

SF 2535 LECTURE 18

ON THE NEAREST NEIGHBOR
IN VERY LARGE METRIC
SPACES

T. Kuski

2017

1 NOTATIONS

1/23

$$\underline{x} = (x_1, \dots, x_d)^T \quad \underline{x} \in \mathbb{R}^d$$

$$D(\underline{x}) = \sum_{i=1}^d x_i^2 = \underline{x}^T \underline{x} = \text{SQUARED DISTANCE TO THE ORIGIN } 0 \in \mathbb{R}^d$$

2 SKEWNORMAL DISTRIBUTION

$$\underline{x} \sim \mathcal{S} \phi_{\mathcal{Q}}(\underline{x}) \pi(\underline{x})$$

$$\phi_{\mathcal{Q}}(\underline{x}) = \frac{1}{(2\pi)^{d/2}} e^{-\underline{x}^T \underline{x} / 2} \quad (= \mathcal{N}(0, \mathbb{I}_{\mathcal{Q}}))$$

$\pi(\underline{x}) =$ A SKEWING FUNCTION

i) $0 \leq \pi(\underline{x}) \leq 1$

ii) $\pi(-\underline{x}) = 1 - \pi(\underline{x})$

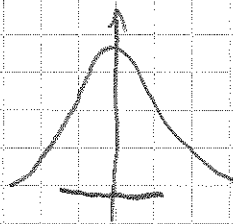
EXAMPLE $d=1$

$$\pi(x) = \underline{\Phi}(\lambda x)$$

$$\rightarrow \lambda < x < +\infty$$

$$\Phi(x) = \int_{-\infty}^x \phi_1(u) du$$

$$2 \cdot \underline{\Phi}(\lambda x) \phi_{\mathcal{Q}}(x)$$



$$\lambda = 0$$

$$2 \cdot \underline{\Phi}(0, x) \phi_{\mathcal{Q}}(x) = 2 \cdot \frac{1}{2} \phi_{\mathcal{Q}}(x) = \phi_{\mathcal{Q}}(x)$$

□

THM 1

$$X \sim N(\Sigma) \varphi_d(X) \Leftrightarrow D = X^T X \sim \chi_d^2$$

Proof: [1] ■

χ_d^2 = CHISQUARED WITH d DEGREES
OF FREEDOM

$$f_D(d) = \frac{1}{\Gamma\left(\frac{d}{2}\right)} d^{\frac{d}{2}-1} \left(\frac{1}{2}\right)^{d/2} e^{-d/2}, \quad 0 \leq d$$

FACTS

$$i) E[D] = d \quad ii) \text{Var}(D) = 2d$$

see [2] ■

3) CONVERGENCE IN PROBABILITY [3]

a) X_1, \dots, X_n, \dots A SEQUENCE OF RANDOM VECTORS CONVERGES IN PROBABILITY TO $a \in \mathbb{R}^d$

IFF

$$\lim_{n \rightarrow \infty} P(\|X_n - a\| > \epsilon) = 0$$

FOR ANY $\epsilon > 0$.

WE WRITE

$$X_n \xrightarrow{P} a, \text{ as } n \rightarrow \infty$$

b) CRAMÉR - SLUTSKY'S THEOREM

a) IF $X_n \xrightarrow{P} a$ as $n \rightarrow \infty$ AND

$g(x)$ IS A CONTINUOUS FUNCTION

$$g(x) : \mathbb{R}^d \rightarrow \mathbb{R}$$

AT a , THEN

$$g(X_n) \rightarrow g(a)$$

as $n \rightarrow \infty$.

