Context McAllester's pioneering work The different PAC-Bayes bounds Main ideas in Chap.1

Application to linear least squares

◆□▶ ◆□▶ ▲□▶ ▲□▶ □ のQ@

PAC-Bayesian bounds

Jean-Yves Audibert^{1,2}

1. Imagine - Université Paris Est, 2. Willow - CNRS/ENS/INRIA

July 2009

Context	McAllester's	pioneering	work
00000	00000		

The different PAC-Bayes bounds Main ideas in Chap.1

Application to linear least squares

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Outline

- Context
- McAllester's pioneering work
- 3 The different PAC-Bayes bounds
 - Seeger's PAC Bayesian bound
 - Catoni's old PAC Bayesian bound
 - Audibert's PAC Bayesian bound
 - Zhang's PAC Bayesian bound
- Main ideas in Chap.1
- 5 Application to linear least squares
 - Framework
 - Variants of known results
 - New results
 - Ridge regression and empirical risk minimization
 - A simple tight risk bound for a sophisticated PAC-Bayes algorithm

 Context
 McAllester's pioneering work
 The different PAC-Bayes bounds
 Main ideas in Chap.1
 Application to linear least squares

 •0000
 •0000
 •00000
 •000000
 •0000000
 •000000000

Supervised learning

Training data = n input-output pairs :

$$Z_1=(X_1,Y_1),\ldots,Z_n=(X_n,Y_n)$$

- A new input X comes.
- Goal: predict the corresponding output Y.
- Probabilistic assumption (batch setting):

$$Z = (X, Y), Z_1, ..., Z_n$$
 i.i.d.

from some unknown distribution P

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

Application to linear least squares

Measuring the quality of prediction

- $\ell(y, y') = \text{loss incurred for predicting } y'$ while the true output is y
- Typical losses are:
 - the least square loss for real outputs

$$\ell(y,y')=(y-y')^2$$

the classification loss for discrete outputs

$$\ell(y,y') = \mathbf{1}_{y \neq y'}$$

- Prediction function: $f: \mathcal{X} \to \mathcal{Y}$
- Risk: $R(f) = \mathbb{E} \ell[Y, f(X)]$

McAllester's pioneering work The different PAC-Bayes bounds Context 00000

Main ideas in Chap.1

Application to linear least squares

Statistical learning theory (SLT)

- Achievable goal for an estimator \hat{f} : predict as well as the best function in a set of prediction functions \mathcal{F} (provided that \mathcal{F} is not too large)
- Central goal of SLT: study $R(\hat{f})$ (whatever \hat{f} is)
- Prominent tool of SLT: probabilistic analysis of the supremum

$$\sup_{f\in\mathcal{F}} \left| R(f) - r(f) \right|$$

with

$$r(f) = \frac{1}{n} \sum_{i=1}^{n} \ell[Y_i, f(X_i)].$$

(日) (日) (日) (日) (日) (日) (日)

Application to linear least squares

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Kullback-Leibler (KL) divergence

$$\mathcal{K}(\rho, \pi) = \begin{cases} \mathbb{E}_{\rho(df)} \log(\frac{\rho}{\pi}(f)) & \text{if } \rho \ll \pi \\ +\infty & \text{otherwise} \end{cases}$$

- 1 If $\rho \ll \pi$, then we have $K(\rho, \pi) = \mathbb{E}_{\pi(df)}\chi(\frac{\rho}{\pi}(f))$ with $\chi: u \mapsto u \log(u) + 1 - u$ convex and nonnegative
- **2** $K(\rho, \pi) > 0$
- If \mathcal{F} is finite and π is the uniform distribution on \mathcal{F} , let $H(\rho) = -\sum_{f \in \mathcal{F}} \rho(f) \log \rho(f)$, then

$$K(\rho,\pi) = \log(|\mathcal{F}|) - H(\rho) \le \log |\mathcal{F}|.$$

Legendre transform of the KL divergence

Let
$$h : \mathcal{F} \to \mathbb{R}$$
 s.t. $\mathbb{E}_{\pi(df)} e^{h(f)} < +\infty$. Define

$$\pi_h(df) = rac{e^{h(f)}}{\mathbb{E}_{\pi(df')}e^{h(f')}} \cdot \pi(df)$$

•
$$\mathcal{K}(\rho, \pi_h) = \mathcal{K}(\rho, \pi) - \mathbb{E}_{\rho(df)}h(f) + \log \mathbb{E}_{\pi(df)}e^{h(f)}$$

• $\sup_{\rho} \{\mathbb{E}_{\rho(df)}h(f) - \mathcal{K}(\rho, \pi)\} = \log \mathbb{E}_{\pi(df)}e^{h(f)}$
• $\operatorname{argmax}_{\rho}\{\mathbb{E}_{\rho(df)}h(f) - \mathcal{K}(\rho, \pi)\} = \pi_h$
• $\lambda \mapsto \mathcal{K}(\pi_{\lambda h}, \pi)$ is nondecreasing on $[0, +\infty)$.

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

Main ideas in Chap.1

Application to linear least squares

PAC-Bayesian analysis

 Study E_{ρ̂(df)}R(f) for any distribution ρ̂ on F depending on the training data

 \longrightarrow similar to the study of $R(\hat{f})$ (whatever \hat{f} is)

 Uses a (prior) distribution to evaluate the complexity of the data-dependent (or posterior) distribution

 \longrightarrow different from VC bounds where the complexity is a global quantity characterizing the model ${\cal F}$

- The bound holds for any prior and posterior
 - \longrightarrow different from the usual Bayesian approach

Application to linear least squares

McAllester's bound (1998,1999)

We assume $0 \le \ell(y, y') \le 1$ for any y, y'.

