



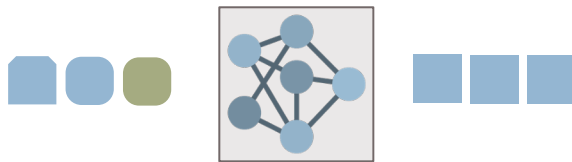
Cheat Sheet – Imbalanced Data in Classification

Blue: Label 1 : 

Green: Label 0 : 

$$\text{Accuracy} = \frac{\text{Correct Predictions}}{\text{Total Predictions}}$$

Classifier that always **predicts label blue** yields prediction accuracy of 90%



Accuracy doesn't always give the correct insight about your trained model

Accuracy: %age correct prediction	Correct prediction over total predictions	One value for entire network
Precision: <u>Exactness</u> of model	From the detected cats, how many were actually cats	Each class/label has a value
Recall: <u>Completeness</u> of model	Correctly detected cats over total cats	Each class/label has a value
F1 Score: Combines Precision/Recall	Harmonic mean of Precision and Recall	Each class/label has a value

Performance metrics associated with Class 1

		Actual Labels	
		1	0
Predicted Labels	1	True Positive	False Positive
	0	False Negative	True Negative

(Is your prediction correct?) (What did you predict)

True Negative

(Your prediction is correct) (You predicted 0)

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

$$\text{False +ve rate} = \frac{FP}{TN + FP}$$

$$\text{F1 score} = 2x \frac{(\text{Prec} \times \text{Rec})}{(\text{Prec} + \text{Rec})}$$

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


$$\text{Specificity} = \frac{TN}{TN + FP}$$


$$\text{Recall, Sensitivity} = \frac{TP}{TP + FN}$$


True +ve rate


Possible solutions

- Data Replication:** Replicate the available data until the number of samples are comparable
- Synthetic Data:** Images: Rotate, dilate, crop, add noise to existing input images and create new data
- Modified Loss:** Modify the loss to reflect greater error when misclassifying smaller sample set
- Change the algorithm:** Increase the model/algorithm complexity so that the two classes are perfectly separable (Con: Overfitting)

Blue: Label 1 : 

Green: Label 0 : 

Blue: Label 1 : 

Green: Label 0 : 

$$\text{loss} = a * \text{loss}_{\text{green}} + b * \text{loss}_{\text{blue}} \quad a > b$$

No straight line ($y=ax$) passing through origin can perfectly separate data. **Best solution:** line $y=0$, predict all labels blue

→ Increase model complexity →

Straight line ($y=ax+b$) can perfectly separate data. Green class will no longer be predicted as blue

