

Class Notes

Topics in Big Data Analytics

Data Reduction

Theorem: Let S be any set of n points in \mathbb{R}^d . Then, \exists a projection f such that $\forall u, v \in S, \|u - v\|^2(1 - \epsilon) \leq \|f(u) - f(v)\|^2 \leq (1 + \epsilon)\|u - v\|^2$. We can find such a project in random polytime.

$$f: \mathbb{R}^d \rightarrow \mathbb{R}^k \quad \text{for } k > 4 \left(\frac{\epsilon^2}{2} - \frac{\epsilon^3}{3} \right)^{-1} \ln n$$

Proof: Consider a random unit vector $Y = \frac{1}{\|X\|}(X_1, \dots, X_d)$ where $X_i = N(0, 1)$, all independent. Let $Z = \frac{1}{\|X\|}(X_1, \dots, X_k)$. The expected length of $Z = k/d = \mu$.

Claim: $P[\|Z\| \leq \frac{\beta k}{d}] \leq e^{\frac{k}{2}(1-\beta+\ln \beta)} \forall \beta < 1$ and $P[\|Z\| > \frac{\beta k}{d}] \leq e^{\frac{k}{2}(1-\beta+\ln \beta)} \forall \beta > 1$.

Fact: $E[e^{sX^2}] = \frac{1}{\sqrt{1-2s}}$ for any $-\infty < s < 1/2, X = N(0, 1)$.

Fact: (Markov's inequality) If X is a non-negative random variable with $E[X] = \mu$ then $P[X \geq a\mu] \leq a^{-1} \Rightarrow P[X \geq 1] \leq \mu$.

Proof: $P[d(X_1^2 + \dots + X_k^2) \leq \beta k(X_1^2 + \dots + X_d^2)] = P[\beta k(X_1^2 + \dots + X_d^2) - d(X_1^2 + \dots + X_k^2) \geq 0] = P\left[e^{t(\beta k(X_1^2 + \dots + X_d^2) - d(X_1^2 + \dots + X_k^2))} \geq 1\right]$.

By Markov's inequality the previous probability is $\leq E\left[e^{t(\beta k(X_1^2 + \dots + X_d^2) - d(X_1^2 + \dots + X_k^2))}\right]$. We convert it to,

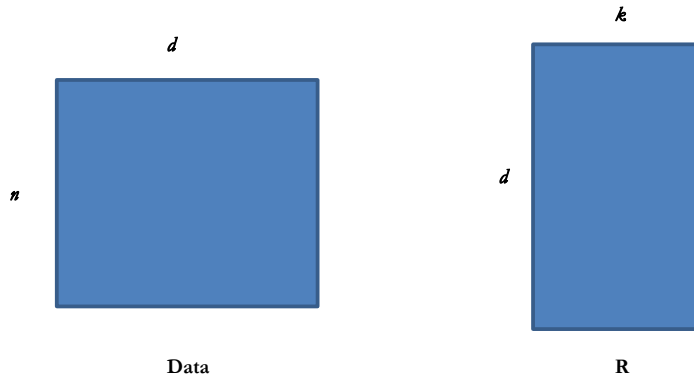
$$E\left[(e^{tk\beta X^2})^{d-k} (e^{t(\beta k-d)X^2})^k\right] = \left(\frac{1}{\sqrt{1-2tk\beta}}\right)^{d-k} \left(\frac{1}{\sqrt{1-2t(\beta k-d)}}\right)^k \text{ where } tk\beta < 1/2.$$

To get the best probability we differentiate it with respect to t , equate it to zero, then we substitute this value of t in it and simplify.

Let $\beta = (1 - \epsilon)$. Then, $P[\|Z\| \leq (1-\epsilon)\frac{k}{d}] \leq \frac{1}{n^2}$ and $P[\|Z\| \geq (1-\epsilon)\frac{k}{d}] \leq \frac{1}{n^2}$. Then, $P[(1-\epsilon)\frac{k}{d} \leq \|Z\| \leq (1+\epsilon)\frac{k}{d}] \leq \frac{2}{n^2} \Rightarrow$ For a fixed pair, the distance between them is more than a factor β away from their distance in \mathbb{R}^d is $\leq \frac{2}{n^2} \Rightarrow$ The probability that this happens for at least one pair is $\leq \binom{n}{2} \frac{2}{n^2} = 1 - \frac{1}{n}$. Repeat this process $\alpha n \ln n$ times. Then, the probability of failure in all of them is $\leq (1 - \frac{1}{n})^{\alpha n \ln n} \leq \left[(1 - \frac{1}{n})^n\right]^{\frac{1}{n} \alpha n \ln n} \leq e^{-\alpha \ln n} = n^{-\alpha}$. If $u \in \mathbb{R}^d$ and $u = (X_1, \dots, X_d)$ then $f(u) = \sqrt{\frac{d}{k}}(X_{i_1}, \dots, X_{i_k})$ where i_1, \dots, i_k are picked randomly.

Runtime: $O(kn^3 \ln n)$.

Achleoptas Procedure



Consider the above matrices. Each row is a point and each element of R will be $= \begin{cases} \sqrt{3} & \text{with prob } \frac{1}{6} \\ 0 & \text{with prob } \frac{2}{3} \\ -\sqrt{3} & \text{with prob } \frac{1}{6} \end{cases}$

Claim: This projection works with high probability.

Runtime: $O(ndk)$

Learning Algorithm

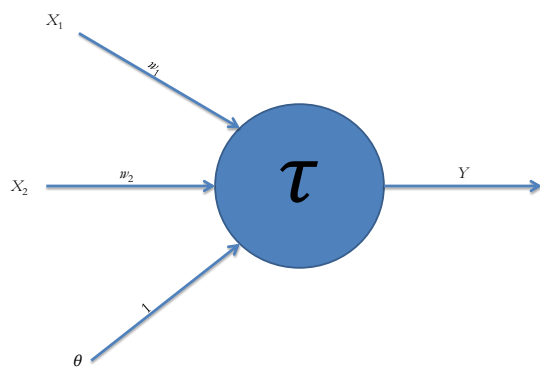
Given positive and negative examples the goal is to learn a concept.

Performance Criteria:

- Sample Complexity \square # of examples
- Time Complexity \square Learning time

Artificial Neural Networks:

It's a directed graph $G(V, E)$ where V are the processors and E are edges with weights.



$$\begin{aligned} Y &= 0 && \text{if } w_1X_1 + w_2X_2 + \theta < 0 \\ Y &= 1 && \text{otherwise} \end{aligned}$$