Model selection theory: a tutorial with applications to learning

Pascal Massart Université Paris-Sud, Orsay

ALT 2012, October 29

- Asymptotic approach to model selection
- Idea of using some penalized empirical criterion goes back to the seminal works of Akaike ('70).

- Akaike celebrated criterion (AIC) suggests to penalize the log-likelihood by the number of parameters of the parametric model.

- This criterion is based on some asymptotic approximation that essentially relies on Wilks' Theorem

Wilks' Theorem: under some proper regularity conditions the log-likelihood $L_n(\theta)$ based on n i.i.d. observations with distribution belonging to a parametric model with D parameters obeys to the following weak convergence result

$$2\left(L_n\left(\hat{\theta}\right) - L_n\left(\theta_0\right)\right) \to \chi^2(D)$$

where $\hat{\theta}$ denotes the MLE and θ_0 is the true value of the parameter.

• Non asymptotic Theory

In many situations, it is usefull to make the size of the models tend to infinity or make the list of models depend on n. In these situations, classical asymptotic analysis breaks down and one needs to introduce an alternative approach that we call non asymptotic.

We still like



But the size of the models as well the size of the list of models should be authorized to be large too.

Functional estimation

• The basic problem

Construct estimators of some function s, using as few prior information on s as possible. Some typical frameworks are the following.

• Density estimation

- $(X_1,...,X_n)$ i.i.d. sample with unknown density s with respect to some given measure μ .
- Regression framework

One observes $(X_1, Y_1), ..., (X_n, Y_n)$ With $Y_i = s(X_i) + \varepsilon_i$

The explanatory variables X_i are fixed or i.i.d. The errors \mathcal{E}_i are i.i.d. with $E[\varepsilon_i | X_i] = 0$

• Binary classification

We consider an i.i.d. regression framework where the response variable Y is a « label » :0 or 1. A basic problem in *statistical learning* is to estimate the best classifier $s(x) = 1_{\eta(x) \ge 1/2}$, where η denotes the regression function

 $\eta(x) = E[Y|X = x]$

• Gaussian white noise

Let **s** be a numerical function on $\begin{bmatrix} 0,1 \end{bmatrix}$. One observes the process $Y^{(n)}$ on $\begin{bmatrix} 0,1 \end{bmatrix}$ defined by

$$dY^{(n)}(x) = s(x)dx + \frac{1}{\sqrt{n}}dB(x), Y^{(n)}(0) = 0$$

Where **B** is a Brownian motion. The level of noise is written as $1/\sqrt{n}$ by allow an easy comparison.

Empirical Risk Minimization (ERM)

A classical strategy to estimate s consists of taking a set of functions S (a « model ») and consider some *empirical criterion* (based on the data) γ_n such that

 $t \to E[\gamma_n(t)]$

achieves a minimum at point t = s. The ERM estimator \hat{s} of s minimizes γ_n over s. One can hope that \hat{s} is close to s, if the target s belongs to model s(or at least is not far from s). This approach is most popular in the parametric case (i.e. when s is defined by a finite number of parameters and one assumes that $s \in S$). • Maximum likelihood estimation (MLE)

Context:density estimation (i.i.d. setting to be simple) $(X_1, ..., X_n)$ i.i.d. sample with distribution $sd\mu$

with

$$\gamma_n(t) = -\frac{1}{n} \sum_{i=1}^n \log t(X_i)$$

 $E\left[\gamma_n\left(t\right)-\gamma_n\left(s\right)\right]=K\left(s,t\right)\geq 0$

Kullback Leibler information

Least squares
 Regression

with

$$\gamma_n(t) = \frac{1}{n} \sum_{i=1}^n \left(Y_i - t(X_i) \right)^2$$

Nhite noise

$$E[\gamma_n(t) - \gamma_n(s)] = \frac{1}{n} \sum_{i=1}^n E[(t-s)^2(X_i)] \ge 0$$

with

$$\gamma_n(t) = \left\| t \right\|^2 - 2 \int t(x) dY^{(n)}(x)$$

Density

$$E\left[\gamma_{n}\left(t\right)-\gamma_{n}\left(s\right)\right]=\left\|t-s\right\|^{2}\geq0$$

with

$$\gamma_n(t) = \left\| t \right\|^2 - \frac{2}{n} \sum_{i=1}^n t(X_i)$$
$$E\left[\gamma_n(t) - \gamma_n(s)\right] = \left\| t - s \right\|^2 \ge 0$$

