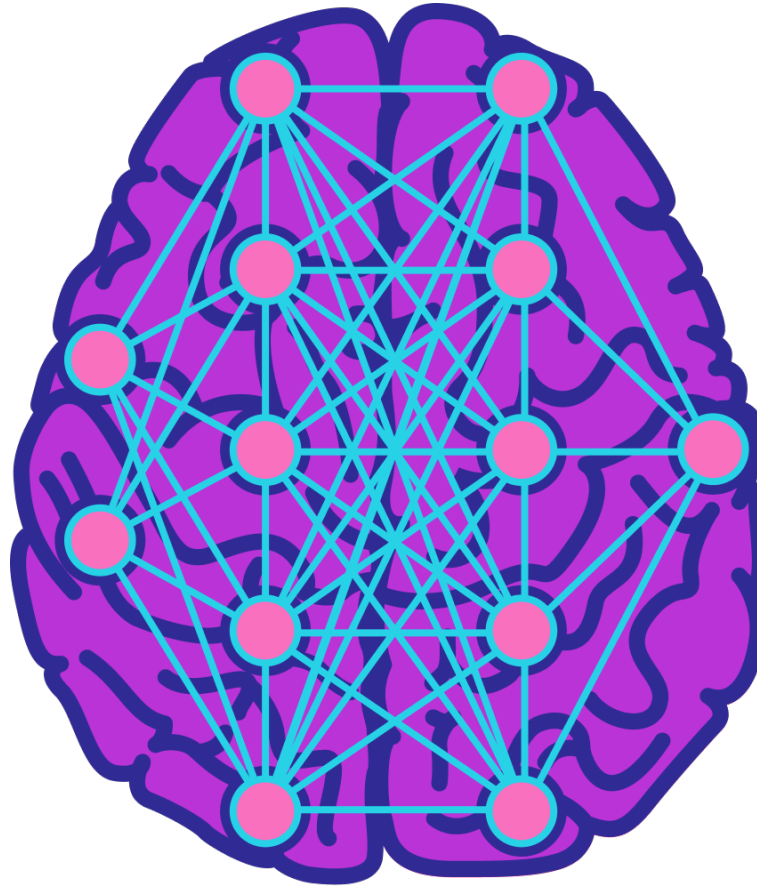# DS 5220

# Supervised Machine Learning and Learning Theory

Alina Oprea

Associate Professor, CCIS

Northeastern University

September 4 2019

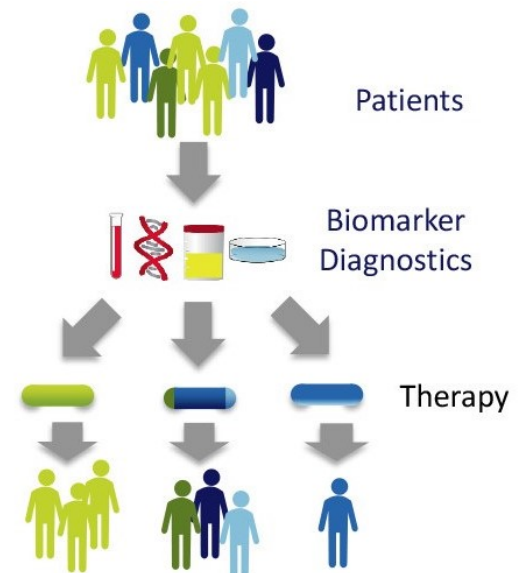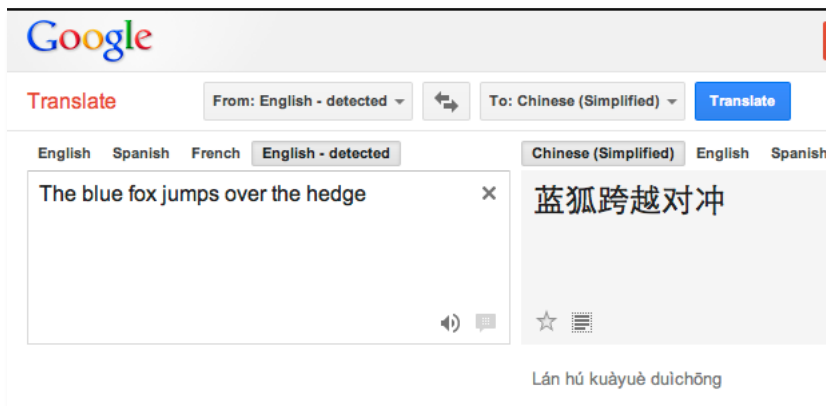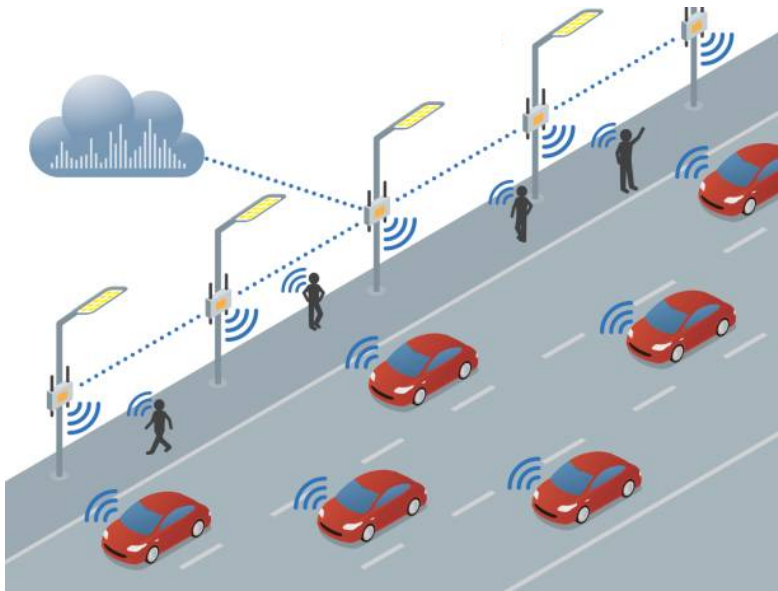# Welcome to DS 5220!



