

Heuristique de pente pour des M-estimateurs à contraste régulier

Journées MAS 2010, Bordeaux

Adrien Saumard

Université Rennes 1, IRMAR

Sept 3, 2010

- 1 The Slope Phenomenon, introduction and first heuristics
- 2 Optimal control of the excess risks when the contrast is "regular", fixed model case

1 - The Slope Phenomenon, introduction and first heuristics

Some general notations

- Unknown law P on a measurable space $(\mathcal{Z}, \mathcal{T})$, generic random variable Z of law P .
- We are given (Z_1, \dots, Z_n) i.i.d. sample of law $P^{\otimes n}$ (also independent of Z).
- Empirical measure associated to the sample

$$P_n = \frac{1}{n} \sum_{i=1}^n \delta_{Z_i}$$

- Expectations :

$$P(s) = \mathbb{E}[s(Z)] ,$$
$$P_n(s) = \frac{1}{n} \sum_{i=1}^n s(Z_i) .$$

- Norms :

$$\|s\|_{2,\mu} = \sqrt{\mu(s^2)} \quad ; \quad \|s\|_2 := \|s\|_{L_2(P)}$$

$$\|s\|_\infty = \operatorname{ess\,sup}_{z \in \mathcal{Z}} |s(z)| .$$

- Positive and negative parts :

$$(x)_+ := \max\{x; 0\} \quad ; \quad (x)_- := \max\{-x; 0\} \geq 0 \quad \forall x \in \mathbb{R}$$

$$(f)_\pm : x \in \mathcal{D}_f \longmapsto (f(x))_\pm .$$

- A functional space : (not a vector space !)

$$L_1^-(P) := \{f : (\mathcal{Z}, \mathcal{T}) \rightarrow \overline{\mathbb{R}}, P(f)_- < +\infty\} (\supset L_1(P)) .$$

Expectation is well-defined on $L_1^-(P)$,

$$Pf := P(f)_+ - P(f)_- \in (-\infty; +\infty] .$$

Definitions (Contrast, Target, Risk)

Given $(\mathcal{Z}, \mathcal{T}, P)$, a **Contrast** is a functional K defined from a set \mathcal{S} of functions to $L_1^-(P)$,

$$K : \begin{cases} \mathcal{S} \longrightarrow L_1^-(P) := \{f : (\mathcal{Z}, \mathcal{T}) \rightarrow \mathbb{R}, P(f)_- < +\infty\} \\ s \longmapsto (Ks : z \longmapsto (Ks)(z)) \end{cases},$$

such that the **risk** function (for any $s \in \mathcal{S}$, $P(Ks)$ is called the **risk** of s)

$$PK : \begin{cases} \mathcal{S} \longrightarrow (-\infty; +\infty] \\ s \longmapsto P(Ks) := \mathbb{E}[(Ks)(Z)] \end{cases}$$

is proper (i.e. not identically equal to $+\infty$) and admits a unique minimum. The argument of this minimum is called the **target**, denoted by s_* .

Definition (M-estimator)

Let $K : \mathcal{S} \rightarrow L_1^-(P)$ be a contrast and let $M \subset \mathcal{S}$ such the restriction of the risk function PK to M is proper. M is called a **model**. We call **M-estimator** associated to the contrast K and to the model M , a random variable $s_n(M)$ satisfying

$$s_n(M) \in \arg \min_{s \in M} P_n(Ks) \quad , \quad |P_n(Ks_n(M))| < +\infty \quad a.s.$$

- Quality of a M-estimator : measured by its **excess risk**,

$$\ell(s_*, s_n(M)) := P(Ks_n(M)) - P(Ks_*) = P(Ks_n(M) - Ks_*) \geq 0 .$$

- Maximum likelihood estimation of density (MLE) :

$$s_* = \frac{dP}{d\mu} ; K(s) = -\ln s$$

Excess risk : **Kullback-Leibler divergence** of s w.r.t. s_* .

$$\ell(s_*, s) = \mathcal{K}(s_*, s) = \int_{\mathcal{Z}} s_* \ln \left(\frac{s_*}{s} \right) d\mu$$

- Maximum likelihood estimation of density (MLE) :

$$s_* = \frac{dP}{d\mu} ; K(s) = -\ln s$$

Excess risk : **Kullback-Leibler divergence** of s w.r.t. s_* .

$$\ell(s_*, s) = \mathcal{K}(s_*, s) = \int_{\mathcal{Z}} s_* \ln \left(\frac{s_*}{s} \right) d\mu$$

- Least-square estimation of density (LSE):

$$s_* = \frac{dP}{d\mu}.$$

Goal : $\ell(s_*, s) = P(Ks - Ks_*) = \|s - s_*\|_{2,\mu}^2$ (excess risk given by the **quadratic norm** of $L_2(\mu)$).