Exact calculations in the linear case In the white noise or the density frameworks, when S is a finite dimensional subspace of $L_2(\mu)$ (where μ denotes the Lebesgue measure in the white noise case), the LSE can be explicitly computed. Let $(\phi_{\lambda})_{\lambda \in \Lambda}$ be some orthonormal basis of S, then

$$\hat{s} = \sum_{\lambda \in \Lambda} \hat{\beta}_{\lambda} \phi_{\lambda}$$

$$\hat{\beta}_{\lambda} = \int \phi_{\lambda}(x) dY^{(n)}(x) \quad \text{or} \quad \hat{\beta}_{\lambda} = \frac{1}{n} \sum_{i=1}^{n} \phi_{\lambda}(X_{i})$$
White noise Density

- The model choice paradigm
- If a model S is defined by a « small » number of parameters (as compared to n), then the target s can happen to be far from the model.
- If the number of parameters is taken too large then \hat{s} will be a poor estimator of s even if s truly belongs to S.

Illustration (white noise)

One takes S as a linear space with dimension D, the expected *quadratic risk* of the LSE can be easily computed

$$E\left\|\hat{s}-s\right\|^{2} = d^{2}\left(s,S\right) + \frac{D}{n}$$

Of course, since we do not know $d^2(s,S)$ the quadratic risk cannot be used as a model choice criterion but just as a benchmark.

• First Conclusions

- It is safer to play with several possible models rather than with *a single* one given in advance.
- The notion of expected risk allows to compare the candidates and can serve as a benchmark.
- According to the risk minimization criterion, S is a « good » model *does not mean* that the target s belongs to S.
- Since the minimization of the risk cannot be used as a selection criterion, one needs to introduce some empirical version of it.

Model selection via penalization

Consider some empirical criterion γ_n .

- Framework: Consider some (at most countable) collection of models $(S_m)_{m \in \mathfrak{M}}$. Represent each model S_m by the ERM S_m on S_m .
- Purpose: select the « best » estimator among the collection $(\hat{s}_m)_{m \in \mathfrak{M}}$.
- Procedure: Given some penalty function pen: $\mathfrak{M} \to \mathbb{R}_+$, we take \hat{m} minimizing

 $\gamma_n(\hat{s}_m) + \operatorname{pen}(m)$

over **m** and define

$$\tilde{s} = \hat{s}_{\hat{m}}.$$

• The classical asymptotic approach

Origin: Akaike (log-likelihood), Mallows (least squares)

1 The penalty function is proportional to the number of parameters D_m of the model S_m .

Akaike: D_m / n Mallows' C_p : $2D_m / n$,

where the variance of the errors of the regression framework is assumed to be equal to 1 by the sake of simplicity.

2 The heuristics (Akaike ('73)) leading to the choice of the penalty function D_n / n relies on the assumption: the dimensions and the number of the models are bounded w.r.t. n and n tends to infinity.

BIC (log-likelihood) criterion Schwartz ('78) :

- aims at selecting a « true » model rather than mimicking an oracle

- also asymptotic, with a penalty which is proportional to the number of parameters: $ln(n)D_m / n$

• The non asymptotic approach Barron,Cover ('91) for discrete models, Birgé, Massart ('97) and Barron, Birgé, Massart ('99)) for general models. Differs from the asymptotic approach on the following points

- The number as well as the dimensions of the models may depend on n.
- One can choose a list of models because of its *approximation properties*:

wavelet expansions, trigonometric or piecewise polynomials, artificial neural networks etc

It may perfectly happen that many models of the list have the same dimension and in our view, the « complexity » of the list of models is typically taken into account. Shape of the penalty

 $C_1 \frac{D_m}{n} + C_2 \frac{x_m}{n}$



Data driven penalization

Practical implementation requires some datadriven calibration of the penalty.

« Recipe »

- 1. Compute the ERM \hat{s}_{D} on the union of models with D parameters
- 2. Use theory to guess the shape of the penalty pen(D), typically pen(D)=aD (but aD(2+ln(n/D)) is another possibility)
- Estimate *a* from the data by multiplying by
 2 the smallest value for which the penalized criterion explodes.