For any distribution π on \mathcal{F} , with probability at least $1 - \epsilon$, for any distribution ρ on \mathcal{F}

$$\left|\mathbb{E}_{
ho(df)}R(f)-\mathbb{E}_{
ho(df)}r(f)\right|\leq \sqrt{rac{K(
ho,\pi)+\log(4n\epsilon^{-1})}{2n-1}}$$

Equivalently (measurability problems set aside), for any data-dependent (posterior) distribution $\hat{\rho}$, with probability at least $1 - \epsilon$.

$$\left|\mathbb{E}_{\hat{
ho}(df)}R(f)-\mathbb{E}_{\hat{
ho}(df)}r(f)\right|\leq\sqrt{rac{K(\hat{
ho},\pi)+\log(4n\epsilon^{-1})}{2n-1}}$$

Application to linear least squares

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

Seeger's proof (slightly revisited)

The PAC lemma

Let V be a real-valued random variable s.t. $\mathbb{E}e^{V} < 1$, then with probability at least $1 - \epsilon$, we have

$$V \le \log(\epsilon^{-1}).$$

McAllester's bound:

$$V = \sup_{\rho} \left\{ (2n-1) \big[\mathbb{E}_{\rho(df)} R(f) - \mathbb{E}_{\rho(df)} r(f) \big]^2 - K(\rho, \pi) - \log(4n) \right\} \leq \log(\epsilon^{-1}).$$

First step: Jensen's ineg. + Legendre transform of KL

$$V \le \sup_{\rho} \left\{ (2n-1)\mathbb{E}_{\rho(df)} [R(f) - r(f)]^2 - K(\rho, \pi) - \log(4n) \right\}$$

= $-\log(4n) + \log \mathbb{E}_{\pi(df)} e^{(2n-1)[R(f) - r(f)]^2}$

Context McAllester's pioneering work The different PAC-Bayes bounds Main ideas in Chap.1 00000 00000

Application to linear least squares

Seeger's proof (second step)

$$\begin{split} \mathbb{E}e^{V} &\leq \frac{1}{4n} \mathbb{E}\mathbb{E}_{\pi(df)} e^{(2n-1)[R(f)-r(f)]^{2}} \\ &= \frac{1}{4n} \mathbb{E}_{\pi(df)} \Big(1 + \mathbb{E} \Big\{ e^{(2n-1)[R(f)-r(f)]^{2}} - 1 \Big\} \Big) \\ &= \frac{1}{4n} \mathbb{E}_{\pi(df)} \Big(1 + \int_{0}^{+\infty} \mathbb{P}(e^{(2n-1)[R(f)-r(f)]^{2}} - 1 > t) dt \Big) \\ &= \frac{1}{4n} \mathbb{E}_{\pi(df)} \Big(1 + \int_{0}^{+\infty} \mathbb{P}(|R(f) - r(f)| > \sqrt{\frac{\log(t+1)}{2n-1}}) dt \Big) \\ &\leq \frac{1}{4n} \mathbb{E}_{\pi(df)} \Big(1 + \int_{0}^{+\infty} 2e^{-2n\frac{\log(t+1)}{2n-1}} dt \Big) \\ &= \frac{1}{4n} \mathbb{E}_{\pi(df)} \Big(1 + 2\int_{1}^{+\infty} (t+1)^{-\frac{2n}{2n-1}} dt \Big) \\ &= \frac{4n-1}{4n} \leq 1 \end{split}$$

Main ideas in Chap.1

Minimizing McAllester's bound and Gibbs estimator

Let $B(\rho) = \mathbb{E}_{\rho(df)} r(f) + \sqrt{\frac{K(\rho,\pi) + \log(4ne^{-1})}{2n-1}}$. McAllester's bound implies: for any distribution ρ

 $\mathbb{E}_{\rho(df)}R(f) \leq B(\rho).$

Theorem

There exists
$$\hat{\lambda} \in [\lambda_1, \lambda_2]$$
 s.t. $B(\pi_{-\hat{\lambda}r}) = \min_{\rho} B(\rho)$ with $\lambda_1 = \sqrt{4(2n-1)\log(4n\epsilon^{-1})}$ and $\lambda_2 = 2\lambda_1 + 4(2n-1)$. Besides, we have

$$\hat{\lambda} = \sqrt{4(2n-1)[K(\pi_{-\hat{\lambda}r},\pi) + \log(4n\epsilon^{-1})]}$$

$$\hat{\lambda} \in \operatorname*{argmin}_{\lambda>0} \left\{ -\frac{1}{\lambda} \log \mathbb{E}_{\pi(df)} e^{-\lambda r(f)} + \frac{\lambda}{4(2n-1)} + \frac{\log(4n\epsilon^{-1})}{\lambda} \right\}$$

Context McAllester's pioneering work The different PAC-Bayes bounds Main ideas in Chap.1

000000

Application to linear least squares

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Seeger's PAC Bayesian bound

Seeger's bound for classification (2002)

slightly revisited

•
$$\mathcal{K}(p||q) = \mathcal{K}(Be(p), Be(q)) = p \log\left(\frac{p}{q}\right) + (1-p) \log\left(\frac{(1-p)}{1-q}\right)$$

Theorem

With probability at least $1 - \epsilon$, for any distribution ρ on \mathcal{F} ,

$$K(\mathbb{E}_{\rho(df)}r(f)||\mathbb{E}_{\rho(df)}R(f)) \leq \frac{K(\rho,\pi) + \log(2\sqrt{n}\epsilon^{-1})}{n}$$

Context	McAllester's pioneering work	The different PAC-Bayes bounds	Main ideas in Chap.1	Application to linear least squares
00000	00000	000000	0000000	0000000000

Seeger's PAC Bayesian bound

This time, it suffices to prove $V = \sup_{\rho} \left\{ n \mathcal{K}(\mathbb{E}_{\rho(df)} r(f)) || \mathbb{E}_{\rho(df)} \mathcal{R}(f)) - \mathcal{K}(\rho, \pi) - \log(2\sqrt{n}) \right\} \leq \log(\epsilon^{-1}).$

We have

$$\mathbb{E}\boldsymbol{e}^{V} \leq \mathbb{E}\boldsymbol{e}^{\sup_{\rho} \left\{ n\mathbb{E}_{\rho(df)}K(r(f)||R(f)) - K(\rho,\pi) - \log(2\sqrt{n}) \right\}} = \frac{1}{2\sqrt{n}}\mathbb{E}\mathbb{E}_{\pi(df)}\boldsymbol{e}^{nK(r(f)),R(f))}$$

$$= \frac{1}{2\sqrt{n}}\mathbb{E}_{\pi(df)}\sum_{k=0}^{n}\mathbb{P}(nr(f) = k)\left(\frac{k}{nR(f)}\right)^{k}\left(\frac{n-k}{n[1-R(f)]}\right)^{n-k}$$

$$= \frac{1}{2\sqrt{n}}\mathbb{E}_{\pi(df)}\sum_{k=0}^{n}\binom{n}{k}\left(\frac{k}{n}\right)^{k}\left(\frac{n-k}{n}\right)^{n-k}$$

$$\leq 1,$$

where the last inequality is obtained from computations using Stirling's approximation.

