# MACHINE LEARNING BASICS

David Kauchak
CS159 Spring 2023

1

## Admin

Assignment 6

pre-pre enrollment

No office hours Friday

2

## Machine Learning is...

Machine learning, a branch of artificial intelligence, concerns the construction and study of systems that can learn from data.

WIKIPEDIA
The Free Encyclopedia

3

## Machine Learning is...

Machine learning is programming computers to optimize a performance criterion using example data or past experience.

-- Ethem Alpaydin

The goal of machine learning is to develop methods that can automatically detect patterns in data, and then to use the uncovered patterns to predict future data or other outcomes of interest.

-- Kevin P. Murphy

The field of pattern recognition is concerned with the automatic discovery of regularities in data through the use of computer algorithms and with the use of these regularities to take actions.

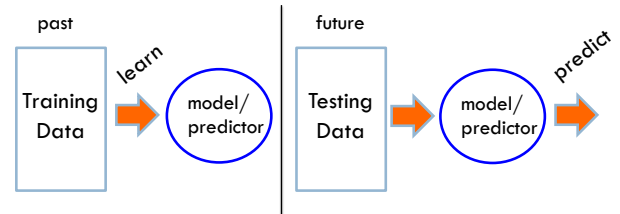-- Christopher M. Bishop

4

## Machine Learning is...

Machine learning is about predicting the future based on the past.
-- Hal Daume III

5

## Machine Learning is...

Machine learning is about predicting the future based on the past.
-- Hal Daume III

| past | | future | |
|------|------|--------|------|
| Training Data | → learn → model/predictor | Testing Data | → model/predictor → predict |

6

## Why machine learning?

Lot's of data

Hand-written rules just don't do it

Performance can be much better than what people can do
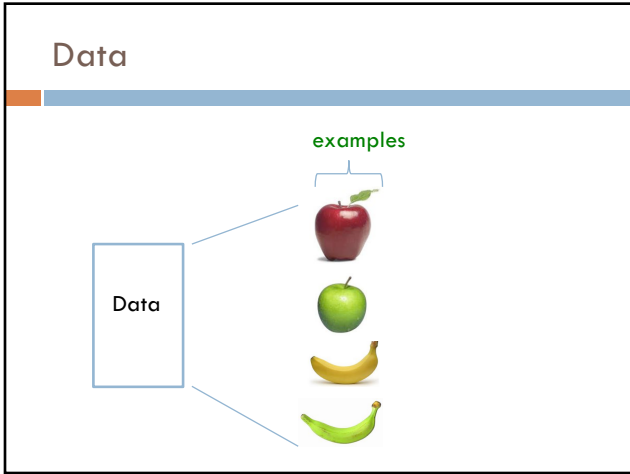
Why not just study machine learning?

☐ Domain knowledge/expertise is still very important
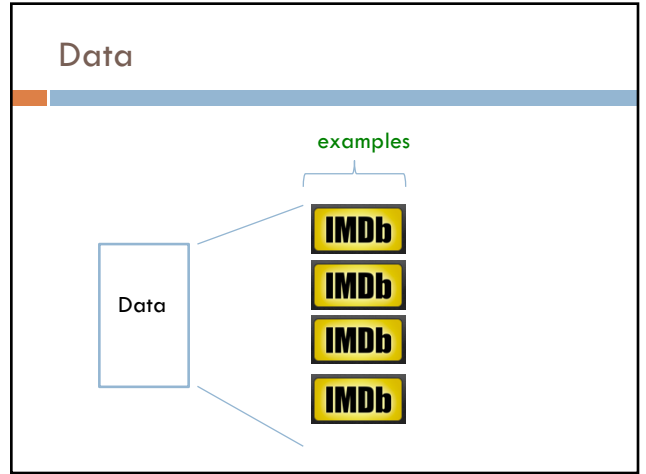☐ What types of features to use
☐ What models are important