Implemented first by Lebarbier ('05) for multiple change points detection

Celeux, Martin, Maugis '07

- Gene expression data: 1020 genes and 20 experiments
- Mixture models $S_{K} = \left\{ x \in \mathbb{R}^{20} \mapsto \sum_{k=1}^{K} p_{k} \Phi(x|\mu_{k}, \Sigma) \right\}$
- Choice of K? Slope heuristics: K=17 BIC: K=17 ICL: K=15



Akaike's heuristics revisited

The main issue is to remove the asymptotic approximation argument in Akaike's heuristics

$$\gamma_{n}(\hat{s}_{D}) = \gamma_{n}(s_{D}) - [\gamma_{n}(s_{D}) - \gamma_{n}(\hat{s}_{D})]$$
variance term
$$\hat{v}_{D}$$
minimizing
$$\gamma_{n}(\hat{s}_{D}) + \text{pen}(D) \quad \text{, is equivalent to}$$
minimizing
$$\gamma_{n}(s_{D}) - \gamma_{n}(s) - \hat{v}_{D} + \text{pen}(D)$$
Fair estimate of $\ell(s, s_{D})$

Ideally: $\operatorname{pen}_{id}(D) = \hat{v}_D + \ell(s_D, \hat{s}_D)$ In order to (approximately) minimize $\ell(s, \hat{s}_D) = \ell(s, s_D) + \ell(s_D, \hat{s}_D)$

The key : Evaluate the excess risks

$$\hat{\boldsymbol{v}}_{D} = \gamma_{n} \left(\boldsymbol{s}_{D} \right) - \gamma_{n} \left(\hat{\boldsymbol{s}}_{D} \right)$$

$$\ell(\boldsymbol{s}_{_{\boldsymbol{D}}}, \hat{\boldsymbol{s}}_{_{\boldsymbol{D}}})$$

This the very point where the various approaches diverge. Akaike's criterion relies on the asymptotic approximation

$$\ell(s_D, \hat{s}_D) \approx \hat{v}_D \approx \frac{D}{2n}$$

The method initiated in Birgé, Massart ('97) relies on upper bounds for the sum of the excess risks which can be written as

$$\hat{\mathbf{v}}_{D} + \ell(\mathbf{s}_{D}, \hat{\mathbf{s}}_{D}) = \left[\overline{\gamma}_{n}(\mathbf{s}_{D}) - \overline{\gamma}_{n}(\hat{\mathbf{s}}_{D})\right]$$

where $\overline{\gamma}_n$ denotes the empirical process

$$\overline{\gamma}_{n}(t) = \gamma_{n}(t) - E[\gamma_{n}(t)]$$

These bounds derive from concentration inequalities for the supremum of the appropriately weighted empirical process = (1) = (1)

$$\frac{\overline{\gamma}_{n}(t) - \overline{\gamma}_{n}(u)}{\omega(t, u)}, t \in S_{D}$$

The prototype being Talagrand's inequality ('96) for empirical processes.

This approach has been fruitfully used in several works. Among others: Baraud ('00) and ('03) for least squares in the regression framework, Castellan ('03) for log-splines density estimation, Patricia Reynaud ('03) for poisson processes, etc...

Pascal Massart Lecture Note: In Mathematics Mathematics Thiosenes reports on here developments to mathematical research and teaching - guilding informally and ana high level. The type of insterial conscient for publication inductor Research menoarcelys 2. Let have on a new field on preventations of a new angle how davated field. 3. Summariset policient intensive courses de topics of current research. Feats that are partial paint for still in demandiment also be providened ecture Notes in The time ness of a manuscript is sometimes more important than its form, which may Concentration transferensiste peloritory of briefly-Deals of the entropy policy can be bound on the holds from prove of a current action. the score section to the publice or the sets with a consultation of your Inequalities project. Addresses are given on the laside back-cover. Manuscripts doubling prepared to colling to Springer-Weike's standard specifications. and Model Selection LCpX tryle Bleatney be found at Rol/Tepson specie/pub/tec/stawimetricy/impro/ (for monographs) and Hav/ftp.doi/rgotdo/pub/tec/late vimatricg//muit//ikitiaummerischopis/tuto faist-Additional 5 technical instructions, if necessary, the aveilable on request tion: inmeso ingencom. LNM 1996 1896 Since the interesting works of Talgarand, concentration, executings have been Ecole d'Eté de Probabilités wagate day backgood back to were block or and an every allowed. spload or rendom sumounationer. They also runn to be external tools to develop of de Saint-Flour XXXIII - 2003 too comptotic theory in protects, exactly config canthal little theorem and large de Actions are intern to play a certifial part in the element block theory. In play wear of accuracy, against the fileway for model to be then by from here, and we consider to a Editor: Jean Picard methodown in soundly whether, sharps paties deterious and risthic of issuing. perdiscover) this experience is the templed of the conversion to 2 Magnet In S. Daw 2007. It to make well-combined and should be weaked to be grand note-STUDIET'S 15835 0/00-6586 ISSN puri Sasa KINN 975-3-546-45497-4 vailable online at pringer Link.com Springer }springercom

