

# Classification

December 18, 2017

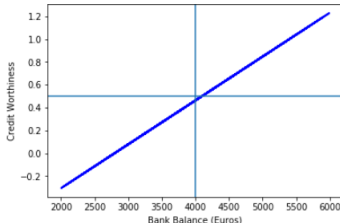
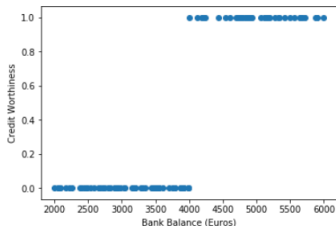
# Classification

- ▶ Linear regression: the response variable  $y$  is quantitative.
- ▶ Classification:  $y$  is qualitative (takes on a number of discrete values).
- ▶ Classification problems seem to occur more often than regression problems:
  - ▶ spam classifiers (spam or ham)
  - ▶ classifying whether a bank transaction is fraudulent or not
  - ▶ given a set of symptoms, determining which medical condition a person has
  - ▶ classifying whether a video is suitable or unsuitable for children
  - ▶ MNIST: given a handwritten digit, determine which digit it actually is

# Why not Linear Regression?

## Example 1

Consider the following (simplified) problem: given the bank balance  $x$  of an individual, determine whether they are credit worthy or not.

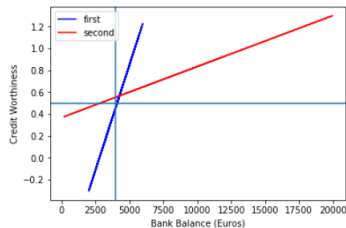
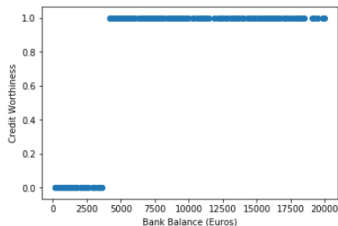


- ▶ Turns out that anyone with a balance of €4000 or more is credit worthy
- ▶ Classification: if  $y(x) \geq 0.5$ , then “credit worthy”; else “not”
- ▶ Slope of the regression line depends on the how many data points are in each of the two buckets

# Why not Linear Regression?

## Example 1

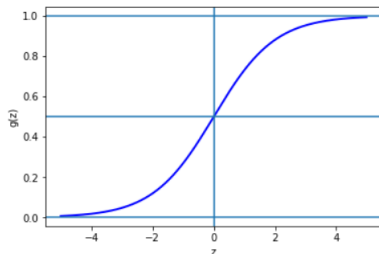
- ▶ With more data points in the “positive” bucket, the slope of the regression line is less steep.
- ▶ The threshold predicted also changes.



# Binary Classification: Logistic Regression

- ▶ Two classes: 0 and 1
- ▶ Logistic regression models the probability that the response variable  $y$  belongs to a particular class:  $P(y = 1 \mid x)$
- ▶  $P(y = 1 \mid x; \theta) = g(\theta^\top x) = \frac{1}{1 + e^{-\theta^\top x}}$
- ▶  $g(z) = \frac{1}{1 + e^{-\theta^\top x}}$  is the sigmoid function

# The Sigmoid Function



- ▶  $g(z) \rightarrow 1$  as  $z \rightarrow \infty$
- ▶  $g(z) \rightarrow 0$  as  $z \rightarrow -\infty$
- ▶  $g'(z) = g(z)(1 - g(z))$

# Logistic Regression: Learning the Model Parameters

In logistic regression, the probabilities are modeled as follows:

$$P(y = 1 \mid x; \theta) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

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

# Logistic Regression: Learning the Model Parameters

In logistic regression, the probabilities are modeled as follows:

$$P(y = 1 \mid x; \theta) = g(\theta^\top x) = \frac{1}{1 + e^{-\theta^\top x}}$$

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

$$P(y \mid x; \theta) = \left(g(\theta^\top x)\right)^y \cdot \left(1 - g(\theta^\top x)\right)^{1-y}$$



# Logistic Regression: Learning the Model Parameters

In logistic regression, the probabilities are modeled as follows:

$$P(y = 1 \mid x; \theta) = g(\theta^\top x) = \frac{1}{1 + e^{-\theta^\top x}}$$

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

$$P(y \mid x; \theta) = \left(g(\theta^\top x)\right)^y \cdot \left(1 - g(\theta^\top x)\right)^{1-y}$$

## The Likelihood Function

Assume that the  $m$  training examples  $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$  were generated *independently*.

$$\begin{aligned} L(\theta; \{(x^{(i)}, y^{(i)})\}_{i=1}^m) &= \prod_{i=1}^m P(y^{(i)} \mid x^{(i)}; \theta) \\ &= \prod_{i=1}^m \left(g(\theta^\top x^{(i)})\right)^{y^{(i)}} \cdot \left(1 - g(\theta^\top x^{(i)})\right)^{(1-y^{(i)})} \end{aligned}$$

