



POLYTECH[®]
NICE-SOPHIA

Classifieur de Bayes et approximation softmax

ANDRIEU Grégoire
GILLE Cyprien
NEGRE Philippe
YOUSSEF Alan



I. Introduction

Sujet :

L'approximation **softmax**, utilisée à la place de la fonction **argmax**, dans un **classifieur de Bayes**

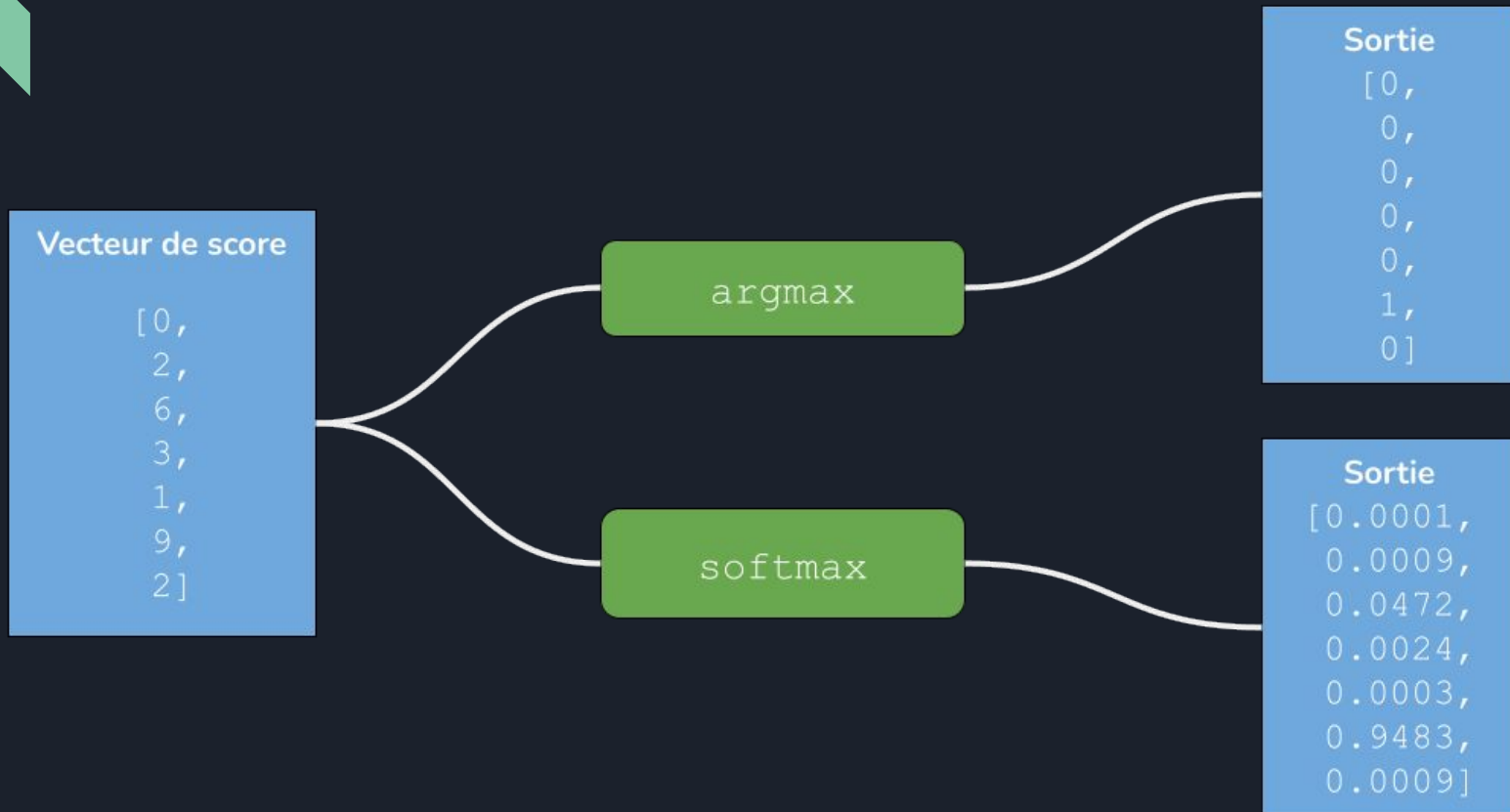
II. Présentation des objets

- Nos fonctions d'étude :