Main drawback: typically involve some unkown multiplicative constante which may depend on the unknown distribution (variance of the regression errors, supremum of the density, classification noise etc...). Needs to be calibrated...

Slope heuristics : one looks for some approximation of \hat{v}_D (typically) of the form *aD* with *a* unknown. When *D* is large, $\gamma_n(s_D)$ is almost constant, it suffices to « read » *a* as a slope on the graph of $\gamma_n(\hat{s}_D)$. On chooses the final penalty as

$$pen(D) = 2 \times aD$$

In fact $pen_{min}(D) = \hat{v}_{D}$ is a minimal penalty and the slope heuristics provides a way of approximating it. The factor 2 which is finally used reflects our hope that the excess risks

$$\hat{\boldsymbol{v}}_{_{D}} = \boldsymbol{\gamma}_{_{n}} \left(\boldsymbol{s}_{_{D}} \right) - \boldsymbol{\gamma}_{_{n}} \left(\hat{\boldsymbol{s}}_{_{D}} \right)$$

$$\ell(\mathbf{s}_{_{D}}, \hat{\mathbf{s}}_{_{D}})$$

are of the same order of magnitude. If this the case then

« optimal » penalty=2 * « minimal » penalty

Recent advances

• Justification of the slope heuristics:

Arlot and Massart (JMLR'08) for histograms in the regression framework. Phd of Saumard (2010) regular parametric models. Boucheron and Massart (PTRF'11) for concentration of the empirical excess risk (Wilks phenomenon)

Calibration of regularization

Linear estimators Arlot and Bach (2010)

Lasso type algorithms. Thesis: Connault (2010)

and Meynet (work in progress...)

High dimensional Wilks' phenomenon

Wilks' Theorem: under some proper regularity conditions the log-likelihood $L_n(\theta)$ based on n i.i.d. observations with distribution belonging to a parametric model with D parameters obeys to the following weak convergence result

$$2\left(L_n\left(\hat{\theta}\right) - L_n\left(\theta_0\right)\right) \to \chi^2(D)$$

where $\hat{\theta}$ denotes the MLE and θ_0 is the true value of the parameter.

Question: what's left if we consider possibly irregular empirical risk minimization procedures and let the dimension of the model tend to infinity? Obviously one cannot expect similar asymptotic results. However it is still possible to exhibit some kind of Wilks' phenomenon. Motivation: modification of Akaike's

heuristics for model selection



Data-driven penalties

• A statistical learning framework

We consider the i.i.d. framework where one observes independent copies $\xi_1,...,\xi_n$ of a random variable ξ with distribution P. We have in mind the *regression framework* for which $\xi = (X, Y)$. X is an explanatory variable and Y is the response variable.

Let *s* be some target function to be estimated. For instance, if η denotes the *regression* function $\eta(x) = E[Y|X = x]$ The function of interest *s* may be the regression function η itself. In the *binary classification* case where the response variable takes only the two values 0 and 1, it may be the Bayes classifier

We consider some criterion γ , such that the target function s achieves the minimum of

 $s = \Pi_{\{\eta \ge 1/2\}}$

 $t \to P\gamma(t,.)$

over some set \mathcal{S} . For example

 $\gamma(t,(X,Y)) = (Y - t(X))^{2}$

- with $S = L_2$ leads to the regression function as a minimizer
- with $S = \{t : X \rightarrow \{0,1\}\}$ leads to the Bayes classifier

Introducing the empirical criterion

$$\gamma_n(t) = P_n(\gamma(t,.)) = \frac{1}{n} \sum_{i=1}^n \gamma(t,\xi_i)$$

in order to estimate *s* one considers some subset *S* of *S* (a « model ») and defines the empirical risk minimizer \hat{s} as a minimizer of γ_n over *S*.