Contrast : $K(s) = \|s\|_{2,\mu}^2 - 2Ps$.

- Least-squares heteroscedastic Regression :

$$Y = s_*(X) + \sigma(X)\varepsilon, \quad \mathbb{E}[\varepsilon | X] = 0 \quad \text{et} \quad \mathbb{E}[\varepsilon^2 | X] = 1 .$$

If $\mathcal{S} = L_2(P^X)$ and $K : \mathcal{S} \mapsto L_1(P) (\subset L_1^-(P))$, with

$$Ks : z = (x, y) \mapsto (Ks)(z) = (Ks)(x, y) = (y - s(x))^2 .$$

$$\ell(s_*, s) = P(Ks - Ks_*) = P^X (s - s_*)^2 = \|s - s_*\|_2^2 .$$

Excess risk : **quadratic norm** of $L_2(P^X)$.

- Binary Classification : $Z = (X, Y) \in \mathcal{X} \times \{-1, 1\}$.

Target : **Bayes classifier**.

$$s_* = \arg \min_{s \in \mathcal{S}} P(Ks) = \arg \min_{s \in \mathcal{S}} P(Y \neq s(X)) .$$

Contrast : $K(s)(x, y) = \mathbf{1}_{y \cdot s(x) \geq 0}$.

- Least-squares **heteroscedastic Regression** :

$$Y = s_*(X) + \sigma(X)\varepsilon, \quad \mathbb{E}[\varepsilon | X] = 0 \quad \text{et} \quad \mathbb{E}[\varepsilon^2 | X] = 1 .$$

If $\mathcal{S} = L_2(P^X)$ and $K : \mathcal{S} \mapsto L_1(P) (\subset L_1^-(P))$, with

$$Ks : z = (x, y) \mapsto (Ks)(z) = (Ks)(x, y) = (y - s(x))^2 .$$

$$\ell(s_*, s) = P(Ks - Ks_*) = P^X (s - s_*)^2 = \|s - s_*\|_2^2 .$$

Excess risk : **quadratic norm** of $L_2(P^X)$.

- Binary Classification : $Z = (X, Y) \in \mathcal{X} \times \{-1, 1\}$.

Target : **Bayes classifier**.

$$s_* = \arg \min_{s \in \mathcal{S}} P(Ks) = \arg \min_{s \in \mathcal{S}} P(Y \neq s(X)) .$$

Contrast : $K(s)(x, y) = \mathbf{1}_{y \cdot s(x) \geq 0}$.

- Convex binary classification : SVM, Boosting, Logistic regression, etc...

Contrasts : **convex surrogate** of the Bayes contrast.

$$K_\phi : s \mapsto [K_\phi(s) : z = (x, y) \mapsto K_\phi(s)(z) = \phi(y \cdot s(x))] .$$

Model Selection in M-estimation, via penalization.

- Contrast : $K : \mathcal{S} \rightarrow L_1^-(P)$, target $s_* = \arg \min_{s \in \mathcal{S}} P(Ks)$.

Model Selection in M-estimation, via penalization.

- Contrast : $K : \mathcal{S} \rightarrow L_1^-(P)$, target $s_* = \arg \min_{s \in \mathcal{S}} P(Ks)$.
- Collection of models : \mathcal{M}_n . Associated collection of M-estimators : $\{s_n(M) ; M \in \mathcal{M}_n\}$,

$$s_n(M) \in \arg \min_{s \in M} P_n(Ks) , \quad \forall M \in \mathcal{M}_n .$$

Model Selection in M-estimation, via penalization.