McAllester's pioneering work The different PAC-Bayes bounds

000000

Main ideas in Chap.1

Application to linear least squares

Seeger's PAC Bayesian bound

McAllester's bound vs Seeger's bound

•
$$\left|\mathbb{E}_{\rho(df)}R(f) - \mathbb{E}_{\rho(df)}r(f)\right| \leq \sqrt{\frac{K(\rho,\pi) + \log(4n\epsilon^{-1})}{2n-1}}$$
 (1)

- $\mathcal{K}(\mathbb{E}_{\rho(df)}r(f)||\mathbb{E}_{\rho(df)}R(f)) \leq \frac{\mathcal{K}(\rho,\pi) + \log(2\sqrt{n}\epsilon^{-1})}{n}$ (2)
- (2) \Rightarrow (1) up to constant since from Pinsker's inequality:

$$\left|\mathbb{E}_{
ho(df)}R(f)-\mathbb{E}_{
ho(df)}r(f)
ight|\leq \sqrt{\mathcal{K}ig(\mathbb{E}_{
ho(df)}r(f)||\mathbb{E}_{
ho(df)}R(f)ig)}.$$

• (2) \gg (1) when $\mathbb{E}_{\rho(df)}R(f)$ is close to 0 since (2) implies

$$\left|\mathbb{E}_{\rho(df)}R(f)-\mathbb{E}_{\rho(df)}r(f)\right| \leq \sqrt{\frac{2\mathbb{E}_{\rho(df)}r(f)[1-\mathbb{E}_{\rho(df)}r(f)]\mathcal{K}}{n}} + \frac{4\mathcal{K}}{3n}$$

with

$$\mathcal{K} = \mathcal{K}(\rho, \pi) + \log(2\sqrt{n}\epsilon^{-1}).$$

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

Context McAllester's pioneering work The different PAC-Bayes bounds Main ideas in Chap.1

0000000

Application to linear least squares

Catoni's old PAC Bayesian bound

Catoni's old bound for classification (2002)

• Let
$$\Psi(\lambda) = \frac{e^t - 1 - t}{t^2} \xrightarrow[t \to 0]{} \frac{1}{2}$$
.

Theorem

For $\lambda > 0$, with proba. at least $1 - \epsilon$, for any distribution ρ on \mathcal{F} ,

$$\mathbb{E}_{\rho(df)} R(f) \leq \frac{\mathbb{E}_{\rho(df)} r(f)}{1 - \frac{\lambda}{n} \Psi(\frac{\lambda}{n})} + \frac{K(\rho, \pi) + \log(\epsilon^{-1})}{\lambda [1 - \frac{\lambda}{n} \Psi(\frac{\lambda}{n})]}$$

Since typical values of λ are in $[C\sqrt{n}; Cn]$, we roughly have

$$\mathbb{E}_{\rho(df)} \mathcal{R}(f) \lesssim \mathbb{E}_{\rho(df)} r(f) + \frac{\lambda}{2n} \mathbb{E}_{\rho(df)} r(f) + \frac{\mathcal{K}(\rho, \pi) + \log(\epsilon^{-1})}{\lambda}$$

$$\approx \sum_{\text{choice of } \lambda} \mathbb{E}_{\rho(df)} r(f) + \sqrt{2\mathbb{E}_{\rho(df)} r(f) \frac{\mathcal{K}(\rho, \pi) + \log(\epsilon^{-1})}{n}}$$

McAllester's pioneering work The different PAC-Bayes bounds

0000000

Main ideas in Chap.1

Application to linear least squares

Audibert's PAC Bayesian bound

Audibert's bound (2004)

• Let
$$\Psi(\lambda) = \frac{e^t - 1 - t}{t^2} \xrightarrow[t \to 0]{} \frac{1}{2}$$
.

Theorem

For $\lambda > 0$, with proba. at least $1 - \epsilon$, for any distribution ρ on \mathcal{F} ,

$$\mathbb{E}_{
ho(df)} R(f) \leq \mathbb{E}_{
ho(df)} r(f) + rac{\lambda}{n} \Psi\left(rac{\lambda}{n}
ight) \mathbb{E}_{
ho(df)} \operatorname{Var}_{Z} \ell(Y, f(X)) + rac{K(
ho, \pi) + \log(\epsilon^{-1})}{\lambda}.$$

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

Context McAllester's pioneering work The different PAC-Bayes bounds Main ideas in Chap.1

0000000

Application to linear least squares

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Zhang's PAC Bayesian bound

Zhang's bound (2005)

Theorem

For $\lambda > 0$, with proba. at least $1 - \epsilon$, for any distribution ρ on \mathcal{F} ,

$$-\frac{n}{\lambda}\mathbb{E}_{\rho(df)}\log\mathbb{E}_{Z}\boldsymbol{e}^{-\frac{\lambda}{n}\ell(Y,f(X))}\leq\mathbb{E}_{\rho(df)}\boldsymbol{r}(f)+\frac{K(\rho,\pi)+\log(\epsilon^{-1})}{\lambda}.$$

Since we have

$$-\frac{1}{t}\log \mathbb{E}_{Z}e^{-t\ell(Y,f(X))} = R(f) - \frac{t}{2}\operatorname{Var}_{Z}\ell(Y,f(X)) + O(t^{2}),$$

we have

I.h.s.
$$\approx \mathbb{E}_{\rho(df)} R(f) - \frac{\lambda}{2n} \mathbb{E}_{\rho(df)} \operatorname{Var}_{Z} \ell(Y, f(X))$$

Context McAllester's pioneering work The different PAC-Bayes bounds Main ideas in Chap.1

000000

Application to linear least squares

(ロ) (同) (三) (三) (三) (三) (○) (○)

Zhang's PAC Bayesian bound

Comparison of the bounds in classification

Shang and Audibert:

$$\mathbb{E}_{\rho(df)}R(f) \lessapprox \mathbb{E}_{\rho(df)}r(f) + \sqrt{2\mathbb{E}_{\rho(df)}(R(f)[1-R(f)])\frac{K(\rho,\pi) + \log(\epsilon^{-1})}{n}}$$

Catoni:

$$\mathbb{E}_{\rho(df)}R(f) \lesssim \mathbb{E}_{\rho(df)}r(f) + \sqrt{2\mathbb{E}_{\rho(df)}R(f)\frac{K(\rho,\pi) + \log(\epsilon^{-1})}{n}}$$

Seeger:

$$\mathbb{E}_{\rho(df)}R(f) \leq \mathbb{E}_{\rho(df)}r(f) + \sqrt{\frac{2\mathbb{E}_{\rho(df)}R(f)[1-\mathbb{E}_{\rho(df)}R(f)]\mathcal{K}}{n}} + \frac{2\mathcal{K}}{3n}$$

with $\mathcal{K} = \mathcal{K}(\rho, \pi) + \log(2\sqrt{n}\epsilon^{-1})$. Besides, we have $\mathbb{E}_{\rho(df)}R(f)[1-\mathbb{E}_{\rho(df)}R(f)] \geq \mathbb{E}_{\rho(df)}R(f)[1-R(f)]$

