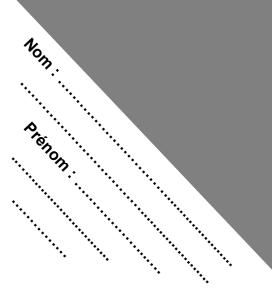
Numéro d'étudiant :



Université de Bordeaux - UF Mathématiques et Interactions

M2 Modélisation statistique et stochastique (4TSS901S, S9) M2 CMI Statistique et informatique (4TST901S, S9)

« Statistique et analyse de données en grande dimension »

CONTROLE CONTINU

27/10/2022 14h-17h

Enseignante: Marta Avalos-Fernandez

Déroulement du devoir :

- Devoir individuel. Vous n'avez pas le droit de communiquer entre vous.
- Vous pouvez utiliser votre ordinateur portable ou un ordinateur de la salle.
- Vous avez droit à tous vos documents et codes, ainsi que à internet. Vous pouvez utiliser un traducteur anglais/français pour l'article, si vous en avez besoin. Le lien vers l'article associé au sujet vous a été envoyé par mail le 19/10/2022.
- Vous pouvez répondre aux questions en anglais ou en français, à la suite des questions (l'espace laissé pour les réponses n'est pas proportionnel à la longueur de réponse attendue) ou sur un document vierge (en précisant le numéro de la question) ou sur papier libre (que vous devrez ensuite scanner avec une application sur téléphone portable). Le fichier doit être transformé au <u>format pdf</u> avant l'envoi, si plusieurs documents ont été créés, ils devront être fusionnés avant l'envoi.
- L'ensemble de réponses de la Partie I sont à envoyer par mail à <u>marta.avalos-fernandez@inria.fr</u> dans un seul fichier pdf : **devoir271022_ nometudiant.pdf** <u>avant 15h15</u>.
- Vous recevrez pour 15h la Partie II du devoir. Vous pouvez y travailler dès que vous avez envoyé votre Partie I.

Partie I

Ce devoir s'appuie sur des extraits de l'article Shah et al. Blood-Based Fingerprint of Cardiorespiratory Fitness and Long-Term Health Outcomes in Young Adulthood. *Journal of the American Heart Association*, 2022;11:e026670.

La forme cardio-vasculaire (CRF) est un bon prédicteur de l'état de santé. La CRF est mesurée par le pic de la consommation d'oxygène (peak VO2) à l'aide d'un test d'effort cardiopulmonaire. Cette mesure est difficile d'obtenir en routine et donc la CRF est sous-utilisée en prévention primaire. La CRF devrait toutefois pouvoir être prédite sans l'utilisation d'un test d'exercice à partir de facteurs de risque cliniques et de biomarqueurs (ou métabolites) sanguins. L'objectif du travail décrit dans l'article était d'obtenir une bonne prédiction du pic de la consommation d'oxygène à partir de facteurs de risque démographiques et cliniques (âge, sexe, indice de masse corporelle, pression artérielle systolique, statut tabagique, cholestérol total, diabètes, antécédent de maladie cardio-vasculaire, ...) et de métabolites sanguins en utilisant les données Framingham Heart Study (FHS). A partir de cette prédiction, les auteurs ont construit un score et l'ont évalué sur les données d'une autre étude.

L'objectif du devoir est de repérer dans vos cours les éléments qui permettent de comprendre la démarche des auteurs (pourquoi, comment) et d'avoir un regard critique vis-à-vis de cette démarche, de leurs résultats et de leurs conclusions et cela sans forcément avoir des connaissances du domaine d'application.

MÉTHODES-STATISTICAL ANALYSIS-DEVELOPMENT OF A METABOLOMIC FITNESS SCORE IN THE FHS

- 1) Il est indiqué : "LASSO regression was selected to
 - avoid overfitting,
 - address collinearity across metabolites, and
 - provide a parsimonious multimetabolite model for peak VO2."

Expliquez avez vos propres mots ces arguments pour utiliser la méthode Lasso.

2) Il est indiqué: "Our aim was to develop a metabolite score independent of standard risk factors, so we forced (unpenalized) adjustment for age, sex, body mass index (BMI), systolic blood pressure, hypertension treatment status, total cholesterol, high-density lipoprotein cholesterol, prevalent CVD, diabetes, and smoking status in the LASSO model."

_ ''	1 /		
Expliquez ce parag	ranne (nour	anoi et com	menti
ENDINGUEL CE PUI US	i apric (pour	gaoi et com	11101107.

3)	Il est indiqué: "Continuous covariates and peak VO2 were log-transformed; forced adjustments in LASSO were not standardized." Pourquoi, à votre avis,
	 les variables explicatives continues ont été log-transformées ? la variable réponse, peak VO2, a été log-transformée ? les variables d'ajustement ne sont pas standardisées ?
	Est-ce que, à votre avis, les métabolites sont standardisés ? Pourquoi ?
4)	Détaillez le modèle qui a été utilisé par les auteurs ainsi que les hypothèses sur lesquelles repose ce modèle.
	Donnez ensuite le problème d'optimisation qui est résolu en utilisant la méthode Lasso (tenez en compte les choix effectués par les auteurs mentionnés dans les questions précédentes).
	Soyez rigoureux dans les notations mathématiques.
D.ź.	Davis Construction Construction of Management France Construction File
KE	SULTATS- DEVELOPMENT AND CHARACTERIZATION OF A METABOLOMIC FITNESS SCORE IN THE FHS
5)	La figure 2A montre les coefficients de régression (d'une part) des métabolites estimés avec la méthode Lasso.
	Combien de métabolites (nombre total de variables explicatives mesurées, moins les facteurs de risque démographiques et cliniques) pourraient être représentés ?
	A votre avis, pourquoi seulement ~70 métabolites sont présentés dans la figure 2A ?
6)	Il n'est pas précisé comment ces estimations de la figure 2A sont obtenues. Indiquez une méthode avec laquelle ces estimations sont obtenues usuellement. Expliquez brièvement le raisonnement de cette méthode.

7)	Les auteurs ont utilisé une partie de la cohorte FHS pour estimer les coefficients de régression (derivation subsample N=451) et la partie restante (validation subsample N=914) pour évaluer le coefficient de corrélation entre la partie du modèle correspondant aux métabolites (estimée à partir de la derivation subsample, soit le score que les auteurs développent) et la réponse. Pourquoi, à votre avis, les auteurs ont suivi cette démarche ? Pourquoi, à votre avis, le coefficient de corrélation obtenu avec la derivation subsamble est meilleur que celui obtenu avec la validation subsample (figure 2B) ?				
8)	La table 2 montre les résultats du modèle linéaire estimé en utilisant la méthode des moindres carrés (sans sélection de variables) à partir de la validation subsample. Onze variables sont utilisées : 10 variables d'ajustement et le score des métabolites (dont les coefficients ont été estimés à partir de la derivation subsample). Donnez l'interprétation de $\hat{\beta}$ associé à Smoking ($\hat{\beta}$ =-0.193), sachant que Smoking=1 si fumeur ou ancien fumeur et Smoking=0 pour ceux qui n'ont jamais fumé, et celui associé à Metabolomic fitness score ($\hat{\beta}$ =0.303). Pour cela,				
	 Soit vous tenez en compte dans votre interprétation de la transformation de la réponse. Soit si vous préférez (car plus facile) supposez que la réponse peak VO2 n'a pas été transformée, et qu'elle est mesurée en mL/kg per min. 				