Acquiring and Understanding Real Data

https://tinyurl.com/cis5190-10-26-2022

Osbert Bastani and Zachary G. Ives CIS 4190/5190 – Fall 2022



- Homework 4 due Wednesday November 2, 8pm
- Quiz 7 due tomorrow, Oct 27, 8pm

Our Machine Learning Toolkit in Practice

So far, we've assumed the data is ready for us to apply ML techniques

- **1.** Available as an "X matrix": instances as rows, features as columns
- **2.** We know the classes we want to predict
- **3.** We are given training labels for the classes

But this isn't always the case!

- Data may need to be wrangled and integrated
- We may need to understand our data and tasks
- We may need to understand what the data "tells us"

For more depth, see CIS 5450 – but here we'll briefly discuss some of the major techniques and ideas



https://www.base-4.com/open-apartments-faster-with-modular



https://dreamsmeaning.site/house-under-construction/



Need to Do Work to Prepare Data for ML

Data is rarely "clean": Real data is messy, fragmented

80% + of human time is spent on data wrangling, according to practitioners

- Depends on who you ask, the application, the data source, etc...
- And "80%+" can be an underestimate!

Having good data with appropriate features is absolutely critical to success

From Raw to Usable Data



May need to iterate many times, and clean at each step!

Challenges of Acquiring Relevant Data

In many data science/ML applications: identifying and *getting access to data sources* is a challenge itself

- "Data fiefdoms"
- Lack of clarity on what is needed / need for data discovery
- Proprietary data
- Privacy constraints on access: HIPAA, FERPA, GDPR, ...



Once you have permissions, may need to *wrangle* the data into Pandas tables

- From popular data sharing formats: CSV, JSON, XML
- From HTML tables
- From structured files: DICOM, HDF5, Excel, MatLab, ...
- Text \rightarrow information extraction and NLP
- Database connections: SQLAlchemy etc

https://www.drivethrurpg.com/product/190343/Fiefdom-Manors-Ruling-manorial-fiefs

Merging Data: Pandas or SQL merge / join

Data may be split across different files

trooko

Requires doing a join based on a key to combine data into one table

liachs							albuills					ai 11515						
• • •	● Insert Pag	• ట్రి ⊽ e Layout Formulas	: Data Review	🗈 Track v View			Q - Search Sheet	© ▼ ≗ + Share ~	● ● ■ Home	Insert Pa	२ ऍ रू ोो Album Q×Searc ge Layout Formulas Data Re	h Sheet view View 💄	+ Share ↓	Hom	● Insert Pag	> ₽ Q y Search ge Layout Formulas	Sheet Data >> よ	© ▼ Share ▼
	A	В	С	DE	F	G	Н			A	В	С	D		A	В	С	D
1 i	d	name	album_id	media_type_id genre_id	composer	milliseconds b	oytes	unit_price	ı id	d	title	artist_id		-	id	name		
2	1	For Those Ab	1	1	1 Angus Young	343719	11170334	0.99	2		1 For Those About To Ro	-	L	2	:	1 AC/DC		
3	2	Balls to the V	2	2 2	1	342562	5510424	0.99	3	:	2 Balls to the Wall	2	2	3	:	2 Accept		
4	3	Fast As a Sha	3	3 2	1 F. Baltes, S. K	230619	3990994	0.99	4	:	3 Restless and Wild	2	2	4	3	3 Aerosmith		
5	4	Restless and	3	3 2	1 F. Baltes, R.A	252051	4331779	0.99	5		4 Let There Be Rock	-	L	5	4	4 Alanis Moriss	ette	
6	5	Princess of th	3	3 2	1 Deaffy & R.A	375418	6290521	0.99	6		5 Big Ones	3	3	6	!	5 Alice In Chair	S	
7	6	Put The Finge	1	. 1	1 Angus Young	205662	6713451	0.99	7		6 Jagged Little Pill	4	1	7	-	7 Apocalyptica		
8	7	Let's Get It U	1	. 1	1 Angus Young	233926	7636561	0.99	8	•	7 Facelift		5	8	į	8 Audioslave		
9	8	Inject The Ve	1	1	1 Angus Young	210834	6852860	0.99	9	!	9 Plays Metallica By Four	-	7	9	(BackBeat		
10	9	Snowballed	1	. 1	1 Angus Young	203102	6599424	0.99	10	1	O Audioslave	5	3	10	10	Billy Cobham		
11	10	Evil Walks	1	1	1 Angus Young	263497	8611245	0.99	11	1	1 Out Of Exile	8	3	11	1	1 Black Label S	ociety	
12	11	C.O.D.	1	1	1 Angus Young	199836	6566314	0.99	12	1	2 BackBeat Soundtrack	ç	9	12	12	2 Black Sabbat	ı	
13	12	Breaking The	1	1	1 Angus Young	263288	8596840	0.99	13	1	3 The Best Of Billy Cobha	10	<mark>)</mark>	13	13	Body Count		
14	13	Night Of The	1	1	1 Angus Young	205688	6706347	0.99	14	14	4 Alcohol Fueled Brewta	11	L		Artist +	Durree Dieldine		
15	14	Spellbound	1	1	1 Angus Young	270863	8817038	0.99	1	Album +	Alashal Fusiad Drouts	1 4		Read	y III	─ -	+	200%
Read	Track +							+ 200%	Ready			+	200%					

