

UCB - CS189  
Introduction to Machine Learning  
Fall 2015

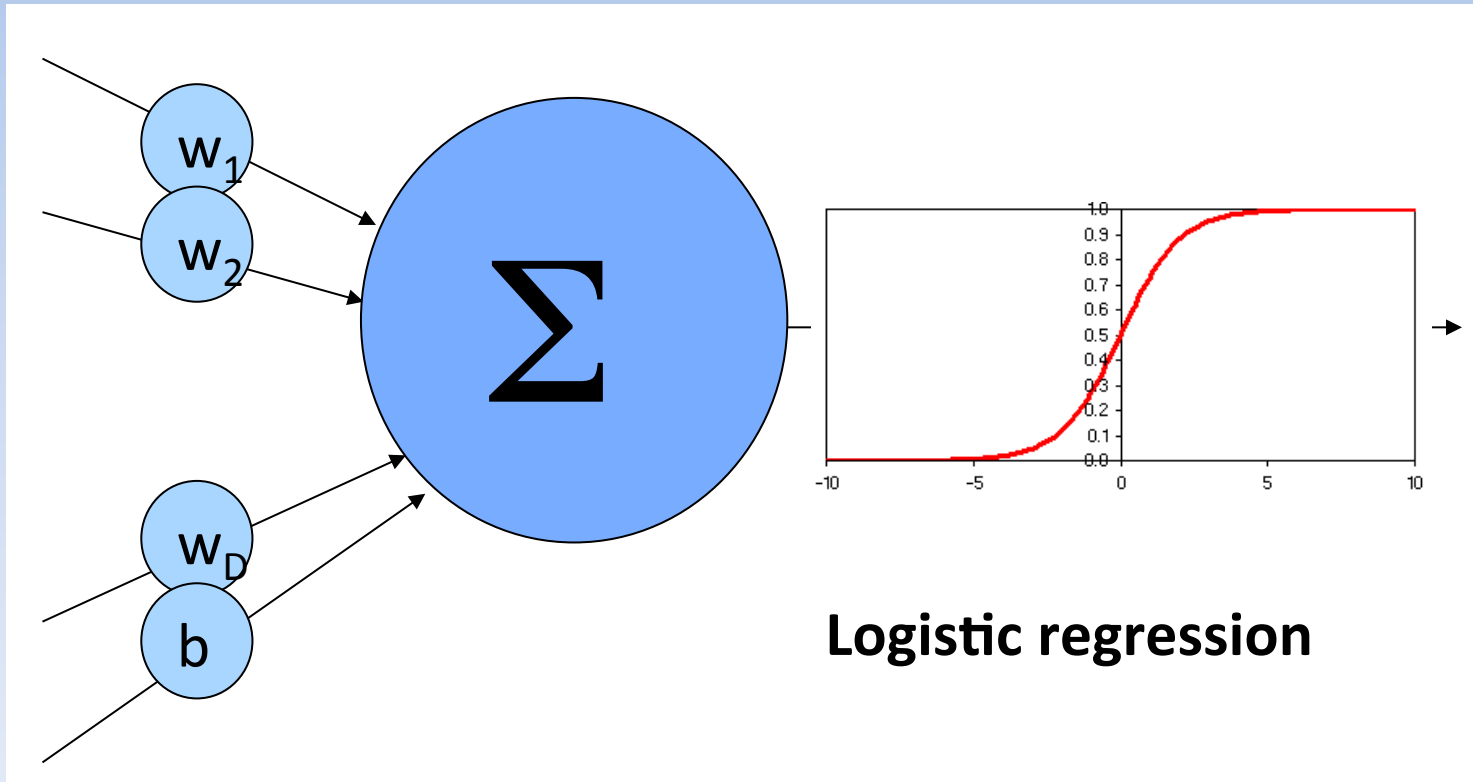
Lecture 7: Ridge regression

Isabelle Guyon  
ChaLearn

Come to my office hours...

