

# Predictive Modelling

## Classification - Logistic Regression

Jonathan Mwaura

Khoury College of Computer Sciences

July 23, 2024

# Introduction

## Textbook

Reading: Chapter 4 of: Gareth James et al (2021) . An Introduction to Statistical Learning (2nd Edition) .

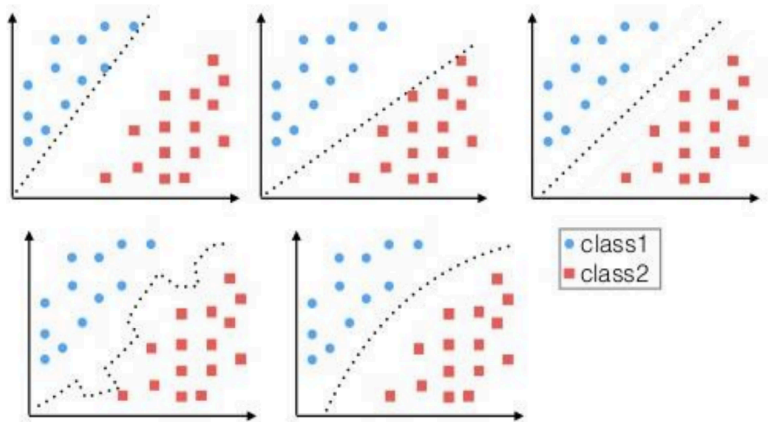
<https://www.statlearning.com/>

## Acknowledgements

These slides have been adapted from the following Professors:

- 1) Andrew Ng - Stanford
- 2) Eric Eaton - UPenn
- 3) David Sontag - MIT
- 4) Alina Oprea - Northeastern

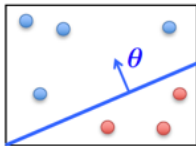
# Linear vs Non-Linear Classifiers



# Linear Classifiers

- **Linear classifiers:** represent decision boundary by hyperplane

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \quad x^\top = \begin{bmatrix} 1 & x_1 & \dots & x_d \end{bmatrix}$$





$h_\theta(x) = f(\theta^T x)$  linear function

- If  $\theta^T x > 0$  classify 1
- If  $\theta^T x < 0$  classify 0

All the points  $x$  on the hyperplane satisfy:  $\theta^T x = 0$

# Linear Classifiers

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

- Examples: perceptron, LDA
- Pros 
  - Very compact model (size  $d$ )
- Cons of linear models studied so far 
  - Perceptron depend on the order of training data and it could take many steps for convergence
  - LDA assumes normal distribution of features

# Classification Based on Probability

- Instead of just predicting the class, give the *probability of the instance being in that class*
  - Learn  $P(Y|X)$
- Consider binary classifier with classes 0 and 1
  - $P(Y = 1|X) + P(Y = 0|X) = 1$
  - Sufficient to learn  $P(Y = 1|X)$
- Advantages: interpretability and confidence of output

# Logistic Regression

- Setup

- Training data:  $\{x_i, y_i\}$ , for  $i = 1, \dots, N$
- Labels:  $y_i \in \{0, 1\}$

- Goals

- Learn  $P(Y = 1|X = x)$

- Highlights

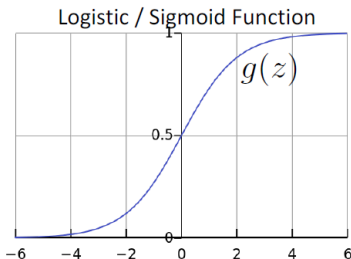
- Probabilistic output
- At the basis of more complex models (e.g., neural networks)
- Supports regularization (Ridge, Lasso)
- Can be trained with Gradient Descent

# Logistic Regression

- Takes a probabilistic approach to learning discriminative functions (i.e., a classifier)
- $h_{\theta}(x)$  should give  $P(Y = 1|X; \theta)$ 
  - Want  $0 \leq h_{\theta}(x) \leq 1$
- Logistic regression model:

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

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





# Interpretation of Model Output

$$h_{\theta}(\mathbf{x}) = \text{estimated} \quad P(Y = 1|X; \theta)$$

Example: Cancer diagnosis from tumor size

$$\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorSize} \end{bmatrix}$$

$$h_{\theta}(\mathbf{x}) = 0.7$$

→ Tell patient that 70% chance of tumor being malignant

Note that:  $P(Y = 0|X; \theta) + P(Y = 1|X; \theta) = 1$

Therefore,  $P(Y = 0|X; \theta) = 1 - P(Y = 1|X; \theta)$

# LR is a Linear Classifier!

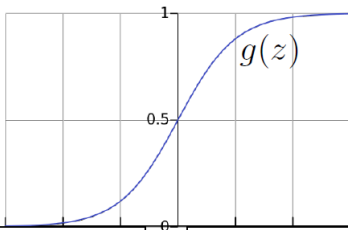
- Predict  $Y = 1$  if:
  - $P[Y = 1|X = x; \theta] > P[Y = 0|X = x; \theta]$
  - $P[Y = 1|X = x; \theta] > \frac{1}{2}$
  - $$\frac{1}{1 + e^{-\theta^T x}} > \frac{1}{2}$$
- Equivalent to:
  - $e^{\theta_0 + \sum_{j=1}^d \theta_j x_j} > 1$
  - $\theta_0 + \sum_{j=1}^d \theta_j x_j > 0$

Logistic Regression is a linear classifier!

# Logistic Regression

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

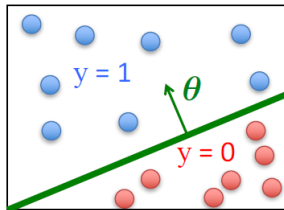
$$g(z) = \frac{1}{1 + e^{-z}}$$



$\theta^T x$  should be large negative values for negative instances

$\theta^T x$  should be large positive values for positive instances

- Assume a threshold and...
  - Predict  $Y = 1$  if  $h_{\theta}(x) \geq 0.5$
  - Predict  $Y = 0$  if  $h_{\theta}(x) < 0.5$



Logistic Regression is a linear classifier!

# Logistic Regression Objective

- Can't just use squared loss as in linear regression:

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N (h_{\theta}(x_i) - y_i)^2$$

- Using the logistic regression model

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

results in a non-convex optimization

# Maximum Likelihood Estimation (MLE)

Given training data  $X = \{x_1, \dots, x_N\}$  with labels  $Y = \{y_1, \dots, y_N\}$

What is the likelihood of training data for parameter  $\theta$ ?

Define **likelihood function**

$$\text{Max}_{\theta} L(\theta) = P[Y|X; \theta]$$

Assumption: training points are independent

$$L(\theta) = \prod_{i=1}^n P[Y = y_i | X = x_i; \theta]$$

# Log Likelihood

- Max likelihood is equivalent to maximizing log of likelihood

$$L(\theta) = \prod_{i=1}^N P[Y = y_i | X = x_i; \theta]$$

$$\log L(\theta) = \sum_{i=1}^N \log P[Y = y_i | X = x_i; \theta]$$

- They both have the same maximum  $\theta_{MLE}$

# MLE for Logistic Regression

$$P(Y = y_i | X = x_i; \theta) = h_{\theta}(x_i)^{y_i} (1 - h_{\theta}(x_i))^{1-y_i}$$

$$\begin{aligned}\theta_{MLE} &= \operatorname{argmax}_{\theta} \sum_{i=1}^N \log P[Y = y_i | X = x_i; \theta] \\ &= \operatorname{argmax}_{\theta} \sum_{i=1}^N y_i \log h_{\theta}(x_i) + (1 - y_i) \log (1 - h_{\theta}(x_i))\end{aligned}$$

Logistic regression objective

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

$$J(\theta) = - \sum_{i=1}^N [y_i \log h_{\theta}(x_i) + (1 - y_i) \log (1 - h_{\theta}(x_i))]$$

# Cross-Entropy Objective

$$J(\theta) = - \sum_{i=1}^N [y_i \log h_{\theta}(x_i) + (1 - y_i) \log (1 - h_{\theta}(x_i))]$$

- Cost of a single instance:

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

- Can re-write objective function as

$$J(\theta) = \sum_{i=1}^n \underbrace{\text{cost}(h_{\theta}(x_i), y_i)}_{\text{Cross-entropy loss}}$$

