

Logistic Regression

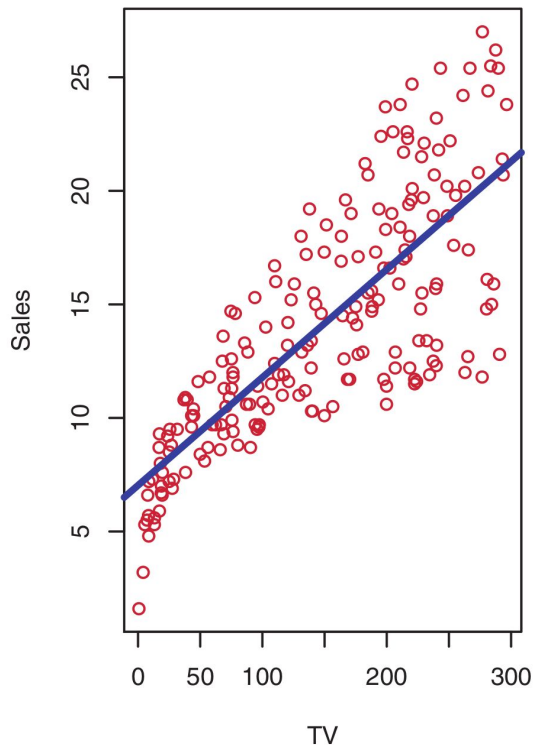
Objectives

After this lecture you should be able to:

- Explain the difference between regression and classification
- Decide which class to make positive in binary classification
- Describe how classifications are made using predicted probabilities and a threshold
- Describe how logistic regression determines probabilities
- Interpret logistic regression coefficients
- Calculate (and define) TP, FP, FN, and TN
- Construct a confusion matrix
- Define metrics useful in classification, and give examples where each might be useful
- Describe what a ROC curve is, how it gets constructed, and how it describes the overall performance of a classifier

Regression vs Classification

Sales as a $f(\text{TV advertising})$



The target, **Sales**, is numerical.
This makes it a **regression** problem.

The features could be numerical or categorical (of different classes), but *what determines if it's a regression or classification problem is the target.*

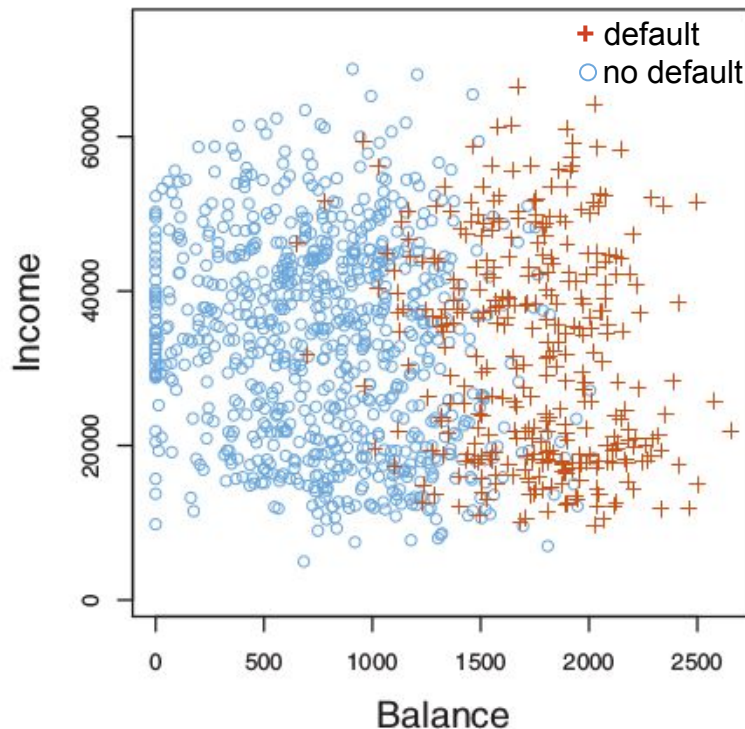
Regression vs Classification

In classification the target is qualitative, not quantitative.
(e.g., no default, default)

In this case there are two categories, or classes. This makes it a **binary classification** problem.