# Maximizing the Likelihood

$$\begin{aligned}l(\theta) &= \log L(\theta) \\&= \sum_{i=1}^m y^{(i)} \log g(\theta^\top x^{(i)}) + (1 - y^{(i)}) \log(1 - g(\theta^\top x^{(i)}))\end{aligned}$$

# Stochastic Gradient Descent

In batch gradient descent, the update step for the  $j$ th component is:

$$\theta_j := \theta_j + \alpha \cdot \sum_{i=1}^m \left( y^{(i)} - \sum_{j=0}^n x_j^{(i)} \theta_j \right) \cdot x_j^{(i)}.$$

- ▶ Has to scan through the entire training set for a *single* update
- ▶ Costly operation if  $m$  is large

# Stochastic Gradient Descent

In batch gradient descent, the update step for the  $j$ th component is:

$$\theta_j := \theta_j + \alpha \cdot \sum_{i=1}^m \left( y^{(i)} - \sum_{j=0}^n x_j^{(i)} \theta_j \right) \cdot x_j^{(i)}.$$

- ▶ Has to scan through the entire training set for a *single* update
- ▶ Costly operation if  $m$  is large
- ▶ **Stochastic Gradient Descent:** for every training instance  $(x, y)$ ,  $x = (x_0, x_1, \dots, x_n)^T$ , update the parameters:

$$\theta_j := \theta_j + \alpha \cdot \left( y - \sum_{j=0}^n x_j \theta_j \right) \cdot x_j.$$

# Stochastic Gradient Descent: Features and Issues

- ▶ Doesn't have to look at the entire training set to make progress.
- ▶ Often gets close to the optimum much faster than batch gradient descent.
- ▶ May never converge to the optimum (can keep on oscillating between values near the optimum). This problem is alleviated by choosing  $\alpha$  to be very small.

# Analytic Solution

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m \left( y^{(i)} - \sum_{j=0}^n x_j^{(i)} \theta_j \right)^2$$

Want to find in closed-form a value of  $\theta$  that minimizes  $J(\theta)$

- Write  $J(\theta)$  in matrix-vector form.

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m \left( y^{(i)} - \sum_{j=0}^n x_j^{(i)} \theta_j \right)^2$$

Want to find in closed-form a value of  $\theta$  that minimizes  $J(\theta)$

- ▶ Write  $J(\theta)$  in matrix-vector form.
- ▶ **Design matrix.** An  $m \times (n + 1)$  matrix  $X$  defined by:

$$X = \begin{bmatrix} - & (x^{(1)})^T & - \\ - & (x^{(2)})^T & - \\ \vdots & \vdots & \vdots \\ - & (x^{(m)})^T & - \end{bmatrix}$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m \left( y^{(i)} - \sum_{j=0}^n x_j^{(i)} \theta_j \right)^2$$

Want to find in closed-form a value of  $\theta$  that minimizes  $J(\theta)$

- ▶ Write  $J(\theta)$  in matrix-vector form.
- ▶ **Design matrix.** An  $m \times (n + 1)$  matrix  $X$  defined by:

$$X = \begin{bmatrix} - & (x^{(1)})^T & - \\ - & (x^{(2)})^T & - \\ \vdots & \vdots & \vdots \\ - & (x^{(m)})^T & - \end{bmatrix}$$

- ▶  $y = (y^{(1)}, \dots, y^{(m)})^T$



$$y - X \cdot \theta = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} - \begin{bmatrix} (x^{(1)})^\top \cdot \theta \\ (x^{(2)})^\top \cdot \theta \\ \vdots \\ (x^{(m)})^\top \cdot \theta \end{bmatrix}$$

Matrix-form of  $J(\theta)$ :

$$\begin{aligned} J(\theta) &= \frac{1}{2} \sum_{i=1}^m \left( y^{(i)} - \sum_{j=0}^n x_j^{(i)} \theta_j \right)^2 \\ &= \frac{1}{2} (y - X\theta)^\top (y - X\theta) \end{aligned}$$

Minimize  $J(\theta)$  w.r.t  $\theta$ :

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \frac{1}{2} (y - X\theta)^{\top} (y - X\theta) \\ &= \text{see Andrew Ng's notes} \\ &= X^{\top} X \theta - X^{\top} y\end{aligned}$$

yielding:

$$\theta = (X^{\top} X)^{-1} \cdot X^{\top} y.$$

# Analytic Solution ...

Minimize  $J(\theta)$  w.r.t  $\theta$ :

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \frac{1}{2} (y - X\theta)^{\top} (y - X\theta) \\ &= \text{see Andrew Ng's notes} \\ &= X^{\top} X \theta - X^{\top} y\end{aligned}$$

yielding:

$$\theta = (X^{\top} X)^{-1} \cdot X^{\top} y.$$

► **Assumption:**  $X$  has full column rank so that  $X^{\top} X$  is invertible.