7

## Machine learning problems

What high-level machine learning problems and algorithms have you seen or heard of before?
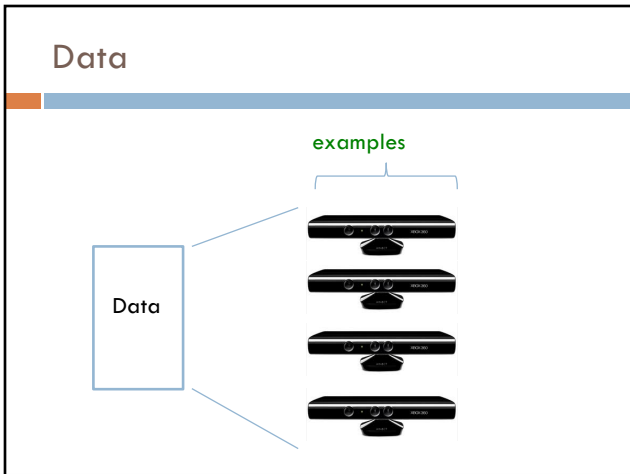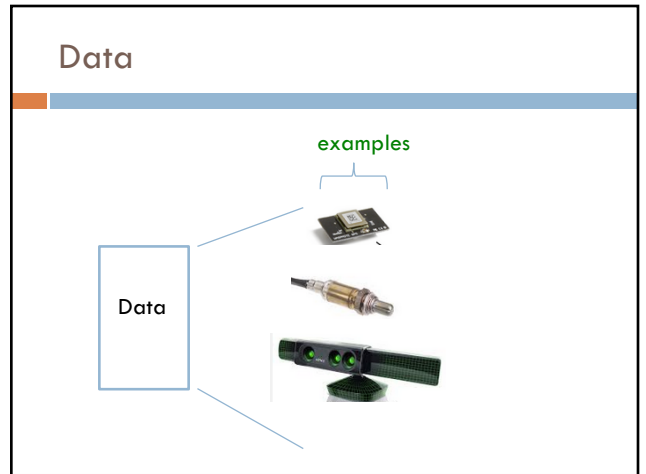
8

## Data

examples

Data

9

## Data

examples

Data

IMDb
IMDb
IMDb
IMDb

10

## Data

examples

Data

11

## Data

examples

Data

12

## Supervised learning

examples

label

label1

label3

label4

label5

labeled examples

Supervised learning: given labeled examples

13

## Supervised learning

label

label1

label3

label4

label5

model/
predictor

Supervised learning: given labeled examples

14

## Supervised learning

model/
predictor

predicted label

Supervised learning: learn to predict new example

15

## Supervised learning: classification

label

apple

apple

banana

banana

Classification: a finite set of labels

Supervised learning: given labeled examples

16

## NLP classification applications

Document classification
- spam
- sentiment analysis
- topic classification

Turn SafeSearch on or off
https://support.google.com/websearch/answer/510
1. Visit the Search Settings page.
2. In the "SafeSearch filters" section, select or unselect **Filter explicit results**.
3. Click **Save** at the bottom of the page.

Does linguistics phenomena X occur in text Y?

Digit recognition

Grammatically correct or not?

Word sense disambiguation

Any question you can pose as to have a discrete set of labels/answers!

17

## Supervised learning: regression

label

-4.5

10.1

Regression: label is real-valued

3.2

4.3
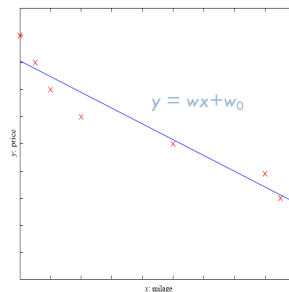
Supervised learning: given labeled examples

18

## Regression Example

Price of a used car

$x$ : car attributes
(e.g. mileage)
$y$ : price

$y = wx + w_0$

y: price

x: mileage

19

## Regression applications

How many clicks will a particular website, ad, etc. get?

Predict the readability level of a document

Predict pause between spoken sentences?

Economics/Finance: predict the value of a stock

Car/plane navigation: angle of the steering wheel, acceleration, …

…

20

## Supervised learning: ranking

label

1

4

2

3

Ranking: label is a ranking

Supervised learning: given labeled examples

21

## NLP Ranking Applications

reranking N-best output lists (e.g. parsing, machine translation, …)

Rank possible simplification options

flight search (search in general)

…

22

## Ranking example

Given a query and a set of web pages, rank them according to relevance

23

## Unsupervised learning

Unsupervised learning: given data, i.e. examples, but no labels

24

## Unsupervised learning applications

learn clusters/groups without any label

- cluster documents
- cluster words (synonyms, parts of speech, …)

compression

bioinformatics: learn motifs

…

25

## Reinforcement learning

| | |
|---|---|
| left, right, straight, left, left, left, straight | GOOD |
| left, straight, straight, left, right, straight, straight | BAD |
| left, right, straight, left, left, left, straight | 18.5 |
| left, straight, straight, left, right, straight, straight | -3 |

Given a *sequence* of examples/states and a *reward* after completing that sequence, learn to predict the action to take in for an individual example/state

26

## Reinforcement learning example

Backgammon



WIN!

