

MAPSI – Examen Réparti – 60 pts

Durée : 2h

*Seuls documents autorisés : Calculatrice, antisèche recto-verso,
tables de lois de probabilité*

Le barème sur 60 pts (1 pt = 2 minutes)

– Barème indicatif –

Exercice 1 (6pts) – Probabilités conditionnelles

Soit le tableau suivant, partiellement rempli, correspondant à la distribution conditionnelle $p(B|A)$.

		A	
		1	2
B	1	0.1	?
	2	0.2	?
	3	?	?

Q 1.1 (1.5pt) Compléter le tableau ci-dessus sachant que $P(B = k|A = 2) \propto k$, c'est à dire que $P(B = k|A = 2) = pk$ où p est une constante à déterminer.

		A	
		1	2
B	1	0.1	$\frac{1}{6}$
	2	0.2	$\frac{2}{6} = \frac{1}{3}$
	3	0.7	$\frac{3}{6} = \frac{1}{2}$

Q 1.2 (1.5pt) Déterminer le tableau de la loi jointe $P(A, B)$ sachant que la marginale en A est uniforme. Déterminer la marginale en B .

		A	
		1	2
B	1	0.05	$\frac{1}{12} = 0.8333$
	2	0.1	$\frac{2}{6} = \frac{1}{6} = 0.1666$
	3	0.35	$\frac{3}{6} = \frac{1}{4} = 0.25$

Marginale $P(B) = [0.1333, 0.2666, 0.6]$

Q 1.3 (1.5pt) Les deux variables A et B sont-elles indépendantes ?

NON
On voit que les lignes ne sont pas proportionnelles. En particulier, la troisième ligne n'est pas proportionnelle aux deux précédentes.

Q 1.4 (1.5pt) Soit le tableau p_{AB} obtenu à la question 1.2 (quelles que soient les valeurs à l'intérieur). Donner les lignes de code permettant de vérifier (ou pas) l'indépendance entre les variables.

On ne vérifie pas les reshape et on accepte évidemment les tests sur le max ou la somme...
-0.5 pour un test exact.

```

1 pAB_hat = pAB.sum(1).reshape(3,1) @ pAB.sum(0).reshape(1,2)
2 if np.abs(pAB - pAB_hat).max() < 1e-5:
3     print('A et B sont independantes')
4 else:
5     print('A et B NE sont PAS independantes')

```

Exercice 2 – Attente au distributeur de billets (6.5pts)

En supposant que le traitement d'une demande d'un client au distributeur est réalisé en temps constant (égal à 1 unité de temps) et que la distribution des arrivées des clients au distributeur suit une loi de Poisson de paramètre μ (unités de temps), la distribution de probabilité des durées X pendant lesquelles le distributeur est occupé suit une distribution de Borel de paramètre μ définie par :

$$P(X = n \text{ unités de temps}) = \frac{e^{-\mu n} (\mu n)^{n-1}}{n!}$$

où le paramètre μ appartient à l'intervalle $[0, 1]$. Sur une période de quelques jours, on a observé les durées d'occupation suivantes :

2	3	1	5	2	7
---	---	---	---	---	---

Q 2.1 (2.5pts) Déterminez par maximum de vraisemblance la valeur du paramètre μ .

La vraisemblance est égale à :

$$L(\mathbf{x}, \mu) = \prod_{i=1}^6 \frac{e^{-\mu x_i} (\mu x_i)^{x_i-1}}{x_i!},$$

où les x_i sont les observations de la variable X . Par conséquent,

$$\log L(\mathbf{x}, \mu) = \sum_{i=1}^6 [-\mu x_i + (x_i - 1) \log \mu + (x_i - 1) \log x_i - \log(x_i!)].$$

Donc la dérivée de la log-vraisemblance est égale à :

$$\frac{\partial \log L(\mathbf{x}, \mu)}{\partial \mu} = \sum_{i=1}^6 \left[-x_i + \frac{x_i - 1}{\mu} \right].$$

En remplaçant les x_i par leurs valeurs, on obtient :

$$\frac{\partial \log L(\mathbf{x}, \mu)}{\mu} = -20 + \frac{14}{\mu}.$$

Cette expression est égale à 0 lorsque $-20\mu + 14 = 0$, autrement dit $\mu = 14/20 = 0,7$.

