

# Les types de variables

- variable qualitative
- variable ordinales
- variables quantitatives (discrètes : thé & café)
- variables quantitatives (continues : taille, poids et âge)

## Quand les données (quantitatives) ne sont pas comparables

- Centrage et réduction => favorise les algo qui utilisent la descente de gradient. Pour pouvoir comparer les choux et les carottes, il faut centrer et réduire les variables quantitatives
- Analyse en composantes principales, ne peut être faite que sur les variables quantitatives. Pour faire des traitements, il va devoir calculer des moyennes, médianes, écarts types, mais cela n'a aucun sens mathématique sur les variables qualitatives.

## Quand les variables sont qualitatives

(moyen de paiement - état de santé)

- Donner à la place de chaque modalité une valeur numérique. Jour de la semaine (1 - 7). Les valeurs n'ont aucun impact, ce sont des modalités, (1 - 7) == (15 - 22).
- On peut transformer ces modalités en bits, chaque modalité devient une colonne Samedi (0 / 1)

## Clustering

comparer un individu avec un autre => clustering, call non supervisé

Non supervisé, j'ai un jeu de données et je n'ai aucune idée de comment comparer ces individus. C'est une première étape pour classer un jeu de données précis.

## Classification

déterminer qu'un indiv. appartient à une classe => classification

Ici on a des données par le métier, on a une colonne qui devra être déterminée, c'est ce que l'individu appartient à la classe A ou B, malade pas malade, toxique, comestible.

## Régression

déterminer les caractéristiques d'un individu => régression

La colonne à déterminer est réelle avec une infinité de valeurs, elle n'est pas sur deux modalités.

Classification => variable qualitative

Régression => variable quantitative

## Analyse en composantes principales

C'est une transformation linéaire, il n'y a pas de déformation des axes. l'ACP calcule la valeur par laquelle on va devoir multiplier chaque axe pour réaliser cette opération de transformation linéaire.

Quand le jeu de données n'est pas linéaire, ce n'est pas applicable. D'en le pourcentage d'inertie, s'il n'y a pas d'éboulement, cela signifie que l'analyse en composant principal n'a pas d'impact. Il faut observer des coudes dans l'éboulement des valeurs propres.

# Algorithmes supervisés et non supervisés

*e.g. donnez moi le nom d'un algorithme supervisés et dites moi comment ça fonctionne.*

## Le clustering (non supervisé)

Hiérarchique ascendant & kmeans :

On se base sur la distance entre les différents individus.

### Classification Hiérarchique ascendant

#### Aglomératif

On recherche la distance la plus petite entre les individus plus les groupes.

**La hauteur des branches représente la différence entre les différents clusters.**

*Cette arbre n'a qu'un rôle, permettre de définir un découpage en nombre de classes. Il est donc possible de couper cet arbre à n'importe quelle hauteur et de trouver un nombre de classe. Cela sert à déterminer le nombre de classes qu'il est intéressant d'étudier.*

*E.g. découper toutes la France en deux (ville / météo) n'est pas pertinent.*

#### Divisif

On fait l'inverse, on part d'un groupe d'individu avec une forte hétérogénéité et on essaye de trouver des feuilles et tous les différencier. Pour le clustering cette algorithme n'a jamais été mis en place mais pour la classification oui.

Dans le cas du cancer du sein si une des colonnes de mon tableau peut séparer les individus en malade / pas malade c'est gagné.

#### Inconvénient

Il est impossible d'utiliser ce type d'algorithme sur des grands jeux de données. Il serait trop long à sortir un résultat.

### K-means (Lloyd)

- définir aléatoirement deux points
- calculer les distances des individus entre ces deux points
- déplacer ces points dans les barycentres de ces classes.

Il existe en deux versions soit on commence en lui donnant 2 individus, soit 2 points aléatoires. L'algorithme s'arrête lorsqu'il n'y a plus de variation entre deux itérations, il est possible que l'agglomération n'arrive pas à faire trois catégories et tourne en boucle.

On détermine le nombre de classes en :

- on fait toutes les classifications possibles
- on utilise l'agglomération silhouette pour déterminer le quel est le meilleur
  - pour chaque individu, on calcule la distance entre lui-même contre tout les individus de sa classe, puis de lui-même contre tout les individus d'un autre groupe, on fait la moyenne et on regarde pour un autre groupe.
  - est-ce que même point n'était pas mieux dans la classe X plutôt que Y ?
  - Où se trouve la distance minimale, entre les individus de son groupe ou d'un autre groupe ?

