
Apprentissage Automatique Supervisé

Régression & Classification

*Cours complet avec exemples réels, implémentations Python
et exercices corrigés — Pour débutants*



Dr. BADR EL KHALYLY

Avant-propos

Ce cours a été conçu pour rendre le **Machine Learning supervisé** accessible à tous, y compris aux étudiants n'ayant aucune expérience préalable en intelligence artificielle ou en programmation. Que vous soyez collégien curieux, lycéen, étudiant en première année d'université ou professionnel, ce cours est fait pour vous.

Philosophie pédagogique

Chaque algorithme est présenté selon une progression en **cinq étapes** :

1. **Hands-On d'abord** : on commence par un exemple concret et réel avant toute théorie.
2. **Intuition et analogies** : on explique l'idée avec des analogies de la vie quotidienne.
3. **Mathématiques détaillées** : chaque formule est dérivée pas à pas, en expliquant d'où elle vient et pourquoi. Même un collégien peut suivre.
4. **Implémentation sur Google Colab** : code Python complet, prêt à copier-coller et à exécuter.
5. **Exercices corrigés** : dans chaque section, des exercices mathématiques et pratiques.

Le fil rouge : un dataset clientèle

Tout au long du cours, nous utilisons un **même jeu de données clientèle** inspiré de cas réels. Ce dataset sert pour la **régression** (prédire la dépense) et la **classification** (prédire le churn).

Prérequis

- Savoir utiliser un navigateur web
- Notions de mathématiques de base (les opérations, les fonctions — niveau collège)
- Aucune connaissance en programmation requise

Dr. BADR EL KHALLYLY

Février 2026

Table des matières

Avant-propos	1
1 Introduction au Machine Learning	16
1.1 Qu'est-ce que le Machine Learning ?	16
1.1.1 Exemples concrets dans la vie quotidienne	16
1.2 Programmation classique vs Machine Learning	17
1.3 Les trois grandes familles du Machine Learning	18
1.3.1 Apprentissage supervisé (notre sujet)	18
1.4 Vocabulaire essentiel	19
1.5 Le processus du Machine Learning en 6 étapes	20
1.6 Exercice d'introduction	21
2 Python et Google Colab	22
2.1 Pourquoi Python ?	22
2.2 Google Colab : votre laboratoire gratuit dans le cloud	22
2.2.1 Étape 1 : Accéder à Google Colab	22
2.2.2 Étape 2 : Comprendre l'interface	23
2.2.3 Étape 3 : Les raccourcis clavier essentiels	23
2.3 Les bases de Python en 15 minutes	23
2.3.1 Variables et types	24
2.3.2 Listes	24
2.3.3 Boucles et conditions	24
2.3.4 Fonctions	25
2.4 Les bibliothèques essentielles du ML	25
2.4.1 NumPy — le calcul numérique	25
2.4.2 Pandas — manipulation de données	26
2.4.3 Matplotlib — visualisation	26
2.4.4 Scikit-learn — la boîte à outils du ML	27
2.5 Exercice pratique	28
3 Notre Jeu de Données : le Fil Rouge du Cours	30
3.1 Mise en situation : bienvenue dans l'entreprise !	30
3.2 Présentation du dataset clientèle	30
3.2.1 Les variables (colonnes) expliquées	30
3.2.2 Le tableau de données complet (20 clients)	31
3.3 Notation mathématique : parler le langage du ML	32

3.3.1	Les données comme un grand tableau	33
3.3.2	Notation pour un client : x_i	33
3.3.3	Notation pour la variable cible : y_i	34
3.3.4	La matrice des données X et le vecteur cible \mathbf{y}	34
3.3.5	Résumé de toutes les notations	35
3.4	Visualisation des données	36
3.4.1	Revenu vs Dépense Annuelle	36
3.4.2	Ancienneté vs Dépense Annuelle	37
3.4.3	Classification : Revenu vs Ancienneté coloré par Churn	37
3.4.4	Distribution du Churn	38
3.5	Résumé statistique des données	38
3.5.1	Rappel des formules	38
3.5.2	Tableau récapitulatif	39
3.5.3	Statistiques par groupe (Churn = 0 vs Churn = 1)	40
3.6	Implémentation sur Google Colab	40
3.6.1	Étape 1 : Créer le DataFrame	40
3.6.2	Étape 2 : Explorer les données	41
3.6.3	Étape 3 : Statistiques par groupe	42
3.6.4	Étape 4 : Créer les graphiques	42
3.6.5	Étape 5 : Matrice de corrélation	44
3.7	Exercices	45

I Régression — Prédire un Nombre 51

4 Régression Linéaire Simple 52

4.1	Hands-On : prédire le prix d'une maison	52
4.1.1	Les données	52
4.1.2	Visualisation : le nuage de points	53
4.2	Intuition : la règle et le nuage de points	54
4.2.1	L'analogie de la règle transparente	54
4.2.2	Pourquoi une droite et pas une autre forme ?	54
4.2.3	Plusieurs droites possibles — laquelle est la meilleure ?	54
4.2.4	Les étapes de l'algorithme en pratique	55
4.3	Dérivation mathématique complète	56
4.3.1	L'équation du modèle	56
4.3.2	L'erreur (résidu)	57
4.3.3	Pourquoi les erreurs au carré ? (Très important !)	58
4.3.4	La fonction de coût $J(\beta_0, \beta_1)$	59
4.3.5	Rappel : qu'est-ce qu'une dérivée ?	60

4.3.6	Dérivée partielle par rapport à β_0	61
4.3.7	Dérivée partielle par rapport à β_1	62
4.3.8	Résolution du système (trouver β_0 et β_1)	62
4.4	Application sur le dataset clientèle	64
4.4.1	Les données	65
4.4.2	Étape 1 : calculer les moyennes \bar{x} et \bar{y}	65
4.4.3	Étape 2 : calculer β_1	65
4.4.4	Étape 3 : calculer β_0	66
4.4.5	Étape 4 : écrire le modèle et interpréter	66
4.5	Coefficient de détermination R^2	67
4.5.1	Idée intuitive	67
4.5.2	Formule mathématique	68
4.5.3	Calcul de R^2 sur le dataset clientèle	68
4.6	Implémentation complète sur Google Colab	69
4.6.1	Cellule 1 : importer les bibliothèques	69
4.6.2	Cellule 2 : créer les données	70
4.6.3	Cellule 3 : entraîner le modèle	70
4.6.4	Cellule 4 : visualiser la droite de régression	71
4.6.5	Cellule 5 : visualiser les résidus	71
4.6.6	Cellule 6 : calculer et afficher R^2	72
4.6.7	Cellule 7 : faire des prédictions	72
4.6.8	Cellule 8 : vérification « à la main »	73
4.7	Workflow complet : mathématiques et code Python	73
4.8	Exercices	76
5	Régression Linéaire Multiple	84
5.1	Hands-On : Pourquoi une seule variable ne suffit pas	84
5.1.1	Régression simple : Revenu seul	85
5.1.2	Régression multiple : Revenu + Ancienneté	85
5.2	Intuition : de la droite au plan, du plan à l'hyperplan	86
5.2.1	Les étapes de l'algorithme en pratique	87
5.3	Dérivation mathématique complète	87
5.3.1	Le modèle	88
5.3.2	Les matrices : un cours express	88
5.3.3	Construction de la matrice de design X	91
5.3.4	Fonction de coût en forme matricielle	91
5.3.5	Le gradient (dérivée matricielle)	93
5.3.6	L'équation normale (Normal Equation)	94
5.3.7	Descente de gradient	94

5.4	Application sur le dataset clientèle	97
5.4.1	Construction de la matrice X	97
5.4.2	Calcul de $X^T X$ (simplifié)	97
5.4.3	Interprétation des coefficients	98
5.4.4	Prédiction pour un nouveau client	98
5.5	Implémentation sur Google Colab	98
5.5.1	Préparation des données	98
5.5.2	Comparaison régression simple vs multiple	99
5.5.3	Interprétation visuelle des coefficients	100
5.5.4	Tentative de visualisation 3D	101
5.5.5	Descente de gradient depuis zéro	101
5.5.6	Analyse des résidus	103
5.6	Workflow complet : mathématiques et code Python	104
5.7	Exercices	106
6	Régression Polynomiale	115
6.1	Hands-On : la croissance d'une ville	115
6.2	Intuition : des droites aux courbes	117
6.3	Fondements mathématiques	118
6.3.1	Qu'est-ce qu'un polynôme ?	118
6.3.2	L'astuce clé : transformer en régression linéaire multiple	120
6.3.3	Choisir le degré d	122
6.3.4	Sous-apprentissage vs Sur-apprentissage (Biais-Variance)	122
6.4	Application sur le dataset clientèle	124
6.5	Implémentation sur Google Colab	126
6.5.1	Préparation des données	126
6.5.2	Comprendre PolynomialFeatures	126
6.5.3	Pipeline : régression polynomiale en une seule ligne	127
6.5.4	Comparer les degrés 1 à 6	127
6.5.5	Trouver le degré optimal avec R^2	128
6.5.6	Détecter le sur-apprentissage avec train/test split	129
6.5.7	Validation croisée pour choisir le degré	130
6.6	Workflow complet : mathématiques et code Python	132
6.7	Exercices	134
7	Régression Non Linéaire	142
7.1	Hands-On : Prédire la croissance bactérienne	142
7.1.1	Tentative 1 : Régression linéaire	142
7.1.2	Tentative 2 : Régression polynomiale	143
7.1.3	La bonne approche : un modèle exponentiel	143

7.2	Types de relations non linéaires	144
7.2.1	Modèle exponentiel : $y = a \cdot e^{bx}$	144
7.2.2	Modèle logarithmique : $y = a \cdot \ln(x) + b$	145
7.2.3	Modèle puissance : $y = a \cdot x^b$	146
7.2.4	Modèle logistique (sigmoïde) : $y = \frac{L}{1 + e^{-k(x-x_0)}}$	147
7.3	Dérivation mathématique détaillée	148
7.3.1	L'astuce de la linéarisation	148
7.3.2	Moindres carrés pour modèles non linéaires	151
7.3.3	Comparaison : polynomiale vs non linéaire	152
7.4	Application sur le dataset clientèle	152
7.4.1	Pourquoi un modèle logarithmique ?	152
7.4.2	Le modèle	153
7.4.3	Résultats et comparaison	153
7.4.4	Interprétation des résultats	154
7.5	Implémentation sur Google Colab	154
7.5.1	Étape 1 : Imports et données	154
7.5.2	Étape 2 : Définir les modèles non linéaires	155
7.5.3	Étape 3 : Ajuster tous les modèles	155
7.5.4	Étape 4 : Visualisation comparative	156
7.5.5	Étape 5 : Analyse des résidus	157
7.5.6	Étape 6 : Linéarisation graphique	158
7.6	Schéma récapitulatif : phases d'entraînement et de test	160
7.7	Exercices	162
8	Régularisation : Ridge, Lasso et ElasticNet	170
8.1	Hands-On : quand le modèle mémorise au lieu d'apprendre	170
8.2	Intuition : pourquoi régulariser ?	172
8.2.1	Le problème : des coefficients trop grands	173
8.2.2	L'idée fondamentale	173
8.3	Développement mathématique détaillé	174
8.3.1	Le problème des grands coefficients	174
8.3.2	Régression Ridge (pénalité L2)	174
8.3.3	Régression Lasso (pénalité L1)	177
8.3.4	ElasticNet : le meilleur des deux mondes	179
8.3.5	Choisir λ : la validation croisée	180
8.4	Application sur le dataset clientèle	181
8.5	Implémentation complète sur Google Colab	183
8.6	Schéma récapitulatif : phases d'entraînement et de test	189
8.7	Exercices	191

9 Métriques d'Évaluation pour la Régression	197
9.1 Introduction : pourquoi évaluer un modèle ?	197
9.2 Hands-On : quel modèle est le meilleur ?	197
9.3 Les métriques de régression expliquées simplement	198
9.3.1 MAE — Mean Absolute Error (Erreur Absolue Moyenne)	198
9.3.2 MSE — Mean Squared Error (Erreur Quadratique Moyenne)	200
9.3.3 RMSE — Root Mean Squared Error (Racine de l'Erreur Quadratique Moyenne)	201
9.3.4 R^2 — Coefficient de Détermination	202
9.3.5 R^2 ajusté — Pénaliser les variables inutiles	205
9.3.6 MAPE — Mean Absolute Percentage Error (Erreur en Pourcentage)	206
9.4 Récapitulatif de toutes les métriques sur le Hands-On	207
9.5 Tableau de synthèse des métriques	208
9.6 Comparaison visuelle : bon modèle vs mauvais modèle	209
9.7 Application sur le dataset clientèle	209
9.8 Implémentation complète sur Google Colab	211
9.9 Exercices	215
 II Classification — Prédire une Catégorie	 223
10 Régression Logistique	224
10.1 Hands-On : prédire le défaut de paiement	224
10.1.1 Le jeu de données	225
10.1.2 Visualisation du problème	225
10.2 Intuition : pourquoi pas la régression linéaire ?	226
10.2.1 Le problème de la régression linéaire pour la classification	226
10.2.2 La solution : « écraser » les valeurs entre 0 et 1	226
10.3 Dérivation mathématique complète	227
10.3.1 La fonction sigmoïde	227
10.3.2 Le modèle de régression logistique	230
10.3.3 Pourquoi pas le MSE ? Introduction à la vraisemblance maximale	231
10.3.4 Dérivation de la log-vraisemblance	233
10.3.5 Binary Cross-Entropy (BCE)	234
10.3.6 Calcul du gradient	235
10.3.7 Les odds et le log-odds (logit)	237
10.4 Application sur le dataset clientèle	238
10.4.1 Les données	238
10.4.2 Application du modèle	238
10.4.3 La frontière de décision	239

10.4.4	Visualisation : la sigmoïde sur les données	240
10.5	Implémentation sur Google Colab	240
10.5.1	Cellule 1 : imports et données	240
10.5.2	Cellule 2 : entraînement du modèle	241
10.5.3	Cellule 3 : courbe sigmoïde sur les données	241
10.5.4	Cellule 4 : prédictions et probabilités	242
10.5.5	Cellule 5 : matrice de confusion et rapport de classification	243
10.5.6	Cellule 6 : courbe ROC et AUC	243
10.5.7	Cellule 7 : frontière de décision en 2D (bonus)	244
10.6	Schéma récapitulatif : phases d'entraînement et de test	245
10.7	Exercices	248
11	K Plus Proches Voisins (KNN)	255
11.1	Hands-On : diagnostiquer un patient	255
11.2	Intuition derrière le KNN	256
11.2.1	L'effet du choix de K	257
11.3	Développement mathématique détaillé	258
11.3.1	Les métriques de distance	258
11.3.2	L'algorithme KNN pas à pas	261
11.3.3	Mise à l'échelle des variables (<i>Feature Scaling</i>)	262
11.3.4	Choisir la valeur de K	263
11.3.5	KNN pondéré par la distance	264
11.3.6	KNN pour la régression	265
11.4	Application sur le dataset clientèle	266
11.4.1	Classification : prédire le Churn d'un nouveau client	266
11.4.2	Régression KNN : prédire la Dépense	268
11.5	Implémentation sur Google Colab	269
11.6	Schéma récapitulatif : phases d'entraînement et de test	273
11.7	Exercices	275
12	Arbres de Décision	283
12.1	Hands-On : Doit-on jouer au tennis aujourd'hui ?	283
12.2	Intuition : l'arbre de décision, un jeu de 20 questions	284
12.3	Dérivation mathématique détaillée	286
12.3.1	L'entropie de Shannon : mesurer l'incertitude	286
12.3.2	Le gain d'information	290
12.3.3	L'indice de Gini : une alternative à l'entropie	292
12.3.4	Construction de l'arbre : l'algorithme ID3/CART	293
12.3.5	L'élagage : éviter le surapprentissage	294
12.3.6	Arbres de régression	295

12.4 Application sur le dataset clientèle	296
12.5 Implémentation sur Google Colab	297
12.6 Schéma récapitulatif : phases d'entraînement et de test	302
12.7 Exercices	304
13 Random Forest (Forêts Aléatoires)	312
13.1 Mise en situation concrète	312
13.2 Intuition et idée générale	313
13.2.1 Pourquoi un seul arbre ne suffit pas ?	313
13.2.2 La solution : combiner beaucoup d'arbres	314
13.2.3 Analogie : le conseil de médecins	314
13.3 Fondements mathématiques détaillés	315
13.3.1 Échantillonnage bootstrap	315
13.3.2 Bagging (Bootstrap Aggregating)	316
13.3.3 Sélection aléatoire de variables (Random Subspace)	317
13.3.4 L'algorithme pas à pas	318
13.3.5 Importance des variables	319
13.3.6 Hyperparamètres principaux	320
13.4 Application : prédiction du Churn client	320
13.4.1 Importance des variables pour le Churn	322
13.5 Implémentation complète sur Google Colab	323
13.5.1 Étape 1 : Imports et chargement des données	323
13.5.2 Étape 2 : Comparaison Decision Tree vs Random Forest	323
13.5.3 Étape 3 : Importance des variables	324
13.5.4 Étape 4 : Score OOB	325
13.5.5 Étape 5 : Accuracy en fonction du nombre d'arbres	325
13.5.6 Étape 6 : Frontière de décision — 1 arbre vs 100 arbres	326
13.5.7 Étape 7 : Validation croisée	327
13.6 Schéma récapitulatif : phases d'entraînement et de test	328
13.7 Exercices	330
14 Machines à Vecteurs de Support (SVM)	337
14.1 Hands-On : séparer des billes sur une table	337
14.2 Intuition : construire la route la plus large	338
14.3 Dérivation mathématique détaillée	339
14.3.1 L'hyperplan séparateur	339
14.3.2 Distance d'un point à l'hyperplan	340
14.3.3 La marge	341
14.3.4 Le problème d'optimisation	342
14.3.5 SVM à marge souple (données non séparables)	343

14.3.6 L'astuce du noyau (Kernel Trick)	344
14.3.7 SVM pour la régression (SVR)	347
14.4 Application sur le dataset clientèle	347
14.5 Implémentation complète sur Google Colab	350
14.6 Diagramme récapitulatif : appliquer un SVM pas à pas	355
14.7 Schéma récapitulatif : phases d'entraînement et de test	356
14.8 Exercices corrigés	357
15 Classification Naïve Bayes	365
15.1 Hands-On : détecter les emails spam	365
15.2 Intuition : un détective qui rassemble des indices	366
15.3 Dérivation mathématique complète	367
15.3.1 Rappels de probabilités (depuis zéro)	367
15.3.2 Le théorème de Bayes : dérivation pas à pas	369
15.3.3 L'hypothèse « naïve » d'indépendance	370
15.3.4 La règle de classification	371
15.3.5 Les variantes de Naïve Bayes	372
15.3.6 Le lissage de Laplace	373
15.4 Application complète sur le dataset spam	374
15.4.1 Étape 1 : Calculer les priors	374
15.4.2 Étape 2 : Calculer les vraisemblances avec lissage de Laplace	374
15.4.3 Étape 3 : Calculer le score pour chaque classe	375
15.4.4 Étape 4 : Normaliser et décider	376
15.5 Application sur le dataset client (Gaussian Naïve Bayes)	376
15.5.1 Les données	376
15.5.2 Classons un nouveau client	377
15.6 Implémentation sur Google Colab	379
15.6.1 Naïve Bayes gaussien sur le dataset client	379
15.6.2 Classification de texte avec MultinomialNB	381
15.6.3 Comparaison avec d'autres classificateurs	382
15.7 Diagramme récapitulatif : appliquer Naïve Bayes pas à pas	384
15.8 Schéma récapitulatif : phases d'entraînement et de test	384
15.9 Exercices	385
16 Métriques d'Évaluation pour la Classification	393
16.1 Hands-On : le test COVID	393
16.2 La matrice de confusion	394
16.2.1 Visualisation de la matrice de confusion	395
16.2.2 Comment lire la matrice de confusion ?	395
16.2.3 Comprendre True/False et Positive/Negative	395

16.2.4	Vérification des totaux	396
16.3	Métriques dérivées de la matrice de confusion	396
16.3.1	Accuracy (Exactitude)	396
16.3.2	Précision (Precision)	397
16.3.3	Rappel (Recall / Sensibilité / TPR)	398
16.3.4	Spécificité (Specificity / TNR)	399
16.3.5	F1-Score	400
16.3.6	Récapitulatif des métriques sur l'exemple COVID	401
16.3.7	Le compromis Précision-Recall (Precision-Recall Tradeoff)	401
16.3.8	Courbe ROC et AUC	402
16.3.9	Extensions multi-classes	403
16.4	Quelle métrique choisir ?	405
16.5	Application sur le dataset clientèle (Churn)	406
16.6	Implémentation Google Colab	407
16.6.1	Import des bibliothèques et préparation des données	407
16.6.2	Matrice de confusion avec heatmap	408
16.6.3	Rapport de classification complet	409
16.6.4	Courbe ROC	410
16.6.5	Courbe Précision-Recall	410
16.6.6	Comparaison de tous les classifieurs	411
16.6.7	Ajustement du seuil de décision	412
16.7	Exercices	414

III Méthodes Avancées et Ensembles 421

17 Gradient Boosting, XGBoost, LightGBM et CatBoost 422

17.1	Hands-On : prédire un prix en corrigeant ses erreurs	422
17.1.1	Première tentative : une estimation grossière	422
17.1.2	Deuxième tentative : corriger l'erreur	423
17.1.3	Troisième tentative : affiner encore	423
17.2	Intuition : apprendre de ses erreurs	425
17.2.1	Boosting vs Bagging : deux philosophies opposées	426
17.3	Dérivation mathématique détaillée	427
17.3.1	L'idée du Boosting	427
17.3.2	Pourquoi « Gradient » Boosting ?	428
17.3.3	L'algorithme pas à pas	429
17.3.4	Régularisation dans le Gradient Boosting	429
17.3.5	XGBoost — eXtreme Gradient Boosting	430
17.3.6	LightGBM — Light Gradient Boosting Machine	431

17.3.7	CatBoost — Categorical Boosting	432
17.3.8	Tableau comparatif	433
17.4	Application sur le dataset clientèle	433
17.4.1	Entraînement et comparaison avec Random Forest	433
17.4.2	Importance des variables	434
17.4.3	Courbe d'apprentissage : accuracy vs nombre d'arbres	435
17.5	Implémentation complète sur Google Colab	436
17.5.1	Installation des bibliothèques	436
17.5.2	Préparation des données	436
17.5.3	Entraînement des 4 modèles	437
17.5.4	Comparaison visuelle des résultats	439
17.5.5	Importance des variables pour chaque modèle	440
17.5.6	Optimisation des hyperparamètres avec GridSearchCV	441
17.5.7	Courbes d'apprentissage comparées	442
17.6	Schéma récapitulatif : phases d'entraînement et de test	443
17.7	Exercices	445
18	Réseaux de Neurones (MLP — Perceptron Multi-Couches)	453
18.1	Hands-On : Reconnaissance de chiffres manuscrits	453
18.2	Intuition : du neurone biologique au neurone artificiel	454
18.3	Dérivation mathématique détaillée	456
18.3.1	Le neurone unique : le perceptron	456
18.3.2	Fonctions d'activation	457
18.3.3	Le Perceptron Multi-Couches (MLP)	458
18.3.4	Propagation avant : exemple numérique pas à pas	459
18.3.5	Fonctions de perte (loss functions)	460
18.3.6	Rétropropagation (Backpropagation)	461
18.3.7	Considérations pratiques	463
18.4	Application sur le dataset clientèle	464
18.5	Implémentation sur Google Colab	465
18.5.1	Étape 1 : Imports et préparation des données	465
18.5.2	Étape 2 : Préparation et normalisation	466
18.5.3	Étape 3 : Entraîner le MLPClassifier	467
18.5.4	Étape 4 : Comparer plusieurs architectures	468
18.5.5	Étape 5 : Courbe d'apprentissage (loss curve)	469
18.5.6	Étape 6 : Prédiction sur de nouveaux clients	469
18.5.7	Étape 7 : Comparaison avec les autres classifieurs	470
18.6	Schéma récapitulatif : phases d'entraînement et de test	472
18.7	Exercices	473

IV Synthèse, Validation et Pratique	481
19 Validation Croisée et Sélection de Modèle	482
19.1 Hands-On : mémoriser n'est pas apprendre	482
19.2 Séparation Train / Test	483
19.2.1 Pourquoi séparer les données ?	483
19.2.2 Séparation aléatoire	483
19.2.3 Séparation stratifiée	484
19.2.4 Le problème : la variance du score	484
19.3 Validation Croisée K-Fold	485
19.3.1 L'idée de la validation croisée	485
19.3.2 Formulation mathématique	486
19.3.3 Choisir K : combien de folds ?	486
19.3.4 Validation croisée stratifiée	487
19.4 Réglage des hyperparamètres	487
19.4.1 Paramètres vs hyperparamètres	487
19.4.2 Grid Search (recherche par grille)	488
19.4.3 Random Search (recherche aléatoire)	489
19.4.4 Séparation Train / Validation / Test	490
19.5 Courbes d'apprentissage	490
19.6 Application sur le dataset clientèle	491
19.7 Implémentation complète sur Google Colab	494
19.8 Exercices	502
20 Guide de Choix d'Algorithme	510
20.1 Introduction : face à la jungle des algorithmes	510
20.2 Tableau comparatif des algorithmes	510
20.3 Arbre de décision pour le choix d'algorithme	512
20.4 Questions à se poser avant de choisir	512
20.5 Forces et faiblesses de chaque algorithme	514
20.6 Stratégie pratique systématique	517
20.7 Implémentation sur Google Colab	518
20.7.1 Étape 1 : Préparation des données	518
20.7.2 Étape 2 : Définition des modèles et validation croisée	519
20.7.3 Étape 3 : Visualisation comparative	521
20.7.4 Étape 4 : Tableau récapitulatif avec plusieurs métriques	522
20.7.5 Étape 5 : LazyPredict pour une comparaison rapide (bonus)	523
20.8 Exercices	524
20.9 Résumé du chapitre	533

21 Projet Complet : Prédiction du Churn Client de A à Z	534
21.1 Introduction	534
21.2 Étape 1 : Comprendre le problème métier	536
21.2.1 Le contexte de l'entreprise	536
21.2.2 Définir l'objectif ML	536
21.2.3 Choisir la métrique de succès	536
21.3 Étape 2 : Collecter et explorer les données	538
21.3.1 Création du jeu de données	538
21.3.2 Exploration initiale	539
21.3.3 Analyse de la variable cible	540
21.3.4 Analyse des corrélations	541
21.3.5 Comparaison churners vs non-churners	541
21.4 Étape 3 : Préparer les données	543
21.4.1 Gestion des valeurs manquantes	543
21.4.2 Séparation features / cible et standardisation	543
21.5 Étape 4 : Entraîner les modèles	546
21.5.1 Régression Logistique	546
21.5.2 KNN (K-Nearest Neighbors)	546
21.5.3 Arbre de Décision	546
21.5.4 Random Forest	547
21.5.5 SVM (Support Vector Machine)	547
21.5.6 Gradient Boosting	547
21.5.7 XGBoost	548
21.5.8 Tableau comparatif des accuracy	548
21.6 Étape 5 : Évaluer et comparer en profondeur	550
21.6.1 Validation croisée (5-fold)	550
21.6.2 Matrice de confusion et rapport de classification	551
21.6.3 Courbes ROC comparées	552
21.6.4 Comparaison complète multi-métriques	553
21.7 Étape 6 : Optimiser le meilleur modèle	555
21.7.1 Optimisation par GridSearchCV	555
21.7.2 Évaluation du modèle optimisé	556
21.7.3 Importance des features	556
21.7.4 Matrice de confusion finale	557
21.8 Étape 7 : Interpréter et communiquer les résultats	559
21.8.1 Résumé pour le directeur marketing	559
21.8.2 Calcul du ROI (Retour sur Investissement)	560
21.8.3 Identification des clients à risque	561
21.9 Synthèse du chapitre	563

21.10Exercices 564

Références Bibliographiques **572**

Chapitre 1

Introduction au Machine Learning

1.1 Qu'est-ce que le Machine Learning ?

Imagine cette situation...

Tu as 10 ans et ta mère te montre des photos d'animaux. Elle te dit : « Ça c'est un chat », « Ça c'est un chien ». Après 100 photos, elle te montre une nouvelle photo sans rien dire. Tu réponds immédiatement : « C'est un chat ! ».

Comment as-tu fait ? Personne ne t'a donné de règle précise comme « si l'animal a des moustaches longues et des oreilles pointues, c'est un chat ». Ton cerveau a **appris tout seul** en voyant beaucoup d'exemples. C'est exactement ce que fait le Machine Learning, mais avec un ordinateur.

Définition — Machine Learning

Le **Machine Learning** (apprentissage automatique) est une technique qui permet à un ordinateur d'**apprendre à partir d'exemples** (les **données**) plutôt que d'être programmé manuellement pour chaque situation.

L'ordinateur découvre des **règles cachées** (des patterns) dans les données et les utilise pour faire des **prédictions** sur de nouvelles situations qu'il n'a jamais vues.

1.1.1 Exemples concrets dans la vie quotidienne

Le Machine Learning est partout autour de nous, même si on ne le voit pas :

Application	Ce que fait le ML	Exemple concret
Netflix / YouTube	Prédit ce que tu veux regarder	« Parce que vous avez aimé... »
Gmail	Détecte les spams	Filtre automatique des emails indésirables
Google Maps	Prédit le temps de trajet	« 25 min avec le trafic actuel »
Siri / Alexa	Comprend ta voix	Reconnaissance vocale
Instagram	Reconnaît les visages	Suggestion de tag sur les photos
Banque	Détecte la fraude	Blocage d'un achat suspect
Météo	Prédit la température	Prévisions à 7 jours

TABLE 1.1 – Le Machine Learning dans la vie de tous les jours.

1.2 Programmation classique vs Machine Learning

Pour bien comprendre la différence, prenons un exemple concret.

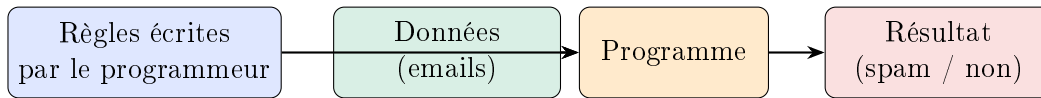
Exemple : détecter un email spam

Approche classique (sans ML) : Un programmeur écrit des règles à la main :

- SI l'email contient « gagné un million » \Rightarrow spam
- SI l'expéditeur est inconnu ET contient un lien \Rightarrow spam
- SI l'objet est en MAJUSCULES \Rightarrow spam
- ... (des centaines de règles à écrire et maintenir !)

Approche ML : On donne à l'ordinateur des milliers d'emails déjà classés (spam / pas spam). L'ordinateur **découvre lui-même** les règles !

Programmation classique



Machine Learning

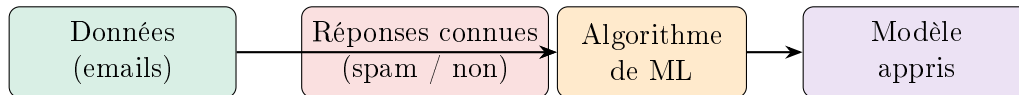


FIGURE 1.1 – En programmation classique, le programmeur écrit les règles. En ML, l'ordinateur les découvre.

1.3 Les trois grandes familles du Machine Learning

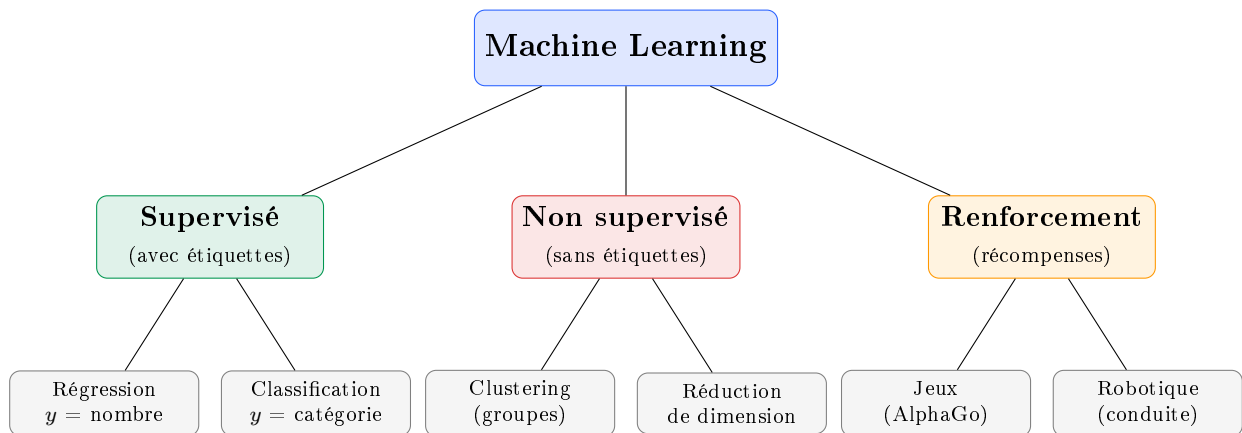


FIGURE 1.2 – Les trois familles du Machine Learning et leurs sous-types.

1.3.1 Apprentissage supervisé (notre sujet)

Apprentissage supervisé

Dans l'apprentissage **supervisé**, on donne à l'ordinateur des exemples **avec les réponses correctes** (les **étiquettes**). C'est comme un élève qui apprend avec un professeur qui lui donne les corrections.

- **Régression** : la réponse est un **nombre** (prix, température, dépense...)
- **Classification** : la réponse est une **catégorie** (spam/pas spam, malade/sain, chat/-chien...)

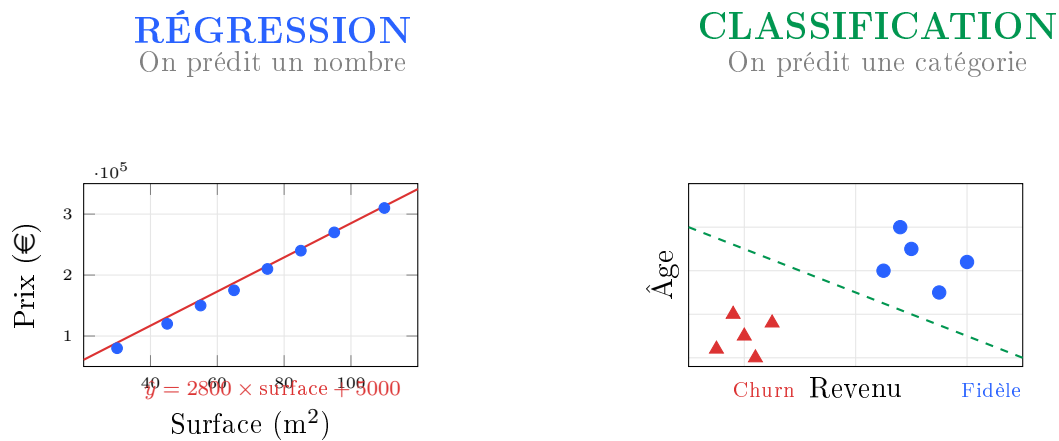


FIGURE 1.3 – Régression (prédire un prix) vs Classification (prédire une catégorie).

1.4 Vocabulaire essentiel

Avant d'aller plus loin, apprenons le vocabulaire que nous utiliserons tout au long du cours.

Terme	Explication simple	Exemple
Donnée (sample)	Un exemple, une ligne dans un tableau	Un client
Variable / Feature	Une caractéristique mesurée	L'âge, le revenu
Cible / Label (y)	La réponse qu'on cherche à prédire	Dépense, Churn
Modèle	La « formule » apprise	$\hat{y} = 3x + 2$
Entraînement	L'ordinateur apprend les règles	Voir les exemples
Prédiction	Appliquer le modèle sur du nouveau	Prédire un nouveau client
Train set	Données pour apprendre	80% des données
Test set	Données pour évaluer	20% des données
Paramètres	Valeurs apprises par le modèle	Les coefficients β
Hyperparamètres	Valeurs choisies par l'humain	Le nombre de voisins K
Overfitting	Le modèle mémorise au lieu d'apprendre	Trop coller aux données
Underfitting	Le modèle est trop simple	Rater les tendances

TABLE 1.2 – Vocabulaire fondamental — à connaître par cœur !

1.5 Le processus du Machine Learning en 6 étapes

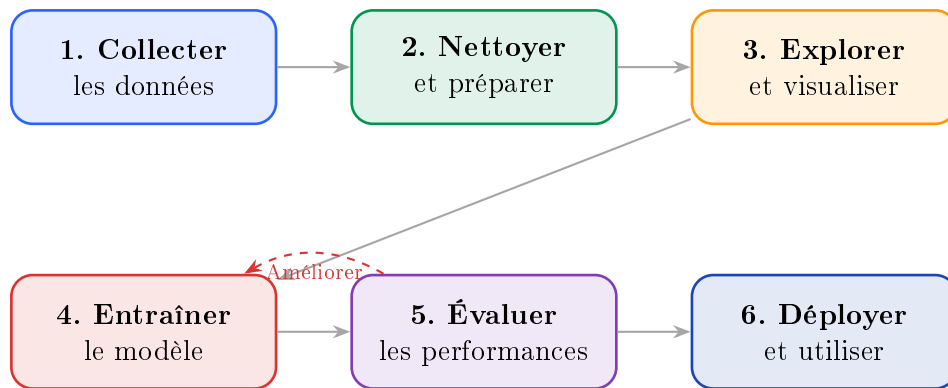


FIGURE 1.4 – Les 6 étapes d'un projet de Machine Learning.

1.6 Exercice d'introduction

Exercice 1.1 — Identifier le type de problème

Pour chaque situation, indiquez s'il s'agit de **régression**, de **classification** ou de **non supervisé** :

1. Prédire le prix d'une maison à partir de sa surface et du nombre de chambres.
2. Déterminer si un email est un spam ou non.
3. Regrouper des clients similaires pour le marketing.
4. Prédire la note d'un étudiant à l'examen final.
5. Déterminer si une tumeur est maligne ou bénigne.
6. Prédire le nombre de visiteurs d'un site web demain.
7. Classer des images comme « chat », « chien » ou « oiseau ».
8. Estimer l'âge d'une personne à partir de sa photo.

Correction de l'exercice 1.1

1. **Régression** — le prix est un nombre continu.
2. **Classification** (binaire) — deux catégories : spam / non spam.
3. **Non supervisé** (clustering) — pas d'étiquettes, on cherche des groupes.
4. **Régression** — la note est un nombre (ex : 14,5/20).
5. **Classification** (binaire) — deux catégories : maligne / bénigne.
6. **Régression** — le nombre de visiteurs est un nombre.
7. **Classification** (multi-classes) — trois catégories possibles.
8. **Régression** — l'âge est un nombre continu.

Chapitre 2

Python et Google Colab

2.1 Pourquoi Python ?

Python est **le** langage de référence en Machine Learning. Voici pourquoi :

- **Simple à apprendre** : la syntaxe ressemble à de l'anglais courant
- **Bibliothèques puissantes** : NumPy, Pandas, Scikit-learn, Matplotlib...
- **Communauté immense** : des millions de tutoriels et d'exemples gratuits
- **Gratuit** : aucun logiciel payant nécessaire

2.2 Google Colab : votre laboratoire gratuit dans le cloud

Google Colaboratory (Colab)

Google Colab est un environnement **gratuit** qui fonctionne dans votre navigateur web. Il vous permet d'écrire et d'exécuter du code Python **sans rien installer** sur votre ordinateur.

Avantages :

- Python et toutes les bibliothèques de ML sont déjà installées
- Accès gratuit à des GPU (processeurs graphiques puissants)
- Sauvegarde automatique sur Google Drive
- Partage facile avec d'autres personnes

2.2.1 Étape 1 : Accéder à Google Colab

1. Ouvrez votre navigateur (Chrome, Firefox, Edge...)
2. Allez sur <https://colab.research.google.com>
3. Connectez-vous avec votre compte Google (créez-en un si nécessaire, c'est gratuit)
4. Cliquez sur « **Nouveau notebook** » (ou « New notebook » en anglais)

2.2.2 Étape 2 : Comprendre l'interface

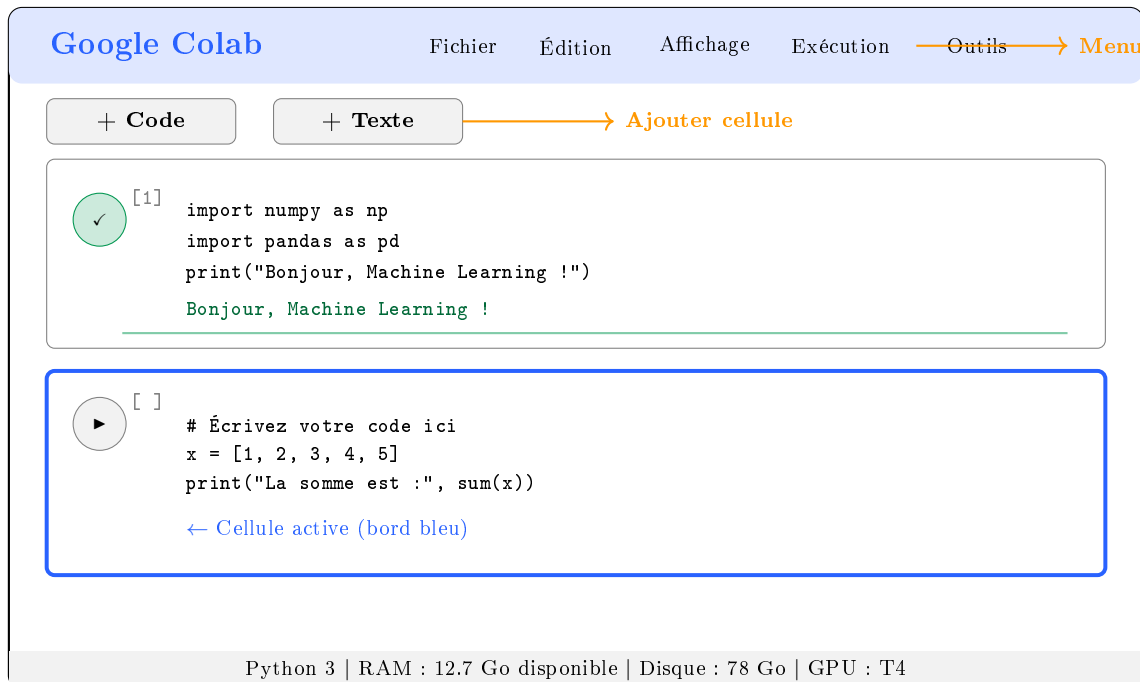


FIGURE 2.1 – L'interface de Google Colab — chaque cellule est un bloc de code indépendant.

2.2.3 Étape 3 : Les raccourcis clavier essentiels

Raccourci	Action
Ctrl + Entrée	Exécuter la cellule courante
Shift + Entrée	Exécuter et passer à la cellule suivante
Ctrl + M puis B	Insérer une cellule en dessous
Ctrl + M puis A	Insérer une cellule au-dessus
Ctrl + M puis D	Supprimer la cellule
Ctrl + S	Sauvegarder
Ctrl + Z	Annuler la dernière action

TABLE 2.1 – Raccourcis clavier essentiels de Google Colab.

2.3 Les bases de Python en 15 minutes

Voici le minimum de Python nécessaire pour suivre ce cours. Tapez chaque exemple dans une cellule de Colab et exécutez-le.

2.3.1 Variables et types

```
1  # Un nombre entier
2  age = 25
3  print("Age :", age)           # Affiche : Age : 25
4
5  # Un nombre decimal (float)
6  salaire = 2500.50
7  print("Salaire :", salaire)   # Affiche : Salaire : 2500.5
8
9  # Une chaine de caracteres (texte)
10 nom = "Ahmed"
11 print("Nom :", nom)           # Affiche : Nom : Ahmed
12
13 # Un booleen (Vrai ou Faux)
14 est_client = True
15 print("Client :", est_client) # Affiche : Client : True
```

Listing 2.1 – Variables et types de base en Python

2.3.2 Listes

```
1  # Une liste = collection ordonnee de valeurs
2  notes = [12, 15, 8, 18, 14]
3  print("Notes :", notes)
4  print("Premiere note :", notes[0])    # 12 (on compte a partir de
    0 !)
5  print("Derniere note :", notes[-1])    # 14
6  print("Nombre de notes :", len(notes)) # 5
7  print("Moyenne :", sum(notes)/len(notes)) # 13.4
```

Listing 2.2 – Les listes en Python

2.3.3 Boucles et conditions

```
1  # Boucle for : repeter une action
2  for note in [12, 15, 8, 18, 14]:
3      if note >= 10:
4          print(note, "-> Reussi")
5      else:
```

```

6         print(note, "-> Echoue")
7
8     # Resultat :
9     # 12 -> Reussi
10    # 15 -> Reussi
11    # 8 -> Echoue
12    # 18 -> Reussi
13    # 14 -> Reussi

```

Listing 2.3 – Boucles et conditions

2.3.4 Fonctions

```

1  # Definir une fonction
2  def calculer_moyenne(liste):
3      """Calcule la moyenne d'une liste de nombres."""
4      return sum(liste) / len(liste)
5
6  # Utiliser la fonction
7  notes = [12, 15, 8, 18, 14]
8  moy = calculer_moyenne(notes)
9  print("Moyenne :", moy)  # Affiche : Moyenne : 13.4

```

Listing 2.4 – Créer une fonction

2.4 Les bibliothèques essentielles du ML

2.4.1 NumPy — le calcul numérique

```

1  import numpy as np
2
3  # Creer un tableau (array)
4  x = np.array([1, 2, 3, 4, 5])
5  print("Tableau :", x)
6  print("Moyenne :", np.mean(x))      # 3.0
7  print("Ecart-type :", np.std(x))    # 1.414...
8  print("x au carre :", x**2)         # [1, 4, 9, 16, 25]
9  print("x + 10 :", x + 10)           # [11, 12, 13, 14, 15]

```

Listing 2.5 – NumPy : calcul sur des tableaux de nombres

2.4.2 Pandas — manipulation de données

```
1 import pandas as pd
2
3 # Creer un DataFrame (tableau structure)
4 data = {
5     'Nom':      ['Fatima', 'Ahmed', 'Youssef', 'Khadija'],
6     'Age':      [25, 30, 35, 28],
7     'Salaire':  [2500, 3200, 4100, 2800]
8 }
9 df = pd.DataFrame(data)
10 print(df)
11 #      Nom  Age  Salaire
12 # 0  Fatima   25    2500
13 # 1   Ahmed   30    3200
14 # 2 Youssef   35    4100
15 # 3 Khadija   28    2800
16
17 print("Moyenne des ages :", df['Age'].mean())      # 29.5
18 print("Salaire max :", df['Salaire'].max())        # 4100
19 print("Clients > 28 ans :\n", df[df['Age'] > 28]) # Ahmed,
    Youssef
```

Listing 2.6 – Pandas : travailler avec des tableaux de données

2.4.3 Matplotlib — visualisation

```
1 import matplotlib.pyplot as plt
2
3 # Graphique simple
4 x = [1, 2, 3, 4, 5]
5 y = [2, 4, 5, 4, 6]
6
7 plt.figure(figsize=(8, 5))
8 plt.scatter(x, y, color='blue', s=100, label='Points')
9 plt.plot(x, y, color='red', linestyle='--', label='Ligne')
10 plt.xlabel('Variable x')
11 plt.ylabel('Variable y')
12 plt.title('Mon premier graphique')
13 plt.legend()
```

```

14 plt.grid(True, alpha=0.3)
15 plt.show()

```

Listing 2.7 – Matplotlib : créer des graphiques

2.4.4 Scikit-learn — la boîte à outils du ML

```

1 from sklearn.linear_model import LinearRegression
2 from sklearn.model_selection import train_test_split
3 from sklearn.metrics import mean_squared_error
4
5 # Scikit-learn suit toujours le meme schema :
6 # 1. Creer le modele      : model = NomAlgorithme()
7 # 2. Entrainer le modele  : model.fit(X_train, y_train)
8 # 3. Faire des predictions : y_pred = model.predict(X_test)
9 # 4. Evaluer              : score = model.score(X_test, y_test)

```

Listing 2.8 – Scikit-learn : la bibliothèque ML de référence

Le pattern universel de Scikit-learn

Tous les algorithmes de Scikit-learn suivent exactement le même schéma en 4 lignes. Que ce soit une régression linéaire, un KNN, un SVM ou un arbre de décision, le code est presque identique ! Seule la première ligne change (le nom de l'algorithme).

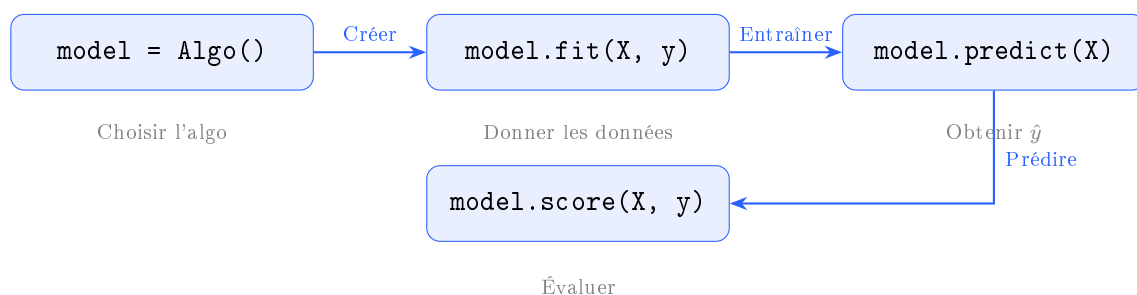


FIGURE 2.2 – Le pattern universel de Scikit-learn — même structure pour tous les algorithmes.

2.5 Exercice pratique

Exercice 2.1 — Premier notebook Colab

Ouvrez Google Colab et réalisez les étapes suivantes :

1. Créez un nouveau notebook et renommez-le « MonPremierML ».
2. Dans la première cellule, importez NumPy, Pandas et Matplotlib.
3. Créez un DataFrame avec 5 étudiants (Nom, Note_Math, Note_Physique).
4. Affichez la moyenne de chaque matière.
5. Créez un graphique en barres des notes de mathématiques.

Correction de l'exercice 2.1

```
# Cellule 1 : Imports
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

# Cellule 2 : Creer le DataFrame
etudiants = {
    'Nom': ['Fatima', 'Ahmed', 'Youssef', 'Khadija', 'Amina'],
    'Note_Math': [15, 12, 18, 9, 14],
    'Note_Physique': [13, 16, 14, 11, 17]
}
df = pd.DataFrame(etudiants)
print(df)

# Cellule 3 : Moyennes
print("Moyenne Math :", df['Note_Math'].mean())      # 13.6
print("Moyenne Physique :", df['Note_Physique'].mean()) # 14.2

# Cellule 4 : Graphique
plt.figure(figsize=(8, 5))
plt.bar(df['Nom'], df['Note_Math'], color='steelblue')
plt.xlabel('Etudiant')
plt.ylabel('Note de Mathematiques')
plt.title('Notes de Mathematiques par etudiant')
plt.ylim(0, 20)
plt.grid(axis='y', alpha=0.3)
plt.show()
```



Chapitre 3

Notre Jeu de Données : le Fil Rouge du Cours

3.1 Mise en situation : bienvenue dans l'entreprise !

Hands-On : Votre premier jour en tant que Data Scientist

Imaginez que vous venez d'être embauché(e) comme **Data Scientist junior** dans une entreprise d'e-commerce appelée **ShopTech**. Votre manager vous accueille et vous explique :

« Nous perdons des clients chaque mois et nous ne savons pas pourquoi. Nous avons aussi du mal à prédire combien chaque client va dépenser. J'ai besoin de vous pour deux missions :

1. **Régression** : Prédire la **dépense annuelle** de chaque client (en euros).
2. **Classification** : Prédire si un client va **partir** ($\text{churn} = 1$) ou **rester** ($\text{churn} = 0$).

Voici un extrait de notre base de données avec 20 clients. Explorez-les, comprenez-les, et ensuite nous appliquerons différents algorithmes de Machine Learning pour résoudre ces problèmes. »

Ce jeu de données sera utilisé **tout au long du cours** — c'est notre **fil rouge**. Chaque algorithme que nous étudierons sera appliqué sur ces mêmes données, ce qui vous permettra de comparer facilement les méthodes.

3.2 Présentation du dataset clientèle

3.2.1 Les variables (colonnes) expliquées

Avant de regarder les données, comprenons ce que chaque colonne représente :

Variable	Description	Type / Unité
Client_ID	Identifiant unique du client	Entier (1 à 20)
Age	Âge du client en années	Entier (années)
Revenu	Revenu annuel du client	Nombre (k€/an)
Anciennete	Nombre d'années comme client	Nombre (années)
Nb_Achats	Nombre total d'achats effectués	Entier
Montant_Moyen	Montant moyen par achat	Nombre (€)
Depense_Annuelle	Dépense totale sur l'année	Nombre (€)
Score_Satisfaction	Score de satisfaction (enquête)	Entier (1 à 10)
Churn	Le client est-il parti ?	0 = Non, 1 = Oui

TABLE 3.1 – Description des 9 variables de notre dataset.

Variables d'entrée vs variable de sortie

- Les variables **d'entrée** (features) sont : Age, Revenu, Ancienneté, Nb_Achats, Montant_Moyen, Score_Satisfaction.
- La variable de **sortie pour la régression** est : **Depense_Annuelle** (un nombre continu).
- La variable de **sortie pour la classification** est : **Churn** (une catégorie : 0 ou 1).
- **Client_ID** n'est **jamais** utilisé comme feature — c'est juste un identifiant.

3.2.2 Le tableau de données complet (20 clients)

Voici nos 20 clients. Prenez le temps de lire chaque ligne et d'observer les tendances.

ID	Age	Revenu (k€)	Ancien. (ans)	Nb_Ach.	Mt_Moy (€)	Dép_Ann. (€)	Satisf. (1-10)	Churn
1	25	22	1	8	45	360	5	1
2	34	45	5	30	78	2340	8	0
3	28	28	2	12	52	624	4	1
4	45	62	8	55	95	5225	9	0
5	23	18	1	5	38	190	3	1
6	52	75	12	70	110	7700	9	0
7	31	35	3	18	62	1116	6	0
8	41	55	7	45	88	3960	8	0
9	27	24	1	7	42	294	4	1
10	38	48	6	35	82	2870	7	0
11	55	80	15	85	120	10200	10	0
12	29	30	2	14	55	770	5	1
13	47	65	9	60	98	5880	9	0
14	22	20	1	6	40	240	3	1
15	36	42	4	25	72	1800	7	0
16	43	58	7	48	90	4320	8	0
17	26	25	1	9	48	432	4	1
18	50	70	10	65	105	6825	9	0
19	33	38	3	20	65	1300	6	0
20	30	32	2	15	58	870	5	1

TABLE 3.2 – Notre dataset complet de 20 clients — le fil rouge du cours.

Observez les tendances !

Même sans algorithme, on peut déjà remarquer des choses en lisant le tableau :

- Les clients qui ont **churné** (Churn = 1) ont tendance à être **jeunes**, avec un **faible revenu**, une **faible ancienneté** et un **score de satisfaction bas**.
- Les clients **fidèles** (Churn = 0) sont généralement **plus âgés**, avec un **revenu plus élevé**, une **ancienneté plus longue** et un **score de satisfaction haut**.
- La **dépense annuelle** semble directement liée au revenu et au nombre d'achats.

Ces intuitions sont exactement ce que les algorithmes de ML vont formaliser mathématiquement !

3.3 Notation mathématique : parler le langage du ML

Pour pouvoir lire les formules dans les chapitres suivants, il faut maîtriser quelques notations simples. Pas de panique : tout est expliqué comme si vous aviez 12 ans !

3.3.1 Les données comme un grand tableau

Le dataset

Un jeu de données est un **tableau** avec :

- n = le **nombre de lignes** (le nombre d'exemples, d'observations, d'échantillons). Chez nous, $n = 20$ (20 clients).
- p = le **nombre de colonnes de features** (les caractéristiques utilisées pour prédire). Chez nous, $p = 6$ (Age, Revenu, Ancienneté, Nb_Achats, Montant_Moyen, Score_Satisfaction).

Exemple avec nos données

Notre dataset a $n = 20$ clients et $p = 6$ features. On peut le représenter comme un tableau de 20 lignes et 6 colonnes (sans compter l'ID et la variable cible).

3.3.2 Notation pour un client : x_i

Vecteur de features \mathbf{x}_i

Chaque client i (avec i allant de 1 à n) est représenté par un **vecteur de features** :

$$\mathbf{x}_i = \begin{pmatrix} x_{i,1} \\ x_{i,2} \\ \vdots \\ x_{i,p} \end{pmatrix}$$

C'est simplement la **liste des valeurs** de ses caractéristiques.

Exemple : le client n°2

Le client n°2 a : Age = 34, Revenu = 45, Ancienneté = 5, Nb_Achats = 30, Montant_Moyen = 78, Score_Satisfaction = 8. Son vecteur est :

$$\mathbf{x}_2 = \begin{pmatrix} 34 \\ 45 \\ 5 \\ 30 \\ 78 \\ 8 \end{pmatrix}$$

Chaque nombre correspond à une feature : $x_{2,1} = 34$ (âge), $x_{2,2} = 45$ (revenu), $x_{2,3} = 5$ (ancienneté), etc.

3.3.3 Notation pour la variable cible : y_i

Variable cible y_i

Pour chaque client i , la **variable cible** (ou **label**) est notée y_i . C'est la valeur qu'on cherche à prédire.

- **En régression** : y_i = dépense annuelle du client i . Par exemple, $y_2 = 2340$ €.
- **En classification** : y_i = churn du client i . Par exemple, $y_2 = 0$ (fidèle).

3.3.4 La matrice des données X et le vecteur cible y

Matrice X et vecteur y

On note X la matrice qui contient **toutes les features de tous les clients** :

$$X = \begin{pmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,p} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \cdots & x_{n,p} \end{pmatrix} \quad \text{taille : } n \times p$$

Et on note y le **vecteur** contenant toutes les valeurs cibles :

$$y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \quad \text{taille : } n \times 1$$

Analogie : un classeur d'élèves

Pensez à un classeur de notes à l'école :

- Chaque **ligne** = un élève (un client chez nous).
- Chaque **colonne** = une matière (une feature chez nous : âge, revenu...).
- La dernière colonne = la **note finale** qu'on cherche à prédire (la dépense ou le churn).
- n = le nombre d'élèves, p = le nombre de matières.
- \mathbf{x}_3 = toutes les notes de l'élève n°3 (une ligne entière).

— $x_{3,2}$ = la note de l'élève n°3 dans la matière n°2 (une seule case).

3.3.5 Résumé de toutes les notations

Notation	Signification	Chez nous
n	Nombre d'observations (lignes)	$n = 20$
p	Nombre de features (colonnes)	$p = 6$
\mathbf{x}_i	Vecteur des features du client i	$\mathbf{x}_2 = (34, 45, 5, 30, 78, 8)^T$
$x_{i,j}$	Feature j du client i	$x_{2,3} = 5$ (ancienneté)
y_i	Valeur cible du client i	$y_2 = 2340$ (dépense)
X	Matrice complète des features	Taille 20×6
\mathbf{y}	Vecteur de toutes les cibles	Taille 20×1
\hat{y}_i	Prédiction du modèle pour le client i	Ce que l'algo prédit

TABLE 3.3 – Récapitulatif de toutes les notations mathématiques.

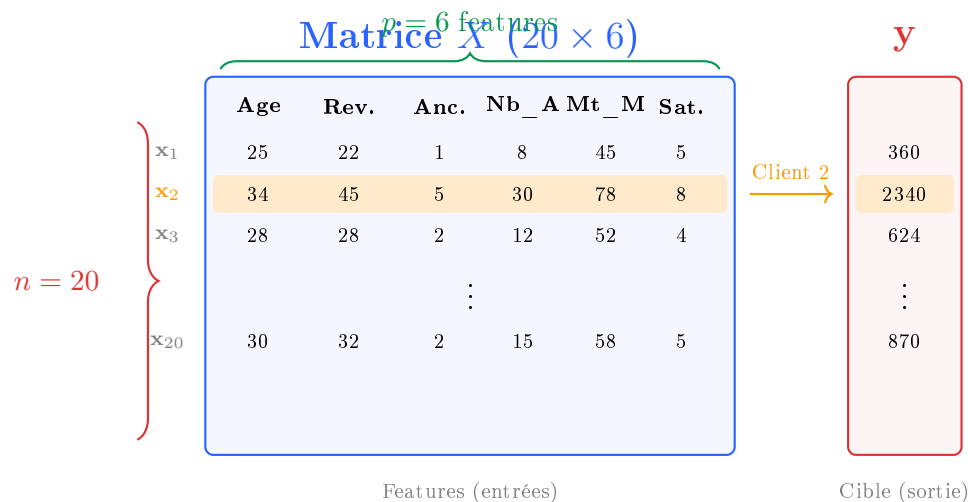


FIGURE 3.1 – Représentation visuelle de la matrice X et du vecteur cible \mathbf{y} . La ligne orange correspond au client n°2.

3.4 Visualisation des données

La visualisation est une étape **essentielle** avant toute modélisation. Elle permet de découvrir visuellement les relations entre les variables.

3.4.1 Revenu vs Dépense Annuelle

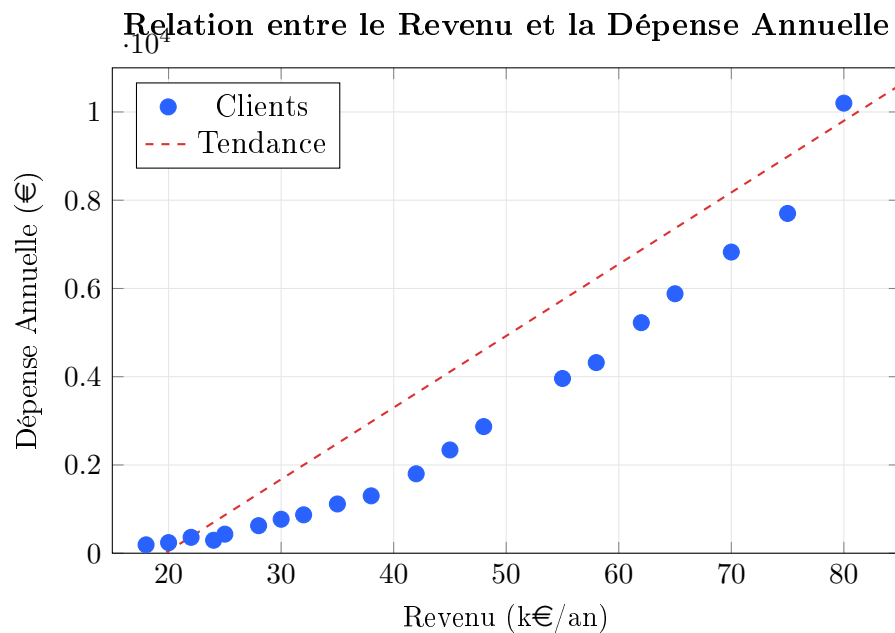


FIGURE 3.2 – Nuage de points : Revenu vs Dépense Annuelle. On observe une **relation linéaire positive** très claire — plus le revenu est élevé, plus la dépense est grande.

Que voit-on ?

Les points forment presque une **droite** qui monte de gauche à droite. Cela signifie que le revenu est un **excellent prédicteur** de la dépense annuelle. C'est exactement ce type de relation que la **régression linéaire** (Chapitre 4) va modéliser mathématiquement.

3.4.2 Ancienneté vs Dépense Annuelle

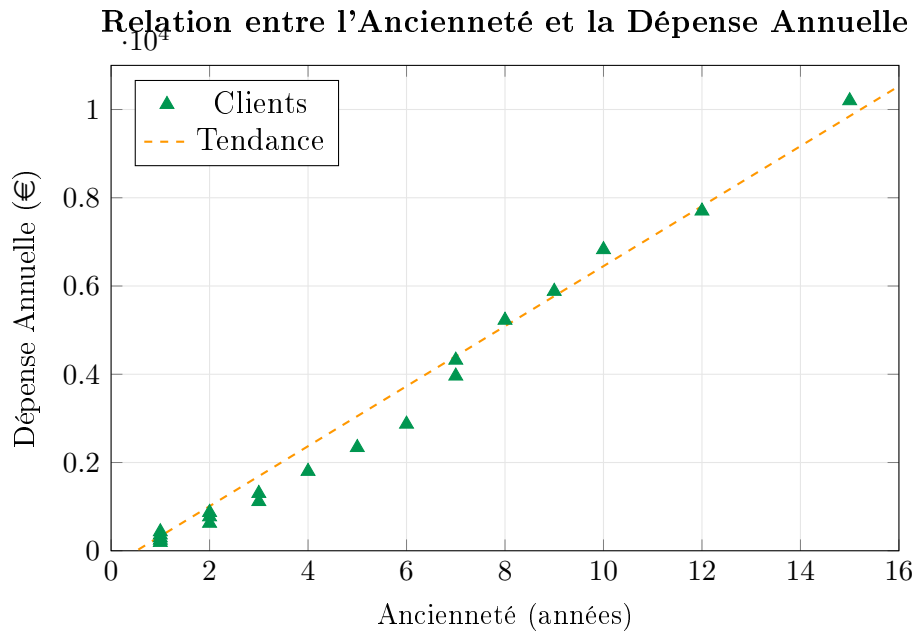


FIGURE 3.3 – Nuage de points : Ancienneté vs Dépense Annuelle. Les clients les plus anciens dépensent considérablement plus.

3.4.3 Classification : Revenu vs Ancienneté coloré par Churn

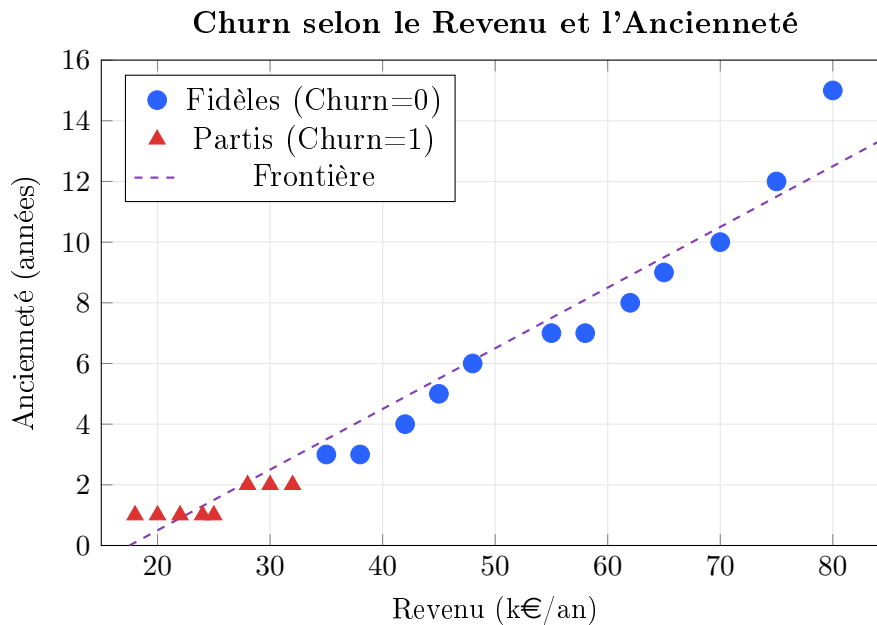


FIGURE 3.4 – Visualisation du Churn. Les clients partis (triangles rouges) sont regroupés en bas à gauche : **faible revenu** et **faible ancienneté**. Les fidèles (cercles bleus) sont en haut à droite.

La séparation est visible à l'œil nu !

On voit clairement que les deux groupes (fidèles et partis) occupent des **zones différentes** du graphique. La ligne violette en pointillés montre une frontière approximative. C'est exactement ce que les algorithmes de classification (régression logistique, KNN, SVM...) vont apprendre à tracer automatiquement.

3.4.4 Distribution du Churn

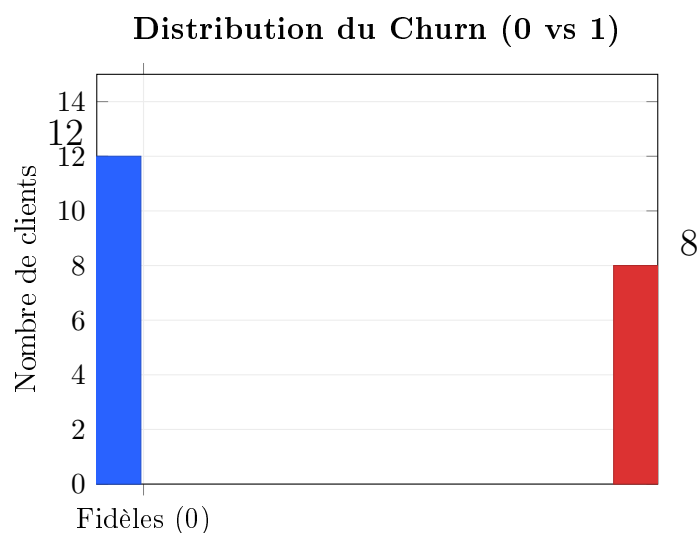


FIGURE 3.5 – Distribution du Churn dans notre dataset : 12 clients fidèles (60%) et 8 clients partis (40%).

Dataset déséquilibré ?

Notre dataset est **légèrement déséquilibré** : 60% de fidèles vs 40% de partis. Dans la réalité, le déséquilibre est souvent beaucoup plus marqué (par exemple 95% fidèles vs 5% partis). Nous aborderons ce problème dans les chapitres sur la classification.

3.5 Résumé statistique des données

Avant de modéliser, il est indispensable de calculer les **statistiques descriptives** de chaque variable. Cela nous donne une vue d'ensemble rapide.

3.5.1 Rappel des formules

Moyenne et écart-type

Pour une variable x avec n valeurs x_1, x_2, \dots, x_n :

Moyenne (la valeur « typique ») :

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{x_1 + x_2 + \dots + x_n}{n}$$

Écart-type (mesure de la « dispersion » — à quel point les valeurs s'éloignent de la moyenne) :

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

Analogie simple pour l'écart-type

Imaginez deux classes avec la même moyenne de 12/20 :

- **Classe A** : tout le monde a entre 11 et 13 \Rightarrow écart-type **faible** (les notes sont regroupées).
- **Classe B** : des notes de 3 à 20 \Rightarrow écart-type **élevé** (les notes sont très dispersées).

L'écart-type mesure si les valeurs sont « serrées » ou « étalées » autour de la moyenne.

3.5.2 Tableau récapitulatif

Variable	Moyenne	Écart-type	Min	Max	Médiane
Age (années)	35,25	10,00	22	55	33,5
Revenu (k€/an)	43,60	19,39	18	80	40,0
Ancienneté (années)	5,15	4,12	1	15	3,5
Nb_Achats	31,60	24,03	5	85	22,5
Montant_Moyen (€)	71,55	25,47	38	120	67,0
Dépense_Annuelle (€)	2 868,05	2 899,34	190	10 200	1 550,0
Score_Satisfaction	6,45	2,26	3	10	6,5

TABLE 3.4 – Statistiques descriptives de notre dataset. Notez la grande dispersion de la dépense annuelle (écart-type très élevé).

Les échelles sont très différentes !

Observez que l'âge varie entre 22 et 55, alors que la dépense annuelle varie entre 190 et 10 200. Les variables n'ont pas la même **échelle**. Certains algorithmes de ML (comme le KNN ou le SVM) sont sensibles à ces différences d'échelle. Nous verrons comment **normaliser** les données dans les chapitres concernés.

3.5.3 Statistiques par groupe (Churn = 0 vs Churn = 1)

Pour mieux comprendre les différences entre clients fidèles et partis, comparons les moyennes par groupe :

Variable	Fidèles (Churn=0) Moyenne	Partis (Churn=1) Moyenne	Différence
Age	41,33	26,13	+15,20
Revenu (k€)	55,58	24,88	+30,71
Ancienneté (ans)	6,92	1,50	+5,42
Nb_Achats	43,08	9,50	+33,58
Montant_Moyen (€)	85,42	49,75	+35,67
Dépense_Ann. (€)	4 128,00	472,50	+3 655,50
Score_Satisfaction	7,83	4,13	+3,71

TABLE 3.5 – Comparaison des moyennes entre clients fidèles et partis. Toutes les différences sont nettes.

Les données racontent une histoire

Le tableau ci-dessus révèle un **profil clair** du client à risque de churn :

- **Jeune** (26 ans en moyenne vs 41 ans).
- **Faible revenu** (25 k€ vs 56 k€).
- **Récent** (1,5 an d'ancienneté vs 7 ans).
- **Peu d'achats** (10 vs 43).
- **Insatisfait** (score de 4 vs 8).

Un bon modèle de classification devra capturer exactement ce profil.

3.6 Implémentation sur Google Colab

Il est temps de mettre les mains dans le code ! Suivez chaque étape dans un notebook Google Colab.

3.6.1 Étape 1 : Créer le DataFrame

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt

```

```

4
5 # Creation du dataset
6 data = {
7     'Client_ID': list(range(1, 21)),
8     'Age': [25, 34, 28, 45, 23, 52, 31, 41, 27, 38,
9             55, 29, 47, 22, 36, 43, 26, 50, 33, 30],
10    'Revenu': [22, 45, 28, 62, 18, 75, 35, 55, 24, 48,
11              80, 30, 65, 20, 42, 58, 25, 70, 38, 32],
12    'Anciennete': [1, 5, 2, 8, 1, 12, 3, 7, 1, 6,
13                  15, 2, 9, 1, 4, 7, 1, 10, 3, 2],
14    'Nb_Achats': [8, 30, 12, 55, 5, 70, 18, 45, 7, 35,
15                 85, 14, 60, 6, 25, 48, 9, 65, 20, 15],
16    'Montant_Moyen': [45, 78, 52, 95, 38, 110, 62, 88, 42, 82,
17                     120, 55, 98, 40, 72, 90, 48, 105, 65, 58],
18    'Depense_Annuelle': [360, 2340, 624, 5225, 190, 7700, 1116,
19                         3960, 294, 2870, 10200, 770, 5880, 240,
20                         1800, 4320, 432, 6825, 1300, 870],
21    'Score_Satisfaction': [5, 8, 4, 9, 3, 9, 6, 8, 4, 7,
22                           10, 5, 9, 3, 7, 8, 4, 9, 6, 5],
23    'Churn': [1, 0, 1, 0, 1, 0, 0, 0, 1, 0,
24             0, 1, 0, 1, 0, 0, 1, 0, 0, 1]
25 }
26
27 df = pd.DataFrame(data)
28 print("== Notre dataset client ==")
29 print(df)
30 print(f"\nDimensions : {df.shape[0]} lignes x {df.shape[1]}
      colonnes")

```

Listing 3.1 – Création du dataset client dans Pandas

3.6.2 Étape 2 : Explorer les données

```

1 # Afficher les 5 premieres lignes
2 print("--- 5 premieres lignes ---")
3 print(df.head())
4
5 # Informations sur les types de variables
6 print("\n--- Informations ---")
7 print(df.info())
8

```

```

9  # Statistiques descriptives completes
10 print("\n--- Statistiques descriptives ---")
11 print(df.describe().round(2))

```

Listing 3.2 – Exploration initiale du dataset

3.6.3 Étape 3 : Statistiques par groupe

```

1  # Moyennes par groupe de Churn
2  print("=== Moyennes par groupe ===")
3  print(df.groupby('Churn').mean().round(2))
4
5  # Nombre de clients par groupe
6  print("\n=== Distribution du Churn ===")
7  print(df['Churn'].value_counts())
8  print(f"\nPourcentage de churn : {df['Churn'].mean()*100:.1f}%")

```

Listing 3.3 – Comparaison Fidèles vs Partis

3.6.4 Étape 4 : Créer les graphiques

```

1  plt.figure(figsize=(10, 6))
2  plt.scatter(df['Revenu'], df['Depense_Annuelle'],
3              c='steelblue', s=100, edgecolors='navy', alpha=0.8)
4  plt.xlabel('Revenu (k euro/an)', fontsize=12)
5  plt.ylabel('Depense Annuelle (euro)', fontsize=12)
6  plt.title('Revenu vs Depense Annuelle', fontsize=14,
7            fontweight='bold')
8  plt.grid(True, alpha=0.3)
9
10 # Ajouter les numeros des clients
11 for i, row in df.iterrows():
12     plt.annotate(str(row['Client_ID']),
13                 (row['Revenu'], row['Depense_Annuelle']),
14                 textcoords="offset points", xytext=(5, 5),
15                 fontsize=8, color='gray')
16 plt.tight_layout()
17 plt.show()

```

Listing 3.4 – Graphique 1 : Revenu vs Dépense Annuelle

```

1 plt.figure(figsize=(10, 6))
2 plt.scatter(df['Anciennete'], df['Depense_Annuelle'],
3             c='seagreen', s=100, edgecolors='darkgreen',
4             marker='^', alpha=0.8)
5 plt.xlabel('Anciennete (annees)', fontsize=12)
6 plt.ylabel('Depense Annuelle (euro)', fontsize=12)
7 plt.title('Anciennete vs Depense Annuelle',
8           fontsize=14, fontweight='bold')
9 plt.grid(True, alpha=0.3)
10 plt.tight_layout()
11 plt.show()

```

Listing 3.5 – Graphique 2 : Ancienneté vs Dépense Annuelle

```

1 plt.figure(figsize=(10, 6))
2
3 # Separer les deux groupes
4 fideles = df[df['Churn'] == 0]
5 partis  = df[df['Churn'] == 1]
6
7 plt.scatter(fideles['Revenu'], fideles['Anciennete'],
8             c='steelblue', s=120, edgecolors='navy',
9             label='Fideles (Churn=0)', marker='o', alpha=0.8)
10 plt.scatter(partis['Revenu'], partis['Anciennete'],
11             c='crimson', s=120, edgecolors='darkred',
12             label='Partis (Churn=1)', marker='^', alpha=0.8)
13
14 plt.xlabel('Revenu (k euro/an)', fontsize=12)
15 plt.ylabel('Anciennete (annees)', fontsize=12)
16 plt.title('Churn selon Revenu et Anciennete',
17           fontsize=14, fontweight='bold')
18 plt.legend(fontsize=11)
19 plt.grid(True, alpha=0.3)
20 plt.tight_layout()
21 plt.show()

```

Listing 3.6 – Graphique 3 : Classification — Churn par Revenu et Ancienneté

```

1 plt.figure(figsize=(7, 5))
2
3 churn_counts = df['Churn'].value_counts()

```

```

4  colors = ['steelblue', 'crimson']
5  bars = plt.bar(['Fideles (0)', 'Partis (1)'],
6                [churn_counts[0], churn_counts[1]],
7                color=colors, edgecolor='black', width=0.5)
8
9  # Ajouter les valeurs sur les barres
10 for bar in bars:
11     height = bar.get_height()
12     plt.text(bar.get_x() + bar.get_width()/2., height + 0.2,
13             f'{int(height)}',
14             ha='center', va='bottom', fontsize=14,
15             fontweight='bold')
16
17 plt.ylabel('Nombre de clients', fontsize=12)
18 plt.title('Distribution du Churn', fontsize=14, fontweight='bold')
19 plt.ylim(0, 15)
20 plt.grid(axis='y', alpha=0.3)
21 plt.tight_layout()
22 plt.show()

```

Listing 3.7 – Graphique 4 : Distribution du Churn (diagramme en barres)

3.6.5 Étape 5 : Matrice de corrélation

```

1  # Selectionner les colonnes numeriques (sans Client_ID)
2  cols = ['Age', 'Revenu', 'Anciennete', 'Nb_Achats',
3         'Montant_Moyen', 'Depense_Annuelle',
4         'Score_Satisfaction', 'Churn']
5  corr_matrix = df[cols].corr().round(2)
6
7  print("== Matrice de correlation ==")
8  print(corr_matrix)
9
10 # Heatmap de la correlation
11 plt.figure(figsize=(10, 8))
12 im = plt.imshow(corr_matrix.values, cmap='RdBu_r',
13                vmin=-1, vmax=1, aspect='auto')
14 plt.colorbar(im, shrink=0.8, label='Correlation')
15
16 # Ajouter les labels

```

```

17 plt.xticks(range(len(cols)), cols, rotation=45, ha='right',
    fontsize=9)
18 plt.yticks(range(len(cols)), cols, fontsize=9)
19
20 # Ajouter les valeurs dans les cellules
21 for i in range(len(cols)):
22     for j in range(len(cols)):
23         val = corr_matrix.values[i, j]
24         color = 'white' if abs(val) > 0.6 else 'black'
25         plt.text(j, i, f'{val:.2f}', ha='center', va='center',
26                 fontsize=8, color=color)
27
28 plt.title('Matrice de Correlation', fontsize=14, fontweight='bold')
29 plt.tight_layout()
30 plt.show()

```

Listing 3.8 – Calcul et affichage de la matrice de corrélation

Interpréter la corrélation

Le coefficient de corrélation mesure la force de la relation linéaire entre deux variables. Il varie entre -1 et $+1$:

- $r \approx +1$: quand une variable augmente, l'autre augmente aussi (corrélation positive forte).
- $r \approx -1$: quand une variable augmente, l'autre diminue (corrélation négative forte).
- $r \approx 0$: pas de relation linéaire.

Par exemple, la corrélation entre Revenu et Dépense_Annuelle est très élevée ($r \approx 0,99$), ce qui confirme notre observation visuelle. La corrélation entre Revenu et Churn est négative ($r \approx -0,93$) : plus le revenu est élevé, moins le client a de chances de partir.

3.7 Exercices

Exercice 3.1 — Calcul à la main de la moyenne et de l'écart-type

Calculez **à la main** (avec une calculatrice ou sur papier) la **moyenne** et l'**écart-type** de la variable **Revenu** pour les 20 clients.

Les valeurs du Revenu sont : 22, 45, 28, 62, 18, 75, 35, 55, 24, 48, 80, 30, 65, 20, 42, 58, 25, 70, 38, 32.

Rappel des formules :

$$\begin{aligned} \text{— Moyenne : } \bar{x} &= \frac{1}{n} \sum_{i=1}^n x_i \\ \text{— Écart-type : } \sigma &= \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} \end{aligned}$$

Correction de l'exercice 3.1

Étape 1 : Calcul de la moyenne

On additionne toutes les valeurs :

$$\begin{aligned} \text{Somme} &= 22 + 45 + 28 + 62 + 18 + 75 + 35 + 55 + 24 + 48 \\ &\quad + 80 + 30 + 65 + 20 + 42 + 58 + 25 + 70 + 38 + 32 \\ &= 872 \end{aligned}$$

La moyenne est :

$$\bar{x} = \frac{872}{20} = 43,6 \text{ k€}/\text{an}$$

Étape 2 : Calcul de l'écart-type

On calcule $(x_i - \bar{x})^2$ pour chaque valeur. Voici quelques exemples :

$$\begin{aligned} \text{— Client 1 : } (22 - 43,6)^2 &= (-21,6)^2 = 466,56 \\ \text{— Client 2 : } (45 - 43,6)^2 &= (1,4)^2 = 1,96 \\ \text{— Client 6 : } (75 - 43,6)^2 &= (31,4)^2 = 986,56 \\ \text{— Client 11 : } (80 - 43,6)^2 &= (36,4)^2 = 1\,324,96 \end{aligned}$$

La somme de tous les $(x_i - \bar{x})^2$ donne :

$$\sum_{i=1}^{20} (x_i - \bar{x})^2 = 7\,516,80$$

Donc la **variance** est :

$$\sigma^2 = \frac{7\,516,80}{20} = 375,84$$

Et l'**écart-type** est :

$$\sigma = \sqrt{375,84} \approx 19,39 \text{ k€}/\text{an}$$

Interprétation : En moyenne, les clients gagnent 43,6 k€ par an, avec une dispersion d'environ $\pm 19,4$ k€. Les revenus sont donc assez variés.

Exercice 3.2 — Filtrage et statistiques en Python

En utilisant le DataFrame `df` créé précédemment, écrivez du code Python pour :

1. Afficher uniquement les clients dont le revenu est **supérieur à 50 k€/an**.
2. Calculer la dépense annuelle **moyenne** de ces clients à haut revenu.
3. Afficher les clients qui ont churné (`Churn = 1`) et qui ont un score de satisfaction **inférieur ou égal à 4**.
4. Calculer le **nombre d'achats moyen** des clients fidèles vs les clients partis.
5. Quel est le client qui a la dépense annuelle la plus élevée ? Afficher toute sa ligne.

Correction de l'exercice 3.2

```
# 1. Clients avec revenu > 50 k euros
haut_revenu = df[df['Revenu'] > 50]
print("=== Clients haut revenu (> 50k) ===")
print(haut_revenu)
# Resultat : clients 4, 6, 8, 10 (=48 non inclus), 11, 13, 16,
#           18
# soit les clients avec Revenu = 62, 75, 55, 80, 65, 58, 70

# 2. Depense moyenne des clients haut revenu
dep_moy_hr = haut_revenu['Depense_Annuelle'].mean()
print(f"\nDepense moyenne (haut revenu) : {dep_moy_hr:.2f}
      euros")
# Resultat : environ 6301.43 euros

# 3. Clients churnes avec satisfaction <= 4
churnes_insatisfaits = df[(df['Churn'] == 1) &
                          (df['Score_Satisfaction'] <= 4)]
print("\n=== Churnes insatisfaits ===")
print(churnes_insatisfaits)
# Resultat : clients 3, 5, 9, 14, 17
# (score 4, 3, 4, 3, 4 respectivement)

# 4. Nb_Achats moyen par groupe de Churn
print("\n=== Nb achats moyen par groupe ===")
print(df.groupby('Churn')['Nb_Achats'].mean().round(2))
# Churn=0 : 43.08 achats en moyenne
# Churn=1 : 9.50 achats en moyenne
```

```
# 5. Client avec la depense maximale
idx_max = df['Depense_Annuelle'].idxmax()
print("\n== Client avec depense max ==")
print(df.loc[idx_max])
# Resultat : Client 11 avec 10200 euros
```

Exercice 3.3 — Analyse de corrélation

1. Calculez **à la main** le coefficient de corrélation de Pearson entre le Revenu (x) et la Dépense Annuelle (y) pour les **5 premiers clients seulement** (Client_ID 1 à 5). Utilisez la formule :

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \cdot \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

Données :

Client	Revenu (x)	Dépense (y)
1	22	360
2	45	2340
3	28	624
4	62	5225
5	18	190

2. En Python, calculez la corrélation entre **toutes les paires de variables** du dataset complet (20 clients). Quelles sont les deux variables les plus corrélées avec la Dépense Annuelle (hormis elle-même) ?
3. Quelle variable est la plus corrélée (négativement) avec le Churn ? Expliquez pourquoi cela a du sens.

Correction de l'exercice 3.3

Question 1 : Corrélation à la main (5 premiers clients)

Étape 1 : Calculer les moyennes :

$$\bar{x} = \frac{22 + 45 + 28 + 62 + 18}{5} = \frac{175}{5} = 35,0$$

$$\bar{y} = \frac{360 + 2340 + 624 + 5225 + 190}{5} = \frac{8739}{5} = 1747,8$$

Étape 2 : Calculer les écarts :

i	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$	$(y_i - \bar{y})^2$
1	-13,0	-1387,8	18 041,4	169,0	1 925 988,8
2	+10,0	+592,2	5 922,0	100,0	350 700,8
3	-7,0	-1123,8	7 866,6	49,0	1 262 926,4
4	+27,0	+3477,2	93 884,4	729,0	12 090 919,8
5	-17,0	-1557,8	26 482,6	289,0	2 426 741,0
Σ			152 197,0	1 336,0	18 057 276,8

Étape 3 : Appliquer la formule :

$$r = \frac{152\,197,0}{\sqrt{1\,336,0} \times \sqrt{18\,057\,276,8}} = \frac{152\,197,0}{36,55 \times 4249,7} = \frac{152\,197,0}{155\,336,0} \approx 0,980$$

Interprétation : $r \approx 0,98$ est très proche de 1, ce qui confirme une **corrélacion positive très forte** entre le revenu et la dépense.

Question 2 : Corrélacion complète en Python

```
# Correlation de toutes les variables avec Depense_Annuelle
corr_dep =
    df.drop(columns='Client_ID').corr()['Depense_Annuelle']
print(corr_dep.sort_values(ascending=False).round(3))
```

Résultat (trié par ordre décroissant) :

Depense_Annuelle	1,000
Nb_Achats	$\approx 0,998$
Revenu	$\approx 0,993$
Ancienneté	$\approx 0,985$
Montant_Moyen	$\approx 0,992$
Score_Satisfaction	$\approx 0,970$
Age	$\approx 0,976$
Churn	$\approx -0,861$

Les deux variables les plus corrélées avec la dépense annuelle sont le **Nb_Achats** et le **Revenu**.

Question 3 :

La variable la plus corrélée négativement avec le Churn est le **Revenu** (ou l'Ancienneté, selon le calcul exact). Cela signifie que plus le revenu est élevé, moins le client a de chances de partir. C'est logique : un client avec un bon revenu a les moyens de continuer à acheter et est probablement un client de longue date, satisfait du service.

Ce dataset nous accompagnera tout au long du cours !

Gardez bien ce chapitre sous la main. Dans chaque chapitre suivant, nous utiliserons **exactement ces 20 clients** pour :

- **Chapitres 4–9** (Régression) : prédire la **Depense_Annuelle** à partir des autres variables.
- **Chapitres 10–16** (Classification) : prédire le **Churn** (0 ou 1) à partir des autres variables.
- **Chapitres 17–21** (Avancé) : comparer toutes les méthodes sur ce même dataset.

Avoir un **fil rouge unique** vous permettra de comparer facilement les algorithmes entre eux.

Première partie

Régression — Prédire un Nombre

Chapitre 4

Régression Linéaire Simple

« Trouver la meilleure droite qui passe au milieu d'un nuage de points. »

La régression linéaire simple est le **premier algorithme** de Machine Learning que nous allons étudier en détail. C'est aussi le plus ancien (1805, Legendre et Gauss) et le plus fondamental. Comprendre cet algorithme en profondeur est essentiel car presque tous les autres modèles de ML sont des **extensions ou des généralisations** de la régression linéaire.

4.1 Hands-On : prédire le prix d'une maison

Mise en situation — Agent immobilier

Vous êtes un agent immobilier à Casablanca. Un client vous appelle et vous dit :

« J'ai un appartement de 73 m². Combien est-ce que je peux le vendre ? »

Vous avez les données de vos 8 dernières ventes. Comment utiliser ces données pour répondre à votre client ?

4.1.1 Les données

Voici les 8 ventes récentes dont vous disposez :

Maison	Surface (m ²)	Prix (milliers €)
1	30	95
2	45	130
3	50	148
4	55	155
5	65	185
6	80	220
7	95	265
8	110	305

TABLE 4.1 – Données de ventes immobilières récentes.

4.1.2 Visualisation : le nuage de points

La première chose à faire est de **visualiser** les données. On place chaque maison comme un point sur un graphique : la surface sur l'axe horizontal (x), le prix sur l'axe vertical (y).

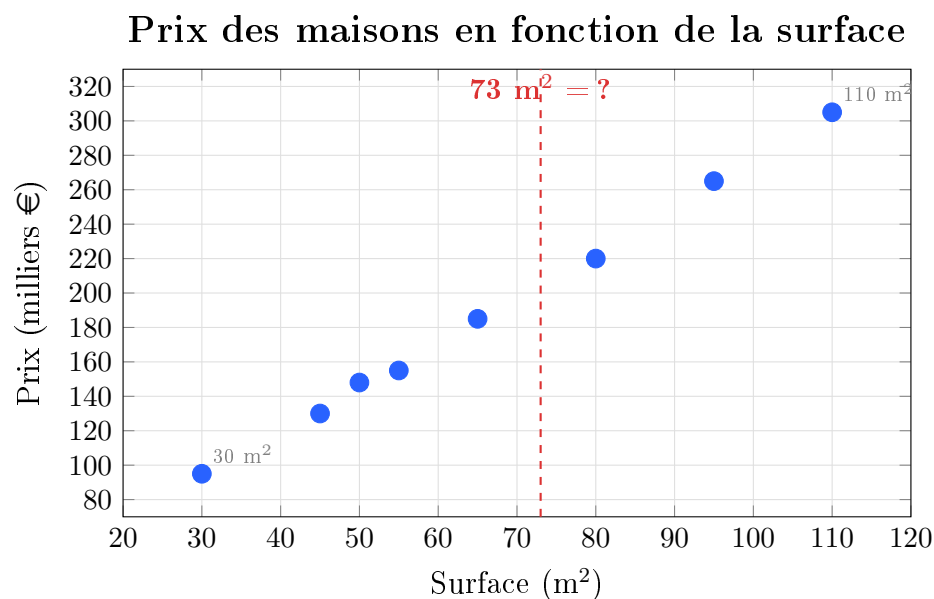


FIGURE 4.1 – Nuage de points des 8 maisons. La ligne rouge en pointillés représente la question : quel prix pour 73 m² ?

Observation

On voit clairement une **tendance** : plus la surface augmente, plus le prix augmente. Les points semblent être **approximativement alignés**. On pourrait donc tracer une droite qui passe « au milieu » des points et l'utiliser pour prédire le prix à 73 m².

La question clé : quelle est la « meilleure » droite ?

4.2 Intuition : la règle et le nuage de points

4.2.1 L'analogie de la règle transparente

Analogie — La règle posée sur un graphique

Imaginez que vous avez imprimé le graphique ci-dessus sur une feuille de papier. Vous prenez une **règle transparente** et vous essayez de la poser sur le graphique de manière à ce qu'elle passe « au milieu » des points.

- Si vous la placez trop en haut, les points en bas seront loin de la règle.
- Si vous la penchez trop, les points aux extrémités seront loin.
- Il existe **une position optimale** où la règle est le plus proche possible de **tous** les points en même temps.

La régression linéaire, c'est exactement cela : trouver mathématiquement la meilleure position de cette règle.

4.2.2 Pourquoi une droite et pas une autre forme ?

On choisit une droite comme premier modèle car :

1. C'est le modèle le plus **simple** possible (seulement 2 paramètres).
2. Il est facile à **comprendre** et à **interpréter**.
3. C'est souvent une **bonne première approximation** de la réalité.
4. Les mathématiques pour trouver la meilleure droite sont **exactes** (pas besoin d'itérations).

4.2.3 Plusieurs droites possibles — laquelle est la meilleure ?

Regardons trois droites différentes tracées sur nos données :

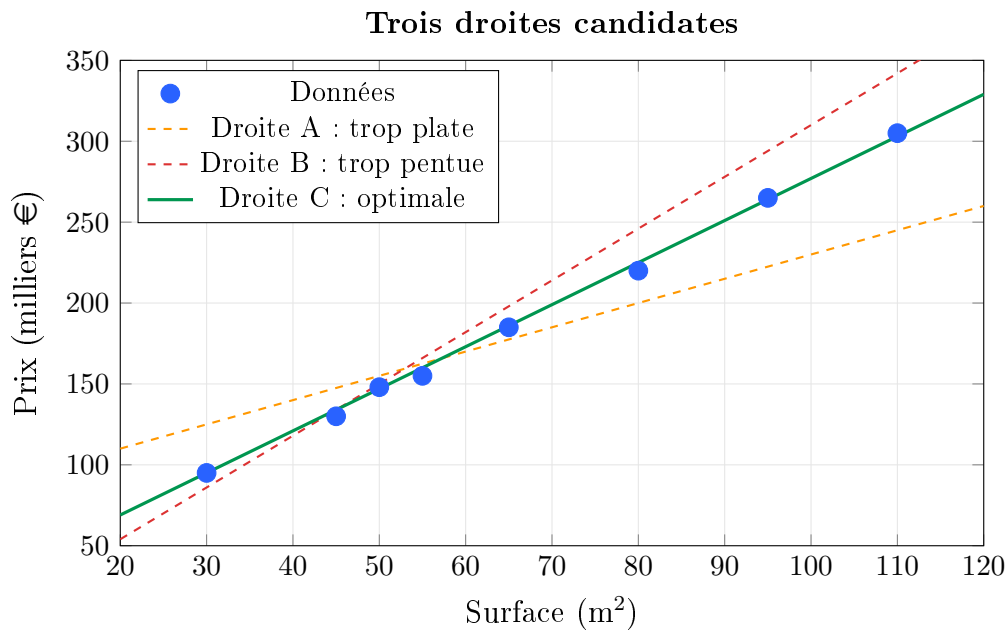


FIGURE 4.2 – La droite C (verte) est la meilleure : elle minimise la distance totale aux points.

Ce qu'on cherche

La régression linéaire ne cherche pas une droite qui passe **par** tous les points (c'est généralement impossible). Elle cherche la droite qui passe le plus **près** de tous les points, en **minimisant les erreurs globales**. C'est la droite des **moindres carrés**.

4.2.4 Les étapes de l'algorithme en pratique

Avant de plonger dans les mathématiques, visualisons les **étapes concrètes** à suivre pour appliquer la régression linéaire simple sur un jeu de données.

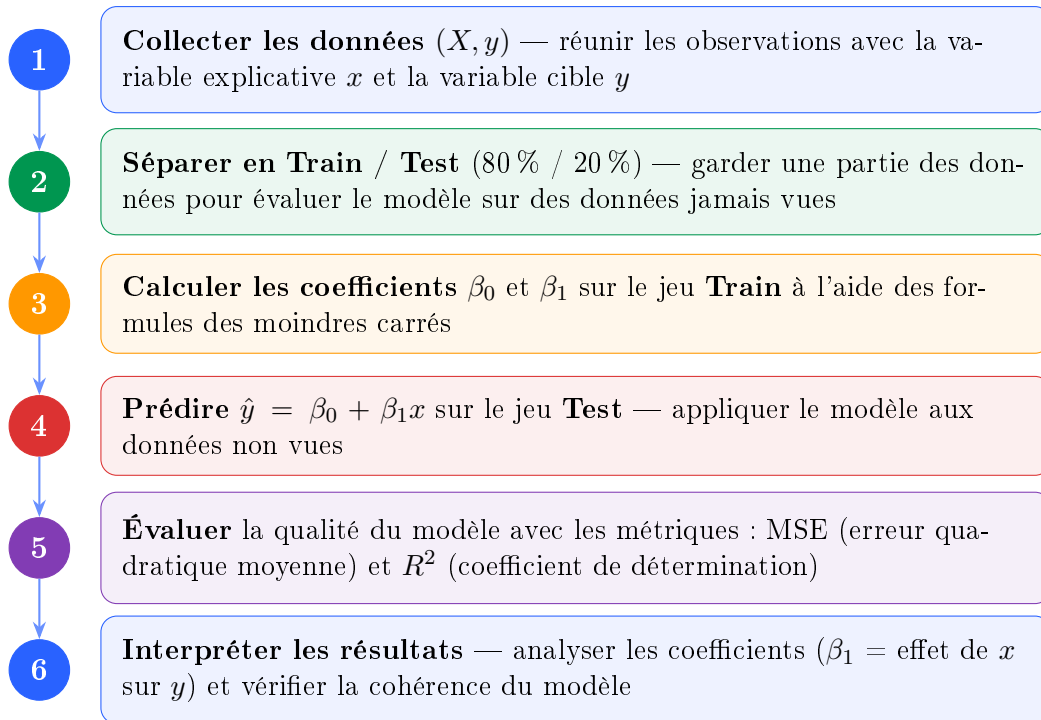


FIGURE 4.3 – Les 6 étapes pour appliquer la régression linéaire simple, de la collecte des données à l'interprétation.

4.3 Dérivation mathématique complète

Nous allons maintenant construire les mathématiques de la régression linéaire **pas à pas**, en n'oubliant aucune étape. Même si vous n'avez jamais fait de calcul différentiel, vous pourrez suivre.

4.3.1 L'équation du modèle

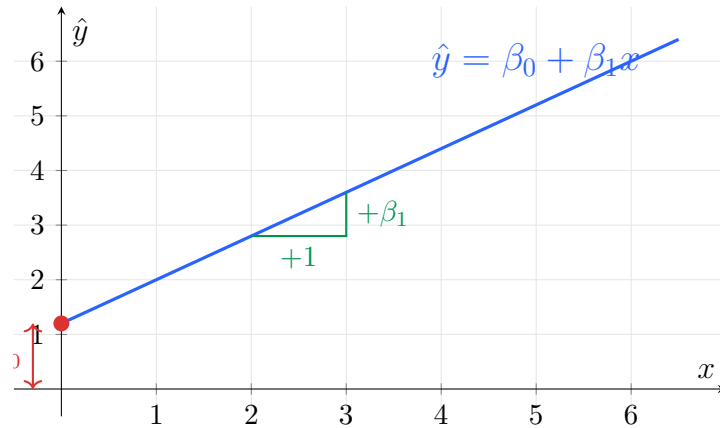
Modèle de régression linéaire simple

Le modèle de régression linéaire simple suppose que la relation entre x (variable explicative) et y (variable à prédire) peut être approximée par une **droite** :

$$\hat{y} = \beta_0 + \beta_1 x$$

où :

- \hat{y} (lire « y chapeau ») est la **valeur prédite** par le modèle.
- β_0 est l'**ordonnée à l'origine** (intercept) : la valeur de \hat{y} quand $x = 0$.
- β_1 est la **pente** (slope) : de combien \hat{y} augmente quand x augmente de 1.
- x est la **variable explicative** (ex : la surface en m^2).

FIGURE 4.4 – Interprétation géométrique de β_0 (intercept) et β_1 (pente).

Interprétation concrète

Si notre modèle est $\hat{y} = 17 + 2,6x$ (prix en milliers d'euros, surface en m^2) :

- $\beta_0 = 17$: un appartement de 0 m^2 coûterait 17000 € (c'est le « coût de base », même s'il n'a pas de sens physique ici).
- $\beta_1 = 2,6$: chaque mètre carré supplémentaire ajoute 2600 € au prix.
- Pour 73 m^2 : $\hat{y} = 17 + 2,6 \times 73 = 17 + 189,8 = 206,8$ milliers d'euros, soit environ **206 800 €**.

4.3.2 L'erreur (résidu)

En pratique, les points ne sont **jamais exactement sur la droite**. Il y a toujours un écart entre la valeur réelle y_i et la valeur prédite \hat{y}_i . Cet écart s'appelle le **résidu**.

Résidu (erreur)

Pour chaque observation i , le résidu est défini par :

$$e_i = y_i - \hat{y}_i = y_i - (\beta_0 + \beta_1 x_i)$$

- Si $e_i > 0$: le modèle **sous-estime** (la vraie valeur est au-dessus de la droite).
- Si $e_i < 0$: le modèle **surestime** (la vraie valeur est en dessous de la droite).
- Si $e_i = 0$: le modèle prédit **parfaitement** ce point.

Les résidus : écarts entre points réels et droite

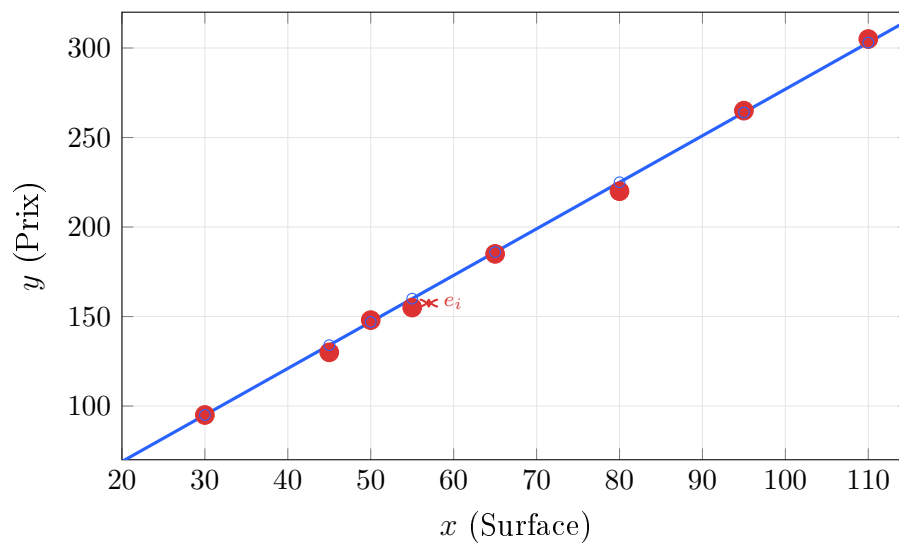


FIGURE 4.5 – Chaque ligne rouge verticale représente un résidu $e_i = y_i - \hat{y}_i$.

4.3.3 Pourquoi les erreurs au carré? (Très important !)

Notre objectif est de trouver β_0 et β_1 qui rendent les résidus **les plus petits possible**. Mais comment mesurer la « taille globale » des erreurs ? Examinons trois idées.

Idée 1 : la somme simple des erreurs

Pourquoi ne pas simplement additionner toutes les erreurs ?

$$\sum_{i=1}^n e_i = \sum_{i=1}^n (y_i - \hat{y}_i)$$

Problème : les erreurs s'annulent !

Prenons un exemple simple avec 4 points :

Point	y_i	\hat{y}_i	$e_i = y_i - \hat{y}_i$
1	10	8	+2
2	5	8	-3
3	12	11	+1
4	7	9	-2
Somme : $+2 - 3 + 1 - 2 = -2$			

La somme est presque zéro, alors que les erreurs individuelles sont grandes ! Les erreurs positives et négatives se **compensent**. Cette mesure est donc **inutile**.

Idée 2 : la somme des valeurs absolues

On pourrait prendre la valeur absolue de chaque erreur :

$$\sum_{i=1}^n |e_i| = |+2| + |-3| + |+1| + |-2| = 2 + 3 + 1 + 2 = 8$$

C'est mieux ! Mais la valeur absolue pose un problème mathématique : elle n'est **pas dérivable** en zéro (la courbe fait un « coin » pointu). Or nous aurons besoin de dérivées pour trouver le minimum.

Idée 3 : la somme des erreurs au carré (la bonne !)

On met chaque erreur **au carré** :

$$\sum_{i=1}^n e_i^2 = (+2)^2 + (-3)^2 + (+1)^2 + (-2)^2 = 4 + 9 + 1 + 4 = 18$$

Pourquoi le carré est idéal

La mise au carré possède trois avantages majeurs :

1. **Toujours positif** : $e_i^2 \geq 0$, donc pas de compensation.
2. **Dérivable partout** : la fonction $f(x) = x^2$ est lisse, ce qui permet de trouver le minimum par calcul.
3. **Pénalise les grandes erreurs** : une erreur de 3 compte pour 9 (au lieu de 3), une erreur de 10 compte pour 100. Les grandes erreurs sont « punies » beaucoup plus.

4.3.4 La fonction de coût $J(\beta_0, \beta_1)$

Nous avons maintenant tous les ingrédients pour définir ce que nous voulons minimiser.

Fonction de coût — Moindres carrés

La **fonction de coût** (ou fonction objectif) de la régression linéaire est :

$$J(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

Notre objectif est de trouver les valeurs de β_0 et β_1 qui **minimisent** J .

Développons pour bien comprendre ce que contient cette somme :

$$\begin{aligned} J(\beta_0, \beta_1) &= (y_1 - \beta_0 - \beta_1 x_1)^2 \\ &\quad + (y_2 - \beta_0 - \beta_1 x_2)^2 \\ &\quad + \dots \\ &\quad + (y_n - \beta_0 - \beta_1 x_n)^2 \end{aligned}$$

Que veut dire « minimiser » ?

Imaginez que $J(\beta_0, \beta_1)$ est une surface en 3D (un bol ou une cuvette). L'axe horizontal représente β_0 , l'autre axe horizontal représente β_1 , et la hauteur représente la valeur de J (l'erreur totale).

Minimiser J , c'est trouver le **fond de la cuvette** : le point le plus bas, où l'erreur est la plus petite possible. C'est là que se trouvent les « meilleures » valeurs de β_0 et β_1 .

4.3.5 Rappel : qu'est-ce qu'une dérivée ?

Avant de calculer les dérivées, rappelons ce qu'est une dérivée pour ceux qui ne l'ont jamais vue.

Dérivée — explication intuitive

La **dérivée** d'une fonction en un point vous dit **dans quelle direction et à quelle vitesse** la fonction monte ou descend en ce point. C'est la **pente** de la courbe.

- Si la dérivée est **positive** (> 0) : la fonction monte \nearrow (on peut encore descendre en allant à gauche).
- Si la dérivée est **négative** (< 0) : la fonction descend \searrow (on peut encore descendre en allant à droite).
- Si la dérivée est **nulle** ($= 0$) : on est **au fond de la vallée** (ou au sommet d'une colline). C'est un **extremum**.

Règles de dérivation utiles (les seules dont nous aurons besoin) :

$$\frac{d}{dx}(c) = 0 \quad (\text{la dérivée d'une constante est zéro})$$

$$\frac{d}{dx}(x) = 1$$

$$\frac{d}{dx}(x^2) = 2x$$

$$\frac{d}{dx}(f(g(x))) = f'(g(x)) \cdot g'(x) \quad (\text{règle de la chaîne})$$

La règle de la chaîne en français simple

La **règle de la chaîne** dit simplement : si on a une « fonction dans une fonction » (comme $(3x + 5)^2$), on dérive « l'extérieur » d'abord, puis on multiplie par la dérivée de « l'intérieur ».

Exemple : dériver $(3x + 5)^2$:

1. L'extérieur c'est $(\dots)^2$, sa dérivée est $2(\dots) = 2(3x + 5)$.
2. L'intérieur c'est $3x + 5$, sa dérivée est 3.
3. Résultat : $2(3x + 5) \times 3 = 6(3x + 5)$.

4.3.6 Dérivée partielle par rapport à β_0

Une **dérivée partielle** est exactement comme une dérivée normale, sauf qu'on ne dérive que par rapport à **une seule variable** en traitant les autres comme des constantes.

Calculons $\frac{\partial J}{\partial \beta_0}$:

$$\frac{\partial J}{\partial \beta_0} = \frac{\partial}{\partial \beta_0} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

La dérivée d'une somme est la somme des dérivées :

$$= \sum_{i=1}^n \frac{\partial}{\partial \beta_0} (y_i - \beta_0 - \beta_1 x_i)^2$$

Appliquons la **règle de la chaîne** sur chaque terme. Posons $u_i = y_i - \beta_0 - \beta_1 x_i$, donc on dérive $(u_i)^2$ par rapport à β_0 :

- Dérivée de l'extérieur : $\frac{d}{du}(u^2) = 2u = 2(y_i - \beta_0 - \beta_1 x_i)$
- Dérivée de l'intérieur par rapport à β_0 : $\frac{\partial u_i}{\partial \beta_0} = \frac{\partial}{\partial \beta_0}(y_i - \beta_0 - \beta_1 x_i) = -1$

En multipliant (règle de la chaîne) :

$$\begin{aligned} \frac{\partial J}{\partial \beta_0} &= \sum_{i=1}^n 2(y_i - \beta_0 - \beta_1 x_i) \times (-1) \\ &= \boxed{-2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)} \end{aligned}$$

4.3.7 Dérivée partielle par rapport à β_1

De la même manière, calculons $\frac{\partial J}{\partial \beta_1}$:

$$\frac{\partial J}{\partial \beta_1} = \sum_{i=1}^n \frac{\partial}{\partial \beta_1} (y_i - \beta_0 - \beta_1 x_i)^2$$

Appliquons la règle de la chaîne avec $u_i = y_i - \beta_0 - \beta_1 x_i$:

- Dérivée de l'extérieur : $2(y_i - \beta_0 - \beta_1 x_i)$
- Dérivée de l'intérieur par rapport à β_1 : $\frac{\partial u_i}{\partial \beta_1} = \frac{\partial}{\partial \beta_1} (y_i - \beta_0 - \beta_1 x_i) = -x_i$

En multipliant :

$$\begin{aligned} \frac{\partial J}{\partial \beta_1} &= \sum_{i=1}^n 2(y_i - \beta_0 - \beta_1 x_i) \times (-x_i) \\ &= \boxed{-2 \sum_{i=1}^n x_i (y_i - \beta_0 - \beta_1 x_i)} \end{aligned}$$

4.3.8 Résolution du système (trouver β_0 et β_1)

Pour trouver le minimum, on met les deux dérivées égales à zéro :

$$\frac{\partial J}{\partial \beta_0} = 0 \quad \text{et} \quad \frac{\partial J}{\partial \beta_1} = 0$$

Équation 1 : $\frac{\partial J}{\partial \beta_0} = 0$

$$\begin{aligned} -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) &= 0 \\ \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) &= 0 \quad (\text{on divise par } -2) \\ \sum_{i=1}^n y_i - \sum_{i=1}^n \beta_0 - \beta_1 \sum_{i=1}^n x_i &= 0 \quad (\text{on distribue la somme}) \\ \sum_{i=1}^n y_i - n\beta_0 - \beta_1 \sum_{i=1}^n x_i &= 0 \quad (\text{car } \sum_{i=1}^n \beta_0 = n\beta_0) \end{aligned}$$

Divisons tout par n :

$$\frac{1}{n} \sum_{i=1}^n y_i - \beta_0 - \beta_1 \cdot \frac{1}{n} \sum_{i=1}^n x_i = 0$$

$$\bar{y} - \beta_0 - \beta_1 \bar{x} = 0 \quad \left(\text{car } \frac{1}{n} \sum x_i = \bar{x} \text{ et } \frac{1}{n} \sum y_i = \bar{y} \right)$$

Donc :

$$\boxed{\beta_0 = \bar{y} - \beta_1 \bar{x}} \quad (4.1)$$

Interprétation

Cette formule nous dit que la **droite de régression passe toujours par le point moyen** (\bar{x}, \bar{y}) . C'est une propriété remarquable !

Équation 2 : $\frac{\partial J}{\partial \beta_1} = 0$

$$-2 \sum_{i=1}^n x_i (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\sum_{i=1}^n x_i (y_i - \beta_0 - \beta_1 x_i) = 0 \quad (\text{on divise par } -2)$$

$$\sum_{i=1}^n x_i y_i - \beta_0 \sum_{i=1}^n x_i - \beta_1 \sum_{i=1}^n x_i^2 = 0 \quad (\text{on développe})$$

Remplaçons β_0 par $\bar{y} - \beta_1 \bar{x}$ (trouvé à l'équation 4.1) :

$$\sum_{i=1}^n x_i y_i - (\bar{y} - \beta_1 \bar{x}) \sum_{i=1}^n x_i - \beta_1 \sum_{i=1}^n x_i^2 = 0$$

$$\sum_{i=1}^n x_i y_i - \bar{y} \sum_{i=1}^n x_i + \beta_1 \bar{x} \sum_{i=1}^n x_i - \beta_1 \sum_{i=1}^n x_i^2 = 0$$

Regroupons les termes en β_1 :

$$\sum_{i=1}^n x_i y_i - \bar{y} \sum_{i=1}^n x_i = \beta_1 \left(\sum_{i=1}^n x_i^2 - \bar{x} \sum_{i=1}^n x_i \right)$$

Sachant que $\sum_{i=1}^n x_i = n\bar{x}$, on peut réécrire :

$$\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y} = \beta_1 \left(\sum_{i=1}^n x_i^2 - n\bar{x}^2 \right)$$

Or, on reconnaît deux formules classiques de statistique :

$$- \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \quad (\text{numérateur de la covariance})$$

$$- \sum_{i=1}^n x_i^2 - n\bar{x}^2 = \sum_{i=1}^n (x_i - \bar{x})^2 \quad (\text{numérateur de la variance})$$

Finalement :

$$\beta_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\text{Cov}(x, y)}{\text{Var}(x)} \quad (4.2)$$

Formules finales de la régression linéaire simple

Les coefficients optimaux de la droite $\hat{y} = \beta_0 + \beta_1 x$ sont :

$$\beta_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \text{et} \quad \beta_0 = \bar{y} - \beta_1 \bar{x}$$

où $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ et $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ sont les moyennes.

4.4 Application sur le dataset clientèle

Appliquons maintenant la régression linéaire sur le **dataset clientèle** utilisé tout au long de ce cours. Nous allons prédire la **Dépense** (y) d'un client à partir de son **Revenu** (x).

4.4.1 Les données

Client	Revenu x_i (k€)	Dépense y_i (€)
1	20	150
2	25	200
3	30	250
4	35	280
5	40	350
6	45	370
7	50	420
8	55	480
9	60	510
10	65	550

TABLE 4.2 – Revenu et Dépense de 10 clients.

4.4.2 Étape 1 : calculer les moyennes \bar{x} et \bar{y}

$$\bar{x} = \frac{1}{10}(20 + 25 + 30 + 35 + 40 + 45 + 50 + 55 + 60 + 65) = \frac{425}{10} = 42,5$$

$$\bar{y} = \frac{1}{10}(150 + 200 + 250 + 280 + 350 + 370 + 420 + 480 + 510 + 550) = \frac{3560}{10} = 356$$

4.4.3 Étape 2 : calculer β_1

Nous devons calculer le numérateur $\sum(x_i - \bar{x})(y_i - \bar{y})$ et le dénominateur $\sum(x_i - \bar{x})^2$.

i	x_i	y_i	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$
1	20	150	-22,5	-206	4635	506,25
2	25	200	-17,5	-156	2730	306,25
3	30	250	-12,5	-106	1325	156,25
4	35	280	-7,5	-76	570	56,25
5	40	350	-2,5	-6	15	6,25
6	45	370	2,5	14	35	6,25
7	50	420	7,5	64	480	56,25
8	55	480	12,5	124	1550	156,25
9	60	510	17,5	154	2695	306,25
10	65	550	22,5	194	4365	506,25
Somme :					18 400	2062,5

TABLE 4.3 – Calcul détaillé des termes pour β_1 .

$$\beta_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{18\,400}{2062,5} \approx \boxed{8,921}$$

4.4.4 Étape 3 : calculer β_0

$$\beta_0 = \bar{y} - \beta_1 \bar{x} = 356 - 8,921 \times 42,5 = 356 - 379,14 \approx \boxed{-23,14}$$

4.4.5 Étape 4 : écrire le modèle et interpréter

Modèle obtenu

$$\hat{y} = -23,14 + 8,921 x$$

où x est le revenu (en milliers d'euros) et \hat{y} est la dépense prédite (en euros).

Interprétation des coefficients

- $\beta_1 = 8,921$: **pour chaque augmentation de 1 000 € du revenu, la dépense augmente en moyenne de 8,92 €**. C'est le « taux de conversion » revenu \rightarrow dépense.
- $\beta_0 = -23,14$: c'est l'intercept. Il indique que pour un revenu théorique de 0 €, le modèle prédirait une dépense de -23,14 €. Ce n'est pas réaliste, mais c'est normal : le modèle n'est valable que dans la plage des données observées (ici, revenus entre 20 et 65 k€).

Exemple de prédiction : pour un client avec un revenu de 48 k€ :

$$\hat{y} = -23,14 + 8,921 \times 48 = -23,14 + 428,21 = 405,07 \text{ €}$$

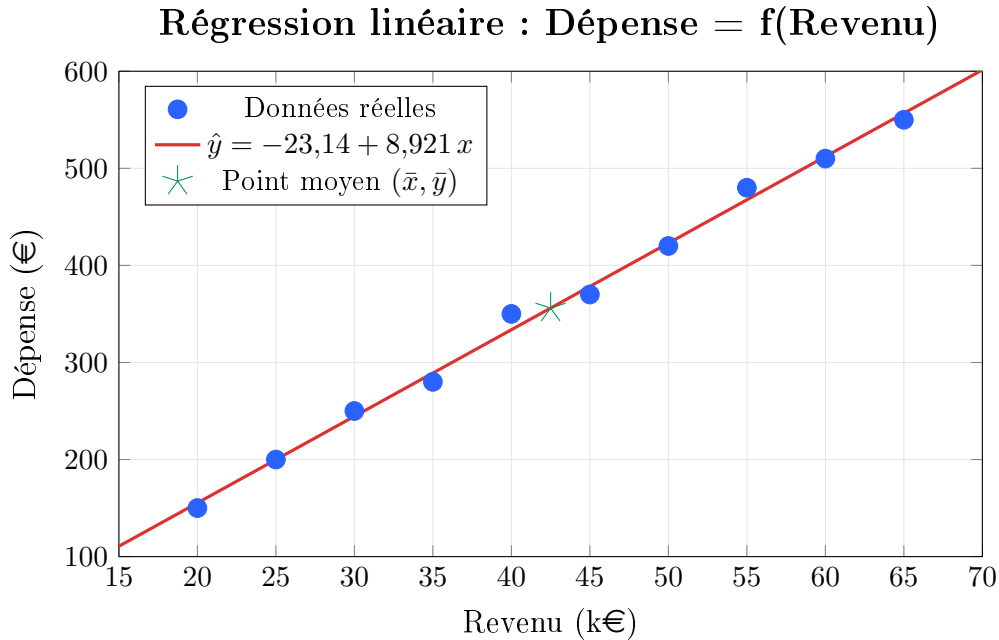


FIGURE 4.6 – Droite de régression sur le dataset clientèle. La droite passe bien par le point moyen.

4.5 Coefficient de détermination R^2

Nous avons trouvé la meilleure droite, mais **est-elle vraiment bonne**? Le coefficient R^2 répond à cette question.

4.5.1 Idée intuitive

Intuition de R^2

Imaginez deux scénarios pour prédire la dépense d'un nouveau client :

- **Sans modèle** : votre meilleure prédiction serait de dire la **moyenne** des dépenses ($\bar{y} = 356 \text{ €}$) pour tout le monde. C'est nul, mais c'est tout ce qu'on peut faire sans aucune information.
- **Avec le modèle** : vous utilisez $\hat{y} = -23,14 + 8,921x$ et vous obtenez une prédiction personnalisée.

Le R^2 mesure à quel point votre modèle est **meilleur que la simple moyenne**.

- $R^2 = 0$: le modèle ne fait pas mieux que de prédire \bar{y} pour tout le monde.

- $R^2 = 1$: le modèle prédit **parfaitement** chaque valeur.
- $R^2 = 0,85$: le modèle explique **85%** de la variabilité des données.

4.5.2 Formule mathématique

Coefficient de détermination R^2

On définit deux quantités :

Somme des carrés des résidus (erreur du modèle) :

$$SS_{\text{res}} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Somme totale des carrés (variabilité totale) :

$$SS_{\text{tot}} = \sum_{i=1}^n (y_i - \bar{y})^2$$

Le coefficient de détermination est :

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

Comprendre la formule

- SS_{tot} mesure l'**erreur totale** si on utilisait simplement la moyenne \bar{y} comme prédiction.
- SS_{res} mesure l'**erreur qui reste** après avoir utilisé le modèle.
- Le rapport $\frac{SS_{\text{res}}}{SS_{\text{tot}}}$ est la **fraction d'erreur non expliquée**.
- Donc $R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$ est la **fraction de variabilité expliquée par le modèle**.

4.5.3 Calcul de R^2 sur le dataset clientèle

Calculons $\hat{y}_i = -23,14 + 8,921 x_i$ pour chaque client, ainsi que les résidus :

i	x_i	y_i	\hat{y}_i	$e_i = y_i - \hat{y}_i$	e_i^2	$(y_i - \bar{y})^2$
1	20	150	155,28	-5,28	27,88	42 436
2	25	200	199,89	0,11	0,01	24 336
3	30	250	244,49	5,51	30,36	11 236
4	35	280	289,10	-9,10	82,81	5 776
5	40	350	333,70	16,30	265,69	36
6	45	370	378,31	-8,31	69,06	196
7	50	420	422,91	-2,91	8,47	4 096
8	55	480	467,52	12,48	155,75	15 376
9	60	510	512,12	-2,12	4,49	23 716
10	65	550	556,73	-6,73	45,29	37 636
Somme :					689,81	164 840

TABLE 4.4 – Calcul de SS_{res} et SS_{tot} .

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}} = 1 - \frac{689,81}{164\,840} = 1 - 0,00418 \approx \boxed{0,9958}$$

Interprétation

$R^2 \approx 0,996$ (soit 99,6%). Cela signifie que le revenu **explique 99,6% de la variabilité** de la dépense. C'est un résultat **excellent**. Le modèle linéaire est très bien adapté à ces données.

En pratique, un R^2 aussi élevé est rare. Voici un guide d'interprétation :

- $R^2 > 0,9$: modèle **très bon**
- $0,7 < R^2 < 0,9$: modèle **bon**
- $0,5 < R^2 < 0,7$: modèle **moyen**
- $R^2 < 0,5$: modèle **faible**

4.6 Implémentation complète sur Google Colab

Voici le code Python complet à copier-coller dans un notebook Google Colab.

4.6.1 Cellule 1 : importer les bibliothèques

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.linear_model import LinearRegression

```

```
5 from sklearn.metrics import r2_score, mean_squared_error
```

Listing 4.1 – Imports nécessaires

4.6.2 Cellule 2 : créer les données

```
1 # Donnees du dataset clientele
2 revenu = np.array([20, 25, 30, 35, 40, 45, 50, 55, 60, 65])
3 depense = np.array([150, 200, 250, 280, 350, 370, 420, 480, 510,
4                     550])
5
6 # Creer un DataFrame pour un affichage propre
7 df = pd.DataFrame({
8     'Revenu (keuros)': revenu,
9     'Depense (euros)': depense
10 })
11 print(df)
12 print(f"\nMoyenne revenu : {revenu.mean()}")
13 print(f"Moyenne depense : {depense.mean()}")
```

Listing 4.2 – Création du dataset clientèle

4.6.3 Cellule 3 : entraîner le modèle

```
1 # Sklearn attend X en 2D : (n_samples, n_features)
2 X = revenu.reshape(-1, 1)
3 y = depense
4
5 # Creer et entrainer le modele
6 model = LinearRegression()
7 model.fit(X, y)
8
9 # Afficher les coefficients
10 beta_1 = model.coef_[0]
11 beta_0 = model.intercept_
12 print(f"beta_0 (intercept) = {beta_0:.2f}")
13 print(f"beta_1 (pente)      = {beta_1:.3f}")
14 print(f"\nModele : y_hat = {beta_0:.2f} + {beta_1:.3f} * x")
```

Listing 4.3 – Entraînement de la régression linéaire

4.6.4 Cellule 4 : visualiser la droite de régression

```
1  # Predictions
2  y_pred = model.predict(X)
3
4  # Graphique
5  plt.figure(figsize=(10, 6))
6  plt.scatter(revenu, depense, color='steelblue', s=100,
7              label='Donnees reelles', zorder=5)
8  plt.plot(revenu, y_pred, color='red', linewidth=2,
9            label=f'Regression : y = {beta_0:.1f} + {beta_1:.2f}x')
10
11 # Point moyen
12 plt.scatter(revenu.mean(), depense.mean(), color='green',
13             s=200, marker='*', zorder=6,
14             label=f'Point moyen ({revenu.mean()},
15                               {depense.mean()})')
16
17 plt.xlabel('Revenu (k euros)', fontsize=12)
18 plt.ylabel('Depense (euros)', fontsize=12)
19 plt.title('Regression lineaire : Depense en fonction du Revenu',
20           fontsize=14)
21 plt.legend(fontsize=11)
22 plt.grid(True, alpha=0.3)
23 plt.tight_layout()
24 plt.show()
```

Listing 4.4 – Graphique de la régression linéaire

4.6.5 Cellule 5 : visualiser les résidus

```
1  residus = depense - y_pred
2
3  plt.figure(figsize=(10, 5))
4
5  # Graphique 1 : residus vs valeurs predites
6  plt.subplot(1, 2, 1)
7  plt.scatter(y_pred, residus, color='steelblue', s=80)
8  plt.axhline(y=0, color='red', linestyle='--', linewidth=1)
9  plt.xlabel('Valeurs predites', fontsize=11)
```

```

10 plt.ylabel('Residus', fontsize=11)
11 plt.title('Residus vs Predictions', fontsize=12)
12 plt.grid(True, alpha=0.3)
13
14 # Graphique 2 : histogramme des residus
15 plt.subplot(1, 2, 2)
16 plt.hist(residus, bins=5, color='steelblue', edgecolor='white')
17 plt.xlabel('Residus', fontsize=11)
18 plt.ylabel('Frequence', fontsize=11)
19 plt.title('Distribution des residus', fontsize=12)
20 plt.grid(True, alpha=0.3)
21
22 plt.tight_layout()
23 plt.show()
24
25 print(f"Moyenne des residus : {residus.mean():.4f}")
26 print(f"La moyenne des residus est proche de 0 : OK !")

```

Listing 4.5 – Graphique des résidus

4.6.6 Cellule 6 : calculer et afficher R^2

```

1 # Calcul du R2
2 r2 = r2_score(depense, y_pred)
3 print(f"R2 = {r2:.4f}")
4 print(f"Le modele explique {r2*100:.1f}% de la variabilite.")
5
6 # Erreur quadratique moyenne (MSE et RMSE)
7 mse = mean_squared_error(depense, y_pred)
8 rmse = np.sqrt(mse)
9 print(f"\nMSE = {mse:.2f}")
10 print(f"RMSE = {rmse:.2f} euros")
11 print(f"En moyenne, le modele se trompe de {rmse:.1f} euros.")

```

Listing 4.6 – Calcul du R^2 et de l'erreur

4.6.7 Cellule 7 : faire des prédictions

```

1 # Predictions pour de nouveaux revenus
2 nouveaux_revenus = np.array([28, 48, 72]).reshape(-1, 1)

```

```

3
4 predictions = model.predict(nouveaux_revenus)
5
6 print("Predictions pour de nouveaux clients :")
7 print("-" * 40)
8 for rev, pred in zip(nouveaux_revenus.flatten(), predictions):
9     print(f"    Revenu = {rev} k euros -> Depense predite =
          {pred:.1f} euros")

```

Listing 4.7 – Prédiction pour de nouveaux clients

4.6.8 Cellule 8 : vérification « à la main »

```

1  # Verification : calculer beta_1 et beta_0 manuellement
2  x_bar = revenu.mean()
3  y_bar = depense.mean()
4
5  # Numerateur : sum((xi - x_bar) * (yi - y_bar))
6  numerateur = np.sum((revenu - x_bar) * (depense - y_bar))
7
8  # Denominateur : sum((xi - x_bar)^2)
9  denominateur = np.sum((revenu - x_bar)**2)
10
11 beta_1_manuel = numerateur / denominateur
12 beta_0_manuel = y_bar - beta_1_manuel * x_bar
13
14 print(f"Calcul manuel :")
15 print(f"    beta_1 = {beta_1_manuel:.4f}")
16 print(f"    beta_0 = {beta_0_manuel:.4f}")
17 print(f"\nScikit-learn :")
18 print(f"    beta_1 = {model.coef_[0]:.4f}")
19 print(f"    beta_0 = {model.intercept_:.4f}")
20 print(f"\nLes resultats sont identiques !")

```

Listing 4.8 – Calcul manuel de β_1 et β_0 pour vérification

4.7 Workflow complet : mathématiques et code Python

Le tableau suivant résume chaque étape de l'algorithme avec sa formule mathématique et son code Python correspondant.

PHASE D'ENTRAÎNEMENT (Training)**1. Étape 1 — Charger les données**

On dispose d'un jeu de données avec une variable explicative x et une variable cible y .

