

Introduction to Logistic Regression

Ivan Corneillet

Data Scientist

Learning Objectives

After this lesson, you should be able to:

- Build a logistic regression classification model using *sklearn*
- Describe the logit and sigmoid functions, odds, and odds ratios as well as how they relate to logistic regression
- Evaluate a model using metrics such as classification accuracy/error



DS

Announcements and Exit Tickets

A black circle containing the white text "DS".

DS

Q & A



DS

Review

KNN | Pros and cons

▸ Pros

- Intuitive and simple to explain
- Training phase is fast
- Non-parametric (does not presume a “form” of the “decision boundary”)
- Easily capture non-linearity

▸ Cons

- Not interpretable
- Prediction phase can be slow when n (number of observations) is large
- Very sensitive to feature scaling; need to standardize the data
- Sensitive to irrelevant features
- Cannot be used if you have sparse data and feature space with dimension ≥ 4

Linear Regression | Interactions

features x_1 and x_2

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2$$

$f(x_1)$ x_2 is a constant

$y \rightarrow y + \Delta y$ $x_2 \rightarrow x_2 + \Delta x_2$

$y_{\text{old}} \quad y_{\text{new}} \quad x_{2,\text{old}} \quad x_{2,\text{new}}$

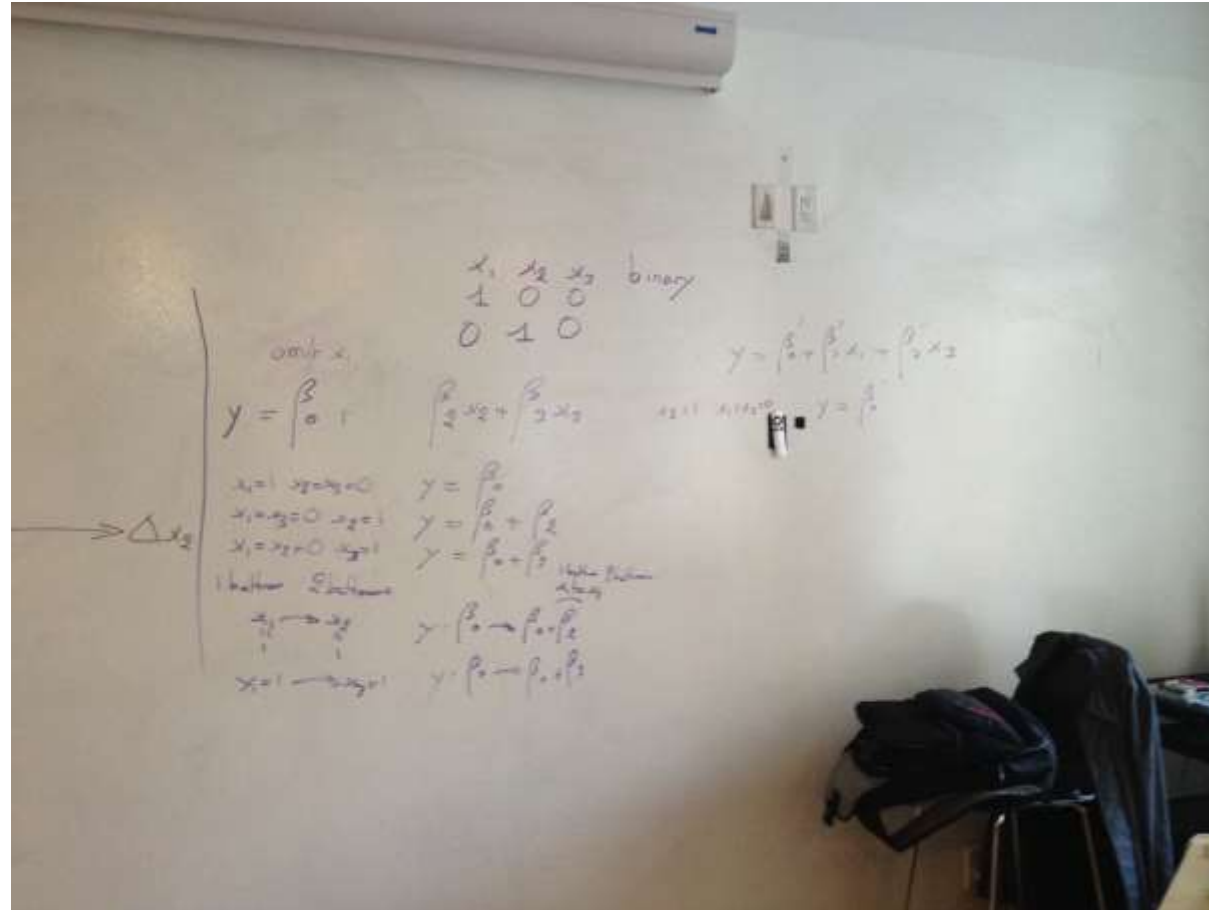
$$\Delta y = y_{\text{new}} - y_{\text{old}} = (\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2) - (\beta_0 + \beta_1 x_1 + \beta_2 x_{2,\text{old}} + \beta_3 x_1 x_{2,\text{old}})$$
$$= \beta_2 \Delta x_2 + \beta_3 x_1 \Delta x_2$$
$$\Delta y = (\beta_2 + \beta_3 x_1) \Delta x_2$$

