Supervised Learning

COSC 410: Applied Machine Learning

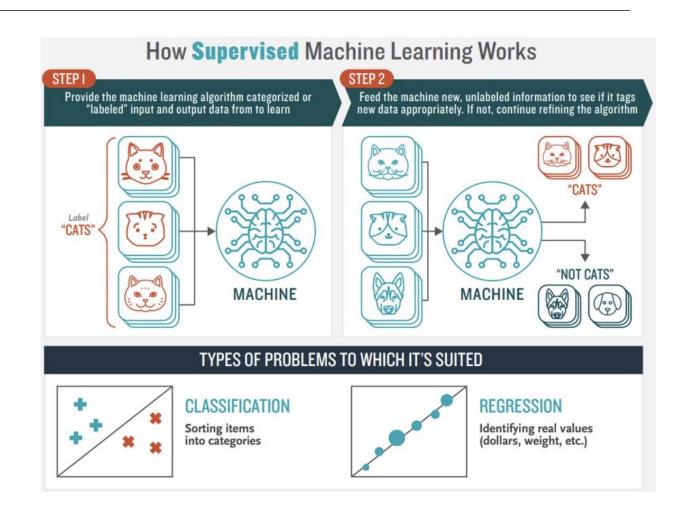
Spring 2022

Prof. Apthorpe

2/1/22

Outline

- Classification & Regression
- Training, Validation, & Test Sets
- Golden Rule of Supervised ML
- Cross-Validation
- Performance Metrics
 - Regression Error Functions
 - Accuracy & Imbalanced Classes
 - ROC curve & AUC
 - Precision/Recall & F₁ Scores
 - Confusion Matrices

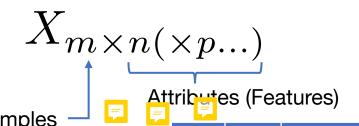


Classification & Regression

Regression: Continuous labels

Classification: Binary or multi-class labels

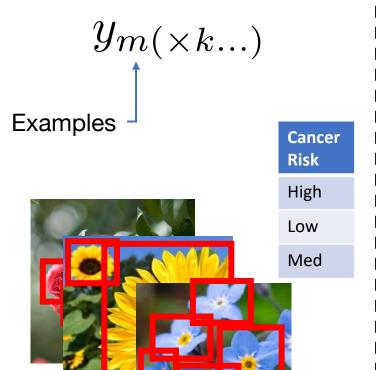
<u>Dataset</u>



Age	His.	Sillokes
43	Υ	Υ
29	N	N
76	N	N



Labels



Ultimate Goal

Train model h() such that for data X and correct labels $\mathcal Y$

$$h(X) \approx y$$

"The model's predictions are close" to the real labels"

*Need to define "close"

Inherent Problem

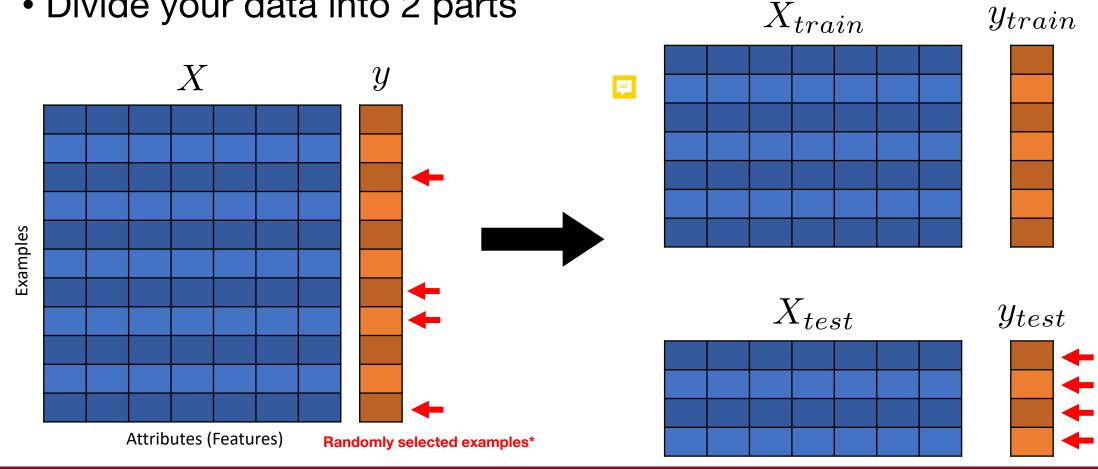
- If you don't know the correct labels for new data,
 how do you evaluate whether your model generalizes?
 - Why is this important?