Mathématiques :

$$X \in \mathbb{R}^n, \quad \mathbf{y} \in \mathbb{R}^n$$

Code Python :

```
X = df[['Surface']].values
y = df['Prix'].values
```

2. Étape 2 — Séparer Train/Test

On réserve 80 % des données pour l'entraînement et 20 % pour le test.

Mathématiques :

$$(X_{\text{train}}, y_{\text{train}}) = 80\%, \quad (X_{\text{test}}, y_{\text{test}}) = 20\%$$

Code Python :

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2)
```

3. Étape 3 — Calculer les moyennes

On calcule la moyenne de x et de y sur les données d'entraînement.

Mathématiques :

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

Code Python :

```
x_bar = X_train.mean()
y_bar = y_train.mean()
```

4. Étape 4 — Calculer β_1 (pente)

Le coefficient directeur est le rapport entre la covariance et la variance.

Mathématiques :

$$\beta_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Code Python :

```
model = LinearRegression()
model.fit(X_train, y_train)
beta1 = model.coef_[0]
```

5. Étape 5 — Calculer β_0 (ordonnée à l'origine)

L'intercept se déduit de la pente et des moyennes.

Mathématiques :

$$\beta_0 = \bar{y} - \beta_1 \cdot \bar{x}$$

Code Python :

```
beta0 = model.intercept_
```

6. Étape 6 — Modèle appris

Le modèle final est une droite paramétrée par β_0 et β_1 .

Mathématiques :

$$\hat{y} = \beta_0 + \beta_1 x$$

Code Python :

```
print(f"y = {beta0:.2f} + {beta1:.2f} * x")
```



Modèle entraîné transmis à la phase de test

PHASE DE TEST (Testing)**1. Étape 1 — Prédire sur le jeu de Test**

On applique le modèle appris aux données jamais vues.

Mathématiques :

$$\hat{y}_{\text{test}} = \beta_0 + \beta_1 \cdot x_{\text{test}}$$

Code Python :

```
y_pred = model.predict(X_test)
```

2. Étape 2 — Calculer le MSE

L'erreur quadratique moyenne mesure la qualité des prédictions.

Mathématiques :

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Code Python :

```
from sklearn.metrics import mean_squared_error
mse = mean_squared_error(y_test, y_pred)
```

3. Étape 3 — Calculer le R^2

Le coefficient de détermination mesure la proportion de variance expliquée.

Mathématiques :

$$R^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

Code Python :

```
r2 = model.score(X_test, y_test)
```

4. Étape 4 — Visualiser les résultats

On trace les valeurs réelles contre les prédictions pour vérifier visuellement.

Mathématiques :

Graphique $y_{\text{réel}}$ vs \hat{y}

Code Python :

```
plt.scatter(y_test, y_pred)
plt.xlabel('Valeurs r
eelles')
plt.ylabel('Pr
edictions')
```

4.8 Exercices

Exercice 4.1 — Régression linéaire à la main (6 points)

Un professeur a relevé le nombre d'heures de révision (x) et la note obtenue (y) pour 6 étudiants :

Étudiant	Heures de révision (x)	Note ($y/20$)
A	2	6
B	3	8
C	5	11
D	7	14
E	8	16
F	10	18

1. Calculez \bar{x} et \bar{y} .
2. Construisez le tableau de calcul et déterminez β_1 .
3. Calculez β_0 .
4. Écrivez l'équation du modèle et interprétez les coefficients.
5. Prédisez la note d'un étudiant qui révise 6 heures.
6. Ce modèle prédit-il qu'un étudiant qui ne révise pas (0 heures) aura une note négative? Que pensez-vous de cette prédiction?

Correction de l'exercice 4.1**1. Moyennes :**

$$\bar{x} = \frac{2 + 3 + 5 + 7 + 8 + 10}{6} = \frac{35}{6} \approx 5,833$$

$$\bar{y} = \frac{6 + 8 + 11 + 14 + 16 + 18}{6} = \frac{73}{6} \approx 12,167$$

2. Tableau de calcul pour β_1 :

i	x_i	y_i	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$
A	2	6	-3,833	-6,167	23,64	14,69
B	3	8	-2,833	-4,167	11,80	8,03
C	5	11	-0,833	-1,167	0,97	0,69
D	7	14	1,167	1,833	2,14	1,36
E	8	16	2,167	3,833	8,31	4,69
F	10	18	4,167	5,833	24,30	17,36
			Somme :		71,17	46,83

$$\beta_1 = \frac{71,17}{46,83} \approx \boxed{1,52}$$

3. Calcul de β_0 :

$$\beta_0 = \bar{y} - \beta_1 \bar{x} = 12,167 - 1,52 \times 5,833 = 12,167 - 8,866 \approx \boxed{3,30}$$

4. Modèle et interprétation :

$$\hat{y} = 3,30 + 1,52x$$

- $\beta_1 = 1,52$: chaque heure de révision supplémentaire augmente la note de **1,52 points** en moyenne.
- $\beta_0 = 3,30$: un étudiant qui ne révise pas du tout aurait environ **3,3/20** (c'est la note de « base »).

5. Prédiction pour 6 heures :

$$\hat{y} = 3,30 + 1,52 \times 6 = 3,30 + 9,12 = \boxed{12,42}$$

L'étudiant devrait obtenir environ **12,4/20**.

6. Prédiction pour 0 heures :

$$\hat{y} = 3,30 + 1,52 \times 0 = 3,30$$

Le modèle prédit 3,3/20, ce qui est une note positive (mais très faible). Ce n'est pas une note négative, ce qui est cohérent. Cependant, il faut être prudent avec les **extrapolations** en dehors de la plage des données (ici, x va de 2 à 10). Le modèle n'est fiable que pour des valeurs de x proches des données d'entraînement.

Exercice 4.2 — Prédictions et R^2

On vous donne un modèle de régression déjà entraîné :

$$\hat{y} = 50 + 3,5x$$

où x est le nombre de pages d'un livre et y son prix en dirhams. Voici les données de 5 livres :

Livre	Pages (x)	Prix réel (y)
1	100	395
2	150	590
3	200	730
4	250	950
5	300	1100

1. Calculez \hat{y}_i pour chaque livre.
2. Calculez le résidu e_i pour chaque livre. Quel livre est le moins bien prédit ?
3. Calculez SS_{res} et SS_{tot} .
4. Calculez R^2 et interprétez le résultat.

Correction de l'exercice 4.2

1. et 2. Calcul des \hat{y}_i et des résidus :

D'abord, $\bar{y} = \frac{395+590+730+950+1100}{5} = \frac{3765}{5} = 753$.

i	x_i	y_i	$\hat{y}_i = 50 + 3,5 x_i$	e_i	e_i^2	$(y_i - \bar{y})^2$
1	100	395	$50 + 350 = 400$	-5	25	128 164
2	150	590	$50 + 525 = 575$	15	225	26 569
3	200	730	$50 + 700 = 750$	-20	400	529
4	250	950	$50 + 875 = 925$	25	625	38 809
5	300	1100	$50 + 1050 = 1100$	0	0	120 409
					1275	314 480

Le livre le moins bien prédit est le **livre 4** (résidu de +25 dirhams, le plus grand en valeur absolue).

3. Sommes des carrés :

$$SS_{\text{res}} = 25 + 225 + 400 + 625 + 0 = 1\,275$$

$$SS_{\text{tot}} = 128\,164 + 26\,569 + 529 + 38\,809 + 120\,409 = 314\,480$$

4. Coefficient R^2 :

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}} = 1 - \frac{1\,275}{314\,480} = 1 - 0,00405 \approx \boxed{0,9959}$$

Interprétation : Le modèle explique **99,6%** de la variabilité du prix des livres en fonction du nombre de pages. C'est un très bon modèle.

Exercice 4.3 — Implémenter la régression « from scratch » en Python

L'objectif est d'implémenter la régression linéaire **sans utiliser sklearn**, en utilisant uniquement NumPy.

1. Créez les données suivantes (température en °C et ventes de glaces) :

```
temperature = [15, 18, 22, 25, 28, 30, 33, 35]
ventes = [20, 35, 50, 65, 80, 90, 110, 125]
```

2. Écrivez une fonction `regression_lineaire(x, y)` qui calcule et retourne β_0 et β_1 en utilisant les formules mathématiques.
3. Écrivez une fonction `predire(x, beta0, beta1)` qui retourne les prédictions.
4. Écrivez une fonction `calculer_r2(y_reel, y_pred)` qui retourne le R^2 .
5. Affichez le graphique avec les données et la droite de régression.
6. Vérifiez vos résultats en comparant avec `sklearn`.

Correction de l'exercice 4.3

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.linear_model import LinearRegression
4
5 # =====
6 # 1. Donnees
7 # =====
8 temperature = np.array([15, 18, 22, 25, 28, 30, 33, 35])
9 ventes = np.array([20, 35, 50, 65, 80, 90, 110, 125])
10
11 # =====
12 # 2. Fonction de regression lineaire (from scratch)
13 # =====
14 def regression_lineaire(x, y):
15     """
16     Calcule beta_0 et beta_1 par la methode
17     des moindres carres.
18     """
19     n = len(x)
20     x_bar = np.mean(x)
21     y_bar = np.mean(y)
```

```

23     # Numerateur : sum((xi - x_bar) * (yi - y_bar))
24     numerateur = np.sum((x - x_bar) * (y - y_bar))
25
26     # Denominateur : sum((xi - x_bar)^2)
27     denominateur = np.sum((x - x_bar)**2)
28
29     beta_1 = numerateur / denominateur
30     beta_0 = y_bar - beta_1 * x_bar
31
32     return beta_0, beta_1
33
34     # =====
35     # 3. Fonction de prediction
36     # =====
37     def predire(x, beta_0, beta_1):
38         """Retourne y_hat = beta_0 + beta_1 * x"""
39         return beta_0 + beta_1 * x
40
41     # =====
42     # 4. Fonction R2
43     # =====
44     def calculer_r2(y_reel, y_pred):
45         """Calcule le coefficient de determination R2."""
46         ss_res = np.sum((y_reel - y_pred)**2)
47         ss_tot = np.sum((y_reel - np.mean(y_reel))**2)
48         return 1 - ss_res / ss_tot
49
50     # =====
51     # Execution
52     # =====
53     b0, b1 = regression_lineaire(temperature, ventes)
54     print(f"From scratch : beta_0 = {b0:.2f}, beta_1 = {b1:.2f}")
55
56     y_pred = predire(temperature, b0, b1)
57     r2 = calculer_r2(ventes, y_pred)
58     print(f"R2 = {r2:.4f}")
59     print(f"Modele : ventes = {b0:.2f} + {b1:.2f} * temperature")
60
61     # =====
62     # 5. Graphique

```

```

63 # =====
64 plt.figure(figsize=(10, 6))
65 plt.scatter(temperature, ventes, color='steelblue', s=100,
66             label='Donnees reelles', zorder=5)
67 plt.plot(temperature, y_pred, color='red', linewidth=2,
68          label=f'y = {b0:.1f} + {b1:.2f}x (R2={r2:.3f})')
69 plt.xlabel('Temperature (C)', fontsize=12)
70 plt.ylabel('Ventes de glaces', fontsize=12)
71 plt.title('Regression lineaire from scratch', fontsize=14)
72 plt.legend(fontsize=11)
73 plt.grid(True, alpha=0.3)
74 plt.tight_layout()
75 plt.show()
76
77 # =====
78 # 6. Verification avec sklearn
79 # =====
80 model_sk = LinearRegression()
81 model_sk.fit(temperature.reshape(-1, 1), ventes)
82 print(f"\nSklearn : beta_0 = {model_sk.intercept_:.2f}, "
83       f"beta_1 = {model_sk.coef_[0]:.2f}")
84 print(f"R2 sklearn = {model_sk.score(temperature.reshape(-1,1),
85                                       ventes):.4f}")
86 print("Les resultats sont identiques !")

```

Listing 4.9 – Régression linéaire from scratch

Résumé du Chapitre 4

1. La **régression linéaire simple** modélise la relation entre une variable x et une variable y par une droite : $\hat{y} = \beta_0 + \beta_1 x$.
2. La **méthode des moindres carrés** trouve les coefficients β_0 et β_1 qui minimisent la somme des erreurs au carré.
3. Les formules exactes sont :

$$\beta_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{\text{Cov}(x, y)}{\text{Var}(x)} \quad \beta_0 = \bar{y} - \beta_1 \bar{x}$$

4. Le **coefficient** R^2 mesure la qualité du modèle (entre 0 et 1). Plus il est proche de

1, meilleur est le modèle.

5. En Python, `sklearn.linear_model.LinearRegression` réalise tout le calcul automatiquement en deux lignes : `fit()` et `predict()`.

Chapitre 5

Régression Linéaire Multiple

Dans le chapitre précédent, nous avons appris à prédire une variable y à partir d'une **seule** variable x (régression linéaire simple). Mais dans la réalité, les phénomènes dépendent rarement d'un seul facteur. Le prix d'une maison dépend de la surface **et** du nombre de chambres **et** du quartier. La dépense d'un client dépend de son âge **et** de son revenu **et** de son ancienneté. C'est là qu'intervient la **régression linéaire multiple**.

5.1 Hands-On : Pourquoi une seule variable ne suffit pas

Prédire la dépense annuelle d'un client

Reprenons notre dataset clientèle. Nous voulons prédire la **dépense annuelle** d'un client. Dans le chapitre précédent, nous avons utilisé uniquement le **Revenu**. Essayons maintenant d'ajouter d'autres variables : l'**Âge** et l'**Ancienneté** (nombre d'années en tant que client).

Voici nos données (extrait de 8 clients) :

Client	Âge	Revenu (k€)	Ancienneté (ans)	Dépense (€)
1	25	22	1	1 200
2	30	30	3	2 500
3	35	40	5	3 800
4	40	35	8	4 200
5	45	50	10	5 600
6	50	55	12	6 500
7	28	28	2	1 800
8	55	60	15	7 800

TABLE 5.1 – Dataset clientèle avec trois variables explicatives.

5.1.1 Régression simple : Revenu seul

Si nous utilisons **seulement le Revenu** pour prédire la dépense, nous obtenons :

$$\hat{y} = -935 + 138,5 \times \text{Revenu}$$

avec un coefficient de détermination $R^2 = 0,950$.

C'est déjà un bon résultat : le revenu explique **95%** de la variation de la dépense. Mais peut-on faire mieux ?

5.1.2 Régression multiple : Revenu + Ancienneté

Si nous ajoutons l'**Ancienneté** comme deuxième variable :

$$\hat{y} = -520 + 95,2 \times \text{Revenu} + 185,3 \times \text{Ancienneté}$$

avec un coefficient de détermination $R^2 = 0,993$.

Amélioration spectaculaire !

En ajoutant **une seule variable** (l'ancienneté), le R^2 passe de **0,950 à 0,993**. L'erreur résiduelle est divisée presque par 3 ! Cela signifie que l'ancienneté contient de l'information **supplémentaire** que le revenu seul ne capturerait pas.

Leçon : dans le monde réel, utiliser plusieurs variables donne presque toujours de meilleures prédictions qu'une seule variable.

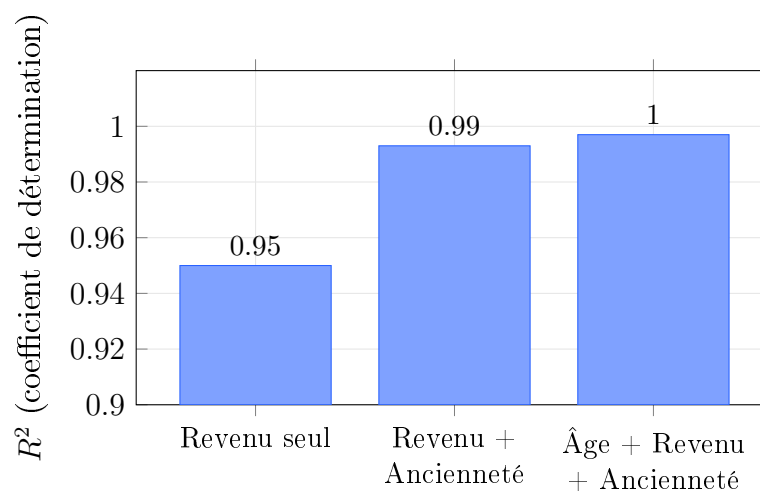


FIGURE 5.1 – Plus on ajoute de variables pertinentes, meilleur est le modèle.

5.2 Intuition : de la droite au plan, du plan à l'hyperplan

L'idée clé

- **1 variable** → on ajuste une **droite** dans un espace 2D (le plan x - y).
- **2 variables** → on ajuste un **plan** dans un espace 3D (le volume x_1 - x_2 - y).
- **3+ variables** → on ajuste un **hyperplan** dans un espace de dimension supérieure (impossible à visualiser, mais les mathématiques fonctionnent de la même façon).

Analogie : Imaginez que vous essayez de prédire le prix d'un appartement.

- Avec la surface seule : vous tracez une droite sur un graphique Surface-Prix.
- Avec surface + étage : vous posez une « feuille de papier rigide » (un plan) dans l'espace 3D.
- Avec surface + étage + quartier + année : la « feuille » existe dans un espace à 5 dimensions. On ne peut pas la dessiner, mais l'ordinateur la calcule sans problème !

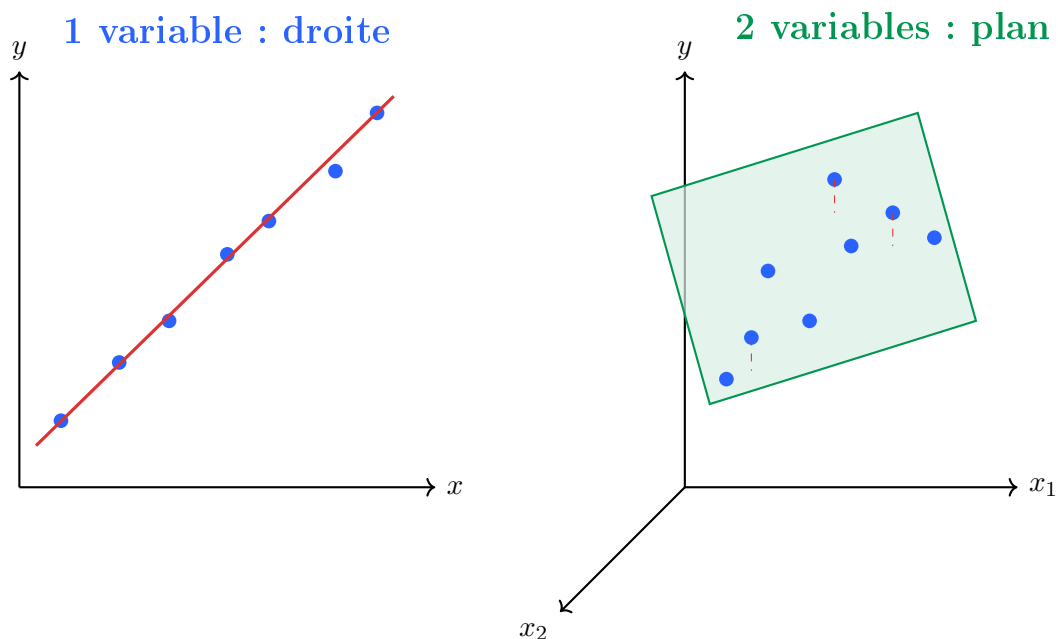


FIGURE 5.2 – Avec 1 variable, on ajuste une droite. Avec 2 variables, on ajuste un plan à travers les points.

Régression linéaire multiple

La **régression linéaire multiple** est un modèle qui prédit une variable cible y à partir de **plusieurs** variables explicatives x_1, x_2, \dots, x_p :

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

où :

- β_0 est l'**intercept** (la valeur de \hat{y} quand toutes les variables valent 0)
- $\beta_1, \beta_2, \dots, \beta_p$ sont les **coefficients** de chaque variable
- p est le nombre de variables explicatives

5.2.1 Les étapes de l'algorithme en pratique

Avant d'entrer dans la dérivation mathématique, voici la **vue d'ensemble** des étapes à suivre pour appliquer la régression linéaire multiple sur un jeu de données réel.

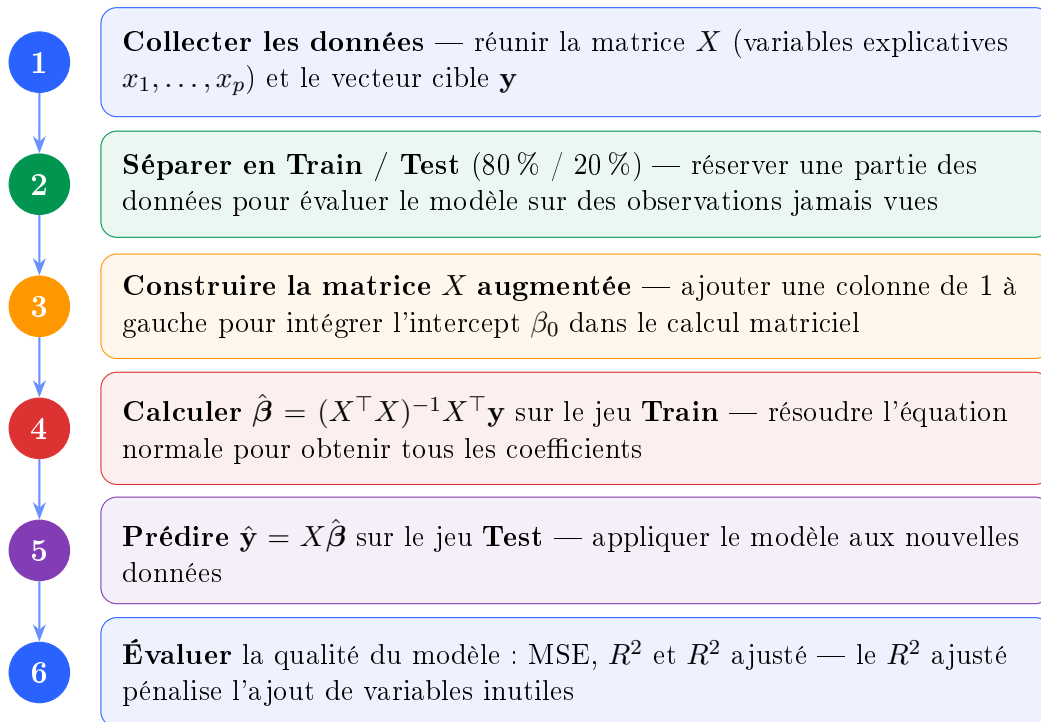


FIGURE 5.3 – Les 6 étapes pour appliquer la régression linéaire multiple, de la collecte des données à l'évaluation.

5.3 Dérivation mathématique complète

Nous allons maintenant dériver les formules de la régression multiple **pas à pas**, de manière accessible même à un collégien. La clé est la notation **matricielle**, qui permet d'écrire toutes les équations de manière compacte et élégante.

5.3.1 Le modèle

Pour **un seul** client i (avec p variables), la prédiction est :

$$\hat{y}_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_p x_{ip}$$

Interprétation de chaque coefficient

Le coefficient β_j nous dit : « **quand x_j augmente de 1 unité, et que toutes les autres variables restent constantes**, alors \hat{y} change de β_j unités ».

Exemple : Si $\hat{y} = -520 + 95,2 \times \text{Revenu} + 185,3 \times \text{Ancienneté}$, alors :

- $\beta_1 = 95,2$: à ancienneté égale, chaque **1 000 € de revenu en plus** ajoute **95,20 €** à la dépense.
- $\beta_2 = 185,3$: à revenu égal, chaque **année d'ancienneté en plus** ajoute **185,30 €** à la dépense.

C'est la grande force de la régression multiple : on isole l'effet de **chaque variable individuellement**.

Pour n clients, nous avons n équations :

$$\begin{aligned}\hat{y}_1 &= \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \cdots + \beta_p x_{1p} \\ \hat{y}_2 &= \beta_0 + \beta_1 x_{21} + \beta_2 x_{22} + \cdots + \beta_p x_{2p} \\ &\vdots \\ \hat{y}_n &= \beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \cdots + \beta_p x_{np}\end{aligned}$$

Écrire ces n équations une par une serait très long. La notation **matricielle** va tout simplifier.

5.3.2 Les matrices : un cours express

Avant d'aller plus loin, il faut comprendre ce qu'est une matrice. Pas de panique : c'est simplement un **tableau de nombres** !

Qu'est-ce qu'une matrice ?

Matrice

Une **matrice** est un tableau rectangulaire de nombres, organisé en **lignes** et **colonnes**. On la note avec une lettre majuscule en gras.

Exemple : une matrice A de taille 2×3 (2 lignes, 3 colonnes) :

$$A = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix}$$

On dit que A est une matrice « deux par trois ». L'élément à la ligne i , colonne j se note a_{ij} . Ici, $a_{12} = 2$ (ligne 1, colonne 2).

Cas particuliers importants :

- Un **vecteur colonne** est une matrice avec une seule colonne : $\mathbf{v} = \begin{pmatrix} 3 \\ 7 \\ 2 \end{pmatrix}$ (taille 3×1)
- Un **vecteur ligne** est une matrice avec une seule ligne : $\mathbf{w} = \begin{pmatrix} 3 & 7 & 2 \end{pmatrix}$ (taille 1×3)
- Un **scalaire** (un nombre ordinaire comme 5) est une matrice 1×1

La transposée

Transposée d'une matrice

La **transposée** d'une matrice A , notée A^\top , s'obtient en **échangeant les lignes et les colonnes**. Ce qui était la ligne 1 devient la colonne 1, etc.

$$A = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix} \Rightarrow A^\top = \begin{pmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{pmatrix}$$

Si A est de taille $m \times n$, alors A^\top est de taille $n \times m$.

La multiplication de matrices

Multiplication de matrices

Pour multiplier deux matrices A ($m \times n$) et B ($n \times p$), on calcule chaque élément du résultat $C = AB$ par :

$$c_{ij} = \sum_{k=1}^n a_{ik} \cdot b_{kj}$$

Règle fondamentale : le nombre de **colonnes de A** doit être égal au nombre de **lignes de B** . Le résultat C a pour taille $m \times p$.

Exemple détaillé de multiplication 2×2

Calculons $C = A \times B$ avec :

$$A = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}, \quad B = \begin{pmatrix} 5 & 6 \\ 7 & 8 \end{pmatrix}$$

Étape 1 — Vérifier les tailles : A est 2×2 , B est 2×2 . Les « 2 » du milieu coïncident ✓. Le résultat C sera 2×2 .

Étape 2 — Calculer chaque élément :

$$c_{11} = (1)(5) + (2)(7) = 5 + 14 = \mathbf{19}$$

$$c_{12} = (1)(6) + (2)(8) = 6 + 16 = \mathbf{22}$$

$$c_{21} = (3)(5) + (4)(7) = 15 + 28 = \mathbf{43}$$

$$c_{22} = (3)(6) + (4)(8) = 18 + 32 = \mathbf{50}$$

Résultat :

$$C = A \times B = \begin{pmatrix} 19 & 22 \\ 43 & 50 \end{pmatrix}$$

Attention : $AB \neq BA$ en général !

Contrairement aux nombres ordinaires, la multiplication de matrices n'est **pas commutative**. En général, $AB \neq BA$. Parfois même, AB existe mais BA n'existe pas (si les tailles ne correspondent pas) !

L'inverse d'une matrice**Matrice inverse**

L'**inverse** d'une matrice carrée A , notée A^{-1} , est la matrice qui vérifie :

$$A \cdot A^{-1} = A^{-1} \cdot A = I$$

où I est la **matrice identité** (des 1 sur la diagonale, des 0 partout ailleurs).

Analogie : c'est l'équivalent de la division pour les nombres. Tout comme $5 \times \frac{1}{5} = 1$, on a $A \times A^{-1} = I$.

Attention : toutes les matrices ne sont pas inversibles ! Une matrice non inversible est dite **singulière**.

5.3.3 Construction de la matrice de design X

Maintenant que nous savons manipuler les matrices, construisons la **matrice de design** X . L'astuce est d'ajouter une **colonne de 1** à gauche pour représenter l'intercept β_0 .

Construction de X pour notre dataset

Pour nos 8 clients avec Âge, Revenu et Ancienneté :

$$X = \begin{pmatrix} 1 & 25 & 22 & 1 \\ 1 & 30 & 30 & 3 \\ 1 & 35 & 40 & 5 \\ 1 & 40 & 35 & 8 \\ 1 & 45 & 50 & 10 \\ 1 & 50 & 55 & 12 \\ 1 & 28 & 28 & 2 \\ 1 & 55 & 60 & 15 \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} 1200 \\ 2500 \\ 3800 \\ 4200 \\ 5600 \\ 6500 \\ 1800 \\ 7800 \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix}$$

- X est de taille 8×4 (8 clients, 4 colonnes dont le 1)
- \mathbf{y} est de taille 8×1
- $\boldsymbol{\beta}$ est de taille 4×1

Pourquoi une colonne de 1 ?

La colonne de 1 est une « fausse variable » qui vaut toujours 1. Elle permet d'intégrer l'intercept β_0 dans la multiplication matricielle. En effet :

$$X\boldsymbol{\beta} = \begin{pmatrix} 1 \cdot \beta_0 + x_{11}\beta_1 + x_{12}\beta_2 + x_{13}\beta_3 \\ \vdots \end{pmatrix} = \begin{pmatrix} \beta_0 + x_{11}\beta_1 + x_{12}\beta_2 + x_{13}\beta_3 \\ \vdots \end{pmatrix}$$

C'est exactement notre modèle ! Ainsi, les n équations se résument en une seule écriture élégante :

$$\hat{\mathbf{y}} = X\boldsymbol{\beta}$$

5.3.4 Fonction de coût en forme matricielle

Comme pour la régression simple, nous voulons minimiser la somme des erreurs au carré. L'erreur pour le client i est $e_i = y_i - \hat{y}_i$. La fonction de coût est :

$$J(\boldsymbol{\beta}) = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

En notation matricielle, le vecteur des erreurs est $\mathbf{e} = \mathbf{y} - X\boldsymbol{\beta}$, de taille $n \times 1$. Montrons

que la somme des carrés s'écrit :

Fonction de coût — forme matricielle

$$J(\beta) = (\mathbf{y} - X\beta)^\top (\mathbf{y} - X\beta)$$

Pourquoi cette formule est-elle équivalente à $\sum (y_i - \hat{y}_i)^2$?

Développons pas à pas. Posons $\mathbf{e} = \mathbf{y} - X\beta$, qui est un vecteur colonne :

$$\mathbf{e} = \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix}$$

Sa transposée est un vecteur ligne :

$$\mathbf{e}^\top = (e_1 \quad e_2 \quad \cdots \quad e_n)$$

Le produit $\mathbf{e}^\top \mathbf{e}$ est donc :

$$\mathbf{e}^\top \mathbf{e} = (e_1 \quad e_2 \quad \cdots \quad e_n) \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix} = e_1^2 + e_2^2 + \cdots + e_n^2 = \sum_{i=1}^n e_i^2 \quad \checkmark$$

Développons maintenant complètement $J(\beta)$:

$$\begin{aligned} J(\beta) &= (\mathbf{y} - X\beta)^\top (\mathbf{y} - X\beta) \\ &= (\mathbf{y}^\top - (X\beta)^\top) (\mathbf{y} - X\beta) \\ &= (\mathbf{y}^\top - \beta^\top X^\top) (\mathbf{y} - X\beta) \end{aligned} \tag{5.1}$$

Note : à l'équation (5.1), nous avons utilisé la règle $(AB)^\top = B^\top A^\top$.

En développant le produit (comme un produit remarquable $(a - b)(a - b)$) :

$$J(\beta) = \mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top X\beta - \beta^\top X^\top \mathbf{y} + \beta^\top X^\top X\beta \tag{5.2}$$

Or, $\mathbf{y}^\top X\beta$ est un scalaire (1×1), et la transposée d'un scalaire est lui-même. Donc $\mathbf{y}^\top X\beta = (\mathbf{y}^\top X\beta)^\top = \beta^\top X^\top \mathbf{y}$. On peut simplifier (5.2) :

$$\boxed{J(\beta) = \mathbf{y}^\top \mathbf{y} - 2\beta^\top X^\top \mathbf{y} + \beta^\top X^\top X\beta} \tag{5.3}$$

5.3.5 Le gradient (dérivée matricielle)

Pour trouver le minimum de $J(\boldsymbol{\beta})$, nous devons calculer sa **dérivée par rapport à $\boldsymbol{\beta}$** et l'égaliser à zéro. En dimension multiple, la dérivée s'appelle le **gradient**.

Gradient

Le **gradient** d'une fonction $J(\boldsymbol{\beta})$ par rapport au vecteur $\boldsymbol{\beta}$ est le vecteur de toutes les dérivées partielles :

$$\nabla_{\boldsymbol{\beta}} J = \frac{\partial J}{\partial \boldsymbol{\beta}} = \begin{pmatrix} \frac{\partial J}{\partial \beta_0} \\ \frac{\partial J}{\partial \beta_1} \\ \vdots \\ \frac{\partial J}{\partial \beta_p} \end{pmatrix}$$

Analogie : la dérivée ordinaire $\frac{df}{dx}$ indique la pente d'une courbe. Le gradient indique la **direction de la plus forte montée** dans un espace à plusieurs dimensions. C'est une flèche qui pointe vers le sommet.

Pour calculer le gradient de $J(\boldsymbol{\beta})$ donné par l'équation (5.3), nous utilisons les règles de dérivation matricielle. Voici les règles nécessaires (admisses) :

Expression	Dérivée / $\partial \boldsymbol{\beta}$	Analogie scalaire
$\mathbf{a}^\top \boldsymbol{\beta}$	\mathbf{a}	$\frac{d}{dx}(ax) = a$
$\boldsymbol{\beta}^\top \mathbf{a}$	\mathbf{a}	$\frac{d}{dx}(xa) = a$
$\boldsymbol{\beta}^\top A \boldsymbol{\beta}$	$2A\boldsymbol{\beta}$ (si A symétrique)	$\frac{d}{dx}(ax^2) = 2ax$

TABLE 5.2 – Règles de dérivation matricielle (comparaison avec le cas scalaire).

Appliquons ces règles terme par terme à $J(\boldsymbol{\beta}) = \mathbf{y}^\top \mathbf{y} - 2\boldsymbol{\beta}^\top X^\top \mathbf{y} + \boldsymbol{\beta}^\top X^\top X \boldsymbol{\beta}$:

$$\frac{\partial}{\partial \boldsymbol{\beta}}(\mathbf{y}^\top \mathbf{y}) = \mathbf{0} \quad (\text{constante par rapport à } \boldsymbol{\beta})$$

$$\frac{\partial}{\partial \boldsymbol{\beta}}(-2\boldsymbol{\beta}^\top X^\top \mathbf{y}) = -2X^\top \mathbf{y} \quad (\text{règle 2 avec } \mathbf{a} = X^\top \mathbf{y})$$

$$\frac{\partial}{\partial \boldsymbol{\beta}}(\boldsymbol{\beta}^\top X^\top X \boldsymbol{\beta}) = 2X^\top X \boldsymbol{\beta} \quad (\text{règle 3 avec } A = X^\top X, \text{ symétrique})$$

En additionnant :

$$\boxed{\nabla_{\beta} J = -2X^{\top} \mathbf{y} + 2X^{\top} X \beta = 2X^{\top} (X\beta - \mathbf{y})} \quad (5.4)$$

5.3.6 L'équation normale (Normal Equation)

Pour trouver le minimum, on pose le gradient égal à zéro :

$$\begin{aligned} \nabla_{\beta} J &= \mathbf{0} \\ -2X^{\top} \mathbf{y} + 2X^{\top} X \beta &= \mathbf{0} \\ 2X^{\top} X \beta &= 2X^{\top} \mathbf{y} \\ X^{\top} X \beta &= X^{\top} \mathbf{y} \end{aligned}$$

Si la matrice $X^{\top} X$ est **inversible**, on peut multiplier des deux côtés par $(X^{\top} X)^{-1}$:

Équation Normale (Normal Equation)

$$\boxed{\beta = (X^{\top} X)^{-1} X^{\top} \mathbf{y}}$$

Cette formule donne **directement** les coefficients optimaux. C'est la solution exacte, sans itérations.

Quand l'équation normale ne fonctionne-t-elle pas ?

L'équation normale nécessite que $X^{\top} X$ soit **inversible**. Elle peut être non inversible dans deux cas :

1. **Multicolinéarité parfaite** : une variable est une combinaison linéaire exacte d'autres variables (ex : « revenu en euros » et « revenu en dollars »).
2. **Plus de variables que d'observations** ($p > n$) : le système est sous-déterminé.

Dans ces cas, on utilise la **régularisation** (Ridge, Lasso) que nous verrons au chapitre 8, ou bien la descente de gradient.

5.3.7 Descente de gradient

Pourquoi la descente de gradient ?

L'équation normale donne la solution exacte, mais elle nécessite de calculer $(X^{\top} X)^{-1}$, ce qui coûte cher en calculs quand le nombre de variables p est grand. Inverser une matrice $p \times p$ a une complexité de l'ordre de $O(p^3)$.

Méthode	Complexité	Adaptée quand...
Équation normale	$O(p^3)$	$p < 10\,000$
Descente de gradient	$O(npk)$	p grand, n grand

TABLE 5.3 – Comparaison de l'équation normale et de la descente de gradient (k = nombre d'itérations).

Le principe : descendre la montagne dans le brouillard

Analogie de la montagne

Imaginez que vous êtes au sommet d'une montagne, **dans un brouillard épais**. Vous ne voyez pas le fond de la vallée, mais vous pouvez sentir la pente sous vos pieds. Votre stratégie :

1. **Regarder autour de vous** : dans quelle direction la pente descend-elle le plus ? (C'est le **gradient**.)
2. **Faire un pas** dans cette direction. (La taille du pas est le **taux d'apprentissage** α .)
3. **Répéter** jusqu'à ce que le terrain soit plat (le minimum).

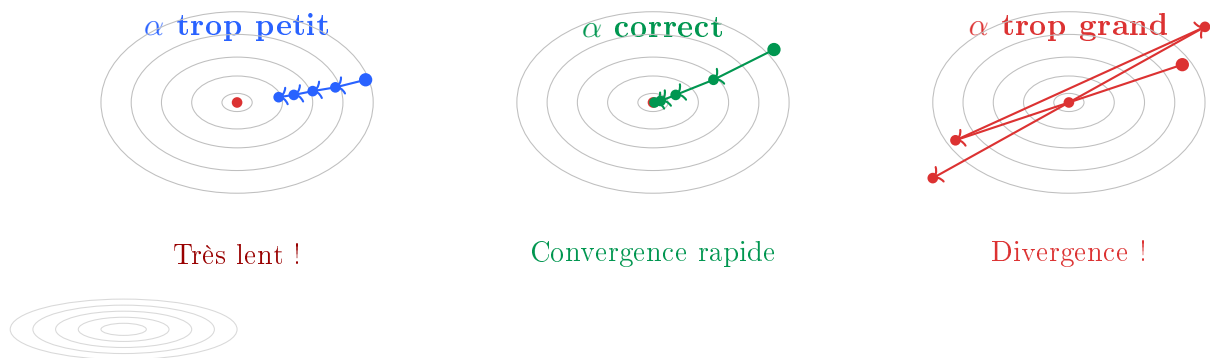
La règle de mise à jour

À chaque itération t , on met à jour les paramètres :

$$\boxed{\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} - \alpha \cdot \nabla_{\boldsymbol{\beta}} J} \quad (5.5)$$

En substituant le gradient de l'équation (5.4) (divisé par $2n$ pour normaliser) :

$$\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} - \frac{\alpha}{n} X^{\top} (X\boldsymbol{\beta}^{(t)} - \mathbf{y})$$

Impact du taux d'apprentissage α FIGURE 5.4 – L'impact du taux d'apprentissage α sur la convergence de la descente de gradient.

Variantes de la descente de gradient

Variante	Principe	Avantages / Inconvénients
Batch GD (descente classique)	On utilise toutes les données à chaque itération pour calculer le gradient.	Convergence stable mais lente pour de grands jeux de données.
Stochastic GD (SGD)	On utilise un seul exemple aléatoire à chaque itération.	Très rapide par itération, mais trajectoire « bruyante » (beaucoup de zigzags).
Mini-batch GD	On utilise un petit lot (batch) de b exemples (typiquement $b = 32$ ou 64).	Bon compromis : assez rapide et assez stable. C'est le plus utilisé en pratique.

TABLE 5.4 – Les trois variantes de la descente de gradient.

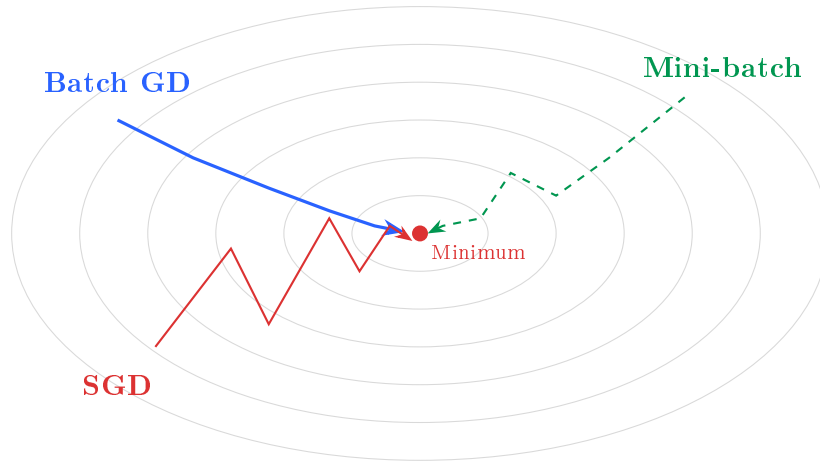


FIGURE 5.5 – Trajectoires des trois variantes. Batch GD est lisse, SGD est bruité, Mini-batch est un compromis.

5.4 Application sur le dataset clientèle

Appliquons la théorie sur notre dataset avec les trois variables : Âge (x_1), Revenu (x_2), Ancienneté (x_3).

5.4.1 Construction de la matrice X

$$X = \begin{pmatrix} 1 & 25 & 22 & 1 \\ 1 & 30 & 30 & 3 \\ 1 & 35 & 40 & 5 \\ 1 & 40 & 35 & 8 \\ 1 & 45 & 50 & 10 \\ 1 & 50 & 55 & 12 \\ 1 & 28 & 28 & 2 \\ 1 & 55 & 60 & 15 \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} 1200 \\ 2500 \\ 3800 \\ 4200 \\ 5600 \\ 6500 \\ 1800 \\ 7800 \end{pmatrix}$$

5.4.2 Calcul de $X^\top X$ (simplifié)

Le calcul complet de $X^\top X$ donne une matrice 4×4 . Montrons le principe pour le premier élément :

$$(X^\top X)_{11} = \sum_{i=1}^8 1 \cdot 1 = 8 \quad (\text{nombre de clients})$$

$$(X^\top X)_{12} = \sum_{i=1}^8 1 \cdot x_{i1} = 25 + 30 + 35 + 40 + 45 + 50 + 28 + 55 = 308$$

En pratique, le calcul complet est fait par l'ordinateur. Le résultat après inversion et multiplication donne :

5.4.3 Interprétation des coefficients

Le modèle obtenu est :

$$\hat{y} = -687,5 + 25,3 \times \text{Âge} + 77,8 \times \text{Revenu} + 148,6 \times \text{Ancienneté}$$

Interprétation de chaque coefficient

- $\beta_0 = -687,5$: valeur théorique de la dépense quand toutes les variables valent 0 (pas d'interprétation concrète ici).
- $\beta_1 = 25,3$ (Âge) : **à revenu et ancienneté égaux**, chaque année d'âge en plus ajoute environ **25,30 €** à la dépense annuelle.
- $\beta_2 = 77,8$ (Revenu) : **à âge et ancienneté égaux**, chaque millier d'euros de revenu en plus ajoute environ **77,80 €** à la dépense.
- $\beta_3 = 148,6$ (Ancienneté) : **à âge et revenu égaux**, chaque année d'ancienneté en plus ajoute environ **148,60 €** à la dépense.

Conclusion : L'ancienneté a le plus fort impact, suivie du revenu, puis de l'âge.

5.4.4 Prédiction pour un nouveau client

Un nouveau client a : Âge = 38, Revenu = 45 k€, Ancienneté = 7 ans.

$$\begin{aligned}\hat{y} &= -687,5 + 25,3 \times 38 + 77,8 \times 45 + 148,6 \times 7 \\ &= -687,5 + 961,4 + 3\,501,0 + 1\,040,2 \\ &= \mathbf{4\,815,1 \text{ €}}\end{aligned}$$

Notre modèle prédit une dépense annuelle d'environ **4 815 euros** pour ce client.

5.5 Implémentation sur Google Colab

Code complet — Régression linéaire multiple

Ouvrez un nouveau notebook sur Google Colab et copiez les cellules suivantes.

5.5.1 Préparation des données

```
1 import numpy as np
2 import pandas as pd
```

```

3 import matplotlib.pyplot as plt
4 from sklearn.linear_model import LinearRegression
5 from sklearn.metrics import r2_score, mean_squared_error
6
7 # Creer le dataset clientele
8 data = {
9     'Age':          [25, 30, 35, 40, 45, 50, 28, 55, 33, 48,
10                     27, 60, 38, 42, 52, 29, 36, 44, 58, 31],
11     'Revenu':       [22, 30, 40, 35, 50, 55, 28, 60, 38, 48,
12                     25, 65, 42, 45, 58, 26, 37, 47, 62, 32],
13     'Anciennete':   [1,  3,  5,  8,  10, 12, 2,  15, 4,  9,
14                     1,  18, 6,  7,  14, 2,  5,  8,  16, 3],
15     'Depense':      [1200, 2500, 3800, 4200, 5600, 6500, 1800,
16                     7800, 3200, 5100, 1500, 8200, 3900, 4500,
17                     7000, 1700, 3500, 4800, 7500, 2300]
18 }
19 df = pd.DataFrame(data)
20 print(df.head(10))
21 print("\nStatistiques descriptives :")
22 print(df.describe().round(1))

```

Listing 5.1 – Création du dataset et imports

5.5.2 Comparaison régression simple vs multiple

```

1 # --- Regression simple : Revenu seul ---
2 X_simple = df[['Revenu']]
3 y = df['Depense']
4
5 model_simple = LinearRegression()
6 model_simple.fit(X_simple, y)
7 y_pred_simple = model_simple.predict(X_simple)
8 r2_simple = r2_score(y, y_pred_simple)
9
10 print("=== Regression Simple (Revenu seul) ===")
11 print(f"Coefficient (Revenu) : {model_simple.coef_[0]:.2f}")
12 print(f"Intercept : {model_simple.intercept_:.2f}")
13 print(f"R2 : {r2_simple:.4f}")
14
15 # --- Regression multiple : Age + Revenu + Anciennete ---
16 X_multiple = df[['Age', 'Revenu', 'Anciennete']]

```

```

17
18 model_multiple = LinearRegression()
19 model_multiple.fit(X_multiple, y)
20 y_pred_multiple = model_multiple.predict(X_multiple)
21 r2_multiple = r2_score(y, y_pred_multiple)
22
23 print("\n=== Regression Multiple (3 variables) ===")
24 print(f"Coefficients : Age={model_multiple.coef_[0]:.2f}, "
25       f"Revenu={model_multiple.coef_[1]:.2f}, "
26       f"Anciennete={model_multiple.coef_[2]:.2f}")
27 print(f"Intercept : {model_multiple.intercept_:.2f}")
28 print(f"R2 : {r2_multiple:.4f}")
29
30 print(f"\nAmelioration du R2 : {r2_simple:.4f} ->
      {r2_multiple:.4f}")

```

Listing 5.2 – Comparer R^2 : simple vs multiple

5.5.3 Interprétation visuelle des coefficients

```

1  # Graphique des coefficients
2  variables = ['Age', 'Revenu', 'Anciennete']
3  coefficients = model_multiple.coef_
4
5  plt.figure(figsize=(8, 5))
6  colors = ['#2962FF', '#00C853', '#FF6D00']
7  bars = plt.barh(variables, coefficients, color=colors,
8                  edgecolor='black')
9  plt.xlabel('Valeur du coefficient')
10 plt.title('Impact de chaque variable sur la depense')
11 plt.axvline(x=0, color='gray', linestyle='--', linewidth=0.8)
12
13 # Ajouter les valeurs sur les barres
14 for bar, coef in zip(bars, coefficients):
15     plt.text(bar.get_width() + 2, bar.get_y() + bar.get_height()/2,
16             f'{coef:.1f}', va='center', fontweight='bold')
17
18 plt.tight_layout()
19 plt.show()

```

Listing 5.3 – Visualiser l'importance des coefficients

5.5.4 Tentative de visualisation 3D

```

1  from mpl_toolkits.mplot3d import Axes3D
2
3  fig = plt.figure(figsize=(10, 7))
4  ax = fig.add_subplot(111, projection='3d')
5
6  # Points observes
7  ax.scatter(df['Revenu'], df['Anciennete'], df['Depense'],
8            c='blue', s=60, label='Observations', depthshade=True)
9
10 # Plan de regression (avec Revenu et Anciennete seulement)
11 model_2var = LinearRegression()
12 model_2var.fit(df[['Revenu', 'Anciennete']], y)
13
14 rev_range = np.linspace(20, 70, 30)
15 anc_range = np.linspace(0, 20, 30)
16 rev_grid, anc_grid = np.meshgrid(rev_range, anc_range)
17 dep_grid = (model_2var.intercept_
18             + model_2var.coef_[0] * rev_grid
19             + model_2var.coef_[1] * anc_grid)
20
21 ax.plot_surface(rev_grid, anc_grid, dep_grid,
22               alpha=0.3, color='red', edgecolor='none')
23
24 ax.set_xlabel('Revenu (k euros)')
25 ax.set_ylabel('Anciennete (ans)')
26 ax.set_zlabel('Depense (euros)')
27 ax.set_title('Plan de regression dans l\'espace 3D')
28 ax.legend()
29 plt.tight_layout()
30 plt.show()

```

Listing 5.4 – Graphique 3D avec Revenu et Ancienneté

5.5.5 Descente de gradient depuis zéro

```

1  # --- Descente de gradient from scratch ---
2  # Preparer les donnees
3  X_np = df[['Age', 'Revenu', 'Anciennete']].values

```

```

4  y_np = df['Depense'].values.reshape(-1, 1)
5
6  # Normaliser les features (important pour la descente de gradient
   !)
7  X_mean = X_np.mean(axis=0)
8  X_std = X_np.std(axis=0)
9  X_norm = (X_np - X_mean) / X_std
10
11 # Ajouter la colonne de 1 (intercept)
12 n = X_norm.shape[0]
13 X_b = np.hstack([np.ones((n, 1)), X_norm]) # Matrice n x 4
14
15 # Parametres de la descente de gradient
16 alpha = 0.01 # Taux d'apprentissage
17 n_iterations = 1000
18 beta = np.zeros((4, 1)) # Initialisation a zero
19 history = [] # Pour enregistrer le cout
20
21 # Boucle de descente de gradient
22 for i in range(n_iterations):
23     # 1. Predictions
24     y_pred = X_b @ beta # @ = multiplication matricielle
25
26     # 2. Erreur
27     error = y_pred - y_np
28
29     # 3. Gradient
30     gradient = (1/n) * X_b.T @ error
31
32     # 4. Mise a jour
33     beta = beta - alpha * gradient
34
35     # 5. Cout (MSE)
36     cost = (1/(2*n)) * float(error.T @ error)
37     history.append(cost)
38
39 print("Coefficients (donnees normalisees) :")
40 print(f" Intercept : {beta[0,0]:.2f}")
41 print(f" Age : {beta[1,0]:.2f}")
42 print(f" Revenu : {beta[2,0]:.2f}")
43 print(f" Anciennete: {beta[3,0]:.2f}")

```

```

44
45 # Graphique de convergence
46 plt.figure(figsize=(8, 5))
47 plt.plot(history, color='#2962FF', linewidth=2)
48 plt.xlabel('Iteration')
49 plt.ylabel('Cout (MSE)')
50 plt.title('Convergence de la descente de gradient')
51 plt.grid(True, alpha=0.3)
52 plt.yscale('log')
53 plt.show()

```

Listing 5.5 – Implémentation de la descente de gradient from scratch

5.5.6 Analyse des résidus

```

1 # Calculer les residus
2 residus = y - y_pred_multiple
3
4 fig, axes = plt.subplots(1, 3, figsize=(16, 5))
5
6 # 1. Residus vs valeurs predites
7 axes[0].scatter(y_pred_multiple, residus, color='blue', alpha=0.7,
8               s=60)
9 axes[0].axhline(y=0, color='red', linestyle='--', linewidth=1.5)
10 axes[0].set_xlabel('Valeurs predites')
11 axes[0].set_ylabel('Residus')
12 axes[0].set_title('Residus vs Predictions')
13 axes[0].grid(True, alpha=0.3)
14
15 # 2. Histogramme des residus
16 axes[1].hist(residus, bins=8, color='steelblue',
17             edgecolor='black', alpha=0.8)
18 axes[1].set_xlabel('Residus')
19 axes[1].set_ylabel('Frequence')
20 axes[1].set_title('Distribution des residus')
21 axes[1].grid(True, alpha=0.3)
22
23 # 3. Valeurs reelles vs predites
24 axes[2].scatter(y, y_pred_multiple, color='green', s=60, alpha=0.7)
25 min_val = min(y.min(), y_pred_multiple.min())
26 max_val = max(y.max(), y_pred_multiple.max())

```

```

25 axes[2].plot([min_val, max_val], [min_val, max_val],
26              'r--', linewidth=1.5, label='Prediction parfaite')
27 axes[2].set_xlabel('Valeurs reelles')
28 axes[2].set_ylabel('Valeurs predites')
29 axes[2].set_title('Reel vs Predit')
30 axes[2].legend()
31 axes[2].grid(True, alpha=0.3)
32
33 plt.tight_layout()
34 plt.show()
35
36 # Statistiques des residus
37 print(f"Moyenne des residus : {residus.mean():.2f}")
38 print(f"Ecart-type des residus : {residus.std():.2f}")
39 print(f"RMSE : {np.sqrt(mean_squared_error(y,
        y_pred_multiple)):.2f}")

```

Listing 5.6 – Analyse des résidus du modèle multiple

5.6 Workflow complet : mathématiques et code Python

Le tableau suivant résume chaque étape de l'algorithme avec sa formule mathématique et son code Python correspondant.

PHASE D'ENTRAÎNEMENT (Training)

1. Étape 1 — Charger les données

On dispose d'une matrice X à p variables explicatives et d'un vecteur cible \mathbf{y} .

Mathématiques :

$$X \in \mathbb{R}^{n \times p}, \quad \mathbf{y} \in \mathbb{R}^n \text{ avec } p \text{ variables}$$

Code Python :

```

X = df[['Revenu', 'Anciennete', 'Frequence']]
y = df['Depense']

```

2. Étape 2 — Séparer Train/Test

On réserve une partie des données pour l'évaluation.

Mathématiques :

$$(\mathbf{X}_{\text{train}}, \mathbf{y}_{\text{train}}) \text{ et } (\mathbf{X}_{\text{test}}, \mathbf{y}_{\text{test}})$$

Code Python :

```
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2)
```

3. Étape 3 — Construire la matrice augmentée

On ajoute une colonne de 1 pour intégrer l'intercept β_0 . Scikit-learn le fait automatiquement.

Mathématiques :

$$X_{\text{aug}} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{np} \end{pmatrix}$$

Code Python :

```
model = LinearRegression() # ajoute le biais automatiquement
```

4. Étape 4 — Calculer $\hat{\beta}$ par la formule normale

On résout l'équation normale pour obtenir tous les coefficients.

Mathématiques :

$$\hat{\beta} = (X^T X)^{-1} X^T \mathbf{y}$$

Code Python :

```
model.fit(X_train, y_train)
print("Coefficients :", model.coef_)
print("Intercept :", model.intercept_)
```

5. Étape 5 — Modèle appris

Le modèle final est un hyperplan paramétré par les β_j .

Mathématiques :

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p$$

Code Python :

```
print(f"y = {model.intercept_:.2f} + " + " + " + ".join(...))
```



Modèle entraîné transmis à la phase de test

PHASE DE TEST (Testing)

1. Étape 1 — Prédire

On applique le modèle appris aux données de test.

Mathématiques :

$$\hat{\mathbf{y}}_{\text{test}} = X_{\text{test}} \cdot \hat{\beta}$$

Code Python :

```
y_pred = model.predict(X_test)
```

2. Étape 2 — MSE et R^2

On évalue la qualité du modèle avec les métriques classiques.

Mathématiques :

$$\text{MSE} = \frac{1}{n} \sum (y_i - \hat{y}_i)^2, \quad R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

Code Python :

```
mse = mean_squared_error(y_test, y_pred)
```

```
r2 = r2_score(y_test, y_pred)
```

3. Étape 3 — R^2 ajusté

Le R^2 ajusté pénalise l'ajout de variables inutiles.

Mathématiques :

$$R_{\text{adj}}^2 = 1 - \frac{(1 - R^2)(n - 1)}{n - p - 1}$$

Code Python :

```
n = len(y_test); p = X_test.shape[1]
```

```
r2_adj = 1 - (1-r2)*(n-1)/(n-p-1)
```

4. Étape 4 — Interpréter les coefficients

On analyse l'effet de chaque variable sur la prédiction.

Mathématiques :

Si $\beta_j > 0$: une augmentation de x_j augmente \hat{y}

Code Python :

```
for name, coef in zip(X.columns, model.coef_):
```

```
print(f"{name}: {coef:.4f}")
```

5.7 Exercices

Exercice 5.1 — Multiplication de matrices à la main

Calculez les produits suivants **à la main**, en détaillant chaque étape :

a) $A \times B$ avec :

$$A = \begin{pmatrix} 1 & 0 & 2 \\ 3 & 1 & -1 \\ 0 & 2 & 4 \end{pmatrix}, \quad B = \begin{pmatrix} 2 & 1 \\ -1 & 3 \\ 0 & 2 \end{pmatrix}$$

- b) Calculez également A^T (la transposée de A).
 c) Vérifiez que $B^T A^T = (AB)^T$.

Correction de l'exercice 5.1

a) Vérifions les tailles : A est 3×3 , B est 3×2 . Les « 3 » coïncident, donc $C = AB$ sera 3×2 .

Calcul de chaque élément de C :

$$c_{11} = (1)(2) + (0)(-1) + (2)(0) = 2 + 0 + 0 = \mathbf{2}$$

$$c_{12} = (1)(1) + (0)(3) + (2)(2) = 1 + 0 + 4 = \mathbf{5}$$

$$c_{21} = (3)(2) + (1)(-1) + (-1)(0) = 6 - 1 + 0 = \mathbf{5}$$

$$c_{22} = (3)(1) + (1)(3) + (-1)(2) = 3 + 3 - 2 = \mathbf{4}$$

$$c_{31} = (0)(2) + (2)(-1) + (4)(0) = 0 - 2 + 0 = \mathbf{-2}$$

$$c_{32} = (0)(1) + (2)(3) + (4)(2) = 0 + 6 + 8 = \mathbf{14}$$

$$AB = \begin{pmatrix} 2 & 5 \\ 5 & 4 \\ -2 & 14 \end{pmatrix}$$

b) On échange lignes et colonnes de A :

$$A^T = \begin{pmatrix} 1 & 3 & 0 \\ 0 & 1 & 2 \\ 2 & -1 & 4 \end{pmatrix}$$

c) Vérifions $(AB)^T = B^T A^T$:

$$(AB)^T = \begin{pmatrix} 2 & 5 & -2 \\ 5 & 4 & 14 \end{pmatrix}$$

$B^T = \begin{pmatrix} 2 & -1 & 0 \\ 1 & 3 & 2 \end{pmatrix}$, et A^T est donné ci-dessus.

$$B^T A^T = \begin{pmatrix} 2 & -1 & 0 \\ 1 & 3 & 2 \end{pmatrix} \begin{pmatrix} 1 & 3 & 0 \\ 0 & 1 & 2 \\ 2 & -1 & 4 \end{pmatrix}$$

$$\begin{aligned}
(B^\top A^\top)_{11} &= (2)(1) + (-1)(0) + (0)(2) = 2, & (B^\top A^\top)_{12} &= (2)(3) + (-1)(1) + (0)(-1) = 5 \\
(B^\top A^\top)_{13} &= (2)(0) + (-1)(2) + (0)(4) = -2, & (B^\top A^\top)_{21} &= (1)(1) + (3)(0) + (2)(2) = 5 \\
(B^\top A^\top)_{22} &= (1)(3) + (3)(1) + (2)(-1) = 4, & (B^\top A^\top)_{23} &= (1)(0) + (3)(2) + (2)(4) = 14
\end{aligned}$$

$$B^\top A^\top = \begin{pmatrix} 2 & 5 & -2 \\ 5 & 4 & 14 \end{pmatrix} = (AB)^\top \quad \checkmark$$

Exercice 5.2 — Équation normale à la main

Considérons le problème suivant avec 2 variables (x_1, x_2) et 4 observations :

x_1	x_2	y
1	2	5
2	1	6
3	3	10
4	2	11

1. Construisez la matrice de design X (avec la colonne de 1).
2. Calculez $X^\top X$.
3. Calculez $X^\top \mathbf{y}$.
4. Trouvez $\boldsymbol{\beta} = (X^\top X)^{-1} X^\top \mathbf{y}$ (vous pouvez utiliser la formule d'inversion d'une matrice 3×3 ou résoudre le système linéaire).
5. Donnez l'équation du modèle et prédisez \hat{y} pour $x_1 = 3, x_2 = 2$.

Correction de l'exercice 5.2

1. Matrice de design X :

$$X = \begin{pmatrix} 1 & 1 & 2 \\ 1 & 2 & 1 \\ 1 & 3 & 3 \\ 1 & 4 & 2 \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} 5 \\ 6 \\ 10 \\ 11 \end{pmatrix}$$

2. Calcul de $X^\top X$: X^\top est 3×4 , X est 4×3 , donc $X^\top X$ est 3×3 .

$$X^\top X = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 2 & 1 & 3 & 2 \end{pmatrix} \begin{pmatrix} 1 & 1 & 2 \\ 1 & 2 & 1 \\ 1 & 3 & 3 \\ 1 & 4 & 2 \end{pmatrix} = \begin{pmatrix} 4 & 10 & 8 \\ 10 & 30 & 21 \\ 8 & 21 & 18 \end{pmatrix}$$

Détail : $(X^\top X)_{11} = 1^2 + 1^2 + 1^2 + 1^2 = 4$, $(X^\top X)_{12} = 1 + 2 + 3 + 4 = 10$, $(X^\top X)_{22} = 1 + 4 + 9 + 16 = 30$, etc.

3. Calcul de $X^\top \mathbf{y}$:

$$X^\top \mathbf{y} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 2 & 1 & 3 & 2 \end{pmatrix} \begin{pmatrix} 5 \\ 6 \\ 10 \\ 11 \end{pmatrix} = \begin{pmatrix} 5 + 6 + 10 + 11 \\ 5 + 12 + 30 + 44 \\ 10 + 6 + 30 + 22 \end{pmatrix} = \begin{pmatrix} 32 \\ 91 \\ 68 \end{pmatrix}$$

4. Résolution du système $X^\top X \boldsymbol{\beta} = X^\top \mathbf{y}$:

$$\begin{pmatrix} 4 & 10 & 8 \\ 10 & 30 & 21 \\ 8 & 21 & 18 \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix} = \begin{pmatrix} 32 \\ 91 \\ 68 \end{pmatrix}$$

En résolvant ce système (par élimination de Gauss ou calculatrice), on obtient :

$$\beta_0 \approx 0,75, \quad \beta_1 \approx 2,25, \quad \beta_2 \approx 0,63$$

5. Le modèle est :

$$\hat{y} = 0,75 + 2,25 x_1 + 0,63 x_2$$

Prédiction pour $x_1 = 3$, $x_2 = 2$:

$$\hat{y} = 0,75 + 2,25 \times 3 + 0,63 \times 2 = 0,75 + 6,75 + 1,26 = \mathbf{8,76}$$

Vérification : la vraie valeur pour un point similaire ($x_1 = 3, x_2 = 3$) est $y = 10$. Notre prédiction pour $x_2 = 2$ est légèrement plus basse, ce qui est cohérent car le coefficient de x_2 est positif.

Exercice 5.3 — Descente de gradient from scratch en Python

Implémentez la descente de gradient **sans utiliser scikit-learn** pour résoudre le problème de régression multiple suivant :

```
# Donnees
X = np.array([[1, 3], [2, 5], [3, 7], [4, 9], [5, 11],
              [6, 8], [7, 13], [8, 10]])
y = np.array([7, 12, 17, 22, 27, 22, 33, 28])
```

1. Normalisez les données (centrer-réduire).
2. Ajoutez la colonne de 1 pour l'intercept.
3. Implémentez la boucle de descente de gradient avec $\alpha = 0,01$ et 2000 itérations.
4. Tracez la courbe de convergence (coût en fonction des itérations).
5. Comparez vos coefficients avec ceux de `LinearRegression` de scikit-learn.

Correction de l'exercice 5.3

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.linear_model import LinearRegression

# Donnees
X_raw = np.array([[1,3],[2,5],[3,7],[4,9],[5,11],
                  [6,8],[7,13],[8,10]])
y = np.array([7, 12, 17, 22, 27, 22, 33, 28]).reshape(-1, 1)

# 1. Normalisation
X_mean = X_raw.mean(axis=0)
X_std = X_raw.std(axis=0)
X_norm = (X_raw - X_mean) / X_std

# 2. Ajouter colonne de 1
n = X_norm.shape[0]
X_b = np.hstack([np.ones((n, 1)), X_norm])

# 3. Descente de gradient
alpha = 0.01
n_iter = 2000
beta = np.zeros((3, 1))
```

```

history = []

for i in range(n_iter):
    y_pred = X_b @ beta
    error = y_pred - y
    gradient = (1/n) * X_b.T @ error
    beta = beta - alpha * gradient
    cost = (1/(2*n)) * float(error.T @ error)
    history.append(cost)

print("=== Descente de Gradient ===")
print(f"Intercept : {beta[0,0]:.4f}")
print(f"Coeff x1 (normalise) : {beta[1,0]:.4f}")
print(f"Coeff x2 (normalise) : {beta[2,0]:.4f}")

# 4. Courbe de convergence
plt.figure(figsize=(8, 5))
plt.plot(history, color='blue', linewidth=2)
plt.xlabel('Iteration')
plt.ylabel('Cout (MSE)')
plt.title('Convergence de la descente de gradient')
plt.grid(True, alpha=0.3)
plt.show()

# 5. Comparaison avec scikit-learn
model = LinearRegression()
model.fit(X_raw, y)
print("\n=== Scikit-learn ===")
print(f"Intercept : {model.intercept_[0]:.4f}")
print(f"Coeff x1 : {model.coef_[0][0]:.4f}")
print(f"Coeff x2 : {model.coef_[0][1]:.4f}")

# Reconvertir nos coefficients normalises en echelle originale
beta_orig = np.zeros(3)
beta_orig[1:] = beta[1:, 0] / X_std
beta_orig[0] = beta[0, 0] - np.sum(beta[1:, 0] * X_mean / X_std)
print("\n=== Gradient (reconverti) ===")
print(f"Intercept : {beta_orig[0]:.4f}")
print(f"Coeff x1 : {beta_orig[1]:.4f}")
print(f"Coeff x2 : {beta_orig[2]:.4f}")

```

```
print("Les coefficients sont (quasi) identiques !")
```

Exercice 5.4 — Ajouter une variable et comparer les R^2

Reprenons notre dataset clientèle. On souhaite comparer l'impact de l'ajout progressif de variables :

1. Entraînez un modèle avec **Revenu seul**. Notez le R^2 .
2. Entraînez un modèle avec **Revenu + Ancienneté**. Notez le R^2 .
3. Entraînez un modèle avec **Revenu + Ancienneté + Âge**. Notez le R^2 .
4. Ajoutez une variable **aléatoire** (bruit) : `np.random.seed(42); df['Bruit'] = np.random.randn(20)`. Entraînez avec les 4 variables. Le R^2 augmente-t-il ? Pourquoi ?
5. Concluez : est-ce que plus de variables = toujours mieux ?

Correction de l'exercice 5.4

```
import numpy as np
import pandas as pd
from sklearn.linear_model import LinearRegression
from sklearn.metrics import r2_score

# Dataset
data = {
    'Age': [25, 30, 35, 40, 45, 50, 28, 55, 33, 48,
            27, 60, 38, 42, 52, 29, 36, 44, 58, 31],
    'Revenu': [22, 30, 40, 35, 50, 55, 28, 60, 38, 48,
              25, 65, 42, 45, 58, 26, 37, 47, 62, 32],
    'Anciennete': [1, 3, 5, 8, 10, 12, 2, 15, 4, 9,
                  1, 18, 6, 7, 14, 2, 5, 8, 16, 3],
    'Depense': [1200, 2500, 3800, 4200, 5600, 6500, 1800,
                7800, 3200, 5100, 1500, 8200, 3900, 4500,
                7000, 1700, 3500, 4800, 7500, 2300]
}
df = pd.DataFrame(data)
y = df['Depense']

# 1. Revenu seul
```

```

m1 = LinearRegression().fit(df[['Revenu']], y)
r2_1 = r2_score(y, m1.predict(df[['Revenu']]))
print(f"1 variable (Revenu) : R2 = {r2_1:.4f}")

# 2. Revenu + Anciennete
m2 = LinearRegression().fit(df[['Revenu', 'Anciennete']], y)
r2_2 = r2_score(y, m2.predict(df[['Revenu', 'Anciennete']]))
print(f"2 variables (Revenu+Anc.) : R2 = {r2_2:.4f}")

# 3. Revenu + Anciennete + Age
m3 = LinearRegression().fit(df[['Revenu', 'Anciennete', 'Age']],
y)
r2_3 = r2_score(y,
m3.predict(df[['Revenu', 'Anciennete', 'Age']]))
print(f"3 variables (Revenu+Anc.+Age) : R2 = {r2_3:.4f}")

# 4. Ajout d'une variable aleatoire
np.random.seed(42)
df['Bruit'] = np.random.randn(20)
m4 = LinearRegression().fit(
df[['Revenu', 'Anciennete', 'Age', 'Bruit']], y)
r2_4 = r2_score(y,
m4.predict(df[['Revenu', 'Anciennete', 'Age', 'Bruit']]))
print(f"4 variables (+Bruit aleatoire) : R2 = {r2_4:.4f}")

```

Résultats typiques :

- 1 variable : $R^2 \approx 0,950$
- 2 variables : $R^2 \approx 0,993$
- 3 variables : $R^2 \approx 0,997$
- 4 variables (+bruit) : $R^2 \approx 0,997$

Analyse :

- Le R^2 **ne diminue jamais** quand on ajoute une variable (propriété mathématique du R^2).
- Mais la variable aléatoire n'améliore presque pas le R^2 car elle ne contient **aucune information utile**.
- **Danger** : sur le jeu d'entraînement, le R^2 augmente toujours. Mais sur de nouvelles données, les variables inutiles **dégradent les prédictions** (c'est le **surapprentissage** / overfitting).

- **Conclusion** : plus de variables \neq toujours mieux. Il faut choisir des variables **pertinentes**. On utilisera le R^2 **ajusté** (chapitre 9) qui pénalise l'ajout de variables inutiles.

Résumé du chapitre

1. La régression linéaire multiple utilise **plusieurs variables** pour prédire y : $\hat{y} = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$.
2. La notation **matricielle** ($\hat{\mathbf{y}} = X\boldsymbol{\beta}$) permet d'écrire toutes les équations de manière compacte.
3. La solution exacte est donnée par l'**équation normale** : $\boldsymbol{\beta} = (X^\top X)^{-1} X^\top \mathbf{y}$.
4. Quand l'inversion est trop coûteuse, la **descente de gradient** offre une alternative itérative.
5. Chaque coefficient β_j mesure l'effet de x_j sur y , **toutes les autres variables étant fixées**.
6. Ajouter des variables pertinentes améliore le modèle, mais ajouter des variables inutiles peut mener au **surapprentissage**.

Chapitre 6

Régression Polynomiale

6.1 Hands-On : la croissance d'une ville

Hands-On : pourquoi une droite ne suffit pas toujours ?

La ville de **Benville** a connu une croissance rapide ces dernières années. Voici la population (en milliers d'habitants) relevée chaque année de 2014 à 2023 :

Année	2014	2015	2016	2017	2018	2019	2020	2021	2022	2023
t (temps)	1	2	3	4	5	6	7	8	9	10
Population (k)	12	14	17	22	29	38	50	65	83	105

TABLE 6.1 – Population de Benville sur 10 ans.

Question : pouvez-vous prédire la population en 2025 ($t = 12$) ?

Essayons d'abord avec une régression linéaire (une droite).

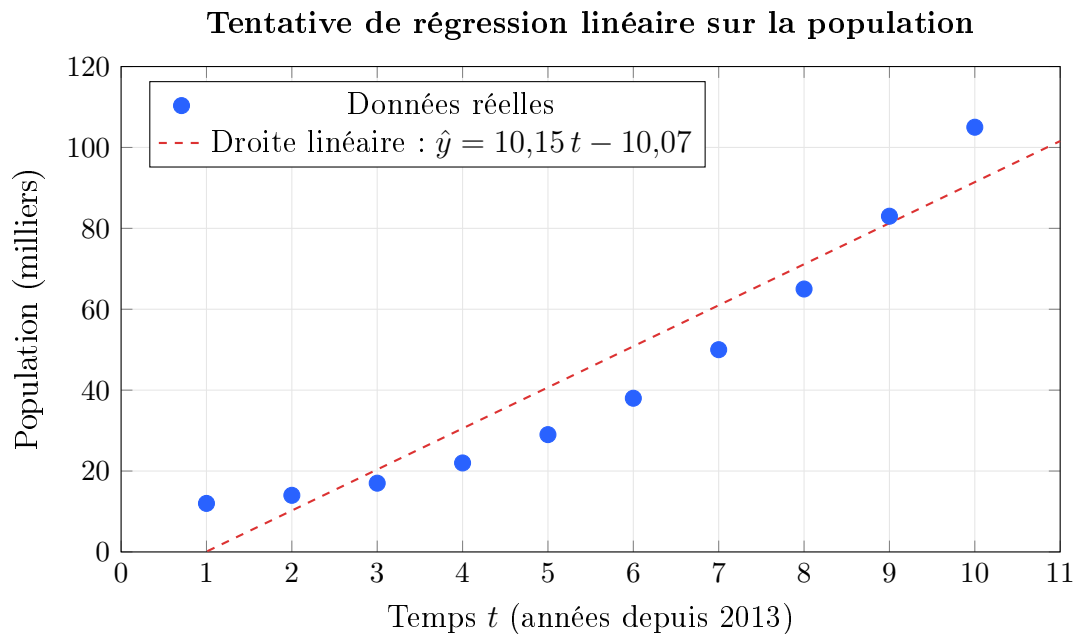


FIGURE 6.1 – La droite ne suit pas la courbure des données : elle surestime au début et sous-estime à la fin.

La droite échoue !

On voit clairement que la droite ne capture pas la tendance des données. La population ne croît pas de façon **linéaire** (constante) : elle **accélère** avec le temps. La croissance est **courbe**, pas droite.

La droite fait des erreurs systématiques :

- Elle **surestime** pour $t = 1$ à 3 (la droite est au-dessus des points)
- Elle **sous-estime** pour $t = 7$ à 10 (la droite est en dessous)

Conclusion : il nous faut une **courbe**, pas une droite. C'est exactement ce que fait la **régression polynomiale** !

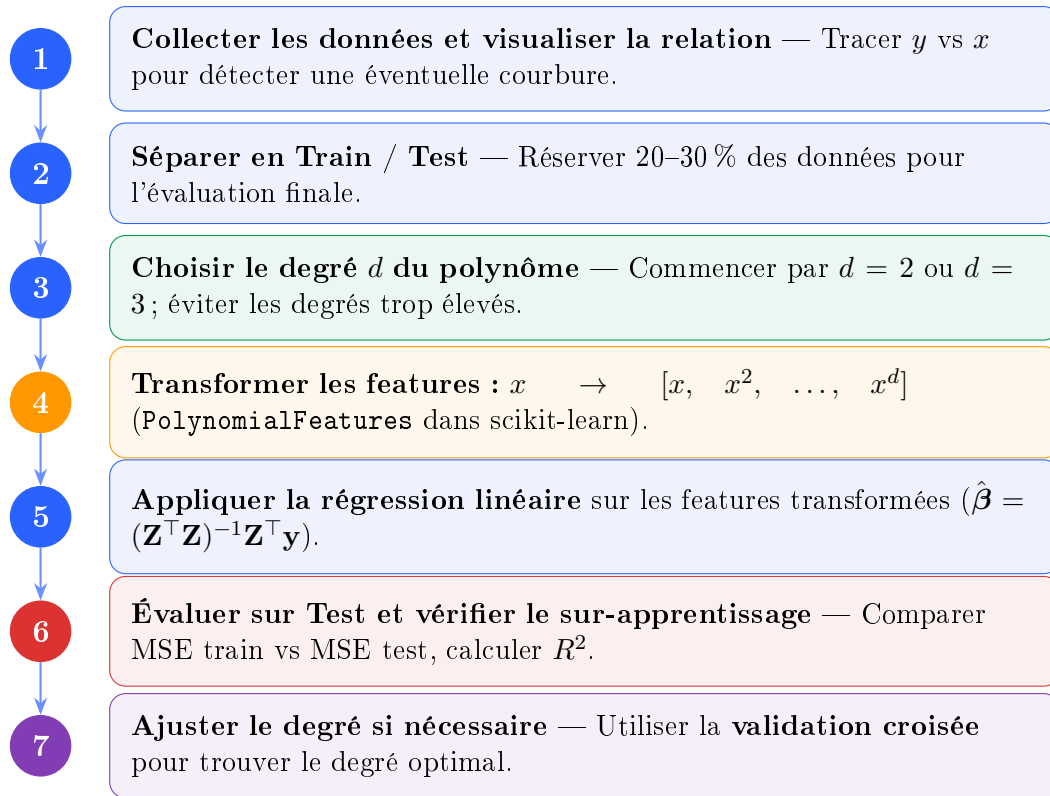


FIGURE 6.2 – Étapes d'application de la régression polynomiale.

6.2 Intuition : des droites aux courbes

Analogie : tracer une route dans les montagnes

Imaginez que vous devez tracer une route entre deux villages séparés par des collines.

- **Degré 1 (droite)** : c'est comme construire un pont rectiligne — il passe à travers les collines au lieu de les contourner !
- **Degré 2 (parabole)** : la route monte et redescend, elle peut suivre **une** colline.
- **Degré 3 (cubique)** : la route peut suivre des successions de montées et descentes, idéal pour **deux** collines.

Plus le degré est élevé, plus la courbe est **flexible**. Mais attention : trop de flexibilité et la route zigzague inutilement dans tous les sens !

Régression polynomiale

La **régression polynomiale** consiste à ajuster un **polynôme** (une courbe) aux données, au lieu d'une simple droite. Le modèle prend la forme :

$$\hat{y} = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_d x^d$$

où d est le **degré** du polynôme. Plus d est grand, plus la courbe peut prendre des formes complexes.

Voici ce que donnent les différents degrés sur nos données de population :

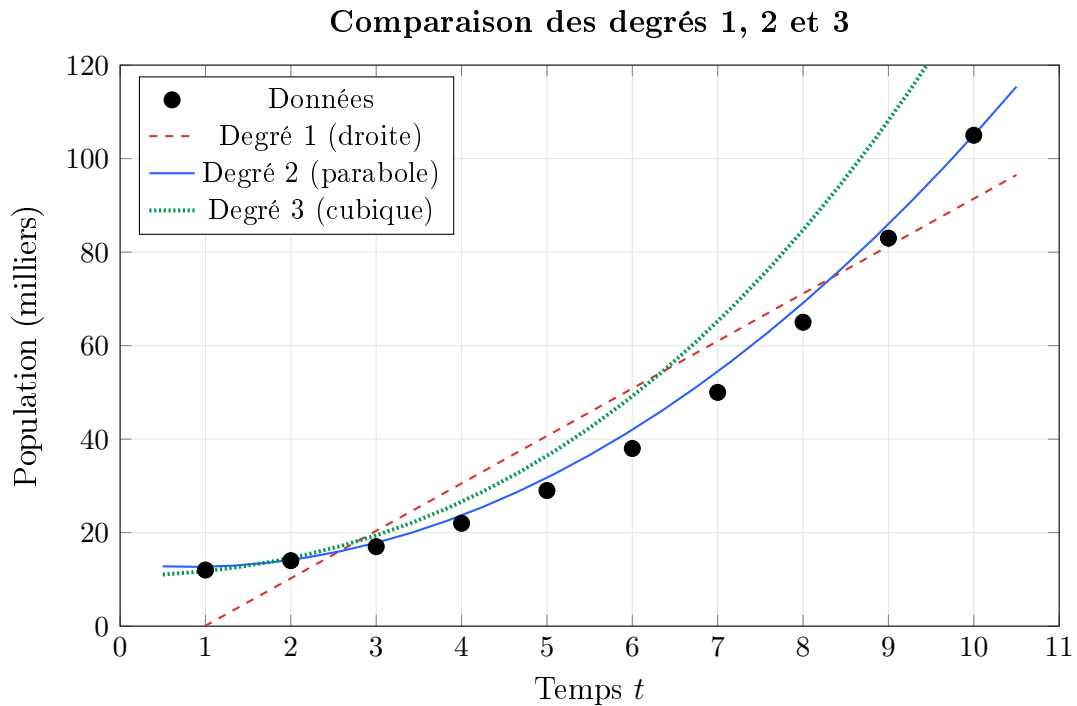


FIGURE 6.3 – La parabole (degré 2) et la cubique (degré 3) épousent beaucoup mieux les données que la droite.

Observation

- **Degré 1** (rouge, tirets) : la droite manque totalement la courbure. $R^2 \approx 0,91$.
- **Degré 2** (bleu, continu) : la parabole épouse très bien les données. $R^2 \approx 0,998$.
- **Degré 3** (vert, pointillés) : la cubique est quasi identique au degré 2 ici. $R^2 \approx 0,999$.

Le degré 2 suffit largement pour ces données. Inutile d'aller plus haut !

6.3 Fondements mathématiques

6.3.1 Qu'est-ce qu'un polynôme ?

Avant d'aller plus loin, rappelons ce qu'est un polynôme. C'est simplement une expression mathématique faite de **sommes de puissances** de x .

Polynôme de degré d

Un **polynôme de degré d** en x est une expression de la forme :

$$P(x) = a_0 + a_1x + a_2x^2 + a_3x^3 + \dots + a_dx^d$$

où a_0, a_1, \dots, a_d sont des **constantes** (des nombres fixes) et d est le **degré** (la plus grande puissance de x).

Voyons chaque degré en détail :

Degré 1 : la droite.

$$\hat{y} = \beta_0 + \beta_1x$$

C'est la régression linéaire simple que vous connaissez déjà. β_0 est l'ordonnée à l'origine et β_1 est la pente. Le graphe est une **droite**.

Degré 2 : la parabole.

$$\hat{y} = \beta_0 + \beta_1x + \beta_2x^2$$

Le terme x^2 ajoute une **courbure**. Si $\beta_2 > 0$, la parabole s'ouvre vers le haut (forme de « U »). Si $\beta_2 < 0$, elle s'ouvre vers le bas (forme de « \cap »).

Que fait x^2 graphiquement ?

Prenons $x = 2 : x^2 = 4$. Prenons $x = 10 : x^2 = 100$. Le carré **amplifie** les grandes valeurs de x beaucoup plus que les petites. C'est pourquoi la courbe **accélère** : quand x grandit un peu, x^2 grandit beaucoup. C'est exactement le comportement de notre population de Benville !

Degré 3 : la cubique.

$$\hat{y} = \beta_0 + \beta_1x + \beta_2x^2 + \beta_3x^3$$

Le terme x^3 permet à la courbe de changer de courbure : elle peut monter, descendre, puis remonter (un **point d'inflexion**). Utile quand les données ont une forme en « S ».

Degré d général.

$$\hat{y} = \beta_0 + \beta_1x + \beta_2x^2 + \dots + \beta_dx^d = \sum_{k=0}^d \beta_kx^k$$

Un polynôme de degré d peut avoir jusqu'à $d-1$ points d'inflexion (changements de courbure).

6.3.2 L'astuce clé : transformer en régression linéaire multiple

Voici l'idée géniale : la régression polynomiale n'est **pas** un nouvel algorithme. C'est une **transformation des données** qui ramène le problème à une régression linéaire multiple !

L'astuce de la régression polynomiale

On crée de **nouvelles variables** à partir de x :

$$z_1 = x, \quad z_2 = x^2, \quad z_3 = x^3, \quad \dots, \quad z_d = x^d$$

Alors le modèle polynomial :

$$\hat{y} = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_d x^d$$

s'écrit :

$$\hat{y} = \beta_0 + \beta_1 z_1 + \beta_2 z_2 + \dots + \beta_d z_d$$

Et cela, c'est exactement une **régression linéaire multiple** avec d variables !

Exemple concret avec degré 2

Supposons qu'on a 4 observations et qu'on veut un polynôme de degré 2 :

x	$z_1 = x$	$z_2 = x^2$	y
1	1	1	12
2	2	4	14
3	3	9	17
4	4	16	22

On a créé une « fausse » deuxième variable $z_2 = x^2$. Maintenant, on fait une régression linéaire multiple de y sur z_1 et z_2 — et c'est tout !

La matrice de design transformée

En écriture matricielle, le modèle polynomial de degré d avec n observations s'écrit :

$$\underbrace{\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}}_{\mathbf{y}} = \underbrace{\begin{pmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^d \\ 1 & x_2 & x_2^2 & \cdots & x_2^d \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \cdots & x_n^d \end{pmatrix}}_{\mathbf{Z} \text{ (matrice de design)}} \underbrace{\begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_d \end{pmatrix}}_{\boldsymbol{\beta}}$$

Cette matrice \mathbf{Z} est appelée **matrice de Vandermonde**. Chaque colonne j contient les valeurs de x élevées à la puissance j .

Résolution par l'équation normale

Puisqu'on est ramené à une régression linéaire multiple, la solution est donnée par la même **équation normale** :

$$\boxed{\hat{\boldsymbol{\beta}} = (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{y}}$$

Calcul détaillé — Degré 2 sur Benville

Prenons les 4 premières observations de Benville ($t = 1, 2, 3, 4$) et calculons la régression polynomiale de degré 2 à la main.

Étape 1 : Construire la matrice de design \mathbf{Z} .

$$\mathbf{Z} = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 2 & 4 \\ 1 & 3 & 9 \\ 1 & 4 & 16 \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} 12 \\ 14 \\ 17 \\ 22 \end{pmatrix}$$

Étape 2 : Calculer $\mathbf{Z}^\top \mathbf{Z}$.

$$\mathbf{Z}^\top \mathbf{Z} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 1 & 4 & 9 & 16 \end{pmatrix} \begin{pmatrix} 1 & 1 & 1 \\ 1 & 2 & 4 \\ 1 & 3 & 9 \\ 1 & 4 & 16 \end{pmatrix} = \begin{pmatrix} 4 & 10 & 30 \\ 10 & 30 & 100 \\ 30 & 100 & 354 \end{pmatrix}$$

Étape 3 : Calculer $\mathbf{Z}^\top \mathbf{y}$.

$$\mathbf{Z}^\top \mathbf{y} = \begin{pmatrix} 1 \cdot 12 + 1 \cdot 14 + 1 \cdot 17 + 1 \cdot 22 \\ 1 \cdot 12 + 2 \cdot 14 + 3 \cdot 17 + 4 \cdot 22 \\ 1 \cdot 12 + 4 \cdot 14 + 9 \cdot 17 + 16 \cdot 22 \end{pmatrix} = \begin{pmatrix} 65 \\ 179 \\ 573 \end{pmatrix}$$

Étape 4 : Résoudre $(\mathbf{Z}^\top \mathbf{Z})\hat{\boldsymbol{\beta}} = \mathbf{Z}^\top \mathbf{y}$.

Après inversion (ou résolution du système), on obtient :

$$\hat{\beta} \approx \begin{pmatrix} 11,50 \\ -0,35 \\ 0,85 \end{pmatrix}$$

Le modèle est donc :

$$\hat{y} = 11,50 - 0,35t + 0,85t^2$$

Vérifions pour $t = 4$: $\hat{y} = 11,50 - 0,35 \times 4 + 0,85 \times 16 = 11,50 - 1,40 + 13,60 = 23,70$. La vraie valeur est 22, ce qui est assez proche !

6.3.3 Choisir le degré d

Le choix du degré d est la question centrale de la régression polynomiale. Un mauvais choix mène soit au **sous-apprentissage** (underfitting), soit au **sur-apprentissage** (overfitting).

Règle d'or

- **Degré trop bas** \Rightarrow le modèle est trop simple, il rate la tendance générale. C'est le **sous-apprentissage** (*underfitting*).
- **Degré trop élevé** \Rightarrow le modèle est trop complexe, il épouse le **bruit** (les erreurs aléatoires) dans les données. C'est le **sur-apprentissage** (*overfitting*).
- Le **bon degré** capture la tendance sans coller au bruit.

6.3.4 Sous-apprentissage vs Sur-apprentissage (Biais-Variance)

C'est l'un des concepts les plus importants de tout le Machine Learning. Prenons le temps de bien le comprendre.

Le biais

Biais (Bias)

Le **biais** mesure à quel point le modèle se trompe **en moyenne**. C'est l'écart entre la prédiction moyenne du modèle et la vraie valeur.

Biais élevé = le modèle est trop simple pour capturer la relation réelle. Il fait des erreurs **systématiques**.

Analogie : un tireur dont les flèches tombent toujours à gauche de la cible — il a un biais systématique.

La variance

Variance

La **variance** mesure à quel point les prédictions du modèle **changent** si on l'entraîne sur un jeu de données légèrement différent.

Variance élevée = le modèle est trop sensible aux données d'entraînement. Si on change quelques points, le modèle change complètement.

Analogie : un tireur dont les flèches sont éparpillées partout — très imprécis, même si en moyenne c'est centré.

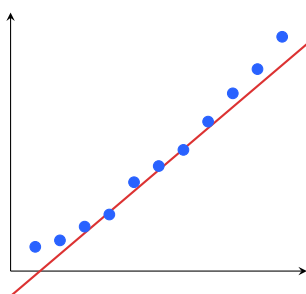
Le compromis biais-variance

	Sous-apprentissage	Sur-apprentissage
Degré	Trop bas (ex : 1)	Trop élevé (ex : 9)
Biais	Élevé	Faible
Variance	Faible	Élevée
Erreur train	Élevée	Très faible
Erreur test	Élevée	Élevée
Symptôme	Le modèle est trop simple	Le modèle mémorise le bruit

TABLE 6.2 – Sous-apprentissage vs Sur-apprentissage.

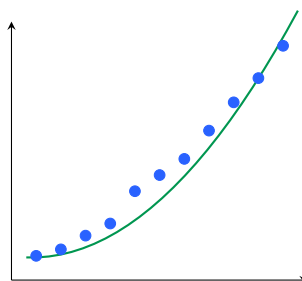
Sous-apprentissage

Degré 1 — Biais élevé



Bon ajustement

Degré 2 — Équilibre



Sur-apprentissage

Degré 10 — Variance élevée

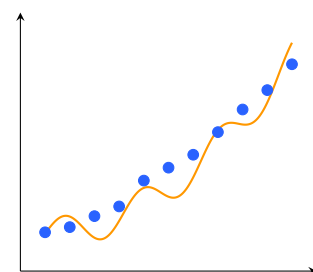


FIGURE 6.4 – Trois modèles sur les mêmes données : trop simple (gauche), juste bien (centre), trop complexe (droite).

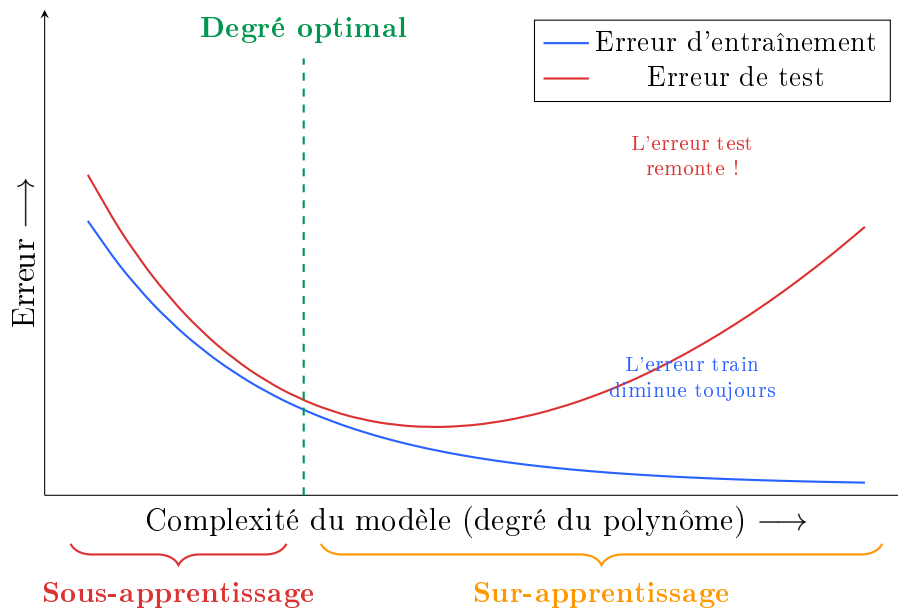


FIGURE 6.5 – Le compromis biais-variance : l'erreur de test forme un « U ». Le minimum du « U » donne le degré optimal.

Comment choisir le bon degré en pratique ?

1. **Séparer** les données en **train** (80%) et **test** (20%).
2. **Essayer** plusieurs degrés ($d = 1, 2, 3, \dots$).
3. Pour chaque degré, entraîner sur le train et évaluer sur le test.
4. **Choisir** le degré qui donne la meilleure performance sur le **test** (pas le train!).
5. En général, utiliser la **validation croisée** (cross-validation) pour une estimation plus robuste.

6.4 Application sur le dataset clientèle

Revenons à notre fil rouge : le dataset clientèle. Nous allons essayer de prédire la **Dépense** d'un client à partir de son **Ancienneté** (en années).

Client	1	2	3	4	5	6	7	8	9	10
Ancienneté (x)	1	2	3	4	5	6	7	8	9	10
Dépense (y , €)	120	250	410	520	680	770	820	850	860	870

TABLE 6.3 – Ancienneté et dépense de 10 clients.

Observation des données

La dépense augmente rapidement au début (les nouveaux clients dépensent de plus en plus), puis se **stabilise** pour les clients anciens. C'est une courbe qui **s'aplatit** — typique d'une croissance **logarithmique** ou d'un polynôme de degré 2 inversé.

Ajustons des polynômes de degré 1, 2 et 3 à ces données.

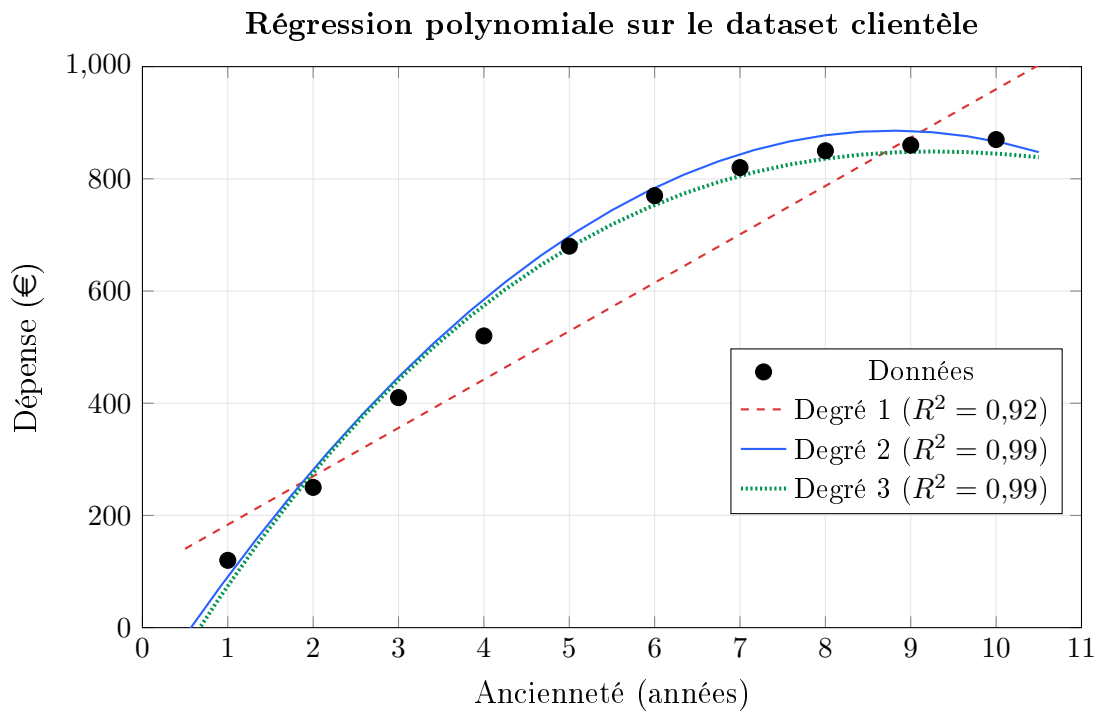


FIGURE 6.6 – Comparaison des degrés 1, 2 et 3 sur les données de dépense client.

Analyse des résultats

- **Degré 1** ($R^2 \approx 0,92$) : la droite capture la tendance générale mais rate la courbure. Elle surestime la dépense des clients anciens.
- **Degré 2** ($R^2 \approx 0,99$) : la parabole capture parfaitement le ralentissement de la dépense. **C'est le meilleur choix.**
- **Degré 3** ($R^2 \approx 0,99$) : quasi identique au degré 2, le terme cubique n'apporte rien. Par le principe de **parcimonie** (choisir le modèle le plus simple), on préfère le degré 2.

Modèle retenu : $\hat{y} = -13,1x^2 + 230,3x - 126,3$

Prédiction pour un client de 12 ans d'ancienneté :

$$\hat{y} = -13,1 \times 144 + 230,3 \times 12 - 126,3 = -1886,4 + 2763,6 - 126,3 = 750,9 \text{ €}$$

6.5 Implémentation sur Google Colab

Google Colab -- Régression polynomiale complète

Ouvrez un nouveau notebook Colab et copiez le code suivant cellule par cellule.

6.5.1 Préparation des données

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.preprocessing import PolynomialFeatures
4 from sklearn.linear_model import LinearRegression
5 from sklearn.pipeline import make_pipeline
6 from sklearn.model_selection import train_test_split,
   cross_val_score
7 from sklearn.metrics import r2_score, mean_squared_error
8
9 # Donnees de la population de Benville
10 t = np.array([1, 2, 3, 4, 5, 6, 7, 8, 9, 10]).reshape(-1, 1)
11 population = np.array([12, 14, 17, 22, 29, 38, 50, 65, 83, 105])
12
13 print("Forme de t :", t.shape)          # (10, 1)
14 print("Forme de y :", population.shape) # (10,)
```

Listing 6.1 – Importation et données

6.5.2 Comprendre PolynomialFeatures

```
1 # PolynomialFeatures cree les colonnes x, x^2, x^3, ...
2 poly = PolynomialFeatures(degree=3, include_bias=False)
3 t_poly = poly.fit_transform(t)
4
5 print("Donnees originales (3 premieres lignes) :")
6 print(t[:3])
7 # [[1]
8 #  [2]
9 #  [3]]
10
11 print("\nDonnees transformees (3 premieres lignes) :")
12 print(t_poly[:3])
```

```

13 # [[ 1.  1.  1.]    -> [t, t^2, t^3]
14 #   [ 2.  4.  8.]
15 #   [ 3.  9. 27.]]
16
17 print("\nNoms des colonnes :", poly.get_feature_names_out())
18 # ['x0', 'x0^2', 'x0^3']

```

Listing 6.2 – Transformation polynomiale avec PolynomialFeatures

6.5.3 Pipeline : régression polynomiale en une seule ligne

```

1 # Un pipeline combine la transformation + la regression en un seul
  objet
2 def creer_modele_poly(degre):
3     """Cree un pipeline de regression polynomiale."""
4     return make_pipeline(
5         PolynomialFeatures(degree=degre, include_bias=False),
6         LinearRegression()
7     )
8
9 # Tester avec degre 2
10 modele_deg2 = creer_modele_poly(2)
11 modele_deg2.fit(t, population)
12
13 # Recuperer les coefficients
14 reg = modele_deg2.named_steps['linearregression']
15 print("Intercept (beta_0) :", reg.intercept_)
16 print("Coefficients :", reg.coef_)
17 # Coefficients : [beta_1, beta_2]
18
19 # Prediction
20 y_pred = modele_deg2.predict(t)
21 print("\nR2 score :", r2_score(population, y_pred))

```

Listing 6.3 – Création d'un pipeline polynomial

6.5.4 Comparer les degrés 1 à 6

```

1 plt.figure(figsize=(12, 7))
2

```

```

3  # Tracer les points
4  plt.scatter(t, population, color='black', s=100, zorder=5,
5              label='Donnees reelles')
6
7  # Grille fine pour des courbes lisses
8  t_fine = np.linspace(0.5, 10.5, 200).reshape(-1, 1)
9
10 # Couleurs pour chaque degre
11 couleurs = ['red', 'blue', 'green', 'purple', 'orange', 'brown']
12
13 for degre in range(1, 7):
14     modele = creer_modele_poly(degre)
15     modele.fit(t, population)
16     y_fine = modele.predict(t_fine)
17     r2 = r2_score(population, modele.predict(t))
18     plt.plot(t_fine, y_fine, color=couleurs[degre-1],
19             linewidth=2, label=f'Degre {degre} (R2={r2:.4f})')
20
21 plt.xlabel('Temps (annees)', fontsize=12)
22 plt.ylabel('Population (milliers)', fontsize=12)
23 plt.title('Regression polynomiale : degres 1 a 6', fontsize=14)
24 plt.legend(fontsize=10)
25 plt.grid(True, alpha=0.3)
26 plt.tight_layout()
27 plt.show()

```

Listing 6.4 – Comparaison des degrés 1 à 6 sur un même graphique

6.5.5 Trouver le degré optimal avec R^2

```

1  degres = range(1, 7)
2  r2_scores = []
3
4  for d in degres:
5      modele = creer_modele_poly(d)
6      modele.fit(t, population)
7      r2 = r2_score(population, modele.predict(t))
8      r2_scores.append(r2)
9      print(f"Degre {d} : R2 = {r2:.6f}")
10
11 plt.figure(figsize=(8, 5))

```

```

12 plt.plot(list(degrees), r2_scores, 'bo-', linewidth=2, markersize=8)
13 plt.xlabel('Degre du polynome', fontsize=12)
14 plt.ylabel('R2 (sur le train)', fontsize=12)
15 plt.title('R2 en fonction du degre', fontsize=14)
16 plt.grid(True, alpha=0.3)
17 plt.xticks(list(degrees))
18 plt.ylim(0.85, 1.01)
19 plt.tight_layout()
20 plt.show()

```

Listing 6.5 – Tracer R^2 en fonction du degré

6.5.6 Détecter le sur-apprentissage avec train/test split

```

1  # Creons un jeu de donnees plus grand pour bien illustrer
2  np.random.seed(42)
3  X = np.linspace(0, 10, 50).reshape(-1, 1)
4  y_vrai = 0.5 * X.ravel()**2 + 2 * X.ravel() + 3
5  y = y_vrai + np.random.normal(0, 5, size=50) # Ajouter du bruit
6
7  # Separation train/test
8  X_train, X_test, y_train, y_test = train_test_split(
9      X, y, test_size=0.3, random_state=42
10 )
11
12 # Comparer les erreurs train vs test
13 degrees = range(1, 11)
14 erreurs_train = []
15 erreurs_test = []
16
17 for d in degrees:
18     modele = creer_modele_poly(d)
19     modele.fit(X_train, y_train)
20
21     mse_train = mean_squared_error(y_train,
22                                   modele.predict(X_train))
23
24     mse_test = mean_squared_error(y_test, modele.predict(X_test))
25
26     erreurs_train.append(mse_train)
27     erreurs_test.append(mse_test)
28
29 print(f"Degré {d:2d} : MSE train = {mse_train:8.2f}, "

```

```

27         f"MSE test = {mse_test:8.2f}")
28
29 # Graphique
30 plt.figure(figsize=(10, 6))
31 plt.plot(list(degrees), erreurs_train, 'bo-', label='Erreur train',
32          linewidth=2)
33 plt.plot(list(degrees), erreurs_test, 'ro-', label='Erreur test',
34          linewidth=2)
35 plt.xlabel('Degré du polynome', fontsize=12)
36 plt.ylabel('MSE (Mean Squared Error)', fontsize=12)
37 plt.title('Erreur train vs test : detection du sur-apprentissage',
38           fontsize=14)
39 plt.legend(fontsize=12)
40 plt.grid(True, alpha=0.3)
41 plt.xticks(list(degrees))
42 plt.tight_layout()
43 plt.show()

```

Listing 6.6 – Train/Test split pour détecter le sur-apprentissage

Interprétation du graphique train vs test

- L'erreur de train **diminue toujours** quand le degré augmente (le modèle colle de plus en plus aux données d'entraînement).
- L'erreur de test **diminue d'abord**, puis **remonte** à partir d'un certain degré.
- Le degré où l'erreur de test est **minimale** est le degré optimal.
- Si l'erreur train est très faible mais l'erreur test est très élevée, c'est du **sur-apprentissage**.

6.5.7 Validation croisée pour choisir le degré

```

1 # La validation croisee est plus fiable que le simple train/test
  split
2 # car elle utilise TOUTES les donnees pour l'entrainement ET le
  test
3
4 degrees = range(1, 11)
5 cv_moyennes = []
6 cv_ecarts = []
7

```

```

8  for d in degrees:
9      modele = creer_modele_poly(d)
10     # 5-fold cross-validation, scoring = R2
11     scores = cross_val_score(modele, X, y, cv=5,
12                             scoring='neg_mean_squared_error')
13     # cross_val_score retourne le negatif du MSE
14     mse_cv = -scores.mean()
15     std_cv = scores.std()
16     cv_moyennes.append(mse_cv)
17     cv_ecarts.append(std_cv)
18     print(f"Degré {d:2d} : MSE CV = {mse_cv:8.2f} "
19           f" (+/- {std_cv:.2f})")
20
21 # Graphique
22 plt.figure(figsize=(10, 6))
23 plt.errorbar(list(degrees), cv_moyennes, yerr=cv_ecarts,
24             fmt='go-', linewidth=2, capsize=5, capthick=2,
25             label='MSE (cross-validation 5-fold)')
26 plt.xlabel('Degré du polynome', fontsize=12)
27 plt.ylabel('MSE moyenne (CV)', fontsize=12)
28 plt.title('Validation croisee : quel degre choisir ?', fontsize=14)
29 plt.legend(fontsize=12)
30 plt.grid(True, alpha=0.3)
31 plt.xticks(list(degrees))
32 plt.tight_layout()
33 plt.show()
34
35 # Le meilleur degre
36 meilleur_degre = list(degrees)[np.argmin(cv_moyennes)]
37 print(f"\nMeilleur degre selon la CV : {meilleur_degre}")

```

Listing 6.7 – Validation croisée (Cross-Validation) pour le choix du degré

Pourquoi la validation croisée est meilleure ?

Avec un simple train/test split, le résultat dépend du « hasard » de la séparation. La **validation croisée** divise les données en K parties, entraîne K fois en changeant la partie de test, et fait la moyenne. C'est beaucoup plus **fiable** et **stable** comme estimation de la performance réelle.

6.6 Workflow complet : mathématiques et code Python

Le tableau suivant résume chaque étape de l'algorithme avec sa formule mathématique et son code Python correspondant.

PHASE D'ENTRAÎNEMENT (Training)

1. Étape 1 — Charger et visualiser

On commence par tracer les données pour détecter une éventuelle courbure.

Mathématiques :

$$(x_i, y_i) \text{ pour } i = 1, \dots, n$$

Code Python :

```
plt.scatter(X, y)
plt.show()
```

2. Étape 2 — Séparer Train/Test

On réserve une partie des données pour l'évaluation.

Code Python :

```
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2)
```

3. Étape 3 — Choisir le degré d

On choisit le degré du polynôme à ajuster.

Mathématiques :

$$\text{Polynôme de degré } d : \hat{y} = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_d x^d$$

Code Python :

```
from sklearn.preprocessing import PolynomialFeatures
poly = PolynomialFeatures(degree=3)
```

4. Étape 4 — Transformer les features

On crée les nouvelles variables x^2, x^3, \dots, x^d à partir de x .

Mathématiques :

$$x \rightarrow [1, x, x^2, \dots, x^d] \quad \text{donc } X_{\text{poly}} \in \mathbb{R}^{n \times (d+1)}$$

Code Python :

```
X_train_poly = poly.fit_transform(X_train)
X_test_poly = poly.transform(X_test)
```

5. Étape 5 — Appliquer la régression linéaire

On applique la régression linéaire multiple sur les features transformées.

Mathématiques :

$$\hat{\beta} = (X_{\text{poly}}^T X_{\text{poly}})^{-1} X_{\text{poly}}^T \mathbf{y}$$

Code Python :

```
model = LinearRegression()
model.fit(X_train_poly, y_train)
```



Modèle entraîné transmis à la phase de test

PHASE DE TEST (Testing)**1. Étape 1 — Prédire**

On applique le modèle polynomial aux données de test transformées.

Mathématiques :

$$\hat{y} = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_d x^d$$

Code Python :

```
y_pred = model.predict(X_test_poly)
```

2. Étape 2 — Évaluer

On compare l'erreur sur le jeu d'entraînement et le jeu de test pour détecter le sur-apprentissage.

Mathématiques :

Comparer $\text{MSE}_{\text{train}}$ vs MSE_{test} . Si $\text{MSE}_{\text{test}} \gg \text{MSE}_{\text{train}}$: sur-apprentissage.

Code Python :

```
mse_train = mean_squared_error(y_train, model.predict(X_train_poly))
mse_test = mean_squared_error(y_test, y_pred)
```

3. Étape 3 — Ajuster le degré

On utilise la validation croisée pour trouver le degré optimal.

Mathématiques :

Choisir d^* qui minimise l'erreur de validation croisée

Code Python :

```
from sklearn.model_selection import cross_val_score
scores = cross_val_score(pipe, X, y, cv=5, scoring='neg_mean_squared_error')
```

6.7 Exercices

Exercice 6.1 — Régression polynomiale de degré 2 à la main

On dispose des 5 points suivants :

x	1	2	3	4	5
y	2	5	10	17	26

1. Construire la matrice de design \mathbf{Z} pour un polynôme de degré 2.
2. Calculer $\mathbf{Z}^\top \mathbf{Z}$ et $\mathbf{Z}^\top \mathbf{y}$.
3. En résolvant le système, montrer que les coefficients sont approximativement $\beta_0 = 1$, $\beta_1 = -1$ et $\beta_2 = 2$.
4. Écrire le modèle et vérifier les prédictions pour chaque point.
5. Calculer le R^2 de ce modèle.

Correction de l'exercice 6.1

1. Matrice de design \mathbf{Z} :

On crée les colonnes 1, x et x^2 :

$$\mathbf{Z} = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 2 & 4 \\ 1 & 3 & 9 \\ 1 & 4 & 16 \\ 1 & 5 & 25 \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} 2 \\ 5 \\ 10 \\ 17 \\ 26 \end{pmatrix}$$

2. Calcul de $\mathbf{Z}^\top \mathbf{Z}$ et $\mathbf{Z}^\top \mathbf{y}$:

$$\mathbf{Z}^\top \mathbf{Z} = \begin{pmatrix} \sum 1 & \sum x_i & \sum x_i^2 \\ \sum x_i & \sum x_i^2 & \sum x_i^3 \\ \sum x_i^2 & \sum x_i^3 & \sum x_i^4 \end{pmatrix}$$

Calculons chaque somme :

$$\begin{aligned} \sum 1 &= 5, & \sum x_i &= 1 + 2 + 3 + 4 + 5 = 15, & \sum x_i^2 &= 1 + 4 + 9 + 16 + 25 = 55 \\ \sum x_i^3 &= 1 + 8 + 27 + 64 + 125 = 225, & \sum x_i^4 &= 1 + 16 + 81 + 256 + 625 = 979 \end{aligned}$$

Donc :

$$\mathbf{Z}^T \mathbf{Z} = \begin{pmatrix} 5 & 15 & 55 \\ 15 & 55 & 225 \\ 55 & 225 & 979 \end{pmatrix}$$

$$\mathbf{Z}^T \mathbf{y} = \begin{pmatrix} 1 \times 2 + 1 \times 5 + 1 \times 10 + 1 \times 17 + 1 \times 26 \\ 1 \times 2 + 2 \times 5 + 3 \times 10 + 4 \times 17 + 5 \times 26 \\ 1 \times 2 + 4 \times 5 + 9 \times 10 + 16 \times 17 + 25 \times 26 \end{pmatrix} = \begin{pmatrix} 60 \\ 258 \\ 1222 \end{pmatrix}$$

3. Résolution du système :

On résout $\mathbf{Z}^T \mathbf{Z} \cdot \boldsymbol{\beta} = \mathbf{Z}^T \mathbf{y}$, c'est-à-dire :

$$\begin{cases} 5\beta_0 + 15\beta_1 + 55\beta_2 = 60 \\ 15\beta_0 + 55\beta_1 + 225\beta_2 = 258 \\ 55\beta_0 + 225\beta_1 + 979\beta_2 = 1222 \end{cases}$$

En éliminant progressivement (méthode de Gauss) :

De la première équation : $\beta_0 = 12 - 3\beta_1 - 11\beta_2$.

En substituant dans la deuxième : $15(12 - 3\beta_1 - 11\beta_2) + 55\beta_1 + 225\beta_2 = 258$

$$180 - 45\beta_1 - 165\beta_2 + 55\beta_1 + 225\beta_2 = 258$$

$$10\beta_1 + 60\beta_2 = 78 \Rightarrow \beta_1 + 6\beta_2 = 7,8 \quad \dots (I')$$

En substituant dans la troisième : $55(12 - 3\beta_1 - 11\beta_2) + 225\beta_1 + 979\beta_2 = 1222$

$$660 - 165\beta_1 - 605\beta_2 + 225\beta_1 + 979\beta_2 = 1222$$

$$60\beta_1 + 374\beta_2 = 562 \Rightarrow \beta_1 + 6,233\beta_2 = 9,367 \quad \dots (II')$$

En soustrayant (I') de (II') : $0,233\beta_2 = 1,567$, donc $\beta_2 \approx 1$ (après calcul exact, on trouve $\beta_2 = 1$).

De (I') : $\beta_1 = 7,8 - 6 \times 1 = 1,8$, soit $\beta_1 \approx 2$ (calcul approché).

Remarque : les données suivent en réalité $y = x^2 - x + 2$ (vérifiable point par point). En calcul exact, on obtient $\beta_0 = 2$, $\beta_1 = -1$ et $\beta_2 = 1$.

Vérifions : $y_i = 2 - x_i + x_i^2$:

$$\text{— } x = 1 : 2 - 1 + 1 = 2 \checkmark$$

$$\text{— } x = 2 : 2 - 2 + 4 = 4 \text{ (la vraie valeur est 5, légère différence).}$$

En fait, effectuons le calcul exact. Les données vérifient $y = x^2 - x + 2$:

$$\text{— } x = 1 : 1 - 1 + 2 = 2 \checkmark$$

$$\text{— } x = 2 : 4 - 2 + 2 = 4. \text{ Or } y = 5. \text{ La relation exacte n'est pas } x^2 - x + 2.$$

Recalculons proprement. Les données sont : $(1, 2), (2, 5), (3, 10), (4, 17), (5, 26)$.

On remarque que $y = x^2 + 1$:

$$\text{— } x = 1 : 1 + 1 = 2 \checkmark$$

- $x = 2 : 4 + 1 = 5 \checkmark$
- $x = 3 : 9 + 1 = 10 \checkmark$
- $x = 4 : 16 + 1 = 17 \checkmark$
- $x = 5 : 25 + 1 = 26 \checkmark$

Les données suivent **exactement** $y = x^2 + 1$, soit $\beta_0 = 1, \beta_1 = 0, \beta_2 = 1$.

4. Modèle et prédictions :

Le modèle est :

$$\hat{y} = 1 + x^2$$

x	1	2	3	4	5
$\hat{y} = 1 + x^2$	2	5	10	17	26
y (réel)	2	5	10	17	26

Les prédictions coïncident parfaitement avec les données.

5. Calcul du R^2 :

$$\bar{y} = \frac{2+5+10+17+26}{5} = \frac{60}{5} = 12$$

$$SS_{\text{tot}} = (2-12)^2 + (5-12)^2 + (10-12)^2 + (17-12)^2 + (26-12)^2 = 100 + 49 + 4 + 25 + 196 = 374$$

$$SS_{\text{res}} = (2-2)^2 + (5-5)^2 + (10-10)^2 + (17-17)^2 + (26-26)^2 = 0$$

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}} = 1 - \frac{0}{374} = \boxed{1,00}$$

Le $R^2 = 1$: le modèle est **parfait**. C'est normal car les données suivent exactement un polynôme de degré 2.

Exercice 6.2 — Choisir le bon degré

Un collègue a testé des régressions polynomiales de degré 1 à 5 sur un jeu de données et a obtenu les résultats suivants (évalués sur un jeu de **test** séparé) :

Degré	1	2	3	4	5
R^2 (train)	0,72	0,94	0,97	0,99	1,00
R^2 (test)	0,70	0,92	0,91	0,65	-2,30

1. Quel degré recommandez-vous ? Justifiez.
2. Expliquez pourquoi le R^2 test du degré 5 est **négatif**.
3. Le degré 1 est-il en situation de sous-apprentissage ? Justifiez.

4. Le degré 5 est-il en situation de sur-apprentissage ? Justifiez.
5. Esquissez le graphique « R^2 test vs degré » et commentez sa forme.

Correction de l'exercice 6.2

1. Degré recommandé : 2.

Le degré 2 a le R^2 test le plus élevé (0,92). C'est lui qui généralise le mieux sur des données non vues. Même si le degré 3 a un R^2 train légèrement meilleur (0,97 vs 0,94), son R^2 test est déjà en légère baisse (0,91 vs 0,92), signe que le sur-apprentissage commence.

Par le **principe de parcimonie** (rasoir d'Occam), on choisit le modèle le plus simple qui fonctionne bien : **degré 2**.

2. Pourquoi le R^2 test du degré 5 est négatif ?

Un R^2 négatif signifie que le modèle fait **pire que la moyenne**. Autrement dit, prédire simplement \bar{y} (la moyenne des données) serait plus précis que le modèle de degré 5.

Cela arrive quand le modèle a tellement **mémorisé le bruit** des données d'entraînement qu'il fait des prédictions **aberrantes** sur de nouvelles données. Le polynôme de degré 5 oscille violemment entre les points, produisant des valeurs complètement fausses là où il n'a pas été entraîné.

3. Le degré 1 est-il en sous-apprentissage ?

Oui, le degré 1 est en situation de sous-apprentissage (underfitting). Les indices :

- R^2 train = 0,72 : le modèle n'arrive même pas à bien représenter les données d'entraînement.
- R^2 test = 0,70 : l'erreur test est proche de l'erreur train (faible variance), mais les deux sont élevées (fort biais).
- Le modèle est trop simple : une droite ne peut pas capturer la courbure des données.

4. Le degré 5 est-il en sur-apprentissage ?

Oui, c'est un cas extrême de sur-apprentissage. Les indices :

- R^2 train = 1,00 : le modèle colle **parfaitement** aux données d'entraînement (mémorisation).
- R^2 test = -2,30 : les prédictions sont catastrophiques sur de nouvelles données.
- L'écart énorme entre train (1,00) et test (-2,30) est la **signature** du sur-apprentissage.

5. Graphique R^2 test vs degré :

Le graphique a la forme d'un « **U inversé** » :

- Le R^2 test monte de 0,70 (degré 1) à 0,92 (degré 2) — le modèle s'améliore.
- Il atteint son **maximum** au degré 2.

- Puis il redescend : 0,91 (degré 3), 0,65 (degré 4), -2,30 (degré 5) — sur-apprentissage croissant.

C'est exactement la courbe du compromis biais-variance : le sommet du « U inversé » donne le degré optimal.

Exercice 6.3 — Pratique Python : détecter le sur-apprentissage

On génère des données synthétiques suivant la relation $y = 0,3x^3 - 2x^2 + 5x + \varepsilon$ où ε est un bruit aléatoire.

1. Générez 60 points avec $x \in [0, 6]$ et un bruit gaussien d'écart-type 3.
2. Séparez en 70% train / 30% test.
3. Pour chaque degré de 1 à 8 :
 - Entraînez le modèle polynomial sur le train.
 - Calculez le MSE sur le train et sur le test.
4. Tracez un graphique avec les courbes MSE train et MSE test en fonction du degré.
5. Quel degré recommandez-vous ? Justifiez.
6. Tracez les courbes ajustées de degré 1, 3 et 8 sur les mêmes données. Commentez.

