



ΕΘΝΙΚΟ ΜΕΤΣΟΒΙΟ ΠΟΛΥΤΕΧΝΕΙΟ  
ΣΧΟΛΗ ΗΛΕΚΤΡΟΛΟΓΩΝ ΜΗΧΑΝΙΚΩΝ ΚΑΙ ΜΗΧΑΝΙΚΩΝ ΥΠΟΛΟΓΙΣΤΩΝ  
ΤΟΜΕΑΣ ΤΕΧΝΟΛΟΓΙΑΣ ΠΛΗΡΟΦΟΡΙΚΗΣ ΚΑΙ ΥΠΟΛΟΓΙΣΤΩΝ  
ΕΡΓΑΣΤΗΡΙΟ ΛΟΓΙΚΗΣ ΚΑΙ ΕΠΙΣΤΗΜΗΣ ΥΠΟΛΟΓΙΣΜΩΝ

Ευσταθή και ασταθή σημεία ισορροπίας σε  
εκμάθηση χωρίς regret και ενθόρυβα μοντέλα

ΔΙΠΛΩΜΑΤΙΚΗ ΕΡΓΑΣΙΑ

Αγγελική Γιάννου

Επιβλέπων: Δημήτριος Φωτάκης  
Αναπληρωτής Καθηγητής Ε.Μ.Π.

Αθήνα, 06/06/2021

Αγγελική Γιάννου





ΕΘΝΙΚΟ ΜΕΤΣΟΒΙΟ ΠΟΛΥΤΕΧΝΕΙΟ  
ΣΧΟΛΗ ΗΛΕΚΤΡΟΛΟΓΩΝ ΜΗΧΑΝΙΚΩΝ ΚΑΙ ΜΗΧΑΝΙΚΩΝ ΥΠΟΛΟΓΙΣΤΩΝ  
ΤΟΜΕΑΣ ΤΕΧΝΟΛΟΓΙΑΣ ΠΛΗΡΟΦΟΡΙΚΗΣ ΚΑΙ ΥΠΟΛΟΓΙΣΤΩΝ

Ευσταθή και ασταθή σημεία ισορροπίας σε εκμάθηση  
χωρίς regret και ενθόρυβα μοντέλα

ΔΙΠΛΩΜΑΤΙΚΗ ΕΡΓΑΣΙΑ

Αγγελική Γιάννου

Επιβλέπων: Δημήτριος Φωτάκης  
Αναπληρωτής Καθηγητής Ε.Μ.Π.

Εγκρίθηκε από την τριμελή εξεταστική επιτροπή την 06/06/2021.

.....  
Δημήτριος Φωτάκης  
Αναπληρωτής Καθηγητής Ε.Μ.Π.

.....  
Παναγιώτης Μερτικόπουλος  
Principal Researcher CNRS

.....  
Αριστέιδης Παγουρτζής  
Αναπληρωτής Καθηγητής Ε.Μ.Π.

.....  
**Αγγελική Γιάννου**

(Διπλωματούχος Ηλεκτρολόγος Μηχανικός & Μηχανικός Υπολογιστών Ε.Μ.Π.)

Οι απόψεις που εκφράζονται σε αυτό το κείμενο είναι αποκλειστικά του συγγραφέα και δεν αντιπροσωπεύουν απαραίτητα την επίσημη θέση του Εθνικού Μετσόβιου Πολυτεχνείου.

Απαγορεύεται η χρήση της παρούσας εργασίας για εμπορικούς σκοπούς.

This work is licensed under a [Creative Commons “Attribution-NonCommercial-ShareAlike 4.0 International”](https://creativecommons.org/licenses/by-nc-sa/4.0/) license.



Αγγελική Γιάννου, 2021

## Περίληψη

Σε αυτή τη διπλωματική εργασία θα εξετάσουμε τη σύγκλιση no-regret διακριτών αλγορίθμων σε σημεία Nash ισορροπίας, μελετώντας πεπερασμένα παίγνια  $N$  παικτών. Παρά το αυξανόμενο ενδιαφέρον μελέτης των no-regret αλγορίθμων, λόγω των πολυποίκιλων εφαρμογών τους, λίγα είναι γνωστά για την πραγματική συμπεριφορά τους (long-run behavior) σε περιβάλλοντα όπου εμπλέκονται πολλοί παίκτες· ενώ συνήθως τα μέχρι τώρα γνωστά αποτελέσματα αφορούν συγκεκριμένες κλάσεις παιγνίων. Σε αυτή τη διπλωματική αντί να εστιάσουμε σε μία συγκεκριμένη κλάση παιγνίων, θα εστιάσουμε στις διαφορετικές κατηγορίες σημείων Nash ισορροπίας. Συγκεκριμένα, θα μελετήσουμε το σύνολο των αλγορίθμων "Follow the Regularized Leader" (FTRL) με διαφορετικές θορυβώδεις ανατροφοδοτήσεις σήματος - από την περίπτωση όπου οι παίκτες έχουν πρόσβαση σε ένα oracle, μέχρι και την bandit περίπτωση, όπου οι παίκτες έχουν πρόσβαση μόνο σε μία τιμή. Σε αυτό το πλαίσιο, θα εδραιώσουμε την εξής ισοδυναμία: ένα σημείο Nash ισορροπίας είναι ευσταθές αν και μόνο αν είναι ένα *strict* σημείο Nash ισορροπίας.

**Λέξεις κλειδιά:** Θεωρία παιγνίων, Follow the Regularized Leader, Εχμάθηση σε περιβάλλοντα πολλών παικτών, Bandits.

## Abstract

In this diploma thesis, we examine the Nash equilibrium convergence properties of no-regret learning in general  $N$ -player games. Despite the importance and widespread applications of no-regret algorithms, their long-run behavior in multi-agent environments is still far from understood, and most of the literature has focused by necessity on certain, specific classes of games (typically zero-sum or congestion games). Instead of focusing on a fixed class of *games*, we instead take a structural approach and examine different classes of *equilibria* in generic games. For concreteness, we focus on the archetypal "follow the regularized leader" (FTRL) class of algorithms, and we consider the full spectrum of information uncertainty that the players may encounter – from noisy, oracle-based feedback, to bandit, payoff-based information. In this general context, we establish a comprehensive equivalence between the stability of a Nash equilibrium and its support: a Nash equilibrium is stable and attracting with arbitrarily high probability if and only if it is strict (i.e., each equilibrium strategy has a unique best response). This result extends existing continuous-time versions of the "folk theorem" of evolutionary game theory to a bona fide discrete-time learning setting, and provides an important link between the literature on multi-armed bandits and the equilibrium refinement literature.

**Keywords:** Online Learning, Follow the Regularized Leader, Game Theory, Multi-agent Learning, Bandits.

## Ευχαριστίες

Με την ολοκλήρωση αυτής της διπλωματικής εργασίας ολοκληρώνεται και ένας κύκλος της ζωής μου, αυτός των φοιτητικών μου χρόνων. Θα ήθελα λοιπόν να ευχαριστήσω όλους αυτούς που με στήριξαν σε όλη μου την πορεία. Αρχικά, θα ήθελα να ευχαριστήσω την οικογένεια μου αλλά και τους φίλους μου, οι οποίοι αποτελούν για μένα μία ευρύτερη οικογένεια. Εν συνεχεία, θα ήθελα να ευχαριστήσω τον κ. Φωτάκη για την εμπιστοσύνη που μου έδειξε και έγινε επιβλέπωντας αυτής της διπλωματικής εργασίας. Τελευταίους θα ήθελα να ευχαριστήσω το Μανώλη Βλατάκη και τον Παναγιώτη Μερτικόπουλο που χωρίς τη βοήθεια τους αυτή η διπλωματική εργασία δε θα είχε πραγματοποιηθεί. Πιο συγκεκριμένα, θα ήθελα να τους ευχαριστήσω για όλες τις πολύωρες και διαφωτιστικές συζητήσεις που είχαμε, τη γνώση που μου μετέδωσαν σε ένα εντελώς άγνωστο σε μένα πεδίο έρευνας, την εμπιστοσύνη που μου έδειξαν και την στήριξη που μου παρείχαν τον τελευταίο ενάμιση χρόνο.





# Table of contents

<b>1 Εκτεταμένη Ελληνική Περίληψη</b>	<b>1</b>
1.1 Παίγνια N-παικτών & Σημεία ισορροπίας	1
1.1.1 Πεπερασμένα παίγνια σε κανονική μορφή	1
1.1.2 Σημεία Nash ισορροπίας	2
1.2 Ελαχιστοποίηση regret και Εξομάλυνση	3
1.2.1 Μοντέλο ανατροφοδότησης	3
1.2.2 Εξομάλυνση	4
1.3 Αποτελέσματα	4
1.3.1 Ασυμπτωτική Ευστάθεια	4
1.3.2 Θεωρήματα	5
<b>2 Introduction</b>	<b>7</b>
2.1 Our contributions	8
2.2 Related work	8
2.3 Proof techniques	9
<b>3 Preliminaries</b>	<b>11</b>
3.1 Finite games in normal form	11
3.2 Solution concepts	12
3.2.1 Dominated strategies	13
3.2.2 Nash equilibrium	13
3.3 No regret learning and Regularization	14
3.3.1 Regret	14
3.3.2 Follow the Regularized Leader (FTRL)	14
<b>4 Analysis and Results</b>	<b>21</b>
4.1 The algorithm	21
4.1.1 The feedback model	21
4.1.2 Regularization	23
4.2 Asymptotic Stability	24
4.3 Main Results	25
4.4 Our Techniques	25
4.4.1 The Stochastic Asymptotic Stability of Strict Nash Equilibria	25
4.4.2 The Stochastic Instability of Mixed Nash Equilibria	27
<b>5 Future work</b>	<b>29</b>

---

<b>A' Theoretical Basis</b>	<b>35</b>
A'.1 Elements of martingale limit theory . . . . .	35
A'.1.1 Basic definitions . . . . .	35
A'.1.2 Conditional Expectation . . . . .	36
A'.1.3 Martingales . . . . .	37
A'.1.4 Martingale limit theorems . . . . .	37
A'.2 Elements of Convex Analysis . . . . .	38
A'.2.1 Basic definitions . . . . .	38
A'.2.2 Convexity & Duality . . . . .	39
A'.2.3 Convexity and Smoothness . . . . .	41
<b>B' Deferred Proofs</b>	<b>49</b>
B'.1 Bregman Divergence and Fenchel Coupling . . . . .	49
B'.1.1 Bregman Divergence . . . . .	49
B'.1.2 Steep vs non-steep . . . . .	51
B'.1.3 Polar Cone . . . . .	51
B'.1.4 Fenchel Coupling . . . . .	52
B'.2 Variational stability . . . . .	54
B'.3 Proofs of assumptions for Model 1, Model 2 . . . . .	55
B'.4 Proofs of Stability . . . . .	55
B'.4.1 Deferred Proof of theorem 4.3.1 . . . . .	55
B'.5 Proofs of Instability . . . . .	61

# List of figures

3.1	Illustration of the simplex in different dimensions. . . . .	11
3.2	Regularizers . . . . .	15
A'.1	Depiction of the epigraph of two functions . . . . .	38
A'.2	Representing the blue line with parameters $\rho, \theta$ . . . . .	40
A'.3	Geometric meaning of the conjugate transform . . . . .	41
B'.1	The level sets of KL-divergence . . . . .	49
B'.2	The polar cone corresponding to different points of the simplex. For an interior point this is a line perpendicular to the simplex. For a point of the boundary, it is a plane perpendicular to the simplex tangential to the point of the boundary. For an edge the polar cone corresponds to a cone. . . . .	52
B'.3	(VS) states that the payoff vectors are pointing "towards" the equilibrium . . . .	54

# Chapter 1

## Εκτεταμένη Ελληνική Περίληψη

Αυτό το κεφάλαιο περιλαμβάνει μία περιληπτική παρουσίαση των περιεχομένων αυτής της διπλωματικής εργασίας στα ελληνικά. Εισάγουμε όλες τις βασικές έννοιες που παρουσιάζονται στο κύριο μέρος του κειμένου της διπλωματικής αυτής στα αγγλικά. Ωστόσο, δεν δίνουμε ούτε αποδείξεις ούτε τεχνικές λεπτομέρειες. Αυτές δίνονται εκτενώς στα appendices.

### 1.1 Παίγνια N-παικτών & Σημεία ισορροπίας

#### 1.1.1 Πεπερασμένα παίγνια σε κανονική μορφή

Σε αυτή τη διπλωματική εργασία θα εστιάσουμε σε πεπερασμένα παίγνια σε κανονική μορφή.

**Ορισμός 1.** Ένα τέτοιο παίγνιο ορίζεται σαν μία τούπλα  $\Gamma = \Gamma(\mathcal{N}, \mathcal{A}, u)$  με τα εξής στοιχεία:

- Ένα πεπερασμένο σύνολο παικτών που απαριθμείται ως  $i \in \mathcal{N} = \{1, \dots, N\}$ .
- Ένα πεπερασμένο σύνολο από αμγείς στρατηγικές που απαριθμείται ως  $\alpha_i \in \mathcal{A}_i = \{1, \dots, A_i\}$ ,  $i \in \mathcal{N}$ . Οι παίκτες μπορούν επίσης να παίζουν μικτές στρατηγικές, οι οποίες αναπαριστούν πιθανοτικές κατανομές πάνω στο σύνολο των αμγών στρατηγικών τους και συμβολίζονται ως  $x_i \in \mathcal{X}_i := \Delta(\mathcal{A}_i)$ <sup>1</sup>. σε αυτή την περίπτωση, με  $x_{i\alpha_i}$  συμβολίζουμε την πιθανότητα με την οποία ο παίκτης  $i \in \mathcal{N}$  selects  $\alpha_i \in \mathcal{A}_i$ . Επιπλέον, αναφερόμενοι στο σύνολο των παικτών, θα γράφουμε  $x = (x_1, \dots, x_N)$  για ένα μικτό προφίλ στρατηγικών και  $\mathcal{X} := \prod_i \mathcal{X}_i$  για το σύνολο στο οποίο ανήκουν όλα τα προφίλ αυτά. Τέλος, όταν θέλουμε να εστιάσουμε στην στρατηγική ενός μόνο παίκτη  $i \in \mathcal{N}$ , θα χρησιμοποιούμε τη συντομογραφία  $(x_i; x_{-i}) := (x_1, \dots, x_i, \dots, x_N)$  – και αντιστοίχως,  $(\alpha_i; \alpha_{-i})$  για αμγείς στρατηγικές.
- Ένα σύνολο συναρτήσεων πληρωμής  $u_i: \mathcal{A} \rightarrow \mathbb{R}$  όπου  $\mathcal{A} := \prod_i \mathcal{A}_i$  είναι ο χώρος όλων των προφίλ αμγών στρατηγικών. Η αναμενόμενη πληρωμή του παίκτη  $i$  ως προς ένα προφίλ μικτών στρατηγικών  $x \in \mathcal{X}$  είναι

$$u_i(x) \equiv u_i(x_i; x_{-i}) = \sum_{\alpha_1 \in \mathcal{A}_1} \cdots \sum_{\alpha_N \in \mathcal{A}_N} u_i(\alpha_1, \dots, \alpha_N) \cdot x_{1,\alpha_1} \cdots x_{N,\alpha_N} \quad (1.1)$$

<sup>1</sup>Με  $\Delta$  συμβολίζουμε το simplex; μία αναπαράσταση του οποίου φαίνεται στην εικόνα 3.1

όπου  $u_i(\alpha_1, \dots, \alpha_N)$  είναι η πληρωμή του παίκτη  $i$  στο προφίλ αμιγών στρατηγικών  $\alpha = (\alpha_1, \dots, \alpha_N) \in \mathcal{A}$ .

Επιπλέον θα γράφουμε  $v_{i\alpha_i}(x) = u_i(\alpha_i; x_{-i})$  για την πληρωμή που ο παίκτης  $i$  θα εκλάμβανε αν επέλεγε να παίξει τη στρατηγική  $\alpha_i \in \mathcal{A}_i$  εναντίον του προφίλ μικτών στρατηγικών  $x_{-i}$  όλων των άλλων παικτών. Έτσι, το μικτό διάνυσμα πληρωμών του  $i$ -στου παίκτη είναι

$$v_i(x) = (v_{i\alpha_i}(x))_{\alpha_i \in \mathcal{A}_i} \quad (1.2)$$

και θα γράφουμε  $v(x) = (v_1(x), \dots, v_N(x))$  για το σύνολο αυτών. Για απλότητα στο συμβολισμό, θα ορίσουμε  $\mathcal{Y}_i = \mathbb{R}^{\mathcal{A}_i}$  και  $\mathcal{Y} = \prod_i \mathcal{Y}_i$  για το χώρο των διανυσμάτων πληρωμών και των προφίλ αυτών αντιστοίχως. Τέλος, θα αναγνωρίζουμε την αμιγή στρατηγική  $\alpha_i$  ως τη μικτή στρατηγική που αναθέτει πιθανότητα 1 στην  $\alpha_i$ , και θα ορίσουμε το αντίστοιχο αμιγές διάνυσμα πληρωμής ως  $v_i(\alpha) = (u_i(\alpha_i; \alpha_{-i}))_{\alpha_i \in \mathcal{A}_i}$ . Η διαφορά μεταξύ αμιγών και μικτών διανυσμάτων πληρωμής θα αποκτήσει σημασία στη συνέχεια.

### 1.1.2 Σημεία Nash ισορροπίας

Η πιο γνωστή έννοια λύσης σε παίγνια είναι αυτή του σημείου Nash ισορροπίας, δηλαδή ένα προφίλ μικτής στρατηγικής το οποίο αποθαρρύνει τους παίκτες μονομερώς να επιλέξουν κάποια άλλη στρατηγική. Ο Nash απέδειξε στο [1] ότι όλα τα πεπερασμένα παίγνια  $N$ -παικτών επιδέχονται τουλάχιστον ένα σημείο Nash ισορροπίας.

**Ορισμός 2.** Ένα σημείο  $x^*$  είναι ένα σημείο Nash ισορροπίας του παιγνίου  $\Gamma$  αν

$$u_i(x^*) \geq u_i(x_i; x_{-i}^*) \quad \text{για κάθε } x_i \in \mathcal{X}_i \text{ και για κάθε } i \in \mathcal{N}. \quad (\text{NE})$$

Το σύνολο των αμιγών στρατηγικών που υποστηρίζονται στη συνιστώσα του σημείου ισορροπίας  $x_i^* \in \mathcal{X}_i$  για κάθε παίκτη, θα συμβολίζεται ως  $\text{supp}(x_i^*) = \{\alpha_i \in \mathcal{A}_i : x_{i\alpha_i}^* > 0\}$ . Αντιστοίχως, τα σημεία Nash ισορροπίας μπορούν να χαρακτηριστούν μέσω της ανισότητας

$$v_{i\alpha_i^*}(x^*) \geq v_{i\alpha_i}(x^*) \quad \text{για κάθε } \alpha_i^* \in \text{supp}(x_i^*) \text{ και για κάθε } \alpha_i \in \mathcal{A}_i, i \in \mathcal{N}. \quad (1.3)$$

Απόρροια του παραπάνω χαρακτηρισμού των σημείων Nash ισορροπίας είναι η ακόλουθη ταξινόμηση αυτών:

- $x^*$  είναι ένα αμιγές σημείο ισορροπίας αν το  $\text{supp}(x_i^*)$  περιέχει μία μόνο στρατηγική για κάθε παίκτη  $i \in \mathcal{N}$ .
- $x^*$  είναι ένα σημείο μικτής ισορροπίας σε οποιαδήποτε άλλη περίπτωση. Συγκεκριμένα, αν το  $\text{supp}(x_i^*) = \mathcal{A}_i$  για κάθε  $i \in \mathcal{N}$ , τότε το  $x^*$  καλείται πλήρως μικτό.

Εξ ορισμού, τα σημεία αμιγής ισορροπίας αντιστοιχούν σε κορυφές του χώρου. By definition  $\mathcal{X}$ , ένω τα πλήρως μικτά σημεία ισορροπίας βρίσκονται στο σχετικό εσωτερικό  $\text{ri}(\mathcal{X})$  του χώρου  $\mathcal{X}$ , και γενικότερα τα σημεία μικτής ισορροπίας βρίσκονται στο σχετικό εσωτερικό του πορτρέτου που γεννάται από το  $\text{supp}(x_i^*)$  του κάθε παίκτη.

Μία άλλη ταξινόμηση των σημείων Nash ισορροπίας πηγάζει από την ανισότητα (1.3) και έχει ως εξής: αν αυτή η ανισότητα 1.3 είναι αυστηρή για κάθε  $\alpha_i \in \mathcal{A}_i \setminus \text{supp}(x_i^*)$ ,  $i \in \mathcal{N}$ , το αντίστοιχο σημείο ισορροπίας καλείται σχεδόν-αυστηρό [6]. Τα σχεδόν-αυστηρά σημεία ισορροπίας έχουν την ιδιότητα ότι όλες οι καλύτερες στρατηγικές επιλέγονται με θετική πιθανότητα. Σημειώνεται ότι τα σχεδόν-αυστηρά σημεία ισορροπίας μπορεί να είναι είτε μικτά είτε αμιγή. Τα αμιγή σχεδόν-αυστηρά σημεία ισορροπίας θα καλούνται απλώς αυστηρά.

## 1.2 Ελαχιστοποίηση regret και Εξομάλυνση

Μία βασική απαίτηση στο πεδίο της ενεργούς εκμάθησης είναι η ελαχιστοποίηση του regret των παικτών, δηλαδή της διαφοράς των συσσωρευμένων πληρωμών μεταξύ της στρατηγικής ενός παίκτη και της καλύτερης στρατηγικής που θα μπορούσε να έχει διαλέξει εκ των υστέρων σε βάθος ενός χρονικού ορίζοντα  $T$ . Αυστηρά μιλώντας, δοθείσας μίας ακολουθίας του παιγνίου  $X_n \in \mathcal{X}$ ,  $n = 1, 2, \dots$ , το (εξωτερικό) regret του κάθε παίκτη  $i \in \mathcal{N}$  ορίζεται ως

$$\text{Reg}_i(T) = \max_{x_i \in \mathcal{X}_i} \sum_{n=1}^T [u_i(x_i; X_{-i,n}) - u_i(X_{i,n}; X_{-i,n})] \quad (1.4)$$

και θα λέμε ότι ο παίκτης  $i$  δεν έχει regret αν  $\text{Reg}_i(T) = o(T)$ .

Ένα από τα πιο γνωστά χρησιμοποιούμενα σχήματα ενεργής εκμάθησης για να επιτευχθεί αυτή η απαίτηση είναι η οικογένεια αλγορίθμων *Follow the Regularized Leader* (FTRL) [35, 22]. Συγκεκριμένα, σε κάθε βήμα της διαδικασίας εκμάθησης ο αλγόριθμος (FTRL) αποδίδει τη μικτή στρατηγική που μεγιστοποιεί τη συσσωρευμένη πληρωμή του παίκτη σε συνδυασμό με έναν εξομαλυντή. Έχουμε λοιπόν της εξής βήμα προς βήμα απεικόνιση

$$\begin{aligned} X_{i,n} &= Q_i(Y_{i,n}) \\ Y_{i,n+1} &= Y_{i,n} + \gamma_n \hat{v}_{i,n} \end{aligned} \quad (\text{FTRL})$$

όπου  $Q_i: \mathcal{Y}_i \rightarrow \mathcal{X}_i$  είναι η συνάρτηση επιλογής του παίκτη  $i \in \mathcal{N}$ ,  $\gamma_n > 0$  είναι ο ρυθμός εκμάθησης, τέτοιος ώστε  $\sum_n \gamma_n = \infty$ , και  $\hat{v}_{i,n}$  είναι ένα "σήμα πληρωμής" που παρέχει μία εκτίμηση των μικτών πληρωμών του παίκτη  $i$  στο βήμα  $n$ . Παρακάτω συζητάμε αναλυτικά όλες αυτές τις συνιστώσες.

### 1.2.1 Μοντέλο ανατροφοδότησης

Έχοντας ως στόχο να συμπεριλάβουμε διαφορετικού τύπου ανατροφοδοτήσεις στο μοντέλο μας, κάνουμε τις παρακάτω συνηθισμένες υποθέσεις το σήμα πληρωμής:

$$\hat{v}_n = v(X_n) + \xi_n \quad (1.5)$$

για κάποια γενική διαδικασία λάθους  $\xi_n = (\xi_{i,n})_{i \in \mathcal{N}}$ . Για να διαχωρίσουμε μεταξύ του μηδενικής μέσης τιμής και μη μηδενικής μέσης τιμής λάθους, αναλύουμε περαιτέρω το  $\xi_n$  σε  $\xi_n = Z_n + b_n$ , όπου

$$b_n = \mathbb{E}[\xi_n | \mathcal{F}_n] \quad \text{and} \quad \mathbb{E}[Z_n | \mathcal{F}_n] = 0 \quad (1.6)$$

όπου  $\mathcal{F}_n$  περιλαμβάνει γνώση για όλα τα  $X_n$  μέχρι και το βήμα  $n$ <sup>2</sup>. Έτσι για το σήμα ανατροφοδότησης  $\hat{v}_n$  χαρακτηρίζεται μέσω των παρακάτω στατιστικών

$$a) \text{ Συστηματικό σφάλμα: } \quad \mathbb{E}[\|b_n\|_* | \mathcal{F}_n] \leq B_n \quad (1.7\alpha')$$

$$b) \text{ Απόκλιση: } \quad \mathbb{E}[\|\hat{v}_n\|_*^2 | \mathcal{F}_n] \leq M_n^2 \quad (1.7\beta')$$

όπου  $B_n$  και  $M_n$  είναι ντετερμινιστικά φράγματα του συστηματικού σφάλματος και της απόκλισης του σήματος ανατροφοδότησης  $\hat{v}_n$ . Επιπλέον, θεωρούμε ως δεδομένες τις παρακάτω υποθέσεις:

$$(A1) \text{ Έλεγχος του συστηματικού σφάλματος: } \lim_{n \rightarrow \infty} B_n = 0 \text{ and } \sum_n \gamma_n B_n < \infty.$$

$$(A2) \text{ Έλεγχος της απόκλισης: } \sum_n \gamma_n^2 M_n^2 < \infty.$$

<sup>2</sup>Φυσικά, αφού το σήμα ανατροφοδότησης γεννάται μετά την επιλογή στρατηγικής από τους παίκτες,  $\hat{v}_n$  δεν είναι  $\mathcal{F}_n$ -μετρήσιμο στη γενική περίπτωση.

- (A3) *Κουτόντες παρατηρήσεις λάθους στο σημείο ισορροπίας:* Για κάθε μικτό σημείο ισορροπίας Nash  $x^*$  του  $\Gamma$  και για κάθε  $n = 1, 2, \dots$ , υπάρχει ένας παίκτης  $i \in \mathcal{N}$  και στρατηγικές  $a, b \in \text{supp}(x_i^*)$  τέτοια ώστε

$$\mathbb{P}(|\hat{v}_{ia,n} - \hat{v}_{ib,n}| \geq \beta \mid \mathcal{F}_n) > 0 \quad \text{για κάθε επαρκώς μικρό } \beta > 0. \quad (1.8)$$

Αυτές οι υποθέσεις είναι αρκετά γενικές και επιτρέπουν ένα μεγάλο εύρος διαφορετικών μοντέλων ανατροφοδότησης.

