

# Machine Learning 2:

## k-nearest neighbors, deviation bounds

Master 2 Computer Science

---

Aurélien Garivier

2019-2020



# 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.**  $\eta$  is  $c$ -Lipschitz continuous:

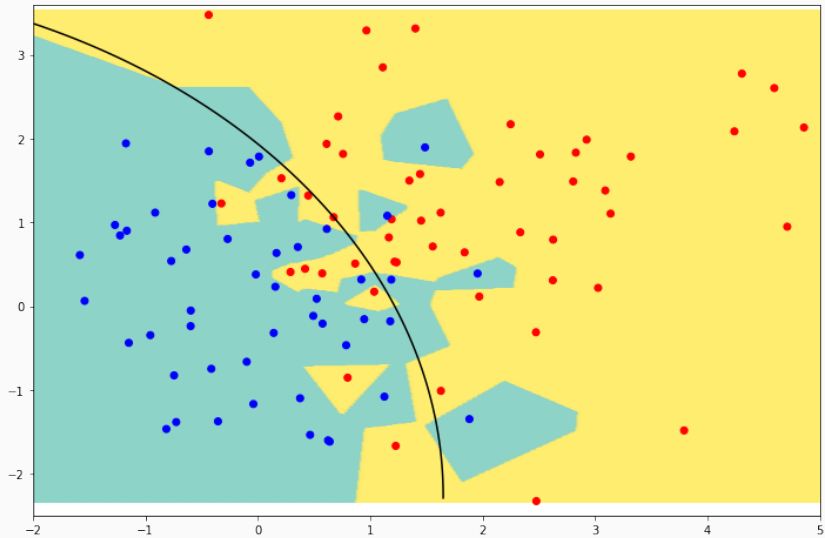
$$\forall x, x' \in \mathcal{X}, |\eta(x) - \eta(x')| \leq c \|x - x'\| .$$

## Theorem

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

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

# Numerically

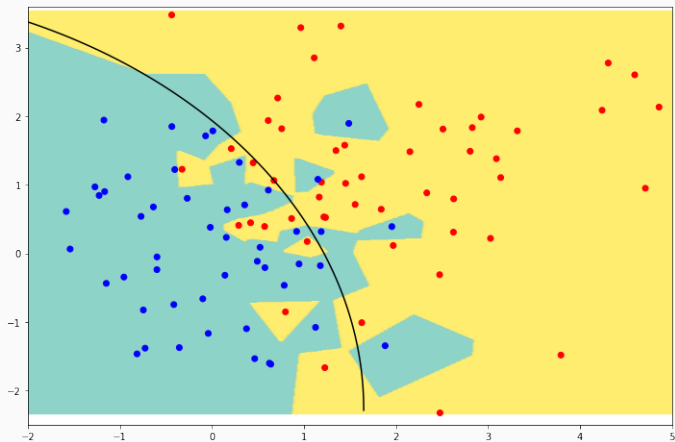


## What does the analysis say?

- 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...)
  - $\implies$  *curse of dimensionality*
- Is it better than a simple grid approach?
  - $\implies$  *adaptivity* to the dimension of manifold supporting data
- How to improve the classifier?
  - $\implies$  k-nearest neighbors

# More neighbors are better?

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

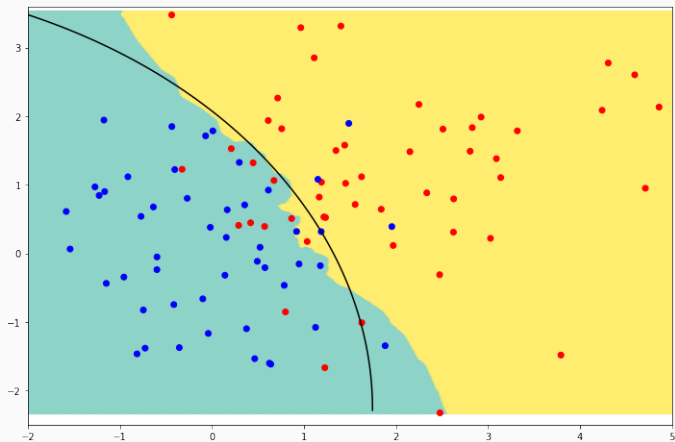


But one can find counterexamples:  $\forall k \geq 3, \forall m \geq k,$

$$R_m(\hat{h}_m^{kNN}) \geq R_m(\hat{h}_m^{NN}).$$

# More neighbors are better?

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



But one can find counterexamples:  $\forall k \geq 3, \forall m \geq k,$

$$R_m(\hat{h}_m^{kNN}) \geq R_m(\hat{h}_m^{NN}).$$

# 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) Chernoff's 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  $\mu$ .

