

CS5234 Cheat Sheet Gabriel Yeo

Inequalities

Markov:	$P(X \geq \alpha) \leq \frac{E(X)}{\alpha}$
Chebyshev:	$P(X - \mu \geq k) \leq \frac{Var(X)}{k^2}$
Chernoff:	$P(X \leq (1 - \delta)\mu) \leq e^{-\frac{\delta^2 \mu}{2}} \leq e^{-\frac{\delta^2 \mu}{3}}$ $P(X \geq (1 + \delta)\mu) \leq e^{-\frac{\delta^2 \mu}{2 + \delta}} \leq e^{-\frac{\delta^2 \mu}{3}}$ $X = X_1 + \dots + X_n$, $X_i \in \{0, 1\}$, (ind. Bernoulli)
Chernoff:	$P(X - \mu \geq \delta\mu) \leq 2e^{-\frac{\delta^2 \mu}{3}}$
Hoeffding:	$P(X - \mu \geq t) \leq 2e^{-\frac{2t^2}{n}}$ $X = X_1 + \dots + X_n$, $X_i \in [0, 1]$
Hoeffding:	$P(\bar{X} - E[\bar{X}] \geq t) \leq 2e^{-2nt^2}$ $\bar{X} = \frac{1}{n}(X_1 + \dots + X_n)$
Hoeffding (gen):	$P(\bar{X} - E[\bar{X}] \geq t) \leq 2e^{-\frac{2t^2}{\sum (a_i - b_i)^2}}$ $X_i \in [a_i, b_i]$

Facts

Taylor expansion:	$e^x = 1 + x + \frac{x^2}{2!} + \dots = \sum_{n=0}^{\infty} \frac{x^n}{n!}$
For $0 < x \leq 1$:	$1 - x \leq e^{-x} \leq 1 - \frac{x}{2}$
For $0 < x \leq 1$:	$\frac{1}{e^2} \leq (1 - x)^{1/x} \leq \frac{1}{e}$
Approx $\binom{a}{b}$:	$\left(\frac{a}{b}\right)^b \leq \binom{a}{b} \leq \left(\frac{ea}{b}\right)^b$
For $0 < x < 1$:	$\sum_{i=0}^{\infty} a^i = \frac{1}{1-a}$
Natural log:	$\ln(n-1) \leq \sum_{i=1}^n \frac{1}{i} \leq \ln(n) + 1$
Powers of 2:	$\sum_{i=0}^n 2^i = 2^{n+1} - 1$
$1/e$:	$= 0.36787944117$
$1/e^2$:	$= 0.13533528323 \leq 1/6$
$1/e^3$:	$= 0.04978706836 \leq 1/20$

Probability

Expectation:	$E[X] = \sum_{v \in D} v \cdot Pr[X = v]$
Variance:	$Var[X] = E[(X - E[X])^2]$ $= E[X^2] - E[X]^2$
More variance:	$Var[aX] = a^2 Var[X]$
Linearity of expectation:	$E[\sum X_i] = \sum E[X_i]$
If independent X_i :	$Var[\sum X_i] = \sum Var[X_i]$
Conditional:	$P[X \text{ and } Y] = P[X Y] \cdot P[Y]$
Independence:	$P[X \text{ and } Y] = P[X] \cdot P[Y]$
Union bound:	$P[\cup E_i] \leq \sum P[E_i]$

Complexities

Algorithm	Classical	Approx
BFS	$O(n + m)$	
Connected?	$O(n + m)$	$O(\frac{1}{\epsilon^2 d})$
# CCs	$O(n + m)$	$O(\frac{d}{\epsilon^3})$
MST weight	$O(m \log n)$	$O(\frac{dW^4 \log W}{\epsilon^3})$
Prim's	$O((m + n) \cdot \log n)$	
Prim's (Fib heap)	$O(m + n \log n)$	
Kruskal's (MST)	$O(m \log n)$	
Dijkstra	$O((m + n) \cdot \log n)$	
Dijkstra (Fib heap)	$O(m + n \log n)$	

Lecture 1 All 0s

Give an array A with n elements, $A[i] \in \{0, 1\}$:
 (1) If array has $\geq \epsilon n$ 1's, return False with probability at least $1 - \delta$:
 Assume $\geq \epsilon n$ 1's, then for sample i :
 $Pr[A[i] = 1] \geq \epsilon n / n \geq \epsilon$.

