

Convex methods for classification

Peter Bartlett

Department of Statistics and Computer Science Division
UC Berkeley

Joint work with

Sylvain Arlot, Mike Jordan, Jon McAuliffe, Mikhail Traskin.

slides at <http://www.stat.berkeley.edu/~bartlett>

The Pattern Classification Problem

- i.i.d. $(X, Y), (X_1, Y_1), \dots, (X_n, Y_n)$ from $\mathcal{X} \times \{\pm 1\}$.
- Use data $(X_1, Y_1), \dots, (X_n, Y_n)$ to choose $f_n : \mathcal{X} \rightarrow \mathbb{R}$ with small risk,

$$R(f) = \Pr(\text{sign}(f(X)) \neq Y) = \mathbf{E}\ell_f(X, Y),$$

where ℓ_f is the 0-1 loss:

$$\ell_f(x, y) = \begin{cases} 1 & \text{if } y \neq \text{sign}(f(x)), \\ 0 & \text{otherwise.} \end{cases}$$

The Pattern Classification Problem

- Natural approach: minimize empirical risk,

$$\hat{R}(f) = \hat{\mathbf{E}}\ell_f = \frac{1}{n} \sum_{i=1}^n \ell_f(X_i, Y_i).$$

- Often computationally intractable...
- An alternative approach:
Replace 0-1 loss, ℓ , with a convex surrogate, ϕ .

Large Margin Algorithms

- Consider the margins, $Y f(X)$.
- Define a margin cost function $\phi : \mathbb{R} \rightarrow \mathbb{R}^+$.
- Define the ϕ -risk of $f : \mathcal{X} \rightarrow \mathbb{R}$ as $R_\phi(f) = \mathbf{E}\phi(Y f(X))$.
- Choose $f \in \mathcal{F}$ to minimize ϕ -risk.
(e.g., use data, $(X_1, Y_1), \dots, (X_n, Y_n)$, to minimize **empirical ϕ -risk**,

$$\hat{R}_\phi(f) = \hat{\mathbf{E}}\phi(Y f(X)) = \frac{1}{n} \sum_{i=1}^n \phi(Y_i f(X_i)),$$

or a regularized version.)

Large Margin Algorithms

- Adaboost:
 - $\mathcal{F} = \text{span}(\mathcal{G})$ for a VC-class \mathcal{G} ,
 - $\phi(\alpha) = \exp(-\alpha)$,
 - Minimizes $\hat{R}_\phi(f)$ using greedy basis selection, line search:

$$f_{t+1} = f_t + \alpha_{t+1}g_{t+1},$$

$$\hat{R}_\phi(f_t + \alpha_{t+1}g_{t+1}) = \min_{\alpha \in \mathbb{R}, g \in \mathcal{G}} \hat{R}_\phi(f_t + \alpha g).$$

Large Margin Algorithms

- Support vector machines:

- \mathcal{F} = ball in reproducing kernel Hilbert space, \mathcal{H} .
- $\phi(\alpha) = \max(0, 1 - \alpha)$.
- Algorithm minimizes $\hat{R}_\phi(f) + \lambda \|f\|_{\mathcal{H}}^2$.

This is equivalent to a quadratic program:

