

Projet 6 — 24 mai 2021



Détection de faux billets

Algorithme de classification avec régression logistique

→ 17 faux billets d'euros pour 1 million d'authentiques

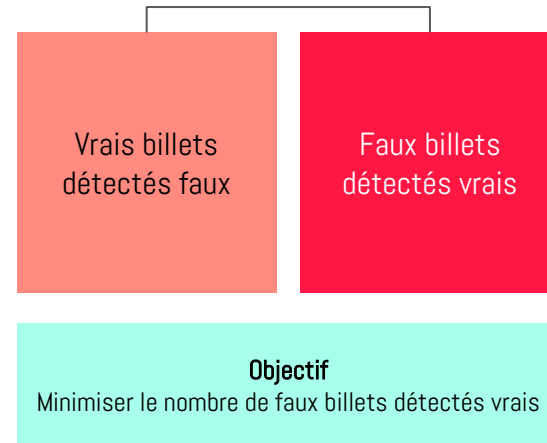
- ◆ Dans la zone euro
- ◆ Niveau historiquement bas, en baisse chaque année
- ◆ 80% des contrefaçons sont des coupures de 20 et 50€
- ◆ 97% saisis dans la zone euro

→ 30 à 40% émis depuis la France

→ Caractéristiques discriminantes

- ◆ Impression sur coton pur
- ◆ Impressions en relief
- ◆ Mesures géométriques (adapté au machine learning)

2 types d'erreur de classification



Sources :

[National Geographic](#)

[Banque centrale européenne](#)

1. Exploration

- Jeu de données
- Vrais et faux billets
- Distributions par authenticité
- Impact des variables
- Analyse des corrélations

2. ACP

- Objectifs de l'ACP
- Standardisation
- Axes d'inertie
- Covariance
- Analyse des éboulis
- Représentation en 2D
- Cercle des corrélations

3. Clustering

- K-means
- Clusters
- Matrice de confusion

4. Modélisation

- Régression linéaire
- Régression logistique
- Split des données
- Classifieur idiot
- Résultats de la modélisation
- Contributions au modèle
- Programme de détection

Exploration

01

170 billets, 7 variables

| | is_genuine | diagonal | height_left | height_right | margin_low | margin_up | length |
|-------|------------|-------------|--------------|--------------|------------|------------|--------|
| 0 | True | 171.81 | 104.86 | 104.95 | 4.52 | 2.89 | 112.83 |
| 1 | True | 171.81 | 104.86 | 104.95 | 4.52 | 2.89 | 112.83 |
| | diagonal | height_left | height_right | margin_low | margin_up | length | |
| count | 170.000000 | 170.000000 | 170.000000 | 170.000000 | 170.000000 | 170.000000 | |
| mean | 171.940588 | 104.066353 | 103.928118 | 4.612118 | 3.170412 | 112.570412 | |
| std | 0.305768 | 0.298185 | 0.330980 | 0.702103 | 0.236361 | 0.924448 | |
| min | 171.040000 | 103.230000 | 103.140000 | 3.540000 | 2.270000 | 109.970000 | |
| 25% | 171.730000 | 103.842500 | 103.690000 | 4.050000 | 3.012500 | 111.855000 | |
| 50% | 171.945000 | 104.055000 | 103.950000 | 4.450000 | 3.170000 | 112.845000 | |
| 75% | 172.137500 | 104.287500 | 104.170000 | 5.127500 | 3.330000 | 113.287500 | |
| max | 173.010000 | 104.860000 | 104.950000 | 6.280000 | 3.680000 | 113.980000 | |
| 169 | False | 171.96 | 104.00 | 103.95 | 5.63 | 3.26 | 110.96 |

170 rows × 7 columns

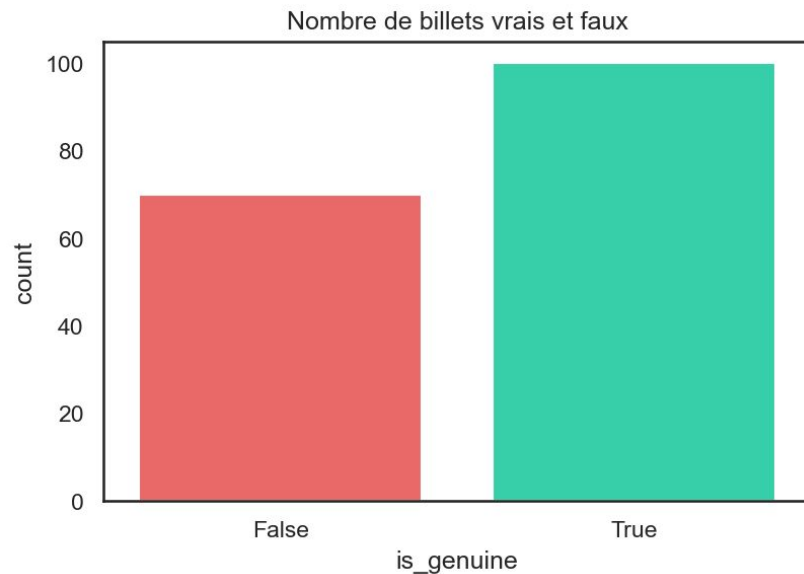
1 variable target (**vrai** ou **faux**)

6 mesures en millimètres

- Diagonale
- Hauteur à gauche
- Hauteur à droite
- Marge basse
- Marge haute
- Longueur

