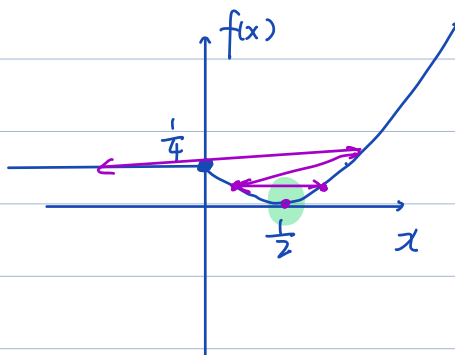1. Devise the update rule for minimizing $f(x) = (\max(0,x) - \frac{1}{2})^2$ using gradient descent. Roughly, when will gradient descent succeed or fail? (based on where you start and step size.)

① Draw the plot for $f(x) = (\max(0,x) - \frac{1}{2})^2$

$$= \begin{cases} \frac{1}{4} & , x \leq 0 \\ (x - \frac{1}{2})^2 & , x > 0 \end{cases}$$

Continuous. $\frac{1}{4} = (0 - \frac{1}{2})^2$



② Compute

$$f'(x) = \begin{cases} 0 & , x \leq 0 \\ 2(x - \frac{1}{2}) = 2x - 1 & , x > 0 \end{cases}$$

③ $x_0 \leq 0$ will fail because it always get 0 gradient.

④ $x_{n+1} = x_n - \eta \cdot \nabla$

$x_0 > 0 \qquad x_{n+1} = x_n - \eta \cdot (2x - 1)$.

Step size need to be reasonably small. Otherwise, it will diverge.

Let's say $x_0 = \frac{1}{4}$  $\eta = 1$  $x_1 = \frac{1}{4} - 1 \cdot (-\frac{1}{2}) = \frac{3}{4}$,

$x_2 = \frac{3}{4} - 1 \cdot (\frac{3}{4} \times 2 - 1) = \frac{1}{4}$

$x$ will be trapped in $\frac{1}{4}$ & $\frac{3}{4}$.

$\eta = 0.01$

2. Given a <u>data set</u> A, under what circumstances will its projection onto its principal components equal its projection on

    1. Right singular vectors      A    $m > n$

    2. Left singular vectors.      A    $n > m$

In PCA, we find the principal components by calculating the eigendecomposition of $A^T A$ or $A A^T$ given a data set A.

If A is $m \times n$, $A^T A \longrightarrow$ sample covariance matrix.

    sample      features      $n \times n$    (i,j : How similar feature i to feature j)

If A is $n \times m$.    $A A^T$.    $n \times n$      Q

Data   left   right

$A = U S V^T$    $A^T A = V S (U^T U) S V^T = V S^2 V^T$    right singular vector

$A A^T = U S (V^T V) S U^T = U S^2 U^T$    left . . . .

$i \in T$      $j \in S$

3. Let S be a set of documents and let T be set of terms. Suppose that C is a binary term-document incidence matrix. (So entry $\langle i,j \rangle$ is 1 if term i appear in document j and 0 otherwise).

What do the entries of $C^T C$ represent?

    ① Dimension of C:   # term x # docs.
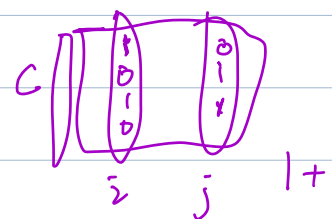
    ② Dimension of $C^T C$:   # docs x # docs

$C^T C \langle i,j \rangle$ = ith row of $C^T$ . jth column of C

             = ith column of C . jth column of C.

column k of C means if each term is in doc k.

             = # shared terms in doc i & j.

If i=j, # terms in doc i/j.

4. How did we derive the equations for simple linear regression from class.

Page 19-20 from CS 378H.

Single variable $\beta_1 = \dfrac{\overline{xy} - \overline{x}\,\overline{y}}{\overline{x^2} - (\overline{x})^2} = \dfrac{Cov(x,y)}{Var(x)}$

$\beta_0 = \overline{y} - \beta_1 \overline{x}$

Multiple - Variable Case.

