

Homework 2 Solution key

Name: 2 Solution key

Problem 1

1. Once there are $n+1$ servers in this setup, the expected number of items on the $(n+1)^{\text{st}}$ server is $\frac{m}{n+1}$, by symmetry. All of these items (and only these items) must have been relocated when the $(n+1)^{\text{st}}$ server was added. So the expected number of items that move is $\frac{m}{n+1}$.

2. For a server S to own more than a $c \log n/n$ fraction of the interval, it would need to be that *no other server* falls within distance $c \log n/n$ to the left of the server. We can choose the random location of server S first. Then the probability of any one server landing within distance $c \log n/n$ from S 's left is $c \log n/n$. So the probability *no servers* land that close is:

$$(1 - c \log n/n)^{n-1} \leq \frac{1}{10n},$$

as long as we choose c to be a large enough constant (same analysis as homework 1). By a *union bound*, we thus have that no server owns more than an $O(\log n/n)$ fraction of the interval with probability $\geq 1 - n \frac{1}{10n} = \frac{9}{10}$ which proves the claim.

3. From Part 2, we could have equivalently proven that no server owns more than a $c \log n/n$ fraction of the interval with probability $19/20$ (by choosing c larger). For the rest of the problem, assume that this event is happening.

For servers S_1, \dots, S_n let $Y_i^{(j)}$ be the indicator random variable that item j lands within distance $c \log n/n$ to S_i 's left. Let X_i equal $X_i = \sum_{j=1}^m Y_i^{(j)}$. Since we assumed that no server owns more than a $c \log n/n$ fraction of the interval, X_i is an *upper bound* on the number of items assigned to server i . So it suffices to show that X_i is not too large for all i .

To do so, note that, for a fixed i , $Y_i^{(1)}, Y_i^{(2)}, \dots, Y_i^{(m)}$ are an independent $\{0, 1\}$ random variables, where each is 1 with probability exactly $c \log n/n$. So they are just biased coin flips!

Let $c > 2$ be a sufficiently large constant. Using the Chernoff bound from class with $\epsilon = c$, we get that:

$$\mathbb{P}[X_i \geq 2c \cdot \frac{m \log n}{n}] \leq e^{\frac{-c^2 m \log n/n}{2+c}} \leq e^{\frac{-c \log n}{2}} \leq \frac{1}{20n},$$

for large enough c . The last inequality uses that $m > n$ (as specified in the problem).

We conclude via a union bound that no server is assigned more than $O(m \log n/n)$ items with probability $\frac{19}{20}$.

There's one last step – we needed two events to hold for our proof to go through: 1) no server owns more than a $c \log n/n$ fraction of the interval and 2) no server was assigned too many items. Since each holds with probability $19/20$, by another union bound, both hold with probability $9/10$.

Problem 2 (a)

1. **Expectation Calculation.** As in class, we have that $\mathbb{E}[\|\Pi x\|_2^2] = \mathbb{E}[\langle \pi, x \rangle^2]$, where π is a single unscaled row from the matrix Π . I.e. π has length n and contains random ± 1 entries. We have:

$$\begin{aligned} \mathbb{E}[\langle \pi, x \rangle^2] &= \mathbb{E} \left[\left(\sum_{j=1}^n \pi_j x_j \right)^2 \right] = \mathbb{E} \left[\sum_{j=1}^n \pi_j^2 x_j^2 \right] + \mathbb{E} \left[\sum_{i \neq j}^n \pi_i \pi_j x_j x_i \right] \\ &= \sum_{j=1}^n \mathbb{E}[\pi_j^2] x_j^2 + \sum_{i \neq j}^n \mathbb{E}[\pi_i \pi_j] x_j x_i. \end{aligned}$$

The last equality follows from linearity of expectation. Since π_i is independent of π_j , we have that for $j \neq i$, $\mathbb{E}[\pi_i \pi_j] = \mathbb{E}[\pi_i] \mathbb{E}[\pi_j] = 0$. On the other hand $\pi_j^2 = 1$ deterministically, so we have $\mathbb{E}[\pi_j^2] = 1$. Plugging in above, we find that

$$\mathbb{E}[\langle \pi, x \rangle^2] = \sum_{j=1}^n x_j^2 + \sum_{i \neq j}^n 0 \cdot x_j x_i = \sum_{j=1}^n x_j^2 = \|x\|_2^2,$$

as desired.