This commonly used procedure includes LSE and also MLE for density estimation.

In this presentation we shall assume that $s \in S$

(makes life simpler but not necessary) and also that

 $0 \le \gamma \le 1$ (boundedness is necessary).

Introducing the natural loss function

$$\ell(s,t) = P\gamma(t,.) - P\gamma(s,.)$$

We are dealing with two « dual » estimation errors

• the excess loss

$$\ell(s,\hat{s}) = P\gamma(\hat{s},.) - P\gamma(s,.)$$

• the empirical excess loss

$$\ell_n(s,\hat{s}) = P_n \gamma(s,.) - P_n \gamma(\hat{s},.)$$

Note that for MLE $\gamma(t,.) = -\log t(.)$ and Wilks' theorem provides the asymptotic behavior of $\ell_n(s,\hat{s})$ when S is a regular parametric model.

Crucial issue:

Concentration of the empirical excess loss: connected to empirical processes theory because

$$\ell_n(s,\hat{s}) = \sup_{t \in S} P_n(\gamma(s,.) - \gamma(t,.))$$

Difficult problem: Talagrand's inequality does not make directly the job (the 1/n rate is hard to gain). Let us begin with the related but easier question:

What is the order of magnitude of the excess loss and the empirical excess loss?

Risk bounds for the excess loss

We need to relate the variance of $\gamma(t,.) - \gamma(s,.)$ with the excess loss

$$\ell(s,t) = P(\gamma(t,.) - \gamma(s,.))$$

Introducing some pseudo-metric d such that

$$P(\gamma(t,.)-\gamma(s,.))^2 \leq d^2(s,t)$$

We assume that for some convenient function W $d(s,t) \le w(\sqrt{\ell(s,t)})$

In the regression or the classification case d is simply the L_2 distance and w is either identity for regression or is related to a margin condition for classification. • Tsybakov's margin condition (AOS 2004)

$$\ell(s,t) \geq hd^{2\kappa}(s,t)$$

where $\kappa \ge 1$ and $h \in [0,1]$ with $d^{2}(s,t) = E\left[\left|s(X) - t(X)\right|\right]$

Since for binary classification

$$\ell(s,t) = E\left[\left|2\eta(X)-1\right|\left|s(X)-t(X)\right|\right]$$

this condition is closely related to the behavior of $\eta(X)$ around 1/2. For example margin condition $\kappa = 1$ is achieved whenever

$$\left|2\eta-1\right|\geq h$$

• Heuristics

Let us introduce

$$\overline{\gamma_n}(t) = (P_n - P)\gamma(t,.)$$

Then

$$\ell(s,\hat{s}) + \ell_n(s,\hat{s}) = \overline{\gamma_n}(s) - \overline{\gamma_n}(\hat{s}) \implies \ell(s,\hat{s}) \leq \overline{\gamma_n}(s) - \overline{\gamma_n}(\hat{s})$$

Now the variance of $\overline{\gamma_n}(s) - \overline{\gamma_n}(t)$ is bounded by $d^2(s,t)/n$ hence empirical process theory tells you that the uniform fluctuation of

$$\sqrt{n}\left[\overline{\gamma_n}(s)-\overline{\gamma_n}(t)\right]$$

remains under control in the ball $d(s,t) \le \sigma$

(Massart, Nédélec AOS 2006)

Theorem : Let ϕ , $w \nearrow$ such that $\phi(x)/x, w(x)/x \searrow$ with $\phi(1) \ge 1$ and $w(1) \ge 1$. Assume that $d(s,t) \le w(\sqrt{\ell(s,t)})$

and

$$\mathbb{E}\left[\sup_{t\in S,d(s,t)\leq\sigma}\sqrt{n}\left(\overline{\gamma_n}(s)-\overline{\gamma_n}(t)\right)\right]\leq\phi(\sigma)$$

for every σ such that $\sqrt{n\sigma^2} \ge \phi(\sigma)$. Then, defining \mathcal{E}_* as

$$\sqrt{n\varepsilon_*^2} = \phi \circ w(\varepsilon_*)$$

one has

		(^)	
\mathbb{E}	ℓ	S,S	$ \leq C\mathcal{E}_*^2 $
		\ /	

where C is an absolute constant.

