

# Introduction

CV\_M1\_L01

# **Computer Science 5970-008**

## **Machine Learning Practice**

**Andrew H. Fagg**

**Symbiotic Computing Laboratory  
School of Computer Science**



*The UNIVERSITY of OKLAHOMA*

# Constructing Models

- Start with observations (data) drawn from the world
  - Motion of an object, force applied to that object
- Models relate different types of observations to one-another

$$F = m \times a$$

# What Makes a Good Model?

A good model:

- Is simple
- Explains the observations that have already been made
- Is predictive of future observations

# Machine Learning

# Machine Learning

Fundamentally: ML is about using data to automatically construct a model. We would like:

- The model to produce meaningful output given novel situations
- The model to give us insights into the problem



- End section



# Example: Brain-Machine Interfaces

CV\_M1\_L02

# Example: Brain-Machine Interfaces

# Example: Brain-Machine Interfaces

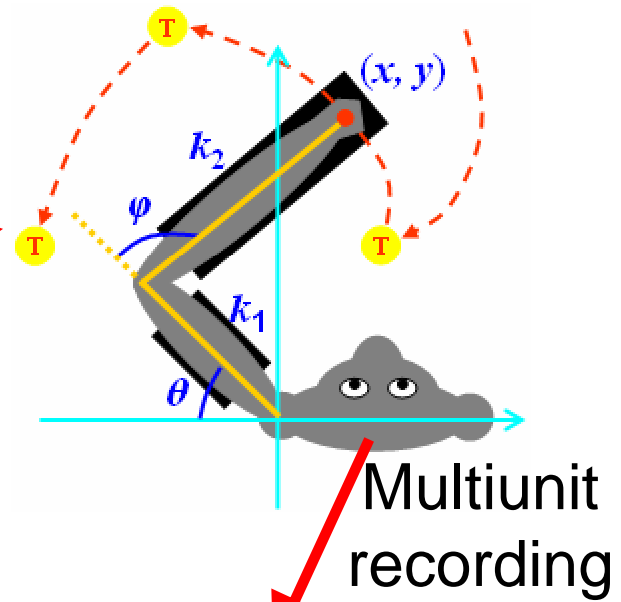
- Goal: to develop a direct connection from the brain to an advanced prosthetic device
- Approach:
  - Electrodes in the primary motor cortex “listen” to individual neurons or small clusters of neurons
  - Cortical neurons communicate by emitting sequences of pulses (“spikes” or “action potentials”) at different rates
  - Use a model to decode these pulses in terms of the intent to move the arm

