

CSE 547 - Assignment 1

Philip Pham

May 16, 2018

Problem 0

List of collaborators: I have not collaborated with anyone.

List of acknowledgements: None.

Certify that you have read the instructions: I have read and understood these policies.

Problem 1: Gaussian Random Projections and Inner Products

Let $\phi(x) = \frac{1}{\sqrt{m}}Ax$ represent our random projection of $x \in \mathbb{R}^d$, with A an $m \times d$ projection matrix with each entry sampled i.i.d from $N(0, 1)$. Note that each row of A is a random projection vector, $v^{(i)}$.

The *norm preservation theorem* states that for all $x \in \mathbb{R}^d$, the norm of the random projection $\phi(x)$ approximately maintains the norm of the original x with high probability:

$$\mathbb{P}\left((1 - \epsilon)\|x\|^2 \leq \|\phi(x)\|^2 \leq (1 + \epsilon)\|x\|^2\right) \geq 1 - 2\exp\left(-\left(\epsilon^2 - \epsilon^3\right)m/4\right), \quad (1)$$

where $\epsilon \in (0, 1/2)$.

Using the norm preservation theorem, prove that for any $u, v \in \mathbb{R}^d$ such that $\|u\| \leq 1$ and $\|v\| \leq 1$,

$$\mathbb{P}(|u \cdot v - \phi(u) \cdot \phi(v)| \geq \epsilon) \leq 4\exp\left(-\left(\epsilon^2 - \epsilon^3\right)m/4\right) \quad (2)$$

Solution

Proof. First note that

$$(1 - \epsilon)\|u + v\|^2 \leq \|\phi(u + v)\|^2 \leq (1 + \epsilon)\|u + v\|^2$$

implies that

$$\|u + v\|^2 - 2\epsilon \leq \|\phi(u + v)\|^2 \leq \|u + v\|^2 + 2\epsilon \quad (3)$$

by triangle inequality and the assumption of the norms of u and v .

Thus, the probability of the event in Equation 3 is less than that of Equation 1.

Using this and taking the additive inverse, we have that

$$\mathbb{P}\left(\|\phi(u + v)\|^2 \notin [\|u + v\|^2 - 2\epsilon, \|u + v\|^2 + 2\epsilon]\right) \leq 2\exp\left(-\left(\epsilon^2 - \epsilon^3\right)m/4\right) \quad (4)$$

$$\mathbb{P}\left(\|\phi(u - v)\|^2 \notin [\|u - v\|^2 - 2\epsilon, \|u - v\|^2 + 2\epsilon]\right) \leq 2\exp\left(-\left(\epsilon^2 - \epsilon^3\right)m/4\right). \quad (5)$$

By the countable sub-additivity property of probability distributions, we have that the probability of both these events occurring is at most $4 \exp(-(\epsilon^2 - \epsilon^3) m/4)$. Thus, we are done if we can show $\{|u \cdot v - \phi(u) \cdot \phi(v)| \geq \epsilon\}$ subsets these two conditions.

If we have the pair

$$\|\phi(u+v)\|^2 \leq \|u+v\|^2 - 2\epsilon \Rightarrow \|u+v\|^2 - \|\phi(u+v)\|^2 \geq 2\epsilon \quad (6)$$

$$\|\phi(u-v)\|^2 \geq \|u-v\|^2 + 2\epsilon \Rightarrow \|u-v\|^2 - \|\phi(u-v)\|^2 \geq 2\epsilon, \quad (7)$$

we can use the linearity of ϕ and the expansion $\|u \pm v\|^2 = \|u\|^2 + \|v\|^2 \pm 2u \cdot v$, we can add the two inequalities to obtain

$$4(u \cdot v - \phi(u) \cdot \phi(v)) \geq 4\epsilon.$$

Thus, we have that the conditions in Equations 4 and 5 imply $u \cdot v - \phi(u) \cdot \phi(v) \geq \epsilon$.

