

— Classification Metrics I

Data Science Process

1. Define the problem
2. Gather data
3. Explore data
4. Model with data
5. Evaluate model
6. Answer problem



Framing

Remember the regression metrics lesson from last week, where we explored different methods for evaluating the performance of **regression models**.

We'll do the same thing today, but for **classification models**.

- In regression, we quantify the performance of our model by comparing predicted and observed values in some capacity.
- We'll do the same thing in classification... but predicted and observed are categories, so it's slightly different.

We're going to focus on **binary classification problems**.

Evaluating Our Model

Suppose you build a logistic regression model predicting whether or not people will vote. You make predictions for 100 people, then check to see if they actually voted.

- There are 40 people you predicted to vote who did vote.
- There are 20 people you predicted to vote who didn't vote.
- There are 15 people you predicted to stay home who did vote.
- There are 25 people you predicted to stay home who didn't vote.



Evaluating Our Model

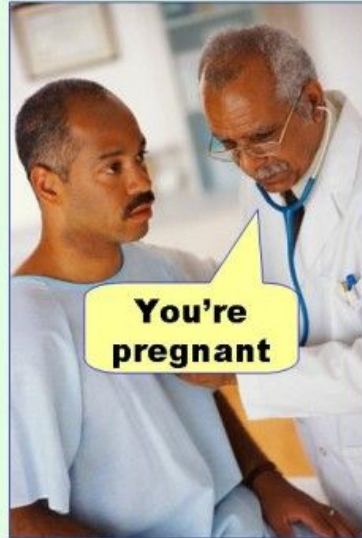
Suppose you build a logistic regression model predicting whether or not people will vote. You make predictions for 100 people, then check to see if they actually voted.

- There are 40 people you predicted to vote who did vote.
 - These are called **true positives**.
- There are 20 people you predicted to vote who didn't vote.
 - These are called **false positives**.
- There are 15 people you predicted to stay home who did vote.
 - These are called **false negatives**.
- There are 25 people you predicted to stay home who didn't vote.
 - These are called **true negatives**.

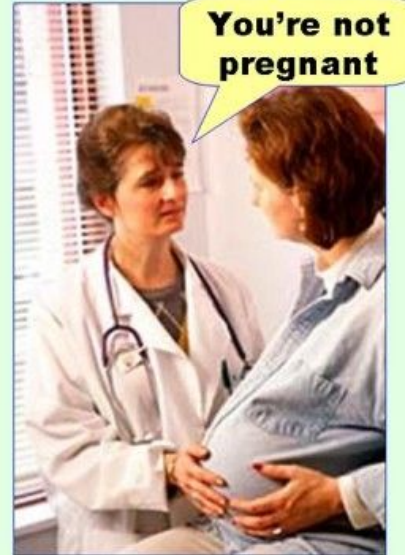


Evaluating Our Model

Type I error
(false positive)



Type II error
(false negative)



Evaluating Our Model

How do I keep true positives/true negatives/false positives/false negatives straight?

- First word (true/false): Was I right?
- Second word (positive/negative): What did I predict?



Evaluating Our Model

How do I keep true positives/true negatives/false positives/false negatives straight?

- First word (true/false): Was I right?
- Second word (positive/negative): What did I predict?

What is it called if I correctly predicted that someone does not vote?

Evaluating Our Model

How do I keep true positives/true negatives/false positives/false negatives straight?

- First word (true/false): Was I right?
- Second word (positive/negative): What did I predict?

What is it called if I incorrectly predicted that someone does vote?

Confusion Matrix

It's helpful for us to list out the number of each category in a 2x2 grid called a **confusion matrix**.

	Actual Positive	Actual Negative
Predicted Positive		
Predicted Negative		

Confusion Matrix

It's helpful for us to list out the number of each category in a 2x2 grid called a **confusion matrix**.

	Actual Positive	Actual Negative
Predicted Positive		
Predicted Negative		

The axes or ordering of “Yes” vs. “No” may be rearranged!

Be clear what “Yes” / “Positive” means.

Confusion Matrix

A confusion matrix is a convenient way for us to visualize how our model performs.

However, there are metrics that can help us to summarize performance with one number.

- Accuracy
- Misclassification Rate
- Sensitivity
- Specificity
- Precision



Accuracy

Interpretation: What percentage of observations did I **correctly** predict?

$$\text{Accuracy} = \frac{\text{All Correct}}{\text{All Predictions}} = \frac{TP + TN}{TP + FP + TN + FN}$$

	Actual Positive	Actual Negative
Predicted Positive	40	20
Predicted Negative	15	25

Misclassification Rate

Interpretation: What percentage of observations did I **incorrectly** predict?

$$\text{Misclassification Rate} = \frac{\text{All Incorrect}}{\text{All Predictions}} = \frac{FP + FN}{TP + FP + TN + FN} = 1 - \text{Acc}$$

	Actual Positive	Actual Negative
Predicted Positive	40	20
Predicted Negative	15	25

Sensitivity

Interpretation: Among those who will vote, how many did I get correct?

$$\text{Sensitivity} = \frac{\text{True Positives}}{\text{All Positives}} = \frac{TP}{TP + FN} = \frac{TP}{P}$$

a.k.a. True Positive Rate, Recall

	Actual Positive	Actual Negative
Predicted Positive	40	20
Predicted Negative	15	25

Specificity

Interpretation: Among those who will not vote, how many did I get correct?

$$\text{Specificity} = \frac{\text{True Negatives}}{\text{All Negatives}} = \frac{TN}{TN + FP} = \frac{TN}{N}$$

a.k.a. True Negative Rate

	Actual Positive	Actual Negative
Predicted Positive	40	20
Predicted Negative	15	25

Precision

Interpretation: Among those I predicted to vote, how many did I get correct?

$$\text{Precision} = \frac{\text{True Positives}}{\text{Predicted Positives}} = \frac{TP}{TP + FP}$$

a.k.a. Positive Predictive Value

	Actual Positive	Actual Negative
Predicted Positive	40	20
Predicted Negative	15	25

Example

Suppose I'm working on a fraud analytics team and our goal is to detect fraudulent credit card transactions. I take a random sample of 500 transactions. Of these transactions, 50 are fraudulent. I predict 100 total fraudulent transaction, 45 of which are correct.

1. Identify the TP, TN, FP, FN and construct a confusion matrix.
2. Calculate the accuracy, misclassification rate, positive predictive value, recall, and true negative rate.



Example

Suppose I'm working on a fraud analytics team and our goal is to detect fraudulent credit card transactions. I take a random sample of 500 transactions. Of these transactions, 50 are fraudulent. I predict 100 total fraudulent transaction, 45 of which are correct.

When building my classification model, I want to optimize one of the above metrics. Given the use-case of identifying fraudulent transactions, which metric should I optimize as I build my model?

Final Notes

We explored binary classification problems today.

We can construct confusion matrices for 3+ categories and calculate a lot of these metrics (accuracy, misclassification error, etc.), but they get a lot more complicated.

These get *especially* complicated when working with **ordinal data**.