albuma

ortioto

Integration May be Hard

Merged table may be too large for memory

- Incrementally load and join data, using SGD or mini-batches
- Use online learning techniques

Encoding issues

- Inconsistent data formats or terminology
- Key aspects mentioned in cell comments or auxiliary files

Record linking problem (next)

The Record Linking Problem

Ins ID	Name
203342	J Smith
123452	Mao Y



Student ID	Name
3432432	Jon Smithee
9734783	Jane Smyth
8273737	Ying Mao

Huge literature. Some popular ideas:

- String similarity above a threshold
 - String edit distance ("J Smith" → "Jon Smithee" with 4 edits)
 - String overlap (n-grams)
- May want to tokenize and compare tokens, not just strings
 - Or consider "soundex", common mis-substitutions, ...
- Often combine similarities of multiple fields (e.g., addresses, employer)

Once Data Is Integrated: Need to Deal with Non-Numeric Feature Types



Encoding Features

Encode *categorical* features

• Use **one-hot encoding:** Expand $X_i \in \{1,2,3\}$ into [1,0,0] or [0,1,0] or [0,0,1]

Encode *ordinal* features

- Convert to a number, preserving the order (e.g. [low, medium, high] \rightarrow [1, 2, 3])
- Encoding may not capture relative differences, so may still want one-hot encoding

_					•
HouseStyle	FullBath	RoofMatl	BsmtCond	KitchenQual	
1Story	2	CompShg	ТА	TA	
SLvl	1	CompShg	ТА	TA	I
2Story	2	CompShg	TA	Gd	1
1Story	2	CompShg	Gd	Ex	1
2Story	2	CompShg	TA	Gd	
SLvl	1	WdShngl	TA	TA	
2Story	2	CompShg	TA	Gd	
SLvl	1	CompShg	TA	TA	l
2Story	2	CompShg	TA	TA	
2Story	2	CompShg	ТА	Gd	

HouseStyle	FullBath	RoofMatl	BsmtCond	KitchenQual
1Story	2	CompShg	3	3
SLvl	1	CompShg	3	3
2Story	2	CompShg	3	4 🛛
1Story	2	CompShg	4	5
2Story	2	CompShg	3	4
SLvl	1	WdShngl	3	3
2Story	2	CompShg	3	4
SLvl	1	CompShg	3	3
2Story	2	CompShg	3	3
2Story	2	CompShg	3	4∥

Data Quality Issues

Missing feature values

- **Delete features** with mostly missing values
- **Delete instances** with (many) missing features, if rare
- Impute via mean (numeric) or mode (categorical or ordinal)
- Learn to **predict the missing values** (i.e., a kind of model stacking)
- Flag missing values using **binary variables**
 - Data might not be "missing at random"
 - It might be meaningful that instances have missing features!
 e.g., history might be missing from an unconscious trauma patient

Missing Values

ID	Last_Visit
1234	2018-03-05
4567	0
8910	2019-12-12

Rather than removing the case where Last_Visit is unknown – flag it with a separate feature!

We will learn a weight for the feature if it's there, and also a different weight if it isn't!

