10-301/601: Introduction to Machine Learning Lecture 2 – ML as Function Approximation

Henry Chai & Matt Gormley 8/31/22

Q & A

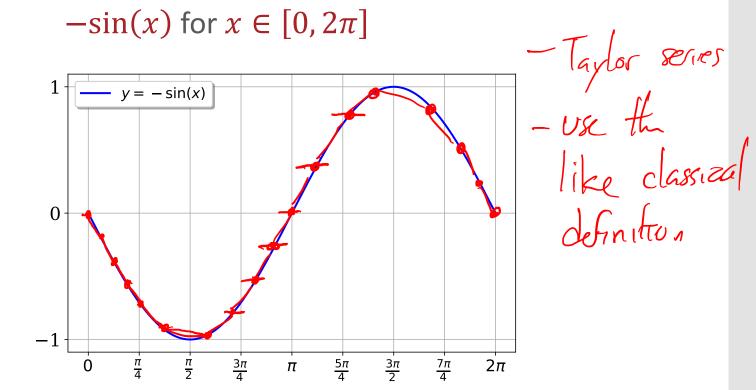
- In Lecture 1, why did we use the term experience instead of just data?
- Because our concern isn't just the data itself, but also where the data comes from (e.g., an agent interacting with the world vs. knowledge from a book). As well, the word experience better aligns with the notion of what humans require in order to learn.

Front Matter

- Announcements:
 - HW1 released 8/29, due 9/7 at 11:59 PM
 - Two components: written and programming
 - Separate submissions on Gradescope
 - Unique policies specific to HW1:
 - Two submissions for the written portion (see writeup for details)
 - Unlimited submissions for the programming portion (really, just keep submitting until you get 100%)
 - We will grant (almost) any extension request

Function Approximation: Example

Challenge: implement a function that computes

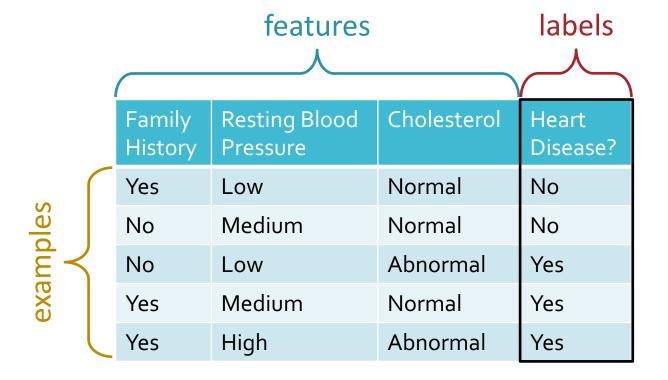


- You may not call any trigonometric functions
- You may call an existing implementation of sin(x) a few times (e.g., 100) to check your work

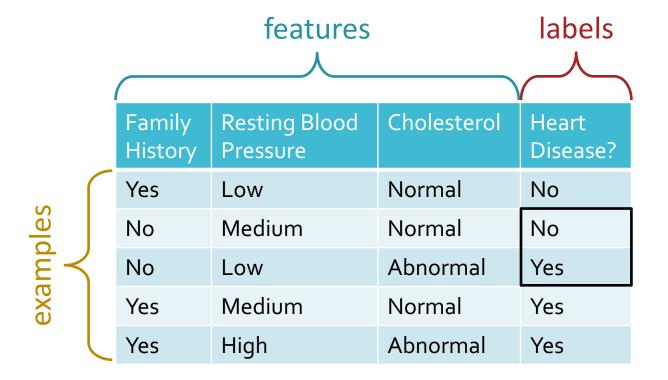
Learning to diagnose heart disease
 as a (supervised) binary classification task

	features			labels	
		Family History	Resting Blood Pressure	Cholesterol	Heart Disease?
examples		Yes	Low	Normal	No
		No	Medium	Normal	No
	<i>,</i>	No	Low	Abnormal	Yes
еха		Yes	Medium	Normal	Yes
		Yes	High	Abnormal	Yes
eXe					

Learning to diagnose heart disease
 as a (<u>supervised</u>) binary classification task