ii)

$$\begin{matrix} \bar{X}_n \\ \bar{Y}_n \end{matrix} \xrightarrow{P} \begin{matrix} a \\ b \end{matrix} \quad (a, b \in \mathbb{R}) \quad \bar{X}_n \in \mathbb{R},$$

$$\bar{Y}_n \in \mathbb{R} \quad b \neq 0 \quad \bar{Y}_n \in \mathbb{R}$$

THEN $\begin{matrix} \bar{X}_n \\ \bar{Y}_n \end{matrix} \xrightarrow{P} \frac{a}{b}, \quad a, n \rightarrow \infty.$

3) CHEBYSHEV'S INEQUALITY

$$E[\bar{X}] = \mu \quad \text{Var}(\bar{X}) = \sigma^2$$

(\bar{X} REAL VALUED ST. V.)

$$\mathbb{P}(|\bar{X} - \mu| > \varepsilon) \leq \frac{1}{\varepsilon^2} \sigma^2$$

4) THE THEOREM [4]

$$X_1, \dots, X_n \quad X_i \text{ i.i.d.} \sim \pi(x) \varphi_d(x)$$

$$D_i = D(X_i) \quad i = 1, \dots, n$$

$$D_{\min} = \min_{1 \leq i \leq n} D_i \quad D_{\max} = \max_{1 \leq i \leq n} D_i$$

THM 2. For any $\epsilon > 0$

$$\lim_{d \rightarrow \infty} P(D_{\max} \leq (1 + \epsilon) D_{\min}) = 1$$

PROOF: SET $\mu = E[D_i]$ (AS X_i
ARE I.I.D.)

AND
$$\tilde{y}_i = \frac{D_i}{\mu}$$

THEN
$$E[\tilde{y}_i] = \frac{1}{\mu} E[D_i] = 1.$$

$$\text{Var}[\tilde{y}_i] = \frac{1}{\mu^2} \text{Var}[D_i] = \frac{2\sigma^2}{\mu^2} = \frac{2}{\mu^2}$$

THEN CHEBYSJEV'S INEQUALITY
TELLS THAT FOR EVERY i

$$\mathbb{P}(|y_i - 1| > \varepsilon) \leq \frac{1}{\varepsilon^2} \text{Var}(y_i) = \frac{2}{d\varepsilon}.$$

$$\therefore \mathbb{P}(|y_i - 1| > \varepsilon) \leq \frac{2}{d\varepsilon} \rightarrow 0, \text{ as } d \rightarrow \infty$$

$$\text{i.e. } p_i \xrightarrow{P} 1, \text{ as } d \rightarrow \infty$$

d = DIMENSION OF THE EUCLIDEAN SPACE

HENCE, BY CONVERGENCE SLUTSKY (i)

$$\max(y_1, \dots, y_n) \xrightarrow{P} 1, \text{ as } d \rightarrow \infty$$

$$\min(y_1, \dots, y_n) \xrightarrow{P} 1, \text{ as } d \rightarrow \infty$$

SINCE \max AND \min ARE CONTINUOUS
FUNCTIONS,

THEN

$$\frac{D_{\max}}{D_{\min}} = \frac{\max(D_1, \dots, D_n)}{\min(D_1, \dots, D_n)}$$

$$= \frac{\frac{1}{\mu} \max(D_1, \dots, D_n)}{\frac{1}{\mu} \min(D_1, \dots, D_n)}$$

$$= \max\left(\frac{D_1}{\mu}, \dots, \frac{D_n}{\mu}\right)$$

$$= \frac{\max\left(\frac{D_1}{\mu}, \dots, \frac{D_n}{\mu}\right)}{\min\left(\frac{D_1}{\mu}, \dots, \frac{D_n}{\mu}\right)}$$

$$= \frac{\max(y_1, \dots, y_n)}{\min(y_1, \dots, y_n)}$$

$$\stackrel{P}{\rightarrow} 1, \text{ as } d \rightarrow d$$

BY CRAMÉR-SLUTSKY

THUS:

$$P(D_{\max} \leq (1+\epsilon) D_{\min})$$

$$= P\left(\frac{D_{\max}}{D_{\min}} \leq 1+\epsilon\right) = P\left(\frac{D_{\max}}{D_{\min}} - 1 \leq \epsilon\right)$$

$$= P\left(\left|\frac{D_{\max}}{D_{\min}} - 1\right| \leq \epsilon\right), \text{ as } D_{\max} > D_{\min}$$