Supervised Machine Learning and Learning Theory

# Introduction
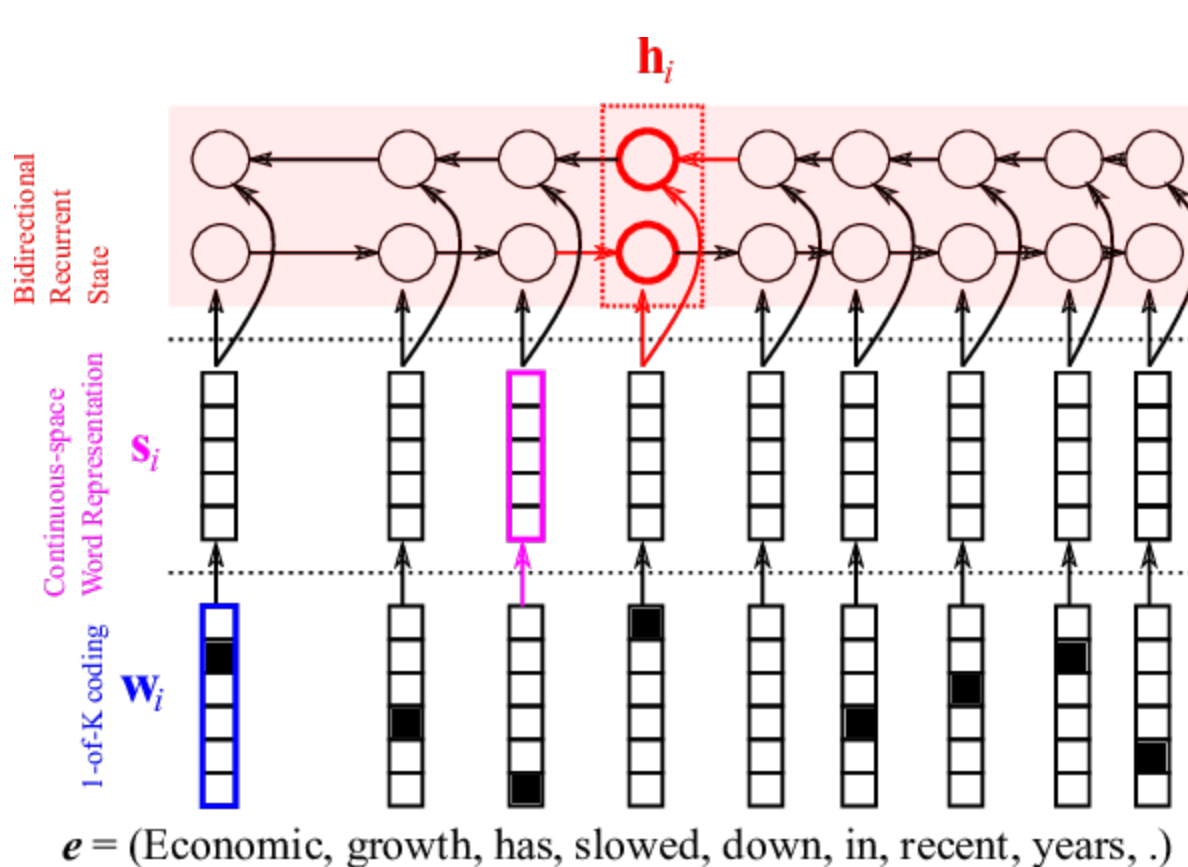
- **Ph.D. at CMU**
  - Research in storage security & cryptographic file systems
- **RSA Laboratories**
  - Cloud security, applied cryptography
  - Security analytics (ML in security)
- **NEU CCIS – since Fall 2016**
  - ML for security applications (attack detection, IoT, connected car security)
  - Adversarial ML (study the vulnerabilities of ML in face of attacks and design defenses)

# Machine learning is everywhere
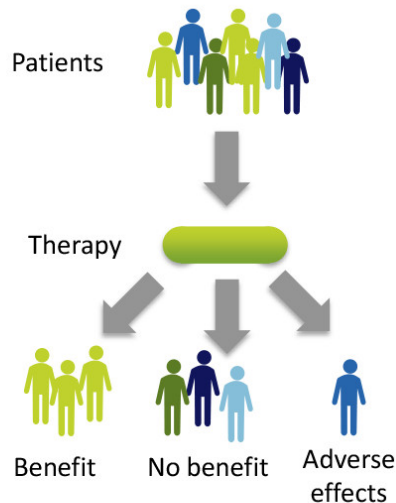
# Natural Language Processing (NLP)



$e$ = (Economic, growth, has, slowed, down, in, recent, years, .)

- Understand language semantics
- Real-time translation, speech recognition

# Personalized Medicine



**Without Personalized Medicine:**
Some Benefit, Some Do Not

Patients

Therapy

Benefit   No benefit   Adverse effects

**With Personalized Medicine:**
Each Patient Receives the Right Medicine For Them

Patients

Biomarker Diagnostics

Therapy

Each Patient Benefits From Individualized Treatment

- Treatment adjusted to individual patients
- Predictive models using a variety of features
- Better outcome and reduced cost

# Playing games



**a** Selection    **b** Expansion    **c** Evaluation    **d** Backup
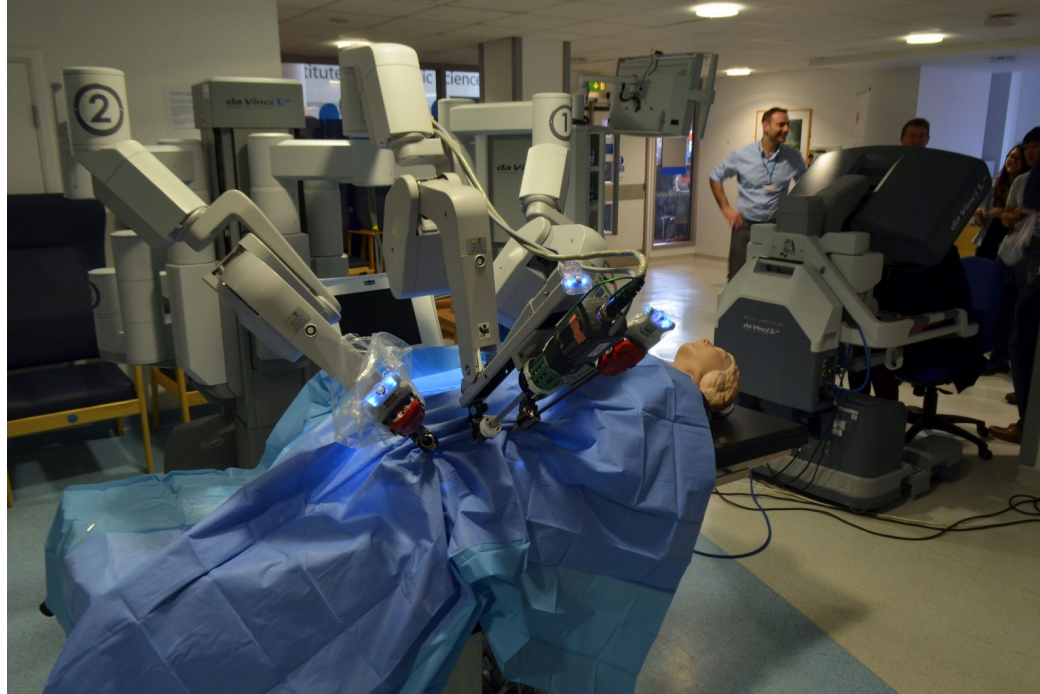
Reinforcement learning
- AlphaGo
- Chess

