

An optimal aggregation type classifier

Alejandro Cholaquidis

CMAT-Facultad de Ciencias, UdelaR

Seminario de Probabilidad y Estadística
Abril 2014

1 Introduction

- On the classification problem
 - k -NN

2 Aggregation type classifiers

- Other approaches
- Our proposal

On the classification problem

Goal:

From a training sample $\mathcal{D}_n = \{(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)\}$ i.i.d. of (X, Y) with values in $\mathcal{F} \times \{0, 1\}$ we want a predictor $g : \mathcal{F} \rightarrow \{0, 1\}$, which *minimize* $\mathbb{P}(g(X) \neq Y)$.

Let us denote

- 1) $\eta(x) = \mathbb{E}(Y|X = x) = \mathbb{P}(Y = 1|X = x)$ the regression function.
- 2) $g^*(x) = \mathbb{I}_{\{\eta(x) > 1/2\}}$ the Bayes rule.
- 3) $L^* = \mathbb{P}(g^*(X) \neq Y)$ the optimal Bayes risk.

The classical estimator of $\eta(x)$ is

$$\eta_n(X) = \sum_{i=1}^n W_{ni}(X) Y_i,$$

for some weights $W_{ni}(X) = W_{ni}(X, X_1, \dots, X_n)$.

On the classification problem: Stone Theorem

Theorem

Stone 1977: If $\mathcal{F} = \mathbb{R}^d$, and the weights fulfills:

- 1) There exists a constant $C \geq 1$ such that for every nonnegative Borel function

$$\mathbb{E}\left(\sum_{i=1}^n |W_{ni}|f(X_i)\right) \leq C\mathbb{E}(f(X)) \quad \forall n \geq 1$$

- 2) there exists a $D \geq 1$ such that $\mathbb{P}\left(\sum_{i=1}^n |W_{ni}(X)| \leq D\right) = 1 \quad \forall n \geq 1$

- 3) $\sum_{i=1}^n |W_{ni}(X)|\mathbb{I}_{\{\|X_i - X\| > a\}} \rightarrow 0$ in probability, $\forall a > 0$

- 4) $\sum_{i=1}^n W_{ni}(X) \rightarrow 1$ in probability, and

- 5) $\max_i |W_{ni}(X)| \rightarrow 0$

then

$$\mathbb{E}\left(\mathbb{P}\left(\mathbb{I}_{\{\eta_n(X) \geq 1/2\}} \neq Y \mid D_n\right)\right) \rightarrow \mathbb{P}\left(\mathbb{I}_{\{\eta(X) \geq 1/2\}} \neq Y\right).$$

1 Introduction

- On the classification problem
- k -NN

2 Aggregation type classifiers

- Other approaches
- Our proposal

k-NN and the infinite dimensional case

k-NN

For $W_{ni}(X) = \frac{1}{k} \mathbb{I}_{\{X_i \in k_n(X)\}}$, where $X_i \in k_n(X)$ if X_i is one of the k nearest neighbours of X , conditions 1 to 5 holds if $n \rightarrow \infty$ and $k_n/n \rightarrow 0$.

Theorem

Let $(X, Y) \in \mathcal{F} \times \{0, 1\}$ where (\mathcal{F}, d) is a separable metric space. Suppose that the probability measure μ of X is a Borel measure. If η satisfies the Besicovitch Condition:

$$\lim_{\delta \rightarrow 0} \mu \left\{ x \in \mathcal{F} : \frac{1}{\mu(B(x, \delta))} \int_{B(x, \delta)} |\eta(y) - \eta(x)| d\mu(y) > \varepsilon \right\} = 0, \quad (1)$$

then, taking $k_n \rightarrow \infty$ and $k_n/n \rightarrow 0$, k-NN is consistent:

$$\mathbb{E} \left(\mathbb{P} \left(\mathbb{I}_{\{\eta_n(X) \geq 1/2\}} \neq Y | \mathcal{D}_n \right) \right) \rightarrow \mathbb{P} \left(\mathbb{I}_{\{\eta(X) \geq 1/2\}} \neq Y \right).$$

