

Manuel des Mnémotechniques pour Examens Data Analyst

Guide de Révision Ultra-Rapide avec Techniques de Mémorisation

1. TYPES DE VARIABLES - “NOIR”

Mnémotechnique: N-O-I-R

Lettre	Type	Mémoire	Exemple
N	Nominale	Noms sans ordre	Type de compte
O	Ordinal	Ordre naturel	Rating AAA>AA>A
I	Intervalle	Intervalles égaux	Température
R	Ratio	Ratios possibles	Montant en HTG

Sous-catégories:

“**Binaire = Bi = 2 choix**” → Oui/Non, Défaut/Non-défaut

“**Polytomique = Poly = Plusieurs**” → Type compte: Épargne/Courant/DAT

Règle des statistiques:

“**Le Mode pour les Mots, la Médiane pour les rangs Ordonnés**” - Nominale → Mode - Ordinale → Médiane - Quantitative → Moyenne (si symétrique)

2. STATISTIQUES DESCRIPTIVES - “MMV-VEC-QP”

Tendance Centrale: M-M-M

- **Moyenne** (sensible aux outliers)
- **Médiane** (robuste)
- **Mode** (pour catégories)

Dispersion: V-E-C-I

- **Variance**
- **Écart-type**
- **Coefficient de Variation**
- **IQR** (InterQuartile Range)

Position: Q-P-D

- **Quartiles** (Q1, Q2, Q3)
- **Percentiles**
- **Déciles**

Règle 68-95-99.7:

“68 à 1, 95 à 2, 99.7 à 3” - 68% dans $\mu \pm 1\sigma$ - 95% dans $\mu \pm 2\sigma$ - 99.7% dans $\mu \pm 3\sigma$

3. TESTS STATISTIQUES - “2 GROUPES = t, 3+ = ANOVA”

Tableau des Tests:

Situation	Test	Mnémotechnique
2 moyennes	t-test	“Two groups = T-test”
3+ moyennes	ANOVA	“Analyse de la Variance pour Autres groupes”
Catégories	Chi-carré	“Chi pour Choix catégoriels”
Corrélation	Pearson/Spearman	Pearson = Paramétrique, Spearman = Sans normalité”

p-value:

“P petit = Preuve Présente” - $p < 0.05 \rightarrow$ Rejeter $H_0 \rightarrow$ Significatif

4. PROBABILITÉS - “BICU”

Formule de Bayes: “Bayes Inverse les Conditionnelles”

$$P(A|B) = P(B|A) \times P(A) / P(B)$$

Distributions - “BPNE”

- Bernoulli → Binaire (0 ou 1)
 - Binomiale → Beaucoup d'essais
 - Poisson → Peu fréquent (événements rares)
 - Normale → Naturelle (la plus commune)
 - Exponentielle → Espacement entre événements
-

5. SQL - “CTE avant WINDOW avant GROUP”

Window Functions: “RRDN-LLF”

- ROW_NUMBER → Unique
- RANK → Saute les rangs
- DENSE_RANK → Ne saute pas
- NTILE → Divise en groupes
- LAG → Ligne précédente
- LEAD → Ligne suivante
- FIRST/LAST_VALUE

Pattern CTE:

“WITH avant WHERE”

```
WITH cte AS (SELECT ...)  
SELECT * FROM cte WHERE ...
```

Optimisation - “SIEWL”

- **SELECT** colonnes spécifiques (pas *)
 - **INDEX** sur colonnes filtrées
 - **EXISTS** plutôt qu'**IN**
 - **WHERE** avant **HAVING**
 - **LIMIT** pour tester
-

6. MACHINE LEARNING - “SNC-CRC”

Types d'apprentissage:

- **Supervisé** → avec **Solutions** (labels)
- **Non supervisé** → **No labels**
- **Semi-supervisé** → **Some labels**

Modèles de Classification: “L-A-R-G-E”

- **Logistic Regression**
- **Arbre de décision**
- **Random Forest**
- **Gradient Boosting**
- **Ensemble methods**

Métriques: “PAR-FAR”

- **Precision** = $TP / (TP + FP)$ → **Parmi les prédicts positifs**
- **Recall** = $TP / (TP + FN)$ → **Récupérer tous les vrais positifs**
- **Accuracy** = $(TP + TN) / \text{Total}$
- **F1** = $2 \times P \times R / (P + R)$

AUC et Gini:

“Gini = 2 AUC moins 1”

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

7. RISQUE DE CRÉDIT - “PLE” (Perte Liée à l'Exposition)

Expected Loss:

“EL = PLD” (Perte Liée au Défaut)

$$\text{Expected Loss} = \text{PD} \times \text{LGD} \times \text{EAD}$$

Composant	Signification	Mnémotechnique
PD	Probability of Default	Probabilité
LGD	Loss Given Default	Loss (perte)
EAD	Exposure at Default	Exposition

Métriques de Scoring: “GAK”

- **Gini** = $2 \times \text{AUC} - 1$ (discrimination)
 - **AUC** = Aire sous ROC
 - **KS** = Kolmogorov-Smirnov (séparation)
-

8. KPIs BANCAIRES - “RRNN-NCC-CLL”

Rentabilité: “ROE-ROA-NIM-CIR”

- **ROE** = Résultat / Equity (capitaux propres)
- **ROA** = Résultat / Assets
- **NIM** = Net Interest Margin
- **CIR** = Coût / Income

Qualité: “NPC”

- **NPL** = Non-Performing Loans / Total
- Provision Coverage = Provisions / NPL
- Cost of Risk = Dotations / Encours

