# Introduction

Machine Learning Fundamentals

Simple Cells: Response to light orientation

Complex Cells: Response to light orientation & movement

Hypercomplex Cells: Response to movement with an end point

No response    Response (end point)

Electrical signal from brain

Recording electrode

Visual area of brain

Stimulus

Stimulus    http://hubel.med.harvard.edu/book/70.jpg

Cell response

Hubel & Wiesel, 1959

# Machine Learning

Field of study that gives computers the ability to learn without being explicitly programmed.

By: Arthur Samuel (1959)

# Applications

- Database mining:
    - Web click data, medical records, biology, engineering etc.

- Application that can't be programmed by hand
    - Autonomous helicopter, handwriting recognition, NLP, computer vision

- Self Customizing Programs Recommendation
    - Amazon, Netflix

- Understanding human learning
    - Brain, real AI

# Learning Problem

- Well-posed learning problem:

  A computer program is said to *learn* from *experience* E with respect to some *task* T and some *performance measure* P, if its performance on T, as measured by P, improves with *experience* E.

  - By Tome Mitchell (1998)

# Example

Your Email program watches which emails you do or do not mark as a spam and based on that learns how to better filter spam.
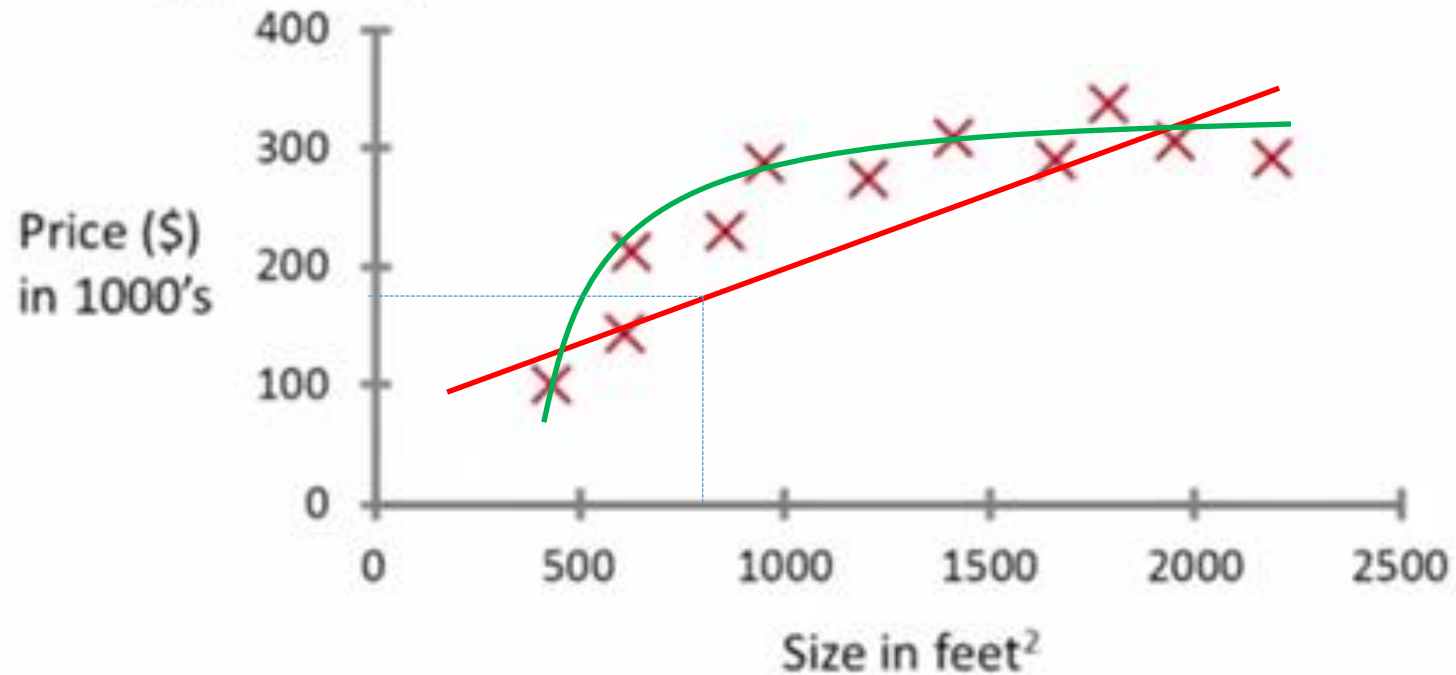
- Task (T): Classifying the emails as spam or not

