Lecture 5 — Extensions

Recall:

Any LSH function satisfies the triangle inequality. Proof via

$$\delta_{x,x''} \leq \delta_{xx'} + \delta_{x'x''}$$

Lemma: Jaccard mean is a kernel (and so are other similarities)

Proof: $\delta_{xx'} = \langle e_x, e_{x'} \rangle$ □

Lemma: Any $sim(x,x') = \mathbb{E}_h\left[\delta_{h(x),h(x')}\right]$ can be embedded as binary mapping

Proof: Construct new hash function

$$G(h(x)) \in \{0,1\}$$

$$\mathbb{E}_{h,b}\left[\delta_{G(h(x)),G(h(x'))}\right] = \frac{1}{2} + \frac{1}{2} sim(x,x')$$

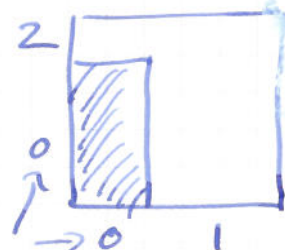Lemma: not all kernels can be embedded as a sim hash:

Proof: $\langle (1,0), (-1,0) \rangle = -1$ ⨏

Extension $\overline{sim}(x,x') = \mathbb{E}_{h,b}\left\{ G(x) G(x') \delta_{h(x),h(x')} \right\}$ ①

Lemma: $\overline{sim}(x,x')$ embeds all kernels $k$ and to app. map $a \cdot b$ to binary

→ Sign is easy → sgn$(a) \to \delta(\cdot)$



this is the only chance to match

← need absolute upper bound per coordinate

Condition 1) trace class $\sum_i \lambda_i < \infty$

2) $L_\infty$ on eigenfcts $\|\varphi_i\|_\infty < \infty$

or rather $\boxed{\sum_i \lambda_i \|\varphi_i\|_\infty^2 < \infty}$

but all kernels satisfy this since we have

$$k(x,x') = \sum_i \lambda_i \varphi_i(x) \varphi_i(x')$$

So, no summand must be $\infty$ □

Properties of the kernel-with-hash

We can use this to approximate kernels, too via

$$k(x,x') \approx \frac{1}{n} \sum_i (h_i(x) = h_i(x'))$$
$$\cdot G_i(x) G_i(x')$$

⇒ we need only 2 bit per dimension ✓

# lecture 6 - Properties of the Min-Wise Hash

## Min-wise hash (recall)

$$\text{sim}(x, x') = \underset{\pi}{\mathbb{E}}\left[\min(\pi(x)) \equiv \min(\pi(x'))\right]$$

$$= \frac{|x \cap x'|}{|x \cup x'|}$$

in practice take $\min_k$ rather than $k$-times min

## Definition:

$$\boxed{\begin{array}{l} \text{min-wise indep. permutation family} \\ Pr\left(\min(\pi(X)) = \pi(x)\right) = \frac{1}{X} \\ \text{for } x \in X \end{array}}$$

## Surprising theorem: (Broder & Mitzenmacher, 2001)

For any mapping $f: \mathbb{P}(x)^{\mathbb{E}} \to \mathcal{R}$ there exists
a permutation $\pi_f$ such that

$$f(x) = f\left(\pi_f^{-1}\left(\min \pi_f(x)\right)\right)$$

and (obviously) there is a distribution over $\pi$
such that

$$\underset{u}{\mathbb{E}}\left[\min \pi(x) = \min \pi(x')\right]$$

$$= \underset{f}{\mathbb{E}}\left[f(x) = f(x')\right]$$

## Advance Preview:

- $f$ is invertible (do slow)
- extract $x_i$ via $f^{-1}(x \setminus \{x_1 \cdots x_{i-1}\})$

## Auxiliary Results

### Theorem to prove:

Assume that there is a mapping $f$, $P(f)$ s.t.

$$Pr_f \{f(A) = f(B)\} = \frac{|A \cap B|}{|A \cup B|}$$

then there exist some $\pi_f$ such that

$$f(x) = f\left(\{\pi_f^{-1}(\min(\pi_f(x)))\}\right)$$

Lemma 1) $Pr_f\left(f(x) = f(\{x\})\right) = \frac{1}{|x|}$

(by definition)

