

# Machine Learning Overview

Yale

CPSC/AMTH 663





1. What is machine learning?
2. Supervised Learning
  - Regression
  - Complexity, overfitting, and underfitting
  - Regularization
  - Classification
  - Validation techniques
  - Anomaly detection
3. Unsupervised Learning
  - Principal components analysis (PCA)
  - Manifold learning
  - Clustering



THIS IS YOUR MACHINE LEARNING SYSTEM?

| YUP! YOU POUR THE DATA INTO THIS BIG  
PILE OF LINEAR ALGEBRA, THEN COLLECT  
THE ANSWERS ON THE OTHER SIDE.

| WHAT IF THE ANSWERS ARE WRONG?

| JUST STIR THE PILE UNTIL  
THEY START LOOKING RIGHT.





- Machine learning algorithm = algorithm that can learn from data
- What do we mean by “learn from data”?
- “A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ .” – Mitchell (1997)
- What are  $E$ ,  $T$ , and  $P$ ?
  - It depends on the problem!



- Machine learning helps us to solve problems that are too difficult for fixed programs
- Learning is the process of attaining the ability to perform the task
- Tasks are typically described in terms of how to process an ***example***
  - An example is a collection of measured ***features***
  - Often represented as a vector  $x \in \mathbb{R}^d$  of features
  - E.g. an example includes a single image with the features comprised of the pixels in the image



- Classification
    - Output: a function  $f: \mathbb{R}^d \rightarrow \{1, \dots, k\}$
  - Regression
    - Output: a function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$
  - Transcription
  - Machine translation
  - Anomaly detection
  - Synthesis and sampling
  - Imputation of missing values
  - Denoising
  - Density estimation



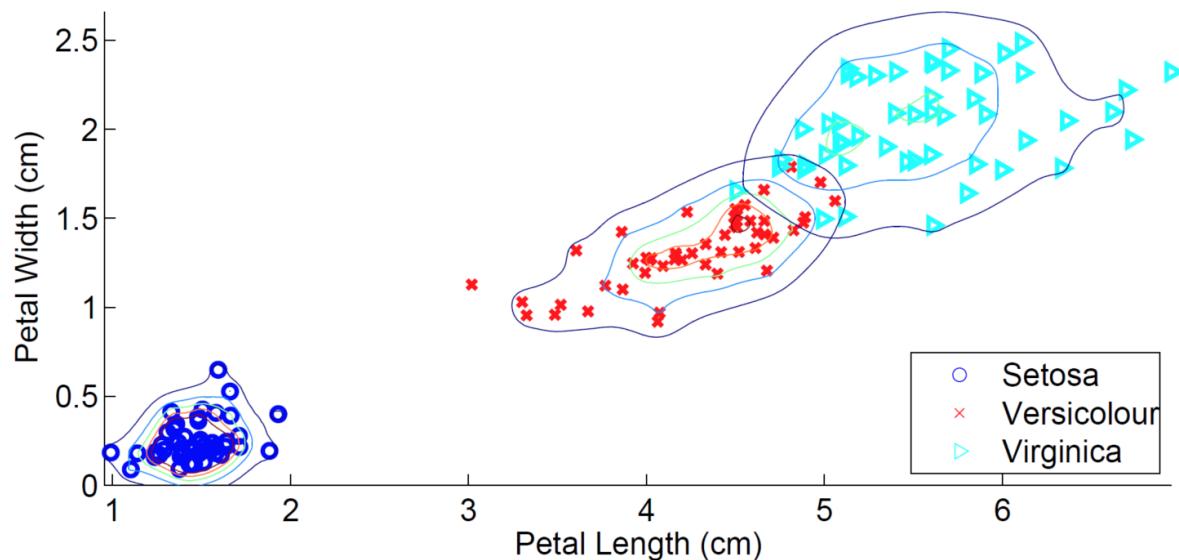
			
"man in black shirt is playing guitar."	"construction worker in orange safety vest is working on road."	"two young girls are playing with lego toy."	"boy is doing backflip on wakeboard."
			
"girl in pink dress is jumping in air."	"black and white dog jumps over bar."	"young girl in pink shirt is swinging on swing."	"man in blue wetsuit is surfing on wave."
			
"little girl is eating piece of cake."	"baseball player is throwing ball in game."	"woman is holding bunch of bananas."	"black cat is sitting on top of suitcase."



- Need a quantitative measure to evaluate a machine learning algorithm
- Provides a measure to guide the selection of an algorithm
- Examples:
  - **Accuracy** (or equivalently, the ***error rate***) in a classification problem
  - Prediction error
- Performance is typically ***estimated***/evaluated based on a *test set* of data (different from training data)
- Choosing a performance measure isn't always straightforward
  - Transcription
  - Regression

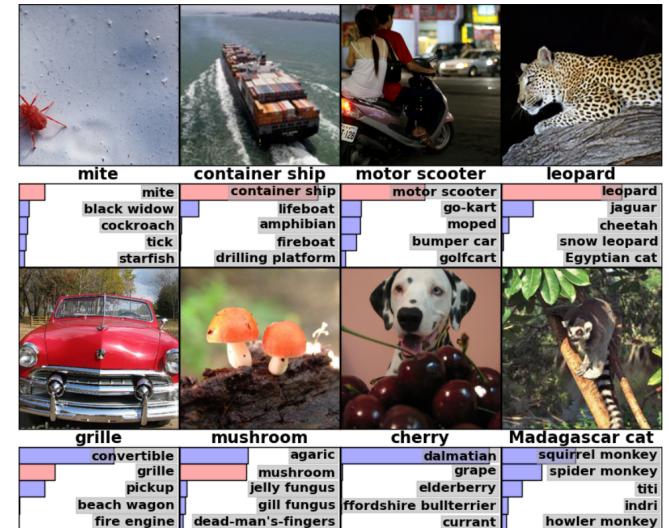
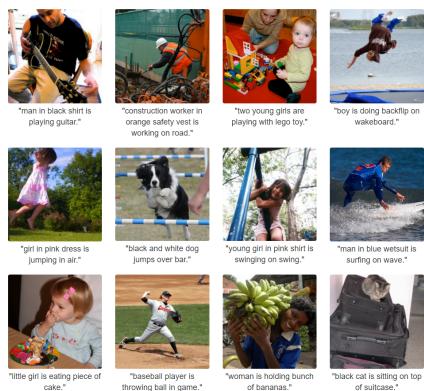
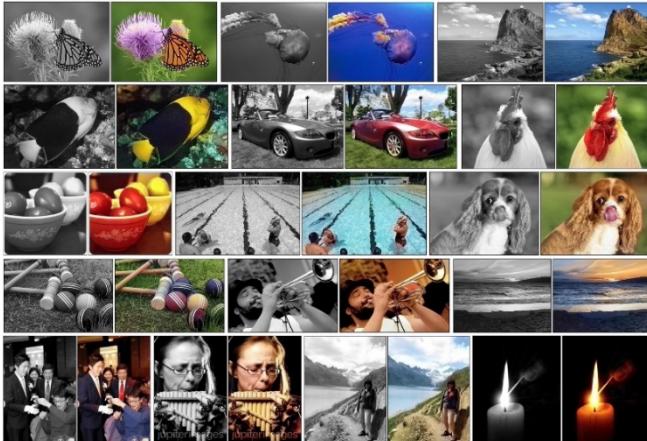


- Most algorithms experience an entire **dataset** consisting of many examples or **data points**
- Example: Iris data set (Fisher, 1936)
  - Four features: sepal length, sepal width, petal length, petal width
  - 3 species/classes, 50 data points per class





- Goal is to build a program that learns from experience  $E$  to do well at task  $T$  as measured by performance measure  $P$



# Supervised Learning



- Dataset contains features and a *label* or target
    - E.g. classify iris plants into the three different species
  - Mathematically, supervised learning observes examples of a random vector  $x$  and an associated label  $y$ 
    - Goal is to learn to predict  $y$  from  $x$
  - Term comes from view of the label  $y$  coming from an instructor or teacher

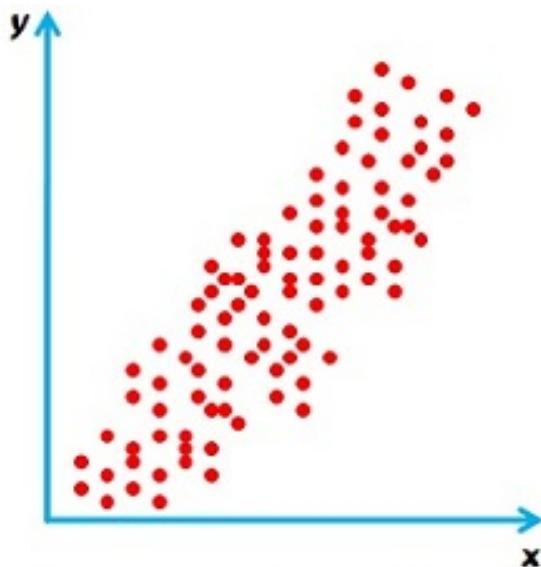
0000000000  
1111111111  
2222222222  
3333333333  
4444444444  
555555535555  
666666666666  
777777777777  
888888888888  
999999999999



