

Week 6 MLT

Lecture 1 :

$$w^* = \hat{w}^{ML} \quad \leftarrow \text{max likelihood} = (x x^T)^+ x y$$

$$y_{\hat{x}} = w^T x + \epsilon \quad \leftarrow \text{noise } N(0, \sigma^2)$$

can be thought of
as a Gaussian $\rightarrow N(w^T x, \sigma^2)$

How good is \hat{w}^{ML} as a guess for the true w ?

$$w \in \mathbb{R}^d ; \hat{w}^{ML} \in \mathbb{R}^d$$

derived out of an optimization problem

i.e., we maximize the likelihood & get \hat{w}^{ML} as
the solution of an optimization problem !!.

We want a way to understand how good \hat{w}^{ML} is in
estimating w ?

$$E \left(\| \hat{w}^{ML} - w \|^2 \right) = \sigma^2 \cdot \text{trace} \left((x x^T)^{-1} \right)$$

expectation over randomness in y

\Rightarrow the avg. deviations between your estimated \hat{w}^{ML} & w
is product of 'variance of the noise that we add to $w^T x$ '
 (σ^2)

and 'trace of inverse of cor. matrix'.
 $\text{trace}((\mathbf{x}\mathbf{x}^T)^{-1})$

It also means that it depends on the data features via the trace of inverse of our covariance matrix, i.e.,

→ depending on how the features themselves are related to each other will also affect Estimators Goodness !!

LECTURE : 2 Cross Validation for Minimizing MSE

trace $(\mathbf{x}\mathbf{x}^T)^{-1}$

$$\mathbf{A} = \begin{bmatrix} a_1 & & \\ & a_2 & \\ & & \ddots \\ & & & a_d \end{bmatrix}$$

sum of diagonal entries of matrix

$$\text{tr}(\mathbf{A}) = \sum_{i=1}^d a_i = \sum_{i=1}^d \lambda_i$$

Let the eigenvalues of $(\mathbf{x}\mathbf{x}^T)^{-1}$ be $\{\lambda_1, \dots, \lambda_d\}$

∴ Eigenvalues of $(\mathbf{x}\mathbf{x}^T)^{-1} = \{\frac{1}{\lambda_1}, \dots, \frac{1}{\lambda_d}\}$

Mean Sq Error ($\hat{\mathbf{w}}_{ML}$)

$$E\left(\|\hat{\mathbf{w}}_{ML} - \mathbf{w}\|^2\right) = \sigma^2 \left(\sum_{i=1}^d \frac{1}{\lambda_i} \right)$$

1 Consider the following estimator :

$$\hat{w}_{\text{new}} = (x x^T + \lambda I)^{-1} x y$$

$\in \mathbb{R}_d$ $\in \mathbb{R}^{d \times d}$

For some matrix A , let eigenvalues be $\{\lambda_1, \dots, \lambda_d\}$

- What are Eigenvalues of $A + \lambda I$?
 $\{\lambda_1 + \lambda, \dots, \lambda_d + \lambda\}$

$$\text{trace} [(x x^T + \lambda I)^{-1}] = \left(\sum_{i=1}^d \frac{1}{\lambda_i + \lambda} \right)$$

EXISTENCE Theorem :

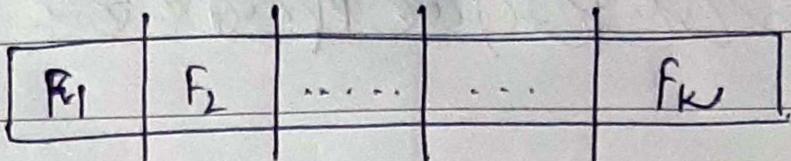
$\exists \lambda \in \mathbb{R}$ such that $\hat{w}_{\text{new}} = (x x^T + \lambda I)^{-1} x y$
 has lesser mean sq. error than \hat{w}_{ML} .

In practice, find λ by cross validation

80%	20%
Training set	Validation set

- Train on the training set & check for error on validation set.

- Pick λ that gives least error.

K-fold Cross Validation

- Train on folds $\{F_1, \dots, F_{i-1}, F_{i+1}, \dots, F_k\}$
 - Validate on F_i
- Pick i that gives least average error !!

Leave One Out Cross Validation

→ not good for large data computationally !!

LECTURE : 3 Bayesian Modelling for Linear RegressionBAYESIAN MODELLING

Need a prior on w , i.e., $P(w) \xrightarrow{R^d}$

Likelihood $y/x \sim N(w^T x, \sigma^2 I)$

(for simplicity
can use σ^2 as well)

Choice for Prior :

$$w \sim N(\vec{0}, \gamma^2 I)$$

$\vec{0} \in R^d$

covariance matrix $R^{d \times d}$

$$\begin{bmatrix} \gamma^2 & 0 & 0 & 0 \\ 0 & \gamma^2 & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & \gamma^2 \end{bmatrix}$$