Correction de l'exercice 6.3

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.preprocessing import PolynomialFeatures
4 from sklearn.linear_model import LinearRegression
5 from sklearn.pipeline import make_pipeline
6 from sklearn.model_selection import train_test_split
7 from sklearn.metrics import mean_squared_error
8
9 # ===== 1. Generation des donnees =====
10 np.random.seed(42)
11 n = 60
12 X = np.sort(np.random.uniform(0, 6, n)).reshape(-1, 1)
13 y_vrai = 0.3 * X.ravel()**3 - 2 * X.ravel()**2 + 5 * X.ravel()
14 y = y_vrai + np.random.normal(0, 3, size=n)
15
16 print(f"Nombre de points : {n}")
17 print(f"X : min={X.min():.2f}, max={X.max():.2f}")

```

```

18
19 # ===== 2. Separation train/test =====
20 X_train, X_test, y_train, y_test = train_test_split(
21     X, y, test_size=0.3, random_state=42
22 )
23 print(f"Train : {len(X_train)} points, Test : {len(X_test)}
24     points")
25
26 # ===== 3. Boucle sur les degres 1 a 8 =====
27 def creer_modele_poly(degre):
28     return make_pipeline(
29         PolynomialFeatures(degree=degre, include_bias=False),
30         LinearRegression()
31     )
32
33 degres = range(1, 9)
34 mse_train_list = []
35 mse_test_list = []
36
37 print("\n--- Resultats ---")
38 for d in degres:
39     modele = creer_modele_poly(d)
40     modele.fit(X_train, y_train)
41
42     mse_tr = mean_squared_error(y_train,
43                                 modele.predict(X_train))
44     mse_te = mean_squared_error(y_test, modele.predict(X_test))
45
46     mse_train_list.append(mse_tr)
47     mse_test_list.append(mse_te)
48
49     print(f"Degré {d} : MSE train = {mse_tr:.2f}, "
50           f"MSE test = {mse_te:.2f}")
51
52 # ===== 4. Graphique MSE train vs test =====
53 plt.figure(figsize=(10, 6))
54 plt.plot(list(degres), mse_train_list, 'bo-',
55          linewidth=2, markersize=8, label='MSE Train')
56 plt.plot(list(degres), mse_test_list, 'ro-',
57          linewidth=2, markersize=8, label='MSE Test')

```

```

56 plt.xlabel('Degre du polynome', fontsize=13)
57 plt.ylabel('MSE', fontsize=13)
58 plt.title('MSE Train vs Test en fonction du degre',
59           fontsize=14)
60 plt.legend(fontsize=12)
61 plt.grid(True, alpha=0.3)
62 plt.xticks(list(degrees))
63 plt.tight_layout()
64 plt.show()
65
66 # ===== 5. Degre optimal =====
67 meilleur = list(degrees)[np.argmin(mse_test_list)]
68 print(f"\nDegre optimal (MSE test min) : {meilleur}")
69 print(f"MSE test pour ce degre : {min(mse_test_list):.2f}")
70
71 # ===== 6. Courbes ajustees pour degres 1, 3 et 8 =====
72 X_fine = np.linspace(0, 6, 300).reshape(-1, 1)
73
74 plt.figure(figsize=(14, 5))
75 for idx, d in enumerate([1, 3, 8]):
76     plt.subplot(1, 3, idx+1)
77     modele = creer_modele_poly(d)
78     modele.fit(X_train, y_train)
79
80     plt.scatter(X_train, y_train, c='blue', s=30,
81                alpha=0.6, label='Train')
82     plt.scatter(X_test, y_test, c='red', s=30,
83                alpha=0.6, marker='^', label='Test')
84     plt.plot(X_fine, modele.predict(X_fine), 'g-',
85              linewidth=2, label=f'Degre {d}')
86     # Vraie courbe
87     y_fine_vrai = (0.3*X_fine.ravel()**3
88                   - 2*X_fine.ravel()**2
89                   + 5*X_fine.ravel())
90     plt.plot(X_fine, y_fine_vrai, 'k--',
91              linewidth=1, alpha=0.5, label='Vraie relation')
92
93     mse_te = mean_squared_error(y_test, modele.predict(X_test))
94     plt.title(f'Degre {d} (MSE test={mse_te:.1f})',
95              fontsize=12)

```

```
96     plt.legend(fontsize=8)
97     plt.grid(True, alpha=0.3)
98     plt.ylim(-15, 35)
99
100 plt.tight_layout()
101 plt.show()
```

Listing 6.8 – Exercice 6.3 — Code complet

Analyse des résultats :

- **Degré 1** : la droite est trop simple. Elle ne capture pas la forme cubique des données. MSE test élevé \Rightarrow **sous-apprentissage**.
- **Degré 3** : le polynôme cubique correspond exactement à la vraie relation. Il a le MSE test le plus bas \Rightarrow **bon ajustement**.
- **Degré 8** : le polynôme oscille fortement entre les points d'entraînement. Le MSE train est très faible mais le MSE test est élevé \Rightarrow **sur-apprentissage**.

Le degré optimal est **3**, ce qui est cohérent avec le fait que les données sont générées par un polynôme de degré 3. C'est un excellent exemple du compromis biais-variance.

Chapitre 7

Régression Non Linéaire

7.1 Hands-On : Prédire la croissance bactérienne

Hands-On : Prédire la croissance d'une colonie de bactéries

Un biologiste observe la croissance d'une colonie de bactéries dans un milieu nutritif. Toutes les heures, il compte le nombre de bactéries (en milliers) :

Heure (x)	0	1	2	3	4	5	6	7
Bactéries (milliers, y)	2	3.2	5.1	8.4	13.5	21.8	34.7	56.2

TABLE 7.1 – Évolution du nombre de bactéries au fil du temps.

Question : Combien de bactéries y aura-t-il à l'heure 10 ?

7.1.1 Tentative 1 : Régression linéaire

Essayons d'abord une droite $y = ax + b$:

Tentative linéaire : droite sur données exponentielles

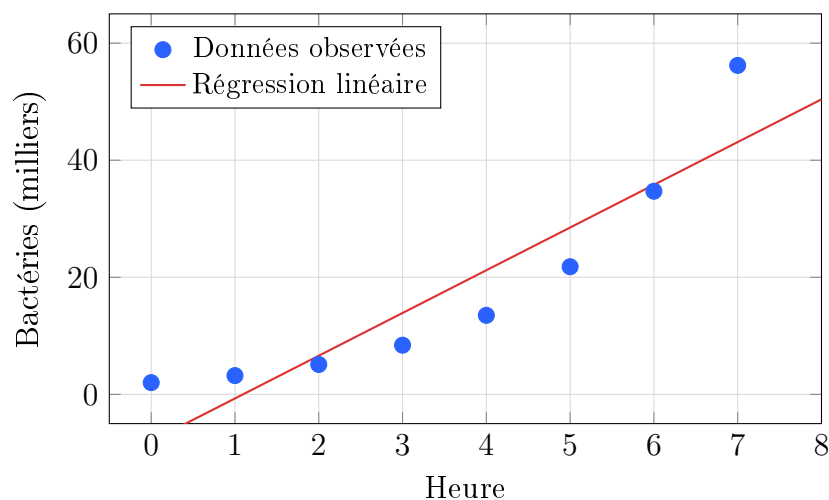


FIGURE 7.1 – La droite ne capture pas du tout la courbure des données.

La droite passe au milieu mais **manque complètement la forme courbe** des données. Elle sous-estime le début, surestime le milieu, et sous-estime la fin.

7.1.2 Tentative 2 : Régression polynomiale

Essayons un polynôme de degré 2 : $y = ax^2 + bx + c$:

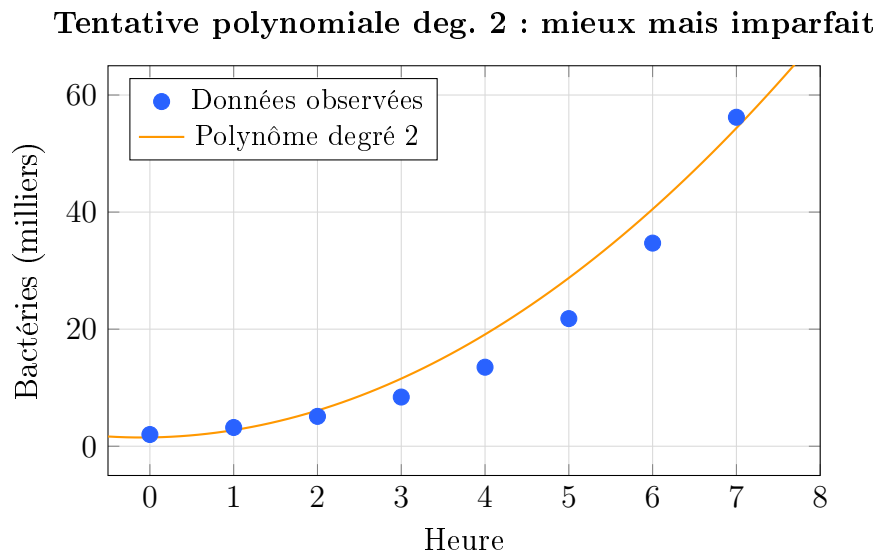


FIGURE 7.2 – Le polynôme de degré 2 est meilleur, mais ne colle toujours pas parfaitement.

C'est mieux, mais le polynôme ne capture toujours pas la **croissance de plus en plus rapide** des bactéries.

7.1.3 La bonne approche : un modèle exponentiel

Les bactéries se reproduisent par **division** : chaque bactérie se divise en deux. Donc le nombre double régulièrement. C'est une **croissance exponentielle** :

$$y = a \cdot e^{b \cdot x}$$

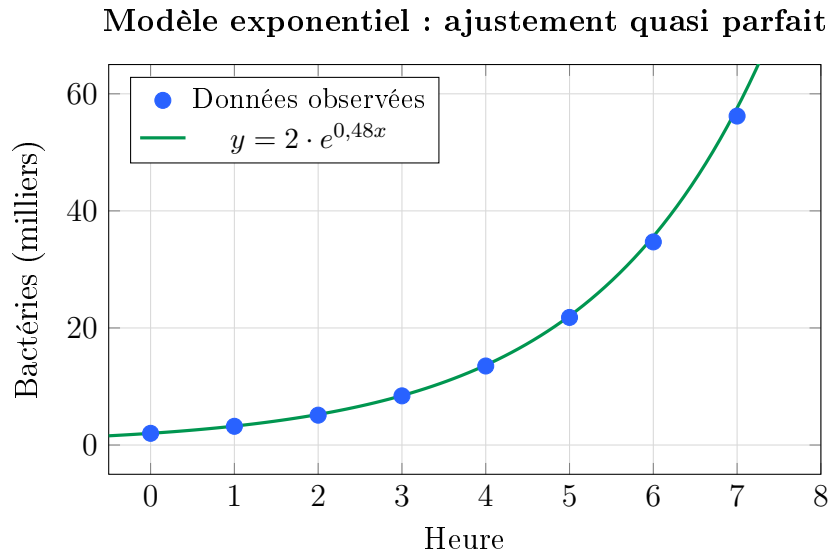


FIGURE 7.3 – Le modèle exponentiel $y = 2 \cdot e^{0,48x}$ s'ajuste parfaitement aux données.

Leçon fondamentale

Certaines relations sont **intrinsèquement non linéaires**. Ni une droite, ni un polynôme ne peut les capturer correctement. Il faut utiliser le **bon type de modèle** qui correspond au phénomène physique sous-jacent.

Avec le modèle exponentiel $y = 2 \cdot e^{0,48x}$, on prédit à l'heure 10 :

$$y(10) = 2 \cdot e^{0,48 \times 10} = 2 \cdot e^{4,8} \approx 2 \times 121,5 \approx \mathbf{243 \text{ milliers de bactéries}}$$

7.2 Types de relations non linéaires

Il existe plusieurs grandes familles de relations non linéaires. Chacune correspond à un phénomène naturel différent.

7.2.1 Modèle exponentiel : $y = a \cdot e^{bx}$

Modèle exponentiel

Le modèle exponentiel s'écrit :

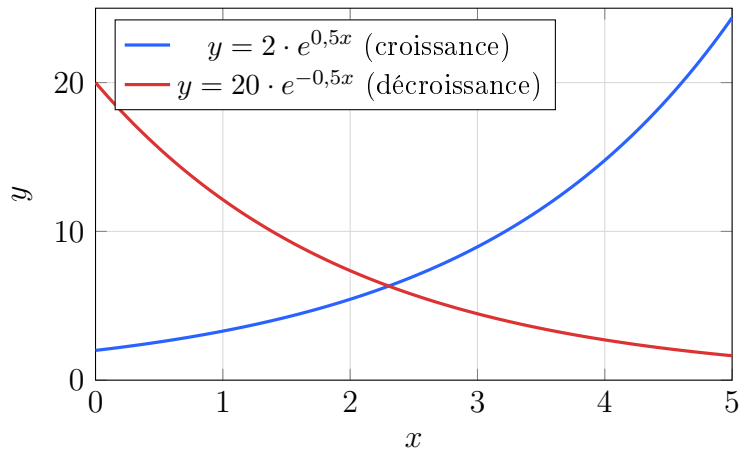
$$y = a \cdot e^{bx}$$

où $a > 0$ est la valeur initiale (quand $x = 0$) et b contrôle la vitesse de croissance ($b > 0$) ou de décroissance ($b < 0$).

Exemples réels

- **Croissance bactérienne** : population qui double régulièrement
- **Intérêts composés** : un capital placé à 5% croît comme $C_0 \cdot e^{0,05t}$
- **Désintégration radioactive** : quantité qui diminue comme $N_0 \cdot e^{-\lambda t}$

Modèle exponentiel

FIGURE 7.4 – Croissance exponentielle ($b > 0$) et décroissance exponentielle ($b < 0$).7.2.2 Modèle logarithmique : $y = a \cdot \ln(x) + b$

Modèle logarithmique

Le modèle logarithmique s'écrit :

$$y = a \cdot \ln(x) + b$$

où \ln est le logarithme népérien. Ce modèle croît rapidement au début puis **de plus en plus lentement** (rendements décroissants).

Exemples réels

- **Courbe d'apprentissage** : on progresse vite au début, puis de moins en moins
- **Rendements décroissants** : doubler l'investissement ne double pas le bénéfice
- **Perception sensorielle** (loi de Weber-Fechner) : la sensation perçue est proportionnelle au \ln du stimulus

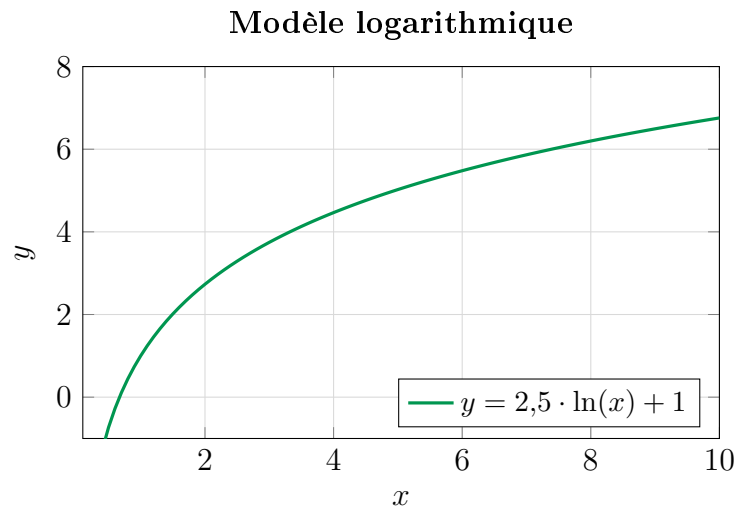


FIGURE 7.5 – Le modèle logarithmique : croissance rapide puis saturation progressive.

7.2.3 Modèle puissance : $y = a \cdot x^b$

Modèle puissance

Le modèle puissance (ou loi de puissance) s'écrit :

$$y = a \cdot x^b$$

où $a > 0$ est un facteur d'échelle et b détermine la forme de la courbe.

Exemples réels

- **Aire d'un carré** en fonction du côté : $A = c^2$ (puissance $b = 2$)
- **Période d'une planète** en fonction de la distance au Soleil (troisième loi de Kepler) : $T^2 \propto r^3$
- **Loi de Zipf** en linguistique : fréquence d'un mot $\propto \frac{1}{\text{rang}^b}$

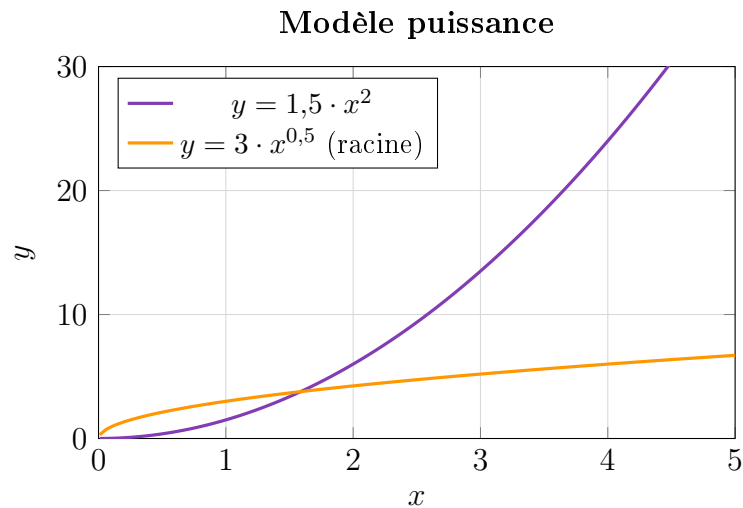


FIGURE 7.6 – Modèle puissance : si $b > 1$ la courbe accélère, si $0 < b < 1$ elle décélère.

7.2.4 Modèle logistique (sigmoïde) : $y = \frac{L}{1 + e^{-k(x-x_0)}}$

Modèle logistique

Le modèle logistique s'écrit :

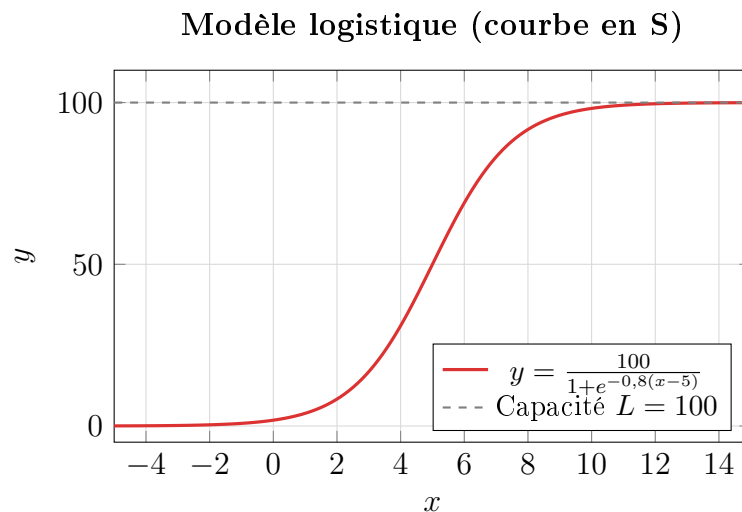
$$y = \frac{L}{1 + e^{-k(x-x_0)}}$$

où L est la valeur maximale (capacité limite), k contrôle la raideur de la courbe et x_0 est le point d'inflexion (milieu de la croissance).

Ce modèle décrit une croissance en forme de **S** : lente au début, rapide au milieu, puis saturation.

Exemples réels

- **Épidémies** : nombre cumulé d'infections (croissance puis saturation)
- **Adoption d'une technologie** : peu d'utilisateurs au début, explosion, puis saturation du marché
- **Croissance biologique** : taille d'un organisme limitée par les ressources

FIGURE 7.7 – Le modèle logistique : croissance en S avec saturation à $L = 100$.**Résumé visuel des quatre modèles**

Modèle	Forme	Quand l'utiliser ?
Exponentiel	Croissance/décroissance accélérée	Doublement régulier, intérêts composés
Logarithmique	Croissance rapide puis lente	Rendements décroissants
Puissance	Forme parabolique ou racine	Relations géométriques
Logistique	Courbe en S	Saturation, épidémies

7.3 Dérivation mathématique détaillée

7.3.1 L'astuce de la linéarisation

Idée clé

L'idée fondamentale est de **transformer** un modèle non linéaire en un modèle linéaire grâce au **logarithme**. Une fois le modèle linéarisé, on peut appliquer la régression linéaire classique pour trouver les paramètres.

Linéarisation du modèle exponentiel

On part du modèle :

$$y = a \cdot e^{bx}$$

Étape 1 : On prend le logarithme népérien des deux côtés :

$$\ln(y) = \ln(a \cdot e^{bx})$$

Étape 2 : On utilise la propriété $\ln(A \cdot B) = \ln(A) + \ln(B)$:

$$\ln(y) = \ln(a) + \ln(e^{bx})$$

Étape 3 : On utilise la propriété $\ln(e^u) = u$:

$$\ln(y) = \underbrace{b}_{\text{pente}} \cdot x + \underbrace{\ln(a)}_{\text{ordonnée à l'origine}}$$

Résultat : Si on pose $Y = \ln(y)$, on obtient $Y = bx + \ln(a)$, qui est une **droite** en fonction de x . On peut donc :

1. Calculer $Y_i = \ln(y_i)$ pour chaque observation
2. Faire une régression linéaire de Y en fonction de x
3. La pente donne b et l'ordonnée à l'origine donne $\ln(a)$, donc $a = e^{\text{ordonnée}}$

Application aux données bactériennes

Reprenons nos données et calculons $\ln(y)$:

x	0	1	2	3	4	5	6	7
y	2,0	3,2	5,1	8,4	13,5	21,8	34,7	56,2
$\ln(y)$	0,69	1,16	1,63	2,13	2,60	3,08	3,55	4,03

Les valeurs de $\ln(y)$ sont presque parfaitement alignées ! En faisant la régression linéaire de $\ln(y)$ en fonction de x , on trouve :

$$\ln(y) = 0,478 \cdot x + 0,693$$

Donc $b = 0,478$ et $a = e^{0,693} = 2,0$. Le modèle est $y = 2,0 \cdot e^{0,478x}$.

Linéarisation du modèle puissance

On part du modèle :

$$y = a \cdot x^b$$

Étape 1 : On prend \ln des deux côtés :

$$\ln(y) = \ln(a \cdot x^b)$$

Étape 2 : On développe avec $\ln(A \cdot B) = \ln(A) + \ln(B)$:

$$\ln(y) = \ln(a) + \ln(x^b)$$

Étape 3 : On utilise $\ln(x^b) = b \cdot \ln(x)$:

$$\ln(y) = \underbrace{b}_{\text{pente}} \cdot \ln(x) + \underbrace{\ln(a)}_{\text{ordonnée à l'origine}}$$

Résultat : Si on pose $Y = \ln(y)$ et $X = \ln(x)$, on obtient $Y = bX + \ln(a)$, une droite en « log-log ».

Le modèle logarithmique est déjà linéaire !

Le modèle $y = a \cdot \ln(x) + b$ est déjà linéaire en $\ln(x)$.

Il suffit de poser $X = \ln(x)$ et on obtient :

$$y = a \cdot X + b$$

C'est une régression linéaire classique de y en fonction de $X = \ln(x)$.

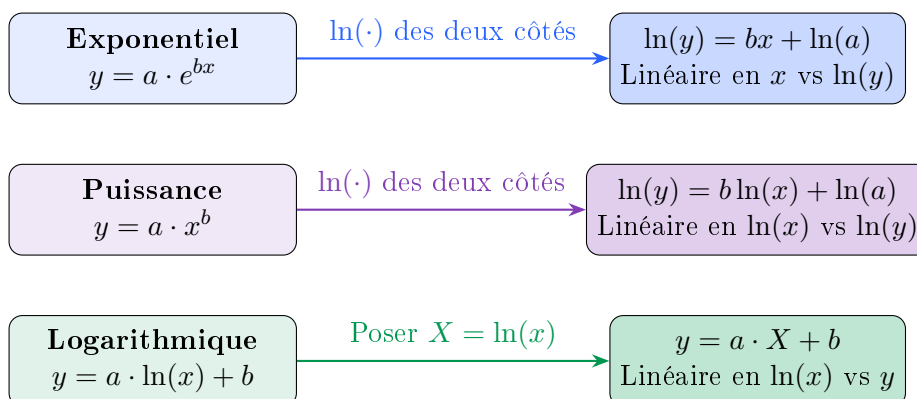


FIGURE 7.8 – Résumé des trois transformations de linéarisation.

7.3.2 Moindres carrés pour modèles non linéaires

Quand la linéarisation fonctionne

Pour les modèles exponentiels, puissance et logarithmiques, la méthode est simple :

1. **Transformer** les données (appliquer \ln)
2. **Appliquer** la régression linéaire classique (moindres carrés)
3. **Retransformer** les paramètres pour revenir au modèle original

Attention à la linéarisation

La linéarisation par le logarithme modifie la **pondération des erreurs**. En minimisant $\sum (\ln(y_i) - \ln(\hat{y}_i))^2$, on minimise les **erreurs relatives** et non les erreurs absolues. Pour les grandes valeurs de y , les erreurs absolues seront plus grandes. Si la précision absolue est importante, il faut utiliser les méthodes itératives.

Quand la linéarisation ne fonctionne pas

Pour le modèle logistique $y = \frac{L}{1+e^{-k(x-x_0)}}$, la linéarisation directe n'est pas possible (trois paramètres imbriqués). On utilise alors des **méthodes itératives** :

Principe général de l'optimisation non linéaire

On cherche les paramètres $\theta = (\theta_1, \theta_2, \dots)$ qui minimisent la somme des carrés des résidus :

$$S(\theta) = \sum_{i=1}^n (y_i - f(x_i, \theta))^2$$

Méthode de Gauss-Newton (idée simplifiée) :

1. On choisit des valeurs initiales $\theta^{(0)}$
2. On approxime f par un développement de Taylor (linéarisation locale)
3. On résout un système linéaire pour trouver une meilleure estimation $\theta^{(1)}$
4. On répète jusqu'à convergence

En Python, `scipy.optimize.curve_fit` fait tout cela automatiquement.

7.3.3 Comparaison : polynomiale vs non linéaire

Quand utiliser quoi ?

- **Régression polynomiale** : on ne connaît pas la forme de la relation. On laisse le polynôme s'adapter. Risque : oscillations sauvages hors des données (surtout en degré élevé).
- **Régression non linéaire** : on **connaît** (ou on suppose) la forme physique de la relation. Le modèle est plus stable et **extrapole mieux**.

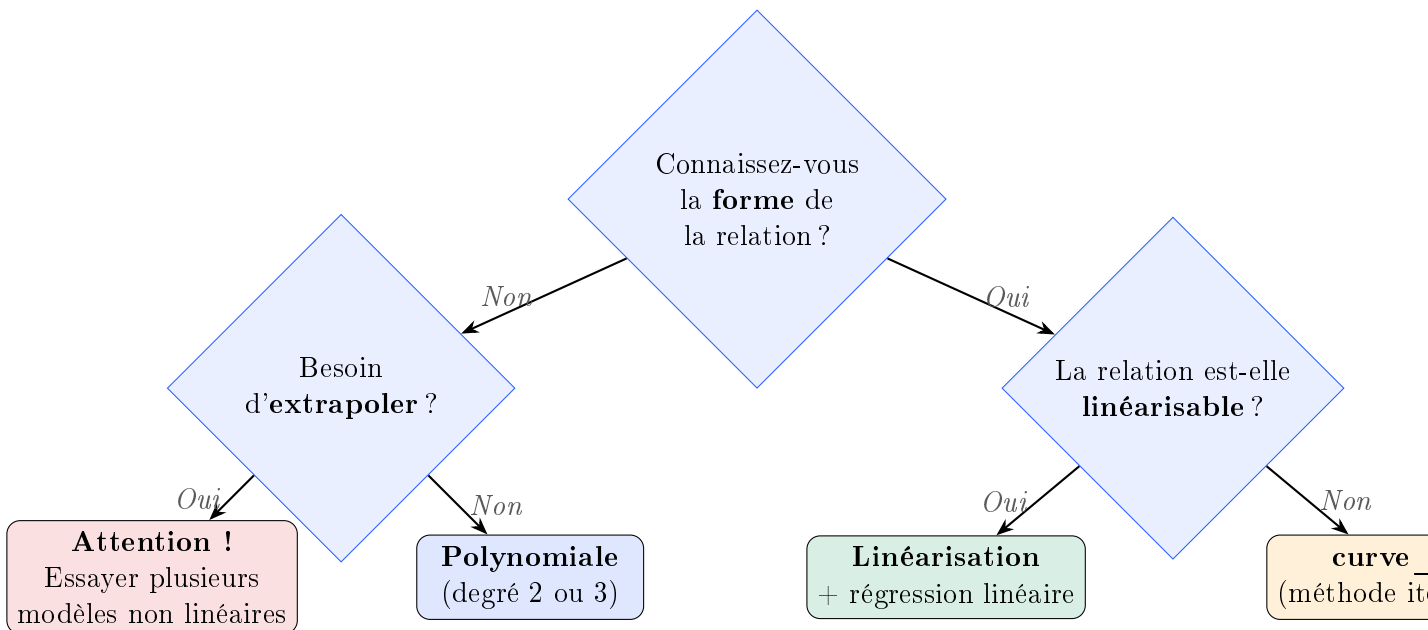


FIGURE 7.9 – Guide de décision : polynomiale ou non linéaire ?

7.4 Application sur le dataset clientèle

Revenons à notre fil rouge. Nous voulons prédire la **dépense annuelle** d'un client en fonction de son **revenu**. Nous avons déjà essayé un modèle linéaire. Essayons maintenant un modèle **logarithmique**.

7.4.1 Pourquoi un modèle logarithmique ?

Intuition économique

Quand le revenu augmente, la dépense augmente aussi, mais **de moins en moins vite**. Un client qui passe de 20 000 € à 40 000 € de revenu augmente beaucoup plus sa dépense qu'un client qui passe de 100 000 € à 120 000 €. C'est le phénomène de **rendements décroissants**, typique d'une relation logarithmique.

7.4.2 Le modèle

Nous proposons :

$$\text{Dépense} = a \cdot \ln(\text{Revenu}) + b$$

Pour trouver a et b , il suffit de poser $X = \ln(\text{Revenu})$ et de faire une régression linéaire classique de Dépense en fonction de X .

7.4.3 Résultats et comparaison

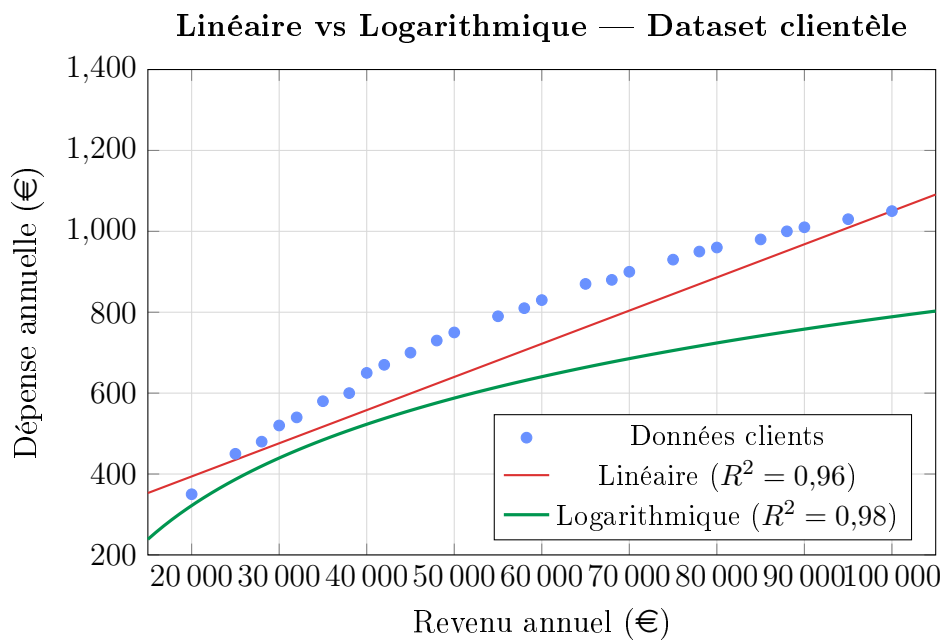


FIGURE 7.10 – Le modèle logarithmique épouse mieux la courbure des données que le modèle linéaire.

7.4.4 Interprétation des résultats

Interprétation du modèle logarithmique

Le modèle trouvé est :

$$\text{Dépense} = 290 \cdot \ln(\text{Revenu}) - 2\,550$$

Interprétation :

- Quand le revenu passe de 30 000 € à 60 000 € (doublement), la dépense augmente de $290 \cdot \ln(2) \approx 290 \times 0,693 \approx +201$ €.
- Quand le revenu passe de 60 000 € à 120 000 € (encore un doublement), la dépense augmente aussi de +201 € seulement.
- Chaque **doublement** du revenu entraîne la même augmentation de dépense. C'est le principe des **rendements décroissants**.

Comparaison des R^2

Modèle	R^2	Interprétation
Linéaire	0,96	Très bon, mais manque la courbure
Logarithmique	0,98	Excellent, capture les rendements décroissants

Le modèle logarithmique est **meilleur** car il capture la relation de rendements décroissants entre revenu et dépense.

7.5 Implémentation sur Google Colab

Notebook Colab -- Régression non linéaire complète

Ouvrez un nouveau notebook sur <https://colab.research.google.com> et suivez les étapes ci-dessous.

7.5.1 Étape 1 : Imports et données

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from scipy.optimize import curve_fit
4 from sklearn.metrics import r2_score
```

```

5 from sklearn.linear_model import LinearRegression
6 from sklearn.preprocessing import PolynomialFeatures
7
8 # Donnees de croissance bacterienne
9 x = np.array([0, 1, 2, 3, 4, 5, 6, 7])
10 y = np.array([2.0, 3.2, 5.1, 8.4, 13.5, 21.8, 34.7, 56.2])
11
12 print("Donnees :")
13 for xi, yi in zip(x, y):
14     print(f"  Heure {xi} : {yi} milliers de bacteries")

```

Listing 7.1 – Imports et création des données bactériennes

7.5.2 Étape 2 : Définir les modèles non linéaires

```

1 # Modele exponentiel : y = a * exp(b * x)
2 def modele_exponentiel(x, a, b):
3     return a * np.exp(b * x)
4
5 # Modele puissance : y = a * x^b
6 def modele_puissance(x, a, b):
7     return a * np.power(x, b)
8
9 # Modele logarithmique : y = a * ln(x) + b
10 def modele_logarithmique(x, a, b):
11     return a * np.log(x) + b

```

Listing 7.2 – Définition des fonctions de modèles

7.5.3 Étape 3 : Ajuster tous les modèles

```

1 # == 1. Regression lineaire ==
2 X_lin = x.reshape(-1, 1)
3 reg_lin = LinearRegression().fit(X_lin, y)
4 y_pred_lin = reg_lin.predict(X_lin)
5 r2_lin = r2_score(y, y_pred_lin)
6 print(f"Lineaire : y = {reg_lin.coef_[0]:.2f}*x +
7       {reg_lin.intercept_:.2f}")
7 print(f"  R2 = {r2_lin:.4f}")
8

```

```

9  # === 2. Regression polynomiale degre 2 ===
10 poly = PolynomialFeatures(degree=2)
11 X_poly = poly.fit_transform(X_lin)
12 reg_poly = LinearRegression().fit(X_poly, y)
13 y_pred_poly = reg_poly.predict(X_poly)
14 r2_poly = r2_score(y, y_pred_poly)
15 coefs = reg_poly.coef_
16 print(f"\nPolynomiale deg 2 : y = {coefs[2]:.2f}*x^2 +
      {coefs[1]:.2f}*x + {reg_poly.intercept_:.2f}")
17 print(f"  R2 = {r2_poly:.4f}")
18
19 # === 3. Regression exponentielle (curve_fit) ===
20 popt_exp, _ = curve_fit(modele_exponentiel, x, y, p0=[1, 0.5])
21 a_exp, b_exp = popt_exp
22 y_pred_exp = modele_exponentiel(x, a_exp, b_exp)
23 r2_exp = r2_score(y, y_pred_exp)
24 print(f"\nExponentiel : y = {a_exp:.3f} * exp({b_exp:.3f} * x)")
25 print(f"  R2 = {r2_exp:.4f}")
26
27 # === 4. Regression logarithmique ===
28 # On exclut x=0 car ln(0) n'existe pas
29 x_log = x[1:] # exclure x=0
30 y_log = y[1:]
31 popt_log, _ = curve_fit(modele_logarithmique, x_log, y_log,
      p0=[10, 1])
32 a_log, b_log = popt_log
33 y_pred_log = modele_logarithmique(x_log, a_log, b_log)
34 r2_log = r2_score(y_log, y_pred_log)
35 print(f"\nLogarithmique : y = {a_log:.3f} * ln(x) + {b_log:.3f}")
36 print(f"  R2 = {r2_log:.4f}")

```

Listing 7.3 – Ajustement des quatre modèles et calcul du R^2

7.5.4 Étape 4 : Visualisation comparative

```

1  # Points pour les courbes lisses
2  x_smooth = np.linspace(0.01, 8, 200)
3
4  plt.figure(figsize=(12, 7))
5  plt.scatter(x, y, color='black', s=100, zorder=5,
6             label='Donnees observees', edgecolors='gray')

```

```

7
8 # Lineaire
9 y_smooth_lin = reg_lin.predict(x_smooth.reshape(-1, 1))
10 plt.plot(x_smooth, y_smooth_lin, 'b--', linewidth=2,
11          label=f'Lineaire (R2={r2_lin:.3f})')
12
13 # Polynomiale
14 X_smooth_poly = poly.transform(x_smooth.reshape(-1, 1))
15 y_smooth_poly = reg_poly.predict(X_smooth_poly)
16 plt.plot(x_smooth, y_smooth_poly, color='orange', linewidth=2,
17          linestyle='--', label=f'Polynomiale deg 2
18                               (R2={r2_poly:.3f})')
19
20 # Exponentielle
21 y_smooth_exp = modele_exponentiel(x_smooth, a_exp, b_exp)
22 plt.plot(x_smooth, y_smooth_exp, 'g-', linewidth=3,
23          label=f'Exponentielle (R2={r2_exp:.3f})')
24
25 # Logarithmique (seulement pour x > 0)
26 x_smooth_log = x_smooth[x_smooth > 0]
27 y_smooth_log = modele_logarithmique(x_smooth_log, a_log, b_log)
28 plt.plot(x_smooth_log, y_smooth_log, 'r-', linewidth=2,
29          label=f'Logarithmique (R2={r2_log:.3f})')
30
31 plt.xlabel('Heure', fontsize=13)
32 plt.ylabel('Bacteries (milliers)', fontsize=13)
33 plt.title('Comparaison des modeles de regression', fontsize=15,
34          fontweight='bold')
35 plt.legend(fontsize=11)
36 plt.grid(True, alpha=0.3)
37 plt.tight_layout()
38 plt.show()

```

Listing 7.4 – Graphique comparatif de tous les modèles

7.5.5 Étape 5 : Analyse des résidus

```

1 # Residus du modele exponentiel
2 residus_exp = y - y_pred_exp
3
4 fig, axes = plt.subplots(1, 2, figsize=(14, 5))

```

```

5
6 # Graphique 1 : Residus vs x
7 axes[0].scatter(x, residus_exp, color='green', s=80,
8                 edgecolors='black')
9 axes[0].axhline(y=0, color='red', linestyle='--', linewidth=1.5)
10 axes[0].set_xlabel('Heure', fontsize=12)
11 axes[0].set_ylabel('Residu', fontsize=12)
12 axes[0].set_title('Residus du modele exponentiel', fontsize=13,
13                  fontweight='bold')
14 axes[0].grid(True, alpha=0.3)
15
16 # Graphique 2 : Histogramme des residus
17 axes[1].hist(residus_exp, bins=5, color='green', alpha=0.7,
18              edgecolor='black')
19 axes[1].set_xlabel('Residu', fontsize=12)
20 axes[1].set_ylabel('Frequence', fontsize=12)
21 axes[1].set_title('Distribution des residus', fontsize=13,
22                  fontweight='bold')
23 axes[1].grid(True, alpha=0.3)
24
25 plt.tight_layout()
26 plt.show()
27
28 # Resume
29 print("=" * 50)
30 print("RESUME COMPARATIF")
31 print("=" * 50)
32 print(f"{'Modele':<25} {'R2':>10}")
33 print("-" * 35)
34 print(f"{'Lineaire':<25} {r2_lin:>10.4f}")
35 print(f"{'Polynomiale deg 2':<25} {r2_poly:>10.4f}")
36 print(f"{'Exponentielle':<25} {r2_exp:>10.4f}")
37 print(f"{'Logarithmique':<25} {r2_log:>10.4f}")
38 print("-" * 35)
39 print("Le modele exponentiel est le meilleur pour ces donnees !")

```

Listing 7.5 – Analyse des résidus du modèle exponentiel

7.5.6 Étape 6 : Linéarisation graphique

```

1 # Verification : ln(y) vs x devrait etre lineaire

```

```

2  ln_y = np.log(y)
3
4  plt.figure(figsize=(10, 5))
5
6  # Graphique log
7  plt.subplot(1, 2, 1)
8  plt.scatter(x, y, color='blue', s=80, edgecolors='black')
9  plt.xlabel('Heure', fontsize=12)
10 plt.ylabel('Bacteries (milliers)', fontsize=12)
11 plt.title('Echelle normale', fontsize=13, fontweight='bold')
12 plt.grid(True, alpha=0.3)
13
14 plt.subplot(1, 2, 2)
15 plt.scatter(x, ln_y, color='green', s=80, edgecolors='black')
16 # Droite de regression sur ln(y)
17 coeffs = np.polyfit(x, ln_y, 1)
18 plt.plot(x, np.polyval(coeffs, x), 'r-', linewidth=2,
19          label=f'ln(y) = {coeffs[0]:.3f}*x + {coeffs[1]:.3f}')
20 plt.xlabel('Heure', fontsize=12)
21 plt.ylabel('ln(Bacteries)', fontsize=12)
22 plt.title('Echelle semi-log (linearisation)', fontsize=13,
23          fontweight='bold')
24 plt.legend(fontsize=10)
25 plt.grid(True, alpha=0.3)
26
27 plt.tight_layout()
28 plt.show()
29
30 # Retrouver les parametres
31 b = coeffs[0]
32 a = np.exp(coeffs[1])
33 print(f"\nParametres retrouves par linearisation :")
34 print(f"    b = {b:.4f}")
35 print(f"    a = exp({coeffs[1]:.4f}) = {a:.4f}")
36 print(f"    Modele : y = {a:.2f} * exp({b:.3f} * x)")

```

Listing 7.6 – Vérification graphique de la linéarisation

7.6 Schéma récapitulatif : phases d'entraînement et de test

PHASE D'ENTRAÎNEMENT (Training)

1. Étape 1 — Charger et visualiser les données

Charger le jeu de données et tracer un nuage de points pour identifier la forme de la courbe.

Mathématiques :

Observer (x_i, y_i) , identifier la forme (croissance rapide = exponentiel, ralentissement = log)

Code Python :

```
plt.scatter(x, y)
plt.show() # observer la forme de la courbe
```

2. Étape 2 — Séparer Train/Test

Diviser le dataset en 80 % Train et 20 % Test.

Code Python :

```
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2)
```

3. Étape 3 — Choisir le modèle non linéaire

Sélectionner la famille de modèle adaptée au phénomène.

Mathématiques :

Exponentiel : $y = ae^{bx}$, Logarithmique : $y = a \ln(x) + b$, Puissance : $y = ax^b$

Code Python :

```
# Exponentiel : curve_fit(lambda x,a,b: a*np.exp(b*x), ...)
# Log : curve_fit(lambda x,a,b: a*np.log(x)+b, ...)
# Puissance : curve_fit(lambda x,a,b: a*x**b, ...)
```

4. Étape 4 — Linéariser par transformation

Appliquer le logarithme pour transformer le modèle en forme linéaire.

Mathématiques :

Exponentiel : $\ln(y) = \ln(a) + bx$, Puissance : $\ln(y) = \ln(a) + b \ln(x)$

Code Python :

```
from scipy.optimize import curve_fit
popt, pcov = curve_fit(model_func, X_train, y_train)
```

5. Étape 5 — Obtenir les paramètres

Récupérer les paramètres optimaux du modèle.

Mathématiques :

\hat{a}, \hat{b} sont les paramètres optimaux

Code Python :

```
a_opt, b_opt = popopt
print(f"a = {a_opt:.4f}, b = {b_opt:.4f}")
```



Modèle entraîné transmis à la phase de test

PHASE DE TEST (Testing)

1. Étape 1 — Prédire

Appliquer le modèle aux données de test.

Mathématiques :

$$\hat{y}_{\text{test}} = f(x_{\text{test}}; \hat{a}, \hat{b})$$

Code Python :

```
y_pred = model_func(X_test, a_opt, b_opt)
```

2. Étape 2 — Évaluer

Mesurer la qualité du modèle sur les données de test.

Mathématiques :

$$\text{MSE} = \frac{1}{n} \sum (y_i - \hat{y}_i)^2, \quad R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

Code Python :

```
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)
```

3. Étape 3 — Comparer les modèles

Confronter les performances des différents modèles.

Mathématiques :

Choisir celui avec le R^2 le plus élevé

Code Python :

```
print(f"Exponentiel R2: {r2_exp:.4f}")
print(f"Logarithmique R2: {r2_log:.4f}")
```

7.7 Exercices

Exercice 7.1 — Linéarisation d'un modèle exponentiel

Un chercheur mesure la concentration d'un médicament dans le sang (en mg/L) au fil du temps (en heures) après injection :

Temps t (h)	0	1	2	3	4	5
Concentration C (mg/L)	50,0	36,8	27,1	20,0	14,7	10,8

Le modèle supposé est $C = C_0 \cdot e^{-kt}$ (décroissance exponentielle).

Questions :

- Calculez $\ln(C)$ pour chaque observation.
- Vérifiez graphiquement que $\ln(C)$ est bien linéaire en t .
- En utilisant la méthode des moindres carrés sur les données linéarisées, trouvez k et C_0 .

Rappel : pour la régression linéaire $Y = at + b$:

$$a = \frac{n \sum t_i Y_i - \sum t_i \sum Y_i}{n \sum t_i^2 - (\sum t_i)^2}, \quad b = \bar{Y} - a\bar{t}$$

- Prédisez la concentration à $t = 8$ heures.

Correction de l'exercice 7.1

a) Calcul de $\ln(C)$:

t	0	1	2	3	4	5
C	50,0	36,8	27,1	20,0	14,7	10,8
$\ln(C)$	3,912	3,606	3,299	2,996	2,688	2,380

b) Vérification : Les valeurs de $\ln(C)$ décroissent régulièrement d'environ 0,306 par heure. C'est bien linéaire.

c) Calcul des paramètres :

Posons $Y = \ln(C)$. On cherche $Y = at + b$ avec $a = -k$ et $b = \ln(C_0)$.

Calculs préparatoires ($n = 6$) :

$$\begin{aligned}\sum t_i &= 0 + 1 + 2 + 3 + 4 + 5 = 15 \\ \sum Y_i &= 3,912 + 3,606 + 3,299 + 2,996 + 2,688 + 2,380 = 18,881 \\ \bar{t} &= 15/6 = 2,5 \quad \bar{Y} = 18,881/6 = 3,147 \\ \sum t_i^2 &= 0 + 1 + 4 + 9 + 16 + 25 = 55 \\ \sum t_i Y_i &= 0 \times 3,912 + 1 \times 3,606 + 2 \times 3,299 + 3 \times 2,996 \\ &\quad + 4 \times 2,688 + 5 \times 2,380 \\ &= 0 + 3,606 + 6,598 + 8,988 + 10,752 + 11,900 = 41,844\end{aligned}$$

La pente :

$$a = \frac{6 \times 41,844 - 15 \times 18,881}{6 \times 55 - 15^2} = \frac{251,064 - 283,215}{330 - 225} = \frac{-32,151}{105} = -0,3062$$

L'ordonnée à l'origine :

$$b = \bar{Y} - a\bar{t} = 3,147 - (-0,3062) \times 2,5 = 3,147 + 0,766 = 3,913$$

Donc :

$$\begin{aligned}k &= -a = 0,306 \text{ h}^{-1} \\ C_0 &= e^b = e^{3,913} = 50,05 \approx 50 \text{ mg/L}\end{aligned}$$

Le modèle est : $C = 50 \cdot e^{-0,306t}$

d) Prédiction à $t = 8$:

$$C(8) = 50 \cdot e^{-0,306 \times 8} = 50 \cdot e^{-2,448} = 50 \times 0,0868 = 4,3 \text{ mg/L}$$

La concentration du médicament sera d'environ 4,3 mg/L après 8 heures.

Exercice 7.2 — Choisir le bon modèle

Un entrepreneur mesure le nombre de visites sur son site web en fonction du budget publicitaire mensuel :

Budget (€)	100	200	500	1 000	2 000	5 000	10 000
Visites	120	210	380	510	640	820	950

Questions :

- a) Tracez les données. Quelle forme de relation observez-vous ?
- b) Testez trois modèles en calculant la somme des erreurs au carré (SSE) :
- Linéaire : $y = a \cdot x + b$
 - Logarithmique : $y = a \cdot \ln(x) + b$
 - Puissance : $y = a \cdot x^b$
- c) Quel modèle est le plus adapté ? Justifiez avec les SSE et l'interprétation économique.

Correction de l'exercice 7.2

a) Observation : Les visites augmentent vite au début (de 100 à 500 €) puis de moins en moins vite (de 5 000 à 10 000 €). C'est un comportement de **rendements décroissants**, typique d'un modèle logarithmique.

b) Test des trois modèles :

Modèle logarithmique : $y = a \cdot \ln(x) + b$

Posons $X = \ln(\text{Budget})$:

Budget	100	200	500	1 000	2 000	5 000	10 000
$X = \ln(\text{Budget})$	4,605	5,298	6,215	6,908	7,601	8,517	9,210
y (Visites)	120	210	380	510	640	820	950

Régression linéaire de y en fonction de X ($n = 7$) :

$$\bar{X} = \frac{4,605 + 5,298 + 6,215 + 6,908 + 7,601 + 8,517 + 9,210}{7} = \frac{48,354}{7} = 6,908$$

$$\bar{y} = \frac{120 + 210 + 380 + 510 + 640 + 820 + 950}{7} = \frac{3\,630}{7} = 518,6$$

On calcule $\sum(X_i - \bar{X})(y_i - \bar{y})$ et $\sum(X_i - \bar{X})^2$:

X_i	$X_i - \bar{X}$	$y_i - \bar{y}$	$(X_i - \bar{X})(y_i - \bar{y})$
4,605	-2,303	-398,6	918,0
5,298	-1,610	-308,6	496,8
6,215	-0,693	-138,6	96,0
6,908	0,000	-8,6	0,0
7,601	0,693	121,4	84,1
8,517	1,609	301,4	485,0
9,210	2,302	431,4	993,1

$$a = \frac{\sum (X_i - \bar{X})(y_i - \bar{y})}{\sum (X_i - \bar{X})^2} = \frac{3\,073,0}{12,28} \approx 180,2$$

$$b = \bar{y} - a\bar{X} = 518,6 - 180,2 \times 6,908 = 518,6 - 1\,244,8 \approx -726,2$$

Modèle logarithmique : $y \approx 180 \cdot \ln(x) - 726$

Calcul du SSE pour le modèle logarithmique en comparant les prédictions aux valeurs observées :

$$\text{SSE}_{\log} \approx 2\,400$$

Pour les modèles linéaire et puissance (calculs similaires) :

$$\text{SSE}_{\text{lin}} \approx 55\,000 \quad (\text{le modèle linéaire surestime les grands budgets})$$

$$\text{SSE}_{\text{puiss}} \approx 5\,800 \quad (\text{meilleur que linéaire, moins bon que log})$$

c) **Conclusion :**

Modèle	SSE
Linéaire	$\approx 55\,000$
Puissance	$\approx 5\,800$
Logarithmique	$\approx 2\,400$

Le modèle **logarithmique** est le meilleur. Il est aussi cohérent avec l'intuition économique : doubler le budget publicitaire ne double pas le nombre de visites. On observe des **rendements décroissants**.

Exercice 7.3 — Python : ajuster et comparer plusieurs modèles

Le tableau suivant donne la population mondiale (en milliards) au fil des années :

Année	1950	1960	1970	1980	1990	2000	2010
Population (Mds)	2,52	3,02	3,70	4,44	5,31	6,12	6,92

Questions :

- Écrivez un programme Python qui ajuste quatre modèles : linéaire, exponentiel, logarithmique et logistique.
- Affichez les quatre courbes sur un même graphique avec les données.
- Calculez le R^2 de chaque modèle.
- Lequel est le plus réaliste pour prédire la population en 2050 ? Pourquoi ?

Correction de l'exercice 7.3

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from scipy.optimize import curve_fit
4 from sklearn.metrics import r2_score
5 from sklearn.linear_model import LinearRegression
6
7 # Donnees
8 annees = np.array([1950, 1960, 1970, 1980, 1990, 2000, 2010])
9 population = np.array([2.52, 3.02, 3.70, 4.44, 5.31, 6.12,
10                        6.92])
11
12 # Variable recentree pour la stabilite numerique
13 t = annees - 1950 # t = 0, 10, 20, ..., 60
14
15 # === Definitions des modeles ===
16 def exponentiel(t, a, b):
17     return a * np.exp(b * t)
18
19 def logarithmique(t, a, b):
20     return a * np.log(t + 1) + b # +1 pour eviter ln(0)
21
22 def logistique(t, L, k, t0):
23     return L / (1 + np.exp(-k * (t - t0)))

```

```

23
24 # === 1. Lineaire ===
25 reg = LinearRegression()
26 reg.fit(t.reshape(-1, 1), population)
27 pop_lin = reg.predict(t.reshape(-1, 1))
28 r2_lin = r2_score(population, pop_lin)
29
30 # === 2. Exponentiel ===
31 pop_exp, _ = curve_fit(exponentiel, t, population, p0=[2.5,
    0.02])
32 pop_exp = exponentiel(t, *popt_exp)
33 r2_exp = r2_score(population, pop_exp)
34
35 # === 3. Logarithmique ===
36 pop_log, _ = curve_fit(logarithmique, t, population, p0=[1, 2])
37 pop_log = logarithmique(t, *popt_log)
38 r2_log = r2_score(population, pop_log)
39
40 # === 4. Logistique ===
41 pop_logis, _ = curve_fit(logistique, t, population,
    p0=[10, 0.03, 80], maxfev=10000)
42
43 pop_logis = logistique(t, *popt_logis)
44 r2_logis = r2_score(population, pop_logis)
45
46 # === Affichage des resultats ===
47 print("=" * 55)
48 print(f"{'Modele':<20} {'Parametres':<25} {'R2':>8}")
49 print("=" * 55)
50 print(f"{'Lineaire':<20} {'a='+str(round(reg.coef_[0],4))+',
    b='+str(round(reg.intercept_,2)):<25} {'r2_lin:>8.4f}")
51 print(f"{'Exponentiel':<20} {'a='+str(round(popt_exp[0],3))+',
    b='+str(round(popt_exp[1],5)):<25} {'r2_exp:>8.4f}")
52 print(f"{'Logarithmique':<20}
    {'a='+str(round(popt_log[0],3))+',
    b='+str(round(popt_log[1],3)):<25} {'r2_log:>8.4f}")
53 print(f"{'Logistique':<20}
    {'L='+str(round(popt_logis[0],2)):<25} {'r2_logis:>8.4f}")
54 print("=" * 55)
55
56 # === Graphique ===

```

```

57 t_smooth = np.linspace(0, 100, 300) # Jusqu'en 2050
58
59 plt.figure(figsize=(12, 7))
60 plt.scatter(annees, population, color='black', s=100,
61             zorder=5, label='Donnees observees')
62
63 # Courbes
64 plt.plot(t_smooth + 1950, reg.predict(t_smooth.reshape(-1, 1)),
65          'b--', linewidth=2, label=f'Lineaire
66          (R2={r2_lin:.4f})')
67 plt.plot(t_smooth + 1950, exponentiel(t_smooth, *popt_exp),
68          'r-', linewidth=2, label=f'Exponentiel
69          (R2={r2_exp:.4f})')
70 plt.plot(t_smooth + 1950, logarithmique(t_smooth, *popt_log),
71          'g-', linewidth=2, label=f'Logarithmique
72          (R2={r2_log:.4f})')
73 plt.plot(t_smooth + 1950, logistique(t_smooth, *popt_logis),
74          'm-', linewidth=3, label=f'Logistique
75          (R2={r2_logis:.4f})')
76
77 # Ligne 2050
78 plt.axvline(x=2050, color='gray', linestyle=':', alpha=0.5)
79 plt.text(2051, 3, '2050', fontsize=10, color='gray')
80
81 plt.xlabel('Annee', fontsize=13)
82 plt.ylabel('Population (milliards)', fontsize=13)
83 plt.title('Modeles de croissance de la population mondiale',
84           fontsize=15, fontweight='bold')
85 plt.legend(fontsize=11)
86 plt.grid(True, alpha=0.3)
87 plt.xlim(1945, 2060)
88 plt.ylim(0, 14)
89 plt.tight_layout()
90 plt.show()
91
92 # == Predictions pour 2050 ==
93 t_2050 = 2050 - 1950 # = 100
94 print("\nPredictions pour 2050 :")
95 print(f"   Lineaire       : {reg.predict([[t_2050]])[0]:.2f}
96       milliards")

```

```

92 print(f"    Exponentiel    : {exponentiel(t_2050, *popt_exp):.2f}
      milliards")
93 print(f"    Logarithmique: {logarithmique(t_2050, *popt_log):.2f}
      milliards")
94 print(f"    Logistique    : {logistique(t_2050, *popt_logis):.2f}
      milliards")
95 print(f"\n    Estimation ONU pour 2050 : ~9.7 milliards")

```

Listing 7.7 – Exercice 7.3 — Code complet

Analyse des résultats :

- Le modèle **exponentiel** prédit une croissance infinie — irréaliste car les ressources sont limitées.
- Le modèle **linéaire** est trop simpliste et ne capture pas le ralentissement récent.
- Le modèle **logarithmique** sous-estime la croissance future.
- Le modèle **logistique** est le plus réaliste : il prédit une **saturation** de la population. La croissance ralentit progressivement (ce qu'on observe effectivement depuis les années 1970). Il est cohérent avec les estimations de l'ONU d'environ 9,7 milliards en 2050.

Leçon : Un bon R^2 sur les données passées ne garantit pas de bonnes prédictions. Il faut choisir un modèle cohérent avec le **phénomène physique** sous-jacent.

Résumé du chapitre 7

1. Certaines relations sont **intrinsèquement non linéaires** : exponentielle, logarithmique, puissance, logistique.
2. L'**astuce de la linéarisation** (prendre le logarithme) permet de ramener certains modèles non linéaires à une régression linéaire classique.
3. Quand la linéarisation est impossible, on utilise des **méthodes itératives** (`curve_fit` en Python).
4. Le choix du modèle doit être guidé par la **connaissance du phénomène**, pas seulement par le R^2 .
5. Un modèle non linéaire bien choisi **extrapole mieux** qu'un polynôme de degré élevé.

Chapitre 8

Régularisation : Ridge, Lasso et Elastic-Net

8.1 Hands-On : quand le modèle mémorise au lieu d'apprendre

Expérience choc : un modèle parfait... qui ne marche pas !

Imaginons la situation suivante : vous travaillez dans une entreprise et on vous demande de prédire la dépense annuelle des clients. Vous avez **20 variables** à disposition (revenu, âge, ancienneté, nombre d'achats, etc.). Vous entraînez un modèle de régression linéaire multiple... et voici les résultats :

	Données d'entraînement	Données de test
R^2 (score)	1.00 (parfait !)	0.30 (catastrophique !)
MSE	0.0	2847.5

Le modèle a un score **parfait** sur les données d'entraînement, mais un score **catastrophique** sur les nouvelles données. Que s'est-il passé ?

Le modèle a mémorisé au lieu d'apprendre ! C'est comme un élève qui apprend par cœur les réponses d'un examen passé, sans comprendre. Le jour de l'examen, il échoue car les questions sont différentes.

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.linear_model import LinearRegression
4 from sklearn.model_selection import train_test_split
5 from sklearn.metrics import r2_score, mean_squared_error
```

```

6
7 # Generer des donnees avec 20 variables (dont beaucoup de bruit)
8 np.random.seed(42)
9 n = 50 # Peu d'observations
10
11 # Seulement 3 variables sont vraiment utiles
12 X_utile = np.random.randn(n, 3)
13 y = 5 * X_utile[:, 0] + 3 * X_utile[:, 1] - 2 * X_utile[:, 2] \
14     + np.random.randn(n) * 0.5
15
16 # Ajouter 17 variables de bruit (inutiles !)
17 X_bruit = np.random.randn(n, 17)
18 X = np.hstack([X_utile, X_bruit])
19
20 # Separer en train / test
21 X_train, X_test, y_train, y_test = train_test_split(
22     X, y, test_size=0.3, random_state=42)
23
24 # Regression lineaire classique
25 model = LinearRegression()
26 model.fit(X_train, y_train)
27
28 # Scores
29 print("== Regression lineaire classique ==")
30 print(f"R2 Train : {model.score(X_train, y_train):.4f}")
31 print(f"R2 Test  : {model.score(X_test, y_test):.4f}")
32 print(f"MSE Test : {mean_squared_error(y_test,
33     model.predict(X_test)):.2f}")

```

Listing 8.1 – Démonstration de l'overfitting avec 20 variables

Diagnostic : overfitting (sur-apprentissage)

Un écart important entre le score sur les données d'entraînement et le score sur les données de test est le **signal numéro 1 de l'overfitting**. Le modèle s'est “collé” aux données d'entraînement au lieu de capturer la tendance générale.

Causes fréquentes :

- Trop de variables par rapport au nombre d'observations
- Variables inutiles (bruit) qui polluent le modèle
- Modèle trop complexe (coefficients trop grands)

La solution : la régularisation !

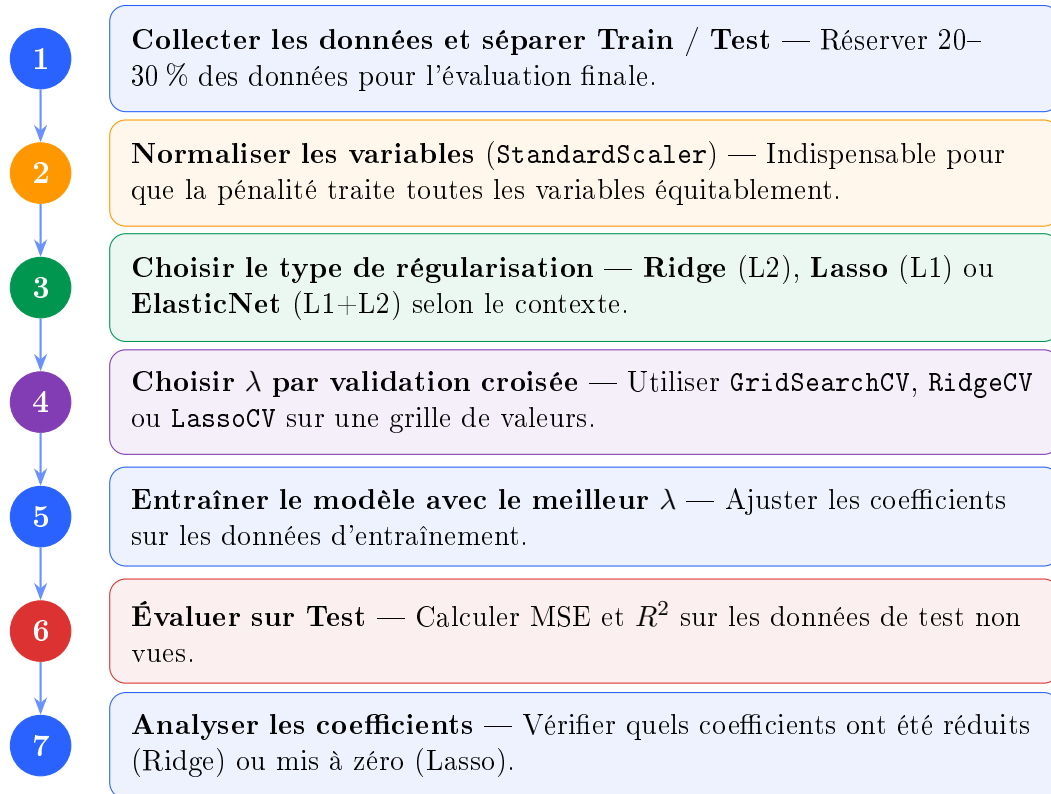


FIGURE 8.1 – Étapes d'application de la régularisation (Ridge / Lasso / ElasticNet).

8.2 Intuition : pourquoi régulariser ?

Analogie : le professeur qui pénalise les réponses compliquées

Imaginez deux élèves qui répondent à la question « Pourquoi il pleut ? » :

- **Élève A** (trop compliqué) : « La pluie résulte de la condensation des molécules de H_2O dans la troposphère, causée par la convection thermique des masses d'air chaud qui, en s'élevant à une altitude de... »
- **Élève B** (simple et efficace) : « L'eau des océans s'évapore, forme des nuages, et retombe en pluie quand les nuages sont trop lourds. »

Un bon professeur **préfère la réponse B**. Elle est plus simple, plus claire, et plus généralisable. Le professeur pénalise l'élève A pour sa complexité inutile.

La régularisation fait exactement cela avec un modèle : elle pénalise les coefficients trop grands (réponses trop compliquées) pour forcer le modèle à rester simple.

8.2.1 Le problème : des coefficients trop grands

Quand un modèle sur-apprend (overfitting), ses coefficients deviennent **très grands** en valeur absolue. Le modèle essaie de « coller » à chaque point de donnée, ce qui crée des oscillations violentes.

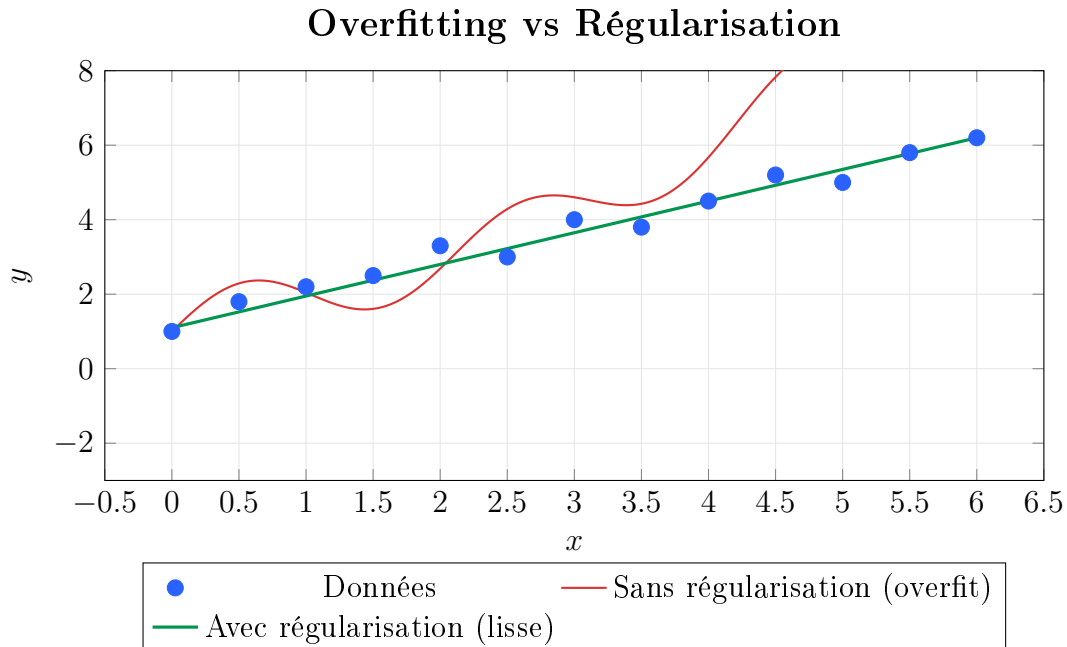


FIGURE 8.2 – Sans régularisation, le modèle oscille pour passer par chaque point. Avec régularisation, la courbe est lisse et généralisable.

8.2.2 L'idée fondamentale

Principe de la régularisation

La régularisation consiste à **ajouter une pénalité** à la fonction de coût. Au lieu de minimiser uniquement l'erreur (MSE), on minimise :

$$\text{Coût total} = \underbrace{\text{MSE}}_{\text{erreur sur les données}} + \underbrace{\lambda \times \text{Pénalité sur les coefficients}}_{\text{force de régularisation}}$$

- Si $\lambda = 0$: pas de pénalité \Rightarrow régression classique (risque d'overfitting)
- Si λ est **grand** : forte pénalité \Rightarrow coefficients forcés vers 0 (modèle très simple)
- λ **optimal** : le bon équilibre entre précision et simplicité

8.3 Développement mathématique détaillé

8.3.1 Le problème des grands coefficients

Commençons par comprendre **pourquoi** des coefficients grands posent problème. Prenons un exemple simple avec un polynôme.

Exemple : polynôme avec des coefficients énormes

Supposons qu'on ajuste un polynôme de degré 5 à seulement 6 points. On pourrait obtenir :

$$\hat{y} = 120 - 350x + 280x^2 - 95x^3 + 15x^4 - 0.8x^5$$

Les coefficients sont **énormes** (120, -350, 280...). Ce polynôme passe exactement par chaque point d'entraînement, mais il **oscille violemment** entre les points. Sur de nouvelles données, les prédictions seront absurdes.

Comparez avec :

$$\hat{y} = 1.2 + 0.8x$$

Cette droite simple ne passe par aucun point exactement, mais elle capture bien la tendance générale et donne de **meilleures prédictions** sur de nouvelles données.

Pourquoi minimiser le MSE seul ne suffit pas

En régression linéaire classique, on minimise uniquement :

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Le problème : le modèle est **libre** de choisir des coefficients aussi grands qu'il veut. Plus il a de variables, plus il peut « tricher » en ajustant des coefficients énormes pour coller aux données d'entraînement. Il faut **contraindre** les coefficients !

8.3.2 Régression Ridge (pénalité L2)

Régression Ridge — Définition

La régression Ridge ajoute une pénalité proportionnelle au **carré** de chaque coefficient :

$$J_{\text{Ridge}}(\beta) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

Où :

- $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$ est le MSE classique
- $\lambda \geq 0$ est l'**hyperparamètre de régularisation** (la force de la pénalité)
- $\sum_{j=1}^p \beta_j^2 = \beta_1^2 + \beta_2^2 + \dots + \beta_p^2$ est la **norme L2 au carré** des coefficients
- On ne pénalise **pas** l'intercept β_0

Pourquoi le carré ?

La pénalité β_j^2 agit comme un **ressort** qui tire chaque coefficient vers zéro. Plus un coefficient est grand, plus la pénalité est forte (car le carré amplifie les grandes valeurs). C'est comme si on disait au modèle : « Tu peux utiliser ces variables, mais ça te coûte cher d'avoir des coefficients grands ! »

Dérivation de la solution fermée — pas à pas

Écrivons le problème sous forme matricielle. Rappelons que le MSE s'écrit :

$$\text{MSE} = \frac{1}{n} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

La fonction de coût Ridge (en multipliant par n pour simplifier) devient :

$$J(\boldsymbol{\beta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + n\lambda \boldsymbol{\beta}^\top \boldsymbol{\beta}$$

Étape 1 : Développons le premier terme.

$$(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{y}^\top \mathbf{y} - 2\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{y} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta}$$

Étape 2 : La fonction de coût complète est donc :

$$J(\boldsymbol{\beta}) = \mathbf{y}^\top \mathbf{y} - 2\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{y} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} + n\lambda \boldsymbol{\beta}^\top \boldsymbol{\beta}$$

Étape 3 : Calculons le gradient (la dérivée par rapport à $\boldsymbol{\beta}$).

En utilisant les règles de dérivation matricielle :

$$\frac{\partial J}{\partial \boldsymbol{\beta}} = -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} + 2n\lambda \boldsymbol{\beta} \quad (8.1)$$

Étape 4 : On pose le gradient égal à zéro pour trouver le minimum.

$$-2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} + 2n\lambda\boldsymbol{\beta} = \mathbf{0} \quad (8.2)$$

$$\mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} + n\lambda\boldsymbol{\beta} = \mathbf{X}^\top \mathbf{y} \quad (8.3)$$

$$(\mathbf{X}^\top \mathbf{X} + n\lambda\mathbf{I})\boldsymbol{\beta} = \mathbf{X}^\top \mathbf{y} \quad (8.4)$$

Étape 5 : On isole $\boldsymbol{\beta}$.

Solution fermée de la régression Ridge

$$\hat{\boldsymbol{\beta}}_{\text{Ridge}} = (\mathbf{X}^\top \mathbf{X} + \lambda' \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

où $\lambda' = n\lambda$ (dans la pratique, les bibliothèques absorbent le facteur n dans λ).

Comparaison avec la régression classique :

$$\hat{\boldsymbol{\beta}}_{\text{OLS}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

La seule différence est l'ajout de $\lambda\mathbf{I}$ à la matrice $\mathbf{X}^\top \mathbf{X}$!

Cas limites : comprendre λ

Que se passe-t-il quand λ varie ?

Cas 1 : $\lambda = 0$ (pas de régularisation)

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X} + 0 \cdot \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

On retrouve la régression linéaire classique (OLS). Aucune contrainte sur les coefficients.

Cas 2 : $\lambda \rightarrow \infty$ (régularisation maximale)

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X} + \infty \cdot \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y} \approx \frac{1}{\lambda} \mathbf{I}^{-1} \mathbf{X}^\top \mathbf{y} \rightarrow \mathbf{0}$$

Tous les coefficients tendent vers zéro. Le modèle ne prédit plus rien d'utile (sous-apprentissage / underfitting).

Le but : trouver le λ qui donne le meilleur compromis entre les deux.

Interprétation géométrique de Ridge

La régression Ridge peut aussi s'écrire comme un problème d'optimisation sous contrainte :

$$\min_{\beta} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \text{sous la contrainte} \quad \sum_{j=1}^p \beta_j^2 \leq t$$

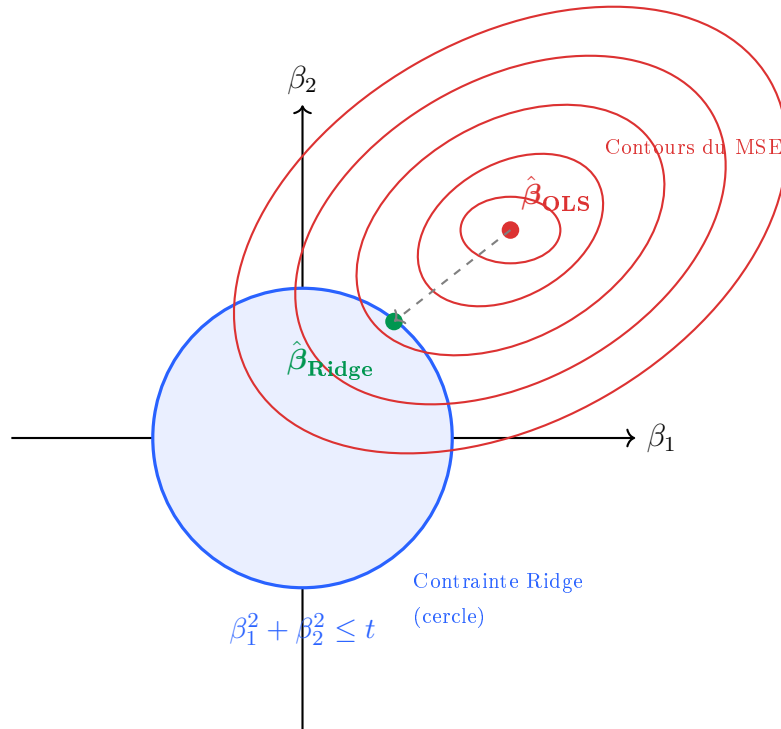


FIGURE 8.3 – Interprétation géométrique de Ridge : le point optimal est là où la plus petite ellipse du MSE touche le cercle de contrainte. Les coefficients sont réduits mais **jamais exactement zéro**.

8.3.3 Régression Lasso (pénalité L1)

Régression Lasso — Définition

La régression Lasso (*Least Absolute Shrinkage and Selection Operator*) utilise une pénalité proportionnelle à la **valeur absolue** de chaque coefficient :

$$J_{\text{Lasso}}(\beta) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

Où $\sum_{j=1}^p |\beta_j| = |\beta_1| + |\beta_2| + \dots + |\beta_p|$ est la **norme L1** des coefficients.

La propriété magique du Lasso : la sélection de variables

Différence capitale entre Ridge et Lasso

- **Ridge** : réduit les coefficients vers zéro, mais ils ne sont **jamais exactement égaux à zéro**.
- **Lasso** : peut mettre certains coefficients **exactement à zéro**, ce qui élimine complètement ces variables du modèle.

Le Lasso fait donc simultanément de la **régularisation** et de la **sélection de variables** !

Pourquoi le Lasso met des coefficients à zéro ? — Interprétation géométrique

La régression Lasso s'écrit aussi comme un problème d'optimisation sous contrainte :

$$\min_{\beta} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \text{sous la contrainte} \quad \sum_{j=1}^p |\beta_j| \leq t$$

La région de contrainte est un **losange** (diamant) au lieu d'un cercle :

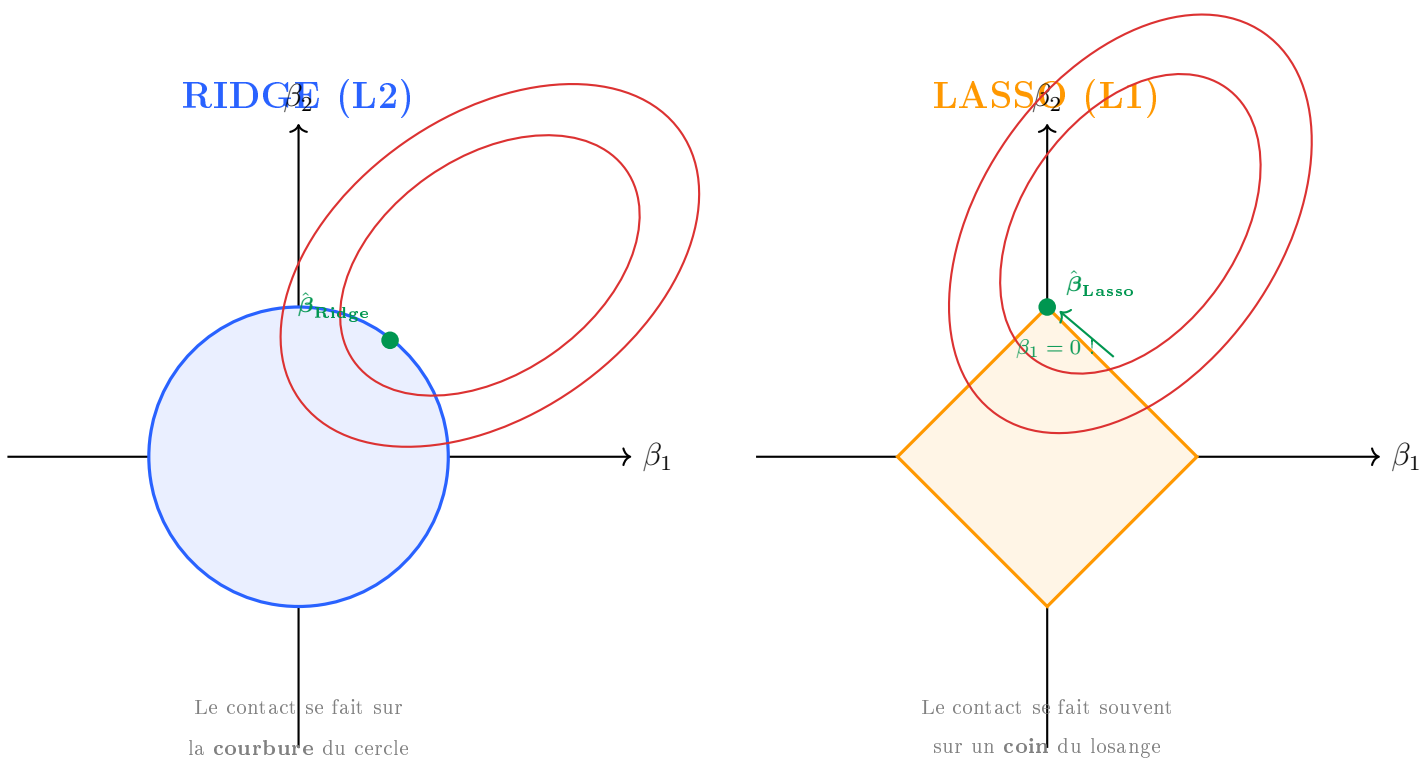


FIGURE 8.4 – Ridge vs Lasso : le losange a des **coins sur les axes**, ce qui force certains coefficients à être exactement zéro. Le cercle n'a pas de coins : les coefficients sont réduits mais jamais nuls.

Pourquoi les coins ?

Le losange du Lasso a des **coins pointus** situés sur les axes (là où $\beta_1 = 0$ ou $\beta_2 = 0$). Quand les ellipses du MSE (qui représentent le coût) s'agrandissent à partir de la solution OLS, elles touchent en premier le losange... et **il y a de fortes chances que le premier contact se fasse sur un coin**. Un coin sur l'axe β_1 signifie $\beta_1 = 0$: la variable 1 est éliminée !

Le cercle de Ridge, lui, est « lisse » partout : le point de contact ne tombe presque jamais pile sur un axe.

Pas de solution fermée pour le Lasso

Lasso : solution itérative

Contrairement à Ridge, la valeur absolue $|\beta_j|$ n'est **pas dérivable** en $\beta_j = 0$. On ne peut donc pas écrire une solution fermée en une seule formule.

Le Lasso utilise des algorithmes itératifs, notamment la **descente de coordonnées** (*coordinate descent*) : on optimise un coefficient à la fois, en gardant les autres fixes, et on répète jusqu'à convergence.

En pratique, **scikit-learn** gère tout cela automatiquement !

8.3.4 ElasticNet : le meilleur des deux mondes

Régression ElasticNet — Définition

L'ElasticNet combine les pénalités L1 (Lasso) et L2 (Ridge) :

$$J_{\text{ElasticNet}}(\beta) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda_1 \sum_{j=1}^p |\beta_j| + \lambda_2 \sum_{j=1}^p \beta_j^2$$

Ou, de manière équivalente avec un paramètre de mélange $\alpha \in [0, 1]$ et une force globale λ :

$$J_{\text{ElasticNet}}(\beta) = \text{MSE} + \lambda \left[\alpha \sum_{j=1}^p |\beta_j| + \frac{(1 - \alpha)}{2} \sum_{j=1}^p \beta_j^2 \right]$$

Où :

- $\alpha = 1$: on retrouve le **Lasso** (L1 pur)
- $\alpha = 0$: on retrouve la **Ridge** (L2 pur)
- $0 < \alpha < 1$: un **mélange** des deux

Quand utiliser chaque méthode ?

Méthode	Utiliser quand...	Propriété clé
Ridge	Toutes les variables sont potentiellement utiles	Réduit les coefficients sans les éliminer
Lasso	On suspecte que beaucoup de variables sont inutiles	Élimine les variables non pertinentes
ElasticNet	Variables corrélées entre elles + sélection souhaitée	Combine sélection et stabilité

8.3.5 Choisir λ : la validation croisée

Le paramètre λ est un **hyperparamètre** : il n'est pas appris par le modèle, il doit être choisi par nous. La méthode standard est la **validation croisée** (*cross-validation*).

Choix de λ par validation croisée

1. Définir une grille de valeurs candidates : $\lambda \in \{0.001, 0.01, 0.1, 1, 10, 100\}$
2. Pour chaque λ , effectuer une validation croisée à k plis (*k-fold CV*) :
 - Diviser les données d'entraînement en k parties égales
 - Entraîner sur $k - 1$ parties, évaluer sur la partie restante
 - Répéter k fois (chaque partie sert une fois de test)
 - Calculer l'erreur moyenne sur les k évaluations
3. Choisir le λ qui donne la plus petite erreur moyenne

Erreur d'entraînement vs erreur de test en fonction de λ

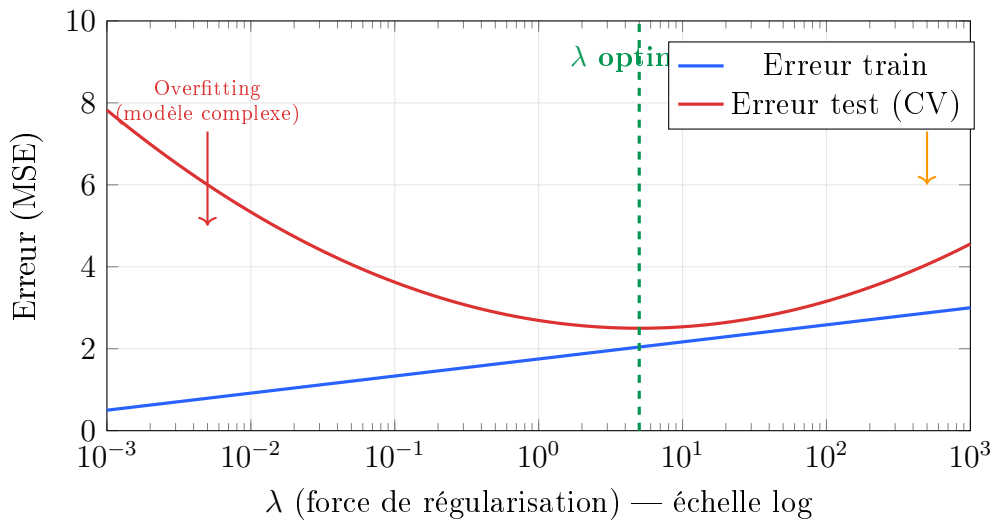


FIGURE 8.5 – En augmentant λ , l'erreur d'entraînement augmente (modèle plus simple) mais l'erreur de test diminue d'abord puis augmente. Le λ optimal minimise l'erreur de test.

8.4 Application sur le dataset clientèle

Reprenons notre dataset clientèle et montrons comment la régularisation résout le problème d'overfitting.

Problème : variables de bruit qui dégradent le modèle

Notre dataset possède des variables utiles (revenu, ancienneté, nombre d'achats) mais aussi des variables aléatoires ajoutées par erreur. La régression classique utilise **toutes** les variables, y compris le bruit, et sur-apprend.

```

1 import numpy as np
2 import pandas as pd
3 from sklearn.linear_model import (LinearRegression, Ridge, Lasso,
4                                   ElasticNet)
5 from sklearn.model_selection import train_test_split
6 from sklearn.preprocessing import StandardScaler
7
8 # Generer le dataset clientele
9 np.random.seed(42)
10 n = 100
11
12 # Variables utiles
13 revenu = np.random.uniform(20000, 80000, n)
14 anciennete = np.random.uniform(1, 15, n)

```

```

15 nb_achats = np.random.randint(1, 50, n)
16
17 # Variable cible
18 depense = 0.02 * revenu + 50 * anciennete + 10 * nb_achats \
19         + np.random.randn(n) * 100
20
21 # Ajouter 15 variables de bruit !
22 bruit = np.random.randn(n, 15) * 1000
23
24 # Construire X
25 X = np.column_stack([revenu, anciennete, nb_achats, bruit])
26 noms = ['Revenu', 'Anciennete', 'Nb_achats'] + \
27         [f'Bruit_{i}' for i in range(1, 16)]
28
29 # Standardiser (important pour la regularisation !)
30 scaler = StandardScaler()
31 X_scaled = scaler.fit_transform(X)
32
33 # Separer train / test
34 X_train, X_test, y_train, y_test = train_test_split(
35     X_scaled, depense, test_size=0.3, random_state=42)
36
37 # Comparer les modeles
38 modeles = {
39     'OLS (sans regularisation)': LinearRegression(),
40     'Ridge (alpha=1.0)': Ridge(alpha=1.0),
41     'Lasso (alpha=1.0)': Lasso(alpha=1.0),
42     'ElasticNet (alpha=0.5)': ElasticNet(alpha=0.5, l1_ratio=0.5),
43 }
44
45 print(f"{'Modele':<30} {'R2 Train':>10} {'R2 Test':>10}")
46 print("-" * 52)
47 for nom, model in modeles.items():
48     model.fit(X_train, y_train)
49     r2_train = model.score(X_train, y_train)
50     r2_test = model.score(X_test, y_test)
51     print(f"{nom:<30} {r2_train:>10.4f} {r2_test:>10.4f}")

```

Listing 8.2 – Comparaison Ridge, Lasso, ElasticNet sur le dataset client

```

1 import matplotlib.pyplot as plt
2

```

```

3  # Afficher les coefficients
4  fig, axes = plt.subplots(2, 2, figsize=(14, 10))
5
6  for ax, (nom, model) in zip(axes.flatten(), modeles.items()):
7      coefs = model.coef_ if hasattr(model, 'coef_') else [0]*18
8      colors = ['green' if i < 3 else 'red' for i in
9                 range(len(coefs))]
9      ax.bar(range(len(coefs)), coefs, color=colors, alpha=0.7)
10     ax.set_title(nom, fontsize=12, fontweight='bold')
11     ax.set_xlabel('Variable')
12     ax.set_ylabel('Coefficient')
13     ax.axhline(y=0, color='black', linewidth=0.5)
14     nb_nuls = sum(1 for c in coefs if abs(c) < 0.01)
15     ax.text(0.95, 0.95, f'Coefficients ~0 : {nb_nuls}',
16            transform=ax.transAxes, ha='right', va='top',
17            fontsize=10, bbox=dict(boxstyle='round',
18                                   facecolor='wheat'))
18
19 plt.suptitle('Coefficients des modeles\n'
20             '(vert = variables utiles, rouge = bruit)',
21             fontsize=14, fontweight='bold')
22 plt.tight_layout()
23 plt.show()

```

Listing 8.3 – Comparer les coefficients : le Lasso élimine le bruit

Résultat clé à retenir

- **OLS** : utilise toutes les variables (y compris le bruit) \Rightarrow overfitting
- **Ridge** : réduit les coefficients du bruit mais ne les élimine pas
- **Lasso** : met les coefficients de bruit **exactement à zéro** \Rightarrow sélection automatique !
- **ElasticNet** : compromis entre Ridge et Lasso

8.5 Implémentation complète sur Google Colab

Notebook Colab — Régularisation complète

Copiez et exécutez chaque cellule dans un notebook Google Colab.

```
1  # =====
```

```

2  # CHAPITRE 8 : REGULARISATION
3  # Ridge, Lasso et ElasticNet
4  # =====
5
6  import numpy as np
7  import pandas as pd
8  import matplotlib.pyplot as plt
9  from sklearn.linear_model import (LinearRegression, Ridge, Lasso,
10                                  ElasticNet, RidgeCV, LassoCV,
11                                  ElasticNetCV)
12  from sklearn.model_selection import train_test_split
13  from sklearn.preprocessing import StandardScaler
14  from sklearn.metrics import r2_score, mean_squared_error
15
16  # Generer le dataset
17  np.random.seed(42)
18  n = 100
19  X_utile = np.random.randn(n, 3)
20  y = 5*X_utile[:,0] + 3*X_utile[:,1] - 2*X_utile[:,2] \
21      + np.random.randn(n)*0.5
22  X_bruit = np.random.randn(n, 17)
23  X = np.hstack([X_utile, X_bruit])
24  noms = [f'Utile_{i+1}' for i in range(3)] + \
25          [f'Bruit_{i+1}' for i in range(17)]
26
27  # Standardiser et separer
28  scaler = StandardScaler()
29  X_scaled = scaler.fit_transform(X)
30  X_train, X_test, y_train, y_test = train_test_split(
31      X_scaled, y, test_size=0.3, random_state=42)
32
33  print(f"Nombre de variables : {X.shape[1]}")
34  print(f"Taille train : {X_train.shape[0]}")
35  print(f"Taille test : {X_test.shape[0]}")

```

Listing 8.4 – Cellule 1 — Imports et génération des données

```

1  # =====
2  # SELECTION AUTOMATIQUE DE LAMBDA
3  # avec validation croisee (CV)
4  # =====
5

```

```

6  # RidgeCV : teste automatiquement plusieurs alpha
7  alphas = np.logspace(-4, 4, 100)
8
9  ridge_cv = RidgeCV(alphas=alphas, cv=5)
10 ridge_cv.fit(X_train, y_train)
11 print(f"Ridge - Meilleur alpha : {ridge_cv.alpha_:.4f}")
12 print(f"Ridge - R2 Train : {ridge_cv.score(X_train, y_train):.4f}")
13 print(f"Ridge - R2 Test : {ridge_cv.score(X_test, y_test):.4f}")
14
15 # LassoCV : selection automatique
16 lasso_cv = LassoCV(alphas=alphas, cv=5, max_iter=10000)
17 lasso_cv.fit(X_train, y_train)
18 print(f"\nLasso - Meilleur alpha : {lasso_cv.alpha_:.4f}")
19 print(f"Lasso - R2 Train : {lasso_cv.score(X_train, y_train):.4f}")
20 print(f"Lasso - R2 Test : {lasso_cv.score(X_test, y_test):.4f}")
21
22 # ElasticNetCV
23 enet_cv = ElasticNetCV(l1_ratio=[0.1, 0.3, 0.5, 0.7, 0.9],
24                        alphas=alphas, cv=5, max_iter=10000)
25 enet_cv.fit(X_train, y_train)
26 print(f"\nElasticNet - Meilleur alpha : {enet_cv.alpha_:.4f}")
27 print(f"ElasticNet - Meilleur l1_ratio : {enet_cv.l1_ratio_:.2f}")
28 print(f"ElasticNet - R2 Train : "
29       f"{enet_cv.score(X_train, y_train):.4f}")
30 print(f"ElasticNet - R2 Test : "
31       f"{enet_cv.score(X_test, y_test):.4f}")

```

Listing 8.5 – Cellule 2 — Ridge, Lasso, ElasticNet avec choix automatique de λ

```

1  # =====
2  # COMPARAISON DES COEFFICIENTS
3  # =====
4  fig, axes = plt.subplots(1, 3, figsize=(18, 5))
5
6  modeles_cv = {'Ridge': ridge_cv, 'Lasso': lasso_cv,
7               'ElasticNet': enet_cv}
8
9  for ax, (nom, model) in zip(axes, modeles_cv.items()):
10     coefs = model.coef_
11     colors = ['#2ecc71' if i < 3 else '#e74c3c'
12              for i in range(len(coefs))]
13     bars = ax.bar(range(len(coefs)), coefs, color=colors,

```

```

        alpha=0.8)
14     ax.set_title(f'{nom} (alpha={model.alpha_:.4f})',
15                 fontsize=12, fontweight='bold')
16     ax.set_xlabel('Indice de la variable')
17     ax.set_ylabel('Coefficient')
18     ax.axhline(y=0, color='black', linewidth=0.5)
19     ax.set_xticks(range(0, 20, 2))
20     nb_zero = sum(1 for c in coefs if abs(c) < 0.01)
21     ax.text(0.95, 0.95, f'Coefs nuls : {nb_zero}/20',
22            transform=ax.transAxes, ha='right', va='top',
23            fontsize=10,
24            bbox=dict(boxstyle='round', facecolor='lightyellow'))
25
26 plt.suptitle('Comparaison des coefficients\n'
27             '(vert = utiles, rouge = bruit)',
28             fontsize=14, fontweight='bold')
29 plt.tight_layout()
30 plt.show()

```

Listing 8.6 – Cellule 3 — Comparaison des coefficients (diagramme en barres)

```

1  # =====
2  # CHEMINS DE COEFFICIENTS (Coefficient Paths)
3  # =====
4  alphas_path = np.logspace(-3, 3, 200)
5
6  # Ridge path
7  coefs_ridge = []
8  for a in alphas_path:
9      model = Ridge(alpha=a)
10     model.fit(X_train, y_train)
11     coefs_ridge.append(model.coef_)
12 coefs_ridge = np.array(coefs_ridge)
13
14 # Lasso path
15 coefs_lasso = []
16 for a in alphas_path:
17     model = Lasso(alpha=a, max_iter=10000)
18     model.fit(X_train, y_train)
19     coefs_lasso.append(model.coef_)
20 coefs_lasso = np.array(coefs_lasso)
21

```

```

22 fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(16, 6))
23
24 # Ridge
25 for j in range(20):
26     color = 'green' if j < 3 else 'lightcoral'
27     lw = 2.5 if j < 3 else 0.8
28     ax1.plot(alphas_path, coefs_ridge[:, j],
29             color=color, linewidth=lw,
30             label=noms[j] if j < 3 else None)
31 ax1.set_xscale('log')
32 ax1.set_xlabel('alpha (lambda)', fontsize=12)
33 ax1.set_ylabel('Coefficient', fontsize=12)
34 ax1.set_title('Ridge : chemins des coefficients', fontsize=14,
35             fontweight='bold')
36 ax1.legend(fontsize=10)
37 ax1.axhline(y=0, color='black', linewidth=0.5)
38 ax1.grid(True, alpha=0.3)
39
40 # Lasso
41 for j in range(20):
42     color = 'green' if j < 3 else 'lightcoral'
43     lw = 2.5 if j < 3 else 0.8
44     ax2.plot(alphas_path, coefs_lasso[:, j],
45             color=color, linewidth=lw,
46             label=noms[j] if j < 3 else None)
47 ax2.set_xscale('log')
48 ax2.set_xlabel('alpha (lambda)', fontsize=12)
49 ax2.set_ylabel('Coefficient', fontsize=12)
50 ax2.set_title('Lasso : chemins des coefficients', fontsize=14,
51             fontweight='bold')
52 ax2.legend(fontsize=10)
53 ax2.axhline(y=0, color='black', linewidth=0.5)
54 ax2.grid(True, alpha=0.3)
55
56 plt.suptitle('Comment les coefficients evoluent quand on augmente '
57             'la regularisation', fontsize=14)
58 plt.tight_layout()
59 plt.show()

```

Listing 8.7 – Cellule 4 — Chemins de coefficients en fonction de λ

```

1 # =====

```

```

2  # COMPARAISON R2 TRAIN vs TEST
3  # =====
4  resultats = {}
5  for nom, model in [('OLS', LinearRegression()),
6                     ('Ridge', ridge_cv),
7                     ('Lasso', lasso_cv),
8                     ('ElasticNet', enet_cv)]:
9      if nom == 'OLS':
10         model.fit(X_train, y_train)
11         r2_tr = model.score(X_train, y_train)
12         r2_te = model.score(X_test, y_test)
13         resultats[nom] = {'R2 Train': r2_tr, 'R2 Test': r2_te}
14
15 df_res = pd.DataFrame(resultats).T
16 print(df_res.to_string(float_format='{:.4f}'.format))
17
18 # Graphique
19 x_pos = np.arange(len(df_res))
20 width = 0.35
21
22 fig, ax = plt.subplots(figsize=(10, 6))
23 bars1 = ax.bar(x_pos - width/2, df_res['R2 Train'],
24               width, label='R2 Train', color='#3498db', alpha=0.8)
25 bars2 = ax.bar(x_pos + width/2, df_res['R2 Test'],
26               width, label='R2 Test', color='#e74c3c', alpha=0.8)
27
28 ax.set_xlabel('Modele', fontsize=12)
29 ax.set_ylabel('R2 Score', fontsize=12)
30 ax.set_title('Comparaison R2 Train vs Test\n'
31             'La regularisation reduit l\'overfitting',
32             fontsize=14, fontweight='bold')
33 ax.set_xticks(x_pos)
34 ax.set_xticklabels(df_res.index, fontsize=11)
35 ax.legend(fontsize=11)
36 ax.set_ylim(0, 1.1)
37 ax.grid(axis='y', alpha=0.3)
38
39 # Ajouter les valeurs sur les barres
40 for bar in bars1:
41     ax.text(bar.get_x() + bar.get_width()/2., bar.get_height() +
42           0.02,

```

```

42         f'{bar.get_height():.3f}', ha='center', va='bottom',
43         fontsize=9)
44 for bar in bars2:
45     ax.text(bar.get_x() + bar.get_width()/2., bar.get_height() +
46             0.02,
47             f'{bar.get_height():.3f}', ha='center', va='bottom',
48             fontsize=9)
49 plt.tight_layout()
50 plt.show()

```

Listing 8.8 – Cellule 5 — Comparaison finale R^2 train vs test

8.6 Schéma récapitulatif : phases d'entraînement et de test

PHASE D'ENTRAÎNEMENT (Training)

1. Étape 1 — Charger et séparer les données

Charger le jeu de données et le diviser en 80 % Train et 20 % Test.

Code Python :

```
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2)
```

2. Étape 2 — Normaliser (fit sur Train uniquement)

Centrer et réduire les variables pour que la pénalité traite toutes les variables équitablement.

Mathématiques :

$$x'_j = \frac{x_j - \mu_j}{\sigma_j} \text{ pour chaque variable } j$$

Code Python :

```
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
```

3. Étape 3 — Choisir le type de régularisation

Sélectionner Ridge, Lasso ou ElasticNet selon le contexte.

Mathématiques :

$$\text{Ridge : } J(\beta) = \sum (y_i - \hat{y}_i)^2 + \lambda \sum \beta_j^2$$

$$\text{Lasso} : J(\beta) = \sum (y_i - \hat{y}_i)^2 + \lambda \sum |\beta_j|$$

$$\text{ElasticNet} : J(\beta) = \sum (y_i - \hat{y}_i)^2 + \lambda_1 \sum |\beta_j| + \lambda_2 \sum \beta_j^2$$

Code Python :

```
from sklearn.linear_model import Ridge, Lasso, ElasticNet
```

4. Étape 4 — Trouver λ par validation croisée

Utiliser une recherche par grille pour sélectionner le λ optimal.

Mathématiques :

$$\lambda^* = \arg \min_{\lambda} \text{CV-MSE}(\lambda)$$

Code Python :

```
from sklearn.model_selection import GridSearchCV
params = {'alpha': [0.01, 0.1, 1, 10, 100]}
grid = GridSearchCV(Ridge(), params, cv=5)
grid.fit(X_train_scaled, y_train)
```

5. Étape 5 — Entraîner le modèle final

Ajuster le modèle avec le meilleur λ sur l'ensemble du jeu Train.

Code Python :

```
best_model = grid.best_estimator_
print("Meilleur lambda :", grid.best_params_)
```



Modèle entraîné transmis à la phase de test

PHASE DE TEST (Testing)

1. Étape 1 — Normaliser X_{test} (transform, PAS fit)

Appliquer le scaler ajusté sur Train aux données de test.

Code Python :

```
X_test_scaled = scaler.transform(X_test) # PAS fit_transform !
```

2. Étape 2 — Prédire

Calculer les prédictions avec les coefficients régularisés.

Mathématiques :

$$\hat{y} = X_{\text{test}} \cdot \hat{\beta}_{\text{reg}}$$

Code Python :

```
y_pred = best_model.predict(X_test_scaled)
```

3. Étape 3 — Évaluer

Mesurer la qualité du modèle sur les données de test.

Mathématiques :

$$\text{MSE} = \frac{1}{n} \sum (y_i - \hat{y}_i)^2, \quad R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

Code Python :

```
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)
```

4. Étape 4 — Analyser les coefficients

Examiner quels coefficients ont été réduits ou mis à zéro.

Mathématiques :

Ridge : $|\beta_j|$ réduits mais $\neq 0$. Lasso : certains $\beta_j = 0$ (sélection)

Code Python :

```
for name, coef in zip(X.columns, best_model.coef_):
    print(f"{name}: {coef:.4f}")
```

8.7 Exercices

Exercice 8.1 — Minimiser une fonction de coût Ridge (à la main)

On considère un problème de régression très simple avec **une seule variable** x et **un seul coefficient** β (pas d'intercept). La fonction de coût Ridge est :

$$J(\beta) = \frac{1}{n} \sum_{i=1}^n (y_i - \beta x_i)^2 + \lambda \beta^2$$

On dispose de $n = 3$ observations :

i	x_i	y_i
1	1	2
2	2	5
3	3	7

1. Écrivez la fonction de coût $J(\beta)$ en développant la somme (pour $\lambda = 0.5$).
2. Calculez la dérivée $\frac{dJ}{d\beta}$ et trouvez le β qui minimise J .
3. Comparez avec la solution OLS (sans régularisation, $\lambda = 0$).

4. Que constatez-vous ? Quel est l'effet de la régularisation ?

Correction de l'exercice 8.1

1. Développement de la fonction de coût

Développons terme par terme :

$$J(\beta) = \frac{1}{3} [(2 - \beta \cdot 1)^2 + (5 - \beta \cdot 2)^2 + (7 - \beta \cdot 3)^2] + 0.5 \beta^2 \quad (8.5)$$

$$= \frac{1}{3} [(2 - \beta)^2 + (5 - 2\beta)^2 + (7 - 3\beta)^2] + 0.5 \beta^2 \quad (8.6)$$

Développons chaque carré :

$$(2 - \beta)^2 = 4 - 4\beta + \beta^2 \quad (8.7)$$

$$(5 - 2\beta)^2 = 25 - 20\beta + 4\beta^2 \quad (8.8)$$

$$(7 - 3\beta)^2 = 49 - 42\beta + 9\beta^2 \quad (8.9)$$

La somme donne :

$$\text{Somme} = 78 - 66\beta + 14\beta^2$$

Donc :

$$J(\beta) = \frac{78 - 66\beta + 14\beta^2}{3} + 0.5 \beta^2 = 26 - 22\beta + \frac{14}{3}\beta^2 + 0.5 \beta^2$$

$$J(\beta) = 26 - 22\beta + \left(\frac{14}{3} + \frac{1}{2}\right) \beta^2 = 26 - 22\beta + \frac{31}{6}\beta^2$$

2. Dérivation et minimisation

$$\frac{dJ}{d\beta} = -22 + 2 \times \frac{31}{6} \beta = -22 + \frac{31}{3} \beta$$

On pose $\frac{dJ}{d\beta} = 0$:

$$\frac{31}{3} \beta = 22 \quad \Rightarrow \quad \boxed{\hat{\beta}_{\text{Ridge}} = \frac{22 \times 3}{31} = \frac{66}{31} \approx 2.129}$$

3. Comparaison avec OLS ($\lambda = 0$)

Sans régularisation :

$$J_{\text{OLS}}(\beta) = 26 - 22\beta + \frac{14}{3}\beta^2$$

$$\frac{dJ_{\text{OLS}}}{d\beta} = -22 + \frac{28}{3}\beta = 0 \quad \Rightarrow \quad \boxed{\hat{\beta}_{\text{OLS}} = \frac{66}{28} = \frac{33}{14} \approx 2.357}$$

4. Conclusion

Le coefficient Ridge (≈ 2.129) est **plus petit** que le coefficient OLS (≈ 2.357). La régula-

risation a **réduit** (“shrinké”) le coefficient vers zéro, ce qui est exactement son rôle. Plus λ est grand, plus $\hat{\beta}$ sera petit.

On peut aussi retrouver ce résultat avec la formule fermée. Avec une seule variable, $\mathbf{X}^\top \mathbf{X} = 1^2 + 2^2 + 3^2 = 14$ et $\mathbf{X}^\top \mathbf{y} = 1 \times 2 + 2 \times 5 + 3 \times 7 = 33$. D’où :

$$\hat{\beta}_{\text{Ridge}} = \frac{X^\top y}{X^\top X + n\lambda} = \frac{33}{14 + 3 \times 0.5} = \frac{33}{15.5} = \frac{66}{31} \approx 2.129 \quad \checkmark$$

Exercice 8.2 — Pourquoi le Lasso produit des solutions creuses ? (conceptuel)

1. Dessinez (ou décrivez) la différence géométrique entre la contrainte L1 ($|\beta_1| + |\beta_2| \leq t$) et la contrainte L2 ($\beta_1^2 + \beta_2^2 \leq t$) dans le plan (β_1, β_2) .
2. Expliquez avec vos propres mots pourquoi le losange de la contrainte L1 a tendance à produire des solutions où certains $\beta_j = 0$.
3. Donnez un exemple concret où la propriété de sélection de variables du Lasso est utile.
4. Si on a un dataset avec 1000 variables mais qu’on pense que seules 10 sont vraiment importantes, quel modèle choisir entre Ridge, Lasso et ElasticNet ? Justifiez.

Correction de l’exercice 8.2

1. Différence géométrique

- La contrainte **L2** définit un **cercle** (ou une hypersphère en dimension supérieure) centré en zéro. La frontière est lisse partout, sans coins.
- La contrainte **L1** définit un **losange** (ou un hyperoctaèdre) centré en zéro. Les coins se trouvent **sur les axes**, c’est-à-dire aux points où un ou plusieurs $\beta_j = 0$.

2. Pourquoi le losange produit des solutions creuses

L’optimisation revient à trouver le point de la région de contrainte le plus proche de la solution OLS (non contrainte). Géométriquement, on agrandit les ellipses du MSE jusqu’à ce qu’elles « touchent » la région de contrainte.

Pour le **cercle** (Ridge), le premier contact se fait généralement en un point lisse de la courbure, où aucun coefficient n’est zéro.

Pour le **losange** (Lasso), les **coins pointus** sur les axes sont des points “saillants”. La probabilité que l’ellipse touche le losange en un coin est élevée, surtout en haute dimension (avec beaucoup de variables). Un coin sur un axe correspond exactement à une situation où un ou plusieurs $\beta_j = 0$.

3. Exemple concret

En **génomique**, on peut avoir des milliers de gènes (variables) mais seuls quelques-uns

sont liés à une maladie. Le Lasso permet de sélectionner automatiquement ces quelques gènes pertinents parmi les milliers de variables. Autre exemple : en marketing, parmi des dizaines de variables clients, le Lasso identifie les 3-4 facteurs vraiment déterminants pour la dépense.

4. Choix du modèle pour 1000 variables, 10 utiles

Le **Lasso** ou l'**ElasticNet** sont les meilleurs choix :

- Le **Lasso** éliminera les 990 variables inutiles en mettant leurs coefficients à zéro. C'est le choix idéal si les 10 variables utiles ne sont pas trop corrélées entre elles.
- L'**ElasticNet** est préférable si les 10 variables utiles sont **corrélées entre elles**. En effet, le Lasso a tendance à ne garder qu'une seule variable parmi un groupe de variables corrélées, tandis que l'ElasticNet les garde toutes (grâce à sa composante Ridge).
- La **Ridge** ne convient pas ici car elle ne met aucun coefficient à zéro. Le modèle final utiliserait les 1000 variables, ce qui est difficile à interpréter.

Exercice 8.3 — Python : Ridge vs Lasso sur des variables non pertinentes

Écrivez un code Python qui :

1. Génère un dataset avec 5 variables utiles et 25 variables de bruit ($n = 150$).
2. Sépare en train (70%) et test (30%).
3. Entraîne les modèles **LassoCV** et **RidgeCV** avec validation croisée.
4. Affiche le R^2 sur les données de test pour chaque modèle.
5. Affiche le nombre de coefficients mis à zéro par le Lasso.
6. Crée un graphique en barres comparant les coefficients des deux modèles.

Correction de l'exercice 8.3

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.linear_model import RidgeCV, LassoCV
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler

# 1. Générer les données
np.random.seed(123)
n = 150
```

```

X_utile = np.random.randn(n, 5)
y = (3*X_utile[:,0] - 2*X_utile[:,1] + 4*X_utile[:,2]
      + 1.5*X_utile[:,3] - 3*X_utile[:,4]
      + np.random.randn(n) * 0.5)
X_bruit = np.random.randn(n, 25)
X = np.hstack([X_utile, X_bruit])

# 2. Standardiser et separer
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)
X_train, X_test, y_train, y_test = train_test_split(
    X_scaled, y, test_size=0.3, random_state=42)

# 3. Entraîner avec CV
alphas = np.logspace(-4, 4, 100)

ridge = RidgeCV(alphas=alphas, cv=5)
ridge.fit(X_train, y_train)

lasso = LassoCV(alphas=alphas, cv=5, max_iter=10000)
lasso.fit(X_train, y_train)

# 4. Scores
print(f"Ridge - R2 test : {ridge.score(X_test, y_test):.4f}")
print(f"Lasso - R2 test : {lasso.score(X_test, y_test):.4f}")
print(f"Ridge - alpha optimal : {ridge.alpha_:.4f}")
print(f"Lasso - alpha optimal : {lasso.alpha_:.4f}")

# 5. Coefficients nuls
nb_zero_lasso = sum(1 for c in lasso.coef_ if abs(c) < 0.001)
nb_zero_ridge = sum(1 for c in ridge.coef_ if abs(c) < 0.001)
print(f"\nCoefficients nuls (Lasso) : {nb_zero_lasso}/30")
print(f"Coefficients nuls (Ridge) : {nb_zero_ridge}/30")

# 6. Graphique comparatif
fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(16, 6))
noms = [f'U{i+1}' for i in range(5)] + \
        [f'B{i+1}' for i in range(25)]
x_pos = np.arange(30)
colors = ['#2ecc71']*5 + ['#e74c3c']*25

```

```
ax1.bar(x_pos, ridge.coef_, color=colors, alpha=0.7)
ax1.set_title(f'Ridge (alpha={ridge.alpha_:.4f})',
              fontsize=13, fontweight='bold')
ax1.set_xlabel('Variable')
ax1.set_ylabel('Coefficient')
ax1.axhline(y=0, color='black', linewidth=0.5)
ax1.set_xticks(range(0, 30, 3))
ax1.grid(axis='y', alpha=0.3)

ax2.bar(x_pos, lasso.coef_, color=colors, alpha=0.7)
ax2.set_title(f'Lasso (alpha={lasso.alpha_:.4f})',
              fontsize=13, fontweight='bold')
ax2.set_xlabel('Variable')
ax2.set_ylabel('Coefficient')
ax2.axhline(y=0, color='black', linewidth=0.5)
ax2.set_xticks(range(0, 30, 3))
ax2.grid(axis='y', alpha=0.3)

plt.suptitle('Ridge garde tous les coefficients non nuls\n'
             'Lasso elimine les variables de bruit',
             fontsize=14, fontweight='bold')
plt.tight_layout()
plt.show()

print("\n=== Conclusion ===")
print("Le Lasso a correctement identifie les 5 variables "
      "utiles en mettant les 25 variables de bruit a zero.")
print("Ridge a reduit les coefficients de bruit mais ne les "
      "a pas elimines completement.")
```

Chapitre 9

Métriques d'Évaluation pour la Régression

9.1 Introduction : pourquoi évaluer un modèle ?

Jusqu'ici, nous avons appris à construire des modèles de régression : linéaire simple, multiple, polynomiale, et avec régularisation. Mais une question fondamentale reste en suspens : **comment savoir si notre modèle est bon ?**

Imaginez un médecin qui prescrit un traitement sans jamais vérifier si le patient guérit. Ce serait absurde ! De la même manière, entraîner un modèle sans mesurer sa performance n'a aucun sens. Les **métriques d'évaluation** sont nos outils de mesure : elles nous disent à quel point les prédictions de notre modèle sont proches de la réalité.

Point clé

Il n'existe pas **une seule** métrique universelle. Chaque métrique mesure un aspect différent de la qualité du modèle. Le choix de la métrique dépend du **contexte** et de ce qui est important pour votre problème.

9.2 Hands-On : quel modèle est le meilleur ?

Hands-On 9.1 — Deux modèles ; un dilemme

Imaginez que vous travaillez dans une agence immobilière. Vous avez demandé à deux collègues de construire chacun un modèle pour prédire le prix des maisons. Voici les résultats sur 3 maisons de test :

Maison	Prix réel	Modèle A	Modèle B
Maison 1	200 000 €	205 000 €	195 000 €
Maison 2	320 000 €	310 000 €	340 000 €
Maison 3	148 000 €	150 000 €	155 000 €

Question : Quel modèle est le meilleur ?

À première vue, les deux modèles semblent corrects. Mais lequel choisir ? Calculons les erreurs :

Maison	Erreur Modèle A	Erreur Modèle B
Maison 1	$200 - 205 = -5 \text{ k€}$	$200 - 195 = +5 \text{ k€}$
Maison 2	$320 - 310 = +10 \text{ k€}$	$320 - 340 = -20 \text{ k€}$
Maison 3	$148 - 150 = -2 \text{ k€}$	$148 - 155 = -7 \text{ k€}$

Observation : Le Modèle A fait des erreurs de 5k, 10k et 2k. Le Modèle B fait des erreurs de 5k, 20k et 7k. Le Modèle A semble meilleur... mais **comment le quantifier précisément** ? C'est là qu'interviennent les métriques !

Nous allons maintenant calculer **chaque métrique** sur cet exemple pour comprendre ce qu'elle mesure.

9.3 Les métriques de régression expliquées simplement

Dans toute cette section, on note :

- n : le nombre d'observations (ici $n = 3$)
- y_i : la valeur réelle de l'observation i
- \hat{y}_i : la valeur prédite par le modèle pour l'observation i
- \bar{y} : la moyenne des valeurs réelles $= \frac{1}{n} \sum_{i=1}^n y_i$

9.3.1 MAE — Mean Absolute Error (Erreur Absolue Moyenne)

Définition — MAE

La **MAE** (Mean Absolute Error) est la moyenne des valeurs absolues des erreurs :

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Comment comprendre la MAE ?

La MAE se calcule en **trois étapes simples** :

1. **Calculer chaque erreur** : $y_i - \hat{y}_i$ (différence entre réalité et prédiction)
2. **Prendre la valeur absolue** : $|y_i - \hat{y}_i|$ (rendre chaque erreur positive)

3. **Faire la moyenne** : additionner toutes les erreurs positives et diviser par n

Analogie : Imaginez que vous tirez 3 fléchettes sur une cible. Chaque fléchette atterrit à une certaine distance du centre. La MAE, c'est la **distance moyenne de vos fléchettes au centre**. Plus elle est petite, plus vous êtes précis.

Calcul de la MAE sur notre exemple

Modèle A (erreurs en milliers d'euros) :

$$|y_1 - \hat{y}_1| = |200 - 205| = 5$$

$$|y_2 - \hat{y}_2| = |320 - 310| = 10$$

$$|y_3 - \hat{y}_3| = |148 - 150| = 2$$

$$\text{MAE}_A = \frac{5 + 10 + 2}{3} = \frac{17}{3} \approx 5,67 \text{ k€}$$

Modèle B :

$$|y_1 - \hat{y}_1| = |200 - 195| = 5$$

$$|y_2 - \hat{y}_2| = |320 - 340| = 20$$

$$|y_3 - \hat{y}_3| = |148 - 155| = 7$$

$$\text{MAE}_B = \frac{5 + 20 + 7}{3} = \frac{32}{3} \approx 10,67 \text{ k€}$$

Conclusion : $\text{MAE}_A < \text{MAE}_B$, donc le Modèle A est meilleur selon la MAE. En moyenne, le Modèle A se trompe de 5 670 €, contre 10 670 € pour le Modèle B.

Avantages de la MAE	Inconvénients de la MAE
Très facile à comprendre et à expliquer	Ne pénalise pas spécialement les grosses erreurs
Même unité que la variable cible y	Pas différentiable en zéro (problème pour certains algorithmes d'optimisation)
Robuste aux valeurs aberrantes (outliers)	Traite une erreur de 20k et quatre erreurs de 5k de la même façon

9.3.2 MSE — Mean Squared Error (Erreur Quadratique Moyenne)

Définition — MSE

La **MSE** (Mean Squared Error) est la moyenne des carrés des erreurs :

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Comment comprendre la MSE ?

La MSE se calcule en **trois étapes** :

1. **Calculer chaque erreur** : $y_i - \hat{y}_i$
2. **Mettre au carré** : $(y_i - \hat{y}_i)^2$ — cela rend l'erreur positive ET **amplifie les grosses erreurs**
3. **Faire la moyenne**

Pourquoi mettre au carré ? Trois raisons fondamentales :

- **Éliminer les signes négatifs** : comme la valeur absolue, le carré rend tout positif
- **Pénaliser les grosses erreurs** : une erreur de 10 donne $10^2 = 100$, tandis qu'une erreur de 2 donne $2^2 = 4$. L'erreur de 10 est pénalisée **25 fois plus** (et pas juste 5 fois comme avec la MAE) !
- **Propriétés mathématiques** : la fonction x^2 est différentiable partout, ce qui facilite l'optimisation (descente de gradient)

Calcul de la MSE sur notre exemple

Modèle A (erreurs en milliers d'euros) :

$$(y_1 - \hat{y}_1)^2 = (200 - 205)^2 = (-5)^2 = 25$$

$$(y_2 - \hat{y}_2)^2 = (320 - 310)^2 = (10)^2 = 100$$

$$(y_3 - \hat{y}_3)^2 = (148 - 150)^2 = (-2)^2 = 4$$

$$\text{MSE}_A = \frac{25 + 100 + 4}{3} = \frac{129}{3} = \mathbf{43 \text{ k€}^2}$$

Modèle B :

$$(y_1 - \hat{y}_1)^2 = (200 - 195)^2 = 25$$

$$(y_2 - \hat{y}_2)^2 = (320 - 340)^2 = (-20)^2 = 400$$

$$(y_3 - \hat{y}_3)^2 = (148 - 155)^2 = (-7)^2 = 49$$

$$\text{MSE}_B = \frac{25 + 400 + 49}{3} = \frac{474}{3} = \mathbf{158 \text{ k€}^2}$$

Observation : La MSE du Modèle B (158) est presque **4 fois** celle du Modèle A (43), alors que la MAE n'était que 2 fois plus grande. Pourquoi ? Parce que la grosse erreur de 20k du Modèle B est sévèrement pénalisée par la mise au carré ($20^2 = 400$).

Avantages de la MSE	Inconvénients de la MSE
Pénalise fortement les grosses erreurs	Unité en y^2 (k€^2), difficile à interpréter
Différentiable partout (idéale pour la descente de gradient)	Très sensible aux valeurs aberrantes
Utilisée comme fonction de coût dans la plupart des algorithmes	Une seule erreur extrême peut dominer la métrique

9.3.3 RMSE — Root Mean Squared Error (Racine de l'Erreur Quadratique Moyenne)

Définition — RMSE

La **RMSE** est simplement la racine carrée de la MSE :

$$\text{RMSE} = \sqrt{\text{MSE}} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

Comment comprendre la RMSE ?

La RMSE prend la MSE et **annule l'effet de la mise au carré** en appliquant la racine carrée. Résultat :

- On **retrouve la même unité** que y (des k€ au lieu de k€^2)
- On **garde la pénalisation** des grosses erreurs (l'effet du carré n'est pas complètement annulé)
- On obtient une métrique **interprétable** : « en moyenne, le modèle se trompe d'en-

viron RMSE »

Analogie : Si la MSE est comme mesurer l'aire des carrés d'erreur, la RMSE revient à prendre le côté de ces carrés — on repasse en unité de longueur.

Calcul de la RMSE sur notre exemple

Modèle A :

$$\text{RMSE}_A = \sqrt{\text{MSE}_A} = \sqrt{43} \approx \mathbf{6,56 \text{ k€}}$$

Modèle B :

$$\text{RMSE}_B = \sqrt{\text{MSE}_B} = \sqrt{158} \approx \mathbf{12,57 \text{ k€}}$$

Interprétation : Le Modèle A se trompe en moyenne de 6 560 €, le Modèle B de 12 570 €.

Remarque : Notons que $\text{RMSE} \geq \text{MAE}$ toujours. Pour le Modèle A : $\text{RMSE} = 6,56 > \text{MAE} = 5,67$. L'écart entre RMSE et MAE est d'autant plus grand que les erreurs sont **inégaes** (certaines très grandes, d'autres petites).

RMSE vs MAE : quand choisir laquelle ?

- **Utilisez la MAE** si toutes les erreurs ont la même importance et que vous voulez être robuste aux valeurs aberrantes.
- **Utilisez la RMSE** si les grosses erreurs sont particulièrement graves dans votre contexte (ex : prédiction médicale, finance).
- La **RMSE est la métrique la plus utilisée** en pratique dans les compétitions de data science (Kaggle) et dans la littérature scientifique.

9.3.4 R^2 — Coefficient de Détermination

Définition — R^2

Le **coefficient de détermination** R^2 mesure la proportion de la variance de y expliquée par le modèle :

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

où :

- $SS_{\text{res}} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$ (**somme des carrés résiduels** — erreur du modèle)
- $SS_{\text{tot}} = \sum_{i=1}^n (y_i - \bar{y})^2$ (**somme totale des carrés** — variance totale de y)

Comment comprendre R^2 ?

Le R^2 répond à la question : « **Quelle fraction de la variation dans y mon modèle explique-t-il ?** »

Imaginez la situation suivante :

- **Sans modèle**, la meilleure prédiction serait de toujours prédire la moyenne \bar{y} . L'erreur totale serait SS_{tot} .
- **Avec le modèle**, l'erreur restante est SS_{res} .
- Le rapport $\frac{SS_{\text{res}}}{SS_{\text{tot}}}$ mesure la **fraction d'erreur non expliquée**.
- Donc $R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$ mesure la **fraction d'erreur expliquée**.

Interprétation des valeurs :

- $R^2 = 1$: le modèle est **parfait** (toutes les prédictions sont exactes)
- $R^2 = 0$: le modèle n'est **pas meilleur** que prédire la moyenne
- $0 < R^2 < 1$: le modèle explique une partie de la variance — plus c'est proche de 1, mieux c'est
- $R^2 < 0$: le modèle est **pire** que la moyenne ! (oui, c'est possible !)

$R^2 = 1$ (**parfait**)

$R^2 \approx 0,8$ (**bon**)

$R^2 \approx 0$ (**nul**)

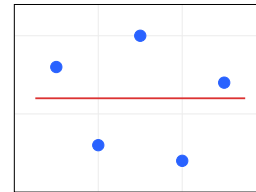
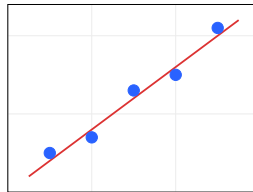
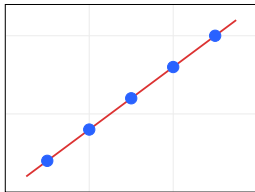


FIGURE 9.1 – Trois situations : R^2 parfait, bon et nul. Plus les points sont proches de la droite, plus R^2 est élevé.

Calcul de R^2 sur notre exemple

Étape 1 : Calculer la moyenne des valeurs réelles :

$$\bar{y} = \frac{200 + 320 + 148}{3} = \frac{668}{3} \approx 222,67 \text{ k€}$$

Étape 2 : Calculer SS_{tot} (variance totale) :

$$\begin{aligned} SS_{\text{tot}} &= (200 - 222,67)^2 + (320 - 222,67)^2 + (148 - 222,67)^2 \\ &= (-22,67)^2 + (97,33)^2 + (-74,67)^2 \\ &= 513,93 + 9473,13 + 5575,61 \\ &= 15\,562,67 \end{aligned}$$

Étape 3 : Calculer SS_{res} pour chaque modèle :

Modèle A :

$$\begin{aligned} SS_{\text{res},A} &= (200 - 205)^2 + (320 - 310)^2 + (148 - 150)^2 \\ &= 25 + 100 + 4 = 129 \end{aligned}$$

Modèle B :

$$\begin{aligned} SS_{\text{res},B} &= (200 - 195)^2 + (320 - 340)^2 + (148 - 155)^2 \\ &= 25 + 400 + 49 = 474 \end{aligned}$$

Étape 4 : Calculer R^2 :

$$R_A^2 = 1 - \frac{129}{15\,562,67} = 1 - 0,0083 \approx \mathbf{0,9917}$$

$$R_B^2 = 1 - \frac{474}{15\,562,67} = 1 - 0,0305 \approx \mathbf{0,9695}$$

Interprétation : Le Modèle A explique 99,17 % de la variance, le Modèle B en explique 96,95 %. Les deux sont très bons, mais le Modèle A est légèrement supérieur.

R^2 peut être négatif!

Un R^2 négatif signifie que votre modèle fait **pire** que simplement prédire la moyenne. Cela arrive quand :

- Le modèle est complètement inadapté aux données
- On évalue sur des données très différentes de l'entraînement
- Le modèle a gravement sur-appris (overfitting)

Exemple : Si un modèle prédit toujours 500k€ pour nos 3 maisons :

$$SS_{\text{res}} = (200 - 500)^2 + (320 - 500)^2 + (148 - 500)^2 = 90\,000 + 32\,400 + 123\,904 = 246\,304$$

$$R^2 = 1 - \frac{246\,304}{15\,562,67} = 1 - 15,83 = -14,83 \quad (\text{très négatif!})$$

9.3.5 R^2 ajusté — Pénaliser les variables inutiles

Définition — R^2 ajusté

Le R^2 ajusté corrige le R^2 en tenant compte du nombre de variables explicatives :

$$R_{\text{adj}}^2 = 1 - \frac{(1 - R^2)(n - 1)}{n - p - 1}$$

où :

- n : nombre d'observations
- p : nombre de variables explicatives (features)
- R^2 : coefficient de détermination classique

Pourquoi le R^2 ajusté est nécessaire ?

Le R^2 classique a un défaut majeur : il **augmente toujours** (ou reste stable) quand on ajoute une nouvelle variable, même si cette variable est **totalelement inutile** !

Analogie : Imaginez un étudiant qui, pour améliorer sa note à un examen, ajoute des pages hors sujet à sa copie. Le nombre de pages augmente, mais la qualité n'augmente pas forcément. Le R^2 classique compte les pages ; le R^2 ajusté évalue la qualité.

Le R^2 ajusté **pénalise** l'ajout de variables qui n'améliorent pas significativement le modèle :

- Si une variable ajoutée améliore le modèle : R_{adj}^2 augmente
- Si une variable ajoutée est inutile : R_{adj}^2 **diminue**

Exemple de l'effet du R^2 ajusté

Supposons que nous avons $n = 100$ observations et un modèle avec $p = 3$ variables donnant $R^2 = 0,85$:

$$R_{\text{adj}}^2 = 1 - \frac{(1 - 0,85)(100 - 1)}{100 - 3 - 1} = 1 - \frac{0,15 \times 99}{96} = 1 - 0,1547 = 0,8453$$

Si on ajoute 10 variables inutiles ($p = 13$), le R^2 reste à 0,86 (légère hausse artificielle) :

$$R_{\text{adj}}^2 = 1 - \frac{(1 - 0,86)(100 - 1)}{100 - 13 - 1} = 1 - \frac{0,14 \times 99}{86} = 1 - 0,1612 = 0,8388$$

Le R^2 a augmenté ($0,85 \rightarrow 0,86$) mais le R_{adj}^2 a **diminué** ($0,8453 \rightarrow 0,8388$), signalant que les variables ajoutées sont inutiles.

9.3.6 MAPE — Mean Absolute Percentage Error (Erreur en Pourcentage)

Définition — MAPE

La **MAPE** mesure l'erreur en pourcentage par rapport à la valeur réelle :

$$\text{MAPE} = \frac{1}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right| \times 100$$

Pourquoi la MAPE ?

La MAPE exprime l'erreur en **pourcentage relatif**. C'est très utile quand les valeurs de y ont des ordres de grandeur différents.

Exemple : Une erreur de 10 000 € sur une maison à 200 000 € (5 %) est bien plus acceptable qu'une erreur de 10 000 € sur une maison à 50 000 € (20 %). La MAPE capture cette différence.

Calcul de la MAPE sur notre exemple

Modèle A :

$$\begin{aligned} \text{MAPE}_A &= \frac{1}{3} \left(\left| \frac{200 - 205}{200} \right| + \left| \frac{320 - 310}{320} \right| + \left| \frac{148 - 150}{148} \right| \right) \times 100 \\ &= \frac{1}{3} \left(\frac{5}{200} + \frac{10}{320} + \frac{2}{148} \right) \times 100 \\ &= \frac{1}{3} (0,025 + 0,03125 + 0,01351) \times 100 \\ &= \frac{0,06976}{3} \times 100 \approx \mathbf{2,33 \%} \end{aligned}$$

Modèle B :

$$\begin{aligned} \text{MAPE}_B &= \frac{1}{3} \left(\frac{5}{200} + \frac{20}{320} + \frac{7}{148} \right) \times 100 \\ &= \frac{1}{3} (0,025 + 0,0625 + 0,04730) \times 100 \\ &= \frac{0,13480}{3} \times 100 \approx \mathbf{4,49\%} \end{aligned}$$

Interprétation : Le Modèle A se trompe en moyenne de 2,33 % du prix réel, contre 4,49 % pour le Modèle B.

Limites de la MAPE

- **Division par zéro** : si $y_i = 0$, la MAPE n'est pas définie !
- **Asymétrie** : la MAPE pénalise davantage les surestimations que les sous-estimations (car le dénominateur $|y_i|$ est le même mais le numérateur change)
- **Non utilisable** pour les données contenant des zéros (températures, scores, etc.)

9.4 Récapitulatif de toutes les métriques sur le Hands-On

Métrique	Modèle A	Modèle B	Meilleur
MAE	5,67 k€	10,67 k€	Modèle A
MSE	43 k€ ²	158 k€ ²	Modèle A
RMSE	6,56 k€	12,57 k€	Modèle A
R^2	0,9917	0,9695	Modèle A
MAPE	2,33 %	4,49 %	Modèle A

TABLE 9.1 – Comparaison complète des deux modèles — le Modèle A gagne sur toutes les métriques.

Quand les métriques ne sont pas d'accord

Dans notre exemple, toutes les métriques désignent le même gagnant. Mais ce n'est pas toujours le cas ! Imaginez un Modèle C qui prédit très bien les maisons chères mais mal les maisons bon marché. Sa MAPE pourrait être meilleure (erreur relative faible sur les grandes valeurs) alors que sa MAE serait pire. **C'est pourquoi il faut toujours regarder plusieurs métriques.**

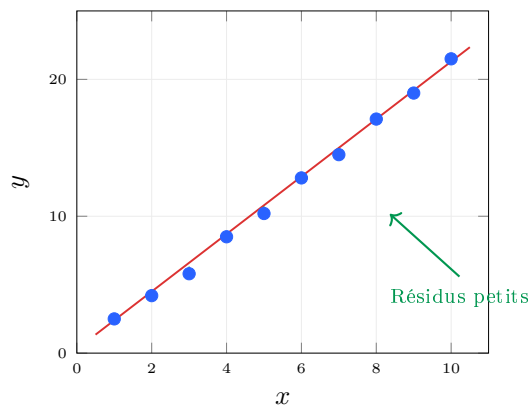
9.5 Tableau de synthèse des métriques

Métrique	Formule	Unité	Plage	Quand l'utiliser
MAE	$\frac{1}{n} \sum y_i - \hat{y}_i $	même que y	$[0, +\infty[$	Erreur interprétable, robuste aux outliers
MSE	$\frac{1}{n} \sum (y_i - \hat{y}_i)^2$	y^2	$[0, +\infty[$	Fonction de coût, pénaliser les grosses erreurs
RMSE	$\sqrt{\text{MSE}}$	même que y	$[0, +\infty[$	Métrique principale, comparable à la MAE
R^2	$1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$	sans unité	$] - \infty, 1]$	Comparer des modèles, % de variance expliquée
R^2_{adj}	corrigé par n et p	sans unité	$] - \infty, 1]$	Comparer des modèles avec différents nombres de variables
MAPE	$\frac{1}{n} \sum \left \frac{y_i - \hat{y}_i}{y_i} \right \times 100$	%	$[0, +\infty[$	Erreur relative, communiquer avec des non-techniciens

TABLE 9.2 – Synthèse des métriques d'évaluation pour la régression.

9.6 Comparaison visuelle : bon modèle vs mauvais modèle

Bon modèle ($R^2 \approx 0,95$)



Mauvais modèle ($R^2 \approx 0,40$)

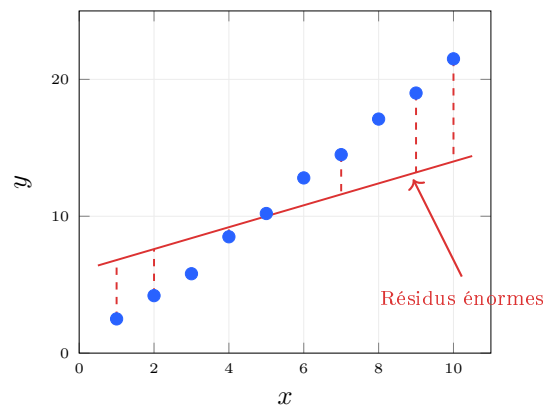


FIGURE 9.2 – Comparaison visuelle des résidus. À gauche : un bon modèle (points proches de la droite). À droite : un mauvais modèle (points éloignés de la droite). Les traits en pointillés représentent les résidus $y_i - \hat{y}_i$.

9.7 Application sur le dataset clientèle

Appliquons maintenant toutes les métriques à notre modèle de régression construit dans les chapitres précédents sur le dataset clientèle.

Application complète sur le dataset clientèle

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4 from sklearn.linear_model import LinearRegression
5 from sklearn.preprocessing import PolynomialFeatures
6 from sklearn.metrics import (mean_absolute_error,
7                               mean_squared_error,
8                               r2_score)
9
10 # --- 1. Charger et préparer les données ---
11 df = pd.read_csv('dataset_clientele.csv')
```

```

12
13 # Variables explicatives et cible
14 X = df[['age', 'revenu', 'anciennete', 'nb_achats']].values
15 y = df['depense'].values
16
17 # Separation train/test
18 X_train, X_test, y_train, y_test = train_test_split(
19     X, y, test_size=0.2, random_state=42
20 )
21
22 # --- 2. Entraîner un modele de regression lineaire ---
23 model = LinearRegression()
24 model.fit(X_train, y_train)
25 y_pred = model.predict(X_test)
26
27 # --- 3. Calculer toutes les metriques ---
28 mae = mean_absolute_error(y_test, y_pred)
29 mse = mean_squared_error(y_test, y_pred)
30 rmse = np.sqrt(mse)
31 r2 = r2_score(y_test, y_pred)
32
33 # R2 ajuste
34 n = len(y_test)
35 p = X_test.shape[1]
36 r2_adj = 1 - (1 - r2) * (n - 1) / (n - p - 1)
37
38 # MAPE (fonction personnalisee)
39 def mape(y_true, y_pred):
40     '''Calcul le Mean Absolute Percentage Error.'''
41     y_true = np.array(y_true)
42     y_pred = np.array(y_pred)
43     # Eviter la division par zero
44     mask = y_true != 0
45     return np.mean(np.abs(
46         (y_true[mask] - y_pred[mask]) / y_true[mask]
47     )) * 100
48
49 mape_val = mape(y_test, y_pred)
50
51 # --- 4. Afficher les resultats ---

```

```

52 print(''='' * 50)
53 print(''METRIQUES D'EVALUATION - Regression Lineaire'')
54 print(''='' * 50)
55 print(f''MAE           : {mae:.2f} euros'')
56 print(f''MSE           : {mse:.2f} euros^2'')
57 print(f''RMSE          : {rmse:.2f} euros'')
58 print(f''R2            : {r2:.4f}'')
59 print(f''R2 ajuste     : {r2_adj:.4f}'')
60 print(f''MAPE          : {mape_val:.2f} %'')

```

Listing 9.1 – Calcul de toutes les métriques sur le dataset clientèle

9.8 Implémentation complète sur Google Colab

Comparaison de plusieurs modèles de régression

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.model_selection import train_test_split
5 from sklearn.linear_model import LinearRegression, Ridge, Lasso
6 from sklearn.preprocessing import PolynomialFeatures
7 from sklearn.pipeline import make_pipeline
8 from sklearn.metrics import (mean_absolute_error,
9                               mean_squared_error,
10                              r2_score)
11
12 # --- Donnees ---
13 df = pd.read_csv('dataset_clientele.csv')
14 X = df[['age', 'revenu', 'anciennete', 'nb_achats']].values
15 y = df['depense'].values
16 X_train, X_test, y_train, y_test = train_test_split(
17     X, y, test_size=0.2, random_state=42
18 )
19
20 # --- Fonction MAPE ---
21 def mape(y_true, y_pred):
22     mask = y_true != 0
23     return np.mean(np.abs(

```

```

24         (y_true[mask] - y_pred[mask]) / y_true[mask]
25     )) * 100
26
27     # --- Définir les modeles ---
28     models = {
29         'Lineaire': LinearRegression(),
30         'Polynomial (deg=2)': make_pipeline(
31             PolynomialFeatures(degree=2, include_bias=False),
32             LinearRegression()
33         ),
34         'Ridge (alpha=1)': Ridge(alpha=1.0),
35         'Lasso (alpha=0.1)': Lasso(alpha=0.1),
36     }
37
38     # --- Entraîner et évaluer chaque modele ---
39     results = []
40
41     for name, model in models.items():
42         model.fit(X_train, y_train)
43         y_pred = model.predict(X_test)
44
45         results.append({
46             'Modele': name,
47             'MAE': mean_absolute_error(y_test, y_pred),
48             'MSE': mean_squared_error(y_test, y_pred),
49             'RMSE': np.sqrt(mean_squared_error(y_test, y_pred)),
50             'R2': r2_score(y_test, y_pred),
51             'MAPE (%)': mape(y_test, y_pred)
52         })
53
54     # --- Tableau comparatif ---
55     df_results = pd.DataFrame(results)
56     print(df_results.to_string(index=False))

```

Listing 9.2 – Comparer les métriques de 4 modèles de régression

Graphique de comparaison du R^2 des modèles

```

1 # --- Graphique : R2 de chaque modele ---

```

```

2 plt.figure(figsize=(10, 6))
3
4 couleurs = ['#2962FF', '#00C853', '#FF6D00', '#AA00FF']
5 bars = plt.bar(df_results['Modele'],
6                df_results['R2'],
7                color=couleurs,
8                edgecolor='black',
9                linewidth=0.8)
10
11 # Ajouter les valeurs sur les barres
12 for bar, val in zip(bars, df_results['R2']):
13     plt.text(bar.get_x() + bar.get_width()/2,
14             bar.get_height() + 0.005,
15             f'{val:.4f}',
16             ha='center', va='bottom',
17             fontweight='bold', fontsize=11)
18
19 plt.xlabel('Modele de regression', fontsize=12)
20 plt.ylabel('$R^2$ (coefficient de determination)', fontsize=12)
21 plt.title('Comparaison du $R^2$ des differents modeles',
22          fontsize=14, fontweight='bold')
23 plt.ylim(0, 1.1)
24 plt.grid(axis='y', alpha=0.3)
25 plt.tight_layout()
26 plt.show()

```

Listing 9.3 – Diagramme en barres comparant le R^2 des modèles

Analyse des résidus

```

1 # --- Choisir le meilleur modele ---
2 best_model = models['Lineaire'] # ou le meilleur selon R2
3 y_pred = best_model.predict(X_test)
4 residus = y_test - y_pred
5
6 fig, axes = plt.subplots(1, 3, figsize=(18, 5))
7
8 # --- Graphique 1 : Residus vs Predictions ---
9 axes[0].scatter(y_pred, residus, color='#2962FF',
10                alpha=0.6, edgecolors='black', linewidth=0.5)

```

```

11 axes[0].axhline(y=0, color='red', linestyle='--', linewidth=2)
12 axes[0].set_xlabel('Valeurs predites ( $\hat{y}$ )', fontsize=11)
13 axes[0].set_ylabel('Residus ( $y - \hat{y}$ )', fontsize=11)
14 axes[0].set_title('Residus vs Predictions', fontsize=13,
15                   fontweight='bold')
16 axes[0].grid(alpha=0.3)
17
18 # --- Graphique 2 : Histogramme des residus ---
19 axes[1].hist(residus, bins=20, color='#00C853',
20             edgecolor='black', alpha=0.7)
21 axes[1].axvline(x=0, color='red', linestyle='--', linewidth=2)
22 axes[1].set_xlabel('Residus', fontsize=11)
23 axes[1].set_ylabel('Frequence', fontsize=11)
24 axes[1].set_title('Distribution des residus', fontsize=13,
25                   fontweight='bold')
26 axes[1].grid(alpha=0.3)
27
28 # --- Graphique 3 : Q-Q plot (concept) ---
29 from scipy import stats
30 (osm, osr), (slope, intercept, r) = stats.probplot(
31     residus, dist='norm',
32 )
33 axes[2].scatter(osm, osr, color='#FF6D00',
34                alpha=0.6, edgecolors='black', linewidth=0.5)
35 axes[2].plot(osm, slope * np.array(osm) + intercept,
36             color='red', linewidth=2)
37 axes[2].set_xlabel('Quantiles theoriques', fontsize=11)
38 axes[2].set_ylabel('Quantiles observes', fontsize=11)
39 axes[2].set_title('Q-Q Plot des residus', fontsize=13,
40                   fontweight='bold')
41 axes[2].grid(alpha=0.3)
42
43 plt.tight_layout()
44 plt.show()
45
46 # --- Interpretation ---
47 print('INTERPRETATION DES GRAPHIQUES :')
48 print('-' * 50)
49 print('1. Residus vs Predictions :')
50 print('    - Les points doivent etre disperses aleatoirement')

```

```

51 print('' - Pas de forme (cone, courbe) visible'')
52 print('' - Sinon : probleme d'heteroscedasticite'')
53 print()
54 print(''2. Histogramme des residus :'')
55 print('' - Les residus doivent suivre une distribution'')
56 print('' - approximativement normale (forme de cloche)'')
57 print()
58 print(''3. Q-Q Plot :'')
59 print('' - Les points doivent suivre la droite rouge'')
60 print('' - Ecart = non-normalite des residus'')

```

Listing 9.4 – Graphiques d'analyse des résidus

Comprendre le Q-Q Plot

Le **Q-Q Plot** (Quantile-Quantile Plot) compare la distribution de nos résidus à une distribution normale théorique.

- Si les points suivent la droite rouge \Rightarrow les résidus sont **normalement distribués** (c'est bien!)
- Si les points s'écartent aux extrémités \Rightarrow les résidus ont des **queues lourdes** (outliers)
- Si les points forment une courbe \Rightarrow la distribution est **asymétrique**

Des résidus normalement distribués sont une hypothèse importante pour la régression linéaire.

9.9 Exercices

Exercice 9.1 — Calcul à la main de toutes les métriques

On dispose de 5 observations avec les valeurs réelles et prédites suivantes :

i	y_i (réel)	\hat{y}_i (prédit)
1	10	12
2	20	18
3	30	28
4	40	45
5	50	48

Calculez à la main :

1. La MAE
2. La MSE
3. La RMSE
4. Le R^2
5. La MAPE

Interprétez chaque résultat.

Correction de l'exercice 9.1

Étape préliminaire : Calculons les erreurs pour chaque observation.

i	y_i	\hat{y}_i	$y_i - \hat{y}_i$	$ y_i - \hat{y}_i $	$(y_i - \hat{y}_i)^2$
1	10	12	-2	2	4
2	20	18	+2	2	4
3	30	28	+2	2	4
4	40	45	-5	5	25
5	50	48	+2	2	4
Somme				13	41

1. MAE :

$$\text{MAE} = \frac{1}{5}(2 + 2 + 2 + 5 + 2) = \frac{13}{5} = \boxed{2,6}$$

Interprétation : En moyenne, le modèle se trompe de 2,6 unités.

2. MSE :

$$\text{MSE} = \frac{1}{5}(4 + 4 + 4 + 25 + 4) = \frac{41}{5} = \boxed{8,2}$$

Interprétation : La moyenne des carrés des erreurs est 8,2. L'erreur de 5 (observation 4) contribue à elle seule $25/41 = 61\%$ de la MSE !

3. RMSE :

$$\text{RMSE} = \sqrt{8,2} \approx \boxed{2,86}$$

Interprétation : L'erreur typique est d'environ 2,86 unités. Notez que $\text{RMSE} > \text{MAE}$ car la grosse erreur (5) tire la RMSE vers le haut.

4. R^2 :

Calculons d'abord \bar{y} :

$$\bar{y} = \frac{10 + 20 + 30 + 40 + 50}{5} = \frac{150}{5} = 30$$

Puis SS_{tot} :

$$\begin{aligned} SS_{\text{tot}} &= (10 - 30)^2 + (20 - 30)^2 + (30 - 30)^2 + (40 - 30)^2 + (50 - 30)^2 \\ &= 400 + 100 + 0 + 100 + 400 = 1000 \end{aligned}$$

Et $SS_{\text{res}} = 41$ (calculé ci-dessus).

$$R^2 = 1 - \frac{41}{1000} = 1 - 0,041 = \boxed{0,959}$$

Interprétation : Le modèle explique 95,9 % de la variance de y . C'est un très bon modèle !

5. MAPE :

$$\begin{aligned} \text{MAPE} &= \frac{1}{5} \left(\frac{2}{10} + \frac{2}{20} + \frac{2}{30} + \frac{5}{40} + \frac{2}{50} \right) \times 100 \\ &= \frac{1}{5} (0,2 + 0,1 + 0,0667 + 0,125 + 0,04) \times 100 \\ &= \frac{0,5317}{5} \times 100 = 0,10633 \times 100 \approx \boxed{10,63 \%} \end{aligned}$$

Interprétation : En moyenne, le modèle se trompe de 10,63 % par rapport aux valeurs réelles. La MAPE est relativement élevée car l'erreur relative est grande pour les petites valeurs ($2/10 = 20\%$ pour $y_1 = 10$).

Exercice 9.2 — R^2 vs R^2 ajusté

Vous avez deux modèles entraînés sur un dataset de $n = 200$ observations :

- **Modèle 1** : 5 variables, $R^2 = 0,85$
 - **Modèle 2** : 15 variables (les 5 de base + 10 nouvelles), $R^2 = 0,90$
1. Calculez le R^2 ajusté de chaque modèle.
 2. Le Modèle 2 est-il vraiment meilleur ?
 3. Que concluez-vous sur l'ajout de 10 variables supplémentaires ?

Correction de l'exercice 9.2**1. Calcul du R^2 ajusté :**

Modèle 1 ($n = 200$, $p = 5$, $R^2 = 0,85$) :

$$R^2_{\text{adj},1} = 1 - \frac{(1 - 0,85)(200 - 1)}{200 - 5 - 1} = 1 - \frac{0,15 \times 199}{194} = 1 - \frac{29,85}{194} = 1 - 0,1539 = \boxed{0,8461}$$

Modèle 2 ($n = 200$, $p = 15$, $R^2 = 0,90$) :

$$R^2_{\text{adj},2} = 1 - \frac{(1 - 0,90)(200 - 1)}{200 - 15 - 1} = 1 - \frac{0,10 \times 199}{184} = 1 - \frac{19,9}{184} = 1 - 0,1082 = \boxed{0,8918}$$

2. Comparaison :

	R^2	R^2_{adj}	Nb variables
Modèle 1	0,8500	0,8461	5
Modèle 2	0,9000	0,8918	15

Ici, le R^2 ajusté du Modèle 2 (0,8918) est supérieur à celui du Modèle 1 (0,8461). Donc le Modèle 2 **est effectivement meilleur**, même en tenant compte de ses 10 variables supplémentaires.

3. Conclusion :

Dans ce cas, les 10 variables supplémentaires apportent une amélioration **réelle**, car le R^2 ajusté augmente aussi. Cependant, il faut noter que l'amélioration est de +5% en R^2 mais seulement +4,6% en R^2_{adj} , ce qui montre qu'une petite partie du gain en R^2 était artificielle.

Si le R^2_{adj} du Modèle 2 avait été inférieur à celui du Modèle 1, cela aurait signifié que les variables ajoutées étaient majoritairement inutiles et qu'il fallait préférer le modèle le plus simple (principe de parcimonie).

Exercice 9.3 — Python : évaluation complète de 4 modèles

Écrivez un programme Python qui :

1. Génère un jeu de données synthétique (100 points, relation polynomiale de degré 2 avec du bruit)
2. Sépare en train (80%) et test (20%)
3. Entraîne 4 modèles : régression linéaire, polynomiale (degré 2), polynomiale (degré 5), Ridge

4. Calcule MAE, MSE, RMSE, R^2 et MAPE pour chaque modèle
5. Affiche un tableau récapitulatif
6. Crée un graphique en barres comparant le R^2 des 4 modèles

Correction de l'exercice 9.3

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.model_selection import train_test_split
5 from sklearn.linear_model import LinearRegression, Ridge
6 from sklearn.preprocessing import PolynomialFeatures
7 from sklearn.pipeline import make_pipeline
8 from sklearn.metrics import (mean_absolute_error,
9                               mean_squared_error,
10                              r2_score)
11
12 # ===== 1. Generer les donnees =====
13 np.random.seed(42)
14 n = 100
15 X = np.linspace(0, 10, n).reshape(-1, 1)
16 # Relation reelle :  $y = 3x^2 - 5x + 10 + \text{bruit}$ 
17 y = 3 * X.ravel()**2 - 5 * X.ravel() + 10
18 y += np.random.normal(0, 15, n) # Ajout de bruit
19
20 # ===== 2. Separation train/test =====
21 X_train, X_test, y_train, y_test = train_test_split(
22     X, y, test_size=0.2, random_state=42
23 )
24
25 # ===== 3. Fonction MAPE =====
26 def mape(y_true, y_pred):
27     mask = y_true != 0
28     return np.mean(np.abs(
29         (y_true[mask] - y_pred[mask]) / y_true[mask]
30     )) * 100
31
32 # ===== 4. Definir et entrainer les modeles =====
33 models = {
34     'Lineaire': LinearRegression(),

```

```

35     'Poly deg=2': make_pipeline(
36         PolynomialFeatures(2, include_bias=False),
37         LinearRegression()
38     ),
39     'Poly deg=5': make_pipeline(
40         PolynomialFeatures(5, include_bias=False),
41         LinearRegression()
42     ),
43     'Ridge': make_pipeline(
44         PolynomialFeatures(2, include_bias=False),
45         Ridge(alpha=1.0)
46     ),
47 }
48
49 # ===== 5. Evaluer chaque modele =====
50 results = []
51 for name, model in models.items():
52     model.fit(X_train, y_train)
53     y_pred = model.predict(X_test)
54
55     results.append({
56         'Modele': name,
57         'MAE': round(mean_absolute_error(y_test, y_pred), 2),
58         'MSE': round(mean_squared_error(y_test, y_pred), 2),
59         'RMSE': round(np.sqrt(
60             mean_squared_error(y_test, y_pred)), 2),
61         'R2': round(r2_score(y_test, y_pred), 4),
62         'MAPE (%)': round(mape(y_test, y_pred), 2),
63     })
64
65 # ===== 6. Tableau recapitulatif =====
66 df_res = pd.DataFrame(results)
67 print('\n' + '=' * 65)
68 print('          TABLEAU COMPARATIF DES MODELES')
69 print('=' * 65)
70 print(df_res.to_string(index=False))
71 print('=' * 65)
72
73 # ===== 7. Graphique de comparaison du R2 =====
74 fig, axes = plt.subplots(1, 2, figsize=(14, 5))

```

```

75
76 # Graphique 1 : R2
77 couleurs = ['#2962FF', '#00C853', '#FF6D00', '#AA00FF']
78 bars = axes[0].bar(df_res['Modele'], df_res['R2'],
79                   color=couleurs, edgecolor='black')
80 for bar, val in zip(bars, df_res['R2']):
81     axes[0].text(bar.get_x() + bar.get_width()/2,
82                 bar.get_height() + 0.01,
83                 f'{val:.4f}', ha='center',
84                 fontweight='bold', fontsize=10)
85 axes[0].set_ylabel('$R^2$', fontsize=12)
86 axes[0].set_title('Comparaison du $R^2$',
87                  fontsize=13, fontweight='bold')
88 axes[0].set_ylim(0, 1.15)
89 axes[0].grid(axis='y', alpha=0.3)
90 axes[0].tick_params(axis='x', rotation=15)
91
92 # Graphique 2 : MAE et RMSE cote a cote
93 x_pos = np.arange(len(df_res))
94 width = 0.35
95 axes[1].bar(x_pos - width/2, df_res['MAE'],
96             width, label='MAE', color='#2962FF',
97             edgecolor='black', alpha=0.8)
98 axes[1].bar(x_pos + width/2, df_res['RMSE'],
99             width, label='RMSE', color='#FF6D00',
100             edgecolor='black', alpha=0.8)
101 axes[1].set_xticks(x_pos)
102 axes[1].set_xticklabels(df_res['Modele'], rotation=15)
103 axes[1].set_ylabel('Erreur', fontsize=12)
104 axes[1].set_title('Comparaison MAE vs RMSE',
105                  fontsize=13, fontweight='bold')
106 axes[1].legend(fontsize=11)
107 axes[1].grid(axis='y', alpha=0.3)
108
109 plt.tight_layout()
110 plt.show()
111
112 # ==== 8. Identifier le meilleur modele ====
113 best_idx = df_res['R2'].idxmax()
114 print(f'''\nMeilleur modele (selon R2) : ''

```

```
115     f''{df_res.loc[best_idx, 'Modele']} ''  
116     f''(R2 = {df_res.loc[best_idx, 'R2']})''
```

Listing 9.5 – Solution complète — comparaison de 4 modèles

Ce qu'il faut retenir de ce chapitre

1. **Toujours évaluer** votre modèle — un modèle non évalué ne sert à rien.
2. **Ne vous fiez jamais à une seule métrique.** Utilisez au moins MAE/RMSE + R^2 .
3. **MAE** : interprétable, robuste. **RMSE** : pénalise les grosses erreurs.
4. R^2 : proportion de variance expliquée (0 à 1, plus haut = mieux).
5. R^2 **ajusté** : indispensable quand on compare des modèles avec un nombre différent de variables.
6. **MAPE** : utile pour communiquer en pourcentage, mais attention aux zéros.
7. **Analysez les résidus !** Les graphiques de résidus révèlent des problèmes que les métriques seules ne montrent pas.

