Lectures on learning theory

Gábor Lugosi

ICREA and Pompeu Fabra University

Barcelona

what is learning theory?

A mathematical theory to understand the behavior of learning algorithms and assist their design.

what is learning theory?

A mathematical theory to understand the behavior of learning algorithms and assist their design.

Ingredients:

- ℜ Probability;
- ₭ (Linear) algebra;
- # Optimization;
- ₭ Complexity of algorithms;
- ₭ High-dimensional geometry;
- Statistics-hypothesis testing, regression, Bayesian methots, etc....

learning theory

Statistical learning

supervised-classification, regression, ranking, ...

unsupervised-clustering, density estimation, ...

semi-unsupervised learning

active learning

online learning

statistical learning

How is it different from "classical" statistics?

- ✤ Focus is on prediction rather than inference;
- ✤ Distribution-free approach;
- * Non-asymptotic results are preferred;
- ₭ High-dimensional problems;
- # Algorithmic aspects play a central role.

statistical learning

How is it different from "classical" statistics?

- Focus is on prediction rather than inference;
- ✤ Distribution-free approach;
- * Non-asymptotic results are preferred;
- ₭ High-dimensional problems;
- * Algorithmic aspects play a central role.

Here we focus on concentration inequalities.

a binary classification problem

(X, Y) is a pair of random variables.

 $X \in X$ represents the observation

 $\mathbf{Y} \in \{-1,1\}$ is the (binary) label.

A classifier is a function $\mathcal{X} \to \{-1,1\}$ whose risk is

 $\mathsf{R}(g) = \mathbb{P}\{g(\mathsf{X}) \neq \mathsf{Y}\} \ .$

a binary classification problem

(X, Y) is a pair of random variables.

 $X \in X$ represents the observation

 $\mathbf{Y} \in \{-1,1\}$ is the (binary) label.

A classifier is a function $\mathcal{X} \to \{-1,1\}$ whose risk is

 $\mathsf{R}(\mathsf{g}) = \mathbb{P}\{\mathsf{g}(\mathsf{X}) \neq \mathsf{Y}\} \ .$

training data: **n** i.i.d. observation/label pairs:

 $\mathsf{D}_n = ((\mathsf{X}_1,\mathsf{Y}_1),\ldots,(\mathsf{X}_n,\mathsf{Y}_n))$

The risk of a data-based classifier \boldsymbol{g}_n is

 $\mathsf{R}(g_n) = \mathbb{P}\{g_n(\mathsf{X}) \neq \mathsf{Y} | \mathsf{D}_n\} \ .$

Given a class $\boldsymbol{\mathcal{C}}$ of classifiers, choose one that minimizes the empirical risk:

$$g_n = \operatorname*{argmin}_{g \in \mathcal{C}} R_n(g) = \operatorname*{argmin}_{g \in \mathcal{C}} \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{g(X_i) \neq Y_i}$$

Given a class \mathcal{C} of classifiers, choose one that minimizes the empirical risk:

$$g_n = \operatorname*{argmin}_{g \in \mathcal{C}} \mathsf{R}_n(g) = \operatorname*{argmin}_{g \in \mathcal{C}} \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{g(X_i) \neq Y_i}$$

Fundamental questions:

- # How close is $R_n(g)$ to R(g)?
- # How close is $R(g_n)$ to $\min_{g \in C} R(g)$?
- # How close is $R(g_n)$ to $R_n(g_n)$?

To understand $|\mathbf{R}_n(\mathbf{g}) - \mathbf{R}(\mathbf{g})|$, we need to study deviations of empirical averages from their means.

For the other two, note that

$$|\mathsf{R}(g_n)-\mathsf{R}_n(g_n)| \leq \sup_{g\in\mathcal{C}} |\mathsf{R}(g)-\mathsf{R}_n(g)|$$

and

$$\begin{split} \mathsf{R}(\mathbf{g}_n) - \min_{\mathbf{g} \in \mathcal{C}} \mathsf{R}(\mathbf{g}) &= (\mathsf{R}(\mathbf{g}_n) - \mathsf{R}_n(\mathbf{g}_n)) + \left(\mathsf{R}_n(\mathbf{g}_n) - \min_{\mathbf{g} \in \mathcal{C}} \mathsf{R}(\mathbf{g})\right) \\ &\leq 2 \sup_{\mathbf{g} \in \mathcal{C}} |\mathsf{R}(\mathbf{g}) - \mathsf{R}_n(\mathbf{g})| \end{split}$$

To understand $|\mathbf{R}_n(\mathbf{g}) - \mathbf{R}(\mathbf{g})|$, we need to study deviations of empirical averages from their means.

For the other two, note that

$$|\mathsf{R}(g_n)-\mathsf{R}_n(g_n)| \leq \sup_{g\in\mathcal{C}} |\mathsf{R}(g)-\mathsf{R}_n(g)|$$

and

$$\begin{split} \mathsf{R}(g_n) - \min_{g \in \mathcal{C}} \mathsf{R}(g) &= (\mathsf{R}(g_n) - \mathsf{R}_n(g_n)) + \left(\mathsf{R}_n(g_n) - \min_{g \in \mathcal{C}} \mathsf{R}(g)\right) \\ &\leq 2 \sup_{g \in \mathcal{C}} |\mathsf{R}(g) - \mathsf{R}_n(g)| \end{split}$$

We need to understand uniform deviations of empirical averages from their means.

$\begin{array}{l} \mbox{markov's inequality} \\ \mbox{If } \mathbf{Z} \geq \mathbf{0}, \mbox{ then} \end{array}$

 $\mathbb{P}\{\mathsf{Z} > \mathsf{t}\} \leq \frac{\mathbb{E}\mathsf{Z}}{\mathsf{t}} \; .$

$\begin{array}{l} \mbox{markov's inequality} \\ \mbox{If ${\sf Z} \geq 0$, then} \end{array}$

$$\mathbb{P}\{\mathsf{Z} > \mathsf{t}\} \leq \frac{\mathbb{E}\mathsf{Z}}{\mathsf{t}} \; .$$

This implies Chebyshev's inequality: if Z has a finite variance $Var(Z) = \mathbb{E}(Z - \mathbb{E}Z)^2$, then

$$\mathbb{P}\{|\mathsf{Z}-\mathbb{E}\mathsf{Z}|>t\}=\mathbb{P}\{(\mathsf{Z}-\mathbb{E}\mathsf{Z})^2>t^2\}\leq \frac{\operatorname{Var}(\mathsf{Z})}{t^2}\;.$$

$\begin{array}{l} \mbox{markov's inequality} \\ \mbox{If ${\sf Z} \geq 0$, then} \end{array}$

$$\mathbb{P}\{\mathsf{Z} > \mathsf{t}\} \leq \frac{\mathbb{E}\mathsf{Z}}{\mathsf{t}} \; .$$

This implies Chebyshev's inequality: if Z has a finite variance $Var(Z) = \mathbb{E}(Z - \mathbb{E}Z)^2$, then

$$\mathbb{P}\{|\mathsf{Z}-\mathbb{E}\mathsf{Z}|>t\}=\mathbb{P}\{(\mathsf{Z}-\mathbb{E}\mathsf{Z})^2>t^2\}\leq \frac{\operatorname{Var}(\mathsf{Z})}{t^2}\;.$$



Andrey Markov (1856–1922)

sums of independent random variables

Let X_1, \ldots, X_n be independent real-valued and let $Z = \sum_{i=1}^n X_i$. By independence, $\operatorname{Var}(Z) = \sum_{i=1}^n \operatorname{Var}(X_i)$. If they are identically distributed, $\operatorname{Var}(Z) = n\operatorname{Var}(X_1)$, so

$$\mathbb{P}\left\{ \left|\sum_{i=1}^n X_i - n\mathbb{E}X_1 \right| > t \right\} \leq \frac{n\mathrm{Var}(X_1)}{t^2} \; .$$

Equivalently,

$$\mathbb{P}\left\{ \left| \sum_{i=1}^n X_i - n \mathbb{E} X_1 \right| > t \sqrt{n} \right\} \leq \frac{\operatorname{Var}(X_1)}{t^2} \; .$$

Typical deviations are at most of the order \sqrt{n} .

sums of independent random variables

Let X_1, \ldots, X_n be independent real-valued and let $Z = \sum_{i=1}^n X_i$. By independence, $\operatorname{Var}(Z) = \sum_{i=1}^n \operatorname{Var}(X_i)$. If they are identically distributed, $\operatorname{Var}(Z) = n\operatorname{Var}(X_1)$, so

$$\mathbb{P}\left\{ \left|\sum_{i=1}^n X_i - n \mathbb{E} X_1 \right| > t \right\} \leq \frac{n \mathrm{Var}(X_1)}{t^2} \; .$$

Equivalently,

$$\mathbb{P}\left\{ \left| \sum_{i=1}^n X_i - n \mathbb{E} X_1 \right| > t \sqrt{n} \right\} \leq \frac{\operatorname{Var}(X_1)}{t^2} \; .$$

Typical deviations are at most of the order \sqrt{n} .



Pafnuty Chebyshev (1821–1894)

By the central limit theorem,

$$\begin{split} \lim_{n \to \infty} \mathbb{P} \left\{ \sum_{i=1}^n X_i - n \mathbb{E} X_1 > t \sqrt{n} \right\} & = 1 - \Psi(t/\sqrt{\operatorname{Var}(X_1)}) \\ & \leq e^{-t^2/(2\operatorname{Var}(X_1))} \end{split}$$

so we expect an exponential decrease in $t^2/\mathrm{Var}(\mathsf{X}_1).$

By the central limit theorem,

$$\begin{split} \lim_{n \to \infty} \mathbb{P} \left\{ \sum_{i=1}^n X_i - n \mathbb{E} X_1 > t \sqrt{n} \right\} & = 1 - \Psi(t/\sqrt{\operatorname{Var}(X_1)}) \\ & \leq e^{-t^2/(2\operatorname{Var}(X_1))} \end{split}$$

so we expect an exponential decrease in $t^2/Var(X_1)$. Trick: use Markov's inequality in a more clever way: if $\lambda > 0$,

$$\mathbb{P}\{\mathsf{Z}-\mathbb{E}\mathsf{Z}>t\}=\mathbb{P}\left\{e^{\lambda(\mathsf{Z}-\mathbb{E}\mathsf{Z})}>e^{\lambda t}\right\}\leq\frac{\mathbb{E}e^{\lambda(\mathsf{Z}-\mathbb{E}\mathsf{Z})}}{e^{\lambda t}}$$

Now derive bounds for the moment generating function $\mathbb{E}e^{\lambda(Z-\mathbb{E}Z)}$ and optimize λ .

