

Linear and Logistic Regression

David I. Inouye

Outline

- Linear regression
 - Formalization
 - Intuitions
 - Solution in closed-form
- Logistic regression
 - Intuitions
 - Formalization
 - Solution requires numerical algorithms

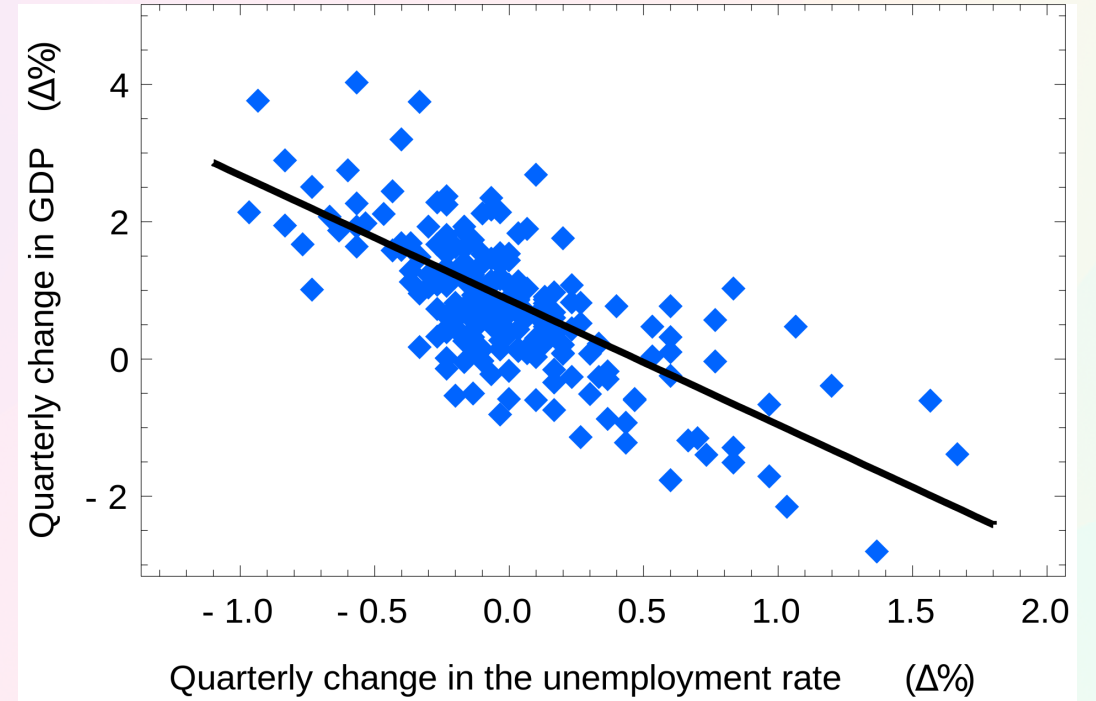
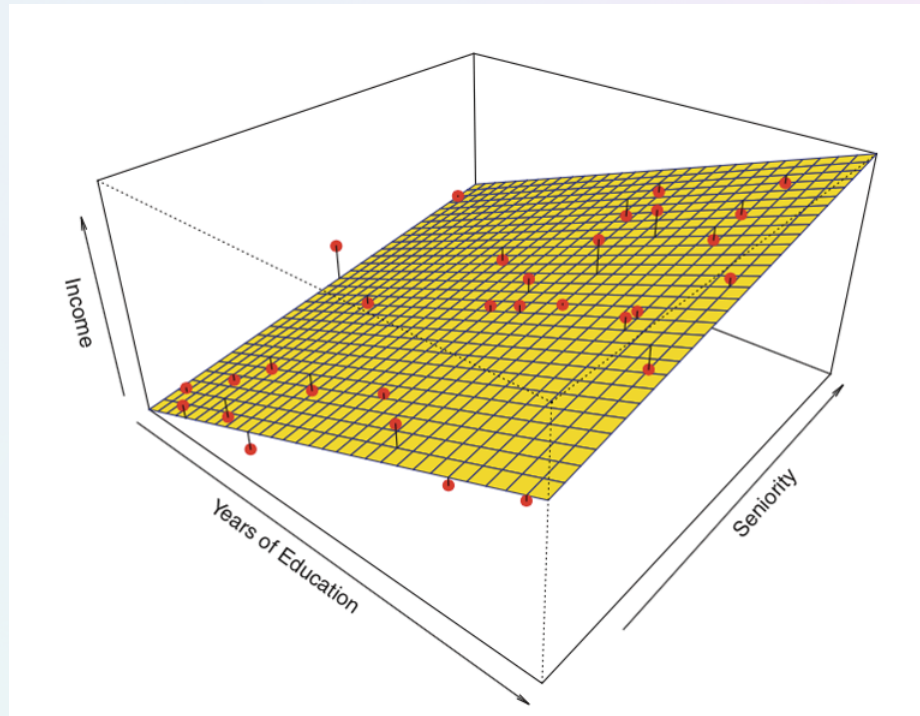
The Linear Regression Model Is Defined by the Coefficients (or **Parameters**) for Each Feature

