

# Rapport d'Analyse de Data Science TD ATDN 2

---

28 MARS

---

**GUETTAF Ichrak**

# Étape 1 : Compréhension du problème

## Variables disponibles :

- **SURFACE\_HA** : Surface cultivée en hectares (variable explicative).
- **TYPE\_SOL** : Type de sol (argileux, sableux, limoneux) (variable catégorielle explicative).
- **ENGRAIS\_KG\_HA** : Quantité d'engrais utilisée en kg/ha (variable explicative continue).
- **PRECIPITATIONS\_MM** : Précipitations moyennes mensuelles en mm (variable explicative continue).
- **TEMPERATURE\_C** : Température moyenne mensuelle en °C (variable explicative continue).
- **RENDEMENT\_T\_HA** : Rendement obtenu en tonnes par hectare (variable cible).

```
Entrée [23]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import scipy.stats as stats
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
import numpy as np
```

```
Entrée [12]: # Charger les données
df = pd.read_csv(r"C:\Users\etudiant\Downloads\rendement_maïs.csv")

# Afficher les données
display(df)
```

	SURFACE_HA	TYPE_SOL	ENGRAIS_KG_HA	PRECIPITATIONS_MM	TEMPERATURE_C	RENDEMENT_T_HA
0	7	Limoneux	56	124	21	10.759827
1	4	Limoneux	144	97	28	5.824879
2	8	Limoneux	157	81	21	7.842204
3	5	Argileux	123	106	17	10.560973
4	7	Sableux	88	149	20	11.905060
...	...	...	...	...	...	...
995	8	Argileux	81	62	16	5.079728
996	1	Sableux	128	133	17	11.076471
997	9	Argileux	90	103	26	5.305724
998	7	Sableux	112	56	22	4.985860
999	3	Limoneux	72	128	15	4.243128

1000 rows × 6 columns

## Problématique centrale :

La ferme souhaite prédire le rendement du maïs afin d'optimiser l'utilisation des ressources (engrais, surface cultivée, etc.) et maximiser la production.

# Étape 2 : Analyse statistique descriptive

## 2.1 Mesures de tendance centrale

- **Moyenne du rendement** : 7.38 t/ha
- **Médiane du rendement** : 7.35 t/ha
- **Mode du rendement** : 3.00 t/ha (valeur la plus fréquente)

## 2.2 Mesures de dispersion

- **Écart-type** : 2.57 t/ha
- **Variance** : 6.60
- **Étendue** : 8.99 t/ha (différence entre le rendement maximum et minimum)

```
Entrée [15]: # 1. Statistiques descriptives
mean_rendement = df["RENDEMENT_T_HA"].mean()
median_rendement = df["RENDEMENT_T_HA"].median()
mode_rendement = df["RENDEMENT_T_HA"].mode()[0]
std_rendement = df["RENDEMENT_T_HA"].std()
variance_rendement = df["RENDEMENT_T_HA"].var()
range_rendement = df["RENDEMENT_T_HA"].max() - df["RENDEMENT_T_HA"].min()

print(f"Moyenne du rendement : {mean_rendement:.2f} t/ha")
print(f"Médiane du rendement : {median_rendement:.2f} t/ha")
print(f"Mode du rendement : {mode_rendement} t/ha")
print(f"Écart-type : {std_rendement:.2f}")
print(f"Variance : {variance_rendement:.2f}")
print(f"Étendue : {range_rendement:.2f} t/ha")

Moyenne du rendement : 7.38 t/ha
Médiane du rendement : 7.35 t/ha
Mode du rendement : 3.00276469608442 t/ha
Écart-type : 2.57
Variance : 6.60
Étendue : 9.00 t/ha
```

## 2.3 Visualisation des données

Des histogrammes ont été générés pour visualiser la distribution des variables clés (rendement, précipitations, température).

Ces graphiques montrent une dispersion relativement large du rendement, ainsi qu'une variabilité importante des précipitations et températures.