Variance Calculation. Since $\|\Pi x\|_2^2 = \frac{1}{k} \sum_{i=1}^k \langle \pi^i, x \rangle^2$, where π^1, \dots, π^k are the unscaled rows of Π , we first observe that $\text{Var}[\|\Pi x\|_2^2] = \frac{1}{k} \text{Var}[\langle \pi, x \rangle^2]$ for a single random ± 1 vector π . So we just need to bound $\text{Var}[\langle \pi, x \rangle^2]$. This gets a bit tricky! There are many ways to do it, but I think the easiest way is to take advantage of linearity of variance by writing:

$$\langle \pi, x \rangle^2 = \sum_{j=1}^n \pi_j^2 x_j^2 + 2 \sum_{i > j} \pi_i \pi_j x_i x_j.$$

The terms in the first part of the sum are actually deterministic, since $\pi_j^2 = 1$. The terms in the second part of the sum are random, but they are *pairwise independent* since $\pi_i \pi_j$ is random ± 1 and independent from any $\pi_i \pi_k$, $\pi_k \pi_j$, or $\pi_k \pi_\ell$. They are not mutually independent, but we only need pairwise independence to apply linearity of variance. Note that to make this claim it's important that I used the form $2 \sum_{i > j}$ instead of $\sum_{i \neq j}$. If I did the later, there would be repeated random variables in the sum ($\pi_i \pi_j x_i x_j$ and $\pi_j \pi_i x_j x_i$). Writing the other way removes duplicates.

$$\text{Var}[\langle \pi, x \rangle^2] = \sum_{j=1}^n \text{Var}[\pi_j^2 x_j^2] + 4 \sum_{i > j} \text{Var}[\pi_i \pi_j x_i x_j] = 0 + 4 \sum_{i > j} x_j^2 x_i^2.$$

Then finally we observe that:

$$\|x\|_2^4 = \|x\|_2^2 \cdot \|x\|_2^2 = (x_1^2 + \dots + x_n^2) \cdot (x_1^2 + \dots + x_n^2) \geq 2 \sum_{i > j} x_j^2 x_i^2.$$

Putting this together we have that $\text{Var}[\langle \pi, x \rangle^2] \leq 2\|x\|_2^4$ and the result follows since $\text{Var}[\|\Pi x\|_2^2] = \frac{1}{k} \text{Var}[\langle \pi, x \rangle^2]$ as claimed above.

2. This just follows directly from Chebyshev's.

3. It's almost the same analysis as in part 1. The first thing to observe is that:

$$\langle \Pi x, \Pi y \rangle = \frac{1}{k} \sum_{i=1}^k \langle \pi^i, x \rangle \langle \pi^i, y \rangle.$$

So we have that $\mathbb{E}[\langle \Pi x, \Pi y \rangle] = \mathbb{E}[\langle \pi, x \rangle \langle \pi, y \rangle]$ and $\text{Var}[\langle \Pi x, \Pi y \rangle] = \frac{1}{k} \text{Var}[\langle \pi, x \rangle \langle \pi, y \rangle]$, where π is a single random ± 1 vector. We also have that

$$\langle \pi, x \rangle \langle \pi, y \rangle = \left(\sum_{j=1}^n \pi_j x_j \right) \cdot \left(\sum_{j=1}^n \pi_j y_j \right) = \sum_{i=1}^n \pi_i^2 x_i y_i + \sum_{j \neq i} \pi_i \pi_j x_i y_j.$$

From this it's clear that

$$\mathbb{E}[\langle \Pi x, \Pi y \rangle] = \mathbb{E}[\langle \pi, x \rangle \langle \pi, y \rangle] = \sum_{i=1}^n x_i y_i = \langle x, y \rangle,$$

as desired.

