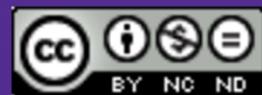


DS-GA 1003 Machine Learning

Week 1: Lecture 1
Model Selection - Clustering and Regression



How can machine learning
help us make guesses or
detect patterns?

DS-GA 1003 Machine Learning

Week 1: Lecture 1 Model Selection - Clustering and Regression

Adapted from Rosenberg, Rudin, Rangan, Steinhardt



What is machine learning?

Age.	Per-	Age.	Per-	Age.	Per-	Age.	Per-	Age.	Per-	Age.	Per-	Age.	Per-
Curt.	sons.	Curt.	sons										
1	1000	8	680	15	628	22	585	29	539	36	481	42	7
2	855	9	670	16	622	23	579	30	531	37	472	21	14
3	798	10	661	17	616	24	573	31	523	38	463	28	21
4	760	11	653	18	610	25	567	32	515	39	454	35	28
5	732	12	646	19	604	26	560	33	507	40	445	42	33
6	710	13	640	20	598	27	553	34	499	41	436	49	42
7	692	14	634	21	592	28	546	35	490	42	427	56	49
													3198
													2709
													2194
													1694
													1204
43	417	50	346	57	272	64	202	71	131	78	58	77	692
44	407	51	335	58	262	65	192	72	120	79	49	84	253
45	397	52	324	59	252	65	182	73	109	80	41	100	107
46	387	53	313	60	242	67	172	74	98	81	34		
47	377	54	302	61	232	68	162	75	88	82	28		34000
48	367	55	292	62	222	69	152	76	78	83	23		
49	357	56	282	63	212	70	142	77	68	84	20		Sum Total.

Based on data collected
between 1687 and 1691 in
Wroclaw, Poland

What is machine learning?

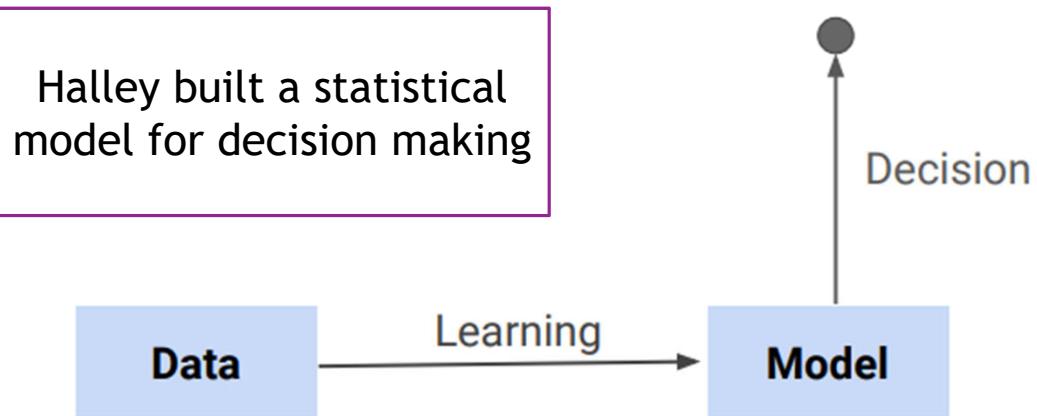
Age.	Per- sons.	Age.	Per- sons												
Curt.		Curt.		Curt.		Curt.		Curt.		Curt.		Curt.		Curt.	
1	1000	8	680	15	628	22	585	29	539	36	481	44	36	14	7
2	855	9	670	16	622	23	579	30	531	37	472	21	38	28	
3	798	10	661	17	616	24	573	31	523	38	463	33	454	35	
4	760	11	653	18	610	25	567	32	515	39	454	42	3198		
5	732	12	646	19	604	26	560	33	507	40	445	49	2709		
6	710	13	640	20	598	27	553	34	499	41	436	56	2194		
7	692	14	634	21	592	28	546	35	490	42	427	63	1694		
Age.	Per- sons.	Age.	Per- sons												
Curt.		Curt.		Curt.		Curt.		Curt.		Curt.		Curt.		Curt.	
43	417	50	346	57	272	64	202	71	131	78	58	77	692		
44	407	51	335	58	262	65	192	72	120	79	49	84	253		
45	397	52	324	59	252	65	182	73	109	80	41	100	107		
46	387	53	313	60	242	66	172	74	98	81	34				
									88	82	28		34000		
									78	83	23				
									68	84	20		Sum Total.		

Based on data collected between 1687 and 1691 in Wroclaw, Poland

$$\text{Price at age } x = \sum_i p[\text{death at age } x+i] 0.95^i \text{ (annual payout)}$$

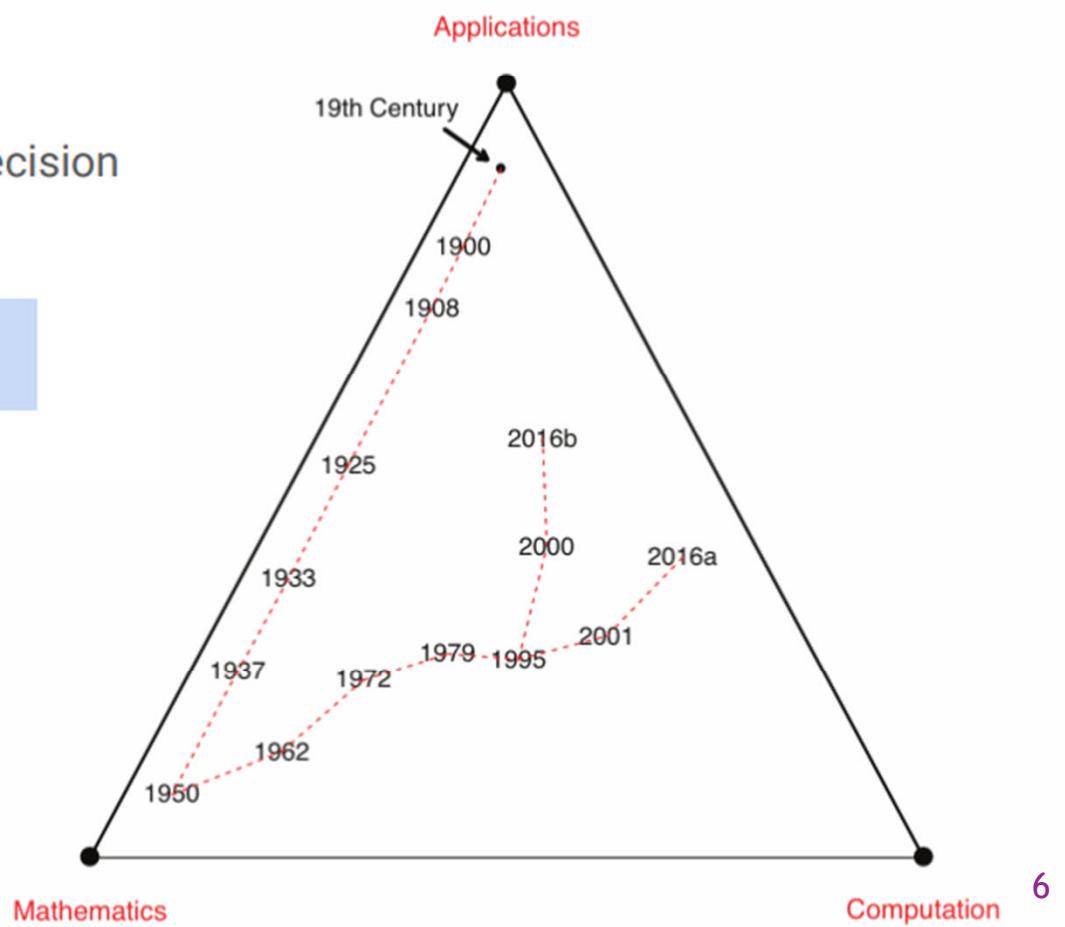
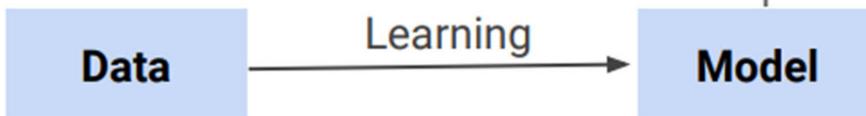
Halley's life expectancy model

What is machine learning?