Data Quality Issues

Incorrect feature values

- Typos: e.g., color = {"bleu", "green", "gren", "red"}
- Inconsistent spelling (e.g., "color", "colour")
- Inconsistent abbreviations (e.g., "Oak St.", "Oak Street")
- Garbage: e.g., color = "w_l r--śij"
- Often: compare against a dictionary (~ spell-check)

Missing instance labels

- Delete instances if only a few are missing labels
- (Can also use "semi-supervised learning" techniques that can exploit unlabeled data samples alongside labeled data samples)

Example: Amazon Product Knowledge Graph



2021 Apple 12.9-inch iPad Pro (Wi-Fi, 256GB) - Space Gray

 \star

Add to Cart

\$99900 FREE delivery Sunday, October 30

New

Adorama

Adorama

★★★★★ (575279 ratings) 90% positive over last 12 months

See more

Ships from

Sold by

5 other options sorted by price + deliver	Filter 🗸	
Used - Like New \$1,078 ⁵⁵	FREE delivery October 27 - November 1. Details	Add to Cart
Ships from Sold by	Expercom - Apple Premier Partner Expercom - Apple Premier Partner ************************************	
Used - Like New ^{\$} 1,079 ⁰⁰	FREE delivery October 27 - 28 . Details	Add to Cart
Condition	Open Box Item # 1586413 What's in the box: Apple iPa Chip (Space Gray), 20W USB Type-C PoMore	d Pro 12.9" M1

Amazon wants to rank + recommend products by relevance – to search or on the front page

Q: How does it know whether a product in a vendor catalog is a 2021 iPad Pro vs a 2019 version? The right size?

A: Aligning with a *product knowledge* graph!

Amazon Product Knowledge Graph



AWS knowledge graph tools [Dong et al KDD 18]:

- Distant supervision to take knowledge graph, use NNs to learn patterns from open Web page product descriptions apply to model catalog entries
- Random forest for record linking based on terms, values

Script Your Data Preprocessing!

Don't manually edit via a spreadsheet program

- No history of changes
- Very easy to introduce mistakes
- Hard to correct earlier decisions



Instead, write a script that loads the raw data and does all preprocessing

- Documents all steps
- Incremental debugging
- Easy to make changes to earlier steps
- Repeatable



Summary: Data Integration

Key tasks as a data science consultant or a data scientist:

- Understand the data, important aspects of the domain
- Identify relevant data sources
- Import and wrangle the data
- Integrate and clean

Integration across sources may happen before or after cleaning – *why?*

- Record linking
- Value imputation

Data processing should often be done in a pipeline or script

Understanding Your Data

What's Interesting? Important? Representative?

Machine learning is algorithmic – but its use should be guided by an understanding of the data, underlying hypothesis space, etc. e.g., two features may be heavily correlated, or there may be odd outliers, or unexpected scaling, or ...

Thus: we should always get a sense of the data before we start trying to build models!

Simplest Starting Point: Dataframe.describe()

			r	Note the nissing va	e lues	No missing target values					
	Id	MSSubClass	LotFrontage	LotArea	OverallQual	OverallCond	YearBuilt	YearRemodAdd	MasVnrArea		SalePrice
count	1022.000000	1022.000000	832.000000	1022.000000	1022.000000	1022.000000	1022.000000	1022.000000	1019.000000		1022.000000
mean	732.338552	57.059687	70.375000	10745.437378	6.128180	5.564579	1970.995108	1984.757339	105.261040		181312.692759
std	425.860402	42.669715	25.533607	11329.753423	1.371391	1.110557	30.748816	20.747109	172.707705		77617.461005
min	1.000000	20.000000	21.000000	1300.000000	1.000000	1.000000	1872.000000	1950.000000	0.000000	•••	34900.000000
25%	367.500000	20.000000	59.000000	7564.250000	5.000000	5.000000	1953.000000	1966.000000	0.000000		130000.000000
50%	735.500000	50.000000	70.000000	9600.000000	6.000000	5.000000	1972.000000	1994.000000	0.000000		165000.000000
75%	1100.500000	70.000000	80.000000	11692.500000	7.000000	6.000000	2001.000000	2004.000000	170.000000		215000.000000
max	1460.000000	190.000000	313.000000	215245.000000	10.000000	9.000000	2010.000000	2010.000000	1378.000000		745000.000000

Potential outliers

Feature Correlation Matrix



1.0

- 0.8

0.6

0.4

- 0.2

- 0.0

-0.2

-0.4

Plot Features Most Correlated with Target Variable (Seaborn pairplot and similar)



25

Handling Outliers: Causes of Outliers

Errors

- Human error in data collection or data entry
- Measurement/instrumentation errors
- Experimental errors
- Data merge errors
 - e.g., merging datasets with different scales
- Data preprocessing errors

Natural

• Novelties in the data – not mistakes!

