# Supervised Learning : Linear  Regression and Logistic Regression

Dr. Srijith P K

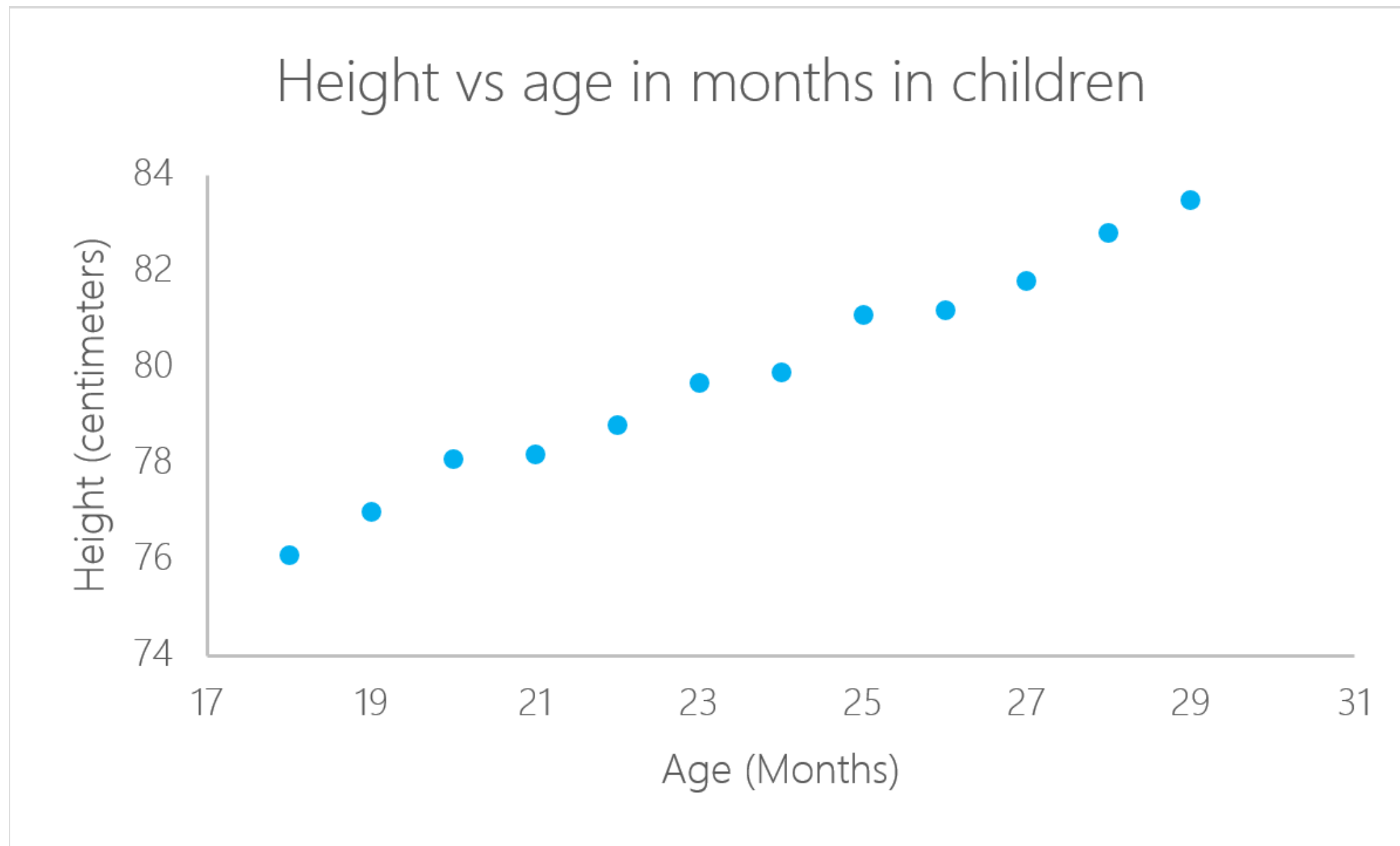Computer Science and Engineering

IIT Hyderabad
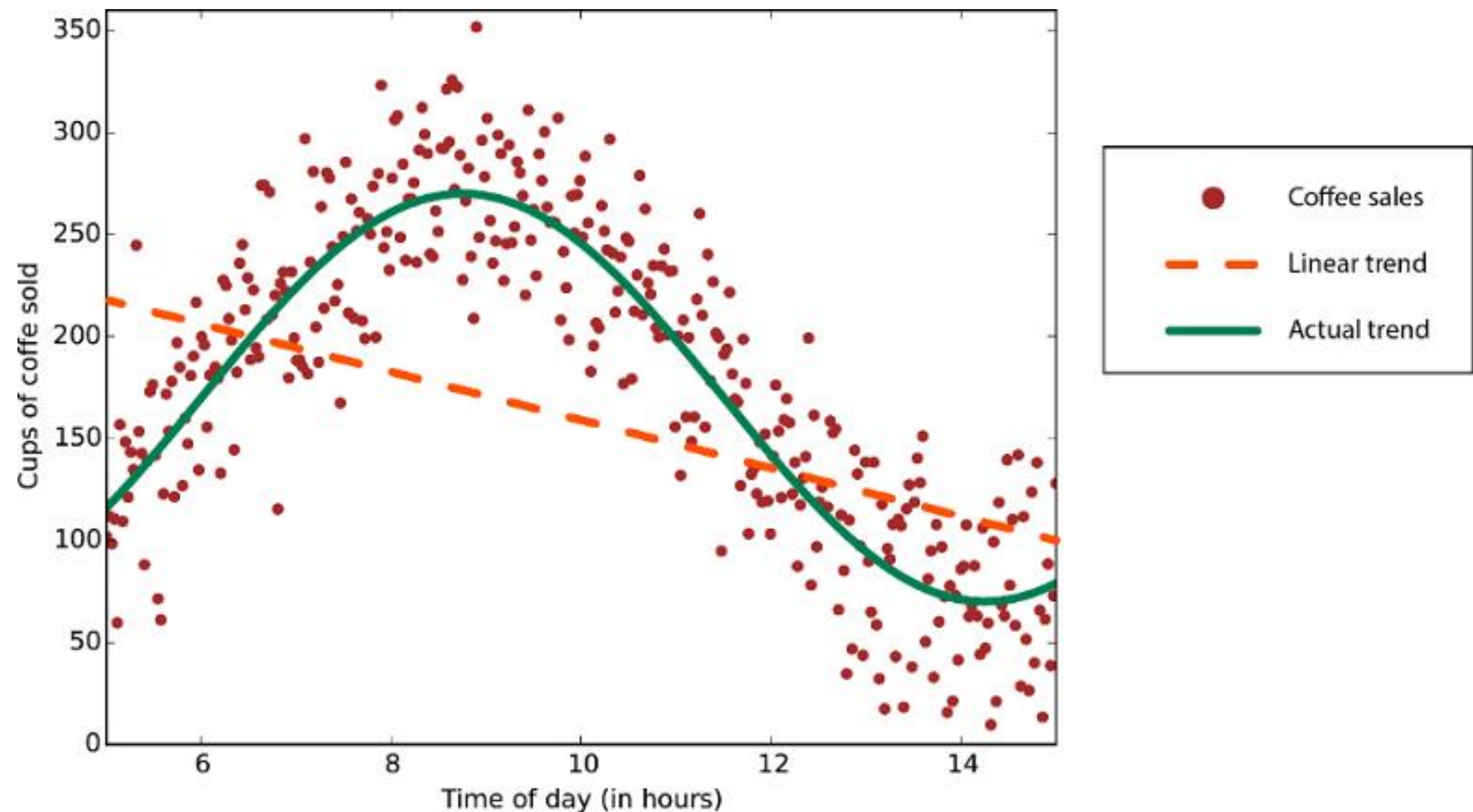
# Supervised learning



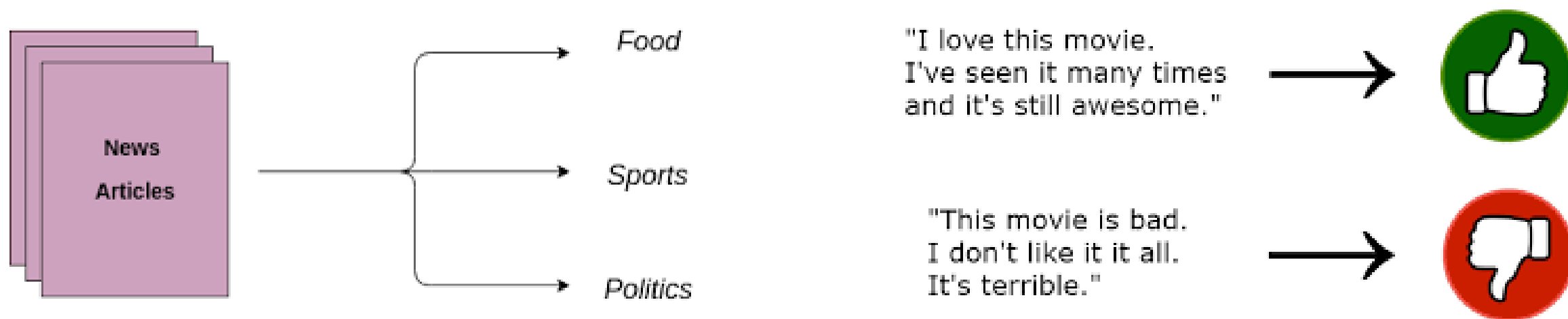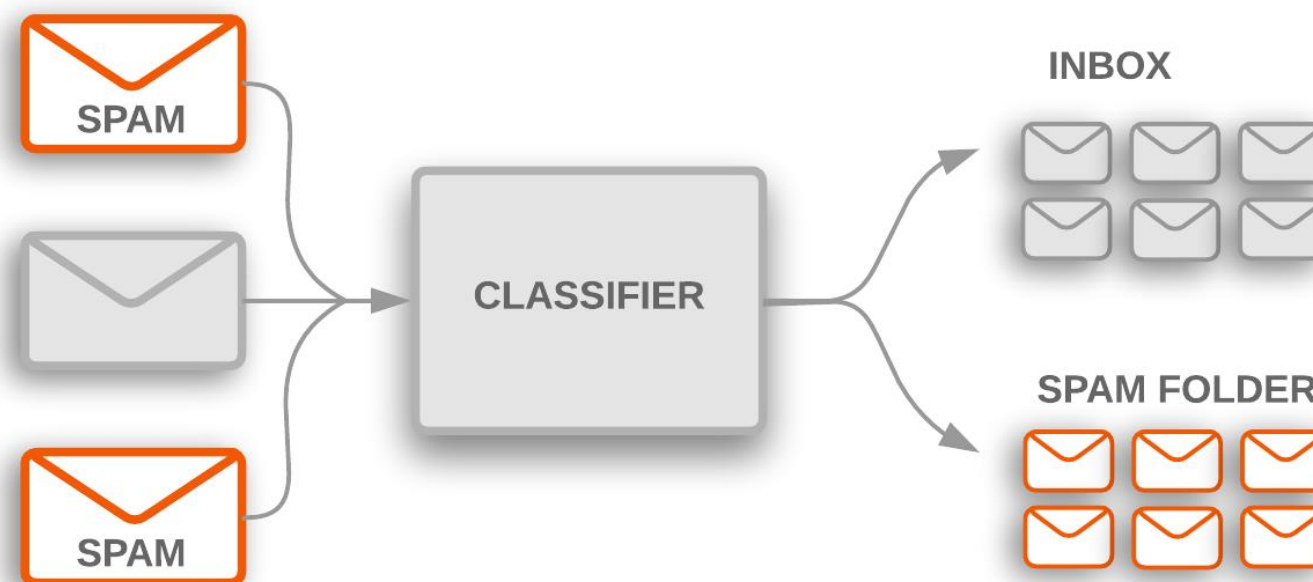Height vs age in months in children