- Contrast : $K : \mathcal{S} \rightarrow L_1^-(P)$, target $s_* = \arg \min_{s \in \mathcal{S}} P(Ks)$.
- Collection of models : \mathcal{M}_n . Associated collection of M-estimators : $\{s_n(M) ; M \in \mathcal{M}_n\}$,

$$s_n(M) \in \arg \min_{s \in M} P_n(Ks) , \quad \forall M \in \mathcal{M}_n .$$

- Oracle model (target of the model selection procedure) :

$$\begin{aligned} M_* &\in \arg \min_{M \in \mathcal{M}_n} P(Ks_n(M)) \\ &= \arg \min_{M \in \mathcal{M}_n} P(Ks_n(M) - Ks_*) \\ &= \arg \min_{M \in \mathcal{M}_n} \{P_n(Ks_n(M) - Ks_*) + (P - P_n)(Ks_n(M) - Ks_*)\} . \end{aligned}$$

Model Selection in M-estimation, via penalization.

- Contrast : $K : \mathcal{S} \rightarrow L_1^-(P)$, target $s_* = \arg \min_{s \in \mathcal{S}} P(Ks)$.
- Collection of models : \mathcal{M}_n . Associated collection of M-estimators : $\{s_n(M) ; M \in \mathcal{M}_n\}$,

$$s_n(M) \in \arg \min_{s \in \mathcal{M}} P_n(Ks) , \quad \forall M \in \mathcal{M}_n .$$

- Oracle model (target of the model selection procedure) :

$$\begin{aligned} M_* &\in \arg \min_{M \in \mathcal{M}_n} P(Ks_n(M)) \\ &= \arg \min_{M \in \mathcal{M}_n} P(Ks_n(M) - Ks_*) \\ &= \arg \min_{M \in \mathcal{M}_n} \{P_n(Ks_n(M) - Ks_*) + (P - P_n)(Ks_n(M) - Ks_*)\} . \end{aligned}$$

- If $\text{Card}(\mathcal{M}_n) \leq c_{\mathcal{M}} n^{\alpha_{\mathcal{M}}}$, then we set the **ideal penalty** (S. Arlot, PhD Thesis, 2007),

$$\text{pen}_{\text{id}} : M \in \mathcal{M}_n \mapsto \text{pen}_{\text{id}}(M) = (P - P_n)(Ks_n(M) - Ks_*) \geq 0 .$$

Hence,

$$M_* \in \arg \min_{M \in \mathcal{M}_n} \{P_n(K_{S_n}(M) - K_{S_*}) + \text{pen}_{\text{id}}(M)\} .$$

- Selected model : Choose $\text{pen} : M \in \mathcal{M}_n \mapsto \text{pen}(M) \geq 0$ and select

$$\begin{aligned} \hat{M} &\in \arg \min_{M \in \mathcal{M}_n} \{P_n(K_{S_n}(M)) + \text{pen}(M)\} \\ &= \arg \min_{M \in \mathcal{M}_n} \{P_n(K_{S_n}(M) - K_{S_*}) + \text{pen}(M)\} . \end{aligned}$$

Hence,

$$M_* \in \arg \min_{M \in \mathcal{M}_n} \{P_n(Ks_n(M) - Ks_*) + \text{pen}_{\text{id}}(M)\} .$$

- Selected model : Choose $\text{pen} : M \in \mathcal{M}_n \mapsto \text{pen}(M) \geq 0$ and select

$$\begin{aligned} \hat{M} &\in \arg \min_{M \in \mathcal{M}_n} \{P_n(Ks_n(M)) + \text{pen}(M)\} \\ &= \arg \min_{M \in \mathcal{M}_n} \{P_n(Ks_n(M) - Ks_*) + \text{pen}(M)\} . \end{aligned}$$

- Quality of the procedure : measured by an **oracle inequality**. With large probability (of order $1 - Ln^{-2}$),

$$\ell\left(s_*, s_n\left(\hat{M}\right)\right) \leq C \times \ell\left(s_*, s_n\left(M_*\right)\right) .$$

The smaller is $C \geq 1$ (under a fixed probability), the better is the model selection procedure in terms of **prediction** (measured by the excess risk).