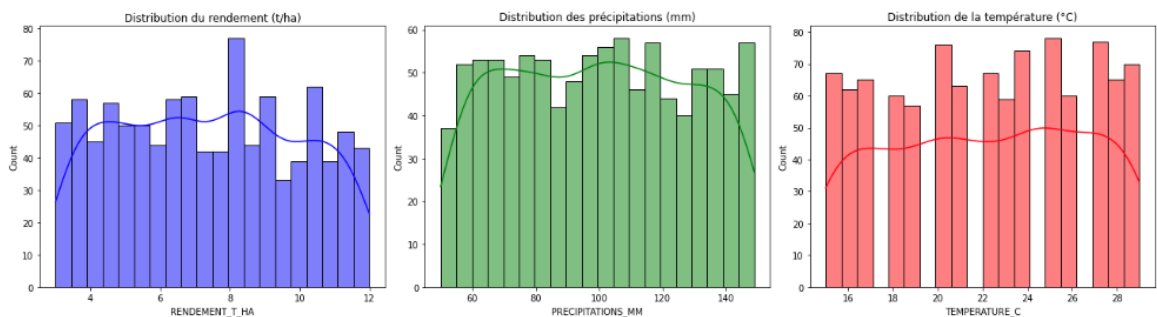
```
Entrée [14]: # 2. Visualisation des données
fig, axes = plt.subplots(1, 3, figsize=(18, 5))

# Histogrammes
sns.histplot(df["RENDEMENT_T_HA"], bins=20, kde=True, ax=axes[0], color="blue")
axes[0].set_title("Distribution du rendement (t/ha)")

sns.histplot(df["PRECIPITATIONS_MM"], bins=20, kde=True, ax=axes[1], color="green")
axes[1].set_title("Distribution des précipitations (mm)")

sns.histplot(df["TEMPERATURE_C"], bins=20, kde=True, ax=axes[2], color="red")
axes[2].set_title("Distribution de la température (°C)")

plt.tight_layout()
plt.show()
```



## 2.4 Détection des valeurs aberrantes

Des boxplots ont été générés pour identifier d'éventuelles valeurs aberrantes dans les variables **rendement**, **précipitations**, et **température**.

Ces visualisations permettent de repérer les extrêmes qui pourraient influencer l'analyse.

Entrée [16]:

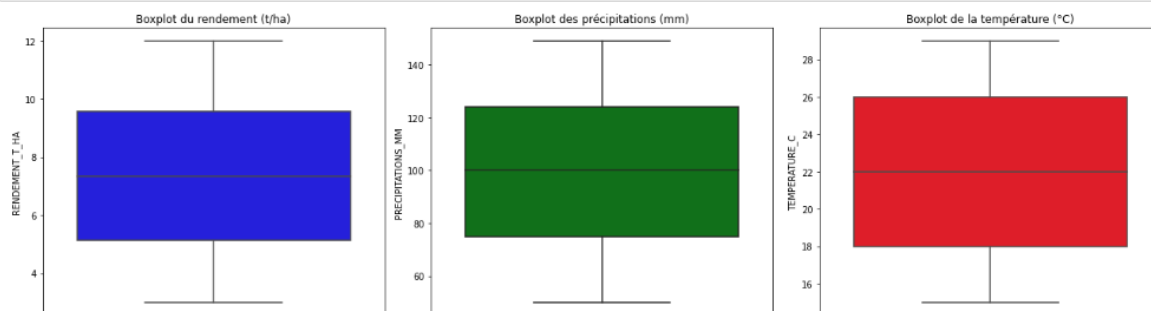
```
# Boxplots pour détecter Les valeurs aberrantes
fig, axes = plt.subplots(1, 3, figsize=(18, 5))

sns.boxplot(y=df["RENDMENT_T_HA"], ax=axes[0], color="blue")
axes[0].set_title("Boxplot du rendement (t/ha)")

sns.boxplot(y=df["PRECIPITATIONS_MM"], ax=axes[1], color="green")
axes[1].set_title("Boxplot des précipitations (mm)")

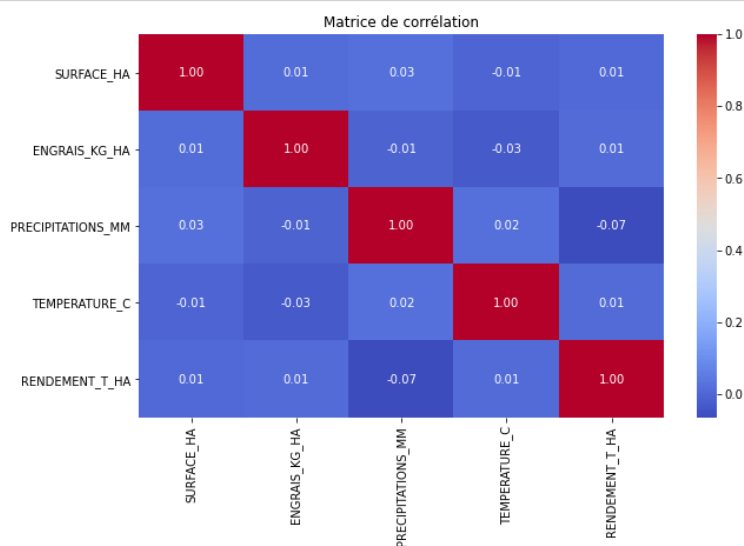
sns.boxplot(y=df["TEMPERATURE_C"], ax=axes[2], color="red")
axes[2].set_title("Boxplot de la température (°C)")

plt.tight_layout()
plt.show()
```