Q 2.2 (4pts) Un expert de la banque affirme qu'*a priori* la distribution des paramètres μ dans la région où est situé le distributeur suit une loi Beta de paramètres $\alpha = 7$ et $\beta = 21$. On rappelle que la loi Beta est définie de la manière suivante :

$$\text{Beta}(x; \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}.$$

Où Γ désigne la loi usuelle dont la formulation est volontairement omise. En utilisant cet *a priori*, calculez par *maximum a posteriori* la valeur du paramètre μ .

Pour le *maximum a posteriori*, on calcule $L(\mathbf{x}, \mu) \times \pi(\mu)$, où π est la loi Beta indiquée ci-dessus. En passant au log, on obtient donc :

$$\log(L(\mathbf{x}, \mu) \times \pi(\mu)) = \sum_{i=1}^6 [-\mu x_i + (x_i - 1) \log \mu + (x_i - 1) \log x_i - \log(x_i!)] + \log\left(\frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)}\right) + (\alpha - 1) \log \mu + (\beta - 1) \log(1 - \mu).$$

Par conséquent, en remplaçant les x_i par leurs valeurs et en dérivant, on obtient :

$$\begin{aligned} \frac{\partial \log(L(\mathbf{x}, \mu) \times \pi(\mu))}{\mu} &= -20 + \frac{14}{\mu} + \frac{\alpha - 1}{\mu} - \frac{\beta - 1}{1 - \mu} = -20 + \frac{14}{\mu} + \frac{6}{\mu} - \frac{20}{1 - \mu} = -20 + \frac{20}{\mu} - \frac{20}{1 - \mu} \\ &= \frac{20 \times (-\mu(1 - \mu) + (1 - \mu) - \mu)}{\mu(1 - \mu)} = \frac{20}{\mu(1 - \mu)}(\mu^2 - 3\mu + 1). \end{aligned}$$

La dérivée s'annule lorsque $\mu^2 - 3\mu + 1 = 0$, autrement dit lorsque $\mu = (3 \pm \sqrt{5})/2$. Mais comme μ doit appartenir à l'intervalle $[0, 1]$, la seule solution possible est $\mu = (3 - \sqrt{5})/2 \approx 0,382$.

Exercice 3 (11 pts) – Pour quelques patates de plus

Nous nous intéressons ici à l'optimisation de la culture de pommes de terres. Pour commencer à travailler, nous avons fait un prélèvement de 100 éléments qui se répartissent comme suit en 4 catégories :

Diamètre moyen (en cm)	3	5	8	12
Nombres de pommes de terres	30	20	30	20

Q 3.1 (2.5pts) L'agriculteur considère que toutes les classes sont grossièrement équi-probable... Qu'en pensez-vous du point de vue statistique –avec un niveau de confiance 0.95– ?

Il s'agit d'un test du chi2 par rapport à la loi uniforme.

$$D^2 = \sum_i \frac{(d_i - d_i^*)^2}{d_i^*}, d_i^* = 25$$

$$D^2 = 4$$

Avec $ddl = 3$, $limite = 7.81$

On doit donc accepter l'hypothèse que cette distribution est uniforme.

Q 3.2 (2.5pts) L'année précédente, le diamètre moyen des pommes de terre était de 6.4cm sur l'ensemble de la récolte. Depuis l'agriculteur a tenté d'optimiser la production par différents moyens. Peut-on conclure que le diamètre moyen des pommes de terre a augmenté (avec une confiance de 95%, en prenant un écart type sur le diamètre de 3 cm) ?

$$\mu = 6.4cm$$

$$p\left(\frac{\bar{D} - \mu}{\sigma/\sqrt{N}} \geq z\right) = 0.05 \iff \bar{D} \geq 1.645 * \sigma/\sqrt{N} + \mu = 6.89$$

$$\bar{D}_{lim} = 6.89$$

Calcul de la moyenne = 6.7cm

NON, on peut pas être sûr qu'il s'agit bien d'une augmentation, ce n'est pas significatif