On the Besicovitch condition

If (\mathcal{F}, d) is complete and separable, Besicovitch condition holds for all measurable function $f : \mathcal{F} \rightarrow \mathbb{R}$ if and only if:

$$\lim_{\delta \rightarrow 0} \mu \left\{ x \in \mathcal{F} : \left| \frac{1}{\mu(B(x, \delta))} \int_{B(x, \delta)} f(y) d\mu(y) - f(x) \right| > \varepsilon \right\} = 0. \quad (2)$$

Condition (2) holds if any of the following conditions holds:

- f is a continuous function.
- μ is doubling: $\exists C \geq 1$ such that $\forall r > 0, \mu(B(x, 2r)) \leq C\mu(B(x, r))$.
- (\mathcal{F}, d) is a Vitali space: for every \mathcal{G} collection of closed balls such that for all $a \in \mathcal{F}$, $\inf\{r > 0 : B(a, r) \in \mathcal{G}\} = 0$ then there exists a disjoint covering of \mathcal{F} with balls in \mathcal{G} , a.s. with respect to μ .
- \mathcal{F} is finite dimensional: $\exists n \in \mathbb{N}$ and $s \in (0, +\infty) : \bigcap_{x \in \mathcal{G}} B(x, r(x)) = \emptyset$,

$\forall \mathcal{G} \subset \mathcal{F}$ finite set such that:

$$1) \#\mathcal{G} > n \quad 2) r(x) \in (0, s) \quad 3) x \notin B(y, r(y)) \quad \forall x \neq y \in \mathcal{G}$$

1 Introduction

- On the classification problem
- k -NN

2 Aggregation type classifiers

- **Other approaches**
- Our proposal

Mojirsheibani 1999

- 1) With \mathcal{D}_n build up $g_{ni} : \mathbb{R}^d \rightarrow \{0, 1\}$, $i = 1, \dots, M$
- 2) Taking $\mathbf{g}_n(x) \doteq (g_{n1}(x), \dots, g_{nM}(x))$, we define the classifier:

$$g_T(x) = \mathbb{I}_{\{T_n(\mathbf{g}_n(x)) > 1/2\}},$$

where

$$T_n(\mathbf{g}_n(x)) = \sum_{j=1}^n W_{n,j}(x) Y_j, \quad x \in \mathcal{F}, \quad (3)$$

and the weights $W_{n,j}(x)$ given by

$$W_{n,j}(x) = \frac{\mathbb{I}_{\{\mathbf{g}_n(x) = \mathbf{g}_n(X_j)\}}}{\sum_{j=1}^n \mathbb{I}_{\{\mathbf{g}_n(x) = \mathbf{g}_n(X_j)\}}}. \quad (4)$$

Here, $0/0$ is assumed to be 0.

Mojirsheibani 1999, consistency results

Notation:

- Taking $\mathcal{T}_n \in (\mathbb{R}^d \times \{0, 1\})^n$, the sets: $\Pi_n(\mathcal{T}_n) = \{A_{n,1}, \dots, A_{n,2^M}\}$ define a partition of \mathbb{R}^d where

$$A_{n,i} = \left\{ x \in \mathbb{R}^d : (g_{n,1}(x), \dots, g_{n,M}(x)) = \nu_i \text{ where } \nu_i \in \{0, 1\}^M \right\}.$$

The family of all partitions induced by $g_{n,1}, \dots, g_{n,M}$ is:

$$\Omega_n = \bigcup_{\mathcal{T}_n \in (\mathbb{R}^d \times \{0,1\})^n} \Pi_n(\mathcal{T}_n).$$

- Let $B = \{x_1, \dots, x_n\}$, define $\Delta_n(\Omega_n)$ the number of distinct partitions:

$$\left\{ \Pi_n(\mathcal{T}_n^1) \cap B, \dots, \Pi_n(\mathcal{T}_n^r) \cap B \right\},$$

of the finite set B that are induced by partitions $\{\Pi_n(\mathcal{T}_n^1), \dots, \Pi_n(\mathcal{T}_n^r)\} \in \Omega_n$.