*Es ce que la silhouette fonctionne dans tous les cas ? Non quand il y a des groupe qui entoure un autre (des cercles imbriquer les un dans les autres) la silhouette ne peut pas fonctionner, car un point à une extrémité du cercle sera forcément plus loins de l'autre extrémité du cercle que du cercle étant situé à l'intérieur de lui.*



K-means peut aussi être utilisé pour diminuer le nombre de valeur, diminuer la complexité des données. En faisant des clustering, je peut rajouter des individus qui vont être dans un cluster particulier. Je peux ainsi rajouter artificiellement le nombre d'individu, je ne crée pas d'information mais on peut avoir plus de valeur ce qui permet au algo de classification de mieux se comporter par la suite.

On en fait k-means que sur les colonnes centré réduite.

## DB Scan

On donne le rayon de la boule & combien d'individu on veut avoir dans cette boule pour que l'individu qui a crée cette boule fasse partie de la même classe. On a A plus 3 individu et l'algo s'arrête avec C car il n'a qu'un seul individu proche.

- A la fin il reste des individus qui ne sont dans aucun groupe, qui représente les cas extrême.
- Le calcul de ce rayon n'est jamais simple, une faible variation entre le nb d'individu proche et le rayon du cercle peut donner des résultats très différents



## Clustering hiérarchique ascendant

C'est la création d'un arbre. Tous les individus sont uniques et on essaye de faire des groupes de plus en plus grands jusqu'à n'en avoir plus qu'un seul cluster. Si on joue de données à N individus, on peut avoir N-1 classes. La hauteur des branches entre deux groupes symbolise leurs éloignements.

## Les plus proches voisins (supervisé)

On continue de faire rechercher les proximités mais on donne un nombre de voisins à trouver pour aiguiller la classification.

## Apprentissage

Il faut découper notre jeu de données en apprentissage & test pour contrôler les prédictions de notre algorithme.

**La complexité** permet de découper les groupes par rapport à notre jeu de données de manière de plus en plus précise (une ligne ou un dessin qui englobe l'intégralité des points). En augmentant la complexité de l'algo, on peut avoir un taux de réussite de 100% sur notre jeu de données.

**généralisation** Cependant, notre modèle une fois entraîné doit rester général. Il ne doit pas servir juste à prédire notre jeu de données initial. Si la complexité est trop élevée, notre algorithme ne sera plus général et sera utilisable que pour son jeu d'entraînement.

Surapprentissage => faible biais, variance élevée

Sous-apprentissage => biais élevé, variance faible

## L'échantillonnage

**stratification** => permet de s'assurer que le découpage de données soit équilibré. Il faut faire attention, car la variable a stratifier doit être qualitative.

**prédiction croisé** => Comment faire pour faire une prédiction à partir de 5 modèle. les algo donne deux métriques la réponse (O/1)+ la probabilité. Il faut donc additionner les probabilités données par chaque modèle.

# Choix des métrique de performance

## Classification (qualitatif)

### Matrice de confusion

Cette matrice permet d'analyser les performances d'un modèle de classification en regardant où il fait des erreurs et en évaluant la pertinence des prédictions



Exemple d'application :

- Diagnostic médical
- Détection de spam
- Reconnaissance faciale

## La courbe Roc

### Définition

La courbe ROC (Receiver Operating Characteristic) est un outil permettant d'évaluer la performance d'un modèle de classification binaire.

☐ Pourquoi l'utiliser ?

Elle mesure la capacité du modèle à bien classer les classes positives et négatives. Elle permet d'analyser l'impact du seuil de décision sur les performances du modèle. L'aire sous la courbe (AUC - Area Under Curve) donne un indicateur global de la qualité du modèle

☐ Les Axes de la Courbe ROC

Axe	Explication
<b>Axe Y (ordonnée) = Sensibilité (Recall)</b>	Proportion de <b>vrais positifs détectés</b> par rapport à tous les positifs réels.
<b>Axe X (abscisse) = Taux de faux positifs (1 - Spécificité)</b>	Proportion de <b>faux positifs détectés</b> par rapport à tous les négatifs réels.

☐ Interprétation d'un Point sur la Courbe ROC

- Chaque point de la courbe ROC correspond à un seuil de classification différent.
- Le seuil de probabilité permet de décider à partir de quelle probabilité on classe une observation comme "positive".
- Si on diminue le seuil, on détecte plus de vrais positifs mais aussi plus de faux positifs.
- Si on augmente le seuil, on réduit les faux positifs, mais on risque de manquer des vrais positifs. Exemple ☐ :

Seuil = 50% → Le modèle classe "positif" si la probabilité est  $\geq 50\%$ . Seuil = 70% → On devient plus strict, donc moins de faux positifs mais plus de faux négatifs.

