Introduction
Pac-Bayesian oracle inequalities
Application to French GDP and quantile prediction

# PAC-Bayesian bounds for dependent observations
## Joint work with Pierre Alquier & Olivier Wintenberger

Xiaoyin LI

Xiaoyin.li@u-cergy.fr

**University of Cergy-Pontoise**

TEST February 2013

**Introduction**
Pac-Bayesian oracle inequalities
Application to French GDP and quantile prediction

Context
Estimators

## Outline

**Introduction**
Pac-Bayesian oracle inequalities
Application to French GDP and quantile prediction

**Context**
Estimators

## A problem of statistical inference

Let $(X_t)_{t \in \mathbb{Z}}$ be a stationary time series in $\mathbb{R}^p$. We want to learn to forecast the series from the observations $X_1, \ldots, X_n$.

We have a family of predictors

$$\mathcal{F} = \left\{ f_\theta : (\mathbb{R}^p)^k \to \mathbb{R}^p \text{ mesurable}, \theta \in \Theta \right\}$$

We consider a model-selection type approach:

$$\Theta = \bigcup_{j=1}^m \Theta_j.$$

Objective: find a $\theta \in \Theta$ such that $f_\theta(X_{t-1}, \ldots, X_{t-k})$ is a good prediction of $X_t$.

**Introduction**
Pac-Bayesian oracle inequalities
Application to French GDP and quantile prediction

**Context**
Estimators

## Familes of classical predictor

### Definition

For any $\theta \in \Theta$,

$$\hat{X}_t^\theta = f_\theta(X_{t-1}, \ldots, X_{t-k}).$$

Linear auto-regressive class of predictors:

$$f_\theta(X_{t-1}, \ldots, X_{t-k}) = \theta_0 + \sum_{j=1}^{k} \theta_j X_{t-j}.$$

Non-parametric auto-regression predictors:

$$f_\theta(X_{t-1}, \ldots, X_{t-k}) = \sum_{i=1}^{j} \theta_i \varphi_i(X_{t-1}, \ldots, X_{t-k}).$$

**Introduction**
Pac-Bayesian oracle inequalities
Application to French GDP and quantile prediction

**Context**
Estimators

## Measure of risk

Let $\ell$ be a loss function: $\ell(\hat{X}_t^\theta, X_t) \geq 0$ measures the forecasting error of predictor $\theta$ at time $t$.

The prevision risk is defined as:

$$R(\theta) = \mathbb{E}\left[\ell\left(\hat{X}_t^\theta, X_t\right)\right] \text{ is } \textit{unknown}.$$

On the other hand, we observe the empirical risk:

$$r_n(\theta) = \frac{1}{n-k} \sum_{i=k+1}^{n} \ell\left(\hat{X}_i^\theta, X_i\right).$$

with $\mathbb{E}\left[r_n(\theta)\right] = R(\theta)$.

**Introduction**
Pac-Bayesian oracle inequalities
Application to French GDP and quantile prediction

**Context**
Estimators

## Objective

Build a parameter $\theta$ on the basis of observations $X_1, \ldots, X_n$ such that :

$$R\left(\theta\right) \simeq \inf_{\theta \in \Theta} R(\theta).$$

More precisely,

$$R\left(\theta\right) \leq \left\{ \inf_{\theta \in \Theta} R(\theta) + \Delta(n, \Theta) \right\}$$

**Introduction**
Pac-Bayesian oracle inequalities
Application to French GDP and quantile prediction

Context
**Estimators**

## Prior and Estimators

Let $\mathcal{M}_+^1(\Theta)$ denote the set of all probability measures on $(\Theta, \mathcal{T})$. Let us take $\pi \in \mathcal{M}_+^1(\Theta)$, the prior.

