

Machine Learning 2:

k-nearest neighbors, deviation bounds

Master 2 Computer Science

Aurélien Garivier

2018-2019



Table of contents

1. Deviation Bound for Bernoulli Variables
2. k -nearest neighbours

The result of last lecture on the nearest-neighbor classifier

A1. $\mathcal{Y} = \{0, 1\}$.

A2. $\mathcal{X} = [0, 1]^d$.

A3. η is c -Lipschitz continuous:

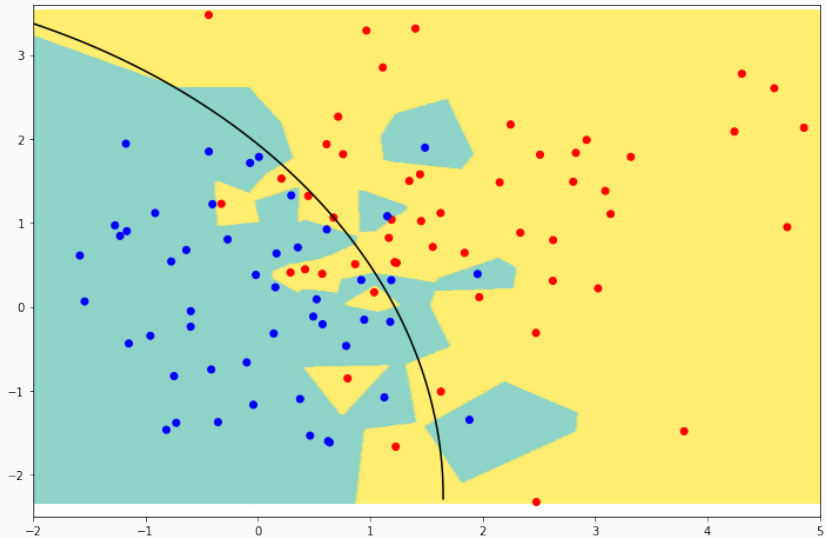
$$\forall x, x' \in \mathcal{X}, |\eta(x) - \eta(x')| \leq c \|x - x'\| .$$

Theorem

Under the previous assumptions, for all distributions D and all $m \geq 1$

$$L_D(\hat{h}_m^{NN}) \leq 2L_D^* + \frac{3c\sqrt{d}}{m^{1/(d+1)}}$$

Numerically

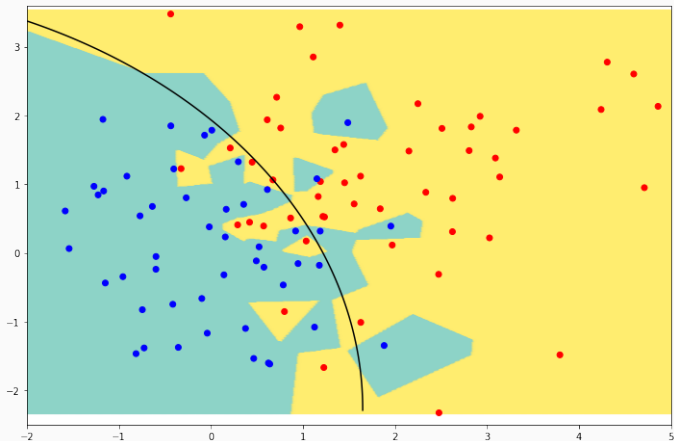


What does the analysis say?

- Where is the analysis loose? (sanity check: uniform \mathcal{D}_X)
- *finite sample* bound: explicit, non-asymptotic
- The second term $\frac{3c\sqrt{d}}{m^{1/(d+1)}}$ is *distribution-free*
- Does not give the trajectorial decreasing rate of the risk
- Exponential bound d (cannot be avoided...)
 - ⇒ *curse of dimensionality*
- Is it better than a simple grid approach?
 - ⇒ *adaptivity* to the dimension of manifold supporting data
- How to improve the classifier?
 - ⇒ k-nearest neighbors

More neighbors are better?

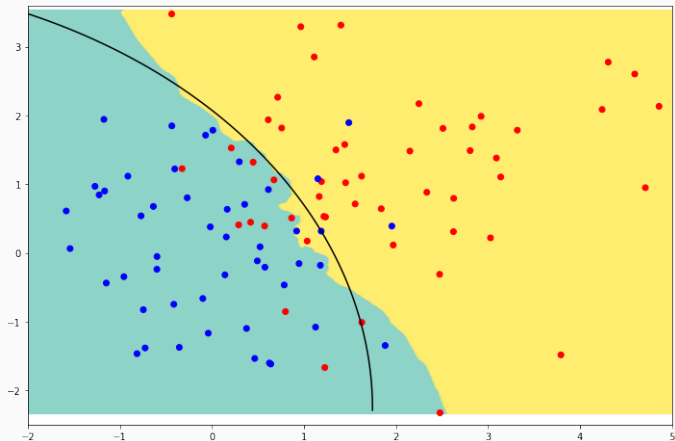
In general, yes in the sense that for m large enough, larger k is better.



