

---

# Machine Learning

## CSE 6363 (Fall 2019)

### Lecture 3 Classification, Linear Regression

Dajiang Zhu, Ph.D.

Department of Computer Science and Engineering

---

*Slides of this course (CSE6363) courtesy: Dr. Heng Huang,  
Dr. Aarti Singh, Dr. Tibshirani*

# Supervised Learning

---

**Data:**  $D = \{d_1, d_2, \dots, d_n\}$     a set of  $n$  examples

$$d_i = \langle \mathbf{x}_i, y_i \rangle$$

$\mathbf{x}_i$  is input vector, and  $y$  is desired output (given by a teacher)

**Objective:** learn the mapping  $f : X \rightarrow Y$

$$\text{s.t. } y_i \approx f(x_i) \quad \text{for all } i = 1, \dots, n$$

**Two types of problems:**

- **Regression:**  $X$  discrete or continuous  $\rightarrow$   
 $Y$  is **continuous**
- **Classification:**  $X$  discrete or continuous  $\rightarrow$   
 $Y$  is **discrete**

# Discrete to Continuous Labels

---

- Classification



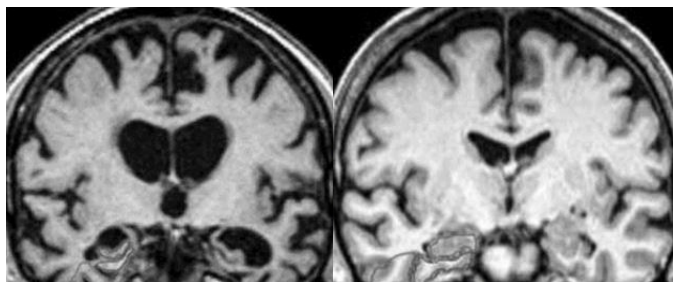
$X = \text{Document}$



Science?  
Sports?  
News?

...

$Y = \text{Topic}$



$X = \text{Brain image}$

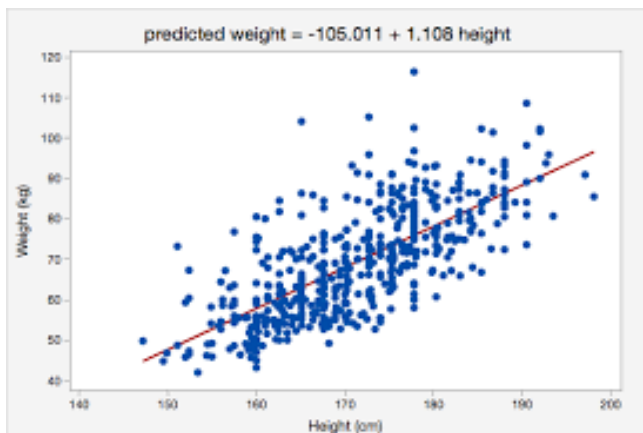


Alzheimer's Disease (AD) patient?  
Normal people?

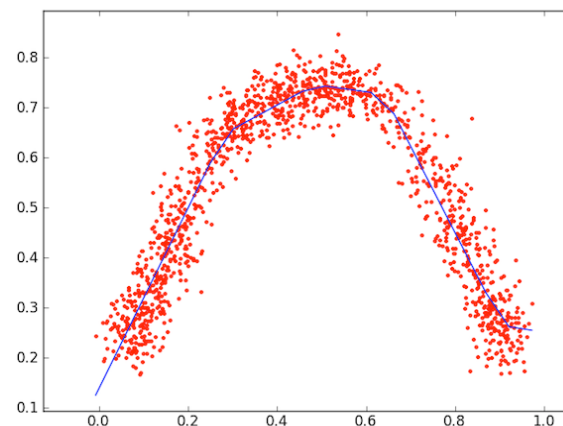
$Y = \text{Patient or not}$

# Discrete to Continuous Labels

- Regression

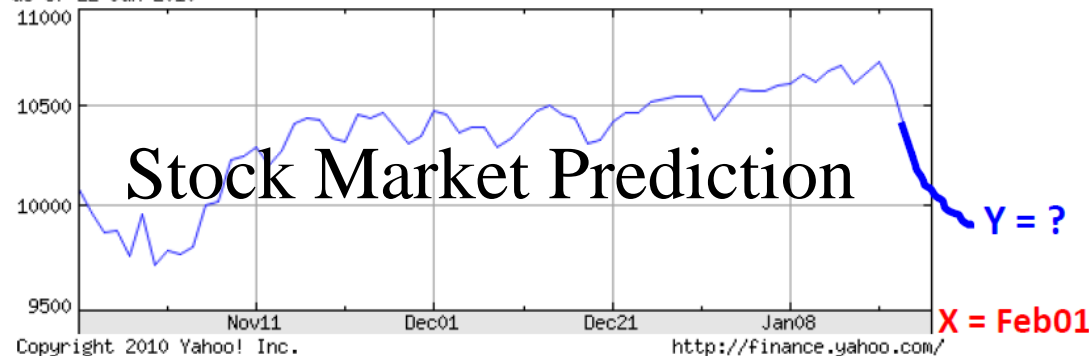


<https://onlinecourses.science.psu.edu/stat200/node/81>



<https://blog.biolab.si/tag/regression/>

DJ INDU AVERAGE (DOW JONES & CO  
as of 22-Jan-2010

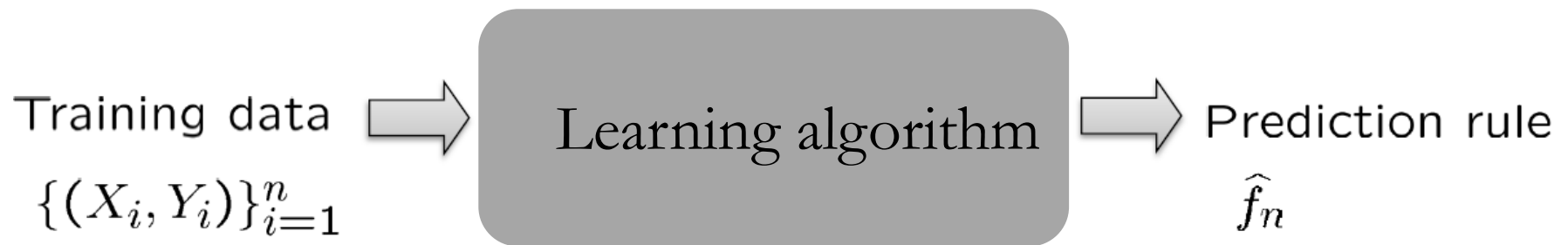


Copyright 2010 Yahoo! Inc.

