

Devoir 1 - Partie Théorique

- Ce devoir doit être fait et envoyé sur Gradescope individuellement. Vous pouvez discuter avec d'autres étudiants mais les réponses que vous soumettez doivent être les vôtres. A noter que nous utiliserons l'outil de détection de plagiat de Gradescope. Tous les cas suspectés de plagiat seront enregistrés et transmis à l'Université pour vérification.
- Vous devez soumettre vos solutions au format pdf sur Gradescope en utilisant le devoir intitulé **Devoir 1 - Partie Théorique**.

1. **Rappels de probabilités: probabilité conditionnelle et règle de Bayes [5 points]**

- (a) Donnez la définition de la probabilité conditionnelle de la variable aléatoire discrète X sachant la variable aléatoire discrète Y

Answer. $P(X|Y) = P(X, Y)/P(Y)$

- (b) Soit une pièce déséquilibrée dont la probabilité d'obtenir face est $2/3$ et la probabilité d'obtenir pile est $1/3$. Cette pièce est lancée à trois reprises. Quelle est la probabilité d'obtenir exactement deux faces (parmi les trois lancers), sachant que le premier lancer a fait face ?

Answer.

$$\begin{aligned} P(2 \text{ H. in } 3 \text{ tosses} | \text{H. in } 1 \text{ toss}) &= P(1 \text{ H. in } 2 \text{ tosses}) \\ &= 2 * P(H.) * P(T.) = 2 * \frac{2}{3} * \frac{1}{3} = \frac{4}{9} \end{aligned}$$

- (c) Donnez deux expressions équivalentes de $P(X, Y)$:

- (i) en fonction de $\mathbb{P}(X)$ et $\mathbb{P}(Y|X)$

Answer. $P(X, Y) = P(Y|X)P(X)$

- (ii) en fonction de $\mathbb{P}(Y)$ et $\mathbb{P}(X|Y)$

Answer. $P(X, Y) = P(X|Y)P(Y)$

(d) Prouvez le théorème de Bayes:

$$\mathbb{P}(X|Y) = \frac{\mathbb{P}(Y|X)\mathbb{P}(X)}{\mathbb{P}(Y)}.$$

Answer. $P(X, Y) = P(X|Y) * P(Y) = P(Y|X) * P(X) \implies P(X|Y) = \frac{P(Y|X) * P(X)}{P(Y)}$

(e) Un sondage des étudiants Montréalais est fait, où 60% des élèves sondés sont affiliés à l'UdeM alors que les autres sont affiliés à McGill. Un étudiant est choisi aléatoirement parmi ce groupe.

i. Quelle est la probabilité que l'étudiant soit affilié à McGill?

Answer. 40%

ii. Considérons maintenant que l'étudiant est bilingue, et que 60% des étudiants de l'UdeM sont bilingues alors que seulement 30% des étudiants de McGill le sont. Étant donné cette information, quelle est la probabilité que cet étudiant soit affilié à McGill ?

Answer. $P(\text{McGill}|\text{bilingual}) = \frac{P(\text{bilingual}|\text{McGill})P(\text{McGill})}{P(\text{bilingual})} = \frac{0.3 * 0.4}{0.3 * 0.4 + 0.6 * 0.6} = 0.25$

2. Bag of words (sac de mots) et modèle de sujet unique [12 points]

On s'intéresse à un problème de classification où l'on veut prédire le sujet d'un document d'un certain corpus (ensemble de documents). Le sujet de chaque document peut être soit *sport*, soit *politique*. 1/3 des documents du corpus sont sur le *sport*, et 2/3 sont sur la *politique*.

On va utiliser un modèle très simple où on ignore l'ordre des mots apparaissant dans le document et l'on considère que les mots dans un document sont indépendants les uns des autres, étant donné le sujet du document.

De plus, nous allons utiliser des statistiques très simples des documents: les probabilités qu'un mot choisi au hasard dans un document soit "goal", "kick", "congress", "vote", ou n'importe quel autre mot (dénnoté par *other*). Nous appelons ces cinq catégories le vocabulaire ou dictionnaire pour les documents: $V = \{\text{"goal"}, \text{"kick"}, \text{"congress"}, \text{"vote"}, \text{"other"}\}$.