## 1.2.2 Εξομάλυνση

Η δεύτερη συνιστώσα του (FTRL) είναι οι συναρτήσεις επιλογής των παικτών  $Q_i: \mathcal{Y}_i \rightarrow \mathcal{X}_i$ . Με στόχο την αποφυγή πρόωρης προσκόλησης σε μία συγκεκριμένη στρατηγική  $Q_i$  ορίζεται ως μία “εξομαλυσμένη” εκδοχή της καλύτερης απόκρισης  $y_i \mapsto \arg \max_{x_i \in \mathcal{X}_i} \langle y_i, x_i \rangle$ . Έτσι, επικεντρωνόμαστε στις εξομαλυσμένες καλύτερες αποκρίσεις που ορίζονται ως

$$Q_i(y_i) = \arg \max_{x_i \in \mathcal{X}_i} \langle y_i, x_i \rangle - h_i(x_i). \quad (1.9)$$

Ο εξομαλυντής κάθε παίκτη  $h_i: \mathcal{X}_i \rightarrow \mathbb{R}$  ορίζεται ως  $h_i(x_i) = \sum_{\alpha_i \in \mathcal{A}_i} \theta_i(x_i)$  για κάποια συνάρτηση πυρήνα  $\theta_i: [0, 1] \rightarrow \mathbb{R}$  που έχει τις εξής ιδιότητες:

- (i)  $\theta_i$  είναι συνεχής στο  $[0, 1]$
- (ii)  $C^2$ -ομαλή στο  $(0, 1]$
- (iii)  $\inf_{[0,1]} \theta_i'' > 0$ .

Φυσικά διαφορετικοί εξομαλυντές συνεπάγονται διαφορετικές εκδοχές του (FTRL). Παρακάτω παρουσιάζουμε δυο χαρακτηριστικά παραδείγματα.

**Example 1.2.1** (Multiplicative/Exponential weights update). Μία γνωστή επιλογή εξομαλυντή είναι η αρνητική εντροπία  $h_i(x) = \sum_i x_i \log x_i$ , που οδηγεί στη συνάρτηση επιλογής  $\Lambda_i(y) = \exp(y_i) / \sum_j \exp(y_j)$  και ακολούθως στον αλγόριθμο γνωστό ως *multiplicative weights update* (MWU), cf. [19, 50, 21, 20, 22].

**Example 1.2.2** (Euclidean projection). Μία άλλη συνηθισμένη επιλογή εξομαλυντή είναι η τετραγωνική  $h_i(x) = \sum_i x_i^2 / 2$ , η οποία συνεπάγεται τη συνάρτηση επιλογής  $\boxtimes_i(y) = \arg \min_{x \in \Delta} \|y - x\|^2$ , cf. [23, 25].

## 1.3 Αποτελέσματα

Στόχος μας είναι να μελετήσουμε τη σε βάθος χρόνου συμπεριφορά του (FTRL). Η ερώτηση που θα μας απασχολήσει είναι *Ποια σημεία Nash ισορροπίας έχουν ιδιότητες σύγκλισης και ευστάθειας οι οποίες δεν επηρεάζονται από την αβεβαιότητα που περιλαμβάνεται στο σήμα ανατροφοδότησης;*

### 1.3.1 Ασυμπτωτική Ευστάθεια

Αρχικά, αξίζει να σημειωθεί ότι σε γενικά παίγνια μπορεί να υπάρχουν πάνω από ένα σημεία Nash ισορροπίας, είτε μικτά είτε αμιγή ή και τα δύο. Ως εκ τούτου τα αποτελέσματα μας είναι λογικό να ισχύουν τοπικά· έτσι θα εστιάσουμε στην έννοια της στοχαστικής ασυμπτωτικής ευστάθειας [7, 8, 9]. Ευριστικά, ένα σημείο ισορροπίας είναι στοχαστικά ευσταθές αν οποιαδήποτε ακολουθία του παιγνίου, η οποία ξεκινά αρκετά κοντά στο σημείο ισορροπίας παραμένει κοντά με μεγάλη πιθανότητα· επιπροσθέτως, αν η ακολουθία συγκλίνει εν τέλει στο σημείο ισορροπίας το σημείο καλείται στοχαστικά ασυμπτωτικά ευσταθές. Αυστηρά μιλώντας:

**Ορισμός 3.** Έστω  $x^* \in \mathcal{X}$  ένα σημείο Nash ισορροπίας. Καθορίζοντας κάποιο αυθαίρετο επίπεδο εμπιστοσύνης  $\delta > 0$  και μία γειτονιά  $U$  του  $x^*$ , τότε το  $x^* \in \mathcal{X}$  καλείται

1. **Στοχαστικά ευσταθές** αν υπάρχει γειτονιά  $U_0$  του  $x^*$  τέτοια ώστε οποτεδήποτε ισχύει  $X_0 = Q(Y_0) \in U_0$ , έχουμε ότι

$$\mathbb{P}(X_n \in U \text{ για κάθε } n = 0, 1, \dots) \geq 1 - \delta \quad (1.10)$$

2. **Συγκλίνον** αν υπάρχει γειτονιά  $U_0$  του  $x^*$  τέτοια ώστε

$$\mathbb{P}(\lim_{n \rightarrow \infty} X_n = x^*) \geq 1 - \delta \quad (1.11)$$

οποτεδήποτε  $X_0 = Q(Y_0) \in U_0$ .

3. **Στοχαστικά ασυμπτωτικά ευσταθές** αν είναι στοχαστικά ευσταθές και συγκλίνον.

Ο ορισμός 3 είναι σημαντικός για την ανάλυση μας και για αυτό το λόγο παραθέτουμε κάποιες παρατηρήσεις.

**Παρατήρηση 1.** Μία πρώτη λεπτομέρεια που αξίζει να σημειωθεί στον παραπάνω ορισμό είναι αυτή της μεγάλης πιθανότητας: πράγματι, υπό την επήρεια της αβεβαιότητας, μία και μόνο λάθος εκτίμηση των διανυσμάτων πληρωμής των παικτών θα μπορούσε να οδηγήσει την ακολουθία  $X_n$  εκτός της γειτονιάς του  $x^*$ , πιθανότατα χωρίς να επιστρέψει ποτέ. Έχοντας αυτό στα υπόψη μας είναι αναμενόμενο τα αποτελέσματά μας να μην ισχύουν με πιθανότητα 1, αλλά με αυθαίρετα μεγάλη πιθανότητα.

**Παρατήρηση 2.** Μία άλλη παρατήρηση που αξίζει να γίνει είναι πως η απαίτηση  $X_0 = Q(Y_0) \in U_0$  υπονοεί πως κάποιες στρατηγικές στο χώρο  $\mathcal{X}$  δεν είναι επιτρεπτές ως αρχικές συνθήκες. Επιστρέφοντας πίσω στα δύο χαρακτηριστικά παραδείγματα του (FTRL), MWU 1.2.1, Projection GD 1.2.2, υπάρχει μία διχοτομία σε ότι αφορά τις ιδιότητες των αντίστοιχων συναρτήσεων επιλογής. Από τη μία πλευρά, ο πυρήνας του Ευκλείδειου εξομαλυντή είναι παντού παραγωγίσιμος σε όλο το διάστημα  $[0, 1]$ . Από την άλλη πλευρά, η παράγωγος του πυρήνα της αρνητικής Shannon-εντροπίας πάει στο  $-\infty$  καθώς το  $x$  πάει στο 0. Αυτό σημαίνει ότι στη δεύτερη περίπτωση τα σύνορα δεν είναι επιτρεπτά και έτσι κάποιες αρχικές συνθήκες δεν ανήκουν στην εικόνα  $\text{im } Q$ . Αυτή η διχοτομία αναλύεται εκτενώς στην ενότητα B.1.2.

### 1.3.2 Θεωρήματα

Έχοντας ορίσει όλα τα παραπάνω είμαστε σε θέση να παρουσιάσουμε τα αποτελέσματά μας.

**Κύριο θεώρημα.** Αν οι υποθέσεις (A1)–(A3) ισχύουν, τότε:  
 $x^*$  είναι ένα αυστηρό σημείο Nash ισορροπίας  $\iff x^*$  είναι στοχαστικά ασυμπτωτικά ευσταθές για τον (FTRL)

**Θεώρημα 1.** Έστω  $x^* \in \mathcal{X}$  ένα αυστηρό σημείο Nash ισορροπίας του  $\Gamma$ . Αν ο αλγόριθμος (FTRL) τρέχει με ημιτελές σήμα ανατροφοδότησης που ικανοποιεί τις υποθέσεις (A1) και (A2), τότε το σημείο  $x^*$  είναι στοχαστικά ασυμπτωτικά ευσταθές.

**Θεώρημα 2.** Έστω  $x^*$  ένα σημείο μκτής Nash ισορροπίας του  $\Gamma$ . Αν ο αλγόριθμος (FTRL) τρέχει με ημιτελές σήμα ανατροφοδότησης που ικανοποιεί την υπόθεση (A3), τότε το σημείο  $x^*$  δεν είναι στοχαστικά ασυμπτωτικά ευσταθές.





## Chapter 2

# Introduction

The prototypical framework for online learning in games can be summarized as follows:

1. At each stage of the process, every participating agent chooses an action from some finite set.
2. All agents receive a reward based on the actions of all other players and their individual payoff functions (assumed a priori unknown).
3. The players record their rewards and any other feedback generated during the payoff phase, and the process repeats.

This multi-agent framework has both important similarities and major differences with *single-agent* online learning. Indeed, if we isolate a single, focal player and abstract away all others, we essentially recover a multi-armed bandit (MAB) problem – stochastic or adversarial, depending on the assumptions for the non-focal players [13, 14]. In this case, the most widely used figure of merit is the agent’s *regret*, i.e., the difference between the agent’s cumulative payoff and that of the best fixed action in hindsight. Accordingly, much of the literature on online learning has focused on deriving regret bounds that are min-max optimal, both in terms of the horizon  $T$  of the process, as well as the number of actions  $A$  available to the focal player.

On the other hand, from a game-theoretic standpoint, the main question that arises is whether players eventually settle on an equilibrium profile from which no player has an incentive to deviate. In this regard, a “folk” result states that the empirical frequency of play under no-regret play converges to the game’s set of *coarse correlated equilibria* (CCE) [27, 28]. However, there are two key caveats with this result. First, CCE are considerably weaker than Nash equilibria, to the extent that they fail even the most basic postulates of rationalizability [24]: as was shown by [29], CCE may be supported *exclusively* on *strictly dominated* strategies, even in simple, symmetric two-player games. Second, the convergence of the empirical mean does not carry any tangible guarantees for the players’ day-to-day behavior: under this type of convergence, the player’s best payoff over time could be close to that of a Nash equilibrium, but the players might otherwise be spending arbitrarily long periods of time on dominated strategies.

The above is just a well-known example of the convergence failures of no-regret learning in games with a possibly exotic equilibrium structure. More to the point, even when the underlying game admits a *unique* Nash equilibrium, recent works have shown that no-regret algorithms – such as the popular multiplicative weights update (MWU) method – could still lead to chaotic [30, 31, 32] or Poincaré recurrent / cycling behavior [33, 16, 34]. From a convergence viewpoint, all these results can be seen as instances of a much more general impossibility result at play:

there are no uncoupled dynamics leading to Nash equilibrium in all games [Hart and Mas-Colell, [42]].<sup>1</sup> Since no-regret dynamics are by definition unilateral, they are *a fortiori* uncoupled, so this result shatters any hope of obtaining a universal Nash equilibrium convergence result for the players’ day-to-day behavior.

## 2.1 Our contributions

In view of the above, a critical question that arises is the following: *Is there a class of Nash equilibria that consistently attract no-regret processes?* Conversely, *are all Nash equilibria equally likely to emerge as outcomes of a no-regret learning process?*

To address these questions in as general a setting as possible, we focus on the “follow the regularized leader” (FTRL) family of algorithms: this is arguably the most widely used class of dynamics for no-regret learning in games, and it includes as special cases the seminal multiplicative weights / EXP3 algorithms [22, 35, 20]. In terms of feedback, we also consider a flexible, context-agnostic template in which players are only assumed to have access to an inexact model of their payoff vectors at a given stage. This model for the players’ feedback covers a broad range of modeling assumptions, such as ( $\hat{a}$ ) the case where players can retroactively compute – or otherwise observe – their full payoff vectors (e.g., as in routing games); and ( $\beta$ ) the *bandit* case, where players only observe their in-game payoffs and have no other information on the game being played.

The range of modeling assumptions covered by our framework is quite extensive, so one would likewise expect different, context-specific answers to these questions – presumably with equilibria becoming “less stable” as information becomes “more scarce”. This expectation is justified by the behavior of no-regret learning in single-agent environments: there, the type of information available to the learner has a dramatic effect on the achieved regret minimization rate. Nevertheless, we show that this conjecture is *false*: as far as the algorithms’ equilibrium convergence properties are concerned, the learning dynamics described above are all *equivalent*.

In more detail, we show that all FTRL algorithms under study enjoy the following properties:

- $\hat{a}$ ) *Strict Nash equilibria are stochastically asymptotically stable* – i.e., they are stable and attracting with arbitrarily high probability.
- $\beta$ ) *Only strict Nash equilibria have this property*: mixed Nash equilibria supported on more than one strategies are inherently unstable from a learning viewpoint.

We are not aware of a similar result in the literature at this level of generality (i.e., including models with bandit feedback), and we believe that this equivalence represents an important refinement criterion for the prediction of the day-to-day behavior of no-regret learners in the face of uncertainty and lack of perfect information.

## 2.2 Related work

To put our contributions in the proper context, we provide below an account of relevant works in the literature, classified along the two directions of our main result: “strictness  $\implies$  stability” and “stability  $\implies$  strictness”.

---

<sup>1</sup>“Uncoupled” means here that each player’s update rule does not depend explicitly on the payoffs of other players.

**I. Strictness  $\implies$  Stability.** Analyzing the convergence of game-theoretic learning dynamics has generated a vast corpus of literature that is impossible to survey here. Nonetheless, an emerging theme in this literature is the focus on specific classes of games (such as potential games or  $2^N$  games). As a purely indicative – and highly incomplete – list, we cite here the works of Leslie and Collins [43] and Leslie [44], Cominetti et al. [45], Kleinberg et al. [41], , Coucheney et al. [46], Syrgkanis et al. [40], and d Cohen et al. [39], who provide a range of equilibrium convergence results in potential,  $2^N$ , and  $(\lambda, \mu)$ -smooth games, under different feedback assumptions – from payoff vector observations [41, 40] to bandit [43, 44, 45, 39]. By contrast, our focus is determining the stochastic stability of a class of *equilibria* – not *games*.

As far as we are aware, the only comparable results in this literature concern an idealized continuous-time, deterministic, full-information version of our setting, which is common in applications to population biology and evolutionary game theory. In this context, building on earlier results on the replicator dynamics [36, 7], the authors of [16] showed that strict Nash equilibria are asymptotically stable under the continuous-time dynamics of FTRL. However, we stress here again that these results only concern continuous-time, deterministic dynamical systems with an inherent full-information assumption; we are not aware of a result providing convergence to strict Nash equilibria with bandit feedback.

**II. Stability  $\implies$  Strictness.** In the converse direction, a related result in the literature on evolutionary games is that only strict Nash equilibria are asymptotically stable under the (multi-population) replicator dynamics [36, 8, 37], a continuous-time, deterministic dynamical system which can be seen as the “mean-field” limit of the exponential weights algorithm [38, 47, 33]. In a much more recent paper [48], this implication was extended to the dynamics of FTRL, but always in a deterministic, full-information, continuous-time setting. In this regard, our results are aligned with [48]; however, other than this high-level conceptual link, there is no precise connection, either at the level of implications or at the level of proofs. Specifically, the analysis of [48] relies crucially on volume-conservation arguments that are neither applicable nor relevant in a discrete-time stochastic setting – where the various processes involved could jump around stochastically without any regard for volume contraction or expansion.

## 2.3 Proof techniques

Learning with partial information is an inherently stochastic process, so our results are also stochastic in nature – hence the requirement for asymptotic stability with arbitrarily high probability. This constitutes a major point of departure from continuous-time models of learning [16, 48], so our proof techniques are also radically different as a result. The principal challenge in our proof of stability of strict Nash equilibria comes in controlling the aggregation of error terms with possibly unbounded variance (coming from inverse propensity scoring of bandit-type observations). Because of this, stochastic approximation techniques that have been used to show convergence with  $L^2$ -bounded feedback [49] cannot be applied in this setting; we achieve this control by applying a sharp version of the Doob-Kolmogorov maximal inequality to control equilibrium deviations with high probability. In the converse direction, the crucial argument in the proof of the *instability* of mixed equilibria is a direct probabilistic estimate which leverages a non-degeneracy argument for the noise entering the process; we are not aware of other works using a similar technique.



# Chapter 3

## Preliminaries

### 3.1 Finite games in normal form

Throughout this diploma thesis we will focus on normal form games with a finite number of players and a finite number of actions per player.

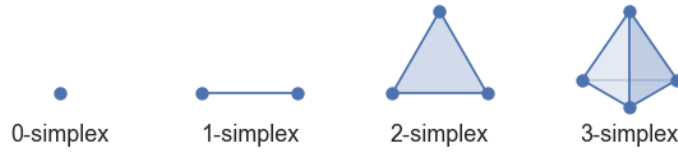


Figure 3.1: Illustration of the simplex in different dimensions.

**Definition 3.1.1.** Such a game is defined as a tuple  $\Gamma = \Gamma(\mathcal{N}, \mathcal{A}, u)$  with the following primitives:

- A finite set of *players* – or *agents* – indexed by  $i \in \mathcal{N} = \{1, \dots, N\}$ .
- A finite set of *actions* – or *pure strategies* – indexed by  $\alpha_i \in \mathcal{A}_i = \{1, \dots, A_i\}$ ,  $i \in \mathcal{N}$ . Players can also play *mixed strategies*, which represent probability distributions  $x_i \in \mathcal{X}_i := \Delta(\mathcal{A}_i)$ <sup>1</sup>; in this case, we will write  $x_{i\alpha_i}$  for the probability that player  $i \in \mathcal{N}$  selects  $\alpha_i \in \mathcal{A}_i$ . Aggregating over all players, we will also write  $x = (x_1, \dots, x_N)$  for the players' *mixed strategy profile* and  $\mathcal{X} := \prod_i \mathcal{X}_i$  for the set thereof. Finally, when we want to focus on the strategy (or action) of a particular player  $i \in \mathcal{N}$ , we will use the shorthand  $(x_i; x_{-i}) := (x_1, \dots, x_i, \dots, x_N)$  – and, similarly,  $(\alpha_i; \alpha_{-i})$  for pure strategies.
- An ensemble of *payoff functions*  $u_i: \mathcal{A} \rightarrow \mathbb{R}$  where  $\mathcal{A} := \prod_i \mathcal{A}_i$  is the space of all pure strategy profiles. The expected payoff of player  $i$  in a mixed strategy profile  $x \in \mathcal{X}$  is then given by

$$u_i(x) \equiv u_i(x_i; x_{-i}) = \sum_{\alpha_1 \in \mathcal{A}_1} \cdots \sum_{\alpha_N \in \mathcal{A}_N} u_i(\alpha_1, \dots, \alpha_N) \cdot x_{1,\alpha_1} \cdots x_{N,\alpha_N} \quad (3.1)$$

<sup>1</sup>With  $\Delta$  we symbolize the simplex; an illustration of the simplex is provided in figure 3.1

where  $u_i(\alpha_1, \dots, \alpha_N)$  is the payoff of player  $i$  in the action profile  $\alpha = (\alpha_1, \dots, \alpha_N) \in \mathcal{A}$ . For posterity, we will also write  $v_{i\alpha_i}(x) = u_i(\alpha_i; x_{-i})$  for the payoff that player  $i$  would have gotten by playing  $\alpha_i \in \mathcal{A}_i$  against the mixed strategy profile  $x_{-i}$  of all other players. In this way, the *mixed payoff vector* of the  $i$ -th player will be

$$v_i(x) = (v_{i\alpha_i}(x))_{\alpha_i \in \mathcal{A}_i} \quad (3.2)$$

and we will write  $v(x) = (v_1(x), \dots, v_N(x))$  for the ensemble thereof. For notational convenience, we will also set  $\mathcal{Y}_i = \mathbb{R}^{\mathcal{A}_i}$  and  $\mathcal{Y} = \prod_i \mathcal{Y}_i$  for the space of payoff vectors and profiles respectively. Finally, in a slight abuse of notation, we will identify  $\alpha_i$  with the mixed strategy that assigns all probability to  $\alpha_i$ , and we will denote the corresponding *pure payoff vector* as  $v_i(\alpha) = (u_i(\alpha_i; \alpha_{-i}))_{\alpha_i \in \mathcal{A}_i}$ . The distinction between pure and mixed payoff vectors will become important later on, when we discuss the information at each player's disposal.

This class of games includes any type of games with finite players and finite action sets, for example zero sum games and potential games. Below we present two well-known examples of such games:

**Example 3.1.1** (Matching pennies). In this game each player flips a coin, if both players' coins are heads or tails player one wins one coin; in any other case player two wins a coin. This game is finite game in normal form consisting of 2 players with action sets  $\mathcal{A}_1 \equiv \mathcal{A}_2 \equiv \{H, T\}$ ; while the payoffs can be seen in the matrix below:

	$H$	$T$
$H$	1/-1	-1/1
$T$	-1/1	1/-1

where player 1 is the "row" player and player 2 is the "column" player. This is an example of a zero-sum game since for any (pure or mixed) players' strategies the sum of the payoffs is always zero.

**Example 3.1.2** (Prisoners' dilemma). In this game, each one of two prisoners who were working together has two choices; either to confess and thus betray the other or to remain silent. The years of sentence depend on what both players will do. Let  $\mathcal{A}_1 \equiv \mathcal{A}_2 \equiv \{B, S\}$ , where  $B$  symbolizes betrayal and  $S$  silence, then the payoff matrix of this game is:

	$S$	$B$
$S$	1/1	9/0
$B$	0/9	6/6

in which again prisoner/player 1 is the "row" player and prisoner/player 2 is the "column" player. The numbers in the matrix symbolize the years of sentence. In this game the smaller the number the better thus the players' payoffs represent losses. In this case we can simply multiply the matrix by  $-1$  and turn the losses into gains leaving intact the structure of the game.

## 3.2 Solution concepts

In terms of solution players need somehow to evaluate the actions they chose to play. It is taken for granted that the players are rational and thus they chose strategies that result in the best possible outcome for them. However this concept is susceptible to many interpretations; below we present two of them.

### 3.2.1 Dominated strategies

A *dominated* strategy is a strategy that results in strictly worst payoff than at least some other strategy no matter what the opponents do. Formally, a strategy  $\alpha \in \mathcal{A}_i$  of player  $i \in \mathcal{N}$  is said to be dominated by a strategy  $b \in \mathcal{A}_i$  if it holds that

$$u_i(\alpha; x_{-i}) < u_i(b; x_{-i}) \text{ for all } x_{-i} \in \mathcal{X}_{-i} \quad (3.3)$$

If the inequality is not strict then we say that the strategy is *weakly dominated*. For example looking at example 3.1.2 the strategy  $S$  is dominated by the strategy  $B$ , since it always results in a worst payoff.

It is reasonable to expect that players will not choose to play dominated strategies and thus these strategies can be eliminated. This elimination could lead either in the existence of only one pure strategy (such as in the example of Prisoners' dilemma) or to a reduced version of the game. Iteratively, players continue to eliminate dominated strategies until there are *no* dominated strategies.

### 3.2.2 Nash equilibrium

The most widely used solution concept is that of a Nash equilibrium, i.e., a mixed strategy profile that discourages unilateral deviations. Nash proved in [1] that all  $N$ -player finite games have at least one *Nash equilibrium*.

**Definition 3.2.1.** A point  $x^*$  is a *Nash equilibrium* of  $\Gamma$  if

$$u_i(x^*) \geq u_i(x_i; x_{-i}^*) \text{ for all } x_i \in \mathcal{X}_i \text{ and all } i \in \mathcal{N}. \quad (\text{NE})$$

The set of pure strategies supported at the equilibrium component  $x_i^* \in \mathcal{X}_i$  of each player will be denoted by  $\text{supp}(x_i^*) = \{\alpha_i \in \mathcal{A}_i : x_{i\alpha_i}^* > 0\}$ . Accordingly, Nash equilibria can be equivalently characterized by means of the variational inequality

$$v_{i\alpha_i^*}(x^*) \geq v_{i\alpha_i}(x^*) \text{ for all } \alpha_i^* \in \text{supp}(x_i^*) \text{ and all } \alpha_i \in \mathcal{A}_i, i \in \mathcal{N}. \quad (3.4)$$

The above characterization gives rise to the following classification of Nash equilibria:

- $x^*$  is a *pure equilibrium* if  $\text{supp}(x_i^*)$  only contains a single strategy for all  $i \in \mathcal{N}$ .
- $x^*$  is a *mixed equilibrium* in any other case; in particular, if  $\text{supp}(x_i^*) = \mathcal{A}_i$  for all  $i \in \mathcal{N}$ , we say that  $x^*$  is *fully mixed*.

By definition, pure equilibria correspond to vertices of  $\mathcal{X}$ , fully mixed equilibria lie in the relative interior  $\text{ri}(\mathcal{X})$  of  $\mathcal{X}$ , and, more generally, mixed equilibria lie in the relative interior of the face of the simplex spanned by the support of each player's equilibrium component.

A further distinction between Nash equilibria that is inherited by the inequality (3.4) is as follows: if the inequality 3.4 holds as a strict inequality for all  $\alpha_i \in \mathcal{A}_i \setminus \text{supp}(x_i^*)$ ,  $i \in \mathcal{N}$ , the equilibrium in question is said to be *quasi-strict* [6]. Quasi-strict equilibria have the defining property that *all pure best responses* are played with positive probability; it is also well known that all Nash equilibria in all but a measure-zero set of games are quasi-strict. For this reason, the property of having a quasi-strict equilibrium is generic, and games that enjoy this property are called themselves *generic*.<sup>2</sup>

<sup>2</sup>Specifically, the set of games with Nash equilibria that are not quasi-strict is *meager* in the Baire category sense.



We stress here by looking at examples 3.1.1, 3.1.2 that quasi-strict equilibria could be either mixed or pure. The equilibrium of Matching Pennies is if both players play each one of their strategies with probability 1/2. Thus it is fully mixed and quasi-strict since if any of the two players unilaterally deviates from the equilibrium point results in a strictly worst payoff. Whereas the equilibrium of the Prisoner's dilemma is quasi-strict and pure and it is the pure strategy profile  $(B, B)$ . In this last case, when a quasi-strict equilibrium is pure, it will be called *strict*: any deviation from an equilibrium strategy results in a strictly worse payoff.

