

Statistical Methods in AI (CS7.403)

Lecture-2: ML Workflow, Data Representations,
Basic Data Transformations, Data Visualization

Ravi Kiran (ravi.kiran@iiit.ac.in)

<https://ravika.github.io>



Center for Visual Information Technology (CVIT)
IIIT Hyderabad

Announcements

- Tutorial (11.40a – 12.40p Saturday, H-205)
 - Python, Pandas, Jupyter notebook, Plotting tools.
 - **Bring your laptops.**
- Ask questions.

Announcements

- **IMPORTANT:** All assignments/projects will need to submitted via Github Classroom
- Tutorial
 - Git
 - Github
- Ask questions.

Announcements

- TAs will share SMAI Course Calendar on Moodle
- You can add it to your Teams Calendar
- Will contain assignment release/due/eval dates
- Will contain exam paper showing dates
- Will contain project-related dates

Queries

- Post queries on Moodle
- Helps all (many may have same question)
- Do not DM TAs !

Announcements

- Do not use Python libraries for assignments unless explicitly allowed/specied.

Additionally ...

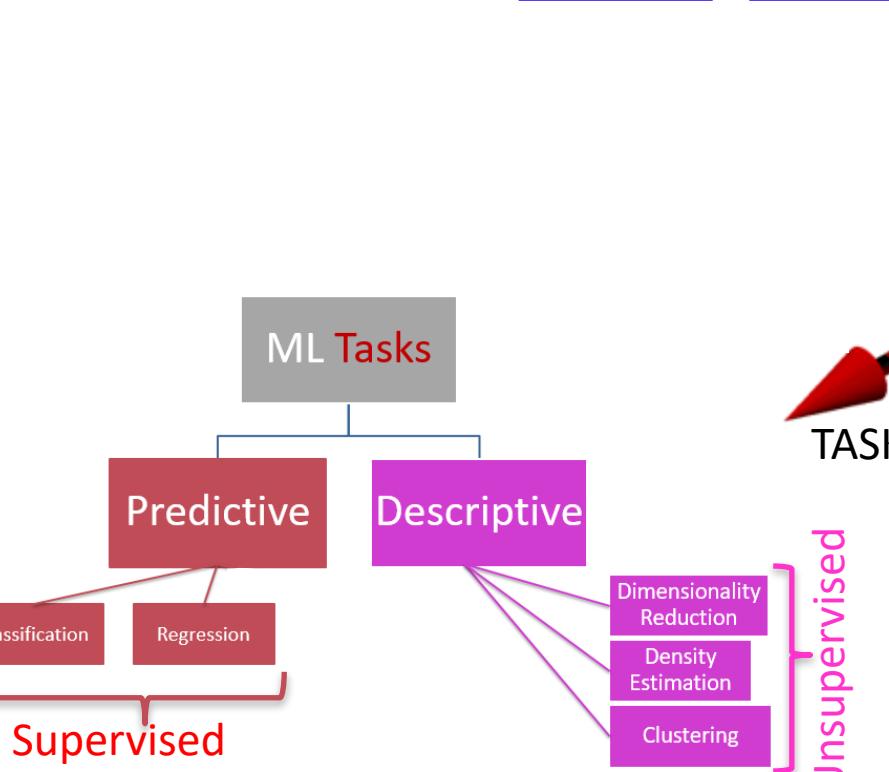
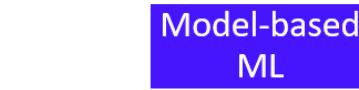
- Spending time everyday on material covered in class helps
 - Take notes
 - Revise
 - Reflect
- **Ask if you wish to take something down, but slide is no longer on screen**



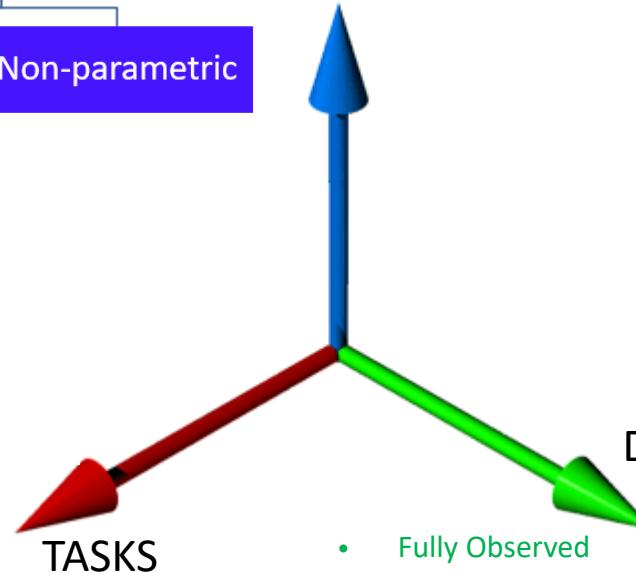
Course TAs

1. Khushi Agarwal
2. Naimeesh Narayan Tiwari
3. Tathagato Roy
4. Nachiket Patil
5. Nitin Shrinivas
6. Vaibhav Agarwal
7. Naraharisetti Siddik Ayyappa
8. Kawshik Manikantan

Recap



ALGORITHMS



Representation
DATA

- Fully Observed
- Partially Observed
 - Some variables systematically not observed (e.g. 'topic' of a document)
 - Some variables missing some of the time (e.g. 'faulty sensor' readings)

Lecture Outline

- ML Workflow
- Data Representations
- Basic Data Transformations
- Data Visualization

Machine Learning



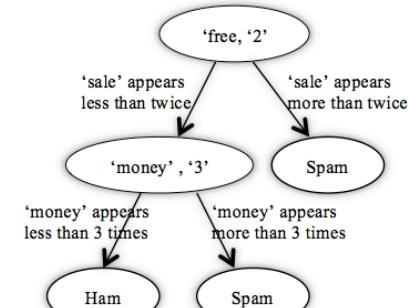
Algorithmic methods that use **data** to improve their **knowledge** of a **task**

Task: Detect spam email



Data: Labelled emails
(in inboxes of other users as well !)

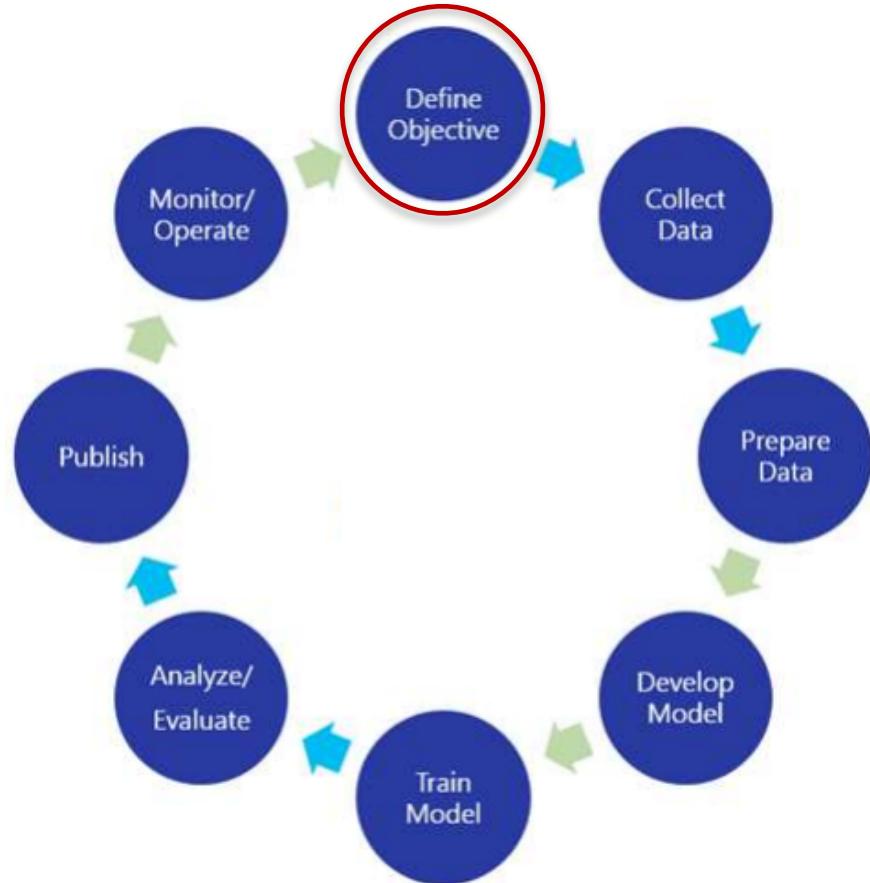
Knowledge:



Improve → 85% reduction of spam emails in Inbox over 3 months

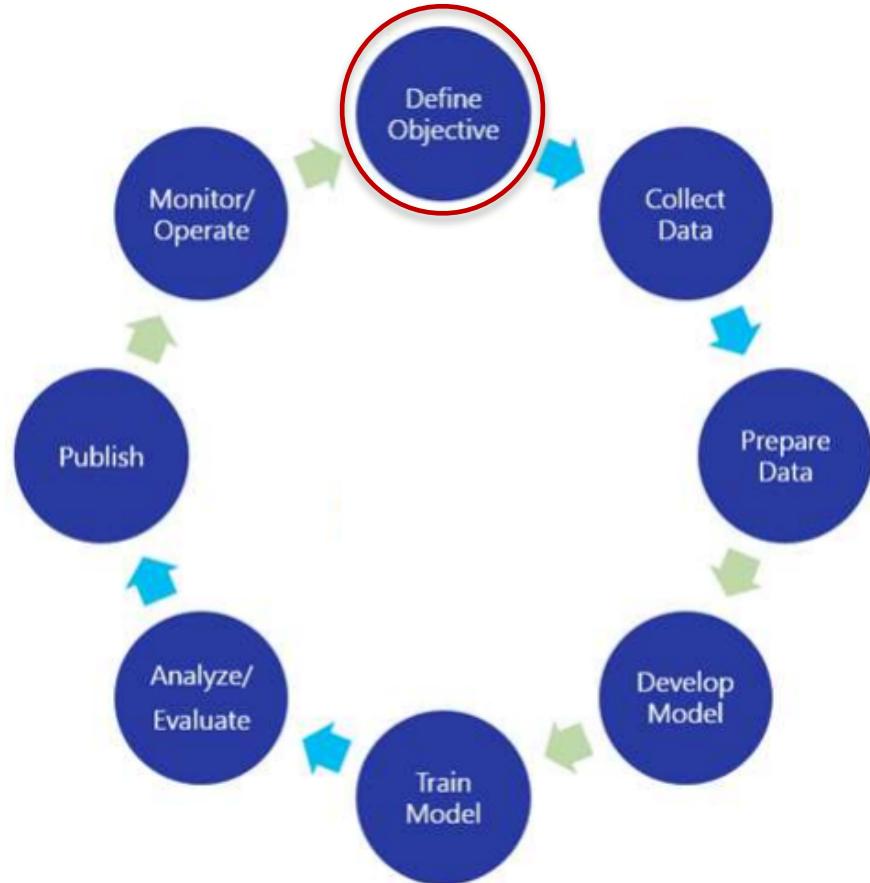
Algorithmic method: Decision Tree

Workflow of a Machine Learning Problem

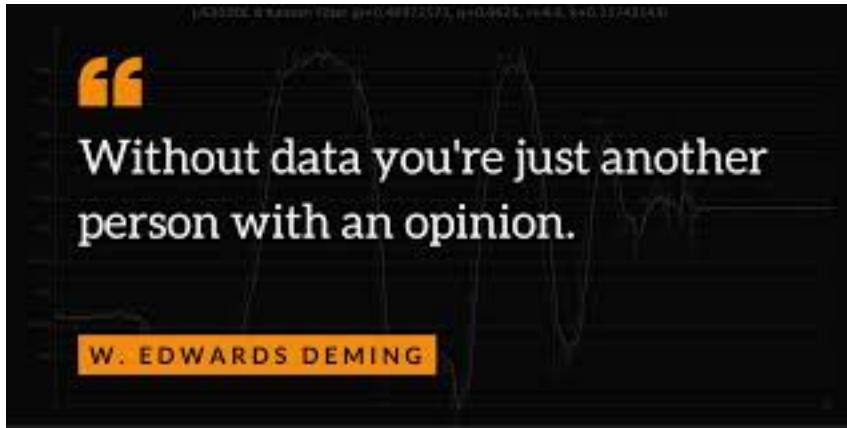


- Detect spam email
- Predict value of a stock
- Predict effect of advertising on sales
- Drive car 'safely' without human intervention
- Translate text from one language to another
- Sentiment Analysis
- ...

Workflow of a Machine Learning Problem



No Data, no ML !



Sources of data

- Detect spam email

Mail - (no subject) 3/22/13 12:59 PM

 [View in browser](#)

(no subject)
1 message

LC Johnson <cgconfidential@gmail.com> Tue, Mar 12, 2013 at 9:39 AM

To: naomi

Hi Naomi,

Hope this email finds you well! It's LC from Colored Girl Confidential.

I won't take up too much of your time as I'm sure you're busy with the epic 73% off sale on IttyBiz. (I'm currently deciding between How to Launch The **** Out of an Ebook and the How Not to Screw Up Bundle!) Anyway, I'm emailing because I was reading through some of your older blog posts and it occurred to me that you might appreciate my recently launched manifesto: **The Red Lipstick Manifesta**.