- Solution:
 - Use the data and labels you do have for both training and evaluation

Training Sets & Test Sets

Divide your data into 2 parts



Training Sets & Test Sets

- Use the **training set** to train the model h()
 - \bigstar Minimize the training error $E(h(X_{train}), y_{train})$
 - Choice of error (loss) function
 - Square error, absolute error, hinge loss, cross-entropy loss, ...

- Use the test set to measure the model's performance and ability to generalize to new data
 - Measure the test error $E(h(X_{test}), y_{test})$

The Golden Rule of Supervised ML

Never train on the test set!!!



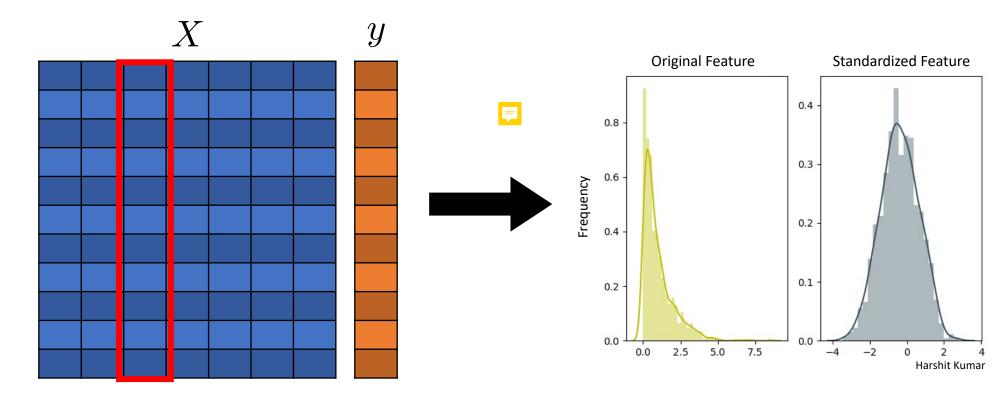
Nothing in your model should be influenced by the examples in the test set. This includes preprocessing steps, training examples, and hyperparameter settings!

The test set is **only** for reporting your final performance

Imagine you don't even have the test set when preparing your model

Common Ways to Break the Golden Rule

- Preprocessing based on entire dataset
 - E.g. scaling the variance of each feature to a normal distribution



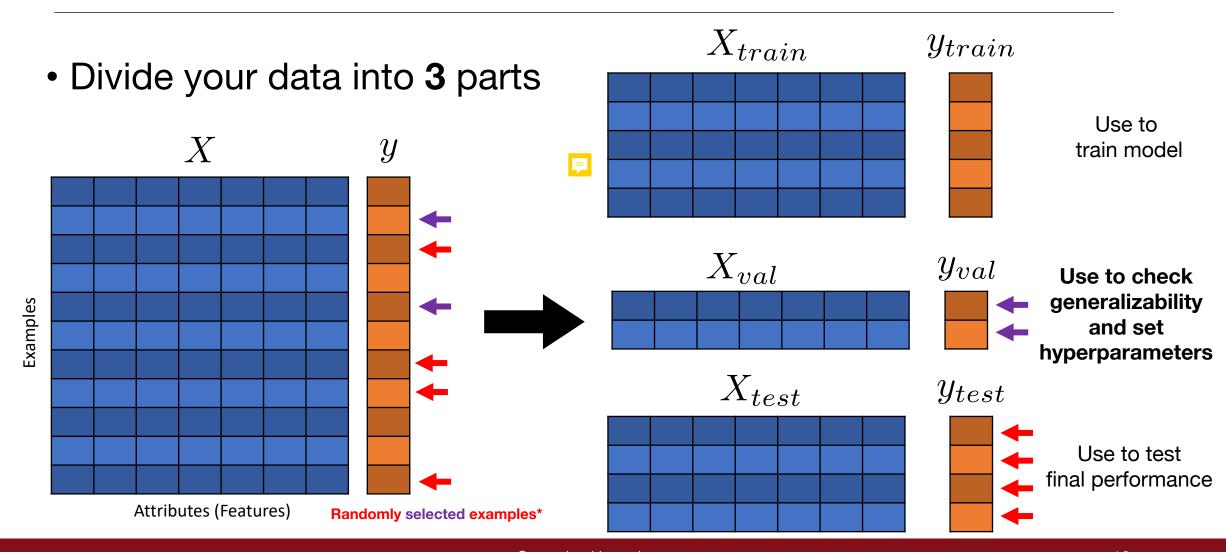
Common Ways to Break the Golden Rule

Modifying model based on test error

- Train on the training set (good)
- Test on the test set (good)
- Think "that worked poorly, let me modify some parameters..." (BAD)

