

CLUSTERING

2020-05-05
Mags Giust



APRIL 30, 2020 CLUSTERING

"Life comes in clusters, clusters of solitude, then a cluster when there is hardly time to breathe"

- May Sarton (Eleanor Marie Sarton)
Poet, Novelist, Memoirist
1912-1995

"Over time, this growing tendency of like marrying like will only reinforce clustering and geographic sorting along class lines, giving the emerging map of social, economic, and cultural segregation even greater permanence."

- Richard Florida
Urban Studies - Social & Economic Theory
Prof. @ Univ. of Toronto
1957 - present

Unsupervised Learning

DIMENSIONALITY REDUCTION

PCA
t-SNE
Factor Analysis

CLUSTERING

AD, NLP may be solved with clustering.

Algorithms

KMeans

DBScan

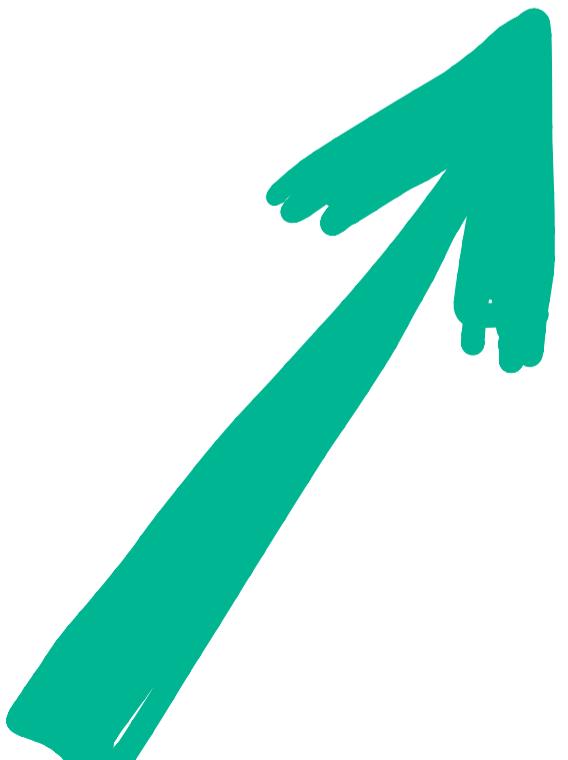
Hierarchical

MACHINE LEARNING

Reinforcement Learning

AD, NLP, RS may be solved with reinforcement learning

Collaborative Filtering
Content Filtering
Neural Networks (DL)



Supervised Learning

REGRESSION

• Continuous Target
[TSA, AD may be regression problems.]

Linear Regression
GLM
Polynomial Reg.
SVR
Decision Tree Regr.
Neural Networks

CLASSIFICATION

• Discrete, Categorical Target
[TSA, NLP, AD may be classification problems]

Logistic Regression
Decision Tree
Random Forest
SVC
KNN
Neural Networks

SUB-METHODOLOGIES

TSA: time series analysis

AD: anomaly detection

NLP: natural language processing

RS: recommendation system



MAGGIE, customer intel analyst
RACKSPACE, c. 2010

I hear you can do Magic with
numbers... so ...
Who are our Cloud
Customers?



LANHAM, CEO
RACKSPACE

MAGGIE,
aspiring data scientist

...They love ...
what the f...
uh-oh!