- Experience (E): Watching you label emails as spam or not spam

- Performance (P): The number of emails correctly classified as spam / not spam

# Machine Learning Algorithms

- Supervised Learning
- Un-supervised Learning

- Others: Reinforcement learning, recommender system

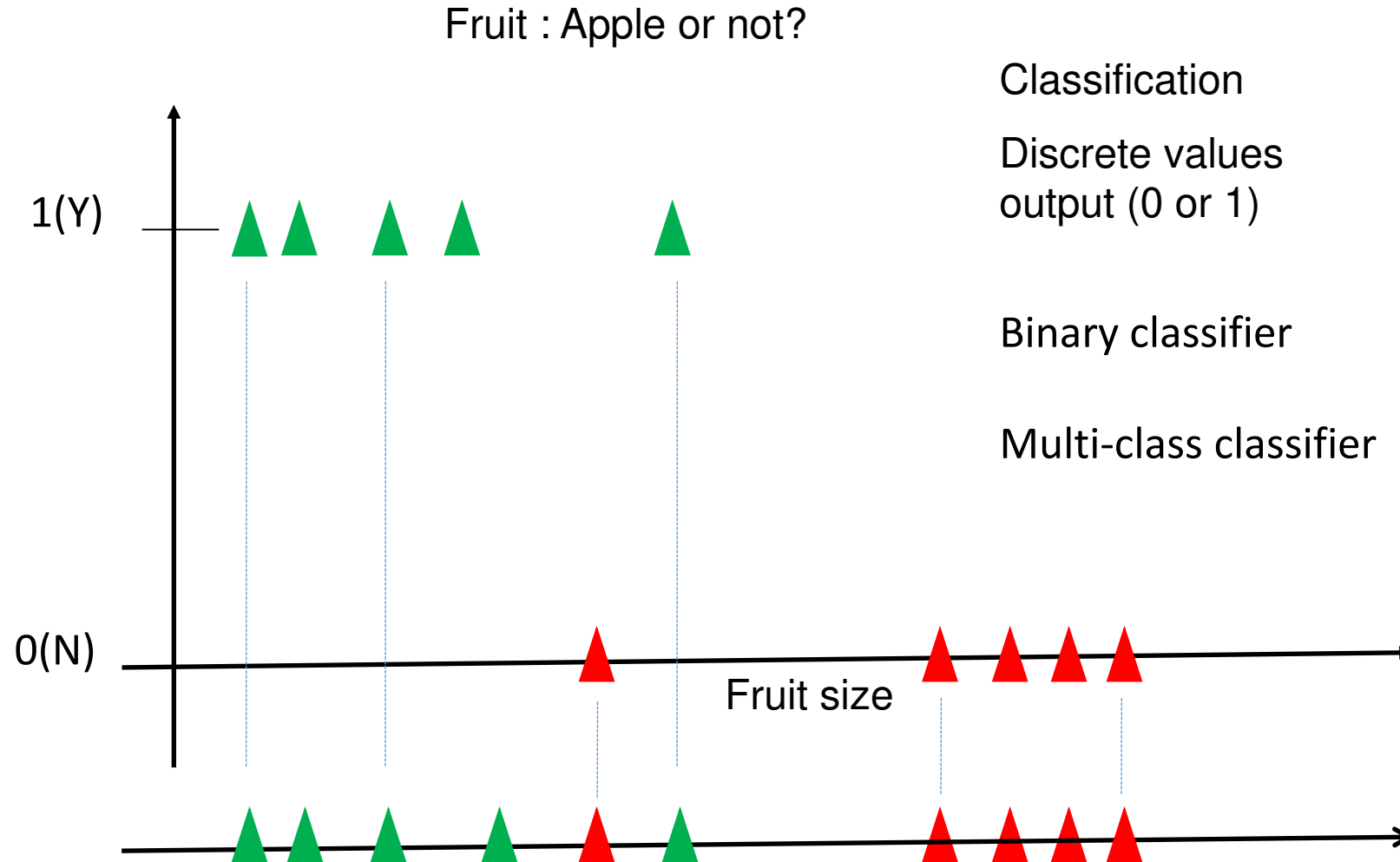- Where to apply which algorithms

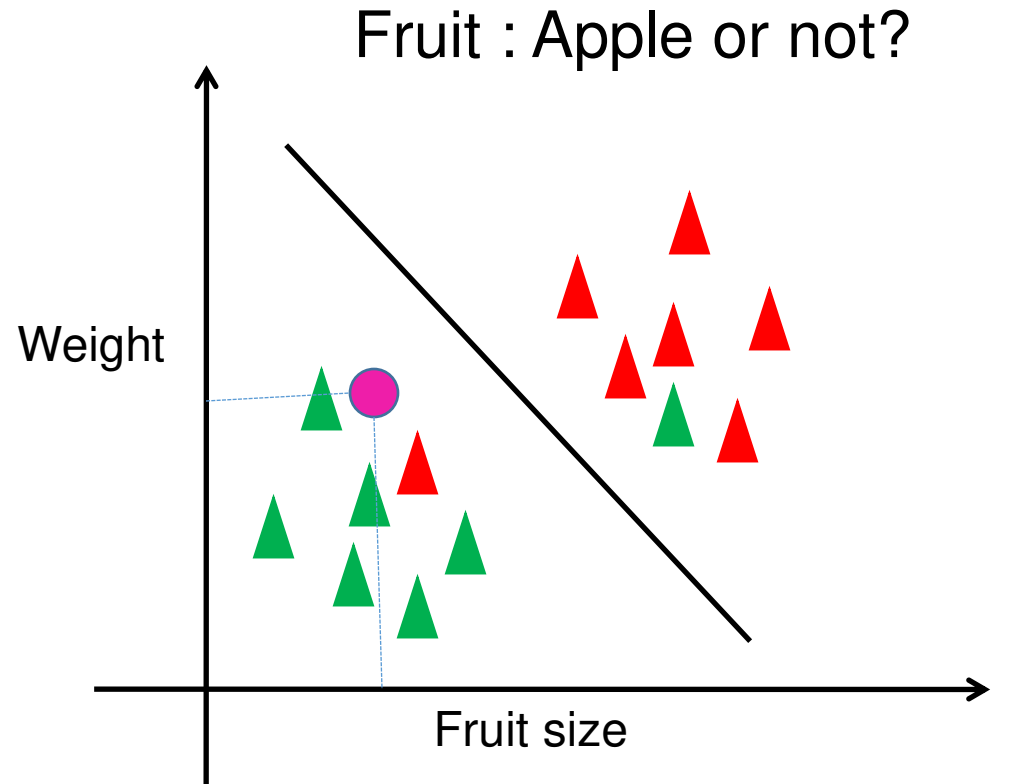# Supervised Learning



Housing price prediction.

Supervised Learning:
    Right answer given

Regression:
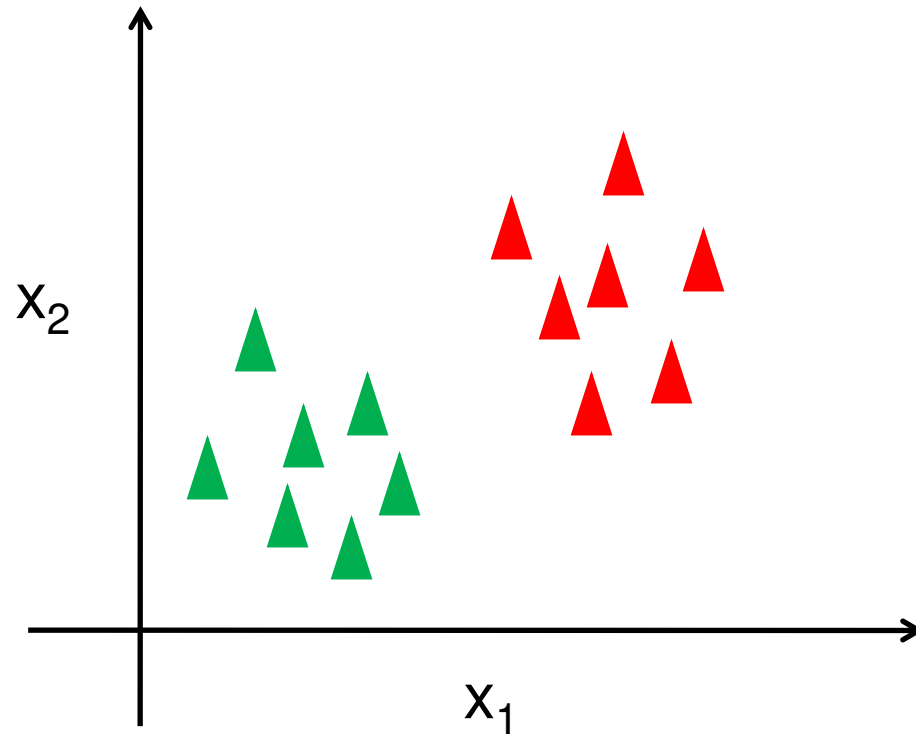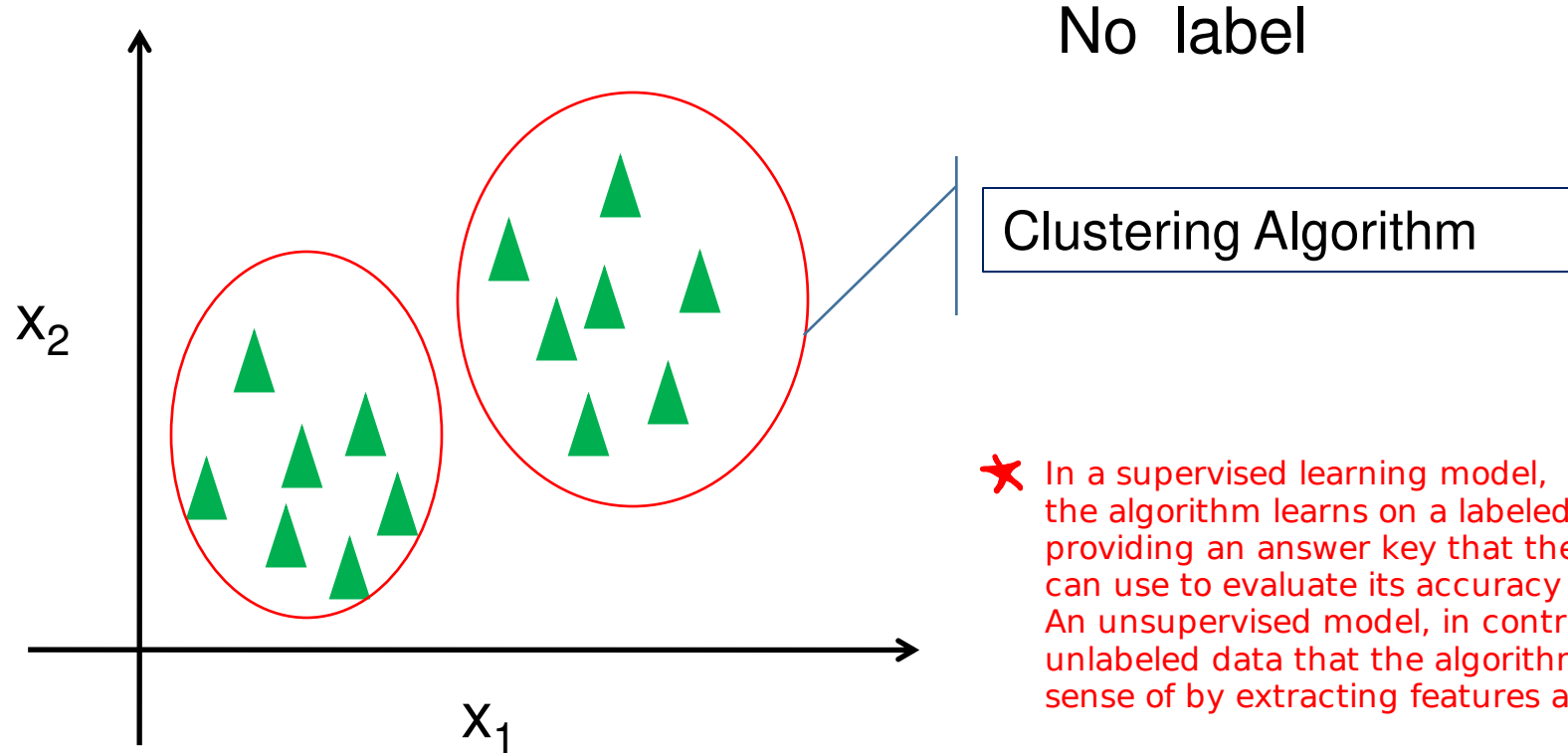Predict continuous valued output

# Supervised Learning

Fruit : Apple or not?

