

MACHINE LEARNING POUR LA LUTTE CONTRE LA FRAUDE

VOIR GRAND ET TOUT PETIT, MAIS JUSTE

Théo Lopès-Quintas

BPCE Payment Services,
Université Paris Dauphine

9 décembre 2024

1	Contexte	1
2	Hyper déséquilibre de classe : voir tout petit	2
2.1	Sur-échantillonnage aléatoire	2
2.2	Sous-échantillonnage aléatoire	4
2.3	SMOTE	6
3	Travail en grande dimension : voir grand	11
4	Calibration : voir juste	16

CONTEXTE

FORMALISATION DU PROBLÈME MACHINE LEARNING

Dans le cadre supervisé, nous avons accès à un dataset \mathcal{D} défini comme :

$$\mathcal{D} = \left\{ (x_i, y_i) \mid \forall i \leq n, x_i \in \mathbb{R}^{d'}, y_i \in \mathcal{Y} \right\} \quad (1)$$

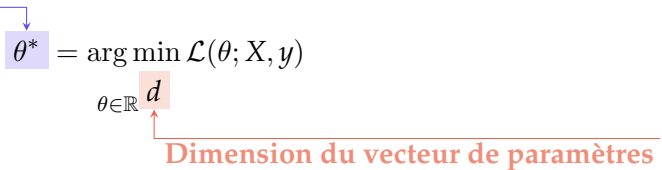
Nombre d'observations

Nombre d'informations

Avec $\mathcal{Y} \subseteq \mathbb{R}$ pour un problème de régression et $\mathcal{Y} \subset \mathbb{N}$ dans le cadre d'une classification. Les problèmes de Machine Learning supervisés peuvent souvent s'écrire sous la forme d'une optimisation d'une fonction de perte $\mathcal{L} : \mathbb{R}^d \times \mathcal{M}_{n,d'} \times \mathbb{R}^n \rightarrow \mathbb{R}_+$ comme :

Vecteur des paramètres optimaux

$$\theta^* = \arg \min_{\theta \in \mathbb{R}^d} \mathcal{L}(\theta; X, y)$$



Dimension du vecteur de paramètres

Dans la suite, pour simplifier les notations, nous omettrons la dépendance de \mathcal{L} en X et y . Notons qu'en général, nous avons $d \neq d'$ et dans le cas du deep learning, très souvent $d \gg d'$.

HYPER DÉSÉQUILIBRE DE CLASSE : VOIR TOUT PETIT

SUR-ÉCHANTILLONNAGE ALÉATOIRE

Une première méthode pour équilibrer le dataset est de dupliquer aléatoirement des observations de la classe minoritaire, jusqu'à ce qu'on atteigne la proportion de la classe minoritaire voulue p_f .

HYPER DÉSÉQUILIBRE DE CLASSE : VOIR TOUT PETIT

SUR-ÉCHANTILLONNAGE ALÉATOIRE

Une première méthode pour équilibrer le dataset est de dupliquer aléatoirement des observations de la classe minoritaire, jusqu'à ce qu'on atteigne la proportion de la classe minoritaire voulue p_f .

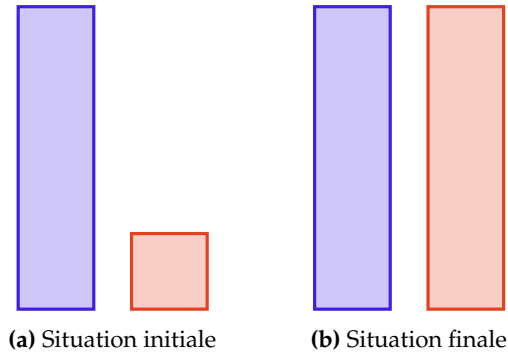


Figure – Illustration du sur-échantillonnage

HYPER DÉSÉQUILIBRE DE CLASSE : VOIR TOUT PETIT

SUR-ÉCHANTILLONNAGE ALÉATOIRE

Une première méthode pour équilibrer le dataset est de dupliquer aléatoirement des observations de la classe minoritaire, jusqu'à ce qu'on atteigne la proportion de la classe minoritaire voulue p_f .