Outlier Detection: Z-score

- Assume feature values are **Gaussian-distributed**
- Discard points more than k standard deviations away from the mean
 Good values for k: 2.5, 3, 3.5+

Cautions:

- Mostly for low-*d* feature spaces on reasonably small-to-medium data sets
- Incorrect if parametric assumption doesn't hold



Feature Standardization

Rescales features to have zero mean and unit variance

• Let
$$\mu_j$$
, σ_j be: $\mu_j = \frac{1}{N} \sum_{i=1}^N x_{ij}$, $\sigma_j^2 = \sum_{i=1}^N (x_{ij} - \mu_j)^2 / N$

• Replace each value with
$$x_{ij} \leftarrow \frac{x_{ij} - \mu_j}{\sigma_j}$$
 for $j = 1 \dots D$ (not x_0)
Same as the "Z score"

Could also rescale features to lie in [0, 1] (for special cases where there's a clear scale, e.g., pixel values)

• Replace each value with
$$x_{ij} \leftarrow \frac{x_{ij} - \min(x_j)}{\max(x_j) - \min(x_j)}$$

Must use the <u>same transformation</u> for both training and prediction $(\mu_j \text{ and } \sigma_j \text{ are computed on training data and also used on the test data)$



- Important to look at data, identify missing values, correlations, and outliers
- Remove outliers
- Impute or drop missing values
- Use regularization or drop correlated features

Beyond describe(): Understanding Structure in Data Domains

- We can visualize data, and call describe(), to get a sense of univariate and bivariate distributions
- But: suppose the data naturally falls into different groupings perhaps even suggesting which *classes* we would like to learn?

This motivates a study of additional techniques for extracting structure present in the data

Unsupervised Machine Learning

Types of Learning

Supervised learning

 Given: training data + desired outputs (labels)

Unsupervised learning

• Given: training data (without labels)

Semi-supervised learning

• Given: training data + some labels

Reinforcement learning

• Given: observations and occasional rewards as the agent performs sequential actions in an environment

From Supervised to Unsupervised Learning

Supervised Learning: Classification



Supervised Learning: Regression

- Given $(x_1, y_1), (x_2, y_2), ..., (x_n, y_n)$
- Learn a function f(x) to predict y given x







Given $x_1, x_2, ..., x_N$ (without labels) Learn "hidden structure" in the data



Types and Uses of Unsupervised Learning

Dimensionality Reduction

Map samples $x_i \in \mathbb{R}^D$ to $f(x_i) \in \mathbb{R}^{D' \ll D}$

Special case: clustering, $f(x_i) \in \mathbb{N}$

- **Feature Learning:** For preprocessing inputs to an ML algorithm, since lowerdimensional features permit smaller models and fewer data samples.
- **Compression (for storage):** e.g. JPEG standard for images is now adopting unsupervised ML approaches https://jpeg.org/items/20190327_press.html
- Visualization: Exploring a dataset, or an ML model's outputs



"Here's some data, could you do something with it?"

Unsupervised learning helps identify uses for data:

- Visualize the data, find clusters

 e.g. "based on our polling data, there are three main voting blocs, based
 on age, race, education level, income, political beliefs, and home ownership. Features like marital status and # children are irrelevant."
- Identify interesting supervised learning problems within your dataset
 e.g. "do our company's profits y_i actually vary systematically based on
 the weather x_i?"
- Generate new data

e.g. "given all of Bach's work, I could generate new music that would sound like Bach."

Popular Clustering Algorithms

- **1.** One of the simplest, scalable techniques: k-Means clustering
- 2. Greedy: hierarchical clustering

The Clustering Setting

Task:

Input: $\mathcal{D} = \{x_i\}_{i=1}^N$

Want to discover a mapping $f(x_i) \in \{1, 2, 3, ..., K\}$ that discovers natural groupings in the data.

Performance Metric / Objective Function: What is a good mapping f(.)?

Somewhat loosely defined, and different clustering algorithms differ in their definition of a good clustering f(.)



© 2019-22 D. Jayaraman, O. Bastani, Z. Ives Figures: Dan Roth



Guess The Cluster Assignment



Cluster Assignment in K-Means

- Define a centroid μ_k for each cluster k
- For any new sample, assign the cluster whose centroid/mean is closest!

Equivalent to 1-nearest neighbor classification over the cluster means.

Note: Certainly not the only answer we could have come up with!



Can we expand this into a full, consistent clustering algorithm?