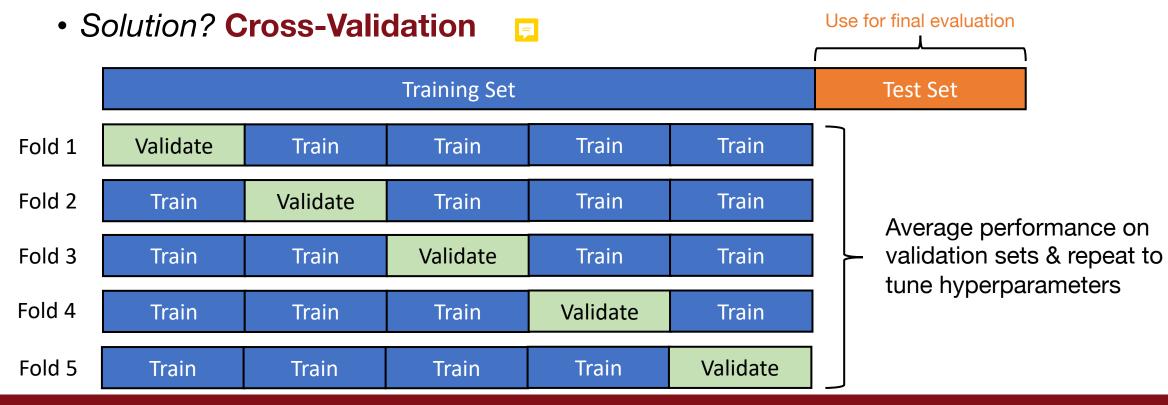
- Why is this bad?
- Why is this tempting?
 - Without testing on data you didn't use to train, how do you improve generalization?

Training / Validation / Test Sets



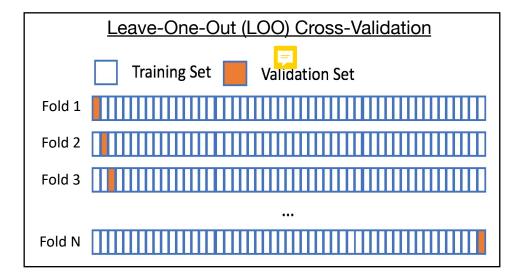
Cross-Validation

- Train / validate / test split leaves much less data for training
 - Can reduce final performance (more data usually improves models)



Cross-Validation

- How many folds?
 - More folds
 - Better estimation of model performance
 - Computation intensive
 - Need to train model from scratch for each fold
 - Fewer folds (5 or 10 standard)
 - Worse estimation of model performance
 - Fewer computations (practical for deep learning)





Performance Metrics

How to Measure ML Success

Performance Metrics

- Performance metric functions
 - Produce value indicating the success of the model
 - Compare predicted labels to actual labels

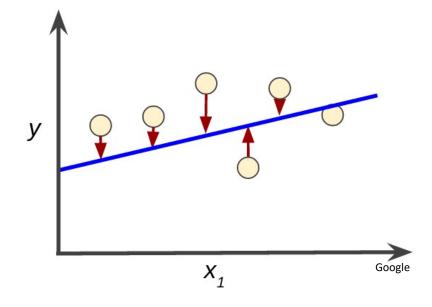
Regression Loss Functions

Mean square error

$$MSE = \frac{\sum_{i=1}^{n} (y_i - \hat{y}_i)^2}{n}$$

Mean absolute error

$$MAE = \frac{\sum_{i=1}^{n} |y_i - \hat{y}_i|}{n}$$



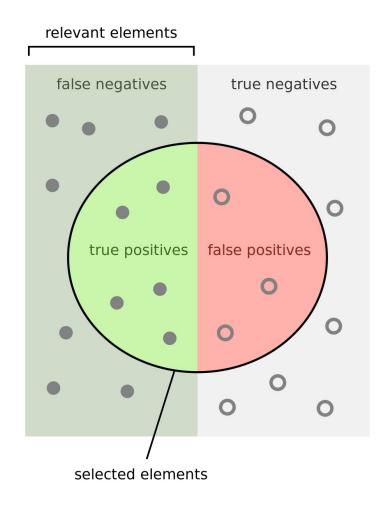
Classification Accuracy

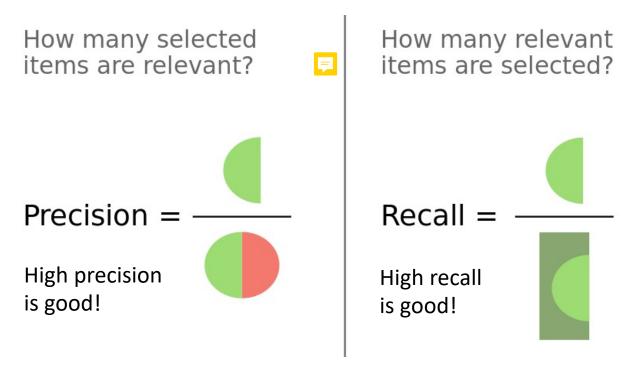
 $Accuracy = \frac{Number of Correct Predictions}{Number of Total Predictions}$

