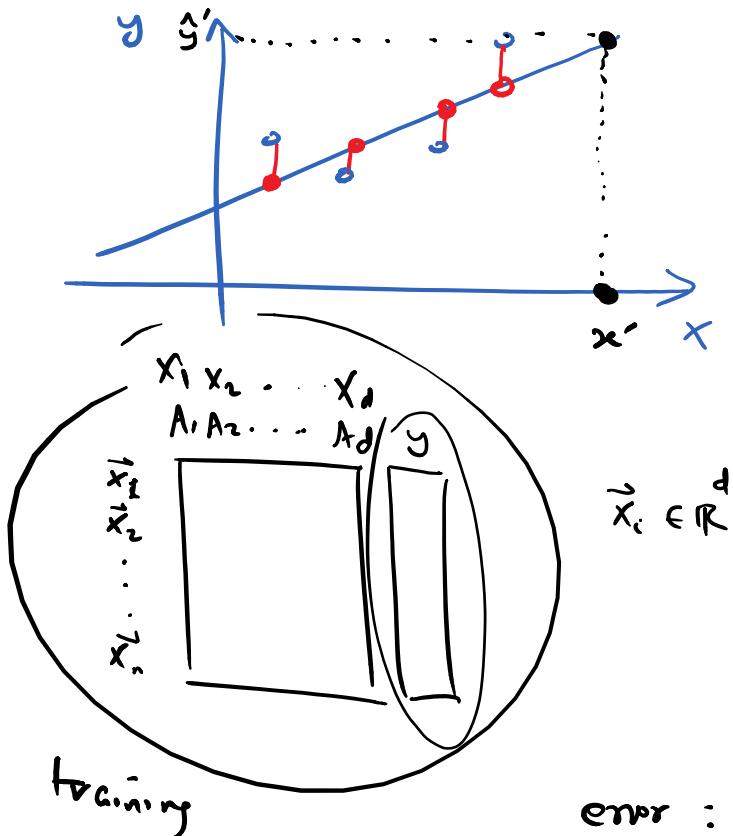
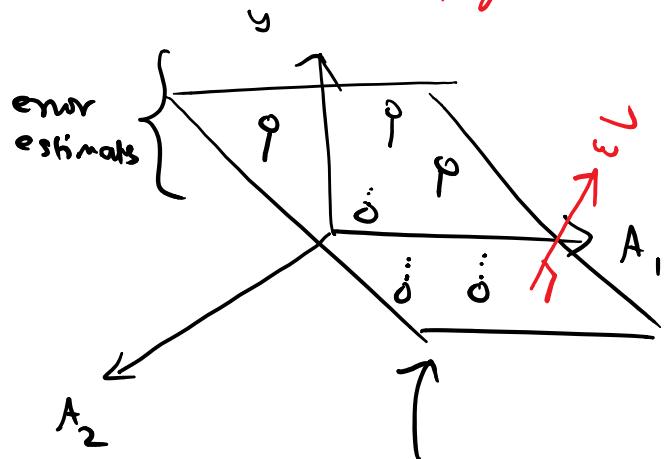


## Linear Regression



dependent var.  
Independent Var  
 $\hat{y} = w x + b$  |  
equation of a line  
slope  
bias/offset



error :  $\sum_{i=1}^n (y_i - \hat{y}_i)^2$

$$\hat{y}_i = \sum_{j=1}^d w_j x_{ij} + w_0$$

bias       $w_0 = b$

$$\vec{x}_i = \begin{pmatrix} 1 \\ x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{pmatrix} \quad \vec{w} = \begin{pmatrix} w_0 = b \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{pmatrix} \quad \left. \right\} \mathbb{R}^{d+1} \text{ space}$$

$$\hat{y}_i = \vec{w}^T \vec{x}_i$$

$$E = \sum_{i=1}^n (y_i - \vec{w}^T \vec{x}_i)^2$$

n diff components

find the best  $\vec{w}$  to minimize the squared errors

*n* diff Components

$$\vec{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \quad X = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \quad \vec{\omega}$$

$$E = \|\vec{y} - X\vec{\omega}\|^2$$

$$E = \|\vec{y} - \vec{g}\|^2$$

$$\vec{g} = \begin{pmatrix} g_1 \\ g_2 \\ \vdots \\ g_n \end{pmatrix}$$

$$E = (\vec{y} - X\vec{\omega})^T (\vec{y} - X\vec{\omega})$$

$$\frac{\partial E}{\partial \vec{\omega}} \Rightarrow (X^T X) \vec{\omega} = X^T \vec{y}$$