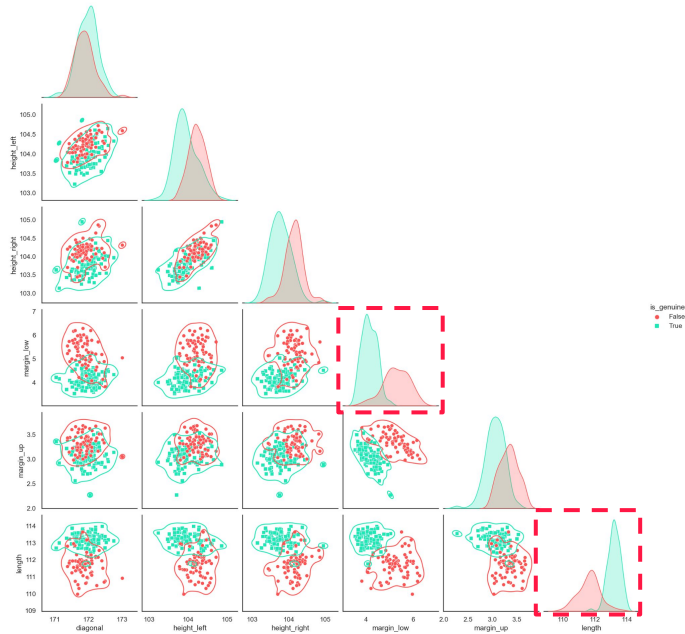
Aucun nettoyage requis

- Aucune valeur manquante
- Aucune valeur aberrante (outlier)
- Aucun doublon



100 vrais billets — 70 contrefaçons

Oversampling inutile



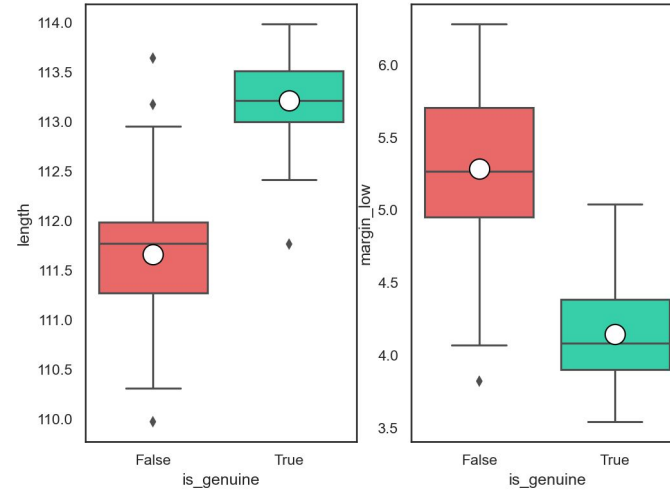
Distribution de chaque variable, par authenticité

→ Fort impact de 2 variables sur l'authenticité

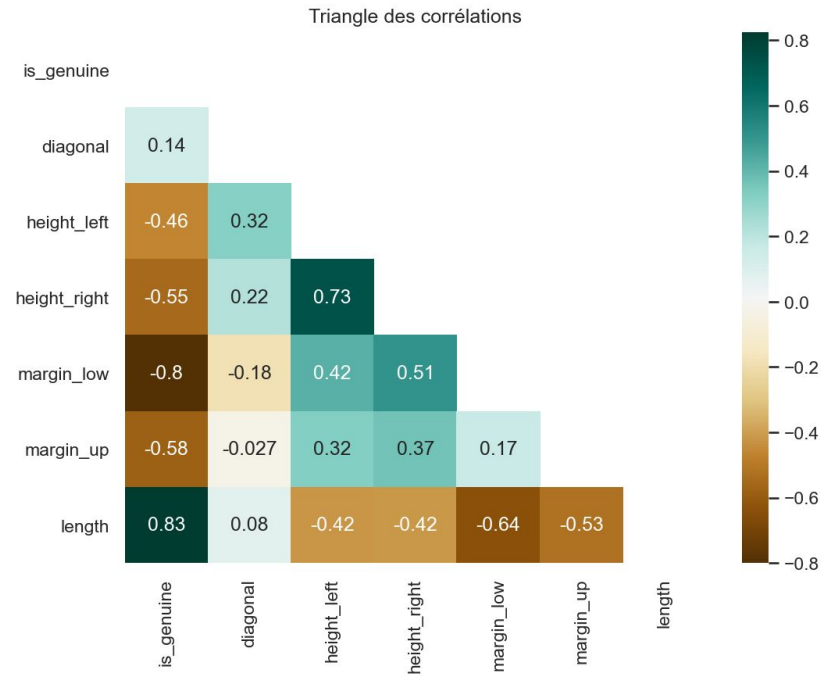
◆ longueur
◆ marge basse

→ Les autres variables ont peu d'impact

- Vrais billets plus longs que les faux
- Marge basse plus courte chez les faux billets



