

Évaluation / sélection de modèles

I / Cadre statistique

- notre but est d'approcher f , défini de la façon suivante:

→ Soit P une mesure de probabilité sur \mathbb{R}^p :

$$\min_{f \in \mathcal{F}} \mathbb{E}[L(Z, f)] = \min_{f \in \mathcal{F}} \int L(Z, f) dP(Z)$$

→ L : fonction de perte : "loss function"

→ (f, \mathcal{F}) : une fonction de régression
un classifieur

une densité associée à une certaine loi

- exemple : régression (linéaire)

→ $L(Z, f) = (y - f(x))^2$ avec $Z = (y, x)$

$$\mathcal{F} = \{x \mapsto \beta^T x : \beta \in \mathbb{R}^d\}$$

- on peut considérer d'autres "loss functions"

→ ex : $L(y, f) = |y - f(x)|$

non dérivable, on pose problème pour résoudre problème d'optim, mais moins sensible aux outliers (utile lorsque la queue de distribution est plus importante qu'une gaussienne)

→ d'autres classes de fonctions :

$$\mathcal{F} = \{x \mapsto \sum \alpha_i K(x, x_i) : \alpha \in \mathbb{R}^n\}$$

SVR : support vector regression, version SVM généralisé à la régression

→ $\mathcal{F} = \{x \mapsto f_\theta(x) : f_\theta \text{ la fonction d'un certain réseau de neurones}\}$

→ classification (similaire à la régression)

→ estimation de densité

$$\mathcal{F} = \{f_\theta : \theta \in \Theta\} : \text{une collection de densités}$$

$$L(Z, f) = -\log(f(Z)) \quad (\text{voir distance de Kullback})$$

min_f $-E[\log(f(x))]$ \Leftrightarrow max_f $E[\log(f) - \log(f_0)]$
 avec f_0 la densité de \mathbb{P}

$$\Leftrightarrow \max_f E[\log(f(x))] \Leftrightarrow \max_f \int \log \frac{f(x)}{f_0(x)} f_0(x) dx$$

On a le résultat suivant:

$$\forall f, \int \log(f(x)/f_0(x)) f_0(x) dx \leq 0$$

avec égalité si et seulement si $f = f_0$.

On suppose observer un échantillon $(z_i)_{i=1 \dots n}$ de variables indépendantes et identiquement distribuées sous \mathbb{P} .

→ on applique la "minimisation du risque empirique"

$$\hat{f} \in \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n l(z_i, f)$$

→ def :

$$R(f) = E[l(z, f)]$$

$$\bar{R}(f) = \frac{1}{n} \sum_{i=1}^n l(z_i, f)$$

Théorie de l'estimation: $\hat{f} \approx f$?

→ $E(\hat{f}) \approx f$

→ $\text{Var}(\hat{f}) \leq ?$

→ questions classiques, mais ici on veut poser une question différente: celle de l'évaluation de la qualité de \hat{f} .

II) Correction du biais

Appréhension : la qualité de \hat{f} sera évaluée par :

- $\hat{R}(\hat{f})$? non car on valorise juste le plus grand modèle, celui avec toutes les variables actives.
- $R(\hat{f})$: oui !

problème : $R(\hat{f})$ est une quantité difficile à évaluer.

→ l'estimateur naturel $\hat{R}(\hat{f})$ est biaisé.

→ intuition : pour tout $f \in \mathcal{F}$: $\mathbb{E}[\hat{R}(f)] = R(f)$

pour tout $f \in \mathcal{F}$, $\hat{R}(f)$ est non biaisé

mais $\hat{R}(\hat{f}) = \min_{f \in \mathcal{F}} \hat{R}(f)$, on attend le résultat suivant :
 $\hat{R}(\hat{f}) \leq R(\hat{f})$ la tendance

Il est difficile de vérifier cela dans un cadre général. Vérifions le pour les moindres carrés :