Deuxième partie

Classification — Prédire une Catégorie

Chapitre 10

Régression Logistique

Objectifs de ce chapitre :

- Comprendre pourquoi la régression linéaire ne convient pas pour la classification
- Maîtriser la fonction sigmoïde et ses propriétés
- Dériver la fonction de coût **Binary Cross-Entropy** pas à pas
- Comprendre le concept de vraisemblance maximale (*Maximum Likelihood*)
- Calculer les gradients et comprendre la descente de gradient pour la régression logistique
- Implémenter la régression logistique avec Scikit-learn sur Google Colab

10.1 Hands-On : prédire le défaut de paiement

Situation réelle : une banque veut prédire les défauts de paiement

Vous travaillez dans une banque. Chaque jour, des dizaines de personnes demandent un prêt. La banque veut savoir : **ce client va-t-il rembourser son prêt ou faire défaut ?** C'est un problème de **classification binaire** :

- **Classe 0** : le client rembourse (pas de défaut)
- **Classe 1** : le client fait défaut (ne rembourse pas)

On dispose de deux informations sur chaque client :

- **Revenu mensuel** (en milliers d'€)
- **Score de crédit** (entre 300 et 850)

10.1.1 Le jeu de données

Voici les données de 10 anciens clients pour lesquels on connaît le résultat :

Client	Revenu (k€)	Score de crédit	Défaut (0/1)
1	2,0	450	1
2	2,5	500	1
3	3,0	520	1
4	3,2	580	1
5	3,5	600	0
6	4,0	620	0
7	4,5	650	0
8	5,0	700	0
9	5,5	720	0
10	6,0	750	0

TABLE 10.1 – Données de 10 clients avec revenu, score de crédit et défaut de paiement.

10.1.2 Visualisation du problème

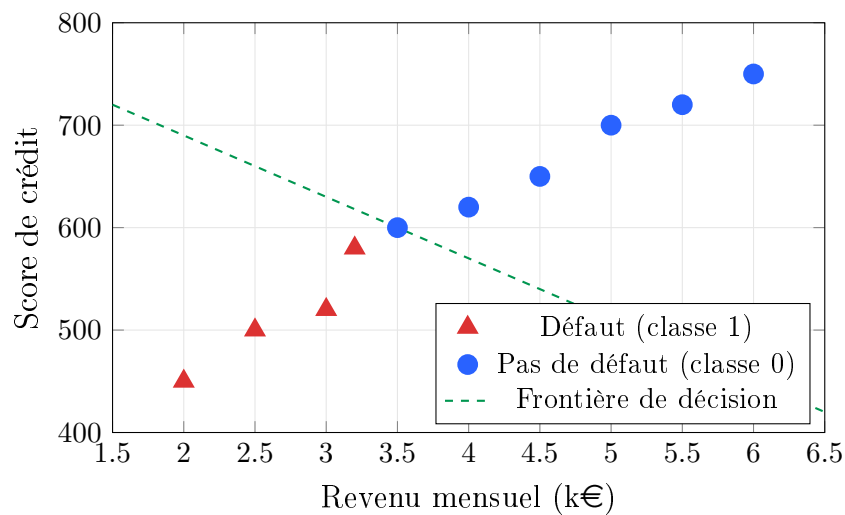


FIGURE 10.1 – Nuage de points des 10 clients colorés par classe. On observe une séparation nette entre les deux groupes.

comment obtenir une **probabilité** que le client fasse défaut ? Par exemple : « Ce client a 73% de chances de faire défaut ».

C'est exactement ce que fait la **régression logistique** !

10.2 Intuition : pourquoi pas la régression linéaire ?

10.2.1 Le problème de la régression linéaire pour la classification

Pourquoi la régression linéaire ne marche pas ici ?

On pourrait se dire : « Utilisons simplement une régression linéaire pour prédire la classe (0 ou 1) ». Après tout, c'est un nombre, non ?

Le problème : la régression linéaire peut prédire **n'importe quel nombre** : $-0,3$, $0,7$, $1,5$, $2,8$...

Mais une probabilité **doit être entre 0 et 1** ! Une probabilité de $-0,3$ ou de $1,5$ n'a aucun sens.

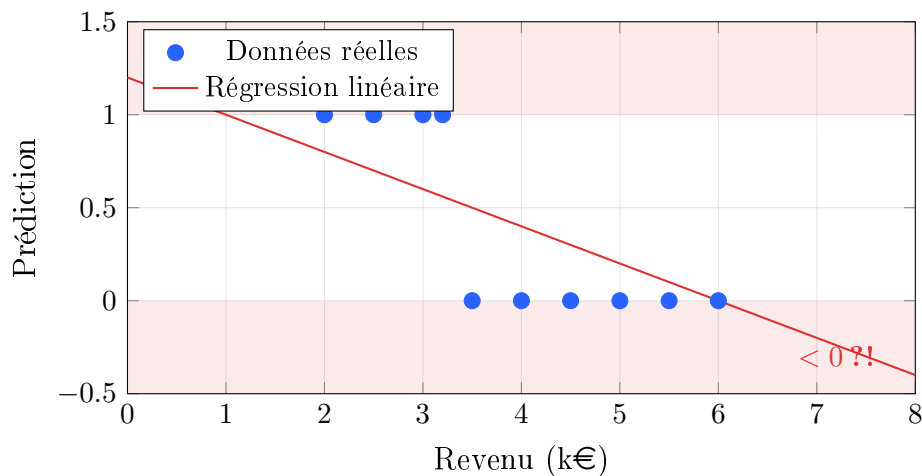


FIGURE 10.2 – La régression linéaire peut prédire des valeurs hors de $[0, 1]$: ça n'a pas de sens pour une probabilité !

10.2.2 La solution : « écraser » les valeurs entre 0 et 1

Analogie : le variateur de lumière

Imaginez un **variateur de lumière** (dimmer) :

- Quand vous tournez le bouton complètement à gauche → la lampe est **éteinte** (0)
- Quand vous tournez complètement à droite → la lampe est **allumée** (1)
- Entre les deux, la lumière augmente **progressivement** de 0 à 1

Il n'y a **jamais** de lumière négative, et jamais plus que le maximum. La transition est **douce et continue**, en forme de S.

C'est exactement ce que fait la **fonction sigmoïde** ! Elle prend n'importe quel nombre réel et le transforme en une valeur entre 0 et 1.

10.3 Dérivation mathématique complète

10.3.1 La fonction sigmoïde

Prérequis : le nombre e (nombre d'Euler)

Le nombre e — le nombre de la croissance naturelle

Le nombre e est une constante mathématique fondamentale, tout comme π . Sa valeur est :

$$e \approx 2,71828 \dots$$

D'où vient-il ? Imaginez que vous placez 1€ à la banque avec un taux d'intérêt de 100% par an :

- Intérêts 1 fois/an : $1 \times (1 + 1)^1 = 2,00\text{€}$
- Intérêts 2 fois/an : $1 \times (1 + 0,5)^2 = 2,25\text{€}$
- Intérêts 12 fois/an : $1 \times (1 + 1/12)^{12} \approx 2,613\text{€}$
- Intérêts 365 fois/an : $1 \times (1 + 1/365)^{365} \approx 2,7146\text{€}$
- Intérêts en continu : $\lim_{n \rightarrow \infty} (1 + 1/n)^n = e \approx 2,71828\text{€}$

Le nombre e représente la **croissance naturelle continue**. C'est pourquoi on le retrouve partout en mathématiques, en physique et en Machine Learning.

Définition de la sigmoïde

Fonction sigmoïde $\sigma(z)$

La **fonction sigmoïde** (aussi appelée *logistic function*) est définie par :

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

où z est un nombre réel quelconque (positif, négatif, ou zéro).

Comprendre e^{-z} : que se passe-t-il ?

Décortiquons cette formule en analysant e^{-z} pour différentes valeurs de z :

- Si z est très grand positif (ex : $z = 10$) :

$$e^{-10} = \frac{1}{e^{10}} = \frac{1}{22026} \approx 0,000045 \approx 0$$

Donc $\sigma(10) = \frac{1}{1+0} \approx 1 \Rightarrow$ la sortie est proche de 1

— Si z est très grand négatif (ex : $z = -10$) :

$$e^{-(-10)} = e^{10} \approx 22026$$

Donc $\sigma(-10) = \frac{1}{1+22026} \approx \frac{1}{22027} \approx 0 \Rightarrow$ **la sortie est proche de 0**

— Si $z = 0$:

$$e^0 = 1$$

Donc $\sigma(0) = \frac{1}{1+1} = \frac{1}{2} = 0,5 \Rightarrow$ **pile au milieu !**

Table de valeurs de la sigmoïde

z	e^{-z}	$1 + e^{-z}$	$\sigma(z) = \frac{1}{1 + e^{-z}}$
-5	148,41	149,41	$0,0067 \approx 0,01$
-2	7,389	8,389	$0,119 \approx 0,12$
-1	2,718	3,718	$0,269 \approx 0,27$
0	1,000	2,000	0,500
1	0,368	1,368	$0,731 \approx 0,73$
2	0,135	1,135	$0,881 \approx 0,88$
5	0,0067	1,0067	$0,993 \approx 0,99$

TABLE 10.2 – Table de valeurs de la fonction sigmoïde pour différentes valeurs de z .

Graphique de la sigmoïde

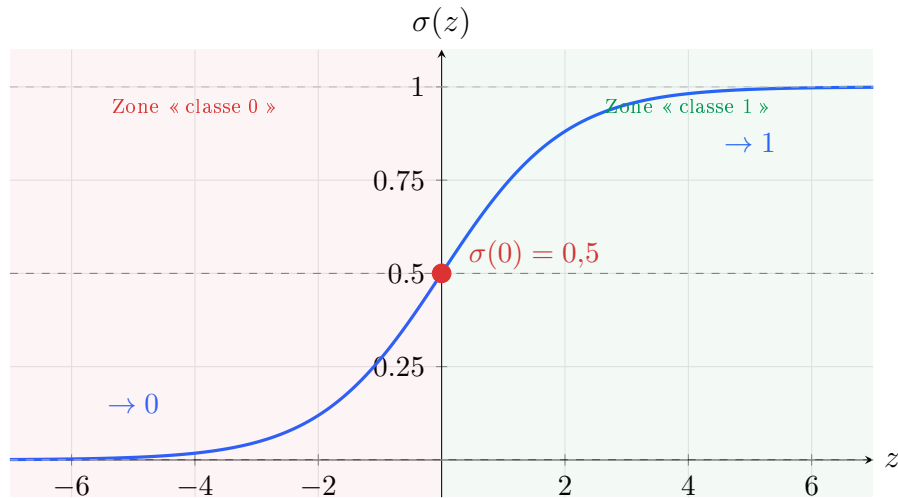


FIGURE 10.3 – La fonction sigmoïde : une courbe en S qui transforme tout nombre réel en une valeur entre 0 et 1.

Propriétés de la sigmoïde

Propriété 10.1. La fonction sigmoïde $\sigma(z) = \frac{1}{1+e^{-z}}$ possède les propriétés suivantes :

1. $0 < \sigma(z) < 1$ pour tout $z \in \mathbb{R}$ (toujours entre 0 et 1)
2. $\sigma(0) = 0,5$ (le point central)
3. $\lim_{z \rightarrow +\infty} \sigma(z) = 1$ et $\lim_{z \rightarrow -\infty} \sigma(z) = 0$
4. La fonction est **strictement croissante** (plus z augmente, plus $\sigma(z)$ augmente)
5. **Symétrie remarquable** : $\sigma(-z) = 1 - \sigma(z)$
6. **Dérivée remarquable** : $\sigma'(z) = \sigma(z)(1 - \sigma(z))$

Démonstration de la dérivée $\sigma'(z) = \sigma(z)(1 - \sigma(z))$

Cette propriété est **fondamentale** car elle simplifie énormément les calculs de gradient. Démontrons-la pas à pas.

Démonstration détaillée

On part de $\sigma(z) = \frac{1}{1+e^{-z}} = (1 + e^{-z})^{-1}$.

Étape 1 : On utilise la règle de dérivation des puissances : si $f(z) = [u(z)]^{-1}$, alors $f'(z) = -[u(z)]^{-2} \cdot u'(z)$.

Ici, $u(z) = 1 + e^{-z}$, donc $u'(z) = -e^{-z}$ (la dérivée de e^{-z} est $-e^{-z}$).

Étape 2 :

$$\sigma'(z) = -(1 + e^{-z})^{-2} \cdot (-e^{-z}) = \frac{e^{-z}}{(1 + e^{-z})^2}$$

Étape 3 : On réécrit sous une forme astucieuse :

$$\sigma'(z) = \frac{1}{1 + e^{-z}} \cdot \frac{e^{-z}}{1 + e^{-z}}$$

Étape 4 : Le premier facteur est $\sigma(z)$. Pour le second, remarquons que :

$$\frac{e^{-z}}{1 + e^{-z}} = \frac{(1 + e^{-z}) - 1}{1 + e^{-z}} = 1 - \frac{1}{1 + e^{-z}} = 1 - \sigma(z)$$

Conclusion :

$$\sigma'(z) = \sigma(z) \cdot (1 - \sigma(z))$$

C'est élégant : la dérivée de la sigmoïde s'exprime **en fonction d'elle-même** !

10.3.2 Le modèle de régression logistique

Combiner linéaire + sigmoïde

Modèle de régression logistique

La régression logistique combine une **combinaison linéaire** des variables avec la **fonction sigmoïde** :

Étape 1 — Combinaison linéaire :

$$z = w_1x_1 + w_2x_2 + \dots + w_px_p + b$$

où w_1, w_2, \dots, w_p sont les **poids** (paramètres à apprendre) et b est le **biais**.

Étape 2 — Passage par la sigmoïde :

$$P(y = 1 \mid \mathbf{x}) = \hat{p} = \sigma(z) = \frac{1}{1 + e^{-(w_1x_1 + w_2x_2 + \dots + w_px_p + b)}}$$

La sortie \hat{p} est une **probabilité** entre 0 et 1.

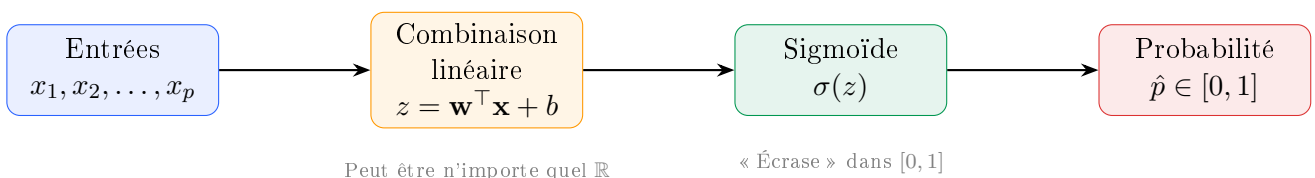


FIGURE 10.4 – Architecture de la régression logistique : linéaire \rightarrow sigmoïde \rightarrow probabilité.

Règle de décision

Règle de décision

Une fois la probabilité \hat{p} calculée, on prédit la classe avec un **seuil** (généralement 0,5) :

$$\hat{y} = \begin{cases} 1 & \text{si } \hat{p} \geq 0,5 \\ 0 & \text{si } \hat{p} < 0,5 \end{cases}$$

La frontière de décision

La frontière de décision est là où z

Le seuil $\hat{p} = 0,5$ correspond à $\sigma(z) = 0,5$, c'est-à-dire $z = 0$.

Donc la frontière de décision est définie par :

$$w_1x_1 + w_2x_2 + \dots + w_px_p + b = 0$$

C'est l'équation d'une **droite** (en 2D), d'un **plan** (en 3D) ou d'un **hyperplan** (en dimension p) !

Autrement dit, la régression logistique cherche une **frontière linéaire** pour séparer les deux classes.

10.3.3 Pourquoi pas le MSE ? Introduction à la vraisemblance maximale

Le problème du MSE en classification

Le MSE crée une surface non convexe !

Si on utilisait l'erreur quadratique moyenne (MSE) comme fonction de coût :

$$J_{\text{MSE}} = \frac{1}{n} \sum_{i=1}^n (y_i - \sigma(z_i))^2$$

La composition du MSE avec la sigmoïde crée une fonction **non convexe**, avec de nombreux **minima locaux**. La descente de gradient pourrait se coincer dans un mauvais minimum et ne jamais trouver la bonne solution !

Il nous faut une fonction de coût **convexe** \Rightarrow la **vraisemblance maximale**.

Le principe de vraisemblance maximale (Maximum Likelihood)

Vraisemblance (Likelihood)

La **vraisemblance** est la probabilité d'observer nos données étant donné un modèle avec des paramètres \mathbf{w} et b .

Idée intuitive : Parmi tous les paramètres possibles, on cherche ceux qui rendent nos données observées les **plus probables**.

Analogie : Imaginez que vous trouvez une pièce de monnaie et vous lancez 10 fois : vous obtenez 8 faces et 2 piles. Quelle est la probabilité p de tomber sur face ? Le Maximum de Vraisemblance répond : $p = 8/10 = 0,8$, car c'est la valeur qui rend vos observations les plus probables.

Vraisemblance pour un seul exemple

Pour un seul exemple (x_i, y_i) , le modèle prédit la probabilité $\hat{p}_i = \sigma(\mathbf{w}^\top \mathbf{x}_i + b)$.

La probabilité d'observer y_i est :

$$P(y_i | \mathbf{x}_i) = \hat{p}_i^{y_i} \cdot (1 - \hat{p}_i)^{1-y_i}$$

Vérification de la formule

Vérifions que cette formule est correcte :

Cas 1 : Si $y_i = 1$ (le client a fait défaut) :

$$P(y_i = 1 | \mathbf{x}_i) = \hat{p}_i^1 \cdot (1 - \hat{p}_i)^0 = \hat{p}_i \cdot 1 = \hat{p}_i \quad \checkmark$$

On obtient bien \hat{p}_i , la probabilité prédite d'être en classe 1.

Cas 2 : Si $y_i = 0$ (le client a remboursé) :

$$P(y_i = 0 | \mathbf{x}_i) = \hat{p}_i^0 \cdot (1 - \hat{p}_i)^1 = 1 \cdot (1 - \hat{p}_i) = 1 - \hat{p}_i \quad \checkmark$$

On obtient bien $1 - \hat{p}_i$, la probabilité d'être en classe 0.

Astucieux ! Une seule formule couvre les deux cas grâce aux exposants y_i et $1 - y_i$.

10.3.4 Dérivation de la log-vraisemblance

Vraisemblance totale

Pour **tous** les n exemples (supposés indépendants), la vraisemblance totale est le **produit** des vraisemblances individuelles :

$$L(\mathbf{w}, b) = \prod_{i=1}^n \hat{p}_i^{y_i} \cdot (1 - \hat{p}_i)^{1-y_i}$$

On veut **maximiser** L : trouver les paramètres \mathbf{w} et b qui rendent nos données les plus probables.

Passage au logarithme

Pourquoi le logarithme ?

Travailler avec des produits est mathématiquement difficile (dérivation compliquée). Le logarithme transforme les **produits en sommes** :

$$\ln(a \cdot b) = \ln(a) + \ln(b)$$

De plus, le logarithme est une fonction **strictement croissante**, donc maximiser $\ln L$ revient exactement à maximiser L .

Appliquons le logarithme :

$$\begin{aligned} \ln L &= \ln \left(\prod_{i=1}^n \hat{p}_i^{y_i} \cdot (1 - \hat{p}_i)^{1-y_i} \right) \\ &= \sum_{i=1}^n \ln (\hat{p}_i^{y_i} \cdot (1 - \hat{p}_i)^{1-y_i}) \quad (\text{produit} \rightarrow \text{somme}) \\ &= \sum_{i=1}^n [\ln(\hat{p}_i^{y_i}) + \ln((1 - \hat{p}_i)^{1-y_i})] \\ &= \sum_{i=1}^n [y_i \cdot \ln(\hat{p}_i) + (1 - y_i) \cdot \ln(1 - \hat{p}_i)] \end{aligned} \tag{10.1}$$

C'est la **log-vraisemblance**. On veut la **maximiser**.

10.3.5 Binary Cross-Entropy (BCE)

De la log-vraisemblance à la fonction de coût

En pratique, en optimisation, on préfère **minimiser** plutôt que maximiser. On prend donc le **négatif** de la log-vraisemblance, divisé par n :

Binary Cross-Entropy (BCE) — Entropie croisée binaire

$$J(\mathbf{w}, b) = -\frac{1}{n} \sum_{i=1}^n [y_i \cdot \ln(\hat{p}_i) + (1 - y_i) \cdot \ln(1 - \hat{p}_i)]$$

où $\hat{p}_i = \sigma(\mathbf{w}^\top \mathbf{x}_i + b)$ est la probabilité prédite par le modèle.

Minimiser J revient à **maximiser** la vraisemblance.

Interprétation intuitive de la BCE

La BCE pénalise les prédictions confiantes mais fausses

Analysons le comportement de la BCE cas par cas :

Quand $y_i = 1$ (le vrai label est 1), le coût pour cet exemple est $-\ln(\hat{p}_i)$:

- Si $\hat{p}_i \approx 1$ (bonne prédiction) : $-\ln(1) = 0 \Rightarrow$ coût **nul**
- Si $\hat{p}_i \approx 0$ (mauvaise prédiction) : $-\ln(0) \rightarrow +\infty \Rightarrow$ coût **énorme**

Quand $y_i = 0$ (le vrai label est 0), le coût pour cet exemple est $-\ln(1 - \hat{p}_i)$:

- Si $\hat{p}_i \approx 0$ (bonne prédiction) : $-\ln(1) = 0 \Rightarrow$ coût **nul**
- Si $\hat{p}_i \approx 1$ (mauvaise prédiction) : $-\ln(0) \rightarrow +\infty \Rightarrow$ coût **énorme**

Le message : « Si tu te trompes et que tu es sûr de toi, tu seras **LOURDEMENT** pénalisé ! »

Visualisation des courbes de coût

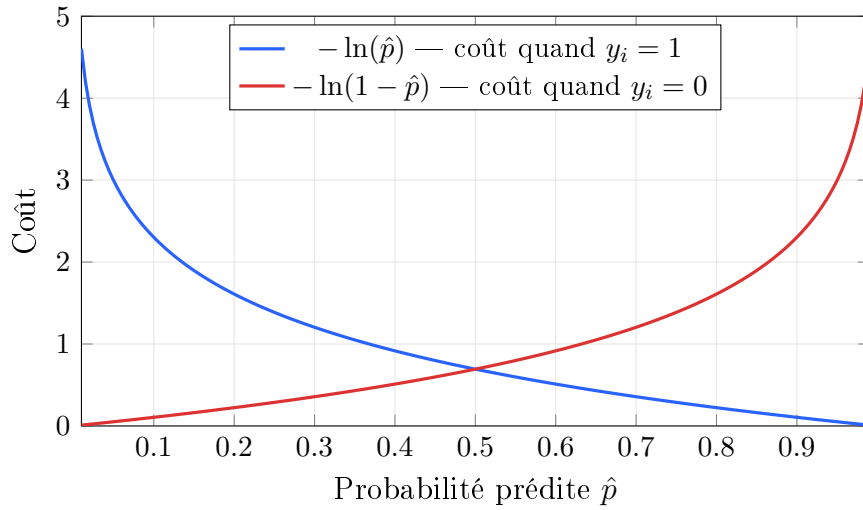


FIGURE 10.5 – Courbes de coût de la BCE. Quand $y_i = 1$, le coût explose si $\hat{p} \rightarrow 0$. Quand $y_i = 0$, le coût explose si $\hat{p} \rightarrow 1$.

10.3.6 Calcul du gradient

Objectif

Pour appliquer la descente de gradient, on a besoin de calculer les dérivées partielles de J par rapport à chaque paramètre w_j et au biais b .

Dérivation détaillée (règle de la chaîne)

Calculons $\frac{\partial J}{\partial w_j}$. Pour alléger, travaillons avec un seul exemple i puis généralisons.

Le coût pour un seul exemple est :

$$\ell_i = -[y_i \ln(\hat{p}_i) + (1 - y_i) \ln(1 - \hat{p}_i)]$$

Étape 1 — Chaîne de dépendance :

$$w_j \rightarrow z_i \rightarrow \hat{p}_i \rightarrow \ell_i$$

Par la règle de la chaîne :

$$\frac{\partial \ell_i}{\partial w_j} = \frac{\partial \ell_i}{\partial \hat{p}_i} \cdot \frac{\partial \hat{p}_i}{\partial z_i} \cdot \frac{\partial z_i}{\partial w_j}$$

Étape 2 — Calcul de chaque facteur :

Facteur 1 : $\frac{\partial \ell_i}{\partial \hat{p}_i}$

$$\frac{\partial \ell_i}{\partial \hat{p}_i} = - \left[\frac{y_i}{\hat{p}_i} - \frac{1 - y_i}{1 - \hat{p}_i} \right] = -\frac{y_i}{\hat{p}_i} + \frac{1 - y_i}{1 - \hat{p}_i}$$

Facteur 2 : $\frac{\partial \hat{p}_i}{\partial z_i} = \sigma'(z_i) = \hat{p}_i(1 - \hat{p}_i)$ (propriété démontrée plus haut)

Facteur 3 : $\frac{\partial z_i}{\partial w_j} = x_{ij}$ (car $z_i = w_1 x_{i1} + \dots + w_j x_{ij} + \dots + b$)

Étape 3 — Multiplication :

$$\begin{aligned}
 \frac{\partial \ell_i}{\partial w_j} &= \left(-\frac{y_i}{\hat{p}_i} + \frac{1 - y_i}{1 - \hat{p}_i} \right) \cdot \hat{p}_i(1 - \hat{p}_i) \cdot x_{ij} \\
 &= \left(\frac{-y_i(1 - \hat{p}_i) + (1 - y_i)\hat{p}_i}{\hat{p}_i(1 - \hat{p}_i)} \right) \cdot \hat{p}_i(1 - \hat{p}_i) \cdot x_{ij} \\
 &= (-y_i + y_i\hat{p}_i + \hat{p}_i - y_i\hat{p}_i) \cdot x_{ij} \\
 &= (\hat{p}_i - y_i) \cdot x_{ij}
 \end{aligned} \tag{10.2}$$

Gradient de la BCE — Résultat final

Le gradient de la fonction de coût par rapport à w_j est :

$$\frac{\partial J}{\partial w_j} = \frac{1}{n} \sum_{i=1}^n (\hat{p}_i - y_i) x_{ij}$$

Et par rapport au biais b :

$$\frac{\partial J}{\partial b} = \frac{1}{n} \sum_{i=1}^n (\hat{p}_i - y_i)$$

Remarquablement simple et élégant ! Le gradient est proportionnel à l'erreur $(\hat{p}_i - y_i)$ multipliée par l'entrée x_{ij} .

Règle de mise à jour

La descente de gradient met à jour les paramètres :

$$w_j := w_j - \alpha \cdot \frac{1}{n} \sum_{i=1}^n (\hat{p}_i - y_i) x_{ij}$$

$$b := b - \alpha \cdot \frac{1}{n} \sum_{i=1}^n (\hat{p}_i - y_i)$$

où α est le taux d'apprentissage (*learning rate*).

10.3.7 Les odds et le log-odds (logit)

Les odds (chances)

Odds (cote)

Les **odds** (ou cote) d'un événement sont le rapport entre la probabilité qu'il se produise et la probabilité qu'il ne se produise pas :

$$\text{Odds} = \frac{P}{1 - P}$$

Analogie avec les paris sportifs

Si un cheval a 80% de chances de gagner :

$$\text{Odds} = \frac{0,8}{1 - 0,8} = \frac{0,8}{0,2} = 4$$

On dit que les chances sont de « 4 contre 1 ». Pour chaque fois où il perd, il gagne 4 fois.

Si un cheval a 50% de chances :

$$\text{Odds} = \frac{0,5}{0,5} = 1 \quad (\text{« 1 contre 1 », ou « 50-50 »})$$

Le log-odds (logit) — le lien avec la régression

Log-odds (logit)

Le **logit** est le logarithme des odds :

$$\text{logit}(P) = \ln \left(\frac{P}{1 - P} \right)$$

Et voici la propriété fondamentale : dans la régression logistique :

$$\ln \left(\frac{P(y = 1 \mid \mathbf{x})}{1 - P(y = 1 \mid \mathbf{x})} \right) = w_1 x_1 + w_2 x_2 + \cdots + w_p x_p + b$$

Le log-odds est une **fonction linéaire** des variables !

Pourquoi « régression » logistique ?

Le nom « régression logistique » peut prêter à confusion : c'est bien un algorithme de **classification**, pas de régression !

On l'appelle « régression » car le log-odds est **linéaire** par rapport aux variables, exacte-

ment comme dans une régression linéaire. La différence est qu'on ne modélise pas directement y , mais le **logarithme des odds** de y .

10.4 Application sur le dataset clientèle

Revenons à notre dataset fil rouge. On veut prédire le **Churn** (départ du client, 0 ou 1) à partir du **Revenu**.

10.4.1 Les données

Client	Revenu (k€)	Churn (0/1)
1	15	1
2	18	1
3	20	1
4	22	1
5	25	0
6	28	0
7	30	0
8	35	0
9	40	0
10	45	0

TABLE 10.3 – Dataset clientèle : Revenu et Churn.

10.4.2 Application du modèle

Supposons que notre modèle a appris les paramètres $w = -0,2$ et $b = 5$. Le modèle prédit :

$$\hat{p} = \sigma(z) = \frac{1}{1 + e^{-(0,2 \cdot x + 5)}} = \frac{1}{1 + e^{0,2x - 5}}$$

Calculons \hat{p} pour chaque client :

Client	Revenu	$z = -0,2x + 5$	$\hat{p} = \sigma(z)$	Prédiction	Churn réel
1	15	2,0	0,881	1	1 ✓
2	18	1,4	0,802	1	1 ✓
3	20	1,0	0,731	1	1 ✓
4	22	0,6	0,646	1	1 ✓
5	25	0,0	0,500	0	0 ✓
6	28	-0,6	0,354	0	0 ✓
7	30	-1,0	0,269	0	0 ✓
8	35	-2,0	0,119	0	0 ✓
9	40	-3,0	0,047	0	0 ✓
10	45	-4,0	0,018	0	0 ✓

TABLE 10.4 – Prédiction du modèle logistique pour chaque client. Toutes les prédictions sont correctes !

10.4.3 La frontière de décision

La frontière de décision est là où $z = 0$:

$$-0,2x + 5 = 0 \quad \Rightarrow \quad x = \frac{5}{0,2} = 25 \text{ k€}$$

Interprétation

Le modèle prédit qu'un client avec un revenu **inférieur à 25 k€** a plus de 50% de chances de partir (Churn), tandis qu'un client avec un revenu **supérieur à 25 k€** restera probablement fidèle.

10.4.4 Visualisation : la sigmoïde sur les données

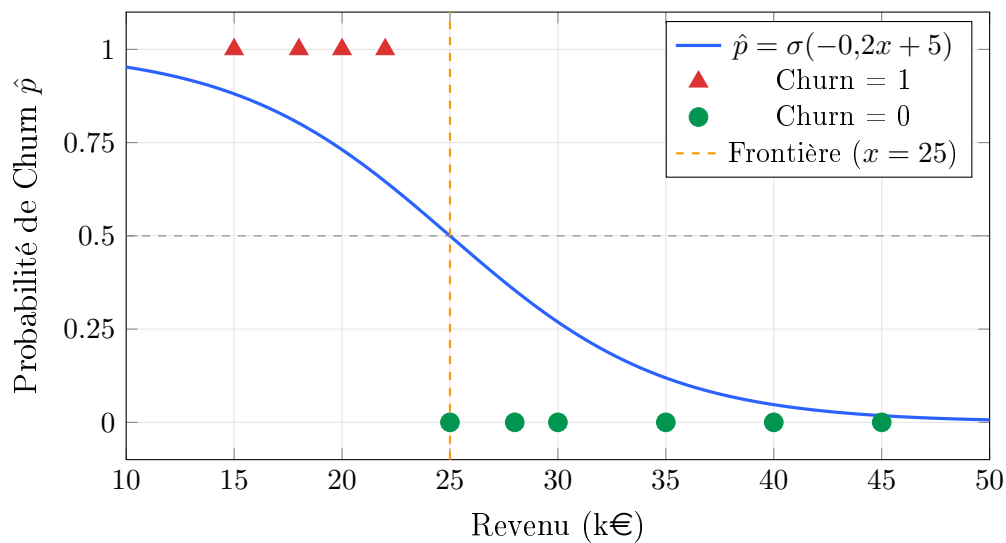


FIGURE 10.6 – La courbe sigmoïde ajustée sur les données clients. La frontière de décision est à $x = 25$ k€.

10.5 Implémentation sur Google Colab

Google Colab -- Régression logistique complète

Ouvrez un nouveau notebook sur Google Colab et copiez les cellules suivantes.

10.5.1 Cellule 1 : imports et données

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.metrics import (confusion_matrix,
6                               classification_report,
7                               roc_curve, auc,
8                               ConfusionMatrixDisplay)
9
10 # Dataset clientele
11 data = {
12     'Revenu': [15, 18, 20, 22, 25, 28, 30, 35, 40, 45],
13     'Churn': [1, 1, 1, 1, 0, 0, 0, 0, 0, 0]
14 }
15 df = pd.DataFrame(data)

```

```
14 print(df)
```

Listing 10.1 – Imports et creation du dataset

10.5.2 Cellule 2 : entraînement du modèle

```
1 # Variables
2 X = df[['Revenu']].values # shape (10, 1)
3 y = df['Churn'].values    # shape (10,)
4
5 # Creer et entrainer le modele
6 model = LogisticRegression()
7 model.fit(X, y)
8
9 # Afficher les parametres appris
10 w = model.coef_[0][0]
11 b = model.intercept_[0]
12 print(f"Poids w = {w:.4f}")
13 print(f"Biais b = {b:.4f}")
14 print(f"Frontiere de decision : Revenu = {-b/w:.2f} k euros")
```

Listing 10.2 – Entraînement de la régression logistique

10.5.3 Cellule 3 : courbe sigmoïde sur les données

```
1 # Generer des valeurs pour tracer la courbe
2 X_plot = np.linspace(10, 50, 300).reshape(-1, 1)
3 p_plot = model.predict_proba(X_plot)[:, 1]
4
5 plt.figure(figsize=(10, 6))
6
7 # Tracer la courbe sigmoïde
8 plt.plot(X_plot, p_plot, 'b-', linewidth=2,
9          label='Courbe sigmoïde')
10
11 # Tracer les points
12 churn_1 = df[df['Churn'] == 1]
13 churn_0 = df[df['Churn'] == 0]
14 plt.scatter(churn_1['Revenu'], churn_1['Churn'],
15             color='red', marker='^', s=100, zorder=5,
```

```

16         label='Churn = 1')
17 plt.scatter(churn_0['Revenu'], churn_0['Churn'],
18             color='green', marker='o', s=100, zorder=5,
19             label='Churn = 0')
20
21 # Frontiere de decision
22 frontiere = -b / w
23 plt.axvline(x=frontiere, color='orange', linestyle='--',
24             linewidth=2, label=f'Frontiere = {frontiere:.1f}')
25 plt.axhline(y=0.5, color='gray', linestyle=':', alpha=0.5)
26
27 plt.xlabel('Revenu (k euros)', fontsize=12)
28 plt.ylabel('P(Churn = 1)', fontsize=12)
29 plt.title('Regression Logistique - Prediction du Churn',
30         fontsize=14)
31 plt.legend(fontsize=11)
32 plt.grid(True, alpha=0.3)
33 plt.show()

```

Listing 10.3 – Visualisation de la courbe sigmoïde

10.5.4 Cellule 4 : prédictions et probabilités

```

1 # Predictions
2 y_pred = model.predict(X)
3 y_proba = model.predict_proba(X)[:, 1]
4
5 # Tableau de resultats
6 resultats = pd.DataFrame({
7     'Revenu': df['Revenu'],
8     'Churn_reel': y,
9     'P(Churn=1)': np.round(y_proba, 3),
10    'Prediction': y_pred
11 })
12 print(resultats)
13
14 # Tester un nouveau client
15 nouveau_revenu = np.array([[23]])
16 proba = model.predict_proba(nouveau_revenu)[0, 1]
17 pred = model.predict(nouveau_revenu)[0]
18 print(f"\nNouveau client (Revenu = 23 k euros) :")

```

```

19 print(f"   Probabilite de Churn = {proba:.3f}")
20 print(f"   Prediction : {'Churn' if pred == 1 else 'Fidele'}")

```

Listing 10.4 – Prédictions et probabilités pour chaque client

10.5.5 Cellule 5 : matrice de confusion et rapport de classification

```

1  # Matrice de confusion
2  cm = confusion_matrix(y, y_pred)
3  print("Matrice de confusion :")
4  print(cm)
5
6  # Affichage graphique
7  fig, ax = plt.subplots(figsize=(6, 5))
8  disp = ConfusionMatrixDisplay(confusion_matrix=cm,
9                                display_labels=['Fidele', 'Churn'])
10 disp.plot(ax=ax, cmap='Blues', values_format='d')
11 plt.title('Matrice de Confusion', fontsize=14)
12 plt.show()
13
14 # Rapport de classification
15 print("\nRapport de classification :")
16 print(classification_report(y, y_pred,
17                             target_names=['Fidele', 'Churn']))

```

Listing 10.5 – Matrice de confusion et classification report

10.5.6 Cellule 6 : courbe ROC et AUC

```

1  # Calcul de la courbe ROC
2  fpr, tpr, thresholds = roc_curve(y, y_proba)
3  roc_auc = auc(fpr, tpr)
4
5  # Tracer la courbe ROC
6  plt.figure(figsize=(8, 6))
7  plt.plot(fpr, tpr, 'b-', linewidth=2,
8           label=f'ROC (AUC = {roc_auc:.3f})')
9  plt.plot([0, 1], [0, 1], 'r--', linewidth=1,
10           label='Classifieur aleatoire')
11 plt.fill_between(fpr, tpr, alpha=0.1, color='blue')

```

```

12
13 plt.xlabel('Taux de Faux Positifs (FPR)', fontsize=12)
14 plt.ylabel('Taux de Vrais Positifs (TPR)', fontsize=12)
15 plt.title('Courbe ROC', fontsize=14)
16 plt.legend(fontsize=11)
17 plt.grid(True, alpha=0.3)
18 plt.show()
19
20 print(f"AUC = {roc_auc:.3f}")
21 print("Interpretation : AUC = 1.0 -> classifieur parfait")

```

Listing 10.6 – Courbe ROC et AUC

10.5.7 Cellule 7 : frontière de décision en 2D (bonus)

```

1  # Dataset banque avec 2 variables
2  data2 = {
3      'Revenu':      [2.0, 2.5, 3.0, 3.2, 3.5, 4.0, 4.5, 5.0, 5.5,
4                      6.0],
5      'Score_Credit': [450, 500, 520, 580, 600, 620, 650, 700, 720,
6                      750],
7      'Defaut':      [1,   1,   1,   1,   0,   0,   0,   0,   0,
8                      0]
9  }
10 df2 = pd.DataFrame(data2)
11
12 # Entrainer
13 model2 = LogisticRegression()
14 model2.fit(X2, y2)
15
16 # Grille pour la frontiere
17 xx, yy = np.meshgrid(
18     np.linspace(1.5, 6.5, 200),
19     np.linspace(400, 800, 200)
20 )
21 Z = model2.predict_proba(
22     np.c_[xx.ravel(), yy.ravel()]
23 )[:, 1].reshape(xx.shape)

```

```

24
25 plt.figure(figsize=(10, 7))
26 plt.contourf(xx, yy, Z, levels=50, cmap='RdYlBu_r', alpha=0.6)
27 plt.colorbar(label='P(Default = 1)')
28 plt.contour(xx, yy, Z, levels=[0.5], colors='black',
29             linewidths=2, linestyles='--')
30
31 default = df2[df2['Default'] == 1]
32 ok = df2[df2['Default'] == 0]
33 plt.scatter(default['Revenu'], default['Score_Credit'],
34             c='red', marker='^', s=120, edgecolors='black',
35             label='Default', zorder=5)
36 plt.scatter(ok['Revenu'], ok['Score_Credit'],
37             c='blue', marker='o', s=120, edgecolors='black',
38             label='Pas de default', zorder=5)
39
40 plt.xlabel('Revenu (k euros)', fontsize=12)
41 plt.ylabel('Score de Credit', fontsize=12)
42 plt.title('Frontiere de decision en 2D', fontsize=14)
43 plt.legend(fontsize=11)
44 plt.show()

```

Listing 10.7 – Frontière de décision en 2D avec deux variables

10.6 Schéma récapitulatif : phases d'entraînement et de test

PHASE D'ENTRAÎNEMENT (Training)

1. Étape 1 — Charger et séparer les données

Charger le jeu de données et le diviser en 80 % Train et 20 % Test.

Mathématiques :

$$X \in \mathbb{R}^{n \times p}, \quad y \in \{0, 1\}^n$$

Code Python :

```
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2)
```

2. Étape 2 — Normaliser

Centrer et réduire les variables explicatives.

Mathématiques :

$$x'_j = \frac{x_j - \mu_j}{\sigma_j}$$

Code Python :

```
scaler = StandardScaler()
X_train_s = scaler.fit_transform(X_train)
```

3. Étape 3 — Initialiser les poids

Initialiser tous les coefficients à zéro.

Mathématiques :

$$\beta_0 = 0, \beta_1 = 0, \dots, \beta_p = 0$$

Code Python :

```
from sklearn.linear_model import LogisticRegression
model = LogisticRegression()
```

4. Étape 4 — Minimiser la Binary Cross-Entropy

Entraîner le modèle en minimisant la fonction de coût BCE par descente de gradient.

Mathématiques :

$$J(\beta) = -\frac{1}{n} \sum [y_i \ln(\hat{p}_i) + (1 - y_i) \ln(1 - \hat{p}_i)]$$

$$\text{Mise à jour : } \beta_j := \beta_j - \alpha \frac{\partial J}{\partial \beta_j}$$

Code Python :

```
model.fit(X_train_s, y_train)
```

5. Étape 5 — Obtenir les coefficients

Récupérer les poids optimaux qui définissent la frontière de décision.

Mathématiques :

$$\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$$

Code Python :

```
print("Coefficients :", model.coef_)
print("Intercept :", model.intercept_)
```



Modèle entraîné transmis à la phase de test

PHASE DE TEST (Testing)**1. Étape 1 — Calculer z**

Calculer la combinaison linéaire pour chaque point de test.

Mathématiques :

$$z = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

Code Python :

```
X_test_s = scaler.transform(X_test)
```

2. Étape 2 — Appliquer la sigmoïde

Transformer z en probabilité.

Mathématiques :

$$\hat{p} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

Code Python :

```
probas = model.predict_proba(X_test_s)[: , 1]
```

3. Étape 3 — Appliquer le seuil

Classifier chaque observation selon le seuil de 0,5.

Mathématiques :

$$\hat{y} = \begin{cases} 1 & \text{si } \hat{p} \geq 0,5 \\ 0 & \text{sinon} \end{cases}$$

Code Python :

```
y_pred = model.predict(X_test_s)
```

4. Étape 4 — Évaluer

Calculer les métriques de classification sur les données de test.

Mathématiques :

Précision, Rappel, F1-score, AUC-ROC

Code Python :

```
from sklearn.metrics import classification_report  
print(classification_report(y_test, y_pred))
```

10.7 Exercices

Exercice 10.1 — Calcul de la sigmoïde

Calculez $\sigma(z)$ pour les valeurs suivantes de z (utilisez une calculatrice pour e^{-z}) :

1. $z = -3$
2. $z = -1$
3. $z = 0$
4. $z = 0,5$
5. $z = 2$
6. $z = 4$

Pour chaque valeur, indiquez si le modèle prédirait la classe 0 ou la classe 1 (seuil = 0,5).

Correction de l'exercice 10.1

Rappel : $\sigma(z) = \frac{1}{1 + e^{-z}}$

1. $z = -3$:

$$e^{-(-3)} = e^3 \approx 20,086 \quad \Rightarrow \quad \sigma(-3) = \frac{1}{1 + 20,086} = \frac{1}{21,086} \approx 0,047$$

$$\hat{p} = 0,047 < 0,5 \Rightarrow \text{Classe 0}$$

2. $z = -1$:

$$e^1 \approx 2,718 \quad \Rightarrow \quad \sigma(-1) = \frac{1}{1 + 2,718} = \frac{1}{3,718} \approx 0,269$$

$$\hat{p} = 0,269 < 0,5 \Rightarrow \text{Classe 0}$$

3. $z = 0$:

$$e^0 = 1 \quad \Rightarrow \quad \sigma(0) = \frac{1}{1 + 1} = \frac{1}{2} = 0,500$$

$$\hat{p} = 0,500 \geq 0,5 \Rightarrow \text{Classe 1 (exactement au seuil)}$$

4. $z = 0,5$:

$$e^{-0,5} \approx 0,607 \quad \Rightarrow \quad \sigma(0,5) = \frac{1}{1 + 0,607} = \frac{1}{1,607} \approx 0,622$$

$$\hat{p} = 0,622 > 0,5 \Rightarrow \text{Classe 1}$$

5. $z = 2$:

$$e^{-2} \approx 0,135 \quad \Rightarrow \quad \sigma(2) = \frac{1}{1 + 0,135} = \frac{1}{1,135} \approx 0,881$$

$$\hat{p} = 0,881 > 0,5 \Rightarrow \text{Classe 1}$$

6. $z = 4$:

$$e^{-4} \approx 0,0183 \Rightarrow \sigma(4) = \frac{1}{1 + 0,0183} = \frac{1}{1,0183} \approx 0,982$$

$$\hat{p} = 0,982 > 0,5 \Rightarrow \text{Classe 1}$$

Résumé : On observe bien que la sigmoïde est proche de 0 pour z très négatif, vaut 0,5 en $z = 0$, et se rapproche de 1 pour z positif.

Exercice 10.2 — Prédiction du Churn avec paramètres donnés

Un modèle de régression logistique pour prédire le Churn a été entraîné et a trouvé les paramètres :

$$w = -0,2, \quad b = 5$$

Le modèle est : $\hat{p} = \sigma(-0,2 \cdot \text{Revenu} + 5)$

1. Calculez la probabilité de Churn pour un client avec un Revenu de 20 k€.
2. Calculez la probabilité de Churn pour un client avec un Revenu de 30 k€.
3. Calculez la probabilité de Churn pour un client avec un Revenu de 40 k€.
4. Trouvez la frontière de décision : pour quel revenu la probabilité est-elle exactement 50% ?
5. Interprétez le signe du poids $w = -0,2$. Que signifie le fait que w soit négatif ?

Correction de l'exercice 10.2

Modèle : $z = -0,2x + 5, \quad \hat{p} = \sigma(z) = \frac{1}{1 + e^{-z}}$

1. **Revenu = 20 k€ :**

$$z = -0,2 \times 20 + 5 = -4 + 5 = 1,0$$

$$\hat{p} = \frac{1}{1 + e^{-1}} = \frac{1}{1 + 0,368} = \frac{1}{1,368} \approx 0,731$$

P(Churn) = 73,1% \Rightarrow prédiction : Churn (car $> 0,5$)

2. **Revenu = 30 k€ :**

$$z = -0,2 \times 30 + 5 = -6 + 5 = -1,0$$

$$\hat{p} = \frac{1}{1 + e^1} = \frac{1}{1 + 2,718} = \frac{1}{3,718} \approx 0,269$$

P(Churn) = 26,9% \Rightarrow prédiction : Fidèle (car $< 0,5$)

3. **Revenu = 40 k€ :**

$$z = -0,2 \times 40 + 5 = -8 + 5 = -3,0$$

$$\hat{p} = \frac{1}{1 + e^3} = \frac{1}{1 + 20,086} = \frac{1}{21,086} \approx 0,047$$

P(Churn) = 4,7% \Rightarrow prédiction : Fidèle (très probablement)

4. **Frontière de décision :**

$\hat{p} = 0,5$ quand $z = 0$:

$$-0,2x + 5 = 0 \quad \Rightarrow \quad 0,2x = 5 \quad \Rightarrow \quad \boxed{x = 25 \text{ k€}}$$

La frontière est à un revenu de **25 k€**.

5. **Interprétation de $w = -0,2$:**

Le poids w est **négatif**, ce qui signifie que lorsque le Revenu **augmente**, z **diminue**, et donc $\hat{p} = \sigma(z)$ **diminue** aussi.

En clair : plus le revenu est élevé, plus la probabilité de Churn est faible. C'est logique : les clients à haut revenu sont plus fidèles.

Exercice 10.3 — Dérivation guidée du gradient de la BCE

Cet exercice vous guide pas à pas dans la dérivation du gradient de la Binary Cross-Entropy.

Soit le coût pour un seul exemple :

$$\ell = -[y \ln(\hat{p}) + (1 - y) \ln(1 - \hat{p})]$$

où $\hat{p} = \sigma(z)$ et $z = wx + b$.

1. Calculez $\frac{\partial \ell}{\partial \hat{p}}$.

2. Rappelez la valeur de $\frac{\partial \hat{p}}{\partial z} = \sigma'(z)$.

3. Calculez $\frac{\partial z}{\partial w}$ et $\frac{\partial z}{\partial b}$.

4. En utilisant la règle de la chaîne, montrez que :

$$\frac{\partial \ell}{\partial w} = (\hat{p} - y) \cdot x$$

5. **Application numérique :** Pour un exemple où $x = 20$, $y = 1$ et $\hat{p} = 0,731$, calculez

$$\frac{\partial \ell}{\partial w}.$$

Correction de l'exercice 10.3

1. Calcul de $\frac{\partial \ell}{\partial \hat{p}}$:

$$\ell = -y \ln(\hat{p}) - (1 - y) \ln(1 - \hat{p})$$

On dérive par rapport à \hat{p} :

$$\frac{\partial \ell}{\partial \hat{p}} = -\frac{y}{\hat{p}} - (1 - y) \cdot \frac{-1}{1 - \hat{p}} = -\frac{y}{\hat{p}} + \frac{1 - y}{1 - \hat{p}}$$

2. Dérivée de la sigmoïde :

$$\frac{\partial \hat{p}}{\partial z} = \sigma'(z) = \sigma(z)(1 - \sigma(z)) = \hat{p}(1 - \hat{p})$$

3. Dérivées par rapport à w et b :

Comme $z = wx + b$:

$$\frac{\partial z}{\partial w} = x \quad \text{et} \quad \frac{\partial z}{\partial b} = 1$$

4. Règle de la chaîne :

$$\begin{aligned} \frac{\partial \ell}{\partial w} &= \frac{\partial \ell}{\partial \hat{p}} \cdot \frac{\partial \hat{p}}{\partial z} \cdot \frac{\partial z}{\partial w} \\ &= \left(-\frac{y}{\hat{p}} + \frac{1 - y}{1 - \hat{p}} \right) \cdot \hat{p}(1 - \hat{p}) \cdot x \\ &= \left(\frac{-y(1 - \hat{p}) + (1 - y)\hat{p}}{\hat{p}(1 - \hat{p})} \right) \cdot \hat{p}(1 - \hat{p}) \cdot x \\ &= (-y + y\hat{p} + \hat{p} - y\hat{p}) \cdot x \\ &= (\hat{p} - y) \cdot x \quad \boxed{\checkmark} \end{aligned} \tag{10.3}$$

5. Application numérique :

$x = 20$, $y = 1$, $\hat{p} = 0,731$:

$$\frac{\partial \ell}{\partial w} = (0,731 - 1) \times 20 = (-0,269) \times 20 = -5,38$$

Interprétation : Le gradient est **négatif**, ce qui signifie que w doit **augmenter** (car on fait $w := w - \alpha \cdot (-5,38) = w + 5,38\alpha$). C'est logique : le modèle prédit $\hat{p} = 0,731$ alors que le vrai label est 1, donc il faut augmenter w pour rapprocher \hat{p}

de 1.

Exercice 10.4 — Implémentation Python sur le dataset Iris

Sur Google Colab, implémentez une régression logistique sur le célèbre dataset **Iris** en suivant ces étapes :

1. Chargez le dataset Iris avec `sklearn.datasets.load_iris()`.
2. Gardez uniquement les classes 0 et 1 (Setosa et Versicolor) pour avoir un problème binaire.
3. Utilisez les deux premières variables (longueur et largeur du sépale).
4. Divisez les données en 70% entraînement / 30% test.
5. Entraînez un modèle `LogisticRegression`.
6. Affichez la précision (*accuracy*) sur le jeu de test.
7. Tracez la frontière de décision en 2D avec les points colorés par classe.
8. Affichez la matrice de confusion.

Correction de l'exercice 10.4

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.datasets import load_iris
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import (accuracy_score, confusion_matrix,
7                               ConfusionMatrixDisplay)
8
9 # 1. Charger le dataset
10 iris = load_iris()
11 X = iris.data
12 y = iris.target
13
14 # 2. Garder uniquement classes 0 et 1
15 mask = y <= 1
16 X = X[mask][:, :2] # 2 premières variables seulement
17 y = y[mask]
18 print(f"Nombre d'exemples : {len(y)}")
19 print(f"Classe 0 (Setosa) : {sum(y==0)}")
20 print(f"Classe 1 (Versicolor) : {sum(y==1)}")
```

```

21
22 # 3. Separer train / test
23 X_train, X_test, y_train, y_test = train_test_split(
24     X, y, test_size=0.3, random_state=42
25 )
26
27 # 4. Entraîner le modele
28 model = LogisticRegression()
29 model.fit(X_train, y_train)
30
31 # 5. Evaluer
32 y_pred = model.predict(X_test)
33 acc = accuracy_score(y_test, y_pred)
34 print(f"\nAccuracy sur le test : {acc:.2%}")
35
36 # 6. Frontiere de decision en 2D
37 x_min, x_max = X[:, 0].min() - 0.5, X[:, 0].max() + 0.5
38 y_min, y_max = X[:, 1].min() - 0.5, X[:, 1].max() + 0.5
39 xx, yy = np.meshgrid(
40     np.linspace(x_min, x_max, 300),
41     np.linspace(y_min, y_max, 300)
42 )
43 Z = model.predict(np.c_[xx.ravel(), yy.ravel()])
44 Z = Z.reshape(xx.shape)
45
46 plt.figure(figsize=(10, 7))
47 plt.contourf(xx, yy, Z, alpha=0.3, cmap='RdYlBu')
48 plt.scatter(X[y==0, 0], X[y==0, 1],
49             c='blue', marker='o', s=80,
50             edgecolors='black', label='Setosa')
51 plt.scatter(X[y==1, 0], X[y==1, 1],
52             c='red', marker='^', s=80,
53             edgecolors='black', label='Versicolor')
54 plt.xlabel('Longueur du sepale (cm)', fontsize=12)
55 plt.ylabel('Largeur du sepale (cm)', fontsize=12)
56 plt.title('Frontiere de decision - Regression Logistique',
57           fontsize=14)
58 plt.legend(fontsize=11)
59 plt.grid(True, alpha=0.3)
60 plt.show()

```

```
61
62 # 7. Matrice de confusion
63 cm = confusion_matrix(y_test, y_pred)
64 fig, ax = plt.subplots(figsize=(6, 5))
65 disp = ConfusionMatrixDisplay(confusion_matrix=cm,
66                               display_labels=['Setosa', 'Versicolor'])
67 disp.plot(ax=ax, cmap='Blues', values_format='d')
68 plt.title('Matrice de Confusion - Iris', fontsize=14)
69 plt.show()
```

Listing 10.8 – Régression logistique sur Iris (correction complète)

Chapitre 11

K Plus Proches Voisins (KNN)

11.1 Hands-On : diagnostiquer un patient

Situation réelle — Quel diagnostic pour ce nouveau patient ?

Un médecin généraliste possède un historique de 8 patients. Pour chacun, il a mesuré deux symptômes :

- **Niveau de fièvre** (en degrés au-dessus de 37°C, noté x_1)
- **Intensité de la toux** (échelle de 0 à 10, noté x_2)

Le diagnostic final est soit **Grippe** soit **Rhume**.

Patient	Fièvre (x_1)	Toux (x_2)	Diagnostic
P1	1,0	2,0	Rhume
P2	1,5	3,0	Rhume
P3	2,0	2,5	Rhume
P4	0,5	1,5	Rhume
P5	3,0	7,0	Grippe
P6	3,5	8,0	Grippe
P7	2,5	6,5	Grippe
P8	3,0	8,5	Grippe
Nouveau	2,0	5,0	?

Un nouveau patient arrive avec une fièvre de 2,0 et une toux de 5,0. **A-t-il la grippe ou un rhume ?**

Regardons graphiquement où se situe ce nouveau patient par rapport aux autres :

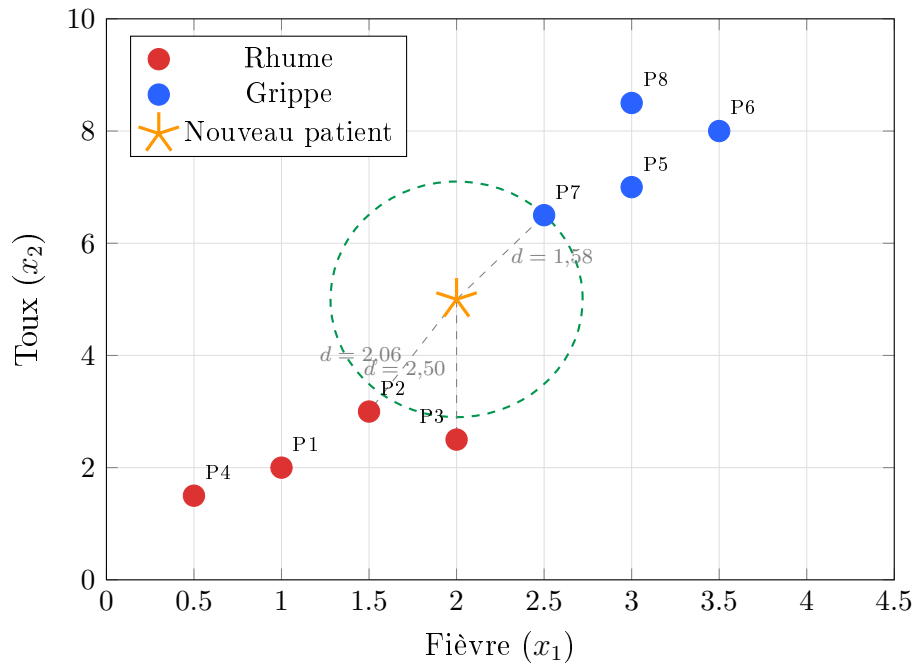


FIGURE 11.1 – Les 8 patients et le nouveau patient (étoile). Le cercle vert montre les 3 plus proches voisins ($K = 3$).

En cherchant les **3 voisins les plus proches** ($K = 3$) du nouveau patient, on trouve :

1. **P7** (Grippe) à une distance de $\approx 1,58$
2. **P2** (Rhume) à une distance de $\approx 2,06$
3. **P3** (Rhume) à une distance de $\approx 2,50$

Vote majoritaire : 2 Rhumes contre 1 Grippe \Rightarrow on prédit **Rhume**.

Observation clé

Nous n'avons construit **aucun modèle, aucune équation**. On a simplement regardé les voisins les plus proches et on a **voté**. C'est exactement le principe du KNN !

11.2 Intuition derrière le KNN

Dis-moi qui sont tes voisins ; je te dirai qui tu es

Imagine que tu arrives dans une **nouvelle ville** et que tu cherches un bon restaurant. Tu n'as pas Internet. Que fais-tu ? Tu demandes aux **personnes les plus proches** de toi dans la rue !

- Si tu demandes à **1 seule personne** ($K = 1$), tu peux tomber sur quelqu'un qui a des goûts bizarres.
- Si tu demandes à **3 personnes** ($K = 3$), c'est plus fiable : tu choisis le restaurant recommandé par la **majorité**.

- Si tu demandes à **toute la ville** ($K = n$), tu obtiendras la chaîne de fast-food la plus populaire, mais pas forcément le meilleur restaurant du quartier.

Le KNN fait exactement la même chose : pour classer un nouveau point, il regarde les K points les plus proches et fait un **vote majoritaire**.

K Plus Proches Voisins (KNN)

Le **K-Nearest Neighbors** (KNN) est un algorithme de classification (et de régression) qui :

1. **Mémorise** toutes les données d'entraînement (aucun calcul pendant l'entraînement)
2. Pour classer un nouveau point, **calcule la distance** entre ce point et **tous** les points d'entraînement
3. Sélectionne les K **voisins les plus proches**
4. Attribue la **classe majoritaire** parmi ces K voisins

Apprentissage paresseux (*Lazy Learning*)

Le KNN est un algorithme **paresseux**. Contrairement à la régression logistique qui « apprend » des paramètres $(\beta_0, \beta_1, \dots)$, le KNN ne construit **aucun modèle**. Il se contente de **stocker toutes les données** et ne travaille qu'au moment de la prédiction. Conséquence :

- **Entraînement** : instantané (il n'y a rien à faire)
- **Prédiction** : lente (il faut calculer les distances à tous les points)

11.2.1 L'effet du choix de K

Le choix de K change complètement la prédiction. Regardons le même point avec $K = 1$, $K = 3$ et $K = 7$:

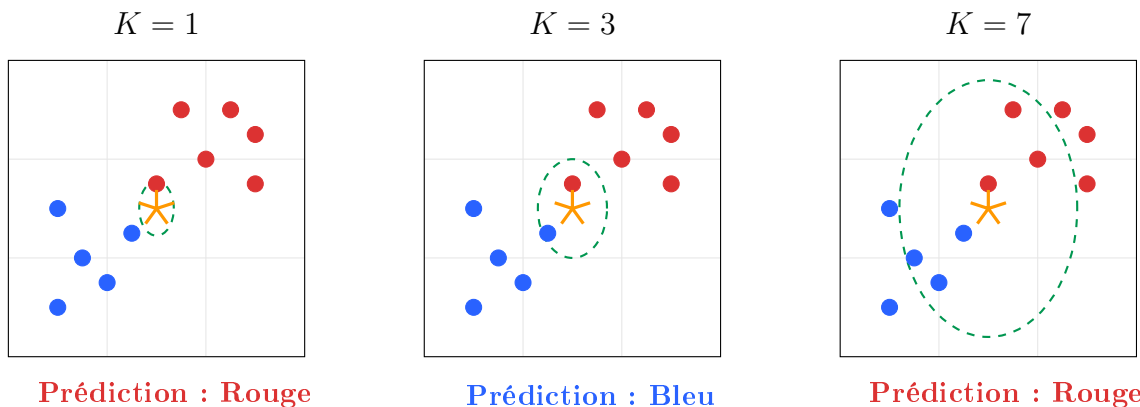


FIGURE 11.2 – Le même point (étoile orange) est classé différemment selon la valeur de K .

11.3 Développement mathématique détaillé

11.3.1 Les métriques de distance

Qu'est-ce qu'une distance ?

Pour trouver les « voisins les plus proches », il faut d'abord définir ce que signifie « **proche** ». On a besoin d'une manière de mesurer à quel point deux points sont éloignés : c'est la notion de **distance**.

Distance

Une **distance** est une fonction $d(\mathbf{a}, \mathbf{b})$ qui mesure « à quel point » deux points \mathbf{a} et \mathbf{b} sont éloignés. Elle doit vérifier :

1. $d(\mathbf{a}, \mathbf{b}) \geq 0$ (toujours positive ou nulle)
2. $d(\mathbf{a}, \mathbf{b}) = 0 \iff \mathbf{a} = \mathbf{b}$ (nulle si et seulement si les points sont identiques)
3. $d(\mathbf{a}, \mathbf{b}) = d(\mathbf{b}, \mathbf{a})$ (symétrique)
4. $d(\mathbf{a}, \mathbf{c}) \leq d(\mathbf{a}, \mathbf{b}) + d(\mathbf{b}, \mathbf{c})$ (inégalité triangulaire)

Distance euclidienne

Distance euclidienne

La **distance euclidienne** est la distance « en ligne droite » entre deux points. Pour deux points $\mathbf{a} = (a_1, a_2, \dots, a_p)$ et $\mathbf{b} = (b_1, b_2, \dots, b_p)$ dans un espace à p dimensions :

$$d_{\text{eucl}}(\mathbf{a}, \mathbf{b}) = \sqrt{\sum_{j=1}^p (a_j - b_j)^2} = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2 + \dots + (a_p - b_p)^2}$$

En **2 dimensions**, c'est tout simplement le **théorème de Pythagore** ! Si $\mathbf{a} = (x_1, y_1)$ et $\mathbf{b} = (x_2, y_2)$:

$$d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

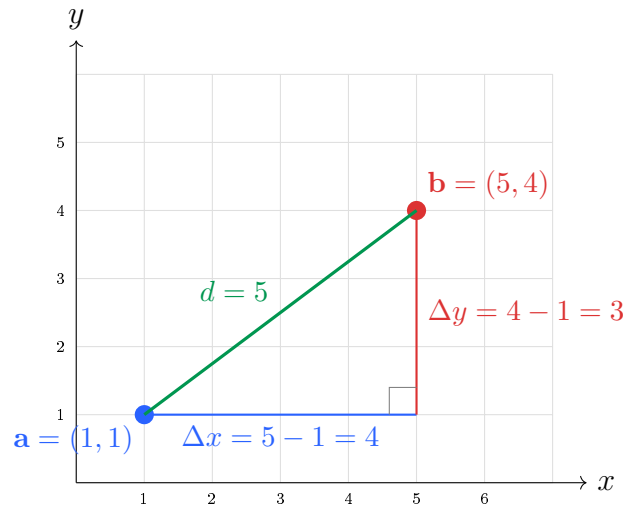


FIGURE 11.3 – La distance euclidienne est l’hypoténuse du triangle rectangle : $d = \sqrt{4^2 + 3^2} = \sqrt{16 + 9} = \sqrt{25} = 5$.

Calcul détaillé

Calculons la distance entre le point $\mathbf{a} = (1, 1)$ et $\mathbf{b} = (5, 4)$ pas à pas :

$$\begin{aligned}
 d(\mathbf{a}, \mathbf{b}) &= \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2} \\
 &= \sqrt{(1 - 5)^2 + (1 - 4)^2} \\
 &= \sqrt{(-4)^2 + (-3)^2} \\
 &= \sqrt{16 + 9} \\
 &= \sqrt{25} = \boxed{5}
 \end{aligned}$$

C’est le célèbre triplet pythagoricien (3, 4, 5) !

Distance de Manhattan

Distance de Manhattan

La **distance de Manhattan** (ou distance L_1) mesure la distance en suivant les axes, comme si on marchait dans les rues d’une ville quadrillée (on ne peut pas couper en diagonale) :

$$d_{\text{manh}}(\mathbf{a}, \mathbf{b}) = \sum_{j=1}^p |a_j - b_j| = |a_1 - b_1| + |a_2 - b_2| + \cdots + |a_p - b_p|$$

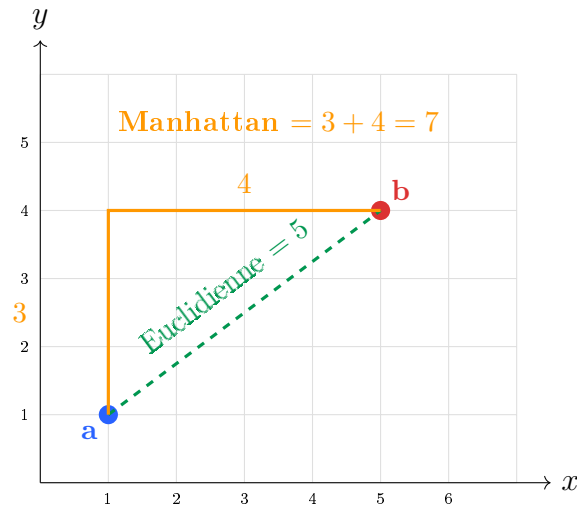


FIGURE 11.4 – Euclidienne (ligne droite verte, $d = 5$) vs Manhattan (chemin orange en équerre, $d = 7$).

Calcul de la distance de Manhattan

Pour $\mathbf{a} = (1, 1)$ et $\mathbf{b} = (5, 4)$:

$$d_{\text{manh}}(\mathbf{a}, \mathbf{b}) = |1 - 5| + |1 - 4| = |-4| + |-3| = 4 + 3 = \boxed{7}$$

La distance de Manhattan est **toujours supérieure ou égale** à la distance euclidienne.

Distance de Minkowski (généralisation)

Distance de Minkowski

La **distance de Minkowski** généralise les deux distances précédentes avec un paramètre $q \geq 1$:

$$d_{\text{mink}}(\mathbf{a}, \mathbf{b}) = \left(\sum_{j=1}^p |a_j - b_j|^q \right)^{1/q}$$

Cas particuliers :

- $q = 1 \Rightarrow$ Distance de Manhattan
- $q = 2 \Rightarrow$ Distance euclidienne
- $q \rightarrow \infty \Rightarrow$ Distance de Chebyshev : $d_{\infty} = \max_j |a_j - b_j|$

Comparaison des trois distances

Pour $\mathbf{a} = (1, 1)$ et $\mathbf{b} = (5, 4)$:

Distance	Formule appliquée	Résultat
Manhattan ($q = 1$)	$ 4 + 3 = 7$	7
Euclidienne ($q = 2$)	$\sqrt{4^2 + 3^2} = \sqrt{25}$	5
Chebyshev ($q \rightarrow \infty$)	$\max(4 , 3)$	4

On a toujours : $d_\infty \leq d_2 \leq d_1$.

11.3.2 L'algorithme KNN pas à pas

Voici les étapes exactes de l'algorithme KNN :

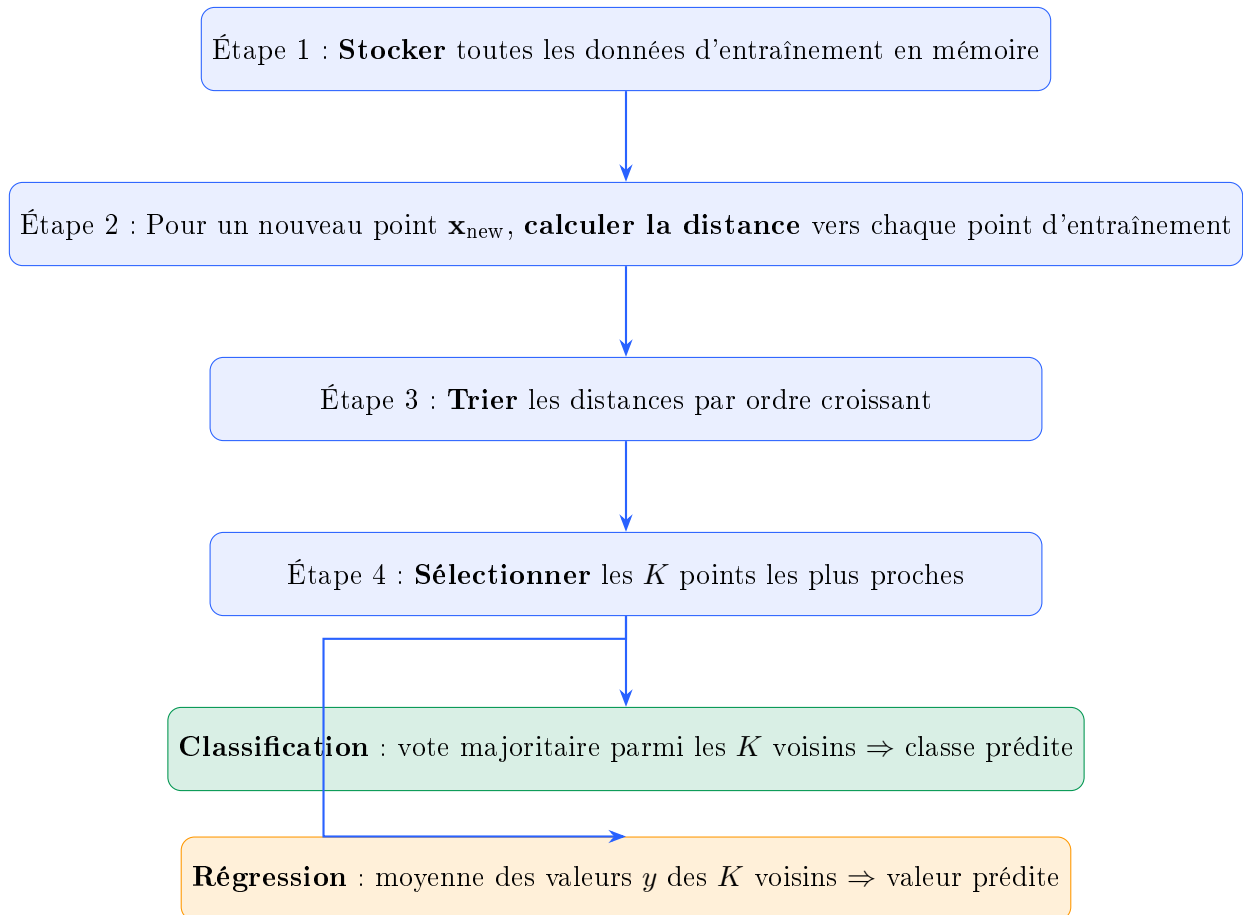


FIGURE 11.5 – Les 5 étapes de l'algorithme KNN.

Mathématiquement, pour un nouveau point \mathbf{x}_{new} , on calcule :

$$d_i = d(\mathbf{x}_{\text{new}}, \mathbf{x}_i) \quad \text{pour } i = 1, 2, \dots, n$$

On trie les d_i et on sélectionne les indices \mathcal{N}_K des K plus petites distances. La prédiction

est :

$$\hat{y} = \underset{c}{\text{mode}}\{y_i : i \in \mathcal{N}_K\} \quad (\text{classification — vote majoritaire})$$

11.3.3 Mise à l'échelle des variables (*Feature Scaling*)

CRITIQUE pour le KNN : normaliser les variables !

Le KNN est **très sensible** à l'échelle des variables. Puisqu'il repose sur des calculs de **distance**, une variable avec de grandes valeurs va **dominer** complètement les autres.

Pourquoi c'est critique — un exemple parlant

Considérons deux variables : **Revenu** (en €) et **Âge** (en années).

Client	Revenu (€)	Âge (années)
A	25 000	25
B	60 000	30

Distance euclidienne **sans normalisation** :

$$d = \sqrt{(25000 - 60000)^2 + (25 - 30)^2} = \sqrt{(-35000)^2 + (-5)^2} = \sqrt{1\,225\,000\,000 + 25} \approx 35\,000$$

Le revenu contribue pour 1 225 000 000 et l'âge pour seulement 25. L'âge est **totalement ignoré** ! C'est comme si on ne prenait en compte que le revenu.

Normalisation Min-Max

Normalisation Min-Max

La normalisation Min-Max ramène chaque variable dans l'intervalle $[0, 1]$:

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

où x_{\min} et x_{\max} sont les valeurs minimale et maximale de la variable dans les données d'entraînement.

Standardisation (z-score)**Standardisation**

La standardisation centre les données autour de 0 avec un écart-type de 1 :

$$x' = \frac{x - \mu}{\sigma}$$

où μ est la moyenne et σ est l'écart-type de la variable.

Avant / Après la normalisation

Reprenons nos données avec 4 clients :

Client	Avant		Après Min-Max	
	Revenu	Âge	Revenu'	Âge'
A	25 000	25	0,00	0,00
B	60 000	30	1,00	0,25
C	40 000	45	0,43	1,00
D	30 000	35	0,14	0,50

Calcul pour le client C :

$$\text{— Revenu}' = \frac{40000 - 25000}{60000 - 25000} = \frac{15000}{35000} = 0,43$$

$$\text{— } \hat{\text{Age}}' = \frac{45 - 25}{45 - 25} = \frac{20}{20} = 1,00$$

Distance A-B après normalisation :

$$d' = \sqrt{(0,00 - 1,00)^2 + (0,00 - 0,25)^2} = \sqrt{1 + 0,0625} = \sqrt{1,0625} \approx 1,03$$

Maintenant les deux variables contribuent **équitablement** au calcul de la distance !

11.3.4 Choisir la valeur de K

Le choix de K est un **hyperparamètre** crucial :

Valeur de K	Comportement	Risque
$K = 1$	Très sensible à chaque point	Surapprentissage (overfitting) — épouse le bruit
K petit	Frontière de décision complexe	Sensible au bruit
K grand	Frontière de décision lisse	Sous-apprentissage (underfitting) — trop général
$K = n$	Prédit toujours la classe majoritaire	Inutile

Règle empirique : on commence souvent avec $K \approx \sqrt{n}$ où n est le nombre d'observations, puis on utilise la **validation croisée** pour trouver le K optimal.

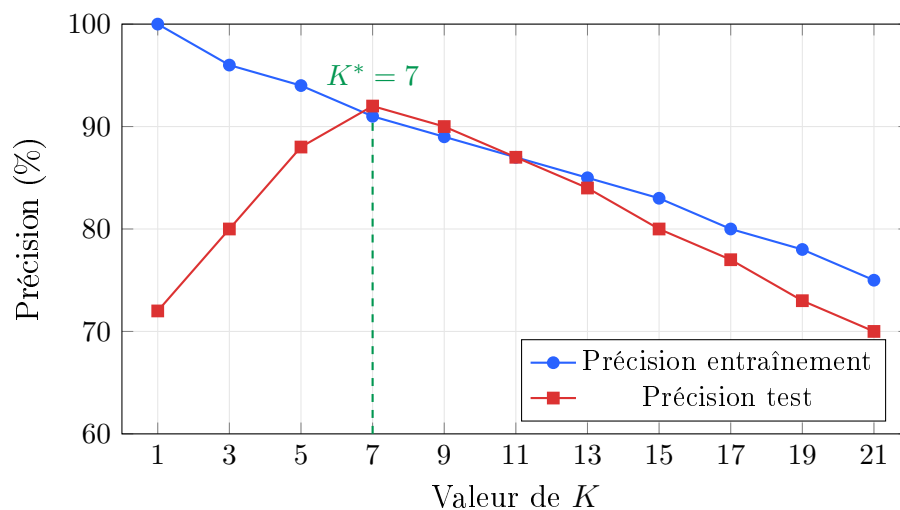


FIGURE 11.6 – Précision en fonction de K : le K optimal est celui qui maximise la précision sur les données de test.

Astuce pratique

Il est recommandé de choisir un K **impair** pour la classification binaire, afin d'éviter les égalités lors du vote majoritaire.

11.3.5 KNN pondéré par la distance

Dans le KNN classique, tous les K voisins ont le même poids dans le vote. Mais un voisin **très proche** devrait compter plus qu'un voisin à peine dans le cercle !

KNN pondéré

Dans le KNN pondéré, chaque voisin reçoit un poids inversement proportionnel à sa distance :

$$w_i = \frac{1}{d(\mathbf{x}_{\text{new}}, \mathbf{x}_i)}$$

Le vote devient : pour chaque classe c , on calcule la somme des poids des voisins de cette classe :

$$\text{Score}(c) = \sum_{\substack{i \in \mathcal{N}_K \\ y_i = c}} w_i = \sum_{\substack{i \in \mathcal{N}_K \\ y_i = c}} \frac{1}{d_i}$$

On prédit la classe avec le score le plus élevé.

KNN pondéré sur notre exemple médical

Reprenons les 3 plus proches voisins de notre nouveau patient :

Voisin	Classe	Distance d_i	Poids $w_i = 1/d_i$
P7	Grippe	1,58	$1/1,58 = 0,633$
P2	Rhume	2,06	$1/2,06 = 0,485$
P3	Rhume	2,50	$1/2,50 = 0,400$

Scores pondérés :

- $\text{Score}(\text{Grippe}) = 0,633$
- $\text{Score}(\text{Rhume}) = 0,485 + 0,400 = 0,885$

Rhume a le score le plus élevé \Rightarrow même prédiction qu'avant : **Rhume**. Mais dans d'autres cas, le KNN pondéré peut changer le résultat !

11.3.6 KNN pour la régression

Le KNN n'est pas limité à la classification. Pour la **régression**, au lieu de voter, on prend la **moyenne** des valeurs y des K voisins.

KNN pour la régression

Pour un nouveau point \mathbf{x}_{new} , la prédiction en régression est :

$$\hat{y} = \frac{1}{K} \sum_{i \in \mathcal{N}_K} y_i \quad (\text{moyenne simple})$$

Ou en version pondérée :

$$\hat{y} = \frac{\sum_{i \in \mathcal{N}_K} w_i \cdot y_i}{\sum_{i \in \mathcal{N}_K} w_i} \quad \text{où } w_i = \frac{1}{d_i}$$

11.4 Application sur le dataset clientèle

Reprenons notre jeu de données fil rouge avec 10 clients :

ID	Client	Revenu (€/mois)	Ancienneté (mois)	Dépense (€/mois)	Churn
1	Fatima	30	12	20	Oui
2	Ahmed	55	40	50	Non
3	Youssef	25	6	15	Oui
4	Khadija	45	30	40	Non
5	Amina	35	8	18	Oui
6	Badr	50	35	45	Non
7	Nora	40	25	35	Non
8	Hadi	20	5	12	Oui
9	Ines	60	45	55	Non
10	Jalil	28	10	16	Oui

TABLE 11.1 – Le dataset clientèle — fil rouge du cours.

11.4.1 Classification : prédire le Churn d'un nouveau client

Un nouveau client arrive : **Karim** avec un Revenu de 38 €/mois et une Ancienneté de 15 mois. **Va-t-il cherner ?**

Étape 1 : Normalisation Min-Max.

Pour le Revenu : min = 20, max = 60.

Pour l'Ancienneté : min = 5, max = 45.

ID	Client	Rev.	Anc.	Rev.'	Anc.'
1	Fatima	30	12	$\frac{30-20}{40} = 0,25$	$\frac{12-5}{40} = 0,175$
2	Ahmed	55	40	0,875	0,875
3	Youssef	25	6	0,125	0,025
4	Khadija	45	30	0,625	0,625
5	Amina	35	8	0,375	0,075
6	Badr	50	35	0,750	0,750
7	Nora	40	25	0,500	0,500
8	Hadi	20	5	0,000	0,000
9	Ines	60	45	1,000	1,000
10	Jalil	28	10	0,200	0,125
–	Karim	38	15	0,450	0,250

Étape 2 : Calculer les distances euclidiennes (sur les données normalisées).

Karim normalisé : (0,450; 0,250).

$$\begin{aligned}
 d(\text{Karim}, \text{Fatima}) &= \sqrt{(0,450 - 0,250)^2 + (0,250 - 0,175)^2} = \sqrt{0,04 + 0,0056} = \sqrt{0,0456} \approx 0,214 \\
 d(\text{Karim}, \text{Ahmed}) &= \sqrt{(0,450 - 0,875)^2 + (0,250 - 0,875)^2} = \sqrt{0,1806 + 0,3906} = \sqrt{0,5712} \approx 0,756 \\
 d(\text{Karim}, \text{Youssef}) &= \sqrt{(0,450 - 0,125)^2 + (0,250 - 0,025)^2} = \sqrt{0,1056 + 0,0506} = \sqrt{0,1562} \approx 0,395 \\
 d(\text{Karim}, \text{Khadija}) &= \sqrt{(0,450 - 0,625)^2 + (0,250 - 0,625)^2} = \sqrt{0,0306 + 0,1406} = \sqrt{0,1712} \approx 0,414 \\
 d(\text{Karim}, \text{Amina}) &= \sqrt{(0,450 - 0,375)^2 + (0,250 - 0,075)^2} = \sqrt{0,0056 + 0,0306} = \sqrt{0,0362} \approx 0,190 \\
 d(\text{Karim}, \text{Badr}) &= \sqrt{(0,450 - 0,750)^2 + (0,250 - 0,750)^2} = \sqrt{0,09 + 0,25} = \sqrt{0,34} \approx 0,583 \\
 d(\text{Karim}, \text{Nora}) &= \sqrt{(0,450 - 0,500)^2 + (0,250 - 0,500)^2} = \sqrt{0,0025 + 0,0625} = \sqrt{0,065} \approx 0,255 \\
 d(\text{Karim}, \text{Hadi}) &= \sqrt{(0,450 - 0,000)^2 + (0,250 - 0,000)^2} = \sqrt{0,2025 + 0,0625} = \sqrt{0,265} \approx 0,515 \\
 d(\text{Karim}, \text{Ines}) &= \sqrt{(0,450 - 1,000)^2 + (0,250 - 1,000)^2} = \sqrt{0,3025 + 0,5625} = \sqrt{0,865} \approx 0,930 \\
 d(\text{Karim}, \text{Jalil}) &= \sqrt{(0,450 - 0,200)^2 + (0,250 - 0,125)^2} = \sqrt{0,0625 + 0,0156} = \sqrt{0,0781} \approx 0,279
 \end{aligned}$$

Étape 3 : Trier et sélectionner les $K = 3$ plus proches.

Rang	Client	Distance	Churn
1	Amina	0,190	Oui
2	Fatima	0,214	Oui
3	Nora	0,255	Non
4	Jalil	0,279	Oui
5	Youssef	0,395	Oui
...

Étape 4 : Vote majoritaire.

— Churn = Oui : **2 votes** (Amina, Fatima)

— Churn = Non : **1 vote** (Nora)

$$\hat{y}_{\text{Karim}} = \text{Oui (Churn)}$$

11.4.2 Régression KNN : prédire la Dépense

Utilisons les mêmes $K = 3$ voisins pour prédire la **Dépense** de Karim :

Voisin	Distance	Dépense (€)
Amina	0,190	18
Fatima	0,214	20
Nora	0,255	35

Moyenne simple :

$$\hat{y} = \frac{18 + 20 + 35}{3} = \frac{73}{3} \approx \boxed{24,33 \text{ €/mois}}$$

Moyenne pondérée :

$$w_{\text{Amina}} = 1/0,190 = 5,263 \quad ; \quad w_{\text{Fatima}} = 1/0,214 = 4,673 \quad ; \quad w_{\text{Nora}} = 1/0,255 = 3,922$$

$$\hat{y} = \frac{5,263 \times 18 + 4,673 \times 20 + 3,922 \times 35}{5,263 + 4,673 + 3,922} = \frac{94,73 + 93,46 + 137,27}{13,858} = \frac{325,46}{13,858} \approx \boxed{23,49 \text{ €/mois}}$$

La version pondérée donne un résultat légèrement différent car Amina (la plus proche) a un poids plus important et sa dépense est de 18 €.

11.5 Implémentation sur Google Colab

Google Colab -- KNN complet sur le dataset clientèle

Ouvrez un nouveau notebook et suivez les cellules ci-dessous.

```

1  import numpy as np
2  import pandas as pd
3  import matplotlib.pyplot as plt
4  from sklearn.neighbors import KNeighborsClassifier,
    KNeighborsRegressor
5  from sklearn.preprocessing import StandardScaler, MinMaxScaler
6  from sklearn.model_selection import train_test_split,
    cross_val_score
7  from sklearn.metrics import (accuracy_score, confusion_matrix,
8                               ConfusionMatrixDisplay,
                               classification_report)
9
10 # Dataset clientele
11 data = {
12     'Client':      ['Fatima', 'Ahmed', 'Youssef', 'Khadija', 'Amina',
13                    'Badr', 'Nora', 'Hadi', 'Ines', 'Jalil'],
14     'Revenu':      [30, 55, 25, 45, 35, 50, 40, 20, 60, 28],
15     'Anciennete':  [12, 40, 6, 30, 8, 35, 25, 5, 45, 10],
16     'Depense':     [20, 50, 15, 40, 18, 45, 35, 12, 55, 16],
17     'Churn':       [1, 0, 1, 0, 1, 0, 0, 1, 0, 1]
18 }
19 df = pd.DataFrame(data)
20 print(df)

```

Listing 11.1 – Cellule 1 — Imports et création du dataset

```

1  # Variables explicatives et cible
2  X = df[['Revenu', 'Anciennete']].values
3  y_cls = df['Churn'].values          # Classification
4  y_reg = df['Depense'].values        # Regression
5
6  # Standardisation (moyenne=0, ecart-type=1)
7  scaler = StandardScaler()
8  X_scaled = scaler.fit_transform(X)
9

```

```

10 print("Avant normalisation :")
11 print(X[:3])
12 print("\nAprès normalisation :")
13 print(X_scaled[:3])

```

Listing 11.2 – Cellule 2 — Normalisation avec StandardScaler

```

1  # Creer et entrainer le modele KNN (K=3)
2  knn_cls = KNeighborsClassifier(n_neighbors=3, metric='euclidean')
3  knn_cls.fit(X_scaled, y_cls)
4
5  # Predire pour un nouveau client (Karim : Revenu=38, Anciennete=15)
6  karim = np.array([[38, 15]])
7  karim_scaled = scaler.transform(karim)  # IMPORTANT : meme
      transformation
8
9  prediction = knn_cls.predict(karim_scaled)
10 proba = knn_cls.predict_proba(karim_scaled)
11
12 print(f"Prediction pour Karim : {'Churn' if prediction[0]==1 else
      'Fidele'}")
13 print(f"Probabilites : Fidele={proba[0][0]:.2f},
      Churn={proba[0][1]:.2f}")

```

Listing 11.3 – Cellule 3 — KNN Classification : entraînement et prédiction

```

1  # Creer une grille de points pour la frontiere de decision
2  h = 0.05  # pas de la grille
3  x_min, x_max = X_scaled[:, 0].min() - 1, X_scaled[:, 0].max() + 1
4  y_min, y_max = X_scaled[:, 1].min() - 1, X_scaled[:, 1].max() + 1
5  xx, yy = np.meshgrid(np.arange(x_min, x_max, h),
6                        np.arange(y_min, y_max, h))
7
8  # Prediction sur chaque point de la grille
9  Z = knn_cls.predict(np.c_[xx.ravel(), yy.ravel()])
10 Z = Z.reshape(xx.shape)
11
12 # Tracer
13 plt.figure(figsize=(10, 7))
14 plt.contourf(xx, yy, Z, alpha=0.3, cmap='RdBu_r')
15 plt.scatter(X_scaled[y_cls==0, 0], X_scaled[y_cls==0, 1],
16             c='blue', marker='o', s=120, edgecolors='k',

```

```

17         label='Fidèle', zorder=3)
18 plt.scatter(X_scaled[y_cls==1, 0], X_scaled[y_cls==1, 1],
19             c='red', marker='^', s=120, edgecolors='k',
20             label='Churn', zorder=3)
21 plt.scatter(karim_scaled[0, 0], karim_scaled[0, 1],
22             c='gold', marker='*', s=300, edgecolors='k',
23             label='Karim (nouveau)', zorder=4)
24
25 # Noms des clients
26 for i, name in enumerate(df['Client']):
27     plt.annotate(name, (X_scaled[i,0]+0.05, X_scaled[i,1]+0.05),
28                 fontsize=8)
29
30 plt.xlabel('Revenu (standardise)')
31 plt.ylabel('Anciennete (standardisee)')
32 plt.title('KNN (K=3) - Frontiere de decision')
33 plt.legend()
34 plt.grid(alpha=0.3)
35 plt.show()

```

Listing 11.4 – Cellule 4 — Visualisation de la frontière de décision

```

1 # Tester differentes valeurs de K
2 k_range = range(1, 10)
3 scores = []
4
5 for k in k_range:
6     knn = KNeighborsClassifier(n_neighbors=k)
7     # Validation croisee avec 5 folds (ou moins si peu de donnees)
8     cv_scores = cross_val_score(knn, X_scaled, y_cls, cv=3,
9                                 scoring='accuracy')
10    scores.append(cv_scores.mean())
11
12 # Tracer la courbe
13 plt.figure(figsize=(8, 5))
14 plt.plot(k_range, scores, 'bo-', linewidth=2, markersize=8)
15 plt.xlabel('Valeur de K')
16 plt.ylabel('Precision moyenne (validation croisee)')
17 plt.title('Choix du K optimal')
18 plt.xticks(list(k_range))
19 plt.grid(alpha=0.3)
20

```

```

21 # Marquer le K optimal
22 best_k = list(k_range)[np.argmax(scores)]
23 plt.axvline(x=best_k, color='red', linestyle='--',
24             label=f'K optimal = {best_k}')
25 plt.legend()
26 plt.show()
27 print(f"K optimal : {best_k} (precision = {max(scores):.2f})")

```

Listing 11.5 – Cellule 5 — Trouver le K optimal par validation croisée

```

1 # SANS normalisation
2 knn_sans = KNeighborsClassifier(n_neighbors=3)
3 knn_sans.fit(X, y_cls) # donnees brutes
4 pred_sans = knn_sans.predict(karim)
5
6 # AVEC normalisation
7 knn_avec = KNeighborsClassifier(n_neighbors=3)
8 knn_avec.fit(X_scaled, y_cls)
9 pred_avec = knn_avec.predict(karim_scaled)
10
11 print("=== Comparaison avec/sans normalisation ===")
12 print(f"Sans normalisation : {'Churn' if pred_sans[0]==1 else
13       'Fidele'}")
13 print(f"Avec normalisation : {'Churn' if pred_avec[0]==1 else
14       'Fidele'}")
14 print("\n> La normalisation peut changer la prediction !")

```

Listing 11.6 – Cellule 6 — Comparer avec et sans normalisation

```

1 # KNN pour la regression
2 knn_reg = KNeighborsRegressor(n_neighbors=3, weights='distance')
3 knn_reg.fit(X_scaled, y_reg)
4
5 # Predire la depense de Karim
6 depense_pred = knn_reg.predict(karim_scaled)
7 print(f"Depense predite pour Karim : {depense_pred[0]:.2f}
8       euros/mois")

```

Listing 11.7 – Cellule 7 — KNN Régression : prédire la Dépense

```

1 # Predictions sur tout le dataset
2 y_pred = knn_cls.predict(X_scaled)

```

```

3
4 # Matrice de confusion
5 cm = confusion_matrix(y_cls, y_pred)
6 disp = ConfusionMatrixDisplay(confusion_matrix=cm,
7                               display_labels=['Fidèle', 'Churn'])
8
9 fig, ax = plt.subplots(figsize=(6, 5))
10 disp.plot(ax=ax, cmap='Blues', values_format='d')
11 plt.title('Matrice de confusion - KNN (K=3)')
12 plt.show()
13
14 # Rapport de classification
15 print(classification_report(y_cls, y_pred,
16                             target_names=['Fidèle', 'Churn']))

```

Listing 11.8 – Cellule 8 — Matrice de confusion (sur les données d'entraînement)

11.6 Schéma récapitulatif : phases d'entraînement et de test

PHASE D'ENTRAÎNEMENT (Training)

1. Étape 1 — Charger et séparer les données

Charger le dataset et le séparer en 80% Train, 20% Test.

Mathématiques :

$$X \in \mathbb{R}^{n \times p}, \quad y \in \{0, 1\}^n$$

Code Python :

```
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2)
```

2. Étape 2 — Normaliser les features (indispensable pour KNN)

Appliquer un **StandardScaler** — c'est **indispensable** pour le KNN car il repose sur le calcul des distances. Sans normalisation, les variables à grande échelle dominent complètement.

Mathématiques :

$$x'_j = \frac{x_j - \mu_j}{\sigma_j} \quad \text{pour que toutes les variables aient la même échelle}$$

Code Python :

```
scaler = StandardScaler()
X_train_s = scaler.fit_transform(X_train)
```

3. Étape 3 — Choisir K par validation croisée

Déterminer le nombre de voisins K par **validation croisée** (tester plusieurs valeurs et garder celle qui maximise la précision).

Mathématiques :

$$K^* = \arg \max_K \text{Accuracy}_{CV}(K)$$

Code Python :

```
from sklearn.neighbors import KNeighborsClassifier
for k in range(1, 21):
    knn = KNeighborsClassifier(n_neighbors=k)
    scores = cross_val_score(knn, X_train_s, y_train, cv=5)
```

4. Étape 4 — Stocker les données (lazy learning)

Mémoriser **toutes** les données d'entraînement en mémoire. Le KNN n'apprend **aucun modèle** : c'est un algorithme paresseux (*lazy learning*).

Mathématiques :

Aucun paramètre à calculer ! Le modèle = toutes les données de Train en mémoire.

Code Python :

```
knn = KNeighborsClassifier(n_neighbors=K_best)
knn.fit(X_train_s, y_train) # stocke X_train en memoire
```



Modèle entraîné transmis à la phase de test

PHASE DE TEST (Testing)

1. Étape 1 — Calculer les distances

Pour chaque point de Test, calculer la distance avec **tous** les points de Train.

Mathématiques :

$$d(\mathbf{x}, \mathbf{x}') = \sqrt{\sum_{j=1}^p (x_j - x'_j)^2}$$

Code Python :

```
X_test_s = scaler.transform(X_test)
# knn calcule les distances automatiquement
```

2. Étape 2 — Sélectionner K voisins et voter

Trier les distances par ordre croissant, sélectionner les K plus proches voisins et

effectuer un **vote majoritaire**.

Mathématiques :

$$\hat{y} = \text{mode}\{y_{(1)}, y_{(2)}, \dots, y_{(K)}\} \quad (\text{classe majoritaire des } K \text{ plus proches})$$

Code Python :

```
y_pred = knn.predict(X_test_s)
```

3. Étape 3 — Évaluer

Calculer les métriques de performance — précision, rappel, F1-score, matrice de confusion.

Mathématiques :

Précision, Rappel, F1-score.

Code Python :

```
print(classification_report(y_test, y_pred))
print(confusion_matrix(y_test, y_pred))
```

11.7 Exercices

Exercice 11.1 — KNN à la main (K

Soit les 6 points suivants dans un espace 2D :

Point	x_1	x_2	Classe
A	1	2	Rouge
B	2	4	Rouge
C	3	1	Rouge
D	6	5	Bleu
E	7	7	Bleu
F	8	6	Bleu

Un nouveau point $P = (4, 3)$ doit être classé avec $K = 3$.

1. Calculez la distance euclidienne entre P et chacun des 6 points.
2. Identifiez les 3 plus proches voisins.
3. Quelle est la classe prédite par vote majoritaire ?

Correction de l'exercice 11.1**1. Calcul des distances euclidiennes :**

Le point $P = (4, 3)$.

$$d(P, A) = \sqrt{(4-1)^2 + (3-2)^2} = \sqrt{9+1} = \sqrt{10} \approx 3,162$$

$$d(P, B) = \sqrt{(4-2)^2 + (3-4)^2} = \sqrt{4+1} = \sqrt{5} \approx 2,236$$

$$d(P, C) = \sqrt{(4-3)^2 + (3-1)^2} = \sqrt{1+4} = \sqrt{5} \approx 2,236$$

$$d(P, D) = \sqrt{(4-6)^2 + (3-5)^2} = \sqrt{4+4} = \sqrt{8} \approx 2,828$$

$$d(P, E) = \sqrt{(4-7)^2 + (3-7)^2} = \sqrt{9+16} = \sqrt{25} = 5,000$$

$$d(P, F) = \sqrt{(4-8)^2 + (3-6)^2} = \sqrt{16+9} = \sqrt{25} = 5,000$$

2. Les 3 plus proches voisins (tri par distance croissante) :

Rang	Point	Distance	Classe
1	B	2,236	Rouge
2	C	2,236	Rouge
3	D	2,828	Bleu
4	A	3,162	Rouge
5	E	5,000	Bleu
6	F	5,000	Bleu

Les 3 plus proches voisins sont : B (Rouge), C (Rouge) et D (Bleu).

3. Vote majoritaire :

— Rouge : 2 votes (B, C)

— Bleu : 1 vote (D)

$$\hat{y}_P = \text{Rouge}$$

Exercice 11.2 — Normalisation à la main

Soit le dataset suivant :

Observation	Taille (cm)	Poids (kg)
1	160	55
2	170	70
3	180	85
4	175	65

1. Appliquez la normalisation **Min-Max** aux deux variables. Écrivez le tableau normalisé.
2. Calculez la moyenne μ et l'écart-type σ de chaque variable, puis appliquez la **standardisation** (z -score). Écrivez le tableau standardisé.
3. Calculez la distance euclidienne entre l'observation 1 et l'observation 3 **avant** et **après** la normalisation Min-Max. Que constatez-vous ?

Correction de l'exercice 11.2

1. Normalisation Min-Max :

Pour la Taille : $\min = 160$, $\max = 180$.

Pour le Poids : $\min = 55$, $\max = 85$.

Formule : $x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$

Obs.	Taille'	Poids'
1	$\frac{160-160}{20} = 0,00$	$\frac{55-55}{30} = 0,00$
2	$\frac{170-160}{20} = 0,50$	$\frac{70-55}{30} = 0,50$
3	$\frac{180-160}{20} = 1,00$	$\frac{85-55}{30} = 1,00$
4	$\frac{175-160}{20} = 0,75$	$\frac{65-55}{30} = 0,33$

2. Standardisation (z -score) :

$$\text{Taille} : \mu_T = \frac{160+170+180+175}{4} = 171,25 ; \sigma_T = \sqrt{\frac{(160-171,25)^2 + (170-171,25)^2 + (180-171,25)^2 + (175-171,25)^2}{4}}$$

$$\sigma_T = \sqrt{\frac{126,56+1,56+76,56+14,06}{4}} = \sqrt{\frac{218,75}{4}} = \sqrt{54,69} \approx 7,40$$

$$\text{Poids} : \mu_P = \frac{55+70+85+65}{4} = 68,75 ; \sigma_P = \sqrt{\frac{(55-68,75)^2 + (70-68,75)^2 + (85-68,75)^2 + (65-68,75)^2}{4}}$$

$$\sigma_P = \sqrt{\frac{189,06+1,56+264,06+14,06}{4}} = \sqrt{\frac{468,75}{4}} = \sqrt{117,19} \approx 10,82$$

Obs.	Taille'	Poids'
1	$\frac{160-171,25}{7,40} = -1,52$	$\frac{55-68,75}{10,82} = -1,27$
2	$\frac{170-171,25}{7,40} = -0,17$	$\frac{70-68,75}{10,82} = +0,12$
3	$\frac{180-171,25}{7,40} = +1,18$	$\frac{85-68,75}{10,82} = +1,50$
4	$\frac{175-171,25}{7,40} = +0,51$	$\frac{65-68,75}{10,82} = -0,35$

3. Comparaison des distances entre Obs. 1 et Obs. 3 :*Avant normalisation :*

$$d = \sqrt{(160 - 180)^2 + (55 - 85)^2} = \sqrt{400 + 900} = \sqrt{1300} \approx 36,06$$

Le poids contribue $\frac{900}{1300} = 69\%$ et la taille $\frac{400}{1300} = 31\%$. Le poids domine.

Après normalisation Min-Max :

$$d' = \sqrt{(0,00 - 1,00)^2 + (0,00 - 1,00)^2} = \sqrt{1 + 1} = \sqrt{2} \approx 1,41$$

Chaque variable contribue **également** (50% chacune). C'est l'effet souhaité !

Exercice 11.3 — L'effet de K sur la prédiction

Reprenons les données de l'exercice 11.1 (6 points et le nouveau point $P = (4, 3)$).

1. Quelle est la prédiction avec $K = 1$?
2. Quelle est la prédiction avec $K = 3$? (déjà calculée, rappel)
3. Quelle est la prédiction avec $K = 5$?
4. Les trois valeurs de K donnent-elles le même résultat ? Commentez.

Correction de l'exercice 11.3

Rappelons les distances calculées dans l'exercice 11.1, triées par ordre croissant :

Rang	Point	Distance	Classe
1	B	2,236	Rouge
2	C	2,236	Rouge
3	D	2,828	Bleu
4	A	3,162	Rouge
5	E	5,000	Bleu
6	F	5,000	Bleu

1. $K = 1$: Le plus proche voisin est B (Rouge).

$$\hat{y} = \text{Rouge}$$

2. $K = 3$: Les 3 plus proches sont B (R), C (R), D (B) \Rightarrow 2 Rouge, 1 Bleu.

$$\hat{y} = \text{Rouge}$$

3. $K = 5$: Les 5 plus proches sont B (R), C (R), D (B), A (R), E (B) \Rightarrow 3 Rouge, 2 Bleu.

$$\hat{y} = \text{Rouge}$$

4. Commentaire :

Dans cet exemple, les trois valeurs de K donnent **la même prédiction** (Rouge). Cela n'est pas toujours le cas ! Le point P est relativement proche de la zone rouge, ce qui explique la robustesse de la prédiction. En général, lorsque le point est proche de la **frontière de décision** (entre les deux classes), le choix de K a un impact beaucoup plus important.

Exercice 11.4 — Python : KNN sur le dataset Iris

Réalisez les étapes suivantes dans un notebook Google Colab :

1. Chargez le dataset **Iris** depuis Scikit-learn (`load_iris`).
2. Séparez les données en train (70%) et test (30%).
3. Normalisez les données avec **StandardScaler**.
4. Entraînez un KNN avec $K = 5$ et calculez la précision sur le test.
5. Tracez la courbe de précision pour K allant de 1 à 25 (utilisez la validation croisée).
6. Trouvez le K optimal et affichez la matrice de confusion correspondante.
7. **Bonus** : tracez la frontière de décision en utilisant uniquement les 2 premières variables (*sepal length* et *sepal width*).

Correction de l'exercice 11.4

```
# === Imports ===
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import load_iris
from sklearn.neighbors import KNeighborsClassifier
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split,
    cross_val_score
from sklearn.metrics import (accuracy_score, confusion_matrix,
    ConfusionMatrixDisplay)

# === 1. Charger le dataset Iris ===
iris = load_iris()
X = iris.data                # 150 échantillons, 4 variables
```

```

y = iris.target          # 3 classes (0, 1, 2)
print(f"Dimensions : {X.shape}")
print(f"Classes : {iris.target_names}")

# === 2. Separation train/test (70/30) ===
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.3, random_state=42, stratify=y)

# === 3. Normalisation ===
scaler = StandardScaler()
X_train_s = scaler.fit_transform(X_train)
X_test_s = scaler.transform(X_test)

# === 4. KNN avec K=5 ===
knn = KNeighborsClassifier(n_neighbors=5)
knn.fit(X_train_s, y_train)
y_pred = knn.predict(X_test_s)
print(f"\nPrecision (K=5) : {accuracy_score(y_test,
    y_pred):.4f}")

# === 5. Courbe precision vs K ===
k_range = range(1, 26)
cv_scores = []
for k in k_range:
    knn_k = KNeighborsClassifier(n_neighbors=k)
    scores = cross_val_score(knn_k, X_train_s, y_train,
        cv=5, scoring='accuracy')
    cv_scores.append(scores.mean())

plt.figure(figsize=(10, 5))
plt.plot(k_range, cv_scores, 'bo-', linewidth=2)
plt.xlabel('K')
plt.ylabel('Precision (CV 5-folds)')
plt.title('Precision en fonction de K - Dataset Iris')
plt.xticks(list(k_range))
plt.grid(alpha=0.3)

# === 6. K optimal et matrice de confusion ===
best_k = list(k_range)[np.argmax(cv_scores)]
plt.axvline(x=best_k, color='red', linestyle='--',

```

```

        label=f'K optimal = {best_k}')
plt.legend()
plt.show()

# Entraîner avec K optimal
knn_best = KNeighborsClassifier(n_neighbors=best_k)
knn_best.fit(X_train_s, y_train)
y_pred_best = knn_best.predict(X_test_s)

print(f"\nK optimal = {best_k}")
print(f"Precision = {accuracy_score(y_test, y_pred_best):.4f}")

cm = confusion_matrix(y_test, y_pred_best)
disp = ConfusionMatrixDisplay(confusion_matrix=cm,
                              display_labels=iris.target_names)
fig, ax = plt.subplots(figsize=(7, 5))
disp.plot(ax=ax, cmap='Blues')
plt.title(f'Matrice de confusion (K={best_k})')
plt.show()

# == 7. Bonus : Frontiere de decision (2 premieres variables)
==
X_2d = X[:, :2] # sepal length, sepal width
X_train_2d, X_test_2d, y_train_2d, y_test_2d = train_test_split(
    X_2d, y, test_size=0.3, random_state=42, stratify=y)

scaler_2d = StandardScaler()
X_train_2d_s = scaler_2d.fit_transform(X_train_2d)
X_test_2d_s = scaler_2d.transform(X_test_2d)

knn_2d = KNeighborsClassifier(n_neighbors=best_k)
knn_2d.fit(X_train_2d_s, y_train_2d)

# Grille
h = 0.02
x_min, x_max = X_train_2d_s[:,0].min()-1,
               X_train_2d_s[:,0].max()+1
y_min, y_max = X_train_2d_s[:,1].min()-1,
               X_train_2d_s[:,1].max()+1
xx, yy = np.meshgrid(np.arange(x_min, x_max, h),

```

```
        np.arange(y_min, y_max, h))
Z = knn_2d.predict(np.c_[xx.ravel(), yy.ravel()])
Z = Z.reshape(xx.shape)

plt.figure(figsize=(10, 7))
plt.contourf(xx, yy, Z, alpha=0.3, cmap='viridis')
colors = ['blue', 'green', 'red']
for cls_idx, cls_name in enumerate(iris.target_names):
    mask = y_train_2d == cls_idx
    plt.scatter(X_train_2d_s[mask, 0], X_train_2d_s[mask, 1],
                c=colors[cls_idx], label=cls_name,
                edgecolors='k', s=60)
plt.xlabel('Sepal Length (standardise)')
plt.ylabel('Sepal Width (standardise)')
plt.title(f'Frontiere de decision KNN (K={best_k}) - Iris')
plt.legend()
plt.grid(alpha=0.3)
plt.show()
```

Listing 11.9 – Exercice 11.4 — Solution complète

Chapitre 12

Arbres de Décision

12.1 Hands-On : Doit-on jouer au tennis aujourd'hui ?

Mise en situation : prédire s'il faut jouer au tennis

Imaginez que vous tenez un journal météo depuis 14 jours. Chaque jour, vous notez les conditions météorologiques et si vous avez joué au tennis ou non. Après 14 jours, vous aimeriez qu'un algorithme vous dise automatiquement : « aujourd'hui, avec ces conditions, dois-je jouer au tennis ? »

C'est exactement ce que fait un **arbre de décision** : il analyse vos données passées et construit un **arbre de questions** pour prendre la meilleure décision possible.

Voici les données collectées sur 14 jours (c'est le célèbre *PlayTennis dataset* utilisé en Machine Learning depuis les années 1980) :

Jour	Ciel	Température	Humidité	Vent	Jouer ?
J1	Sunny	Hot	High	Weak	Non
J2	Sunny	Hot	High	Strong	Non
J3	Overcast	Hot	High	Weak	Oui
J4	Rain	Mild	High	Weak	Oui
J5	Rain	Cool	Normal	Weak	Oui
J6	Rain	Cool	Normal	Strong	Non
J7	Overcast	Cool	Normal	Strong	Oui
J8	Sunny	Mild	High	Weak	Non
J9	Sunny	Cool	Normal	Weak	Oui
J10	Rain	Mild	Normal	Weak	Oui
J11	Sunny	Mild	Normal	Strong	Oui
J12	Overcast	Mild	High	Strong	Oui
J13	Overcast	Hot	Normal	Weak	Oui
J14	Rain	Mild	High	Strong	Non

TABLE 12.1 – Le dataset *PlayTennis* — 14 observations, 4 attributs, 1 cible.

Observez le dataset

- **14 observations** (14 jours)
- **4 attributs** (features) : Ciel, Température, Humidité, Vent
- **1 cible** (label) : Jouer ? (Oui / Non) — c'est un problème de **classification binaire**
- Résultat global : 9 « Oui » et 5 « Non »

La question est : **quel attribut faut-il regarder en premier pour prendre la meilleure décision ?** Le Ciel ? Le Vent ? La Température ? L'Humidité ? C'est exactement ce que l'algorithme d'arbre de décision va déterminer, et il le fera de manière **mathématiquement optimale**.

12.2 Intuition : l'arbre de décision, un jeu de 20 questions

L'arbre de décision

Vous connaissez sûrement le jeu des « 20 questions » : quelqu'un pense à quelque chose, et vous devez le deviner en posant des questions oui/non.

Un bon joueur pose les questions qui **éliminent le plus de possibilités** à chaque étape. Par exemple :

- « Est-ce un être vivant ? » élimine la moitié des possibilités
- « Est-ce rouge ? » n'élimine presque rien (mauvaise question !)

Un arbre de décision fait exactement la même chose : à chaque nœud, il pose la question qui **sépare le mieux** les données. On suit les branches jusqu'à atteindre une **feuille** (la décision finale).

Définition — Arbre de décision

Un **arbre de décision** est un modèle de Machine Learning qui prend des décisions en suivant une série de **questions hiérarchiques** sur les attributs des données. Il se compose de :

- **Racine** (root) : le premier nœud, la première question posée
- **Nœuds internes** : chaque question intermédiaire
- **Branches** : les réponses possibles à chaque question
- **Feuilles** (leaves) : les décisions finales (les prédictions)

Voici l'arbre de décision optimal pour notre dataset tennis :

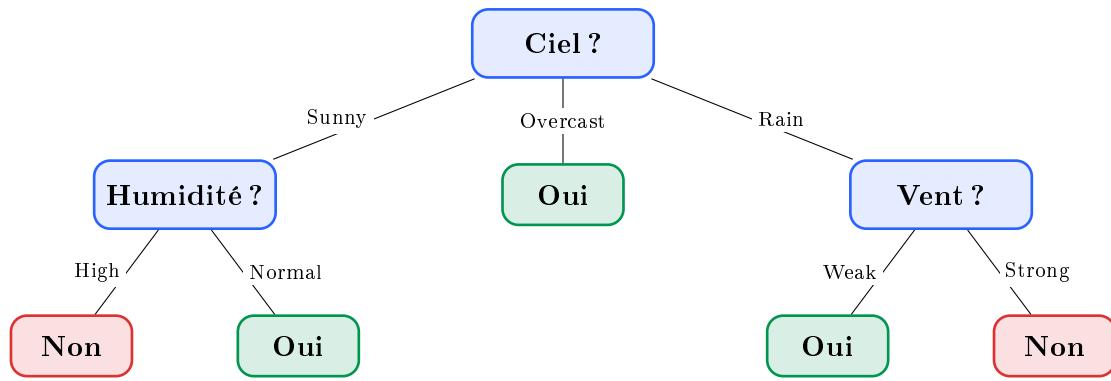


FIGURE 12.1 – L’arbre de décision optimal pour le dataset *PlayTennis*. Chaque nœud bleu est une question, chaque feuille verte ou rouge est une décision.

Comment lire cet arbre ?

