

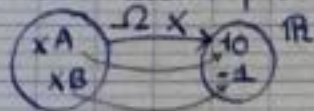
TD: ML

TD1 Exercice 1

- 1) 1) événement élémentaire = résultat de l'expérience
 univers = ensemble des événements élémentaires

- 2) $P: \mathcal{P}(\Omega)$ ensemble des parties de $\Omega \rightarrow [0, 1]$
 • $P(\Omega) = 1$ et $P(\emptyset) = 0$
 • $P(A) \geq 0$
 • $A \cap B = \emptyset$ alors $P(A \cup B) = P(A) + P(B)$

- 3) variable aléatoire : permet d'attribuer un résultat à un événement élémentaire



$$X: \Omega \rightarrow \mathbb{R}$$

$$X(\omega) = 10$$

$$P(X=10) = P(X^{-1}(\{10\}))$$

$$P(X \in [a, b]) = P(X^{-1}([a, b]))$$

→ passerelle entre év et valeurs

- 2) loi géométrique : $p^b(1-p)$

- 3) Les événements sont indépendants car $P(\text{Roi}) = 4/52$ et $P(\text{Pique}) = 13/52$ et $P(\text{Roi de Pique}) = 1/52$
 Ils ne sont pas incompatibles car on peut tirer le roi de pique.
 $P(\text{Roi} \cup \text{Pique}) = 4/52 + 13/52 - 1/52 = 16/52$

Exercice 2

- 1) \mathbb{R}^n univers est continu (les événements sont infinis et non dénombrables)
 1) Comment évaluer la proba autour d'un point de Ω sachant que cette proba = 0 ?
 ⇒ densité de probabilité = $\int_a^b p(x) dx = P_X([a, b])$
 densité jointe comme dans le cas discret
 densité marginale : $P(X) = \int_{-\infty}^{+\infty} p(x, y) dy$

2) $E(X) = \int_{\mathbb{R}} x p(x) dx = \int_{\Omega} X(\omega) p(\omega) d\mu$

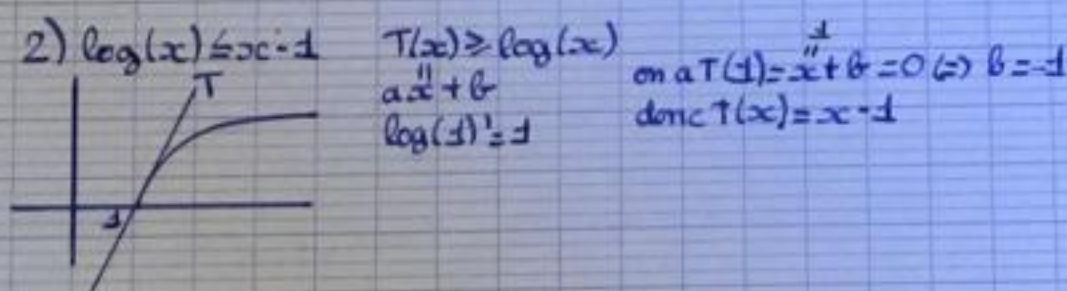
3) $E(X|Y=y) = \int_{\mathbb{R}} x p(x|Y=y) dx$

4) $E(X) \leftarrow \frac{1}{|E|} \sum_i x_i$
 $|E| \rightarrow +\infty$

iid: échantillon les uns des autres et choisis dans le même ensemble
 → loi des grands nombres

Exercice 3

$$2) 1) H(p) \geq 0 \text{ car } p_i \in]0; 1[\Rightarrow \log p_i \leq 0 \Rightarrow -\log p_i \geq 0 \Rightarrow -\sum p_i \log p_i \geq 0 = H(p)$$



$$\begin{aligned} -\sum p_i \log(p_i) &\leq -\sum p_i \log(q_i) \\ \Leftrightarrow \sum p_i \log(p_i) - \sum p_i \log(q_i) &= \sum p_i \log\left(\frac{p_i}{q_i}\right) \geq 0 \\ \Leftrightarrow \sum p_i \log\left(\frac{q_i}{p_i}\right) &\leq 0 \\ \sum p_i \log\left(\frac{q_i}{p_i}\right) &\leq \sum p_i \left(\frac{q_i}{p_i} - 1\right) \quad \text{car } \log(x) \leq x-1 \\ &\leq \sum q_i - \sum p_i = 0 \quad \text{ici } x = \frac{q_i}{p_i} \end{aligned}$$

$$H(p) \leq -\sum p_i \log q_i$$

$$3) \text{ avec } q_i = 1/m, -\sum p_i \log(1/m) = \sum p_i \log(m) = \log(m)$$

3) 1) égalité montroale par récurrence en faisant un changement de variable

$$2) f(\mathbb{E}(X)) \leq \mathbb{E}(f(X))$$

$\Rightarrow \sum_{i=1}^N x_i p_i \dots$

$$3) KL(p_i \| q_i) = \sum p_i \log \frac{p_i}{q_i} = -\sum p_i \log \frac{q_i}{p_i} \geq -\log\left(\sum p_i \frac{q_i}{p_i}\right) \geq -\log\left(\sum q_i\right) = 0$$

