

Statistique en grande dimension

Lecturer : Dalalyan A., Scribe : Thomas F.-X.

First lecture

1 Introduction

1.1 Statistique classique

Statistique *paramétriques* : Z_1, \dots, Z_n iid, avec une loi commune \mathbb{P}_θ

On fait l'hypothèse $\theta \in \Theta \subset \mathbb{R}^d$

- Connu : Z_1, \dots, Z_n et Θ
- Inconnu : θ ou \mathbb{P}_θ

Hypothèse importante : d est fixe et $n \rightarrow +\infty$

On sait dans ce cas que l'estimateur du MV est asymptotiquement (le plus) efficace (convergent) : $\hat{\theta}_{\text{MV}}$ vérifie quand $n \rightarrow +\infty$:

$$\mathbb{E}_P \left[\|\hat{\theta}_{\text{MV}} - \theta\|^2 \right] = \frac{C}{n} (1 + o(1))$$

On estime θ à une vitesse $\frac{1}{\sqrt{n}}$ (vitesse paramétrique)

Constat. Si $d = d_n$ t.q. $\lim_{n \rightarrow +\infty} d_n = +\infty$, alors toute la théorie paramétrique est inutilisable. De plus, l'estimateur du MV n'est **plus** le meilleur estimateur !

1.2 Statistique non paramétrique

On observe Z_1, \dots, Z_n iid de loi \mathbb{P} , inconnue, telle que $\mathbb{P} \in \{\mathbb{P}_\theta, \theta \in \Theta\}$, mais avec Θ soit de dimension infinie, soit de dimension $d = d_n$ finie mais $\rightarrow +\infty$ avec la taille de l'échantillon.

Exemples:

$$\Theta = \left\{ f : [0, 1] \rightarrow \mathbb{R}, f \text{ Lipschitz de constante } L \right\} \quad (1)$$

$$= \left\{ f : [0, 1] \rightarrow \mathbb{R}, \forall x, y, |f(x) - f(y)| \leq L|x - y| \right\} \quad (2)$$

$$\Theta = \left\{ \theta = (\theta_1, \theta_2, \dots), \sum_{j=1}^{\infty} \theta_j^2 < +\infty \right\} = \ell_2 \quad (3)$$

Démarche générale:

On approche Θ par une suite croissante $\{\Theta_k\}$ de sous-ensembles de Θ telle que Θ_k est de

dimension d_k . En procédant comme si θ appartenait à Θ_k (ce n'est pas nécessairement le cas), on utilise une méthode paramétrique pour définir un estimateur $\hat{\theta}_k$ de θ . Cela nous donne une famille d'estimateurs $\{\hat{\theta}_k\}$.

Question principale. Comment choisir k pour minimiser le risque de $\hat{\theta}_k$?

- Si k est petit, on est face à un phénomène de sous-apprentissage (*underfitting*)
- Inversement, si k est grand, phénomène de sur-apprentissage (*overfitting*)

1.3 Principal models in non-parametric statistics

Density model. We have X_1, \dots, X_n iid with a density f defined on \mathbb{R}^p , and :

$$\mathbb{P}(X_1 \in A) = \int_A f(x) dx$$

The assumptions imposed on f are very weak as opposed to the parametric setting. For instance, a typical assumption in parametric setting is that f is the Gaussian density :

$$f(x) = \frac{\det(\Sigma^{-1})}{(2\pi)^{p/2}} \exp\left[-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right],$$

whereas a common assumption on f in nonparametric framework is : f is smooth, say, twice continuously differentiable with a bounded second derivative.

Regression model. We observe $Z_i = (X_i, Y_i)$, with input X_i , output Y_i and error ε_i :

$$Y_i = f(X_i) + \varepsilon_i.$$

The function f is called the regression function. Here, the goal is to estimate f without assuming any parametric structure on it.

Practical examples.

Marketing.

- Each i represents a consumer
- X_i are the features of the consumer

A typical question is "how do I estimate different relevant groups of consumers". A typical answer is then to use clustering algorithms. We assume that X_1, \dots, X_n are iid with density f . Then, we estimate f in a non-parametric manner by \hat{f} . The clusters are defined as regions around the local maxima of the function \hat{f} .

