

lec 12

Aside if X rv $\in \mathbb{R}$ $E(X) = \sum_k k \cdot P(X=k)$
 if $X \geq 0 = \sum_{k=1}^{\infty} P(X \geq k)$

$X \in \mathbb{R}$ then $E(X) = \int_{-\infty}^{\infty} u \cdot P(u) du$

if $x > 0 = \int_0^{\infty} \Pr(X \geq u) du$

$\text{Var}(x) = E(x^2) - E(x)^2$

$x = 0, 1, p$
 $\rightarrow \text{Var}(x) = p(1-p)$

if X, Y ind $\Rightarrow \text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y)$

Document duplicate detection

- Hashing for exact-duplicate detection:

- Pick $h: \text{Docs} \rightarrow \{0,1\}^b$ from 2-wise independent family.
- $\Pr(h(A) = h(B)) = \begin{cases} 1 & \text{if doc. } A \text{ is identical to doc. } B \\ 2^{-b} & \text{if } A \neq B \end{cases}$

\rightarrow only for exact matches

Consider \rightarrow a document is a sequence of words. Let's transform into set that captures size of this

("quick brown fox jumped") \rightarrow of "quick brow", "brown fox", ...]

\rightarrow so now we care abt set similarity

\rightarrow similarity \rightarrow if $A = B \Rightarrow \text{sim}(A, B) = 1$
 if $A \cap B = \emptyset \Rightarrow \text{sim}(A, B) = 0$

\rightarrow we will use Jaccard Similarity: $\text{sim}(A, B) = \frac{|A \cap B|}{|A \cup B|}$

MinHash

- Pick $h: [u] \rightarrow [m]$ (want $m \gg$ size of any set)
 - The basic sketch: $sk(S) = \min\{h(x) \mid x \in S\}$.
 - To compare sketches of two sets S_1, S_2 :
- $$Y = \begin{cases} 0 & \text{if } sk(S_1) \neq sk(S_2) \\ 1 & \text{if } sk(S_1) = sk(S_2) \end{cases}$$
- What is $E(Y)$?
 - $- Y = 1$ if
 - The element in $x \in S_1 \cup S_2$ with minimum $h(x)$ lies in $S_1 \cap S_2$.
 - The elements $x_1 \in S_1, x_2 \in S_2$ with minimum hash values have $h(x_1) = h(x_2)$ but $x_1 \neq x_2$.
 - $- E(Y) \leq \frac{|S_1 \cap S_2|}{|S_1 \cup S_2|} + \frac{|S_2|}{m} = \text{sim}(S_1, S_2) + o(1)$.

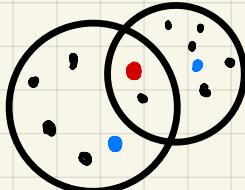
$E(Y) = \Pr(Y=1) \Rightarrow \Pr(\text{sk}(S_1) = \text{sk}(S_2))$

\Rightarrow we want minimal val in $S_1 \cup S_2$ be in $S_1 \cap S_2 \rightarrow$ std way red

$\text{CD} \Pr = \text{sim}(A, B)$

or we could just map 2 different things that hash to same min val...

$\hookrightarrow \Pr \leftarrow \frac{|S_2|}{m}$



in big news

as in expectation

But the above isn't great \rightarrow y is only 0,1 so we will get its value in expectation by reducing variance

MinHash

- Y is an unbiased estimate of $\text{sim}(S_1, S_2)$ but it has a large variance. How can variance be reduced?
- The sketch: $SK(S) = (sk_1(S), sk_2(S), \dots, sk_t(S))$
- Estimate is $\tilde{\text{sim}}(SK(S_1), SK(S_2)) = (\sum_{i=1}^t Y_i)/t$, where:
$$Y_i = \begin{cases} 0 & \text{if } sk_i(S_1) \neq sk_i(S_2), \\ 1 & \text{if } sk_i(S_1) = sk_i(S_2). \end{cases}$$
- $E(\tilde{\text{sim}}(SK(S_1), SK(S_2))) = \text{sim}(S_1, S_2) + o(1)$.
- $\text{Var}(\tilde{\text{sim}}(SK(S_1), SK(S_2))) = \frac{1}{t} \text{sim}(S_1, S_2)(1 - \text{sim}(S_1, S_2)) + o(1)$
- Space?
- $t \cdot \log m$ bits.
- Output per sketch

$$\text{Var}\left(\frac{1}{t} \sum_{i=1}^t Y_i\right) = \frac{1}{t^2} \sum_{i=1}^t \text{Var}(Y_i)$$