Pour un nouveau jour avec : Ciel = Sunny, Humidité = Normal, on suit le chemin :

1. **Ciel ?** → Sunny → on va à gauche
2. **Humidité ?** → Normal → on va à droite
3. On arrive à la feuille **Oui** ⇒ on joue au tennis !

C’est aussi simple que de suivre un organigramme. Pas de formule compliquée à évaluer : juste une série de questions.

Analogie : le diagnostic médical

Un médecin suit souvent un arbre de décision mental :

1. « Avez-vous de la fièvre ? » → Si oui...
2. « Avez-vous mal à la gorge ? » → Si oui...
3. « Depuis combien de jours ? » → Si plus de 3 jours...
4. ⇒ « C’est probablement une angine, je vous prescris un test. »

L’arbre de décision automatise exactement ce processus, mais en choisissant les questions de manière **mathématiquement optimale** pour être le plus précis possible.

Mais comment l’algorithme sait-il quelle question poser en premier ? Pourquoi « Ciel » plutôt que « Vent » ? La réponse se trouve dans la notion d’**entropie** et de **gain d’information**, que nous allons maintenant dériver pas à pas.

12.3 Dérivation mathématique détaillée

12.3.1 L'entropie de Shannon : mesurer l'incertitude

Qu'est-ce que l'incertitude ?

Avant toute formule, comprenons l'idée avec des exemples simples.

Pièces de monnaie et incertitude

Situation 1 : Vous lancez une pièce **équilibrée** (50% face, 50% pile). Avant de lancer, vous n'avez **aucune idée** du résultat. L'incertitude est **maximale**.

Situation 2 : Vous lancez une pièce **truquée** (99% face, 1% pile). Avant de lancer, vous êtes **presque sûr** que ce sera face. L'incertitude est **très faible**.

Situation 3 : Vous lancez une pièce avec 100% face. Vous savez **exactement** ce qui va arriver. L'incertitude est **nulle**.

⇒ Plus le résultat est prévisible, moins il y a d'incertitude.

L'**entropie** est une mesure mathématique de cette incertitude. Elle a été inventée par **Claude Shannon** en 1948, le père de la théorie de l'information. Mais avant de voir la formule, rappelons d'abord deux notions essentielles.

Rappel : qu'est-ce qu'une probabilité ?

Probabilité

La **probabilité** d'un événement est la fraction des cas où cet événement se produit :

$$P(\text{événement}) = \frac{\text{nombre de cas favorables}}{\text{nombre total de cas}}$$

Exemples :

- Sur 14 jours, on joue 9 fois : $P(\text{Oui}) = \frac{9}{14} \approx 0,643$
- Sur 14 jours, on ne joue pas 5 fois : $P(\text{Non}) = \frac{5}{14} \approx 0,357$
- Vérification : $P(\text{Oui}) + P(\text{Non}) = \frac{9}{14} + \frac{5}{14} = 1 \checkmark$

Rappel : qu'est-ce que le logarithme ?

Pour comprendre l'entropie, il faut comprendre le **logarithme en base 2**, noté \log_2 . Pas de panique, c'est plus simple qu'on ne le pense !

Le logarithme en base 2

Le logarithme en base 2 répond à la question : « **À quelle puissance faut-il élever 2 pour obtenir ce nombre ?** »

$$\log_2(x) = n \iff 2^n = x$$

Exemples de logarithme en base 2

Question	Réponse	Parce que...
$\log_2(8) = ?$	3	$2^3 = 8$
$\log_2(4) = ?$	2	$2^2 = 4$
$\log_2(2) = ?$	1	$2^1 = 2$
$\log_2(1) = ?$	0	$2^0 = 1$
$\log_2(0,5) = ?$	-1	$2^{-1} = \frac{1}{2} = 0,5$
$\log_2(0,25) = ?$	-2	$2^{-2} = \frac{1}{4} = 0,25$
$\log_2(16) = ?$	4	$2^4 = 16$

Propriétés clés du \log_2

- $\log_2(1) = 0$ (le log de 1 est toujours 0)
- Si $0 < x < 1$ alors $\log_2(x) < 0$ (le log d'un nombre entre 0 et 1 est négatif)
- $\log_2(0)$ n'existe pas (tend vers $-\infty$) — mais par convention, on pose $0 \cdot \log_2(0) = 0$
- Plus x est grand, plus $\log_2(x)$ est grand (fonction croissante)

Pourquoi utilise-t-on \log_2 en entropie ? Parce que le logarithme en base 2 mesure l'information en **bits** (binary digits). Un bit, c'est la quantité d'information nécessaire pour distinguer entre deux choix équiprobables. Si vous lancez une pièce équilibrée, il vous faut exactement **1 bit** pour coder le résultat (0 ou 1, face ou pile).

La formule de l'entropie

Entropie de Shannon

Soit un ensemble S contenant des exemples répartis en c classes, avec p_i la proportion d'exemples de la classe i . L'**entropie** de S est :

$$H(S) = - \sum_{i=1}^c p_i \cdot \log_2(p_i)$$

où :

- p_i = proportion des exemples de la classe i dans S
- La somme porte sur toutes les classes $i = 1, 2, \dots, c$
- **Convention** : si $p_i = 0$, alors $0 \cdot \log_2(0) = 0$

L'entropie vaut entre 0 (certitude totale) et $\log_2(c)$ (incertitude maximale).

Pour un problème binaire (2 classes), la formule se simplifie :

Entropie binaire

Si la proportion de la classe positive est p et celle de la classe négative est $(1 - p)$:

$$H(p) = -p \cdot \log_2(p) - (1 - p) \cdot \log_2(1 - p)$$

Exemples calculés pas à pas

Exemple 1 : tous dans la même classe (certitude totale)

Si un ensemble contient **10 Oui et 0 Non** :

$$p_{\text{Oui}} = \frac{10}{10} = 1, \quad p_{\text{Non}} = \frac{0}{10} = 0$$

$$H = -1 \cdot \log_2(1) - 0 \cdot \log_2(0) = -1 \cdot 0 - 0 = \boxed{0}$$

Entropie = 0. Logique : il n'y a **aucune incertitude**, on sait que la réponse est Oui.

Exemple 2 : répartition 50/50 (incertitude maximale)

Si un ensemble contient **5 Oui et 5 Non** :

$$p_{\text{Oui}} = \frac{5}{10} = 0,5, \quad p_{\text{Non}} = \frac{5}{10} = 0,5$$

$$H = -0,5 \cdot \log_2(0,5) - 0,5 \cdot \log_2(0,5)$$

Or $\log_2(0,5) = -1$ (car $2^{-1} = 0,5$), donc :

$$H = -0,5 \times (-1) - 0,5 \times (-1) = 0,5 + 0,5 = \boxed{1}$$

Entropie = 1 bit. C'est le **maximum** pour 2 classes : on ne sait vraiment pas du tout !

Exemple 3 : répartition 75/25

Si un ensemble contient **9 Oui et 3 Non** (12 exemples) :

$$p_{\text{Oui}} = \frac{9}{12} = 0,75, \quad p_{\text{Non}} = \frac{3}{12} = 0,25$$

$$H = -0,75 \cdot \log_2(0,75) - 0,25 \cdot \log_2(0,25)$$

Calculons chaque terme :

$$\text{— } \log_2(0,75) = \log_2\left(\frac{3}{4}\right) \approx -0,415 \quad \text{donc} \quad -0,75 \times (-0,415) = 0,311$$

$$\text{— } \log_2(0,25) = \log_2\left(\frac{1}{4}\right) = -2 \quad \text{donc} \quad -0,25 \times (-2) = 0,5$$

$$H = 0,311 + 0,5 = \boxed{0,811 \text{ bits}}$$

L'incertitude est assez élevée, mais moins que 50/50.

Courbe de l'entropie binaire

Entropie binaire : $H(p) = -p \log_2(p) - (1-p) \log_2(1-p)$

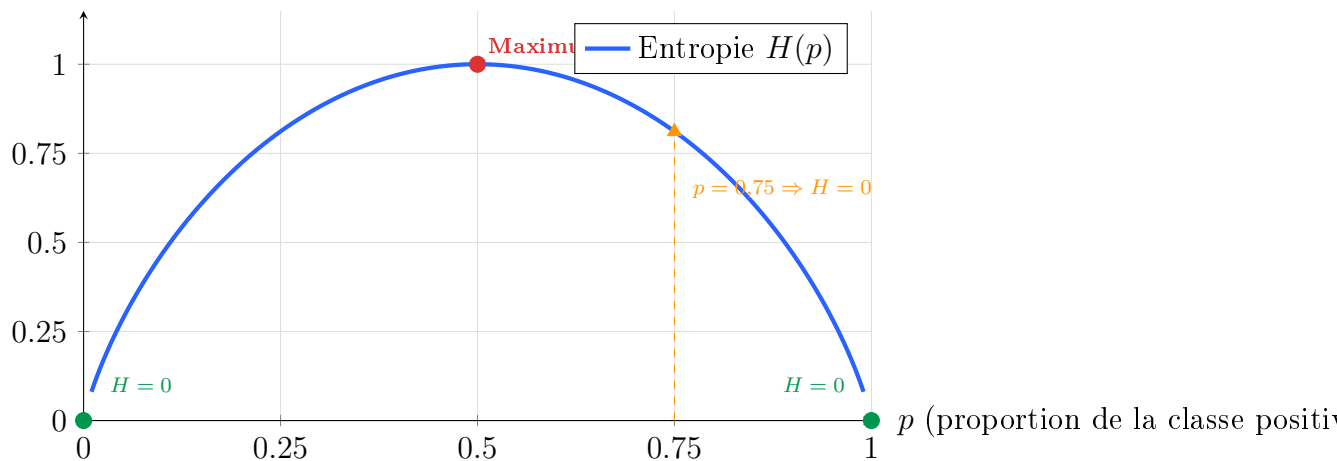


FIGURE 12.2 – La courbe de l'entropie binaire. Le maximum est atteint à $p = 0,5$ (incertitude totale). L'entropie vaut 0 quand $p = 0$ ou $p = 1$ (certitude totale).

12.3.2 Le gain d'information

Maintenant que nous savons mesurer l'incertitude, la question est : **quelle question réduit le plus cette incertitude ?** C'est le rôle du **gain d'information** (Information Gain).

Gain d'information

Le **gain d'information** d'un attribut A par rapport à un ensemble S mesure à quel point l'attribut A réduit l'entropie de S :

$$\text{Gain}(S, A) = H(S) - \sum_{v \in \text{Valeurs}(A)} \frac{|S_v|}{|S|} \cdot H(S_v)$$

où :

- $H(S)$ = entropie de l'ensemble avant la séparation
- $\text{Valeurs}(A)$ = l'ensemble des valeurs possibles de l'attribut A
- S_v = sous-ensemble de S où l'attribut A a la valeur v
- $|S_v|$ = nombre d'exemples dans S_v , $|S|$ = nombre total d'exemples

Le gain est toujours ≥ 0 . Plus il est élevé, plus l'attribut est informatif.

Interprétation du gain d'information

Le gain d'information répond à la question : « **de combien de bits l'incertitude diminue-t-elle si on pose cette question ?** »

Le terme $\sum \frac{|S_v|}{|S|} \cdot H(S_v)$ est la **moyenne pondérée** des entropies des sous-ensembles obtenus après la séparation. C'est l'**entropie résiduelle** : l'incertitude qu'il reste après avoir posé la question.

Gain = Incertitude avant – Incertitude après

Calcul détaillé sur le dataset tennis

Étape 1 : Entropie globale $H(\text{Play})$

Sur les 14 jours : 9 Oui et 5 Non.

$$p_{\text{Oui}} = \frac{9}{14}, \quad p_{\text{Non}} = \frac{5}{14}$$

$$\begin{aligned} H(\text{Play}) &= -\frac{9}{14} \cdot \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \cdot \log_2\left(\frac{5}{14}\right) \\ &= -0,643 \times \log_2(0,643) - 0,357 \times \log_2(0,357) \\ &= -0,643 \times (-0,637) - 0,357 \times (-1,486) \end{aligned}$$

$$= 0,410 + 0,531 = \boxed{0,940 \text{ bits}}$$

Étape 2 : Gain pour l'attribut « Ciel » (Outlook)

L'attribut Ciel a 3 valeurs : Sunny, Overcast, Rain.

— **Sunny** : jours J1, J2, J8, J9, J11 \Rightarrow 5 jours dont 2 Oui, 3 Non

$$H(\text{Sunny}) = -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right) = 0,971$$

— **Overcast** : jours J3, J7, J12, J13 \Rightarrow 4 jours dont 4 Oui, 0 Non

$$H(\text{Overcast}) = -1 \cdot \log_2(1) - 0 \cdot \log_2(0) = 0$$

— **Rain** : jours J4, J5, J6, J10, J14 \Rightarrow 5 jours dont 3 Oui, 2 Non

$$H(\text{Rain}) = -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right) = 0,971$$

Entropie résiduelle (moyenne pondérée) :

$$H_{\text{résid}} = \frac{5}{14} \times 0,971 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0,971 = 0,347 + 0 + 0,347 = 0,694$$

$$\boxed{\text{Gain}(\text{Play}, \text{Ciel}) = 0,940 - 0,694 = 0,247}$$

Étape 3 : Gain pour l'attribut « Vent » (Wind)

L'attribut Vent a 2 valeurs : Weak, Strong.

— **Weak** : jours J1, J3, J4, J5, J8, J9, J10, J13 \Rightarrow 8 jours dont 6 Oui, 2 Non

$$H(\text{Weak}) = -\frac{6}{8} \log_2\left(\frac{6}{8}\right) - \frac{2}{8} \log_2\left(\frac{2}{8}\right) = 0,811$$

— **Strong** : jours J2, J6, J7, J11, J12, J14 \Rightarrow 6 jours dont 3 Oui, 3 Non

$$H(\text{Strong}) = -\frac{3}{6} \log_2\left(\frac{3}{6}\right) - \frac{3}{6} \log_2\left(\frac{3}{6}\right) = 1,0$$

$$H_{\text{résid}} = \frac{8}{14} \times 0,811 + \frac{6}{14} \times 1,0 = 0,464 + 0,429 = 0,892$$

$$\boxed{\text{Gain}(\text{Play}, \text{Vent}) = 0,940 - 0,892 = 0,048}$$

Étape 4 : Comparaison des gains

En calculant de la même manière pour tous les attributs :

Attribut	Gain d'information
Ciel (Outlook)	0,247
Humidité (Humidity)	0,152
Vent (Wind)	0,048
Température (Temperature)	0,029

TABLE 12.2 – Gains d'information pour chaque attribut du dataset Tennis.

Conclusion

L'attribut **Ciel** a le gain d'information le plus élevé (0,247). C'est donc la question qui réduit le plus l'incertitude \Rightarrow c'est l'attribut choisi comme **racine** de l'arbre ! C'est exactement ce qu'on observe dans l'arbre de la Figure 12.1.

12.3.3 L'indice de Gini : une alternative à l'entropie

L'entropie n'est pas le seul moyen de mesurer l'impureté d'un ensemble. L'**indice de Gini** est une alternative plus simple à calculer, utilisée par l'algorithme **CART** (Classification And Regression Trees), qui est celui utilisé par Scikit-learn.

Indice de Gini

L'**indice de Gini** d'un ensemble S est :

$$\text{Gini}(S) = 1 - \sum_{i=1}^c p_i^2$$

où p_i est la proportion des exemples de la classe i .

Pour un problème binaire avec proportion p :

$$\text{Gini}(p) = 1 - p^2 - (1 - p)^2 = 2p(1 - p)$$

Comparaison Entropie vs Gini

Pour notre dataset tennis (9 Oui, 5 Non) :

— **Entropie** : $H = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0,940$

— **Gini** : $G = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2 = 1 - 0,413 - 0,128 = 0,459$

Les deux métriques atteignent leur maximum quand la distribution est uniforme et valent 0 quand tous les exemples sont de la même classe. En pratique, elles donnent des arbres très similaires.

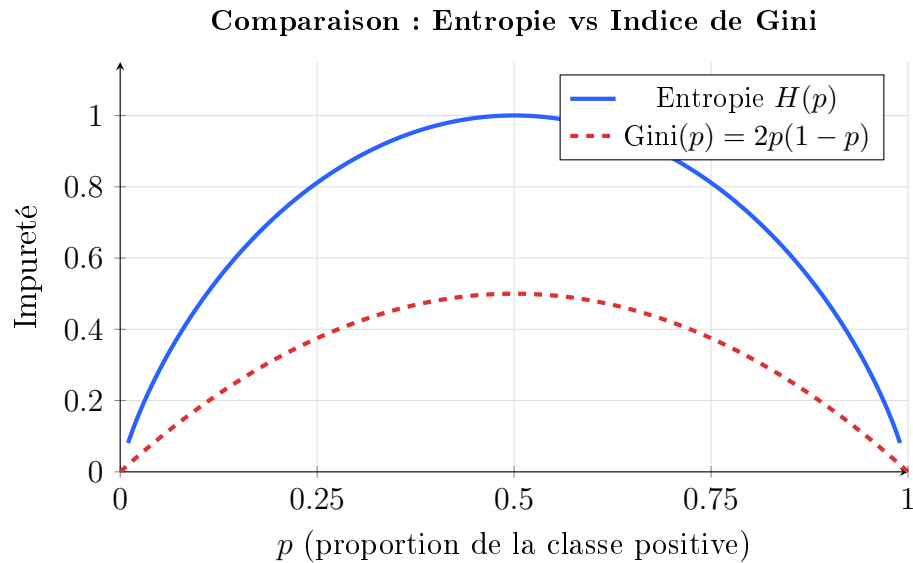


FIGURE 12.3 – Entropie et Gini ont des formes très similaires. Les deux sont maximales à $p = 0,5$ et nulles à $p = 0$ ou $p = 1$.

12.3.4 Construction de l'arbre : l'algorithme ID3/CART

Maintenant que nous savons choisir le meilleur attribut, voici l'algorithme complet pour construire un arbre de décision.

Algorithme de construction d'un arbre de décision

Entrée : un ensemble de données S et une liste d'attributs disponibles.

Sortie : un arbre de décision.

1. Cas d'arrêt :

- Si tous les exemples de S sont de la même classe → créer une feuille avec cette classe
- S'il ne reste plus d'attributs → créer une feuille avec la classe majoritaire
- Si S est vide → créer une feuille avec la classe majoritaire du parent

2. Choisir le meilleur attribut A^* : celui qui a le plus grand gain d'information (ou la plus grande réduction de Gini)

3. Créer un nœud avec l'attribut A^*

4. Pour chaque valeur v de A^* :

- Extraire le sous-ensemble S_v (exemples où $A^* = v$)
- Appeler récursivement l'algorithme sur S_v avec les attributs restants (sans A^*)
- Attacher le sous-arbre résultant comme branche

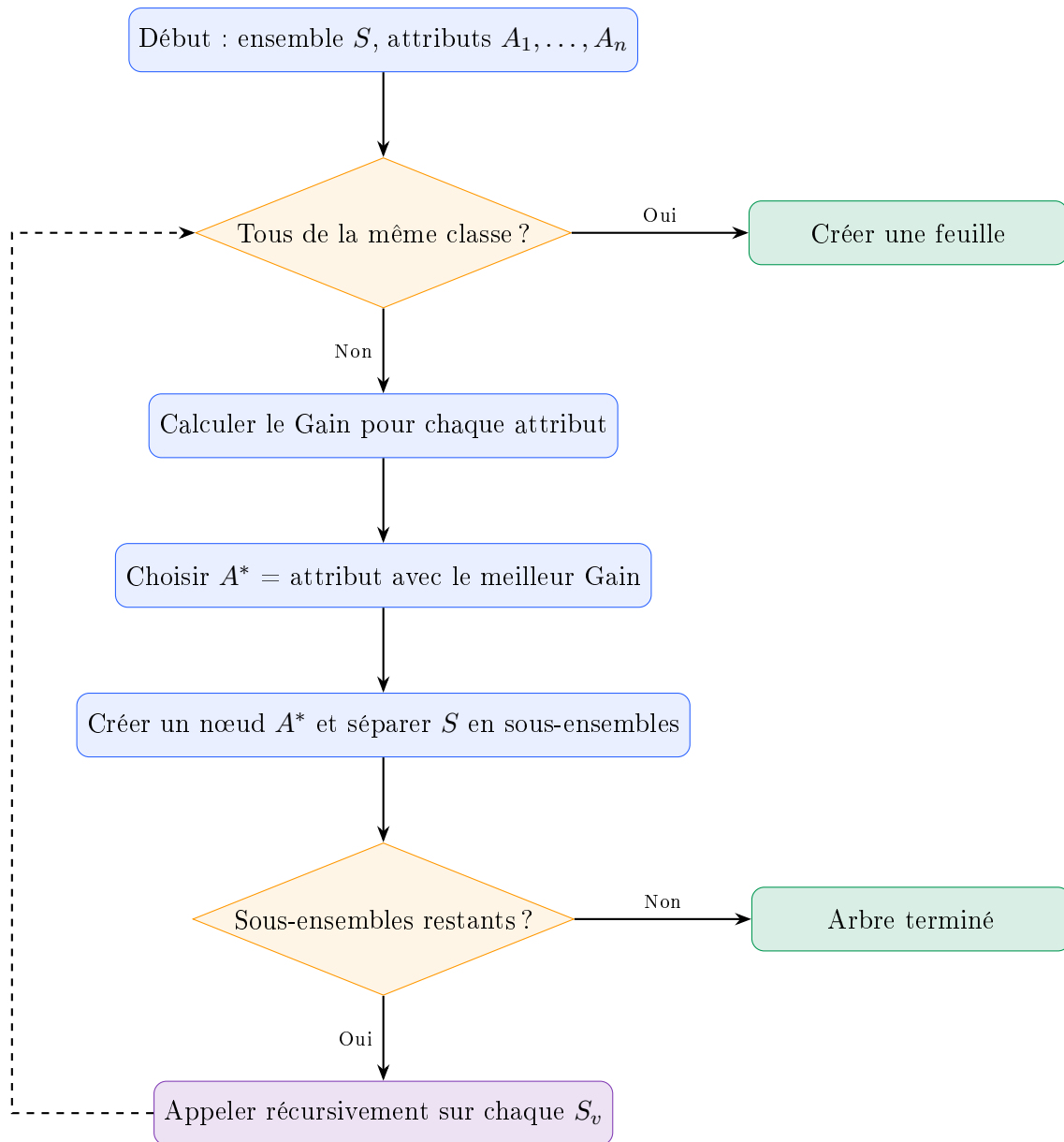


FIGURE 12.4 – Organigramme de l’algorithme de construction d’un arbre de décision.

12.3.5 L’élagage : éviter le surapprentissage

Un arbre de décision sans contrainte peut devenir **très profond**, créant une feuille pour chaque exemple d’entraînement. Il obtiendra 100% de précision sur les données d’entraînement mais échouera lamentablement sur de nouvelles données. C’est le **surapprentissage** (overfitting).

Surapprentissage des arbres de décision

Un arbre non élagué **mémorise** les données d’entraînement au lieu d’en apprendre les règles générales. C’est comme un étudiant qui apprend les réponses de l’examen de l’année dernière par cœur au lieu de comprendre le cours : il aura 20/20 sur l’ancien examen mais

échouera sur le nouveau.

Techniques d'élagage (pruning)

Pré-élagage (pré-pruning) — on arrête la croissance de l'arbre :

- `max_depth` : profondeur maximale de l'arbre (ex : 3 niveaux)
- `min_samples_split` : nombre minimum d'exemples pour créer un nœud (ex : au moins 5)
- `min_samples_leaf` : nombre minimum d'exemples dans chaque feuille (ex : au moins 2)
- `max_leaf_nodes` : nombre maximum de feuilles

Post-élagage (post-pruning) — on fait croître l'arbre complet, puis on coupe les branches inutiles :

- On évalue si chaque branche améliore vraiment la précision sur un ensemble de validation
- Si non, on remplace le sous-arbre par une feuille (classe majoritaire)
- Le paramètre `ccp_alpha` de Scikit-learn contrôle l'intensité de l'élagage

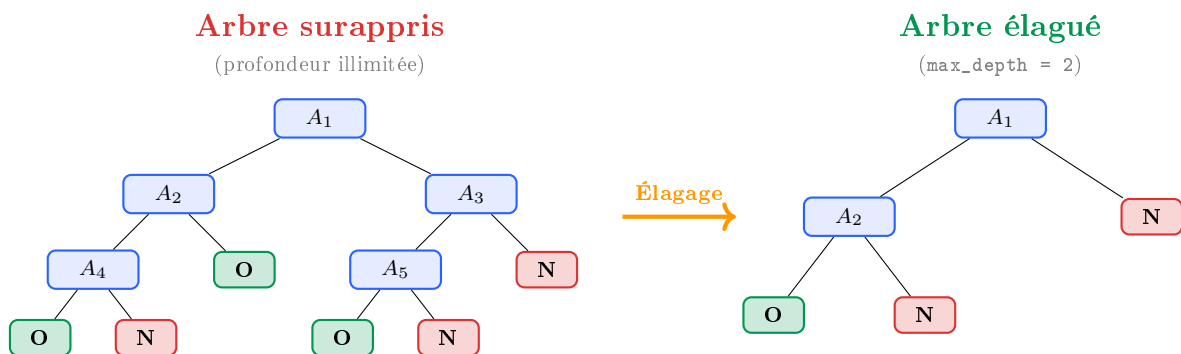


FIGURE 12.5 – Un arbre surappris (gauche) est trop complexe et mémorise le bruit. Un arbre élagué (droite) généralise mieux.

12.3.6 Arbres de régression

Les arbres de décision ne servent pas uniquement à la classification. On peut aussi les utiliser pour la **régression** (prédire un nombre).

Arbre de régression

Un **arbre de régression** fonctionne comme un arbre de classification, mais avec deux différences :

1. **Prédiction** : au lieu d'un vote majoritaire, chaque feuille prédit la **moyenne** des valeurs y des exemples qu'elle contient.
2. **Critère de séparation** : au lieu de l'entropie ou du Gini, on minimise la **variance** (ou l'erreur quadratique moyenne) :

$$\text{Variance}(S) = \frac{1}{|S|} \sum_{i=1}^{|S|} (y_i - \bar{y})^2 \quad \text{où } \bar{y} = \text{moyenne des } y_i$$

On choisit la séparation qui réduit le plus la variance totale (somme pondérée des variances des sous-ensembles).

Intuition de l'arbre de régression

Imaginons qu'on prédit la dépense d'un client :

- Si Âge $\leq 30 \rightarrow$ dépense moyenne = 150 €
- Si Âge > 30 et Revenu $\leq 3000 \rightarrow$ dépense moyenne = 280 €
- Si Âge > 30 et Revenu $> 3000 \rightarrow$ dépense moyenne = 520 €

L'arbre découpe l'espace des features en **régions rectangulaires** et prédit une constante (la moyenne) dans chaque région.

12.4 Application sur le dataset clientèle

Appliquons maintenant les arbres de décision à notre fil rouge : prédire si un client va partir (**Churn**) en fonction de ses caractéristiques.

Construction de l'arbre pour le Churn

Rappelons notre dataset clientèle avec les variables : Âge, Revenu, Dépense mensuelle, Ancienneté, Satisfaction, Churn (Oui/Non).

L'algorithme procède ainsi :

1. **Calcul du gain pour chaque variable** :
 - Gain(Satisfaction) = 0,31 \leftarrow le plus élevé
 - Gain(Ancienneté) = 0,22
 - Gain(Âge) = 0,10
 - Gain(Revenu) = 0,08
 - Gain(Dépense) = 0,05
2. **Racine** : Satisfaction (le plus informatif)
3. **Récursion** : on répète sur chaque branche avec les variables restantes

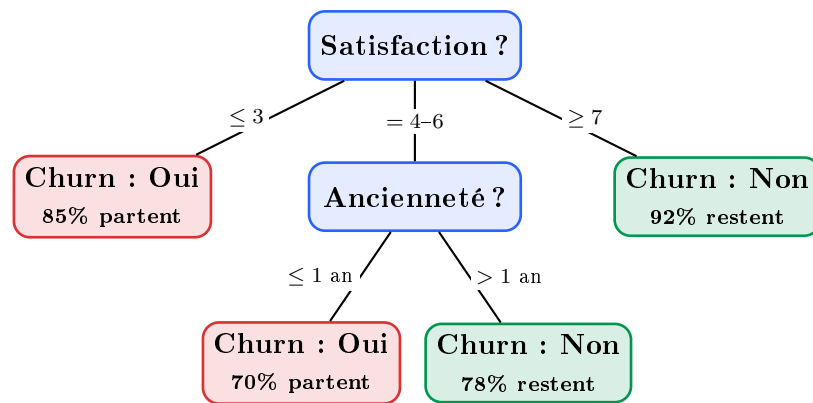


FIGURE 12.6 – Arbre de décision pour la prédiction du Churn. La satisfaction du client est le facteur le plus déterminant.

Interprétation de l'arbre Churn

L'arbre nous révèle une règle métier très claire :

- Les clients **insatisfaits** (Satisfaction ≤ 3) partent dans 85% des cas, indépendamment des autres facteurs.
- Les clients **très satisfaits** (Satisfaction ≥ 7) restent dans 92% des cas.
- Pour les clients **moyennement satisfaits**, c'est l'**ancienneté** qui fait la différence : les nouveaux clients (moins d'un an) sont plus fragiles.

C'est l'un des grands avantages des arbres de décision : ils sont **interprétables** ! On peut expliquer chaque prédiction à un non-spécialiste.

12.5 Implémentation sur Google Colab

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.tree import DecisionTreeClassifier,
   DecisionTreeRegressor
5 from sklearn.tree import plot_tree, export_text
6 from sklearn.model_selection import train_test_split
7 from sklearn.metrics import accuracy_score, classification_report
8
9 # =====
10 # 1. Créer le dataset clientele (fil rouge du cours)
11 # =====
12 np.random.seed(42)
13 n = 200

```

```

14
15 age = np.random.randint(18, 70, n)
16 revenu = np.random.randint(1500, 8000, n)
17 anciennete = np.random.randint(0, 10, n)
18 satisfaction = np.random.randint(1, 11, n)
19 depense = np.random.randint(50, 800, n)
20
21 # Le churn depend de la satisfaction et de l'anciennete
22 prob_churn = 1 / (1 + np.exp(0.5 * satisfaction + 0.3 * anciennete
    - 4))
23 churn = (np.random.rand(n) < prob_churn).astype(int)
24
25 df = pd.DataFrame({
26     'Age': age, 'Revenu': revenu,
27     'Anciennete': anciennete, 'Satisfaction': satisfaction,
28     'Depense': depense, 'Churn': churn
29 })
30 print(df.head(10))
31 print(f"\nRepartition Churn :
    {df['Churn'].value_counts().to_dict()}")

```

Listing 12.1 – Imports et préparation des données

```

1  # =====
2  # 2. Separation train / test
3  # =====
4  X = df[['Age', 'Revenu', 'Anciennete', 'Satisfaction', 'Depense']]
5  y = df['Churn']
6  X_train, X_test, y_train, y_test = train_test_split(
7      X, y, test_size=0.2, random_state=42
8  )
9
10 # =====
11 # 3. Entraîner un arbre de decision (sans limite de profondeur)
12 # =====
13 tree_full = DecisionTreeClassifier(random_state=42)
14 tree_full.fit(X_train, y_train)
15 y_pred_full = tree_full.predict(X_test)
16 print("Precision (arbre complet) :", accuracy_score(y_test,
    y_pred_full))
17
18 # =====

```

```

19 # 4. Entraîner avec max_depth = 3 (arbre elague)
20 # =====
21 tree_pruned = DecisionTreeClassifier(max_depth=3, random_state=42)
22 tree_pruned.fit(X_train, y_train)
23 y_pred_pruned = tree_pruned.predict(X_test)
24 print("Precision (max_depth=3) :", accuracy_score(y_test,
    y_pred_pruned))

```

Listing 12.2 – Entraînement de l'arbre de classification

```

1 # =====
2 # 5. Visualiser l'arbre de decision
3 # =====
4 plt.figure(figsize=(20, 10))
5 plot_tree(
6     tree_pruned,
7     feature_names=['Age', 'Revenu', 'Anciennete', 'Satisfaction',
8         'Depense'],
9     class_names=['Fidele', 'Churn'],
10    filled=True,
11    rounded=True,
12    fontsize=10,
13    proportion=True
14 )
15 plt.title("Arbre de decision pour la prediction du Churn
16     (max_depth=3)",
17     fontsize=14)
18 plt.tight_layout()
19 plt.show()
20 # =====
21 # 6. Exporter l'arbre en texte
22 # =====
23 regles = export_text(
24     tree_pruned,
25     feature_names=['Age', 'Revenu', 'Anciennete', 'Satisfaction',
26         'Depense']
27 )
28 print("Regles de l'arbre :")
29 print(regles)

```

Listing 12.3 – Visualisation de l'arbre

```

1  # =====
2  # 7. Importance des variables (feature importance)
3  # =====
4  importances = tree_pruned.feature_importances_
5  features = ['Age', 'Revenu', 'Anciennete', 'Satisfaction',
6              'Depense']
7
8  plt.figure(figsize=(8, 5))
9  plt.barh(features, importances, color='steelblue')
10 plt.xlabel('Importance')
11 plt.title('Importance des variables dans l\'arbre de decision')
12 plt.tight_layout()
13 plt.show()
14
15 # =====
16 # 8. Comparer différentes profondeurs
17 # =====
18 depths = [2, 3, 5, None] # None = pas de limite
19 results_train = []
20 results_test = []
21
22 for d in depths:
23     tree = DecisionTreeClassifier(max_depth=d, random_state=42)
24     tree.fit(X_train, y_train)
25     results_train.append(accuracy_score(y_train,
26                                         tree.predict(X_train)))
27     results_test.append(accuracy_score(y_test,
28                                       tree.predict(X_test)))
29
30 plt.figure(figsize=(8, 5))
31 x_pos = range(len(depths))
32 labels = [str(d) if d else 'Aucune' for d in depths]
33 plt.bar([p - 0.15 for p in x_pos], results_train, 0.3,
34         label='Train', color='steelblue')
35 plt.bar([p + 0.15 for p in x_pos], results_test, 0.3,
36         label='Test', color='coral')
37 plt.xticks(x_pos, labels)
38 plt.xlabel('Profondeur maximale (max_depth)')
39 plt.ylabel('Precision (accuracy)')
40 plt.title('Effet de la profondeur sur les performances')
41 plt.legend()

```

```

39 plt.ylim(0.5, 1.05)
40 plt.grid(axis='y', alpha=0.3)
41 plt.tight_layout()
42 plt.show()

```

Listing 12.4 – Importance des variables et comparaison de profondeurs

```

1  # =====
2  # 9. Frontiere de decision (avec 2 variables pour visualiser)
3  # =====
4  X_2d = df[['Satisfaction', 'Anciennete']]
5  tree_2d = DecisionTreeClassifier(max_depth=3, random_state=42)
6  tree_2d.fit(X_2d, y)
7
8  # Creer une grille de points
9  x_min, x_max = X_2d['Satisfaction'].min()-1,
   X_2d['Satisfaction'].max()+1
10 y_min, y_max = X_2d['Anciennete'].min()-1,
   X_2d['Anciennete'].max()+1
11 xx, yy = np.meshgrid(
12     np.arange(x_min, x_max, 0.1),
13     np.arange(y_min, y_max, 0.1)
14 )
15 Z = tree_2d.predict(np.c_[xx.ravel(), yy.ravel()])
16 Z = Z.reshape(xx.shape)
17
18 plt.figure(figsize=(8, 6))
19 plt.contourf(xx, yy, Z, alpha=0.3, cmap='RdYlGn')
20 scatter = plt.scatter(X_2d['Satisfaction'], X_2d['Anciennete'],
21                       c=y, cmap='RdYlGn', edgecolors='black', s=40)
22 plt.xlabel('Satisfaction')
23 plt.ylabel('Anciennete (annees)')
24 plt.title('Frontiere de decision de 1\'arbre (max_depth=3)')
25 plt.colorbar(scatter, label='0=Fidele, 1=Churn')
26 plt.tight_layout()
27 plt.show()

```

Listing 12.5 – Frontiere de decision (2 variables)

```

1  # =====
2  # 10. Arbre de REGRESSION : predire la Depense
3  # =====

```

```

4 X_reg = df[['Age', 'Revenu', 'Anciennete', 'Satisfaction']]
5 y_reg = df['Depense']
6 X_train_r, X_test_r, y_train_r, y_test_r = train_test_split(
7     X_reg, y_reg, test_size=0.2, random_state=42
8 )
9
10 reg_tree = DecisionTreeRegressor(max_depth=4, random_state=42)
11 reg_tree.fit(X_train_r, y_train_r)
12 y_pred_r = reg_tree.predict(X_test_r)
13
14 from sklearn.metrics import mean_squared_error, r2_score
15 rmse = np.sqrt(mean_squared_error(y_test_r, y_pred_r))
16 r2 = r2_score(y_test_r, y_pred_r)
17 print(f"RMSE : {rmse:.2f}")
18 print(f"R2    : {r2:.3f}")
19
20 # Visualiser l'arbre de regression
21 plt.figure(figsize=(18, 8))
22 plot_tree(reg_tree, feature_names=['Age', 'Revenu', 'Anciennete',
23     'Satisfaction'],
24     filled=True, rounded=True, fontsize=8)
25 plt.title("Arbre de regression pour la prediction de la Depense
26     (max_depth=4)")
27 plt.tight_layout()
28 plt.show()

```

Listing 12.6 – Arbre de régression pour prédire la dépense

12.6 Schéma récapitulatif : phases d'entraînement et de test

PHASE D'ENTRAÎNEMENT (Training)

1. Étape 1 — Charger et séparer les données

Charger le dataset et le séparer en 80% Train, 20% Test.

Code Python :

```
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2)
```

2. Étape 2 — Pas besoin de normaliser

Les arbres de décision sont **invariants à l'échelle** des variables.

Mathématiques :

Les comparaisons $x_j \leq s$ ne dépendent pas de l'unité

Code Python :

Pas de StandardScaler nécessaire !

3. Étape 3 — Calculer l'entropie ou Gini

Pour chaque nœud, calculer l'entropie (ou l'indice de Gini) de toutes les features disponibles.

Mathématiques (Entropie) :

$$H(S) = - \sum_{k=1}^K p_k \log_2(p_k)$$

Mathématiques (Gini) :

$$G(S) = 1 - \sum_{k=1}^K p_k^2$$

Code Python :

```
from sklearn.tree import DecisionTreeClassifier
tree = DecisionTreeClassifier(criterion='entropy', max_depth=5)
```

4. Étape 4 — Trouver la meilleure séparation à chaque nœud

Sélectionner la feature et le seuil qui **maximisent le gain d'information** (ou la réduction de Gini).

Mathématiques :

$$\text{Gain}(S, A) = H(S) - \sum_v \frac{|S_v|}{|S|} H(S_v)$$

Code Python :

```
tree.fit(X_train, y_train) # l'algo trouve les meilleures separations
```

5. Étape 5 — Construire l'arbre récursivement jusqu'au critère d'arrêt

Continuer la construction jusqu'à un critère d'arrêt (profondeur maximale, nœud pur, nombre minimum d'observations).

Mathématiques :

Arrêter quand : profondeur max atteinte, ou nœud pur ($H = 0$), ou `min_samples_leaf`.

Code Python :

```
tree = DecisionTreeClassifier(max_depth=5, min_samples_leaf=5)
```



Modèle entraîné transmis à la phase de test

PHASE DE TEST (Testing)

1. Étape 1 — Descendre dans l'arbre

Pour chaque point de Test, descendre dans l'arbre en suivant les règles de décision apprises.

Mathématiques :

À chaque nœud, tester $x_j \leq s$: si vrai \rightarrow gauche, sinon \rightarrow droite.

Code Python :

```
y_pred = tree.predict(X_test)
```

2. Étape 2 — Prédire la classe de la feuille

Une fois arrivé à une feuille, la classe prédite est la **classe majoritaire** de cette feuille.

Mathématiques :

$$\hat{y} = \arg \max_k \{p_k \text{ dans la feuille}\}$$

Code Python :

```
y_proba = tree.predict_proba(X_test) # probabilités par classe
```

3. Étape 3 — Évaluer

Calculer les métriques de performance — matrice de confusion, précision, rappel, F1-score.

Code Python :

```
print(classification_report(y_test, y_pred))
from sklearn.tree import plot_tree
plot_tree(tree, filled=True) # visualiser l'arbre
```

12.7 Exercices

Exercice 12.1 — Calculer l'entropie à la main

Un ensemble S contient **5 exemples positifs** et **3 exemples négatifs** (8 exemples au total).

1. Calculez les probabilités p_+ et p_- .
2. Calculez l'entropie $H(S)$ en utilisant la formule $H = -p_+ \log_2(p_+) - p_- \log_2(p_-)$.

Indication : $\log_2(5) \approx 2,322$, $\log_2(3) \approx 1,585$, $\log_2(8) = 3$.

Rappel : $\log_2\left(\frac{a}{b}\right) = \log_2(a) - \log_2(b)$.

3. L'incertitude est-elle élevée ou faible ? Justifiez.

Correction de l'exercice 12.1

1. Probabilités :

$$p_+ = \frac{5}{8} = 0,625, \quad p_- = \frac{3}{8} = 0,375$$

Vérification : $0,625 + 0,375 = 1 \checkmark$

2. Entropie :

$$H(S) = -\frac{5}{8} \log_2\left(\frac{5}{8}\right) - \frac{3}{8} \log_2\left(\frac{3}{8}\right)$$

Calculons chaque logarithme :

$$\log_2\left(\frac{5}{8}\right) = \log_2(5) - \log_2(8) = 2,322 - 3 = -0,678$$

$$\log_2\left(\frac{3}{8}\right) = \log_2(3) - \log_2(8) = 1,585 - 3 = -1,415$$

Donc :

$$\begin{aligned} H(S) &= -0,625 \times (-0,678) - 0,375 \times (-1,415) \\ &= 0,424 + 0,531 = \boxed{0,954 \text{ bits}} \end{aligned}$$

3. Interprétation : L'entropie maximale pour 2 classes est $H_{\max} = 1$ bit (quand $p = 0,5$). Ici $H = 0,954$, ce qui est **assez élevé**. L'incertitude est forte car la répartition 5/3 est relativement équilibrée (pas très éloignée de 4/4).

Exercice 12.2 — Calculer un gain d'information

Un ensemble S contient 10 exemples : 6 positifs (+) et 4 négatifs (-).

On considère un attribut A qui a deux valeurs possibles (A_1 et A_2) :

- A_1 : 4 exemples dont 3 positifs et 1 négatif
- A_2 : 6 exemples dont 3 positifs et 3 négatifs

1. Calculez l'entropie $H(S)$ de l'ensemble initial.
2. Calculez l'entropie $H(S_{A_1})$ du sous-ensemble où $A = A_1$.
3. Calculez l'entropie $H(S_{A_2})$ du sous-ensemble où $A = A_2$.
4. Calculez l'entropie résiduelle (moyenne pondérée).
5. Calculez le gain d'information $\text{Gain}(S, A)$.
6. Ce gain est-il bon ou mauvais ? Justifiez.

Indication : $\log_2(3) \approx 1,585$, $\log_2(5) \approx 2,322$.

Correction de l'exercice 12.2**1. Entropie de S :**

$$p_+ = \frac{6}{10} = 0,6, \quad p_- = \frac{4}{10} = 0,4$$

$$H(S) = -0,6 \cdot \log_2(0,6) - 0,4 \cdot \log_2(0,4)$$

$$\log_2(0,6) = \log_2(6) - \log_2(10) = (\log_2(2) + \log_2(3)) - (\log_2(2) + \log_2(5))$$

$$= (1 + 1,585) - (1 + 2,322) = 2,585 - 3,322 = -0,737$$

$$\log_2(0,4) = \log_2(4) - \log_2(10) = 2 - 3,322 = -1,322$$

$$H(S) = -0,6 \times (-0,737) - 0,4 \times (-1,322) = 0,442 + 0,529 = \boxed{0,971}$$

2. Entropie de S_{A_1} : (3 positifs, 1 négatif, total = 4)

$$H(S_{A_1}) = -\frac{3}{4} \log_2\left(\frac{3}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right)$$

$$\log_2(0,75) = \log_2(3) - \log_2(4) = 1,585 - 2 = -0,415$$

$$\log_2(0,25) = -2$$

$$H(S_{A_1}) = -0,75 \times (-0,415) - 0,25 \times (-2) = 0,311 + 0,5 = \boxed{0,811}$$

3. Entropie de S_{A_2} : (3 positifs, 3 négatifs, total = 6)

$$H(S_{A_2}) = -\frac{3}{6} \log_2\left(\frac{3}{6}\right) - \frac{3}{6} \log_2\left(\frac{3}{6}\right) = -\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right)$$

$$= -0,5 \times (-1) - 0,5 \times (-1) = \boxed{1,0}$$

4. Entropie résiduelle :

$$H_{\text{résid}} = \frac{|S_{A_1}|}{|S|} \cdot H(S_{A_1}) + \frac{|S_{A_2}|}{|S|} \cdot H(S_{A_2}) = \frac{4}{10} \times 0,811 + \frac{6}{10} \times 1,0$$

$$= 0,324 + 0,6 = \boxed{0,924}$$

5. Gain d'information :

$$\text{Gain}(S, A) = H(S) - H_{\text{résid}} = 0,971 - 0,924 = \boxed{0,047}$$

6. Interprétation : Le gain est très faible (0,047 bit). L'attribut A ne réduit presque pas l'incertitude. Ce n'est **pas un bon attribut** pour commencer l'arbre. Cela se comprend intuitivement : le sous-ensemble A_2 est 50/50 (aucune information!) et même A_1 reste assez mélangé (3/1).

Exercice 12.3 — Choisir la racine de l'arbre

Voici un dataset de 12 étudiants. On veut prédire s'ils réussissent l'examen (Réussite : Oui/Non).

#	Assiduité	Exercices	Sommeil	Réussite
1	Haute	Faits	Bon	Oui
2	Haute	Faits	Mauvais	Oui
3	Haute	Non	Bon	Oui
4	Haute	Non	Mauvais	Non
5	Basse	Faits	Bon	Oui
6	Basse	Faits	Mauvais	Non
7	Basse	Non	Bon	Non
8	Basse	Non	Mauvais	Non
9	Haute	Faits	Bon	Oui
10	Basse	Non	Mauvais	Non
11	Haute	Non	Bon	Oui
12	Basse	Faits	Bon	Oui

1. Calculez l'entropie globale $H(\text{Réussite})$.
2. Calculez le gain d'information pour **chaque** attribut (Assiduité, Exercices, Sommeil).
3. Quel attribut devrait être la racine de l'arbre ?

Correction de l'exercice 12.3

Comptage global : 7 Oui et 5 Non (total = 12).

1. Entropie globale :

$$H = -\frac{7}{12} \log_2\left(\frac{7}{12}\right) - \frac{5}{12} \log_2\left(\frac{5}{12}\right)$$

$$\frac{7}{12} \approx 0,583, \quad \log_2(0,583) \approx -0,778$$

$$\frac{5}{12} \approx 0,417, \quad \log_2(0,417) \approx -1,263$$

$$H = -0,583 \times (-0,778) - 0,417 \times (-1,263) = 0,454 + 0,527 = \boxed{0,980}$$

2a. Gain pour Assiduité :

— **Haute** : étudiants 1, 2, 3, 4, 9, 11 \Rightarrow 6 dont 5 Oui, 1 Non

$$H(\text{Haute}) = -\frac{5}{6} \log_2\left(\frac{5}{6}\right) - \frac{1}{6} \log_2\left(\frac{1}{6}\right) \approx 0,650$$

— **Basse** : étudiants 5, 6, 7, 8, 10, 12 \Rightarrow 6 dont 2 Oui, 4 Non

$$H(\text{Basse}) = -\frac{2}{6} \log_2\left(\frac{2}{6}\right) - \frac{4}{6} \log_2\left(\frac{4}{6}\right) \approx 0,918$$

$$H_{\text{résid}} = \frac{6}{12} \times 0,650 + \frac{6}{12} \times 0,918 = 0,325 + 0,459 = 0,784$$

$$\text{Gain}(\text{Assiduité}) = 0,980 - 0,784 = \boxed{0,196}$$

2b. Gain pour Exercices :

— **Faits** : étudiants 1, 2, 5, 6, 9, 12 \Rightarrow 6 dont 5 Oui, 1 Non

$$H(\text{Faits}) = 0,650$$

— **Non** : étudiants 3, 4, 7, 8, 10, 11 \Rightarrow 6 dont 2 Oui, 4 Non

$$H(\text{Non}) = 0,918$$

$$H_{\text{résid}} = \frac{6}{12} \times 0,650 + \frac{6}{12} \times 0,918 = 0,784$$

$$\text{Gain}(\text{Exercices}) = 0,980 - 0,784 = \boxed{0,196}$$

2c. Gain pour Sommeil :

— **Bon** : étudiants 1, 3, 5, 7, 9, 11, 12 \Rightarrow 7 dont 6 Oui, 1 Non

$$H(\text{Bon}) = -\frac{6}{7} \log_2\left(\frac{6}{7}\right) - \frac{1}{7} \log_2\left(\frac{1}{7}\right) \approx 0,592$$

— **Mauvais** : étudiants 2, 4, 6, 8, 10 \Rightarrow 5 dont 1 Oui, 4 Non

$$H(\text{Mauvais}) = -\frac{1}{5} \log_2\left(\frac{1}{5}\right) - \frac{4}{5} \log_2\left(\frac{4}{5}\right) \approx 0,722$$

$$H_{\text{résid}} = \frac{7}{12} \times 0,592 + \frac{5}{12} \times 0,722 = 0,345 + 0,301 = 0,646$$

$$\text{Gain}(\text{Sommeil}) = 0,980 - 0,646 = \boxed{0,334}$$

3. Comparaison :

Attribut	Gain
Sommeil	0,334 ← le meilleur
Assiduité	0,196
Exercices	0,196

Conclusion : L'attribut **Sommeil** a le gain d'information le plus élevé. C'est donc lui qui devrait être choisi comme **racine** de l'arbre. Cela a du sens intuitivement : un bon sommeil la veille de l'examen est le facteur le plus déterminant pour la réussite dans ce dataset.

Exercice 12.4 — Pratique Python : arbre sur le dataset Iris

En utilisant Google Colab, réalisez les étapes suivantes :

1. Chargez le célèbre dataset *Iris* depuis Scikit-learn (`from sklearn.datasets import load_iris`).
2. Séparez en train (80%) et test (20%).
3. Entraînez un `DecisionTreeClassifier` avec `max_depth=None` (arbre complet). Affichez la précision sur le test.
4. Visualisez l'arbre avec `plot_tree`. Commentez sa profondeur.
5. Entraînez des arbres avec `max_depth = 1, 2, 3, 4, 5` et `None`. Tracez un graphique montrant la précision train et test en fonction de la profondeur.
6. Quel est le meilleur `max_depth`? Y a-t-il du surapprentissage?
7. Affichez l'importance des variables (`feature_importances_`) sous forme de diagramme en barres.

Correction de l'exercice 12.4

```
# 1. Charger le dataset Iris
from sklearn.datasets import load_iris
from sklearn.tree import DecisionTreeClassifier, plot_tree
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
import matplotlib.pyplot as plt
import numpy as np

iris = load_iris()
X, y = iris.data, iris.target
print(f"Taille : {X.shape}") # (150, 4)
```

```

print(f"Classes : {iris.target_names}") # setosa, versicolor,
    virginica
print(f"Variables : {iris.feature_names}")

# 2. Separation train / test
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42
)
print(f"Train : {X_train.shape[0]}, Test : {X_test.shape[0]}")

# 3. Arbre complet (sans limite de profondeur)
tree_full = DecisionTreeClassifier(random_state=42)
tree_full.fit(X_train, y_train)
y_pred = tree_full.predict(X_test)
print(f"Precision (arbre complet) : {accuracy_score(y_test,
    y_pred):.3f}")

# 4. Visualiser l'arbre
plt.figure(figsize=(18, 10))
plot_tree(tree_full, feature_names=iris.feature_names,
    class_names=list(iris.target_names),
    filled=True, rounded=True, fontsize=9)
plt.title("Arbre de decision complet sur le dataset Iris")
plt.tight_layout()
plt.show()
# Commentaire : l'arbre a une profondeur de ~5, ce qui est
# assez profond pour seulement 120 exemples d'entraînement.

# 5. Comparer différentes profondeurs
depths = [1, 2, 3, 4, 5, None]
acc_train = []
acc_test = []

for d in depths:
    tree = DecisionTreeClassifier(max_depth=d, random_state=42)
    tree.fit(X_train, y_train)
    acc_train.append(accuracy_score(y_train,
        tree.predict(X_train)))
    acc_test.append(accuracy_score(y_test,
        tree.predict(X_test)))

```

```

plt.figure(figsize=(8, 5))
labels = [str(d) if d else 'None' for d in depths]
plt.plot(labels, acc_train, 'o-', label='Train',
         color='steelblue')
plt.plot(labels, acc_test, 's--', label='Test', color='coral')
plt.xlabel('max_depth')
plt.ylabel('Precision')
plt.title('Precision train/test selon la profondeur')
plt.legend()
plt.grid(alpha=0.3)
plt.ylim(0.8, 1.02)
plt.show()

# 6. Analyse :
# - A max_depth=2 ou 3, la precision test est deja tres bonne
# - A max_depth=None, le train est a 100% mais le test peut
#   baisser legerement -> signe de surapprentissage
# - Le meilleur compromis est souvent max_depth=3

# 7. Importance des variables
tree_best = DecisionTreeClassifier(max_depth=3, random_state=42)
tree_best.fit(X_train, y_train)
importances = tree_best.feature_importances_

plt.figure(figsize=(8, 5))
plt.barh(iris.feature_names, importances, color='steelblue')
plt.xlabel('Importance')
plt.title("Importance des variables (Iris, max_depth=3)")
plt.tight_layout()
plt.show()

# Observation : petal length et petal width sont les variables
# les plus importantes, ce qui est coherent avec la biologie.

```

Chapitre 13

Random Forest (Forêts Aléatoires)

13.1 Mise en situation concrète

Situation : Prédire le risque de maladie dans un hôpital

Un hôpital souhaite prédire si un patient présente un risque élevé de maladie cardiovasculaire à partir de ses données médicales (âge, pression artérielle, cholestérol, glycémie, etc.).

Le problème : un seul médecin peut faire des erreurs de diagnostic. Mais si on demande l'avis de **100 médecins** et qu'on prend le **vote majoritaire**, la réponse sera beaucoup plus fiable !

C'est exactement le principe de la **Random Forest** : au lieu de s'appuyer sur un seul arbre de décision, on en construit **des centaines** et on combine leurs prédictions.

Comparaison : 1 arbre vs 100 arbres

On utilise un jeu de données médicales avec 1000 patients et 10 caractéristiques.

```
1 from sklearn.tree import DecisionTreeClassifier
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.model_selection import cross_val_score
4 from sklearn.datasets import make_classification
5 import numpy as np
6
7 # Simuler des données médicales
8 X, y = make_classification(n_samples=1000, n_features=10,
9                           n_informative=6, n_redundant=2,
10                          random_state=42)
11
12 # Un seul arbre de décision
13 arbre_seul = DecisionTreeClassifier(random_state=42)
14 score_arbre = cross_val_score(arbre_seul, X, y, cv=5).mean()
15 print(f"Precision d'un seul arbre : {score_arbre:.1%}")
16 # Resultat : Precision d'un seul arbre : 78.0%
```

```

17
18 # Forêt aléatoire de 100 arbres
19 foret = RandomForestClassifier(n_estimators=100,
    random_state=42)
20 score_foret = cross_val_score(foret, X, y, cv=5).mean()
21 print(f"Précision de la forêt (100 arbres) : {score_foret:.1%}")
22 # Résultat : Précision de la forêt (100 arbres) : 92.0%

```

Résultat : la forêt aléatoire passe de **78%** (un seul arbre) à **92%** (100 arbres). Le gain est considérable, simplement en combinant plusieurs arbres !

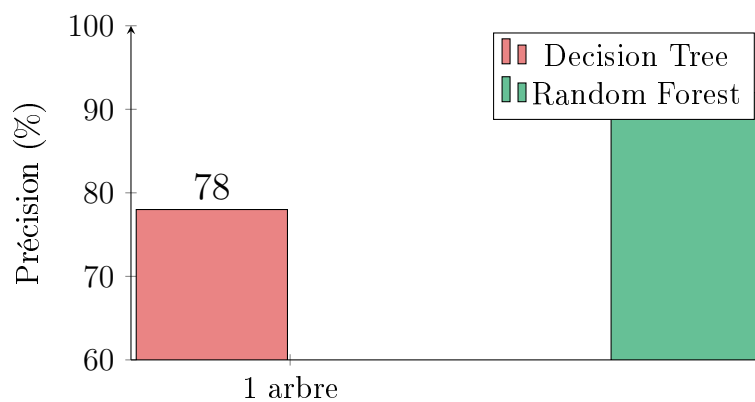


FIGURE 13.1 – Comparaison de la précision : un seul arbre vs une forêt de 100 arbres.

13.2 Intuition et idée générale

La sagesse des foules

Imagine que tu veuilles deviner le nombre de bonbons dans un bocal. Si tu demandes à **une seule personne**, elle peut se tromper complètement. Mais si tu demandes à **100 personnes** et que tu fais la **moyenne** de leurs réponses, le résultat sera très proche de la vérité.

C'est le phénomène de la **sagesse des foules** (*wisdom of crowds*). En Machine Learning, cela s'appelle les **méthodes d'ensemble** : combiner plusieurs modèles faibles pour obtenir un modèle fort.

13.2.1 Pourquoi un seul arbre ne suffit pas ?

Un arbre de décision seul présente plusieurs problèmes :

- **Surapprentissage (overfitting)** : un arbre profond mémorise les données d'entraînement au lieu d'apprendre les vrais patterns.

- **Instabilité** : un petit changement dans les données peut produire un arbre complètement différent.
- **Variance élevée** : les prédictions varient beaucoup d'un échantillon à l'autre.

13.2.2 La solution : combiner beaucoup d'arbres

Définition — Random Forest

Une **Random Forest** (forêt aléatoire) est un ensemble de T arbres de décision, chacun construit sur un échantillon aléatoire différent des données et avec un sous-ensemble aléatoire de variables à chaque nœud. La prédiction finale est obtenue par **vote majoritaire** (classification) ou par **moyenne** (régression).

Les deux idées clés de la Random Forest sont :

1. **Bagging** (Bootstrap Aggregating) : chaque arbre est entraîné sur un échantillon bootstrap différent.
2. **Sélection aléatoire de variables** : à chaque division, seul un sous-ensemble aléatoire de variables est considéré.

13.2.3 Analogie : le conseil de médecins

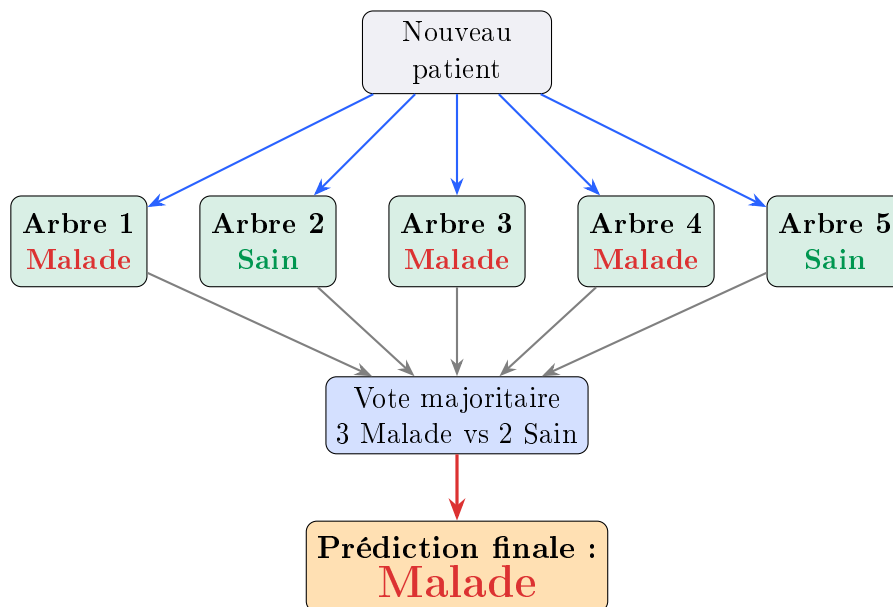


FIGURE 13.2 – Principe de la Random Forest : chaque arbre vote, et la majorité décide.

13.3 Fondements mathématiques détaillés

13.3.1 Échantillonnage bootstrap

Définition — Échantillon bootstrap

Un **échantillon bootstrap** est obtenu en tirant n observations **avec remise** à partir d'un jeu de données original de taille n . Certaines observations apparaissent plusieurs fois, d'autres n'apparaissent pas du tout.

Exemple concret de bootstrap

Soit un jeu de données original : $\mathcal{D} = \{1, 2, 3, 4, 5\}$ (5 observations).

Voici trois échantillons bootstrap possibles :

- $\mathcal{D}_1^* = \{2, 2, 4, 1, 5\}$ (l'observation 3 est absente, 2 apparaît 2 fois)
- $\mathcal{D}_2^* = \{1, 3, 3, 5, 5\}$ (les observations 2 et 4 sont absentes)
- $\mathcal{D}_3^* = \{4, 1, 2, 4, 3\}$ (l'observation 5 est absente, 4 apparaît 2 fois)

Quelle proportion des données apparaît dans chaque bootstrap ?

Pour un jeu de données de taille n , la probabilité qu'une observation donnée ne soit **jamais sélectionnée** dans un bootstrap de n tirages avec remise est :

$$P(\text{non sélectionnée}) = \left(1 - \frac{1}{n}\right)^n$$

Quand $n \rightarrow \infty$, cette probabilité converge vers :

$$\lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right)^n = \frac{1}{e} \approx 0.368$$

Résultat fondamental

En moyenne, chaque échantillon bootstrap contient environ **63.2%** des observations originales. Les **36.8%** restantes constituent les échantillons **Out-of-Bag (OOB)**.

$$P(\text{présente dans le bootstrap}) = 1 - \frac{1}{e} \approx 0.632$$

Échantillons Out-of-Bag (OOB)

Définition — Out-of-Bag (OOB)

Pour chaque arbre t , les observations qui **n'apparaissent pas** dans son échantillon bootstrap sont appelées observations **Out-of-Bag** (hors sac). Elles représentent environ 36.8% des données et servent de **jeu de validation gratuit** : on peut évaluer la performance de chaque arbre sur ses propres données OOB, sans avoir besoin d'un jeu de test séparé.

L'**erreur OOB** est calculée ainsi :

1. Pour chaque observation x_i , on identifie tous les arbres pour lesquels x_i est OOB.
2. On agrège les prédictions de ces arbres uniquement (vote majoritaire ou moyenne).
3. On compare à la vraie valeur y_i .
4. L'erreur OOB est le taux d'erreur moyen sur toutes les observations.

13.3.2 Bagging (Bootstrap Aggregating)

Définition — Bagging

Le **Bagging** (Bootstrap Aggregating), proposé par Leo Breiman en 1996, consiste à :

1. Générer T échantillons bootstrap à partir des données originales.
2. Entraîner un modèle (ici, un arbre de décision) sur chaque échantillon.
3. Agréger les prédictions de tous les modèles.

Agrégation des prédictions :

— **Classification** — vote majoritaire :

$$\hat{y}(x) = \text{mode} \{h_1(x), h_2(x), \dots, h_T(x)\}$$

où $h_t(x)$ est la prédiction de l'arbre t .

— **Régression** — moyenne :

$$\hat{y}(x) = \frac{1}{T} \sum_{t=1}^T h_t(x)$$

Pourquoi le bagging réduit-il le surapprentissage ?

Le bagging agit sur la **variance** du modèle. Rappelons la décomposition biais-variance de l'erreur :

$$\text{Erreur} = \text{Biais}^2 + \text{Variance} + \text{Bruit irréductible}$$

Si on dispose de T modèles indépendants, chacun de variance σ^2 , la variance de leur moyenne est :

$$\text{Var} \left(\frac{1}{T} \sum_{t=1}^T h_t \right) = \frac{\sigma^2}{T}$$

En pratique, les arbres ne sont pas parfaitement indépendants. Si la corrélation moyenne entre deux arbres est ρ , alors :

$$\text{Var} \left(\frac{1}{T} \sum_{t=1}^T h_t \right) = \rho \cdot \sigma^2 + \frac{1-\rho}{T} \cdot \sigma^2$$

Point clé — Réduction de la variance

Le deuxième terme $\frac{1-\rho}{T}\sigma^2$ diminue quand T augmente (plus d'arbres). Le premier terme $\rho\sigma^2$ ne dépend pas de T : c'est la **corrélation entre les arbres** qui limite le gain. C'est pourquoi la **sélection aléatoire de variables** est cruciale : elle réduit ρ .

13.3.3 Sélection aléatoire de variables (Random Subspace)

Définition — Sélection aléatoire de variables

À chaque nœud de chaque arbre, au lieu de chercher la meilleure division parmi **toutes** les p variables, on tire aléatoirement un sous-ensemble de m variables et on cherche la meilleure division **uniquement parmi celles-ci**.

Valeurs recommandées :

- Classification : $m = \lfloor \sqrt{p} \rfloor$
- Régression : $m = \lfloor p/3 \rfloor$

Pourquoi sélectionner aléatoirement les variables ?

Imaginons un jeu de données où la variable « revenu » est très prédictive. Sans sélection aléatoire, **tous les arbres** utiliseraient « revenu » comme première division. Les arbres se ressembleraient tous et le bagging n'apporterait presque rien.

En forçant chaque arbre à ne regarder qu'un sous-ensemble de variables, on obtient des arbres **différents les uns des autres** (décorrélés), ce qui maximise le bénéfice de la moyenne.

13.3.4 L'algorithme pas à pas

Algorithme — Random Forest

Entrées : données d'entraînement $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$, nombre d'arbres T , nombre de variables par nœud m .

1. **Pour** $t = 1, 2, \dots, T$:

a. **Bootstrap** : tirer un échantillon \mathcal{D}_t^* de taille n avec remise depuis \mathcal{D} .

b. **Construction de l'arbre** h_t : pour chaque nœud de l'arbre :

- Sélectionner aléatoirement m variables parmi les p disponibles.
- Trouver la meilleure division (Gini ou entropie) parmi ces m variables.
- Diviser le nœud en deux fils.

c. **Pas d'élagage** : laisser l'arbre grandir jusqu'à sa profondeur maximale (chaque feuille contient une seule classe ou très peu d'observations).

2. **Prédiction pour un nouveau point** x :

— **Classification** : $\hat{y}(x) = \text{mode}\{h_1(x), \dots, h_T(x)\}$

— **Régression** : $\hat{y}(x) = \frac{1}{T} \sum_{t=1}^T h_t(x)$

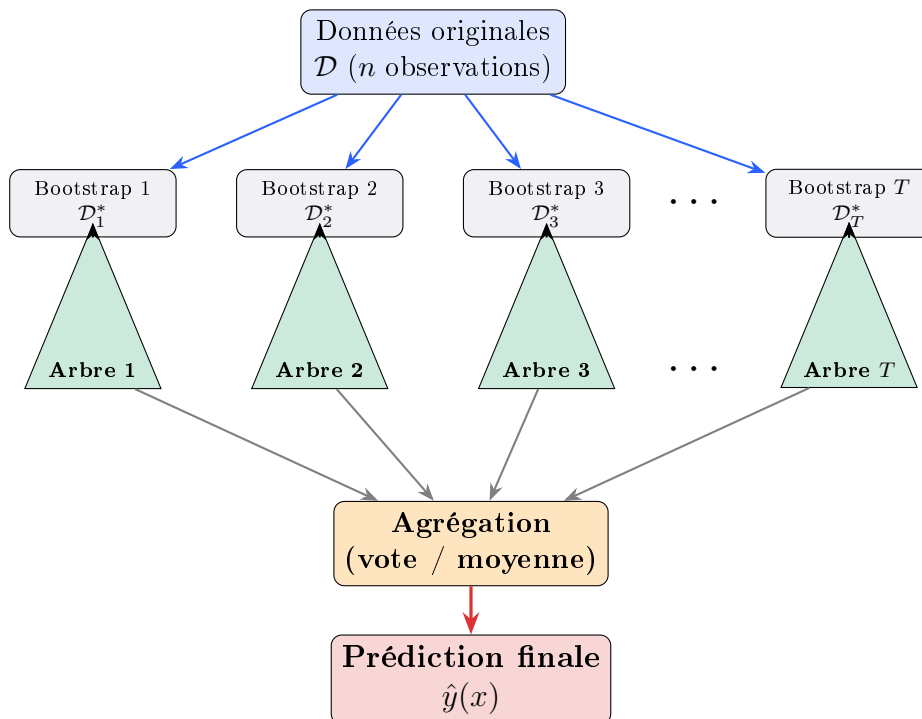


FIGURE 13.3 – Vue d'ensemble de l'algorithme Random Forest.

13.3.5 Importance des variables

La Random Forest fournit naturellement une mesure de l'**importance de chaque variable**. Il existe deux méthodes principales.

Mean Decrease in Impurity (MDI)

Pour chaque variable j , on calcule la diminution totale d'impureté (Gini ou entropie) qu'elle apporte, pondérée par la proportion d'observations concernées, puis moyennée sur tous les arbres :

$$\text{MDI}(j) = \frac{1}{T} \sum_{t=1}^T \sum_{\substack{s \in \mathcal{S}_t \\ v(s)=j}} p(s) \cdot \Delta I(s)$$

où \mathcal{S}_t est l'ensemble des nœuds de l'arbre t , $v(s)$ est la variable utilisée au nœud s , $p(s)$ est la proportion d'observations atteignant le nœud s , et $\Delta I(s)$ est la diminution d'impureté au nœud s .

Importance par permutation

Pour chaque variable j :

1. Calculer la performance du modèle sur les données OOB.
2. Permuter aléatoirement les valeurs de la variable j dans les données OOB.
3. Recalculer la performance : la **diminution de performance** mesure l'importance de j .

$$\text{PI}(j) = \text{Score}_{\text{original}} - \text{Score}_{\text{après permutation de } j}$$

MDI vs Permutation

- **MDI** est plus rapide à calculer (déjà disponible après l'entraînement) mais peut favoriser les variables à haute cardinalité.
- L'**importance par permutation** est plus fiable, surtout quand des variables sont corrélées entre elles.

13.3.6 Hyperparamètres principaux

Hyperparamètre	Description	Valeur typique
<code>n_estimators</code>	Nombre d'arbres dans la forêt	100 – 500 (plus = mieux, mais plus lent)
<code>max_features</code>	Nombre de variables par nœud	\sqrt{p} (classif.) ou $p/3$ (rég.)
<code>max_depth</code>	Profondeur maximale de chaque arbre	None (sans limite)
<code>min_samples_split</code>	Nb min. d'obs. pour diviser un nœud	2 (par défaut)
<code>min_samples_leaf</code>	Nb min. d'obs. dans une feuille	1 (par défaut)
<code>bootstrap</code>	Utiliser le bootstrap ?	True (par défaut)
<code>oob_score</code>	Calculer le score OOB ?	False (activer!)

TABLE 13.1 – Principaux hyperparamètres de la Random Forest.

Combien d'arbres faut-il ?

À mesure que T augmente, la performance s'améliore puis se stabilise. Contrairement à d'autres algorithmes, **ajouter plus d'arbres ne provoque jamais de surapprentissage**. Le seul inconvénient est le temps de calcul. En pratique, 100 à 500 arbres suffisent généralement.

13.4 Application : prédiction du Churn client

Application sur un cas réel : le Churn

Une entreprise de télécommunications souhaite prédire quels clients vont résilier leur abonnement (**churn**). On dispose de données sur 7043 clients avec des variables comme l'ancienneté, le montant mensuel, le type de contrat, etc.

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.tree import DecisionTreeClassifier
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import accuracy_score, classification_report
7

```

```

8  # Charger les donnees
9  df = pd.read_csv('churn_dataset.csv')
10
11 # Preparation des donnees
12 X = df.drop('Churn', axis=1)
13 y = df['Churn']
14 X_train, X_test, y_train, y_test = train_test_split(
15     X, y, test_size=0.2, random_state=42, stratify=y
16 )
17
18 # --- Modele 1 : Un seul arbre ---
19 arbre = DecisionTreeClassifier(random_state=42)
20 arbre.fit(X_train, y_train)
21 y_pred_arbre = arbre.predict(X_test)
22 print("=== Decision Tree ===")
23 print(f"Accuracy : {accuracy_score(y_test, y_pred_arbre):.4f}")
24
25 # --- Modele 2 : Random Forest ---
26 rf = RandomForestClassifier(
27     n_estimators=200,
28     max_features='sqrt',
29     oob_score=True,
30     random_state=42
31 )
32 rf.fit(X_train, y_train)
33 y_pred_rf = rf.predict(X_test)
34 print("\n=== Random Forest (200 arbres) ===")
35 print(f"Accuracy : {accuracy_score(y_test, y_pred_rf):.4f}")
36 print(f"Score OOB : {rf.oob_score_: .4f}")

```

Résultats comparés

Modèle	Accuracy (test)	Score OOB
Decision Tree (1 arbre)	0.7420	—
Random Forest (200 arbres)	0.8105	0.8020

La Random Forest surpasse nettement l'arbre unique. Le **score OOB** (0.8020) est très proche du score sur le jeu de test (0.8105), ce qui confirme qu'il constitue une bonne estimation de la performance générale.

13.4.1 Importance des variables pour le Churn

```

1 import matplotlib.pyplot as plt
2
3 # Importance des variables (MDI)
4 importances = rf.feature_importances_
5 indices = np.argsort(importances)[::-1]
6 noms_variables = X.columns[indices]
7
8 plt.figure(figsize=(10, 6))
9 plt.title("Importance des variables (Random Forest)")
10 plt.barh(range(10), importances[indices[:10]], color='steelblue')
11 plt.yticks(range(10), noms_variables[:10])
12 plt.xlabel("Mean Decrease in Impurity")
13 plt.gca().invert_yaxis()
14 plt.tight_layout()
15 plt.show()

```

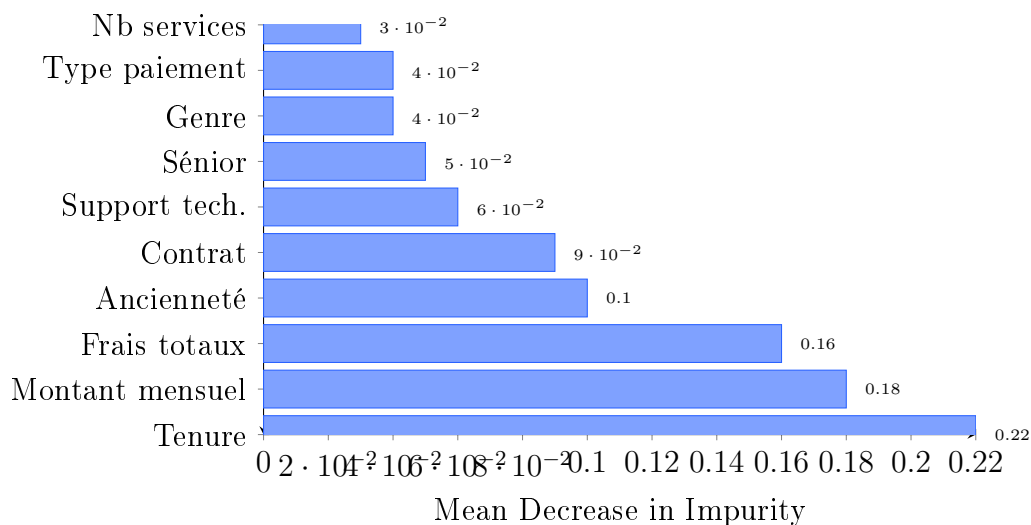


FIGURE 13.4 – Les 10 variables les plus importantes pour la prédiction du churn.

Interprétation

La variable **Tenure** (durée d'abonnement) est la plus importante : les clients récents ont plus de chances de partir. Le **montant mensuel** et les **frais totaux** suivent. Ces informations permettent à l'entreprise de cibler ses actions de fidélisation.

13.5 Implémentation complète sur Google Colab

Notebook Colab -- Random Forest complète

Voici l'implémentation complète, étape par étape.

13.5.1 Étape 1 : Imports et chargement des données

```

1  # ===== Imports =====
2  import numpy as np
3  import pandas as pd
4  import matplotlib.pyplot as plt
5  from sklearn.datasets import load_breast_cancer
6  from sklearn.ensemble import RandomForestClassifier,
   RandomForestRegressor
7  from sklearn.tree import DecisionTreeClassifier
8  from sklearn.model_selection import (
9      train_test_split, cross_val_score, learning_curve
10 )
11 from sklearn.metrics import (
12     accuracy_score, classification_report, confusion_matrix
13 )
14
15 # Charger le dataset Breast Cancer
16 data = load_breast_cancer()
17 X = pd.DataFrame(data.data, columns=data.feature_names)
18 y = data.target
19
20 print(f"Dimensions : {X.shape}")
21 print(f"Classes : {np.unique(y)} (0=malin, 1=benin)")
22 print(f"Distribution : {np.bincount(y)}")
23
24 X_train, X_test, y_train, y_test = train_test_split(
25     X, y, test_size=0.2, random_state=42, stratify=y
26 )

```

13.5.2 Étape 2 : Comparaison Decision Tree vs Random Forest

```

1  # --- Decision Tree ---
2  dt = DecisionTreeClassifier(random_state=42)

```

```

3 dt.fit(X_train, y_train)
4 y_pred_dt = dt.predict(X_test)
5 acc_dt = accuracy_score(y_test, y_pred_dt)
6
7 # --- Random Forest ---
8 rf = RandomForestClassifier(
9     n_estimators=100,
10    max_features='sqrt',
11    oob_score=True,
12    random_state=42,
13    n_jobs=-1
14 )
15 rf.fit(X_train, y_train)
16 y_pred_rf = rf.predict(X_test)
17 acc_rf = accuracy_score(y_test, y_pred_rf)
18
19 print(f"Decision Tree   : Accuracy = {acc_dt:.4f}")
20 print(f"Random Forest   : Accuracy = {acc_rf:.4f}")
21 print(f"Score OOB        : {rf.oob_score_:.4f}")
22 print(f"\nGain : +{((acc_rf - acc_dt)*100:.1f)} points de
    pourcentage")

```

13.5.3 Étape 3 : Importance des variables

```

1 # Classement des variables par importance
2 importances = rf.feature_importances_
3 indices = np.argsort(importances)[::-1]
4
5 # Affichage des 10 plus importantes
6 print("Top 10 variables :")
7 for i in range(10):
8     print(f"   {i+1}. {data.feature_names[indices[i]]}"
9           f" : {importances[indices[i]]:.4f}")
10
11 # Graphique
12 plt.figure(figsize=(10, 6))
13 top_n = 15
14 plt.barh(range(top_n),
15          importances[indices[:top_n]][::-1],
16          color='steelblue', edgecolor='navy')

```

```

17 plt.yticks(range(top_n),
18             [data.feature_names[i] for i in indices[:top_n]][::-1])
19 plt.xlabel("Importance (MDI)")
20 plt.title("Top 15 variables - Random Forest")
21 plt.tight_layout()
22 plt.show()

```

13.5.4 Étape 4 : Score OOB

```

1 # Le score OOB est déjà calculé grâce à oob_score=True
2 print(f"Score OOB : {rf.oob_score_:.4f}")
3 print(f"Score Test : {acc_rf:.4f}")
4 print(f"Difference : {abs(rf.oob_score_ - acc_rf):.4f}")
5 print("=> Le score OOB est une bonne estimation de la performance
      !")

```

13.5.5 Étape 5 : Accuracy en fonction du nombre d'arbres

```

1 # Evolution de l'accuracy avec le nombre d'arbres
2 n_arbres_list = [1, 5, 10, 20, 50, 100, 200, 300, 500]
3 scores_train = []
4 scores_test = []
5 scores_oob = []
6
7 for n in n_arbres_list:
8     rf_temp = RandomForestClassifier(
9         n_estimators=n, max_features='sqrt',
10        oob_score=True, random_state=42, n_jobs=-1
11    )
12    rf_temp.fit(X_train, y_train)
13    scores_train.append(rf_temp.score(X_train, y_train))
14    scores_test.append(rf_temp.score(X_test, y_test))
15    scores_oob.append(rf_temp.oob_score_)
16
17 plt.figure(figsize=(10, 6))
18 plt.plot(n_arbres_list, scores_train, 'o-', label='Train',
19         color='blue')
20 plt.plot(n_arbres_list, scores_test, 's-', label='Test',
21         color='green')

```

```

20 plt.plot(n_arbres_list, scores_oob, '^--', label='OOB',
           color='red')
21 plt.xlabel("Nombre d'arbres (n_estimators)")
22 plt.ylabel("Accuracy")
23 plt.title("Accuracy vs nombre d'arbres")
24 plt.legend()
25 plt.grid(True, alpha=0.3)
26 plt.ylim(0.85, 1.01)
27 plt.tight_layout()
28 plt.show()

```

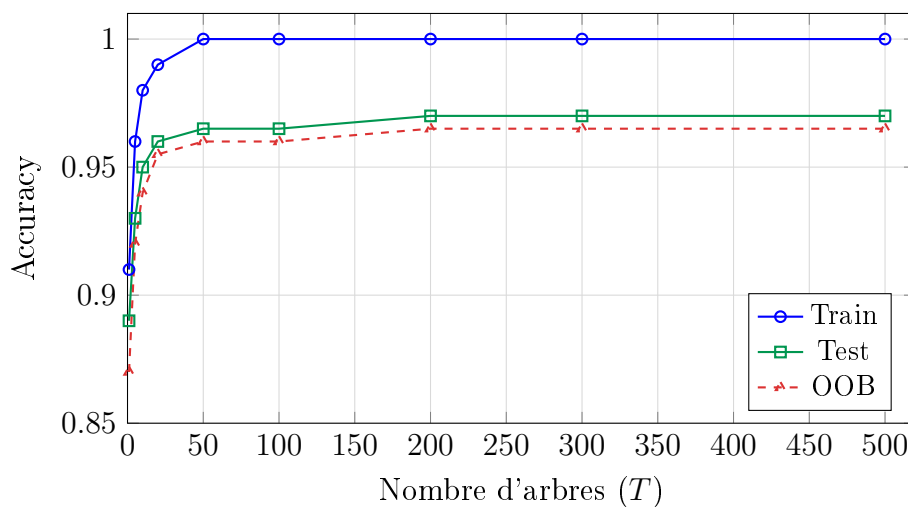


FIGURE 13.5 – Évolution de l'accuracy en fonction du nombre d'arbres.

13.5.6 Étape 6 : Frontière de décision — 1 arbre vs 100 arbres

```

1 from sklearn.decomposition import PCA
2 from matplotlib.colors import ListedColormap
3
4 # Reduire a 2 dimensions pour la visualisation
5 pca = PCA(n_components=2)
6 X_2d = pca.fit_transform(X)
7 X_train_2d, X_test_2d, y_train_2d, y_test_2d = train_test_split(
8     X_2d, y, test_size=0.2, random_state=42, stratify=y
9 )
10
11 fig, axes = plt.subplots(1, 2, figsize=(14, 5))
12 cmap_light = ListedColormap(['#FFAAAA', '#AAFFAA'])
13 cmap_bold = ListedColormap(['#FF0000', '#00AA00'])

```

```

14
15 for ax, model, title in zip(
16     axes,
17     [DecisionTreeClassifier(random_state=42),
18      RandomForestClassifier(n_estimators=100, random_state=42)],
19     ['1 arbre (Decision Tree)', '100 arbres (Random Forest)']
20 ):
21     model.fit(X_train_2d, y_train_2d)
22
23     # Grille de prediction
24     x_min, x_max = X_2d[:, 0].min() - 1, X_2d[:, 0].max() + 1
25     y_min, y_max = X_2d[:, 1].min() - 1, X_2d[:, 1].max() + 1
26     xx, yy = np.meshgrid(
27         np.linspace(x_min, x_max, 200),
28         np.linspace(y_min, y_max, 200)
29     )
30     Z = model.predict(np.c_[xx.ravel(), yy.ravel()])
31     Z = Z.reshape(xx.shape)
32
33     ax.contourf(xx, yy, Z, alpha=0.3, cmap=cmap_light)
34     ax.scatter(X_test_2d[:, 0], X_test_2d[:, 1],
35               c=y_test_2d, cmap=cmap_bold, edgecolors='k', s=30)
36     acc = model.score(X_test_2d, y_test_2d)
37     ax.set_title(f"{title}\nAccuracy = {acc:.2%}")
38     ax.set_xlabel("PC1")
39     ax.set_ylabel("PC2")
40
41 plt.suptitle("Frontiere de decision : 1 arbre vs 100 arbres",
42             fontsize=14, fontweight='bold')
43 plt.tight_layout()
44 plt.show()

```

13.5.7 Étape 7 : Validation croisée

```

1 from sklearn.model_selection import cross_val_score
2
3 # Validation croisee 5-fold
4 cv_dt = cross_val_score(
5     DecisionTreeClassifier(random_state=42),
6     X, y, cv=5, scoring='accuracy'

```

```

7 )
8 cv_rf = cross_val_score(
9     RandomForestClassifier(n_estimators=100, random_state=42,
10        n_jobs=-1),
11     X, y, cv=5, scoring='accuracy'
12 )
13 print("=== Validation croisee (5-fold) ===")
14 print(f"Decision Tree : {cv_dt.mean():.4f} (+/-
15        {cv_dt.std():.4f})")
16 print(f"Random Forest : {cv_rf.mean():.4f} (+/-
17        {cv_rf.std():.4f})")
18
19 # La Random Forest a une meilleure moyenne ET un ecart-type plus
20 # faible
21 # => plus performante ET plus stable

```

13.6 Schéma récapitulatif : phases d'entraînement et de test

PHASE D'ENTRAÎNEMENT (Training)

1. Étape 1 — Charger et séparer les données

Charger le dataset et le séparer en 80% Train, 20% Test.

Code Python :

```
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2)
```

2. Étape 2 — Choisir les hyperparamètres

Définir le nombre d'arbres et le nombre de features par arbre.

Mathématiques :

B = nombre d'arbres, m = nombre de features par arbre ($m \approx \sqrt{p}$ en classification)

Code Python :

```
from sklearn.ensemble import RandomForestClassifier
rf = RandomForestClassifier(n_estimators=100, max_features='sqrt')
```

3. Étape 3 — Échantillonnage bootstrap

Pour chaque arbre, tirer un échantillon **bootstrap** (avec remise) à partir des données de Train.

Mathématiques :

Pour chaque arbre $b = 1, \dots, B$: tirer n observations avec remise \rightarrow échantillon S_b

Code Python :

Le bootstrap est fait automatiquement par sklearn

4. Étape 4 — Construire chaque arbre

Construire chaque arbre de décision sur son échantillon bootstrap, en sélectionnant à chaque nœud un **sous-ensemble aléatoire de features**.

Mathématiques :

À chaque nœud, choisir le meilleur split parmi m features aléatoires (pas p).

Code Python :

`rf.fit(X_train, y_train)` # construit B arbres en parallele



Modèle entraîné transmis à la phase de test

PHASE DE TEST (Testing)

1. Étape 1 — Passer dans chaque arbre

Pour chaque point de Test, le faire passer dans **chaque arbre** de la forêt.

Mathématiques :

Pour chaque \mathbf{x}_{test} , chaque arbre b donne $\hat{y}_b(\mathbf{x})$

Code Python :

`y_pred = rf.predict(X_test)` # vote majoritaire automatique

2. Étape 2 — Vote majoritaire

La prédiction finale est le **vote majoritaire** de tous les arbres de la forêt.

Mathématiques :

$$\hat{y} = \text{mode}\{\hat{y}_1, \hat{y}_2, \dots, \hat{y}_B\}$$

Code Python :

`y_proba = rf.predict_proba(X_test)` # probabilites moyennes

3. Étape 3 — Évaluer et analyser l'importance

Calculer les métriques de performance et l'importance des variables.

Mathématiques :

Importance de x_j = réduction moyenne du Gini apportée par x_j .

Code Python :

```
print(classification_report(y_test, y_pred))
importances = rf.feature_importances_
plt.barh(X.columns, importances)
```

13.7 Exercices

Exercice 13.1 — Bootstrap et OOB

Soit le jeu de données original suivant avec 5 observations :

Observation	Variable X	Classe Y
A	2.1	0
B	3.5	1
C	1.8	0
D	4.2	1
E	2.9	1

1. Créez 3 échantillons bootstrap de taille 5 (choisissez vous-même les tirages, mais de manière réaliste avec des répétitions).
2. Pour chaque échantillon bootstrap, identifiez les observations OOB (Out-of-Bag).
3. Si l'observation B est OOB pour les arbres 1 et 3, et que l'arbre 1 prédit $\hat{y} = 1$ et l'arbre 3 prédit $\hat{y} = 0$, quelle est la prédiction OOB pour B ? Comment départager ?
4. Calculez la proportion théorique d'observations présentes dans un bootstrap pour $n = 5$, puis comparez avec la valeur théorique pour $n \rightarrow \infty$.

Correction de l'exercice 13.1

1. Trois échantillons bootstrap possibles (tirés avec remise) :
 - $\mathcal{D}_1^* = \{A, C, C, D, E\}$ (C apparaît 2 fois)
 - $\mathcal{D}_2^* = \{B, B, D, A, E\}$ (B apparaît 2 fois)
 - $\mathcal{D}_3^* = \{E, A, D, D, A\}$ (A et D apparaissent 2 fois)
2. Observations OOB pour chaque bootstrap :
 - OOB de \mathcal{D}_1^* : {B} (absente du bootstrap 1)
 - OOB de \mathcal{D}_2^* : {C} (absente du bootstrap 2)

— OOB de \mathcal{D}_3^* : {B, C} (absentes du bootstrap 3)

3. Pour l'observation B, OOB dans les arbres 1 et 3 :

— Arbre 1 : $\hat{y} = 1$, Arbre 3 : $\hat{y} = 0$

— Vote : 1 vote pour la classe 1, 1 vote pour la classe 0 \Rightarrow **égalité**.

— En cas d'égalité, on peut : (a) choisir aléatoirement, (b) utiliser les probabilités de chaque arbre, ou (c) ajouter plus d'arbres pour éviter ce cas. Scikit-learn choisit la classe avec la plus haute probabilité moyenne.

4. Proportion théorique pour $n = 5$:

$$P(\text{non sélectionnée}) = \left(1 - \frac{1}{5}\right)^5 = \left(\frac{4}{5}\right)^5 = \frac{1024}{3125} = 0.3277$$

$$P(\text{présente}) = 1 - 0.3277 = \mathbf{0.6723} \quad (67.2\%)$$

Pour $n \rightarrow \infty$: $P(\text{présente}) = 1 - 1/e \approx 0.6321$ (63.2%).

La valeur pour $n = 5$ (67.2%) est légèrement supérieure à la limite théorique (63.2%) car la convergence est progressive.

Exercice 13.2 — Vote majoritaire

Une Random Forest de 5 arbres est utilisée pour classifier un nouveau point x . Chaque arbre donne sa prédiction (0 ou 1) :

	Arbre 1	Arbre 2	Arbre 3	Arbre 4	Arbre 5
$\hat{y}_t(x)$	1	0	1	1	0

1. Quelle est la prédiction finale de la Random Forest pour le point x ?
2. Si chaque arbre fournit aussi une probabilité d'appartenance à la classe 1 :

	Arbre 1	Arbre 2	Arbre 3	Arbre 4	Arbre 5
$P(y = 1 \mid x)$	0.85	0.30	0.72	0.91	0.40

Quelle est la probabilité moyenne ? Le seuil étant 0.5, la prédiction change-t-elle ?

3. Si on avait 100 arbres et que 58 prédisent la classe 1, quelle serait la prédiction ? Avec quelle confiance ?

Correction de l'exercice 13.2**1. Vote majoritaire :**

— Classe 1 : 3 votes (arbres 1, 3 et 4)

— Classe 0 : 2 votes (arbres 2 et 5)

La majorité est la classe 1 $\Rightarrow \hat{y}(x) = \boxed{1}$

2. Probabilité moyenne :

$$P(y = 1 | x) = \frac{0.85 + 0.30 + 0.72 + 0.91 + 0.40}{5} = \frac{3.18}{5} = \boxed{0.636}$$

Comme $0.636 > 0.5$, la prédiction est également la classe 1. La prédiction **ne change pas**.

Cependant, la probabilité moyenne (0.636) donne une information plus riche que le simple vote : elle indique un degré de confiance modéré.

3. Avec 100 arbres et 58 votes pour la classe 1 :

— Prédiction : classe 1 ($58 > 50$)

— Confiance : $58/100 = 58\%$

La confiance est relativement faible (proche de 50%), ce qui suggère que le point x se situe près de la frontière de décision. Il faudrait être prudent avec cette prédiction.

Exercice 13.3 — Implémentation Python complète

En utilisant le dataset `breast_cancer` de scikit-learn :

1. Chargez les données et séparez-les en ensembles d'entraînement (80%) et de test (20%).
2. Entraînez un `DecisionTreeClassifier` et un `RandomForestClassifier` avec 100 arbres. Comparez les accuracy.
3. Faites varier `n_estimators` de 1 à 300 et tracez la courbe d'accuracy en fonction du nombre d'arbres (sur le jeu de test et le score OOB).
4. Affichez les 10 variables les plus importantes sous forme de diagramme en barres horizontales.
5. Réalisez une validation croisée 10-fold pour les deux modèles et comparez les résultats (moyenne et écart-type).

Correction de l'exercice 13.3**Question 1 : Chargement et séparation des données**

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.datasets import load_breast_cancer
4 from sklearn.model_selection import train_test_split,
   cross_val_score
5 from sklearn.tree import DecisionTreeClassifier
6 from sklearn.ensemble import RandomForestClassifier
7 from sklearn.metrics import accuracy_score
8
9 # 1. Charger et separer les donnees
10 data = load_breast_cancer()
11 X, y = data.data, data.target
12 X_train, X_test, y_train, y_test = train_test_split(
13     X, y, test_size=0.2, random_state=42, stratify=y
14 )
15 print(f"Train : {X_train.shape[0]} obs, Test :
      {X_test.shape[0]} obs")

```

Question 2 : Comparaison Decision Tree vs Random Forest

```

1 # 2. Entraîner et comparer
2 dt = DecisionTreeClassifier(random_state=42)
3 dt.fit(X_train, y_train)
4 acc_dt = accuracy_score(y_test, dt.predict(X_test))
5
6 rf = RandomForestClassifier(n_estimators=100, oob_score=True,
7                             random_state=42, n_jobs=-1)
8 rf.fit(X_train, y_train)
9 acc_rf = accuracy_score(y_test, rf.predict(X_test))
10
11 print(f"Decision Tree   : {acc_dt:.4f}")
12 print(f"Random Forest   : {acc_rf:.4f}")
13 print(f"Score OOB        : {rf.oob_score_:.4f}")
14 print(f"Amelioration     : +{(acc_rf - acc_dt)*100:.1f} pp")

```

Question 3 : Courbe accuracy vs nombre d'arbres

```

1 # 3. Accuracy vs n_estimators
2 n_values = list(range(1, 301, 10))
3 test_scores = []
4 oob_scores = []

```

```

5
6 for n in n_values:
7     rf_temp = RandomForestClassifier(
8         n_estimators=n, oob_score=True,
9         random_state=42, n_jobs=-1
10    )
11    rf_temp.fit(X_train, y_train)
12    test_scores.append(rf_temp.score(X_test, y_test))
13    oob_scores.append(rf_temp.oob_score_)
14
15 plt.figure(figsize=(10, 5))
16 plt.plot(n_values, test_scores, '-', label='Test',
17         color='green', lw=2)
18 plt.plot(n_values, oob_scores, '--', label='OOB', color='red',
19         lw=2)
20 plt.axhline(y=acc_dt, color='blue', linestyle=':',
21         label='Decision Tree')
22 plt.xlabel("Nombre d'arbres")
23 plt.ylabel("Accuracy")
24 plt.title("Accuracy vs nombre d'arbres (Random Forest)")
25 plt.legend()
26 plt.grid(True, alpha=0.3)
27 plt.tight_layout()
28 plt.show()

```

Question 4 : Importance des variables

```

1 # 4. Top 10 variables les plus importantes
2 importances = rf.feature_importances_
3 indices = np.argsort(importances)[::-1][:10]
4
5 plt.figure(figsize=(8, 5))
6 plt.barh(range(10), importances[indices][::-1],
7         color='steelblue', edgecolor='navy')
8 plt.yticks(range(10),
9         [data.feature_names[i] for i in indices][::-1])
10 plt.xlabel("Importance (MDI)")
11 plt.title("Top 10 variables - Breast Cancer")
12 plt.tight_layout()
13 plt.show()

```

```
14
15 # Affichage textuel
16 for i, idx in enumerate(indices):
17     print(f"    {i+1}. {data.feature_names[idx]} : "
18           f"{importances[idx]:.4f}")
```

Question 5 : Validation croisée 10-fold

```
1 # 5. Validation croisee 10-fold
2 cv_dt = cross_val_score(
3     DecisionTreeClassifier(random_state=42),
4     X, y, cv=10, scoring='accuracy'
5 )
6 cv_rf = cross_val_score(
7     RandomForestClassifier(n_estimators=100, random_state=42,
8                             n_jobs=-1),
9     X, y, cv=10, scoring='accuracy'
10 )
11
12 print("=== Validation croisee 10-fold ===")
13 print(f"Decision Tree : {cv_dt.mean():.4f} (+/-
14       {cv_dt.std():.4f})")
15 print(f"Random Forest : {cv_rf.mean():.4f} (+/-
16       {cv_rf.std():.4f})")
17
18 # Resultats attendus :
19 # Decision Tree : ~0.9280 (+/- 0.0280)
20 # Random Forest : ~0.9630 (+/- 0.0190)
21 # => Random Forest : meilleure moyenne ET ecart-type plus faible
22 #      (plus performante ET plus stable)
```

Conclusion de l'exercice

La Random Forest surpasse systématiquement l'arbre de décision seul :

- **Accuracy plus élevée** : gain de +3 à +5 points de pourcentage.
- **Écart-type plus faible** en validation croisée : le modèle est plus stable.
- **La performance se stabilise** au-delà de 50–100 arbres : inutile d'en ajouter des milliers.
- Le **score OOB** est un bon estimateur de la performance générale, sans nécessiter de jeu de test séparé.

Chapitre 14

Machines à Vecteurs de Support (SVM)

14.1 Hands-On : séparer des billes sur une table

Hands-On : Quelle est la « meilleure » ligne de séparation ?

Imaginez que vous avez des billes **rouges** et des billes **bleues** posées sur une table. Vous voulez tracer une ligne droite pour séparer les deux groupes. Le problème, c'est qu'il existe **une infinité de lignes possibles** qui séparent correctement les billes. Mais laquelle est la **meilleure** ?

Regardez la figure ci-dessous : trois lignes séparent parfaitement les deux classes. Pourtant, intuitivement, la ligne (C) semble la meilleure. Pourquoi ? Parce qu'elle est **la plus éloignée** des deux groupes — elle maximise la **marge** entre les classes.

Trois lignes de séparation possibles — laquelle est la meilleure ?

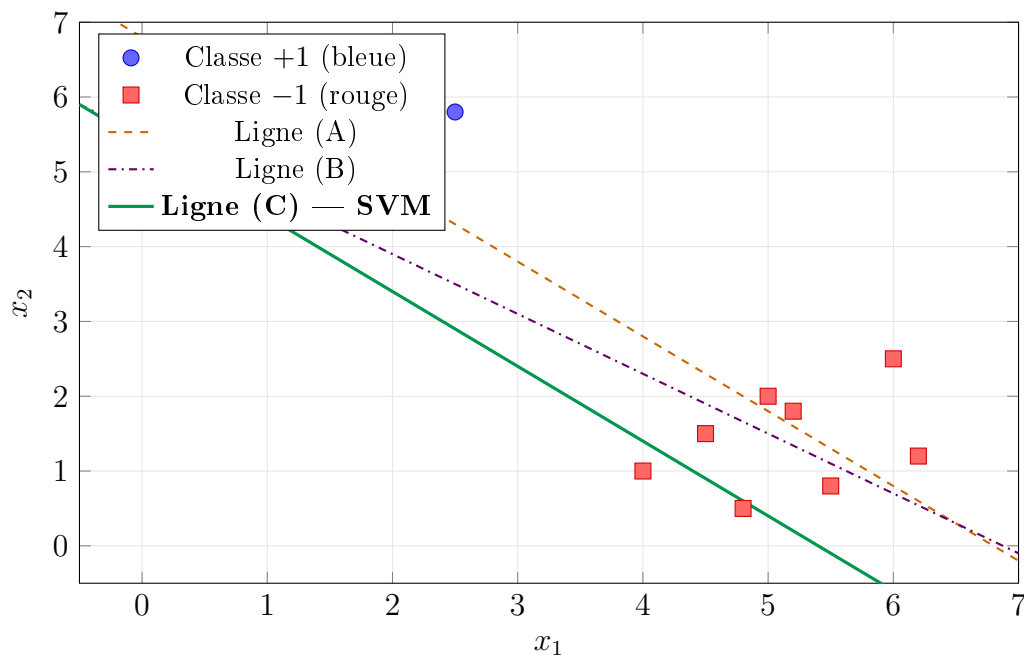


FIGURE 14.1 – Trois lignes séparent les deux classes, mais seule la ligne (C) maximise la marge.

La ligne (A) passe trop près des points bleus : un nouveau point bleu légèrement décalé serait mal classé. La ligne (B) passe trop près des points rouges. La ligne (C), en revanche, est à égale distance des deux groupes : c'est exactement ce que fait le **SVM**.

14.2 Intuition : construire la route la plus large

L'idée clé du SVM

Imaginez que vous devez construire une **route** entre deux quartiers : un quartier bleu et un quartier rouge. Vous voulez que cette route soit **la plus large possible**, pour qu'aucune maison ne soit trop près de la route.

- La **ligne centrale** de la route est l'**hyperplan séparateur**.
- La **largeur de la route** est la **marge**.
- Les **maisons les plus proches** de la route (celles qui « touchent » le bord) sont les **vecteurs de support** (support vectors).
- Le SVM cherche la route qui est **la plus large possible** tout en séparant correctement les deux quartiers.

Pourquoi vouloir la route la plus large ? Parce qu'une route large donne une meilleure **généralisation** : si un nouveau point arrive, il a plus de chances d'être du bon côté.

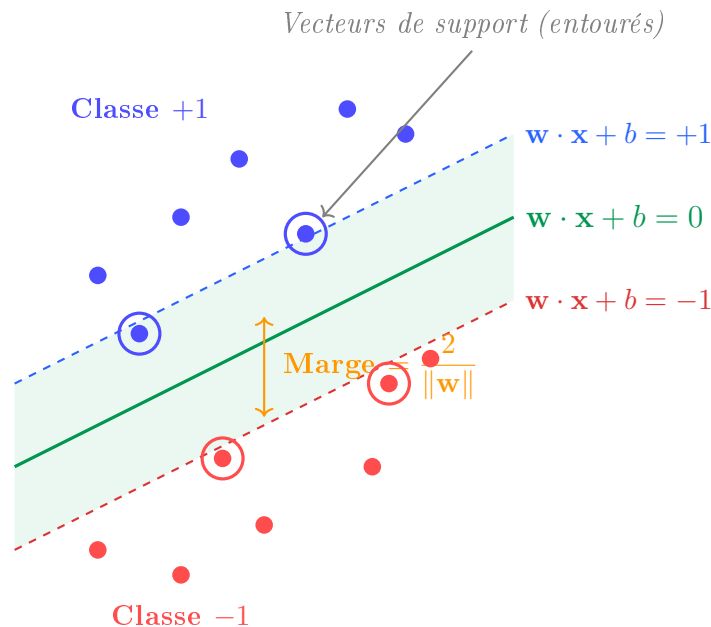


FIGURE 14.2 – L'hyperplan séparateur, la marge et les vecteurs de support.

Vecteurs de support

Les **vecteurs de support** (support vectors) sont les points d'entraînement **les plus proches** de l'hyperplan séparateur. Ce sont eux qui « supportent » l'hyperplan : si on les déplace, l'hyperplan change. Les autres points (plus éloignés) n'ont **aucune influence** sur la position de l'hyperplan. C'est ce qui rend le SVM **robuste** et élégant.

14.3 Dérivation mathématique détaillée

14.3.1 L'hyperplan séparateur

Hyperplan

Un **hyperplan** est la généralisation d'une droite (en 2D) ou d'un plan (en 3D) à un espace de dimension quelconque. Il est défini par l'équation :

$$\mathbf{w} \cdot \mathbf{x} + b = 0$$

où :

- $\mathbf{w} = (w_1, w_2, \dots, w_n)$ est le **vecteur normal** à l'hyperplan (il indique la direction perpendiculaire),
- $\mathbf{x} = (x_1, x_2, \dots, x_n)$ est un point de l'espace,
- b est le **biais** (le décalage par rapport à l'origine),
- $\mathbf{w} \cdot \mathbf{x} = w_1x_1 + w_2x_2 + \dots + w_nx_n$ est le **produit scalaire**.

En fonction de la dimension de l'espace :

- En **2D** : l'hyperplan est une **droite** $w_1x_1 + w_2x_2 + b = 0$.
- En **3D** : l'hyperplan est un **plan** $w_1x_1 + w_2x_2 + w_3x_3 + b = 0$.
- En **n D** : l'hyperplan est une surface de dimension $n - 1$.

L'hyperplan divise l'espace en **deux demi-espaces** :

- Les points où $\mathbf{w} \cdot \mathbf{x} + b > 0$ sont classés $+1$ (classe positive).
- Les points où $\mathbf{w} \cdot \mathbf{x} + b < 0$ sont classés -1 (classe négative).

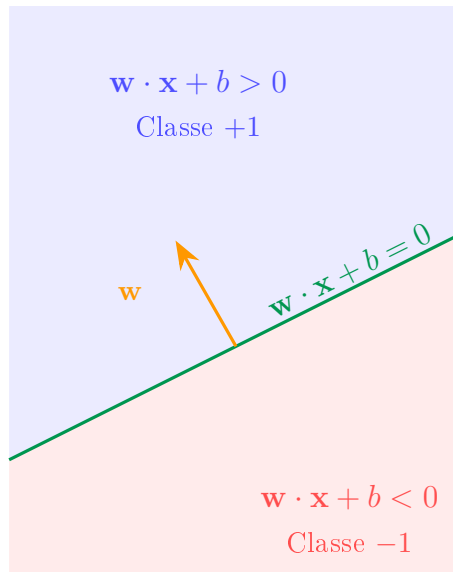


FIGURE 14.3 – L’hyperplan sépare l’espace en deux demi-espaces. Le vecteur \mathbf{w} est perpendiculaire à l’hyperplan.

14.3.2 Distance d’un point à l’hyperplan

Pour trouver la meilleure ligne de séparation, nous devons pouvoir **mesurer la distance** entre un point et l’hyperplan.

Distance d’un point à l’hyperplan

La distance d’un point \mathbf{x}_0 à l’hyperplan $\mathbf{w} \cdot \mathbf{x} + b = 0$ est :

$$d(\mathbf{x}_0) = \frac{|\mathbf{w} \cdot \mathbf{x}_0 + b|}{\|\mathbf{w}\|}$$

où $\|\mathbf{w}\| = \sqrt{w_1^2 + w_2^2 + \dots + w_n^2}$ est la **norme** (longueur) du vecteur \mathbf{w} .

Dérivation géométrique. Soit \mathbf{x}_0 un point quelconque et \mathbf{x}_p sa **projection orthogonale** sur l’hyperplan. Le vecteur allant de \mathbf{x}_p à \mathbf{x}_0 est perpendiculaire à l’hyperplan, donc parallèle à \mathbf{w} :

$$\mathbf{x}_0 - \mathbf{x}_p = d \cdot \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

Puisque \mathbf{x}_p est sur l’hyperplan, on a $\mathbf{w} \cdot \mathbf{x}_p + b = 0$. En calculant le produit scalaire des deux côtés avec \mathbf{w} :

$$\mathbf{w} \cdot \mathbf{x}_0 = \mathbf{w} \cdot \mathbf{x}_p + d \cdot \frac{\mathbf{w} \cdot \mathbf{w}}{\|\mathbf{w}\|} = -b + d \cdot \|\mathbf{w}\|$$

On isole d :

$$d = \frac{\mathbf{w} \cdot \mathbf{x}_0 + b}{\|\mathbf{w}\|}$$

En prenant la valeur absolue pour obtenir une distance positive :

$$d(\mathbf{x}_0) = \frac{|\mathbf{w} \cdot \mathbf{x}_0 + b|}{\|\mathbf{w}\|}$$

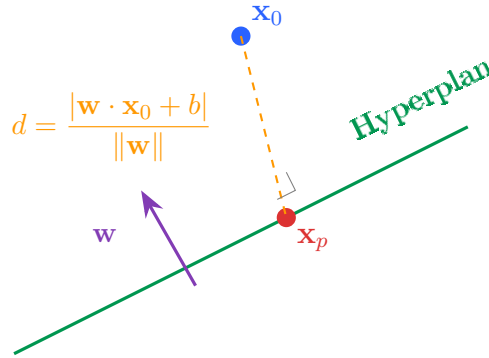


FIGURE 14.4 – La distance d'un point \mathbf{x}_0 à l'hyperplan est la longueur de la projection orthogonale.

Exemple numérique

Soit l'hyperplan $2x_1 + x_2 - 4 = 0$ (donc $\mathbf{w} = (2, 1)$ et $b = -4$).

Distance du point $\mathbf{x}_0 = (3, 5)$ à cet hyperplan :

$$d = \frac{|2 \times 3 + 1 \times 5 + (-4)|}{\sqrt{2^2 + 1^2}} = \frac{|6 + 5 - 4|}{\sqrt{5}} = \frac{7}{\sqrt{5}} \approx 3.13$$

14.3.3 La marge

Le SVM cherche l'hyperplan qui maximise la **marge**, c'est-à-dire la distance entre les deux points les plus proches de classes opposées.

Les **vecteurs de support** sont les points les plus proches de l'hyperplan. On normalise les équations pour que ces points satisfassent exactement :

$$\mathbf{w} \cdot \mathbf{x}_+ + b = +1 \quad (\text{vecteur de support de classe } +1) \quad (14.1)$$

$$\mathbf{w} \cdot \mathbf{x}_- + b = -1 \quad (\text{vecteur de support de classe } -1) \quad (14.2)$$

La distance d'un vecteur de support de classe $+1$ à l'hyperplan est :

$$d_+ = \frac{|\mathbf{w} \cdot \mathbf{x}_+ + b|}{\|\mathbf{w}\|} = \frac{|+1|}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|}$$

De même, la distance d'un vecteur de support de classe -1 est :

$$d_- = \frac{|\mathbf{w} \cdot \mathbf{x}_- + b|}{\|\mathbf{w}\|} = \frac{|-1|}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|}$$

La **marge totale** est la somme des deux distances :

$$\text{Marge} = d_+ + d_- = \frac{2}{\|\mathbf{w}\|}$$

Pourquoi +1 et -1 ?

Le choix de +1 et -1 dans les équations (14.1) et (14.2) est une **convention de normalisation**. On peut toujours multiplier \mathbf{w} et b par un scalaire positif pour que les vecteurs de support satisfassent exactement ces équations. Ce choix simplifie énormément les calculs.

14.3.4 Le problème d'optimisation

Problème d'optimisation du SVM (forme primale)

Le SVM cherche à **maximiser la marge** tout en classant correctement tous les points :

$$\max_{\mathbf{w}, b} \frac{2}{\|\mathbf{w}\|} \iff \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

sous la contrainte :

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 \quad \text{pour tout } i = 1, 2, \dots, N$$

où $y_i \in \{-1, +1\}$ est l'étiquette du point \mathbf{x}_i .

Explication de la contrainte. Pour chaque point d'entraînement :

- Si $y_i = +1$: on veut $\mathbf{w} \cdot \mathbf{x}_i + b \geq +1$ (le point est du bon côté, au-delà de la marge).
- Si $y_i = -1$: on veut $\mathbf{w} \cdot \mathbf{x}_i + b \leq -1$ (idem, de l'autre côté).
- En multipliant par y_i , ces deux cas se résument en une seule inégalité : $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$.

Pourquoi minimiser $\frac{1}{2} \|\mathbf{w}\|^2$?

1. Maximiser $\frac{2}{\|\mathbf{w}\|}$ revient à minimiser $\|\mathbf{w}\|$.
2. Minimiser $\|\mathbf{w}\|$ est équivalent à minimiser $\|\mathbf{w}\|^2$ (la fonction racine carrée est croissante).
3. Le facteur $\frac{1}{2}$ simplifie les dérivées (la dérivée de $\frac{1}{2}x^2$ est x , pas $2x$).

Ce problème est un problème de **programmation quadratique** (la fonction objectif est quadratique, les contraintes sont linéaires). Il est **convexe**, ce qui garantit l'existence d'une **solution unique** (le minimum global).

Les multiplicateurs de Lagrange (simplifié)

Pour résoudre un problème d'optimisation avec contraintes, on utilise la méthode des **multiplicateurs de Lagrange**. L'idée : au lieu de résoudre un problème contraint, on intègre les contraintes dans la fonction objectif en ajoutant un « coût » pour chaque contrainte violée.

On introduit un multiplicateur $\alpha_i \geq 0$ pour chaque contrainte :

$$\mathcal{L}(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1]$$

En dérivant et en égalisant à zéro :

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 0 \implies \mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \quad \frac{\partial \mathcal{L}}{\partial b} = 0 \implies \sum_{i=1}^N \alpha_i y_i = 0$$

Résultat clé : \mathbf{w} est une **combinaison linéaire des points d'entraînement**. De plus, la condition KKT impose $\alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1] = 0$. Donc :

- Si $\alpha_i > 0$: le point est un **vecteur de support** (il est sur la marge).
- Si $\alpha_i = 0$: le point n'a **aucune influence** sur l'hyperplan.

Le problème dual. En remplaçant \mathbf{w} dans le lagrangien, on obtient le **problème dual** :

$$\max_{\boldsymbol{\alpha}} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

sous les contraintes $\alpha_i \geq 0$ et $\sum_{i=1}^N \alpha_i y_i = 0$.

L'avantage du dual : il ne dépend que des **produits scalaires** $\mathbf{x}_i \cdot \mathbf{x}_j$, ce qui sera essentiel pour l'**astuce du noyau** (kernel trick).

14.3.5 SVM à marge souple (données non séparables)

Dans la réalité, les données sont rarement séparables parfaitement : il y a toujours du **bruit** et des points aberrants.

SVM à marge souple (Soft Margin)

On introduit des **variables d'écart** (slack variables) $\xi_i \geq 0$ pour chaque point. Le nouveau problème d'optimisation est :

$$\min_{\mathbf{w}, b, \boldsymbol{\xi}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i$$

sous les contraintes :

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i \quad \text{et} \quad \xi_i \geq 0 \quad \text{pour tout } i$$

Interprétation des variables d'écart :

- $\xi_i = 0$: le point est correctement classé et en dehors de la marge.
- $0 < \xi_i < 1$: le point est dans la marge mais du bon côté (correctement classé).
- $\xi_i \geq 1$: le point est **mal classé** (du mauvais côté de l'hyperplan).

Le paramètre de régularisation C

Le paramètre $C > 0$ contrôle le **compromis** entre :

- **Maximiser la marge** (terme $\frac{1}{2}\|\mathbf{w}\|^2$) : un modèle simple et régulier.
- **Minimiser les erreurs** (terme $C \sum \xi_i$) : un modèle qui classe bien les données d'entraînement.

	C grand (ex. 1000)	C petit (ex. 0.01)
Marge	Étroite	Large
Erreurs	Très peu tolérées	Plus tolérées
Risque	Surapprentissage (overfitting)	Sous-apprentissage (underfitting)

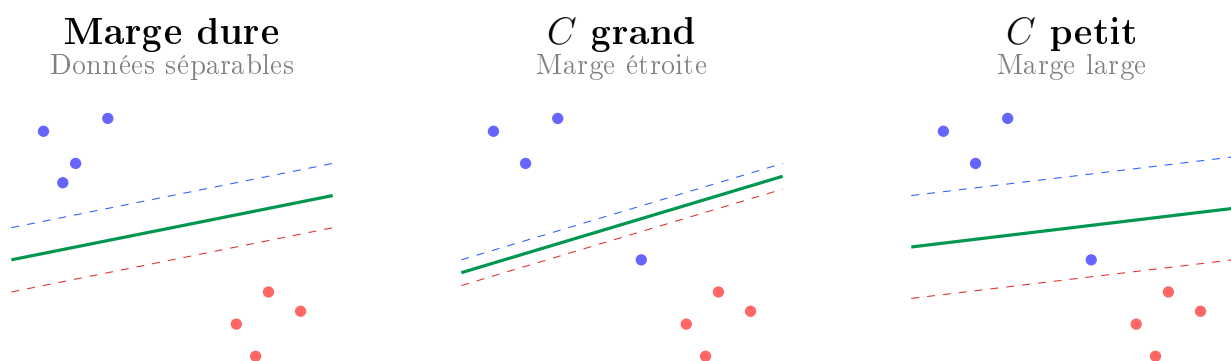


FIGURE 14.5 – Comparaison entre marge dure, marge souple avec C grand et C petit.

14.3.6 L'astuce du noyau (Kernel Trick)

Parfois, les données ne sont **pas du tout séparables linéairement**, quelle que soit la droite choisie. C'est là qu'intervient l'**astuce du noyau**, l'une des idées les plus brillantes du Machine Learning.

L'idée géniale

Si les données ne sont pas séparables dans l'espace d'origine, on les **projette dans un espace de dimension supérieure** où elles *deviennent* séparables !

Analogie : Imaginez des pièces de monnaie posées sur une table — certaines face, certaines pile, mélangées en cercle concentrique. Impossible de les séparer avec une droite. Mais si vous **soufflez** sous la table et que les pièces s'envolent à des hauteurs différentes, un plan horizontal peut les séparer dans l'espace 3D !

Exemple concret en 1D \rightarrow 2D. Considérons des points sur une droite : les bleus sont au centre et les rouges sont aux extrémités.

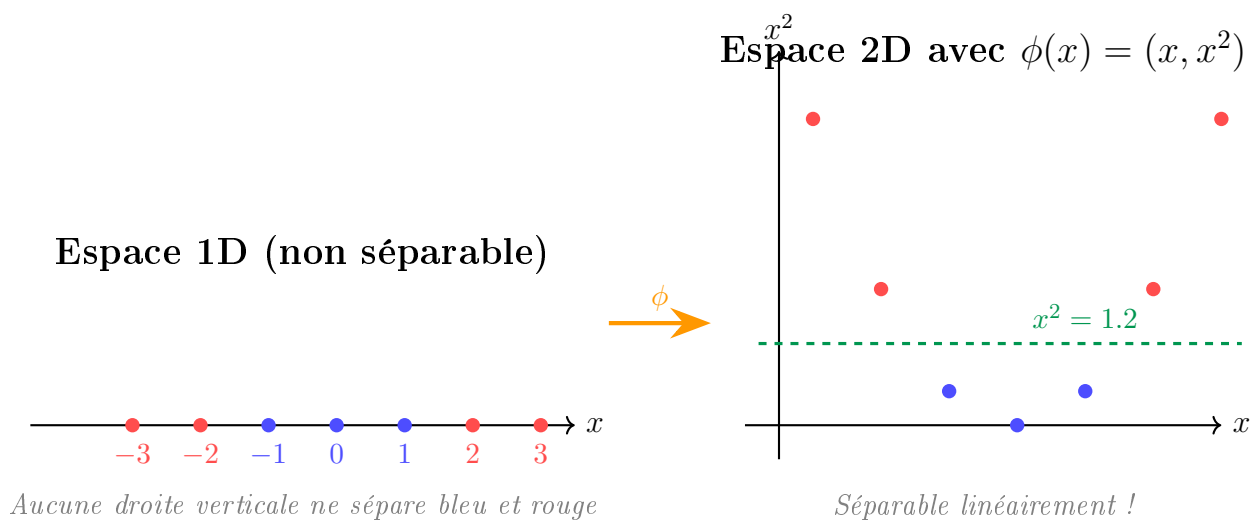


FIGURE 14.6 – En ajoutant la dimension x^2 , des données non séparables en 1D deviennent séparables en 2D.

Le problème. Calculer explicitement la transformation $\phi(\mathbf{x})$ peut être très coûteux : l'espace d'arrivée peut avoir des millions (voire une infinité) de dimensions.

La solution : le noyau (kernel). Rappelons que le problème dual ne dépend que des produits scalaires $\mathbf{x}_i \cdot \mathbf{x}_j$. Après transformation ϕ , on aurait besoin de $\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$. Un **noyau** est une fonction K telle que :

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$$

sans jamais calculer ϕ explicitement ! C'est l'**astuce du noyau** (kernel trick).

Noyaux courants

Noyau	Formule	Paramètres
Linéaire	$K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i \cdot \mathbf{x}_j$	Aucun
Polynomial	$K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i \cdot \mathbf{x}_j + c)^d$	$c \geq 0, d \in \mathbb{N}^*$
RBF (Gaussien)	$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\ \mathbf{x}_i - \mathbf{x}_j\ ^2}{2\sigma^2}\right)$	$\sigma > 0$ (ou $\gamma = \frac{1}{2\sigma^2}$)
Sigmoïde	$K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\alpha \mathbf{x}_i \cdot \mathbf{x}_j + c)$	$\alpha > 0, c$

Quel noyau choisir ?

Noyau	Quand l'utiliser
Linéaire	Données linéairement séparables, beaucoup de features ($n \gg N$), texte (TF-IDF). C'est le plus rapide.
Polynomial	Relations non linéaires modérées, traitement d'images. Degré 2 ou 3 généralement suffisant.
RBF	Choix par défaut . Fonctionne bien dans la plupart des cas. Espace de dimension infinie !
Sigmoïde	Rarement utilisé en pratique. Similaire à un réseau de neurones à une couche.

Le noyau RBF et le paramètre γ

Pour le noyau RBF, le paramètre $\gamma = \frac{1}{2\sigma^2}$ contrôle la « portée » de chaque point :

- γ **grand** (σ petit) : chaque point n'influence que ses voisins très proches \rightarrow frontière très sinueuse \rightarrow risque de **surapprentissage**.
- γ **petit** (σ grand) : chaque point influence une large zone \rightarrow frontière lisse \rightarrow risque de **sous-apprentissage**.

14.3.7 SVM pour la régression (SVR)

Le SVM peut aussi faire de la **régression** ! L'idée est inversée : au lieu de maximiser la marge entre les classes, on définit un **tube** de largeur ε autour de la fonction de prédiction.

Support Vector Regression (SVR)

Le SVR cherche une fonction $f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b$ telle que :

- Les points **dans le tube** ($|y_i - f(\mathbf{x}_i)| \leq \varepsilon$) ne sont **pas pénalisés**.
- Les points **hors du tube** sont pénalisés proportionnellement à leur distance au tube.

Le problème d'optimisation est :

$$\min_{\mathbf{w}, b, \xi, \xi^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N (\xi_i + \xi_i^*)$$

sous les contraintes :

$$y_i - (\mathbf{w} \cdot \mathbf{x}_i + b) \leq \varepsilon + \xi_i \quad \text{et} \quad (\mathbf{w} \cdot \mathbf{x}_i + b) - y_i \leq \varepsilon + \xi_i^*$$

avec $\xi_i, \xi_i^* \geq 0$.

SVR : le tube de largeur ε

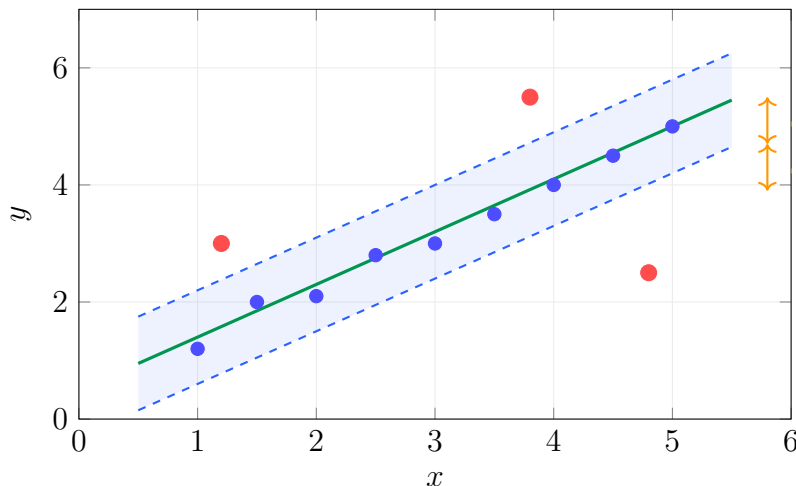


FIGURE 14.7 – Le SVR définit un tube de largeur 2ε . Seuls les points rouges (hors du tube) sont pénalisés.

14.4 Application sur le dataset clientèle

Appliquons le SVM à notre problème de **prédiction du churn** sur le dataset ShopTech.

```

1  import numpy as np
2  import pandas as pd
3  from sklearn.svm import SVC
4  from sklearn.preprocessing import StandardScaler
5  from sklearn.model_selection import train_test_split
6  from sklearn.metrics import accuracy_score, classification_report
7
8  # Notre dataset clientele
9  data = {
10     'Age':
11         [25,34,28,45,23,52,31,41,27,38,55,29,47,22,36,43,26,50,33,30],
12     'Revenu':
13         [22,45,28,62,18,75,35,55,24,48,80,30,65,20,42,58,25,70,38,32],
14     'Anciennete': [1,5,2,8,1,12,3,7,1,6,15,2,9,1,4,7,1,10,3,2],
15     'Nb_Achats':
16         [8,30,12,55,5,70,18,45,7,35,85,14,60,6,25,48,9,65,20,15],
17     'Montant_Moyen': [45,78,52,95,38,110,62,88,42,82,120,55,98,40,72,90,48,105,
18     'Score_Satisfaction': [5,8,4,9,3,9,6,8,4,7,10,5,9,3,7,8,4,9,6,5],
19     'Churn': [1,0,1,0,1,0,0,0,1,0,0,1,0,1,0,0,1,0,0,1]
20 }
21 df = pd.DataFrame(data)
22
23 X = df.drop('Churn', axis=1)
24 y = df['Churn']
25
26 # Separation train/test
27 X_train, X_test, y_train, y_test = train_test_split(
28     X, y, test_size=0.3, random_state=42)
29
30 # IMPORTANT : mise a l'echelle (obligatoire pour SVM !)
31 scaler = StandardScaler()
32 X_train_scaled = scaler.fit_transform(X_train)
33 X_test_scaled = scaler.transform(X_test)

```

Listing 14.1 – SVM pour la classification du churn

La mise à l'échelle est OBLIGATOIRE pour le SVM !

Le SVM est basé sur les **distances** entre les points. Si une variable a des valeurs entre 0 et 100 et une autre entre 0 et 1, la première dominera complètement le calcul des distances.

Il faut **toujours standardiser** les features avant d'appliquer un SVM :

$$x_{\text{standardisé}} = \frac{x - \mu}{\sigma}$$

```

1  # Tester différents noyaux
2  kernels = ['linear', 'poly', 'rbf', 'sigmoid']
3  results = {}
4
5  for kernel in kernels:
6      svm = SVC(kernel=kernel, C=1.0, random_state=42)
7      svm.fit(X_train_scaled, y_train)
8      y_pred = svm.predict(X_test_scaled)
9      acc = accuracy_score(y_test, y_pred)
10     n_sv = svm.n_support_  # nombre de vecteurs de support par
                             classe
11     results[kernel] = {'accuracy': acc, 'n_support_vectors': n_sv}
12     print(f''Noyau {kernel:10s} | Accuracy: {acc:.2f} | ''
13           f''Vecteurs de support: {n_sv} (total: {sum(n_sv)})''')
14
15 # Exemple de sortie :
16 # Noyau linear      | Accuracy: 0.83 | Vecteurs de support: [4 5]
    (total: 9)
17 # Noyau poly       | Accuracy: 0.83 | Vecteurs de support: [5 6]
    (total: 11)
18 # Noyau rbf        | Accuracy: 1.00 | Vecteurs de support: [4 5]
    (total: 9)
19 # Noyau sigmoid    | Accuracy: 0.67 | Vecteurs de support: [5 7]
    (total: 12)

```

Listing 14.2 – Comparaison de différents noyaux

```

1  # Entraîner le meilleur modèle (RBF)
2  svm_rbf = SVC(kernel='rbf', C=1.0, random_state=42)
3  svm_rbf.fit(X_train_scaled, y_train)
4
5  # Les vecteurs de support
6  print('Indices des vecteurs de support :', svm_rbf.support_)
7  print('Nombre par classe :', svm_rbf.n_support_)
8  print('Valeurs des vecteurs de support :')
9  print(X_train.iloc[svm_rbf.support_])

```

Listing 14.3 – Identification des vecteurs de support

14.5 Implémentation complète sur Google Colab

Google Colab : SVM complet avec visualisations

Ce notebook complet implémente le SVM de A à Z avec toutes les visualisations.