Mojirsheibani 1999, consistency results

Theorem

If, as $n \rightarrow \infty$,

- $M_n / \log(n) \rightarrow 0$
- $n^{-1} \log(\Delta_n) \rightarrow 0$
- for every $\gamma > 0$ and $\delta \in (0, 1)$

$$\inf_{S \subset \mathbb{R}^d: \mu(S) \geq 1-\delta} \mu \left\{ x : \text{diam}(A_n[x] \cap S) > \gamma \right\} \rightarrow 0 \quad \text{a.s.}$$

if we denote $g_T^*(x) = \mathbb{I}_{\{\mathbb{P}(Y=1|g_n(x)) > 1/2\}}$ the optimal combined classifier, then

$$\mathbb{P}(g_T(X) \neq Y | \mathcal{D}_n) - \mathbb{P}(g_T^*(X) \neq Y | \mathcal{D}_n) \rightarrow 0 \quad \text{a.s.}$$

1 Introduction

- On the classification problem
- k -NN

2 Aggregation type classifiers

- Other approaches
- **Our proposal**

The classifier

- 1) Split \mathcal{D}_n into
 - $\mathcal{D}_k = \{(X_1, Y_1), \dots, (X_k, Y_k)\}$
 - $\mathcal{E}_l = \{(X_{k+1}, Y_{k+1}), \dots, (X_n, Y_n)\}$ with $l = n - k \geq 1$.
- 2) With \mathcal{D}_k build up $g_{ki} : \mathcal{F} \rightarrow \{0, 1\}$, $i = 1, \dots, M$
- 3) Taking $\mathbf{g}_k(x) \doteq (g_{k1}(x), \dots, g_{kM}(x))$, we define the classifier:

$$gT(x) = \mathbb{I}_{\{T_n(\mathbf{g}_k(x)) > 1/2\}},$$

where

$$T_n(\mathbf{g}_k(x)) = \sum_{j=k+1}^n W_{n,j}(x) Y_j, \quad x \in \mathcal{F}, \quad (5)$$

and the weights $W_{n,j}(x)$ given by

$$W_{n,j}(x) = \frac{\mathbb{I}_{\{\mathbf{g}_k(x) = \mathbf{g}_k(X_j)\}}}{\sum_{j=k+1}^n \mathbb{I}_{\{\mathbf{g}_k(x) = \mathbf{g}_k(X_j)\}}}. \quad (6)$$

Here, $0/0$ is assumed to be 0.

The classifier, example

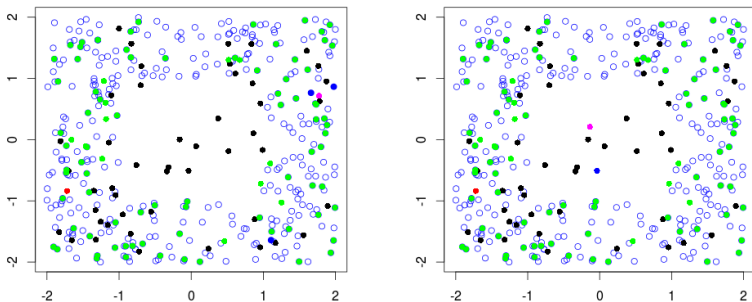


Figure : An example in \mathbb{R}^2 where the 0-class points are shown in black dots and the 1-class points are shown in blue circles. In both pictures the voters of the point in red are shown as green dots while the voters of the point in magenta are shown with blue dots.