Optimal Model Selection, Slope Heuristics

- A model selection procedure is **optimal** - or nearly optimal - if, with probability at least $1 - Ln^{-2}$, we have

$$\ell \left(s_*, s_n \left(\widehat{M} \right) \right) \leq (1 + \varepsilon_n) \times \ell \left(s_*, s_n \left(M_* \right) \right) , \quad \varepsilon_n \rightarrow 0 .$$

Slope Heuristics : (Birgé & Massart, 2007, extended by Arlot & Massart, 2009) There exists a penalty, called **minimal penalty** and denoted pen_{\min} , such that:

(I) If a penalty $\text{pen} : \mathcal{M}_n \longrightarrow \mathbb{R}_+$ is such that, for all model $M \in \mathcal{M}_n$,

$$\text{pen}(M) \leq (1 - \delta) \text{pen}_{\min}$$

with $\delta > 0$, then the dimension of the selected model \widehat{M} is “very large” and the excess risk of the selected estimator $s_n \left(\widehat{M} \right)$ is “much larger” than the excess risk of the oracle.

(II) If $\text{pen} \approx (1 + \delta) \text{pen}_{\min}$ with $\delta > 0$, then the corresponding model selection procedure satisfies an oracle inequality with a leading constant $C(\delta) < +\infty$ and the dimension of the selected model is “not too large”.

(III) Moreover,

$$\text{pen}_{\text{opt}} \approx 2 \text{pen}_{\min}$$

is a (quasi)-optimal penalty.

If the **projection** s_M of the target s_* exists and is unique, i.e.

$$s_M = \arg \min_{s \in M} P(Ks), \quad Ks_M \in L_1(P),$$

then

$$P_n(Ks_M - Ks_n(M)) \geq 0$$

is called the **empirical excess risk on M** . In this case, Arlot & Massart (09) conjecture that the following equality holds with great generality,

$$\text{for all } M \in \mathcal{M}_n, \quad \text{pen}_{\min}(M) = \mathbb{E}[P_n(Ks_M - Ks_n(M))] .$$

- **Question** : In what extend the conjecture of Arlot & Massart is true ? Find a general positive answer, find (nontrivial) counter-examples...

Practice ? Listen to the talk of J-P. Baudry

- Baudry, Maugis & Michel (10) : Survey. Overview of the theoretical and practical results about the Slope heuristics. Logiciel CAPUSHE. **Already many conclusive empirical study (simulations and real data) !**
- When it is possible, use the Slope Heuristics to calibrate your penalty in practice, it seems to work quite well !...

One Heuristic on the slope phenomenon

If we take $\text{pen} \approx 2 \times \text{pen}_{\min} = 2\mathbb{E}[P_n(Ks_M - Ks_n(M))]$. \hat{M} minimizes

$$\begin{aligned} & P_n(Ks_n(M) - Ks_*) + \text{pen}(M) \\ & \approx \ell(s_*, s_M) + P_n(Ks_M - Ks_n(M)) + (P_n - P)(Ks_M - Ks_*) \\ & \quad + 2(\underbrace{\mathbb{E}[P_n(Ks_M - Ks_n(M))] - P_n(Ks_M - Ks_n(M))}_{\text{Boucheron \& Massart, 2010.}}) \\ & \approx \ell(s_*, s_M) + P_n(Ks_M - Ks_n(M)) . \end{aligned}$$

If

$$\boxed{P_n(Ks_M - Ks_n(M)) \sim P(Ks_n(M) - Ks_M)} \quad (*)$$

then

$$\begin{aligned} P(Ks_M - Ks_*) + P_n(Ks_M - Ks_n(M)) &\approx \ell(s_*, s_M) + P(Ks_n(M) - Ks_M) \\ &\approx \ell(s_*, s_n(M)) . \end{aligned}$$

Hence,

$$\ell(s_*, s_n(\widehat{M})) \approx \ell(s_*, s_n(M_*))$$

and the procedure is nearly optimal.

The **keystone of the slope heuristics** is the equivalence (*) with high probability between the true and empirical excess risk, for the model of interest.

2 - Optimal control of the excess risks when the contrast is "regular", fixed model case

The notion of regular contrast

Definition (Regular Contrast w.r.t. a model)

Let $K : \mathcal{S} \rightarrow L_1^-(P)$ be a contrast with \mathcal{S} . Take $M \subset \mathcal{S}$ a convex model. Then K is said to be regular w.r.t. M if there exists a projection s_M of the target s_* onto M ,

$$s_M \in \arg \min_{s \in M} P(Ks),$$

if the restriction $PK|_M : M \rightarrow (-\infty, +\infty]$ is strictly convex and if there exists $c > 0$ such that, by denoting

$$B_c := \{s \in \text{Aff}(M) ; \|s - s_M\|_\infty < c\}$$

we have

$$B_c \subset M$$

and the restriction $K|_{B_c} : B_c \rightarrow L_\infty(P)$ is \mathcal{C}^3 in the sense of the Fréchet-differentiability.

