

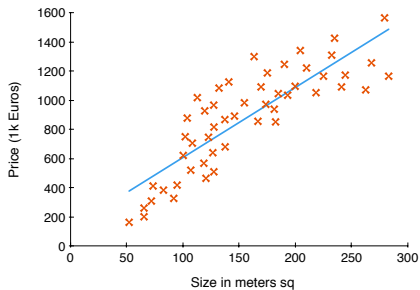
Classification

Lecture 2 - DAMLF | ML1

Regression Analysis

Univariate Linear Regression:

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



Regression Analysis

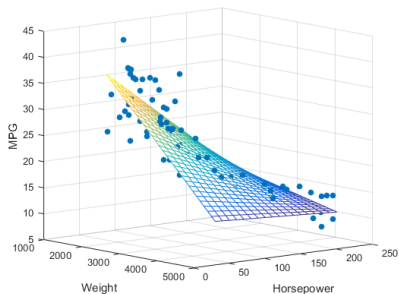
Univariate Linear Regression:

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

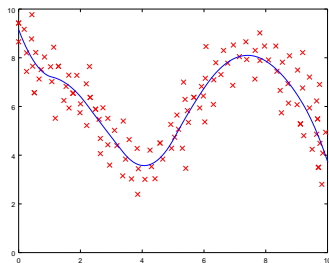
Multiple Linear Regression:

$$h(x; \theta) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 \cdots + \theta_n x_n$$

where $x_0 = 1$



Regression Analysis



Univariate Linear Regression:

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

Multiple Linear Regression:

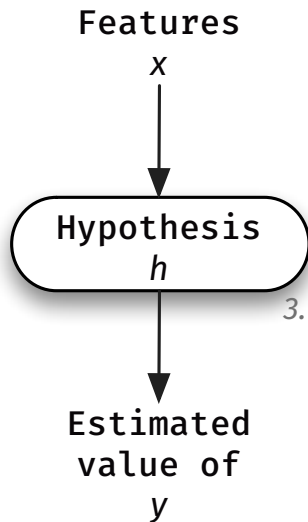
$$h(x; \theta) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 \cdots + \theta_n x_n$$

where $x_0 = 1$

Polynomial Regression:

$$h(x; \theta) = \theta_0 + \theta_1 x + \theta_2 x^2 \cdots + \theta_n x^n$$

Regression Analysis



Univariate Linear Regression:

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

Multiple Linear Regression:

$$h(x; \theta) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 \cdots + \theta_n x_n$$

where $x_0 = 1$

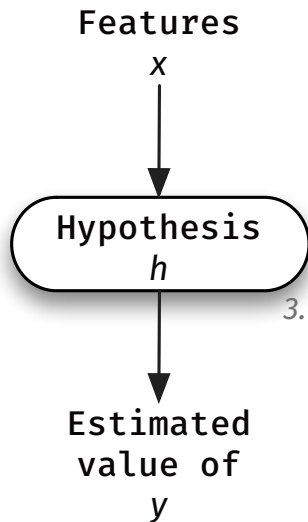
Polynomial Regression:

$$h(x; \theta) = \theta_0 + \theta_1 x + \theta_2 x^2 \cdots + \theta_n x^n$$

Cost Function $J(\theta)$:

$$MSE = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - h(x^{(i)}; \theta))^2$$

Regression Analysis



Univariate Linear Regression:

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

$$\text{Normal Equation: } \hat{\theta} = (X^T X)^{-1} X^T y$$

Multiple Linear Regression:

$$h(x; \theta) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 \cdots + \theta_n x_n$$

where $x_0 = 1$

$$\text{Normal Equation: } \hat{\theta} = (X^T X)^{-1} X^T y$$

Polynomial Regression:

$$h(x; \theta) = \theta_0 + \theta_1 x + \theta_2 x^2 \cdots + \theta_n x^n$$

$$\text{Normal Equation: } \hat{\theta} = (X^T X)^{-1} X^T y$$

Cost Function $J(\theta)$:

$$MSE = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - h(x^{(i)}; \theta))^2$$

Normal Equation & the MSE Estimator

The **Ordinary Least Squares** or **Mean Squared Error Estimator** $\hat{\theta}$

$$\hat{\theta} := \min_{\theta} J(\theta) = \min_{\theta_0, \theta_1, \dots, \theta_n} J(\theta_0, \theta_1, \dots, \theta_n)$$

where $\hat{\theta}$ is solved¹ by

$$(X\hat{\theta} - y)^T X = 0$$

$$y^T X = \hat{\theta}^T (X^T X)$$

$$X^T y = (X^T X) \hat{\theta}$$

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

¹In statistics, β and $\hat{\beta}$ are sometimes used instead of θ and $\hat{\theta}$, respectively.

WHAT IS A LEARNING ALGORITHM?

"A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E "

–Tom Mitchell

Ingredients:

Task *to perform*

Type of **Experience**

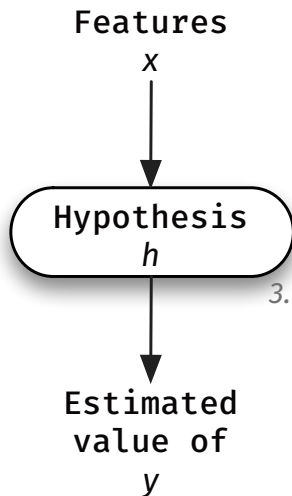
Performance *measure*

WHAT IS A LEARNING ALGORITHM?

*"A computer program is said to learn from **experience** E with respect to some class of tasks T and **performance** measure P , if its performance at **tasks** in T , as measured by P , improves with experience E "*

–Tom Mitchell

Gradient Descent for Multiple Linear Regression



Multiple Linear Regression

1. Hypothesis

$$h(x; \theta) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

2. Parameters

θ ($n + 1$ -dimensional scalar-valued vector)

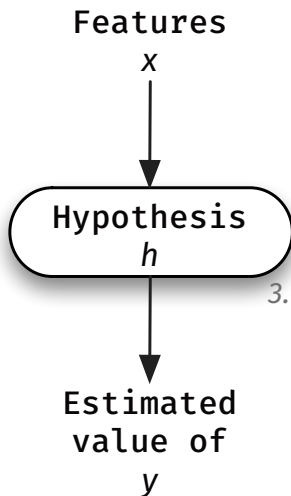
3. Cost Function

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}; \theta) - y^{(i)})^2$$

4. Goal

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

Gradient Descent for Multiple Linear Regression



For $n \geq 1$:
» Enter loop
Repeat:

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

simultaneously update θ_j
for $j = 0, \dots, n$
Stop at convergence
» Exit loop

Note: remember that $x_0 = 1$;
that is., $x_0^{(i)} = 1$ for all $i = 1 \dots m$.

Gradient Descent for Multiple Linear Regression

Tips:

How to choose a good learning rate, α

Importance of feature scaling

Iterations to convergence can vary enormously from problem-to-problem

For $n \geq 1$:
» Enter loop
Repeat:

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

simultaneously update θ_j
for $j = 0, \dots, n$

Stop at convergence
» Exit loop

Note: remember that $x_0 = 1$;
that is., $x_0^{(i)} = 1$ for all $i = 1 \dots m$.