Distributions de *length* et *margin_low*, par authenticité



Matrice des corrélations

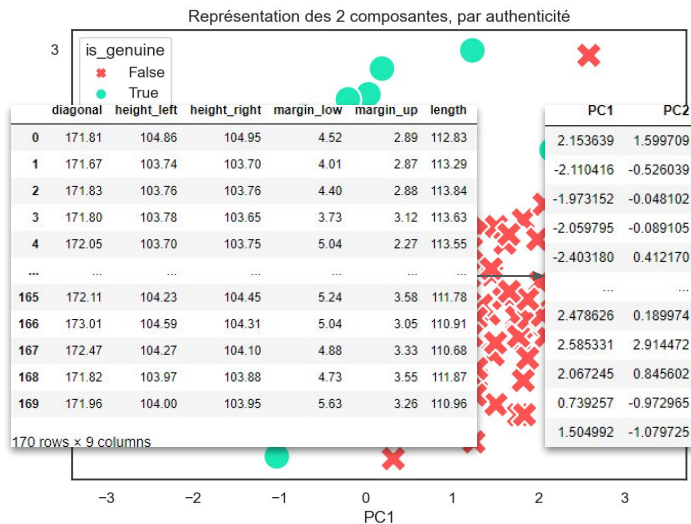
Résume les dépendances entre variables

- 0 pas de corrélation
- 1 corrélation positive
- -1 corrélation négative

- Confirmation des fortes corrélations
- Toutes les variables apportent de l'information

ACP

02



- Réduction des variables à n dimensions
- Transformation linéaire
Préserve les rapports de colinéarité
- Perte d'informations
Valeurs obtenues non interprétables
- Valeurs projetables sur un plan à 2 ou 3 dimensions
Requis : Préserver le maximum d'information dans les 2 premières composantes
Objectif : visualiser la ressemblance (variabilité) des individus et la linéarité
- Sensible à la variance

→ Différences d'échelle entre variables

Poids trop fort de certaines variables sur d'autres
Contributions inégales à l'ACP

→ Objectif : harmoniser les **écarts type**

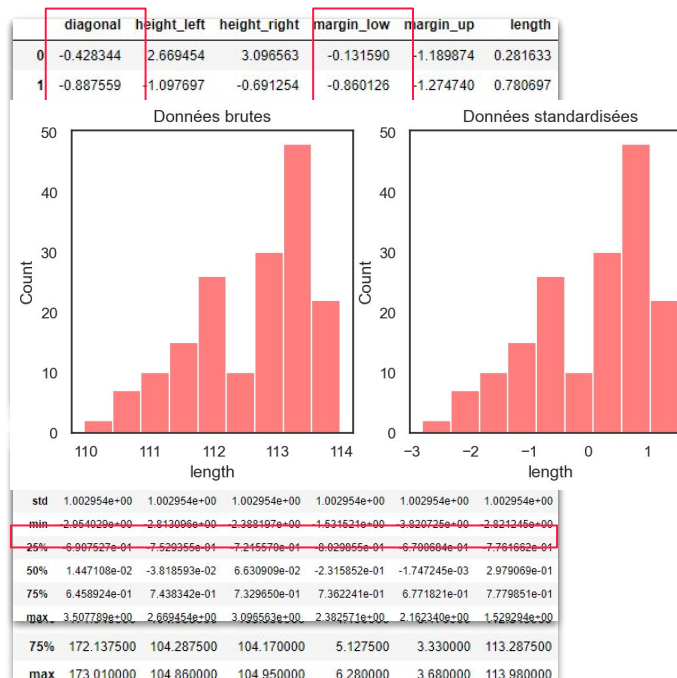
Racine carrée de la variance (moyenne des écarts à la moyenne)

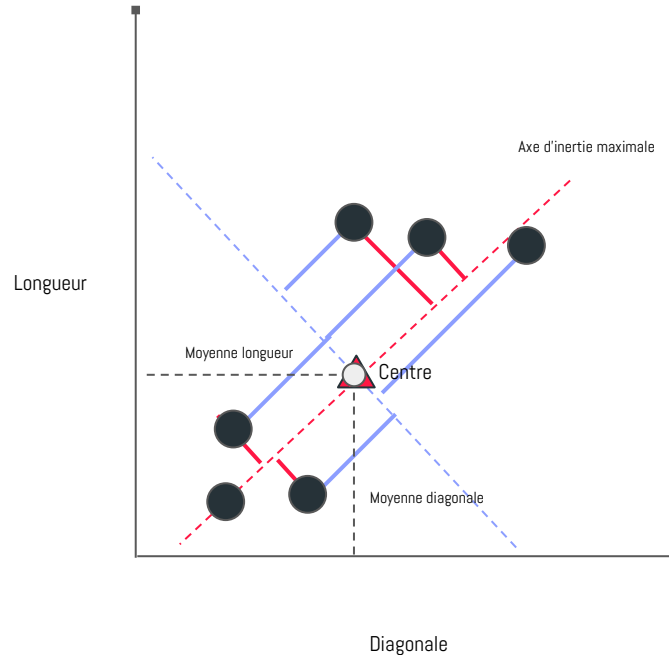
→ Standardisation

- ◆ Moyenne = 0 (centrage autour de cette valeur)
- ◆ Soustraction de la moyenne à toutes les observations
- ◆ Écart-type = 1
- ◆ Conserve la forme de la distribution

→ Adapté aux algorithmes linéaires

ACP, régression logistique ...



