# Overfitting, Model Selection

CSC 461: Machine Learning
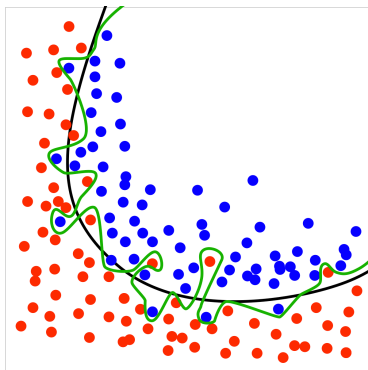
Fall 2021

Prof. Marco Alvarez
University of Rhode Island
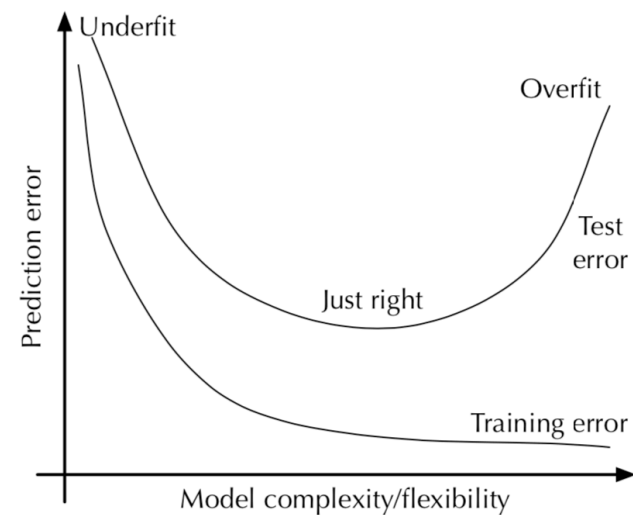
---

# Overfitting
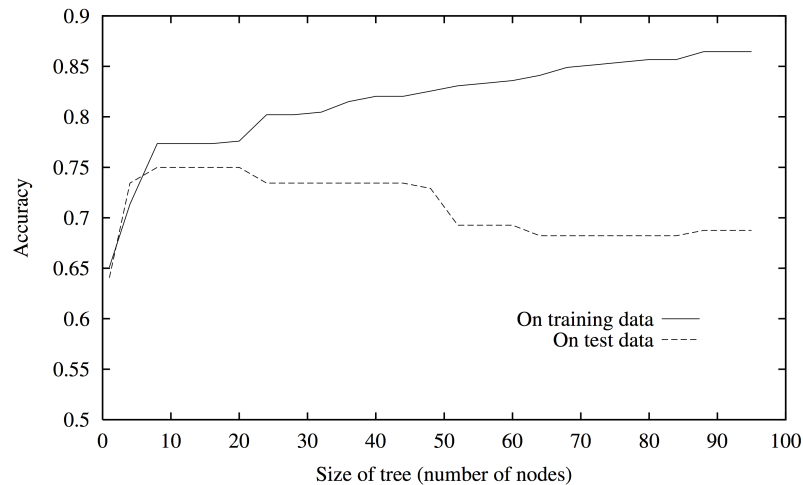
---

# Overfitting

‣ Learning a **model** that "knows" the training data very well but does not **generalize**



---

# Model complexity

# Model complexity (DTs)



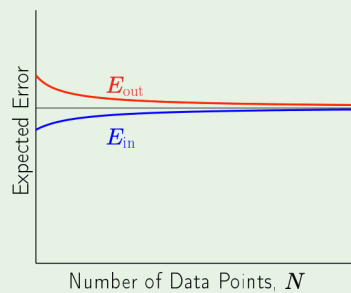Machine Learning, Tom Mitchell, McGraw Hill, 1997

# Overfitting

‣ Reasons

  ✓ model is **too complex**

  ✓ model is **fitting noise** present in the training data

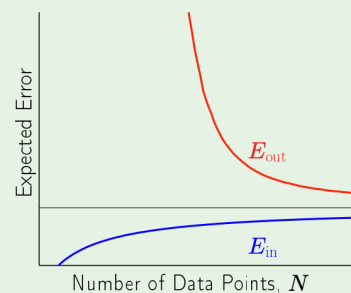  ✓ training data is **not a representative sample** of the distribution

‣ How to prevent?

