

SMAI-M20-L18:Roundup Session

C. V. Jawahar

IIIT Hyderabad

September 21, 2020

Class Review

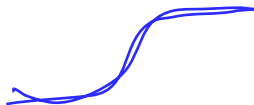
$$\begin{bmatrix} 0 \\ 1 \end{bmatrix} \rightarrow \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

$$\frac{1}{He^{-z}} \times \frac{e^z}{e^z}$$

Consider the sigmoid function $g(\alpha z) = \frac{1}{1+e^{-\alpha z}}$

- When α varies what happens to this function?
- Where is its range, what is its value at zero, max value, min value,
- How is this related to tanh? $\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$
- What is the derivative of sigmoid?
- What is $1 - g(z)$?

Sigm
tanh



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

This Lecture:

① Plans and Preparation for the Quiz

② Topic 1:

- Feature vector, Data Matrix, Bayesian Optimal Classifier, MLE etc.
- Eigen Decomposition, SVD, LSI etc.

③ Topic 2: Supervised Learning

- Training, Testing, Validation, Performance Metrics, Overfitting, Regularization etc.
- Loss Functions, Optimization, Bias and Variance

④ Topic 3: Algorithms

- PCA, Linear Regression, Nearest Neighbour, Naive Bayes etc.
- Perceptrons, Gradient Descent, Logistic Regression (not for Quiz)

6.30-7.10

Questions? Comments?

Quiz Preparation

- Go through the class material well enough
 - Micro-Lecture Videos
 - Lecture Session (slides/recordings) and some references there in.
 - Class Reviews
 - Home Works
- Read the text book related chapters, Solve problems and practice
- Not expected for Quiz1:
 - Additional references/material list available (moodle)
 - Lecture notes and additional references
 - (<https://www.dropbox.com/sh/h91lhc0xpmh2ekw/AAC-FuNgqOO0-txb3FvJqns-a?dl=0>)
- Use computer; Keep phone next to you.
- Keep pen, paper and a calculator for your use.
- Strictly avoid any communication with classmates during this period.

Example 5 Questions; Need not be the same

- 1 Test whether you are a student of SMAI 2020 (like CAPTCHA) (may be 10 questions that you can solve fast if you are a student in SMAI-2020) ¹
- 2 Test whether you not only enrol but also attempt to follow² SMAI-2020.
- 3 Numerical Problems
- 4 A question set very similar to the regular class review one.
- 5 A bit more involved. Keep couple of A4 sheets and a paper ready.

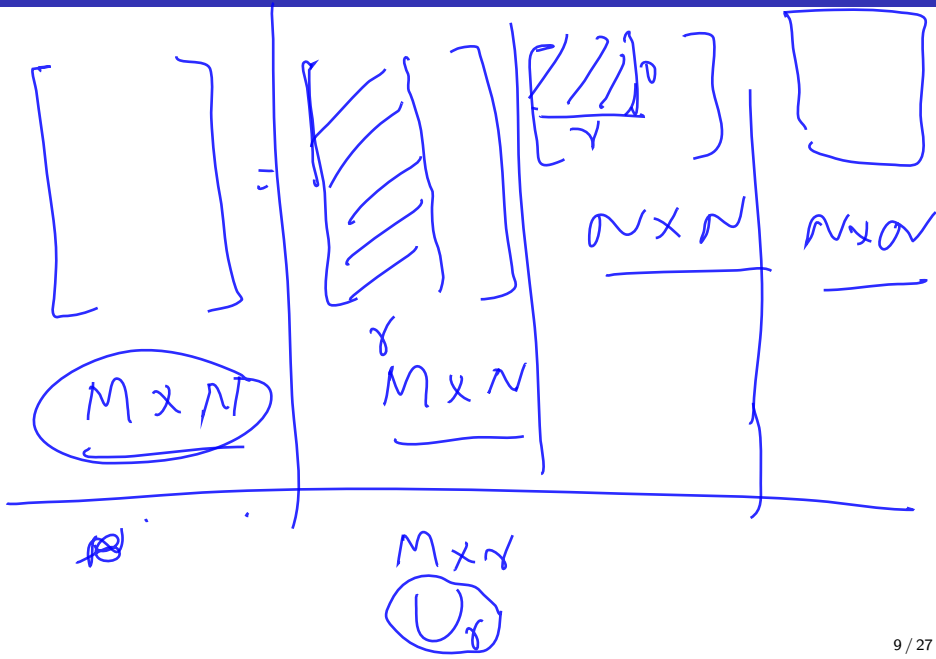
¹CAPTCHA stands for Completely Automated Public Turing test to tell Computers and Humans Apart.

²Follow Class Material

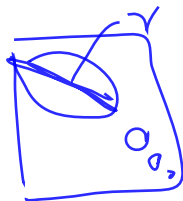
Topic - I

$$\{(c_i, d_i)\}$$

- Feature vector, Data Matrix, Bayesian Optimal Classifier, MLE etc.
- Eigen Decomposition, SVD, LSI etc.
- Q: In posterior probability $P(w|x)$, does x denote class label or feature in the vector x ?
 $P(w|x)$ $P(x|w)$
- Yes;
- In SVD we do $A = UDV^T$. The dimensions of U are given as $m \times r$, but we know that U is the matrix of eigen vectors of AA^T , and AA^T is symmetric so it is supposed to have independent eigen vectors. If we assume a to be $M \times N$ matrix, the dimension of you must be $M \times M$ not $M \times R$, since all the eigen vectors are linearly independent.