Regression vs Classification

Regression Problems

Input: Features x are **continuous** valued *data*

Output: *Predicts* a **continuous** value y .

Classification Problems

Input: Features x are **continuous** valued *data*

Output: *Predicts* a **discrete** value y .

Binary classification problem

Predicts output variables y in $\{0, 1\}$, where

0 denotes the **negative** class

1 denotes the **positive** class

Binary classification problem

Predicts output variables y in $\{0, 1\}$, where

0 denotes the **negative** class

1 denotes the **positive** class

Examples:

- **Email spam filter**: 1 = spam; 0 = not spam
- **Tumor classification**: 1 = malignant; 0 = benign
- **University Admissions**: 1 = accept; 0 = reject
- **Credit score**: 1 = credit-worthy; 0 = not credit-worthy
- \vdots

Classification is Different from Regression

Example: Credit Scoring



Classification is Different from Regression

Example: Credit Scoring



Suppose we applied **linear regression** to this example:

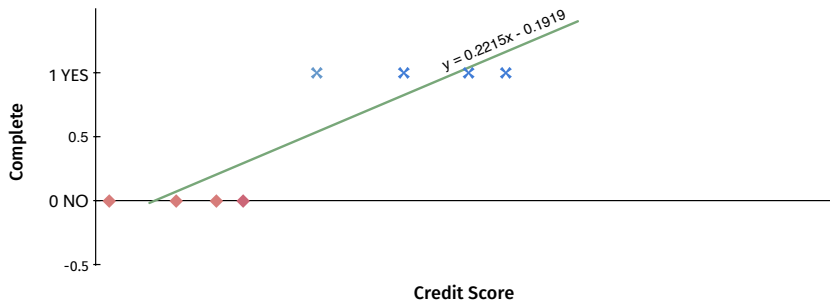
$$h(x^{(i)}; \theta) = \theta^T x^{(i)}$$

where

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

Classification is Different from Regression

Example: $h(x; \theta) = \theta^T x$



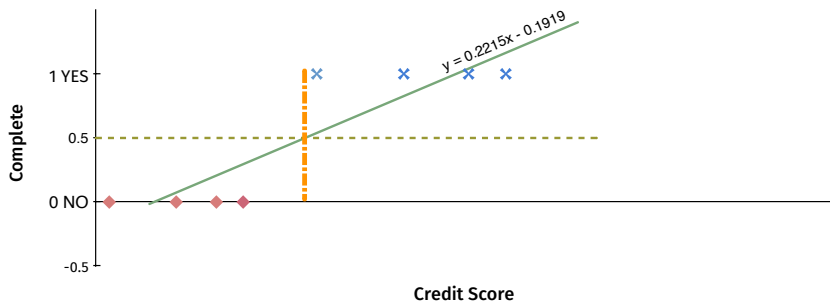
Threshold Classification

If $h(x; \theta) \geq \frac{1}{2}$, predict that $y = 1$

If $h(x; \theta) < \frac{1}{2}$, predict that $y = 0$

Classification is Different from Regression

Example: $h(x; \theta) = \theta^T x$



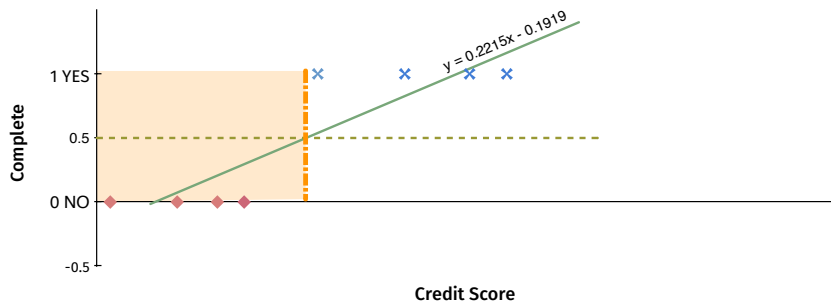
Threshold Classification

If $h(x; \theta) \geq \frac{1}{2}$, predict that $y = 1$

If $h(x; \theta) < \frac{1}{2}$, predict that $y = 0$

Classification is Different from Regression

Example: $h(x; \theta) = \theta^T x$



Threshold Classification

If $h(x; \theta) \geq \frac{1}{2}$, predict that $y = 1$

If $h(x; \theta) < \frac{1}{2}$, predict that $y = 0$

Classification is Different from Regression

Example: $h(x; \theta) = \theta^T x$



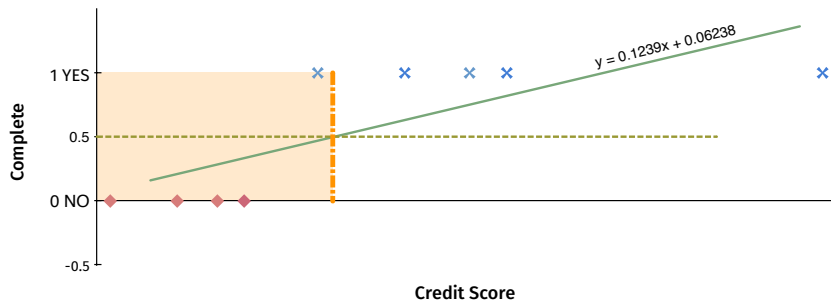
Threshold Classification

If $h(x; \theta) \geq \frac{1}{2}$, predict that $y = 1$

If $h(x; \theta) < \frac{1}{2}$, predict that $y = 0$

Classification is Different from Regression

Example: $h(x; \theta) = \theta^T x$



Threshold Classification

If $h(x; \theta) \geq \frac{1}{2}$, predict that $y = 1$

If $h(x; \theta) < \frac{1}{2}$, predict that $y = 0$

Classification is Different from Regression

Binary Classification Output:

$$y = 0 \quad \text{or} \quad y = 1$$

Classification is Different from Regression

Binary Classification Output:

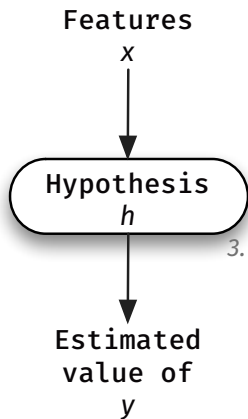
$$y = 0 \quad \text{or} \quad y = 1$$

However,

$$h(x; \theta) = \theta^T x < 0 \quad \text{and} \quad h(x; \theta) = \theta^T x > 1,$$

is possible.

Hypothesis Representation



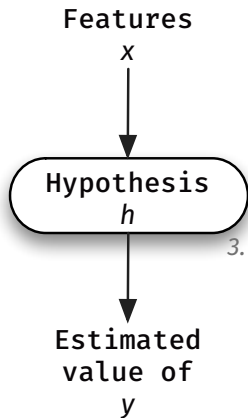
Recall:

When selecting a learning algorithm, you must decide how to **represent** h .