<http://finance.yahoo.com/>

# Regression Algorithms

---



Linear Regression

Lasso, Ridge regression (Regularized Linear Regression)

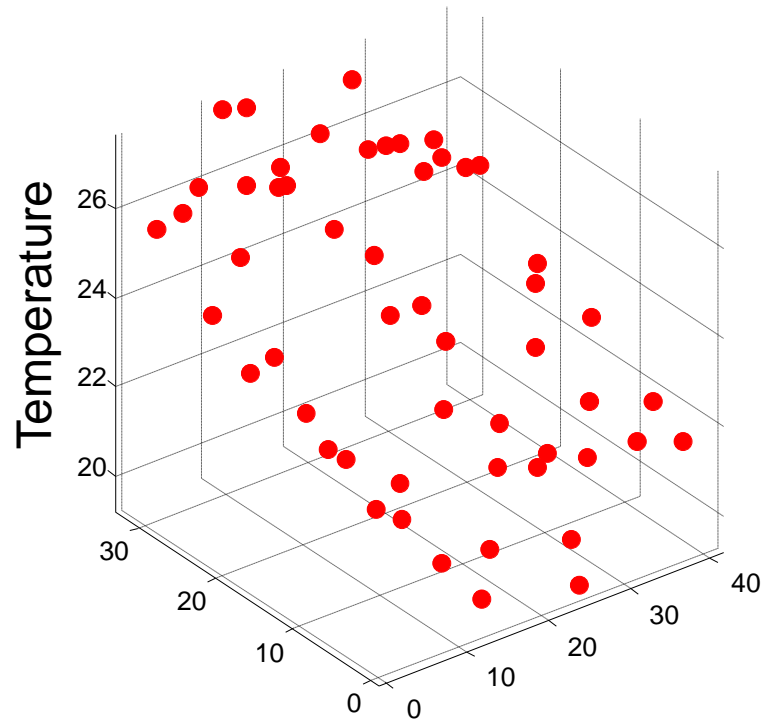
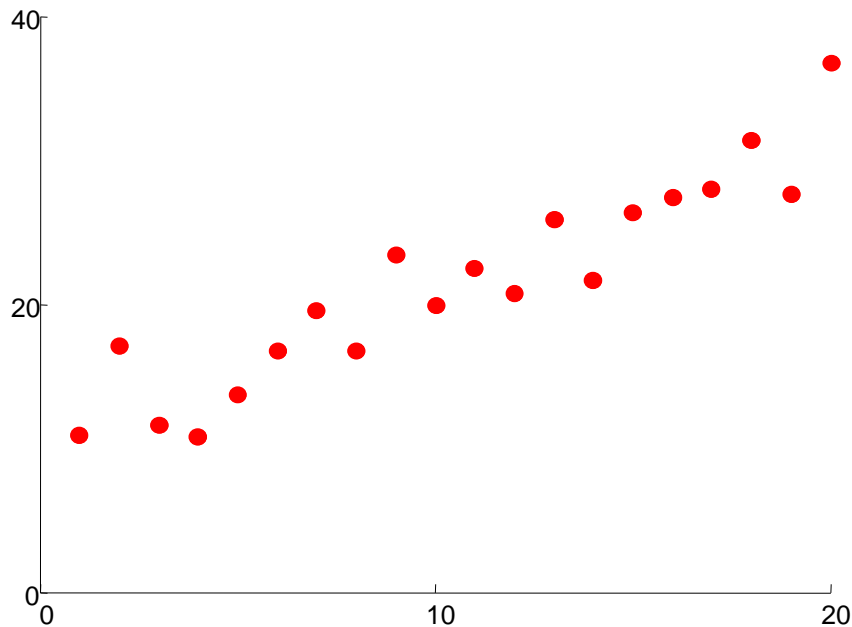
Nonlinear Regression

Kernel Regression

Regression Trees, Splines, Wavelet estimators, ...

# Linear Regression

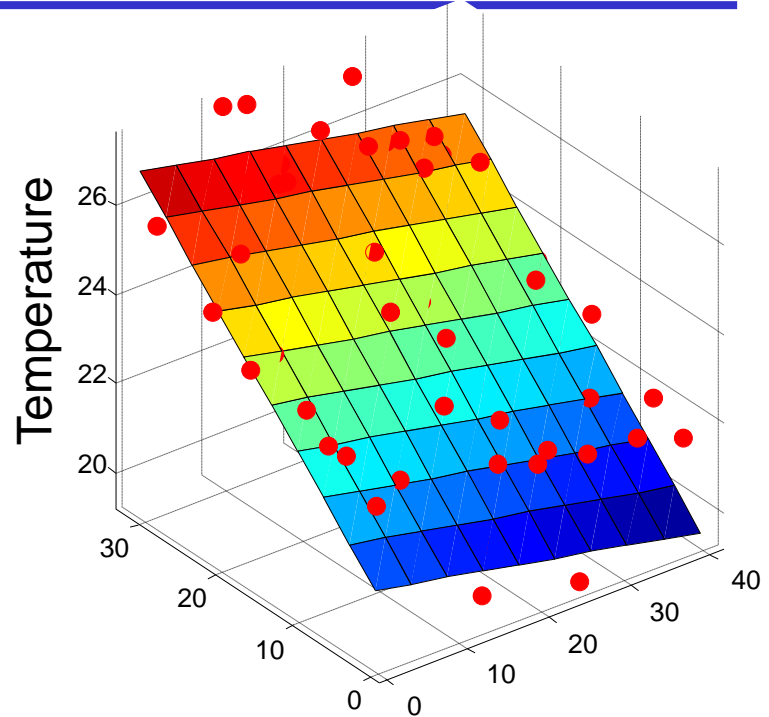
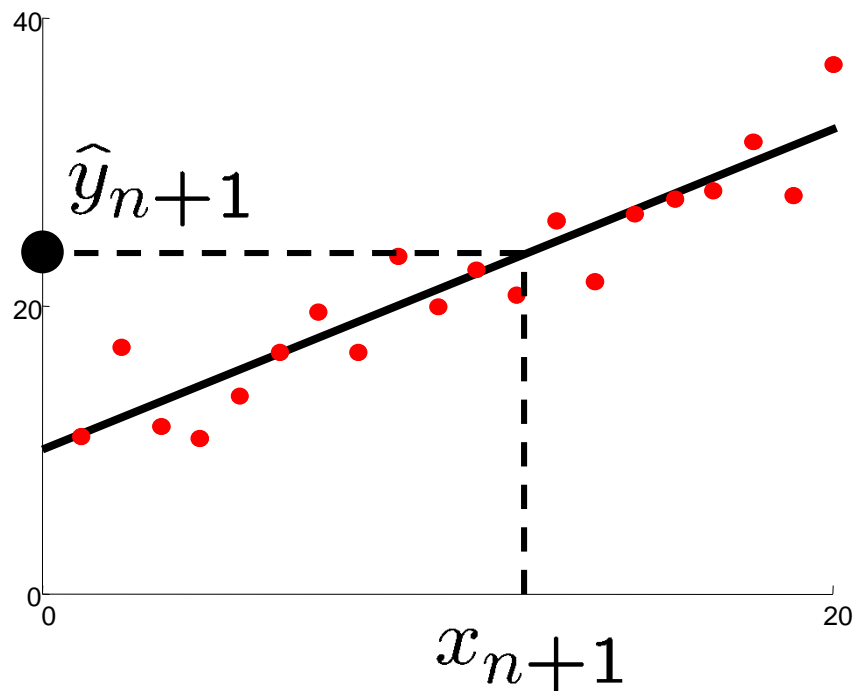
---



