

# SMAI-M20-L12:PCA

C. V. Jawahar

IIIT Hyderabad

September 4, 2020

Topics for Quiz 1 (10 days from now):

- Upto (including PCA) (yet to see some more details)
  - Mathematical Foundations (LA, Prob, Opt.)
  - Supervised Learning (formulation, Simple Algorithms, Bayesian Optimal, Related Concepts, Performance Metrics and Evaluation)
  - Matrix Factorization and Applications
  - Linear Regression
  - PCA
- New topics (Not for Q1): Perceptrons and Gradient Descent. (even if you see these in the class before the Q1)

# Class Review

(use notations and conventions from the class) Consider the problem of linear regression where we minimize the loss

$$\mathcal{L}_1 = \frac{1}{N} \sum_{i=1}^N \alpha_i (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \lambda_1 g(\mathbf{w})$$

where  $g()$  is a regularization term. We also write the loss in matrix form as

$$\mathcal{L}_2 = \frac{1}{N} [Y - \mathbf{X}\mathbf{w}]^T A [Y - \mathbf{X}\mathbf{w}] + \lambda_2 g(\mathbf{w}).$$

- What can we say about the matrix  $A$ ? size? properties? elements?
- What can we say about the optima, equivalence?
- What can we say when the objective is regularized?
- When regularizing what happens to the objective/optimal solution?



# Recap:

- Problem Space:
  - Learn a function  $y = f(\mathbf{W}, \mathbf{x})$  from the data.
  - Learn useful features
- Supervised Learning:
  - Notion of Training, Validation and Testing
  - Loss Function and Optimization
  - Performance Metrics, Estimating error using validation set.
  - Need of Generalization, overfitting, Occam's razor, model complexity, Bias and Variance, Regularization.
- Classification Algorithms:
  - Nearest Neighbour Algorithm
  - Linear Classification; Linear Regression
  - Decide as  $\omega_1$  if  $P(\omega_1|\mathbf{x}) \geq P(\omega_2|\mathbf{x})$  else  $\omega_2$
- Mathematical Foundations: Linear Algebra, Probability, Optimization
  - SVD, Eigen Decomposition, MLE
  - Matrix Factorization in LSI, Recommendation Systems

# This Lecture:

$$\log a + \log b - \log c - \log d$$
$$\log \frac{a}{c} + \log \frac{b}{d}$$

## Micro-Lecture Videos

- 1 PCA: Dimensions that preserve maximum variance

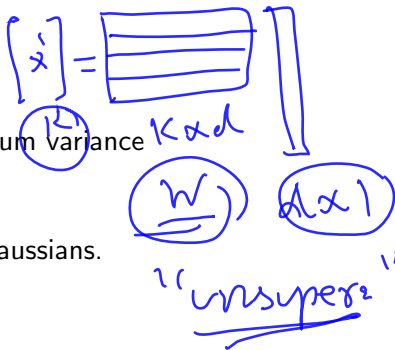
- $\mathbf{x}' = \mathbf{W}\mathbf{x}$
- How to find  $\mathbf{W}$ ?

- 2 Decision boundaries for Multivariate Gaussians.

- Analysis under ideal situations.

- 3 Deep/Neural Embeddings.

- What are good features?
- How data can help in discovering features?



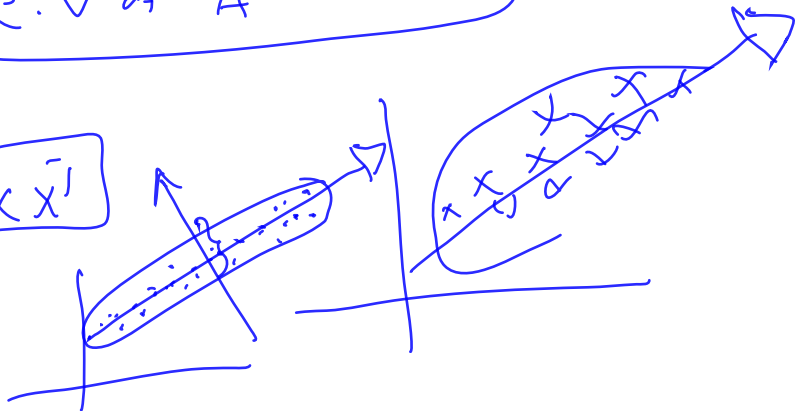
Questions? Comments?

$$\log a + \log b - \log c - \log d$$
$$\log \frac{a}{c} + \log \frac{b}{d}$$
$$\log(p \times q) - \log p - \log q$$

$w^T A w$  s.t.  $c^T w = 1$   
 e.v. of  $A$

$$p(w_1|x) \geq p(w_2|x)$$

$|x x^T|$



$$\underline{P(w_1|x)} \geq \underline{P(w_2|x)}$$

$$\frac{a \cdot b}{\cancel{\times}} < \frac{d \cdot e}{\cancel{\times}}$$

$$a \cdot b < d \cdot e$$

$$\log a + \log b < \log d + \log e$$

$$\log a + \log b - \log d - \log e \geq 0$$



## Discussions Point - I

$$\underline{w_1}x_1 + \underline{w_2}x_2 = 10$$

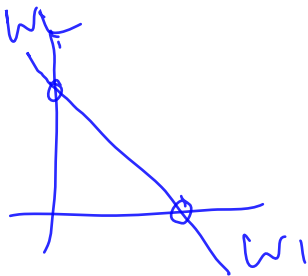
Regularization with Lp norm is popular. Why does L1 norm induce sparsity?