If $\textbf{Z} = \sum_{i=1}^n \textbf{X}_i$ is a sum of independent random variables,

$$\mathbb{E} e^{\lambda Z} = \mathbb{E} \prod_{i=1}^{n} e^{\lambda X_{i}} = \prod_{i=1}^{n} \mathbb{E} e^{\lambda X_{i}}$$

by independence. Now it suffices to find bounds for $\mathbb{E}e^{\lambda X_i}$.

If $\textbf{Z} = \sum_{i=1}^n \textbf{X}_i$ is a sum of independent random variables,

$$\mathbb{E} e^{\lambda Z} = \mathbb{E} \prod_{i=1}^{n} e^{\lambda X_{i}} = \prod_{i=1}^{n} \mathbb{E} e^{\lambda X_{i}}$$

by independence. Now it suffices to find bounds for $\mathbb{E}e^{\lambda X_i}$.



Serguei Bernstein (1880-1968)

Herman Chernoff (1923-)

hoeffding's inequality

If $X_1,\ldots,X_n\in[0,1]$, then

 $\mathbb{E} e^{\lambda(X_i - \mathbb{E} X_i)} \leq e^{\lambda^2/8}$.

hoeffding's inequality

If $X_1,\ldots,X_n\in[0,1]$, then

$$\mathbb{E} \mathrm{e}^{\lambda (\mathsf{X}_{\mathrm{i}} - \mathbb{E} \mathsf{X}_{\mathrm{i}})} \leq \mathrm{e}^{\lambda^2/8}$$
 .

We obtain

$$\mathbb{P}\left\{ \left| \frac{1}{n} \sum_{i=1}^n X_i - \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n X_i \right] \right| > t \right\} \leq 2e^{-2nt^2}$$



Wassily Hoeffding (1914–1991)

bernstein's inequality

Hoeffding's inequality is distribution free. It does not take variance information into account.

Bernstein's inequality is an often useful variant:

Let X_1,\ldots,X_n be independent such that $X_i\leq 1.$ Let $v=\sum_{i=1}^n\mathbb{E}\left[X_i^2\right].$ Then

$$\mathbb{P}\left\{\sum_{i=1}^n \left(X_i - \mathbb{E} X_i\right) \geq t\right\} \leq exp\left(-\frac{t^2}{2(\nu + t/3)}\right) \;.$$

a maximal inequality

Suppose $\boldsymbol{Y}_1,\ldots,\boldsymbol{Y}_N$ are sub-Gaussian in the sense that

 $\mathbb{E} \mathbf{e}^{\lambda \mathbf{Y}_{\mathsf{i}}} \leq \mathbf{e}^{\lambda^2 \sigma^2 / 2}$.

Then

$$\mathbb{E}\max_{i=1,...,\mathsf{N}}\mathsf{Y}_{\mathsf{i}} \leq \sigma\sqrt{2\log\mathsf{N}}\;.$$

a maximal inequality

Suppose Y_1,\ldots,Y_N are sub-Gaussian in the sense that

$$\mathbb{E}\mathbf{e}^{\lambda \mathbf{Y}_{\mathsf{i}}} \leq \mathbf{e}^{\lambda^2 \sigma^2/2}$$
 .

Then

$$\mathbb{E} \max_{\mathbf{i}=1,...,\mathsf{N}} \mathsf{Y}_{\mathbf{i}} \leq \sigma \sqrt{2 \log \mathsf{N}} \;.$$

Proof:

$$e^{\lambda \mathbb{E} \max_{i=1,\ldots,N} \mathsf{Y}_i} \leq \mathbb{E} e^{\lambda \max_{i=1,\ldots,N} \mathsf{Y}_i} \leq \sum_{i=1}^{\mathsf{N}} \mathbb{E} e^{\lambda \mathsf{Y}_i} \leq \mathsf{N} e^{\lambda^2 \sigma^2/2}$$

Take logarithms, and optimize in λ .

uniform deviations-finite classes

Let $A_1,\ldots,A_N\subset \mathcal{X}$ and let X_1,\ldots,X_n be i.i.d. random points in $\mathcal{X}.$ Let

$$\mathsf{P}(\mathsf{A}) = \mathbb{P}\{\mathsf{X}_1 \in \mathsf{A}\} \quad \text{and} \quad \mathsf{P}_n(\mathsf{A}) = \frac{1}{n}\sum_{i=1}^n \mathbb{1}_{\mathsf{X}_i \in \mathsf{A}}$$

By Hoeffding's inequality, for each A,

$$\begin{split} \mathbb{E} e^{\lambda(\mathsf{P}(\mathsf{A})-\mathsf{P}_n(\mathsf{A}))} &= \mathbb{E} e^{(\lambda/n)\sum_{i=1}^n (\mathsf{P}(\mathsf{A})-\mathbbm{1}_{X_i\in\mathsf{A}})} \\ &= \prod_{i=1}^n \mathbb{E} e^{(\lambda/n)(\mathsf{P}(\mathsf{A})-\mathbbm{1}_{X_i\in\mathsf{A}})} \leq e^{\lambda^2/(8n)} \;. \end{split}$$

By the maximal inequality,

$$\mathbb{E} \max_{j=1,\dots,\mathsf{N}} (\mathsf{P}(\mathsf{A}_j) - \mathsf{P}_\mathsf{n}(\mathsf{A}_j)) \leq \sqrt{\frac{\log\mathsf{N}}{2\mathsf{n}}} \; .$$

Suppose $A = \{a_1, \dots, a_n\} \subset \mathbb{R}^D$ is a finite set, D is large. We would like to embed A in \mathbb{R}^d where $d \ll D$.

Suppose $A = \{a_1, \ldots, a_n\} \subset \mathbb{R}^D$ is a finite set, D is large. We would like to embed A in \mathbb{R}^d where $d \ll D$. Is this possible? In what sense?

Suppose $A = \{a_1, \ldots, a_n\} \subset \mathbb{R}^D$ is a finite set, D is large. We would like to embed A in \mathbb{R}^d where $d \ll D$. Is this possible? In what sense? Given $\varepsilon > 0$, a function $f : \mathbb{R}^D \to \mathbb{R}^d$ is an ε -isometry if for all $a, a' \in A$,

 $\left(1-arepsilon
ight)\left\|\mathsf{a}-\mathsf{a}'
ight\|^2\leq \left\|\mathsf{f}(\mathsf{a})-\mathsf{f}(\mathsf{a}')
ight\|^2\leq \left(1+arepsilon
ight)\left\|\mathsf{a}-\mathsf{a}'
ight\|^2~.$

Suppose $A = \{a_1, \ldots, a_n\} \subset \mathbb{R}^D$ is a finite set, D is large. We would like to embed A in \mathbb{R}^d where $d \ll D$. Is this possible? In what sense? Given $\varepsilon > 0$, a function $f : \mathbb{R}^D \to \mathbb{R}^d$ is an ε -isometry if for all $a, a' \in A$,

$$\left(1-arepsilon
ight)\left\| \mathsf{a}-\mathsf{a}'
ight\|^2 \leq \left\| \mathsf{f}(\mathsf{a})-\mathsf{f}(\mathsf{a}')
ight\|^2 \leq \left(1+arepsilon
ight)\left\| \mathsf{a}-\mathsf{a}'
ight\|^2 \;.$$

Johnson-Lindenstrauss lemma: If $\mathbf{d} \geq (\mathbf{c}/\varepsilon^2) \log \mathbf{n}$, then there exists an ε -isometry $\mathbf{f} : \mathbb{R}^{\mathsf{D}} \to \mathbb{R}^{\mathsf{d}}$.

Suppose $A = \{a_1, \ldots, a_n\} \subset \mathbb{R}^D$ is a finite set, D is large. We would like to embed A in \mathbb{R}^d where $d \ll D$. Is this possible? In what sense? Given $\varepsilon > 0$, a function $f : \mathbb{R}^D \to \mathbb{R}^d$ is an ε -isometry if for all $a, a' \in A$,

$$\left(1-arepsilon
ight)\left\| \mathsf{a}-\mathsf{a}'
ight\|^2 \leq \left\| \mathsf{f}(\mathsf{a})-\mathsf{f}(\mathsf{a}')
ight\|^2 \leq \left(1+arepsilon
ight)\left\| \mathsf{a}-\mathsf{a}'
ight\|^2 \;.$$

Johnson-Lindenstrauss lemma: If $\mathbf{d} \geq (\mathbf{c}/\varepsilon^2) \log \mathbf{n}$, then there exists an ε -isometry $\mathbf{f} : \mathbb{R}^{\mathsf{D}} \to \mathbb{R}^{\mathsf{d}}$.

Independent of D!

We take **f** to be linear. How? At random!

We take f to be linear. How? At random! Let $f = (W_{i,j})_{d \times D}$ with

$$\mathsf{W}_{\mathsf{i},\mathsf{j}} = rac{\mathsf{I}}{\sqrt{\mathsf{d}}}\mathsf{X}_{\mathsf{i},\mathsf{j}}$$

where the $X_{i,j}$ are independent standard normal.

We take f to be linear. How? At random! Let $f=(W_{i,j})_{d\times D}$ with

$$\mathsf{W}_{\mathsf{i},\mathsf{j}} = \frac{1}{\sqrt{\mathsf{d}}}\mathsf{X}_{\mathsf{i},\mathsf{j}}$$

where the $X_{i,j}$ are independent standard normal.