# Brain-Machine Interfaces

Estimate of  
intended  
movement

Command  
prosthetic arm

Predictive  
model

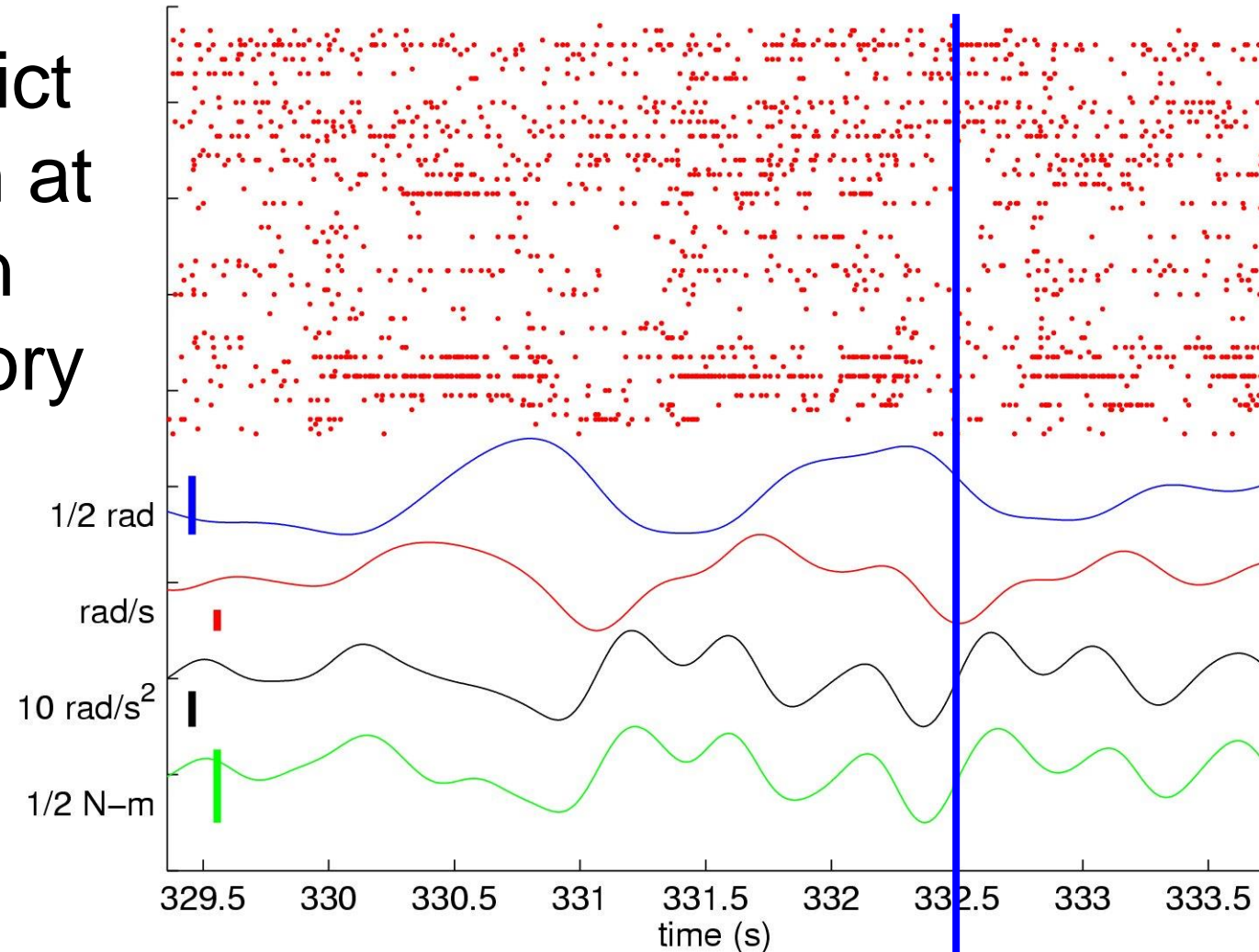


The UNIVERSITY of OKLAHOMA