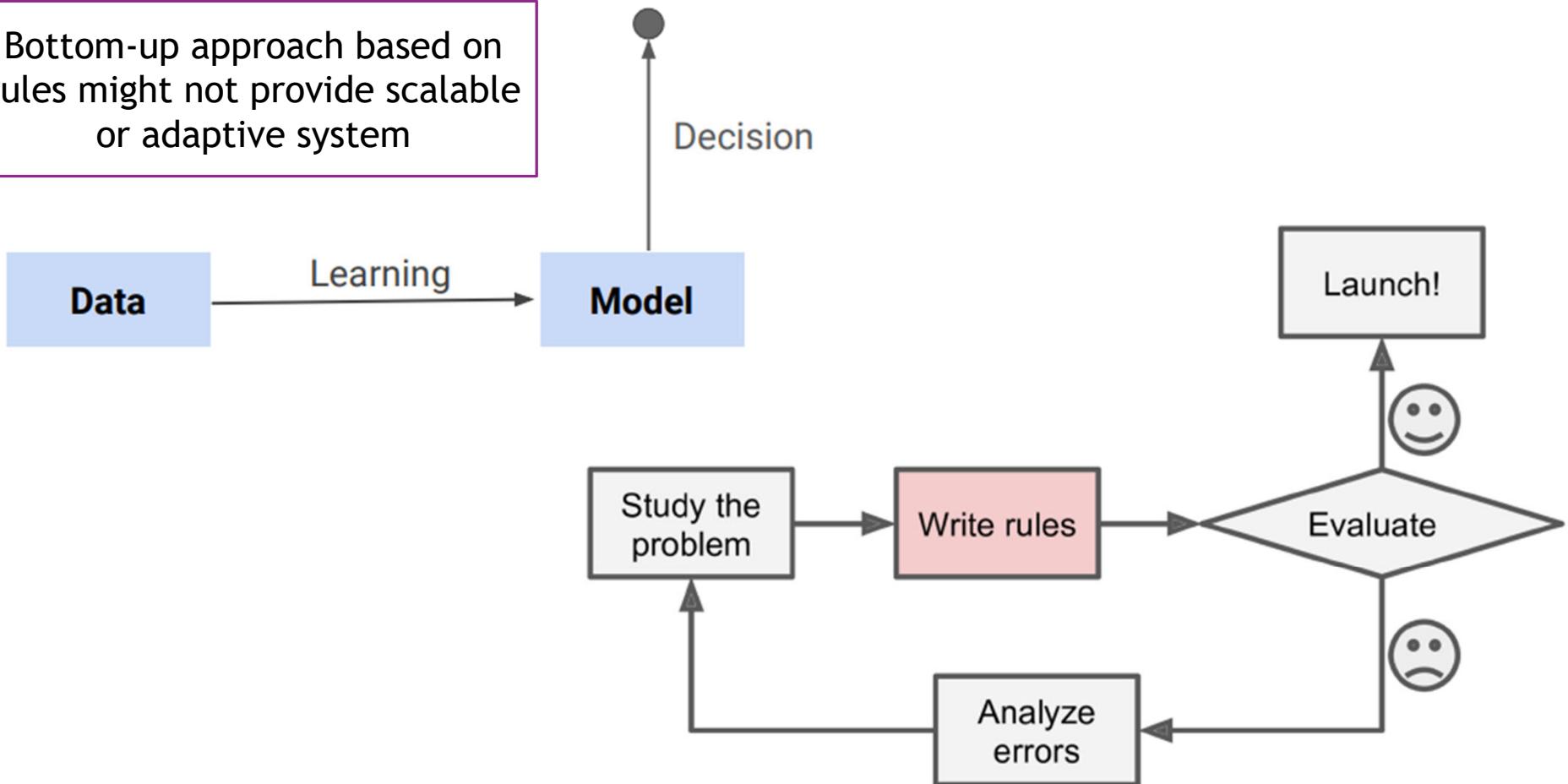
What is machine learning?

Advances in computing
and storage informed the
methods

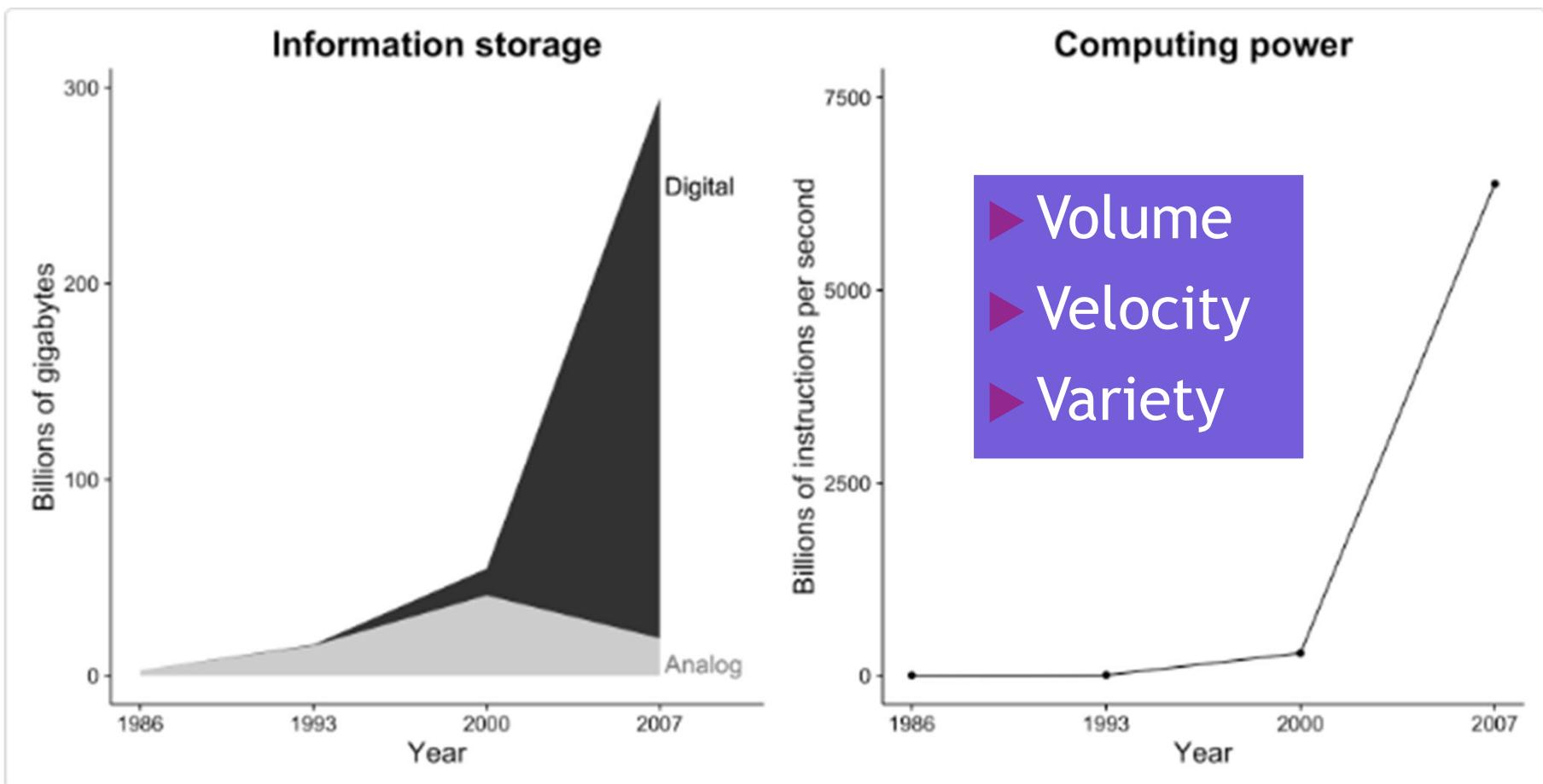


What is machine learning?