For any $\mathbf{a} = (\alpha_1, \ldots, \alpha_D) \in \mathbb{R}^D$,

$$\mathbb{E} \| \mathbf{f}(\mathbf{a}) \|^2 = \frac{1}{\mathsf{d}} \sum_{\mathsf{i}=1}^{\mathsf{d}} \sum_{\mathsf{j}=1}^{\mathsf{D}} \alpha_\mathsf{j}^2 \mathbb{E} \mathsf{X}_{\mathsf{i},\mathsf{j}}^2 = \| \mathsf{a} \|^2 \; .$$

The expected squared distances are preserved!

We take f to be linear. How? At random! Let $f=(W_{i,j})_{d\times D}$ with

$$\mathsf{W}_{\mathsf{i},\mathsf{j}} = \frac{1}{\sqrt{\mathsf{d}}}\mathsf{X}_{\mathsf{i},\mathsf{j}}$$

where the $X_{i,j}$ are independent standard normal.

For any $\mathbf{a} = (\alpha_1, \ldots, \alpha_D) \in \mathbb{R}^D$,

$$\mathbb{E} \| \mathsf{f}(\mathsf{a}) \|^2 = \frac{1}{\mathsf{d}} \sum_{\mathsf{i}=1}^{\mathsf{d}} \sum_{\mathsf{j}=1}^{\mathsf{D}} \alpha_\mathsf{j}^2 \mathbb{E} \mathsf{X}_{\mathsf{i},\mathsf{j}}^2 = \| \mathsf{a} \|^2 \; .$$

The expected squared distances are preserved! $\|\mathbf{f}(\mathbf{a})\|^2 / \|\mathbf{a}\|^2$ is a weighted sum of squared normals.
random projections

Let
$$\mathbf{b} = \mathbf{a}_i - \mathbf{a}_j$$
 for some $\mathbf{a}_i, \mathbf{a}_j \in \mathbf{A}$. Then

$$\mathbb{P}\left\{\exists \mathbf{b} : \left|\frac{\|\mathbf{f}(\mathbf{b})\|^2}{\|\mathbf{b}\|^2} - \mathbf{1}\right| > \sqrt{\frac{8\log(n/\sqrt{\delta})}{d}} + \frac{8\log(n/\sqrt{\delta})}{d}\right\}$$

$$\leq {\binom{n}{2}}\mathbb{P}\left\{\left|\frac{\|\mathbf{f}(\mathbf{b})\|^2}{\|\mathbf{b}\|^2} - \mathbf{1}\right| > \sqrt{\frac{8\log(n/\sqrt{\delta})}{d}} + \frac{8\log(n/\sqrt{\delta})}{d}\right\}$$

$$\leq \delta \quad \text{(by a Bernstein-type inequality)}.$$
If $\mathbf{d} \geq (\mathbf{c}/\varepsilon^2)\log(n/\sqrt{\delta})$, then

$$\sqrt{rac{8\log(\mathsf{n}/\sqrt{\delta})}{\mathsf{d}}} + rac{8\log(\mathsf{n}/\sqrt{\delta})}{\mathsf{d}} \leq arepsilon$$

and **f** is an ε -isometry with probability $\geq 1 - \delta$.

martingale representation

 $\textbf{X}_1,\ldots,\textbf{X}_n$ are independent random variables taking values in some set $\mathcal{X}.$ Let $f:\mathcal{X}^n\to\mathbb{R}$ and

 $\mathsf{Z}=\mathsf{f}(\mathsf{X}_1,\ldots,\mathsf{X}_n)\;.$

Denote $\mathbb{E}_i[\cdot] = \mathbb{E}[\cdot|X_1, \dots, X_i]$. Thus, $\mathbb{E}_0 Z = \mathbb{E} Z$ and $\mathbb{E}_n Z = Z$.

martingale representation

 X_1,\ldots,X_n are independent random variables taking values in some set $\mathcal{X}.$ Let $f:\mathcal{X}^n\to\mathbb{R}$ and

 $\mathsf{Z}=\mathsf{f}(\mathsf{X}_1,\ldots,\mathsf{X}_n)\;.$

Denote $\mathbb{E}_i[\cdot] = \mathbb{E}[\cdot|X_1, \dots, X_i]$. Thus, $\mathbb{E}_0 Z = \mathbb{E} Z$ and $\mathbb{E}_n Z = Z$. Writing

$$\Delta_{\mathsf{i}} = \mathbb{E}_{\mathsf{i}}\mathsf{Z} - \mathbb{E}_{\mathsf{i}-1}\mathsf{Z} ,$$

we have

$$\mathsf{Z} - \mathbb{E}\mathsf{Z} = \sum_{i=1}^n \Delta_i$$

This is the Doob martingale representation of Z.

martingale representation

 $\begin{array}{l} \textbf{X}_1,\ldots,\textbf{X}_n \text{ are independent random variables taking values in}\\ \text{some set } \mathcal{X}. \text{ Let } f: \mathcal{X}^n \to \mathbb{R} \text{ and} \end{array}$

 $\mathsf{Z}=\mathsf{f}(\mathsf{X}_1,\ldots,\mathsf{X}_n)\;.$

Denote $\mathbb{E}_i[\cdot] = \mathbb{E}[\cdot|X_1, \dots, X_i]$. Thus, $\mathbb{E}_0 Z = \mathbb{E} Z$ and $\mathbb{E}_n Z = Z$. Writing

$$\Delta_{\mathsf{i}} = \mathbb{E}_{\mathsf{i}}\mathsf{Z} - \mathbb{E}_{\mathsf{i}-1}\mathsf{Z} ,$$

we have

$$Z - \mathbb{E} Z = \sum_{i=1}^n \Delta_i$$

This is the Doob martingale representation of **Z**.



Joseph Leo Doob (1910-2004)

martingale representation: the variance

$$\operatorname{Var}\left(\mathsf{Z}\right) = \mathbb{E}\left[\left(\sum_{i=1}^{n} \Delta_{i}\right)^{2}\right] = \sum_{i=1}^{n} \mathbb{E}\left[\Delta_{i}^{2}\right] + 2\sum_{j>i} \mathbb{E}\Delta_{i}\Delta_{j} \ .$$

Now if j > i, $\mathbb{E}_i \Delta_j = 0$, so

$$\mathbb{E}_i \Delta_j \Delta_i = \Delta_i \mathbb{E}_i \Delta_j = 0 \ ,$$

We obtain

$$\mathrm{Var}\left(\mathsf{Z}\right) = \mathbb{E}\left[\left(\sum_{i=1}^{n} \Delta_{i}\right)^{2}\right] = \sum_{i=1}^{n} \mathbb{E}\left[\Delta_{i}^{2}\right] \ .$$

martingale representation: the variance

$$\operatorname{Var}\left(\mathsf{Z}\right) = \mathbb{E}\left[\left(\sum_{i=1}^{n} \Delta_{i}\right)^{2}\right] = \sum_{i=1}^{n} \mathbb{E}\left[\Delta_{i}^{2}\right] + 2\sum_{j>i} \mathbb{E}\Delta_{i}\Delta_{j} \ .$$

Now if $j>i,\ \mathbb{E}_i\Delta_j=0,$ so

$$\mathbb{E}_i \Delta_j \Delta_i = \Delta_i \mathbb{E}_i \Delta_j = 0 \ ,$$

We obtain

$$\mathrm{Var}\left(\boldsymbol{Z}\right) = \mathbb{E}\left[\left(\sum_{i=1}^n \boldsymbol{\Delta}_i\right)^2\right] = \sum_{i=1}^n \mathbb{E}\left[\boldsymbol{\Delta}_i^2\right] \;.$$

From this, using independence, it is easy derive the Efron-Stein inequality.

Let X_1,\ldots,X_n be independent random variables taking values in $\mathcal{X}.$ Let $f:\mathcal{X}^n\to\mathbb{R}$ and $\mathsf{Z}=f(\mathsf{X}_1,\ldots,\mathsf{X}_n).$ Then

$$\operatorname{Var}(\mathsf{Z}) \leq \mathbb{E} \sum_{i=1}^n (\mathsf{Z} - \mathbb{E}^{(i)}\mathsf{Z})^2 = \mathbb{E} \sum_{i=1}^n \operatorname{Var}^{(i)}(\mathsf{Z}) \; .$$

where $\mathbb{E}^{(i)} Z$ is expectation with respect to the i-th variable X_i only.

Let X_1,\ldots,X_n be independent random variables taking values in $\mathcal{X}.$ Let $f:\mathcal{X}^n\to\mathbb{R}$ and $\mathsf{Z}=f(X_1,\ldots,X_n).$ Then

$$\operatorname{Var}(\mathsf{Z}) \leq \mathbb{E} \sum_{i=1}^{n} (\mathsf{Z} - \mathbb{E}^{(i)}\mathsf{Z})^{2} = \mathbb{E} \sum_{i=1}^{n} \operatorname{Var}^{(i)}(\mathsf{Z}) \; .$$

where $\mathbb{E}^{(i)} Z$ is expectation with respect to the i-th variable X_i only.

We obtain more useful forms by using that

$$\operatorname{Var}(\mathsf{X}) = rac{1}{2} \mathbb{E}(\mathsf{X} - \mathsf{X}')^2$$
 and $\operatorname{Var}(\mathsf{X}) \leq \mathbb{E}(\mathsf{X} - \mathsf{a})^2$

for any constant **a**.

If $\textbf{X}_1',\ldots,\textbf{X}_n'$ are independent copies of $\textbf{X}_1,\ldots,\textbf{X}_n,$ and

$$\mathsf{Z}'_i = f(\mathsf{X}_1, \dots, \mathsf{X}_{i-1}, \mathsf{X}'_i, \mathsf{X}_{i+1}, \dots, \mathsf{X}_n),$$

then

$$\operatorname{Var}(\mathsf{Z}) \leq \frac{1}{2} \mathbb{E} \left[\sum_{i=1}^n (\mathsf{Z} - \mathsf{Z}'_i)^2 \right]$$

Z is concentrated if it doesn't depend too much on any of its variables.

If X_1',\ldots,X_n' are independent copies of $X_1,\ldots,X_n,$ and

$$\mathsf{Z}'_i = f(\mathsf{X}_1, \dots, \mathsf{X}_{i-1}, \mathsf{X}'_i, \mathsf{X}_{i+1}, \dots, \mathsf{X}_n),$$

then

$$\operatorname{Var}(\mathsf{Z}) \leq \frac{1}{2} \mathbb{E} \left[\sum_{i=1}^n (\mathsf{Z} - \mathsf{Z}_i')^2 \right]$$

Z is concentrated if it doesn't depend too much on any of its variables.