- A simple linear combination where θ are the parameters

$$f_{\theta}(x) = \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_d x_d + \theta_{d+1}$$

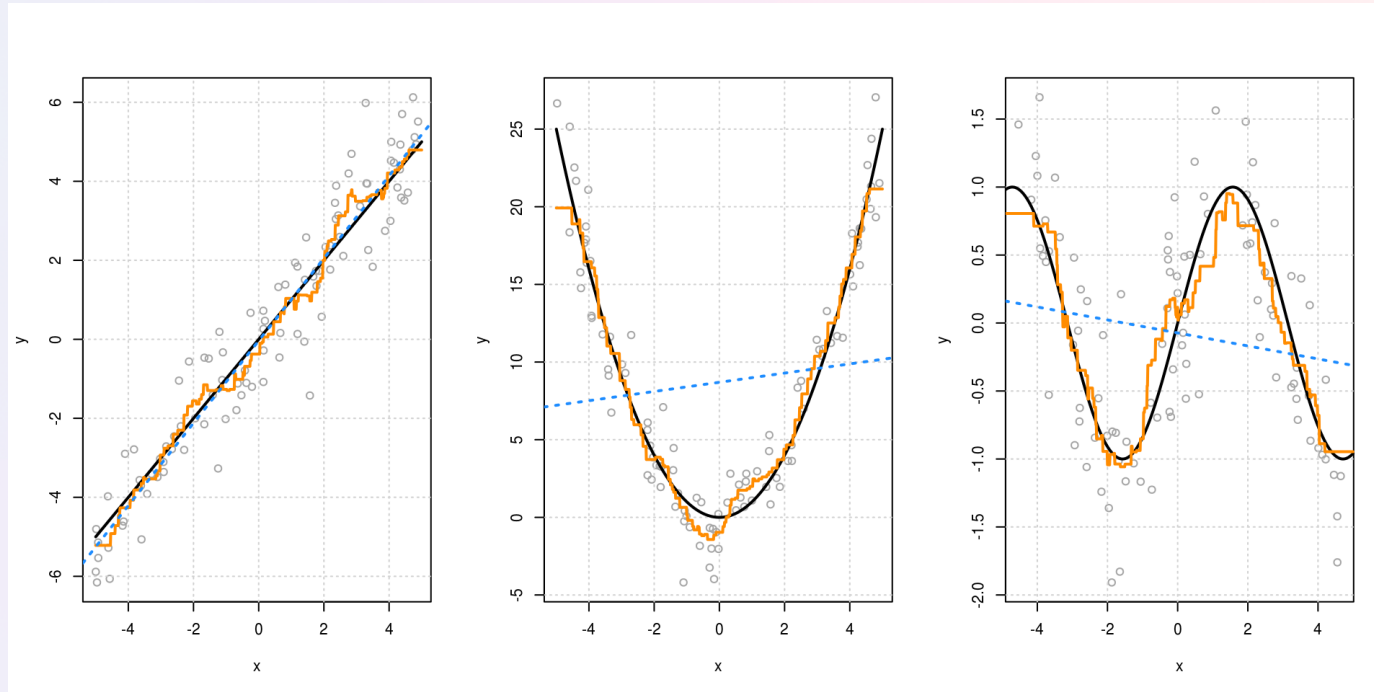
- Letting $x = [x_1, x_2, \dots, x_d, 1]$, we can write as $f_{\theta}(x) = \theta^T x$
- This is known as a *parametric model*

