

SF 2935 LECTURE 20

k-NN REGRESSION IN
HIGH DIMENSIONS

2017-12-14

T. Kussai

0. INTRODUCTION

1/19

$$Y = f(x) + \varepsilon$$

$$\hat{f}_{\text{opt}} = \underset{f}{\text{argmin}} E[(Y - f(x))^2]$$

$$\hat{f}_{\text{opt}}(x) = E[Y | x] = \int_{-\infty}^{\infty} y p_{Y|X}(y|x) dy$$

We discuss learning of \hat{f}_{opt}

IF WE DO NOT KNOW THE DISTRIBUTION

$$P_{X,Y}(y|x)$$

WE STUDY NEAREST-NEIGHBOR

REGRESSION & HIGH DIMENSIONAL

SPACES

$$\hat{f}_{\text{opt}}(x) \approx \hat{f}_k(x) = \text{Average}(y_i | x) \cdot m_k(x)$$

$$m_k(x) = k \text{ nearest neighbors}$$

=

HERE x LIES IN THE UNIT HYPERCUBE

2/19

C.F. EXERCISES 4.7

EXERCISE 4 pp. 168 - 169

in JAMES - WITEN - HARTIG
TIOSHITANI

An Introduction to
Statistical Learning

(Project 1)

1. BASIC ASSUMPTIONS AND NOTATIONS

$\mathcal{X} = \text{DATA SPACE} = [0, 1]^d$
= A HYPERCUBE IN \mathbb{R}^d

$x \in \mathcal{X}, y \in \mathcal{X}$ x, y points

$$d(x, y) = \|x - y\|_2 = \sqrt{\sum_{i=1}^d (x_i - y_i)^2}$$

= L₂-distance

$\mathcal{D} = \{x_1, \dots, x_N\}$ A DATA SET $x_i \in \mathcal{X}$

ASSUMPTION

THE DATA POINTS x_1, \dots, x_N ARE I.I.D
(EXAMPLE) OF $\bar{X}_1, \dots, \bar{X}_N$ RESPECTIVELY

$$\bar{X}_i \in \mathcal{U}([0, 1]^d) = \text{UNIFORM DIST'N}$$

$$\boxed{\text{Vol}(\mathcal{X}) = 1}$$

$$\bar{X} = (\bar{X}_1, \dots, \bar{X}_d)$$

$$A \subseteq [0, 1]^d$$

$$\boxed{\bar{X}_i \sim \mathcal{U}(0, 1)}$$

$$P(\bar{X} \in A) = \int \dots \int_A dx_1 \dots dx_d$$

= Vol(A).

HYPER BALL (in \mathbb{R}^d)

4/19

$$S^d(x, r) = \{y \in \mathbb{R}^d \mid \|y - x\|_2 \leq r\}$$

$s \in \mathbb{R}^d$ $nn(s) =$ NEAREST NEIGHBOR
TO s IN \mathcal{D} .

IS DEFINED AS

$$nn(s) = \{x_i \in \mathcal{D} \mid \forall x' \in \mathcal{D} : \|x_i - s\|_2 \leq \|x' - s\|_2\}$$

$$nn_{dist}(s) = \|nn(s) - s\|_2$$

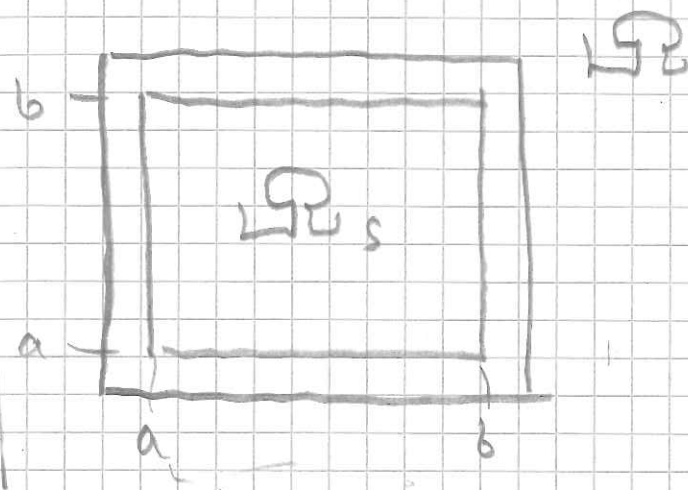
2. FACTS

FACT 1. PROBABILITY OF A
SUBCUBE

5/19

TAKE $s < 1$. AND $b - a = s$

$$\Omega_s = \{x \in \Omega \mid a \leq x_i \leq b \quad i=1, \dots, d\}$$



$$\mathbb{P}((X_1, \dots, X_d) \in \Omega_s)$$

$$= \mathbb{P}(a \leq X_1 \leq b, a \leq X_2 \leq b, \dots, a \leq X_d \leq b)$$

ASSUMPTION

$$\stackrel{\text{ASSUMPTION}}{=} \mathbb{P}(a \leq X_1 \leq b) \cdot \mathbb{P}(a \leq X_2 \leq b) \cdot \dots \cdot \mathbb{P}(a \leq X_d \leq b)$$