$$\mathbb{R}^n \rightarrow \{0, 1\}^n$$
$$\text{argmax} : V = [v_i] \mapsto W = [w_i] \text{ avec } \begin{cases} w_i = 1 & \text{si } v_i = \max_i(V) \\ 0 & \text{sinon} \end{cases}$$

$$\mathbb{R}^n \rightarrow [0, 1]^n$$
$$\text{softmax} : V = [v_i] \mapsto W = [w_i] \text{ avec } w_i = \frac{\exp(v_i)}{\sum_{k=1}^n \exp(v_k)}$$

II. Présentation des objets





II. Présentation des objets

- Classification

Associer à une observation X une classe Y :

$$\mathcal{D}(X) = Y_i$$

- Règle de Bayes

$$\hat{y} = \underset{k \in \{1, \dots, K\}}{\operatorname{argmax}} p(Y = k) \prod_{i=1}^n p(X_i \mid Y = k)$$

III. Etude d'un cas particulier

3 Classes :

- Excellente
- Ok
- Honteuse



4 features :

- Degré de cuisson
- Nombre de pepperonis
- Quantité de fromage
- Rayon

III. Etude d'un cas particulier

- Répartition dans les classes :

P(Y)	0 : Excellente	1 : Ok	2 : Honteuse
	0.30	0.50	0.20

- Exemple d'une feature : le degré de cuisson (1 = pas cuite, 5 = en cendre)

DdC selon acceptabilité	1	2	3	4	5
Excellente	0.0	0.1	0.8	0.1	0.0
Ok	0.0	0.3	0.4	0.3	0.0
Honteuse	0.4	0.1	0.0	0.1	0.4



III. Etude d'un cas particulier

```
print("\nClassification d'une pizza la plus mauvaise possible")
mauvaise_pizza = [1, 0, 1, 10]
pizza_bc.classify(mauvaise_pizza)
pizza_bc.bayes_softmax(mauvaise_pizza)
```

```
Classification d'une pizza la plus mauvaise possible
Classe: 2
Classe 0. Confiance: 0.2119
Classe 1. Confiance: 0.2119
Classe 2. Confiance: 0.5761
[Finished in 0.7s]
```


III. Le risque d'erreur global

- Expressions matricielles:

Notons N le nombre d'observations (i.e. de données), K restant le nombre de classes. Notre matrice de probabilités A' est alors de dimension (K, N) :

$$A' = \begin{pmatrix} P(Y = 1 | X_1) & \dots & P(Y = 1 | X_N) \\ \vdots & \ddots & \vdots \\ P(Y = K | X_1) & \dots & P(Y = K | X_N) \end{pmatrix}$$

On applique le théorème de Bayes, et on obtient la matrice A des probabilités postérieures, où apparaît la première dépendance en π .

$$A = \begin{pmatrix} P(X_1 | Y = 1)\pi_1 & \dots & P(X_N | Y = 1)\pi_1 \\ \vdots & \ddots & \vdots \\ P(X_1 | Y = K)\pi_k & \dots & P(X_N | Y = K)\pi_k \end{pmatrix}$$

III. Le risque d'erreur global

- Expressions matricielles:
- Argmax:
- Softmax:

$$D_j^* = \underset{i}{\operatorname{argmax}}[A_j]$$

$$d_{ij}^* = \left(\underset{k}{\operatorname{argmax}} P(X_j | Y = k) \pi_k \right)_i$$

$$D^* = \underset{D \in \mathcal{D}}{\operatorname{argmin}} R(D)$$

$$\tilde{D}_j = \operatorname{softmax}[A_j]$$

$$\tilde{d}_{ij} = \frac{\exp(\pi_i P(X_j | Y = i))}{\sum_{k=1}^K \exp(\pi_k P(X_j | Y = k))}$$