If there are more than two classes, it's a **multi-class** or **multinomial** classification problem.

Credit card balance default
as $f(\text{Income}, \text{Balance})$



Regression vs Classification

You'd like to predict the winner of the Kentucky Derby.



How would you approach it? Regression? Classification?

Regression vs Classification

Could calculate either of these numerical values for each horse and sort:

- Finish time
- Lengths behind winner



Regression vs Classification

Could calculate either of these numerical values for each horse and sort:

- Finish time
- Lengths behind winner

Binary classification

- Win or lose

Multinomial classification

- First, Second, Third, Fourth, etc.



Regression vs Classification

Could calculate either of these numerical values for each horse and sort:

- Finish time
- Lengths behind winner

Binary classification

- Win or lose

Multinomial classification

- First, Second, Third, Fourth, etc.

But what about:

- 1, 2, 3, 4 etc. regression? (does -1, or 3.5 make sense?)
- 1, 2, 3, 4 etc. classification? (1 better than 2, 2 better than 3, etc.)

Example of [ordinal regression](#), in Python [mord](#) package.



In binary classification, picking the + class

- In binary classification, there are two classes, e.g.:
 - negative, positive
 - False, True
 - no default, default
 - no cancer, cancer
 - no churn, churn
- These options are encoded into an indicator variable, taking value 0 or 1
- Usually you care about one of them more than the other, and that's the one you make 1, or the positive class
 - default
 - cancer
 - churn
- This matters when you pick an evaluation metric for your model (later)

Deciding whether a data point is the + class

A sklearn classification model (like logistic regression) will predict which class each datapoint is in:

```
>>> model = LogisticRegression()  
>>> model.fit(X_train, y_train)  
>>> y_hat = model.predict(X_test)  
>>> y_hat  
array([[0],  
       [1],  
       [0],  
       ...,  
       [1]])
```

But - and this is important - **this is not actually that useful.**

Deciding whether a data point is the + class

The *probability* that each row of data belongs to the positive class - that's much more useful:

```
>>> model = LogisticRegression()
>>> model.fit(X_train, y_train)
>>> y_hat_probs = model.predict_proba(X_test)[: , 1]
>>> y_hat_probs
array([[0.33],
       [0.78],
       [0.05],
       ...,
       [0.94]])
```

Now, you can pick a probability threshold, say 0.4, that you can use to transform the probabilities into classifications:

```
>>> threshold = 0.4
>>> y_hat = (y_hat_probs >= threshold).astype(int)
```

Breakout 1

Use the predicted probability of the + class and the given threshold to fill in the predictions in each of the tables below. One is filled out for you.

threshold:	1.0
y_hat_prob	y_hat
0.67	
0.20	
0.98	
0.03	

threshold:	0.75
y_hat_prob	y_hat
0.67	
0.20	
0.98	
0.03	

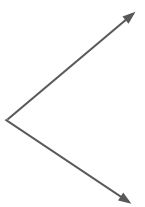
threshold:	0.5
y_hat_prob	y_hat
0.67	1
0.20	0
0.98	1
0.03	0

threshold:	0.25
y_hat_prob	y_hat
0.67	
0.20	
0.98	
0.03	

Determining probabilities in logistic regression

- Question: if logistic regression does classification, why is it called regression?
- Answer: It's a model, linear in its coefficients, that regresses values on to a numerical quantity that scales from -infinity to +infinity (and that's exactly what linear regression does).

both can
take value
from
 $-\infty$ to $+\infty$


$$\begin{aligned} ? &= \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots && \text{Linear regression} \\ ? &= \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots && \text{Logistic regression} \end{aligned}$$

Determining probabilities in logistic regression

- In Linear Regression, the sum of the products of the coefficients and features is the target, \hat{y}
- In Logistic regression, the sum of the products of the coefficients and features is the target, the *log odds*

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots \quad \text{Linear regression}$$

$$\ln(odds) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots \quad \text{Logistic regression}$$

Determining probabilities in logistic regression

The *log odds*, otherwise known as the [logit](#), links values that range from $[-\infty, \infty]$ to a probability in the range $[0, 1]$.

