

Lecture Notes for **Machine Learning in Python**



Professor Eric Larson
Data Quality and Imputation

Class Logistics and Agenda

- Logistics:
 - Need help? canvas has links to various resources
 - the class GitHub is also a resource!
 - TA hours!
 - Team Forming Discussion sections
- Agenda:
 - Data Quality
 - Data Representations
 - Imputation methods

Course Github Page:	https://github.com/eclarson/MachineLearningNotebooks ↗
Other Useful Guides:	Helpful Links and Guides for Semester
Participation For Distance Students	Turn in answers to questions here: Participation

Using the SMU HPC

- Tutorials available for various types of analysis:
 - <https://www.smu.edu/oit/research>

Events This Month

September 2025						
S	M	T	W	T	F	S
	1	2	3	4	5	6
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30				
Upcoming Events:						
Introduction to Python						
1:30pm - 2:30pm, Wednesday, September 17, 2025						
Data Science, AI/ML Office hour						
10:00am - 1:00pm, Wednesday, September 24, 2025						
Data Science workflow with Python						
1:30pm - 2:30pm, Wednesday, September 24, 2025						
HPC and ColdFront drop in support						
10:00am - 12:00pm, Monday, September 29, 2025						
Data Science, AI/ML Office hour						
10:00am - 1:00pm, Wednesday, October 1, 2025						
Introduction to Neural Network with Pytorch						
1:30pm - 2:30pm, Wednesday, October 1, 2025						
View the calendar of events at SMU						

Class Overview, by topic

Table Data
Visualization

Numpy, Pandas, Seaborn
Overviews with some in-depth discussion

Dimension
Reduction and
Image Processing

Scikit-learn, Scikit Image,
Intuition only, Some mathematics

Linear and
Logistic
Regression

Numpy, Recreate API for Scikit-learn
Detailed mathematics for simple optimization
intuition for advanced optimization

Neural Networks
and Back Prop.

Numpy
Detailed mathematics for NN operations

Wide and Deep
Networks

Convolutional
Networks

Recurrent
Networks

Keras, Tensorflow
Intuition, Detailed implement.

Ethics in
Language Models

ConceptNet
Case studies

Last Time

Data Quality Problems

- Missing
 - Easy to find, NaNs
- Duplicated
 - Easy to find, hard to verify
- Noise or Outlier
 - Hard to define
 - Hard to catch

Discrete	Nominal or Categorical	Variable could be one value in a set of categories. No ordering of values. <i>Example: Employee ID</i>	Allowed Transforms permuting values boolean, one hot encoding, or hash function
	Ordinal	Variable could be one value in a set of categories. Ordering matters. <i>Example: Star Ratings, 1-5</i>	Allowed Transforms $V_{new} = f_{mono}(V_{old}) + b$ integer (or boolean)
Continuous	Interval or Numeric	Value is continuous numeric value. Could be in specified range. <i>Example: BMI, Temperature, etc.</i>	Allowed Transforms $V_{new} = f_{mono}(V_{old}) + b$ float
	Ratio or Numeric	Value is continuous numeric value. Zero is meaningful. Often not treated differently than interval. <i>Example: Length, Elevation</i>	Allowed Transforms $V_{new} = f_{mono}(V_{old})$ float

Split-Impute-Combine

TID	Pregnant	BMI	Age	Diabetes
1	Y	33.6	41-50	positive
2	N	26.6	31-40	negative
3	Y	23.3	?	positive
4	N	28.1	21-30	negative
5	N	43.1	31-40	positive
6	Y	25.8	21-30	negative
7	Y	31.0	21-30	positive
8	Y	35.3	?	negative
9	N	30.5	51-60	positive
10	Y	37.6	51-60	positive



split: pregnant
split: BMI > 32

TID	Pregnant	BMI	Age	Diabetes
1	Y	>32	41-50	positive
8	Y	>32	?	negative
10	Y	>32	51-60	positive

Mode: none, can't impute

TID	Pregnant	BMI	Age	Diabetes
3	Y	<32	?	positive
6	Y	<32	21-30	negative
7	Y	<32	21-30	positive

Mode: 21-30

TID	Hair Color	Height	Age	Arrested
1	Brown	5'2"	23	no
2	Hazel	1.5m	12	no
3	Bl	5	999	no
4	Brown	5'2"	28	no

Self Test, Missing data

- Can all missing data be found by searching for NaNs?
 - A. Yes. Missing data should always be a NaN.
 - B. Yes. Pandas defaults all missing data to NaN.
 - C. No. This only works for floats, because that is the data type for NaN.
 - D. No. NaN only represents missing data that is already found.