slope = $f(x_1)$

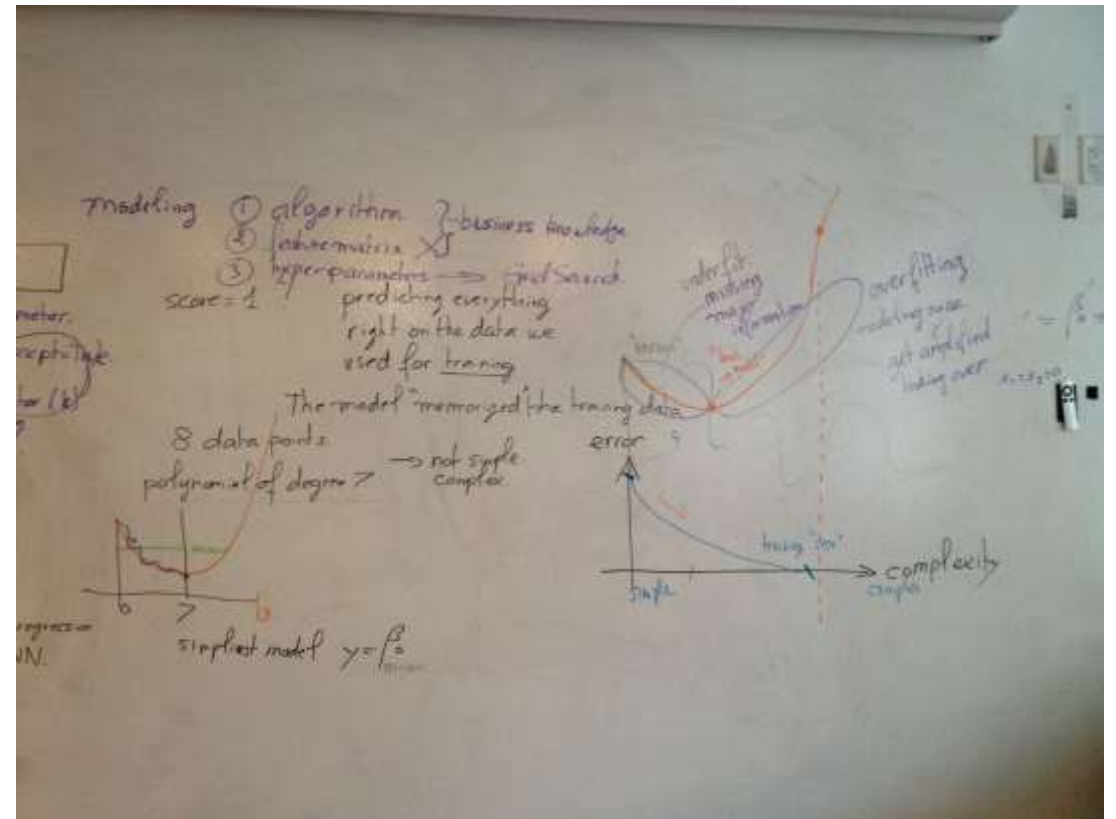
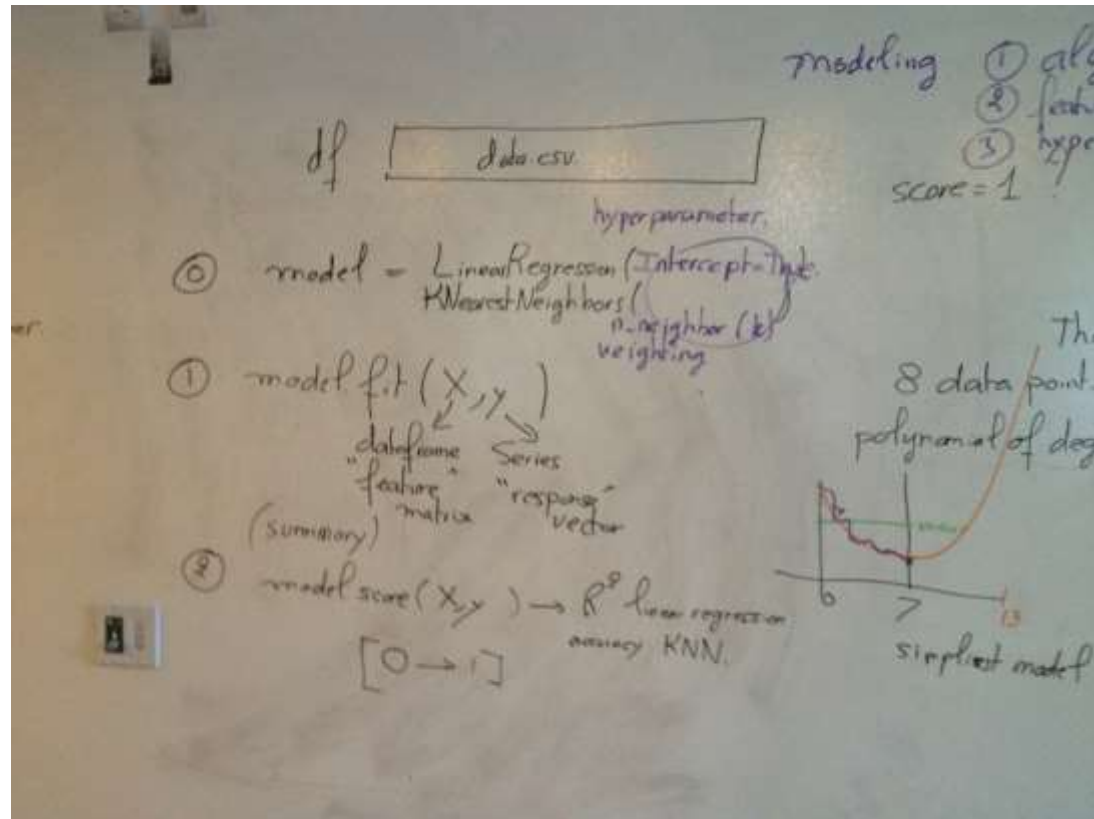
Δy

1000 = x_1 TV
1000-1000 = x_2 Newspaper
0.0100 = y Sales
 $\Delta y = (f(x_1) - 1000)100$