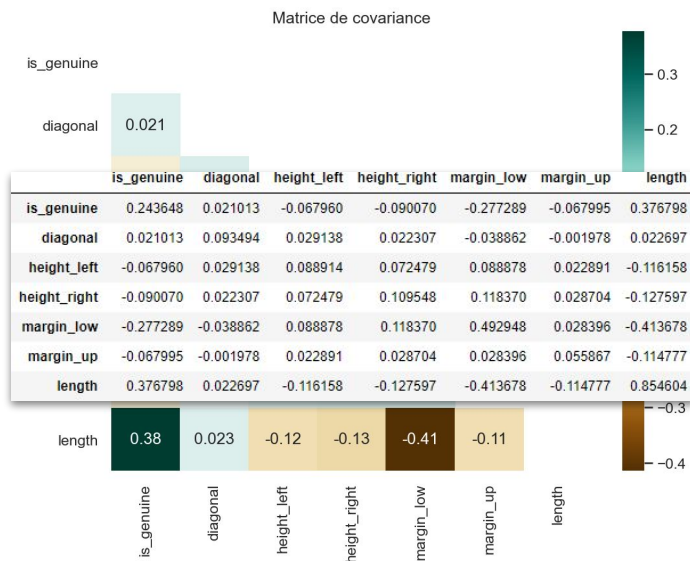


→ **Objectif** : trouver l'axe exprimant le maximum de variance

- ◆ En partant du centroïde
- ◆ Conservation du maximum d'information
(= variance : moyenne des distances au carré l'axe)

→ Axes orthogonaux aux précédents
Composantes décorrélées

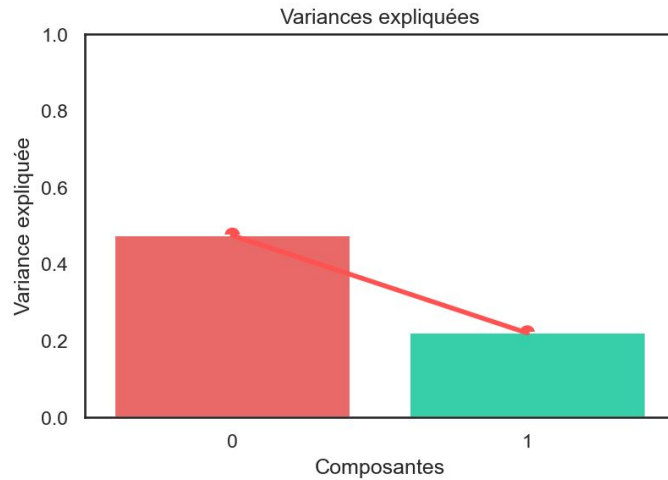
→ Autant d'axes (composantes) que de variables



→ 7 axes (vecteurs propres)
 Combinaisons linéaires de chaque variables avec les autres

→ Classement par somme des valeurs propres
 La 1ère composante est celle cumulant le plus de variance

$$\text{Axe 7} = 0.37 * \text{is_genuine} + 0.02 * \text{diagonal} - 0.11 * \text{height_left} - 0.12 * \text{height_right} - 0.41 * \text{margin_low} - 0.11 * \text{margin_up} + 0.85 * \text{length}$$



Analyse des éboulis des 2 premières composantes

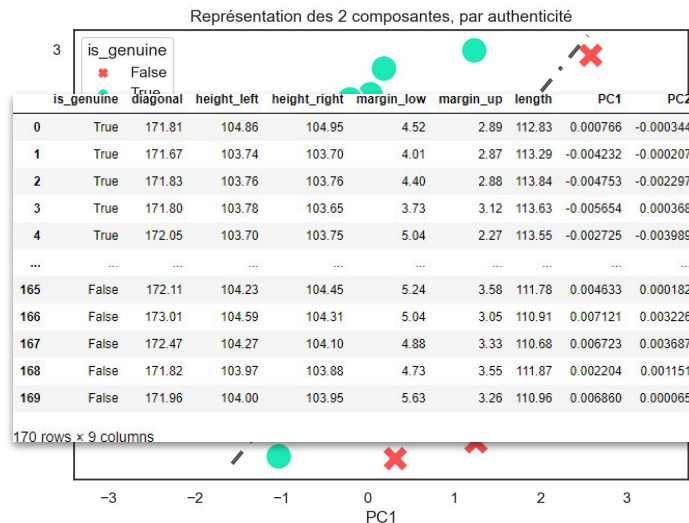
- Le maximum d'informations restantes conservé à chaque itération
- Essentiel de l'information stocké dans les 2 premières composantes

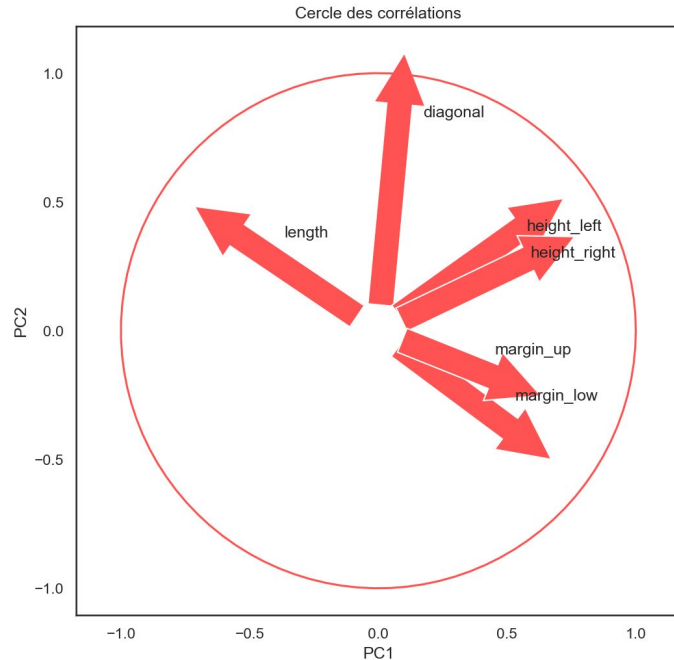
Variance expliquée : **69%**

- PC1 : 47% des données
- PC2 : 22% des données

→ Jointure des 2 composantes avec les données

→ Les données sont linéaires
Les vrais et faux billets forment 2 groupes distincts