Q 3.3 L'agriculteur nous fait remarquer qu'il peut y avoir un biais sur le diamètre D selon les champs C d'où viennent les pommes de terre. L'échantillon considéré vient de 3 champs, chaque diamètre de patate correspondant la distribution suivante :

Diamètre moyen (en cm)	3	5	8	12
Champ 1	0.80	0.30	0.20	0.00
Champ 2	0.10	0.30	0.20	0.50
Champ 3	0.10	0.40	0.60	0.50

Q 3.3.1 (1pt) A quoi correspond la distribution du tableau ci-dessus par rapport aux variables aléatoires Champ C et Diamètre D ?

$$p(C|D)$$

Q 3.3.2 (2.5pts) Est-il possible de calculer le diamètre moyen des patates du champ 1 à partir des informations présentes dans les questions précédentes ? Dans l'affirmative, faites le calcul, sinon, justifier de l'impossibilité de le faire.

Oui, il faut utiliser les probas a priori de la première question

```
pCsD = np.array([[0.8, 0.3, 0.2, 0],[0.1, 0.3, 0.2, 0.5], [0.1, 0.4, 0.6, 0.5]])
pD = np.array([30,20,30,20])
pD = pD/pD.sum()

pCD = pCsD * pD
pDsC = pCD / pCD.sum(1).reshape(3,1)

print(pCD)
print(pDsC)

print(pDsC[0]@ech[0])

[[0.24 0.06 0.06 0. ]
 [0.03 0.06 0.06 0.1 ]
 [0.03 0.08 0.18 0.1 ]]
[[0.66666667 0.16666667 0.16666667 0. ]
 [0.12      0.24      0.24      0.4      ]
 [0.07692308 0.20512821 0.46153846 0.25641026]]
4.166666666666666

diam moy = 4.16 cm
```

Q 3.4 (2.5pts) Chaque pomme de terre engendre des coûts fixe de 5cts (plantation, arrosage, récolte). Elles sont ensuite vendues en moyenne selon leur calibre :

Diamètre moyen (en cm)	3	5	8	12
Prix unitaire	3cts	6cts	10cts	15cts

Est ce que l'un des champs est déficitaire ? Donner les étapes principales analytiquement puis l'application numérique.

!!!! ATTENTION effectif changé par rapport au copier-coller de sujet!!!!
 (1) il faut calculer $p(\text{diam}|\text{champ}) = p(\text{champ}|\text{diam}) * p(\text{diam})/p(\text{champ})$
 puis $E[\text{prix}|\text{champ}] = \sum_{\text{diam}} \text{prix}(\text{diam})p(\text{diam}|\text{champ})$
 avec $\text{prix}(\text{diam})$ minoré des 5ts de coût
 (2) calcul

```
# gain
prix = np.array([-2, 1, 5, 10])
print((pDsC * prix).sum(1))
# [-0.33333333  5.2          4.92307692]
```

Exercice 4 (4pts) – [CODE] Multiplication matricielle

Q 4.1 (2.5pts) Donner le code de la fonction `multmat(A,B)` qui prend en argument deux matrices et retourne la matrice résultant du produit matriciel.

Vous travaillerez obligatoirement avec des boucles `for`, après avoir extrait les dimensions des matrices d'entrées mais vous considèrerez que les dimensions des matrices sont de dimensions compatibles sans faire de test.

```
1 def multmat(A, B):
2     ligA, colA = A.shape
3     ligB, colB = B.shape
4     C = np.zeros((ligA, colB))
5     for i in range(ligA):
6         for j in range(colB):
7             for k in range(colA):
8                 C[i, j] += A[i, k] * B[k, j]
9     return C
```

Q 4.2 (1.5pt) Vous donnerez les lignes de script permettant l'invocation de la méthode et la vérification du résultat par rapport à l'opérateur `@`. On supposera A et B pré-existantes.

```
1 ligA, colA, ligB, colB = 2,3,3,5 # arbitraire
2 A = np.random.randn(ligA, colA)
3 B = np.random.randn(ligB, colB)
4
5 # seules les lignes suivantes sont demandees:
6 C = multmat(A, B)
7 if np.abs(C-A@B).max() < 1e-5:
8     print('OK')
9 else
10    print('KO')
```