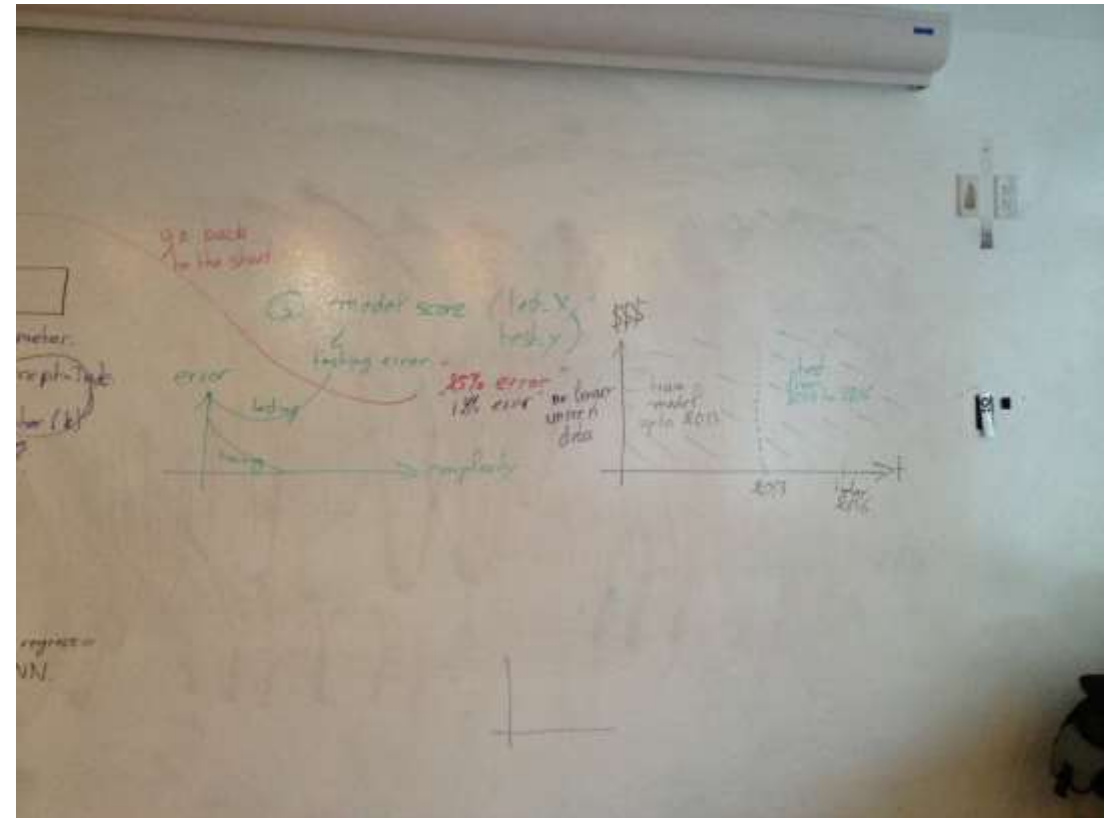
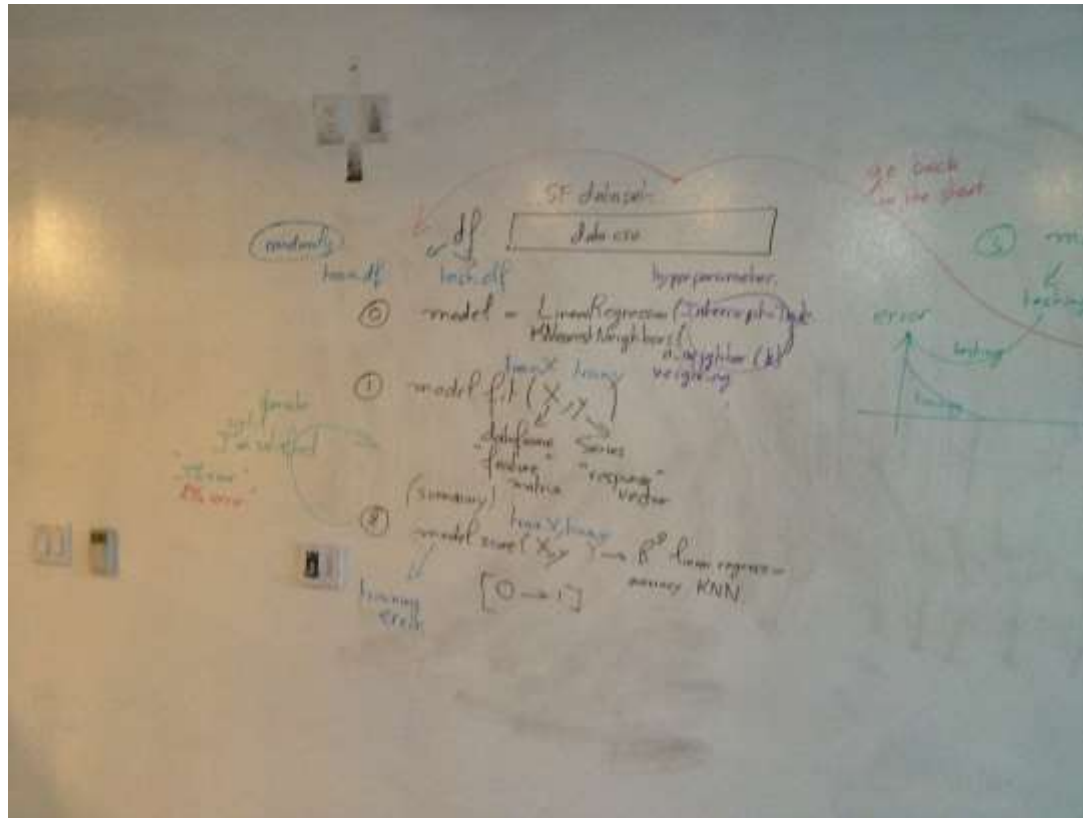
Linear Regression | Binary Variables



