# Introduction to Machine Learning for Social Scientists

## Class 6: Classification

Edgar Franco Vivanco

Stanford University
Department of Political Science

*edgarf1@stanford.edu*

Summer 2018

# Where are you struggling?

Mini survey results:

- ▶ Functions
- ▶ Subsetting (using [])
- ▶ Difference between linear regression and logistic regression
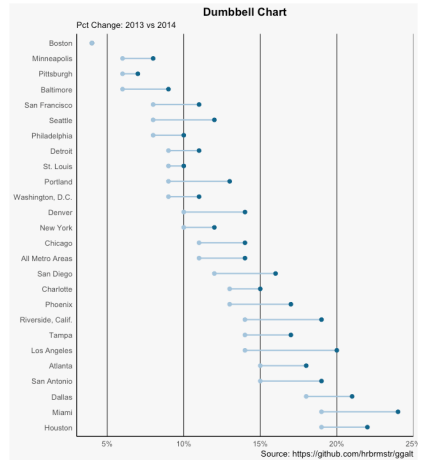
# Where are you struggling?

Mini survey results:

- ► Functions
- ► Subsetting (using [])
- ► Difference between linear regression and logistic regression

**Tutorials available before midterm**

# Extra workshops

- ggplot!!
- data manipulation
- text analysis

# Other petitions

Mini survey results:

- ▶ Connection with Machine Learning:

# Other petitions

Mini survey results:

- ▶ Connection with Machine Learning:
- ▶ Next class will study an application of these methods
- ▶ Other fields:

# Other petitions

Mini survey results:

- ▶ Connection with Machine Learning:
- ▶ Next class will study an application of these methods
- ▶ Other fields:
- ▶ Fake news, Psychology, Sociology, etc.
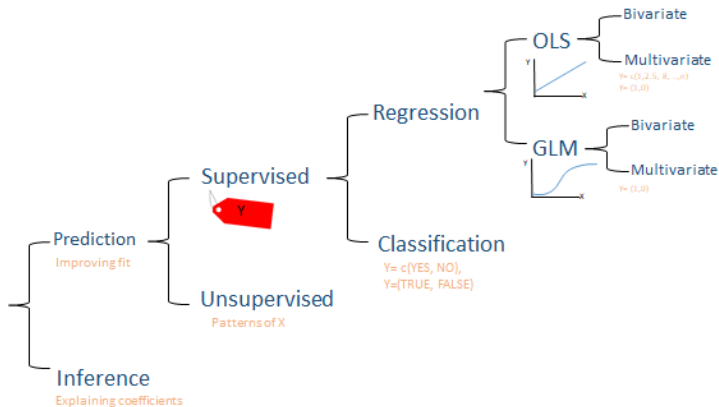
# Today's Goals