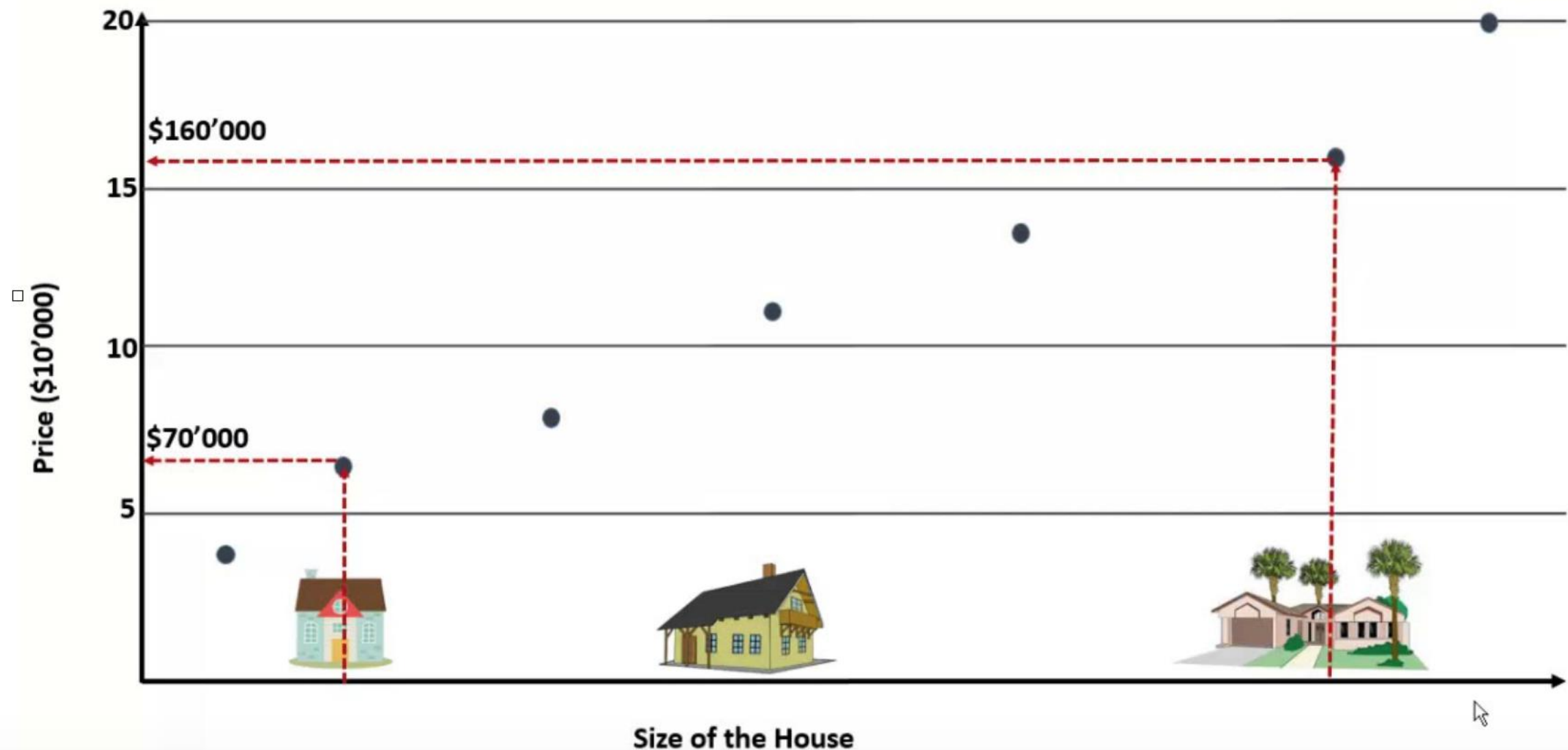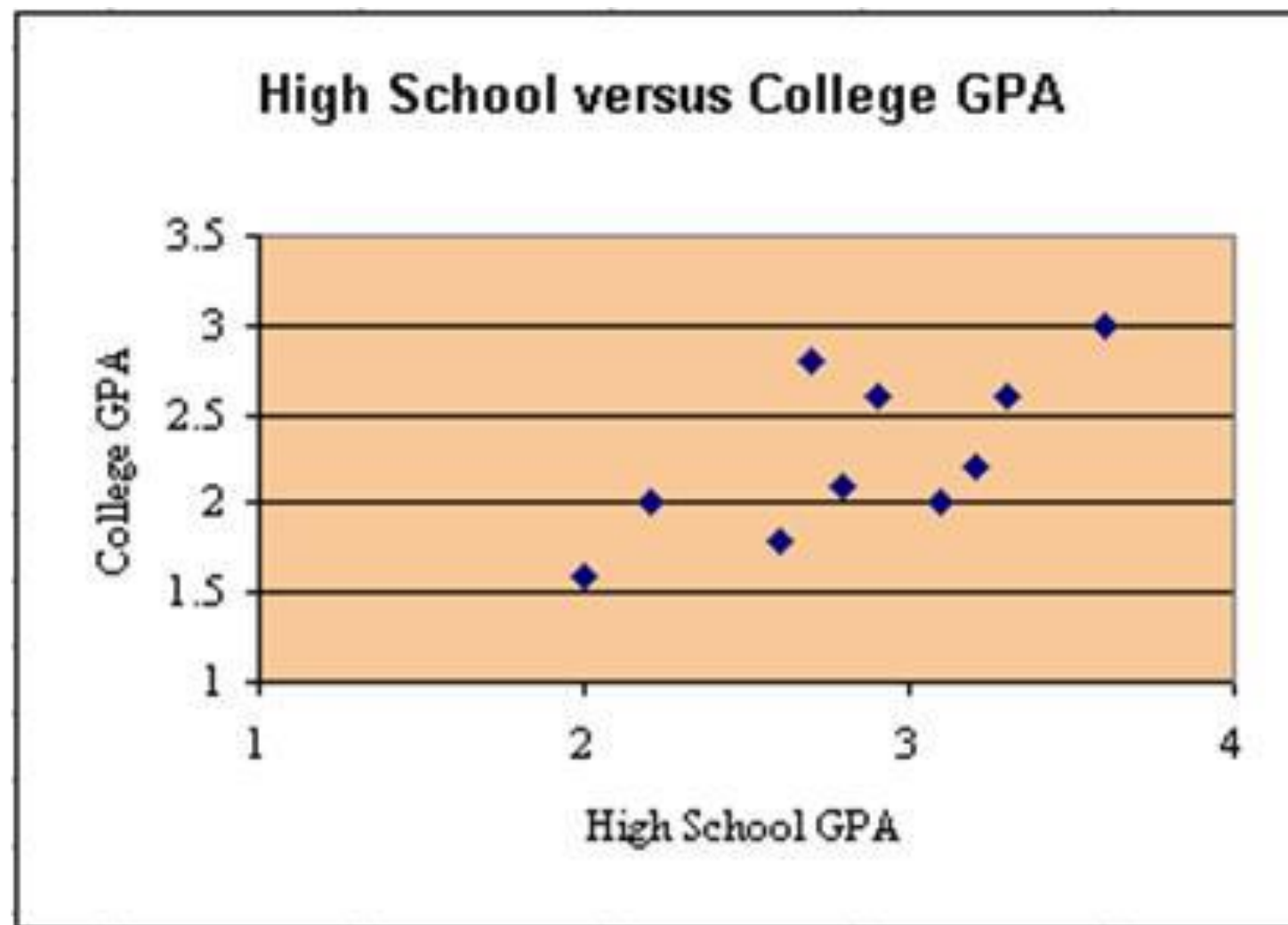# Supervised learning

# Supervised Learning

# Outline

- Supervised Learning

- Regression

- Classification

- Linear Regression

- Logistic Regression

- Poisson Regression

# Supervised learning : Regression

**Estimating Price of a house**

# Supervised learning : Regression



High School versus College GPA

# Supervised learning : Regression



High School versus College GPA

Real valued targets (outputs)

**Generalization performance**

Goal is to learn a function which maps inputs to outputs so that it will predict well on future data points
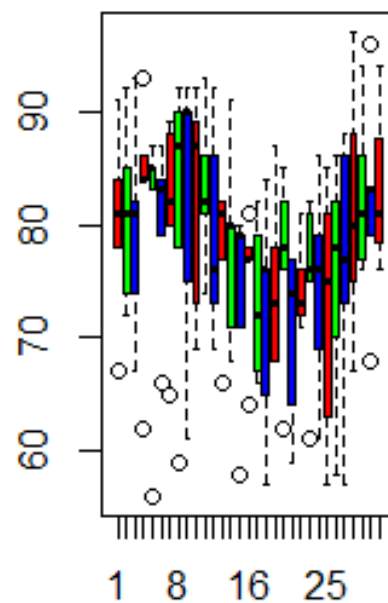
# Airquality data.