$$= \left(\mathbb{P}(a \leq u \leq b) \right)^d \quad \underline{u \in U(0,1)}$$

$$= (b - a)^d = s^d$$

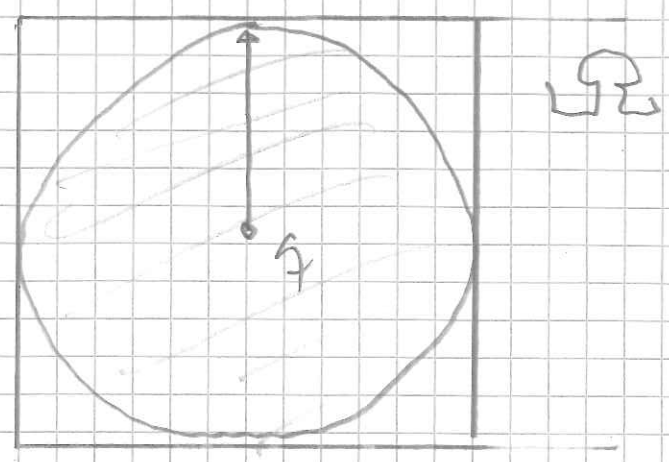
EX. $d=100, \quad s=0.55$

$$s^d = 0.55^{100} \approx 0.0055 = \underline{\underline{0.55\%}}$$

FACT 2

THE DATA SPACE IS VERY
SPARSELY POPULATED.

FACT 2 THE LARGEST HYPERBALL THAT FITS ENTIRELY IN Ω



$$sp^d(q, \frac{1}{2}) = \{ x \in \Omega \mid \|x - q\|_2 \leq \frac{1}{2} \}$$

SOUGHT: THE PROBABILITY THAT ONE DATA POINT LIES IN $sp^d(q, \frac{1}{2})$.

$$P(X \in sp^d(q, \frac{1}{2})) = \int_{sp^d(q, \frac{1}{2})} \dots \int dx_1 \dots dx_d$$

$$= \underbrace{K_d}_{\rho} \left(\frac{1}{2}\right)^{d+2}$$

volume of the unit ball

$$sp^d(q, 1) = \{ x \in \mathbb{R}^d \mid \|x - q\|_2 \leq 1 \}$$

7/19
CHAPTER 2 in FOUNDATIONS OF
DATA SCIENCE SHOWS THAT

$$K_d = \frac{\pi^{d/2}}{\Gamma\left(\frac{d}{2} + 1\right)} = \frac{\pi^{d/2}}{\frac{d}{2} \Gamma\left(\frac{d}{2}\right)}$$

$\Gamma(x+1) = x \Gamma(x)$ Euler's Gamma Function

STIRLING'S FORMULA

$$\Gamma(x) \sim \sqrt{2\pi} e^{-x} \cdot x^{x-1/2}$$

$$\Gamma\left(\frac{d}{2}\right) \sim \sqrt{2\pi} e^{-\frac{d}{2}} \cdot \left(\frac{d}{2}\right)^{\frac{d}{2}-1/2}$$

$$\therefore \Gamma(\bar{X} \in \text{Sp}^d(\mathbb{F}, 1/2))$$

$$\approx \frac{\pi^{d/2}}{\sqrt{2\pi}} \frac{e^{d/2}}{\left(\frac{d}{2}\right)^{d/2-1/2}} \cdot \left(\frac{d}{2}\right)^{\frac{d}{2}}$$

$$\frac{\pi^{d/2}}{\sqrt{2\pi}} = e^{\frac{d}{2} \ln \pi} \quad \frac{d/2-1}{d/2} = e^{\left(\frac{d}{2}-1\right) \ln \frac{d}{2}}$$

$$\approx \frac{1}{\sqrt{2\pi}} e^{\frac{d}{2} (\ln \pi - 1) - \left(\frac{d}{2}-1\right) \ln \frac{d}{2}} \cdot \left(\frac{d}{2}\right)^{d/2}$$

