

SMAI-M20-L17: Perceptron and Loss Functions

C. V. Jawahar

IIIT Hyderabad

September 18, 2020

Announcements

- **Quiz 1** Already announced.
- **Some Observations:**
 - Number of views of the video before the lecture session is low. (some times even after the session :-)). May be most of you are very comfortable.
 - Number of attempts for CR is high. Good enough. Performance is OK. No concerns.
 - Number of people who ask queries to TAs for office hours is low. Increased a bit.
 - Ask: (i) Are you putting effort (ii) Are you asking questions at the right forums (iii) Are you coming prepared?
 - Higher Education and Online education expects you to put effort beyond the class rooms.
- **Round up Session on Monday 21.** (no new content)
 - See all videos (and lecture sessions?), if required more than once.
 - Ask specific queries (a form will be shared)
 - Do post by Saturday mid night. (latest by Sunday noon)

Class Review

$$M \leq x^T A x$$

$$s.t. \|x\| = 1$$

$$x^*$$

$$\|w^{k+1} - v^k\| < \epsilon$$

Appreciating PCA

- Is PCA good for compression?
- How does it compare with linear regression?
- Can we have a GD based computation for PCA? (there is a small extra step of norm being unity)
- How are eigen values, vectors related to the covariances?



$\nabla_{\text{over } x} \in \text{subset}$

All samples
batch



one of e^i

Single s

~~stock~~
Min bdr

$$\omega^{k+1} \leftarrow \omega^k - \eta \nabla J$$

Recap:

- Supervised Learning:
 - Notions of Training, Validation and Testing; Loss Function and Optimization, Generalization, Overfitting, Occam's razor, Model Complexity, Bias and Variance, Regularization.
 - Performance Metrics, Estimating error using validation set.
- Approaches:
 - Optimal Decision as ω_1 if $P(\omega_1|\mathbf{x}) \geq P(\omega_2|\mathbf{x})$ else ω_2
 - MLE
 - Dimensionality Reduction and Representation (Feature Selection, PCA, Neural Embeddings)
 - Application of PCA: Eigen Face
 - Matrix Factorization for Data Matrices (SVD, Eigen Decomposition)
 - Application of Matrix Factorization: LSI, Matrix Completion, Recommendation Systems)
 - Nearest Neighbour, Linear Discriminants
 - Gradient Descent
 - Linear Regression: Closed form, GD, Regularization Optimization
 - Perceptron Algorithm and Neuron Model

This Lecture:

$$y_n = +1$$
$$y_n = -1$$

$$\text{sign}(\tilde{w}^T x)$$

1 Perceptron -II

- Analysis of Perceptron Algorithm

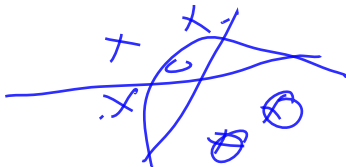
2 Logistic Regression

- Probabilistic view in defining loss/goal. (MLE next)

3 Loss Functions

- MSE, MAE, Hinge Loss, Cross Entropy

Questions? Comments?



Discussions Point - I

Look at three popular functions we know:

A 0-1 loss function

B Hinge loss

C $P(\omega|x)$ in Logistic regression

Comment on:

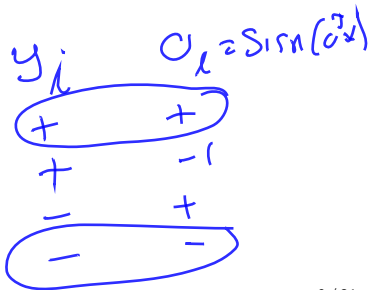
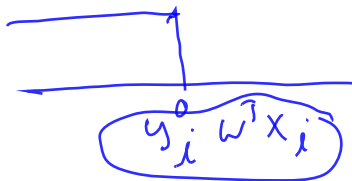
① Smoothness