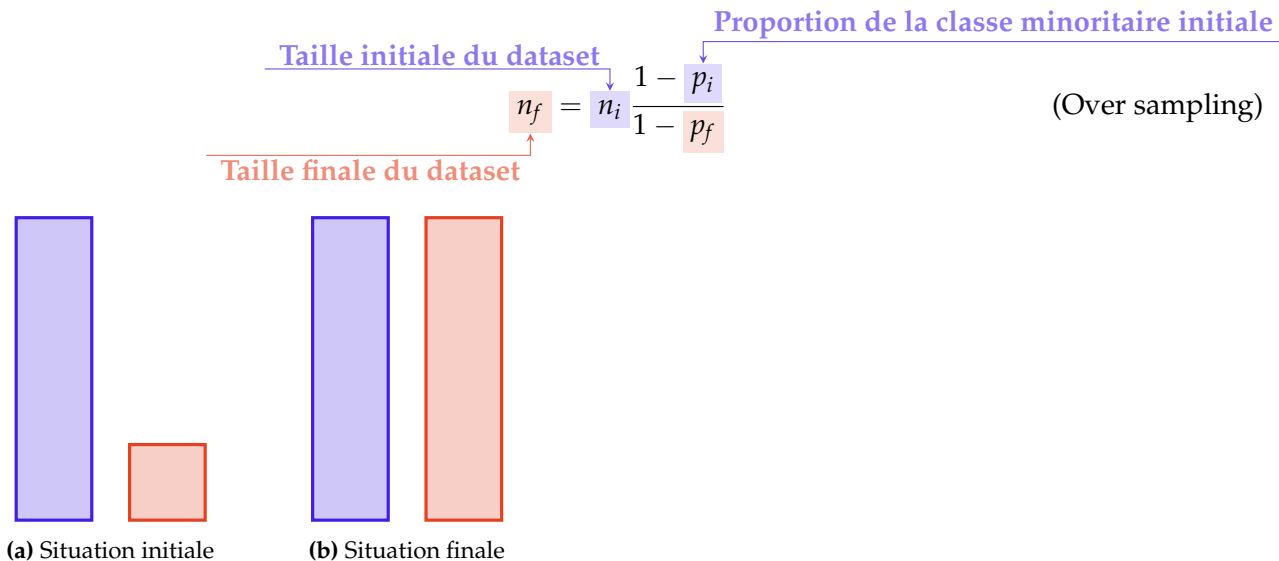


Figure – Illustration du sur-échantillonnage

HYPER DÉSÉQUILIBRE DE CLASSE : VOIR TOUT PETIT

SUR-ÉCHANTILLONNAGE ALÉATOIRE

Une première méthode pour équilibrer le dataset est de dupliquer aléatoirement des observations de la classe minoritaire, jusqu'à ce qu'on atteigne la proportion de la classe minoritaire voulue p_f .

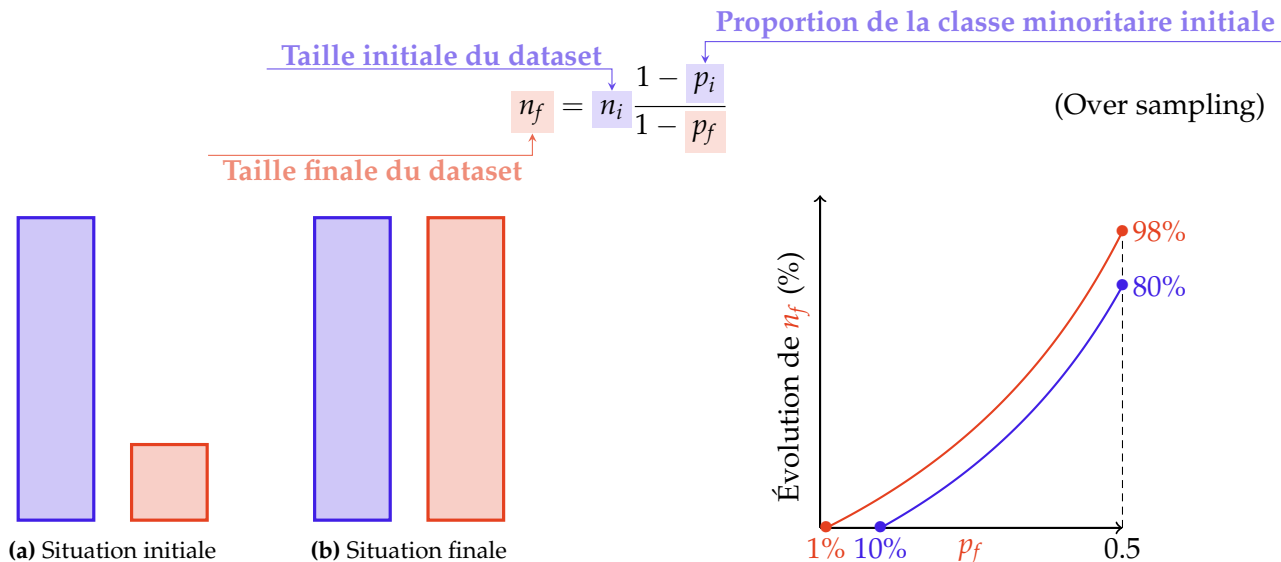


Figure – Illustration du sur-échantillonnage

Figure – Évolution de la taille du dataset final en fonction de la proportion finale de la classe minoritaire

