

# Support Didactique - Séance 1

## Fondations Data Science

### Formation Machine Learning & MLOps

## Objectifs de la Séance

- Comprendre l'importance de l'analyse exploratoire
- Maîtriser l'interprétation des visualisations
- Gérer efficacement les valeurs manquantes
- Identifier et mitiger les biais dans les données
- Créer des features pertinentes et éthiques

## 1 Philosophie de l'Exploration de Données

### Concept Clé

**L'analyse exploratoire n'est pas une formalité mais une nécessité :**

- **Détecter les problèmes** avant qu'ils ne polluent le modèle
- **Comprendre le domaine** métier à travers les données
- **Éviter les conclusions erronées** causées par des biais cachés
- **Choisir les bonnes techniques** de prétraitement

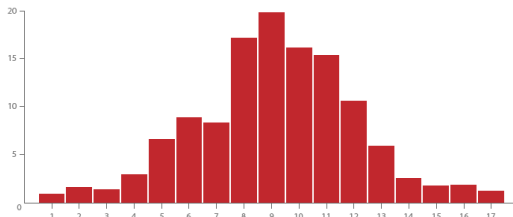
### 1.1 Les 4 Questions Fondamentales

Question	Technique	Piège à éviter
Quelle est la distribution ?	Histogramme, Boxplot	Confondre asymétrie avec outlier
Y a-t-il des relations ?	Scatterplot, Heatmap	Corrélation $\neq$ causation
Les données sont-elles propres ?	.isnull(), .duplicated()	Supprimer sans analyser le pattern
Existe-t-il des sous-groupes ?	GroupBy, FacetGrid	Conclusions sur petits échantillons

TABLE 1 – Questions fondamentales pour l'EDA

## 2 Interprétation des Visualisations

### 2.1 Histogramme : Lire la Distribution



figureExemple d'histogramme (Âge des passagers)

À analyser :

1. **Forme** : Normale, bimodale, asymétrique ?
2. **Centre** : Moyenne  $\neq$  médiane ?
3. **Étalement** : Variance, outliers ?
4. **Gaps** : Problème de collecte ?

#### Concept Clé

**Exemple concret (Âge Titanic) :**

- **Asymétrie positive** : Plus de jeunes passagers
- **Pic autour de 20-30 ans** : Voyageurs en âge de travailler
- **Queue longue droite** : Quelques passagers âgés
- **Implication** : Moyenne  $>$  Médiane, utiliser médiane pour imputation

### 2.2 Boxplot : Détecter les Outliers

La règle des  $1.5 \times IQR$

$$\text{Outlier inférieur} < Q1 - 1.5 \times IQR \quad (1)$$

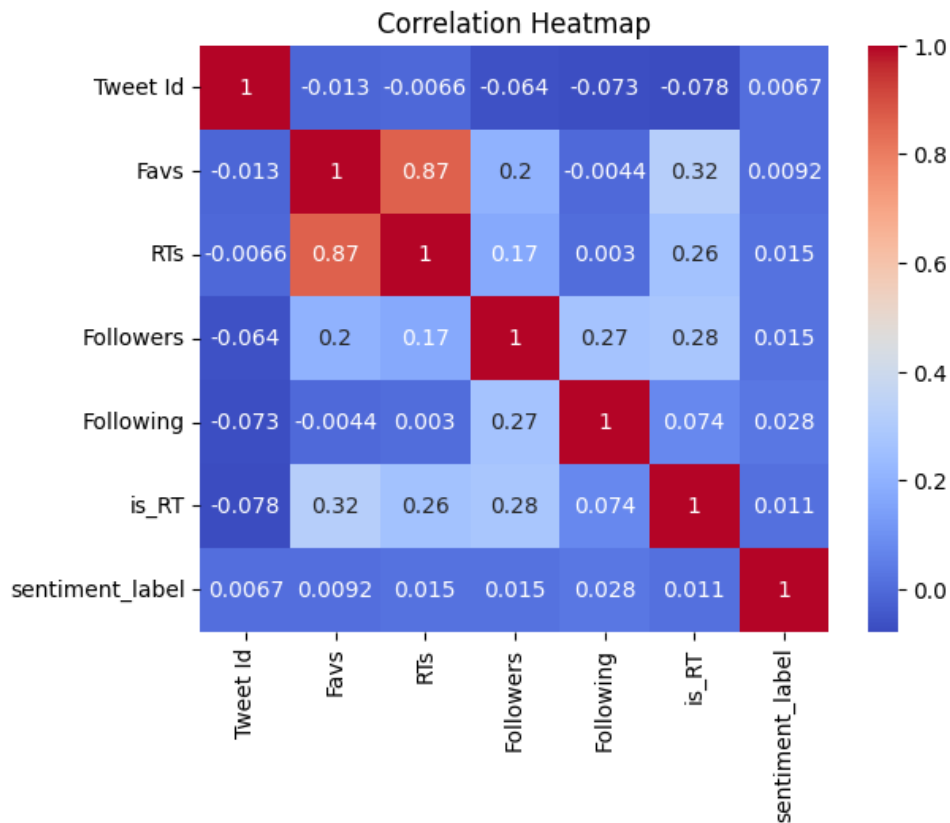
$$\text{Outlier supérieur} > Q3 + 1.5 \times IQR \quad (2)$$