$\hat{\theta} \in \arg \min_{\theta \in \mathbb{R}^p} \|Y - X\theta\|_2^2$ avec $Y \in \mathbb{R}^n$ $X \in \mathbb{R}^{n \times p}$

on a $X^T X \hat{\theta} = X^T Y$

$$\Rightarrow X^T (X\hat{\theta} - Y) = 0$$

$$\Rightarrow X^T (X(\hat{\theta} - \theta^*) + (X\theta^* - Y)) = 0$$

$$\Rightarrow \boxed{X^T (X(\hat{\theta} - \theta^*)) = X^T \varepsilon}$$

Regardons le comportement de $\hat{R}(\hat{f})$ dans ce cadre. = -2 \langle \varepsilon, X(\hat{\theta} - \theta^*) \rangle

$$\hat{R}(\hat{f}) = \frac{\|Y - X\hat{\theta}\|_2^2}{n} = \frac{\|\varepsilon + X(\theta^* - \hat{\theta})\|_2^2}{n} = \frac{\|\varepsilon\|_2^2}{n} + 2 \frac{\langle \varepsilon, X(\theta^* - \hat{\theta}) \rangle}{n} + \frac{\|X(\theta^* - \hat{\theta})\|_2^2}{n}$$

$$\left(\begin{aligned} \text{car } \|X(\theta^* - \hat{\theta})\|_2^2 &= (\theta^* - \hat{\theta})^T X^T X (\theta^* - \hat{\theta}) = -(\theta^* - \hat{\theta})^T X^T \varepsilon = \langle \varepsilon, X(\hat{\theta} - \theta^*) \rangle \\ &= \frac{\|\varepsilon\|_2^2}{n} - \frac{\langle \varepsilon, X(\hat{\theta} - \theta^*) \rangle}{n} \end{aligned} \right)$$

On obtient :

$$= \frac{\|\varepsilon\|_2^2}{n} - \frac{\|X(\theta^* - \hat{\theta})\|_2^2}{n}$$

le risque empirique est biaisé inférieurement

$X, \hat{\theta}$

On cherche à estimer $R(\hat{f}) = \mathbb{E}_Z[\ell(Z, \hat{f})] = \mathbb{E}_Z[(Y - \hat{f}(X))^2]$

$$\begin{aligned} &= \mathbb{E}[(\varepsilon_1 + X_1^T (\theta^* - \hat{\theta}))^2] = \mathbb{E}[\varepsilon_1^2] + 2 \mathbb{E}[\varepsilon_1 X_1^T (\theta^* - \hat{\theta})] + \mathbb{E}[X_1^T (\theta^* - \hat{\theta})^2] \\ &= \mathbb{E}[\varepsilon_1^2] + \mathbb{E}[X_1^T (\theta^* - \hat{\theta})^2] \end{aligned}$$

cadre asymptotique = $n \rightarrow \infty$

$$\frac{\|X(\hat{\theta} - \theta^*)\|_2^2}{n} = (\hat{\theta} - \theta^*)^T \hat{\Sigma}_n (\hat{\theta} - \theta^*) \text{ avec } \hat{\Sigma}_n = \frac{X^T X}{n}$$

$$= \|\hat{\Sigma}_n^{1/2} (\hat{\theta} - \theta^*)\|_2^2$$

Théorème MCO : $\sqrt{n}(\hat{\theta} - \theta^*) \rightarrow \mathcal{N}(0, \sigma^2 \Sigma^{-1})$ $\Sigma = E(X_i X_i^T)$

$$\text{ Slutsky } \Rightarrow \hat{\Sigma}_n^{1/2} \sqrt{n}(\hat{\theta} - \theta^*) \rightarrow \mathcal{N}(0, \sigma^2 I_2)$$

Théorème de continuité :

$$\|\hat{\Sigma}_n^{1/2} \sqrt{n}(\hat{\theta} - \theta^*)\|_2^2 \rightarrow \|\mathcal{N}(0, \sigma^2 I)\|_2^2 = \sigma^2 \chi_d^2$$

