

Statistical Methods in AI (CS7.403)

Lecture-3: Data Visualization, Intro to Performance Measures, Benchmarking

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

<https://ravika.github.io>

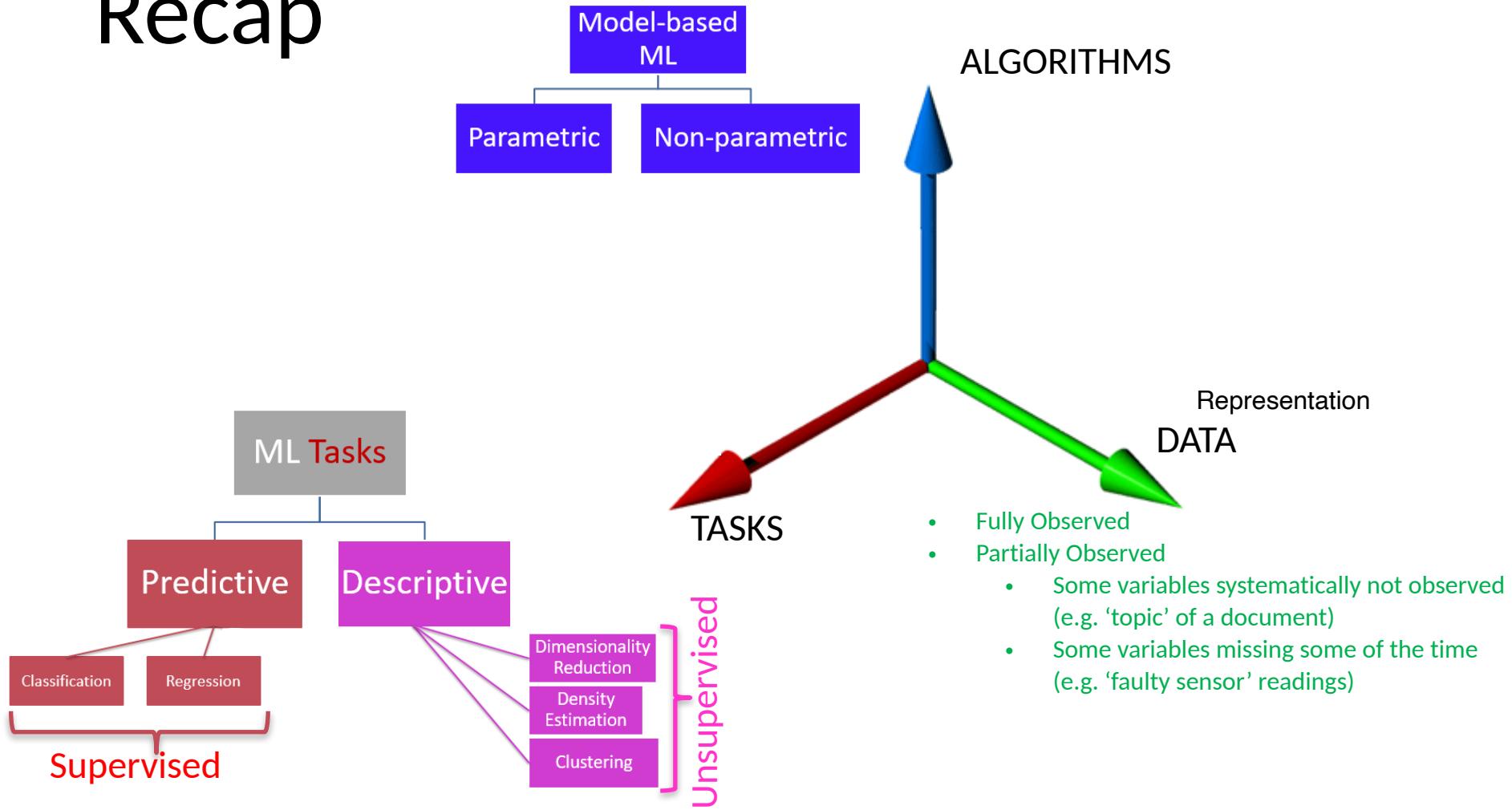


Center for Visual Information Technology (CVIT)

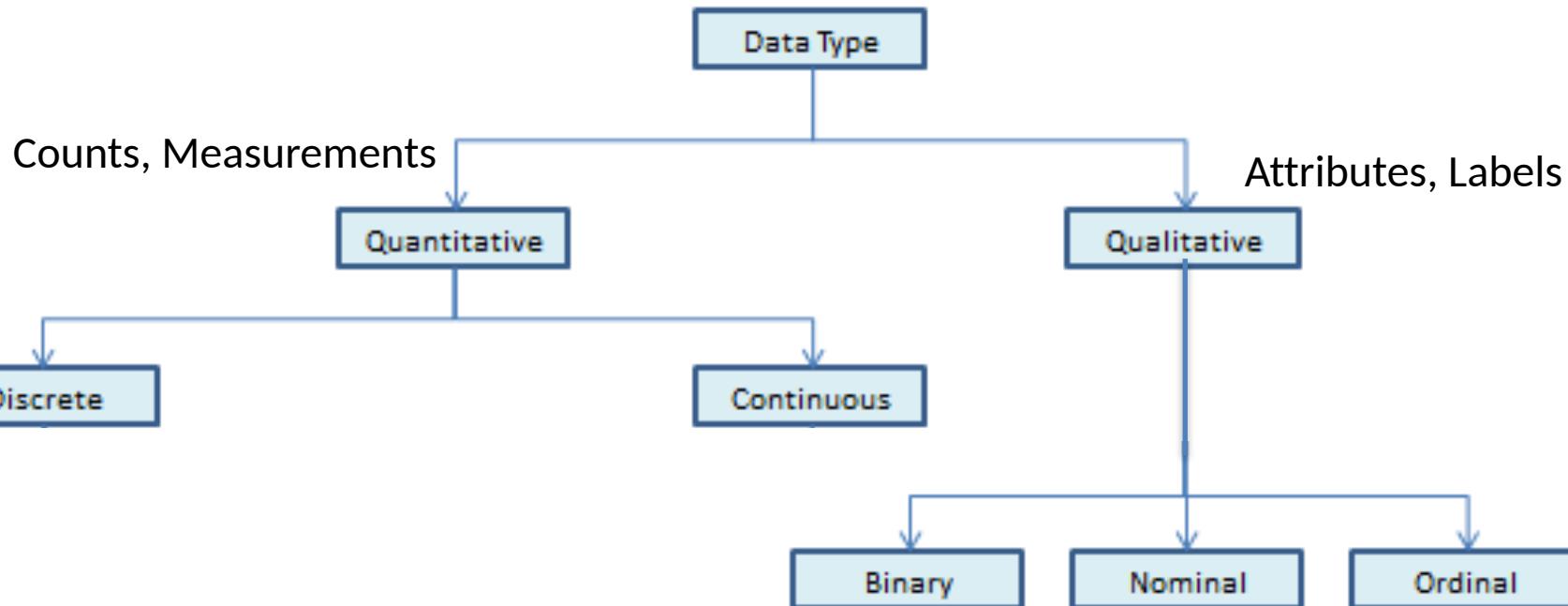
IIIT Hyderabad



Recap



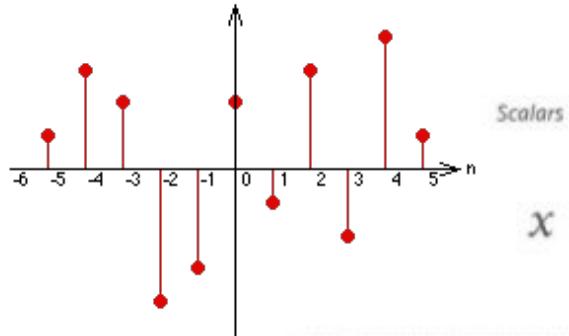
Taxonomy of data



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]

Data Sample Representations



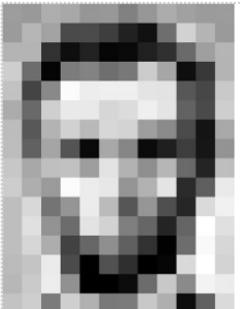
Matrix

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 & \dots & \mathbf{x}_N \end{bmatrix} = \begin{bmatrix} \mathbf{x}_{1,1} & & \mathbf{x}_{N,1} \\ & \ddots & \\ \mathbf{x}_{1,M} & & \mathbf{x}_{N,M} \end{bmatrix}$$

1st dimension

2nd dimension

2-d image



157	163	174	168	160	152	129	151	172	161	155	156
155	182	163	74	75	62	89	17	110	210	180	154
180	180	50	14	54	6	10	33	48	104	159	181
206	109	6	124	131	111	120	204	164	15	56	180
194	83	137	251	237	239	230	228	227	87	71	231
172	104	207	233	233	214	220	220	228	88	74	206
188	81	179	209	186	216	211	158	139	75	22	169
189	97	165	64	15	168	134	11	91	62	22	146
199	168	191	193	158	227	176	143	182	105	36	190
205	174	188	262	236	231	140	178	228	63	95	234
190	216	116	149	236	187	86	150	76	38	218	241
190	234	147	108	227	210	127	102	36	101	255	224
190	214	173	66	113	143	91	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	213
183	202	237	140	6	5	13	108	200	138	143	236
195	206	133	207	177	121	125	200	175	19	36	218

Vectors

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

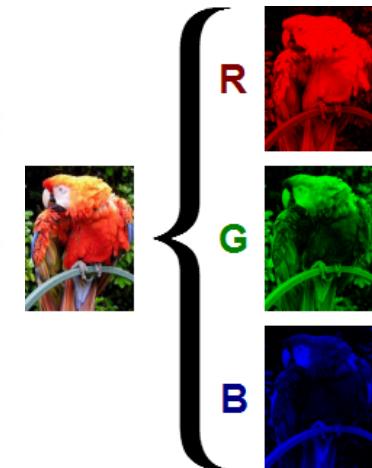
Tensor

$$\mathbf{X} = \{\mathbf{X}_1, \dots, \mathbf{X}_K\} = \begin{bmatrix} \mathbf{x}_{1,1,1} & & \mathbf{x}_{N,1,1} \\ \vdots & \cdots & \vdots \\ \mathbf{x}_{1,M,1} & & \mathbf{x}_{N,M,1} \\ & \ddots & \\ & & \mathbf{x}_{1,1,K} & \cdots & \mathbf{x}_{N,1,K} \\ & & \vdots & \cdots & \vdots \\ & & \mathbf{x}_{1,M,K} & & \mathbf{x}_{N,M,K} \end{bmatrix}$$