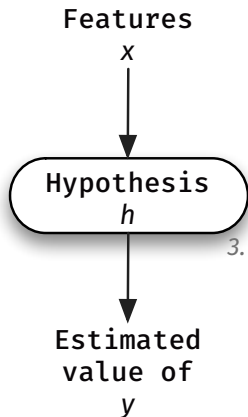
Hypothesis Representation

For binary classification, we want $h(x; \theta)$ to take values between 0 and 1:

$$0 \leq h(x; \theta) \leq 1$$



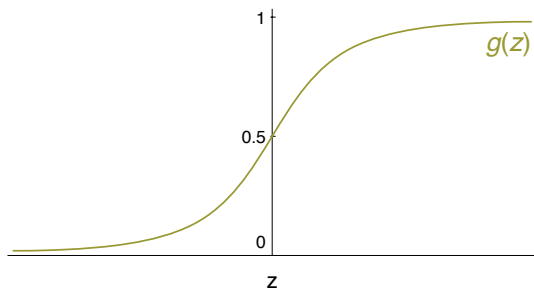
Hypothesis Representation



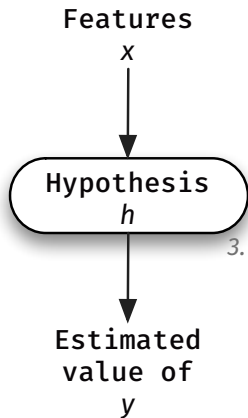
For binary classification, we want $h(x; \theta)$ to take values between 0 and 1:

$$0 \leq h(x; \theta) \leq 1$$

Sigmoid Function $g(\cdot)$



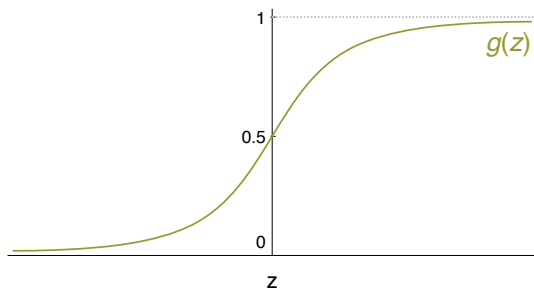
Hypothesis Representation



For binary classification, we want $h(x; \theta)$ to take values between 0 and 1:

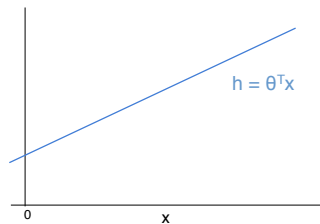
$$0 \leq h(x; \theta) \leq 1$$

Sigmoid Function $g(\cdot)$



Hypothesis Representation

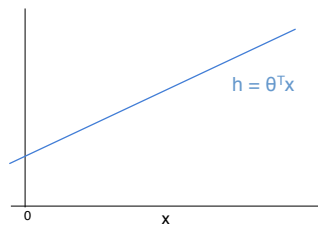
Linear Regression



$$\begin{aligned} h(x; \theta) &= \theta_0 + \theta_1 x \\ &= \theta^T x \end{aligned}$$

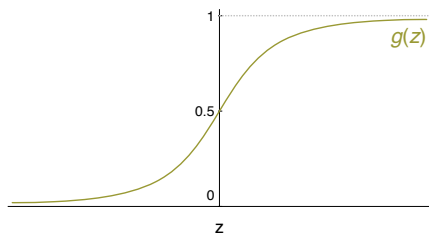
Hypothesis Representation

Linear Regression



$$\begin{aligned} h(x; \theta) &= \theta_0 + \theta_1 x \\ &= \theta^T x \end{aligned}$$

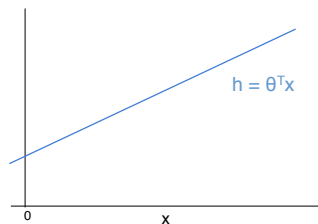
Logistic Regression



$$h(x; \theta) = g(\theta^T x)$$

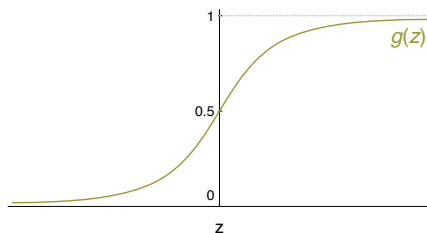
Hypothesis Representation

Linear Regression



$$\begin{aligned} h(x; \theta) &= \theta_0 + \theta_1 x \\ &= \theta^T x \end{aligned}$$

Logistic Regression



$$h(x; \theta) = g(\theta^T x)$$

where

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

So:

$$h(x; \theta) = \frac{1}{1 + \exp^{-\theta^T x}}$$

Hypothesis Interpretation

Logistic Regression

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

= estimated probability that $y = 1$ given data x parameterized by θ

$$= P(y = 1 \mid x; \theta)$$

Hypothesis Interpretation

Logistic Regression

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

= estimated probability that $y = 1$ given data x parameterized by θ

$$= P(y = 1 \mid x; \theta)$$

Example:

$h(x; \theta) = .37$ means: 37% chance that applicant will be accepted.

Hypothesis Interpretation

Logistic Regression

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

= estimated probability that $y = 1$ given data x parameterized by θ

$$= P(y = 1 \mid x; \theta)$$

$$= 0.37$$

Question: *What is the probability that $y = 0$?*

Hypothesis Interpretation

Logistic Regression

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

= estimated probability that $y = 1$ given data x parameterized by θ

$$= P(y = 1 \mid x; \theta)$$

$$= 0.37$$

Question: What is the probability that $y = 0$?

Since ($y = 0$ xor $y = 1$), apply the **Law of Total Probability**:

$$P(y = 0 \mid x; \theta) + P(y = 1 \mid x; \theta) = 1 \quad (\text{Law of Total Probability})$$

$$P(y = 0 \mid x; \theta) = 1 - P(y = 1 \mid x; \theta)$$

$$P(y = 0 \mid x; \theta) = 1 - 0.37$$

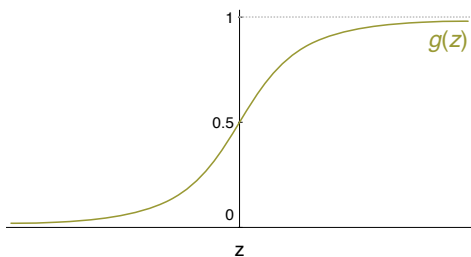
$$P(y = 0 \mid x; \theta) = 0.63$$

Using $g(\cdot)$ for prediction

$$\begin{aligned}h(x; \theta) &= g(\theta^T x) \\ &= P(y = 1 \mid x; \theta)\end{aligned}$$

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

Sigmoid Function $g(z)$



Using $g(\cdot)$ for prediction

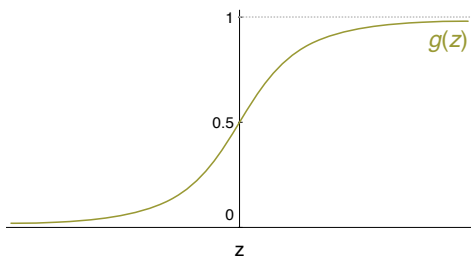
$$\begin{aligned}h(x; \theta) &= g(\theta^T x) \\ &= P(y = 1 \mid x; \theta)\end{aligned}$$