$d+1$  dimensions

$$\vec{\omega} = (X^T X)^{-1} X^T \vec{y}$$

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

$d \times 1$

$$X \in \mathbb{R}^{n \times d}$$

$$X = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}$$

$$\hat{y} = X\vec{\omega}$$

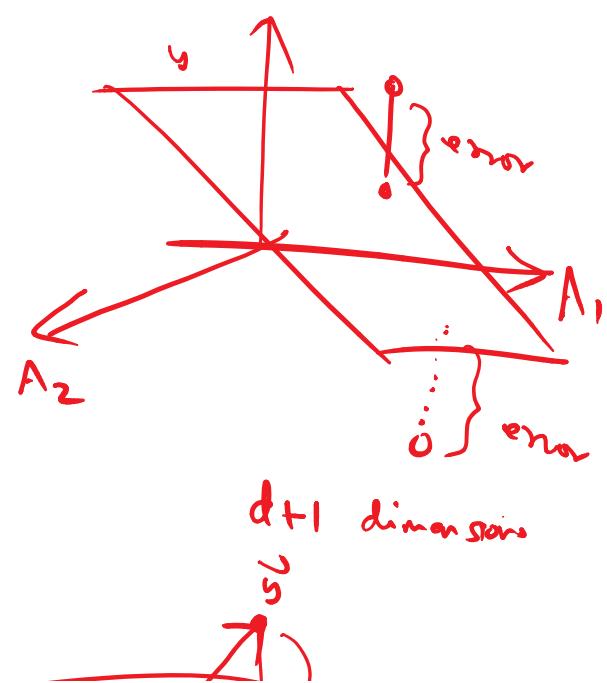
$$= X((X^T X)^{-1} X^T \vec{y})$$

$$\hat{y} = (X(X^T X)^{-1}) X^T \vec{y}$$

Predicted  
values

HGT matrix  
Projection matrix  
 $H$

true  
values

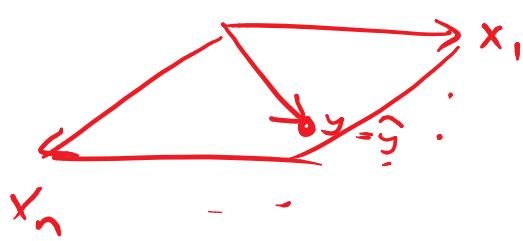
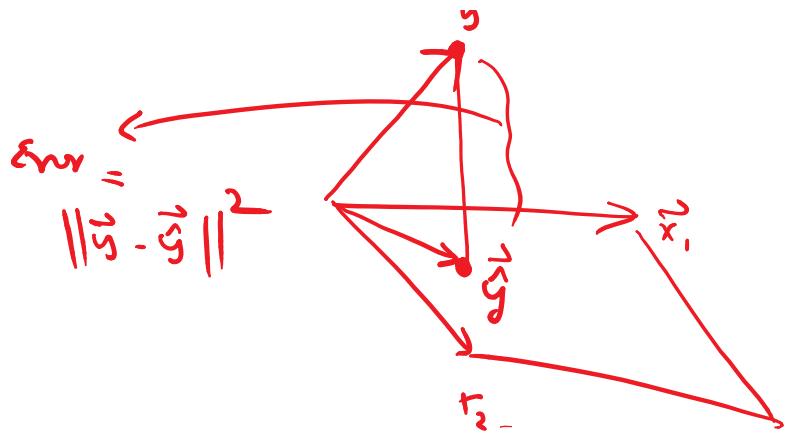