Clustering Data



Clustering Data

K-Means (K, X)

- Randomly choose *K* cluster means
- Loop until convergence, do:
 - Assign each point to the cluster of the closest centroid
 - Re-estimate the cluster centroids based on the data assigned to each cluster

K-Means (K, X)

- Randomly choose *K* cluster center locations (centroids)
- Loop until convergence, do:
 - Assign each point to the cluster of the closest centroid
 - Re-estimate the cluster centroids based on the data assigned to each cluster

K-Means (K, X)

- Randomly choose *K* cluster center locations (centroids)
- Loop until convergence, do:
 - Assign each point to the cluster of the closest centroid
 - Re-estimate the cluster centroids based on the data assigned to each cluster

K-Means (K, X)

- Randomly choose *K* cluster center locations (centroids)
- Loop until convergence, do:
 - Assign each point to the cluster of the closest centroid
 - Re-estimate the cluster centroids based on the data assigned to each cluster

Optimizer: "Alternating Minimization"

K-means finds a local optimum of the following objective function:

"Sum of squared distances" loss function

$$\arg\min_{\boldsymbol{S}}\sum_{k=1}^{K}\sum_{\boldsymbol{x}\in S_{k}}\|\boldsymbol{x}-\boldsymbol{\mu}_{K}\|_{2}^{2}$$

where $S = \{S_1, ..., S_K\}$ are sets corresponding to disjoint clusters, and the clusters together include all samples.

K-Means Clustering Convergence

https://dashee87.github.io/data%20science/general/Clustering-with-Scikit-with-GIFs/

Generalizing to Other Distance Functions

K-Means Objective Function:

$$\arg\min_{\boldsymbol{S}}\sum_{k=1}^{K}\sum_{\boldsymbol{x}\in S_{k}}\|\boldsymbol{x}-\boldsymbol{\mu}_{K}\|_{2}^{2}$$

But it is possible to define k-means with other notions of pairwise distance between samples too. For example:

$$\left(\sum_{d} |x_{1d} - x_{2d}|^{1}\right)^{\frac{1}{1}} Q: What would have to change in the algorithm?$$

$$\ell_{1} \text{ distance}$$

$$\sum_{d} |x_{1d} - x_{2d}|$$

$$\sum_{d} |x_{1d} - x_{2d}|$$

$$(2019-2x_{2d})^{2} (2019-2x_{2d})^{2} (2019-2x_{2d})^{2}$$

K-Means is Too Sensitive to Initialization

Alternative strategies:

- 1. Do many runs of K-Means, each with different initial centroids, and pick the best
- 2. Pick initial centroids using a better method than random choice

K-means+ + initialization

- Choose a data point uniformly at random as the first centroid
- Loop for 2: *K*, do:
 - Let D(x) be the distance from each point x to the closest centroid
 - Choose data point x randomly $\propto D(x)^2$ as the next centroid. Higher chance to pick points that are far from previous centroids.

K-means++ Illustrated

Place the initial centroids **far away from one another**:

- Initialize an empty set M (for storing selected centroids);
 Randomly select the first centroid from the input sample and assign it to M
- 2. For each x_i that is not in M, find the distance $D(x_i)$ to the closest centroid in M
- Choose one new data point at random as a new centroid using probability distribution ~ D(x)²
- 4. Repeat (2) and (3) until K centroids have been chosen

Then do "classic" k-means

How Many Clusters?

Measuring the Performance of k-Means

How can we evaluate how good our clustering is? Some options:

- Evaluation using the k-means objective itself
- Comparing to class labels (for a subset of data) Sometimes possible
- Subjective evaluation by a human domain expert

. . .

"Knee Point" For Selecting K

Elbow Method For Optimal k

https://blog.cambridgespark.com/how-to-determine-the-optimal-number-of-clusters-for-k-means-clustering-14f27070048f

- Non-deterministic (may get different output based on different start values) but guaranteed to converge
- Iterative algorithm with two sub-steps (after random cluster centroids chosen):
 - **1.** Assign points to nearest cluster
 - **2.** Recompute cluster centroid
- Select number of clusters by exploring error (distortion)

Are there other methods of clustering?

Hierarchical Clustering

Hierarchical Clustering

Instead of repeating until convergence based on global measures, like k-Means –

Let's consider a greedy *local* algorithm, that iteratively makes choices

- 1. Start with **single-item clusters**, then build successively bigger and bigger clusters: *agglomerative*
- 2. Or: start with one cluster, break into the most logical sub-clusters, repeat: *divisive*

These are called *hierarchical clustering* approaches...