Similarly, we show that the pair

$$\|\phi(u+v)\|^2 \geq \|u+v\|^2 + 2\epsilon \quad (8)$$

$$\|\phi(u-v)\|^2 \leq \|u-v\|^2 - 2\epsilon \quad (9)$$

implies $u \cdot v - \phi(u) \cdot \phi(v) \leq -\epsilon$, which gives us $|u \cdot v - \phi(u) \cdot \phi(v)| \geq \epsilon$. \square

Problem 2: Locality-Sensitive Hashing (LSH) for Angle Similarity

Suppose our set of n points $D = \{p_1, \dots, p_n\}$ are vectors in d dimensions. Our problem is: given a query point q find a point $p \in D$, which has a small angle with q . Recall that the angle between two vectors a and b is $\cos^{-1}\left(\frac{a \cdot b}{\|a\|\|b\|}\right)$.

As doing this exactly may be computationally expensive, let us try to do this approximately with a fast algorithm. The approximate objective is as follows: suppose there exists a point $p \in D$ which has angle less than θ with q , then our goal is return a point with angle less than $c\theta$, where $c > 1$.

Let us try to do this with LSH. Let us consider the a family of hash functions, where $h(p) = \text{sign}(u \cdot p)$ where we will sample u uniformly at random from a Gaussian (or from a unit sphere).

1. Provide an exact expression for $\mathbb{P}(h(p) = h(p'))$ based on some geometric relation between p and p' .

Solution

Define

$$\text{angle}(u, v) = \cos^{-1}\left(\frac{u \cdot v}{\|u\|\|v\|}\right), \quad (10)$$

which is the angle between two vectors.

Then,

$$\boxed{\mathbb{P}(h(p) = h(p')) = 1 - \frac{\text{angle}(p, p')}{\pi}}. \quad (11)$$

2. Provide an expression for P_1 and P_2 in terms of θ and $c\theta$. Note that since we want a small angle, we should use:

- (a) If $\text{angle}(p, p') < \theta$, then $\mathbb{P}(h(p) = h(p')) \geq P_1$.
- (b) If $\text{angle}(p, p') > c\theta$, then $\mathbb{P}(h(p) = h(p')) \leq P_2$.

Solution

If $\text{angle}(p, p') < \theta$, then

$$\mathbb{P}(h(p) = h(p')) = 1 - \frac{\text{angle}(p, p')}{\pi} \geq 1 - \frac{\theta}{\pi},$$

so $\boxed{P_1 = 1 - \frac{\theta}{\pi}}.$

If $\text{angle}(p, p') > c\theta$, then

$$\mathbb{P}(h(p) = h(p')) = 1 - \frac{\text{angle}(p, p')}{\pi} \leq 1 - \frac{c\theta}{\pi},$$

so $\boxed{P_2 = 1 - \frac{c\theta}{\pi}}.$

3. Provide expressions for query time for point q , the space to store the hash tables, and the construction time of our datastructure.

Solution

Suppose we have L hash functions. If we use the algorithm discussed in class, to query a point, we need to compute L hashes. Then, up to 3 times, we iterate through the buckets: for each bucket, we choose a point and check how close it is to q ; if it is $c\theta$ close, we stop. The worst case is that we decide there exists no point that is θ close to q . In this case, we iterate through the L buckets 3 times, so the time complexity is $O(L)$.

For the space needed to store the hash tables, we need to store L bits for each point, so the space needed is $O(nL)$.

For construction, we need to compute L hashes for each point, so the computational complexity is $O(nL)$ as well.

Problem 3: Dual Coordinate Ascent

Consider the problem

$$\min_w L(x), \text{ where } L(x) = \sum_{i=1}^n (w \cdot x_i - y_i)^2 + \lambda \|w\|^2. \quad (12)$$

1. Show that the solution for Equation 12 is obtained for weights

$$w^* = (X^\top X + \lambda I)^{-1} X^\top Y \quad (13)$$

$$= \frac{1}{\lambda} X^\top \alpha^*, \quad (14)$$

where $\alpha^* = (I + XX^\top/\lambda)^{-1}$.