Classification

Discrete values output (0 or 1)

Binary classifier

Multi-class classifier

1(Y)

0(N)

Fruit size

# Supervised Learning

Fruit : Apple or not?

Weight

Fruit size

Two features examples

- Features:

  - color
  - texture
  - cost
  - …..

# Supervised Learning

# Un-supervised Learning

No  label

$x_2$

Clustering Algorithm

⭐ In a supervised learning model,
the algorithm learns on a labeled dataset,
providing an answer key that the algorithm
can use to evaluate its accuracy on training data.
An unsupervised model, in contrast, provides
unlabeled data that the algorithm tries to make
sense of by extracting features and patterns on its own.

$x_1$

CS 590 Lecture 2

# Application of Clustering Algorithms

- Organizing computer clusters

- Social network analysis

- Market segmentation

- Astronomical image/data analysis

- Speaker recognition and many more…

# Supervised Learning

# Supervised Learning



Price
In 1000's of
dollars

Size (feet$^2$)

- Given the right answer for each example of the data
  - Classification: discrete no. of outputs
  - Regression: Predict real valued data

# Supervised Learning

| | Size in feet² (x) | Price ($) in 1000's (y) |
|---|---|---|
| **Training set of housing prices** | 2104 | 460 |
| | 1416 | 232 |
| | 1534 | 315 |
| | 852 | 178 |
| | ... | ... |

Notation:

**m** = Number of training examples

**x**'s = "input" variable / features

**y**'s = "output" variable / "target" variable

$(x,y)$ → one training example      $x^{(i)} = 2104$

$(x^{(i)}, y^{(i)})$ → $i^{th}$ training example      $y^{(i)} = 460$

# Supervised Learning



Training Set → Learning Algorithm

Size of house → h → Estimated price

x    hypothesis    Estimated value of y

How do we represent h

$h_\theta(x) = h(x) = \theta_0 + \theta_1 x$

$h(x) = \theta_0 + \theta_1 x$

Univariate linear regression: linear regression with one variable

# Cost Function

Training Set

| Size in feet$^2$ (x) | Price ($) in 1000's (y) |
|---|---|
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| ... | ... |

Hypothesis:  $h_\theta(x) = \theta_0 + \theta_1 x$

$\theta_i$'s  $\rightarrow$ Parameters

How to choose $\theta_i$'s

# Cost Function

Hypothesis Function:

$$h_\theta(x) = \theta_0 + \theta_1 x$$



h(x)= 1.5 + 0.x

$\theta_0 = 1.5$

$\theta_1 = 0$

h(x)= 0 + 0.5x

$\theta_0 = 0$

$\theta_1 = 0.5$

h(x)= 1+ 0.5x

$\theta_0 = 1$

$\theta_1 = 0.5$

# Cost Function



**Hypothesis:** $h_\theta(x) = \theta_0 + \theta_1 x$

**Parameters:** $\theta_0, \theta_1$

**Cost Function:** $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$

Squared error function

**Goal:** $\underset{\theta_0, \theta_1}{\text{minimize}} \, J(\theta_0, \theta_1)$

m = No. of training samples

Idea: Choose $\theta_0, \theta_1$ so that $h_\theta(x)$ is close to $y$ for our training examples $(x, y)$

# Cost Function

$h_\theta(x)$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)



$J(\theta_0, \theta_1)$= value of the height of the surface

# Contour Plots / Figures

$h_\theta(x)$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)



J($\theta_0$, $\theta_1$)

$(\theta_0, \theta_1) = (800, - 0.125)$

CS 590 Lecture 3

$$h_\theta(x)$$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0, \theta_1$)



$(\theta_0, \theta_1) = (360, 0)$

CS 590 Lecture 3

$h_\theta(x)$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)

$$h_\theta(x)$$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0, \theta_1$)

# Gradient Descent
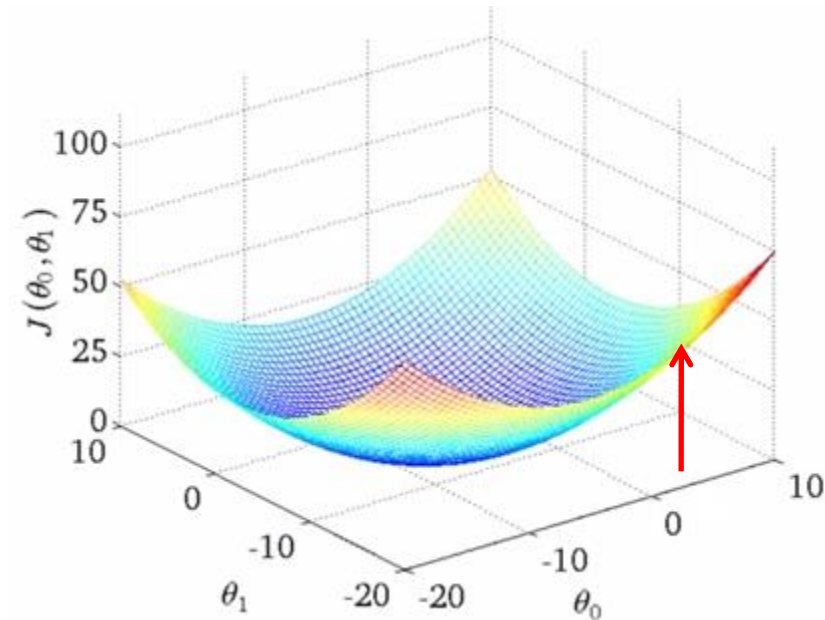
- Let some function $J(\theta_0, \theta_1)$

- We have to find $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

- Start with some $(\theta_0, \theta_1)$ (let say $\theta_0 = 0$, $\theta_1 = 0$)

- Keep changing $\theta_0$, $\theta_1$ to reduce $J(\theta_0, \theta_1)$
  until we hopefully end up at a minimum

# Gradient Descent

# Gradient Descent

# Gradient Descent



$$J(\theta_0, \theta_1)$$

# Gradient Descent

# Gradient Descent Algorithm

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \qquad (\text{for } j = 0 \text{ and } j = 1)$$

}

$\alpha$ = learning rate

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \qquad j \in [0, 1]$$

Implication of $\alpha$ = it controls how bigger steps we are taking over gradient descent

✓ Correct: Simultaneous update    ✗ Incorrect:

$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$
$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$
$$\theta_0 := \text{temp0}$$
$$\theta_1 := \text{temp1}$$

$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$
$$\theta_0 := \text{temp0}$$
$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$
$$\theta_1 := \text{temp1}$$

# Gradient Descent Algorithm

- Let take a single variable

- we have to minimize $\min\limits_{\theta_1} J(\theta_1)$

  where $\theta_1 \in R$

- So the GD algorithm becomes

$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

# Gradient Descent Algorithm



$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

≥ 0
(slope is positive)

Local minima ⟵ $\theta_1$ is reduced

$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

≤ 0
(slope is negative)

$\theta_1$ ⟶ Local minima

# Gradient Descent Algorithm

$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

If α is too small, gradient descent can be slow.

If α is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.