Minimize  $J(\theta)$  w.r.t  $\theta$ :

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \frac{1}{2} (y - X\theta)^{\top} (y - X\theta) \\ &= \text{see Andrew Ng's notes} \\ &= X^{\top} X \theta - X^{\top} y\end{aligned}$$

yielding:

$$\theta = (X^{\top} X)^{-1} \cdot X^{\top} y.$$

- **Assumption:**  $X$  has full column rank so that  $X^{\top} X$  is invertible.  
*Proof.* Show that  $X^{\top} X z = 0$  implies  $z = 0$ .

Minimize  $J(\theta)$  w.r.t  $\theta$ :

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \frac{1}{2} (y - X\theta)^{\top} (y - X\theta) \\ &= \text{see Andrew Ng's notes} \\ &= X^{\top} X \theta - X^{\top} y\end{aligned}$$

yielding:

$$\theta = (X^{\top} X)^{-1} \cdot X^{\top} y.$$

- ▶ **Assumption:**  $X$  has full column rank so that  $X^{\top} X$  is invertible.  
*Proof.* Show that  $X^{\top} X z = 0$  implies  $z = 0$ .
- ▶ If  $X$  does not have full column rank, the usual strategy is to remove redundant columns.

# Probabilistic Interpretation

## Assumptions

$$y^{(i)} = \theta^T \cdot x^{(i)} + \epsilon^{(i)},$$

where the error terms  $\epsilon^{(i)}$

- ▶ capture unmodeled effects and/or random noise
- ▶ are independent and identically distributed as  $N(0, \sigma^2)$

# Probabilistic Interpretation

## Assumptions

$$y^{(i)} = \theta^T \cdot x^{(i)} + \epsilon^{(i)},$$

where the error terms  $\epsilon^{(i)}$

- ▶ capture unmodeled effects and/or random noise
- ▶ are independent and identically distributed as  $N(0, \sigma^2)$

Given  $x^{(i)}$ ,

- ▶  $E(y^{(i)} \mid x^{(i)}) = \theta^T \cdot x^{(i)}$

# Probabilistic Interpretation

## Assumptions

$$y^{(i)} = \theta^T \cdot x^{(i)} + \epsilon^{(i)},$$

where the error terms  $\epsilon^{(i)}$

- ▶ capture unmodeled effects and/or random noise
- ▶ are independent and identically distributed as  $N(0, \sigma^2)$

Given  $x^{(i)}$ ,

- ▶  $E(y^{(i)} \mid x^{(i)}) = \theta^T \cdot x^{(i)}$
- ▶  $\text{Var}(y^{(i)} \mid x^{(i)}) = \text{Var}(\epsilon^{(i)}) = \sigma^2$



# Probabilistic Interpretation

## Assumptions

$$y^{(i)} = \theta^T \cdot x^{(i)} + \epsilon^{(i)},$$

where the error terms  $\epsilon^{(i)}$

- ▶ capture unmodeled effects and/or random noise
- ▶ are independent and identically distributed as  $N(0, \sigma^2)$

Given  $x^{(i)}$ ,

- ▶  $E(y^{(i)} \mid x^{(i)}) = \theta^T \cdot x^{(i)}$
- ▶  $\text{Var}(y^{(i)} \mid x^{(i)}) = \text{Var}(\epsilon^{(i)}) = \sigma^2$

Thus  $y^{(i)} \mid x^{(i)} \sim N(\theta^T x^{(i)}, \sigma^2)$ :

$$p(y^{(i)} \mid x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T \cdot x^{(i)})^2}{2\sigma^2}\right)$$

# Maximum Likelihood Estimation

- ▶ Given  $x = (x^{(1)}, \dots, x^{(m)})^\top$  and  $\theta$ , what is the joint distribution of the  $y = (y^{(1)}, \dots, y^{(m)})^\top$ ?