### 3.3 No regret learning and Regularization

Suppose that a person goes everyday to her work and has two possible routes to follow. How will she decide which one to chose? Suppose that the criterion based on which the choice is made is the time spend i.e., the fastest route is the best one. However, the person is not able to know a priori which route will be fastest each day (let aside applications such as google maps). This problem can be considered as a problem of online learning; the person at each day  $T = 1, 2, \dots$  observes the loss incurred in the previous  $T - 1$  days for each one of the two routes and takes a decision based on this knowledge.

#### 3.3.1 Regret

In this context of online learning, a key requirement is the minimization of the players' regret, i.e., the cumulative payoff difference between each player's chosen action and the best possible action in hindsight over a given horizon of play  $T$ . Formally, given a sequence of play  $X_n \in \mathcal{X}$ ,  $n = 1, 2, \dots$ , the (external) *regret* of player  $i \in \mathcal{N}$  is defined as

$$\text{Reg}_i(T) = \max_{x_i \in \mathcal{X}_i} \sum_{n=1}^T [u_i(x_i; X_{-i,n}) - u_i(X_{i,n}; X_{-i,n})] \quad (3.5)$$

and we will say that player  $i$  has no regret if  $\text{Reg}_i(T) = o(T)$ . This implies that in the long run the player does not regret not to have chosen a fixed action. One may think the sequence  $X_{-i,n}$  in the example mentioned above as the others' people choices of routes.

This definition of regret constitutes the minimum requirement that players would like to satisfy while it takes for granted that there exists a strategy that performs well for the whole window of time. For this reason this type of regret is also known as *static* regret.

Another type of regret, the *dynamic* regret can be also defined as

$$\text{Reg}_i(T) = \sum_{n=1}^T \max_{x_{i,n} \in \mathcal{X}_i} [u_i(x_{i,n}; X_{-i,n}) - u_i(X_{i,n}; X_{-i,n})] \quad (3.6)$$

in which the action chosen as a baseline changes at each round. In the rest of this work we will not focus on this type of regret.

#### 3.3.2 Follow the Regularized Leader (FTRL)

As we have already mentioned, each player needs somehow to choose the strategy that she will play on each round  $T = 1, 2, \dots$ . A simple idea is to chose the strategy that so far has the best cumulative payoff

$$X_{i,T} = \arg \max_{x \in \mathcal{X}_i} \left\{ \sum_{n=0}^{T-1} u_i(x; X_{-i,n}) \right\} \quad (\text{FTL})$$

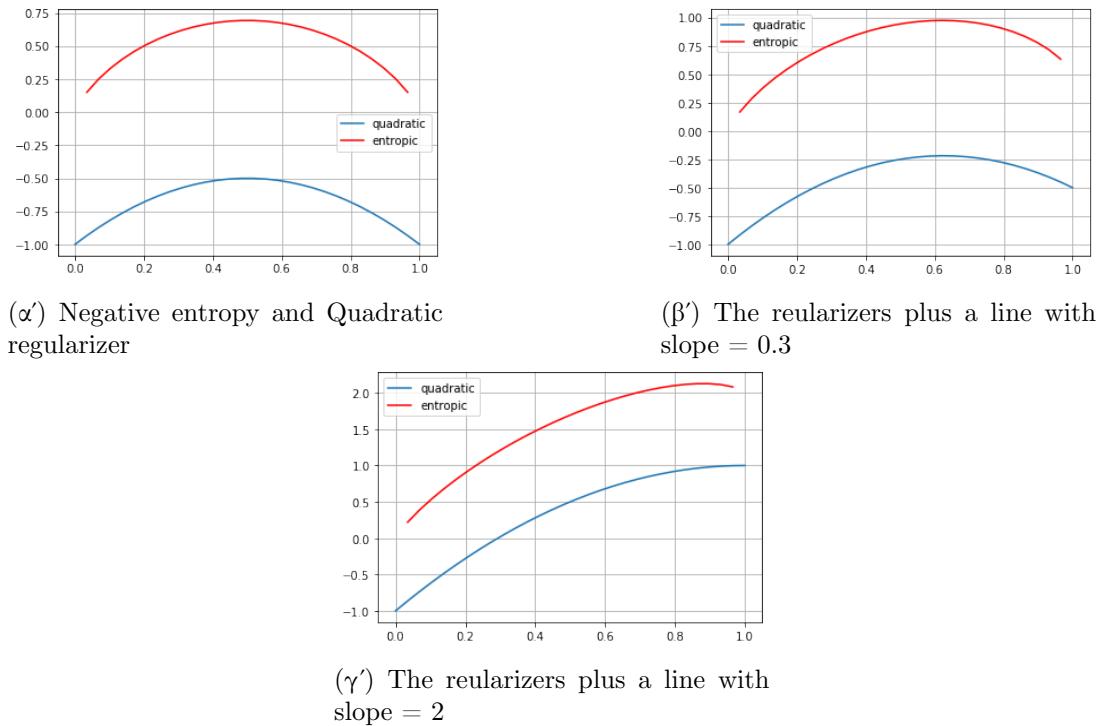


Figure 3.2: Regularizers

This algorithm is known as *Follow the Leader* (FTL). However, this algorithm has a linear worst case regret. It is easy to construct an example to prove this claim.

**Example 3.3.1.** Suppose that the player (learner) has two strategies  $H1, L1$  and that there exists an adversary that also has two strategies  $H2, L2$ . Let the payoff matrix of the learner to be

	H2	L2
H1	0	1
L1	$1 - \epsilon$	0

Suppose now that she starts with the strategy  $H1$  (without loss of generality) and that the adversary chooses  $H2$ . For the next round, she will choose  $L1$  (since  $L1$  has a cumulative payoff of 1 and  $H1$  has a cumulative payoff of 0) while the adversary chooses  $L2$ . Now the learner has the incentive to chose again  $L1$  while the adversary plays  $L2$ . We continue this game with the adversary playing in a way always detrimental to the learner. One can easily verify that indeed in this case  $\text{Reg}(T) = T$ .

This regret is due to the construction of (FTL), which permits to the player to abruptly change her decisions. One solution that ensures the desired regret is to add a regularization penalty, which "smooths out" the transitions between two different states. Intuitively, looking at figure 3.2 one can compare the maximum of a linear function (which is normally presented in the

corners) and the maximum when a regularizer i.e., a strongly convex function is added<sup>3</sup>. Thus this idea gives rise to a new algorithm known as *Follow the Regularized Leader* (FTRL), which can be represented as

$$X_{i,T} = \arg \max_{x \in \mathcal{X}_i} \left\{ \sum_{n=1}^{T-1} u_i(x; X_{-i,n}) - \frac{1}{\eta} h_i(x) \right\} \quad (\text{FTRL})$$

if  $\eta$  is chosen appropriately, no-regret properties of this algorithm can be ensured.

We will present the proof of this statement in the simplest case possible, focusing on one player that has only two strategies. Below we first prove some auxiliary results.

**Lemma 3.3.1** (Closeness of minima). *Consider two strongly convex functions  $f: [0, 1] \rightarrow \mathbb{R}$  and  $g: [0, 1] \rightarrow \mathbb{R}$ , such that  $f''(x) \geq \frac{1}{\eta}$  and  $g''(x) \geq \frac{1}{\eta}$  for all  $x \in [0, 1]$ , and such that  $h(x) = g(x) - f(x)$  is an  $L$ -Lipchitz function, i.e.  $|h(x) - h(x')| \leq L|x - x'|$ . Then, if  $x_f = \arg \min_{x \in [0,1]} f(x)$  and  $x_g = \arg \min_{x \in [0,1]} g(x)$  it holds that:  $|x_f - x_g| \leq \eta L$ .*

*Proof.* First, define the functions

$$\begin{aligned} f_1(x) &= f'(x) - \frac{1}{\eta}x \\ g_1(x) &= g'(x) - \frac{1}{\eta}x \end{aligned}$$

These two functions are apparently increasing. Suppose without loss of generality that  $x_g < x_f$ , then from the *Mean Value theorem* there exists  $x_0 \in (x_g, x_f)$  such that

$$h'(x_0)(x_g - x_f) = h(x_g) - h(x_f) \Rightarrow (g'(x_0) - f'(x_0))(x_f - x_g) = h(x_f) - h(x_g) \quad (3.7)$$

We also have

$$x_g \leq x_0 \leq x_f \Rightarrow f'(x_g) - \frac{1}{\eta}x_g \leq f'(x_0) - \frac{1}{\eta}x_0 \leq -\frac{1}{\eta}x_f \quad (3.8)$$

$$x_g \leq x_0 \leq x_f \Rightarrow -\frac{1}{\eta}x_g \leq g'(x_0) - \frac{1}{\eta}x_0 \leq g'(x_f) - \frac{1}{\eta}x_f \quad (3.9)$$

Of course  $f'(x_f) = 0$  and  $g'(x_g) = 0$  since  $x_f, x_g$  are minimizers of  $f, g$  equivalently. By using (3.8),(3.9) we get

$$\frac{1}{\eta}(x_f - x_g) \leq g'(x_0) - f'(x_0) \leq -\frac{1}{\eta}(x_f - x_g) + g'(x_f) - f'(x_g) \quad (3.10)$$

Combing the above equation with (3.7) and the Lipschitz continuity of  $h$  we have

$$h(x_f) - h(x_g) = (g'(x_0) - f'(x_0))(x_f - x_g) \geq \frac{1}{\eta}(x_f - x_g)^2 \quad (3.11)$$

$$\frac{1}{\eta}(x_f - x_g)^2 \leq L|x_f - x_g| \quad (3.12)$$

$$|x_f - x_g| \leq L\eta \quad (3.13)$$

■

---

<sup>3</sup>The exact assumptions of the regularizers are presented in the next section, but for now think that  $h'' \geq 1$

**Proposition 3.3.1.** *Let  $i \in \mathcal{N}$  be a player that has only two strategies  $H, L$ . Let  $1_n$  be the probability of her first strategy at each round  $n$ . Then under (FTRL) it holds that*

$$|x_{n+1} - x_n| \leq 2\eta \max_{\alpha \in \mathcal{A}} |u_i(\alpha)| \quad (3.14)$$

*Proof.* We will present the steps to reach in the desired result.

- Player  $i$  chooses strategy  $H$  at round  $n = 0, 1, \dots$  with probability  $x_n$  and receives a payoff  $u_{H,n}$ , while with probability  $1 - x_n$  chooses strategy  $L$  and receives a payoff  $u_{L,n}$  at each round  $n = 0, 1, \dots$
- Since the player chooses based on the (FTRL) algorithm it holds that

$$x_{n+1} = \arg \max_{x \in [0,1]} \left\{ x \sum_{k=1}^n u_{H,k} + (1-x) \sum_{k=1}^n u_{L,k} - \frac{1}{\eta} h_i(x) \right\} \quad (3.15)$$

$$= \arg \min_{x \in [0,1]} \left\{ -x \sum_{k=1}^n u_{H,k} - (1-x) \sum_{k=1}^n u_{L,k} + \frac{1}{\eta} h_i(x) \right\} \quad (3.16)$$

- Notice now that the function  $H_n(x) = -x \sum_{k=1}^n u_{H,k} - (1-x) \sum_{k=1}^n u_{L,k} + \frac{1}{\eta} h_i(x)$  has second derivative

$$H''_n(x) = \frac{1}{\eta} h''_i(x) \geq \frac{1}{\eta} \text{ for all } n = 0, 1, \dots \quad (3.17)$$

- Applying Lemma 3.3.1 with  $f = H_n$  and  $g = H_{n-1}$  we have that

$$|x_{n+1} - x_n| \leq \eta L \quad (3.18)$$

where  $L = 2 \max_{\alpha \in \mathcal{A}} |u_i(\alpha)|$ .

Indeed let  $G_n(x) = H_n(x) - H_{n-1}(x) = -x u_{H,n} - (1-x) u_{L,n}$  then

$$|G_n(x) - G_n(x')| = |x(u_{L,n} - u_{H,n}) - x'(u_{L,n} - u_{H,n})| \quad (3.19)$$

$$\leq |u_{L,n} - u_{H,n}| |x - x'| \quad (3.20)$$

$$\leq 2 \max_{\alpha \in \mathcal{A}} |u_i(\alpha)| \quad (3.21)$$

■

Our goal is to prove that (FTRL) is no-regret. For convenience of symbolism, we will also define the following algorithm known as *Be the Regularized Leader* (BTRL). This is an idealized case of (FTRL); suppose that player has access to the induced payoffs for all rounds  $n = 0, 1, \dots, T$  in order to make a decision at round  $T$  then

$$X_{i,T}^* = \arg \max_{x \in \mathcal{X}_i} \left\{ \sum_{n=1}^T u_i(x; X_{-i,n}) - \frac{1}{\eta} h_i(x) \right\} \quad (\text{BTRL})$$

We will now prove that in the simple case, in which the player has only two strategies  $H, L$ , (FTRL) is indeed no-regret.

**Theorem 3.3.2.** *The expected regret of (FTRL) is upper bounded. Specifically,*

$$\text{Reg}(T) \leq \frac{2 \max_{x \in [0,1]} |h(x)|}{\eta} + 2\eta \max_{\alpha \in \mathcal{A}} |u_i(\alpha)| T \quad (3.22)$$

*Proof.* Remember that we focus on the case that player  $i \in \mathcal{N}$  has only two strategies  $H, L$ . For convenience we will adopt the following symbolism. Let

$$f_n(x) = xu_{H,n} + (1-x)u_{L,n} \quad (3.23)$$

$$F_T(x) = \sum_{n=1}^T xu_{H,n} + (1-x)u_{L,n} \quad (3.24)$$

and

$$X_T = \arg \max_{x \in \mathcal{X}_i} \left\{ F_{T-1}(x) - \frac{1}{\eta} h(x) \right\} \quad (3.25)$$

$$\tilde{X}_T = \arg \max_{x \in \mathcal{X}_i} \{ F_{T-1}(x) \} \quad (3.26)$$

$$X_T^* = \arg \max_{x \in \mathcal{X}_i} \left\{ F_T(x) - \frac{1}{\eta} h(x) \right\} \quad (3.27)$$

$$\tilde{X}_T^* = \arg \max_{x \in \mathcal{X}_i} \{ F_T(x) \} \quad (3.28)$$

We first focus on the regret of (BTRL) which we will symbolize as  $\text{Reg}_{\text{BTRL}}(T)$

$$\text{Reg}_{\text{BTRL}}(T) = \max_{x \in [0,1]} \sum_{n=1}^T f_n(x) - \sum_{n=1}^T f_n(X_n^*) \quad (3.29)$$

$$= F_T(\tilde{X}_T^*) - \sum_{n=1}^T (F_n(X_n^*) - F_{n-1}(X_n^*)) \quad (3.30)$$

$$= F_T(\tilde{X}_T^*) - \sum_{n=1}^T (F_n(X_n^*) - \frac{1}{\eta} h(X_n^*) - F_{n-1}(X_n^*) + \frac{1}{\eta} h(X_n^*)) \quad (3.31)$$

$$= F_T(\tilde{X}_T^*) - F_T(X_T^*) + \frac{1}{\eta} h(X_T^*) + F_0(X_1^*) - \frac{1}{\eta} h(X_T^*) \quad (3.32)$$

$$\leq F_T(\tilde{X}_T^*) - F_T(\tilde{X}_T^*) + \frac{1}{\eta} h(\tilde{X}_T^*) - \frac{1}{\eta} h(X_T^*) \quad (3.33)$$

$$\leq \frac{2 \max_{x \in [0,1]} |h(x)|}{\eta} \quad (3.34)$$

We now continue to upper bound the regret of (FTRL). Simply notice that

$$\text{Reg}_{\text{FTRL}}(T) - \text{Reg}_{\text{BTRL}}(T) = \sum_{n=1}^T f_n(X_n^*) - \sum_{n=1}^T f_n(X_n) \quad (3.35)$$

$$\sum_{n=1}^T f_n(X_{n+1}) - \sum_{n=1}^T f_n(X_n) \quad (3.36)$$

Using lemma 3.3.1 and by rearranging we have

$$\text{Reg}_{\text{FTRL}}(T) \leq \text{Reg}_{\text{BTRL}}(T) + 2\eta \max_{\alpha \in \mathcal{A}} |u_i(\alpha)| T \quad (3.37)$$

$$\leq \frac{2 \max_{x \in [0,1]} |h(x)|}{\eta} + 2\eta \max_{\alpha \in \mathcal{A}} |u_i(\alpha)| T \quad (3.38)$$

■

*Remark 1.* By choosing  $\eta$  appropriately ( $\eta = 1/\sqrt{T}$ ), no-regret guarantees are achieved for (FTRL).

*Remark 2.* All these results can be extended for the general case, in which player has  $A > 2$  strategies. The proof follows the steps above but leverages tools from convex analysis presented in appendix A', section A'.2.



# Chapter 4

## Analysis and Results

### 4.1 The algorithm

For the analysis of our results we use an alternative (but equivalent) form of (FTRL). Formally, we have the round-by-round recursive rule

$$\begin{aligned} X_{i,n} &= Q_i(Y_{i,n}) \\ Y_{i,n+1} &= Y_{i,n} + \gamma_n \hat{v}_{i,n} \end{aligned} \tag{FTRL}$$

where  $Q_i: \mathcal{Y}_i \rightarrow \mathcal{X}_i$  denotes the “choice map” of player  $i \in \mathcal{N}$ ,  $\gamma_n > 0$  is a “learning rate” parameter such that  $\sum_n \gamma_n = \infty$ , and  $\hat{v}_{i,n}$  is a “payoff signal” that provides an estimate for the mixed payoffs of player  $i$  at stage  $n$ . We discuss each of these components in detail below.

#### 4.1.1 The feedback model

Depending on the specific framework at play, the modeling details concerning the feedback received by the players may vary wildly. For example, when modeling congestion in a city, it is reasonable to assume that commuters can estimate the time it would have taken them to get to their destination via a different route – e.g., by means of a GPS service or an app like GoogleMaps or Waze. By contrast, in applications of online learning to auctions and online advertising, it is not clear how a player could estimate the payoff of actions they did not play. To account for as broad a range of feedback models as possible, we will take a context-agnostic approach and assume that each player receives a “black-box” model of their payoff vector of the form

$$\hat{v}_n = v(X_n) + \xi_n \tag{4.1}$$

for some abstract error process  $\xi_n = (\xi_{i,n})_{i \in \mathcal{N}}$ . To differentiate between random (zero-mean) and systematic (non-zero-mean) errors, we will further decompose  $\xi_n$  as  $\xi_n = Z_n + b_n$ , where

$$b_n = \mathbb{E}[\xi_n | \mathcal{F}_n] \quad \text{and} \quad \mathbb{E}[Z_n | \mathcal{F}_n] = 0 \tag{4.2}$$

with  $\mathcal{F}_n$  denoting the history of  $X_n$  up to stage  $n$  (inclusive)<sup>1</sup>. We may then characterize the input signal  $\hat{v}_n$  by means of the following statistics:

---

<sup>1</sup>Of course, since the feedback signal is generated only *after* the player chooses a strategy,  $\hat{v}_n$  is not  $\mathcal{F}_n$ -measurable in general.



$$a) \text{ Bias: } \quad \mathbb{E}[\|b_n\|_* | \mathcal{F}_n] \leq B_n \quad (4.3\alpha')$$

$$b) \text{ Variance: } \quad \mathbb{E}[\|Z_n\|_*^2 | \mathcal{F}_n] \leq M_n^2 \quad (4.3\beta')$$

In the above,  $B_n$  and  $M_n$  represent deterministic bounds on the bias and variance of the feedback signal  $\hat{v}_n$ . For concreteness, we will also make the following blanket assumptions:

(A1) *Bias control*:  $\lim_{n \rightarrow \infty} B_n = 0$  and  $\sum_n \gamma_n B_n < \infty$ .

(A2) *Variance control*:  $\sum_n \gamma_n^2 M_n^2 < \infty$ .

(A3) *Generic observation errors at equilibrium*: For every mixed Nash equilibrium  $x^*$  of  $\Gamma$  and for all  $n = 1, 2, \dots$ , there exists a player  $i \in \mathcal{N}$  and strategies  $a, b \in \text{supp}(x_i^*)$  such that

$$\mathbb{P}(|\hat{v}_{ia,n} - \hat{v}_{ib,n}| \geq \beta | \mathcal{F}_n) > 0 \quad \text{for all sufficiently small } \beta > 0. \quad (4.4)$$

The formulation of these hypotheses has been kept intentionally abstract because we have not made any modeling assumptions for how the players' payoff signals are generated. In this regard, they are to be construed as an "inexact model" that allows for a wide variety of settings; as an application, we illustrate below how these assumptions are verified in two widely used learning frameworks.

**Model 1** (Oracle-based feedback). Assume that each player chooses an action based on a given mixed strategy. Then, once this procedure has been completed, an oracle reveals to each player the payoffs corresponding to their pure strategies given the other players' chosen strategies (in the congestion example, this oracle could be Waze or a GPS device). Formally, at each round  $n$ , every player  $i \in \mathcal{N}$  picks an action  $\alpha_{i,n} \in \mathcal{A}_i$  based on  $X_{i,n} \in \mathcal{X}_i$  and observes the pure payoff vector  $v_i(\alpha_n) \equiv (u_i(\alpha_i; \alpha_{-i,n}))_{\alpha_i \in \mathcal{A}_i}$ . Then the player's feedback signal is  $\hat{v}_{i,n} = v_i(\alpha_n)$ , which is a special case of the model (4.1) with  $\xi_n = v(X_n) - v(\alpha_n)$  and  $b_n = 0$ . In more detail, we have:

- (A1) is trivial because  $\mathbb{E}[\hat{v}_n | \mathcal{F}_n] = \mathbb{E}_{X_n}[v(\alpha_n)] = v(X_n)$ , i.e.,  $b_n = 0$ .
- (A2) is satisfied as long as  $\sum_n \gamma_n^2 < \infty$ , since  $\|Z_n\|_* = \|\hat{v}_n - v(X_n)\|_* \leq 2 \max_X \|v(X)\|_*$ .
- (A3) is proved in B'5.

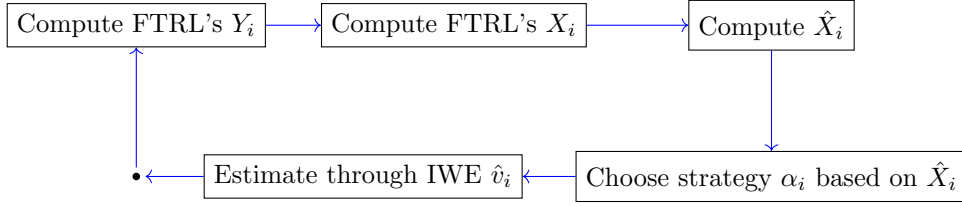
**Model 2** (Payoff-based feedback). Assume that each player picks an action based on some mixed strategy as above; however, players now only observe their realized payoffs  $u_i(\alpha_{i,n}; \alpha_{-i,n})$ . This is the standard model for multi-armed bandits [13, 14], and it is also known as the "bandit feedback" setting. In this case, players can estimate their payoff vectors by means of the importance-weighted estimator:

$$\hat{v}_{i\alpha_i,n} = \frac{\mathbb{1}\{\alpha_{i,n} = \alpha_i\}}{\hat{X}_{i\alpha_i,n}} u_i(\alpha_n) \quad (\text{IWE})$$

where  $\hat{X}_{i,n} = (1 - \varepsilon_n) X_{i,n} + \varepsilon_n / |\mathcal{A}_i|$  is the mixed strategy of the  $i$ -th player at stage  $n$ . Compared to  $X_{i,n}$ , the player's actual sampling strategy is recalibrated by an explicit exploration parameter  $\varepsilon_n \rightarrow 0$  whose role is to stabilize the learning process by controlling the variance of (IWE). The idea is that even if a strategy has zero probability to be chosen under  $X_n$ , it will still be sampled with positive probability thanks to the mixing factor  $\varepsilon_n$ . Schematically players act the following actions:

A standard calculation (that we defer to B'5) shows that (IWE) can be recast in the general form (4.1) with  $B_n = \mathcal{O}(\varepsilon_n)$  and  $M_n^2 = \mathcal{O}(1/\varepsilon_n)$ . We then have:

- (A1) is satisfied as long as  $\varepsilon_n \rightarrow 0$  and  $\sum_n \gamma_n \varepsilon_n < \infty$ .



- (A2) is satisfied as long as  $\sum_n \gamma_n^2 \varepsilon_n^{-1} < \infty$ .
- (A3) is proved in B'5.

*Remark.* The above conditions for the method's learning rate and exploration parameters can be achieved by using schedules of the form  $\gamma_n \propto 1/n^p$  and  $\varepsilon_n \propto 1/n^q$  with  $p + q > 1$  and  $2p - q > 1$ . A popular choice is  $p = 2/3 + \delta$  and  $q = 1/3 + \delta$  for some arbitrarily small  $\delta > 0$  – or  $\delta = 0$  and including an extra logarithmic factor, cf. [15] and references therein.

### 4.1.2 Regularization

The second component of the FTRL method is the players' "choice map"  $Q_i: \mathcal{Y}_i \rightarrow \mathcal{X}_i$ . Because the players' score variables  $Y_{i,n}$  essentially represent an estimate of each strategy's cumulative payoff over time,  $Q_i$  is defined as a "regularized" version of the best-response correspondence  $y_i \mapsto \arg \max_{x_i \in \mathcal{X}_i} \{ \langle y_i, x_i \rangle \}$  (the regularization being necessary to avoid prematurely committing to a strategy). On that account, we will consider *regularized best responses* of the general form

$$Q_i(y_i) = \arg \max_{x_i \in \mathcal{X}_i} \{ \langle y_i, x_i \rangle - h_i(x_i) \}. \quad (4.5)$$

In the above, each player's *regularizer*  $h_i: \mathcal{X}_i \rightarrow \mathbb{R}$  is defined as  $h_i(x_i) = \sum_{\alpha_i \in \mathcal{A}_i} \theta_i(x_i)$  for some "kernel function"  $\theta_i: [0, 1] \rightarrow \mathbb{R}$  with the following properties:

- (i)  $\theta_i$  is *continuous* on  $[0, 1]$ ;
- (ii)  $C^2$ -smooth on  $(0, 1]$ ; and
- (iii)  $\inf_{[0,1]} \theta_i'' > 0$ .

Of course, different regularizers give rise to different instances of (FTRL); for concreteness, we present below two prototypical examples thereof.

**Example 4.1.1** (Multiplicative/Exponential weights update). A popular choice of regularizer is the (negative) entropy  $h_i(x) = \sum_i x_i \log x_i$ , which leads to the *logit choice* map  $\Lambda_i(y) = \exp(y_i) / \sum_j \exp(y_j)$  and the algorithm known as *multiplicative weights update* (MWU), cf. [19, 50, 21, 20, 22].

**Example 4.1.2** (Euclidean projection). Another popular regularizer is the quadratic penalty  $h_i(x) = \sum_i x_i^2 / 2$ , which yields the *payoff projection* choice map  $\boxtimes_i(y) = \arg \min_{x \in \Delta} \|y - x\|^2$ , cf. [23, 25].

To understand the long-run behavior of (FTRL), we will focus on the following overarching question: *Which Nash equilibria hold convergence and stability properties and how are these properties affected by the uncertainty in the players' feedback model?*

We provide the technical groundwork for our answers in 4.2 below; subsequently, we state our results in section 4.3, and present the technical analysis in section 4.4.