Ainsi, il y a bien un biais asymptotique car :

$$E[\chi_d^2] = E\left[\sum_{i=1}^d Z_i^2\right] = d$$

Le C_p de Mallows est défini par le biais $\hat{R}(\hat{f}) - R(\hat{f})$:

$$E[\hat{R}(\hat{f}) - R(\hat{f})] \approx E[\|e\|_2^2 - \|X(\hat{\theta} - \theta^*)\|_2^2 - \|e\|_2^2 - \|X^T(\hat{\theta} - \theta^*)\|_2^2] = \frac{-2\|X^T(\hat{\theta} - \theta^*)\|_2^2}{n}$$

$$C_p = \frac{\|Y - X\hat{\theta}\|_2^2}{n} + \frac{2d}{n} \hat{\sigma}^2$$

$$= \frac{-2\chi_d^2 \sigma^2}{n}$$

Autre critère similaire : AIC, BIC, GIC

- AIC : similaire à Mallows mais dans un contexte plus général de vraisemblance :

$$\hat{\theta}_n \in \arg\max_{\theta} \text{Lib}(\theta) = \arg\max_{\theta} \sum_{i=1}^n \log(f_{\theta}(Z_i))$$

$$\Rightarrow \text{AIC} = 2 \text{Lib}(\hat{\theta}) - 2 \text{length}(\hat{\theta})$$

- BIC, GIC : tout comme le Mallows, BIC et GIC ont un terme correcteur du biais.

$$\text{Lib}(\theta) = \exp\left(-\frac{1}{2}(y - \beta^T x)(y - \beta^T x) / \sigma^2\right)$$

III) Petite parenthèse sur l'utilisation de ces critères/méthodes

• $R(\hat{f}) = \text{comment } \hat{f} \text{ résout notre problème}$) model assessment
→ quantités importantes en elle-même

• Les estimations de $R(f)$ permettent de faire de la sélection en comparant différents modèles. On préfère \hat{f}_1 à \hat{f}_2 si $R(\hat{f}_1) \leq R(\hat{f}_2)$

• Exemple: → sélection forward; on compare chaque modèle correspondant à un jeu de variables sélectionnées.

→ dès qu'on a un hyperparamètre à régler ou définir
par exemple Lasso $\|Y - X\hat{\beta}\|_2^2 + \lambda \|\hat{\beta}\|_1$

support (variable active) qui permet de calculer le Mallows

△ le choix / du critère

de la méthode d'estimation de $R(\hat{f})$ est importante

Esc: $\hat{R}(f) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}_2(x_i))$ \hat{f}_2 = estimateur à plus petits carrés
car $\hat{R}(\hat{f}) = 1$ car ici $\hat{f}(x_i) = x_i$

IV) Validation croisée

• le principe: séparer les données en 2 parties:

→ une pour apprendre le modèle \hat{f}

→ une pour évaluer le risque $R(\hat{f})$

• utilisation: plus large que AIC, BIC, etc... car libre de tout modèle. Ici on peut régler n'importe quel hyperparamètre.

a) le hold-out

Soit $S \subset \{1, \dots, n\}$ l'échantillon de validation
 S^c le complémentaire de S dans $\{1, \dots, n\}$: l'échantillon d'apprentissage

on définit: $\hat{f}_{S^c} \in \arg \min_{f \in \mathcal{F}} \sum_{i \in S^c} \ell(z_i, f)$

l'estimateur "Hold-out" du risque de \hat{f} , $R(\hat{f})$, avec
 $\hat{R}_{Ho} = \frac{1}{|S|} \sum_{i \in S} \ell(z_i, \hat{f}_{S^c})$ $|S|$: cardinalité de S .