Solution

Proof. We can take the derivative of L in Equation 12 directly. Note that $D(x \mapsto Ax)(x) = A$ and $D(x \mapsto x^\top x)(x) = 2x^\top$. Therefore by the chain rule,

$$D(x \mapsto (Ax)^\top (Ax))(x) = 2x^\top A^\top A. \quad (15)$$

We can reformulate Equation 12 as a function of w

$$\begin{aligned} l_{X,y}(w) &= (Aw - y)^\top (Aw - y) + \lambda w^\top w \\ &= (Aw)^\top (Aw) - 2y^\top Aw + y^\top y + \lambda w^\top w. \end{aligned} \quad (16)$$

Taking the derivative, we have that

$$D(l_{X,y})(w) = 2w^\top X^\top X - 2y^\top X + 2\lambda w^\top. \quad (17)$$

Setting Equation 17 to 0 and solving for w , we have

$$\begin{aligned} 0 &= 2w^\top X^\top X - 2y^\top X + 2w^\top \\ w^\top (X^\top X + \lambda I) &= y^\top X \\ (X^\top X + \lambda I) w &= X^\top y \\ w &= (X^\top X + \lambda I)^{-1} X^\top y. \end{aligned}$$

Since Equation 16 is a quadratic form, the problem is convex, and

$$w^* = (X^\top X + \lambda I)^{-1} X^\top Y$$

minimizes Equation 12.

Now, note that

$$(X^\top X + \lambda I) X^\top = X^\top X X^\top + \lambda X^\top = X^\top (X X^\top + \lambda I).$$

Multiplying on the left by $(X^\top X + \lambda I)^{-1}$ and on the right by $(X X^\top + \lambda I)^{-1}$, we have that

$$X^\top (X X^\top + \lambda I)^{-1} = (X^\top X + \lambda I)^{-1} X^\top.$$

Substituting this into Equation 13, we obtain

$$\begin{aligned} w^* &= X^\top (X X^\top + \lambda I)^{-1} y \\ &= X^\top \left(\lambda \left(I + \frac{X X^\top}{\lambda} \right) \right)^{-1} y \\ &= \frac{1}{\lambda} X^\top \left(I + \frac{X X^\top}{\lambda} \right)^{-1} y, \end{aligned}$$

which gives us the desired result. \square

If $\lambda = 0$, in general, this is not true since $XX^\top + \lambda I$ may not be invertable when $n > d$. However, if $d \geq n$, and $\text{rank}(X) \geq n$, Equation 14 may still be well-defined.

2. Define

$$G(\alpha_1, \alpha_2, \dots, \alpha_n) = \frac{1}{2} \alpha^\top (I + XX^\top / \lambda) \alpha - Y^\top \alpha. \quad (18)$$

Start with $\alpha = 0$. Choose coordinate i randomly, and update

$$\alpha_i = \arg \min_z G(\alpha_1, \dots, \alpha_{i-1}, z, \alpha_{i+1}, \dots, \alpha_n). \quad (19)$$

Show that the solution to the inner optimization problem for α_i is:

$$\alpha_i = \frac{y_i - \frac{1}{\lambda} \left(\sum_{j \neq i} \alpha_j x_j \right) \cdot x_i}{1 + \|x_i\|^2 / \lambda}. \quad (20)$$

Solution

Proof. We can take the partial derivative of Equation 18 directly to obtain

$$\begin{aligned} \frac{\partial G}{\partial \alpha_i} &= \alpha_i + \frac{1}{\lambda} (\alpha^\top XX^\top)_i - y_i \\ &= \alpha_i + \frac{1}{\lambda} \left(\alpha_i \|x_i\|^2 + \sum_{j \neq i} \alpha_j (x_j \cdot x_i) \right) - y_i. \end{aligned} \quad (21)$$

Setting Equation 21 to 0, solving for α_i , and taking advantage of convexity, we find

$$\alpha_i = \frac{y_i - \left(\sum_{j \neq i} \alpha_j x_j \right) \cdot x_i}{1 + \|x_i\|^2 / \lambda} \quad (22)$$

minimizes Equation 18 as a function of α_i , and solves the inner optimization problem. \square

3. What is the computational complexity of this update, as it is stated?

Solution

The complexity of updating α_i with Equation 20 is $O(nd)$ since we need to iterate over the n rows of X , and take the d -dimensional dot product of each row with x_i .

