

STAT 37710 / CAAM 37710 / CMSC 35400
Machine Learning

Introduction

Cong Ma

What is machine learning?

- Wiki's definition of machine learning (ML):
 - Machine learning (ML) is a field of inquiry devoted to understanding and building methods that 'learn', that is, methods that **leverage data** to improve performance on some set of **tasks** ---adapted from Tom Mitchell

Examples of machine learning

- **spam filter**

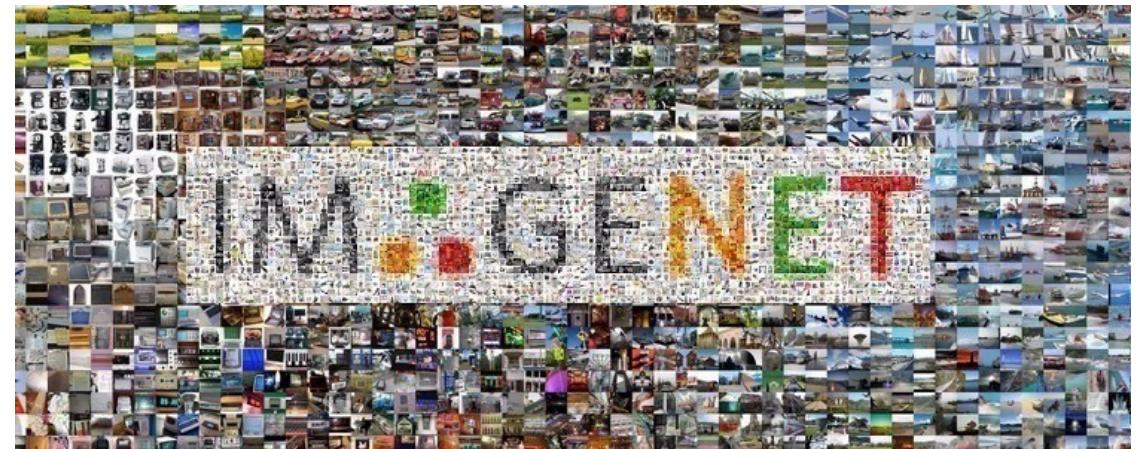
- Data: emails together with labels
- Task: label emails to either spam or non-spam
- Performance: accuracy in labelling emails



Examples of machine learning

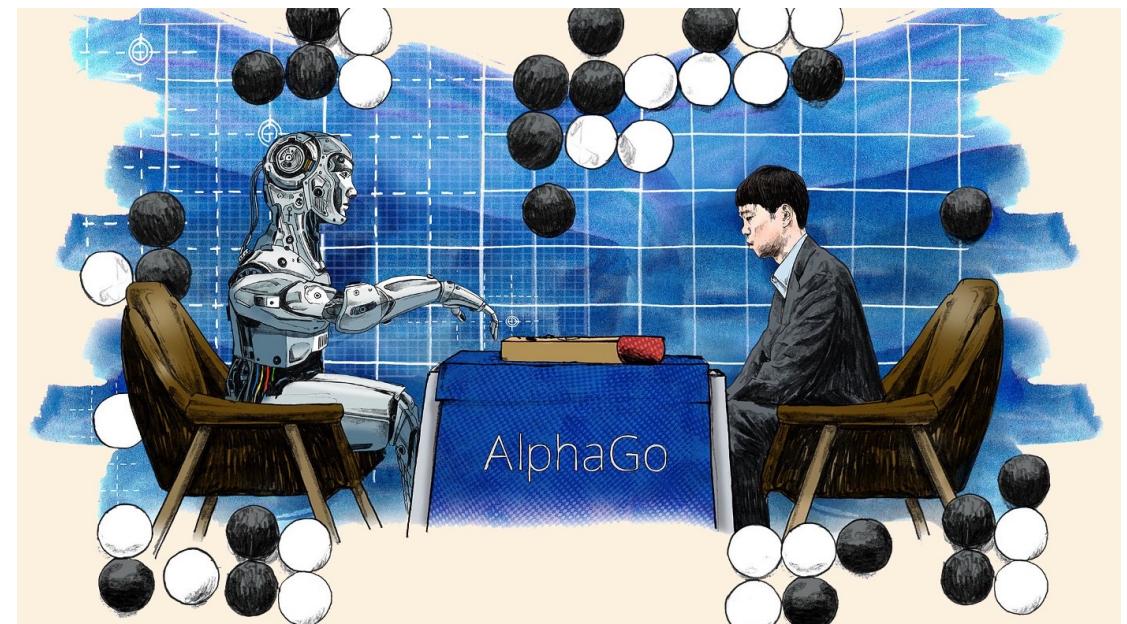
- **image classification**

- Data: images together with labels (ImageNet)
- Task: label images to categories (e.g., cat, dog)
- Performance: accuracy in labelling images



