

Classifier Evaluation

Zach Gulde

Why Evaluate?

Plan → Acquire → Prepare → Explore → Model → Deliver

You are here

- Quantifying model performance allows to compare models (ML or otherwise!)
- How we quantify performance is key
- Many different ways to quantify depending on what we want to optimize for

Vocab


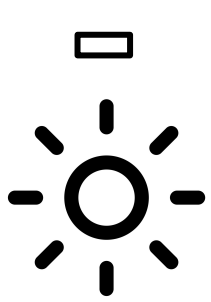
- Classifier
 - Binary
 - Multi-Class
- Evaluation Metric
- Label / target / outcome
- Actual and Predicted Values



Classification Outcomes



- True Positive (TP)
- True Negative (TN)
- False Positive (FP)
- False Negative (FN)

 rain

 no rain

		Actual	
Predicted	+		
	-		

  rain

  no rain

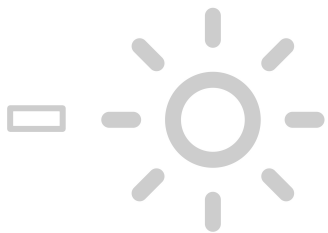


Actual




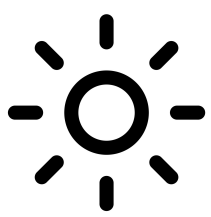


Note that the choice of +/- is arbitrary!

Predic



 rain

 no rain


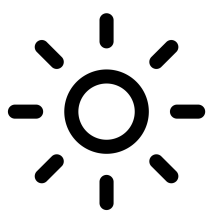




		Actual	
Predicted	+		
	-	<div>TP</div>	



rain



no rain


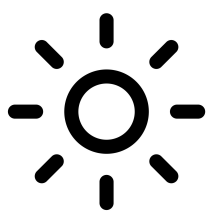
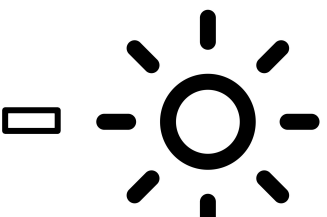

		Actual	
Predicted	+		
	-	<div>TP</div>  	<div>FP</div>  





rain





no rain

		Actual	
Predicted	+		
	-		


TP

FP

FN


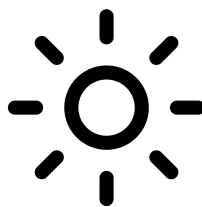
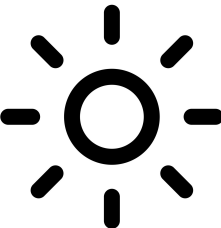







rain





no rain

		Actual	
Predicted	+		
	-		


TP


FP

FN



TN



Depending on the domain / business case, not all outcomes are equally important!

Different classifier evaluation metrics look at or ignore certain outcomes.






