

## Plan

Logistics

Review

Similarity Search

Problem set

Spm Tomorrow Gradescope

↳ Please select which question  
you're answering

## Proposal

Spm Monday Gradescope

## JL Lemma

$$x_1, \dots, x_n \in \mathbb{R}^d$$

$$\Pi \in \mathbb{R}^{k \times d} \quad \Pi_{ij} = \text{random variable}$$

$$(1-\epsilon) \|x_i - x_j\|_2^2 \leq \|\Pi x_i - \Pi x_j\|_2^2 \leq (1+\epsilon) \|x_i - x_j\|_2^2$$

wp 9/10

$$k = O\left(\frac{\log n}{\epsilon^2}\right) \leftarrow \text{no dimension dependence}$$

## Distributional JL Lemma

$$(1-\epsilon) \|x\|_2^2 \leq \|\Pi x\|_2^2 \leq (1+\epsilon) \|x\|_2^2$$

wp  $1-\delta$

$$k = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$$

what if we want to preserve inner products?

$$|\langle x_i, x_j \rangle - \langle \Pi x_i, \Pi x_j \rangle| \leq \frac{\epsilon}{2} (\|x_i\|_2^2 + \|x_j\|_2^2)$$

$$\hookrightarrow \text{using JL} \quad \hookrightarrow \|x-y\|_2^2 = \|x\|_2^2 + \|y\|_2^2 - 2\langle x, y \rangle$$

Application: Fast Set Size Estimation

$X$  = people in class

$Y$  = people who climb

$$x = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \end{bmatrix} \leftarrow \begin{array}{l} \text{Sujay} \\ \text{Iris} \end{array}$$

$$y = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} \leftarrow \text{Aidan}$$

$$\langle x, y \rangle = |X \cap Y|$$

$$\|x\|_2^2 = |X|$$

## Similarity Estimation

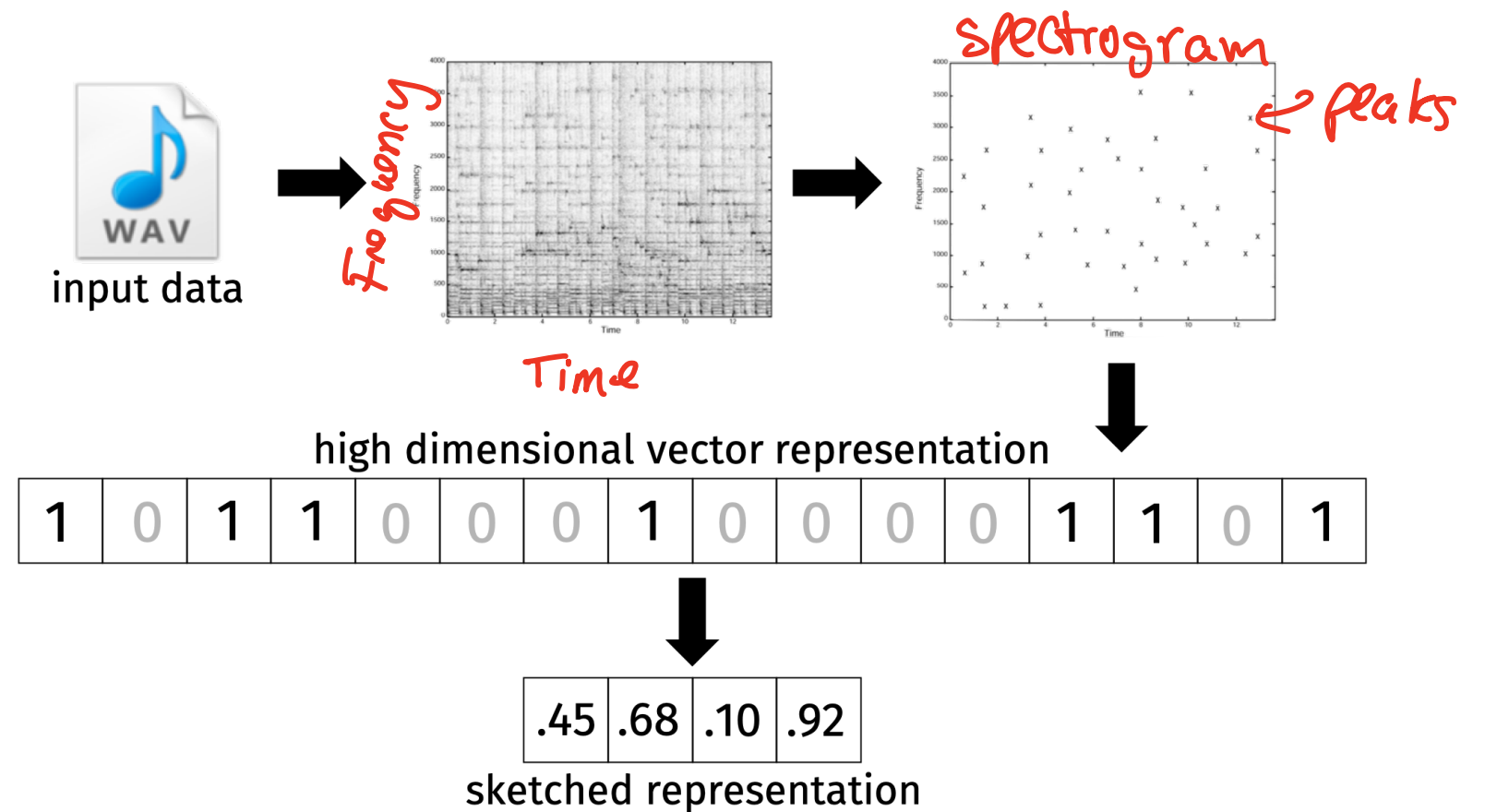
JL preserves distance.

How about "similarity"?

Shazam matches short,  
noisy clips against  
huge database

Problem: Given query  $x \in \{0,1\}^d$ , find similar song  $y \in \{0,1\}^d$

Naive:  $O(nd)$  space to store,  $O(nd)$  time to search



"Sketch" into lower dimension

$$c: \{0,1\}^d \rightarrow \mathbb{R}^k \quad k \ll d$$

$$c(x) \approx c(y) \quad \text{if } x \approx y$$

↑  
similarity

Jaccard Similarity

$$J(x,y) = \frac{|X \cap Y|}{|X \cup Y|}$$

$$= \frac{\# \text{ non-zero in common}}{\# \text{ non-zero total}}$$

e.g.:

$$x = \begin{bmatrix} 1 \\ -1 \\ 0 \\ 0 \\ 1 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ -1 \\ 0 \\ -1 \\ 1 \end{bmatrix}$$

$$J(x,y) = \frac{2}{5}$$

Also useful for

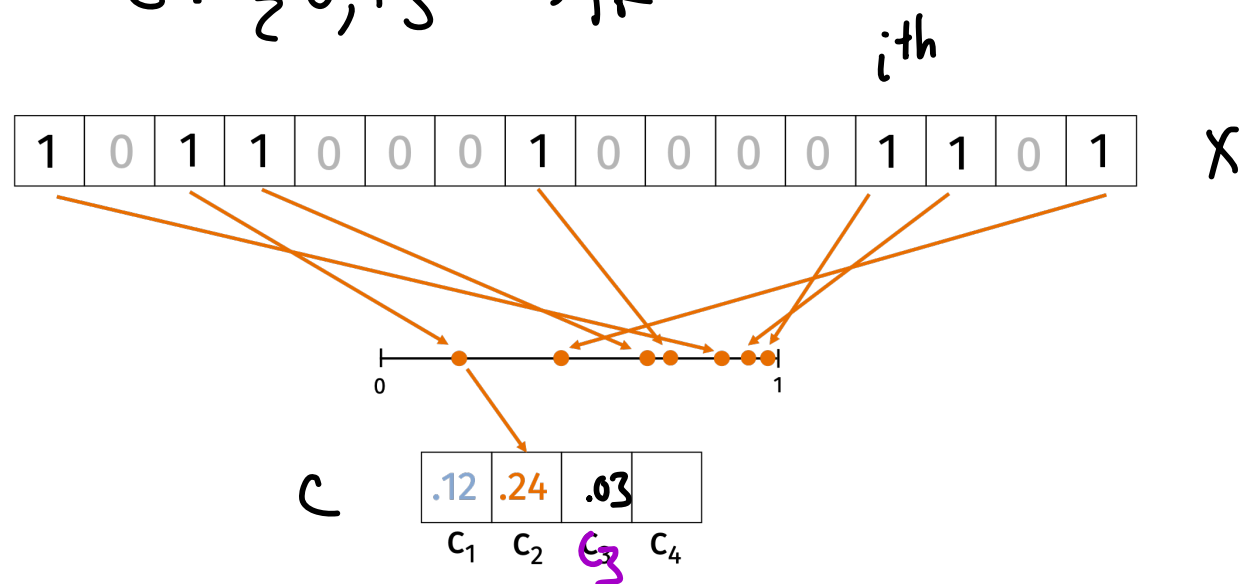
↳ "bag of words" document

↳ cached webpages

↳ earthquake detection

# Minhash

$$C: \{0,1\}^d \rightarrow \mathbb{R}^k$$



$$h_i: \{1, \dots, d\} \rightarrow [0, 1]$$

$$c_i = \min_{j: x_j = 1} h_i(j)$$

$$\Pr(c_i(x) = c_i(y)) = \frac{|X \cap Y|}{|X \cup Y|} = J(x, y)$$

Estimate  $J(x, y)$  using  $C$

$$\hat{J}(x, y) = \frac{1}{k} \sum_{i=1}^k \mathbb{1}[c_i(x) = c_i(y)]$$

$$\begin{aligned} \mathbb{E}[\hat{J}(x, y)] &= \frac{1}{k} \sum_{i=1}^k \mathbb{E}[\mathbb{1}[c_i(x) = c_i(y)]] \\ &= J(x, y) \end{aligned}$$

$$\begin{aligned} \text{Var}(\hat{J}(x, y)) &= \frac{1}{k^2} \sum_{i=1}^k \text{Var}(\mathbb{1}[c_i(x) = c_i(y)]) \\ &\leq \frac{1}{k^2} \sum_{i=1}^k J(x, y) = \frac{J(x, y)}{k} \end{aligned}$$

$$\begin{aligned} \text{Var}(\mathbb{1}) &= \mathbb{E}[\mathbb{1}^2] - \mathbb{E}[\mathbb{1}]^2 \\ &= J(x, y) - J(x, y)^2 \end{aligned}$$

## Chebyshev's

$$\Pr(|\hat{J} - J| \geq \alpha \cdot \sigma) \leq \frac{1}{\alpha^2}$$

$$\sigma = \sqrt{\text{Var}(\hat{J})} \approx \sqrt{\frac{J}{k}}$$

$$\Pr(|\hat{J} - J| \geq \alpha \cdot \frac{\sqrt{J}}{\sqrt{k}}) \leq \frac{1}{\alpha^2} \stackrel{\text{want}}{=} \delta$$

$$\Rightarrow \Pr(|\hat{J} - J| \geq \underbrace{\epsilon}_{\epsilon} \cdot \frac{1}{\sqrt{k}}) \leq \frac{1}{\alpha^2}$$

$$\Pr(|\hat{J} - J| \geq \epsilon) \leq \delta$$

$$k = O\left(\frac{1/\delta}{\epsilon^2}\right)$$

$$\epsilon = \alpha \frac{1}{\sqrt{k}}$$

$$\delta = \frac{1}{\alpha^2}$$

$$\epsilon = \frac{1}{\sqrt{\delta}} \cdot \frac{1}{\sqrt{k}}$$

$$\alpha = \frac{1}{\sqrt{\delta}}$$

$$\epsilon^2 \delta = k$$

using biased coin theorem

↳ holds if  $c_i(x) = c_i(y)$

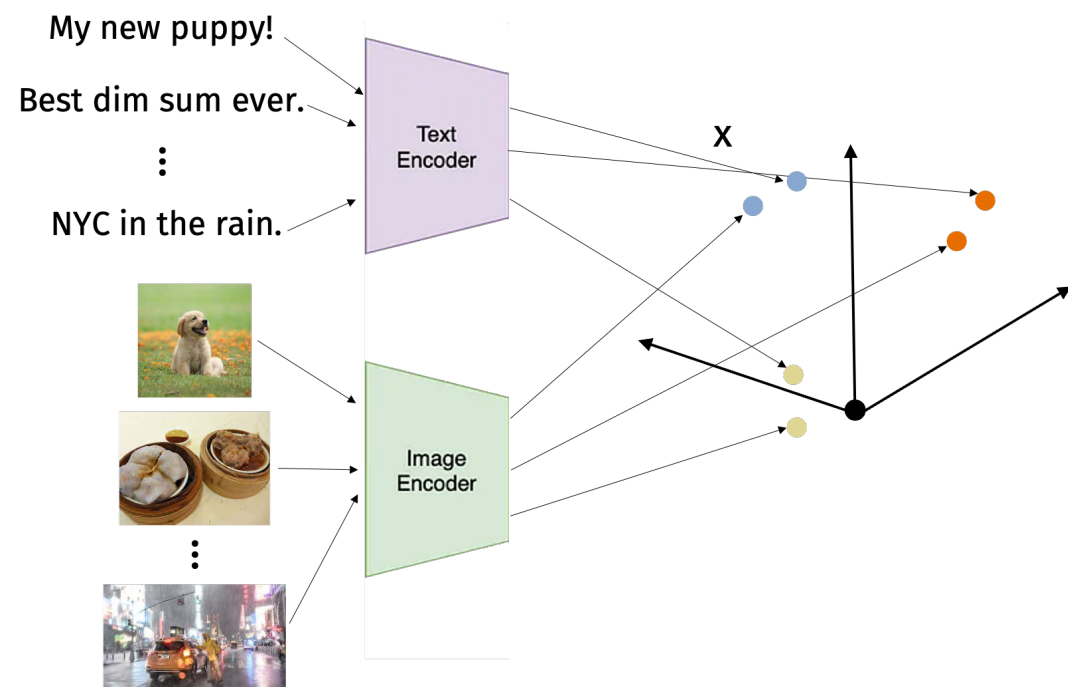
↳ bias  $b = J(x, y)$

$$k = O\left(\frac{\log 1/\delta}{\epsilon^2}\right)$$

$O(d)$  <sup>compression</sup>  $\rightarrow O(k)$  compute (approx) similarity

$O(dn)$   $\rightarrow O(k \cdot n)$  naive search

How do we find similar points faster?



## Locality Sensitive Hashing

↳ Hash function  $h: \{0,1\}^d \rightarrow \{1, \dots, m\}$

↳ Similarity function  $c$  e.g. Jaccard

↳  $h$  is locally sensitive if

$$\Pr(h(x) = h(y)) = \begin{cases} \text{large} & \text{when } x \approx y \\ \text{small} & \text{when } x \not\approx y \end{cases}$$

Our approach:

$c: \{0,1\}^d \rightarrow [0,1]$  single MinHash

$g: \mathbb{R} \rightarrow \{1, \dots, m\}$

$$h(x) = g(c(x))$$

$h(x) = h(y)$  when

(1)  $c(x) = c(y)$  or

(2)  $c(x), c(y)$  happen

to hash to the same cell

$$\Pr(h(x) = h(y))$$

$$= \Pr(c(x) = c(y)) \cdot 1$$

$$+ (1 - \Pr(c(x) = c(y))) \cdot \frac{1}{m}$$

$$\approx \Pr(c(x) = c(y))$$

↑ negligible

$$= J(x, y)$$



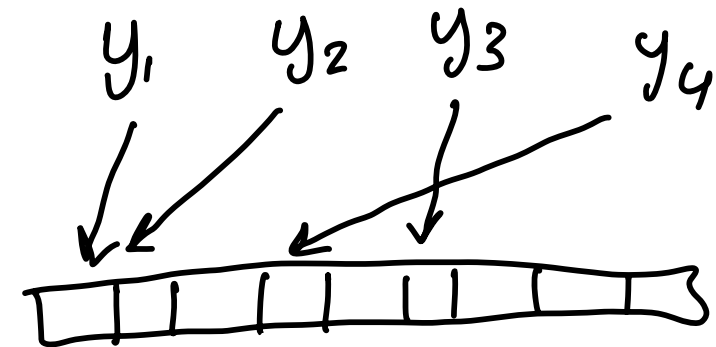
## Pre processing

Choose  $h$  by choosing  $g, c$

Create a table with  $m$  slots

For each vector, we compute

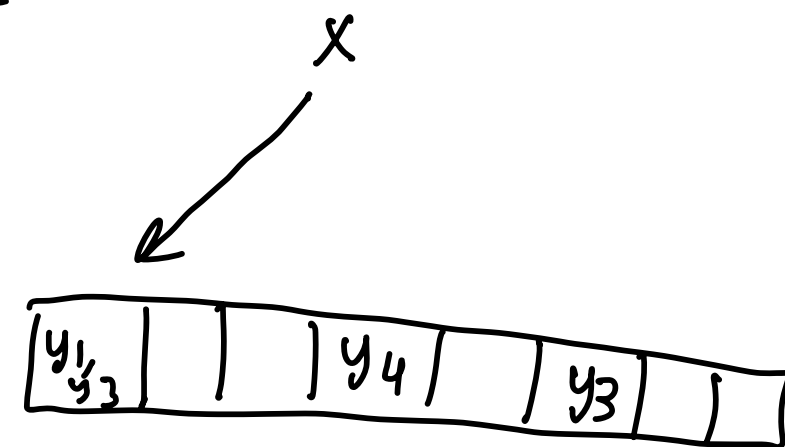
$h(y)$  and store in corresponding slot



## Query

Compute  $h(x)$  and look in

corresponding slot



Repeat with  $t$  tables

Two questions:

↳ False negative: what's the probability we don't find similar vectors?

↳ False positive: what's the probability we do find a non-similar vector?

$$\begin{aligned} \Pr(\text{find } y) &= 1 - \Pr(y \text{ not in slot of table})^t \\ &= 1 - \Pr(h_i(x) \neq h_i(y))^t = 1 - (1 - J(x, y))^t \end{aligned}$$

When  $J(x, y) = .4$ ,  $t = 10$ ,  $\Pr(\text{find } y) = 1 - (1 - .4)^{10} \approx .99 \quad \ddot{\smile}$

When  $J(x, y) = .2$ ,  $t = 10$ ,  $\Pr(\text{find } y) = 1 - (1 - .2)^{10} \approx .89 \quad \ddot{\smile}$

## Our Approach

$$c_1, \dots, c_r : \{0, 1\}^d \rightarrow [0, 1]$$

$$g : [0, 1]^r \rightarrow \{1, \dots, m\}$$

$$h(x) = g(c_1(x), \dots, c_r(x))$$

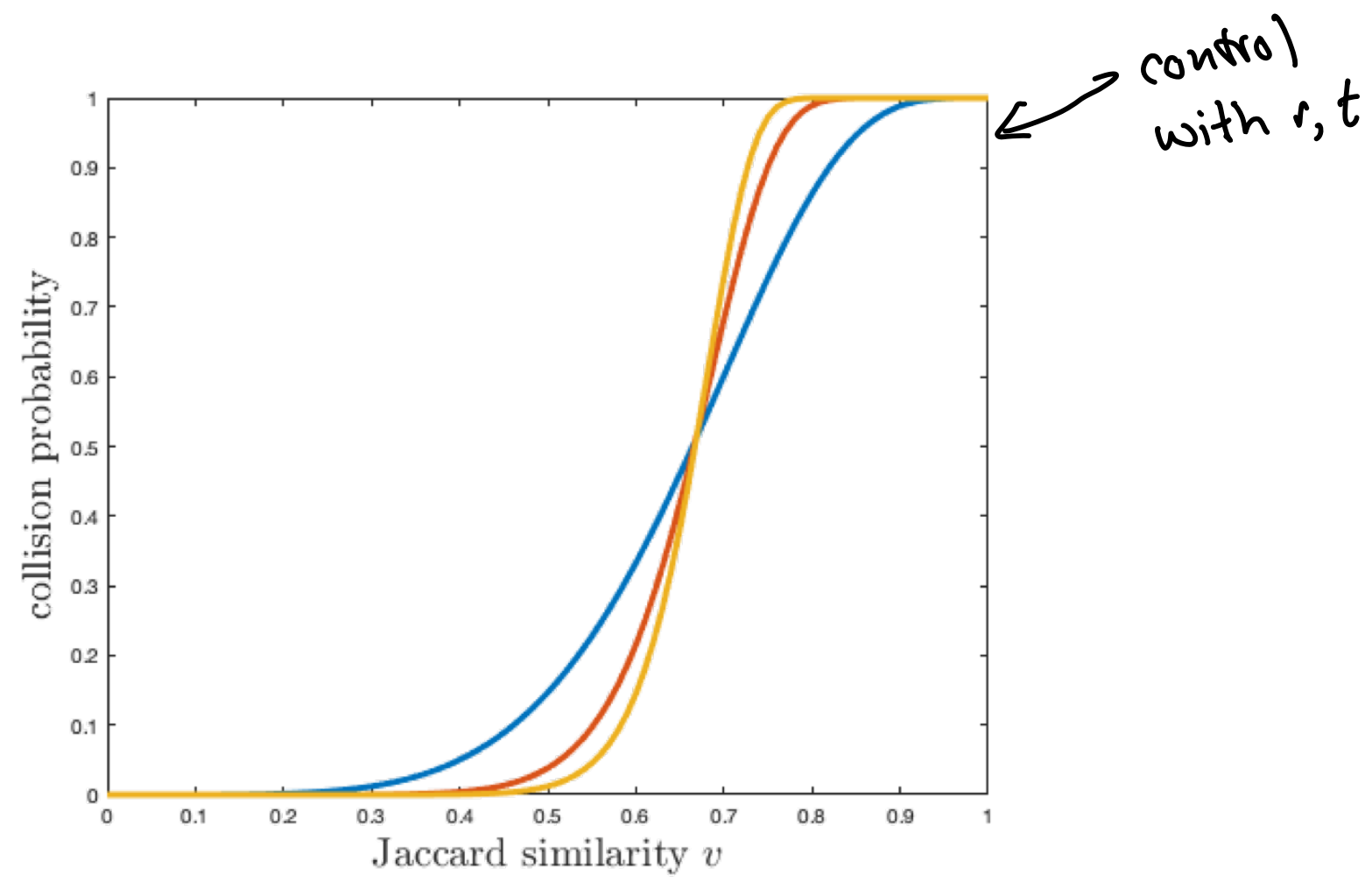
$$\Pr(h(x) = h(y)) = \Pr(c_i(x) = c_i(y) \quad \forall i) + \frac{1}{m} \swarrow \text{negligible}$$

$$\stackrel{\text{indep}}{=} \Pr(c_1(x) = c_1(y)) \Pr(c_2(x) = c_2(y)) \dots \Pr(c_r(x) = c_r(y))$$

$$\stackrel{=}{=} J(x, y)^r$$

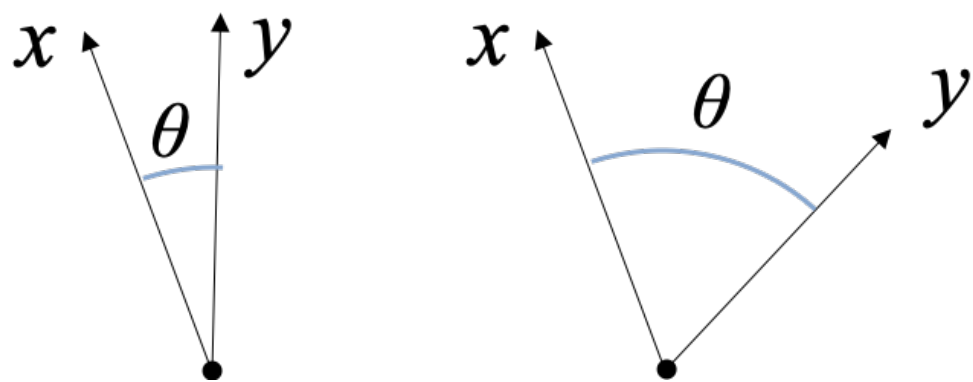
$$\Pr(\text{find } y) = 1 - \Pr(\text{not find } y \text{ in table})^t$$

$$\stackrel{=}{=} 1 - (1 - J(x, y)^r)^t$$



## Cosine Similarity

$$\cos(\theta(x,y)) = \frac{\langle x, y \rangle}{\|x\|_2 \|y\|_2}$$



"inverse to distance"

$$\begin{aligned} \|x-y\|_2^2 &= \|x\|_2^2 - 2\langle x, y \rangle + \|y\|_2^2 \\ &= \|x\|_2^2 + \|y\|_2^2 - 2\|x\|_2 \|y\|_2 \cos\theta(x,y) \end{aligned}$$

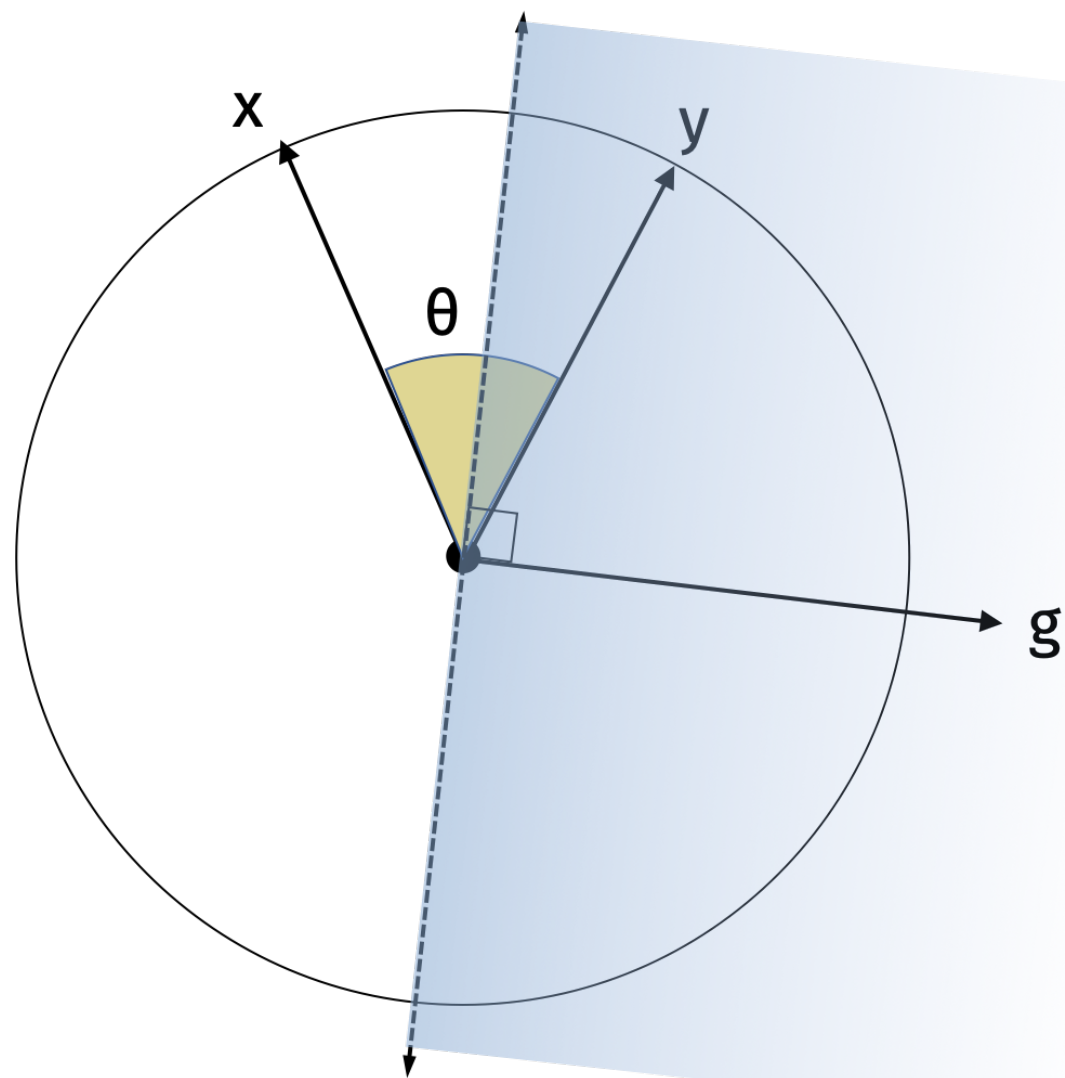
Let  $g_1, \dots, g_r \in \mathbb{R}^d$   
random vectors with  $\mathcal{N}(0,1)$

$$f: \{-1, 1\}^r \rightarrow \{1, \dots, m\}$$

$$h: \mathbb{R}^d \rightarrow \{1, \dots, m\}$$

$$h(x) = f(\text{sign}(\langle g_1, x \rangle), \dots, \text{sign}(\langle g_r, x \rangle))$$

$$\begin{aligned} \Pr(\text{sign}(\langle g_i, x \rangle) = \text{sign}(\langle g_i, y \rangle)) \\ = \end{aligned}$$



$\text{sign}(\langle g, x \rangle)$  = which side of hyperplane  
 $\text{sign}(\langle g, y \rangle)$  = which side of hyperplane

$$\text{Pr}(\text{different sides}) = \frac{2\theta}{2\pi} = \frac{\theta}{\pi}$$

$$\text{Pr}(\text{same}) = 1 - \frac{\theta}{\pi}$$

$$= \text{Pr}[\text{sign}(\langle g, x \rangle) = \text{sign}(\langle g, y \rangle)]$$

$$\begin{aligned}
 \text{Pr}(\text{find } y) &= 1 - \text{Pr}(y \text{ not in table})^t = 1 - \text{Pr}(h_i(x) \neq h_i(y))^t \\
 &= 1 - (1 - \text{Pr}(h(x) = h(y)))^t = 1 - (1 - \text{Pr}(\text{same})^r)^t = 1 - (1 - (1 - \frac{\theta}{\pi})^r)^t
 \end{aligned}$$