Projection des variables sur le plan factoriel

Cosinus de l'angle entre 2 flèches = coefficient de corrélation entre les 2 variables

- 0° : corrélation positive
- 90° : absence de corrélation
- 180° : corrélation négative

- Corrélations préservées
- Bonne inertie
Longueurs des flèches homogènes et proches de 1
- Représentation fiable sur les 2 premières dimensions

Clustering

03

→ Regrouper les billets en K groupes homogènes

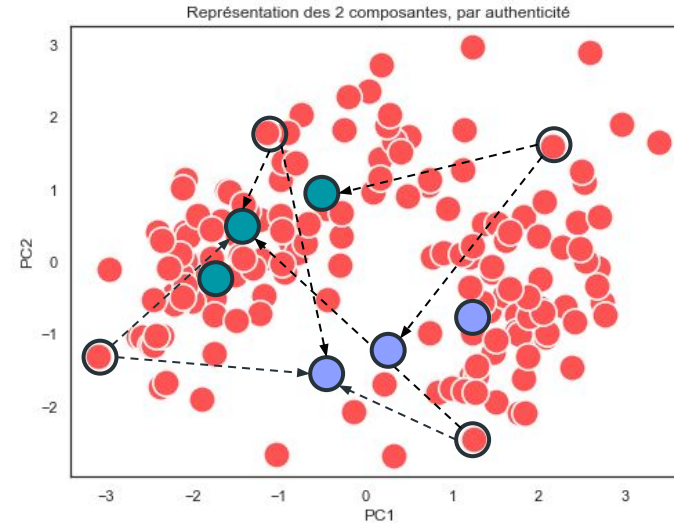
→ Classification non supervisée

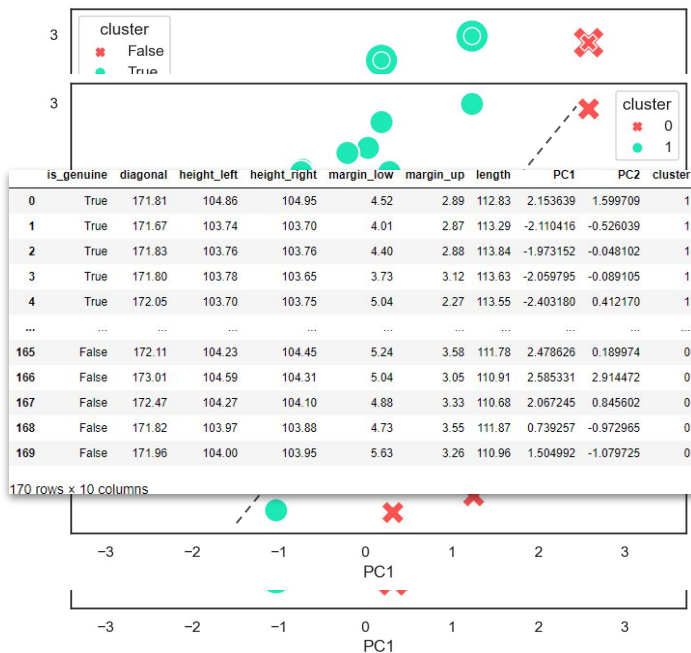
Postulat : on ne connaît pas les groupes auxquels appartiennent les billets

