

68 Best Interview Questions - Data Analyst

UniBank Haiti - Commercial Banking Context

Édition Complète avec Régression, Séries Temporelles et Tests Non-Paramétriques

Niveau: Intermédiaire à Senior

Format: Question + Réponse complète avec introduction

GLOSSAIRE DES DÉFINITIONS CLÉS

SQL et Bases de Données

Terme	Définition
SQL	Structured Query Language - langage standardisé pour interroger et manipuler des bases de données relationnelles
JOIN	Opération qui combine des lignes de deux ou plusieurs tables basées sur une colonne liée
CTE (Common Table Expression)	Requête nommée temporaire définie dans une clause WITH, utilisable dans la requête principale
Window Function	Fonction qui effectue un calcul sur un ensemble de lignes liées à la ligne courante
Index	Structure de données qui améliore la vitesse des opérations de récupération des données
N+1 Problem	Anti-pattern où une requête initiale génère N requêtes supplémentaires, une par résultat
Aggrégation	Opération qui combine plusieurs valeurs en une seule (SUM, AVG, COUNT, etc.)

Machine Learning et Modélisation

Terme	Définition
Apprentissage supervisé	Modèle entraîné avec des données étiquetées (features + target)
Apprentissage non supervisé	Modèle qui découvre des patterns sans étiquettes (clustering)
Classification	Prédiction d'une catégorie (défaut oui/non, segment client)
Régression	Prédiction d'une valeur continue (montant, LTV)
AUC-ROC	Aire sous la courbe ROC - mesure de discrimination du modèle
Gini	Coefficient = $2 \times \text{AUC} - 1$, mesure le pouvoir discriminant (scoring)
Overfitting	Modèle trop adapté aux données d'entraînement, mauvaise généralisation
Cross-validation	Technique de validation avec découpage en k folds

Types de Variables

Terme	Définition
Variable nominale	Catégorielle sans ordre naturel (type de compte, région)
Variable ordinale	Catégorielle avec ordre (rating AAA > AA > A)
Variable discrète	Valeurs entières dénombrables (nb de transactions)
Variable continue	Valeurs décimales dans un intervalle (montant, taux)
MCAR	Missing Completely At Random - manquants aléatoires
MAR	Missing At Random - manquants liés à d'autres variables
MNAR	Missing Not At Random - manquants liés à la valeur elle-même

Statistiques

Terme	Définition
Moyenne (Mean)	Somme des valeurs divisée par le nombre d'observations
Médiane	Valeur centrale qui sépare un ensemble de données en deux moitiés égales
Mode	Valeur la plus fréquente dans un ensemble de données
Écart-type	Mesure de la dispersion des données autour de la moyenne
Variance	Carré de l'écart-type, mesure la dispersion des données
Skewness (Asymétrie)	Mesure de l'asymétrie d'une distribution par rapport à sa moyenne
Kurtosis	Mesure de l'épaisseur des queues d'une distribution
p-value	Probabilité d'obtenir un résultat au moins aussi extrême que celui observé, si H_0 est vraie
Corrélation	Mesure statistique de la force et de la direction de la relation entre deux variables
Intervalle de confiance	Plage de valeurs qui contient probablement le paramètre de population réel

KPIs Bancaires

Terme	Définition
NPL (Non-Performing Loans)	Prêts en défaut ou proches du défaut (généralement > 90 jours de retard)
ROE (Return on Equity)	Rendement des capitaux propres - mesure la rentabilité pour les actionnaires
ROA (Return on Assets)	Rendement des actifs - mesure l'efficacité d'utilisation des actifs
NIM (Net Interest Margin)	Marge nette d'intérêt - différence entre revenus et charges d'intérêts
CAR (Capital Adequacy Ratio)	Ratio de solvabilité - fonds propres / actifs pondérés par les risques
LDR (Loan-to-Deposit Ratio)	Ratio prêts/dépôts - mesure la liquidité et l'utilisation des dépôts
CIR (Cost-to-Income Ratio)	Ratio d'efficacité opérationnelle - charges / revenus

Terme	Définition
AML (Anti-Money Laundering)	Ensemble des procédures pour prévenir le blanchiment d'argent

Data Analysis

Terme	Définition
EDA (Exploratory Data Analysis)	Analyse exploratoire des données - processus d'investigation initiale des données
Outlier	Valeur anormalement éloignée des autres observations
Data Wrangling	Processus de nettoyage, transformation et préparation des données
ETL	Extract, Transform, Load - processus d'intégration de données
DataFrame	Structure de données tabulaire en Pandas avec lignes et colonnes indexées

SECTION 1: SQL (15 Questions)

Q1. Quelle est la différence entre WHERE et HAVING?

Réponse:

Cette question teste votre compréhension de l'ordre d'exécution des requêtes SQL. WHERE et HAVING sont tous deux des clauses de filtrage, mais ils s'appliquent à des moments différents dans l'exécution de la requête.

- **WHERE** filtre les lignes **AVANT** l'agrégation (GROUP BY)
- **HAVING** filtre les groupes **APRÈS** l'agrégation

```
-- WHERE: filtre les transactions > 1000 avant de grouper
SELECT client_id, SUM(montant)
FROM transactions
WHERE montant > 1000
GROUP BY client_id;

-- HAVING: filtre les clients dont le total > 50000
SELECT client_id, SUM(montant) as total
FROM transactions
GROUP BY client_id
HAVING SUM(montant) > 50000;
```

Q2. Expliquez la différence entre ROW_NUMBER, RANK et DENSE_RANK.

Réponse:

Ces trois fonctions sont des Window Functions qui attribuent un rang aux lignes, mais elles diffèrent dans leur traitement des valeurs égales. Comprendre cette différence est essentiel pour les classements et la pagination.

Fonction	Comportement avec égalité	Exemple
ROW_NUMBER	Toujours unique, arbitraire pour égalités	1, 2, 3, 4, 5
RANK	Saute les rangs après égalité	1, 2, 2, 4, 5
DENSE_RANK	Ne saute pas les rangs	1, 2, 2, 3, 4

SELECT

```

client_id,
montant,
ROW_NUMBER() OVER (ORDER BY montant DESC) as row_num,
RANK() OVER (ORDER BY montant DESC) as rank,
DENSE_RANK() OVER (ORDER BY montant DESC) as dense_rank
FROM transactions;

```

Q3. Qu'est-ce qu'une CTE et quand l'utiliser?

Réponse:

Une CTE (Common Table Expression) est une requête nommée temporaire qui améliore la lisibilité et la maintenabilité du code SQL. Elle est définie dans une clause WITH et peut être référencée comme une table dans la requête principale.

```

WITH clients_actifs AS (
    SELECT client_id, COUNT(*) as nb_tx
    FROM transactions
    WHERE date >= CURRENT_DATE - INTERVAL '30 days'
    GROUP BY client_id
),
clients_vip AS (
    SELECT client_id FROM clients WHERE segment = 'VIP'
)
SELECT c.client_id, ca.nb_tx
FROM clients_vip c
JOIN clients_actifs ca ON c.client_id = ca.client_id;

```

Utiliser quand: requêtes complexes avec sous-requêtes multiples, besoin de réutiliser un résultat, ou pour améliorer la lisibilité.

Q4. Comment calculer un running total (cumul) en SQL?

Réponse:

Le calcul d'un total cumulatif est une opération courante en analyse financière, notamment pour suivre l'évolution des soldes ou des ventes. En SQL moderne, on utilise les Window Functions avec une clause de frame.

```

SELECT
    date,
    montant,
    SUM(montant) OVER (ORDER BY date ROWS UNBOUNDED PRECEDING) as cumul
FROM transactions;

```

Explication: ROWS UNBOUNDED PRECEDING indique que la somme inclut toutes les lignes depuis le début du résultat jusqu'à la ligne courante.

Q5. Quelle est la différence entre INNER JOIN, LEFT JOIN et FULL OUTER JOIN?

Réponse:

Les différents types de JOIN déterminent comment les lignes des deux tables sont combinées et quelles lignes sont incluses dans le résultat final. Le choix du type de JOIN dépend des besoins métier.