  ✓ use **more training data**

  ✓ use **fewer features**

  ✓ **regularize** your model

# Number of data instances



Simple Model          Complex Model

https://work.caltech.edu/lectures.html
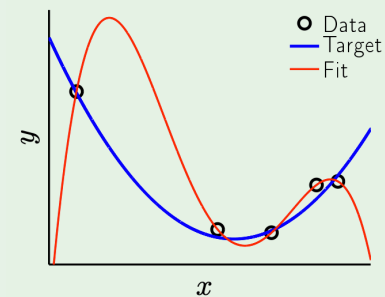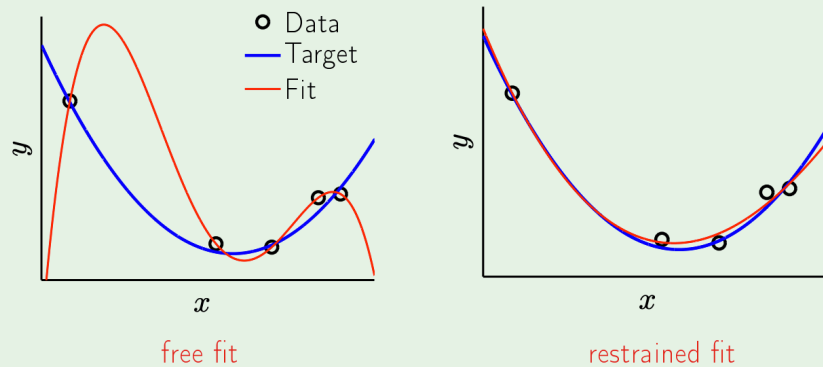
# Restricting the model

‣ Imagine the target function below …

  ✓ 5 noisy data points and a 4th order polynomial fit

  ✓ what can you say about training error? test error?



https://work.caltech.edu/lectures.html

## Restricting the model



free fit

restrained fit

https://work.caltech.edu/lectures.html

## Model Evaluation

## Confusion matrix (2 classes)

|  | POSITIVE (1) | NEGATIVE (0) |
|---|---|---|
| POSITIVE (1) | TP | FN |
| NEGATIVE (0) | FP | TN |

Actual values? Predicted values?

## Evaluation metrics (2 classes)

**sensitivity, recall, hit rate, or true positive rate (TPR)**

$$\text{TPR} = \frac{\text{TP}}{\text{P}} = \frac{\text{TP}}{\text{TP} + \text{FN}} = 1 - \text{FNR}$$

**specificity, selectivity or true negative rate (TNR)**

$$\text{TNR} = \frac{\text{TN}}{\text{N}} = \frac{\text{TN}}{\text{TN} + \text{FP}} = 1 - \text{FPR}$$

**precision or positive predictive value (PPV)**

$$\text{PPV} = \frac{\text{TP}}{\text{TP} + \text{FP}} = 1 - \text{FDR}$$

**negative predictive value (NPV)**

$$\text{NPV} = \frac{\text{TN}}{\text{TN} + \text{FN}} = 1 - \text{FOR}$$

**miss rate or false negative rate (FNR)**

$$\text{FNR} = \frac{\text{FN}}{\text{P}} = \frac{\text{FN}}{\text{FN} + \text{TP}} = 1 - \text{TPR}$$

**fall-out or false positive rate (FPR)**

$$\text{FPR} = \frac{\text{FP}}{\text{N}} = \frac{\text{FP}}{\text{FP} + \text{TN}} = 1 - \text{TNR}$$

https://en.wikipedia.org/wiki/Confusion_matrix

# Evaluation metrics (2 classes)

**accuracy (ACC)**

$$\text{ACC} = \frac{\text{TP} + \text{TN}}{\text{P} + \text{N}} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