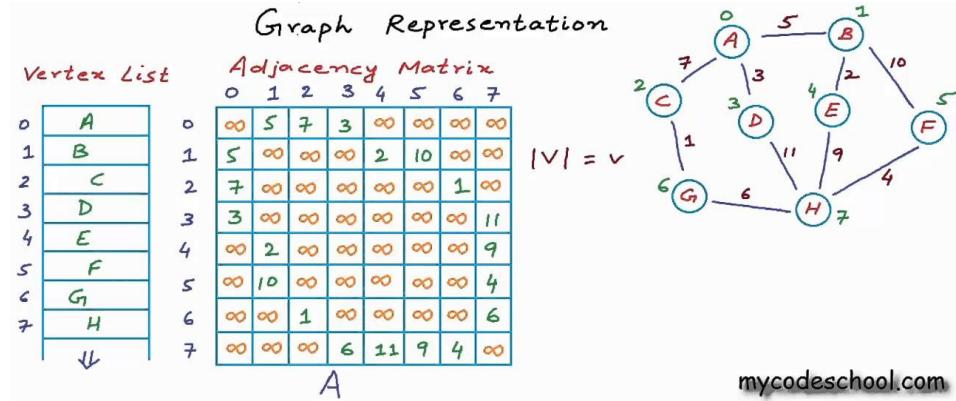
1st dimension

2nd dimension

3rd dimension

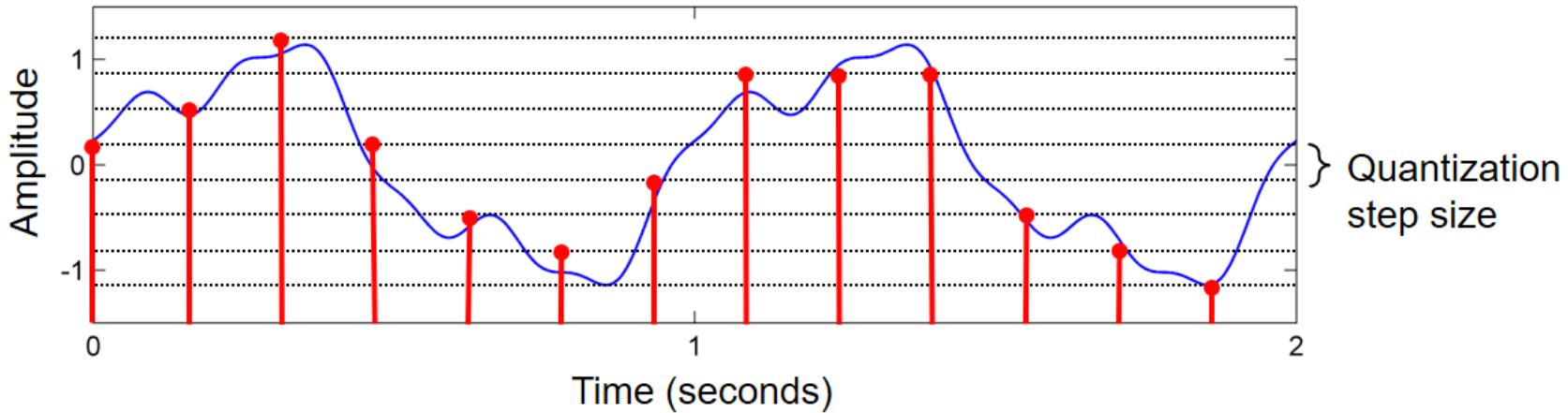


Data Sample Representations



Quantization

1. Continuous → Discrete ('Rounding off')



2. Binary Quantization ('Thresholding')

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

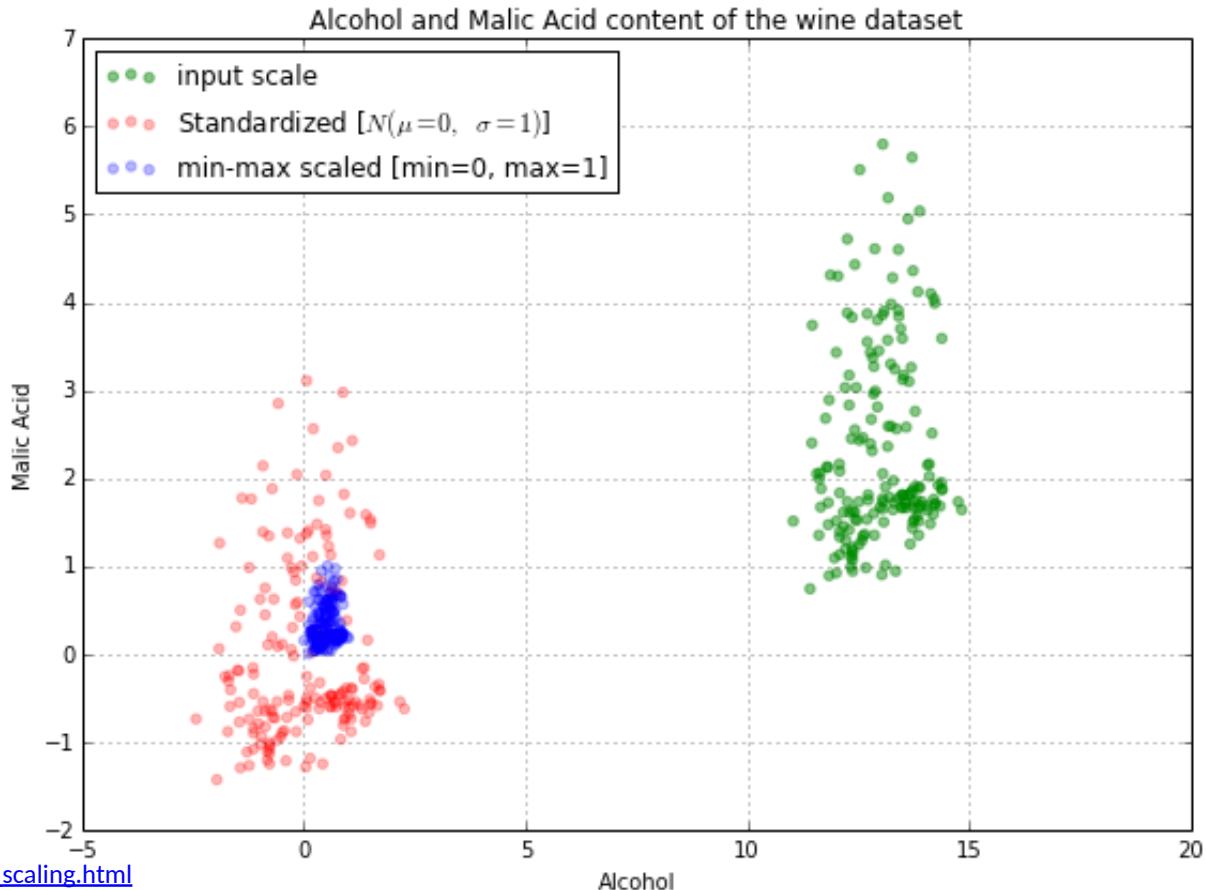
-
-
-
-

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

MinMax Scaling

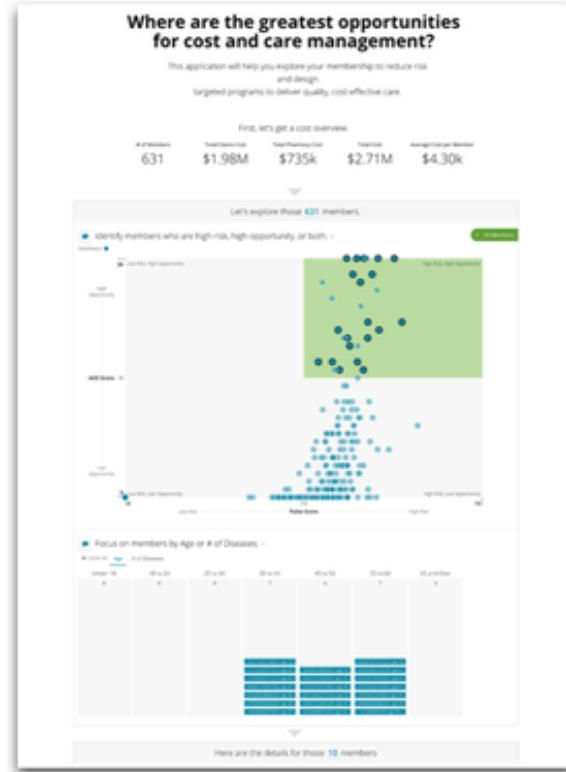
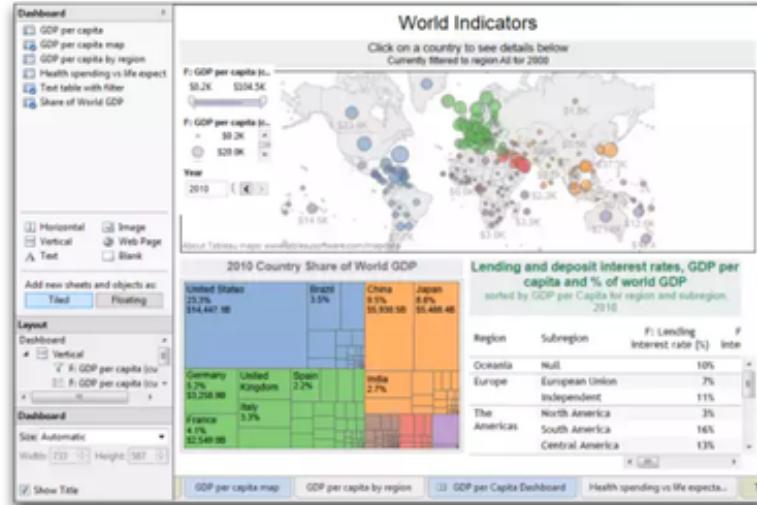
Standardization
(Unit Normal Scaling)