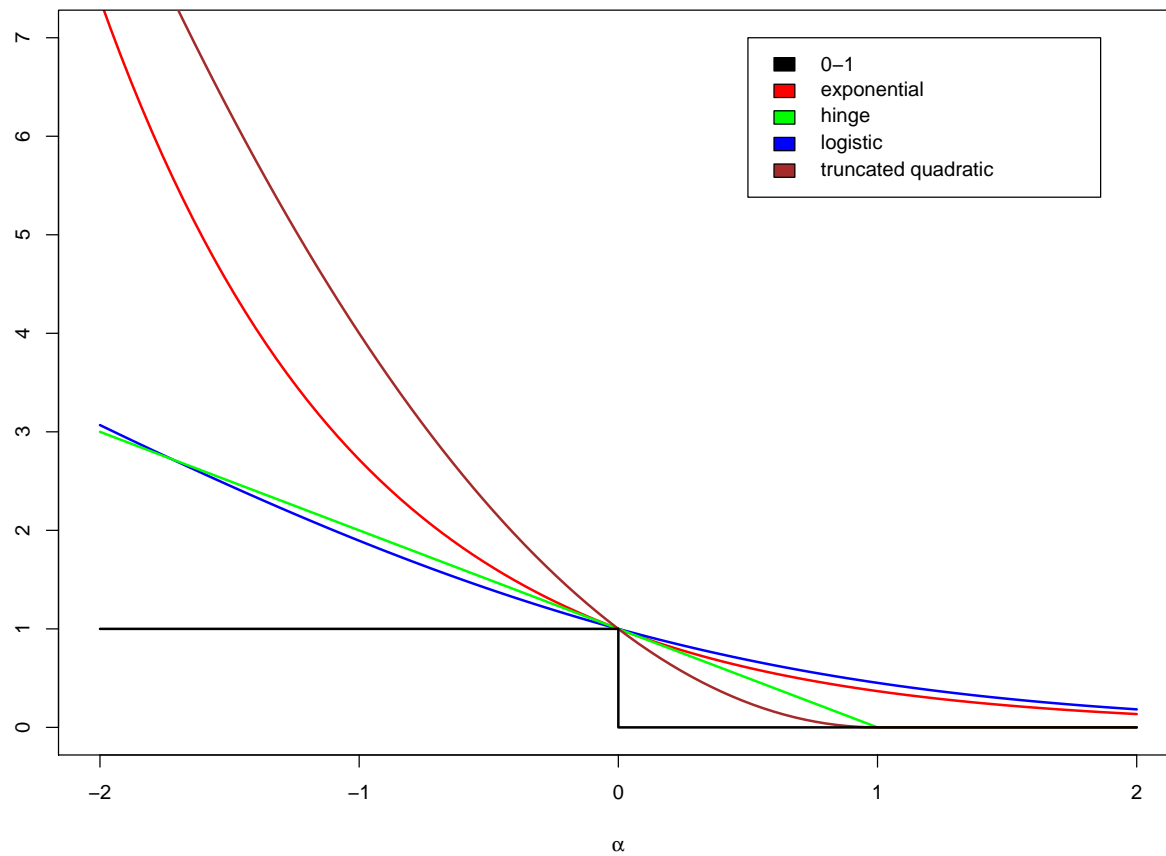
$$\begin{aligned} \min \quad & \xi' \mathbf{1} + \lambda \alpha' K \alpha \\ \text{s.t.} \quad & 1 - \xi \leq \text{diag}(\mathbf{y}) K \alpha, \\ & \xi \geq 0, \end{aligned}$$

where $\mathbf{y} = (Y_1, \dots, Y_n)$,
 $K_{i,j} = k(X_i, X_j)$,
 $\hat{f}(x) = \sum_{i=1}^n \alpha_i k(X_i, x)$,
and k is the reproducing kernel of \mathcal{H} .

Large Margin Algorithms

- Many other variants
 - Neural net classifiers
 - L2Boost, LS-SVMs
 - Logistic regression

Large Margin Algorithms



Overview

1. Classification with convex loss.
2. Universal consistency of large margin algorithms.
3. Classification problems with low noise.

Overview

1. Classification with convex loss.
 - Impact of replacing 0-1 loss with convex loss?
 - What ϕ are suitable for classification?
 - Relationship between risk and ϕ -risk?
2. Universal consistency of large margin algorithms.
3. Classification problems with low noise.

Definitions and Facts

$$R(f) = \Pr(\text{sign}(f(X)) \neq Y)$$

$$R^* = \inf_f R(f)$$

risk

$$R_\phi(f) = \mathbb{E}\phi(Y f(X))$$

$$R_\phi^* = \inf_f R_\phi(f)$$

ϕ -risk

$$\eta(x) = \Pr(Y = 1|X = x)$$

conditional probability

$$f^*(x) = \text{sign}(2\eta(x) - 1)$$

Bayes decision rule.