$[-\infty, \infty]$	$\ln(odds) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots$	logistic regression
---------------------	-----------------------------------------------------------	------------------------

$[0, \infty]$	$odds = \frac{p}{1 - p}$	$\frac{\text{prob. of + class}}{\text{prob. of - class}}$
---------------	--------------------------	-----------------------------------------------------------

$[0, 1]$	p	prob. of + class
----------	-----	------------------

Odds - some intuition

Statistical odds are the long run ratio of the probability of an event occurring to it not occurring.

For example, the odds of rolling a 3 on a fair, 6 sided die are 1 to 5.

$$\text{odds of } 3 = 1 \text{ to } 5 = \frac{1}{5} = \frac{p_3}{1 - p_3}$$

$$p_3 = \frac{1}{6}$$

Rearranging the logistic regression equation (1 of 4)

Solve for p

$$\ln(odds) = \beta_0 + \beta_1 x_1 + \dots$$

*logistic
regression*

$$odds = e^{\beta_0 + \beta_1 x_1 + \dots}$$

*exponentiate
both sides*

Rearranging the logistic regression equation (2 of 4)

Solve for p

$$odds = e^{\beta_0 + \beta_1 x_1 + \dots}$$

odds

$$\frac{p}{1 - p} = e^{\beta_0 + \beta_1 x_1 + \dots}$$

*substitute p
in odds*

$$p = e^{\beta_0 + \beta_1 x_1 + \dots} - e^{\beta_0 + \beta_1 x_1 + \dots} p$$

*multiply both
sides by $1 - p$*

Rearranging the logistic regression equation (3 of 4)

Solve for p

$$p = e^{\beta_0 + \beta_1 x_1 + \dots} - e^{\beta_0 + \beta_1 x_1 + \dots} p$$

*multiply both
sides by 1 - p*

$$p + e^{\beta_0 + \beta_1 x_1 + \dots} p = e^{\beta_0 + \beta_1 x_1 + \dots}$$

*get p on the
same side*

$$p(1 + e^{\beta_0 + \beta_1 x_1 + \dots}) = e^{\beta_0 + \beta_1 x_1 + \dots}$$

factor out p

Rearranging the logistic regression equation (4 of 4)