# Fast Forward in the Near Future



AI Transportation in Cities of the Future (10-20 years)

# Fast Forward in the Near Future



AI Robots in Medicine of the Future (10-20 years)

# Short History

- Legendre and Gauss – linear regression, 1805
  - Astronomy applications
- Bayes and Laplace - Bayes Theorem, 1812
- Markov chains, 1913
- Fisher – linear discriminant analysis, 1936
  - Logistic regression, 1940
- Widrow and Hoff ADELINE neural network, 1959
- Nelder, Wedderburn, generalized linear models, 1970
- "AI winter", limitations of  perceptron, 1970
- Breiman, Friedman, Olshen, Stone, decision trees, 1980
- More work on neural networks, 1980
- Cortes and Vapnik , SVM with kernels, 1990
- IBM Deep Blue beats Kasparov at chess, 1996
- Geoffrey Hinton, Deep learning, back propagation, 2006

# DS-5220

- What is *machine learning*?
  - The science of teaching machines how to learn
  - Design predictive algorithms that learn from data
  - Replace humans in critical tasks
  - Subset of Artificial Intelligence (AI)
- Machine learning very successful in:
  - Machine translation
  - Precision medicine
  - Recommendation systems
  - Self-driving cars
- Why the hype?
  - Availability: data created/reproduced in 2010 reached 1,200 exabytes
  - Reduced cost of storage
  - Computational power (cloud, multi-core CPUs, GPUs)

# DS-5220 Course objectives

- Become familiar with main machine learning tasks
  - Supervised learning vs unsupervised learning
  - Classification vs Regression
- Study most well-known algorithms
  - Regression (linear regression, spline regression)
  - Classification  (SVM, decision trees, Naïve Bayes, ensembles, etc.)
  - Deep learning (different neural network architectures)
- Learn the theory and foundation behind ML algorithms and learn to apply them to real datasets
- Learn about security challenges of ML
  - Introduction to adversarial ML

http://www.ccs.neu.edu/home/alina/classes/Fall2019

# Class Outline

- Introduction
  - Probability and linear algebra review
- Regression - 2 weeks
  - Linear regression, polynomial, spline regression
- Classification - 4 weeks
  - Linear classification (logistic regression, LDA)
  - Non-linear models (decision trees, SVM, Naïve Bayes)
  - Ensembles (random forest, AdaBoost)
  - Model selection, regularization, cross validation
- Neural networks and deep learning – 2 weeks
  - Back-propagation, gradient descent
  - NN architectures (feed-forward, convolutional, recurrent)
- Adversarial ML – 1 lecture
  - Security of ML at testing and training time

# Resources

• Gareth James, Daniela Witten, Trevor Hastie and Robert Tibshirani. [An Introduction to Statistical Learning with Applications in R](#)

• Trevor Hastie, Rob Tibshirani, and Jerry Friedman, [Elements of Statistical Learning](#), Second Edition, Springer, 2009.

• Christopher Bishop. [Pattern Recognition and Machine Learning](#). Springer, 2006.

• A. Zhang, Z. Lipton, and A. Smola. [Dive into Deep Learning](#)

# Policies

- Instructors
  - Alina Oprea
  - TAs: Yuxuan (Ewen) Wang; Christopher Gomes
- Schedule
  - Mon, Wed 2:50-4:30pm
  - West Village H 108
  - Office hours:
    - Alina: Wed 4:30 – 6:00 pm (ISEC 625)
    - Christopher : Monday 5:00-6:00pm (ISEC 605)
    - Ewen: Thursday 5:00-6:00pm (ISEC 605)
- Online resources
  - Slides will be posted after each lecture on public website
  - Piazza for questions and discussion
  - Gradescope for homework and project submission

# Policies, cont.

- Your responsibilities
  - Please be on time, attend classes, and take notes
  - Participate in interactive discussion in class
  - Submit assignments/ programming projects on time
- Late days for assignments
  - 5 total late days, after that loose 20% for every late day
  - Assignments are due at 11:59pm on the specified date
  - No need to email for late days, Gradescope shows submission time

# Grading

- Assignments – 20%
  - 4-5 assignments and programming exercises based on studied material in class
- Midterm – 25%
- Final exam – 25%
- Final project – 25%
  - Select your own project, public dataset
  - Submit short project proposal and milestone
  - Project pitch and final presentation
  - Project report
- Class participation – 5%
  - Participate in class discussion and on Piazza

# Assignments

- Programming exercises, and theory questions
  - Prefer Latex/Word/… write up
- <span style="color:red">Language</span>
  - Use R or Python
  - Jupyter notebooks recommended
- <span style="color:red">Submission</span>
  - Submit PDF report in Gradescope
  - Includes all the results, as well as link to code and instructions to run it

# Final project

- Goal: work on a larger data science project
  - Build your portfolio and increase your experience
- Requirements
  - Large dataset: at least 10,000 records (public source)
  - Not recommended to collect your own data
  - Pick application of interest, but instructor will also provide potential list of projects
  - Experiment with at least 3 ML models
  - Perform in-depth analysis (which features contribute mostly to prediction, which model performs best)
- Timeline
  - Proposal: mid class; milestone 2-3 weeks after (Instructors will provide early feedback)
  - Final presentation (10 mins) and report (5-6 pages)

