

# L7: Approximate Nearest Neighbors

Jeff M. Phillips

January 30, 2019

# Edit Distance

defined between <sup>short</sup> 1 string

# operations (delete, add, replace)  
to turn one string  
into another

edit (mines, smiles)

- 1 mines
- 2 smiles
- 3 smiles

works only  
if dist  
small

No LSH

data  $a, b \in \mathbb{R}^d$

$\rightarrow L_p$

$\rightarrow \tilde{a} \leftarrow \frac{a}{\|a\|_2}$

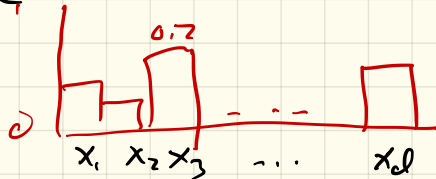
Cosine, Angular

$\rightarrow \tilde{a} \leftarrow \frac{a}{\|a\|_1}$

ensure each coord  $a_i \geq 0$

$$\Delta^{d-1} = \{x \in \mathbb{R}^d \mid \|x\|_1 = 1, \text{ each } x_i \geq 0\}$$

$$\Delta_0^{d-1} = \{x \in \mathbb{R}^d \mid \|x\|_1 = 1, \text{ each } x_i > 0\}$$



$$\sum_{i=1}^d x_i = 1$$

probabilities dist.

$$\mathbb{P}[x=3] = 0.2$$

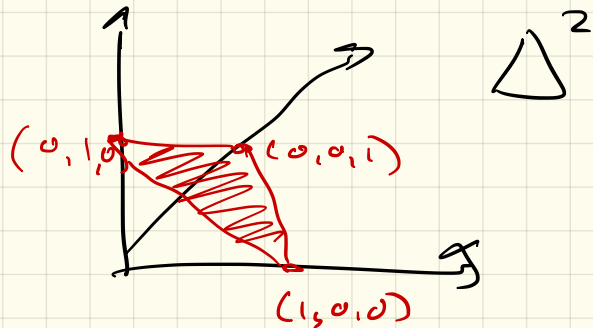
# KL Divergence

$$g \text{ of } P \in \Delta_0^{d-1}$$

$$\text{no } p_i = 0$$

$$\begin{aligned} d_{KL}(P, g) &= d_{KL}(P \parallel g) \\ &= \sum_{i=1}^d p_i \ln \left( \frac{p_i}{g_i} \right) \end{aligned}$$

Jensen-Shannon



# Word Vector Embeddings

large corpus text (entire web)

"The story said, ..."

$v_{\text{story}} \in \mathbb{R}^{300}$

Use Euclidean dist.

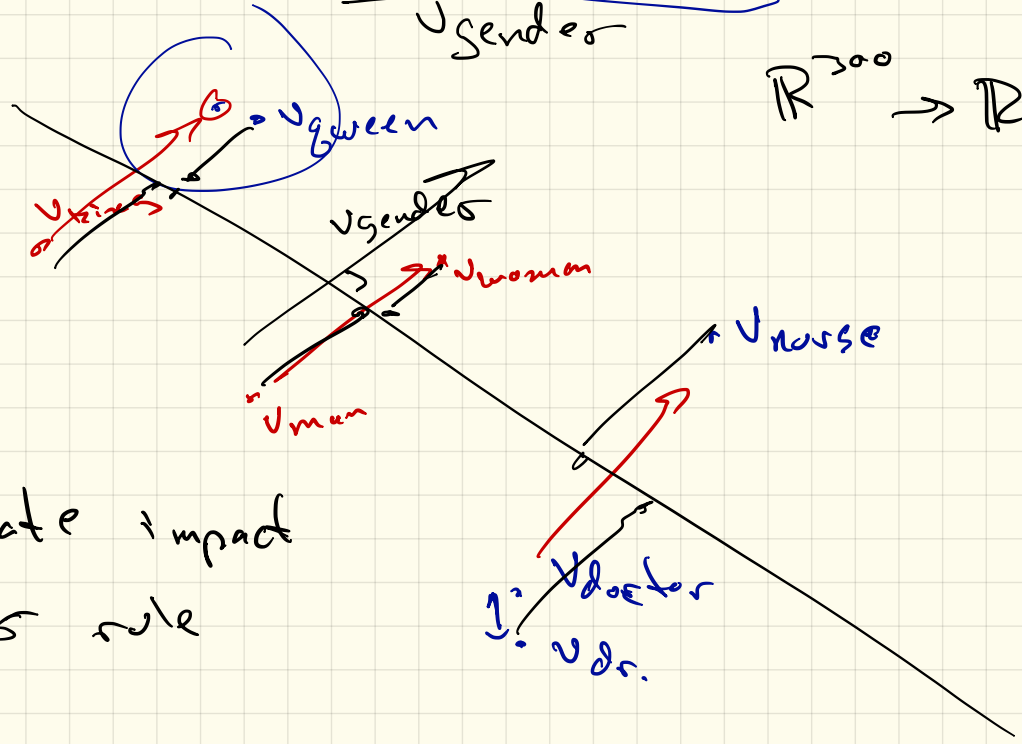
Use Cosine Sim

text  $\rightarrow$  vectors  
dim = 1,000,000  
word count  
in neighborhood

$\rightarrow$  define PPMI  $\rightarrow$  maps  $v_i$   
to  $\mathbb{R}^{300}$   
 $\max\left(0, \frac{P(i,j)}{P(i)P(j)}\right)$

$$v_{\text{king}} + \frac{(v_{\text{woman}} - v_{\text{man}})}{v_{\text{gender}}}$$