Soit les distributions suivantes des mots du vocabulaire, par sujet:

	$\mathbb{P}(\text{mot} \mid \text{sujet} = \textit{sport})$	$\mathbb{P}(\text{mot} \mid \text{sujet} = \textit{politique})$
mot = “ <i>goal</i> ”	1/100	5/1000
mot = “ <i>kick</i> ”	2/100	1/1000
mot = “ <i>congress</i> ”	1/1000	1/100
mot = “ <i>vote</i> ”	3/1000	4/100
mot = <i>other</i>	966/1000	944/1000

Table 1

Cette table nous dit par exemple que la probabilité qu’un mot choisi aléatoirement dans un document soit textit“vote” n’est que de 3/1000 si le sujet du document est le *sport*, mais est de 4/100 si le sujet est la *politique*.

- (a) Quelle est la probabilité qu’un mot aléatoire dans un document soit “goal” étant donné que le sujet est la *politique* ? **Answer.** 5/1000
- (b) Quelle est l’espérance du nombre de fois où le mot “congress” apparait dans un document de 2000 mots dont le sujet est le *sport*? **Answer.** $2000 * \frac{1}{1000} = 2$
- (c) On tire aléatoirement un document du corpus. Quelle est la probabilité qu’un mot aléatoire de ce document soit “goal”? **Answer.**

$$\begin{aligned}
 P(\textit{“goal”}) &= P(\textit{“goal”}|\textit{sports})P(\textit{sports}) + P(\textit{“goal”}|\textit{politics})P(\textit{politics}) \\
 &= \frac{1}{100} * \frac{1}{3} + \frac{5}{1000} * \frac{2}{3} \\
 &= \frac{2}{300} = \frac{1}{150} = 0.0067
 \end{aligned}$$

- (d) Supposons que l’on tire aléatoirement un mot d’un document et que ce mot est “kick”. Quelle est la probabilité que le sujet du

document soit le *sport*? **Answer.**

$$\begin{aligned}
 P(sports|"kick") &= \frac{P("kick"|sports)P(sports)}{P("kick")} \\
 &= \frac{2/100 * 1/3}{P("kick"|sports)P(sports) + P("kick"|politics)P(politics)} \\
 &= \frac{2/100 * 1/3}{2/100 * 1/3 + 1/1000 * 2/3} \\
 &= \frac{20}{22}
 \end{aligned}$$

- (e) Supposons que l'on tire aléatoirement deux mots d'un document et que le premier soit "kick". Quelle est la probabilité que le second mot soit "vote"? **Answer.** Assume drawing with replacement

$$\begin{aligned}
 P("vote"|"kick") &= P("vote"|sports)P(sports|"kick") \\
 &\quad + P("vote"|politics)P(politics|"kick")
 \end{aligned}$$

we already calculated $P(sports|"kick")$, and we could recalculate $P(politics|"kick")$ the same way or just realize that

$$\begin{aligned}
 P(politics|"kick") &= 1 - P(sports|"kick") \\
 &= \frac{2}{22}
 \end{aligned}$$

putting it all together

$$\begin{aligned}
 P("vote"|"kick") &= 3/1000 * 20/22 + 4/100 * 2/22 \\
 &\approx 0.006
 \end{aligned}$$

- (f) Pour en revenir à l'apprentissage, supposons que l'on ne connaisse pas les probabilités conditionnelles étant donné chaque sujet ni les probabilités de chaque sujet (i.e. nous n'avons pas accès aux informations de la table 1 où aux proportions de chaque sujet), mais nous avons un jeu de données de N documents où chaque document est annoté avec un des sujet *sport* ou *politique*. Comment estimeriez vous les probabilités conditionnelles (e.g., $\mathbb{P}(\text{mot} = "goal" \mid \text{sujet} = \text{politique})$) et les probabilités des sujets (e.g., $\mathbb{P}(\text{sujet} = \text{politique})$) à partir de ce jeu de données ? **Answer.** Nous estimons les probabilités de sujet en supposant

que vos données sont i.i.d et en utilisant des bayes naïfs $P(topic) = N_{topic}/N$. Nous estimons les probabilités de mots en utilisant des bayes naïfs $P(word|topic) = Count_{word}/Counts_{total}$.