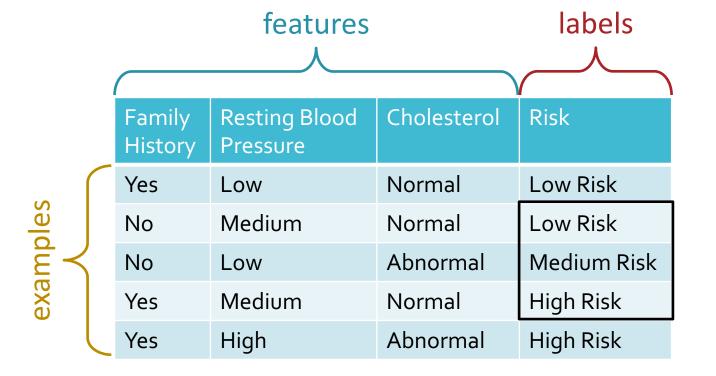
Learning to diagnose heart disease
 as a (supervised) binary classification task



Learning to diagnose heart disease

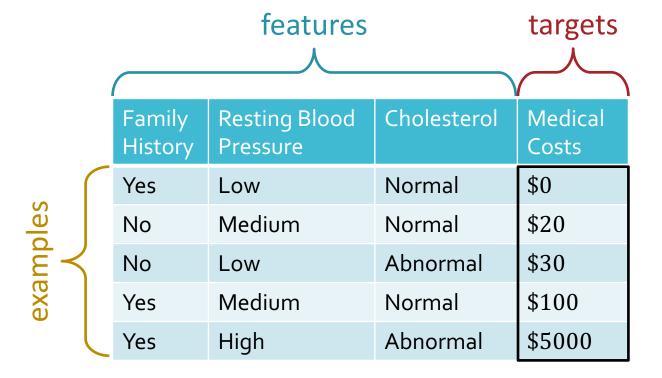
as a (supervised)

classification task



Learning to diagnose heart disease

as a (supervised) regression task



- A classifier is a function that takes feature values as input and outputs a label
- Majority vote classifier: always predict the most common label in the training dataset

features			labels
Familia	Deaties Diesel	Chalastanal	
History	Pressure	Cholesterol	Heart Disease?
Yes	Low	Normal	No
No	Medium	Normal	No
No	Low	Abnormal	Yes
Yes	Medium	Normal	Yes
Yes	High	Abnormal	Yes
	Yes No No Yes	Family Resting Blood Pressure Yes Low No Medium No Low Yes Medium	Family Resting Blood Cholesterol Pressure Yes Low Normal No Medium Normal No Low Abnormal Yes Medium Normal

- A classifier is a function that takes feature values as input and outputs a label
- Majority vote classifier: always predict the most common label in the training dataset

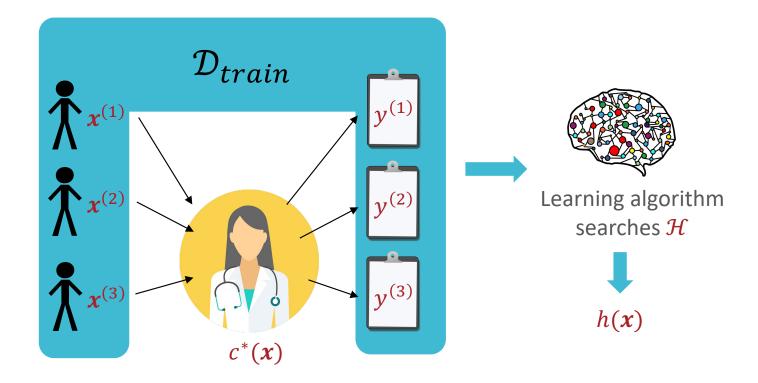
		features			labels	
	ı					
		Family History	Resting Blood Pressure	Cholesterol	Heart Disease?	Predictions
		Yes	Low	Normal	No	Yes
<u> es</u>		No	Medium	Normal	No	Yes
examples \(\lambda\)	,	No	Low	Abnormal	Yes	Yes
exa		Yes	Medium	Normal	Yes	Yes
		Yes	High	Abnormal	Yes	Yes