### Definition (Gibbs Estimators)

*We put, for any $\lambda > 0$,*

$$\hat{\theta}_\lambda = \int_\Theta \theta \hat{\rho}_\lambda(\mathrm{d}\theta)$$

*where*

$$\hat{\rho}_\lambda(\mathrm{d}\theta) = \frac{e^{-\lambda r_n(\theta)} \pi(\mathrm{d}\theta)}{\int e^{-\lambda r_n(\theta')} \pi(\mathrm{d}\theta')}.$$

Introduction
**Pac-Bayesian oracle inequalities**
Application to French GDP and quantile prediction

A basic PAC-Bayesian Bound
$\theta$-weakly dependent
Some examples

# Outline

Introduction
**Pac-Bayesian oracle inequalities**
Application to French GDP and quantile prediction

A basic PAC-Bayesian Bound
$\theta$-weakly dependent
Some examples

## Overview of the results

The idea is that the risk of the Gibbs estimator will be close to $\inf_{\theta} R(\theta)$ up to a small remainder term $\Delta$ called the rate of convergence.

For the sake of simplicity, let $\overline{\theta} \in \Theta$ be such that

$$R(\overline{\theta}) = \inf_{\theta} R(\theta).$$

We want to prove that our estimators satisfy, for any $\varepsilon > 0$,

$$\mathbb{P}\left(R\left(\hat{\theta}\right) \leq R(\overline{\theta}) + \Delta(n, \lambda, \pi, \varepsilon)\right) \geq 1 - \varepsilon$$

where $\Delta(n, \lambda, \pi, \varepsilon) \to 0$ as $n \to \infty$ for some $\lambda = \lambda(n)$.

Introduction
**Pac-Bayesian oracle inequalities**
Application to French GDP and quantile prediction

A basic PAC-Bayesian Bound
$\theta$-weakly dependent
Some examples

## Assumptions

We assume:

- $(X_t)_{t \in \mathbb{Z}}$ is bounded and $\theta$-dependent, *ie* a.s. $\|X_0\|_\infty \leq \mathcal{B} < \infty$ and $\theta_{\infty,n}(1) \leq \mathcal{C} < \infty$.
- $\ell(x, x') = g(x - x')$ for some convex function $g$ and $g$ is $K$-Lipschitz.
- for any $f$,

$$\|f_\theta (x_1, \ldots, x_k) - f_\theta (y_1, \ldots, y_k)\| \leq \sum_{j=1}^{k} a_j (\theta) \|x_j - y_j\|.$$

$$L := \sup_{\theta \in \Theta} \sum_{j=1}^{k} a_j (\theta)$$

- $k = k(n) \leq n/2$.

Introduction
Pac-Bayesian oracle inequalities
Application to French GDP and quantile prediction
A basic PAC-Bayesian Bound
$\theta$-weakly dependent
Some examples

# A basic PAC-Bayesian Bound

## Theorem

*for any $\lambda > 0$, for any $\varepsilon > 0$,*

$$\mathbb{P}\left\{\int R(\hat{\theta}_\lambda) \leq \inf_\rho \left[\int R\mathrm{d}\rho + \frac{2\lambda\kappa_n^2}{n} + 2\frac{\mathcal{K}(\rho,\pi) + \log\left(\frac{2}{\varepsilon}\right)}{\lambda}\right]\right\}$$

$$\geq 1 - \varepsilon$$

*where $\kappa_n := \sqrt{2}K(1+L)(\mathcal{B} + \theta_{\infty,n}(1))$.*

Introduction
**Pac-Bayesian oracle inequalities**
Application to French GDP and quantile prediction

**A basic PAC-Bayesian Bound**
$\theta$-weakly dependent
Some examples

# Example of loss function

1. Absolute loss $\ell(x, x') = \|x - x'\|$ with $K = 1$ (Alquier and Wintenberger).
2. Quadratic loss $\ell(x, x') = \|x - x'\|^2$ with $K = 4\mathcal{B}$ (Meir).
3. Quantile loss

$$\ell_\tau(x, y) = \begin{cases} \tau (x - y), & \text{if } x - y > 0 \\ - (1 - \tau) (x - y), & \text{otherwise} \end{cases}.$$

with $K = \max(\tau, 1 - \tau) \leq 1$(Alquier and Li).