I consider it one of the most important things that I've ever created and very much inspired by your constant reminders that successful women don't just blindly follow the rules. They are rebels, challenging the assumptions of what everyone thinks is possible, chasing down their "unrealistic" dreams, and eventually creating lives and careers they love!

If you have a few minutes I hope you'll check it out. Let me know what you think and have a great week!

Much love,

LC

LC Johnson
Founding Editor, Colored Girl Confidential
Latest Post: [The Red Lipstick Manifesta](#)

Check us out at www.coloredgirlconfidential.com.
Join the conversation on [Facebook](#) and [Twitter](#).

<https://mail.google.com/mail/u/0/?hl=zhd4f9ab&view=pb&search=send&h=13d5ed40fb1d1c14>

A screenshot of a Gmail compose screen. The top navigation bar shows 'Gmail' and a search bar. Below it, the 'Compose' button is highlighted in red. The message header includes 'SEND', 'Preview', 'Send', and 'Labels'. The recipient field contains 'Salvador Faria <salvador.mrf@gmail.com>'. The 'To' field has 'Salvador Faria' and his email address. The 'Subject' field contains 'lorem ipsum'. Below the subject are buttons for 'Attach a file', 'Insert: Invitation', and 'Canned responses'. The message body starts with 'Lorem ipsum dolor sit amet, sit nostro utamur qualisque ne, no tantas electram est. Per legimus ludicribus omnitarium eu, el has antipopem neglegentibus philosophia. Viderer luvaret vis eum, mel legimus vivendo ad. Eum no atqui nullam, harum solet pericula quo te, facer ludus partem an nec. Ex ius habeo mnesarchum, ne nisl augue sadipsicing vis, sumo doming patroique nec at.' A toolbar below the message body includes icons for bold, italic, underline, strikethrough, font size, alignment, and a check spelling button. The bottom of the screen shows a portion of the Gmail inbox with several messages listed.

Business Email Sample

To: "Anna Jones" <annajones@buzzle.com.>
CC: All Staff
From: "James Brown"
Subject: Welcome to our Hive!

Dear Anna,

Welcome to our Hive!

It is a pleasure to welcome you to the team of _____. We are excited to have you join our team, and we hope that you will enjoy Working with our Company.

On the last Saturday of each month we hold a special staff party to welcome any new employees. Please be sure to come next Week to meet all of our senior staff and any other new staff members who have joined _____ this month. You will receive an e-mail regarding the same with further details.

If you have any questions during your training period, please do not hesitate to contact me. You can reach me at my email address or on my office line at 000-0001.

Warm regards,
James

Jackie Brown, Manager, Staff
jamesbrown@abcd.com
Tel: 000-0001

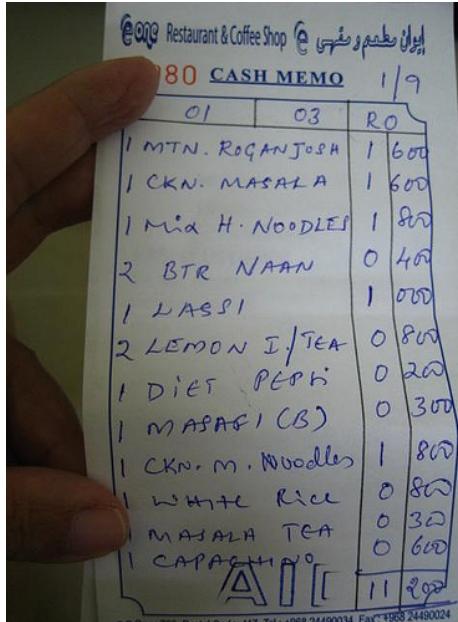
Sources of data

- Predict value of a stock



Sources of data

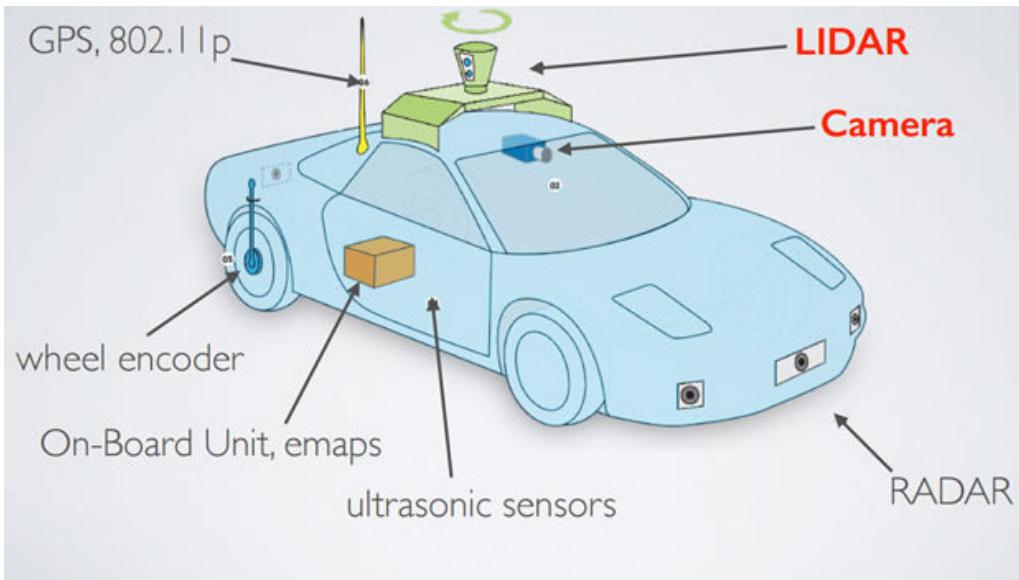
- Predict effect of advertising on sales



Raw Data may not always be digital in nature !

Sources of data

- Drive car safely without human intervention



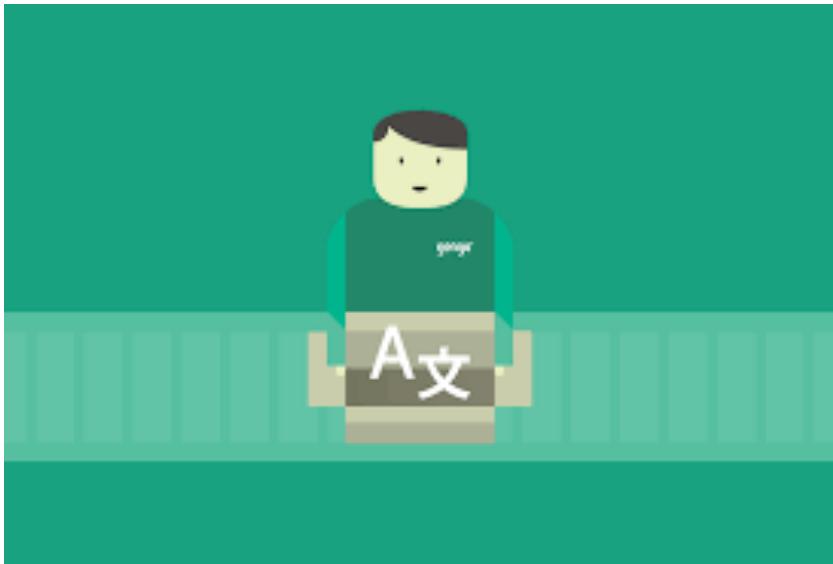
<https://inai.iiit.ac.in/bodhyaan.html>



Data can be multi-modal and may
need to be 'synchronized'

Sources of data

- Translate text from one language to another



A human domain expert
may be required to obtain
raw data

Two fundamental questions

- What data to collect ?
- How to collect ?

Raw data

- May be too little in quantity

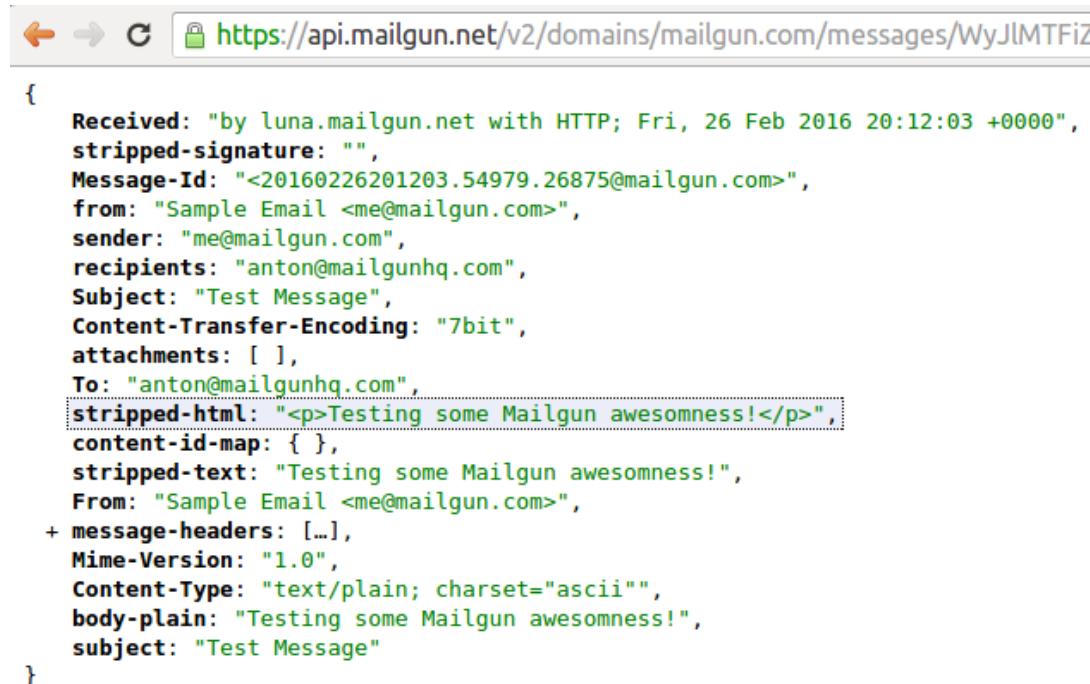
Raw data

- May be **too much** in quantity
 - Limitations on system end (compute, storage)



Raw data

- Not all of it relevant



A screenshot of a web browser displaying a JSON object. The URL in the address bar is <https://api.mailgun.net/v2/domains/mailgun.com/messages/WyJlMTFiZ>. The JSON object represents a single email message with various fields like Received, Message-ID, From, To, Subject, Content-Type, and body.

```
{
  Received: "by luna.mailgun.net with HTTP; Fri, 26 Feb 2016 20:12:03 +0000",
  stripped-signature: "",
  Message-Id: "<20160226201203.54979.26875@mailgun.com>",
  from: "Sample Email <me@mailgun.com>",
  sender: "me@mailgun.com",
  recipients: "anton@mailgunhq.com",
  Subject: "Test Message",
  Content-Transfer-Encoding: "7bit",
  attachments: [ ],
  To: "anton@mailgunhq.com",
  stripped-html: "<p>Testing some Mailgun awesomness!</p>",
  content-id-map: { },
  stripped-text: "Testing some Mailgun awesomness!",
  From: "Sample Email <me@mailgun.com>",
  + message-headers: [...],
  Mime-Version: "1.0",
  Content-Type: "text/plain; charset='ascii'",
  body-plain: "Testing some Mailgun awesomness!",
  subject: "Test Message"
}
```

Raw data

- Often not directly usable
 - Filter (needed data)
 - Transform (to numerical data)