Notice: $R_\phi(f) = \mathbb{E}(\mathbb{E}[\phi(Y f(X))|X])$, and **conditional ϕ -risk** is:

$$\mathbb{E}[\phi(Y f(X))|X = x] = \eta(x)\phi(f(x)) + (1 - \eta(x))\phi(-f(x)).$$

Definitions

$$H(\eta) = \inf_{\alpha \in \mathbb{R}} (\eta\phi(\alpha) + (1 - \eta)\phi(-\alpha))$$

$$H^-(\eta) = \inf_{\alpha: \alpha(2\eta-1) \leq 0} (\eta\phi(\alpha) + (1 - \eta)\phi(-\alpha)).$$

Definition: We say that ϕ is **classification-calibrated** if, for $\eta \neq 1/2$,

$$H^-(\eta) > H(\eta).$$

i.e., pointwise optimization of conditional ϕ -risk leads to the correct sign.

The ψ transform

Definition: Given convex ϕ , define $\psi : [0, 1] \rightarrow [0, \infty)$ by

$$\psi(\theta) = H^{-} \left(\frac{1 + \theta}{2} \right) - H \left(\frac{1 + \theta}{2} \right).$$

(The definition is a little more involved for non-convex ϕ .)

The Relationship between Excess Risk and Excess ϕ -risk

Theorem: [with Mike Jordan and Jon McAuliffe]

1. For any P and f , $\psi(R(f) - R^*) \leq R_\phi(f) - R_\phi^*$.
2. For $|\mathcal{X}| \geq 2$, $\epsilon > 0$ and $\theta \in [0, 1]$, there is a P and an f with

$$R(f) - R^* = \theta$$

$$\psi(\theta) \leq R_\phi(f) - R_\phi^* \leq \psi(\theta) + \epsilon.$$

3. The following conditions are equivalent:
 - (a) ϕ is classification calibrated.
 - (b) $\psi(\theta_i) \rightarrow 0$ iff $\theta_i \rightarrow 0$.
 - (c) $R_\phi(f_i) \rightarrow R_\phi^*$ implies $R(f_i) \rightarrow R^*$.

Classification-calibrated ϕ

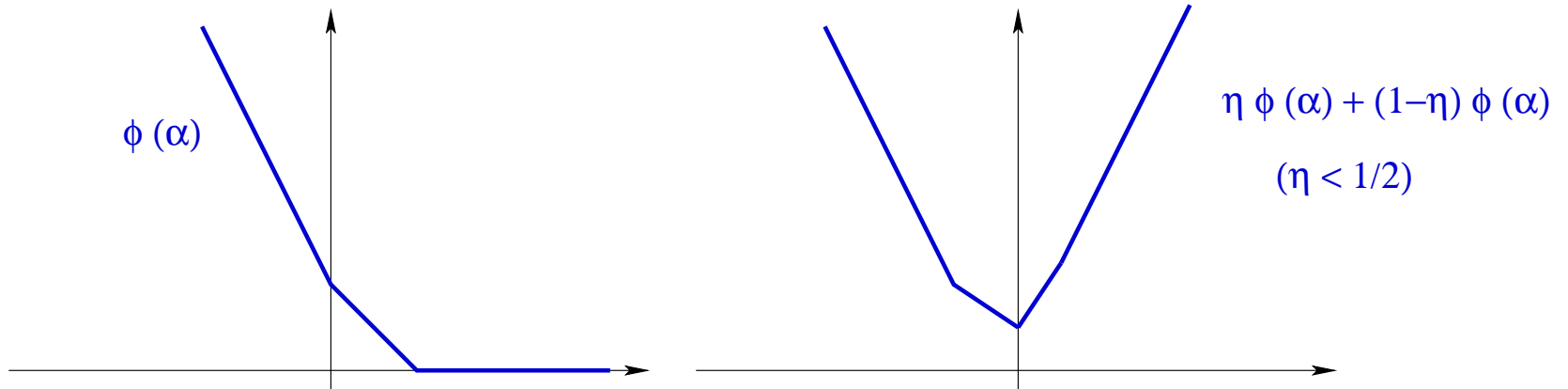
Theorem: If ϕ is convex,

$$\phi \text{ is classification calibrated} \Leftrightarrow \begin{cases} \phi \text{ is differentiable at } 0 \\ \phi'(0) < 0. \end{cases}$$