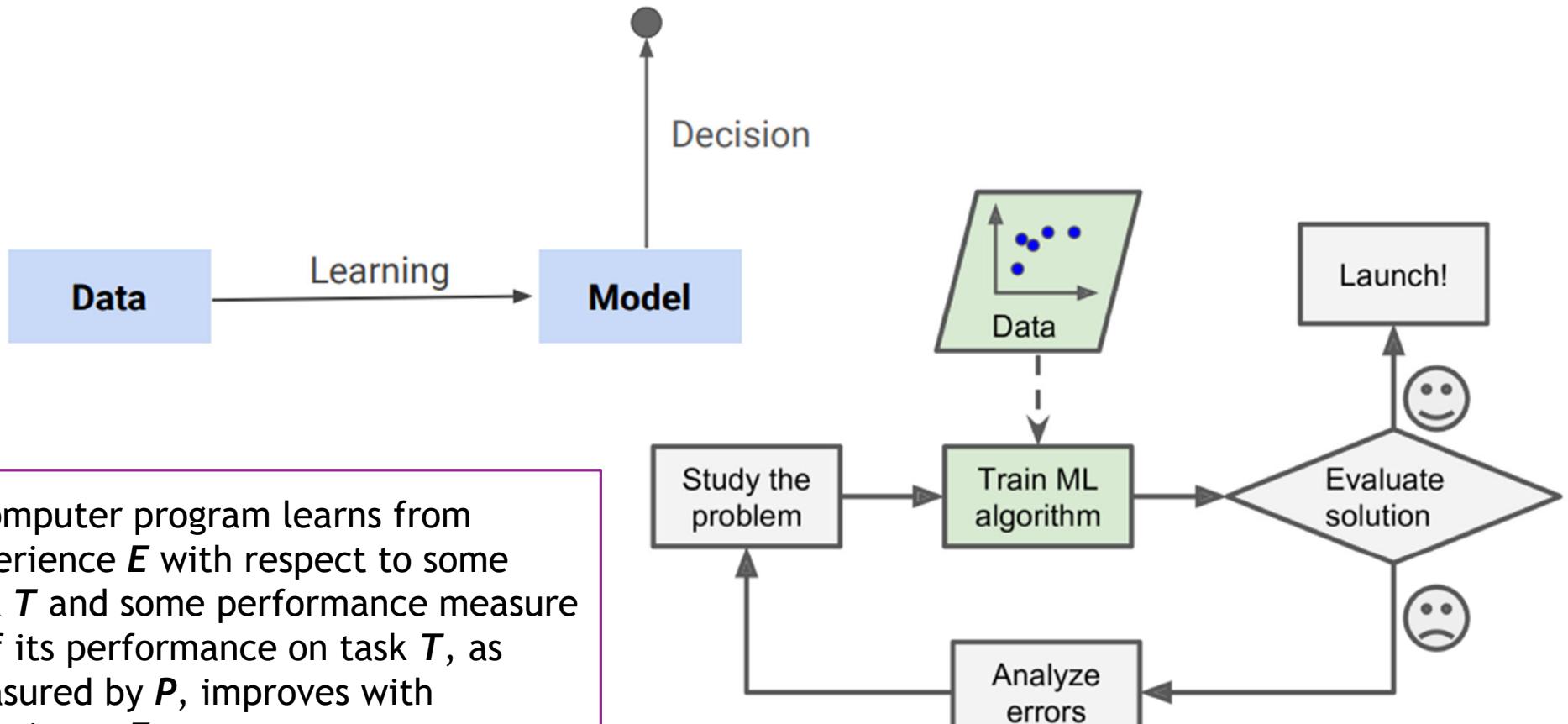
Bottom-up approach based on rules might not provide scalable or adaptive system



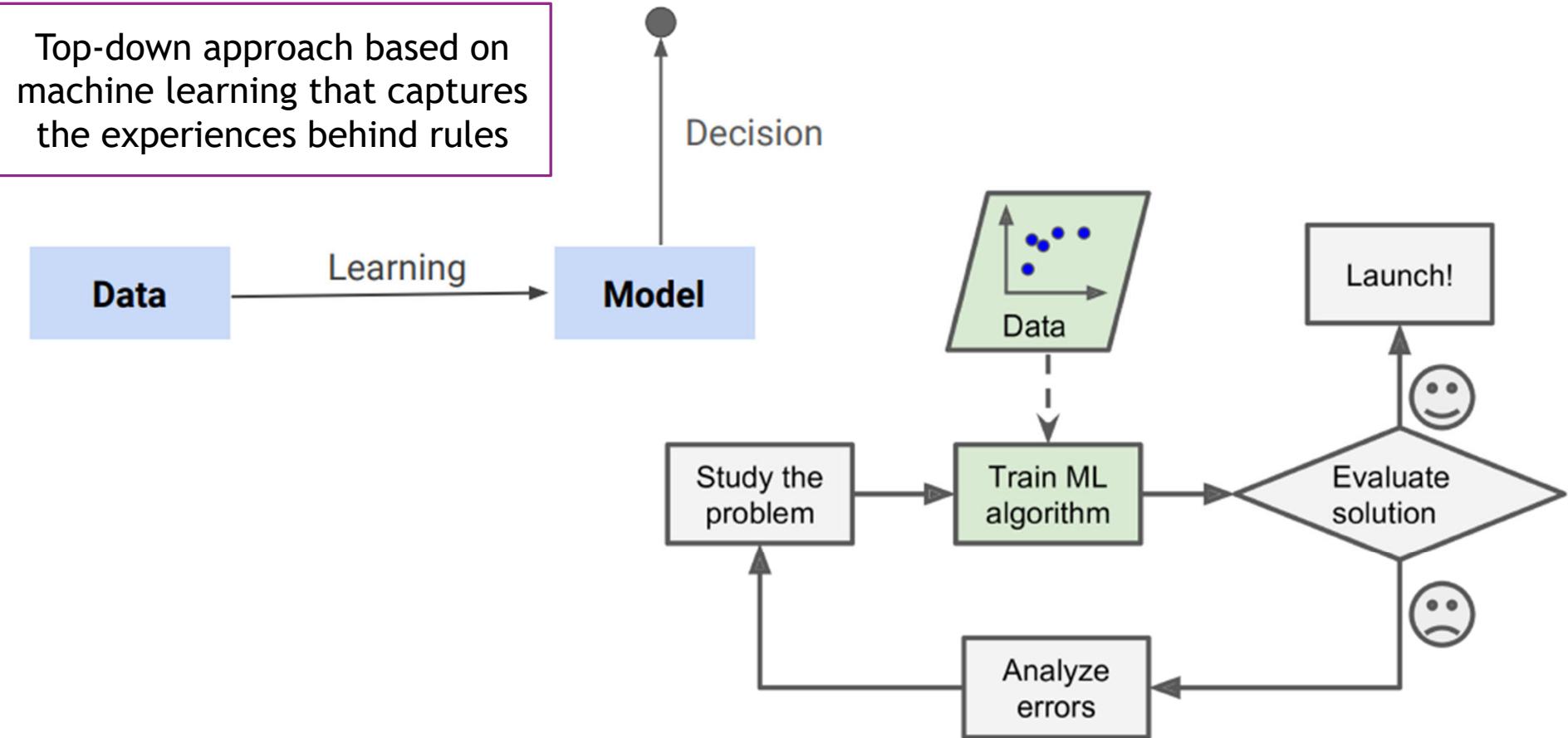
What is machine learning?



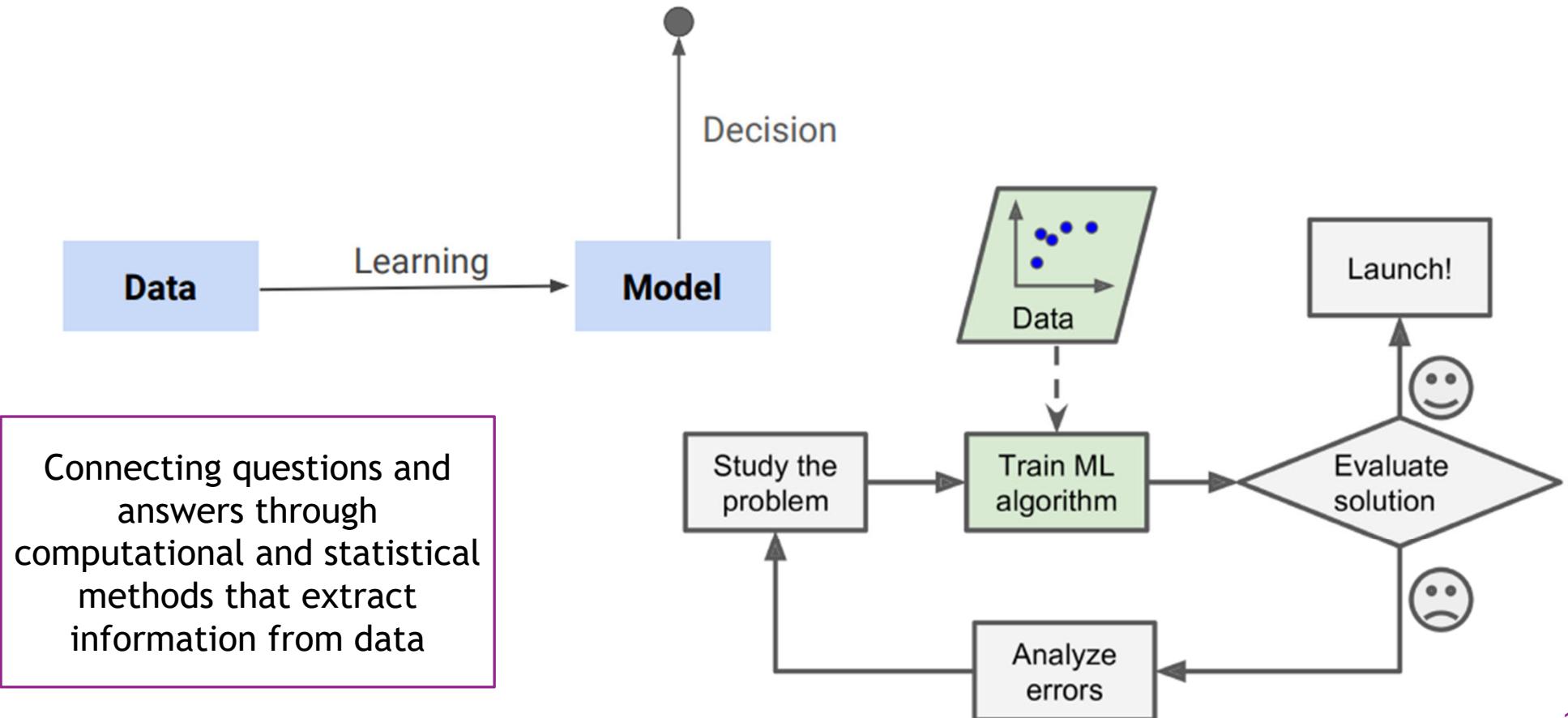
What is machine learning?