Linear Regression Models the Output as a Line (1D) or Hyperplane (>1D)



<https://towardsdatascience.com/linear-regression-detailed-view-ea73175f6e86>

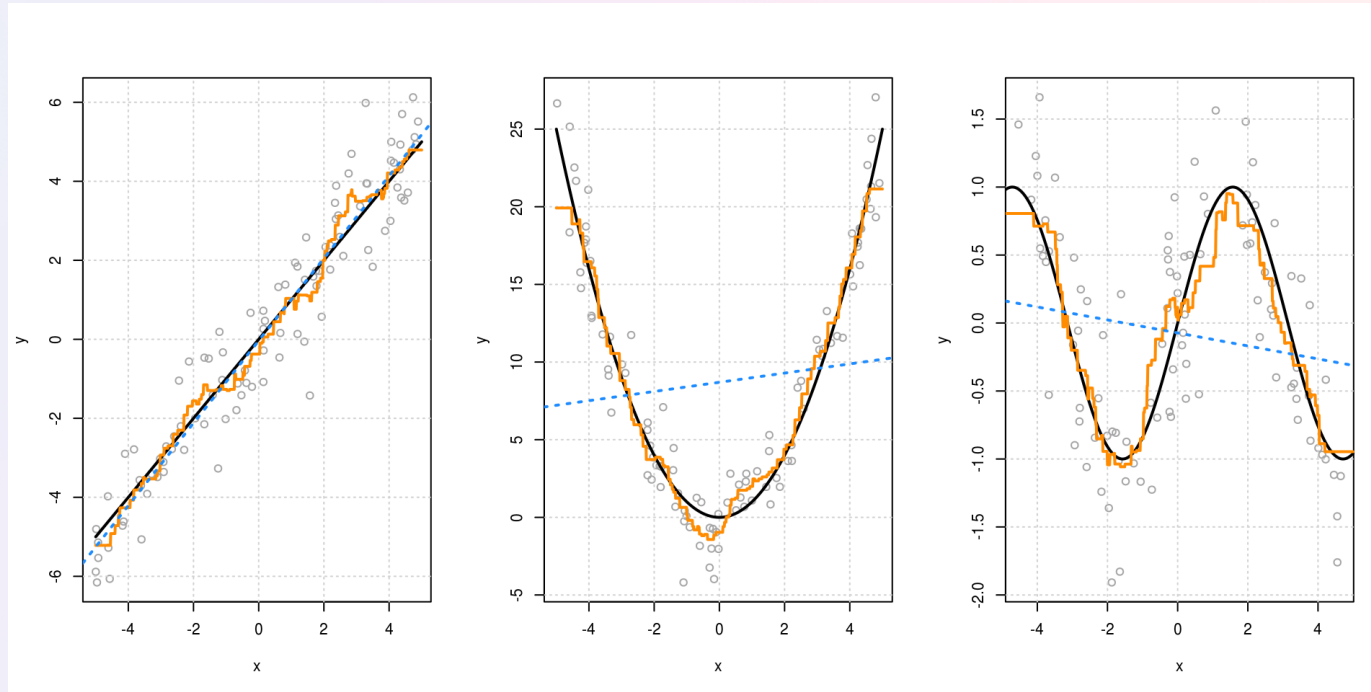
How Does This Compare to KNN Regression?



If KNN regression is a much more flexible model, is there ever a reason to choose linear models for 1D regression? — Discuss with partner.

<https://daviddalpiaz.github.io/r4sl/knn-reg.html>

How Does This Compare to KNN Regression?



If true phenomena is linear (i.e., *assumption matches reality*), linear regression will do the best (left). However, if true phenomena is not linear, KNN regression will perform better. (Black line is true function, **dotted blue line** is best linear approximation, and **orange line** is KNN regression.)