## 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

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

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

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

# Sub-Gaussian inequalities

## 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  $\sigma^2$  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  $x > 0$ ,

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

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  $b > 0$ ,

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

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

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

# Parenthesis: a nice proof for the technicalities of Bernstein

From [Pollard, MiniEmpirical ex.14, <http://www.stat.yale.edu/~pollard/Books/Mini/Basic.pdf>]

For any sufficiently smooth real-valued function  $g$  defined at least in a neighborhood of 0 let

$$G(x) = \frac{g(x) - g(0) - xg'(0)}{x^2/2} \text{ if } x \neq 0, \text{ and } G(0) = g''(0) .$$

By Taylor's integral formula

$$g(x) - g(0) - xg'(0) = \int_0^x g''(u)(x-u)du = x^2 \int_0^1 g''(sx)(1-s)ds .$$

Thus,  $G(x) = \int g''(sx)d\nu(s)$ , where  $d\nu(s) = 2(1-s)\mathbb{1}\{0 \leq s \leq 1\}ds$ .

Hence, if  $g$  is convex then  $g'' \geq 0$  and  $G \geq 0$ . Moreover, if  $g''$  is increasing then the functions  $x \mapsto g''(sx)$  for  $s \in [0, 1]$  are all increasing and  $G$  is also increasing as an average of increasing functions. For  $g(u) = \exp(u)$ , this yields that  $(\exp(u) - u - 1)/u^2$  is increasing, as required for the proof of Bernstein's inequality.

Similarly, if  $g''$  is convex then  $G$  is also convex as an average of convex functions  $(x \mapsto g''(sx))$ . Moreover, by Jensen's inequality applied to convex function  $\psi(s) = g''(xs)$  with the probability measure  $d\nu(s) = 2(1-s)\mathbb{1}\{0 \leq s \leq 1\}ds$

$$G(x) = \int_0^1 g''(xs) 2(1-s)ds \geq g'' \left( x \int_0^1 s \times 2(1-s)ds \right) = g'' \left( \frac{x}{3} \right) .$$

For  $g(u) = (1+u) \log(1+u) - u$ ,  $g''(u) = 1/(1+u)$  and this yields:

$$\frac{g(u)}{u^2/2} \geq g'' \left( \frac{u}{3} \right) = \frac{1}{1+u/3} .$$

## Exercise: for $X_i \stackrel{iid}{\sim} \mathcal{B}(\mu)$ , $\mathbb{P}(\bar{X}_m \geq 2\mu) \leq \exp(-m \times ?)$

**Chernoff + Taylor:** since  $\log(u) \geq (u - 1)/u$ ,

$$\text{kl}(2\mu, \mu) = 2\mu \log(2) + (1 - 2\mu) \log \frac{1 - 2\mu}{1 - 2\mu} \geq 2\mu \log(2) - \mu = \mu(2 \log(2) - 1) \approx 0.386 \mu .$$

**Chernoff with convexity:**

$$\text{kl}(2\mu, \mu) \geq \frac{(2\mu - \mu)^2/2}{4/3\mu} = \frac{3}{8} \mu = 0.375 \mu .$$

**Improved Hoeffding:**

$$\text{kl}(2\mu, \mu) \geq \frac{(2\mu - \mu)^2/2}{\max_{\mu \leq u \leq 2\mu} u(1-u)} \geq \frac{\mu^2/2}{2\mu} = \frac{1}{4} \mu = 0.25 \mu .$$

**Bennett:**

$$2\mu \log \frac{2\mu}{\mu} - (2\mu - \mu) = \mu(2 \log(2) - 1) \approx 0.386 \mu .$$