Solve for p

$$p(1 + e^{\beta_0 + \beta_1 x_1 + \dots}) = e^{\beta_0 + \beta_1 x_1 + \dots}$$

factor out p

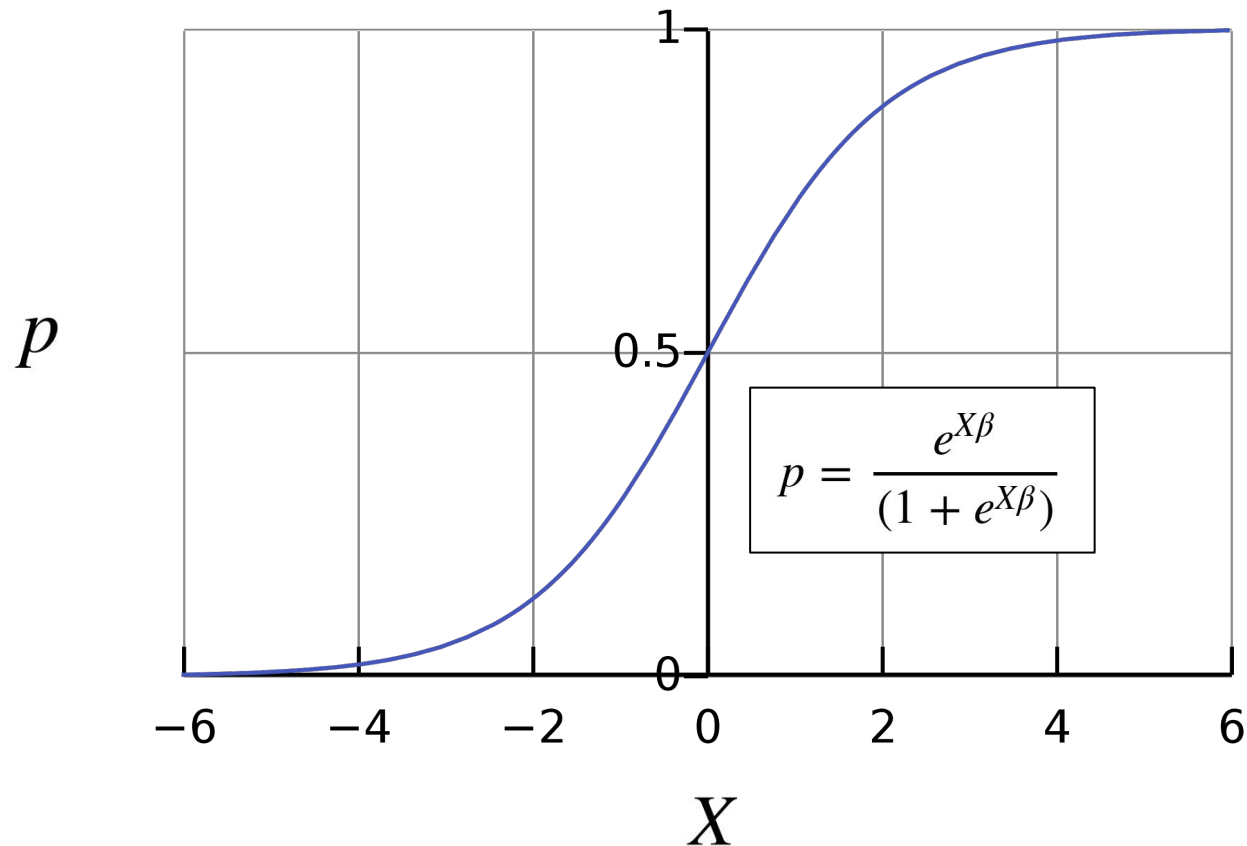
$$p = \frac{e^{\beta_0 + \beta_1 x_1 + \dots}}{(1 + e^{\beta_0 + \beta_1 x_1 + \dots})}$$

*divide by term
to isolate p*

Whew! Does this look familiar?

It's the logistic (aka sigmoid) function.

Logistic function - determining probability given $X\beta$



Determining probabilities in logistic regression - recap

- In logistic regression prediction, $X\beta$ is computed and then placed in the logistic function to determine the probability that each row of data belongs to the positive class.
- The logistic function naturally bounds the probability between 0 and 1.
- In logistic regression training, maximum likelihood is used to find the coefficients β that maximize the likelihood of the existing classifications given the data X .

You'll solve for these coefficients in the Gradient Descent assignment.

Interpreting logistic regression coefficients

How to interpret coefficient β_1 ?

- In linear regression, for a 1 unit increase in x_1 the response \hat{y} will increase by β_1 assuming all other values are held constant.

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots \quad \text{Linear regression}$$

- In logistic regression, for a 1 unit increase in x_1 the *log odds* will increase by β_1 assuming all other values are held constant.

$$\ln(odds) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots \quad \text{Logistic regression}$$

Interpreting logistic regression coefficients

How to interpret coefficient β_1 ?

- In linear regression, for a 1 unit increase in x_1 the response \hat{y} will increase by β_1 assuming all other values are held constant.

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots \quad \text{Linear regression}$$

- In logistic regression, for a 1 unit increase in x_1 the *log odds* will increase by β_1 assuming all other values are held constant.

$$\ln(odds) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots \quad \text{Logistic regression}$$

“... the log odds will increase ...” <- True, but what does that mean? Speak odds instead...

Interpreting logistic regression coefficients

How to interpret coefficient β_1 ?

Recall:

$$odds = e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2}$$

Using the [exponential power rule](#):

$$odds = e^{\beta_0} e^{\beta_1 x_1} e^{\beta_2 x_2}$$

If everything else held constant:

$$odds = C \cdot e^{\beta_1 x_1}$$

In logistic regression, for a 1 unit increase in x_1 the *odds* will increase by e^{β_1} assuming all other values are held constant.