3. Estimateur de maximum de vraisemblance [10 points]

Soit la fonction de densité de probabilité suivante:

$$f_{\theta}(x) = 2\theta x e^{-\theta x^2}$$

où θ est un paramètre et x est un nombre réel positif.

Supposons que n points $D = \{x_1, \dots, x_n\}$ sont tirés aléatoirement indépendemment selon $f_{\theta}(x)$.

- (a) Soit $f_{\theta}(x_1, x_2, \dots, x_n)$ la fonction de densité de probabilité jointe de n points indépendemment et identiquement distribué (i.i.d) selon $f_{\theta}(x)$. Exprimez $f_{\theta}(x_1, x_2, \dots, x_n)$ en fonction de $f_{\theta}(x_1)$, $f_{\theta}(x_2), \dots, f_{\theta}(x_n)$

Answer. $f_{\theta}(x_1, x_2, \dots, x_n) = f_{\theta}(x_1)f_{\theta}(x_2) \dots f_{\theta}(x_n) = \prod_{i=1}^n f_{\theta}(x_i)$

- (b) On définit l'estimateur du maximum de vraisemblance comme la valeur de θ qui maximise la vraisemblance de générer le jeu de donnée D à partir de la distribution $f_{\theta}(x)$. Formellement,

$$\theta_{MLE} = \arg \max_{\theta \in \mathbb{R}} f_{\theta}(x_1, x_2, \dots, x_n),$$

calculez l'estimateur du maximum de vraisemblance de θ .

Answer.

$$\begin{aligned} L(D; \theta) &= \log f_{\theta}(x_1, x_2, \dots, x_n) \\ &= \log \prod_{i=1}^n f_{\theta}(x_i) \\ &= \sum_{i=1}^n \log f_{\theta}(x_i) \\ &= \sum_{i=1}^n \log(2\theta x_i e^{-\theta x_i^2}) \\ &= \sum_{i=1}^n (\log 2 + \log \theta + \log x_i - \theta x_i^2) \end{aligned}$$

En prenant la dérivée par rapport à θ , on obtient:

$$\begin{aligned}
\frac{\partial L(D; \theta)}{\partial \theta} &= \frac{\partial}{\partial \theta} \left(\sum_{i=1}^n (\log 2 + \log \theta + \log x_i - \theta x_i^2) \right) \\
&= \sum_{i=1}^n \left(\frac{\partial}{\partial \theta} \log 2 + \frac{\partial}{\partial \theta} \log \theta + \frac{\partial}{\partial \theta} \log x_i - \frac{\partial}{\partial \theta} \theta x_i^2 \right) \\
&= \sum_{i=1}^n \left(0 + \frac{1}{\theta} + 0 - x_i^2 \right) \\
&= \sum_{i=1}^n \left(\frac{1}{\theta} - x_i^2 \right) \\
&= \frac{n}{\theta} - \sum_{i=1}^n x_i^2
\end{aligned}$$

En mettant ceci à 0 et en résolvant pour θ nous obtenons:

$$\theta = \frac{n}{\sum_{i=1}^n x_i^2}$$

4. Estimateur du maximum de vraisemblance et histogramme [10 points]

Soit le jeu de données X_1, X_2, \dots, X_n sont n données i.i.d tirées d'une distribution de probabilité constante par morceaux. Il y a N morceaux de longueur égale entre 0 et 1 (B_1, B_2, \dots, B_N), où les constantes sont $\theta_1, \theta_2, \dots, \theta_N$.

$$p(x; \theta_1, \dots, \theta_N) = \begin{cases} \theta_j & \frac{j-1}{N} \leq x < \frac{j}{N} \text{ for } j \in \{1, 2, \dots, N\} \\ 0 & \text{otherwise} \end{cases}$$

On défini μ_j pour $j \in \{1, 2, \dots, N\}$ en tant que $\mu_j := \sum_{i=1}^n \mathbb{1}(X_i \in B_j)$.

- (a) Utilisant le fait que l'aire totale sous une distribution de probabilité est 1, exprimez θ_N en fonction de $\theta_1, \theta_2, \dots, \theta_{N-1}$.

