

Classification

Bernhard Stankewitz
bernhard.stankewitz@posteo.de

March 5, 2019

Abstract

Summary of elementary results for classification.

Contents

1	Elementary definitions and the Bayes classifier	1
2	The KNN-classifier	2

1 Elementary definitions and the Bayes classifier

Definition 1.1 (Classifier).

- (a) For i.i.d. training data $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathbb{R}^d \times \{1, \dots, K\}$, a classifier is a measurable function $C : \mathbb{R}^d \rightarrow \{1, \dots, K\}$. The classification error is given by

$$R(C) := \mathbb{P}\{C(X) \neq Y\} = \mathbb{E}\mathbf{1}_{C(X) \neq Y}.$$

- (b) In case that the labels are given by $\{0, 1\}$, the classification error

$$R(C) = \mathbb{E}(Y - C(X))^2 = \int \mathbf{1}_{y \neq C(x)} \mathbb{P}^{X,Y}(d(x, y)).$$

In case all theoretical quantities are known, a classification problem has an optimal solution.

Proposition 1.2 (Bayes-Classifier).

- (i) In the situation of Definition 1.1, the classification error is minimised by the Bayes classifier

$$C^{Bayes}(x) := \operatorname{argmax}_{k=1, \dots, K} \mathbb{P}\{Y = k | X = x\}.$$

- (ii) If the labels are given by $\{0, 1\}$, we have

$$C^{Bayes}(x) = \mathbf{1}_{\{\eta(x) \geq 1/2\}} \quad \text{with} \quad \eta(x) := \mathbb{P}\{Y = 1 | X = x\}.$$

Proof. For any classifier C , we have

$$R(C) = 1 - \mathbb{E}\mathbb{E}\mathbf{1}_{C=Y} | X = 1 - \mathbb{E} \sum_{k=1}^K \mathbb{E}(\mathbf{1}_{C=k} \mathbf{1}_{Y=k} | X) = 1 - \mathbb{E} \sum_{k=1}^K \mathbf{1}_{C=k} \mathbb{E}(\mathbf{1}_{Y=k} | X).$$

□

2 The KNN-classifier

Definition 2.1 (KNN-classifier).

- (a) Let $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathbb{R}^d \times \{1, \dots, J\}$ be a training sample and $K \in \mathbb{N}$. For $x \in \mathbb{R}^d$, let $N_K(x)$ be the set of the K nearest neighbours of x with respect to the euclidean distance. Then, the KNN-classifier is given by

$$\hat{C}^{\text{KNN}}(x) := \operatorname{argmax}_{j=1, \dots, J} \frac{1}{K} \sum_{X_i \in N_K(x)} \mathbf{1}_{Y_i=j}.$$

- (b) In case, the labels are given by $\{0, 1\}$, we have

$$\hat{C}^{\text{KNN}}(x) = \mathbf{1}\{\hat{\eta}(x) \geq 1/2\} \quad \text{with} \quad \hat{\eta}(x) := \frac{1}{K} \sum_{X_i \in N_K(x)} \mathbf{1}_{Y_i=1} =: \sum_{i=1}^n w_i(x) Y_i,$$

where $w_i := \mathbf{1}_{X_i \in N_K(x)} / K$ with $\sum_{i=1}^n w_i = 1$.

Lemma 2.2 (Reduction to the regression function). *In the situation of Definition 2.1 (b), we have that*

$$|\mathbb{E}_{\leq n} R(\hat{C}^{\text{KNN}}) - R(C^{\text{Bayes}})| \leq 2\sqrt{\mathbb{E}_{n+1} |\hat{\eta}(X) - \eta(X)|^2}$$

Proof. For any classifier C , we have

$$\mathbb{P}\{C(X) = Y|X\} = \mathbf{1}_{C=1}\eta + \mathbf{1}_{C=0}(1 - \eta) = \eta + \mathbf{1}_{C=0}(1 - 2\eta).$$