JOIN Type	Résultat
INNER	Seulement les lignes qui ont une correspondance dans les deux tables
LEFT	Toutes les lignes de la table gauche + correspondances de droite (NULL si absent)
RIGHT	Toutes les lignes de la table droite + correspondances de gauche
FULL OUTER	Toutes les lignes des deux tables, avec NULL où il n'y a pas de correspondance

-- Tous les clients avec leurs transactions (même sans transactions)

```
SELECT c.nom, t.montant  
FROM clients c  
LEFT JOIN transactions t ON c.id = t.client_id;
```

Q6. Comment trouver les doublons dans une table?

Réponse:

Identifier les doublons est une étape cruciale du nettoyage des données. Il existe plusieurs méthodes, chacune adaptée à des situations spécifiques.

-- Méthode 1: GROUP BY + HAVING - identifie les valeurs dupliquées

```
SELECT email, COUNT(*) as occurrences  
FROM clients  
GROUP BY email  
HAVING COUNT(*) > 1;
```

-- Méthode 2: Window function - retourne les lignes dupliquées elles-mêmes

```
SELECT * FROM (   
    SELECT *, ROW_NUMBER() OVER (PARTITION BY email ORDER BY id) as rn  
    FROM clients  
) t WHERE rn > 1;
```

Q7. Comment calculer la variation mois sur mois (MoM)?

Réponse:

Le calcul de la variation période sur période est fondamental en analyse financière pour mesurer la croissance ou la décroissance. La fonction LAG permet d'accéder aux valeurs des lignes précédentes.

```

WITH monthly AS (
    SELECT
        DATE_TRUNC('month', date) AS mois,
        SUM(montant) AS total
    FROM transactions
    GROUP BY DATE_TRUNC('month', date)
)
SELECT
    mois,
    total,
    LAG(total) OVER (ORDER BY mois) AS mois_precedent,
    ROUND((total - LAG(total) OVER (ORDER BY mois)) /
          NULLIF(LAG(total) OVER (ORDER BY mois), 0) * 100, 2) AS variation_pct
FROM monthly;

```

Q8. Qu'est-ce que le problème N+1 et comment le résoudre?

Réponse:

Le problème N+1 est un anti-pattern de performance critique qui se produit principalement avec les ORM. Il survient quand une requête initiale est suivie de N requêtes supplémentaires pour charger les données liées.

Définition: 1 requête pour la liste principale + N requêtes pour les détails = N+1 total.

```

# MAUVAIS - N+1 (génère 1001 requêtes pour 1000 clients)
clients = session.query(Client).all()
for client in clients:
    print(client.transactions) # 1 requête par client!

# BON - Eager loading (génère seulement 2 requêtes)
clients = session.query(Client).options(selectinload(Client.transactions)).all()

```

Solutions: JOIN (joinedload), batch loading avec IN (selectinload), ou requête SQL directe avec agrégation.

Q9. Comment optimiser une requête lente?

Réponse:

L'optimisation des requêtes est une compétence essentielle pour tout analyste de données travaillant avec des volumes importants. Voici une approche méthodique en 6 étapes.

1. **EXPLAIN ANALYZE** - Comprendre le plan d'exécution et identifier les goulots d'étranglement
 2. **Ajouter des index** - Sur les colonnes utilisées dans WHERE, JOIN et ORDER BY
 3. **Sélectionner les colonnes nécessaires** - Éviter SELECT * qui charge toutes les colonnes
 4. **Éviter les fonctions dans WHERE** - Elles empêchent l'utilisation des index
 5. **Partitionner les grandes tables** - Généralement par date pour les données temporelles
 6. **Utiliser EXISTS plutôt que IN** - Plus efficace pour les sous-requêtes volumineuses
-

Q10. Expliquez COALESCE et son utilisation.

Réponse:

COALESCE est une fonction SQL standard qui retourne la première valeur non-NULL parmi ses arguments. Elle est essentielle pour gérer les valeurs manquantes et fournir des valeurs par défaut.

```
-- Retourne le premier numéro de téléphone disponible
SELECT
    client_id,
    COALESCE(telephone_mobile, telephone_fixe, 'N/A') AS telephone
FROM clients;

-- Évite les NULL dans les agrégations
SELECT COALESCE(SUM(montant), 0) AS total FROM transactions;
```

Q11. Comment faire un pivot en SQL?

Réponse:

Le pivotement transforme les valeurs d'une colonne en colonnes distinctes, permettant une vue matricielle des données. Cette technique est courante pour les rapports financiers et les tableaux croisés.

```
-- Transactions par client et type (pivot manuel avec CASE)
SELECT
    client_id,
    SUM(CASE WHEN type = 'DEPOT' THEN montant ELSE 0 END) AS depots,
    SUM(CASE WHEN type = 'RETRAIT' THEN montant ELSE 0 END) AS retraits,
    SUM(CASE WHEN type = 'TRANSFERT' THEN montant ELSE 0 END) AS transferts
FROM transactions
GROUP BY client_id;
```

Q12. Quelle est la différence entre DELETE, TRUNCATE et DROP?

Réponse:

Ces trois commandes suppriment des données, mais avec des implications très différentes en termes de performance, récupération et structure. Le choix dépend du cas d'usage.

Commande	Action	Rollback possible	Clause WHERE	Performance
DELETE	Supprime des lignes	Oui	Oui	Lent (log chaque ligne)
TRUNCATE	Vide la table entière	Non*	Non	Rapide (pas de log par ligne)
DROP	Supprime la table et sa structure	Non	N/A	Instantané

*Note: TRUNCATE est plus rapide car il ne journalise pas chaque suppression individuellement.

Q13. Comment calculer le percentile en SQL?

Réponse:

Les percentiles sont essentiels pour comprendre la distribution des données et identifier les seuils. SQL offre des fonctions dédiées pour ce calcul.

-- Médiane (50e percentile) et P95 des montants

```
SELECT  
    PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY montant) as mediane,  
    PERCENTILE_CONT(0.95) WITHIN GROUP (ORDER BY montant) as p95  
FROM transactions;
```

-- Créer des quartiles avec NTILE

```
SELECT  
    *,  
    NTILE(4) OVER (ORDER BY montant) as quartile  
FROM transactions;
```

Q14. Comment gérer les NULL dans les comparaisons?

Réponse:

NULL représente une valeur inconnue en SQL et ne peut pas être comparé avec les opérateurs standards. Cette particularité est source de nombreuses erreurs pour les débutants.

-- NE FONCTIONNE PAS - retourne toujours faux

```
SELECT * FROM clients WHERE telephone = NULL;
```

-- CORRECT - utiliser IS NULL / IS NOT NULL

```
SELECT * FROM clients WHERE telephone IS NULL;  
SELECT * FROM clients WHERE telephone IS NOT NULL;
```

-- Comparaison sécurisée avec COALESCE

```
SELECT * FROM clients  
WHERE COALESCE(telephone, '') = COALESCE(@param, '');
```

Q15. Écrivez une requête pour trouver les clients inactifs depuis 90 jours.

Réponse:

Identifier les clients inactifs est crucial pour les stratégies de rétention. Cette requête illustre l'utilisation de LEFT JOIN pour trouver les absences de correspondance.

-- Méthode 1: LEFT JOIN avec condition de date

```
SELECT c.*  
FROM clients c  
LEFT JOIN transactions t ON c.id = t.client_id  
    AND t.date >= CURRENT_DATE - INTERVAL '90 days'  
WHERE t.id IS NULL;
```

```
-- Méthode 2: NOT EXISTS (souvent plus performant)
SELECT c.*
FROM clients c
WHERE NOT EXISTS (
    SELECT 1 FROM transactions t
    WHERE t.client_id = c.id
    AND t.date >= CURRENT_DATE - INTERVAL '90 days'
);
```

SECTION 2: Types de Variables et Machine Learning (12 Questions)

Q16. Quelle est la différence entre variable nominale et ordinale?

Réponse:

Cette distinction est fondamentale pour choisir les bonnes techniques d'analyse et d'encodage. Les deux sont catégorielles mais diffèrent par la notion d'ordre.