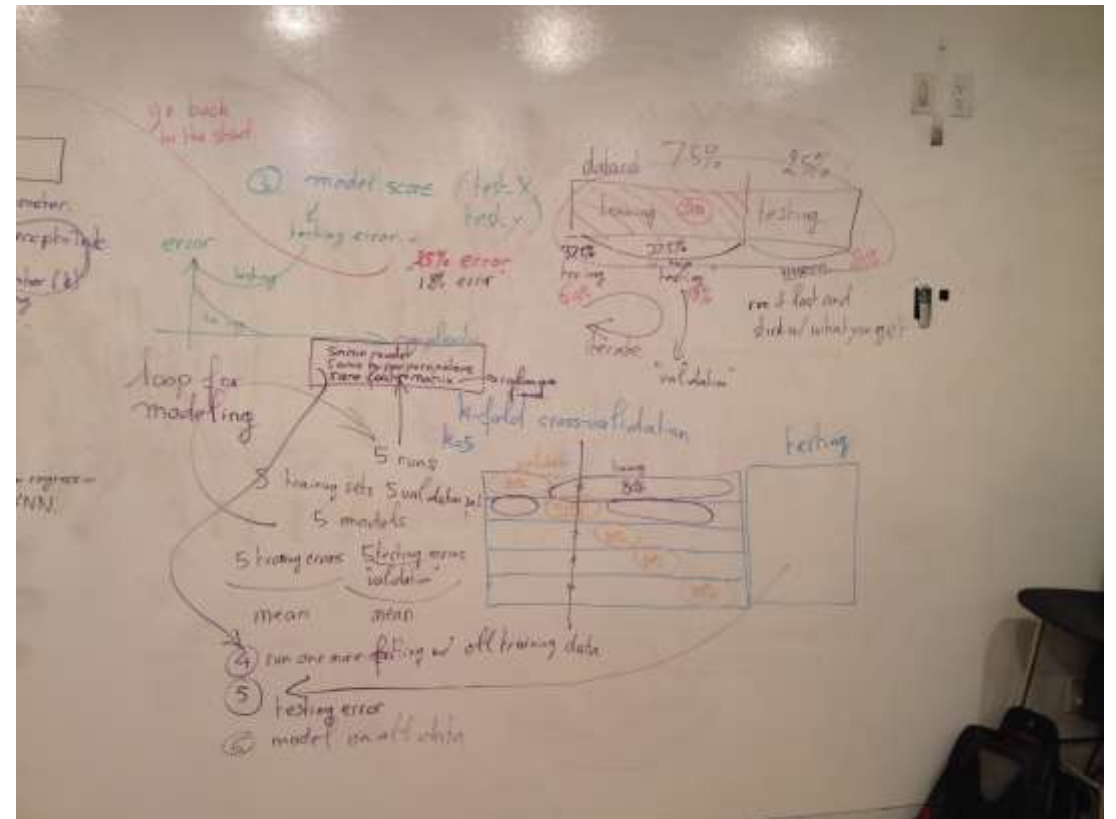
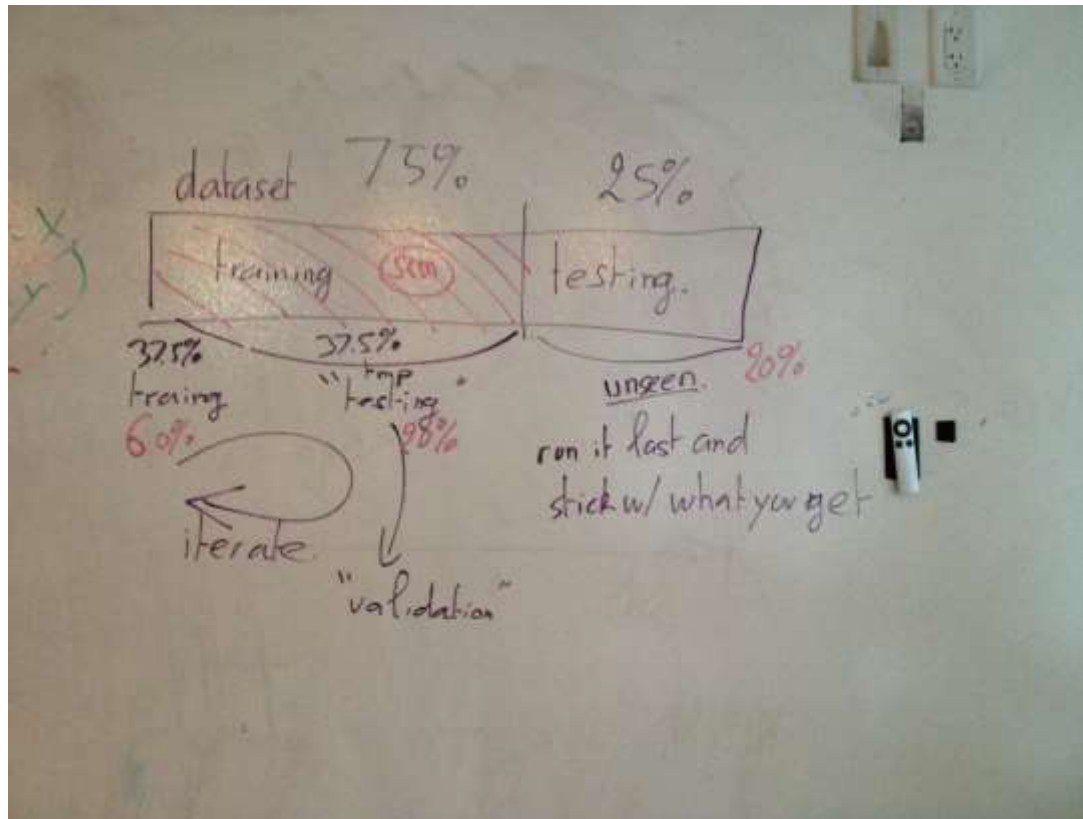
Modeling | You need a testing set ("unseen" data) to prevent overfitting



Modeling, Take 2 | Once you use your testing set, you can't go back and remodel. If you do, you are incorporating a learned knowledge from your "unseen" data which is no longer unseen...



Modeling, Take 3 | Use train and validation sets for training ("seen" data) and keep your test set as "unseen" data. Rotate "seen" data with k-fold cross-validation. Finally, use all your available data for the final model





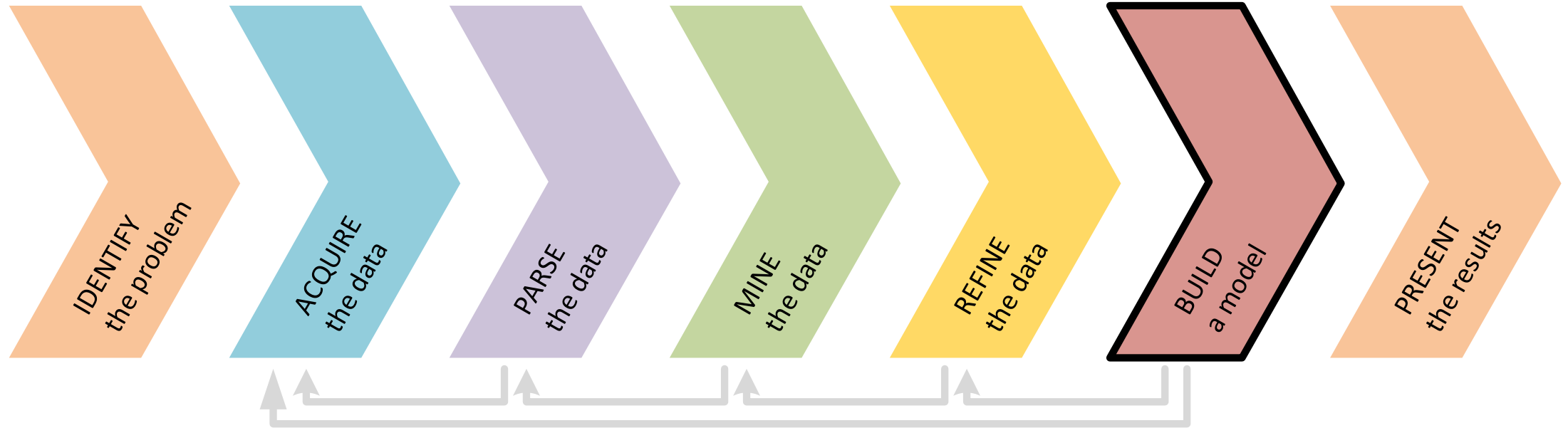
DS

Today

Today, we are focusing on logistic regression

Research Design and Data Analysis	Research Design	Data Visualization in <i>pandas</i>	Statistics	Exploratory Data Analysis in <i>pandas</i>
Foundations of Modeling	Linear Regression	Classification Models	Evaluating Model Fit	Presenting Insights from Data Models
Data Science in the Real World	Decision Trees and Random Forests	Time Series Data	Natural Language Processing	Databases

Today, we keep our focus on the **BUILD** a model step but with a focus on logistic regression



Here's what's happening today:

- Announcements and Exit Tickets
- Review
- ⑥ Build a Model | Logistic Regression
 - How logistic regression relates to linear regression
 - “Retrofitting” linear regression into logistic regression
 - Interpreting the logistic regression coefficients
- Iris dataset and Codealong on the Iris dataset
- Lab – Introduction to Logistic Regression
- Review
- **Unit Project 3 (due next session on 6/7)**



DS

Pre-Work

Pre-Work

Before this lesson, you should already be able to:

- Implement a linear model (`LinearRegression`) with *sklearn*
- Define the concept of coefficients
- Recall metrics for accuracy and misclassification

DS

Logistic Regression

Why is logistic regression so valuable to know?