This yields

$$\begin{aligned} |\mathbb{P}\{\hat{C}^{\text{KNN}}(X) \neq Y|X\} - \mathbb{P}\{\hat{C}^{\text{Bayes}}(X) \neq Y|X\}| &= |\mathbb{P}\{\hat{C}^{\text{KNN}}(X) = Y|X\} - \mathbb{P}\{\hat{C}^{\text{Bayes}}(X) = Y|X\}| \\ &= |\mathbf{1}_{\hat{C}^{\text{KNN}}=0}(1 - 2\eta) - \mathbf{1}_{\hat{C}^{\text{Bayes}}=0}(1 - 2\eta)| = \mathbf{1}_{\hat{C}^{\text{KNN}} \neq \hat{C}^{\text{Bayes}}} |1 - 2\eta| \leq 2|\hat{\eta} - \eta|. \end{aligned}$$

For the last step, use that either $\hat{\eta} \geq 1/2 > \eta$ or reverse. For the second to last step, use that one of η and $\hat{\eta}$ has to be above and one below $1/2$. By conditioning on X and Jensen's inequality, this gives

$$|\mathbb{E}_{\leq n} R(\hat{C}^{\text{KNN}}) - R(C^{\text{Bayes}})|^2 = |\mathbb{E}_{\leq n+1} (\mathbf{1}_{\hat{C}^{\text{KNN}} \neq Y} - \mathbf{1}_{C^{\text{Bayes}} \neq Y})|^2 \leq 4\mathbb{E}_{\leq n+1} |\hat{\eta} - \eta|^2.$$

□

Theorem 2.3 (Consistency of KNN). *In the situation of Definition 2.1 (b), let $k \rightarrow \infty$, $k/n \rightarrow 0$ and let $x \mapsto \eta(x)$ be uniformly continuous. Then, the KNN-classifier \hat{C}^{KNN} is consistent, i.e.*

$$|\mathbb{E}_{\leq n} R(\hat{C}^{\text{KNN}}) - R(C^{\text{Bayes}})| \xrightarrow{n \rightarrow \infty} 0.$$

Proof. From Lemma 2.2, we obtain with the triangle inequality

$$\begin{aligned} |\mathbb{E}_{\leq n} R(\hat{C}^{\text{KNN}}) - R(C^{\text{Bayes}})|/2 &\leq \sqrt{\mathbb{E}_{\leq n+1} |\hat{\eta}(X) - \eta(X)|^2} = \sqrt{\mathbb{E}_{\leq n+1} \left| \sum_{i=1}^n w_i(X) (Y_i - \eta(X)) \right|^2} \\ &\leq \sqrt{\mathbb{E}_{\leq n+1} \left| \sum_{i=1}^n w_i(X) (Y_i - \eta(X_i)) \right|^2} + \sqrt{\mathbb{E}_{\leq n+1} \left| \sum_{i=1}^n w_i(X) (\eta(X_i) - \eta(X)) \right|^2} \end{aligned}$$

In the following, we consider the two terms separately.

For the first term, we have by independence

$$\begin{aligned} \mathbb{E}_{\leq n+1} \left| \sum_{i=1}^n w_i(X) (Y_i - \eta(X_i)) \right|^2 &= \mathbb{E}_{\leq n+1} \sum_{i=1}^n w_i(X)^2 (Y_i - \eta(X_i))^2 \\ &\leq \mathbb{E}_{\leq n+1} \left(\max_{i \leq n} w_i(X) \sum_{i=1}^n w_i(X) \right) \leq 1/K \xrightarrow{n \rightarrow \infty} 0. \end{aligned}$$

For the second term, we have for any $\varepsilon > 0$ and corresponding $\delta > 0$, that

$$\begin{aligned} \mathbb{E}_{\leq n+1} \left| \sum_{i=1}^n w_i(X) (\eta(X_i) - \eta(X)) \right|^2 &\leq \mathbb{E}_{\leq n+1} \left| \sum_{i=1}^n w_i(X) \mathbf{1}_{|X_i - X| \geq \delta} + \varepsilon \right|^2 \\ &\leq \mathbb{E}_{\leq n+1} \frac{1}{K} \sum_{i=1}^K w_i(X) \mathbf{1} \left\{ \sum_{i=1}^n \mathbf{1}_{|X_i - X| < \delta} \leq i \right\} \leq \mathbb{P} \left\{ \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{|X_i - X| < \delta} \leq K/n \right\} + \varepsilon. \end{aligned}$$

The last term is eventually smaller than 2ε by dominated convergence.

□