K-Nearest Neighbors Imputation

TID	Pregnant	BMI	Age	Diabetes
1	Y	33.6	31-40	positive
2	N	26.6	31-40	negative
3	Y	23.3	?	positive
4	?	28.1	21-30	negative
5	N	43.1	31-40	positive
6	Y	25.6	21-30	negative
7	Y	31.0	21-30	positive
8	Y	35.3	?	negative
9	N	30.5	51-60	positive
10	Y	37.6	21-30	positive

$f_i^{(unk)}$

For $k = 3$, find 3 closest neighbors

TID	Preg.	BMI	Age	Diabetes	Distance d_k
3	Y	23.3	?	positive	0
6	Y	25.6	21-30	negative	$(0 + 2.3 + 1)/3$
2	N	26.6	31-40	negative	$(1 + 3.3 + 1)/3$
4	?	28.1	21-30	negative	$(4.8 + 1)/2$

... repeat for all rows, select 3 closest ...

Imputed Age: 21-30

Distance can be calculated differently:

- Difference for valid features only
- May need to normalize ranges
- Weight neighbors differently?
- Have min # of valid features?
- Type: Euclidean, city-block, etc.

$$d_k = \frac{1}{|F_{valid}|} \sum_{i \in F_{valid}} \|f_i^{(unk)} - f_i^{(k)}\|$$