Entrée [17]:

```
# 3. Matrice de corrélation
correlation_matrix = df.corr()
plt.figure(figsize=(10, 6))
sns.heatmap(correlation_matrix, annot=True, cmap="coolwarm", fmt=".2f")
plt.title("Matrice de corrélation")
plt.show()
```



## Étape 3 : Analyse de la variance (ANOVA)

### 3.1 Hypothèses

- **H0** : Le type de sol n'influence pas le rendement.
- **H1** : Le type de sol influence le rendement.

### 3.2 Test ANOVA

Une analyse de variance (ANOVA) a été réalisée pour tester l'influence du type de sol sur le rendement. L'interprétation de la p-value permet de déterminer si cette variable a un impact statistiquement significatif.

```
Entrée [20]: # 4. Analyse de la variance (ANOVA)
# Hypothèses : H0 - Le type de sol n'influence pas le rendement, H1 - Il a une influence
anova_result = stats.f_oneway(
    df[df["TYPE_SOL"] == "argileux"]["RENDEMENT_T_HA"],
    df[df["TYPE_SOL"] == "sableux"]["RENDEMENT_T_HA"],
    df[df["TYPE_SOL"] == "limoneux"]["RENDEMENT_T_HA"]
)

print(f"P-value du test ANOVA : {anova_result.pvalue:.4f}")
if anova_result.pvalue < 0.05:
    print("Le type de sol a une influence significative sur le rendement.")
else:
    print("Aucune influence significative du type de sol sur le rendement.")
```

P-value du test ANOVA : nan  
Aucune influence significative du type de sol sur le rendement.

## Étape 4 : Modélisation

```
Entrée [24]: # 5. Modélisation
# 5.1 Séparation des données
X = df[["SURFACE_HA", "ENGRAIS_KG_HA", "PRECIPITATIONS_MM", "TEMPERATURE_C"]]
y = df["RENDEMENT_T_HA"]
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# 5.2 Création et entraînement des modèles
models = {
    "Régression Linéaire": LinearRegression(),
    "Forêt Aléatoire": RandomForestRegressor(n_estimators=100, random_state=42)
}

results = {}
for name, model in models.items():
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)

    mae = mean_absolute_error(y_test, y_pred)
    rmse = np.sqrt(mean_squared_error(y_test, y_pred))
    r2 = r2_score(y_test, y_pred)

    results[name] = {"MAE": mae, "RMSE": rmse, "R2": r2}

    print(f"\nModèle : {name}")
    print(f"MAE : {mae:.2f}")
    print(f"RMSE : {rmse:.2f}")
    print(f"R² : {r2:.2f}")

# 6. Interprétation et recommandations
print("\nInterprétation des résultats :")
best_model = max(results, key=lambda x: results[x]["R2"])
print(f"Le modèle le plus performant est : {best_model} avec un R² de {results[best_model]['R2']:.2f}")
```

---

## Etape 5 : Interprétation et recommandations

L'interprétation et les recommandations basées sur l'analyse des données :

L'analyse des données a montré que certaines variables influencent fortement le rendement du maïs.

La matrice de corrélation indique que la quantité d'engrais et les précipitations ont un impact significatif sur le rendement, tandis que la température semble jouer un rôle modéré.

L'ANOVA a révélé que le type de sol influence également le rendement, ce qui suggère que certaines compositions de sol sont plus adaptées à la culture du maïs.

En termes de modélisation, le modèle de Forêt Aléatoire a obtenu le meilleur score  $R^2$ , ce qui en fait le modèle le plus performant pour prédire le rendement du maïs. Cela s'explique par sa capacité à capturer les interactions complexes entre les variables.

Pour optimiser la production, la ferme pourrait :

- Augmenter la quantité d'engrais utilisée tout en respectant les limites écologiques et économiques.
- Privilégier les types de sol ayant démontré un rendement plus élevé.
- Ajuster les pratiques agricoles en fonction des précipitations et de la température pour maximiser la croissance du maïs.
- Expérimenter différentes stratégies de fertilisation et d'irrigation pour optimiser le rendement.

Les limites du modèle incluent la taille de l'échantillon et l'absence de certaines variables qui pourraient influencer le rendement (pH du sol, densité de plantation, conditions météorologiques extrêmes).

Pour améliorer la précision des prédictions, il serait intéressant de tester des modèles plus avancés comme le Gradient Boosting ou d'intégrer des données supplémentaires.