LOSE!

Given sequences of moves and whether or not the player won at the end, learn to make good moves

27

## Reinforcement learning example

https://www.youtube.com/watch?v=tXlM99xPQC8

28

## Other learning variations

What data is available:
- Supervised, unsupervised, reinforcement learning
- semi-supervised, active learning, …

How are we getting the data:
- online vs. offline learning

Type of model:
- generative vs. discriminative
- parametric vs. non-parametric

29

## Text classification

label

spam

For this class, I'm mostly going to focus on classification

not spam

I'll use text classification as a running example

not spam

30

## Representing examples

examples

What is an example?
How is it represented?

31

## Features

examples        features

$f_1, f_2, f_3, …, f_n$

$f_1, f_2, f_3, …, f_n$

$f_1, f_2, f_3, …, f_n$

$f_1, f_2, f_3, …, f_n$

How our algorithms actually "view" the data

Features are the questions we can ask about the examples

32

## Features

examples

features

red, round, leaf, 3oz, …

green, round, no leaf, 4oz, …

yellow, curved, no leaf, 4oz, …

green, curved, no leaf, 5oz, …

How our algorithms actually "view" the data

Features are the questions we can ask about the examples

33

## Text: raw data

Raw data

Features?

34

## Feature examples

Raw data

Features

Clinton said banana repeatedly last week on tv, "banana, banana, banana"

(1, 1, 1, 0, 0, 1, 0, 0, ...)

banana clinton said california across tv wrong capital

Occurrence of words (unigrams)

35

## Feature examples

Raw data

Features

Clinton said banana repeatedly last week on tv, "banana, banana, banana"

(4, 1, 1, 0, 0, 1, 0, 0, ...)

banana clinton said california across tv wrong capital

Frequency of word occurrence (unigram frequency)

36

## Feature examples

Raw data

Features

Clinton said banana
repeatedly last week on tv,
"banana, banana, banana"

(1, 1, 1, 0, 0, 1, 0, 0, ...)

banana repeatedly
clinton said
said banana
california schools
across the
tv banana
wrong way
capital city

Occurrence of bigrams

37

## Feature examples

Raw data

Features

Clinton said banana
repeatedly last week on tv,
"banana, banana, banana"

(1, 1, 1, 0, 0, 1, 0, 0, ...)

banana repeatedly
clinton said
said banana
california schools
across the
tv banana
wrong way
capital city

Other features?

38

## Lots of other features

POS: occurrence, counts, sequence

Constituents

Whether 'V1agra' occurred 15 times

Whether 'banana' occurred more times than 'apple'

If the document has a number in it

…

Features are very important, but we're going to focus
on the model

39

## Classification revisited

examples          label

red, round, leaf, 3oz, …        apple

green, round, no leaf, 4oz, …    apple     learn

yellow, curved, no leaf, 4oz, …   banana          model/
classifier

green, curved, no leaf, 5oz, …    banana

During learning/training/induction, learn a model of what
distinguishes apples and bananas *based on the features*

40

## Classification revisited

red, round, no leaf, 4oz, … ➡ model/classifier ➡ *predict* **Apple or banana?**

The model can then classify a new example *based on the features*

41

## Classification revisited

red, round, no leaf, 4oz, … ➡ model/classifier ➡ *predict* **Apple**

**Why?**

The model can then classify a new example *based on the features*

42

## Classification revisited

| Training data | | Test set |
|---|---|---|
| examples | label | |
| red, round, leaf, 3oz, … | apple | |
| green, round, no leaf, 4oz, … | apple | red, round, no leaf, 4oz, … **?** |
| yellow, curved, no leaf, 4oz, … | banana | |
| green, curved, no leaf, 5oz, … | banana | |

43

## Classification revisited

| Training data | | Test set |
|---|---|---|
| examples | label | |
| red, round, leaf, 3oz, … | apple | |
| green, round, no leaf, 4oz, … | apple | red, round, no leaf, 4oz, … **?** |
| yellow, curved, no leaf, 4oz, … | banana | |
| green, curved, no leaf, 5oz, … | banana | |

Learning is about **generalizing** from the training data

What does this assume about the training and test set?

44

## Past predicts future

Training data          Test set



45

## Past predicts future

Training data          Test set



**Not always the case, but we'll often assume it is!**

46

## Past predicts future

Training data          Test set



**Not always the case, but we'll often assume it is!**

47

## More technically…

We are going to use the *probabilistic model* of learning

There is some probability distribution over example/label pairs called the *data generating distribution*

**Both** the training data **and** the test set are generated based on this distribution

48

## data generating distribution

Training data          Test set

data generating distribution

49

## data generating distribution

Training data          Test set

data generating distribution

50

## data generating distribution

Training data          Test set

data generating distribution

51

## 5 weeks left

Now that you know more about NLP, anything left that you'd like to know more about?