- It addresses many commercially valuable classification problems, such as:
 - Fraud detection (e.g., payments, e-commerce)
 - Churn prediction (marketing)
 - Medical diagnoses (e.g., is the test positive or negative?)
 - and many, many others...

A black circle containing the white text "DS".

DS

Logistic Regression

How logistic regression relates to linear regression

Logistic regression is a generalization of the linear regression model to classification problems

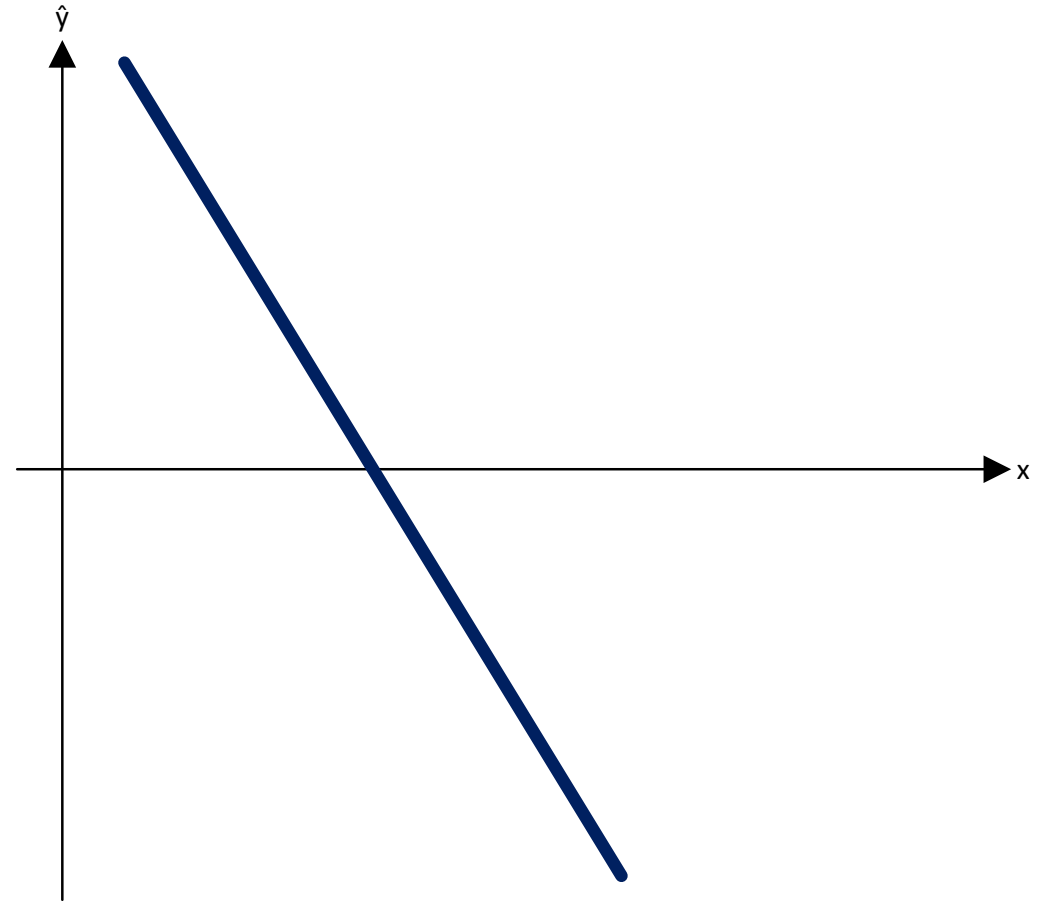
- The name is somewhat misleading
 - “Regression” comes from fact that we fit a linear model to the feature space
 - But it is really a technique for classification, not regression
- We use a linear model, similar to linear regression, in order to solve if an item *belongs* or *does not* belong to a class model
 - It is a binary classification technique: $y = \{0, 1\}$
 - Our goal is to classify correctly two types of examples:
 - Class 0, labeled as 0, e.g., “*belongs*”
 - Class 1, labeled as 1, e.g., “*does not belong*”

With linear regression, \hat{y} is in $]-\infty; +\infty[$, not $[0; 1]$. How do we fix this for logistic regression?

- ▶ The key variable in any regression problem is the outcome variable \hat{y} given the covariate x

$$\hat{y} = \hat{\beta}x$$

- ▶ With linear regression, \hat{y} takes values in $]-\infty; +\infty[$
- ▶ However, with logistic regression, \hat{y} takes values in the unit interval $[0; 1]$



DS

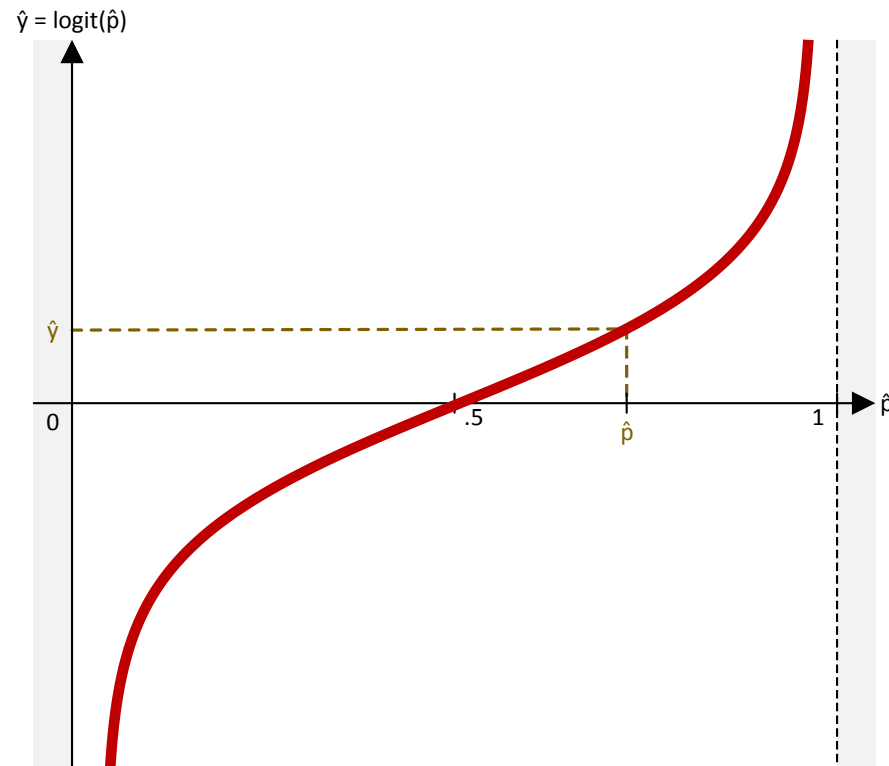
Logistic Regression

“Retrofitting” linear regression into logistic regression

We “retrofit” linear regression in logistic regression with a transformation called the *logit* function (a.k.a., the *log-odds* function) and its inverse, the *logistic* function (a.k.a., *sigmoid* function)