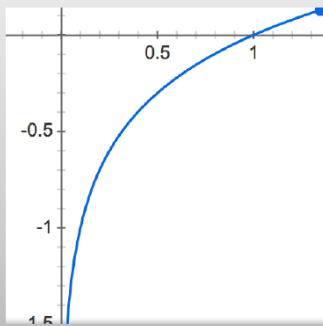
Cross-entropy loss



# Intuition

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

Aside: Recall the plot of  $\log(z)$

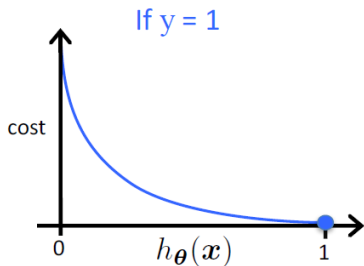


# Intuition

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

If  $y = 1$

- Cost = 0 if prediction is correct
- As  $h_{\theta}(x) \rightarrow 0$ ,  $\text{cost} \rightarrow \infty$
- Captures intuition that larger mistakes should get larger penalties
  - e.g., predict  $h_{\theta}(x) = 0$ , but  $y = 1$

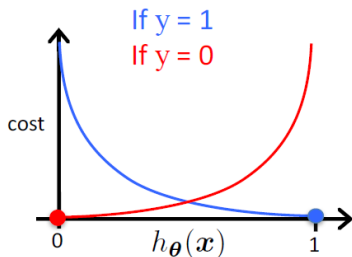


# Intuition

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

If  $y = 0$

- Cost = 0 if prediction is correct
- As  $(1 - h_{\theta}(\mathbf{x})) \rightarrow 0$ ,  $\text{cost} \rightarrow \infty$
- Captures intuition that larger mistakes should get larger penalties



# Gradient Descent for Logistic Regression

$$J(\theta) = - \sum_{i=1}^N [y_i \log h_{\theta}(x_i) + (1 - y_i) \log (1 - h_{\theta}(x_i))]$$

$$J(\theta) = - \sum_{i=1}^n C_i$$

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

- Initialize  $\theta$
- Repeat until convergence

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

simultaneous update  
for  $j = 0 \dots d$

# Computing Gradients

- Derivative of sigmoid

$$- g(z) = \frac{1}{1+e^{-z}}; g'(z) = \frac{e^{-z}}{(1+e^{-z})^2} = g(z)(1 - g(z))$$

- Derivative of hypothesis

$$- h_{\theta}(x) = g(\theta^T x) = g(\theta_j x_j + \sum_{k \neq j} \theta_k x_k)$$

$$- \frac{\partial h_{\theta}(x)}{\partial \theta_j} = \frac{\partial g(\theta^T x)}{\partial \theta_j} x_j = g(\theta^T x)(1 - g(\theta^T x))x_j$$

- Derivation of  $C_i$

$$\begin{aligned} - \frac{\partial C_i}{\partial \theta_j} &= y_i \frac{1}{h_{\theta}(x_i)} g(\theta^T x_i)(1 - g(\theta^T x_i))x_{ij} - \\ &\quad (1 - y_i) \frac{1}{1 - h_{\theta}(x_i)} g(\theta^T x_i)(1 - g(\theta^T x_i))x_{ij} \\ &= (y_i - h_{\theta}(x_i))x_{ij} \end{aligned}$$

# Gradient Descent for Logistic Regression

$$J(\theta) = - \sum_{i=1}^N [y_i \log h_{\theta}(x_i) + (1 - y_i) \log (1 - h_{\theta}(x_i))]$$

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

- Initialize  $\theta$
- Repeat until convergence (simultaneous update for  $j = 0 \dots d$ )

$$\theta_0 \leftarrow \theta_0 - \alpha \sum_{i=1}^N (h_{\theta}(x_i) - y_i)$$

$$\theta_j \leftarrow \theta_j - \alpha \sum_{i=1}^N (h_{\theta}(x_i) - y_i) x_{ij}$$

# Gradient Descent for Logistic Regression

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

- Initialize  $\theta$
- Repeat until convergence (simultaneous update for  $j = 0 \dots d$ )

$$\theta_0 \leftarrow \theta_0 - \alpha \sum_{i=1}^N (h_{\theta}(x_i) - y_i)$$

$$\theta_j \leftarrow \theta_j - \alpha \sum_{i=1}^N (h_{\theta}(x_i) - y_i) x_{ij}$$

**This looks IDENTICAL to Linear Regression!**

- However, the form of the model is very different:

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

# Regularized Logistic Regression

$$J(\theta) = - \sum_{i=1}^N [y_i \log h_{\theta}(x_i) + (1 - y_i) \log (1 - h_{\theta}(x_i))]$$

- We can regularize logistic regression exactly as before:

$$\begin{aligned} J_{\text{regularized}}(\theta) &= J(\theta) + \lambda \sum_{j=1}^d \theta_j^2 \\ &= J(\theta) + \lambda \|\theta_{[1:d]}\|_2^2 \end{aligned}$$

L2 regularization