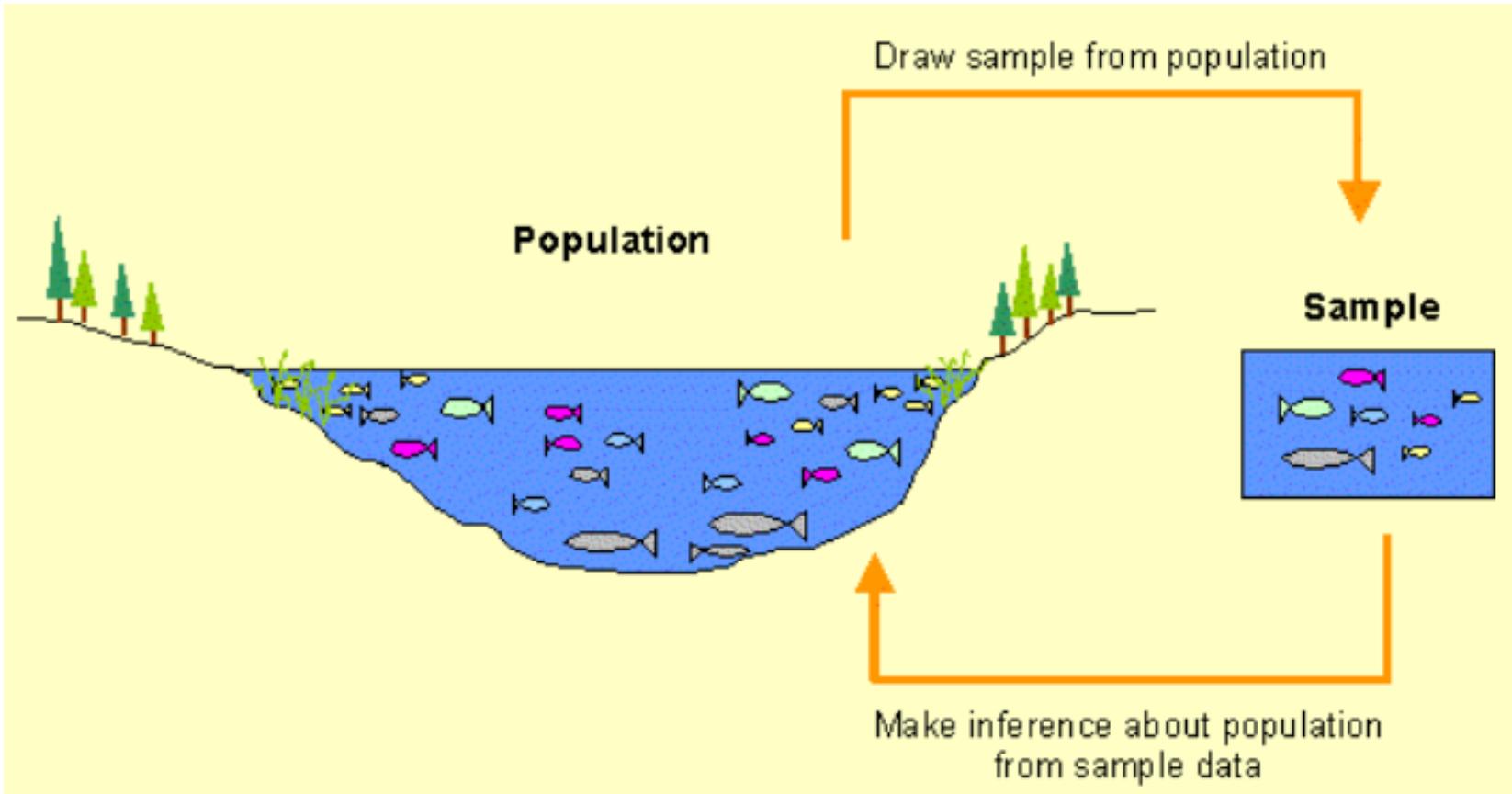
The screenshot shows a browser window with the URL <https://api.mailgun.net/v2/domains/mailgun.com/messages/WyJlMTFiZ>. The page displays a JSON object representing an email message. The JSON structure includes fields like Received, stripped-signature, Message-Id, from, sender, recipients, Subject, Content-Transfer-Encoding, attachments, To, stripped-html, content-id-map, stripped-text, From, + message-headers, Mime-Version, Content-Type, body-plain, and subject. The stripped-html field contains the text "Testing some Mailgun awesomness!". The browser's address bar and navigation buttons are visible at the top.

```
{  
  Received: "by luna.mailgun.net with HTTP; Fri, 26 Feb 2016 20:12:03 +0000",  
  stripped-signature: "",  
  Message-Id: "<20160226201203.54979.26875@mailgun.com>",  
  from: "Sample Email <me@mailgun.com>",  
  sender: "me@mailgun.com",  
  recipients: "anton@mailgunhq.com",  
  Subject: "Test Message",  
  Content-Transfer-Encoding: "7bit",  
  attachments: [ ],  
  To: "anton@mailgunhq.com",  
  stripped-html: "<p>Testing some Mailgun awesomness!</p>",  
  content-id-map: { },  
  stripped-text: "Testing some Mailgun awesomness!",  
  From: "Sample Email <me@mailgun.com>",  
+ message-headers: [...],  
  Mime-Version: "1.0",  
  Content-Type: "text/plain; charset='ascii'",  
  body-plain: "Testing some Mailgun awesomness!",  
  subject: "Test Message"  
}
```

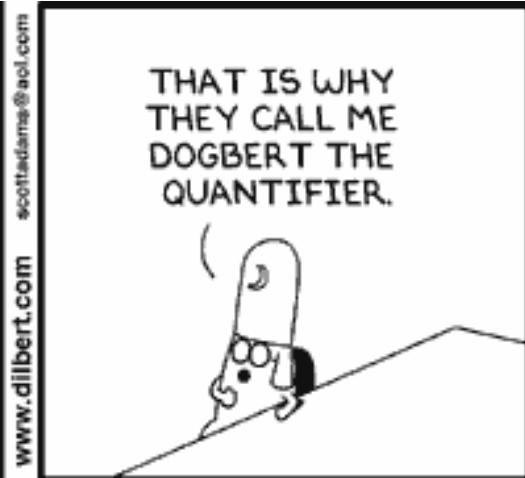
Two fundamental questions

- What data to collect ?
- How (much) to collect ?

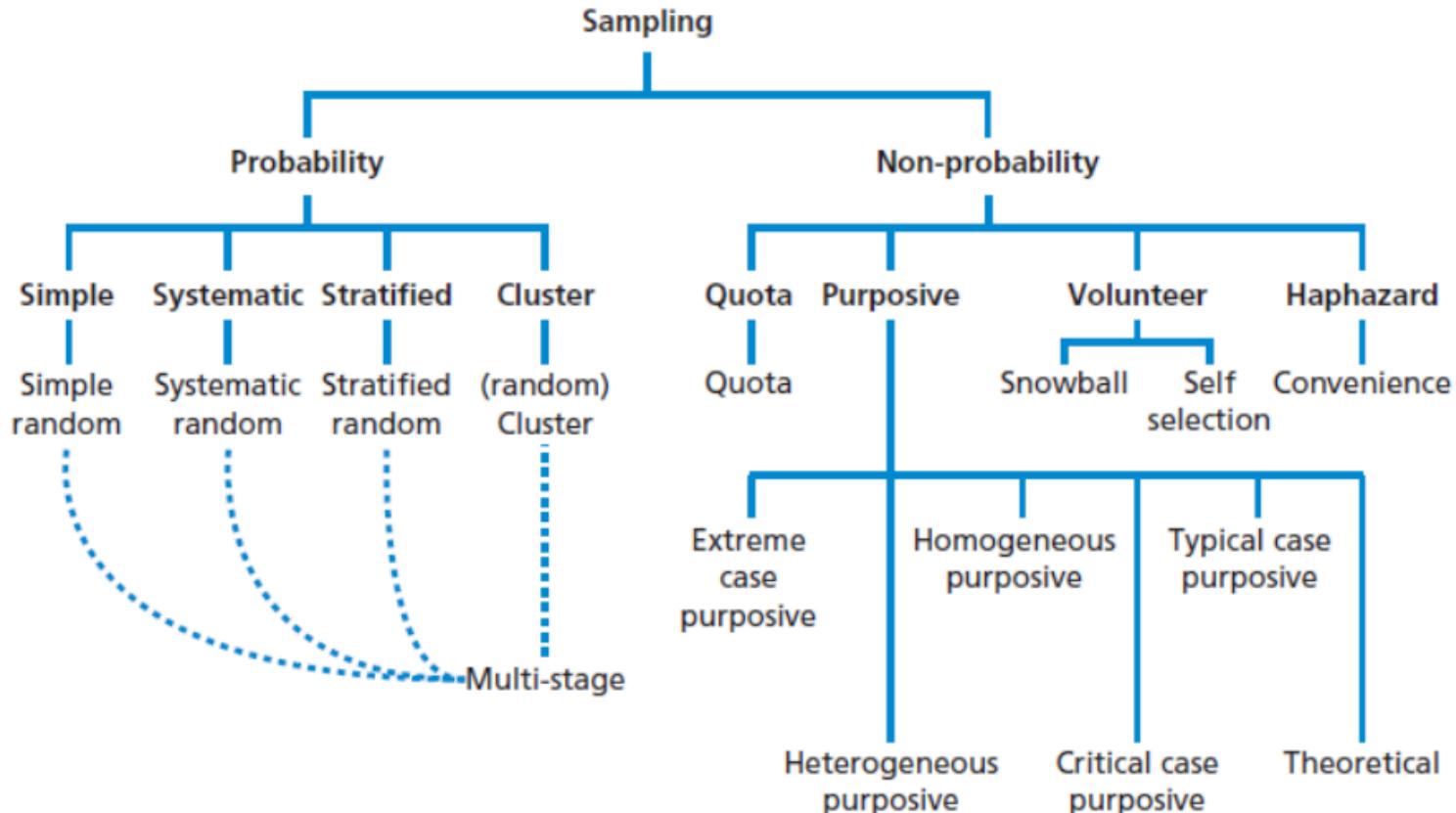
The Research Method



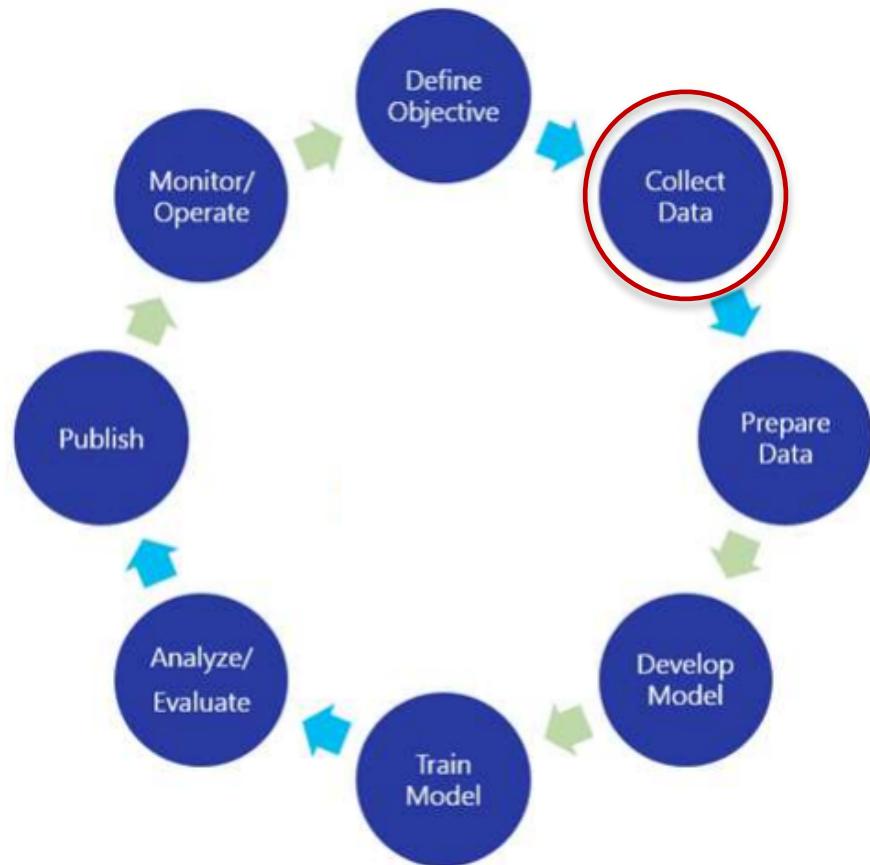
Are our samples ‘representative’ of the population?