Examples of machine learning

- Playing Go
 - Data: history of game playing
 - Task: playing Go well
 - Performance: winning rate against world champion



Machine learning is pervasive...

A large-scale crowd-sourced analysis of abuse against women journalists and politicians on Twitter

Laur Ele Archy Ele

We re between to stud curate for the release the tec aware

The us ing an disaster fast an and it's most c and of frame satelli find a novel of two frame on the

From Satellite Imagery to Disaster Insights

jigaz

1C

sgur jinyongch

Po the com reso cons nega Fort

Wildlife Poaching Prediction with Data and Human Knowledge *

A Scalable, Flexible Augmentation of the Student Education Process

Bhairav Mehta
Mila, Université de Montréal
bhairav.mehta@umontreal.ca

Adithya Ramanathan
University of Michigan
adithram@umich.edu

Abstract

We present a novel intelligent tutoring system which builds upon well-established

Helping Bees and Beekeepers with AI

Honeybee Identification with Machine Learning Using Augmented Microscopy

Peter He (Department of Computing, Imperial College London) · Alexis Gkantrigas (Department of Molecular Biology, University College London) · Gerard Glowacki (Imperial College London)

Towards a Sustainable Food Supply Chain Powered By Artificial Intelligence

Volodymyr Kuleshov, Marjan Seymour, Danny Nemer, Nathan Fenner, Matthew Schwartz
Afresh Technologies and Stanford University

Introduction

Improving Traffic Safety in Jakarta Through video Analysis

João Caldeira, Alex Fout, Aniket Kesari, Raesetje Sefala, Katy Dupre, Joe Walsh

University of Chicago, Colorado State University

- **Problem:** ~2,000 people die in traffic accidents in Jakarta, Indonesia
- The city of Jakarta invests in traffic cameras to capture data on traffic behavior, but this data does not scale with an increasing population.

THE
UNIVERSITY
OF IOWA

Next Hit Predictor - Self-exciting Risk Modeling for Predicting Next Locations of Serial Crimes

Yunyi Li
The University of Iowa
yunyi-li@uiowa.edu

Tong Wang
The University of Iowa

ML is interdisciplinary

statistics

information theory

Machine Learning

optimization

algorithms

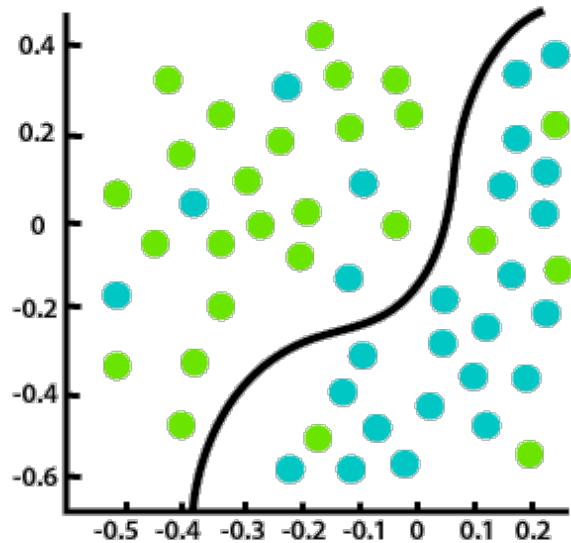
Machine learning tasks

- *Depending on the feedbacks, machine learning can be decomposed into*
 - *Supervised Learning*
 - Regression
 - Classification
 - *Unsupervised Learning*
 - Clustering
 - Dimension reduction Anomaly detection, ...
- *Many other specialized tasks*