Local minima ←

Local minima

$\theta_1$ →

# Multivariate Linear Regression

Univariate Hypothesis function:

$$h_\theta(x) = \theta_0 + \theta_1 x$$

Multivariate Hypothesis function:

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$$

$$X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\theta^T x = \begin{bmatrix} \theta_0 & \theta_1 & \dots & \theta_n \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

where $x_0 = 1$

$$\boxed{h_\theta(x) = \theta^T x}$$

# Multivariate Gradient Descent

Hypothesis: $h_\theta(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$

Parameters: $\theta_0, \theta_1, \ldots, \theta_n$ → **ɵ** : n+1 dimensional vector

Cost function:

$$J(\theta_0, \theta_1, \ldots, \theta_n) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$$

**J(ɵ)**

Gradient descent:

Repeat $\{$

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \ldots, \theta_n)$$

$\}$        (simultaneously update for every $j = 0, \ldots, n$)

**J(ɵ)**

# Multivariate Gradient Descent

$$J(\theta) = \frac{1}{2m} \sum (h_\theta(x^i) - y^i)^2$$

**Gradient Descent**

Previously (n=1):

Repeat $\{$

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad}_{\frac{\partial}{\partial \theta_0} J(\theta)}$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x^{(i)}$$

(simultaneously update $\theta_0, \theta_1$)

$\}$

New algorithm $(n \geq 1)$:

$$\frac{\partial}{\partial \theta_j} J(\theta)$$

Repeat $\{$

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^{m} \boxed{(h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)}}$$

(simultaneously update $\theta_j$ for $j = 0, \ldots, n$)

$\}$

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x_0^{(i)}$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x_1^{(i)}$$

$$\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x_2^{(i)}$$

# … to continue

# ✱ Feature Scaling

**Feature Scaling**
Idea: Make sure features are on a similar scale.    $x_1 = \frac{\text{size (feet}^2)}{2000}$

E.g. $x_1$ = size (0-2000 feet$^2$)

   $x_2$ = number of bedrooms (1-5)    $x_2 = \frac{\text{number of bedrooms}}{5}$

$$0 \leq x_1 \leq 1 \qquad 0 \leq x_2 \leq 1$$

θ₂                                    J(θ)          θ₂                                    J(θ)

θ₁                                                  θ₁

Get every feature into approximately a $-1 \leq x_i \leq 1$ range.

# Logistic Regression: Classification

$h_\theta(x) = \theta^T x$



$h_\theta(x) = \theta^T x$

Threshold classifier output $h_\theta(x)$ at 0.5:

- If $h_\theta(x) \geq 0.5$, predict "y = 1"

If $h_\theta(x) < 0.5$, predict "y = 0"

# Logistic Regression



Linear regression for classification problem is not always good

Classification:  y  =  0  or  1

$h_\theta(x)$ can be > 1 or < 0

Logistic Regression:  $0 \leq h_\theta(x) \leq 1$

CS 590 Lecture 3

# Logistic Regression Model

Logistic Regression:    $0 \leq h_\theta(x) \leq 1$

Linear Regression:    $h_\theta(x) = \theta^T x$

Logistic  Regression:

$h_\theta(x) = {\color{red} g}(\theta^T x)$          $h_\theta(x) = \dfrac{1}{1 + e^{-\theta^T x}}$

$* \quad g(z) = \dfrac{1}{1 + e^{-z}}$



Sigmoid Function or Logistic function

# Hypothesis Representation

$h_\Theta(x)$ = estimated probability that y=1 on input x

Example: if $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ Object\ Size \end{bmatrix}$

$h_\Theta(x)$ = 0.7
There is 70% chance that the object is salient

$h_\Theta(x)$ = p(y=1|x, Ө)
i.e. "probability that y=1, given x, parameterized by Ө"

p(y=0|x; Ө) + p(y=1|x; Ө) = 1

p(y=0|x; Ө) = 1 - p(y=1|x; Ө)

# Decision Boundary

**Logistic regression**

$$h_\theta(x) = g(\theta^T x)$$

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



Suppose predict "$y = 1$" if $h_\theta(x) \geq 0.5$

$$g(z) \geq 0.5 \quad when \quad z \geq 0$$

i.e.  $\theta^T x \geq 0$

predict "$y = 0$" if $h_\theta(x) < 0.5$

$$h_\theta(x) = g(\theta^T x) \geq 0.5$$

whenever  $\theta^T x \geq 0$

$$h_\theta(x) = g(\theta^T x)$$

i.e.  $\theta^T x < 0$

$$z$$

# Decision Boundary



$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

$x_1 + x_2 = 3$

**Decision Boundary**

$$h_\theta(x) = g(\theta^T x) = 0.5$$

Predict "$y = 1$" if $-3 + x_1 + x_2 \geq 0$

Predict "y = 0" if

i.e. $\quad \theta^T x \geq 0$

$x_1 + x_2 < 3$

$x_1 + x_2 \geq 3$

✱ <u>Decision boundary is a property of hypothesis function</u> NOT of a data set

# Non-Linear Decision Boundary

$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

Let $\theta^T = \begin{bmatrix} -1 & 0 & 0 & 1 & 1 \end{bmatrix}$

Predict "$y = 1$" if $-1 + x_1^2 + x_2^2 \geq 0$

$$x_1^2 + x_2^2 \geq 1$$

$$x_1^2 + x_2^2 = 1$$

Decision Boundary

Again, decision boundary is a property of hypothesis function NOT of a data set

# Cost Function

- Optimization objective of the cost function

Training set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \cdots, (x^{(m)}, y^{(m)})\}$

m examples $\qquad x \in \begin{bmatrix} x_0 \\ x_1 \\ \dots \\ x_n \end{bmatrix} \qquad x_0 = 1, y \in \{0, 1\}$

$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

How to choose parameters $\theta$ ?

# Cost Function

**Cost function**

✱ Linear regression:  $J(\theta) = \frac{1}{m} \sum\limits_{i=1}^{m} \frac{1}{2} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$

Let,   $\text{Cost}\left( h_\theta(x^{(i)}), y^{(i)} \right) = \frac{1}{2} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$

So, $J(\theta) = \frac{1}{m} \sum\limits_{i=1}^{m} \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$

where, for logistic regression

$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

# Cost Function



**Logistic Regression**

Non convex  **J(Ɵ)**

**Ɵ**

**Linear Regression**

Convex  **J(Ɵ)**

**Ɵ**

# Cost Function: Logistic Regression

$$\star \quad \text{Cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y = 1 \\ -\log(1 - h_\theta(x)) & \text{if } y = 0 \end{cases}$$

If  y = 1

$\longrightarrow$  Cost $= 0$ if $y = 1, h_\theta(x) = 1$
But as $\quad h_\theta(x) \to 0$
$\quad\quad\quad$ $Cost \to \infty$

$\longrightarrow$ Captures intuition that if $h_\theta(x) = 0$,
(predict $P(y = 1|x; \theta) = 0$), but $y = 1$,
we'll penalize learning algorithm by a very
large cost.

Cost

0

$h_\theta(x)$

1

# Cost Function: Logistic Regression

$$\text{Cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y = 1 \\ -\log(1 - h_\theta(x)) & \text{if } y = 0 \end{cases}$$

If  y = 0

Cost

0          $h_\theta(x)$          1

⟶ Cost =0 if y=0, $h_\Theta(x) = 0$
But as $h_\Theta(x) \rightarrow 1$
    Cost → ∞
Captures intuition that if $h_\Theta(x) = 1$,
(predict P(y=0|x; Θ) = 1), but y = 0,
We will penalize learning algorithm
by a very large cost.

✱ It can be shown that the overall cost function is convex function and local optimum free. But details of such convexity analysis is beyond of the scope of this course.