**F1 score**

is the harmonic mean of precision and sensitivity

$$F_1 = 2 \cdot \frac{\text{PPV} \cdot \text{TPR}}{\text{PPV} + \text{TPR}} = \frac{2\text{TP}}{2\text{TP} + \text{FP} + \text{FN}}$$

**Matthews correlation coefficient (MCC)**

$$\text{MCC} = \frac{\text{TP} \times \text{TN} - \text{FP} \times \text{FN}}{\sqrt{(\text{TP} + \text{FP})(\text{TP} + \text{FN})(\text{TN} + \text{FP})(\text{TN} + \text{FN})}}$$
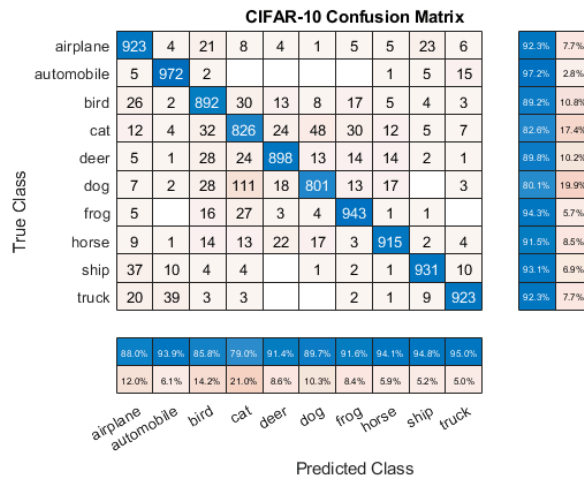
https://en.wikipedia.org/wiki/Confusion_matrix

# Confusion matrix



https://scikit-learn.org/stable/auto_examples/model_selection/plot_confusion_matrix.html