- Input data space  $\mathcal{X}$
- Output (label) space  $\mathcal{Y}$
- Unknown function  $f: \mathcal{X} \rightarrow \mathcal{Y}$
- Dataset  $\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\}$  with  $x_i \in \mathcal{X}$  and  $y_i \in \mathcal{Y}$
- Finite  $\mathcal{Y} \Rightarrow$  classification
- Continuous  $\mathcal{Y} \Rightarrow$  regression



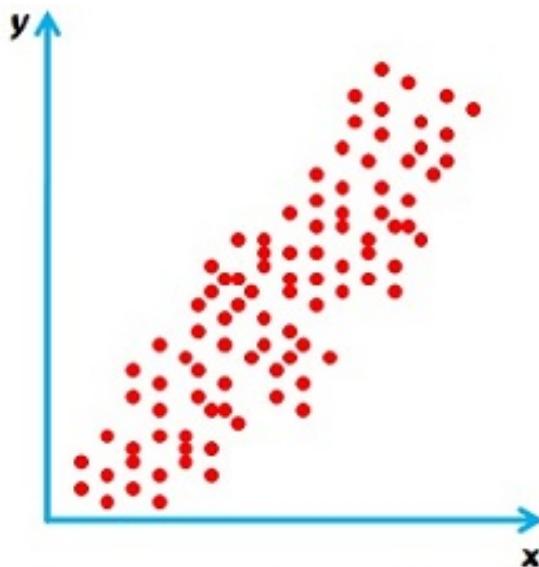
- We have a set of  $n$  observations  $(\mathbf{x}_i, y_i)$  with  $y_i \in \mathbb{R}$
- Goal is to learn a function that predicts label  $y$  from sample  $\mathbf{x}$



Surbhi S, 2016



1. Choose a *model class* of functions
2. Choose a performance measure to guide the selection of a function from the model class



Surbhi S, 2016

- We'll begin with a simple model class: linear functions



- We want to fit a linear function to dataset  $\mathcal{D} = \{(x_1, y_1), \dots (x_n, y_n)\}$
- $\hat{y} = \mathbf{w} \cdot \mathbf{x} = \mathbf{w}^T \mathbf{x}$ 
  - $\mathbf{w} \in \mathbb{R}^d$  is a vector of parameters
  - Line if  $d = 1$
  - Plane if  $d = 2$
  - Hyperplane in general  $d$
- How do we determine  $\mathbf{w}$ ?
- Based on the performance measure  $P$



- Labels are in  $\mathcal{Y}$ 
  - Discrete (classification)
  - Continuous (regression)
- ***Loss function***  $L: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$
- $L$  maps decisions/predictions to a cost.
  - $L(\hat{y}, y)$  is the penalty for predicting  $\hat{y}$  when the true label is  $y$
- Standard choice for regression is ***mean squared error*** (MSE)
  - $L(\hat{y}, y) = (\hat{y} - y)^2$
- Absolute loss:  $L(\hat{y}, y) = |\hat{y} - y|$



- Parametric function  $f(\mathbf{x}; \mathbf{w})$ 
  - E.g.  $f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x}$  for regression
- The *empirical loss* on a dataset  $\mathcal{D}$  is

$$L(\mathbf{w}; \mathcal{D}) = \frac{1}{n} \sum_i L(f(\mathbf{x}_i; \mathbf{w}), y_i)$$

- Least squares minimizes the empirical MSE
- We care about predicting labels for *new* samples
- Does empirical loss minimization help with that?



- Assumption: sample and label pairs  $(\mathbf{x}, y)$  are drawn from the unknown joint distribution  $p(\mathbf{x}, y)$
- Data are sampled i.i.d.
- Empirical loss measured on training data

$$L(\mathbf{w}; \mathcal{D}) = \frac{1}{n} \sum_i L(f(\mathbf{x}_i; \mathbf{w}), y_i)$$

- Goal: minimize the ***expected loss*** (aka ***risk***)

$$R(\mathbf{w}) = \mathbb{E}[L(f(\mathbf{x}; \mathbf{w}), y)]$$



- Empirical loss

$$L(\mathbf{w}; \mathcal{D}) = \frac{1}{n} \sum_i L(f(\mathbf{x}_i; \mathbf{w}), y_i)$$

- Risk

$$R(\mathbf{w}) = \mathbb{E}[L(f(\mathbf{x}; \mathbf{w}), y)]$$

- If training set is representative of  $p(\mathbf{x}, y)$ , then empirical loss is a suitable proxy for the risk
- Essentially, we estimate  $p(\mathbf{x}, y)$  and  $R(\mathbf{w})$  from the empirical distribution of the data



1. Choose a model class  $\mathcal{H}$  of functions  $f: \mathcal{X} \rightarrow \mathcal{Y}$ 
  - Regression example: linear functions parameterized by  $\mathbf{w}$ :  
$$f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x}$$
2. Select the function in  $\mathcal{H}$  that minimizes the empirical risk
  - Linear regression example: minimize empirical squared loss:  
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2$$
  - How do we find  $\mathbf{w}^*$ ?
  - Use calculus



- Add an offset  $w_0$ :  $f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x} + w_0$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i + w_0 - y_i)^2$$

$$= \arg \min_{\mathbf{w}} L(\mathbf{w}; \mathcal{D})$$

- Set  $\frac{\partial L(\mathbf{w}; \mathcal{D})}{\partial w_i} = 0$  for each  $i$
- Results in  $d + 1$  equations and  $d + 1$  unknowns



- Switching to vector/matrix notation

- $X$  = data matrix with an additional column of 1's (for the offset)
- $\mathbf{y}$  = vector of labels
- Prediction:  $\hat{\mathbf{y}} = X\mathbf{w}$

- Resulting solution is

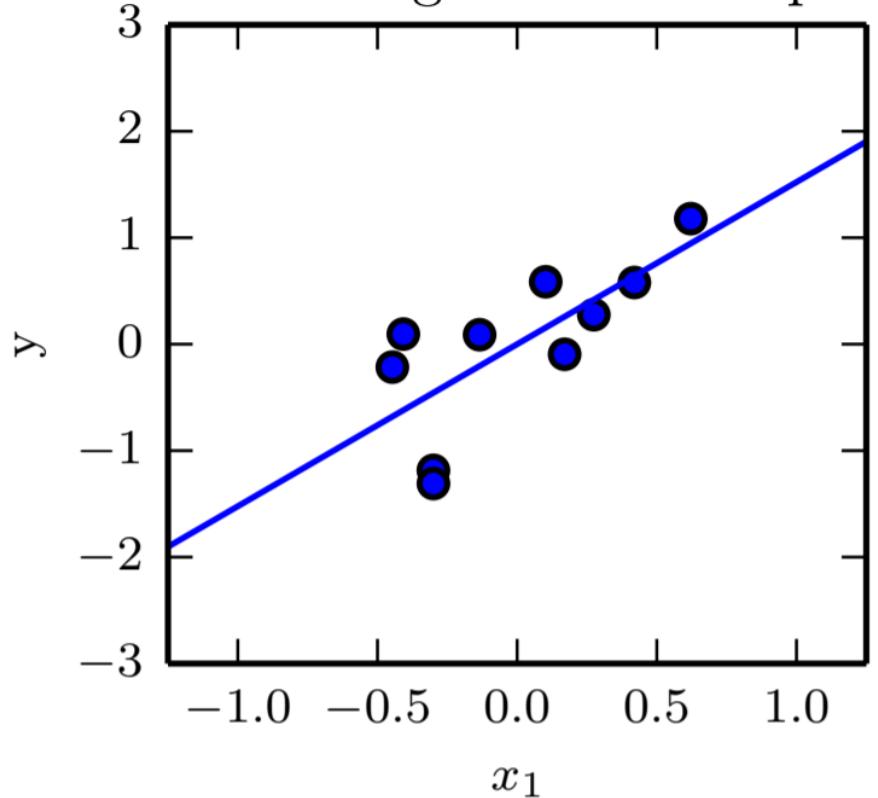
$$\mathbf{w}^* = (X^T X)^{-1} X^T \mathbf{y}$$

- Closed form solution!