HYPER DÉSÉQUILIBRE DE CLASSE : VOIR TOUT PETIT

SUR-ÉCHANTILLONNAGE ALÉATOIRE

$$n_f = n_i \frac{1 - p_i}{1 - p_f}$$

Diagram illustrating the formula for over-sampling to balance a dataset:

- Taille initiale du dataset** points to n_i .
- Taille finale du dataset** points to n_f .
- Proportion de la classe minoritaire initiale** points to p_i .
- Proportion de la classe minoritaire finale** points to p_f .

(Over sampling)

Exercice 1 (Paramétrer un arbre)

Nous avons accès à un dataset d'un million de lignes avec une proportion initiale de la classe minoritaire de 1%. Après avoir équilibré le dataset avec du sur-échantillonnage, vous utilisez un arbre ou une méthode à base d'arbres.

1. Combien y avait-il d'observations de la classe minoritaire au début ? A la fin ?
2. En moyenne, combien de fois une même observation de la classe minoritaire est présente dans le dataset final ?
3. Comment ajuster en conséquence les paramètres d'un arbre ?

HYPER DÉSÉQUILIBRE DE CLASSE : VOIR TOUT PETIT

SOUS-ÉCHANTILLONNAGE ALÉATOIRE

Une deuxième méthode pour équilibrer le dataset est de supprimer aléatoirement des observations de la classe majoritaire, jusqu'à ce qu'on atteigne la proportion de la classe minoritaire voulue p_f .

HYPER DÉSÉQUILIBRE DE CLASSE : VOIR TOUT PETIT

SOUS-ÉCHANTILLONNAGE ALÉATOIRE

Une deuxième méthode pour équilibrer le dataset est de supprimer aléatoirement des observations de la classe majoritaire, jusqu'à ce qu'on atteigne la proportion de la classe minoritaire voulue p_f .

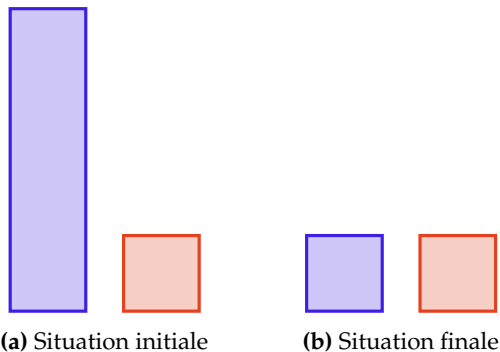


Figure – Illustration du sous-échantillonnage

HYPER DÉSÉQUILIBRE DE CLASSE : VOIR TOUT PETIT

SOUS-ÉCHANTILLONNAGE ALÉATOIRE

Une deuxième méthode pour équilibrer le dataset est de supprimer aléatoirement des observations de la classe majoritaire, jusqu'à ce qu'on atteigne la proportion de la classe minoritaire voulue p_f .