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



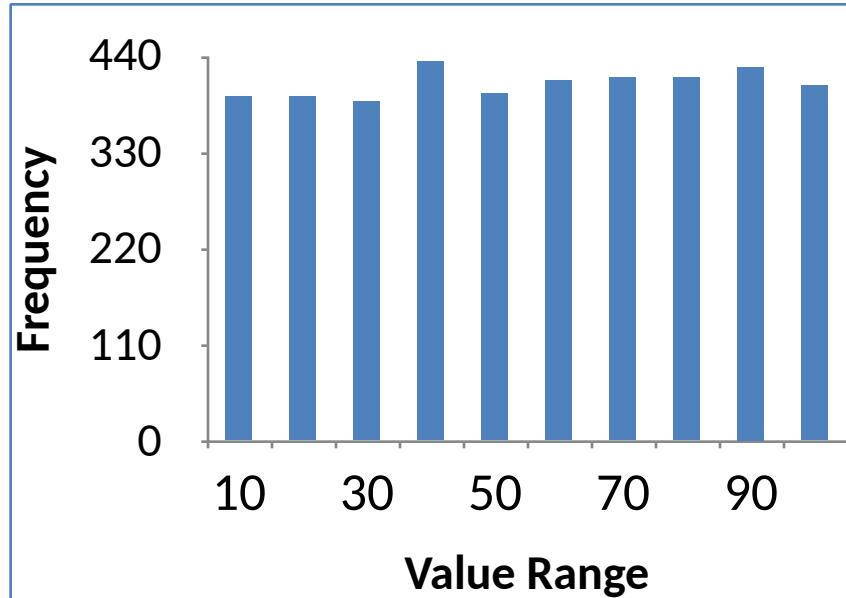
Gazing at Data: Data visualization

data exploration data presentation



Why Plot Data?

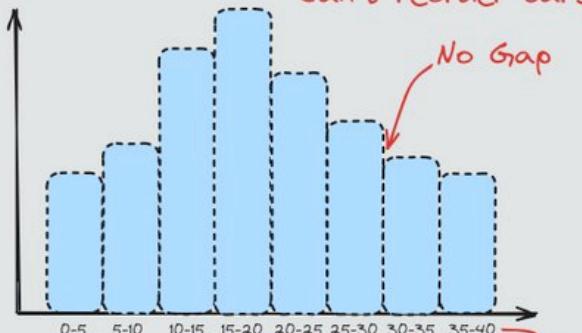
22.65	42.12	67.24	59.13	81.49
23.03	53.42	40.54	89.97	21.85
12.07	93.43	51.93	49.30	43.76
47.68				
86.20	51.91	13.12	73.88	60.29
20.02	41.28	66.24	62.15	46.87
48.38				
55.23	92.09	26.50	83.53	70.99
65.30	46.21	10.85	29.61	62.15
75.71				
69.73	84.90	15.37	35.00	83.23
77.95	26.56	5.78	72.59	12.47
58.90				
6.45	93.15	3.67	49.80	43.05
32.39	53.77	82.80	43.59	32.35
74.35				
14.94	63.71	9.30	1.31	
82.87				
47.05	42.53	62.74	99.91	53.17
34.26				
12.54	46.29	67.34	32.65	23.94
	57.39	10.61	54.07	53.28
	60.10	2.25	77.55	12.05
	17.02	80.73	29.60	9.96
	97.01	19.84	76.59	45.90
	86.80	19.11	4.80	1.24
	30.40	67.94	55.53	58.25
			73.13	0.23



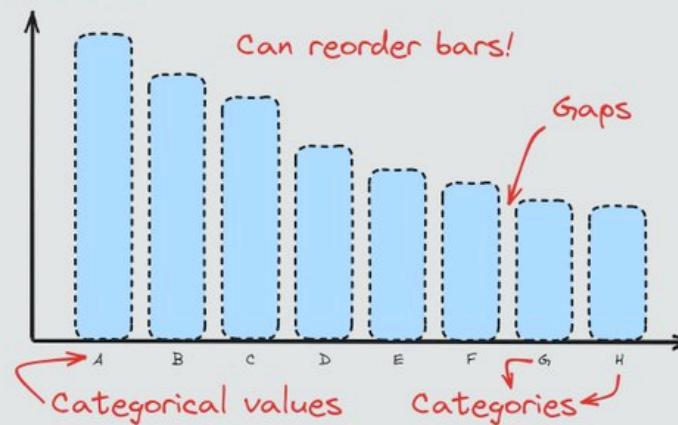
- Visualization of data provides specific insights into the nature of the data.
- Depending on the plot, we gain different insights

Barchart Vs Histogram

Histogram:

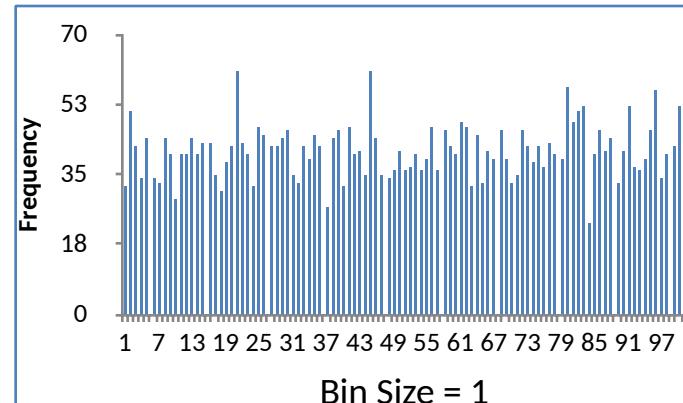
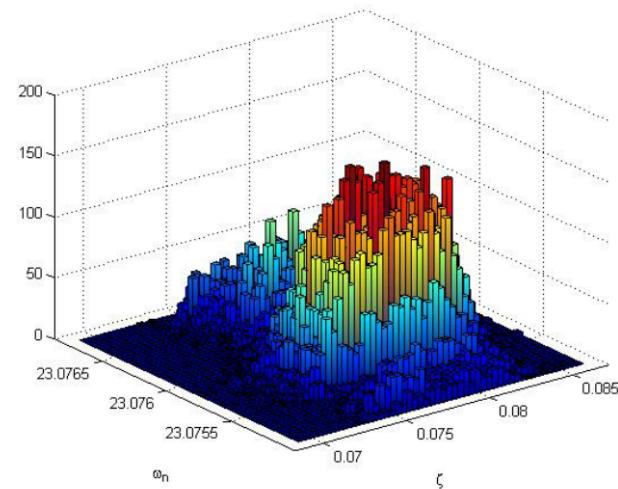
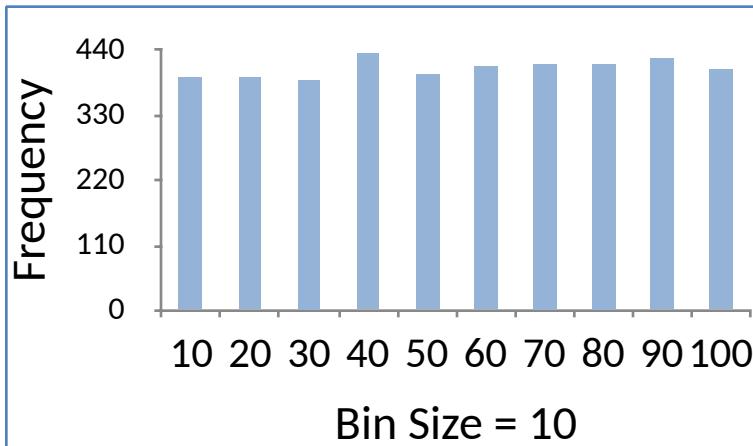


Barchart:



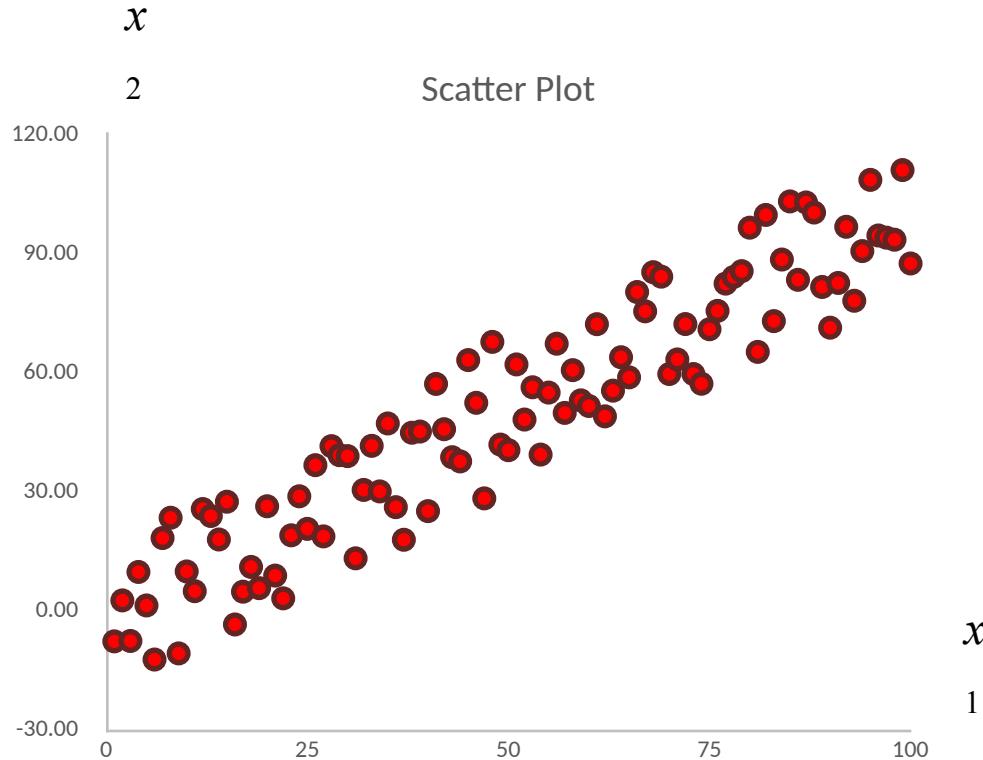
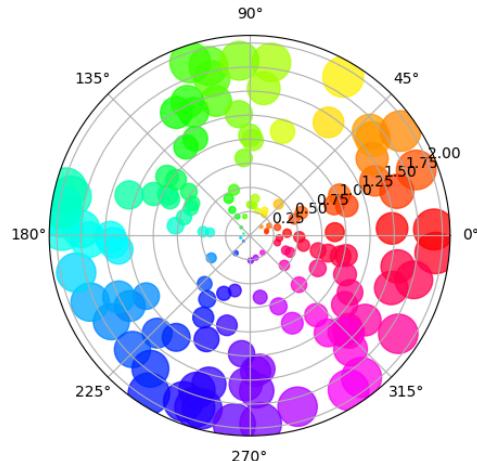
Histogram