Example:-Sin(x)

- Feature space, $x = \begin{bmatrix} 0,2\pi \end{bmatrix}$
- · Label space, y R, (-1,1) function definition
- (Unknown) Target function, $c^*: \mathcal{X} \to \mathcal{Y}$
- Training dataset: calling the observe implementation $\mathcal{D} = \{ (x^{(1)}, c^*(x^{(1)}) = y^{(1)}), (x^{(2)}, y^{(2)}) \dots, (x^{(N)}, y^{(N)}) \}$
- Example: $(x^{(n)}, y^{(n)}) = (x_1^{(n)}, x_2^{(n)}, \dots, x_D^{(n)}, y^{(n)})$
- · Hypothesis space: H

 precewise linear functions step functions
- Goal: find a classifier, $h \in \mathcal{H}$, that best approximates c^*

Notation



 Majority vote classifier: always predict the most common label in the training dataset

	Family History	Resting Blood Pressure	Cholesterol	Heart Disease?	Predictions
	Yes	Low	Normal	No	Yes
$\boldsymbol{x}^{(2)}$	No	Medium	Normal	No	Yes
•	No	Low	Abnormal	Yes	Yes
	Yes	Medium	Normal	Yes	Yes
	Yes	High	Abnormal	Yes	Yes

•
$$N = 5$$
 and $D = 3$

•
$$\mathbf{x}^{(2)} = (x_1^{(2)} = \text{"No"}, x_2^{(2)} = \text{"Medium"}, x_3^{(2)} = \text{"Normal"})$$

Evaluation

- Loss function, $\ell: \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$
 - Defines how "bad" predictions, $\hat{y} = h(x)$, are compared to the true labels, $y = c^*(x)$
 - Common choices

absolute loss: /y-//

- 1. Squared loss (for regression): $\ell(y, \hat{y}) = (y \hat{y})^{\frac{1}{2}}$
- 2. Binary or 0-1 loss (for classification):

$$\ell(y, \hat{y}) = \begin{cases} 1 & \text{if } y \neq \hat{y} \\ 0 & \text{otherwise} \end{cases}$$

Evaluation

- Loss function, $\ell: \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$
 - Defines how "bad" predictions, $\hat{y} = h(x)$, are compared to the true labels, $y = c^*(x)$
 - Common choices
 - 1. Squared loss (for regression): $\ell(y, \hat{y}) = (y \hat{y})^2$
 - 2. Binary or 0-1 loss (for classification):

$$\ell(y, \hat{y}) = 1(y \neq \hat{y})$$
indicator function

• Error rate:

$$err(h, \mathcal{D}) = \frac{1}{N} \sum_{n=1}^{N} \mathbb{1}(y^{(n)} \neq \hat{y}^{(n)})$$

Different Kinds of Error

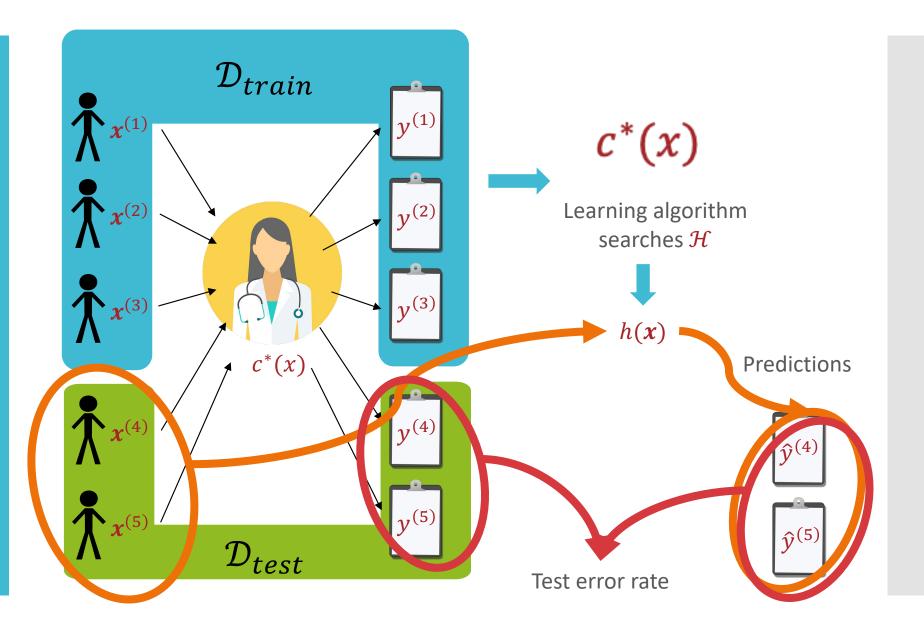
• Training error rate = $err(h, \mathcal{D}_{train})$