Exercice 5 (14pts) – Max de vraisemblance

Un fabricant d'ampoules *basse consommation* affirme que la durée de vie de ses ampoules est de 2 ans et demi. Un échantillon de 1000 ampoules a été testé et les durées de vie constatées (en dizaines d'années) ont été reportées dans le tableau ci-dessous :

durée (dizaines d'années)	0.05	0.1	0.2	0.25	0.3	0.35
nombre d'ampoules	200	100	300	100	200	100

Q 5.1 (2.5pts) Sachant que la durée de vie d'une ampoule de l'entreprise est modélisée par une loi exponentielle de densité $p(x) = \lambda e^{-\lambda x}$, estimez le paramètre λ par maximum de vraisemblance en vous appuyant sur l'échantillon ci-dessous. Vous détaillerez vos calculs.

La vraisemblance $L(\mathbf{x}, \lambda)$ est égale à $\prod_{i=1}^{1000} p(x_i|\lambda)$, où x_i correspond aux réalisations du tableau ci-dessus. Autrement dit, $L(\mathbf{x}, \lambda) = \prod_{i=1}^{1000} \lambda e^{-\lambda x_i} = \lambda^{1000} e^{-\lambda \sum_{i=1}^{1000} x_i} = \lambda^{1000} e^{-200\lambda}$ avec $200 * 0.05 + 100 * 0.1 + 300 * 0.2 + 100 * 0.25 + 200 * 0.3 + 100 * 0.35 = 200$.
 D'où $\frac{\partial L}{\partial \lambda} = 1000\lambda^{999} e^{-200\lambda} - 200\lambda^{1000} e^{-200\lambda}$ et, donc, $\frac{\partial L}{\partial \lambda}(\lambda_{ML}^*) = 0 \iff 1000 - 200\lambda_{ML}^* = 0$. Par conséquent, $\lambda_{ML}^* = 5$.

Q 5.2 Un expert de l'entreprise intervient alors pour vous expliquer qu'il existe deux gammes de produits : la gamme basique, qui représente approximativement 60% des ventes et la gamme premium, plus résistante, qui représente 40% des ventes. Les deux gammes sont modélisables par des lois exponentielles.

Q 5.2.1 (2pts) Identifier les paramètres initiaux des deux modèles. Afin d'éviter un nouveau calcul de max de vraisemblance, vous pourrez utiliser directement le fait que lorsque X suit une loi exponentielle, $E[X] = 1/\lambda$

On va entrer dans une boucle EM... ET il faut des valeurs de paramètres initiaux (ou des affectations, mais c'est plus rare).
 $\pi_1 = 0.6, \pi_2 = 0.4$
 On divise l'échantillon en deux selon ces fractions et on en déduit :

```

1 mu1 = (0.05*200 + 0.1*100 + 0.2 *300 ) / 600
2 mu2 = (0.25*100 + 0.3*200 + 0.35 *100 ) / 400
3
4 print (mu1 ,mu2)
5 print ( 1/mu1 , 1/mu2)

```

$\mu_1 = 0.1333 \quad \mu_2 = 0.3 \quad \lambda_1 = 7.5 \quad \lambda_2 = 3.333$
 On notera la classe $C \in \{1, 2\}$ dans la suite

Q 5.2.2 (1pt) Donner les formules permettant de calculer la probabilité d'appartenance à un modèle sachant la durée de vie à partir des paramètres initiaux (i.e. la probabilité des variables cachées sachant les observations).

Note : donner la formulation sans faire les applications numériques

$$Q_i(j) = P(C = j|x_i, \lambda) \propto P(x_i|C = j, \lambda)\pi_j = \lambda_j e^{-\lambda_j x_i} \pi_j$$

Il faudra impérativement normaliser le tableau des Q de sorte que $\sum_j Q_i(j) = 1$.