Answer. By integrating the probability density function over

its domain, we obtain:

$$\begin{aligned} 1 &= \int_{\bigcup_{j=1}^N U_j} p(x; \theta_1, \dots, \theta_N) dx = \sum_{i=1}^N \int_{U_j} p(x; \theta_1, \dots, \theta_N) dx \\ &= \sum_{i=1}^N \int_{U_j} \theta_j dx = \sum_{j=1}^N \frac{\theta_j}{N} = \frac{1}{N} \sum_{j=1}^N \theta_j = \frac{1}{N} \sum_{j=1}^{N-1} \theta_j + \frac{\theta_N}{N}. \end{aligned}$$

Thus, we have that:

$$1 = \frac{1}{N} \sum_{j=1}^{N-1} \theta_j + \frac{\theta_N}{N} \implies \theta_N = N - \sum_{j=1}^{N-1} \theta_j.$$

- (b) Écrivez le log-vraisemblance du jeu de données en fonction de $\theta_1, \theta_2, \dots, \theta_{N-1}$ et $\mu_1, \mu_2, \dots, \mu_{N-1}$.

Answer.

$$\begin{aligned} \mathcal{L}(x_{1\dots n}; \theta_{1\dots N}, \mu_{1\dots N}) &= \log \left(\prod_{i=1}^n p(x_i; \theta_{1\dots N}, \mu_{1\dots N}) \right) \\ &= \log \left(\prod_{i=1}^n \prod_{j=1}^N \theta_j^{\mathbb{1}(x_i \in B_j)} \right) \\ &= \log \left(\prod_{j=1}^N \prod_{i=1}^n \theta_j^{\mathbb{1}(x_i \in B_j)} \right) \\ &= \log \left(\prod_{j=1}^N \theta_j^{\mu_j} \right) \\ &= \log \left(\left(\prod_{j=1}^{N-1} \theta_j^{\mu_j} \right) (\theta_N^{\mu_N}) \right) \\ &= \log \left(\left(\prod_{j=1}^{N-1} \theta_j^{\mu_j} \right) \left((N - \sum_{j=1}^{N-1} \theta_j)^{(n - \sum_{j=1}^{N-1} \mu_j)} \right) \right) \\ &= \sum_{j=1}^{N-1} \mu_j \log \theta_j + (n - \sum_{j=1}^{N-1} \mu_j) \log (N - \sum_{j=1}^{N-1} \theta_j). \end{aligned}$$

- (c) Trouvez l'estimateur du maximum de vraisemblance de θ_j , $j \in \{1, 2, \dots, N\}$. **Answer.** The expression for the likelihood, as

derived in the previous points, is given by:

$$\mathcal{L}(\theta_{1\dots N}; \mu_{1\dots N}) = \sum_{j=1}^N \mu_j \log \theta_j.$$

We are looking for the arguments $\hat{\theta}_j$ which maximize this function, with the constraint that $\sum_{i=1}^N \frac{\theta_i}{N} = 1$ for the underlying density to be representing a valid probability distribution.

To find the maximum in this set, we can use the method of Lagrange multipliers, which translates the constrained optimization problem in an unconstrained one. The goal becomes to find the stationary point of the Lagrangian function:

$$\mathcal{L}_\lambda(\theta_{1\dots N}; \mu_{1\dots N}) = \sum_{j=1}^N \mu_j \log \theta_j - \lambda \left(\frac{1}{N} \sum_{j=1}^N \theta_j - 1 \right).$$

To do it, let us compute its gradient with respect to a particular θ_j and equate it to zero:

$$\nabla_{\theta_j} \mathcal{L}_\lambda(\theta_{1\dots j\dots N}; \mu_{1\dots j\dots N}) = \frac{\mu_j}{\theta_j} - \frac{\lambda}{N} = 0.$$

We can thus express θ_j in terms of μ_j as:

$$\theta_j = \frac{\mu_j N}{\lambda}.$$

If we now remember the constraint, we can find an expression for λ :

$$\frac{1}{N} \sum_{j=1}^N \theta_j = 1 \implies \frac{1}{N} \sum_{j=1}^N \frac{\mu_j N}{\lambda} = 1 \implies \lambda = \sum_{j=1}^N \mu_j = n.$$

Thus, we obtain

$$\theta_j = \frac{\mu_j N}{n}$$

as the maximum likelihood estimate for the parameters θ_j .