Classification-calibrated ϕ

Theorem: If ϕ is convex,

$$\phi \text{ is classification calibrated} \Leftrightarrow \begin{cases} \phi \text{ is differentiable at } 0 \\ \phi'(0) < 0. \end{cases}$$



Classification with convex loss

Bartlett, Jordan and McAuliffe, *Convexity, classification, and risk bounds*.

See also:

Zhang, Statistical behavior and consistency of classification methods based on convex risk minimization.

Steinwart, How to compare different loss functions and their risks.

Overview

1. Classification with convex loss.
2. Universal consistency of large margin algorithms.
 - AdaBoost.
3. Classification problems with low noise.

Universal Consistency

- Assume: **i.i.d. data**, $(X, Y), (X_1, Y_1), \dots, (X_n, Y_n)$ from $\mathcal{X} \times \mathcal{Y}$ (with $\mathcal{Y} = \{\pm 1\}$).
- Consider a method $f_n = A((X_1, Y_1), \dots, (X_n, Y_n))$,
e.g., $f_n = \text{AdaBoost}((X_1, Y_1), \dots, (X_n, Y_n), t_n)$.

Definition: We say that the method is **universally consistent** if, for all distributions P ,

$$R(f_n) \xrightarrow{a.s.} R^*,$$

where R is the risk and R^* is the Bayes risk:

$$R(f) = \Pr(Y \neq \text{sign}(f(X))), \quad R^* = \inf_f R(f).$$

The Approximation/Estimation Decomposition

Consider an algorithm that chooses

$$f_n = \arg \min_{f \in \mathcal{F}_n} \hat{R}_\phi(f) \quad \text{or} \quad f_n = \arg \min_{f \in \mathcal{F}} \left(\hat{R}_\phi(f) + \lambda_n \Omega(f) \right).$$

($\hat{R}_\phi(f)$ is empirical ϕ -risk, $\mathcal{F}_1 \subseteq \mathcal{F}_2 \subseteq \dots \subseteq \mathcal{F}$, and Ω is regularization.)

We can decompose the excess risk estimate as

$$\begin{aligned} \psi(R(f_n) - R^*) &\leq R_\phi(f_n) - R_\phi^* \\ &= \underbrace{R_\phi(f_n) - \inf_{f \in \mathcal{F}_n} R_\phi(f)}_{\text{estimation error}} + \underbrace{\inf_{f \in \mathcal{F}_n} R_\phi(f) - R_\phi^*}_{\text{approximation error}}. \end{aligned}$$

The Approximation/Estimation Decomposition

$$\begin{aligned} \psi(R(f_n) - R^*) &\leq R_\phi(f_n) - R_\phi^* \\ &= \underbrace{R_\phi(f_n) - \inf_{f \in \mathcal{F}_n} R_\phi(f)}_{\text{estimation error}} + \underbrace{\inf_{f \in \mathcal{F}_n} R_\phi(f) - R_\phi^*}_{\text{approximation error}}. \end{aligned}$$

- Approximation and estimation errors are in terms of R_ϕ , not R .
- Like a regression problem.
- With a rich class and appropriate regularization, $R_\phi(f_n) \rightarrow R_\phi^*$.
(e.g., \mathcal{F}_n gets large slowly, or $\lambda_n \rightarrow 0$ slowly.)
- Universal consistency ($R(f_n) \rightarrow R^*$) iff ϕ is classification calibrated.