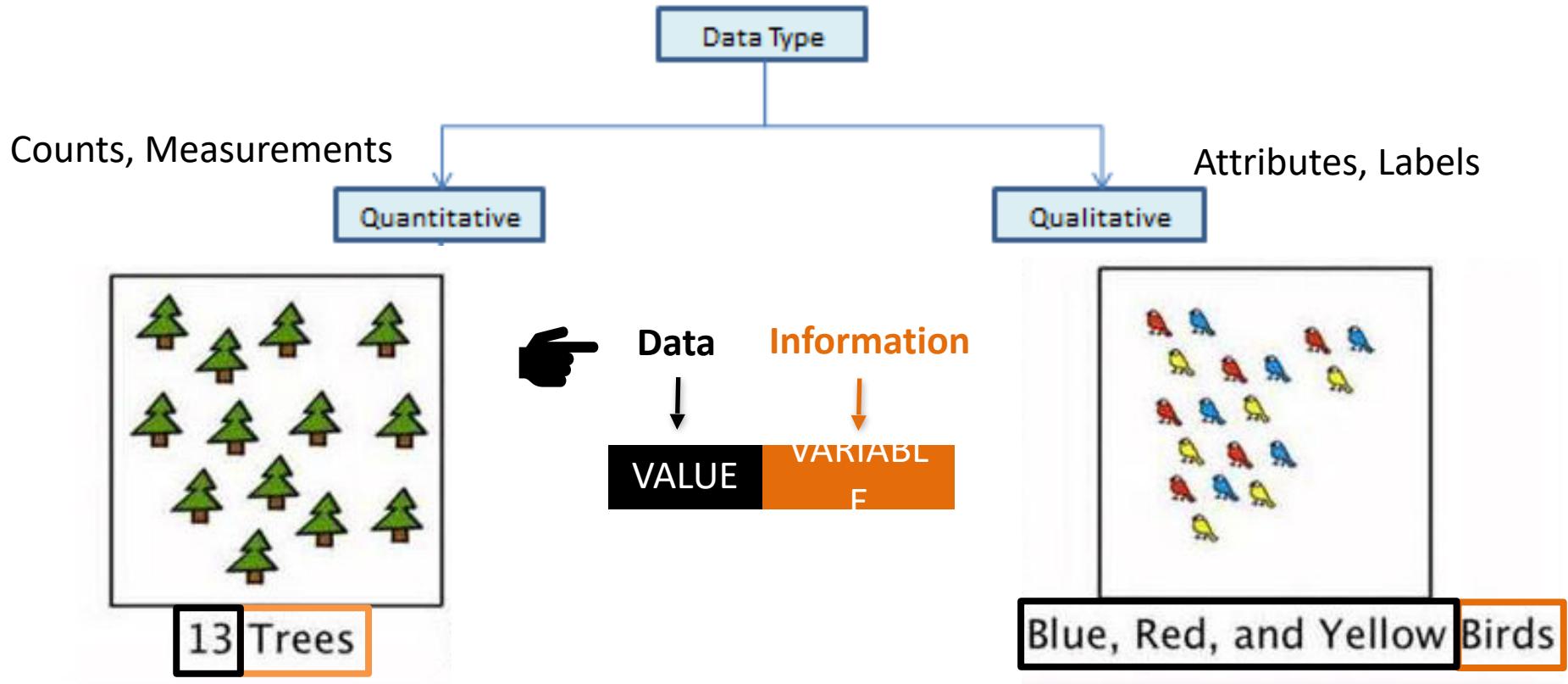
Sampling Techniques



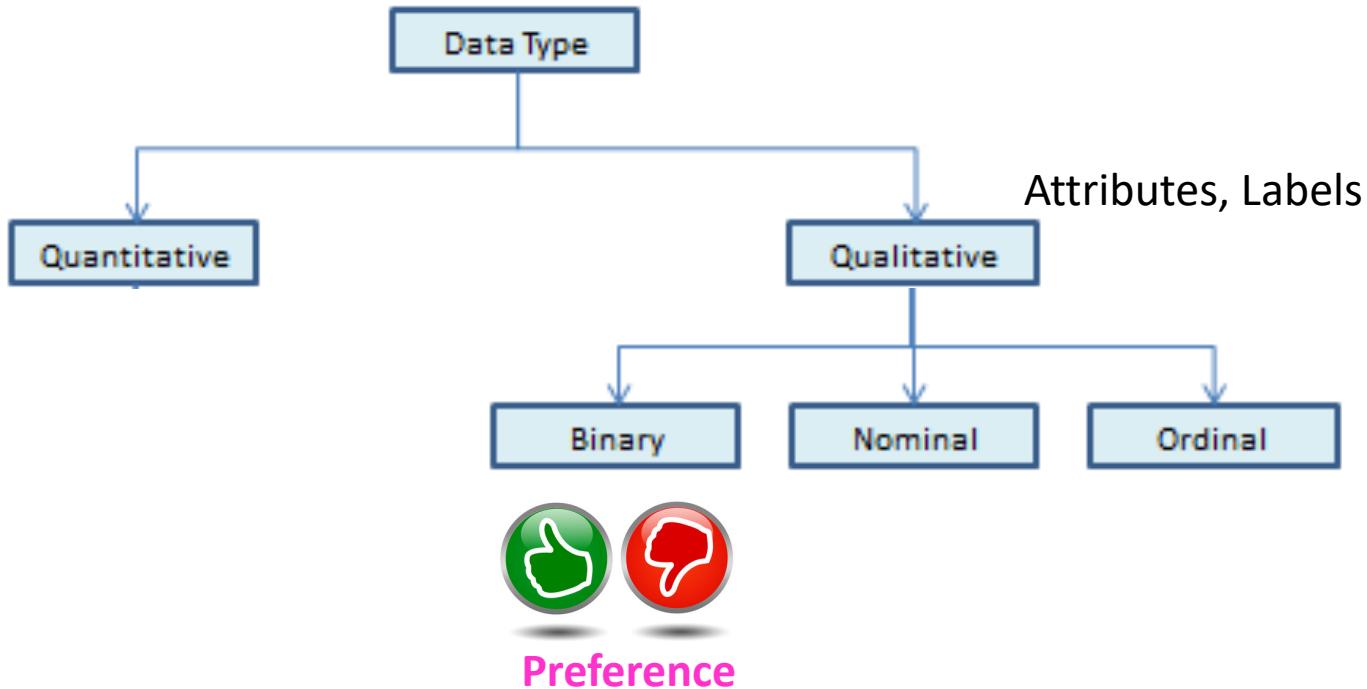
Workflow of a Machine Learning Problem



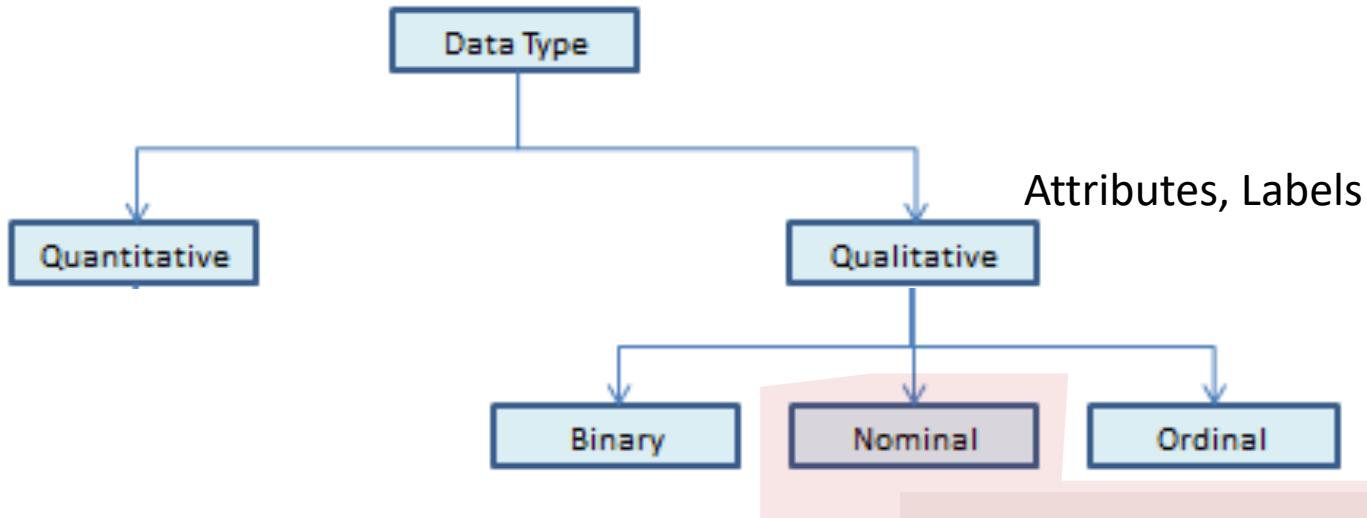
Taxonomy of data variables



Taxonomy of data

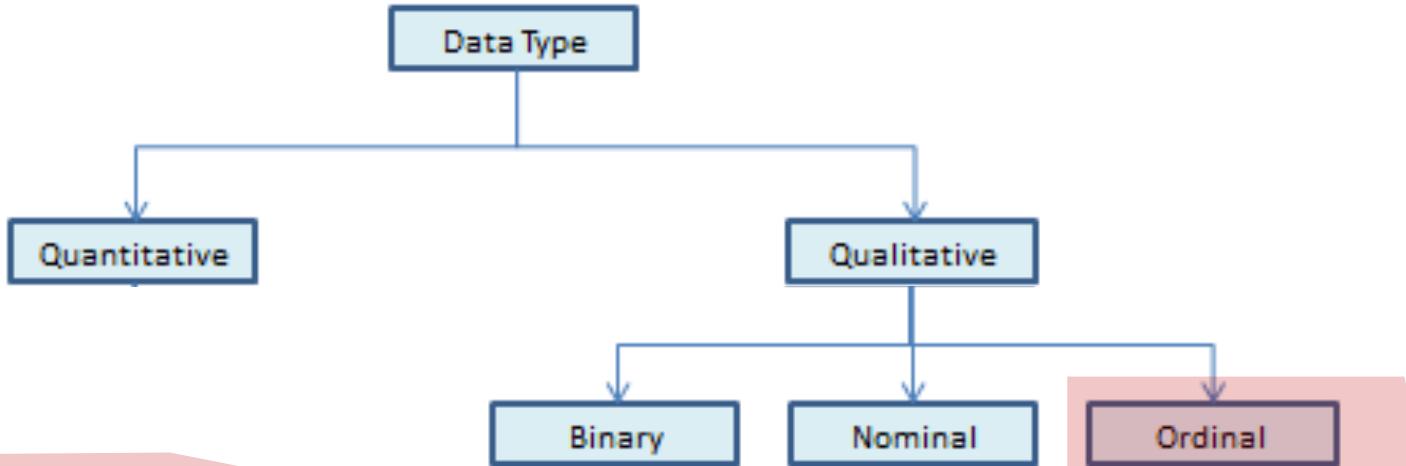


Taxonomy of data



Pin Code





How comfortable are you with Python *

No knowledge ... Very comfortable

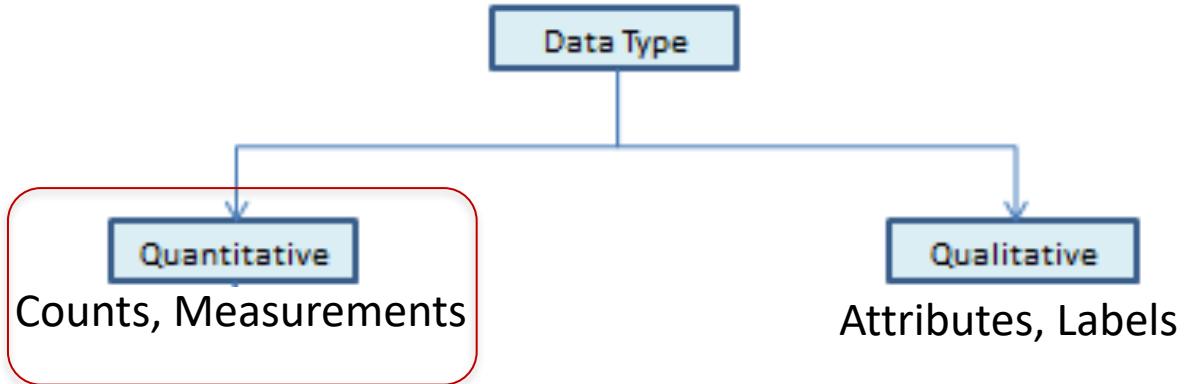
Very comfortable

XS **S** **M** **L** **XL** **XXL**

Letter grade
A +
A
A -
B +
B
B -
C +
C
C -
D +
D
E



Taxonomy of data



QUANTITATIVE DATA:



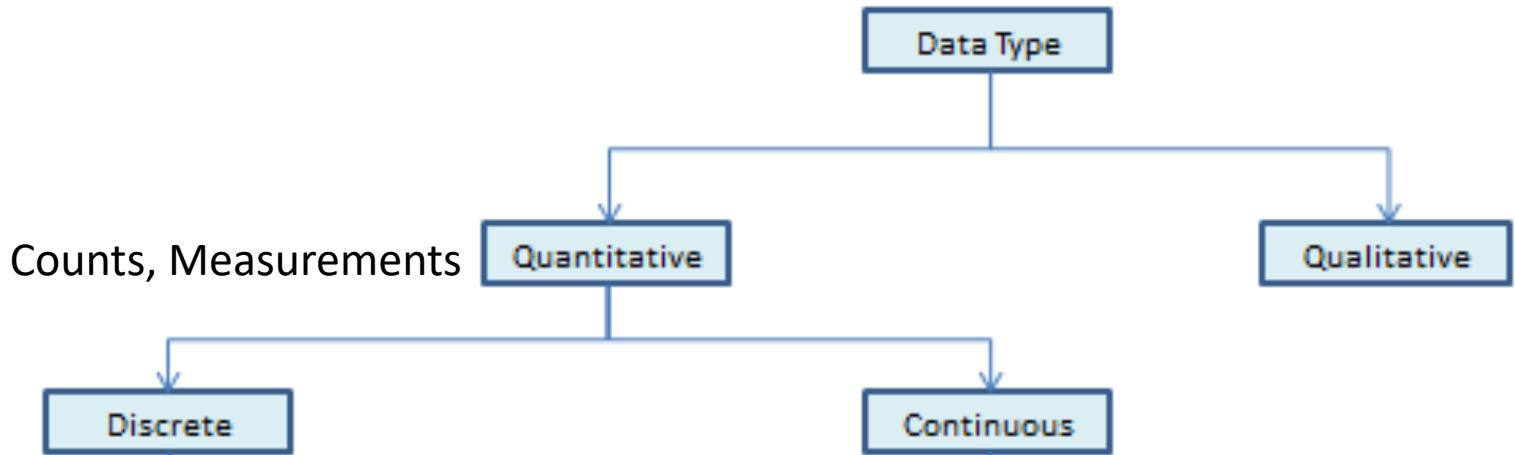
Discrete data:

- There are 3 cones
- Cone 1 has 2 scoops

Continuous data:

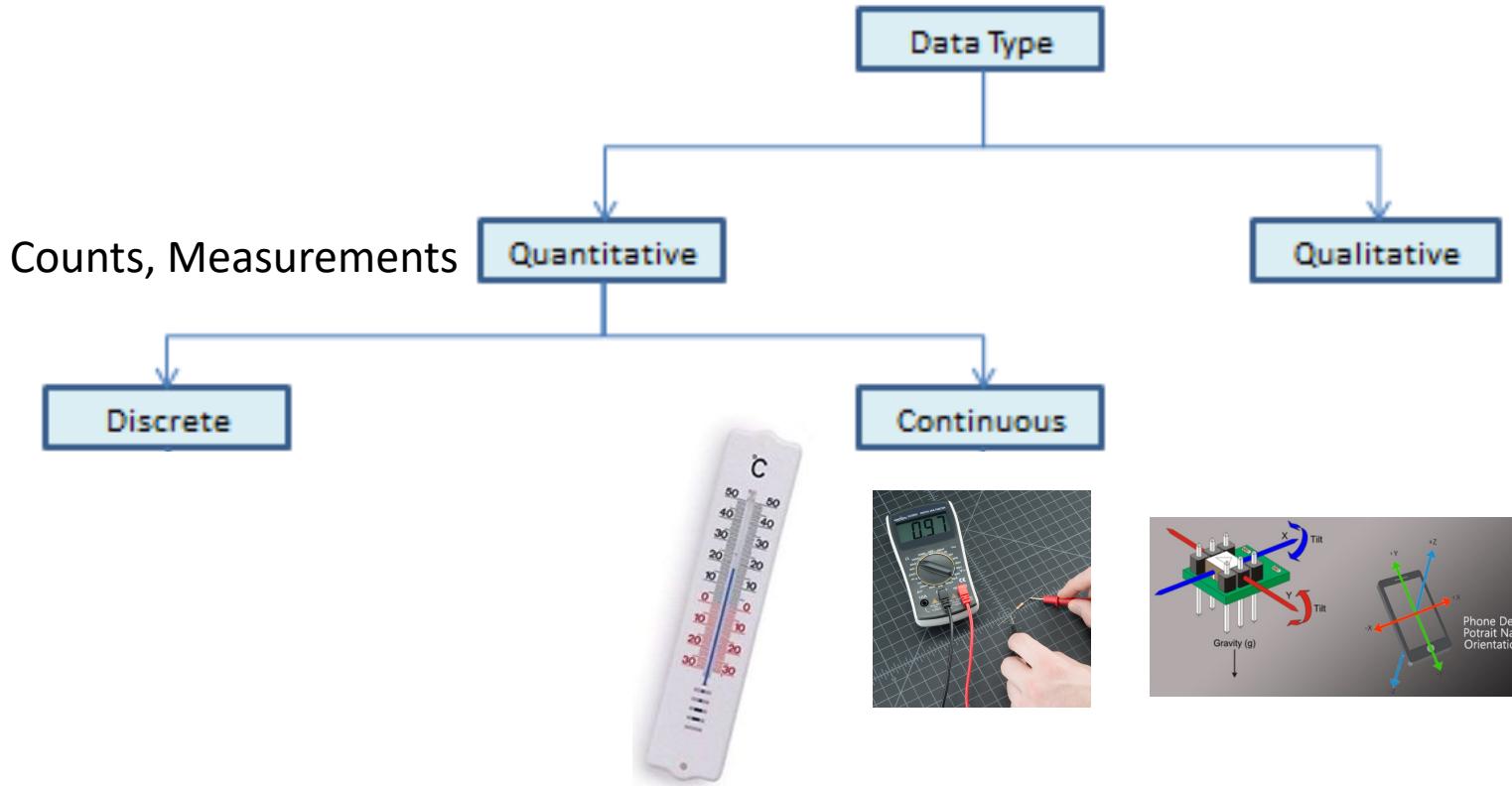
- Cone 3 weighs 79.4 grams
- Cone 2 ice cream is at 8.3°F

Taxonomy of data



- # of CPU cores
- # of courses taken in a semester
- # of times word 'sale' appears in a doc

Taxonomy of data



Samples and Features

Dataset

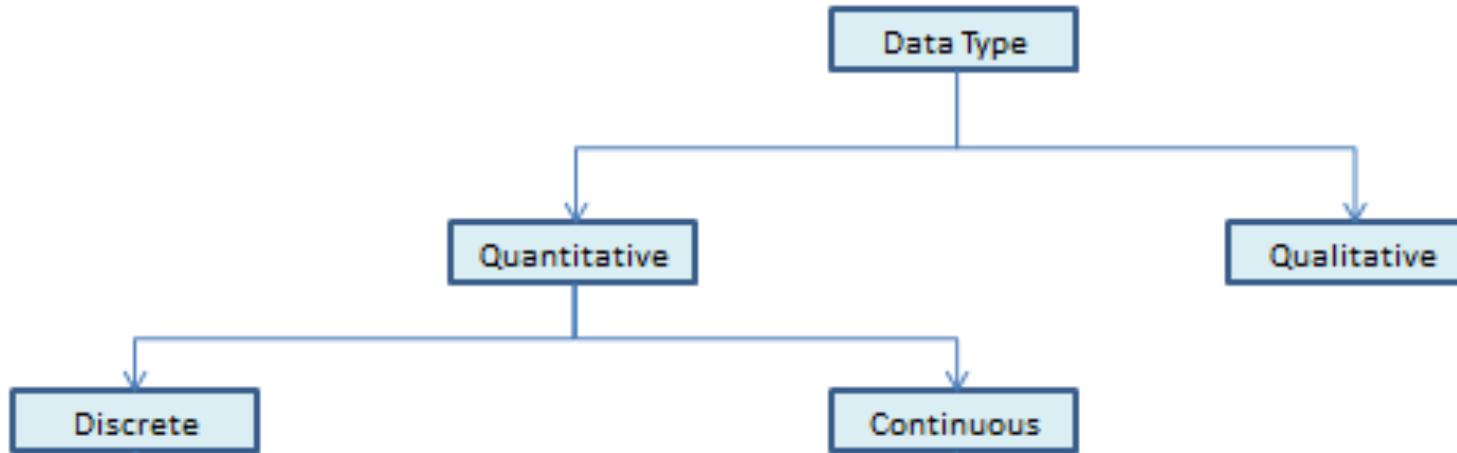
Feature / Attribute

Sample

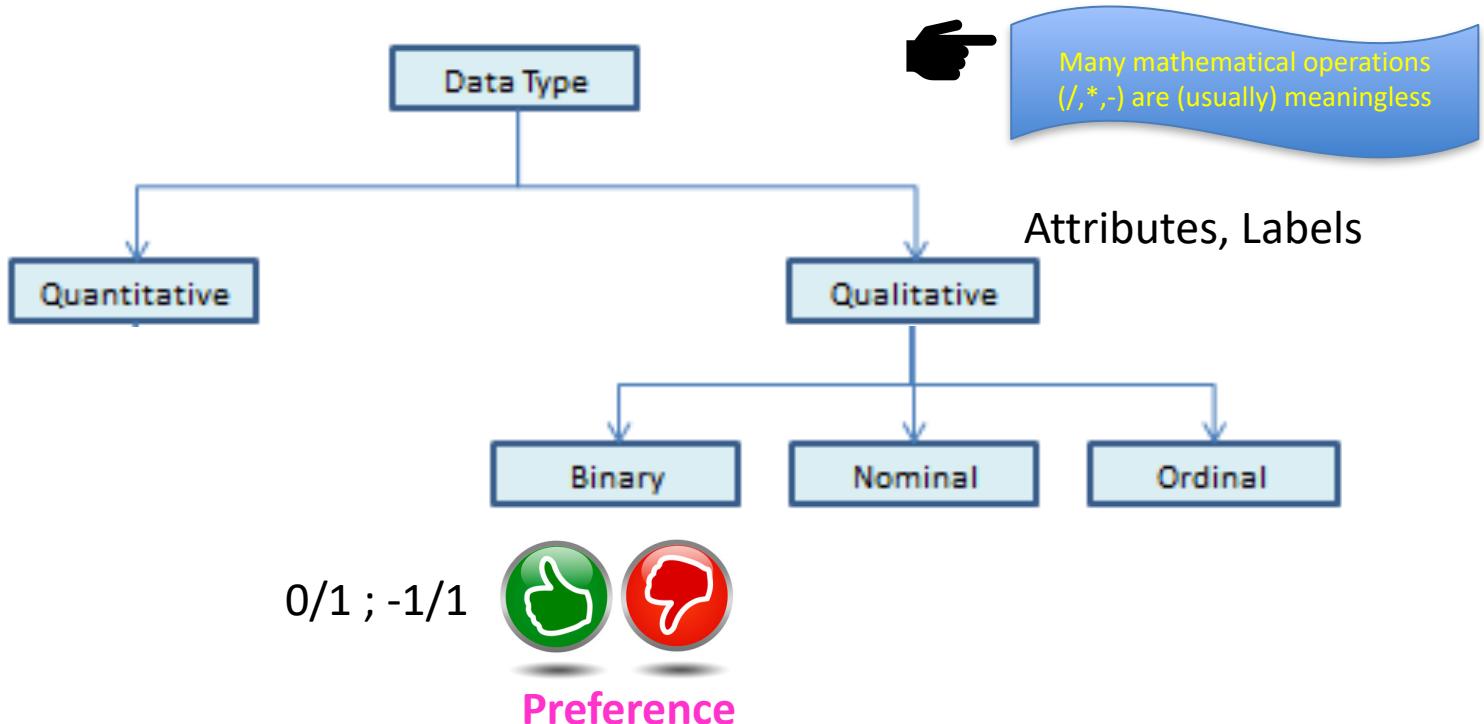
The diagram illustrates a dataset as a table. A large blue arrow points from the word "Dataset" on the left to the table. Another blue arrow points from the word "Sample" on the left to the second column of the table. Above the table, the text "Feature / Attribute" is centered, with a vertical blue arrow pointing downwards to the first row of the table headers.

B	C	D	E	F	G	H	I
Quality	Usage	Dioxane [mol%]	Toluene [mol%]	Cyclohexane [mol%]	Temperature [°C]	{Instrument}	Timestamp
Good	train	18.238	59.40672	22.3555	22.1	RXN1	2019-11-14
Good	train	23.315	37.88732	38.7977	22.2	RXN1	2019-11-14
Good	train	16.405	56.02367	27.5714	22.0	RXN1	2019-11-14
Good	train	41.196	3.06438	55.7395	22.1	RXN1	2019-11-14
Bad	ignore		51.75047		22.2	LTT-R	2019-11-15
Good	test	13.476	67.81965	18.7039	22.5	LTT-R	2019-11-15
Good	test	16.802	13.56112	69.6365	21.9	LTT-R	2019-11-15

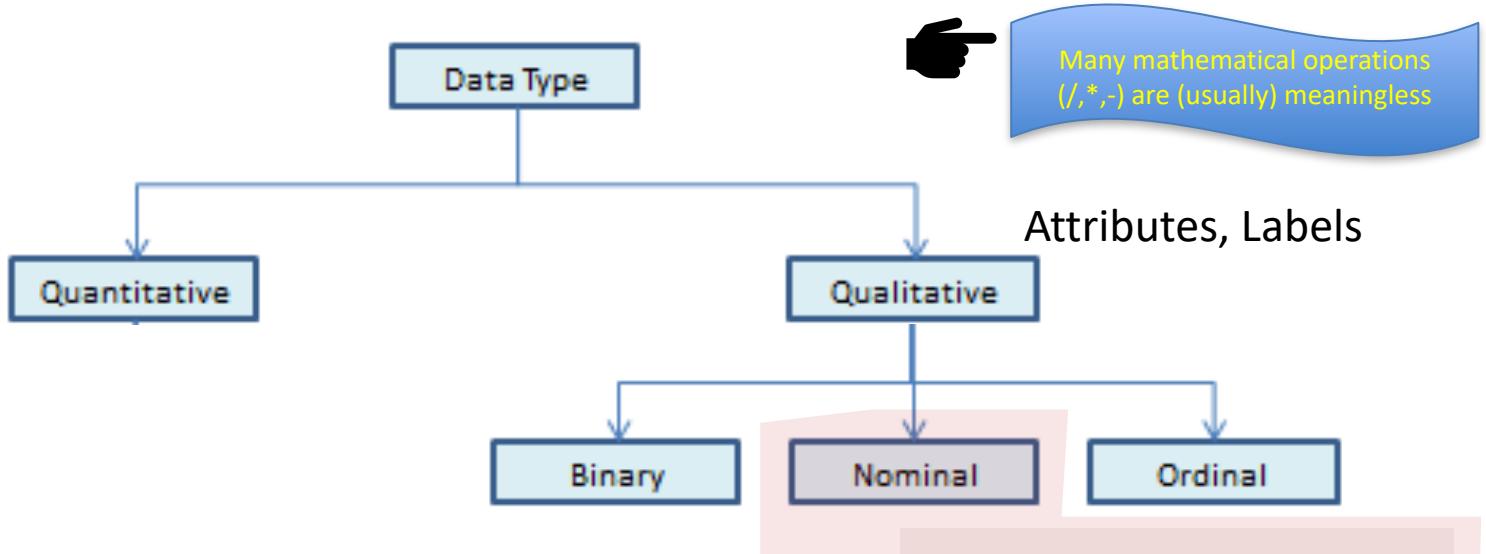
Ultimately, all data needs to be quantitative