Ce n'est pas symétrique \Rightarrow on peut le transformer en distance en prenant la moyenne entre $p \| q$ et $q \| p$

estimation ponctuelle

$$4) \hat{p}(x) = \frac{1}{N} \sum_{i=1}^N \delta(x - x_i) \text{ on essaie de l'approcher par } p_\theta$$

$$KL(\hat{p} \| p_\theta) = \sum_{x \in \mathcal{X}} \hat{p}(x) \log \frac{\hat{p}(x)}{p_\theta(x)} = \sum_x \hat{p} \log \hat{p}(x) - \sum_x \hat{p} \log p_\theta(x)$$

$\underbrace{\sum_x \hat{p} \log \hat{p}(x)}_{=H(\hat{p})}$

$$= -\sum \hat{p} \log p_\theta(x) = -\frac{1}{N} \sum_{i=1}^N \sum_{x \in \mathcal{X}} \delta(x - x_i) \log p_\theta(x)$$

$$= \frac{1}{N} \sum \log p_\theta(x_i)$$

1) 1) 3^{46} valeurs différentes sont possibles

2) $\frac{P(D|V=v)}{P(R|V=v)}$ on ne peut pas calculer ça car on a pas tous les votes étiquetés (on a juste les votes agrégés par partie et pas les votes individuels pour chaque voteur)

\Rightarrow on utilise Bayes

$$\frac{P(V=v|D)P(D)P(V=v)}{P(V=v|R)P(R)P(V=v)} = \frac{P(V=v|D)P(D)}{P(V=v|R)P(R)} \xrightarrow{\frac{|D|}{|D|+|R|}} \xrightarrow{\frac{|R|}{|D|+|R|}}$$

On ne peut pas estimer $P(V=v|D)$ car on a pas de votes individuels. On est obligé de faire l'hypothèse d'indépendance.

Rmq: ici l'hypothèse d'ind n'est pas vraiment juste car les votes dépendent les uns des autres. Mais on a pas le choix de faire autrement

$$P(V=v|D) = \prod_{i=1}^n P(v_i|D) \quad \text{on passe au log}$$

$$\log \frac{P(D|V)}{P(R|V)} = \sum \log P(v_i|D) - \sum \log P(v_i|R) + \log \left(\frac{|D|}{|R|} \right)$$

$$P(v_i=0|D) = 1 - P(v_i=1|D) - P(v_i=NS|D)$$

2) 1) $G \sim \text{Ber} \left(\frac{|R|}{|D|+|R|} \right)$

2) $P(t) = \underbrace{p(t|G=h)}_{\mathcal{N}(\mu_h, \sigma_h)} p(G=h) + \underbrace{p(t|G=f)}_{\mathcal{N}(\mu_f, \sigma_f)} p(G=f) \Rightarrow$ mixture gaussienne

$$P(G=h|t) = \frac{p(t|G=h)p(G=h)}{P(t)}$$

3) classifieur bayésien = $\frac{P(G=h|t)}{P(G=f|t)} = \frac{p(t|G=h)p(h)}{p(t|G=f)p(f)}$

3) 5 paramètres : 2 par gaussienne (esp + variance)
1 param pour le rapport h/f

Exercice 4

1) 1) $p(x_1, \dots, x_m) = \prod_{i=1}^m P(x_i) = \prod_{i=1}^m p_i^{m_i}$
car indépendance

2) $m \rightarrow mp_k$ (loi des grands nb) donc $\prod_{i=1}^6 p_i^{m_i} \rightarrow \prod_{i=1}^6 p_i^{mp_i}$
 $= \prod_{i=1}^6 2^{m p_i \log p_i} = 2^{m \sum_{i=1}^6 p_i \log p_i} = 2^{-m H(p)}$

faible entropie = le nb de séquences est très diversifié donc d'aléatoire
forte entropie = on tire presque toujours la même séquence

Td2 Exercice 1

$p(y)$: a priori
 $p(x)$: évidence
 $p(x|y)$: vraisemblance
 $p(y|x)$: a posteriori

- 1) On se positionne dans un monde idéal.
 X dans \mathbb{R}^d et $Y = \{y_1, \dots, y_k\}$ $X = \{x_1, \dots, x_n\}$ $x_i \in \mathbb{R}^d$
 Classifieur $X \rightarrow Y$ qui utilise des algo fondés sur Bayes
 $\hat{y} = \arg \max_y p(y|x)$