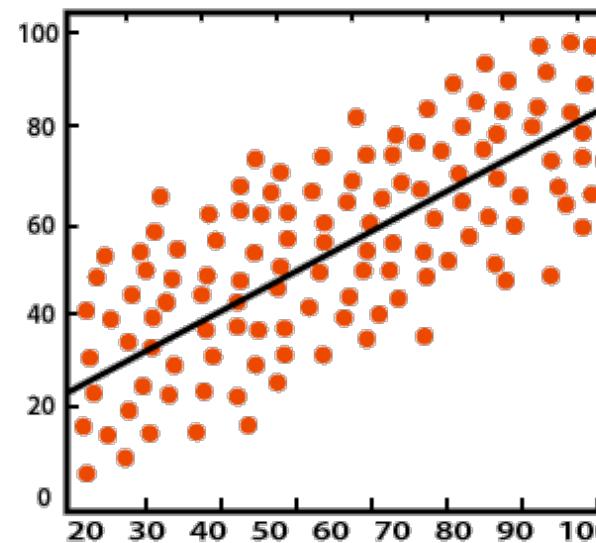
Supervised learning

$$f : X \rightarrow Y$$

Supervised learning

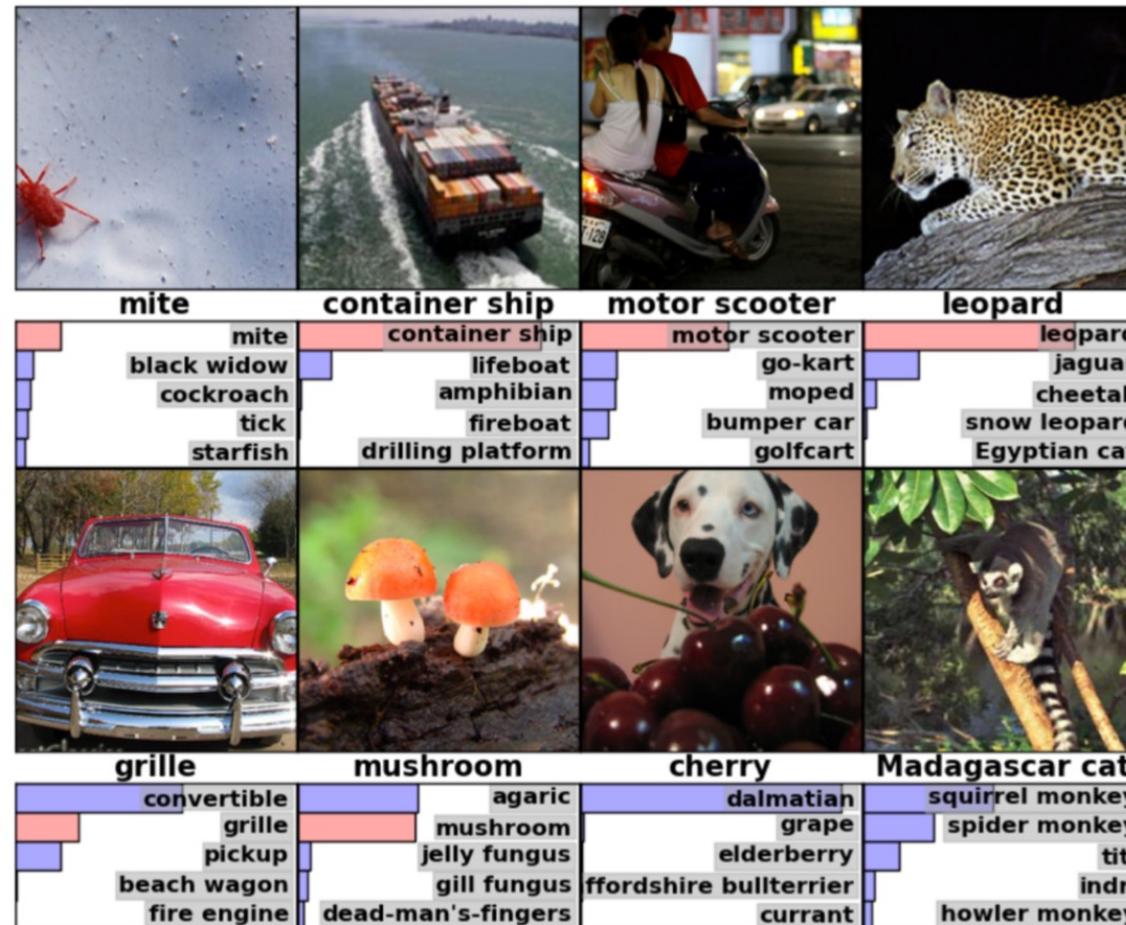


Classification



Regression

Example: Image classification



Regression

- **Goal:** Predict **real valued** labels (possibly vectors)
- Examples:

X

Flight route

Real estate objects

Patient & drug

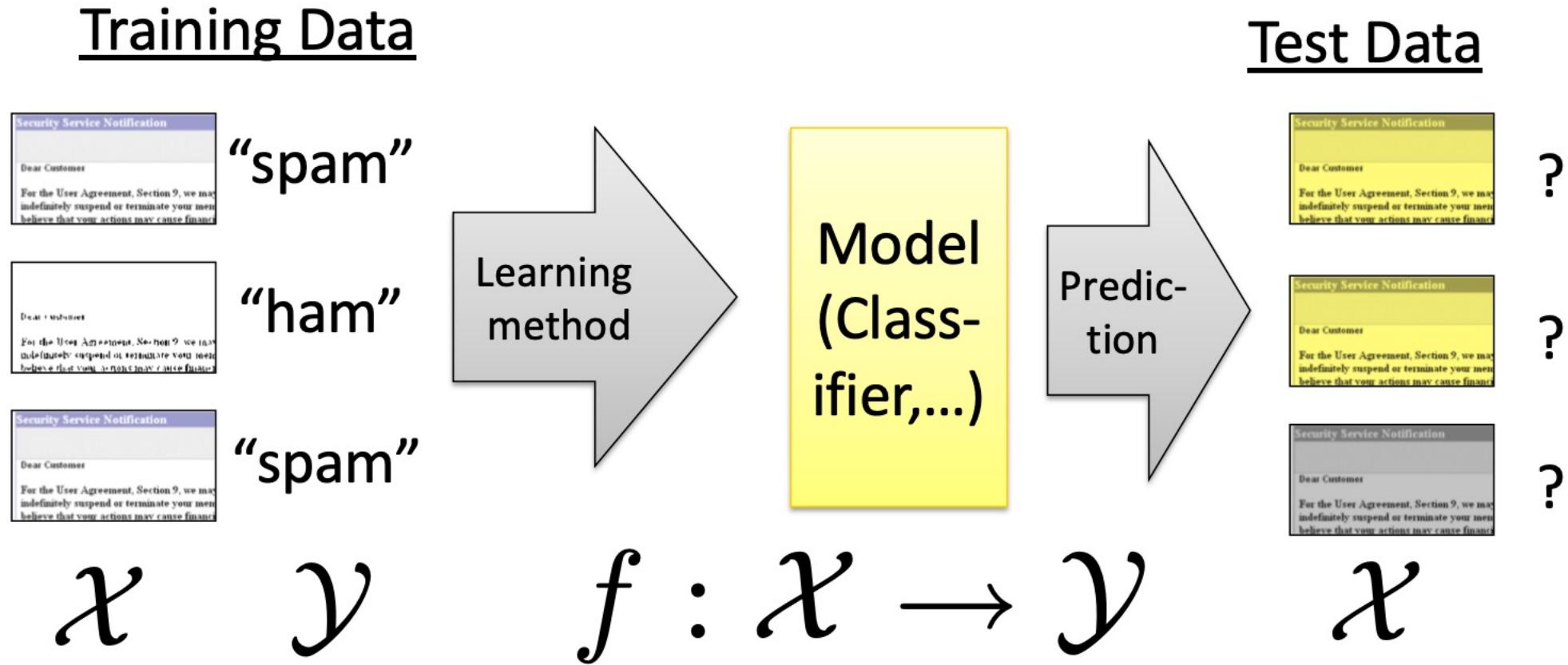
Y

Delay (minutes)

Price

Treatment effectiveness ...

Basic supervised learning pipeline



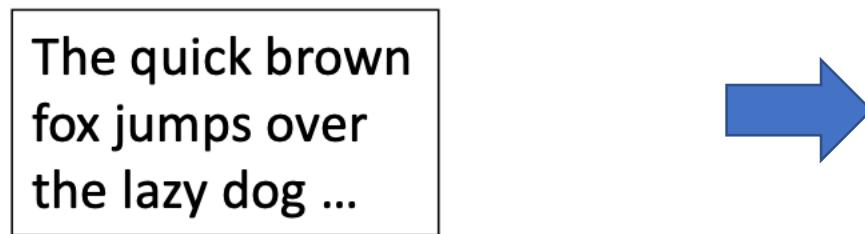
Example: Classifying documents

- **Input:**
 - Documents with labels, but how to represent documents
- **Goal:**



Representing data

- Learning methods expect standardized representation of data
 - (e.g., Points in vector spaces, nodes in a graph, similarity matrices ...)