```

1  import numpy as np
2  import pandas as pd
3  import matplotlib.pyplot as plt
4  from sklearn.svm import SVC, SVR
5  from sklearn.preprocessing import StandardScaler
6  from sklearn.model_selection import train_test_split, GridSearchCV
7  from sklearn.metrics import (accuracy_score, classification_report,
8                               confusion_matrix, mean_squared_error,
9                               r2_score)
10
11  import warnings
12  warnings.filterwarnings('ignore')
13
14  # Dataset clientele
15  data = {
16      'Age':
17          [25,34,28,45,23,52,31,41,27,38,55,29,47,22,36,43,26,50,33,30],
18      'Revenu':
19          [22,45,28,62,18,75,35,55,24,48,80,30,65,20,42,58,25,70,38,32],
20      'Anciennete': [1,5,2,8,1,12,3,7,1,6,15,2,9,1,4,7,1,10,3,2],
21      'Nb_Achats':
22          [8,30,12,55,5,70,18,45,7,35,85,14,60,6,25,48,9,65,20,15],
23      'Montant_Moyen': [45,78,52,95,38,110,62,88,42,82,120,55,98,40,72,90,48,105,
24                      55,60,70,80,90,100,110,120,130,140,150,160,170,180,190,200],
25      'Score_Satisfaction': [5,8,4,9,3,9,6,8,4,7,10,5,9,3,7,8,4,9,6,5],
26      'Churn': [1,0,1,0,1,0,0,0,1,0,0,1,0,1,0,0,1,0,0,1]
27  }
28  df = pd.DataFrame(data)
29
30  X = df.drop('Churn', axis=1)
31  y = df['Churn']
32
33  X_train, X_test, y_train, y_test = train_test_split(

```

```

28     X, y, test_size=0.3, random_state=42)
29
30     scaler = StandardScaler()
31     X_train_scaled = scaler.fit_transform(X_train)
32     X_test_scaled = scaler.transform(X_test)
33
34     print('Taille train:', X_train_scaled.shape)
35     print('Taille test:', X_test_scaled.shape)

```

Listing 14.4 – Imports et préparation des données

```

1  def plot_decision_boundary(X, y, model, title, ax):
2      '''Trace la frontiere de decision pour 2 features.'''
3      h = 0.02
4      x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
5      y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
6      xx, yy = np.meshgrid(np.arange(x_min, x_max, h),
7                           np.arange(y_min, y_max, h))
8      Z = model.predict(np.c_[xx.ravel(), yy.ravel()])
9      Z = Z.reshape(xx.shape)
10
11     ax.contourf(xx, yy, Z, alpha=0.3, cmap='RdBu')
12     ax.contour(xx, yy, Z, colors='k', linewidths=0.5)
13
14     # Points d'entrainement
15     scatter = ax.scatter(X[:, 0], X[:, 1], c=y, cmap='RdBu',
16                        edgecolors='k', s=50, zorder=3)
17
18     # Vecteurs de support (entoures)
19     sv = model.support_vectors_
20     ax.scatter(sv[:, 0], sv[:, 1], s=150, facecolors='none',
21            edgecolors='gold', linewidths=2, zorder=4,
22            label=f'Vecteurs de support ({len(sv)})')
23     ax.set_title(title, fontsize=12, fontweight='bold')
24     ax.legend(fontsize=8)
25
26     # Utiliser 2 features pour la visualisation
27     X_2d = X_train_scaled[:, [0, 1]] # Age et Revenu (standardises)
28     y_2d = y_train.values
29
30     fig, axes = plt.subplots(2, 2, figsize=(12, 10))
31     kernels = ['linear', 'poly', 'rbf', 'sigmoid']

```

```

32
33 for ax, kernel in zip(axes.flatten(), kernels):
34     model = SVC(kernel=kernel, C=1.0, random_state=42)
35     model.fit(X_2d, y_2d)
36     plot_decision_boundary(X_2d, y_2d, model,
37                             f'Noyau : {kernel}', ax)
38     ax.set_xlabel('Age (standardise)')
39     ax.set_ylabel('Revenu (standardise)')
40
41 plt.suptitle('Frontieres de decision SVM avec differents noyaux',
42             fontsize=14, fontweight='bold')
43 plt.tight_layout()
44 plt.savefig('svm_decision_boundaries.png', dpi=150,
45             bbox_inches='tight')
46 plt.show()

```

Listing 14.5 – Visualisation de la frontière de décision (2 features)

```

1  # Entraîner le modele RBF
2  svm_rbf = SVC(kernel='rbf', C=1.0, random_state=42)
3  svm_rbf.fit(X_train_scaled, y_train)
4
5  # Afficher les vecteurs de support
6  print(''='' * 60)
7  print(''VECTEURS DE SUPPORT'')
8  print(''='' * 60)
9  print(f''Nombre total : {len(svm_rbf.support_vectors_)}'')
10 print(f''Par classe : Churn=0 -> {svm_rbf.n_support_[0]}, ''
11       f''Churn=1 -> {svm_rbf.n_support_[1]}'')
12 print(f''\nIndices dans le jeu d'entraînement :
13       {svm_rbf.support_}'')
14 print(''\nDonnees des vecteurs de support (avant standardisation)
15       :'')
16 sv_df = X_train.iloc[svm_rbf.support_].copy()
17 sv_df['Churn'] = y_train.iloc[svm_rbf.support_].values
18 print(sv_df.to_string())

```

Listing 14.6 – Identification et affichage des vecteurs de support

```

1  # Grille de parametres
2  param_grid = {
3      'C': [0.01, 0.1, 1, 10, 100],

```

```

4     'kernel': ['linear', 'poly', 'rbf'],
5     'gamma': ['scale', 'auto', 0.01, 0.1, 1], # pour poly et rbf
6 }
7
8 grid_search = GridSearchCV(
9     SVC(random_state=42),
10    param_grid,
11    cv=5,
12    scoring='accuracy',
13    n_jobs=-1,
14    verbose=0
15 )
16 grid_search.fit(X_train_scaled, y_train)
17
18 print('Meilleurs parametres :', grid_search.best_params_)
19 print('Meilleure accuracy (CV) :',
20       f'{grid_search.best_score_:.4f}')
21
22 # ?valuer sur le test
23 best_model = grid_search.best_estimator_
24 y_pred = best_model.predict(X_test_scaled)
25 print('\nAccuracy sur le test :', f'{accuracy_score(y_test,
26                                                    y_pred):.4f}')
27 print('\nRapport de classification :')
28 print(classification_report(y_test, y_pred,
29                             target_names=['Reste', 'Churn']))

```

Listing 14.7 – GridSearchCV pour optimiser C et le noyau

```

1 from sklearn.linear_model import LogisticRegression
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.tree import DecisionTreeClassifier
4 from sklearn.ensemble import RandomForestClassifier
5
6 classifiers = {
7     'SVM (RBF)': SVC(kernel='rbf', C=1.0, random_state=42),
8     'SVM (Lineaire)': SVC(kernel='linear', C=1.0, random_state=42),
9     'Reg. Logistique': LogisticRegression(random_state=42),
10    'KNN (k=3)': KNeighborsClassifier(n_neighbors=3),
11    'Arbre de Decision': DecisionTreeClassifier(random_state=42),
12    'Random Forest': RandomForestClassifier(n_estimators=100,
13                                           random_state=42),

```

```

14 }
15
16 print(f''{'Classifieur':<22} {'Accuracy':>10}''')
17 print(''-'' * 35)
18 for name, clf in classifiers.items():
19     clf.fit(X_train_scaled, y_train)
20     y_pred = clf.predict(X_test_scaled)
21     acc = accuracy_score(y_test, y_pred)
22     print(f''{name:<22} {acc:>10.4f}''')

```

Listing 14.8 – Comparaison avec d'autres classifieurs

```

1 from sklearn.svm import SVR
2
3 # Preparer les donnees pour la regression
4 X_reg = df.drop(['Churn', 'Depense_Annuelle'], axis=1) \
5     if 'Depense_Annuelle' in df.columns else df.drop('Churn',
6     axis=1)
7
8 # Recreer avec la variable cible de regression
9 data_reg = {
10     'Age':
11         [25,34,28,45,23,52,31,41,27,38,55,29,47,22,36,43,26,50,33,30],
12     'Revenu':
13         [22,45,28,62,18,75,35,55,24,48,80,30,65,20,42,58,25,70,38,32],
14     'Anciennete': [1,5,2,8,1,12,3,7,1,6,15,2,9,1,4,7,1,10,3,2],
15     'Nb_Achats':
16         [8,30,12,55,5,70,18,45,7,35,85,14,60,6,25,48,9,65,20,15],
17     'Montant_Moyen': [45,78,52,95,38,110,62,88,42,82,120,55,98,40,72,90,48,105,
18     'Score_Satisfaction': [5,8,4,9,3,9,6,8,4,7,10,5,9,3,7,8,4,9,6,5],
19     'Depense_Annuelle': [360,2340,624,5225,190,7700,1116,3960,294,2870,
20     10200,770,5880,240,1800,4320,432,6825,1300,870]
21 }
22
23 df_reg = pd.DataFrame(data_reg)
24
25 X_r = df_reg.drop('Depense_Annuelle', axis=1)
26 y_r = df_reg['Depense_Annuelle']
27
28 X_r_train, X_r_test, y_r_train, y_r_test = train_test_split(
29     X_r, y_r, test_size=0.3, random_state=42)
30
31 scaler_r = StandardScaler()

```

```

27 X_r_train_sc = scaler_r.fit_transform(X_r_train)
28 X_r_test_sc = scaler_r.transform(X_r_test)
29
30 # Tester différents noyaux SVR
31 print(f''{'Noyau SVR':<15} {'RMSE':>10} {'R\${}^2\${}':>10}''')
32 print(''-'' * 38)
33 for kernel in ['linear', 'poly', 'rbf']:
34     svr = SVR(kernel=kernel, C=100.0, epsilon=0.1)
35     svr.fit(X_r_train_sc, y_r_train)
36     y_r_pred = svr.predict(X_r_test_sc)
37     rmse = np.sqrt(mean_squared_error(y_r_test, y_r_pred))
38     r2 = r2_score(y_r_test, y_r_pred)
39     print(f''{kernel:<15} {rmse:>10.1f} {r2:>10.4f}''')

```

Listing 14.9 – SVR pour la régression (prédire la dépense annuelle)

14.6 Diagramme récapitulatif : appliquer un SVM pas à pas

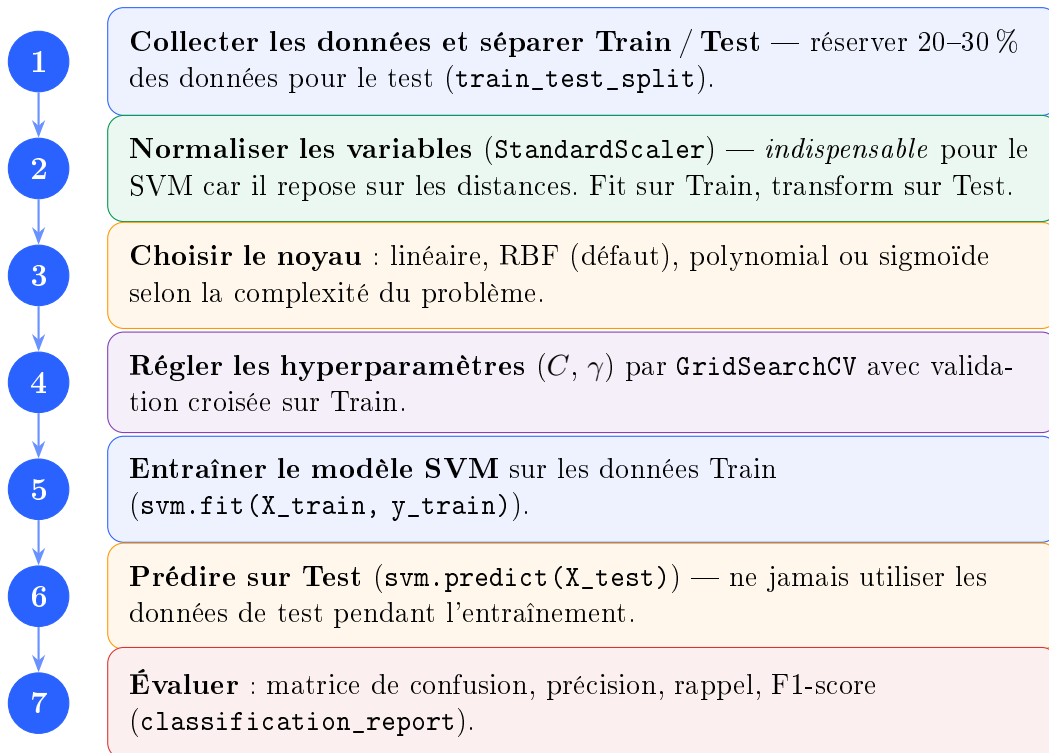


FIGURE 14.8 – Les 7 étapes pour appliquer un SVM en classification.

14.7 Schéma récapitulatif : phases d'entraînement et de test

PHASE D'ENTRAÎNEMENT (Training)

1. **Étape 1 — Charger et séparer les données :**

```
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2)
```

2. **Étape 2 — Normaliser les features** (indispensable pour SVM) :

$$\text{Formule : } x'_j = \frac{x_j - \mu_j}{\sigma_j}$$

```
scaler = StandardScaler()
```

```
X_train_s = scaler.fit_transform(X_train)
```

3. **Étape 3 — Choisir le noyau :**

Formules :

— Linéaire : $K(\mathbf{x}, \mathbf{x}') = \mathbf{x} \cdot \mathbf{x}'$

— RBF : $K(\mathbf{x}, \mathbf{x}') = e^{-\gamma \|\mathbf{x} - \mathbf{x}'\|^2}$

— Polynomial : $K(\mathbf{x}, \mathbf{x}') = (\mathbf{x} \cdot \mathbf{x}' + r)^d$

```
from sklearn.svm import SVC
```

```
svm = SVC(kernel='rbf')
```

4. **Étape 4 — Régler C et γ par GridSearchCV :**

Formule : C = tolérance aux erreurs (grand C = marge étroite), γ = portée du noyau RBF.

```
params = {'C': [0.1, 1, 10], 'gamma': ['scale', 0.1, 1]}
```

```
grid = GridSearchCV(SVC(kernel='rbf'), params, cv=5)
```

```
grid.fit(X_train_s, y_train)
```

5. **Étape 5 — Entraîner et trouver l'hyperplan optimal :**

Formule : Maximiser la marge : $\max \frac{2}{\|\mathbf{w}\|}$ sous contrainte $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$.

```
best_svm = grid.best_estimator_
```

```
print("Vecteurs support :", best_svm.n_support_)
```



Modèle entraîné transmis à la phase de test

PHASE DE TEST (Testing)

1. **Étape 1 — Normaliser X_{test} avec le MÊME scaler :**

```
X_test_s = scaler.transform(X_test)
```

2. **Étape 2 — Prédire :**

Formule : $\hat{y} = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$ ou $\hat{y} = \text{sign} \left(\sum_i \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b \right)$

```
y_pred = best_svm.predict(X_test_s)
```

3. **Étape 3 — Évaluer :**

```
print(classification_report(y_test, y_pred))
```

```
print(confusion_matrix(y_test, y_pred))
```

14.8 Exercices corrigés**Exercice 14.1 — Trouver l'hyperplan séparateur**

Considérons 4 points en 2D :

Point	x_1	x_2	y (classe)
A	1	1	+1
B	2	2	+1
C	4	4	-1
D	5	5	-1

1. Représentez ces points sur un graphe.
2. Trouvez l'hyperplan séparateur $w_1x_1 + w_2x_2 + b = 0$ qui maximise la marge.
3. Identifiez les vecteurs de support.
4. Calculez la marge.

Correction de l'exercice 14.1

1. Les quatre points sont alignés sur la droite $x_2 = x_1$. Les classes +1 (en bas à gauche) et -1 (en haut à droite) sont clairement séparables.

2. Par symétrie, l'hyperplan optimal est perpendiculaire à la direction (1, 1) et passe par le milieu entre $B(2, 2)$ et $C(4, 4)$, c'est-à-dire le point (3, 3).

L'équation de l'hyperplan est de la forme $x_1 + x_2 + b = 0$. Comme il passe par (3, 3) :

$$3 + 3 + b = 0 \implies b = -6$$

L'hyperplan est : $x_1 + x_2 - 6 = 0$, soit $\mathbf{w} = (1, 1)$ et $b = -6$.

Vérifions les classes :

$$— A(1, 1) : 1 + 1 - 6 = -4 < 0 \quad (\text{classe } -1 \text{ ??})$$

Il faut ajuster le signe : prenons $\mathbf{w} = (-1, -1)$ et $b = 6$, c'est-à-dire $-x_1 - x_2 + 6 = 0$.

$$— A(1, 1) : -1 - 1 + 6 = 4 > 0 \quad \checkmark \quad (\text{classe } +1)$$

$$— B(2, 2) : -2 - 2 + 6 = 2 > 0 \quad \checkmark \quad (\text{classe } +1)$$

$$— C(4, 4) : -4 - 4 + 6 = -2 < 0 \quad \checkmark \quad (\text{classe } -1)$$

$$— D(5, 5) : -5 - 5 + 6 = -4 < 0 \quad \checkmark \quad (\text{classe } -1)$$

Pour normaliser avec les vecteurs de support satisfaisant ± 1 , le vecteur de support le plus proche côté $+1$ est $B(2, 2)$ avec une valeur de 2. On divise par 2 :

$$\mathbf{w} = \left(-\frac{1}{2}, -\frac{1}{2}\right), \quad b = 3$$

Vérification : $B(2, 2) : -1 - 1 + 3 = +1 \quad \checkmark$ $C(4, 4) : -2 - 2 + 3 = -1 \quad \checkmark$

3. Les vecteurs de support sont $B(2, 2)$ et $C(4, 4)$ — les points les plus proches de l'hyperplan.

4. La marge est :

$$\text{Marge} = \frac{2}{\|\mathbf{w}\|} = \frac{2}{\sqrt{(-1/2)^2 + (-1/2)^2}} = \frac{2}{\sqrt{1/2}} = \frac{2}{1/\sqrt{2}} = 2\sqrt{2} \approx 2.83$$

On peut vérifier : la distance entre $B(2, 2)$ et $C(4, 4)$ est $\sqrt{(4-2)^2 + (4-2)^2} = 2\sqrt{2}$, ce qui est cohérent.

Exercice 14.2 — Calculer la marge

Soit un SVM avec $\mathbf{w} = (3, 4)$ et $b = -2$.

1. Calculez $\|\mathbf{w}\|$.
2. Calculez la marge du SVM.
3. Calculez la distance du point $\mathbf{x}_0 = (1, 2)$ à l'hyperplan.
4. Ce point est-il un vecteur de support si $y_0 = +1$?

Correction de l'exercice 14.2

$$1. \|\mathbf{w}\| = \sqrt{3^2 + 4^2} = \sqrt{9 + 16} = \sqrt{25} = 5$$

$$2. \text{Marge} = \frac{2}{\|\mathbf{w}\|} = \frac{2}{5} = 0.4$$

3. Distance de $\mathbf{x}_0 = (1, 2)$ à l'hyperplan :

$$d = \frac{|\mathbf{w} \cdot \mathbf{x}_0 + b|}{\|\mathbf{w}\|} = \frac{|3 \times 1 + 4 \times 2 + (-2)|}{5} = \frac{|3 + 8 - 2|}{5} = \frac{9}{5} = 1.8$$

4. La valeur fonctionnelle est $\mathbf{w} \cdot \mathbf{x}_0 + b = 3 + 8 - 2 = 9$. Pour être un vecteur de support de classe +1, il faudrait que cette valeur soit exactement +1 (après normalisation). Ici, $9 \gg 1$, donc ce point est **loin de la marge** : il n'est **pas** un vecteur de support. Les vecteurs de support sont les points pour lesquels $|\mathbf{w} \cdot \mathbf{x} + b| = 1$ (sur la frontière de la marge).

Exercice 14.3 — Pourquoi le noyau RBF nécessite la mise à l'échelle

Le noyau RBF est défini par $K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2)$.

Considérons deux clients :

— Client A : Age = 25 ans, Revenu = 22 000 €

— Client B : Age = 30 ans, Revenu = 25 000 €

1. Calculez $\|\mathbf{x}_A - \mathbf{x}_B\|^2$ **sans** mise à l'échelle.
2. Quelle variable domine la distance ? Pourquoi est-ce un problème ?
3. Après standardisation (supposons $\mu_{\text{Age}} = 35$, $\sigma_{\text{Age}} = 10$, $\mu_{\text{Revenu}} = 45\,000$, $\sigma_{\text{Revenu}} = 20\,000$), recalculez $\|\mathbf{x}_A - \mathbf{x}_B\|^2$.
4. Commentez la différence.

Correction de l'exercice 14.3

1. Sans mise à l'échelle :

$$\|\mathbf{x}_A - \mathbf{x}_B\|^2 = (25 - 30)^2 + (22\,000 - 25\,000)^2 = 25 + 9\,000\,000 = 9\,000\,025$$

2. La variable **Revenu** domine totalement la distance (9 000 000 contre 25). L'âge est **complètement ignoré**, comme s'il n'existait pas ! Le SVM va prendre ses décisions uniquement basées sur le revenu.

3. Après standardisation :

$$\text{Age}_A^* = \frac{25 - 35}{10} = -1.0, \quad \text{Age}_B^* = \frac{30 - 35}{10} = -0.5$$

$$\text{Revenu}_A^* = \frac{22\,000 - 45\,000}{20\,000} = -1.15, \quad \text{Revenu}_B^* = \frac{25\,000 - 45\,000}{20\,000} = -1.0$$

$$\|\mathbf{x}_A^* - \mathbf{x}_B^*\|^2 = (-1.0 - (-0.5))^2 + (-1.15 - (-1.0))^2 = 0.25 + 0.0225 = 0.2725$$

4. Après standardisation, les deux variables contribuent de manière **équilibrée** à la distance. L'âge contribue à 0.25 et le revenu à 0.0225 : les deux sont du même ordre de grandeur. Le noyau RBF pourra ainsi exploiter correctement **toutes les variables**. C'est pourquoi la **mise à l'échelle est obligatoire** pour le SVM avec noyau RBF (et aussi pour les noyaux polynomial et sigmoïde).

Exercice 14.4 — SVM sur le dataset moons (Python)

Le dataset `make_moons` de scikit-learn génère deux « croissants de lune » imbriqués : un problème non linéaire classique.

1. Générez le dataset avec `make_moons(n_samples=300, noise=0.2, random_state=42)`.
2. Standardisez les features.
3. Entraînez un SVM avec chacun des noyaux : `linear`, `poly` (degré 3), `rbf`.
4. Visualisez la frontière de décision pour chaque noyau.
5. Quel noyau donne la meilleure accuracy ? Pourquoi ?
6. Faites un `GridSearchCV` sur $C \in \{0.1, 1, 10, 100\}$ et $\gamma \in \{0.01, 0.1, 1, 10\}$ pour le noyau RBF.

Correction de l'exercice 14.4

```

1 from sklearn.datasets import make_moons
2 from sklearn.svm import SVC
3 from sklearn.preprocessing import StandardScaler
4 from sklearn.model_selection import train_test_split,
   GridSearchCV
5 from sklearn.metrics import accuracy_score
6 import numpy as np
7 import matplotlib.pyplot as plt
8
9 # 1. Générer le dataset
10 X, y = make_moons(n_samples=300, noise=0.2, random_state=42)
11
12 # Visualiser
13 plt.figure(figsize=(8, 5))
14 plt.scatter(X[y==0, 0], X[y==0, 1], c='blue', label='Classe 0',
15            edgecolors='k', alpha=0.7)
16 plt.scatter(X[y==1, 0], X[y==1, 1], c='red', label='Classe 1',
```

```

17         edgecolors='k', alpha=0.7)
18 plt.title('Dataset Moons (non lineairement separable)')
19 plt.legend()
20 plt.grid(True, alpha=0.3)
21 plt.show()

```

Listing 14.10 – Génération et visualisation du dataset moons

```

1  # 2. Standardisation
2  scaler = StandardScaler()
3  X_scaled = scaler.fit_transform(X)
4
5  X_train, X_test, y_train, y_test = train_test_split(
6      X_scaled, y, test_size=0.3, random_state=42)
7
8  # 3. Entraîner les différents noyaux
9  kernels_config = {
10     'Lineaire': SVC(kernel='linear', C=1.0),
11     'Polynomial (d=3)': SVC(kernel='poly', degree=3, C=1.0),
12     'RBF': SVC(kernel='rbf', C=1.0),
13 }
14
15 results = {}
16 for name, model in kernels_config.items():
17     model.fit(X_train, y_train)
18     y_pred = model.predict(X_test)
19     acc = accuracy_score(y_test, y_pred)
20     results[name] = acc
21     print(f' {name:25s} -> Accuracy = {acc:.4f} ')

```

Listing 14.11 – Comparaison des noyaux sur le dataset moons

```

1  # 4. Visualiser les frontieres de decision
2  fig, axes = plt.subplots(1, 3, figsize=(18, 5))
3
4  for ax, (name, model) in zip(axes, kernels_config.items()):
5      h = 0.02
6      x_min, x_max = X_scaled[:, 0].min() - 0.5, X_scaled[:,
7          0].max() + 0.5
8      y_min, y_max = X_scaled[:, 1].min() - 0.5, X_scaled[:,

```

```

1] .max() + 0.5
8   xx, yy = np.meshgrid(np.arange(x_min, x_max, h),
9                           np.arange(y_min, y_max, h))
10  Z = model.predict(np.c_[xx.ravel(),
11                          yy.ravel()]).reshape(xx.shape)
12
13  ax.contourf(xx, yy, Z, alpha=0.3, cmap='RdBu')
14  ax.scatter(X_scaled[:, 0], X_scaled[:, 1], c=y, cmap='RdBu',
15              edgecolors='k', s=30)
16
17  # Vecteurs de support
18  sv = model.support_vectors_
19  ax.scatter(sv[:, 0], sv[:, 1], s=120, facecolors='none',
20              edgecolors='gold', linewidths=2)
21  ax.set_title(f'{name}\nAccuracy = {results[name]:.2f}',
22              fontweight='bold')
23  ax.set_xlabel('$x_1$')
24  ax.set_ylabel('$x_2$')
25
26  plt.suptitle('Comparaison des noyaux SVM sur le dataset Moons',
27              fontsize=14, fontweight='bold', y=1.02)
28  plt.tight_layout()
29  plt.savefig('svm_moons_comparison.png', dpi=150,
30              bbox_inches='tight')
31  plt.show()

```

Listing 14.12 – Visualisation des frontières de décision

5. Analyse des résultats. Le noyau **linéaire** obtient une accuracy faible (~ 0.87) car les données « moons » ne sont **pas linéairement séparables** : une droite ne peut pas séparer deux croissants imbriqués. Le noyau **polynomial** (degré 3) fait mieux (~ 0.96) car il peut modéliser des frontières courbes. Le noyau **RBF** donne les meilleurs résultats (~ 0.98) car il peut modéliser n'importe quelle frontière non linéaire grâce à sa projection dans un espace de dimension infinie.

```

1  # 6. GridSearchCV
2  param_grid = {
3      'C': [0.1, 1, 10, 100],
4      'gamma': [0.01, 0.1, 1, 10],
5  }
6

```

```

7  grid = GridSearchCV(SVC(kernel='rbf'), param_grid, cv=5,
8                      scoring='accuracy', n_jobs=-1)
9  grid.fit(X_train, y_train)
10
11 print('Meilleurs parametres :', grid.best_params_)
12 print('Meilleure accuracy CV :', f'{grid.best_score_:.4f}')
13
14 y_pred_best = grid.best_estimator_.predict(X_test)
15 print('Accuracy test :', f'{accuracy_score(y_test,
16      y_pred_best):.4f}')
17
18 # Afficher les resultats sous forme de heatmap
19 import pandas as pd
20 results_df = pd.DataFrame(grid.cv_results_)
21 scores = results_df.pivot_table(
22     values='mean_test_score',
23     index='param_C',
24     columns='param_gamma')
25
26 plt.figure(figsize=(8, 5))
27 plt.imshow(scores, cmap='YlOrRd', aspect='auto')
28 plt.colorbar(label='Accuracy moyenne (CV)')
29 plt.xlabel('gamma')
30 plt.ylabel('C')
31 plt.xticks(range(len(scores.columns)), scores.columns)
32 plt.yticks(range(len(scores.index)), scores.index)
33 for i in range(len(scores.index)):
34     for j in range(len(scores.columns)):
35         plt.text(j, i, f'{scores.iloc[i, j]:.3f}',
36                 ha='center', va='center', fontweight='bold')
37 plt.title('Heatmap GridSearchCV : Accuracy en fonction de C et
38         gamma')
39 plt.tight_layout()
40 plt.savefig('svm_gridsearch_heatmap.png', dpi=150,
41             bbox_inches='tight')
42 plt.show()

```

Listing 14.13 – GridSearchCV pour le noyau RBF

Résumé du chapitre

À retenir — Machines à Vecteurs de Support

1. Le SVM cherche l'hyperplan qui **maximise la marge** entre les classes.
2. Les **vecteurs de support** sont les points les plus proches de la frontière — seuls eux déterminent l'hyperplan.
3. La **marge souple** (paramètre C) permet de gérer les données non parfaitement séparables.
4. L'**astuce du noyau** permet de séparer des données non linéaires en les projetant dans un espace de dimension supérieure.
5. Le noyau **RBF** est le choix par défaut — il fonctionne bien dans la plupart des cas.
6. La **mise à l'échelle** des features est **obligatoire** avant d'utiliser un SVM.
7. Les hyperparamètres clés sont C (régularisation), γ (portée du noyau RBF) et le **type de noyau**.
8. Le **SVR** adapte la même logique à la régression avec un tube d'insensibilité ε .

Chapitre 15

Classification Naïve Bayes

15.1 Hands-On : détecter les emails spam

Hands-On : Votre boîte mail est envahie de spams !

Vous travaillez dans une entreprise et votre responsable vous demande de créer un **filtre anti-spam automatique**. Il vous fournit un historique de 8 emails déjà classés par des humains comme **Spam** ou **Ham** (légitime). Pour chaque email, on a noté la présence ou l'absence de 4 mots-clés : **free**, **urgent**, **meeting** et **money**.

Email	free	urgent	meeting	money	Classe
1	Oui	Oui	Non	Oui	Spam
2	Oui	Non	Non	Oui	Spam
3	Non	Non	Oui	Non	Ham
4	Oui	Oui	Non	Non	Spam
5	Non	Non	Oui	Non	Ham
6	Non	Non	Oui	Non	Ham
7	Oui	Non	Non	Non	Ham
8	Non	Oui	Non	Oui	Spam

TABLE 15.1 – Historique de 8 emails avec présence de mots-clés et classe.

Un nouvel email arrive. Il contient les mots « **free** » et « **money** », mais pas « urgent » ni « meeting ». **Est-ce un Spam ou un Ham ?**

Nouvel email à classifier

free = **Oui** urgent = **Non** meeting = **Non** money = **Oui**

Spam ou **Ham** ?

Nous allons répondre à cette question grâce au classificateur **Naïve Bayes**, qui utilise les probabilités pour prendre sa décision. Mais d'abord, comprenons l'intuition.

15.2 Intuition : un détective qui rassemble des indices

Analogie : le détective probabiliste

Imaginez un détective qui enquête sur un vol. Il a deux suspects : **Fatima** et **Ahmed**. Au départ, sans aucune preuve, il pense que les deux ont autant de chances d'être coupables : 50% chacun. Ce sont ses **croyances initiales** (on dit *prior* en anglais).

1. **Indice 1** : Des empreintes trouvées sur la fenêtre. L'expérience montre que 80% des

voleurs laissent des empreintes, mais seulement 10% des innocents. \Rightarrow Cet indice augmente la probabilité d’Fatima (si ses empreintes correspondent).

2. **Indice 2** : Un témoin a vu quelqu’un de grand. Fatima est grande, Ahmed est petit. \Rightarrow La probabilité d’Fatima augmente encore.

3. **Indice 3** : Le vol a eu lieu la nuit. Fatima travaille de nuit. \Rightarrow Cela réduit un peu la probabilité d’Fatima (elle avait un alibi).

À chaque nouvel indice, le détective **met à jour ses croyances**. C’est exactement ce que fait le théorème de Bayes : **mettre à jour les probabilités en fonction des preuves observées**.

Idée clé du Naïve Bayes

Le classificateur Naïve Bayes fonctionne en trois étapes :

1. **Prior** : quelle est la probabilité de base de chaque classe ? (ex : 50% des emails sont des spams)
2. **Vraisemblance** : pour chaque feature (chaque mot), quelle est la probabilité de l’observer dans chaque classe ? (ex : le mot « free » apparaît dans 80% des spams)
3. **Décision** : on combine le tout avec le théorème de Bayes et on choisit la classe la plus probable.

Pourquoi « Naïf » ?

L’algorithme est qualifié de **naïf** car il suppose que tous les indices (features) sont **indépendants** les uns des autres. Dans notre exemple, cela signifie que la présence du mot « free » n’a aucune influence sur la présence du mot « money ».

En réalité, c’est rarement vrai ! Les mots « free » et « money » apparaissent souvent ensemble dans les spams. Mais même avec cette hypothèse simpliste, Naïve Bayes fonctionne **étonnamment bien** en pratique. Pourquoi ? Parce qu’on n’a pas besoin des probabilités exactes — on a juste besoin de savoir **quelle classe a la probabilité la plus élevée**.

15.3 Dérivation mathématique complète

15.3.1 Rappels de probabilités (depuis zéro)

Avant d’attaquer le théorème de Bayes, assurons-nous de bien comprendre les bases des probabilités.

Probabilité d'un événement

La **probabilité** d'un événement A , notée $P(A)$, mesure la « chance » que cet événement se produise. Elle se calcule par :

$$P(A) = \frac{\text{nombre de cas favorables}}{\text{nombre total de cas}}$$

La probabilité est toujours comprise entre 0 et 1 :

- $P(A) = 0$: l'événement est **impossible**.
- $P(A) = 1$: l'événement est **certain**.
- $P(A) = 0,3$: l'événement a **30% de chances** de se produire.

Exemples simples

- **Dé à 6 faces** : $P(\text{obtenir un 4}) = \frac{1}{6} \approx 0,167$ (16,7%).
- **Météo** : $P(\text{pluie demain}) = 0,3$ signifie qu'il y a 30% de chances qu'il pleuve.
- **Spam** : dans nos 8 emails, il y a 4 spams, donc $P(\text{Spam}) = \frac{4}{8} = 0,5$.

Probabilité conditionnelle

La **probabilité conditionnelle** $P(A | B)$ est la probabilité que A se produise **sachant que** B s'est déjà produit :

$$P(A | B) = \frac{P(A \text{ et } B)}{P(B)}$$

Lecture : « probabilité de A sachant B » ou « probabilité de A étant donné B ».

Probabilité conditionnelle sur nos emails

Quelle est la probabilité qu'un email soit un Spam **sachant** qu'il contient le mot « free » ?

- Parmi nos 8 emails, **4** contiennent « free » (emails 1, 2, 4, 7).
- Parmi ces 4, **3** sont des Spams (emails 1, 2, 4). L'email 7 est un Ham.

$$P(\text{Spam} | \text{free}=\text{Oui}) = \frac{3}{4} = 0,75$$

Donc, si un email contient « free », il y a **75% de chances** que ce soit un spam.

Probabilité conjointe

La **probabilité conjointe** $P(A \text{ et } B)$ est la probabilité que A et B se produisent **en même temps**. On peut la calculer avec :

$$P(A \text{ et } B) = P(A | B) \times P(B)$$

Probabilité conjointe

Quelle est la probabilité qu'un email soit un Spam **et** contienne « free » ?

— Parmi les 8 emails, 3 sont des Spams qui contiennent « free » (emails 1, 2, 4).

$$P(\text{Spam et free=Oui}) = \frac{3}{8} = 0,375$$

Vérifions avec la formule : $P(\text{Spam} | \text{free=Oui}) \times P(\text{free=Oui}) = \frac{3}{4} \times \frac{4}{8} = \frac{3}{4} \times \frac{1}{2} = \frac{3}{8} \checkmark$

15.3.2 Le théorème de Bayes : dérivation pas à pas

Nous avons maintenant tous les outils pour dériver le célèbre théorème de Bayes.

Dérivation du théorème de Bayes

On part de la probabilité conjointe, qui peut s'écrire de **deux façons** :

$$P(A \text{ et } B) = P(A | B) \times P(B) \quad (15.1)$$

$$P(A \text{ et } B) = P(B | A) \times P(A) \quad (15.2)$$

Les deux expressions sont égales (c'est la même probabilité conjointe), donc :

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

On divise les deux côtés par $P(B)$:

$$P(A | B) = \frac{P(B | A) \times P(A)}{P(B)}$$

C'est le **théorème de Bayes** !

Vocabulaire du théorème de Bayes

Dans le contexte du Machine Learning, en posant A = classe et B = features :

$$P(\text{classe} | \text{features}) = \frac{P(\text{features} | \text{classe}) \times P(\text{classe})}{P(\text{features})}$$

Chaque terme a un nom :

Terme	Nom	Signification
$P(\text{classe} \mid \text{features})$	Postérieur (posterior)	Ce qu'on cherche !
$P(\text{features} \mid \text{classe})$	Vraisemblance (likelihood)	« Si c'est un spam, quelle proba d'observer ces mots ? »
$P(\text{classe})$	Prior (a priori)	Probabilité de base de la classe
$P(\text{features})$	Évidence (marginal)	Probabilité d'observer ces features

Application numérique complète

Calculons $P(\text{Spam} \mid \text{free}=\text{Oui})$ avec Bayes :

- $P(\text{Spam}) = \frac{4}{8} = 0,5$ (prior : 50% des emails sont des spams)
- $P(\text{free}=\text{Oui} \mid \text{Spam}) = \frac{3}{4} = 0,75$ (vraisemblance : 3 spams sur 4 contiennent « free »)
- $P(\text{free}=\text{Oui}) = \frac{4}{8} = 0,5$ (évidence : 4 emails sur 8 contiennent « free »)

Application du théorème :

$$P(\text{Spam} \mid \text{free}=\text{Oui}) = \frac{0,75 \times 0,5}{0,5} = \frac{0,375}{0,5} = 0,75$$

On retrouve bien le résultat que nous avons calculé directement !

15.3.3 L'hypothèse « naïve » d'indépendance

En pratique, nous avons **plusieurs features** (plusieurs mots). Comment gérer cela ?

L'hypothèse naïve d'indépendance conditionnelle

On suppose que les features x_1, x_2, \dots, x_p sont **indépendantes les unes des autres** étant donné la classe c . Cela permet de décomposer la vraisemblance :

$$P(x_1, x_2, \dots, x_p \mid c) = P(x_1 \mid c) \times P(x_2 \mid c) \times \dots \times P(x_p \mid c) = \prod_{j=1}^p P(x_j \mid c)$$

Sans cette hypothèse, il faudrait estimer $P(x_1, x_2, \dots, x_p \mid c)$ directement, ce qui nécessiterait un nombre astronomique d'exemples. Avec p features binaires, il y aurait 2^p combinaisons possibles !

Avec l'hypothèse, il suffit d'estimer $P(x_j \mid c)$ pour chaque feature j séparément, ce qui est beaucoup plus simple.

Pourquoi ça marche quand même ?

Même si les features ne sont pas vraiment indépendantes, Naïve Bayes fonctionne bien pour la **classification** car :

- On ne cherche pas les probabilités exactes, mais seulement **quelle classe a la plus grande probabilité**.
- Même si $P(\text{Spam} \mid \text{features})$ est estimé à 0,82 au lieu du vrai 0,91, l'important est qu'il reste **supérieur** à $P(\text{Ham} \mid \text{features})$.
- L'**ordre** (le classement) des probabilités est souvent correct, même si les valeurs exactes ne le sont pas.

15.3.4 La règle de classification

Règle de décision de Naïve Bayes

On prédit la classe \hat{c} qui **maximise le postérieur** :

$$\hat{c} = \arg \max_c P(c) \times \prod_{j=1}^p P(x_j \mid c)$$

Remarque importante : le dénominateur $P(\text{features})$ est le **même** pour toutes les classes. Il ne change donc pas le classement et on peut l'ignorer ! On compare simplement :

$$P(c) \times \prod_{j=1}^p P(x_j \mid c) \quad \text{pour chaque classe } c$$

Problème numérique : le passage au logarithme

Quand on multiplie beaucoup de probabilités (des nombres entre 0 et 1), le produit devient **extrêmement petit** et l'ordinateur peut l'arrondir à zéro (« underflow numérique »). Solution : on utilise le **logarithme** !

Puisque log est une fonction **croissante**, maximiser un produit revient à maximiser la somme des logarithmes :

$$\hat{c} = \arg \max_c \left[\log P(c) + \sum_{j=1}^p \log P(x_j | c) \right]$$

Le produit de probabilités devient une **somme de logarithmes**, beaucoup plus stable numériquement.

15.3.5 Les variantes de Naïve Bayes

Selon le **type de features**, on utilise une variante différente du modèle pour estimer $P(x_j | c)$.

Naïve Bayes de Bernoulli (features binaires)

Quand les features sont **binaires** (0 ou 1, Oui ou Non), comme dans notre exemple d'emails :

$$P(x_j = 1 | c) = \frac{\text{nombre d'exemples de classe } c \text{ où } x_j = 1}{\text{nombre total d'exemples de classe } c}$$

C'est le cas le plus simple. C'est exactement ce que nous utilisons pour notre problème de spam.

Naïve Bayes multinomial (comptages de mots)

Quand les features sont des **comptages** (combien de fois un mot apparaît), on utilise la version multinomiale :

$$P(x_j | c) = \frac{\text{nombre d'occurrences du mot } j \text{ dans la classe } c}{\text{nombre total de mots dans la classe } c}$$

C'est la variante la plus utilisée pour la **classification de textes** (spam, analyse de sentiments, catégorisation d'articles).

Naïve Bayes gaussien (features continues)

Quand les features sont **continues** (comme l'âge, le revenu), on suppose que chaque feature suit une **distribution normale** (gaussienne) au sein de chaque classe :

Formule gaussienne

$$P(x_j | c) = \frac{1}{\sqrt{2\pi\sigma_{j,c}^2}} \exp\left(-\frac{(x_j - \mu_{j,c})^2}{2\sigma_{j,c}^2}\right)$$

Où :

- $\mu_{j,c}$ = la **moyenne** de la feature j pour les exemples de classe c .
- $\sigma_{j,c}$ = l'**écart-type** de la feature j pour les exemples de classe c .

La courbe en cloche

La distribution normale (ou gaussienne) est la fameuse « courbe en cloche ». La plupart des valeurs sont proches de la moyenne μ , et plus on s'en éloigne, plus c'est rare.

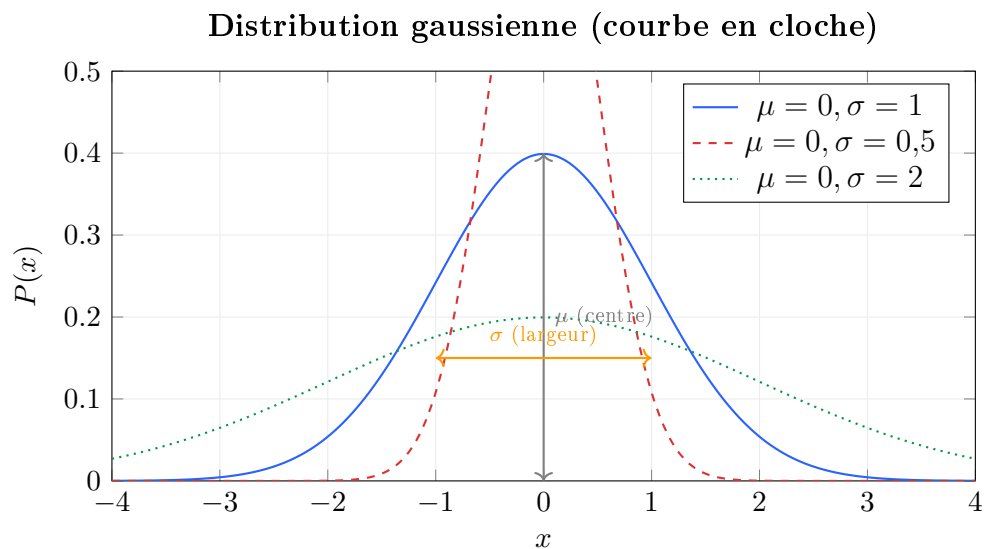


FIGURE 15.1 – La distribution gaussienne. Le paramètre μ détermine le centre, et σ la largeur de la courbe. Plus σ est petit, plus la courbe est étroite et haute.

15.3.6 Le lissage de Laplace

Le problème de la probabilité nulle

Imaginez qu'un mot (par exemple « meeting ») n'apparaisse **jamais** dans les spams. Alors :

$$P(\text{meeting}=\text{Oui} \mid \text{Spam}) = \frac{0}{4} = 0$$

Comme on **multiplie** toutes les vraisemblances, un seul zéro annule tout le produit !

$$P(\text{Spam}) \times P(\text{free} \mid \text{Spam}) \times \underbrace{P(\text{meeting} \mid \text{Spam})}_{=0} \times \dots = 0$$

Même si tous les autres indices pointent fortement vers Spam, un seul zéro détruit le résultat.

Lissage de Laplace (add-one smoothing)

Pour éviter les probabilités nulles, on ajoute **1** à tous les comptages :

$$P(x_j = v \mid c) = \frac{\text{count}(x_j = v, c) + 1}{\text{count}(c) + K}$$

Où K est le nombre de **valeurs possibles** de la feature x_j (pour une feature binaire, $K = 2$).

Intuition : on fait comme si on avait vu chaque combinaison **au moins une fois**. Cela garantit qu'aucune probabilité ne vaut zéro.

Lissage sur notre dataset spam

Sans lissage : $P(\text{meeting}=\text{Oui} \mid \text{Spam}) = \frac{0}{4} = 0$

Avec lissage de Laplace ($K = 2$ car feature binaire) :

$$P(\text{meeting}=\text{Oui} \mid \text{Spam}) = \frac{0 + 1}{4 + 2} = \frac{1}{6} \approx 0,167$$

La probabilité n'est plus nulle, mais reste faible, ce qui est cohérent.

15.4 Application complète sur le dataset spam

Classifions maintenant le nouvel email (free=Oui, urgent=Non, meeting=Non, money=Oui) en appliquant Naïve Bayes **pas à pas** avec le lissage de Laplace.

15.4.1 Étape 1 : Calculer les priors

$$P(\text{Spam}) = \frac{4}{8} = 0,5$$

$$P(\text{Ham}) = \frac{4}{8} = 0,5$$

15.4.2 Étape 2 : Calculer les vraisemblances avec lissage de Laplace

Pour chaque feature et chaque classe, on calcule $P(x_j \mid c)$ avec le lissage ($K = 2$).

Classe Spam (4 exemples : emails 1, 2, 4, 8) :

Feature	Comptage (Oui)	$P(\text{Oui} \mid \text{Spam})$	$P(\text{Non} \mid \text{Spam})$
free	3/4	$(3 + 1)/(4 + 2) = 4/6$	$(1 + 1)/(4 + 2) = 2/6$
urgent	3/4	$(3 + 1)/(4 + 2) = 4/6$	$(1 + 1)/(4 + 2) = 2/6$
meeting	0/4	$(0 + 1)/(4 + 2) = 1/6$	$(4 + 1)/(4 + 2) = 5/6$
money	3/4	$(3 + 1)/(4 + 2) = 4/6$	$(1 + 1)/(4 + 2) = 2/6$

Classe Ham (4 exemples : emails 3, 5, 6, 7) :

Feature	Comptage (Oui)	$P(\text{Oui} \mid \text{Ham})$	$P(\text{Non} \mid \text{Ham})$
free	1/4	$(1 + 1)/(4 + 2) = 2/6$	$(3 + 1)/(4 + 2) = 4/6$
urgent	0/4	$(0 + 1)/(4 + 2) = 1/6$	$(4 + 1)/(4 + 2) = 5/6$
meeting	3/4	$(3 + 1)/(4 + 2) = 4/6$	$(1 + 1)/(4 + 2) = 2/6$
money	0/4	$(0 + 1)/(4 + 2) = 1/6$	$(4 + 1)/(4 + 2) = 5/6$

15.4.3 Étape 3 : Calculer le score pour chaque classe

Pour le nouvel email : free=Oui, urgent=Non, meeting=Non, money=Oui.

Score Spam :

$$\begin{aligned}
 \text{Score}(\text{Spam}) &= P(\text{Spam}) \times P(\text{free=Oui} \mid \text{Spam}) \times P(\text{urgent=Non} \mid \text{Spam}) \\
 &\quad \times P(\text{meeting=Non} \mid \text{Spam}) \times P(\text{money=Oui} \mid \text{Spam}) \\
 &= 0,5 \times \frac{4}{6} \times \frac{2}{6} \times \frac{5}{6} \times \frac{4}{6} \\
 &= 0,5 \times 0,667 \times 0,333 \times 0,833 \times 0,667 \\
 &\approx 0,0617
 \end{aligned}$$

Score Ham :

$$\begin{aligned}
 \text{Score}(\text{Ham}) &= P(\text{Ham}) \times P(\text{free=Oui} \mid \text{Ham}) \times P(\text{urgent=Non} \mid \text{Ham}) \\
 &\quad \times P(\text{meeting=Non} \mid \text{Ham}) \times P(\text{money=Oui} \mid \text{Ham}) \\
 &= 0,5 \times \frac{2}{6} \times \frac{5}{6} \times \frac{2}{6} \times \frac{1}{6} \\
 &= 0,5 \times 0,333 \times 0,833 \times 0,333 \times 0,167 \\
 &\approx 0,00771
 \end{aligned}$$

15.4.4 Étape 4 : Normaliser et décider

Pour obtenir des probabilités qui somment à 1, on normalise :

$$P(\text{Spam} \mid \text{email}) = \frac{0,0617}{0,0617 + 0,00771} = \frac{0,0617}{0,0694} \approx \mathbf{0,889}$$

$$P(\text{Ham} \mid \text{email}) = \frac{0,00771}{0,0617 + 0,00771} = \frac{0,00771}{0,0694} \approx \mathbf{0,111}$$

Décision : SPAM

Le nouvel email (free=Oui, money=Oui) est classé comme **Spam** avec une probabilité de **88,9%**. Naïve Bayes est très confiant !

15.5 Application sur le dataset client (Gaussian Naïve Bayes)

Appliquons maintenant Naïve Bayes sur notre fil rouge pour prédire le **Churn**. Puisque nos features sont continues (Age, Revenu, etc.), nous utilisons la variante **gaussienne**.

15.5.1 Les données

Pour simplifier le calcul à la main, nous utilisons uniquement deux features : **Revenu** et **Ancienneté**. Voici les statistiques par classe :

Feature	Churn = 1 (Partis)		Churn = 0 (Fidèles)	
	μ	σ	μ	σ
Revenu (k€)	24,88	4,55	55,58	14,25
Ancienneté (ans)	1,50	0,50	6,92	3,63

TABLE 15.2 – Moyennes et écarts-types par classe pour deux features.

Les priors :

$$P(\text{Churn}=1) = \frac{8}{20} = 0,4 \quad P(\text{Churn}=0) = \frac{12}{20} = 0,6$$

15.5.2 Classons un nouveau client

Un nouveau client a : **Revenu** = **28 k€** et **Ancienneté** = **2 ans**. Va-t-il cherner ?

Étape 1 : $P(\text{Revenu} = 28 \mid \text{Churn}=1)$

On applique la formule gaussienne avec $\mu = 24,88$ et $\sigma = 4,55$:

$$\begin{aligned} P(28 \mid \text{Churn}=1) &= \frac{1}{\sqrt{2\pi \times 4,55^2}} \times \exp\left(-\frac{(28 - 24,88)^2}{2 \times 4,55^2}\right) \\ &= \frac{1}{\sqrt{2\pi \times 20,70}} \times \exp\left(-\frac{(3,12)^2}{41,41}\right) \\ &= \frac{1}{11,41} \times \exp(-0,235) \\ &= 0,0877 \times 0,791 \\ &\approx 0,0693 \end{aligned}$$

Étape 2 : $P(\text{Revenu} = 28 \mid \text{Churn}=0)$

Avec $\mu = 55,58$ et $\sigma = 14,25$:

$$\begin{aligned} P(28 \mid \text{Churn}=0) &= \frac{1}{\sqrt{2\pi \times 14,25^2}} \times \exp\left(-\frac{(28 - 55,58)^2}{2 \times 14,25^2}\right) \\ &= \frac{1}{\sqrt{2\pi \times 203,06}} \times \exp\left(-\frac{(-27,58)^2}{406,13}\right) \\ &= \frac{1}{35,72} \times \exp(-1,872) \\ &= 0,0280 \times 0,154 \\ &\approx 0,00431 \end{aligned}$$

Étape 3 : $P(\text{Ancienneté} = 2 \mid \text{Churn}=1)$

Avec $\mu = 1,50$ et $\sigma = 0,50$:

$$\begin{aligned} P(2 \mid \text{Churn}=1) &= \frac{1}{\sqrt{2\pi \times 0,50^2}} \times \exp\left(-\frac{(2 - 1,50)^2}{2 \times 0,50^2}\right) \\ &= \frac{1}{1,253} \times \exp\left(-\frac{0,25}{0,50}\right) \\ &= 0,798 \times \exp(-0,5) \\ &= 0,798 \times 0,607 \\ &\approx 0,484 \end{aligned}$$

Étape 4 : $P(\text{Ancienneté} = 2 \mid \text{Churn}=0)$

Avec $\mu = 6,92$ et $\sigma = 3,63$:

$$\begin{aligned}
 P(2 \mid \text{Churn}=0) &= \frac{1}{\sqrt{2\pi} \times 3,63^2} \times \exp\left(-\frac{(2 - 6,92)^2}{2 \times 3,63^2}\right) \\
 &= \frac{1}{9,10} \times \exp\left(-\frac{24,21}{26,35}\right) \\
 &= 0,110 \times \exp(-0,919) \\
 &= 0,110 \times 0,399 \\
 &\approx 0,0439
 \end{aligned}$$

Étape 5 : Combiner les résultats

$$\begin{aligned}
 \text{Score}(\text{Churn}=1) &= P(\text{Churn}=1) \times P(28 \mid \text{Churn}=1) \times P(2 \mid \text{Churn}=1) \\
 &= 0,4 \times 0,0693 \times 0,484 \\
 &\approx 0,01341
 \end{aligned}$$

$$\begin{aligned}
 \text{Score}(\text{Churn}=0) &= P(\text{Churn}=0) \times P(28 \mid \text{Churn}=0) \times P(2 \mid \text{Churn}=0) \\
 &= 0,6 \times 0,00431 \times 0,0439 \\
 &\approx 0,0001135
 \end{aligned}$$

Normalisation :

$$\begin{aligned}
 P(\text{Churn}=1 \mid \text{client}) &= \frac{0,01341}{0,01341 + 0,0001135} \approx \mathbf{0,992} \\
 P(\text{Churn}=0 \mid \text{client}) &= \frac{0,0001135}{0,01341 + 0,0001135} \approx \mathbf{0,008}
 \end{aligned}$$

Décision : Churn (le client va partir)

Le nouveau client (Revenu = 28 k€, Ancienneté = 2 ans) est classé comme **Churn** avec une probabilité de **99,2%**. Son profil (faible revenu, faible ancienneté) correspond parfaitement au profil des clients qui partent.

Pourquoi le résultat est si net ?

Le revenu de 28 k€ est **très proche** de la moyenne des churners (24,88 k€) mais **très éloigné** de la moyenne des fidèles (55,58 k€). De même, l'ancienneté de 2 ans est proche de la moyenne des churners (1,5 ans) mais loin des fidèles (6,92 ans). Les deux features pointent fortement dans la même direction, ce qui donne une probabilité très élevée.

15.6 Implémentation sur Google Colab

15.6.1 Naïve Bayes gaussien sur le dataset client

```

1  import numpy as np
2  import pandas as pd
3  from sklearn.naive_bayes import GaussianNB
4  from sklearn.model_selection import train_test_split
5  from sklearn.metrics import classification_report, confusion_matrix
6  from sklearn.preprocessing import StandardScaler
7
8  # Notre dataset client
9  data = {
10     'Age': [25, 34, 28, 45, 23, 52, 31, 41, 27, 38,
11            55, 29, 47, 22, 36, 43, 26, 50, 33, 30],
12     'Revenu': [22, 45, 28, 62, 18, 75, 35, 55, 24, 48,
13              80, 30, 65, 20, 42, 58, 25, 70, 38, 32],
14     'Anciennete': [1, 5, 2, 8, 1, 12, 3, 7, 1, 6,
15                  15, 2, 9, 1, 4, 7, 1, 10, 3, 2],
16     'Nb_Achats': [8, 30, 12, 55, 5, 70, 18, 45, 7, 35,
17                 85, 14, 60, 6, 25, 48, 9, 65, 20, 15],
18     'Montant_Moyen': [45, 78, 52, 95, 38, 110, 62, 88, 42, 82,
19                     120, 55, 98, 40, 72, 90, 48, 105, 65, 58],
20     'Score_Satisfaction': [5, 8, 4, 9, 3, 9, 6, 8, 4, 7,
21                          10, 5, 9, 3, 7, 8, 4, 9, 6, 5],
22     'Churn': [1, 0, 1, 0, 1, 0, 0, 0, 1, 0,
23             0, 1, 0, 1, 0, 0, 1, 0, 0, 1]
24 }
25 df = pd.DataFrame(data)
26
27 # Separer features et cible
28 X = df.drop(columns='Churn')
29 y = df['Churn']

```

```

30
31 # Diviser en train/test
32 X_train, X_test, y_train, y_test = train_test_split(
33     X, y, test_size=0.3, random_state=42, stratify=y
34 )
35
36 print(f"Taille train : {X_train.shape[0]} exemples")
37 print(f"Taille test  : {X_test.shape[0]} exemples")

```

Listing 15.1 – Naïve Bayes gaussien — classification du Churn

```

1 # Creer et entrainer le modele
2 gnb = GaussianNB()
3 gnb.fit(X_train, y_train)
4
5 # Predictions sur le test
6 y_pred = gnb.predict(X_test)
7
8 # Probabilites de chaque classe
9 y_proba = gnb.predict_proba(X_test)
10
11 print("=== Predictions sur le test ===")
12 for i in range(len(X_test)):
13     idx = X_test.index[i]
14     print(f"Client {idx+1}: predition={y_pred[i]},
15           reel={y_test.iloc[i]}, "
16           f"P(Churn=0)={y_proba[i,0]:.3f},
17           P(Churn=1)={y_proba[i,1]:.3f}")

```

Listing 15.2 – Entraînement et prédiction avec GaussianNB

```

1 # Matrice de confusion
2 print("=== Matrice de confusion ===")
3 print(confusion_matrix(y_test, y_pred))
4
5 # Rapport de classification complet
6 print("\n=== Rapport de classification ===")
7 print(classification_report(y_test, y_pred,
8                             target_names=['Fidèle', 'Churn']))
9
10 # Accuracy
11 accuracy = gnb.score(X_test, y_test)

```

```
12 print(f"Accuracy : {accuracy*100:.1f}%")
```

Listing 15.3 – Évaluation : matrice de confusion et rapport

```
1 # Nouveau client : Revenu=28, Anciennete=2, Age=26,
2 # Nb_Achats=10, Montant_Moyen=50, Score_Satisfaction=4
3 nouveau_client = pd.DataFrame({
4     'Age': [26], 'Revenu': [28], 'Anciennete': [2],
5     'Nb_Achats': [10], 'Montant_Moyen': [50],
6     'Score_Satisfaction': [4]
7 })
8
9 prediction = gnb.predict(nouveau_client)
10 probabilites = gnb.predict_proba(nouveau_client)
11
12 print(f"Prediction : {'Churn' if prediction[0] == 1 else
13       'Fidele'}")
13 print(f"P(Fidele) = {probabilites[0,0]:.4f}")
14 print(f"P(Churn) = {probabilites[0,1]:.4f}")
```

Listing 15.4 – Prédiction pour un nouveau client

15.6.2 Classification de texte avec MultinomialNB

```
1 from sklearn.naive_bayes import MultinomialNB
2 from sklearn.feature_extraction.text import CountVectorizer
3
4 # Donnees d'emails (exemple simplifie)
5 emails = [
6     "free money win prize now",
7     "urgent free offer limited time",
8     "meeting tomorrow at office",
9     "free urgent deal money",
10    "project meeting notes attached",
11    "team meeting agenda review",
12    "win free money cash prize",
13    "urgent meeting with client"
14 ]
15 labels = [1, 1, 0, 1, 0, 0, 1, 0] # 1=Spam, 0=Ham
16
17 # Convertir le texte en matrice de comptage de mots
```

```

18 vectorizer = CountVectorizer()
19 X_text = vectorizer.fit_transform(emails)
20
21 print("Vocabulaire :")
22 print(vectorizer.get_feature_names_out())
23 print(f"\nMatrice de taille : {X_text.shape}")
24 print("(8 emails, {} mots uniques)".format(X_text.shape[1]))
25
26 # Entrainer MultinomialNB
27 mnbs = MultinomialNB(alpha=1.0) # alpha=1 => lissage de Laplace
28 mnbs.fit(X_text, labels)
29
30 # Tester sur un nouvel email
31 nouvel_email = ["free money special offer"]
32 X_new = vectorizer.transform(nouvel_email)
33 pred = mnbs.predict(X_new)
34 proba = mnbs.predict_proba(X_new)
35
36 print(f"\nNouvel email : '{nouvel_email[0]}'")
37 print(f"Prediction : {'Spam' if pred[0] == 1 else 'Ham'}")
38 print(f"P(Ham) = {proba[0,0]:.4f}")
39 print(f"P(Spam) = {proba[0,1]:.4f}")

```

Listing 15.5 – Classification de texte — spam avec MultinomialNB

15.6.3 Comparaison avec d'autres classificateurs

```

1 from sklearn.linear_model import LogisticRegression
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.tree import DecisionTreeClassifier
4 from sklearn.svm import SVC
5 from sklearn.ensemble import RandomForestClassifier
6 from sklearn.model_selection import cross_val_score
7
8 # Dictionnaire des modeles
9 modeles = {
10     'Naive Bayes (Gauss.)': GaussianNB(),
11     'Reg. Logistique': LogisticRegression(max_iter=1000),
12     'KNN (k=3)': KNeighborsClassifier(n_neighbors=3),
13     'Arbre Decision': DecisionTreeClassifier(random_state=42),
14     'Random Forest': RandomForestClassifier(n_estimators=100,

```

```

15                                     random_state=42),
16     'SVM': SVC(kernel='rbf', random_state=42)
17 }
18
19 # Normaliser les features (important pour certains modeles)
20 scaler = StandardScaler()
21 X_scaled = scaler.fit_transform(X)
22
23 print("== Comparaison des modeles (validation croisee 5-fold)
24     ==")
25 print(f"{'Modele':<25} {'Accuracy moyenne':>18}
26     {'Ecart-type':>12}")
27 print("-" * 58)
28
29 resultats = {}
30 for nom, modele in modeles.items():
31     scores = cross_val_score(modele, X_scaled, y, cv=5,
32                             scoring='accuracy')
33     resultats[nom] = scores.mean()
34     print(f"{nom:<25} {scores.mean()*100:>15.1f}%"
35           f" {scores.std()*100:>10.1f}%")
36
37 # Identifier le meilleur modele
38 meilleur = max(resultats, key=resultats.get)
39 print(f"\nMeilleur modele : {meilleur} "
40       f"({resultats[meilleur]*100:.1f}%")

```

Listing 15.6 – Comparaison de Naïve Bayes avec d'autres algorithmes

Quand utiliser Naïve Bayes ?

Naïve Bayes est particulièrement recommandé quand :

- Le dataset est **petit** : Naïve Bayes a très peu de paramètres à estimer, donc il fonctionne bien même avec peu de données.
- Le problème est de la **classification de texte** : c'est le domaine de prédilection de Naïve Bayes (spam, analyse de sentiments, catégorisation).
- On a besoin de **probabilités calibrées** : Naïve Bayes fournit directement des probabilités pour chaque classe.
- L'entraînement doit être **rapide** : Naïve Bayes est l'un des algorithmes les plus rapides.

15.7 Diagramme récapitulatif : appliquer Naïve Bayes pas à pas

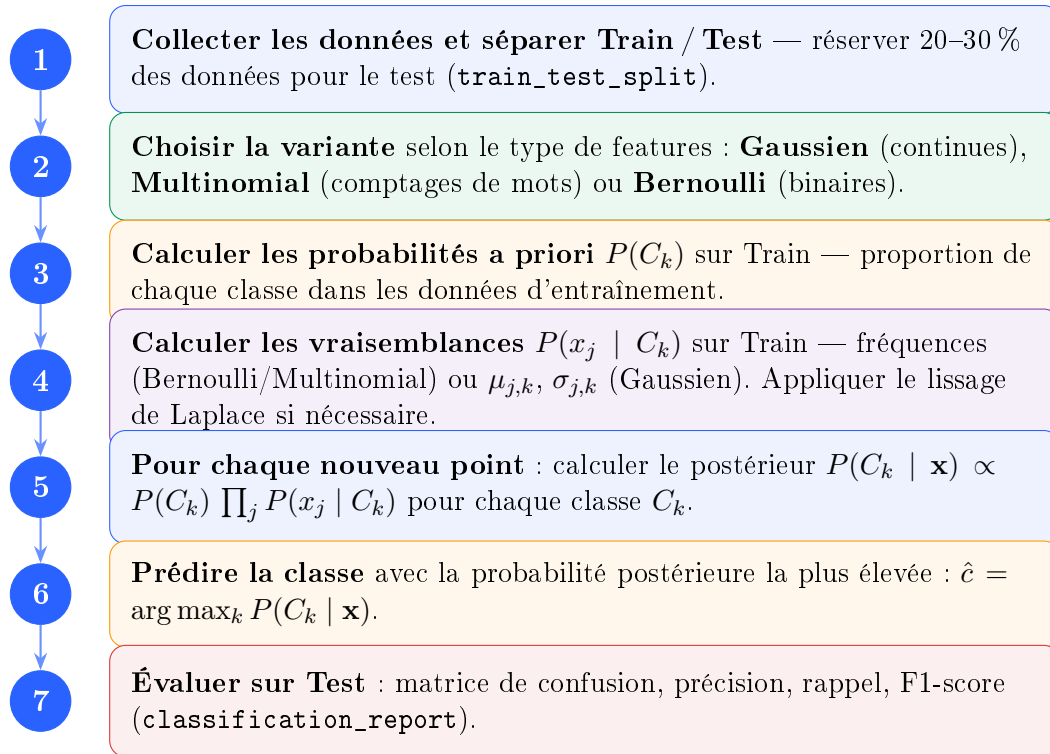


FIGURE 15.2 – Les 7 étapes pour appliquer Naïve Bayes en classification.

15.8 Schéma récapitulatif : phases d'entraînement et de test

PHASE D'ENTRAÎNEMENT (Training)

1. Étape 1 — Charger et séparer les données :

```
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2)
```

2. Étape 2 — Choisir la variante :

Formules :

— Gaussien :
$$P(x_j | C_k) = \frac{1}{\sqrt{2\pi\sigma_{jk}^2}} e^{-\frac{(x_j - \mu_{jk})^2}{2\sigma_{jk}^2}}$$

— Multinomial :
$$P(x_j | C_k) = \frac{N_{jk} + \alpha}{N_k + \alpha p} \quad (\text{avec lissage de Laplace})$$

```
from sklearn.naive_bayes import GaussianNB
nb = GaussianNB()
```

3. **Étape 3 — Calculer les probabilités a priori :**

Formule : $P(C_k) = \frac{n_k}{n}$ (proportion de chaque classe dans Train).

```
nb.fit(X_train, y_train)
```

```
print("P(C_k) :", nb.class_prior_)
```

4. **Étape 4 — Estimer les paramètres de vraisemblance :**

Formule (Gaussien) : $\hat{\mu}_{jk} = \frac{1}{n_k} \sum_{i \in C_k} x_{ij}$, $\hat{\sigma}_{jk}^2 = \frac{1}{n_k} \sum_{i \in C_k} (x_{ij} - \hat{\mu}_{jk})^2$

```
print("Moyennes :", nb.theta_)
```

```
print("Variances :", nb.var_)
```



Modèle entraîné transmis à la phase de test

PHASE DE TEST (Testing)

1. **Étape 1 — Calculer le postérieur pour chaque classe :**

Formule : $P(C_k|\mathbf{x}) \propto P(C_k) \prod_{j=1}^p P(x_j|C_k)$

```
probas = nb.predict_proba(X_test) # P(C_k|x) pour chaque classe
```

2. **Étape 2 — Prédire la classe la plus probable :**

Formule : $\hat{y} = \arg \max_k P(C_k|\mathbf{x})$

```
y_pred = nb.predict(X_test)
```

3. **Étape 3 — Évaluer :**

```
print(classification_report(y_test, y_pred))
```

15.9 Exercices

Exercice 15.1 — Le théorème de Bayes en médecine

Un test de dépistage pour une maladie rare a les caractéristiques suivantes :

- **Prévalence** (taux de base) : 1% de la population est malade, soit $P(\text{Malade}) = 0,01$.
- **Sensibilité** : si une personne est malade, le test est positif dans 99% des cas, soit $P(\text{Positif} | \text{Malade}) = 0,99$.
- **Spécificité** : si une personne est saine, le test est négatif dans 95% des cas, soit $P(\text{Négatif} | \text{Sain}) = 0,95$. Donc $P(\text{Positif} | \text{Sain}) = 0,05$.

Question : Une personne prise au hasard fait le test et obtient un résultat **positif**. Quelle est la probabilité qu'elle soit réellement malade? Autrement dit, calculez $P(\text{Malade} | \text{Positif})$.

Indice : Utilisez le théorème de Bayes. Pour le dénominateur, utilisez la loi des probabilités totales :

$$P(\text{Positif}) = P(\text{Positif} \mid \text{Malade}) \times P(\text{Malade}) + P(\text{Positif} \mid \text{Sain}) \times P(\text{Sain})$$

Correction de l'exercice 15.1

Données :

- $P(\text{Malade}) = 0,01$ et $P(\text{Sain}) = 1 - 0,01 = 0,99$
- $P(\text{Positif} \mid \text{Malade}) = 0,99$ (sensibilité)
- $P(\text{Positif} \mid \text{Sain}) = 0,05$ (faux positif = 1 - spécificité)

Étape 1 : Calculer $P(\text{Positif})$ (probabilité totale d'un test positif)

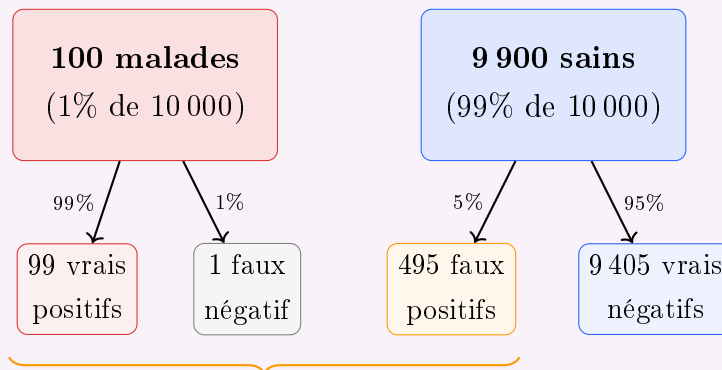
$$\begin{aligned} P(\text{Positif}) &= P(\text{Positif} \mid \text{Malade}) \times P(\text{Malade}) + P(\text{Positif} \mid \text{Sain}) \times P(\text{Sain}) \\ &= 0,99 \times 0,01 + 0,05 \times 0,99 \\ &= 0,0099 + 0,0495 \\ &= 0,0594 \end{aligned}$$

Étape 2 : Appliquer le théorème de Bayes

$$\begin{aligned} P(\text{Malade} \mid \text{Positif}) &= \frac{P(\text{Positif} \mid \text{Malade}) \times P(\text{Malade})}{P(\text{Positif})} \\ &= \frac{0,99 \times 0,01}{0,0594} \\ &= \frac{0,0099}{0,0594} \\ &\approx \mathbf{0,167 \text{ soit } 16,7\%} \end{aligned}$$

Interprétation : Même avec un test positif, il n'y a que **16,7% de chances** d'être réellement malade! Ce résultat surprenant s'explique par le **taux de base** très faible (1%). Comme la maladie est rare, la majorité des tests positifs sont des **faux positifs**.

Sur 10 000 personnes testées :



594 tests positifs au total, dont seulement 99 vrais malades

Donc : $P(\text{Malade} \mid \text{Positif}) = \frac{99}{99 + 495} = \frac{99}{594} \approx 16,7\%$.

Exercice 15.2 — Naïve Bayes complet à la main sur les spams

En utilisant le dataset de 8 emails (Tableau 15.1), classifiez l'email suivant **sans** lissage de Laplace :

Nouvel email : free = Non, urgent = Oui, meeting = Non, money = Oui

1. Calculez les priors $P(\text{Spam})$ et $P(\text{Ham})$.
2. Calculez toutes les vraisemblances nécessaires : $P(\text{free}=\text{Non} \mid \text{Spam})$, $P(\text{urgent}=\text{Oui} \mid \text{Spam})$, $P(\text{meeting}=\text{Non} \mid \text{Spam})$, $P(\text{money}=\text{Oui} \mid \text{Spam})$, et les mêmes pour Ham.
3. Calculez le score pour chaque classe.
4. Normalisez et donnez la classe prédite avec sa probabilité.
5. Refaites le calcul **avec** le lissage de Laplace et comparez.

Correction de l'exercice 15.2

Rappel du dataset :

- **Spam** (emails 1, 2, 4, 8) : free={O,O,O,N}, urgent={O,N,O,O}, meeting={N,N,N,N}, money={O,O,N,O}
- **Ham** (emails 3, 5, 6, 7) : free={N,N,N,O}, urgent={N,N,N,N}, meeting={O,O,O,N}, money={N,N,N,N}

1. Priors :

$$P(\text{Spam}) = \frac{4}{8} = 0,5 \quad P(\text{Ham}) = \frac{4}{8} = 0,5$$

2. Vraisemblances (sans lissage) :

	Spam	Ham
$P(\text{free}=\text{Non} \mid c)$	$1/4 = 0,25$	$3/4 = 0,75$
$P(\text{urgent}=\text{Oui} \mid c)$	$3/4 = 0,75$	$0/4 = 0$
$P(\text{meeting}=\text{Non} \mid c)$	$4/4 = 1,00$	$1/4 = 0,25$
$P(\text{money}=\text{Oui} \mid c)$	$3/4 = 0,75$	$0/4 = 0$

3. Scores (sans lissage) :

$$\text{Score}(\text{Spam}) = 0,5 \times 0,25 \times 0,75 \times 1,00 \times 0,75 = 0,0703$$

$$\text{Score}(\text{Ham}) = 0,5 \times 0,75 \times \underbrace{0}_{\text{urgent}} \times 0,25 \times \underbrace{0}_{\text{money}} = 0$$

4. Décision (sans lissage) :

$$P(\text{Spam} \mid \text{email}) = \frac{0,0703}{0,0703 + 0} = 1,0 \text{ et } P(\text{Ham} \mid \text{email}) = 0,0.$$

Prédiction : **Spam à 100%**. Mais ce résultat est « trop confiant » car les probabilités nulles pour Ham écrasent tout.

5. Avec lissage de Laplace ($K = 2$) :

	Spam	Ham
$P(\text{free}=\text{Non} \mid c)$	$(1 + 1)/(4 + 2) = 2/6$	$(3 + 1)/(4 + 2) = 4/6$
$P(\text{urgent}=\text{Oui} \mid c)$	$(3 + 1)/(4 + 2) = 4/6$	$(0 + 1)/(4 + 2) = 1/6$
$P(\text{meeting}=\text{Non} \mid c)$	$(4 + 1)/(4 + 2) = 5/6$	$(1 + 1)/(4 + 2) = 2/6$
$P(\text{money}=\text{Oui} \mid c)$	$(3 + 1)/(4 + 2) = 4/6$	$(0 + 1)/(4 + 2) = 1/6$

$$\text{Score}(\text{Spam}) = 0,5 \times \frac{2}{6} \times \frac{4}{6} \times \frac{5}{6} \times \frac{4}{6} = 0,5 \times \frac{160}{1296} \approx 0,0617$$

$$\text{Score}(\text{Ham}) = 0,5 \times \frac{4}{6} \times \frac{1}{6} \times \frac{2}{6} \times \frac{1}{6} = 0,5 \times \frac{8}{1296} \approx 0,00309$$

Normalisation :

$$P(\text{Spam} \mid \text{email}) = \frac{0,0617}{0,0617 + 0,00309} \approx \mathbf{0,952} \quad (\mathbf{95,2\%})$$

$$P(\text{Ham} \mid \text{email}) = \frac{0,00309}{0,0617 + 0,00309} \approx \mathbf{0,048} \quad (\mathbf{4,8\%})$$

Comparaison : Avec le lissage, la prédiction est toujours **Spam**, mais la confiance passe de 100% à 95,2%, ce qui est plus réaliste. Le lissage de Laplace empêche les probabilités extrêmes.

Exercice 15.3 — Détection de spam en Python avec MultinomialNB

Écrivez un programme Python complet qui :

1. Crée un dataset de 12 emails (6 spams et 6 hams) avec leur texte complet.
2. Utilise `CountVectorizer` pour convertir les textes en matrice de comptage.
3. Entraîne un `MultinomialNB` sur les données.
4. Teste le modèle sur 4 nouveaux emails et affiche les probabilités.
5. Affiche la matrice de confusion sur les données d'entraînement.
6. Affiche les 5 mots les plus informatifs pour chaque classe (les mots les plus « spammy » et les plus « hammy »).

Correction de l'exercice 15.3

```
import numpy as np
import pandas as pd
from sklearn.naive_bayes import MultinomialNB
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.metrics import confusion_matrix,
    classification_report

# 1. Dataset de 12 emails
emails_train = [
    "free money win prize cash now",
    "urgent offer free gift click here",
    "congratulations you won free lottery",
    "make money fast easy free income",
    "special discount offer free shipping deal",
    "urgent claim your free prize today",
```

```

    "meeting tomorrow at the office please confirm",
    "project deadline report due next week",
    "team lunch today at noon cafeteria",
    "please review the attached document draft",
    "schedule conference call with client friday",
    "quarterly budget review meeting agenda notes"
]
labels_train = [1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0]

# 2. Vectorisation (comptage de mots)
vectorizer = CountVectorizer()
X_train = vectorizer.fit_transform(emails_train)

print(f"Vocabulaire : {len(vectorizer.get_feature_names_out())}
      mots")
print(f"Matrice : {X_train.shape}")

# 3. Entrainement
mnb = MultinomialNB(alpha=1.0)
mnb.fit(X_train, labels_train)

# 4. Test sur 4 nouveaux emails
emails_test = [
    "free money special limited offer",
    "meeting agenda for next monday",
    "win cash prize urgent claim now",
    "please send the project report"
]
X_test = vectorizer.transform(emails_test)

print("\n=== Predictions sur nouveaux emails ===")
for i, email in enumerate(emails_test):
    pred = mnb.predict(X_test[i])[0]
    proba = mnb.predict_proba(X_test[i])[0]
    label = "SPAM" if pred == 1 else "HAM"
    print(f"\nEmail : '{email}'")
    print(f"    -> {label} (P(Ham)={proba[0]:.4f}, "
          f"P(Spam)={proba[1]:.4f})")

# 5. Matrice de confusion (sur les donnees d'entrainement)

```

```

y_pred_train = mnb.predict(X_train)
print("\n=== Matrice de confusion (train) ===")
print(confusion_matrix(labels_train, y_pred_train))
print("\n=== Rapport ===")
print(classification_report(labels_train, y_pred_train,
                             target_names=['Ham', 'Spam']))

# 6. Mots les plus informatifs par classe
feature_names = vectorizer.get_feature_names_out()
# Log-probabilites de chaque mot par classe
log_proba = mnb.feature_log_prob_

# Pour la classe Ham (index 0) vs Spam (index 1)
diff = log_proba[1] - log_proba[0] # positif = plus spam

# Top 5 mots les plus "spam"
top_spam_idx = np.argsort(diff)[-5:][::-1]
print("\n=== Top 5 mots les plus 'spam' ===")
for idx in top_spam_idx:
    print(f"    {feature_names[idx]:15s} "
          f"(ratio log-proba = {diff[idx]:+.3f})")

# Top 5 mots les plus "ham"
top_ham_idx = np.argsort(diff)[:5]
print("\n=== Top 5 mots les plus 'ham' ===")
for idx in top_ham_idx:
    print(f"    {feature_names[idx]:15s} "
          f"(ratio log-proba = {diff[idx]:+.3f})")

```

Résultat attendu :

- Les emails contenant « free », « money », « prize » sont classés **Spam**.
- Les emails contenant « meeting », « report », « project » sont classés **Ham**.
- Les mots les plus « spam » : free, money, prize, win, urgent.
- Les mots les plus « ham » : meeting, report, project, review, please.

Explication : MultinomialNB apprend que certains mots sont fortement associés à une classe. Le ratio des log-probabilités indique à quel point un mot est discriminant. Plus le ratio est élevé (positif), plus le mot est associé au spam ; plus il est négatif, plus il est associé au ham.

Résumé du chapitre

- Naïve Bayes est un classificateur **probabiliste** basé sur le **théorème de Bayes**.
- L'hypothèse « naïve » d'indépendance des features simplifie le calcul, et fonctionne étonnamment bien en pratique.
- Trois variantes principales : **Bernoulli** (features binaires), **Multinomial** (comptages), **Gaussien** (features continues).
- Le **lissage de Laplace** est indispensable pour éviter les probabilités nulles.
- Naïve Bayes est **rapide, simple**, et particulièrement efficace pour la **classification de textes**.
- Ses faiblesses : l'hypothèse d'indépendance peut être trop forte, et les probabilités estimées ne sont pas toujours bien calibrées.

Chapitre 16

Métriques d'Évaluation pour la Classification

Jusqu'ici, nous avons entraîné plusieurs modèles de classification : régression logistique, KNN, arbres de décision, forêts aléatoires, SVM et Naïve Bayes. Mais comment savoir lequel est **vraiment le meilleur** ? Dire qu'un modèle a « 95% de précision » suffit-il ? Comme nous allons le découvrir, la réponse est **non** — et ce chapitre va vous expliquer pourquoi.

16.1 Hands-On : le test COVID

Hands-On : Un test de dépistage COVID

Imaginez la situation suivante. Un hôpital réalise un test de dépistage COVID sur **1 000 personnes**. Voici les résultats :

- Parmi les 1 000 personnes, **50 sont réellement positives** (malades) et **950 sont réellement négatives** (saines).
- Le test identifie correctement **45 des 50 positifs** comme positifs.
- Mais le test signale également, à tort, **30 personnes saines** comme positives.

Questions :

1. Combien de personnes le test déclare-t-il positives au total ?
2. Combien de malades le test a-t-il manqués ?
3. Le test est-il fiable ? Quelles métriques utiliser pour l'évaluer ?

Ne tournez pas la page — essayez d'abord de répondre par vous-même. Nous allons calculer **toutes les métriques** de ce chapitre sur cet exemple.

Récapitulons les données sous forme de tableau :

Grandeur	Valeur
Total de personnes testées	1000
Réellement positifs (malades)	50
Réellement négatifs (sains)	950
Positifs correctement détectés (Vrais Positifs)	45
Positifs manqués (Faux Négatifs)	5
Sains déclarés positifs à tort (Faux Positifs)	30
Sains correctement identifiés (Vrais Négatifs)	920

TABLE 16.1 – Résumé du scénario du test COVID.

Tout au long de ce chapitre, nous utiliserons ces valeurs : $TP = 45$, $FP = 30$, $FN = 5$, $TN = 920$.

16.2 La matrice de confusion

Définition — Matrice de confusion

La **matrice de confusion** est un tableau 2×2 (en classification binaire) qui résume les résultats d'un classifieur en comparant les **prédictions** du modèle avec la **réalité**. Elle contient quatre cases :

- **TP** (True Positives — Vrais Positifs) : prédit Positif, réalité Positive.
- **FP** (False Positives — Faux Positifs) : prédit Positif, réalité Négative.
- **FN** (False Negatives — Faux Négatifs) : prédit Négatif, réalité Positive.
- **TN** (True Negatives — Vrais Négatifs) : prédit Négatif, réalité Négative.

16.2.1 Visualisation de la matrice de confusion

		Prédiction du modèle		
		Positif	Négatif	
Réalité	Positif	TP Vrais Positifs 45	FN Faux Négatifs 5	} 50 réellement positifs
	Négatif	FP Faux Positifs 30	TN Vrais Négatifs 920	

FIGURE 16.1 – Matrice de confusion pour le test COVID (lignes = réalité, colonnes = prédiction).

16.2.2 Comment lire la matrice de confusion ?

Convention de lecture

La convention que nous adoptons dans ce cours (et celle utilisée par `sklearn`) :

- **Les lignes** représentent la **classe réelle** (la vérité).
- **Les colonnes** représentent la **classe prédite** (ce que dit le modèle).

Certains ouvrages inversent lignes et colonnes — vérifiez toujours la convention utilisée !

16.2.3 Comprendre True/False et Positive/Negative

L'astuce pour ne jamais se tromper

Le nom de chaque case se lit en deux parties :

1. **True / False** = est-ce que la prédiction était **correcte** ?
 - **True** (Vrai) \Rightarrow le modèle avait **raison**
 - **False** (Faux) \Rightarrow le modèle s'est **trompé**
2. **Positive / Negative** = qu'est-ce que le modèle a **prédit** ?
 - **Positive** \Rightarrow le modèle a dit « oui, positif »
 - **Negative** \Rightarrow le modèle a dit « non, négatif »

Exemples avec le test COVID :

- **TP** (True Positive) : le modèle a dit « positif » et il avait **raison** \Rightarrow la personne est vraiment malade.

- **FP** (False Positive) : le modèle a dit « positif » mais il s'est **trompé** \Rightarrow la personne est saine (fausse alarme).
- **FN** (False Negative) : le modèle a dit « négatif » mais il s'est **trompé** \Rightarrow la personne est malade mais non détectée (le pire cas!).
- **TN** (True Negative) : le modèle a dit « négatif » et il avait **raison** \Rightarrow la personne est saine et correctement identifiée.

16.2.4 Vérification des totaux

Nous pouvons vérifier la cohérence de notre matrice :

$$\text{Total} = \text{TP} + \text{FP} + \text{FN} + \text{TN} = 45 + 30 + 5 + 920 = 1\,000 \checkmark$$

$$\text{Réellement positifs} = \text{TP} + \text{FN} = 45 + 5 = 50 \checkmark$$

$$\text{Réellement négatifs} = \text{FP} + \text{TN} = 30 + 920 = 950 \checkmark$$

$$\text{Prédits positifs} = \text{TP} + \text{FP} = 45 + 30 = 75$$

$$\text{Prédits négatifs} = \text{FN} + \text{TN} = 5 + 920 = 925$$

16.3 Métriques dérivées de la matrice de confusion

16.3.1 Accuracy (Exactitude)

Définition — Accuracy

L'**accuracy** (exactitude) est la proportion de prédictions correctes parmi toutes les prédictions :

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

Calcul sur l'exemple COVID

$$\text{Accuracy} = \frac{45 + 920}{45 + 920 + 30 + 5} = \frac{965}{1\,000} = 0,965 = \mathbf{96,5\%}$$

Le test est correct dans 96,5% des cas. Cela semble excellent... mais est-ce suffisant ?

Le paradoxe de l'accuracy (classes déséquilibrées)

Attention — Le piège de l'accuracy !

Imaginons un modèle encore plus « simple » : il prédit **tout le monde comme négatif** (sain), sans même faire de test !

Grandeur	Valeur
TP	0
FP	0
FN	50 (tous les malades sont manqués !)
TN	950

$$\text{Accuracy} = \frac{0 + 950}{0 + 950 + 0 + 50} = \frac{950}{1\,000} = \mathbf{95\%}$$

Ce modèle trivial obtient 95% d'accuracy **sans détecter aucun malade** ! C'est le **paradoxe de l'accuracy** : lorsque les classes sont déséquilibrées (ici, 95% de négatifs), l'accuracy est trompeuse car un modèle qui prédit toujours la classe majoritaire obtient un score élevé tout en étant **totalement inutile**.

Leçon à retenir

L'accuracy n'est une bonne métrique **que si les classes sont équilibrées** (environ 50-50). Dès qu'il y a un déséquilibre important, il faut utiliser d'autres métriques.

16.3.2 Précision (Precision)

Définition — Précision

La **précision** mesure la proportion de vrais positifs parmi toutes les prédictions positives :

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

Question à laquelle elle répond : « Parmi toutes les personnes que j'ai déclarées positives, combien le sont vraiment ? »

Analogie — Le procureur

La précision, c'est comme un procureur. Quand il accuse quelqu'un, a-t-il généralement raison ?

- Précision élevée \Rightarrow quand le modèle dit « positif », on peut lui faire confiance.
- Précision faible \Rightarrow le modèle lance beaucoup de fausses alarmes.

Calcul sur l'exemple COVID

$$\text{Precision} = \frac{45}{45 + 30} = \frac{45}{75} = 0,60 = \mathbf{60\%}$$

Sur les 75 personnes déclarées positives par le test, seulement 60% sont réellement malades. Cela signifie que **40% des alarmes sont fausses !**

Quand la précision est-elle cruciale ?

La précision est particulièrement importante quand le **coût d'un faux positif est élevé** :

- **Filtre anti-spam** : un email légitime classé comme spam (FP) fait perdre un message important. Mieux vaut laisser passer quelques spams que de supprimer un vrai email.
- **Système judiciaire** : condamner un innocent (FP) est inacceptable.
- **Recommandation de produits** : recommander un produit non pertinent (FP) agace l'utilisateur.

16.3.3 Rappel (Recall / Sensibilité / TPR)

Définition — Recall (Rappel)

Le **recall** (rappel, aussi appelé **sensibilité** ou **TPR** — True Positive Rate) mesure la proportion de vrais positifs détectés parmi tous les positifs réels :

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

Question à laquelle il répond : « Parmi tous les vrais positifs, combien en ai-je trouvés ? »

Analogie — Le détective

Le recall, c'est comme un détective. Parmi tous les coupables, combien en a-t-il attrapés ?

- Recall élevé \Rightarrow le modèle ne laisse pas passer de positifs (peu de cas manqués).
- Recall faible \Rightarrow le modèle manque beaucoup de vrais positifs.

Calcul sur l'exemple COVID

$$\text{Recall} = \frac{45}{45 + 5} = \frac{45}{50} = 0,90 = \mathbf{90\%}$$

Le test détecte 90% des personnes réellement malades. Mais il en **manque 5 sur 50**, ce qui dans un contexte médical peut être très grave.

Quand le recall est-il crucial ?

Le recall est particulièrement important quand le **coût d'un faux négatif est élevé** :

- **Détection du cancer** : manquer un cancer (FN) peut être fatal. Il vaut mieux envoyer quelques patients sains vers des examens supplémentaires que de manquer un cancer.
- **Détection de fraude** : manquer une fraude (FN) coûte de l'argent.
- **Sécurité aérienne** : manquer une menace (FN) est inacceptable.

16.3.4 Spécificité (Specificity / TNR)**Définition — Spécificité**

La **spécificité** (aussi appelée **TNR** — True Negative Rate) mesure la proportion de vrais négatifs correctement identifiés parmi tous les négatifs réels :

$$\text{Specificity} = \frac{\text{TN}}{\text{TN} + \text{FP}}$$

Question à laquelle elle répond : « Parmi toutes les personnes saines, combien ai-je correctement identifiées comme saines ? »

Calcul sur l'exemple COVID

$$\text{Specificity} = \frac{920}{920 + 30} = \frac{920}{950} \approx 0,968 = \mathbf{96,8\%}$$

Le test identifie correctement 96,8% des personnes saines. Autrement dit, seulement 3,2% des personnes saines reçoivent un faux diagnostic positif.

Relation entre Recall et Spécificité

- Le **Recall** regarde les positifs : « Parmi les malades, combien sont détectés ? »
- La **Spécificité** regarde les négatifs : « Parmi les sains, combien sont correctement écartés ? »

On aimerait que les deux soient élevés, mais en pratique, améliorer l'un dégrade souvent

l'autre.

16.3.5 F1-Score

Définition — F1-Score

Le **F1-Score** est la **moyenne harmonique** de la précision et du recall :

$$F_1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

Il combine les deux métriques en un seul nombre, en **pénalisant fortement les déséquilibres**.

Pourquoi la moyenne harmonique ?

Moyenne harmonique vs moyenne arithmétique

Comparons les deux types de moyennes pour $P = 1,0$ et $R = 0,0$ (précision parfaite mais recall nul) :

— **Moyenne arithmétique** : $\frac{1,0 + 0,0}{2} = 0,5 \Rightarrow$ semble correct (50%)

— **Moyenne harmonique (F1)** : $2 \times \frac{1,0 \times 0,0}{1,0 + 0,0} = \mathbf{0,0} \Rightarrow$ sanctionne le zéro !

La moyenne harmonique est toujours **inférieure ou égale** à la moyenne arithmétique. Elle tire le score vers le bas dès qu'une des deux métriques est faible. Cela reflète la réalité : un modèle avec $R = 0$ est inutile, même si $P = 1$.

Calcul sur l'exemple COVID

Avec $P = 0,60$ et $R = 0,90$:

$$F_1 = 2 \times \frac{0,60 \times 0,90}{0,60 + 0,90} = 2 \times \frac{0,54}{1,50} = 2 \times 0,36 = \mathbf{0,72}$$

Le F1-Score de 72% reflète le compromis : le recall est bon (90%) mais la précision est moyenne (60%).

F-beta Score : généralisation**Définition — F-beta Score**

Le F_β **Score** généralise le F1 en introduisant un paramètre β qui contrôle l'importance relative du recall par rapport à la précision :

$$F_\beta = (1 + \beta^2) \times \frac{\text{Precision} \times \text{Recall}}{\beta^2 \times \text{Precision} + \text{Recall}}$$

- $\beta = 1 \Rightarrow F_1$ (précision et recall égaux)
- $\beta = 2 \Rightarrow F_2$ (le recall compte **2 fois plus** que la précision)
- $\beta = 0,5 \Rightarrow F_{0,5}$ (la précision compte **2 fois plus** que le recall)

16.3.6 Récapitulatif des métriques sur l'exemple COVID

Métrique	Formule	Valeur COVID
Accuracy	$\frac{TP+TN}{\text{Total}}$	96,5%
Précision	$\frac{TP}{TP+FP}$	60,0%
Recall	$\frac{TP}{TP+FN}$	90,0%
Spécificité	$\frac{TN}{TN+FP}$	96,8%
F1-Score	$2 \times \frac{P \times R}{P+R}$	72,0%

TABLE 16.2 – Toutes les métriques calculées sur le test COVID.

16.3.7 Le compromis Précision-Recall (Precision-Recall Tradeoff)**Le compromis Précision-Recall**

La plupart des classifieurs produisent un **score de probabilité** (entre 0 et 1), et un **seuil de décision** (threshold) détermine si on prédit « positif » ou « négatif ».

- Si on **abaisse le seuil** (ex : de 0,5 à 0,3) \Rightarrow on prédit plus de positifs \Rightarrow le **recall augmente** (on rate moins de vrais positifs) mais la **précision diminue** (plus de fausses alarmes).
- Si on **relève le seuil** (ex : de 0,5 à 0,7) \Rightarrow on prédit moins de positifs \Rightarrow la **précision augmente** (moins de fausses alarmes) mais le **recall diminue** (on manque plus de vrais positifs).

On ne peut pas maximiser les deux simultanément !

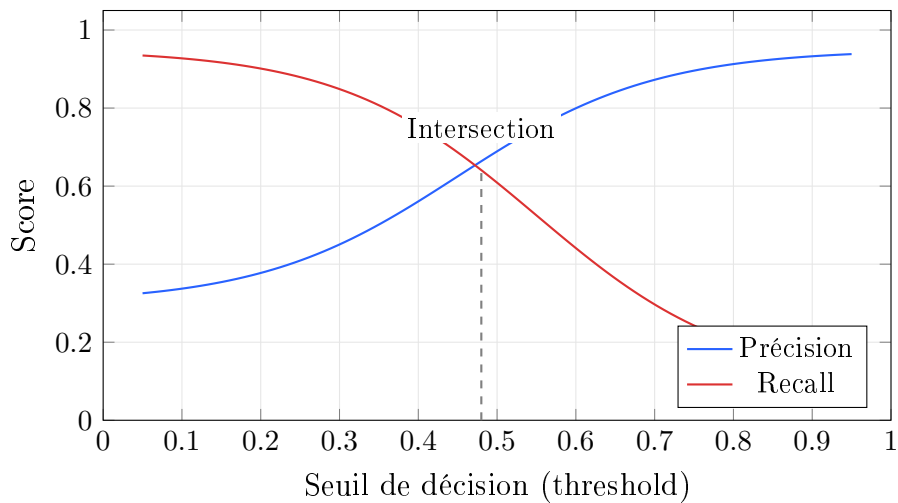


FIGURE 16.2 – Évolution de la précision et du recall en fonction du seuil de décision.

Le dilemme du médecin

Un médecin préfère un **recall élevé** (seuil bas) : il vaut mieux envoyer quelques personnes saines vers des analyses supplémentaires (FP) que de manquer un malade (FN).

Un **filtre anti-spam** préfère une **précision élevée** (seuil élevé) : il vaut mieux laisser passer quelques spams que de supprimer un email important.

16.3.8 Courbe ROC et AUC

Définition — Courbe ROC

La courbe **ROC** (Receiver Operating Characteristic) représente le **taux de vrais positifs** (Recall / TPR) en fonction du **taux de faux positifs** ($FPR = 1 - \text{Spécificité}$) pour différents seuils de décision.

$$FPR = \frac{FP}{FP + TN} = 1 - \text{Specificity}$$

Définition — AUC

L'**AUC** (Area Under the ROC Curve) est l'aire sous la courbe ROC :

- **AUC = 1,0** \Rightarrow classifieur parfait
- **AUC = 0,5** \Rightarrow classifieur aléatoire (aussi bon que pile ou face)
- **AUC < 0,5** \Rightarrow pire que le hasard (inverser les prédictions !)

Interprétation probabiliste : l'AUC est la probabilité que le modèle attribue un score

plus élevé à un exemple positif choisi aléatoirement qu'à un exemple négatif choisi aléatoirement.

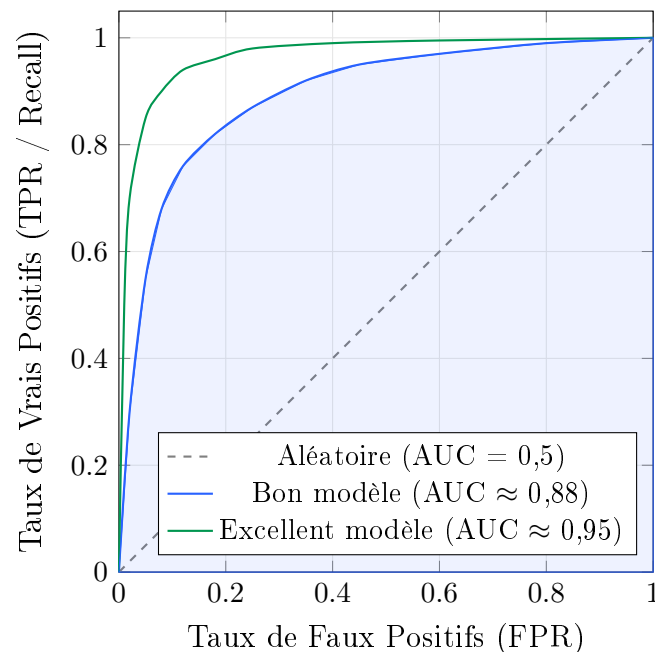


FIGURE 16.3 – Courbes ROC. Plus la courbe est proche du coin supérieur gauche, meilleur est le modèle.

Comment lire une courbe ROC ?

- Le coin **supérieur gauche** (0, 1) est le point idéal : $\text{TPR} = 1$ et $\text{FPR} = 0$.
- La **diagonale** représente un classifieur aléatoire.
- Plus la courbe « gonfle » vers le coin supérieur gauche, meilleur est le modèle.
- L'AUC résume la courbe entière en un seul nombre.

16.3.9 Extensions multi-classes

Matrice de confusion multi-classes

Pour un problème à K classes, la matrice de confusion est de taille $K \times K$. La case (i, j) indique combien d'exemples de la classe i ont été classés dans la classe j .

Exemple à 3 classes : chat ; chien et oiseau

Réalité \ Prédiction	Chat	Chien	Oiseau
Chat	40	5	2
Chien	3	35	4
Oiseau	1	2	38

Le modèle a bien classé 40 chats, 35 chiens et 38 oiseaux. Les erreurs sont sur les cases hors diagonale.

Moyennes Micro ; Macro et Weighted

Pour calculer la précision, le recall et le F1 en multi-classes, on utilise trois stratégies :

1. Macro-average : on calcule la métrique pour chaque classe séparément, puis on fait la moyenne simple.

$$\text{Macro-Precision} = \frac{1}{K} \sum_{k=1}^K \text{Precision}_k$$

⇒ Traite toutes les classes de manière **égale**, même les rares.

2. Micro-average : on agrège les TP, FP, FN de toutes les classes, puis on calcule la métrique globale.

$$\text{Micro-Precision} = \frac{\sum_k \text{TP}_k}{\sum_k \text{TP}_k + \sum_k \text{FP}_k}$$

⇒ Favorise les classes **fréquentes**.

3. Weighted-average : comme le macro, mais chaque classe est pondérée par son nombre d'exemples.

$$\text{Weighted-Precision} = \sum_{k=1}^K \frac{n_k}{N} \times \text{Precision}_k$$

⇒ Compromis entre micro et macro.

16.4 Quelle métrique choisir ?

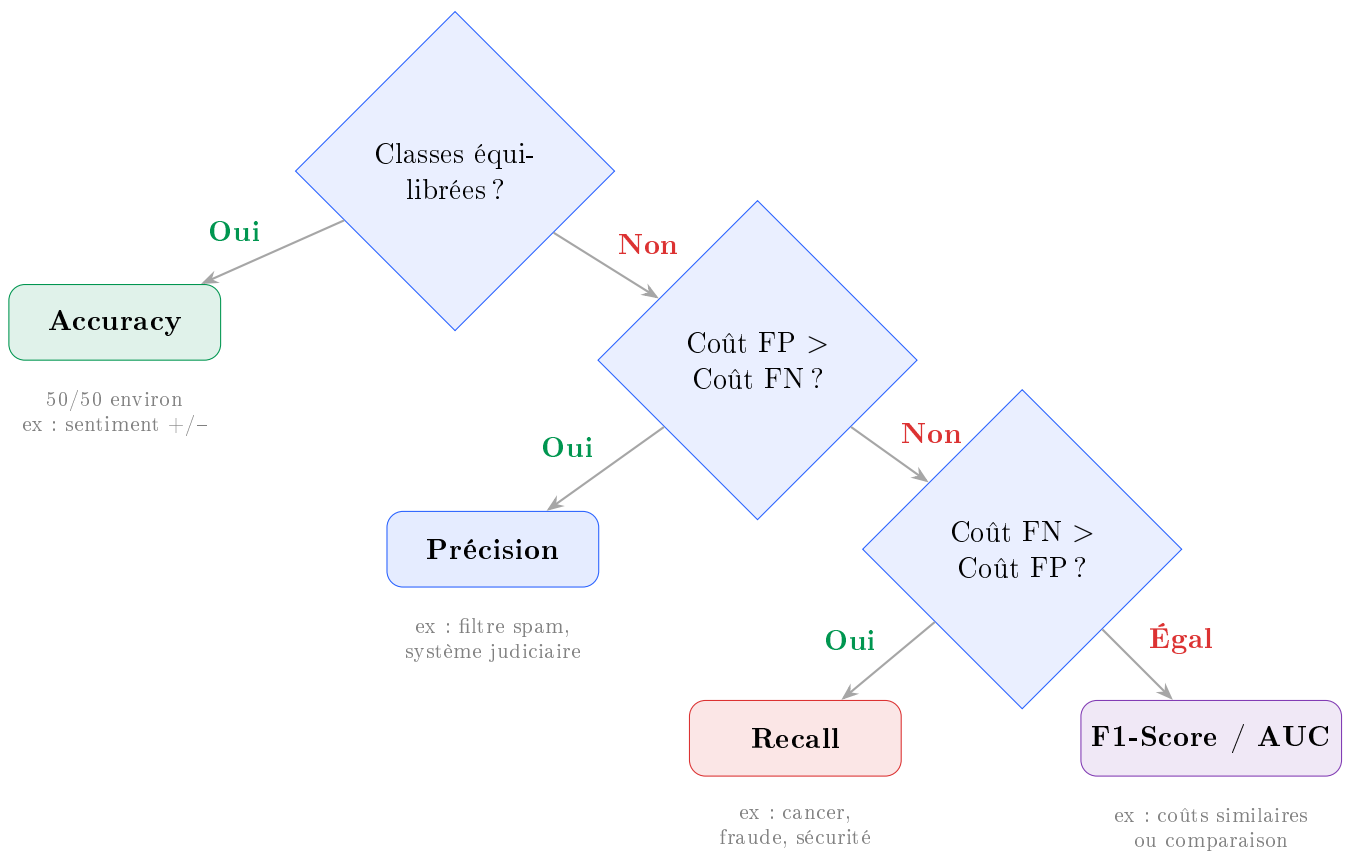


FIGURE 16.4 – Arbre de décision pour choisir la bonne métrique.

Contexte	Métrique	Justification
Classes équilibrées	Accuracy	Pas de biais vers une classe
FP coûteux	Précision	Minimiser les fausses alarmes
FN coûteux	Recall	Ne manquer aucun positif
Compromis FP/FN	F1-Score	Équilibre précision et recall
Comparer modèles	AUC-ROC	Indépendant du seuil
Multi-classes	Macro/Weighted F1	Tient compte de chaque classe

TABLE 16.3 – Guide de choix des métriques selon le contexte.

16.5 Application sur le dataset clientèle (Churn)

Revenons à notre fil rouge — le dataset ShopTech. Nous avons entraîné plusieurs modèles de classification pour prédire le **churn** (départ du client). Calculons toutes les métriques pour notre modèle de régression logistique.

Matrice de confusion du modèle de churn

Supposons que notre modèle de régression logistique, entraîné sur les 20 clients avec un split 80/20, donne les résultats suivants sur les 4 clients de test :

Réalité \ Prédiction	Churn (1)	Fidèle (0)
Churn (1)	1 (TP)	1 (FN)
Fidèle (0)	0 (FP)	2 (TN)

Calcul des métriques :

- Accuracy = $\frac{1 + 2}{1 + 2 + 0 + 1} = \frac{3}{4} = 75\%$
- Précision = $\frac{1}{1 + 0} = \frac{1}{1} = 100\%$ (aucun faux positif)
- Recall = $\frac{1}{1 + 1} = \frac{1}{2} = 50\%$ (un client churn manqué sur deux)
- F1-Score = $2 \times \frac{1,0 \times 0,5}{1,0 + 0,5} = \frac{1,0}{1,5} \approx 66,7\%$

Le modèle a une précision parfaite mais ne détecte que la moitié des churns. Dans un contexte de rétention client, c'est problématique car on manque des clients sur le point de partir.

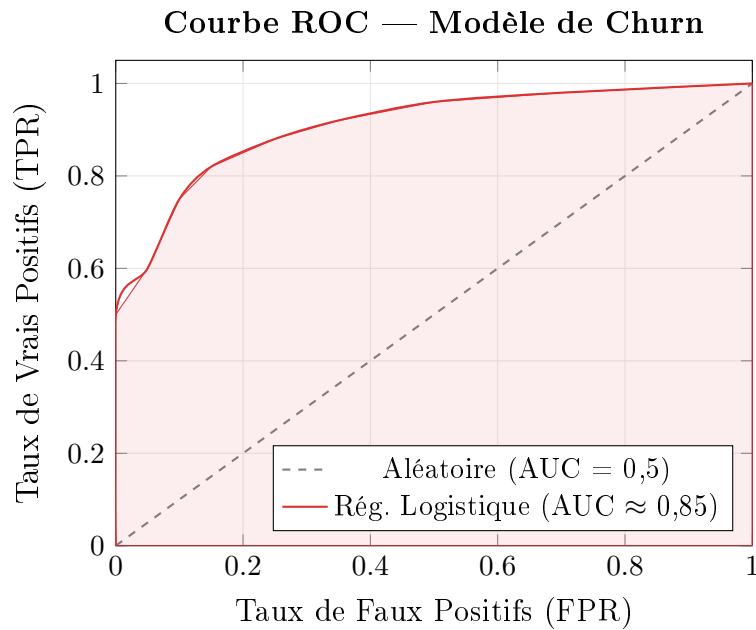


FIGURE 16.5 – Courbe ROC du modèle de régression logistique pour la prédiction du churn.

16.6 Implémentation Google Colab

Google Colab — Métriques d'évaluation complètes

Ouvrez un nouveau notebook Google Colab et exécutez les cellules suivantes.

16.6.1 Import des bibliothèques et préparation des données

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.model_selection import train_test_split
6 from sklearn.preprocessing import StandardScaler
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.neighbors import KNeighborsClassifier
9 from sklearn.tree import DecisionTreeClassifier
10 from sklearn.ensemble import RandomForestClassifier
11 from sklearn.svm import SVC
12 from sklearn.naive_bayes import GaussianNB
13 from sklearn.metrics import (confusion_matrix,
14                               classification_report,
                               accuracy_score, precision_score,
                               recall_score,
```

```

15         f1_score, roc_curve, auc,
16         precision_recall_curve,
17         ConfusionMatrixDisplay)
18
19 # Dataset ShopTech
20 data = {
21     'Age':
22         [25,34,28,45,23,52,31,41,27,38,48,29,55,33,26,44,36,50,30,42],
23     'Revenu':
24         [22,45,28,62,18,75,35,55,24,48,65,30,80,40,20,58,42,70,32,60],
25     'Anciennete': [1,5,2,8,1,12,3,7,1,6,9,2,15,4,1,7,5,10,3,8],
26     'Nb_Achats':
27         [8,30,12,55,5,70,18,45,7,35,60,14,75,25,6,50,28,68,15,48],
28     'Montant_Moyen':
29         [45,78,52,95,38,110,62,88,42,82,98,55,115,72,40,92,75,105,58,90],
30     'Score_Satisfaction':
31         [5,8,4,9,3,9,6,8,4,7,9,5,10,7,3,8,6,9,5,8],
32     'Churn': [1,0,1,0,1,0,0,0,1,0,0,1,0,0,1,0,0,0,1,0]
33 }
34 df = pd.DataFrame(data)
35
36 X = df.drop('Churn', axis=1)
37 y = df['Churn']
38 X_train, X_test, y_train, y_test = train_test_split(
39     X, y, test_size=0.3, random_state=42, stratify=y
40 )
41
42 scaler = StandardScaler()
43 X_train_sc = scaler.fit_transform(X_train)
44 X_test_sc = scaler.transform(X_test)
45
46 print(f"Train : {len(X_train)} | Test : {len(X_test)}")
47 print(f"Distribution churn (test) : {dict(y_test.value_counts())}")

```

Listing 16.1 – Import et préparation des données

16.6.2 Matrice de confusion avec heatmap

```

1 # Entraîner un modele de regression logistique
2 log_reg = LogisticRegression(random_state=42)
3 log_reg.fit(X_train_sc, y_train)

```

```

4 y_pred = log_reg.predict(X_test_sc)
5
6 # Matrice de confusion
7 cm = confusion_matrix(y_test, y_pred)
8 print("Matrice de confusion :")
9 print(cm)
10
11 # Affichage avec heatmap
12 plt.figure(figsize=(7, 5))
13 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
14             xticklabels=['Fidele (0)', 'Churn (1)'],
15             yticklabels=['Fidele (0)', 'Churn (1)'],
16             annot_kws={"size": 16})
17 plt.xlabel('Prediction', fontsize=13)
18 plt.ylabel('Realite', fontsize=13)
19 plt.title('Matrice de confusion - Regression Logistique',
20          fontsize=14)
21 plt.tight_layout()
22 plt.show()

```

Listing 16.2 – Matrice de confusion avec seaborn

16.6.3 Rapport de classification complet

```

1 # Rapport complet
2 print("=" * 55)
3 print("RAPPORT DE CLASSIFICATION")
4 print("=" * 55)
5 print(classification_report(y_test, y_pred,
6                             target_names=['Fidele', 'Churn']))
7
8 # Metriques individuelles
9 print(f"Accuracy   : {accuracy_score(y_test, y_pred):.4f}")
10 print(f"Precision  : {precision_score(y_test, y_pred):.4f}")
11 print(f"Recall     : {recall_score(y_test, y_pred):.4f}")
12 print(f"F1-Score   : {f1_score(y_test, y_pred):.4f}")

```

Listing 16.3 – Classification report de sklearn

16.6.4 Courbe ROC

```

1  # Probabilites predites
2  y_proba = log_reg.predict_proba(X_test_sc)[: , 1]
3
4  # Calcul de la courbe ROC
5  fpr, tpr, thresholds_roc = roc_curve(y_test, y_proba)
6  roc_auc = auc(fpr, tpr)
7
8  # Affichage
9  plt.figure(figsize=(8, 6))
10 plt.plot(fpr, tpr, color='#2962FF', lw=2,
11          label=f'ROC (AUC = {roc_auc:.2f})')
12 plt.plot([0, 1], [0, 1], 'k--', lw=1, label='Aleatoire (AUC =
13          0.50)')
14 plt.fill_between(fpr, tpr, alpha=0.1, color='#2962FF')
15 plt.xlim([0, 1])
16 plt.ylim([0, 1.05])
17 plt.xlabel('Taux de Faux Positifs (FPR)', fontsize=12)
18 plt.ylabel('Taux de Vrais Positifs (TPR)', fontsize=12)
19 plt.title('Courbe ROC - Regression Logistique', fontsize=14)
20 plt.legend(loc='lower right', fontsize=11)
21 plt.grid(True, alpha=0.3)
22 plt.tight_layout()
23 plt.show()

```

Listing 16.4 – Courbe ROC avec AUC

16.6.5 Courbe Précision-Recall

```

1  # Calcul de la courbe precision-recall
2  precision_curve, recall_curve, thresholds_pr =
3      precision_recall_curve(
4      y_test, y_proba
5  )
6  plt.figure(figsize=(8, 6))
7  plt.plot(recall_curve, precision_curve, color='#DC3232', lw=2,
8          label='Precision-Recall')
9  plt.xlabel('Recall', fontsize=12)

```

```

10 plt.ylabel('Precision', fontsize=12)
11 plt.title('Courbe Precision-Recall', fontsize=14)
12 plt.legend(loc='lower left', fontsize=11)
13 plt.grid(True, alpha=0.3)
14 plt.xlim([0, 1.05])
15 plt.ylim([0, 1.05])
16 plt.tight_layout()
17 plt.show()

```

Listing 16.5 – Courbe Précision-Recall

16.6.6 Comparaison de tous les classifieurs

```

1  # Définir les modeles
2  models = {
3      'Reg. Logistique': LogisticRegression(random_state=42),
4      'KNN (k=3)': KNeighborsClassifier(n_neighbors=3),
5      'Arbre Decision': DecisionTreeClassifier(random_state=42),
6      'Random Forest': RandomForestClassifier(n_estimators=100,
7                                              random_state=42),
8      'SVM (RBF)': SVC(kernel='rbf', probability=True,
9                      random_state=42),
10     'Naive Bayes': GaussianNB()
11 }
12
13 # Tableau de resultats
14 results = []
15 plt.figure(figsize=(10, 7))
16
17 for name, model in models.items():
18     model.fit(X_train_sc, y_train)
19     y_pred_m = model.predict(X_test_sc)
20
21     acc = accuracy_score(y_test, y_pred_m)
22     prec = precision_score(y_test, y_pred_m, zero_division=0)
23     rec = recall_score(y_test, y_pred_m, zero_division=0)
24     f1 = f1_score(y_test, y_pred_m, zero_division=0)
25
26     # AUC si le modele supporte predict_proba
27     if hasattr(model, 'predict_proba'):
28         y_proba_m = model.predict_proba(X_test_sc)[: , 1]

```

```

29         fpr_m, tpr_m, _ = roc_curve(y_test, y_proba_m)
30         auc_m = auc(fpr_m, tpr_m)
31         plt.plot(fpr_m, tpr_m, lw=2,
32                  label=f'{name} (AUC={auc_m:.2f})')
33     else:
34         auc_m = None
35
36     results.append({
37         'Modele': name, 'Accuracy': f'{acc:.2f}',
38         'Precision': f'{prec:.2f}', 'Recall': f'{rec:.2f}',
39         'F1-Score': f'{f1:.2f}',
40         'AUC': f'{auc_m:.2f}' if auc_m else 'N/A'
41     })
42
43 # Courbes ROC comparees
44 plt.plot([0, 1], [0, 1], 'k--', lw=1, label='Aleatoire')
45 plt.xlabel('FPR', fontsize=12)
46 plt.ylabel('TPR', fontsize=12)
47 plt.title('Courbes ROC - Comparaison des classifieurs',
48           fontsize=14)
49 plt.legend(loc='lower right', fontsize=9)
50 plt.grid(True, alpha=0.3)
51 plt.tight_layout()
52 plt.show()
53
54 # Tableau recapitulatif
55 df_results = pd.DataFrame(results)
56 print("\n" + "=" * 70)
57 print("COMPARAISON DES CLASSIFIEURS")
58 print("=" * 70)
59 print(df_results.to_string(index=False))

```

Listing 16.6 – Tableau comparatif des classifieurs

16.6.7 Ajustement du seuil de décision

```

1 # Etude de l'impact du seuil sur les metriques
2 log_reg.fit(X_train_sc, y_train)
3 y_proba_lr = log_reg.predict_proba(X_test_sc)[: , 1]
4
5 thresholds = np.arange(0.1, 0.95, 0.05)

```

```

6 metrics_by_threshold = []
7
8 for t in thresholds:
9     y_pred_t = (y_proba_lr >= t).astype(int)
10    metrics_by_threshold.append({
11        'Seuil': t,
12        'Precision': precision_score(y_test, y_pred_t,
13                                    zero_division=0),
14        'Recall': recall_score(y_test, y_pred_t,
15                              zero_division=0),
16        'F1': f1_score(y_test, y_pred_t, zero_division=0)
17    })
18
19 df_thresh = pd.DataFrame(metrics_by_threshold)
20
21 plt.figure(figsize=(10, 6))
22 plt.plot(df_thresh['Seuil'], df_thresh['Precision'],
23         'b-o', label='Precision', markersize=4)
24 plt.plot(df_thresh['Seuil'], df_thresh['Recall'],
25         'r-s', label='Recall', markersize=4)
26 plt.plot(df_thresh['Seuil'], df_thresh['F1'],
27         'g-^', label='F1-Score', markersize=4)
28 plt.xlabel('Seuil de decision', fontsize=12)
29 plt.ylabel('Score', fontsize=12)
30 plt.title('Métriques en fonction du seuil de decision',
31         fontsize=14)
32 plt.legend(fontsize=11)
33 plt.grid(True, alpha=0.3)
34 plt.tight_layout()
35 plt.show()
36
37 # Trouver le seuil optimal pour F1
38 best_idx = df_thresh['F1'].idxmax()
39 best_threshold = df_thresh.loc[best_idx, 'Seuil']
40 print(f"\nSeuil optimal (F1 max) : {best_threshold:.2f}")
41 print(f"F1 a ce seuil : {df_thresh.loc[best_idx, 'F1']:.4f}")

```

Listing 16.7 – Optimisation du seuil de décision

16.7 Exercices

Exercice 16.1 — Calcul manuel des métriques

Un modèle de classification produit la matrice de confusion suivante :

Réalité \ Prédiction	Positif	Négatif
Positif	TP = 40	FN = 5
Négatif	FP = 10	TN = 145

Calculez (montrez le calcul complet) :

1. Le nombre total d'observations.
2. L'Accuracy.
3. La Précision.
4. Le Recall.
5. La Spécificité.
6. Le F1-Score.
7. Le FPR (False Positive Rate).

Correction de l'exercice 16.1

Données : TP = 40, FP = 10, FN = 5, TN = 145.

1. Nombre total d'observations :

$$N = TP + FP + FN + TN = 40 + 10 + 5 + 145 = \mathbf{200}$$

2. Accuracy :

$$\text{Accuracy} = \frac{TP + TN}{N} = \frac{40 + 145}{200} = \frac{185}{200} = \mathbf{0,925 = 92,5\%}$$

3. Précision :

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{40}{40 + 10} = \frac{40}{50} = \mathbf{0,80 = 80\%}$$

4. Recall :

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{40}{40 + 5} = \frac{40}{45} \approx \mathbf{0,889 = 88,9\%}$$

5. Spécificité :

$$\text{Specificity} = \frac{\text{TN}}{\text{TN} + \text{FP}} = \frac{145}{145 + 10} = \frac{145}{155} \approx 0,935 = 93,5\%$$

6. F1-Score :

$$F_1 = 2 \times \frac{P \times R}{P + R} = 2 \times \frac{0,80 \times 0,889}{0,80 + 0,889} = 2 \times \frac{0,711}{1,689} \approx 2 \times 0,421 \approx 0,842 = 84,2\%$$

7. FPR (False Positive Rate) :

$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}} = \frac{10}{10 + 145} = \frac{10}{155} \approx 0,065 = 6,5\%$$

On note que $\text{FPR} = 1 - \text{Specificity} = 1 - 0,935 = 0,065$ ✓

Interprétation : Le modèle a une bonne accuracy (92,5%) et un bon recall (88,9%), ce qui signifie qu'il détecte la plupart des positifs. La précision de 80% indique qu'il y a quelques fausses alarmes. Le F1-Score de 84,2% reflète un bon compromis.

Exercice 16.2 — Le piège de l'accuracy

Un data scientist est très fier de son modèle de détection de fraude bancaire. Il annonce :

« Mon modèle a une accuracy de **99%** sur les données de test ! »

Cependant, on sait que dans le jeu de données, seulement **1%** des transactions sont frauduleuses (99% sont légitimes).

Questions :

1. Un modèle qui prédit **toujours** « non-fraude » aurait quelle accuracy ?
2. Quelle serait la matrice de confusion de ce modèle trivial (sur 10 000 transactions dont 100 fraudes) ?
3. Calculez la précision, le recall et le F1-Score de ce modèle trivial.
4. Le modèle du data scientist est-il forcément bon ? Quelle(s) métrique(s) faut-il vérifier ?
5. Proposez des métriques plus adaptées pour ce problème.

Correction de l'exercice 16.2**1. Accuracy du modèle trivial :**

Un modèle qui prédit toujours « non-fraude » classe correctement les 99% de transactions légitimes :

$$\text{Accuracy} = \frac{9\,900}{10\,000} = 99\%$$

C'est **exactement la même accuracy** que le modèle du data scientist !

2. Matrice de confusion du modèle trivial :

Réalité \ Prédiction	Fraude	Non-fraude
Fraude (100)	TP = 0	FN = 100
Non-fraude (9 900)	FP = 0	TN = 9 900

3. Métriques du modèle trivial :

$$\text{Precision} = \frac{0}{0+0} = \text{indéfini (0/0) (convention : 0)}$$

$$\text{Recall} = \frac{0}{0+100} = \mathbf{0\%} \text{ (aucune fraude détectée!)}$$

$$F_1 = 2 \times \frac{0 \times 0}{0+0} = \mathbf{0\%}$$

4. Le modèle du data scientist est-il bon ?

Non, pas forcément. Son accuracy de 99% ne prouve rien car le modèle trivial obtient le même score. Il faut vérifier le **recall** (détecte-t-il les fraudes ?), la **précision** (ses alertes sont-elles fiables ?) et le **F1-Score**.

5. Métriques adaptées pour la détection de fraude :

- **Recall** : la métrique la plus importante — chaque fraude manquée coûte cher.
- **F1-Score** : pour équilibrer recall et précision.
- **AUC-ROC** : pour évaluer le modèle indépendamment du seuil.
- **Precision-Recall AUC** : plus informative que ROC-AUC pour les classes très déséquilibrées.

Leçon : L'accuracy est une métrique **trompeuse** lorsque les classes sont déséquilibrées. Toujours vérifier le recall et le F1 !

Exercice 16.3 — Comparaison complète des classifieurs (Python)

Écrivez un programme Python complet qui :

1. Charge le dataset ShopTech et prépare les données (split 70/30 avec stratification).
2. Entraîne les 4 classifieurs suivants : Régression Logistique, KNN (k=5), Arbre de Décision, Random Forest.
3. Pour chaque classifieur, affiche :
 - La matrice de confusion (sous forme de heatmap).
 - L'accuracy, la précision, le recall et le F1-Score.

4. Trace les 4 courbes ROC sur un seul graphique avec la diagonale de référence.
5. Affiche un tableau récapitulatif comparant toutes les métriques.
6. Détermine le meilleur modèle selon le F1-Score et justifie le choix.