## 4.2 Asymptotic Stability

The first thing to note in this general context is that a game may admit several Nash equilibria, both mixed and pure. As a result, global convergence to an equilibrium from all initializations is not possible; for this reason, we will focus on the notion of (*stochastic*) *asymptotic stability* [7, 8, 9]. Heuristically, an equilibrium is *stochastically stable* if any sequence of play that begins close enough to the equilibrium in question, remains close enough with high probability; in addition, if the sequence of play eventually converges to said equilibrium, then we say that it is stochastically asymptotically stable. Formally, we have the following definition.

**Definition 4.2.1.** Let  $x^* \in \mathcal{X}$  be a Nash equilibrium. Fix some arbitrary confidence level  $\delta > 0$  and a neighborhood  $U$  of  $x^*$ . Then  $x^* \in \mathcal{X}$  is said to be

1. **Stochastically stable** if, there exists a neighborhood  $U_0$  of  $x^*$  such that whenever  $X_0 = Q(Y_0) \in U_0$ , we have

$$\mathbb{P}(X_n \in U \text{ for all } n = 0, 1, \dots) \geq 1 - \delta \quad (4.6)$$

whenever  $X_0 = Q(Y_0) \in U_0$ .

2. **Attracting** if there exists a neighborhood  $U_0$  of  $x^*$  such that

$$\mathbb{P}(\lim_{n \rightarrow \infty} X_n = x^*) \geq 1 - \delta \quad (4.7)$$

whenever  $X_0 = Q(Y_0) \in U_0$ .

3. **Stochastically asymptotically stable** if it is stochastically stable and attracting.

Definition 4.2.1 will be the mainstay of our analysis and results, so some remarks are in order.

*Remark 3.* A first intricate detail in the above definition is the high probability requirement: indeed, under uncertainty, a single unlucky estimation of the players' payoff vector could drive  $X_n$  away from any neighborhood of  $x^*$ , possibly never to return. In this regard, local stability results cannot be expected to hold with probability 1, hence the requirement to hold with some arbitrary confidence level in the definition above.

*Remark 4.* Another remark worth making is the requirement  $X_0 = Q(Y_0) \in U_0$  that indicates that some strategies in  $\mathcal{X}$  are not admissible as initial states. Going back to the two archetypal examples of (FTRL), MWU 4.1.1, Projection GD 4.1.2, there is a dichotomy in the properties of the corresponding mirror maps. On the one hand, the kernel of the Euclidean/quadratic regularizer is differentiable on all of  $[0, 1]$ . On the other hand, the derivative of the kernel of the negative Shannon-entropy goes to  $-\infty$  as  $x$  goes to 0. This means that in the latter the boundaries are off the limits and inevitably some initial conditions do not belong in  $\text{im } Q$ . We discuss this dichotomy extensively in section B.1.2.

### 4.3 Main Results

We are now in a position to state our main results. The informal version is as follows.

**Main Theorem.** Suppose that Assumptions (A1)–(A3) hold. Then:  
 $x^*$  is a strict Nash equilibrium  $\iff x^*$  is stochastically asymptotically stable under (FTRL)

Formally, we get the following precise statements and corollaries for the specific feedback models described in section 4.1.1.

**Theorem 4.3.1.** *Let  $x^* \in \mathcal{X}$  be a strict Nash equilibrium of  $\Gamma$ . If (FTRL) is run with inexact payoff feedback satisfying Assumptions (A1) and (A2), then  $x^*$  is stochastically asymptotically stable.*

**Theorem 4.3.2.** *Let  $x^*$  be a mixed Nash equilibrium of  $\Gamma$ . If (FTRL) is run with inexact payoff feedback satisfying assumption (A3), then  $x^*$  is not stochastically asymptotically stable.*

**Corollary 4.3.1.** *Suppose that (FTRL) is run in a generic game with oracle-based feedback as in model 1 and a sufficiently small step-size  $\gamma_n$  with  $\sum_n \gamma_n^2 < \infty$ . Then, a Nash equilibrium is stochastically asymptotically stable if and only if it is strict.*

**Corollary 4.3.2.** *Suppose that (FTRL) is run in a generic game with bandit feedback as in model 2 and sufficiently small step-size and explicit exploration parameters with  $\sum_n \gamma_n^2 / \varepsilon_n < \infty$ ,  $\sum_n \gamma_n \varepsilon_n < \infty$ . Then, a Nash equilibrium is stochastically asymptotically stable if and only if it is strict.*

These results – and, in particular, the implications for the bandit case – provide a learning justification to the abundance of arguments that have been made in the refinement literature against selecting mixed Nash equilibria [17, 6, 24]. In the rest of this work, we present an outline of the main proof ideas and defer the details to the appendix.

## 4.4 Our Techniques

### 4.4.1 The Stochastic Asymptotic Stability of Strict Nash Equilibria

At a high level, the standard tool in FTRL dynamics for questions pertaining to asymptotic stability of strict Nash equilibria is the construction of a potential – or *Lyapunov* – function. However, the analysis and the underlying structural results are considerably more involved when we shift from the continuous dynamics to discrete algorithms and more importantly in a stochastic framework with incomplete feedback information. Still, to build intuition we first recall the continuous and deterministic analogue.

**The continuous-time case.** In prior work [10, 11, 12], multiple instantiations of Bregman functions, like the KL-divergence have been employed as a potent tool for understanding replicator & population dynamics, which are the continuous analogues of MWU/EW (4.1.1). Unfortunately, Bregman functions are insufficient to cover the full spectrum of regularizers

studied in this work. This limitation has been sidestepped in [16] by exploiting the information of the dual space  $\mathcal{Y}$  of the payoff scores, via the Fenchel coupling:

$$F_h(x, y) = h(x) + h^*(y) - \langle y, x \rangle \text{ for all } x \in \mathcal{X}, y \in \mathcal{Y} \quad (4.8)$$

where  $h^* : \mathcal{Y} \rightarrow \mathbb{R}$  is the convex conjugate of  $h$ :  $h^*(y) = \sup_{x \in \mathcal{X}} \{\langle y, x \rangle - h(x)\}$ . Indeed,  $F_h(x^*, y) \geq 0$  where equality holds if and only if  $x^* = Q(y)$  (Proposition B.1.4). Therefore, for the continuous FTRL dynamics  $\dot{y}(t) = v(x(t))$ ,  $x(t) = Q(y(t))$ , it remains to show that the time derivative of the Lyapunov-candidate-function  $L_{x^*}(y(t)) = F_h(x^*, y(t))$  is negative. This last key ingredient for the strict Nash equilibria is derived by their *variational stability* property. Formally, a point  $x^*$  is *variationally stable* if there exists a neighborhood  $U$  of  $x^*$  such that

$$\langle v(x), x - x^* \rangle \leq 0 \text{ for all } x \in U \quad (\text{VS})$$

with equality if and only if  $x = x^*$ . Roughly speaking, this property states that the payoff vectors are pointing “towards” the equilibrium in question since in a neighborhood of  $x^*$ , it strictly dominates over all other strategies. Thus by applying the chain rule, (VS) implies that  $dL_{x^*}(y(t))/dt \leq 0$ <sup>2</sup>. Given their usefulness also in the discrete time stochastic case, we present all the aforementioned properties in detail in the paper’s supplement (B.1-B.2).

**The discrete time.** The core elements of the continuous time proof do not trivially extend to the discrete time case. Even though we are not able to show that  $(F_h(x^*, Y_k))_{k=1}^\infty$  is a decreasing sequence, due to the discretization and the uncertainty involved, we prove that  $F_h(x^*, Y_k) \rightarrow 0$ . This immediately implies that FTRL algorithm converges to  $x^*$ , since from proposition B.1.4  $F_h(x^*, Y_k) \geq \frac{1}{2K_h} \|x^* - X_k\|$ .

To exploit again the Fenchel coupling as a Lyapunov function, successive differences have to be taken among  $F_h(x^*, Y_{n+1}), \dots, F_h(x^*, Y_0)$ . In contrast to the continuous time analysis, since the chain rule no longer applies, we can only do a second order Taylor expansion of the Fenchel coupling. Additionally, let us recall that in our stochastic feedback model, the payoff vector  $\hat{v}_n = v(X_n) + Z_n + b_n$  including possibly either random zero-mean noise or systematic biased noise. Combining proposition B.1.4, definition of  $\hat{v}_n$  and (FTRL), we can create the following upper-bound of Fenchel coupling at each round:

$$F_h(x^*, Y_{n+1}) \leq F_h(x^*, Y_0) + \sum_{k=0}^n \gamma_k (\text{drift}_k + \text{noise}_k + \text{bias}_k) + \frac{1}{2K_h} \sum_{k=0}^n \gamma_k^2 \|\hat{v}_k\|_*^2 \quad (\star)$$

where  $\text{drift}_k = \langle v(X_k), X_k - x^* \rangle$ ,  $\text{noise}_k = \langle Z_k, X_k - x^* \rangle$ ,  $\text{bias}_k = \langle b_k, X_k - x^* \rangle$  are the related terms with the drift of the actual payoff, the zero-mean noise and the bias correspondingly. When  $X_n$  lies in a variationally stable region  $U_{VS}$  of  $x^*$ , the first-order term of  $\text{drift}_k$ , which also appears in the continuous time, corresponds actually to the negative “drift” of the variational stability which attracts Fenchel coupling to zero.

Having settled the basic framework, we split the proof sketch of theorem 4.3.1 into two parts: *stochastic stability & convergence*. Our analysis relies heavily on tools from the convex analysis and martingale limit theory to control the influence of the stochastic terms in the aforementioned bound.

**Step 1: Stability.** Let  $U_\varepsilon = \{x : D_h(x^*, x) < \varepsilon\}$  and  $U_\varepsilon^* = \{y \in \mathcal{Y} : F_h(x^*, y) < \varepsilon\}$  be the  $\varepsilon$ -sublevel sets of Bregman function and Fenchel coupling respectively. Our first observation

<sup>2</sup>Analytically,  $\frac{dL_{x^*}(y(t))}{dt} = \frac{dh^*(y(t))}{dt} - \langle \dot{y}(t), x^* \rangle = \langle \dot{y}(t), \nabla h^*(y) \rangle - \langle \dot{y}(t), x^* \rangle = \langle v(x(t)), x(t) - x^* \rangle \leq 0$ .

is that for all “natural” decomposable regularizers, it holds the so-called “reciprocity condition” (B'.1.1, B'.1.5): essentially, this posits that  $U_\varepsilon$  and  $Q(U_\varepsilon^*)$  are neighborhoods of  $x^*$  in  $\mathcal{X}$ . Additionally, since  $F_h(x^*, y) = D_h(x^*, x)$  whenever  $Q(y) = x$  and  $\text{supp}(x)$  contains  $\text{supp}(x^*)$ , from proposition B'.1.4, it holds that  $Q(U_\varepsilon^*) \subseteq U_\varepsilon$  and  $Q^{-1}(U_\varepsilon) = U_\varepsilon^*$ . Thus, we conclude that whenever  $y \in U_\varepsilon^*$ ,  $x = Q(y) \in U_\varepsilon$ .

To proceed, fix a confidence level  $\delta$  and  $\varepsilon$  sufficiently small such that (VS) holds for all  $x \in U_\varepsilon$ . Using Doob’s maximal inequalities for (sub)martingales (A'.1.6, A'.1.5) we can prove that with probability at least  $1 - \delta$ ,

$$\begin{aligned} (\alpha') & \left\{ \sum_{k=0}^n \gamma_k \text{noise}_k \right\}, \\ (\beta') & \left\{ \sum_{k=0}^n \gamma_k \text{bias}_k \right\} \text{ and} \\ (\gamma') & \left\{ \frac{1}{2K_h} \sum_{k=0}^n \gamma_k^2 \|\hat{v}_k\|_*^2 \right\} \end{aligned}$$

are less than  $\varepsilon/4$  for all  $n \geq 0$ . For concision, we defer the full proof to the supplement of the paper in section B'.4.1. For the rest of this part, we condition on this event and rewrite ( $\star$ ) as  $F_h(x^*, Y_{n+1}) < \sum_{k=0}^n \gamma_k \text{drift}_k + \varepsilon$ .

Following the definition of stability (4.2.1), we prove inductively that if  $X_0$  belongs a smaller neighborhood, namely if  $X_0 \in U_{\varepsilon/4} \cap \text{im } Q$ , then  $X_n$  never escapes  $U_\varepsilon$ ,  $X_n \in U_\varepsilon$  for all  $n \geq 0$ .

- Induction Basis/Hypothesis: Since  $X_0 \in U_{\varepsilon/4} \cap \text{im } Q$ , apparently  $F_h(x^*, Y_0) < \varepsilon/4$  and  $X_0 \in U_\varepsilon$ . Assume that  $X_k \in U_\varepsilon$  for all  $0 \leq k \leq n$ .
- Induction Step: We will prove that  $Y_{n+1} \in U_\varepsilon^*$  and consequently  $X_{n+1} \in U_\varepsilon$ . Since  $U_\varepsilon$  is a neighborhood of  $x^*$  in which (VS) holds we have that  $\text{drift}_k \leq 0$  for all  $0 \leq k \leq n$ . Consequently  $F_h(x^*, Y_{n+1}) < \varepsilon$  which implies that  $Y_{n+1} \in U_\varepsilon^*$  or equivalently  $X_{n+1} \in U_\varepsilon$ .

**Step 2: Convergence.** A tandem combination of stochastic Lyapunov and variational stability is the following lemma:

**Lemma 4.4.1** (Informal statement of Lemma B'.4.1). *Let  $x^* \in \mathcal{A}$  be a strict Nash equilibrium. If  $X_n$  does not exit a neighborhood  $R$  of  $x^*$ , in which variational stability holds, then there exists a subsequence  $X_{n_k}$  of  $X_n$  that converges to  $x^*$  almost surely.*

Indeed, if  $X_n$  is entrapped in a variationally stable region  $U_\varepsilon$  of  $x^*$  without converging to  $x^*$ , we can show that  $\sum_{k=0}^{\infty} \gamma_k \text{drift}_k \rightarrow -\infty$ , while comparatively by the law of the large numbers for martingales (A'.1.3), the contribution of  $(\alpha'), (\beta'), (\gamma')$  is negligible. Thus, in limit ( $\star$ ) implies that  $0 \leq \liminf F_h(x^*, Y_n) \leq -\infty$ , which is a contradiction.

Our final ingredient to complete the proof is that  $(F_h(x^*, Y_k))_{k=1}^{\infty}$  behaves like an almost supermartingale when it is entrapped in a variationally stable region  $U_\varepsilon$  of  $x^*$ . So, by convergence theorem for (sub)-martingales (A'.1.4),  $(F_h(x^*, Y_k))_{k=1}^{\infty}$  actually converges to a random finite variable. Inevitably though,  $\liminf_{n \rightarrow \infty} F_h(x^*, Y_n) = \lim_{n \rightarrow \infty} F_h(x^*, Y_n) = 0$  and by the properties of Fenchel coupling B'.1.4,  $Q(Y_n) = X_n \rightarrow x^*$ .

## 4.4.2 The Stochastic Instability of Mixed Nash Equilibria

For the proof of theorem 4.3.2, it is worth mentioning that in this case stability fails for any choice of step-size. We start by focusing on the assumption of non-degeneracy (A3) of theorem’s statement.

- From a game-theoretic perspective, (A3) actually demands that with non-zero probability, when players receive the payoffs corresponding to pure strategy profiles, there exists at least one player for whom at least two strategies of the equilibrium have distinct payoff signal.

Note that if for each player, the payoffs corresponding to two different strategies of  $\text{supp}(x^*)$  were all equal <sup>3</sup> immediately implies a non-generic game with pure Nash equilibria.

- To illustrate this assumption in our generic feedback model, suppose that this error term  $\xi_n$  is standard normal random noise  $\xi_n$ . Indeed, the requirement of (A3) is satisfied since  $\mathbb{P}(|v_{i,a}(X_n) + \xi_{i,a,n} - v_{i,b}(X_n) - \xi_{i,b,n}| \geq 1/|\mathcal{N}|) > 1 - \mathcal{O}(\exp(-1/|\mathcal{N}|^2))$ . Such kind of property can be derived actually for any per-coordinate independent noise since actually the event of two independent coordinates to be exactly equal has zero measure.

For the bandit models 1, 2 of the previous section, we show that (A3) is satisfied in corollaries B'.3.1, B'.3.2 of B'.5.

Moving on to the proof of theorem 4.3.2, we start our analysis by connecting the difference of the payoff signal between two pure strategies, with the difference of the changes in the output of the regularizers' kernels,  $\theta_i$ :

**Lemma 4.4.2** (Informal Statement of lemma B'.5.1). *Let  $X_{i,n}$  be the sequence of play in (FTRL) i.e.,  $X_{i,n} = Q(Y_{i,n}) \in \mathcal{X}_i$  of player  $i \in \mathcal{N}$ ; and for some round  $n \geq 0$  let  $a, b \in \text{supp}(X_{i,n})$  be two pure strategies of player  $i \in \mathcal{N}$ . Then it holds:*

$$(\theta'_i(X_{ia,n+1}) - \theta'_i(X_{ia,n})) - (\theta'_i(X_{ib,n+1}) - \theta'_i(X_{ib,n})) = \gamma_n(\hat{v}_{ia,n} - \hat{v}_{ib,n})$$

To proceed with the proof of theorem 4.3.2 assume ad absurdum that a mixed Nash equilibrium  $x^*$  is stochastically asymptotically stable. Since  $x^*$  is mixed, there exist  $a, b \in \text{supp}(x^*)$ . Second, the stochastic stability implies that for all  $\varepsilon, \delta > 0$  if  $X_0$  belongs to an initial neighborhood  $U_\varepsilon$ , then  $\|X_n - x^*\| < \varepsilon$  for all  $n \geq 0$ , with probability at least  $1 - \delta$ . Third, by the triangle inequality for two consecutive instances of the sequence of play  $X_{i,n}, X_{i,n+1}$  for any player  $i \in \mathcal{N}$  it holds:

$$|X_{ia,n+1} - X_{ia,n}| + |X_{ib,n+1} - X_{ib,n}| < \mathcal{O}(\varepsilon) \text{ with probability } 1 - \delta \quad (4.9)$$

Consider  $\varepsilon$  sufficiently small, such that the probabilities of the strategies that belong to the support of the equilibrium are bounded away from 0, for all the points of the neighborhood. Since  $\theta_i$  is continuously differentiable in  $(0, 1]$ , the differences described in 4.4.2 are bounded from  $\mathcal{O}(\varepsilon)$  due to (4.9). Thus, if the sequence of play  $X_n$  is contained to an  $\varepsilon$ -neighborhood of  $x^*$ , then the difference of the feedback, for any player  $i \in \mathcal{N}$ , to two strategies of the equilibrium is  $\mathcal{O}(\varepsilon/\gamma_n)$  with probability at least  $1 - \delta$ :

$$\mathbb{P}(|\hat{v}_{ia,n} - \hat{v}_{ib,n}| = \mathcal{O}(\varepsilon/\gamma_n) \mid \mathcal{F}_n) \geq 1 - \delta$$

However, from assumption (A3) for a fixed round  $n$  and some player  $i \in \mathcal{N}$ , there exist  $\beta, \pi > 0$  such that:  $\mathbb{P}(|\hat{v}_{ia,n} - \hat{v}_{ib,n}| \geq \beta \mid \mathcal{F}_n) = \pi > 0$ . Thus by choosing  $\varepsilon = \mathcal{O}(\beta\gamma_n)$  and  $\delta = \pi/2$ , we obtain a contradiction and our proof is complete.

---

<sup>3</sup>when all other players' also employ strategies of the equilibrium

## Chapter 5

# Future work

The equivalence between strict Nash equilibria and stable attracting states of feedback-limited (FTRL) implies that any equilibrium that exhibits payoff-indifference between different strategies is inherently unstable. This fragility has already been remarked from an epistemic viewpoint [17], and our results provide a complementary justification based on realistic models of learning. In the converse direction, the generality of the feedback models considered also provides a template for proving stochastic asymptotic stability results in more demanding learning environments. A particular case of interest arises in online ad auctions where payoffs are observed with delay (or are dropped completely): depending on the delay, the estimation of the player's payoff could exhibit a bias relative to the sampling strategy, and our generic conditions provide an estimate of how large the delays can be before convergence breaks down. This opens the door to an array of fruitful research directions that we intend to pursue in the future.





# Βιβλιογραφία

- [1] J. Nash. Non-Cooperative Games. *Annals of Mathematics*, 54(2), second series, 286-295, 1951.
- [2] Hui-Hsiung Kuo. *Introduction to Stochastic Integration*. Springer-Verlag New York, 2006.
- [3] Bernt Øksendal (2003). *Stochastic Differential Equations*. Springer-Verlag New York. doi:10.1007/978-3-642-14394-6
- [4] Χελιώτης, Δ., 2015. Εισαγωγή στον στοχαστικό λογισμό. [ηλεκτρ. βιβλ.] Αθήνα:Σύνδεσμος Ελληνικών Ακαδημαϊκών Βιβλιοθηκών. Διαθέσιμο στο: <http://hdl.handle.net/11419/4143>
- [5] P. Hall & C.C. Heyde . *Martingale Limit Theory and its Applications*. Academic Press, 1980.
- [6] Drew Fudenberg & Jean Tirole. *Game Theory*. The MIT Press, 1991.
- [7] Josef Hofbauer and Karl Sigmund. *Evolutionary Games and Population Dynamics*. Cambridge University Press, Cambridge, UK, 1998.
- [8] William H. Sandholm. *Population Games and Evolutionary Dynamics*. MIT Press, Cambridge, MA, 2010.
- [9] Rafail Z. Khasminskii. *Stochastic Stability of Differential Equations*. Number 66 in *Stochastic Modelling and Applied Probability*. Springer-Verlag, Berlin, 2 edition, 2012.
- [10] Edwin Hewitt and Karl Stromberg. *Real and abstract analysis*. Graduate Texts in Mathematics. Springer-Verlag, New York, NY, 1975.
- [11] Steven Weinberg. *Quantum Field Theory*. Cambridge University Press, Cambridge, UK, 2000.
- [12] Marc Harper. Escort evolutionary game theory. *Physica D: Nonlinear Phenomena*, 240(18):1411–1415, September 2011.
- [13] Nicolò Cesa-Bianchi and Gábor Lugosi. *Prediction, Learning, and Games*. Cambridge University Press, 2006.
- [14] Sébastien Bubeck and Nicolò Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends in Machine Learning*, 5(1):1–122, 2012.
- [15] Aleksandrs Slivkins. Introduction to multi-armed bandits. *Foundations and Trends in Machine Learning*, 12(1-2): 1–286, November 2019.

- [16] Panayotis Mertikopoulos and William H. Sandholm. Learning in games via reinforcement and regularization. *Mathematics of Operations Research*, 41(4):1297–1324, November 2016.
- [17] Eric van Damme. *Stability and perfection of Nash equilibria*. Springer-Verlag, Berlin, 1987.
- [18] Ralph Tyrrell Rockafellar and Roger J. B. Wets. *Variational Analysis*, volume 317 of *A Series of Comprehensive Studies in Mathematics*. Springer-Verlag, Berlin, 1998.
- [19] Vladimir G. Vovk. Aggregating strategies. In *COLT '90: Proceedings of the 3rd Workshop on Computational Learning Theory*, pages 371–383, 1990.
- [20] Sanjeev Arora, Elad Hazan, and Satyen Kale. The multiplicative weights update method: A meta-algorithm and applications. *Theory of Computing*, 8(1):121–164, 2012.
- [21] Peter Auer, Nicolò Cesa-Bianchi, Yoav Freund, and Robert E. Schapire. Gambling in a rigged casino: The adversarial multi-armed bandit problem. In *Proceedings of the 36th Annual Symposium on Foundations of Computer Science*, 1995.
- [22] Shai Shalev-Shwartz. Online learning and online convex optimization. *Foundations and Trends in Machine Learning*, 4(2):107–194, 2011.
- [23] Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *ICML '03: Proceedings of the 20th International Conference on Machine Learning*, pages 928–936, 2003.
- [24] Eddie Dekel and Drew Fudenberg. Rational behavior with payoff uncertainty. *Journal of Economic Theory*, 52: 243–267, 1990.
- [25] Ratul Lahkar and William H. Sandholm. The projection dynamic and the geometry of population games. *Games and Economic Behavior*, 64:565–590, 2008.
- [26] Shai Shalev-Shwartz and Yoram Singer. Convex repeated games and Fenchel duality. In *NIPS' 06: Proceedings of the 19th Annual Conference on Neural Information Processing Systems*, pages 1265–1272. MIT Press, 2006.
- [27] James Hannan. Approximation to Bayes risk in repeated play. In Melvin Dresher, Albert William Tucker, and P. Wolfe, editors, *Contributions to the Theory of Games, Volume III*, volume 39 of *Annals of Mathematics Studies*, pages 97–139. Princeton University Press, Princeton, NJ, 1957.
- [28] Sergiu Hart and Andreu Mas-Colell. A simple adaptive procedure leading to correlated equilibrium. *Econometrica*, 68(5):1127–1150, September 2000.
- [29] Yannick Viossat and Andriy Zapechelnjuk. No-regret dynamics and fictitious play. *Journal of Economic Theory*, 148(2):825–842, March 2013.
- [30] Gerasimos Palaiopoulos, Ioannis Panageas, and Georgios Piliouras. Multiplicative weights update with constant stepsize in congestion games: Convergence, limit cycles and chaos. In *NIPS '17: Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017.
- [31] Yun Kuen Cheung and Georgios Piliouras. Vortices instead of equilibria in minmax optimization: Chaos and butterfly effects of online learning in zero-sum games. In *COLT '19: Proceedings of the 32nd Annual Conference on Learning Theory*, 2019.
- [32] Barnabé Monnot and Georgios Piliouras. Limits and limitations of no-regret learning in games. *The Knowledge Engineering Review*, 32, 2017.