# Cost Function: Logistic Regression

**Logistic regression cost function**

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$\text{Cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y = 1 \\ -\log(1 - h_\theta(x)) & \text{if } y = 0 \end{cases}$$
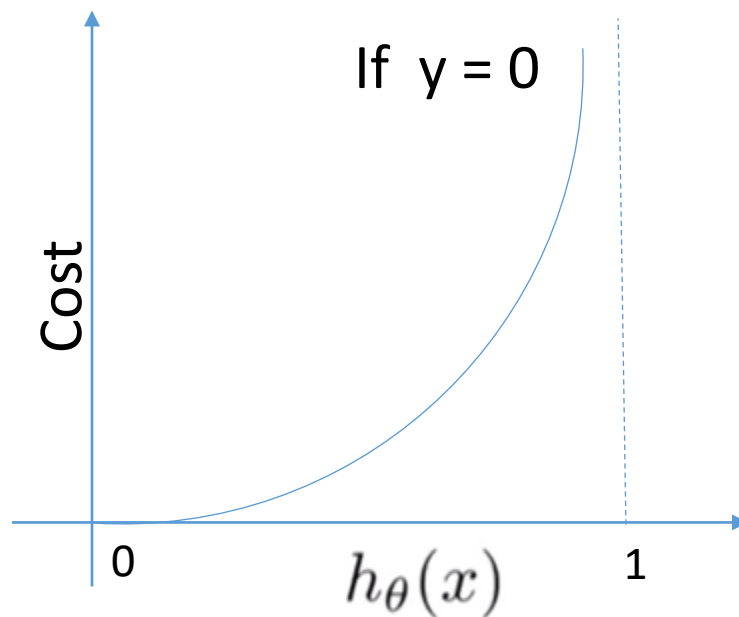
Note: $y = 0$ or $1$ always

$$\text{Cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1 - y)\log(1 - h_\theta(x))$$

$$\text{If } y = 1: \ \text{Cost}(h_\theta(x), y) = -\log(h_\theta(x)))$$

$$\text{If } y = 0: \ \text{Cost}(h_\theta(x), y) = -\log(1 - h_\theta(x))$$

# Cost Function: Logistic Regression

**Logistic regression cost function**

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$\longrightarrow \; J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^{m} y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)})) \right]$$

* Principle of Maximum Likelihood Estimation

To fit parameters $\theta$ :

Obtain $\min_\theta J(\theta)$

and get $\theta$

To make a prediction given new $x$:

Output: $h_\theta(x) = \dfrac{1}{1 + e^{-\theta^T x}}$

For $p(y=1|x; \theta)$

**How to minimize J(θ) ?**

# Cost Function and Gradient Descent

**Gradient Descent**

$$J(\theta) = -\frac{1}{m}[\sum_{i=1}^{m} y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_\theta(x^{(i)}))]$$

Want $\min_\theta J(\theta)$:

  Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

                                 (simultaneously update all $\theta_j$)

  }

$$\frac{\partial}{\partial \theta_j} J(\Theta) = \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

# Cost Function and Gradient Descent

**Gradient Descent**

$$J(\theta) = -\frac{1}{m}\left[\sum_{i=1}^{m} y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))\right]$$

Want $\min_\theta J(\theta)$:

Repeat {

$$\theta_j := \theta_j - \alpha \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)}$$

(simultaneously update all $\theta_j$)

}

For Linear Regression:
$$h_\theta(\text{x}) = \theta^T x$$

For Logistic Regression:
$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

Algorithm looks identical to linear regression!

# Gradient descent optimization

- Problems:
  - Choosing step size
    - too small $\rightarrow$ convergence is slow and inefficient
    - too large $\rightarrow$ may not converge
  - Can get stuck on "flat" areas of function
  - Easily trapped in local minima

# Stochastic gradient descent

Stochastic (definition):

1. involving a random variable
2. involving chance or probability; probabilistic

# Stochastic gradient descent

- Application to training a machine learning model:
    1. Choose one sample from training set
    2. Calculate loss function for that single sample
    3. Calculate gradient from loss function
    4. Update model parameters a single step based on gradient and learning rate
    5. Repeat from 1) until stopping criterion is satisfied

- Typically entire training set is processed multiple times before stopping.

- Order in which samples are processed can be fixed or random.

# Multi Class Classification

# One vs. All (One vs. Rest)

**One-vs-all (one-vs-rest):**



Class 1: △
Class 2: □
Class 3: ✗

$$h_\theta^{(i)}(x) = P(y = i | x; \theta) \qquad (i = 1, 2, 3)$$

$h_\theta^{(1)}(x)$

$h_\theta^{(2)}(x)$

$h_\theta^{(3)}(x)$

# One vs. All (One vs. Rest)

★ Train a logistic regression classifier $h_\theta^{(i)}(x)$ for each class $i$ to predict the probability that $y = i$.

★ On a new input $x$, to make a prediction, pick the class $i$ that maximizes

$$\max_i h_\theta^{(i)}(x)$$

# Overfitting

- A hypothesis function h is said to overfit the training data if there is another hypothesis h' such that h' has more error than h on training data but h' has less error than h on testing data.

- Learning a classifier that classifies a training data perfectly may not lead to the classifier with best generalization performance
    - There may be noise in training data
    - Training data set is too small

- Simplistically, overfitting occurs when model is too complex whether underfitting occurs when model is too simple.

- Note: Training error is not a good predictor for the testing error.

# The problem of overfitting

Example: Linear regression (housing prices)



$$\theta_0 + \theta_1 x$$

$$\theta_0 + \theta_1 x + \theta_2 x^2$$

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

✱ Under fit or High bias

Low varience

Over fit or High variance

✱ **Overfitting:** If we have too many features, the learned hypothesis may fit the training set very well ($J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 \approx 0$), but fail to generalize to new examples (predict prices on new examples).

# The problem of overfitting

Example: Logistic regression



$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$

( $g$ = sigmoid function)

$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2$
$+\theta_3 x_1^2 + \theta_4 x_2^2$
$+\theta_5 x_1 x_2)$

$g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2$
$+\theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2$
$+\theta_5 x_1^2 x_2^3 + \theta_6 x_1^3 x_2 + \ldots)$

Under fit or High bias

Over fit or High variance

# The problem of overfitting

- **Let's consider D, the entire distribution of data, and T, the training set.**

- **Hypothesis h $\in$ H overfits D if**

  **$\exists$ h'$\neq$ h $\in$ H such that**

  **(1) error$_T$(h) < error$_T$(h')** [i.e. doing well on training set] but

  **(2) error$_D$(h) > error$_D$(h')**

- **What do we care about most (1) or (2)?**

- *Estimate error on full distribution by using test data set.*
  *Error on test data: Generalization error (want it low!!)*

- *Generalization to unseen examples/data is what we care about.*

# The problem of overfitting

- Data overfitting is the arguably the most common pitfall in machine learning.

- **Why?**

- Temptation to use as much data as possible to train on. ("Ignore test till end." Test set too small.) Data "peeking" not noticed.

- *Temptation to fit very complex hypothesis* (e.g. large decision tree). In general, the larger the tree, the better the fit to the training data.

- It's hard to think of a better fit to the training data as a "worse" result. Often difficult to fit training data well, so it seems that "a good fit to the training data means a good result."

**Key figure in machine learning**



Legend:
Validation Set Error ———
Training Set Error ––✗––

**Error rate** (y-axis, values 0, 10, 20, 30, 40, 50, 60)

Optimal tree size

We set tree size as a parameter in our DT learning alg.

**Overfitting kicks in…**

Tree size (x-axis, values 1, 2, 3, 4, 5, 6, 7, 8, 9, 10)

**Tree size**

$error_T(h) < error_T(h')$ but $error_D(h) > error_D(h')$

**Note: with larger and larger trees, we just do better and better on the training set!**

But note the performance on the validation set degrades!

Note: Similar curves can happen when training too long in complex hypothesis space with lots of parameters to set.

