

8th Jan

## TUTORIAL

'n' points

$$x_i \in \mathbb{R}^{d \times 1}$$

$$w \in \mathbb{R}^{d \times 1}$$

$$y_i \in \mathbb{R} \text{ (a single number)}$$

$$Y \in \mathbb{R}^{n \times 1}$$

$$X = \begin{bmatrix} \text{--- } x_1^T \text{ ---} \\ \vdots \\ \text{--- } x_n^T \text{ ---} \end{bmatrix}$$

← d →

## Matrix algebra

derivative of inner product (vector  $w$ ):  $\frac{\partial}{\partial w} [a^T w] = \frac{\partial}{\partial w} \left[ \sum_{i=1}^d a_i w_i \right]$

derivative wrt each component of the vector & then stack them up

$$= \begin{bmatrix} \frac{\partial}{\partial w_1} \left[ \sum_{i=1}^d a_i w_i \right] \\ \vdots \\ \frac{\partial}{\partial w_d} \left[ \sum_{i=1}^d a_i w_i \right] \end{bmatrix} = \begin{bmatrix} a_1 \\ \vdots \\ a_d \end{bmatrix} = a$$

$$\frac{d^2 f}{d w^2} = \frac{\partial}{\partial w} \left[ \frac{d f}{d w} \right] \Rightarrow \frac{\partial^2 [w^T A w]}{\partial w^2} = \frac{\partial}{\partial w} [(A + A^T) w] = A + A^T$$

$$\nabla^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial w_1^2} & \frac{\partial^2 f}{\partial w_1 \partial w_2} & \dots & \frac{\partial^2 f}{\partial w_1 \partial w_d} \\ \vdots & \ddots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial w_d^2} \end{bmatrix}$$

← d →

$$\frac{d}{dx} [x^T A x] = (A + A^T)$$

$$\frac{d}{dx} [a^T x b] = a b^T$$

Matrix cookbook for other rules.

$$\text{tr}(A) = \sum_i A_{ii} = \sum_i \lambda_i$$

Trace is sum of diagonal elements / eigen values.

$$\lambda_i = \text{eig}(A)$$

$$|A| = \prod \lambda_i$$

$$\text{tr}(ABC) = \text{tr}(CAB) = \text{tr}(BCA) \quad [\text{cyclically equal}]$$

$$\nabla f(w) = 0$$

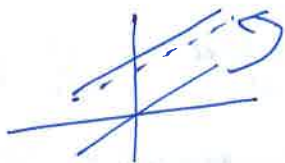
$$\Rightarrow w = (X^T X)^{-1} X^T y$$

In MATLAB:

$$A x = b$$

$$x = A \setminus b$$

Adding bias:



to be able to

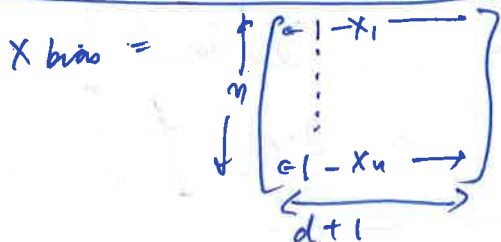
increase expressive power, that is to incorporate the  $y$  intercept term, we add a bias variable.

$$\hat{y}_i = w^T x_i + w_0$$

$$= \begin{bmatrix} w_0 \\ w \end{bmatrix}^T \begin{bmatrix} 1 \\ x_i \end{bmatrix}$$

$$\hat{y}_i = w^T \begin{bmatrix} 1 \\ x_i \end{bmatrix}$$

$$w = \text{argmin} \| X_{\text{bias}} w - y \|^2 \quad \textcircled{*}$$



$$X = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^p \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^p \end{bmatrix} \rightarrow \text{if } d=1, X_{poly} \in \mathbb{R}^{n \times (p+1)}$$

$$\text{if } d=2 \begin{bmatrix} 1 & x_{1,1} & x_{1,2} & x_{1,1}^2 & x_{1,2}^2 & x_{1,1}x_{1,2} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n,1} & x_{n,2} & x_{n,1}^2 & x_{n,2}^2 & x_{n,1}x_{n,2} \end{bmatrix}$$