1. Key concepts:
   - ▶ Linear Probability Model vs. Generalized Linear Model
   - ▶ Classification
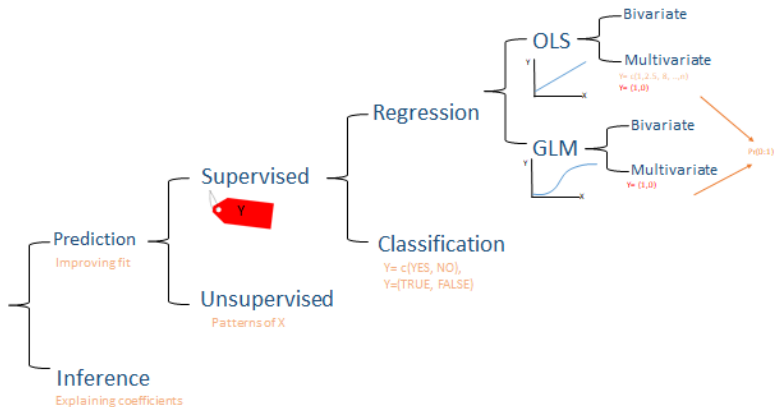   - ▶ Confusion Matrix
   - ▶ Performance measures
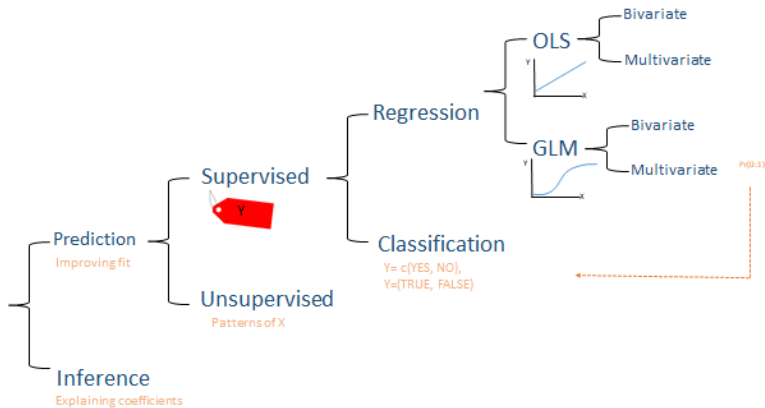2. Key techniques and R functions:
   - ▶ ifelse
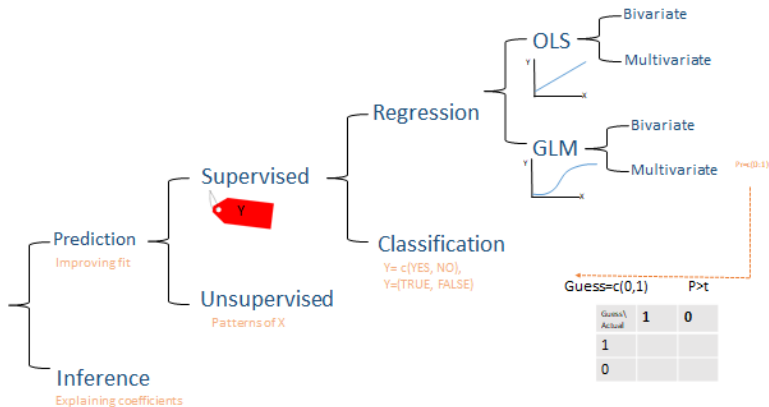   - ▶ table

# Our Mental Map: OLS and GLM

# Our Mental Map: Predicting probabilities
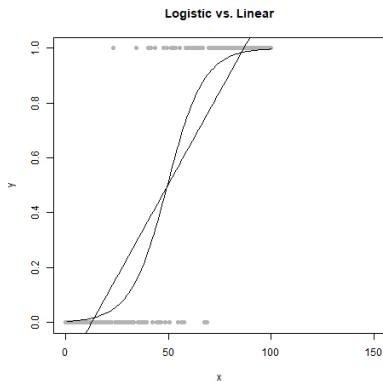
# Our Mental Map: Classify

# Our Mental Map: Test our model

# Overview

Logistics

From prediction to classification

Performance Measures

# From Prediction to classification



Logistic vs. Linear

▶ If we have a qualitative
  outcome (Y=0 or Y=1) we can
  predict probabilities using a
  linear or a logistic model.

# From Prediction to classification



- If we have a qualitative outcome (Y=0 or Y=1) we can predict probabilities using a linear or a logistic model.

- In our example:
  - Y: Vote for Iraq War (YES=1, NO=0)
  - rep: Senator is Republican
  - gorevote: Percentage of vote for Al Gore in Senator's state

# From Prediction to classification

```
> fit <- lm(y ~ rep + gorevote, data = iraqvote)
> summary(fit)

Call:
lm(formula = y ~ rep + gorevote, data = iraqvote)

Residuals:
    Min      1Q  Median      3Q     Max
-0.7654 -0.1533  0.0509  0.2904  0.5707

Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept)  1.174458   0.236256   4.971 2.87e-06 ***
repTRUE      0.316933   0.080493   3.937 0.000155 ***
gorevote    -0.012376   0.004715  -2.625 0.010072 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.3603 on 97 degrees of freedom
Multiple R-squared:  0.2888,    Adjusted R-squared:  0.2742
F-statistic:  19.7 on 2 and 97 DF,  p-value: 6.617e-08
```

▶ We can run a linear model

▶ $p(Y = 1|X) = \beta_0 + \beta_1 rep + \beta_2 gorevote$

# From Prediction to classification

```
> fit <- lm(y ~ rep + gorevote, data = iraqVote)
> summary(fit)

Call:
lm(formula = y ~ rep + gorevote, data = iraqVote)

Residuals:
    Min      1Q  Median      3Q     Max
-0.7654 -0.1533  0.0509  0.2904  0.5707

Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept)  1.174458   0.236256   4.971 2.87e-06 ***
repTRUE      0.316933   0.080493   3.937 0.000155 ***
gorevote    -0.012376   0.004715  -2.625 0.010072 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.3603 on 97 degrees of freedom
Multiple R-squared:  0.2888,     Adjusted R-squared:  0.2742
F-statistic:  19.7 on 2 and 97 DF,  p-value: 6.617e-08
```