- Count of items in each bin
 - Not a bar chart of Data
 - Approximation of Distribution
- Visualize one feature at a time
- Possible to extend to two
- Dependency on bins (\sqrt{n} ; $2\sqrt[3]{n}$)



Scatter Diagram

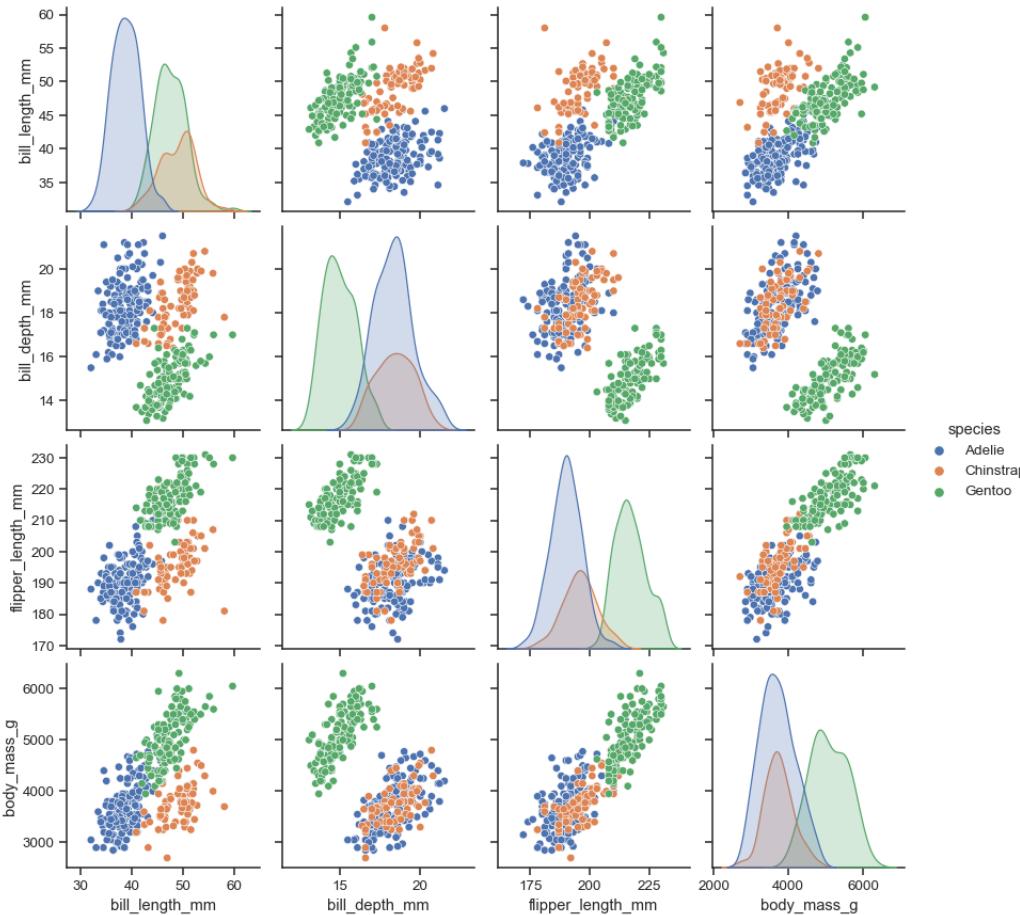
- Plots two features at a time
- Captures the correlation between the two
- Other formats possible



Polar Plot Courtesy: Scatterplot Documentation [matplotlib.org]

Pair Plot

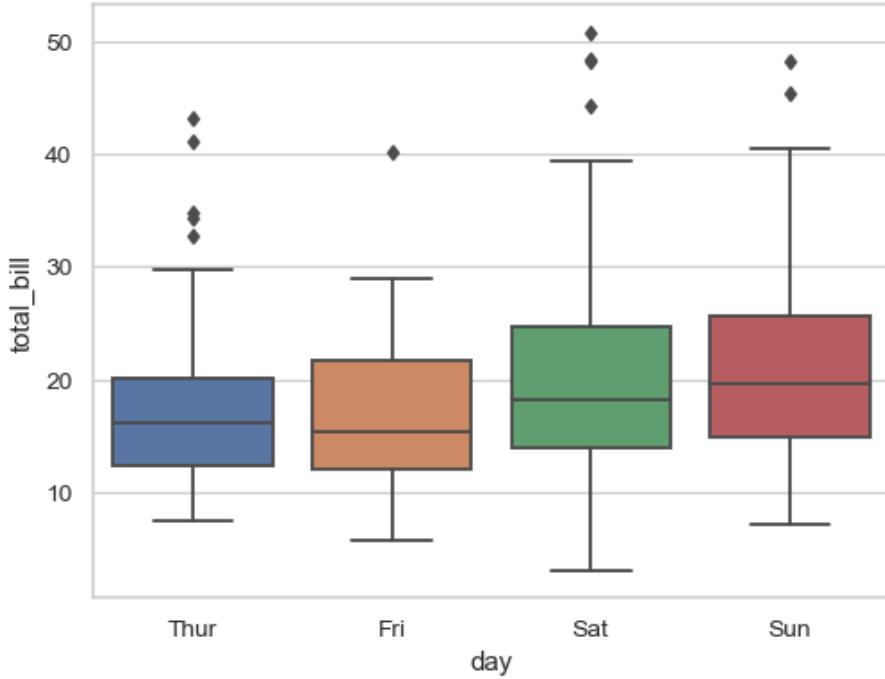
- Plot each pair of features as a matrix
 - Diagonal entries are histograms (densities)
 - Off-diagonal entries are scatter plots
- Can use other plots at each cell.



Plot Courtesy: Seaborn Documentation [seaborn.pydata.org]

Box Plot

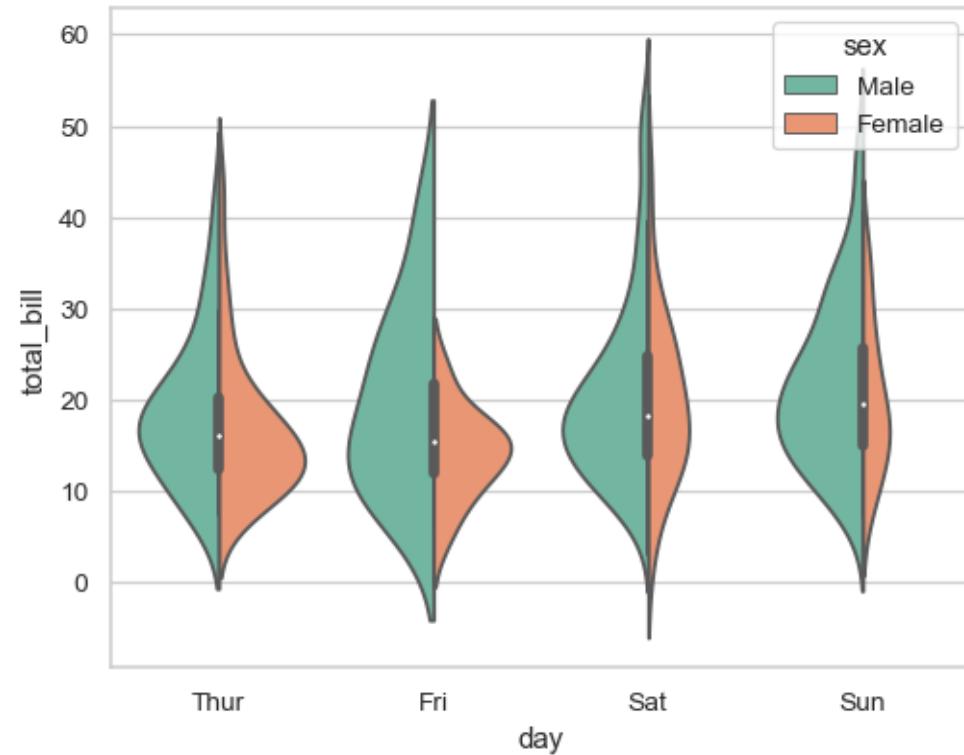
- Show median and quartiles of each feature
 - Outliers are removed
 - Box-and-whisker plot
- Whiskers can represent other percentiles/data
- Simpler than histograms of each feature



Plot Courtesy: Seaborn Documentation [seaborn.pydata.org]

Violin Plot

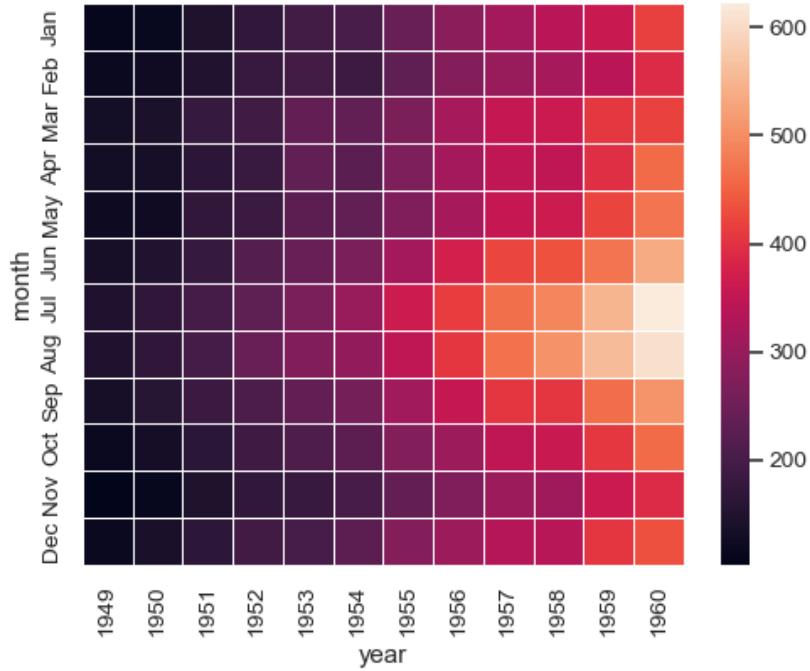
- Shows the **density plot** of each feature
 - Similar to Box Plot
- Either side can represent different densities
- Densities are smoothed estimates from data