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

Predict that ' $y = 1$ ' if $h(x; \theta) \geq 0.5$

Predict that ' $y = 0$ ' if $h(x; \theta) < 0.5$

Sigmoid Function $g(z)$



Using $g(\cdot)$ for prediction

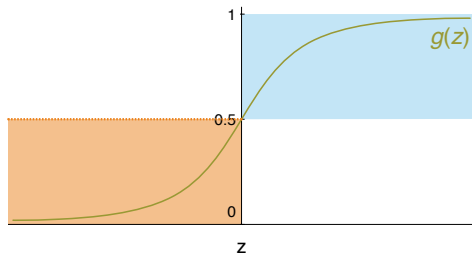
$$\begin{aligned}h(x; \theta) &= g(\theta^T x) \\ &= P(y = 1 \mid x; \theta)\end{aligned}$$

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

Predict that ' $y = 1$ ' if $h(x; \theta) \geq 0.5$

Predict that ' $y = 0$ ' if $h(x; \theta) < 0.5$

Sigmoid Function $g(z)$



Using $g(\cdot)$ for prediction

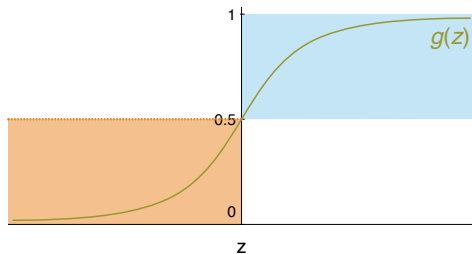
$$\begin{aligned}h(x; \theta) &= g(\theta^T x) \\ &= P(y = 1 \mid x; \theta)\end{aligned}$$

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

Predict that ' $y = 1$ ' if $h(x; \theta) \geq 0.5$
when $z \geq 0$

Predict that ' $y = 0$ ' if $h(x; \theta) < 0.5$
when $z < 0$

Sigmoid Function $g(z)$



Using $g(\cdot)$ for prediction

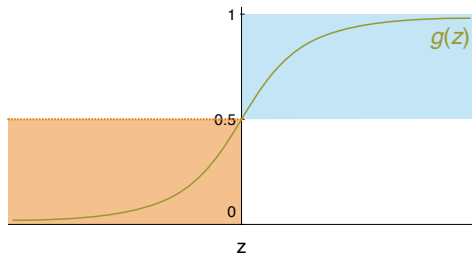
$$\begin{aligned}h(x; \theta) &= g(\theta^T x) \\ &= P(y = 1 \mid x; \theta)\end{aligned}$$

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

Predict that ' $y = 1$ ' if $h(x; \theta) \geq 0.5$
when $z \geq 0$
when $\theta^T x \geq 0$

Predict that ' $y = 0$ ' if $h(x; \theta) < 0.5$
when $z < 0$
when $\theta^T x < 0$

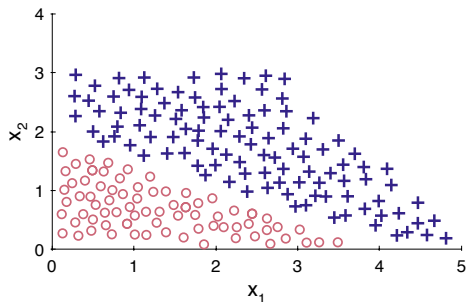
Sigmoid Function $g(z)$



Decision Boundaries

$$\begin{aligned}h(x; \theta) &= g(\theta^T x) \\&= g(\theta_0 + \theta_1 x_1 + \theta_2 x_2) \\&= \\&= \end{aligned}$$

Decision Boundary

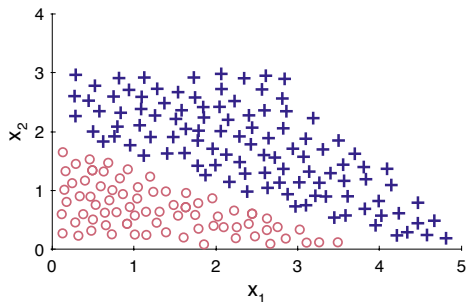


Decision Boundaries

$$\begin{aligned}h(x; \theta) &= g(\theta^T x) \\&= g(\theta_0 + \theta_1 x_1 + \theta_2 x_2) \\&= g(-4 + x_1 + 2x_2), \text{ where}\end{aligned}$$

$$\theta = \begin{bmatrix} -4 \\ 1 \\ 2 \end{bmatrix}$$

Decision Boundary



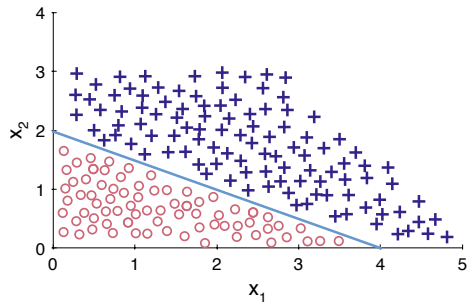
Decision Boundaries

$$\begin{aligned}h(x; \theta) &= g(\theta^T x) \\&= g(\theta_0 + \theta_1 x_1 + \theta_2 x_2) \\&= g(-4 + x_1 + 2x_2), \text{ where}\end{aligned}$$

$$\theta = \begin{bmatrix} -4 \\ 1 \\ 2 \end{bmatrix}$$

Predict that ' $y = 1$ ' if $h(x; \theta) \geq 0.5$

Decision Boundary



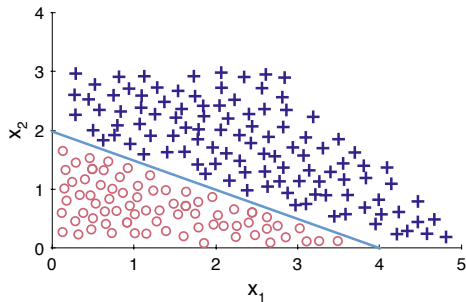
Decision Boundaries

$$\begin{aligned}h(x; \theta) &= g(\theta^T x) \\&= g(\theta_0 + \theta_1 x_1 + \theta_2 x_2) \\&= g(-4 + x_1 + 2x_2), \text{ where}\end{aligned}$$

$$\theta = \begin{bmatrix} -4 \\ 1 \\ 2 \end{bmatrix}$$

Predict that ' $y = 1$ ' if $h(x; \theta) \geq 0.5$
when $-4 + x_1 + 2x_2 \geq 0$

Decision Boundary



Decision Boundaries

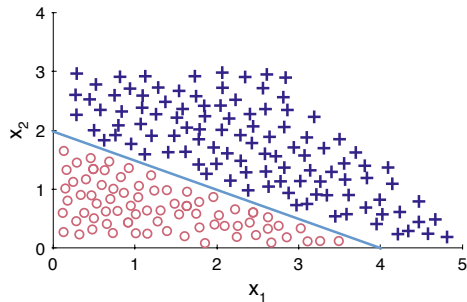
$$\begin{aligned}h(x; \theta) &= g(\theta^T x) \\&= g(\theta_0 + \theta_1 x_1 + \theta_2 x_2) \\&= g(-4 + x_1 + 2x_2), \text{ where}\end{aligned}$$

$$\theta = \begin{bmatrix} -4 \\ 1 \\ 2 \end{bmatrix}$$

Predict that ' $y = 1$ ' if $h(x; \theta) \geq 0.5$
when $-4 + x_1 + 2x_2 \geq 0$

Predict that ' $y = 0$ ' if $h(x; \theta) < 0.5$
when $-4 + x_1 + 2x_2 < 0$

Decision Boundary



Decision Boundaries

$$\begin{aligned}h(x; \theta) &= g(\theta^T x) \\&= g(\theta_0 + \theta_1 x_1 + \theta_2 x_2) \\&= g(-4 + x_1 + 2x_2), \text{ where}\end{aligned}$$