Application to classification

• Tsybakov's framework

Tsybakov's margin condition means that $w(\varepsilon) \approx \varepsilon^{1/\kappa}$ An entropy with bracketing condition implies that one can take $\phi(\sigma) \approx \sigma^{1-\rho}$ and we recover Tsybakov's rate $\varepsilon_*^2 \approx n^{-\kappa/(2\kappa+\rho-1)}$

• VC-classes under margin condition $|2\eta - 1| \ge h$ one has $w(\varepsilon) = \varepsilon / \sqrt{h}$. If S is a VC-class with VC-dimension D

$$\phi(\sigma) \approx \sigma \sqrt{D(1 + \log(1 / \sigma))}$$

so that (whenever $nh^2 \ge D$)

$$\varepsilon_*^2 = \frac{C}{nh} D\left(1 + \log\left(\frac{nh^2}{D}\right)\right)$$

Main points

- Local behavior of the empirical process
- Connection between ℓ and d.



These rates are optimal (minimax)

Concentration of the empirical excess loss

Joint work with Boucheron (PTRF 2011). Since

 $\ell(s,\hat{s}) + \ell_n(s,\hat{s}) = \overline{\gamma_n}(s) - \overline{\gamma_n}(\hat{s})$

the proof of the above Theorem also leads to an upper bound for the empirical excess risk for the same price. In other words $\mathbb{E}\left[\overline{\gamma_n}(s) - \overline{\gamma_n}(\hat{s})\right]$ is also bounded by ε_*^2 up to some absolute constant. Concentration: start from identity

$$\ell_n(s,\hat{s}) = \sup_{t \in S} P_n(\gamma(s,.) - \gamma(t,.))$$

• Concentration tools

Let $\xi'_1, ..., \xi'_n$ be some independent copy of $\xi_1, ..., \xi_n$. Defining $Z'_i = \zeta(\xi_1, ..., \xi_{i-1}, \xi'_i, \xi_{i+1}, ..., \xi_n)$ and setting

$$V^{+} = \mathbb{E}\left[\sum_{i=1}^{n} \left(Z - Z_{i}^{'}\right)_{+}^{2} \left|\xi\right]\right]$$

Efron-Stein's inequality asserts that

 $\operatorname{Var}\left[Z\right] \leq \mathbb{E}\left[V^+\right]$

A Burkholder-type inequality (BBLM, AOP2005) For every $q \ge 2$ such that $|Z|^q$ is integrable, one has

$$\left\| \left(Z - \mathbb{E} \left[Z \right] \right)_{+} \right\|_{q} \leq \sqrt{3q \left\| V^{+} \right\|_{q/2}}$$

Comments

This inequality can be compared to Burkholder's martingale inequality

$$\left\|Z - \mathbb{E}\left[Z\right]\right\|_{q} \leq (q-1)\sqrt{\left\|\langle Z \rangle\right\|_{q/2}}$$

where $\langle Z \rangle = \sum_{k=1}^{n} \left(\mathbb{E} \left[Z | \mathbf{F}_{k} \right] - \mathbb{E} \left[Z | \mathbf{F}_{k-1} \right] \right)^{2}$ denotes the quadratic variation w.r.t. Doob's filtration $\mathbf{F}_{k} = \sigma \left(X_{1}, ..., X_{k} \right), 1 \leq k \leq n$ and \mathbf{F}_{0} trivial σ -field. It can also be compared with Marcinkiewicz Zygmund's inequality which asserts that in the case where $Z = X_{1} + ... + X_{n}$

$$\left|Z - \mathbb{E}\left[Z\right]\right|_{q} \leq \sqrt{q \left\|\left\langle Z\right\rangle\right\|_{q/2}}$$

Note that the constant q-1 in Burkhölder's inequality cannot generally be improved. Our inequality is therefore somewhere « between » Burkholder's and Marcinkiewicz-Zygmund's inequalities. We always get the \sqrt{q} factor instead of q-1 which turns out to be a crucial gain if one wants to derive sub-Gaussian inequalities. The price is to make the substitution $V^+ \longleftrightarrow \langle Z \rangle$

which is absolutely painless in the Marcinkiewicz-Zygmund case. More generally, for the examples that we have in view, V^+ will turn to be a quite manageable quantity and explicit moment inequalities will be obtained by applying iteratively the preceding one.