Taxonomy of data: Qualitative → Quantitative



Taxonomy of data: Qualitative → Quantitative



Pin Code



Numerical encoding of categorical variables

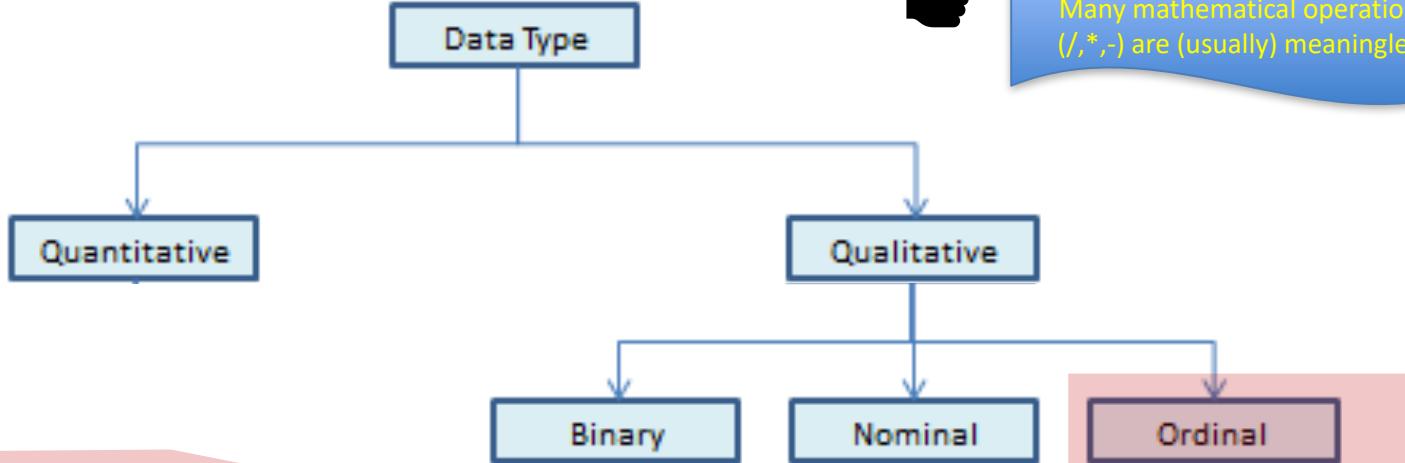
Original data:	
id	Color
1	White
2	Red
3	Black
4	Purple
5	Gold

Numerical encoding of categorical variables

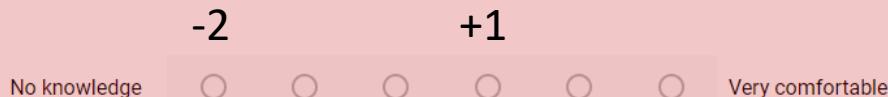
Original data:		One-hot encoding format:					
id	Color	id	White	Red	Black	Purple	Gold
1	White	1	1	0	0	0	0
2	Red	2	0	1	0	0	0
3	Black	3	0	0	1	0	0
4	Purple	4	0	0	0	1	0
5	Gold	5	0	0	0	0	1

Numerical encoding of categorical variables

Rome = [1, 0, 0, 0, 0, 0, ..., 0]
Paris = [0, 1, 0, 0, 0, 0, ..., 0]
Italy = [0, 0, 1, 0, 0, 0, ..., 0]
France = [0, 0, 0, 1, 0, 0, ..., 0]



How comfortable are you with Python *



Letter grade
A +
A
A -
B +
B
B -
C +
C
C -
D +
D
E



Example: Contact Lenses dataset

No patient id

Age is not a number !

Age	Spectacle prescription	Astigmatism	Tear production rate	Recommended lenses
Young	Myope	No	Reduced	None
Young	Myope	No	Normal	Soft
Young	Myope	Yes	Reduced	None
Young	Myope	Yes	Normal	Hard
Young	Hypermetropic	No	Reduced	None
Young	Hypermetropic	No	Normal	Soft
Young	Hypermetropic	Yes	Reduced	None
Young	Hypermetropic	Yes	Normal	hard
Pre-presbyopic	Myope	No	Reduced	None
Pre-presbyopic	Myope	No	Normal	Soft
Pre-presbyopic	Myope	Yes	Reduced	None
Pre-presbyopic	Myope	Yes	Normal	Hard
Pre-presbyopic	Hypermetropic	No	Reduced	None
Pre-presbyopic	Hypermetropic	No	Normal	Soft
Pre-presbyopic	Hypermetropic	Yes	Reduced	None
Pre-presbyopic	Hypermetropic	Yes	Normal	None
Presbyopic	Myope	No	Reduced	None
Presbyopic	Myope	No	Normal	None
Presbyopic	Myope	Yes	Reduced	None
Presbyopic	Myope	Yes	Normal	Hard
Presbyopic	Hypermetropic	No	Reduced	None
Presbyopic	Hypermetropic	No	Normal	Soft
Presbyopic	Hypermetropic	Yes	Reduced	None
Presbyopic	Hypermetropic	Yes	Normal	None

Example: PlayTennis dataset

Outlook	Temperature	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	Normal	False	Yes
...

Outlook	Temperature	Humidity	Windy	Play
Sunny	85	85	False	No
Sunny	80	90	True	No
Overcast	83	86	False	Yes
Rainy	75	80	False	Yes
...

Sometimes data can be missing

Outlook	Temperature	Humidity	Windy	Play
Sunny	85	85	False	No
Sunny	80	<input type="text"/>	True	No
Overcast	83	86	False	Yes
Rainy	75	80	False	Yes
...

→ Unknown or unrecorded

... or incorrect

	DBAName	AKAName	Address	City	State	Zip
t1	John Veliotis Sr.	Johnnyo's	3465 S Morgan ST	Chicago	IL	60608
t2	John Veliotis Sr.	Johnnyo's	3465 S Morgan ST	Chicago	IL	60609
t3	John Veliotis Sr.	Johnnyo's	3465 S Morgan ST	Chicago	IL	60609
t4	Johnnyo's	Johnnyo's	3465 S Morgan ST	Cicago	IL	60608

Conflicts

Does not obey data distribution

Conflict

Data imputation

- Approaches that aim to estimate missing data
- Options
 - Remove sample
 - Fill with 0
 - Fill with constant
 - Fill with a statistical measure (mean, median, mode)
 - Do nothing. Use a learning method which can handle missing data.

Lecture Outline

- *ML Workflow*
- Data sample Representations
- Basic Data Transformations
- Data Visualization

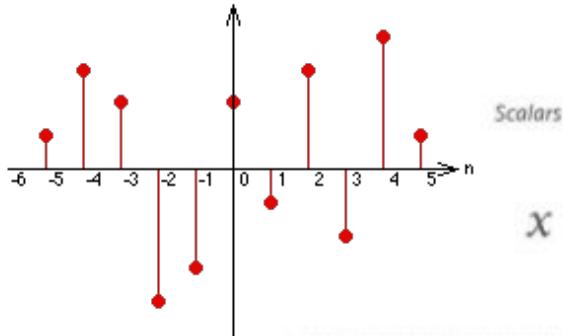
Samples, Features, Labels

The diagram illustrates the components of a dataset:

- Label:** A red box labeled "Label" with a red arrow pointing to the "Quality" column.
- Feature / Attribute:** A blue arrow pointing to the "Dioxane [mol%]" column.
- Sample:** A blue arrow pointing to the first row of the table.
- Table:** A screenshot of a spreadsheet application showing a dataset with columns: B (Quality), C (Usage), D (Dioxane [mol%]), E (Toluene [mol%]), F (Cyclohexane [mol%]), G (Temperature [°C]), H (Instrument), and I (Timestamp). The rows represent individual samples, categorized by Quality (Good or Bad) and Usage (train or test).

B	C	D	E	F	G	H	I
Quality	Usage	Dioxane [mol%]	Toluene [mol%]	Cyclohexane [mol%]	Temperature [°C]	{Instrument}	Timestamp
Good	train	18.238	59.40672	22.3555	22.1	RXN1	2019-11-14
Good	train	23.315	37.88732	38.7977	22.2	RXN1	2019-11-14
Good	train	16.405	56.02367	27.5714	22.0	RXN1	2019-11-14
Good	train	41.196	3.06438	55.7395	22.1	RXN1	2019-11-14
Bad	ignore		51.75047		22.2	LTT-R	2019-11-15
Good	test	13.476	67.81965	18.7039	22.5	LTT-R	2019-11-15
Good	test	16.802	13.56112	69.6365	21.9	LTT-R	2019-11-15

Data Sample Representations



Matrix

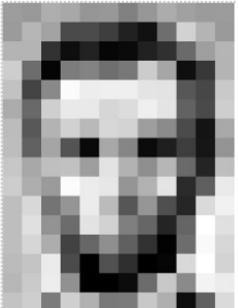
$$X = \begin{bmatrix} x & \cdots & x_N \end{bmatrix} =$$



$$\begin{bmatrix} X_{1,1} & & X_{N,1} \\ \vdots & \ddots & \vdots \\ X_{1,M} & & X_{N,M} \end{bmatrix}$$

1st dimension

2-d image



167	133	174	148	150	162	129	151	72	161	155	156	
168	182	143	74	75	62	81	17	110	210	184	185	
169	180	60	54	54	6	30	48	46	106	159	181	
206	109	9	124	131	131	111	120	204	166	15	96	180
164	187	257	237	239	238	238	227	87	87	71	201	202
170	116	267	203	239	214	220	229	238	238	238	238	238
172	88	179	209	185	215	211	158	139	75	20	169	170
187	97	185	84	10	168	134	11	91	62	22	188	189
166	181	193	183	157	227	174	182	186	186	36	190	191
204	174	286	236	231	241	228	228	245	95	234	234	234
216	116	146	236	187	86	50	79	76	218	241	241	241
204	214	147	109	227	210	177	121	36	255	254	254	254
214	174	173	66	133	143	95	50	2	249	219	219	219
187	196	236	73	8	1	47	0	6	217	256	231	231
183	202	237	146	0	0	128	200	136	243	238	238	238
206	123	207	177	121	120	176	175	175	175	175	175	175

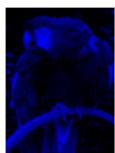
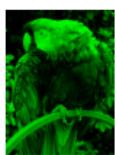
Vectors

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix}$$

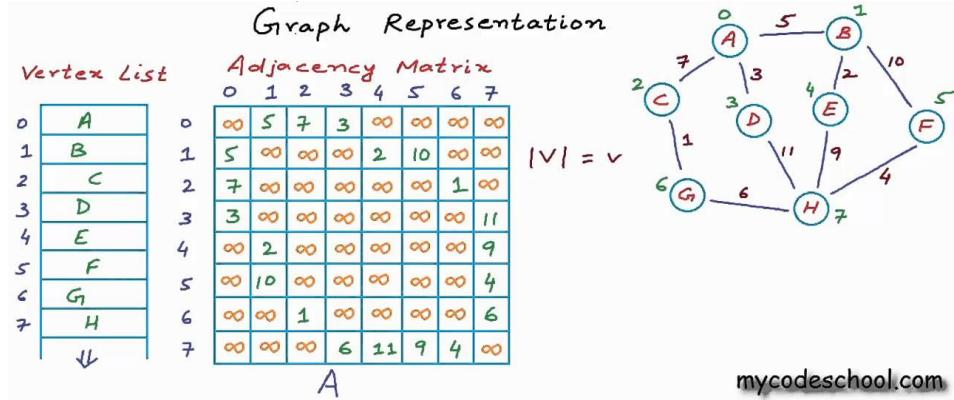
Tensor

$$X = \{X_1, \dots, X_K\} = \begin{bmatrix} X_{1,1,1} & & X_{N,1,1} \\ \vdots & \ddots & \vdots \\ X_{1,M,1} & & X_{N,M,1} \\ & \ddots & \\ & & X_{1,K,K} & & X_{N,K,K} \\ & & \vdots & \ddots & \vdots \\ & & X_{1,M,K} & & X_{N,M,K} \end{bmatrix}$$

row dimension column dimension



Data Representations



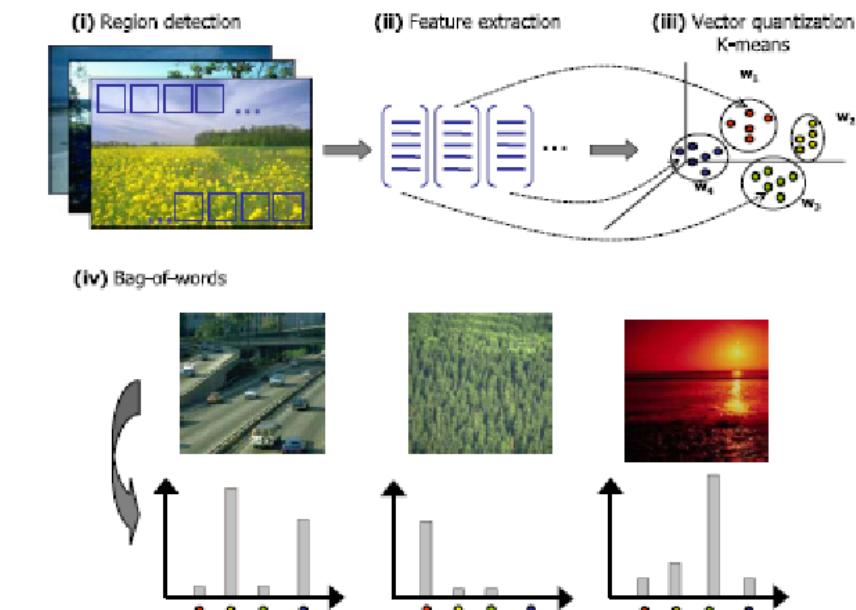
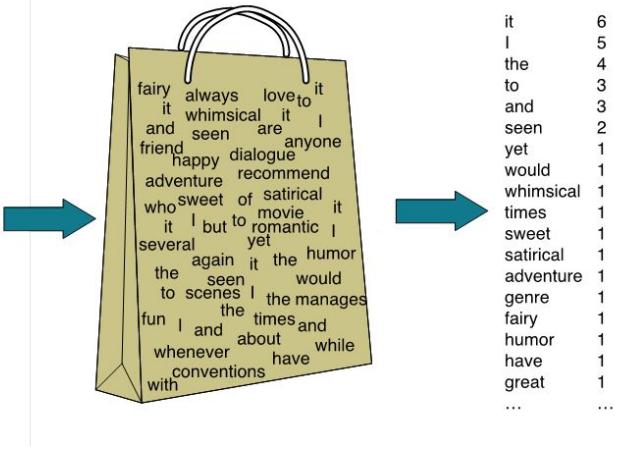
Feature Extraction (FE)