III. Le risque d'erreur global

- L'expression du risque :

$$R(D, \pi) = \sum_{i \in \{1, \dots, K\}} \pi_i R_i(D)$$

- Risque conditionnel:

$$\begin{aligned} R_i(D) &= \sum_j \sum_{k \neq i} P(Y = k | X_j) \\ &= \sum_j C_j^i D \end{aligned}$$

Avec C^i un vecteur $(1, K)$ tel que :

$$C_j^i = \begin{cases} 1 & \text{si } j \neq i \\ 0 & \text{sinon} \end{cases}$$

III. Le risque d'erreur global

- Expressions du risque:

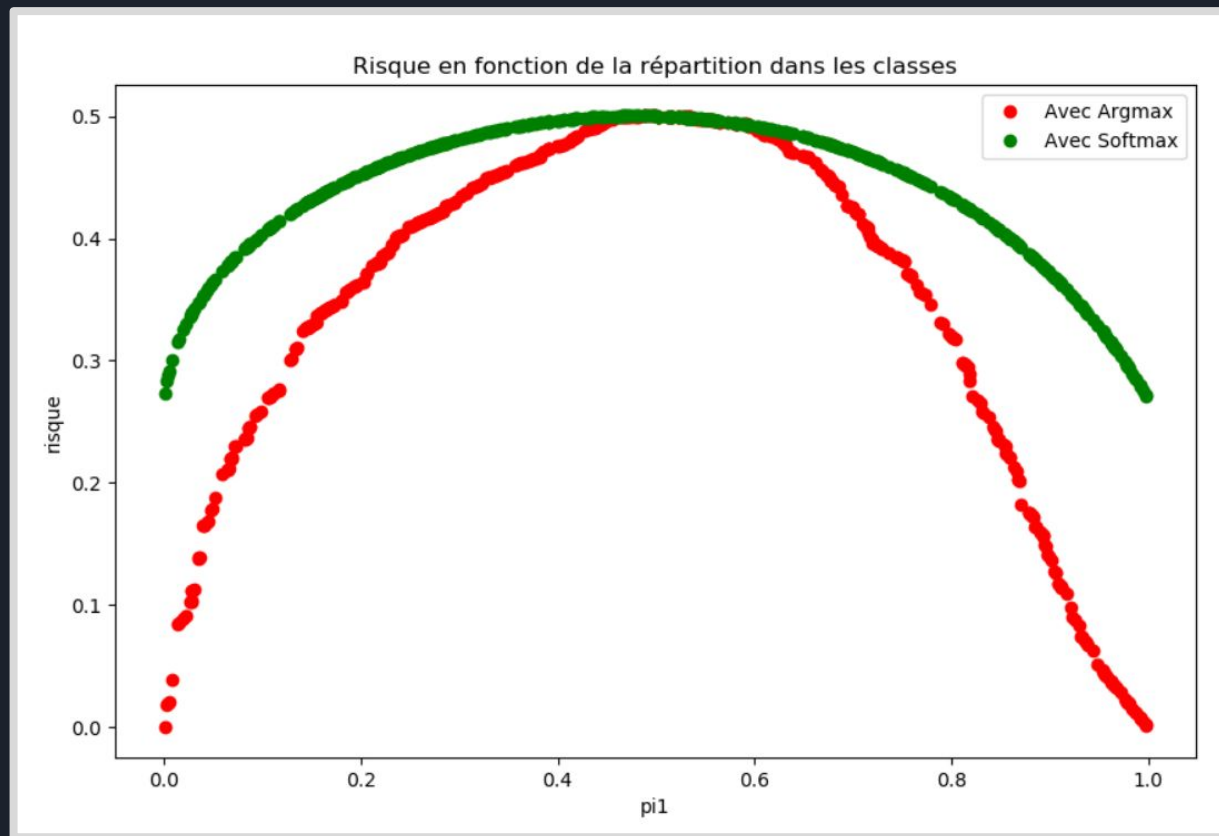
- Argmax:

$$\begin{aligned} R(D^*, \pi) &= \sum_{i \in \{1, \dots, K\}} \pi_i R_i(D^*) \\ &= \sum_i \sum_j \pi_i C^i D^* \end{aligned}$$