But one can find counter-examples: $\forall k \geq 3, \forall m \geq k, L(\hat{h}_m^{kNN}) \geq L(\hat{h}_m^{NN})$.

More neighbors are better?

In general, yes in the sense that for m large enough, larger k is better.



But one can find counter-examples: $\forall k \geq 3, \forall m \geq k, L(\hat{h}_m^{kNN}) \geq L(\hat{h}_m^{NN})$.

Deviation Bound for Bernoulli Variables

Remember: Jensen's Inequality

Let \mathcal{X} be a convex set and $\phi : \mathcal{X} \rightarrow \mathbb{R}$ be a convex function.

Basic: For all $x, x' \in \mathcal{X}$, $\phi(tx + (1-t)x') \leq t\phi(x) + (1-t)\phi(x')$.

Probabilistic version: If $\phi : \mathcal{X} \rightarrow \mathbb{R}$ is convex and if X is a random variable with range in \mathcal{X} , then $\phi(\mathbb{E}[X]) \leq \mathbb{E}[\phi(X)]$.

Conditional version: If X and Y are random variables and the range of X is included in \mathcal{X} , if $\phi(X)$ is integrable then $\phi(\mathbb{E}[X|Y]) \leq \mathbb{E}[\phi(X)|Y]$.

Example: For a real-valued random variable X with finite expectation, $\mathbb{E}[X^2] \geq \mathbb{E}[X]^2$ and thus $\text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 \geq 0$.

Make a picture. Think about equality case.

Chernoff's Bound

Theorem (Chernoff-Hoeffding Deviation Bound)

Let $\mu \in (0, 1)$. $X_1, \dots, X_n \stackrel{iid}{\sim} \mathcal{B}(\mu)$, and let $x \in (\mu, 1]$.

(i) Chernoffs' bound for Bernoulli variables:

$$\mathbb{P}(\bar{X}_n \geq x) \leq \exp(-n \text{kl}(x, \mu)), \quad (1)$$

where $\text{kl}(p, q) = p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q}$. Same for left deviations.

(ii) If $\phi(x) = \text{kl}(x, \mu)$, then $\phi''(x) = 1/[x(1-x)]$ and

$$\begin{aligned} \text{kl}(x, \mu) &= \frac{(x - \mu)^2}{2} \int_0^1 \phi''(\mu + s(x - \mu)) 2(1-s) ds \\ &\geq \frac{(x - \mu)^2}{2\tilde{x}(1 - \tilde{x})} \quad \text{with } \tilde{x} = \frac{2\mu + x}{3} \text{ by Jensen, since } \phi'' \text{ is convex and } \int_0^1 s 2(1-s) ds = \frac{1}{3} \\ &\geq \frac{1}{2 \max_{x \leq u \leq p} u(1-u)} (x - \mu)^2 \geq 2(x - \mu)^2. \end{aligned}$$

(iii) Hoeffding's bound for Bernoulli variables:

$$\mathbb{P}(\bar{X}_n \geq x) \leq \exp(-2n(x - \mu)^2). \quad (2)$$

(iv) Inequalities (1) and (2) hold for arbitrary independent random variables with range $[0, 1]$ and expectation μ .

Examples

- If $\mu < 1/2$,

$$\mathbb{P}\left(\bar{X}_k > \frac{1}{2}\right) \leq \exp\left(-\frac{k}{2}(1-2\mu)^2\right).$$

(Consequence of Chernoff or direct computation with $(1-u)^k \leq \exp(-ku)$, or of Hoeffding).

- For all $\mu \in [0, 1]$, Chernoff's bound with $\log(u) \geq (u-1)/u$ yields

$$\mathbb{P}\left(\bar{X}_k < \frac{\mu}{2}\right) \leq \exp\left(-\frac{1-\log(2)}{2} k\mu\right) \approx \exp(-0.153 k\mu) \leq \exp\left(-\frac{k\mu}{7}\right).$$

Hoeffding yields a very poor result, but (ii) gives:

$$\mathbb{P}\left(\bar{X}_k < \frac{\mu}{2}\right) \leq \exp\left(-\frac{3}{20} k\mu\right) = \exp(-0.15 k\mu) \leq \exp\left(-\frac{k\mu}{8}\right).$$

Bennett's and Bernstein's inequalities

Let $(X_i)_{1 \leq i \leq n}$ be independent random variables upper-bounded by 1, let $\bar{\mu} = (\mathbb{E}[X_1] + \dots + \mathbb{E}[X_n])/n$, let σ^2 be such that $\mathbb{E}[X_i^2] \leq \sigma^2$ for all i and let $\phi(u) = (1+u) \log(1+u) - u$. Then, for all $x > 0$,

$$\mathbb{P}(\bar{X} \geq \bar{\mu} + x) \leq \exp\left(-n\sigma^2\phi\left(\frac{x}{\sigma^2}\right)\right) \leq \exp\left(-\frac{nx^2/2}{1+x/3}\right).$$

Bernstein from Bennett: $\phi(x) \geq \frac{x^2}{2(1+\frac{x}{3})}$ since $\psi(x) = 2(1+\frac{x}{3})\phi(x) - x^2 \geq 0$.