$$\text{avec } IQR = Q3 - Q1 \quad (3)$$

Ne pas supprimer automatiquement !	Questions à se poser
Point de données légitime	Est-ce une erreur de mesure ?
Représente une population rare	Vient-il du même processus génératif ?
Information importante	Impacte-t-il la métrique business ?

TABLE 2 – Décision concernant les outliers

## 2.3 Heatmap de Corrélation



figureExemple de  
heatmap de corrélations

Valeur	Force	Interprétation
$ r  = 1$	Parfaite	Relation linéaire exacte
$ r  > 0.7$	Forte	Prédiction possible
$0.3 <  r  < 0.7$	Modérée	Relation détectable
$ r  < 0.3$	Faible	Peu ou pas de relation

TABLE 3 – Interprétation des coefficients de corrélation

### Attention

#### Avertissements critiques :

1. **Corrélation  $\neq$  Causation** : Le chant du coq ne fait pas lever le soleil
2. **Sensibilité aux outliers** : Un outlier peut créer une corrélation illusoire
3. **Relations non-linéaires** :  $r = 0$  n'implique pas indépendance
4. **Biais d'agrégation** : Corrélation écologique    corrélation individuelle

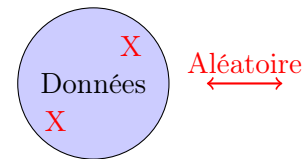
### 3 Gestion des Valeurs Manquantes

#### 3.1 Classification des Mécanismes

##### MCAR : Missing Completely At Random

**Définition :** Absence indépendante des valeurs observées ET non observées

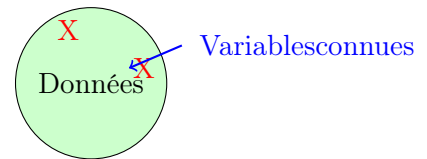
- Exemple : Questionnaire perdu aléatoirement
- **Solution :** Suppression ou imputation simple



##### MAR : Missing At Random

**Définition :** Absence dépendante des valeurs observées

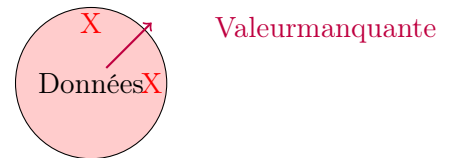
- Exemple : Âge manquant pour jeunes participants (sexe connu)
- **Solution :** Imputation conditionnelle



##### MNAR : Missing Not At Random

**Définition :** Absence dépendante des valeurs non observées

- Exemple : Salaire manquant pour hauts revenus
- **Solution :** Modélisation explicite ou flag



#### 3.2 Méthodes d'Imputation Comparées

Méthode	Avantages	Inconvénients	Utilisation
<b>Suppression</b>	Simple, pas de biais ajouté	Perte d'information	MCAR, <5%
<b>Moyenne/Médiane</b>	Rapide, stable	Réduit la variance	Variables numériques
<b>KNN</b>	Utilise similarité entre observations	Comput. coûteux, choix de k	MAR, dataset riche
<b>MICE</b>	Gère interdépendances complexes	Complexe, convergence ?	Données structurées
<b>Flag + Imputation</b>	Préserve l'information du manque	Augmente dimensionnalité	MNAR suspecté