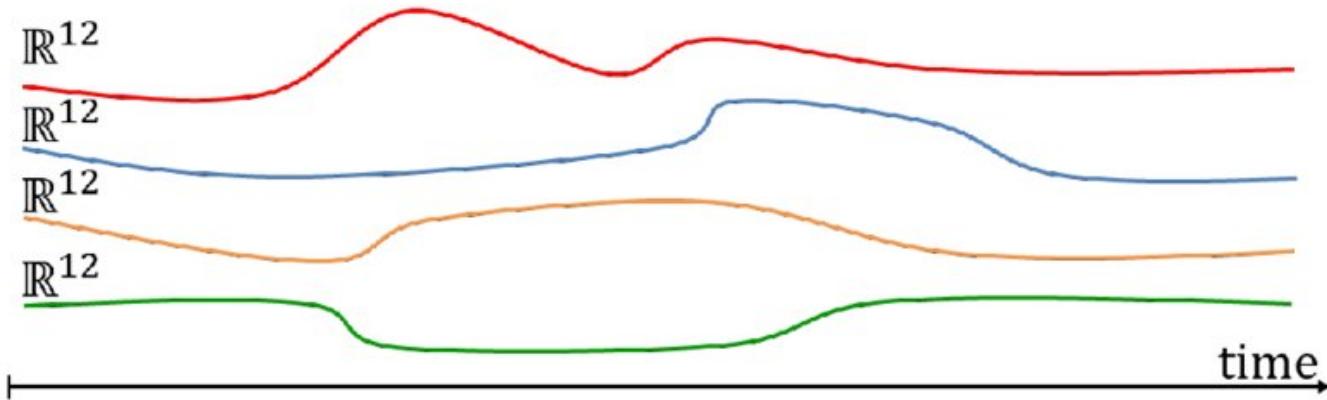
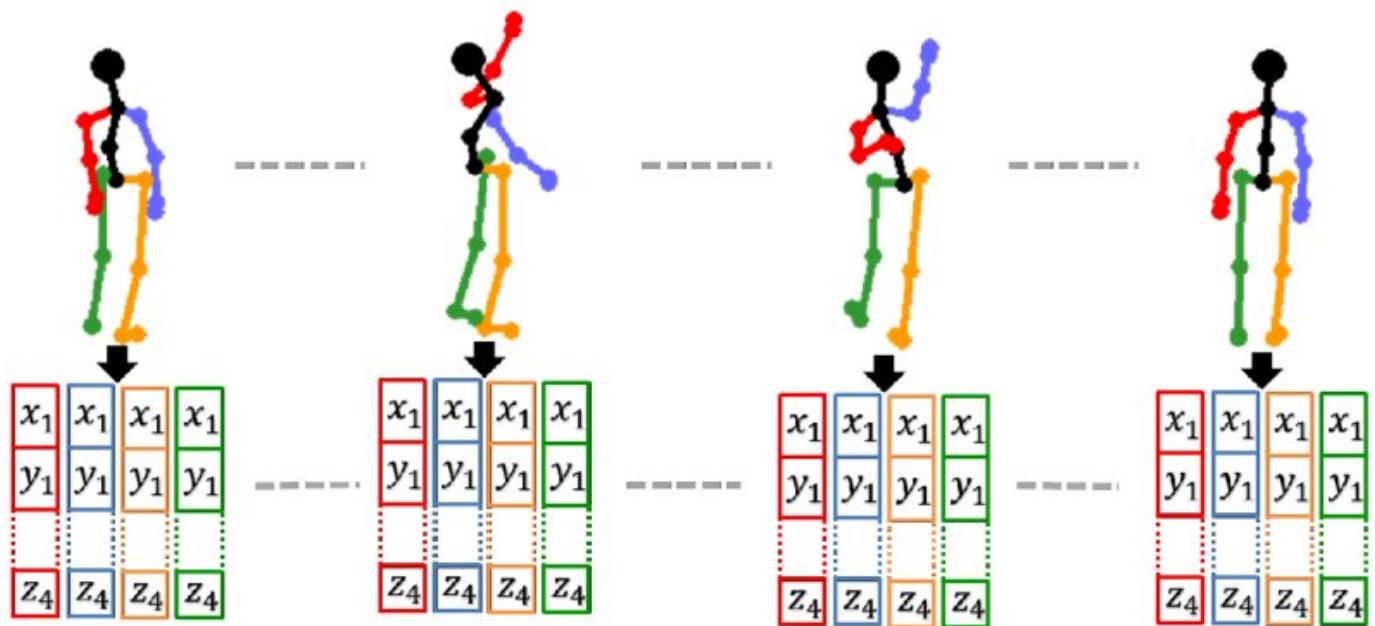
- **Def:** Feature Extraction (FE) is any algorithm that transformation raw data into features that can be used as an input for a learning algorithm.

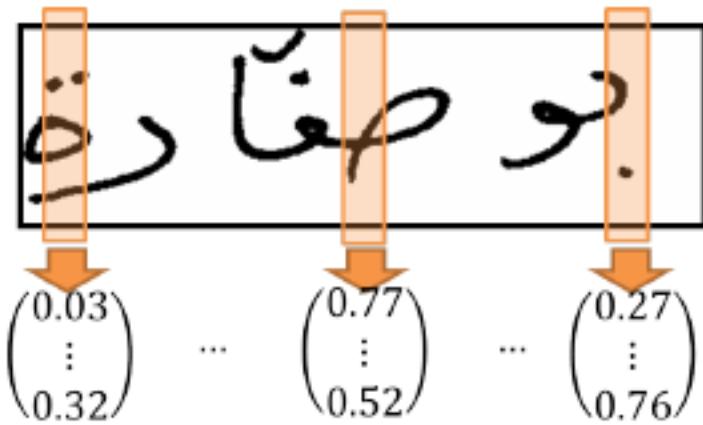
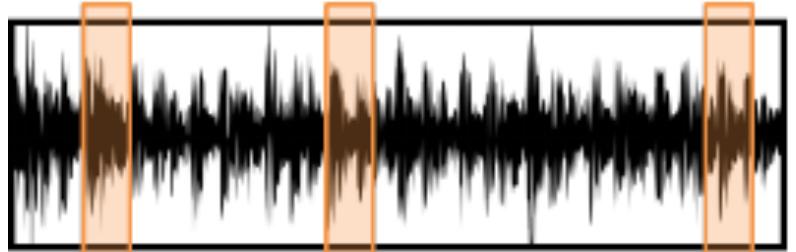
The Bag of Words Representation

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!

15







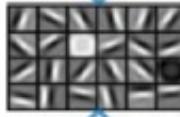
Feature-based, Hierarchical Data Representations



3rd layer
“Objects”



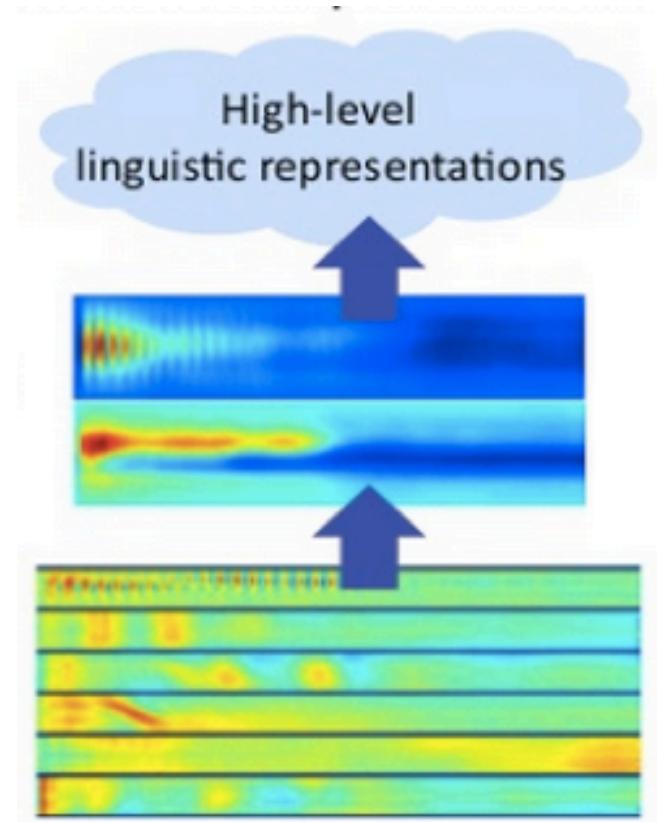
2nd layer
“Object parts”



1st layer
“Edges”

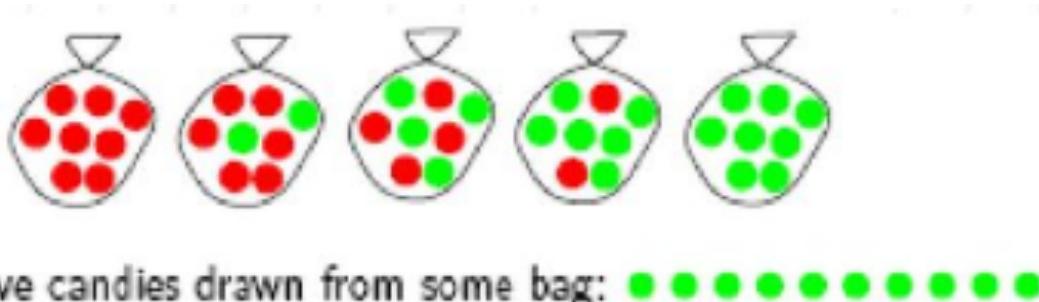


Pixels



Data – a probability-based perspective

- The basis for Statistical Learning Theory



- Then we observe candies drawn from some bag:
• • • • • • • •
- Domain described by random variables (r.v.)
 - $X = \{\text{apple, grape}\}$
 - $b_i \in [1,5]$
- Data = Instantiation of some or all r.v.'s in the domain

Data: a probabilistic perspective

Output					
Proposed Cleaned Dataset					
	DBAName	Address	City	State	Zip
t1	John Veliotis Sr.	Johnnyo's	3465 S Morgan ST	Chicago	IL 60608
t2	John Veliotis Sr.	Johnnyo's	3465 S Morgan ST	Chicago	IL 60609
t3	John Veliotis Sr.	Johnnyo's	3465 S Morgan ST	Chicago	IL 60609
t4	John Veliotis Sr.	Johnnyo's	3465 S Morgan ST	Cicago	IL 60608

Marginal Distribution of Cell Assignments		
Cell	Possible Values	Probability
t2.Zip	60608	0.84
	60609	0.16
t4.City	Chicago	0.95
	Cicago	0.05
t4.DBAName	John Veliotis Sr.	0.99
	Johnnyo's	0.01

DBAName	AKAName	Address	City	State	Zip
t1	John Veliotis Sr.	Johnnyo's	3465 S Morgan ST	Chicago	IL 60608
t2	John Veliotis Sr.	Johnnyo's	3465 S Morgan ST	Chicago	IL 60609
t3	John Veliotis Sr.	Johnnyo's	3465 S Morgan ST	Chicago	IL 60609
t4	Johnnyo's	Johnnyo's	3465 S Morgan ST	Cicago	IL 60608

Conflicts



Does not obey
data distribution

Conflict

Other important aspects of data

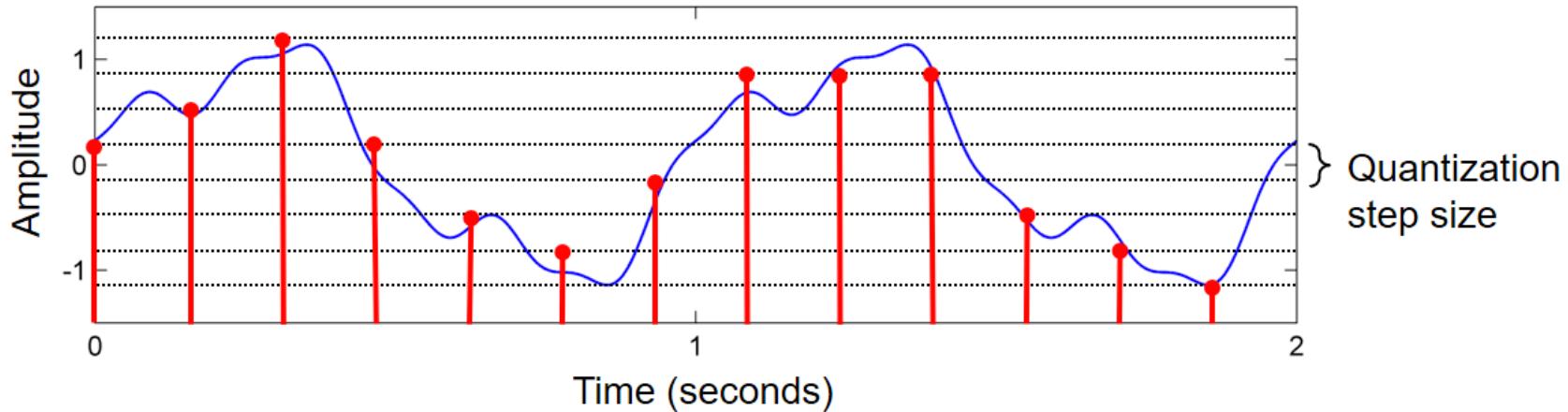
- Mode of collection
 - Passive ('sense')
 - Active ('explore, sense, repeat')
- Statistical assumptions on data
 - i.i.d (independent and identically distributed)
 - Online (e.g. time-series data)