Example: Universal Consistency of SVMs

For a Reproducing Kernel Hilbert Space \mathcal{H} , choose

$$f_n = \arg \min_{f \in \mathcal{H}} \left(\hat{R}_\phi(f) + \lambda_n \|f\|_{\mathcal{H}}^2 \right),$$

or $f_n = \arg \min_{f \in \mathcal{H}_n} \hat{R}_\phi(f)$ with $\mathcal{H}_n = \{f \in \mathcal{H} : \lambda_n \|f\|_{\mathcal{H}}^2 \leq 1\}$.

$$\text{Then } \psi(R(f_n) - R^*) \leq \underbrace{R_\phi(f_n) - \inf_{f \in \mathcal{H}_n} R_\phi(f)}_{\text{estimation error}} + \underbrace{\inf_{f \in \mathcal{H}_n} R_\phi(f) - R_\phi^*}_{\text{approximation error}}.$$

If \mathcal{H} is large (e.g., a Gaussian kernel on \mathbb{R}^d), $\inf_{f \in \mathcal{H}_n} R_\phi(f) \rightarrow R_\phi^*$.

For $\lambda_n \rightarrow 0$ (suitably slowly), $|\hat{R}_\phi(f_n) - R_\phi(f_n)| \xrightarrow{a.s.} 0$.

In that case, $R_\phi(f_n) \xrightarrow{a.s.} R_\phi^*$, and universal consistency follows.

(Steinwart, 2005)

Universal Consistency: AdaBoost?

- For SVMs, the regularization term keeps f_n small, which is essential for the uniform convergence result: $|\hat{R}_\phi(f_n) - R_\phi(f_n)| \xrightarrow{a.s.} 0$.
- AdaBoost?

AdaBoost

Sample, $S_n = ((x_1, y_1), \dots, (x_n, y_n)) \in (X \times \{\pm 1\})^n$

Number of iterations, T

Class of basis functions, \mathcal{G}

function AdaBoost(S_n, T):

$f_0 := 0$

for t from $1, \dots, T$

$$(\alpha_t, g_t) := \arg \min_{\alpha \in \mathbb{R}, g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \exp(-y_i (f_{t-1}(x_i) + \alpha g(x_i)))$$

$f_t := f_{t-1} + \alpha_t g_t$

return f_T

Previous results: Regularized versions

Instead, we could consider a regularized version of AdaBoost:

1. Minimize $\hat{R}_\phi(f)$ over $\mathcal{F}_n = \gamma_n \text{co}(\mathcal{G})$, the scaled convex hull of \mathcal{G} .
2. Minimize

$$\hat{R}_\phi(f) + \lambda_n \|f\|_*,$$

over $\text{span}(\mathcal{G})$, where $\|f\|_* = \inf\{\gamma : f \in \gamma \text{co}(\mathcal{G})\}$.

For suitable choices of the parameters (γ_n and λ_n), these algorithms are universally consistent. (Lugosi and Vayatis, 2004), (Zhang, 2004)

Also **bounded step size**. (Zhang and Yu, 2005), (Bickel, Ritov, Zakai, 2006)

Previous results: ‘Process consistency’

Theorem: [Jiang, 2004]

For a (suitable) basis class defined on \mathbb{R}^d , and for all probability distributions P satisfying certain smoothness assumptions, there is a sequence t_n such that $f_n = \text{AdaBoost}(S_n, t_n)$ satisfies

$$R(f_n) \xrightarrow{a.s.} R^*.$$

Universal consistency of AdaBoost

Theorem: [with Mikhail Traskin]

If

$$d_{VC}(\mathcal{G}) < \infty,$$
$$R_{\phi}^* = \liminf_{\lambda \rightarrow \infty} \{R_{\phi}(f) : f \in \lambda \text{co}(\mathcal{G})\},$$
$$t_n \rightarrow \infty$$
$$t_n = O(n^{1-\alpha}) \quad \text{for some } \alpha > 0,$$

then AdaBoost is universally consistent.