- Data set has various air quality parameters in New York city.

- These are the parameters in the data set:

- Daily temperature from May to August

- Solar radiation data

- Ozone data

- Wind data

- Goal : predict the temperature for a particular month in New York using solar radiation, ozone and wind data.

# Airquality data

```
##       Ozone Solar.R Wind Temp Month Day
## 1       41     190  7.4   67     5   1
## 2       36     118  8.0   72     5   2
## 3       12     149 12.6   74     5   3
## 4       18     313 11.5   62     5   4
## 5       NA      NA 14.3   56     5   5
## 6       28      NA 14.9   66     5   6
## 7       23     299  8.6   65     5   7
## 8       19      99 13.8   59     5   8
## 9        8      19 20.1   61     5   9
## 10      NA     194  8.6   69     5  10
```
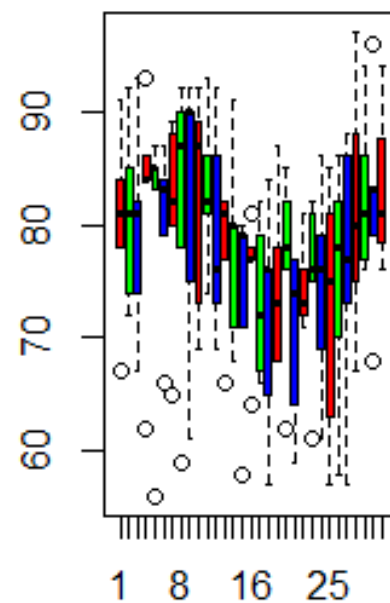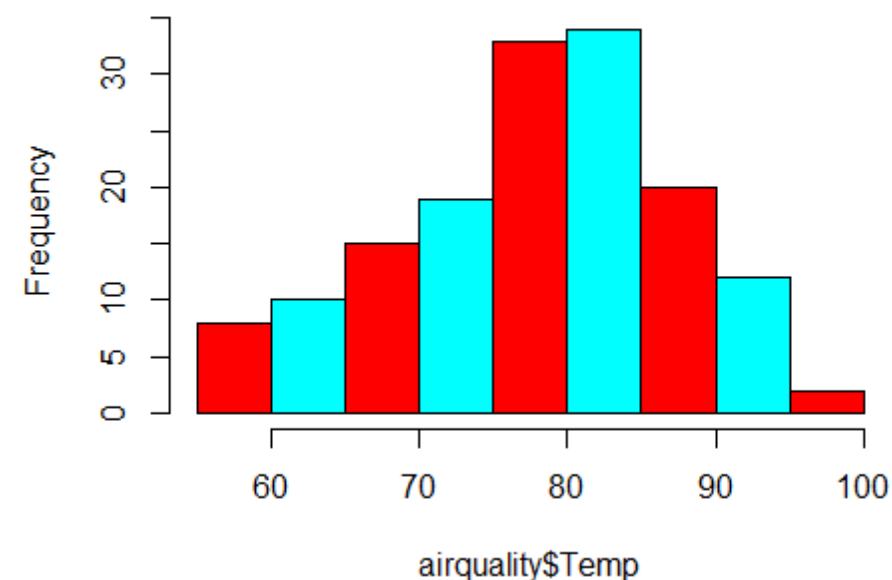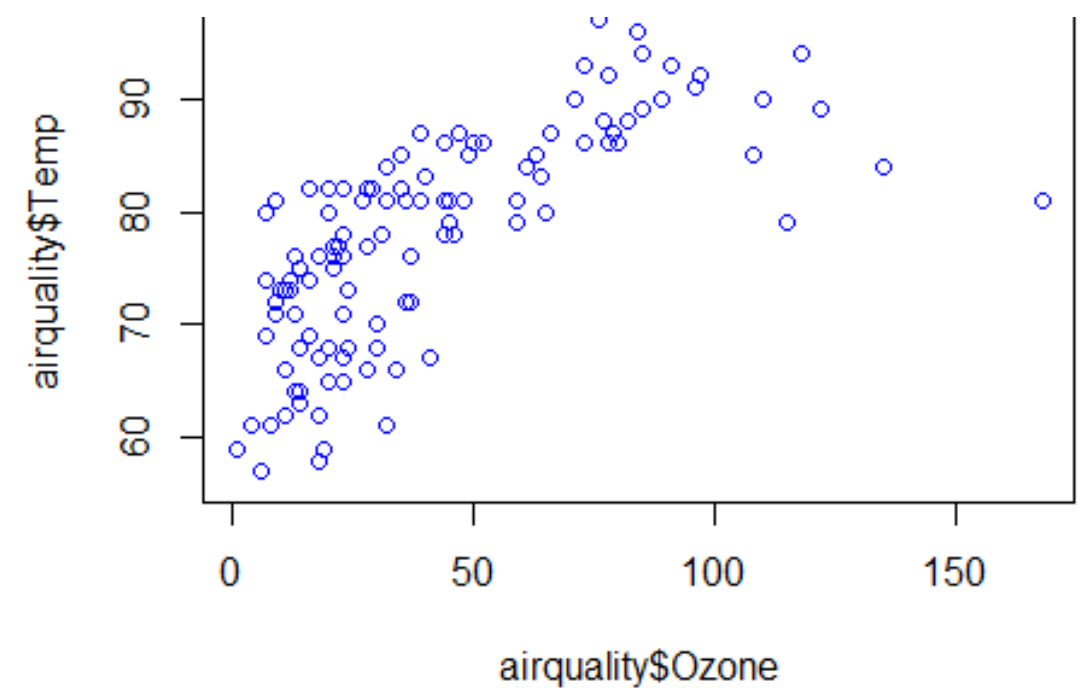


**Month 5**



**Month 6**

**Histogram of airquality$Temp**

# Airquality data

# Linear regression

- Temp=w1.Solar.R +w2.Ozone + w3.Wind + error.

- Temperature of house depends on ozone, wind and solar radiations

- linear regression helps to discover relation between dependent and independent variables

# Linear Regression

- Observations need not lie on a line

  - Observations are not generated by a linear line

  - Observations are noisy, due to measurement errors

# Linear regression

- Learn a function which maps input to output f : X -> Y

Regression Output is real and scalar, $y \in \mathbb{R}$

- Consider a Linear function

Estimated (or predicted) Y value for observation i

Estimate of the regression intercept

Estimate of the regression slope

Value of X for observation i

$$\hat{Y}_i = w_0 + w_1 X_i$$

# Linear regression

- Learn a function which maps input to output f : X -> Y

Regression Output is real and scalar, $y \in \mathbb{R}$
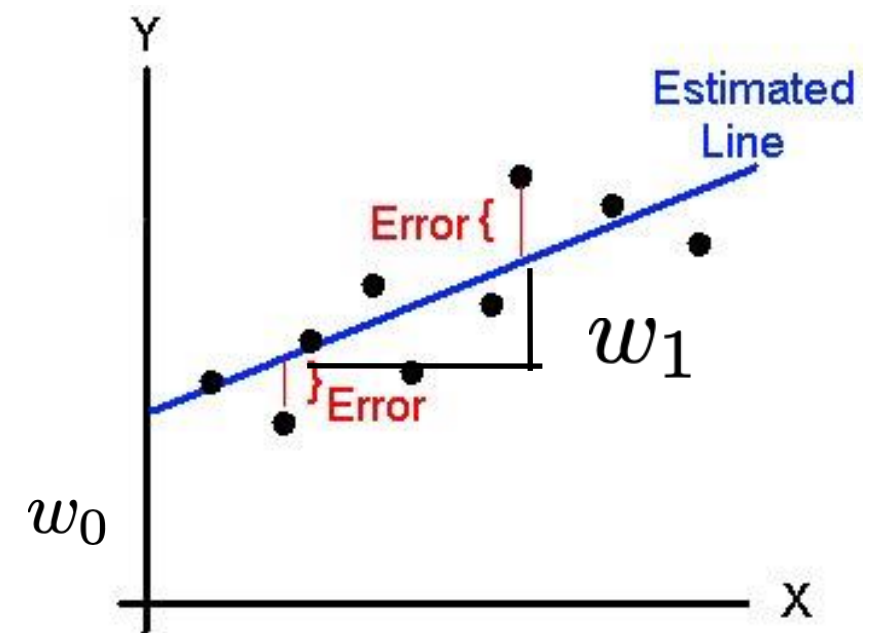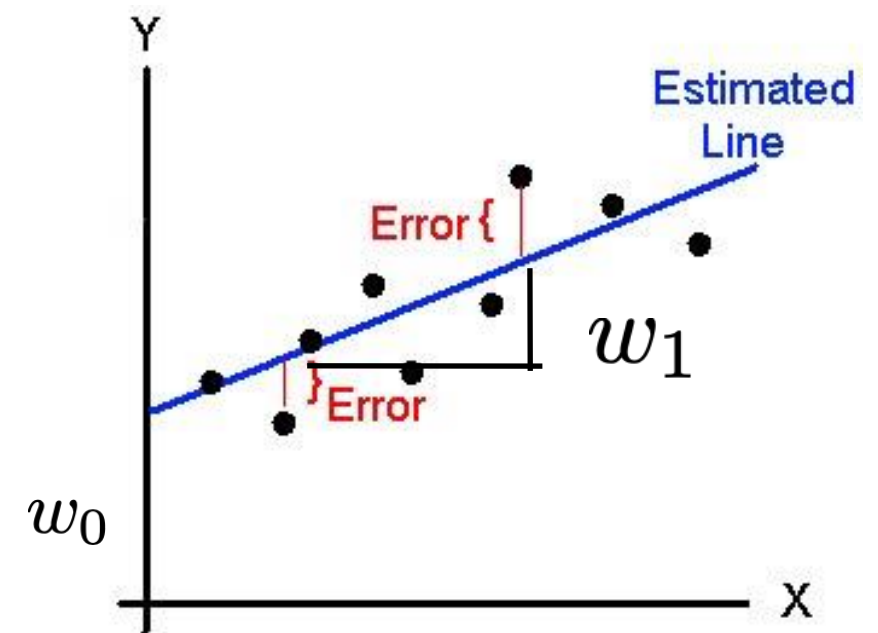
- Consider a Linear function