Lemma 2) $Pr_f\left(f(\{x\}) = f(\{x'\})\right) = \delta_{xx'}$

(by definition)

Lemma 3) $Pr_f\left(f(x) \in \{f(\{x_i\}) \mid \forall i\}\right) = 1$

$$\approx \sum_i Pr_f\left(f(x) = f(\{x_i\})\right) = 1$$

(by Lemma 1)

Lemma 4) For $x \subseteq y$ we have

$$Pr_f\left(f(x) = f(y)\right) = \frac{|x|}{|y|}$$

And if $f(y) \in \{f(x_i) \mid i\}$

then $f(y) = f(x)$

Proof: Part 1 is trivial

Corollary: averaging reduces variance

$x_i$ with $\mathbb{E}\{x_i\} = \rho$

$$Var[x_i] = \mathbb{E}[x^2] - (\mathbb{E}[x])^2 = \sigma^2$$

define $X := \frac{1}{n} \sum_{i=1}^{u} x_i$

$\Rightarrow$ $\mathbb{E}[X] = \rho$ ; $Var[X] = \frac{1}{n}\sigma^2$

(since variances add up indep.)

Corollary: $\Pr\{|X - \rho| > \gamma\sigma\} \le \frac{1}{8\sqrt{n}}$

$\rightarrow$ Terrible scaling behavior in $\sigma$, good in $\gamma$, $u$

$\rightarrow$ Want logarithmic in $\sigma$, but cannot use Bernstein / Chernoff since the moments are not bounded ...   Similar to $\min\{x_i\}$ trick
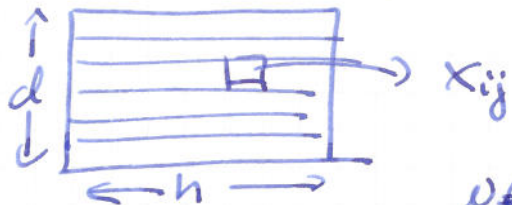
Chernoff bound: for $x_i \in \{a_i, b_i\}$ we have

$$\Pr\left\{\sum_i x_i > \sum_i \mathbb{E}[x_i] + \varepsilon\right\} \le \exp\left(-2\frac{\varepsilon^2}{c^2}\right)$$
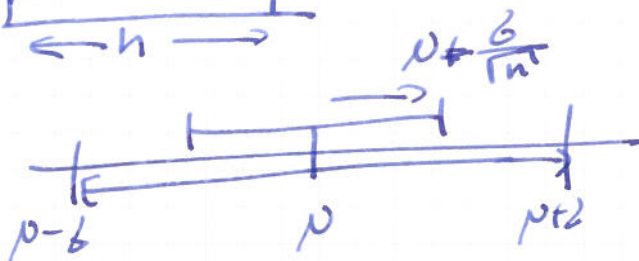
when $c^2 = \sum_i (b_i - a_i)^2$

key idea: Controlling variance & sample size separately



$\hat{x}_i := \frac{1}{n} \sum_{j=1}^{u} x_{ij}$

$\hat{x} := \text{median}_i \hat{x}_i$

$\Pr\left\{|\hat{x}_i - \rho| > \gamma \cdot \frac{\sigma}{\sqrt{n}}\right\} < \frac{1}{8\sqrt{n}}$

$\Pr\left\{|\text{med}_i\{\hat{x}_i\} - \rho| > \gamma\cdot\sigma\right\}$

$\le \Pr\left\{\sum_{i=1}^{d}\{|\hat{x}_i - \rho| > \gamma\cdot\sigma\} > \frac{d}{2}\right\}$

Lots of things must go bad before this happens

this is a R.V. with $\mathbb{E}[\ ] \le \frac{1}{8\sqrt{n}}$

$\Rightarrow$ $\varepsilon = d\left(\frac{1}{2} - \frac{1}{8\sqrt{n}}\right)$

$c^2 = d$

$\Rightarrow$ $\Pr\{\text{failure}\} \le \exp\left(-2d\left(\frac{1}{2} - \frac{1}{8\sqrt{n}}\right)^2\right)$