If $\textbf{Z} = \sum_{i=1}^n \textbf{X}_i$ then we have an equality. Sums are the "least concentrated" of all functions!

If for some arbitrary functions \mathbf{f}_i

$$\mathsf{Z}_i = f_i(\mathsf{X}_1, \ldots, \mathsf{X}_{i-1}, \mathsf{X}_{i+1}, \ldots, \mathsf{X}_n) \ ,$$

then

$$\operatorname{Var}(\mathsf{Z}) \leq \mathbb{E}\left[\sum_{i=1}^n (\mathsf{Z}-\mathsf{Z}_i)^2\right]$$

efron, stein, and steele



Bradley Efron



Charles Stein



Mike Steele

example: uniform deviations

Let \mathcal{A} be a collection of subsets of \mathcal{X} , and let X_1, \ldots, X_n be n random points in \mathcal{X} drawn i.i.d. Let

$$\begin{split} \mathsf{P}(\mathsf{A}) &= \mathbb{P}\{\mathsf{X}_1 \in \mathsf{A}\} \quad \text{and} \quad \mathsf{P}_\mathsf{n}(\mathsf{A}) = \frac{1}{\mathsf{n}}\sum_{i=1}^{\mathsf{n}}\mathbbm{1}_{\mathsf{X}_i \in \mathsf{A}} \\ \text{If } \mathsf{Z} &= \mathsf{sup}_{\mathsf{A} \in \mathcal{A}} \, |\mathsf{P}(\mathsf{A}) - \mathsf{P}_\mathsf{n}(\mathsf{A})|, \\ &\quad \mathrm{Var}(\mathsf{Z}) \leq \frac{1}{2\mathsf{n}} \end{split}$$

example: uniform deviations

Let \mathcal{A} be a collection of subsets of \mathcal{X} , and let X_1, \ldots, X_n be n random points in \mathcal{X} drawn i.i.d. Let

$$\begin{split} \mathsf{P}(\mathsf{A}) &= \mathbb{P}\{\mathsf{X}_1 \in \mathsf{A}\} \quad \text{and} \quad \mathsf{P}_\mathsf{n}(\mathsf{A}) = \frac{1}{\mathsf{n}}\sum_{i=1}^{\mathsf{n}}\mathbbm{1}_{\mathsf{X}_i \in \mathsf{A}} \\ \text{If } \mathsf{Z} &= \mathsf{sup}_{\mathsf{A} \in \mathcal{A}} \, |\mathsf{P}(\mathsf{A}) - \mathsf{P}_\mathsf{n}(\mathsf{A})|, \\ &\quad \mathrm{Var}(\mathsf{Z}) \leq \frac{1}{2\mathsf{n}} \end{split}$$

regardless of the distribution and the richness of \mathcal{A} .

example: kernel density estimation

Let X_1, \ldots, X_n be i.i.d. real samples drawn according to some density ϕ . The kernel density estimate is

$$\phi_{n}(\mathbf{x}) = \frac{1}{nh} \sum_{i=1}^{n} K\left(\frac{\mathbf{x} - \mathbf{X}_{i}}{h}\right),$$

where h>0, and ${\sf K}$ is a nonnegative "kernel" $\int {\sf K}=1.$ The ${\sf L}_1$ error is

$$\mathsf{Z} = \mathsf{f}(\mathsf{X}_1, \dots, \mathsf{X}_\mathsf{n}) = \int |\phi(\mathsf{x}) - \phi_\mathsf{n}(\mathsf{x})| \mathsf{d}\mathsf{x} \; .$$

example: kernel density estimation

Let X_1, \ldots, X_n be i.i.d. real samples drawn according to some density ϕ . The kernel density estimate is

$$\phi_{n}(\mathbf{x}) = \frac{1}{nh} \sum_{i=1}^{n} K\left(\frac{\mathbf{x} - \mathbf{X}_{i}}{h}\right),$$

where h>0, and ${\sf K}$ is a nonnegative "kernel" $\int {\sf K}=1.$ The ${\sf L}_1$ error is

$$\mathsf{Z} = \mathsf{f}(\mathsf{X}_1, \dots, \mathsf{X}_n) = \int |\phi(\mathsf{x}) - \phi_\mathsf{n}(\mathsf{x})| \mathsf{d}\mathsf{x} \; .$$

It is easy to see that

$$\begin{split} |f(x_1,\ldots,x_n)-f(x_1,\ldots,x'_i,\ldots,x_n)| \\ &\leq \ \frac{1}{nh}\int \left|\mathsf{K}\left(\frac{x-x_i}{h}\right)-\mathsf{K}\left(\frac{x-x'_i}{h}\right)\right|dx \leq \frac{2}{n} \ , \\ & \text{ so we get } \ \mathbf{Var}(\mathsf{Z}) \leq \frac{2}{n} \ . \end{split}$$

bounding the expectation

Let $P_n'(A) = \frac{1}{n} \sum_{i=1}^n \mathbbm{1}_{X_i' \in A}$ and let \mathbb{E}' denote expectation only with respect to X_1', \ldots, X_n' .

$$\begin{split} \mathbb{E} \sup_{A \in \mathcal{A}} |\mathsf{P}_n(A) - \mathsf{P}(A)| &= \mathbb{E} \sup_{A \in \mathcal{A}} |\mathbb{E}'[\mathsf{P}_n(A) - \mathsf{P}'_n(A)]| \\ &\leq \mathbb{E} \sup_{A \in \mathcal{A}} |\mathsf{P}_n(A) - \mathsf{P}'_n(A)| &= \frac{1}{n} \mathbb{E} \sup_{A \in \mathcal{A}} \left| \sum_{i=1}^n (\mathbbm{1}_{X_i \in A} - \mathbbm{1}_{X'_i \in A}) \right| \end{split}$$

bounding the expectation

Let $P_n'(A) = \frac{1}{n} \sum_{i=1}^n \mathbbm{1}_{X_i' \in A}$ and let \mathbb{E}' denote expectation only with respect to X_1', \ldots, X_n' .

$$\mathbb{E} \sup_{A \in \mathcal{A}} |\mathsf{P}_n(A) - \mathsf{P}(A)| = \mathbb{E} \sup_{A \in \mathcal{A}} |\mathbb{E}'[\mathsf{P}_n(A) - \mathsf{P}'_n(A)]|$$

$$\leq \mathbb{E} \sup_{A \in \mathcal{A}} |\mathsf{P}_n(A) - \mathsf{P}'_n(A)| = \frac{1}{n} \mathbb{E} \sup_{A \in \mathcal{A}} \left| \sum_{i=1}^n (\mathbb{1}_{X_i \in A} - \mathbb{1}_{X'_i \in A}) \right|$$

Second symmetrization: if $\varepsilon_1, \ldots, \varepsilon_n$ are independent Rademacher variables, then

$$= \frac{1}{n} \mathbb{E} \sup_{A \in \mathcal{A}} \left| \sum_{i=1}^{n} \varepsilon_{i} (\mathbb{1}_{X_{i} \in A} - \mathbb{1}_{X_{i}' \in A}) \right| \leq \frac{2}{n} \mathbb{E} \sup_{A \in \mathcal{A}} \left| \sum_{i=1}^{n} \varepsilon_{i} \mathbb{1}_{X_{i} \in A} \right|$$

conditional rademacher average

lf

If

$$\begin{split} R_n &= \mathbb{E}_{\varepsilon} \sup_{A \in \mathcal{A}} \left| \sum_{i=1}^n \varepsilon_i \mathbb{1}_{X_i \in A} \right| \\ \end{split}$$
then

$$\begin{split} \mathbb{E} \sup_{A \in \mathcal{A}} |P_n(A) - P(A)| &\leq \frac{2}{n} \mathbb{E} R_n \end{split}$$

٠

conditional rademacher average

If
$$\mathsf{R}_n = \mathbb{E}_{\varepsilon} \sup_{\mathsf{A} \in \mathcal{A}} \left| \sum_{i=1}^n \varepsilon_i \mathbb{1}_{\mathsf{X}_i \in \mathsf{A}} \right|$$
 then

then

$$\mathbb{E}\sup_{A\in\mathcal{A}}|\mathsf{P}_n(A)-\mathsf{P}(A)|\leq \frac{2}{n}\mathbb{E}\mathsf{R}_n\;.$$

 $\mathbf{R}_{\mathbf{n}}$ is a data-dependent quantity!

concentration of conditional rademacher average

Define

$$\mathsf{R}_{\mathsf{n}}^{(\mathsf{i})} = \mathbb{E}_{\varepsilon} \sup_{\mathsf{A} \in \mathcal{A}} \left| \sum_{\mathsf{j} \neq \mathsf{i}} \varepsilon_{\mathsf{j}} \mathbb{1}_{\mathsf{X}_{\mathsf{j}} \in \mathsf{A}} \right|$$

One can show easily that

$$0 \leq \mathsf{R}_n - \mathsf{R}_n^{(i)} \leq 1 \quad \text{and} \quad \sum_{i=1}^n (\mathsf{R}_n - \mathsf{R}_n^{(i)}) \leq \mathsf{R}_n \ .$$

By the Efron-Stein inequality,

$$\operatorname{Var}(\mathsf{R}_n) \leq \mathbb{E} \sum_{i=1}^n (\mathsf{R}_n - \mathsf{R}_n^{(i)})^2 \leq \mathbb{E} \mathsf{R}_n \; .$$

concentration of conditional rademacher average

Define

$$\mathsf{R}_{\mathsf{n}}^{(\mathsf{i})} = \mathbb{E}_{\varepsilon} \sup_{\mathsf{A} \in \mathcal{A}} \left| \sum_{\mathsf{j} \neq \mathsf{i}} \varepsilon_{\mathsf{j}} \mathbb{1}_{\mathsf{X}_{\mathsf{j}} \in \mathsf{A}} \right|$$

