

# Examen Réparti 1 - MAPSI

Durée : 2 heures

*Seuls documents autorisés : Calculatrice, antisèche recto-verso, tables de lois.  
– Barème indicatif –*

## Exercice 1 – Indépendance

Soit deux dés à six faces non pipés, un de couleur blanc et un de couleur noir. Les deux sont jetés une fois. On définit les événements suivants :

- le dé blanc donne 1, 2 ou 3.
- le dé blanc donne 2, 3 ou 6.
- la somme des deux dés est égal à 9.
- les deux dés donnent deux nombres égaux, dont la somme est inférieure à 9.

**Q 1.1 (0.5 pt)** Quel est la probabilité des ces événements ?

**Q 1.2 (1 pt)** Quels événements sont deux-à-deux indépendants ?

**Q 1.3 (0.5 pt)** Sont-ils mutuellement indépendants ? Si non, trouvez les groupes (à trois ou quatre variables) qui sont mutuellement indépendants.

## Exercice 2 – Une page top news ?

Les gérants d'un site internet de news veulent garder l'intérêt des internautes. Pour cela, ils ont pensé contrôler en temps réel les articles qui sont en première page par le nombre de clics qu'ils attirent. Si le nombre de clics par heure descend en dessous d'un certain seuil donné (par catégorie), un article n'est plus considéré comme intéressant et est enlevé de la première page pour être archivé sur le site. Pour évaluer les effets de cette règle on veut étudier statistiquement le comportement du site.

On a trois catégories d'articles (Politique, Culture, Sport). On suppose que la décroissance de l'intérêt en fonction du temps suit une loi exponentielle de paramètre  $\lambda$  dont la fonction de densité est  $f(x|\lambda) = \lambda e^{-\lambda x}$ , pour  $x \geq 0$ , et que chaque catégorie suit une loi d'intérêt différente ( $\lambda_P, \lambda_C, \lambda_S$ ).

Le tableau suivant recense le temps d'affichage en première page (en heures) pour un échantillon de 13 articles :

P	C	C	P	S	S	S	C	P	P	C	C	S
35	95	75	23	13	25	35	59	55	34	73	50	14

**Q 2.1 (1.5 pt)** On suppose que la distribution de l'attention sur les articles est effectivement une loi exponentielle. Estimer par maximum de vraisemblance la valeur de  $\lambda$  pour chaque catégorie.

**Q 2.2 (0.5 pt)** Après un certain moment, on s'est aperçu que le mécanisme de mise à jour ne permet pas de garder l'hypothèse que ces paramètres soient fixes. Partagez-vous cette idée ? Pourquoi ?

## Exercice 3 – Production des oeufs

Dans un élevage de poules pondeuses, on veut calculer le poids moyen des oeufs, en sachant qu'il suit une loi normale (en grammes)  $\mathcal{N}(\mu, 6^2)$ . Le tableau ci-dessous donne le poids de 9 oeufs pris au hasard :

43	67	52	66	46	71	65	48	64
----	----	----	----	----	----	----	----	----

**Q 3.1 (0.5 pt)** Donnez une estimation ponctuelle du poids moyen d'un oeuf dans cet élevage, et la médiane de cet échantillon.

**Q 3.2 (1.5 pt)** Le responsable de l'élevage vous dit que le poids moyen normalement est de 64 gr, mais la moyenne de cet échantillon a baissé. En utilisant un test de niveau de confiance de 95%, pouvez vous confirmer que la moyenne donnée par le responsable n'a pas baissée ?

#### Exercice 4 – Test d'indépendance

Le responsable des opérations d'une usine de production veut déterminer s'il y a des différences concernant la qualité entre deux groupes de travail qui travaillent alternativement. Il sélectionne au hasard un certain nombre de produits pour chaque groupe et les inspecte avec soin, en suivant un protocole qui classe chaque produit comme défectueux ou satisfaisant. Ce tableau donne une synthèse des données recueillies :

	Défectueux	Satisfaisant
Groupe 1	6	102
Groupe 2	5	83

Il se demande maintenant, au vu des données, si il existe une différence de qualité entre les deux groupes.

**Q 4.1 (0.25 pt)** Construisez le tableau de contingence avec les effectifs reels par catégorie et marginaux. **Q 4.2 (0.75 pt)** Calculez les effectifs attendus, en supposant que les deux variables aléatoires *groupe* et *qualité* sont indépendantes.

**Q 4.3 (1.25 pt)** Tester, au sens du chi2, si l'hypothèse d'indépendance est correcte, avec un niveau de confiance du 95 %.

#### Exercice 5 – Publicité dynamique sur les sites web

Nous nous intéressons à la création d'un système de personnalisation de la publicité affichée sur un site d'information important (comme `www.lemonde.fr` par exemple). Le site en question a établi  $K = 20$  catégories d'articles (e.g. sport, culture, économie nationale...) et utilise des cookies pour déterminer l'intérêt des utilisateurs pour telle ou telle catégorie. Ainsi, pour le site, un utilisateur  $\mathbf{x} = [0, 1, 0, \dots, 0]$  est un vecteur binaire de taille  $K$ .

La régie publicitaire du site distingue  $M = 10$  types de publicités (articles de luxe, films, ...). Pour savoir si un utilisateur est intéressé par une publicité, ils enregistrent les *clics* effectués (afin d'économiser de la place, seuls les clics sont enregistrés et pas les simples affichages).

Pendant une longue période de référence, le site utilise une politique d'affichage aléatoire et enregistre les comportements des utilisateurs. Il constitue ainsi une base de donnée  $X, Y$  de  $N$  entrées où chaque ligne désigne un couple utilisateur/catégorie de la publicité cliquée :  $\mathbf{x}_i = \{x_{i1}, \dots, x_{ik}, \dots, x_{iK}\}$ ,  $x_{ik} \in \{0, 1\}$ ,  $y_i \in \{1, \dots, M\}$ .