What is machine learning?



What is machine learning?



What is machine learning?

	id	subject	email	spam
0	0	Subject: A&L Daily to be auctioned in bankrupt...	url: http://boingboing.net/#85534171\n date: n...	0
1	1	Subject: Wired: "Stronger ties between ISPs an...	url: http://scriptingnews.userland.com/backiss...	0
2	2	Subject: It's just too small ...	<html>\n <head>\n </head>\n <body>\n <font siz...	1
3	3	Subject: liberal definitions\n	depends on how much over spending vs. how much...	0
4	4	Subject: RE: [ILUG] Newbie seeks advice - Suse...	hehe sorry but if you hit caps lock twice the ...	0

url: http://boingboing.net/#85534171
date: not supplied

arts and letters daily, a wonderful and dense blog, has folded up its tent due to the bankruptcy of its parent company. a&l daily will be auctioned off by the receivers. link[1] discuss[2] (_thanks, misha!_)

[1] <http://www.aldaily.com/>
[2] <http://www.quicktopic.com/boing/h/zlftejnd6jf>

Homework 1

What is machine learning?

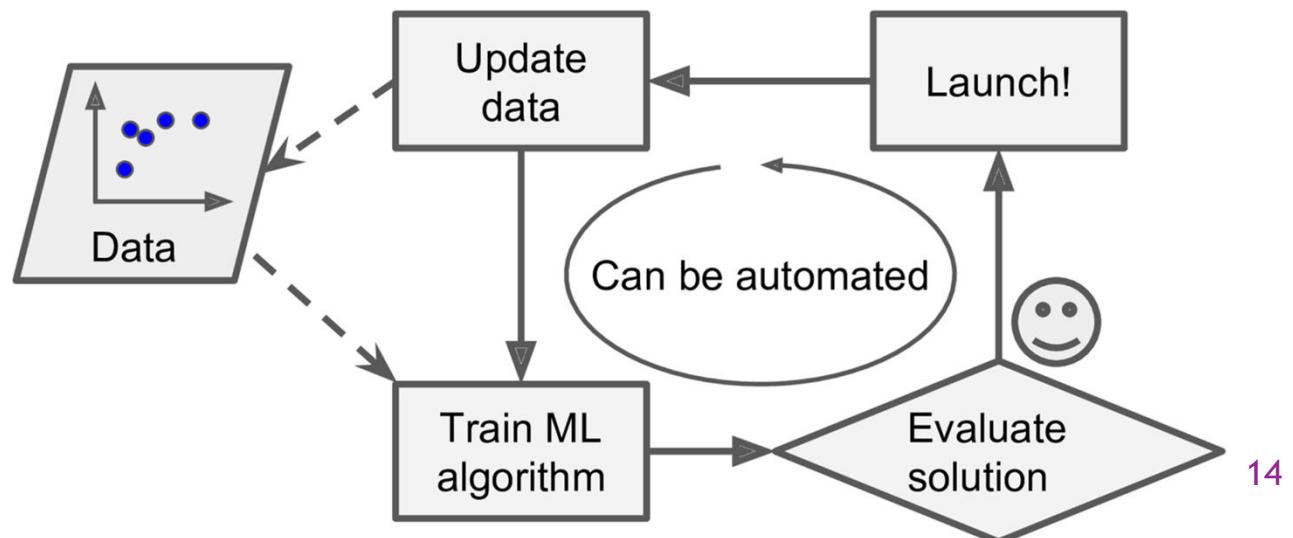
	id	subject	email	spam
0	0	Subject: A&L Daily to be auctioned in bankrupt...	url: http://boingboing.net/#85534171\n date: n...	0
1	1	Subject: Wired: "Stronger ties between ISPs an...	url: http://scriptingnews.userland.com/backiss...	0
2	2	Subject: It's just too small ...	<html>\n <head>\n </head>\n <body>\n <font siz...	1
3	3	Subject: liberal definitions\n	depends on how much over spending vs. how much...	0
4	4	Subject: RE: [ILUG] Newbie seeks advice - Suse...	hehe sorry but if you hit caps lock twice the ...	0

Homework 1

```
dear ricardo1 ,\n\n<html>\n<body>\n<center>\n<b><font color = "red" size = "+2.5">cost effective direct email advertising</font><br>\n<font color = "blue" size = "+2">promote your business for as low as </font><br>\n<font color = "red" size = "+2">$50</font> <font color = "blue" size = "+2">per<br>\n<font color = "red" size = "+2">1 million</font>\n<font color = "blue" size = "+2"> email addresses</font></font><p>\n<b><font color = "#44c300" size ="+2">maximize your marketing dollars!<p></font></b>\n<font size = "+2">complete and fax this information form to 309-407-7378.<br>\na consultant will contact you to discuss your marketing needs.<br>
```

What is machine learning?

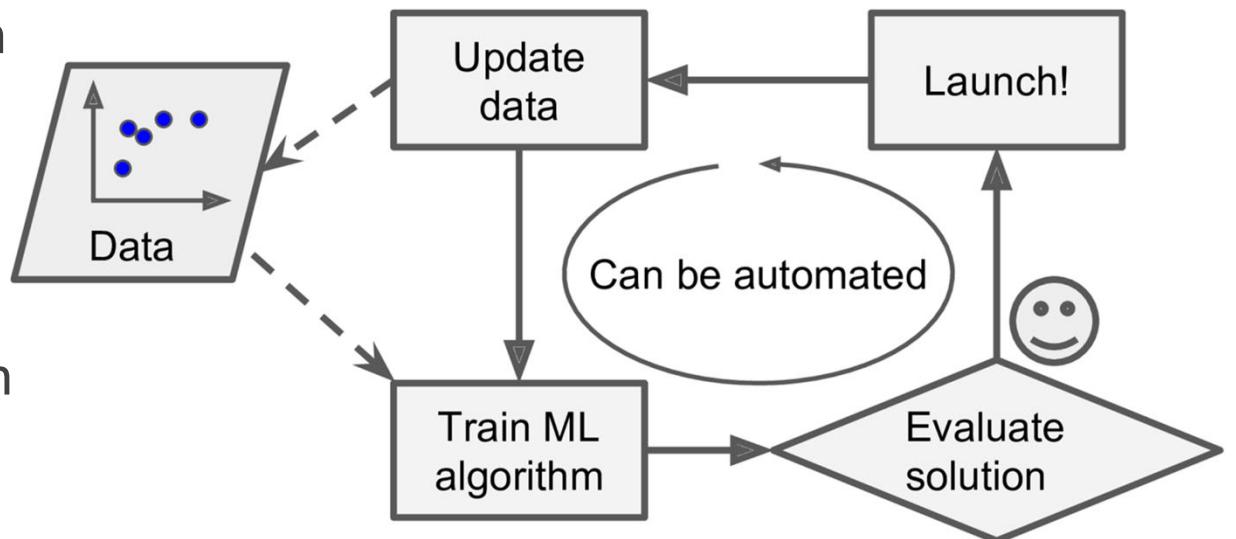
- ▶ Rules not Appropriate
- ▶ Possess Data
- ▶ Pattern in the Data



What is machine learning?