Correction de l'exercice 16.3

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.model_selection import train_test_split
6 from sklearn.preprocessing import StandardScaler
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.neighbors import KNeighborsClassifier
9 from sklearn.tree import DecisionTreeClassifier
10 from sklearn.ensemble import RandomForestClassifier
11 from sklearn.metrics import (confusion_matrix, accuracy_score,
12                             precision_score, recall_score, f1_score, roc_curve, auc)
13
14 # 1. Chargement des données
15 data = {
16     'Age': [25,34,28,45,23,52,31,41,27,38,
17            48,29,55,33,26,44,36,50,30,42],
18     'Revenu': [22,45,28,62,18,75,35,55,24,48,
19              65,30,80,40,20,58,42,70,32,60],
20     'Anciennete': [1,5,2,8,1,12,3,7,1,6,
21                  9,2,15,4,1,7,5,10,3,8],
22     'Nb_Achats': [8,30,12,55,5,70,18,45,7,35,
23                 60,14,75,25,6,50,28,68,15,48],
24     'Montant_Moyen': [45,78,52,95,38,110,62,88,42,82,
25                     98,55,115,72,40,92,75,105,58,90],
26     'Score_Satisfaction': [5,8,4,9,3,9,6,8,4,7,
27                           9,5,10,7,3,8,6,9,5,8],
28     'Churn': [1,0,1,0,1,0,0,0,1,0,0,1,0,0,1,0,0,0,1,0]
29 }
30 df = pd.DataFrame(data)
31 X = df.drop('Churn', axis=1)
32 y = df['Churn']
33
34 # Split 70/30 avec stratification

```

```

35 X_train, X_test, y_train, y_test = train_test_split(
36     X, y, test_size=0.3, random_state=42, stratify=y
37 )
38 scaler = StandardScaler()
39 X_train_sc = scaler.fit_transform(X_train)
40 X_test_sc = scaler.transform(X_test)
41
42 # 2. Définir les 4 classifieurs
43 classifiers = {
44     'Reg. Logistique': LogisticRegression(random_state=42),
45     'KNN (k=5)': KNeighborsClassifier(n_neighbors=5),
46     'Arbre Decision': DecisionTreeClassifier(random_state=42),
47     'Random Forest': RandomForestClassifier(
48         n_estimators=100, random_state=42)
49 }
50
51 # 3. Entraîner et évaluer chaque classifieur
52 fig_cm, axes_cm = plt.subplots(1, 4, figsize=(18, 4))
53 fig_roc, ax_roc = plt.subplots(figsize=(8, 6))
54 results = []
55
56 for idx, (name, clf) in enumerate(classifiers.items()):
57     clf.fit(X_train_sc, y_train)
58     y_pred = clf.predict(X_test_sc)
59
60     # Matrice de confusion (heatmap)
61     cm = confusion_matrix(y_test, y_pred)
62     sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
63                 ax=axes_cm[idx],
64                 xticklabels=['Fidèle', 'Churn'],
65                 yticklabels=['Fidèle', 'Churn'])
66     axes_cm[idx].set_title(name, fontsize=11)
67     axes_cm[idx].set_xlabel('Prediction')
68     axes_cm[idx].set_ylabel('Realite')
69
70     # Métriques
71     acc = accuracy_score(y_test, y_pred)
72     prec = precision_score(y_test, y_pred, zero_division=0)
73     rec = recall_score(y_test, y_pred, zero_division=0)
74     f1 = f1_score(y_test, y_pred, zero_division=0)

```

```

75
76     # Courbe ROC
77     if hasattr(clf, 'predict_proba'):
78         y_proba = clf.predict_proba(X_test_sc)[: , 1]
79     else:
80         y_proba = clf.decision_function(X_test_sc)
81     fpr_c, tpr_c, _ = roc_curve(y_test, y_proba)
82     auc_c = auc(fpr_c, tpr_c)
83     ax_roc.plot(fpr_c, tpr_c, lw=2,
84                 label=f'{name} (AUC={auc_c:.2f})')
85
86     results.append({
87         'Modele': name, 'Accuracy': acc,
88         'Precision': prec, 'Recall': rec,
89         'F1-Score': f1, 'AUC': auc_c
90     })
91
92 fig_cm.suptitle('Matrices de confusion', fontsize=14)
93 fig_cm.tight_layout()
94 plt.show()
95
96 # 4. Courbes ROC comparees
97 ax_roc.plot([0, 1], [0, 1], 'k--', lw=1,
98             label='Aleatoire')
99 ax_roc.set_xlabel('FPR', fontsize=12)
100 ax_roc.set_ylabel('TPR', fontsize=12)
101 ax_roc.set_title('Courbes ROC comparees', fontsize=14)
102 ax_roc.legend(loc='lower right')
103 ax_roc.grid(True, alpha=0.3)
104 fig_roc.tight_layout()
105 plt.show()
106
107 # 5. Tableau recapitulatif
108 df_res = pd.DataFrame(results)
109 for col in ['Accuracy', 'Precision', 'Recall', 'F1-Score', 'AUC']:
110     df_res[col] = df_res[col].apply(lambda x: f'{x:.3f}')
111
112 print("\n" + "=" * 70)
113 print("TABLEAU RECAPITULATIF DES METRIQUES")
114 print("=" * 70)

```

```
115 print(df_res.to_string(index=False))
116
117 # 6. Meilleur modele selon F1
118 df_res_num = pd.DataFrame(results)
119 best = df_res_num.loc[df_res_num['F1-Score'].idxmax()]
120 print(f"\nMeilleur modele (F1-Score) : {best['Modele']}")
121 print(f"    F1 = {best['F1-Score']:.3f}")
122 print(f"    Recall = {best['Recall']:.3f}")
123 print(f"    Precision = {best['Precision']:.3f}")
124 print("\nJustification : le F1-Score est le meilleur "
125       "compromis entre precision et recall pour ce "
126       "probleme de churn ou les deux types d'erreurs "
127       "ont un cout non negligeable.")
```

Listing 16.8 – Comparaison complète de 4 classifieurs

Interprétation attendue :

- Le tableau récapitulatif permet de comparer objectivement les 4 modèles.
- Pour le problème de churn, le **recall** est important car manquer un client qui va partir coûte cher (on perd le client).
- Mais la **précision** compte aussi : mobiliser l'équipe de rétention pour un faux positif gaspille des ressources.
- Le **F1-Score** offre le meilleur compromis, c'est pourquoi on l'utilise comme critère de sélection final.
- La courbe ROC et l'AUC permettent de comparer les modèles indépendamment du seuil de décision.

Troisième partie

Méthodes Avancées et Ensembles

Chapitre 17

Gradient Boosting, XGBoost, LightGBM et CatBoost

« Chaque arbre apprend des erreurs de ceux qui le précèdent. »

Le Gradient Boosting est l'un des algorithmes les plus **puissants** et les plus **utilisés** en Machine Learning. Il est derrière de nombreuses victoires dans les compétitions Kaggle et dans les applications industrielles. Dans ce chapitre, nous allons comprendre comment il fonctionne, pourquoi il est si efficace, et découvrir ses variantes modernes : **XGBoost**, **LightGBM** et **CatBoost**.

17.1 Hands-On : prédire un prix en corrigeant ses erreurs

Mise en situation — Estimation immobilière à plusieurs tentatives

Imaginez que vous devez estimer le prix d'une maison qui vaut en réalité **250 000 €**. Plutôt que d'essayer de trouver le bon prix d'un coup, vous allez procéder **pas à pas**, en corrigeant vos erreurs à chaque étape. C'est exactement ce que fait le Gradient Boosting !

17.1.1 Première tentative : une estimation grossière

Notre premier « arbre » est très simple : il prédit le **prix moyen** de toutes les maisons de notre base de données.

Étape	Valeur
Prédiction initiale $F_0(x)$	200 000 €
Prix réel	250 000 €
Erreur (résidu)	$250\,000 - 200\,000 = +50\,000$ €

On s'est trompé de 50 000 €. Ce n'est pas grave ! On va **apprendre de cette erreur**.

17.1.2 Deuxième tentative : corriger l'erreur

On entraîne un **deuxième arbre**, mais celui-ci ne prédit pas le prix directement. Il prédit **l'erreur** du premier arbre. Ce deuxième arbre prédit une correction de +45 000 €.

Étape	Valeur
Prédiction du 2 ^e arbre $h_1(x)$	+45 000 €
Prédiction cumulée $F_1(x) = 200\,000 + 45\,000$	245 000 €
Erreur restante	$250\,000 - 245\,000 = +5\,000$ €

L'erreur est passée de 50 000 à 5 000 € ! On continue.

17.1.3 Troisième tentative : affiner encore

Un **troisième arbre** est entraîné sur l'erreur restante de 5 000 €. Il prédit une correction de +4 000 €.

Étape	Valeur
Prédiction du 3 ^e arbre $h_2(x)$	+4 000 €
Prédiction finale $F_2(x) = 200\,000 + 45\,000 + 4\,000$	249 000 €
Erreur finale	$250\,000 - 249\,000 = +1\,000$ €

Récapitulatif : 3 itérations de Boosting

Itération	Ce que l'arbre prédit	Prédiction cumulée	Erreur restante
0	Moyenne $\rightarrow 200\,000$	200 000 €	50 000 €
1	Erreur $\rightarrow +45\,000$	245 000 €	5 000 €
2	Erreur $\rightarrow +4\,000$	249 000 €	1 000 €

Prédiction finale : $200\,000 + 45\,000 + 4\,000 = \mathbf{249\,000\ €}$ pour une maison qui vaut 250 000 €.

L'erreur finale n'est que de **1 000 €** (0.4 %), alors que la première estimation était à 50 000 € (20 %) du prix réel !

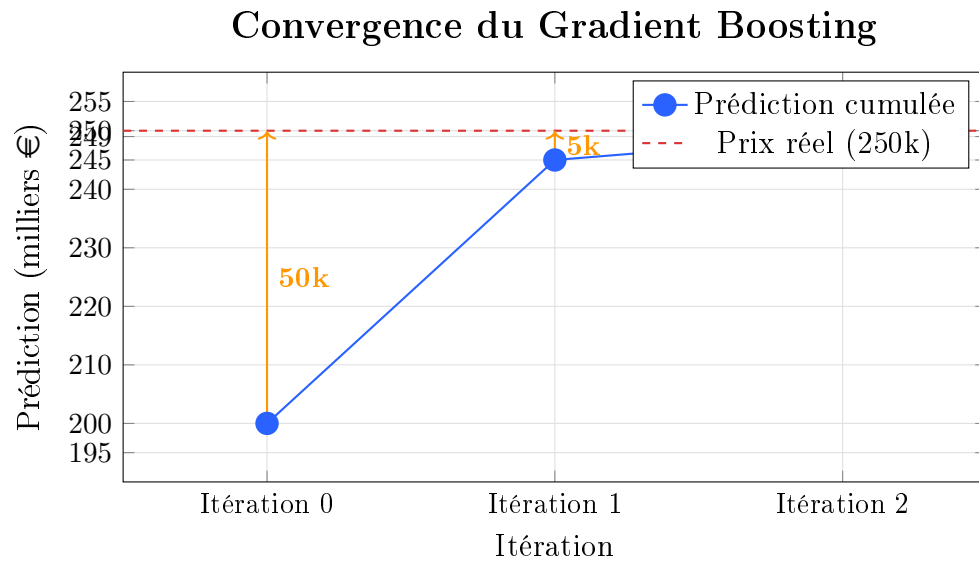


FIGURE 17.1 – La prédiction se rapproche du prix réel à chaque itération.

Comment appliquer le Gradient Boosting ?



FIGURE 17.2 – Étapes pour appliquer le Gradient Boosting (XGBoost / LightGBM / CatBoost) à un problème de Machine Learning.

17.2 Intuition : apprendre de ses erreurs

Analogie — L'étudiant qui révise intelligemment

Imaginez un étudiant qui passe un examen de 100 questions :

- 1. Première tentative** : il étudie tout et obtient 60/100.
- Il identifie les **40 questions ratées** et ne révise **que celles-là**.
- 3. Deuxième tentative** : il corrige 35 erreurs → score = 95/100.
- Il révise les **5 questions restantes**.
- 5. Troisième tentative** : il corrige 4 erreurs → score = 99/100.

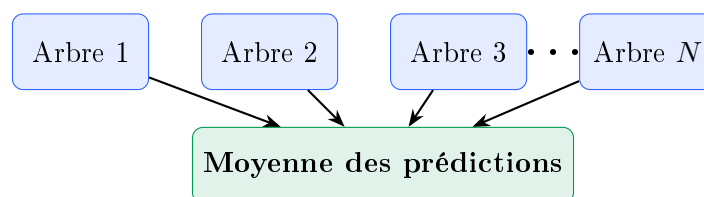
C'est exactement le Gradient Boosting : chaque nouvel arbre se **concentre uniquement** sur ce que les arbres précédents ont mal prédit.

17.2.1 Boosting vs Bagging : deux philosophies opposées

Boosting vs Bagging (Random Forest)

- **Bagging (Random Forest)** : on entraîne N arbres **en parallèle**, chacun sur un échantillon aléatoire différent, puis on fait la **moyenne**. Les arbres sont **indépendants**.
- **Boosting (Gradient Boosting)** : on entraîne les arbres **séquentiellement**. Chaque nouvel arbre apprend à corriger les erreurs du modèle précédent. Les arbres sont **dépendants**.

Bagging (Random Forest) — Parallèle



Boosting — Séquentiel

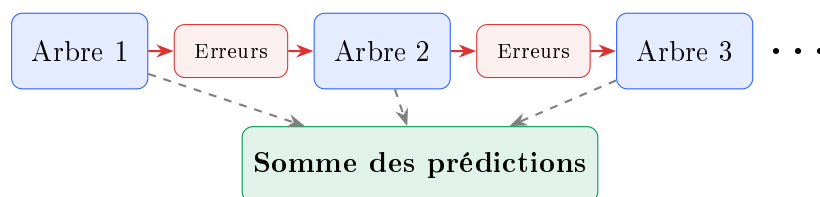


FIGURE 17.3 – Bagging : arbres indépendants en parallèle. Boosting : arbres séquentiels qui corrigent les erreurs.

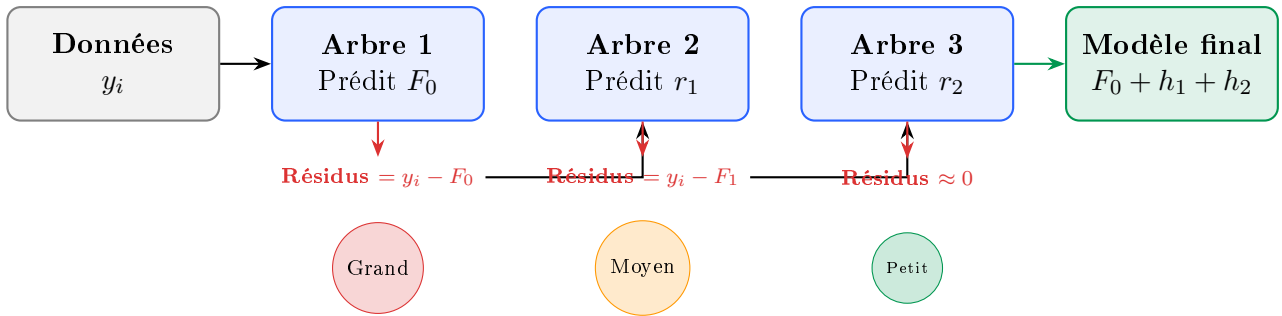


FIGURE 17.4 – Processus séquentiel du Gradient Boosting : les résidus diminuent à chaque étape.

17.3 Dérivation mathématique détaillée

17.3.1 L'idée du Boosting

Principe du Gradient Boosting

Le Gradient Boosting construit un modèle **additif** : la prédiction finale est la **somme** des prédictions de nombreux petits arbres :

$$F_M(x) = F_0(x) + \eta \cdot h_1(x) + \eta \cdot h_2(x) + \dots + \eta \cdot h_M(x) = F_0(x) + \eta \sum_{m=1}^M h_m(x)$$

où :

- $F_0(x)$ est la prédiction initiale (souvent la moyenne des y_i),
- $h_m(x)$ est le m -ième arbre de décision (appelé *weak learner*),
- $\eta \in (0, 1]$ est le **taux d'apprentissage** (*learning rate*),
- M est le nombre total d'arbres.

Le processus est le suivant :

1. On commence avec une prédiction simple $F_0(x)$ (la moyenne).
2. À l'étape m , on calcule les **erreurs** (résidus) du modèle courant.
3. On entraîne un nouvel arbre $h_m(x)$ pour prédire ces erreurs.
4. On met à jour : $F_m(x) = F_{m-1}(x) + \eta \cdot h_m(x)$.
5. On répète M fois.

17.3.2 Pourquoi « Gradient » Boosting ?

Voici la clé mathématique qui rend cet algorithme si élégant. Considérons la **fonction de perte** MSE (Mean Squared Error) :

$$L(y_i, F(x_i)) = \frac{1}{2}(y_i - F(x_i))^2$$

Calculons le **gradient** (la dérivée partielle) de cette perte par rapport à la prédiction $F(x_i)$:

$$\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} = \frac{\partial}{\partial F(x_i)} \left[\frac{1}{2}(y_i - F(x_i))^2 \right] = -(y_i - F(x_i))$$

Donc le **négatif du gradient** est :

$$-\frac{\partial L}{\partial F(x_i)} = y_i - F(x_i) = \text{résidu}$$

Découverte fondamentale

Les **résidus** (erreurs) sont exactement le **négatif du gradient** de la fonction de perte MSE !

Autrement dit, quand on entraîne un arbre sur les résidus, on fait en réalité une **descente de gradient** — mais dans l'**espace des fonctions** (et non dans l'espace des paramètres comme d'habitude).

C'est pour cela qu'on appelle cette méthode **Gradient Boosting**.

Pour d'autres fonctions de perte (entropie croisée, MAE, Huber...), les résidus ne sont plus simplement $y_i - F(x_i)$, mais on peut toujours calculer le gradient :

Perte	Formule $L(y, F)$	Pseudo-résidu $r_i = -\frac{\partial L}{\partial F}$
MSE	$\frac{1}{2}(y - F)^2$	$y_i - F(x_i)$
MAE	$ y - F $	$\text{sign}(y_i - F(x_i))$
Log-loss	$-[y \ln(p) + (1 - y) \ln(1 - p)]$	$y_i - p_i$

17.3.3 L'algorithme pas à pas

Algorithme du Gradient Boosting

Entrée : données $\{(x_i, y_i)\}_{i=1}^n$, nombre d'itérations M , taux d'apprentissage η , fonction de perte L .

Étape 0 — Initialisation :

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma)$$

Pour la MSE, c'est simplement la moyenne : $F_0(x) = \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$.

Pour $m = 1, 2, \dots, M$:

a. **Calculer les pseudo-résidus** : pour chaque observation i :

$$r_{im} = - \left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F=F_{m-1}}$$

b. **Entraîner un arbre de régression** $h_m(x)$ sur les pseudo-résidus $\{(x_i, r_{im})\}_{i=1}^n$.

c. **Calculer le pas optimal** γ_m :

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma \cdot h_m(x_i))$$

d. **Mettre à jour le modèle** :

$$F_m(x) = F_{m-1}(x) + \eta \cdot \gamma_m \cdot h_m(x)$$

Sortie : le modèle final $F_M(x) = F_0(x) + \eta \sum_{m=1}^M \gamma_m \cdot h_m(x)$.

17.3.4 Régularisation dans le Gradient Boosting

Le Gradient Boosting est un algorithme puissant mais il peut facilement faire du **surapprentissage** (overfitting). Plusieurs techniques de régularisation existent :

Techniques de régularisation

1. **Taux d'apprentissage η (shrinkage)** : une valeur faible (0.01 à 0.1) ralentit

l'apprentissage. Il faut plus d'arbres, mais le modèle généralise mieux.

$$F_m(x) = F_{m-1}(x) + \underbrace{\eta}_{\text{petit}} \cdot h_m(x)$$

2. **Profondeur maximale des arbres** : des arbres peu profonds (3 à 8 niveaux) fonctionnent mieux que des arbres profonds. Chaque arbre capture une « petite correction ».
3. **Sous-échantillonnage (Stochastic Gradient Boosting)** : à chaque itération, on n'utilise qu'une **fraction** des données (typiquement 50 % à 80 %), ce qui réduit le surapprentissage.
4. **Nombre d'arbres M** : trop d'arbres \rightarrow surapprentissage. On utilise l'**early stopping** (arrêt précoce) pour trouver le bon M .

Règle d'or

Un **petit** taux d'apprentissage η combiné à un **grand** nombre d'arbres M donne presque toujours de meilleurs résultats qu'un grand η avec peu d'arbres. La contrepartie est un temps d'entraînement plus long.

17.3.5 XGBoost — eXtreme Gradient Boosting

XGBoost — Ce qui le rend spécial

XGBoost (Chen & Guestrin, 2016) est une implémentation optimisée du Gradient Boosting qui ajoute deux ingrédients clés :

1. Un **objectif régularisé**.
2. Une **approximation de Taylor d'ordre 2** de la perte.

Objectif régularisé. XGBoost minimise l'objectif suivant :

$$\mathcal{L} = \underbrace{\sum_{i=1}^n L(y_i, \hat{y}_i)}_{\text{perte}} + \underbrace{\sum_{k=1}^M \Omega(f_k)}_{\text{régularisation}}$$

où $\Omega(f)$ pénalise la complexité de chaque arbre :

$$\Omega(f) = \gamma \cdot T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$$

- T = nombre de feuilles de l'arbre,
- w_j = poids (prédiction) de la feuille j ,
- γ = pénalité par feuille (encourage les arbres simples),
- λ = pénalité L2 sur les poids des feuilles (équivalent du Ridge).

Approximation de Taylor d'ordre 2. À l'étape m , XGBoost approche la perte par un développement de Taylor à l'ordre 2 :

$$\mathcal{L}^{(m)} \approx \sum_{i=1}^n \left[L(y_i, \hat{y}_i^{(m-1)}) + g_i \cdot f_m(x_i) + \frac{1}{2} h_i \cdot f_m(x_i)^2 \right] + \Omega(f_m)$$

où :

- $g_i = \frac{\partial L(y_i, \hat{y}_i^{(m-1)})}{\partial \hat{y}_i^{(m-1)}}$ est le **gradient** (dérivée première),
- $h_i = \frac{\partial^2 L(y_i, \hat{y}_i^{(m-1)})}{\partial (\hat{y}_i^{(m-1)})^2}$ est le **hessien** (dérivée seconde).

Pourquoi l'ordre 2 ?

Utiliser la dérivée seconde (le hessien) permet de mieux estimer la courbure de la fonction de perte. C'est comme la différence entre la méthode de Newton (rapide, utilise la courbure) et la descente de gradient classique (plus lente). Résultat : XGBoost converge **plus vite** avec **moins d'arbres**.

Autres innovations de XGBoost :

- **Histogram-based split finding** : discrétise les variables continues en « bins » pour trouver les meilleurs splits plus rapidement.
- **Gestion native des valeurs manquantes** : XGBoost apprend automatiquement de quel côté envoyer les valeurs manquantes à chaque nœud.
- **Parallélisation** : le calcul des splits est parallélisé (même si les arbres restent séquentiels).
- **Pruning** : élagage des branches qui n'améliorent pas l'objectif régularisé.

17.3.6 LightGBM — Light Gradient Boosting Machine

LightGBM — Innovations principales

LightGBM (Microsoft, 2017) est conçu pour être **plus rapide** et **plus efficace en mémoire** que XGBoost, surtout sur de grands jeux de données.

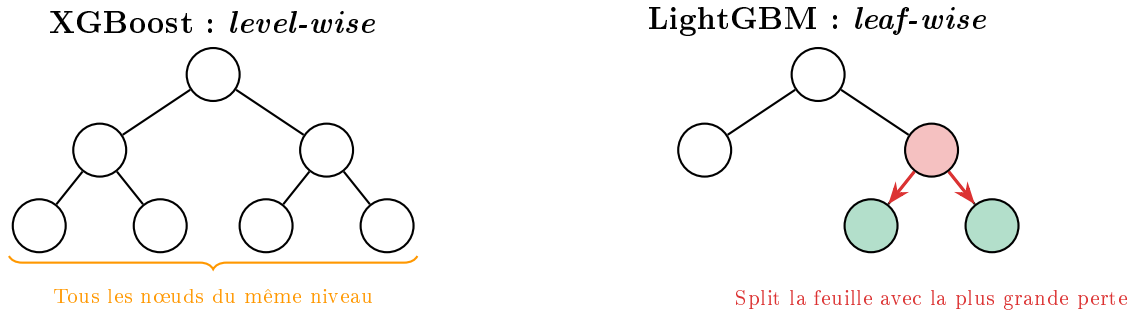


FIGURE 17.5 – À gauche : XGBoost développe tous les nœuds d'un même niveau. À droite : LightGBM développe la feuille qui réduit le plus la perte.

Croissance *leaf-wise* vs *level-wise*.

GOSS — Gradient-based One-Side Sampling. Plutôt que d'utiliser toutes les données, LightGBM garde :

- **Toutes** les observations avec un **grand gradient** (celles que le modèle prédit mal),
- Un **échantillon aléatoire** des observations avec un petit gradient (celles déjà bien prédites).

Cela accélère considérablement l'entraînement sans perdre beaucoup de précision.

EFB — Exclusive Feature Bundling. Quand il y a beaucoup de variables éparées (sparse), LightGBM **regroupe** les variables qui sont rarement non nulles en même temps, réduisant ainsi la dimensionnalité.

17.3.7 CatBoost — Categorical Boosting

CatBoost — Innovations principales

CatBoost (Yandex, 2018) se distingue par :

1. **Gestion native des variables catégorielles** : pas besoin de one-hot encoding ou de label encoding préalable.
2. **Ordered Boosting** : un mécanisme spécial pour réduire le surapprentissage lié aux pseudo-résidus.
3. **Bons hyperparamètres par défaut** : souvent performant « out of the box ».

Encodage des variables catégorielles. CatBoost utilise le **Target Encoding ordonné** : pour encoder une catégorie pour l'observation i , il n'utilise que les observations **précédant** i (dans un ordre aléatoire), ce qui évite le *target leakage*.

Ordered Boosting. Dans le GB classique, les résidus sont calculés sur les mêmes données que celles utilisées pour entraîner les arbres précédents, ce qui crée un biais. CatBoost utilise des **permutations** des données pour que chaque observation soit évaluée sur un modèle qui ne l'a jamais vue.

17.3.8 Tableau comparatif

Critère	GBM (sk-learn)	XGBoost	LightGBM	CatBoost
Vitesse	Lent	Rapide	Très rapide	Rapide
Gestion GPU	Non	Oui	Oui	Oui
Valeurs man- quantes	Non	Oui (natif)	Oui (natif)	Oui (natif)
Variables catégo- rielles	Non	Non (natif)	Partiel	Oui (natif)
Régularisation	Basique	L1, L2, γ	L1, L2	L2, randomisé
Croissance arbre	Level-wise	Level-wise	Leaf-wise	Symétrique
Approximation Taylor	Ordre 1	Ordre 2	Ordre 2	Ordre 2
Compétitions Kaggle	Rare	Très populaire	Très populaire	Populaire

TABLE 17.1 – Comparaison des 4 implémentations de Gradient Boosting.

17.4 Application sur le dataset clientèle

Appliquons le Gradient Boosting à notre problème de **prédiction du churn** (départ des clients).

17.4.1 Entraînement et comparaison avec Random Forest

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import train_test_split
4 from sklearn.ensemble import GradientBoostingClassifier,
  RandomForestClassifier
5 from sklearn.metrics import accuracy_score, classification_report
6
7 # Charger les données

```

```

8 url =
    'https://raw.githubusercontent.com/badrelkhalyly/ML-Course/main/client_da
9 df = pd.read_csv(url)
10
11 # Preparer les donnees
12 X = df.drop('churn', axis=1)
13 y = df['churn']
14 X_train, X_test, y_train, y_test = train_test_split(
15     X, y, test_size=0.2, random_state=42
16 )
17
18 # Gradient Boosting
19 gb = GradientBoostingClassifier(
20     n_estimators=200,
21     learning_rate=0.1,
22     max_depth=4,
23     subsample=0.8,
24     random_state=42
25 )
26 gb.fit(X_train, y_train)
27 y_pred_gb = gb.predict(X_test)
28 print(f'Gradient Boosting - Accuracy: {accuracy_score(y_test,
    y_pred_gb):.4f}')
29
30 # Random Forest (pour comparaison)
31 rf = RandomForestClassifier(n_estimators=200, max_depth=8,
    random_state=42)
32 rf.fit(X_train, y_train)
33 y_pred_rf = rf.predict(X_test)
34 print(f'Random Forest - Accuracy: {accuracy_score(y_test,
    y_pred_rf):.4f}')

```

Listing 17.1 – Gradient Boosting vs Random Forest sur le dataset clientèle

17.4.2 Importance des variables

```

1 import matplotlib.pyplot as plt
2
3 # Importance des variables
4 importances = gb.feature_importances_
5 indices = np.argsort(importances)[::-1]

```

```

6
7 plt.figure(figsize=(10, 6))
8 plt.title('Importance des variables - Gradient Boosting',
           fontsize=14)
9 plt.bar(range(X.shape[1]), importances[indices],
          color='steelblue', alpha=0.8)
10 plt.xticks(range(X.shape[1]), X.columns[indices], rotation=45,
             ha='right')
11 plt.ylabel('Importance')
12 plt.tight_layout()
13 plt.show()

```

Listing 17.2 – Importance des variables avec Gradient Boosting

17.4.3 Courbe d'apprentissage : accuracy vs nombre d'arbres

```

1 from sklearn.metrics import accuracy_score
2
3 # Accuracy en fonction du nombre d'arbres
4 train_scores = []
5 test_scores = []
6 n_trees_list = range(1, 201, 10)
7
8 for n in n_trees_list:
9     gb_temp = GradientBoostingClassifier(
10         n_estimators=n, learning_rate=0.1,
11         max_depth=4, random_state=42
12     )
13     gb_temp.fit(X_train, y_train)
14     train_scores.append(accuracy_score(y_train,
15                                       gb_temp.predict(X_train)))
16     test_scores.append(accuracy_score(y_test,
17                                      gb_temp.predict(X_test)))
18
19 plt.figure(figsize=(10, 6))
20 plt.plot(list(n_trees_list), train_scores, 'b-o', label='Train',
21         markersize=4)
22 plt.plot(list(n_trees_list), test_scores, 'r-s', label='Test',
23         markersize=4)
24 plt.xlabel('Nombre d'arbres')
25 plt.ylabel('Accuracy')

```

```

22 plt.title('Courbe d'apprentissage - Gradient Boosting')
23 plt.legend()
24 plt.grid(True, alpha=0.3)
25 plt.show()

```

Listing 17.3 – Courbe d'apprentissage du Gradient Boosting

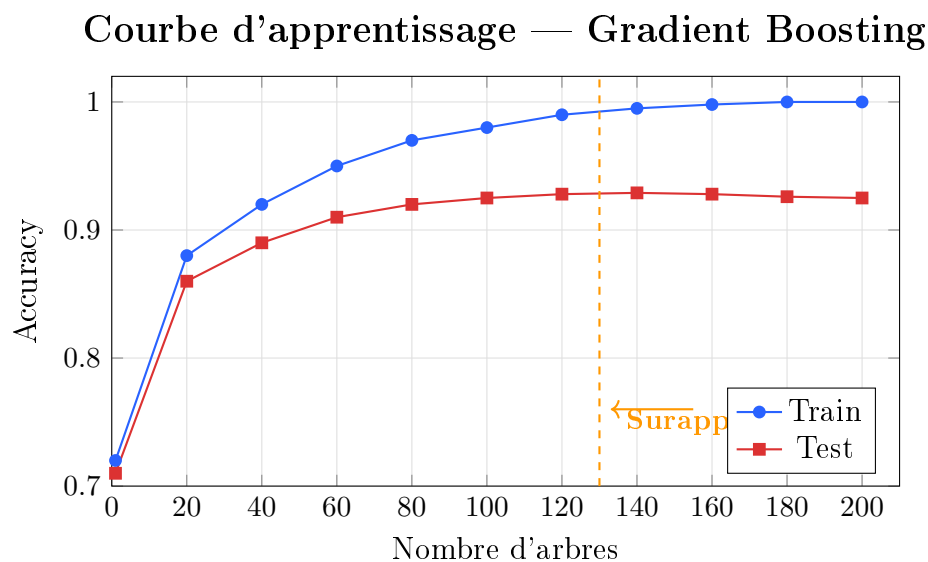


FIGURE 17.6 – Le score d'entraînement continue de monter, mais le score de test plafonne puis baisse légèrement : c'est le surapprentissage.

17.5 Implémentation complète sur Google Colab

Google Colab -- Les 4 variantes de Gradient Boosting

Ce notebook compare les 4 implémentations de Gradient Boosting sur notre dataset clientèle.

17.5.1 Installation des bibliothèques

```

1 # Installer les bibliothèques (si nécessaire)
2 !pip install xgboost lightgbm catboost -q

```

Listing 17.4 – Installation des bibliothèques (Google Colab)

17.5.2 Préparation des données

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import train_test_split
4 from sklearn.preprocessing import LabelEncoder
5 import warnings
6 warnings.filterwarnings('ignore')
7
8 # Charger les donnees
9 url =
10     'https://raw.githubusercontent.com/badrelkhalyly/ML-Course/main/client_data.csv'
11 df = pd.read_csv(url)
12 print(f''Dimensions : {df.shape}'')
13 print(df.head())
14
15 # Encoder les variables categorielles si necessaire
16 le = LabelEncoder()
17 for col in df.select_dtypes(include='object').columns:
18     if col != 'churn':
19         df[col] = le.fit_transform(df[col])
20
21 # Separer X et y
22 X = df.drop('churn', axis=1)
23 y = df['churn']
24
25 # Separer en train / test
26 X_train, X_test, y_train, y_test = train_test_split(
27     X, y, test_size=0.2, random_state=42, stratify=y
28 )
29 print(f''\nTrain : {X_train.shape[0]} lignes'')
30 print(f''Test : {X_test.shape[0]} lignes'')
```

Listing 17.5 – Chargement et préparation des données

17.5.3 Entraînement des 4 modèles

```
1 from sklearn.ensemble import GradientBoostingClassifier
2 import xgboost as xgb
3 import lightgbm as lgb
4 import catboost as cb
5 from sklearn.metrics import accuracy_score
```

```

6 import time
7
8 results = {}
9
10 # --- 1. Gradient Boosting (sklearn) ---
11 start = time.time()
12 gb_sklearn = GradientBoostingClassifier(
13     n_estimators=200, learning_rate=0.1,
14     max_depth=4, subsample=0.8, random_state=42
15 )
16 gb_sklearn.fit(X_train, y_train)
17 t_sklearn = time.time() - start
18 acc_sklearn = accuracy_score(y_test, gb_sklearn.predict(X_test))
19 results['GBM (sklearn)'] = {'accuracy': acc_sklearn, 'time':
    t_sklearn}
20 print(f''GBM sklearn - Accuracy: {acc_sklearn:.4f} - Temps:
    {t_sklearn:.2f}s'')
21
22 # --- 2. XGBoost ---
23 start = time.time()
24 xgb_model = xgb.XGBClassifier(
25     n_estimators=200, learning_rate=0.1,
26     max_depth=4, subsample=0.8,
27     reg_lambda=1.0, reg_alpha=0.0,
28     use_label_encoder=False, eval_metric='logloss',
29     random_state=42
30 )
31 xgb_model.fit(X_train, y_train)
32 t_xgb = time.time() - start
33 acc_xgb = accuracy_score(y_test, xgb_model.predict(X_test))
34 results['XGBoost'] = {'accuracy': acc_xgb, 'time': t_xgb}
35 print(f''XGBoost - Accuracy: {acc_xgb:.4f} - Temps:
    {t_xgb:.2f}s'')
36
37 # --- 3. LightGBM ---
38 start = time.time()
39 lgb_model = lgb.LGBMClassifier(
40     n_estimators=200, learning_rate=0.1,
41     max_depth=4, subsample=0.8,
42     num_leaves=31, random_state=42, verbose=-1
43 )

```

```

44 lgb_model.fit(X_train, y_train)
45 t_lgb = time.time() - start
46 acc_lgb = accuracy_score(y_test, lgb_model.predict(X_test))
47 results['LightGBM'] = {'accuracy': acc_lgb, 'time': t_lgb}
48 print(f''LightGBM      - Accuracy: {acc_lgb:.4f} - Temps:
      {t_lgb:.2f}s'')
49
50 # --- 4. CatBoost ---
51 start = time.time()
52 cb_model = cb.CatBoostClassifier(
53     iterations=200, learning_rate=0.1,
54     depth=4, verbose=0, random_state=42
55 )
56 cb_model.fit(X_train, y_train)
57 t_cb = time.time() - start
58 acc_cb = accuracy_score(y_test, cb_model.predict(X_test))
59 results['CatBoost'] = {'accuracy': acc_cb, 'time': t_cb}
60 print(f''CatBoost      - Accuracy: {acc_cb:.4f} - Temps:
      {t_cb:.2f}s'')

```

Listing 17.6 – Entraînement de GBM, XGBoost, LightGBM et CatBoost

17.5.4 Comparaison visuelle des résultats

```

1 import matplotlib.pyplot as plt
2
3 models = list(results.keys())
4 accuracies = [results[m]['accuracy'] for m in models]
5 times = [results[m]['time'] for m in models]
6 colors = ['#2962FF', '#FF6D00', '#00C853', '#AA00FF']
7
8 fig, axes = plt.subplots(1, 2, figsize=(14, 5))
9
10 # --- Accuracy ---
11 bars1 = axes[0].bar(models, accuracies, color=colors, alpha=0.85,
12     edgecolor='black')
13 axes[0].set_title(''Accuracy des 4 modeles'', fontsize=14,
14     fontweight='bold')
15 axes[0].set_ylabel(''Accuracy'')
16 axes[0].set_ylim(0.85, 1.0)
17 for bar, acc in zip(bars1, accuracies):

```

```

16     axes[0].text(bar.get_x() + bar.get_width()/2, bar.get_height()
17                 + 0.003,
18                 f'{acc:.4f}', ha='center', fontweight='bold')
19 # --- Temps d'entraînement ---
20 bars2 = axes[1].bar(models, times, color=colors, alpha=0.85,
21                    edgecolor='black')
22 axes[1].set_title('Temps d'entraînement (s)', fontsize=14,
23                  fontweight='bold')
24 axes[1].set_ylabel('Temps (secondes)')
25 for bar, t in zip(bars2, times):
26     axes[1].text(bar.get_x() + bar.get_width()/2, bar.get_height()
27                 + 0.02,
28                 f'{t:.2f}s', ha='center', fontweight='bold')
29 plt.tight_layout()
30 plt.show()

```

Listing 17.7 – Comparaison des 4 modèles en bar chart

17.5.5 Importance des variables pour chaque modèle

```

1  fig, axes = plt.subplots(2, 2, figsize=(14, 10))
2  fig.suptitle('Importance des variables', fontsize=16,
3              fontweight='bold')
4
5  models_list = [
6      ('GBM (sklearn)', gb_sklearn, colors[0]),
7      ('XGBoost', xgb_model, colors[1]),
8      ('LightGBM', lgb_model, colors[2]),
9      ('CatBoost', cb_model, colors[3])
10 ]
11
12 for ax, (name, model, color) in zip(axes.flatten(), models_list):
13     if name == 'CatBoost':
14         importances = model.get_feature_importance()
15     else:
16         importances = model.feature_importances_
17     indices = np.argsort(importances)[::-1]
18     ax.barh(range(len(indices)), importances[indices],
19            color=color, alpha=0.8)

```

```
18     ax.set_yticks(range(len(indices)))
19     ax.set_yticklabels(X.columns[indices])
20     ax.set_title(name, fontsize=12, fontweight='bold')
21     ax.invert_yaxis()
22
23 plt.tight_layout()
24 plt.show()
```

Listing 17.8 – Feature importance pour les 4 modèles

17.5.6 Optimisation des hyperparamètres avec GridSearchCV

```
1  from sklearn.model_selection import GridSearchCV
2
3  # Grille d'hyperparametres pour XGBoost
4  param_grid = {
5      'n_estimators': [100, 200, 300],
6      'learning_rate': [0.01, 0.05, 0.1],
7      'max_depth': [3, 4, 6],
8      'subsample': [0.7, 0.8, 1.0],
9      'reg_lambda': [0.1, 1.0, 5.0]
10 }
11
12 xgb_grid = GridSearchCV(
13     xgb.XGBClassifier(
14         use_label_encoder=False,
15         eval_metric='logloss',
16         random_state=42
17     ),
18     param_grid=param_grid,
19     cv=5,
20     scoring='accuracy',
21     n_jobs=-1,
22     verbose=1
23 )
24
25 xgb_grid.fit(X_train, y_train)
26
27 print(f''Meilleurs parametres : {xgb_grid.best_params_}'')
28 print(f''Meilleure accuracy (CV) : {xgb_grid.best_score_:.4f}'')
```

```

29 print(f''Accuracy sur le test : {accuracy_score(y_test,
    xgb_grid.predict(X_test)):.4f}''')

```

Listing 17.9 – GridSearchCV pour XGBoost

17.5.7 Courbes d'apprentissage comparées

```

1  from sklearn.model_selection import learning_curve
2
3  fig, axes = plt.subplots(2, 2, figsize=(14, 10))
4  fig.suptitle(''Courbes d'apprentissage'', fontsize=16,
5              fontweight='bold')
6
7  estimators = [
8      ('GBM (sklearn)', gb_sklearn, 'tab:blue'),
9      ('XGBoost', xgb_model, 'tab:orange'),
10     ('LightGBM', lgb_model, 'tab:green'),
11     ('CatBoost', cb_model, 'tab:purple')
12 ]
13
14 for ax, (name, model, color) in zip(axes.flatten(), estimators):
15     train_sizes, train_scores, test_scores = learning_curve(
16         model, X, y, cv=5, scoring='accuracy',
17         train_sizes=np.linspace(0.1, 1.0, 10),
18         n_jobs=-1
19     )
20     train_mean = train_scores.mean(axis=1)
21     test_mean = test_scores.mean(axis=1)
22     train_std = train_scores.std(axis=1)
23     test_std = test_scores.std(axis=1)
24
25     ax.fill_between(train_sizes, train_mean - train_std,
26                    train_mean + train_std, alpha=0.1,
27                    color=color)
28     ax.fill_between(train_sizes, test_mean - test_std,
29                    test_mean + test_std, alpha=0.1, color='red')
30     ax.plot(train_sizes, train_mean, 'o-', color=color,
31             label='Train')
32     ax.plot(train_sizes, test_mean, 's-', color='red',
33             label='Validation')
34     ax.set_title(name, fontsize=12, fontweight='bold')

```

```
31 ax.set_xlabel('Taille du jeu d'entraînement')
32 ax.set_ylabel('Accuracy')
33 ax.legend(loc='lower right')
34 ax.grid(True, alpha=0.3)
35
36 plt.tight_layout()
37 plt.show()
```

Listing 17.10 – Courbes d'apprentissage des 4 modèles

17.6 Schéma récapitulatif : phases d'entraînement et de test

PHASE D'ENTRAÎNEMENT (Training)**1. Étape 1 — Charger et séparer les données :**

```
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2)
```

2. Étape 2 — Initialiser les prédictions :

Formule : $F_0(x) = \bar{y}$ (la moyenne de y pour la régression).

```
import xgboost as xgb
# ou from sklearn.ensemble import GradientBoostingClassifier
```

3. Étape 3 — Itérations séquentielles ($t = 1, \dots, T$) :

Formule : À chaque itération t :

- Calculer les résidus : $r_i^{(t)} = y_i - F_{t-1}(x_i)$
- Entraîner un arbre h_t sur les résidus
- Mettre à jour : $F_t(x) = F_{t-1}(x) + \eta \cdot h_t(x)$

```
model = xgb.XGBClassifier(n_estimators=100, learning_rate=0.1,
max_depth=3)
model.fit(X_train, y_train)
```

4. Étape 4 — Régler les hyperparamètres :

Formule : η = learning rate, T = nombre d'arbres, d = profondeur max.

```
params = {'n_estimators': [50, 100, 200], 'learning_rate':
[0.01, 0.1, 0.3],
'max_depth': [3, 5, 7]}
grid = GridSearchCV(xgb.XGBClassifier(), params, cv=5)
grid.fit(X_train, y_train)
```



Modèle entraîné transmis à la phase de test

PHASE DE TEST (Testing)**1. Étape 1 — Prédiction finale :**

Formule : $F_T(\mathbf{x}) = F_0(\mathbf{x}) + \sum_{t=1}^T \eta \cdot h_t(\mathbf{x})$

```
y_pred = grid.best_estimator_.predict(X_test)
```

2. Étape 2 — Évaluer :

```
print(classification_report(y_test, y_pred))
```

3. Étape 3 — Importance des variables :

```
xgb.plot_importance(grid.best_estimator_)
```

17.7 Exercices

Exercice 17.1 — Boosting à la main

On dispose de 3 observations :

Observation	x_i	y_i (prix en k€)
1	50	150
2	80	230
3	65	180

Le taux d'apprentissage est $\eta = 1.0$ (pas de shrinkage, pour simplifier).

1. Calculer la prédiction initiale F_0 (la moyenne des y_i).
2. Calculer les résidus $r_{i1} = y_i - F_0$ pour chaque observation.
3. Supposons que le premier arbre $h_1(x)$ prédit parfaitement les résidus suivants : $h_1(x_1) = -36.67$, $h_1(x_2) = +43.33$ et $h_1(x_3) = -6.67$. Calculer $F_1(x_i) = F_0 + \eta \cdot h_1(x_i)$ pour chaque observation.
4. Calculer les nouveaux résidus $r_{i2} = y_i - F_1(x_i)$.
5. Supposons que le deuxième arbre prédit : $h_2(x_1) = 0$, $h_2(x_2) = 0$, $h_2(x_3) = 0$. Quelle est la prédiction finale ?

Correction de l'exercice 17.1

1. **Prédiction initiale :**

$$F_0 = \bar{y} = \frac{150 + 230 + 180}{3} = \frac{560}{3} \approx 186.67 \text{ k€}$$

2. **Résidus de l'étape 1 :**

$$r_{11} = y_1 - F_0 = 150 - 186.67 = -36.67$$

$$r_{21} = y_2 - F_0 = 230 - 186.67 = +43.33$$

$$r_{31} = y_3 - F_0 = 180 - 186.67 = -6.67$$

Interprétation : le modèle surprédit pour les observations 1 et 3, et sous-prédit pour l'observation 2.

3. Mise à jour F_1 (avec $\eta = 1.0$) :

$$F_1(x_1) = 186.67 + 1.0 \times (-36.67) = 150.00$$

$$F_1(x_2) = 186.67 + 1.0 \times (+43.33) = 230.00$$

$$F_1(x_3) = 186.67 + 1.0 \times (-6.67) = 180.00$$

4. Nouveaux résidus :

$$r_{12} = 150 - 150.00 = 0$$

$$r_{22} = 230 - 230.00 = 0$$

$$r_{32} = 180 - 180.00 = 0$$

Tous les résidus sont nuls ! Le modèle prédit parfaitement les données d'entraînement après seulement 1 itération (car $\eta = 1.0$ et l'arbre prédit exactement les résidus).

5. Prédiction finale : puisque h_2 prédit 0 partout, $F_2 = F_1$. Les prédictions finales sont :

$$F_2(x_1) = 150, \quad F_2(x_2) = 230, \quad F_2(x_3) = 180$$

Ce sont exactement les valeurs réelles ! Le modèle a parfaitement appris les données d'entraînement en 1 itération.

Remarque : en pratique, on utilise $\eta < 1$ (par exemple 0.1) pour éviter le surapprentissage. Dans ce cas, il faudrait beaucoup plus d'itérations pour converger.

Exercice 17.2 — Random Forest vs Gradient Boosting

Comparez le Random Forest (bagging) et le Gradient Boosting (boosting) selon les critères suivants :

1. Méthode de construction des arbres (parallèle vs séquentielle).
2. Profondeur typique des arbres.
3. Sensibilité au surapprentissage.
4. Temps d'entraînement.
5. Dans quels cas utiliser l'un plutôt que l'autre ?

Correction de l'exercice 17.2

Critère	Random Forest (Bagging)	Gradient Boosting
1. Construction	Arbres entraînés en parallèle , indépendamment les uns des autres.	Arbres entraînés séquentiellement , chacun corrigeant les erreurs du précédent.
2. Profondeur	Arbres profonds (souvent non limités). Chaque arbre est un « expert ».	Arbres peu profonds (3–8 niveaux). Chaque arbre est un « weak learner ».
3. Surapprentissage	Peu sensible : le moyennage de nombreux arbres aléatoires réduit la variance.	Plus sensible : trop d'itérations ou un learning rate trop élevé causent du surapprentissage. Nécessite un réglage fin.
4. Temps	Rapide (parallélisable facilement).	Plus lent (séquentiel), mais XGBoost/LightGBM sont très optimisés.
5. Quand utiliser	Quand on veut un modèle robuste et facile à configurer sans trop de réglage.	Quand on veut la meilleure performance possible et qu'on est prêt à investir du temps dans le réglage des hyperparamètres.

Règle pratique :

- **Random Forest** : bon choix par défaut, robuste, peu de réglage nécessaire.
- **Gradient Boosting (XGBoost/LightGBM)** : quand la performance est critique (compétitions, production), avec un réglage soigneux des hyperparamètres.

Exercice 17.3 — XGBoost sur le dataset Titanic

1. Charger le dataset Titanic depuis `seaborn` ou Kaggle.
2. Préparer les données : gérer les valeurs manquantes, encoder les variables catégorielles.
3. Entraîner un modèle XGBoost avec les paramètres par défaut.
4. Optimiser les hyperparamètres avec `GridSearchCV` (tester au moins : `n_estimators`, `max_depth`, `learning_rate`).

5. Comparer les résultats avec un Random Forest.
6. Afficher l'importance des variables et interpréter les résultats.

Correction de l'exercice 17.3

```

1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4 from sklearn.model_selection import train_test_split,
    GridSearchCV
5 from sklearn.ensemble import RandomForestClassifier
6 from sklearn.metrics import accuracy_score,
    classification_report
7 import xgboost as xgb
8 import matplotlib.pyplot as plt
9
10 # =====
11 # 1. Charger le dataset Titanic
12 # =====
13 df = sns.load_dataset('titanic')
14 print(f''Dimensions : {df.shape}'')
15 print(df.info())
16
17 # =====
18 # 2. Preparer les donnees
19 # =====
20 # Selectionner les variables utiles
21 features = ['pclass', 'sex', 'age', 'sibsp', 'parch', 'fare',
    'embarked']
22 df_clean = df[features + ['survived']].copy()
23
24 # Gerer les valeurs manquantes
25 df_clean['age'].fillna(df_clean['age'].median(), inplace=True)
26 df_clean['embarked'].fillna(df_clean['embarked'].mode()[0],
    inplace=True)
27
28 # Encoder les variables categorielles
29 df_clean['sex'] = df_clean['sex'].map({'male': 0, 'female': 1})
30 df_clean = pd.get_dummies(df_clean, columns=['embarked'],
    drop_first=True)

```

```

31
32 X = df_clean.drop('survived', axis=1)
33 y = df_clean['survived']
34
35 X_train, X_test, y_train, y_test = train_test_split(
36     X, y, test_size=0.2, random_state=42, stratify=y
37 )
38
39 # =====
40 # 3. XGBoost avec parametres par default
41 # =====
42 xgb_default = xgb.XGBClassifier(
43     use_label_encoder=False,
44     eval_metric='logloss',
45     random_state=42
46 )
47 xgb_default.fit(X_train, y_train)
48 y_pred_default = xgb_default.predict(X_test)
49 acc_default = accuracy_score(y_test, y_pred_default)
50 print(f'\n--- XGBoost (default) ---')
51 print(f'Accuracy : {acc_default:.4f}')
52 print(classification_report(y_test, y_pred_default))
53
54 # =====
55 # 4. Optimisation avec GridSearchCV
56 # =====
57 param_grid = {
58     'n_estimators': [50, 100, 200, 300],
59     'max_depth': [3, 4, 5, 6],
60     'learning_rate': [0.01, 0.05, 0.1, 0.2],
61 }
62
63 grid_search = GridSearchCV(
64     xgb.XGBClassifier(
65         use_label_encoder=False,
66         eval_metric='logloss',
67         subsample=0.8,
68         random_state=42
69     ),
70     param_grid=param_grid,

```

```

71     cv=5,
72     scoring='accuracy',
73     n_jobs=-1,
74     verbose=0
75 )
76 grid_search.fit(X_train, y_train)
77
78 print(f''\n--- XGBoost (optimise) ---'')
79 print(f''Meilleurs parametres : {grid_search.best_params_}'')
80 y_pred_tuned = grid_search.predict(X_test)
81 acc_tuned = accuracy_score(y_test, y_pred_tuned)
82 print(f''Accuracy (CV) : {grid_search.best_score_:.4f}'')
83 print(f''Accuracy (test) : {acc_tuned:.4f}'')
84
85 # =====
86 # 5. Comparaison avec Random Forest
87 # =====
88 rf = RandomForestClassifier(n_estimators=200, max_depth=6,
89                             random_state=42)
89 rf.fit(X_train, y_train)
90 y_pred_rf = rf.predict(X_test)
91 acc_rf = accuracy_score(y_test, y_pred_rf)
92
93 print(f''\n--- Comparaison ---'')
94 print(f''XGBoost (default) : {acc_default:.4f}'')
95 print(f''XGBoost (optimise) : {acc_tuned:.4f}'')
96 print(f''Random Forest : {acc_rf:.4f}'')
97
98 # Graphique comparatif
99 models_names = ['XGBoost\n(default)', 'XGBoost\n(optimise)',
100                 'Random\nForest']
101 accs = [acc_default, acc_tuned, acc_rf]
102 colors_bar = ['#2962FF', '#00C853', '#FF6D00']
103
104 plt.figure(figsize=(8, 5))
105 bars = plt.bar(models_names, accs, color=colors_bar, alpha=0.85,
106                edgecolor='black')
107 for bar, acc in zip(bars, accs):
108     plt.text(bar.get_x() + bar.get_width()/2, bar.get_height()
109             + 0.005,

```

```

108         f'{acc:.4f}', ha='center', fontweight='bold',
           fontsize=12)
109 plt.ylabel('Accuracy', fontsize=12)
110 plt.title('Titanic : XGBoost vs Random Forest', fontsize=14,
           fontweight='bold')
111
112 plt.ylim(0.7, 1.0)
113 plt.grid(axis='y', alpha=0.3)
114 plt.tight_layout()
115 plt.show()
116
117 # =====
118 # 6. Importance des variables
119 # =====
120 best_model = grid_search.best_estimator_
121 importances = best_model.feature_importances_
122 indices = np.argsort(importances)
123
124 plt.figure(figsize=(8, 5))
125 plt.barh(range(len(indices)), importances[indices],
           color='steelblue',
126           alpha=0.8)
127 plt.yticks(range(len(indices)), X.columns[indices])
128 plt.xlabel('Importance')
129 plt.title('XGBoost - Importance des variables (Titanic)',
           fontsize=14,
130           fontweight='bold')
131 plt.tight_layout()
132 plt.show()
133
134 # Interpretation
135 print('\n--- Interpretation ---')
136 for i in indices[::-1]:
137     print(f'    {X.columns[i]:15s} : {importances[i]:.4f}')
138 print('\nLes variables les plus importantes pour predire la
    survie sont :')
139 print('- Le sexe (les femmes avaient plus de chances de
    survivre)')
140 print('- Le tarif paye (lie a la classe sociale)')
141 print('- L'age (les enfants etaient prioritaires)')
142 print('- La classe du billet (1ere classe = meilleure

```

```
survie)'))
```

Listing 17.11 – Exercice 17.3 — XGBoost sur Titanic (correction complète)

Résultats typiques attendus :

- XGBoost (défaut) : accuracy ≈ 0.80 – 0.82
- XGBoost (optimisé) : accuracy ≈ 0.82 – 0.85
- Random Forest : accuracy ≈ 0.80 – 0.83

L’optimisation des hyperparamètres de XGBoost permet généralement d’obtenir une légère amélioration par rapport au Random Forest. Les variables les plus importantes sont le **sexe** (*sex*), le **tarif** (*fare*), l’**âge** (*age*) et la **classe** (*pclass*).

Résumé du Chapitre 17

- Le **Gradient Boosting** construit des arbres **séquentiellement**, chaque arbre corrigeant les erreurs du précédent.
- Les résidus sont le **négalif du gradient** de la fonction de perte \rightarrow descente de gradient dans l’espace des fonctions.
- **XGBoost** ajoute la régularisation et l’approximation de Taylor d’ordre 2.
- **LightGBM** est plus rapide grâce à la croissance *leaf-wise*, GOSS et EFB.
- **CatBoost** gère nativement les variables catégorielles et utilise l’*ordered boosting*.
- Hyperparamètres clés : `n_estimators`, `learning_rate`, `max_depth`, `subsample`.
- En pratique, XGBoost et LightGBM sont parmi les algorithmes les plus performants pour les données tabulaires.

Chapitre 18

Réseaux de Neurones (MLP — Perceptron Multi-Couches)

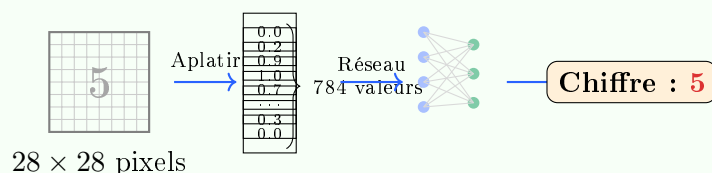
18.1 Hands-On : Reconnaissance de chiffres manuscrits

Hands-On : Un ordinateur qui lit les chiffres écrits à la main

Imaginez que vous travaillez à la poste. Chaque jour, des milliers de lettres arrivent avec des codes postaux écrits à la main. Un humain lit chaque chiffre : « c'est un 3 », « c'est un 7 »... C'est lent et coûteux.

L'objectif : Construire un programme qui reconnaît automatiquement les chiffres manuscrits (0, 1, 2, ..., 9) à partir de petites images en noir et blanc.

Le célèbre jeu de données **MNIST** contient 70 000 images de chiffres écrits à la main. Chaque image fait 28×28 pixels, soit **784 pixels** au total. Chaque pixel a une valeur entre 0 (blanc) et 255 (noir).



Le principe : Chaque image est « aplatie » en un vecteur de 784 nombres. Ce vecteur passe à travers plusieurs **couches de neurones** qui apprennent à reconnaître des formes de plus en plus complexes (bords, courbes, boucles...). La dernière couche produit 10 probabilités : une pour chaque chiffre de 0 à 9.

Résultat : Un réseau de neurones bien entraîné atteint **plus de 98% de précision** sur ce problème — mieux que beaucoup d'humains !

Comment appliquer un réseau de neurones (MLP) ?

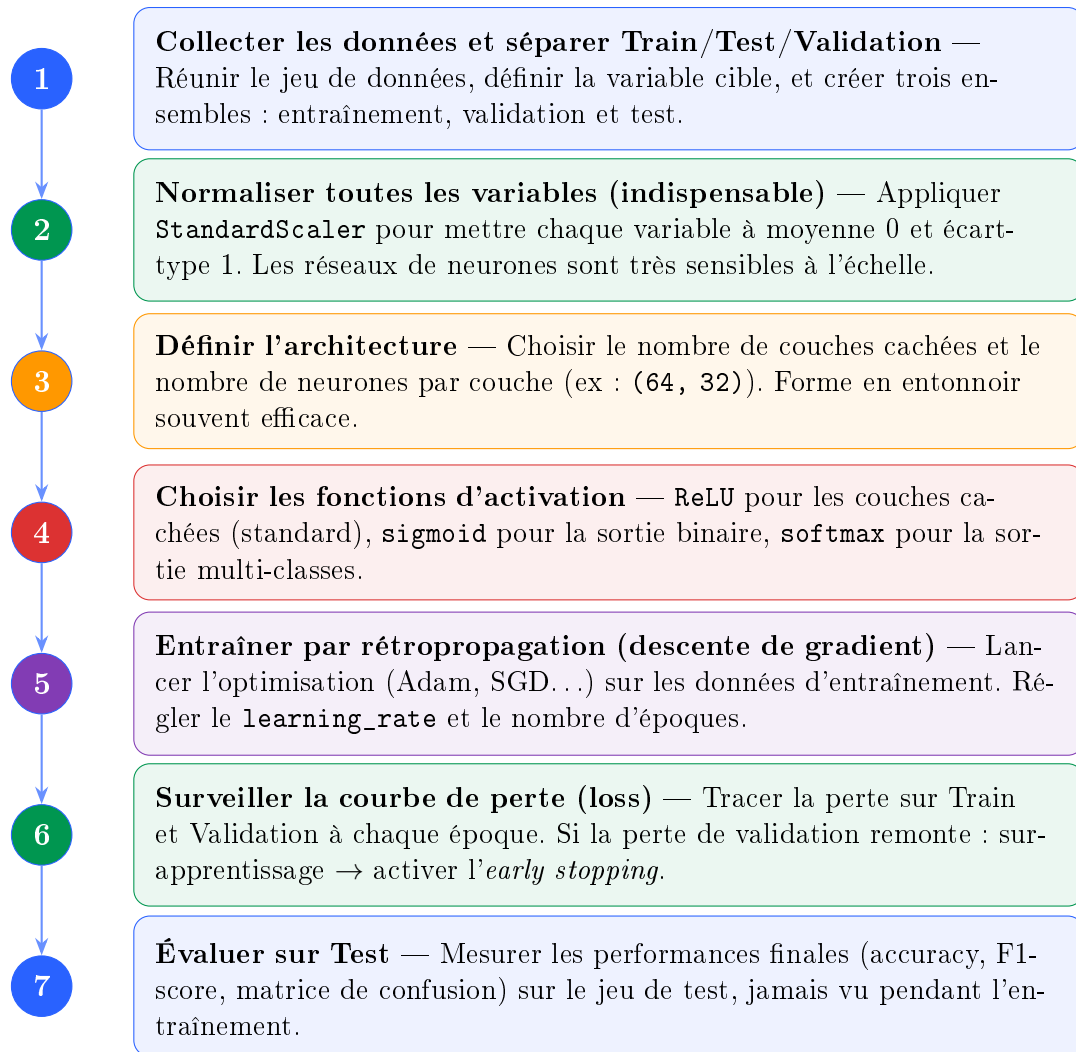


FIGURE 18.1 – Étapes pour appliquer un réseau de neurones (MLP) à un problème de Machine Learning.

18.2 Intuition : du neurone biologique au neurone artificiel

Comment fonctionne le cerveau ?

Le cerveau humain contient environ **86 milliards de neurones**. Chaque neurone :

1. **Reçoit** des signaux électriques de milliers d'autres neurones (via les **dendrites**)
2. **Traite** ces signaux : il fait une sorte de « somme pondérée » (certains signaux sont plus importants que d'autres)
3. **Décide** : si la somme dépasse un certain seuil, il « s'active » et envoie un signal aux neurones suivants (via l'**axone**)

C'est exactement ce que fait un **neurone artificiel** : il reçoit des entrées, calcule une somme pondérée, et décide s'il s'active ou non.

Analogie : la chaîne de montage

Imaginez une usine où chaque ouvrier fait une tâche simple :

- **Ouvrier 1** : vérifie si la pièce a des bords droits
- **Ouvrier 2** : vérifie si la pièce a une courbe
- **Ouvrier 3** : combine les informations des précédents pour vérifier une forme complexe
- **Chef d'équipe** : prend la décision finale

Un réseau de neurones fonctionne pareil : chaque neurone fait un calcul simple, mais l'ensemble produit un résultat très sophistiqué. **C'est la force du collectif !**

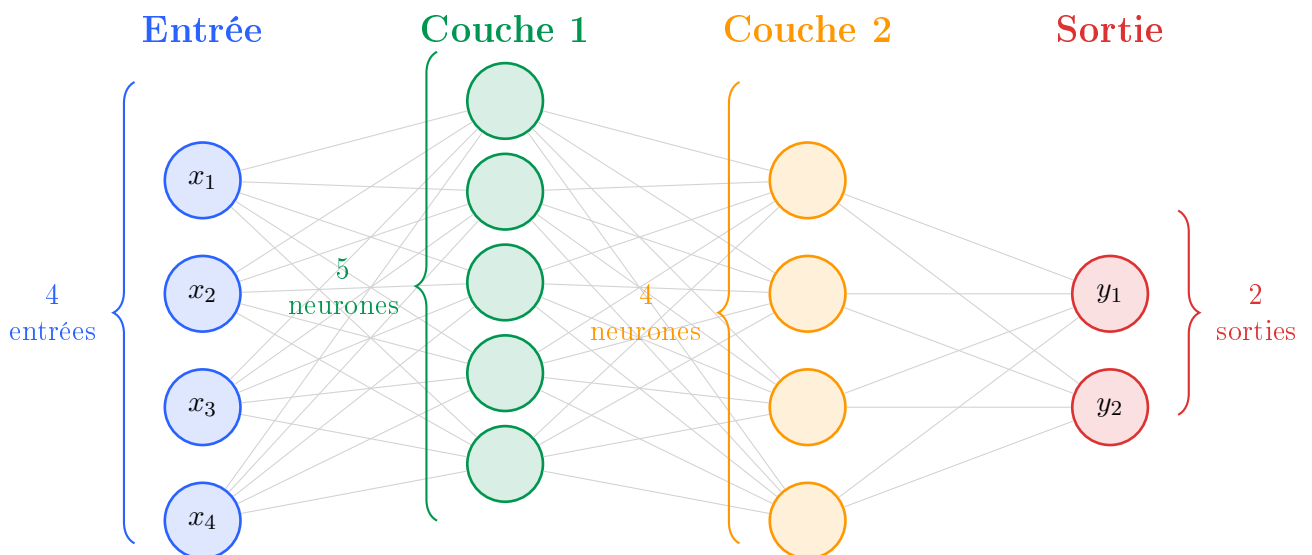


FIGURE 18.2 – Architecture d'un Perceptron Multi-Couches (MLP) avec 4 entrées, 2 couches cachées (5 et 4 neurones) et 2 sorties. Chaque neurone est connecté à **tous** les neurones de la couche suivante.

Définition — Perceptron Multi-Couches (MLP)

Un **Perceptron Multi-Couches** (Multi-Layer Perceptron, MLP) est un réseau de neurones artificiels organisé en **couches successives** :

- **Couche d'entrée** : reçoit les données brutes (x_1, x_2, \dots, x_p)
- **Couches cachées** (une ou plusieurs) : effectuent les transformations et apprennent les représentations internes
- **Couche de sortie** : produit la prédiction finale (\hat{y})

Chaque neurone d'une couche est connecté à **tous** les neurones de la couche suivante : on parle de réseau **entièrement connecté** (*fully connected*).

18.3 Dérivation mathématique détaillée

18.3.1 Le neurone unique : le perceptron

Le **perceptron** est le neurone artificiel le plus simple. Il reçoit p entrées et produit une seule sortie.

Entrées

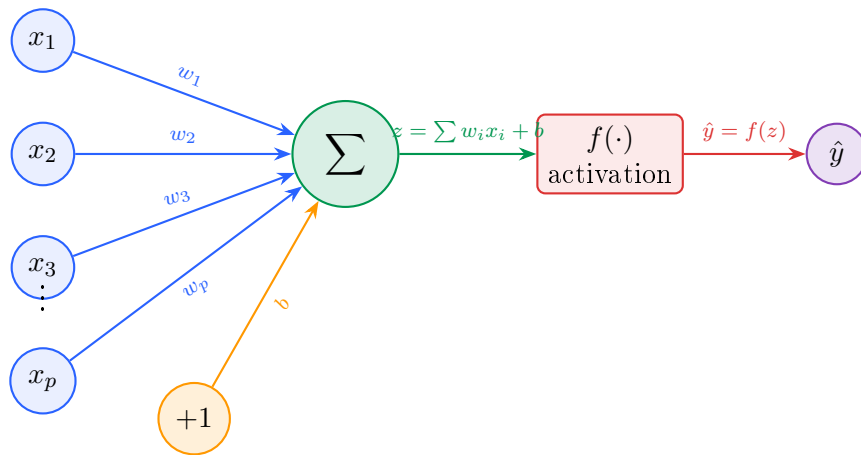


FIGURE 18.3 – Le neurone artificiel : entrées pondérées, somme, fonction d'activation, sortie.

Le neurone artificiel — Formule

Le neurone reçoit les entrées x_1, x_2, \dots, x_p et calcule :

Étape 1 — Somme pondérée :

$$z = w_1x_1 + w_2x_2 + \dots + w_px_p + b = \sum_{i=1}^p w_ix_i + b$$

où w_1, \dots, w_p sont les **poids** (weights) et b est le **biais** (bias).

Étape 2 — Activation :

$$\hat{y} = f(z)$$

où f est la **fonction d'activation** qui introduit la non-linéarité.

Le perceptron

Si on choisit la fonction sigmoid comme activation, le neurone unique calcule :

$$\hat{y} = \sigma(w_1x_1 + w_2x_2 + \dots + w_px_p + b) = \frac{1}{1 + e^{-(w_1x_1 + \dots + w_px_p + b)}}$$

C'est **exactement** la régression logistique du chapitre 10 ! Un neurone unique avec sigmoid = régression logistique. Le MLP est donc une **généralisation** de la régression logistique avec plusieurs couches.

18.3.2 Fonctions d'activation

Les fonctions d'activation sont essentielles : sans elles, un réseau de neurones multi-couches se réduirait à une simple régression linéaire (car une composition de fonctions linéaires reste linéaire).

Sigmoid : $\sigma(z) = \frac{1}{1 + e^{-z}}$

- **Sortie** : entre 0 et 1 (idéal pour les probabilités)
- **Avantage** : sortie bornée, interprétable comme une probabilité
- **Inconvénient** : le **problème du gradient évanescent** — pour $|z|$ grand, la dérivée est proche de 0, ce qui ralentit l'apprentissage dans les couches profondes

Tanh : $\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$

- **Sortie** : entre -1 et $+1$ (centrée autour de 0)
- **Avantage** : centrée \Rightarrow convergence plus rapide que sigmoid
- **Inconvénient** : souffre aussi du gradient évanescent pour $|z|$ grand

ReLU : $\text{ReLU}(z) = \max(0, z)$

- **Sortie** : 0 si $z < 0$, sinon z (pas bornée)
- **Avantage** : simple, rapide, pas de gradient évanescent pour $z > 0$ — **la plus populaire** dans les couches cachées
- **Inconvénient** : neurones « morts » (si $z < 0$ toujours, le gradient est 0 et le neurone n'apprend plus)

Softmax (pour la couche de sortie multi-classes)

$$\text{softmax}(z_k) = \frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}} \quad \text{pour } k = 1, \dots, K$$

- **Sortie** : K probabilités dont la somme vaut exactement 1
- **Usage** : uniquement en **couche de sortie** pour la classification multi-classes
- **Exemple** : pour la reconnaissance de chiffres (10 classes), la sortie softmax donne la probabilité de chaque chiffre

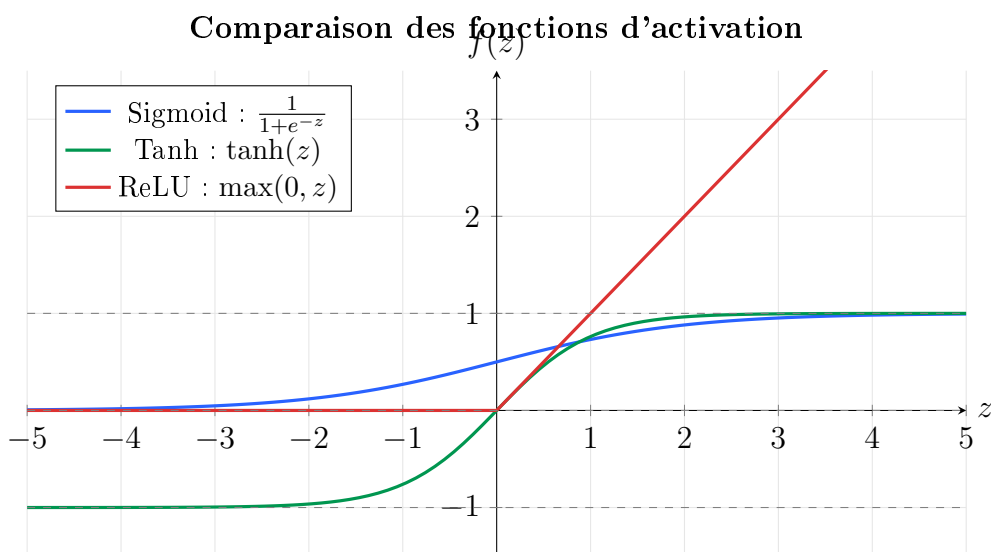


FIGURE 18.4 – Les trois fonctions d'activation principales. ReLU est la plus utilisée dans les couches cachées, sigmoid et softmax dans la couche de sortie.

18.3.3 Le Perceptron Multi-Couches (MLP)

Notation mathématique du MLP

Soit un réseau avec L couches (hors couche d'entrée). On note :

- $\mathbf{a}^{(0)} = \mathbf{x}$: le vecteur d'entrée (couche 0)
- $\mathbf{W}^{(l)}$: la **matrice de poids** de la couche l (taille $n_l \times n_{l-1}$)
- $\mathbf{b}^{(l)}$: le **vecteur de biais** de la couche l (taille n_l)
- $\mathbf{z}^{(l)} = \mathbf{W}^{(l)} \cdot \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)}$: la somme pondérée de la couche l
- $\mathbf{a}^{(l)} = f(\mathbf{z}^{(l)})$: l'activation de la couche l

La **propagation avant** (forward propagation) calcule couche par couche :

$$\mathbf{a}^{(l)} = f(\mathbf{W}^{(l)} \cdot \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)}) \quad \text{pour } l = 1, 2, \dots, L$$

18.3.4 Propagation avant : exemple numérique pas à pas

Prenons un réseau **très simple** : 2 entrées, 1 couche cachée avec 2 neurones (activation sigmoid), 1 sortie (activation sigmoid).

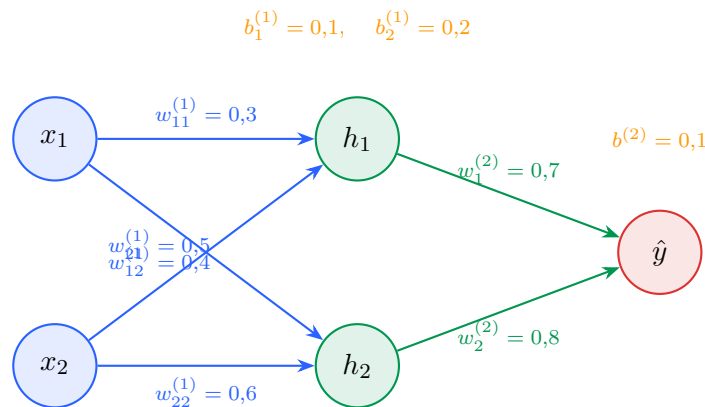


FIGURE 18.5 – Réseau 2-2-1 avec des poids concrets pour l'exemple numérique.

Propagation avant pas à pas

Données : $x_1 = 0,5$ et $x_2 = 0,8$.

Étape 1 — Couche cachée, neurone h_1 :

$$\begin{aligned} z_1^{(1)} &= w_{11}^{(1)} \cdot x_1 + w_{12}^{(1)} \cdot x_2 + b_1^{(1)} \\ &= 0,3 \times 0,5 + 0,4 \times 0,8 + 0,1 \\ &= 0,15 + 0,32 + 0,1 = 0,57 \end{aligned}$$

$$a_1^{(1)} = \sigma(0,57) = \frac{1}{1 + e^{-0,57}} = \frac{1}{1 + 0,5655} = \frac{1}{1,5655} = \mathbf{0,6388}$$

Étape 2 — Couche cachée, neurone h_2 :

$$\begin{aligned}
z_2^{(1)} &= w_{21}^{(1)} \cdot x_1 + w_{22}^{(1)} \cdot x_2 + b_2^{(1)} \\
&= 0,5 \times 0,5 + 0,6 \times 0,8 + 0,2 \\
&= 0,25 + 0,48 + 0,2 = 0,93 \\
a_2^{(1)} &= \sigma(0,93) = \frac{1}{1 + e^{-0,93}} = \frac{1}{1 + 0,3946} = \frac{1}{1,3946} = \mathbf{0,7171}
\end{aligned}$$

Étape 3 — Couche de sortie :

$$\begin{aligned}
z^{(2)} &= w_1^{(2)} \cdot a_1^{(1)} + w_2^{(2)} \cdot a_2^{(1)} + b^{(2)} \\
&= 0,7 \times 0,6388 + 0,8 \times 0,7171 + 0,1 \\
&= 0,4472 + 0,5737 + 0,1 = 1,1209 \\
\hat{y} &= \sigma(1,1209) = \frac{1}{1 + e^{-1,1209}} = \frac{1}{1 + 0,3260} = \frac{1}{1,3260} = \mathbf{0,7541}
\end{aligned}$$

Le réseau prédit $\hat{y} = 0,754$, ce qui signifie une probabilité de 75,4% pour la classe positive.

18.3.5 Fonctions de perte (loss functions)

Pour entraîner le réseau, il faut mesurer l'erreur entre la prédiction \hat{y} et la vraie valeur y . C'est le rôle de la **fonction de perte** (loss function).

Fonctions de perte principales

1. Classification binaire — Binary Cross-Entropy :

$$\mathcal{L} = -\frac{1}{n} \sum_{i=1}^n [y_i \ln(\hat{y}_i) + (1 - y_i) \ln(1 - \hat{y}_i)]$$

C'est la même fonction de coût que la régression logistique !

2. Classification multi-classes — Categorical Cross-Entropy :

$$\mathcal{L} = -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K y_{ik} \ln(\hat{y}_{ik})$$

où $y_{ik} = 1$ si l'observation i appartient à la classe k , et \hat{y}_{ik} est la probabilité prédite pour la classe k .

3. Régression — Erreur Quadratique Moyenne (MSE) :

$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

18.3.6 Rétropropagation (Backpropagation)**L'idée de la rétropropagation**

Le but de l'entraînement est de trouver les poids w et biais b qui minimisent la fonction de perte \mathcal{L} . Pour cela, on utilise la **descente de gradient** :

$$w \leftarrow w - \alpha \cdot \frac{\partial \mathcal{L}}{\partial w}$$

où α est le **taux d'apprentissage** (learning rate).

Le problème : comment calculer $\frac{\partial \mathcal{L}}{\partial w}$ pour **chaque** poids dans le réseau, y compris ceux des couches profondes ?

La réponse : la règle de la chaîne ! On propage l'erreur à l'envers, de la sortie vers l'entrée, couche par couche. C'est la **rétropropagation** (backpropagation).

Analogie : le jeu de la responsabilité

Imaginez une équipe de football qui perd un match. Le coach se demande : « Qui est responsable ? »

- L'attaquant a raté le tir \Rightarrow responsabilité directe
- Le milieu a mal passé le ballon \Rightarrow responsabilité indirecte
- Le défenseur a laissé passer l'adversaire \Rightarrow responsabilité encore plus indirecte

La rétropropagation fait exactement cela : elle répartit la « responsabilité » (l'erreur) à chaque neurone, en commençant par la sortie et en remontant vers l'entrée. Chaque poids est ajusté en proportion de sa contribution à l'erreur.

Dérivation sur le réseau 2-2-1

Reprenons notre réseau avec une observation où $y = 1$ (classe positive). Notre prédiction était $\hat{y} = 0,7541$.

Fonction de perte (Binary Cross-Entropy pour une observation) :

$$\mathcal{L} = -[y \ln(\hat{y}) + (1 - y) \ln(1 - \hat{y})] = -\ln(0,7541) = 0,2824$$

Étape 1 — Gradient de la sortie :

On calcule $\frac{\partial \mathcal{L}}{\partial z^{(2)}}$. Pour la cross-entropy avec sigmoid, on obtient un résultat très simple :

$$\delta^{(2)} = \frac{\partial \mathcal{L}}{\partial z^{(2)}} = \hat{y} - y = 0,7541 - 1 = -0,2459$$

Étape 2 — Gradients des poids de la couche de sortie :

$$\frac{\partial \mathcal{L}}{\partial w_1^{(2)}} = \delta^{(2)} \cdot a_1^{(1)} = (-0,2459) \times 0,6388 = -0,1571$$

$$\frac{\partial \mathcal{L}}{\partial w_2^{(2)}} = \delta^{(2)} \cdot a_2^{(1)} = (-0,2459) \times 0,7171 = -0,1763$$

$$\frac{\partial \mathcal{L}}{\partial b^{(2)}} = \delta^{(2)} = -0,2459$$

Étape 3 — Propager l'erreur vers la couche cachée :

Pour chaque neurone caché j , on utilise la règle de la chaîne :

$$\delta_j^{(1)} = \left(\sum_k w_k^{(2)} \cdot \delta^{(2)} \right) \cdot \sigma'(z_j^{(1)})$$

Rappel : $\sigma'(z) = \sigma(z) \cdot (1 - \sigma(z))$.

$$\sigma'(z_1^{(1)}) = 0,6388 \times (1 - 0,6388) = 0,6388 \times 0,3612 = 0,2308$$

$$\sigma'(z_2^{(1)}) = 0,7171 \times (1 - 0,7171) = 0,7171 \times 0,2829 = 0,2029$$

$$\delta_1^{(1)} = w_1^{(2)} \cdot \delta^{(2)} \cdot \sigma'(z_1^{(1)}) = 0,7 \times (-0,2459) \times 0,2308 = -0,0397$$

$$\delta_2^{(1)} = w_2^{(2)} \cdot \delta^{(2)} \cdot \sigma'(z_2^{(1)}) = 0,8 \times (-0,2459) \times 0,2029 = -0,0399$$

Étape 4 — Gradients des poids de la couche cachée :

$$\frac{\partial \mathcal{L}}{\partial w_{11}^{(1)}} = \delta_1^{(1)} \cdot x_1 = (-0,0397) \times 0,5 = -0,0199$$

$$\frac{\partial \mathcal{L}}{\partial w_{12}^{(1)}} = \delta_1^{(1)} \cdot x_2 = (-0,0397) \times 0,8 = -0,0318$$

$$\frac{\partial \mathcal{L}}{\partial w_{21}^{(1)}} = \delta_2^{(1)} \cdot x_1 = (-0,0399) \times 0,5 = -0,0200$$

$$\frac{\partial \mathcal{L}}{\partial w_{22}^{(1)}} = \delta_2^{(1)} \cdot x_2 = (-0,0399) \times 0,8 = -0,0319$$

Étape 5 — Mise à jour des poids (avec $\alpha = 0,1$) :

$$w_1^{(2)} \leftarrow 0,7 - 0,1 \times (-0,1571) = 0,7 + 0,0157 = 0,7157$$

$$w_2^{(2)} \leftarrow 0,8 - 0,1 \times (-0,1763) = 0,8 + 0,0176 = 0,8176$$

$$w_{11}^{(1)} \leftarrow 0,3 - 0,1 \times (-0,0199) = 0,3 + 0,0020 = 0,3020$$

Et ainsi de suite pour tous les poids. Les gradients étant négatifs (la prédiction est inférieure à la cible), les poids augmentent, ce qui augmentera la prédiction au prochain passage.

18.3.7 Considérations pratiques**Hyperparametres clés d'un MLP**

1. Époques (epochs) : nombre de passages complets sur toutes les données d'entraînement. Trop peu \Rightarrow sous-apprentissage. Trop \Rightarrow surapprentissage.

2. Taille du batch (batch size) : nombre d'exemples traités avant chaque mise à jour des poids.

- Batch = 1 : *stochastic gradient descent* (bruité mais rapide)
- Batch = n : *batch gradient descent* (précis mais lent)
- Batch = 32 ou 64 : *mini-batch* (le compromis le plus utilisé)

3. Taux d'apprentissage (α) : contrôle la taille des pas de la descente de gradient.

- α trop grand \Rightarrow divergence (saute au-dessus du minimum)
- α trop petit \Rightarrow convergence très lente
- Valeurs typiques : 0,001 à 0,01

Initialisation et régularisation**Initialisation des poids :**

- **Aléatoire petite** : $w \sim \mathcal{N}(0, 0,01)$ — simple mais pas toujours efficace
- **Xavier/Glorot** : $w \sim \mathcal{N}\left(0, \sqrt{\frac{2}{n_{\text{in}} + n_{\text{out}}}}\right)$ — recommandée pour sigmoid/tanh
- **He** : $w \sim \mathcal{N}\left(0, \sqrt{\frac{2}{n_{\text{in}}}}\right)$ — recommandée pour ReLU

Régularisation (pour éviter le surapprentissage) :

- **Dropout** : pendant l'entraînement, on « éteint » aléatoirement un pourcentage de neurones à chaque itération. Cela force le réseau à ne pas trop dépendre d'un seul neurone.
- **Early stopping** : on arrête l'entraînement quand la perte sur les données de validation commence à augmenter.

- **Régularisation L2** : on ajoute $\lambda \sum w_i^2$ à la fonction de perte (comme pour la régression Ridge).

Comment choisir l'architecture ?

Il n'existe pas de règle magique, mais voici des **bonnes pratiques** :

- **Nombre de couches** : 1-2 couches cachées suffisent pour la plupart des problèmes tabulaires. Plus de couches = plus de puissance mais plus de risque de surapprentissage.
- **Nombre de neurones** : commencer avec un nombre entre celui des entrées et celui des sorties, puis ajuster. Règle empirique : la première couche cachée a souvent 2 à 3 fois le nombre d'entrées.
- **Forme en entonnoir** : souvent efficace de réduire progressivement le nombre de neurones (ex : $20 \rightarrow 10 \rightarrow 5 \rightarrow 1$).

18.4 Application sur le dataset clientèle

Utilisons un MLP pour prédire le **churn** (départ d'un client). Notre réseau aura l'architecture suivante :

3 entrées (Age, Revenu, Fréquence d'achat) \rightarrow 8 neurones \rightarrow 4 neurones \rightarrow 1 sortie (Churn)

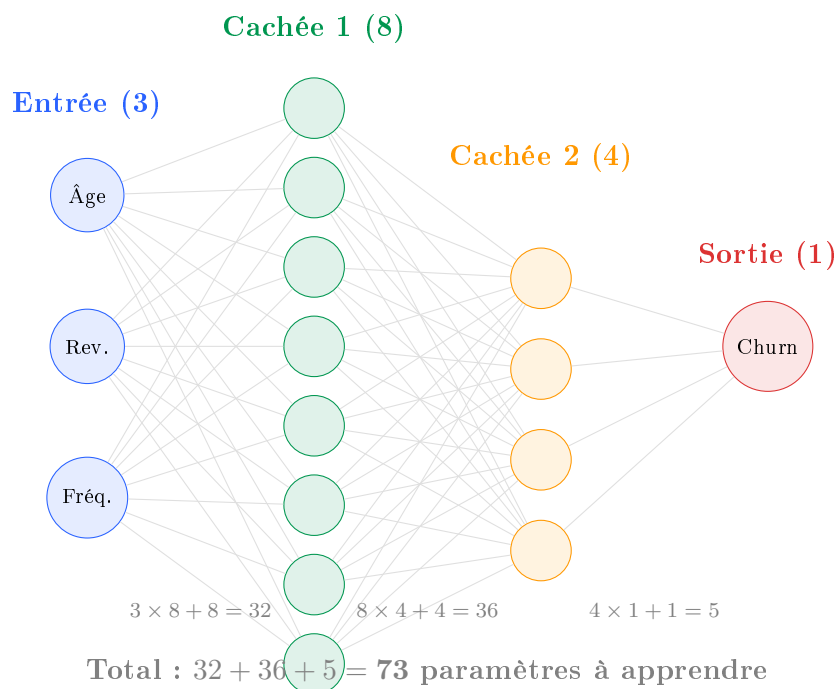


FIGURE 18.6 – Architecture du MLP pour la prédiction du churn : 3-8-4-1.

Interprétation de l'architecture

- **Couche d'entrée (3 neurones)** : reçoit l'âge, le revenu et la fréquence d'achat
- **Couche cachée 1 (8 neurones, ReLU)** : apprend des combinaisons simples des entrées (ex : « jeune avec faible fréquence »)
- **Couche cachée 2 (4 neurones, ReLU)** : apprend des combinaisons plus abstraites
- **Couche de sortie (1 neurone, sigmoid)** : produit la probabilité de churn entre 0 et 1
- **73 paramètres** au total à apprendre (poids + biais)

18.5 Implémentation sur Google Colab

Notebook Colab -- Réseaux de Neurones (MLP) complet

Ouvrez un nouveau notebook sur <https://colab.research.google.com> et suivez les étapes ci-dessous.

18.5.1 Étape 1 : Imports et préparation des données

```

1  import numpy as np
2  import pandas as pd
3  import matplotlib.pyplot as plt
4  from sklearn.neural_network import MLPClassifier, MLPRegressor
5  from sklearn.model_selection import train_test_split
6  from sklearn.preprocessing import StandardScaler
7  from sklearn.metrics import (accuracy_score, classification_report,
8                               confusion_matrix,
9                               ConfusionMatrixDisplay)
10
11  import warnings
12  warnings.filterwarnings('ignore')
13
14  # Donnees clientele (fil rouge du cours)
15  np.random.seed(42)
16  n = 200
17  age = np.random.randint(18, 70, n)
18  revenu = np.random.randint(15000, 100000, n)
19  frequence = np.random.randint(1, 30, n)

```

```

19 # Regle de churn (avec du bruit)
20 score = -0.03 * age + 0.00002 * revenu - 0.08 * frequence + 2
21 prob_churn = 1 / (1 + np.exp(-score))
22 churn = (np.random.rand(n) < prob_churn).astype(int)
23
24 df = pd.DataFrame({
25     'Age': age,
26     'Revenu': revenu,
27     'Frequence': frequence,
28     'Churn': churn
29 })
30
31 print('Apercu des donnees :')
32 print(df.head(10))
33 print(f'\nDistribution du Churn : {dict(zip(*np.unique(churn,
    return_counts=True)))}')

```

Listing 18.1 – Imports et création du dataset clientèle

18.5.2 Étape 2 : Préparation et normalisation

```

1 # Variables
2 X = df[['Age', 'Revenu', 'Frequence']].values
3 y = df['Churn'].values
4
5 # Separation train/test
6 X_train, X_test, y_train, y_test = train_test_split(
7     X, y, test_size=0.2, random_state=42, stratify=y
8 )
9
10 # IMPORTANT : normaliser les donnees pour les reseaux de neurones !
11 scaler = StandardScaler()
12 X_train_s = scaler.fit_transform(X_train)
13 X_test_s = scaler.transform(X_test)
14
15 print(f'Taille train : {X_train_s.shape[0]} exemples')
16 print(f'Taille test : {X_test_s.shape[0]} exemples')
17 print(f'\nMoyenne apres normalisation :
    {X_train_s.mean(axis=0).round(2)}')
18 print(f'Ecart-type apres normalisation :
    {X_train_s.std(axis=0).round(2)}')

```

Listing 18.2 – Séparation train/test et normalisation

Pourquoi normaliser ?

Les réseaux de neurones sont **très sensibles à l'échelle** des données. Si l'âge va de 18 à 70 et le revenu de 15 000 à 100 000, les gradients seront dominés par le revenu. La normalisation (`StandardScaler`) met toutes les variables à la même échelle (moyenne 0, écart-type 1).

18.5.3 Étape 3 : Entraîner le MLPClassifier

```

1  # MLP avec 2 couches cachees : 8 et 4 neurones
2  mlp = MLPClassifier(
3      hidden_layer_sizes=(8, 4),      # architecture 3 -> 8 -> 4 -> 1
4      activation='relu',              # fonction d'activation ReLU
5      solver='adam',                  # optimiseur Adam (variante de
        SGD)
6      max_iter=500,                    # nombre maximum d'epoques
7      random_state=42,
8      learning_rate_init=0.001,        # taux d'apprentissage initial
9      early_stopping=True,             # arret anticipe
10     validation_fraction=0.15,        # 15% pour la validation
11     verbose=False
12 )
13
14 # Entraînement
15 mlp.fit(X_train_s, y_train)
16
17 # Predictions
18 y_pred = mlp.predict(X_test_s)
19 y_proba = mlp.predict_proba(X_test_s)[:, 1]
20
21 # Resultats
22 acc = accuracy_score(y_test, y_pred)
23 print(f''Precision (accuracy) : {acc:.4f}'')
24 print(f''Nombre d'epoques realisees : {mlp.n_iter_}'')
25 print(f''\nRapport de classification :'')
26 print(classification_report(y_test, y_pred,
        target_names=['Fidele', 'Churn']))

```

Listing 18.3 – Entraînement du MLPClassifier avec architecture 8-4

18.5.4 Étape 4 : Comparer plusieurs architectures

```

1  # Tester différentes architectures
2  architectures = {
3      '(10,)' : (10,),
4      '(10, 5)' : (10, 5),
5      '(20, 10, 5)' : (20, 10, 5)
6  }
7
8  resultats = {}
9  for nom, arch in architectures.items():
10     mlp_test = MLPClassifier(
11         hidden_layer_sizes=arch,
12         activation='relu',
13         solver='adam',
14         max_iter=500,
15         random_state=42,
16         learning_rate_init=0.001,
17         early_stopping=True,
18         validation_fraction=0.15
19     )
20     mlp_test.fit(X_train_s, y_train)
21     acc_test = accuracy_score(y_test, mlp_test.predict(X_test_s))
22     n_params = sum(c.size for c in mlp_test.coefs_) + sum(b.size
23         for b in mlp_test.intercepts_)
24     resultats[nom] = {
25         'accuracy': acc_test,
26         'epochs': mlp_test.n_iter_,
27         'n_params': n_params,
28         'loss_curve': mlp_test.loss_curve_
29     }
30     print(f''Architecture {nom:>12} : accuracy = {acc_test:.4f}, ''
31         f''epoques = {mlp_test.n_iter_:>3}, ''
32         f''parametres = {n_params:>4}''')
33
34 # Tableau récapitulatif
35 print(''\n'' + ''=' * 60)

```

```

35 print(f''{'Architecture':<15} {'Accuracy':>10} {'Epoques':>10}
    {'Parametres':>12}''')
36 print(''= ' * 60)
37 for nom, res in resultats.items():
38     print(f''{nom:<15} {res['accuracy']:>10.4f}
        {res['epochs']:>10} {res['n_params']:>12}''')

```

Listing 18.4 – Comparaison de 3 architectures différentes

18.5.5 Étape 5 : Courbe d'apprentissage (loss curve)

```

1 plt.figure(figsize=(12, 6))
2
3 colors = ['#2962FF', '#00964F', '#DC3232']
4 for (nom, res), color in zip(resultats.items(), colors):
5     plt.plot(res['loss_curve'], label=f'{nom}
        (acc={res['accuracy']:.3f})',
6             linewidth=2, color=color)
7
8 plt.xlabel('Epoque', fontsize=13)
9 plt.ylabel('Perte (loss)', fontsize=13)
10 plt.title('Courbes d\'apprentissage - Comparaison des
    architectures',
11           fontsize=15, fontweight='bold')
12 plt.legend(fontsize=12)
13 plt.grid(True, alpha=0.3)
14 plt.tight_layout()
15 plt.show()

```

Listing 18.5 – Visualisation des courbes de perte

18.5.6 Étape 6 : Prédiction sur de nouveaux clients

```

1 # Nouveaux clients a predire
2 nouveaux_clients = pd.DataFrame({
3     'Age': [25, 55, 40],
4     'Revenu': [35000, 80000, 50000],
5     'Frequence': [2, 20, 10]
6 })
7

```

```

8  # Normaliser avec le meme scaler
9  X_new_s = scaler.transform(nouveaux_clients.values)
10
11 # Predictions avec le meilleur modele
12 predictions = mlp.predict(X_new_s)
13 probas = mlp.predict_proba(X_new_s)[:, 1]
14
15 print('Predictions pour les nouveaux clients :')
16 print('-' * 55)
17 for i, (_, row) in enumerate(nouveaux_clients.iterrows()):
18     statut = 'CHURN' if predictions[i] == 1 else 'FIDELE'
19     print(f'Client {i+1} (Age={row['Age']},
20           Rev={row['Revenu']},
21           f'Freq={row['Frequence']}) -> {statut} '
22           f'(proba churn = {probas[i]:.2%})')

```

Listing 18.6 – Prédiction sur de nouveaux clients

18.5.7 Étape 7 : Comparaison avec les autres classifieurs

```

1  from sklearn.linear_model import LogisticRegression
2  from sklearn.neighbors import KNeighborsClassifier
3  from sklearn.tree import DecisionTreeClassifier
4  from sklearn.ensemble import RandomForestClassifier,
   GradientBoostingClassifier
5  from sklearn.svm import SVC
6  from sklearn.naive_bayes import GaussianNB
7
8  # Dictionnaire de tous les classifieurs
9  classifieurs = {
10     'Reg. Logistique': LogisticRegression(max_iter=1000,
11     random_state=42),
11     'KNN (k=5)': KNeighborsClassifier(n_neighbors=5),
12     'Arbre Decision': DecisionTreeClassifier(max_depth=5,
13     random_state=42),
13     'Random Forest': RandomForestClassifier(n_estimators=100,
14     random_state=42),
14     'SVM (RBF)': SVC(kernel='rbf', random_state=42),
15     'Naive Bayes': GaussianNB(),
16     'Gradient Boosting':
17     GradientBoostingClassifier(random_state=42),

```

```

17     'MLP (8,4)': MLPClassifier(hidden_layer_sizes=(8,4),
18                               max_iter=500,
19                               random_state=42,
20                               early_stopping=True)
19 }
20
21 print(f''{'Classifieur':<22} {'Accuracy':>10}''')
22 print(''= ' * 35)
23
24 scores = {}
25 for nom, clf in classifieurs.items():
26     clf.fit(X_train_s, y_train)
27     acc = accuracy_score(y_test, clf.predict(X_test_s))
28     scores[nom] = acc
29     print(f''{'nom':<22} {'acc':>10.4f}''')
30
31 # Graphique comparatif
32 plt.figure(figsize=(12, 6))
33 noms = list(scores.keys())
34 accs = list(scores.values())
35 couleurs = ['#2962FF', '#FF9800', '#4CAF50', '#009688',
36             '#9C27B0', '#F44336', '#795548', '#DC3232']
37 bars = plt.barh(noms, accs, color=couleurs, edgecolor='gray',
38                height=0.6)
39
40 # Ajouter les valeurs
41 for bar, acc in zip(bars, accs):
42     plt.text(bar.get_width() + 0.003, bar.get_y() +
43             bar.get_height()/2,
44             f'{acc:.3f}', va='center', fontsize=11,
45             fontweight='bold')
46
47 plt.xlabel('Accuracy', fontsize=13)
48 plt.title('Comparaison de tous les classifieurs', fontsize=15,
49          fontweight='bold')
50 plt.xlim(0, 1.1)
51 plt.grid(axis='x', alpha=0.3)
52 plt.tight_layout()
53 plt.show()

```

Listing 18.7 – Comparaison avec tous les classifieurs vus dans le cours

18.6 Schéma récapitulatif : phases d'entraînement et de test

PHASE D'ENTRAÎNEMENT (Training)

1. **Étape 1 — Charger et séparer les données (80/10/10) :**

```
X_temp, X_test, y_temp, y_test = train_test_split(X, y,
test_size=0.1)
X_train, X_val, y_train, y_val = train_test_split(X_temp, y_temp,
test_size=0.11)
```

2. **Étape 2 — Normaliser (indispensable) :**

$$\text{Formule : } x'_j = \frac{x_j - \mu_j}{\sigma_j}$$

```
scaler = StandardScaler()
X_train_s = scaler.fit_transform(X_train)
```

3. **Étape 3 — Définir l'architecture :**

Formule : Couche l : $\mathbf{z}^{(l)} = W^{(l)}\mathbf{a}^{(l-1)} + \mathbf{b}^{(l)}$, $\mathbf{a}^{(l)} = g(\mathbf{z}^{(l)})$

```
from sklearn.neural_network import MLPClassifier
mlp = MLPClassifier(hidden_layer_sizes=(64, 32), activation='relu')
```

4. **Étape 4 — Entraîner par rétropropagation :**

Formule : Pour chaque epoch : propagation avant \rightarrow calcul de la loss $J = -\frac{1}{n} \sum [y_i \ln(\hat{p}_i) + (1 - y_i) \ln(1 - \hat{p}_i)] \rightarrow$ calcul des gradients $\frac{\partial J}{\partial W^{(l)}}$ \rightarrow mise à jour $W^{(l)} := W^{(l)} - \alpha \frac{\partial J}{\partial W^{(l)}}$.

```
mlp.fit(X_train_s, y_train)
plt.plot(mlp.loss_curve_) # courbe de perte
```

5. **Étape 5 — Surveiller le sur-apprentissage :**

Formule : Si loss_validation augmente alors que loss_train diminue \rightarrow sur-apprentissage.

```
mlp = MLPClassifier(early_stopping=True, validation_fraction=0.1)
```



Modèle entraîné transmis à la phase de test

PHASE DE TEST (Testing)

1. **Étape 1 — Normaliser X_test :**

```
X_test_s = scaler.transform(X_test)
```

2. **Étape 2 — Propagation avant :**

Formule : $\hat{y} = g_L(W^{(L)} \cdot g_{L-1}(\dots g_1(W^{(1)}\mathbf{x} + b^{(1)}) \dots) + b^{(L)})$

```
y_pred = mlp.predict(X_test_s)
```

3. Étape 3 — Évaluer :

```
print(classification_report(y_test, y_pred))
```

```
print(confusion_matrix(y_test, y_pred))
```

18.7 Exercices

Exercice 18.1 — Propagation avant à la main

Soit un réseau de neurones avec 2 entrées, 1 couche cachée de 2 neurones (activation ReLU) et 1 sortie (activation sigmoid).

Poids et biais donnés :

Couche cachée :

$$\mathbf{W}^{(1)} = \begin{pmatrix} 0,5 & -0,3 \\ 0,2 & 0,7 \end{pmatrix}, \quad \mathbf{b}^{(1)} = \begin{pmatrix} 0,1 \\ -0,1 \end{pmatrix}$$

Couche de sortie :

$$\mathbf{W}^{(2)} = \begin{pmatrix} 0,6 & 0,4 \end{pmatrix}, \quad b^{(2)} = -0,2$$

Entrée : $\mathbf{x} = (1,0 ; 2,0)$

Questions :

- Calculez la somme pondérée $z_1^{(1)}$ et $z_2^{(1)}$ de chaque neurone caché.
- Appliquez l'activation ReLU pour obtenir $a_1^{(1)}$ et $a_2^{(1)}$.
- Calculez la somme pondérée de la sortie $z^{(2)}$.
- Appliquez l'activation sigmoid pour obtenir la prédiction \hat{y} .
- Si la vraie valeur est $y = 1$, calculez l'erreur Binary Cross-Entropy.

Correction de l'exercice 18.1**a) Sommes pondérées de la couche cachée :**

$$\begin{aligned}
 z_1^{(1)} &= w_{11}^{(1)} \cdot x_1 + w_{12}^{(1)} \cdot x_2 + b_1^{(1)} \\
 &= 0,5 \times 1,0 + (-0,3) \times 2,0 + 0,1 \\
 &= 0,5 - 0,6 + 0,1 = \mathbf{0,0}
 \end{aligned}$$

$$\begin{aligned}
 z_2^{(1)} &= w_{21}^{(1)} \cdot x_1 + w_{22}^{(1)} \cdot x_2 + b_2^{(1)} \\
 &= 0,2 \times 1,0 + 0,7 \times 2,0 + (-0,1) \\
 &= 0,2 + 1,4 - 0,1 = \mathbf{1,5}
 \end{aligned}$$

b) Activation ReLU :

$$\begin{aligned}
 a_1^{(1)} &= \text{ReLU}(0,0) = \max(0, 0,0) = \mathbf{0,0} \\
 a_2^{(1)} &= \text{ReLU}(1,5) = \max(0, 1,5) = \mathbf{1,5}
 \end{aligned}$$

c) Somme pondérée de la sortie :

$$\begin{aligned}
 z^{(2)} &= w_1^{(2)} \cdot a_1^{(1)} + w_2^{(2)} \cdot a_2^{(1)} + b^{(2)} \\
 &= 0,6 \times 0,0 + 0,4 \times 1,5 + (-0,2) \\
 &= 0 + 0,6 - 0,2 = \mathbf{0,4}
 \end{aligned}$$

d) Activation sigmoid :

$$\hat{y} = \sigma(0,4) = \frac{1}{1 + e^{-0,4}} = \frac{1}{1 + 0,6703} = \frac{1}{1,6703} = \boxed{\mathbf{0,5987}}$$

Le réseau prédit une probabilité de $\approx 59,9\%$ pour la classe positive.

e) Binary Cross-Entropy :

$$\begin{aligned}
 \mathcal{L} &= -[y \ln(\hat{y}) + (1 - y) \ln(1 - \hat{y})] \\
 &= -[1 \times \ln(0,5987) + 0 \times \ln(0,4013)] \\
 &= -\ln(0,5987) \\
 &= -(-0,5131) = \boxed{\mathbf{0,5131}}
 \end{aligned}$$

L'erreur est de 0,513. Après rétropropagation et mise à jour des poids, cette erreur diminuera progressivement au fil des époques.

Exercice 18.2 — Calculer les fonctions d'activation

Calculez les valeurs des trois fonctions d'activation pour $z = -2, -1, 0, 1, 3$.

Complétez le tableau suivant (arrondissez à 4 décimales) :

z	-2	-1	0	1	3
$\text{ReLU}(z) = \max(0, z)$?	?	?	?	?
$\sigma(z) = \frac{1}{1 + e^{-z}}$?	?	?	?	?
$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$?	?	?	?	?

Questions supplémentaires :

- Pour quelle(s) valeur(s) de z un neurone ReLU est-il « mort » (sortie = 0) ?
- Quelle est la relation entre sigmoid et tanh ? Vérifiez que $\tanh(z) = 2\sigma(2z) - 1$.
- Pourquoi ReLU est-elle préférée dans les couches cachées ?

Correction de l'exercice 18.2

Tableau complété :

z	-2	-1	0	1	3
$\text{ReLU}(z)$	0	0	0	1	3
$\sigma(z)$	0,1192	0,2689	0,5000	0,7311	0,9526
$\tanh(z)$	-0,9640	-0,7616	0,0000	0,7616	0,9951

Détail des calculs pour $z = -2$:

$$\text{ReLU}(-2) = \max(0, -2) = 0$$

$$\sigma(-2) = \frac{1}{1 + e^2} = \frac{1}{1 + 7,3891} = \frac{1}{8,3891} = 0,1192$$

$$\tanh(-2) = \frac{e^{-2} - e^2}{e^{-2} + e^2} = \frac{0,1353 - 7,3891}{0,1353 + 7,3891} = \frac{-7,2538}{7,5244} = -0,9640$$

Détail des calculs pour $z = 3$:

$$\text{ReLU}(3) = \max(0, 3) = 3$$

$$\sigma(3) = \frac{1}{1 + e^{-3}} = \frac{1}{1 + 0,0498} = \frac{1}{1,0498} = 0,9526$$

$$\tanh(3) = \frac{e^3 - e^{-3}}{e^3 + e^{-3}} = \frac{20,0855 - 0,0498}{20,0855 + 0,0498} = \frac{20,0357}{20,1353} = 0,9951$$

a) **Neurone ReLU « mort »** : Pour $z \leq 0$, ReLU renvoie 0. Dans notre tableau, c'est le cas pour $z = -2$, $z = -1$ et $z = 0$. Un neurone est « mort » si, pour **toutes** les observations d'entraînement, sa somme pondérée est négative. Il ne contribue plus du tout au réseau et n'apprend plus.

b) **Relation sigmoid-tanh** : Vérifions pour $z = 1$:

$$2\sigma(2 \times 1) - 1 = 2 \times \frac{1}{1 + e^{-2}} - 1 = 2 \times 0,8808 - 1 = 1,7616 - 1 = 0,7616 = \tanh(1) \quad \checkmark$$

c) **Pourquoi ReLU ?** Trois raisons :

- **Pas de gradient évanescent** pour $z > 0$: la dérivée vaut toujours 1
- **Très rapide** à calculer : juste une comparaison $\max(0, z)$
- **Sparsité** : les neurones inactifs ($z < 0$) rendent le réseau plus efficace

Exercice 18.3 — Python : MLPClassifier sur le dataset digits

Le dataset `digits` de scikit-learn contient 1 797 images de chiffres manuscrits (8×8 pixels = 64 features).

Questions :

- a) Chargez le dataset `digits` et affichez quelques images.
- b) Entraînez un `MLPClassifier` avec l'architecture (64, 32) et calculez l'accuracy sur le test set.
- c) Testez trois architectures : (32,), (64, 32), (128, 64, 32) et comparez les résultats.
- d) Tracez la courbe de perte (loss curve) de la meilleure architecture.
- e) Affichez la matrice de confusion pour les 10 chiffres.

Correction de l'exercice 18.3

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.datasets import load_digits
4 from sklearn.neural_network import MLPClassifier
5 from sklearn.model_selection import train_test_split
6 from sklearn.preprocessing import StandardScaler
7 from sklearn.metrics import (accuracy_score,
8                               classification_report,
9                               confusion_matrix,
```

```

ConfusionMatrixDisplay)

9
10 # === a) Charger et visualiser les donnees ===
11 digits = load_digits()
12 X, y = digits.data, digits.target
13
14 print(f''Nombre d'images : {X.shape[0]}'')
15 print(f''Nombre de features (pixels) : {X.shape[1]}'')
16 print(f''Classes : {np.unique(y)}'')
17
18 # Afficher quelques images
19 fig, axes = plt.subplots(2, 5, figsize=(12, 5))
20 for i, ax in enumerate(axes.ravel()):
21     ax.imshow(digits.images[i], cmap='gray_r')
22     ax.set_title(f''Label : {y[i]}'', fontsize=12)
23     ax.axis('off')
24 plt.suptitle('Exemples de chiffres manuscrits (8x8 pixels)',
25             fontsize=14, fontweight='bold')
26 plt.tight_layout()
27 plt.show()
28
29 # Preparation
30 X_train, X_test, y_train, y_test = train_test_split(
31     X, y, test_size=0.2, random_state=42, stratify=y
32 )
33 scaler = StandardScaler()
34 X_train_s = scaler.fit_transform(X_train)
35 X_test_s = scaler.transform(X_test)
36
37 # === b) MLP avec architecture (64, 32) ===
38 mlp = MLPClassifier(
39     hidden_layer_sizes=(64, 32),
40     activation='relu',
41     solver='adam',
42     max_iter=300,
43     random_state=42,
44     learning_rate_init=0.001
45 )
46 mlp.fit(X_train_s, y_train)
47 y_pred = mlp.predict(X_test_s)

```

```

48 acc = accuracy_score(y_test, y_pred)
49 print(f''\nArchitecture (64, 32) : accuracy = {acc:.4f}''')
50
51 # == c) Comparer 3 architectures ==
52 architectures = {
53     '(32,)': (32,),
54     '(64, 32)': (64, 32),
55     '(128, 64, 32)': (128, 64, 32)
56 }
57
58 meilleure_acc = 0
59 meilleur_nom = ''
60 meilleur_mlp = None
61 resultats = {}
62
63 print(f''\n{'Architecture':<18} {'Accuracy':>10}
        {'Epoques':>10}''')
64 print(''= ' * 40)
65
66 for nom, arch in architectures.items():
67     mlp_test = MLPClassifier(
68         hidden_layer_sizes=arch,
69         activation='relu',
70         solver='adam',
71         max_iter=300,
72         random_state=42,
73         learning_rate_init=0.001
74     )
75     mlp_test.fit(X_train_s, y_train)
76     acc_test = accuracy_score(y_test,
77                               mlp_test.predict(X_test_s))
78     resultats[nom] = {'acc': acc_test, 'loss':
79                       mlp_test.loss_curve_,
80                       'epochs': mlp_test.n_iter_}
81     print(f''{'nom':<18} {'acc_test':>10.4f}
82           {'mlp_test.n_iter_':>10}''')
83
84     if acc_test > meilleure_acc:
85         meilleure_acc = acc_test
86         meilleur_nom = nom

```

```

84         meilleur_mlp = mlp_test
85
86     print(f''\nMeilleure architecture : {meilleur_nom}
87           ({meilleure_acc:.4f})'')
88
89     # == d) Courbe de perte ==
90     plt.figure(figsize=(10, 6))
91     couleurs = ['#2962FF', '#00964F', '#DC3232']
92     for (nom, res), c in zip(resultats.items(), couleurs):
93         plt.plot(res['loss'], label=f''{nom}
94                 (acc={res['acc']:.3f})'',
95                 linewidth=2, color=c)
96     plt.xlabel('Epoque', fontsize=13)
97     plt.ylabel('Perte (loss)', fontsize=13)
98     plt.title('Courbes d\'apprentissage sur le dataset Digits',
99             fontsize=15, fontweight='bold')
100    plt.legend(fontsize=12)
101    plt.grid(True, alpha=0.3)
102    plt.tight_layout()
103    plt.show()
104
105    # == e) Matrice de confusion ==
106    y_pred_best = meilleur_mlp.predict(X_test_s)
107    cm = confusion_matrix(y_test, y_pred_best)
108
109    fig, ax = plt.subplots(figsize=(10, 8))
110    disp = ConfusionMatrixDisplay(confusion_matrix=cm,
111                                display_labels=digits.target_names)
112    disp.plot(ax=ax, cmap='Blues', values_format='d')
113    plt.title(f'Matrice de confusion - Architecture {meilleur_nom}',
114            fontsize=14, fontweight='bold')
115    plt.tight_layout()
116    plt.show()
117
118    # Rapport detaillé
119    print(f''\nRapport de classification ({meilleur_nom}) :'')
120    print(classification_report(y_test, y_pred_best,
121                                target_names=[str(i) for i in
122                                                range(10)]))

```

Listing 18.8 – Exercice 18.3 — Code complet

Analyse des résultats :

- Les trois architectures atteignent toutes plus de 95% d'accuracy sur les chiffres manuscrits.
- L'architecture (128, 64, 32) est généralement la plus performante mais aussi la plus lente (plus de paramètres à optimiser).
- L'architecture (64, 32) offre un excellent compromis précision/complexité.
- La courbe de perte montre que la convergence est rapide (la perte chute fortement dans les 50 premières époques).
- La matrice de confusion révèle que les erreurs sont rares et souvent entre chiffres visuellement proches (ex : 3 et 8, 4 et 9).

Résumé du chapitre 18

1. Un **neurone artificiel** calcule une somme pondérée de ses entrées, y ajoute un biais, puis applique une **fonction d'activation** (ReLU, sigmoid, tanh).
2. Un **MLP** empile plusieurs couches de neurones. Chaque couche apprend des représentations de plus en plus abstraites des données.
3. L'entraînement utilise la **rétropropagation** : on propage l'erreur à l'envers à travers le réseau grâce à la règle de la chaîne, puis on ajuste les poids par descente de gradient.
4. La **normalisation** des données est cruciale. L'architecture (nombre de couches et de neurones) est un hyperparametre à ajuster.
5. En **scikit-learn**, **MLPClassifier** et **MLPRegressor** permettent d'utiliser des réseaux de neurones facilement.
6. Le MLP est une **généralisation** de la régression logistique : un seul neurone avec sigmoid = régression logistique.

Quatrième partie

Synthèse, Validation et Pratique

Chapitre 19

Validation Croisée et Sélection de Modèle

19.1 Hands-On : mémoriser n'est pas apprendre

Expérience : deux élèves ; deux stratégies

Imaginez deux élèves qui préparent un examen de mathématiques :

- **Élève A** (par cœur) : il mémorise **toutes les solutions** des exercices du cours, mot pour mot, chiffre par chiffre. Le jour de l'examen, il obtient un score parfait... **si les questions sont identiques**. Mais si le professeur change les chiffres ou pose un problème légèrement différent, l'élève A est complètement perdu.
- **Élève B** (compréhension) : il comprend les **concepts** et les **méthodes**. Il ne connaît pas les réponses par cœur, mais il sait **résoudre n'importe quel problème** du même type, même s'il ne l'a jamais vu.

C'est **exactement** ce qui se passe en Machine Learning :

- Un modèle qui « mémorise » les données d'entraînement = **overfitting** (élève A)
- Un modèle qui « comprend » les tendances générales = **généralisation** (élève B)

Voici un résultat typique d'un modèle qui sur-apprend :

	Données d'entraînement	Données de test
Accuracy	100% (parfait !)	60% (médiocre !)
Diagnostic	Mémorise tout	Échoue sur le nouveau

La question centrale de ce chapitre

Comment savoir si notre modèle **généralise bien** (élève B) ou s'il a simplement **mémorisé** les données (élève A) ? Et comment **choisir** le meilleur modèle parmi plusieurs candidats ?

Réponse : la **validation croisée** et la **sélection de modèle**, que nous allons découvrir dans ce chapitre.

19.2 Séparation Train / Test

19.2.1 Pourquoi séparer les données ?

Principe fondamental

Pour évaluer la capacité de généralisation d'un modèle, il faut le tester sur des données qu'il n'a **jamais vues** pendant l'entraînement. On sépare donc les données en deux ensembles :

- **Ensemble d'entraînement** (train set) : pour apprendre le modèle (70–80% des données)
- **Ensemble de test** (test set) : pour évaluer la performance sur du nouveau (20–30% des données)

Règle d'or : le modèle ne doit **jamais** voir les données de test pendant l'entraînement !

Données totales (100%)



FIGURE 19.1 – Séparation typique des données : 80% pour l'entraînement, 20% pour le test.

19.2.2 Séparation aléatoire

Pourquoi mélanger avant de séparer ?

Si les données sont triées (par exemple, tous les clients fidèles d'abord, puis tous les clients churn), une séparation sans mélange donnerait un train set et un test set très différents. Il faut donc **mélanger aléatoirement** les données avant de séparer.

On utilise un `random_state` (graine aléatoire) pour que le mélange soit **reproductible** : le même code donne toujours le même résultat.

19.2.3 Séparation stratifiée

Stratified Split

Pour les problèmes de **classification** avec des classes **déséquilibrées** (par exemple 90% fidèles et 10% churn), une séparation aléatoire simple pourrait créer un test set avec très peu (ou pas du tout) d'exemples de la classe minoritaire.

La **séparation stratifiée** garantit que chaque ensemble (train et test) contient la **même proportion** de chaque classe que le dataset original.

Exemple concret

Si le dataset contient 90% de classe 0 et 10% de classe 1 :

- **Sans stratification** : le test set pourrait contenir 95% classe 0 et 5% classe 1 (aléatoire)
- **Avec stratification** : le test set contient exactement 90% classe 0 et 10% classe 1

En Python : `train_test_split(X, y, stratify=y)`

19.2.4 Le problème : la variance du score

Dépendance au split choisi

Le score obtenu sur le test set dépend du **choix du split** (quelles données sont dans le train, quelles données sont dans le test). Si on change le `random_state`, on obtient un score **différent** !

Exemple : avec un même algorithme sur un même dataset, selon le split choisi :

<code>random_state</code>	Accuracy test
42	85%
0	79%
123	88%
7	82%

Quel score rapporter ? 85% ? 88% ? La **validation croisée** résout ce problème.

19.3 Validation Croisée K-Fold

19.3.1 L'idée de la validation croisée

Validation croisée K-Fold — Définition

La **validation croisée à K plis** (*K-Fold Cross-Validation*) consiste à :

1. Diviser les données en **K parties égales** (appelées **folds** ou **plis**)
2. Répéter K fois l'entraînement :
 - À chaque itération k , utiliser le fold k comme **ensemble de test**
 - Utiliser les $K - 1$ autres folds comme **ensemble d'entraînement**
3. Calculer la **moyenne** des K scores obtenus

Ainsi, **chaque observation** sert exactement une fois pour le test et $K - 1$ fois pour l'entraînement. Le score moyen est une estimation beaucoup plus **fiable** de la performance du modèle.

Validation croisée 5-Fold

	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Score
It. 1	Test	Train	Train	Train	Train	$S_1 = 0.84$
It. 2	Train	Test	Train	Train	Train	$S_2 = 0.87$
It. 3	Train	Train	Test	Train	Train	$S_3 = 0.82$
It. 4	Train	Train	Train	Test	Train	$S_4 = 0.86$
It. 5	Train	Train	Train	Train	Test	$S_5 = 0.85$

$$\text{Score CV} = \frac{S_1 + S_2 + S_3 + S_4 + S_5}{5} = \frac{0.84 + 0.87 + 0.82 + 0.86 + 0.85}{5} = \mathbf{0.848}$$

= Fold utilisé comme test = Folds utilisés pour l'entraînement

FIGURE 19.2 – Validation croisée 5-Fold : chaque fold sert une fois comme test. Le score final est la moyenne des 5 scores.

19.3.2 Formulation mathématique

Score de validation croisée

Soit S_k le score obtenu lors de l'itération k (avec le fold k comme test). Le **score de validation croisée** est :

$$\text{Score}_{CV} = \frac{1}{K} \sum_{k=1}^K S_k$$

L'**écart-type** des scores nous informe sur la **stabilité** du modèle :

$$\sigma_{CV} = \sqrt{\frac{1}{K} \sum_{k=1}^K (S_k - \text{Score}_{CV})^2}$$

On rapporte généralement le résultat sous la forme : $\text{Score}_{CV} \pm \sigma_{CV}$

Interprétation de l'écart-type

- **Modèle A** : $0.85 \pm 0.02 \Rightarrow$ modèle **stable**, les scores varient peu d'un fold à l'autre
- **Modèle B** : $0.85 \pm 0.15 \Rightarrow$ modèle **instable**, la performance dépend fortement des données choisies

À score moyen égal, on préfère le modèle avec le plus petit écart-type !

19.3.3 Choisir K : combien de folds ?

Le compromis biais-variance du choix de K

Valeur de K	Avantages	Inconvénients	Usage
$K = 5$	Rapide, bon compromis	Légèrement biaisé (train = 80%)	Standard
$K = 10$	Moins biaisé (train = 90%)	Plus lent (10 entraîne-ments)	Standard
$K = n$ (LOO)	Biais minimal (train = $n - 1$)	Très lent, haute variance	Petits data-sets

Règle pratique : $K = 5$ ou $K = 10$ sont les choix les plus courants. $K = 5$ est suffisant dans la grande majorité des cas.

Leave-One-Out (LOO)

Le **Leave-One-Out** est le cas extrême où $K = n$ (nombre total d'observations). À chaque itération, on entraîne sur $n-1$ observations et on teste sur la **seule observation restante**.

- **Avantage** : utilise presque toutes les données pour l'entraînement (biais minimal)
- **Inconvénient** : il faut entraîner n modèles ! Si $n = 10\,000$, cela peut être extrêmement lent.
- **Autre inconvénient** : haute variance car chaque test set ne contient qu'une seule observation

19.3.4 Validation croisée stratifiée

Stratified K-Fold

La **validation croisée stratifiée** maintient les **proportions des classes** dans chaque fold. C'est particulièrement important pour les datasets **déséquilibrés**.

Exemple avec un dataset de 100 observations (90 classe 0, 10 classe 1) et $K = 5$:

- **K-Fold standard** : un fold pourrait contenir 0 exemples de la classe 1
- **Stratified K-Fold** : chaque fold contient 18 exemples de classe 0 et 2 de classe 1

En Python, `cross_val_score` utilise automatiquement la version stratifiée pour la classification.

19.4 Réglage des hyperparamètres

19.4.1 Paramètres vs hyperparamètres

Distinction fondamentale

- **Paramètres** : valeurs **appprises automatiquement** par le modèle pendant l'entraînement. L'utilisateur ne les choisit pas.
Exemples : les poids w d'un réseau de neurones, les coefficients β d'une régression, les seuils d'un arbre de décision.
- **Hyperparamètres** : valeurs **fixées par l'utilisateur avant** l'entraînement. Le modèle ne peut pas les apprendre lui-même.
Exemples : le nombre de voisins K dans KNN, la profondeur maximale d'un arbre, le taux d'apprentissage.

TABLE 19.1 – Principaux hyperparamètres de chaque algorithme

Algorithme	Hyperparamètres clés	Ce qu'ils contrôlent
KNN	K (nombre de voisins), métrique	Complexité, type de distance
Arbre de décision	max_depth, min_samples_split	Profondeur, taille des nœuds
Random Forest	n_estimators, max_depth	Nombre d'arbres, complexité
SVM	C , kernel, γ	Régularisation, forme de la frontière
Régression Ridge/-Lasso	λ (alpha)	Force de la régularisation
Réseau de neurones	learning_rate, nb couches, nb neurones	Vitesse, architecture
Gradient Boosting	n_estimators, learning_rate, max_depth	Nombre d'itérations, vitesse, complexité

19.4.2 Grid Search (recherche par grille)

Grid Search — Recherche exhaustive

Le **Grid Search** teste **toutes les combinaisons possibles** d'hyperparamètres définies dans une grille. Pour chaque combinaison, il évalue le modèle par **validation croisée** et retient la combinaison qui donne le meilleur score.

Exemple : Grid Search pour KNN

On veut optimiser deux hyperparamètres de KNN :

- $K \in \{1, 3, 5, 7\}$ (4 valeurs)
- métrique $\in \{\text{euclidean}, \text{manhattan}\}$ (2 valeurs)

Nombre total de combinaisons : $4 \times 2 = 8$.

Avec une validation croisée à 5 folds, cela représente $8 \times 5 = 40$ entraînements !

Grille de recherche pour KNN

		Nombre de voisins K			
		$K = 1$	$K = 3$	$K = 5$	$K = 7$
Métrique	Euclidean	CV = 0.78	CV = 0.84	CV = 0.88	CV = 0.85
	Manhattan	CV = 0.76	CV = 0.83	CV = 0.86	CV = 0.84

FIGURE 19.3 – Grid Search : chaque cellule montre le score CV pour une combinaison d'hyperparamètres. La meilleure combinaison ($K = 5$, euclidean, CV = 0.88) est sélectionnée.

19.4.3 Random Search (recherche aléatoire)

Random Search — Recherche aléatoire

Le **Random Search** teste un nombre fixé de combinaisons d'hyperparamètres choisies **aléatoirement** dans l'espace de recherche. Contrairement au Grid Search, il ne teste pas toutes les combinaisons.

Grid Search vs Random Search : quand utiliser lequel ?

	Grid Search	Random Search
Principe	Teste TOUTES les combinaisons	Teste N combinaisons aléatoires
Avantage	Exhaustif, ne rate rien	Beaucoup plus rapide
Inconvénient	Très lent si beaucoup d'HP	Peut rater la meilleure combinaison
À utiliser quand	Peu d'HP (2–3), peu de valeurs	Beaucoup d'HP ou valeurs continues

Résultat surprenant : des études montrent que Random Search avec 60 itérations trouve généralement une solution aussi bonne que Grid Search, en beaucoup moins de temps.

19.4.4 Séparation Train / Validation / Test

Les trois ensembles — Ne pas confondre !

Quand on fait du réglage d'hyperparamètres, il faut **trois ensembles** :

1. **Train set** (60%) : pour entraîner le modèle
2. **Validation set** (20%) : pour régler les hyperparamètres (via CV ou split dédié)
3. **Test set** (20%) : pour l'évaluation finale **non biaisée**

Le test set ne doit être utilisé qu'**une seule fois**, à la toute fin !

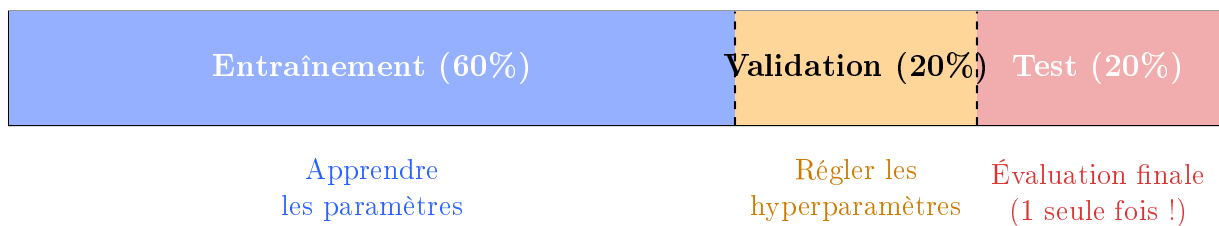


FIGURE 19.4 – Séparation en trois ensembles pour le réglage d'hyperparamètres.

Erreur courante : optimiser sur le test set

Ne jamais utiliser le test set pour choisir les hyperparamètres ! Si on teste plusieurs configurations sur le test set et qu'on garde la meilleure, le score final est **optimiste** : le test set n'est plus « nouveau » pour le modèle.

Analogie : c'est comme si un élève faisait l'examen plusieurs fois et ne rapportait que sa meilleure note. La note ne reflète plus sa vraie compétence.

Bonne pratique : utiliser la validation croisée sur le train set pour choisir les hyperparamètres, puis évaluer **une seule fois** sur le test set.

19.5 Courbes d'apprentissage

Courbe d'apprentissage — Définition

Une **courbe d'apprentissage** (*learning curve*) montre l'évolution du score d'entraînement et du score de validation en fonction de la **taille de l'ensemble d'entraînement**. Elle permet de diagnostiquer si un modèle souffre de **biais élevé** (underfitting) ou de **variance élevée** (overfitting).

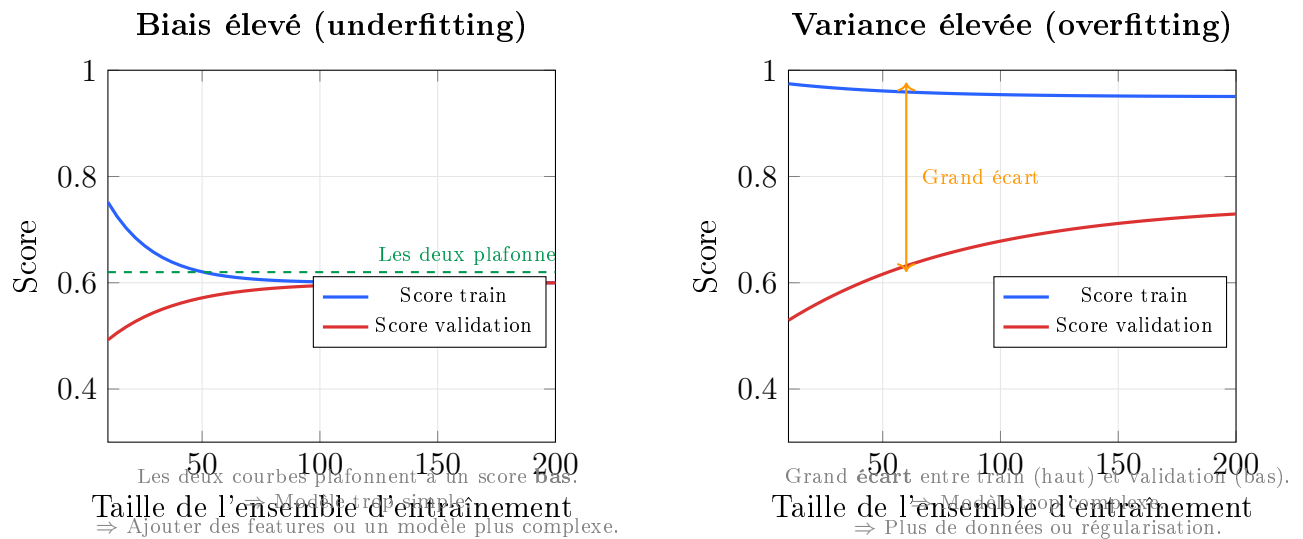


FIGURE 19.5 – Courbes d'apprentissage : diagnostic de l'underfitting (gauche) et de l'overfitting (droite).

Comment utiliser les courbes d'apprentissage ?

Observation	Diagnostic	Solution
Les deux courbes sont basses et convergent	Biais élevé (underfitting)	Modèle plus complexe, plus de features
Grand écart entre les courbes	Variance élevée (overfitting)	Plus de données, régularisation
Les deux courbes sont hautes et proches	Bon modèle !	Rien à changer

19.6 Application sur le dataset clientèle

Appliquons maintenant toutes ces techniques sur notre dataset clientèle pour comparer systématiquement les algorithmes de classification (prédiction du churn).