Type	Définition	Exemple bancaire	Statistiques possibles
Nominale	Catégories SANS ordre naturel	Type de compte (Épargne, Courant, DAT)	Mode, Chi-carré
Ordinal	Catégories AVEC ordre naturel	Rating crédit (AAA > AA > A > BBB)	Médiane, Mann-Whitney

Encodage approprié:

```
# Nominale → One-Hot Encoding
pd.get_dummies(df, columns=['type_compte'])

# Ordinale → Label Encoding avec ordre
ordre = {'AAA': 4, 'AA': 3, 'A': 2, 'BBB': 1}
df['rating_encoded'] = df['rating'].map(ordre)
```

Q17. Quelle est la différence entre variable discrète et continue?

Réponse:

Cette distinction détermine les visualisations et tests statistiques appropriés.

Type	Définition	Exemple	Test/Visualisation
Discrète	Valeurs entières dénombrables	Nb transactions, Nb produits	Bar chart, Poisson
Continue	Valeurs décimales dans un intervalle	Montant, Taux d'intérêt	Histogramme, t-test

Règle simple: Si 2.5 n'a pas de sens (2.5 transactions?), c'est discret.

```

# Identifier le type
if df['col'].dtype == 'int64' and df['col'].nunique() < 20:
    print("Probablement discrète")
elif df['col'].dtype == 'float64':
    print("Continue")

```

Q18. Quels sont les niveaux de mesure des variables (NOIR)?

Réponse:

Les quatre niveaux de mesure déterminent quelles opérations mathématiques sont valides.
Mnémotechnique: **NOIR**.

Niveau	Opérations	Exemple	Zéro
Nominale	= ≠	Type de compte	N/A
Ordinal	= ≠ < >	Rating crédit	N/A
Intervalle	+ -	Température °C	Arbitraire
Ratio	× ÷	Montant HTG	Absolu (0 = rien)

Important: Le revenu est Ratio (0 = pas de revenu), la température est Intervalle (0°C n'est pas "pas de température").

Q19. Comment traiter les valeurs manquantes (MCAR, MAR, MNAR)?

Réponse:

Le type de manquant détermine la stratégie d'imputation appropriée.

Type	Signification	Diagnostic	Traitement
MCAR	Manquant complètement aléatoire	Indépendant de toutes les variables	Imputation simple (médiane)
MAR	Manquant aléatoire	Dépend d'autres variables observées	Imputation par groupe
MNAR	Manquant non aléatoire	Dépend de la valeur elle-même	Modèle spécialisé

```

# Règle des seuils
# < 5% → Suppression acceptable
# 5-20% → Imputation nécessaire
# > 20% → Envisager suppression de la variable

# Imputation médiane
df['col'].fillna(df['col'].median(), inplace=True)

# Imputation par groupe (MAR)
df['revenu'] = df.groupby('type_emploi')['revenu'].transform(
    lambda x: x.fillna(x.median()))
)

```

```
# IMPORTANT: Fit sur train seulement!
scaler.fit_transform(X_train)
scaler.transform(X_test) # Pas de fit!
```

Q20. Quelle est la différence entre apprentissage supervisé et non supervisé?

Réponse:

Cette distinction fondamentale en ML détermine le type de problème et les algorithmes à utiliser.

Aspect	Supervisé	Non Supervisé
Données	Avec labels (X, Y)	Sans labels (X seulement)
Objectif	Prédire Y	Découvrir des patterns
Exemples	Classification, Régression	Clustering, Réduction dim.
Bancaire	Scoring crédit, Fraude	Segmentation clients

```
# Supervisé - Classification
from sklearn.linear_model import LogisticRegression
model = LogisticRegression()
model.fit(X_train, y_train) # Avec target!
y_pred = model.predict(X_test)

# Non supervisé - Clustering
from sklearn.cluster import KMeans
kmeans = KMeans(n_clusters=5)
labels = kmeans.fit_predict(X) # Pas de target!
```

Q21. Quelle est la différence entre Classification et Régression?

Réponse:

Les deux sont supervisés mais diffèrent par le type de variable cible.

Aspect	Classification	Régression
Target	Catégorielle (classes)	Continue (valeurs)
Output	Probabilités/Classes	Valeur numérique
Métriques	Accuracy, AUC, F1	MAE, RMSE, R ²
Bancaire	Défaut oui/non	Montant de perte

Cas bancaires: - **Classification:** Scoring crédit, Détection fraude, Prédiction churn - **Régression:** Estimation LTV, Prédiction montant de prêt, LGD

Q22. Qu'est-ce que le Gini et l'AUC en scoring?

Réponse:

Ces métriques mesurent le pouvoir discriminant d'un modèle de scoring - sa capacité à séparer les bons des mauvais clients.

Formule: $\text{Gini} = 2 \times \text{AUC} - 1$

Métrique	Plage	Interprétation
AUC	0.5 - 1.0	0.5 = aléatoire, 1.0 = parfait
Gini	0 - 1.0	> 0.4 acceptable, > 0.5 bon

```
from sklearn.metrics import roc_auc_score

# Calcul
y_proba = model.predict_proba(X_test)[:, 1]
auc = roc_auc_score(y_test, y_proba)
gini = 2 * auc - 1

print(f"AUC: {auc:.4f}")
print(f"Gini: {gini:.4f}")
```

Q23. Expliquez Precision, Recall et F1-Score.

Réponse:

Ces métriques sont cruciales pour les problèmes de classification déséquilibrée comme la détection de fraude.

Métrique	Formule	Question	Priorité si...
Precision	$\text{TP}/(\text{TP}+\text{FP})$	Parmi les prédicts positifs, combien sont vrais?	Coût des faux positifs élevé
Recall	$\text{TP}/(\text{TP}+\text{FN})$	Parmi les vrais positifs, combien détectés?	Coût des faux négatifs élevé
F1	$2 \times (\text{P} \times \text{R}) / (\text{P} + \text{R})$	Équilibre P et R	Besoin d'équilibre

En détection de fraude: Privilégier le **Recall** car manquer une fraude coûte plus cher que vérifier une fausse alerte.

```
from sklearn.metrics import precision_score, recall_score, f1_score

precision = precision_score(y_test, y_pred)
recall = recall_score(y_test, y_pred)
f1 = f1_score(y_test, y_pred)
```

Q24. Comment gérer un dataset déséquilibré (1% fraude)?

Réponse:

Le déséquilibre de classes est commun en banque (peu de défauts, peu de fraudes). Plusieurs techniques existent.

Technique	Description	Quand utiliser
SMOTE	Suréchantillonnage synthétique	Générer plus de cas minoritaires
Class weights	Pénaliser plus les erreurs sur minorité	Simple, dans le modèle
Undersampling	Réduire la classe majoritaire	Beaucoup de données
Ajuster le seuil	Seuil < 0.5	Privilégier le recall

```
# SMOTE
from imblearn.over_sampling import SMOTE
smote = SMOTE(random_state=42)
X_res, y_res = smote.fit_resample(X_train, y_train)

# Class weights
from sklearn.linear_model import LogisticRegression
model = LogisticRegression(class_weight='balanced')

# Ajuster le seuil
y_pred = (y_proba >= 0.3).astype(int) # Seuil plus bas
```

Q25. Qu'est-ce que l'overfitting et comment l'éviter?

Réponse:

L'overfitting se produit quand le modèle apprend trop bien les données d'entraînement, y compris le bruit, et généralise mal aux nouvelles données.

Signes: Train accuracy » Test accuracy

Solutions:

Technique	Comment
Plus de données	Réduit le risque naturellement
Régularisation	L1 (Lasso), L2 (Ridge)
Validation croisée	cv=5 ou cv=10
Early stopping	Arrêter quand test error augmente
Simplifier	Moins de features, max_depth plus petit

```
from sklearn.model_selection import cross_val_score

# Validation croisée
scores = cross_val_score(model, X, y, cv=5, scoring='roc_auc')
print(f"AUC: {scores.mean():.4f} (+/- {scores.std():.4f})")
```

Q26. Comment calculer l'Expected Loss (EL) en risque de crédit?

Réponse:

L'Expected Loss est la perte moyenne anticipée sur un portefeuille de crédit. C'est un concept fondamental en gestion des risques.

