# Logistic Regression for Predicting Binary Responses

#### **Dr Rebecca Barter**



#### Labeled data

#### Outcome of interest: Amount spent on healthcare in next year

#### Labeled data

Age	Sex	Weight	Diabetes	Healthare expenses
54	Μ	132	Ν	4,300
76	F	155	Y	3,467
49	Μ	166	Y	103
39	F	129	Ν	844
47	Μ	177	Ν	6,591
70	F	192	Ν	8,089



#### Labeled data

Outcome of interest: Amount More than 1K spent on healthcare in next year

Age	Sex	Weight	Diabetes	Healthare expenses
54	Μ	132	Ν	4,300
76	F	155	Y	3,467
49	Μ	166	Y	103
39	F	129	Ν	844
47	Μ	177	Ν	6,591
70	F	192	Ν	8,089



#### Labeled data

**Outcome of interest:** 

#### More than 1K spent on healthcare in next year

Age	Sex	Weight	Diabetes	Healthcare expenses
54	Μ	132	Ν	1
76	F	155	Y	1
49	Μ	166	Y	0
39	F	129	Ν	0
47	Μ	177	Ν	1
70	F	192	Ν	1



### (New) Unlabeled data

**Outcome of interest:** 

More than 1K spent on healthcare in next year

Age	Sex	Weight	Diabetes	Healthcare expenses
44	М	165	Ν	?
69	F	161	Y	?
78	М	170	N	?
66	Μ	191	Ν	?



#### For continuous responses...





#### For binary responses...



#### For binary responses...





#### **Logistic regression**



Y-axis  
P(hlthx1k) = 
$$\frac{e^{b_0 + b_1 age}}{1 + e^{b_0 + b_1 age}}$$

#### This is the logistic function

- It is always between 0 and 1
- Values between 0 and 1 are interpreted as the *probability* of >1K healthcare expenditure

### **Logistic regression**



Y-axis  

$$P(hlthx1k) = \frac{e^{b_0 + b_1 age}}{1 + e^{b_0 + b_1 age}}$$

#### This is the logistic function

- It is always between 0 and 1
- Values between 0 and 1 are interpreted as the *probability* of >1K healthcare expenditure

#### Training an algorithm and generating a prediction

#### Labeled data

**Unlabeled data** 

**Prediction** 



# **Evaluating binary response predictions**

Dr Rebecca Barter



# **Training and testing sets**

Since the only labeled data is usually the data we have, we need to **split** our data into training and testing sets



# **Evaluating binary response predictions**

When evaluating binary response predictions, we are comparing 0's and 1's.

(Observed) Hlthx1k		Predicted Hlthx1k
1	$\longleftrightarrow$	1
1	$\longleftrightarrow$	1
0	$\leftrightarrow$	0
0	←→	1
1	$\leftrightarrow$	1
1	$\longleftrightarrow$	1

#### Accuracy

Proportion of responses where observed == predicted

Accuracy = 5/6 = 88.3%



# Accuracy is not everything

A skin cancer image dataset contains images for which:

5% are melanoma and 95% are non-melanoma



What is the accuracy of an algorithm that always predicts "non-melanoma"?95%How accurate is the algorithm for the non-melanoma patients?100%How accurate is the algorithm for the melanoma patients?0%



### **Confusion Matrix**



	Observed Yes hlthx1k (1)	Observed No hlthx1k (0)
Predicted Yes Hlthx1k (1)	4	1
Predicted No Hlthx1k (0)	0	1

Accuracy 
$$= \frac{A+D}{A+B+C+D} = 5/6 = 0.83$$

Sensitivity/Recall (True Positive) = 
$$\frac{A}{A+C} = 4/4 = 1$$

Specificity (True Negative) = 
$$\frac{D}{B+D} = 1/2 = 0.5$$

	Observed Yes Hlthx1k (1)	Observed No Hlthx1k (0)
Predicted Yes Hlthx1k (1)	А	В
Predicted No Hlthx1k (0)	С	D



#### **Confusion Matrix**

Accuracy = 
$$\frac{A+D}{A+B+C+D} = \frac{23+45}{23+10+6+45} = 0.81$$
  
Sensitivity (TP) =  $\frac{A}{A+C} = \frac{23}{23+6} = 0.79$   
Specificity (TN) =  $\frac{D}{B+D} = \frac{45}{45+10} = 0.82$   
 $\frac{Ves hithx1k}{(1)} = \frac{Ves hithx1k}{23} = 0.81$   
 $\frac{Ves hithx1k}{(1)} = \frac{Ves hithx1k}{(1)} = 0.82$ 

Observed

Observed

# **ROC curves for binary** predictions

Dr Rebecca Barter



#### **Converting probability predictions to binary predictions**



We can choose the probability threshold.

Predict "Spent 1K on healthcare" if:

Threshold: P(hlthx1k) > 0.5



#### **Converting probability predictions to binary predictions**



We can choose the probability threshold.

Predict "Spent 1K on healthcare" if:

Threshold: P(hlthx1k) > 0.75



#### **Converting probability predictions to binary predictions**



We can choose the probability threshold.

Predict "Spent 1K on healthcare" if:

Threshold: P(hlthx1k) > 0.25

Each choice will have different accuracy/sensitivity/specificity

### **ROC curve**

Predict "healthcare expenses over 1K" if predicted probability is greater than a threshold value





### **ROC curve**

Predict "healthcare expenses over 1K" if predicted probability is greater than a threshold value



Often the "best" threshold value is the proportion of positive class obs

#### **ROC curve**

The "best" ROC curve is the one that gets the closest to the top-left corner





# Area under the ROC curve (AUC)

The "best" model is the one that has the highest Area Under the Curve (AUC)





# Area under the ROC curve (AUC)

The "best" model is the one that has the highest Area Under the Curve (AUC)





# Area under the ROC curve (AUC)

The "best" model is the one that has the highest Area Under the Curve (AUC)



# Class imbalance in binary prediction problems

Dr Rebecca Barter



#### **Class imbalance**

Class imbalance occurs when you have many more observations in one class than the other

Actual "observed" class



Problem: algorithms will tend to "over predict" the "majority" class

**Predicted class** 





### Sampling methods for class imbalance

# Original imbalanced data

#### Upsampling

When training predictive algorithms, use a version of the training data where the **minority** class observations are randomly **repeated** so their number matches the majority class

Repeated data observations



#### Sampling methods for class imbalance

# Original imbalanced data

#### Downsampling

When training predictive algorithms, use a version of the training data where only a random **subsample** of the **majority** class is used so their number matches the minority class





#### Sampling methods for class imbalance

# Original imbalanced data

#### **SMOTE (Synthetic Minority Over-sampling Technique)**

When training predictive algorithms, use a version of the training data where **new** "synthetic" minority class data points are created based on "interpolating" existing minority class observations

"Synthetic" (made-up) data observations

#### **Choosing your threshold for class imbalance**

When converting a binary probability prediction to a binary 0/1 value, we use a **threshold**.

The default threshold is predict positive class when **P > 0.5** 

Improve predictive performance by setting the **threshold** to be the **proportion** of positive (minority) class observations in training data

Proportional threshold =  $\frac{6}{34} = 0.176$ 

Predict positive class when **P > 0.176** 