Basic Intuition

- Agglomerative: In each iteration, we find the closest clusters and merge them
- Divisive: In each iteration, we find the two most distant sub-clusters and split there

But: we know how to compute differences in points – need to generalize this to computing distances among clusters

Potential Cluster Differences

Single Linkage: Compute distances between the **most similar** members for each pair of clusters

Merge the clusters with the smallest **min-distance**

Complete Linkage: Compute distance between the **most dissimilar** members for each pair of clusters

Merge the clusters with the smallest **max-distance**

How Do We Do this Efficiently? Considering Agglomerative Case

Really inefficient to iterate over each point and compute its distance to every other point

• O(n²) computations, which is bad for big data!

Idea: precompute and memorize:

- Compute a distance matrix where distance[i,j] is the distance between nodes i,j
- We'll update this matrix every time we merge

Pseudocode: Agglomerative Clustering with Complete Linkage Single

1. Compute distance matrix **dist** between all pairs of points (a,b)

2. Repeat:

Iterate over pairs of clusters A, B, compute their distance: Look at all pairwise distances dist[a,b] between $a \in A, b \in B$ Merge the pair of clusters with *min* distance between most distant members least Update the distance matrix

Until a single cluster remains

Example: Reconstructing Phylogenetic Trees

Q: Is a panda a bear?

Or is its closest relative the red panda, which is related to the raccoon?

https://www.nwf.org/Educational-Resources/Wildlife-Gu

https://www.smithsonianmag.com/science-nature/eight-amazing-facts-about-red-pandas-180979708/

https://towardsdatascience.com/hierarchical-clustering-and-its-applications-41c1ad4441a6

Example: Reconstructing Phylogenetic Trees

https://towardsdatascience.com/hierarchical-clustering-and-its-applications-41c1ad4441a6

Summary of Hierarchical Clustering

- Hierarchical clustering is often easier to visualize and interpret with a taxonomy
 - "Dendrogram" plots
- We don't need to specify the number of clusters up front!

Limitation: Doesn't scale well to big problems

- In studying Clustering techniques, we assume that we are given a matrix of distances between all pairs of data points.
- We assume that the input to the problem is:

- In studying Clustering techniques, we assume that we are given a matrix of distances between all pairs of data points.
- A distance measure (metric) is a function $d: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ that satisfies:

$$1. d(x, y) \ge 0, d(x, y) = 0 \iff x = y$$

$$2. d(x, y) + d(y, z) \ge d(x, z)$$

$$3. d(x, y) = d(y, x)$$

- For the purpose of clustering, sometimes the distance (similarity) is not required to be a metric
 - No Triangle Inequality
 - No Symmetry

Examples:

Euclidean Distance:

$$d(x,y) = \sqrt{(x-y)^2} = \sqrt{(x-y)^T (x-y)} = \sqrt{\sum_{i=1}^d (x_i - y_i)^2}$$

Manhattan Distance:

$$d(x,y) = |x - y| = \sum_{i=1}^{d} |x_i - y_i|$$

Infinity (Sup) Distance:

$$d(x, y) = \max_{1 \le i \le d} |x_i - y_i|$$

- Notice that if d(x, y) is the Euclidean metric, $d^2(x, y)$ is not a metric
- But can be used as a measure (no triangle inequality)

© 2019-22 D. Jayaraman, O. Bastani, Z. Ives

- Examples:
 - Euclidean Distance:

$$d(x,y) = \sqrt{(x-y)^2} = \sqrt{(x-y)^T (x-y)} = \sqrt{\sum_{i=1}^d (x_i - y_i)^2}$$

Manhattan Distance:

$$d(x, y) = |x - y| = \sum_{i=1}^{d} |x_i - y_i|$$

Infinity (Sup) Distance:

$$d(x, y) = \max_{1 \le i \le d} |x_i - y_i|$$

Euclidean: $(4^2 + 2^2)^{\frac{1}{2}} =$ 4.47 Manhattan: 4 + 2 = 6Sup: Max(4,2) = 4

- Notice that:
 - Infinity (Sup) Distance < Euclidean <u>Distance < Manhattan Distance</u>:

$$L_{\infty} = \max_{1 \le i \le d} |x_i - y_i| \quad L_2 = \sqrt{(x - y)^2} = \sqrt{\sum_{i=1}^d (x_i - y_i)^2} \quad L_1 = |x - y| = \sum_{i=1}^d |x_i - y_i|$$

• But different distances do not induce same order on pairs of points

https://scikit-learn.org/stable/modules/clustering.html#clustering

Summary of Clustering

- Critical to understanding the structure of our data
- Often useful for creating high-level features useful for supervised learning
- We saw two approaches: k-Means vs hierarchical clustering