One can show easily that

$$0 \leq \mathsf{R}_n - \mathsf{R}_n^{(i)} \leq 1 \quad \text{and} \quad \sum_{i=1}^n (\mathsf{R}_n - \mathsf{R}_n^{(i)}) \leq \mathsf{R}_n \ .$$

By the Efron-Stein inequality,

$$\operatorname{Var}(\mathsf{R}_n) \leq \mathbb{E} \sum_{i=1}^n (\mathsf{R}_n - \mathsf{R}_n^{(i)})^2 \leq \mathbb{E} \mathsf{R}_n \; .$$

Standard deviation is at most $\sqrt{\mathbb{E}R_n}$!

concentration of conditional rademacher average

Define

$$\mathsf{R}_{\mathsf{n}}^{(\mathsf{i})} = \mathbb{E}_{\varepsilon} \sup_{\mathsf{A} \in \mathcal{A}} \left| \sum_{\mathsf{j} \neq \mathsf{i}} \varepsilon_{\mathsf{j}} \mathbb{1}_{\mathsf{X}_{\mathsf{j}} \in \mathsf{A}} \right|$$

One can show easily that

$$0 \leq \mathsf{R}_n - \mathsf{R}_n^{(i)} \leq 1 \quad \text{and} \quad \sum_{i=1}^n (\mathsf{R}_n - \mathsf{R}_n^{(i)}) \leq \mathsf{R}_n \ .$$

By the Efron-Stein inequality,

$$\operatorname{Var}(\mathsf{R}_n) \leq \mathbb{E} \sum_{i=1}^n (\mathsf{R}_n - \mathsf{R}_n^{(i)})^2 \leq \mathbb{E} \mathsf{R}_n \; .$$

Standard deviation is at most $\sqrt{\mathbb{E}R_n}$!

Such functions are called self-bounding.

bounding the conditional rademacher average

If $S(X_1^n, \mathcal{A})$ is the number of different sets of form

 $\{X_1,\ldots,X_n\}\cap A:A\in\mathcal{A}$

then R_n is the maximum of $S(X_1^n, A)$ sub-Gaussian random variables. By the maximal inequality,

$$\frac{1}{2}\mathsf{R}_\mathsf{n} \leq \sqrt{\frac{\log\mathsf{S}(\mathsf{X}_1^\mathsf{n},\mathcal{A})}{2\mathsf{n}}} \; .$$

bounding the conditional rademacher average

If $S(X_1^n, \mathcal{A})$ is the number of different sets of form

 $\{X_1,\ldots,X_n\}\cap A:A\in\mathcal{A}$

then R_n is the maximum of $S(X_1^n, \mathcal{A})$ sub-Gaussian random variables. By the maximal inequality,

$$rac{1}{2}\mathsf{R}_{\mathsf{n}} \leq \sqrt{rac{\mathsf{log}}{\mathsf{S}}(\mathsf{X}_{1}^{\mathsf{n}},\mathcal{A})}{2\mathsf{n}}} \; .$$

In particular,

$$\mathbb{E}\sup_{\mathsf{A}\in\mathcal{A}}|\mathsf{P}_\mathsf{n}(\mathsf{A})-\mathsf{P}(\mathsf{A})|\leq 2\mathbb{E}\sqrt{\frac{\log\mathsf{S}(\mathsf{X}_1^\mathsf{n},\mathcal{A})}{2\mathsf{n}}}$$

random VC dimension

Let $V = V(x_1^n, A)$ be the size of the largest subset of $\{x_1, \ldots, x_n\}$ shattered by A. By Sauer's lemma,

 $\log \mathsf{S}(\mathsf{X}_1^n,\mathcal{A}) \leq \mathsf{V}(\mathsf{X}_1^n,\mathcal{A}) \log(n+1)$.

random VC dimension

Let $V = V(x_1^n, A)$ be the size of the largest subset of $\{x_1, \ldots, x_n\}$ shattered by A. By Sauer's lemma,

$\log \mathsf{S}(\mathsf{X}_1^n,\mathcal{A}) \leq \mathsf{V}(\mathsf{X}_1^n,\mathcal{A}) \log(n+1)$.

V is also self-bounding:

$$\sum_{i=1}^n (\mathsf{V}-\mathsf{V}^{(i)})^2 \leq \mathsf{V}$$

so by Efron-Stein,

 $\operatorname{Var}(\mathsf{V}) \leq \mathbb{E}\mathsf{V}$

vapnik and chervonenkis



Vladimir Vapnik



Alexey Chervonenkis

beyond the variance

 X_1, \ldots, X_n are independent random variables taking values in some set \mathcal{X} . Let $f : \mathcal{X}^n \to \mathbb{R}$ and $Z = f(X_1, \ldots, X_n)$. Recall the Doob martingale representation:

$$\mathsf{Z} - \mathbb{E}\mathsf{Z} = \sum_{i=1}^n \Delta_i \quad \text{where} \quad \Delta_i = \mathbb{E}_i\mathsf{Z} - \mathbb{E}_{i-1}\mathsf{Z} \ ,$$

with $\mathbb{E}_i[\cdot] = \mathbb{E}[\cdot | X_1, \dots, X_i]$.

To get exponential inequalities, we bound the moment generating function $\mathbb{E}e^{\lambda(Z-\mathbb{E}Z)}$.

azuma's inequality

Suppose that the martingale differences are bounded: $|\Delta_i| \leq c_i.$ Then

$$\begin{split} \mathbb{E} \mathbf{e}^{\lambda(\mathsf{Z}-\mathbb{E}\mathsf{Z})} &= \mathbb{E} \mathbf{e}^{\lambda\left(\sum_{i=1}^{n} \Delta_{i}\right)} = \mathbb{E} \mathbb{E}_{n} \mathbf{e}^{\lambda\left(\sum_{i=1}^{n-1} \Delta_{i}\right) + \lambda \Delta_{n}} \\ &= \mathbb{E} \mathbf{e}^{\lambda\left(\sum_{i=1}^{n-1} \Delta_{i}\right)} \mathbb{E}_{n} \mathbf{e}^{\lambda \Delta_{n}} \\ &\leq \mathbb{E} \mathbf{e}^{\lambda\left(\sum_{i=1}^{n-1} \Delta_{i}\right)} \mathbf{e}^{\lambda^{2} c_{n}^{2}/2} \text{ (by Hoeffding)} \\ & \cdots \\ &< \mathbf{e}^{\lambda^{2} \left(\sum_{i=1}^{n} c_{i}^{2}\right)/2} \text{ .} \end{split}$$

This is the Azuma-Hoeffding inequality for sums of bounded martingale differences.

bounded differences inequality If $Z = f(X_1, ..., X_n)$ and f is such that

 $|f(x_1,\ldots,x_n)-f(x_1,\ldots,x_i',\ldots,x_n)|\leq c_i$

then the martingale differences are bounded.

bounded differences inequality If $Z = f(X_1, ..., X_n)$ and f is such that

 $|f(x_1,\ldots,x_n)-f(x_1,\ldots,x_i',\ldots,x_n)|\leq c_i$

then the martingale differences are bounded.

Bounded differences inequality: if X_1,\ldots,X_n are independent, then

 $\mathbb{P}\{|\mathsf{Z}-\mathbb{E}\mathsf{Z}|>t\}\leq 2e^{-2t^2/\sum_{i=1}^n c_i^2}\;.$

bounded differences inequality If $Z = f(X_1, ..., X_n)$ and f is such that

 $|f(x_1,\ldots,x_n)-f(x_1,\ldots,x_i',\ldots,x_n)|\leq c_i$

then the martingale differences are bounded.

Bounded differences inequality: if X_1,\ldots,X_n are independent, then

$$\mathbb{P}\{|\mathsf{Z}-\mathbb{E}\mathsf{Z}|>t\}\leq 2e^{-2t^2/\sum_{i=1}^n c_i^2}\;.$$

McDiarmid's inequality.



Colin McDiarmid

hoeffding in a hilbert space

Let X_1,\ldots,X_n be independent zero-mean random variables in a separable Hilbert space such that $||X_i|| \leq c/2$ and denote $v=nc^2/4.$ Then, for all $t\geq \sqrt{v},$

$$\mathbb{P}\left\{ \left\|\sum_{i=1}^n X_i\right\| > t \right\} \leq e^{-(t-\sqrt{\nu})^2/(2\nu)} \; .$$

hoeffding in a hilbert space

Let X_1,\ldots,X_n be independent zero-mean random variables in a separable Hilbert space such that $||X_i|| \leq c/2$ and denote $v = nc^2/4$. Then, for all $t \geq \sqrt{v}$,

$$\mathbb{P}\left\{ \left\|\sum_{i=1}^n X_i\right\| > t \right\} \leq e^{-(t-\sqrt{\nu})^2/(2\nu)} \; .$$

Proof: By the triangle inequality, $\left\|\sum_{i=1}^n X_i\right\|$ has the bounded differences property with constants c, so

$$\begin{split} \mathbb{P}\left\{ \left\|\sum_{i=1}^{n}X_{i}\right\| > t\right\} &= \mathbb{P}\left\{ \left\|\sum_{i=1}^{n}X_{i}\right\| - \mathbb{E}\left\|\sum_{i=1}^{n}X_{i}\right\| > t - \mathbb{E}\left\|\sum_{i=1}^{n}X_{i}\right\|\right\} \\ &\leq exp\left(-\frac{\left(t - \mathbb{E}\left\|\sum_{i=1}^{n}X_{i}\right\|\right)^{2}}{2v}\right) \,. \end{split}$$

Also,

$$\mathbb{E} \left\| \sum_{i=1}^{n} X_{i} \right\| \leq \sqrt{\mathbb{E} \left\| \sum_{i=1}^{n} X_{i} \right\|^{2}} = \sqrt{\sum_{i=1}^{n} \mathbb{E} \left\| X_{i} \right\|^{2}} \leq \sqrt{\nu} \; .$$

bounded differences inequality

- ₭ Easy to use.
- ✤ Distribution free.
- # Often close to optimal (e.g., \textbf{L}_1 error of kernel density estimate).
- ✤ Does not exploit "variance information."
- ✤ Often too rigid.
- * Other methods are necessary.
shannon entropy