Formule: EL = PD × LGD × EAD

Composante	Définition	Valeur typique
PD	Probability of Default	1-15%
LGD	Loss Given Default	30-60%
EAD	Exposure at Default	Montant exposé

Exemple: - PD = 5%, LGD = 40%, EAD = 100,000 HTG - EL = $0.05 \times 0.40 \times 100,000 = 2,000$ HTG

Q27. Quels algorithmes pour quels cas d'usage bancaires?

Réponse:

Le choix de l'algorithme dépend du cas d'usage et des contraintes (interprétabilité, performance).

Cas d'usage	Algorithme	Raison
Scoring crédit	Régression logistique	Interprétabilité réglementaire
Détection fraude	Random Forest, XGBoost	Performance sur déséquilibre
Segmentation clients	K-Means	Non supervisé, clustering
Prédiction churn	Gradient Boosting	Bon équilibre perf/interprétabilité
Détection anomalies	Isolation Forest	Non supervisé, peu de labels

```
# Scoring crédit - Interprétable
from sklearn.linear_model import LogisticRegression
model = LogisticRegression()

# Fraude - Performance
from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier(class_weight='balanced')

# Segmentation - Clustering
from sklearn.cluster import KMeans
kmeans = KMeans(n_clusters=5)
```

SECTION 3: Statistiques (10 Questions)

Q28. Quelles sont les hypothèses de la régression linéaire (LINE)?

Réponse:

La régression linéaire repose sur quatre hypothèses fondamentales, résumées par l'acronyme **LINE**. Violer ces hypothèses peut invalider les résultats.

Hypothèse	Description	Diagnostic
Linéarité	Relation linéaire entre X et Y	Graphique résidus vs fitted
Indépendance	Observations indépendantes	Test de Durbin-Watson
Normalité	Résidus normalement distribués	QQ-plot, Shapiro-Wilk
Égalité variances	Homoscédasticité (variance constante)	Test de Breusch-Pagan

```

import statsmodels.api as sm
from statsmodels.stats.stattools import durbin_watson

# Ajuster le modèle
X_const = sm.add_constant(X)
model = sm.OLS(y, X_const).fit()

# Vérifier les hypothèses
dw = durbin_watson(model.resid) # DW ≈ 2 = OK
print(f"Durbin-Watson: {dw:.3f}")

```

Q29. Comment détecter et traiter la multicolinéarité?

Réponse:

La multicolinéarité survient quand des variables explicatives sont fortement corrélées entre elles. Elle rend les coefficients instables et difficiles à interpréter.

Détection avec VIF (Variance Inflation Factor):

VIF	Interprétation	Action
1	Pas de corrélation	Idéal
1-5	Corrélation modérée	Acceptable
5-10	Corrélation élevée	Attention
> 10	Multicolinéarité sévère	Intervention requise

```

from statsmodels.stats.outliers_influence import variance_inflation_factor

def calculate_vif(X):
    vif_data = pd.DataFrame()
    vif_data["Variable"] = X.columns
    vif_data["VIF"] = [variance_inflation_factor(X.values, i)
                      for i in range(X.shape[1])]
    return vif_data.sort_values('VIF', ascending=False)

# Si VIF > 10: supprimer une variable ou créer un indice composite

```

Q30. Quelle est la différence entre R² et R² ajusté?

Réponse:

Les deux mesurent le pouvoir explicatif du modèle, mais avec une différence cruciale pour la sélection de variables.

Métrique	Formule	Comportement
R^2	$1 - \text{SSE/SST}$	Augmente TOUJOURS avec plus de variables
R^2 ajusté	$1 - [(1-R^2)(n-1)/(n-p-1)]$	Peut DIMINUER si variable inutile

Règle: Utiliser R^2 ajusté pour comparer des modèles avec différents nombres de variables.

```
print(f"R²: {model.rsquared:.4f}")
print(f"R² ajusté: {model.rsquared_adj:.4f}")
```

Si R^2 ajusté baisse en ajoutant une variable → ne pas l'inclure

Q31. Comment interpréter un coefficient de régression?

Réponse:

Le coefficient β mesure l'effet marginal de X sur Y, toutes autres choses égales par ailleurs.

Exemple: Dans le modèle $\text{Prêt} = 50000 + 0.25 \times \text{Revenu} + 5000 \times \text{Ancienneté}$:

Coefficient	Interprétation
$\beta_0 = 50000$	Prêt de base si revenu et ancienneté = 0
$\beta_1 = 0.25$	Pour +1 HTG de revenu, le prêt augmente de 0.25 HTG
$\beta_2 = 5000$	Pour +1 an d'ancienneté, le prêt augmente de 5000 HTG

Significativité: - p-value < 0.05 → Coefficient significatif - IC ne contient pas 0 → Significatif

Q32. Comment tester la stationnarité d'une série temporelle?

Réponse:

Une série stationnaire a une moyenne et variance constantes dans le temps. C'est une condition préalable pour la plupart des modèles de prévision.

Test ADF (Augmented Dickey-Fuller): - H_0 : Série non stationnaire (racine unitaire) - H_1 : Série stationnaire

```
from statsmodels.tsa.stattools import adfuller

result = adfuller(series)
print(f"ADF Statistic: {result[0]:.4f}")
print(f"P-value: {result[1]:.4f}")

# Si p < 0.05 → Stationnaire
# Sinon → Différencier: series_diff = series.diff().dropna()
```

Si non stationnaire: Appliquer différenciation ($d=1, d=2$) avant ARIMA.

Q33. Qu'est-ce qu'un modèle ARIMA(p,d,q)?

Réponse:

ARIMA est le modèle de référence pour les séries temporelles. Il combine trois composantes.

Paramètre	Signification	Comment choisir
p	Ordre autorégressif (AR)	PACF coupe au lag p
d	Ordre de différenciation	Nb de diff. pour stationnarité
q	Ordre moyenne mobile (MA)	ACF coupe au lag q

```
from statsmodels.tsa.arima.model import ARIMA
```

```
# ARIMA(1,1,1) pour une série avec tendance
model = ARIMA(series, order=(1, 1, 1))
results = model.fit()
```

```
# Prévision
forecast = results.forecast(steps=12)
```

SARIMA: ARIMA + saisonnalité → SARIMA(p,d,q)(P,D,Q,s)

Q34. Quand utiliser Mann-Whitney au lieu du t-test?

Réponse:

Mann-Whitney U est l'alternative non-paramétrique au t-test pour comparer deux groupes indépendants.

Utiliser Mann-Whitney quand: - Distribution non normale (Shapiro-Wilk $p < 0.05$) - Échantillon petit ($n < 30$) - Données ordinaires - Présence d'outliers importants

```
from scipy.stats import mannwhitneyu, shapiro

# Vérifier normalité d'abord
_, p_normal = shapiro(groupe1)

if p_normal < 0.05: # Non normal
    stat, p = mannwhitneyu(groupe1, groupe2)
    print(f"Mann-Whitney U: p = {p:.4f}")
else: # Normal
    stat, p = ttest_ind(groupe1, groupe2)
    print(f"t-test: p = {p:.4f}")
```

Q35. Quelle est la différence entre Pearson et Spearman?

Réponse:

Les deux mesurent la corrélation mais avec des hypothèses différentes.

Aspect	Pearson (r)	Spearman (ρ)
Type de relation	Linéaire	Monotone
Type de données	Continues, normales	Ordinales ou non-normales
Sensibilité outliers	Élevée	Faible (basé sur rangs)
Interprétation	Variance partagée	Concordance des rangs

```
from scipy.stats import pearsonr, spearmanr

r_pearson, p_pearson = pearsonr(x, y)
r_spearman, p_spearman = spearmanr(x, y)

# Utiliser Spearman si:
# - Distribution asymétrique (revenus bancaires)
# - Présence d'outliers
# - Variables ordinaires
```

Règle: Pour les données bancaires (souvent asymétriques), préférer Spearman.

SECTION 4: Statistiques (10 Questions)

Q36. Quelle est la différence entre moyenne et médiane? Quand utiliser chacune?

Réponse:

La moyenne et la médiane sont deux mesures de tendance centrale qui résument différemment la position centrale des données. Le choix entre les deux dépend de la distribution des données.