Lecture Outline

- *ML Workflow*
- *Data Representations*
- Basic Data Transformations
- Data Visualization

Quantization

1. Continuous → Discrete ('Rounding off')

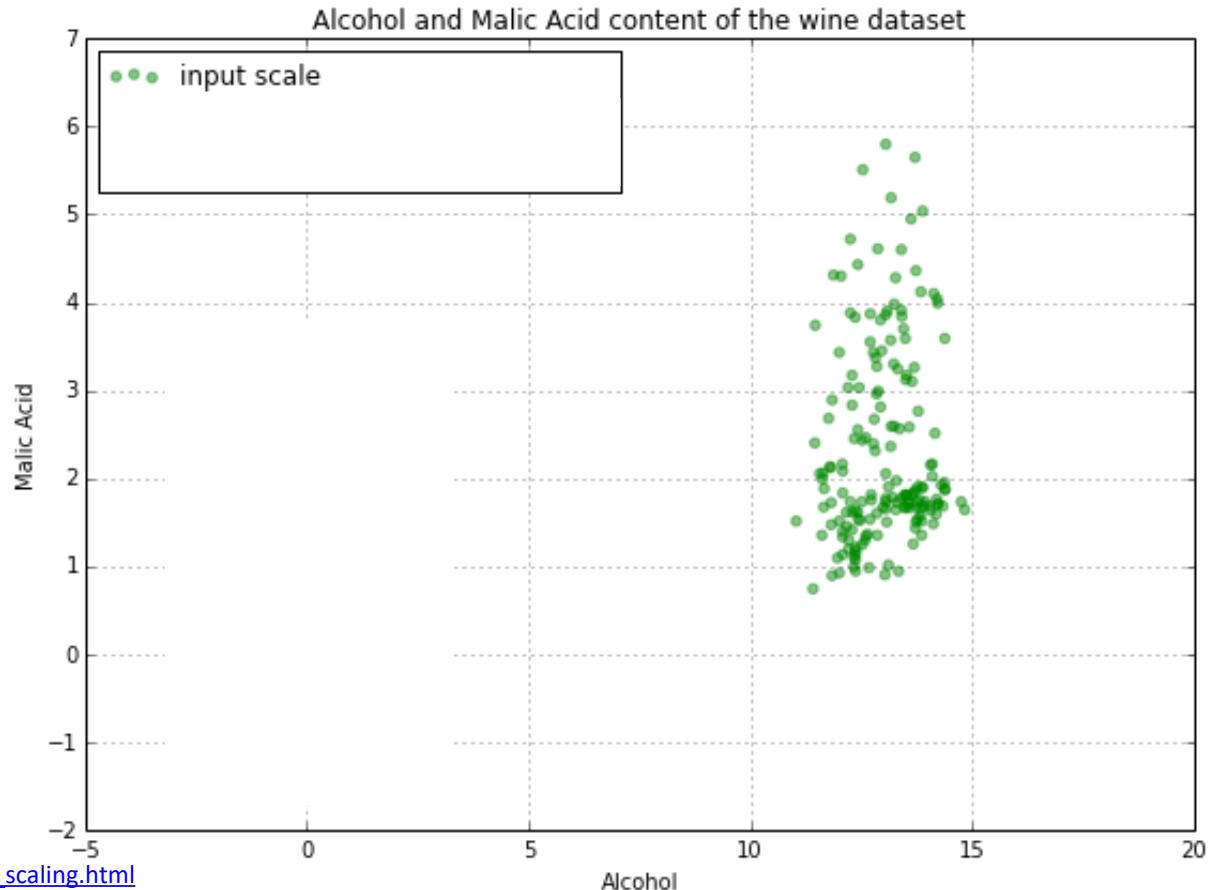


2. Binary Quantization ('Thresholding')

Data Normalization

	Class label	Alcohol	Malic acid
0	1	14.23	1.71
1	1	13.20	1.78
2	1	13.16	2.36
3	1	14.37	1.95
4	1	13.24	2.59

●
●
●
●



Popular normalization approaches

$$X_{norm} = \frac{X - X_{min}}{X_{max} - X_{min}}$$

MinMax Scaling

$$z = \frac{x - \mu}{\sigma}$$

Standardization
(Unit Normal Scaling)

$$\mu = \frac{1}{N} \sum_{i=1}^N (x_i)$$

$$\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2}$$

	Class label	Alcohol	Malic acid
0	1	14.23	1.71
1	1	13.20	1.78
2	1	13.16	2.36
3	1	14.37	1.95
4	1	13.24	2.59

Data Normalization (applied to each feature)

	Class label	Alcohol	Malic acid
0	1	14.23	1.71
1	1	13.20	1.78
2	1	13.16	2.36
3	1	14.37	1.95
4	1	13.24	2.59

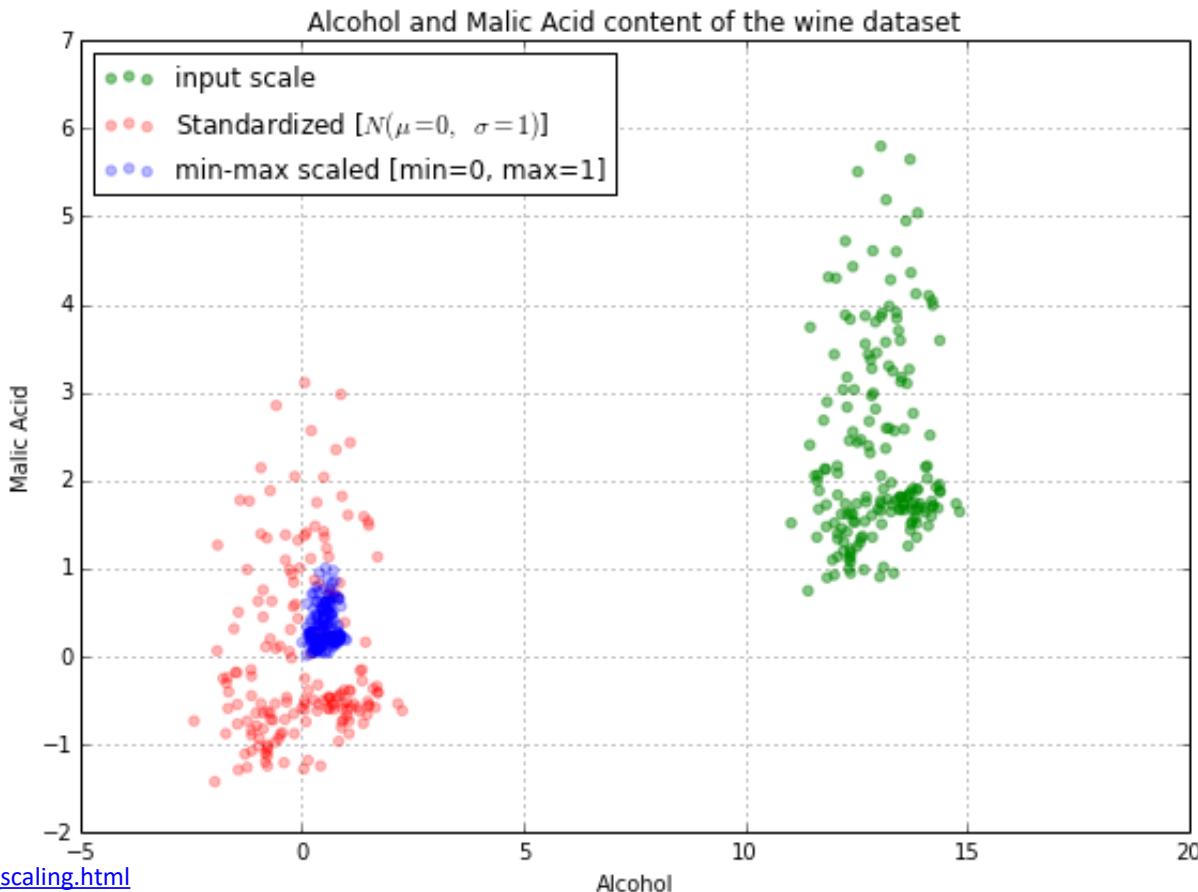
-
-
-
-

$$X_{norm} = \frac{X - X_{min}}{X_{max} - X_{min}}$$

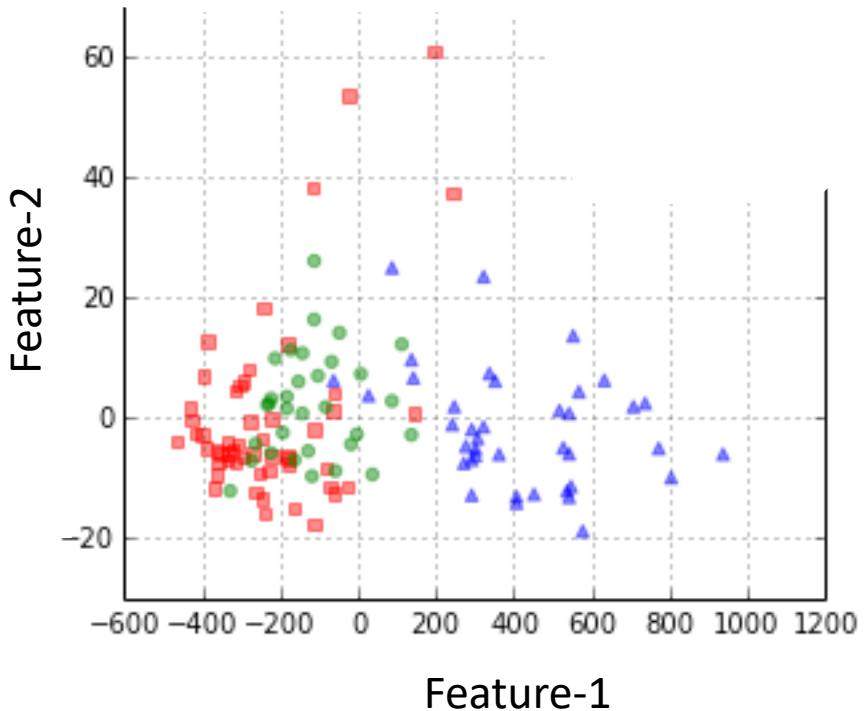
MinMax Scaling

Standardization
(Unit Normal Scaling)

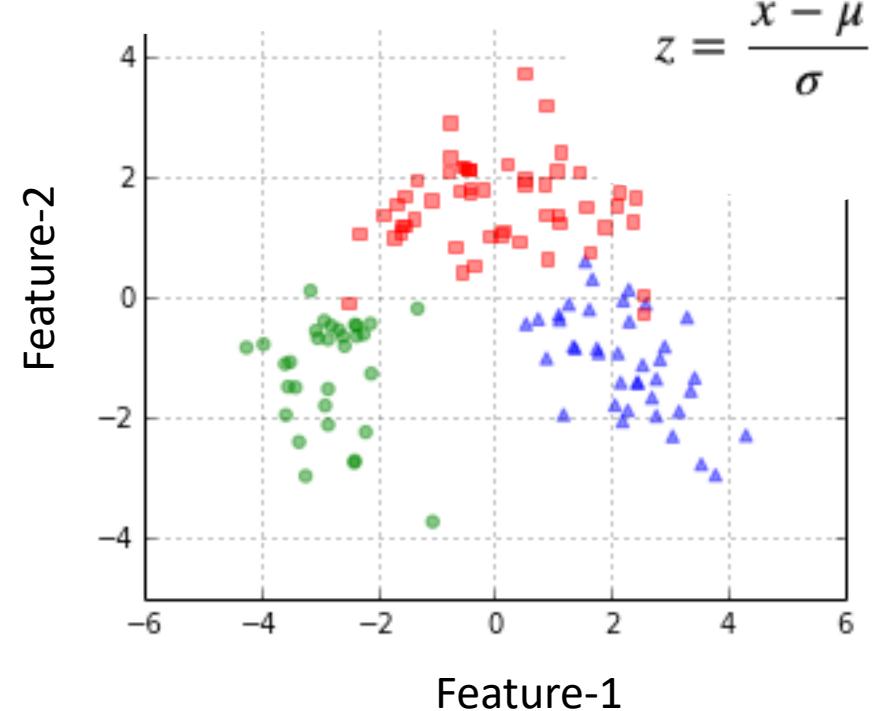
$$z = \frac{x - \mu}{\sigma}$$



Before standardization



After standardization



Why normalize data ?

- Uniform treatment of all features
- (Empirically) Helps stabilize optimization and lead to faster convergence
- Disadvantages?

Workflow of a Machine Learning Problem