- Linear Regression Is a Much Simpler Function

<https://daviddalpiaz.github.io/r4sl/knn-reg.html>

The Goal of Linear Regression Is to Find the Parameters θ That Minimize the Prediction Error

Using mean squared error (MSE) this means:

$$\theta^* = \operatorname{argmin}_{\theta} \frac{1}{n} \sum_{i=1}^n (y_i - f_{\theta}(x_i))^2$$

Or equivalently

$$\theta^* = \operatorname{argmin}_{\theta} \sum_{i=1}^n (y_i - \theta^T x_i)^2$$

Or in matrix form

$$\theta^* = \operatorname{argmin}_{\theta} \|y - X\theta\|_2^2$$

Known as **Ordinary Least Squares (OLS)**

The Solution for OLS Can Be Computed in Closed Form

How do you find maximum or minimum in calculus?

- Calculate gradient

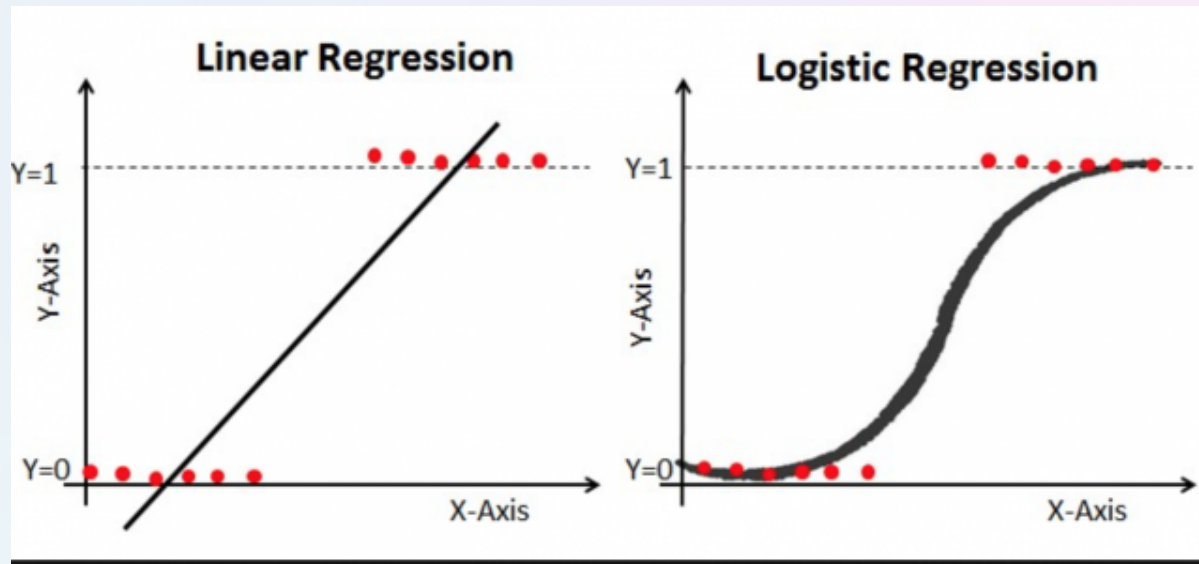
$$\begin{aligned}\nabla_{\theta} \|y - X\theta\|_2^2 &= (2(y - X\theta)^T (-X))^T \\ &= 2(-X^T)(y - X\theta) \\ &= 2(-X^T y + X^T X\theta)\end{aligned}$$

- Set equal to zero and solve (**Known as normal equations**)

$$\begin{aligned}2(-X^T y + X^T X\theta) &= 0 \\ X^T X\theta &= X^T y \\ \theta^* &= (X^T X)^{-1} X^T y\end{aligned}$$

Derivation hints: Use equivalence of $\|v\|_2^2 = v^T v$. Then use matrix calculus ([wikipedia reference](#)).