Introduction
**Pac-Bayesian oracle inequalities**
Application to French GDP and quantile prediction

A basic PAC-Bayesian Bound
$\theta$-weakly dependent
Some examples

## $\theta$ weak dependent coefficient

Introduced by Doukhan and Louhichi (SPA, 1999). $\theta$ coefficient is as follow:

$$\theta_\infty(\mathfrak{S}, Z) = \sup_{f \in \Lambda_1^q} \left\| \mathbb{E}\left[f(Z)|\mathfrak{S}\right] - \mathbb{E}\left[f(Z)\right] \right\|_\infty$$

where

$$\Lambda_1^q = \left\{ f : (\mathbb{R}^p)^q \to \mathbb{R}, \quad \frac{|f(z_1, \ldots, z_q) - f(z_1', \ldots, z_q')|}{\sum_{j=1}^q \|z_j - z_j'\|} \leq 1 \right\},$$

and that

$$\theta_{\infty,k}(1) := \sup_{p < j_1 < \ldots < j_\ell, 1 \leq \ell \leq k} \left\{ \theta_\infty(\sigma(X_t, t \leq p), (X_{j_1}, \ldots, X_{j_\ell})) \right\}.$$

Introduction    A basic PAC-Bayesian Bound
**Pac-Bayesian oracle inequalities**    $\theta$-weakly dependent
Application to French GDP and quantile prediction    Some examples

## Example of $\theta$-weakly dependent series

For any series

$$X_t = F(\varepsilon_t, \varepsilon_{t-1}, \varepsilon_{t-2}, \ldots)$$

with $\varepsilon_i$ iid, upper bounded by $b$, and

$$\|F(x_1, x_2, \ldots) - F(x_1', x_2', \ldots)\| \leq \sum_{j=1}^{\infty} a_j \|x_j - x_j'\|$$

we have

$$\theta_{\infty,n}(1) \leq 2b \sum_{j=1}^{\infty} j a_j.$$

Introduction
**Pac-Bayesian oracle inequalities**
Application to French GDP and quantile prediction

A basic PAC-Bayesian Bound
**θ-weakly dependent**
Some examples

# Examples of $\theta$-weakly dependent series

Let us remind the $\phi$-mixing coefficient:

$$\phi(r) = \sup_{A \in \sigma(X_t, t \leq 0), B \in \sigma(X_t, t \geq r)} |P(B|A) - P(B)|$$

Then, for $(X_t)$ upper bounded by $\mathcal{B}$, we have

$$\theta_{\infty,n}(1) \leq 2\mathcal{B} \sum_{r=1}^{n} \phi(r).$$

Introduction
**Pac-Bayesian oracle inequalities**
Application to French GDP and quantile prediction

A basic PAC-Bayesian Bound
$\theta$-weakly dependent
**Some examples**

## Reminder

### Theorem

*for any $\lambda > 0$, for any $\varepsilon > 0$,*

$$\mathbb{P}\left\{ \int R(\hat{\theta}_\lambda) \leq \inf_\rho \left[ \int R \mathrm{d}\rho + \frac{2\lambda\kappa_n^2}{n} + \frac{2\mathcal{K}(\rho,\pi) + 2\log\left(\frac{2}{\varepsilon}\right)}{\lambda} \right] \right\}$$

$$\geq 1 - \varepsilon$$

*where $\kappa_n := \sqrt{2}K(1+L)(\mathcal{B} + \theta_{\infty,n}(1))$.*

Introduction    A basic PAC-Bayesian Bound
**Pac-Bayesian oracle inequalities**    $\theta$-weakly dependent
Application to French GDP and quantile prediction    **Some examples**

# Toy example: $\mathrm{card}(\Theta) = M < \infty$

We take $\pi$ as the uniform distribution:

$$
\begin{aligned}
R(\hat{\theta}_\lambda) &\leq \inf_\rho \left\{ \int R \mathrm{d}\rho + \frac{2\lambda\kappa_n^2}{n} + \frac{2\mathcal{K}(\rho,\pi) + 2\log\left(\frac{2}{\varepsilon}\right)}{\lambda} \right\} \\
&\leq \inf_\theta \left\{ R(\theta) + \frac{2\lambda\kappa_n^2}{n} + 2\frac{\log(M) + \log\left(\frac{2}{\varepsilon}\right)}{\lambda} \right\}
\end{aligned}
$$