ind iid

$$= \frac{1}{t} \text{Var}(Y_i)$$

$$= \frac{1}{t} \text{sim}(S_1, S_2)(1 - \text{sim}(S_1, S_2))$$

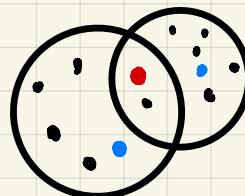
Can reduce by $\uparrow t$

How do we reduce space

MinHash

- Reducing space: pick hash functions
 - $h: [u] \rightarrow [u]$
 - $g: [u] \rightarrow \{0,1\}$
 - Basic sketch: $sk(S) = g(x)$, where $x = \text{argmin}_{x \in S}(h(x))$.
- $SK(S) = (sk_1(S), sk_2(S), \dots, sk_t(S))$
- $Y_i = \begin{cases} 0 & \text{if } sk_i(S_1) \neq sk_i(S_2), \\ 1 & \text{if } sk_i(S_1) = sk_i(S_2). \end{cases}$
- $E(Y_i) = \text{sim}(S_1, S_2) + (1 - \text{sim}(S_1, S_2))/2 = (1 + \text{sim}(S_1, S_2))/2$.
- $\tilde{\text{sim}}(SK(S_1), SK(S_2)) = (\sum_{i=1}^t (2Y_i - 1))/t$.
- $E(\tilde{\text{sim}}(SK(S_1), SK(S_2))) = \text{sim}(S_1, S_2)$.
- $\text{Var}(\tilde{\text{sim}}(SK(S_1), SK(S_2))) = \frac{4}{t} \cdot \frac{(1 + \text{sim}(S_1, S_2))}{2} \cdot \frac{(1 - \text{sim}(S_1, S_2))}{2}$.
- Space?
- t bits.

Case 1: if min x in intersection then
 $\Pr(Y_i = 1) = \text{sim}(S_1, S_2)$



Case 2: if outside its 50/50 chance of 0.

$$(1 - \text{sim}(S_1, S_2)) \frac{1}{2}$$

$$\Rightarrow E(Y_i) = \text{sim}(S_1, S_2) + \frac{(1 - \text{sim}(S_1, S_2))}{2}$$

$$= \frac{(1 + \text{sim}(S_1, S_2))}{2}$$

$$\Rightarrow Y = \frac{1}{t} \sum_{i=1}^t (2Y_i - 1)$$

in exp it is $\text{sim}(S_1, S_2)$

$$\text{Var}(Y) = \frac{4}{t} \cdot \left(\frac{1 + \text{sim}(S_1, S_2)}{2} \right) \left(\frac{1 - \text{sim}(S_1, S_2)}{2} \right)$$

Spin practice $(sk_1, \dots, sk_t) \rightarrow$ min t high values not min for + hash function

Distinct Elements in Data Streams

- Input is a stream (multiset) $x_1, x_2, \dots, x_N \in [u]^N$.
- How many distinct elements were seen?
 - Estimate $n = |\{x_1, x_2, \dots, x_N\}|$: return some $\hat{n} \approx n$.
- Assume we have an ideal hash function $h: [u] \rightarrow [0,1]$.
- At time t , maintain $X = \min(h(x_1), \dots, h(x_t))$.

$$\text{If } n=1 \rightarrow E(X) = \frac{1}{2}$$

$$n=2 \rightarrow E(X) = \frac{1}{3}$$

:

$$n=100 \rightarrow E(X) = \frac{1}{100}$$

Distinct Elements in Data Streams

- What is $E(X) = ?$
- What is $\text{Var}(X) = E((X - E(X))^2) = E(X^2) - E(X)^2 = ?$
- $E(X^2) = \int_0^1 \Pr(X^2 > u) du = \int_0^1 u^n du = \frac{1}{n+1}$
- $= \int_0^1 (1 - \sqrt{u})^n du = \int_0^1 (1 - v)^n 2v dv = 2 \int_0^1 v^n (1 - v) dv$
- $= 2 \left(\frac{1}{n+1} - \frac{1}{n+2} \right) = \frac{2}{(n+1)(n+2)}$
- $\text{Var}(X) = \frac{2}{(n+1)(n+2)} - \frac{1}{(n+1)^2} < \frac{1}{(n+1)^2}$
- $Y = \frac{(X_1 + \dots + X_m)}{m}$ (avg. of m independent sketches)
- $\text{Var}(Y) = \frac{1}{m^2} \text{Var}(X_1 + \dots + X_m) = \frac{1}{m} \text{Var}(X_1) < \frac{1}{m} \cdot \frac{1}{(n+1)^2}$.

(Markov's Ineq.: $X \geq 0$) $\Pr(X \geq t) \leq \frac{E(X)}{t}$.

(Chebychev: apply Markov to 2nd moment)

$$\Pr(|X - E(X)| \geq t) \leq \frac{E((X - E(X))^2)}{t^2} = \frac{\text{Var}(X)}{t^2}$$

Estimate $\frac{1}{Y} \in \left[\frac{n+1}{1+\epsilon}, \frac{n+1}{1-\epsilon} \right]$ whenever $Y \in \frac{1}{n+1} \pm \frac{\epsilon}{n+1}$.

$$\Pr\left(\left|Y - \frac{1}{n+1}\right| > \frac{\epsilon}{n+1}\right) < \frac{\text{Var}(Y)}{\left(\frac{\epsilon}{n+1}\right)^2} < \frac{1}{\left(\frac{\epsilon}{n+1}\right)^2} = \frac{1}{m\epsilon^2}$$

So m grows quad with ϵ

for exp

$$\Pr\left(\hat{n} \in \left[\frac{1}{1-\epsilon}, \frac{1}{1+\epsilon}\right]\right)$$

So no big as \hat{n} not too far away, we're good