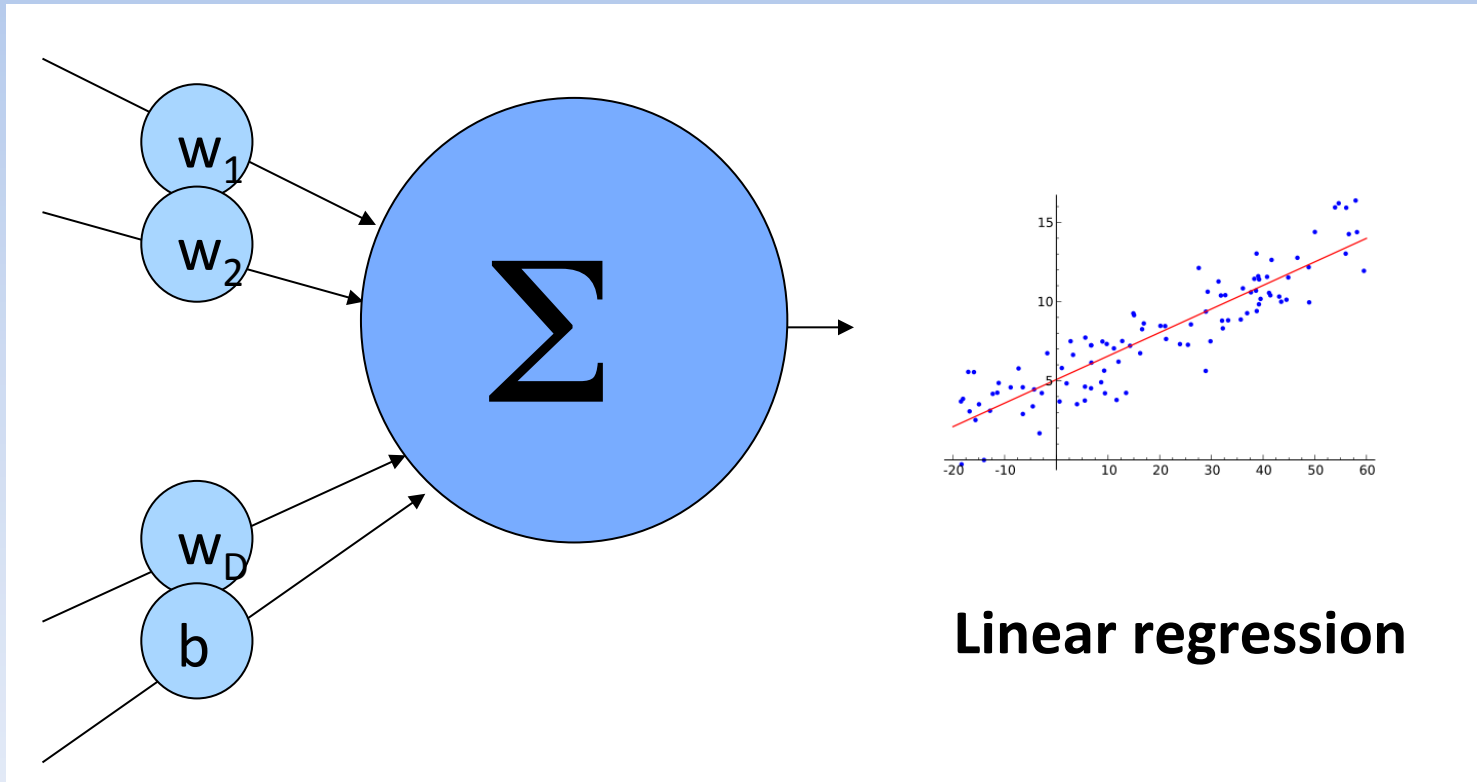
Wed 2:30-4:30 Soda 329

Last time



Come to my office hours...  
Wed 2:30-4:30 Soda 329

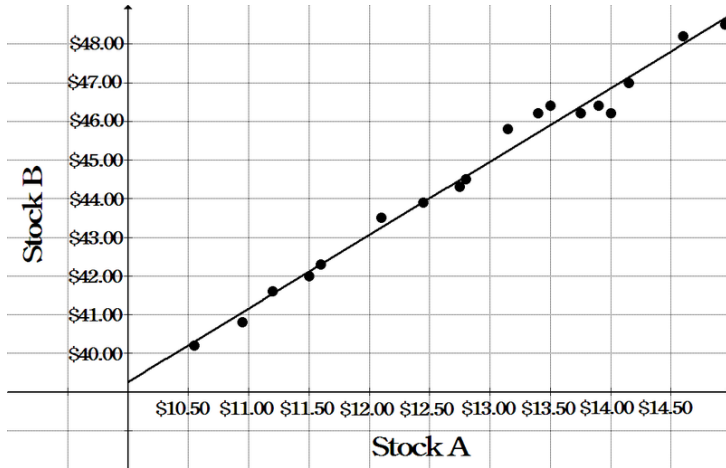
## Today



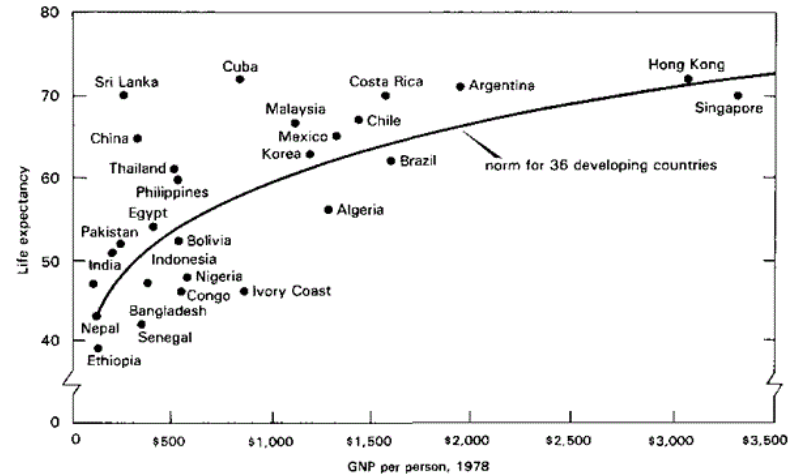
# Math prerequisites

- Vectors and matrices
- Matrix multiplication
- Matrix inverse, determinant
- Matrix diagonalization, eigen vectors, eigen values, rank of a matrix
- Pseudo-inverse