► We can run a linear model

► $p(Y = 1|X) = \beta_0 + \beta_1 rep + \beta_2 gorevote$

► And calculate predictions:

► $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 rep + \hat{\beta}_2 gorevote$

# From Prediction to classification

```
> fit <- lm(y ~ rep + gorevote, data = iraqvote)
> summary(fit)

Call:
lm(formula = y ~ rep + gorevote, data = iraqvote)

Residuals:
    Min      1Q  Median      3Q     Max
-0.7654 -0.1533  0.0509  0.2904  0.5707

Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept)  1.174458   0.236256   4.971 2.87e-06 ***
repTRUE      0.316933   0.080493   3.937 0.000155 ***
gorevote    -0.012376   0.004715  -2.625 0.010072 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.3603 on 97 degrees of freedom
Multiple R-squared:  0.2888,    Adjusted R-squared:  0.2742
F-statistic:  19.7 on 2 and 97 DF,  p-value: 6.617e-08
```
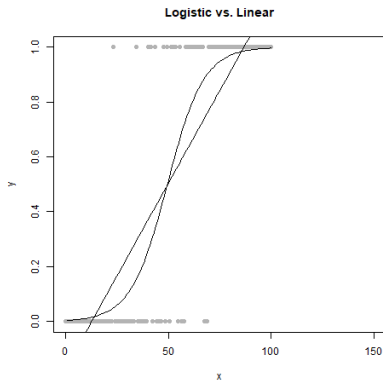
- We can run a linear model

- $p(Y = 1|X) = \beta_0 + \beta_1 rep + \beta_2 gorevote$

- And calculate predictions:

- $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 rep + \hat{\beta}_2 gorevote$

- $\hat{Y} = 1.144 + 0.3169 rep - 0.0123 gorevote$
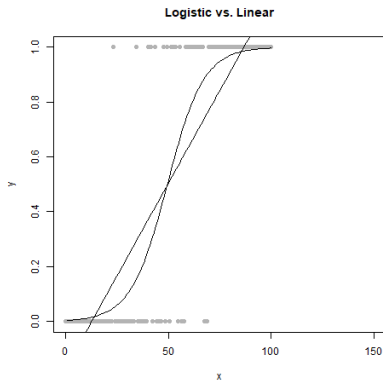
# From Prediction to classification



- We can run a linear model

- $p(Y = 1|X) = \beta_0 + \beta_1 rep + \beta_2 gorevote$

- And calculate predictions:

- $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 rep + \hat{\beta}_2 gorevote$

- $\hat{Y} = 1.1744 + 0.3169 rep - 0.0123 gorevote$

- $0.9766 = 1.1744 + 0.3169 - 0.0123 * 41.59$

# From Prediction to classification



Logistic vs. Linear

▶ A logistic model will produce predictions between 0 and 1.

# From Prediction to classification



Logistic vs. Linear

▶ A logistic model will produce predictions between 0 and 1.

# From Prediction to classification

```
> rep_reg_glm <- glm(y~rep+gorevote, family = binomial, data = iraqvote)
> summary(rep_reg_glm)

Call:
glm(formula = y ~ rep + gorevote, family = binomial, data = iraqvote)

Deviance Residuals:
     Min       1Q   Median       3Q      Max
 -2.12054  0.07761  0.19676  0.59926  1.59277

Coefficients:
             Estimate Std. Error z value Pr(>|z|)
(Intercept)   5.87859    2.27506   2.584  0.00977 **
repTRUE       3.01881    1.07138   2.818  0.00484 **
gorevote     -0.11322    0.04508  -2.512  0.01201 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 107.855  on 99  degrees of freedom
Residual deviance:  71.884  on 97  degrees of freedom
AIC: 77.884

Number of Fisher Scoring iterations: 6
```

▶ A logistic model will produce predictions between 0 and 1.

▶ Because it models a relationship:

$$p(X) = \frac{1}{1 + exp^{-\beta X}}$$

# From Prediction to classification

```
> rep_reg_glm <- glm(y~rep+gorevote, family = binomial, data = iraqvote)
> summary(rep_reg_glm)

Call:
glm(formula = y ~ rep + gorevote, family = binomial, data = iraqvote)

Deviance Residuals:
     Min       1Q   Median       3Q      Max
-2.12054  0.07761  0.19676  0.59926  1.59277

Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept)  5.87859    2.27506   2.584  0.00977 **
repTRUE      3.01881    1.07138   2.818  0.00484 **
gorevote    -0.11322    0.04508  -2.512  0.01201 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 107.855  on 99  degrees of freedom
Residual deviance:  71.884  on 97  degrees of freedom
AIC: 77.884

Number of Fisher Scoring iterations: 6
```

- A logistic model will produce predictions between 0 and 1.