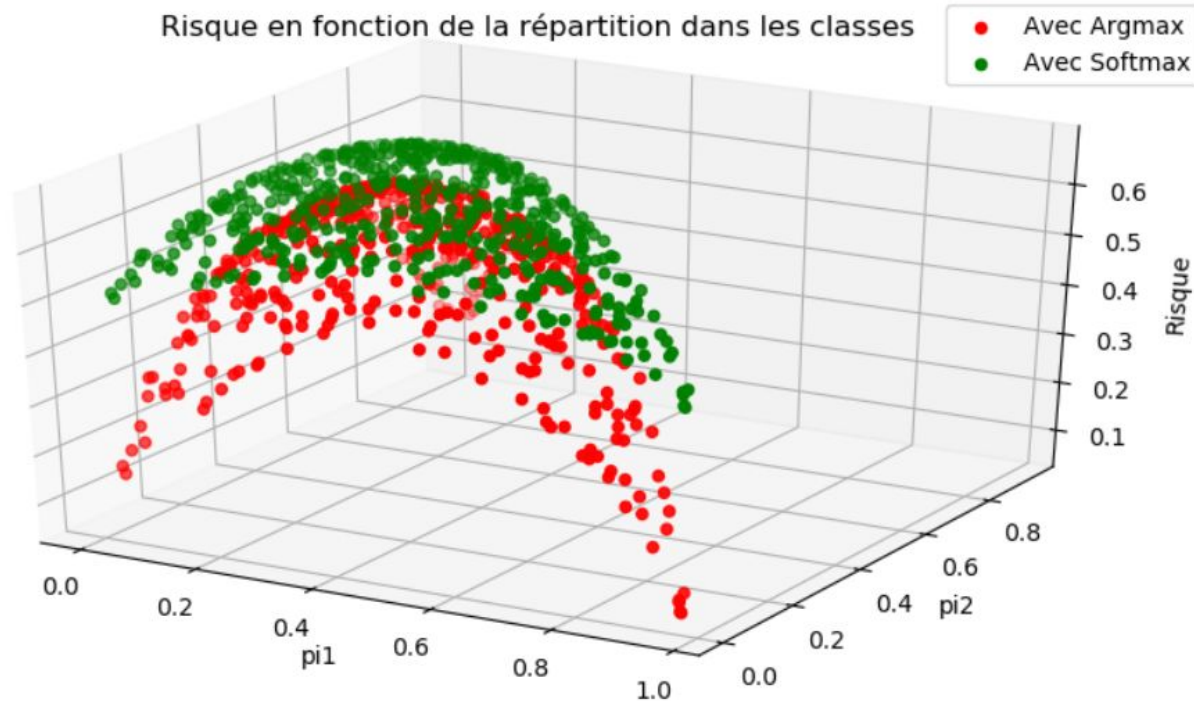
- Softmax:

$$R(\tilde{D}, \pi) = \sum_{i=1}^K \sum_{j=1}^N \pi_i (1 - \delta_{ij}) \frac{\exp(\pi_i P(X_j | Y = i))}{\sum_{k=1}^K \exp(\pi_k P(X_j | Y = k))}$$

III. Le risque d'erreur global



III. Le risque d'erreur global





III. Le risque d'erreur global

- Concavité
 - Risque Argmax
 - Risque Softmax
- Gradient

$$\frac{\partial R(\tilde{D}, \pi)}{\partial \pi_i} = \sum_{j=1}^N \left[(1 - \delta_{ij}) (S_i + \pi_i S_i (1 - S_i)) + \sum_{k \neq i}^K (1 - \delta_{kj}) (-S_k S_i) \right]$$