Plot Courtesy: Seaborn Documentation [seaborn.pydata.org]

Heat Map

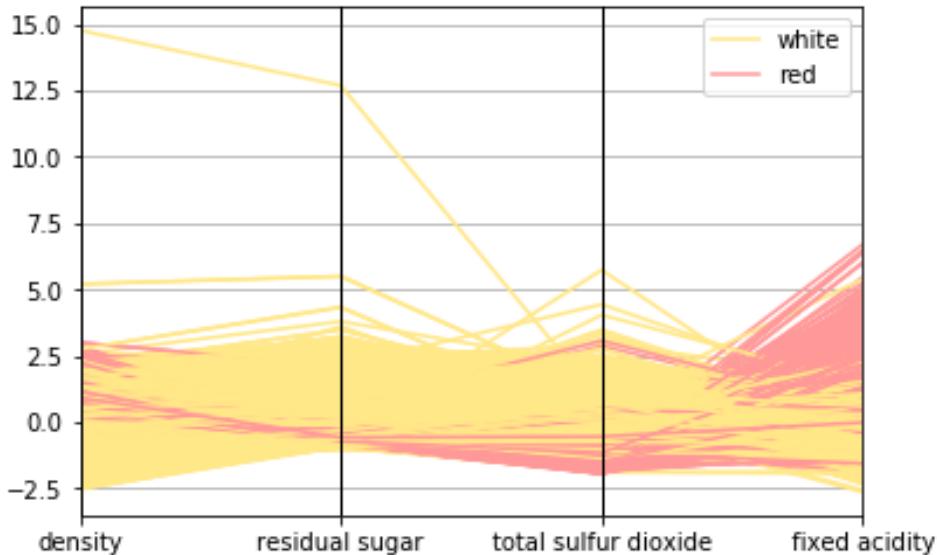
- A color-coded representation of 2D data
- Can be raw data, 2D histogram or any other function of 2 variables
- A color map accompanies the heat map
- We will learn other metrics in future that may be visualized as a heat map



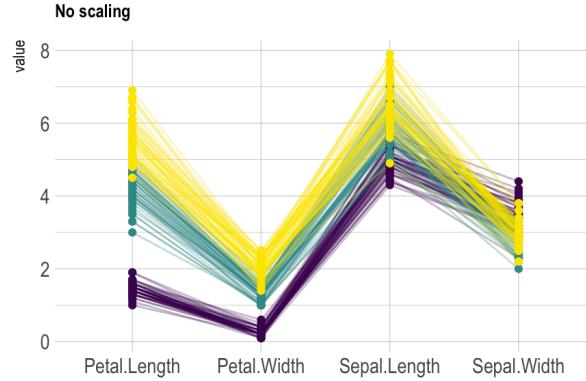
Plot Courtesy: Seaborn Documentation [seaborn.pydata.org]

Direct Visualization

- Parallel Co-ordinates
 - Each vertical line is a dimension
 - A data item is connected by line segments
 - Large number of samples clutters the visualization



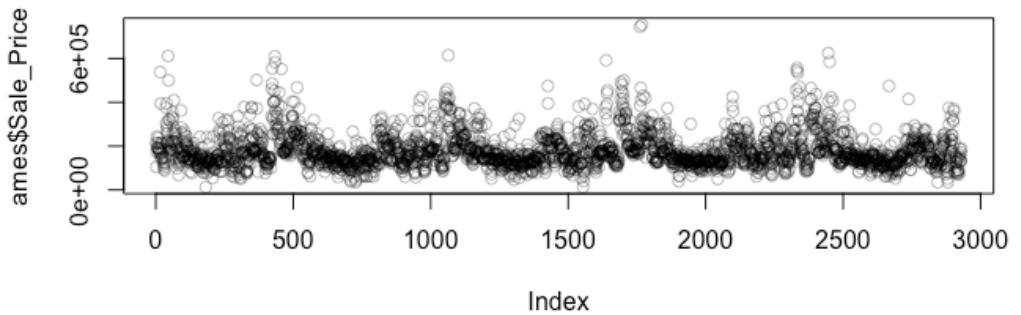
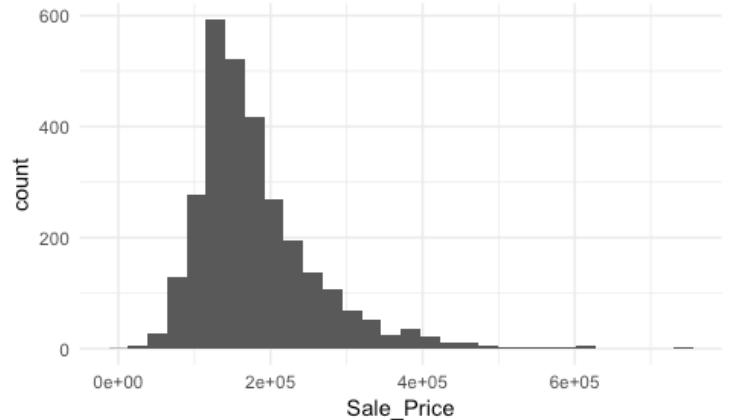
Parallel Coordinates



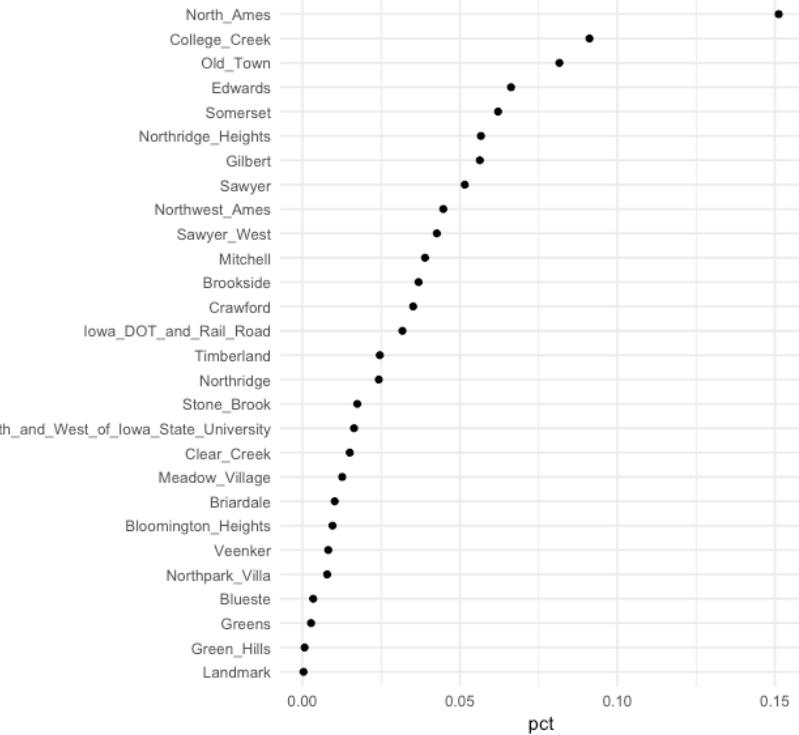
Leaderboard

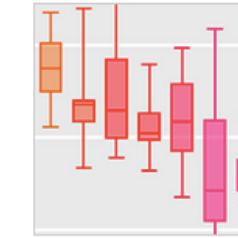
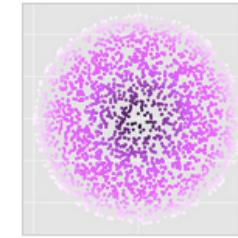
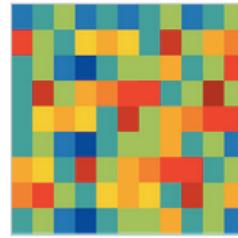
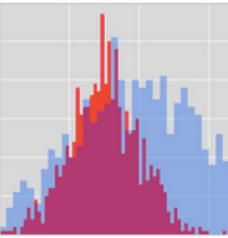
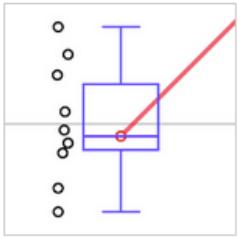
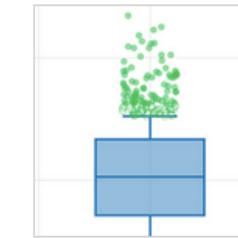
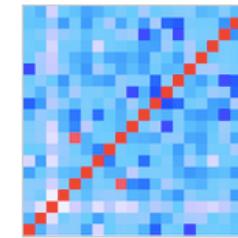
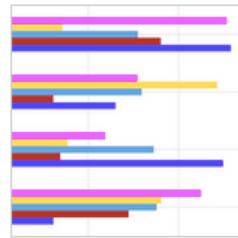
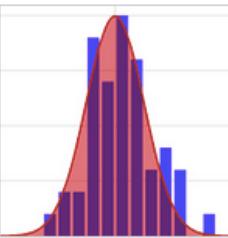
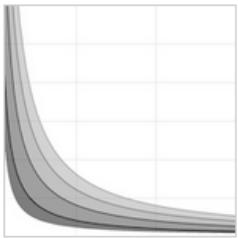
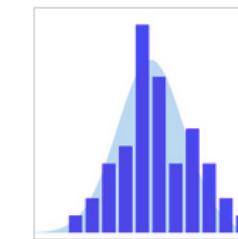
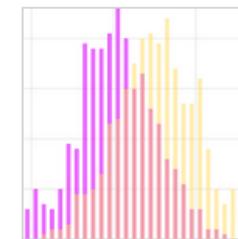
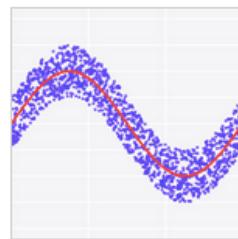
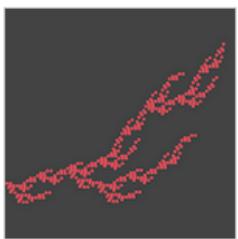
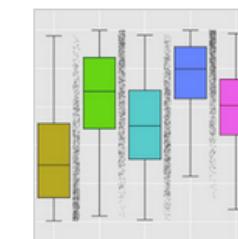
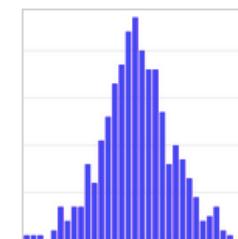
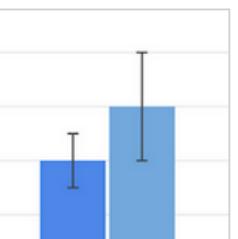
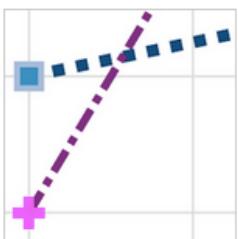
Proposals	Proposal Dollars	Awards	Award Dollars
1st Computer Science 10	Recreation & Tourism \$9.18M	1st Biology 20	1st Computer Scienc \$8.92M
2nd Biology 9	Health Sciences \$8.95M	2nd Computer Science 17	2nd Biology \$8.84M
2nd Electrical & Comp Engr 9	3rd Computer Scien \$8.92M	Health Sciences 12	3rd Electrical & Com \$4.56M
4th Chemistry and Biochem: 8	4th Biology \$8.84M	3rd Electrical & Comp Eng 12	4th Chemistry and B \$3.67M
Journalism 6	Psychology \$6.53M	5th Chemistry and Biochem 10	Health Sciences \$3.55M
Mathematics 4	Undergraduate Studi \$5.18M	Secondary Education 7	Recreation & Tourism \$3.06M
Deaf Studies 4	7th Chemistry and B \$3.67M	Elementary Education 7	Secondary Education \$2.85M
Elementary Education 4	8th Electrical & Com \$3.65M	Journalism 5	Elementary Education \$2.67M
Art 4	Secondary Education \$2.85M	Recreation & Tourism 5	Communication Stud \$2.39M
Health Sciences 3	Elementary Education \$2.67M	Mechanical Engineering 5	Psychology \$2.17M