- ▶ Natural Language Processing
 - ▶ Text classification
 - ▶ Part of speech recognition
- ▶ Computer Vision
 - ▶ Character recognition
 - ▶ Image retrieval
- ▶ Speech Recognition
 - ▶ Source separation
 - ▶ Speaker Identification

- ▶ Rules not Appropriate
- ▶ Possess Data
- ▶ Pattern in the Data



What is machine learning?

Make Your Job Application Robot-Proof

It takes planning to make sure AI gatekeepers don't bounce your résumé—sometimes for arbitrary reasons—before a human can make a call

“Often a job candidate doesn’t even know a system is in use,” and employers aren’t required to disclose it, says Sarah Myers West, a researcher at the AI Now Institute, a New York University research group.



What is machine learning?

 ProPublica

Machine Bias

Machine Bias ... In 2014, then U.S. Attorney General Eric Holder warned that the risk scores might be injecting bias into the courts. ... So ProPublica did, as part of a larger examination of the powerful, largely hidden effect of May 23, 2016



A commercial tool COMPAS automatically predicts some categories of future crime to assist in bail and sentencing decisions. It is used in courts in the US.

What is machine learning?

IEEE Spectrum

Are Your Students Bored? This AI Could Tell You

Qu and his colleagues tested their AI system in two classrooms ... chief architect at Squirrel AI Learning, who was not involved in the paper.

2 weeks ago



The system “provides teachers with a quick and convenient measure of the students’ engagement level in a class,” says Huamin Qu, a computer scientist at the Hong Kong University of Science and Technology, who co-authored the paper. “Knowing whether the lectures are too hard and when students get bored can help improve teaching.”

What is machine learning?



Memorandum

Date: January 24, 2020

To: THE NYU COMMUNITY

From: Carlo Ciotoli, MD, Associate Vice President for Student Health and Executive Director of the Student Health Center

Re: The Emergence of the Novel Coronavirus

An AI Epidemiologist Sent the First Warnings of the Wuhan Virus

WIRED · 2 days ago



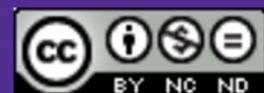
Logistics

For more information please refer to the syllabus on NYU Classes.

Please complete **Homework 0** and **Survey 1** linked to the Weekly Agenda on NYU Classes



A word cloud graphic centered on data science, containing words like learn, data, science, python, application, etc.



Logistics

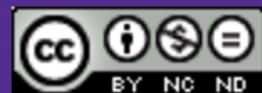
- Instructors: He and Chris
 - Office Hours:
Thursdays 9:30-10:30AM at 60 Fifth Avenue, Room 650



Logistics

- ▶ Section Leaders:
Shubham, Joshua,
Yiqiu
- ▶ Graders: Raghav,
Sarthak, Peeyush,
Aniket, Ieshan,
Chinmay

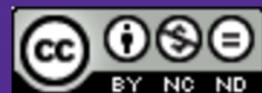
A word cloud centered on the words "data" and "science". Other prominent words include "learning", "program", "work", "good", "project", "real", "python", "application", and "method". Smaller words surrounding the center include "algorithm", "statistics", "expect", "gain", "knowledge", "basic", "making", "experience", "practical", "analyze", "library", "create", "expand", "actual", "hope", "large", "tool", "skill", "job", "code", "world", "large", "idea", "lot", "field", "clean", "don't", "understanding", "deep", "interest", "apply", and "algorithm".



Logistics

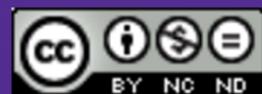
- ▶ NYU Classes
 - ▶ Weekly Agenda, Zoom Conference, Syllabus
- ▶ JupyterHub
 - ▶ Class Materials, Submission Labs/Homework/Projects
- ▶ Piazza
 - ▶ Announcements, Discussion
- ▶ Gradescope
 - ▶ Submission Homework/Projects, Retrieve Exams

applied
algorithm
don't
interest
understanding
deep
field
learning
clean
program
lot
idea
expect
skill
job
code
world
large
hope
data
science
python
application
method
knowledge
real
basic
making
practical
analyze
experience
library
help
create
expand
actual
method
hope
large
code
job
skill
lot
idea
field
deep
don't
understanding
interest
apply
algorithm
statistics
model
function
set
expect
gain
work
good
project
knowledge
real
python
application
method
class
hand
making
practical
analyze
experience
library
help
create
expand
actual

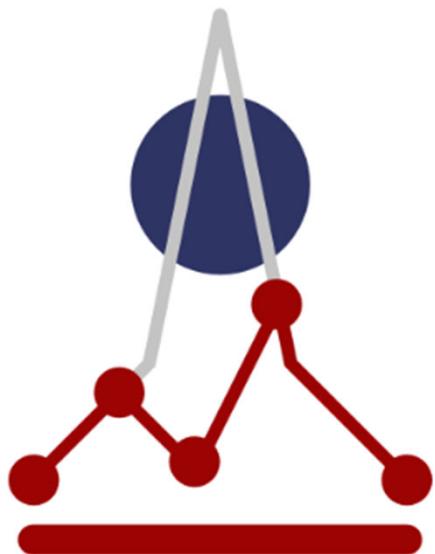


Logistics

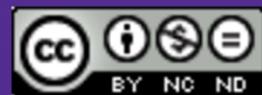
- ▶ Homework 40%
 - ▶ Project 15%
 - ▶ Exams
 - ▶ Midterm 20%
 - ▶ Final 25%
 - ▶ Extra Credit
 - ▶ Optional Problems
 - ▶ Office Hours
 - ▶ Piazza



Logistics

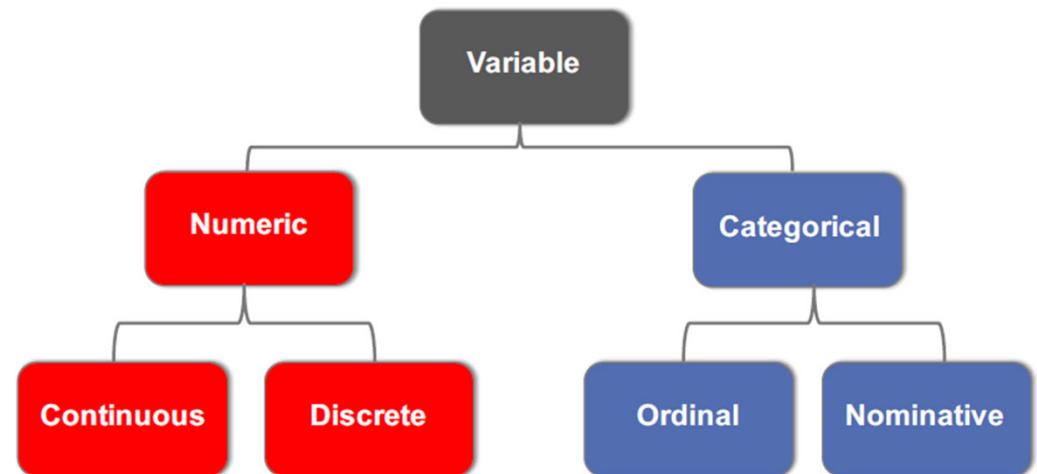


- ▶ Project Goals
 - ▶ Formulate
 - ▶ Process
 - ▶ Access
 - ▶ Visualize
 - ▶ Fit
 - ▶ Evaluate