Given examples  $(x_i, y_i)_{i=1\dots n}$

Predict  $y_{n+1}$  given a new point  $x_{n+1}$

# Linear Regression

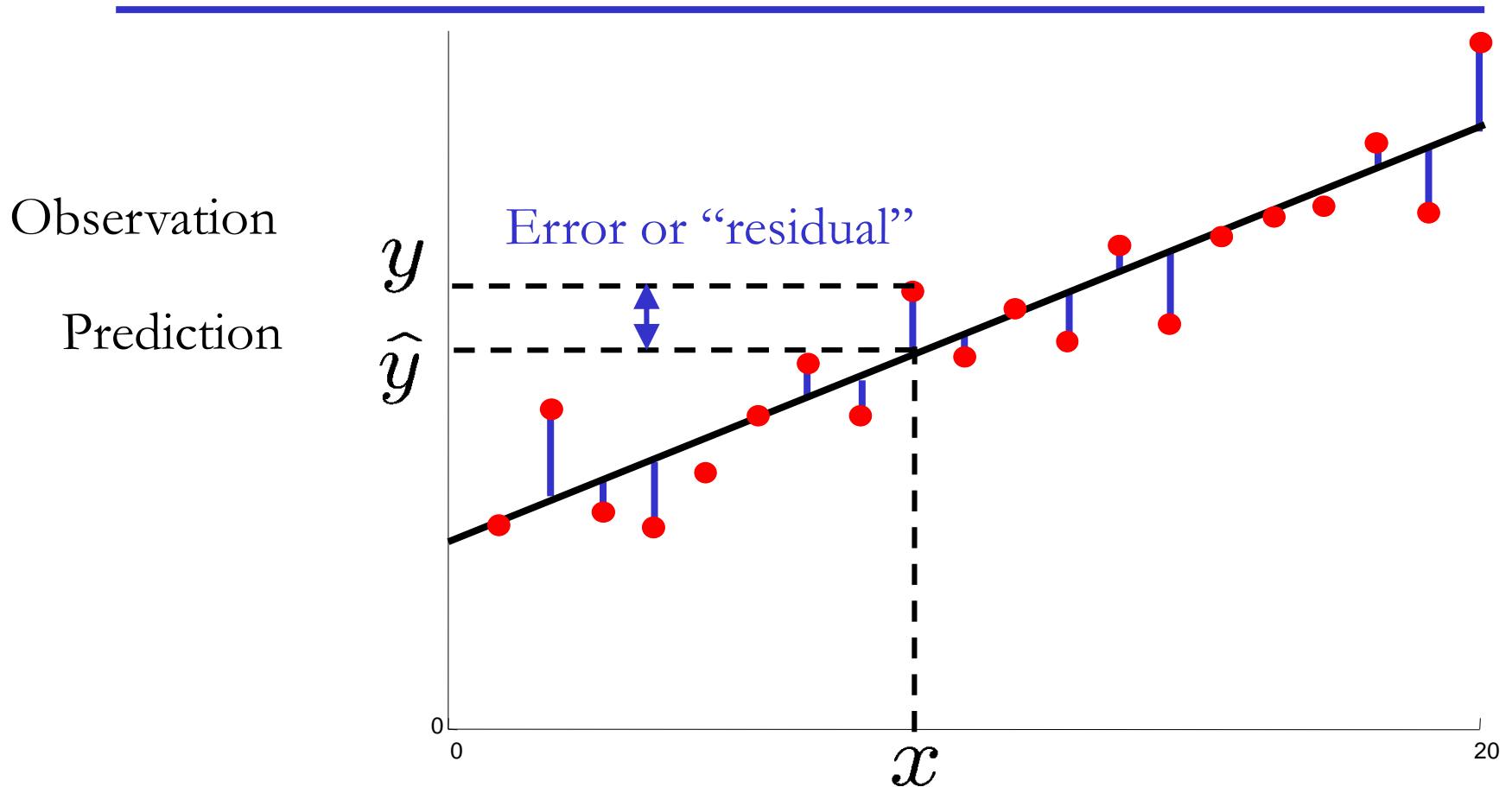


Prediction

$$\hat{y}_i = w_0 + w_1 x_i$$

$$\begin{aligned} \text{Pre}_{\hat{y}_i} &= \begin{pmatrix} 1 & x_{i,1} & x_{i,2} \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ w_2 \end{pmatrix} \\ &= X_i^\top w \end{aligned}$$

# Ordinary Least Squares (OLS)



Sum squared error  $\sum_i (X_i^\top w - y_i)^2$



# Least Squares Estimator

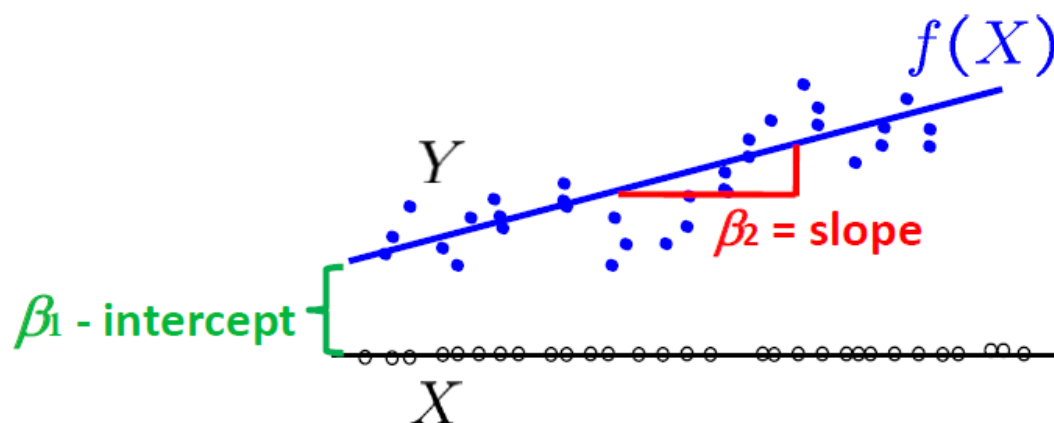
---

$$\hat{f}_n^L = \arg \min_{f \in \mathcal{F}_L} \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2$$

$\mathcal{F}_L$  - Class of Linear functions

Uni-variate case:

$$Y = \beta_1 + \beta_2 X$$



# Least Squares Estimator – Example

---

Goal: To minimize

$$D = \sum_{i=1}^n d_i^2 = \sum_{i=1}^n (y_i - \beta_1 - \beta_2 x_i)^2$$

# Least Squares Estimator

---

Multi-variate case:

$$\begin{aligned}\hat{\beta} &= \arg \min_{\beta} \frac{1}{n} \sum_{i=1}^n (X_i \beta - Y_i)^2 \\ &= \arg \min_{\beta} \frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y})\end{aligned}$$

$$\mathbf{A} = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix} = \begin{bmatrix} X_1^{(1)} & \dots & X_1^{(p)} \\ \vdots & \ddots & \vdots \\ X_n^{(1)} & \dots & X_n^{(p)} \end{bmatrix} \quad \mathbf{Y} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}$$