Extension: if $X_i \leq b$ with $b > 0$,

$$\mathbb{P}(\bar{X} \geq \bar{\mu} + x) \leq \exp\left(-\frac{n\sigma^2}{b^2}\phi\left(\frac{bx}{\sigma^2}\right)\right) \leq \exp\left(-\frac{nx^2/2}{\sigma^2 + bx/3}\right).$$

Example: for X with range in $[0, 1]$,

$$\mathbb{P}\left(\bar{X}_k < \frac{\mu}{2}\right) \leq \exp\left(-k\left(\frac{3}{2}\log\frac{3}{2} - \frac{1}{2}\right)\mu\right) \leq \exp\left(-\frac{3k\mu}{28}\right).$$

***k*-nearest neighbours**

Definition

Let \mathcal{X} be a (pre-compact) metric space with distance d .

k-NN classifier

$h^{kNN} : x \mapsto \mathbb{1}\{\hat{\eta}(x) \geq 1/2\}$ = plugin for Bayes classifier with estimator

$$\hat{\eta}(x) = \frac{1}{k} \sum_{j=1}^k Y_{(j)}(X)$$

where

$$d(X_{(1)}(X), X) \leq d(X_{(2)}(X), X) \leq \dots \leq d(X_{(m)}(X), X) .$$

Risk bound

Let \mathcal{C}_ϵ be an ϵ -covering of \mathcal{X} :

$$\forall x \in \mathcal{X}, \exists x' \in \mathcal{C}_\epsilon : d(x, x') \leq \epsilon .$$

Excess risk for k-nearest-neighbours

If η is c -Lipschitz continuous: $\forall x, x' \in \mathcal{X}, |\eta(x) - \eta(x')| \leq c d(x, x')$,
then for all $k \geq 2$ and all $m \geq 1$:

$$\begin{aligned} L(\hat{h}^{kNN}) - L(h^*) &\leq \frac{1}{\sqrt{ke}} + \frac{2k|\mathcal{C}_\epsilon|}{m} + 4c\epsilon \\ &\leq \frac{1}{\sqrt{ke}} + (2 + 4c) \left(\frac{\alpha k}{m}\right)^{\frac{1}{d+1}} \begin{cases} \text{for } \epsilon = \left(\frac{\alpha k}{m}\right)^{\frac{1}{d+1}}, \\ \text{if } |\mathcal{C}_\epsilon| \leq \alpha \epsilon^{-d} \end{cases} \\ &\leq (3 + 4c) \left(\frac{\alpha}{m}\right)^{\frac{1}{d+3}} \quad \text{for } k = \left(\frac{m}{\alpha}\right)^{\frac{2}{d+3}} . \end{aligned}$$

Sketch of the analysis

$$\begin{aligned}L(\hat{h}_m^{kNN}) - L(h^*) &= \mathbb{E} \left[|2\eta(X) - 1| \mathbb{1} \{ \hat{h}_m^{kNN} \neq h^*(x) \} \right] \\ &\leq \mathbb{P}(d(X, X_{(k)}) > 2\epsilon) + \mathbb{E} \left[|2\eta(X) - 1| \mathbb{1} \{ \hat{h}_m^{kNN} \neq h^*(x) \} \mathbb{1} \{ d(X, X_{(k)}) \leq 2\epsilon \} \right]\end{aligned}$$

- $\mathbb{P}(d(X, X_{(k)}) > 2\epsilon) \leq \sum_{c \in \mathcal{C}_\epsilon} \mathbb{P}(X \in c, N_c < k) \leq \frac{2k|\mathcal{C}_\epsilon|}{m}$

- For x such that $\eta(x) \leq 1/2 - 2c\epsilon$,

$$P(\hat{h}_m^{kNN}(x) = 1 | X = x, d(X, X_{(k)}) \leq 2\epsilon) \leq \exp\left(-\frac{k}{2}(2\eta(x) + 4c\epsilon - 1)^2\right).$$

Same for $\eta(x) \geq 1/2 + 2c\epsilon$. And for $1/2 - 2c\epsilon \leq \eta(x) \leq 1/2 + 2c\epsilon$ the probability is upper-bounded by 1. In all cases, on $\{d(X, X_{(k)}) \leq 2\epsilon\}$:

$$|2\eta(X) - 1| P(\hat{h}_m^{kNN}(X) \neq h^*(X)) \leq 4c\epsilon + \sup_{u \geq 0} u \exp(-ku^2/2) = 4c\epsilon + \frac{1}{\sqrt{ke}}.$$