H

$y \in \mathbb{R}^n$

$A_1, A_2, \dots, A_d \in \mathbb{R}^n$

$\vdots$



$$y = Xw$$

$y$  could be expressed as  
a linear combination of  
attributes

$$y = A_1 w_1 + A_2 w_2 + \dots + A_d w_d$$

$\vec{y}$  is the  
orthogonal projection  
of  $\vec{y}$  onto  
 $\{A_1, A_2, \dots, A_d\}$

$b_{i,j}, w_1, w_2, \dots, w_j$

$\vec{y}$  has been projected onto the space  
spanned by the columns of  $X$   
attributes

the projected point  $\vec{y}$  is the  
best we can do

$$\vec{y} = Xw$$

$$\sum_{i=1}^n \underbrace{x_i^T w}_{\sum_{j=1}^d x_{ij} w_j}$$

$$\sum_{i=1}^n \left( \sum_{j=1}^d x_{ij} w_j \right)$$

↑  
weights

## Regularized linear regression

$w_i$  should not be too large

$$\sum_{i=1}^d w_i^2 = \|w\|_2^2$$

### Regularized objective

$$\min_w E = \underbrace{\|\vec{y} - \vec{y}\|_2^2}_{\text{error}} + C \underbrace{\|w\|_2^2}_{\substack{\text{size/norm of} \\ \text{weight vector}}} \quad C \geq 0$$

Small ones with small weights

$$\begin{aligned} E &= (\vec{y} - X\vec{w})^T (\vec{y} - X\vec{w}) + \vec{w}^T (C \cdot I) \vec{w} \\ &= \vec{y}^T \vec{y} - 2(\vec{y}^T X)\vec{w} + \vec{w}^T (X^T X) \vec{w} + \vec{w}^T (C \cdot I) \vec{w} \end{aligned}$$

$$E = \vec{y}^T \vec{y} - 2(\vec{y}^T X)\vec{w} + \vec{w}^T (X^T X + C I) \vec{w}$$

$$\frac{\partial E}{\partial \vec{w}} = 0 \Rightarrow \cancel{\vec{w}^T (X^T X + C I)} = \cancel{2 X^T \vec{y}}$$

ridge regression

$$\hat{\omega} = (X^T X + cI)^{-1} X^T y$$

$$X^T X$$

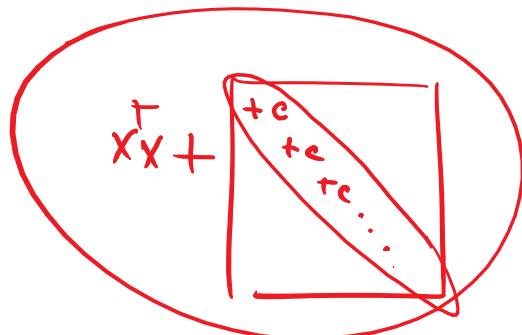
Covariance

$d \times d$

Symmetric

PSD

$$\lambda_i > 0$$



Inverse always exists

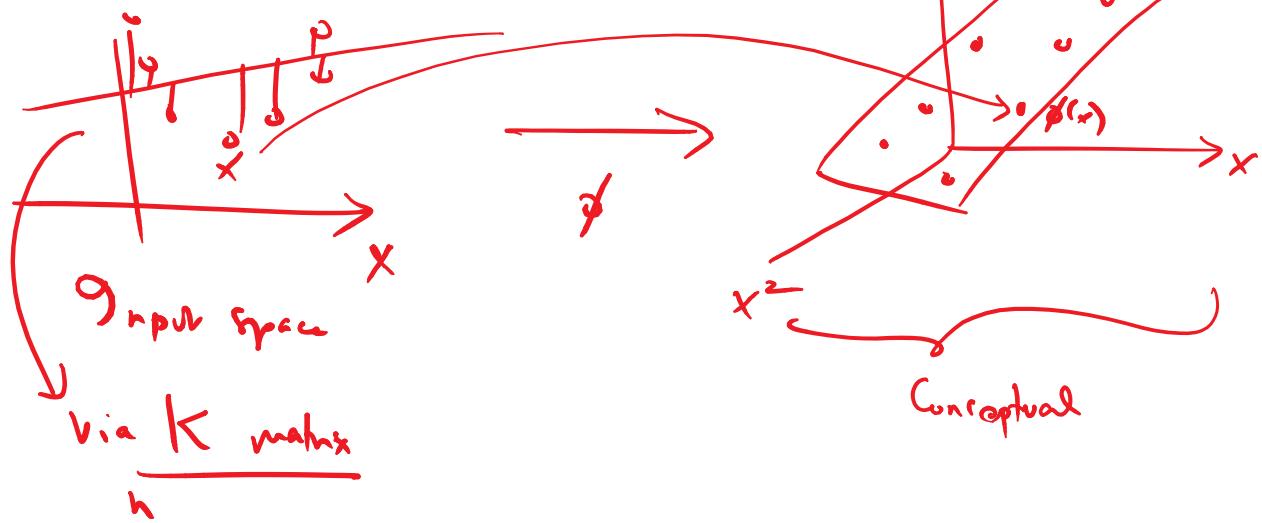
• if  $\lambda_i$  is eigenvalue of  $X^T X$