$$S_i = \frac{\exp(a_i)}{\sum_{k=1}^K \exp(a_k)}$$

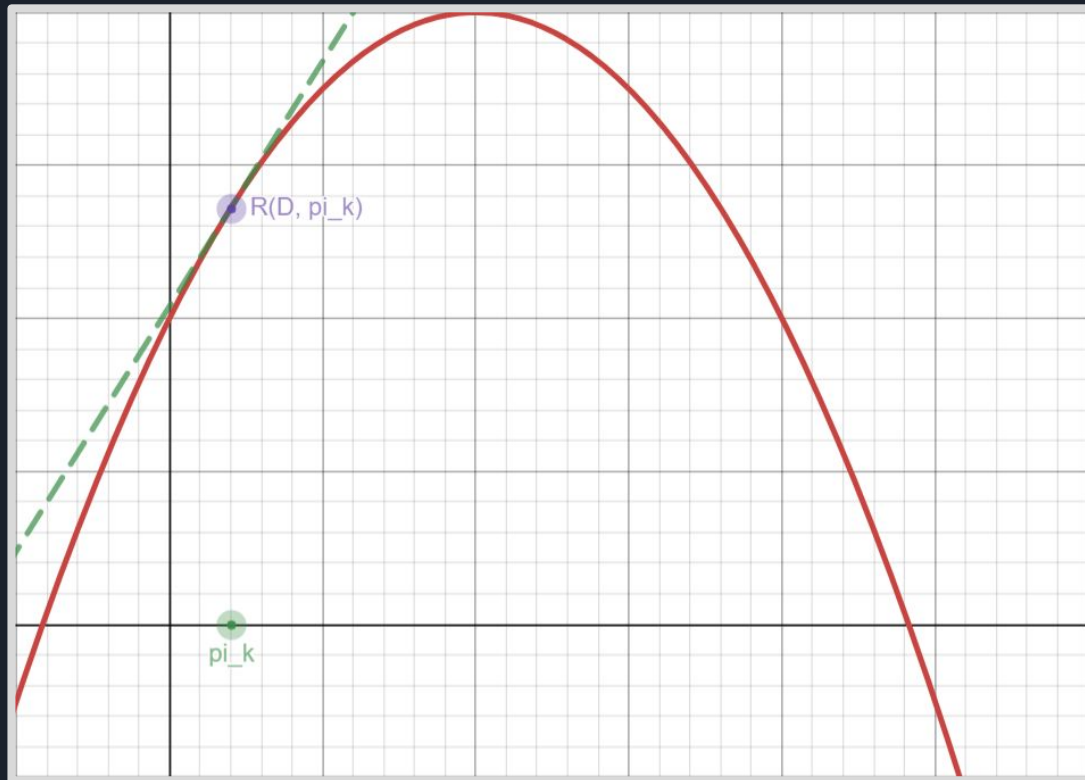


IV. Maximisation du risque : méthode du gradient projeté

- Montée du gradient :
 - Paramètres :
 - Nombre d'étapes
 - π_0 (Point de départ)
 - η (learning rate, ou facteur multiplicatif du pas) .