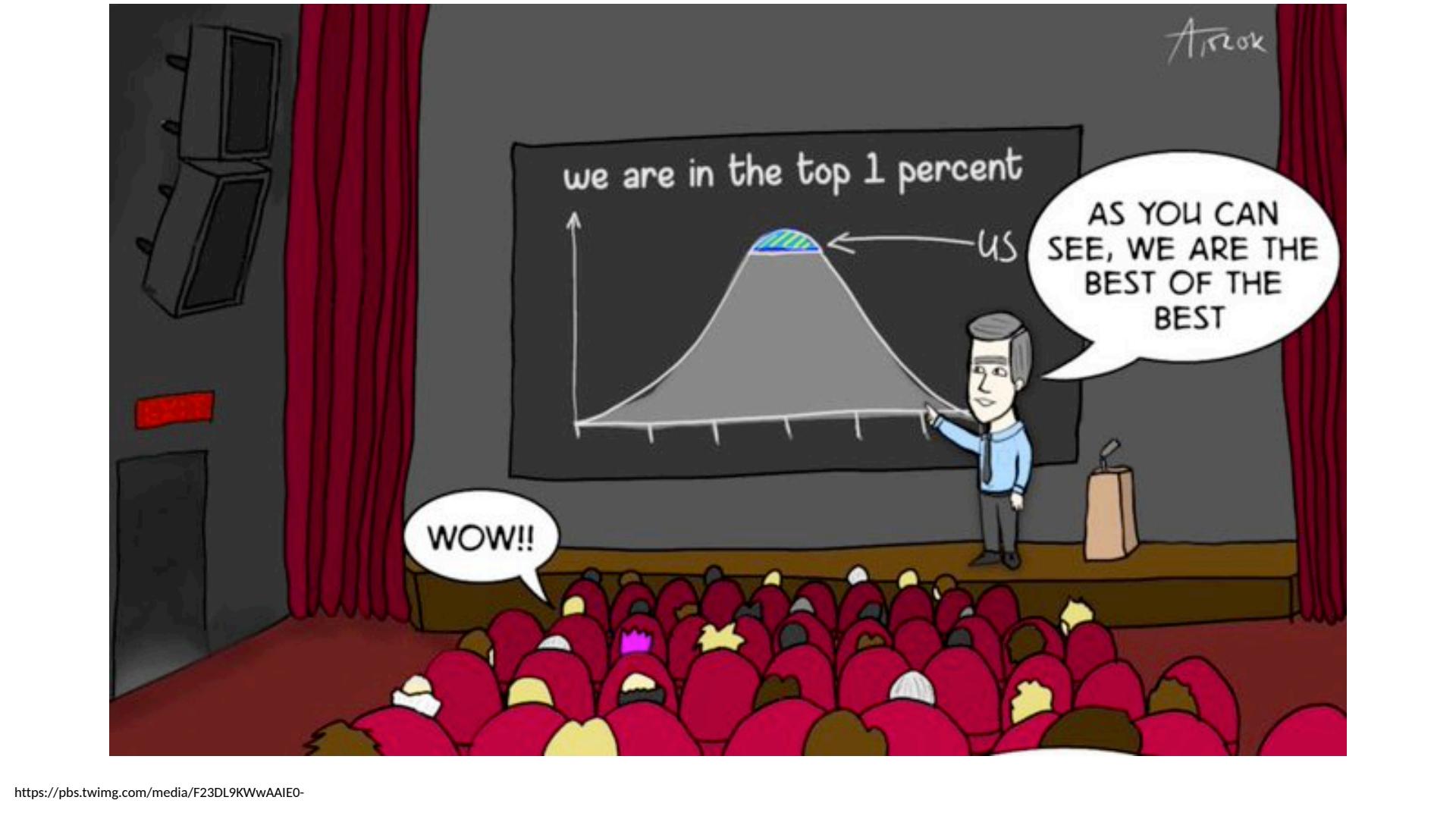
Top 10



reorder(Neighborhood, pct)





A political cartoon by Arrok. A man in a blue shirt and tie stands on a stage, pointing at a large screen displaying a bell curve. The top right of the curve is shaded green and labeled 'US'. Above the graph, the text reads 'we are in the top 1 percent'. A speech bubble from the man says 'AS YOU CAN SEE, WE ARE THE BEST OF THE BEST'. In the audience, a person says 'WOW!!' in a speech bubble.

Arrok

we are in the top 1 percent

US

AS YOU CAN
SEE, WE ARE THE
BEST OF THE
BEST

WOW!!

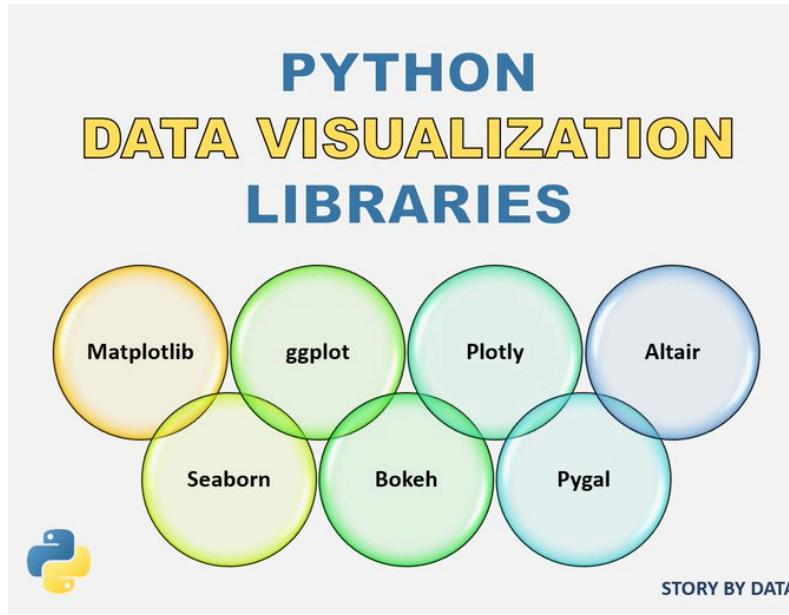
“In good information visualization, there are no rules, no guidelines, no templates, no standard technologies, no stylebooks ... You must simply do whatever it takes.**”**

—Edward Tufte



Resources

- <https://towardsdatascience.com/5-quick-and-easy-data-visualizations-in-python-with-code-a2284bae952f>

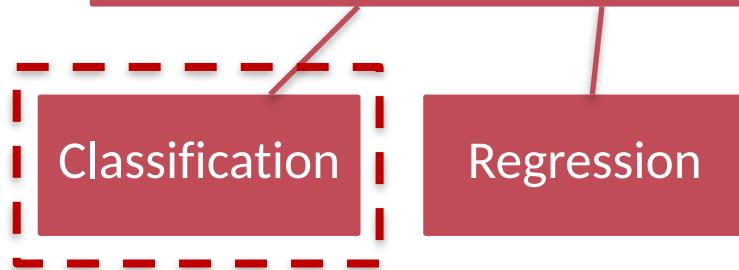


[https://twitter.com/
storybydata/status/
1166337648341991424](https://twitter.com/storybydata/status/1166337648341991424)