4. What is the computational complexity of one stochastic gradient descent update?

Solution

The complexity of one stochastic gradient descent update is $O(d)$. We computed the derivative in Equation 17 for the full matrix X . In stochastic gradient descent we'd replace X by a vector by randomly sampling a row from X . Then, to compute the gradient we have to do some dot products along with scalar operations.

5. Now consider the procedure.

- Start with $\alpha = 0$, $w = \frac{1}{\lambda} X^\top \alpha = 0$.
- Choose coordinate i randomly and perform the following update:
 - Compute the differences:

$$\Delta\alpha_i = \frac{(y_i - w \cdot x_i) - \alpha_i}{1 + \|x_i\|^2 / \lambda} \quad (23)$$

- Update the parameters as follows:

$$\begin{aligned} \alpha_i &\leftarrow \alpha_i + \Delta\alpha_i \\ w &\leftarrow w + \frac{\Delta\alpha_i}{\lambda} x_i. \end{aligned} \quad (24)$$

Prove that the update rule in Equation 24 is valid.

Solution

Proof. Let α' and w' be the result of updating coordinate i of α . Assume that $w = \frac{1}{\lambda} X^\top \alpha$. This is true when $\alpha = 0$. We will show that this invariant holds as α is updated.

To see that, the update rule for w is valid, we can rewrite

$$w = \frac{\alpha_1}{\lambda} x_1 + \cdots + \frac{\alpha_i}{\lambda} x_i + \cdots + \frac{\alpha_n}{\lambda} x_n, \quad (25)$$

so

$$\begin{aligned} w' &= w + \frac{\Delta\alpha_i}{\lambda} x_i \\ &= \frac{\alpha_1}{\lambda} x_1 + \cdots + \frac{\alpha_i + \Delta\alpha_i}{\lambda} x_i + \cdots + \frac{\alpha_n}{\lambda} x_n \\ &= \frac{\alpha_1}{\lambda} x_1 + \cdots + \frac{\alpha'_i}{\lambda} x_i + \cdots + \frac{\alpha_n}{\lambda} x_n \\ &= \frac{1}{\lambda} X^\top \alpha'. \end{aligned}$$

Thus, the w update is valid.

To see that the α update is valid, we show that Equations 20 and 24 are equivalent. Both algorithms initiate $\alpha = 0$, so they are equivalent at the initial step.

By using the definition $w = \frac{1}{\lambda} X^\top \alpha$,

$$\begin{aligned} \alpha'_i &= \alpha_i + \Delta\alpha_i \\ &= \frac{(y_i - w \cdot x_i) - \alpha_i}{1 + \|x_i\|^2 / \lambda} + \frac{\alpha_i + \alpha_i \|x_i\|^2 / \lambda}{1 + \|x_i\|^2 / \lambda} \\ &= \frac{1}{1 + \|x_i\|^2 / \lambda} \left(y_i - \frac{1}{\lambda} \left(\sum_{j \neq i} \alpha_j x_j \right) \cdot x_i - \frac{1}{\lambda} \alpha_i \|x_i\|^2 + \frac{1}{\lambda} \alpha_i \|x_i\|^2 \right) \\ &= \frac{y_i - \frac{1}{\lambda} \left(\sum_{j \neq i} \alpha_j x_j \right) \cdot x_i}{1 + \|x_i\|^2 / \lambda}, \end{aligned}$$

so both update rules are equivalent. \square

6. What is the computation complexity of the update defined by Equations 23 and 24?

Solution

The computation complexity is $O(d)$. Computing the dot product when computing $\Delta\alpha_i$ and updating w are both $O(d)$ operations. Everywhere else, we do scalar operations.

This is much faster than the $O(nd)$ update for Equation 20.

Problem 4: Project Milestone

Build a simple object detection model. Our object detector will comprise of a binary classifier per category: given features of an image patch corresponding to a bounding box, does this patch contain an object of the category of interest?

Solution