Estimated (or predicted) Y value for observation i

Estimate of the regression intercept

Estimate of the regression slope

Value of X for observation i

$$\hat{Y}_i = w_0 + w_1 X_i$$

$$= X_i^\top w$$



1 dim input  $X_i = [1, X_i]^\top$  $w = [w_0, w_1]^\top$

D dim input $X_i = [1, X_{i1}, \ldots, X_{iD}]^\top$  $w = [w_0, w_1, \ldots, w_D]^\top$

# Linear Regression - Learning Parameters

- Learn the function which passes through as many points as possible : Minimize the **Least Squares Error**

$$X_i = [1, X_{i1}, \ldots, X_{iD}]^\top$$

Design matrix

$$X = \begin{pmatrix} X_{1,}\, X_{2,}\, \ldots\, X_N \end{pmatrix}$$

(D+1) × N

$$E(w) = \frac{1}{2} \sum_{i=1}^{N} (y_i - \hat{y}_i)^2$$

$$= \frac{1}{2} \sum_{i=1}^{N} (y_i - X_i^\top w)^2$$
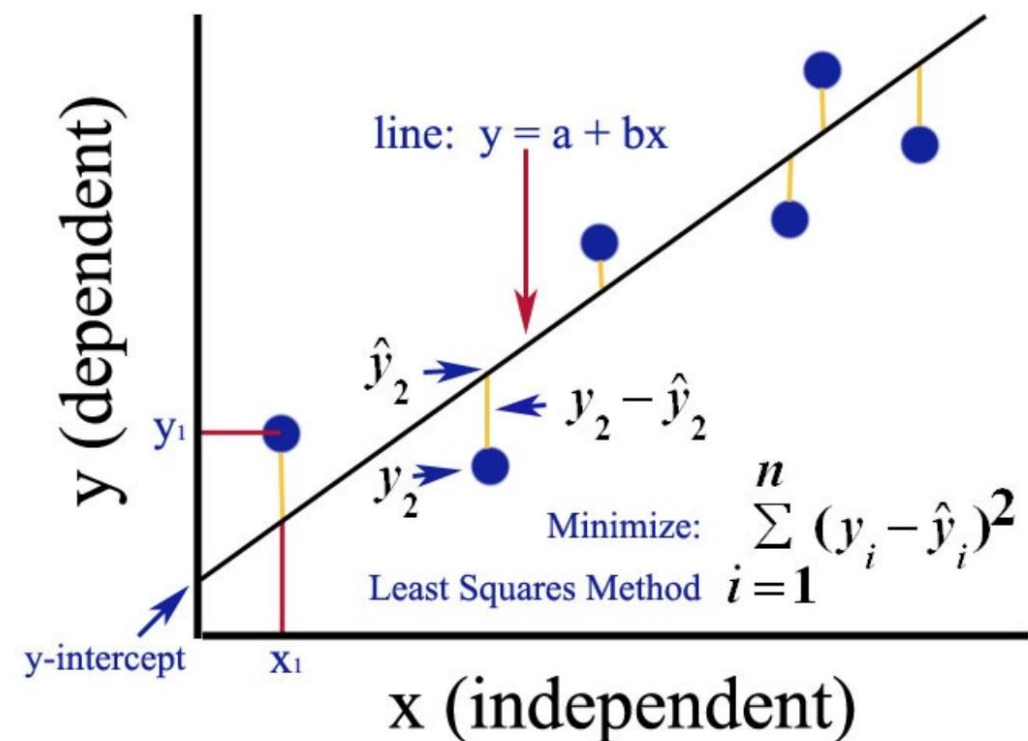
$$\frac{1}{2} \|(y - X^\top w)\|^2$$

# Linear Regression - Learning Parameters

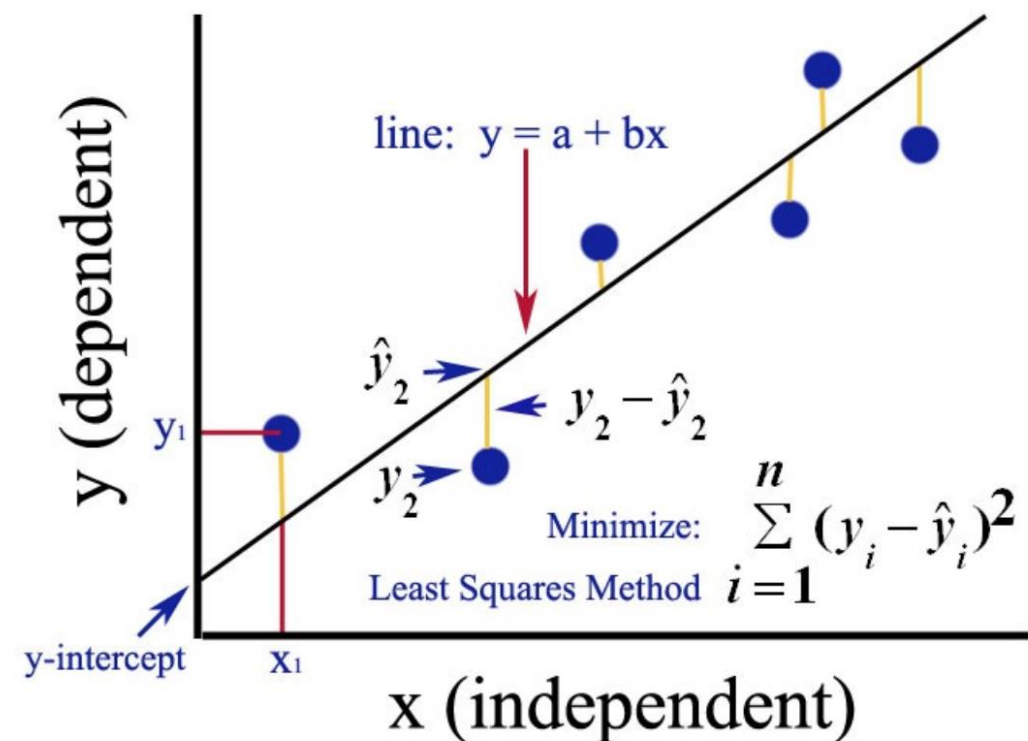- Learn the function which passes through as many points as possible : Minimize the **Least Squares Error**

$X_i = [1, X_{i1}, \ldots, X_{iD}]^\top$

Design matrix $\quad X = \begin{pmatrix} X_{1,} X_{2, \ldots} X_N \end{pmatrix}$

$(D+1) \times N$

$$E(w) = \frac{1}{2} \sum_{i=1}^{N} (y_i - \hat{y}_i)^2$$

$$= \frac{1}{2} \sum_{i=1}^{N} (y_i - X_i^\top w)^2$$

$$\frac{1}{2} \| (y - X^\top w) \|^2$$



line: y = a + bx

$\hat{y}_2$

$y_2 - \hat{y}_2$

$y_1$

$y_2$

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

Least Squares Method

y-intercept $\quad$ X₁

y (dependent)

x (independent)