- [33] Josef Hofbauer, Sylvain Sorin, and Yannick Viossat. Time average replicator and best reply dynamics. *Mathematics of Operations Research*, 34(2):263–269, May 2009.
- [34] Panayotis Mertikopoulos, Christos H. Papadimitriou, and Georgios Piliouras. Cycles in adversarial regularized learning. In *SODA '18: Proceedings of the 29th annual ACM-SIAM Symposium on Discrete Algorithms*, 2018.
- [35] Shai Shalev-Shwartz and Yoram Singer. Convex repeated games and Fenchel duality. In *NIPS' 06: Proceedings of the 19th Annual Conference on Neural Information Processing Systems*, pages 1265–1272. MIT Press, 2006.
- [36] Jörgen W. Weibull. *Evolutionary Game Theory*. MIT Press, Cambridge, MA, 1995.
- [37] Josef Hofbauer and Karl Sigmund. Evolutionary game dynamics. *Bulletin of the American Mathematical Society*, 40(4):479–519, July 2003.
- [38] Aldo Rustichini. Optimal properties of stimulus-response learning models. *Games and Economic Behavior*, 29 (1-2):244–273, 1999.
- [39] Johanne Cohen, Amélie Héliou, and Panayotis Mertikopoulos. Learning with bandit feedback in potential games. In *NIPS '17: Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017.
- [40] Vasilis Syrgkanis, Alekh Agarwal, Haipeng Luo, and Robert E. Schapire. Fast convergence of regularized learning in games. In *NIPS '15: Proceedings of the 29th International Conference on Neural Information Processing Systems*, pages 2989–2997, 2015.
- [41] Robert David Kleinberg, Georgios Piliouras, and Éva Tardos. Load balancing without regret in the bulletin board model. *Distributed Computing*, 24(1):21–29, 2011.
- [42] Sergiu Hart and Andreu Mas-Colell. Uncoupled dynamics do not lead to Nash equilibrium. *American Economic Review*, 93(5):1830–1836, 2003.
- [43] David S. Leslie and E. J. Collins. Individual Q-learning in normal form games. *SIAM Journal on Control and Optimization*, 44(2):495–514, 2005.
- [44] David S. Leslie. Generalised weakened fictitious play. *Games and Economic Behavior*, 56(2):285–298, August 2006.
- [45] Roberto Cominetti, Emerson Melo, and Sylvain Sorin. A payoff-based learning procedure and its application to traffic games. *Games and Economic Behavior*, 70(1):71–83, 2010.
- [46] Pierre Coucheney, Bruno Gaujal, and Panayotis Mertikopoulos. Penalty-regulated dynamics and robust learning procedures in games. *Mathematics of Operations Research*, 40(3):611–633, August 2015.
- [47] Sylvain Sorin. Exponential weight algorithm in continuous time. *Mathematical Programming*, 116(1):513–528, 2009.
- [48] Lampros Flokas, Emmanouil Vasileios Vlatakis-Gkaragkounis, Thanasis Lianas, Panayotis Mertikopoulos, and Georgios Piliouras. No-regret learning and mixed Nash equilibria: They do not mix. In *NeurIPS '20: Proceedings of the 34th International Conference on Neural Information Processing Systems*, 2020.
- [49] Panayotis Mertikopoulos and Zhengyuan Zhou. Learning in games with continuous action sets and unknown payoff functions. *Mathematical Programming*, 173(1-2):465–507, January 2019.

- [50] Nick Littlestone and Manfred K. Warmuth. The weighted majority algorithm. *Information and Computation*, 108 (2):212–261, 1994.
- [51] Xingyu Zhou. On the Fenchel Duality between Strong Convexity and Lipschitz Continuous Gradient. March, 2018.

# Appendix A'

## Theoretical Basis

### A'.1 Elements of martingale limit theory

#### A'.1.1 Basic definitions

In this part we provide some basic definitions necessary for the rest of this thesis.

**Definition A'.1.1.** Let  $\Omega$  be a given set, then a  $\sigma$ -algebra  $\mathcal{F}$  on  $\Omega$  is a family  $\mathcal{F}$  of subsets of  $\Omega$  with the following properties

1.  $\emptyset \in \mathcal{F}$
2.  $F \in \mathcal{F} \Rightarrow F^c \in \mathcal{F}$ , where  $F^c = \Omega \setminus F$  is the complement of  $F$  in  $\Omega$
3.  $A_1, A_2, \dots \in \mathcal{F} \Rightarrow A := \bigcup_{i=1}^{\infty} A_i \in \mathcal{F}$

The pair  $(\Omega, \mathcal{F})$  is called a measurable space. A probability measure  $P$  on a measurable space  $(\Omega, \mathcal{F})$  is a function  $P : \mathcal{F} \rightarrow [0, 1]$  such that

1.  $P(\emptyset) = 0, P(\Omega) = 1$
2. If  $A_1, A_2, \dots \in \mathcal{F}$  and  $\{A_i\}_{i=1}^{\infty}$  is disjoint (i.e.,  $A_i \cap A_j = \emptyset$  for  $i \neq j$ ) then

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i) \quad (\text{A'.1})$$

The triple  $(\Omega, \mathcal{F}, P)$  is called a probability space.

The subsets  $F$  of  $\Omega$  which belong to  $\mathcal{F}$  are called  $\mathcal{F}$ -measurable sets. In a probability context these sets are called *events* and  $P(F)$  is the probability that the event  $F$  occurs. Given a set  $\mathcal{V}$  which contains some subsets of  $\Omega$ , there is a smallest  $\sigma$ -algebra  $\mathcal{H}_{\mathcal{V}}$  containing  $\mathcal{V}$ :

$$\mathcal{H}_{\mathcal{V}} = \bigcap \{ \mathcal{H}; \mathcal{H} \text{ } \sigma\text{-algebra of } \Omega, \mathcal{V} \subset \mathcal{H} \} \quad (\text{A'.2})$$

We call  $\mathcal{H}_{\mathcal{V}}$  the  $\sigma$ -algebra generated by  $\mathcal{V}$

A special case of the above definition emerges if we consider  $\mathcal{V}$  the set, containing all the open sets of  $\Omega = \mathbb{R}^n$ . The  $\sigma$ -algebra generated by  $\mathcal{V}$  is called *Borel*  $\sigma$  algebra on  $\mathbb{R}^n$ . Consider now a probability space  $(\Omega, \mathcal{F}, P)$ ; a random variable  $X$  is an  $\mathcal{F}$ -measurable function  $X : \Omega \rightarrow \mathbb{R}^n$ . Every random variable induces a probability measure  $\mu_X$  on  $\mathbb{R}^n$ , defined by

$$\mu_X(B) = P(X^{-1}(B)) \quad (\text{A'.3})$$

We call  $\mu_X$  the distribution of  $X$ . The number

$$\mathbb{E}[X] := \int_{\Omega} X(\omega) dP(\omega) = \int_{\mathbb{R}^n} x d\mu_x(x) \quad (\text{A'.4})$$

is called the expectation of  $X$ , if  $\int_{\Omega} X(\omega) dP(\omega) < \infty$ .

Equivalently, if  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is Borel measurable and  $\int_{\Omega} |f(X(\omega))| dP(\omega) < \infty$  then the expectation of the random variable  $f(X)$  is

$$\mathbb{E}[f(X)] := \int_{\Omega} f(X(\omega)) dP(\omega) \quad (\text{A'.5})$$

Notice that in the finite case, in which  $X$  is a random variable and  $x_1, \dots, x_n$  are the possible outcomes of  $X$ , occurring with probabilities  $p_1, \dots, p_n$  the expectation of  $X$  is

$$\mathbb{E}[X] = \sum_{i=1}^n x_i p_i \quad (\text{A'.6})$$

**Definition A'.1.2.** Let  $(\Omega, \mathcal{F}, P)$  be a probability space, then a stochastic process is a collection of random variables

$$\{X_n\}_{n \in T} \quad (\text{A'.7})$$

for some set  $T \subseteq [0, \infty)$ .

### A'.1.2 Conditional Expectation

Let  $(\Omega, \mathcal{F}, P)$  be a probability space and let  $X : \Omega \rightarrow \mathbb{R}^n$  be a random variable with finite expectation i.e.,  $\mathbb{E}[|X|] < \infty$ . If  $\mathcal{H} \subset \mathcal{F}$  is a  $\sigma$ -algebra then the conditional expectation of  $X$  given  $\mathcal{H}$ , which is denoted by  $\mathbb{E}[X | \mathcal{H}]$  is:

**Definition A'.1.3.**  $\mathbb{E}[X | \mathcal{H}]$  is the almost surely (a.s.) unique function from  $\Omega$  to  $\mathbb{R}^n$  satisfying:

- $\mathbb{E}[X | \mathcal{H}]$  is  $\mathcal{H}$ -measurable
- $\int_{H \in \mathcal{H}} \mathbb{E}[X | \mathcal{H}] dP = \int_{H \in \mathcal{H}} X dP$

The existence and uniqueness of this function can be proven using Radon-Nikodym theorem. A proof can be found in [2],[3],[4]. Below we present some basic properties of the conditional expectation:

**Theorem A'.1.1.** Let  $X : \Omega \rightarrow \mathbb{R}^n$  and  $Y : \Omega \rightarrow \mathbb{R}^n$  be two random variables with finite expectations and  $a, b \in \mathbb{R}$ . Then

1.  $\mathbb{E}[aX + bY | \mathcal{H}] = a \mathbb{E}[X | \mathcal{H}] + b \mathbb{E}[Y | \mathcal{H}]$
2.  $\mathbb{E}[\mathbb{E}[X | \mathcal{H}]] = \mathbb{E}[X]$
3. If  $X$  is  $\mathcal{H}$ -measurable then  $\mathbb{E}[X | \mathcal{H}] = X$
4. If  $X$  is independent of  $\mathcal{H}$  then  $\mathbb{E}[X | \mathcal{H}] = \mathbb{E}[X]$
5. If  $Y$  is  $\mathcal{H}$ -measurable, then  $\mathbb{E}[\langle X, Y \rangle | \mathcal{H}] = \langle Y, \mathbb{E}[X | \mathcal{H}] \rangle$ .

**Theorem A'.1.2.** Let  $\mathcal{G}_1, \mathcal{G}_2$  be two  $\sigma$ -algebras such that  $\mathcal{G}_1 \subset \mathcal{G}_2 \subset \mathcal{F}$ . Then

1.  $\mathbb{E}[\mathbb{E}[X | \mathcal{G}_2] | \mathcal{G}_1] = \mathbb{E}[X | \mathcal{G}_1]$
2.  $\mathbb{E}[\mathbb{E}[X | \mathcal{G}_1] | \mathcal{G}_2] = \mathbb{E}[X | \mathcal{G}_1]$

### A'.1.3 Martingales

Let  $(\Omega, \mathcal{F}, P)$  be a probability space. We call *filtration* in this space an increasing sequence  $(\mathcal{F}_n)_{n \geq 0}$  of  $\sigma$ -algebras which are all subsets of  $\mathcal{F}$  i.e.,  $\mathcal{F}_n \subset \mathcal{F}_{n+1} \subset \mathcal{F}$  for all  $n \geq 0$ . A sequence of random variables  $(X_n)_{n \geq 0}$  is attached to the filtration  $(\mathcal{F}_n)_{n=0}^\infty$ , if for all  $n \geq 0$   $X_n$  is  $\mathcal{F}_n$ -measurable.

**Definition A'.1.4.** A sequence of random variables  $X = (X_n)_{n \geq 0}$  that satisfies the following properties

1.  $(X_n)_{n \geq 0}$  is attached to  $(\mathcal{F}_n)_{n \geq 0}$
2.  $\mathbb{E}[X_n] < \infty$  for all  $n \geq 0$
3.  $\mathbb{E}[X_{n+1} | \mathcal{F}_n] = X_n$  for all  $n \geq 0$

is called a *martingale* with respect to the filtration  $(\mathcal{F}_n)_{n \geq 0}$ . If *iii*) holds as an inequality then

1.  $X_n$  is called a *submartingale* if  $\mathbb{E}[X_{n+1} | \mathcal{F}_n] \geq X_n$  for all  $n \geq 0$
2.  $X_n$  is called a *supermartingale* if  $\mathbb{E}[X_{n+1} | \mathcal{F}_n] \leq X_n$  for all  $n \geq 0$

Actually the filtration  $\mathcal{F}_n$  includes all the information up to round  $n$ . As though, if  $X_n$  is  $\mathcal{F}_n$ -measurable, it holds that  $\mathbb{E}[X_n | \mathcal{F}_n] = X_n$ .

### A'.1.4 Martingale limit theorems

Below we first present a simple fact for the reader to keep in mind, followed by the main theorems that we utilize in the main body of our proofs presented in the next chapters.

**Fact 1.** Let  $R_n = \sum_{k=1}^n r_k$ , where  $r_k$  is a positive random variable for all  $k = 0, 1, \dots$  attached to the filtration  $\mathcal{F}_{k-1}$ . Then  $R_n$  is a submartingale.

We begin with the strong law of large numbers for martingale difference sequences:

**Theorem A'.1.3.** Let  $R_n = \sum_{k=1}^n r_k$  be a martingale with respect to an underlying stochastic basis  $(\Omega, \mathcal{F}, (\mathcal{F}_n)_{n=1}^\infty, \mathbb{P})$  and let  $(\tau_n)_{n=1}^\infty$  be a nondecreasing sequence of positive numbers with  $\lim_{n \rightarrow \infty} \tau_n = \infty$ . If  $\sum_{n=1}^\infty \tau_n^{-p} \mathbb{E}[|r_n|^p | \mathcal{F}_{n-1}] < \infty$  for some  $p \in [1, 2]$  almost surely, then

$$\lim_{n \rightarrow \infty} \tau_n^{-1} R_n = 0 \text{ almost surely} \quad (\text{A'.8})$$

The second important result for our analysis is Doob's martingale convergence theorem:

**Theorem A'.1.4.** If  $R_n$  is a submartingale that is bounded in  $L_1$  (i.e.,  $\sup_n \mathbb{E}[|R_n|] < \infty$ ),  $R_n$  converges almost surely to a random variable  $R$  with  $\mathbb{E}[R] < \infty$ .

Finally, we use the known as Doob's maximal inequality and one of its variants, presented below:

**Theorem A'.1.5.** Let  $R_n$  be a non-negative submartingale and fix some  $\varepsilon > 0$ . Then:

$$\mathbb{P}(\sup_n R_n \geq \varepsilon) \leq \frac{\mathbb{E}[R_n]}{\varepsilon} \quad (\text{A'.9})$$

**Theorem A'.1.6.** Let  $R_n$  be a martingale and fix some  $\varepsilon > 0$ . Then:

$$\mathbb{P}(\sup_n |R_n| \geq \varepsilon) \leq \frac{\mathbb{E}[R_n^2]}{\varepsilon^2} \quad (\text{A'.10})$$

Proofs of all these results can be found in [5].



## A'.2 Elements of Convex Analysis

In this section we provide basic definitions and results from convex analysis. Many of these are implicitly used in our proofs

### A'.2.1 Basic definitions

**Definition A'.2.1** (Convex Set). A subset  $C$  of  $\mathbb{R}^d$  is said to be **convex** when for every pair  $x, y \in C \subseteq \mathbb{R}^d$  and every  $\lambda \in \mathbb{R}$  for which  $0 \leq \lambda \leq 1$  the following holds:

$$z = (1 - \lambda)x + \lambda y \in C$$

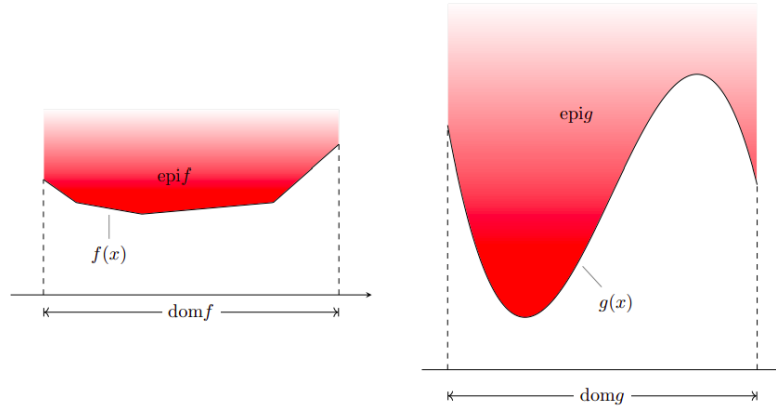


Figure A'.1: Depiction of the epigraph of two functions

Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be a function. We can imagine  $f$  as defining a hyper-surface in the joint space of its input space and its output space,  $\mathbb{R}^n \times \mathbb{R}$ . The points above that surface whose perpendicular projections on  $\mathbb{R}^n$  remain in  $\text{dom } f$  form the epigraph of the given function. More formally:

**Definition A'.2.2** (Epigraph). An **epigraph** of a function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is said to be the set of points  $(x, \mu)$  such that  $\mu \geq f(x)$  and it is noted as:

$$\text{epi } f = \{(x, \mu) \mid \mu \geq f(x)\}$$

An illustration of the epigraph can be viewed in figure A'.1.

**Definition A'.2.3** (Proper function). A function  $f$  is called proper if its epigraph is non-empty and contains no vertical lines.

**Definition A'.2.4** (Lipschitz continuity). Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  be a vector-valued function over some open set  $\mathcal{X} \subset \mathbb{R}^n$ ; we say that  $f$  is  **$(a, b)$ -Lipschitz continuous** if there exists a constant  $L$  for norms  $\|\cdot\|_a, \|\cdot\|_b$  such that for all  $x, y \in \mathcal{X}$ :

$$\|f(x) - f(y)\|_b \leq L \cdot \|x - y\|_a$$

The Lipschitz constant,  $L^{(a,b)}(f, \mathcal{X})$ , is the infimum over all such all such  $L$ . Equivalently, one can define  $L^{(a,b)}(f, \mathcal{X})$  as

$$L^{(a,b)}(f, \mathcal{X}) = \sup_{x, y \in \mathcal{X}, x \neq y} \frac{\|f(x) - f(y)\|_b}{\|x - y\|_a}$$

**Definition A'.2.5** (Differentiability). We say that  $f$  is **differentiable** at  $x$  if there exists some linear operator  $\nabla f(x) \in \mathbb{R}^{n \times m}$  such that

$$\lim_{h \rightarrow 0} \frac{\|f(x+h) - f(x) - \nabla f(x)^T h\|}{\|h\|} = 0$$

A linear operator such that the above equation holds is defined as the Jacobian.

**Definition A'.2.6** (Directional derivative). Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  be a vector-valued function, then the **directional derivative** of  $f$  along a direction  $v \in \mathbb{R}^n$  is defined as

$$\delta_v f(x) := \lim_{t \rightarrow 0} \frac{f(x+tv) - f(x)}{t}$$

We now add the following facts:

**Fact 2.** Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  be a vector-valued function. Then the following hold:

1. If  $f$  is Lipschitz continuous, then it is absolutely continuous.
2. If  $f$  is differentiable at a point  $x \in \mathbb{R}^n$ , all directional derivatives exist at  $x$ . The converse is not true, however.
3. If  $f$  is differentiable at a point  $x \in \mathbb{R}^n$ , then for any vector  $v \in \mathbb{R}^n$ ,  $\delta_v f(x) = \nabla f(x)^T v$ .
4. (**Rademacher's Theorem**): If  $f$  is Lipschitz continuous, then  $f$  is everywhere differentiable except for a set of measure zero (under the standard Lebesgue measure in  $\mathbb{R}^n$ ).

**Definition A'.2.7** (Subgradient). Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be a function, then a vector  $s \in \mathbb{R}^n$  is a subgradient of  $f$  at  $x \in \text{dom} f$  if for all  $y \in \text{dom} f$  it holds

$$f(y) \geq f(x) + s^T(y-x) \tag{A'.11}$$

*Remark 5.* If  $f$  is convex and differentiable, then its gradient at a point  $x$  is also a subgradient. But a subgradient can exist even when  $f$  is not differentiable at  $x$ .

**Definition A'.2.8** (Subdifferential). The subdifferential of a function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  at a point  $x \in \mathbb{R}^n$ , denoted by  $\partial f(x)$ , is the set of subgradients of  $f$  at that point  $x$ . A function  $f$  is called subdifferential at a point  $x$  if there exists at least one subgradient of  $f$  at  $x$ . If  $f$  is subdifferential at all  $x \in \text{dom} f$  then  $f$  is called subdifferential.

**Definition A'.2.9** (Lipschitz Continuous Gradient). A function  $f$  is said to have a  $L$ -Lipschitz continuous gradient if there exists  $L > 0$  such that for all  $x, y$  in its domain, it holds:

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$$

## A'.2.2 Convexity & Duality

In this section we are going to discuss the conjugate transform of functions. It is a transform that maps the parameters of hyper-planes tangent to the curve of a function to a certain value. It may not be the first time one sees such a transform, one that shifts our attention to a parameter space. For example, in traditional computer vision a certain transform, known as Hough Transform, is used in order to map whole lines of the 2-D space to tuples  $(\rho, \theta)$ ;  $\theta$  being

the angle that the line perpendicular to the line in question forms with the horizontal axis and  $\rho$  being the distance of the line from the origin. Although this only vaguely resembles the subject of our discussion – and we regret causing any confusion – we mention it in order to motivate more ways of thinking of lines than just as a set of points. There are various implementations based on this premise that help us detect and recognize not only lines but also regular geometric shapes. Our subject revolves around tangent lines (or hyper-planes for

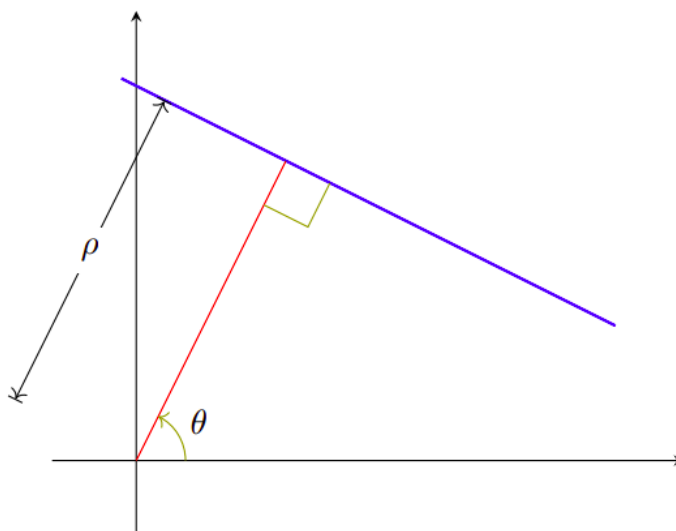


Figure A'.2: Representing the blue line with parameters  $\rho, \theta$

function domains with dimension greater than 1) on a convex function. We will demonstrate a way that has been devised in order to represent elegantly the whole set of these tangent lines.

Frankly the definition seems a bit awkward. Considering its geometric interpretation could maybe shed some light as to what this is supposed to mean.

The conjugate transform of a function  $f$  is merely a function  $f^*$  that maps slopes  $\alpha$  to the maximum available offset  $\beta$  such that the given line  $\alpha x + \beta$  will be tangent to the curve defined by  $f$ .

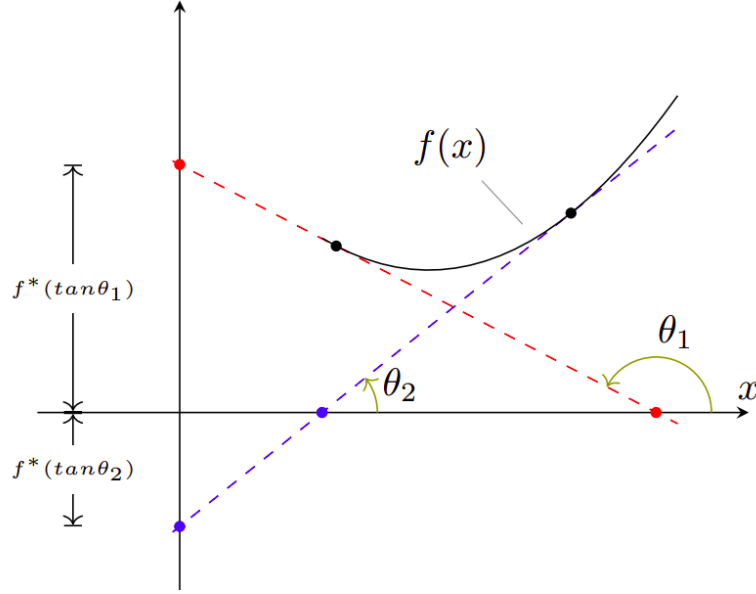


Figure A'.3: Geometric meaning of the conjugate transform

**Definition A'.2.10** (Dual space). Given any vector space  $V$  over a field  $F$ , the (algebraic) dual space  $V^*$  is the set of all linear maps  $\phi : V \rightarrow F$ .

**Definition A'.2.11** (Fenchel conjugate). Let  $X$  be a real topological vector space and  $X^*$  its dual space. Then for a function  $f : X \rightarrow \mathbb{R}$ , the convex conjugate  $f^* : X^* \rightarrow \mathbb{R}$  is defined as

$$f^*(x^*) := \sup_{x \in X} \{\langle x^*, x \rangle - f(x)\} \quad (\text{A'.12})$$

**Theorem A'.2.1** (Fenchel's inequality). For any subgradient vector  $p \in f^*(\text{dom } f^*)$  and any  $x \in \text{dom } f$  the following inequality stands:

$$f^*(p) + f(x) \geq \langle p, x \rangle$$

*Proof.* By the definition of the conjugate transform:  $f^*(p) = \sup_{x \in \text{dom } f} \{\langle p, x \rangle - f(x)\}$  we can decide that:

$$f^*(p) \geq \langle p, x \rangle - f(x), \text{ for all } x \quad (\text{A'.13})$$

$$f^*(p) + f(x) \geq \langle p, x \rangle \quad (\text{A'.14})$$

■

### A'.2.3 Convexity and Smoothness

**Definition A'.2.12** (Convexity in  $\mathbb{R}^n$ ). A function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is convex, if its domain  $A$  is a convex set and for all  $x, y \in A$  and for all  $\lambda \in [0, 1]$  it holds:

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$$

**Definition A'.2.13** (Effective Domain of a Convex Function). The effective domain of a convex function,  $\text{dom } f$ , is the set of  $x$  such that:

$$\text{dom } f = \{x \mid f(x) < -\infty\}$$

**Lemma A'.2.2** (Equivalence for Convexity). Suppose  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  with the extended value extension. Then, the following statements are equivalent:

[1] (Jensen's inequality):  $f$  is convex.

[2] (First order):  $f(y) \geq f(x) + s_x^T(y - x)$  for all  $x, y$  and any  $s_x \in \partial f(x)$ .

[3] (Monotonicity of subgradient):  $(s_y - s_x)^T(y - x) \geq 0$  for all  $x, y$  and any  $s_x \in \partial f(x), s_y \in \partial f(y)$ .

*Proof.* [1]  $\Rightarrow$  [2] By definition of convexity we have

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) \quad (\text{A'.15})$$

$$\frac{f(\lambda x + (1 - \lambda)y) - f(y)}{\lambda} \leq f(x) - f(y) \quad (\text{A'.16})$$

But, if  $s_x \in \partial f(x)$  then

$$f(y) \geq f(x) + s_x^T(y - x) \text{ for all } y \in \text{dom } f \quad (\text{A'.17})$$

$$f(\lambda x + (1 - \lambda)y) \geq f(y) + s_y^T(\lambda x + (1 - \lambda)y - y) \quad (\text{A'.18})$$

$$f(\lambda x + (1 - \lambda)y) - f(y) \geq \lambda s_y^T(x - y) \quad (\text{A'.19})$$

where for the first we have substitute  $y = \lambda x + (1 - \lambda)y, x = y$ . Combining the above two results we get

$$f(x) \geq f(y) + s_y^T(x - y) \text{ or} \quad (\text{A'.20})$$

$$f(y) \geq f(x) + s_x^T(y - x) \quad (\text{A'.21})$$

[2]  $\Rightarrow$  [1]

$$f(y) \geq f(x) + s_x^T(y - x) \text{ for all } x, y \quad (\text{A'.22})$$

By substituting  $x = x', x' = \lambda x + (1 - \lambda)y$  in the above inequality we get

$$f(y) \geq f(x') + s_{x'}^T(y - x') \quad (\text{A'.23})$$

$$(1 - \lambda)f(y) \geq (1 - \lambda)f(x') + (1 - \lambda)\lambda s_{x'}^T(y - x) \quad (\text{A'.24})$$

Furthermore, let  $x' = \lambda x + (1 - \lambda)y$  by renaming we have

$$f(x) \geq f(y) + s_y^T(x - y) \quad (\text{A'.25})$$

$$f(x) \geq f(x') + s_{x'}^T(x - x') \quad (\text{A'.26})$$

$$\lambda f(x) \geq \lambda f(x') + \lambda(1 - \lambda)s_{x'}^T(x - y) \quad (\text{A'.27})$$

By adding the last two inequalities we get

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) \quad (\text{A'.28})$$

[2]  $\Rightarrow$  [3]

$$f(y) \geq f(x) + s_x^T(y-x) \quad (\text{A'.29})$$

$$f(x) \geq f(y) + s_y^T(x-y) + \quad (\text{A'.30})$$

$$(s_y - s_x)^T(y-x) \geq 0 \quad (\text{A'.31})$$

[3]  $\Rightarrow$  [2]

Since the set of subgradients is non-empty for all  $x, y$  then [2] is just the definition of the subgradient.  $\blacksquare$

**Lemma A'.2.3** (Equivalence for Strong convexity). *Suppose  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  with the extended-value extension. Then the following statements are equivalent:*

[1]  $f$  is strongly convex with parameter  $\mu$ .