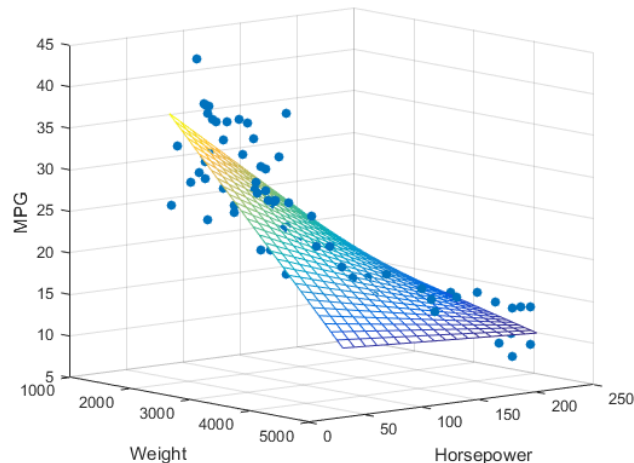
# Uses of regression



Economy and Finance



Epidemiology and medicine

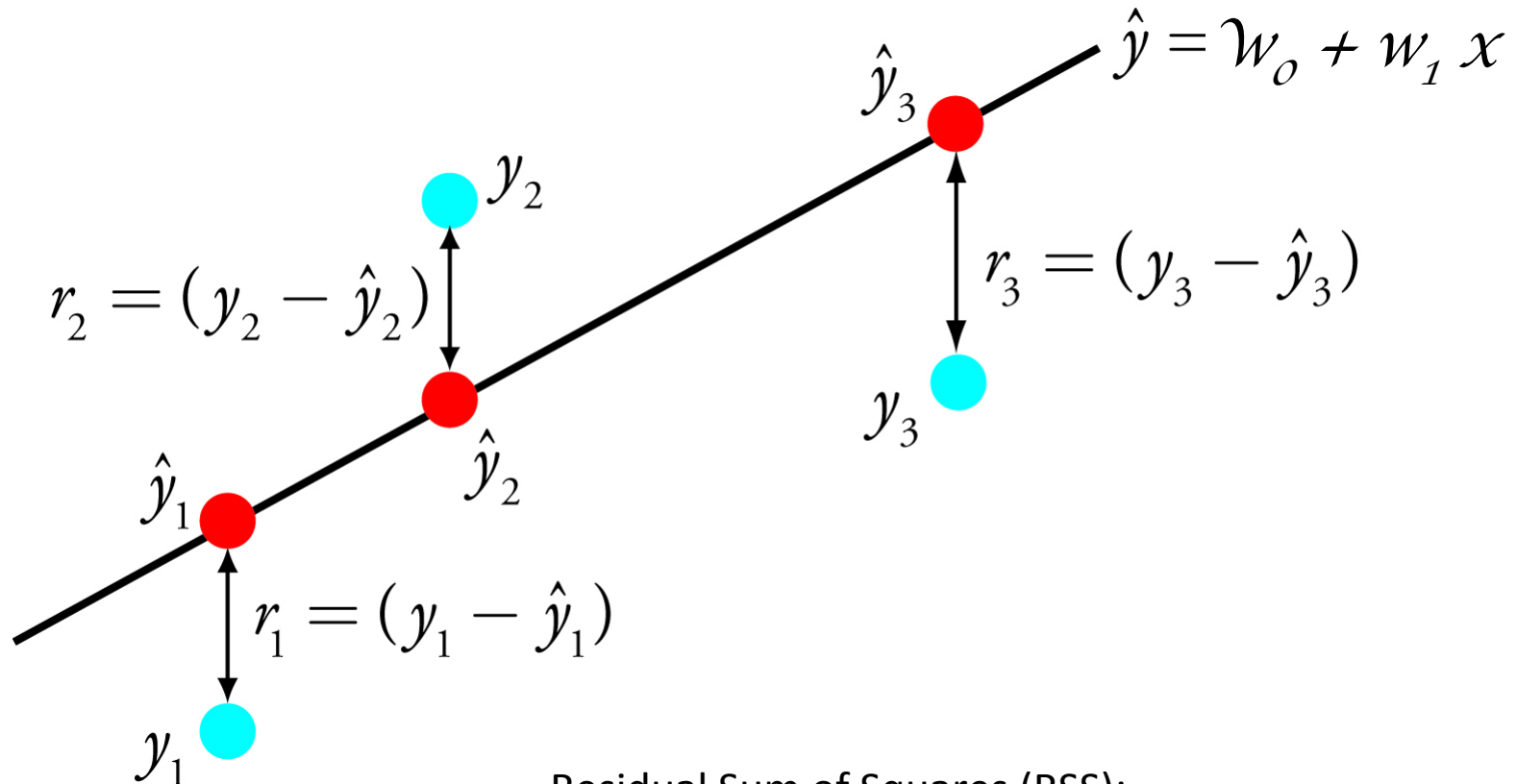


Engineering



Apparent age estimation: <http://gesture.chalearn.org>