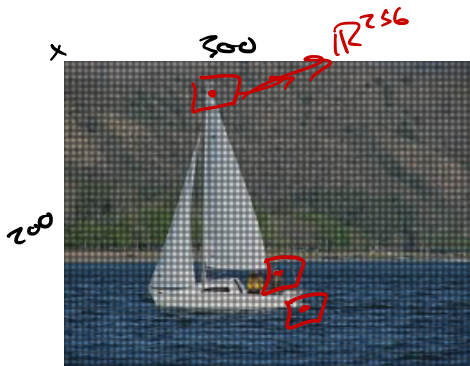
$$\mathbb{R}^{300} \rightarrow \mathbb{R}^{299}$$



disparate impact  
3/5 rule

# Images and SIFT Features

Generate in  $\mathbb{R}^d$



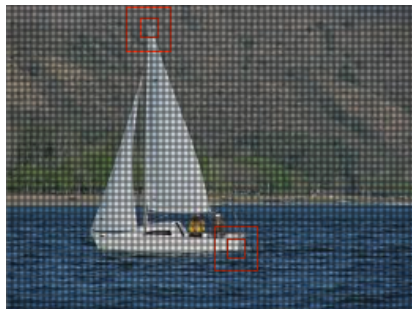
N1	N2	N3
N8	X	N4
N7	N6	N5

$\rightarrow \mathbb{R}^{200 \cdot 300 \cdot 3}$  20 yrs ago

$\rightarrow$  SIFT vector  $\rightarrow \mathbb{R}^{256}$

Euclidean 25 yrs ago

# Images and SIFT Features



N1	N2	N3
N8	X	N4
N7	N6	N5



# Quickly Find Nearest Neighbors

Large  $P \subset \mathbb{R}^d$   $|P| = n \leftarrow \text{large}$

→ Given query  $g \in \mathbb{R}^d$

↳ quickly find

$$\phi_P(g) = \arg \min_{P \in P} \|p - g\|$$

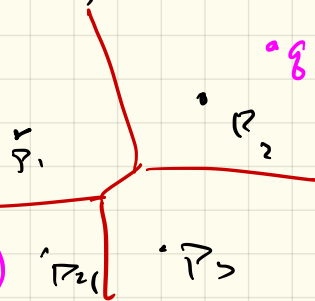
• LSH binary search on  $T$

•  $d=1$  binary tree  $O(\log n)$

•  $d=2$  Voronoi Diagram

$$O(\log n)$$

$$\text{size VD } O\left(n^{\lfloor d/2 \rfloor}\right)$$



$d \geq 3$       Aprox    NN

Find  $p \in \mathcal{P}$  s.t.

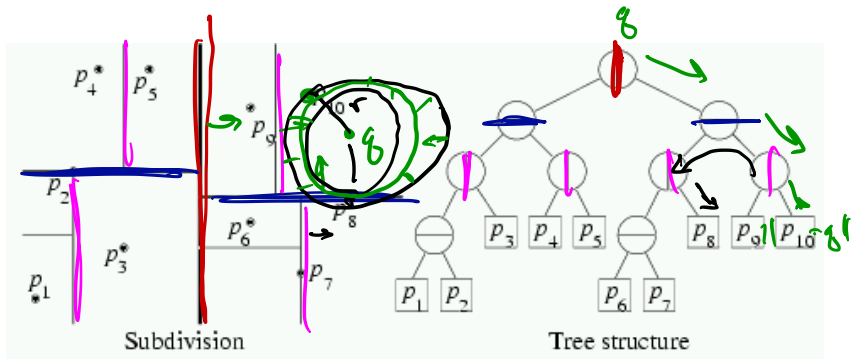
$$\|p - g\| \leq (1 + \epsilon) \|p^* - g\|$$

$$p^* = \underset{\mathcal{P}}{\operatorname{arg\,min}} (g)$$

$$d = 3 \dots 12$$

$1/2d$ -free

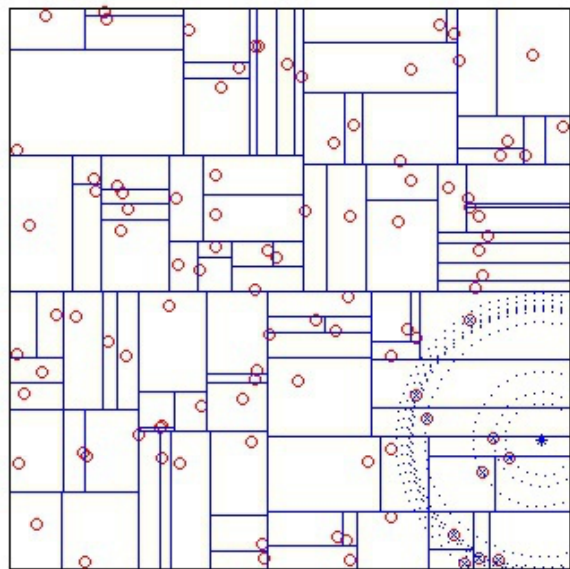
# kD-Tree



dim 2 ... 12, 20

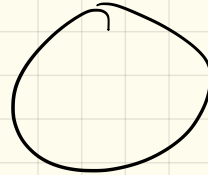
Data-adaptive kd-tree  $\rightarrow d=200$

## Approximate Queries on $k$ D-Tree

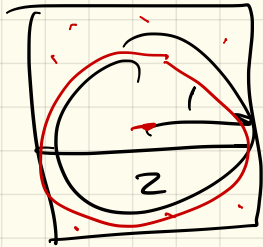




vs.



ball rad | vol



$$\frac{\pi r^d}{\Gamma(d/2 + 1)} \approx \frac{\pi r^d}{(d/2)!} \rightarrow 0$$

as  $d \rightarrow \infty$   
 $1 \cdot 2 \cdot 3 \cdots d/2$

box side length 2 vol

$$2^d$$