\Rightarrow similar PAC-Bayes bounds

000000

Application to linear least squares

(ロ) (同) (三) (三) (三) (三) (○) (○)

Explicit Laplace transform in classification

Instead of using

$$\log \mathbb{E} e^{-\frac{\lambda}{n}\ell(Y,f(X))} \leq -\frac{\lambda}{n}R(f) + \frac{\lambda^2}{n^2}\Psi\left(\frac{\lambda}{n}\right)R(f),$$

use

$$\log \mathbb{E} e^{-\frac{\lambda}{n}\ell(Y,f(X))} = \log \left(1 - R(f)(1 - e^{-\frac{\lambda}{n}})\right)$$
$$= -\frac{\lambda}{n} \Phi_{\frac{\lambda}{n}}(R(f)).$$

with

$$\Phi_a(p) = -a^{-1}\log[1 - (1 - e^{-a})p] = p - \frac{a}{2}p(1 - p) + O(a^2)$$

 Zhang's bound can be used to obtain exactly the same basic bound as Theorem 1.2.6.

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

Concentration of the risk w.r.t. the posterior distribution

 All PAC-Bayes bounds can be stated as: for any posterior distribution $\hat{\rho}$, with probability at least $1 - \epsilon$ w.r.t. to the joint probability $P^{\otimes n}\hat{\rho}$ of the training set and the randomized prediction function $\hat{f} \sim \hat{\rho}$,

$$R(\hat{f}) \leq r(\hat{f}) + ext{ terms with } \log\left(rac{\hat{
ho}}{\pi}(\hat{f})
ight)$$
 instead of $\mathcal{K}(\hat{
ho},\pi)$

For instance, Seeger's bound becomes:

$$K(r(\hat{f})||R(\hat{f})) \leq rac{\log(rac{\hat{
ho}}{\pi}(\hat{f})) + \log(2\sqrt{n}\epsilon^{-1})}{n}$$

Main ideas in Chap.1 000000

Application to linear least squares

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Unbiased empirical bounds

- Known problem of PAC bounds: pessimistic constants
- Proposed solution: find an empirical quantity $B(\rho)$ s.t.

$$\mathbb{E}[\mathbb{E}_{\rho(df)}R(f)] \leq \mathbb{E}[B(\rho)],$$

and choose the estimator or the parameters by minimizing $B(\rho)$.

Relative bounds

- Main idea: the difference of empirical risks of two close prediction functions has much smaller variations around its mean than the empirical risk of one of these functions.
- Typically, we start with

$$\begin{split} \mathbb{E}_{\rho_1(df)} R(f) &- \mathbb{E}_{\rho_2(df)} R(f) \leq \mathbb{E}_{\rho_1(df)} r(f) - \mathbb{E}_{\rho_2(df)} r(f) \\ &+ \frac{\lambda}{n} \Psi\left(\frac{2\lambda}{n}\right) \mathbb{E}_{\rho_1(df)} \mathbb{E}_{\rho_2(df)} \mathsf{Var}_Z \big[\ell(Y, f_1(X)) - \ell(Y, f_2(X)) \big] \\ &+ \frac{K(\rho_1, \pi_1) + K(\rho_2, \pi_2) + \log(\epsilon^{-1})}{\lambda}. \end{split}$$

instead of

$$\mathbb{E}_{\rho(df)}R(f) \leq \mathbb{E}_{\rho(df)}r(f) + \frac{\lambda}{n}\Psi\left(\frac{\lambda}{n}\right)\mathbb{E}_{\rho(df)}\operatorname{Var}_{Z}\ell(Y, f(X)) + \frac{K(\rho, \pi) + \log(\epsilon^{-1})}{\lambda}.$$

 Context
 McAllester's pioneering work
 The different PAC-Bayes bounds
 Main ideas in Chap.1

 00000
 0000000
 0000000
 0000000

Application to linear least squares

Fast rates under margin assumptions

- in classification
 - Mammen and Tsybakov's assumption: for a reference prediction function *t̃* ∈ *G*, for any *f* ∈ *G*,

$$\mathbb{P}[f(X) \neq \tilde{f}(X)] \leq C[R(f) - R(\tilde{f})]^{1/\kappa}$$

• Catoni's margin functions:

$$\varphi(t) = \sup_{f \in \mathcal{F}} \left\{ \mathbb{E} \left| \mathbf{1}_{Y \neq f(X)} - \mathbf{1}_{Y \neq \tilde{f}(X)} \right| - t[R(f) - R(\tilde{f})] \right\}$$

$$\bar{\varphi}(t) = \sup_{f \in \mathcal{F}} \left\{ \bar{\mathbb{E}} \left| \mathbf{1}_{Y \neq f(X)} - \mathbf{1}_{Y \neq \tilde{f}(X)} \right| - t[r(f) - r(\tilde{f})] \right\}$$

in least squares regression, under reasonable assumptions,

$$\operatorname{Var}_{Z}\left[\ell(Y, f(X)) - \ell(Y, \tilde{f}(X))\right] \leq c[R(f) - R(\tilde{f})]$$

Algorithm design by successive improvement

A better variance control in classification

• Classification : $|\mathcal{Y}| < +\infty$ and $L(y, y') \triangleq \mathbb{1}_{y \neq y'}$

Transductive setting : we are given the training set Z_1^N and N points to classify X_{N+1}, \ldots, X_{2N} . **Target** : predict unknown labels Y_{N+1}, \ldots, Y_{2N}

$$\begin{cases} \bar{\mathbb{P}} & \triangleq \frac{1}{N} \sum_{i=1}^{N} \delta_{(X_{i},Y_{i})} \\ \bar{\mathbb{P}}' & \triangleq \frac{1}{N} \sum_{i=N+1}^{2N} \delta_{(X_{i},Y_{i})} \\ \bar{\mathbb{P}} & \triangleq \frac{1}{2N} \sum_{i=1}^{2N} \delta_{(X_{i},Y_{i})} \\ r(f) & \triangleq \frac{1}{N} \sum_{i=1}^{N} \mathbb{1}_{Y_{i} \neq f(X_{i})} = \bar{\mathbb{P}}[Y \neq f(X)] \\ r'(f) & \triangleq \frac{1}{N} \sum_{i=N+1}^{2N} \mathbb{1}_{Y_{i} \neq f(X_{i})} = \bar{\mathbb{P}}'[Y \neq f(X)] \\ \bar{\mathbb{P}}_{f_{1},f_{2}} & \triangleq \bar{\mathbb{P}}[f_{1}(X) \neq f_{2}(X)] \end{cases}$$