- Because it models a relationship:

$$p(X) = \frac{1}{1 + exp^{-\beta X}}$$

  Optimized via Maximum Likelihood

- $4.18 =$
  5.88+3.021-0.113*41.59

# From Prediction to classification

```
> rep_reg_glm <- glm(y~rep+gorevote, family = binomial, data = iraqvote)
> summary(rep_reg_glm)

Call:
glm(formula = y ~ rep + gorevote, family = binomial, data = iraqvote)

Deviance Residuals:
    Min       1Q    Median       3Q      Max
-2.12054  0.07761  0.19676  0.59926  1.59277

Coefficients:
             Estimate Std. Error z value Pr(>|z|)
(Intercept)  5.87859    2.27506   2.584  0.00977 **
repTRUE      3.01881    1.07138   2.818  0.00484 **
gorevote    -0.11322    0.04508  -2.512  0.01201 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 107.855  on 99  degrees of freedom
Residual deviance:  71.884  on 97  degrees of freedom
AIC: 77.884

Number of Fisher Scoring iterations: 6
```

▶ A logistic model will produce predictions between 0 and 1.

▶ Because it models a relationship:

$$p(X) = \frac{1}{1 + exp^{-\beta X}}$$

Optimized via Maximum Likelihood

▶ 4.18 = 5.88+3.021-0.113*41.59

▶ 0.985 $= \dfrac{1}{1 + exp^{-4.18}}$

Call $p_i = \Pr(\text{Vote}_i = 1 | \boldsymbol{x}_i)$

$$
\begin{aligned}
\text{Vote}_i &\sim \text{Bernoulli}(p_i) \\
p_i &= f(\boldsymbol{\beta} \cdot \boldsymbol{x}_i) \\
\log\left(\frac{p_i}{1 - p_i}\right) &= \boldsymbol{\beta} \cdot \boldsymbol{x}_i \\
p_i &= \frac{\exp(\boldsymbol{\beta} \cdot \boldsymbol{x}_i)}{1 + \exp(\boldsymbol{\beta} \cdot \boldsymbol{x}_i)} \\
&= \frac{1}{1 + \exp(-\boldsymbol{\beta} \cdot \boldsymbol{x}_i)}
\end{aligned}
$$

Important functions:

$$
\begin{aligned}
\text{odds}(p) &= \frac{p}{1 - p} \\
\log \text{ odds or logit}(p) &= \log\left(\frac{p}{1 - p}\right) \\
\text{logistic function or logit}^{-1}(a) &= \frac{1}{1 + \exp(-a)}
\end{aligned}
$$

# From Prediction to classification



- A logistic model will produce predictions between 0 and 1.

- Because it models a relationship:

$$p(X) = \frac{1}{1 + exp^{-\beta X}}$$

- 4.18 = 5.88+3.021-0.113*41.59

- 0.985 = $\dfrac{1}{1 + exp^{-4.18}}$

# How to create classifications?



Logistic vs. Linear

▶ We can choose a threshold such as:

$$Pr(Y \hat{=} 1 | X) >= t$$

Then clas=1, and 0 otherwise

# How to create classifications?



- ▶ We can choose a threshold such as:

$$Pr(Y \stackrel{\hat{}}{=} 1|X) >= t$$

  Then clas=1, and 0 otherwise

- ▶ We can do this by using the function 'ifelse()'

# How to create classifications?



- ▶ We can choose a threshold such as:

$$Pr(Y \hat{=} 1 | X) >= t$$

  Then clas=1, and 0 otherwise

- ▶ We can do this by using the function 'ifelse()'

- ▶ And now we can start comparing our models with the observed values

# Errors

# Confusion Matrix

To asses the quality of our data we compare our classifications with the real data or the "gold standard".

| Actual \\ Guess | Yes | No |
|:---:|:---:|:---:|
| Yes | | |
| No | | |

# Confusion Matrix

| Actual \ Guess | Yes | No |
|---|---|---|
| Yes | True positive | False Negative |
| No | False Positive | True Negative |