**Q 5.1 (0.5 pt) Modélisation.** Quelle loi utiliser pour modéliser chaque centre d'intérêt  $k$  (connaissant la catégorie de publicité  $m$ ) ? Montrer que la probabilité d'un centre d'intérêt  $k$  pour un utilisateur  $i$  peut s'écrire :

$$p(x_{ik}|y_m) = \theta_{k,m}^{x_{ik}} (1 - \theta_{k,m})^{(1-x_{ik})}$$

Pour un utilisateur, nous faisons l'hypothèse que ses centres d'intérêt sont indépendants les uns des autres. Exprimer  $p(\mathbf{x}_i|y_m)$  en fonction des  $\theta_{k,m}$ .

**Q 5.2 (0.5 pt)** Donner le protocole pour créer une politique d'affichage adaptée aux centres d'intérêt de l'utilisateur à partir de ces données. Donner le nombre de paramètres à apprendre pour entraîner notre modèle.

**Q 5.3 (1.5 pt) Apprentissage.** Donner l'expression des paramètres maximisant la vraisemblance des données pour chaque catégorie de publicité. Formaliser le problème et donner les détails des calculs.

**Q 5.4 (1 pt) Inférence.** Pour un  $\mathbf{x}$  donné, quel type de publicité la régie doit-elle afficher pour maximiser le taux de clic (ie, au sens du maximum a posteriori) ?

Donner la formulation détaillée du problème à poser et une estimation des paramètres manquants.

La rapidité du système est un facteur clé, le chargement des pages devant rester (quasi-)instantané... Comment jugez-vous l'efficacité du système proposé ?

**Q 5.5 (0.5 pt) Amélioration du système.** Afin de palier le risque de non-affichage de certaines catégories de publicité, nous utilisons les capacités génératives de notre modèle. Expliquer sur quel critère et comment tirer des publicités reflétant l'ensemble des goûts de l'utilisateur et pas seulement son meilleur score d'affinité.

**Q 5.6 (0.5 pt)** L'évaluation de ce type d'approche est délicate, nous procédons en général par A-B *testing*, en comparant 2 systèmes mis en ligne en parallèle. Sur quel critère départager les deux systèmes ? Décrire brièvement l'implémentation de l'A-B *testing*.

**Q 5.7 (0.25pt)** Quelle entreprise française a fait fortune grâce à ce type d'approche ?

## Exercice 6 – Modélisation des usagers RATP par des lois normales

Nous nous intéressons à la modélisation des usagers du métro parisien pour les départs au travail, sur la plage [6h 10h]. Les logs d'une population contiennent simplement les temps auxquels ont eu lieu les validations des usagers. Ils ont été constitués sur 13 semaines et le temps a été discrétisé par minutes (soit 240 intervalles entre 6h et 10h).

Nous modélisons cette situation à l'aide d'une loi normale  $\mathcal{N}(\mu, \sigma)$ , en prenant un temps discrétisé  $t \in [0, 239]$

$$p(t) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(t - \mu)^2\right)$$

Note : pour simplifier, nous faisons l'hypothèse que la loi est toujours normalisée (même si  $\mu$  se rapproche d'une des bornes et que la distribution est en réalité tronquée).

**Q 6.1 (2 pts)** Il y a deux façons de voir un usager :

1. Comme un ensemble de  $N$  validations :  $u = \{t_1, \dots, t_N\}$ ,  $t \in [0, 239]$
2. Comme le comptage de ces validations par créneau horaire :  $u = \{c_0, \dots, c_t, \dots, c_{239}\}$ ,  $c_t \in \mathbb{N}$

Un usager avec une validation unique à  $t = 10$  est donc représenté par :  $u = \{t_1 = 10\}$  ou par  $u = \{c_0 = 0, c_1 = 0, \dots, c_9 = 0, c_{10} = 1, c_{11} = 0, \dots, c_{239} = 0\}$ .

Exprimer la vraisemblance d'un usager en explicitant l'hypothèse que vous faites pour simplifier le calcul. Vous exprimerez la vraisemblance en fonction des  $c_t$ . Passer ensuite la formule au log et développer l'expression.

**Q 6.2 (2 pts)** Pour distinguer les comptages de 2 utilisateurs  $u$  et  $u'$  au pas de temps  $t$ , nous introduisons un nouvel index :  $c_{u,t}$  (qui est donc différent de  $c_{u',t}$ ). Trouver les paramètres optimaux de la loi normale pour un ensemble d'usager au sens du maximum de vraisemblance en fonction des  $c_{u,t}$ .

NB : donner la formulation du problème et les détails du calcul.

**Q 6.3 (0.5 pt)** A partir des données précédentes, un expert a constitué 3 ensembles d'utilisateurs correspondant à des classes distinctes : les *lève-tôt*, les *tardifs*, les *variables*. Donner le protocole expérimental pour créer un classifieur d'usager, au sens du maximum de vraisemblance, capable de ranger un nouvel utilisateur (identifié par ses logs) dans l'une des catégories.

Donner la formulation du critère de décision.

**Q 6.4 (0.5 pt)** Comment construire un classifieur au sens du maximum a posteriori à partir de ces mêmes données.

**Q 6.5 (0.5 pt)** Est-il possible que notre modèle de classification de voyageur soit en désaccord avec certaines étiquettes proposées par l'expert sur le jeu de données original ? Dans l'affirmative, comment faire évoluer le modèle ?