

**Exercice 1 – Perceptron (examen 2009)**

On s'intéresse à des classifieurs linéaires pour faire de la discrimination à deux classes. On note  $D = \{x^i, y^i\}_{i=1, \dots, N}$  un ensemble d'apprentissage avec  $y^i = 1$  si  $x^i \in C_1$  et  $y^i = -1$  si  $x^i \in C_2$ .  $\mathbf{w}$  est le vecteur de poids du classifieur linéaire.

**Q 1.1** Rappeler l'algorithme du perceptron

1. INPUT : xapp, yapp, epsilon, OPT : nbIterations
2. Initialiser  $\mathbf{w}$  (random, zeros...)
3. Tant que (critère insatisfait)
  - (a) Tirer un échantillon  $\mathbf{x}^i$  aléatoirement
  - (b) Si Erreur sur  $\mathbf{x}^i$ , i.e.  $\mathbf{x}^i \mathbf{w} y^i \leq 0$   
Alors :  $\mathbf{w} = \mathbf{w} + \varepsilon * y^i * x^{iT}$

**Q 1.2** On suppose que l'algorithme est initialisé avec  $\mathbf{w}_0 = 0$

**Q 1.2.1** Montrer qu'à une étape de l'algorithme, il existe des coefficients  $\alpha_i$  tels que  $\mathbf{w}_t = \sum_{i=1}^N \alpha_i y^i x^i$

$$\mathbf{w}_1 = \varepsilon * y^i * x^{iT} \text{ car } \mathbf{w}_0 = 0$$

$$\mathbf{w}_2 = \mathbf{w}_1 + \varepsilon * y^i * x^{iT}$$

$$\mathbf{w}_t = \sum_{i=1}^N \alpha_i * y^i * x^{iT}, \quad \alpha_i = k_i \varepsilon$$

$k_i$  est le nombre de fois où  $x^i$  a été tiré et a été mal classé.

**Q 1.2.2** Exprimer en fonction des  $\alpha_i$  la condition qui indique que le perceptron fait une erreur sur  $x_i$ .

$$f(\mathbf{x}_i) = \sum_{k=1}^N \alpha_k y^k x^i x^{kT} = \sum_{k=1}^N \alpha_k y^k \langle x^i, x^k \rangle$$

Critère :

$$f(\mathbf{x}_i) y_i \leq 0$$

**Q 1.2.3** Reformuler l'algorithme du perceptron avec les seuls  $\alpha_i$  comme paramètres.

1. INPUT : xapp, yapp, epsilon, OPT : nbIterations
2. Initialiser w (random, zeros...)
3. Tant que (critère insatisfait)
  - (a) Tirer un échantillon  $\mathbf{x}^i$  aléatoirement
  - (b) Si Erreur sur  $\mathbf{x}^i$ , i.e.  $\sum_{k=1}^N \alpha_k y^k \langle x^i, x^k \rangle y^i \leq 0$   
Alors :  $\alpha_i \leftarrow \alpha_i + \varepsilon$

**Q 1.2.4** Exprimer la fonction de décision en fonction des  $\alpha_i$

**Q 1.3** On considère le critère d'apprentissage  $Q(\mathbf{w}) = \sum_{i \in Z} -y^i x^i \mathbf{w}$  avec  $Z = \{x^i | y^i x^i \mathbf{w} < 0\}$  l'ensemble des points mal classés par  $\mathbf{w}$ .

**Q 1.3.1** Donner l'algorithme de gradient qui minimise ce critère.

$$\nabla_{\mathbf{w}} Q = \sum_{i \in Z} -y^i x^{iT}$$

**Q 1.3.2** Donner l'algorithme de gradient adaptatif (=stochastique) qui minimise ce critère. Quel algorithme reconnaissez vous ?