Mesure	Calcul	Sensibilité aux outliers	Quand utiliser
Moyenne	Somme / N	Très sensible	Distributions symétriques
Médiane	Valeur centrale	Robuste	Distributions asymétriques, présence d'outliers

Exemple concret - Revenus: [30K, 35K, 40K, 45K, 500K] - Moyenne = 130K (fortement influencée par le 500K) - Médiane = 40K (représente mieux la situation typique)

Règle pratique: Utiliser la médiane pour les revenus, prix immobiliers, et toute variable avec des valeurs extrêmes.

Q29. Qu'est-ce que le skewness et comment l'interpréter?

Réponse:

Le skewness (asymétrie) mesure la déformation d'une distribution par rapport à une distribution symétrique. Il indique si les données sont concentrées d'un côté avec une queue étendue de l'autre.

Valeur	Interprétation	Relation entre mesures centrales
= 0	Distribution symétrique	Moyenne = Médiane = Mode
> 0	Asymétrie positive (queue à droite)	Moyenne > Médiane > Mode
< 0	Asymétrie négative (queue à gauche)	Moyenne < Médiane < Mode

Exemple bancaire: La distribution des soldes de compte a généralement un skewness > 0, car il y a beaucoup de petits soldes et quelques très gros comptes.

Q30. Expliquez la p-value et comment l'interpréter.

Réponse:

La p-value est un concept fondamental en statistique inférentielle. Elle quantifie la force de l'évidence contre l'hypothèse nulle dans un test d'hypothèse.

Définition: La p-value est la probabilité d'observer un résultat au moins aussi extrême que celui obtenu, en supposant que l'hypothèse nulle (H_0) est vraie.

p-value	Interprétation	Décision
< 0.01	Très forte évidence contre H_0	Rejeter H_0
< 0.05	Forte évidence contre H_0	Rejeter H_0 (seuil standard)
< 0.10	Évidence modérée	Marginalement significatif
≥ 0.10	Évidence faible	Ne pas rejeter H_0

Attention importante: Significatif statistiquement ne signifie pas important pratiquement. Une différence de 0.1% peut être significative avec un grand échantillon mais sans importance business.

Q31. Quelle est la différence entre corrélation et causalité?

Réponse:

Cette distinction est cruciale pour éviter de tirer des conclusions erronées à partir des données. La corrélation est une observation, la causalité est un mécanisme.

- **Corrélation:** Deux variables varient ensemble de manière prévisible
- **Causalité:** Une variable provoque directement le changement de l'autre

Exemple classique: - Observation: Les ventes de glaces et les noyades sont corrélées positivement - Fausse conclusion: Les glaces causent les noyades - Réalité: La chaleur (variable confondante) augmente les deux

Pour établir la causalité: Expérience contrôlée randomisée, analyse temporelle (la cause précède l'effet), élimination des variables confondantes.

Q32. Comment détecter les outliers?

Réponse:

Les outliers sont des observations qui s'écartent significativement du reste des données. Les détecter est essentiel avant toute analyse car ils peuvent fausser les résultats.

Méthode IQR (Interquartile Range):

$$\begin{aligned} Q1 &= 25\text{e percentile} \\ Q3 &= 75\text{e percentile} \\ IQR &= Q3 - Q1 \end{aligned}$$

Outlier si: $x < Q1 - 1.5 \times IQR$ ou $x > Q3 + 1.5 \times IQR$

Méthode Z-score:

$z = (x - \text{moyenne}) / \text{écart-type}$
Outlier si: $|z| > 3$ (plus de 3 écarts-types de la moyenne)

```
# Implémentation Python
Q1, Q3 = df['montant'].quantile([0.25, 0.75])
IQR = Q3 - Q1
outliers = df[(df['montant'] < Q1 - 1.5*IQR) | (df['montant'] > Q3 + 1.5*IQR)]
```

Q33. Qu'est-ce qu'un intervalle de confiance à 95%?

Réponse:

L'intervalle de confiance quantifie l'incertitude autour d'une estimation statistique. Il fournit une plage de valeurs plausibles pour le paramètre réel de la population.

Interprétation: Si on répétait l'échantillonnage 100 fois, 95 des intervalles calculés contiendraient la vraie valeur de la population.

Formule pour la moyenne:

$$\begin{aligned} IC &= \bar{x} \pm z \times (s/\sqrt{n}) \\ IC \text{ 95\%} &= \bar{x} \pm 1.96 \times (s/\sqrt{n}) \end{aligned}$$

Exemple bancaire: Montant moyen des transactions = 500 HTG avec IC 95% = [480, 520]. On est 95% confiant que le vrai montant moyen de la population se situe entre 480 et 520 HTG.

Q34. Quel test statistique utiliser pour comparer deux moyennes?

Réponse:

Le choix du test dépend des caractéristiques des données (normalité, indépendance) et du type de comparaison souhaitée. Voici un guide décisionnel.

Situation	Test recommandé
2 groupes indépendants, données normales	t-test pour échantillons indépendants
2 groupes indépendants, données non-normales	Mann-Whitney U (test non-paramétrique)

Situation	Test recommandé
2 mesures sur les mêmes sujets (avant/après)	t-test apparié
3+ groupes indépendants, données normales	ANOVA
3+ groupes indépendants, données non-normales	Kruskal-Wallis

Exemple bancaire: Comparer le montant moyen des prêts entre l'agence A et l'agence B → t-test pour échantillons indépendants (si normalité vérifiée).

Q35. Qu'est-ce que le coefficient de variation (CV)?

Réponse:

Le coefficient de variation est une mesure de dispersion relative qui permet de comparer la variabilité de variables ayant des unités ou des échelles différentes.

Formule: $CV = (\text{Écart-type} / \text{Moyenne}) \times 100$

Avantage: Exprimé en pourcentage, il est indépendant de l'unité de mesure.

Exemple: - Revenus: moyenne = 50,000 HTG, écart-type = 10,000 HTG → CV = 20% - Âge: moyenne = 40 ans, écart-type = 8 ans → CV = 20%

Les deux variables ont la même dispersion relative (20%) malgré des unités et des échelles complètement différentes.

Q36. Comment interpréter un coefficient de corrélation (r)?

Réponse:

Le coefficient de corrélation de Pearson (r) mesure la force et la direction de la relation linéaire entre deux variables quantitatives. Sa valeur est toujours comprise entre -1 et +1.

Valeur de r	Interprétation
0.9 à 1.0	Très forte
0.7 à 0.9	Forte
0.5 à 0.7	Modérée
0.3 à 0.5	Faible
0 à 0.3	Très faible ou inexiste

Le signe indique la direction: Positif (+) = les deux variables augmentent ensemble. Négatif (-) = quand une augmente, l'autre diminue.

R² (R-squared): Le carré de r représente le pourcentage de variance expliquée. Si $r = 0.7$, alors $r^2 = 0.49$, signifiant que 49% de la variance de Y est expliquée par X.

Q37. Qu'est-ce que l'erreur de Type I et Type II?

Réponse:

Ces deux types d'erreurs représentent les deux façons dont un test d'hypothèse peut se tromper. Comprendre ce compromis est essentiel pour interpréter les résultats de tests.

Erreur	Description	Analogie	Conséquence
Type I (α)	Rejeter H_0 alors qu'elle est vraie	Faux positif	Déclarer une fraude inexistante
Type II (β)	Ne pas rejeter H_0 alors qu'elle est fausse	Faux négatif	Manquer une vraie fraude

Relation importante: Il existe un compromis entre les deux types d'erreurs. Réduire le risque de Type I augmente le risque de Type II, et vice versa. Le choix du seuil α (souvent 0.05) reflète ce compromis.

SECTION 5: KPIs Bancaires (8 Questions)

Q38. Qu'est-ce que le NPL Ratio et pourquoi est-il important?

Réponse:

Le NPL Ratio (Non-Performing Loans Ratio) est l'indicateur clé de la qualité du portefeuille de crédit d'une banque. Il mesure la proportion des prêts en difficulté par rapport au total des prêts accordés.

Définition: NPL Ratio = Prêts non-performants / Total des prêts × 100

Qu'est-ce qu'un prêt non-performant? Un prêt est classé NPL quand le paiement est en retard de plus de 90 jours ou quand il est peu probable que l'emprunteur rembourse intégralement.