Another way of controlling the variance term

Reminder

•
$$\rho_2 r' - \rho_1 r' + \rho_1 r - \rho_2 r \leq \frac{2\lambda}{N} (\rho_1 \otimes \rho_2) \overline{\mathbb{P}}_{\cdot,\cdot} + \frac{\mathcal{K}_{1,2}}{\lambda}$$

• Target : use the bounds to design efficient estimators **Basic approach** : consider $(\rho_2, \pi_2, \rho_1, \pi_1) = (\rho, \pi, \delta_{\tilde{f}}, \delta_{\tilde{f}}).$ $\rightsquigarrow \rho r' - r'(\tilde{f}) \leq \rho r - r(\tilde{f}) + \frac{2\lambda}{N}\rho \bar{\mathbb{P}}_{\cdot,\tilde{f}} + \frac{K(\rho,\pi) + \log(\epsilon^{-1})}{\lambda}$

Main problem : control the variance term

Non localized estimator

Theorem. Let
$$L \triangleq \log \left[\log(eN)\epsilon^{-1} \right]$$
 and
 $S(\rho', \rho'') \triangleq \min_{\lambda \in [\sqrt{N};N]} \left\{ \frac{2\lambda}{N} \left(\rho' \otimes \rho'' \right) \overline{\mathbb{P}}_{\cdot,\cdot} + \sqrt{e} \frac{K(\rho',\pi) + K(\rho'',\pi) + L}{\lambda} \right\}.$
With $\mathbb{P}^{\otimes N}$ -proba at least $1 - \epsilon, \forall \rho', \rho'' \in \mathcal{M}^1_+(\mathcal{F}),$
 $\rho''r' - \rho'r' \leq \rho''r - \rho'r + S(\rho', \rho'')$

Algorithm. Let $\rho_0 = \pi$. For any $k \ge 1$, define ρ_k as the distribution with the smallest complexity $K(\rho_k, \pi)$ such that

 $\rho_k r - \rho_{k-1} r + S(\rho_{k-1}, \rho_k) \leq 0$. Classify using a function drawn according to the last posterior distribution ρ_K .

Non localized estimator

Theorem. Let

$$\mathbb{G}(\lambda) \triangleq -\frac{1}{\lambda} \log \pi \exp\left(-\lambda r'\right) + \frac{1}{2\lambda} \log \pi_{-\lambda r'} \exp\left(\frac{72\sqrt{e\lambda^2}}{N} \pi_{-\lambda r'} \bar{\mathbb{P}}_{\cdot,\cdot}\right) + \frac{L}{2\lambda}$$

With $\mathbb{P}^{\otimes 2N}$ -probability at least $1 - \epsilon$, for any $k \in \{1, \ldots, K\}$,

•
$$\rho_k r - \rho_{k-1} r + S(\rho_k, \rho_{k-1}) = 0, \ \rho_k r < \rho_{k-1} r \text{ and } \rho_k r' \le \rho_{k-1} r',$$

• $K(\rho_k, \pi) \ge K(\rho_{k-1}, \pi),$

• $\rho_K r' \leq \min_{\substack{\sqrt{N} \\ \frac{\sqrt{N}}{6\sqrt{e}} \leq \lambda \leq \frac{N}{6\sqrt{e}}}} \mathbb{G}(\lambda).$

Tsybakov's type assumptions:

- there exists C' > 0 and 0 < q < 1 such that the covering entropy of the model \mathcal{F} for the distance $\mathbb{P}_{\cdot,\cdot}$ satisfies for any u > 0, $H(u, \mathcal{F}, \mathbb{P}_{\cdot,\cdot}) \leq C' u^{-q}$,
- there exist c'', C'' > 0 and $\kappa \ge 1$ such that for any function $f \in \mathcal{F}$,

$$c'' \big[R(f) - R(\tilde{f}) \big]^{\frac{1}{\kappa}} \le \mathbb{P}_{f,\tilde{f}} \le C'' \big[R(f) - R(\tilde{f}) \big]^{\frac{1}{\kappa}},$$

 \Rightarrow with $\mathbb{P}^{\otimes 2N}$ -probability at least $1 - \epsilon$,

$$\mathbb{G}(\lambda) \le r'(\tilde{f}) + \log(e\epsilon^{-1}) \mathbf{O}\left(N^{-\frac{\kappa}{2\kappa-1+q}}\right)$$

provided that $\lambda = N^{\frac{\kappa}{2\kappa-1+q}} (\in [\sqrt{N}; N])$ and π is appropriately chosen.

Context McAllester's pioneering work The different PAC-Bayes bounds

Main ideas in Chap.1 0000000

Application to linear least squares

PAC-Bayesian localization

• For a given $\hat{\rho}$, the prior minimizing the expected value of the bound for $\hat{\rho}$ is

$$\pi = \operatorname{argmin}_{\pi'} \mathbb{E} \mathcal{K}(\hat{
ho}, \pi') = \mathbb{E}[\hat{
ho}]$$

since $\mathbb{E}K(\hat{\rho},\pi) = \mathbb{E}K(\hat{\rho},\mathbb{E}[\hat{\rho}]) + K(\mathbb{E}[\hat{\rho}],\pi).$

- Problem: $\mathbb{E}[\hat{\rho}]$ is not observable
- First solution (Catoni, 2003): apply basic PAC bound to $\pi_{-\beta R}$, expand $K(\hat{\rho}, \pi_{-\beta R})$ and develop additional empirical bounds to control the non observable terms
 - Zhang (2005) uses $\pi_{\alpha \log \mathbb{E}_{\mathcal{Z}} e^{-\lambda \ell(Y, f(X))}}$.
 - Ambroladze, P.-H. and S.-T. (2006) localizes by cutting the training set into two parts
 - Catoni (2007) uses π_{-βΦ_β}[R(f)].
 - Alquier (2007,2008) also uses $\pi_{-\beta R}$ but for general unbounded losses (regression, density estimation)

Main ideas in Chap.1 000000

Application to linear least squares

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

Properties of PAC-Bayesian localization

- Advantages
 - allow to replace $K(\rho, \pi)$ with $K(\rho, \pi_{-\lambda r})$
 - gain of logarithmic factor in parametric convergence rates
- Disadvantages = increase of the constant factors
- Open question = useful to build linear classifiers ? (Herbrich, Graepel, 2001; Langford, Shawe-Taylor, 2002; Germain, Lacasse, Laviolette, Marchand, 2009)

A better variance control in classification

• Classification : $|\mathcal{Y}| < +\infty$ and $L(y, y') \triangleq \mathbb{1}_{y \neq y'}$