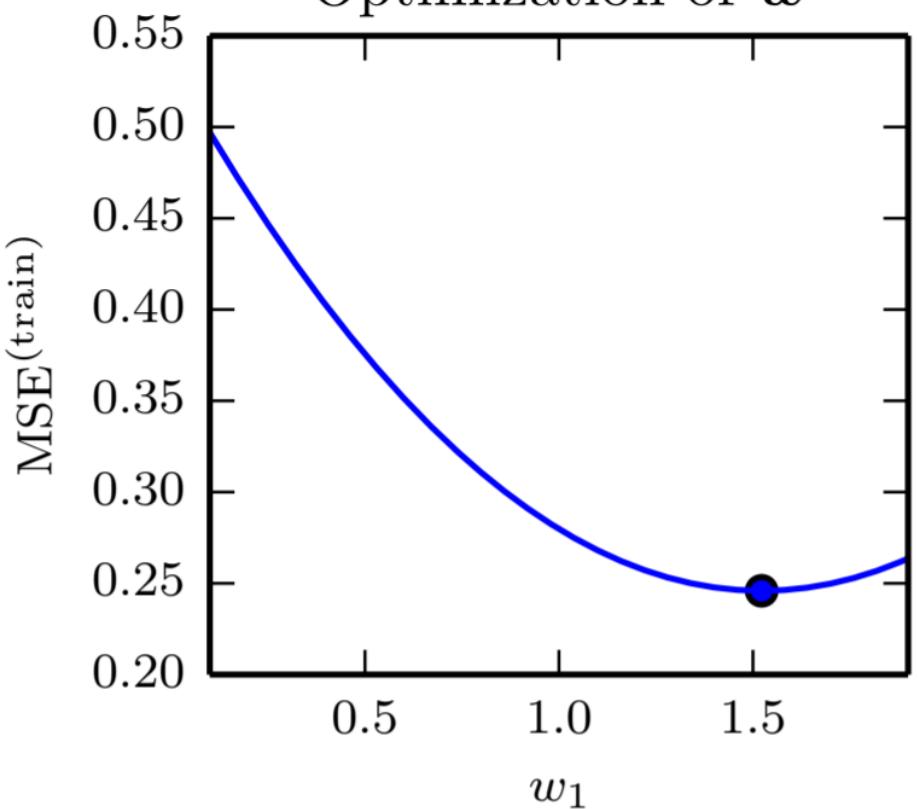
- $(X^T X)^{-1} X^T$  is the Moore-Penrose pseudoinverse of  $X$



Linear regression example



Optimization of  $w$

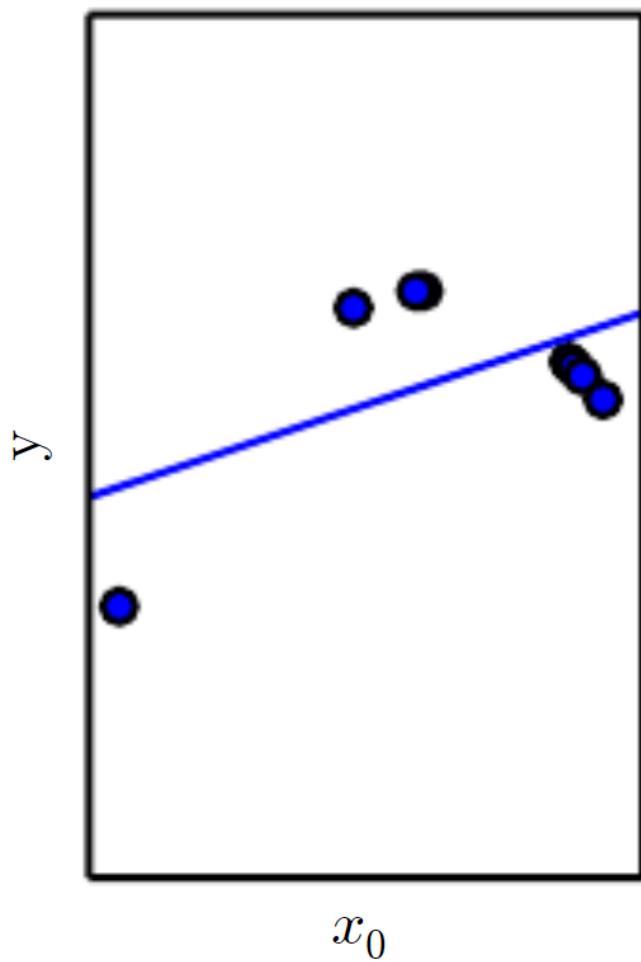




- Challenge of machine learning: algorithm must perform well on previously unseen inputs
  - Not just training inputs
- In other words, we want the algorithm to ***generalize*** well to unseen data
  - That is, we want the expected error rate on previously unseen inputs to be low
  - This error rate is called the ***generalization error*** or the ***test error***
- We also need to estimate the generalization error
  - Done typically by dividing the total amount of data into two sets: ***training data*** and ***test data***
    - More on how this is done specifically later
    - Model is trained on the training data and then tested on the test data
- Expected test error  $\geq$  expected training error



- Many processes are not well modeled by linear functions
  - Biological systems
  - Stock market prices
  - Periodic signals
- In these cases, linear functions can ***underfit*** the data
  - Training error is very high
  - Test error will also be high

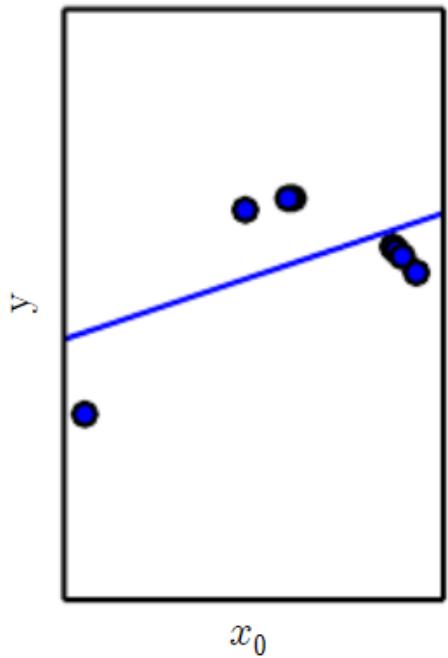




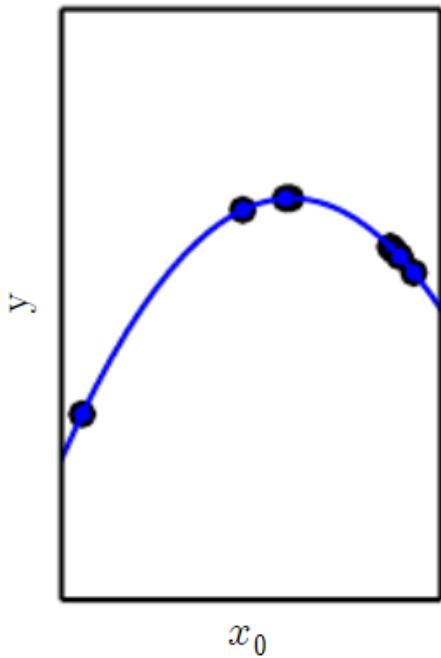
- Solution: increase ***complexity*** of the model
  - Model complexity is the # of independent parameters in the function model (i.e. “degrees of freedom”)
- We can increase the complexity of our linear function by adding polynomial terms
  - Example:  $d = 1$ 
$$f(x; \mathbf{w}) = \sum_{j=0}^m w_j x^j$$
  - Still linear in  $\mathbf{w} \Rightarrow$  we can use linear regression!
  - Can generalize further by using other nonlinear functions of  $x$  (e.g. log or exp)