[]



$n > n$
 $n > n$

\max

~~$\text{rank} - r$~~

$$A = \sum_{i=1}^r x_i x_i^T$$

$n \times n$ $n \times 1$

$$\begin{matrix} m > n \\ n \neq m \end{matrix}$$

$$A = U D V^T$$

$$\hat{A} = V \underline{D} U^T$$

Topic - II

$$\begin{matrix} \times & A & \times \\ \times & & \times \end{matrix}$$

eigen vect

PGD

$$\frac{\partial}{\partial} = 0$$

- Training, Testing, Validation, Performance Metrics, Overfitting, Regularization etc.
- Loss Functions, Optimization, Bias and Variance
- Q: Given a regularization problem, how do we minimise the loss?

$$L = \text{Emp } L + \text{Res Terms}$$
$$L = \sum_{i=1}^n \underline{\underline{l(y_i, f(x_i))}} + ||w||_2^2$$

$$L = L_1 + \lambda L_2$$

$$\frac{\partial L}{\partial w} = \frac{\partial L_1}{\partial w} + \lambda \frac{\partial L_2}{\partial w}$$

$$L = \min_w \left[\sum (y_n - w^T x_n)^2 \right]$$

$$\max_w \left[w^T A w \right] \quad \text{subject to } w^T w = 1$$

$2Ac$

~~$\frac{\partial L}{\partial w}$~~

$$w^{k+1} \leftarrow w^k - \eta \nabla L$$

$$L = \omega^T A \omega$$

$$(\|\omega\| = 1)$$

$$\frac{\partial L}{\partial \omega} = 2A\omega$$

$\omega^0 \leftarrow \text{random}$

$$\textcircled{1} \omega^{k+1} \leftarrow \omega^k + \eta \frac{\partial L}{\partial \omega}$$

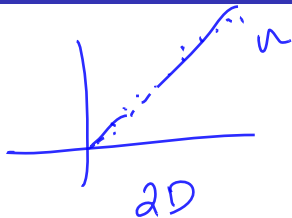
$$\textcircled{2} \text{normalize } \omega^{k+1} \text{ s.t. } (\|\omega^{k+1}\| = 1)$$

- PCA, Linear Regression, Nearest Neighbour, Naive Bayes etc.
- Perceptrons, Gradient Descent, Logistic Regression (not for Quiz)

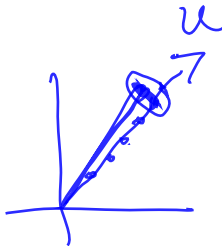
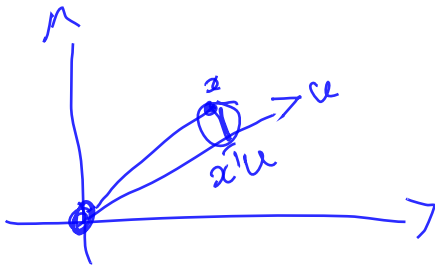
Blank

$$x_i \rightarrow x_i^T u$$

α_n



- Q: Steps of minimizing loss functions to obtain $\min ||X - X'||$



Blank

Let $\mathbf{u}_1, \dots, \mathbf{u}_d$ be d orthonormal vectors. We can represent the vectors \mathbf{x} as $\sum_{i=1}^d \alpha_i \mathbf{u}_i$. Where the scalar α_i is $\mathbf{x}^T \mathbf{u}_i$. However, if we use smaller than d basis vectors, there could be some loss or reconstruction error. Let us consider the loss when we use only one \mathbf{u} . i.e.,

$$\mathbf{x} - \mathbf{u}\mathbf{u}^T \mathbf{x}$$

— ①

Sum of the reconstruction loss for all the N data samples is now:

$$\sum_{i=1}^N \|\mathbf{x}_i - \mathbf{u}\mathbf{u}^T \mathbf{x}_i\|^2$$

$$\mathbf{u}^T \mathbf{u} = 1$$

②

$$= \sum_{i=1}^N \left(\mathbf{x}_i^T \mathbf{x}_i + (\mathbf{u}\mathbf{u}^T \mathbf{x}_i)^T (\mathbf{u}\mathbf{u}^T \mathbf{x}_i) - 2\mathbf{x}_i^T \mathbf{u}\mathbf{u}^T \mathbf{x}_i \right)$$

③

We would like to minimize this. First term is positive (non negative). It is independent of \mathbf{u} . Therefore, we would like to minimize:

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

④

We also know that $\mathbf{u}^T \mathbf{u} = 1$.

min ||x - x' ||

$$= \sum_{i=1}^N -\mathbf{x}_i^T \mathbf{u} \mathbf{u}^T \mathbf{x}_i = \sum_{i=1}^N -\mathbf{u}^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{u} = \underline{\underline{-\mathbf{u}^T \Sigma \mathbf{u}}}$$

Minimizing the reconstruction error now becomes that of Maximizing

$$\text{Max } \mathbf{u}^T \Sigma \mathbf{u}$$

with our familiar constraint of $\mathbf{u}^T \mathbf{u} = 1$. This reduces the solution as the eigen vectors corresponding to the largest eigen values.

$$\boxed{\text{Max } \mathbf{u}^T \Sigma \mathbf{u} \text{ s.t. } \mathbf{u}^T \mathbf{u} = 1}$$

Blank

Diagram illustrating the structure of the input x and the output z in the context of the proposed architecture:

- The input x is represented as a vertical rectangle.
- The output z is represented as a horizontal rectangle.
- The output z is composed of two parts, each of size $d \times d$.
- The output z is also composed of two parts, each of size $n \times d$.

X
x x d

F is vec
of x_2

$$\frac{1}{N} \sum_{i=1}^N [x_i \phi] [x_i \phi]^T$$

Ex. of
Covariation

What Next:? (next three)

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