• Test error rate = $err(h, \mathcal{D}_{test})$

• True error rate = err(h)

= the error rate of h on all possible examples

• In machine learning, this is the quantity that we care about but, in most cases, it is unknowable.



Majority vote classifier:

def train (D_{train}).

store
$$V = mode(x^{(1)}, y^{(2)}, ..., y^{(N)})$$

def predict(x'):

return V

Test your understanding

x_1	x_2	y
1	0	-
1	0	-
1	0	+
1	0	+
1	1	+
1	1	+
1	1	+
1	1	+

 What is the training error of the majority vote
 classifier on this dataset?

- A classifier is a function that takes feature values as input and outputs a label
- Majority vote classifier: always predict the most common label in the training dataset



labels	ı
Heart Disease?	Predictions
No	Yes
No	Yes
Yes	Yes
Yes	Yes
Yes	Yes

ماما

• This classifier completely ignores the features...

 Memorizer: if a set of features exists in the training dataset, predict its corresponding label; otherwise, predict a random label.

Family History	Resting Blood Pressure	Cholesterol	Heart Disease?
Yes	Low	Normal	No
No	Medium	Normal	No
No	Low	Abnormal	Yes
Yes	Medium	Normal	Yes
Yes	High	Abnormal	Yes

Memorizer:

def train (Dtrain):

store Dtrain

def predict (x'):

if
$$\exists x^{(n)} \in D_{train} \text{ s.t. } x^{(n)} == x!$$

return $y^{(n)}$

else:

return $y \in Y$ at random

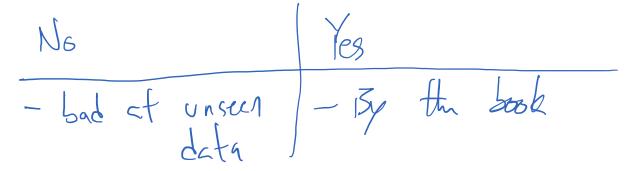
 Memorizer: if a set of features exists in the training dataset, predict its corresponding label; otherwise, predict a random label.

Family History	Resting Blood Pressure	Cholesterol	Heart Disease?	Predictions
Yes	Low	Normal	No	No
No	Medium	Normal	No	No
No	Low	Abnormal	Yes	Yes
Yes	Medium	Normal	Yes	Yes
Yes	High	Abnormal	Yes	Yes

The training error rate is always 0!

8/31/22 **24**

Is the Memorizer learning?



 Memorizer: if a set of features exists in the training dataset, predict its corresponding label; otherwise, predict a random label.

Family History	Resting Blood Pressure	Cholesterol	Heart Disease?	Predictions
Yes	Low	Normal	No	No
No	Medium	Normal	No	No
No	Low	Abnormal	Yes	Yes
Yes	Medium	Normal	Yes	Yes
Yes	High	Abnormal	Yes	Yes

• The training error rate is always 0!

- Memorizer: if a set of features exists in the training dataset, predict its corresponding label; otherwise, predict a random label.
- The memorizer (typically) does not **generalize** well, i.e., it does not perform well on unseen data points
- In some sense, good generalization, i.e., the ability to make accurate predictions given a small training dataset, is the whole point of machine learning!