Q 5.2.3 (3pts) Rappeler la formalisation du problème d'optimisation pour la mise à jour des paramètres des deux modèles en détaillant le calcul de la log-vraisemblance en fonction des $Q_i(j)$.

Calculer la valeur des nouveaux paramètres optimaux en fonction des $Q_i(j)$.

$$\pi^*, \lambda^* = \operatorname{argmax}_{\pi, \lambda} \log -\mathcal{L}$$

Avec Jensen :

$$\log -\mathcal{L} = \sum_i \log(\sum_j p(x_i, C = j)) = \sum_{i,j} Q_i(j) \log\left(\frac{P(x_i, C = j)}{Q_i(j)}\right)$$

Or [tjs dans le cas générique] :

$$\operatorname{argmax}_{\pi, \lambda} \log -\mathcal{L} = \operatorname{argmax}_{\pi, \lambda} \sum_{i,j} Q_i(j) \log(P(x_i, C = j))$$

En développant la somme sur j + remplacement de variable de π_2 [comme en TD pour éviter les contraintes lors de l'optimisation]

$$\log -\mathcal{L}' = \sum_i Q_i(1) (\log(\lambda_1) - \lambda_1 x_i + \log(\pi_1)) + Q_i(2) (\log(\lambda_2) - \lambda_2 x_i + \log(1 - \pi_1))$$

$$\frac{\partial \log -\mathcal{L}'}{\partial \pi_1} = \sum_i \frac{Q_i(1)}{\pi_1} - \frac{Q_i(2)}{1 - \pi_1} = 0$$

$$[\dots] \Leftrightarrow \pi_1 = \frac{\sum_i Q_i(1)}{\sum_i Q_i(1) + Q_i(2)}$$

$$\frac{\partial \log -\mathcal{L}'}{\partial \lambda_1} = \sum_i Q_i(1) \left(\frac{1}{\lambda_1} - x_i \right) = 0$$

$$[\dots] \Leftrightarrow \lambda_1 = \frac{\sum_i Q_i(1)}{\sum_i Q_i(1) x_i}$$

Q 5.2.4 (1.5pt) Donner l'algorithme général et proposer un critère d'arrêt.

On arrête lorsque les paramètres ne bougent plus d'une itération à l'autre :

$$(\lambda_1^{(t)} - \lambda_1^{(t+1)})^2 + (\lambda_2^{(t)} - \lambda_2^{(t+1)})^2 < \varepsilon$$

Algo :

— Tant que le critère n'est pas atteint

(E) Calculer les $Q_i(j)$ à t

(M) Optimiser les paramètres à $t + 1$: $\pi^*, \lambda^* = \operatorname{argmax}_{\pi, \lambda} \log -\mathcal{L}$

Q 5.3 (4pts) Donner une implémentation en python de l'algorithme ci-dessus. Vous préciserez en une phrase vos hypothèses sur la forme des données en entrée et vous utiliserez obligatoirement au moins deux fonctions (permettant de bien identifier les arguments dont vous avez besoin pour le calcul).

Exercice 6 (10.5 pts) – Tests divers

Q 6.1 (1.5 pt) Rappeler l'espérance d'une variable de Bernoulli B puis démontrer que la variance vaut $V[B] = p(1 - p)$.

$$E(B) = p$$

$$V(B) = p(1 - p)^2 + (1 - p)(0 - p)^2 = p(1 - p)(1 - p + p) = p(1 - p)$$

Q 6.2 (0.5 pt) *Espérance et variance d'une loi binomiale*

En remarquant qu'une variable X suivant une loi binomiale peut s'écrire comme la somme de variables X_i indépendantes suivant toutes une même loi de Bernoulli, calculer l'espérance et la variance de X .