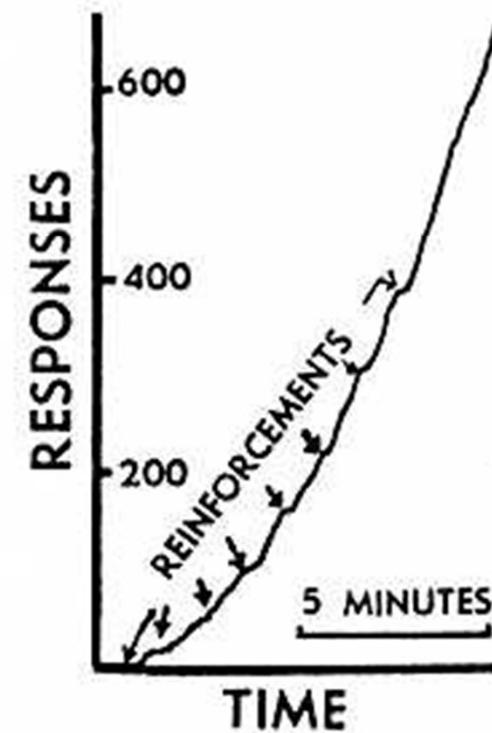
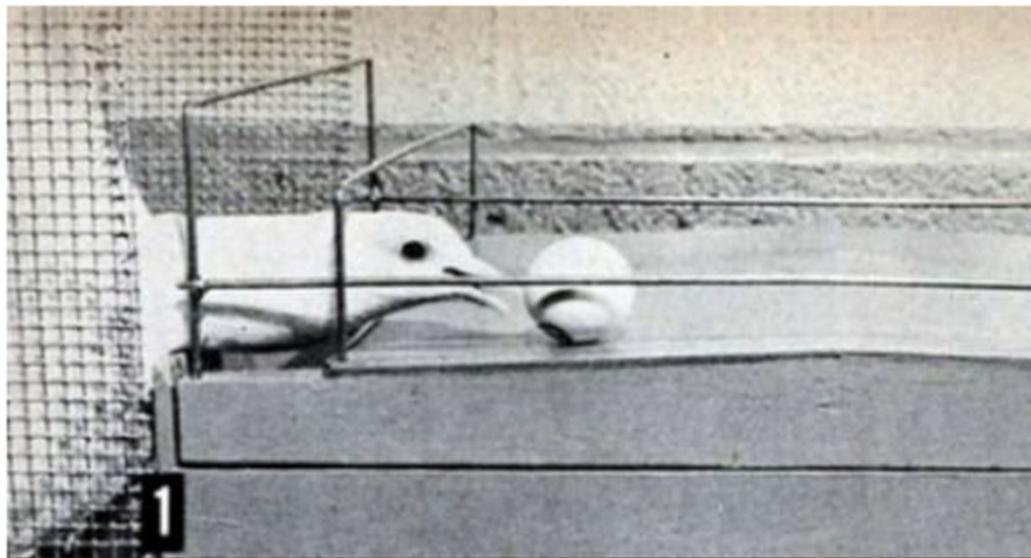
Types of machine learning

- ▶ Supervised Learning
 - ▶ Classification
 - ▶ Regression
- ▶ Unsupervised Learning
 - ▶ Clustering
 - ▶ Dimension Reduction
- ▶ Reinforcement Learning

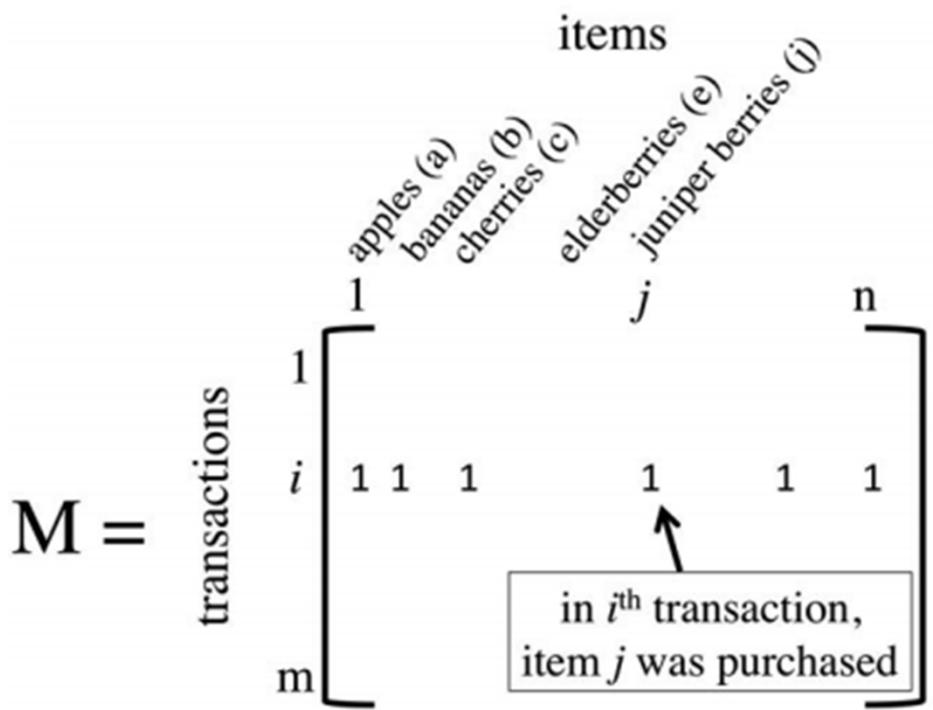


Types of machine learning

- ▶ B.F. Skinner (1950)
- ▶ Garcia and Koelling (1996)



Clustering



Clustering

- **Support** of an itemset: number of transactions containing it,

$$\text{Supp}(\text{bananas, cherries, elderberries}) = \sum_{i=1}^m M_{i,2} \cdot M_{i,3} \cdot M_{i,5}.$$

- **Confidence** of rule $a \rightarrow b$: the fraction of times itemset b is purchased when itemset a is purchased.

$$\begin{aligned} \text{Conf}(a \rightarrow b) &= \frac{\text{Supp}(a \cup b)}{\text{Supp}(a)} = \frac{\#\text{times } a \text{ and } b \text{ are purchased}}{\#\text{times } a \text{ is purchased}} \\ &= \hat{P}(b|a). \end{aligned}$$

- **Itemset**: a subset of items, e.g., (bananas, cherries, elderberries), indexed by $\{2, 3, 5\}$.

$$M =$$

transactions

	1	<i>i</i>	<i>j</i>	<i>n</i>
1	1	1	1	1
<i>i</i>	1	1	1	1
<i>j</i>			1	1
<i>n</i>			1	1

in i^{th} transaction,
 item j was purchased

Clustering

- **Support** of an itemset: number of transactions containing it,

$$\text{Supp}(\text{bananas, cherries, elderberries}) = \sum_{i=1}^m M_{i,2} \cdot M_{i,3} \cdot M_{i,5}.$$

- **Confidence** of rule $a \rightarrow b$: the fraction of times itemset b is purchased when itemset a is purchased.

$$\begin{aligned} \text{Conf}(a \rightarrow b) &= \frac{\text{Supp}(a \cup b)}{\text{Supp}(a)} = \frac{\#\text{times } a \text{ and } b \text{ are purchased}}{\#\text{times } a \text{ is purchased}} \\ &= \hat{P}(b|a). \end{aligned}$$

- **Itemset**: a subset of items, e.g., (bananas, cherries, elderberries), indexed by $\{2, 3, 5\}$.

$\text{Supp}(a \cup b) \geq \theta$, and $\text{Conf}(a \rightarrow b) \geq \text{minconf}$.

$M =$

items

	apples (a)	bananas (b)	cherries (c)	elderberries (e)	juniper berries (j)
1	1				
i	1	1	1	1	1
j				1	1
m					

n

in i^{th} transaction,
 item j was purchased

Clustering

- For each frequent itemset ℓ :
 - Find all nonempty subsets of ℓ
 - For each subset a , output $a \rightarrow \{\ell \setminus a\}$ whenever

$$\frac{\text{Supp}(\ell)}{\text{Supp}(a)} \geq \text{minconf.}$$

$M =$

Confidence

		items				
		apples (a)	bananas (b)	cherries (c)	elderberries (e)	juniper berries (j)
		1				
		i	1 1	1	1	j
		m				n