Valeur	Interprétation
< 3%	Excellent - Portefeuille de très bonne qualité
3-5%	Acceptable - Niveau de risque gérable
5-10%	Préoccupant - Nécessite une attention particulière
> 10%	Critique - Risque significatif pour la banque

Importance: Un NPL élevé signale des problèmes de qualité de crédit, peut nécessiter des provisions importantes, et impacte directement la rentabilité.

Q39. Expliquez ROE, ROA et leur différence.

Réponse:

Le ROE et le ROA sont deux indicateurs fondamentaux de rentabilité bancaire qui mesurent différents aspects de la performance financière.

Indicateur	Formule	Ce qu'il mesure	Valeur typique
ROE	Résultat Net / Capitaux Propres	Rentabilité pour les actionnaires	12-18%
ROA	Résultat Net / Total Actifs	Efficacité d'utilisation des actifs	1-2%

Relation clé: ROE = ROA × Levier financier (Actifs / Capitaux Propres)

Interprétation: Un ROE élevé peut résulter soit d'un bon ROA (bonne gestion opérationnelle), soit d'un fort levier financier (beaucoup de dettes). Il est important d'analyser les deux pour comprendre la source de la performance.

Q40. Qu'est-ce que le NIM (Net Interest Margin)?

Réponse:

Le NIM (Net Interest Margin) est l'indicateur principal de la rentabilité des activités d'intermédiation bancaire. Il mesure l'écart entre ce que la banque gagne sur ses prêts et ce qu'elle paie sur ses dépôts.

Définition: NIM = (Revenus d'intérêts - Charges d'intérêts) / Actifs productifs moyens × 100

Valeur typique: 3-5% pour les banques commerciales

Interprétation: - NIM élevé = Bonne marge sur les activités de prêt - NIM en baisse = Compression des marges (concurrence accrue, environnement de taux bas) - NIM stable = Bonne gestion de la tarification des produits

Q41. Qu'est-ce que le ratio LDR (Loan-to-Deposit)?

Réponse:

Le LDR (Loan-to-Deposit Ratio) mesure la proportion des dépôts clients utilisée pour financer les prêts. C'est un indicateur clé de liquidité et d'efficacité de l'utilisation des ressources.

Définition: LDR = Total des prêts / Total des dépôts × 100

Valeur	Interprétation
< 70%	Sous-utilisation des dépôts - La banque pourrait prêter plus
70-90%	Optimal - Bon équilibre entre liquidité et rentabilité
90-100%	Agressif - Utilisation maximale, moins de marge de sécurité
> 100%	Risqué - La banque dépend de financements externes

Q42. Qu'est-ce que le CAR (Capital Adequacy Ratio)?

Réponse:

Le CAR (Capital Adequacy Ratio) est le ratio réglementaire principal qui mesure la capacité d'une banque à absorber des pertes et à protéger les déposants. Il est étroitement surveillé par les régulateurs.

Définition: CAR = Fonds propres réglementaires / Actifs pondérés par les risques (RWA) × 100

Seuils réglementaires: - Minimum BRH (Haïti): 12% - Minimum Bâle III international: 8% - Coussin de conservation: +2.5%

Importance: Un CAR insuffisant peut entraîner des restrictions sur les activités de la banque, voire une intervention du régulateur. C'est un indicateur de solvabilité et de stabilité financière.

Q43. Comment calculer l'Expected Loss (EL)?

Réponse:

L'Expected Loss (Perte Attendue) est la perte moyenne anticipée sur un portefeuille de crédit. C'est un concept fondamental en gestion des risques de crédit.

Formule: EL = PD × LGD × EAD

Où: - **PD (Probability of Default):** Probabilité que l'emprunteur fasse défaut sur une période donnée - **LGD (Loss Given Default):** Pourcentage du montant exposé perdu en cas de défaut (1 - taux de récupération) - **EAD (Exposure at Default):** Montant total exposé au moment du défaut

Exemple de calcul: - PD = 5% (probabilité de défaut) - LGD = 40% (la banque récupère 60% en cas de défaut) - EAD = 100,000 HTG - EL = $0.05 \times 0.40 \times 100,000 = 2,000 \text{ HTG}$

Q44. Qu'est-ce que le CIR (Cost-to-Income Ratio)?

Réponse:

Le CIR (Cost-to-Income Ratio) mesure l'efficacité opérationnelle d'une banque en comparant ses charges d'exploitation à ses revenus. C'est un indicateur clé de productivité.

Définition: CIR = Charges d'exploitation / Produit Net Bancaire (PNB) × 100

Valeur	Interprétation
< 45%	Excellent - Banque très efficace
45-55%	Bon - Efficacité dans la norme
55-65%	Moyen - Potentiel d'amélioration
> 65%	Inefficace - Charges trop élevées par rapport aux revenus

Levier d'amélioration: Digitalisation des processus, optimisation des agences, automatisation.

Q45. Quels sont les red flags AML (Anti-Money Laundering)?

Réponse:

L'AML (Anti-Money Laundering) vise à détecter et prévenir le blanchiment d'argent. En tant que Data Analyst, vous devez connaître les signaux d'alerte qui peuvent indiquer des activités suspectes.

Principaux red flags à surveiller:

1. **Structuration (Smurfing):** Transactions délibérément fractionnées juste sous le seuil de déclaration (ex: plusieurs dépôts de 9,999 USD au lieu d'un seul de 30,000 USD)
 2. **Activité anormale:** Volume de transactions > moyenne historique du client sans justification économique
 3. **Pays à risque:** Transactions avec des juridictions à haut risque (liste GAFI)
 4. **Cash intensif:** Dépôts ou retraits cash fréquents et importants sans activité économique justifiant ce besoin
 5. **Shell companies:** Entreprises sans activité économique visible mais avec des flux financiers importants
 6. **Montants ronds répétitifs:** Transferts de montants exactement identiques de manière répétitive (ex: exactement 10,000 USD chaque semaine)
-

SECTION 6: Python / Pandas (8 Questions)

Q46. Comment gérer les valeurs manquantes en Pandas?

Réponse:

La gestion des valeurs manquantes est une étape cruciale du nettoyage des données. Pandas offre plusieurs méthodes pour identifier et traiter ces valeurs.

```
# ÉTAPE 1: Identifier les valeurs manquantes
df.isnull().sum()           # Compte par colonne
df.info()                    # Vue d'ensemble avec types et non-null count
df.isnull().mean() * 100     # Pourcentage de valeurs manquantes

# ÉTAPE 2: Traiter - Option A: Supprimer
df.dropna()                  # Supprime toutes les lignes avec au moins un NA
df.dropna(subset=['col1', 'col2']) # Supprime si NA dans ces colonnes spécifiques
df.dropna(thresh=3)          # Garde les lignes avec au moins 3 valeurs non-NA

# ÉTAPE 2: Traiter - Option B: Imputer
df['col'].fillna(df['col'].median())    # Médiane (recommandé pour données numériques asymétriques)
df['col'].fillna(df['col'].mean())       # Moyenne (pour données symétriques)
df['col'].fillna(df['col'].mode()[0])    # Mode (pour données catégorielles)
df['col'].fillna(method='ffill')        # Forward fill (pour séries temporelles)
```

Q47. Comment faire une agrégation groupée en Pandas?

Réponse:

L'agrégation groupée (groupby) est l'équivalent Pandas du GROUP BY en SQL. Elle permet de calculer des statistiques par catégorie.

```
# Agrégation simple - une mesure
df.groupby('region')['montant'].sum()

# Agrégations multiples sur une colonne
df.groupby('region')['montant'].agg(['sum', 'mean', 'count', 'std'])

# Agrégations différentes par colonne
df.groupby('region').agg({
    'montant': ['sum', 'mean'],
    'client_id': 'nunique',      # Nombre de clients uniques
    'date': ['min', 'max']       # Période d'activité
})

# Avec reset_index pour obtenir un DataFrame propre
result = df.groupby('region')['montant'].sum().reset_index()
result.columns = ['region', 'total_montant']
```

Q48. Comment créer une colonne conditionnelle?