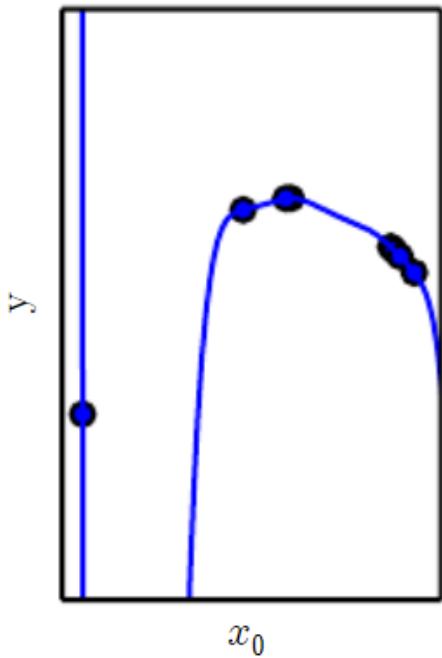
Underfitting



Appropriate capacity

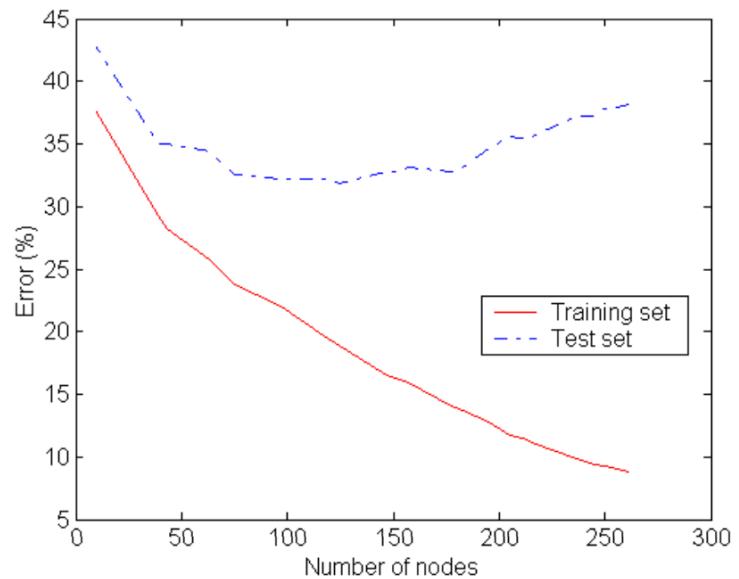


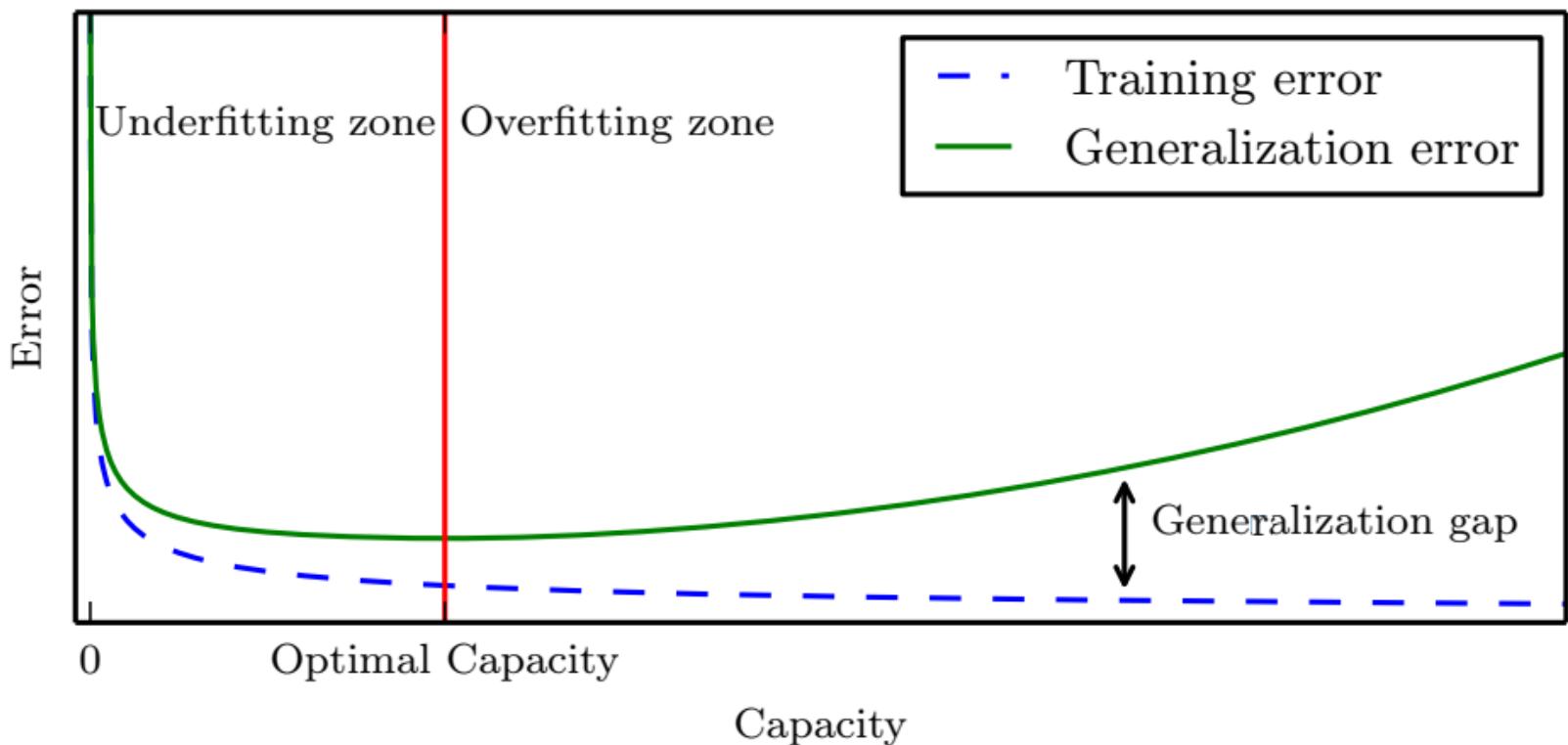
Overfitting





- Overfitting occurs when the gap between training error and test error is too large
  - The model is capturing unique characteristics of the training set
  - The model does not generalize well to new samples
  - Occurs when the complexity is too large
- Avoiding overfitting
  - Choose model complexity based on test error
  - Generally requires the use of another test set for evaluation







- Include a penalty on model complexity directly in the cost function:

$$\sum_i L(f(\mathbf{x}_i; \mathbf{w}), y_i) + \#\text{parameters}$$

- The addition of the penalty is known as ***regularization***
- Regularization: “any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error.” (Goodfellow et al., 2016)
  - Lots of different penalties can be included in regularization
  - Regularization is useful/important in deep learning (and machine learning in general)



- Ridge regression: penalize with L2 norm

$$\mathbf{w}^* = \arg \min \sum_i L(f(\mathbf{x}_i; \mathbf{w}), y_i) + \lambda \sum_{j=1}^m w_j^2$$

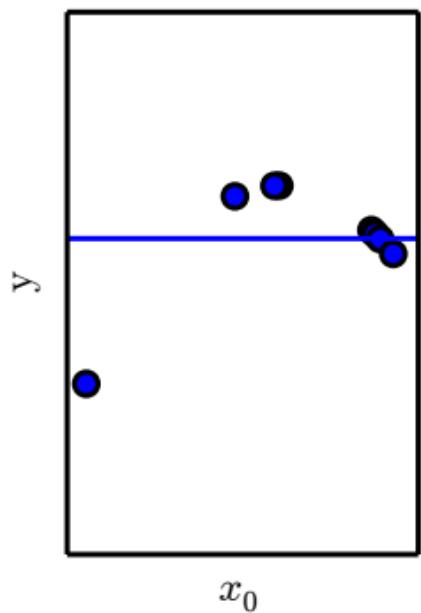
- Closed form solution exists  $\mathbf{w}^* = (\lambda I + X^T X)^{-1} X^T \mathbf{y}$
- LASSO regression: penalize with L1 norm

$$\mathbf{w}^* = \arg \min \sum_i L(f(\mathbf{x}_i; \mathbf{w}), y_i) + \lambda \sum_{j=1}^m |w_j|$$

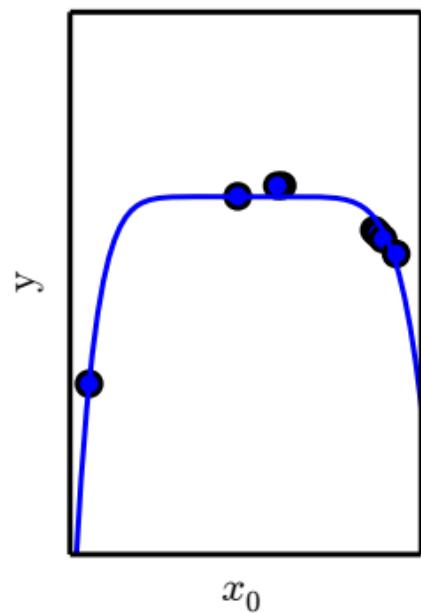
- No closed form solution but still convex (optimal solution can be found)



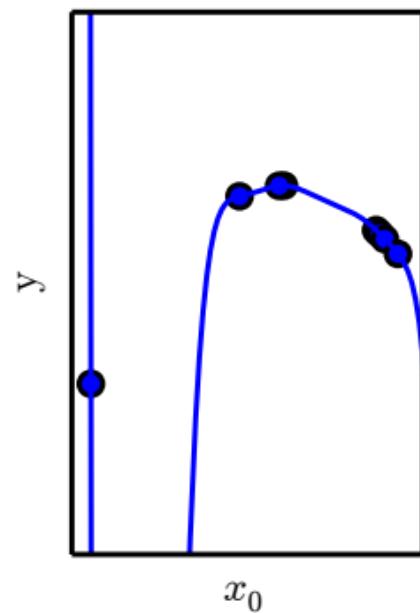
Underfitting  
(Excessive  $\lambda$ )



Appropriate weight decay  
(Medium  $\lambda$ )



Overfitting  
( $\lambda \rightarrow 0$ )



