

# Week 3: Least Squares/Ridge Regn. MLE / over-underrfitting

## Least Squares:

We have seen the linear regression model with the MSE objective function:  $J_W(x) = XW = Y$

and

$$L(W) = \frac{1}{2N} \sum_{n=1}^N (y_n - x_n^T W)^2 = \frac{1}{2N} \|Y - XW\|^2 = \frac{1}{2N} e^T e$$

$$\nabla L(W) = -\frac{1}{N} X^T e$$

- We have seen that we can find  $W^*$  with gradient methods (GD, SGD, ...) but this problem has also a closed form solution.

We know by the **Principle of optimality principle**:

if  $L(W)$  is convex then  $\nabla L(W^*) = 0$

given that  $L(W)$  is convex (sum of convex fct.)

$$\text{then } \nabla L(W^*) = 0 \Leftrightarrow -\frac{1}{N} X^T (Y - XW^*) = 0$$

$$\Leftrightarrow X^T X W^* = X^T Y$$

$$\text{(normal equations)} \Leftrightarrow W^* = (X^T X)^{-1} X^T Y$$

geom. F. help:  $X^T C = 0$  with  $X^T$  rows being the feature vector

- each feature is  $\perp\!\!\!\perp$  to the  $e$
- $W^*$  is given by Proj of  $y$  on  $\text{Span}(X)$

Note:  $X^T X$  is invertible iff  
 $\text{rank}(X) = D$   
 (full column rank)

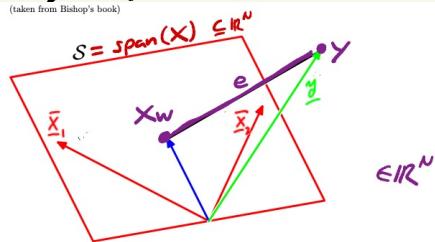
But however,  $X$  is ill-cond.  
 and rank deficient. (column)

- non-linear solver

cond. number =  $\frac{\lambda_{\max}}{\lambda_{\min}}$   $\rightarrow \infty$   
 if ill-cond.

## Linear Regression Summary:

- Grid Search
- Iterative Optim. alg. (GD, SGD ... )  $O(DN)$  per iteration
- Least Squares (closed form sol<sup>1/2</sup>)  
 $\geq O(D^2 N + D^3)$



Maximum Likelihood:  
Probabilistic approach for Lin. Reg.

We assume that our data is generated by:

$$Y_n = X_n^T w + \underbrace{E_n}_{\text{noise}} \sim \mathcal{N}(0, \sigma^2)$$

→ given N inputs, The likelihood of  $Y$  given  $X$ :

$$\underline{P}(Y|X, w) = \prod_{n=1}^N P(Y_n|X_n, w)$$

$$(X_n^T w + E_n) \stackrel{iid}{\sim} \prod_{n=1}^N \mathcal{N}(Y_n | X_n^T w, \sigma^2)$$

⇒ maximize this likelihood

⇒ maximize the log likelihood

$$\mathcal{L}_{LL}(w) = \log \underline{P}(Y|X, w)$$

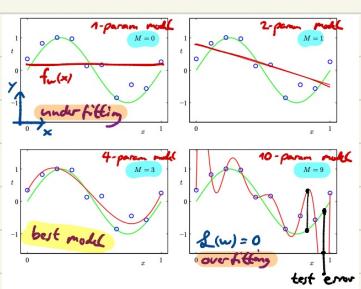
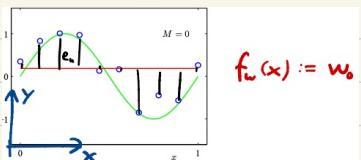
$$= -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - x_n^T w)^2 + \text{const}$$

$$(\mathcal{L}_{MSE}(w) = \frac{1}{2N} \sum_{n=1}^N (y_n - x_n^T w)^2)$$

$$\min_w \mathcal{L}_{MSE}(w) = \max_w \mathcal{L}_{LL}(w) = w^*$$

MLE is consistent:  $w_{MLE} \xrightarrow{P} w_{True}$

## Over-Underfitting:



A common problem in GL is underfitting and overfitting.

Assume we observe a noisy version of samples of raw  $f(x)$ :

$$y_n = g(x_n) + \varepsilon_n, \text{ where } \varepsilon_n \text{ is the noise. (meas. noise, outliers...)}$$

- a simple linear model (1-param) may overly underfit.

→ One way to overcome this is feature expansion

$$\phi(x_n) = [1, x_n, x_n^2, \dots, x_n^M]$$

$\phi(x_n), \log(x_n) \dots$

- add polynomial basis

M: model capacity

$$f_w(x) = \phi_w(\phi(x)) \quad (\text{Linear Model on "augmented" input})$$

$$= \phi(x) w$$

→ To avoid overfitting, we can increase the data overfitting might reduce. (1/N, #Param)

We will see other ways too (Regularization: Ridge-Reg)

# Regularization: Ridge Regression & Lasso

We have seen that by augmenting the feature vector we can make linear models as powerful as we want. However this leads to the problem of overfitting.

Regularization is one way to mitigate this undesirable behavior.

(Occam's Razor)  
Simple models are preferred

Regularization: penalize complex models

$$\min_w \mathcal{L}(w) + \frac{\lambda}{2} \|w\|^2$$

$\downarrow$  complexity of model  $w$   
 $\downarrow$  regularizer

$L_2$ -Regularization: (most frequent)

$$\mathcal{L}(w) = \lambda \|w\|_2^2$$

$\uparrow$  Trade-off par.  $\lambda > 0$

$\lambda \rightarrow 0$	No regd. potential overfit
$\lambda \rightarrow \infty$	underfitting $\ w\  \rightarrow 0$

• When  $\mathcal{L}(w)$  is MSE this is called:

$$\text{Ridge Regression} = \min_w \frac{1}{2N} \sum (y_i - x_i^T w + \lambda \|w\|^2)^2$$

→ Explicit sol:  $p.w$  ( $\lambda = 0 \rightarrow \text{Least Squ.}$ )

$$w^* = (X^T X + \lambda' I)^{-1} X^T y \quad (\lambda' = 2N\lambda)$$

Find  $w$  at optimality

- $\mathcal{L}(w)$  convex
- $\nabla \mathcal{L}(w) = 0$   
 $\Rightarrow w$  optimal

Ridge Regression to fight ill-conditioning: lifting the eigenvalues  
 $(X^T X + \lambda' I)$  eigenvalues are at least  $\lambda' > 0$  so the inverse exists.

$L_1$ -Regularization (The Lasso) ( $\|w\|_1 = \sum_{j=1}^D |w_j|$ )

In combination with GSE for lin. Reg this is called

The Lasso:

$$\min_w \frac{1}{2N} \sum_1^N (y_n - x_n^T w)^2 + \lambda \|w\|_1$$

penalized objective

$\Rightarrow$  this encourages sparse  $w^*$