Réponse:

Créer des colonnes conditionnelles permet de catégoriser ou transformer les données selon des règles métier. Pandas et NumPy offrent plusieurs approches.

```
import numpy as np

# Méthode 1: np.where - pour conditions binaires (2 valeurs possibles)
df['flag'] = np.where(df['montant'] > 1000, 'High', 'Low')

# Méthode 2: np.select - pour conditions multiples (3+ valeurs possibles)
conditions = [
    df['montant'] < 100,
    df['montant'] < 1000,
    df['montant'] >= 1000
]
choices = ['Low', 'Medium', 'High']
df['category'] = np.select(conditions, choices, default='Unknown')

# Méthode 3: apply avec fonction lambda (plus flexible, mais plus lent)
df['category'] = df['montant'].apply(lambda x: 'High' if x > 1000 else 'Low')

# Méthode 4: pd.cut pour créer des bins
df['bin'] = pd.cut(df['montant'], bins=[0, 100, 1000, np.inf], labels=['Low', 'Medium', 'High'])
```

Q49. Comment fusionner deux DataFrames?

Réponse:

La fusion de DataFrames est l'équivalent Pandas des JOIN SQL. C'est une opération fondamentale pour combiner des données de différentes sources.

```

# Merge (équivalent JOIN SQL) - basé sur des colonnes communes
df = pd.merge(clients, transactions, on='client_id', how='left')

# Merge avec colonnes de noms différents
df = pd.merge(df1, df2, left_on='id', right_on='client_id', how='inner')

# Types de merge (how):
# 'inner' = seulement les correspondances
# 'left' = tout de gauche + correspondances de droite
# 'right' = tout de droite + correspondances de gauche
# 'outer' = tout des deux côtés

# Concat - empiler des DataFrames
df = pd.concat([df1, df2], axis=0)    # Verticalement (ajouter des lignes)
df = pd.concat([df1, df2], axis=1)    # Horizontalement (ajouter des colonnes)

```

Q50. Comment créer un pivot table?

Réponse:

Les pivot tables restructurent les données pour faciliter l'analyse. Elles transforment les valeurs uniques d'une colonne en nouvelles colonnes.

```

# Pivot table avec agrégation
pivot = df.pivot_table(
    values='montant',                      # Colonne à agréger
    index='region',                        # Lignes
    columns='type',                         # Colonnes
    aggfunc='sum',                          # Fonction d'agrégation
    fill_value=0,                           # Remplacer NA par 0
    margins=True                           # Ajouter les totaux
)

# Crosstab - pour les comptages et pourcentages
pd.crosstab(df['region'], df['type'])           # Comptages
pd.crosstab(df['region'], df['type'], normalize='index') # % par ligne
pd.crosstab(df['region'], df['type'], normalize='columns') # % par colonne
pd.crosstab(df['region'], df['type'], normalize='all')   # % du total

```

Q51. Comment optimiser la mémoire d'un DataFrame?

Réponse:

L'optimisation mémoire est cruciale pour travailler avec de grands datasets. Pandas utilise par défaut des types de données généreux qui peuvent être réduits.

```

# Vérifier l'usage mémoire actuel
df.info(memory_usage='deep')
print(f"Mémoire: {df.memory_usage(deep=True).sum() / 1024**2:.2f} MB")

# Optimiser les types numériques
df['col_int'] = df['col_int'].astype('int32')          # int64 → int32 (divide par 2)
df['col_int'] = df['col_int'].astype('int16')           # Encore plus petit si valeurs < 32767

```

```

df['col_float'] = pd.to_numeric(df['col_float'], downcast='float')

# Convertir les colonnes texte à faible cardinalité en catégories
df['col_cat'] = df['col_cat'].astype('category') # Très efficace si peu de valeurs uniques

# Lire seulement les colonnes nécessaires
df = pd.read_csv('file.csv', usecols=['col1', 'col2', 'col3'])

# Lire en chunks pour fichiers très volumineux
for chunk in pd.read_csv('big_file.csv', chunksize=10000):
    process(chunk)

```

Q52. Comment calculer une moyenne mobile?

Réponse:

Les moyennes mobiles lissent les données temporelles pour révéler les tendances sous-jacentes. Elles sont essentielles en analyse financière et de séries temporelles.

```

# Moyenne mobile simple (SMA - Simple Moving Average)
df['MA_7'] = df['montant'].rolling(window=7).mean()      # Fenêtre de 7 périodes
df['MA_30'] = df['montant'].rolling(window=30).mean()     # Fenêtre de 30 périodes

# Moyenne mobile exponentielle (EMA - donne plus de poids aux valeurs récentes)
df['EMA_7'] = df['montant'].ewm(span=7).mean()

# Moyenne mobile par groupe (ex: par client)
df['MA_7_client'] = df.groupby('client_id')['montant'].transform(
    lambda x: x.rolling(7, min_periods=1).mean()
)

# Somme mobile (pour volumes cumulés sur une période)
df['sum_30d'] = df['montant'].rolling(window=30).sum()

```

Q53. Comment travailler avec les dates?

Réponse:

La manipulation des dates est omniprésente en analyse de données, surtout pour les données bancaires. Pandas offre des outils puissants via l'accessor .dt.

```

# Conversion en datetime
df['date'] = pd.to_datetime(df['date'])
df['date'] = pd.to_datetime(df['date'], format='%Y-%m-%d') # Format spécifique

# Extractions de composantes
df['year'] = df['date'].dt.year
df['month'] = df['date'].dt.month
df['day'] = df['date'].dt.day
df['day_of_week'] = df['date'].dt.dayofweek      # 0=Lundi, 6=Dimanche
df['week'] = df['date'].dt.isocalendar().week
df['quarter'] = df['date'].dt.quarter
df['is_weekend'] = df['date'].dt.dayofweek >= 5

```

```

# Filtrage par date
df[df['date'] >= '2024-01-01']
df[df['date'].between('2024-01-01', '2024-03-31')]

# Calcul de différences
df['days_since'] = (pd.Timestamp.now() - df['date']).dt.days
df['days_between'] = (df['date2'] - df['date1']).dt.days

# Resampling (agrégation temporelle)
df.set_index('date').resample('M')['montant'].sum() # Par mois
df.set_index('date').resample('W')['montant'].mean() # Par semaine

```

SECTION 7: Data Visualization & EDA (7 Questions)

Q54. Quel graphique utiliser pour quelle situation?

Réponse:

Le choix du bon graphique est essentiel pour communiquer efficacement vos insights. Chaque type de graphique est optimisé pour un type de message spécifique.

Objectif	Graphique recommandé	Quand éviter
Distribution d'une variable continue	Histogramme, KDE	Si peu de données (<30)
Comparer des catégories	Bar chart (horizontal si noms longs)	Si trop de catégories (>10)
Évolution temporelle	Line chart	Si données non ordonnées
Relation entre 2 variables	Scatter plot	Si variables catégorielles
Composition d'un tout	Pie chart (< 5 parts), Stacked bar	Si > 5-6 catégories
Comparer distributions, détecter outliers	Box plot	Si audience non-technique
Visualiser corrélations multiples	Heatmap	Si pas de pattern clair

Règle d'or: Un graphique = Un message clair. Ne pas surcharger d'informations.

Q55. Quelles sont les étapes d'une EDA?

Réponse:

L'EDA (Exploratory Data Analysis) est un processus systématique d'investigation des données avant toute analyse approfondie. Voici un framework en 7 étapes.

1. COMPRENDRE LE CONTEXTE BUSINESS
→ Quel problème essayons-nous de résoudre?

- Quelles décisions seront prises avec ces données?
 - 2. EXPLORER LA STRUCTURE
 - df.shape, df.columns, df.dtypes
 - df.head(), df.tail(), df.sample(10)
 - 3. STATISTIQUES DESCRIPTIVES
 - df.describe() pour variables numériques
 - df.describe(include='object') pour catégorielles
 - Distribution de chaque variable clé
 - 4. QUALITÉ DES DONNÉES
 - Valeurs manquantes: df.isnull().sum()
 - Doubloons: df.duplicated().sum()
 - Outliers: IQR ou Z-score
 - 5. RELATIONS ENTRE VARIABLES
 - Corrélations: df.corr()
 - Scatter plots pour paires intéressantes
 - Analyses bivariées
 - 6. SEGMENTATION ET PATTERNS
 - Patterns par groupe (région, segment client)
 - Tendances temporelles
 - Anomalies
 - 7. DOCUMENTER
 - Insights clés découverts
 - Questions pour la suite
 - Prochaines étapes recommandées
-

Q56. Comment présenter des données à des non-techniciens?

