

- Je suis face à un problème de régression avec $n = 10000$ individus et $p = 100$ variables hétérogènes (quantitatives et qualitatives), et j'aimerais savoir quelles sont les variables explicatives pertinentes :

- ☐ je réalise une régression linéaire avec sélection de variables backward
- ☐ je discrétise les variables quantitatives en variables catégorielles
- ☐ j'utilise la méthode de régression des forêts aléatoires
- ☐ je transforme les variables catégorielles en variables linéaires
- ☐ je réalise une régression PLS
- ☐ je réalise une régression LASSO
- ☐ je réalise une régression ridge
- ☐ je réalise une régression avec sélection de variables forward.
- ☐ je commence par trier les variables en enlevant toutes celles non significativement corrélées avec la réponse à prédire.

- Je dispose d'une base de données d'imagerie médicale, pour lesquelles je sais dans quelles images il y a une tumeur. Je veux construire un modèle de scoring qui m'indique la probabilité de présence d'une tumeur.
Je peux utiliser

- ☐ une régression linéaire
- ☐ une régression logistique
- ☐ une forêt aléatoire
- ☐ un SVM
- ☐ un KPPV
- ☐ un réseau de neurones

- Je dispose d'une base de données génomique, avec 10 000 gènes, et pour chaque individu la connaissance de si ou non ils ont une certaine maladie d'intérêt. Pour détecter les gènes d'intérêt, je peux réaliser :

- ☐ une régression linéaire
- ☐ une régression logistique
- ☐ une régression logistique avec pénalité Lasso
- ☐ un SVM
- ☐ un Random Forest
- ☐ une Artificial Neural Network

- Exercice : sur les données breast tumors, pour laquelle vous avez un défaut de connection :
 - tester une autre technique d'imputation
 - évaluer par validation croisée
 - tester une régression logistique