$\rightarrow 0$ as $d \rightarrow \infty$

<u>2</u>	$P(\mathcal{E} \in \mathcal{S}_T^d(a, 1/2))$
2	0.785
4	0.308
10	0.002
20	$2.46 \cdot 10^{-8}$
40	$3.278 \cdot 10^{-21}$
100	$1.868 \cdot 10^{-70}$

FACT 2:

THE VOLUME OF $\mathcal{S}_T^d(a, 1/2)$
 SHRINKS STRONGLY AS d GROWS
 AND IT IS INCREASINGLY
UNPROBABLE THAT ANY POINT
WILL BE FOUND IN THIS BALL
AT ALL

REMARKS:

3/19

$$\text{vol}(\mathbb{S}^d) = 2$$

$$\begin{aligned} \text{diameter}(\mathbb{S}^d) &= \|\underline{1} - \underline{0}\|_2 \\ &= \sqrt{d} \end{aligned}$$

$$\begin{aligned} \underline{1} &= (\underbrace{1, \dots, 1}_{d \text{ ones}}) \\ \underline{0} &= (\underbrace{0, \dots, 0}_{d \text{ zeros}}) \end{aligned}$$

$$\text{vol}(\mathbb{S}^d(\frac{1}{2}, \frac{1}{2})) \rightarrow 0, \text{ as } d \rightarrow +\infty,$$

$$\text{diameter}(\mathbb{S}^d) \rightarrow +\infty, \text{ as } d \rightarrow +\infty$$

||

FACTS

10/19

$$p \equiv P(X \in S_p^d(q, 1/2))$$

X_1, \dots, X_N i.i.d., N GIVEN IN ADVANCE

SUCCESS

$$\text{SUCCESS} = P(X_i \in S_p^d(q, 1/2))$$

Y = NUMBER OF SUCCESSSES
IN X_1, \dots, X_N

$$Y \sim \text{Bin}(N, p)$$

$$E[Y] = Np$$

$$E[Y] \geq 1 \Leftrightarrow$$

$$N \geq \frac{1}{p} = \frac{\Gamma(\frac{d}{2} + 1)}{\pi^{d/2} \cdot (\frac{1}{2})^d}$$

$d = 10$, $N = 401.5$ SO THAT WE
MAY EXPECT IN AVERAGE
AT LEAST ONE POINT IN
 $S_p^d(q, 1/2)$

$$d = 20 \quad N = 40631624$$

etc.

FACT 4 : WE COMPUTE

$$P(Q, r) = \mathbb{P}(\text{nn}(q) \in \mathcal{S}_r^d(S, r))$$

= PROBABILITY THAT $\text{nn}^{dist}(q) \leq r$.

$$P(Q, r) = 1 - \mathbb{P}(\text{nn}^{dist}(q) > r)$$

$$= 1 - \mathbb{P}(\|x_1 - q\|_2 > r, \dots, \|x_N - q\|_2 > r)$$

$$= 1 - \mathbb{P}(\|x_1 - q\|_2 > r) \cdot \dots \cdot \mathbb{P}(\|x_N - q\|_2 > r)$$

$$= 1 - \left(\mathbb{P}(\|x - q\|_2 > r) \right)^N$$

$$= 1 - \left(1 - \mathbb{P}(\|x - q\|_2 \leq r) \right)^N$$

$$= 1 - \left(1 - \mathbb{P}(\{x \in \mathcal{S} \mid \|x - q\|_2 \leq r\}) \right)$$

$$= 1 - \left(1 - \text{Vol}(\mathcal{S}_r^d(q, r) \cap \mathcal{S}) \right)^N$$

$r = 1/2$ $\mathcal{S}_r^d(q, 1/2) \cap \mathcal{S} = \mathcal{S}_r^d(q, 1/2)$

$d \rightarrow \infty$ $\mathbb{P}(Q, r) \approx 0$

$\therefore \text{nn}^{dist}(q)$ GROWS WITH d !

$$P(\text{nn}(z) \in S_{r^d}(z, r)) =$$

12/15

I.E. PROBABILITY THAT $\text{nn}(z) \in S_{r^d}(z, r) \leq r =$

$$= \frac{1 - (1 - \text{Vol}(S_{r^d}(z, r) \cap \Omega))^N}{N}$$

1) $P(\text{nn}(z) \in S_{r^d}(z, r)) \rightarrow 1,$
as $N \rightarrow \infty$ (LARGE DATA SETS)
FOR GIVEN d