5. Méthodes à base d'histogrammes [10 points]

Soit le jeu de données $\{x_j\}_{j=1}^n$ où chaque point $x \in [0, 1]^d$. La fonction $f(x)$ représente la vraie distribution inconnue des données. Vous décidez d'utiliser une méthode à base d'histogramme pour estimer $f(x)$. Chaque dimension est divisée en m régions.

- (a) Démontrez que pour un ensemble mesurable S , $\mathbb{E}[\mathbb{1}_{\{x \in S\}}] = \mathbb{P}(x \in S)$.

Answer. Soit p une distribution définie sur un ensemble \mathcal{X} :

$$\begin{aligned}\mathbb{E}_{x \sim p}[\mathbb{1}_{\{x \in S\}}] &= \int_{x \in \mathcal{X}} \mathbb{1}_{\{x \in S\}} p(x) dx \\ &= \int_{x \in S} \mathbb{1}_{\{x \in S\}} p(x) dx + \int_{x \notin S} \mathbb{1}_{\{x \in S\}} p(x) dx \\ &= \int_{x \in S} p(x) dx \\ &= \mathbb{P}_{x \sim p}(x \in S)\end{aligned}$$

- (b) En combinant le résultat de la question précédente avec la loi des grands nombres, démontrez que la probabilité estimée d'être dans la région i , telle que donnée par les méthodes à base d'histogramme, tend vers $\int_{V_i} f(x) dx$, la vraie probabilité d'être dans la région i , lorsque $n \rightarrow \infty$. V_i est le volume occupée par la région i .

Answer. On note $f_{\text{hist},n}$ le p.d.f du modèle d'histogramme avec n points de données. Par définition de l'histogramme on obtient $\mathbb{P}_{x \sim f_{\text{hist},n}}(x \in V_i) = \frac{1}{n} \sum_{j=1}^n \mathbb{1}_{\{x_j \in V_i\}}$. La loi des grands nombres nous dit que la moyenne empirique converge vers l'espérance lorsque le nombre d'échantillons tend vers l'infini (presque sûrement) :

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n \mathbb{1}_{\{x_j \in V_i\}} \stackrel{\text{as}}{=} \mathbb{E}_{x \sim f} \mathbb{1}_{\{x \in V_i\}} = \mathbb{P}_{x \sim f}(x \in V_i)$$

où nous avons obtenu la dernière égalité en utilisant (a). On peut ainsi conclure

$$\lim_{n \rightarrow \infty} \mathbb{P}_{x \sim f_{\text{hist},n}}(x \in V_i) \stackrel{\text{as}}{=} \mathbb{P}_{x \sim f}(x \in V_i)$$

.

- (c) Soit le jeu de données MNIST, où chaque point vit en 784 dimensions. Nous divisons chaque dimensions en 2 régions. Combien de chiffres (en base 10) contient le nombre total de régions?

Answer. $2^{784} \approx 1.02 \times 10^{236} \Rightarrow 237$ digits

- (d) Supposons l'existence d'un classificateur MNIST idéalisé basé sur un histogramme de $m = 2$ bins par dimension. La précision du classificateur augmente de $\epsilon = 5\%$ (à partir de 10% et jusqu'à un

maximum de 100%) chaque fois que $k = 4$ de nouveaux points sont ajoutés à chaque bin. Quel est le plus petit nombre de points dont le classificateur aurait besoin pour atteindre une précision de 90% ?

Answer. $4 * (90 - 10) / 5 * 2^{784} = 2^{790}$.

- (e) Assumant une distribution uniforme sur tous les régions, quelle est la probabilité qu’une région en particulier est vide, en fonction de d , m et n ?