BUT $\frac{D_{\max}}{D_{\min}} \xrightarrow{P} 1$, as $d \rightarrow \infty$

AND THENCE

$$P\left(\left|\frac{D_{\max}}{D_{\min}} - 1\right| \leq \epsilon\right) = 1 - P\left(\left|\frac{D_{\max}}{D_{\min}} - 1\right| > \epsilon\right)$$

WE GET

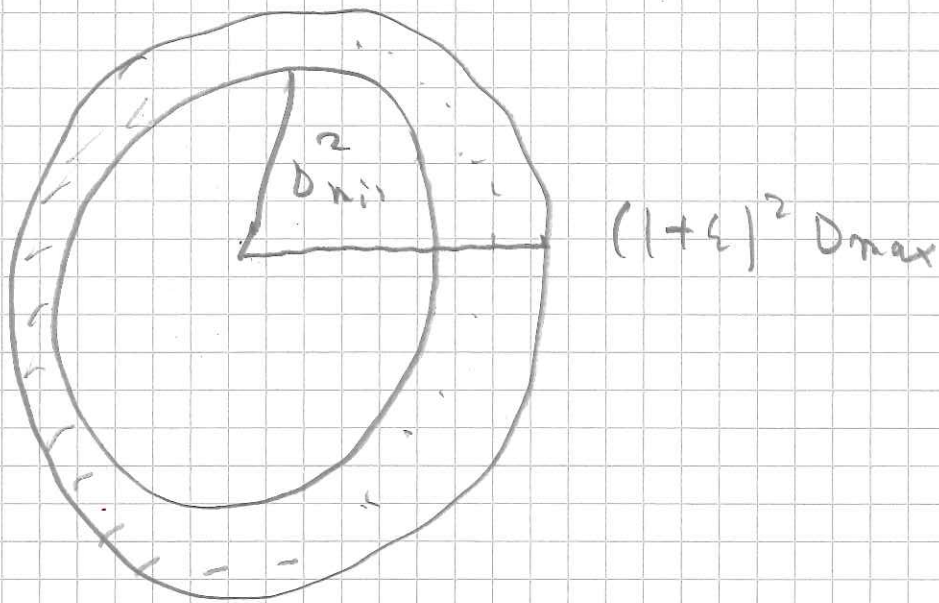
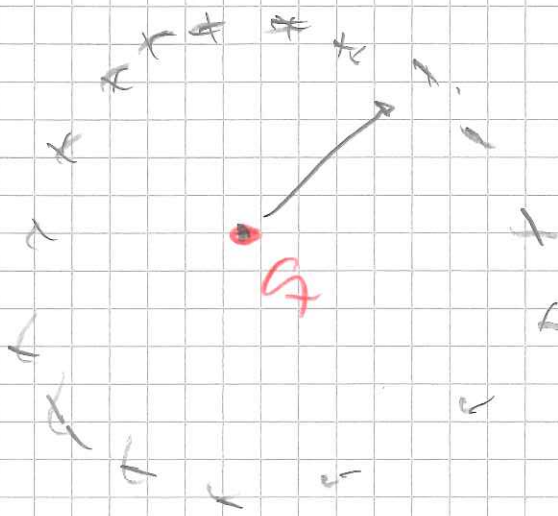
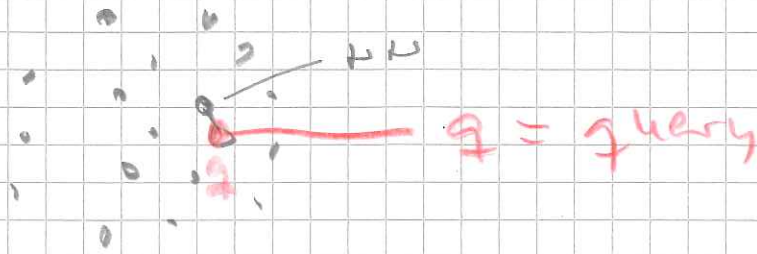
$$\lim_{d \rightarrow \infty} P\left(\left|\frac{D_{\max}}{D_{\min}} - 1\right| \leq \epsilon\right) = 1 - \lim_{d \rightarrow \infty} P\left(\left|\frac{D_{\max}}{D_{\min}} - 1\right| > \epsilon\right)$$