### Theorem

*Assume that $\mathrm{card}(\Theta) = M$ and let $\pi$ be the uniform probability distribution on $\Theta$. Then the oracle inequality is satisfied for any $\lambda > 0$, $\varepsilon > 0$ with probability at least $1 - \varepsilon$*

$$
R(\hat{\theta}_\lambda) \leq \inf_\theta \left\{ R(\theta) + \frac{2\lambda\kappa_n^2}{n} + 2\frac{\log(M) + \log\left(\frac{2}{\varepsilon}\right)}{\lambda} \right\}
$$

Introduction    A basic PAC-Bayesian Bound
**Pac-Bayesian oracle inequalities**    $\theta$-weakly dependent
Application to French GDP and quantile prediction    **Some examples**

# Toy example: $\text{card}(\Theta) = M < \infty$

The choice $\lambda = \sqrt{n \log(M)}/\kappa_n$ yields the oracle inequality:

$$R(\hat{\theta}_\lambda) \leq \inf_\theta R + 2\kappa_n \sqrt{\frac{2 \log(M)}{n}} + \frac{2\kappa_n \log\left(\frac{2}{\varepsilon}\right)}{n \log(M)}$$

Bad news: the optimal $\lambda$ depends on $\theta_{\infty,n}(1)$, unknown in practice.

Introduction
**Pac-Bayesian oracle inequalities**
Application to French GDP and quantile prediction

A basic PAC-Bayesian Bound
$\theta$-weakly dependent
**Some examples**

## Linear autoregressive predictors

Consider the linear autoregressive model of $\mathrm{AR}(k)$ predictors

$$f_\theta(x_{t-1}, \ldots, x_{t-k}) = \sum_{j=1}^{k} \theta_j x_{t-j}$$

with $\theta \in \Theta = \{\theta \in \mathbb{R}^k, \|\theta\| \leq L\}$
We take $\pi$ uniform and we restrict $\rho$ to the uniform distributions on
$\{\theta' : \|\theta - \theta'\| \leq \delta\}$.

$$R(\hat{\theta}_\lambda) \leq \inf_\rho \left\{ \overbrace{\int R \mathrm{d}\rho_{\delta,\theta}}^{\leq R(\theta) + \delta \mathcal{B}} + \frac{2\lambda \kappa_n^2}{n} + \overbrace{\frac{2\mathcal{K}(\rho_{\delta,\theta}, \pi) + 2\log\left(\frac{2}{\varepsilon}\right)}{\lambda}}^{2\frac{k \log\left(\frac{L+1}{\delta}\right) + \log\left(\frac{2}{\varepsilon}\right)}{\lambda}} \right\}.$$

$$\leq \inf_\theta \left\{ R(\theta) + \delta \mathcal{B} + \frac{2\lambda \kappa_n^2}{n} + 2\frac{k \log\left(\frac{L+1}{\delta}\right) + \log\left(\frac{2}{\varepsilon}\right)}{\lambda} \right\}.$$

Introduction
**Pac-Bayesian oracle inequalities**
Application to French GDP and quantile prediction

A basic PAC-Bayesian Bound
$\theta$-weakly dependent
**Some examples**

## Linear autoregressive predictors

We can now take $\delta = \frac{k}{\lambda \mathcal{B}}$,

$$R(\hat{\theta}_\lambda) \leq \inf_\theta R + \frac{2\lambda\kappa_n^2}{n} + 2\frac{k \log\left(\frac{K\mathcal{B}(L+1)\sqrt{e}\lambda}{k}\right) + \log\left(2/\varepsilon\right)}{\lambda}.$$

The optimal inverse temperature parameter is $\lambda = \frac{\sqrt{nk}}{\kappa_n}$.