Supervised Learning



Supervised Learning





Classification

Binary

$\{0,1\}$

Multi-class

1-of-K

Multi-label

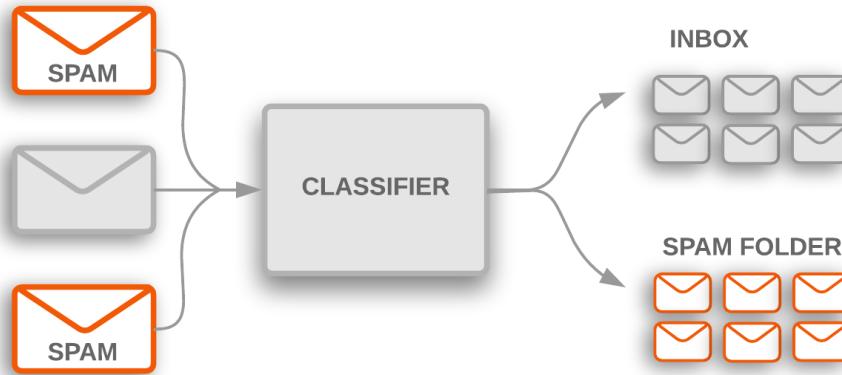
n-of-K

Structure



E.g. graph/sequence

Binary Classification



Performance Measures - Accuracy

$$Accuracy = \frac{(100 + 50)}{165} = 0.91$$

$$Misclassification = \frac{(10 + 5)}{165} = 0.09$$

n=165	Predicted:	
	NO	YES
Actual: NO	50	10
Actual: YES	5	100

- Pool of 100 patients' data used for validation of a cancer prediction ML model
- Prediction:
 - 3 have cancer
 - Rest ($100-3=97$) are healthy.
- Reality:
 - 1 of the 3 did not actually have cancer !
 - 3 from 97 predicted healthy actually have cancer
- Accuracy =

n=	Predicted: NO	Predicted: YES
Actual: NO		
Actual: YES		

- Pool of 100 patients' data used for validation of a cancer prediction ML model
- Prediction:
 - 3 have cancer
 - Rest ($100-3=97$) are healthy.
- Reality:
 - 1 of the 3 did not actually have cancer !
 - 3 from 97 predicted healthy actually have cancer
- Accuracy = $(100 - 4) / 100 = 96\%$!

n=	Predicted: NO	Predicted: YES
Actual: NO		
Actual: YES		

- Pool of 100 patients' data used for validation of a cancer prediction ML model
- Prediction:
 - 3 have cancer → selected for chemotherapy
 - Rest ($100 - 3 = 97$) are healthy.
- Reality:
 - 1 of the 3 did not actually have cancer !
 - 3 from 97 predicted healthy actually have cancer → should have been selected for chemotherapy
- Accuracy = $(100 - 4) / 100 = 96\%$!

n=	Predicted: NO	Predicted: YES
Actual: NO		
Actual: YES		

Performance Measures - Accuracy

$$Accuracy = \frac{(100 + 50)}{165} = 0.91$$

$$Misclassification = \frac{(10 + 5)}{165} = 0.09$$

n=165	Predicted:	
	NO	YES
Actual: NO	50	10
Actual: YES	5	100

Performance Measures - Accuracy, TPR, FPR

$$Accuracy = \frac{(100 + 50)}{165} = 0.91$$

$$Misclassification = \frac{(10 + 5)}{165} = 0.09$$

$$FalsePositiveRate(FP) = \frac{(10)}{60} = 0.17$$

$$FalseNegativeRate(FN) = \frac{(5)}{105} = 0.048$$

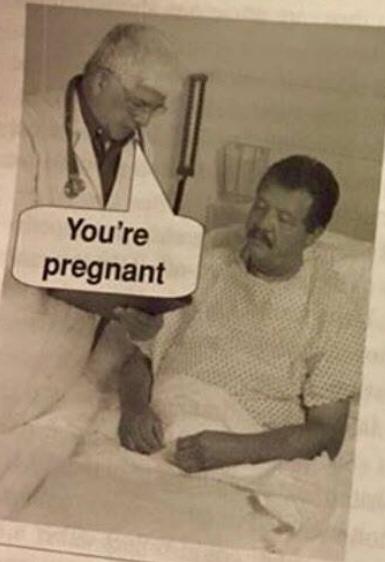
n=165	Predicted:	
	NO	YES
Actual: NO	TN = 50	FP = 10
Actual: YES	FN = 5	TP = 100
	55	110

$$TrueNegativeRate(TN) = \frac{(50)}{60} = 0.833$$

$$TruePositiveRate(TP) = \frac{(100)}{105} = 0.95$$

n=165	Predicted: NO	Predicted: YES	
Actual: NO	TN = 50	FP = 10	60
Actual: YES	FN = 5	TP = 100	105
	55	110	

Type I error
(false positive)



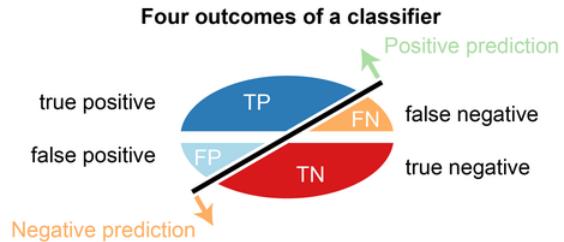
Type II error
(false negative)



Figure 3.1 Type I and Type II errors

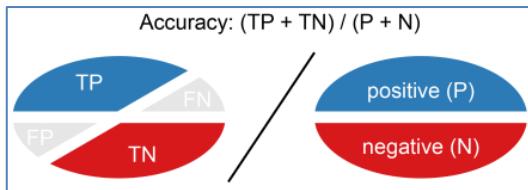
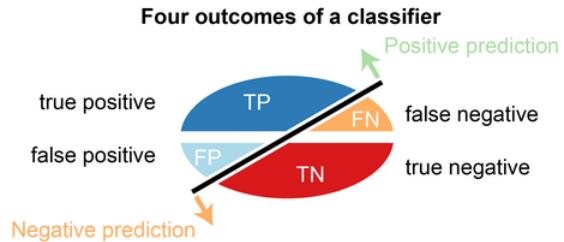
levels to .01 or even .001

Summary of Measures



n=165	Predicted: NO	Predicted: YES	
Actual: NO	TN = 50	FP = 10	60
Actual: YES	FN = 5	TP = 100	105
	55	110	

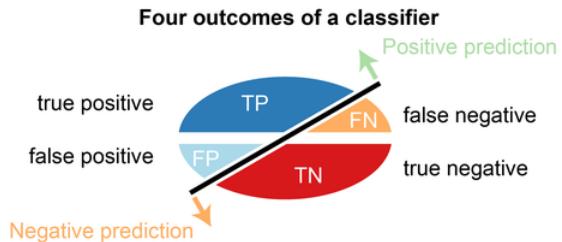
Summary of Measures



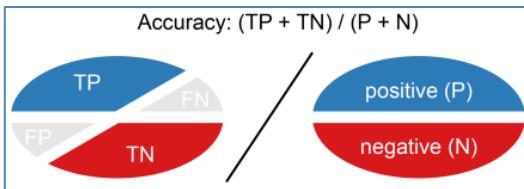
% of correct predictions

n=165	Predicted: NO	Predicted: YES	
Actual: NO	TN = 50	FP = 10	60
Actual: YES	FN = 5	TP = 100	105
	55	110	

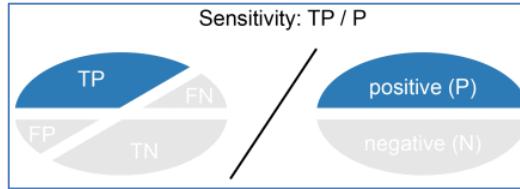
Summary of Measures



n=165	Predicted: NO	Predicted: YES	
Actual: NO	TN = 50	FP = 10	60
Actual: YES	FN = 5	TP = 100	105
	55	110	

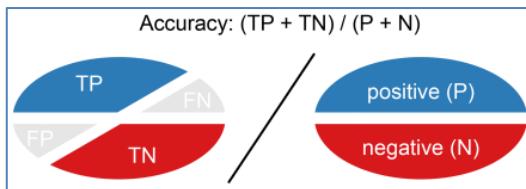
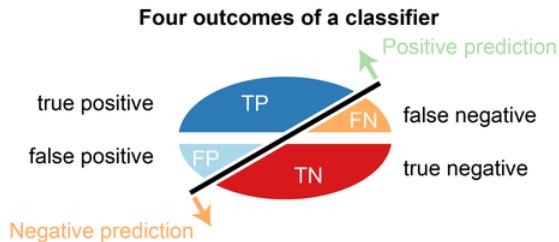


% of correct predictions

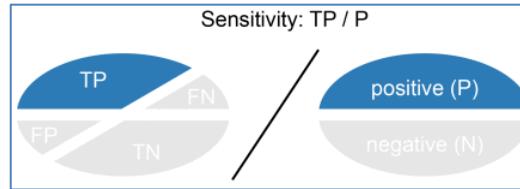


% of + class correctly predicted
[aka Recall / TPR]

Summary of Measures

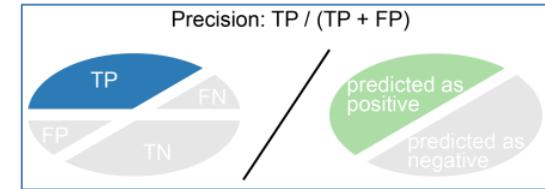


% of correct predictions



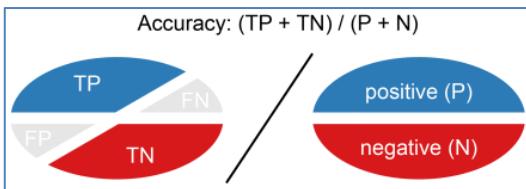
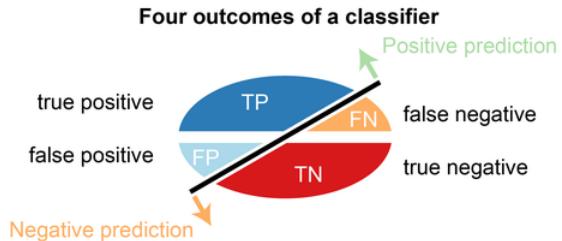
% of + class correctly predicted
[aka Recall / TPR]

n=165	Predicted: NO	Predicted: YES	
Actual: NO	TN = 50	FP = 10	60
Actual: YES	FN = 5	TP = 100	105
	55	110	

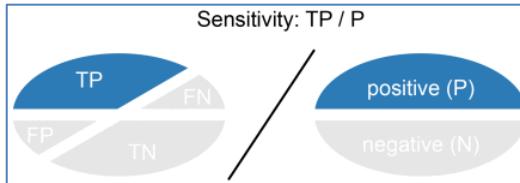


correct prediction of + class
[aka Precision]

Summary of Measures

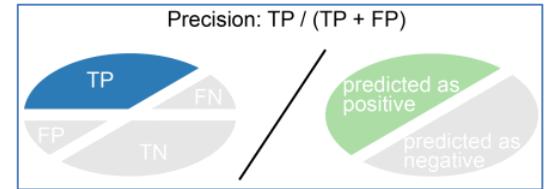


% of correct predictions



% of + class correctly predicted
[aka Recall / TPR]

n=165	Predicted: NO	Predicted: YES	
Actual: NO	TN = 50	FP = 10	60
Actual: YES	FN = 5	TP = 100	105
	55	110	



correct prediction of + class
[aka Precision]



% of - class incorrectly predicted

- **Cancer-Prediction System**
- Precision =
- Recall =
- Accuracy =

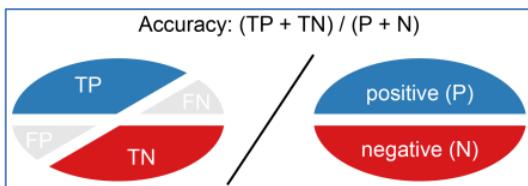
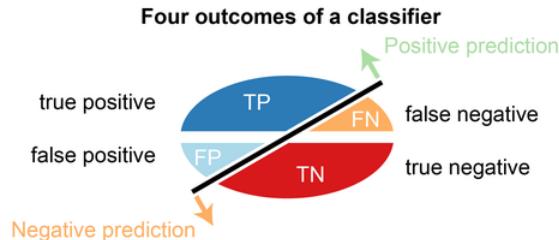
- **Cancer-Prediction System**
- Precision = $2/(2+1) = 67\%$
- Recall = $2/(2+3) = 40\%$
- Accuracy = $(94+2)/100 = 96\%$

What measure are we optimizing for ?