In collaboration with Nicholas G. Hatsopoulos and Lee E. Miller

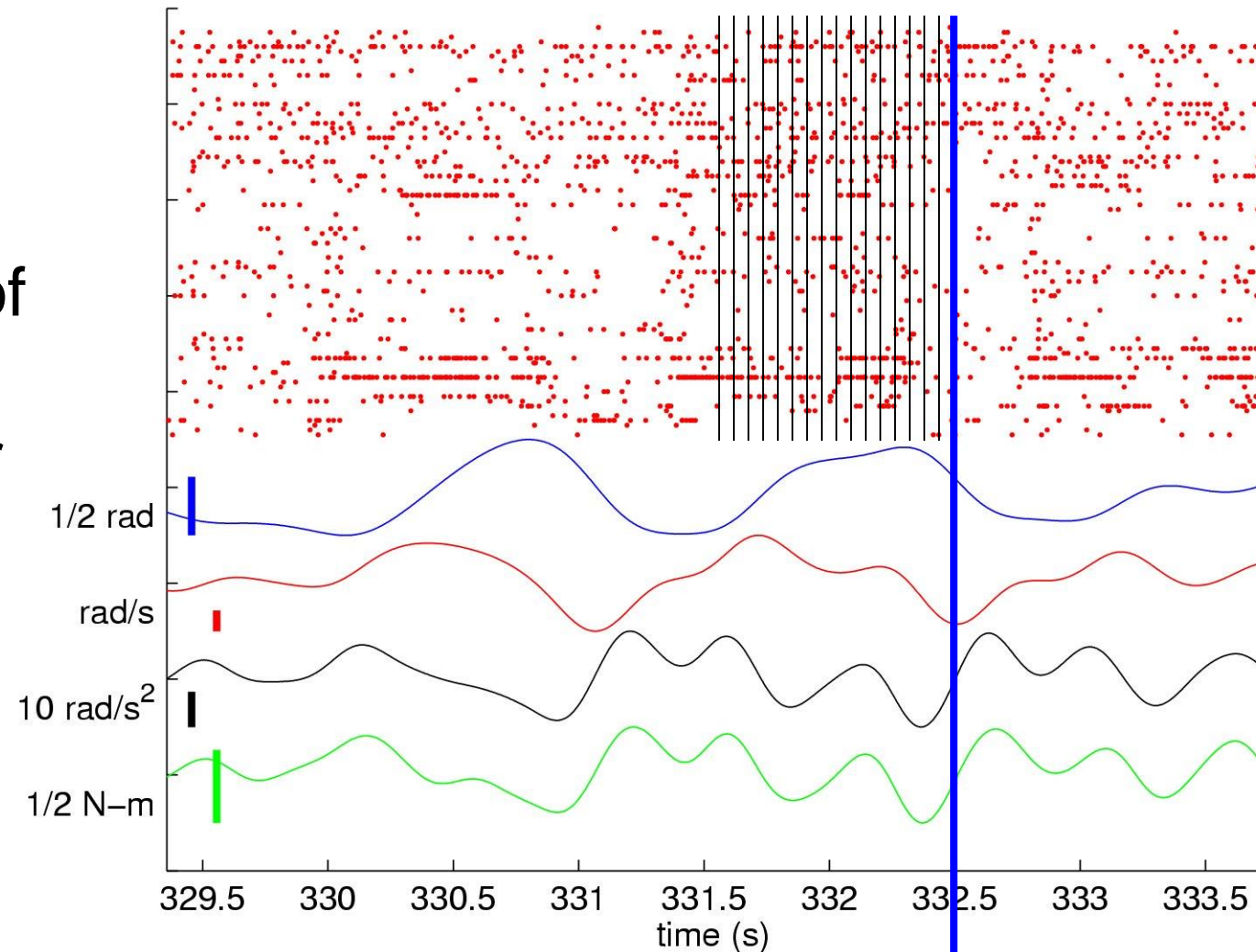
# Decoding Arm State

Want to predict  
arm motion at  
time  $t$  given  
recent history  
of spiking  
behavior