If \mathbf{X}, \mathbf{Y} are random variables taking values in a set of size \mathbf{N} ,

$$H(X) = -\sum_{x} p(x) \log p(x)$$

$$\begin{aligned} H(X|Y) &= H(X,Y) - H(Y) \\ &= -\sum_{x,y} p(x,y) \log p(x|y) \end{aligned}$$

 $\mathsf{H}(\mathsf{X}) \leq \mathsf{log}\,\mathsf{N} \quad \mathsf{and} \quad \mathsf{H}(\mathsf{X}|\mathsf{Y}) \leq \mathsf{H}(\mathsf{X})$



Claude Shannon (1916–2001)

han's inequality

If
$$X = (X_1, \dots, X_n)$$
 and
 $X^{(i)} = (X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n)$, then

$$\sum_{i=1}^n \left(H(X) - H(X^{(i)}) \right) \le H(X)$$



Proof:

$$\begin{split} \mathsf{H}(\mathsf{X}) &= \mathsf{H}(\mathsf{X}^{(i)}) + \mathsf{H}(\mathsf{X}_{i} | \mathsf{X}^{(i)}) \\ &\leq \mathsf{H}(\mathsf{X}^{(i)}) + \mathsf{H}(\mathsf{X}_{i} | \mathsf{X}_{1}, \dots, \mathsf{X}_{i-1}) \end{split}$$

Te Sun Han

Since $\sum_{i=1}^n H(X_i|X_1,\ldots,X_{i-1}) = H(X)$, summing the inequality, we get

$$(\mathsf{n}-1)\mathsf{H}(\mathsf{X}) \leq \sum_{\mathsf{i}=1}^{\mathsf{n}}\mathsf{H}(\mathsf{X}^{(\mathsf{i})}) \; .$$

subadditivity of entropy

The entropy of a random variable $\mathbf{Z} \geq \mathbf{0}$ is

$\operatorname{Ent}(\mathsf{Z}) = \mathbb{E}\Phi(\mathsf{Z}) - \Phi(\mathbb{E}\mathsf{Z})$

where $\Phi(x) = x \log x$. By Jensen's inequality, $Ent(Z) \ge 0$.

subadditivity of entropy

The entropy of a random variable $\mathbf{Z} \geq \mathbf{0}$ is

$\operatorname{Ent}(\mathsf{Z}) = \mathbb{E}\Phi(\mathsf{Z}) - \Phi(\mathbb{E}\mathsf{Z})$

where $\Phi(x)=x\log x$. By Jensen's inequality, $\operatorname{Ent}(Z)\geq 0$. Han's inequality implies the following sub-additivity property. Let X_1,\ldots,X_n be independent and let $Z=f(X_1,\ldots,X_n)$, where $f\geq 0$. Denote

$$\operatorname{Ent}^{(i)}(\mathsf{Z}) = \mathbb{E}^{(i)} \Phi(\mathsf{Z}) - \Phi(\mathbb{E}^{(i)}\mathsf{Z})$$

Then

$$\operatorname{Ent}(\mathsf{Z}) \leq \mathbb{E} \sum_{i=1}^{n} \operatorname{Ent}^{(i)}(\mathsf{Z})$$
 .

a logarithmic sobolev inequality on the hypercube

Let $X=(X_1,\ldots,X_n)$ be uniformly distributed over $\{-1,1\}^n.$ If $f:\{-1,1\}^n\to\mathbb{R}$ and Z=f(X),

$$\operatorname{Ent}(\mathsf{Z}^2) \leq \frac{1}{2} \mathbb{E} \sum_{i=1}^n (\mathsf{Z} - \mathsf{Z}'_i)^2$$

The proof uses subadditivity of the entropy and calculus for the case $\mathbf{n} = \mathbf{1}$.

Implies Efron-Stein and the edge-isoperimetric inequality.

herbst's argument: exponential concentration

If $f : \{-1, 1\}^n \to \mathbb{R}$, the log-Sobolev inequality may be used with $g(x) = e^{\lambda f(x)/2}$ where $\lambda \in \mathbb{R}$. If $F(\lambda) = \mathbb{E}e^{\lambda Z}$ is the moment generating function of Z = f(X), $\operatorname{Ent}(g(X)^2) = \lambda \mathbb{E}\left[Ze^{\lambda Z}\right] - \mathbb{E}\left[e^{\lambda Z}\right] \log \mathbb{E}\left[Ze^{\lambda Z}\right]$ $= \lambda F'(\lambda) - F(\lambda) \log F(\lambda)$.

Differential inequalities are obtained for $F(\lambda)$.

herbst's argument

As an example, suppose f is such that $\sum_{i=1}^n (Z-Z_i')_+^2 \leq v.$ Then by the log-Sobolev inequality,

$$\lambda \mathsf{F}'(\lambda) - \mathsf{F}(\lambda) \log \mathsf{F}(\lambda) \leq rac{\mathsf{v}\lambda^2}{4}\mathsf{F}(\lambda)$$

If $G(\lambda) = \log F(\lambda)$, this becomes

$$\left(rac{\mathsf{G}(\lambda)}{\lambda}
ight)'\leqrac{\mathsf{v}}{4}\;.$$

This can be integrated: $\mathsf{G}(\lambda) \leq \lambda \mathbb{E}\mathsf{Z} + \lambda \mathsf{v}/4$, so

 $\mathsf{F}(\lambda) \leq \mathrm{e}^{\lambda \mathbb{E}\mathsf{Z} - \lambda^2 \mathsf{v}/4}$

This implies

$$\mathbb{P}\{\mathsf{Z} > \mathbb{E}\mathsf{Z} + \mathsf{t}\} \leq e^{-\mathsf{t}^2/\mathsf{v}}$$

herbst's argument

As an example, suppose f is such that $\sum_{i=1}^n (Z-Z_i')_+^2 \leq v.$ Then by the log-Sobolev inequality,

$$\lambda \mathsf{F}'(\lambda) - \mathsf{F}(\lambda) \log \mathsf{F}(\lambda) \leq rac{\mathsf{v}\lambda^2}{4}\mathsf{F}(\lambda)$$

If $G(\lambda) = \log F(\lambda)$, this becomes

$$\left(rac{\mathsf{G}(\lambda)}{\lambda}
ight)'\leqrac{\mathsf{v}}{4}\;.$$

This can be integrated: $\mathsf{G}(\lambda) \leq \lambda \mathbb{E}\mathsf{Z} + \lambda \mathsf{v}/4$, so

 $\mathsf{F}(\lambda) \leq \mathrm{e}^{\lambda \mathbb{E}\mathsf{Z} - \lambda^2 \mathsf{v}/4}$

This implies

$$\mathbb{P}\{\mathsf{Z} > \mathbb{E}\mathsf{Z} + \mathsf{t}\} \leq e^{-\mathsf{t}^2/\mathsf{v}}$$

Stronger than the bounded differences inequality!

gaussian log-sobolev and concentration inequalities

Let $X = (X_1, \dots, X_n)$ be a vector of i.i.d. standard normal If $f : \mathbb{R}^n \to \mathbb{R}$ and Z = f(X),

$$\operatorname{Ent}(\mathsf{Z}^2) \leq 2\mathbb{E}\left[\|\nabla f(\mathsf{X})\|^2\right]$$

This can be proved using the central limit theorem and the Bernoulli log-Sobolev inequality.

gaussian log-sobolev and concentration inequalities

Let $X = (X_1, \dots, X_n)$ be a vector of i.i.d. standard normal If $f : \mathbb{R}^n \to \mathbb{R}$ and Z = f(X),

$$\operatorname{Ent}(\mathsf{Z}^2) \leq 2\mathbb{E}\left[\|\nabla f(\mathsf{X})\|^2\right]$$

This can be proved using the central limit theorem and the Bernoulli log-Sobolev inequality. It implies the Gaussian concentration inequality: Suppose **f** is Lipschitz: for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{n}$,

 $|f(x) - f(y)| \le L \|x - y\| \ .$

Then, for all $\mathbf{t} > \mathbf{0}$,

 $\mathbb{P}\left\{f(X) - \mathbb{E}f(X) \geq t\right\} \leq e^{-t^2/(2L^2)} \;.$

an application: supremum of a gaussian process

Let $(X_t)_{t\in\mathcal{T}}$ be an almost surely continuous centered Gaussian process. Let $\mathsf{Z}=sup_{t\in\mathcal{T}}\,\mathsf{X}_t.$ If

$$\sigma^2 = \sup_{\mathbf{t}\in\mathcal{T}} \left(\mathbb{E}\left[\mathbf{X}_{\mathbf{t}}^2
ight]
ight) \;,$$

then

$$\mathbb{P}\left\{|\mathsf{Z} - \mathbb{E}\mathsf{Z}| \geq \mathsf{u}\right\} \leq 2\mathsf{e}^{-\mathsf{u}^2/(2\sigma^2)}$$

an application: supremum of a gaussian process

Let $(X_t)_{t\in\mathcal{T}}$ be an almost surely continuous centered Gaussian process. Let $\mathsf{Z}=sup_{t\in\mathcal{T}}\,\mathsf{X}_t.$ If

$$\sigma^2 = \sup_{\mathbf{t} \in \mathcal{T}} \left(\mathbb{E} \left[\mathbf{X}_{\mathbf{t}}^2 \right] \right) \;,$$

then

$$\mathbb{P}\left\{ \left| \mathsf{Z} - \mathbb{E}\mathsf{Z} \right| \geq \mathsf{u}
ight\} \leq 2\mathsf{e}^{-\mathsf{u}^2/(2\sigma^2)}$$

Proof: We may assume $\mathcal{T} = \{1, ..., n\}$. Let Γ be the covariance matrix of $X = (X_1, \ldots, X_n)$. Let $A = \Gamma^{1/2}$. If Y is a standard normal vector, then

$$f(\mathbf{Y}) = \max_{i=1,\dots,n} (\mathbf{AY})_i \stackrel{\text{distr.}}{=} \max_{i=1,\dots,n} X_i$$

By Cauchy-Schwarz,

$$\begin{split} |(\mathsf{A}\mathsf{u})_{\mathsf{i}} - (\mathsf{A}\mathsf{v})_{\mathsf{i}}| &= \left| \sum_{\mathsf{j}} \mathsf{A}_{\mathsf{i},\mathsf{j}} \left(\mathsf{u}_{\mathsf{j}} - \mathsf{v}_{\mathsf{j}} \right) \right| \leq \left(\sum_{\mathsf{j}} \mathsf{A}_{\mathsf{i},\mathsf{j}}^2 \right)^{1/2} \|\mathsf{u} - \mathsf{v}\| \\ &\leq \sigma \|\mathsf{u} - \mathsf{v}\| \end{split}$$

beyond bernoulli and gaussian: the entropy method

For general distributions, logarithmic Sobolev inequalities are not available.