- 2) Erreur: Probab de se tromper = $P(\hat{y} \neq y|x) = 1 - p(\hat{y} = y|x) = 1 - \max_y p(y|x)$ (ou argmax)
* vérité d'entraînement * spécialisation du théorème

On ne peut pas trouver mieux. Pour le prouver on utilise un raisonnement par l'absurde et on suppose qu'il existe un classifieur \hat{y}_0 tq $E_{Bay} > E_0$
 \Rightarrow contradiction. $1 - \max_y p(y|x)$ $\Rightarrow 1 - p(y = \hat{y}_0|x)$

- 3) $\lambda(y_i, y_j)$ coût de prédiction de y_i plutôt que y_j

Dans le cas 0-1: $\lambda(y_i, y_j) = \mathbb{1}_{\{y_i \neq y_j\}} = \delta(y_i \neq y_j) = \begin{cases} 1 & \text{si } y_i \neq y_j \\ 0 & \text{sinon} \end{cases}$

Dans le cas asymétrique: \rightarrow cas des classes "proches" (non binaire)

ex: chien plus proche de chat que camion $\lambda(\text{chien}, \text{chat}) < \lambda(\text{chien}, \text{camion})$

\rightarrow cas binaire

ex: texte médical/fake $\lambda(+, -) > \lambda(-, +)$

- 4) $R(y_i|x)$ en fonction de λ et $P(y|x)$

$$R(y_i|x) = \mathbb{E}_{y \sim p(y|x)} [\lambda(y_i, y)] = \sum_{y \in Y} p(y|x) \lambda(y_i, y)$$

espérance du coût

- 5) $D = \{(x_i, y_i)\}_{i=1}^m$ $x_i \in X$ $y_i \in Y$

$$R(f) = \sum_{x, y \in X \times Y} \lambda(f(x), y) p(x, y) \text{ avec } dy = \mathbb{E}_{x \sim p} [\lambda(f(x), y)]$$

coût moyen distribution de proba des données

Dans la pratique on suppose:

$\rightarrow D$ est un échantillon iid suivant D la distribution

\rightarrow loi des grands nombres $Z = \lambda(f(x), y)$

$$E(Z) \approx \frac{1}{m} \sum_{i=1}^m Z_i = \frac{1}{m} \sum_{i=1}^m \lambda(f(x_i), y_i) = R_D(f)$$

statistique empirique

- 6) $\lambda_{+-} \triangleq \lambda(x, y)$

quel critère de décision faut-il prendre?

\rightarrow on veut minimiser l'erreur

$$E(\hat{y}) = \lambda(\hat{y}, +) p(y = +|x) + \lambda(\hat{y}, -) p(y = -|x)$$

Décision: on prédit plus de $E(+)$ $E(+)$ $E(-)$

$$\lambda_{++} p(y = +|x) + \lambda_{+-} p(y = -|x) < \lambda_{-+} p(y = +|x) + \lambda_{--} p(y = -|x)$$

But: arriver à un critère du type $p(y = +|x) > \text{seuil}$ (vérifier que $\lambda_{+-} > \lambda_{--}$) cas 0-1

\rightarrow on va supposer que $\lambda_{+-} > \lambda_{--}$ et $\lambda_{-+} > \lambda_{++}$

5) x_i' est le plus proche voisin de x parmi $\{x_i\}_{i=1}^m$
 $\mathbb{P}(x_i' \in B) \xrightarrow{m \rightarrow \infty} 1$
 $\mathbb{P}(x_i' \in B) \xrightarrow{m \rightarrow \infty} 1$
 Hyp: $p(x) > 0$ (densité)
 $\rightarrow \exists \delta p(x) > 0 \forall B > 0$ et $\mathbb{P}(X_{m,i} \in B) > 0$ pour m'
 $\rightarrow \mathbb{P}(x_i' \in B) > 0$ pour $i \geq m'$

$$6) q_B(x) = P(y=B|x) \\ \pi(x, x'm) = P(y=y'm|x, x'm) \\ = 1 - P(y \neq y'm|x, x'm) \\ = 1 - \sum_{y \neq y'm} p(y=B|x, x'm) \\ = 1 - \sum_{y \neq y'm} q_B(x) q_B(x'm)$$

$$7) \pi(x, x'm) \xrightarrow{m \rightarrow \infty} 1 - q_B(x)^2 \quad \text{car } x'm \xrightarrow{\text{prob}} x \quad q_B(x'm) \xrightarrow{\text{prob}} q_B(x)$$

$$8) \pi(x) = 1 - q_B(x)^2$$

9) $\mathbb{P}(n(x) \leq n_f(x) (2 - \frac{K}{K-1} n_f(x)))$
 En utilisant $(\sum v_i)^2 \leq (\sum v_i^2) (\sum 1)$
 $n_f(x) = 1 - \max_y p(y|x) = 1 - p(y|x)$ pour $y = \arg \max_y p(y|x)$
 avec $v_i = (\sum v_i)^2 \leq (\sum v_i^2) (\sum 1)$
 $(\sum q_i)^2 \leq (\sum q_i^2) \times (K-1) \Rightarrow (1 - n_f)^2 \leq \dots$