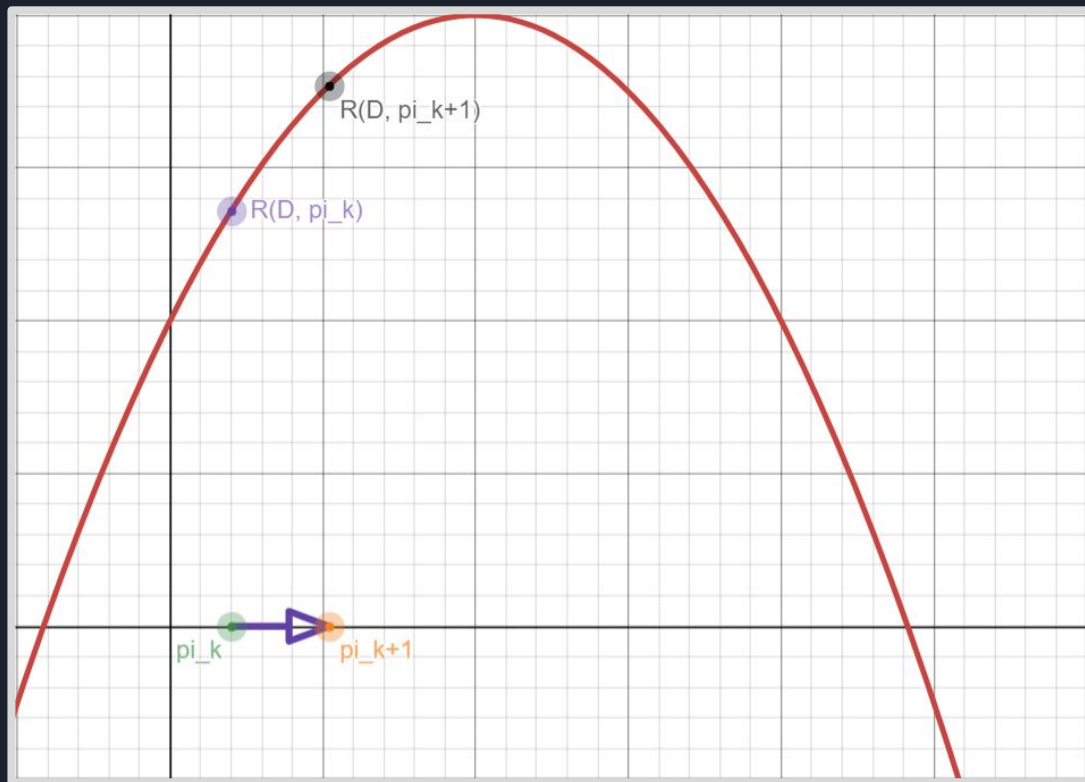
IV. Maximisation du risque : méthode du gradient projeté

- Déroulé d'une étape:
 1. Evaluation du gradient en π_k



IV. Maximisation du risque : méthode du gradient projeté

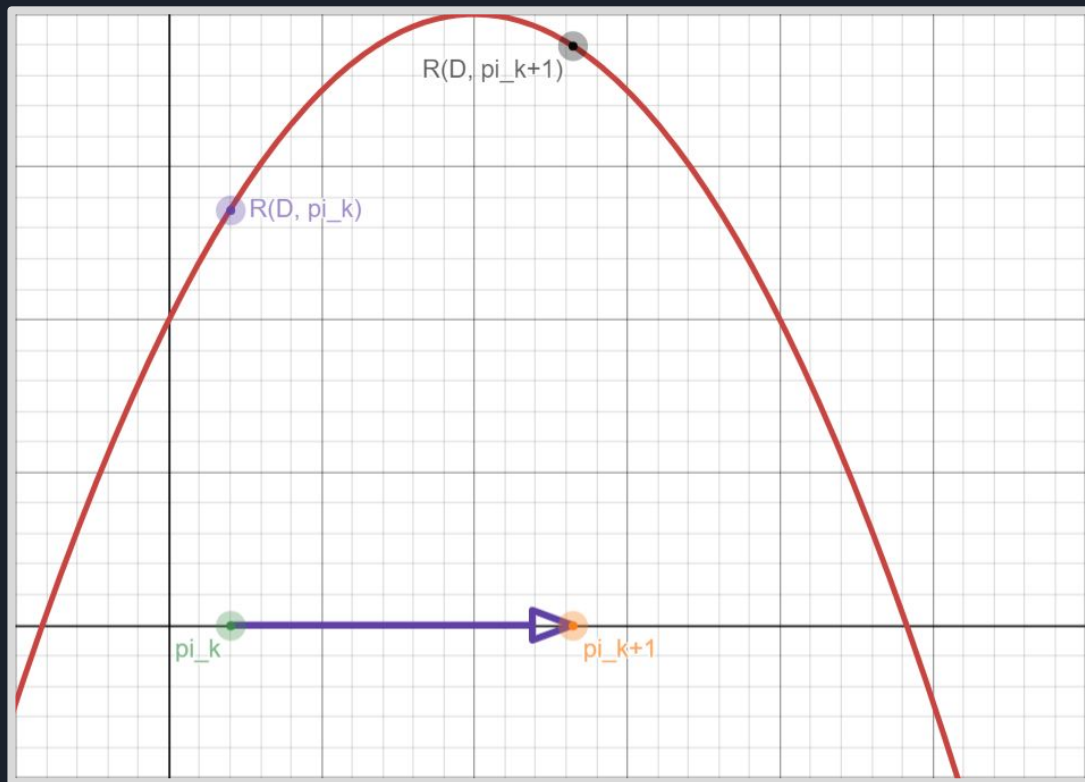
- Déroulé d'une étape:
 1. Evaluation du gradient en π_k
 2. $\pi_{k+1} = \pi_k + \eta * \text{gradient}$