A more flexible version

For $0 \leq \alpha < 1$ a more flexible version of the classifier, called $g_T(x, \alpha)$, can be defined replacing the weights in (6) by

$$W_{n,j}(x) = \frac{\mathbb{I}_{\{\frac{1}{M} \sum_{i=1}^M \mathbb{I}_{\{g_{ki}(x)=g_{ki}(X_j)\}} \geq 1-\alpha\}}}{\sum_{j=k+1}^n \mathbb{I}_{\{\frac{1}{M} \sum_{i=1}^M \mathbb{I}_{\{g_{ki}(x)=g_{ki}(X_j)\}} \geq 1-\alpha\}}} . \quad (7)$$

Observe that if we choose $\alpha = 0$ in (7) we obtain original weights given in (6).

Consistency

Theorem

Consistency: Suppose that for every $i = 1, \dots, M$ the classifier g_{ki} converges in probability to g^* as $k \rightarrow \infty$. If $\mathbb{P}(Y = 1|g^*(X) = 1) > 1/2$ and $\mathbb{P}(Y = 0|g^*(X) = 0) > 1/2$, then

$$\lim_{k \rightarrow +\infty} \lim_{l \rightarrow +\infty} \mathbb{P}_{\mathcal{D}_k}(g_T(X) \neq Y) - L^* = 0.$$

Theorem

Optimality: Let $\mathbb{C} \doteq \{0, 1\}^M$. For $\nu \in \mathbb{C}$ we define:

$$A_\nu^0 \doteq \bigcap_{i=1}^M g_{ki}^{-1}(\nu(i)) \cap \{Y = 0\}, \quad A_\nu^1 \doteq \bigcap_{i=1}^M g_{ki}^{-1}(\nu(i)) \cap \{Y = 1\}$$

$\forall \nu \in \mathbb{C}$, we assume that $\mathbb{P}_{\mathcal{D}_k}((X, Y) \in A_\nu^1) \neq \mathbb{P}_{\mathcal{D}_k}((X, Y) \in A_\nu^0)$ a.s., then

$$\lim_{l \rightarrow \infty} \mathbb{P}_{\mathcal{D}_k}(g_T(X) \neq Y) \leq \mathbb{P}_{\mathcal{D}_k}(g_{ki}(X) \neq Y),$$

for each $i = 1, \dots, M$. Therefore, $\lim_{l \rightarrow \infty} \mathbb{P}_{\mathcal{D}_k}(g_T(X) \neq Y) \leq \min_{1 \leq i \leq M} \mathbb{P}_{\mathcal{D}_k}(g_{ki}(X) \neq Y)$.

Generating the sample \mathcal{D}_n

- 1) Generate Z_1, \dots, Z_N iid uniform in $[0, 1]$,
 - if $Z_i > p$ we generate $X_i \sim \mu_1$ and $Y_i = 1$,
 - if $Z_i \leq p$ we generate $X_i \sim \mu_0$ and set $Y_i = 0$.
- 2) The first $n = N/2$ pairs (X_i, Y_i) for the training sample, the remaining for the testing sample.
- 3) We consider M nearest neighbour classifiers with the numbers of neighbours taken as follows:
 - We fix M consecutive odd numbers.
 - We choose at random M different odd integers between 1 and $\min\{\sum_{i=1}^k Y_i, k - \sum_{i=1}^k Y_i\}$.