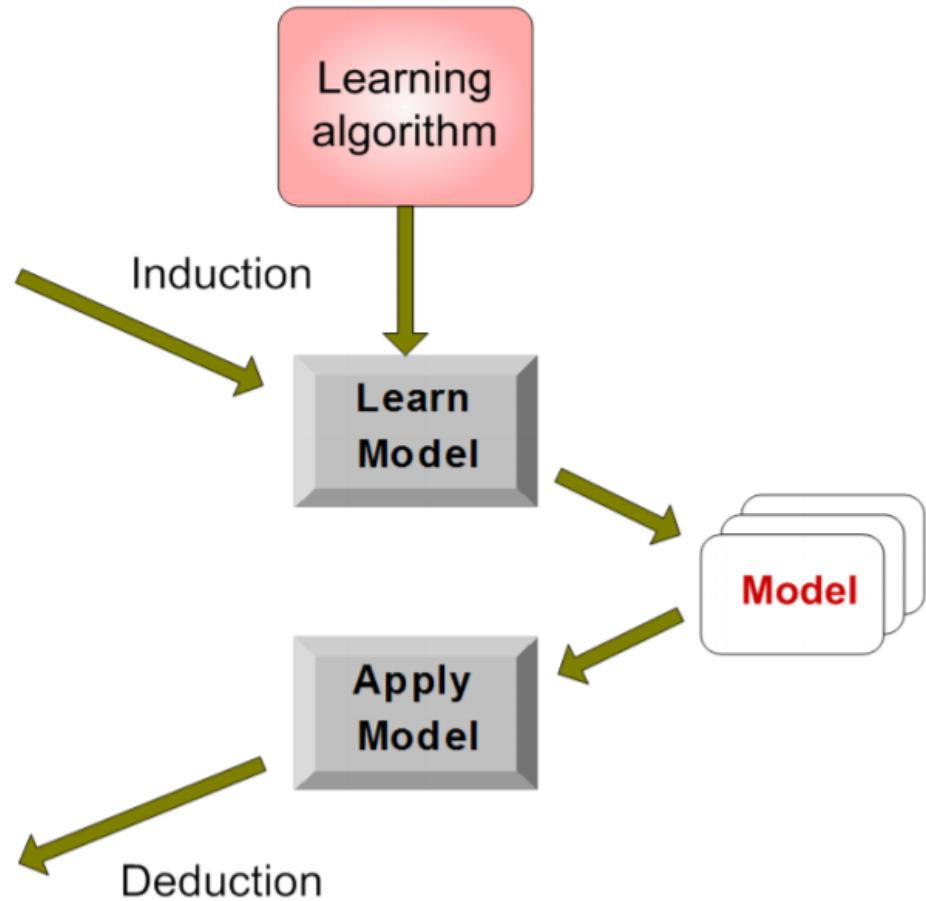


Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

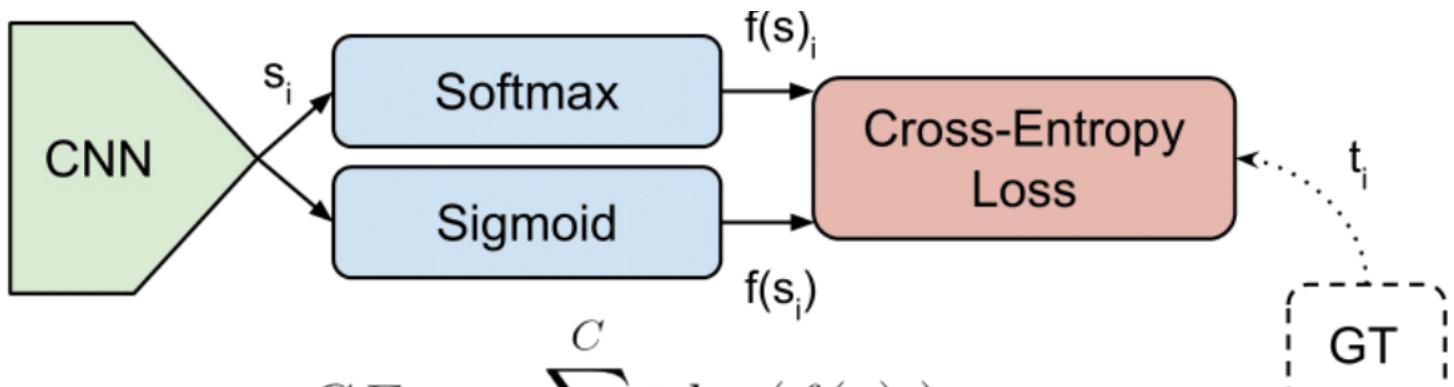
Test Set





- Let  $f: \mathcal{X} \rightarrow \mathcal{Y}$  be a classifier
- Common loss function is the 0/1 loss:

$$L(f(x), y) = \begin{cases} 0, & f(x) = y \\ 1, & f(x) \neq y \end{cases}$$



$$CE = - \sum_i^C t_i \log(f(s)_i)$$

$$CE = - \sum_{i=1}^{C'=2} t_i \log(f(s_i)) = -t_1 \log(f(s_1)) + (1 - t_1) \log(1 - f(s_1))$$



- For binary case, we get that

$$f(x) = 1 \text{ iff } \log \frac{p(y = 1|x)}{p(y = 0|x)} \geq 0$$

- We can model the decision boundary with a linear function:

$$\log \frac{p(y = 1|x)}{p(y = 0|x)} = 0 = w_0 + \mathbf{w}^T \mathbf{x}$$

- It can be shown that

$$p(y = 1|x) = \frac{1}{1 + \exp(-w_0 - \mathbf{w}^T \mathbf{x})}$$

- Known as the *logistic function*



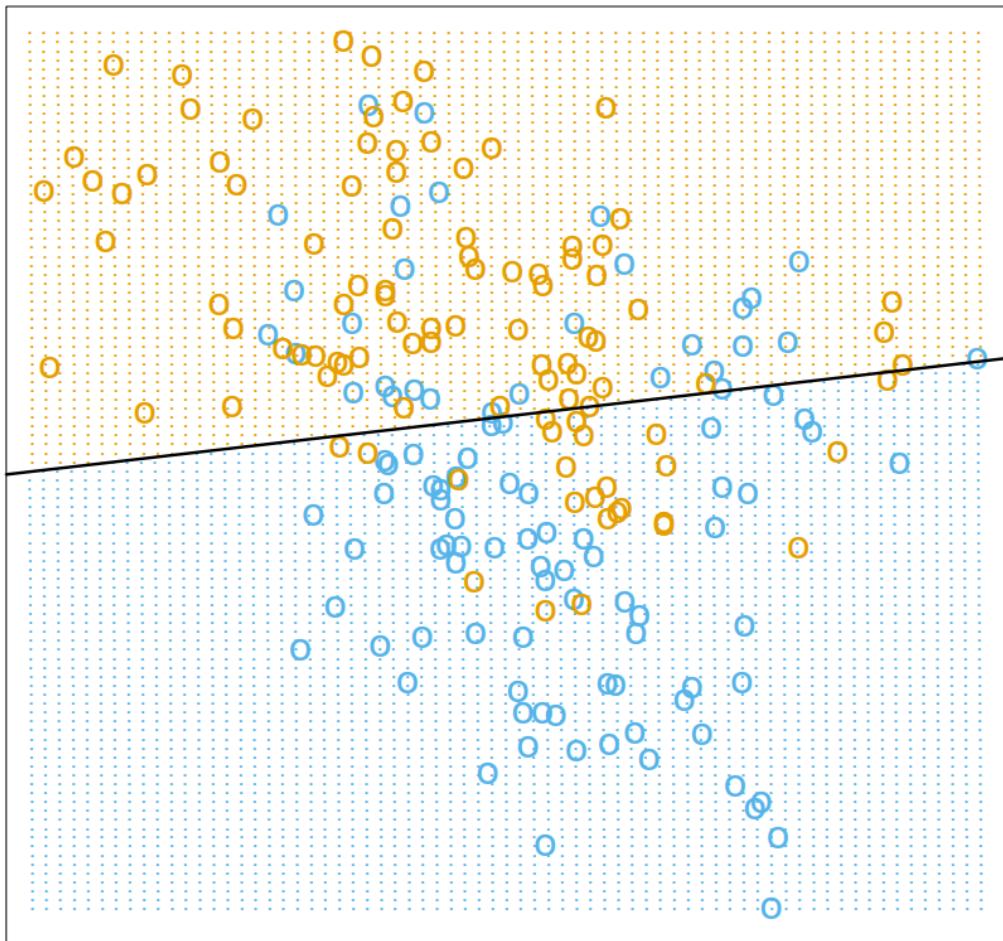
- Logistic function

$$p(y = 1|x) = \frac{1}{1 + \exp(-w_0 - \mathbf{w}^T \mathbf{x})}$$

- Optimal parameters can be found using ***Maximum Likelihood***
- Constructs a linear decision boundary with linear model
  - Can generalize to nonlinear boundaries
- This classifier is known as ***Logistic Regression***