$$= 1 - 0 = 1.$$



NN

Hence: Nearest Neighbor is UNSTABLE
IN HIGH DIMENSION

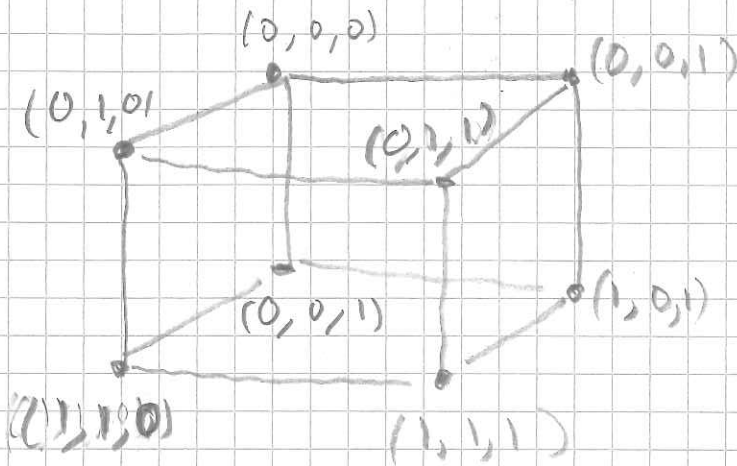


5) BINARY HYPERCUBE

10/23

\mathcal{U} = THE BINARY HYPERCUBE

$$= \{ x \mid x = (x_1, \dots, x_d), x_i \in \{0, 1\} \}$$



$$d = 3$$

$$2^d = 8$$

$$|\mathcal{U}| = 2^d = \text{CARDINALITY OF } \mathcal{U}$$

I.E. THERE ARE 2^d VERTICES x in \mathcal{U}

TAKE $x_i \sim \text{Be}(1/2)$ i.e. $\frac{1}{2} P(x_i=1) = P(x_i=0)$

FOR ALL $i=1, \dots, d$ INDEPENDENTLY

$$\therefore x \in \mathcal{U} : \mathbb{P}(x) = \frac{1}{2^d}$$

HAMMING DISTANCE

x	y	$x \oplus y$
1	1	0
1	0	1
0	1	1
0	0	0

$$x \in U, y \in U$$

$$d(x, y) = \sum_{i=1}^d x_i \oplus y_i \quad \left(= \sum_{i=1}^d \overbrace{|x_i - y_i|}^{\text{real number}} \right)$$

= NUMBER OF POSITIONS

WHERE x_i AND y_i DISAGREE
OR WHERE $x \neq y$.

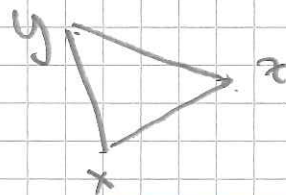
$d(x, y)$ IS A METRIC ON $U \times U$

1) $d(x, y) \geq 0$, $d(x, y) = 0 \Leftrightarrow x = y$

2) $d(x, y) = d(y, x)$

3) $d(x, y) \leq d(x, z) + d(z, y)$

$d(x, y)$ IS CALLED THE HAMMING
DISTANCE.



