

Debriefing du TP n°1.

1/ Soit η la variable aléatoire de Bernouilli qui vaut 1 quand la discrimination est réussie et 0 sinon. Soit τ_{gen} la valeur moyenne de η obtenue en calculant la moyenne statistique sur différentes réalisations des vecteurs de la base de généralisation. La moyenne empirique τ_g (si besoin, revoir sa définition dans l'énoncé du TP) est un estimateur de τ_{gen} .

2/ Il est possible d'estimer la barre d'erreur de l'estimation τ_g avec $\hat{\sigma}_{\tau_g} = \sqrt{\frac{\tau_g(1-\tau_g)}{P_{gen}}}$, où P_{gen} est la taille de la base de généralisation.

3/ Sur la courbe $P_{app} \mapsto \tau_g^{(PI)}$, on observe avec un minimum local en $P_{app} = N$. On explique ci-dessous l'origine de ce phénomène, mais avant, on a besoin de quelques rappels mathématiques. Comme expliqué en cours, pour PI la frontière \mathbf{w} est donnée par

$$\mathbf{w} = (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{d}$$

où $(\mathbf{X} \mathbf{X}^T)^{-1}$ est remplacée par $\mathbf{U} \mathbf{\Delta}_{1/\lambda_{1:r}} \mathbf{U}^T$ quand $\mathbf{X} \mathbf{X}^T$ n'est pas inversible. On rappelle que

$$\mathbf{X} \mathbf{X}^T = \mathbf{U} \mathbf{\Delta}_{\lambda_{1:r}} \mathbf{U}^T$$

est une décomposition en valeurs propres. Ainsi $\mathbf{U} \mathbf{U}^T = \mathbf{Id}$ et $\mathbf{\Delta}_{\lambda_{1:r}}$ est une matrice diagonale avec sur les r premiers termes de la diagonale $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$.

Si la matrice est inversible, que peut-on dire de r ?

Remarque. Si, quand la matrice $\mathbf{X} \mathbf{X}^T$ est singulière, on complète la suite de $(\lambda_n)_{n=1,\dots,r}$ avec des zéros, alors on peut aussi s'écrire

$$\mathbf{X} \mathbf{X}^T = \sum_{n=1}^N \lambda_n \mathbf{u}_n \mathbf{u}_n^T \quad (42)$$

où $(\mathbf{u}_n)_{n=1,\dots,N}$ sont des vecteurs propres qui forment une base orthonormée.

Ainsi, quand $(\mathbf{X} \mathbf{X}^T)^{-1}$ est remplacée par $\mathbf{U} \mathbf{\Delta}_{1/\lambda_{1:r}} \mathbf{U}^T$, tout se passe comme si on n'inverse *uniquement* les valeurs propres strictement positives.

A présent, expliquons le phénomène observé. Par définition $\mathbf{X} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(P_{app})})$ donc

$$\mathbf{X} \mathbf{X}^T = \sum_{p=1}^{P_{app}} \mathbf{x}^{(p)} (\mathbf{x}^{(p)})^T$$

que l'on peut comparer à l'équation (42).

A priori, les P_{app} vecteurs $\mathbf{x}^{(p)}$ forment une famille libre, donc

- si $P_{app} \gg N$, alors $\lambda_N > 0$. Dans ce cas, la matrice $\mathbf{X} \mathbf{X}^T$ est inversible. RAS.
- si $P_{app} < N$, alors $r = P_{app}$. La matrice $\mathbf{X} \mathbf{X}^T$ n'est pas inversible et $(\mathbf{X} \mathbf{X}^T)^{-1}$ est remplacée par $\mathbf{U} \mathbf{\Delta}_{1/\lambda_{1:r}} \mathbf{U}^T$. RAS.
- si $P_{app} = N$, alors $r = N = P_{app}$ car, a priori, la famille des \mathbf{x} de la base d'apprentissage est libre. Néanmoins, quand N est grand et qu'on tire de manière aléatoire une famille de taille N , on a une probabilité forte que l'un des vecteurs soit proche d'une combinaison linéaire des autres. Dans ce cas, $\lambda_N \mapsto 0$, sans l'atteindre. Par conséquent, $\lambda_N^{-1} \mapsto +\infty$. Ceci est un *gros* problème, car on obtient alors une estimation *très bruitée* de l'inverse de $\mathbf{X} \mathbf{X}^T$. C'est pour cette raison que l'on a une chute de performance en $P_{app} = N$.

Une solution possible à ce problème est de *régulariser* l'inversion de la matrice $\mathbf{X} \mathbf{X}^T$.

Mathématiquement, cela peut revenir à inverser les valeurs propres de $\mathbf{X} \mathbf{X}^T$ uniquement si celle-ci sont plus grandes qu'un certain seuil, par exemple $\lambda_1/100$.

4/ Avec la ridge approximation (RA), on s'intéresse à

$$\mathbf{w} = (\mathbf{X} \mathbf{X}^T + P_{app} \sigma^2 \mathbf{Id})^{-1} \mathbf{X} \mathbf{d}$$

et donc à l'inverse de $\mathbf{X} \mathbf{X}^T + P\sigma^2 \mathbf{I}_d$. On a

$$\begin{aligned} w_{RA} &\mapsto w_{PI} && \text{quand } \sigma \mapsto 0 \\ w_{RA} &\mapsto w_{Hebb} && \text{quand } \sigma \mapsto \infty \end{aligned}$$

Plus intéressant, on a $\mathbf{X} \mathbf{X}^T = U \Delta_{\lambda_{1:N}} U^T$ et aussi

$$\mathbf{X} \mathbf{X}^T + P_{app} \sigma^2 \mathbf{I}_d = U \Delta_{\lambda_{1:N} + P_{app} \sigma^2} U^T$$

Que peut-on dire de la plus petite valeur propre de $(\mathbf{X} \mathbf{X}^T + P_{app} \sigma^2 \mathbf{I}_d)$?

Ainsi fixer la valeur de σ permet de régulariser l'inversion de la matrice.

Il reste néanmoins à trouver la "bonne" valeur du paramètre σ . Pour cela, on ne doit pas utiliser la base de généralisation, car la base de généralisation est réservée à l'analyse des performances. En revanche, empiriquement, il semble que fixer σ quand le taux d'apprentissage commence à baisser semble une solution intéressante.

Synthèse du TP n°1 : on a comparé une approche empirique (Hebb) avec une approche mathématique (PI). Finalement, la méthode mathématique s'avère intéressante, mais cela a nécessité de faire attention à l'inversion de matrice.