Linear Regression of 0/1 Response



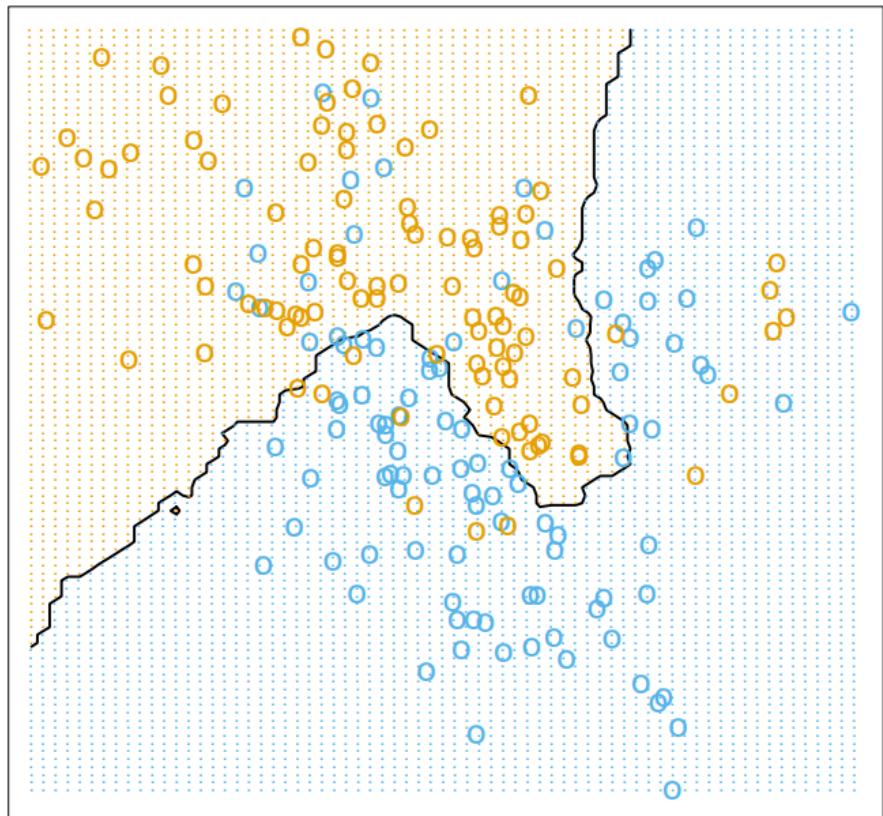
Hastie et al., 2009



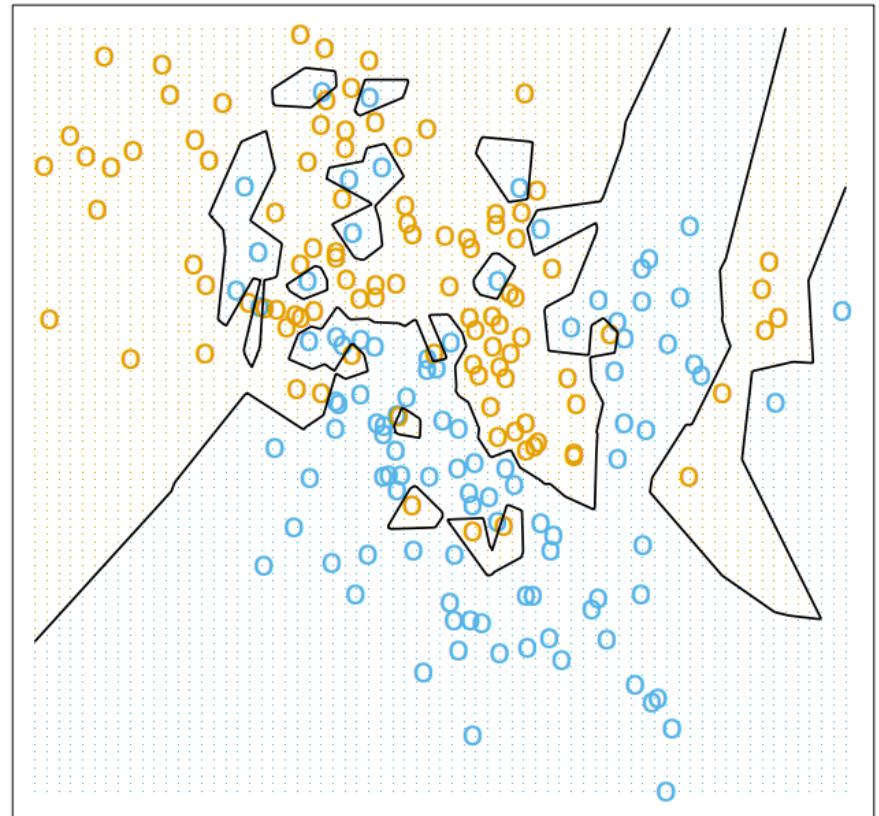
- No training required
- 1. Find the  $k$ -nearest neighbors of the test point  $x$  in the training data
  - Call this set  $\mathcal{N}_k(x)$
- 2. Choose  $\hat{y}$  according to a majority vote of the corresponding labels of the neighbors in  $\mathcal{N}_k(x)$



15-Nearest Neighbor Classifier



1-Nearest Neighbor Classifier





- Decision trees and random forests
- Bayesian classification (Naïve Bayes and Bayesian Belief Networks)
- Deep learning!
- Support Vector Machines
  - Kernel methods



- Recall that we consider both training error and test error
- How do we select the data for estimating the two?
- Common approaches:
  - **Holdout validation:** Randomly split the data into two disjoint sets. Train the model on one set and compute test error on the other.
  - **Random subsampling:** Repeat the holdout validation  $k$  times independently and compute the mean test error
  - Typically these approaches are not very robust



- Cross validation averages evaluation scores over several different splits to provide a more realistic assessment
- Types of cross validation
  - **Leave-one-out:** Use all data points except one for training and validate on the remaining point.
  - **2-fold cross validation:** Randomly split the data into two halves and run two validation iterations, each using one half for training and the other for validation
  - **k-fold cross validation:** Randomly partition the data into  $k$  equal size parts and run  $k$  validation iterations, each using one part for validation and the rest for training
- In all cases, the validation scores are averaged
- Each data point is only used once in testing but many times in training



- Goal: detect significant deviations from normal behavior
- Examples
  - Credit card fraud detection
  - Cybersecurity intrusion detection
  - Detecting bot traffic in online advertising
  - Malfunction detection in process monitoring
- Approaches
  - 1-class classification
  - Density estimation
  - Entropy estimation
  - Others...

# Unsupervised learning



- Dataset contains only features and no labels
- Goal is to find hidden patterns in the unlabeled data
- Examples
  - Density estimation
  - Data generation
  - Clustering
  - Data denoising
  - Representation learning
- Aside: ***Semi-supervised learning*** attempts to combine information from both labeled and unlabeled data to deduce information



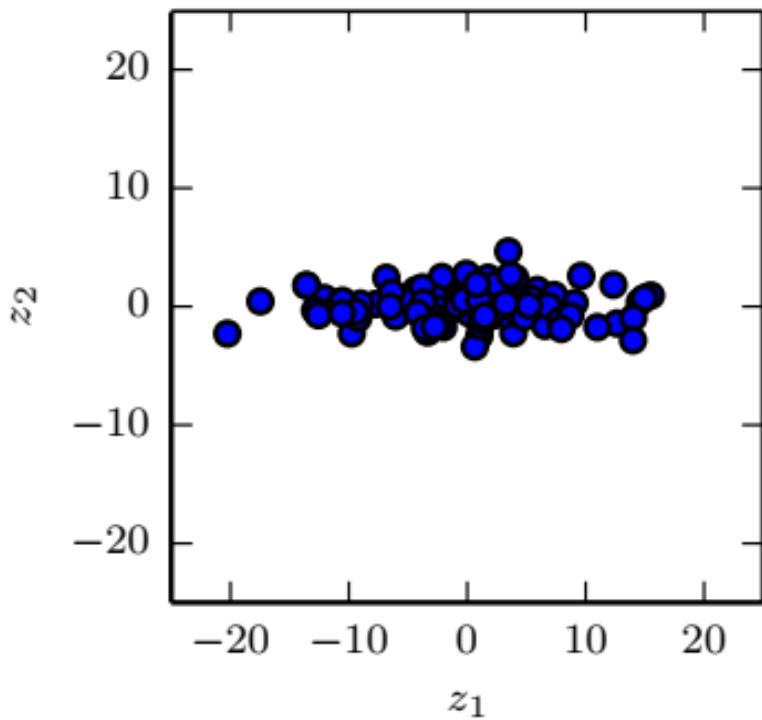
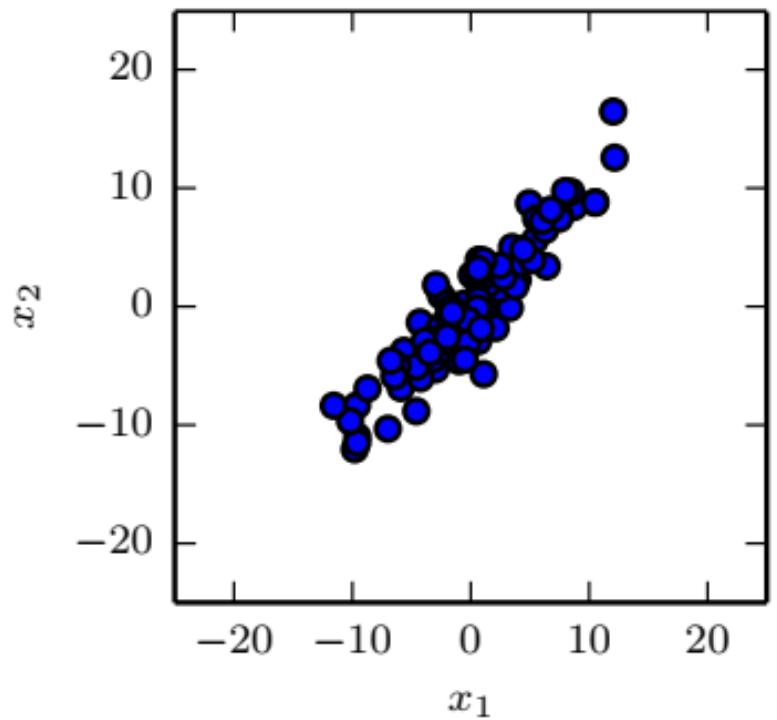
- Find the “best” representation of the data
  - I.e., find a representation of the data that preserves as much information as possible while obeying some penalty or constraint that simplifies the representation
- How do we define simple?
  - Low-dimensional
  - Sparse
  - Independent components
- Above criteria aren’t necessarily mutually exclusive
  - E.g., low-dimensional representations often have independent or weakly dependent components



