

SMAI-M20-L20: LDA

C. V. Jawahar

IIIT Hyderabad

September 25, 2020

Announcements

- ① **Quiz 1** on Next Wed.
- ② **Topics:** Topics remain the same.
- ③ Most-Likely the same time.
- ④ Any other announcements: by Monday.

Class Review

Consider the following three samples and their labels $((x_1, x_2), y)$:

$$\{((1, 1), +), ((2, 2), -), ((0, 0), +)\}$$

Look at the perceptron update rule with $\eta = 0.1$

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

Classify as + ve if $\mathbf{w}^T \mathbf{x} \geq 0$ else - ve.

Given \mathbf{w}^0 . What do we know about \mathbf{w}^1 and \mathbf{w}^2 ?

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
 - Logistic Regression
 - LDA
 - Multi-Class Classification Architectures
- Fisher's Face*

This Lecture:

① Logistic Regression - III

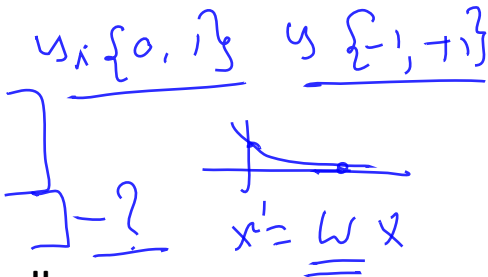
- Insight into LR objective

② LDA - II

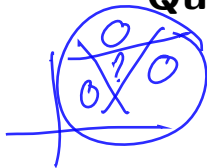
- LDA solution

③ Multi-Class Classification - II

- D-DAG



Questions? Comments?



- compact
- well separated

$$P(y_1 = +1 | x, w) =$$

$$P(y = -1 | x, w)$$

$$\frac{1}{1 + e^{-w^T x}}$$

$$\frac{1}{1 + e^{w^T x}}$$

$$\left(1 - \frac{1}{1 + e^{w^T x}}\right)$$

$$\frac{1}{1 + e^{-y_1 w^T x}}$$

$$y_1 = -1$$

$$1 - \frac{1}{1 + e^{-w^T x}} = \frac{1 + e^{-w^T x}}{1 + e^{-w^T x}} - \frac{1}{1 + e^{-w^T x}} = \frac{e^{-w^T x}}{1 + e^{-w^T x}} = \frac{1}{1 + e^{w^T x}}$$

Discussions Point - I

We know the solution to LDA as

$$\mathbf{w}^* = \alpha \mathbf{S}_W^{-1} [\mu_A - \mu_B]$$

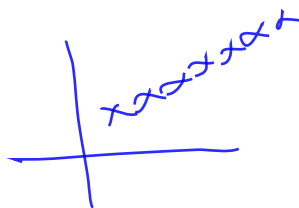
A potential worry is "If \mathbf{S}_W is singular? "

- Suggest a configuration of the data when \mathbf{S}_W can be singular?
- Suggest solutions to handling this singularity problem while computing \mathbf{w}^* ?

N d

$N > d$
 $d \not> N$

$S_c = d \times d$



$$\begin{aligned}
 S_v &= \underbrace{\sum_{x_i \in A} [x_i - \mu_A][x_i - \mu_A]^T}_{N_1} + \sum_{x_i \in B} [x_i - \mu_B \dots]^T \\
 &\quad N_1 + N_2
 \end{aligned}$$

N

$$\begin{aligned}
 &\underbrace{100 \times 100}_{\sum_{i=1}^5 x_i x_i^T} \\
 &\quad \min(d, N)
 \end{aligned}$$

$$\textcircled{1} \quad \underline{S_w} \leftarrow [S_u + \underset{\substack{\uparrow \\ \text{small}}}{\times} I]$$

$$\textcircled{2} \quad \text{Pseudo inverse.}^{-1} \\ S_u = (U D V^T)^{-1} = \underbrace{V^{-1} D^{-1} U^T}_{S_u^{-1}}$$

$$\textcircled{3} \quad \text{Step 1 } \underbrace{1000 \rightarrow 200}_d \text{ Do PCA} \quad \left(\begin{array}{l} d = 1000 \\ N = 500 \end{array} \right)$$

$$\text{Step 2 Do LDA } \underline{S_u} \quad \underbrace{d \times d}_{200 \times 200}$$