~~1. (Ans)~~ ~~Ans~~ ~~Ques 1~~ ~~Ques 1~~ ~~Ques 1~~ ~~Ques 1~~ ~~Ques 1~~

As usual,

$$P\left(\omega \mid \{(x_1, y_1), \dots, (x_n, y_n)\}\right) \propto P\left(\{(x_1, y_1), \dots, (x_n, y_n)\} \mid \omega\right) \cdot P(\omega)$$

$$\propto \left(\prod_{i=1}^n e^{-\frac{(y_i - \omega^T x_i)^2}{2}} \right) \cdot e^{-\frac{\| \omega \|_2^2}{2\gamma^2}}$$

$$\propto \left(\prod_{i=1}^n e^{-\frac{(y_i - \omega^T x_i)^2}{2}} \right) \cdot e^{-\frac{\| \omega \|_2^2}{2\gamma^2}}$$

How will the Max. A posterior estimator look like?

$$\hat{\omega}_{MAP} = \arg \min_{\omega} \sum_{i=1}^n \frac{(y_i - \omega^T x_i)^2}{2} + \frac{\| \omega \|_2^2}{2\gamma^2}$$

$$\hat{\omega}_{MAP} = \arg \min_{\omega} \frac{1}{2} \sum_{i=1}^n (y_i - \omega^T x_i)^2 + \frac{1}{2\gamma^2} \| \omega \|_2^2$$

Take gradient, set it to 0 to solve for $\hat{\omega}_{MAP}$

$$\nabla f(\omega) = (2x^T) \omega - 2y + \frac{\omega}{\gamma^2}$$

$$\hat{w}_{MAP} = \left[(x^T x) + \frac{1}{\gamma^2} I \right]^{-1} x^T y$$

→ cross validated in practice

Conclusion : MAP estimation for linear regression with a Gaussian prior $N(0, \gamma^2 I)$ for w is equivalent to "NEW" estimator we used earlier.

LECTURE 4 : Ridge Regression

Linear Regression $\hat{w}_{ML} = \arg \min_w \sum_{i=1}^n (w^T x_i - y_i)^2$

Ridge Regression $\hat{w}_R = \arg \min_w \underbrace{\sum_{i=1}^n (w^T x_i + \gamma_i)^2}_{\text{Loss}} + \lambda \|w\|^2$ λ Regularizer

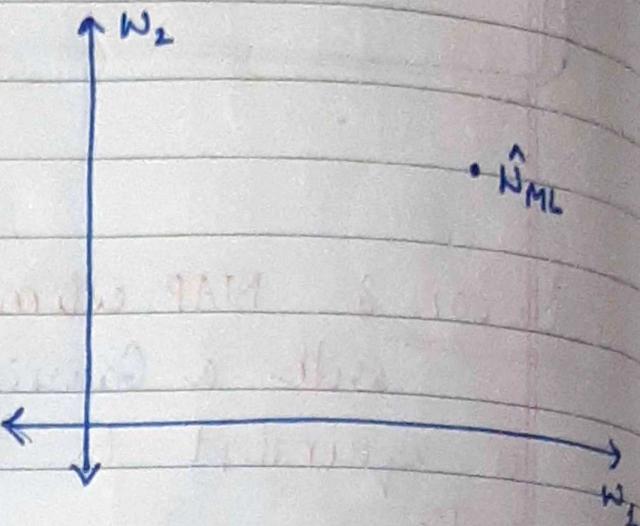
If there are multiple w 's which have the same loss, then the w that you would prefer is the one that has the least length — the one that has many small components!

1 LECTURE : 5

Prelation Between solⁿ of L.R & R.R

PARAMETER SPACE

Where is \hat{w}_R - solⁿ of the ridge regression problem?



Ridge Regression :

$$(A) \rightarrow \arg \min_{w \in \mathbb{R}^d} \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|^2$$

is equivalent to

\equiv

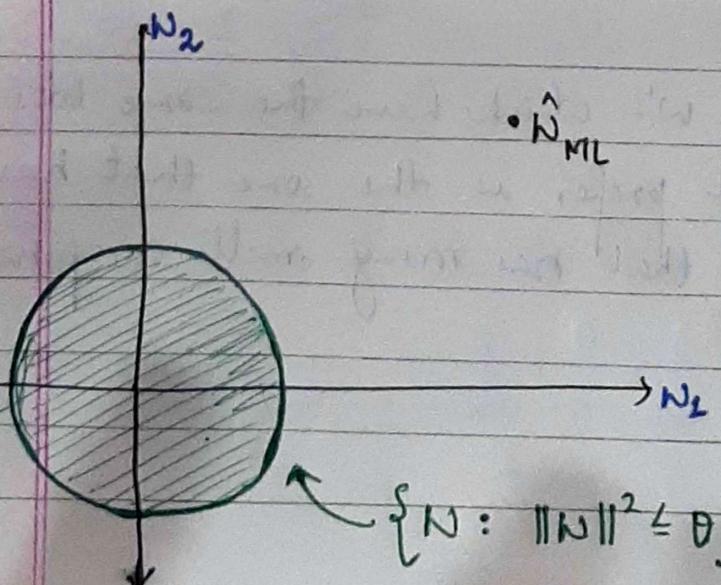
$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^n (w^T x_i - y_i)^2 \quad (B)$$

such that $\|w\|^2 \leq \theta$

$\leftarrow (B)$

depends on λ

For every choice of $\lambda > 0$, $\exists \theta$ such that the optimal solutions of problems (A) & (B) coincide.