TD3 Exercice 1

- 1) f convexe $\Leftrightarrow f'' \geq 0 \Leftrightarrow f'$ est croissant
 $f(x) = x \cos(x)$ la fonction trigonométrique + de 2 fois l'axe des abscisses \Rightarrow pas convexe
 $g(x) = -\log(x) + x^2$ $g''(x) = \frac{1}{x^2} + 2 \geq 0$ (2 fonctions convexes) \Rightarrow convexe
 $h(x) = x\sqrt{x}$ $h'(x) = \frac{3}{2}\sqrt{x}$ $h''(x) = \frac{3}{4\sqrt{x}} \geq 0 \Rightarrow$ convexe
 $t(x) = -\log(x) - \log(10-x) \rightarrow -\frac{1}{x} + \frac{1}{10-x} \rightarrow \frac{1}{x^2} + \frac{1}{(10-x)^2} > 0 \Rightarrow$ convexe

2) $\nabla_x f = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{pmatrix}$ ici $f(x) = 2x_1 + x_2^2 + x_2 x_3$
 $\nabla f = \begin{pmatrix} 2 \\ 2x_2 + x_3 \\ x_2 \end{pmatrix}$ $\nabla_x f$ est de même dimension que x

$$3) \nabla_x (f(x) + g(x)) = \begin{pmatrix} \frac{\partial (f+g)}{\partial x_1} \\ \vdots \\ \frac{\partial (f+g)}{\partial x_d} \end{pmatrix} = \begin{pmatrix} \frac{\partial f}{\partial x_1} + \frac{\partial g}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_d} + \frac{\partial g}{\partial x_d} \end{pmatrix} = \nabla f + \nabla g$$

$$\nabla_x (t f(x)) = \begin{pmatrix} \frac{\partial (t f)}{\partial x_1} \\ \vdots \\ \frac{\partial (t f)}{\partial x_d} \end{pmatrix} = \begin{pmatrix} t \frac{\partial f}{\partial x_1} \\ \vdots \\ t \frac{\partial f}{\partial x_d} \end{pmatrix} = t \nabla f$$

Exercice 2



1) estimation de la densité P_B de X pour B de volume V

On a $D = \{x_i\}_{i=1}^m$ $P(X \in B)$

$P(X \in B) = \sum_{x_i \in B} p(x_i) \Delta x$ + il s'annule

hyp: $x \in B$ est suffisamment petit, $p(x) \approx p_B$ passage à la limite p_B constante

$$P(X \in B) \approx \sum_{x_i \in B} \Delta x = P_B V$$

$X \in B$ loi de Bernoulli success ou échec

Si X_i iid et de même loi que X , $\sum_{i=1}^m X_i \sim \text{Binomiale}(p(X \in B), m)$

$$E(\sum_{i=1}^m X_i) = m p(X \in B) \approx \sum_{i=1}^m \mathbb{1}_{\{x_i \in B\}}$$

$$\rightarrow P(X \in B) \approx \frac{k}{m}$$

$$\rightarrow P_B V \approx \frac{k}{m} \quad P_B \approx \frac{k}{mV}$$

2) $X \in X$ On a $X_0 = \{x_i\}_{i=1}^m$

$p(x)$?

x_0, y_0						x_1, y_1
0						
1	x				xyy	
2		x		xyy		
3	x				x	
x_0, y_0	0	1	2	3	4	x_1, y_1

résolution Δ
en x : $m_x = \lceil \frac{x_1 - x_0}{\Delta} \rceil$ no d'intervalle

$$m_y = \lceil \frac{y_1 - y_0}{\Delta} \rceil$$

$$\Delta x = \frac{x_1 - x_0}{m_x}$$

$$\Delta y = \frac{y_1 - y_0}{m_y}$$

On suppose que la densité est constante dans les hypercubes rectangles

Pour chaque rectangle π_{ij} on estime $P_{ij} = \frac{k_{ij}}{mV} = \frac{k_{ij}}{m \Delta x \Delta y}$

3) P_B de la méthode des histo: dépend de la manière dont on découpe les données
 $\phi(x)$ $p(x) \approx \frac{1}{m} \sum \phi(x - x_i)$

ϕ \rightarrow symétrique la direction n'a pas d'importance

ϕ \rightarrow doublement monotone avec un max en 0 plus un point est proche, plus il est important

ϕ $\rightarrow \int \phi(x) dx = 1$ densité de probabilité

On utilise souvent un facteur d'échelle $\phi'(x) = \phi(\frac{x}{\sigma})$ pour avoir une résolution plus fine ou plus grossière