# Least square regression

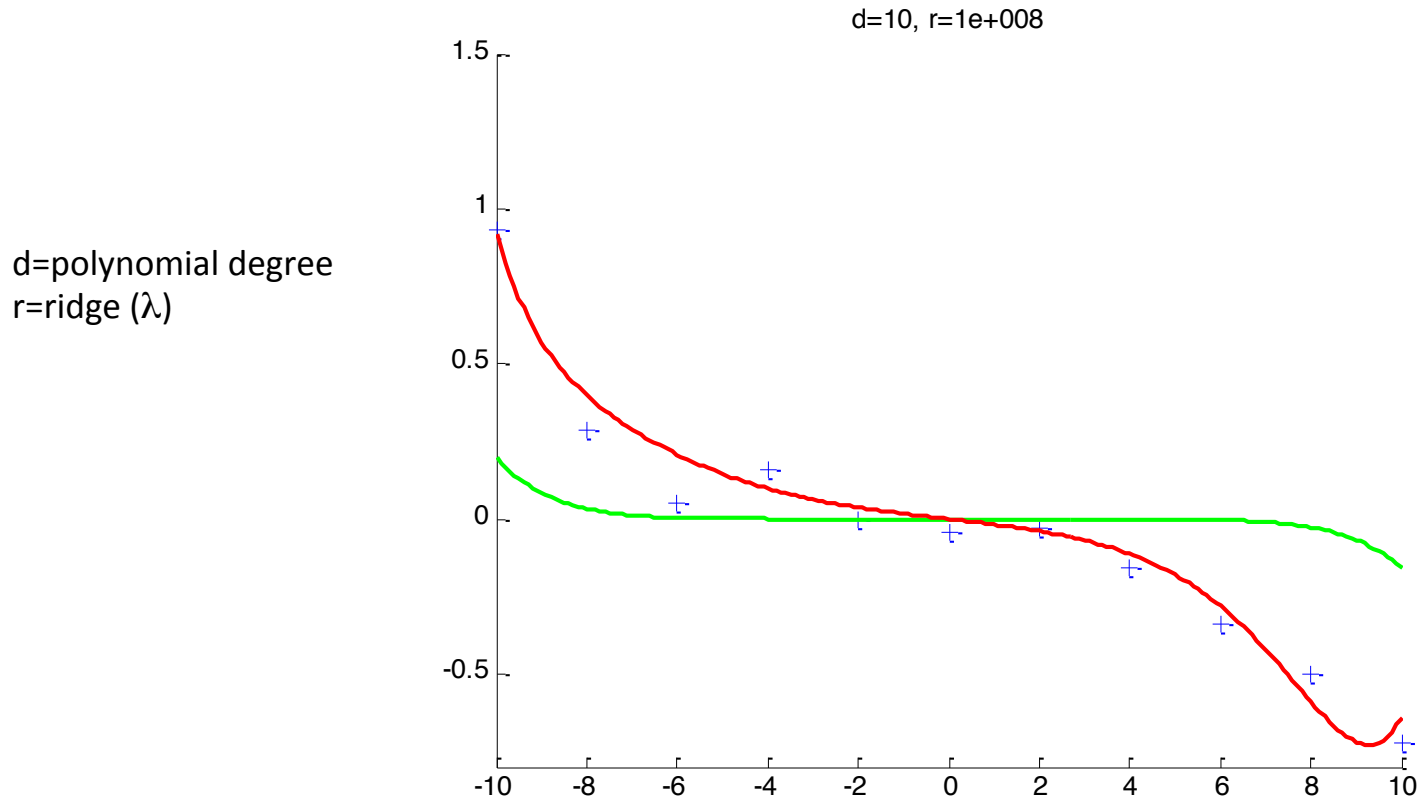


Residual Sum of Squares (RSS):

$$R[\mathbf{w}] = \sum_k (w_0 + w_1 x_k - y_k)^2 = \sum_k (\mathbf{w} \cdot \mathbf{x}_k - y_k)^2$$