→ Classification non hiérarchique

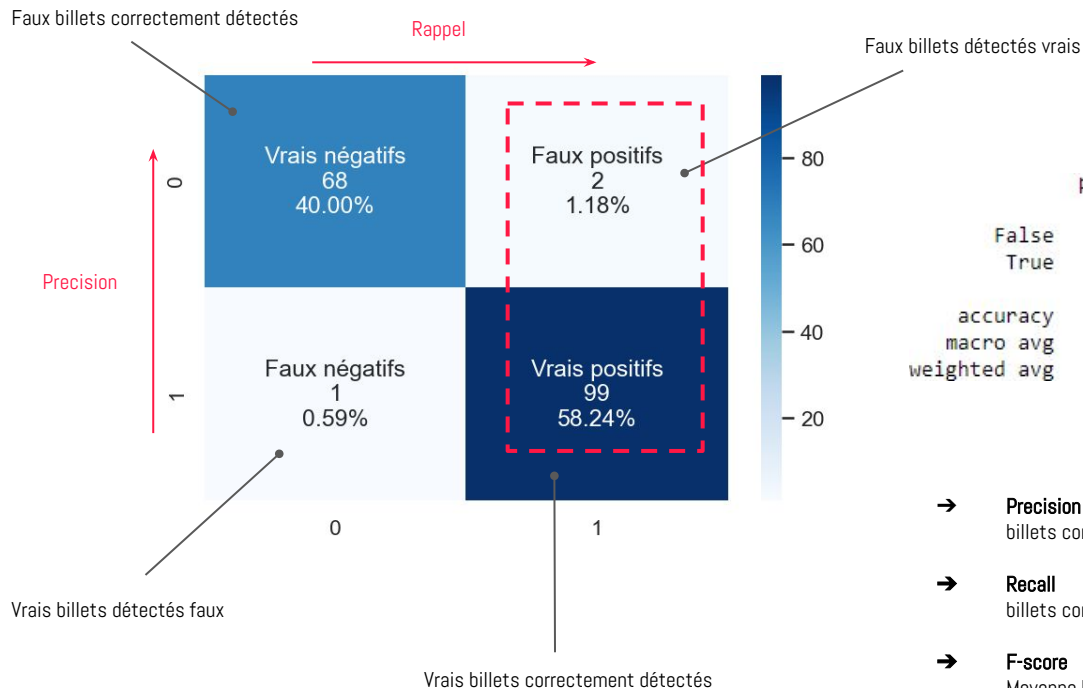
Méthode des k-moyennes

1. Sélectionner le nombre clusters à identifier (K)
2. Sélectionner K points aléatoires (clusters initiaux)
3. Mesurer la distance euclidienne entre le 1er point et les K clusters
4. Assigner le 1er point au cluster le plus proche
5. Répéter pour tous les points
6. Calculer le centroïde de chaque cluster
7. Itérer jusqu'à ce que les centres ne bougent plus





- Jointure des clusters
- 2 clusters distincts
Données linéaires
- Clusters presque identiques aux classes réelles
- 3 erreurs de classification



| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| False | 0.99 | 0.97 | 0.98 | 70 |
| True | 0.98 | 0.99 | 0.99 | 100 |
| accuracy | | | 0.98 | 170 |
| macro avg | 0.98 | 0.98 | 0.98 | 170 |
| weighted avg | 0.98 | 0.98 | 0.98 | 170 |

- **Precision** : à optimiser pour réduire les faux positifs
billets correctement prédits / billets prédits dans la classe
- **Recall**
billets correctement prédits / nbre de billets dans la classe d'origine
- **F-score**
Moyenne harmonique de la précision et du rappel

Modélisation

04

→ Calcul d'une variable quantitative d'après d'autres quantitatives

X = variable(s) indépendante(s) explicatives (features, inputs, paramètres)
 y = variable dépendante expliquée

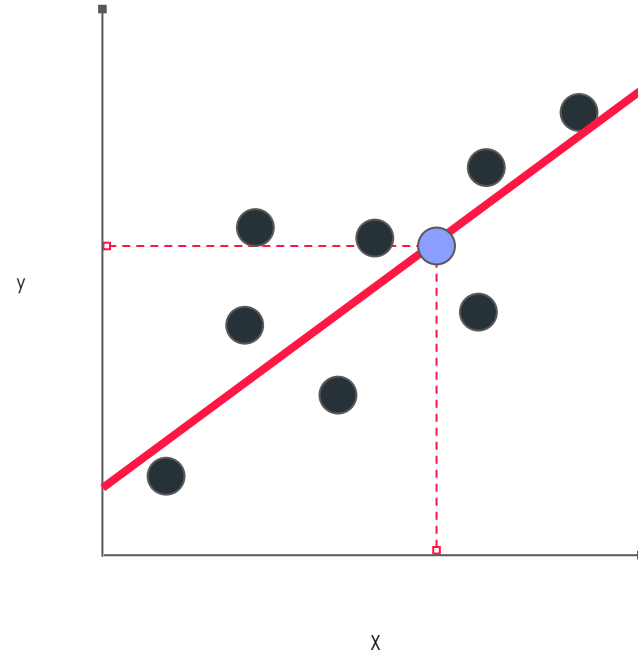
→ Apprentissage supervisé

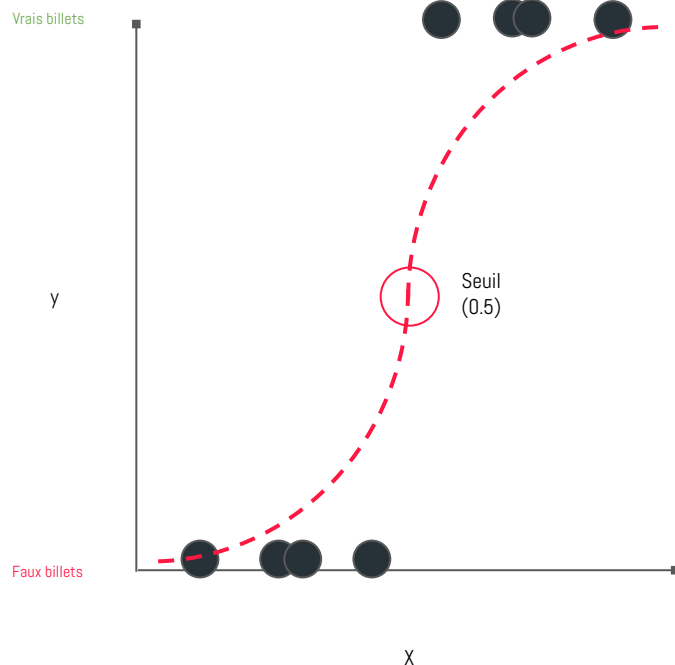
On connaît la valeur réelle de y pour chaque groupe d'explicatives

→ Postulat : on peut aligner les points sur une droite

Relation linéaire décrivant le mieux la relation entre X et y
 $= f(X) = \text{pente} \times X + \text{constante}$ (point sur l'ordonnée quand $X = 0$)