$$Pr(\text{all samples are 0}) \leq (1 - \epsilon)^s$$

$$\leq (1 - \epsilon)^{\frac{\ln(1/\delta)}{\epsilon}}$$

$$\leq e^{-\ln(1/\delta)}$$

$$\leq \delta$$

$$\text{Fix } s = \frac{\ln(1/\delta)}{\epsilon}$$

Lecture 1 Number of 1s

Give an array A with n elements, $A[i] \in \{0, 1\}$:
 (1) Find number of 1's $\pm \epsilon$ with probability at least $1 - \delta$.
 Let Y_i be s independent samples in $\{0, 1\}$.
 Output $= Z = 1/s \sum Y_i$
 Probability of failure:

$$Pr(|Z - E[Z]| \geq \epsilon) \leq 2e^{-2s\epsilon^2}$$

$$\leq 2e^{-2\frac{\ln(2/\delta)}{2\epsilon^2}\epsilon^2}$$

$$\leq 2e^{-\ln(2/\delta)}$$

$$\leq 2e^{\ln(\delta/2)}$$

$$\leq 2 \cdot \delta/2$$

$$\leq \delta$$

$$\text{Fix } s = \frac{\ln(2/\delta)}{2\epsilon^2}$$

Lecture 2 Connectivity

G is ϵ -close to connected if you can modify at most ϵnd entries in the adjacency list to make it connected. If $G(V, E)$ is connected, output True, else if ϵ -far from connected, output False.

```

Connected(G, n, d,  $\epsilon$ )
  Repeat  $16/\epsilon d$  times:
    - Choose a random node  $u$ 
    - Do a BFS from  $u$ , stopping after  $8/\epsilon d$  nodes seen.
    - If CC of  $u$  has  $\leq 8/\epsilon d$  nodes, return FALSE.
  return TRUE

```

BFS cost $= \frac{8}{\epsilon d} \cdot d$
 Total complexity $= O(\frac{1}{\epsilon^2 d})$

Lecture 2 Connected Components

Output CC such that: $|CC(G) - C| \leq \epsilon n$, w.p. $> 1 - \delta$
 Define $n(u)$ = num nodes in the CC of u .
 $cost(u) = 1/n(u)$.
 $\sum_{u \in CC_i} cost(u) = 1$.

```

sum = 0
for  $j = 1$  to  $s$ :
  - Sample  $u$  randomly
  - BFS from  $u$ , stop when see up to  $2/\epsilon$  nodes
  - If BFS found  $> 2/\epsilon$  node:
    - sum := sum +  $\epsilon/2$ 
  - Else:
    - sum := sum + cost( $u$ )
return  $n \cdot (\text{sum}/s)$ 

```

Let $\bar{C} = \sum cost(u_j)$
 Let $Y_j = cost(u_j)$ of our sample j
 $|CC(G) - \bar{C}| \leq \epsilon n/2$
 $E[Y_j] = \sum \frac{1}{n} cost(u_i) = \frac{1}{n} \bar{C}$
 $E[\sum Y_j] = s \cdot E[Y_j] = \frac{s}{n} \bar{C}$
 Since we output $\frac{n}{s} \sum Y_j$, we get $E[\frac{n}{s} \sum Y_j] = \bar{C}$

$$P(|\bar{C} - \frac{n}{s} \sum Y_j| > \epsilon n/2)$$

$$= P(|E[\sum Y_j] - \sum Y_j| > \frac{s}{n} \epsilon n/2)$$

$$= P(|E[\sum Y_j] - \sum Y_j| > s\epsilon/2)$$

$$\leq 2e^{-2\epsilon^2 s^2 / 4s}$$

$$\leq 2e^{-\epsilon^2 s/2} \leq \delta$$

Set $s = \frac{2}{\epsilon^2} \ln(2/\delta)$.
 Complexity $= 2d/\epsilon \cdot O(\frac{1}{\epsilon^2} \ln(1/\delta)) = O(\frac{d}{\epsilon^3} \ln(1/\delta))$.

Lecture 3 MST weight

Output M such that: $M = MST(G)(1 \pm \epsilon)$ w.p. $> 1 - \delta$.
 Let G_j be the graph with edge weights j and below.
 Let C_j be the connected components in G_j .
 $MST(G)$ contains $C_j - 1$ edges of weight $> j$.
 Therefore:

$$MST(G) = (n - C_1) + \sum_{j=1}^{W-1} (j+1)(C_j - C_{j+1})$$

$$= n - W + \sum_{j=1}^{W-1} C_j$$

```
sum := n - W
for j = 1 to W - 1:
    -  $X_j = \text{ApproxCC}(G_j, d, \epsilon' \delta')$ 
    - sum := sum +  $X_j$ 
return sum
```

Sum of errors: $(\epsilon' n)(W - 1)$
 Set $\epsilon' = \epsilon/W$, then sum of errors $\leq \epsilon n$.
 Set $\delta' = \delta/W$. Then $P(\text{any fail}) \leq \sum_{i=1}^{W-1} \delta/W \leq \delta$.
 $MST(G) - \epsilon n \leq \text{sum} \leq MST(G) + \epsilon n$
 Since $MST(G) \geq n - 1 \geq n/2$, $n \leq 2MST(G)$,
 $MST(G)(1 - 2\epsilon) \leq \text{sum} \leq MST(G)(1 + 2\epsilon)$

Lecture 3 Maximal Matching

Return the size of the Maximal Matching.

```
query(e):
    for all neighbours  $e'$  of  $e$ :
        if hash( $e'$ ) < hash( $e$ )
            if query( $e'$ ) = TRUE
                return FALSE
    return TRUE

sum := 0
for j = 1 to s:
    - Choose edge  $e$  uniformly at random.
    - if (query( $e$ ) = True)
        sum := sum + 1
return m.sum/s
```

$$E[\text{cost}] = 2 \sum_{k=1}^{\infty} \frac{d^k}{k!} = O(e^d).$$

$$\text{sum} = MM(G) \pm \epsilon m.$$

$$\text{Complexity} = O\left(\frac{e^d}{\epsilon^2}\right).$$

Can do better: $O(\frac{d^4}{\epsilon^2})$, even $O(\frac{d^2}{\epsilon^2})$, and reduce error to $\pm \epsilon n$.