TABLE 4 – Choix de la méthode d'imputation

#### 3.3 Détection de Pattern en Python

```

1 import pandas as pd
2 import numpy as np
3 from scipy.stats import chi2_contingency
4 import seaborn as sns
5
6 # 1. Pattern simple
7 missing_age = df['Age'].isnull()
8 print("Survie quand Age manquant :")
9 print(df.loc[missing_age, 'Survived'].mean())
10 print("Survie quand Age pr sent :")
11 print(df.loc[~missing_age, 'Survived'].mean())
12

```

```

13 # 2. Test statistique (test du chi2)
14 contingency = pd.crosstab(missing_age, df['Survived'])
15 chi2, p, dof, expected = chi2_contingency(contingency)
16 print(f"p-value : {p:.4f}") # p < 0.05      pattern significatif
17
18 # 3. Visualisation du pattern
19 df['Age_missing'] = missing_age
20 sns.countplot(data=df, x='Pclass', hue='Age_missing')
21 plt.title('R partition des valeurs manquantes par Classe')
22 plt.show()

```

Listing 1 – Analyse des patterns de valeurs manquantes

**Concept Clé****Interprétation des résultats :**

- **p-value < 0.05** : Pattern détecté → risque de biais MNAR
- **Différence de survie** : Si significative, le manque n'est pas aléatoire
- **Action** : Créer une variable "Age\_was\_missing" + imputation conservatrice

## 4 Types de Biais et Mitigation

### 4.1 Taxonomie des Biais

Type de biais	Description	Exemple Titanic	Impact
Biais de mesure	Erreur systématique collecte	"Cabin" non-enregistré 3e classe	Variables incomplètes
Biais de sélection	Échantillon non représentatif	Seulement passagers survivants	Généralisation impossible
Biais d'étiquetage	Labels incorrects/incomplets	Survie influencée par normes sociales	Modèle apprend stéréotypes
Biais de confirmation	Chercher ce qu'on veut voir	Focus sur "femmes et enfants d'abord"	Ignorer contre-exemples
Biais algorithmique	Algorithme amplifie les biais	Sur-optimisation sur classe majoritaire	Discrimination

TABLE 5 – Types de biais courants en data science

### 4.2 Stratégies de Mitigation

#### Avant la modélisation

##### 1. Analyse des marges de distribution :

$$\text{Écart} = \left| \frac{\text{Proportion échantillon} - \text{Proportion population}}{\text{Proportion population}} \right| \quad (4)$$

##### 2. Pondération des observations :

$$w_i = \frac{N_{\text{pop}}/K_{\text{pop}}}{N_{\text{sample}}/K_{\text{sample}}} \quad (5)$$

##### 3. Stratification de l'échantillon : Garantir représentation minimale

### Pendant la modélisation

- **Regularisation** : Pénaliser les coefficients extrêmes
- **Métriques équitables** : F1-score par classe, AUC par sous-groupe
- **Contraintes de fairness** :

$$\text{s.t. } |P(\hat{Y} = 1|A = a) - P(\hat{Y} = 1|A = b)| < \epsilon \quad (6)$$

### Après la modélisation

- **Audit des prédictions** : Analyse des erreurs par sous-groupe
- **Interprétabilité** : SHAP values pour comprendre les décisions
- **Monitoring continu** : Détection de dérive dans le temps

## 4.3 Détection de Biais de Sélection

```

1 def detect_selection_bias(df, reference_population):
2     """
3     Detecte les ecarts de representation entre l'échantillon
4     et la population de reference
5     """
6     bias_report = {}
7
8     for col in ['Pclass', 'Sex', 'AgeGroup']:
9         # Calcul des proportions
10        sample_prop = df[col].value_counts(normalize=True)
11        ref_prop = reference_population[col].value_counts(normalize=True)
12
13        # Calcul cart relatif
14        common_categories = set(sample_prop.index) & set(ref_prop.index)
15        max_bias = 0
16        for cat in common_categories:
17            bias = abs(sample_prop[cat] - ref_prop[cat]) / ref_prop[cat]
18            max_bias = max(max_bias, bias)
19
20        bias_report[col] = {
21            'max_relative_bias': max_bias,
22            'warning': max_bias > 0.3 # Seuil de 30%
23        }
24
25    return bias_report
26
27 # Utilisation
28 reference_data = pd.read_csv('titanic_population_reference.csv')
29 bias_results = detect_selection_bias(df, reference_data)
30 print(bias_results)

