# Internet of Things

Lecture 6: Logistic Regression with PySpark

# LOGISTIC REGRESSION

# Background

- We want to learn about Logistic Regression as a method for **Classification.**
- Some examples of classification problems:
  - Spam versus "Ham" emails
  - Loan Default (yes/no)
  - Disease Diagnosis
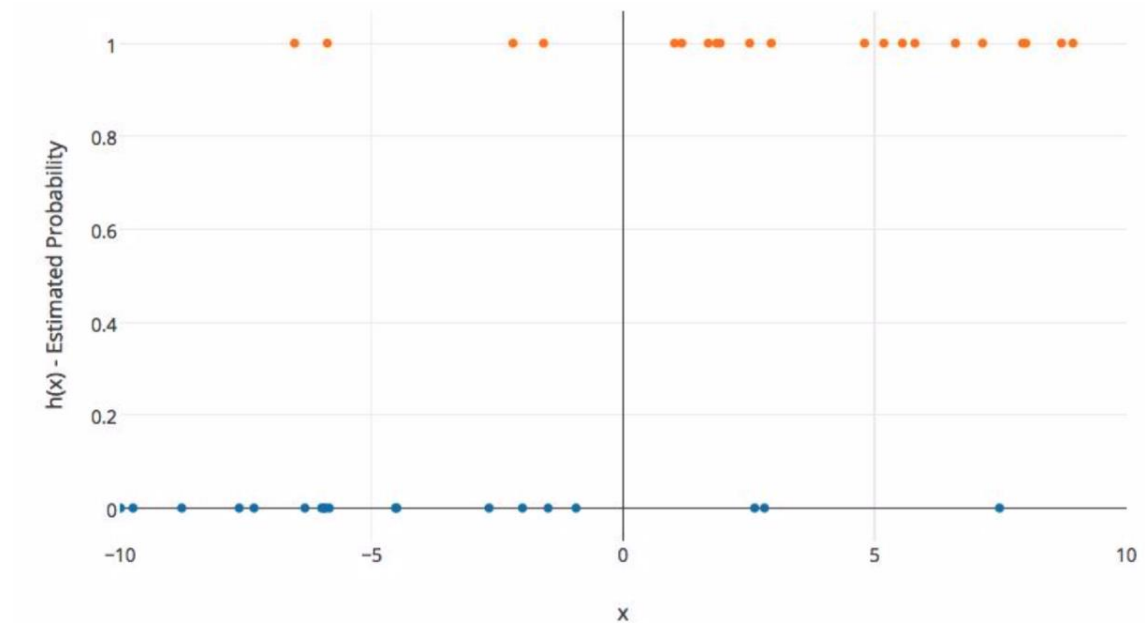- Above were all examples of Binary Classification

# Background

- The convention for binary classification is to have two classes 0 and 1.
- Let's walk through the basic idea for logistic regression.
- We'll also explain why it has the term regression in it, even though it's used for classification!
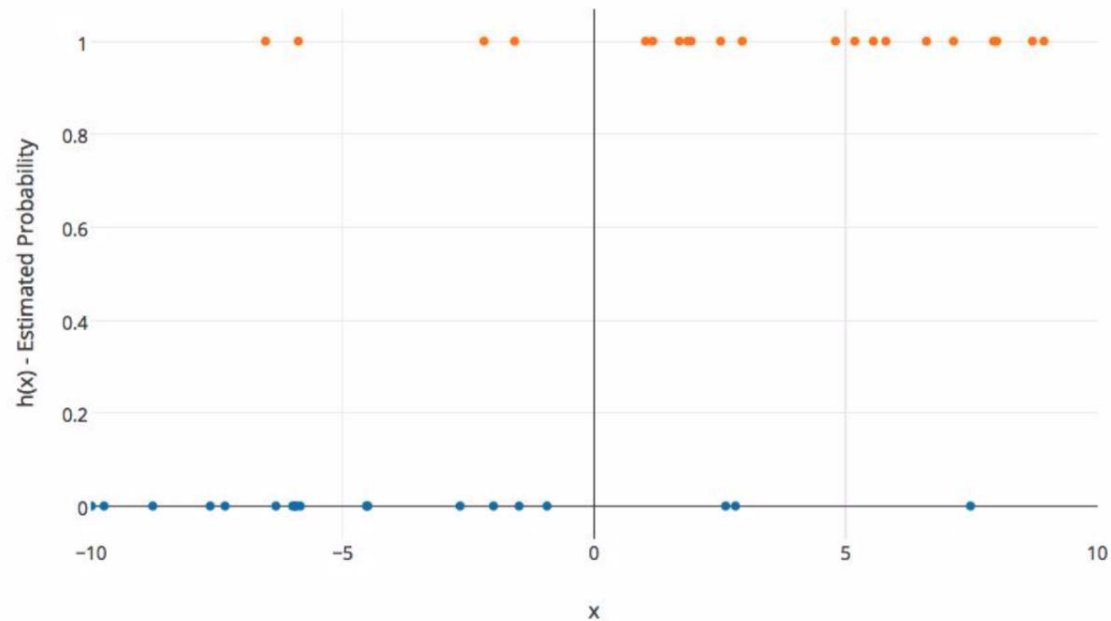
# Background