Evaluating logistic regression (and other classification models)

- In a test or hold-out set, you should have data that has true classifications (y_{true}) that you can compare to your predictions (\hat{y})
- Comparing predictions to true values, you can label each prediction a True Positive (TP), a False Positive (FP), a False Negative (FN), and a True Negative (TN):

y_{true}	\hat{y}	label
1	1	TP
0	1	FP
1	0	FN
0	0	TN

Evaluating logistic regression (and other classification models)

- The counts of the number of TP, FP, FN, and TN are typically summarized in a table called, appropriately, a [confusion matrix](#).

		Actual class	
		Cat	Non-cat
Predicted class	Cat	5 True Positives	2 False Positives
	Non-cat	3 False Negatives	17 True Negatives

- Useful metrics can be calculated from the counts in the confusion matrix.

Classification notation and metrics

condition positive (P)

the number of real positive cases in the data

condition negative (N)

the number of real negative cases in the data

true positive (TP)

eqv. with hit

true negative (TN)

eqv. with correct rejection

false positive (FP)

eqv. with false alarm, Type I error

false negative (FN)

eqv. with miss, Type II error

sensitivity, recall, hit rate, or true positive rate (TPR)

$$\text{TPR} = \frac{\text{TP}}{P} = \frac{\text{TP}}{\text{TP} + \text{FN}} = 1 - \text{FNR}$$

precision or positive predictive value (PPV)

$$\text{PPV} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

accuracy (ACC)

$$\text{ACC} = \frac{\text{TP} + \text{TN}}{P + N} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

F1 score

is the harmonic mean of precision and sensitivity

$$F_1 = 2 \cdot \frac{\text{PPV} \cdot \text{TPR}}{\text{PPV} + \text{TPR}} = \frac{2\text{TP}}{2\text{TP} + \text{FP} + \text{FN}}$$

fall-out or false positive rate (FPR)

$$\text{FPR} = \frac{\text{FP}}{N} = \frac{\text{FP}}{\text{FP} + \text{TN}} = 1 - \text{TNR}$$

Breakout 2

Situation	What are trying to minimize? FP or FN? Both?	Therefore you should maximize? (Accuracy, Precision, Recall, F1 score)
Checking for a disease		
Detecting spam in your email		
Seal of quality test result for parachute manufacturing		
Identifying people to target with a marketing campaign		
Deciding to convict someone of a crime		

Breakout 3

Use the true value, the predicted probability of the + class and the given threshold to fill in the predictions in the tables below. Then, construct a confusion matrix and calculate the accuracy, precision, and recall.

An example is filled out for you.

thresh:	0.5		
y_h_prob	y_hat	y_true	label
0.67	1	0	FP
0.20	0	1	FN
0.98	1	1	TP
0.03	0	0	TN

		True	
		1	0
Predicted	1	1 (TP)	1 (FP)
	0	1 (FN)	1 (TN)