② Convexity

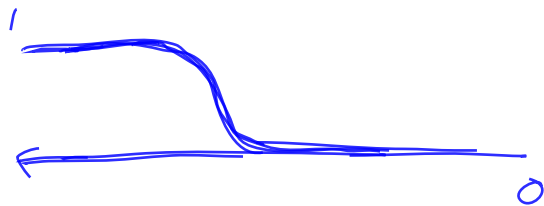
③ Continuity

④ Differentiability

0-1 loss
count

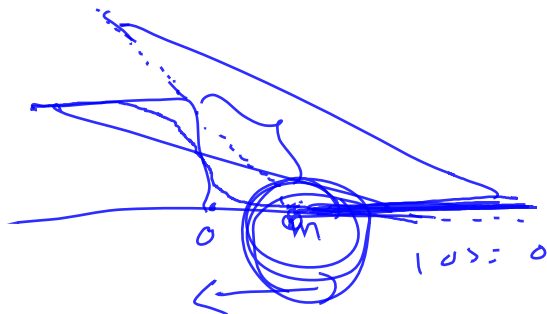


Blank



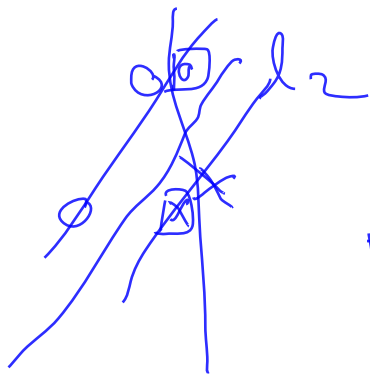
$$\frac{1}{1 + e^{cTx}}$$

$$\frac{1}{1 + e^{-\alpha Wx}}$$



$$\max(a, m - y_{i,j})$$

correct class small penalty



$$w^T x > m_1 \quad l_1$$

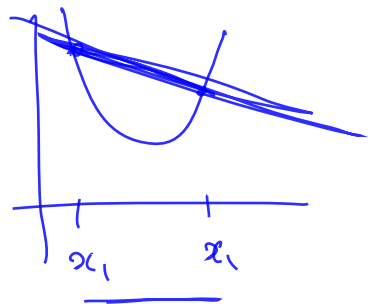
$$w^T x < -m_1$$

l_1, l_2 are
active for x

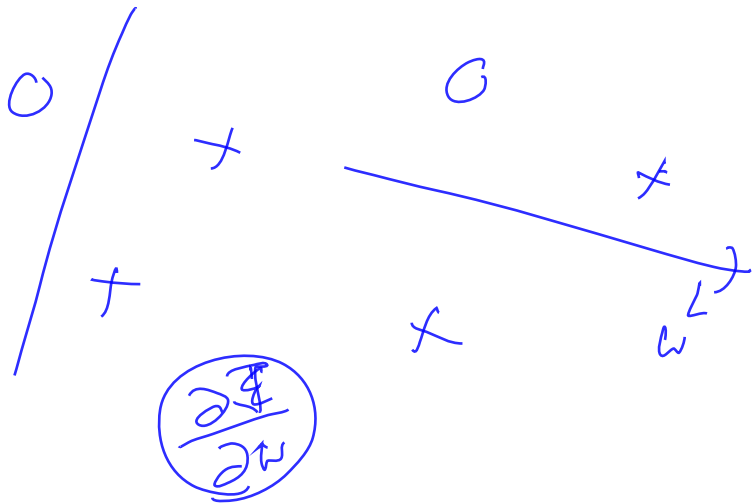
$$\text{sing}(w^T x) \geq 0 \quad x$$

$$w^T x > \underline{m}_1$$

$$w^T x < -\underline{m}_1$$



Convex opt



Discussion Point - II: Are the perceptrons doing GD?

$$J' = \sum_{i=1}^N \underline{(t_i - o_i)^2}$$

$$t_i = y_i \\ o_i = \text{sign}(w^T x_i)$$

with $o_i = \text{sign}(\mathbf{w}^T \mathbf{x}_i)$? There is a problem. With small changes in the \mathbf{w} or the line, the J is not changing. Let us now define:

$$J = \sum_{i=1}^N (t_i - o_i)^2 \underline{(-\mathbf{w}^T \mathbf{x}_i)}$$

JP

This additional terms pulls (or pushes) proportionally. Let us now rewrite this J as sum of two parts one over \mathcal{E} and the other on not in \mathcal{E} .