2) $P(\text{nn}(z) \in S_{r^d}(z, r)) \rightarrow 0,$
as $d \rightarrow \infty$

3) $\left. \begin{array}{l} \text{Both } N \rightarrow \infty \\ \text{and } d \rightarrow \infty \end{array} \right\} \Rightarrow ?$

3. EXPECTED PREDICTION ERROR & BIAS-VARIANCE TRADE-OFF

13/19

2.1 IN GENERAL

$$Y = f(X) + \varepsilon \quad E[\varepsilon] = 0, E[\varepsilon^2] = \sigma_\varepsilon^2$$

\hat{f} = A MODEL OF f ESTIMATED FROM TRAINING DATA

$$D = (x_i, y_i)_{i=1}^N \quad \text{i.e.,}$$

$$y_i = f(x_i) + \varepsilon_i$$

TAKE A FIXED x AND CONSIDER

$$\text{EPE}(x) = E[(Y - \hat{f}(x))^2]$$

E = EXPECTATION W.R.T $P(X, Y)$ GENERATING D

Mean square error of \hat{f} at x

$$\begin{aligned} E[\hat{f}] &= E[f(x) + \varepsilon] \\ &= f(x) + E[\varepsilon] = f(x) \end{aligned}$$

$$\begin{aligned} \text{Var}[\hat{f}] &= E[(Y - f(x))^2] \\ &= E[\varepsilon^2] = \sigma_\varepsilon^2 \end{aligned}$$

EPE = EXPECTED PREDICTION ERROR

14/19

$$E[(Y - \hat{f}(x))^2]$$

$$= E[Y^2 - 2Y\hat{f}(x) + \hat{f}^2(x)]$$

$$= E[Y^2] - 2E[Y\hat{f}(x)] + E[\hat{f}^2(x)]$$

$$\begin{aligned} \bullet E[Y^2] &= \text{Var}(Y) + (E(Y))^2 \\ &= \sigma_\varepsilon^2 + f^2(x) \end{aligned}$$

$$\bullet E[Y\hat{f}(x)] = E[(f(x) + \varepsilon)\hat{f}(x)]$$

$$= E[f(x)\hat{f}(x)] + E[\varepsilon\hat{f}(x)]$$

$$= E[f(x)\hat{f}(x)] + E[\hat{f}(x)] \cdot \underbrace{E(\varepsilon)}_{=0}$$

$$= E[f(x)\hat{f}(x)] = f(x) E[\hat{f}(x)]$$

$$\bullet E[\text{MSE}(x)] = \sigma_\varepsilon^2 + f^2(x) - 2f(x)E[\hat{f}(x)] + E[\hat{f}^2(x)]$$

$$\bullet E[\hat{f}^2(x)] = \text{Var}(\hat{f}(x)) + (E[\hat{f}(x)])^2$$



$$E \text{ ERROR}^2(x) =$$

$$\sigma_{\epsilon}^2 + \text{Var}(f(x)) + \left(E[\hat{f}(x)] - f(x) \right)^2$$

Variance σ_{ϵ}

\hat{f} AS

ESTIMATED

FROM TRAINING
DATA

BIAS²:

MODEL COMPLEXITY

ERROR DUE TO

APPROXIMATION

3.2 EPC & NEAREST NEIGHBOR REGRESSION

$$Y = f(x) + \epsilon \quad \text{AS ABOVE}$$

$$(x_i, y_i)_{i=1}^N \quad \text{TRAINING DATA}$$

$$\hat{f}(x) = \hat{f}_k(x) = \frac{1}{k} \sum_{i: x_i \in \text{nn}_k(x)} y_i = \frac{1}{k} \sum_{\text{nn}_k(x)} y_i$$

$\text{nn}_k(x)$ = THE SET OF k WITH NEAREST
NEIGHBORS

$$\|x_{(1)} - x\|_2 \leq \|x_{(2)} - x\|_2 \leq \dots \leq \|x_{(k)} - x\|_2$$

\uparrow
nearest neighbor

$$\hat{f}_k(x) = y_i \quad \text{DECOMPOSITION:}$$

IF $x_i \in \text{nn}_k(x) \quad x_i \in \mathcal{D}$

1/6/19

AS $N, k \rightarrow \infty$ SO THAT $\frac{k}{N} \rightarrow 0$

IT CAN BE SHOWN THAT

$$\hat{f}_k(x) \rightarrow E[Y|X]$$

RATE OF CONVERGENCE
 $O(1/N^{1/d})$

(Recall that

$$E[(\bar{Y} - g(\bar{X}))^2]$$

is minimized by $g(\bar{X}) = E[Y|\bar{X}]$)

UNIVERSAL ESTIMATOR!