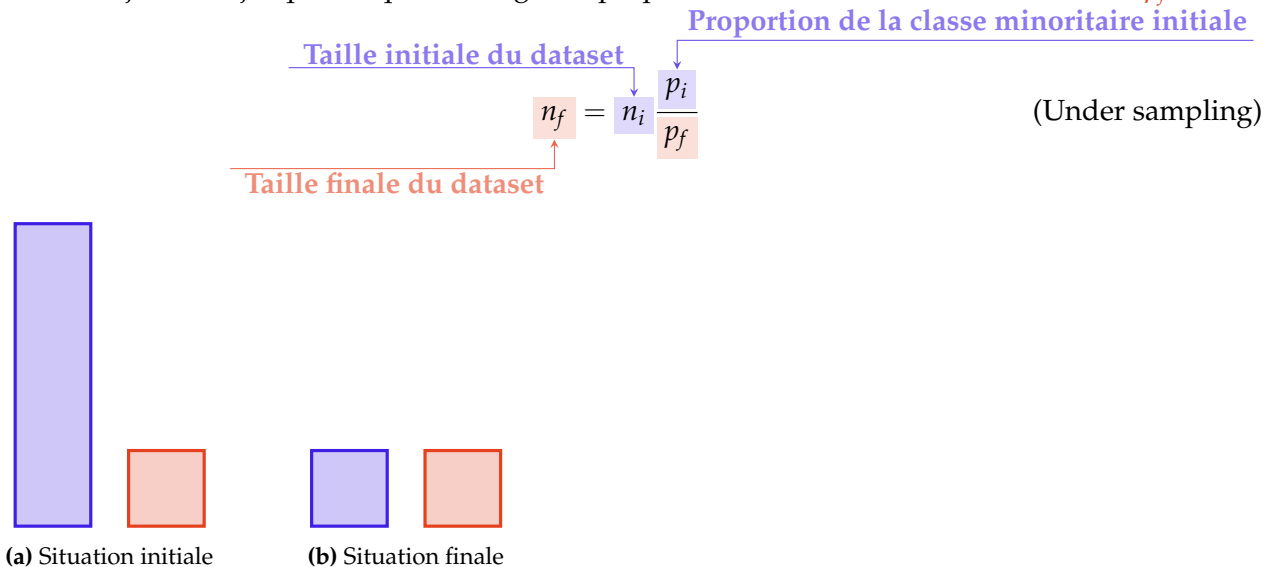


Figure – Illustration du sous-échantillonnage

HYPER DÉSÉQUILIBRE DE CLASSE : VOIR TOUT PETIT

SOUS-ÉCHANTILLONNAGE ALÉATOIRE

Une deuxième méthode pour équilibrer le dataset est de supprimer aléatoirement des observations de la classe majoritaire, jusqu'à ce qu'on atteigne la proportion de la classe minoritaire voulue p_f .

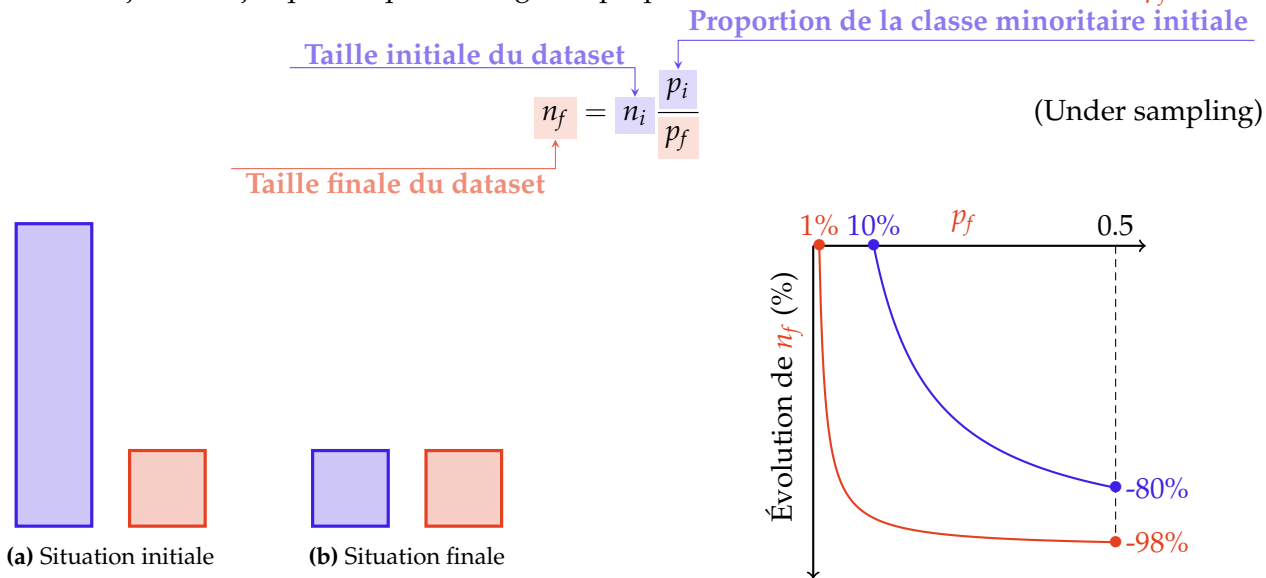
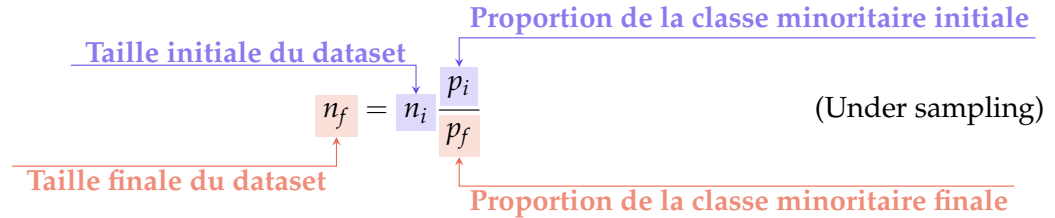


Figure – Illustration du sous-échantillonnage

Figure – Évolution de la taille du dataset final en fonction de la proportion finale de la classe minoritaire

HYPER DÉSÉQUILIBRE DE CLASSE : VOIR TOUT PETIT

SOUS-ÉCHANTILLONNAGE ALÉATOIRE



Exercice 2 (Conserver la représentativité)

Nous avons accès à un dataset d'un million de lignes avec une proportion initiale de la classe minoritaire de 1%. Après avoir équilibré le dataset avec du sous-échantillonnage, vous explorez les données.

