

## L7: Approximate Nearest Neighbors

Data  $\rightarrow \mathbb{R}^d$

Jeff M. Phillips

{Words}  
{Pictures}

January 31, 2018

# Word Vector Embeddings

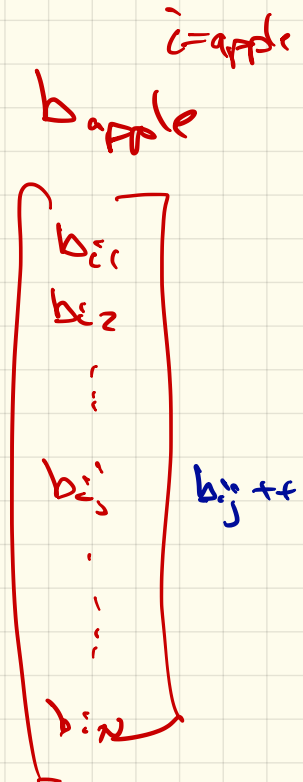
2013

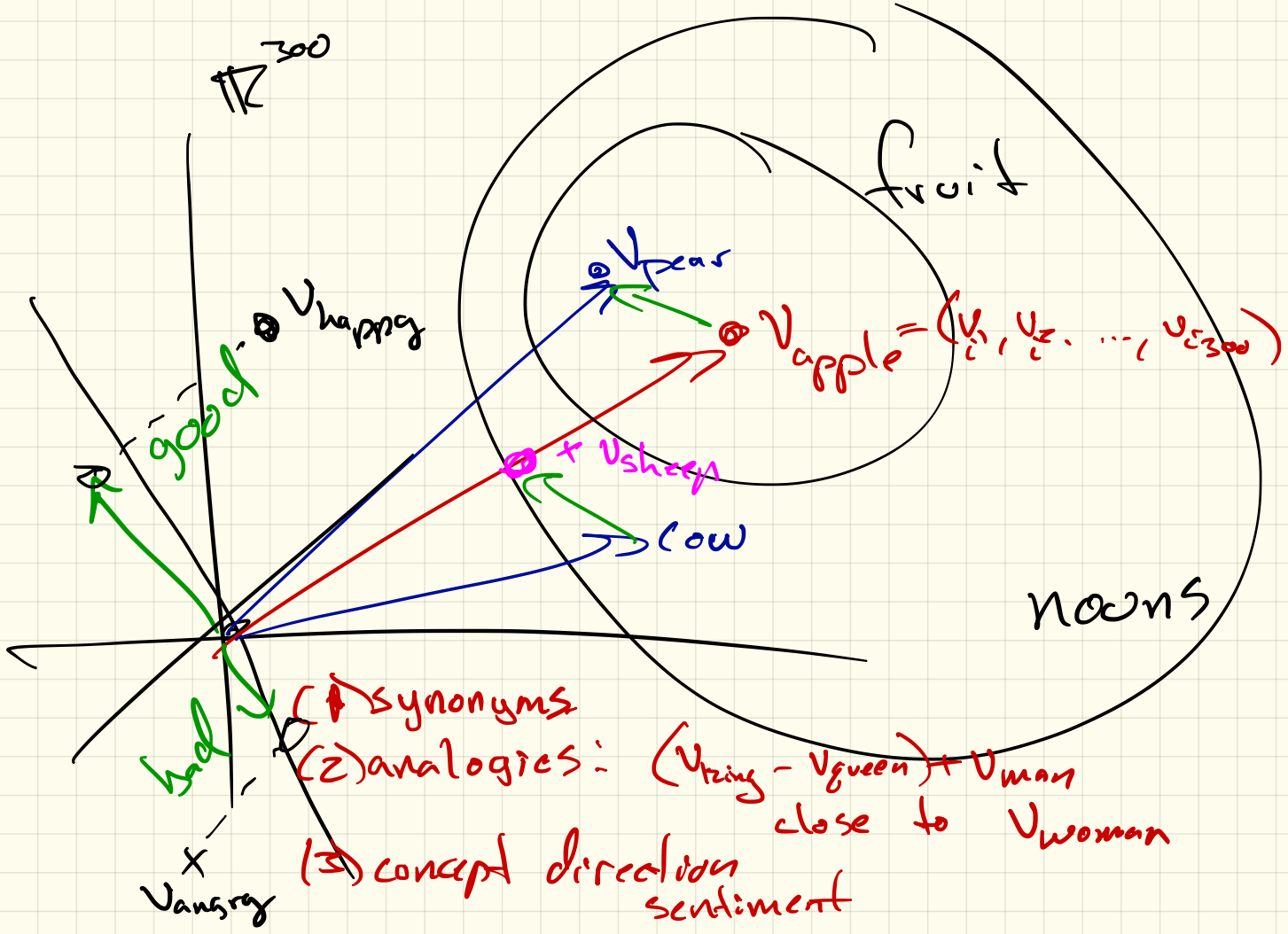
Text  
(all of Wikipedia)