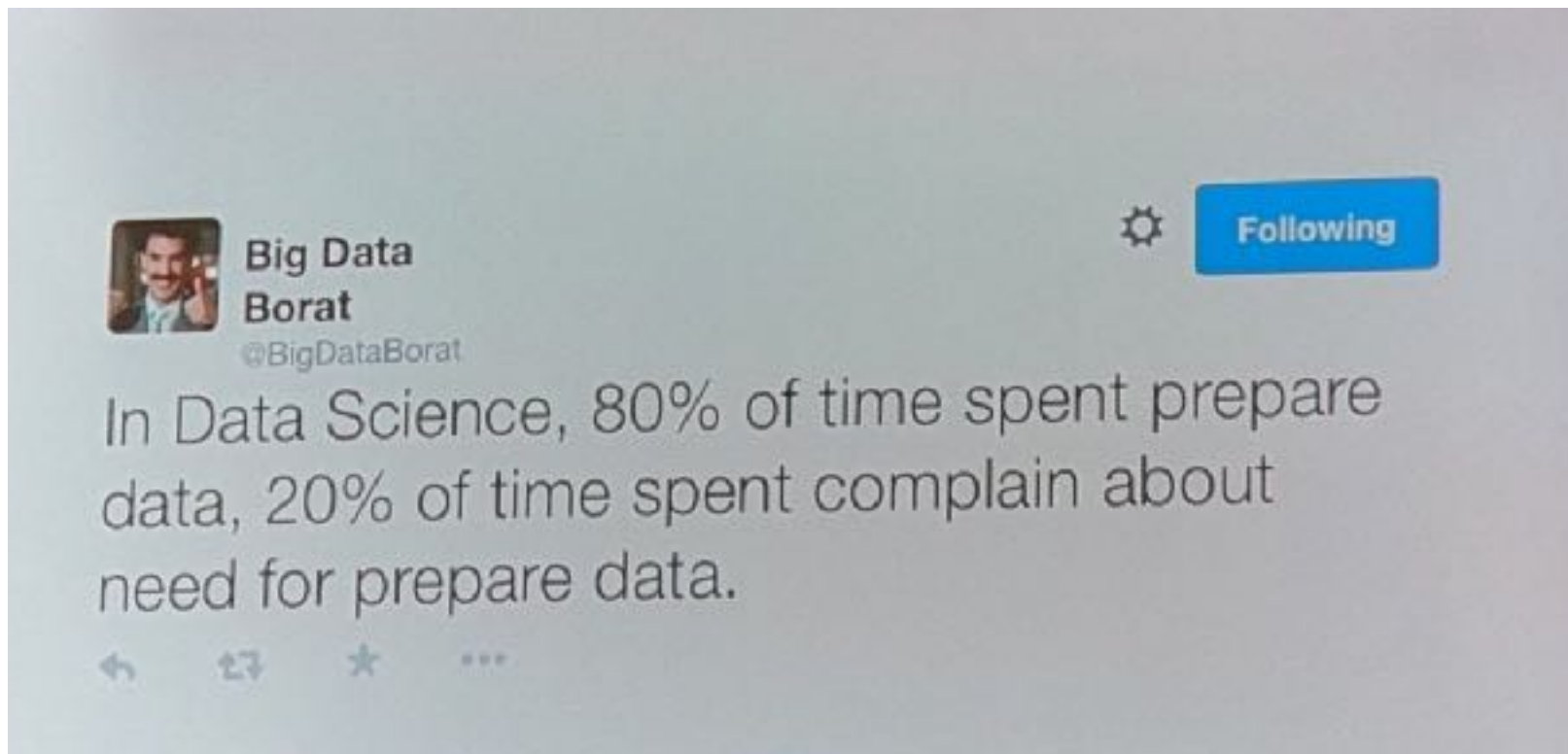
\uparrow
 i^{th} feature, f , in row

Pandas and Imputation
Scikit-Learn



03. Data Visualization.ipynb

Data Representation and Documents



Data Tables as Variable Representations

Table

<i>TID</i>	<i>Pregnant</i>	<i>BMI</i>	<i>Age</i>	<i>Eye Color</i>	<i>Diabetes</i>
1	Y	33.6	41-50	brown	positive
2	N	26.6	31-40	hazel	negative
3	Y	23.3	31-40	blue	positive
4	N	28.1	21-30	brown	inconclusive
5	N	43.1	31-40	blue	positive
6	Y	25.6	21-30	hazel	negative

Internal Rep.

<i>TID</i>
1
2
3
4
5
6

Data Tables as Variable Representations

Table

<i>TID</i>	<i>Pregnant</i>	<i>BMI</i>	<i>Age</i>	<i>Eye Color</i>	<i>Diabetes</i>
1	Y	33.6	41-50	brown	positive
2	N	26.6	31-40	hazel	negative
3	Y	23.3	31-40	blue	positive
4	N	28.1	21-30	brown	inconclusive
5	N	43.1	31-40	blue	positive
6	Y	25.6	21-30	hazel	negative

Internal Rep.

<i>TID</i>	<i>Binary</i>	<i>Float</i>	<i>Ordinal</i>	<i>Object</i>	<i>Diabetes</i>
1	1	33.6	2	hash(0)	1
2	0	26.6	1	hash(1)	0
3	1	23.3	1	hash(2)	1
4	0	28.1	0	hash(0)	2
5	0	43.1	1	hash(2)	1
6	1	25.6	0	hash(1)	0

Bag of words model

<i>TID</i>	<i>Pregnant</i>	<i>BMI</i>	<i>Chart Notes</i>	<i>Diabetes</i>
1	Y	33.6	Complaints of fatigue wh...	positive
2	N	26.6	Sleeplessness and some...	negative
3	Y	23.3	First saw signs of rash o...	positive
4	N	28.1	Came in to see Dr. Steve...	inconclusive
5	N	43.1	First diagnosis for hospit...	positive
6	Y	25.6	N/A	negative