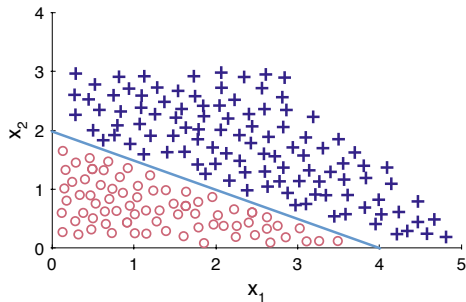
$$\theta = \begin{bmatrix} -4 \\ 1 \\ 2 \end{bmatrix}$$

Predict that ' $y = 1$ ' if $h(x; \theta) \geq 0.5$
when $-4 + x_1 + 2x_2 \geq 0$

— **The Decision Boundary** — $x_1 + 2x_2 = 4$

Predict that ' $y = 0$ ' if $h(x; \theta) < 0.5$
when $-4 + x_1 + 2x_2 < 0$

Decision Boundary

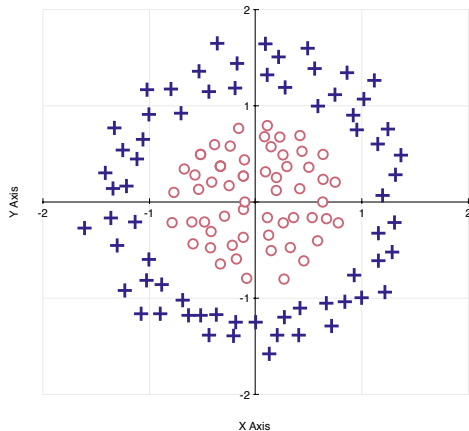


Non-linear Decision Boundaries

$$\begin{aligned}h(x; \theta) &= g(\theta^T x) \\&= g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2) \\&= g(-1 + x_1^2 + x_2^2), \text{ where}\end{aligned}$$

$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

Decision Boundary

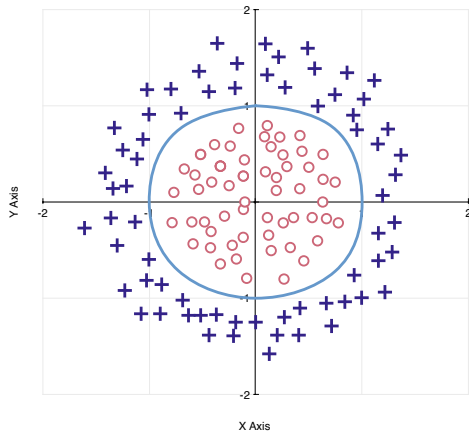


Non-linear Decision Boundaries

$$\begin{aligned}h(x; \theta) &= g(\theta^T x) \\&= g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2) \\&= g(-1 + x_1^2 + x_2^2), \text{ so:}\end{aligned}$$

Decision Boundary: $x_1^2 + 2x_2^2 = 1$

Decision Boundary

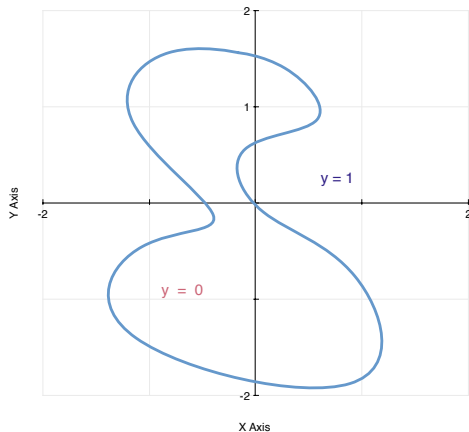


Non-linear Decision Boundaries

Note:

Decision Boundaries are functions of **parameters** rather than functions of data + parameters.

Decision Boundary



How to choose θ ?

How to choose θ ?

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

How to choose θ ?

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

$n + 1$ Features: $x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$, where $x_0 = 1$.

How to choose θ ?

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

$n + 1$ Features: $x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$, where $x_0 = 1$.

Target variable: $y \in \{0, 1\}$.

How to choose θ ?

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

$n + 1$ Features: $x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$, where $x_0 = 1$.

Target variable: $y \in \{0, 1\}$.

Form of the Hypothesis: $h(x; \theta) = \frac{1}{1 + e^{-\theta^T x}}$

How to choose θ ?

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

$n + 1$ Features: $x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$, where $x_0 = 1$.

Target variable: $y \in \{0, 1\}$.

Form of the Hypothesis: $h(x; \theta) = \frac{1}{1 + e^{-\theta^T x}}$

How are the values of the vector θ chosen?

How to choose θ

Recall how θ are chosen for Linear Regression:

Linear Regression **Cost Function:**

$$\begin{aligned} J(\theta) &= \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h(x^{(i)}; \theta) - y^{(i)})^2 \\ &= \frac{1}{m} \sum_{i=1}^m \text{Cost}(h(x^{(i)}; \theta), y^{(i)}) \end{aligned}$$

Simplifying, the cost function for Linear Regression:

$$\text{Cost}(h(x; \theta), y) = \frac{1}{2} (h(x; \theta) - y)^2$$

Logistic Regression Cost Function

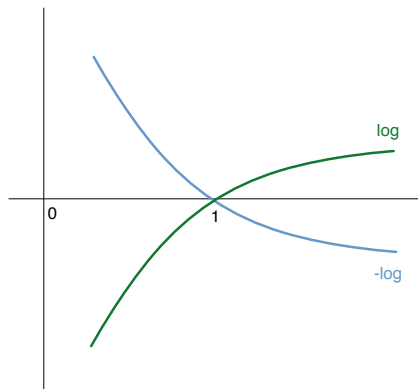
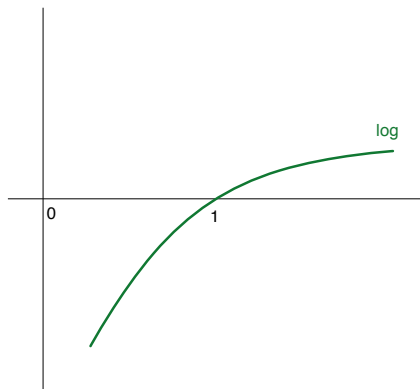
$$\text{Cost}(h(x; \theta), y) = \begin{cases} -\log(h(x; \theta)) & \text{if } y = 1 \\ -\log(1 - h(x; \theta)) & \text{if } y = 0 \end{cases}$$

Intuition:

Logistic Regression Cost Function

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

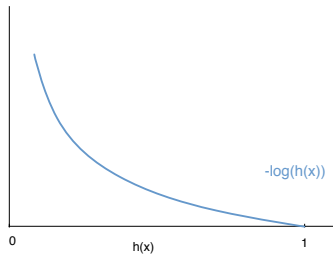
Intuition:



Logistic Regression Cost Function

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

If $y = 1$:



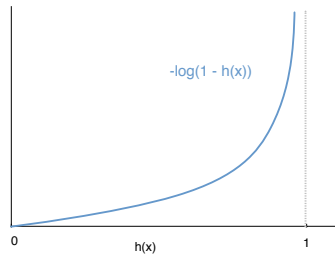
Cost = 0 if $y = 1$ and $h(x; \theta) = 1$;