$\sigma \in \mathbb{R}^+$

σ le voisinage est plus grand et la résolution est plus fine

σ le voisinage est plus petit et la résolution est plus grossière

$$\text{cas 1D: } \int \phi'(x) dx = \int \phi(\frac{x}{\sigma}) dx$$

\rightarrow changement de variable $y: x \rightarrow \sigma y$

$$= \int \phi(y) \left| \det \frac{dy}{dx} \right| dx = \sigma \int \phi(y) dy$$

$$\Rightarrow \phi(\sigma x) = \frac{1}{\sigma} \phi(\frac{x}{\sigma})$$

Exercice 3

1) Parzen \rightarrow prends un k suffisamment grand
 \rightarrow classe choisie = classe majoritaire de l'hypercube centré en x
 \rightarrow prend en compte la localité des points

\rightarrow classe choisie = classe majoritaire parmi les k plus proches voisins
K-NN: \rightarrow avec un k majoritaire, prédit la classe majoritaire

2) et 3) voir feuille (en bleu)

4) voir feuille (en gris)

$k \rightarrow \infty \Rightarrow$ classe majoritaire

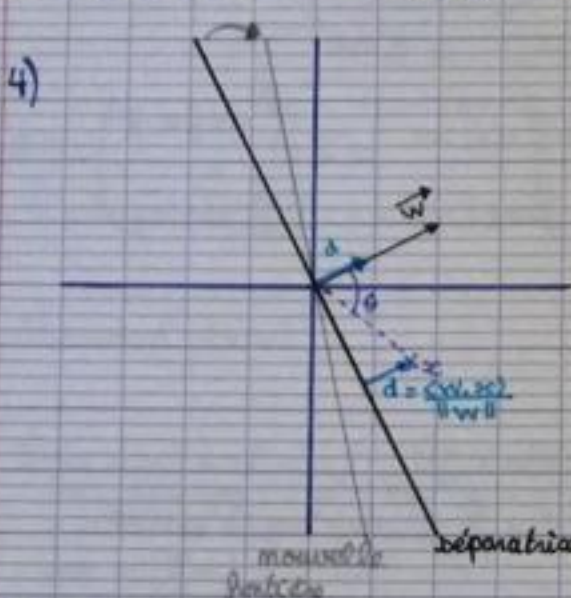
TD4 Exercice 1

$$f_w(x) = \sum_{i=1}^d x_i w_i + b = \langle x, w \rangle + b$$

$$1) HSE(f_w(x), y) = (f_w(x) - y)^2 \text{ si } y \in \mathcal{B}^d : \|f_w(x) - y\|^2$$

$$2) L(f_w(x), y) = \max(0, -y f_w(x))$$

- si $y > 0$ et $f_w(x) > 0$ } $-y f_w(x) < 0 \Rightarrow L = 0$
- ou $y < 0$ et $f_w(x) < 0$
- sinon $L > 0$



$$\langle w, x \rangle = \|w\| \|x\| \cos(\theta)$$

$$w' = w + yx$$

$$\langle w', x \rangle = \langle w + yx, x \rangle = \langle w, x \rangle + y \|x\|^2$$

5) w^1, w^2 sont les mêmes que $w = (2, 1)$ car ils sont colinéaires
 w' est le symétrique de $w \Rightarrow$ inversement des classes

$$6) \nabla_w L = \begin{cases} -y f_w(x) < 0 \rightarrow L = 0 \rightarrow 0 \\ -y (\langle w, x \rangle + b) > 0 \rightarrow -y x \end{cases}$$

$$w' = w - \sum \nabla_w L = w + \begin{cases} \sum y x + L > 0 \\ 0 \leftarrow L = 0 \end{cases}$$

$$\max(0, x - y f_w(x))$$

$$y \cdot f_w(x) = 0.5 \quad -y \cdot f_w(x) = -0.5 \rightarrow L = 0 \quad L = -0.5$$

$$y \cdot f_w(x) \geq d \Rightarrow L = 0$$

7) Si $m = 0$, alors il revient toujours à 0 donc on pose $\max(0, d - y f_w(x))$
 $y \cdot f_w(x) = 0.5$
 $-y \cdot f_w(x) = 0.5 \rightarrow L = 0$

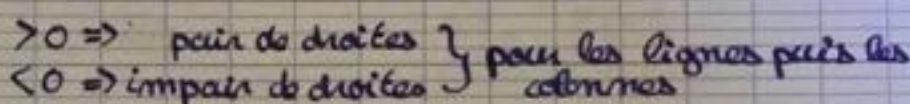
$$8) \text{Batch: } \nabla L(f(x), x) = \frac{1}{N} \sum_{i=1}^N \nabla L(f(x_i), y_i)$$

stochastique: \rightarrow on mélange les points de X
 for (x_i, y_i) :
 $w \leftarrow w - \eta \nabla L(f(x_i), y_i)$

mini-batch: (x_1, \dots, x_D)
 for x_i :
 $w \leftarrow w - \eta \sum \nabla L(f(x_i), y_i)$