Solution: modified logarithmic Sobolev inequalities. Suppose X_1, \ldots, X_n are independent. Let $Z = f(X_1, \ldots, X_n)$ and $Z_i = f_i(X^{(i)}) = f_i(X_1, \ldots, X_{i-1}, X_{i+1}, \ldots, X_n)$.

Let
$$\phi(\mathbf{x}) = \mathbf{e}^{\mathbf{x}} - \mathbf{x} - \mathbf{1}$$
. Then for all $\lambda \in \mathbb{R}$,
 $\lambda \mathbb{E} \left[\mathsf{Z} \mathbf{e}^{\lambda \mathsf{Z}} \right] - \mathbb{E} \left[\mathbf{e}^{\lambda \mathsf{Z}} \right] \log \mathbb{E} \left[\mathbf{e}^{\lambda \mathsf{Z}} \right]$
 $\leq \sum_{i=1}^{n} \mathbb{E} \left[\mathbf{e}^{\lambda \mathsf{Z}} \phi \left(-\lambda (\mathsf{Z} - \mathsf{Z}_{i}) \right) \right].$



Michel Ledoux

the entropy method

Define $\mathsf{Z}_i = \mathsf{inf}_{\mathsf{x}_i'}\,\mathsf{f}(\mathsf{X}_1, \dots, \mathsf{x}_i', \dots, \mathsf{X}_n)$ and suppose

$$\sum_{i=1}^n (\mathsf{Z}-\mathsf{Z}_i)^2 \leq \mathsf{v} \ .$$

Then for all t > 0,

$$\mathbb{P}\left\{\mathsf{Z} - \mathbb{E}\mathsf{Z} > t\right\} \leq e^{-t^2/(2\nu)} \; .$$

the entropy method

Define $\mathsf{Z}_i = \mathsf{inf}_{\mathsf{x}_i'}\,\mathsf{f}(\mathsf{X}_1, \dots, \mathsf{x}_i', \dots, \mathsf{X}_n)$ and suppose

$$\sum_{i=1}^n (\mathsf{Z}-\mathsf{Z}_i)^2 \leq \mathsf{v} \ .$$

Then for all $\mathbf{t} > \mathbf{0}$,

$$\mathbb{P}\left\{\mathsf{Z} - \mathbb{E}\mathsf{Z} > \mathsf{t}\right\} \leq e^{-\mathsf{t}^2/(2\mathsf{v})} \; .$$

This implies the bounded differences inequality and much more.

example: the largest eigenvalue of a symmetric matrix Let $A = (X_{i,j})_{n \times n}$ be symmetric, the $X_{i,j}$ independent $(i \le j)$ with $|X_{i,j}| \le 1$. Let $Z = \lambda_1 = \text{sup } u^T A u$.

u:||u||=1

and suppose \mathbf{v} is such that $\mathbf{Z} = \mathbf{v}^{\mathsf{T}} \mathbf{A} \mathbf{v}$. $\mathbf{A}'_{i,j}$ is obtained by replacing $\mathbf{X}_{i,j}$ by $\mathbf{x}'_{i,j}$. Then

$$\begin{split} (\mathsf{Z} - \mathsf{Z}_{i,j})_+ &\leq \left(\mathsf{v}^\mathsf{T} \mathsf{A} \mathsf{v} - \mathsf{v}^\mathsf{T} \mathsf{A}'_{i,j} \mathsf{v} \right) \mathbbm{1}_{\mathsf{Z} > \mathsf{Z}_{i,j}} \\ &= \left(\mathsf{v}^\mathsf{T} (\mathsf{A} - \mathsf{A}'_{i,j}) \mathsf{v} \right) \mathbbm{1}_{\mathsf{Z} > \mathsf{Z}_{i,j}} \leq 2 \left(\mathsf{v}_i \mathsf{v}_j (\mathsf{X}_{i,j} - \mathsf{X}'_{i,j}) \right)_+ \\ &\leq 4 |\mathsf{v}_i \mathsf{v}_j| \ . \end{split}$$

Therefore,

$$\sum_{1 \leq i \leq j \leq n} (\mathsf{Z} - \mathsf{Z}'_{i,j})_+^2 \leq \sum_{1 \leq i \leq j \leq n} 16 |\mathsf{v}_i \mathsf{v}_j|^2 \leq 16 \left(\sum_{i=1}^n \mathsf{v}_i^2\right)^2 = 16 \; .$$

example: convex lipschitz functions

Let $f:[0,1]^n \to \mathbb{R}$ be a convex function. Let $Z_i = inf_{x_i'} f(X_1, \ldots, x_i', \ldots, X_n)$ and let X_i' be the value of x_i' for which the minimum is achieved. Then, writing $\overline{X}^{(i)} = (X_1, \ldots, X_{i-1}, X_i', X_{i+1}, \ldots, X_n)$,

$$\begin{split} \sum_{i=1}^{n} (\mathsf{Z} - \mathsf{Z}_i)^2 &= \sum_{i=1}^{n} (\mathsf{f}(\mathsf{X}) - \mathsf{f}(\overline{\mathsf{X}}^{(i)})^2 \\ &\leq \sum_{i=1}^{n} \left(\frac{\partial \mathsf{f}}{\partial \mathsf{x}_i}(\mathsf{X}) \right)^2 (\mathsf{X}_i - \mathsf{X}'_i)^2 \\ &\quad (\text{by convexity}) \\ &\leq \sum_{i=1}^{n} \left(\frac{\partial \mathsf{f}}{\partial \mathsf{x}_i}(\mathsf{X}) \right)^2 \\ &= \| \nabla \mathsf{f}(\mathsf{X}) \|^2 \leq \mathsf{L}^2 \;. \end{split}$$

self-bounding functions

Suppose Z satisfies

$$0 \leq \mathsf{Z} - \mathsf{Z}_i \leq 1 \quad \text{and} \quad \sum_{i=1}^n (\mathsf{Z} - \mathsf{Z}_i) \leq \mathsf{Z} \ .$$

Recall that $Var(Z) \leq \mathbb{E}Z$. We have much more:

 $\mathbb{P}\{\mathsf{Z} > \mathbb{E}\mathsf{Z} + t\} \leq e^{-t^2/(2\mathbb{E}\mathsf{Z} + 2t/3)}$

and

$$\mathbb{P}\{\mathsf{Z} < \mathbb{E}\mathsf{Z} - \mathsf{t}\} \leq e^{-\mathsf{t}^2/(2\mathbb{E}\mathsf{Z})}$$

self-bounding functions

Suppose ${\bf Z}$ satisfies

$$0 \leq \mathsf{Z} - \mathsf{Z}_i \leq 1 \quad \text{and} \quad \sum_{i=1}^n (\mathsf{Z} - \mathsf{Z}_i) \leq \mathsf{Z} \ .$$

Recall that $Var(Z) \leq \mathbb{E}Z$. We have much more:

$$\mathbb{P}\{\mathsf{Z} > \mathbb{E}\mathsf{Z} + t\} \leq e^{-t^2/(2\mathbb{E}\mathsf{Z} + 2t/3)}$$

and

$$\mathbb{P}\{\mathsf{Z} < \mathbb{E}\mathsf{Z} - \mathsf{t}\} \leq e^{-\mathsf{t}^2/(2\mathbb{E}\mathsf{Z})}$$

Rademacher averages and the random VC dimension are self bounding.

self-bounding functions

Suppose ${\bf Z}$ satisfies

$$0 \leq \mathsf{Z} - \mathsf{Z}_i \leq 1 \quad \text{and} \quad \sum_{i=1}^n (\mathsf{Z} - \mathsf{Z}_i) \leq \mathsf{Z} \ .$$

Recall that $Var(Z) \leq \mathbb{E}Z$. We have much more:

$$\mathbb{P}\{\mathsf{Z} > \mathbb{E}\mathsf{Z} + t\} \leq e^{-t^2/(2\mathbb{E}\mathsf{Z} + 2t/3)}$$

and

$$\mathbb{P}\{\mathsf{Z} < \mathbb{E}\mathsf{Z} - \mathsf{t}\} \leq e^{-\mathsf{t}^2/(2\mathbb{E}\mathsf{Z})}$$

Rademacher averages and the random VC dimension are self bounding.

Configuration functions.

weakly self-bounding functions

$$\begin{split} &f:\mathcal{X}^n\to [0,\infty) \text{ is weakly } (a,b)\text{-self-bounding if there exist} \\ &f_i:\mathcal{X}^{n-1}\to [0,\infty) \text{ such that for all } x\in\mathcal{X}^n, \end{split}$$

$$\sum_{i=1}^n \left(f(x)-f_i(x^{(i)})\right)^2 \leq af(x)+b\,.$$

weakly self-bounding functions

$$\begin{split} &f:\mathcal{X}^n\to [0,\infty) \text{ is weakly } (a,b)\text{-self-bounding if there exist} \\ &f_i:\mathcal{X}^{n-1}\to [0,\infty) \text{ such that for all } x\in\mathcal{X}^n, \end{split}$$

$$\sum_{i=1}^n \left(f(x)-f_i(x^{(i)})\right)^2 \leq af(x)+b\,.$$

Then

$$\mathbb{P}\left\{\mathsf{Z} \geq \mathbb{E}\mathsf{Z} + t\right\} \leq exp\left(-\frac{t^2}{2\left(a\mathbb{E}\mathsf{Z} + b + at/2\right)}\right) \;.$$

weakly self-bounding functions

$$\begin{split} &f:\mathcal{X}^n\to [0,\infty) \text{ is weakly } (a,b)\text{-self-bounding if there exist} \\ &f_i:\mathcal{X}^{n-1}\to [0,\infty) \text{ such that for all } x\in\mathcal{X}^n, \end{split}$$

$$\sum_{i=1}^n \left(f(x)-f_i(x^{(i)})\right)^2 \leq af(x)+b\,.$$

Then

$$\mathbb{P}\left\{\mathsf{Z} \geq \mathbb{E}\mathsf{Z} + t\right\} \leq \exp\left(-\frac{t^2}{2\left(a\mathbb{E}\mathsf{Z} + b + at/2\right)}\right) \; .$$

If, in addition, $f(x) - f_i(x^{(i)}) \leq 1,$ then for $0 < t \leq \mathbb{E} \mathsf{Z},$

$$\mathbb{P}\left\{\mathsf{Z} \leq \mathbb{E}\mathsf{Z} - t\right\} \leq \exp\left(-\frac{t^2}{2\left(a\mathbb{E}\mathsf{Z} + b + c_-t\right)}\right) \;.$$

where c = (3a - 1)/6.