Theorem

Let $\alpha, A_-, A_+, A_H, A_{\text{cons}} > 0$ and let $K : S \rightarrow L_1^-(P)$, be a regular contrast w.r.t. a model M . Denote by M_0 the underlying vector space of $\text{Aff}(M)$. Assume that

$$0 < A_- (\ln n)^2 \leq \dim(M_0) = D \leq A_+ \frac{n}{(\ln n)^2} < +\infty$$

and that there exists a positive constant $A_H > 0$ such that

$$\text{for all } s \in B_c, \quad \|s - s_M\|_2 \leq A_H P(K''(s_M)(s - s_M, s - s_M)) .$$

Hence the norm defined by

$$\|h\|_{H,M} = \sqrt{P(K''(s_M)(h, h))}, \quad h \in M_0,$$

is an Hilbertian norm on M_0 .

Theorem

Moreover, assume that there exists $R_{n,D,\alpha} \leq A_{\text{cons}} (\ln n)^{-1/2}$ such that for all $n \geq n_1$,

$$\mathbb{P} [\|s_n(M) - s_M\|_\infty > R_{n,D,\alpha}] \leq n^{-\alpha} .$$

Finally, assume that $(M_0, \|\cdot\|_{H,M})$ has a localized basis structure : there exists an orthonormal basis $\varphi = (\varphi_k)_{k=1}^D$ in $(M_0, \|\cdot\|_{H,M})$ that satisfies, for a positive constant $r_M(\varphi)$ and all $\beta = (\beta_k)_{k=1}^D \in \mathbb{R}^D$,

$$\left\| \sum_{k=1}^D \beta_k \varphi_k \right\|_\infty \leq r_M(\varphi) \sqrt{D} |\beta|_\infty ,$$

where $|\beta|_\infty = \max \{|\beta_k| ; k \in \{1, \dots, D\}\}$ is the sup-norm of the D -dimensional vector β .

Theorem

Then there exists $A_0 > 0$ and a positive number n_0 depending on the constants of the problem such that by setting

$$\varepsilon_n = A_0 \max \left\{ \left(\frac{\ln n}{D} \right)^{1/4}, \left(\frac{D \ln n}{n} \right)^{1/4}, \sqrt{R_{n,D,\alpha}} \right\},$$

we have for all $n \geq n_0$, with probability at least $1 - 15n^{-\alpha}$,

$$(1 - \varepsilon_n) \frac{1}{4} \frac{D}{n} \mathcal{K}_{1,M}^2 \leq P(K_{S_n}(M) - K_{S_M}) \leq (1 + \varepsilon_n) \frac{1}{4} \frac{D}{n} \mathcal{K}_{1,M}^2,$$

$$(1 - \varepsilon_n^2) \frac{1}{4} \frac{D}{n} \mathcal{K}_{1,M}^2 \leq P_n(K_{S_M}(M) - K_{S_n}) \leq (1 + \varepsilon_n^2) \frac{1}{4} \frac{D}{n} \mathcal{K}_{1,M}^2,$$

where $\mathcal{K}_{1,M}^2 = D^{-1} \sum_{k=1}^D \text{Var}(K'(s_M)(\varphi_k))$.

- **Conclusion :** In this case, we have proved

$$P(K_{S_n}(M) - K_{S_M}) \sim P_n(K_{S_M} - K_{S_n}(M)).$$

Preprints:

- A.S., Optimal upper and lower bounds for the true and empirical excess risks in heteroscedastic least-squares regression, 2010, hal-00512304, v1.
- A.S., The Slope Heuristics in Heteroscedastic Regression, 2010, hal-00512306, v1.
- A.S., Nonasymptotic quasi-optimality of AIC and the slope heuristics in maximum likelihood estimation of density using histogram models, 2010, hal-00512310, v1.

In preparation:

- Convergence in sup-norm of the least-squares estimator of a regression function with heteroscedastic noise.
- Regular Contrast Estimation on a fixed convex model.
- Regular Contrast Estimation and the Slope Heuristics.