☐ L'Aire Sous la Courbe (AUC - Area Under Curve)

L'AUC (Area Under the Curve) représente la probabilité qu'un modèle classe un exemple positif avant un négatif.

- AUC = 1 → Modèle parfait (sépare parfaitement les classes).
- AUC = 0.5 → Modèle aléatoire (aucune capacité de classification).

- $AUC < 0.5 \rightarrow$  Modèle inversé (prédit le contraire de la vérité).

# Regression (quantitatif)

## Corrélation de Pearson (quantitatif)

Mesure la relation linéaire entre deux variables numériques.

Lorsqu'on entraîne un modèle de régression, on cherche à prédire une variable continue (ex: prix d'une maison, température, chiffre d'affaires, etc.)

La corrélation de Pearson ( $r$ ) mesure le degré de linéarité entre les prédictions et les valeurs réelles :

- Si  $r$  est proche de 1  $\rightarrow$  Le modèle suit bien la tendance des données réelles.
- Si  $r$  est proche de 0  $\rightarrow$  Aucune corrélation linéaire entre les prédictions et les valeurs réelles (le modèle est mauvais).
- Si  $r$  est négatif  $\rightarrow$  Les prédictions sont à l'opposé des valeurs réelles (erreur systématique).

En complément des erreurs absolues (comme RMSE et MAE), la corrélation de Pearson aide à voir si le modèle suit la bonne dynamique générale.

## RSS (Residual Sum of Squares)

Le résidu  $rpz$  l'erreur pour chacune des prédictions.

C'est la somme des erreurs élevées au carré. Plus le RSS est faible, plus le modèle ajuste bien les données.

Ne permet pas de comparer des modèles entre différents jeux de données, car sa valeur dépend de la taille de l'échantillon.

## MSE (Mean Squared Error)

□ Comparaison des performances d'un modèle entre différents jeux de test ou algorithmes.

□ Permet de choisir la meilleure optimisation ou le meilleur kernel (par exemple en SVM ou régression).

□ Problème : l'unité du MSE n'est pas la même que celle de la variable cible  $Y$ .

## RMSE (Root Mean Squared Error)

□ Permet d'avoir une erreur interprétable car elle est dans la même unité que  $Y$ .

□ Très utile pour comparer des modèles de régression et voir l'ampleur des erreurs en unités réelles.

□ Ne fonctionne PAS pour des prédictions qualitatives ! (À corriger dans tes notes)  $\rightarrow$  Il est exclusivement utilisé en régression et non en classification.

# Kernel

Un kernel est une fonction mathématique qui transforme les données pour les rendre séparables linéairement dans un espace de plus haute dimension.

□ Utilisation principale : Il permet d'utiliser des algorithmes linéaires (comme SVM ou régression) sur des données non linéaires, en évitant de calculer explicitement toutes les dimensions.

⚠ Si on choisit le mauvais kernel, la transformation peut rendre les données encore plus complexes, ce qui dégrade la performance du modèle.

□ Cela permet d'économiser du temps de calcul et de conserver des algorithmes optimisés tout en traitant des problèmes complexes.

Concept	Explication
Définition	Le kernel transforme les données pour les rendre <b>séparables linéairement</b>
Pourquoi ?	Quand les données sont <b>non linéaires</b> , il permet d'utiliser un <b>modèle linéaire</b> dans un <b>espace transformé</b>
Exemples de Kernel	<b>Linéaire</b> (pas de transformation) <b>Polynomial</b> (tordre l'espace en puissance) <b>Radial (RBF)</b> (séparer des données circulaires)
Avantage	Évite de <b>calculer toutes les dimensions</b> , permet d'utiliser des <b>algorithmes linéaires</b> sur des <b>données complexes</b>



# La régression linéaire

## Définition

La régression est une méthode de prédiction utilisée en Machine Learning pour estimer une valeur numérique en fonction d'autres variables.

- ☐ Objectif : Trouver une relation entre une ou plusieurs variables d'entrée (features) et une variable cible.

## ☐ Quand utiliser la régression linéaire ?

Cas d'utilisation	Explication
<b>Prédire une valeur continue</b>	Prix d'une maison, salaire, température.
<b>Relation linéaire entre les variables</b>	Si augmenter X augmente Y de manière proportionnelle.
<b>Peu de données, besoin d'un modèle simple</b>	Modèle rapide et interprétable.
<b>Comprendre quelles variables influencent Y</b>	Identifier les variables importantes avant d'utiliser un modèle plus complexe.

- ☐ Elle optimise les autres algorithmes en aidant à :

- ✓ Sélectionner les meilleures variables.
- ✓ Tester si une relation linéaire existe avant d'appliquer un modèle complexe.
- ✓ Servir de base mathématique pour d'autres modèles (logistique, SVM)

## La fonction de coût

La fonction de coût est un outil utilisé en Machine Learning pour mesurer l'erreur d'un modèle.