Picture: <http://chemwiki.ucdavis.edu/>

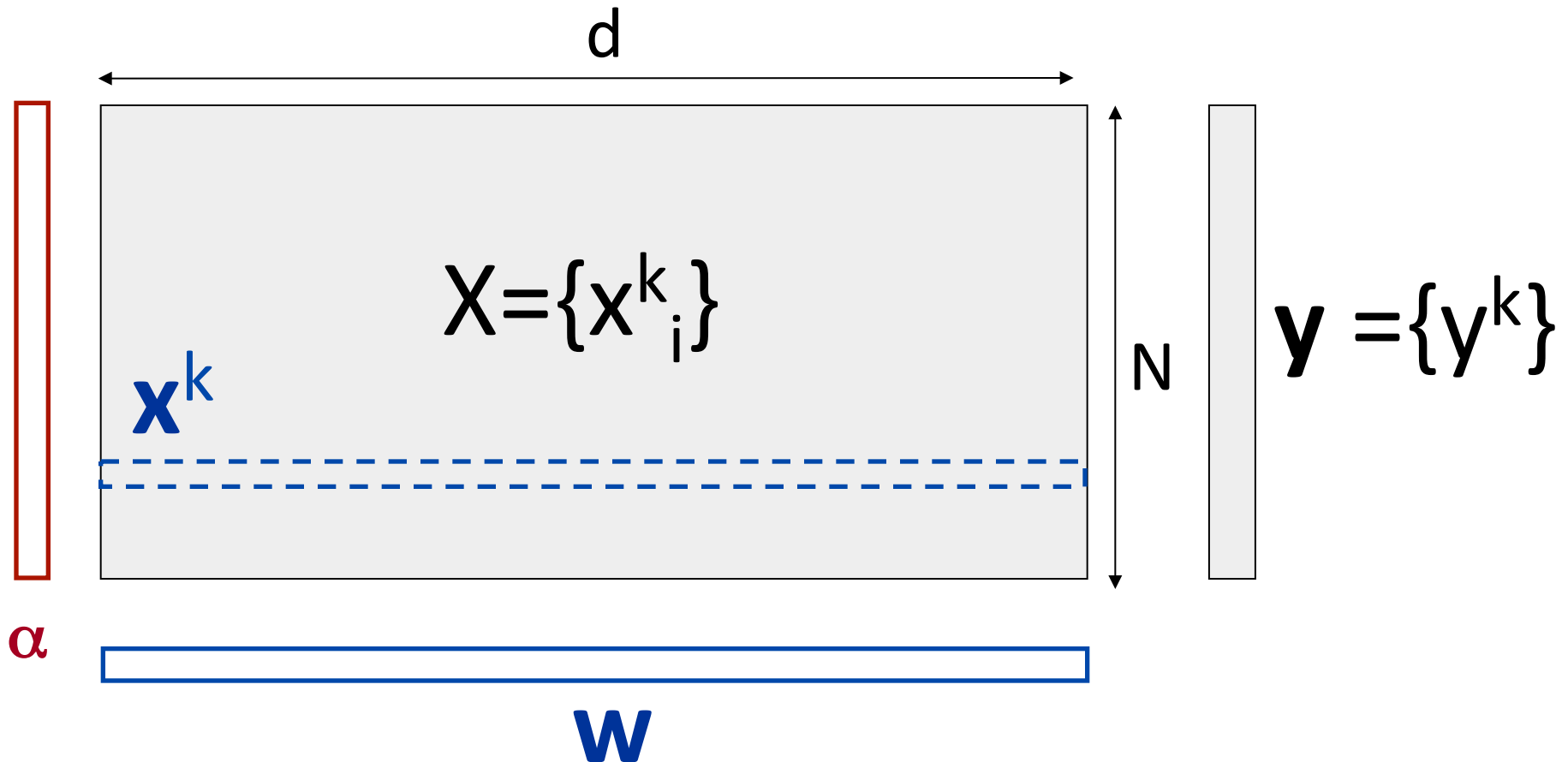
# Polynomial Ridge Regression



$$R[\mathbf{w}] = \sum_k (\mathbf{w} \cdot \Phi(\mathbf{x}_k) - y_k)^2 + \lambda \|\mathbf{w}\|^2$$

$$\Phi(\mathbf{x}) = [1, x, x^2, x^3, \dots, x^d]$$

# Conventions





# Matrix Notations

$$w_i = \sum_k y^k x_{i,k}$$

$$\mathbf{w} = \mathbf{y}^\top \mathbf{X}$$

$$\mathbf{w}^\top = \mathbf{X}^\top \mathbf{y}$$

$$w_i = \sum_k \alpha^k x_{i,k}$$

$$\mathbf{w} = \boldsymbol{\alpha}^\top \mathbf{X}$$

$$\mathbf{w}^\top = \mathbf{X}^\top \boldsymbol{\alpha}$$

$$(1,d) = (1,N)(N,d)$$

$$(d,1) = (d,N)(N,1)$$

$$f(\mathbf{x}) = \sum_i w_i x_i$$

$$f(\mathbf{x}) = \mathbf{x} \mathbf{w}^\top = \mathbf{w} \mathbf{x}^\top$$

$$(1,d)(d,1) \quad (1,d)(d,1)$$

# Linear Regression

- What we want:

$$\sum_i w_i x_i^k = y^k \quad \text{for all examples } k=1\dots m \quad (b=w_0)$$

or for classification,  $y^k=\pm 1$ ,  $\text{sign}(\sum_i w_i x_i^k) = y^k$

- Solve:

$$\mathbf{X}\mathbf{w}^T = \mathbf{y}$$

$$(N,d)(d,1)=(N,1)$$

# Regression: $N > d$

- Solve:

$$X \mathbf{w}^T = \mathbf{y}$$

$(N,d)(d,1) = (N,1)$

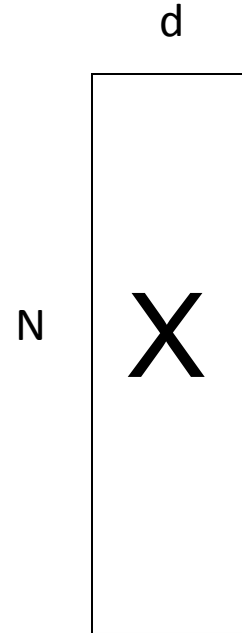
- Normal equations

$$X^T X \mathbf{w}^T = X^T \mathbf{y}$$

$(d,N)(N,d)(d,1) = (d,N)(N,1)$

- Solution:

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



$\text{rank}(X) \leq \min(d, N)$   
assume  $\text{rank}(X) = d$   
implies  $\text{rank}(X^T X) = d$   
 $X^T X$  is invertible

# Pseudo-Inverse

- Solution:

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

$(d,1)$        $(d,N)(N,d)$     $(d,N)(N,1)$        $(d,N)$  pseudo-inverse,  $\mathbf{X}^+ \mathbf{X} = \mathbf{I}$