1.4 Machine Learning

- Essentially the same as non-parametric statistics
- The main focus here is on the algorithms (rather than on the models), their statistical performance and their computational complexity.

2 Main concepts and notations

Observations : Z_1, \dots, Z_n iid, \mathbb{P}

- Non-supervised learning : $Z_i = X_i$
- Supervised learning : $Z_i = (X_i, Y_i)$, where X_i is an example or a feature, and Y_i a label.

Aim. To learn the distribution \mathbb{P} or some properties of it.

Prediction. We assume that a new feature X (from the same prob. distribution as X_1, \dots, X_n) is observed. The aim is to predict the label associated to X .

To measure the quality of a prediction, we need a loss function $\ell(y, \tilde{y})$ (y is the true label, \tilde{y} is the predicted label). In practice, both y and \tilde{y} are random variables, furthermore y and its distribution are unknown, so ℓ is hard to compute!

Risk function. This is the expectation of the loss.

Definition 1 Assume that $Z_i = (X_i, Y_i) \in \mathcal{X} \times \mathcal{Y}$ and $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ is a loss function. A predictor, or prediction algorithm, is any mapping :

$$\hat{g} : (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathcal{Y}^{\mathcal{X}}$$

The risk of the prediction function g is :

$$\mathcal{R}_P[g] = \mathbb{E}_P[\ell(Y, g(X))]$$

The risk of a predictor \hat{g} is $\mathcal{R}_P[\hat{g}]$, which is random since \hat{g} depends on the data.

$$\mathcal{R}_P[\hat{g}] = \int_{\mathcal{X} \times \mathcal{Y}} \ell(y, \hat{g}(x)) d\mathbb{P}(x, y)$$

Examples:

1. Binary classification: $\mathcal{Y} = \{0, 1\}$, with any \mathcal{X}

$$\ell(y, \tilde{y}) = \begin{cases} 0, & \text{if } y = \tilde{y} \\ 1, & \text{otherwise} \end{cases} = \mathbf{1}(y \neq \tilde{y}) = (y - \tilde{y})^2.$$

2. Least-squares regression: $\mathcal{Y} \subset \mathbb{R}$, with any \mathcal{X}

$$\ell(y, \tilde{y}) = (y - \tilde{y})^2.$$

3 Excess risk and Bayes predictor

We have $Z_i = (X_i, Y_i)$

$$\mathcal{R}_P[g] = \int_{\mathcal{X} \times \mathcal{Y}} \ell(y, g(x)) \mathbb{P}(dx, dy)$$

$$\mathbb{P}(dx, dy) = \mathbb{P}_{Y|X}(dy|X=x) \mathbb{P}_X(dx)$$

Definition 2 Given a loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$, the Bayes predictor, or "oracle" is the prediction function minimizing the risk :

$$g^* \in \arg \min_{g \in \mathcal{Y}^{\mathcal{X}}} \mathcal{R}_P[g]$$

Remark 1 In practice, g^* is unavailable, since it depends on \mathbb{P} , which is unknown. The ultimate goal is to do almost as well as the oracle.

A predictor \hat{g}_n will be considered as a good one if :

$$\lim_{n \rightarrow +\infty} \underbrace{\mathcal{R}_P[\hat{g}_n] - \mathcal{R}_P[g^*]}_{\text{excess risk}} = 0$$

Definition 3 We say that the predictor \hat{g}_n is consistent (universally consistent) if $\forall \mathbb{P}$, we have :

$$\lim_{n \rightarrow +\infty} \mathbb{E}_P[\mathcal{R}_P[\hat{g}_n]] - \mathcal{R}_P[g^*] = 0$$

Theorem 1

1. Suppose that $\forall x \in \mathcal{X}$, the infimum of $y \mapsto \mathbb{E}_P[\ell(Y, y) | X = x]$ is reached. Then the function g^* defined by :

$$g^*(x) \in \arg \min_{y \in \mathcal{Y}} \mathbb{E}_P[\ell(Y, y) | X = x]$$

...is a Bayes predictor.