TRIANGLE INEQUALITY IS READILY
PROVED,

TAKE A FIXED $y \in U$

AND DRAW (INDEPENDENTLY OF y)

$$x \sim P(x) = \frac{1}{2^d}$$

THEN IT FOLLOWS

LEMMA

$$Q(x, y) \sim \text{Bin}(d, \frac{1}{2})$$

PROOF: 1) d TRIALS WITH TWO OUTCOMES
 $x_i = 0$ OR $x_i = 1$

SUCCESS = $\{x_i \neq y_i\}$ FOR TRIAL i

$$2) P(x_i \neq y_i) = \frac{1}{2} = p$$

SINCE y IS GIVEN, IF $y_i = 1$, $P(x_i = 0) = \frac{1}{2}$
 $y_i = 0$ $P(x_i = 1) = \frac{1}{2}$

3) p IS THE SAME FOR ALL TRIALS

4) THE NUMBER OF TRIALS, d ,
 IS FIXED IN ADVANCE.

5) TRIALS ARE INDEPENDENT

1) THE SYMPLY THAT, d , 13/23
THE CLAIM THAT
HENCE LEMMA FOLLOWS,

$d(x, y) =$ THE NUMBER OF
SUCCESSSES IN d
TRIALS $\sim \text{Bin}(d, \frac{1}{2})$.

BUT

$$\begin{aligned} X \sim \text{Bin}(n, p) &\Rightarrow E[X] = np \\ \text{Var}[X] &= np(1-p) \end{aligned}$$

HENCE

$$E_x[d(x, y)] = \frac{d}{2} \quad \text{Var}_x[d(x, y)] = \frac{d}{4}$$

Fix $y \in U$ AND TAKE AN I.I.D
SAMPLE x_1, \dots, x_n FROM $P(x)$

$$D_i = d(x_i, y) \quad i=1, \dots, n$$

$$D_{\max} = \max_{1 \leq i \leq n} (D_1, \dots, D_n)$$

$$D_{\min} = \min_{1 \leq i \leq n} (D_1, \dots, D_n)$$

BUT THM 2 ABOVE DEPENDS
CRITICALLY ON THE FACT THM

$$\frac{\text{Var}(D_i)}{(E(D_i))^2} \rightarrow 0, \text{ as } d \rightarrow \infty$$

SHOULD HOLD.

HERE

$$\frac{\text{Var}(D_i)}{(E(D_i))^2} = \frac{d/4}{d^2/4} = \frac{1}{d} \rightarrow 0$$

AS $d \rightarrow \infty$, BUT THEN WE SEE
BY THE PROOF OF THM 2 THAT

THM 3,
IF D_1, \dots, D_n ARE INDEPENDENT
HAMMING DISTANCE TO AN $y \in U$,
THEN
$$\lim_{d \rightarrow \infty} \mathbb{P}(D_{\max} \leq (1+\epsilon)D_{\min}) = 1$$

WE SEE CLEARLY THAT THIS
HOLDS FOR ANY $y \in U$.

SOME CONCEPTS OF \mathcal{U}_n

15/12/23

A VERTEX $x \in \mathcal{U}$ IS ORTHOGONAL
TO A VERTEX $y \in \mathcal{U}$ IF AND ONLY
IF

$$d(x, y) = \frac{d}{2}$$

~~x s.t. $d(x, y) = \frac{d}{2}$ FOR A FIXED y
ARE SAID TO LIE IN THE EQUATOR~~

• $0 \leq d(x, y) \leq d$

• $d(x, y) = d$ IFF $x = (1, 1, \dots, 1)$
 $y = (0, 0, \dots, 0)$

BALL WITH RADIUS r AND
CENTER IN y

$$B_r(y) = \{x \in \mathcal{U} \mid d(x, y) \leq r\}$$

WE CAN ROTATE U SO THAT
ANY OF ITS POINT IS ORIGIN
(A NORTH POLE) $y = 0$

THE POINTS AT THE DISTANCE
 $d/2$
FROM THE POLE 0 ARE CALLED
THE EQUATOR.