Universal consistency of AdaBoost

Theorem:

If

$$d_{VC}(\mathcal{G}) < \infty,$$

$$R_{\phi}^* = \liminf_{\lambda \rightarrow \infty} \{R_{\phi}(f) : f \in \lambda \text{co}(\mathcal{G})\},$$

$$t_n \rightarrow \infty$$

$$t_n = O(n^{1-\alpha}) \quad \text{for some } \alpha > 0,$$

then AdaBoost is universally consistent.

Idea of proof:

Uniform convergence of clipped t_n -combinations. Clipping does not greatly increase \hat{R}_{ϕ} . Then $\hat{R}_{\phi}(f_{t_n})$ approaches best in an ℓ_* -ball. Then uniform convergence over ℓ_* -balls.

Overview

1. Classification with convex loss.
2. Universal consistency of large margin algorithms.
3. Classification problems with low noise.

Low Noise

The difficulty of a binary classification problem is determined by the probability that $\eta(X) = \Pr(Y = 1|X)$ is near $1/2$.

Most favorable case:

for some $c > 0$, $\Pr(0 < |2\eta(X) - 1| < c) = 0$.

Low Noise

Definition: [Tsybakov] The distribution P on $\mathcal{X} \times \{\pm 1\}$ has *noise exponent* $0 \leq \alpha < \infty$ if there is a $c > 0$ such that

$$\Pr(0 < |2\eta(X) - 1| < \epsilon) \leq c\epsilon^\alpha.$$

- Tsybakov considered empirical risk minimization in binary classification.
- Under the noise assumption, if the Bayes classifier is in the function class, the risk of the empirical risk minimizer converges surprisingly quickly to the minimum.

Overview

1. Classification with convex loss.
2. Universal consistency of large margin algorithms.
3. Classification problems with low noise.
 - Large margin classifiers exploit low noise.
 - Adaptivity to low noise.

Risk Bounds with Low Noise: Convex Losses

Low noise improves the comparison inequality:

(Bartlett, Jordan, McAuliffe)

$$c(R(f) - R^*)^\beta \psi \left(\frac{(R(f) - R^*)^{1-\beta}}{2c} \right) \leq R_\phi(f) - R_\phi^*,$$

where $\beta = \frac{\alpha}{1 + \alpha} \in [0, 1]$. (Consider, for example, $\alpha = \infty$.)

- Strictly convex loss ϕ (e.g., AdaBoost's exponential loss)
 $\Rightarrow \psi$ strictly convex \Rightarrow strict improvement.

Risk Bounds with Low Noise: Convex Losses

Example: Suppose that ϕ has quadratic modulus of convexity, $\Pr(0 < |2\eta(X) - 1| < c) = 0$, and \hat{f} minimizes \hat{R}_ϕ over a finite-dimensional function class \mathcal{F} . Then

$$\mathbf{E}R(\hat{f}) - R^* \leq C \left(\inf_{f \in \mathcal{F}} R_\phi(f) - R_\phi^* + \frac{\log n}{n} \right).$$

- Striking: the fluctuations in $\hat{R}(\hat{f})$ are of the order of $1/\sqrt{n}$ in this case.
- Note that the algorithm minimizes the empirical ϕ -risk as before, but now an improvement in the noise exponent gives an improvement in the rate.

Low Noise: Small Variance

Definition: The distribution P on $\mathcal{X} \times \{\pm 1\}$ has *noise exponent* $0 \leq \alpha < \infty$ if there is a $c > 0$ such that

$$\Pr(0 < |2\eta(X) - 1| < \epsilon) \leq c\epsilon^\alpha.$$

- Equivalently, there is a c such that for every $f : \mathcal{X} \rightarrow \{\pm 1\}$,

$$\begin{aligned} \Pr(f(X) \neq f^*(X)) &\leq c(R(f) - R^*)^\beta \\ \Leftrightarrow \mathbf{E}(\ell_f - \ell_{f^*})^2 &\leq c(\mathbf{E}(\ell_f - \ell_{f^*}))^\beta, \end{aligned}$$

where f^* is the Bayes decision rule and $\beta = \frac{\alpha}{1 + \alpha}$.