Summer plans?

Write the name of someone in the class that has their birthday later in the year than you?

52

## Probabilistic Modeling

training data → train → probabilistic model

Model the data with a probabilistic model

specifically, learn p(*features, label*)

p(*features, label*) tells us how likely these features and this example are

## An example: classifying fruit

Training data

| examples | label |
|---|---|
| red, round, leaf, 3oz, … | apple |
| green, round, no leaf, 4oz, … | apple |
| yellow, curved, no leaf, 4oz, … | banana |
| green, curved, no leaf, 5oz, … | banana |

train → probabilistic model: p(*features, label*)

## Probabilistic models

Probabilistic models define a *probability distribution* over features and labels:

yellow, curved, no leaf, 6oz, banana → probabilistic model: p(*features, label*) → 0.004

## Probabilistic model vs. classifier

Probabilistic model:

yellow, curved, no leaf, 6oz, banana → probabilistic model: p(*features, label*) → 0.004

Classifier:

yellow, curved, no leaf, 6oz → probabilistic model: p(*features, label*) → banana

## Probabilistic models: classification

Probabilistic models define a *probability distribution* over features and labels:

yellow, curved, no leaf, 6oz, banana → probabilistic model: p(*features, label*) → 0.004

Given an unlabeled example: yellow, curved, no leaf, 6oz predict the label

**How do we use a probabilistic model for classification/prediction?**

57

## Probabilistic models

Probabilistic models define a *probability distribution* over features and labels:

yellow, curved, no leaf, 6oz, banana → probabilistic model: p(*features, label*) → **0.004**

yellow, curved, no leaf, 6oz, apple → 0.00002

For each label, ask for the probability under the model
Pick the label with the highest probability

58

## Probabilistic model vs. classifier

Probabilistic model:

yellow, curved, no leaf, 6oz, banana → probabilistic model: p(*features, label*) → 0.004

Classifier:

yellow, curved, no leaf, 6oz → probabilistic model: p(*features, label*) → banana

**Why probabilistic models?**

59

## Probabilistic models

Probabilities are nice to work with
- range between 0 and 1
- can combine them in a well understood way
- lots of mathematical background/theory

Provide a strong, well-founded groundwork
- Allow us to make clear decisions about things like smoothing
- Tend to be much less "heuristic"
- Models have very clear meanings

60

## Probabilistic models: big questions

1. Which model do we use, i.e. how do we calculate p(*feature, label*)?

2. How do train the model, i.e. how to we we estimate the probabilities for the model?

3. How do we deal with overfitting (i.e. smoothing)?

61

## Basic steps for probabilistic modeling

Step 1: pick a model

Step 2: figure out how to estimate the probabilities for the model

Step 3 (optional): deal with overfitting

### Probabilistic models

Which model do we use, i.e. how do we calculate p(*feature, label*)?

How do train the model, i.e. how to we we estimate the probabilities for the model?

How do we deal with overfitting?

62

## What was the data generating distribution?

Training data          Test set

data generating distribution

63

## Step 1: picking a model

What we're really trying to do is model the data generating distribution, that is how likely the feature/label combinations are

data generating distribution

64

## Some math

$$p(features, label) = p(x_1, x_2, ..., x_m, y)$$

$$= p(y)p(x_1, x_2, ..., x_m \mid y)$$

What rule?

65

## Some math

$$p(features, label) = p(x_1, x_2, ..., x_m, y)$$

$$= p(y)p(x_1, x_2, ..., x_m \mid y)$$

$$= p(y)p(x_1 \mid y)p(x_2, ..., x_m \mid y, x_1)$$

$$= p(y)p(x_1 \mid y)p(x_2 \mid y, x_1)p(x_3, ..., x_m \mid y, x_1, x_2)$$

$$= p(y)\prod_{j=1}^{m} p(x_i \mid y, x_1, ..., x_{i-1})$$

66

## Step 1: pick a model

$$p(features, label) = p(y)\prod_{j=1}^{m} p(x_i \mid y, x_1, ..., x_{i-1})$$

So, far we have made NO assumptions about the data

$$p(x_m \mid y, x_1, x_2, ..., x_{m-1})$$

How many entries would the probability distribution table have if we tried to represent all possible values and we had 7000 binary features?