Transductive setting : we are given the training set Z_1^N and N points to classify X_{N+1}, \ldots, X_{2N} . **Target** : predict unknown labels Y_{N+1}, \ldots, Y_{2N}

$$\begin{cases} \bar{\mathbb{P}} & \triangleq \frac{1}{N} \sum_{i=1}^{N} \delta_{(X_{i},Y_{i})} \\ \bar{\mathbb{P}}' & \triangleq \frac{1}{N} \sum_{i=N+1}^{2N} \delta_{(X_{i},Y_{i})} \\ \bar{\mathbb{P}} & \triangleq \frac{1}{2N} \sum_{i=1}^{2N} \delta_{(X_{i},Y_{i})} \\ r(f) & \triangleq \frac{1}{N} \sum_{i=1}^{N} \mathbb{1}_{Y_{i} \neq f(X_{i})} = \bar{\mathbb{P}}[Y \neq f(X)] \\ r'(f) & \triangleq \frac{1}{N} \sum_{i=N+1}^{2N} \mathbb{1}_{Y_{i} \neq f(X_{i})} = \bar{\mathbb{P}}'[Y \neq f(X)] \\ \bar{\mathbb{P}}_{f_{1},f_{2}} & \triangleq \bar{\mathbb{P}}[f_{1}(X) \neq f_{2}(X)] \end{cases}$$

Relative PAC-Bayesian bounds

Definitions. • A function Q on \mathbb{Z}^{2N} is said to be exchangeable iff for any permutation σ , $Q_{Z_{\sigma(1)},\dots,Z_{\sigma(2N)}} = Q_{Z_1,\dots,Z_{2N}}$. • $\pi_h \triangleq \frac{\exp(h)}{\pi \exp(h)} \cdot \pi$

Theorem. Let π_1 and π_2 be exchangeable prior distributions. Define $\mathcal{K}_{1,2} \triangleq K(\rho_1, \pi_1) + K(\rho_2, \pi_2) + \log(\epsilon^{-1})$. For any $\epsilon > 0, \lambda > 0$, with $\mathbb{P}^{\otimes 2N}$ -probability at least $1 - \epsilon$, for any $\rho_1, \rho_2 \in \mathcal{M}^1_+(\mathcal{F})$,

$$\rho_2 r' - \rho_1 r' + \rho_1 r - \rho_2 r \le \frac{2\lambda}{N} (\rho_1 \otimes \rho_2) \overline{\mathbb{P}}_{\cdot,\cdot} + \frac{\mathcal{K}_{1,2}}{\lambda}$$

Theorem. For any $\xi \in]0; 1[$ and $\lambda, \lambda_1, \lambda_2 > 0$, define $\mathcal{K}_{1,2}^{\text{loc}} \triangleq K(\rho_1, (\pi_1)_{-\lambda_1 r}) + K(\rho_2, (\pi_2)_{-\lambda_2 r}) + \log(\pi_1)_{-\lambda_1 r} \exp\left(\frac{\lambda_1^2}{2\xi N}\rho_1 \overline{\mathbb{P}}_{\cdot,\cdot}\right) + \log(\pi_2)_{-\lambda_2 r} \exp\left(\frac{\lambda_2^2}{2\xi N}\rho_2 \overline{\mathbb{P}}_{\cdot,\cdot}\right) + (1+\xi)\log(\epsilon^{-1}).$

With $\mathbb{P}^{\otimes 2N}$ -probability at least $1 - \epsilon$, for any $\rho_1, \rho_2 \in \mathcal{M}^1_+(\mathcal{F})$, $\rho_2 r' - \rho_1 r' + \rho_1 r - \rho_2 r \leq \frac{2\lambda}{N} (\rho_1 \otimes \rho_2) \overline{\mathbb{P}}_{\cdot,\cdot} + \frac{\mathcal{K}_{1,2}^{\text{loc}}}{(1-\xi)\lambda}$

Application to VC theory (1/3)

•
$$\mathbb{X} \triangleq X_1^{2N}$$

• $\mathcal{A}(\mathbb{X}) \triangleq \left\{ \left\{ f \in \mathcal{F} : \forall 1 \le i \le N, f(X_i) = \sigma_i \right\}; \sigma_1^{2N} \in \{0; 1\}^{2N} \right\} \right\}$

- $N(\mathbb{X}) \triangleq |\mathcal{A}(\mathbb{X})| = |\{[f(X_k)]_{k=1}^{2N} : f \in \mathcal{F}\}|$
- $\pi_{\mathcal{U}(\mathbb{X})}$: exchangeable distribution uniform on $\mathcal{A}(\mathbb{X})$ to the extent that $\pi_{\mathcal{U}(\mathbb{X})}(A) = \frac{1}{N(\mathbb{X})}$ for any $A \in \mathcal{A}(\mathbb{X})$.

Theorem. With $\mathbb{P}^{\otimes 2N}$ -probability at least $1 - \epsilon$, for any $f_1, f_2 \in \mathcal{F}$, $r'(f_2) - r'(f_1) \leq r(f_2) - r(f_1) + \sqrt{\frac{8\bar{\mathbb{P}}_{f_1,f_2}\left[2\log N(\mathbb{X}) + \log(\epsilon^{-1})\right]}{N}}$. In particular, introducing $\tilde{f}' \triangleq \operatorname{argmin}_{\mathcal{F}} r'$, we obtain $r'(\hat{f}_{\text{ERM}}) - r'(\tilde{f}') \leq r(\hat{f}_{\text{ERM}}) - r(\tilde{f}') + \sqrt{\frac{8\bar{\mathbb{P}}_{\hat{f}_{\text{ERM}},\tilde{f}'}\left[2\log N(\mathbb{X}) + \log(\epsilon^{-1})\right]}{N}}$.