$$\nabla E(w) = Xy - XX^\top w = 0$$
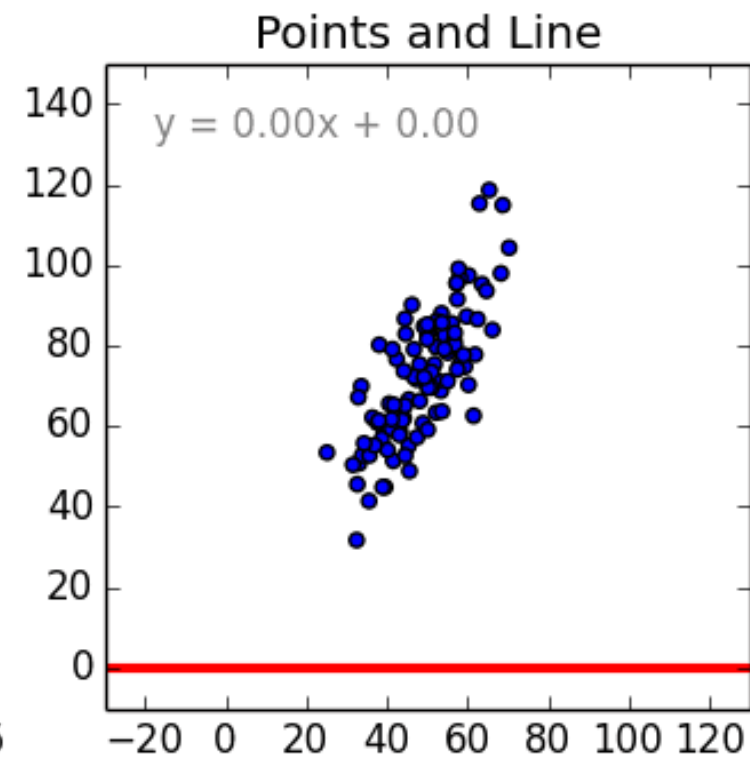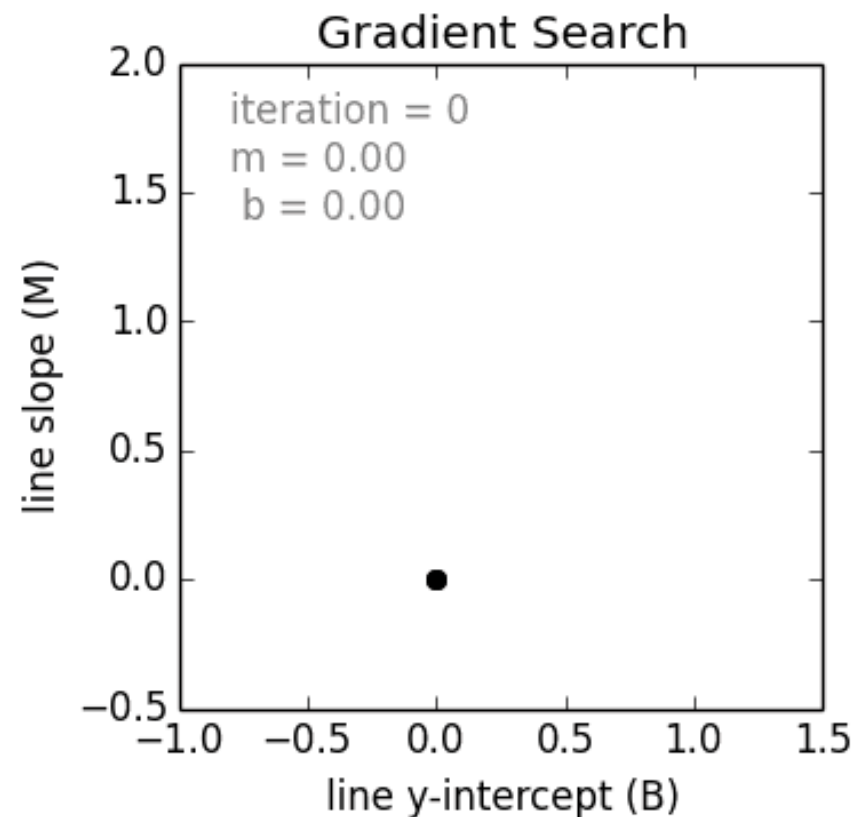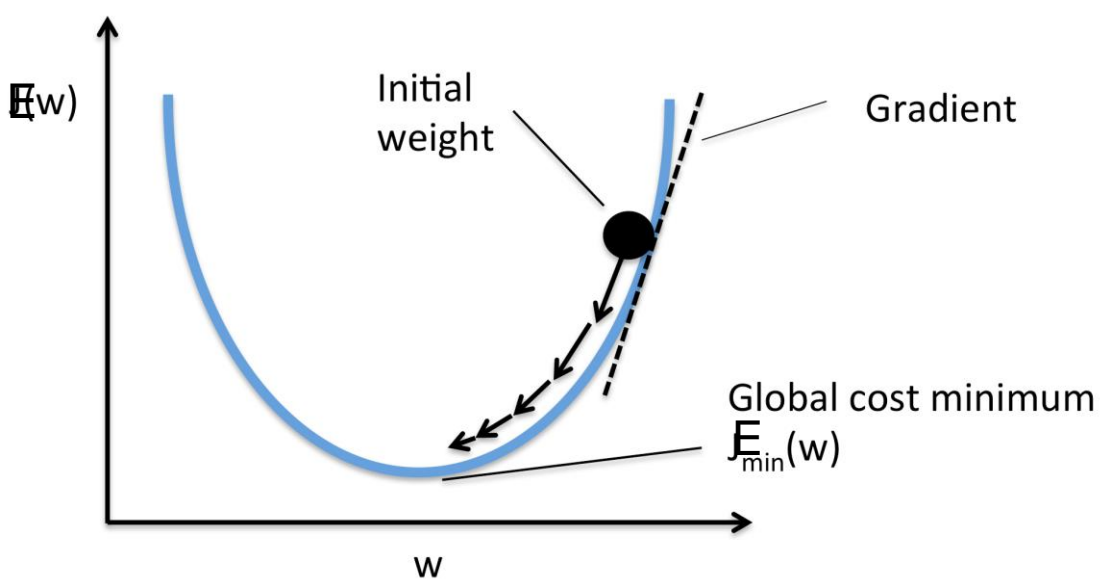
$$\boxed{w_{ML} = (XX^\top)^{-1} Xy}$$

# Linear Regression - Learning Parameters

- Learn the function which passes through as many points as possible
  : Minimize the **least squares error** using **gradient descent**

$$\nabla E(w) = Xy - XX^\top w = 0$$

$$Error_{(m,c)} = \frac{1}{N} \sum_{i=1}^{N} (y_i - (mx_i + c))^2$$

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n$$

# Question

- write down three equations for the line **y = mx+c** to go through **y = 7** at **x = -1**, **y = 7** at **x = 1** and **y = 21** at **x = 2**. Find the least squares solution **(c,m)**?

- Implement In python least squares solution to linear regression
  - Analytical approach
  - Gradient descent approach

# Question
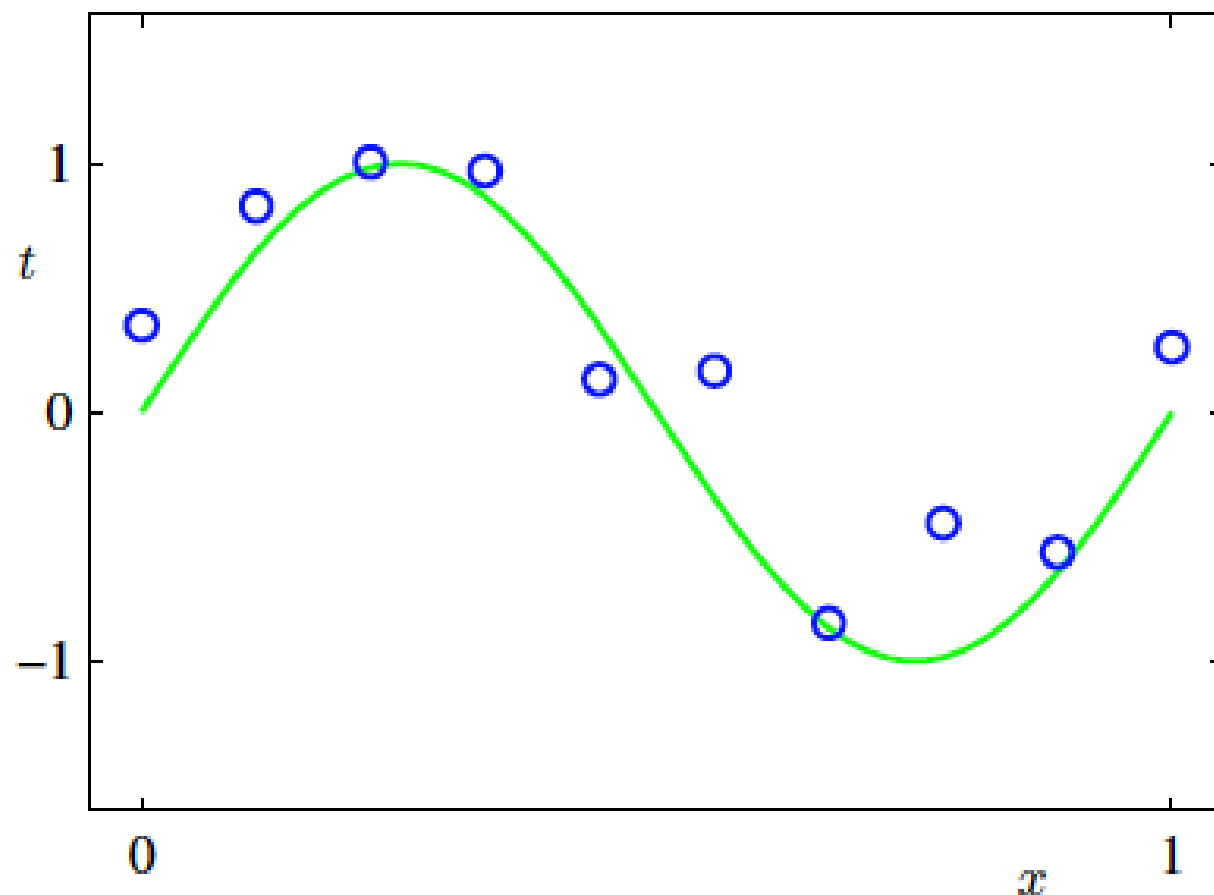
- write down three equations for the line $y = mx+c$ to go through $y = 7$ at $x = -1$, $y = 7$ at $x = 1$ and $y = 21$ at $x = 2$. Find the least squares solution $(c,m)$?
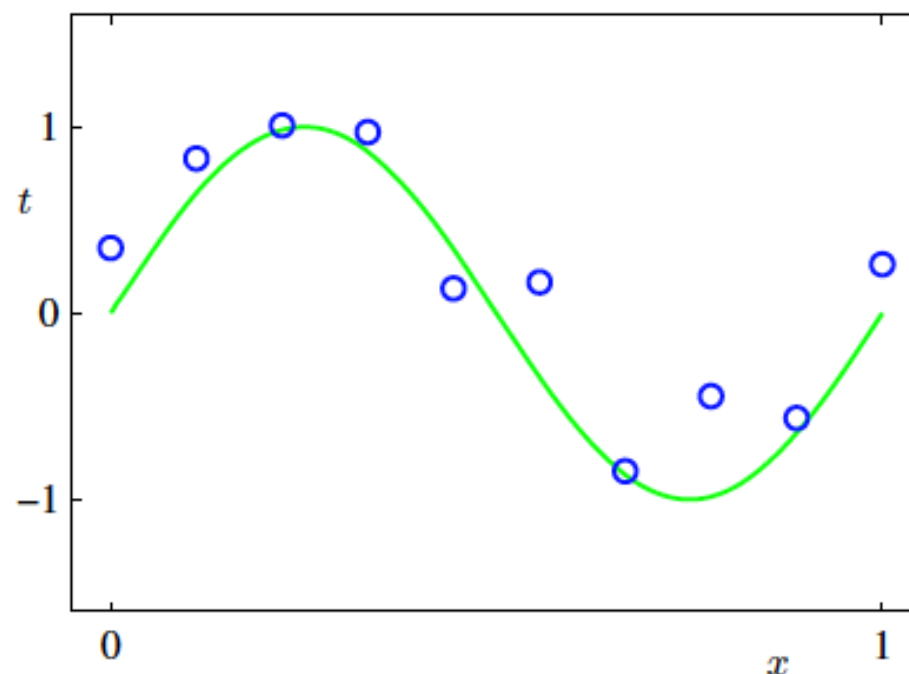
- Answer : (9,4)

# Non Linear Regression - curve fitting

- Remember high school maths !