# Confusion Matrix:

Code approach:
'ifelse()' function: ifelse(condition, yes, no)

```
# Actual yes and guess yes
tp <- ifelse (y ==1 & predicted==1,1,0)
# Actual no and guess n0
tn <- ifelse(y ==0 & predicted==0,1,0)
# Actual no and guess yes
fp <- ifelse(y ==0 & predicted==1,1,0)
# Actual yes and guess no
fn <- ifelse(y ==1 & predicted==0,1,0)
```

# Accuracy

Accuracy is the percentage of observations classified correctly.

$$Accuracy = \frac{TruePositive + TrueNegative}{TruePositive + TrueNegative + FalseNegative + FalsePositive}$$

# Precision

How many items classified as Yes are correctly classified?

$$Precision = \frac{TruePositive}{TruePositive + FalsePositive}$$

It is equal to 1 if all the guesses as Yes are actually Yes.

# Recall

How many items **that are actually** as Yes are correctly classified?
In other words, is the number of correct results divided by the
number of results that should have been returned.

$$Recall = \frac{TruePositive}{TruePositive + FalseNegative}$$

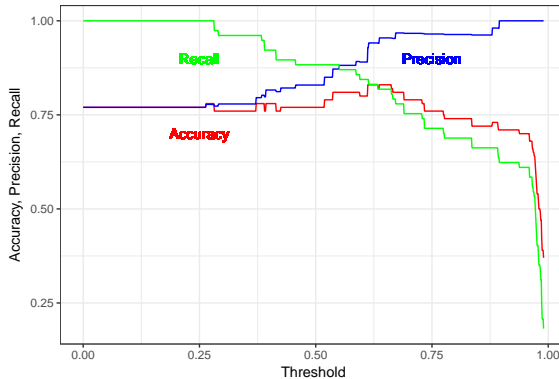It is equal to 1 if all the actual Yes are classified as Yes.

# F-score

Harmonic mean of precision and recall:

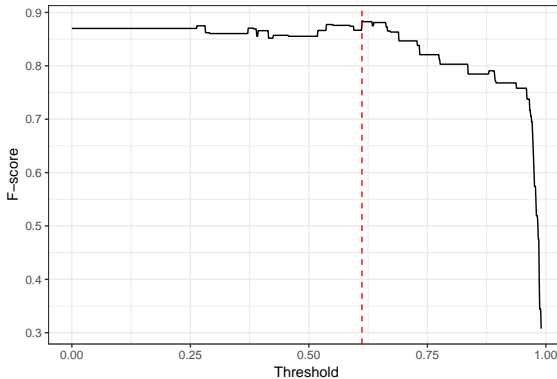$$F = \frac{2 * Precision * Recall}{Precision + Recall}$$

# Performance

All these measures are function of the threshold.

# F-score

We can find the threshold that optimizes the F-score.

# Examples

- **Fraud in bank transactions**: High recall, ie. most of the fraudulent transactions are identified, probably at loss of precision.

- **Twitter:** If we are interested in finding out when a tweet expresses a negative sentiment, we can probably raise precision (to gain certainty).

- **Terrorist attacks:** Given the 800 million average passengers on US flights per year and the 19 (confirmed) terrorists who boarded US flights from 20002017, a very accurate model will predict everyone as non terrorist. Instead, we should focus on recall.

# R!

# Confusion Matrix: LM

| Actual \ Guess | Yes | No |
|:---:|:---:|:---:|
| Yes | 69 | 8 |
| No | 15 | 8 |

# Confusion Matrix: GLM

| Actual \ Guess | Yes | No |
|:---:|:---:|:---:|
| Yes | 68 | 9 |
| No | 14 | 9 |

# Results

- **LM:**
    - Accuracy: 0.77
    - Precision: 0.8214
    - Recall: 0.8961
    - F-score: 0.8571
- **Logistic**
    - Accuracy: 0.77
    - Precision: 0.8293
    - Recall: 0.8831
    - F-score: 0.8554

# NEXT

- ▶ Resampling methods (Crossvalidation)
  - ▶ Training
  - ▶ Test
  - ▶ Validation
- ▶ Midterm guidelines
- ▶ Article

i. White (Seattle, Washington)

ii. Black (Seattle, Washington)

iii. Asian (Seattle, Washington)

iv. Less than High school (Milwaukee, Wisconsin)

v. Graduate school (Milwaukee, Wisconsin)

vi. Income (Tampa, Florida)