2. In the case of the binary classification, $\mathcal{Y} = \{0, 1\}$ and $\ell(y, \tilde{y}) = \mathbf{1}(y \neq \tilde{y})$,

$$g^*(x) = \mathbf{1}\left(\eta^*(x) > \frac{1}{2}\right) \quad \text{where} \quad \eta^*(x) = \mathbb{P}[Y = 1 | X = x].$$

Furthermore, the excess risk can be computed by

$$\mathcal{R}_P[g] - \mathcal{R}_P[g^*] = \mathbb{E}_P[(g(X) - g^*(X))(1 - 2\eta^*(X))]. \quad (4)$$

3. In the case of the least squares regression,

$$g^*(x) = \eta^*(x) \quad \text{where} \quad \eta^*(x) = \mathbb{E}_P[Y | X = x]$$

Furthermore, for any $\eta : \mathcal{X} \rightarrow \mathcal{Y}$, we have :

$$\mathcal{R}_P[\eta] - \mathcal{R}_P[\eta^*] = \mathbb{E}_P[(\eta(X) - \eta^*(X))^2]$$

Proof

1. Let $g \in \mathcal{Y}^{\mathcal{X}}$ and let :

$$g^*(x) \in \arg \min_{y \in \mathcal{Y}} \mathbb{E}_P[\ell(Y, y) | X = x].$$

We have :

$$\begin{aligned} \mathcal{R}_P[g] &= \mathbb{E}_P[\ell(Y, g(X))] \\ &= \int \mathbb{E}_P[\ell(Y, g(X)) | X = x] \mathbb{P}_X(dx) \\ &\geq \int \mathbb{E}_P[\ell(Y, g^*(x)) | X = x] \mathbb{P}_X(dx) \\ &= \mathcal{R}_P[g^*]. \end{aligned}$$

2. Using the first assertion,

$$\begin{aligned}
g^*(x) &\in \arg \min_{y \in \{0,1\}} \mathbb{E}_P [\mathbf{1}(Y \neq y) | X = x] \\
&= \arg \min_{y \in \{0,1\}} \mathbb{P}(Y \neq y | X = x) \\
&= \arg \max_{y \in \{0,1\}} \mathbb{P}(Y = y | X = x) \\
&= \arg \max_{y \in \{0,1\}} \left\{ \eta^*(x) \mathbf{1}(y = 1) + (1 - \eta^*(x)) \mathbf{1}(y = 0) \right\}.
\end{aligned}$$

Therefore,

$$g^*(x) = \begin{cases} 0, & \text{if } \mathbb{P}(Y = 1 | X = x) \leq \frac{1}{2} \\ 1, & \text{otherwise.} \end{cases}$$

To check (4), it suffices to remark that

$$\begin{aligned}
\mathcal{R}_P[g] &= \mathbb{E}_P[(g(X) - Y)^2] = \mathbb{E}_P[g(X)^2] + \mathbb{E}_P[Y^2] - 2\mathbb{E}_P[Yg(X)] \\
&= \mathbb{E}_P[g(X)] + \mathbb{E}_P[Y] - 2\mathbb{E}_P[\mathbb{E}_P(Yg(X) | X)] \\
&= \mathbb{E}_P[g(X)] + \mathbb{E}_P[Y] - 2\mathbb{E}_P[g(X)\mathbb{E}_P(Y | X)] \\
&= \mathbb{E}_P[g(X)] + \mathbb{E}_P[Y] - 2\mathbb{E}_P[g(X)\eta_P^*(X)] \\
&= \mathbb{E}_P[g(X)(1 - 2\eta_P^*(X)) + \mathbb{E}_P[Y].
\end{aligned}$$

Writing the same identity for g_P^* and making the difference of these two identities, we get the desired result.

