

New York University Tandon School of Engineering  
Computer Science and Engineering  
CS-GY 6763: Midterm Practice.

## Logistics

- Exam will be held in class on **Friday, 3/14** starting at 2:00pm sharp. Please arrive on time!
- You will have **1 hour, 15 minutes** to answer a variety of short answer and longer form questions.
- You can bring a one page sheet of paper (two-sided if you want) with notes, theorems, etc. written down for reference.
- Scrap paper will be provided and can be used to write solutions if extra space is needed.
- I will be in the room to answer any questions.

## Concepts to Know

### Random variables and concentration.

- Linearity of expectation and variance.
- Indicator random variables and how to use them.
- Markov's inequality, Chebyshev's inequality (ideally should know from memory so you can apply quickly).
- Union bound (should know from memory).
- Chernoff and Bernstein bounds (don't need to memorize the exact bounds, but can apply if given).
- General idea of law of large numbers and central limit theorem.
- The probability that a normal random variables  $\mathcal{N}(0, \sigma^2)$  falls further than  $k\sigma$  away from its expectation is  $\leq O(e^{-k^2/2})$ .

### Hashing, Dimensionality Reduction, High Dimensional Vectors

- Random hash functions.
- Random hashing for frequency estimation.
- Random hashing for distinct elements estimation.
- MinHash for Jaccard similarity estimation.
- Locality sensitive hash functions.
- MinHash and SimHash for Jaccard Similarity and Cosine Similarity.
- Adjusting false positive rate and false negative rate in an LSH scheme.
- Statement of Johnson-Lindenstrauss lemma (know from memory).
- Statement of *distributional* JL lemma and how it can be used to prove JL.

## High dimensional geometry

- How to draw a random unit vector from the sphere in  $d$  dimensions (draw  $\mathbf{x}$  with all entries i.i.d.  $\mathcal{N}(0, 1)$  and normalize it).
- How does  $\|x - y\|_2^2$  relate to  $\langle x, y \rangle$  if  $x$  and  $y$  are unit vectors?
- How many mutually orthogonal unit vectors are there in  $d$  dimensions?
- There are  $2^{\theta(\epsilon^2 d)}$  nearly orthogonal unit vectors in  $d$  dimensions (with  $\langle x, y \rangle \leq \epsilon$ ). Know roughly how prove this fact using the *probabilistic method*, which required a an exponential *concentration inequality* + *union bound*.
- Know how to prove that all but an  $2^{\theta(-\epsilon d)}$  fraction of a balls volume in  $d$  dimensions lies in a spherical shell of width  $\epsilon$  near its surface.
- The surface area/volume ratio *increases* in high dimensions.
- The cube volume/ball volume ratio *increases* in high dimensions.

## Practice Problems

1. Show that for any random variable  $X$ ,  $\mathbb{E}[X^2] \geq \mathbb{E}[X]^2$ .
2. Show that for independent  $X$  and  $Y$  with  $\mathbb{E}[X] = \mathbb{E}[Y] = 0$ ,  $\text{var}[X \cdot Y] = \text{var}[X] \cdot \text{var}[Y]$ .
3. Given a random variable  $X$ , can we conclude that  $\mathbb{E}[1/X] = 1/\mathbb{E}[X]$ ? If so, prove this. If not, give an example where the equality does not hold.
4. Indicate whether each of the following statements is **always** true, **sometimes** true, or **never** true. Provide a short justification for your choice.
  - (a)  $\Pr[X = s \text{ and } Y = t] > \Pr[X = s]$ .    ALWAYS    SOMETIMES    NEVER
  - (b)  $\Pr[X = s \text{ or } Y = t] \leq \Pr[X = s] + \Pr[Y = t]$ .    ALWAYS    SOMETIMES    NEVER
  - (c)  $\Pr[X = s \text{ and } Y = t] = \Pr[X = s] \cdot \Pr[Y = t]$ .    ALWAYS    SOMETIMES    NEVER

5. Let  $\Pi$  be a random Johnson-Lindenstrauss matrix (e.g. scaled random Gaussians) with  $O(\log(1/\delta)/\epsilon^2)$  rows. With probability  $(1 - \delta)$ ,

$$\min_{\mathbf{x}} \|\Pi \mathbf{A} \mathbf{x} - \Pi \mathbf{b}\|_2^2 \leq (1 + \epsilon) \min_{\mathbf{x}} \|\mathbf{A} \mathbf{x} - \mathbf{b}\|_2^2.$$

Is the statement true:    ALWAYS    SOMETIMES    NEVER

6. Assume there are 1000 registered users on your site  $u_1, \dots, u_{1000}$ , and in a given day, each user visits the site with some probability  $p_i$ . The event that any user visits the site is independent of what the other users do. Assume that  $\sum_{i=1}^{1000} p_i = 500$ .
  - (a) Let  $X$  be the number of users that visit the site on the given day. What is  $E[X]$ ?
  - (b) Apply a Chernoff bound to show that  $\Pr[X \geq 600] \leq .01$ .
  - (c) Apply Markov's inequality and Chebyshev's inequality to bound the same probability. How do they compare?
7. Give an example of a random variable and a deviation  $t$  where Markov's inequality gives a tighter upper bound than Chebyshev's inequality.
8. Suppose there is some unknown vector  $\mu \in \mathbb{R}^d$ . We receive noise perturbed random samples of the form  $\mathbf{Y}_1 = \mu + \mathbf{X}_1, \dots, \mathbf{Y}_k = \mu + \mathbf{X}_k$  where each  $\mathbf{X}_i$  is a random vector with each of its entries distributed as an independent random normal  $\mathcal{N}(0, 1)$ . From our samples  $\mathbf{Y}_1, \dots, \mathbf{Y}_k$  we hope to estimate  $\mu$  by  $\tilde{\mu} = \frac{1}{k} \sum_{i=1}^k \mathbf{Y}_i$ .

- (a) Prove that with  $k = O(\log d/\epsilon^2)$  samples,  $\max_{i=1,\dots,d} |\mu_i - \tilde{\mu}_i| \leq \epsilon$  with probability 9/10.
- (b) Prove that with  $k = O(d \log d/\epsilon)$  samples,  $\|\mu - \tilde{\mu}\|_2 \leq \epsilon$  with probability 9/10.
9. For two length  $d$  binary vectors  $\mathbf{q}, \mathbf{y} \in \{0, 1\}^d$ , consider the hamming similarity:

$$s(\mathbf{q}, \mathbf{y}) = 1 - \frac{\|\mathbf{q} - \mathbf{y}\|_0}{d}.$$

Recall that  $\|\mathbf{q} - \mathbf{y}\|_0 = \sum_{i=1}^d \mathbb{1}[\mathbf{q}_i \neq \mathbf{y}_i]$ . Construct a function  $h$  as follows: define the random function  $c : \{0, 1\}^d \rightarrow \{0, 1\}$  as  $c(\mathbf{x}) = \mathbf{x}[j]$ , where  $j$  is a uniform random integer in  $\{1, \dots, d\}$ . Then, let  $g$  be a uniform random hash function from  $\{0, 1\} \rightarrow \{1, \dots, m\}$ . Finally, let:

$$h(\mathbf{x}) = g(c(\mathbf{x})).$$

Prove that  $h$  is a locality sensitive hash function for hamming similarity.