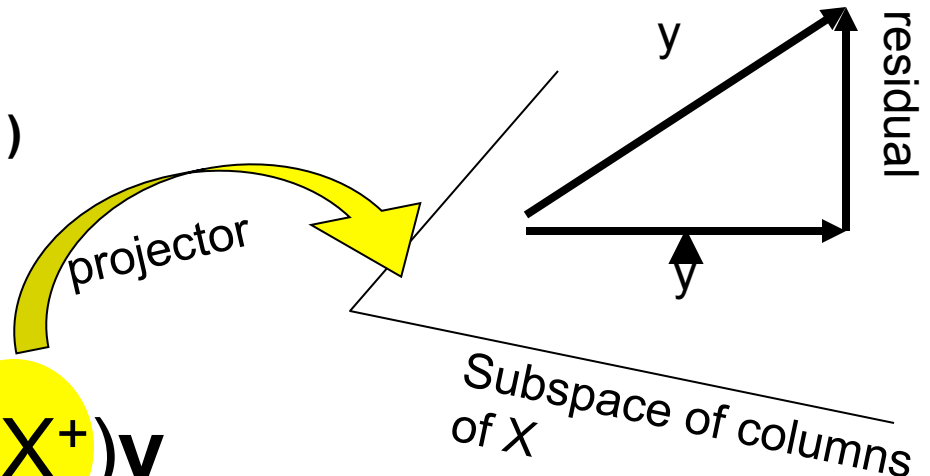
- Predictor:

$$\mathbf{f}(\mathbf{x}) = \mathbf{x} \mathbf{w}^T = \mathbf{x} \mathbf{X}^+ \mathbf{y}$$

$(1,1)$        $(1,d)(d,1)$        $(1,d)(d,N)(N,1)$

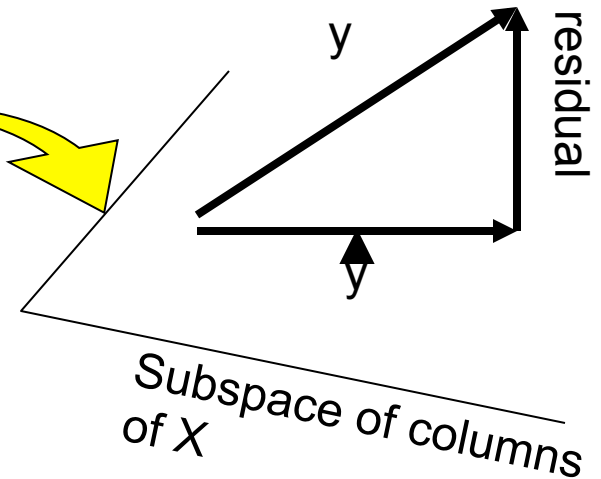
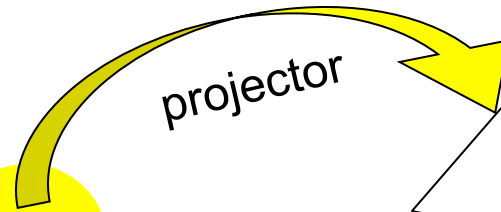
- Residual:

$$\mathbf{y} - \hat{\mathbf{y}} = \mathbf{y} - \mathbf{X} \mathbf{w}^T = (\mathbf{I} - \mathbf{X} \mathbf{X}^+) \mathbf{y}$$



# Least-Squares

$$\mathbf{y} - \hat{\mathbf{y}} = \mathbf{y} - \mathbf{X} \mathbf{w}^T = (\mathbf{I} - \mathbf{X}\mathbf{X}^+) \mathbf{y}$$



The pseudo-inverse solution is optimal in the least-square sense:

$$\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X} \mathbf{w}^T\|^2 = \|(\mathbf{I} - \mathbf{X}\mathbf{X}^+) \mathbf{y}\|^2$$

# Gradient Descent

- Square loss:

$$L_k = (\mathbf{x}^k \mathbf{w}^T - y^k)^2$$

Risk = Residual Sum of Squares (RSS):

$$\begin{aligned} R &= \sum_k (\mathbf{x}^k \mathbf{w}^T - y^k)^2 \\ &= \| X\mathbf{w}^T - \mathbf{y} \|^2 \\ &= \mathbf{w}X^T X \mathbf{w}^T - 2\mathbf{w}X^T \mathbf{y} + \mathbf{y}^T \mathbf{y} \end{aligned}$$

- Gradient:

$$\nabla_{\mathbf{w}} R = 2 (X^T X \mathbf{w}^T - X^T \mathbf{y})$$

# Normal Equations

- At the optimum:

$$\nabla_{\mathbf{w}} R = \mathbf{0}$$

$$2 (X^T X \mathbf{w}^T - X^T \mathbf{y}) = \mathbf{0}$$

- Normal equations (again):

$$X^T X \mathbf{w}^T = X^T \mathbf{y}$$

Solve by inverting  $X^T X$ , if regular.

- What if  $X^T X$ , is singular?

# Regularization

- Normal equations:
- Normal equations:

$$\begin{pmatrix} X^T X & X^T y \\ X^T y & y^T y \end{pmatrix} \begin{pmatrix} w \\ b \end{pmatrix} = \begin{pmatrix} X^T y \\ y^T y \end{pmatrix}$$

)(

$$X^T X \begin{pmatrix} w \\ b \end{pmatrix} = X^T y$$

- $X^T X$  is singular
- Solution:  $(X^T X + \lambda I)^{-1} X^T y$   $\lambda > 0$



# Why it works

- Diagonalization:

$$X^T X = U D U^T$$

- Disambiguation:
  - U orthogonal matrix of eigenvectors ( $U U^T = I$ )

D diagonal matrix of eigenvalues

Singularity: some eigenvalues are zero.

$$X^T X = U D U^T$$

U orthogonal matrix of eigenvectors ( $U U^T = I$ )

no more zero eigenvalue.