- Real-valued target variable t.

- Training set comprising N observations

# Regression - curve fitting

- M is the order of the polynomial, y(x,w) is a nonlinear function of x, it is a linear function of the coefficients w.

- Functions, such as the polynomial, which are linear in the unknown parameters have important properties and are called **linear models**

$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \ldots + w_M x^M = \sum_{j=0}^{M} w_j x^j$$
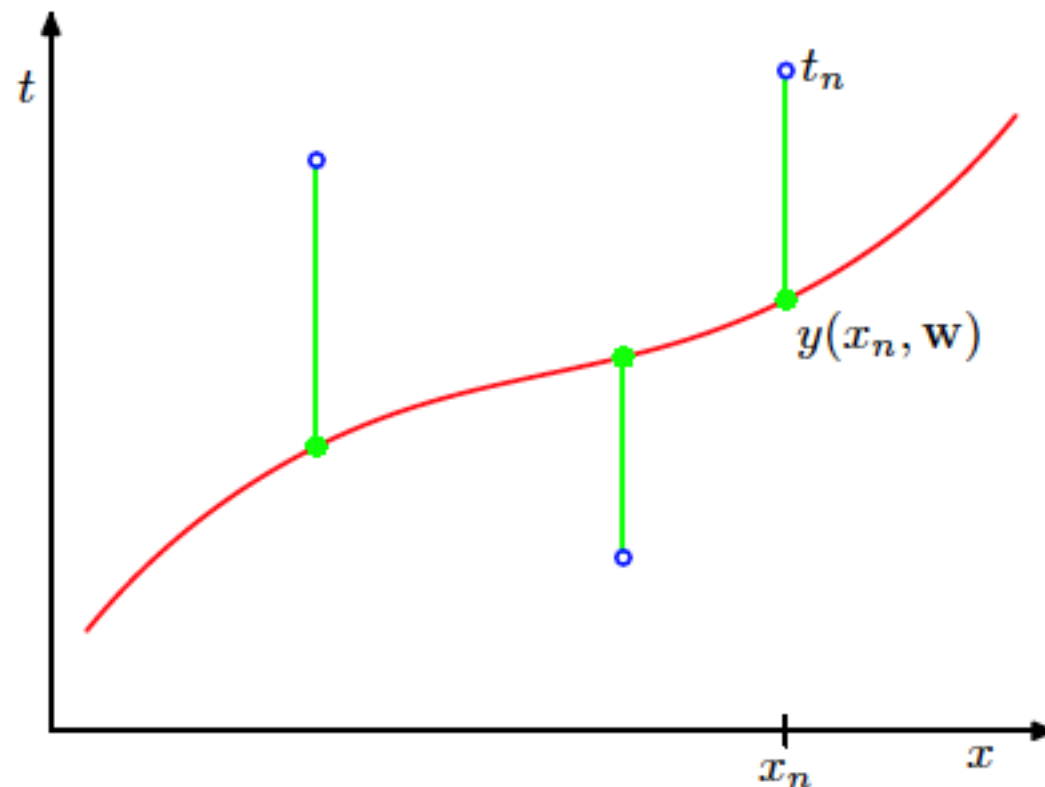
# Regression - curve fitting

$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \ldots + w_M x^M = \sum_{j=0}^{M} w_j x^j$$

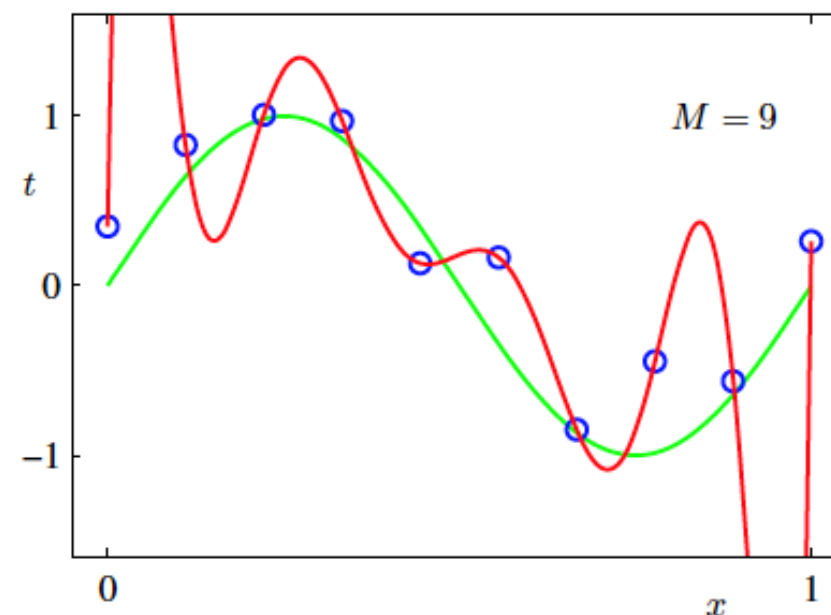- Coefficients will be determined by fitting the polynomial to the training data. This can be done by minimizing an error function

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \mathbf{w}) - t_n\}^2$$

# Regression - curve fitting

- Model selection (choosing M) : higher order polynomial (M = 9), provide excellent fit to the training data but gives a very poor representation of the function

# Regression - curve fitting

- Model selection (choosing M) : hi**Overfitting**omial (M = 9), provide excellent fit to the training ... gives a very poor representation of the function

# Regression - curve fitting

- Generalization performance : root mean square error on test data

- Weights coefficients for M=9 is extremely large !

$$E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^\star)/N}$$

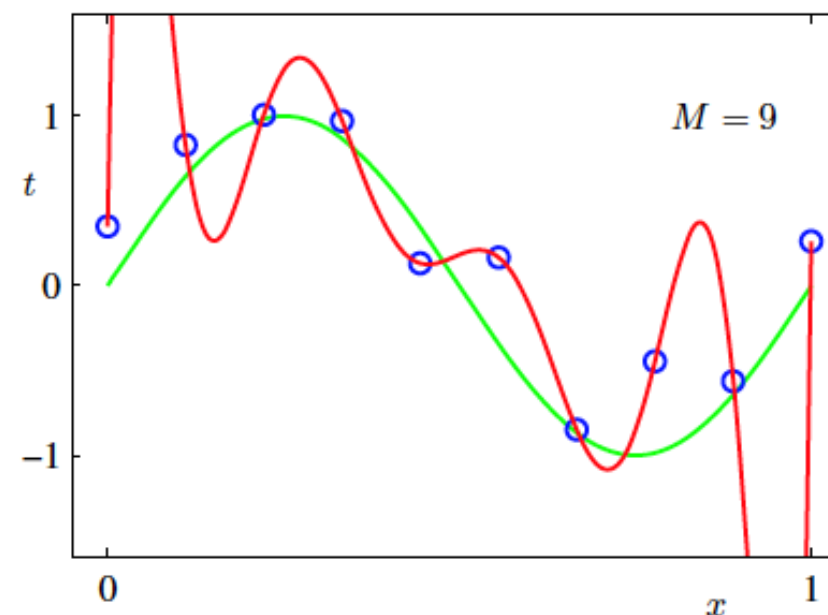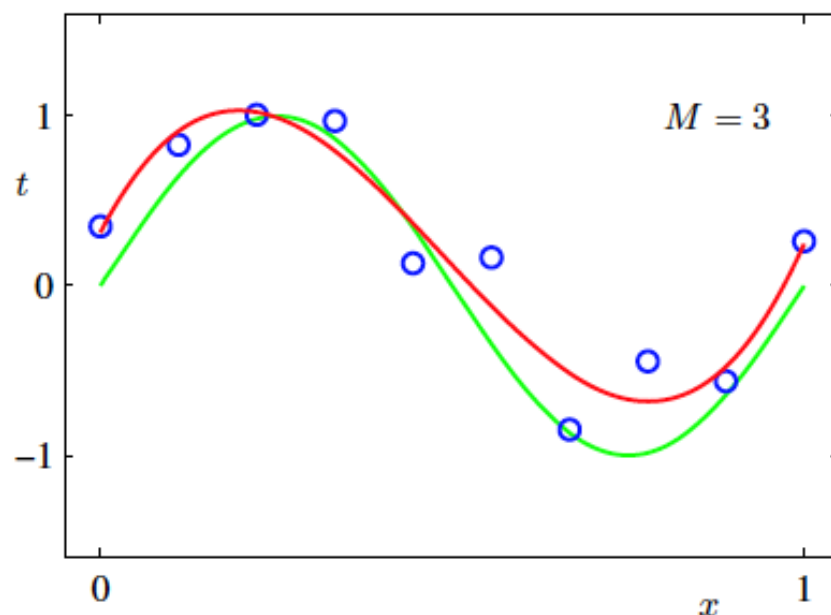$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \mathbf{w}) - t_n\}^2$$

|  | $M = 0$ | $M = 1$ | $M = 6$ | $M = 9$ |
|---|---|---|---|---|
| $w_0^\star$ | 0.19 | 0.82 | 0.31 | 0.35 |
| $w_1^\star$ |  | -1.27 | 7.99 | 232.37 |
| $w_2^\star$ |  |  | -25.43 | -5321.83 |
| $w_3^\star$ |  |  | 17.37 | 48568.31 |
| $w_4^\star$ |  |  |  | -231639.30 |
| $w_5^\star$ |  |  |  | 640042.26 |
| $w_6^\star$ |  |  |  | -1061800.52 |
| $w_7^\star$ |  |  |  | 1042400.18 |
| $w_8^\star$ |  |  |  | -557682.99 |
| $w_9^\star$ |  |  |  | 125201.43 |

# Regression - curve fitting

- Given model complexity, the over-fitting problem become less severe as the size of the data set increases.