In \mathbb{R}^{150}

- $M = 8$ neighbour rules with $k = 5, 7, 9, 11, 13, 15, 17, 19$
- $p = 1/6$
- $\mu_1 \sim U([-2, 2]^{150})$
- $\mu_0 \sim U(\tau_v([-2, 2]^{150}))$ where τ_v is the translation along $(v, \dots, v) \in \mathbb{R}^{150}$, for $v = 1/4$.

n (k)	$g_T(\cdot)$	$g_T(\cdot, 1/4)$	g_{k1}	g_{k2}	g_{k3}	g_{k4}	g_{k5}	g_{k6}	g_{k7}	g_{k8}
400 (300)	.046	.056	.071 (.074)	.067 (.072)	.066 (.072)	.067 (.073)	.068 (.074)	.069 (.077)	.071 (.080)	.073 (.082)
600 (400)	.043	.052	.067 (.072)	.062 (.069)	.061 (.068)	.061 (.068)	.061 (.069)	.062 (.071)	.063 (.073)	.065 (.076)
800 (600)	.037	.045	.062 (.066)	.057 (.061)	.055 (.060)	.055 (.060)	.055 (.060)	.056 (.061)	.056 (.062)	.057 (.064)
1000 (700)	.035	.043	.061 (.065)	.055 (.060)	.053 (.058)	.052 (.057)	.052 (.057)	.052 (.058)	.053 (.059)	.054 (.060)

Table : \mathbb{R}^{150} with fixed number of neighbours

In \mathbb{R}^{150}

We take at random $M = 10$ values for k , and compare with cross validated nearest neighbour classifier.

n	k	$g_T(\cdot)$	$g_T(\cdot, 1/8)$	$g_T(\cdot, 1/4)$	gcv_n	gcv_k
400	300	.052	.065	.077	.068	.073
600	400	.049	.063	.074	.061	.068
800	500	.048	.062	.073	.056	.061
1000	700	.047	.061	.072	.053	.058









Table : \mathbb{R}^{150} with the number of neighbours chosen at random, compared with cross validation

FDA

- $M = 6$ neighbour rules with $k = 3, 9, 15, 21, 27, 33$
- $p = 1/8$
- if $Z_i < p$ then $X_i = B(t)$ where $B(t)$ is the Brownian bridge on a grid of length 100, and $Y_i = 0$
- if $Z_i \geq p$, $X_i = B(t) + 4t(1 - t)$ and $Y_i = 1$.

n (k)	$g_T(\cdot)$	$g_T(\cdot, 1/4)$	g_{k1}	g_{k2}	g_{k3}	g_{k4}	g_{k5}	g_{k6}
400 (300)	.074	.072	.079 (.079)	.071 (.072)	.070 (.070)	.069 (.070)	.070 (.071)	.070 (.072)
600 (400)	.072	.070	.080 (.080)	.071 (.071)	.070 (.070)	.069 (.070)	.069 (.070)	.069 (.070)
800 (500)	.071	.069	.080 (.079)	.071 (.071)	.070 (.070)	.069 (.069)	.069 (.069)	.068 (.069)
1000 700	.071	.068	.080 (.079)	.071 (.071)	.069 (.069)	.069 (.069)	.068 (.068)	.068 (.068)

Table : Functional Data with fixed number of neighbours and the values of g_{ki}

-  Baíllo, A.; Cuevas, A. y Fraiman, R. (2010) - *Classification methods for functional data*. The Oxford Handbook of Functional Data Analysis.
-  Billingsley, P. (1999) *Convergence of probability measures, s.e.*. Wiley.
-  Cérou, F. y Guyader, A. (2005) *Nearest neighbor classification in infinite dimension*. Rapport de Recherche 5536, INRIA.
-  Devroye, L.; Györfi, L. y Lugosi, G. (1996) *A probabilistic theory of pattern recognition.*. Springer–Verlag, New York.
-  Forzani, L.; Fraiman, R. y Llop, P. (2012) *Consistent nonparametric regression for functional data under the Stone-Besicovitch conditions*
-  Mojirsheibani, M. (1999) Combining classifiers via discretization. *JASA*, **94**(446), 600–609.
-  Mojirsheibani, M. (2002) An almost surely optimal combined classification rule. *J. Multivariate Anal.* **81**(1), 28–46.
-  Stone, C. (1977) *Consistent Nonparametric Regression*. The Annals of Statistics.