Comme les X_i sont indépendantes,

$$E(X) = E\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n E(X_i) = \sum_{i=1}^n p = n \cdot p$$

$$V(X) = V\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n V(X_i) = \sum_{i=1}^n p \cdot (1-p) = n \cdot p \cdot (1-p)$$

Q 6.3 Compositions de groupes

Dans une université bien connue, 800 étudiants se sont inscrits à une UE de L1. Ils sont répartis aléatoirement en 200 groupes de travail de 4 étudiants. Sur ces 200 groupes :

- 13 groupes sont composés de garçons uniquement,
- 65 groupes contiennent une unique fille (et donc 3 garçons)
- 72 groupes contiennent exactement 2 filles (et donc 2 garçons),
- 35 groupes contiennent exactement 3 filles (et donc 1 garçon),
- 15 groupes contiennent exactement 4 filles (et donc pas de garçons),

Les enseignants de l'UE se demande si cette répartition est compatible avec une hypothèse de répartitions équilibrées entre filles et garçons (autant de chance pour un membre d'un groupe d'être une fille ou un garçon).

Q 6.3.1 (0.5pt) Quel type de test doit-on donc effectuer ?

Test d'ajustement.

Q 6.3.2 (1pt) Exprimer H_0 . Soit X la variable représentant le nombre de filles dans un groupe, quelle loi doit-elle suivre sous H_0 ?

H_0 : la distribution des 200 groupes est en accord avec la théorie d'équirépartition des filles et des garçons.
Sous H_0 , avec X le nombre de filles parmi 4, X suit une loi binomiale $\mathcal{B}(4, p)$

Q 6.3.3 (1.5pt) Quels seraient alors les effectifs attendus pour les groupes comprenant de 0 à 4 filles ?

Note : $\binom{4}{i} = C_4^i = [1, 4, 6, 4, 1]$ pour i allant de 0 à 4

$$f_k = 200 \cdot p(X = k) = 200 \cdot C_4^k \left(\frac{1}{2}\right)^k \left(\frac{1}{2}\right)^{n-k}$$

D'où : $f_{k \in \{0, \dots, 4\}} = (12.5, 50, 75, 50, 12.5)$

Q 6.3.4 (2.5pt) Vérifier que les conditions d'acceptation du test sont vérifiées puis vérifier si H_0 est acceptable pour un seuil de signification $\alpha = 0.05$.

Conditions d'application : échantillon supposé simple et aléatoire, fréquence de plus de 5 éléments dans chaque case du tableau de contingence.

$$\chi_{obs}^2 = \sum_{i=1}^5 \frac{(obs_i - theo_i)^2}{theo_i} = \frac{(13 - 12.5)^2}{12.5} + \frac{(65 - 50)^2}{50} + \frac{(72 - 75)^2}{75} + \frac{(35 - 50)^2}{50} + \frac{(15 - 12.5)^2}{12.5}$$

Donc $\chi_{obs}^2 = 9.64$. Or dans la table du χ^2 avec 4 degré de liberté et un $\alpha = 0.95$, on lit $c_{4;0.95} = 9.49$.

$$\chi_{obs}^2 > c_{4;0.95} \Rightarrow \text{on refuse } H_0.$$

Q 6.4 (3 pts) On se rend compte rapidement en étudiant la répartition des groupes que l'hypothèse globale d'équilibre entre les sexes n'est pas respectée dans cette UE.

Donner les calculs permettant d'estimer la répartition garçons/filles à partir de cet échantillon puis vérifier votre hypothèse H_0 modifiée par rapport à ce nouveau paramètre, toujours au même niveau de confiance. Si vous n'avez pas le temps de faire les calculs, donnez votre intuition (motivée) sur le résultat que nous allons obtenir.

$$p_G = (13 * 4 + 65 * 3 + 72 * 2 + 35) / 800 = 0.5325 \quad p_F = 0.4675$$

Effectifs théoriques : [16.1, 56.5, 74.4, 43.5, 9.5]

$$\chi_{obs}^2 = 6.91$$

On accepte maintenant H_0 !

[2/3 des points pour une bonne intuition bien motivée]

Intuition : le nouveau calcul de $p = p_F$ va augmenter les effectifs théoriques de la première moitié du tableau et réduire ceux de la seconde moitié. Le plus grand *delta* dans le premier calcul venait des valeurs (65, 35) : ce delta va diminuer et étant donné notre distance à la limite d'acceptation de H_0 , il y a toutes les chances que nous franchissions cette limite et acceptions H_0 dans cette nouvelle configuration.