BUT: ASSUME $X_i \sim U([0,1]^d)$

$$Y_i = f(X_i) + \epsilon_i \quad f \text{ UNKNOWN}$$

• $P(\epsilon \in \text{span}(0.8)) = 0.8^{10} = 0.1$

80% of range in every dimension to

cover 10% of data

- MOST DATA POINTS X_1, \dots, X_N ARE CLOSER TO BOUNDARIES THAN ANY OTHER DATA POINT

WE NEED

17/19

$$\begin{aligned} E[\hat{f}_k(x)] &= \frac{1}{k} \sum_{n \in \mathcal{N}(x)} E[\mathcal{Y}_i] \\ &= \frac{1}{k} \sum_{n \in \mathcal{N}(x)} E[f(x_i) + \varepsilon_i] = \\ &= \frac{1}{k} \sum_{n \in \mathcal{N}(x)} f(x_i) + \underbrace{E[\varepsilon_i]}_{=0} = \frac{1}{k} \sum_{n \in \mathcal{N}(x)} f(x_i) \\ &= \frac{1}{k} \sum_{e=1}^k f(x_{(e)}) \quad \left\{ \begin{array}{l} x_{(1)}, \dots, x_{(k)} \text{ ORDERED BY} \\ \|x_{(1)} - x\|_2 \leq \|x_{(2)} - x\|_2 \leq \dots \leq \\ \leq \|x_{(k)} - x\|_2 \leq \dots \end{array} \right. \\ &\quad \text{IND.} \end{aligned}$$

$$\begin{aligned} \text{Var}(\hat{f}_k(x)) &\stackrel{b}{=} \frac{1}{k^2} \sum_{n \in \mathcal{N}(x)} \text{Var}(\mathcal{Y}_i) \\ &= \frac{1}{k^2} \sum_{n \in \mathcal{N}(x)} \sigma_{\varepsilon}^2 = \frac{1}{k^2} k \cdot \sigma_{\varepsilon}^2 = \frac{1}{k} \sigma_{\varepsilon}^2 \end{aligned}$$

$$\therefore \text{EPE}(x) = \sigma_{\varepsilon}^2 + \left(f(x) - \frac{1}{k} \sum_{e=1}^k f(x_{(e)}) \right)^2 + \frac{\sigma_{\varepsilon}^2}{k}$$

EXPECTED PREDICTION ERROR FOR
k NN - REGRESSION

18/19

IN SMALL DIMENSION, FOR

SMALL h , THE FEW CLOSEST

NEIGHBORS WILL HAVE THEIR

FUNCTION VALUES $f(x_i)$ CLOSE

TO $f(x)$, SO THEIR AVERAGE

SHOULD BE CLOSE TO $f(x)$

IN ADDITION

$$\sup_x |E[Y|X=x] - f(x)| \rightarrow 0$$

WITH PROBABILITY ONE AND WITH

$$\text{RATE} = O\left(\left(\frac{1}{N}\right)^{1/d}\right) \leftarrow \text{SEVERAL DIFFERENT PROOFS}$$

A UNIVERSAL APPROXIMATION!

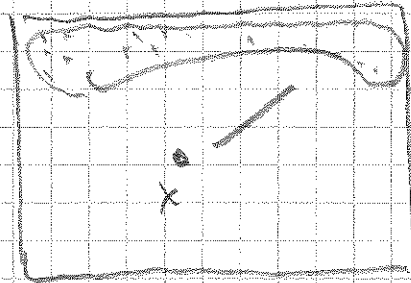
19/19

BUT: IN LARGE DIMENSION

$$D = \{x_1, \dots, x_N\}, N \text{ SAMPLES } \mathcal{U}([0, 1]^d)$$

FOR A HIGH VALUE OF d THE
SAMPLES IN D ARE CLOSE TO

THE BOUNDARY OF $[0, 1]^d$ (As shown
in 2.)



FOR ANY $x \in \mathcal{S}_d$, $\frac{1}{k} \sum_{i=1}^k f(x_{(i)})$

IS AN EXTRAPOLATION FROM NEIGHBORING
SAMPLES RATHER THAN INTERPOLATION
BETWEEN THEM.

$$\left(f(x) - \frac{1}{k} \sum_{i=1}^k f(x_{(i)}) \right)^2$$

(CAN BE VERY LARGE

CONVERGENCE RATE $\mathcal{O}\left(\frac{1}{N^{1/d}}\right)$

IS STILL VALID, BUT IS VERY
SLOW.