Low Noise: Small Variance

Suppose that, for some g^* (think ℓ_{f^*}) and for all g (think ℓ_f),

$$b\left(\sqrt{\text{Var}(g - g^*)}\right) \leq \mathbf{E}(g - g^*),$$

where b is a convex, increasing function.

- The **variance of the excess loss** is bounded in terms of its expectation.
- As the risk, $\mathbf{E}g$, approaches the optimal risk, $\mathbf{E}g^*$, the loss g becomes more correlated with g^* .
- This ensures that the excess risk $\mathbf{E}(\hat{g} - g^*)$ for the empirical minimizer \hat{g} , converges quickly.

Local Low Noise and Model Selection

Suppose that we wish to do model selection over a nested hierarchy,

$$F_1 \subseteq F_2 \subseteq \dots \subseteq F_m \subseteq \dots$$

Tsybakov's low noise assumption bounds the variance of the excess loss of all functions. Instead, suppose we have convex increasing b_m for which

$$\text{for all } m \text{ and all } f \in F_m, \quad b_m \left(\sqrt{\text{Var}(\ell_f - \ell_{f^*})} \right) \leq \mathbf{E}(\ell_f - \ell_{f^*}).$$

This allows for the **local condition** to be favorable, even when the best **global condition** is weak.

For instance, for some small model F_m , we might have small variance of excess loss, that is, the best functions tend to agree with the Bayes rule.

Local Low Noise and Model Selection

Consider penalization-based model selection schemes:

empirical minimizer in F_m :
$$\hat{f}_m = \arg \min_{f \in F_m} \hat{\mathbf{E}} \ell_f,$$

selected model:
$$\hat{m} = \arg \min_m \left(\hat{\mathbf{E}} \ell_{\hat{f}_m} + \text{pen}(m) \right),$$

estimator:
$$\hat{f} = \hat{f}_{\hat{m}}.$$

Ideally, the penalty $\text{pen}(m)$ would approximate the difference between the risk and the empirical risk of \hat{f}_m .

Local Low Noise and Model Selection

Theorem: [with Sylvain Arlot] There is a penalty $\text{pen}(m)$ such that for all b_m , with probability at least $1 - e^{-t}$, $R(\hat{f}) - R^*$ is no more than

$$C \inf_m \left(\inf_{f \in F_m} (R(f) - R^*) + \text{pen}(m) + b_m^* \left(\sqrt{\frac{ct}{n}} + \frac{t}{n} \right) \right),$$

where b_m^* is the convex conjugate of b_m :

$$b_m^*(x) = \sup\{xy - b_m(y) : y \geq 0\}.$$

- The penalties are *local Rademacher averages*.

See (Massart, 2000), (Lugosi and Wegkamp, 2004), (Bartlett, Bousquet and Mendelson, 2005), (Koltchinskii, 2006)

- This model selection scheme *adapts* to the b_m .

Local Low Noise and Model Selection

For example, suppose that

- $\text{VCdim}(F_m) = V_m$,
- every $f \in F_m$ satisfies the **local low noise condition**

$$\Pr(f \neq f^*) \leq \frac{1}{h_m} (R(f) - R^*).$$

Then this model selection scheme satisfies the **oracle inequality**

$$\mathbf{E}R(\hat{f}) - R^* \leq c \inf_m \left(\inf_{f \in F_m} (R(f) - R^*) + \frac{V_m \log n}{h_m n} \right).$$

This is optimal up to a factor of $\log n$.

Convex methods for classification

1. Classification with convex loss.
 - Large margin methods.
 - Classification-calibrated ϕ : minimization of R_ϕ minimizes R .
2. Universal consistency of large margin algorithms.
 - AdaBoost, stopped after $t_n = O(n^{1-\alpha})$, is universally consistent.
3. Classification problems with low noise.
 - Large margin methods exploit low noise.
 - Penalization-based model selection methods that are adaptive to local low noise conditions.

slides at <http://www.stat.berkeley.edu/~bartlett>