1. Combien y avait-il d'observations de la classe majoritaire au début ? A la fin ?
2. Quelle est la probabilité qu'une observation de la classe majoritaire soit présente dans le dataset final ?
3. Supposons qu'un groupe d'observations (par exemple les transactions par téléphone ou mail) représente 5% des observations de la classe majoritaire. Quelle est la probabilité qu'il ne soit plus présent dans le dataset final ?
4. Quel problème cela induit-il sur le modèle à construire en production ? Comment s'en prémunir ?

HYPER DÉSÉQUILIBRE DE CLASSE : VOIR TOUT PETIT

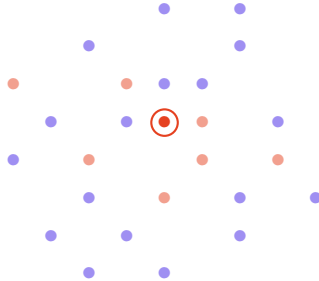
SMOTE : UNE APPROCHE PROMETTEUSE

Les deux précédentes approches dupliquent ou suppriment des observations du dataset. On peut explorer la possibilité de *créer* des observations synthétiques : c'est l'objet de SMOTE.

HYPER DÉSÉQUILIBRE DE CLASSE : VOIR TOUT PETIT

SMOTE : UNE APPROCHE PROMETTEUSE

Les deux précédentes approches dupliquent ou suppriment des observations du dataset. On peut explorer la possibilité de *créer* des observations synthétiques : c'est l'objet de SMOTE.



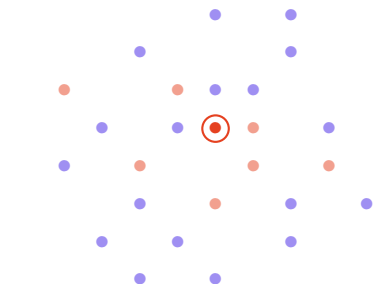
(a) Sélection aléatoire d'une observation de la classe minoritaire

Figure – Fonctionnement de SMOTE

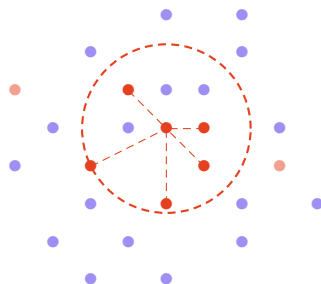
HYPER DÉSÉQUILIBRE DE CLASSE : VOIR TOUT PETIT

SMOTE : UNE APPROCHE PROMETTEUSE

Les deux précédentes approches dupliquent ou suppriment des observations du dataset. On peut explorer la possibilité de *créer* des observations synthétiques : c'est l'objet de SMOTE.



(a) Sélection aléatoire d'une observation de la classe minoritaire



(b) Identification des $k=5$ voisins les plus proches

Figure – Fonctionnement de SMOTE

HYPER DÉSÉQUILIBRE DE CLASSE : VOIR TOUT PETIT

SMOTE : UNE APPROCHE PROMETTEUSE

Les deux précédentes approches dupliquent ou suppriment des observations du dataset. On peut explorer la possibilité de *créer* des observations synthétiques : c'est l'objet de SMOTE.