# Least Squares Estimator

---

Multi-variate case:

$$\hat{\beta} = \arg \min_{\beta} \frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y})$$

$$E(\beta) = (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y})$$



$$\left. \frac{\partial E(\beta)}{\partial \beta} \right|_{\hat{\beta}} = 0$$

# Least Squares Estimator

---

Multi-variate case:

$$\underbrace{(\mathbf{A}^T \mathbf{A})}_{p \times p} \underbrace{\hat{\boldsymbol{\beta}}}_{p \times 1} = \underbrace{\mathbf{A}^T \mathbf{Y}}_{p \times 1}$$

If  $(\mathbf{A}^T \mathbf{A})$  is invertible, we have:

$$\hat{\boldsymbol{\beta}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Y}$$



**You need this for your project!**

# Project 1: Linear regression for classification

---

Data: <http://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data>

## Attribute Information:

1. sepal length in cm
2. sepal width in cm
3. petal length in cm
4. petal width in cm
5. class:
  - Iris Setosa
  - Iris Versicolour
  - Iris Virginica

```
5.1,3.5,1.4,0.2,Iris-setosa
4.9,3.0,1.4,0.2,Iris-setosa
4.7,3.2,1.3,0.2,Iris-setosa
4.6,3.1,1.5,0.2,Iris-setosa
5.0,3.6,1.4,0.2,Iris-setosa
5.4,3.9,1.7,0.4,Iris-setosa
4.6,3.4,1.4,0.3,Iris-setosa
5.0,3.4,1.5,0.2,Iris-setosa
4.4,2.9,1.4,0.2,Iris-setosa
4.9,3.1,1.5,0.1,Iris-setosa
5.4,3.7,1.5,0.2,Iris-setosa
4.8,3.4,1.6,0.2,Iris-setosa
4.8,3.0,1.4,0.1,Iris-setosa
4.3,3.0,1.1,0.1,Iris-setosa
5.8,4.0,1.2,0.2,Iris-setosa
5.7,4.4,1.5,0.4,Iris-setosa
5.4,3.9,1.3,0.4,Iris-setosa
```

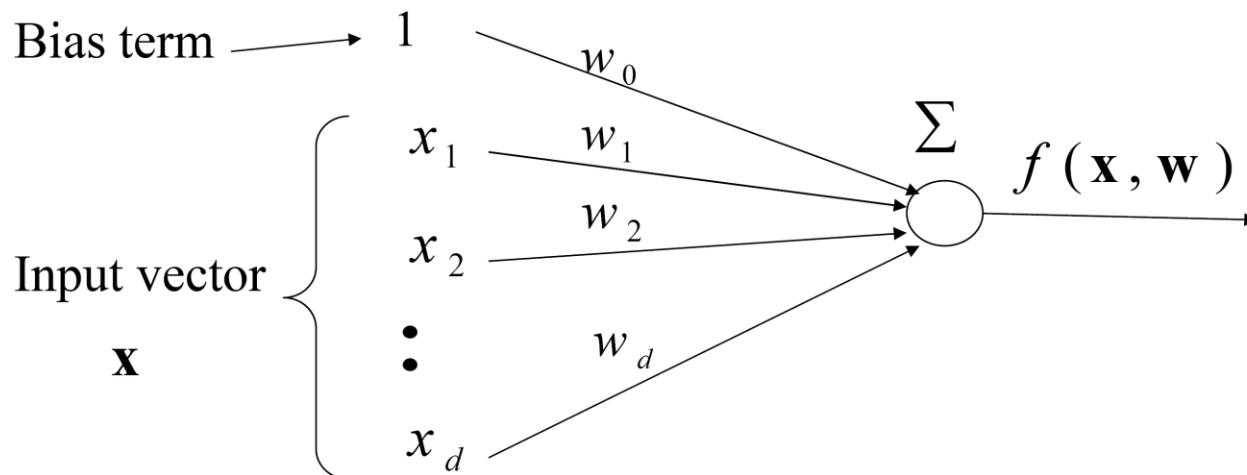
# Linear Regression

---

- **Function**  $f : X \rightarrow Y$  is a linear combination of input components

$$f(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + \dots w_d x_d = w_0 + \sum_{j=1}^d w_j x_j$$

$w_0, w_1, \dots, w_k$  - **parameters (weights)**



# Linear Regression

---

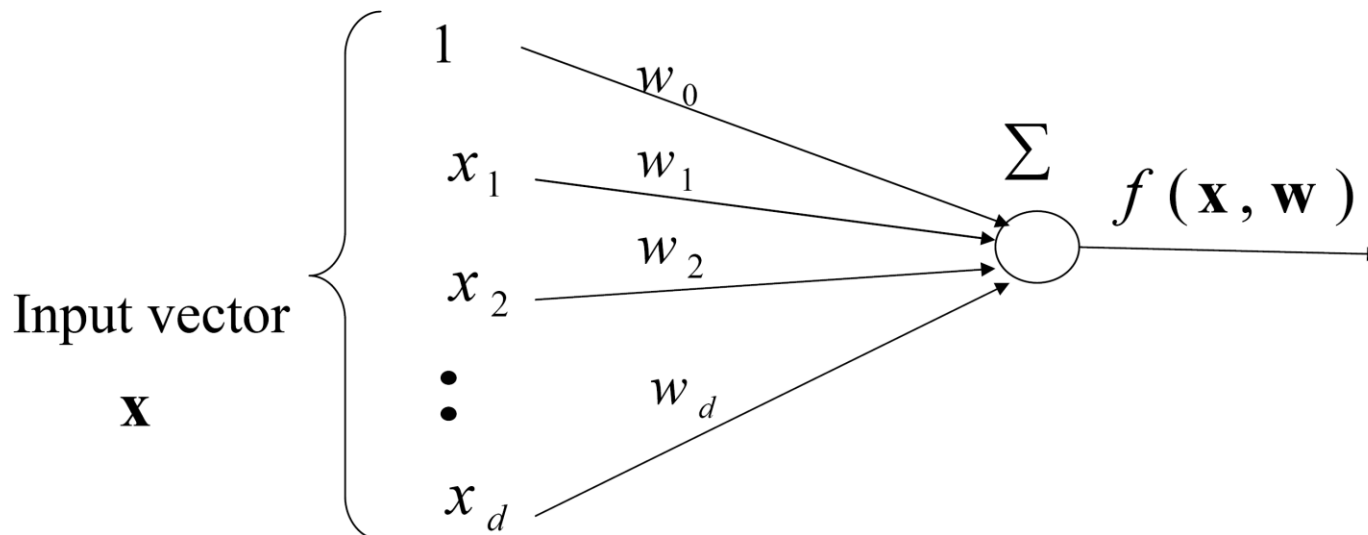
- **Shorter (vector) definition of the model**

- Include bias constant in the input vector

$$\mathbf{x} = (1, x_1, x_2, \dots, x_d)$$

$$f(\mathbf{x}) = w_0 x_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d = \mathbf{w}^T \mathbf{x}$$