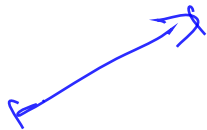
$$S_B u = \lambda S_U u$$

Second w

$$[M_A - M_B] [M_A - M_D]^T w = \lambda S_U u$$

$$\alpha = \text{scalar} / \lambda$$

$$w = \alpha S_U^{-1} [M_A - M_D]$$



$$\|v\| \leq 1$$

$$\underline{v} = \underline{S_U^{-1} [M_A - M_D]}$$

"generalized E.V problem"

PCA

$$Ax = \lambda x \rightarrow \text{E.V. Pr}$$

$$Ax = \lambda Bx \rightarrow \underline{\underline{\text{G. E.V. Pr}}}$$

LDA

$$S_B = [A - \bar{x}][A - \bar{x}]^T$$

S_B is singular

Discussions Point -II

Are there any design considerations in D-DAG?

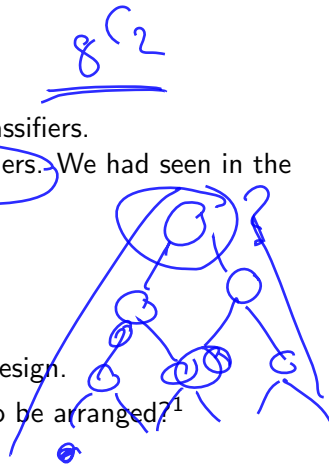
We know the DAG way of arranging pair-wise classifiers.

(Assume we have 4 classes and 6 pairwise classifiers. We had seen in the micro-lecture.)

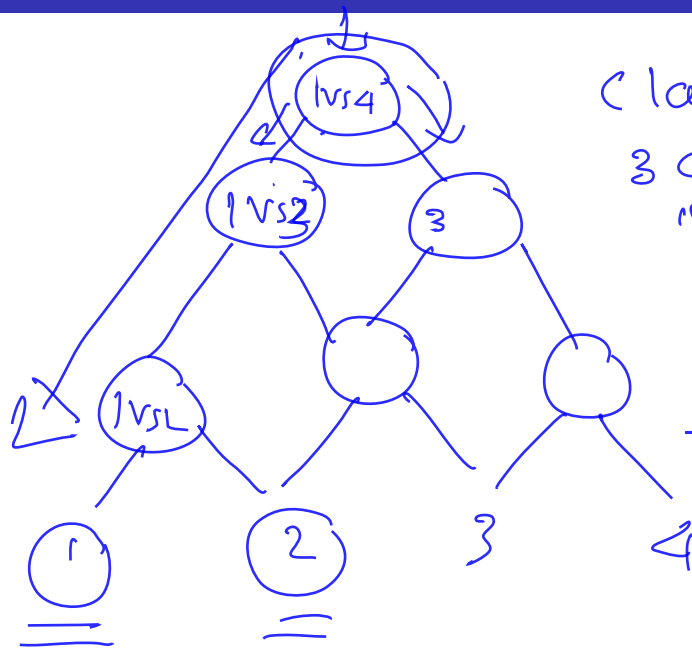
What will you prefer as the root classifier?

- the highest accuracy classifier
- the least accuracy classifier
- any classifier. This does not matter in the design.

Any other insight into how the classifiers/class to be arranged?¹



¹Read later: An old but relevant analysis:
<https://cvit.iiit.ac.in/images/ConferencePapers/2003/pavan03multiclass.pdf>



① Appreciate DAG

- ① how to arrange, how is
- ② Comp Complexity

DP

② There are design concerns

- which class of the dy
- Ass problems?

③ Design an algorithm

Ref

R (

Research

Better user to des?

Discussion Point - III

Comment on the following three different ways of formulating the loss for a binary classification²:

① $\sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2$

② $\sum_{i=1}^N (y_i - g(\mathbf{w}^T \mathbf{x}_i))^2$

③ LR objective

mse

why LR did not?
use this loss

— Non-convex

MLE

?

²Read later: somewhat relevant reference: <http://books.jackson.me/Cross-Entropy-vs-Squared-Error-Training-a-Theoretical-and-Experimental-Comparison.pdf>

What Next:? (next three or even more)

- ① More on LR, Multi-Class Classification, Dimensionality Reduction
- ② Intro to SVMs and Kernels.