Answer. A partir de la probabilité (uniforme) que x se trouve dans la région V : $\mathbb{P}(x \in V) = \frac{1}{m^d}$, nous obtenons la probabilité que ce ne soit pas : $\mathbb{P}(x \notin V) = 1 - \frac{1}{m^d}$. Étant donné n points, la probabilité qu’aucun ne se trouve dans V est alors $\mathbb{P}(x_1 \notin V \text{ and } \dots \text{ and } x_n \notin V) = \prod_{j=1}^n \mathbb{P}(x_j \notin V) = \left(1 - \frac{1}{m^d}\right)^n$. NB qu’à mesure que m et/ou d croissent, cette probabilité devient rapidement très proche de 1, ce qui illustre “the curse of dimensionality”, ce qui signifie que nous aurions besoin d’un nombre exponentielle (en d) de points de données pour qu’une méthode d’histogramme fonctionne.

Notez le contraste entre (b) et (e): même si pour une infinité de points de données, l’histogramme sera arbitrairement précis pour estimer la vraie distribution (b), en pratique le nombre d’échantillons requis pour obtenir même un seul point de données dans chaque région croît de façon exponentielle avec d .

6. Mélange de Gaussiennes [10 points]

Soit $\mu_0, \mu_1 \in \mathbb{R}^d$, et soit Σ_0, Σ_1 deux matrices $d \times d$ positives définies (i.e. symétriques avec des valeurs propres strictement positives). On définit maintenant les deux fonctions de densités de probabilités suivantes sur \mathbb{R}^2 :

$$f_{\mu_0, \Sigma_0}(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{d}{2}} \sqrt{\det(\Sigma_0)}} e^{-\frac{1}{2}(\mathbf{x} - \mu_0)^T \Sigma_0^{-1} (\mathbf{x} - \mu_0)}$$

$$f_{\mu_1, \Sigma_1}(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{d}{2}} \sqrt{\det(\Sigma_1)}} e^{-\frac{1}{2}(\mathbf{x} - \mu_1)^T \Sigma_1^{-1} (\mathbf{x} - \mu_1)}$$

Ces fonctions de densité de probabilités correspondent à la distribution gaussienne multivariée de centre μ_0 et covariance Σ_0 , notée $\mathcal{N}_d(\mu_0, \Sigma_0)$, et à la gaussienne multivariée de centre μ_1 et covariance Σ_1 , notée $\mathcal{N}_d(\mu_1, \Sigma_1)$.

On lance maintenant une pièce équilibrée Y , et on tire une variable aléatoire X dans \mathbb{R}^d , en suivant le procédé suivant : si la pièce atterit sur pile ($Y = 0$), on tire X selon $\mathcal{N}_d(\mu_0, \Sigma_0)$, et si la pièce atterit sur face ($Y = 1$), on tire X selon $\mathcal{N}_d(\mu_1, \Sigma_1)$.

- (a) Calculez $\mathbb{P}(Y = 0|X = \mathbf{x})$, la probabilité que la pièce atterisse sur pile sachant $X = \mathbf{x} \in \mathbb{R}^2$, en fonction de $\mu_0, \mu_1, \Sigma_0, \Sigma_1$, et \mathbf{x} . Montrez toutes les étapes du calcul.

Answer.

$$\begin{aligned}
 P(Y = 0|X = \mathbf{x}) &= \frac{P(X = \mathbf{x}|Y = 0)P(Y = 0)}{P(X = \mathbf{x})} \\
 &= \frac{P(X = \mathbf{x}|Y = 0)P(Y = 0)}{P(X = \mathbf{x}|Y = 0)P(Y = 0) + P(X = \mathbf{x}|Y = 1)P(Y = 1)} \\
 &= \frac{0.5f_{\mu_0, \Sigma_0}(\mathbf{x})}{0.5f_{\mu_0, \Sigma_0}(\mathbf{x}) + 0.5f_{\mu_1, \Sigma_1}(\mathbf{x})} = \frac{f_{\mu_0, \Sigma_0}(\mathbf{x})}{f_{\mu_0, \Sigma_0}(\mathbf{x}) + f_{\mu_1, \Sigma_1}(\mathbf{x})} \\
 &= \frac{\frac{1}{\sqrt{\det(\Sigma_0)}} e^{-\frac{1}{2}(\mathbf{x}-\mu_0)^T \Sigma_0^{-1}(\mathbf{x}-\mu_0)}}{\frac{1}{\sqrt{\det(\Sigma_0)}} e^{-\frac{1}{2}(\mathbf{x}-\mu_0)^T \Sigma_0^{-1}(\mathbf{x}-\mu_0)} + \frac{1}{\sqrt{\det(\Sigma_1)}} e^{-\frac{1}{2}(\mathbf{x}-\mu_1)^T \Sigma_1^{-1}(\mathbf{x}-\mu_1)}} \\
 &= \frac{1}{1 + \frac{\sqrt{\det(\Sigma_0)}}{\sqrt{\det(\Sigma_1)}} e^{-\frac{1}{2}((\mathbf{x}-\mu_1)^T \Sigma_1^{-1}(\mathbf{x}-\mu_1) - (\mathbf{x}-\mu_0)^T \Sigma_0^{-1}(\mathbf{x}-\mu_0))}}.
 \end{aligned}$$