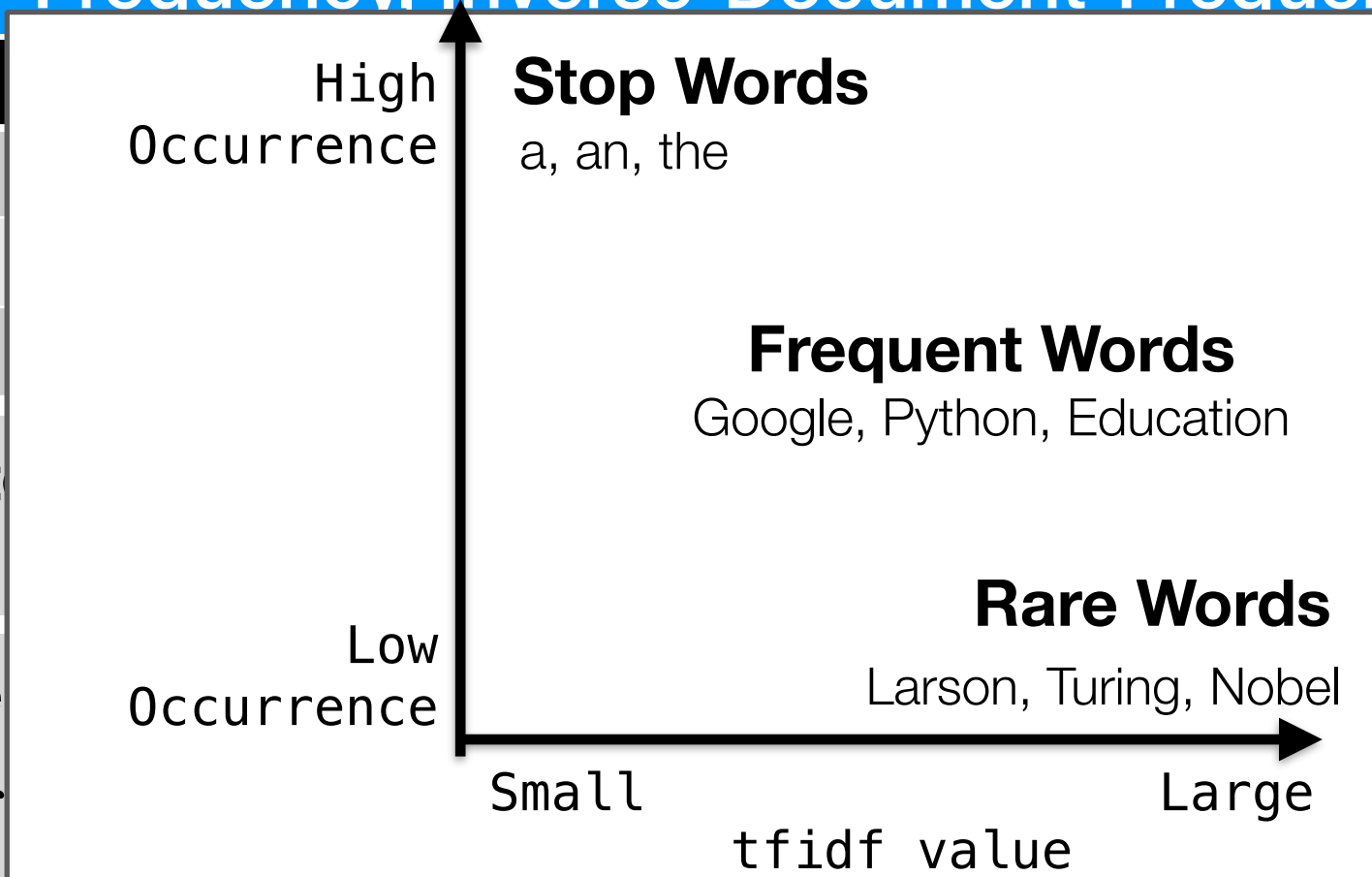
Bag of Words

Vocabulary						
TID	Sleep	Fatigue	Weight	Rash	First	Sight
1	0	1	0	0	2	0
2	1	1	0	0	1	1
3	1	1	0	2	1	1

number of occurrences

Term-Frequency, Inverse-Document-Frequency

TID	Slee
1	0
2	0.1
3	0.1



tc	Stev
.86	0
.02	0
.1	0

This is often used in RAG systems, for Keyword retrieval!!

Want to know more?

Take Natural Language Processing!

For Next Lecture

- Before next class:
 - verify installation of seaborn, plotly, (and/or bokeh if you want)
 - look at pandas table data and additional tutorials
- Next time: Data Visualization

Lecture Notes for **Machine Learning in Python**

Professor Eric Larson
Data Quality and Imputation