IV. Maximisation du risque : méthode du gradient projeté

- Déroulé d'une étape:
 1. Evaluation du gradient en π_k
 2. $\pi_{k+1} = \pi_k + \eta * \text{gradient}$

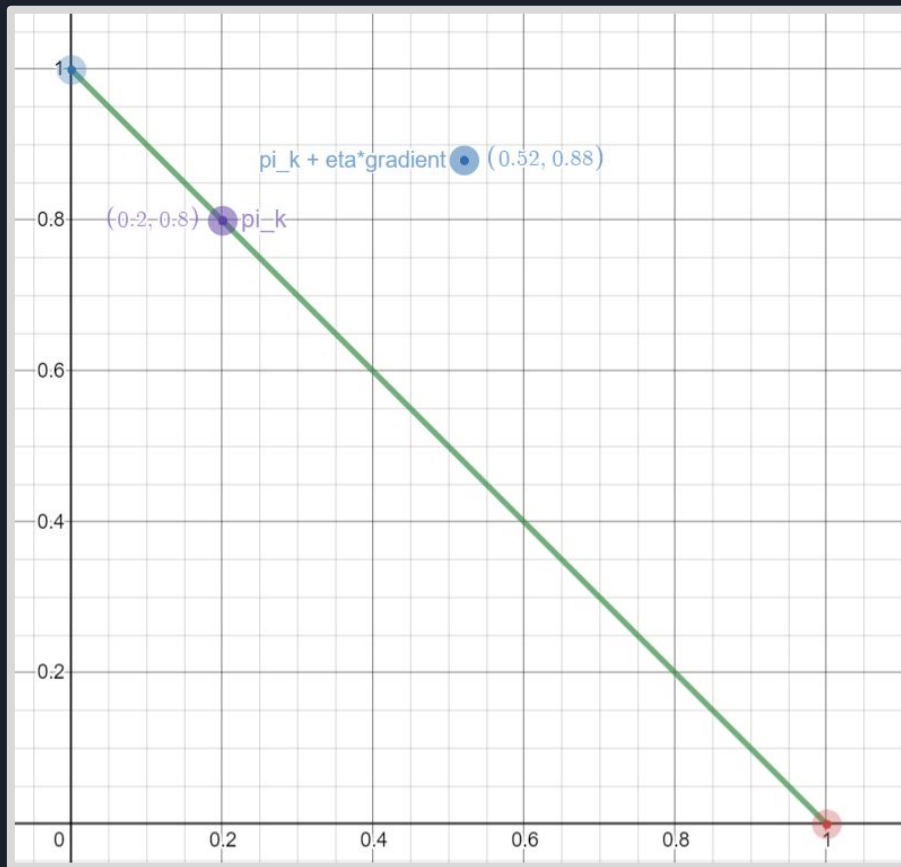
Influence du choix de η :



IV. Maximisation du risque : méthode du gradient projeté

- Problème:

Perte des propriétés
probabilistes!