- (b) Rappelez-vous que le classifieur optimale Bayes est $h_{Bayes}(\mathbf{x}) = \operatorname{argmax}_{y \in \{0,1\}} \mathbb{P}(Y = y|X = \mathbf{x})$. Assumant $\Sigma_0 = \Sigma_1$, Démontrez que le classifieur optimal Bayes est linéaire en \mathbf{x} .

Answer. When the covariances are equal as in $\Sigma = \Sigma_0 = \Sigma_1$,

the expression derived at the previous point simplifies to:

$$\begin{aligned}
P(Y = 0|X = \mathbf{x}) &= \frac{1}{1 + \frac{\sqrt{\det(\Sigma_0)}}{\sqrt{\det(\Sigma_1)}} e^{-\frac{1}{2}((\mathbf{x}-\mu_1)^T \Sigma_1^{-1}(\mathbf{x}-\mu_1) - (\mathbf{x}-\mu_0)^T \Sigma_0^{-1}(\mathbf{x}-\mu_0))}} \\
&= \frac{1}{1 + \frac{\sqrt{\det(\Sigma)}}{\sqrt{\det(\Sigma)}} e^{-\frac{1}{2}((\mathbf{x}-\mu_1)^T \Sigma^{-1}(\mathbf{x}-\mu_1) - (\mathbf{x}-\mu_0)^T \Sigma^{-1}(\mathbf{x}-\mu_0))}} \\
&= \frac{1}{1 + e^{-\frac{1}{2}((\mathbf{x}-\mu_1)^T \Sigma^{-1}(\mathbf{x}-\mu_1) - (\mathbf{x}-\mu_0)^T \Sigma^{-1}(\mathbf{x}-\mu_0))}} \\
&= \frac{1}{1 + e^{-\frac{1}{2}(-\mu_1^T \Sigma^{-1} \mathbf{x} - \mathbf{x}^T \Sigma^{-1} \mu_1 + \mu_1^T \Sigma^{-1} \mu_1 + \mu_0^T \Sigma^{-1} \mathbf{x} + \mathbf{x}^T \Sigma^{-1} \mu_0 - \mu_0^T \Sigma^{-1} \mu_0)}} \\
&= \frac{1}{1 + e^{-\frac{1}{2}(2(\mu_0 - \mu_1)^T \Sigma^{-1} \mathbf{x} + \mu_1^T \Sigma^{-1} \mu_1 - \mu_0^T \Sigma^{-1} \mu_0)}} \\
&= \frac{1}{1 + e^{-(\mu_0 - \mu_1)^T \Sigma^{-1} \mathbf{x} + \frac{1}{2}(\mu_1^T \Sigma^{-1} \mu_1 - \mu_0^T \Sigma^{-1} \mu_0)}} \\
&= \frac{1}{1 + e^{-(wx+b)}} = \sigma(wx+b),
\end{aligned}$$

where we defined $w = (\mu_0 - \mu_1)^T \Sigma^{-1} \mathbf{x}$ and $b = \frac{1}{2} (\mu_1^T \Sigma^{-1} \mu_1 - \mu_0^T \Sigma^{-1} \mu_0)$ and $\sigma(\cdot)$ denotes the sigmoid (logistic) function. The resulting posterior distribution corresponds to the one of a logistic regression classifier, with a linear decision boundary which depends on the Gaussian's means and covariance matrix.