- Concrete choice of representation (“features”) is crucial for successful learning
- This class (typically): **feature vectors** in \mathbb{R}^d

Example: Bag-of-words

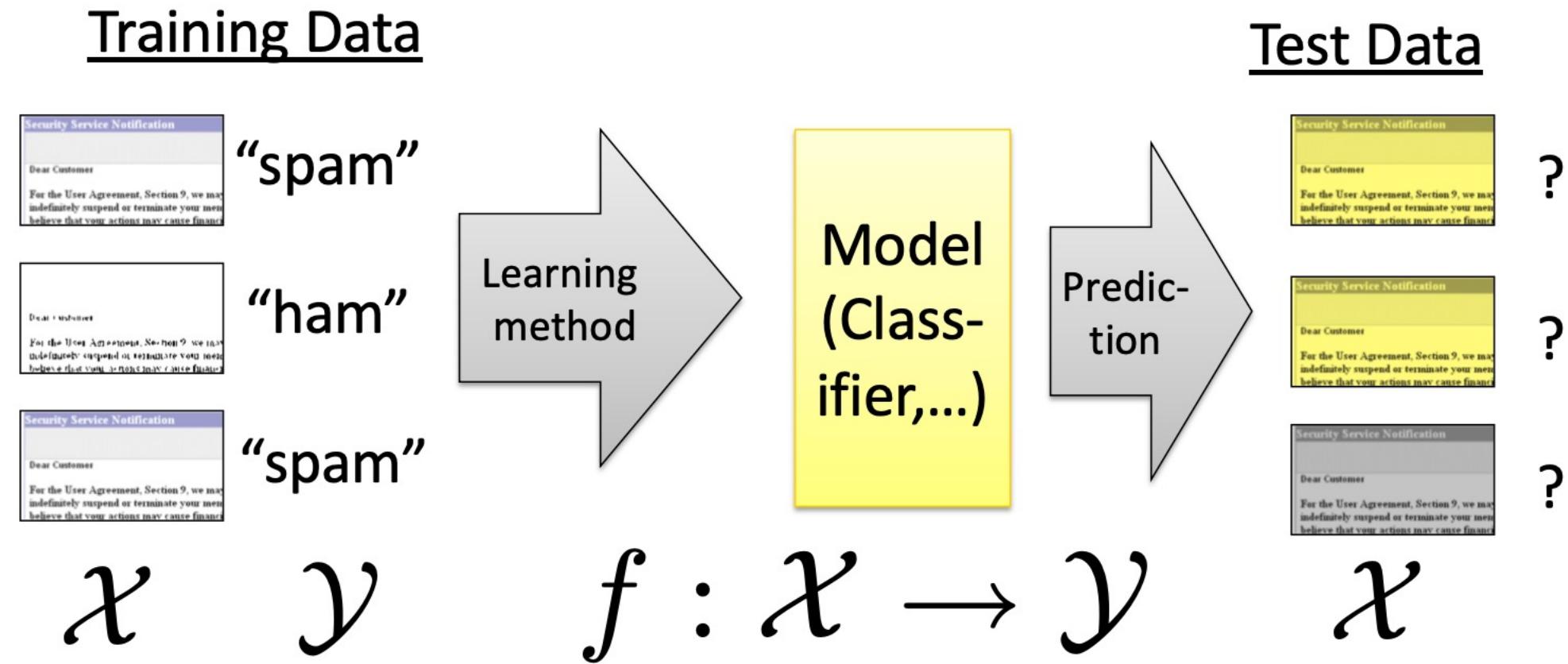
- Suppose language contains at most $d=100000$ words
- Represent each document as a vector \mathbf{x} in \mathbb{R}^d
- i -th component x_i counts occurrence of i -th word

Word	Index
a	1
abandon	2
ability	3
...	
is	578
...	
test	2512
...	
this	2809
....	