... [ I ate an apple. It was good ] ...

continuous bag of words

(BOW)





# PPMI vectors $\mathbb{R}^N$

$N = 1$  million

positive, pointwise mutual information

$b_i$  vector

$$V_i = (V_{i1}, V_{i2}, \dots, V_{ij}, \dots, V_{iN})$$

$$b_i = (b_{i1}, b_{i2}, \dots, b_{ij}, \dots, b_{iN})$$

$b_{ij}$  # occurrences of  $w_j$  in  $\mathbb{R}^N$  of  $b_i$

$M =$  Total # words in corpus.

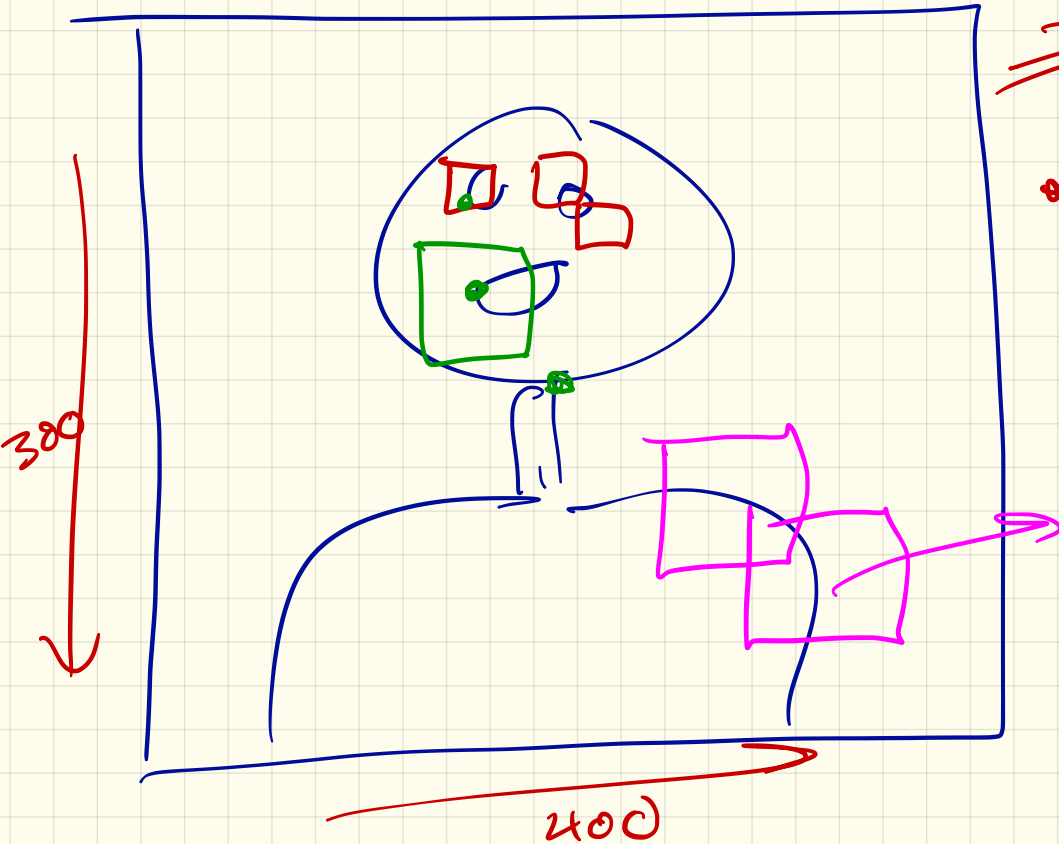
$n_i =$  # times word  $i$  occurs

$$M = \sum_{i=1}^N n_i$$