# Academic Integrity

- Homework is done individually!
- Rules
  - Can discuss with colleagues or instructor
  - Can post and answer questions on Piazza
  - Code cannot be shared with colleagues
  - Cannot use code from the Internet
    - Use python or R packages, but not directly code for ML analysis written by someone else
- NO CHEATING WILL BE TOLERATED!
- Any cheating will automatically result in grade F and report to the university administration
- http://www.northeastern.edu/osccr/academic-integrity-policy/

# Outline

- Supervised learning
  - Classification
  - Regression

- Unsupervised learning
  - Clustering

- Bias-Variance Tradeoff

- Occam's Razor

# Introduction

- What is Machine Learning?
  - Subset of AI
  - Design algorithms that learn from real data and can automate critical tasks
- When can it be applied?
  - It cannot solve any problem!
  - When task can be expressed as learning task
  - When high-quality data is available
    - Labeled data (by human experts) is preferable!
  - When some error is acceptable (can rarely achieve 100% accuracy)
    - Example: recommendation system, advertisement engine

# Example 1
# Handwritten digit recognition



Images are 28 x 28 pixels

Represent input image as a vector $\mathbf{x} \in \mathbb{R}^{784}$
Learn a classifier $f(\mathbf{x})$ such that,

$$f : \mathbf{x} \rightarrow \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$$

MNIST dataset: Predict the digit
Multi-class classifier

# Data Representation

# Model the problem

As a supervised classification problem

Start with training data, e.g. 6000 examples of each digit



- Can achieve testing error of 0.4%

- One of first commercial and widely used ML systems (for zip codes & checks)

# Classification



Decision boundary

- Suppose we are given a training set of N observations

$$(x_1, \ldots, x_N) \text{ and } (y_1, \ldots, y_N), x_i \in \mathbb{R}^d, y_i \in \{0, 1\}$$

- Classification problem is to estimate f(x) from this data such that

$$f(x_i) = y_i$$

Extended to multi-class classification
- Handwritten digit recognition: $y_i \in \{0, 1, \ldots, 9\}$

# Classification

- **Training data**
  - $x_i = [x_{i,1}, \dots x_{i,d}]$: vector of image pixels (features)
  - Size $d = $ 28x28 $= 784$
  - $y_i$: image label
- **Models (hypothesis)**
  - Example: Linear model (parametric model)
    - $f(x) = wx + b$
  - Classify 1 if $f(x) > \mathrm{T}$ ; 0 otherwise
- **Classification algorithm**
  - Training: Learn model parameters $w, b$ to minimize error (number of training examples for which model gives wrong label)
  - Output: "optimal" model
- **Testing**
  - Apply learned model to new data and generate prediction $f(x)$