The variance calculation is also a bit tricky since we need to make sure our sums involve pairwise independent random variables. We have that:

$$\langle \pi, x \rangle \langle \pi, y \rangle = \sum_{i=1}^n \pi_i^2 x_i y_i + \sum_{j>i} \pi_i \pi_j (x_i y_j + x_j y_i).$$

Applying linearity of variance, we find that

$$\begin{aligned} \text{Var}[\langle \pi, x \rangle \langle \pi, y \rangle] &= \sum_{j>i} (x_i y_j + x_j y_i)^2 = \sum_{j>i} x_i^2 y_j^2 + x_j^2 y_i^2 + 2x_i x_j y_i y_j \\ &\leq 2 \sum_{j>i} x_i^2 y_j^2 + x_j^2 y_i^2 \\ &\leq 2(x_1^2 + \dots + x_n^2)(y_1^2 + \dots + y_n^2) \\ &= 2\|x\|_2^2 \|y\|_2^2. \end{aligned}$$

In second to last inequality we have used that for any a, b , $2ab \leq a^2 + b^2$, which follows from the fact that $(a - b)^2 \geq 0$ for all a, b (this is technically called the AM-GM inequality).

Overall, we get a variance bound of:

$$\text{Var}[\langle \Pi x, \Pi y \rangle] \leq \frac{2}{k} \|x\|_2^2 \|y\|_2^2.$$

Once they get the mean and variance, the bound just follows from applying Chebyshev inequality again. .

Problem 2 (b)

1. Construct 2 length U binary vectors x and y where $x_i = 1$ if $i \in X$ and 0 otherwise, and $y_i = 1$ if $i \in Y$ and 0 otherwise. Note that $|X \cap Y|$ is exactly equal to $\langle x, y \rangle$, so we can estimate the quantity using sketches Πx and Πy . If we set $k = O(1/\epsilon^2)$, then with 9/10 probability we will have:

$$|\langle x, y \rangle - \langle \Pi x, \Pi y \rangle| \leq \epsilon \|x\|_2 \|y\|_2$$

Note that $\|x\|_2^2 = |X|$ and $\|y\|_2^2 = |Y|$, which yields the bound.

Problem 3

1. For any vector x , let z be the point on the hyperplane closest to x . Now:

$$\langle x, a \rangle = \langle x - z, a \rangle + \langle z, a \rangle = \langle x - z, a \rangle + c = \|x - z\|_2 + c \geq c + \epsilon.$$

In the second step we used that $\langle z, a \rangle = c$ since z is on the hyperplane. And in the next step we use that $x - z$ must be perpendicular to the hyperplane (for z to be the closest point). And thus $x - z$ is *parallel* to a . Since a is a unit vector, $\langle x - z, a \rangle = \|x - z\|_2$. The proof for any y on the other side of the hyperplane is the same, but in that case, $y - z$ points directly opposite of a

2. To show that there exists a good separating hyperplane for the dimension reduced data, we exhibit one: consider the hyperplane given by parameters $\Pi a / \|\Pi a\|_2, c / \|\Pi a\|_2$.

We can apply Problem 2 to claim that, if Π reduces to $O(\log(1/\delta)/\epsilon^2)$ dimensions, then with probability $(1 - \delta)$ for *any* $x \in X$ or $\forall y \in Y$,

$$\langle \Pi a, \Pi x \rangle \geq \langle a, x \rangle - \epsilon/2 \geq c + \epsilon/2 \quad \text{and} \quad \langle \Pi a, \Pi y \rangle \leq \langle a, y \rangle + \epsilon/2 \leq c - \epsilon/2.$$