Accuracy = $2/4 = 0.5$
Precision = $1/2 = 0.5$
Recall = $1/2 = 0.5$

Example

thresh:	0.20		
y_h_prob	y_hat	y_true	label
0.67		0	
0.20		1	
0.98		1	
0.03		0	

		True	
		1	0
Predicted	1		
	0		

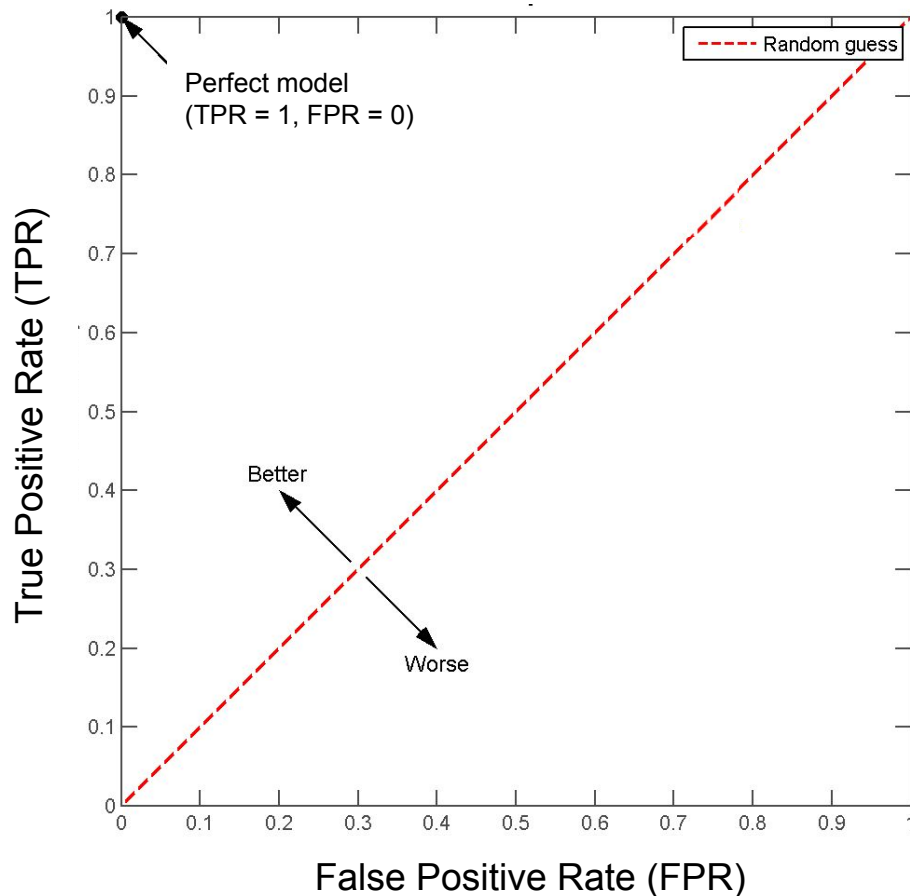
Accuracy =
Precision =
Recall =

Determining the overall performance of a classifier

- You may have noticed that you get different model performance metrics, depending on the threshold you use to determine the classifications from the prediction probabilities.
- This is obviously not ideal - would like to have a metric that incorporates all possible values of the threshold.
- There is a curve and metric that encompass this performance: the [Receiver Operating Characteristic \(ROC curve\)](#).

Receiver Operating Characteristic (ROC) curve

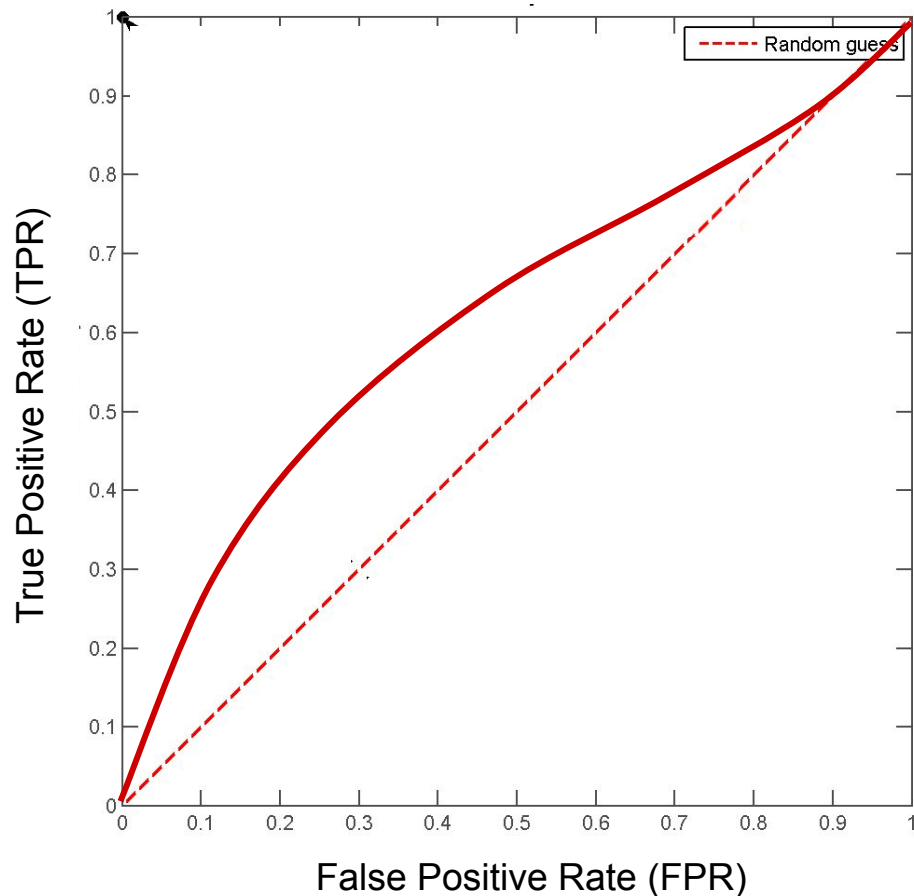
- ROC curve plots the TPR vs the FPR for all thresholds of interest.
- It's easy to get a good TPR with a high FPR (just guessing the positive class all the time).
 - very low threshold usually gives many FPs
- It's difficult to get a good TPR with a low FPR.
 - very high threshold minimizes FPs, but usually miss some TPs
- Total performance quantified by Area Under the Curve (AUC)
 - Perfect: $AUC = 1$
 - Random guessing: $AUC = 0.5$