# Maximum Likelihood Estimation

- ▶ Given  $x = (x^{(1)}, \dots, x^{(m)})^\top$  and  $\theta$ , what is the joint distribution of the  $y = (y^{(1)}, \dots, y^{(m)})^\top$ ?
- ▶ Since the  $\epsilon^{(i)}$ s are independent:

# Maximum Likelihood Estimation

- ▶ Given  $x = (x^{(1)}, \dots, x^{(m)})^\top$  and  $\theta$ , what is the joint distribution of the  $y = (y^{(1)}, \dots, y^{(m)})^\top$ ?
- ▶ Since the  $\epsilon^{(i)}$ s are independent:

$$p(y \mid x; \theta) = \prod_{i=1}^m p(y^{(i)} \mid x^{(i)}; \theta)$$

# Maximum Likelihood Estimation

- ▶ Given  $x = (x^{(1)}, \dots, x^{(m)})^\top$  and  $\theta$ , what is the joint distribution of the  $y = (y^{(1)}, \dots, y^{(m)})^\top$ ?
- ▶ Since the  $\epsilon^{(i)}$ s are independent:

$$p(y \mid x; \theta) = \prod_{i=1}^m p(y^{(i)} \mid x^{(i)}; \theta)$$

- ▶ **Likelihood function.**  $L(\theta) = p(y \mid x; \theta)$

# Maximum Likelihood Estimation

- ▶ Given  $x = (x^{(1)}, \dots, x^{(m)})^\top$  and  $\theta$ , what is the joint distribution of the  $y = (y^{(1)}, \dots, y^{(m)})^\top$ ?
- ▶ Since the  $\epsilon^{(i)}$ s are independent:

$$p(y \mid x; \theta) = \prod_{i=1}^m p(y^{(i)} \mid x^{(i)}; \theta)$$

- ▶ **Likelihood function.**  $L(\theta) = p(y \mid x; \theta)$

$$L(\theta) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(y^{(i)} - \theta^\top x^{(i)})^2}{2\sigma^2} \right\}$$

# Maximum Likelihood Estimation

- ▶ Given  $x = (x^{(1)}, \dots, x^{(m)})^\top$  and  $\theta$ , what is the joint distribution of the  $y = (y^{(1)}, \dots, y^{(m)})^\top$ ?
- ▶ Since the  $\epsilon^{(i)}$ s are independent:

$$p(y \mid x; \theta) = \prod_{i=1}^m p(y^{(i)} \mid x^{(i)}; \theta)$$

- ▶ **Likelihood function.**  $L(\theta) = p(y \mid x; \theta)$

$$L(\theta) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(y^{(i)} - \theta^\top x^{(i)})^2}{2\sigma^2} \right\}$$

- ▶ **Principle of Maximum Likelihood.** Choose the parameters to make the data as likely as possible. Choose  $\theta$  to maximize  $L(\theta)$ .

# Maximum Likelihood Estimation ...

Maximizing  $L(\theta)$  is equivalent to maximizing *any* strictly increasing function of  $L(\theta)$ .



# Maximum Likelihood Estimation ...

Maximizing  $L(\theta)$  is equivalent to maximizing *any* strictly increasing function of  $L(\theta)$ .

- ▶ Usual to maximize the log likelihood  $l(\theta) = \log L(\theta)$ .
- ▶  $l(\theta) = m \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{\sigma^2} \cdot \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^\top x^{(i)})^2$ .
- ▶ Maximizing  $l(\theta)$  is equivalent to minimizing  $\frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^\top x^{(i)})^2$ .

# Maximum Likelihood Estimation ...

Maximizing  $L(\theta)$  is equivalent to maximizing *any* strictly increasing function of  $L(\theta)$ .

- ▶ Usual to maximize the log likelihood  $l(\theta) = \log L(\theta)$ .
- ▶  $l(\theta) = m \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{\sigma^2} \cdot \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$ .
- ▶ Maximizing  $l(\theta)$  is equivalent to minimizing  $\frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$ .

## Summary

Under the previous probabilistic assumptions: **least-squares regression** corresponds to finding the **maximum likelihood estimate** of  $\theta$ .

# The Goodness of Fit

- ▶ Residual Standard Error: standard deviation of the error terms  $\epsilon^{(i)}$ .

$$\text{RSE} = \sqrt{\frac{1}{n-2} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2}$$

# The Goodness of Fit

- ▶ Residual Standard Error: standard deviation of the error terms  $\epsilon^{(i)}$ .

$$\text{RSE} = \sqrt{\frac{1}{n-2} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2}$$

- ▶  $R^2$  Statistic:

$$R^2 = 1 - \frac{\sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2}{\sum_{i=1}^m (y^{(i)} - \bar{y})^2}$$

# The Goodness of Fit

- ▶ Residual Standard Error: standard deviation of the error terms  $\epsilon^{(i)}$ .

$$\text{RSE} = \sqrt{\frac{1}{n-2} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2}$$

- ▶  $R^2$  Statistic:

$$R^2 = 1 - \frac{\sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2}{\sum_{i=1}^m (y^{(i)} - \bar{y})^2}$$

- ▶ Total sum of squares =  $\sum_{i=1}^m (y^{(i)} - \bar{y})^2$ : variability inherent in the response
- ▶ Residual sum of squares =  $\sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$ : variability left unexplained after performing the regression