WE RECALL THAT
 $Bin(n, p) \approx N(np, \underbrace{np(1-p)}_{NR \text{ VARIANCE}})$
IF n IS LARGE.

HENCE WE GET THAT
 $d(x, y)$ IS APPROXIMATELY $N(\frac{d}{2}, \frac{d}{4})$
WHEN d IS LARGE.

SET $D = d(x, y)$

$$D \approx N\left(\frac{d}{2}, \frac{d}{4}\right)$$

WE COMPUTE THE PROBABILITY THAT D LIES IN

$$B_{\frac{d}{2}+k}(y) \setminus B_{\frac{d}{2}-k}(y)$$

$$= \left\{ x \mid \frac{d}{2}-k < d(x, y) \leq \frac{d}{2}+k \right\}$$

OR DIFFER FROM MEAN DISTANCE TO y BY AT MOST k BITS

$$P\left(\frac{d}{2}-k < D < \frac{d}{2}+k\right) =$$

$$= P\left(-k < D - \frac{d}{2} < +k\right)$$

$$= P\left(\frac{-k}{\sqrt{\frac{d}{4}}} < \underbrace{\frac{D - \frac{d}{2}}{\sqrt{\frac{d}{4}}}}_{\sim N(0,1)} < \frac{+k}{\sqrt{\frac{d}{4}}}\right)$$

$D \sim N\left(\frac{d}{2}, \frac{d}{4}\right)$
 $\Phi(-x) = 1 - \Phi(x)$

$$= \Phi\left(\frac{k}{\sqrt{\frac{d}{4}}}\right) - \Phi\left(-\frac{k}{\sqrt{\frac{d}{4}}}\right) =$$