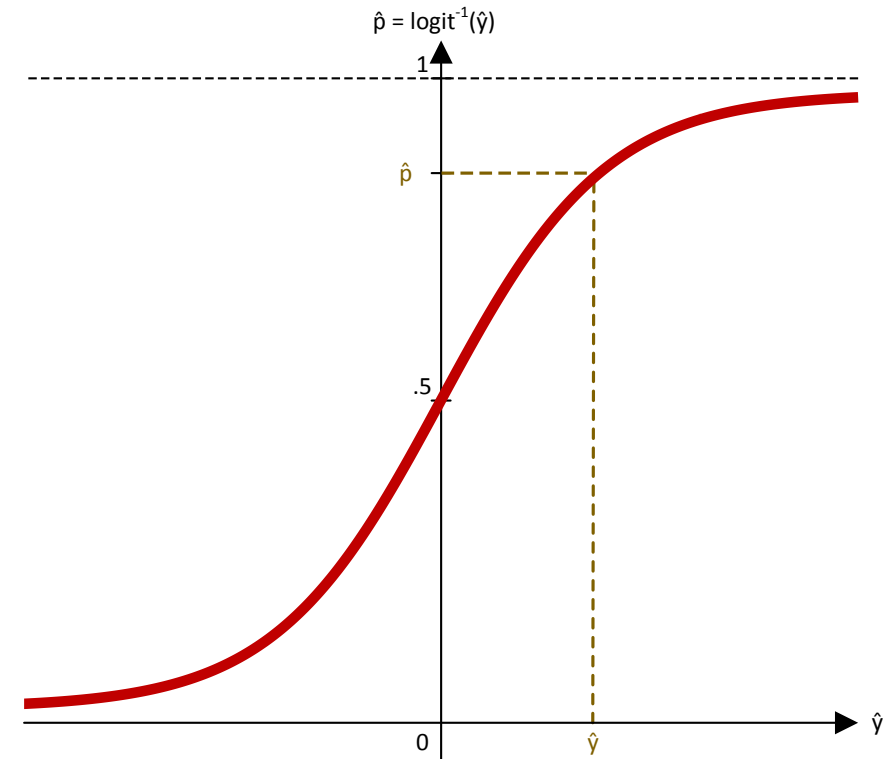
logit maps \hat{p} $([0; 1])$ to \hat{y} $(]-\infty; +\infty[)$

$$\text{logit}(\hat{p}) = \ln\left(\frac{\hat{p}}{1-\hat{p}}\right) = \hat{y}$$



$\pi = \text{logit}^{-1}$ maps \hat{y} $(]-\infty; +\infty[)$ to \hat{p} $([0; 1])$

$$\pi(\hat{y}) = \frac{e^{\hat{y}}}{e^{\hat{y}} + 1} = \hat{p}$$

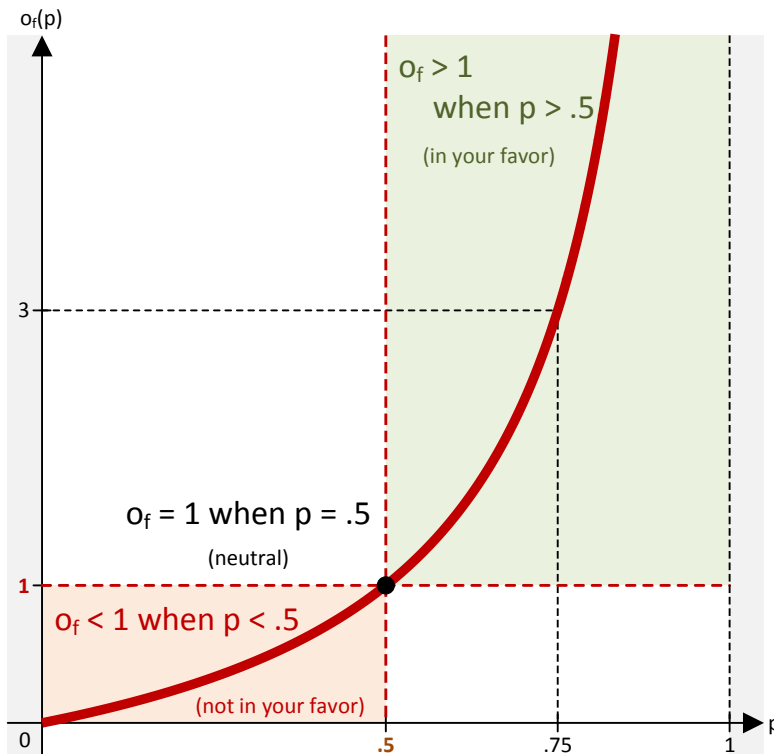


Why is the *logit* function also called the *log-odds* function?

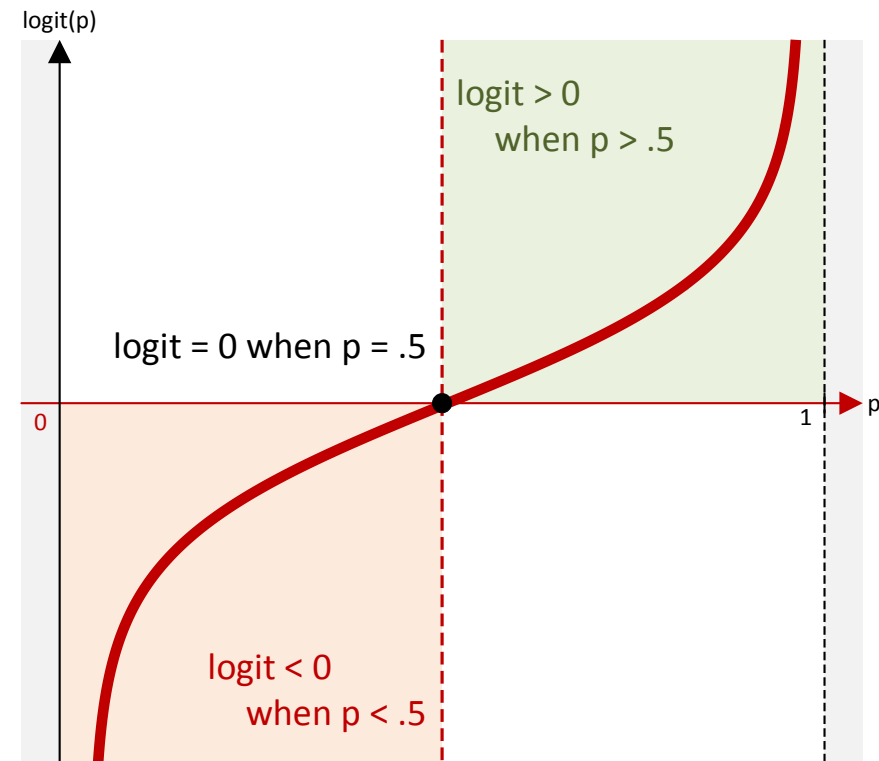
$$o_f = \frac{\text{probability that the event (with probability } p) \text{ happens}}{\text{probability that the event does not happen}}$$