 Accuracy is intuitive, but deceptive if classes are unbalanced Trivial Solution: Always guess "Not a platypus" 90% Accuracy!!!

Paolo Perrotta

Precision/Recall & F₁ Scores



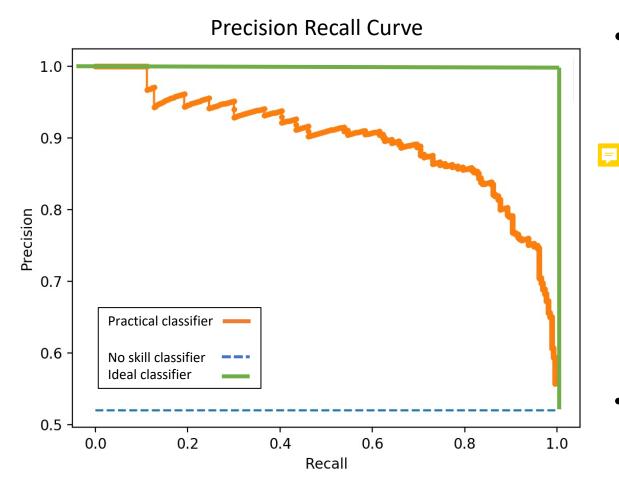


- What does high precision, low recall mean conceptually?
 - In a real-world context? When is this good/bad?
- What does low precision, high recall mean conceptually?
 - In a real-world context? When is this good/bad?

Performance Metrics

- Most ML algorithms are tunable
 - Allow you to trade-off different performance priorities
 - Try a range of values to see performance curve

Precision/Recall & F₁ Scores

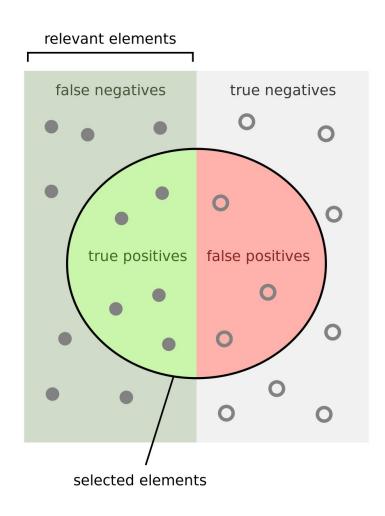


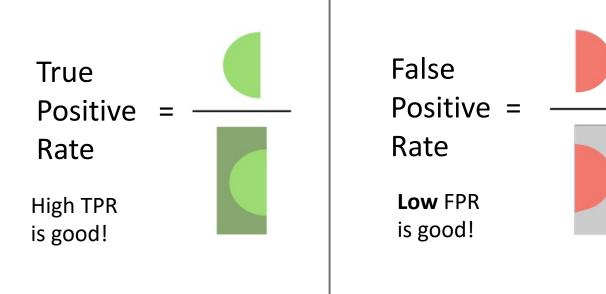
Often report the F₁ score
 (harmonic mean of precision & recall)

 $F_1 = 2 \cdot \frac{precision \cdot recall}{precision + recall}$

 Harmonic mean gives more weight to low values → why does this matter?

Receiver Operating Characteristic (ROC)