D diagonal matrix of eigenvalues

Singularity: some eigenvalues are zero.

# Penalized Risk

$$\sum_k (\mathbf{x}^k \mathbf{w}^T - y^k)^2$$

$$= \| \mathbf{X} \mathbf{w}^T - \mathbf{y} \|^2$$

- Add “

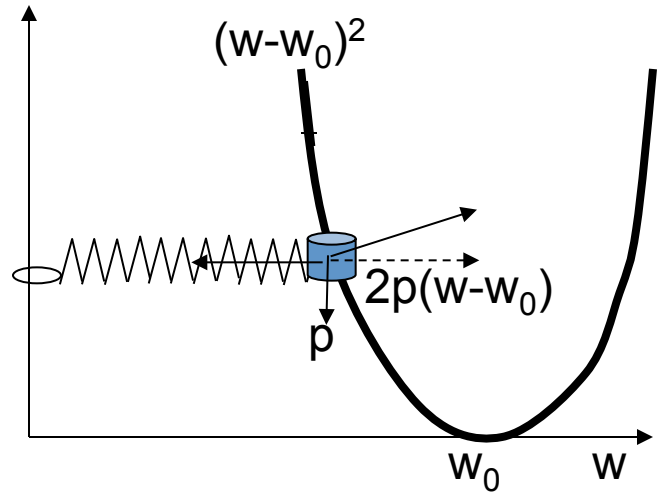
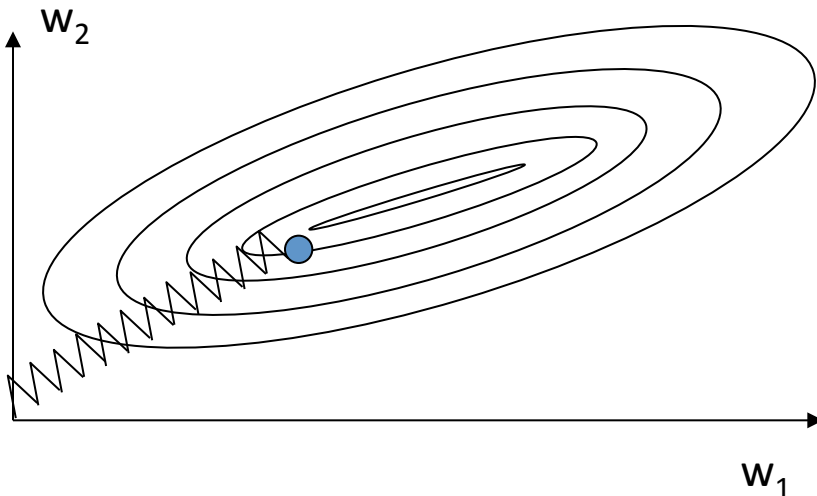
regularizer  $\lambda \|\mathbf{w}\|^2$   $\lambda > 0$

- Objective:

$$\nabla_{\mathbf{w}} \left( \sum_k (\mathbf{x}^k \mathbf{w}^T - y^k)^2 + \lambda \|\mathbf{w}\|^2 \right)$$

# Mechanical Interpretation

- Quadratic form:
  - Quadratic form:  $R = \frac{1}{2} \lambda \| \mathbf{w} - \mathbf{y} \|^2 + \lambda \| \mathbf{w} \|^2$
  - One dimension:
- $$R = p (w - w_0)^2 + \lambda w^2$$
- Two dimensions:  $R = \frac{1}{2} \lambda \| \mathbf{w} - \mathbf{y} \|^2 + \lambda \| \mathbf{w} \|^2$



# Principal Component Analysis

$$\mathbf{x}^k \begin{matrix} X \\ \hline (N, d) \end{matrix} \begin{matrix} U(d, d') \\ \hline \end{matrix} = \begin{matrix} X' \\ \hline (N, d') \end{matrix} \mathbf{f}_j \mathbf{x}'^k$$

$\mathbf{u}_j$

- Problem: Construct features that are linear combinations of the original features, such that the reconstructed patterns are as close as possible to the original in the least square sense.
- $\mathbf{f}_j = X \mathbf{u}_j$  linear combinations of columns of  $X$
- Problem: Construct features that are linear combinations of the
- $\mathbf{x}''^k = \mathbf{x}'^k U^T = \sum_j \mathbf{x}'^k_j \mathbf{u}_j$

$$\mathbf{x}'^k \begin{matrix} X' \\ \hline (N, d') \end{matrix} \begin{matrix} U^T \\ \hline (d', d) \end{matrix} = \begin{matrix} X'' \\ \hline (N, d) \end{matrix} \mathbf{x}''^k$$

# PCA Solution

- $X' = X U$
- $X'' = X' U^T$
- $X' = X U$
- $\min_U \|X U^T X U U^T\|^2$
- Can be brought back to solving and eigenvalue
- Compare:  
 Regularization  $X^T X + \lambda I = U(D + \lambda I)U^T$   
 $U U^T X U U^T$   
 PCA: Remove the dimensions with smallest