$$p(i) = \frac{n_i}{M}, \quad P(i, j) = \frac{b_{ij}}{M}$$

$$V_{ij} = \max \left\{ 0, \log \left( \frac{P(i, j)}{p(i) \cdot p(j)} \right) \right\}$$

# Embed Images

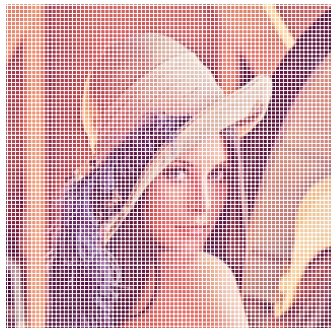


$\mathbb{R}^{120,000}$

SIFT features

$\mathbb{R}^{128}$

# Images and SIFT Features



N1	N2	N3
N8	X	N4
N7	N6	N5

# Approx Nearest Neighbors

$\mathbb{R}^1$

Sort Data  $\rightarrow$  BST

$\mathbb{R}^{2-3}$

Voronoi Diagram

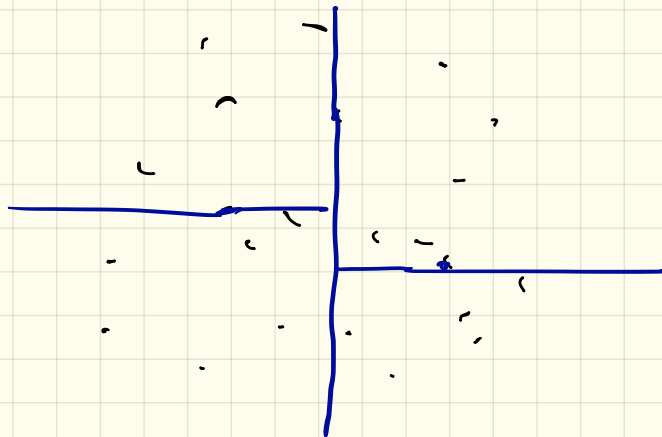
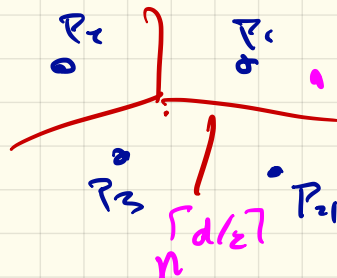
$\mathbb{R}^{3-10}$