# Solutions for Overfitting

- K- fold Cross Validation

- Regularization

- Early stopping

- Drop-out

- Pre or post pruning for decision tree

- Minimum description length (MDL) principle

# K- fold Cross Validation

| Training Set | | Testing Set |
|---|---|---|
| Training Set | Validation Set | Testing Set |

| S1 | S2 | S3 | ... | Sk |
|---|---|---|---|---|

Training Set = S

✗ Average test score = $1/k \left( \sum Si \right)$

| Round | Training Set | Testing Set |
|---|---|---|
| 1 | S1 | S – S1 |
| 2 | S2 | S – S2 |
| | | |
| i | Si | S-Si |

- **Trade-off**:
- Complex hypothesis fit the data well      → may tend to overfitting
- Simple hypothesis may generalize better   → may tend to underfitting
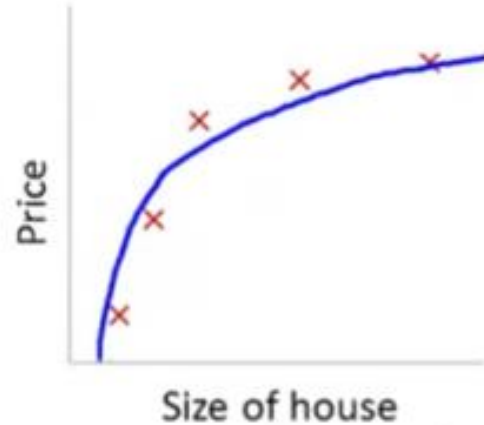- As the training data samples increase, generalization error decreases.
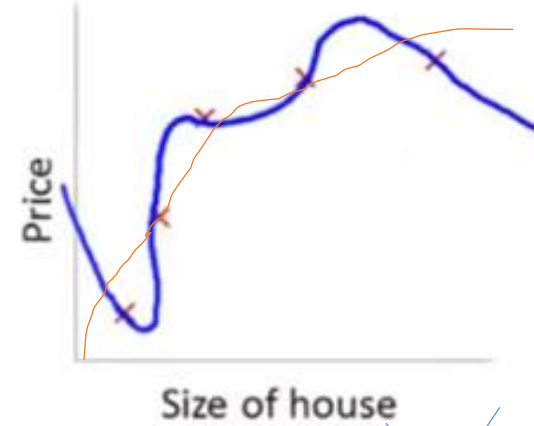
# Regularization

Options:
1. Reduce number of features.
   — Manually select which features to keep.
   — Model selection algorithm

2. Regularization.
   — Keep all the features, but reduce magnitude/values of parameters $\theta_j$.
   — Works well when we have a lot of features, each of which contributes a bit to predicting $y$.

# Regularization

**Intuition**



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Suppose we penalize and make $\theta_3, \theta_4$ really small.

$$\min_\theta \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + 1000\,\theta_3^2 + 1000\,\theta_4^2$$

$$\theta_3 \cong 0 \qquad \theta_4 \cong 0$$

# Regularization

**Regularization.**

Small values for parameters $\theta_0, \theta_1, \ldots, \theta_n$
- — "Simpler" hypothesis
- — Less prone to overfitting

Housing:
- — Features: $x_1, x_2, \ldots, x_{100}$
- — Parameters: $\theta_0, \theta_1, \theta_2, \ldots, \theta_{100}$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{i=1}^{n} \theta_j^2$$
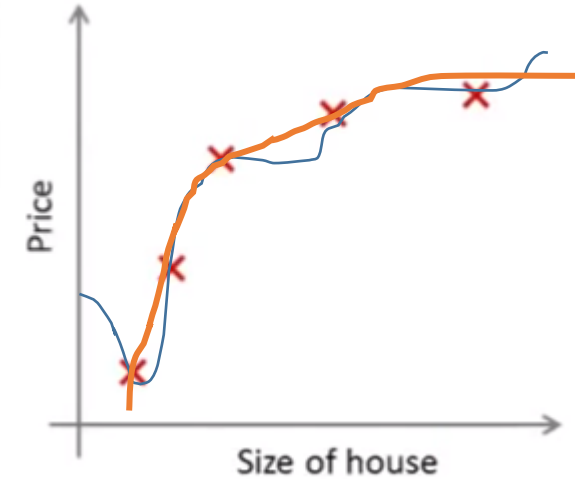
# Regularization

**Regularization.**

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$

$$\min_\theta J(\theta)$$

Regularization parameter

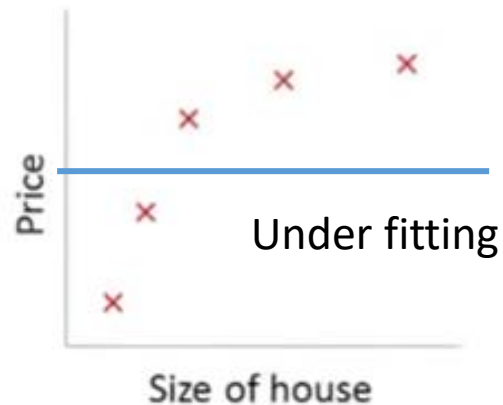- Fitting the data points well

- Keeping the no. of parameters (Ɵs) small



Price

Size of house

# Regularization

In regularized linear regression, we choose $\theta$ to minimize

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$

What if $\lambda$ is set to an extremely large value (perhaps for too large for our problem, say $\lambda = 10^{10}$)?



$\theta_1 \cong 0; \ \theta_2 \cong 0; \theta_3 \cong 0; \ \theta_4 \cong 0;$

Under fitting

$h_\theta(x) = 0$

$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

# Regularized Linear Regression

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$

$$\min_{\theta} J(\theta)$$

## Gradient descent

Repeat {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\frac{\partial}{\partial \theta_j} J(\theta) \Big)$$

$$\theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right]$$

$$(j = 0, 1, 2, 3, \ldots, n)$$

}

# Regularized Linear Regression

$$\theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right]$$
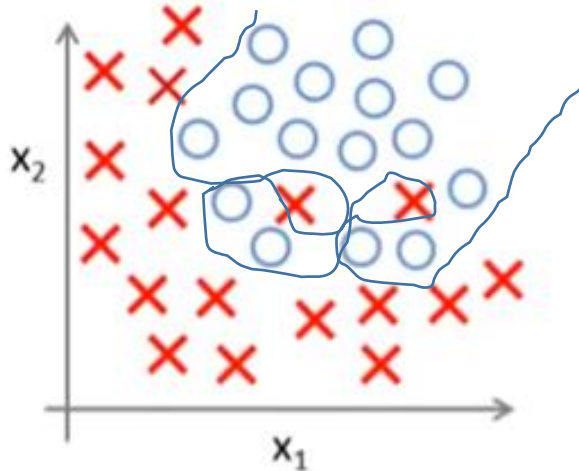
$$(j = 0, 1, 2, 3, \ldots, n)$$

Shrinkage

Parameter updation

$$\theta_j := \theta_j \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$\left(1 - \alpha \frac{\lambda}{m}\right) < 1$$

Shrinkage $\left(1 - \frac{\alpha \lambda}{m}\right)$

# Regularized Logistic Regression



$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2$$
$$+ \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2$$
$$+ \theta_5 x_1^2 x_2^3 + \dots)$$

Cost function:

$$J(\theta) = -\left[ \frac{1}{m} \sum_{i=1}^{m} y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log \left(1 - h_\theta(x^{(i)})\right) \right]$$

$$+ \frac{\lambda}{2m} \sum_{i=1}^{n} \theta_j^2$$

# Regularized Logistic Regression

$$\theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)} + \frac{\lambda}{m}\theta_j \right]$$

$$(j = 0, 1, 2, 3, \ldots, n)$$

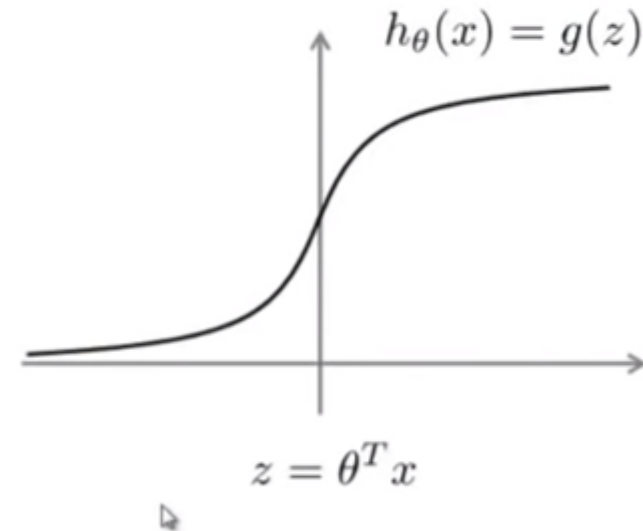$$\frac{\partial}{\partial\theta_j} J(\theta)$$