**Bernstein:**

$$\frac{(2\mu - \mu)^2/2}{\mu(1 - \mu) + (2\mu - \mu)/3} \geq \frac{\mu^2/2}{\mu + \mu/3} \frac{3}{8} \mu = 0.375 \mu .$$

**Hoeffding:**  $2(2\mu - \mu)^2 = 2\mu^2$ , very poor (as expected) when  $\mu$  is small.

## ***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

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

where  $\Sigma_x$  is a random permutation defined by:

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



# Risk bound

Let  $\mathcal{C}_\epsilon$  be an  $\epsilon$ -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  $\eta$  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  $m \geq 1$ :

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

Bias-variance decomposition of the risk.

# Sketch of the analysis

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

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

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

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

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  $\{d(X, X_{\Sigma_X(k)}) \leq 2\epsilon\}$ :

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

# Room for improvement

- Lower bound? in  $m^{-\frac{1}{d}}$ .
- Margin conditions
  - ⇒ fast rates
- More regularity?
  - ⇒ weighted nearest neighbors
- Is regularity required everywhere?
  - ⇒ What matters are the balls of mass  $\approx k/m$  near the decision boundary.

## Classification in general finite dimensional spaces with the $k$ -nearest neighbor rule

by Sébastien Gadat, Thierry Klein, and Clément Marteau

Annals of Statistics Volume 44, Number 3 (2016), 982-1009.

arXiv: arXiv:1506.01171v1

### CLASSIFICATION WITH THE NEAREST NEIGHBOR RULE IN GENERAL FINITE DIMENSIONAL SPACES

BY SÉBASTIEN GADAT AND THIERRY KLEIN AND CLÉMENT MARTEAU

*Toulouse School of Economics, Université Toulouse 1 Capitole  
Institut Mathématiques de Toulouse, Université Paul Sabatier*

Given an  $n$ -sample of random vectors  $(X_i, Y_i)_{1 \leq i \leq n}$  whose joint law is unknown, the long-standing problem of supervised classification aims to *optimally* predict the label  $Y$  of a given a new observation  $X$ . In this context, the nearest neighbor rule is a popular flexible and intuitive method in non-parametric situations. Even if this algorithm is commonly used in the machine learning and statistics communities, less is known about its prediction ability in general finite dimensional spaces, especially when the support of the density of the observations is  $\mathbb{R}^d$ . This paper is devoted to the study of the statistical properties of the nearest neighbor rule in various situations. In particular, attention is paid to the marginal law of  $X$ , as well as the smoothness and margin properties of the regression function  $\eta(X) = \mathbb{E}[Y|X]$ . We identify two necessary and sufficient conditions to obtain uniform consistency rates of classification and to derive sharp estimates in the case of the nearest neighbor rule. Some numerical experiments are proposed at the end of the paper to help illustrate the discussion.

**1. Introduction.** The supervised classification model has been at the core of numerous contributions to statistical literature in recent years. It continues to provide interesting problems, both from the theoretical and practical point of views. The classical task in supervised classification is to predict a feature  $Y \in \mathcal{M}$  when a variable of interest  $X \in \mathbb{R}^d$  is observed, the set  $\mathcal{M}$  being finite. In this paper, we focus on the binary classification problem where  $\mathcal{M} = \{0, 1\}$ .

In order to provide a prediction of the label  $Y$  of  $X$ , it is assumed that a training set  $\mathcal{S}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$  is at our disposal, where  $(X_i, Y_i)$  are i.i.d. and with a common law  $\mathbb{P}_{X,Y}$ . This training set  $\mathcal{S}_n$  makes it possible to retrieve some information on the joint law of  $(X, Y)$  and to provide, depending on some technical conditions, a pertinent prediction. In particular,

*AMS 2000 subject classifications:* Primary 62G05; secondary 62G20

*Keywords and phrases:* Supervised classification, nearest neighbor algorithm, plug in rules, minimax classification rates

## Rates of convergence for nearest neighbor classification

by Kamalika Chaudhuri and Sanjoy Dasgupta

Advances in Neural Information Processing Systems 27 (NIPS 2014)