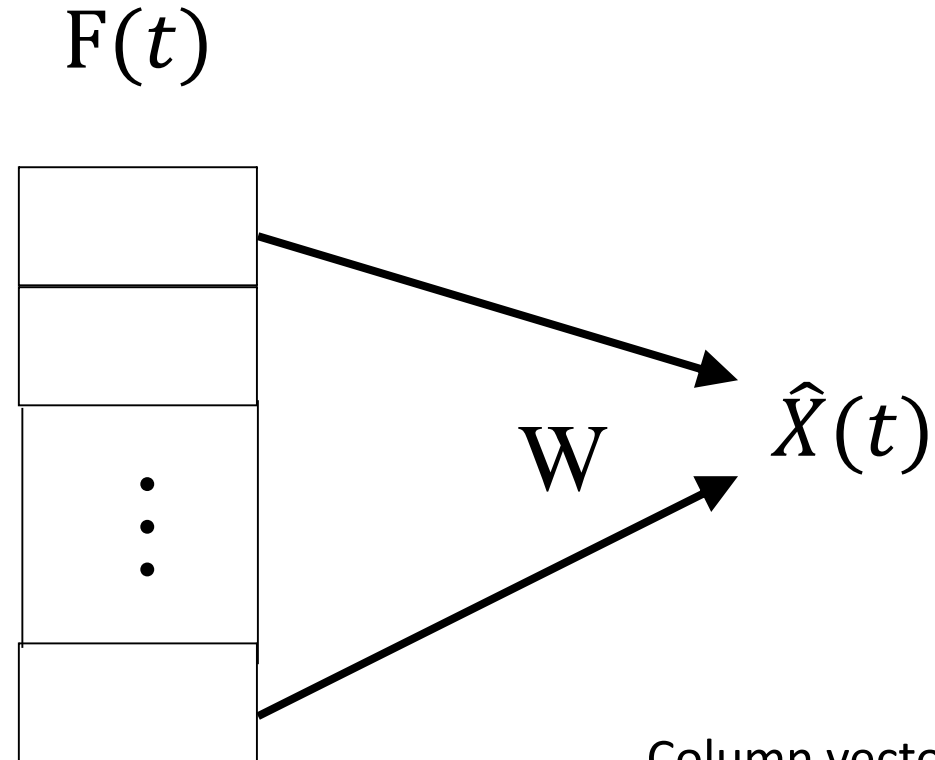
# Decoding Arm State

50ms bins: 20  
descriptors of  
neural  
activation for  
each cell



# Linear Model

Each feature  
( $F_i$ ) is a count  
of spikes by a  
neuron for a  
50 ms bin



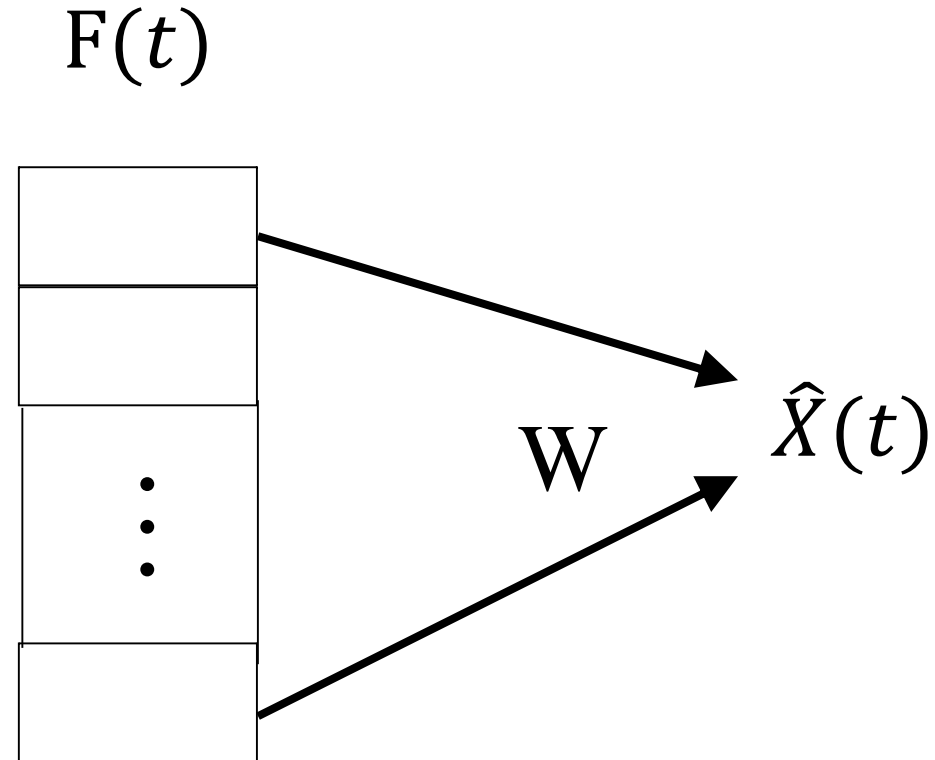
Column vector encoding  
spike counts for  $N$  cells at  
 $T$  taps up to time  $t$

$$\hat{X}(t) = g_W(F(t)) = W^T F(t)$$



# Linear Model

Each feature  
( $F_i$ ) is a count  
of spikes by a  
neuron for a  
50 ms bin



$$\hat{X}(t) = g_W(F(t)) = W^T F(t) = \sum_{i=0}^{N-1} w_i \times F_i(t)$$



# Training a Linear Model

Gathering the data:

