

CS5234 - Algorithms at Scale

Liew Zhao Wei

Semester 1, 2023-2024

1 Probability and Bounds

Lemma 1.1 (Union Bound)

For a countable set of events A_1, A_2, \dots , we have

$$\Pr \left[\bigcup_{i=1}^{\infty} A_i \right] \leq \sum_{i=1}^{\infty} \Pr[A_i] \quad (1)$$

Lemma 1.2 (Linearity of Expectation)

For any random variables X_1, X_2, \dots, X_n , we have

$$\mathbb{E} \left[\sum_{i=1}^n X_i \right] = \sum_{i=1}^n \mathbb{E}[X_i] \quad (2)$$

Lemma 1.3 (Markov's Inequality)

For any *non-negative* random variable X and $t > 0$, we have

$$\Pr[X \geq t] \leq \frac{\mathbb{E}[X]}{t} \quad (3)$$

Lemma 1.4 (Chebyshev's Inequality)

For any random variable X with mean μ and variance σ^2 , we have

$$\Pr[|X - \mu| \geq t] \leq \frac{\sigma^2}{t^2} \quad (4)$$

In fact, this holds for any moment p instead of 2.

Lemma 1.5 (Chernoff-Hoeffding Bounds)

Suppose X_1, \dots, X_n are *independent* random variables with $X_i \in [0, 1]$. Let $X = \sum_{i=1}^n X_i$ and $\mu = \mathbb{E}[X]$ such that $\mu_L \leq \mu \leq \mu_H$. Then, for any $0 \leq \delta \leq 1$,

$$\Pr[X \geq (1 + \delta)\mu] \leq \exp\left\{-\frac{\delta^2\mu}{3}\right\} \quad (5)$$

$$\Pr[X \leq (1 - \delta)\mu] \leq \exp\left\{-\frac{\delta^2\mu}{2}\right\} \quad (6)$$

$$\Pr[|X - \mu| \geq \delta\mu] \leq 2 \exp\left\{-\frac{\delta^2\mu}{3}\right\} \quad (7)$$

$$(8)$$

For any $\delta \geq 0$,

$$\Pr[X \geq \mu + \delta] \leq \exp\left\{-\frac{2\delta^2}{n}\right\} \quad (9)$$

$$\Pr[X \leq \mu - \delta] \leq \exp\left\{-\frac{2\delta^2}{n}\right\} \quad (10)$$

$$\Pr[|X - \mu| \geq \delta] \leq 2 \exp\left\{-\frac{2\delta^2}{n}\right\} \quad (11)$$

More generally, if $a_i \leq X_i \leq b_i$, then

$$\Pr[X \geq \mu + \delta] \leq \exp\left\{-\frac{2\delta^2}{\sum_{i=1}^n (b_i - a_i)^2}\right\} \quad (12)$$

$$\Pr[X \leq \mu - \delta] \leq \exp\left\{-\frac{2\delta^2}{\sum_{i=1}^n (b_i - a_i)^2}\right\} \quad (13)$$

WIP.

3. k-universal hash family (include space analysis)

2 Simple Techniques

WIP. 1. Reservoir sampling 2. Mean trick to drive down variance 3. Median trick to boost success Probability 4. Median of mean trick (2 and 3) to bound concentration

3 Sketches

1. combining sketches, linear sketches 2. Misra-Gries 3. Count-Min-Sketch and Count-Sketch

4 Dimensions and Distances

Lemma 4.1 (Johnson-Lindenstrauss Lemma)

For any set $S \subseteq \mathbb{R}^d$ of n -points, there is an embedding $f: \mathbb{R}^d \rightarrow \mathbb{R}^m$ for $m = O(\epsilon^{-2} \log n)$ such that

$$\forall u, v \in S \quad (1 - \epsilon)\|u - v\|_2^2 \leq \|f(u) - f(v)\|_2^2 \leq (1 + \epsilon)\|u - v\|_2^2 \quad (14)$$

In other words, we can embed S into a lower-dimensional space while approximately preserving ℓ_2 norms.

Some observations:

- The embedding has only a logarithmic dependence on n and *no* dependence on d .

- The embedding is can be generated using a Gaussian distribution.
- The embedding can be represented as a linear transformation, or in other words, a matrix.

Definition 4.2 (Locality Sensitive Hash)

A hash family $\mathcal{H} = \{h: \mathcal{U} \rightarrow S\}$ is a (r_1, r_2, p_1, p_2) -locally sensitive if for all points $p, p' \in \mathcal{U}$,

1. if $d(p, p') \leq r_1$, then $\Pr_{h \in \mathcal{H}}[h(p) = h(p')] \geq p_1$,
2. if $d(p, p') > r_2$, then $\Pr_{h \in \mathcal{H}}[h(p) = h(p')] \leq p_2$.

In other words, a *locality sensitive hash* (LSH) is a hash family where similar items are more likely to collide. Note that the definition makes sense only if $r_1 < r_2$ and $p_1 > p_2$.

WIP. 1. ANN, PLEB, how to solve them