$w_0, w_1, \dots, w_k$  - **parameters (weights)**





# Examples

---

- Voltage  $\rightarrow$  Temperature
- Stock prediction  $\rightarrow$  Money
- Processes, memory  $\rightarrow$  Power consumption
- Protein structure  $\rightarrow$  Energy
- Robot arm controls  $\rightarrow$  Torque at effector
- Location, industry, past losses  $\rightarrow$  Premium

# Linear Regression. Optimization.

---

- We want the **weights minimizing the error**

$$J_n = \frac{1}{n} \sum_{i=1, \dots, n} (y_i - f(\mathbf{x}_i))^2 = \frac{1}{n} \sum_{i=1, \dots, n} (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

- For the optimal set of parameters, derivatives of the error with respect to each parameter must be 0

$$\frac{\partial}{\partial w_j} J_n(\mathbf{w}) = -\frac{2}{n} \sum_{i=1}^n (y_i - w_0 x_{i,0} - w_1 x_{i,1} - \dots - w_d x_{i,d}) x_{i,j} = 0$$

- **Vector of derivatives:**

$$\text{grad}_{\mathbf{w}} (J_n(\mathbf{w})) = \nabla_{\mathbf{w}} (J_n(\mathbf{w})) = -\frac{2}{n} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i = \overline{\mathbf{0}}$$

# Linear Regression. Optimization.

---

- $\text{grad}_{\mathbf{w}}(J_n(\mathbf{w})) = \bar{\mathbf{0}}$  defines a set of equations in  $\mathbf{w}$

$$\frac{\partial}{\partial w_0} J_n(\mathbf{w}) = -\frac{2}{n} \sum_{i=1}^n (y_i - w_0 x_{i,0} - w_1 x_{i,1} - \dots - w_d x_{i,d}) = 0$$

$$\frac{\partial}{\partial w_1} J_n(\mathbf{w}) = -\frac{2}{n} \sum_{i=1}^n (y_i - w_0 x_{i,0} - w_1 x_{i,1} - \dots - w_d x_{i,d}) x_{i,1} = 0$$

...

$$\frac{\partial}{\partial w_j} J_n(\mathbf{w}) = -\frac{2}{n} \sum_{i=1}^n (y_i - w_0 x_{i,0} - w_1 x_{i,1} - \dots - w_d x_{i,d}) x_{i,j} = 0$$

...

$$\frac{\partial}{\partial w_d} J_n(\mathbf{w}) = -\frac{2}{n} \sum_{i=1}^n (y_i - w_0 x_{i,0} - w_1 x_{i,1} - \dots - w_d x_{i,d}) x_{i,d} = 0$$

# Solving Linear Regression

---

$$\frac{\partial}{\partial w_j} J_n(\mathbf{w}) = -\frac{2}{n} \sum_{i=1}^n (y_i - w_0 x_{i,0} - w_1 x_{i,1} - \dots - w_d x_{i,d}) x_{i,j} = 0$$

By rearranging the terms we get a **system of linear equations** with  $d+1$  unknowns

$$\mathbf{A}\mathbf{w} = \mathbf{b}$$

$$\begin{aligned} w_0 \sum_{i=1}^n x_{i,0} 1 + w_1 \sum_{i=1}^n x_{i,1} 1 + \dots + w_j \sum_{i=1}^n x_{i,j} 1 + \dots + w_d \sum_{i=1}^n x_{i,d} 1 &= \sum_{i=1}^n y_i 1 \\ w_0 \sum_{i=1}^n x_{i,0} x_{i,1} + w_1 \sum_{i=1}^n x_{i,1} x_{i,1} + \dots + w_j \sum_{i=1}^n x_{i,j} x_{i,1} + \dots + w_d \sum_{i=1}^n x_{i,d} x_{i,1} &= \sum_{i=1}^n y_i x_{i,1} \\ &\vdots \\ w_0 \sum_{i=1}^n x_{i,0} x_{i,j} + w_1 \sum_{i=1}^n x_{i,1} x_{i,j} + \dots + w_j \sum_{i=1}^n x_{i,j} x_{i,j} + \dots + w_d \sum_{i=1}^n x_{i,d} x_{i,j} &= \sum_{i=1}^n y_i x_{i,j} \\ &\vdots \end{aligned}$$

# Solving Linear Regression

---

- The optimal set of weights satisfies:

$$\nabla_{\mathbf{w}} (J_n(\mathbf{w})) = -\frac{2}{n} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i = \bar{\mathbf{0}}$$

Leads to a **system of linear equations (SLE)** with  $d+1$  unknowns of the form

$$\mathbf{A}\mathbf{w} = \mathbf{b}$$

$$w_0 \sum_{i=1}^n x_{i,0} x_{i,j} + w_1 \sum_{i=1}^n x_{i,1} x_{i,j} + \dots + w_j \sum_{i=1}^n x_{i,j} x_{i,j} + \dots + w_d \sum_{i=1}^n x_{i,d} x_{i,j} = \sum_{i=1}^n y_i x_{i,j}$$

**Solution to SLE:**

$$\mathbf{w} = \mathbf{A}^{-1} \mathbf{b}$$

# Gradient Descent Solution

---

**Goal:** the weight optimization in the linear regression model

$$J_n = \text{Error}(\mathbf{w}) = \frac{1}{n} \sum_{i=1, \dots, n} (y_i - f(\mathbf{x}_i, \mathbf{w}))^2$$

An alternative to SLE solution:

- **Gradient descent**

**Idea:**

- Adjust weights in the direction that improves the Error
- The gradient tells us what is the right direction