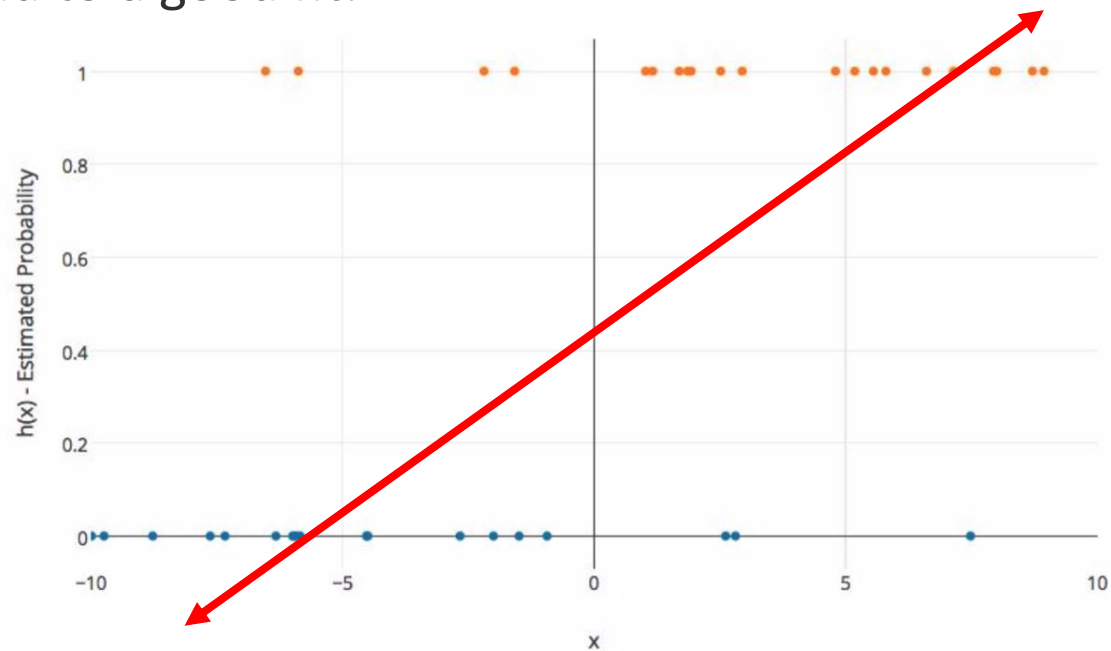- Imagine we plotted out some categorical data against one feature.

# Background

- The X axis represents a feature value and the Y axis represents the probability of belonging to class 1.
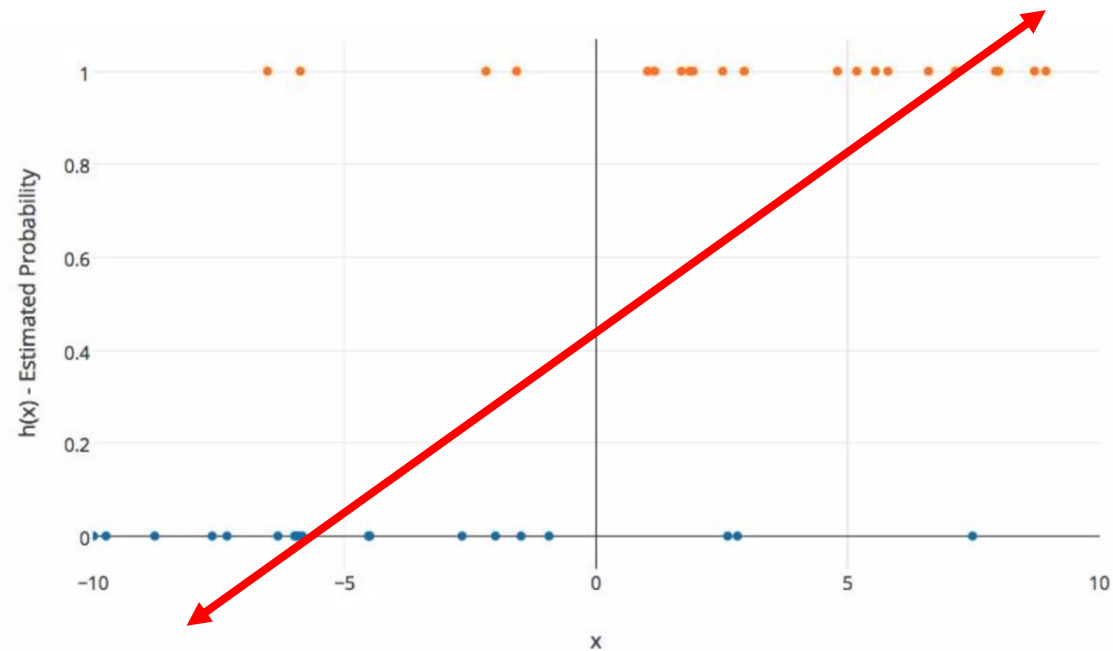
# Background

