

CS-541 A Midterm

Name: Stephen Szemis

1. I find the idea of random projection very interesting (since it is a new topic for me).

I suppose the idea that you can generate a random matrix for some size k (in lecture it's call A), then multiply it by some original data (called X), and get a lower dimensional representation without losing relative distance seems

very counter intuitive. After all, compared to PCA, it's very straight forward.

Of course the proof uses Chebyshev's Inequality to eventually reach Johnson-Lindenstrauss Lemma, so RP does have a serious theoretical backing.

2. Let X be random variable.

- $E[X] = 1$, this does tell us something about $P(X)$. It tells us that $P(X)$ sampled many times, averages to 1. However it isn't very practical to base real world decisions solely on expected value.

- $\text{Var}[X] = 10$, this (along with $E(x)$) gives much stronger information about the behavior of $P(X)$. While before we could only make weak claims about $P(X)$ underlying behavior, now we can give tight bounds around how quickly we can reach accumulated values of X , say for gambling.

2. continue

This all comes from Markov's and Chebyshelv's Inequalities.

If $E(X)$ gives us a "center" for our underlying distribution, $\text{Var}[X]$ gives us how "spread out" our distribution is from that center.

- 3.
- We want probability ≥ 0.99 .
 - No more than 100 dollars (e.g. $n \leq 100$)
 - Let $\Pr(X=1) = p$ where $X=1$ is a correct label. $\Pr(X=-1) = 1-p$ where $X=-1$ is WRONG!

$$\begin{aligned} E(X) &= \Pr(X=1) \cdot 1 + \Pr(X=0) \cdot 0 = p \\ \text{Var}[X] &= \Pr(1-p) \end{aligned}$$

$$S = \Pr\left(\sum_{i=1}^n x_i > 0\right)$$

- Use Chebyshelv's Inequality for high probability argument

$$\Pr(|X - E[X]| \geq t) \leq \frac{\text{Var}[X]}{t^2}$$

$$E(S) = \sum_{i=1}^n (p + p - 1) = n(2p - 1)$$

$$\text{Var}[S] = \sum_{i=1}^n \left(1 - (2p - 1)^2\right) = n(4p - 4p^2) = 4pn(1-p)$$

So...

$$\Pr(|S - n(2p - 1)| \geq t) \leq \frac{4pn(1-p)}{t^2}$$

↳ continue to next page

① 3.

$$n(2p-1) - t \geq 1, \text{ at least majority}$$

$$\frac{4np(1-p)}{t^2} \leq 0.01, \text{ at least } 0.99$$

$$n \leq \frac{0.01t^2}{4p(1-p)} \leq 99, \text{ from } n < 100 \text{ (const)}$$

$$t = \sqrt{0.01 \cdot 99 \cdot 4p(1-p)}$$

$$99(2p-1) - \sqrt{0.01 \cdot 99 \cdot 4p(1-p)} \geq 1$$

$$f(p) = 99(2p-1) - \sqrt{0.01 \cdot 99 \cdot 4p(1-p)} - 1 \geq 0$$

4. • Major advantages:
- Not dependent on input data.
 - Keeps relative distance
 - Not very expensive
- You would not gain a computational boost if your new k dimensions is not suitably smaller than d .
- For some $O(nd)$ algorithm you have computational gain of
- $$nd - nk = n(d - k)$$
- But RF cost some operations, say C then for it to be worth it
- $$n \geq \frac{C}{d-k}$$
- If n is not large enough, then dimension reduction is not worth it.
- Using another distribution would probably work, but we still require a standard deviation of around $1/\sqrt{k}$ and mean of 0. Such a distribution would need to be designed in the discrete case based on k . A bit more work. Though sampling discretely is probably less computationally costly than normal distribution for large matrixes.

5. The main idea of collaborative filtering is to populate a sparse matrix of user data by "group" users by their similar interests. For instance I may have a row in an amazon product database

$$[\dots, 0, 1, 0, 1, \dots]$$

$n \ n+1 \ n+2 \ n+3$

And the professor may have this as his row

$$[\dots, 1, 1, 0, 1 \dots]$$

$n \ n+1 \ n+2 \ n+3$

We have similar $n+1$ and $n+3$ interest, so we may have similar n index interest.

From that assumption we fill may n th item index with some value. And our matrix is slightly less sparse.

Drawbacks for this formulation is that we must assume there exists some low rank representation of our data, which is not always true in all applications. It also isn't clear how to build such a system in a scalable way. For instance amazon must deal with million by million matrix of user items. Very sparse. There are also more complications in a binary case, since we need to recover an actual continuous value in order to base our distance measures.

Some improvements might include clustering data so we can work on smaller subsets of our data. Lowering dimensionality of data. Using a machine learning model to find more arbitrary patterns in our data.