1^{er} résultat: $E[\hat{R}_{Ho}] = \frac{1}{|S|} \sum_{i \in S} E[\ell(z_i, \hat{f}_{S^c})]$

2 sources d'aléatoire: $z_i, \hat{f}_{S^c} = T(S^c)$; indépendantes
 car $S \cap S^c = \emptyset$

$$= \frac{1}{|S|} \sum_{i \in S} E[E[\ell(z_i, \hat{f}_{S^c}) | S^c]]$$

$$= \frac{1}{|S|} \sum_{i \in S} E[E_2[\ell(z_i, \hat{f}_{S^c})]]$$

espérance par rapport à z uniquement

$$= \frac{1}{|S|} \sum_{i \in S} E[R(\hat{f}_{S^c})]$$

$$= E[R(\hat{f}_{S^c})]$$

intuition:

$$\hat{R}_{Ho} - R(\hat{f}) = \underbrace{\hat{R}_{Ho} - E[\hat{R}_{Ho}]}_{\text{borne de variance car on peut calculer:}} + \underbrace{E[\hat{R}_{Ho}] - R(\hat{f})}_{E[R(\hat{f}_{S^c}) - R(\hat{f})] (= 0 \text{ si } S^c = \{1, \dots, n\} \text{ full sample})}$$

(décroissante quand $|S^c| \uparrow$)

$$\text{Var}(\hat{R}_{Ho}) \approx O\left(\frac{1}{|S|}\right)$$

incompatibilité entre les 2 directions car
 $|S| + |S^c| = n$

Il en résulte un trade-off entre $|S|$ et $|S^c|$

En pratique, il semble plus raisonnable de choisir $|S|$ petit par rapport à $|S^c|$ car l'apprentissage de \hat{f} est souvent plus compliqué que celui de $R(\hat{f})$ qui est une simple moyenne empirique.

En général, on choisit $|S| \approx 10\%$ de n

6) Algorithmique

• à partir du simple hold-out, on peut définir les algorithmes suivants, qui sont des classiques:

• K-fold: \rightarrow on suppose que $n/K \in \mathbb{N}$

\rightarrow on définit $S_k = \{Z_{\frac{n}{K}+1}, \dots, Z_{(k+1) \cdot \frac{n}{K}}\}$ $k=0 \dots K-1$

S_k : la $k^{\text{ième}}$ fold; chaque fold contient n/K points et $S_i \cap S_j = \emptyset$

aussi: $\bigcup_{k=1}^K S_k = \{Z_1, \dots, Z_n\}$ l'échantillon entier

$$\rightarrow \hat{R}_{K\text{-fold}} = \frac{1}{K} \sum_{k=1}^K R_{H_0}(S_k)$$

\rightarrow étape préliminaire à l'algo: mélange aléatoire uniforme afin d'éliminer une dépendance possible de classement

Cas particulier: leave one out: $\hat{R}_{LOO} = \frac{1}{n} \sum_{i=1}^n R_{H_0}(Z_i)$

problème: temps de calcul très long bien que dans certains cas comme celui de la régression linéaire, des simplifications existent; on a alors pas besoin de calculer chaque fold.

• leave p-out:

$$\rightarrow \hat{R}_{LPO} = \frac{1}{\binom{n}{p}} \sum_{S \in S_p} R_{H_0}(S) \quad \text{avec } S_p: \text{ensemble des sous-échantillons de taille } p.$$



Dans le LPO, on apprend sur $n-p$ points tout comme avec K-fold si $n - (n-1)(K/n) = p$. Il y a plus d'échantillons par rapport à K-fold.

Monte Carlo: les exemples de validation set sont tirés au hasard.

Ex: si $|S|=100$, et $K=10$ alors un validation set sera de taille 10, et tiré

aléatoirement avec Monte

Carlo, comme par exemple:

2, 45, 79, 53, 33, 87, 31, 28, 48, 74

• validation croisée par Monte Carlo

\rightarrow K-fold + tirage aléatoire sur les autres échantillons

\rightarrow si $(R_k)_{k=1 \dots K}$ est une suite de variables aléatoires indépendamment

distribuées (pas nécessairement indépendantes) alors:

$$\text{Var}\left(\frac{1}{K} \sum_{k=1}^K R_k\right) = E\left[\left(\frac{1}{K} \sum_{k=1}^K R_k - E(R_k)\right)^2\right] \leq E\left[\frac{1}{K} \sum_{k=1}^K (R_k - E(R_k))^2\right] = \frac{1}{K} \sum_{k=1}^K \text{Var}(R_k) = \text{Var}(R_k)$$