Retour sur le perceptron classique

---

## Exercice 2 – Classification d'images binaires (exam 2010)

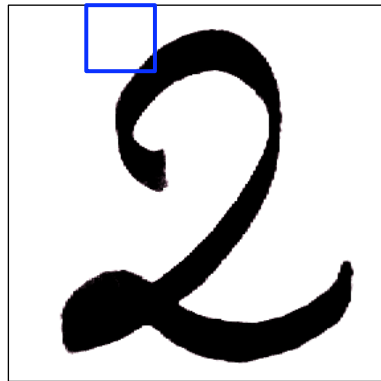
---

Soit une base d'images binaires (pixels noirs ou blancs) de taille 256x256 pixels. Un exemple d'image est fourni Fig. ?? . Nous avons tenté au cours du semestre une analyse bayésienne pixel par pixel à l'aide d'une loi de Bernoulli. Malheureusement, les chiffres n'étaient pas tous positionnés de la même manière et cela à causer une perte de performances. Pour palier ce problème, nous proposons de recommencer l'expérience en utilisant une fenêtre glissante sur l'image. Une telle fenêtre est dessinée sur l'exemple.

La fenêtre peut prendre  $N$  positions dans l'image. Nous indexerons la position en  $i$ ,  $x_i$  désigne la fenêtre dans la position  $i$  avec  $i$  variant de 1 à  $N$ . Pour une image donnée, la variable  $x_i$  prend la valeur  $k_i$ .  $k_i$  est le nombre de pixels noirs dans la fenêtre.  $x_i$  suit une loi Binomiale  $\mathcal{B}(n, p_i)$ . La fenêtre choisie est de taille 2x2 soit 4 pixels. Afin de représenter une image  $\mathbf{x}$ , nous utiliserons l'ensemble des fenêtres :  $\{x_i\}_{i=1, \dots, N}$ .

1. Quelles valeurs peuvent prendre les  $x_i$  ?

5 valeurs de 0 à 5



2. Donner l'expression de la loi binomiale  $p(x_i = k_i)$ ?. Que modélisent  $n$  et  $p_i$ ?

$$\text{Rappelons que pour la loi binomiale : } p(x_i = k_i) = C_n^{k_i} p_i^{k_i} (1 - p_i)^{n - k_i}$$

3. Soit  $\mathbf{x}$  une image,  $\mathbf{x} = \{x_1, \dots, x_N\}$  donner l'expression de  $p(\mathbf{x})$  si la valeur observée pour  $x_i$  est  $k_i$ .

$$p(\mathbf{x}) = \prod_{i=1}^N p(x_i = k_i) = \prod_{i=1}^N C_n^{k_i} p_i^{k_i} (1 - p_i)^{n - k_i}$$

4. On note  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_P\}$  un ensemble d'image  $\mathbf{x}_j = \{x_{j,1}, \dots, x_{j,N}\}$ . Donner la log vraisemblance de  $X$

$$\begin{aligned} \log L &= \sum_{ij} \log p(x_i^j = k_i^j) \\ \log L &= \sum_{ij} \log C_n^{k_i^j} + k_i^j \log p_i + (n - k_i^j) \log(1 - p_i) \end{aligned}$$

5. Montrer que la vraisemblance est maximisée pour  $p_i = \frac{\sum_{j=1}^P k_{j,i}}{nP}$

$$\begin{aligned} \frac{\partial \log L}{\partial p_i} &= \sum_j k_i^j \frac{1}{p_i} + (n - k_i^j) \frac{1}{1 - p_i} \\ \frac{\partial \log L}{\partial p_i} &= 0 \Leftrightarrow p_i = \frac{K_i}{nP} \end{aligned}$$

$$\text{avec } K_i = \sum_j k_i^j$$

6. On veut résoudre un problème à deux classes avec cette approche, comment procède-t-on?

On considère un problème de régression linéaire. La cible est une variable  $t$  qui est une combinaison linéaire des variables d'entrée  $\mathbf{x}$  plus un bruit blanc. Ce qui s'écrit :  $t = y + \varepsilon$ .