# Supervised Learning: Classification

**Training**



Data → Pre-processing → Feature extraction → Learning model

Labeled

$x_i, y_i \in \{0,1\}$

Normalization

Feature Selection

Classification

$f(x)$

**Testing**

New data → Learning model → Predictions

$y' = f(x') \in \{0,1\}$

Unlabeled

$x'$

$f(x)$

Positive
Negative

Classification

# Example Classifiers

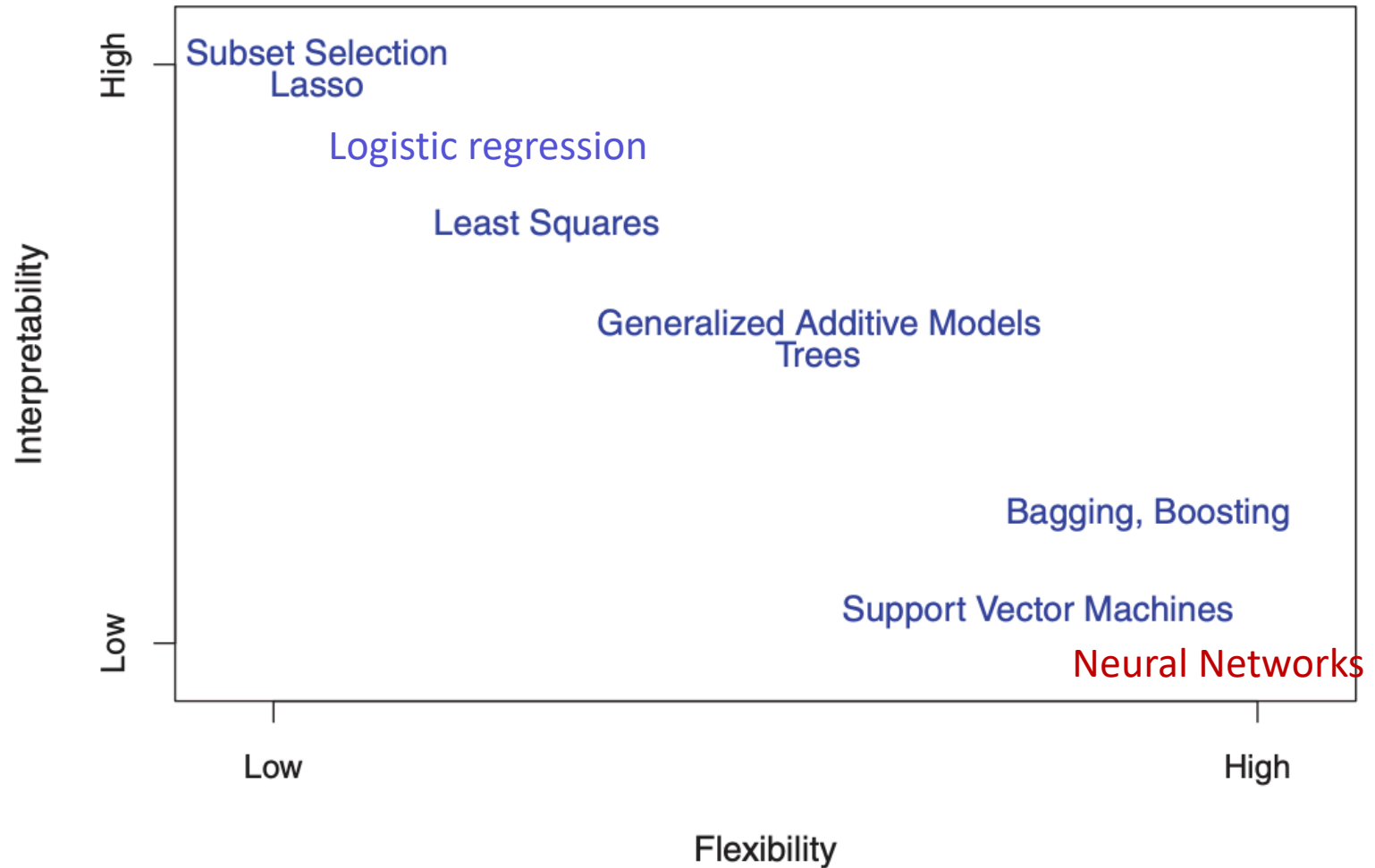

Linear classifiers: logistic regression, LDA (parametric)



Decision trees (non-parametric)

Nearest Neighbors: kNN (non-parametric)

# Interpretability

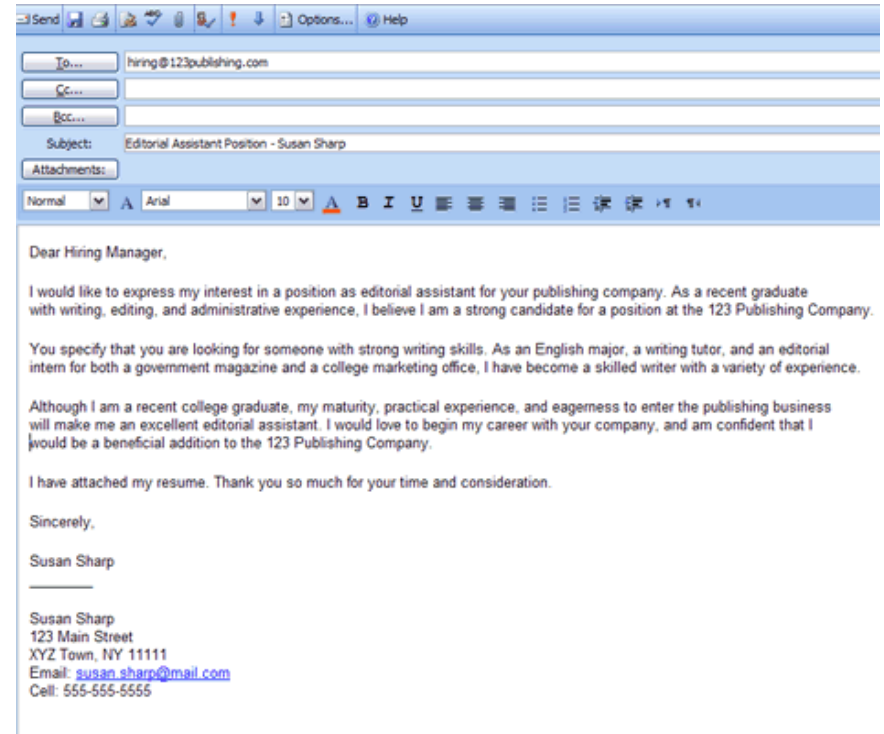# Real-world example: Spam email



SPAM email
- Unsolicited
- Advertisement
- Sent to a large number of people

# Classifying spam email



**Content-related features**
- Use of certain words
- Word frequencies
- Language
- Sentence

**Structural features**
- Sender IP address
- IP blacklist
- DNS information
- Email server
- URL links (non-matching)

# Classifying spam email

SPAM

REGULAR

**Feature extraction**
- Content
- Structural

Numerical

**Classifier**
- Logistic regression
- Decision tree
- SVM

Labeled data
- SPAM
- REGULAR

Training

New email

Model

SPAM → Filter

REGULAR → Allow

Testing

33

# Example 2
# Stock market prediction



S&P/TSX COMPOSITE
as of 4-Apr-2008

Copyright 2008 Yahoo! Inc.          http://finance.yahoo.com/

- Task is to predict stock price at future date
- This is a regression task, as the output is continuous

# Regression



Linear regression
1 dimension

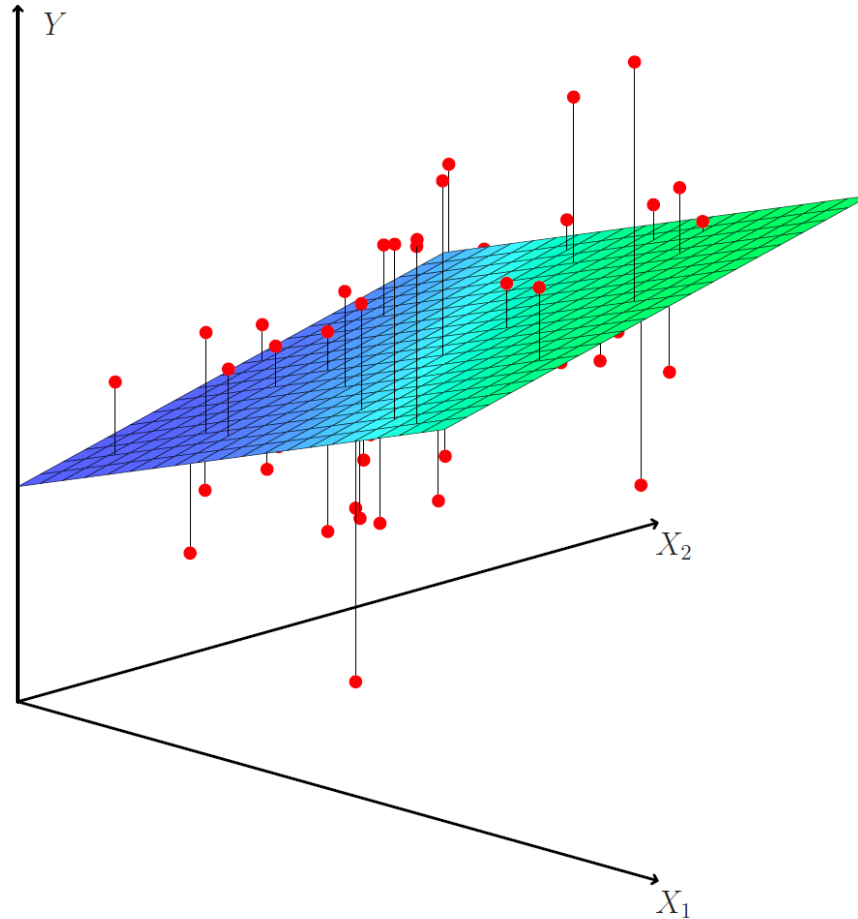- Suppose we are given a training set of N observations