- An unsupervised representation learning algorithm
  - Learns a lower dimensional representation
  - Elements are linearly uncorrelated with each other
  - Linear method
- PCA finds a low-dimensional linear transformation of the data that best reconstructs the original data (in the mean squared error sense)
  - $X$  is the data matrix
  - $\hat{X}_k$  is the reconstructed matrix from the  $k$ -dimensional representation
  - PCA minimizes  $\|X - \hat{X}_k\|_2^2$

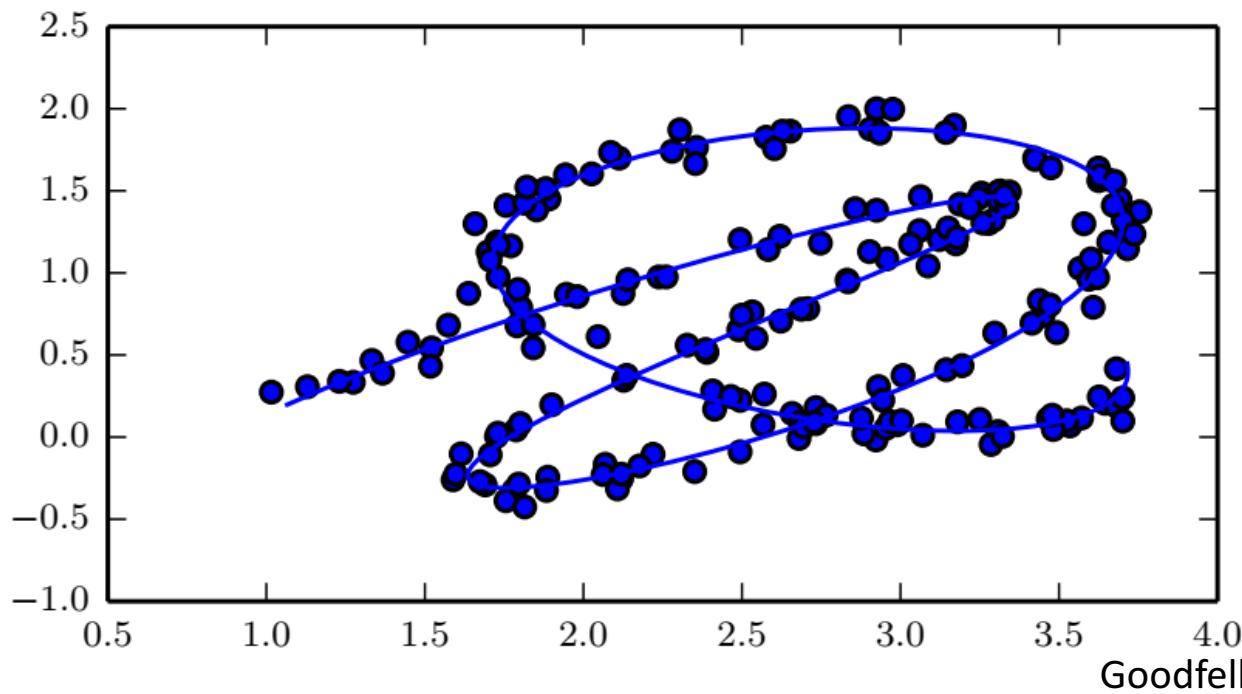


- PCA minimizes  $\|X - \hat{X}_k\|_2^2$
- Corresponds to selecting the first  $k$  eigenvectors and eigenvalues of  $X^T X$ 
  - $X^T X = W \Lambda W^T$
  - $W_k$  = the first  $k$  principal components, or the eigenvectors corresponding to the  $k$  largest eigenvalues
  - $W_k$  gives the low dimensional representation
  - $\hat{X}_k = X W_k W_k^T$  gives the reconstruction

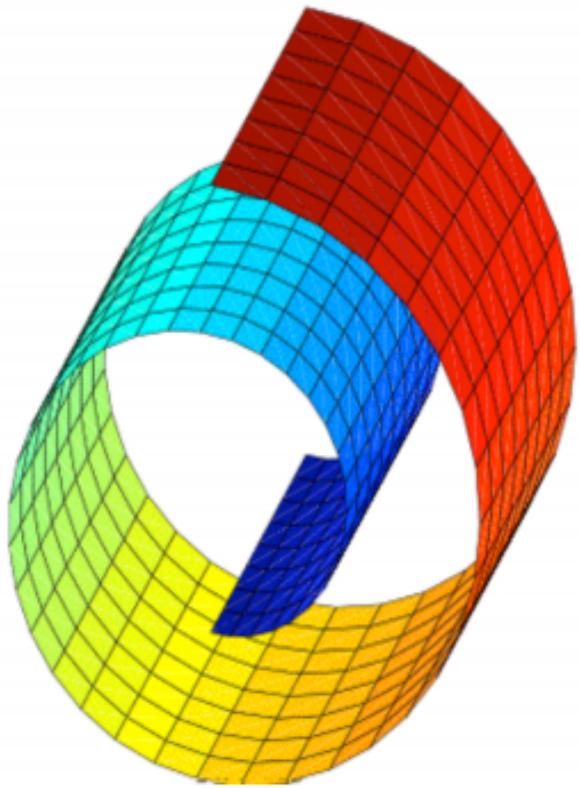




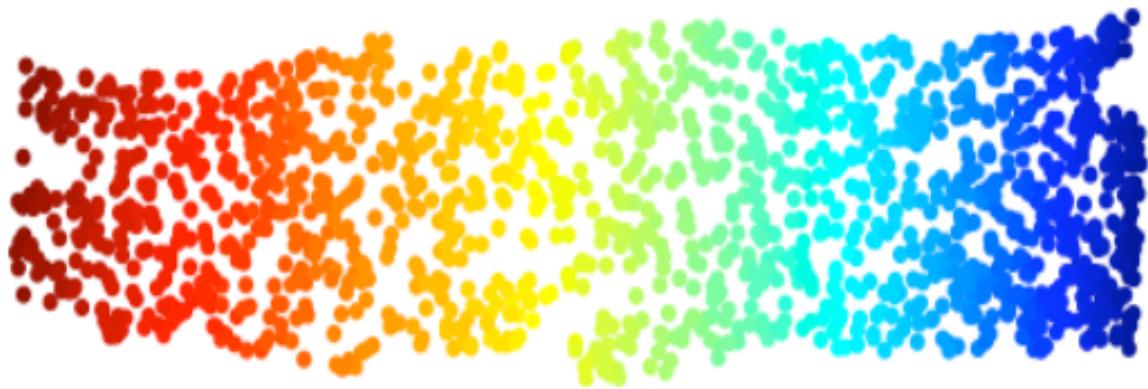
- Informally, a manifold in machine learning loosely means a connected set of points in high dimensions that can be approximated well by a small number of dimensions
- Generally a good assumption in many machine learning tasks (see Goodfellow book chapter 5.11.3)