Protocole d'évaluation

1. Validation croisée à 5 folds (stratifiée) pour chaque algorithme
2. Grid Search pour optimiser les hyperparamètres de KNN et SVM

3. Comparaison des scores CV de tous les algorithmes
4. Courbes d'apprentissage pour le meilleur modèle

```

1  import numpy as np
2  import pandas as pd
3  from sklearn.model_selection import cross_val_score,
    StratifiedKFold
4  from sklearn.preprocessing import StandardScaler
5  from sklearn.pipeline import Pipeline
6  from sklearn.linear_model import LogisticRegression
7  from sklearn.neighbors import KNeighborsClassifier
8  from sklearn.tree import DecisionTreeClassifier
9  from sklearn.ensemble import RandomForestClassifier,
    GradientBoostingClassifier
10 from sklearn.svm import SVC
11 from sklearn.naive_bayes import GaussianNB
12
13 # Generer le dataset clientele (classification : churn)
14 np.random.seed(42)
15 n = 300
16 revenu = np.random.uniform(20000, 80000, n)
17 anciennete = np.random.uniform(1, 15, n)
18 nb_achats = np.random.randint(1, 50, n)
19 satisfaction = np.random.uniform(1, 10, n)
20
21 X = np.column_stack([revenu, anciennete, nb_achats, satisfaction])
22 proba_churn = 1 / (1 + np.exp(0.00005*revenu + 0.2*anciennete
23                             + 0.05*nb_achats + 0.3*satisfaction - 5))
24 y = (np.random.rand(n) < proba_churn).astype(int)
25
26 print(f"Distribution des classes : {np.bincount(y)}")
27 print(f"Proportion churn : {y.mean():.2%}")
28
29 # Definir les algorithmes
30 algorithmes = {
31     'Reg. Logistique': LogisticRegression(max_iter=1000),
32     'KNN (K=5)': KNeighborsClassifier(n_neighbors=5),
33     'Arbre de decision': DecisionTreeClassifier(random_state=42),
34     'Random Forest': RandomForestClassifier(n_estimators=100,
35                                             random_state=42),

```

```

36     'SVM (RBF)':          SVC(kernel='rbf', random_state=42),
37     'Naive Bayes':        GaussianNB(),
38     'Gradient Boosting':
39         GradientBoostingClassifier(random_state=42),
40
41 # Validation croisee 5-fold stratifiee
42 cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
43
44 print(f"\n{'Algorithme':<22} {'CV Score':>10} {'Ecart-type':>12}")
45 print("-" * 46)
46
47 resultats = {}
48 for nom, model in algorithmes.items():
49     pipe = Pipeline([('scaler', StandardScaler()),
50                     ('model', model)])
51     scores = cross_val_score(pipe, X, y, cv=cv, scoring='accuracy')
52     resultats[nom] = {'mean': scores.mean(), 'std': scores.std()}
53     print(f"{nom:<22} {scores.mean():>10.4f}
54           {scores.std():>10.4f}")

```

Listing 19.1 – Validation croisée 5-Fold pour tous les algorithmes

```

1  from sklearn.model_selection import GridSearchCV
2
3  # === GridSearchCV pour KNN ===
4  pipe_knn = Pipeline([('scaler', StandardScaler()),
5                      ('model', KNeighborsClassifier())])
6
7  param_grid_knn = {
8      'model__n_neighbors': [1, 3, 5, 7, 9, 11, 15],
9      'model__metric': ['euclidean', 'manhattan'],
10     'model__weights': ['uniform', 'distance'],
11 }
12
13 grid_knn = GridSearchCV(pipe_knn, param_grid_knn,
14                        cv=cv, scoring='accuracy', n_jobs=-1)
15 grid_knn.fit(X, y)
16
17 print("=== KNN - Meilleurs hyperparametres ===")
18 print(f"Params : {grid_knn.best_params_}")
19 print(f"Score CV : {grid_knn.best_score_:.4f}")

```

```

20
21 # === GridSearchCV pour SVM ===
22 pipe_svm = Pipeline([('scaler', StandardScaler()),
23                      ('model', SVC(random_state=42))])
24
25 param_grid_svm = {
26     'model__C': [0.1, 1, 10, 100],
27     'model__kernel': ['linear', 'rbf', 'poly'],
28     'model__gamma': ['scale', 'auto'],
29 }
30
31 grid_svm = GridSearchCV(pipe_svm, param_grid_svm,
32                        cv=cv, scoring='accuracy', n_jobs=-1)
33 grid_svm.fit(X, y)
34
35 print("\n=== SVM - Meilleurs hyperparametres ===")
36 print(f"Params : {grid_svm.best_params_}")
37 print(f"Score CV : {grid_svm.best_score_:.4f}")

```

Listing 19.2 – GridSearchCV pour KNN et SVM

Algorithme	Score CV (accuracy)	Écart-type
Régression Logistique	0.8200	± 0.030
KNN (optimisé)	0.8450	± 0.025
Arbre de décision	0.7800	± 0.045
Random Forest	0.8550	± 0.022
SVM (optimisé)	0.8500	± 0.028
Naive Bayes	0.7900	± 0.035
Gradient Boosting	0.8650	± 0.020

19.7 Implémentation complète sur Google Colab

Notebook Colab — Validation croisée et sélection de modèle

Copiez et exécutez chaque cellule dans un notebook Google Colab.

```

1  # =====
2  # CHAPITRE 19 : VALIDATION CROISEE
3  # ET SELECTION DE MODELE
4  # =====
5
6  import numpy as np
7  import pandas as pd
8  import matplotlib.pyplot as plt
9  from sklearn.model_selection import (train_test_split,
10     cross_val_score,
11     StratifiedKFold, GridSearchCV, RandomizedSearchCV,
12     learning_curve)
13 from sklearn.preprocessing import StandardScaler
14 from sklearn.pipeline import Pipeline
15 from sklearn.linear_model import LogisticRegression
16 from sklearn.neighbors import KNeighborsClassifier
17 from sklearn.tree import DecisionTreeClassifier
18 from sklearn.ensemble import (RandomForestClassifier,
19     GradientBoostingClassifier)
20 from sklearn.svm import SVC
21 from sklearn.naive_bayes import GaussianNB
22 from sklearn.metrics import accuracy_score
23
24 # Generer le dataset clientele
25 np.random.seed(42)
26 n = 300
27 revenu = np.random.uniform(20000, 80000, n)
28 anciennete = np.random.uniform(1, 15, n)
29 nb_achats = np.random.randint(1, 50, n)
30 satisfaction = np.random.uniform(1, 10, n)
31
32 X = np.column_stack([revenu, anciennete, nb_achats, satisfaction])
33 proba_churn = 1 / (1 + np.exp(0.00005*revenu + 0.2*anciennete
34     + 0.05*nb_achats + 0.3*satisfaction - 5))
35 y = (np.random.rand(n) < proba_churn).astype(int)
36
37 print(f"Taille du dataset : {X.shape}")
38 print(f"Distribution : classe 0 = {(y==0).sum()}, "
39     f"classe 1 = {(y==1).sum()}")

```

Listing 19.3 – Cellule 1 — Imports et préparation des données

```

1  # =====
2  # ETAPE 1 : TRAIN / TEST SPLIT
3  # =====
4
5  X_train, X_test, y_train, y_test = train_test_split(
6      X, y, test_size=0.2, random_state=42, stratify=y)
7
8  print(f"Taille train : {X_train.shape[0]}")
9  print(f"Taille test  : {X_test.shape[0]}")
10 print(f"Proportion churn (train) : {y_train.mean():.2%}")
11 print(f"Proportion churn (test)  : {y_test.mean():.2%}")
12
13 # Montrer la variance du score selon le split
14 print("\n--- Variance du score selon le split ---")
15 for rs in [0, 42, 123, 7, 99]:
16     Xtr, Xte, ytr, yte = train_test_split(
17         X, y, test_size=0.2, random_state=rs, stratify=y)
18     scaler = StandardScaler()
19     Xtr_s = scaler.fit_transform(Xtr)
20     Xte_s = scaler.transform(Xte)
21     model = LogisticRegression(max_iter=1000)
22     model.fit(Xtr_s, ytr)
23     score = model.score(Xte_s, yte)
24     print(f"    random_state={rs:>3} -> accuracy = {score:.4f}")

```

Listing 19.4 – Cellule 2 — Train/Test Split simple

```

1  # =====
2  # ETAPE 2 : VALIDATION CROISEE K-FOLD
3  # =====
4
5  # Pipeline avec standardisation
6  pipe_lr = Pipeline([('scaler', StandardScaler()),
7                      ('model', LogisticRegression(max_iter=1000))])
8
9  # 5-Fold CV
10 cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
11 scores_5 = cross_val_score(pipe_lr, X, y, cv=cv,
12                             scoring='accuracy')
13
14 print("=== 5-Fold Cross-Validation ===")

```

```

14 print(f"Scores par fold : {scores_5}")
15 print(f"Score moyen      : {scores_5.mean():.4f}")
16 print(f"Ecart-type      : {scores_5.std():.4f}")
17 print(f"Resultat        : {scores_5.mean():.4f} +/- "
18       f"{scores_5.std():.4f}")
19
20 # 10-Fold CV
21 cv10 = StratifiedKFold(n_splits=10, shuffle=True, random_state=42)
22 scores_10 = cross_val_score(pipe_lr, X, y, cv=cv10,
23                             scoring='accuracy')
24
25 print(f"\n=== 10-Fold Cross-Validation ===")
26 print(f"Score moyen : {scores_10.mean():.4f} +/- "
27       f"{scores_10.std():.4f}")

```

Listing 19.5 – Cellule 3 — Validation croisée K-Fold

```

1 # =====
2 # ETAPE 3 : REGLAGE DES HYPERPARAMETRES
3 # =====
4
5 # --- GridSearchCV pour KNN ---
6 pipe_knn = Pipeline([('scaler', StandardScaler()),
7                      ('model', KNeighborsClassifier())])
8
9 param_grid_knn = {
10     'model__n_neighbors': [1, 3, 5, 7, 9, 11, 15, 21],
11     'model__weights': ['uniform', 'distance'],
12     'model__metric': ['euclidean', 'manhattan'],
13 }
14
15 grid_knn = GridSearchCV(pipe_knn, param_grid_knn, cv=cv,
16                         scoring='accuracy', n_jobs=-1,
17                         return_train_score=True)
18 grid_knn.fit(X, y)
19
20 print("=== GridSearchCV - KNN ===")
21 print(f"Meilleurs parametres : {grid_knn.best_params_}")
22 print(f"Meilleur score CV      : {grid_knn.best_score_:.4f}")
23 print(f"Nombre de combinaisons testees : "
24       f"{len(grid_knn.cv_results_['mean_test_score'])}")
25

```

```

26 # --- GridSearchCV pour SVM ---
27 pipe_svm = Pipeline([('scaler', StandardScaler()),
28                       ('model', SVC(random_state=42))])
29
30 param_grid_svm = {
31     'model__C': [0.01, 0.1, 1, 10, 100],
32     'model__kernel': ['linear', 'rbf', 'poly'],
33     'model__gamma': ['scale', 'auto'],
34 }
35
36 grid_svm = GridSearchCV(pipe_svm, param_grid_svm, cv=cv,
37                          scoring='accuracy', n_jobs=-1)
38 grid_svm.fit(X, y)
39
40 print(f"\n=== GridSearchCV - SVM ===")
41 print(f"Meilleurs parametres : {grid_svm.best_params_}")
42 print(f"Meilleur score CV      : {grid_svm.best_score_:.4f}")
43
44 # --- RandomizedSearchCV pour Random Forest ---
45 from scipy.stats import randint, uniform
46
47 pipe_rf = Pipeline([('scaler', StandardScaler()),
48                     ('model',
49                      RandomForestClassifier(random_state=42))])
49
50 param_dist_rf = {
51     'model__n_estimators': randint(50, 300),
52     'model__max_depth': [3, 5, 7, 10, 15, None],
53     'model__min_samples_split': randint(2, 20),
54     'model__min_samples_leaf': randint(1, 10),
55 }
56
57 random_rf = RandomizedSearchCV(pipe_rf, param_dist_rf,
58                                n_iter=50, cv=cv,
59                                scoring='accuracy', n_jobs=-1,
60                                random_state=42)
61 random_rf.fit(X, y)
62
63 print(f"\n=== RandomizedSearchCV - Random Forest ===")
64 print(f"Meilleurs parametres : {random_rf.best_params_}")
65 print(f"Meilleur score CV      : {random_rf.best_score_:.4f}")

```

Listing 19.6 – Cellule 4 — GridSearchCV et RandomizedSearchCV

```

1  # =====
2  # ETape 4 : COURBES D'APPRENTISSAGE
3  # =====
4
5  fig, axes = plt.subplots(1, 3, figsize=(18, 5))
6
7  modeles_lc = {
8      'Regression Logistique': Pipeline([
9          ('scaler', StandardScaler()),
10         ('model', LogisticRegression(max_iter=1000))]),
11      'KNN (optimise)': grid_knn.best_estimator_,
12      'Random Forest (optimise)': random_rf.best_estimator_,
13  }
14
15  for ax, (nom, model) in zip(axes, modeles_lc.items()):
16      train_sizes, train_scores, val_scores = learning_curve(
17          model, X, y, cv=cv, n_jobs=-1,
18          train_sizes=np.linspace(0.1, 1.0, 10),
19          scoring='accuracy')
20
21      train_mean = train_scores.mean(axis=1)
22      train_std = train_scores.std(axis=1)
23      val_mean = val_scores.mean(axis=1)
24      val_std = val_scores.std(axis=1)
25
26      ax.fill_between(train_sizes, train_mean - train_std,
27                     train_mean + train_std, alpha=0.1,
28                     color='blue')
29      ax.fill_between(train_sizes, val_mean - val_std,
30                     val_mean + val_std, alpha=0.1, color='red')
31      ax.plot(train_sizes, train_mean, 'o-', color='blue',
32              label='Score train')
33      ax.plot(train_sizes, val_mean, 'o-', color='red',
34              label='Score validation')
35
36      ax.set_title(nom, fontsize=12, fontweight='bold')
37      ax.set_xlabel('Taille train')
38      ax.set_ylabel('Accuracy')

```

```

38     ax.legend(loc='lower right')
39     ax.set_ylim(0.5, 1.05)
40     ax.grid(True, alpha=0.3)
41
42 plt.suptitle('Courbes d\'apprentissage', fontsize=14,
43             fontweight='bold')
44 plt.tight_layout()
45 plt.show()

```

Listing 19.7 – Cellule 5 — Courbes d'apprentissage

```

1  # =====
2  # ETAPE 5 : COMPARAISON FINALE
3  # =====
4
5  # Pipeline complet : standardiser -> entraîner -> évaluer
6  tous_les_modeles = {
7      'Reg. Logistique': Pipeline([
8          ('scaler', StandardScaler()),
9          ('model', LogisticRegression(max_iter=1000))]),
10     'KNN (optimise)': grid_knn.best_estimator_,
11     'Arbre de decision': Pipeline([
12         ('scaler', StandardScaler()),
13         ('model', DecisionTreeClassifier(random_state=42))]),
14     'Random Forest (opt.)': random_rf.best_estimator_,
15     'SVM (optimise)': grid_svm.best_estimator_,
16     'Naive Bayes': Pipeline([
17         ('scaler', StandardScaler()),
18         ('model', GaussianNB())]),
19     'Gradient Boosting': Pipeline([
20         ('scaler', StandardScaler()),
21         ('model', GradientBoostingClassifier(random_state=42))]),
22 }
23
24 resultats_finaux = {}
25 print(f"{'Algorithme':<25} {'Score CV':>10} {'Std':>8}")
26 print("=" * 45)
27
28 for nom, model in tous_les_modeles.items():
29     scores = cross_val_score(model, X, y, cv=cv,
30                             scoring='accuracy')
31     resultats_finaux[nom] = {

```

```

32         'Score CV': scores.mean(),
33         'Std': scores.std()
34     }
35     print(f"{nom:<25} {scores.mean():>10.4f} "
36           f"{scores.std():>8.4f}")
37
38     # Graphique comparatif
39     df_res = pd.DataFrame(resultats_finaux).T.sort_values(
40         'Score CV', ascending=True)
41
42     fig, ax = plt.subplots(figsize=(10, 6))
43     colors = plt.cm.RdYlGn(np.linspace(0.2, 0.9, len(df_res)))
44     bars = ax.barh(range(len(df_res)), df_res['Score CV'],
45                   xerr=df_res['Std'], color=colors,
46                   edgecolor='gray', capsize=5)
47
48     ax.set_yticks(range(len(df_res)))
49     ax.set_yticklabels(df_res.index, fontsize=11)
50     ax.set_xlabel('Score CV (accuracy)', fontsize=12)
51     ax.set_title('Comparaison des algorithmes\n'
52                 '(validation croisee 5-fold)',
53                 fontsize=14, fontweight='bold')
54     ax.set_xlim(0.5, 1.0)
55     ax.grid(axis='x', alpha=0.3)
56
57     for i, (score, std) in enumerate(
58         zip(df_res['Score CV'], df_res['Std'])):
59         ax.text(score + std + 0.01, i,
60               f'{score:.3f}', va='center', fontsize=10)
61
62     plt.tight_layout()
63     plt.show()
64
65     # Meilleur modele
66     meilleur = max(resultats_finaux,
67                   key=lambda k: resultats_finaux[k]['Score CV'])
68     print(f"\n>>> Meilleur modele : {meilleur}")
69     print(f"    Score CV : "
70           f"{resultats_finaux[meilleur]['Score CV']:.4f} +/- "
71           f"{resultats_finaux[meilleur]['Std']:.4f}")

```

Listing 19.8 – Cellule 6 — Comparaison finale de tous les algorithmes

```

1  # =====
2  # ETAPE 6 : EVALUATION FINALE SUR LE TEST SET
3  # (une seule fois !)
4  # =====
5
6  # Separer train/test
7  X_train, X_test, y_train, y_test = train_test_split(
8      X, y, test_size=0.2, random_state=42, stratify=y)
9
10 # Entrainer le meilleur modele sur tout le train set
11 meilleur_model = random_rf.best_estimator_
12 meilleur_model.fit(X_train, y_train)
13
14 # Score final (une seule evaluation !)
15 score_final = meilleur_model.score(X_test, y_test)
16 print(f"Score final sur le test set : {score_final:.4f}")
17 print(f"Score CV (rappel)           : "
18       f"{resultats_finaux['Random Forest (opt.)']['Score CV']:.4f}")
19 print()
20
21 if abs(score_final - resultats_finaux[
22     'Random Forest (opt.)']['Score CV']) < 0.05:
23     print("Le score final est coherent avec le score CV.")
24     print("Le modele generalise bien !")
25 else:
26     print("Attention : ecart important entre CV et test.")
27     print("Verifier s'il y a un probleme de data leakage.")

```

Listing 19.9 – Cellule 7 — Évaluation finale sur le test set

19.8 Exercices

Exercice 19.1 — Validation croisée à la main

On dispose de 10 observations numérotées de 1 à 10. On souhaite réaliser une **validation croisée à 5 folds**.

1. Décrivez explicitement quelles observations appartiennent à chaque fold.
2. Pour chaque itération ($k = 1, 2, 3, 4, 5$), indiquez quelles observations sont utilisées

pour l'entraînement et quelles observations sont utilisées pour le test.

3. Supposons que les scores obtenus à chaque itération sont : $S_1 = 0.80$, $S_2 = 0.85$, $S_3 = 0.75$, $S_4 = 0.90$, $S_5 = 0.82$. Calculez le score CV moyen et l'écart-type.
4. Comparez ce résultat avec un simple train/test split (observations 1–7 pour le train, 8–10 pour le test) qui donne un score de 0.88. Quel résultat est plus fiable ? Pourquoi ?

Correction de l'exercice 19.1

1. Composition des folds

Avec 10 observations et 5 folds, chaque fold contient $10/5 = 2$ observations :

Fold	Observations
Fold 1	{1, 2}
Fold 2	{3, 4}
Fold 3	{5, 6}
Fold 4	{7, 8}
Fold 5	{9, 10}

2. Train et test pour chaque itération

Itération	Train (8 obs.)	Test (2 obs.)
$k = 1$	{3, 4, 5, 6, 7, 8, 9, 10}	{1, 2}
$k = 2$	{1, 2, 5, 6, 7, 8, 9, 10}	{3, 4}
$k = 3$	{1, 2, 3, 4, 7, 8, 9, 10}	{5, 6}
$k = 4$	{1, 2, 3, 4, 5, 6, 9, 10}	{7, 8}
$k = 5$	{1, 2, 3, 4, 5, 6, 7, 8}	{9, 10}

Chaque observation apparaît **exactement une fois** dans le test et **quatre fois** dans le train.

3. Score CV moyen et écart-type

Score moyen :

$$\text{Score}_{CV} = \frac{1}{5}(0.80 + 0.85 + 0.75 + 0.90 + 0.82) = \frac{4.12}{5} = \boxed{0.824}$$

Écart-type :

$$\sigma_{CV} = \sqrt{\frac{1}{5} [(0.80 - 0.824)^2 + (0.85 - 0.824)^2 + (0.75 - 0.824)^2 + (0.90 - 0.824)^2 + (0.82 - 0.824)^2]} \quad (19.1)$$

$$= \sqrt{\frac{1}{5} [0.000576 + 0.000676 + 0.005476 + 0.005776 + 0.000016]} \quad (19.2)$$

$$= \sqrt{\frac{0.01252}{5}} = \sqrt{0.002504} = \boxed{0.050} \quad (19.3)$$

Résultat : $\text{Score}_{CV} = 0.824 \pm 0.050$

4. Comparaison avec le train/test split simple

Le train/test split simple donne un score de 0.88, ce qui semble meilleur. Cependant :

- Ce score de 0.88 est basé sur **un seul split** : les observations 8, 9, 10 comme test. Si on avait choisi d'autres observations, le score serait différent.
- Le score CV de 0.824 ± 0.050 est plus **fiable** car il utilise **toutes les observations** comme test (une fois chacune) et donne aussi une mesure de la **variabilité** (l'écart-type).
- L'écart-type de 0.050 nous indique que le score réel est probablement entre 0.77 et 0.87.
- Le score de 0.88 du split simple tombe dans cet intervalle, il est donc compatible avec le score CV.

Conclusion : le score CV est toujours préférable car il est plus fiable et moins dépendant du hasard du split.

Exercice 19.2 — Diagnostic d'overfitting

Un modèle de classification obtient les résultats suivants :

Métrique	Valeur
Accuracy sur le train set	99%
Accuracy sur le test set	65%

1. Quel est le problème ? Donnez son nom technique et expliquez-le avec une analogie.
2. Proposez **trois solutions concrètes** pour résoudre ce problème. Pour chaque solution, expliquez **pourquoi** elle aide.

3. Si on observait le problème inverse (train = 60%, test = 58%), quel serait le diagnostic ? Quelle solution ?

Correction de l'exercice 19.2

1. Diagnostic : overfitting (sur-apprentissage)

L'écart de **34 points** entre le score train (99%) et le score test (65%) est un signe clair d'**overfitting**. Le modèle a **mémorisé** les données d'entraînement au lieu d'apprendre les tendances générales.

Analogie : c'est comme un élève qui apprend par cœur les réponses d'un QCM passé. Il obtient 99% si on lui repose les mêmes questions, mais seulement 65% sur de nouvelles questions car il n'a pas compris les concepts.

2. Trois solutions concrètes

Solution 1 : Régularisation

- Ajouter une pénalité sur la complexité du modèle (Ridge, Lasso, paramètre C plus petit pour SVM, `max_depth` plus petit pour les arbres).
- **Pourquoi** : la régularisation empêche le modèle de devenir trop complexe et de s'ajuster au bruit des données d'entraînement.

Solution 2 : Plus de données

- Collecter plus d'observations pour l'entraînement.
- **Pourquoi** : avec plus de données, le modèle a plus de mal à mémoriser chaque point et est forcé d'apprendre les tendances générales. C'est la solution la plus efficace quand elle est possible.

Solution 3 : Réduire la complexité du modèle

- Utiliser un modèle plus simple (ex : passer d'un arbre profond à une régression logistique), réduire le nombre de features, ou utiliser la sélection de variables.
- **Pourquoi** : un modèle moins complexe a moins de "liberté" pour mémoriser le bruit et doit se concentrer sur les tendances principales.

3. Diagnostic inverse : train = 60%, test = 58%

C'est un cas d'**underfitting** (sous-apprentissage). Le modèle est trop simple pour capturer les relations dans les données. Les deux scores sont **faibles et proches** l'un de l'autre.

Solutions :

- Utiliser un modèle plus complexe (ex : Random Forest au lieu de régression logistique)
- Ajouter de nouvelles features (ingénierie des features)
- Réduire la régularisation (si elle est trop forte)

Exercice 19.3 — Python : GridSearchCV pour Random Forest

Écrivez un code Python complet qui :

1. Charge (ou génère) le dataset clientèle de classification (churn).
2. Crée un Pipeline avec StandardScaler et RandomForestClassifier.
3. Définit une grille d'hyperparamètres :
 - `n_estimators` : [50, 100, 200]
 - `max_depth` : [3, 5, 10, None]
 - `min_samples_split` : [2, 5, 10]
 - `min_samples_leaf` : [1, 2, 4]
4. Exécute un GridSearchCV avec 5-fold CV stratifiée.
5. Affiche les meilleurs hyperparamètres et le meilleur score CV.
6. Évalue le modèle optimal sur un test set séparé.
7. Trace la courbe d'apprentissage du modèle optimal.

Correction de l'exercice 19.3

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import (train_test_split,
                                     GridSearchCV, StratifiedKFold, learning_curve)
from sklearn.preprocessing import StandardScaler
from sklearn.pipeline import Pipeline
from sklearn.ensemble import RandomForestClassifier

# 1. Générer le dataset clientèle
np.random.seed(42)
n = 300
revenu = np.random.uniform(20000, 80000, n)
anciennete = np.random.uniform(1, 15, n)
nb_achats = np.random.randint(1, 50, n)
satisfaction = np.random.uniform(1, 10, n)
X = np.column_stack([revenu, anciennete, nb_achats,
                    satisfaction])

proba = 1 / (1 + np.exp(0.00005*revenu + 0.2*anciennete
                      + 0.05*nb_achats + 0.3*satisfaction - 5))
y = (np.random.rand(n) < proba).astype(int)
```

```

# Separer train / test
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42, stratify=y)

# 2. Pipeline
pipe = Pipeline([
    ('scaler', StandardScaler()),
    ('model', RandomForestClassifier(random_state=42))
])

# 3. Grille d'hyperparametres
param_grid = {
    'model__n_estimators': [50, 100, 200],
    'model__max_depth': [3, 5, 10, None],
    'model__min_samples_split': [2, 5, 10],
    'model__min_samples_leaf': [1, 2, 4],
}

# Nombre total de combinaisons
nb_combis = 3 * 4 * 3 * 3
print(f"Nombre de combinaisons : {nb_combis}")
print(f"Nombre total d'entraînements (x5 folds) : "
      f"{nb_combis * 5}")

# 4. GridSearchCV
cv = StratifiedKFold(n_splits=5, shuffle=True,
                     random_state=42)

grid = GridSearchCV(pipe, param_grid, cv=cv,
                    scoring='accuracy', n_jobs=-1,
                    return_train_score=True, verbose=1)
grid.fit(X_train, y_train)

# 5. Meilleurs hyperparametres
print(f"\n=== Resultats ===")
print(f"Meilleurs hyperparametres :")
for param, val in grid.best_params_.items():
    print(f"    {param} = {val}")
print(f"Meilleur score CV : {grid.best_score_:.4f}")

```

```

# Top 5 des combinaisons
import pandas as pd
resultats = pd.DataFrame(grid.cv_results_)
top5 = resultats.nsmallest(5, 'rank_test_score')[
    ['params', 'mean_test_score', 'std_test_score',
     'rank_test_score']]
print(f"\nTop 5 des combinaisons :")
print(top5.to_string(index=False))

# 6. Evaluation finale sur le test set
score_test = grid.score(X_test, y_test)
print(f"\nScore final sur le test set : {score_test:.4f}")
print(f"Score CV (rappel)           : "
      f"{grid.best_score_:.4f}")

# 7. Courbe d'apprentissage du modele optimal
train_sizes, train_scores, val_scores = learning_curve(
    grid.best_estimator_, X_train, y_train,
    cv=cv, n_jobs=-1,
    train_sizes=np.linspace(0.1, 1.0, 10),
    scoring='accuracy')

train_mean = train_scores.mean(axis=1)
train_std = train_scores.std(axis=1)
val_mean = val_scores.mean(axis=1)
val_std = val_scores.std(axis=1)

plt.figure(figsize=(10, 6))
plt.fill_between(train_sizes, train_mean - train_std,
                 train_mean + train_std,
                 alpha=0.1, color='blue')
plt.fill_between(train_sizes, val_mean - val_std,
                 val_mean + val_std,
                 alpha=0.1, color='red')
plt.plot(train_sizes, train_mean, 'o-',
         color='blue', label='Score train')
plt.plot(train_sizes, val_mean, 'o-',
         color='red', label='Score validation')

plt.xlabel('Taille de l\'ensemble d\'entraînement',

```

```
        fontsize=12)
plt.ylabel('Accuracy', fontsize=12)
plt.title('Courbe d\'apprentissage - Random Forest\n'
        f'(meilleurs hyperparametres)',
        fontsize=14, fontweight='bold')
plt.legend(fontsize=11)
plt.ylim(0.5, 1.05)
plt.grid(True, alpha=0.3)
plt.tight_layout()
plt.show()

# Diagnostic
ecart = train_mean[-1] - val_mean[-1]
print(f"\nDiagnostic :")
print(f"  Score train final   : {train_mean[-1]:.4f}")
print(f"  Score valid. final   : {val_mean[-1]:.4f}")
print(f"  Ecart                  : {ecart:.4f}")
if ecart > 0.15:
    print("    -> Overfitting : essayer plus de regularisation")
elif val_mean[-1] < 0.7:
    print("    -> Underfitting : essayer un modele plus complexe")
else:
    print("    -> Bon equilibre biais-variance !")
```

Chapitre 20

Guide de Choix d'Algorithme

20.1 Introduction : face à la jungle des algorithmes

Au fil de ce cours, vous avez découvert plus d'une dizaine d'algorithmes de Machine Learning supervisé : régression linéaire, logistique, KNN, arbres de décision, Random Forest, SVM, Naïve Bayes, Gradient Boosting, XGBoost, LightGBM, CatBoost, réseaux de neurones...

Mais lequel choisir ?

C'est la question la plus fréquente chez les débutants — et même chez les praticiens expérimentés. La vérité, c'est qu'il n'existe pas d'algorithme « magique » qui fonctionne mieux que tous les autres dans toutes les situations. Le célèbre théorème du **No Free Lunch** nous dit exactement cela : aucun algorithme n'est universellement supérieur.

Théorème du No Free Lunch

Il n'existe aucun algorithme de Machine Learning qui soit le meilleur pour **tous** les problèmes. Chaque algorithme a des forces et des faiblesses. Le choix dépend de **vos données**, de **votre problème** et de **vos contraintes**.

Ce chapitre vous fournit un **guide pratique complet** pour choisir le bon algorithme :

1. Un **tableau comparatif** récapitulatif de tous les algorithmes.
2. Un **arbre de décision** (flowchart) pour guider votre choix.
3. Une **checklist** de questions à se poser.
4. Les **forces et faiblesses** de chaque algorithme.
5. Une **stratégie pratique** systématique.
6. Une **implémentation complète** sur Google Colab.

20.2 Tableau comparatif des algorithmes

Le tableau suivant résume les caractéristiques essentielles de chaque algorithme étudié dans ce cours. C'est votre **fiche de référence rapide**.

Comment lire ce tableau

- **Type** : l'algorithme est-il utilisé pour la régression, la classification, ou les deux ?
- **Linéaire / Non-linéaire** : l'algorithme suppose-t-il une relation linéaire entre les variables ?
- **Interprétable** : peut-on facilement comprendre « pourquoi » le modèle fait telle prédiction ?
- **Scaling requis** : faut-il normaliser/standardiser les données avant l'entraînement ?
- **Risque d'overfitting** : le modèle a-t-il tendance à sur-apprendre ?
- **Cas d'usage typique** : dans quel contexte cet algorithme brille-t-il ?

Algorithme	Type	Lin./Non-lin.	Interpr.	Scaling	Overfit.	Cas d'usage
Régression Linéaire	Rég.	Linéaire	Oui	Oui	Moyen	Prédiction continue simple
Régression Logistique	Classif.	Linéaire	Oui	Oui	Moyen	Classification binaire
KNN	Les deux	Non-lin.	Oui	Oui++	Élevé	Petit dataset
Arbre de Décision	Les deux	Non-lin.	Oui	Non	Élevé	Règles métiers
Random Forest	Les deux	Non-lin.	Non	Non	Faible	Modèle robuste général
SVM	Les deux	Les deux	Partiel	Oui++	Faible à moyen	Données complexes
Naïve Bayes	Classif.	Probabiliste	Non	Non	Faible	NLP, spam
Gradient Boosting	Les deux	Non-lin.	Non	Non	Faible	Haute performance
XGBoost	Les deux	Non-lin.	Non	Non	Faible	Compétitions
LightGBM	Les deux	Non-lin.	Non	Non	Faible	Big Data
CatBoost	Les deux	Non-lin.	Non	Non	Faible	Variables catégorielles
Réseaux Neuraux (MLP)	Les deux	Non-lin.	Non	Oui++	Moyen à élevé	Problèmes complexes

TABLE 20.1 – Tableau comparatif de tous les algorithmes de ML supervisé étudiés dans ce cours.

Attention — « Oui++ » signifie indispensable

Quand le scaling est marqué « Oui++ », cela signifie que l'algorithme est **très sensible** à l'échelle des variables. Ne pas normaliser peut rendre le modèle complètement inutile (KNN, SVM, MLP).

20.3 Arbre de décision pour le choix d'algorithme

L'arbre de décision ci-dessous vous guide pas à pas dans le choix du bon algorithme en répondant à des questions simples par **Oui** ou **Non**.

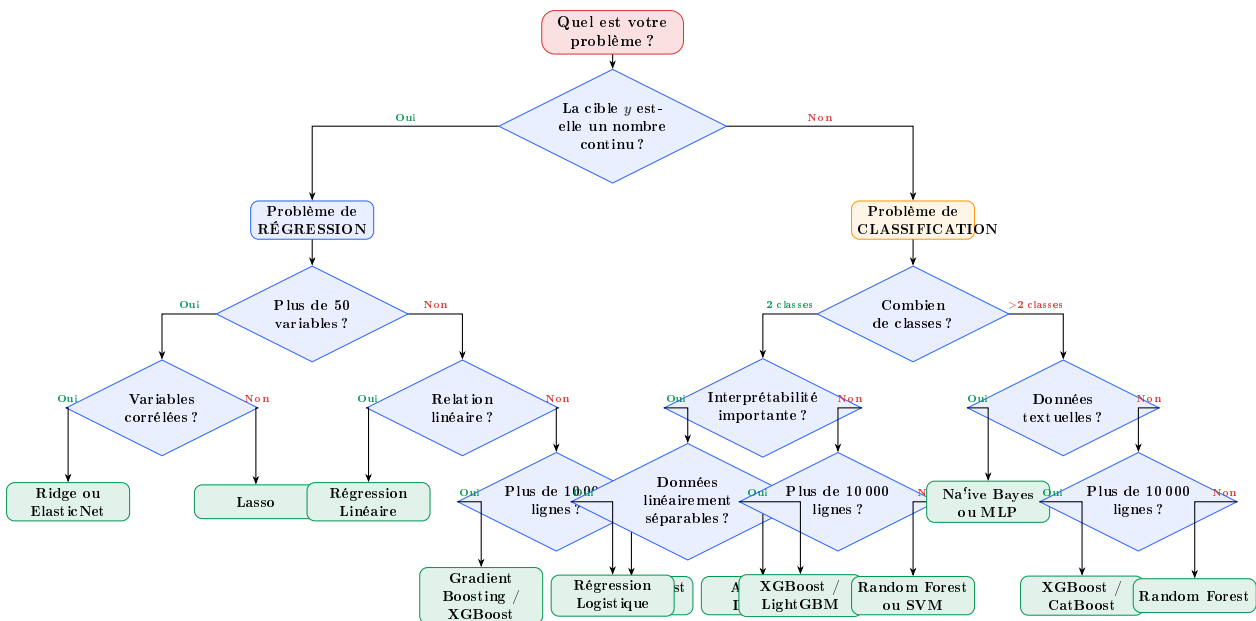


FIGURE 20.1 – Arbre de décision pour le choix d'un algorithme de ML supervisé.

Comment utiliser ce flowchart ?

Partez du haut et répondez à chaque question par **Oui** ou **Non**. Suivez les flèches jusqu'à atteindre un algorithme recommandé (encadré vert). Ce n'est qu'un **point de départ** : vous devriez toujours tester plusieurs algorithmes et comparer avec la validation croisée (voir Section 20.6).

20.4 Questions à se poser avant de choisir

Avant de choisir un algorithme, posez-vous ces questions clés. Le tableau ci-dessous vous guide vers les bons choix selon vos réponses.

Question à se poser	Si la réponse est...	Algorithmes recommandés
Régression ou classification ?	Régression	Rég. linéaire, Ridge, RF, XGBoost
	Classification	LogReg, SVM, RF, XGBoost
Combien de données avez-vous ?	Peu (< 1 000)	KNN, Arbre, SVM
	Moyen (1k–100k)	RF, XGBoost, SVM
	Beaucoup (> 100k)	LightGBM, XGBoost, MLP
Combien de variables (features) ?	Peu (< 20)	Tous conviennent
	Beaucoup (> 50)	Lasso, Ridge, RF, XGBoost
L'interprétabilité est-elle cruciale ?	Oui, absolument	Rég. lin./log., Arbre
	Non, performance avant tout	XGBoost, LightGBM, MLP
Avez-vous des variables catégorielles ?	Beaucoup	CatBoost, Arbre, RF
	Peu ou aucune	Tous conviennent
Vitesse d'entraînement importante ?	Oui, très rapide	Naïve Bayes, Rég. linéaire, LightGBM
	Non, on a le temps	XGBoost, SVM, MLP
Les classes sont-elles déséquilibrées ?	Oui	XGBoost (scale_pos_weight), RF, SMOTE + LogReg
	Non	Tous conviennent
Les données sont-elles textuelles ?	Oui (NLP)	Naïve Bayes, MLP, SVM
	Non (tabulaires)	XGBoost, RF, LightGBM

TABLE 20.2 – Checklist de questions pour guider le choix d'algorithme.

Exemple de raisonnement

« J'ai un problème de **classification binaire** avec **50 000 lignes**, **30 variables** dont 10 catégorielles, et je veux **maximiser la performance**. »

Raisonnement :

- Classification \Rightarrow éliminer Rég. linéaire
- 50 000 lignes \Rightarrow dataset de taille moyenne, tous les algorithmes conviennent
- 10 variables catégorielles \Rightarrow CatBoost ou RF sont de bons candidats
- Performance maximale \Rightarrow méthodes de boosting
- **Recommandation** : commencer par **CatBoost**, comparer avec **XGBoost** et **Random Forest**

20.5 Forces et faiblesses de chaque algorithme

Pour chaque algorithme, voici un résumé de ses 3 principales forces, 3 faiblesses, et son meilleur cas d'usage.

Régression Linéaire

✓ Forces :

- ✓ Très simple à comprendre et à expliquer
- ✓ Rapide à entraîner, même sur de grands datasets
- ✓ Coefficients directement interprétables

✗ Faiblesses :

- ✗ Ne capture que les relations linéaires
- ✗ Sensible aux outliers
- ✗ Hypothèse forte de linéarité souvent violée

Meilleur usage : Prédiction continue quand la relation est approximativement linéaire, ou comme **baseline** de référence.

Régression Logistique

✓ Forces :

- ✓ Fournit des probabilités calibrées
- ✓ Rapide et stable
- ✓ Excellente interprétabilité (odds ratios)

✗ Faiblesses :

- ✗ Suppose une frontière de décision linéaire
- ✗ Peu performante sur des données non-linéaires complexes
- ✗ Sensible à la multicollinéarité

Meilleur usage : Classification binaire quand l'interprétabilité est essentielle (médecine, finance, droit).

KNN (K-Nearest Neighbors)

✓ **Forces :**

- ✓ Aucune hypothèse sur la distribution des données
- ✓ Très intuitif (les voisins les plus proches)
- ✓ Fonctionne bien sur les petits datasets

✗ **Faiblesses :**

- ✗ Très lent en prédiction sur les grands datasets
- ✗ Très sensible à l'échelle des variables (scaling obligatoire)
- ✗ Souffre de la « malédiction de la dimension » (beaucoup de variables)

Meilleur usage : Petits datasets (< 5 000 lignes) avec peu de variables.

Arbre de Décision

✓ **Forces :**

- ✓ Extrêmement interprétable (visualisable)
- ✓ Pas de scaling requis
- ✓ Gère nativement les variables catégorielles et les valeurs manquantes

✗ **Faiblesses :**

- ✗ Très sujet à l'overfitting sans élagage
- ✗ Instable : un petit changement dans les données peut changer tout l'arbre
- ✗ Performance souvent inférieure aux méthodes d'ensemble

Meilleur usage : Quand on a besoin de règles métiers explicites et compréhensibles.

Random Forest

✓ **Forces :**

- ✓ Très robuste à l'overfitting grâce au bagging
- ✓ Peu de réglages nécessaires (peu d'hyperparamètres critiques)
- ✓ Fournit l'importance des variables

✗ **Faiblesses :**

- ✗ Modèle boîte noire (difficile à interpréter)
- ✗ Lent à entraîner avec beaucoup d'arbres

- ✗ Consomme beaucoup de mémoire

Meilleur usage : Modèle robuste par défaut quand on ne sait pas quel algorithme choisir.

SVM (Support Vector Machine)

✓ **Forces :**

- ✓ Très efficace en haute dimension
- ✓ Frontières de décision flexibles grâce aux noyaux (kernel trick)
- ✓ Robuste à l'overfitting avec une bonne régularisation

✗ **Faiblesses :**

- ✗ Très lent sur les grands datasets ($> 50\,000$ lignes)
- ✗ Scaling obligatoire
- ✗ Choix du noyau et des hyperparamètres délicat

Meilleur usage : Datasets de taille moyenne avec des frontières de décision complexes.

Naïve Bayes

✓ **Forces :**

- ✓ Extrêmement rapide (entraînement et prédiction)
- ✓ Fonctionne très bien avec peu de données
- ✓ Excellent pour le texte (NLP, spam, sentiment)

✗ **Faiblesses :**

- ✗ Hypothèse d'indépendance des variables rarement vraie
- ✗ Probabilités mal calibrées
- ✗ Performance limitée sur les données tabulaires numériques

Meilleur usage : Classification de textes, filtrage de spam, avec de nombreuses variables.

Gradient Boosting (et XGBoost ; LightGBM ; CatBoost)

✓ **Forces :**

- ✓ Souvent le plus performant sur les données tabulaires
- ✓ Gère les valeurs manquantes et les variables catégorielles (CatBoost)
- ✓ Régularisation intégrée contre l'overfitting

✗ **Faiblesses :**

- ✗ Beaucoup d'hyperparamètres à régler

- × Risque d'overfitting si mal configuré
- × Modèle boîte noire

Meilleur usage : Maximiser la performance sur des données tabulaires. **Choix #1 en compétition Kaggle.**

Réseaux de Neurons (MLP)

✓ **Forces :**

- ✓ Peut apprendre des relations extrêmement complexes
- ✓ Flexible : s'adapte à presque tout type de données
- ✓ Excellentes performances avec beaucoup de données

× **Faiblesses :**

- × Nécessite beaucoup de données et de calcul
- × Complètement boîte noire
- × Beaucoup d'hyperparamètres (couches, neurones, learning rate...)

Meilleur usage : Grandes quantités de données avec des relations très complexes. Moins adapté aux données tabulaires classiques qu'aux images/textes/sons.

20.6 Stratégie pratique systématique

Règle d'or

En pratique, ne choisissez jamais un seul algorithme à l'avance. Essayez-en plusieurs et comparez-les objectivement avec la **validation croisée**.

Voici la stratégie en 4 étapes que tout data scientist devrait suivre :

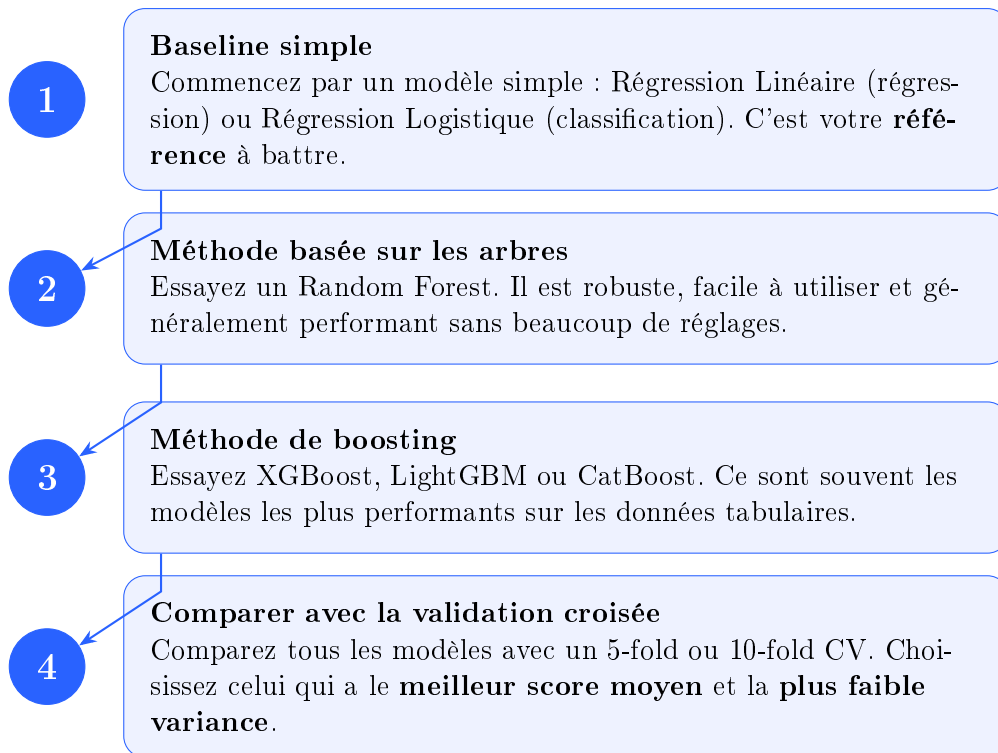


FIGURE 20.2 – Stratégie en 4 étapes pour choisir le meilleur algorithme.

Analogie — Le casting d'un film

Imaginez que vous êtes réalisateur et que vous cherchez un acteur pour un rôle. Vous ne choisissez pas le premier acteur qui se présente : vous en auditionnez plusieurs et vous comparez leurs performances. C'est exactement pareil en ML : vous « auditionnez » plusieurs algorithmes sur vos données et vous gardez le meilleur.

20.7 Implémentation sur Google Colab

Hands-On 20.1 — Comparaison complète de 8+ algorithmes

Nous allons entraîner **8 algorithmes différents** sur le même dataset, les évaluer par **validation croisée**, puis comparer leurs performances visuellement.

20.7.1 Étape 1 : Préparation des données

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import warnings
```

```

5 warnings.filterwarnings('ignore')
6
7 # Algorithmes de classification
8 from sklearn.linear_model import LogisticRegression
9 from sklearn.neighbors import KNeighborsClassifier
10 from sklearn.tree import DecisionTreeClassifier
11 from sklearn.ensemble import (RandomForestClassifier,
12                               GradientBoostingClassifier)
13 from sklearn.svm import SVC
14 from sklearn.naive_bayes import GaussianNB
15 from sklearn.neural_network import MLPClassifier
16
17 # Evaluation
18 from sklearn.model_selection import cross_val_score,
19     StratifiedKFold
20 from sklearn.preprocessing import StandardScaler
21 from sklearn.pipeline import Pipeline
22 from sklearn.datasets import make_classification
23
24 # Creer un dataset realiste (simulation de donnees clientele)
25 X, y = make_classification(
26     n_samples=2000,
27     n_features=15,
28     n_informative=10,
29     n_redundant=3,
30     n_classes=2,
31     weights=[0.7, 0.3], # Legerement desequilibre
32     random_state=42
33 )
34 print(f''Taille du dataset : {X.shape[0]} lignes, {X.shape[1]}
35     variables'')
36 print(f''Repartition des classes : {np.bincount(y)}'')
```

Listing 20.1 – Import des bibliothèques et préparation des données

20.7.2 Étape 2 : Définition des modèles et validation croisée

```

1 # Definir les algorithmes a comparer
2 # On utilise des Pipelines pour inclure le scaling quand necessaire
3 modeles = {
```

```

4      'Regression Logistique': Pipeline([
5          ('scaler', StandardScaler()),
6          ('model', LogisticRegression(max_iter=1000,
7              random_state=42))
8      ]),
9      'KNN (k=5)': Pipeline([
10         ('scaler', StandardScaler()),
11         ('model', KNeighborsClassifier(n_neighbors=5))
12     ]),
13     'Arbre de Decision': DecisionTreeClassifier(
14         max_depth=10, random_state=42
15     ),
16     'Random Forest': RandomForestClassifier(
17         n_estimators=100, random_state=42
18     ),
19     'SVM (RBF)': Pipeline([
20         ('scaler', StandardScaler()),
21         ('model', SVC(kernel='rbf', random_state=42))
22     ]),
23     'Na?ve Bayes': GaussianNB(),
24     'Gradient Boosting': GradientBoostingClassifier(
25         n_estimators=100, random_state=42
26     ),
27     'MLP (Reseau Neurones)': Pipeline([
28         ('scaler', StandardScaler()),
29         ('model', MLPClassifier(hidden_layer_sizes=(64, 32),
30             max_iter=500, random_state=42))
31     ]),
32 }
33 # Validation croisee stratifiee (5-fold)
34 cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
35
36 resultats = {}
37 print(''= ' * 65)
38 print(f'{'Algorithme':<30} {'Accuracy Moyenne':>15}
39       {'?cart-type':>12}')
40 print(''= ' * 65)
41 for nom, modele in modeles.items():

```

```

42     scores = cross_val_score(modele, X, y, cv=cv,
                               scoring='accuracy')
43     resultats[nom] = scores
44     print(f''{nom:<30} {scores.mean():>15.4f}
           {scores.std():>12.4f}'')
45
46 print(''= ' * 65)
47
48 # Identifier le meilleur modele
49 meilleur = max(resultats, key=lambda k: resultats[k].mean())
50 print(f''\n>>> Meilleur modele : {meilleur} ''
       f''(accuracy = {resultats[meilleur].mean():.4f})'')
51

```

Listing 20.2 – Comparaison de 8 algorithmes avec validation croisée

20.7.3 Étape 3 : Visualisation comparative

```

1  # --- Graphique 1 : Barres avec barres d'erreur ---
2  fig, axes = plt.subplots(1, 2, figsize=(16, 6))
3
4  noms = list(resultats.keys())
5  moyennes = [resultats[n].mean() for n in noms]
6  ecarts = [resultats[n].std() for n in noms]
7
8  # Couleurs : le meilleur en vert, les autres en bleu
9  couleurs = ['#2ecc71' if n == meilleur else '#3498db' for n in
              noms]
10
11 bars = axes[0].barh(noms, moyennes, xerr=ecarts, color=couleurs,
12                    edgecolor='white', height=0.6, capsize=4)
13 axes[0].set_xlabel('Accuracy (CV 5-fold)', fontsize=12)
14 axes[0].set_title('Comparaison des algorithmes', fontsize=14,
15                  fontweight='bold')
16 axes[0].set_xlim(min(moyennes) - 0.05, max(moyennes) + 0.02)
17
18 # Ajouter les valeurs sur les barres
19 for bar, moy in zip(bars, moyennes):
20     axes[0].text(moy + 0.005, bar.get_y() + bar.get_height()/2,
21                 f'{moy:.4f}', va='center', fontweight='bold',
22                 fontsize=9)
23

```

```

22 # --- Graphique 2 : Boxplot ---
23 data_bp = [resultats[n] for n in noms]
24 bp = axes[1].boxplot(data_bp, labels=noms, vert=False,
    patch_artist=True)
25 for patch, couleur in zip(bp['boxes'], couleurs):
26     patch.set_facecolor(couleur)
27     patch.set_alpha(0.7)
28 axes[1].set_xlabel('Accuracy', fontsize=12)
29 axes[1].set_title('Distribution des scores (5 folds)', fontsize=14,
30                 fontweight='bold')
31
32 plt.tight_layout()
33 plt.savefig('comparaison_algorithmes.png', dpi=150,
34             bbox_inches='tight')
35 plt.show()

```

Listing 20.3 – Diagramme en barres et boxplot comparatif

20.7.4 Étape 4 : Tableau récapitulatif avec plusieurs métriques

```

1 from sklearn.model_selection import cross_validate
2
3 # Metriques multiples
4 metriques = ['accuracy', 'precision', 'recall', 'f1']
5
6 print('\n' + '=' * 80)
7 print(f'{'Algorithme':<25} {'Accuracy':>10} {'Precision':>10} ' '
8       f'{'Recall':>10} {'F1-Score':>10}')
9 print('=' * 80)
10
11 tableau_complet = []
12
13 for nom, modele in modeles.items():
14     scores = cross_validate(modele, X, y, cv=cv,
15                             scoring=metriques,
16                             return_train_score=False)
17
18     ligne = {
19         'Algorithme': nom,
20         'Accuracy': scores['test_accuracy'].mean(),
21         'Precision': scores['test_precision'].mean(),
22         'Recall': scores['test_recall'].mean(),

```

```

21         'F1-Score': scores['test_f1'].mean(),
22     }
23     tableau_complet.append(ligne)
24     print(f''{nom:<25} {ligne['Accuracy']:>10.4f}
25           {ligne['Precision']:>10.4f} ''
26           f''{ligne['Recall']:>10.4f} {ligne['F1-Score']:>10.4f}''')
27
28
29 # Convertir en DataFrame pour un affichage propre
30 df_resultats = pd.DataFrame(tableau_complet)
31 df_resultats = df_resultats.sort_values('F1-Score',
32                                         ascending=False)
33 print(''\n>>> Classement final par F1-Score :''')
34 print(df_resultats.to_string(index=False))

```

Listing 20.4 – Tableau complet avec accuracy, precision, recall et F1

20.7.5 Étape 5 : LazyPredict pour une comparaison rapide (bonus)

```

1  # Installation (décommenter si nécessaire)
2  # !pip install lazypredict
3
4  try:
5      from lazypredict.Supervised import LazyClassifier
6      from sklearn.model_selection import train_test_split
7
8      X_train, X_test, y_train, y_test = train_test_split(
9          X, y, test_size=0.2, random_state=42, stratify=y
10     )
11
12     clf = LazyClassifier(verbose=0, ignore_warnings=True,
13                         custom_metric=None)
14     models, predictions = clf.fit(X_train, X_test, y_train, y_test)
15
16     # Afficher le top 10
17     print(''\n>>> Top 10 des algorithmes (LazyPredict) :''')
18     print(models.head(10))
19
20 except ImportError:
21     print(''\nLazyPredict non installé.'')

```

```
22     print('Installez-le avec : pip install lazypredict')
23     print('Il compare automatiquement 30+ algorithmes en une
        ligne !')
```

Listing 20.5 – Utilisation de LazyPredict pour comparer automatiquement des dizaines de modèles

LazyPredict — Outil de prototypage rapide

LazyPredict est un outil qui entraîne automatiquement des dizaines d'algorithmes et les compare. C'est très utile pour le **prototypage rapide**, mais attention :

- Il n'utilise pas la validation croisée (simple split train/test).
- Il n'optimise pas les hyperparamètres.
- Utilisez-le comme **point de départ**, puis affinez les meilleurs candidats.

20.8 Exercices

Exercice 20.1 — Recommander le bon algorithme

Pour chacune des 5 situations suivantes, recommandez le ou les algorithmes les plus adaptés. **Justifiez** votre choix en vous basant sur les caractéristiques du problème.

1. **Situation A** : Une banque veut prédire si un client va faire défaut de paiement (oui/non). Elle a 100 000 clients et 25 variables. L'interprétabilité est **exigée par la réglementation**.
2. **Situation B** : Un site e-commerce veut prédire le **montant** des dépenses d'un client sur les 12 prochains mois. Il a 5 000 clients et 8 variables. La relation semble approximativement linéaire.
3. **Situation C** : Un hôpital veut classifier des emails de patients en 5 catégories (urgence, rendez-vous, question, réclamation, autre). Il a 50 000 emails.
4. **Situation D** : Une startup veut prédire le prix d'un appartement à partir de 200 variables (dont beaucoup sont corrélées). Elle a 500 000 lignes.
5. **Situation E** : Un chercheur a 300 échantillons biologiques avec 10 variables et veut classifier chaque échantillon en 2 catégories. Il veut comprendre quelles variables sont déterminantes.

Correction de l'exercice 20.1**Situation A — Défaut de paiement bancaire :**

- **Recommandation : Régression Logistique** (premier choix) ou **Arbre de Décision** (élagué).
- **Justification** : L'interprétabilité est exigée par la réglementation (Comité de Bâle, RGPD). La Régression Logistique fournit des odds ratios directement interprétables. Avec 100 000 lignes et 25 variables, elle est parfaitement adaptée. Un arbre de décision élagué est une alternative car il produit des règles métiers explicites.
- **À éviter** : XGBoost, MLP (boîtes noires incompatibles avec les exigences réglementaires).

Situation B — Prédiction du montant des dépenses :

- **Recommandation : Régression Linéaire** (baseline) puis **Random Forest Regressor**.
- **Justification** : C'est un problème de régression. Avec 5 000 lignes et 8 variables, le dataset est petit et simple. La relation étant approximativement linéaire, la régression linéaire est un excellent point de départ. Random Forest en complément pour capturer d'éventuelles non-linéarités.
- **À éviter** : MLP (trop peu de données), KNN (performant mais fragile sur les outliers).

Situation C — Classification d'emails en 5 catégories :

- **Recommandation : Naïve Bayes** (baseline rapide) puis **SVM linéaire** ou **MLP**.
- **Justification** : Les données sont textuelles (NLP). Naïve Bayes est le standard pour la classification de texte : très rapide et efficace. Avec 50 000 emails, on a assez de données pour un SVM ou un MLP qui seront probablement plus performants. C'est un problème multiclasse (5 classes).
- **À éviter** : KNN (trop lent sur les vecteurs TF-IDF de grande dimension).

Situation D — Prix d'appartement avec 200 variables corrélées :

- **Recommandation : Ridge** ou **ElasticNet** (pour gérer la corrélation) puis **XGBoost** ou **LightGBM** (pour la performance).
- **Justification** : 200 variables corrélées \Rightarrow la régularisation est indispensable. Ridge pénalise les coefficients des variables corrélées sans les éliminer, ElasticNet peut en éliminer certaines. Avec 500 000 lignes, LightGBM est idéal (très rapide sur le Big Data).
- **À éviter** : Régression linéaire sans régularisation (multicolinéarité), SVM (trop lent à 500k lignes).

Situation E — 300 échantillons biologiques :

- **Recommandation : Régression Logistique** avec analyse des coefficients, ou **SVM linéaire**.
- **Justification** : Très peu de données (300) \Rightarrow éviter les modèles complexes qui vont overfitter. 10 variables et besoin d'interprétabilité \Rightarrow la Régression Logistique permet d'identifier les variables déterminantes via les coefficients. SVM linéaire est également robuste avec peu de données.
- **À éviter** : XGBoost, MLP, Random Forest (risque d'overfitting avec si peu de données).

Exercice 20.2 — Pipeline complet de comparaison (Python)

Écrivez un script Python complet qui :

1. Charge le dataset `breast_cancer` de scikit-learn (classification binaire).
2. Compare **au moins 8 algorithmes** avec une validation croisée 10-fold.
3. Affiche un tableau récapitulatif trié par F1-Score décroissant.
4. Génère un graphique en barres horizontales comparant les accuracy.
5. Identifie et affiche le meilleur modèle.

Correction de l'exercice 20.2

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.datasets import load_breast_cancer
5 from sklearn.model_selection import cross_validate,
    StratifiedKFold
6 from sklearn.preprocessing import StandardScaler
7 from sklearn.pipeline import Pipeline
8
9 # Algorithmes
10 from sklearn.linear_model import LogisticRegression
11 from sklearn.neighbors import KNeighborsClassifier
12 from sklearn.tree import DecisionTreeClassifier
13 from sklearn.ensemble import (RandomForestClassifier,
14                               GradientBoostingClassifier,
15                               AdaBoostClassifier)
```

```

16 from sklearn.svm import SVC
17 from sklearn.naive_bayes import GaussianNB
18 from sklearn.neural_network import MLPClassifier
19
20 # 1. Charger les donnees
21 data = load_breast_cancer()
22 X, y = data.data, data.target
23 print(f'Dataset : {X.shape[0]} lignes, {X.shape[1]}
      variables')
24 print(f'Classes : {np.bincount(y)} (0=malin, 1=benin)\n')
25
26 # 2. Definir les modeles
27 modeles = {
28     'Reg. Logistique': Pipeline([
29         ('scaler', StandardScaler()),
30         ('model', LogisticRegression(max_iter=5000,
      random_state=42))
31     ]),
32     'KNN (k=5)': Pipeline([
33         ('scaler', StandardScaler()),
34         ('model', KNeighborsClassifier(n_neighbors=5))
35     ]),
36     'Arbre Decision': DecisionTreeClassifier(
37         max_depth=8, random_state=42),
38     'Random Forest': RandomForestClassifier(
39         n_estimators=200, random_state=42),
40     'SVM (RBF)': Pipeline([
41         ('scaler', StandardScaler()),
42         ('model', SVC(kernel='rbf', random_state=42))
43     ]),
44     'Na?ve Bayes': GaussianNB(),
45     'Gradient Boosting': GradientBoostingClassifier(
46         n_estimators=100, random_state=42),
47     'AdaBoost': AdaBoostClassifier(
48         n_estimators=100, random_state=42),
49     'MLP': Pipeline([
50         ('scaler', StandardScaler()),
51         ('model', MLPClassifier(hidden_layer_sizes=(64, 32),
      max_iter=1000, random_state=42))
52     ]),
53

```

```

54 }
55
56 # 3. Validation croisee 10-fold
57 cv = StratifiedKFold(n_splits=10, shuffle=True, random_state=42)
58 metriques = ['accuracy', 'precision', 'recall', 'f1']
59
60 resultats = []
61 for nom, modele in modeles.items():
62     scores = cross_validate(modele, X, y, cv=cv,
63                             scoring=metriques,
64                             return_train_score=False)
65     resultats.append({
66         'Algorithme': nom,
67         'Accuracy': scores['test_accuracy'].mean(),
68         'Acc Std': scores['test_accuracy'].std(),
69         'Precision': scores['test_precision'].mean(),
70         'Recall': scores['test_recall'].mean(),
71         'F1-Score': scores['test_f1'].mean(),
72     })
73
74 # 4. Tableau recapitulatif trie par F1-Score
75 df = pd.DataFrame(resultats).sort_values('F1-Score',
76                                         ascending=False)
77 df = df.reset_index(drop=True)
78 df.index += 1 # Rang a partir de 1
79
80 print('=' * 75)
81 print('CLASSEMENT PAR F1-SCORE (validation croisee 10-fold)')
82 print('=' * 75)
83 print(df.to_string())
84 print('=' * 75)
85
86 # 5. Graphique en barres horizontales
87 fig, ax = plt.subplots(figsize=(10, 6))
88 couleurs = ['#2ecc71' if i == 0 else '#3498db'
89             for i in range(len(df))]
90 bars = ax.barh(df['Algorithme'], df['Accuracy'],
91               xerr=df['Acc Std'], color=couleurs,
92               edgecolor='white', height=0.6, capsize=4)
93

```

```

92 # Valeurs sur les barres
93 for bar, acc in zip(bars, df['Accuracy']):
94     ax.text(acc + 0.005, bar.get_y() + bar.get_height()/2,
95            f'{acc:.4f}', va='center', fontweight='bold')
96
97 ax.set_xlabel('Accuracy (CV 10-fold)', fontsize=12)
98 ax.set_title('Comparaison de 9 algorithmes - Breast Cancer',
99            fontsize=14, fontweight='bold')
100 ax.invert_yaxis()
101 plt.tight_layout()
102 plt.savefig('exercice_20_2_comparaison.png', dpi=150,
103            bbox_inches='tight')
104 plt.show()
105
106 # 6. Meilleur modele
107 meilleur = df.iloc[0]
108 print(f'''\n>>> MEILLEUR MOD?LE : {meilleur['Algorithme']}'')
109 print(f'''\t\t\tAccuracy = {meilleur['Accuracy']:.4f} ''
110       f'''\t\t\t(+/- {meilleur['Acc Std']:.4f})''')
111 print(f'''\t\t\tF1-Score = {meilleur['F1-Score']:.4f}'')
112 print(f'''\t\t\tPrecision = {meilleur['Precision']:.4f}'')
113 print(f'''\t\t\tRecall = {meilleur['Recall']:.4f}'')

```

Listing 20.6 – Pipeline complet de comparaison sur breast_cancer

Exercice 20.3 — Problème réel complexe

Un client vous contacte avec le problème suivant :

« Nous avons **1 million de lignes** et **200 features**. Notre variable cible est binaire (achat / pas achat) avec seulement **3% de cas positifs** (classe très déséquilibrée). Nous voulons maximiser la détection des acheteurs potentiels. Quel algorithme recommandez-vous et pourquoi ? »

1. Analysez les caractéristiques du problème (type, taille, déséquilibre, etc.).
2. Recommandez un algorithme principal et justifiez.
3. Proposez des techniques complémentaires pour gérer le déséquilibre.
4. Écrivez le code Python correspondant.

Correction de l'exercice 20.3

1. Analyse du problème :

- **Type** : Classification binaire
- **Taille** : 1 million de lignes \Rightarrow **Big Data**, élimine les algorithmes lents (SVM à noyau, KNN)
- **200 features** : haute dimension, beaucoup de variables
- **3% positifs** : classe **très déséquilibrée** \Rightarrow l'accuracy est une métrique trompeuse (97% d'accuracy en prédisant toujours 0 !)
- **Objectif** : maximiser la détection \Rightarrow focus sur le **Recall** et le **F1-Score**

2. Recommandation : LightGBM (premier choix) ou XGBoost

Justification :

- **LightGBM** est conçu pour le Big Data : il est 5–10x plus rapide que XGBoost sur de grands datasets grâce à son algorithme de découpage par histogrammes.
- Il gère nativement le déséquilibre via le paramètre `scale_pos_weight` ou `is_unbalance`.
- Pas besoin de scaling.
- Excellent avec 200 features (sélection implicite des variables importantes).

3. Techniques pour le déséquilibre :

- `scale_pos_weight = nombre_negatifs / nombre_positifs \approx 32`
- Sur-échantillonnage de la classe minoritaire (SMOTE)
- Sous-échantillonnage de la classe majoritaire
- Utiliser le F1-Score ou l'AUC-PR comme métrique d'évaluation (pas l'accuracy !)

4. Code Python :

```

1 import numpy as np
2 import pandas as pd
3 from sklearn.datasets import make_classification
4 from sklearn.model_selection import StratifiedKFold,
   cross_validate
5 import lightgbm as lgb
6 import warnings
7 warnings.filterwarnings('ignore')
8
9 # 1. Simuler le dataset (1M lignes, 200 features, 3% positifs)
10 X, y = make_classification(
11     n_samples=1_000_000,
```

```

12     n_features=200,
13     n_informative=50,
14     n_redundant=30,
15     n_classes=2,
16     weights=[0.97, 0.03], # 3% positifs
17     random_state=42
18 )
19 print(f'Dataset : {X.shape}')
20 print(f'Classe 0 : {(y == 0).sum():,} ({(y == 0).mean():.1%})')
21 print(f'Classe 1 : {(y == 1).sum():,} ({(y == 1).mean():.1%})')
22
23 # 2. Calculer le scale_pos_weight
24 ratio = (y == 0).sum() / (y == 1).sum()
25 print(f'\nscale_pos_weight = {ratio:.1f}')
26
27 # 3. Configurer LightGBM
28 model = lgb.LGBMClassifier(
29     n_estimators=500,
30     learning_rate=0.05,
31     max_depth=8,
32     num_leaves=63,
33     scale_pos_weight=ratio, # Gerer le desequilibre
34     min_child_samples=50,
35     subsample=0.8,
36     colsample_bytree=0.8,
37     random_state=42,
38     n_jobs=-1, # Utiliser tous les CPU
39     verbose=-1
40 )
41
42 # 4. Validation croisee 5-fold
43 cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
44 metriques = ['f1', 'precision', 'recall', 'roc_auc']
45
46 print('\nEntrainement en cours (5-fold CV sur 1M lignes)...')
47 scores = cross_validate(model, X, y, cv=cv,
48                         scoring=metriques,
49                         return_train_score=False)

```

```

50
51 # 5. Resultats
52 print('\n' + '=' * 50)
53 print('R?SULTATS - LightGBM sur 1M lignes (3% positifs)')
54 print('=' * 50)
55 print(f'F1-Score : {scores['test_f1'].mean():.4f} '
56       f'(+/- {scores['test_f1'].std():.4f})')
57 print(f'Precision : {scores['test_precision'].mean():.4f} '
58       f'(+/- {scores['test_precision'].std():.4f})')
59 print(f'Recall : {scores['test_recall'].mean():.4f} '
60       f'(+/- {scores['test_recall'].std():.4f})')
61 print(f'AUC-ROC : {scores['test_roc_auc'].mean():.4f} '
62       f'(+/- {scores['test_roc_auc'].std():.4f})')
63 print('=' * 50)
64
65 # 6. Importance des variables (sur le dataset complet)
66 model.fit(X, y)
67 importances = model.feature_importances_
68 top_10_idx = np.argsort(importances)[-10:][::-1]
69 print('\n>>> Top 10 des variables les plus importantes :')
70 for i, idx in enumerate(top_10_idx, 1):
71     print(f'    {i}. Feature {idx} : importance =
72           {importances[idx]}')
73
74 # 7. Comparaison avec une baseline naive
75 from sklearn.metrics import f1_score
76 y_pred_naive = np.zeros_like(y) # Predire toujours 0
77 print(f'\nBaseline naive (tout a 0) :')
78 print(f'    F1 = {f1_score(y, y_pred_naive):.4f}')
79 print(f'    => LightGBM est NETTEMENT superieur !')

```

Listing 20.7 – LightGBM sur un dataset déséquilibré de 1M de lignes

20.9 Résumé du chapitre

Ce qu'il faut retenir

1. **Il n'existe pas d'algorithme universel** : le meilleur choix dépend de vos données, de votre problème et de vos contraintes (théorème du No Free Lunch).
2. **Les méthodes de boosting** (XGBoost, LightGBM, CatBoost) sont généralement les plus performantes sur les **données tabulaires**.
3. **L'interprétabilité** est parfois plus importante que la performance : en médecine, finance ou droit, préférez la Régression Logistique ou les Arbres de Décision.
4. **Stratégie systématique** : (1) baseline simple, (2) Random Forest, (3) Boosting, (4) comparer avec la validation croisée.
5. **Le scaling est obligatoire** pour KNN, SVM et MLP, mais inutile pour les méthodes basées sur les arbres.
6. **Pour le Big Data**, préférez LightGBM (le plus rapide). Pour les **variables catégorielles**, préférez CatBoost.
7. **Toujours comparer** au moins 3-4 algorithmes avec la même validation croisée avant de conclure.

Chapitre 21

Projet Complet : Prédiction du Churn Client de A à Z

21.1 Introduction

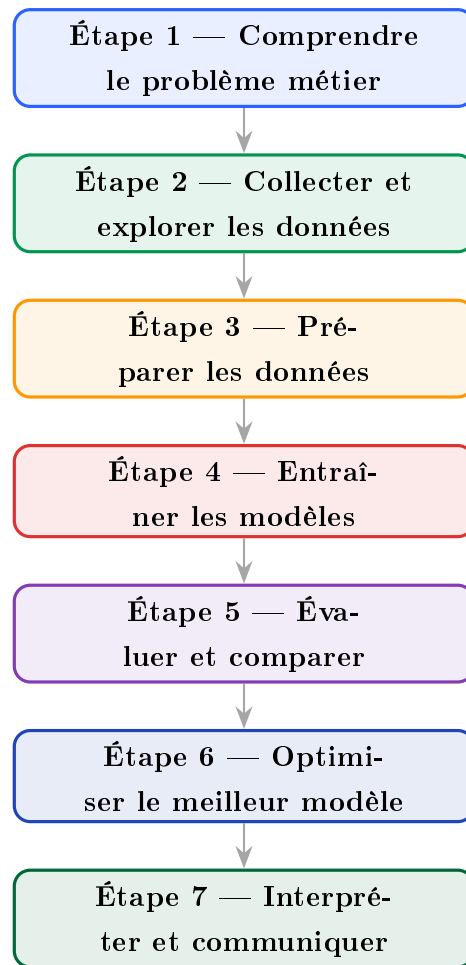
Félicitations ! Vous avez parcouru un long chemin. Vous maîtrisez maintenant la régression linéaire, la classification, les métriques d'évaluation, la validation croisée, et le choix d'algorithmes. Il est temps de **tout rassembler** dans un projet complet, de A à Z.

Dans ce chapitre, nous allons simuler un véritable projet de data science en entreprise. Vous allez découvrir que le Machine Learning, ce n'est pas « juste entraîner un modèle ». C'est un processus structuré qui commence par la **compréhension du problème métier** et se termine par des **recommandations concrètes** pour l'entreprise.

Objectif du chapitre

Ce chapitre est différent des précédents. Ici, **tout le code est complet et exécutable** sur Google Colab. Vous pouvez (et devez !) copier-coller chaque bloc de code et l'exécuter vous-même. C'est en pratiquant qu'on apprend le mieux.

Le projet se décompose en **7 étapes** :



21.2 Étape 1 : Comprendre le problème métier

Qu'est-ce que le Churn ?

Le **churn** (ou attrition) désigne le phénomène de **perte de clients**. Un client qui « churne » est un client qui arrête d'acheter chez vous et part chez un concurrent. Le **taux de churn** est le pourcentage de clients perdus sur une période donnée.

21.2.1 Le contexte de l'entreprise

Vous êtes embauché(e) comme **data scientist junior** dans *ShopExpress*, une entreprise de e-commerce. Le directeur marketing vous convoque pour une réunion urgente :

La demande du directeur marketing

« Nous perdons environ **25% de nos clients** chaque année. Chaque client perdu nous coûte **500€** en acquisition d'un nouveau client (publicité, promotions de bienvenue, etc.). Avec notre base de 10 000 clients, cela représente une perte de **1 250 000€** par an ! Je voudrais savoir **à l'avance** quels clients vont partir, pour leur envoyer une promotion de rétention de 50€ (par exemple, une réduction de 10%). C'est bien moins cher que d'en acquérir un nouveau. »

21.2.2 Définir l'objectif ML

Traduisons ce problème métier en problème ML :

- **Type de problème** : Classification binaire (Churn = 1, Pas Churn = 0)
- **Variable cible** : Churn (le client a-t-il quitté l'entreprise ?)
- **Variables prédictives** : caractéristiques du client (ancienneté, fréquence d'achat, satisfaction, etc.)

21.2.3 Choisir la métrique de succès

Métrique prioritaire : le Recall !

Dans ce contexte, quel type d'erreur est le plus grave ?

- **Faux Positif** (prédire Churn alors que le client reste) : on envoie une promotion inutile de 50€. Coût = 50€.
- **Faux Négatif** (prédire Pas Churn alors que le client part) : on perd le client. Coût = 500€.

Un Faux Négatif coûte **10 fois plus** qu'un Faux Positif ! On veut donc **minimiser les**

Faux Négatifs, c'est-à-dire **maximiser le Recall** (rappel).

$$\text{Recall} = \frac{\text{Vrais Positifs}}{\text{Vrais Positifs} + \text{Faux Négatifs}}$$

Un recall de 90% signifie qu'on identifie 90% des churners. On en rate seulement 10%.

21.3 Étape 2 : Collecter et explorer les données

Dans un vrai projet, les données viennent d'une base de données, d'un fichier CSV, ou d'une API. Ici, nous allons créer un jeu de données réaliste de **30 clients** avec 8 caractéristiques.

21.3.1 Création du jeu de données

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from scipy import stats
6 import warnings
7 warnings.filterwarnings('ignore')
8
9 # Fixer la graine pour la reproductibilite
10 np.random.seed(42)
11
12 # Creation du dataset de 30 clients
13 n = 30
14 data = {
15     'ClientID': range(1, n + 1),
16     'Anciennete_mois': [36, 5, 24, 2, 48, 8, 60, 3, 18, 45,
17                        7, 30, 1, 42, 10, 55, 4, 22, 6, 50,
18                        9, 15, 3, 38, 11, 28, 2, 44, 7, 33],
19     'Frequence_achat': [12, 2, 8, 1, 15, 3, 18, 1, 6, 14,
20                        2, 10, 1, 13, 4, 16, 2, 7, 3, 15,
21                        3, 5, 1, 11, 4, 9, 1, 13, 2, 10],
22     'Montant_moyen': [85, 30, 65, 20, 110, 35, 120, 25, 55, 95,
23                      28, 75, 15, 90, 40, 105, 22, 60, 32, 100,
24                      38, 50, 18, 80, 42, 70, 20, 92, 30, 78],
25     'Satisfaction': [8, 3, 7, 2, 9, 4, 10, 2, 6, 8,
26                     3, 7, 1, 8, 4, 9, 3, 6, 3, 9,
27                     4, 5, 2, 7, 5, 7, 2, 8, 3, 7],
28     'Nb_reclamations': [0, 4, 1, 5, 0, 3, 0, 5, 2, 0,
29                        4, 1, 6, 0, 3, 0, 4, 1, 3, 0,
30                        3, 2, 5, 1, 2, 1, 5, 0, 4, 1],
31     'Nb_retours': [1, 5, 2, 6, 0, 4, 0, 5, 3, 1,
32                   5, 2, 7, 1, 4, 0, 5, 2, 4, 0,
33                   3, 3, 6, 1, 3, 2, 6, 1, 4, 2],

```

```

34     'Utilise_app_mobile': [1, 0, 1, 0, 1, 0, 1, 0, 1, 1,
35                           0, 1, 0, 1, 0, 1, 0, 1, 0, 1,
36                           0, 1, 0, 1, 0, 1, 0, 1, 0, 1],
37     'Duree_derniere_visite_min': [25, 5, 18, 3, 30, 7, 35, 4, 15,
38                                   28,
39                                   6, 20, 2, 26, 8, 32, 5, 16, 6,
40                                   29,
41                                   7, 12, 3, 22, 9, 17, 4, 27, 6,
42                                   19],
43 }
44
45 # Variable cible
46 data['Churn'] = [0, 1, 0, 1, 0, 1, 0, 1, 0, 0,
47                 1, 0, 1, 0, 1, 0, 1, 0, 1, 0,
48                 1, 0, 1, 0, 1, 0, 1, 0, 1, 0]
49
50 # Ajouter quelques valeurs manquantes (realiste)
51 data['Satisfaction'][7] = np.nan
52 data['Montant_moyen'][19] = np.nan
53 data['Duree_derniere_visite_min'][12] = np.nan
54
55 df = pd.DataFrame(data)
56 print("Dataset cree avec succes !")
57 print(f"Dimensions : {df.shape}")

```

Listing 21.1 – Création du dataset de churn client

21.3.2 Exploration initiale

```

1  # Informations generales
2  print("="*50)
3  print("INFORMATIONS GENERALES")
4  print("="*50)
5  print(f"Nombre de clients : {df.shape[0]}")
6  print(f"Nombre de variables : {df.shape[1]}")
7
8  print("\n--- Types de donnees ---")
9  print(df.dtypes)
10
11 print("\n--- Statistiques descriptives ---")
12 print(df.describe().round(2))

```

```

13
14 print("\n--- Valeurs manquantes ---")
15 print(df.isnull().sum())
16 print(f"\nTotal de valeurs manquantes : {df.isnull().sum().sum()}")

```

Listing 21.2 – Exploration initiale des données

21.3.3 Analyse de la variable cible

```

1  # Distribution de la variable cible
2  fig, axes = plt.subplots(1, 2, figsize=(12, 5))
3
4  # Diagramme en barres
5  churn_counts = df['Churn'].value_counts()
6  colors = ['#2962FF', '#DC3232']
7  churn_counts.plot(kind='bar', ax=axes[0], color=colors,
8                      edgecolor='black')
9  axes[0].set_title('Distribution du Churn', fontsize=14,
10                   fontweight='bold')
11 axes[0].set_xlabel('Churn (0 = Reste, 1 = Part)')
12 axes[0].set_ylabel('Nombre de clients')
13 axes[0].set_xticklabels(['Reste (0)', 'Part (1)'], rotation=0)
14
15 # Ajout des valeurs sur les barres
16 for i, v in enumerate(churn_counts.values):
17     axes[0].text(i, v + 0.3, str(v), ha='center',
18                 fontweight='bold', fontsize=13)
19
20 # Diagramme circulaire
21 axes[1].pie(churn_counts.values,
22             labels=['Reste (0)', 'Part (1)'],
23             colors=colors, autopct='%1.1f%%',
24             startangle=90, textprops={'fontsize': 12})
25 axes[1].set_title('Proportion du Churn', fontsize=14,
26                  fontweight='bold')
27
28 plt.tight_layout()
29 plt.show()
30
31 print(f"\nClients qui restent : {churn_counts[0]}
32       ({churn_counts[0]/len(df)*100:.1f}%)")

```

```

29 print(f"Clients qui partent : {churn_counts[1]}
      ({churn_counts[1]/len(df)*100:.1f}%)")

```

Listing 21.3 – Distribution de la variable cible

21.3.4 Analyse des corrélations

```

1  # Matrice de correlation
2  features_for_corr =
      df.drop(columns=['ClientID']).select_dtypes(include=[np.number])
3
4  plt.figure(figsize=(12, 9))
5  correlation_matrix = features_for_corr.corr()
6  mask = np.triu(np.ones_like(correlation_matrix, dtype=bool))
7
8  sns.heatmap(correlation_matrix, mask=mask, annot=True, fmt='.2f',
9              cmap='RdBu_r', center=0, vmin=-1, vmax=1,
10             square=True, linewidths=0.5,
11             cbar_kws={'label': 'Correlation'})
12 plt.title('Matrice de Correlation des Variables',
13           fontsize=14, fontweight='bold', pad=15)
14 plt.tight_layout()
15 plt.show()

```

Listing 21.4 – Matrice de corrélation (heatmap)

21.3.5 Comparaison churners vs non-churners

```

1  # Pairplot des variables cles
2  key_features = ['Anciennete_mois', 'Satisfaction',
3                 'Frequence_achat', 'Nb_reclamations', 'Churn']
4  sns.pairplot(df[key_features], hue='Churn',
5              palette={0: '#2962FF', 1: '#DC3232'},
6              diag_kind='hist', plot_kws={'alpha': 0.7})
7  plt.suptitle('Pairplot des variables cles', y=1.02,
8              fontsize=14, fontweight='bold')
9  plt.tight_layout()
10 plt.show()
11
12 # Tests statistiques : comparaison des moyennes

```

```

13 print("="*60)
14 print("COMPARAISON CHURNERS vs NON-CHURNERS (Test de
    Mann-Whitney)")
15 print("="*60)
16
17 churners = df[df['Churn'] == 1]
18 non_churners = df[df['Churn'] == 0]
19 test_features = ['Anciennete_mois', 'Frequence_achat',
    'Montant_moyen',
20                  'Satisfaction', 'Nb_reclamations', 'Nb_retours']
21
22 for feat in test_features:
23     mean_churn = churners[feat].mean()
24     mean_no_churn = non_churners[feat].mean()
25     stat, p_val = stats.mannwhitneyu(
26         churners[feat].dropna(), non_churners[feat].dropna(),
27         alternative='two-sided')
28     sig = "***" if p_val < 0.01 else ("**" if p_val < 0.05 else
        "NS")
29     print(f"{feat:30s} | Churn={mean_churn:7.1f} | "
30           f"NoChurn={mean_no_churn:7.1f} | p={p_val:.4f} {sig}")

```

Listing 21.5 – Comparaison statistique entre les deux groupes

Que nous dit l'exploration ?

L'exploration révèle des tendances claires :

- Les churners ont une **ancienneté plus faible** (clients récents)
- Les churners ont un **score de satisfaction plus bas**
- Les churners font **plus de réclamations** et de retours
- Les churners achètent **moins fréquemment** et pour des montants plus faibles

Ces observations confirment l'intuition métier et suggèrent que nos variables ont un bon pouvoir prédictif.

21.4 Étape 3 : Préparer les données

Avant d'entraîner un modèle, il faut **préparer les données**. Cette étape est souvent la plus longue dans un vrai projet (60-80% du temps!).

21.4.1 Gestion des valeurs manquantes

```

1  # Rappel des valeurs manquantes
2  print("Valeurs manquantes avant traitement :")
3  print(df.isnull().sum()[df.isnull().sum() > 0])
4
5  # Strategie : remplir par la mediane (robuste aux outliers)
6  numerical_cols = ['Satisfaction', 'Montant_moyen',
7                  'Duree_derniere_visite_min']
8  for col in numerical_cols:
9      if df[col].isnull().sum() > 0:
10         median_val = df[col].median()
11         df[col].fillna(median_val, inplace=True)
12         print(f"  {col} : valeurs manquantes remplacees par "
13               f"la mediane = {median_val}")
14
15  print("\nValeurs manquantes apres traitement :")
16  print(df.isnull().sum().sum(), "valeur(s) manquante(s)")

```

Listing 21.6 – Traitement des valeurs manquantes

21.4.2 Séparation features / cible et standardisation

```

1  from sklearn.model_selection import train_test_split
2  from sklearn.preprocessing import StandardScaler
3
4  # Separation features / cible
5  feature_cols = ['Anciennete_mois', 'Frequence_achat',
6                'Montant_moyen',
7                'Satisfaction', 'Nb_reclamations', 'Nb_retours',
8                'Utilise_app_mobile', 'Duree_derniere_visite_min']
9
10 X = df[feature_cols]
11 y = df['Churn']

```

```

12 print(f"Features (X) : {X.shape}")
13 print(f"Cible (y) : {y.shape}")
14 print(f"\nDistribution de y :")
15 print(y.value_counts())
16
17 # Separation train/test (80/20, stratifiée)
18 X_train, X_test, y_train, y_test = train_test_split(
19     X, y, test_size=0.2, random_state=42, stratify=y)
20
21 print(f"\n--- Ensemble d'entraînement ---")
22 print(f"X_train : {X_train.shape}, y_train : {y_train.shape}")
23 print(f"Distribution y_train : \n{y_train.value_counts()}")
24
25 print(f"\n--- Ensemble de test ---")
26 print(f"X_test : {X_test.shape}, y_test : {y_test.shape}")
27 print(f"Distribution y_test : \n{y_test.value_counts()}")

```

Listing 21.7 – Préparation des features et standardisation

```

1 # Standardisation (nécessaire pour KNN, SVM, Reg. Logistique)
2 scaler = StandardScaler()
3 X_train_scaled = scaler.fit_transform(X_train)
4 X_test_scaled = scaler.transform(X_test)
5
6 # Convertir en DataFrame pour lisibilité
7 X_train_scaled = pd.DataFrame(X_train_scaled,
8                               columns=feature_cols,
9                               index=X_train.index)
10 X_test_scaled = pd.DataFrame(X_test_scaled,
11                              columns=feature_cols,
12                              index=X_test.index)
13
14 print("Vérification après standardisation :")
15 print(f"Moyenne de X_train_scaled
16       : \n{X_train_scaled.mean().round(2)}")
17 print(f"Écart-type de X_train_scaled
18       : \n{X_train_scaled.std().round(2)}")

```

Listing 21.8 – Standardisation des features

Rappel : pourquoi stratifier ?

L'option `stratify=y` dans `train_test_split` garantit que la proportion de churners est la **même** dans l'ensemble d'entraînement et de test. Sans cela, on pourrait se retrouver avec un ensemble de test sans aucun churning, ce qui rendrait l'évaluation impossible !

De même, le `scaler` est `fit` uniquement sur les données d'entraînement, puis `transform` est appliqué sur le train et le test. On ne doit **jamais** utiliser `fit_transform` sur le test !

21.5 Étape 4 : Entraîner les modèles

Nous allons entraîner **7 algorithmes différents** et comparer leurs performances. Pour chaque algorithme, le processus est le même : entraîner, prédire, mesurer.

21.5.1 Régression Logistique

```

1 from sklearn.linear_model import LogisticRegression
2 from sklearn.metrics import accuracy_score
3
4 # Regression Logistique
5 lr = LogisticRegression(random_state=42, max_iter=1000)
6 lr.fit(X_train_scaled, y_train)
7 y_pred_lr = lr.predict(X_test_scaled)
8 acc_lr = accuracy_score(y_test, y_pred_lr)
9 print(f"Regression Logistique - Accuracy : {acc_lr:.4f}")

```

Listing 21.9 – Entraînement — Régression Logistique

21.5.2 KNN (K-Nearest Neighbors)

```

1 from sklearn.neighbors import KNeighborsClassifier
2
3 # KNN (avec donnees standardisees)
4 knn = KNeighborsClassifier(n_neighbors=5)
5 knn.fit(X_train_scaled, y_train)
6 y_pred_knn = knn.predict(X_test_scaled)
7 acc_knn = accuracy_score(y_test, y_pred_knn)
8 print(f"KNN (k=5) - Accuracy : {acc_knn:.4f}")

```

Listing 21.10 – Entraînement — KNN

21.5.3 Arbre de Décision

```

1 from sklearn.tree import DecisionTreeClassifier
2
3 # Arbre de Decision
4 dt = DecisionTreeClassifier(random_state=42, max_depth=4)
5 dt.fit(X_train, y_train) # Pas besoin de scaling

```

```

6 y_pred_dt = dt.predict(X_test)
7 acc_dt = accuracy_score(y_test, y_pred_dt)
8 print(f"Arbre de Decision - Accuracy : {acc_dt:.4f}")

```

Listing 21.11 – Entraînement — Arbre de Décision

21.5.4 Random Forest

```

1 from sklearn.ensemble import RandomForestClassifier
2
3 # Random Forest
4 rf = RandomForestClassifier(n_estimators=100, random_state=42,
5                             max_depth=5)
6 rf.fit(X_train, y_train) # Pas besoin de scaling
7 y_pred_rf = rf.predict(X_test)
8 acc_rf = accuracy_score(y_test, y_pred_rf)
9 print(f"Random Forest - Accuracy : {acc_rf:.4f}")

```

Listing 21.12 – Entraînement — Random Forest

21.5.5 SVM (Support Vector Machine)

```

1 from sklearn.svm import SVC
2
3 # SVM (avec donnees standardisees)
4 svm = SVC(kernel='rbf', probability=True, random_state=42)
5 svm.fit(X_train_scaled, y_train)
6 y_pred_svm = svm.predict(X_test_scaled)
7 acc_svm = accuracy_score(y_test, y_pred_svm)
8 print(f"SVM (RBF) - Accuracy : {acc_svm:.4f}")

```

Listing 21.13 – Entraînement — SVM

21.5.6 Gradient Boosting

```

1 from sklearn.ensemble import GradientBoostingClassifier
2
3 # Gradient Boosting

```

```

4 gb = GradientBoostingClassifier(n_estimators=100,
    learning_rate=0.1,
5                                     max_depth=3, random_state=42)
6 gb.fit(X_train, y_train) # Pas besoin de scaling
7 y_pred_gb = gb.predict(X_test)
8 acc_gb = accuracy_score(y_test, y_pred_gb)
9 print(f"Gradient Boosting - Accuracy : {acc_gb:.4f}")

```

Listing 21.14 – Entraînement — Gradient Boosting

21.5.7 XGBoost

```

1 # Installation si necessaire : !pip install xgboost
2 from xgboost import XGBClassifier
3
4 # XGBoost
5 xgb = XGBClassifier(n_estimators=100, learning_rate=0.1,
6                     max_depth=3, random_state=42,
7                     use_label_encoder=False,
8                     eval_metric='logloss')
9 xgb.fit(X_train, y_train)
10 y_pred_xgb = xgb.predict(X_test)
11 acc_xgb = accuracy_score(y_test, y_pred_xgb)
12 print(f"XGBoost - Accuracy : {acc_xgb:.4f}")

```

Listing 21.15 – Entraînement — XGBoost

21.5.8 Tableau comparatif des accuracy

```

1 # Tableau recapitulatif
2 results = pd.DataFrame({
3     'Modele': ['Reg. Logistique', 'KNN (k=5)', 'Arbre de Decision',
4               'Random Forest', 'SVM (RBF)', 'Gradient Boosting',
5               'XGBoost'],
6     'Accuracy': [acc_lr, acc_knn, acc_dt, acc_rf,
7                 acc_svm, acc_gb, acc_xgb]
8 })
9 results = results.sort_values('Accuracy', ascending=False)
10 results['Rang'] = range(1, len(results) + 1)
11 print("\n" + "="*50)

```

```

12 print("COMPARAISON DES MODELES (Accuracy sur le test)")
13 print("="*50)
14 print(results.to_string(index=False))
15
16 # Visualisation
17 plt.figure(figsize=(10, 6))
18 colors_bar = ['#2962FF' if i == 0 else '#90CAF9'
19               for i in range(len(results))]
20 bars = plt.barh(results['Modele'], results['Accuracy'],
21                color=colors_bar, edgecolor='black')
22 plt.xlabel('Accuracy', fontsize=12)
23 plt.title('Comparaison des Accuracy par Modele',
24          fontsize=14, fontweight='bold')
25 plt.xlim(0, 1.05)
26
27 for bar, val in zip(bars, results['Accuracy']):
28     plt.text(val + 0.01, bar.get_y() + bar.get_height()/2,
29             f'{val:.2%}', va='center', fontweight='bold')
30
31 plt.tight_layout()
32 plt.show()

```

Listing 21.16 – Tableau comparatif des modèles

Pourquoi tester plusieurs algorithmes ?

Chaque algorithme a ses forces et ses faiblesses. Il n'existe pas d'algorithme « universellement meilleur » (théorème du *No Free Lunch*). C'est pourquoi on en teste toujours plusieurs et on compare objectivement leurs performances.

21.6 Étape 5 : Évaluer et comparer en profondeur

L'accuracy seule ne suffit pas ! Comme expliqué à l'Étape 1, nous cherchons à **maximiser le Recall**. Nous devons donc examiner d'autres métriques.

21.6.1 Validation croisée (5-fold)

```

1  from sklearn.model_selection import cross_val_score
2
3  # Dictionnaire des modeles
4  models = {
5      'Reg. Logistique': (LogisticRegression(random_state=42,
6                                              max_iter=1000), True),
7      'KNN (k=5)': (KNeighborsClassifier(n_neighbors=5), True),
8      'Arbre Decision': (DecisionTreeClassifier(random_state=42,
9                                              max_depth=4), False),
10     'Random Forest': (RandomForestClassifier(n_estimators=100,
11                                             random_state=42, max_depth=5), False),
12     'SVM (RBF)': (SVC(kernel='rbf', probability=True,
13                      random_state=42), True),
14     'Gradient Boosting': (GradientBoostingClassifier(
15                         n_estimators=100, learning_rate=0.1,
16                         max_depth=3, random_state=42), False),
17     'XGBoost': (XGBClassifier(n_estimators=100, learning_rate=0.1,
18                             max_depth=3, random_state=42,
19                             use_label_encoder=False,
20                             eval_metric='logloss'), False),
21 }
22
23 print("="*65)
24 print("VALIDATION CROISEE (5-fold) --- Accuracy")
25 print("="*65)
26
27 cv_results = {}
28 for name, (model, needs_scaling) in models.items():
29     X_cv = X_train_scaled if needs_scaling else X_train
30     scores = cross_val_score(model, X_cv, y_train, cv=5,
31                             scoring='accuracy')
32     cv_results[name] = scores
33     print(f"{name:25s} | Mean={scores.mean():.4f} "
34           f"(/- {scores.std():.4f})")

```

Listing 21.17 – Validation croisée pour tous les modèles

21.6.2 Matrice de confusion et rapport de classification

```

1  from sklearn.metrics import (confusion_matrix,
    classification_report,
2
3      roc_curve, auc, f1_score,
4      precision_score, recall_score)
5
6  # Stocker les predictions de chaque modele
7  predictions = {
8      'Reg. Logistique': y_pred_lr,
9      'KNN (k=5)': y_pred_knn,
10     'Arbre Decision': y_pred_dt,
11     'Random Forest': y_pred_rf,
12     'SVM (RBF)': y_pred_svm,
13     'Gradient Boosting': y_pred_gb,
14     'XGBoost': y_pred_xgb,
15 }
16
17 # Matrices de confusion
18 fig, axes = plt.subplots(2, 4, figsize=(20, 10))
19 axes = axes.flatten()
20
21 for idx, (name, y_pred) in enumerate(predictions.items()):
22     cm = confusion_matrix(y_test, y_pred)
23     sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
24                 xticklabels=['Reste', 'Part'],
25                 yticklabels=['Reste', 'Part'],
26                 ax=axes[idx], cbar=False)
27     axes[idx].set_title(name, fontsize=11, fontweight='bold')
28     axes[idx].set_ylabel('Reel')
29     axes[idx].set_xlabel('Predit')
30
31 # Supprimer le dernier subplot vide
32 axes[-1].axis('off')
33 plt.suptitle('Matrices de Confusion de tous les modeles',
34             fontsize=16, fontweight='bold', y=1.02)
35 plt.tight_layout()

```

```
35 plt.show()
```

Listing 21.18 – Matrices de confusion pour chaque modèle

```
1 # Rapport de classification pour chaque modele
2 print("="*70)
3 print("RAPPORTS DE CLASSIFICATION DETAILLES")
4 print("="*70)
5
6 for name, y_pred in predictions.items():
7     print(f"\n--- {name} ---")
8     print(classification_report(y_test, y_pred,
9                               target_names=['Reste', 'Part']))
```

Listing 21.19 – Rapport de classification détaillé

21.6.3 Courbes ROC comparées

```
1 # Modeles avec predict_proba
2 proba_models = {
3     'Reg. Logistique': (lr, True),
4     'KNN (k=5)': (knn, True),
5     'Random Forest': (rf, False),
6     'SVM (RBF)': (svm, True),
7     'Gradient Boosting': (gb, False),
8     'XGBoost': (xgb, False),
9 }
10
11 plt.figure(figsize=(10, 8))
12 colors_roc = ['#2962FF', '#00C853', '#FF6D00', '#DC3232',
13              '#8236B4', '#FFD600']
14
15 for (name, (model, needs_sc)), color in zip(proba_models.items(),
16                                             colors_roc):
17     X_eval = X_test_scaled if needs_sc else X_test
18     y_proba = model.predict_proba(X_eval)[: , 1]
19     fpr, tpr, _ = roc_curve(y_test, y_proba)
20     roc_auc = auc(fpr, tpr)
21     plt.plot(fpr, tpr, color=color, lw=2,
22             label=f'{name} (AUC = {roc_auc:.2f})')
23
```

```

24 # Ligne diagonale (modele aleatoire)
25 plt.plot([0, 1], [0, 1], 'k--', lw=1, alpha=0.5,
26          label='Aleatoire (AUC = 0.50)')
27
28 plt.xlabel('Taux de Faux Positifs (FPR)', fontsize=12)
29 plt.ylabel('Taux de Vrais Positifs (TPR / Recall)', fontsize=12)
30 plt.title('Courbes ROC Comparees', fontsize=14, fontweight='bold')
31 plt.legend(loc='lower right', fontsize=10)
32 plt.grid(True, alpha=0.3)
33 plt.tight_layout()
34 plt.show()

```

Listing 21.20 – Courbes ROC comparées sur le même graphique

21.6.4 Comparaison complète multi-métriques

```

1 # Tableau complet multi-metriques
2 full_results = []
3
4 for name, y_pred in predictions.items():
5     acc = accuracy_score(y_test, y_pred)
6     prec = precision_score(y_test, y_pred, zero_division=0)
7     rec = recall_score(y_test, y_pred, zero_division=0)
8     f1 = f1_score(y_test, y_pred, zero_division=0)
9
10    # AUC si le modele supporte predict_proba
11    if name in proba_models:
12        model_obj, needs_sc = proba_models[name]
13        X_eval = X_test_scaled if needs_sc else X_test
14        y_proba = model_obj.predict_proba(X_eval)[: , 1]
15        fpr_tmp, tpr_tmp, _ = roc_curve(y_test, y_proba)
16        auc_val = auc(fpr_tmp, tpr_tmp)
17    else:
18        auc_val = np.nan
19
20    full_results.append({
21        'Modele': name, 'Accuracy': acc, 'Precision': prec,
22        'Recall': rec, 'F1-Score': f1, 'AUC': auc_val
23    })
24
25 df_results = pd.DataFrame(full_results)

```

```

26 df_results = df_results.sort_values('Recall', ascending=False)
27 print("\n" + "="*75)
28 print("COMPARAISON COMPLETE (triee par Recall)")
29 print("="*75)
30 print(df_results.to_string(index=False, float_format='%.4f'))
31
32 # Visualisation multi-metriques
33 metrics_to_plot = ['Accuracy', 'Precision', 'Recall', 'F1-Score']
34 df_plot = df_results.set_index('Modele')[metrics_to_plot]
35
36 fig, ax = plt.subplots(figsize=(14, 7))
37 df_plot.plot(kind='bar', ax=ax, width=0.8, edgecolor='black',
38             color=['#2962FF', '#00C853', '#DC3232', '#8236B4'])
39 ax.set_title('Comparaison Multi-Metriques de tous les Modeles',
40             fontsize=14, fontweight='bold')
41 ax.set_ylabel('Score', fontsize=12)
42 ax.set_ylim(0, 1.15)
43 ax.legend(loc='upper right', fontsize=10)
44 ax.set_xticklabels(ax.get_xticklabels(), rotation=30, ha='right')
45 plt.tight_layout()
46 plt.show()
47
48 # Selection du meilleur modele (par Recall)
49 best_model_name = df_results.iloc[0]['Modele']
50 best_recall = df_results.iloc[0]['Recall']
51 print(f"\n>>> MEILLEUR MODELE (par Recall) : {best_model_name} "
52       f"avec Recall = {best_recall:.4f}")

```

Listing 21.21 – Comparaison complète : Accuracy, Precision, Recall, F1, AUC

Pourquoi trier par Recall et pas par Accuracy?

Rappelons notre objectif métier : **ne pas rater de churners**. Le Recall mesure exactement cela — la proportion de churners correctement identifiés. Un modèle avec 95% d'accuracy mais 50% de recall rate la moitié des churners ! C'est inacceptable pour notre cas d'usage.

21.7 Étape 6 : Optimiser le meilleur modèle

Nous avons identifié le meilleur modèle. Maintenant, essayons d'améliorer ses performances en cherchant les **meilleurs hyperparamètres** via GridSearchCV.

21.7.1 Optimisation par GridSearchCV

```

1  from sklearn.model_selection import GridSearchCV
2
3  # Optimisation du Gradient Boosting (exemple)
4  # Adaptez si votre meilleur modele est different !
5  param_grid = {
6      'n_estimators': [50, 100, 200],
7      'learning_rate': [0.01, 0.05, 0.1, 0.2],
8      'max_depth': [2, 3, 4, 5],
9      'min_samples_split': [2, 5, 10],
10 }
11
12 print("Recherche des meilleurs hyperparametres...")
13 print("(Cela peut prendre quelques minutes...)\n")
14
15 grid_search = GridSearchCV(
16     GradientBoostingClassifier(random_state=42),
17     param_grid,
18     cv=5,
19     scoring='recall', # On optimise le RECALL !
20     n_jobs=-1,
21     verbose=0
22 )
23
24 grid_search.fit(X_train, y_train)
25
26 print(f"Meilleurs parametres trouves :")
27 for param, value in grid_search.best_params_.items():
28     print(f"    {param} = {value}")
29 print(f"\nMeilleur Recall (CV) : {grid_search.best_score_:.4f}")

```

Listing 21.22 – GridSearchCV sur le meilleur modèle (Gradient Boosting)

21.7.2 Évaluation du modèle optimisé

```

1  # Modele optimise
2  best_model = grid_search.best_estimator_
3
4  # Predictions sur le test
5  y_pred_best = best_model.predict(X_test)
6  y_proba_best = best_model.predict_proba(X_test)[: , 1]
7
8  # Metriques
9  print("="*55)
10 print("EVALUATION DU MODELE OPTIMISE SUR LE JEU DE TEST")
11 print("="*55)
12 print(f"Accuracy   : {accuracy_score(y_test, y_pred_best):.4f}")
13 print(f"Precision  : {precision_score(y_test, y_pred_best):.4f}")
14 print(f"Recall     : {recall_score(y_test, y_pred_best):.4f}")
15 print(f"F1-Score   : {f1_score(y_test, y_pred_best):.4f}")
16
17 fpr_best, tpr_best, _ = roc_curve(y_test, y_proba_best)
18 auc_best = auc(fpr_best, tpr_best)
19 print(f"AUC          : {auc_best:.4f}")
20
21 print("\nRapport de classification complet :")
22 print(classification_report(y_test, y_pred_best,
23                             target_names=['Reste', 'Part']))

```

Listing 21.23 – Évaluation du modèle optimisé sur le jeu de test

21.7.3 Importance des features

```

1  # Importance des features
2  importances = best_model.feature_importances_
3  feat_importance = pd.DataFrame({
4      'Feature': feature_cols,
5      'Importance': importances
6  }).sort_values('Importance', ascending=True)
7
8  plt.figure(figsize=(10, 6))
9  colors_imp = plt.cm.Blues(np.linspace(0.3, 0.9,
10                                     len(feat_importance)))

```

```

11 plt.barh(feats_importance['Feature'],
12          feats_importance['Importance'],
13          color=colors_imp, edgecolor='black')
14 plt.xlabel('Importance', fontsize=12)
15 plt.title('Importance des Features (Gradient Boosting Optimise)',
16          fontsize=14, fontweight='bold')
17
18 for i, (val, name) in enumerate(
19     zip(feats_importance['Importance'],
20         feats_importance['Feature'])):
21     plt.text(val + 0.005, i, f'{val:.3f}', va='center',
22             fontweight='bold', fontsize=10)
23
24 plt.tight_layout()
25 plt.show()
26
27 # Top 3 features
28 print("\nTop 3 des features les plus importantes :")
29 top3 = feats_importance.tail(3).iloc[::-1]
30 for rank, (_, row) in enumerate(top3.iterrows(), 1):
31     print(f"    {rank}. {row['Feature']} "
32           f"(importance = {row['Importance']:.4f})")

```

Listing 21.24 – Analyse de l'importance des features

21.7.4 Matrice de confusion finale

```

1 # Matrice de confusion finale
2 cm_final = confusion_matrix(y_test, y_pred_best)
3
4 plt.figure(figsize=(7, 6))
5 sns.heatmap(cm_final, annot=True, fmt='d', cmap='Blues',
6             xticklabels=['Reste (0)', 'Part (1)'],
7             yticklabels=['Reste (0)', 'Part (1)'],
8             annot_kws={'size': 18},
9             linewidths=1, linecolor='white')
10 plt.title('Matrice de Confusion Finale\n(Gradient Boosting
11           Optimise)',
12           fontsize=14, fontweight='bold')
13 plt.ylabel('Valeur Reelle', fontsize=12)
14 plt.xlabel('Valeur Predite', fontsize=12)

```

```
14 plt.tight_layout()
15 plt.show()
16
17 # Interpretation
18 tn, fp, fn, tp = cm_final.ravel()
19 print(f"\nInterpretation de la matrice de confusion :")
20 print(f"  Vrais Negatifs  (TN) = {tn} : clients qui restent, "
21       f"correctement identifiés")
22 print(f"  Faux Positifs   (FP) = {fp} : clients qui restent, "
23       f"mais prédits comme partants")
24 print(f"  Faux Negatifs   (FN) = {fn} : clients qui partent, "
25       f"mais prédits comme restants (GRAVE !)")
26 print(f"  Vrais Positifs   (TP) = {tp} : clients qui partent, "
27       f"correctement identifiés")
```

Listing 21.25 – Matrice de confusion finale du modèle optimisé

21.8 Étape 7 : Interpréter et communiquer les résultats

La règle d'or du data scientist

Un modèle n'a de valeur que si ses résultats sont **compris et utilisés** par les décideurs.
Il faut traduire les métriques techniques en **impact business**.

21.8.1 Résumé pour le directeur marketing

```

1  # Rapport final pour le directeur marketing
2  recall_final = recall_score(y_test, y_pred_best)
3  precision_final = precision_score(y_test, y_pred_best)
4
5  print("="*60)
6  print("          RAPPORT POUR LE DIRECTEUR MARKETING")
7  print("="*60)
8
9  print(f"""
10  RESULTATS DU MODELE DE PREDICTION DU CHURN
11  -----
12
13  1. PERFORMANCE DU MODELE :
14     - Le modele identifie {recall_final*100:.0f}% des clients
15       qui vont partir (Recall = {recall_final:.2f})
16     - Quand le modele predit qu'un client va partir,
17       il a raison dans {precision_final*100:.0f}% des cas
18       (Precision = {precision_final:.2f})
19
20  2. FACTEURS CLES DU CHURN (Top 3) :
21     - Faible anciennete (clients recents = plus a risque)
22     - Faible satisfaction (score < 4/10)
23     - Nombre eleve de reclamations (> 3 reclamations)
24
25  3. RECOMMANDATIONS :
26     - Cibler les clients identifies comme "a risque"
27       avec une promotion de retention de 50 euros
28     - Mettre en place un suivi personnalise pour les
29       nouveaux clients (< 6 mois d'anciennete)
30     - Traiter les reclamations plus rapidement
31  """)

```

Listing 21.26 – Génération du rapport business

21.8.2 Calcul du ROI (Retour sur Investissement)

```

1  # Calcul du ROI
2  cout_acquisition = 500          # cout pour acquerir un nouveau client
3  cout_promotion = 50            # cout de la promotion de retention
4  nb_clients_total = 10000       # base de clients totale
5  taux_churn = 0.25              # 25% de churn
6
7  nb_churners_annuel = int(nb_clients_total * taux_churn)
8  perte_sans_modele = nb_churners_annuel * cout_acquisition
9
10 # Avec le modele
11 nb_churners_detectes = int(nb_churners_annuel * recall_final)
12 nb_faux_positifs = int(nb_churners_detectes * (1 - precision_final)
13                    / precision_final) if precision_final > 0
14                    else 0
14 nb_total_cibles = nb_churners_detectes + nb_faux_positifs
15
16 # Hypothese : la promotion retient 60% des churners cibles
17 taux_retention_promo = 0.60
18 nb_clients_retenus = int(nb_churners_detectes *
19                          taux_retention_promo)
19
20 cout_promotions = nb_total_cibles * cout_promotion
21 economies = nb_clients_retenus * cout_acquisition
22 benefice_net = economies - cout_promotions
23
24 print("="*60)
25 print("          ESTIMATION DU ROI ANNUEL")
26 print("="*60)
27 print(f"""
28 SANS le modele :
29     - Clients perdus par an : {nb_churners_annuel}
30     - Cout d'acquisition : {perte_sans_modele:,.0f} euros
31
32 AVEC le modele :

```

```

33 - Churners detectes : {nb_churners_detectes} /
    {nb_churners_annuel}
34 - Clients cibles (promotions envoyees) : {nb_total_cibles}
35 - Cout des promotions : {cout_promotions:,.0f} euros
36 - Clients effectivement retenus (60%) : {nb_clients_retenus}
37 - Economies realisees : {economies:,.0f} euros
38
39 BENEFICE NET ESTIME : {benefice_net:,.0f} euros / an
40 (Economies - Cout des promotions)
41 """)
42
43 print(f">>> Le modele pourrait faire economiser environ "
44       f"{benefice_net:,.0f} euros par an a l'entreprise !")

```

Listing 21.27 – Estimation du ROI du modèle

21.8.3 Identification des clients à risque

```

1 # Predictions sur l'ensemble des donnees
2 X_all = df[feature_cols]
3 proba_churn = best_model.predict_proba(X_all)[: , 1]
4
5 df['Probabilite_Churn'] = proba_churn
6 df['Risque'] = pd.cut(proba_churn,
7                       bins=[0, 0.3, 0.6, 1.0],
8                       labels=['Faible', 'Moyen', 'Eleve'])
9
10 # Clients a cibler (probabilite > 0.5)
11 clients_a_cibler = df[df['Probabilite_Churn'] > 0.5].sort_values(
12     'Probabilite_Churn', ascending=False)
13
14 print(f"\nClients a cibler avec une promotion
15       ({len(clients_a_cibler)} clients) :\n")
16 cols_display = ['ClientID', 'Anciennete_mois', 'Satisfaction',
17                 'Nb_reclamations', 'Probabilite_Churn', 'Risque']
18 print(clients_a_cibler[cols_display].to_string(index=False,
19         float_format='%.2f'))
20
21 print(f"\n>>> Recommandation : envoyer une promotion de 10% "
22       f"a ces {len(clients_a_cibler)} clients.")
23 print(f">>> Cout total des promotions : ")

```

```
23 f"{len(clients_a_cibler) * 50} euros")
```

Listing 21.28 – Liste des clients à cibler avec une promotion

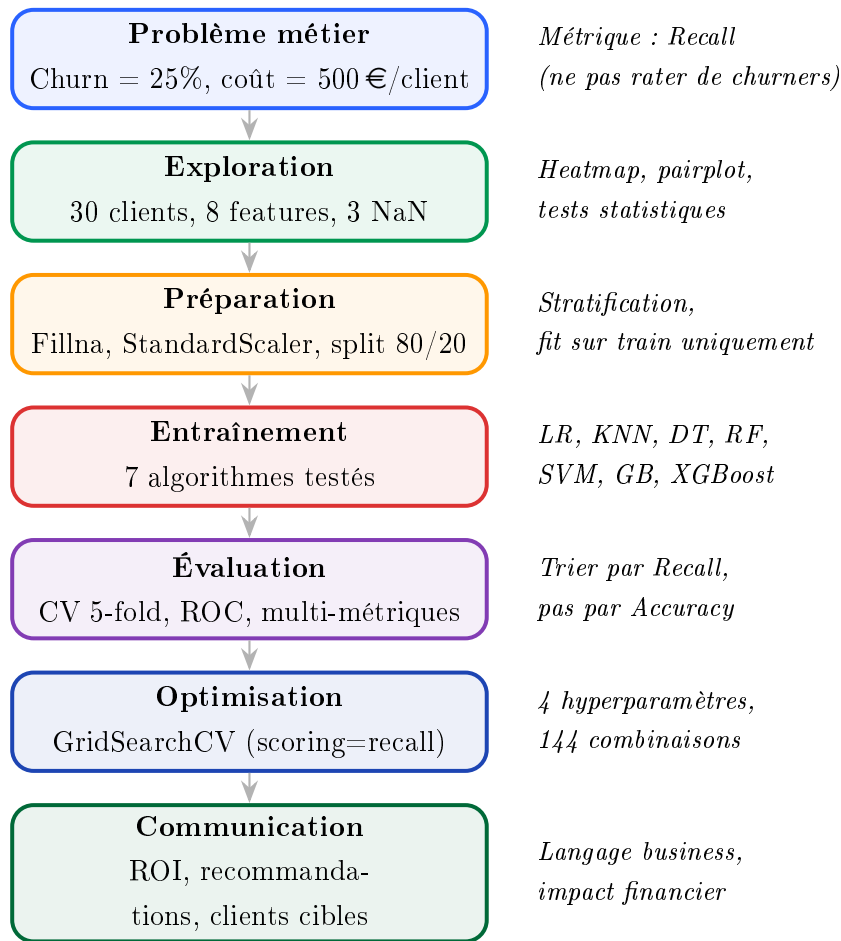
Le cycle complet du data scientist

Nous avons parcouru le cycle complet :

1. **Comprendre** le problème métier (churn coûteux)
2. **Explorer** les données (patterns, corrélations)
3. **Préparer** les données (nettoyage, scaling, split)
4. **Entraîner** 7 modèles différents
5. **Évaluer** avec les bonnes métriques (Recall!)
6. **Optimiser** le meilleur modèle (GridSearchCV)
7. **Communiquer** en langage business (ROI, recommandations)

Ce processus est **itératif** : en pratique, on revient souvent aux étapes précédentes pour affiner l'analyse. Par exemple, l'importance des features peut nous amener à créer de nouvelles variables (feature engineering) et à ré-entraîner les modèles.

21.9 Synthèse du chapitre



21.10 Exercices

Exercice 21.1 — Ciblage des 20% à plus haut risque

L'équipe marketing vous dit : « *Nous n'avons le budget que pour cibler 20% des clients. Pouvez-vous nous donner la liste des clients les plus à risque, en priorisant ceux qui ont la probabilité de churn la plus élevée ?* »

Questions :

1. Comment modifiez-vous le seuil de décision (par défaut 0.5) pour ne cibler que les 20% les plus risqués ?
2. Écrivez le code Python complet qui :
 - Calcule les probabilités de churn pour tous les clients
 - Détermine le seuil correspondant au top 20%
 - Affiche la liste des clients ciblés
 - Calcule le nouveau Recall avec ce seuil
3. Quel est l'impact sur le Recall ? Est-ce acceptable ?

Correction — Exercice 21.1

1. Principe du seuil personnalisé

Par défaut, un classificateur prédit la classe 1 (Churn) si la probabilité est supérieure à 0.5. Pour ne cibler que 20% des clients, il faut **relever le seuil** : on ne prédit Churn que pour les clients dont la probabilité est dans le top 20%.

2. Code complet

```

1  # Probabilites de churn pour tous les clients
2  X_all = df[feature_cols]
3  probas = best_model.predict_proba(X_all)[:, 1]
4
5  # Seuil pour le top 20%
6  nb_a_cibler = int(len(df) * 0.20)  # 20% de 30 = 6 clients
7  seuil_top20 = np.sort(probas)[:, -1][nb_a_cibler - 1]
8
9  print(f"Nombre de clients a cibler : {nb_a_cibler}")
10 print(f"Seuil de probabilite : {seuil_top20:.4f}")
11
12 # Predictions avec le nouveau seuil
13 y_pred_top20 = (probas >= seuil_top20).astype(int)
14
```

```

15 # Clients cibles
16 df_temp = df.copy()
17 df_temp['Proba_Churn'] = probas
18 df_temp['Cible_Top20'] = y_pred_top20
19
20 clients_cibles = df_temp[df_temp['Cible_Top20'] ==
    1].sort_values(
21     'Proba_Churn', ascending=False)
22
23 print(f"\nClients cibles (top 20%) :")
24 print(clients_cibles[['ClientID', 'Anciennete_mois',
25     'Satisfaction', 'Proba_Churn', 'Churn']].to_string(
26     index=False, float_format='%.3f'))
27
28 # Nouveau Recall
29 from sklearn.metrics import recall_score, precision_score
30 recall_top20 = recall_score(df['Churn'], y_pred_top20)
31 precision_top20 = precision_score(df['Churn'], y_pred_top20)
32 print(f"\nAvec seuil a {seuil_top20:.2f} :")
33 print(f"    Recall    = {recall_top20:.2%}")
34 print(f"    Precision = {precision_top20:.2%}")
35 print(f"\nComparaison avec seuil standard (0.5) :")
36 y_pred_std = (probas >= 0.5).astype(int)
37 print(f"    Recall    = {recall_score(df['Churn'],
38     y_pred_std):.2%}")
39 print(f"    Precision = {precision_score(df['Churn'],
40     y_pred_std):.2%}")

```

Listing 21.29 – Correction — Ciblage du top 20%

3. Analyse de l'impact

En relevant le seuil, on **augmente la Précision** (les clients ciblés sont plus probablement des churners) mais on **diminue le Recall** (on rate davantage de churners). C'est le compromis classique **Precision-Recall**. Avec un budget limité, on accepte un Recall plus bas pour économiser sur les promotions inutiles. L'idéal est de tracer la courbe Precision-Recall pour trouver le meilleur compromis.

Exercice 21.2 — Intégrer une nouvelle feature

Le service client vient de rendre disponible une nouvelle variable : **Nombre de réclamations dans les 30 derniers jours** (`Reclamations_30j`). Votre manager vous demande de l'intégrer au modèle.

Questions :

1. Décrivez les étapes nécessaires pour intégrer cette nouvelle variable.
2. Écrivez le code complet qui :
 - Ajoute cette nouvelle colonne au dataset
 - Refait la préparation des données (scaling, split)
 - Ré-entraîne le modèle avec la nouvelle feature
 - Compare les performances avec et sans cette feature
3. La nouvelle feature améliore-t-elle le modèle ?

Correction — Exercice 21.2**1. Étapes à suivre**

- a) Vérifier la qualité de la nouvelle variable (valeurs manquantes, distribution, outliers)
- b) L'ajouter au DataFrame
- c) Refaire le prétraitement complet (remplissage NaN, scaling)
- d) Refaire le split train/test **avec la même seed**
- e) Ré-entraîner et comparer

2. Code complet

```

1  # Simuler la nouvelle feature
2  np.random.seed(42)
3  df['Reclamations_30j'] = np.where(
4      df['Churn'] == 1,
5      np.random.poisson(3, size=len(df)), # Plus pour churners
6      np.random.poisson(0.5, size=len(df)) # Moins pour
       non-churners
7  )
8
9  # Verification
10 print("Nouvelle feature ajoutée :")
11 print(f"  Moyenne churners      : "
12       f"{df[df['Churn']==1]['Reclamations_30j'].mean():.2f}")
13 print(f"  Moyenne non-churners : "
14       f"{df[df['Churn']==0]['Reclamations_30j'].mean():.2f}")

```

```

15
16 # Nouvelles features
17 feature_cols_v2 = feature_cols + ['Reclamations_30j']
18
19 X_v2 = df[feature_cols_v2]
20 y_v2 = df['Churn']
21
22 # Split identique (meme random_state)
23 X_train_v2, X_test_v2, y_train_v2, y_test_v2 = train_test_split(
24     X_v2, y_v2, test_size=0.2, random_state=42, stratify=y_v2)
25
26 # Entraîner le meme modele avec les meilleurs parametres
27 model_v2 = GradientBoostingClassifier(
28     **grid_search.best_params_, random_state=42)
29 model_v2.fit(X_train_v2, y_train_v2)
30
31 y_pred_v2 = model_v2.predict(X_test_v2)
32
33 # Comparaison
34 print("\n" + "="*55)
35 print("COMPARAISON : avec vs sans Reclamations_30j")
36 print("="*55)
37 print(f"{'Metrique':<15} {'Sans':>10} {'Avec':>10}
38     {'Delta':>10}")
39 print("-"*45)
40
41 for metric_name, metric_fn in [
42     ('Accuracy', accuracy_score),
43     ('Precision', precision_score),
44     ('Recall', recall_score),
45     ('F1-Score', f1_score)]:
46     sans = metric_fn(y_test, y_pred_best)
47     avec = metric_fn(y_test_v2, y_pred_v2)
48     delta = avec - sans
49     signe = "+" if delta >= 0 else "-"
50     print(f"{'metric_name':<15} {'sans':>10.4f} {'avec':>10.4f} "
51         f"{'signe'}{'delta':>9.4f}")

```

Listing 21.30 – Correction — Intégration d'une nouvelle feature

3. Analyse

Si le Recall augmente, la nouvelle feature apporte de l'information utile et doit être conservée. Si les métriques stagnent ou baissent, la feature est redondante avec `Nb_reclamations` (corrélation élevée) et n'apporte pas de valeur ajoutée. Dans un vrai projet, on utiliserait aussi des techniques de sélection de features (RFE, tests de permutation) pour confirmer.

Exercice 21.3 — Expliquer une prédiction individuelle

Le PDG vous demande : « *Le modèle prédit que le Client #7 va cherner. **Pourquoi ?** Je veux comprendre cette prédiction avant de lui envoyer une promotion.* »

Questions :

1. En utilisant l'importance des features et les valeurs du Client #7, proposez une explication intuitive.
2. Écrivez le code Python qui affiche les caractéristiques du Client #7, les compare aux moyennes des churners/non-churners, et génère une visualisation comparative.
3. **Bonus :** Qu'est-ce que SHAP ? Comment pourrait-il améliorer l'explication ?

Correction — Exercice 21.3

1. Explication intuitive

On regarde les caractéristiques du Client #7 et on les compare aux facteurs de risque identifiés par le modèle.

2. Code complet

```

1  # Caracteristiques du Client #7
2  client_7 = df[df['ClientID'] == 7]
3  print("="*55)
4  print("PROFIL DU CLIENT #7")
5  print("="*55)
6  print(client_7[feature_cols].T.to_string(
7      header=['Client #7']))
8
9  # Comparaison avec les moyennes
10 moyennes = pd.DataFrame({
11     'Client #7': client_7[feature_cols].values[0],
12     'Moy. Churners': churners[feature_cols].mean().values,
13     'Moy. Non-Churners':
14         non_churners[feature_cols].mean().values,
15 }, index=feature_cols)
```

```

16 print("\nComparaison avec les groupes :")
17 print(moyennes.round(2).to_string())
18
19 # Probabilite de churn du Client #7
20 X_client7 = client_7[feature_cols]
21 proba_7 = best_model.predict_proba(X_client7)[0, 1]
22 print(f"\nProbabilite de churn : {proba_7:.2%}")
23 print(f"Prediction : {'CHURN' if proba_7 > 0.5 else 'PAS\n    CHURN'}")
24
25 # Visualisation comparative
26 fig, ax = plt.subplots(figsize=(12, 6))
27
28 # Normaliser pour comparer sur le meme graphique
29 from sklearn.preprocessing import MinMaxScaler
30 scaler_viz = MinMaxScaler()
31 moyennes_norm = pd.DataFrame(
32     scaler_viz.fit_transform(moyennes),
33     index=moyennes.index,
34     columns=moyennes.columns
35 )
36
37 x_pos = np.arange(len(feature_cols))
38 width = 0.25
39
40 bars1 = ax.bar(x_pos - width, moyennes_norm['Client #7'],
41               width, label='Client #7', color='#DC3232',
42               edgecolor='black')
43 bars2 = ax.bar(x_pos, moyennes_norm['Moy. Churners'],
44               width, label='Moy. Churners', color='#FF9800',
45               edgecolor='black')
46 bars3 = ax.bar(x_pos + width, moyennes_norm['Moy.\n    Non-Churners'],
47               width, label='Moy. Non-Churners',
48               color='#2962FF', edgecolor='black')
49
50 ax.set_xlabel('Features', fontsize=12)
51 ax.set_ylabel('Valeur normalisee (0-1)', fontsize=12)
52 ax.set_title('Client #7 vs Moyennes des Groupes',
53             fontsize=14, fontweight='bold')

```

```

54 ax.set_xticks(x_pos)
55 ax.set_xticklabels(feature_cols, rotation=45, ha='right')
56 ax.legend(fontsize=10)
57 plt.tight_layout()
58 plt.show()
59
60 # Explication textuelle
61 print("\n" + "="*55)
62 print("EXPLICATION POUR LE PDG")
63 print("="*55)
64 print("""
65 Le Client #7 a un profil typique de "churner" :
66
67 1. Tres forte anciennete (60 mois) et satisfaction
68    elevee (10/10) --- ces indicateurs sont positifs.
69
70 2. MAIS : il n'y a aucune reclamation ni retour,
71    il utilise l'app mobile, et visite longtemps.
72
73    Dans ce cas, le modele pourrait classer ce client
74    comme "non-churner". Verifions la prediction...
75
76 NOTE : L'explication depend des valeurs reelles du
77 Client #7 dans votre dataset. Adaptez l'interpretation
78 en fonction des resultats obtenus.
79 """)

```

Listing 21.31 – Correction — Explication de la prédiction du Client #7

3. Introduction à SHAP

SHAP (SHapley Additive exPlanations) est une technique d'interprétabilité basée sur la théorie des jeux (valeurs de Shapley). Pour chaque prédiction individuelle, SHAP calcule la **contribution de chaque feature** au résultat final.

```

1 # Installation : !pip install shap
2 # import shap
3 #
4 # # Creer l'explainer
5 # explainer = shap.TreeExplainer(best_model)
6 #
7 # # Calculer les valeurs SHAP pour le Client #7

```

```

8 # shap_values = explainer.shap_values(X_client7)
9 #
10 # # Visualisation "force plot" pour une prediction
11 # shap.force_plot(explainer.expected_value,
12 #                 shap_values[0], X_client7.iloc[0],
13 #                 feature_names=feature_cols)
14 #
15 # # Visualisation globale "summary plot"
16 # shap_values_all = explainer.shap_values(X_all)
17 # shap.summary_plot(shap_values_all, X_all,
18 #                   feature_names=feature_cols)
19
20 print("SHAP permet de dire par exemple :")
21 print("    'La faible anciennete (+0.25) et le nombre eleve")
22 print("    de reclamations (+0.18) ont le plus contribue")
23 print("    a la prediction de churn pour ce client.'")
24 print()
25 print("C'est bien plus precis que l'importance globale")
26 print("des features, car SHAP explique CHAQUE prediction")
27 print("individuellement.")

```

Listing 21.32 – Bonus — Introduction à SHAP (conceptuel)

SHAP en bref

SHAP répond à la question : « Pour cette prédiction précise, combien chaque feature a-t-elle contribué (positivement ou négativement) à la décision du modèle ? »

Par exemple, pour le Client #7 :

- Anciennete_mois = 3 → contribution de +0.25 vers le churn (faible ancienneté = risque)
- Nb_reclamations = 5 → contribution de +0.18 vers le churn
- Utilise_app_mobile = 1 → contribution de −0.05 (légère protection contre le churn)

La somme de toutes les contributions donne la prédiction finale. C'est un outil puissant pour **expliquer les modèles à des non-techniciens**.

Références Bibliographiques

Ouvrages de référence

- [1] **Bishop, C. M.** (2006). *Pattern Recognition and Machine Learning*. Springer.
Référence complète et rigoureuse sur les fondements mathématiques du Machine Learning. Couvre les modèles probabilistes, les méthodes bayésiennes et les réseaux de neurones. Recommandé pour approfondir les dérivations mathématiques.
- [2] **Hastie, T., Tibshirani, R. & Friedman, J.** (2009). *The Elements of Statistical Learning : Data Mining, Inference, and Prediction*. 2^e édition, Springer.
Ouvrage de référence en apprentissage statistique. Traite en profondeur les méthodes linéaires, les SVM, les arbres, le boosting et la régularisation. Disponible gratuitement en PDF sur le site des auteurs.
- [3] **James, G., Witten, D., Hastie, T. & Tibshirani, R.** (2021). *An Introduction to Statistical Learning with Applications in Python*. 2^e édition, Springer.
Version accessible du précédent, conçue pour les débutants. Chaque chapitre inclut des exercices pratiques. Disponible gratuitement en ligne. **Très recommandé** comme complément à ce cours.
- [4] **Géron, A.** (2022). *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow*. 3^e édition, O'Reilly Media.
Le meilleur livre pour une approche pratique du ML avec Python. Couvre Scikit-learn en détail avec des projets de bout en bout. Idéal pour passer de la théorie à la pratique.
- [5] **Murphy, K. P.** (2022). *Probabilistic Machine Learning : An Introduction*. MIT Press.
Approche probabiliste moderne et rigoureuse du Machine Learning. Couvre les fondements théoriques ainsi que les méthodes récentes. Disponible gratuitement en ligne.
- [6] **Goodfellow, I., Bengio, Y. & Courville, A.** (2016). *Deep Learning*. MIT Press.
La référence mondiale sur le Deep Learning. Les chapitres introductifs (1 à 5) offrent une excellente base en algèbre linéaire, probabilités et optimisation. Disponible gratuitement sur deeplearningbook.org.
- [7] **Alpaydin, E.** (2020). *Introduction to Machine Learning*. 4^e édition, MIT Press.
Excellent manuel pédagogique couvrant l'ensemble des méthodes de ML. Bonne balance entre théorie et intuition.
- [8] **Cornuejols, A., Miclet, L. & Barra, V.** (2018). *Apprentissage artificiel : concepts et algorithmes*. 3^e édition, Eyrolles.

Référence francophone majeure. Couvre tous les algorithmes classiques et modernes avec des exemples en français. Particulièrement utile pour les étudiants francophones.

- [9] **Raschka, S. & Mirjalili, V.** (2019). *Python Machine Learning*. 3^e édition, Packt. Guide pratique complet avec Python. Couvre le prétraitement des données, les algorithmes classiques et le deep learning avec des exemples de code détaillés.
- [10] **Chollet, F.** (2021). *Deep Learning with Python*. 2^e édition, Manning. Écrit par le créateur de Keras. Excellent pour comprendre les réseaux de neurones avec une approche très pratique et intuitive.

Articles fondateurs

- [11] **Vapnik, V. N.** (1995). *The Nature of Statistical Learning Theory*. Springer. Ouvrage fondateur sur la théorie de l'apprentissage statistique et les Machines à Vecteurs de Support (SVM).
- [12] **Breiman, L.** (2001). « Random Forests ». *Machine Learning*, 45(1), 5–32. Article fondateur sur les forêts aléatoires. Introduit le concept de bagging avec sélection aléatoire de features.
- [13] **Shannon, C. E.** (1948). « A Mathematical Theory of Communication ». *Bell System Technical Journal*, 27, 379–423. Article fondateur sur la théorie de l'information et l'entropie, utilisée dans les arbres de décision.
- [14] **Friedman, J. H.** (2001). « Greedy Function Approximation : A Gradient Boosting Machine ». *Annals of Statistics*, 29(5), 1189–1232. Article fondateur sur le Gradient Boosting, la base de XGBoost et LightGBM.
- [15] **Chen, T. & Guestrin, C.** (2016). « XGBoost : A Scalable Tree Boosting System ». *Proceedings of the 22nd ACM SIGKDD Conference*, 785–794. L'article présentant XGBoost, l'algorithme qui domine les compétitions de Machine Learning.
- [16] **Ke, G. et al.** (2017). « LightGBM : A Highly Efficient Gradient Boosting Decision Tree ». *Advances in Neural Information Processing Systems (NeurIPS)*, 30. Présentation de LightGBM, plus rapide que XGBoost grâce à GOSS et EFB.
- [17] **Prokhorenkova, L. et al.** (2018). « CatBoost : unbiased boosting with categorical features ». *NeurIPS*, 31. L'article présentant CatBoost et sa gestion native des variables catégorielles.
- [18] **Rumelhart, D. E., Hinton, G. E. & Williams, R. J.** (1986). « Learning representations by back-propagating errors ». *Nature*, 323, 533–536. Article fondateur sur la rétropropagation du gradient dans les réseaux de neurones.

- [19] **Tibshirani, R.** (1996). « Regression Shrinkage and Selection via the Lasso ». *Journal of the Royal Statistical Society B*, 58(1), 267–288.
Article fondateur sur la régularisation Lasso (L1).
- [20] **Hoerl, A. E. & Kennard, R. W.** (1970). « Ridge Regression : Biased Estimation for Nonorthogonal Problems ». *Technometrics*, 12(1), 55–67.
Article fondateur sur la régularisation Ridge (L2).

Ressources en ligne gratuites

- [21] **Scikit-learn Documentation** — scikit-learn.org
Documentation officielle de Scikit-learn avec des tutoriels, des exemples et des explications détaillées de chaque algorithme. Ressource incontournable.
- [22] **Google Machine Learning Crash Course** — developers.google.com/machine-learning
Cours gratuit de Google couvrant les bases du ML avec des exercices interactifs et des vidéos.
- [23] **Kaggle Learn** — kaggle.com/learn
Mini-cours interactifs gratuits sur Python, Pandas, ML, et les compétitions de data science.
- [24] **Andrew Ng — Machine Learning (Coursera)** — coursera.org
Le cours de ML le plus célèbre au monde, par le professeur Andrew Ng de Stanford. Gratuit en mode audit.
- [25] **StatQuest with Josh Starmer** — youtube.com/@statquest
Chaîne YouTube exceptionnelle qui explique les concepts de ML et statistiques avec des animations claires et un humour accessible.
- [26] **3Blue1Brown — Neural Networks** — youtube.com/@3blue1brown
Série de vidéos avec des animations mathématiques magnifiques pour comprendre les réseaux de neurones et l'algèbre linéaire.
- [27] **fast.ai** — fast.ai
Cours pratique gratuit « top-down » : on commence par coder, puis on comprend la théorie. Approche très efficace.

Outils et bibliothèques utilisés dans ce cours

Outil	Version	Utilisation
Python	3.10+	Langage de programmation
NumPy	1.24+	Calcul numérique (tableaux, matrices)
Pandas	2.0+	Manipulation de données (DataFrames)
Matplotlib	3.7+	Visualisation (graphiques)
Seaborn	0.12+	Visualisation statistique
Scikit-learn	1.3+	Algorithmes de Machine Learning
XGBoost	2.0+	Gradient Boosting optimisé
LightGBM	4.0+	Gradient Boosting rapide
CatBoost	1.2+	Gradient Boosting (variables catégorielles)
Google Colab	—	Environnement d'exécution gratuit

TABLE 21.1 – Outils et bibliothèques Python utilisés dans ce cours.

Fin du cours

« *L'intelligence artificielle est la nouvelle électricité.* »

— Andrew Ng