Lecture 3 Yao's Mini-Max

Every randomized algorithm on a worst-case input is always slower than the best deterministic algorithm on the worst distribution.

$$\forall A \in R : \max_{x \in X} (E[\text{cost}(A, x)]) \geq \min_{B \in D} (E[\text{cost}(B, x \text{ chosen from } \gamma)])$$

Recipe:

- Choose distribution γ .
- Show that the expected cost of every deterministic algorithm from γ is greater than some cost c .
- Conclude that every randomized algorithm has at least one input with expected cost at least as bad as c .

Lecture 4 Misra Gries

$\text{count}(x): N(x) - \epsilon m \leq \text{count}(x) \leq N(x) + \epsilon m$.
 Heavy Hitters: returns

- every item that appears $\geq 2\epsilon m$ times.
- no item that appears $< \epsilon m$ times.

```
Set  $P$  of  $\langle \text{item}, \text{count} \rangle$  pairs
For each  $u$  in stream  $S$ :
    1. if  $\langle u, c \rangle$  is in set  $P$ , increment  $c$ .
    2. else add  $\langle u, 1 \rangle$  to set  $P$ .
    3. if  $|P| > k$ , decrement count  $c$  for each item.
    4. Remove all items from  $P$  with count  $c = 0$ .

Count( $x$ ):
    1. if  $\langle x, c \rangle$  is in  $P$ , return  $c$ .
    2. else return 0.
```

Choose $k = 1/\epsilon$. Then $N(x) - \epsilon m \leq \text{count}(x) \leq N(x)$.
 Space required: $O(k \log m)$.

Proof:

- Count of x is incremented $N(x)$ times in total.
- Total increments is m .
- When $\text{count}(x)$ is decremented, at least k items are also decremented.
- At most m decrements in total.
- So $\text{count}(x)$ is decremented at most m/k times.

For Heavy Hitters: return x if $\text{count}(x) \geq \epsilon m$.

Lecture 5 Spanner

If graph H has $\text{girth}(H) > 2k \rightarrow H$ has $O(n^{1+1/k})$ edges.
 $\log(n)$ -spanner space with $k = \log(n)$: requires $O(n \log n)$ space.

Lecture 6 k-Median Clustering

Cost of each node i : $C_i = \sum_j x_{i,j} d(p_i, p_j)$.

The LP minimizes $\min \sum_i C_i$.

Goal: round fractional LP such that $C'_i \leq 4C_i$.

- If some p_i is within $4C_j$ of p_j , remove p_i from centers.
- In other words: if there is q s.t. $d(p_i, q) \leq 2C_i$ and $d(p_j, q) \leq 2C_j$, delete p_i .
- $C_j \leq C_i$ because of the order of node processing.

Goal: less than $2k$ centers.

- $\sum_{i: d(p_i, p_j) \leq 2C_j} y_i \geq 1/2$
- Since y 's sum to k , if $V(i)$ are disjoint, cannot add more than $2k$ points to S .

Lecture 7 External Memory Model

Cache size = M . Block size = B . Number of lines (cache slots) = M/B .

Assumptions:

- One cache level.
- Only memory access has cost.
- Ideal cache and replacement.

Problem	Classical	EMM
Scan Linked List	$O(N)$	$O(N)$
Scan Array	$O(N)$	$O(N/B)$
Search Linked List	$O(N)$	$O(N)$
Search Red-black tree	$O(\log N)$	$O(\log N)$
Search Array	$O(\log N)$	$O(\log(N/B))$
Search B-tree	$O(\log_B N)$	$O(\log_B N)$
Sort B-tree	$O(N \log N)$	$O(N \log_B N)$
Read/Write B-tree	$O(\log_B N)$	$O(\log_B N)$

Lecture 7 Caching

(a, b) - tree with n keys has height $\leq \log_a(\frac{n}{a}) + 1$.
 At most n/a leaves. Every node except the root has degree $> a$.
 Node at height $\log_a(\frac{n}{a})$ has $\geq a^{\log_a(\frac{n}{a})} \geq \frac{n}{a}$ leaves.
 Corollary: if $a \geq B$, then (a, b) - tree with n keys has height $O(\log_B n)$.
 Amortized cost of split/share/merge is $O(1/B)$, thus $O((1/B) \log_B N)$ per operation.
 With parent points: insert may cost $O(B \log_B N)$ if every level needs to split.