- We can't use a normal linear regression model on binary groups. It won't lead to a good fit:
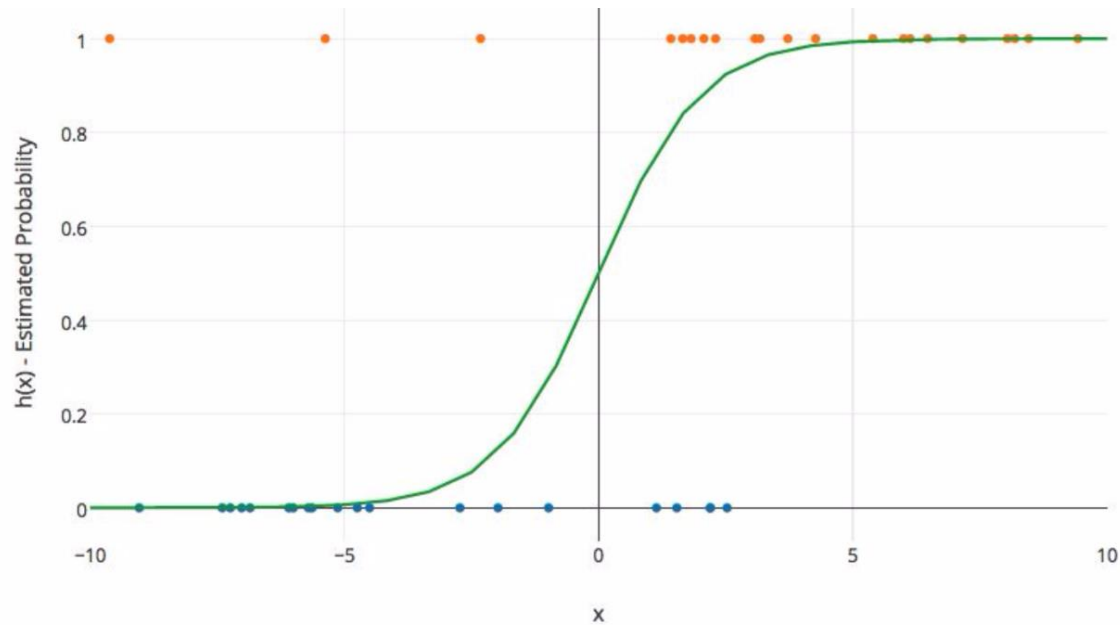
# Background

● We need a function that will fit binary categorical data!
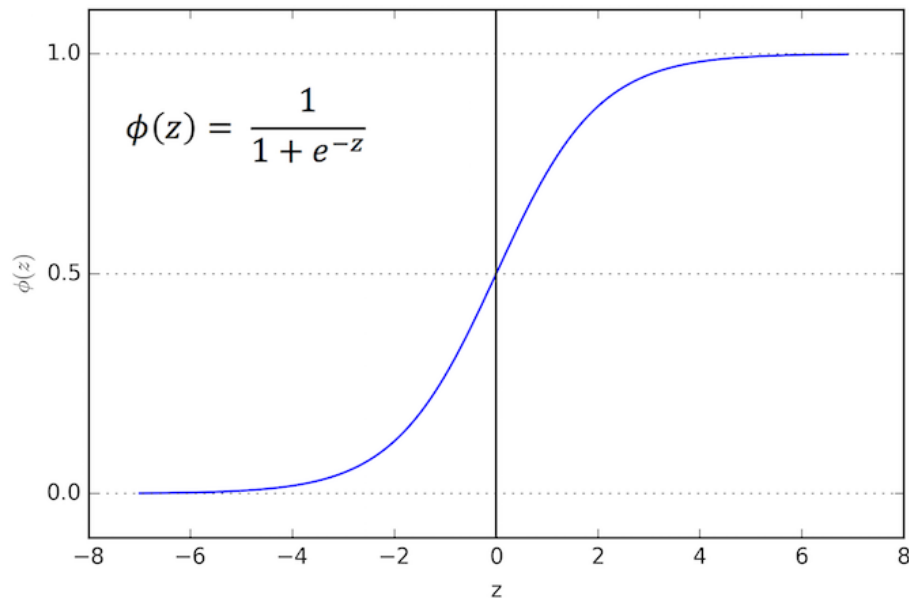
# Background

- It would be great if we could find a function with this sort of behavior:

# Sigmoid Function

- The Sigmoid (aka Logistic) Function takes in any value and outputs it to be between 0 and 1.