# Confusion matrix



https://www.mathworks.com/help/deeplearning/ref/confusionchart.html

# Train, Validation, Test

# Train, validation, and test sets

TRAIN SET

TRAIN SET | TEST SET

TRAIN SET | VALID SET | TEST SET

---

Original set

Training set | Test set

Training set | Validation set | Test set

Training, tuning, and evaluation

Machine learning algorithm

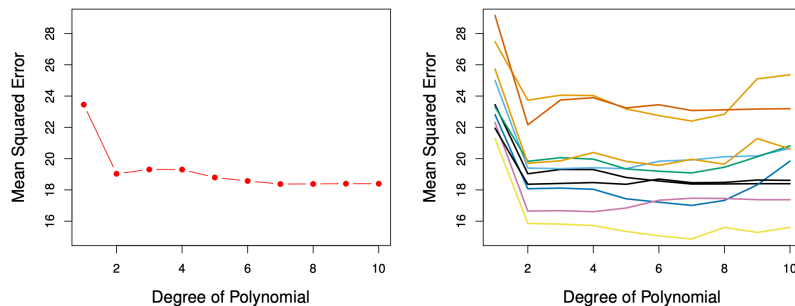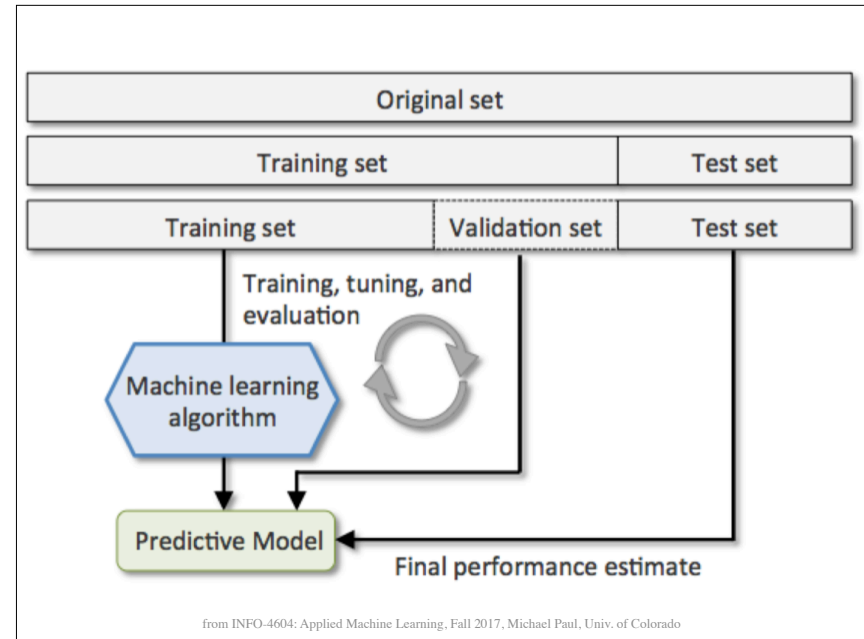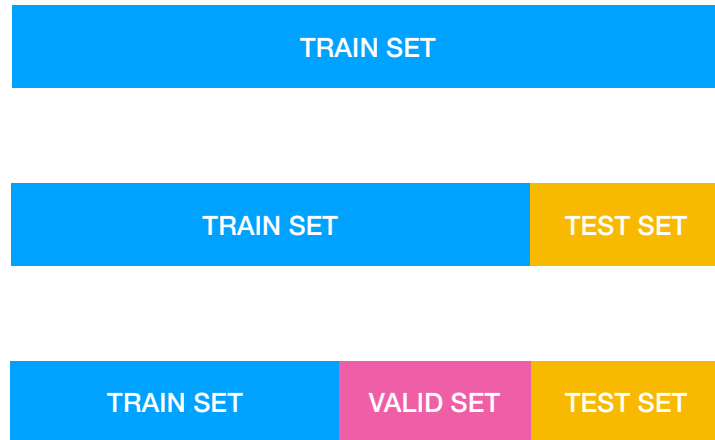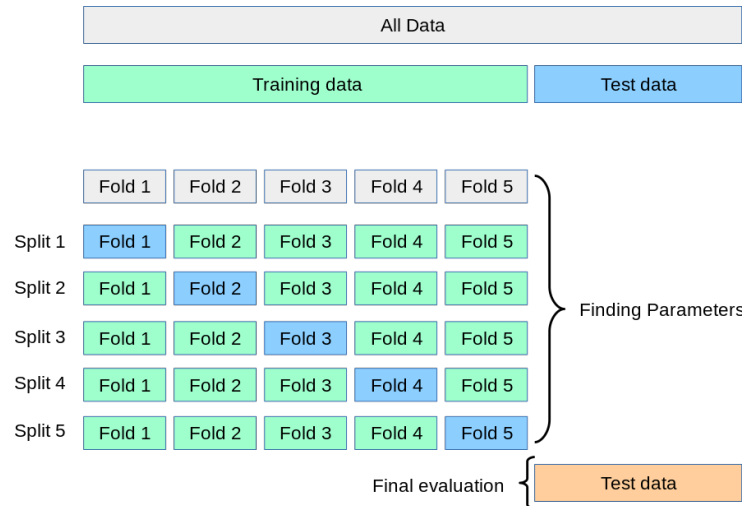Predictive Model

