

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 .
- (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?
- (c) Donnez deux expressions équivalentes de $P(X, Y)$:
 - (i) en fonction de $\mathbb{P}(X)$ et $\mathbb{P}(Y|X)$
 - (ii) en fonction de $\mathbb{P}(Y)$ et $\mathbb{P}(X|Y)$
- (d) Prouvez le théorème de Bayes:

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

- (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?

- 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?

2. Bag of words (sac de mots) et modèle de sujet unique [15 points]

On s'intéresse à un problème de classification où 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 on suppose 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: $V = \{ \text{"goal", "kick", "congress", "vote", other} \}$.

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 = "goal"	1/100	5/1000
mot = "kick"	2/100	1/1000
mot = "congress"	1/1000	1/100
mot = "vote"	3/1000	4/100
mot = other	966/1000	944/1000

Table 1

Cette table nous dit par exemple que la probabilité qu'un mot choisi aléatoirement dans un document soit "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* ?
- (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*?
- (c) On tire aléatoirement un document du corpus. Quelle est la probabilité qu'un mot aléatoire de ce document soit "goal"?
- (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*?
- (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"?
- (f) Pour en revenir à l'apprentissage, supposons que nous ne savons 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 que nous avons un jeu de données de N documents où chaque document est annoté avec un des sujets *sport* ou *politique*. Comment estimeriez-vous les probabilités conditionnelles (e.g., $\mathbb{P}(\text{mot} = \text{"goal"} \mid \text{sujet} = \text{"politique"})$) et les probabilités des sujets (e.g., $\mathbb{P}(\text{sujet} = \text{"politique"})$) à partir de ce jeu de données?

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 $f_{\theta}(x_1)$, $f_{\theta}(x_2), \dots, f_{\theta}(x_n)$

- (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 .

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

Soit le jeu de données X_1, X_2, \dots, X_n , composé de 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 définit 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 $\theta_1, \theta_2, \dots, \theta_{N-1}$.
- (b) Écrivez le log-vraisemblance du jeu de données en fonction de $\theta_1, \theta_2, \dots, \theta_{N-1}$ et $\mu_1, \mu_2, \dots, \mu_{N-1}$.
- (c) Trouvez l'estimateur du maximum de vraisemblance de θ_j , $j \in \{1, 2, \dots, N\}$.

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 , $\mathbb{E}[\mathbb{1}_{\{x \in S\}}] = \mathbb{P}(x \in S)$.

- (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 .
- (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?
- (d) Supposons l'existence d'un classificateur MNIST idéalisé basé sur un histogramme de $m = 2$ régions par dimension. La précision du classificateur augmente de $\epsilon = 5\%$ (à partir de 10% et jusqu'à un maximum de 100%) chaque fois que $k = 4$ nouveaux points sont ajoutés à chaque région. Quel est le plus petit nombre de points dont le classificateur aurait besoin pour atteindre une précision de 90% ?
- (e) En supposant une distribution uniforme sur toutes les régions, quelle est la probabilité qu'une région en particulier est vide, en fonction de d , m et n ?

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