- $$\Rightarrow f_W(x_i) y_i = 1 \Rightarrow x_i = 0 \Rightarrow x_i(1 - f_W(x_i) y_i) = 0 \Rightarrow x_i = 0$$

$d_i = 0 \Rightarrow \beta_i = K \Rightarrow \hat{z}_i = 0 \Rightarrow$ point bien classé

$$f_W(x) = \langle W, x \rangle + b = \left\langle \sum_i d_i y_i \alpha_i, x \right\rangle + b = \sum_i d_i y_i \langle x, \alpha_i \rangle + b$$


A 7x7 grid of numbers 1-6 with a 4x4 grid of numbers 1-6 below it.

$$3) \frac{1}{N} \sum_{i=1}^N (\hat{f}_w(x^i) - y^i)^2 \xrightarrow{n \rightarrow \infty} R(\hat{f}_w) = \int (\hat{f}_w(x-y))^2 p(x,y) dx dy = \mathbb{E} (\hat{f}_w(x) - y)^2$$

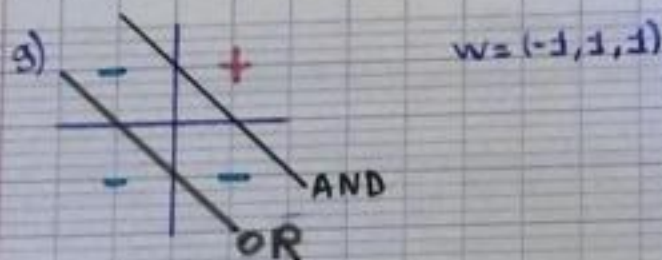
$$\frac{\partial R}{\partial f_w} = 2 S_y (f_w(x) - g) p(y|x) dy = 0$$

$$(\Rightarrow) \int w^*(x) \int p(y|x) dy = \int y p(y|x) dy = 0$$

$$\Rightarrow \int y v'(x) p(y|x) dy = \int y p(y|x) dy$$

$$(2) f_W^*(x) = E[g|x]$$

5) $\begin{pmatrix} 0 \\ 0 \\ 4 \\ 0 \\ 0 \end{pmatrix} \begin{matrix} \beta_1(x) \\ \beta_2(x) \\ \vdots \\ \beta_k(x) \end{matrix}$ $f_k(x) = \sum_{i \neq k} 0 \cdot p(y_i|x) + 1 \cdot p(y_k|x)$



Exercice 3

1) $x \in \mathbb{R}^2 \rightarrow (x_1, x_2)$
 $f_w(x) = \langle x, w \rangle \quad w \in \mathbb{R}^d$

2) $\phi(x) = (1, x_1, x_2, x_1 x_2, x_1^2, x_2^2) \in \mathbb{R}^6$
 $w \rightarrow \mathbb{R}^6$
 $f_w(\phi(x)) = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1 x_2 + w_4 x_1^2 + w_5 x_2^2$
 $x_1^2, w_4, w_5 > 0 \rightarrow \text{ellipsoïde}$
 $x_1^2, w_4, w_5 < 0 \rightarrow \text{hyperboloïde}$
 $\text{sinon} \rightarrow \text{parabole}$

3) Qui c'est plutôt intéressant (on a plus d'expressivité mais on garde des fonctions "propres")
 \rightarrow pas le même coût (on ne peut pas le comparer entre familles de fonctions diff.)
 \Rightarrow ça nous intéresse le taux de bonne classification

4)

On peut le réaliser en utilisant w_3, w_4 et w_5

$$\phi(x) = \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \end{pmatrix} = \begin{pmatrix} S(x, p, \delta) \\ \vdots \end{pmatrix}$$

$$f_w(\phi(x)) = \sum_{i=1}^6 \sum_{j=1}^2 w_{i,j} \phi_{i,j}(x)$$

$$= \sum_i \sum_j w_{i,j} e^{-\frac{\|x - p_i\|^2}{\delta^2}}$$

\Rightarrow on classe x uniquement en fonction de son voisinage.