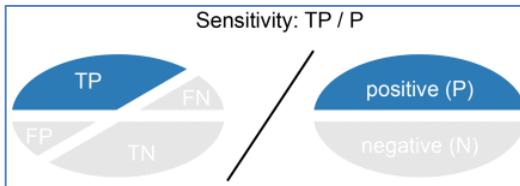
- Screening for a terminal disease
 - Do not want to miss anyone: Maximize Recall
- Classification between apples and oranges
 - Both types of errors are equally imp.: Maximize Accuracy
- Automatic bombing on detecting a target from a drone
 - Should not hurt civilians: Zero False Positives
- Giving access to a secure installation
 - No access to unauthorized personnel: Low False Positives

Accuracy vs Precision vs Recall

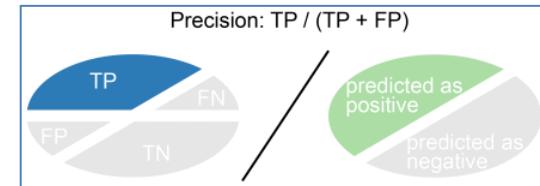
- Accuracy : Performance w.r.t both classes
- Recall : Performance w.r.t '+' class
- Precision : Reliability of predictions w.r.t '+' class



% of correct predictions



% of + class correctly predicted
[aka Recall / TPR]



correct prediction of + class
[aka Precision]

Utility and Cost

- Sometimes, there is a cost for each error
 - E.g. Earthquake prediction
 - False positive: Cost of preventive measures
 - False negative: Cost of recovery
- Detection Cost (Event detection)
 - $\text{Cost} = C_{FP} * FP + C_{FN} * FN$

Farmer Shri MoneyBags and ML-FruitPicker

- MB : I want an automated fruit picker and packer. I will pay an unholy amount for it.
- You (having just finished this lecture) : Sure
- You (*Thinking*): *I love unholy amounts of money* 😎

Farmer Shri MoneyBags and ML-FruitPicker

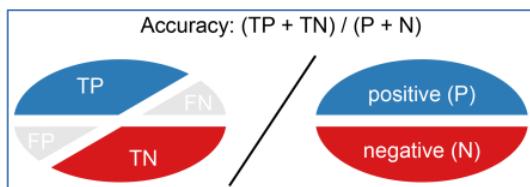
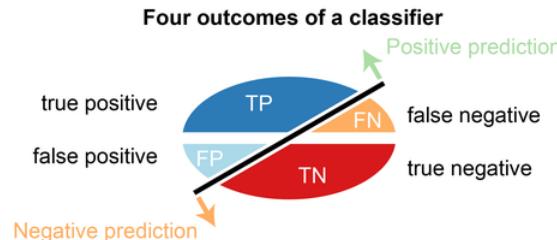
After 6 months ...

- MB : Well ?
- You : I have a High Precision ML-FruitPicker. But its Recall is 20% !
- MB : (confused) Precision ? Recall ?
- You : (*thinking*) Should I go over first 3 lectures of SMAI with MB ? He'll probably run away !
- You : It rejects 80% of good, pickable fruit, but whatever it picks, those fruits are good !
- MB : I'll take your system. How do I transfer unholly amount of money to you ?
- You : 
- MB (seeing your shocked face) : See, in a batch of 100 fruits, 10 fruits are usually bad. Among the 90 good ones, your system will select 18 of them on average. But from any given selection, I pack only 8.

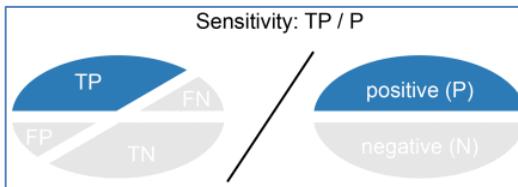


Accuracy vs Precision vs Recall

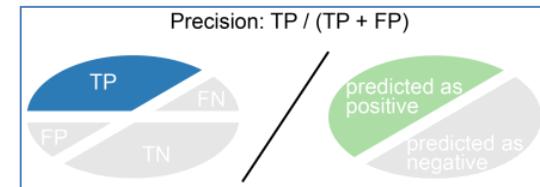
- Monitor **Precision** if a **false positive** carries higher cost.
- Monitor **Recall** if a **false negative** carries higher cost.



% of correct predictions



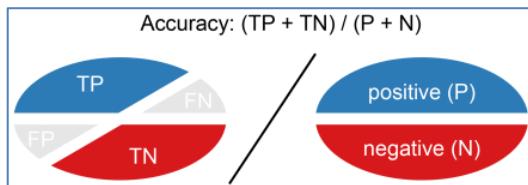
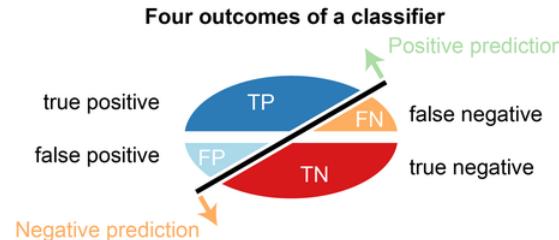
% of + class correctly predicted
[aka Recall / TPR]



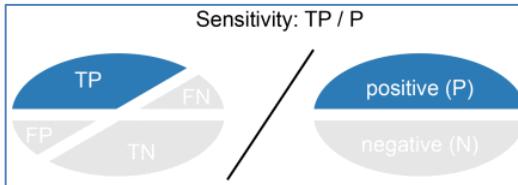
correct prediction of + class
[aka Precision]

Accuracy vs Precision vs Recall

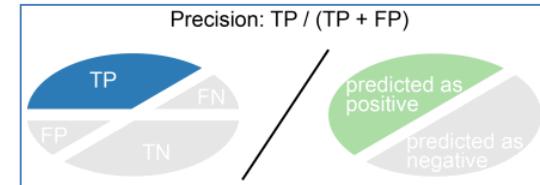
- Precision → Cost of inclusion
- Recall → Cost of exclusion



% of correct predictions

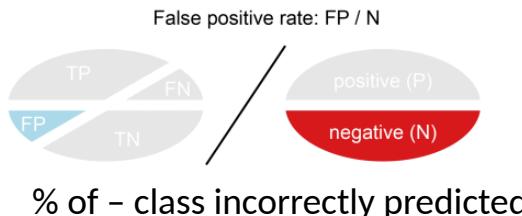
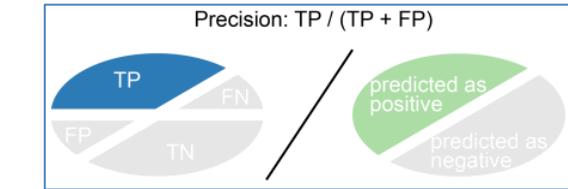
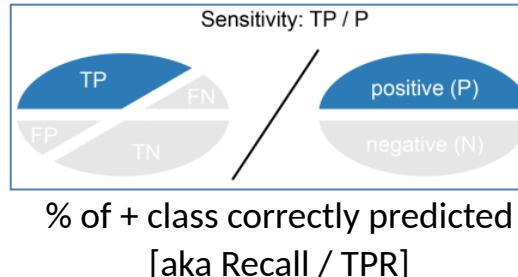
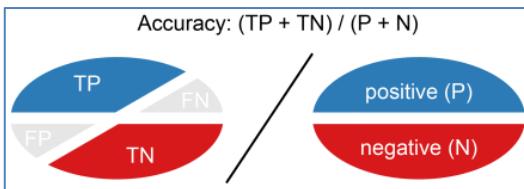
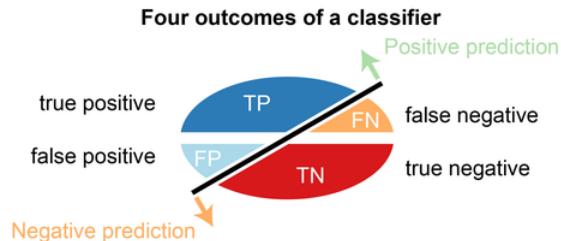


% of + class correctly predicted
[aka Recall / TPR]



correct prediction of + class
[aka Precision]

Summary of Measures



F1-score: A unified measure

- What to do when one classifier has better precision but worse Recall, while other classifier behaves exactly opposite?
 - F-measure (Information Retrieval)

$$\blacksquare \quad F_1 = \frac{2}{\frac{1}{Recall} + \frac{1}{Precision}}$$

Utility and Cost

- What to do when one classifier has better Precision but worse Recall, while other classifier behaves exactly opposite?
 - F-measure (Information Retrieval)
 - $$F_1 = \frac{2}{\frac{1}{Recall} + \frac{1}{Precision}}$$

→ F1 measure punishes extreme values more !

→ Definition of Recall and Precision have same numerator, different denominators. A sensible way to combine them is harmonic mean.