For Logistic Regression:

$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

# Support Vector Machine

# Optimization objective

**Alternative view of logistic regression**

$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

$h_\theta(x) = g(z)$

$z = \theta^T x$

If $y = 1$, we want $h_\theta(x) \approx 1$, $\quad \theta^T x \gg 0$

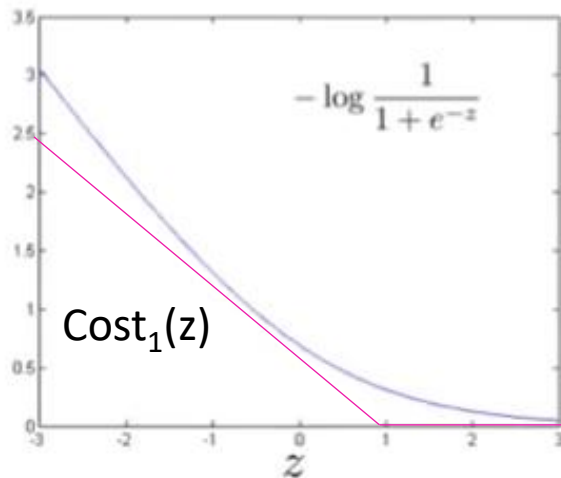If $y = 0$, we want $h_\theta(x) \approx 0$, $\quad \theta^T x \ll 0$

# SVM

**Alternative view of logistic regression**
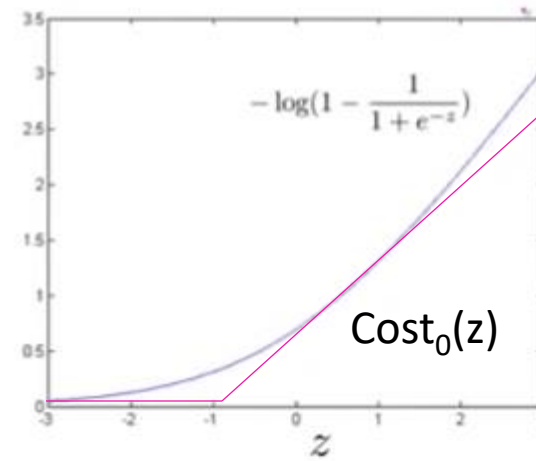
Cost of example:   $-(y \log h_\theta(x) + (1 - y) \log(1 - h_\theta(x)))$

$$= -y \log \frac{1}{1 + e^{-\theta^T x}} - (1 - y) \log(1 - \frac{1}{1 + e^{-\theta^T x}})$$

If $y = 1$ (want $\theta^T x \gg 0$):

$-\log \frac{1}{1 + e^{-z}}$

Cost$_1$(z)

If $y = 0$ (want $\theta^T x \ll 0$):

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

Cost$_0$(z)

# SVM

Logistic regression:

$$\min_{\theta} \frac{1}{m} \left[ \sum_{i=1}^{m} y^{(i)} \left( -\log h_\theta(x^{(i)}) \right) + (1 - y^{(i)}) \left( (-\log(1 - h_\theta(x^{(i)}))) \right) \right] + \frac{\lambda}{2m} \sum_{j=1}^{n} \theta_j^2$$

Support vector machine: $Cost_1(\Theta^T x^{(i)})$      $Cost_0(\Theta^T x^{(i)})$

$$\min_{\theta} \frac{1}{m} \sum_{1}^{m} y^{(i)} Cost_1(\Theta^T x^{(i)}) + (1 - y^{(i)}) Cost_0(\Theta^T x^{(i)}) + \frac{\lambda}{2m} \sum_{i=1}^{n} \theta_j^2$$

A + λ B  =  C A + B     where C = 1/ λ

$$\min_{\theta} C \sum_{i=1}^{m} \left[ y^{(i)} cost_1(\theta^T x^{(i)}) + (1 - y^{(i)}) cost_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{i=1}^{n} \theta_j^2$$

# SVM

**SVM hypothesis**
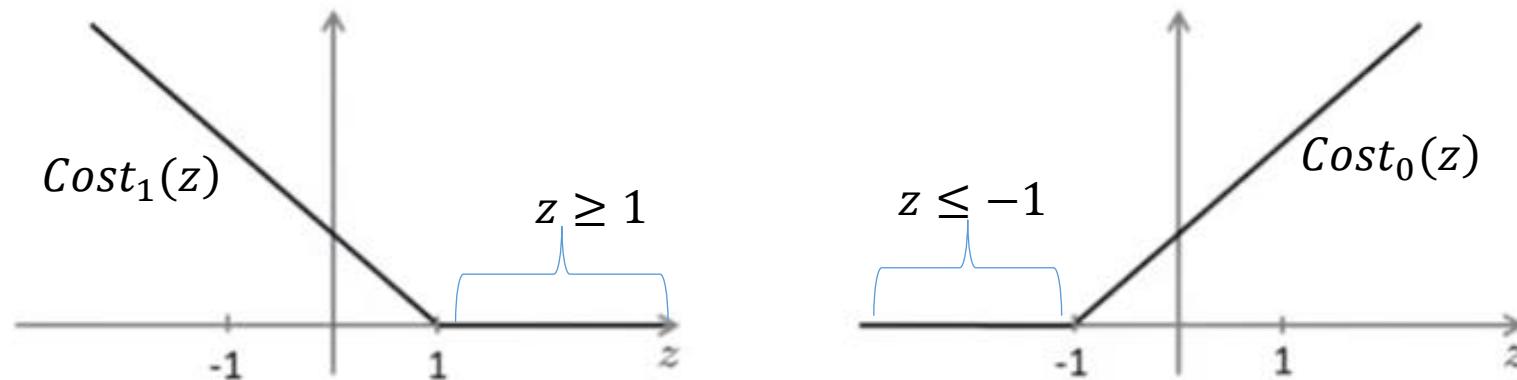
$$\min_{\theta} C \sum_{i=1}^{m} \left[ y^{(i)} cost_1(\theta^T x^{(i)}) + (1 - y^{(i)}) cost_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{i=1}^{n} \theta_j^2$$

Hypothesis:

$$h_{\theta}(x) = \begin{cases} 1 & if \ \theta^T x \geq 0 \\ 0 & otherwise \end{cases}$$

# SVM: As Large Margin Classifier

$$\min_{\theta} C \sum_{i=1}^{m} \left[ y^{(i)} cost_1(\theta^T x^{(i)}) + (1 - y^{(i)}) cost_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{i=1}^{n} \theta_j^2$$

$Cost_1(z)$

$z \geq 1$

$z \leq -1$

$Cost_0(z)$

-1          1          $z$

-1          1          $z$

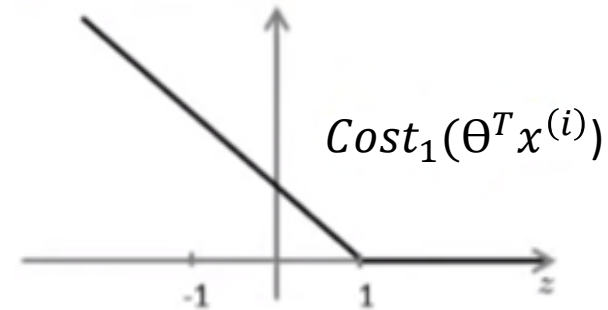If $y = 1$, we want $\theta^T x \geq 1$ (not just $\geq 0$)

If $y = 0$, we want $\theta^T x \leq -1$ (not just $< 0$)

# SVM Decision Boundary

$$\min_{\theta} C \sum_{i=1}^{m} \left[ y^{(i)} cost_1(\theta^T x^{(i)}) + (1 - y^{(i)}) cost_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{i=1}^{n} \theta_j^2$$