Final performance estimate

---



**FIGURE 5.2.** *The validation set approach was used on the* Auto *data set in order to estimate the test error that results from predicting* mpg *using polynomial functions of* horsepower. *Left: Validation error estimates for a single split into training and validation data sets. Right: The validation method was repeated ten times, each time using a different random split of the observations into a training set and a validation set. This illustrates the variability in the estimated test MSE that results from this approach.*
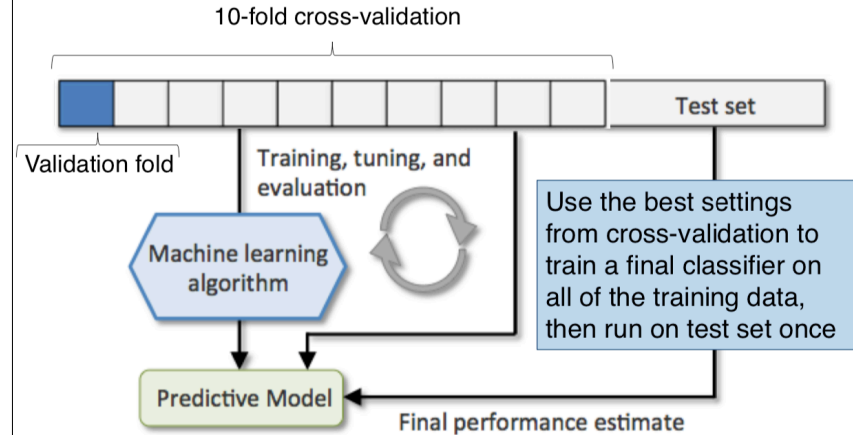
---

# Cross Validation
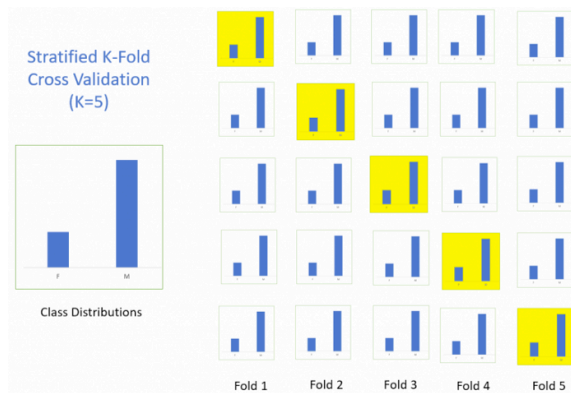
# What is k-fold Cross Validation?



https://scikit-learn.org/stable/modules/cross_validation.html

# Using Cross-Validation



from INFO-4604: Applied Machine Learning, Fall 2017, Michael Paul, Univ. of Colorado

# Stratified cross validation



**Stratified cross validation** aims at having the same class distribution within each fold

https://towardsdatascience.com/cross-validation-explained-evaluating-estimator-performance-e51e5430ff85

# Leave-One-Out CV

‣ Special case of CV when **k = n**

‣ Can be expensive for large **n**

*n = 8*



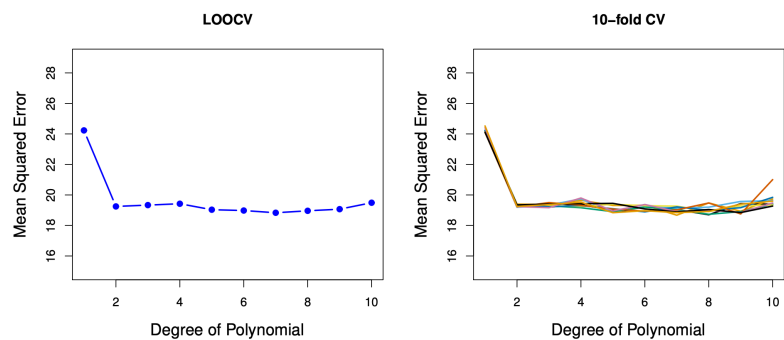https://en.wikipedia.org/wiki/Cross-validation_(statistics)

**FIGURE 5.4.** *Cross-validation was used on the* `Auto` *data set in order to es-timate the test error that results from predicting* `mpg` *using polynomial functions of* `horsepower`. Left: *The LOOCV error curve.* Right: 10-*fold CV was run nine separate times, each with a different random split of the data into ten parts. The figure shows the nine slightly different CV error curves.*