$\Rightarrow$ Set $\frac{1}{8\sqrt{n}} = \frac{1}{4}$ $\Rightarrow$ $\gamma = \frac{1}{4\sqrt{n}}$

$\Rightarrow$ prob is $\exp\left(-\frac{1}{2}d\right)$ ✓

McDiarmid ineq. ($\sim$ 95)

$$\left| f(x \setminus \{x_i\} \cup \{x_i'\}) - f(x) \right| \le c_i$$

$\Rightarrow$ $\Pr\{\mathbb{E}[f(x)] - f(x) > \varepsilon\} \le \exp\left(-2\frac{\varepsilon^2}{c^2}\right)$

when $c^2 = \sum_i c_i^2$

Self-bounding ineq. of McDiarmid & Reed ($\sim$ 05)

$$\sum_i (g(x) - g_i(x)) \le a \cdot g(x) + b$$

where $g_i(x) = \inf_{x^i} g(x \setminus \{x_i\} \cup \{x^i\})$

then $\Pr\{g(x) - \mathbb{E}_x\{g(x)\} \ge \varepsilon\} \le e^{-\frac{\varepsilon^2}{2(a\rho + b + a\varepsilon)}}$

$\le$ _ u _ $\frac{a\varepsilon}{s}$

$$\Pr\{f(x) = f(y)\}$$

$$= \sum_i \underbrace{\Pr\{f(x) = f(x_i) \mid f(y) = f(x_i)\}}_{= 1} \cdot \underbrace{p(f(y)\} = f(\{x_i\})}_{\frac{1}{|Y|}}$$

---

## Proof of main theorem:

(Lemma 3)

$f$ is invertible on its image. Now define

$$x_1 = f^{-1} \circ f(x) \qquad\qquad i := \pi_f(x_i)$$

$$x_i := f^{-1} \circ f(X \setminus \{x_1 \dots x_{i-1}\})$$

now we want to apply this to arbitrary sets $Y$

define $\quad i := \arg\min_i \{x_i \mid x_i \in Y\}$

by construct $\quad Y \in X \setminus \{x_1 \dots x_{i-1}\}$

and $\qquad\qquad x_i = f^{-1} \circ f(X \setminus \{x_1 \dots x_{i-1}\})$

hence by lemma 4 $\quad f(Y) = f(\{x_i\})$ ▢

(i.e. $\pi_f$ satisfies the conditions)

---

Useful result: (will not prove this)

$$\Pr\{\text{rank}_r(\pi(x)) = \pi(x)\} = \frac{1}{|X|}$$

for min-wise independent families.

---

## Tail bounds:

Recall: Gauss - Markov inequality

$$\Pr(x > \rho c) \le \frac{\rho \mathbb{1}}{c} \quad \text{for } x \ge 0$$

Useful corollary

$$\Pr\{(\min_{i \in 1..k} x_i) > \rho c\} \le \left(\frac{\rho \mathbb{1}}{c}\right)^k$$

(will use this for Count Min Sketch)

Quantile trick

$$\text{define } F(\bar{x}) = \int_{-\infty}^{\bar{x}} dp(x)$$

hma for $x \sim p(x)$ we have $F(x) \sim U[0,1]$

$$\boxed{\text{define } F_k(\bar{x}) := \Pr\{\min_{i \in 1..k} x_i \le \bar{x}\} = F^k(\bar{x})}$$

Corollary 59/95 trick

define $\bar{x}_{0.95} = F^{-1}(0.95)$

then $\Pr\{\max_{i \in 1..59} x_i < \bar{x}_{0.95}\} < 0.05$

Proof: $(0.95)^{59} < 0.05$ ▢

---

Chebychev inequality:

Define $z := (x - \mathbb{E}[x])^2$; $\mathbb{E}[z] = \sigma^2$

$$\Rightarrow \Pr\{z > \gamma^2 \sigma^2\} \le \frac{1}{\gamma^2}$$

$$\Pr\{|x - \mathbb{E}[x]| > \gamma \sigma\}$$