Application to VC theory (2/3)

Localized theorem. For any $\lambda \geq 0$, define $C_{\lambda}(f) \triangleq \log \sum_{A \in \mathcal{A}(\mathbb{X})} \exp \left\{ -\lambda \left[(r+r')_{A} - (r+r')(f) \right] \right\}.$ Let $C(f,g) \triangleq \min_{\lambda \geq 0} \left\{ C_{\lambda}(f) + C_{\lambda}(g) \right\}.$ For any $\epsilon > 0$, with $\mathbb{P}^{\otimes 2N}$ -probability at least $1 - \epsilon$,

 $r'(\hat{f}_{\text{ERM}}) - r'(\tilde{f}') \leq r(\hat{f}_{\text{ERM}}) - r(\tilde{f}') + \sqrt{\frac{8\bar{\mathbb{P}}_{\hat{f}_{\text{ERM}},\tilde{f}'}[\mathcal{C}(\hat{f}_{\text{ERM}},\tilde{f}') + \log(\epsilon^{-1})]}{N}}.$

Application to VC theory (2/3)

Localized theorem. For any $\lambda \geq 0$, define $C_{\lambda}(f) \triangleq \log \sum_{A \in \mathcal{A}(\mathbb{X})} \exp \left\{ -\lambda \left[(r+r')_{A} - (r+r')(f) \right] \right\}.$ Let $C(f,g) \triangleq \min_{\lambda \geq 0} \left\{ C_{\lambda}(f) + C_{\lambda}(g) \right\}.$ For any $\epsilon > 0$, with $\mathbb{P}^{\otimes 2N}$ -probability at least $1 - \epsilon$,

$$r'(\hat{f}_{\text{ERM}}) - r'(\tilde{f}') \le r(\hat{f}_{\text{ERM}}) - r(\tilde{f}') + \sqrt{\frac{8\bar{\mathbb{P}}_{\hat{f}_{\text{ERM}},\tilde{f}'}[\mathcal{C}(\hat{f}_{\text{ERM}},\tilde{f}') + \log(\epsilon^{-1})]}{N}}$$

Illustration of localization efficiency by a toy example.

•
$$\mathcal{X} = [0; 1], \mathcal{F} = \{\mathbb{1}_{[\theta; 1]}; \theta \in [0; 1]\}$$

- $Y = \mathbb{1}_{X \ge \tilde{\theta}}$ for some $\tilde{\theta} \in [0; 1]$ and $\mathbb{P}(dX)$ absolutely continuous wrt Lebesgue measure.
- \rightsquigarrow Non localized inequality gives $r'(\hat{f}_{\text{ERM}}) \leq \frac{8\log(2N+1) + 4\log(\epsilon^{-1})}{N}$
- \rightsquigarrow Localized inequality gives $r'(\hat{f}_{\text{ERM}}) \leq \frac{37 + 5\log(\epsilon^{-1})}{N}$

	McAllester's pioneering work	The different PAC-Bayes bounds	Main ideas in Chap.1	Application to linear least squares
Framewo	ork			

Target

•
$$\ell(y, y') = (y - y')^2$$

•
$$R(f) = \mathbb{E}(Y - f(X))^2$$

• $\varphi_1, \ldots, \varphi_d$ functions from \mathcal{X} to \mathbb{R}

$$X \longrightarrow \begin{pmatrix} \varphi_1(X) \\ \vdots \\ \varphi_d(X) \end{pmatrix} = \varphi(X)$$

• $\Theta \subset \mathbb{R}^d$ closed convex

•
$$\mathcal{F} = \left\{ f_{\theta} = \sum_{j=1}^{d} \theta_{j} \varphi_{j}; \theta = (\theta_{1}, \dots, \theta_{d}) \in \Theta \right\}$$

Goal: predict as well as f^{*} ∈ argmin_{f∈F}R(f) (which is possibly different from f^(reg) : x → E(Y|X = x))

Context McAllester's pioneering work The different PAC-Bayes bounds Main ideas in Chap.1

Application to linear least squares

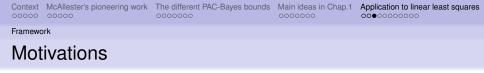
▲□▶▲□▶▲□▶▲□▶ □ のQ@

Framework

Decomposition of the risk

- Gram matrix: $Q = \mathbb{E}[\varphi(X)\varphi^T(X)]$
- The risk is a guadratic form with matrix Q:

$$egin{aligned} & \mathcal{R}(\mathit{f}_{ heta}) = \mathbb{E}(\mathit{Y} - \mathit{ heta}^{\mathsf{T}}arphi(\mathit{X}))^{\mathsf{2}} \ & = \mathbb{E}\mathit{Y}^{\mathsf{2}} - \mathit{2}\mathit{ heta}^{\mathsf{T}}\mathbb{E}[arphi(\mathit{X})\mathit{Y}] + \mathit{ heta}^{\mathsf{T}}\mathit{ extsf{Q}}\mathit{ heta} \end{aligned}$$



Better understanding of the parametric linear least squares regression

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三■ - のへぐ

- Central task for nonparametric regression with linear approximation space
- Two-stage model selection

Context	McAllester's pioneering work	The different PAC-Bayes bounds	Main ideas in Chap.1	Application to linear least squares
00000	00000	000000	000000	000000000

Variants of known results

Ordinary least squares and empirical risk minimization

- Linear aggregation: F = F_{lin} = span{φ₁,...,φ_d} and f^{*}_{lin} = f^{*}
- Let $\hat{f}^{(\text{ols})} \in \operatorname{argmin}_{f \in \mathcal{F}_{\text{lin}}} \frac{1}{n} \sum_{i=1}^{n} [Y_i f(X_i)]^2$.
- $\mathbb{E}R(\hat{f}^{(\text{ols})}) R(f^*_{\text{lin}}) = \mathbb{E}[\hat{f}^{(\text{ols})}(X) f^*_{\text{lin}}(X)]^2.$
- if $\sup_{x \in \mathcal{X}} \operatorname{Var}(Y|X = x) = \sigma^2 < +\infty$ and $f^{(\operatorname{reg})} = f^*_{\operatorname{lin}}$, we have

$$\mathbb{E}\left\{\frac{1}{n}\sum_{i=1}^{n}\left[\hat{f}^{(\text{ols})}(X_{i})-f_{\text{lin}}^{*}(X_{i})\right]^{2}\right\}\leq\sigma^{2}\frac{d}{n}$$

(日) (日) (日) (日) (日) (日) (日)

• It does not imply a $\frac{d}{n}$ upper bound on $\mathbb{E}R(\hat{f}^{(\text{ols})}) - R(f^*_{\text{lin}})$.

	McAllester's pioneering work	The different PAC-Bayes bounds	Main ideas in Chap.1	Application to linear least squares
Varianta	of known rooulto			

Theorem (Györfi, Kohler, Krzyżak, Walk, 2004)

If
$$\sup_{x \in \mathcal{X}} Var(Y|X = x) = \sigma^2 < +\infty$$
 and

$$\|f^{(\mathsf{reg})}\|_{\infty} = \sup_{x \in \mathcal{X}} |f^{(\mathsf{reg})}(x)| \le H$$

for some H > 0, then the truncated estimator $\hat{f}_{H}^{(ols)} = (\hat{f}^{(ols)} \land H) \lor -H$ satisfies

$$\mathbb{E}R(\hat{f}_{H}^{(\text{ols})}) - R(f^{(\text{reg})}) \\ \leq 8[R(f_{\text{lin}}^{*}) - R(f^{(\text{reg})})] + \kappa \frac{(\sigma^{2} \vee H^{2})d\log n}{n}$$

for some numerical constant κ .