TD5 Exercice 1

1) frontière de décision: $f_w(x) = \langle x, w \rangle + b \in \mathbb{R}$
 $f_d(x) = \langle d, x, w \rangle + a, b = d (\langle x, w \rangle + b)$

1) $x_A = x_B + \gamma y_i \frac{w}{\|w\|}$ (projeté)
 $x_B = x_A - \gamma y_i \frac{w}{\|w\|}$ \rightarrow translation du point

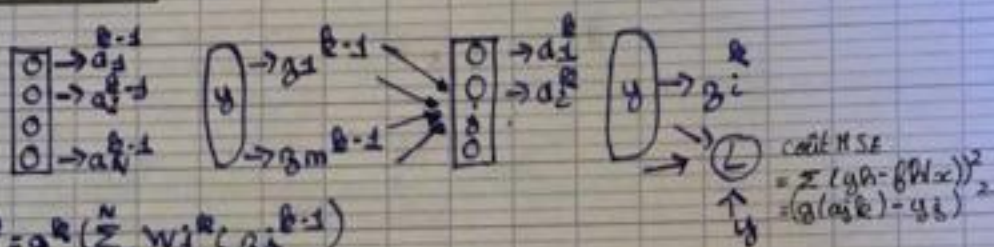
$\langle x_B, w \rangle = 0$ car B est sur la frontière de décision
 $\Rightarrow \langle x_A - \gamma y_i \frac{w}{\|w\|}, w \rangle + b = 0$
 $\Rightarrow \langle x_A, w \rangle + b - \gamma y_i \frac{\langle w, w \rangle}{\|w\|} = 0$
 $\Rightarrow f_w(x_A) = \gamma y_i \frac{\|w\|^2}{\|w\|}$
 $\Rightarrow \gamma = \frac{y_i f_w(x_A)}{\|w\|}$

$$\begin{aligned}
 7) R(fw(x)) &= \int_0^1 (fw(x) - y)^2 p(y|bc) dy \\
 &= E_{y|x} ((fw(x) - y)^2) \\
 &= E_{y|x} ((fw(x) - f^*(x) + f^*(x) - y)^2) \\
 &= E_{y|x} ((fw(x) - f^*(x))^2 + (f^*(x) - y)^2 + 2(fw(x) - f^*(x))(f^*(x) - y))
 \end{aligned}$$

$$\begin{aligned}
 8) E_{y|x} ((fw(x) - f^*(x))(f^*(x) - y)) &= (fw(x) - f^*(x)) E_{y|x} (f^*(x) - y) \\
 &= (fw(x) - f^*(x)) (f^*(x) - E_{y|x}(y)) \\
 &= 0
 \end{aligned}$$

$$R(fw(x)) = E_{y|x} ((fw(x) - f^*(x))^2) + \underbrace{(f^*(x) - y)^2}_{\text{variance}}$$

Exercise 3



$$z_i^0 = g\left(\sum_{j=1}^n w_{j1} z_j^{0-1}\right)$$

$$\frac{\partial L}{\partial w_{ij}^k} = \frac{\partial L}{\partial a_j^k} \frac{\partial a_j^k}{\partial w_{ij}^k}$$

$$\frac{\partial a_j^k}{\partial w_{ij}^k} = z_j^{k-1}$$

$$\frac{\partial L}{\partial a_j^k} = \frac{\partial (g(a_j^k) - y_j)^2}{\partial a_j^k} = 2(g(a_j^k) - y_j) g'(a_j^k)$$

$$\begin{aligned}
 \delta_i^k &= \frac{\partial L}{\partial a_i^k} = \sum_j \frac{\partial L}{\partial a_j^{k+1}} \frac{\partial a_j^{k+1}}{\partial a_i^k} \\
 &= \sum_j \delta_j^{k+1} \frac{\partial w_{ji}^{k+1} g'(a_i^k)}{\partial a_i^k} \\
 &= \sum_j \delta_j^{k+1} w_{ji}^{k+1} g'(a_i^k)
 \end{aligned}$$

$$\frac{\partial L}{\partial w_{ij}^k} = \frac{\partial L}{\partial a_j^k} \frac{\partial a_j^k}{\partial w_{ij}^k} = \delta_j^k z_j^{k-1} = \left(\sum_h \delta_h^{k+1} w_{jh}^{k+1} g'(a_j^k) \right) (z_j^{k-1})$$

TD7 Exercise 1

$$X \xrightarrow{\phi} E \quad \text{Kernel}(bc, y) = \langle \phi(x), \phi(y) \rangle$$

$$\begin{aligned}
 K_1(x, y) &= \langle \phi_1(x), \phi_1(y) \rangle \\
 K_2(x, y) &= \langle \phi_2(x), \phi_2(y) \rangle
 \end{aligned}$$

$$a) cK_1(x, y) \stackrel{?}{=} \langle \phi_3(x), \phi_3(y) \rangle \quad \text{on pose } \phi_3(x) = \sqrt{c} \phi_1(x)$$

$$\begin{aligned}
 \langle \phi_3(x), \phi_3(y) \rangle &= \langle \sqrt{c} \phi_1(x), \sqrt{c} \phi_1(y) \rangle \\
 &= c \langle \phi_1(x), \phi_1(y) \rangle \\
 &= c K_1(x, y)
 \end{aligned}$$

$$b) K_4(x, y) \stackrel{?}{=} K_1(x, y) + K_2(x, y) = \langle \phi_1(x), \phi_1(y) \rangle + \langle \phi_2(x), \phi_2(y) \rangle$$

$$\begin{aligned}
 \text{on pose } \phi_4(x) &= \begin{pmatrix} \phi_1(x) \\ \phi_2(x) \end{pmatrix} \quad \langle \phi_4(x), \phi_4(y) \rangle = \left\langle \begin{pmatrix} \phi_1(x) \\ \phi_2(x) \end{pmatrix}, \begin{pmatrix} \phi_1(y) \\ \phi_2(y) \end{pmatrix} \right\rangle \\
 &= \langle \phi_1(x), \phi_1(y) \rangle + \langle \phi_2(x), \phi_2(y) \rangle
 \end{aligned}$$