[2]  $f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y) - \frac{\mu}{2}\alpha(1-\alpha)\|y-x\|^2$  for any  $x, y$  and  $\alpha \in [0, 1]$ .

[3]  $f(y) \geq f(x) + s_x^T(y-x) + \frac{\mu}{2}\|y-x\|^2$  for all  $x, y$  and any  $s_x \in \partial f(x)$ .

[4]  $(s_y - s_x)^T(y-x) \geq \mu\|y-x\|^2$  for all  $x, y$  and any  $s_x \in \partial f(x), s_y \in \partial f(y)$ .

[5]  $g(x) - \frac{\mu}{2}\|x\|^2$  is convex.

*Proof.* [5]  $\Rightarrow$  [3], [2], [4]

By definition of convexity of  $g$ .

[4]  $\Rightarrow$  [5]

By applying that  $s_x^g = s_x^f - \mu x$  and doing the calculations.

[2]  $\Rightarrow$  [3]

It is

$$f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y) - \frac{\mu}{2}\alpha(1-\alpha)\|y-x\|^2 \quad (\text{A'.32})$$

$$f(\alpha x + (1-\alpha)y) \geq f(x) + s_x^T(\alpha x + (1-\alpha)y - x) \quad (\text{A'.33})$$

By combining them we get

$$f(y) \geq f(x) + s_x^T(y-x) + \frac{\mu}{2}\alpha\|y-x\|^2 \quad (\text{A'.34})$$

$$f(y) \geq f(x) + s_x^T(y-x) \quad (\text{A'.35})$$

where we have set  $\alpha = 1$  [3]  $\Rightarrow$  [2], [3]  $\Rightarrow$  [4]

The proofs are similar to the ones in the previous lemma.  $\blacksquare$

**Lemma A'.2.4** (Implications of Strong Convexity). *Suppose  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  with the extended-value extension. The following conditions are all implied by strong convexity with parameter  $\mu$ :*

[1]  $\frac{1}{2}\|s_x\|^2 \geq \mu(f(x) - f^*)$ .

[2]  $\|s_y - s_x\| \geq \mu\|y-x\|$ .

[3]  $f(y) \leq f(x) + s_x^T(y-x) + \frac{1}{2\mu}\|s_y - s_x\|^2$ .

[4]  $(s_y - s_x)^T(y-x) \leq \frac{1}{\mu}\|s_y - s_x\|^2 \forall x, y$  and any  $s_x \in \partial f(x), s_y \in \partial f(y)$ .

*Remark 6.* A point  $x^*$  is a minimizer of a function  $f$  iff  $f$  is subdifferentiable at  $x^*$  and  $0 \in \partial f(x^*)$

*Proof.* [1] Since  $f$  is strong convex, we have

$$f(y) \geq f(x) + s_x^T(y-x) + \frac{\mu}{2}\|y-x\|^2 \quad (\text{A'.36})$$

By minimizing both parts of this inequality, we get  $f^*$  for the first part and for the second  $y = x - \frac{s_x}{\mu}$  and by substituting we get:

$$f^* \geq f(x) - \frac{\|s_x\|^2}{\mu} + \frac{\|s_x\|^2}{2\mu} \quad (\text{A'.37})$$

$$\frac{1}{2}\|s_x\|^2 \geq \mu(f(x) - f^*) \quad (\text{A'.38})$$

[2] Since  $f$  is strongly convex:

$$(s_x - s_y)^T(x-y) \geq \mu\|x-y\|^2 \quad (\text{A'.39})$$

$$\|s_x - s_y\|\|x-y\| \geq \mu\|x-y\|^2 \quad (\text{A'.40})$$

$$\|s_x - s_y\| \geq \mu\|x-y\| \quad (\text{A'.41})$$

[3] Consider the functions  $h_x(z) = f(z) - s_x^T z$ . Then the subgradient of  $h_x(z)$ , say  $g_z$ , equals  $s_z - s_x$ . So, we have

$$(s_{z_1} - s_{z_2})^T(z_1 - z_2) \geq \mu\|z_1 - z_2\|^2 \quad (\text{A'.42})$$

$$(g_{z_1} - s_{z_2})^T(z_1 - z_2) \geq \mu\|z_1 - z_2\|^2 \quad (\text{A'.43})$$

So  $h_x(z)$  is strong convex and using the [1] of this lemma for  $h_x$  we get

$$h_x^* \geq h_x(y) - \frac{1}{2\mu}\|g_y\|^2 \Rightarrow \quad (\text{A'.44})$$

$$f(x) - s_x^T x \geq f(y) - s^T y - \frac{1}{2\mu}\|s_y - s_x\|^2 \quad (\text{A'.45})$$

$$f(y) \leq f(x) + s_x^T(y-x) + \frac{1}{2\mu}\|s_y - s_x\|^2 \quad (\text{A'.46})$$

[4] It can be derived from [3] with change of variables and adding the two inequalities.  $\blacksquare$

**Lemma A'.2.5.** For a function  $f$  with Lipschitz continuous gradient over  $\mathbb{R}^n$ , the following relations hold:

$$[5] \Leftrightarrow [7] \Rightarrow [6] \Rightarrow [0] \Rightarrow [1] \Leftrightarrow [2] \Leftrightarrow [3] \Leftrightarrow [4] \quad (\text{A'.47})$$

If the function  $f$  is convex, then all the conditions [0]-[7] are equivalent.

$$[0] \|\nabla f(x) - \nabla f(y)\| \leq L\|y - x\|, \forall x, y. \quad (\text{A'.48})$$

$$[1] g(x) = \frac{L}{2}x^T x - f(x) \text{ is convex, for all } x. \quad (\text{A'.49})$$

$$[2] f(y) \leq f(x) + \nabla f(x)^T(y - x) + \frac{L}{2}\|y - x\|^2, \forall x, y. \quad (\text{A'.50})$$

$$[3] (\nabla f(x) - \nabla f(y))^T(y - x) \leq L\|x - y\|^2, \forall x, y. \quad (\text{A'.51})$$

$$[4] f(xy + (1 - x)a) \leq xf(y) + (1 - x)f(a) - \frac{1(1 - a)L}{2}\|x - y\|^2, \forall x, y \text{ and } a \in [0, 1]. \quad (\text{A'.52})$$

$$[5] f(y) \geq f(x) + s_x^T(y - x) + \frac{1}{2L}\|\nabla f(y) - \nabla f(x)\|^2, \forall x, y. \quad (\text{A'.53})$$

$$[6] (\nabla f(x) - \nabla f(y))^T(x - y) \geq \frac{1}{L}\|\nabla f(x) - \nabla f(y)\|^2, \forall x, y. \quad (\text{A'.54})$$

$$[7] f(ax + (1 - a)y) \leq af(x) + (1 - a)f(y) - \frac{a(1 - a)}{2L}\|\nabla f(x) - \nabla f(y)\|^2, \forall x, y \text{ and } a \in [0, 1]. \quad (\text{A'.55})$$

*Proof.* [1]  $\Leftrightarrow$  [2]

It is,

$$\nabla g(x) = Lx - \nabla f(x) \quad (\text{A'.56})$$

Furthermore (all these steps perform equivalences)

$$g(y) \geq g(x) + \nabla g(x)^T(y - x) \quad (\text{A'.57})$$

$$\frac{L}{2}\|y\|^2 - f(y) \geq \frac{L}{2}\|x\|^2 - f(x) + (Lx - \nabla f(x))^T(y - x) \quad (\text{A'.58})$$

$$f(y) \leq f(x) + \frac{L}{2}\|y\|^2 + \frac{L}{2}\|x\|^2 - Lx^T y + \nabla f(x)(y - x) \quad (\text{A'.59})$$

$$f(y) \leq f(x) + \nabla f(x)(y - x) + \frac{L}{2}\|y - x\|^2 \quad (\text{A'.60})$$

[2]  $\Rightarrow$  [3] Interchange  $x, y$  and add.

[3]  $\Rightarrow$  [1] Substitute  $\nabla f(x) = Lx - \nabla g(x)$  and conclude the monotonicity of gradient for  $g$ .

[1]  $\Leftrightarrow$  [4]

$$g(ax + (1 - a)y) \leq ag(x) + (1 - a)g(y) \Leftrightarrow \quad (\text{A'.61})$$

$$\frac{L}{2}(ax + (1 - a)y)^T(ax + (1 - a)y) - f(x') \leq \frac{L}{2}a\|x\|^2 - af(x) + \frac{L}{2}(1 - a)\|y\|^2 - (1 - a)f(y) \Leftrightarrow \quad (\text{A'.62})$$

$$f(x') \geq af(x) + (1 - a)f(y) - \frac{L}{2}[a(1 - a)\|x\|^2 + a(1 - a)\|y\|^2 - 2a(1 - a)x^T y] \Leftrightarrow \quad (\text{A'.63})$$

$$f(x') \geq af(x) + (1 - a)f(y) - \frac{L}{2}a(1 - a)\|x - y\|^2 \quad (\text{A'.64})$$

[0]  $\Rightarrow$  [3]

$$(\nabla f(x) - \nabla f(y))^T(x - y) \leq \|\nabla f(x) - \nabla f(y)\|\|x - y\| \leq L\|y - x\|^2 \quad (\text{A'.65})$$



[5]  $\Rightarrow$  [6] Change  $x, y$  and add.

[5]  $\Rightarrow$  [7] Following the same procedure as to prove that first order implies Jensen's inequality we get

$$f(x') \leq af(x) + (1-a)f(y) - \frac{1-a}{2L} \|\nabla f(x) - \nabla f(x')\| - \frac{a}{2L} \|\nabla f(x') - \nabla f(y)\| \quad (\text{A'.66})$$

It also holds  $\forall x, y$

$$a\|x\|^2 + (1-a)\|y\|^2 \geq a(1-a)\|x+y\|^2 \Leftrightarrow \quad (\text{A'.67})$$

$$a^2\|x\|^2 - 2a(1-a)x^T y + (1-a)^2\|y\|^2 \geq 0 \Leftrightarrow \quad (\text{A'.68})$$

$$\|ax - (1-a)y\|^2 \geq 0 \text{ which is true } \forall x, y \quad (\text{A'.69})$$

Combining the above two results we get [7].

[7]  $\Rightarrow$  [5] We follow the same procedure as the one to prove that Jensen's inequality implies first order criterion. [6]  $\Rightarrow$  [0] again by using Cauchy-Swartz inequality.

Also if  $f$  is convex it is sufficient to show that one of [1], [2], [3], [4] implies [5] or [7].

[3]  $\Rightarrow$  [5]

As before we will define the function  $h_x(z) = f(z) - \nabla f(x)^T z$  and then  $\nabla h_x(z) = \nabla f(z) - \nabla f(x)$ . Then

$$(\nabla f(z_1) - \nabla f(z_2))^T (z_1 - z_2) \leq L\|z_1 - z_2\|^2 \quad (\text{A'.70})$$

$$(\nabla h_x(z_1) - \nabla h_x(z_2))^T (z_1 - z_2) \leq L\|z_1 - z_2\|^2 \quad (\text{A'.71})$$

$$h_x(z) \leq h_x(y) + \nabla h_x(y)^T (z - y) + \frac{L}{2} \|z - y\|^2 \quad (\text{A'.72})$$

where we used the fact that [2]  $\Leftrightarrow$  [3] to go from the second to the third inequality. Since  $f$  is convex from the first order criterion we get that  $h_x(z)$  attains its minimum when  $z = x$ . Also by taking minimization of the right part of the last inequality we get  $z = y - \frac{1}{L} \nabla h_x(y)$  and by applying this minimization two both ends of the inequality we get [5].  $\blacksquare$

**Lemma A'.2.6.** Consider the following conditions for a general function  $f$ :

$$[1] f^*(s) = s^T x - f(x).$$

$$[2] s \in \partial f(x).$$

$$[3] x \in \partial f^*(s).$$

Then, we have

$$[1] \Leftrightarrow [2] \Rightarrow [3] \quad (\text{A'.73})$$

Further, if  $f$  is closed and convex, then all these conditions are equivalent.

*Proof.* [1]  $\Leftrightarrow$  [2]

$$s \in \partial f(x) \Leftrightarrow f(y) \geq f(x) + s^T (y - x) \quad (\text{A'.74})$$

$$\Leftrightarrow s^T x - f(x) \geq s^T y - f(y) = f^*(s) \quad (\text{A'.75})$$

Also by the definition of Fenchel conjugate  $f^*(s) \geq s^T x - f(x)$ . [2]  $\Rightarrow$  [3]

$$s \in \partial f(x) \Rightarrow f(z) \geq f(x) + s^T (z - x) \quad (\text{A'.76})$$

Also, [2]  $\Rightarrow$  [3]

$$f^*(z) \geq z^T x - f(x) \geq z^T x - f(z) + s^T(z - x) = f^*(s) + x^T(z - s) \Rightarrow x \in \partial f^*(s) \quad (\text{A'.77})$$

If  $f$  is closed and convex [3]  $\Rightarrow$  [2]

**Lemma A'.2.7.** *If  $f$  is convex and closed then  $f^{**} = f$ .*

*Proof.* We know that  $f^*(s) = \sup_{x \in \text{dom} f} (s^T x - f(x))$  and  $f^{**}(x) = \sup_{s \in \text{dom} f^*} (x^T s - f^*(s))$ . Now suppose that  $f^*(s) = s^T y - f(y)$  for some  $y$ , we know that this holds due to the closeness of the linear function and  $f$  ( $f$  is closed  $\Rightarrow$  the sublevel sets of  $f$  are closed and so the supremum can be achieved) which implies that  $s \in \partial f(y)$ . Then,

$$f^{**}(x) = \sup_s (s^T x - f^*(s)) \geq s^T x - f^*(s) \quad \forall s \quad (\text{A'.78})$$

If we choose  $s \in \partial f(x) \Rightarrow f^*(s) = s^T x - f(x)$  we have

$$f^{**}(x) \geq s^T x - f^*(s) \geq f(x) \quad (\text{A'.79})$$

Also, from the previous observation we have

$$f^{**}(x) = \sup_{s \in \partial f(y)} (s^T x - s^T y + f(y)) \quad s \in \partial f(y) \quad (\text{A'.80})$$

By convexity of  $f$ ,  $\forall x, y$  and  $s \in \partial f(y)$

$$f(x) \geq f(y) + s^T(x - y) \quad (\text{A'.81})$$

and so  $f^{**}(x) \leq f(x)$ , which means that  $f^{**}(x) = f(x)$  ■

So,

$$f(y) \geq f(x) \geq x^T s - f^* s \quad (\text{A'.82})$$

$$x \in \partial f^*(s) \Rightarrow f^*(z) \geq f^*(s) + x^T(z - s) \quad (\text{A'.83})$$

Finally

$$f^{**}(z) \geq z^T s - f^*(s) \geq z^T s - f^*(z) + x^T(z - s) \quad (\text{A'.84})$$

$$f(z) \geq f(x) + s^T(z - x) \quad (\text{A'.85})$$

$$s \in \partial f(x) \quad (\text{A'.86})$$

■

**Lemma A'.2.8** (Differentiability). *For a closed and strictly convex  $f$ ,  $\nabla f^*(s) = \arg \max_x (s^T x - f(x))$*

*Proof.* Since  $f$  is strictly convex and closed all of the properties from the previous lemma are equivalent. Also,  $f^*(s) = \sup_z \{s^T z - f(z)\}$ . So, we know that the supremum of  $s^T z - f(z)$  is achieved when  $z = x$  and also

$$\nabla f^*(s) = x = \arg \max_x (s^T x - f(x)) \quad (\text{A'.87})$$

Now we only have to show that for two points  $x_1 \neq x_2$  s.t. the  $x_1 \in \partial f^*(s)$  and  $x_2 \in \partial f^*(s)$ , which means that  $s \in \partial f(x_1)$  and  $s \in \partial f(x_2)$ . If there were then  $\forall z$

$$f(z) \geq f(x_1) + s^T(z - x_1) \quad (\text{A'.88})$$

$$f(z) \geq f(x_2) + s^T(z - x_2) \quad (\text{A'.89})$$

With strict inequality for  $z \neq x_1$  and  $z \neq x_2$  equivalently. Then by using  $z = x_2$  and  $z = x_1$  and adding the two inequalities we conclude that  $0 > 0$  which a contradiction.  $\blacksquare$

**Theorem A'.2.9.** *A function  $f$  and its Fenchel conjugate function  $f^*$  satisfy the following assertions:*

1. *if  $f$  is closed and strong convex with parameter  $\mu$ , then  $f^*$  has a Lipschitz continuous gradient with parameter  $\frac{1}{\mu}$ .*
2. *If  $f$  is convex and has Lipschitz continuous gradient with parameter  $L$ , then  $f^*$  is strong convex with parameter  $\frac{1}{L}$ .*

*Proof.* We start by proving 1. If  $f$  is strong convex and closed and  $f^*(s) = s^T x - f(x)$  and  $f^*(p) = p^T y - f(y)$  then

$$\|s_x - s_y\| \geq \mu \|x - y\| \quad (\text{A'.90})$$

$$\frac{1}{\mu} \|s - p\| \geq \|\nabla f^*(s) - \nabla f^*(p)\| \quad (\text{A'.91})$$

We now proceed to the proof of the second claim. Again since  $f$  is convex all the previous properties hold and we have

$$\|\nabla f(x) - \nabla f(y)\| \leq L \|y - x\| \quad (\text{A'.92})$$

$$\|\nabla f^*(s) - \nabla f^*(p)\| \geq \frac{1}{L} \|p - s\| \quad (\text{A'.93})$$

$\blacksquare$

These proofs can also be found in [51].

# Appendix B'

## Deferred Proofs

### B'.1 Bregman Divergence and Fenchel Coupling

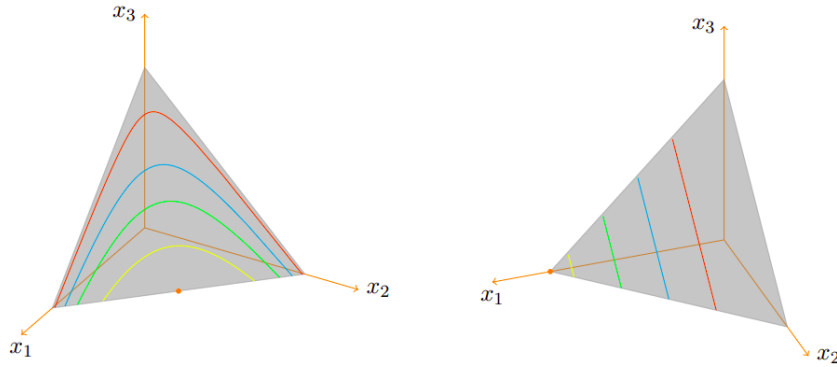


Figure B'.1: The level sets of KL-divergence

#### B'.1.1 Bregman Divergence

Bregman divergence provides a way to measure the distance of two points that belong to the simplex. Its properties render it a useful tool to prove convergence results. Below we state its definition and prove these properties that would be crucial in the establishment of our proof. Given a fixed point  $p \in \mathcal{X}$  then the Bregman divergence of a function  $h$  is defined for all points  $x \in \mathcal{X}$  as

$$D_h(p, x) = h(p) - h(x) - h'(x; p - x) \text{ for all } p, x \in \mathcal{X} \quad (\text{B'.1})$$

where  $h'(x; p - x)$  is the one-sided derivative

$$h'(x; p - x) \equiv \lim_{t \rightarrow 0^+} t^{-1} [h(x + t(p - x)) - h(x)] \quad (\text{B'.2})$$

Notice that this definition of the Bregman divergence permits to work also with points on the boundary. It is possible that the limit of  $D_h$  attains the value of  $+\infty$  if  $h'(x; p - x) = -\infty$ , as

$x \rightarrow p$ , where  $p$  is a point of the boundary. However, the condition below ensures that this is not the case.

$$D_h(p; x) \rightarrow 0 \text{ whenever } x \rightarrow p \quad (\text{Reciprocity})$$

This is known as the reciprocity condition. What this property actually means is that the sublevel sets of  $D(p, \cdot)$  are neighborhoods of  $p$ . This is illustrated in B'.1, when the function employed is the negative Shannon-entropy and the induced Bregman divergence the Kullback–Leibler divergence. Notice that for most decomposable functions  $h$ , this property holds. Below we present a proof of this statement.

**Proposition B'.1.1.** *If  $h(x) = \sum_i \theta(x_i)$ , for some kernel function  $\theta$  having the properties described in (regularizer's properties) and furthermore it holds that  $\theta'(x) = o(1/x)$  for  $x$  close to 0, then  $D_h(p; x) \rightarrow 0$  whenever  $x \rightarrow p$  for all  $x, p \in \mathcal{X}$ .*

*Proof.* It is sufficient to prove that  $\lim_{x \rightarrow 0} (\theta(0) - \theta(x) - \theta'(x)(0 - x)) = 0$ . The difference of the first two terms is obviously gives zero. Now, for the last term notice that if  $\theta'(x) = o(1/x)$  for  $x$  close to 0, then  $\lim_{x \rightarrow 0} x\theta'(x) = 0$  and the proof is completed. ■

Additionally, Bregman divergence satisfies the properties described below.

**Proposition B'.1.2.** *Let  $h$  be a  $K$ -strongly convex function defined on the simplex  $\mathcal{X} = \Delta(\mathcal{A})$ , that has the properties described in regularizer's properties and let  $\Delta_p$  be the union of the relative interiors of the faces of  $\mathcal{X}$  that contain  $p$  i.e.,*

$$\Delta_p = \{x \in \mathcal{X} : \text{supp}(p) \subseteq \text{supp}(x)\} = \{x \in \mathcal{X} : x_a > 0 \text{ whenever } p_a > 0\} \quad (\text{B'.3})$$

Then

1.  $D_h(p, x) < \infty$  whenever  $x \in \Delta_p$ .
2.  $D_h(p, x) \geq 0$  for all  $x \in \mathcal{X}$ , with equality if and only if  $p = x$ , more particularly

$$D_h(p, x) \geq \frac{1}{2}K\|x - p\|^2 \text{ for all } x \in \mathcal{X} \quad (\text{B'.4})$$

*Proof.* For the first part, if  $x \in \Delta_p$  then  $h(x + t(x - p))$  is finite and smooth in a neighborhood of 0 and thus  $D(p, x)$  is also finite.

The second part of the proposition, let  $z = x - p$  then strong convexity yields

$$\begin{aligned} h(x + tz) &\leq th(p) + (1 - t)h(x) - \frac{1}{2}Kt(1 - t)\|x - p\|^2 \\ t^{-1}(h(x + tz) - h(x)) &\leq h(p) - h(x) - \frac{1}{2}(1 - t)K\|x - p\|^2 \\ h(p) - h(x) - t^{-1}(h(x + tz) - h(x)) &\geq \frac{1}{2}(1 - t)K\|x - p\|^2 \end{aligned}$$

And by taking  $t \rightarrow 0$ , we obtain the result. ■

We mention at this point that from (regularizer's properties), since for each  $i \in \mathcal{N}$ :  $\inf_{\epsilon \in [0, 1]} \theta_i'' > 0$ , there exists  $K_i > 0$  such that for all  $x, y \in [0, 1]$  and  $t \in [0, 1]$

$$\theta_i(tx + (1 - t)y) \leq t\theta_i(x) + (1 - t)\theta_i(y) - \frac{K_i}{2}t(1 - t)|x - y|^2 \quad (\text{B'.5})$$

In all the proofs  $h$  symbolizes the aggregate function of all the regularizers i.e.,  $h(x) = \sum_i h_i(x_i)$ , with strong convexity parameter  $K \equiv \min_i K_i$ .

### B'.1.2 Steep vs non-steep

In this section we elaborate in detail the dichotomy of the properties of different regularizers mentioned in the remark 4. As we mentioned players may have different regularizers  $h_i$  employed in their choice maps  $Q_i(y) = \arg \max_{x \in \mathcal{X}_i} \{\langle x, y \rangle - h_i(x)\}$ . Depending on the regularizer chosen, FTRL dynamics may differ significantly. To formally express this difference, it is convenient to consider that  $h$  is an extended-real valued function  $h : \mathcal{V} \rightarrow \mathbb{R} \cup \{\infty\}$  with value  $\infty$  outside of the simplex  $\mathcal{X}$ . Then the subdifferential of  $h$  at  $x \in \mathcal{V}$  is defined as:

$$\partial h(x) = \{y \in \mathcal{V}^* : h(x') \geq h(x) + \langle y, x' - x \rangle \forall x' \in \mathcal{V}\} \quad (\text{B'.6})$$

If  $\partial h(x)$  is nonempty, then  $h$  is called subdifferentiable at  $x \in \mathcal{X}$ . When  $x \in \text{ri}(\mathcal{X})$  then  $\partial h(x)$  is always non-empty or  $\text{ri}(\mathcal{X}) \subseteq \text{dom } \partial h \equiv \{x \in \mathcal{X} : \partial h(x) \neq \emptyset\}$ . Notice that when the gradient of  $h$  exists, then its subgradient always contains it. With these in mind, we present a typical separation between the different regularizers. On the one hand, *steep* regularizers like the negative Shannon-entropy become infinitely steep as  $x$  approaches the boundary or  $\|\nabla h(x)\| \rightarrow \infty$ . On the other hand, *non-steep* are everywhere differentiable, like the Euclidean, allowing the sequence of play to transfer between the different faces of the simplex. In the dual space of payoffs, steepness implies that the choice map is not surjective (since it cannot map all payoff vectors to points of the boundary), it is however injective (it maps a payoff vector plus a multiple of  $(1, 1, \dots, 1)$  to the same strategy). Non-steep regularizers give rise to surjective maps, which are not injective, not even up to a multiple of  $(1, 1, \dots, 1)$ , to the boundary. Focusing on the more simple case of decomposable regularizers, the kernel of a steep one is differentiable on  $(0, 1]$  while for non-steep the kernel is differentiable in all of  $[0, 1]$ . As a result, when a steep regularizer is employed the mirror map  $Q : \mathcal{Y} \rightarrow \mathcal{X}$  cannot return any point of the boundary. In other words, the points of the boundary are infeasible not only as initial conditions but also as part of the sequence of play.