Logistic Regression Generalizes Linear Regression to the Classification Setting (Despite the Name)



- Output is probability of $Y = 1$
- Note that it is always between 0 and 1

The **Logistic Function** Is a Sigmoid Curve With a Simple Form

- $\sigma(a) = \frac{1}{1+e^{-a}}$

- Equivalently

$$\sigma(a) = \frac{e^a}{e^a+1}$$

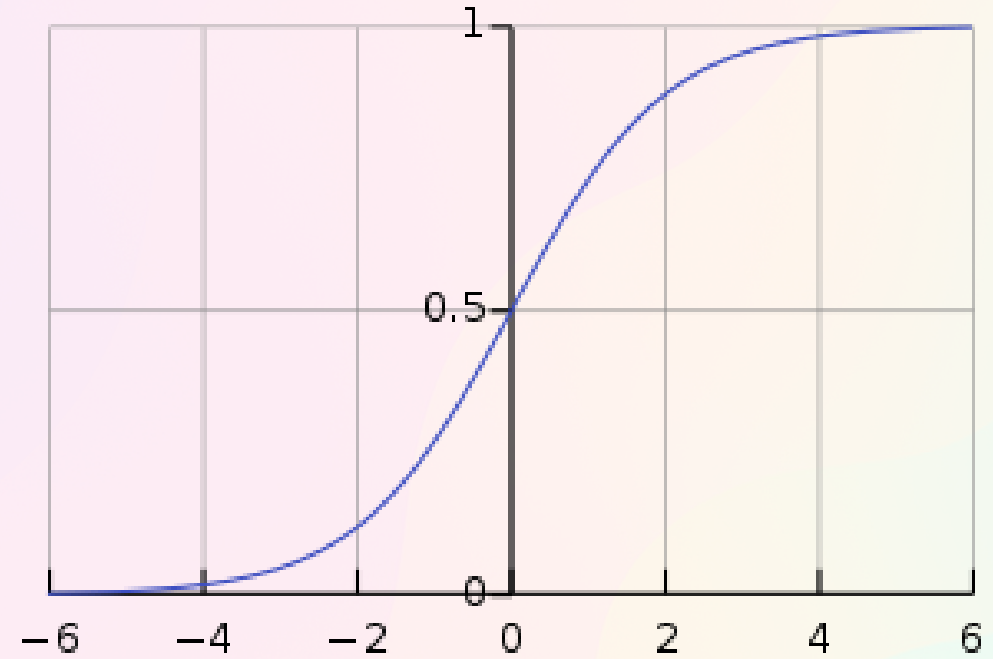
- Bounds

- $a \rightarrow \infty, \sigma(a) \rightarrow 1$

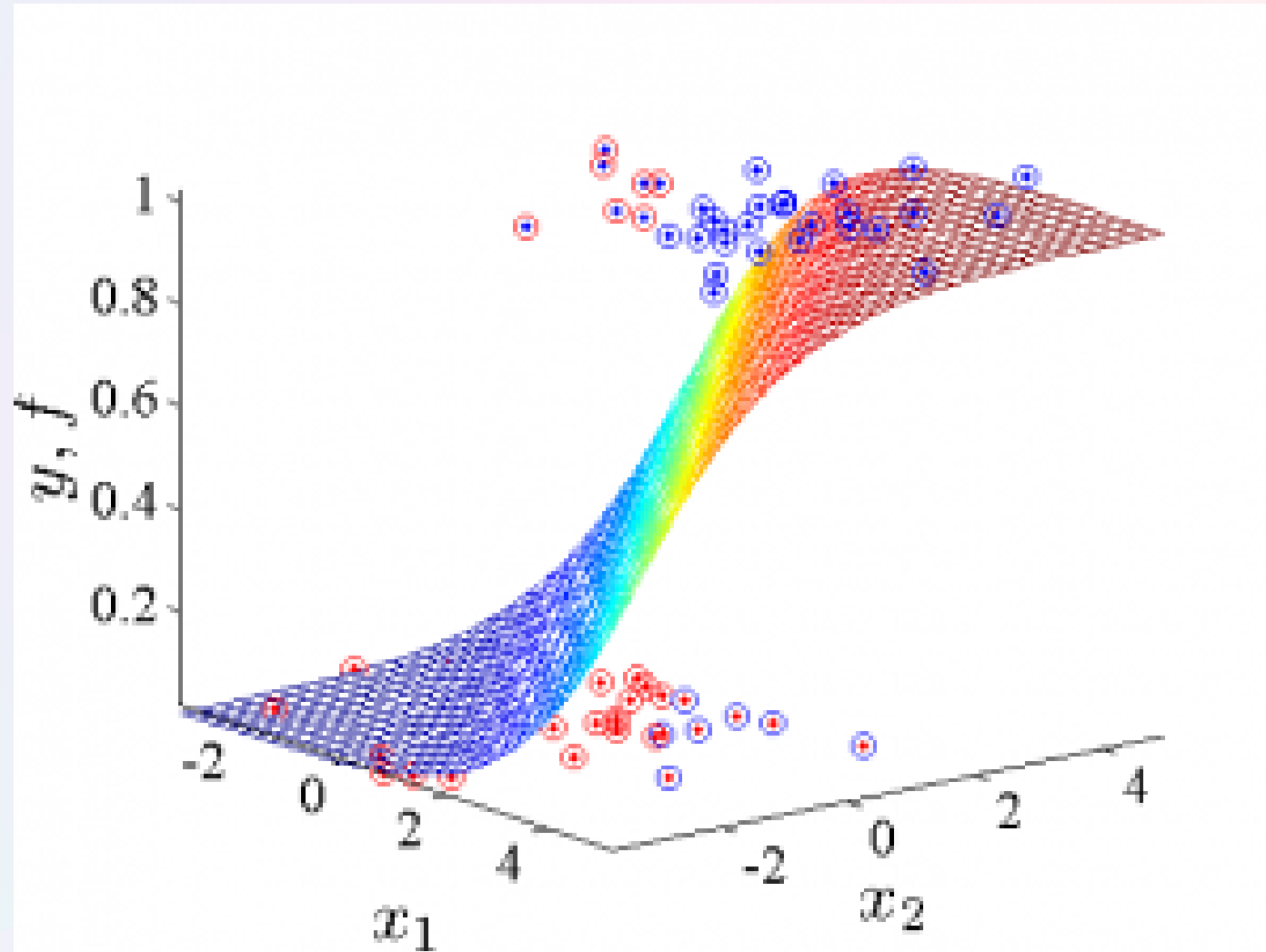
- $a \rightarrow -\infty, \sigma(a) \rightarrow 0$

- 1D logistic model

$$f_{\theta}(x) = \sigma(\theta_1 x + \theta_2)$$



Logistic Regression in Higher Dimensions Is Just the Logistic Curve Along a Single Direction



Multivariate Logistic Regression Merely Applies a Logistic Function to the Output of a Linear Function

- The multivariate logistic regression model is $f_{\theta}(x) = \sigma(\theta^T x)$
- Notice similarity to linear regression model
- However, we can interpret $f_{\theta}(x)$ as the **probability** of $y = 1$ instead of predicting y directly
- Thresholding this probability allows us to predict the class

$$\hat{y} = \begin{cases} 1, & \text{if } f_{\theta}(x) \geq 0.5 \\ 0, & \text{otherwise} \end{cases}$$

The Logistic Regression Optimization Maximizes the Log Likelihood of the Training Data

- In theory, we could use MSE:

$$\theta^* = \operatorname{argmin}_{\theta} \|y - \sigma(X\theta)\|_2^2$$

- However, the true output y is always 0 or 1
- Instead we maximize the *log likelihood* (which is equal to the log probability of the data)

$$\theta^* = \operatorname{argmax}_{\theta} \sum_{i=1}^n [y_i \log \operatorname{Pr}(y_i = 1|x_i) + (1 - y_i) \log \operatorname{Pr}(y_i = 0|x_i)]$$

- Equivalently

$$\theta^* = \operatorname{argmax}_{\theta} \sum_{i=1}^n [y_i \log \sigma(\theta^T x_i) + (1 - y_i) \log(1 - \sigma(\theta^T x_i))]$$

Logistic Regression Does **Not** Have a Closed-Form Solution!

- Must resort to numerical optimization
- Examples: Gradient descent, Newton's method