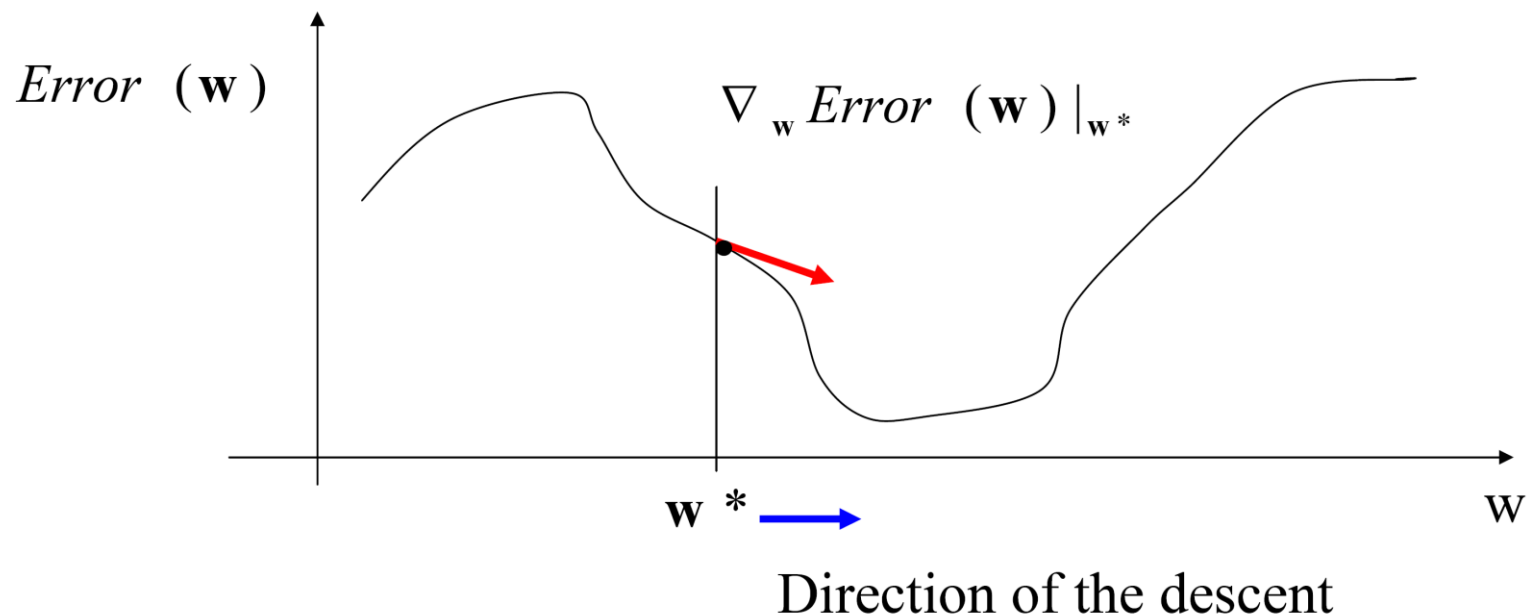
$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla_{\mathbf{w}} \text{Error}_i(\mathbf{w})$$

$\alpha > 0$  - a **learning rate** (scales the gradient changes)

# Gradient Descent Method

---

- Descend using the gradient information

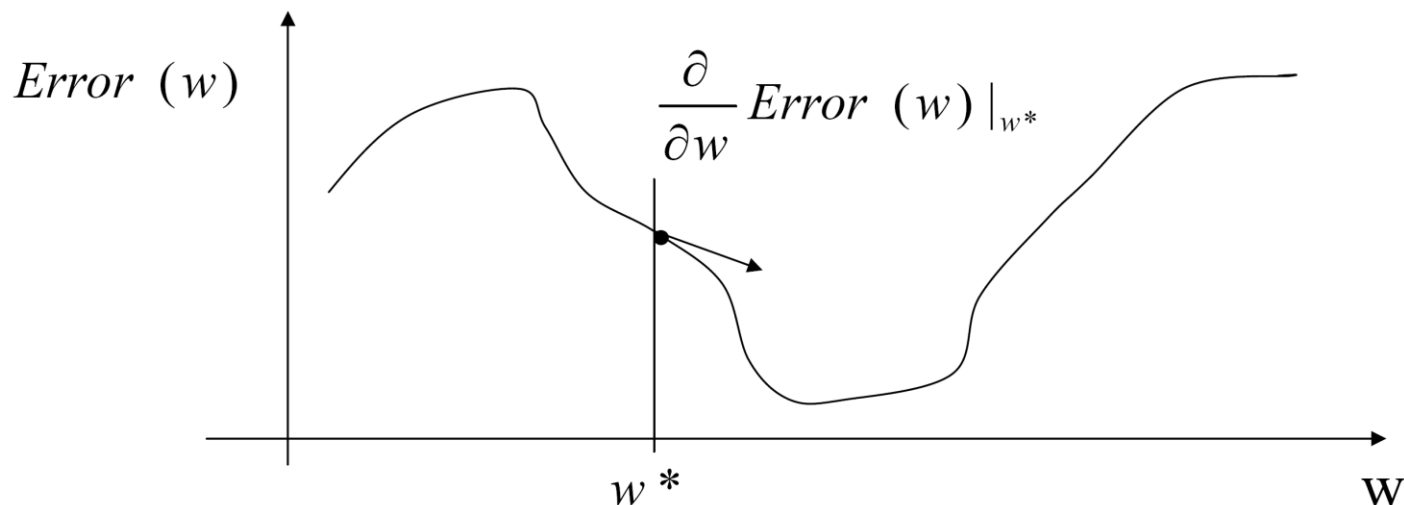


- Change the value of  $w$  according to the gradient

$$w \leftarrow w - \alpha \nabla_w Error_i(w)$$

# Gradient Descent Method

---



- New value of the parameter

$$w_j \leftarrow w_j^* - \alpha \frac{\partial}{\partial w_j} Error (w) |_{w^*} \quad \text{For all } j$$

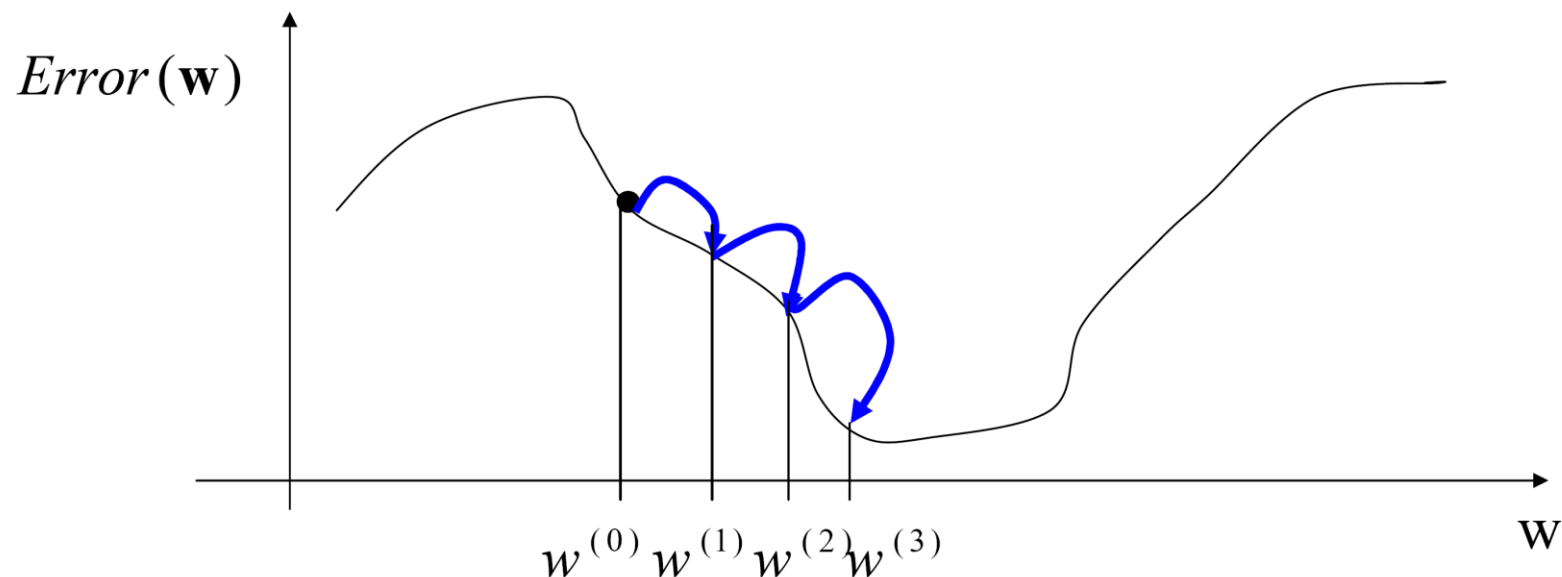
$\alpha > 0$  - a learning rate (scales the gradient changes)



