# L7: Approximate Nearest Neighbors

Jeff M. Phillips

January 29, 2020

Euclidean Distance $\quad d(a,b) = \|a-b\|_2$

$\quad a, b \in \mathbb{R}^d \qquad a = (a_1, a_2, \dots a_d)$

Raw Data $\longrightarrow \mathbb{R}^d \qquad$ usually $d$ large

Example $\Rightarrow$ bag-of-word

$\qquad \Rightarrow$ word-vector embedding

$\qquad \Rightarrow$ images

$\qquad\qquad$ CNN $\longrightarrow$ intermediate layer

$\qquad\qquad$ SIFT vector $\in \mathbb{R}^{128}$

# Images and SIFT Features



| N1 | N2 | N3 |
|----|----|----|
| N8 | X  | N4 |
| N7 | N6 | N5 |

# Images and SIFT Features



| N1 | N2 | N3 |
|----|----|----|
| N8 | X  | N4 |
| N7 | N6 | N5 |

# LSH for Euclidean dist.

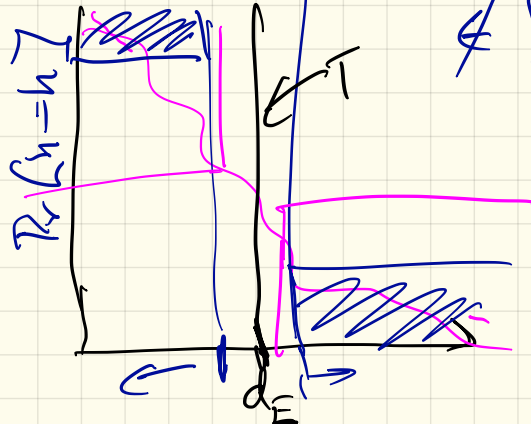$$\Pr_{h \in \mathcal{H}}\left[h(a) = h(b)\right] = S(a,b)$$

<span style="color:red">Jaccard<br>triangle<br>angular</span>

Euclidean

$$S(a,b) = \langle a,b \rangle = \sum_{i=1}^{d} a_i b_i$$
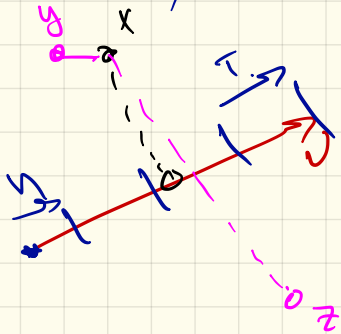
$$\notin [0,1]$$

$$h_{u,\eta}(x)$$

$$\eta \in \text{Unif}(0, T)$$

↳ dist threshold

$$h_{u,\eta}(x) = \left\lfloor \left( \frac{\langle x, u \rangle + \eta}{T} \right) \right\rfloor \quad (\text{mod } m)$$

Gaussian RV.

$$g \in \frac{1}{(2\pi)^{d/2}} e^{-\frac{|M|^2}{2}}$$

$$U \in \frac{g}{\|g\|}$$

$$g \in G_d$$
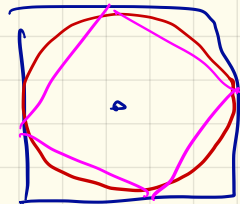
$$g = (g_1, g_2, \cdots g_d)$$

each

$$g_i \in G_1$$

Box-Muller

$$g = \sqrt{-2\ln(\textcircled{$u_1$})} \cos(2\pi \textcircled{$u_2$})$$

$U \in \mathbb{R}^d$

$$\bar{U} = \begin{cases} \bar{U}_1 \in Unif\,(o,1) \\ \bar{U}_2 \in Unif\,(o,1) \\ \phantom{x} \\ \bar{U}_d \in Unif\,(o,1) \end{cases}$$

$$U \in \frac{\bar{U}}{\|\bar{U}\|}$$



box rad1

$2 \cdots 2$ box

$1 = 2^d$

Vol ball radios 1 in $\mathbb{R}^d$

$$\frac{\pi^{d/2}}{\Gamma(d/2+1)} \underset{\widetilde{}}{=} \frac{\pi^{d/2}}{(d/2)!} \leq 1$$

# High-d NN Euclidean

- LSH
- Data Structures Trees
- Graphs

High-D $n^{\lceil d/2 \rceil}$

Voronoi Diagram

# $k$D-Tree

maybe $d = 12$



Subdivision

Tree structure

# Approximate Queries on $k$D-Tree

# Word Vectors & Association

word = nurse $\longrightarrow$ $v_{nurse} \in \mathbb{R}^{300}$

doctor $\quad v_{doctor} \in \mathbb{R}^{300}$

if text context is similar
the two words should have
small cosine distance

woman

man

$$\left( V_{woman} - V_{man} \right) + V_{king}$$

queen

$$0 \in \mathbb{R}^{300}$$

king

nurse

analogy

man : woman :: king : ?

gender

queen

doctor