# Curve fitting - regularization

- Add a penalty term to the error function (1.2) in order to discourage the coefficients from reaching large values

$$\widetilde{E}(\mathbf{w}) = \frac{1}{2}\sum_{n=1}^{N}\{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2}\|\mathbf{w}\|^2$$

- Ridge regression : L2 norm

$$\|\mathbf{w}\|^2 \equiv \mathbf{w}^{\mathrm{T}}\mathbf{w} = w_0^2 + w_1^2 + \ldots + w_M^2$$

**Regularization constant**

# Curve fitting - regularization



| | $\ln \lambda = -\infty$ | $\ln \lambda = -18$ | $\ln \lambda = 0$ |
|---|---|---|---|
| $w_0^\star$ | 0.35 | 0.35 | 0.13 |
| $w_1^\star$ | 232.37 | 4.74 | -0.05 |
| $w_2^\star$ | -5321.83 | -0.77 | -0.06 |
| $w_3^\star$ | 48568.31 | -31.97 | -0.05 |
| $w_4^\star$ | -231639.30 | -3.89 | -0.03 |
| $w_5^\star$ | 640042.26 | 55.28 | -0.02 |
| $w_6^\star$ | -1061800.52 | 41.32 | -0.01 |
| $w_7^\star$ | 1042400.18 | -45.95 | -0.00 |
| $w_8^\star$ | -557682.99 | -91.53 | 0.00 |
| $w_9^\star$ | 125201.43 | 72.68 | 0.01 |

# Regularized Least Squares

- Parameter shrinkage, weight decay

- **Ridge regression** q=2

$$\frac{\lambda}{2} \sum_{j=1}^{M} w_j^2$$

- **Lasso regression** q=1, if λ is sufficiently large, some of the coefficients are driven to zero

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{t_n - \mathbf{w}^\mathrm{T} \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^{M} |w_j|$$

- Elastic net regularization

$$\frac{1}{n} \|Y - \mathbf{X}\, w\|_2^2 + \lambda_1 \sum_{j=1}^{d} |w_j| + \lambda_2 \sum_{j=1}^{d} |w_j|^2$$

$$\frac{1}{2} \sum_{n=1}^{N} \{t_n - \mathbf{w}^\mathrm{T} \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^{M} |w_j|^q$$



$E(\mathbf{w})$, $w^2$, $|w|$



$q = 0.5$     $q = 1$     $q = 2$     $q = 4$

# Ridge Regression

- Regularized Least Squares
$$\frac{1}{2}\sum_{n=1}^{N}\{t_n - \mathbf{w}^\mathrm{T}\phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2}\mathbf{w}^\mathrm{T}\mathbf{w}.$$

$$\mathbf{\Phi} = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

N × M

- Show that the regularized least squares solution is

$$\mathbf{w} = \left(\lambda\mathbf{I} + \mathbf{\Phi}^\mathrm{T}\mathbf{\Phi}\right)^{-1}\mathbf{\Phi}^\mathrm{T}\mathbf{t}.$$

Stable and unique solution

# Regularized Least Squares : Cross Validation

How to choose  λ ?
Use validation data!

$$\frac{1}{2}\sum_{n=1}^{N}\{t_n - \mathbf{w}^{\mathrm{T}}\phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2}\sum_{j=1}^{M}|w_j|^q$$

| TRAIN | VALIDATION | TEST |
|---|---|---|

1. Training set is a set of examples used for learning a model parameters (e.g.,  weight vector w in linear regression

2. Validation set is a set of examples that cannot be used for learning the model parameter  but can help tune model hyper-parameters e.g Regularization constant  in LR.  Validation helps control overfitting.

3. Test set is used to assess the performance of the final model and provide an estimation of the test error.

Example: Split the data randomly into 60% for training, 20% for validation and 20% for testing.

**Note:  Dont use  the  test  set  to  further  tune**

**the  parameters  or  revise  the  model**.

# Cross Validation



https://scikit-learn.org/stable/modules/cross_validation.html

# Probabilistic Interpretation : Least Squares = Maximum likelihood estimation

$$y_i \sim w \cdot x_i + N(0, \sigma^2) = N(w \cdot x_i, \sigma^2)$$

Likelihood

$$\boxed{y_i | x_i \sim N(X_i^\top w, \sigma^2)}$$

$y_i$

$X_i^\top w$

$x_i$

0

20

$$L = \prod_i \exp -\frac{1}{2\sigma^2}(X_i^\top w - y_i)^2 = \exp -\frac{1}{2\sigma^2}\sum_i (X_i^\top w - y_i)^2$$

$$\underset{w}{\arg\max}\, L = \underset{w}{\arg\min}\, E$$

- Similarly Regularized least squares is same as maximum aposteriori estimate assuming p(w) to be a Gaussian :  argmax  p(yi |w, xi) p(w)

# Supervised Learning : Classification

- Binary classification :  y = {0,1}
- Multiclass classification : y = {1,2,...K}



**p(y|x)?**

# Linear regression to Logistic regression

- Y takes value 0 or 1

# Logistic Regression

- A discriminative approach which directly models p(y|x)
- To predict an outcome variable that is categorical from one or more categorical or continuous predictor variables.

- Let X be the data instance, and Y be the class label {0,1} :
  Model P(Y|X) directly using a Sigmoid function:

**Logistic Sigmoid :** $P(Y = 1 | \mathbf{X}) = \dfrac{1}{1 + e^{-\mathbf{wx}}}$

**wx**

**X**

# Logistic Regression

- To predict an outcome variable that is categorical from one or more categorical or continuous predictor variables.

  - Let X be the data instance, and Y be the class label:
    Model P(Y|X) directly using a Sigmoid function:

$$P(Y = 1 \mid \mathbf{X}) = \frac{1}{1 + e^{-\mathbf{wx}}}$$

**[HW] Find derivative of s(w) = p(y=1|X)!**

# Logistic Regression



$$P(Y = 1 | \mathbf{X}) = \frac{1}{1 + e^{-\mathbf{wx}}}$$

P(y|X)

$X_1$ $X_2$

**Decision boundary**



- Male
- Female

Weight (pounds)

Height (inches)

# Logistic Regression

- In logistic regression, we learn the conditional distribution $P(y|x)$
- Let $p_y(x;w)$ be our estimate of $P(y|x)$, where $w$ is a vector of adjustable parameters.
- Assume there are two classes, $y = 0$ and $y = 1$ and

$$p_1(\mathbf{x};\mathbf{w}) = \frac{1}{1+e^{-\mathbf{wx}}} \qquad p_0(\mathbf{x};\mathbf{w}) = 1 - \frac{1}{1+e^{-\mathbf{wx}}} = \frac{1}{1+e^{\mathbf{wx}}}$$

- This is equivalent to

$$\log \frac{p_1(\mathbf{x};\mathbf{w})}{p_0(\mathbf{x};\mathbf{w})} = \mathbf{wx}$$

- That is, the log odds of class 1 is a linear function of x
- Q: How to find **W**?

- Alternate representation of $p(y|x)$ :     $\mathbf{p_y(x;w)} = \dfrac{1}{1+\mathbf{e^{-ywx}}}; \mathbf{y} = \{-1, 1\}$

# Logistic Regression

- Conditional data likelihood - Probability of observed Y values in the training data, conditioned on corresponding X values.

- We choose parameters w that satisfy

$$\mathbf{w} = \arg\max_{\mathbf{w}} \prod_{l} P(y^l \mid \mathbf{x}^l, \mathbf{w})$$

- where
  - $\mathbf{w} = <w_0, w_1, \ldots, w_n>$ is the vector of parameters to be estimated,
  - $y^l$ denotes the observed value of Y in the $l$ th training example, and
  - $\mathbf{x}^l$ denotes the observed value of $\mathbf{X}$ in the $l$ th training example

# Logistic Regression

- Equivalently, we can work with log of conditional likelihood:

$$\mathbf{w} = \arg\max_{\mathbf{w}} \sum_l \ln P(y^l \mid \mathbf{x}^l, \mathbf{w})$$

- Conditional data log likelihood, l(W), can be written as

$$l(\mathbf{w}) = \sum_l y^l \ln P(y^l = 1 \mid \mathbf{x}^l, \mathbf{w}) + (1 - y^l) \ln P(y^l = 0 \mid \mathbf{x}^l, \mathbf{w})$$

- Note here that Y can take only values 0 or 1, so only one of the two terms in the expression will be non-zero for any given $y^l$

# Logistic Regression

- We need to estimate:

$$\mathbf{w} = \arg\max_{\mathbf{w}} \sum_{l} \ln P(y^l \mid \mathbf{x}^l, \mathbf{w})$$

$$l(\mathbf{w}) = \sum_{l} y^l \ln P(y^l = 1 \mid \mathbf{x}^l, \mathbf{w}) + (1 - y^l) \ln P(y^l = 0 \mid \mathbf{x}^l, \mathbf{w})$$

- Equivalently, we can minimize negative log likelihood using gradient descent technique

- No closed-form solution though. Iterative method required.

- [HW] Find the derivative of l(w) !

# Logistic Regression

- Overfitting can arise especially when data has very high dimensions and is sparse.

- One approach -> modified "penalized log likelihood function," which penalizes large values of **w**, as before.

$$\mathbf{w} = \arg\max_{\mathbf{w}} \sum_{l} \ln P(y^l \mid \mathbf{x}^l, \mathbf{w}) - \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- [HW] Find the Derivative !