Réponse:

Communiquer des insights à un public non-technique est une compétence essentielle. Le but est de transmettre le message sans perdre l'audience dans les détails techniques.

Principes clés:

1. **Commencer par le “So What?”** - Présentez l'insight principal et son impact business d'abord, pas la méthodologie
 2. **Simplifier** - Évitez le jargon statistique. Dites “la moitié des clients” plutôt que “la médiane”
 3. **Utiliser des comparaisons concrètes** - “20% de plus que l'an dernier” est plus parlant que “un ratio de 1.2”
 4. **Visualiser** - Un bon graphique remplace souvent des pages d'explication
 5. **Contextualiser** - Expliquez si c'est bon ou mauvais et pourquoi
 6. **Recommander des actions** - Terminez toujours par des recommandations concrètes et actionnables
-

Q57. Qu'est-ce qu'un bon dashboard?

Réponse:

Un bon dashboard permet à l'utilisateur de comprendre rapidement la situation et de prendre des décisions éclairées. Il doit être à la fois informatif et facile à utiliser.

Principes de conception:

1. **Hiérarchie visuelle** - KPIs les plus importants en haut et en grand
2. **Cohérence** - Même palette de couleurs, même style de graphiques
3. **Interactivité appropriée** - Filtres utiles, mais pas trop
4. **Performance** - Chargement rapide (<3 secondes)
5. **Responsive** - Adapté au mobile si nécessaire

Structure typique recommandée:

KPIs principaux (cartes avec indicateurs)	
Tendance (line chart)	Répartition (bar chart / pie)
Tableau détaillé avec tri et filtres	

Q58. Comment gérer les outliers?

Réponse:

Les outliers sont des observations anormalement éloignées du reste des données. Leur traitement dépend du contexte et de l'objectif de l'analyse.

Options de traitement:

1. **Investiguer d'abord** - Sont-ils des erreurs de saisie ou des valeurs légitimes?
2. **Supprimer** - Si ce sont des erreurs confirmées
3. **Cap/Floor (Winsorizing)** - Remplacer les valeurs extrêmes par des percentiles (ex: 1% et 99%)
4. **Transformation** - Appliquer log ou racine carrée pour réduire l'asymétrie
5. **Analyse séparée** - Traiter les outliers comme un segment distinct

Important: En détection de fraude bancaire, les outliers SONT souvent les cas intéressants! Ne les supprimez pas aveuglément.

Q59. Comment valider la qualité des données?

Réponse:

La qualité des données est fondamentale pour des analyses fiables. Un framework de validation couvre plusieurs dimensions complémentaires.

Les 5 dimensions de la qualité:

Dimension	Question	Vérification
Complétude	Toutes les données sont-elles présentes?	Taux de valeurs manquantes
Unicité	Y a-t-il des doublons?	Comptage de duplicates
Exactitude	Les valeurs sont-elles correctes?	Plages attendues, règles métier
Cohérence	Les données sont-elles logiques?	Cross-checks entre colonnes
Fraîcheur	Les données sont-elles à jour?	Date de dernière mise à jour

```
def data_quality_report(df):
    """Génère un rapport de qualité des données"""
    return {
        'rows': len(df),
        'columns': len(df.columns),
        'duplicates': df.duplicated().sum(),
        'missing_pct': df.isnull().mean().to_dict(),
        'dtypes': df.dtypes.astype(str).to_dict()
    }
```

Q60. Comment aborder un nouveau dataset inconnu?

Réponse:

Face à un nouveau dataset, un analyste expérimenté suit un processus systématique pour comprendre rapidement les données. Voici un framework en 10 minutes.

```
# 1. STRUCTURE (30 secondes)
print(f"Dimensions: {df.shape}")
print(f"Colonnes: {df.columns.tolist()}")

# 2. TYPES DE DONNÉES (30 secondes)
print(df.dtypes)
print(df.info())

# 3. APERÇU (1 minute)
display(df.head(10))
display(df.sample(5))

# 4. STATISTIQUES DESCRIPTIVES (2 minutes)
display(df.describe())
display(df.describe(include='object'))

# 5. QUALITÉ DES DONNÉES (2 minutes)
print("Valeurs manquantes:")
print(df.isnull().sum())
print(f"\nDoublons: {df.duplicated().sum()}")

# 6. DISTRIBUTIONS (2 minutes)
df.hist(figsize=(15, 10))
plt.show()
```

```
# 7. CORRÉLATIONS (2 minutes)
plt.figure(figsize=(10, 8))
sns.heatmap(df.corr(), annot=True, cmap='coolwarm')
plt.show()
```

À la fin de ces 10 minutes, vous devez savoir: - Taille et structure du dataset - Types de variables (numériques, catégorielles, dates) - Qualité des données (missing, doublons) - Distributions générales - Premières corrélations intéressantes

Résumé Final

Top 20 des Questions les Plus Probables pour UniBank Haiti

SQL et Python: 1. **Q2 (SQL):** ROW_NUMBER vs RANK vs DENSE_RANK - Classement des clients 2. **Q8 (SQL):** Problème N+1 - Optimisation des requêtes 3. **Q47 (Python):** groupby - Agrégations par catégorie 4. **Q46 (Python):** Gestion des valeurs manquantes

Types de Variables et ML: 5. **Q16:** Variable nominale vs ordinale - Encodage approprié 6. **Q19:** MCAR/MAR/MNAR - Traitement des manquants 7. **Q20:** Supervisé vs non supervisé - Choix d'approche 8. **Q22:** Gini et AUC - Métriques de scoring 9. **Q24:** Dataset déséquilibré - SMOTE, class weights

Régression et Séries Temporelles: 10. **Q28:** Hypothèses LINE de la régression 11. **Q29:** Multicolinéarité et VIF 12. **Q30:** R^2 vs R^2 ajusté 13. **Q32:** Stationnarité et test ADF 14. **Q34:** Mann-Whitney vs t-test 15. **Q35:** Pearson vs Spearman

Statistiques: 16. **Q36:** Moyenne vs Médiane - Analyse des montants 17. **Q38:** p-value - Tests d'hypothèse 18. **Q40:** Détection des outliers - IQR, Z-score

KPIs Bancaires: 19. **Q46 (KPI):** NPL Ratio - Qualité du portefeuille 20. **Q51 (KPI):** Expected Loss = $PD \times LGD \times EAD$

Formules Essentielles à Retenir

$$\begin{aligned} \text{Gini} &= 2 \times \text{AUC} - 1 \\ \text{EL} &= PD \times LGD \times EAD \\ \text{IQR} &= Q3 - Q1 \text{ (outlier si } x < Q1 - 1.5 \times \text{IQR} \text{ ou } x > Q3 + 1.5 \times \text{IQR}) \\ \text{IC 95\%} &= \bar{x} \pm 1.96 \times (s/\sqrt{n}) \\ \text{F1} &= 2 \times (\text{Precision} \times \text{Recall}) / (\text{Precision} + \text{Recall}) \\ \text{R}^2 \text{ ajusté} &= 1 - [(1-R^2)(n-1)/(n-p-1)] \\ \text{VIF} > 10 &= \text{Multicolinéarité sévère} \\ \text{Durbin-Watson} \approx 2 &= \text{Pas d'autocorrélation} \end{aligned}$$

Mnémotechniques Clés

- **NOIR:** Nominale, Ordinale, Intervalle, Ratio
 - **LINE:** Linéarité, Indépendance, Normalité, Égalité variances
 - **MCAR/MAR/MNAR:** Types de valeurs manquantes
 - **TSCI:** Tendance, Saisonnalité, Cycle, Irrégulier
 - **PAR-FAR:** Precision, Accuracy, Recall, F1
 - **PLE:** EL = $PD \times LGD \times EAD$
 - **ARIMA(p,d,q):** AR-Integrated-MA
-

BONNE CHANCE POUR VOTRE ENTRETIEN CHEZ UNIBANK HAITI!