Let $X = (X_1, \dots, X_n)$ have independent components, taking values in \mathcal{X}^n . Let $A \subset \mathcal{X}^n$. The Hamming distance of X to A is

$$d(X,A) = \min_{y \in A} d(X,y) = \min_{y \in A} \sum_{i=1}^{n} \mathbb{1}_{X_i \neq y_i} \ .$$



Michel Talagrand

Let $X = (X_1, \dots, X_n)$ have independent components, taking values in \mathcal{X}^n . Let $A \subset \mathcal{X}^n$. The Hamming distance of X to A is

$$d(X,A) = \min_{y \in A} d(X,y) = \min_{y \in A} \sum_{i=1}^n \mathbb{1}_{X_i \neq y_i} \ .$$



Michel Talagrand

$$\mathbb{P}\left\{\mathsf{d}(\mathsf{X},\mathsf{A})\geq\mathsf{t}+\sqrt{rac{\mathsf{n}}{2}\lograc{1}{\mathbb{P}[\mathsf{A}]}}
ight\}\leq\mathsf{e}^{-2\mathsf{t}^2/\mathsf{n}}\;.$$

Let $X = (X_1, \dots, X_n)$ have independent components, taking values in \mathcal{X}^n . Let $A \subset \mathcal{X}^n$. The Hamming distance of X to A is

$$d(X,A) = \min_{y \in A} d(X,y) = \min_{y \in A} \sum_{i=1}^n \mathbb{1}_{X_i \neq y_i} \; .$$



Michel Talagrand

$$\mathbb{P}\left\{\mathsf{d}(\mathsf{X},\mathsf{A})\geq\mathsf{t}+\sqrt{rac{\mathsf{n}}{2}\lograc{1}{\mathbb{P}[\mathsf{A}]}}
ight\}\leq\mathsf{e}^{-2\mathsf{t}^2/\mathsf{n}}\;.$$

Concentration of measure!

Proof: By the bounded differences inequality,

$$\begin{split} \mathbb{P}\{\mathbb{E}d(\mathsf{X},\mathsf{A})-d(\mathsf{X},\mathsf{A})\geq t\}\leq e^{-2t^2/n}.\\ \text{Taking }t=\mathbb{E}d(\mathsf{X},\mathsf{A})\text{, we get}\\ \mathbb{E}d(\mathsf{X},\mathsf{A})\leq \sqrt{\frac{n}{2}\log\frac{1}{\mathbb{P}\{\mathsf{A}\}}}. \end{split}$$

By the bounded differences inequality again,

$$\mathbb{P}\left\{\mathsf{d}(\mathsf{X},\mathsf{A})\geq\mathsf{t}+\sqrt{\frac{\mathsf{n}}{2}\log\frac{1}{\mathbb{P}\{\mathsf{A}\}}}\right\}\leq\mathsf{e}^{-2\mathsf{t}^2/\mathsf{n}}$$

talagrand's convex distance

The weighted Hamming distance is

$$\mathsf{d}_{\alpha}(\mathsf{x},\mathsf{A}) = \inf_{\mathsf{y}\in\mathsf{A}}\mathsf{d}_{\alpha}(\mathsf{x},\mathsf{y}) = \inf_{\mathsf{y}\in\mathsf{A}}\sum_{\mathsf{i}:\mathsf{x}_{\mathsf{i}}\neq\mathsf{y}_{\mathsf{i}}}|\alpha_{\mathsf{i}}|$$

where $\alpha = (\alpha_1, \ldots, \alpha_n)$. The same argument as before gives

$$\mathbb{P}\left\{\mathsf{d}_{\alpha}(\mathsf{X},\mathsf{A})\geq\mathsf{t}+\sqrt{\frac{\|\alpha\|^{2}}{2}\log\frac{1}{\mathbb{P}\{\mathsf{A}\}}}\right\}\leq\mathsf{e}^{-2\mathsf{t}^{2}/\|\alpha\|^{2}}\;,$$

This implies

 $\sup_{\alpha: \|\alpha\|=1} \min \left(\mathbb{P}\{\mathsf{A}\}, \mathbb{P}\left\{\mathsf{d}_{\alpha}(\mathsf{X},\mathsf{A}) \geq t\right\} \right) \leq e^{-t^{2}/2} \;.$

convex distance inequality

convex distance:

$$\begin{split} \mathsf{d}_\mathsf{T}(\mathsf{x},\mathsf{A}) &= \sup_{\alpha \in [0,\infty)^n : \|\alpha\| = 1} \mathsf{d}_\alpha(\mathsf{x},\mathsf{A}) \ . \\ & \mathbb{P}\{\mathsf{A}\}\mathbb{P}\left\{\mathsf{d}_\mathsf{T}(\mathsf{X},\mathsf{A}) \geq t\right\} \leq e^{-t^2/4} \ . \end{split}$$

convex distance inequality

convex distance:

$$\mathsf{d}_\mathsf{T}(\mathsf{x},\mathsf{A}) = \sup_{lpha \in [0,\infty)^n: \|lpha\| = 1} \mathsf{d}_lpha(\mathsf{x},\mathsf{A}) \;.$$

$$\mathbb{P}{A}\mathbb{P}{d_T(X, A) \ge t} \le e^{-t^2/4}$$
.

Follows from the fact that $d_T(X, A)^2$ is (4, 0) weakly self bounding (by a saddle point representation of d_T).

Talagrand's original proof was different.

convex lipschitz functions For $A \subset [0,1]^n$ and $x \in [0,1]^n$, define $D(x,A) = \inf_{y \in A} \|x - y\| \ .$

If **A** is convex, then

 $\mathsf{D}(x,\mathsf{A}) \leq \mathsf{d}_\mathsf{T}(x,\mathsf{A})$.

convex lipschitz functions For $A \subset [0,1]^n$ and $x \in [0,1]^n$, define $D(x,A) = \inf_{y \in A} ||x - y|| \ .$

If **A** is convex, then

 $\mathsf{D}(x,\mathsf{A}) \leq \mathsf{d}_\mathsf{T}(x,\mathsf{A})$.

Proof:

$$\begin{split} \mathsf{D}(\mathsf{x},\mathsf{A}) &= \inf_{\nu \in \mathcal{M}(\mathsf{A})} \|\mathsf{x} - \mathbb{E}_{\nu}\mathsf{Y}\| \quad (\mathsf{since }\mathsf{A} \text{ is convex}) \\ &\leq \inf_{\nu \in \mathcal{M}(\mathsf{A})} \sqrt{\sum_{j=1}^{n} \left(\mathbb{E}_{\nu}\mathbb{1}_{\mathsf{x}_{j} \neq \mathsf{Y}_{j}}\right)^{2}} \quad (\mathsf{since }\mathsf{x}_{j},\mathsf{Y}_{j} \in [0,1]) \\ &= \inf_{\nu \in \mathcal{M}(\mathsf{A})} \sup_{\alpha: \|\alpha\| \leq 1} \sum_{j=1}^{n} \alpha_{j} \mathbb{E}_{\nu}\mathbb{1}_{\mathsf{x}_{j} \neq \mathsf{Y}_{j}} \quad (\mathsf{by Cauchy-Schwarz}) \\ &= \mathsf{d}_{\mathsf{T}}(\mathsf{x},\mathsf{A}) \quad (\mathsf{by minimax theorem}) \; . \end{split}$$

convex lipschitz functions

Let $X = (X_1, \dots, X_n)$ have independent components taking values in [0, 1]. Let $f : [0, 1]^n \to \mathbb{R}$ be quasi-convex such that $|f(x) - f(y)| \le ||x - y||$. Then

 $\mathbb{P}\{f(X) > \mathbb{M}f(X) + t\} \leq 2e^{-t^2/4}$

and

$$\mathbb{P}\{\mathsf{f}(\mathsf{X}) < \mathbb{M}\mathsf{f}(\mathsf{X}) - \mathsf{t}\} \leq 2\mathrm{e}^{-\mathsf{t}^2/4}$$
 .

convex lipschitz functions

Let $X=(X_1,\ldots,X_n)$ have independent components taking values in [0,1]. Let $f:[0,1]^n\to \mathbb{R}$ be quasi-convex such that $|f(x)-f(y)|\leq \|x-y\|$. Then

 $\mathbb{P}\{f(X) > \mathbb{M}f(X) + t\} \leq 2e^{-t^2/4}$

and

$$\mathbb{P}\{f(\mathsf{X}) < \mathbb{M}f(\mathsf{X}) - t\} \leq 2e^{-t^2/4}$$

Proof: Let $A_s = \{x: f(x) \leq s\} \subset [0,1]^n.$ A_s is convex. Since f is Lipschitz,

$$f(x) \leq s + D(x,A_s) \leq s + d_T(x,A_s) \ ,$$

By the convex distance inequality,

$$\mathbb{P}\{f(\mathsf{X}) \geq s+t\}\mathbb{P}\{f(\mathsf{X}) \leq s\} \leq e^{-t^2/4}$$
 .

Take s = Mf(X) for the upper tail and s = Mf(X) - t for the lower tail.

Stéphane Boucheron Gábor Lugosi Pascal Massart

CONCENTRATION INEQUALITIES



OXFORD