*Remark 7.* This dichotomy is important for our analysis since we study the stochastic asymptotic stability of Nash equilibria, which may lie on the boundary, and we seek a neighborhood of initial conditions such that the equilibrium to be stable and attracting. Thus, instead of demanding the existence of a neighborhood  $U$  of an equilibrium  $x^*$ , such that whenever  $X_0 \in U$ ,  $x^*$  is stable and attracting; we demand the existence of a neighborhood  $U$  of  $x^*$  such that whenever  $X_0 \in U \cap \text{im } Q$  then  $x^*$  is stable and attracting.

### B'.1.3 Polar Cone

The notion of the polar cone is tightly connected with the notion of duality. Given a finite dimensional vector space  $\mathcal{V}$ , a convex set  $\mathcal{C} \subseteq \mathcal{V}$  and a point  $x \in \mathcal{C}$  the tangent cone  $\text{TC}_{\mathcal{C}}(x)$  is the closure of the set of all rays emanating from  $x$  and intersecting  $\mathcal{C}$  in at least one other point. The dual of the tangent cone is the polar cone  $\text{PC}_{\mathcal{C}}(x) = \{y \in \mathcal{V}^* : \langle y, z \rangle \leq 0 \text{ for all } z \in \text{TC}_{\mathcal{C}}(x)\}$ .

When the under consideration convex set is the simplex of the players' strategies, the polar cone corresponding to the boundary differs significantly from the one corresponding to the interior. Formally, the polar cone at a point  $x$  of the simplex is

$$\text{PC}(x) = \{y \in \mathcal{Y} : y_a \geq y_b \text{ for all } a, b \in \mathcal{A}\}^1 \quad (\text{B'.7})$$

An illustration of this is depicted in figure B'.2. When (FTRL) is run, the notion of the polar cone emerges from the choice map  $Q : \mathcal{Y} \rightarrow \mathcal{X}$ , connecting the primal space of the strategies with the dual space of the payoffs. The proposition below presents this exact connection.

<sup>1</sup>It is always  $y_a = y_b$  whenever  $a, b \in \text{supp}(x)$ .

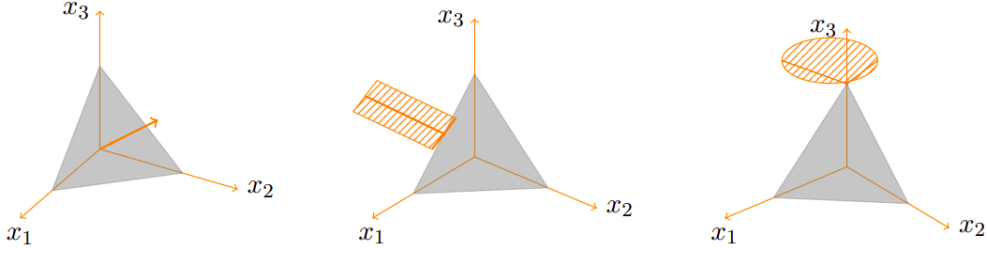


Figure B'.2: The polar cone corresponding to different points of the simplex. For an interior point this is a line perpendicular to the simplex. For a point of the boundary, it is a plane perpendicular to the simplex tangential to the point of the boundary. For an edge the polar cone corresponds to a cone.

**Proposition B'.1.3.** *Let  $h$  be a strong convex regularizer that satisfies the properties described in [regularizer's properties](#) and let  $Q : \mathcal{Y} \rightarrow \mathcal{X}$  be the induced choice map then*

1.  $x = Q(y) \Leftrightarrow y \in \partial h(x)$
2.  $\partial h(x) = \nabla h(x) + \text{PC}(x)$  for all  $x \in \mathcal{X}$ .

### B'.1.4 Fenchel Coupling

Even though Bregman divergence is a useful tool, (FTRL) evolves in the dual space of payoffs. Thus dually to the above the Fenchel coupling<sup>2</sup> is defined,  $F_h : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$

$$F_h(p, y) = h(p) + h^*(y) - \langle y, p \rangle \text{ for all } p \in \mathcal{X}, y \in \mathcal{Y} \quad (\text{B'.8})$$

where  $h^* : \mathcal{Y} \rightarrow \mathbb{R}$  is the convex conjugate of  $h$ :  $h^*(y) = \sup_{x \in \mathcal{X}} \{\langle y, x \rangle - h(x)\}$ . The fenchel conjugate is differentiable on  $\mathcal{Y}$  and it holds that

$$\nabla h^*(y) = Q(y) \text{ for all } y \in \mathcal{Y} \quad (\text{B'.9})$$

Fenchel coupling is also a measure that connects the primal with the dual space. As we mentioned above, (FTRL) evolves in the dual space and thus we use Fenchel coupling to trace its convergence properties. As the next proposition states, whenever Fenchel coupling  $F(p, y)$  is bounded from above so does  $\|Q(y) - p\|$ . This proposition in its entirety, is critical for our proof, since we first need to find a neighborhood  $U$  of attractness (See 4.2.1). For this step, Bregman divergence is necessary in order to define the aforementioned neighborhood since  $\|Q(y) - p\| < c$  for some constant  $c$  is not necessarily a neighborhood of  $p$  (See section B'.1.2).

**Proposition B'.1.4.** *Let  $h$  be a  $K$ -strongly convex function on  $\mathcal{X}$  and has the properties described in [regularizer's properties](#). Let  $p \in \mathcal{X}$ , then*

1.  $F_h(p, y) \geq \frac{1}{2}K\|Q(y) - p\|^2$  for all  $y \in \mathcal{Y}$  and whenever  $F_h(p, y) \rightarrow 0$ ,  $Q(y) \rightarrow p$ .
2.  $F_h(p, y) = D_h(p, x)$  whenever  $Q(y) = x$  and  $x \in \Delta_p$ .
3.  $F_h(p, y') \leq F_h(p, y) + \langle y' - y, Q(y) - p \rangle + \frac{1}{2K}\|y' - y\|_*^2$ .

<sup>2</sup>The term is due to [16].

*Remark 8.* Notice that the first part of the proposition is not implied by the second one, since it is possible that  $\text{im } Q = \text{dom } \partial h$  is not always contained in  $\Delta_p$  (see section [B'.1.2](#)).

*Proof.* For the first part, let  $x = Q(y)$  then  $h^*(y) = \langle y, x \rangle - h(x)$

$$F_h(p, y) = h(p) - h(x) - \langle y, p - x \rangle \quad (\text{B'.10})$$

Since  $y \in \partial h(x)$  (Proposition [B'.1.3](#)), it is

$$h(x + t(p - x)) \geq h(x) + t\langle y, p - x \rangle \quad (\text{B'.11})$$

and by strong convexity of  $h$ , we have

$$h(x + t(p - x)) \leq th(p) + (1 - t)h(x) - \frac{1}{2}Kt(1 - t)\|p - x\|^2 \quad (\text{B'.12})$$

Thus by combining [\(B'.11\)](#),[\(B'.12\)](#) and taking  $t \rightarrow 0$  we get

$$F_h(p, y) \geq h(p) - h(x) - h(p) + h(x) + \frac{K}{2}\|p - x\|^2 \geq \frac{K}{2}\|p - x\|^2 \quad (\text{B'.13})$$

For the second part of the proposition, notice that  $x + t(p - x)$  lies in the relative interior of some face of  $\mathcal{X}$  for  $t$  in a neighborhood of 0 and thus  $h(x + t(p - x))$  is smooth and finite. So,  $h$  admits a two-sided derivative along  $x - p$  and since  $y \in \partial h(x)$ ,  $\langle y, p - x \rangle = h'(x; p - x)$  and our claim naturally follows.

Finally for the last part of the proposition, we have

$$\begin{aligned} F_h(p, y') &= h(p) + h^*(y') - \langle y', p \rangle \\ &\leq h(p) + h^*(y) + \langle y' - y, \nabla h^*(y) \rangle + \frac{1}{2K}\|y' - y\|_*^2 - \langle y', p \rangle \\ &= F_h(p, y) + \langle y' - y, Q(y) - p \rangle + \frac{1}{2K}\|y' - y\|_*^2 \end{aligned}$$

where the second inequality follows from the fact that  $h^*$  is  $1/K$  strongly smooth [\[18\]](#).  $\blacksquare$

In terms of Fenchel coupling our reciprocity assumption can be written as

$$F_h(p, y) \rightarrow 0 \text{ whenever } Q(y) \rightarrow p \quad (\text{Reciprocity})$$

Again for most of  $h$  decomposable, the assumption is turned into a property as we prove below.

**Proposition B'.1.5.** *If  $h(x) = \sum_i \theta(x_i)$ , with  $\theta$  having the properties described in [\(regularizer's properties\)](#) and furthermore it holds that  $\theta'(x) = o(1/x)$  for  $x$  close to 0, then  $F_h(p, y) \rightarrow 0$  whenever  $Q(y) \rightarrow p$  for all  $p \in \mathcal{X}$ .*

*Proof.* Again it is sufficient to prove that whenever  $Q(y) = x \rightarrow 0$  then  $F_h(p, y) \rightarrow 0$ . Notice that from [B'.1.4](#)  $F_h(p, y) = D_h(p, x)$  whenever  $x = Q(y)$  and  $x \in \Delta_p$ . Thus by [B'.1.1](#)  $Q(y) = x \rightarrow 0$  implies that  $F_h(p, y) \rightarrow 0$ .  $\blacksquare$



## B'.2 Variational stability

**Definition B'.2.1** (Variational stability). A point  $x^* \in \mathcal{X}$  is said to be *variationally stable* if there exists neighborhood  $U$  of  $x^*$  such that

$$\langle v(x), x - x^* \rangle \leq 0 \text{ for all } x \in U \quad (\text{VS})$$

with equality if and only if  $x = x^*$ .

What this property actually states is that in a neighborhood of  $x^*$ , it strictly dominates over all other strategies. Interestingly, strict Nash equilibria hold this property:

**Proposition B'.2.1.** *For finite games in normal form, the following are equivalent:*

- i)  $x^*$  is a strict Nash equilibrium.
- ii)  $\langle v(x^*), z \rangle \leq 0$  for all  $z \in \text{TC}(x^*)$  with equality if and only if  $z=0$ .
- iii)  $x^*$  is variationally stable.

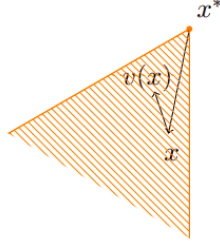


Figure B'.3: (VS) states that the payoff vectors are pointing "towards" the equilibrium

*Proof.* We will first prove that  $i) \Rightarrow ii) \Rightarrow iii) \Rightarrow i)$ .

$i) \Rightarrow ii)$  Since  $x^*$  is a Nash equilibrium by definition it holds for each player  $i$  that

$$\langle v(x^*), x - x^* \rangle \leq 0 \text{ for all } x \in \mathcal{X} \quad (\text{B'.14})$$

For the strict part of the inequality, by definition of strict Nash equilibria it holds that  $\langle v_i(x^*), x_i - x_i^* \rangle < 0$  whenever  $x_i \neq x_i^*$  and thus

$$\langle v(x^*), z \rangle = \sum_{i=1}^N \langle v_i(x^*), x_i - x_i^* \rangle < 0 \text{ if } x_i \neq x_i^* \text{ for some } i \text{ or } z \neq 0 \quad (\text{B'.15})$$

$ii) \Rightarrow iii)$  By definition of the polar cone, we have that  $v(x^*)$  belongs to the interior of  $\text{PC}(x^*)$ <sup>3</sup>. Thus by continuity there exists some neighborhood of  $x^*$  such that  $v(x)$  also belongs to the polar cone of  $\text{PC}(x^*)$  or  $x^*$  is variationally stable.

$iii) \Rightarrow i)$  Assume now that  $x^*$  is variationally stable but not strict, then there exist for some player  $i$   $a, b \in \mathcal{A}_i$  such that  $u_i(a; x_{-i}^*) = u_i(b; x_{-i}^*)$ . Then for  $x_i = x_i^* + \lambda(e_a - e_b)$  and  $x_{-i} = x_{-i}^*$  we have

$$\langle v(x^*), x - x^* \rangle = \langle v_i(x^*), \lambda(e_a - e_b) \rangle = 0 \quad (\text{B'.16})$$

which is a contradiction. ■

<sup>3</sup>Indeed if it belonged to the boundary then the equality in  $ii)$  would not hold only for  $z = 0$ .

### B'.3 Proofs of assumptions for Model 1, Model 2

Below we provide a proof for our claim in Model 2 that  $b_n = \mathcal{O}(\varepsilon_n)$ ,  $M_n^2 = \mathcal{O}(1/\varepsilon_n)$ . Focusing on one player  $i \in \mathcal{N}$ , notice that

$$\mathbb{E}[\hat{v}_{i,n} | \mathcal{F}_n] = \sum_{\alpha_{-i} \in \mathcal{A}_{-i}} \hat{X}_{-i,n}(u_i(\alpha_{i,1}; \alpha_{-i}), \dots, u_i(\alpha_{i,|\mathcal{A}_i|}; \alpha_{-i})) = v_i(\hat{X}_n) \quad (\text{B'.17})$$

Having this in mind  $\hat{v}_{i,n}$  can be viewed as

$$\hat{v}_{i,n} = v_i(X_n) + Z_{i,n} + b_{i,n} \quad (\text{B'.18})$$

where  $Z_{i,n} = \hat{v}_{i,n} - \mathbb{E}[\hat{v}_{i,n} | \mathcal{F}_n] = \hat{v}_{i,n} - v_i(\hat{X}_n)$  and  $b_{i,n} = v_i(\hat{X}_n) - v_i(X_n)$ . Thus, since  $v_i(x)$  is multi-linear in  $x$  and  $\hat{X}_{i,n} = (1 - \varepsilon_n)X_{i,n} + \varepsilon_n/|\mathcal{A}_i|$  it follows that  $b_n = \mathcal{O}(\varepsilon_n)$ . Finally, similarly to (B'.17) we can conclude that  $M_n^2 = \mathcal{O}(1/\varepsilon_n)$ .

We continue by proving that assumption (A3) is indeed satisfied for both Models 1, 2. This is due to the genericity of the game. Actually in the following lemma and corollaries we show that there exist player  $i \in \mathcal{N}$ , strategies  $a, b \in \text{supp}(x_i^*)$  and pure strategy profile  $\alpha_{-i} \in \text{supp}(x_{-i}^*)$ , where  $x^*$  is a mixed Nash equilibrium such that  $|u_i(a; \alpha_{-i}) - u_i(b; \alpha_{-i})| \geq \beta$  for some  $\beta > 0$ . In order to acquire the exact statement of (A3), we have to take into account the round in which the game is evolved. Let  $n > 0$  be this round, then when examining the stochastic asymptotic stability of a mixed Nash equilibrium  $x^*$ , the sequence of play is contained in a neighborhood of  $x^*$  and thus all of the strategies belonging to the support of  $x^*$  have strictly positive probability to be chosen, verifying the statement of (A3).

**Lemma B'.3.1.** *If the game is generic and has a mixed Nash equilibrium  $x^*$ , then there exist player  $i \in \mathcal{N}$ , pure strategies  $a, b \in \text{supp}(x_i^*)$  ( $a \neq b$ ) and pure strategy profile  $\alpha_{-i} \in \text{supp}(x_{-i}^*)$  such that  $u_i(a; \alpha_{-i}) \neq u_i(b; \alpha_{-i})$ .*

*Proof.* Assume that for all players  $i \in \mathcal{N}$ , pure strategy profiles  $\alpha_{-i} \in \text{supp}(x_{-i}^*)$  and pure strategies  $a, b \in \text{supp}(x_i^*)$  it is

$$u_i(a; \alpha_{-i}) = u_i(b; \alpha_{-i}) \quad (\text{B'.19})$$

Then for each player  $i$ , this implies that all of the payoffs corresponding to pure strategy profiles, which consists of the support of the equilibrium, are equal. Then each pure strategy profile  $(\alpha_i; \alpha_{-i}) \in \text{supp}(x^*)$  is a pure Nash equilibrium, which is a contradiction to the genericity of the game. ■

Immediate implications of lemma B'.3.1 are:

**Corollary B'.3.1.** *There exists player  $i \in \mathcal{N}$  and pure strategy profile  $(\alpha_i; \alpha_{-i}) \in \text{supp}(x^*)$ , such that  $u_i(\alpha_i; \alpha_{-i}) \neq 0$ .*

**Corollary B'.3.2.** *There exist  $\beta' > 0$ , player  $i$ , strategies  $a, b \in \text{supp}(x_i^*)$  and pure strategy profile  $\alpha_{-i} \in \text{supp}(x_{-i}^*)$  such that  $|u_i(a; \alpha_{-i}) - u_i(b; \alpha_{-i})| \geq \beta'$ . There also exist  $\beta'' > 0$  and  $(\alpha_i; \alpha_{-i}) \in \text{supp}(x^*)$  such that  $|u_i(\alpha_i; \alpha_{-i})| \geq \beta''$ .*

### B'.4 Proofs of Stability

#### B'.4.1 Deferred Proof of theorem 4.3.1

In the following preliminary result, we focus on the case of (FTRL) with payoff feedback as described in section 4.1.1 and we show that if  $x^*$  is a *strict* Nash equilibrium, there exists a

subsequence of  $(X_n)_{n=0}^\infty$  that converges to it. In order to achieve this convergence result, it is necessary to assume that the sequence  $(X_n)_{n=0}^\infty$  is contained in a neighborhood of  $x^*$ , in which (VS) holds. Here, we outline the basic steps below:

- Step 0: By contradiction, assume that there exists a neighborhood, in which  $X_n$  is not contained for all sufficiently large  $n$  and assume without loss of generality that holds for all  $n = 0, 1, \dots$
- Step 1: We start by showing that the terms of the RHS of the third property described in proposition B'.1.4 are converging almost surely to finite values, except for one. This term, which is a consequence of  $x^*$  being variational stable, goes to  $-\infty$  as  $n \rightarrow \infty$ .
- Step 2: The next crucial observation is that the Fenchel coupling is bounded from below by 0, thanks to the first property in proposition B'.1.4, which gives us the contradiction.

*Remark.* For the interested reader, the assumption (A2),  $\sum_n \gamma_n^2 M_n^2 < \infty$ , that we use in the preliminary lemma and in theorem 4.3.1 could be relaxed by using the Hölder inequality to  $\sum_n \gamma_n^{1+q/2} M_n^q < \infty$  for any  $q \in [2, \infty)$ .

**Lemma B'.4.1.** *Let  $x^* \in \mathcal{A}$  be a strict Nash equilibrium. If (FTRL) is run with payoff feedback of the type (4.1), that satisfies (A1)-(A2) and the sequence of play  $(X_n)_{n=0}^\infty$  does not exit a neighborhood  $\mathcal{R}$  of  $x^*$ , in which variational stability holds, then there exists a subsequence  $X_{n_k}$  of  $X_n$  that converges to  $x^*$  almost surely.*

*Proof.* Suppose that there exists a neighborhood  $U \subseteq \mathcal{R}$  of  $x^*$ , such that  $X_n \notin U$  for all large enough  $n$ . Assume without loss of generality that this is true for all  $n \geq 0$ . Since variational stability holds in  $\mathcal{R}$ , we have

$$\langle v(x), x - x^* \rangle < 0 \text{ for all } x \in \mathcal{R}, x \neq x^* \quad (\text{B'.20})$$

Furthermore, from proposition B'.1.4 we have that for each round  $n$ :

$$F_h(x^*, Y_{n+1}) \leq F_h(x^*, Y_n) + \gamma_n \langle \hat{v}_n, X_n - x^* \rangle + \frac{1}{2K} \gamma_n^2 \|\hat{v}_n\|_*^2 \quad (\text{B'.21})$$

By applying the above inequality for all rounds from 1, ...,  $n$  and creating the telescopic sum we get

$$F_h(x^*, Y_{n+1}) \leq F_h(x^*, Y_0) + \sum_{k=0}^n \gamma_k \langle \hat{v}_k, X_k - x^* \rangle + \frac{1}{2K} \sum_{k=0}^n \gamma_k^2 \|\hat{v}_k\|_*^2 \quad (\text{B'.22})$$

Remember that for the payoff vector holds that

$$\hat{v}_n = v(X_n) + Z_n + b_n$$

We now rewrite (B'.22)

$$\begin{aligned} F_h(x^*, Y_{n+1}) &\leq F_h(x^*, Y_0) + \sum_{k=0}^n \gamma_k \langle v(X_k), X_k - x^* \rangle + \sum_{k=0}^n \gamma_k \langle Z_k, X_k - x^* \rangle \\ &\quad + \sum_{k=0}^n \gamma_k \langle b_k, X_k - x^* \rangle + \frac{1}{2K} \sum_{k=0}^n \gamma_k^2 \|\hat{v}_k\|_*^2 \end{aligned} \quad (\text{B'.23})$$

Let  $\tau_n = \sum_{k=0}^n \gamma_k$  then

$$\begin{aligned} F_h(x^*, Y_{n+1}) &\leq F_h(x^*, Y_0) + \sum_{k=0}^n \gamma_k \langle v(X_k), X_k - x^* \rangle + \tau_n \left( \frac{\sum_{k=0}^n \gamma_k \langle b_k, X_k - x^* \rangle}{\tau_n} \right) \\ &\quad + \tau_n \left( \frac{\sum_{k=0}^n \gamma_k \langle Z_k, X_k - x^* \rangle}{\tau_n} + \frac{\frac{1}{2K} \sum_{k=0}^n \gamma_k^2 \|\hat{v}_k\|_*^2}{\tau_n} \right) \end{aligned} \quad (\text{B'.24})$$

We focus on the asymptotic behavior of each particular term of the previous inequality. We remind that  $\mathcal{F}_n$  denotes the history of  $X_n$  up to stage  $n$  (inclusive) and thus the feedback signal,  $\hat{v}_n$  is not  $\mathcal{F}_n$ -measurable in general.

- Let  $R_n = \sum_{k=0}^n \gamma_k^2 \|\hat{v}_k\|_*^2$ . Then

$$\mathbb{E}[R_n] \leq \sum_{k=0}^n \gamma_k^2 \mathbb{E}[\|\hat{v}_k\|_*^2] = \sum_{k=0}^n \gamma_k^2 \mathbb{E}[\mathbb{E}[\|\hat{v}_k\|_*^2 | \mathcal{F}_k]] \leq \sum_{k=0}^n \gamma_k^2 M_k^2 < \infty \quad (\text{B'.25})$$

where  $\sum_{k=0}^n \gamma_k^2 M_k^2$  is finite by assumption (A2). Hence by 1 and (B'.25)  $R_n$  is an  $L_1$  bounded submartingale while Doob's convergence theorem (A'.1.4) shows that almost surely

$$\lim_{n \rightarrow \infty} \tau_n^{-1} R_n = 0 \quad (\text{B'.26})$$

- Let  $S_n = \sum_{k=0}^n \gamma_k \langle Z_k, X_k - x^* \rangle$  and  $\psi_k = \gamma_k \langle Z_k, X_k - x^* \rangle$ . For the expected value of  $\psi_n$  we have

$$\mathbb{E}[\psi_n | \mathcal{F}_n] = \gamma_n \langle \mathbb{E}[Z_n | \mathcal{F}_n], X_n - x^* \rangle = 0 \quad (\text{B'.27})$$

and so  $S_n$  is a martingale since  $\mathbb{E}[S_n | \mathcal{F}_n] = S_{n-1}$ . Moreover, for the expectation of the absolute value of  $\psi_n$ , Cauchy-Schwarz inequality implies

$$\mathbb{E}[|\psi_n|^2 | \mathcal{F}_n] \leq \gamma_n^2 \mathbb{E}[\|Z_n\|_*^2 \|X_n - x^*\|^2 | \mathcal{F}_n] \quad (\text{B'.28})$$

$$\leq \gamma_n^2 \mathbb{E}[\|Z_n\|_*^2 | \mathcal{F}_n] \|\mathcal{X}\|^2 \quad (\text{B'.29})$$

$$\leq \gamma_n^2 M_n^2 \|\mathcal{X}\|^2 \quad (\text{B'.30})$$

since

$$\mathbb{E}[\|Z_n\|_*^2 | \mathcal{F}_n] = \mathbb{E}[\|\hat{v}_n - \mathbb{E}[\hat{v}_n | \mathcal{F}_n]\|_*^2 | \mathcal{F}_n] \quad (\text{B'.31})$$

$$= \mathbb{E}[\|\hat{v}_n\|_*^2 - 2\langle \hat{v}_n, \mathbb{E}[\hat{v}_n | \mathcal{F}_n] \rangle + \|\mathbb{E}[\hat{v}_n | \mathcal{F}_n]\|_*^2 | \mathcal{F}_n] \quad (\text{B'.32})$$

$$= \mathbb{E}[\|\hat{v}_n\|_*^2 | \mathcal{F}_n] - \|\mathbb{E}[\hat{v}_n | \mathcal{F}_n]\|_*^2 \quad (\text{B'.33})$$

$$\leq \mathbb{E}[\|\hat{v}_n\|_*^2 | \mathcal{F}_n] \leq M_n^2 \quad (\text{B'.34})$$

where  $M_n^2$  is the upper bound of  $\mathbb{E}[\|\hat{v}_n\|_*^2 | \mathcal{F}_n]$  described in section 4.1.1.

Obviously,  $\sum_{n=0}^{\infty} \tau_n^{-2} \mathbb{E}[|\psi_n|^2 | \mathcal{F}_n] < \infty$  and so by the strong law of large number for martingales (A'.1.3) yields that almost surely

$$\lim_{n \rightarrow \infty} \tau_n^{-1} S_n = 0 \quad (\text{B'.35})$$

- Let  $W_n = \sum_{k=0}^n \gamma_k \langle b_k, X_k - x^* \rangle$  then by Cauchy-Schwarz inequality

$$\begin{aligned} |\tau_n^{-1} W_n| &\leq |\tau_n^{-1} \sum_{k=0}^n \gamma_k \langle b_k, X_k - x^* \rangle| \leq \tau_n^{-1} \sum_{k=0}^n \gamma_k |\langle b_k, X_k - x^* \rangle| \\ &\leq \tau_n^{-1} \sum_{k=0}^n \gamma_k \|b_k\|_* \|\mathcal{X}\| \end{aligned} \quad (\text{B'.36})$$

Let  $J_n = \sum_{k=0}^n \gamma_k \|b_k\|_* \|\mathcal{X}\|$ . Notice that  $W_n \leq J_n$  and that from 1  $J_n$  is a submartingale with

$$\mathbb{E}[J_n] = \|\mathcal{X}\| \sum_{k=0}^n \gamma_k \mathbb{E}[\|b_k\|_*] \leq \|\mathcal{X}\| \sum_{k=0}^n \gamma_k \mathbb{E}[\mathbb{E}[\|b_k\|_* | \mathcal{F}_k]] \leq \|\mathcal{X}\| \sum_{k=0}^n \gamma_k B_k < \infty \quad (\text{B'.37})$$

where  $B_n$  is the upper bound of  $\mathbb{E}[\|b_n\|_* | \mathcal{F}_n]$ . Thus,  $J_n$  is a  $L_1$  bounded submartingale and by Doob's convergence theorem (A'.1.4) almost surely

$$\lim_{n \rightarrow \infty} \tau_n^{-1} J_n = 0 \quad (\text{B'.38})$$

As a result,  $\tau_n^{-1} W_n \rightarrow 0$ .