# Gradient Descent Method

---

- Iteratively approaches the optimum of the Error function



# Online Gradient Algorithm

---

Linear model

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

On-line error  $J_{online} = Error_i(\mathbf{w}) = \frac{1}{2}(y_i - f(\mathbf{x}_i, \mathbf{w}))^2$

**On-line algorithm:** generates a sequence of online updates

**(i)-th update step with :**  $D_i = \langle \mathbf{x}_i, y_i \rangle$

**j-th weight:**

$$w_j^{(i)} \leftarrow w_j^{(i-1)} - \alpha(i) \frac{\partial Error_i(\mathbf{w})}{\partial w_j} \Big|_{\mathbf{w}^{(i-1)}}$$

$$w_j^{(i)} \leftarrow w_j^{(i-1)} + \alpha(i)(y_i - f(\mathbf{x}_i, \mathbf{w}^{(i-1)}))x_{i,j}$$

**Fixed learning rate:**  $\alpha(i) = C$

- Use a small constant

**Annealed learning rate:**  $\alpha(i) \approx \frac{1}{i}$

- Gradually rescales changes

# Online Regression Algorithm

---

**Online-linear-regression** ( $D$ , *number of iterations*)

**Initialize** weights  $\mathbf{w} = (w_0, w_1, w_2 \dots w_d)$

**for**  $i=1:1$ : *number of iterations*

**do**        **select** a data point  $D_i = (\mathbf{x}_i, y_i)$  from  $D$

**set** learning rate  $\alpha(i)$

**update** weight vector

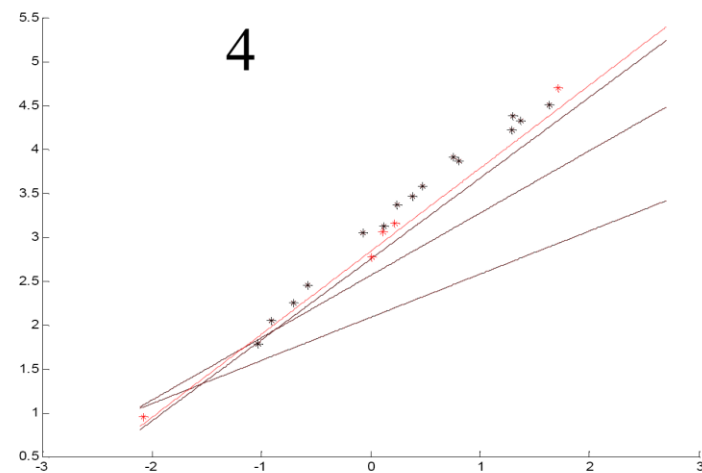
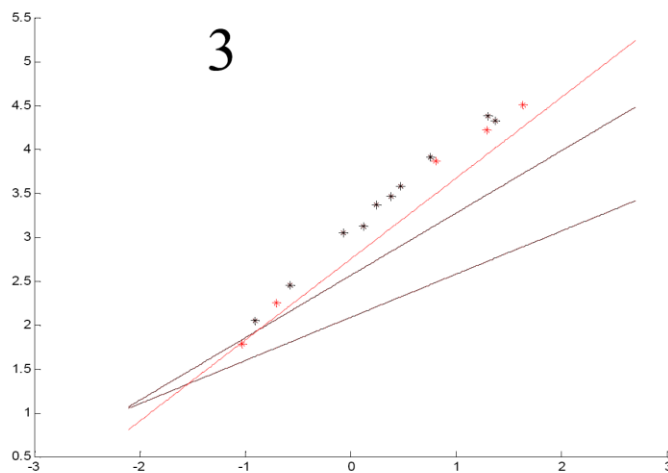
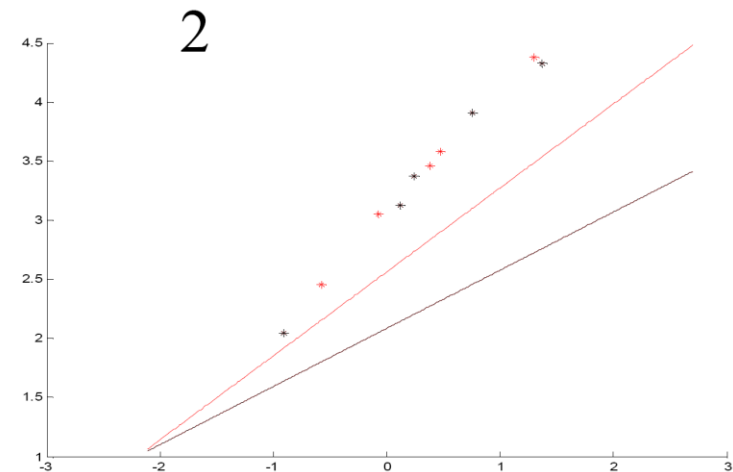
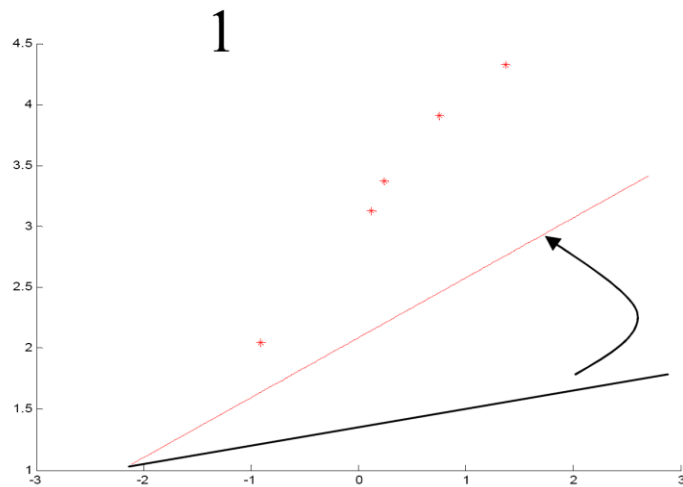
$$\mathbf{w} \leftarrow \mathbf{w} + \alpha(i)(y_i - f(\mathbf{x}_i, \mathbf{w}))\mathbf{x}_i$$

**end for**

**return** weights  $\mathbf{w}$

- **Advantages:** very easy to implement, continuous data streams

# On-line Learning Example



# Linear Regression

---

## Summary

$$\hat{f}_n^L = \arg \min_{f \in \mathcal{F}_L} \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2$$

- is the orthogonal projection of  $Y$  onto the linear subspace spanned by the columns of  $A$

# Linear Regression - Example

---

## Data pre-processing

Age

Sex

Data source

...