### Theorem

*Let $\pi$ be the uniform probability distribution on $\Theta$. Then the oracle inequality is satisfied for any $\lambda > 0$, $\varepsilon > 0$ with high probability at least $1 - \varepsilon$*

$$R(\hat{\theta}_\lambda) \leq \inf_\theta R + \frac{2\lambda\kappa_n^2}{n} + 2\frac{k \log\left(\frac{K\mathcal{B}(L+1)\sqrt{e}\lambda}{k}\right) + \log\left(2/\varepsilon\right)}{\lambda}.$$

Introduction
**Pac-Bayesian oracle inequalities**
Application to French GDP and quantile prediction

A basic PAC-Bayesian Bound
$\theta$-weakly dependent
**Some examples**

## General parametric model

We state a general result about finite-dimensional families of predictors.
In general, one can always consider

$$\rho_\delta(\mathrm{d}\theta) \propto \pi(\mathrm{d}\theta)\mathbf{1}(R(\theta) - \inf_\Theta R \le \delta).$$

Let $\pi$ be uniform and $\rho$ be the uniform distributions on $\{\theta : R(\theta) - \inf_\Theta R(\theta) \le \delta\}$.
we assume $\dim(\Theta, \pi) := \sup \frac{-\log \pi\{\theta : R(\theta) - \inf_\Theta R \le \delta\}}{\log \lambda} = D$

$$R(\hat{\theta}_\lambda) \le \inf_\rho \left\{ \underbrace{\int R \mathrm{d}\rho_{\delta,\theta}}_{\le R(\theta) + \delta} + \frac{2\lambda\kappa_n^2}{n} + \underbrace{\frac{2\mathcal{K}(\rho_{\delta,\theta}, \pi) + 2\log\left(\frac{2}{\varepsilon}\right)}{\lambda}}_{2\frac{-\log\left(\pi\left\{\theta : R(\theta) - \inf_\Theta R \le \delta\right\}\right) + \log\left(\frac{2}{\varepsilon}\right)}{\lambda}} \right\}.$$

$$\le \inf_\Theta R + \delta + \frac{2\lambda\kappa_n^2}{n} + 2\frac{d\log(D/\delta) + \log\left(\frac{2}{\varepsilon}\right)}{\lambda}$$

Introduction
**Pac-Bayesian oracle inequalities**
Application to French GDP and quantile prediction

A basic PAC-Bayesian Bound
$\theta$-weakly dependent
**Some examples**

The infimum is reached for $\delta = d/\lambda$ and we have:

$$R(\hat{\theta}_\lambda) \leq R(\bar{\theta}) + 2\frac{\lambda \kappa_n^2}{n} + 2\frac{d \log(D\sqrt{e}\lambda/d) + \log\left(\frac{2}{\varepsilon}\right)}{\lambda}$$

### Theorem

*Let $\pi$ be the uniform probability distribution on $\Theta$. Then the oracle inequality is satisfied for any $\lambda > 0$, $\varepsilon > 0$ with high probability at least $1 - \varepsilon$*

$$R\left(\hat{\theta}_\lambda\right) \leq R(\bar{\theta}) + \frac{2\lambda \kappa_n^2}{n} + 2\frac{d \log\left(\frac{D\sqrt{e}\lambda}{d}\right) + \log\left(\frac{2}{\varepsilon}\right)}{\lambda}.$$

Introduction
**Pac-Bayesian oracle inequalities**
Application to French GDP and quantile prediction

A basic PAC-Bayesian Bound
$\theta$-weakly dependent
**Some examples**

## Tools used in the proofs

### Lemma (Donsker-Varadhan variational formula)

*For any $\pi \in \mathcal{M}_+^1(E)$, for any measurable upper-bounded function $h$, we have:*

$$\int \exp(h)\mathrm{d}\pi = \exp\left(\sup_{\rho \in \mathcal{M}_+^1(E)} \left(\int h\mathrm{d}\rho - \mathcal{K}(\rho,\pi)\right)\right)$$