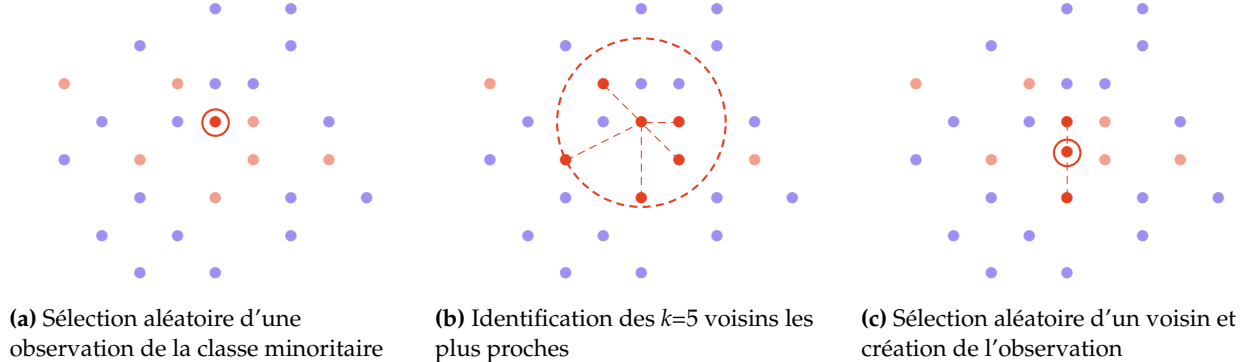


Figure – Fonctionnement de SMOTE

HYPER DÉSÉQUILIBRE DE CLASSE : VOIR TOUT PETIT

SMOTE : OBSERVATIONS SYNTHÉTIQUES AMBIGUËS

SMOTE et One Hot Encoding

SMOTE est une méthode qui ne travaille que les datasets avec des données quantitatives. Elle fonctionne donc théoriquement sur des informations issues d'un One Hot Encoding mais on ne doit pas l'utiliser en pratique : cela n'a aucun sens. Il faut exploiter une variante de SMOTE : SMOTE-NC pour pouvoir traiter des données qualitatives.

Soit \mathcal{D} un dataset comme défini dans (1) dont on reprend les notations, et \mathcal{D}^- le dataset contenant uniquement les observations de la classe majoritaire. On note $n^- = \#\mathcal{D}^-$.

Soit \tilde{x} un point synthétique généré par l'algorithme SMOTE. On dit que \tilde{x} est un **point ambigu** si

$\min_{x \in \mathcal{D}^-} \|\tilde{x} - x\| \leq \delta$ pour $\delta \geq 0$ un paramètre fixé.

On peut montrer, sous l'hypothèse que $\mathcal{D} \sim \mathcal{N}(0_d, I_d)$, que :

$$\mathbb{P} \left(\min_{x \in \mathcal{D}^-} \|\tilde{x} - x\| \leq \delta \right) = 1 - \left(1 - \int_0^{\frac{\delta^2}{2}} \frac{t^{\frac{d}{2}-1} e^{-\frac{t}{2}}}{2^{\frac{d}{2}} \Gamma\left(\frac{d}{2}\right)} dt \right)^{n^-}$$

HYPER DÉSÉQUILIBRE DE CLASSE : VOIR TOUT PETIT

SMOTE : OBSERVATIONS SYNTHÉTIQUES AMBIGUËS - OBSERVATION SUR UN DATASET RÉEL

Nous souhaitons appliquer cela à un dataset de 20 millions de lignes ayant 1% de déséquilibre.

Déséquilibre final (%)	Observation synthétique ambiguë (%)	Taille du dataset final (en millions)
1	0	20.000
2	39.29	20.204
5	62.86	20.842
10	70.72	22.000
20	74.65	24.750
30	75.96	28.285
40	76.61	33.000
50	77.01	39.600

Table – Évolution de la proportion d'observation synthétique ambiguë et de la taille du dataset final en fonction de la proportion de déséquilibre souhaité

Nous avons fixé la valeur de δ comme le quantile à 0.5% de la loi du χ^2 induite par le dataset.

HYPER DÉSÉQUILIBRE DE CLASSE : VOIR TOUT PETIT

SMOTE : OBSERVATIONS SYNTHÉTIQUES AMBIGUËS

$$\mathbb{P} \left(\min_{x \in \mathcal{D}^-} \|\tilde{x} - x\| \leq \delta \right) = 1 - \left(1 - \int_0^{\frac{\delta^2}{2}} \frac{t^{\frac{d}{2}-1} e^{-\frac{t}{2}}}{2^{\frac{d}{2}} \Gamma\left(\frac{d}{2}\right)} dt \right)^{n^-}$$

Exercice 3 (Inducteurs d'observations synthétiques ambiguës)

1. Étudier le comportement asymptotique selon δ . Interpréter.
2. Calculer $\lim_{n^- \rightarrow +\infty} \mathbb{P} \left(\min_{x \in \mathcal{D}^-} \|\tilde{x} - x\| \leq \delta \right)$ et commentez.
3. Que peut-on dire de l'évolution de la probabilité selon la dimension du dataset ?

HYPER DÉSÉQUILIBRE DE CLASSE : VOIR TOUT PETIT

EN RÉSUMÉ

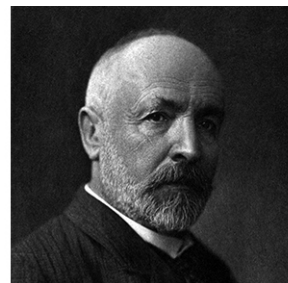
- ▶ L'hyper déséquilibre de classe introduit de nombreux challenges et exacerbe certains phénomènes classiques comme le drift par exemple.
- ▶ Les solutions théoriques sont à explorer et à adapter à chaque cas d'usage. Il est possible qu'aucune ne fonctionne.