→ Prédit des données continues





- Quand la variable expliquée est qualitative
Nombre limité de valeurs possibles
- Classification supervisée
Résultat de la variable dépendante déjà connu
- Renvoie probabilités entre 0 et 1
Probabilité que l'individu appartienne à la classe True
Probabilité que l'individu appartienne à la classe False
- Transformation logistique sur la fonction de régression linéaire
 $S = \text{logit (logarithme) sur } f(x)$
Seuil de probabilité fixé généralement à 0.5

Conditions pour de bons résultats

- Peu de features (risque d'overfitting sinon)
- Données facilement séparables (par une ligne)
- Suffisamment d'individus à disposition
- Pas de valeurs aberrantes
- Données normalisées

train

| | diagonal | height_left | height_right | margin_low | margin_up | length |
|-----|----------|-------------|--------------|------------|-----------|--------|
| 27 | 172.02 | 104.23 | 104.26 | 4.92 | 2.89 | 113.49 |
| 78 | 172.16 | 104.39 | 103.85 | 3.77 | 3.32 | 112.55 |
| 147 | 172.25 | 104.52 | 104.22 | 4.65 | 3.43 | 110.48 |
| 38 | 172.21 | 104.27 | 104.01 | 4.23 | 2.79 | 113.78 |
| 41 | 171.81 | 104.10 | 103.69 | 4.29 | 2.95 | 112.72 |
| ... | ... | ... | ... | ... | ... | ... |
| 71 | 172.17 | 103.93 | 103.62 | 4.06 | 3.08 | 113.10 |
| 106 | 172.22 | 104.17 | 104.07 | 4.52 | 3.67 | 112.13 |
| 14 | 172.04 | 103.94 | 103.76 | 3.81 | 3.24 | 113.41 |
| 92 | 171.86 | 103.47 | 103.59 | 4.04 | 2.97 | 113.22 |
| 102 | 171.94 | 104.21 | 104.10 | 4.28 | 3.47 | 112.23 |

113 rows x 6 columns

```

27    True
78    True
147   False
38    True
41    True
...
71    True
106   False
14    True
92    True
102   False
Name: is_genuine, Length: 113, dtype: bool

```

test

| | diagonal | height_left | height_right | margin_low | margin_up | length |
|-----|----------|-------------|--------------|------------|-----------|--------|
| 139 | 171.60 | 104.37 | 104.20 | 5.82 | 3.08 | 112.84 |
| 30 | 172.19 | 104.05 | 103.81 | 3.90 | 3.22 | 113.52 |
| 119 | 171.51 | 104.13 | 103.90 | 4.99 | 3.60 | 111.23 |
| 29 | 171.84 | 103.75 | 103.38 | 4.08 | 2.70 | 113.72 |
| 144 | 171.56 | 103.80 | 103.87 | 5.66 | 2.98 | 112.95 |
| 163 | 171.78 | 104.07 | 104.16 | 5.77 | 3.30 | 111.27 |
| 166 | 173.01 | 104.59 | 104.31 | 5.04 | 3.05 | 110.91 |
| 51 | 172.22 | 104.48 | 104.06 | 4.59 | 2.91 | 112.82 |
| 105 | 171.99 | 104.18 | 104.20 | 5.26 | 3.23 | 111.83 |
| 60 | 172.11 | 103.67 | 103.43 | 4.19 | 2.98 | 113.09 |

```

139   False
30    True
119   False
29    True
144   False
163   False
166   False
51    True
105   False
60    True
Name: is_genuine, dtype: bool

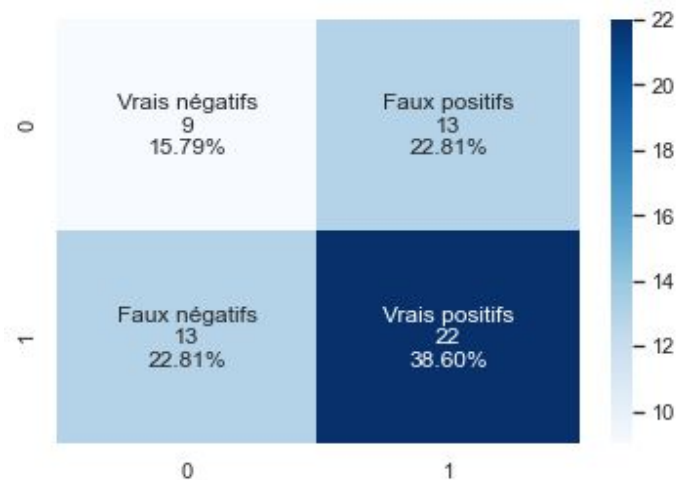
```

→ Séparation du dataset en 2 jeux

- ◆ train : 113 billets
- ◆ test : 57 billets

→ Séparation de chaque jeu entre X et y

- ◆ X = explicatives (mesures)
- ◆ y = expliquée (is_genuine)