\hat{p}
 $1 - \hat{p}$

odds (in favor)



$$\text{logit}(p) = \ln(o_f) = \ln\left(\frac{p}{1-p}\right)$$

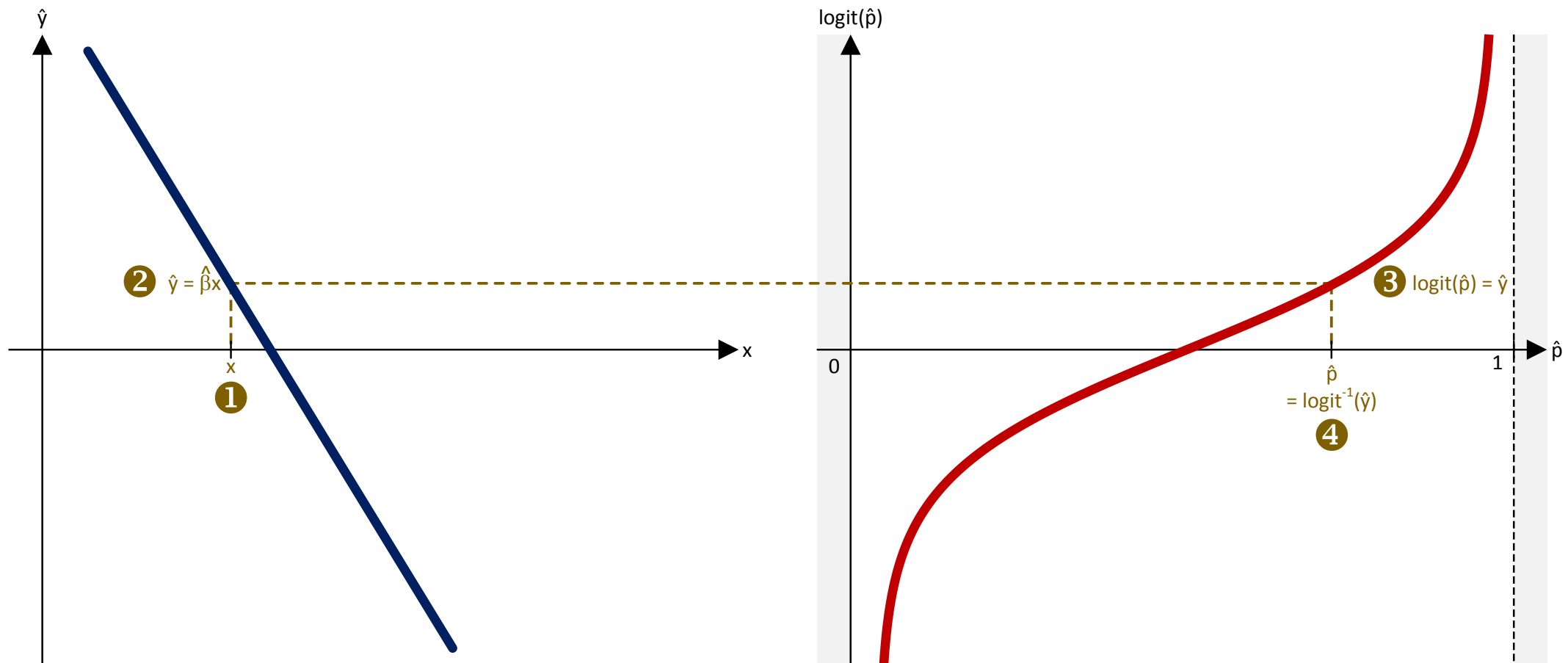


Logistic Regression

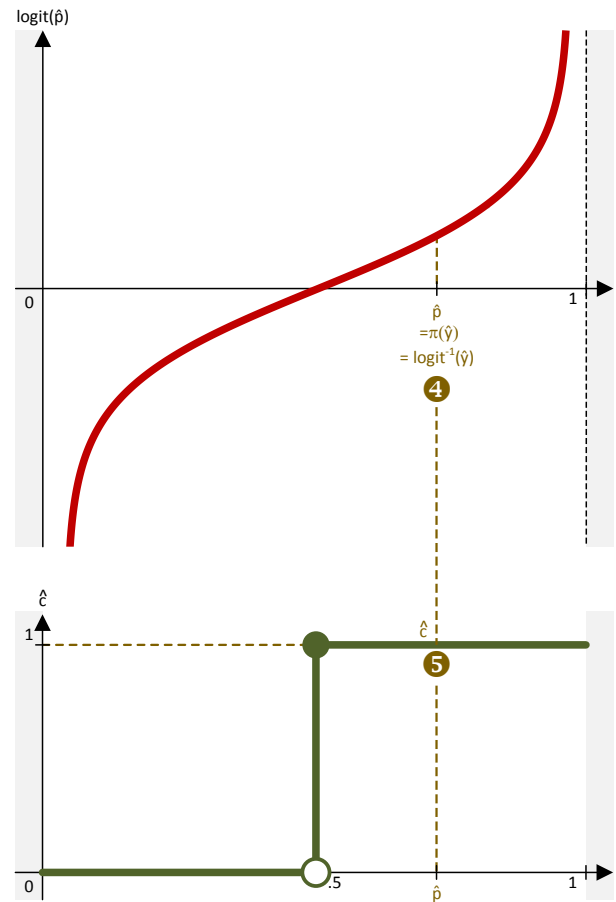
- ▶ Putting together $\hat{y} = \hat{\beta}x$ and $\hat{p} = \pi(\hat{y})$ (really, mapping \hat{y} back to \hat{p}), we get

$$\hat{p} = \pi(\hat{\beta}x) = \frac{e^{\hat{\beta}x}}{e^{\hat{\beta}x} + 1} = \frac{1}{1 + e^{-\hat{\beta}x}}$$

$$\hat{p} = \text{logit}^{-1}(\hat{y}) = \text{logit}^{-1}(\hat{\beta}x) = \frac{1}{1 + e^{-\hat{\beta}x}}$$



Finally, probabilities are “snapped” to class labels (e.g., by thresholding at the 50% level)



A black circle containing the white text "DS".