Alright, let's actually (try to) extract a pattern from the data

x_1 Family History	x_2 Resting Blood Pressure	x_3 Cholesterol	<i>y</i> Heart Disease?
Yes	Low	Normal	No
No	Medium	Normal	No
No	Low	Abnormal	Yes
Yes	Medium	Normal	Yes
Yes	High	Abnormal	Yes

• Decision stump: based on a single feature, x_d , predict the most common label in the **training** dataset among all data points that have the same value for x_d

Alright, let's actually (try to) extract a pattern from the data

x_1 Family History	x_2 Resting Blood Pressure	x_3 Cholesterol	<i>y</i> Heart Disease?
Yes	Low	Normal	No
No	Medium	Normal	No
No	Low	Abnormal	Yes
Yes	Medium	Normal	Yes
Yes	High	Abnormal	Yes

• Decision stump on x_1 :

$$h(x') = h(x'_1, ..., x'_D) = \begin{cases} ??? & \text{if } x'_1 = \text{"Yes"} \\ ??? & \text{otherwise} \end{cases}$$

Alright, let's actually (try to) extract a pattern from the data

x_1 Family History	x_2 Resting Blood Pressure	x_3 Cholesterol	<i>y</i> Heart Disease?
Yes	Low	Normal	No
No	Medium	Normal	No
No	Low	Abnormal	Yes
Yes	Medium	Normal	Yes
Yes	High	Abnormal	Yes

• Decision stump on x_1 :

$$h(x') = h(x'_1, ..., x'_D) = \begin{cases} \text{"Yes" if } x'_1 = \text{"Yes"} \\ ??? \text{ otherwise} \end{cases}$$

Alright, let's actually (try to) extract a pattern from the data

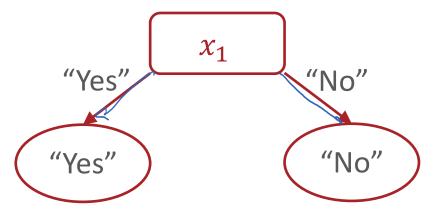
x_1 Family History	x_2 Resting Blood Pressure	x_3 Cholesterol	<i>y</i> Heart Disease?
Yes	Low	Normal	No
No	Medium	Normal	No
No	Low	Abnormal	Yes
Yes	Medium	Normal	Yes
Yes	High	Abnormal	Yes

• Decision stump on x_1 :

$$h(\mathbf{x}') = h(x_1', \dots, x_D') = \begin{cases} \text{"Yes" if } x_1' = \text{"Yes"} \\ \text{"No" otherwise} \end{cases}$$

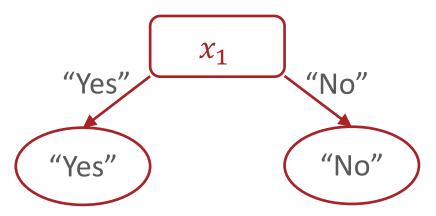
Alright, let's actually (try to) extract a pattern from the data

x_1 Family History	x_2 Resting Blood Pressure	x_3 Cholesterol	<i>y</i> Heart Disease?
Yes	Low	Normal	No
No	Medium	Normal	No
No	Low	Abnormal	Yes
Yes	Medium	Normal	Yes
Yes	High	Abnormal	Yes



• Alright, let's actually (try to) extract a pattern from the data

x_1 Family History	x_2 Resting Blood Pressure	x_3 Cholesterol	<i>y</i> Heart Disease?	\hat{y} Predictions
Yes	Low	Normal	No	Yes
No	Medium	Normal	No	No
No	Low	Abnormal	Yes	No
Yes	Medium	Normal	Yes	Yes
Yes	High	Abnormal	Yes	Yes



8/31/22 **32**

Decision Stumps: Pseudocode

def train (Dtrain):

1. prock a feature,
$$\chi_{3}$$
 (???)

2. split Dtrain according to χ_{3}

— For v in $V(\chi_{3})$, all possible values of $V(\chi_{3})$, all possible values of $V(\chi_{3})$, all possible values of $V(\chi_{3})$.

3. Compute some majority votes

— For V in $V(\chi_{3})$:

store $V(\chi_{3})$:

 $V(\chi_{3})$:

Decision Stumps: Questions

1. How can we pick which feature to split on?

2. Why stop at just one feature?