Let us take an example:

$$\min_{\mathbf{w}} (w_1x_1 + w_2x_2 - 10)^2 + \lambda g(\mathbf{w})$$

Let  $g(\mathbf{w})$  is the Lp norm of  $\mathbf{w}$

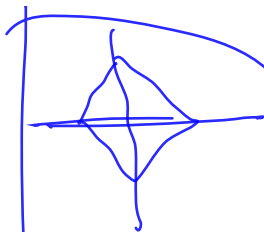
- First term has many solutions and at all these points the term is zero. (locus is a line)
- We need to find a solution (on this line) that minimizes the  $g(\mathbf{w})$ .



min err  
w

min err + 2p norm

Lasso



$$\|w\|_1 = K$$

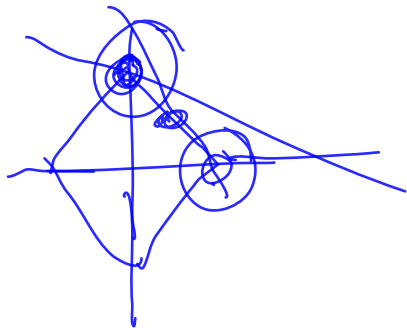
$$|w_1| + |w_2| = K$$



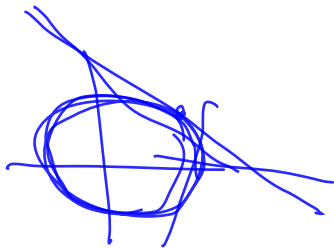
$$\|w\|_2 = 0$$

$$\sqrt{w_1^2 + w_2^2} = K$$

↑  
Ridge



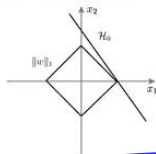
L1 norm  
lead to  
sparsity.



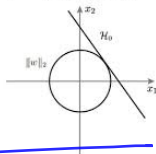
L2 need  
not be

# Some Plots(from Internet)

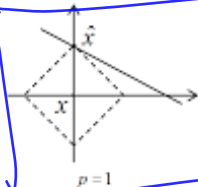
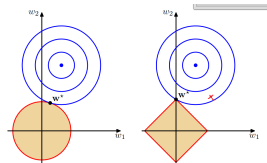
**A** L1 regularization



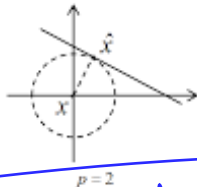
**B** L2 regularization



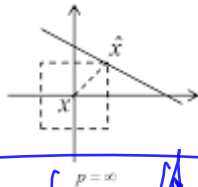
**Figure 3.4** Plot of the contours of the unregularized error function (blue) along with the constraint region (3.30) for the quadratic regularizer  $q = 2$  on the left and the lasso regularizer  $q = 1$  on the right, in which the optimum value for the parameter vector  $w$  is denoted by  $w^*$ . The lasso gives a **sparse** solution in which  $w_1^* = 0$ .



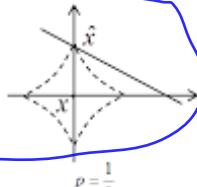
$p = 1$



$p = 2$



$p = \infty$



$p = \frac{1}{2}$



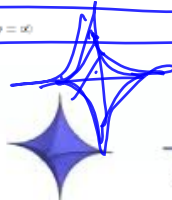
$p = \infty$



$p = 2$



$p = 1$



$0 < p < 1$



$p = 0$



## Discussions Point -II

$$\boxed{w = (X^T X)^{-1} X^T y} \cdot \frac{\partial (X^T A X)}{\partial X} \quad \boxed{2AX}$$

Can we have a closed form expression for ridge regression just like the simple linear regression?

$$w^T X^T X w$$

$$\boxed{\sum_{i=1}^N (y_i - w^T x_i)^2 + \lambda w^T w} \rightarrow \underline{\underline{L_2}}$$

Hint:

$$\underline{\underline{[Y - Xw]^T [Y - Xw] + \lambda w^T w = [Y - Xw]^T [Y - Xw] + \lambda w^T (I^T I) w}}$$
$$\boxed{w = (X^T X + \lambda \underline{\underline{I}})^{-1} X^T y} \quad (X^T X)^{-1} X^T y$$

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

near singular

$$(X^T X + \lambda I)^{-1}$$





## What Next:? (next three)

- ① PCA and Dimensionality Reduction (more)
- ② Bayesian Optimal (Cont.)
- ③ What are good features?
- ④ Perceptron Algorithm
- ⑤ Gradient Descent Optimization