Solvabilité: “CAR-LDR-LCR”

- **CAR** = Capital / RWA ($\geq 12\%$ BRH)
 - **LDR** = Loans / Deposits
 - **LCR** = Liquidity Coverage Ratio
-

9. EDA - “CCANVD”

Checklist EDA:

“Comprendre, Charger, Analyser, Nettoyer, Visualiser, Documenter”

1. Comprendre le business
2. Charger les données
3. Analyser la structure
4. Nettoyer (manquants, doublons)
5. Visualiser
6. Documenter

Traitement des Outliers: “IQR × 1.5”

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

10. SEGMENTATION RFM - “Récent-Fréquent-Montant”

R-F-M:

- Recency → Jours depuis dernière activité (inversé: 5=récent)
- Frequency → Nombre de transactions
- Monetary → Montant total

Segments à retenir:

- **555** = Champions
 - **5XX** = Actifs récents
 - **X5X** = Fréquents
 - **XX5** = Gros dépenseurs
 - **111** = Perdus
-

11. DATA ENGINEERING - “EOLT”

ETL:

“Extraire-Transformer-Charger” - Extract → Sources - Transform → Nettoyage, agrégation
 - Load → Data Warehouse

OLAP vs OLTP:

“OLAP = Analyse, OLTP = Transactions” - OLAP → Requêtes complexes, agrégations - OLTP
 → Transactions rapides, CRUD

12. VISUALISATION - “HBL-SPC”

Choix du graphique:

Type de données	Graphique	Mnémotechnique
Distribution continue	Histogramme	“Histogramme pour Histoire des valeurs”
Comparer groupes	Bar chart	“Bar pour Bandes de catégories”
Tendance temporelle	Line chart	“Ligne pour le Long terme”
Relation 2 variables	Scatter	“Scatter pour Sauter entre points”
Composition	Pie	“Pie = Parts (< 6 catégories)”
Corrélation	Heatmap	“Heatmap pour Hot spots”

13. ENCODAGE - “One-Hot = Noms, Ordinal = Ordre”

Règle simple:

“One-Hot pour les Noms, Ordinal pour l’Ordre”

```

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

# Ordinale → Label Encoding ordonné
ordre = {'Faible': 1, 'Moyen': 2, 'Élevé': 3}
df['risque'] = df['risque'].map(ordre)

```

Data Leakage:

“Fit sur Train, Transform sur Test”

```

scaler.fit_transform(X_train) # FIT + TRANSFORM
scaler.transform(X_test)     # TRANSFORM seulement

```

14. FRAUDE - “Recall > Precision”

Règle d'or:

“En fraude, mieux vaut trop d'alertes que de rater une fraude”

- **Recall élevé** → Détecer toutes les fraudes
- **Precision faible** acceptable si ressources disponibles

Déséquilibre:

“SMOTE pour les PETITS” - SMOTE → Suréchantillonne la classe minoritaire

15. VALIDATION - “TTT-PSI”

Split temporel:

“Train sur le passé, Test sur le présent”

Monitoring:

“PSI > 0.25 = Problème” - $\text{PSI} < 0.10 \rightarrow \text{OK}$ - $0.10 \leq \text{PSI} < 0.25 \rightarrow \text{Attention}$ - $\text{PSI} \geq 0.25 \rightarrow \text{Action requise}$

16. FORMULES ESSENTIELLES

Statistiques:

Moyenne: $\bar{x} = \sum x_i / n$

Variance: $s^2 = \sum (x_i - \bar{x})^2 / (n-1)$

IC 95%: $\bar{x} \pm 1.96 \times (s/\sqrt{n})$

CV: $(s / \bar{x}) \times 100$

ML:

Gini = $2 \times \text{AUC} - 1$
F1 = $2 \times (\text{Precision} \times \text{Recall}) / (\text{Precision} + \text{Recall})$
EL = PD × LGD × EAD

Finance:

ROE = Résultat Net / Capitaux Propres
Variation % = (Nouveau - Ancien) / Ancien × 100
CAGR = $(V_f/V_i)^{(1/n)} - 1$

CHECKLIST FINALE - “12 Points Clés”

- 1. NOIR → Types de variables
 - 2. MMM-VECI → Tendance centrale et dispersion
 - 3. t pour 2, ANOVA pour 3+ → Tests statistiques
 - 4. $p < 0.05$ = Significatif
 - 5. EL = PD × LGD × EAD → Risque crédit
 - 6. Gini = $2 \times \text{AUC} - 1$ → Performance scoring
 - 7. ROE, ROA, NPL, CAR → KPIs bancaires
 - 8. One-Hot/Ordinal → Encodage
 - 9. Fit Train, Transform Test → Data leakage
 - 10. SMOTE pour minorité → Déséquilibre
 - 11. Recall pour fraude → Priorité détection
 - 12. PSI > 0.25 = Drift → Monitoring
-

RÉVISION EXPRESS (15 min avant l'examen)**Top 10 à ne pas oublier:**

1. **Variable nominale** = pas d'ordre (type compte)
 2. **Variable ordinale** = ordre naturel (rating)
 3. **Médiane** pour outliers/asymétrie
 4. **Chi-carré** pour indépendance catégories
 5. **Gini = $2 \times \text{AUC} - 1$**
 6. **EL = PD × LGD × EAD**
 7. **CAR $\geq 12\%$** (exigence BRH)
 8. **NPL** = prêts > 90 jours
 9. **Fit sur TRAIN uniquement**
 10. **Recall** prioritaire en fraude
-

VOUS AVEZ LES OUTILS - CONFIANCE!