$$= 2\Phi\left(\frac{k}{\sqrt{\frac{d}{4}}}\right) - 1 \quad \square$$

$$p_{nd} \hat{=} 2 \Phi\left(\frac{k}{\sqrt{d}}\right) - 1$$

$$d = 100$$

$$k = 3 \quad p_{nd} = 0.4515$$

$$k = 12 \quad p_{nd} = 0.9836$$

$$k = 15 \quad p_{nd} = 0.9973$$

$$d = 1000$$

$$k = 11 \quad p_{nd} = 0.51$$

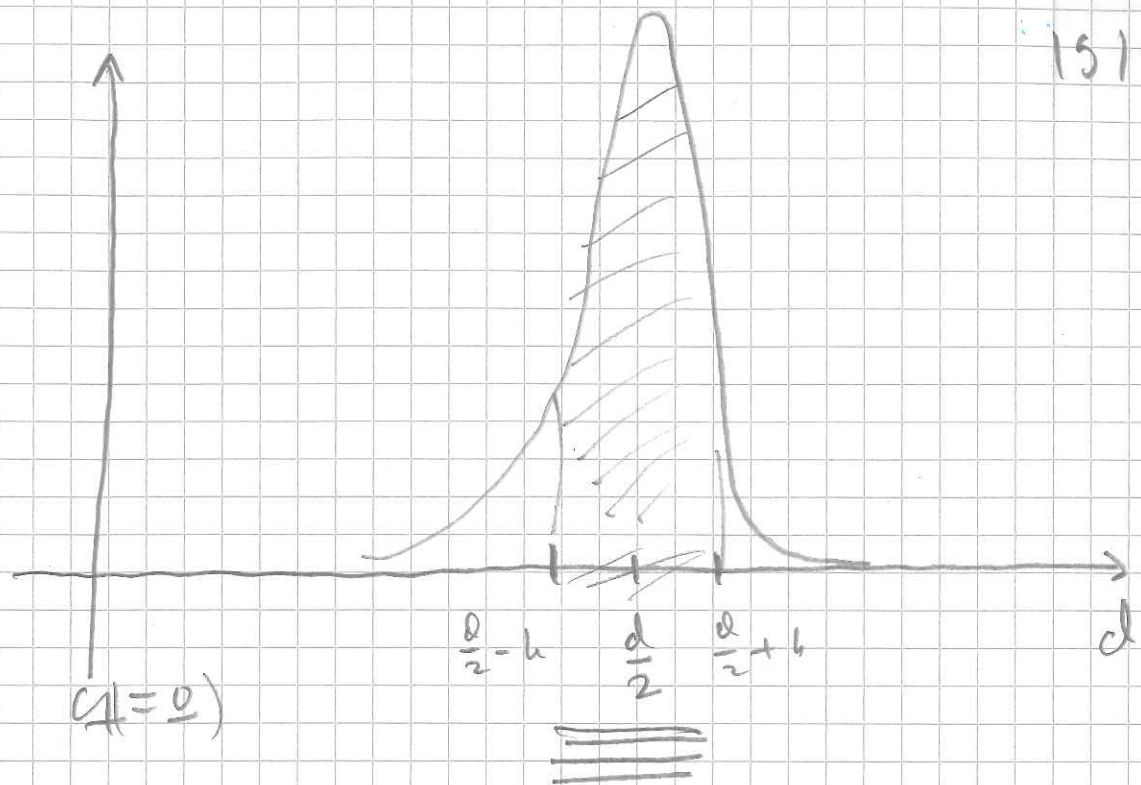
$$k = 12 \quad p_{nd} = 0.55$$

$$k = 55 \quad p_{nd} = 0.9998$$

$$d = 10000$$

$$k = 186 \quad p_{nd} = 0.9999$$

15/23



MOST OF THE PROBABILITY MASS IS
 LOCATED AROUND THE MEAN
 FOR LARGE d .

$f =$ INTRINSIC DIMENSIONALITY (S. BOZINS
et. al.)

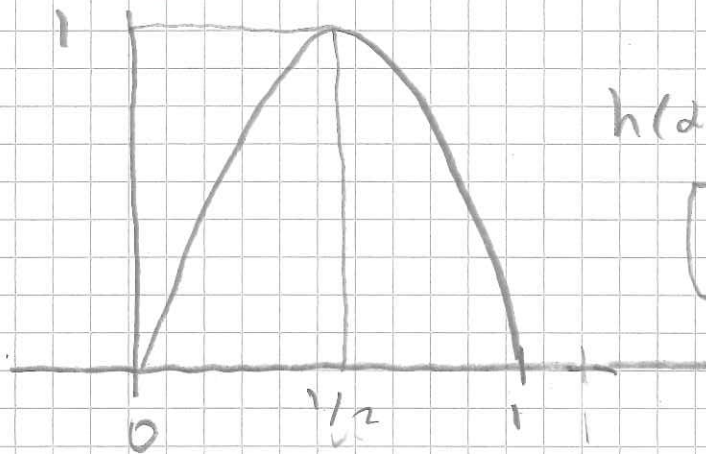
$$f = \frac{(E[d(x,y)])^2}{2 \text{Var}(d(x,y))}$$

ONE MORE OBSERVATION

IF $\alpha < 1/2$ AND

$$\sum_{h=0}^{\lfloor \alpha \rfloor} \binom{\alpha}{h} < 2^{\alpha h(\alpha)}$$

$$\begin{aligned} K &= 2^\alpha \\ \alpha &= K/2 \end{aligned}$$

WHERE $h(\alpha) \triangleq -\alpha \log_2(\alpha) - (1-\alpha) \log_2(1-\alpha)$ For $0 \leq \alpha \leq 1$, IS THE BINARY ENTROPY FUNCTION $h(\alpha), 0 \leq \alpha \leq 1$