- Monkey makes a sequence of reaches
- Simultaneously observe the movement of the monkey's arm and the neural activity
- This provides a set of example input / output examples for our model

# Training a Linear Model

- Linear model works well for this problem:

$$\hat{X}(t) = \sum_{i=0}^{N-1} w_i \times F_i(t)$$

- Cost function:

$$E = \frac{1}{n} \sum_t (X(t) - \hat{X}(t))^2$$

- Learning algorithm: pick the  $w_i$ 's so as to minimize  $E$

# Using Our Model

Given new observations of neural spiking patterns, we can:

- Predict how the monkey will move her arm
- Use these predictions to drive the motion of the prosthesis

- End section

# Machine Learning Taxonomy

- CV\_M1\_L03

# Machine Learning Taxonomy

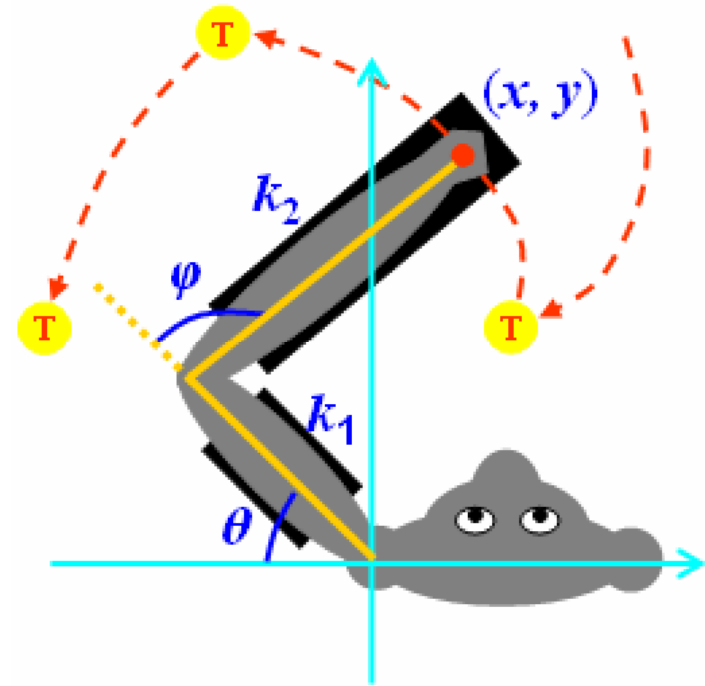
# Classes of Models

Defined by the data type of the output. Very broadly:

- Continuous output: regression-type models
- Categorical output: classifier models