- ☐ Elle sert à évaluer à quel point les prédictions du modèle sont proches ou éloignées des vraies valeurs.

- ☐ Le modèle fait une prédiction.
- ☐ On compare cette prédiction avec la valeur réelle.
- ☐ La fonction de coût attribue une "pénalité" en fonction de l'écart entre les deux.
- ☐ Le modèle ajuste ses paramètres pour réduire cette erreur et améliorer ses futures prédictions.

Concept	Explication
<b>Définition</b>	Mesure l'erreur entre la prédiction du modèle et la valeur réelle.

Concept	Explication
Utilité	Permet d'entraîner le modèle et d'améliorer ses performances.
Fonctionnement	Compare la prédiction avec la vraie valeur et attribue une pénalité (erreur).
Optimisation	Les algorithmes comme la <b>descente de gradient</b> minimisent cette erreur.
Types	Régression : MSE (erreur quadratique) / Classification : Log-Loss (entropie croisée).

## La régression logistique

La régression logistique est un algorithme de classification qui permet de prédire une catégorie à partir de données. Contrairement à la régression linéaire (qui prédit une valeur continue), la régression logistique prédit une probabilité et attribue une classe (ex: "Oui" ou "Non", "Spam" ou "Non-Spam", etc.).

☐ Exemples concrets :

- ✓ Détecter un spam : Est-ce qu'un email est "Spam" ou "Non Spam" ?
- ✓ Diagnostic médical : Un patient est-il malade (1) ou en bonne santé (0) ?
- ✓ Prédiction d'achat : Un utilisateur va-t-il acheter un produit (Oui/Non) ?

☐ Avantages :

- ✓ Simple à comprendre et rapide à entraîner.
- ✓ Fournit une probabilité, utile pour ajuster les seuils.
- ✓ Fonctionne bien sur des petits jeux de données.

☐ Limites :

- ☐ Ne fonctionne que pour deux catégories (mais il existe des extensions pour plusieurs classes).
- ☐ Suppose que les données sont bien séparables (si ce n'est pas le cas, un SVM ou un réseau de neurones est plus adapté).

## l'arbre de décision

Cette algo control colone par colone la capacité de chaque colone à séparer le jeu de données vis à vis de la colone à définir (e.g. espèce).

On décrit jusqu'à quel niveau on va aller en terme de profondeur, je ne veut pas que mon arbre aille au dela de deux niveau de profondeur.

prb :

- il fait trop de coupage à profondeur max pour ne trouver qu'un exemple à la fin
- A chaque étape de l'analyse l'algo ne prend qu'une colone à la fois, les informations de sont pas brassées à l'intérieur d'une seule colone. En traitant variable après variable, il est problématique de transformer ses données en modalité pour les données qualitative.

Réduire le sur apprentissage :

- réduire la profondeur de l'arbre
- augmenter le nombre d'arbre
  - (tree bagging / random forest)
  - random forest crée des arbres en parallèle (c'est rapide)
  - boosting création des arbres en série => plus lent mais plus pertinent
  - vecteur machine

## Présentation

2/3 min

- Expliquer les grandes lignes de la recherche
- Les étapes traversées
- Dire si on a validé / invalidé notre hypothèse

(avoir un plan de ce que l'on va présenter)

Ne pas montrer tout les graphiques, parler de l'analyse, les axes de recherche, nos conclusion. Montrer quelques graph intermédiaire.

# Révision

## typologie d'algorithme

### Non supervisé

Quelle sont les ama d'individu qui vont ensemble.

- Il s'utilise quand on a des jeux de données que l'on veut libélé, mettre en cluster pour faire de l'apprentissage supervisé par la suite.
- Quand on veut connaître une variable qui n'est pas binaire (prix de l'appart)
- Je suis un opérateur internet, je regarde ce que mes clients achète, je propose des publicités ciblé en rapprochant deux individu.

Une fois que j'ai rapprocher / libélé ces informations je peux transporté ce jeu de données dans un algorithme supervisé et rapidement déterminé sa classe.

### Supervisé

Il y a une colonne X ou Y a expliquer.

Il y deux type d'algo :

- **Régression** Y est *quantitatif*
- **Classification** binaire Y est *qualitatif*

Une modalité c'est la représentation de toute les valeur possible d'une valeur qualitative par Vrai / Faux. En somme on crée une colonne par valeur possible (vert / jaune / rouge) et chaque colonne répond par Vrai ou Faux. Chaque colonne est une modalité.