# Kernel “Trick” ( $N < d$ )

- Solve:  $X\mathbf{w}^T = \mathbf{y}$
- Assume:  $\mathbf{w} = \sum_k \alpha^k \mathbf{x}^k$   $\mathbf{w} = \alpha^T X$
- Solve instead:  $X X^T \alpha = \mathbf{y}$ 

(1,d) (1,N) (N,d)

( ) (d,N) (N,1) = (N,1)

Full rank (N,N) matrix
- Solution:  $\alpha = (X X^T)^{-1} \mathbf{y}$ 

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

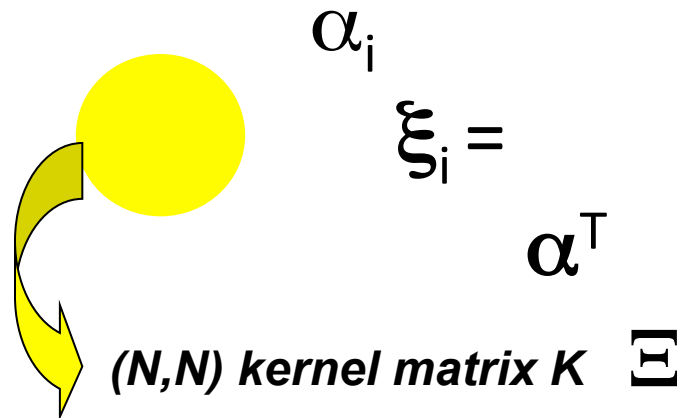
$X^+$

# Kernel Ridge Regression

- $\Xi = \Phi(X)$
- ~~Solve:~~  $\Xi \mathbf{w} = \mathbf{y}$

$$\mathbf{M} = \sum_i \Xi \mathbf{w}^T = \mathbf{y}$$

$$= \mathbf{y} \boldsymbol{\alpha}^T$$



- ~~Regularization:~~ Replace  $K$  by  $K + \lambda \mathbf{I}$

# Regularization and PI

- Case  $N > d$  and  $\text{rank}(X^T X) = d$

$$X^+ = (X^T X)^{-1} X^T$$

- Case  $N < d$  and  $\text{rank}(X^T X) = N$

- Case  $N > d$  and  $\text{rank}(X^T X) = d$

- Either case  $X^+ = (X^T X)^{-1} X^T$

- Case  $N < d$  and  $\text{rank}(X^T X) = N$

- Either case  $X^+ = \lim_{\lambda \rightarrow 0} (X^T X + \lambda I)^{-1} X^T$

$$\lambda I$$



# Summary

- Least square regression for models linear in their parameters can be achieved with:
  - Stochastic gradient (suited for big data, see next lesson)
  - Pseudo-inverse (requires matrix inversion):
    - If  $d < N$ , invert  $X^T X$ , a  $(d, d)$  matrix.
    - If  $N < d$ , invert  $X X^T$ , a  $(N, N)$  matrix.
- Kernelization is easy  $\Xi = \Phi(X)$ 
  - Invert  $\Xi \Xi^T$ , the kernel matrix  $K$ , a  $(N, N)$  matrix.
- Shrinkage for robustness:
  - Add a “ridge” = a small positive value to the diagonal.

Come to my office hours...  
Wed 2:30-4:30 Soda 329

## Next time

### Kernel machines

#### PARAMETRIC (Perceptrons)

$$f(\mathbf{x}) = \mathbf{w} \cdot \Phi(\mathbf{x})$$

$$\mathbf{w} = \sum_k \alpha_k \Phi(\mathbf{x}^k)$$

(Large margin) Perceptron

$$\begin{aligned} \Delta \mathbf{w} &\sim y_k \Phi(\mathbf{x}^k) \quad \text{if } y_k f(\mathbf{x}^k) < 1 \\ &\sim \mathbf{1}(1 - z_k) y_k \Phi(\mathbf{x}^k) \quad z_k = y_k f(\mathbf{x}^k) \end{aligned}$$

(Rosenblatt 1958)

Logistic regression

$$\Delta \mathbf{w} \sim S(-z_k) y_k \Phi(\mathbf{x}^k)$$

(Cox 1958)

LMS regression or classification

$$\Delta \mathbf{w} \sim (y_k - f(\mathbf{x}^k)) \Phi(\mathbf{x}^k) \sim (1 - z_k) y_k \Phi(\mathbf{x}^k)$$

(Widrow-Hoff, 1960)

#### NON PARAMETRIC (Kernel machines)

$$f(\mathbf{x}) = \sum_k \alpha_k k(\mathbf{x}^k, \mathbf{x})$$

$$k(\mathbf{x}^k, \mathbf{x}) = \Phi(\mathbf{x}^k) \cdot \Phi(\mathbf{x})$$

Potential Function algorithm

$$\begin{aligned} \Delta \alpha_k &\sim y_k \quad \text{if } y_k f(\mathbf{x}^k) < 1 \\ &\sim \mathbf{1}(1 - z_k) y_k \end{aligned}$$

(Aizerman et al 1964)

Dual logistic regression

$$\Delta \alpha_k \sim S(-z_k) y_k$$

Dual LMS

$$\Delta \alpha_k \sim (y_k - f(\mathbf{x}^k)) \sim (1 - z_k) y_k$$