Cost $\rightarrow \infty$ as $h(x; \theta) \rightarrow 0$

Logistic Regression Cost Function

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

If $y = 0$:



Cost = 0 if $y = 0$ and $h(x; \theta) = 0$;

Cost $\rightarrow \infty$ as $h(x; \theta) \rightarrow 1$

Logistic Regression Cost Function

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

Simplification:

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

if $y = 1$:

$$\begin{aligned} \text{Cost}(h(x; \theta), y) &= -1 \log(h(x; \theta)) - (1 - 1) \log(1 - h(x; \theta)) \\ &= -\log(h(x; \theta)) - 0 \\ &= -\log(h(x; \theta)) \end{aligned}$$

if $y = 0$:

$$\begin{aligned} \text{Cost}(h(x; \theta), y) &= -0 \log(h(x; \theta)) - (1 - 0) \log(1 - h(x; \theta)) \\ &= 0 - \log(1 - h(x; \theta)) \\ &= -\log(1 - h(x; \theta)) \end{aligned}$$

Logistic Regression Cost Function

$$\begin{aligned} J(\theta) &= \frac{1}{m} \sum_{i=1}^m \text{Cost}(h(x^{(i)}; \theta), y^{(i)}) \\ &= -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(h(x; \theta)) + (1 - y^{(i)}) \log(1 - h(x; \theta)) \right] \end{aligned}$$

Logistic Regression Cost Function

$$\begin{aligned} J(\theta) &= \frac{1}{m} \sum_{i=1}^m \text{Cost}(h(x^{(i)}; \theta), y^{(i)}) \\ &= -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(h(x; \theta)) + (1 - y^{(i)}) \log(1 - h(x; \theta)) \right] \end{aligned}$$

$J(\theta)$ is derived from the **Principle of Maximum Likelihood Estimation** (MLE).

Logistic Regression Cost Function

$$\begin{aligned} J(\theta) &= \frac{1}{m} \sum_{i=1}^m \text{Cost}(h(x^{(i)}; \theta), y^{(i)}) \\ &= -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(h(x; \theta)) + (1 - y^{(i)}) \log(1 - h(x; \theta)) \right] \end{aligned}$$

In order to **fit parameters** θ , we first $\min_{\theta} J(\theta)$

Then, given $\min_{\theta} J(\theta)$ and a **new** x , we **make a prediction** for **unknown** y using $h(x; \theta)$ calculated by:

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

Gradient Descent and Binary Logistic Regression

Algorithm: Gradient Descent

» Enter loop
Repeat:

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

for all θ_j simultaneously
Stop at convergence

» Exit loop

Model: Logistic Regression

Hypothesis:

$$h(x; \theta) = g(\theta^T x)$$

Cost Function:

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

Derivative Term:

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

Gradient Descent and Binary Logistic Regression

Algorithm: Gradient Descent

» Enter loop

Repeat:

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

for all θ_j simultaneously

Stop at convergence

» Exit loop

Remarks:

Gradient Descent and Binary Logistic Regression

Algorithm: Gradient Descent

» Enter loop

Repeat:

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

for all θ_j simultaneously

Stop at convergence

» Exit loop

Remarks:

Compare to Gradient
Descent for Linear Regression:

Gradient Descent and Binary Logistic Regression

Algorithm: Gradient Descent

» Enter loop

Repeat:

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

for all θ_j simultaneously

Stop at convergence

» Exit loop

Remarks:

Compare to Gradient
Descent for Linear Regression:

Same algorithm?!

Gradient Descent and Binary Logistic Regression

Algorithm: Gradient Descent

» Enter loop
Repeat:

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

for all θ_j simultaneously

Stop at convergence

» Exit loop

Remarks:

Compare to Gradient
Descent for Linear Regression:

Same algorithm?!

The Difference:

$h(x; \theta) = \theta^T x$
- linear regression

$h(x; \theta) = \frac{1}{1 + e^{-\theta^T x}}$
- logistic regression

Simplifying the Cost Function

Gradient Descent & Advanced Algorithms

So far, you have learned the fundamentals of regression-based machine learning.

So far, you have learned the fundamentals of regression-based machine learning.

But, to be serious – to deal with thousands (or more) features and millions (or more) training examples – we need more sophisticated algorithms.

Optimization Task

Here's what we know how to do now.

Optimization Task

Here's what we know how to do now.

1. Start with a cost function $J(\theta)$

Optimization Task

Here's what we know how to do now.

1. Start with a cost function $J(\theta)$
2. Goal: $\min_{\theta} J(\theta)$

Optimization Task

Here's what we know how to do now.

1. Start with a cost function $J(\theta)$

2. Goal: $\min_{\theta} J(\theta)$

3. Given θ , compute two things:

$J(\theta)$ the cost of fitting with θ
 $\frac{\partial}{\partial \theta_j} J(\theta)$ the partial derivatives for each $\theta_j \in \theta$

Optimization Task

Here's what we know how to do now.

1. Start with a cost function $J(\theta)$

2. Goal: $\min_{\theta} J(\theta)$

3. Given θ , compute two things:

$J(\theta)$ the cost of fitting with θ
 $\frac{\partial}{\partial \theta_j} J(\theta)$ the partial derivatives for each $\theta_j \in \theta$

4. Compute the partial derivative (well, this has been done for you)

Optimization Task

Here's what we know how to do now.

1. Start with a cost function $J(\theta)$

2. Goal: $\min_{\theta} J(\theta)$

3. Given θ , compute two things:

$J(\theta)$ the cost of fitting with θ

$\frac{\partial}{\partial \theta_j} J(\theta)$ the partial derivatives for each $\theta_j \in \theta$

4. Compute the partial derivative (well, this has been done for you)

5. Plug the derivative term into *Gradient Descent*:

» Enter loop

Repeat:

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

for all θ_j simultaneously

Stop at convergence

» Exit loop

Optimization Task

Now let's look at this optimization task anew.

1. Start with a cost function $J(\theta)$
2. Goal: $\min_{\theta} J(\theta)$
3. Given θ , compute two things:
 - $J(\theta)$ the cost of fitting with θ
 - $\frac{\partial}{\partial \theta_j} J(\theta)$ the partial derivatives for each $\theta_j \in \theta$
4. Optimization algorithms:
 - Gradient Descent

Optimization Task

Now let's look at this optimization task anew.

1. Start with a cost function $J(\theta)$

2. Goal: $\min_{\theta} J(\theta)$

3. Given θ , compute two things:

$J(\theta)$ the cost of fitting with θ

$\frac{\partial}{\partial \theta_j} J(\theta)$ the partial derivatives for each $\theta_j \in \theta$

4. Optimization algorithms:

- Gradient Descent

however, asymptotic rate of convergence is inferior to more advanced algorithms

- **BFGS**

 - *Broyden-Fletcher-Goldfarb-Shanno* Algorithm

Optimization Task

Now let's look at this optimization task anew.