Avec :

$$y = \mathbf{w}^T \mathbf{x} = w_0 x_0 + w_1 x_1 + \dots + w_n x_n$$

$$t \in \mathbb{R}, y \in \mathbb{R}, \mathbf{x} \in \mathbb{R}^{n+1}, \varepsilon \sim \mathcal{N}(0, \beta^{-1})$$

$\beta = \frac{1}{\sigma^2}$ , inverse de la variance est appelé *précision*. De manière générale, on note  $\mathcal{N}(x|\mu, \sigma^2)$  la distribution gaussienne de moyenne  $\mu$  et d'écart type  $\sigma$  d'une variable  $x$ . Dans toute la suite  $x_0 = 1$ , pour les vecteurs  $\mathbf{x}$  (on a augmenté tous les vecteurs avec une première coordonnée 1). Soit  $D = \{\mathbf{x}_i, t_i\}$  un ensemble d'apprentissage. On note  $\mathbf{t} = (t_1, \dots, t_N)^T$  le vecteur colonne des sorties cibles et  $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$  l'ensemble des vecteurs d'entrée pour les données d'apprentissage.

### Q 3.1 Préliminaires

**Q 3.1.1** Montrer que pour des valeurs  $\mathbf{x}, \mathbf{w}, \beta$  données,  $t$  suit une loi normale  $\mathcal{N}(y, \beta^{-1})$ . On notera par la suite  $p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|y, \beta^{-1})$  cette probabilité.

On remarque que  $t$  est la somme d'une valeur constante  $y$  et d'une variable aléatoire gaussienne  $\varepsilon$ .

C'est surtout une histoire de notation. On s'en tient à la deuxième formule : il suffit de démontrer que  $t$  à pour espérance  $y$  et pour écart type  $\beta^{-1}$

$$E[t] = E[\mathbf{w}^T \mathbf{x} + \varepsilon] = E[\mathbf{w}^T \mathbf{x}] + 0 = \mathbf{w}^T \mathbf{x} = y$$

car on travaille pour des valeurs  $\mathbf{x}, \mathbf{w}, \beta$  données.

$$V[t] = V[\mathbf{w}^T \mathbf{x} + \varepsilon] = 0 + \beta^{-1}$$

Donc :  $t \sim \mathcal{N}(y, \beta^{-1})$

**Q 3.1.2** Donner l'expression de  $p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)$

ATTENTION : il s'agit du vecteur  $\mathbf{t}$

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = \prod_{i=1}^N p(t_i|\mathbf{x}, \mathbf{w}, \beta) = \prod_{i=1}^N p(t_i|\mathbf{x}_i, \mathbf{w}, \beta)$$

Définition d'une gaussienne (il faut aussi leur apprendre à lire les questions un peu en avance...)

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{1}{2\sigma^2} \|t_i - y_i\|^2\right)$$

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = \left(\frac{\beta}{2\pi}\right)^{N/2} \exp\left(-\frac{\beta}{2} \sum_i \|t_i - y_i\|^2\right)$$

**Q 3.1.3** Montrer que

$$\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w})$$

Avec :

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (t_i - \mathbf{w}^T \mathbf{x}_i)^2$$

Application de ln sur la formule précédente

**Q 3.2** Estimation du maximum de vraisemblance

**Q 3.2.1** Calculer les gradients  $\frac{\partial \ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)}{\partial \mathbf{w}}$  et  $\frac{\partial \ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)}{\partial \beta}$

Remarquer que cette question est indépendante des précédentes

$$\frac{\partial \ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)}{\partial w_j} = \beta \sum_i -x_{ij}(t_i - \mathbf{w}^T \mathbf{x}_i)$$

$$\frac{\partial \ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)}{\partial \mathbf{w}} = -\beta X^T (\mathbf{t} - X\mathbf{w})$$

$$\frac{\partial \ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)}{\partial \beta} = \frac{N}{2} \frac{1}{\beta} - E_D(\mathbf{w})$$

**Q 3.2.2** Montrer que la solution à  $\frac{\partial \ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)}{\partial \mathbf{w}} = 0$  est le vecteur  $\mathbf{w}_{ML} = (XX^T)^{-1} X\mathbf{t}$

**Q 3.2.3** ...