$p \rightarrow$  degree of polynomial for basis expansion  
 $d \rightarrow$  no. of features in  $X$  [ $x_i \in \mathbb{R}^{d \times 1}$ ]  
 $x_{i,j} \rightarrow$   $j$ th feature of the  $i$ th point.

if  $d > 1 \rightarrow$  use polynomial kernel.  
 RBF kernel  $\rightarrow$   $\infty$  basis expansion.

As  $p \uparrow$ , model complexity  $\uparrow$   
 training error  $\downarrow$   
 generalization  $\rightarrow$  worse  
 test error  $\rightarrow$  high }  $\rightarrow$  overfitting

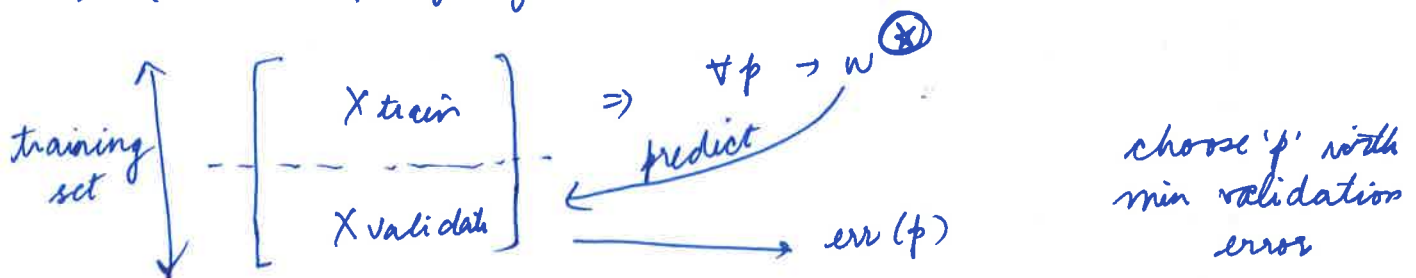
simple model, training error  $\uparrow$   
 generalization  $\rightarrow$  good  
 test error  $\uparrow$  }  $\rightarrow$  underfitting

$\therefore$  we have so much data, our model will not overfit, even the very complex ones (deep neural n/w's for instance have become popular now because we have a lot of data now).

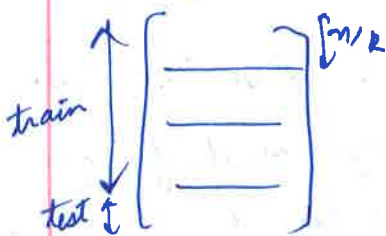
<< REFER or READ : VC-Dimension, Rademacher complexity >>

CROSS VALIDATION - choosing 'p':

split your training set into  $X_{train}$  &  $X_{validate}$  &  $X_{validate}$  will be used as a proxy for the 'test' set.



## k-fold cross validation



⇒ do this  $k$  times. We will get ' $k$ ' errors. Average over all these ' $k$ ' errors  $\Rightarrow \frac{\text{sum of errors}}{k}$

This is to be done for all ' $p$ 's'

Choose the ' $p$ ' with the minimum average error.

$$p^* = \text{argmin} (\text{all avg errors})$$

Once you get  $p^* \rightarrow$  train  $\rightarrow$  get ' $w$ '

## Regularization

$$\min_w \|Xw - y\|^2 + \underbrace{\frac{\lambda}{2} \|w\|^2}_{L_2 \text{ regularization}}$$

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

improves  
conditioning  
number.

$$\lambda = 1e^{-2}$$

condition number tells  
you how easy it is to  
invert a matrix

- unique solution.