1. Start with a cost function $J(\theta)$

2. Goal: $\min_{\theta} J(\theta)$

3. Given θ , compute two things:

$J(\theta)$ the cost of fitting with θ

$\frac{\partial}{\partial \theta_j} J(\theta)$ the partial derivatives for each $\theta_j \in \theta$

4. Optimization algorithms:

- Gradient Descent

however, asymptotic rate of convergence is inferior to more advanced algorithms

- **BFGS**

- *Broyden-Fletcher-Goldfarb-Shanno* Algorithm
- BFGS is a *quasi-Newton* approximation method
(stores histories of dense matrices in memory)
- L-BFGS refers to *limited memory* BFGS
(stores vectors as estimates of full matrices)

Advanced Optimization Algorithms

Hessian based methods: use 2nd derivative information

- Broyden-Fletcher-Goldfarb-Shanno (BFGS) Algorithm
- Limited Memory BFGS (L-BFGS)
- Conjugate gradient

Advantages:

- α is picked and adjusted automatically (at every iteration)
- Often faster than gradient descent (i.e., fewer iterations)

Advanced Optimization Algorithms

Hessian based methods: use 2nd derivative information

- Broyden-Fletcher-Goldfarb-Shanno (BFGS) Algorithm
- Limited Memory BFGS (L-BFGS)
- Conjugate gradient

Advantages:

- α is picked and adjusted automatically (at every iteration)
- Often faster than gradient descent (i.e., fewer iterations)

Disadvantages:

- Each step is more computationally / memory expensive
- More complicated to implement / harder to inspect

How do these advanced algorithms work?

Hessian based methods: use 2nd derivative information

- Broyden-Fletcher-Goldfarb-Shanno (BFGS) Algorithm
- Limited Memory BFGS (L-BFGS)
- Conjugate gradient

Line Search Strategy

- Find a descent direction (pos, neg, 0)
 - by gradient descent or Newton's method, etc.
- Compute a step-size in that direction.
- Repeat.

Simulated annealing

- probabilistic technique for finding a **global minima**
 - *for more, look it up!*

How do these advanced algorithms work?

In Python, there is a built-in implementation of BFGS called `fmin_bfgs`.

Let's see an example.

BFGS Example

using **fmin_bfgs** to compute `theta_new`

```
» theta_new = fmin_bfgs(computeCostFunction,  
:      initial_theta, fprime=computeGradient,  
:      maxiter=100, args=(X,y), disp=True)
```


BFGS Example

using **fmin_bfgs** to compute `theta_new`

```
» theta_new = fmin_bfgs(computeCostFunction,  
:      initial_theta, fprime=computeGradient,  
:      maxiter=100, args=(X,y), disp=True)
```

fmin_bfgs(-1-, -2-, -3-, -4-, -5-, -6-)

BFGS Example

using **fmin_bfgs** to compute theta_new

```
» theta_new = fmin_bfgs(computeCostFunction,  
:                       initial_theta, fprime=computeGradient,  
:                       maxiter=100, args=(X,y), disp=True)
```

fmin_bfgs(-1-, -2-, -3-, -4-, -5-, -6-)

1. Optimization function to be minimized
2. Initial parameter vector
3. fprime = —
optional parameter to
compute gradient with —
4. maxiter = —
maximum number of iterations set to —
5. args = (,)
features
target
6. disp = —
optional parameter to print convergence
message when set to True

BFGS Example

$$\text{initial_theta } \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$$

Compute Cost Function (J):

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

Gradient (grad):

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

using **fmin_bfgs** to compute theta_new

```
» theta_new = fmin_bfgs(computeCostFunction,  
:                       initial_theta, fprime=computeGradient,  
:                       maxiter=100, args=(X,y), disp=True)
```

fmin_bfgs(-1-, -2-, -3-, -4-, -5-, -6-)

1. Optimization function to be minimized
2. Initial parameter vector
3. fprime = —
optional parameter to
compute gradient with —
4. maxiter = —
maximum number of iterations set to —
5. args = (,)
features
target
6. disp = —
optional parameter to print convergence
message when set to True

Multi-class Classification

| y | $y = 1$ | $y = 2$ |
|---------|---------|---------|
| answer: | yes | no |

Multi-class Classification

| y | $y = 1$ | $y = 2$ | $y = 3$ |
|------------------|---------|---------|---------|
| answer: | yes | no | |
| parents' answer: | yes | no | maybe |

Multi-class Classification

| y | $y = 1$ | $y = 2$ | $y = 3$ | $y = 4$ |
|------------------|---------|---------|---------|---------|
| answer: | yes | no | | |
| parents' answer: | yes | no | maybe | |
| cardinal points: | North | South | East | West |

Multi-class Classification

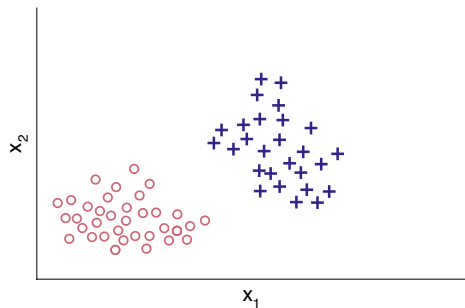
| y | $y = 1$ | $y = 2$ | $y = 3$ | $y = 4$ | $y = 5$ |
|------------------|---------|---------|---------|---------|---------|
| answer: | yes | no | | | |
| parents' answer: | yes | no | maybe | | |
| cardinal points: | North | South | East | West | |
| fruit: | apples | oranges | pears | bananas | mangos |

Multi-class Classification

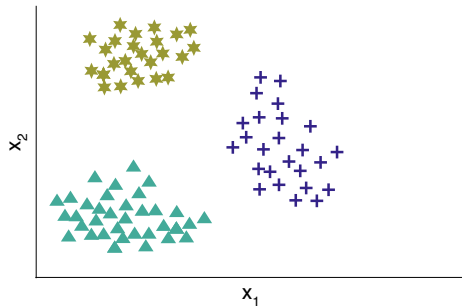
| y | $y = 1$ | $y = 2$ | $y = 3$ | $y = 4$ | $y = 5$ |
|------------------|---------|---------|---------|---------|---------|
| answer: | yes | no | | | |
| parents' answer: | yes | no | maybe | | |
| cardinal points: | North | South | East | West | |
| fruit: | apples | oranges | pears | bananas | mangos |
| ⋮ | | | | | |

binary vs multi-class classification

Binary Classification

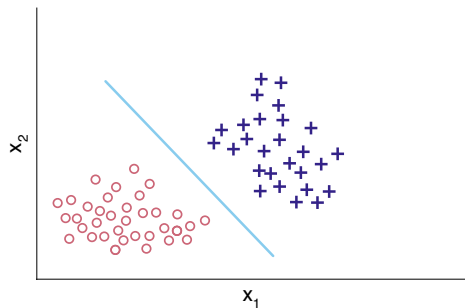


Multi-class Classification

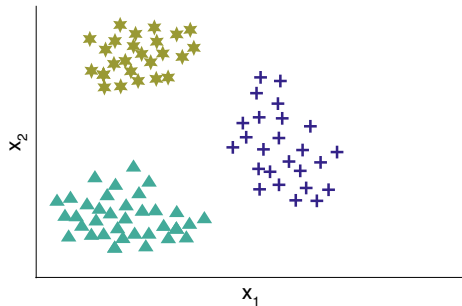


binary vs multi-class classification

Binary Classification

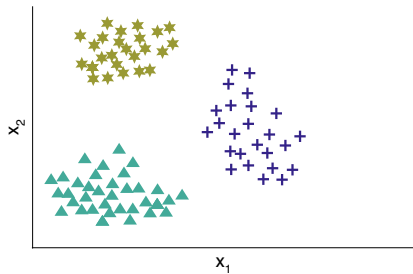


Multi-class Classification



One-vs-Rest (OvR)