Context	McAllester's pioneering work	The different PAC-Bayes bounds	Main ideas in Chap.1	Application to linear least squares
00000	00000	0000000	0000000	000000000

Variants of known results

Theorem (Birgé, Massart, 1998)

Assume that for any f_1, f_2 in \mathcal{F} , $\|f_1 - f_2\|_{\infty} \leq H$ and $\exists f_0 \in \mathcal{F}$ satisfying

for any
$$x \in \mathcal{X}, \quad \mathbb{E}\Big\{ \exp\Big[A^{-1} \big| \, Y - f_0(X) \big| \Big] \, \Big| \, X = x \Big\} \leq M,$$

for some positive constants A and M. Let

$$\tilde{B} = \inf_{\phi_1, \dots, \phi_d} \sup_{\theta \in \mathbb{R}^d - \{0\}} \frac{\|\sum_{j=1}^d \theta_j \phi_j\|_{\infty}^2}{\|\theta\|_{\infty}^2}$$

where the infimum is taken w.r.t. all possible orthonormal basis of \mathcal{F} for $\langle f_1, f_2 \rangle = \mathbb{E} f_1(X) f_2(X)$. Then, with probability at least $1 - \epsilon$:

$$R(\hat{f}^{(ext{erm})}) - R(f^*) \leq \kappa (A^2 + H^2) rac{d \log[2 + (ilde{B}/n) \wedge (n/d)] + \log(\epsilon^{-1})}{n},$$

where κ is a positive constant depending only on M.

 Context
 McAllester's pioneering work
 The different PAC-Bayes bounds
 Main ideas in Chap.1
 Application to linear least squares

 00000
 000000
 0000000
 0000000
 0000000
 00000000

Variants of known results

Projection estimator

Theorem (Tsybakov, 2003)

Let ϕ_1, \ldots, ϕ_d be an o.n.b. of \mathcal{F}_{lin} for $\langle f_1, f_2 \rangle = \mathbb{E}f_1(X)f_2(X)$. The projection estimator on this basis is $\hat{f}^{(\text{proj})} = \sum_{i=1}^d \hat{\theta}_i^{(\text{proj})}\phi_i$, with

$$\hat{\theta}^{(\text{proj})} = \frac{1}{n} \sum_{i=1}^{n} Y_i \phi_i(X_i).$$

	r
h	t.

$$\sup_{x\in\mathcal{X}} \operatorname{Var}(Y|X=x) = \sigma^2 < +\infty$$

and

$$\|f^{(\operatorname{reg})}\|_{\infty} = \sup_{x\in\mathcal{X}} |f^{(\operatorname{reg})}(x)| \le H < +\infty,$$

then we have

$$\mathbb{E}R(\hat{f}^{(\mathsf{proj})}) - R(f^*_{\mathsf{lin}}) \le (\sigma^2 + H^2)\frac{d}{n}.$$

Context	McAllester's pioneering work	The different PAC-Bayes bounds	Main ideas in Chap.1	Application to linear least squares
00000	00000	000000	0000000	00000000000

Variants of known results

Conclusion of the survey

- $R(\hat{f}^{(\text{erm})}) R(f^*) = O(\frac{d \log(2+n/d) + \log(e^{-1})}{n})$ for L_{∞} -bounded \mathcal{F} and exponential moments
- There is no simple d/n which does not require strong assumptions
- Degraded convergence rate when Q is ill-conditioned ?

(日) (日) (日) (日) (日) (日) (日)

	McAllester's pioneering work	The different PAC-Bayes bounds	Main ideas in Chap.1	Application to linear least squares
New resu	ults			

Theorem

Assume $\mathbb{E}[\|\varphi(X)\|^4] < +\infty$ and $\sup \mathbb{E}\{[Y - \tilde{f}(X)]^2 | X\} \le \sigma^2$. For any $\epsilon > 0$, there is n_{ϵ} s.t. for any $n \ge n_{\epsilon}$, with proba. at least $1 - \epsilon$, $R(\hat{f}^{(\text{erm})}) \le R(f^*_{\text{lin}}) + \sigma^2 \frac{30d + 1000 \log(3\epsilon^{-1})}{n}$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

	McAllester's pioneering work	The different PAC-Bayes bounds	Main ideas in Chap.1	Application to linear least squares
Now roci	ulte			

- Θ bounded
- π uniform distribution on \mathcal{F}
- λ > 0

•
$$W_i(f, f') = \frac{\lambda}{n} \left\{ \left[Y_i - f(X_i) \right]^2 - \left[Y_i - f'(X_i) \right]^2 \right\}$$

•
$$\hat{\mathcal{E}}(f) = \log \mathbb{E}_{\pi(df')} \frac{1}{\prod_{i=1}^{n} [1 - W_i(f, f') + \frac{1}{2} W_i(f, f')^2]}$$

- We consider the "posterior" distribution $\hat{\pi} = \pi_{-\hat{\mathcal{E}}(f)}$
- for $\frac{\lambda}{n}$ small enough, $1 W_i(f, f') + \frac{1}{2}W_i(f, f')^2$ is close to $e^{-W_i(f, f')}$, and consequently

$$\hat{\mathcal{E}}(f) \approx \lambda r(f) + \log \mathbb{E}_{\pi(df')} e^{-\lambda r(f')},$$

and

$$\hat{\pi} \approx \pi_{-\lambda r}$$

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

	McAllester's pioneering work	The different PAC-Bayes bounds	Main ideas in Chap.1	Application to linear least squares
New res	ults			

Theorem

Assume $\sup_{f_1, f_2 \in \mathcal{F}} \|f_1 - f_2\|_{\infty} \leq H$ and, for some $\sigma > 0$,

$$\sup_{\mathbf{X}\in\mathcal{X}}\mathbb{E}\big\{[\mathbf{Y}-f^*(\mathbf{X})]^2\big|\mathbf{X}=\mathbf{X}\big\}\leq\sigma^2<+\infty.$$

Let $\lambda = \frac{n}{3(2\sigma+H)^2}$ and \hat{f} be a prediction function drawn from the distribution $\hat{\pi}$. Then for any $\epsilon > 0$, with probability at least $1 - \epsilon$, we have

$$R(\widehat{f}) - R(f^*) \leq 17(2\sigma + H)^2 \, rac{d + \log(2\epsilon^{-1})}{n}.$$

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの