

Reminder: project proposals due Monday

Last time:

Lossy compression:

- Bloom filters
- Heavy hitters & count-min sketch

Key idea:

use magic of hashing
and sacrifice a little bit
of correctness \Rightarrow
significant space savings

Today

- short reviews from probability
 - variance & tail bounds
 - Gaussians & CLT
- Distinct elts
- Similarity search & dimension reduction

Distinct Elements

a_1, a_2, a_3, \dots each $a_i \in U$
 goal: maintain approx count of number of distinct elts seen.

n_t : # distinct elts seen in a_1, a_2, \dots, a_t

do this approximately.

$h: U \rightarrow [0, 1]$

$\tilde{h}: U \rightarrow \{0, 1, \dots, m-1\}$

set $h(x) = \frac{i}{m}$

track $Y = \min_{1 \leq i \leq t} h(a_i)$



	a_1	a_2	a_3
	32	5	17	32	14	5	17
h	0.43	0.06	0.19	0.43	0.19	0.06	0.19

suppose that t distinct elts seen

$E(Y) = E[\min \text{ of } n \text{ hash values}]$

$= \frac{1}{n+1}$

if magically

$Y \approx \frac{1}{n+1}$ not true.

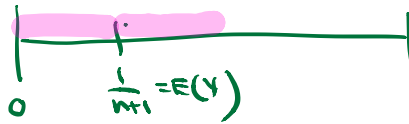
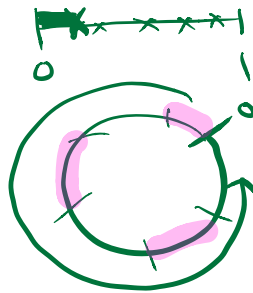
$n+1 = \frac{1}{Y}$

estimate $n = \frac{1}{Y} - 1$

$Var(Y) \approx \frac{1}{(n+1)^2}$

$\sigma(Y) = \frac{1}{n+1}$

- $h(32) = 0.43$
- $h(5) = 0.06$
- $h(17) = 0.19$
- $h(14) = 0.85$



	a_1	a_2	a_3
Y_1	32	5	17	32	14	5	17	5
Y_2	0.43	0.43	0.19	0.19	0.19	0.19	0.19	0.19
Y_k								
								h_1
								h_2
								h_k

$$Y_j = \min_{1 \leq i \leq t} h_j(a_i)$$

$$E[Y_j] = \frac{1}{n+1}$$

$$\text{Var}(Y_j) = \frac{1}{(n+1)^2}$$

$$\bar{Y} = \frac{1}{K} \sum_{j=1}^K Y_j$$

$$E(\bar{Y}) = \frac{1}{n+1}$$

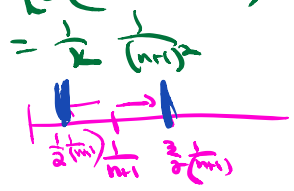
$$\text{Var}(\bar{Y}) = \frac{1}{K} \frac{1}{(n+1)^2}$$

$$\text{Var}(\bar{Y}) = \frac{1}{K^2} (K \text{Var}(Y_j))$$

$$k=16 \quad \sigma^2 = \frac{1}{16} \cdot \frac{1}{(n+1)^2}$$

$$\sigma = \frac{1}{4} \cdot \frac{1}{(n+1)}$$

$$2\sigma = \frac{1}{2} \frac{1}{(n+1)}$$



$$\Pr\left(|\bar{Y} - \frac{1}{n+1}| > 2\sigma\right) \leq \frac{1}{4}$$

w.p. $\geq \frac{3}{4}$

$$\frac{1}{2} \frac{1}{(n+1)} \leq \bar{Y} \leq \frac{3}{2} \frac{1}{(n+1)}$$

$$\frac{1}{2\bar{Y}} \leq n+1 \leq \frac{3}{2} \frac{1}{\bar{Y}}$$

$$\Pr(|Y - n| > c\sigma) \leq \frac{1}{c^2}$$

$$\frac{1}{4} \leq 2(n+1) \leq \frac{3}{2}(n+1) \leq \frac{1}{\bar{Y}}$$

Similarity search:

dataset, notion of similarity between items in set.

Items
documents

web pages

DNA sequences

movie tastes

ML-classification

Goal

plagiarism detection
similar topic

mirror sites

find similar genes

similar

What is an item (data pt)?

vector in high-dimensional space \mathbb{R}^k

k very large.

(2) w_1, w_2, \dots

w_w bag of words

vector representing pixel values image

Similarity

① ℓ_2 distance

$$d(x, y) = \sqrt{\sum_{i=1}^k (x_i - y_i)^2}$$

$x, y \in \mathbb{R}^k$
 $x = (x_1, x_2, \dots, x_k)$

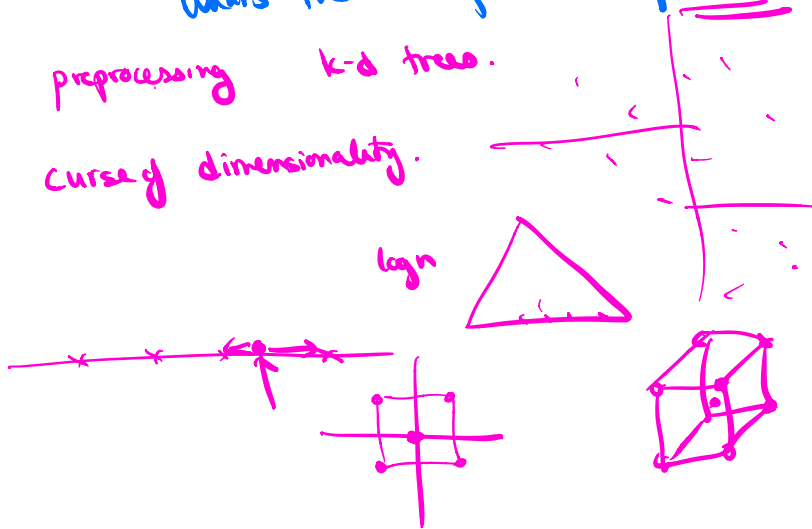
n vectors, $\in \mathbb{R}^k$
compute all distances:

$$O(n^2 \cdot k)$$

Preprocess a data set consisting of n vectors
 & then be able to quickly answer queries
 What's the closest pt to $q \in \mathbb{R}^k$

preprocessing k-d trees.

curse of dimensionality.



Dimension reduction:

$$(x_1, \dots, x_n \in \mathbb{R}^k)$$

$$d \ll k$$

$$f: \mathbb{R}^k \rightarrow \mathbb{R}^d$$

so that distances between pts
 are almost exactly preserved.

- compression of data
- Visualization

$$x_1, \dots, x_n \in \mathbb{R}^k \rightarrow \underbrace{f(x_1), \dots, f(x_n)}_{\in \mathbb{R}^d} \leftarrow \text{dist}(x_i, x_j)$$

$$\forall i, j \quad \text{dist}(x_i, x_j) \approx \text{dist}(f(x_i), f(x_j)) \quad \text{w.h.p.}$$



r at random

$$\mathbb{R}^k \rightarrow \mathbb{R}$$

$$x \rightarrow r \cdot x$$

$$y \rightarrow r \cdot y$$

$$E[(r \cdot x - r \cdot y)^2] = \|x - y\|^2$$

$(r \cdot x - r \cdot y)^2$ is an unbiased estimator for $\|x - y\|^2$

$$f_r(x) = r \cdot x = \sum_{i=1}^k r_i x_i$$

$$r = (r_1, \dots, r_k) \quad x = (x_1, \dots, x_k)$$

$r_i \sim N(0, 1)$ indep.

$$f_r(x) - f_r(y) = r \cdot x - r \cdot y = r \cdot (x - y) = \sum_{i=1}^k r_i (x_i - y_i)$$

$$\sum_{i=1}^k r_i (x_i - y_i)$$

$N(0, 1)$

$$\|x - y\| = \sqrt{\sum_{i=1}^k (x_i - y_i)^2}$$

$$N(0, \sum_{i=1}^k (x_i - y_i)^2)$$

$$= N(0, \|x - y\|^2)$$

$$E[(f_r(x) - f_r(y))^2]$$

$$= \text{Var}(f_r(x) - f_r(y))$$

$$= \|x - y\|^2$$

$$Z = \sum_{i=1}^k z_i^2$$

$Z \sim N(0, 1)$
 $Z \sim N(\sum_{i=1}^k r_i^2, \sum_{i=1}^k z_i^2)$

$$\text{Var}(w) = E[(w - \mu)^2] = E[w^2] - \mu^2$$

$$r^{(1)}, r^{(2)}, \dots, r^{(d)}$$

$$r^{(i)} \in \mathbb{R}^k$$

$$r_j^{(i)} \sim N(0, 1)$$

$$f_r(x) = r \cdot x$$

$$f_{r^{(i)}}(x) = c r^{(i)} \cdot x$$

$$x \rightarrow f(x) = \begin{pmatrix} c r^{(1)} \cdot x \\ c r^{(2)} \cdot x \\ \vdots \\ c r^{(d)} \cdot x \end{pmatrix}$$

$$f: \mathbb{R}^k \rightarrow \mathbb{R}^d$$

$$\|f(x) - f(y)\|^2 = \sum_{i=1}^d (c r^{(i)} \cdot x - c r^{(i)} \cdot y)^2 = c^2 \sum_{i=1}^d (r^{(i)} \cdot (x - y))^2 = \sum_{i=1}^d (f_{r^{(i)}}(x) - f_{r^{(i)}}(y))^2$$

$$f(x) = A x$$

↑

$$\underbrace{\begin{bmatrix} r^{(1)} \\ r^{(2)} \\ \vdots \\ r^{(d)} \end{bmatrix}}_{d \times k \text{ matrix}} \begin{bmatrix} x_1 \\ \vdots \\ x_k \end{bmatrix} = \underbrace{\begin{bmatrix} r^{(1)} x \\ r^{(2)} x \\ \vdots \\ r^{(d)} x \end{bmatrix}}_{\text{vector}}$$

$$f(y) = \begin{pmatrix} c r^{(1)} \cdot y \\ c r^{(2)} \cdot y \\ \vdots \\ c r^{(d)} \cdot y \end{pmatrix}$$

$$E(\|f(x) - f(y)\|^2) = c^2 \sum_{i=1}^d E \left[\underbrace{\left[r^{(i)} \cdot (x-y) \right]^2}_{\|x-y\|^2} \right]$$

$$= c^2 \cdot d \cdot \|x-y\|^2$$

want this = $\|x-y\|^2$

so should set $c = \frac{1}{\sqrt{d}}$

If you choose $d \approx \frac{2 \log n}{\epsilon^2}$

then \forall 2 pts x, y in dataset of size n .

w.h.p.

$$\|f(x) - f(y)\| \leq (1 + \epsilon) \|x - y\|$$