Multi-class Classification



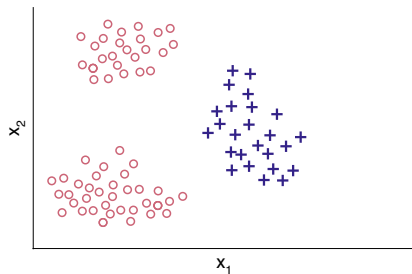
$+$ $y = 1$

\star $y = 2$

\blacktriangle $y = 3$

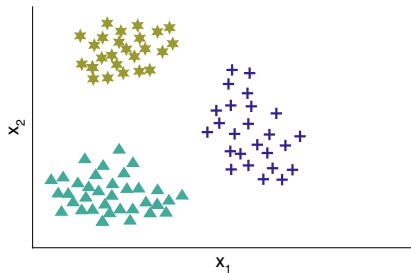
One-vs-Rest (OvR)

Hypothesis for Class 1: $h_{\theta}^{(1)}(x)$



$y = 1$ vs rest (i.e., $y = 2, 3$)

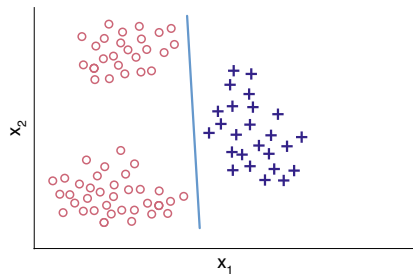
Multi-class Classification



+ $y = 1$
★ $y = 2$
▲ $y = 3$

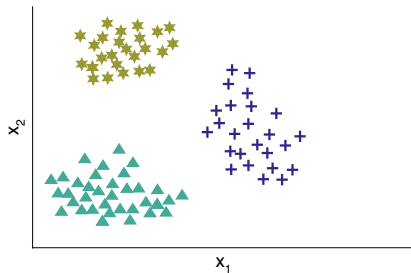
One-vs-Rest (OVR)

Hypothesis for Class 1: $h_{\theta}^{(1)}(x)$



$y = 1$ vs rest (i.e., $y = 2, 3$)

Multi-class Classification



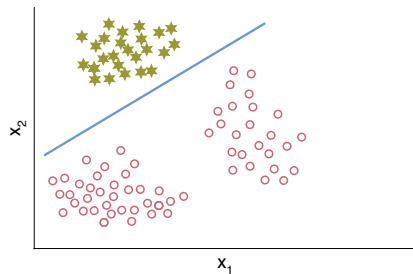
+ $y = 1$

★ $y = 2$

▲ $y = 3$

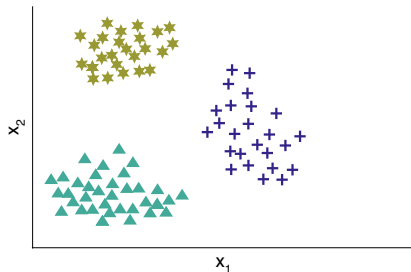
One-vs-Rest (OVR)

Hypothesis for Class 2: $h_{\theta}^{(2)}(x)$



$y = 2$ vs rest (i.e., $y = 1, 3$)

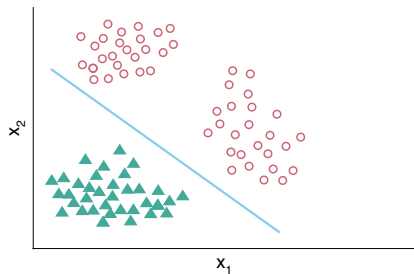
Multi-class Classification



- $+$ $y = 1$
- \star $y = 2$
- \blacktriangle $y = 3$

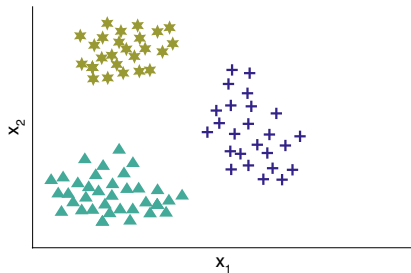
One-vs-Rest (OVR)

Hypothesis for Class 3: $h_{\theta}^{(3)}(x)$



$y = 3$ vs rest (i.e., $y = 1, 2$)

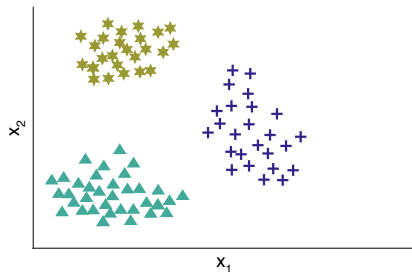
Multi-class Classification



- + $y = 1$
- ★ $y = 2$
- ▲ $y = 3$

One-vs-Rest (OVR)

Multi-class Classification



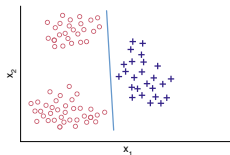
+ $y = 1$

★ $y = 2$

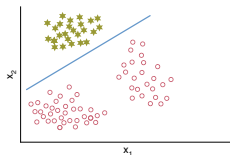
▲ $y = 3$

$$h_{\theta}^{(i)}(x) = P(y = i \mid x; \theta) \quad (\text{for } i = 1, 2, 3)$$

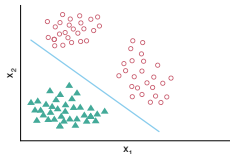
$h_{\theta}^{(1)}(x)$



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



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



One-vs-rest (OVR)

How to implement OvR:

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

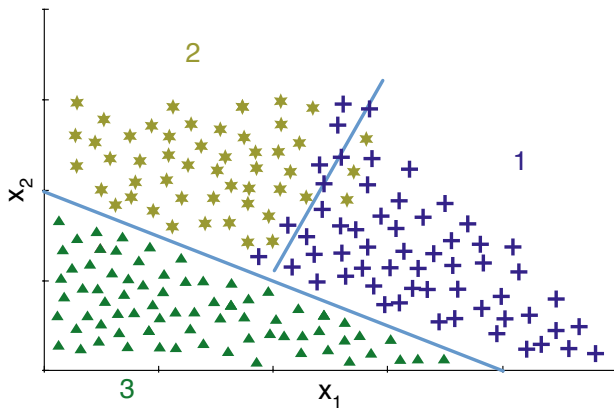
One-vs-rest (OVR)

How to implement OvR:

1. Train a **logistic regression classifier** $h_{\theta}^{(i)}(x)$ for **each** class i to predict the probability that $y = i$.
2. To make a prediction on a **new** input x , pick the class i that **maximizes**:

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

Example



$$h_{\theta}^{(1)}(x) = P(y = 1 \mid x; \theta)$$

$$h_{\theta}^{(2)}(x) = P(y = 2 \mid x; \theta)$$

$$h_{\theta}^{(3)}(x) = P(y = 3 \mid x; \theta)$$

References