$X$ $m \times n$. $y \in R^m$. Goal: minimize $\|X \cdot w - y\|^2$.

Normal Equation $w = (X^T X)^{-1} X^T y$. $O(n^3 + mn^2)$

Take gradient to $\|X \cdot w - y\|_2^2$ with respect to $w$.

$\nabla = 2(X \cdot w - y)^T \cdot X = 0$

$\underset{m \times 1}{\underbrace{\dfrac{m \times n \cdot n \times 1}{\vphantom{1}}}} \quad m \times n$

$2(X \cdot w)^T \cdot X = 2 y^T \cdot X$

$(X^T X)^{-1} X^T X w = X^T y$

$w = (X^T X)^{-1} X^T y$.

5. Why can you assume without loss of generality that a mistake - bounded Learner only updates its state when it makes a mistake?

Re-ordering the samples

No matter you update the state or not when you predict it right.

It always need to make $t$ mistakes to ensure no mistakes in future.

So updating its state when predicting right is redundant.

6. You are given a data set with **m points** and an algorithm that satisfies the **weak-learing condition**. (It always outputs a classifier with accuracy **60%**). Each classifier output by the weak-learning algorithm can be encoded using **two bits**. How can you construct a classifier that can be described by less than **m** bits and **is correct on every data point in the data set**. (You may assume m is very large).

What's the **size** of **your final classifier?**

$T$          Trainning.

$\varepsilon \leq \frac{1}{m} \rightarrow 0$

Adaboost.

After $T$ iterations of the algorithm, the error of $h_{final} = $
MAJ $(h_1, \cdots h_T)$ is at most $e^{-2T\gamma^2}$.   $T \geq \dfrac{\log(\frac{1}{\varepsilon})}{2\gamma^2}$
Majority

makes the error of $h_{final}$ at most $\varepsilon$.   $\gamma, \varepsilon$

Since m is very large   all correct means
$\varepsilon \leq \frac{1}{m}$

$\gamma = 1 - \varepsilon \approx 1$        $T \geq \log(m)$   $O(\log(m))$

$\gamma = 0.1$

$T \geq \dfrac{\log(\frac{1}{\varepsilon})}{\gamma^2}$        $\varepsilon = \frac{1}{m}$

$\frac{1}{\varepsilon} = m$

$= \dfrac{\log(m)}{\gamma^2}$        $\gamma^2 = 1 - \varepsilon = 1$

$= \log(m)$

T. Markov's inequality. [ weakest ]

$$\forall X \geq 0, \quad P(X \geq a) \leq \frac{E(X)}{a}$$

non-negative $\quad a > 0$

$$E(X) = P(X \leq a) \cdot E(X | X < a) + P(X \geq a) \cdot E(X | X \geq a)$$

$$E(X) \geq P(X \geq a) \cdot \underbrace{E(X | X \geq a)}_{\geq a.} \geq \underbrace{a \cdot P(X \geq a)}$$

Chebyshev's Inequality : ( X random variable with mean $\mu$ and variance $\sigma^2$ )

$$P_r(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2} \quad \text{for any } k > 0.$$

$$P_r(|X - \mu| \geq k\sigma) = \underbrace{P_r((X-\mu)^2 \geq k^2\sigma^2)}_{\text{Apply Markov inequality}} \leq \frac{E[(X-\mu)^2]}{k^2\sigma^2} = \frac{\sigma^2}{k^2\sigma^2} = \frac{1}{k^2}$$

] $Y = (X - \mu)^2$ with $a = k^2\sigma^2$

$$X = (x - \mu)^2 \qquad a = k^2 a^2.$$

$$\text{Markov :} \quad P(X \geq a) \leq \frac{E[(X-\mu)^2]}{k^2 a^2} = \frac{\sigma^2}{k^2\sigma^2} = \frac{1}{k^2}$$

Chernoff Bounds (Strongest)

$$P_r(X \geq a) = P_r(e^{t \cdot X} \geq e^{t \cdot a}) \leq \frac{E(e^{t \cdot X})}{e^{t \cdot a}}$$