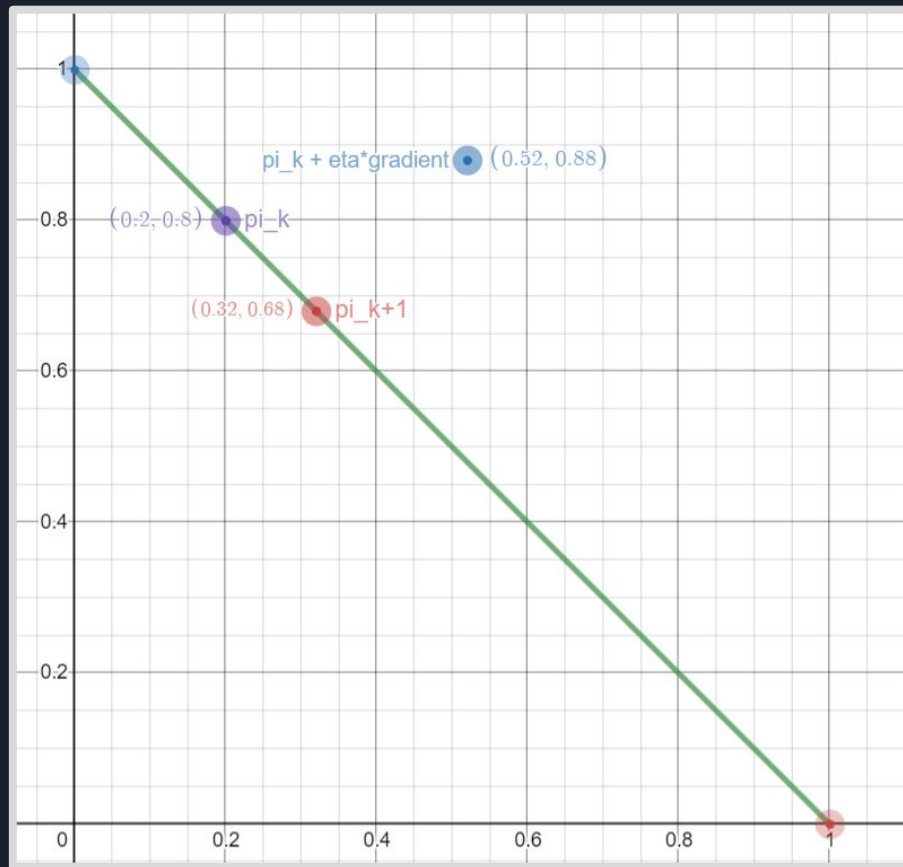
IV. Maximisation du risque : méthode du gradient projeté

- Problème:

Perte des propriétés probabilistes!

- Solution:

Projection sur le simplex



IV. Maximisation du risque : méthode du gradient projeté

- Principe de la projection sur le simplex:

Le point du K-simplex le plus proche de π_i est:

$$t_i = \max_i \{ \pi_i - \Delta_i, 0 \}$$

Avec :

$$\Delta_i = \frac{\left(\sum_{j=i+1}^K \pi_j \right) - 1}{K - i}$$

Avec les π_i triés par ordre croissant.

IV. Maximisation du risque : méthode du gradient projeté

```
Etape: 0 | [0.1, 0.2, 0.7]
Etape: 20 | [0.24258121 0.27990927 0.47750952]
Etape: 40 | [0.29918982 0.31194929 0.38886089]
Etape: 60 | [0.32166766 0.32480849 0.35352385]
Etape: 80 | [0.33059316 0.32997157 0.33943526]
Etape: 100 | [0.33413718 0.33204496 0.33381786]
[Finished in 0.7s]
```



Convergence vers la situation attendue : l'équiprobabilité à priori



V. Pistes futures

- Différence entre les risques argmax et softmax (méthode alpha)
- Influence du choix argmax/softmax sur l'entraînement d'un classifieur non-idéal ?
- Je manque d'idées aled



Bibliographie

Classification bayésienne et risque, expressions et calculs:

[ECE531 Lecture 2a: A Mathematical Model for Hypothesis Testing \(wpi.edu\)](#)

[ECE531 Lecture 2b: Bayesian Hypothesis Testing \(wpi.edu\)](#)

[ECE531 Lecture 3: Minimax Hypothesis Testing \(wpi.edu\)](#)

Projection sur le simplex:

http://www.gipsa-lab.fr/~laurent.condat/download/proj_simplex_l1ball.m

<https://arxiv.org/pdf/1101.6081.pdf>

Concavité de la fonction softmax:

<https://arxiv.org/pdf/1502.04635.pdf>