$$(x_1, \ldots, x_N) \text{ and } (y_1, \ldots, y_N)$$

- Regression problem is to estimate y(x) from this data

$$x_i = (x_{i1}, \ldots, x_{id})$$ - d predictors (features)
$$y_i$$ - response variable, numerical
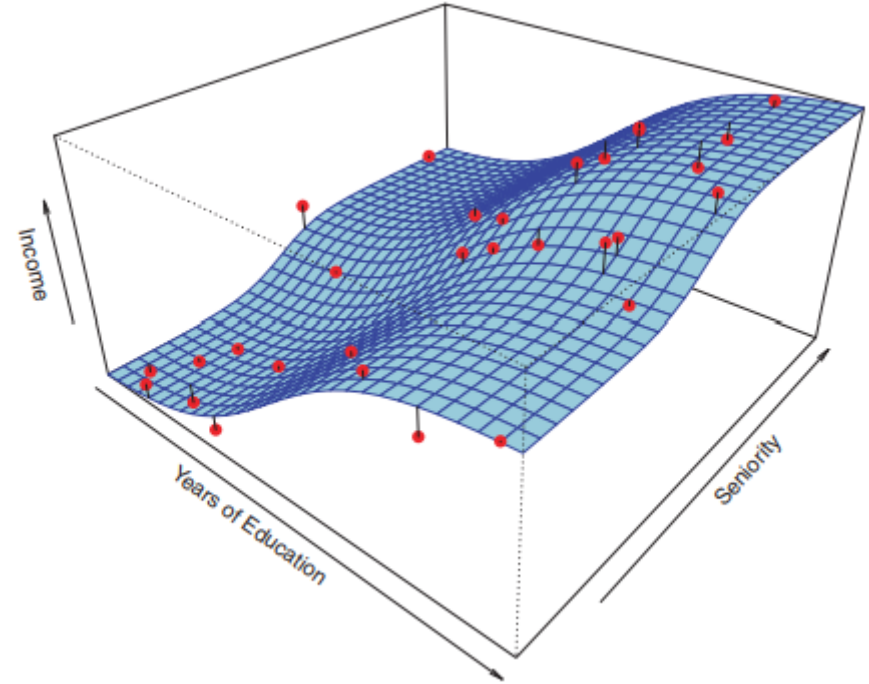
# Multi-dimensional linear regression



Minimize sum of square error
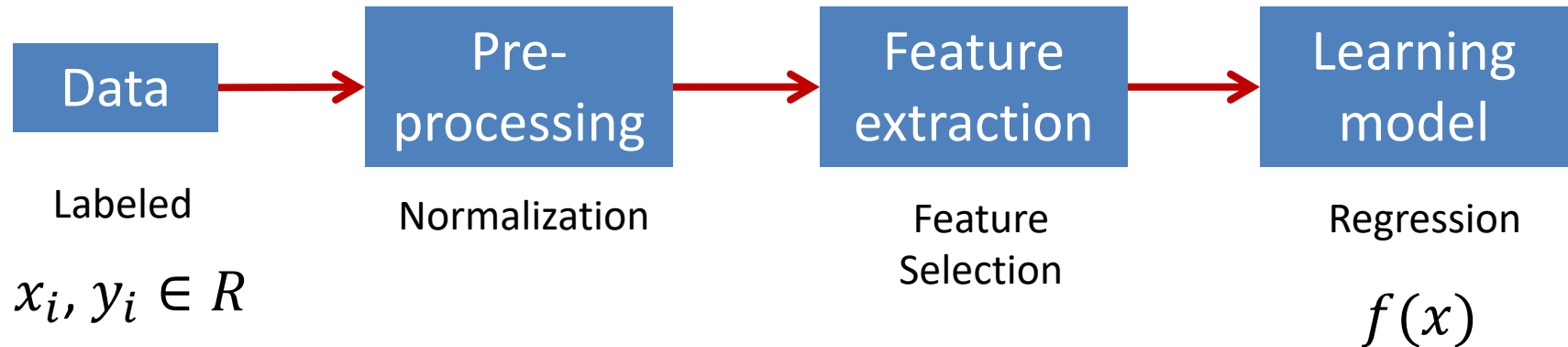(error/loss function)