TRAVAIL EN GRANDE DIMENSION : VOIR GRAND

MOTIVATION

Tant que vous ne m'aurez pas approuvé, je ne puis que dire : je le vois mais je ne le crois pas.

— Georg Cantor (1877)



TRAVAIL EN GRANDE DIMENSION : VOIR GRAND

HYPERSPHÈRE



(a) En dimension 1, volume = 2

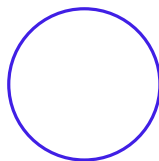
Figure – Représentation et volume d'une hypersphère de rayon 1 dans 3 espaces de dimensions différentes

TRAVAIL EN GRANDE DIMENSION : VOIR GRAND

HYPERSPHÈRE

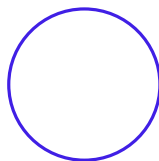


(a) En dimension 1, volume = 2



(b) En dimension 2, volume = π

Figure – Représentation et volume d'une hypersphère de rayon 1 dans 3 espaces de dimensions différentes



(b) Avec la norme 2

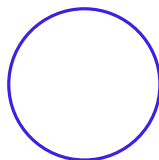
Figure – Représentation d'une hypersphère de rayon 1 en dimension 2 pour 3 normes différentes

TRAVAIL EN GRANDE DIMENSION : VOIR GRAND

HYPERSPHÈRE

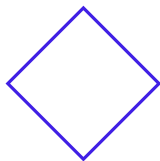


(a) En dimension 1, volume = 2

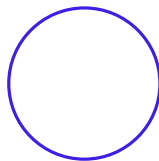


(b) En dimension 2, volume = π

Figure – Représentation et volume d'une hypersphère de rayon 1 dans 3 espaces de dimensions différentes



(a) Avec la norme 1



(b) Avec la norme 2



(c) Avec la norme infinie

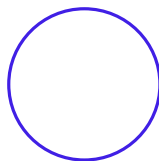
Figure – Représentation d'une hypersphère de rayon 1 en dimension 2 pour 3 normes différentes

TRAVAIL EN GRANDE DIMENSION : VOIR GRAND

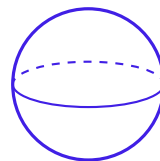
HYPERSPHÈRE



(a) En dimension 1, volume = 2

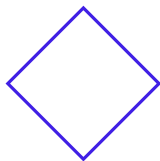


(b) En dimension 2, volume = π

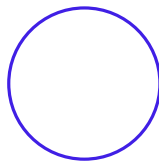


(c) En dimension 3, volume = $\frac{4}{3}\pi$

Figure – Représentation et volume d'une hypersphère de rayon 1 dans 3 espaces de dimensions différentes



(a) Avec la norme 1



(b) Avec la norme 2



(c) Avec la norme infinie

Figure – Représentation d'une hypersphère de rayon 1 en dimension 2 pour 3 normes différentes

TRAVAIL EN GRANDE DIMENSION : VOIR GRAND

VOLUME D'UNE HYPERSPHERE

On appelle *boule* ou hypersphère l'objet défini par :

$$B_n^p(R) = \{u \in \mathbb{R}^n, \|u\|_p^p \leq R^p\}$$

Et son volume par :

$$V_n^p(R) = \int_{B_n^p(R)} \bigotimes_{i=1}^n dx_i$$

Proposition 1 (Volume d'une hypersphere)

Avec les notations précédentes, on a :

$$\begin{aligned} \forall R > 0, \forall n \geq 2, \forall p \geq 1, \quad V_n^p(R) &= \frac{\left(2R\Gamma\left(\frac{1}{p} + 1\right)\right)^n}{\Gamma\left(\frac{n}{p} + 1\right)} \\ &\sim \sqrt{\frac{p}{2\pi n}} \left[2R\Gamma\left(\frac{1}{p} + 1\right) \left(\frac{pe}{n}\right)^{\frac{1}{p}}\right]^n \end{aligned}$$

Avec la fonction Γ définie comme :

$$\Gamma(x) = \int_0^{+\infty} e^{-t} t^{x-1} dt$$

TRAVAIL EN GRANDE DIMENSION : VOIR GRAND

CONCENTRATION DANS UNE HYPERSPHÈRE

On rappelle que :

$$\forall R > 0, \forall n \geq 2, \forall p \geq 1, \quad V_n^p(R) = \frac{\left(2R\Gamma\left(\frac{1}{p} + 1\right)\right)^n}{\Gamma\left(\frac{n}{p} + 1\right)}$$