Receiver Operating Characteristic (ROC) curve

A so-so model

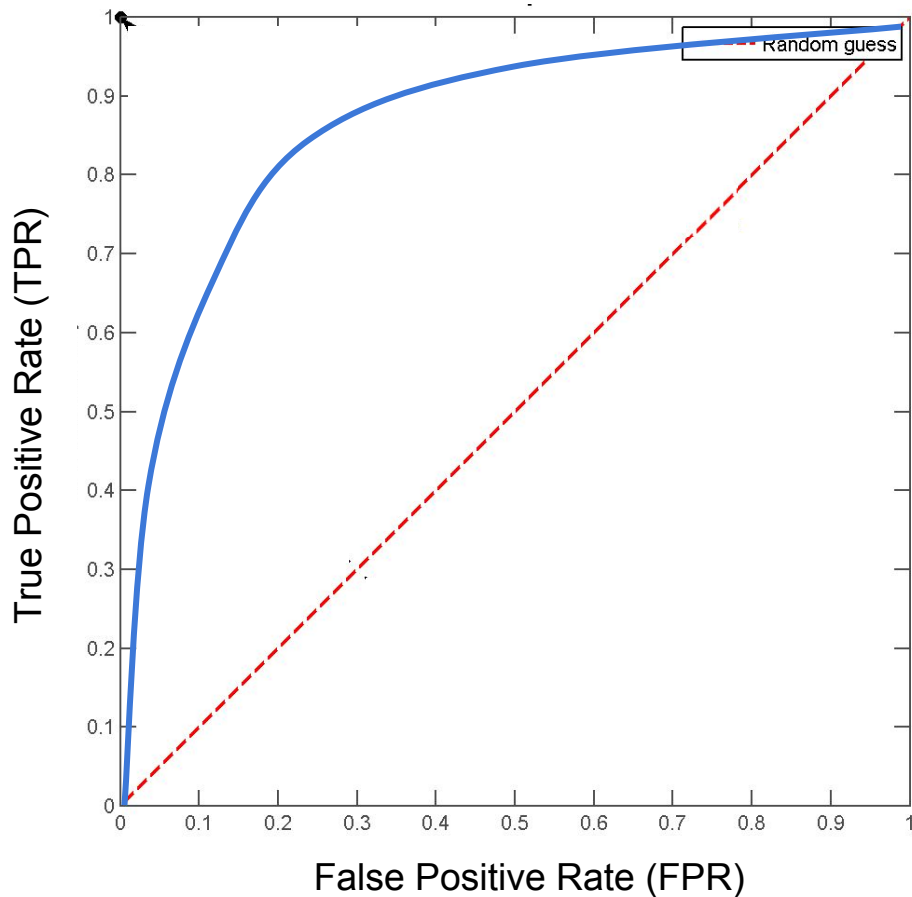
AUC ~ 0.65



Receiver Operating Characteristic (ROC) curve

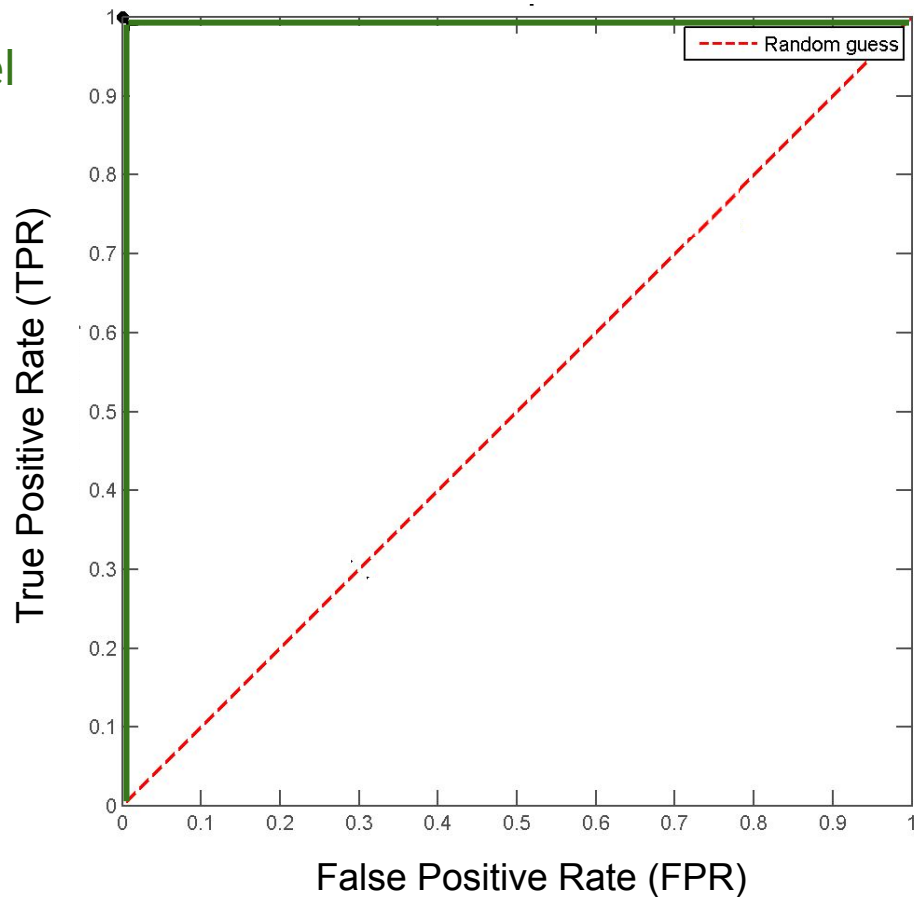
A good model

AUC ~ 0.8



Receiver Operating Characteristic (ROC) curve

An unreasonably good model
AUC ~ 0.999999



Constructing a ROC curve

1. Fit a classifier to your training data and get positive class prediction probabilities.
2. Sort these probabilities - can be low to high or high to low. These probabilities are doubling as your thresholds.
3. Compare all your predicted probabilities to each threshold to get predictions, and from predictions and true values fill out a confusion matrix. You should get a confusion matrix for each threshold.
4. From each confusion matrix, calculate the TPR and FPR.
5. The FPR and TPR for a given threshold and confusion matrix are an (x, y) pair for the ROC - plot them!

Quick demo

`horse_or_dog.ipynb`

Objectives

After this lecture you should be able to:

- Explain the difference between regression and classification
- Decide which class to make positive in binary classification
- Describe how classifications are made using predicted probabilities and a threshold
- Describe how logistic regression determines probabilities
- Interpret logistic regression coefficients
- Calculate (and define) TP, FP, FN, and TN
- Construct a confusion matrix
- Define metrics useful in classification, and give examples where each might be useful
- Describe what a ROC curve is, how it gets constructed, and how it describes the overall performance of a classifier