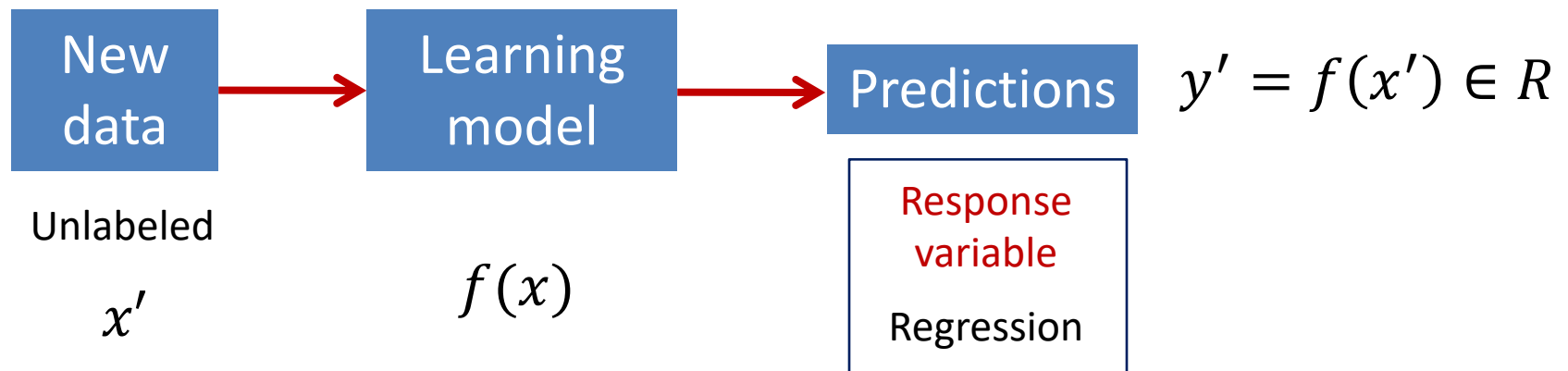# Income Prediction



Linear Regression

Non-Linear Regression
Polynomial/Spline Regression
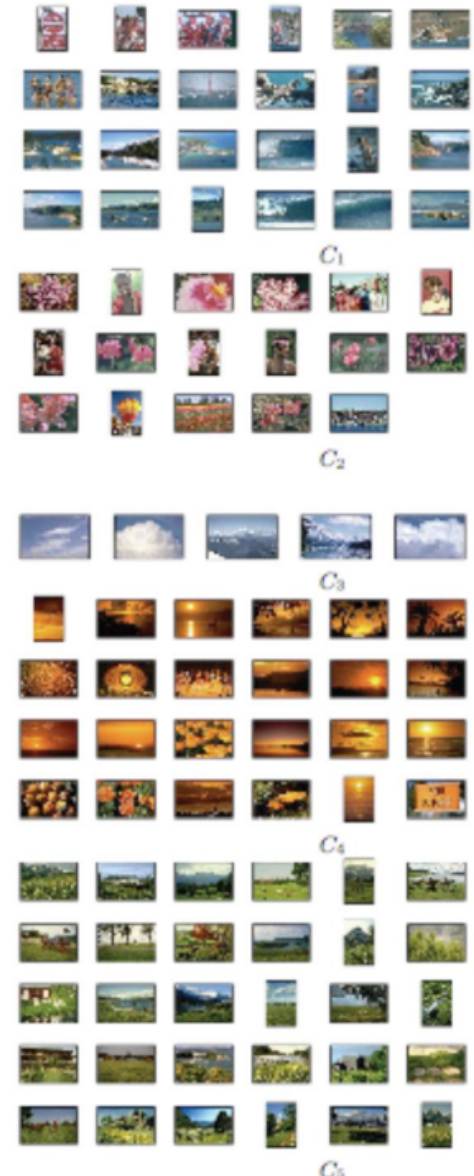
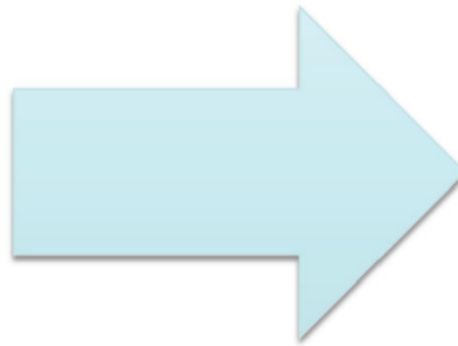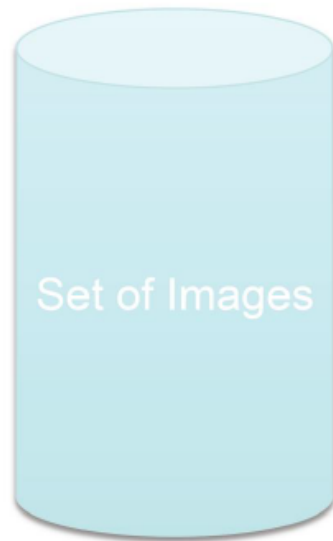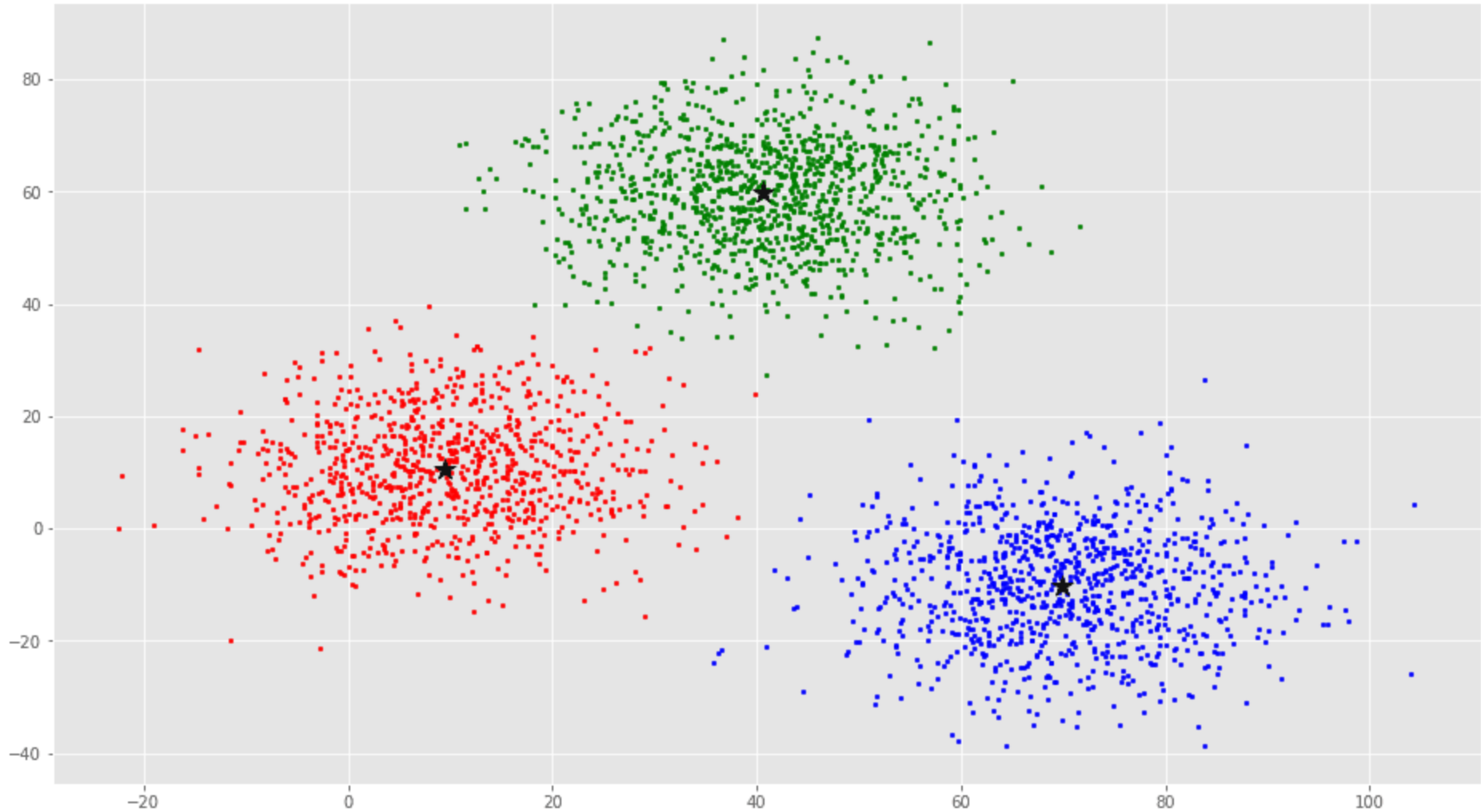# Supervised Learning: Regression