$$0 \log_2 0 = 0$$

$$\begin{aligned} \therefore P(B_K(y)) &= P_x(d(x, u) \leq K) \\ &= \sum_{h=0}^K \binom{\alpha}{h} \left(\frac{1}{2}\right)^h \left(\frac{1}{2}\right)^{\alpha-h} = \frac{1}{2^\alpha} \sum_{h=0}^K \binom{\alpha}{h} \\ &\leq \frac{1}{2^\alpha} 2^{\alpha h(K/\alpha)} \quad (\alpha = K/\alpha) \end{aligned}$$

$$\begin{aligned} P(B_K(y)) &= \frac{|B_K(y)|}{2^\alpha} \quad \begin{array}{l} \text{NUMBER} \\ \text{OF } x \text{ IN} \\ B_K(y) \end{array} \\ &= \frac{\text{VOLUME OF } B_K(y)}{\text{TOTAL VOLUME}} \end{aligned}$$

OR

21/23

$$D(B_k(1)) \leq \frac{1}{2^{d(1-h(k/d))}} = \left(\frac{1}{2}\right)^{d(1-h(k/d))}$$



BECOMES SMALL

WITH d INCREASE

FOR GIVEN k , SINCE

$$d(1-h(k/d)) \rightarrow \infty, \text{ as } d \rightarrow \infty$$

APPENDIX

$$h\left(\frac{k}{d}\right) = -\frac{k}{d} \log_2 \frac{k}{d} - \left(1 - \frac{k}{d}\right) \log_2 \left(1 - \frac{k}{d}\right)$$

$$d(1-h(k/d))$$

$$= d + d \left(\frac{k}{d} \log_2 \left(\frac{k}{d}\right) + \left(1 - \frac{k}{d}\right) \log_2 \left(1 - \frac{k}{d}\right) \right)$$

$$= d + k \log_2 \left(\frac{k}{d}\right) - (d - k) \log_2 \left(1 - \frac{k}{d}\right)$$

$$= d - k \log_2(d) - \underbrace{(d - k) \log_2 \left(1 - \frac{k}{d}\right)}_{= C(d)} + k \log_2 \left(1 - \frac{k}{d}\right)$$

$$\log_2 \left(1 - \frac{k}{d}\right)^d$$

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

$$\rightarrow \log_2 e^{-k} \text{ as } d \rightarrow \infty$$

I.E.

$$d(1-h(k/d)) = d - k \log_2(d) + \underbrace{C(d)}_{\text{BOUNDED}}$$

WHERE d INCREASES,

$$\underline{d - k \log_2(d) \rightarrow +\infty}$$

in d



HENCE WE SEE THAT FOR
ARBITRARY RADIUS $k, < d/2$

$$\mathbb{P}(\mathcal{B}_k^u)$$

WILL BE VERY SMALL IF

d IS SUFFICIENTLY LARGE,

\equiv

ANOTHER ANALYSIS OF THESE [67]
PHENOMENA IN \mathbb{R}^d IS

FOUND IN [6]

REFERENCES

23/23

- [1] WANG, BOYER, LEWTON:
STATISTICS AND PROBABILITY
LETTERS 66, 2004, 355-358
- [2] WEIKENGREN: MADE: BETA
- [3] A. GUT: AN INTERMEDIATE COURSE IN
PROBABILITY, 2ND ED.
2005
- [4] BETER, GOLDSTEIN,
PAMAKRISHWAN, EHOFT.
- When is "NEAREST NEIGHBOR"
MEANINGFUL?
- Lecture Notes in Computer Science
1540; pp. 217-231, 1998
- [5] S. BRIN: NEAR NEIGHBOR SEARCH
IN LARGE METRIC
SPACE
INT'L CONFERENCE ON VERY
LARGE DATABASES: BASEL VLDB
1995
- [6] BLUM, HOPLCRAFT, KANNAN:
CHAPTER 2 IN
FOUNDATIONS OF DATA SCIENCE