Above we use the fact that $\|x\|_2 \|\vec{a}\|_2 = 1$ and $\|y\|_2 \|\vec{a}\|_2 = 1$ since all x and y are specified to be unit vectors. Equivalently, we have:

$$\langle \Pi a / \|\Pi a\|_2, \Pi x \rangle \geq c / \|\Pi a\|_2 + \epsilon/2 \|\Pi a\|_2 \quad \text{and} \quad \langle \Pi a / \|\Pi a\|_2, \Pi y \rangle \leq c / \|\Pi a\|_2 - \epsilon/2 \|\Pi a\|_2. \quad (1)$$

We also have from the distributional JL lemma that, with probability $1 - \delta$, $\|\Pi a\|_2 \leq 2$. And if we set $\delta = 1/99(n + 1)$, by a union bound we have that (1) holds for all n points in our data set and $\|\Pi a\|_2 \leq 2$ simultaneously with probability $99/100$. This proves the claim with margin $\epsilon/4$.

Problem 4

This problem can be solved in a similar way to the Shazam example from the Lecture 4 notes. You need to optimize over values of s and t , where s is the number of independent locality-sensitive hash functions used in your scheme and t is the number of tables used.

Following the analysis in Lecture 5, given a query vector \vec{y} and some database vector \vec{x} , the probability of \vec{x} showing up as a candidate near-duplicate (which will need to be scanned when \vec{y} is issued as a query) is equal to:

$$(1 - (1 - \theta(\vec{x}, \vec{y})/\pi)^s)^t \quad (2)$$

where θ is the angle between vectors \vec{x} and \vec{y} .

Our goal is to find the s, t pair with the smallest value of t which satisfies:

$$1. \text{ If } \cos(\theta(\vec{x}, \vec{y})) \geq .98, (1 - (1 - \theta(\vec{x}, \vec{y})/\pi)^s)^t \geq .99$$

```

1 % brute force search over t values to find smallest t that works
2 - tvals = 1:100;
3 % boudaries and heights of the buckets given
4 - cos_sims = [-1:.25:1];
5 - freq = [.01,1.99,14,34,34,14,1.90,.01];
6 % convert boundaries from cosine similarities to angles
7 - thetas = acos(cos_sims);
8 % will only use top edge of each bucket for a worst case analysis
9 - thetas = thetas(2:end);
10 % we want to find any match with cosine similarity >= .98 with prob. ,99
11 - cutoff = acos(.98);
12 - cprob = .99;
13
14 % keep track of how many candidate matches there are for a given t
15 - ncandidates = zeros(1,length(tvals));
16 - for t = tvals
17     % smallest value of s which ensures we find near-matches
18     - sopt = floor(log(1 - (1-cprob)^(1/t))/log(1-cutoff/pi));
19     % expected number of hits
20     - hitProbs = 1 - (1 - (1-thetas./pi).^sopt).^t;
21     - ncandidates(t) = (hitProbs*freq'/100)*100000000;
22 - end
23 % answer to part (a)
24 - min(find(ncandidates < 1e6))
25 % answer to part (b)
26 - min(find(ncandidates < 2e5))

```

2. Based on the histogram data provided, the expected number of candidate near-duplicates in less than 1 million, or 200k, for parts (1) and (2).

Observe that $(1 - (1 - \theta(\vec{x}, \vec{y})/\pi)^s)^t$ is monotonically decreasing with s and the expected number of duplicates monotonically decreases with s . So, for a given value of t , it suffices to find the largest possible s such that $(1 - (1 - \cos^{-1}(.98)/\pi)^s)^t \geq .99$. I did this by solving for:

$$s = \frac{\log(1 - (1 - .99)^{1/t})}{\log(1 - \cos^{-1}(.98)/\pi)}$$

and taking the floor.

Then an upper bound on the number of expected candidates can be computed for this s . This will be the smallest possible number of expected candidates for the given t .

My code is included below. I obtained solutions of:

- **20 tables** for ≤ 1 million candidates
- **44 tables** for $\leq 200k$ candidates