**Training**



Data → Pre-processing → Feature extraction → Learning model

Labeled

$x_i, y_i \in R$

Normalization

Feature Selection

Regression

$f(x)$

**Testing**

New data → Learning model → Predictions    $y' = f(x') \in R$

Unlabeled

$x'$

$f(x)$

Response variable

Regression

# Example 3: image search

## Clustering images



Set of Images
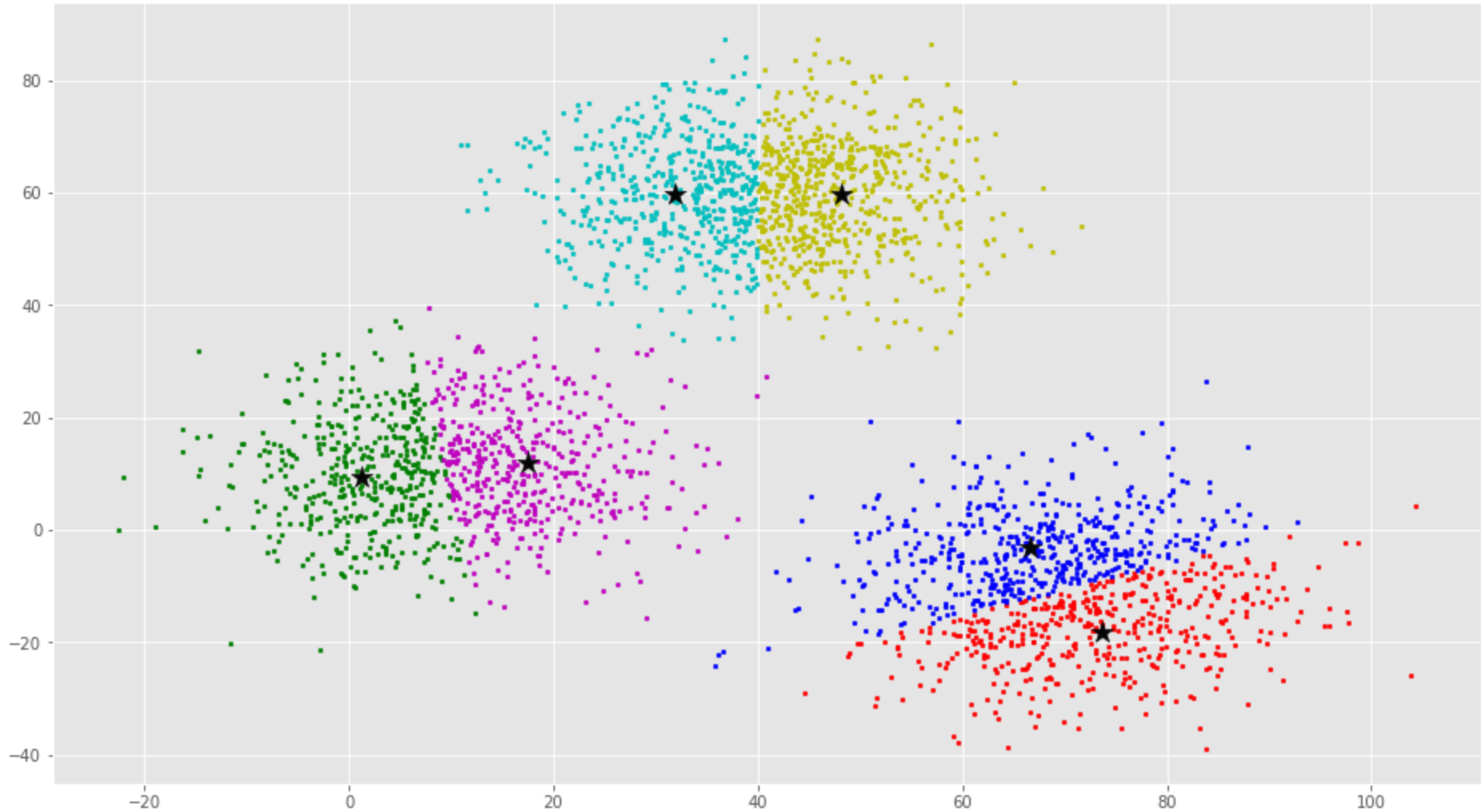
Find similar images to a target one

# K-Means Clustering



K=3

# K-means Clustering



K=6

# Hierarchical Clustering



**Cluster Dendrogram**

d
hclust (*, "ward.D2")

# Unsupervised Learning

- Clustering
  - Group similar data points into clusters
  - Example: k-means, hierarchical clustering
- Dimensionality reduction
  - Project the data to lower dimensional space
  - Example: PCA (Principal Component Analysis)
- Feature learning
  - Find feature representations
  - Example: Autoencoders

# Supervised Learning Tasks

- Classification
  - Learn to predict class (discrete)
  - Minimize error 1/N $\sum_{i=1}^{N}[y^{(i)} \neq f(x^{(i)})]$

- Regression
  - Learn to predict response variable (numerical)
  - Minimize MSE (Mean Square Error between prediction and actual values)

- Both classification and regression
  - Training and testing phase
  - "Optimal" model is learned in training and applied in testing