k-d Tree

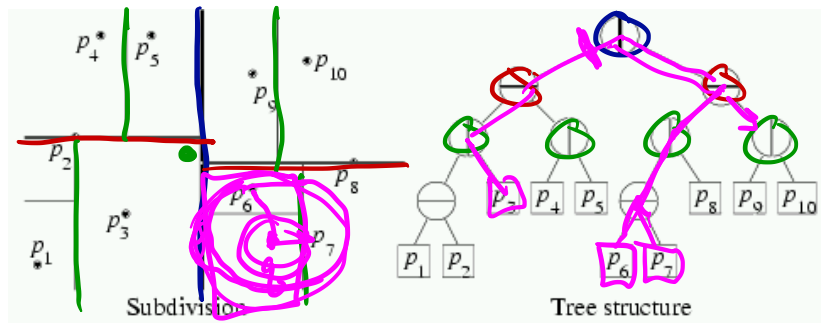
$\mathbb{R}^{10 \rightarrow 100}$

$\downarrow$

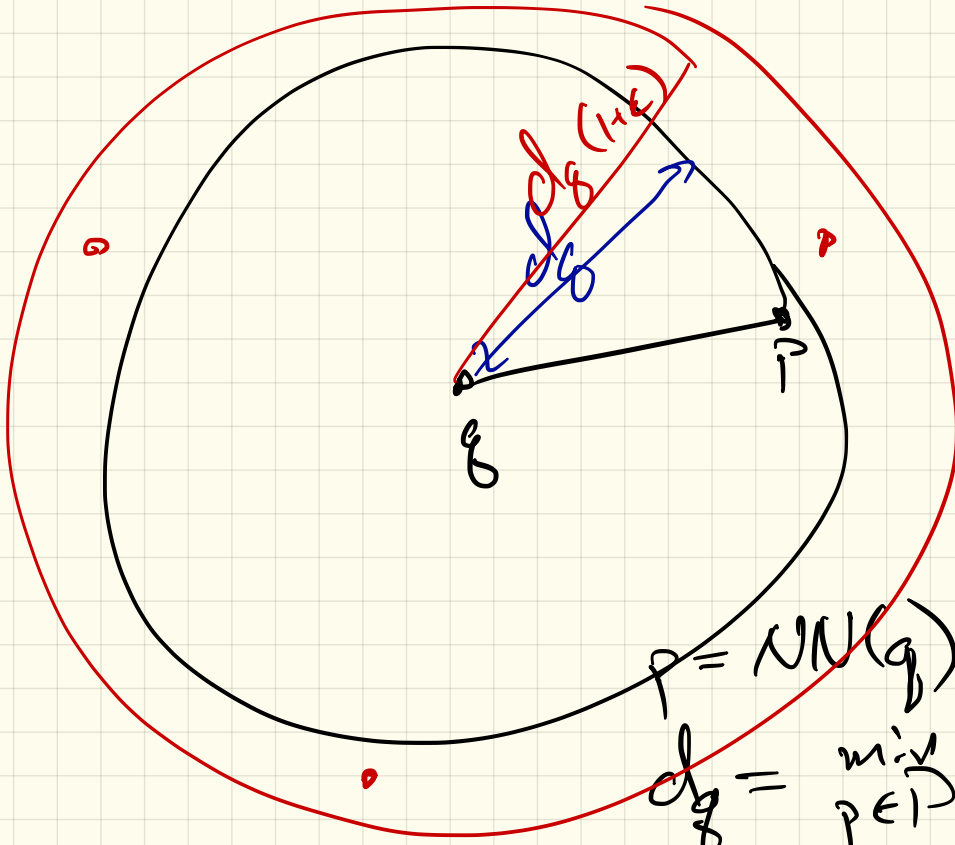
LSH



# kD-Tree



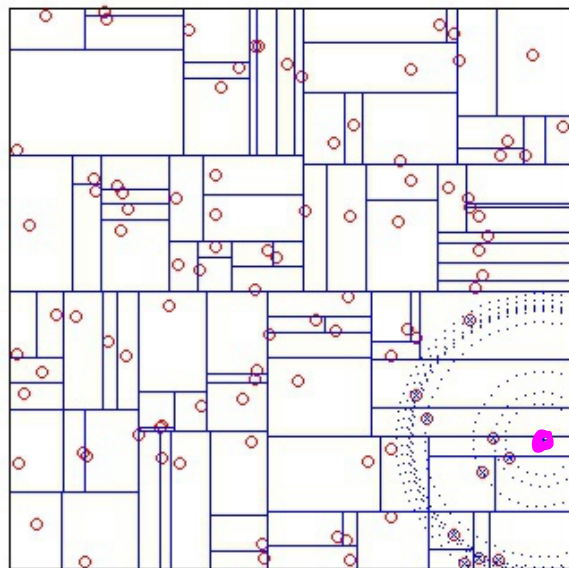




$$p = NN(g)$$

$$d_g = \min_{p \in I} d(p, g)$$

## Approximate Queries on $k$ D-Tree



# k-d-Tree

Split dependent on data.

- PCA  $\rightarrow$  1 dim

- 2-means cluster

# Ethics of Using Word Embedding

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