$\lambda_i + c$  is eigenvalue of  $X^T X + cI$

$$(X^T X + cI)$$

$$\lambda'_i \geq c$$

## Kernel Regression



$$n \quad K$$

## Kernel Ridge Regression

$$(cI)$$

$$\hat{\omega} = \sum \alpha_i \phi(x_i) \in \mathbb{C}$$

$$\vec{\omega} = \sum \alpha_i x_i$$

$$\vec{\omega} = \sum \alpha_i \underline{\phi(x_i)}$$

linear case / kernel

$$\phi(x_i) = x_i$$

we cannot  
reconstruct  
the weight  
vector

$$\vec{\alpha} = \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \end{pmatrix}$$

$$\hat{y}_j = \phi(x_j) \vec{\omega}$$

$$= \sum_{i=1}^n \alpha_i \phi(x_j)^T \phi(x_i)$$

$$= \sum_{i=1}^n \alpha_i k(j,i)$$

$$\hat{y}_j = \vec{\alpha} \cdot \vec{k}_j$$

alpha vector      j-th row of K

linear kernel:  $\phi(\vec{x}_i) = \vec{x}_i$

ridge objective  $E = \|y - X\omega\|^2 + c\omega^T\omega$

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

$$(X^T X + cI)\vec{\omega} = X^T y$$

$$c\vec{\omega} = X^T y - X^T X\omega$$

$$\vec{\omega} = X^T \frac{y - X\vec{\omega}}{c}$$

$$\vec{\omega} = X^T \vec{\alpha} = \sum_{i=1}^n x_i \alpha_i$$

$$\vec{\alpha} = \frac{\vec{y} - X\vec{\omega}}{C} \Rightarrow C\vec{\alpha} = \vec{y} - X\vec{\omega}$$

$$X \in \mathbb{R}^{n \times d}$$

$$\Rightarrow C\vec{\alpha} = \vec{y} - (X^T X)\vec{\alpha}$$

$$C^T \vec{\alpha} + (X^T X)\vec{\alpha} = \vec{y}$$

$$(X^T X + C I)\vec{\alpha} = \vec{y}$$

$$\vec{\alpha} = (X^T X + C I)^{-1} \vec{y}$$

$$\vec{\alpha} = (K + C I)^{-1} \vec{y}$$

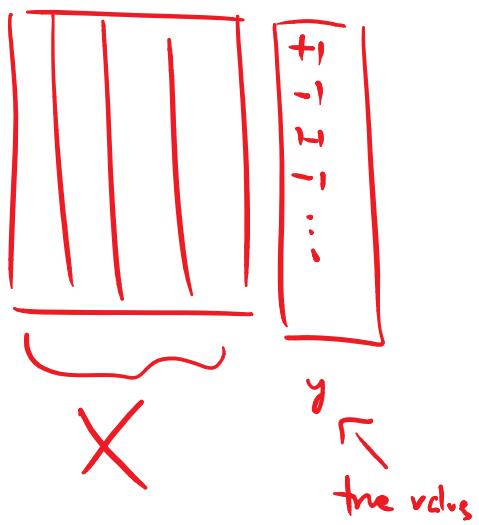
kernel ridge

$$\vec{\omega} = \sum_{i=1}^n \alpha_i \phi(x_i)$$

$d \times d \rightarrow n \times n$  space

$$\hat{y} = \vec{\omega}^T \phi(z) = \sum_{i=1}^n \alpha_i k(i, z) \text{ for test case } z.$$