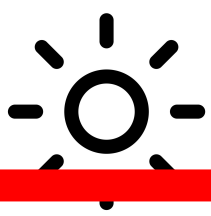



rain



no rain

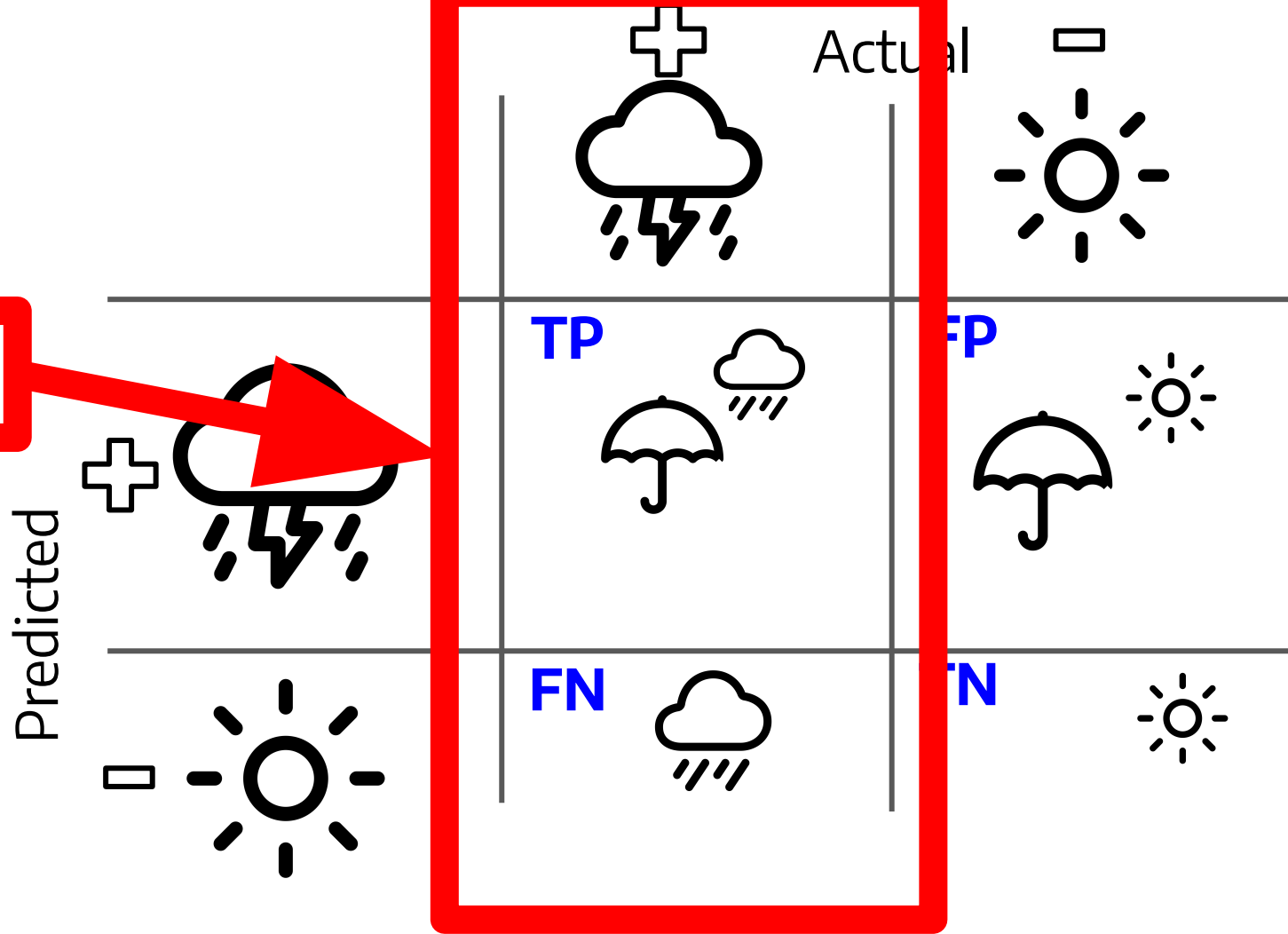
Accuracy

$$\frac{TP + TN}{TP + FP + FN + TN}$$

		Actual	
Predicted	+	<div>TP</div> <div></div>	<div>FP</div> <div></div>
	-	<div>FN</div> <div></div>	<div>TN</div> <div></div>



$$\frac{TP}{TP + FN}$$





Precision

Actual

Predicted

TP

FP

FN

TN

$$\frac{TP}{TP + FP}$$

Actual



Predicted



Let's look at an example...

Actual



Predicted



Actual



Predicted



Actual	+	+	-	+	-	-	+
Predicted	-	+	+	+	+	-	+
	✗	✓	✗	✓	✗	✓	✓















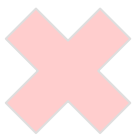






Accuracy
correct / total

$$= 4 / 7 \approx 57\%$$

Actual	+	+	-	+	-	-	+
Predicted	-	+	+	+	+	-	+
	✗	✓	✗	✓	✗	✓	✓

Precision

$$\frac{TP}{TP + FP}$$

Actual							
Predicted							
							

Precision

$$\frac{TP}{TP + FP}$$






















$$= 3 / 5 = 60\%$$

How good are our
positive predictions?

Actual	+	+	-	+	-	-	+
Predicted	-	+	+	+	+	-	+
	✗	✓	✗	✓	✗	✓	✓

Recall

$$\frac{TP}{TP + FN}$$

Actual							
Predicted							
							

Recall

$$\frac{TP}{TP + FN}$$

$$= 3 / 4 = 75\%$$

How many of the
actually positive
cases do we catch?

Actual	+	+	-	+	-	-	+
Predicted	+	+	+	+	+	+	+

Consider a classifier that always predicts positive...

Actual	+	+	-	+	-	-	+
Predicted	+	+	+	+	+	+	+

$$\text{Accuracy} = 4 / 7 = 57\%$$

$$\text{Recall} = 4 / 4 = 100\%$$

$$\text{Precision} = 4 / 7 = 47\%$$

Actual	+	+	-	+	-	-	+
Predicted	-	-	-	-	-	-	-

What if we always predict negative?

Actual	+	+	-	+	-	-	+
Predicted	-	-	-	-	-	-	-

$$\text{Accuracy} = 3 / 7 = 43\%$$

$$\text{Recall} = 0 / 4 = 0\%$$

$$\text{Precision} = 0 / 0 = \text{undef}$$

Classifier Evaluation Metrics

Recap

- Accuracy: Overall Performance
- Recall: When we don't want to "miss out" on an actually positive case
- Precision: When a positive prediction is expensive

Other Metrics

- Sensitivity: aka recall
- Specificity: recall for the negative class
- F1 score: harmonic mean of precision and recall
- ROC / AUC