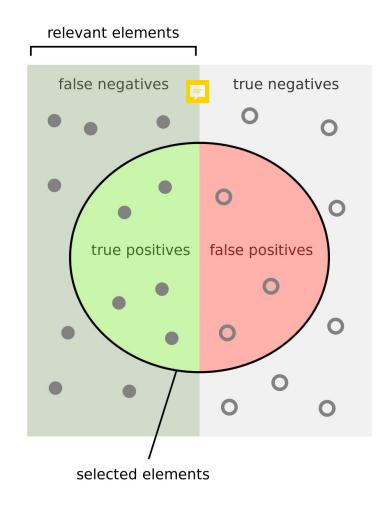


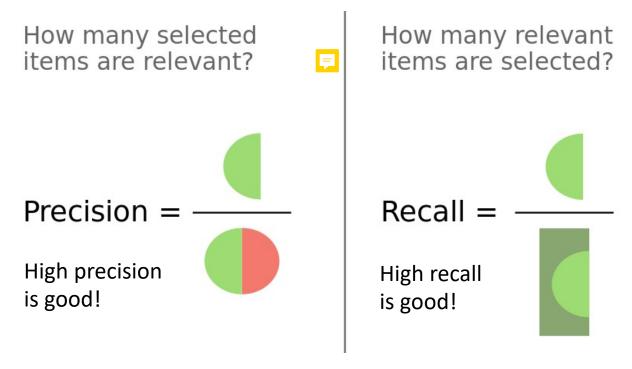


This definition works for binary classification ROC can also be generalized to multiclass classification

Wikipedia

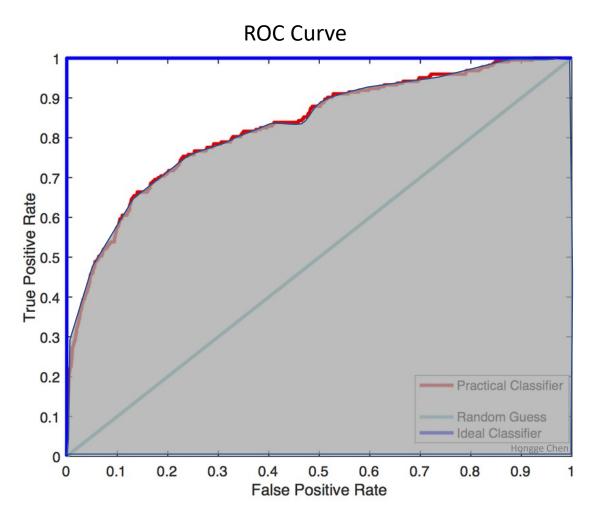
Precision/Recall & F₁ Scores





- What does high precision, low recall mean conceptually?
 - In a real-world context? When is this good/bad?
- What does low precision, high recall mean conceptually?
 - In a real-world context? When is this good/bad?

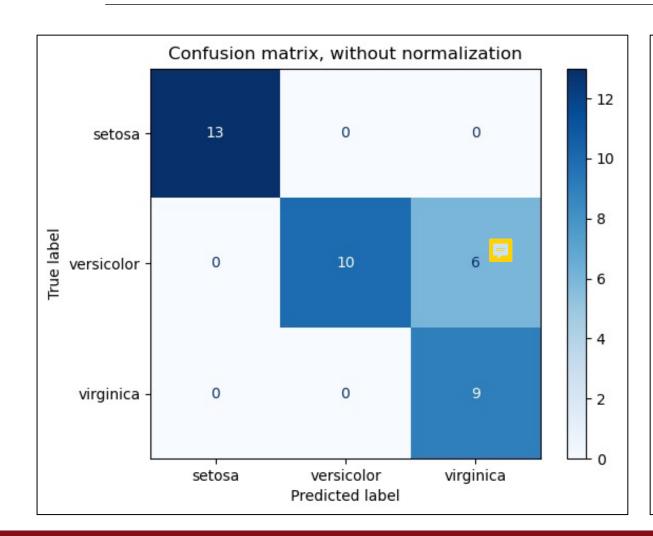
Receiver Operating Characteristic (ROC)

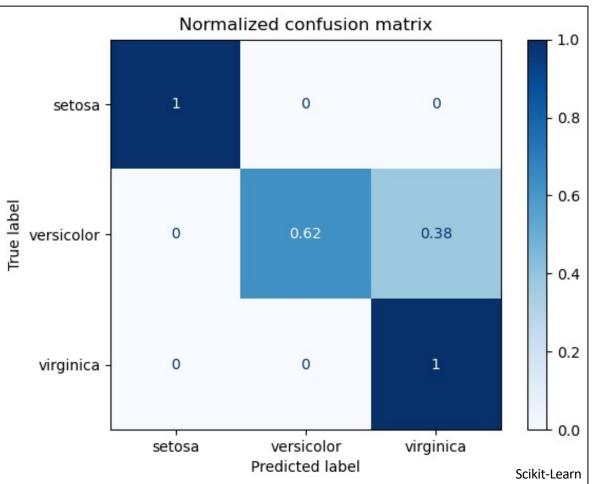


- Often report the area under the curve (AUC) as indicator of model performance
 - Ideal AUC = 1

Confusion Matrices







Programming Practice

SupervisedLearning.ipynb