DS

Logistic Regression

Interpreting the logistic regression coefficients

Interpreting the logistic regression coefficients

- With linear regressions, $\widehat{\beta}_j$ represents the change in y for a change in unit of x_j

$$\ln\left(\frac{\hat{p}}{1-\hat{p}}\right) = \hat{\beta}x = \widehat{\beta}_0 + \widehat{\beta}_1 \cdot x_1 + \cdots + \widehat{\beta}_k \cdot x_k$$

- With logistic regressions, $\widehat{\beta}_j$ represents the **log-odds** change in y for a change in unit of x_j
- This also means that $e^{\widehat{\beta}_j}$ represents the multiplier change in **odds** in y for a change in unit of x_j

$$\frac{\widehat{odds}(x_j + 1)}{\widehat{odds}(x_j)} = \frac{e^{\hat{y}(x_{j+1})}}{e^{\hat{y}(x_j)}} = e^{\hat{y}(x_{j+1}) - \hat{y}(x_j)} = e^{(\boxed{\times} + \widehat{\beta}_j \cdot x_j + \otimes) - (\boxed{\times} + \widehat{\beta}_j \cdot (x_j + 1) + \otimes)} = e^{\widehat{\beta}_j}$$

A black circle containing the white text 'DS' in a bold, sans-serif font.

Logistic Regression

Activity / Interpreting the logistic regression coefficients

Activity | Interpreting the logistic regression coefficients



EXERCISE

DIRECTIONS (5 minutes)

1. Suppose we are interested in mobile purchasing behavior. Let y be a class label denoting purchase/no purchase, and x a feature denoting whether a phone is an iPhone or not. After performing a logistic regression, we get $\beta_1 = .693$. What does this mean?
2. When finished, share your answers with your table

DELIVERABLE

Answers to the above question

Activity | Interpreting the logistic regression coefficients (cont.)



EXERCISE

1. In this case, the odds ratio change is $e^{\beta_1} = e^{.693} = 2$, meaning the likelihood of purchase is twice as high if the phone is an iPhone

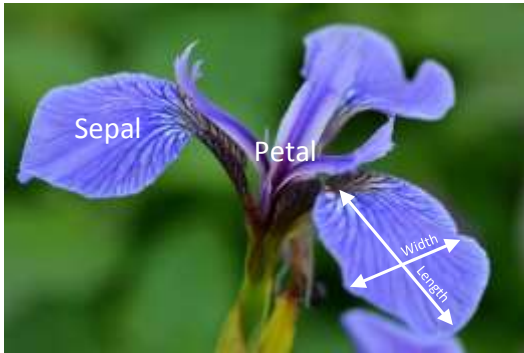
A black circle containing the white text "DS".

DS

Iris Dataset, Take 2

Review | Iris dataset

Iris Setosa



Iris Versicolor



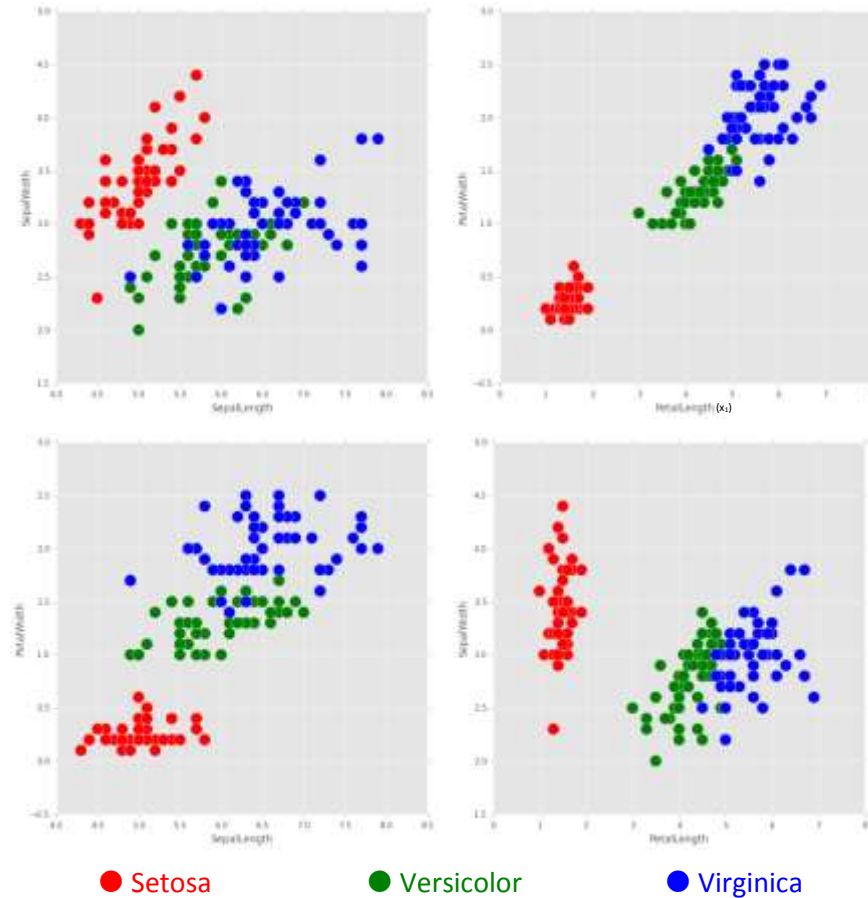
Iris Virginica



Source: Flickr

- 3 classes of Irises (*Setosa*, *Versicolor*, and *Virginica*)
- 4 features
 - Sepal length and width
 - Petal length and width
- 50 instances of each class

Review | Iris dataset (cont.)



- Features (4)

- “SepalLength”, “SepalWidth”, “PetalLength”, and “PetalWidth”

- Classes (3)

- “Setosa” as c_{red} , “Versicolor” as c_{green} , and “Virginica” as c_{blue}

A black circle containing the white text "DS".

DS

Iris Dataset, Take 2

Codealong – Logistic Regression

A black circle containing the white text "DS".

DS

Review

Review

You should now be able to:

- Build a Logistic regression classification model using *sklearn*
- Describe the logit and sigmoid functions, odds, and odds ratios as well as how they relate to logistic regression
- Evaluate a model using metrics such as classification accuracy/error

A black circle containing the white text "DS".

DS

Q & A



DS

Before Next Class

Before Next Class

Before the next lesson, you should already be able to:

- Create and interpret results from a binary classification problem
- Know what a decision line is in logistic regression

Next Class

Advanced Metrics and Communicating Results

Learning Objectives

After the next lesson, you should be able to:

- Evaluate a model using advanced metrics such as confusion matrix and ROC/AUC curves
- Explain the trade-offs between the precision and recall of a model while articulating the cost of false positives vs. false negatives
- Describe the difference between visualization for presentations vs. exploratory data analysis
- Identify the components of a concise and convincing report and how they relate to specific audiences/stakeholders



DS

Exit Ticket

Don't forget to fill out your exit ticket [here](#)

Slides © 2016 Ivan Corneillet Where Applicable
Do Not Reproduce Without Permission