3. In view of the first assertion of the theorem, we have:

$$g^*(x) \in \arg \min_{y \in \mathbb{R}} \mathbb{E}_P [(Y - y)^2 | X = x] = \arg \min_{y \in \mathbb{R}} \varphi(y)$$

where $\varphi(y) = \mathbb{E}_P [Y^2 | X = x] - 2y\mathbb{E}_P [Y | X = x] + y^2$ is a second order polynomial. The minimization of such a polynomial is straightforward and leads to:

$$\arg \min_{y \in \mathbb{R}} \varphi(y) = \mathbb{E}_P [Y | X = x].$$

This shows that the Bayes predictor is equal to the regression function $\eta^*(x)$. The risk of this predictor is:

$$\begin{aligned}
\mathcal{R}_P[\eta] &= \mathbb{E}_P [(Y - \eta(X))^2] \\
&= \mathbb{E}_P \left(\mathbb{E}_P [(Y - \eta(X))^2 | X] \right) \\
&= \mathbb{E}_P \left(\mathbb{E}_P [(Y - \eta^*(X))^2 | X] + 2\mathbb{E}_P [(Y - \eta^*(X))(\eta^* - \eta)(X) | X] + (\eta^* - \eta)^2(X) \right) \\
&= \mathcal{R}_P[\eta^*] + 0 + \mathbb{E}_P [(\eta^* - \eta)^2(X)],
\end{aligned}$$

where the cross-product term vanishes since

$$\mathbb{E}_P [(Y - \eta^*(X))(\eta^* - \eta)(X) | X] = (\eta^* - \eta)(X) \mathbb{E}_P [(Y - \eta^*(X)) | X] = 0.$$

This completes the proof of the theorem. □

3.1 Link between Binary Classification & Regression

Plug-in rule

- We start by estimating $\eta^*(x)$ by $\hat{\eta}_n(x)$,
- We define $\hat{g}_n(x) = \mathbf{1}\left(\hat{\eta}_n > \frac{1}{2}\right)$.

Question: How good the plug-in rule \hat{g}_n is ?

Proposition 1 Let $\hat{\eta}$ be an estimator of the regression function η^* , and let $\hat{g}(x) = \mathbf{1}\left(\hat{\eta}(x) > \frac{1}{2}\right)$. Then, we have :

$$\mathcal{R}_{\text{class}}[\hat{g}] - \mathcal{R}_{\text{class}}[g^*] \leq 2\sqrt{\mathcal{R}_{\text{reg}}[\hat{\eta}] - \mathcal{R}_{\text{reg}}[\eta^*]}$$

Proof Let $\eta : \mathcal{X} \rightarrow \mathcal{Y} \subset \mathbb{R}$, and $g(x) = \mathbf{1}\left(\eta(x) > \frac{1}{2}\right)$, and let's compute the excess risk of g . We have,

$$\mathcal{R}_{\text{class}}[g] - \mathcal{R}_{\text{class}}[g^*] = \mathbb{E}_P[(g(X) - g^*(X))(1 - 2\eta^*(X))].$$

Since g and g^* are both indicator functions and, therefore, take only the values 0 and 1, their difference will be nonzero if and only if one of them is equal to 1 and the other one is equal to 0. This leads to

$$\begin{aligned} \mathcal{R}_{\text{class}}[g] - \mathcal{R}_{\text{class}}[g^*] &\leq \mathbb{E}_P[\mathbf{1}(\eta(X) \leq 1/2 < \eta^*(X))|2\eta^*(X) - 1|] \\ &\quad + \mathbb{E}_P[\mathbf{1}(\eta^*(X) \leq 1/2 < \eta(X))|2\eta^*(X) - 1|] \\ &= 2\mathbb{E}_P[\mathbf{1}(1/2 \in [\eta^*(X), \eta(X)])|\eta^*(X) - 1/2|] \end{aligned}$$

If $\eta(X) \leq 1/2$ and $\eta^*(X) > 1/2$, then $|\eta^*(X) - 1/2| \leq |\eta^*(X) - \eta(X)|$, and thus :

$$\begin{aligned} \mathcal{R}_{\text{class}}[g] - \mathcal{R}_{\text{class}}[g^*] &\leq 2\mathbb{E}_P[\mathbf{1}(1/2 \in [\eta(X), \eta^*(X)])|\eta(X) - \eta^*(X)|] \\ &\leq 2\mathbb{E}_P[|\eta(X) - \eta^*(X)|] \\ &\leq 2\sqrt{\mathbb{E}_P[(\eta(X) - \eta^*(X))^2]} = 2\sqrt{\mathcal{R}_{\text{reg}}(\eta) - \mathcal{R}_{\text{reg}}(\eta^*)}. \end{aligned}$$

Since this inequality is true for every deterministic η , we get the desired property. \square