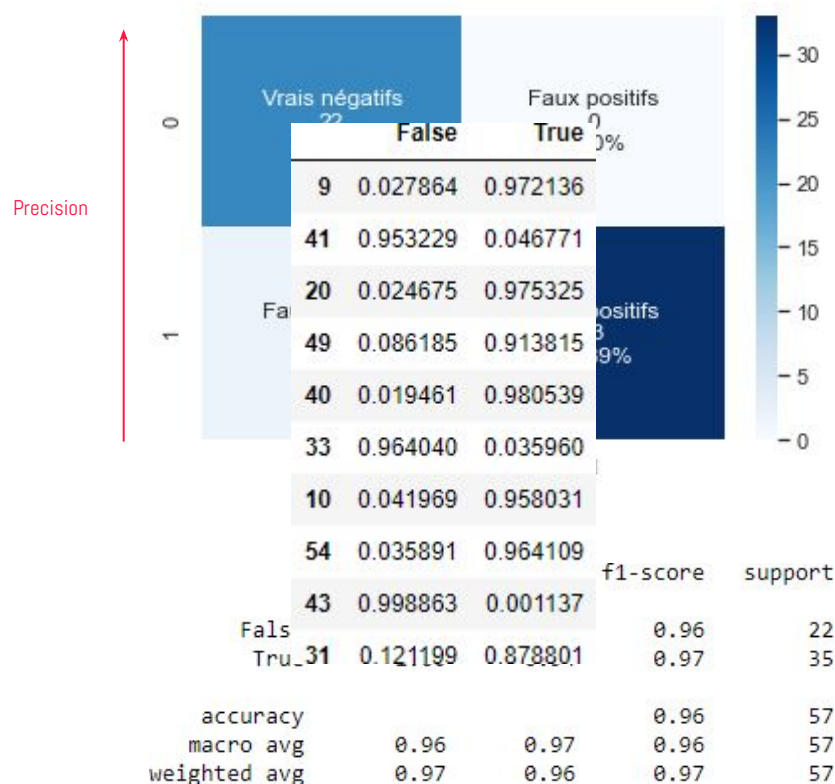
→ Base de comparaison avec le futur modèle

→ Prédictions aléatoires
Donc médiocres

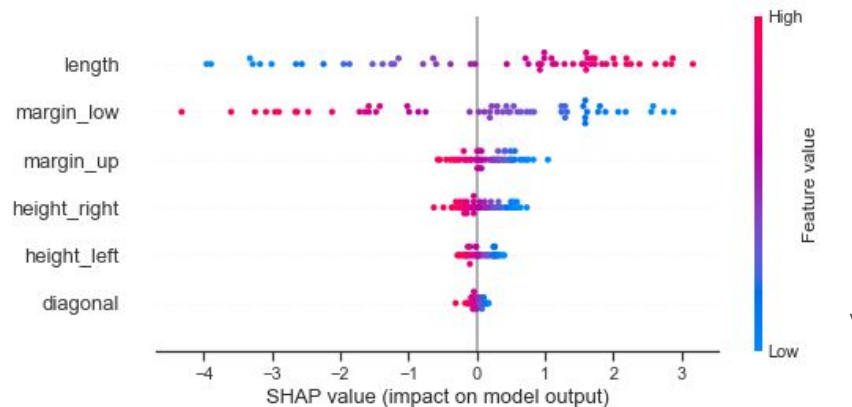
→ Précision très faible
63 % de vrais billets correctement prédits
/ vrais billets prédits

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| False | 0.41 | 0.41 | 0.41 | 22 |
| True | 0.63 | 0.63 | 0.63 | 35 |
| accuracy | | | 0.54 | 57 |
| macro avg | 0.52 | 0.52 | 0.52 | 57 |
| weighted avg | 0.54 | 0.54 | 0.54 | 57 |

- Probabilités entre 0 et 1
Pour chaque classe et chaque billet
- Comparaison avec `y_test`
- Aucun faux billet détecté vrai
Précision de 100% des vrais billets
- 2 vrais billets détectés faux



Explication de la sortie du modèle



Classement par force de contribution au modèle

Rouge : valeurs les + fortes

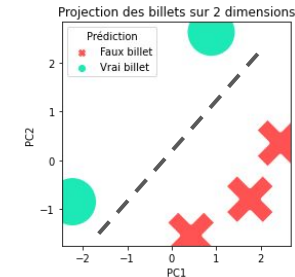
Bleu : valeurs les + faibles

Force de contribution de la variable

- Fonction autonome
- Modèles enregistrés dans un fichier Pickle
 - StandardScaler
 - ACP
 - Régression logistique
- Fichier CSV en entrée
- Résultats en tableau
- Contrôle des clusters sur 2 dimensions

Données en sortie

| | Prédiction | Probabilité de faux | Probabilité de vrai | id | | | |
|---|-------------|---------------------|---------------------|------|------|--------|-----|
| 0 | Faux billet | 0.962899 | 0.037101 | A_1 | | | |
| 1 | Faux billet | 0.994102 | 0.005898 | A_2 | | | |
| 2 | Faux billet | 0.986890 | 0.013110 | A_3 | | | |
| 3 | Vrai billet | 0.058722 | 0.941278 | A_4 | | | |
| 4 | Vrai billet | 0.004059 | 0.995941 | A_5 | | | |
| 2 | 172.00 | 104.56 | 104.29 | 4.99 | 3.39 | 111.97 | A_3 |
| 3 | 172.49 | 104.55 | 104.34 | 4.44 | 3.03 | 113.20 | A_4 |
| 4 | 171.65 | 103.63 | 103.56 | 3.77 | 3.16 | 113.33 | A_5 |



Résumé

→ Données d'entraînement adaptées à la régression logistique

- ◆ Variables à la même échelle
- ◆ Individus séparables
- ◆ Pas d'outliers

→ Modèle à 100% de précision sur les vrais billets

- ◆ Aucun faux billet classé comme vrai
- ◆ 2 vrais billets détectés faux