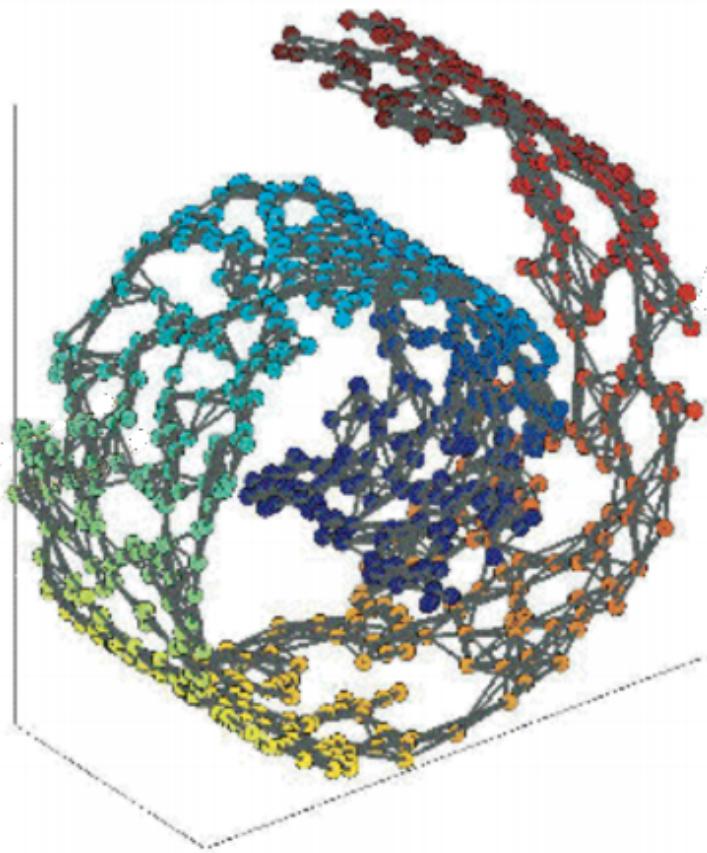


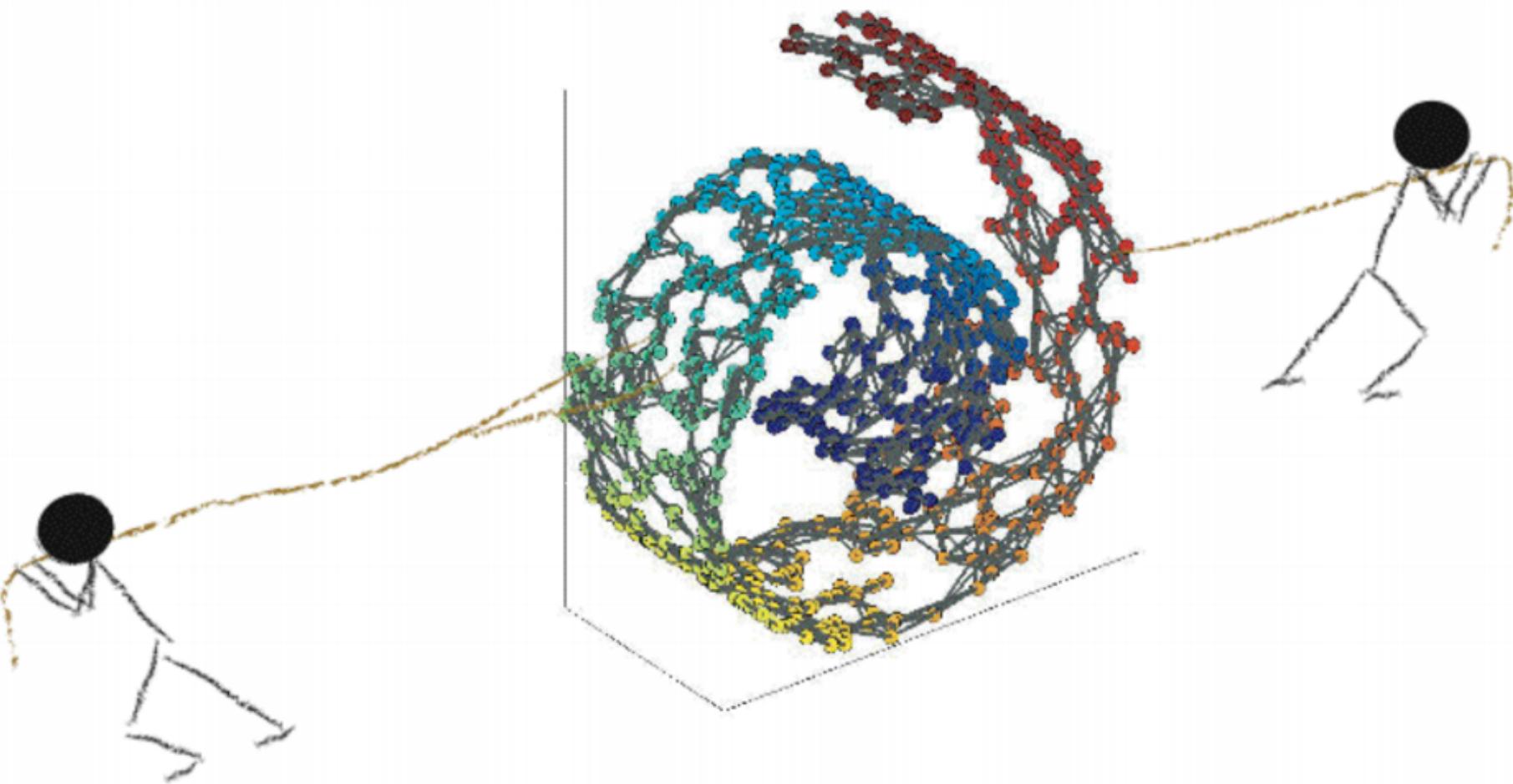
Goodfellow et al., 2016









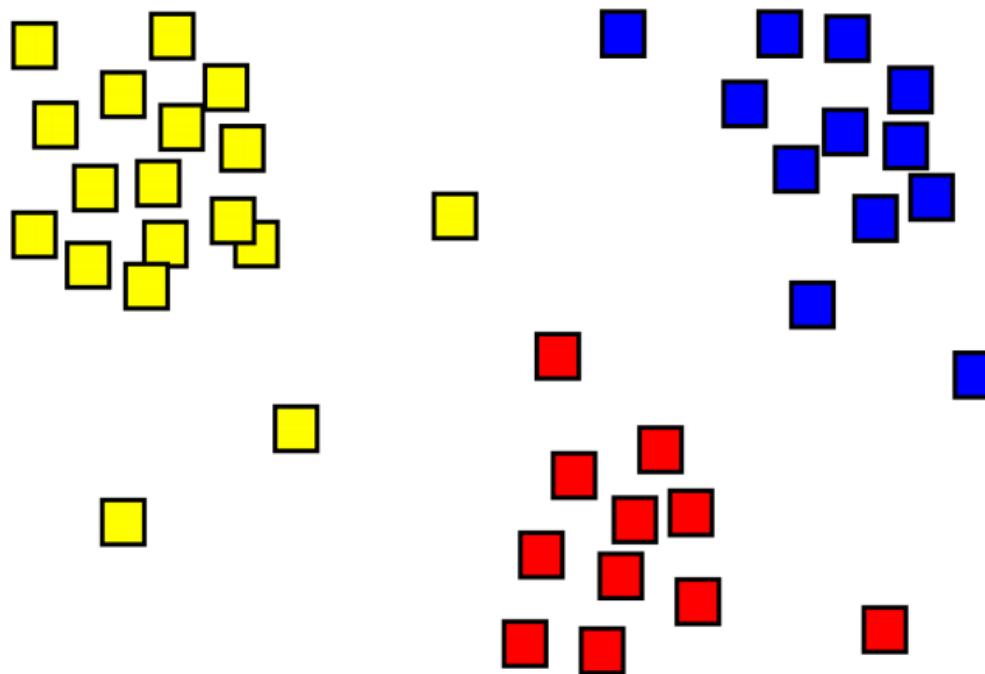




- Isomap
- Diffusion Maps (developed at Yale)
- Locally linear embedding
- Others



- Group objects so that objects within a cluster are similar to each other and dissimilar to objects in another cluster
  - Ill-posed problem  $\Rightarrow$  MANY clustering algorithms





- Suppose we have cars and trucks that are either red or gray
- How do you group these?
  - By color?
  - By vehicle type?
  - By both?
- Different clustering algorithms may give different results
  - Be careful which algorithm you choose
- Clustering quality depends on interpretability



- Examples
  - Clustering stocks to diversify investment
  - Community detection in social networks
  - Clustering genes and cells to uncover activities, reactions, and interactions
  - Network activity profiling by clustering packets/session
- Approaches
  - k-means
  - Hierarchical clustering
  - Spectral clustering
  - Variations on these using different metrics



- Chapter 5 in Goodfellow book
- *Elements of Statistical Learning* by Hastie, Tibshirani, & Friedman
- *Pattern Recognition and Machine Learning* by Bishop
- Online lecture notes from various sources, etc.