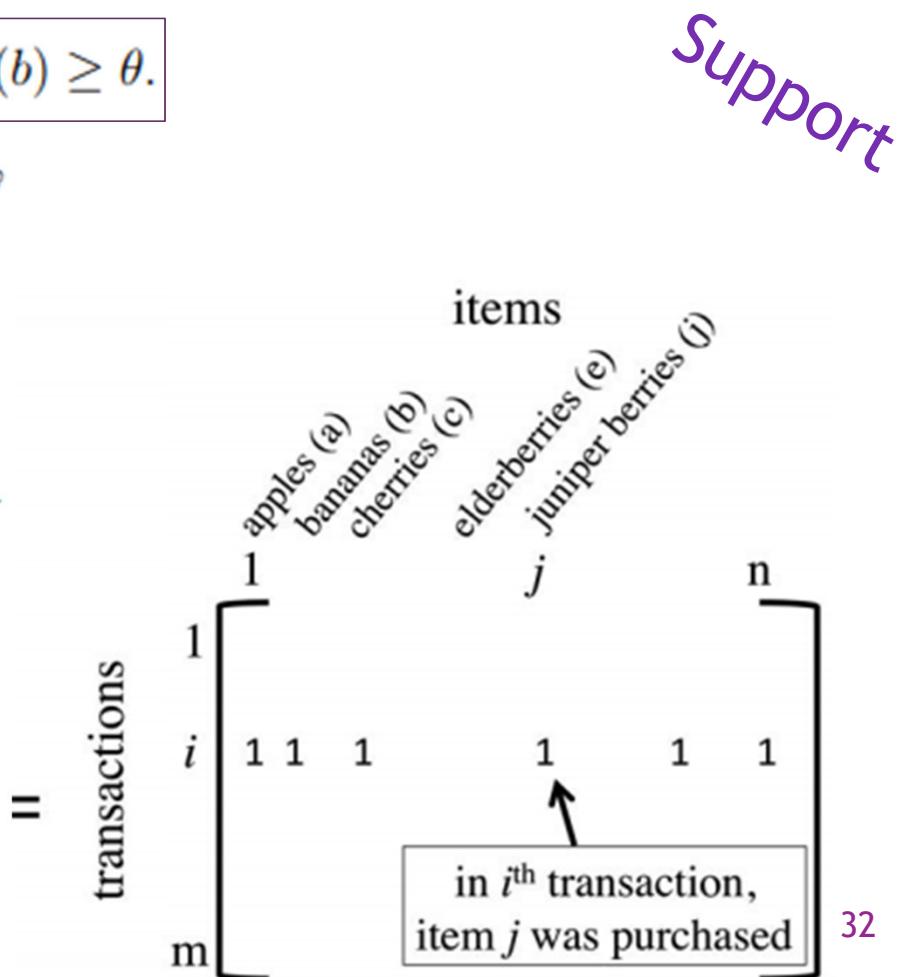
transactions

in i^{th} transaction,
item j was purchased

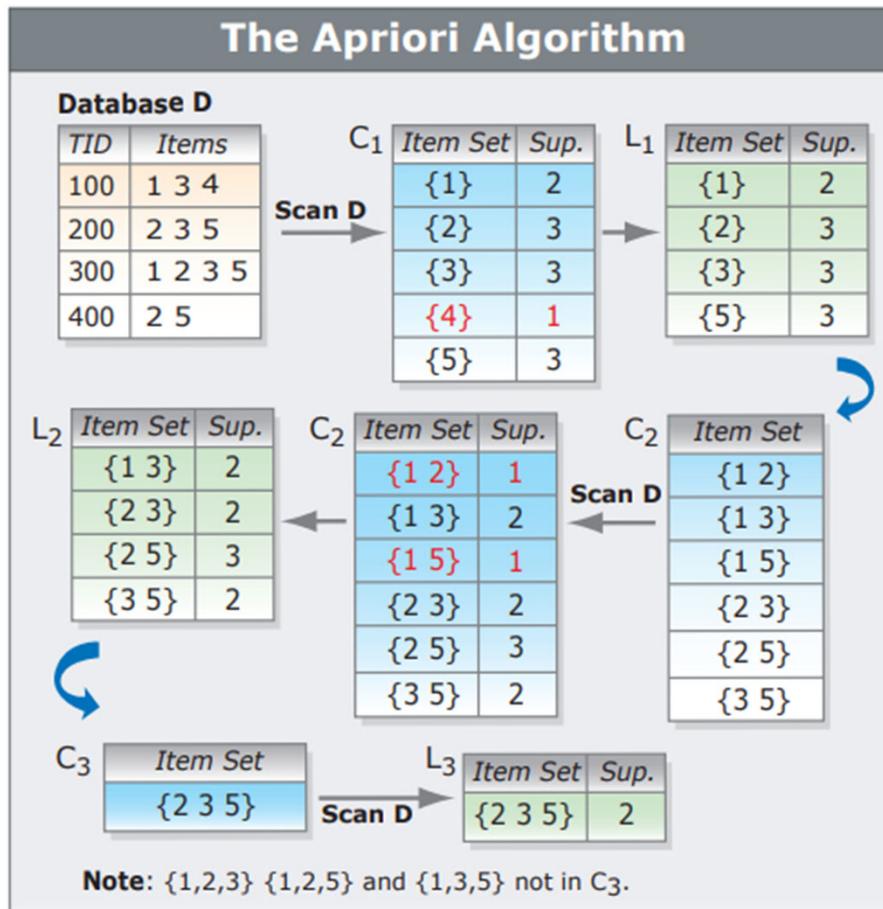
Clustering

If $\text{Supp}(a \cup b) \geq \theta$ then $\text{Supp}(a) \geq \theta$ and $\text{Supp}(b) \geq \theta$.

	<i>apples</i>	<i>bananas</i>	<i>cherries</i>	<i>elderberry</i>	<i>grapes</i>	
1-itemsets:	a	b	c	d	e	f
supp:	25	20	30	45	29	5
2-itemsets:	{a,b}	{a,c}	{a,d}	{a,e}	...	{e,g}
supp:	7	25	15	23	...	3
3-itemsets:	{a,c,d}	{a,c,e}	{b,d,g}	...		
supp:	15	22	15			
4-itemsets:	{a,c,d,e}					
supp:	12					



Clustering



Support

Clustering

The Apriori Algorithm

Database D

TID	Items
100	1 3 4
200	2 3 5
300	1 2 3 5
400	2 5

C ₁	Item Set	Sup.
	{1}	2
	{2}	3
	{3}	2
	{4}	
	{5}	

L ₂	Item Set	Sup.
	{1 3}	2
	{2 3}	2
	{2 5}	3
	{3 5}	2

C ₂	Item Set	Sup.
	{1 2}	1
	{1 3}	2
	{1 5}	1
	{2 3}	2
	{2 5}	3
	{3 5}	2



C ₃	Item Set
	{2 3 5}

L ₃	Item
	{2 3}

Note: {1,2,3} {1,2,5} and {1,3,5} no

Input: Matrix M

$L_1 = \{\text{frequent 1-itemsets}; i \text{ such that } \text{Supp}(i) \geq \theta\}.$

For $k = 2$, while $L_{k-1} \neq \emptyset$ (while there are large $k-1$ -itemsets), $k++$

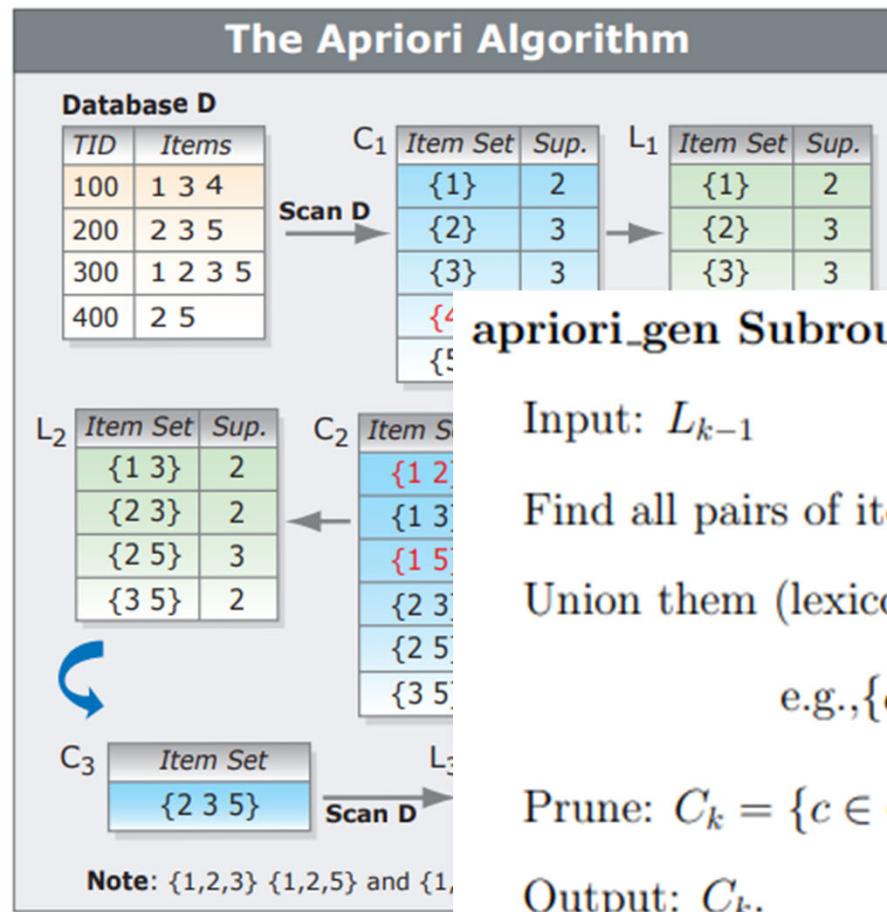
- $C_k = \text{apriori_gen}(L_{k-1})$ generate candidate itemsets of size k
- $L_k = \{c : c \in C_k, \text{Supp}(c) \geq \theta\}$ frequent itemsets of size k (loop over transactions, scan the database)