### Theorem (Rio,2000)

*For any $f$ that is a function $1$-Lipshitz*

$$\forall t \geq 0, \;\; \mathbb{E}[e^{tf(X_1,\ldots,X_n) - t\mathbb{E}[f(X_1,\ldots,X_n)]}] \leq e^{\frac{nt^2(\mathcal{B}+\theta_{\infty,n}(1))^2}{2}}$$

Introduction
Pac-Bayesian oracle inequalities
**Application to French GDP and quantile prediction**

# Outline

Introduction
Pac-Bayesian oracle inequalities
**Application to French GDP and quantile prediction**

## The context

Objective: at each quarter $t$, predict the flash estimate of GDP growth: $\Delta\mathrm{GDP}_t$.

Available information:

- $\Delta\mathrm{GDP}_{t'}$, for all $t' < t$
- $I_{t'}$, for all $t' < t$, $I_{t-1}$ is the climate indicator available to the INSEE at time $t$.
- The observation period is 1988-Q1 (1st quarter of 1988) to 2011-Q3.

Introduction
Pac-Bayesian oracle inequalities
Application to French GDP and quantile prediction

## Quantile loss function

We define $X_t = (\Delta \mathrm{GDP}_t, I_t)' \in \mathbb{R}^2$ . Following Cornec(CIRET conference 2010), we consider predictors of the form :

$$f_\theta(X_{t-1}, X_{t-2}) = \theta_0 + \theta_1 \Delta \mathrm{GDP}_{t-1} + \theta_2 I_{t-1} + \theta_3 (I_{t-1} - I_{t-2})|I_{t-1} - I_{t-2}|$$

These family of predictors allow to obtain a forecasting as precise as the INSEE one.

We use the quantile loss function :

$$\ell_\tau((\Delta \mathrm{GDP}_t, I_t), (\Delta' \mathrm{GDP}_t, I_t'))$$
$$= \begin{cases} \tau \left( \Delta \mathrm{GDP}_t - \Delta' \mathrm{GDP}_t \right), & \text{if } \Delta \mathrm{GDP}_t - \Delta' \mathrm{GDP}_t > 0 \\ -\left( 1 - \tau \right) \left( \Delta \mathrm{GDP}_t - \Delta' \mathrm{GDP}_t \right), & \text{otherwise.} \end{cases}$$

Introduction
Pac-Bayesian oracle inequalities
Application to French GDP and quantile prediction

# Results: GDP forecasting



**Out-of-sample forecasts**

Figure: French GDP online prediction using the quantile loss function with $\tau = 0.5$.

Introduction
Pac-Bayesian oracle inequalities
Application to French GDP and quantile prediction

# Results: Confidence intervals



Figure: French GDP online 50%-confidence intervals (left) and 90%-confidence intervals (right).

Introduction
Pac-Bayesian oracle inequalities
Application to French GDP and quantile prediction

# Reference

Alquier, P., Li, X. and Wintenberger, O.(2012). *Prediction of time series by statistical learning: general losses and fast rates.Preprint arXiv:1211.1847*

Alquier, P. and Li, X. (2012). *Prediction of Quantiles by Statistical Learning and Application to GDP Forecasting. in Proceeding of DS'12, Springer LNAI n.7569. pp.22-36*

Alquier, P. and Wintenberger, O. (2012). *Model selection for weakly dependent time series forecasting. Bernoulli, vol.18, no.3.pp.883-913*

Rio, E.(2000). *Inégalités de Hoeffding pour les fonctions lipshitziennes de suites dépendantes. CRAS série I, vol 330, pp 905-908*

Dedecker, J., Doukhan, P., Lang, G.,Léon, J.R., Louhichi, S. and Prieur, C. (2007). *Weak dependence, examples and applications. Springer Lecture Notes in Mathematics n.190*

Catoni, O.(2004). *Statistical learning theory and stochastic optimization. Saint-Flour 2001 lecture notes. J.Picard Ed. Springer Lecture Notes in Mathematics.*