$$J = J_1 + J_2 = \sum_{\mathbf{x}_i \in \mathcal{E}} (t_i - o_i)^2 (-\mathbf{w}^T \mathbf{x}_i) + \sum_{\mathbf{x}_i \text{ not in } \mathcal{E}} (t_i - o_i)^2 (-\mathbf{w}^T \mathbf{x}_i)$$

$$J = J_1 + J_2 = \sum_{\mathbf{x}_i \in \mathcal{E}} (2 \cdot t_i)^2 (-\mathbf{w}^T \mathbf{x}_i) + \sum_{\mathbf{x}_i \text{ not in } \mathcal{E}} 0 \times (\mathbf{w}^T \mathbf{x}_i)$$

$$t_i = 0$$

We know that J_2 is zero. When $\mathbf{x}_i \in \mathcal{E}$, $(t_i - o_i)$ is $2t_i$ i.e., $2y_i$.

Discussion Point - II: Are the perceptrons doing GD?

- We know either $(t_i - o_i)$ is zero Or we know that $(t_i - o_i)$ is $2t_i$
- And $\frac{\partial(-\mathbf{w}^T \mathbf{x}_i)}{\partial \mathbf{w}} = -\mathbf{x}_i$.

$$\nabla J = -2 \sum_{i=1}^N (t_i - o_i) \mathbf{x}_i$$

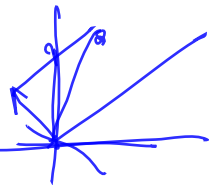
This leads to: $\mathbf{w}^{k+1} \leftarrow \mathbf{w}^k - \eta \nabla J$ as:

$$\mathbf{w}^{k+1} \leftarrow \mathbf{w}^k - \eta \sum_{i=1}^N (t_i - o_i) (-\mathbf{x}_i)$$

$$\mathbf{w}^{k+1} \leftarrow \mathbf{w}^k + \eta \sum_{\mathbf{x} \in \mathcal{E}} 2 \cdot t_i \mathbf{x}_i$$

$$\mathbf{w}^{k+1} \leftarrow \mathbf{w}^k + \eta' \sum_{\mathbf{x}_i \in \mathcal{E}} y_i \mathbf{x}_i$$

t_i	o_i	$t_i - o_i$
+1	+1	0
-1	-1	0
+1	-1	2
-1	+1	-2



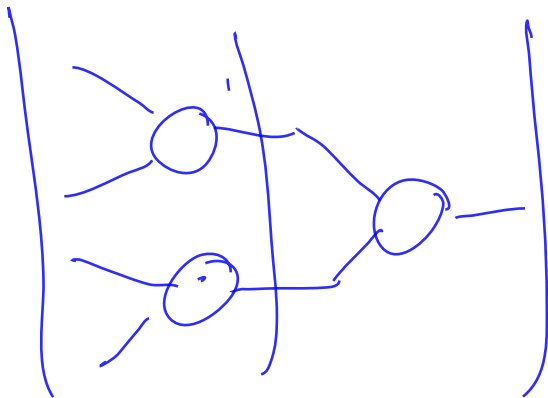
(with changes in scale for learning rate η)

Discussions Point -III: Perceptron Convergence

- **Claim 1** If there exist a set of weights that are consistent with the data, the perceptron algorithm will converge.
- **Claim 2** If the training data is not Linearly Separable, the perceptron algorithm will eventually repeat the same set of weights and thereby enter an infinite loop.
- **Claim 3** If the training data is linearly separable, algorithm will converge in a maximum of N steps. Find N .
- **Claim 4** Every boolean function can be represented by some network of perceptrons only two levels deep.

Read and understand the proofs for the above claims¹

¹Page 229 of Duda, Hart and Stork:
https://cds.cern.ch/record/683166/files/0471056693_TOC.pdf



What Next:? (next three)

- ① Logistic Regression
- ② Multi Class Classification (beyond binary)
- ③ More Dimensionality Reduction Schemes (eg. LDA/Fisher)