67

## Full distribution tables

| $x_1$ | $x_2$ | $x_3$ | ... | y | p( ) |
|---|---|---|---|---|---|
| 0 | 0 | 0 | ... | 0 | * |
| 0 | 0 | 0 | ... | 1 | * |
| 1 | 0 | 0 | ... | 0 | * |
| 1 | 0 | 0 | ... | 1 | * |
| 0 | 1 | 0 | ... | 0 | * |
| 0 | 1 | 0 | ... | 1 | * |
| | | | ... | | |

All possible combination of features!

Table size: $2^{7000} = $ ?

68

17

# $2^{7000}$

162169675566220202646666508547837709519111124303637432562359820841515270231627023529870802378794460004651996019099530984538652557892546513204107022110253564658647431585227076599373340842842722420012281878260072931082617043194484266392077784125099999686016943600666001120981757929667878196255237700655294757256678055809293844627218640216108862600816097132874749204352087401101862690842327501724605231129395523505905454421455477250950909650788947809468359293957411256947343861912152968484743444406741204174020887540371869421701550220735398381224299258743537536161041593435945576665617017909041725970253365266626820218084938928126997095285708906963755754143448760882483699419938024151975145101251270438290872809195384763028578118540240999588959641922776012553604911562403499947144160905730842429313962119953679373012944795600248333570738998392029910322346598038953069042980174009801732521069130797124201696339723021835300758978451952584855371088581956317370007438051674111891346175014845217679842967828422873731274221220225175975359948392570298779077063553347902449354353866605125910795672914312162977887848185522928196541766009803989979916814047493842157435158026038115106828640678973048382922034604277576550737765675475070271446622634876857096212610747627052030494889072089785933689047063428548531668665653732717466065818560906648495080127617546145721617695557519921175075140677751044967285908225585477714472423349007640263217608921135525612411945387026802990440018385850576719369689759366121356888838680023840932567380777501891470304962150996983853975207154939633923720287592041517294937079097785362510832009283960480723795488706954662168804465211249307629009199071774235503913511744153297374793008995583051888413533479846411368000499940373724560035428811232632821866113106455077289922996946915601858083982074170460683212438815202609958469658816137582638292102954734388888321636271223029212297953848683554835357106034077891774170263636562027269554375177807413134551018100094688094078112205738033537112463295891623708958047622459509182530163690923624067141164433165615982805837207834398885623908920284409025538293376

Any problems with this?

---

# Full distribution tables

| $x_1$ | $x_2$ | $x_3$ | ... | $y$ | $p(\ )$ |
|---|---|---|---|---|---|
| 0 | 0 | 0 | ... | 0 | * |
| 0 | 0 | 0 | ... | 1 | * |
| 1 | 0 | 0 | ... | 0 | * |
| 1 | 0 | 0 | ... | 1 | * |
| 0 | 1 | 0 | ... | 0 | * |
| 0 | 1 | 0 | ... | 1 | * |
|  |  |  | ... |  |  |

- Storing a table of that size is impossible!
- How are we supposed to learn/estimate each entry in the table?

---

# Step 1: pick a model

$$p(features, label) = p(y)\prod_{j=1}^{m} p(x_i \mid y, x_1, ..., x_{i-1})$$

So, far we have made NO assumptions about the data

Model selection involves making assumptions about the data

We've done this before, n-gram language model, parsing, etc.

These assumptions allow us to represent the data more compactly and to estimate the parameters of the model

---

# Naïve Bayes assumption

$$p(features, label) = p(y)\prod_{j=1}^{m} p(x_i \mid y, x_1, ..., x_{i-1})$$

$$p(x_i \mid y, x_1, x_2, ..., x_{i-1}) = p(x_i \mid y)$$

What does this assume?

### Naïve Bayes assumption

$$p(features, label) = p(y) \prod_{j=1}^{m} p(x_i \mid y, x_1, ..., x_{i-1})$$

$$p(x_i \mid y, x_1, x_2, ..., x_{i-1}) = p(x_i \mid y)$$

Assumes feature i is independent of the the other features *given the label*

Is this true for text, say, with unigram features?

73

### Naïve Bayes assumption

$$p(x_i \mid y, x_1, x_2, ..., x_{i-1}) = p(x_i \mid y)$$

For most applications, this is not true!

For example, the fact that "San" occurs will probably make it *more likely* that "Francisco" occurs

However, this is often a reasonable approximation:

$$p(x_i \mid y, x_1, x_2, ..., x_{i-1}) \approx p(x_i \mid y)$$

74