**Using residual for analysis!**

# Ridge Regression

---

- Suppose we have an estimator  $\hat{f}(\mathbf{z}) = \mathbf{z}^\top \hat{\beta}$
- To see if  $\hat{f}(\mathbf{z}) = \mathbf{z}^\top \hat{\beta}$  is a good candidate, we can ask ourselves two questions:
  - 1.) Is  $\hat{\beta}$  close to the true  $\beta$ ?
  - 2.) Will  $\hat{f}(\mathbf{z})$  fit future observations well?

# Ridge Regression

---

- Just because  $\hat{f}(\mathbf{z})$  fits our data well, this doesn't mean that it will be a good fit to new data
- In fact, suppose that we take new measurements  $y'_i$  at the same  $\mathbf{z}_i$ 's:

$$(\mathbf{z}_1, y'_1), (\mathbf{z}_2, y'_2), \dots, (\mathbf{z}_n, y'_n)$$

- So if  $\hat{f}(\cdot)$  is a good model, then  $\hat{f}(\mathbf{z}_i)$  should also be close to the new target  $y'_i$

This is the notion of **prediction error** (PE)



# Ridge Regression

---

- If the  $\beta_j$ 's are unconstrained...
  - They can explode
  - And hence are susceptible to very high variance
- To control variance, we might **regularize** the coefficients
  - E.g. control how large the coefficients grow

# Ridge Regression

---

$$\text{minimize } \sum_{i=1}^n (y_i - \boldsymbol{\beta}^\top \mathbf{z}_i)^2 \text{ s.t. } \sum_{j=1}^p \beta_j^2 \leq t$$

$$\Leftrightarrow \text{minimize } (y - \mathbf{Z}\boldsymbol{\beta})^\top (y - \mathbf{Z}\boldsymbol{\beta}) \text{ s.t. } \sum_{j=1}^p \beta_j^2 \leq t$$

By convention:

- $\mathbf{Z}$  is assumed to be standardized (mean 0, unit variance)
- $y$  is assumed to be centered

# Ridge Regression

---

- Can write the ridge constraint as the following **penalized** residual sum of squares (PRSS):

$$\begin{aligned} PRSS(\boldsymbol{\beta})_{\ell_2} &= \sum_{i=1}^n (y_i - \mathbf{z}_i^\top \boldsymbol{\beta})^2 + \lambda \sum_{j=1}^p \beta_j^2 \\ &= (\mathbf{y} - \mathbf{Z}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{Z}\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_2^2 \end{aligned}$$

- Its solution may have smaller average PE than  $\hat{\boldsymbol{\beta}}^{\text{ls}}$
- $PRSS(\boldsymbol{\beta})_{\ell_2}$  is convex, and hence has a unique solution
- Taking derivatives, we obtain:

$$\frac{\partial PRSS(\boldsymbol{\beta})_{\ell_2}}{\partial \boldsymbol{\beta}} = -2\mathbf{Z}^\top (\mathbf{y} - \mathbf{Z}\boldsymbol{\beta}) + 2\lambda\boldsymbol{\beta}$$

# Ridge Regression

---

- The solution to  $PRSS(\hat{\beta})_{\ell_2}$  is now seen to be:

$$\hat{\beta}_{\lambda}^{\text{ridge}} = (\mathbf{Z}^{\top} \mathbf{Z} + \lambda \mathbf{I}_p)^{-1} \mathbf{Z}^{\top} \mathbf{y}$$

- Solution is indexed by the tuning parameter  $\lambda$  (more on this later)
- So for each  $\lambda$ , we have a solution  
Hence, the  $\lambda$ 's trace out a path of solutions
- $\lambda$  is the shrinkage parameter
  - $\lambda$  controls the size of the coefficients
  - $\lambda$  controls amount of **regularization**
  - As  $\lambda \downarrow 0$ , we obtain the least squares solutions
  - As  $\lambda \uparrow \infty$ , we have  $\hat{\beta}_{\lambda=\infty}^{\text{ridge}} = 0$  (intercept-only model)

# Ridge Regression

---

- Need disciplined way of selecting  $\lambda$ :  
That is, we need to “tune” the value of  $\lambda$
- In their original paper, Hoerl and Kennard introduced **ridge traces**:

Plot the components of  $\hat{\beta}_{\lambda}^{\text{ridge}}$  against  $\lambda$   
Choose  $\lambda$  for which the coefficients are not rapidly changing  
and have “sensible” signs  
No objective basis; heavily criticized by many

- Standard practice now is to use cross-validation

# Cross Validation

---

- We need a disciplined way of choosing  $\lambda$
- Obviously want to choose  $\lambda$  that minimizes the mean squared error
- Issue is part of the bigger problem of **model selection**

# Cross Validation

---

- If we have a good model, it should predict well when we have new data
- In machine learning terms, we compute our statistical model  $\hat{f}(\cdot)$  from the **training set**
- A good estimator  $\hat{f}(\cdot)$  should then perform well on a new, independent set of data
- We “test” or assess how well  $\hat{f}(\cdot)$  performs on the new data, which we call the **test set**

# Cross Validation

---

- Ideally, we would separate our available data into both training and test sets

In reality, it is very difficult!

- Hope to come up with the best-trained algorithm that will stand up to the test



# Cross Validation

---

Most common approach is  **$K$ -fold cross validation**:

- Partition the training data  $T$  into  $K$  separate sets of equal size

Suppose  $T = (T_1, T_2, \dots, T_K)$

Commonly chosen  $K$ 's are  $K = 5$  and  $K = 10$

- For each  $k = 1, 2, \dots, K$ , fit the model  $\hat{f}_{-k}^{(\lambda)}(\mathbf{z})$  to the training set excluding the  $k$ th-fold  $T_k$
- Compute the fitted values for the observations in  $T_k$ , based on the training data that excluded this fold
- Compute the cross-validation (CV) error for the  $k$ -th fold:

$$(\text{CV Error})_k^{(\lambda)} = |T_k|^{-1} \sum_{(\mathbf{z}, y) \in T_k} (y - \hat{f}_{-k}^{(\lambda)}(\mathbf{z}))^2$$

# Cross Validation

---

- The model then has overall cross-validation error:

$$(\text{CV Error})^{(\lambda)} = K^{-1} \sum_{k=1}^K (\text{CV Error})_k^{(\lambda)}$$

- Select  $\lambda^*$  as the one with minimum  $(\text{CV Error})^{(\lambda)}$

# Cross Validation

---

What happens when  $K = 1$ ?

This is called **leave-one-out cross validation**

For squared error loss, there is a convenient approximation to  $CV(1)$ , which is the leave one-out CV error

# Linear Regression Summary

---

$$\min_{\beta} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) + \lambda \text{pen}(\beta) = \min_{\beta} J(\beta) + \lambda \text{pen}(\beta)$$

Least Square Solution

$$\text{pen}(\beta) = \|\beta\|_2^2$$

**Ridge Regression**

$$\text{pen}(\beta) = \|\beta\|_1$$

**Lasso Regression**

Lasso (L1 penalty) results in sparse solutions –vector with more zero coordinates. Will come to Lasso later!