<https://papers.nips.cc/paper/5439-rates-of-convergence-for-nearest-neighbor-classification>

### Rates of convergence for nearest neighbor classification

Kamalika Chaudhuri  
Computer Science and Engineering  
University of California, San Diego  
kamalika@cs.ucsd.edu

Sanjoy Dasgupta  
Computer Science and Engineering  
University of California, San Diego  
dasgupta@cs.ucsd.edu

#### Abstract

We analyze the behavior of nearest neighbor classification in metric spaces and provide finite-sample, distribution-dependent rates of convergence under minimal assumptions. These are more general than existing bounds, and enable us, as a by-product, to establish the universal consistency of nearest neighbor in a broader range of data spaces than was previously known. We illustrate our upper and lower bounds by introducing a new smoothness class customized for nearest neighbor classification. We find, for instance, that under the Tsybakov margin condition the convergence rate of nearest neighbor matches recently established lower bounds for nonparametric classification.

#### 1 Introduction

In this paper, we deal with binary prediction in metric spaces. A classification problem is defined by a metric space  $(X, \rho)$  from which instances are drawn, a space of possible labels  $\mathcal{Y} = \{0, 1\}$ , and a distribution  $\mathbb{P}$  over  $X \times \mathcal{Y}$ . The goal is to find a function  $h : X \rightarrow \mathcal{Y}$  that minimizes the probability of error on pairs  $(X, Y)$  drawn from  $\mathbb{P}$ ; this error rate is the risk  $R(h) = \mathbb{P}(h(X) \neq Y)$ . The best such function is easy to specify: if we let  $\mu$  denote the marginal distribution of  $X$  and  $\eta$  the conditional probability  $\eta(x) = \mathbb{P}(Y = 1 | X = x)$ , then the predictor  $1(\eta(x) \geq 1/2)$  achieves the minimum possible risk,  $R^* = \mathbb{E}_X[\min(\eta(X), 1 - \eta(X))]$ . The trouble is that  $\mathbb{P}$  is unknown and thus a prediction rule must instead be based only on a finite sample of points  $(X_1, Y_1), \dots, (X_n, Y_n)$  drawn independently at random from  $\mathbb{P}$ .

Nearest neighbor (NN) classifiers are among the simplest prediction rules. The  $1$ -NN classifier assigns each point  $x \in X$  the label  $Y_i$  of the closest point in  $X_1, \dots, X_n$  (breaking ties arbitrarily, say). For a positive integer  $k$ , the  $k$ -NN classifier assigns  $x$  the majority label of the  $k$  closest points in  $X_1, \dots, X_n$ . In the latter case, it is common to let  $k$  grow with  $n$ , in which case the sequence  $(k_n : n \geq 1)$  defines a  $k_n$ -NN classifier.

The asymptotic consistency of nearest neighbor classification has been studied in detail, starting with the work of Fix and Hodges [2]. The risk of the NN classifier, henceforth denoted  $R_n$ , is a random variable that depends on the data set  $(X_1, Y_1), \dots, (X_n, Y_n)$ ; the usual order of business is to first determine the limiting behavior of the expected value  $\mathbb{E}R_n$ , and to then study stronger modes of convergence of  $R_n$ . Cover and Hart [3] studied the asymptotics of  $\mathbb{E}R_n$  in general metric spaces, under the assumption that every  $x$  in the support of  $\mu$  is either a continuity point of  $\eta$  or has  $\mu(\{x\}) > 0$ . For the  $1$ -NN classifier, they found that  $\mathbb{E}R_n \rightarrow \mathbb{E}_X[2\eta(X)(1 - \eta(X))] \leq 2R^*(1 - R^*)$ ; for  $k_n$ -NN with  $k_n \uparrow \infty$  and  $k_n/n \downarrow 0$ , they found  $\mathbb{E}R_n \rightarrow R^*$ . For points in Euclidean space, a series of results starting with Stone [13] established consistency without any distributional assumptions. For  $k_n$ -NN in particular,  $R_n \rightarrow R^*$  almost surely [8].

These consistency results place nearest neighbor methods in a favored category of nonparametric estimators. But for a fuller understanding it is important to also have rates of convergence. For