# Regression-Type Models

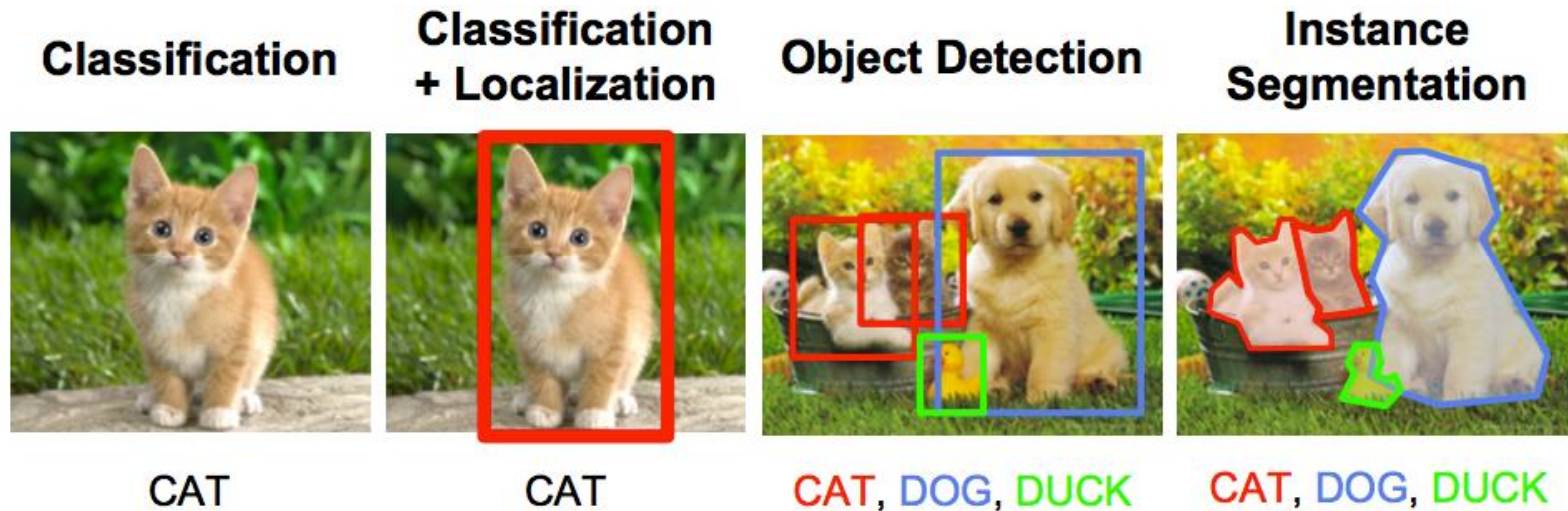
- Continuous output
- In our brain-machine interface example: what velocity should the arm be moving at given the recent history of neural activity patterns?





# Classification-Type Models

- Classification: given an input, which one of several classes does the input belong to?
- Can be crisp (choose exactly one class)
- Or can be probabilistic (each class is assigned a probability)



# Classes of Machine Learning Problems

What information is provide at the time of training?

# Classes of Machine Learning Problems

Supervised learning:

- Training set contains input / output (labels) pairs
- Outputs could be continuous, probabilistic or categorical

# Classes of Machine Learning Problems

Unsupervised learning:

- The training set contains only inputs
- Fundamental question: what is the structure of these inputs?
  - A common case: algorithm assigns categorical labels to each of the inputs (this is clustering)
  - But we can also ask continuous questions. For example: are there linear or nonlinear manifolds that the data live on?

- IPAD HERE

# Classes of Machine Learning Problems

Semi-Supervised learning:

- Part of the training set contains input / output pairs
- The rest of the training set contains only inputs
- Using all of the data can yield a better model than if we only used the labeled data

# Classes of Machine Learning Problems

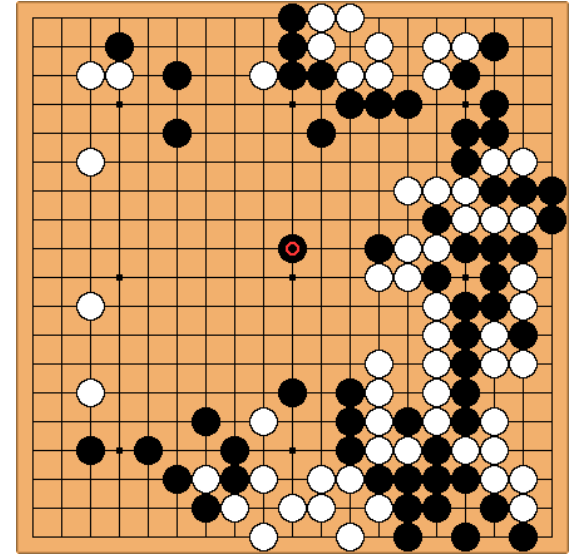
Reinforcement learning:

- Different than direct prediction or classification: RL is about taking sequences of actions in some environment
- At each step:
  - In response to an input, the model (agent) produces some action
  - The feedback signal is an evaluation of the results of this and previous actions

# Classes of Machine Learning Problems

Reinforcement learning:

- Common reward types:
  - How much time did it take to execute an action?
  - How much energy did an action take?
  - Did the agent win the game?
- Learning problem: for a given input, what is the action that maximizes the expected sum of rewards over time?