- Finally, we will examine the term  $\sum_{k=0}^n \gamma_k \langle v(X_k), X_k - x^* \rangle$ . Recall that we had assumed that  $X_n \in \mathcal{R} \setminus U$  for all  $n \geq 0$ , while variational stability holds in  $\mathcal{R}$ , so by continuity there exists  $c > 0$ , such that for all  $n \geq 0$

$$\langle v(X_n), X_n - x^* \rangle \leq -c \quad (\text{B'.39})$$

We return to (B'.24) and we equivalently we have that

$$\begin{aligned} F_h(x^*, Y_{n+1}) &\leq F_h(x^*, Y_0) + \sum_{k=0}^n \gamma_k \langle v(X_k), X_k - x^* \rangle + \tau_n (\tau_n^{-1} W_n + \tau_n^{-1} R_n + \tau_{-1}^n S_n) \\ &\leq F_h(x^*, Y_0) - c \tau_n + \tau_n (\tau_n^{-1} W_n + \tau_n^{-1} R_n + \tau_n^{-1} S_n) \end{aligned} \quad (\text{B'.40})$$

Thus,  $F_h(x^*, Y_{n+1}) \sim -c \sum_{k=0}^{\infty} \gamma_k \rightarrow -\infty$ .

By proposition B'.1.4 we conclude to a contradiction. This implies that some instance of the sequence of play is included to every neighborhood  $U$  of  $x^*$  and thus there exists subsequence  $X_{n_k}$  of  $X_n$  that almost surely converges to  $x^*$ .  $\blacksquare$

**Theorem B'.4.2** (Restatement of theorem 4.3.1). *Let  $x^*$  be a strict Nash equilibrium. If (FTRL) is run with payoff feedback that satisfies (A1)-(A2), then  $x^*$  is stochastically asymptotically stable.*

*Proof.* Fix a confidence level  $\delta$  and let  $U_\varepsilon = \{x : D_h(x^*, x) < \varepsilon\}$  and  $U_\varepsilon^* = \{y \in \mathcal{Y} : F_h(x^*, y) < \varepsilon\}$ .

- By proposition B'.1.2 for all  $x \in U_\varepsilon$  it holds that  $\|x - x^*\|^2 < 2\varepsilon/K$ .
- By proposition B'.1.4 for all  $x = Q(y)$ ,  $y \in U_\varepsilon^*$  it holds that  $\|x - x^*\|^2 < 2\varepsilon/K$ .
- Notice that from proposition B'.1.4  $Q(U_\varepsilon^*) \subseteq U_\varepsilon$  and  $Q^{-1}(U_\varepsilon) = U_\varepsilon^*$ .

Thus we conclude that whenever  $y \in U_\varepsilon^*$ ,  $x = Q(y) \in U_\varepsilon$ . Finally, by (Reciprocity)  $U_\varepsilon$  is a neighborhood of  $x^*$ . Since  $x^*$  is a strict Nash equilibrium, pick  $\varepsilon$  sufficiently small such that (VS) holds for all  $x \in U_{4\varepsilon}$ .

(Stability).

Assume now that  $Y_0 \in U_\varepsilon^*$  and thus  $F_h(x^*, Y_0) < \varepsilon \leq 4\varepsilon$ . We will prove by induction that  $Y_n \in U_{4\varepsilon}^*$  for all  $n \geq 1$  with probability at least  $1 - \delta$ . Suppose that  $F_h(x^*, Y_k) < 4\varepsilon$  for all  $1 \leq k \leq n$  and we will prove that  $Y_{n+1} \in U_{4\varepsilon}^*$  and consequently  $X_{n+1} \in U_{4\varepsilon}$ .

From proposition B'.1.4 we have

$$F_h(x^*, Y_{n+1}) \leq F_h(x^*, Y_n) + \gamma_n \langle \hat{v}_n, X_n - x^* \rangle + \frac{1}{2K} \gamma_n^2 \|\hat{v}_n\|_*^2 \quad (\text{B'.41})$$

For the payoff feedback, it holds  $\hat{v}_n = v(X_n) + Z_n + b_n$ . Then by telescoping the above inequality and substituting we get

$$\begin{aligned} F_h(x^*, Y_{n+1}) &\leq F_h(x^*, Y_0) + \sum_{k=0}^n \gamma_k \langle v(X_k), X_k - x^* \rangle + \sum_{k=0}^n \gamma_k \langle Z_k, X_k - x^* \rangle \\ &\quad + \sum_{k=0}^n \gamma_k \langle b_k, X_k - x^* \rangle + \frac{1}{2K} \sum_{k=0}^n \gamma_k^2 \|\hat{v}_k\|_*^2 \end{aligned} \quad (\text{B'.42})$$

We will study each term of the inequality separately.

- Let  $R_n = \frac{1}{2K} \sum_{k=0}^n \gamma_k^2 \|\hat{v}_k\|_*^2$  and  $F_{n,\varepsilon} = \{\sup_{0 \leq k \leq n} R_k \geq \varepsilon\}$ . As we discussed in lemma B'.4.1,  $R_n$  is a submartingale with  $\mathbb{E}[R_n] \leq \sum_{k=0}^n \gamma_k^2 M_k^2$ . Doob's maximal inequality (A'.1.5) yields

$$\mathbb{P}(F_{n,\varepsilon}) \leq \frac{\mathbb{E}[R_n]}{\varepsilon} \leq \frac{\sum_{k=0}^n \gamma_k^2 M_k^2}{2K\varepsilon} \quad (\text{B'.43})$$

By demanding  $\sum_{k=0}^{\infty} \gamma_k^2 M_k^2 \leq 2K\varepsilon\delta/3$  the event  $F_\varepsilon = \bigcup_{n=0}^{\infty} F_{\varepsilon,n}$  will occur with probability at most  $\delta/3$ .

- Let  $S_n = \sum_{k=0}^n \gamma_k \langle Z_k, X_k - x^* \rangle$  and  $E_{n,\varepsilon} = \{\sup_{0 \leq k \leq n} S_k \geq \varepsilon\}$ . Since  $S_n$  is a martingale, as we discussed in lemma B'.4.1, Doob's maximal inequality (A'.1.6) yields

$$\mathbb{P}(E_{n,\varepsilon}) \leq \frac{\mathbb{E}[S_n^2]}{\varepsilon^2} \leq \frac{\|\mathcal{X}\|^2 \sum_{k=0}^n \gamma_k^2 M_k^2}{\varepsilon^2} \quad (\text{B'.44})$$

In order to calculate the above upper bound, we define  $\psi_k = \langle Z_k, X_k - x^* \rangle$ . Notice that  $S_n^2 = \sum_{k=0}^n |\psi_k|^2 + 2 \sum_{k < \ell} \psi_k \psi_\ell$ . Indeed it holds that

$$\mathbb{E}[|\psi_k|^2] \leq \mathbb{E}[\mathbb{E}[\|Z_k\|_*^2 | X_k - x^*]^2 | \mathcal{F}_k]] \quad (\text{B'.45})$$

$$\leq \mathbb{E}[\mathbb{E}[\|Z_k\|_*^2 | \mathcal{F}_k]] \|\mathcal{X}\|^2 \quad (\text{B'.46})$$

where,

$$\mathbb{E}[\|Z_k\|_*^2 | \mathcal{F}_k] = \mathbb{E}[\|\hat{v}_k - \mathbb{E}[\hat{v}_k | \mathcal{F}_k]\|_*^2 | \mathcal{F}_k] \quad (\text{B'.47})$$

$$= \mathbb{E}[\|\hat{v}_k\|_*^2 - 2\langle \hat{v}_k, \mathbb{E}[\hat{v}_k | \mathcal{F}_k] \rangle + \|\mathbb{E}[\hat{v}_k | \mathcal{F}_k]\|_*^2 | \mathcal{F}_k] \quad (\text{B'.48})$$

$$= \mathbb{E}[\|\hat{v}_k\|_*^2 | \mathcal{F}_k] - \|\mathbb{E}[\hat{v}_k | \mathcal{F}_k]\|_*^2 \leq M_k^2 \quad (\text{B'.49})$$

$$\leq \mathbb{E}[\|\hat{v}_k\|_*^2 | \mathcal{F}_k] \leq M_k^2 \quad (\text{B'.50})$$

Furthermore, for all  $k \neq \ell$  it holds that  $\mathbb{E}[\psi_k \psi_\ell] = \mathbb{E}[\mathbb{E}[\psi_k \psi_\ell | \mathcal{F}_{k \vee \ell}]] = 0$ .

Thus, by demanding  $\sum_{k=0}^{\infty} \gamma_k^2 M_k^2 \leq \frac{\varepsilon^2 \delta}{3\|\mathcal{X}\|^2}$  we ensure that the event  $E_\varepsilon = \bigcup_{n=0}^{\infty} E_{\varepsilon,n}$  will occur with probability at most  $\delta/3$ .

- Let  $W_n = \sum_{k=0}^n \gamma_k \langle b_k, X_k - x^* \rangle$ ,  $J_n = \sum_{k=0}^n \gamma_k \|b_k\|_* \|\mathcal{X}\|$  as we discussed in lemma B'.4.1

$$W_n \leq J_n \quad (\text{B'.51})$$

where  $J_n$  is a submartingale with  $\mathbb{E}[J_n] \leq \|\mathcal{X}\| \sum_{k=0}^n \gamma_k B_k$ . Similarly to the previous steps let  $D_{\varepsilon,n} = \{\sup_{0 \leq k \leq n} J_k \geq \varepsilon\}$ , then Doob's maximal inequality (A'.1.5) yields

$$\mathbb{P}(D_{\varepsilon,n}) \leq \frac{\mathbb{E}[J_n]}{\varepsilon} \leq \frac{\|\mathcal{X}\| \sum_{k=0}^n \gamma_k B_k}{\varepsilon} \quad (\text{B'.52})$$

By demanding  $\sum_{k=0}^{\infty} \gamma_k B_k \leq \frac{\varepsilon \delta}{3\|\mathcal{X}\|}$  then the event  $D_\varepsilon = \bigcup_{n=0}^{\infty} D_{\varepsilon,n}$  will happen with probability at most  $\delta/3$ , which implies that with probability at most  $\delta/3$   $W_n$  will exceed  $\varepsilon$  for all  $n \geq 0$ .

- Furthermore, if  $X_k$  belongs to a neighborhood in which (VS) holds for all  $0 \leq k \leq n$ , we have

$$\langle v(X_k), X_k - x^* \rangle \leq 0 \text{ for all } n \geq 0 \quad (\text{B'.53})$$

By demanding the parameters of the algorithm to satisfy:

$$\sum_{k=0}^{\infty} \gamma_k^2 M_k^2 \leq \min \left\{ \frac{\varepsilon^2 \delta}{3 \|\mathcal{X}\|^2}, \frac{2K\varepsilon\delta}{3} \right\} \quad \& \quad \sum_{k=0}^{\infty} \gamma_k B_k \leq \frac{\varepsilon\delta}{3 \|\mathcal{X}\| \|\mathcal{Y}\|_*}$$

If all of  $\bar{E}_\varepsilon, \bar{F}_\varepsilon, \bar{D}_\varepsilon$  hold, this happens with probability  $\mathbb{P}(\bar{E}_\varepsilon \cap \bar{F}_\varepsilon \cap \bar{D}_\varepsilon) \geq 1 - \delta$  and from (B'.42) we have  $F_h(x^*, Y_{n+1}) < 4\varepsilon$ . This immediately yields that  $Y_{n+1} \in U_{4\varepsilon}^*$  and consequently as we explained in the begin of the proof  $X_{n+1} \in U_{4\varepsilon}$ , in which variational stability holds, with probability at least  $1 - \delta$ .

(Convergence).

By lemma B'.4.1 there exists a subsequence  $X_{n_k}$  that converges to  $x^*$ . By (Reciprocity) we have that  $\liminf_{n \rightarrow \infty} F_h(x^*, Y_n) = 0$ . In order to complete the proof, it is sufficient to prove that the limit of  $F_h(x^*, Y_n)$  exists. Notice that since the sequence of play remains in  $U_{4\varepsilon}$  variational stability holds and thus  $\langle v(X_n), X_n - x^* \rangle \leq 0$ . Again using proposition B'.1.4 we have:

$$F_h(x^*, Y_{n+1}) \leq F_h(x^*, Y_n) + \gamma_n \langle \hat{v}_n, X_n - x^* \rangle + \frac{1}{2K} \gamma_n^2 \|\hat{v}_n\|_*^2 \quad (\text{B'.54})$$

$$\mathbb{E}[F_h(x^*, Y_{n+1}) | \mathcal{F}_n] \leq F_h(x^*, Y_n) + \gamma_n \mathbb{E}[\langle b_n, X_n - x^* \rangle | \mathcal{F}_n] + \frac{1}{2K} \gamma_n^2 \mathbb{E}[\|\hat{v}_n\|_*^2 | \mathcal{F}_n] \quad (\text{B'.55})$$

$$\leq F_h(x^*, Y_n) + \gamma_n \mathbb{E}[\langle b_n, X_n - x^* \rangle | \mathcal{F}_n] + \frac{1}{2K} \gamma_n^2 M_n^2 \quad (\text{B'.56})$$

Notice that since from proposition B'.1.4  $F_h(x^*, Y) \geq 0$ , if we apply absolute values in the above inequality we have

$$\mathbb{E}[F_h(x^*, Y_{n+1}) | \mathcal{F}_n] = |\mathbb{E}[F_h(x^*, Y_{n+1}) | \mathcal{F}_n]| \quad (\text{B'.57})$$

$$\leq |F_h(x^*, Y_n)| + \gamma_n \mathbb{E}[|\langle b_n, X_n - x^* \rangle| | \mathcal{F}_n] + \frac{1}{2K} \gamma_n^2 M_n^2 \quad (\text{B'.58})$$

$$\leq F_h(x^*, Y_n) + \gamma_n \mathbb{E}[\|b_n\|_* | \mathcal{F}_n] \|\mathcal{X}\| + \frac{1}{2K} \gamma_n^2 M_n^2 \quad (\text{B'.59})$$

$$\leq F_h(x^*, Y_n) + \gamma_n B_n \|\mathcal{X}\| + \frac{1}{2K} \gamma_n^2 M_n^2 \quad (\text{B'.60})$$

Let

$$R_n = F_h(x^*, Y_n) + \|\mathcal{X}\| \sum_{k=n}^{\infty} \gamma_k B_k + \frac{1}{2K} \sum_{k=n}^{\infty} \gamma_k^2 M_k^2 \quad (\text{B'.61})$$

Then

$$\mathbb{E}[R_{n+1} | \mathcal{F}_n] \leq \mathbb{E}[F_h(x^*, Y_{n+1}) | \mathcal{F}_n] + \sum_{k=n+1}^{\infty} \gamma_k B_k \|\mathcal{X}\| + \frac{1}{2K} \sum_{k=n+1}^{\infty} \gamma_k^2 M_k^2 \quad (\text{B'.62})$$

$$\leq F_h(x^*, Y_n) + \sum_{k=n}^{\infty} \gamma_k B_k \|\mathcal{X}\| + \frac{1}{2K} \sum_{k=n}^{\infty} \gamma_k^2 M_k^2 \quad (\text{B'.63})$$

$$= R_n \quad (\text{B'.64})$$

Therefore  $R_n$  is a supermartingale and it is also  $L_1$  bounded (each one of the terms is bounded) and so from Doob's convergence theorem (A'.1.4)  $R_n$  converges to a finite random variable and so does  $F_h(x^*, Y_n)$ . Inevitably,  $\liminf_{n \rightarrow \infty} F_h(x^*, Y_n) = \lim_{n \rightarrow \infty} F_h(x^*, Y_n) = 0$  and by proposition B'.1.4,  $Q(Y_n) = X_n \rightarrow x^*$ .

The above analysis shows that whenever  $Y_0 \in U_\varepsilon^*$  and thus  $X_0 \in U_\varepsilon \cap \text{im } Q$ ,  $X_n \in U_{4\varepsilon} \cap \text{im } Q$  and converges to  $x^*$  with arbitrary high probability. Hence,  $x^*$  is stochastically asymptotically stable. ■

## B'.5 Proofs of Instability

Before moving on our proof we first provide some intuition derived from the notion of the polar cone (B'.1.3). Looking at the figure B'.2, the polar cone corresponding to *fully mixed* or *mixed* Nash equilibria has a key difference with the one corresponding to *strict* Nash equilibria. The latter, in contrast to the former, is fully dimensional. Thus intuitively, considering a sufficiently small neighborhood of a *mixed* Nash equilibrium, the slightest perturbation in the dual space of the payoffs, will lead to instability of the system. Our result is based on this intuition; we prove by contradiction that there exists a sufficiently small neighborhood of a *mixed* Nash equilibrium, from which the sequence of play will escape with strictly positive probability. The decomposability assumption of the regularizers ensures that the proof holds also for steep regularizers (See B'.1.2).

Below, leveraging the definition of the polar cone in simplex, we prove a useful property for the difference of the aggregated payoffs of FTRL for a sequence of play that shares common pure strategies.

**Lemma B'.5.1.** *Let  $X_i = Q(Y_i) \in \mathcal{X}_i$  be a mixed strategy profile and  $a, b \in \text{supp}(X_i)$  be two pure strategies, for some player  $i \in \mathcal{N}$ . Then it holds:*

$$\langle Y_i, e_a - e_b \rangle = \langle \nabla h_i(X_i), e_a - e_b \rangle$$

*Additionally, if (FTRL) is run then for a sequence of play  $X_{i,n_1}, \dots, X_{i,n_2}$  that maintains in its support both pure strategies  $a, b \in \mathcal{A}_i$  it holds*

$$\langle Y_{i,k_1} - Y_{i,k_2}, e_a - e_b \rangle = \langle \nabla h_i(X_{i,k_1}) - \nabla h_i(X_{i,k_2}), e_a - e_b \rangle \quad \forall k_1, k_2 \in \{n_1, \dots, n_2\}$$

*Proof.* From proposition B'.1.3,  $Y_i$  can be analyzed as  $Y_i = \nabla h_i(X_i) + G$ ,  $G \in \text{PC}(X_i)$ . Notice that  $\nabla h_i(X_i) = (\theta_i(X_{i,\alpha_1}), \dots, \theta_i(X_{i,\alpha_{|\mathcal{A}_i|}}))$ . Since  $X_i$  assigns positive probability to both  $a, b$ , by definition of the polar cone it is  $G_a = G_b$ . Thus,

$$\langle Y_i, e_a - e_b \rangle = G_a + \theta'_i(X_{i,a}) - G_b - \theta'_i(X_{i,b}) \tag{B'.65}$$

$$= \langle \nabla h_i(X_i), e_a - e_b \rangle \tag{B'.66}$$

For the second part, by applying (B'.66) for both cases of  $Y_{i,k_1}, Y_{i,k_2}$  we have:

$$\langle Y_{i,k_1}, e_a - e_b \rangle = \langle \nabla h_i(X_{i,k_1}), e_a - e_b \rangle \tag{B'.67}$$

$$\langle Y_{i,k_2}, e_a - e_b \rangle = \langle \nabla h_i(X_{i,k_2}), e_a - e_b \rangle \tag{B'.68}$$

From the subtraction of the above equations, we derive the desideratum:

$$\langle Y_{i,k_1} - Y_{i,k_2}, e_a - e_b \rangle = \langle \nabla h_i(X_{i,k_1}) - \nabla h_i(X_{i,k_2}), e_a - e_b \rangle \tag{B'.69}$$

■



**Theorem B'.5.2.** *Let  $x^*$  be a mixed Nash equilibrium. If (FTRL) is run with any feedback model that satisfies (A3), then  $x^*$  cannot be stochastically asymptotically stable for any choice of step-schedules.*

*Proof.* We start by determining all the parameters of the algorithm (FTRL) and we assume ad absurdum that  $x^*$  is a mixed Nash equilibrium, which is stochastically asymptotically stable. Then for all neighborhoods  $U$  of  $x^*$  and  $\delta > 0$ , there exists some neighborhood  $U_1$  such that whenever  $X_0 \in U_1$ , it holds that  $X_n \in U$  for all  $n \geq 0$  with probability at least  $1 - \delta$ . This equivalently implies that for all  $\varepsilon, \delta > 0$  if  $X_0 \in U_1$ ,  $\|X_n - x^*\| < \varepsilon$  for all  $n \geq 0$ , with probability at least  $1 - \delta$ . We leave  $\varepsilon$  to be chosen at the end of our analysis, but we will consider it to be fixed.

For each player  $i \in \mathcal{N}$  and round  $n$  if  $X_{i,n}, X_{i,n+1}$  are two consecutive instances of the sequence of play; then  $\|X_{i,n} - x_i^*\| < \varepsilon$ ,  $\|X_{i,n+1} - x_i^*\| < \varepsilon$  and by the triangle inequality

$$\|X_{i,n+1} - X_{i,n}\| < 2\varepsilon \quad (\text{B'.70})$$

We fix a round  $n$  and focus on player  $i \in \mathcal{N}$  who has the property of (A3); Since for two pure strategies  $a, b \in \text{supp}(x_i^*)$  of player  $i \in \mathcal{N}$ , holds that  $\mathbb{P}(|\hat{v}_{ia,n} - \hat{v}_{ib,n}| \geq \beta \mid \mathcal{F}_n) > 0$  for all  $n \geq 0$ , there exists for each round  $n \geq 0$ ,  $\pi_n > 0$  such that  $\mathbb{P}(|\hat{v}_{ia,n} - \hat{v}_{ib,n}| \geq \beta \mid \mathcal{F}_n) = \pi_n$ . Choose  $\delta$  such that  $\delta < \pi_n$  and consequently

$$1 - \delta > 1 - \pi_n \quad (\text{B'.71})$$

This is possible, since  $\pi_n$  is strictly positive and  $\delta$  can be chosen arbitrarily small.

Consider now the projection of the aggregate payoffs  $Y_{i,n}, Y_{i,n+1}$  in the difference of the directions of these two strategies. From lemma B'.5.1 we have

$$\langle Y_{i,n+1} - Y_{i,n}, e_a - e_b \rangle = \langle \nabla h_i(X_{i,n+1}) - \nabla h_i(X_{i,n}), e_a - e_b \rangle \quad (\text{B'.72})$$

However, by definition of (FTRL)  $Y_{i,n+1} - Y_{i,n} = \gamma_n \hat{v}_{i,n}$  and by taking into consideration that the regularizers used are decomposable, we get

$$(\theta'_i(X_{ia,n+1}) - \theta'_i(X_{ib,n+1}) - (\theta'_i(X_{ia,n}) - \theta'_i(X_{ib,n}))) = \gamma_n \langle \hat{v}_{i,n}, e_a - e_b \rangle \quad (\text{B'.73})$$

By rearranging we have

$$(\theta'_i(X_{ia,n+1}) - \theta'_i(X_{ia,n})) - (\theta'_i(X_{ib,n+1}) - \theta'_i(X_{ib,n})) = \gamma_n (\hat{v}_{ia,n} - \hat{v}_{ib,n}) \quad (\text{B'.74})$$

As a consequence of  $\theta_i$  being continuously differentiable in all of  $(0, 1]$ ,  $\theta'_i$  is continuous in  $[L(\varepsilon), 1]$ , where  $L(\varepsilon)$  is the lower bound of  $X_{ia}, X_{ib}$  whenever  $\|X_i - x_i^*\| < \varepsilon$ .  $L(\varepsilon)$  can be guaranteed to be positive for a sufficiently small  $\varepsilon < \varepsilon'$ , which ensures that all the points of the neighborhood contain the support of the equilibrium for player  $i$ . Therefore, from extreme value theorem in  $\theta'_i$ , there exist finite  $C_a, C_b$  corresponding to  $a, b$  equivalently, such that

$$|\theta'_i(X_{ia,n+1}) - \theta'_i(X_{ia,n})| \leq C_a |X_{ia,n+1} - X_{ia,n}| < 2 \cdot C_a \cdot \varepsilon \quad (\text{B'.75})$$

$$|\theta'_i(X_{ib,n+1}) - \theta'_i(X_{ib,n})| \leq C_b |X_{ib,n+1} - X_{ib,n}| < 2 \cdot C_b \cdot \varepsilon \quad (\text{B'.76})$$

By applying the triangle inequality in (B'.74) and using (B'.75),(B'.76) we get

$$\gamma_n |\hat{v}_{ia,n} - \hat{v}_{ib,n}| < (2 \cdot C_a + 2 \cdot C_b) \cdot \varepsilon \quad (\text{B'.77})$$

Equivalently,

$$|\hat{v}_{ia,n} - \hat{v}_{ib,n}| < \frac{2 \cdot C_a + 2 \cdot C_b}{\gamma_n} \cdot \varepsilon \quad (\text{B'.78})$$

The above inequality holds with probability  $1 - \delta$ . Thus, if the sequence of play  $X_n$  is contained to an  $\varepsilon$ -neighborhood of  $x^*$  i.e.,  $\|X_n - x^*\| < \varepsilon$  for all  $n \geq 0$ , then the difference of the feedback, for some player  $i \in \mathcal{N}$ , to two strategies of the equilibrium is  $O(\varepsilon/\gamma_n)$  with probability at least  $1 - \delta$ .

We now fix  $\varepsilon$  to be

$$\varepsilon < \min \left\{ \varepsilon', \frac{\gamma_n}{2 \cdot C_a + 2 \cdot C_b} \beta \right\} \quad (\text{B'.79})$$

and consequently

$$\mathbb{P}(|\hat{v}_{ia,n} - \hat{v}_{ib,n}| < \beta \mid \mathcal{F}_n) \geq 1 - \delta \quad (\text{B'.80})$$

However, from assumption (A3), it holds that

$$\mathbb{P}(|\hat{v}_{ia,n} - \hat{v}_{ib,n}| \geq \beta) \geq \pi_n \quad (\text{B'.81})$$

Combining (B'.80),(B'.81) we conclude

$$1 = \mathbb{P}[\{|\hat{v}_{ia,n} - \hat{v}_{ib,n}| \geq \beta\} \cup \{|\hat{v}_{ia,n} - \hat{v}_{ib,n}| < \beta\}] \quad (\text{B'.82})$$

$$= \mathbb{P}[|\hat{v}_{ia,n} - \hat{v}_{ib,n}| \geq \beta] + \mathbb{P}[|\hat{v}_{ia,n} - \hat{v}_{ib,n}| < \beta] \quad (\text{B'.83})$$

$$\geq \pi_n + 1 - \delta \quad (\text{B'.84})$$

$$> 1 \quad (\text{B'.85})$$

which is a contradiction.

Thus, a mixed Nash equilibrium cannot be stochastically asymptotically stable, under (FTRL) for types of payoff feedback described in section 4.1.1. Notice that this analysis holds even for the first round. Once the parameters of the algorithm have been determined, asymptotic instability can be derived in whichever finite round. ■