• Fundamental example

The supremum of an empirical process

 $Z = \sup_{t \in S} \sum_{i=1}^{n} f_t(\xi_i),$

provides an important example, both for theory and applications. Assuming that $\sup_{t \in T} \|f_t\|_{\infty} \leq 1$, Talagrand's inequality (Inventiones 1996) ensures the existence of some absolute positive contants *K* and η , such that

$$\mathbb{P}\left[Z - \mathbb{E}\left[Z\right] \ge x\right] \le K \exp\left(-\eta\left(\frac{x^2}{\mathbb{E}\left[W\right]} \land x\right)\right)$$

where
$$W = \sup_{t \in S} \sum_{i=1}^{n} f_t^2(\xi_i)$$

- Moment inequalities in actionWhy?
 - Example: for empirical processes, one can prove that for some absolute positive constant

$$\left\| Z - \mathbb{E} \left[Z \right] \right\|_{q} \le C \left[\sqrt{q \mathbb{E} \left[W \right]} \lor q \right] \quad \forall q \ge 2$$

Refining Talagrand's inequality

The key: In the process of recovering Talagrand's inequality via the moment method above, we may improve on the variance factor. Indeed, setting $f_t = \gamma(s, .) - \gamma(t, .)$ and

$$Z = \sup_{t \in S} \sum_{i=1}^{n} f_t(\xi_i) = \sum_{i=1}^{n} f_{\hat{s}}(\xi_i) = n\ell_n(s,\hat{s})$$

we see that

$$Z - Z'_{i} \leq f_{\hat{s}}\left(\xi_{i}\right) - f_{\hat{s}}\left(\xi'_{i}\right)$$

and therefore

$$V^{+} = \mathbb{E}\left[\sum_{i=1}^{n} \left(Z - Z_{i}^{'}\right)_{+}^{2} \left|\xi\right] \le 2\sum_{i=1}^{n} \left(Pf_{\hat{s}}^{2} + f_{\hat{s}}^{2}\left(\xi_{i}\right)\right)$$

at this stage instead of using the crude bound

$$\frac{V^+}{n} \le 2 \left(\sup_{t \in S} Pf_t^2 + \sup_{t \in S} P_n f_t^2 \right)$$

we can use the refined bound

$$\frac{V^{+}}{n} \le 2Pf_{\hat{s}}^{2} + 2P_{n}(f_{\hat{s}}^{2}) \le 4Pf_{\hat{s}}^{2} + 2(P_{n} - P)(f_{\hat{s}}^{2})$$

Now the point is that on the one hand

$$P(f_{\hat{s}}^2) \le w^2(\sqrt{\ell(s,s)})$$

and on the other hand we can handle the second term $(P_n - P)(f_{\hat{s}}^2)$ by using some kind of *square root trick*. $(P_n - P)(f_{\hat{s}}^2)$ can indeed shown to behave not worse than

$$(P_n - P)(f_{\hat{s}}) = \ell_n(s, \hat{s}) + \ell(s, \hat{s})$$

So finally it can be proved that

$$\operatorname{Var}\left[Z\right] \leq \mathbb{E}\left[V^+\right] \leq Cnw^2\left(\varepsilon_*\right)$$

and similar results for higher moments.

Illustration 1

In the (bounded) regression case. If we consider the regressogram estimator on some partition with *D* pieces, it can be proved that

$$n \left\| \ell_n(s,\hat{s}) - \mathbb{E} \left[\ell_n(s,\hat{s}) \right] \right\|_q \le C \left[\sqrt{qD} + q \right]$$

In this case $n\mathbb{E}\left[\ell_n(s,\hat{s})\right]$ can be shown to be approximately proportional to D. This exemplifies the high dimensional Wilks phenomenon.

Application to model selection with adaptive penalties: Arlot and Massart, JMLR 2009.

Illustration 2

It can be shown that in the classification case, If S is a VC-class with VC-dimension D, under the margin condition $|2\eta - 1| \ge h$

$$nh\left\|\ell_{n}\left(s,\hat{s}\right)-\mathbb{E}\left[\ell_{n}\left(s,\hat{s}\right)\right]\right\|_{q} \leq C\left[\sqrt{qD\left(1+\log\left(\frac{nh^{2}}{D}\right)\right)}+q\right]$$

provided that $nh^2 \ge D$.

Application to model selection: work in progress with Saumard and Boucheron.

Thanks for your attention!