end

Output: $\bigcup_k L_k.$

Support

Clustering



Support

Regression

Assume a linear relation

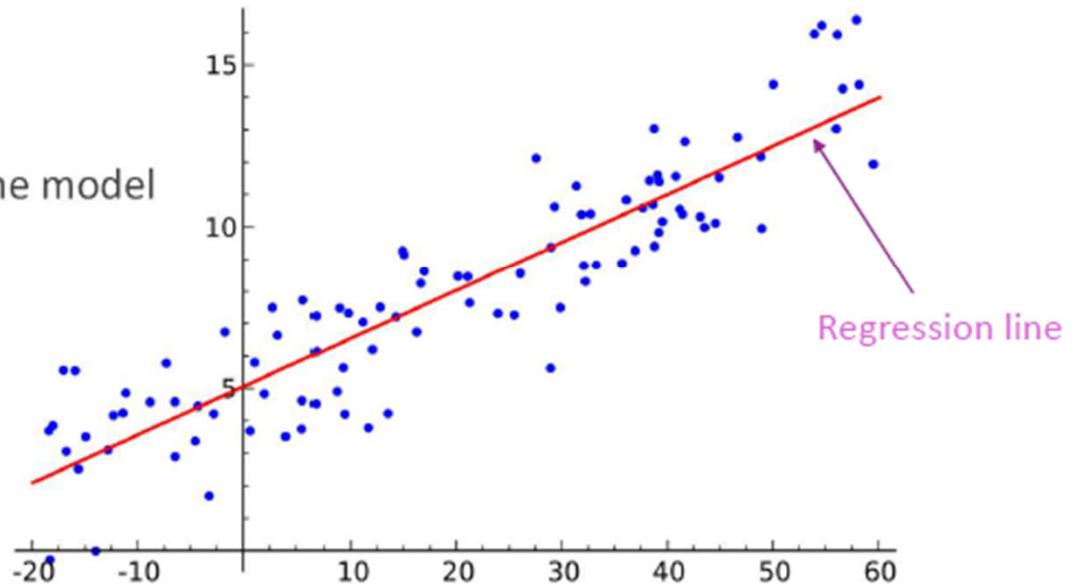
$$y \approx \beta_0 + \beta_1 x$$

- › β_0 = intercept
- › β_1 = slope

$\beta = (\beta_0, \beta_1)$ are the **parameters** of the model

What are the units of β_0, β_1 ?

When is this model good?



Regression

Knowing x does not exactly predict y

- Variation in y due to factors other than x

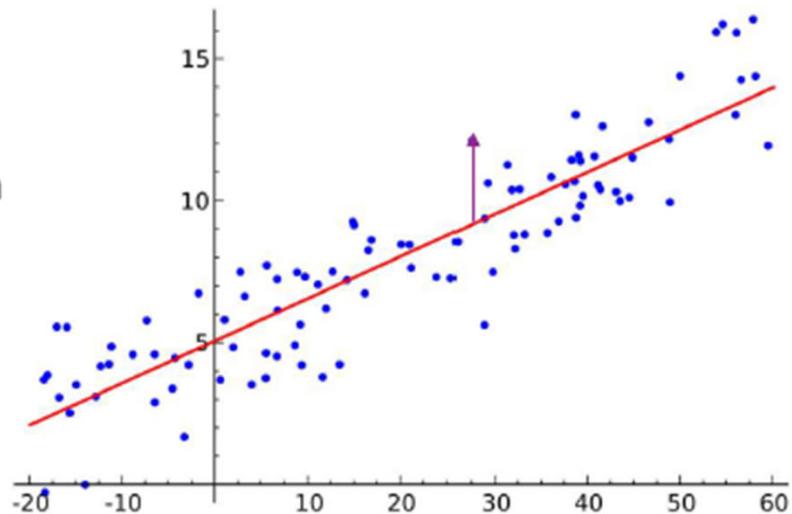
Add a **residual** term

$$y = \beta_0 + \beta_1 x + \epsilon$$

Residual = component the model does not explain

- Predicted value: $\hat{y}_i = \beta_1 x_i + \beta_0$
- Residual: $\epsilon_i = y_i - \hat{y}_i$

Vertical deviation from the regression line



Regression

| How do we select parameters $\beta = (\beta_0, \beta_1)$?

| Define $\hat{y}_i = \beta_1 x_i + \beta_0$

- Predicted value on sample i for parameters $\beta = (\beta_0, \beta_1)$

| Define average residual sum of squares:

$$\text{RSS}(\beta_0, \beta_1) := \sum_{I=1}^n (y_i - \hat{y}_i)^2$$

- Note that \hat{y}_i is implicitly a function of $\beta = (\beta_0, \beta_1)$
- Also called the sum of squared residuals (SSR) and sum of squared errors (SSE)

| Least squares solution: Find (β_0, β_1) to minimize RSS.

- Geometrically, minimizes squared distances of samples to regression line

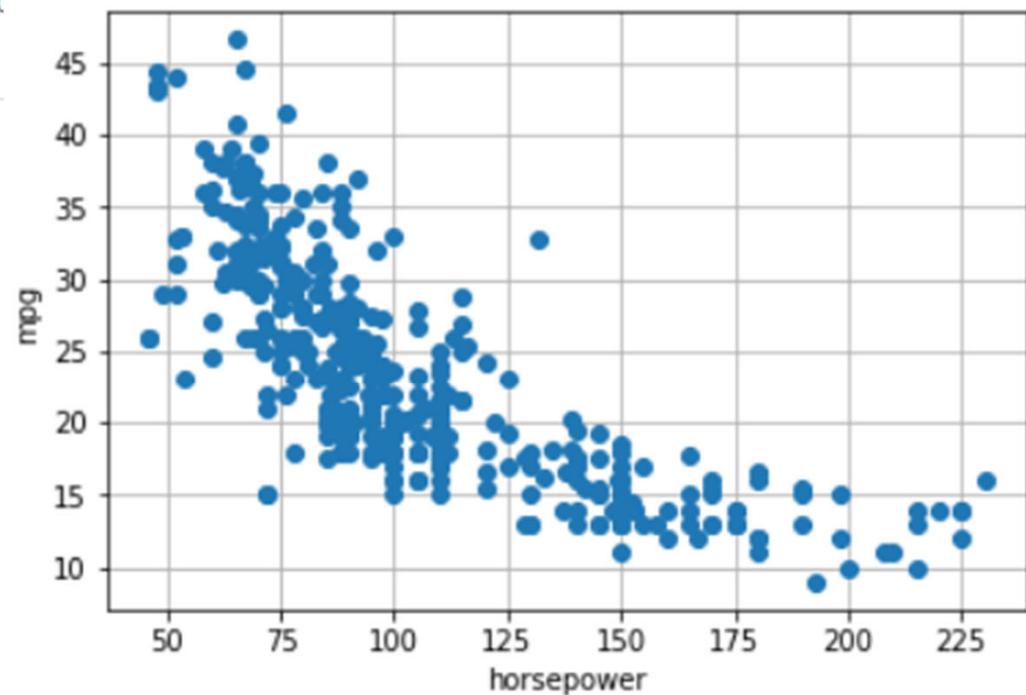
Exercise

SA Scientific American

VW Emissions Cheating Scandal Increased Children's ...

The Volkswagen emissions cheating scandal—also known as ... September 2015, when a team of scientists at West Virginia University discovered that the ... The working paper from the Chicago Fed relied ...

Jul 17, 2019



Regression

General ML problem

|Find a **model** with **parameters**

→ Linear model: $\hat{y} = \beta_0 + \beta_1 x$

|Get **data**

→ Data: $(x_i, y_i), i = 1, 2, \dots, N$

|Pick a **loss function**

- Measures goodness of fit model to data
- Function of the parameters

→ Loss function:

$$RSS(\beta_0, \beta_1) := \sum (y_i - \beta_0 + \beta_1 x_i)^2$$

|Find parameters that **minimizes** loss

→ Select β_0, β_1 to minimize $RSS(\beta_0, \beta_1)$

Questions

- ▶ Questions on Piazza?
- ▶ Question for You!
 - ▶ Lift is Weighted Confidence?

$$\frac{\hat{P}(b|a)}{\hat{P}(b)}$$

- ▶ What is the interpretation of lift?
- ▶ Why would we use lift instead of confidence?
- ▶ Compute lift in terms of the supports of a,b and aUb