Exercice 4 (Concentration dans l'hypersphère)

Soit $\varepsilon > 0$. On considère une hypersphère de rayon R . Montrer que :

$$\frac{V_n^p(R - \varepsilon)}{V_n^p(R)} = \left(1 - \frac{\varepsilon}{R}\right)^n$$

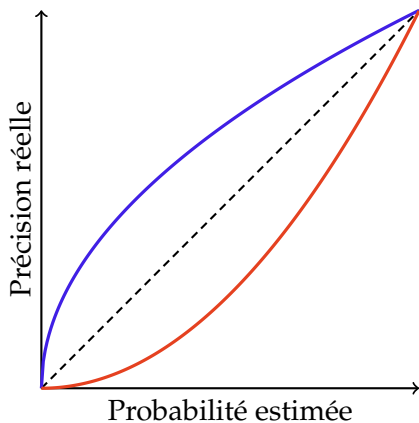
TRAVAIL EN GRANDE DIMENSION : VOIR GRAND

EN RÉSUMÉ

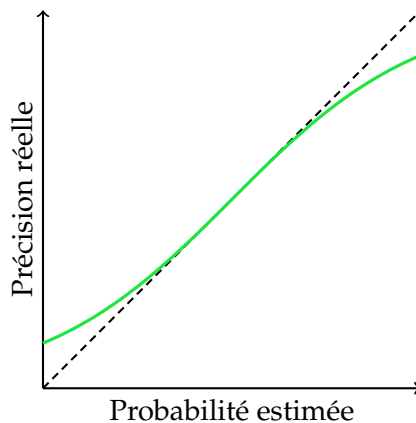
- ▶ Le travail en grande dimension exacerbe l'ensemble des phénomènes présentés jusqu'ici : hyper-déséquilibre ou drift par exemple.
- ▶ Le travail sur la qualité des informations transmises à un modèle ainsi que son architecture est une partie de la réponse à ces problèmes.

CALIBRATION : VOIR JUSTE

PARFAITE CALIBRATION



(a) Prédicteur **trop confiant** et **pas assez confiant**



(b) Prédicteur **non parfaitement calibré**

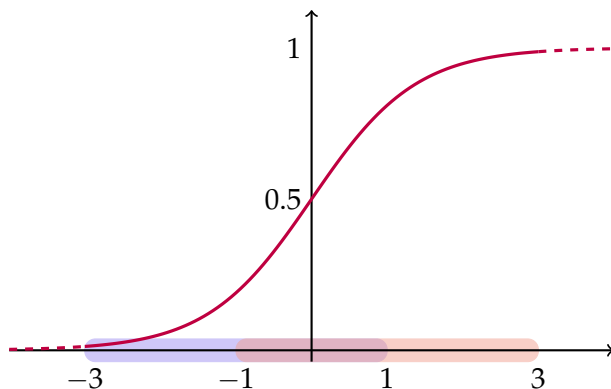
Figure – Comparaison de prédicteurs par rapport à un prédicteur parfaitement calibré (en pointillés)

CALIBRATION : VOIR JUSTE

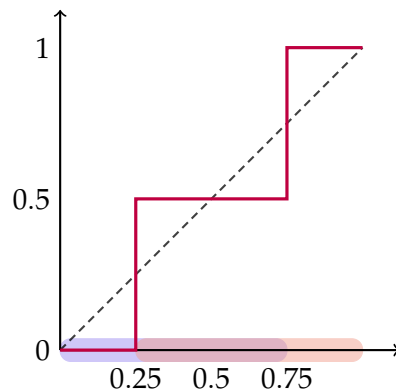
COMMENT INDUIRE LA CALIBRATION ?

Properly regularized logistic regression is well calibrated by default thanks to the use of the log-loss

— Guide utilisateur scikit-learn (2007)



(a) Meilleur modèle obtenu



(b) Diagramme de calibration

Figure – Régression logistique pour la distribution $(x, y) \sim \mathcal{D}$ avec $y = 0$ pour $x \sim \mathcal{U}([-3, 1])$ et $y = 1$ pour $x \sim \mathcal{U}([-1, 3])$

CALIBRATION : VOIR JUSTE

EN RÉSUMÉ

- ▶ La calibration est souhaitable pour tous les modèles destinés à un usage professionnel.
- ▶ Un des moyens de l'induire est de travailler avec des modèles suffisamment puissants pour qu'il *s'auto-calibre* ou d'exploiter un traitement du score après son entraînement.