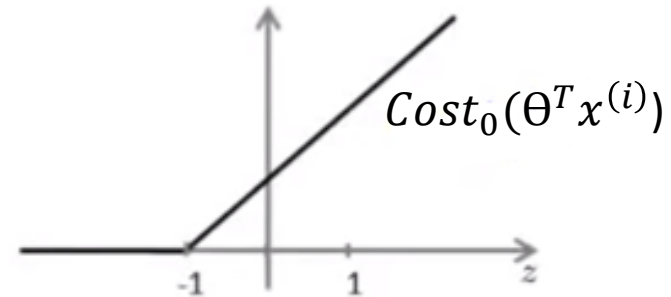**Whenever** $y^{(i)} = 1$:

$$\theta^T x^{(i)} \geq 1$$

$Cost_1(\theta^T x^{(i)})$

**Whenever** $y^{(i)} = 0$:

$$\theta^T x^{(i)} \leq -1$$

$Cost_0(\theta^T x^{(i)})$

# SVM Decision Boundary

$$\min_{\theta} C \sum_{i=1}^{m} \left[ y^{(i)} cost_1(\theta^T x^{(i)}) + (1 - y^{(i)}) cost_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{i=1}^{n} \theta_j^2$$

= 0   as C is very big number

Whenever $y^{(i)} = 1$:

$$\theta^T x^{(i)} \geq 1$$

$$\min_{\theta} C \times 0 + \frac{1}{2} \sum_{i=1}^{n} \theta_j^2 \quad i.e. \quad \min_{\theta} \frac{1}{2} \sum_{i=1}^{n} \theta_j^2$$
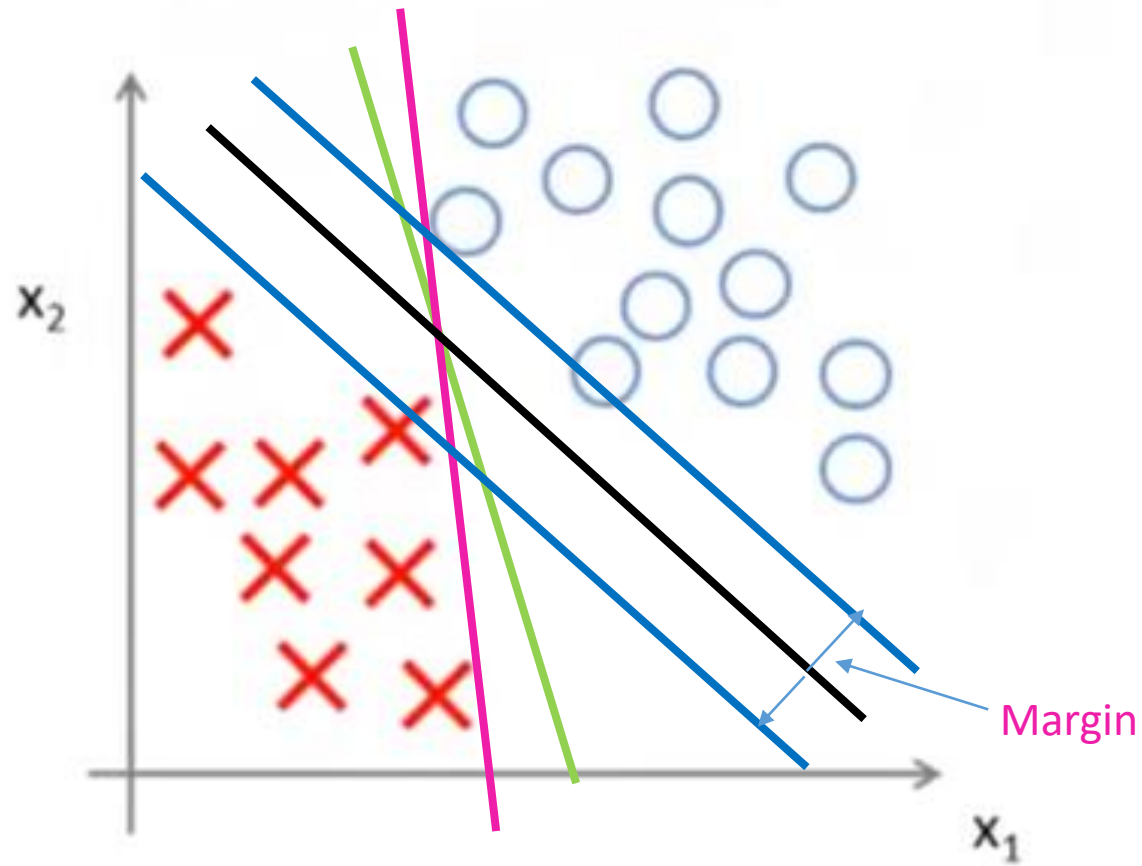
$$s.t. \qquad \theta^T x^{(i)} \begin{cases} \geq 1 & if \quad y^{(i)} = 1 \\ \leq -1 & if \quad y^{(i)} = 0 \end{cases}$$

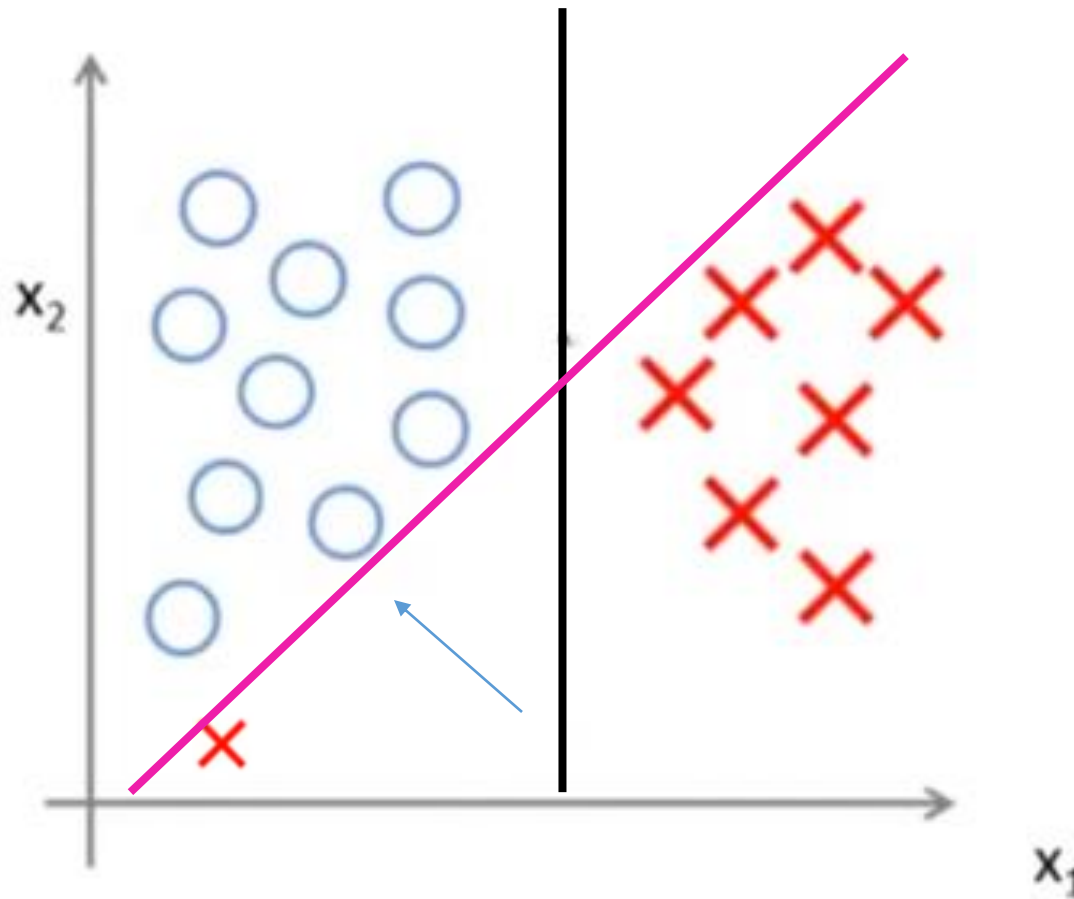Whenever $y^{(i)} = 0$:

$$\theta^T x^{(i)} \leq -1$$

# SVM Decision Boundary
## Linearly separable case

# Large Margin Classifier
## In case of Outliers



C is very large

Sensitive to outliers

# Thanks

CS 590 Lecture 3

# … to continue

CS 590 Lecture 2