# Logistic Regression

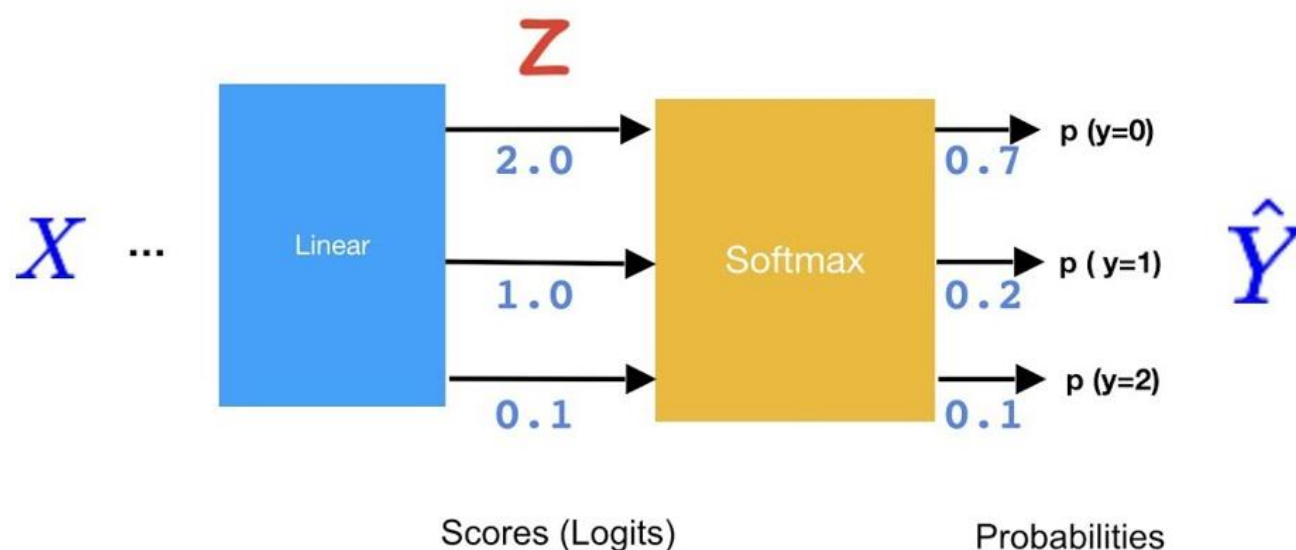- LR: Functional form of P(Y|X), no assumption on P(X|Y)
- LR is a linear classifier
- LR optimized by conditional likelihood
- Extending logistic regression to multiple classes
  - Use softmax for each class k!

$$p(y = k|x) = \frac{\exp(w_k^\top x)}{\sum_{i=1}^{K} \exp(w_i^\top x)}$$

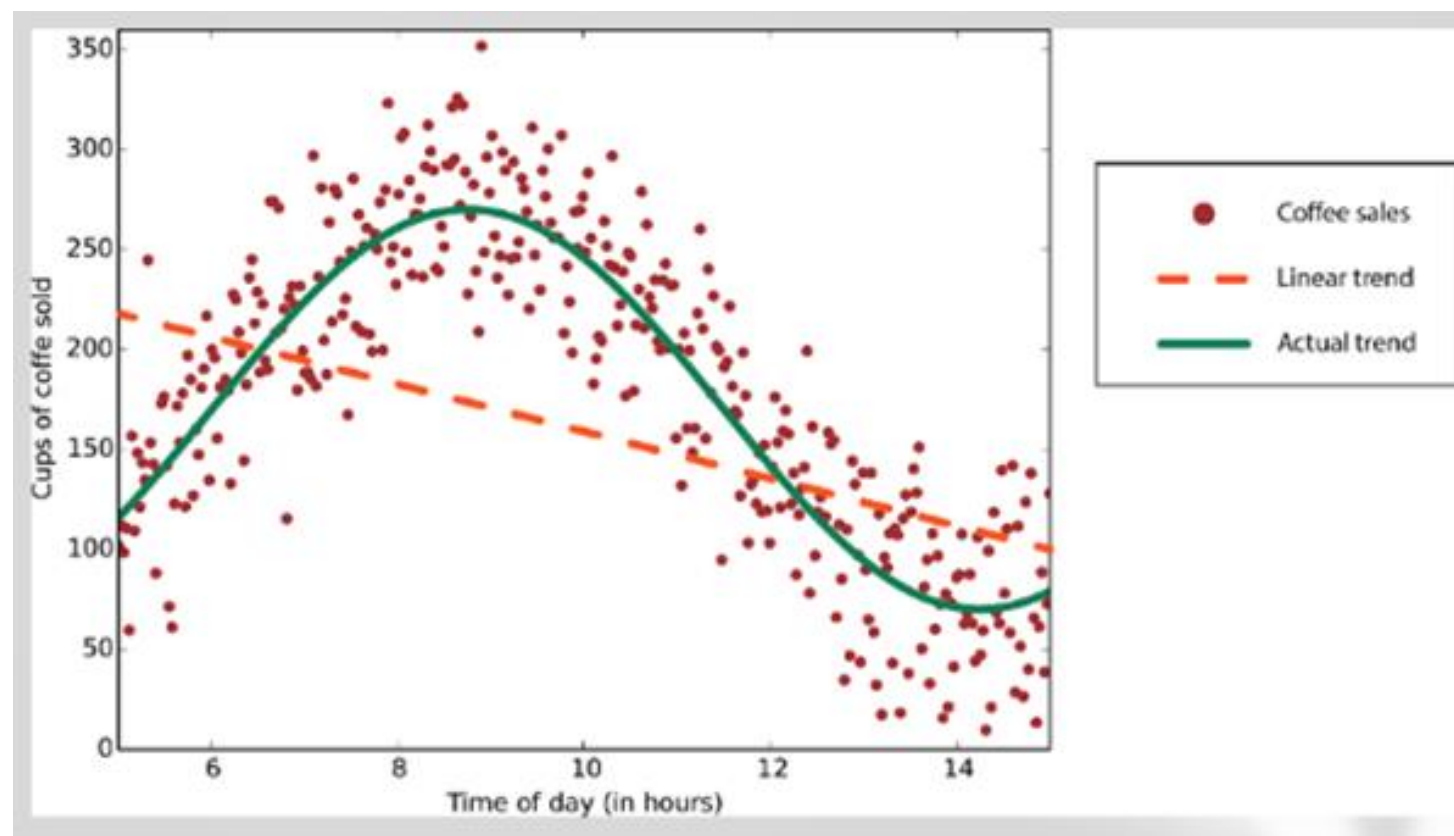Meet Softmax $\qquad \sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^{K} e^{z_k}} \quad$ for $j = 1, ..., K.$

# Classification : Evaluation metrics

**Accuracy  :**   $$\frac{1}{N} \sum_{i=1}^{N} \mathbb{I}[y_i == \hat{y}_i]$$

| | | Actual Label | |
|---|---|---|---|
| | | Positive | Negative |
| **Predicted Label** | Positive | True Positive (TP) | False Positive (FP) |
| | Negative | False Negative (FN) | True Negative (TN) |

| | | |
|---|---|---|
| **Accuracy** | (TP + TN) / (TP + TN + FP + FN) | The percentage of predictions that are correct |
| **Precision** | TP / (TP + FP) | The percentage of positive predictions that are correct |
| **Sensitivity (Recall)** | TP / (TP + FN) | The percentage of positive cases that were predicted as positive |
| **Specificity** | TN / (TN + FP) | The percentage of negative cases that were predicted as negative |

# Supervised learning : Regression



Number of vehicles passing a junction

Y take values 0,1,2,3,.....  but not 2.1, 3.4, 5.55......
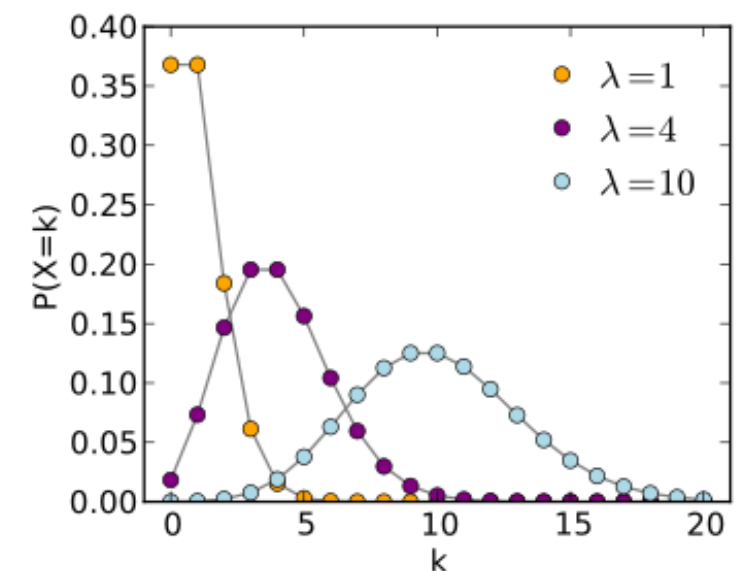
**p(y|x)?**

# Poisson Regression

- Poisson distribution : Model number of events occurring in a fixed interval of time/space

$$P(k \text{ events in interval}) = e^{-\lambda} \frac{\lambda^k}{k!}$$



- λ is the average (mean) number of events

- Y has a Poisson distribution, and assumes the logarithm of its expected value can be modeled by a linear combination of unknown parameters.

$$\lambda := \mathrm{E}(Y \mid x) = e^{\theta' x},$$

$$p(y \mid x; \theta) = \frac{\lambda^y}{y!} e^{-\lambda} = \frac{e^{y\theta' x} e^{-e^{\theta' x}}}{y!}$$

# Poisson Regression : Learning parameters

- Likelihood

$$p(y_1, \ldots, y_m \mid x_1, \ldots, x_m; \theta) = \prod_{i=1}^{m} \frac{e^{y_i \theta' x_i} e^{-e^{\theta' x_i}}}{y_i!}.$$

- Estimate parameters by maximum likelihood estimation

$$\ell(\theta \mid X, Y) = \log L(\theta \mid X, Y) = \sum_{i=1}^{m} \left( y_i \theta' x_i - e^{\theta' x_i} - \log(y_i!) \right).$$

- Use gradient descent to find the optimal value of θ.

# Thank you !

**Reference**

**[1] Christopher Bishop, Pattern Recognition and Machine Learning**
**[2] Kevin Murphy, Machine Learning : A Probabilistic Perspective**