Bag-of-words: Improvements

- Length of the document should not matter
 - Replace counts by binary indicator (yes/no) Normalize to unit length
- Some words more “important” than others
 - Remove “stopwords” (the, a, is, ...)
 - Stemming (learning, learner, learns -> learn)
 - Discount frequent words (tf-idf)
- Bag-of-words ignores order
 - Consider pairs (n-grams) of consecutive words
- Does not differentiate between similar and dissimilar words (ignores semantics)
 - Word embeddings (e.g., word2vec, GloVe)

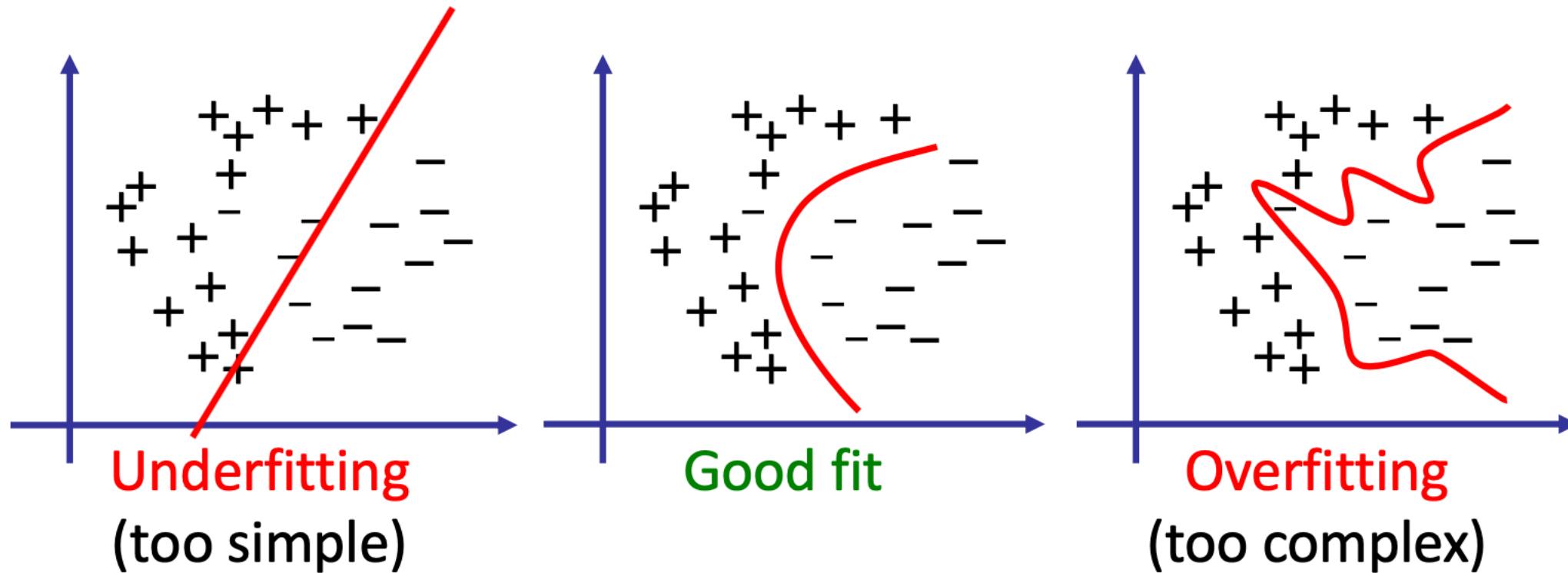
Basic supervised learning pipeline



Model class

- Linear
- decision tree
- random forests
- graphical models
- deep neural networks

Model selection and validation



Other models of learning

- **Unsupervised learning**
 - Learning without labels
- **Semi-supervised learning**
 - Learning from both labeled and unlabeled data
- **Transfer learning**
 - Learn on one domain and test on another
- **Active learning**
 - Acquiring most informative data for learning
- **Online learning**
 - Learning from examples as they arrive over time
- **Reinforcement learning**
 - Learning by interacting with an unknown environment

Summary

- Where we are

- Basic forms of learning:
 - Supervised learning and other modes of machine learning
- Key challenge in ML
 - Trading goodness of fit and model complexity
- Representation of data is of key importance

- What's next

- Formally state machine learning problems
- Estimation theory and bias-variance tradeoff