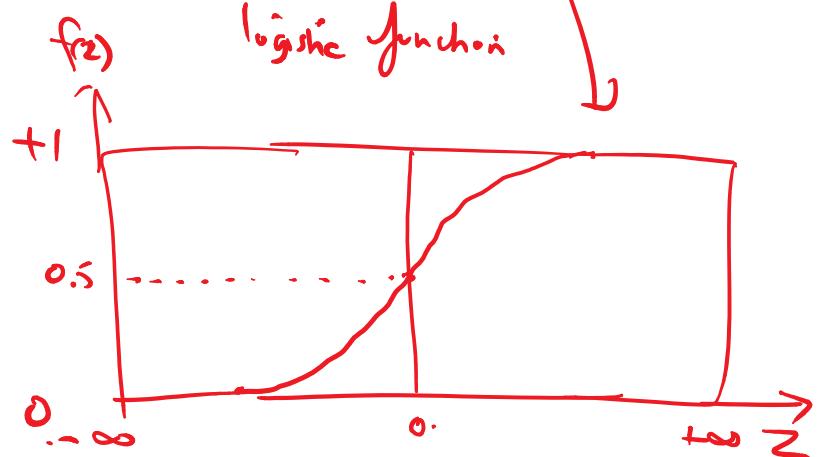
## Logistic Regression



$$z = \sum_{i=1}^d w_i x_i + w_0$$

$$f(z) = \frac{1}{1+e^{-z}}$$

squash the value  $z$  in the range  $(0, 1)$



find the best  $w$   
with logistic function

classes :  $\{+1, -1\}$

$$\begin{aligned} +1 & \quad \hat{y} = P(y=1|x) = f(z) \\ -1 & \quad 1 - \hat{y} = 1 - P(y=1|x) = 1 - f(z) \end{aligned}$$

$$f(-z) = 1 - f(z) \quad \leftarrow \text{property of logistic function}$$

$$1 - f(z)$$

$$1 - \frac{1}{1+e^z} = \frac{1-e^{-z}}{1+e^z} = \frac{e^{-z}/e^z}{(1+e^z)/e^z} = \frac{1}{e^z + 1}$$

$$f(y \cdot z)$$

$$f(z) = \frac{1}{1+e^{-z}} = \frac{1}{1+e^z} \quad ||| \quad +1 \equiv f(z) \\ -1 \equiv f(-z) = 1-f(z)$$

likelihood function for all data points-

$$L = \prod_{i=1}^n f(y_i z_i) = \prod_{i=1}^n f(y_i w^T x_i)$$

Maximize  $\ln L = \sum_{i=1}^n \ln f(y_i w^T x_i)$

finding the parameters  $\tilde{w}$  so that we maximize  
the (log of the likelihood) (MLF)

How probable is  $X, y$  given the parameters  $(\tilde{w})$

likelihood

$$w^T = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{pmatrix}$$

$$L = \sum_{i=1}^n \ln f(y_i w^T x_i)$$

$$f(z_i) = \frac{1}{1+e^{-z_i}}$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial f(z_i)} \cdot \frac{\partial f(z_i)}{\partial z_i} \cdot \frac{\partial z_i}{\partial w} \\ = \frac{1}{f(z_i)} \cdot [f(z_i)(1-f(z_i))] \cdot y_i x_i$$

$\boxed{\nabla_w = f(-z_i) y_i x_i}$

where  $z_i = y_i w^T x_i$

Gradient descent optimization for  $\vec{w}$ , per component  $w_j$

$\vec{w} = (0 \dots 0)$

repeat

For all points  $i$  in random order

$$z_i = y_i w^T x_i$$

$$\nabla_i = f(-z_i) \cdot y_i x_i$$

$$w = w + \eta \nabla_i$$

until  $\|w - w_{\text{prev}}\|^2 \leq \varepsilon$