- What is the loss/error/obj function value of linear regression at \hat{w}_{ML} ?

$$\sum_{i=1}^n (\hat{w}_{ML}^T x_i - y_i)^2 = \underbrace{f(\hat{w}_{ML})}_{\text{LOSS F.N}}$$

↓
Find the least value among all possible w 's.

Consider the set of all w 's such that

$$f(w) = f(\hat{w}_{ML}) + c \quad [c > 0]$$

$$S_c = \{w : f(w) = f(\hat{w}_{ML}) + c\}$$

→ i.e., every $w \in S_c$ satisfies

$$\underbrace{\|x^T w - y\|^2}_{f(w)} = \underbrace{\|x^T \hat{w}_{ML} - y\|^2}_{f(\hat{w}_{ML})} + c$$

On simplification : we get,

$$(w - \hat{w}_{ML})^T (x x^T) (w - \hat{w}_{ML}) = c'$$

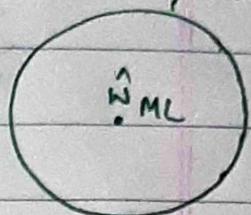
some constant
that depends on
 $c, (x x^T), \hat{w}_{ML}$
and not on w .

If $xx^T = I$,

$$(w - \hat{w}_{ML})^T I (w - \hat{w}_{ML}) = c'$$

$$\Rightarrow \|w - \hat{w}_{ML}\|^2 = c'$$

set of
data point



Conclusion :

- 1. Ridge regression pushes weight values towards 0.
but does not necessarily make it 0.

LECTURE 6 : Relationship b/w L.R & Lasso Regression

- An alternate way to regularize would then be using L_1 norm instead of L_2 norm

$$L_1 \text{ norm} : \|w\|_1 = \sum_{i=1}^d |w_i|$$

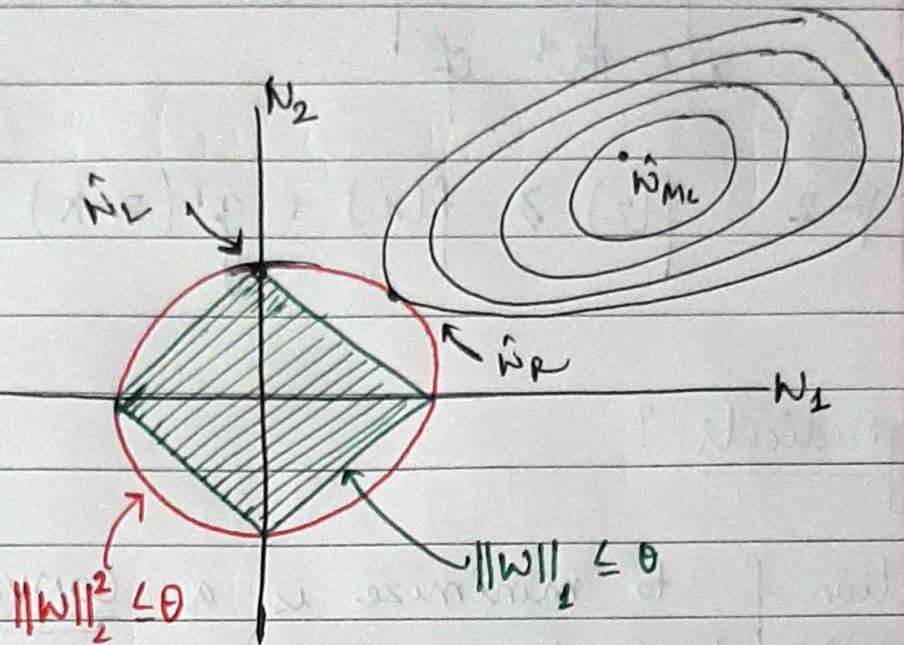
L_1 Regularization :

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|_1$$

It is equivalent to,

JINDAL
PAGE NO.
DATE:

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^n (w^T x_i - y_i)^2 \text{ such that } \|w\|_2 \leq \theta$$



- Instead of L_2 regularizer if you used an L_1 penalty or a regularizer you will get more sparse solution.

↑
This is known as LASSO Regression

LASSO : Least Absolute Shrinkage & Selection Operator

LECTURE 7 : Characteristics of Lasso Regression

- Lasso does not have a closed form solution.
- Sub-gradient methods are usually used to solve LASSO.

Subgradient : A vector $g \in \mathbb{R}^d$ is a subgradient of $f: \mathbb{R}^d \rightarrow \mathbb{R}$ at a point $x \in \mathbb{R}^d$ if

$$\nabla f(z) \geq f(x) + g^T(z-x)$$

Why subgradients?

If function f to minimize is a convex f^n ,
then subgradient descent converges!

- Lasso problems as a convex optimization problem !!.

There are other special purpose methods for LASSO

Eg : IRLS (Iterative reweighted least squares)