2) $\gamma^* = \min \frac{\sum_i f_w(x_i)}{\|w\|}$
 Si on souhaite maximiser, il faut minimiser $\frac{\sum_i f_w(x_i)}{\|w\|}$

2) 1) ≥ 1 car le point le proche est 1 et il faut que les points soient bien classés.
 On veut que les points soient $\geq 1 \Rightarrow \exists i \text{ tq } f_w(x_i) y_i = 1$
 Ça ne doit pas être 0 car sinon on ne considère que la frontière et pas la marge \Rightarrow on retrouve dans le cas du perceptron.

2) Quand on a un pb avec argmin $(C(w))$ et des contraintes tq $g_1(w) = 0, g_2(w) = 0 \dots g_m(w) = 0$ et $c_1(w) \leq 0 \dots c_n(w) \leq 0$.

Au point solution : les gradients sont colinéaires. (sinon on pourrait continuer à descendre dans la descente de gradient et trouver une autre solution)

\Rightarrow On a le Lagrangien $L(w, \beta, \alpha, d) = C(w) + \beta_1 g_1(w) + \dots + \beta_n g_n(w) + \alpha_1 c_1(w) + \dots + \alpha_n c_n(w)$
 $= \min_w \max_{\beta, \alpha} L(w, \beta, \alpha) \Leftrightarrow \max_{\beta, \alpha} \min_w L(w, \beta, \alpha)$

3) Les contraintes sont très difficiles à respecter. (les contraintes d'égalités demandent beaucoup moins de corrections à faire).
 \Rightarrow si les contraintes ne sont pas respectées, la quantité max sera non bornée.

On pose le Lagrangien du SVM. On ajoute m contraintes (car on a m exemples)

$L(w, \alpha) = \frac{1}{2} \|w\|^2 + \sum_{i=1}^m \alpha_i (1 - \langle x_i, w \rangle + b) y_i$

\Rightarrow on dérive par rapport à w
 $\nabla_w L = w - \sum_{i=1}^m \alpha_i y_i x_i = 0 \Rightarrow w^* = \sum_{i=1}^m \alpha_i y_i x_i$

$\frac{\partial L}{\partial b} = \sum_{i=1}^m \alpha_i y_i = 0$

On va remplacer w et b
 $L(w, \alpha) = \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i y_i x_i \right\|^2 + \sum_{i=1}^m \alpha_i \left(1 - \left\langle x_i, \sum_{j=1}^m \alpha_j y_j x_j \right\rangle + b \right) y_i$
 $= \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle + \sum_{i=1}^m \alpha_i (1 - \sum_{j=1}^m \alpha_j y_j \langle x_i, x_j \rangle + b y_i)$
 $= \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle + \sum_{i=1}^m \alpha_i (1 - \sum_{j=1}^m \alpha_j y_j \langle x_i, x_j \rangle + b y_i)$

$L(w^*, \alpha) = \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle + \sum_{i=1}^m \alpha_i$

On veut maximiser par rapport à α

7) argmin $\frac{1}{2} \|w\|^2 + w \sum_{i=1}^N \alpha_i$ tq $\forall i, f_w(x_i) y_i \geq 1 - \alpha_i$ et au min du Lagrangien

$L(w, \alpha, \beta) = \frac{1}{2} \|w\|^2 + \sum_{i=1}^N \alpha_i (1 - \langle x_i, w \rangle + b) y_i - \sum_{i=1}^N \beta_i (1 - \langle x_i, w \rangle + b) y_i - \sum_{i=1}^N \beta_i \alpha_i$

$\frac{\partial L}{\partial \alpha_i} = -\alpha_i + w - \beta_i = 0 \Rightarrow w = \alpha_i + \beta_i \Rightarrow \alpha_i \leq w$

8) Ça s'annule dans l'optimisation du Lagrangien donc on a
 $\max_{\alpha} L(w^*, \alpha) = \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle + \sum_{i=1}^N \alpha_i$

9) $\nabla L = 0$

$1 - \alpha_i - f_w(x_i) y_i \leq 0$

$\alpha_i / \beta_i \geq 0$

$\alpha_i (1 - \langle x_i, w \rangle + b) y_i = 0$