- End section

# Practical Challenges

CV\_M1\_L04

# Practical Challenges

# Practical Challenges

## Modeling Choices:

- Right model and learning algorithm
  - Worry about computational complexity in training or querying a model
- Hyper-parameters
- Selecting a data set to train from
  - Data can be expensive to collect
  - Different algorithms require different amounts of data

# Practical Challenges

## Overfitting

- Model matches the training data set well, but does not perform well on independent data

- INSERT IPAD DRAWING HERE

# Practical Challenges

## Overfitting

- Model matches the training data set well, but does not perform well on independent data
- How do we detect this?
- How do we mitigate this?
  - Some algorithms will handle this automatically
  - In some cases, we have to be careful about how we choose our training set



# Practical Challenges

Comparing models and algorithms

- Measuring performance of a model
- Performance is inherently a random variable
  - Must acknowledge this when we are comparing two models
  - This implies that comparison is an empirical process
  - Also must acknowledge this issue when selecting hyper-parameters



- End section

# Course Topics

- CV\_M1\_L05

# Course Topics

# Course Topics

## Preliminaries:

- Python
- Jupyter
- Pandas
- Numpy
- Scikit-Learn
- Python best practices

# Course Topics

- Classifiers
  - Logistic regression, support vector machines, decision trees
  - Feature importance
- Regression
  - Linear and non-linear
  - Polynomial / kernel regression, support vector regression and decision tree regression
- Decision Trees: ensemble methods and random forests

# Course Topics

## Unsupervised Methods

- Principal component analysis
  - Local linear embeddings
  - Multidimensional scaling
  - ISomap
- 
- Clustering: K-Means, Mixture Models



# Course Topics

## Tuning Models

- Detecting and mitigating overfitting
- Choosing hyperparameters
- Comparing algorithm types in a statistically sound way



# End section

# Course Mechanics

- CV\_M1\_L06

# Course Mechanics

# Course Delivery

- All lecture material is on-line via Canvas
  - Will release the videos and homework assignments at the beginning of the week
- Our lecture time will be used for my office hours:
  - T/Th 9:00 – 10:15am in Sarkey's A0133
  - Will also livestream via Canvas if requested
- TA office hours:
  - To be chosen

# Computing Environment

- All homework assignments will be done in Python
- We are providing a computing server for these assignments (more details to come)
  - Your primary interface is through **Jupyter Lab**
  - Packages pre-installed; data and code skeletons automatically available
  - You are also welcome to work on your local machine, if you wish

# What I am assuming about you...

- Programming background:
  - Experience with object-oriented programming
  - Python is not a necessary prerequisite, but is a bonus
- Statistical Methods:
  - Linear regression
  - Hypothesis testing



# Resources

- Course web page:  
<http://www.cs.ou.edu/~fagg/classes/aml>
- Canvas: grade book, announcements, discussion board, office hours, videos
- Text: Aurélien Géron (2017) **Hands-On Machine Learning with Scikit-Learn and TensorFlow** (Concepts, Tools, and Techniques to Build Intelligent Systems) ISBN-13: 978-1491962299, O'Reilly Media
- Web resources: documentation, tutorials, papers (linked from the schedule or announced on Canvas)

# Grading

## Homework

- 12 assignments (+ one test assignment)
- Explore different ML methods and data sets
- Criteria:
  - Success in solving the problem
  - Cleanliness of the code (yes, we expect documentation)

No final exam or end-of-semester project

# Proper Academic Conduct

Homework assignments are to be done on your own

- No communication of solutions in any form with anyone other than the instructor or TA
- Do not copy code off the net
- General communication or drawing inspiration off of the net is okay

# Keys to Success

- Stay on top of lectures and homework assignments
- Learn to read the documentation
- Most assignments will not be doable in the day before the deadline. Start early
- The net is filled with lots of advice about how to do things
  - Much of the advice is poor or down-right wrong
  - Even when the advice is correct, you should still be able to write your own code
- Ask plenty of questions



- End section

# For Next Time

- For today: chapter 1
- Next time: start of chapter 2
- We will get you started on python and numpy