```

Listing 2 – Détection automatique de biais de représentation

## 5 Feature Engineering Critique

### 5.1 Principes de Création Éthique

#### Concept Clé

Les 5 commandements du feature engineering :

1. **Comprendre avant de créer** : Connaître la signification métier
2. **Documenter la provenance** : Traçabilité des transformations
3. **Éviter la fuite temporelle** : Features calculées sur données futures
4. **Tester l'impact** : Validation croisée pour chaque nouvelle feature
5. **Auditer les biais** : Vérifier l'impact sur les sous-groupes

### 5.2 Exemples de Features Robustes

```

1 # 1. Variables composites (réduisent la colinearité)
2 df['FamilySize'] = df['SibSp'] + df['Parch'] + 1
3 df['IsAlone'] = (df['FamilySize'] == 1).astype(int)
4
5 # 2. Binning intelligent (basé sur la distribution)
6 df['AgeGroup'] = pd.cut(df['Age'],
7                          bins=[0, 12, 18, 30, 50, 100],
8                          labels=['Child', 'Teen', 'YoungAdult',
9                                  'Adult', 'Senior'])
10
11 # 3. Ratio features (plus robustes que les différences)
12 df['FarePerPerson'] = df['Fare'] / df['FamilySize']
13
14 # 4. Interactions non-biaisées
15 df['FemaleInFirstClass'] = ((df['Sex'] == 'female') &
16                             (df['Pclass'] == 1)).astype(int)
17
18 # 5. Encodage target-aware (avec prudence !)
19 from category_encoders import TargetEncoder
20
21 # UNIQUEMENT sur le training set !
22 encoder = TargetEncoder()
23 X_train_encoded = encoder.fit_transform(X_train[['Embarked']], y_train)

```

Listing 3 – Création de features informatives et équitables

## 6 Checklist de Fin de Session

Validation Technique	Validation Éthique
<input type="checkbox"/> Valeurs manquantes analysées et traitées <input type="checkbox"/> Outliers détectés et justifiés <input type="checkbox"/> Distributions visualisées et comprises  <input type="checkbox"/> Corrélations interprétées <input type="checkbox"/> Code documenté et reproductible	<input type="checkbox"/> Biais de mesure identifiés <input type="checkbox"/> Biais de sélection quantifiés <input type="checkbox"/> Features créées sans reproduire stéréotypes <input type="checkbox"/> Impact sur sous-groupes évalué <input type="checkbox"/> Limitations clairement énoncées

TABLE 6 – Checklist de validation de l'analyse exploratoire

**La phrase à retenir :**

*"Un bon data scientist ne se contente pas de trouver des patterns, il comprend lesquels sont réels, lesquels sont des artefacts, et lesquels sont éthiquement problématiques."*

**Ressources Complémentaires**

- **Visualisation** : "The Visual Display of Quantitative Information" - Tufte
- **Biais statistiques** : "Statistical Fallacies" - Altman et Bland
- **Éthique ML** : "Weapons of Math Destruction" - O'Neil
- **Tools** : Aequitas (audit de fairness), SHAP (interprétabilité)
- **Datasets de test** : UCI ML Repository, Kaggle (avec précaution)

**Exercice Pratique****Exercice d'application immédiate :**

Créez une fonction Python qui :

1. Prend un DataFrame en entrée
2. Détecte automatiquement le type de mécanisme de valeurs manquantes (MCAR/-MAR/MNAR)
3. Propose une stratégie de traitement adaptée
4. Génère un rapport d'analyse avec visualisations

**Bonus** : Testez cette fonction sur le dataset Titanic et comparez avec votre approche initiale.

**Document préparé pour la Formation Machine Learning & MLOps - Séance 1**

*Révision : 1.0 | Date : Déc-2025 | Principe pédagogique : "Comprendre pourquoi avant d'apprendre comment"*