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

# Sigmoid Function

- This means we can take our Linear Regression Solution and place it into the Sigmoid Function.

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

# Sigmoid Function

- This results in a probability from 0 to 1 of belonging in the 1 class.

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

# Sigmoid Function

- We can set a cutoff point at 0.5, anything below it results in class 0, anything above is class 1.



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

# Review

- We use the logistic function to output a value ranging from 0 to 1. Based off of this probability we assign a class.

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

# Model Evaluation

- After you train a logistic regression model on some training data, you will evaluate your model's performance on some test data.
- You can use a confusion matrix to evaluate classification models.

# Model Evaluation

- The main point to remember with the confusion matrix and the various calculated metrics is that they are all fundamentally ways of comparing the predicted values versus the true values.
- What constitutes "good" metrics, will really depend on the specific situation!

# Model Evaluation

- We can use a confusion matrix to evaluate our model.
- For example, imagine testing for disease.

| n=165 | Predicted: NO | Predicted: YES |
|---|---|---|
| Actual: NO | 50 | 10 |
| Actual: YES | 5 | 100 |

Example: Test for presence of disease
NO = negative test = False = 0
YES = positive test = True = 1

# Confusion Matrix

| n=165 | Predicted: NO | Predicted: YES | |
|---|---|---|---|
| Actual: NO | TN = 50 | FP = 10 | 60 |
| Actual: YES | FN = 5 | TP = 100 | 105 |
| | 55 | 110 | |

Basic Terminology:
- True Positives (TP)
- True Negatives (TN)
- False Positives (FP)
- False Negatives (FN)

# Confusion Matrix

| n=165 | Predicted: NO | Predicted: YES | |
|---|---|---|---|
| Actual: NO | TN = 50 | FP = 10 | 60 |
| Actual: YES | FN = 5 | TP = 100 | 105 |
| | 55 | 110 | |

Accuracy:
- Overall, how often is it **correct**?
- (TP + TN) / total = 150/165 = 0.91

# Confusion Matrix

| n=165 | Predicted: NO | Predicted: YES | |
|---|---|---|---|
| **Actual: NO** | TN = 50 | FP = 10 | 60 |
| **Actual: YES** | FN = 5 | TP = 100 | 105 |
| | 55 | 110 | |

Misclassification Rate (Error Rate):
- Overall, how often is it **wrong**?
- (FP + FN) / total = 15/165 = 0.09

# Confusion Matrix

# Confusion Matrix

|  | total population | predicted condition | |
|---|---|---|---|
|  |  | prediction positive | prediction negative |
| **true condition** | condition positive | **True Positive (TP)** | **False Negative (FN)** (type II error) |
|  | condition negative | **False Positive (FP)** (Type I error) | **True Negative (TN)** |

# Confusion Matrix

| | predicted condition | | |
|---|---|---|---|
| total population | prediction positive | prediction negative | Prevalence $= \dfrac{\Sigma\,\text{condition positive}}{\Sigma\,\text{total population}}$ |
| **true condition** — condition positive | **True Positive (TP)** | **False Negative (FN)** (type II error) | True Positive Rate (TPR), Sensitivity, Recall, Probability of Detection $= \dfrac{\Sigma\,\text{TP}}{\Sigma\,\text{condition positive}}$ |
| condition negative | **False Positive (FP)** (Type I error) | **True Negative (TN)** | False Positive Rate (FPR), Fall-out, Probability of False Alarm $= \dfrac{\Sigma\,\text{FP}}{\Sigma\,\text{condition negative}}$ |
| Accuracy $= \dfrac{\Sigma\,\text{TP} + \Sigma\,\text{TN}}{\Sigma\,\text{total population}}$ | Positive Predictive Value (PPV), Precision $= \dfrac{\Sigma\,\text{TP}}{\Sigma\,\text{prediction positive}}$ | False Omission Rate (FOR) $= \dfrac{\Sigma\,\text{FN}}{\Sigma\,\text{prediction negative}}$ | Positive Likelihood Ratio (LR+) $= \dfrac{\text{TPR}}{\text{FPR}}$ |
| | False Discovery Rate (FDR) $= \dfrac{\Sigma\,\text{FP}}{\Sigma\,\text{prediction positive}}$ | Negative Predictive Value (NPV) $= \dfrac{\Sigma\,\text{TN}}{\Sigma\,\text{prediction negative}}$ | Negative Likelihood Ratio (LR–) $= \dfrac{\text{FNR}}{\text{TNR}}$ |