Wrangle

Plan

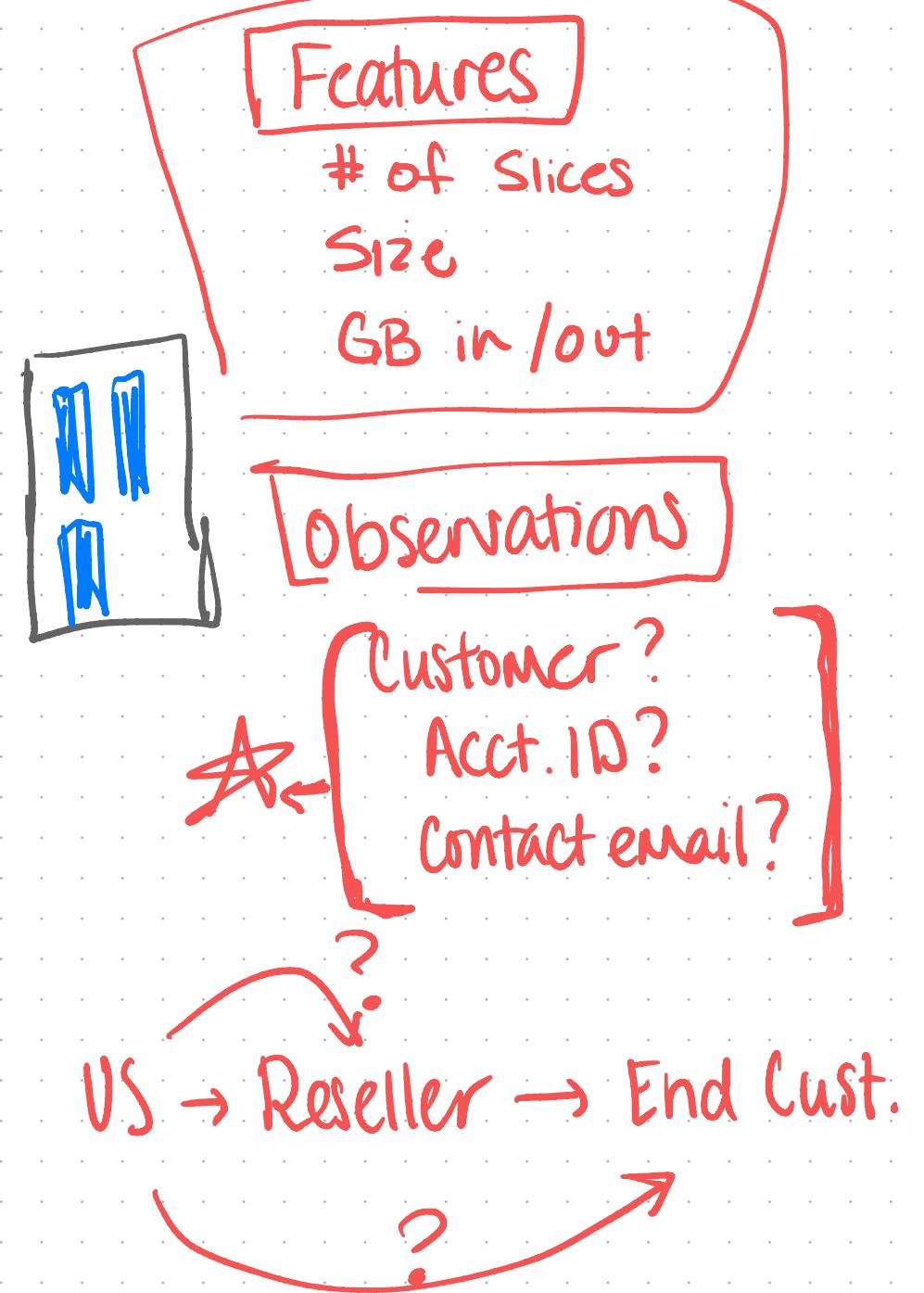
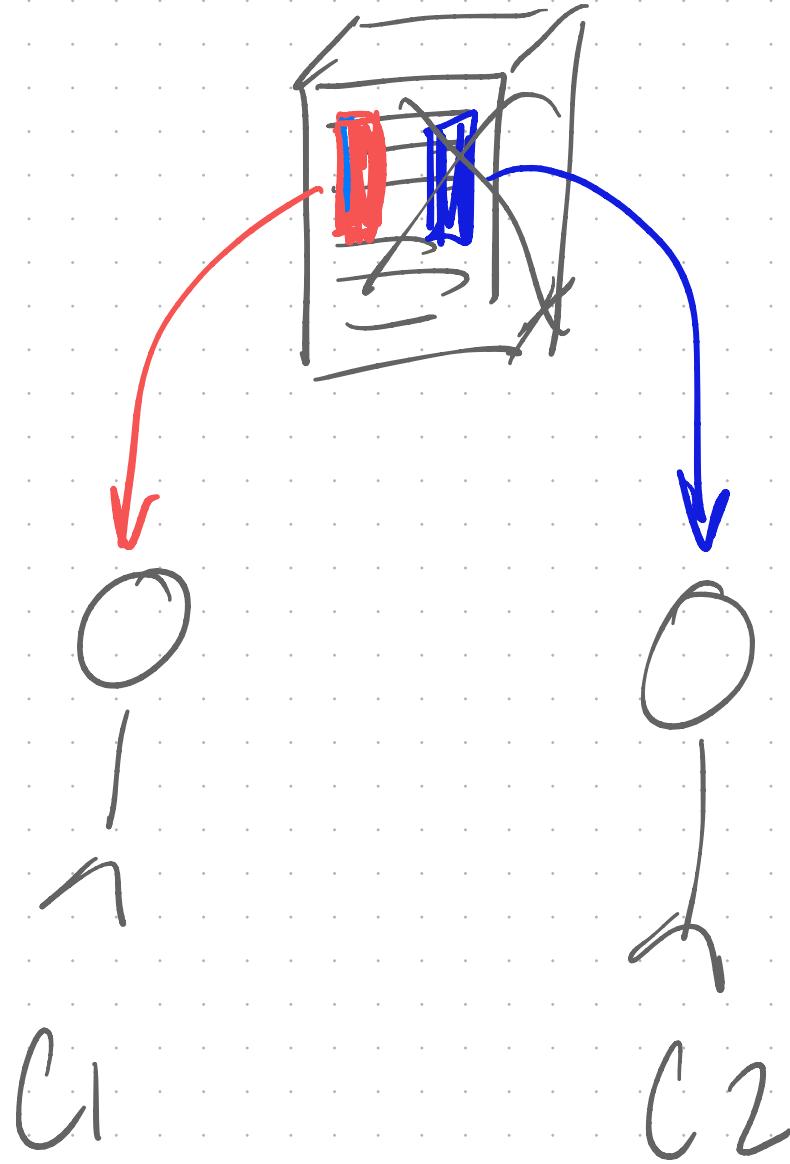
- define the ?, deliverables, goals - CLEARLY
- get necessary domain knowledge •
- early hypotheses •
- data definitions, locations, etc. •
- timeline? how much do I have?
how bad do I want this?

Explore/Analyze

Deliver

Model

This Pipeline heuristic was not available / clear to much of anyone at this time.



Wrangle

Plan

- define the ?, deliverables, goals - CLEARLY
- get necessary domain knowledge
- early hypotheses
- data definitions, locations, etc.
- timeline? how much do I have? how bad do I want this?

Deliver

- Common Usage Patterns of Cloud Customers!
- NTH : Profile of each of those patterns - Who, how Many, \$, ...

How do I define usage?

Explore / Analyze
CLUSTERING

Find customers with
Similar Usage Patterns

~~Model~~
Not yet!

No labels =
Unsupervised

Plan

- define the ?, deliverables, goals - CLEARLY
- get necessary domain knowledge
- early hypotheses
- data definitions, locations, etc.
- timeline? how much do I have? how bad do I want this?

Deliver

- Common Usage Patterns of Cloud Customers!
- NTH : Profile of each of those patterns - Who, how Many, \$, ...

Wrangle

- hundreds of millions of records
- tens of tables
- handful of DB's
- NOT indexed well!
- overnormalized
- the "unique" ID vs. Acct vs. Customer issue

Scale data



- How to aggregate
- How to sample
- How much history?
- What to measure
- What is an "observation"

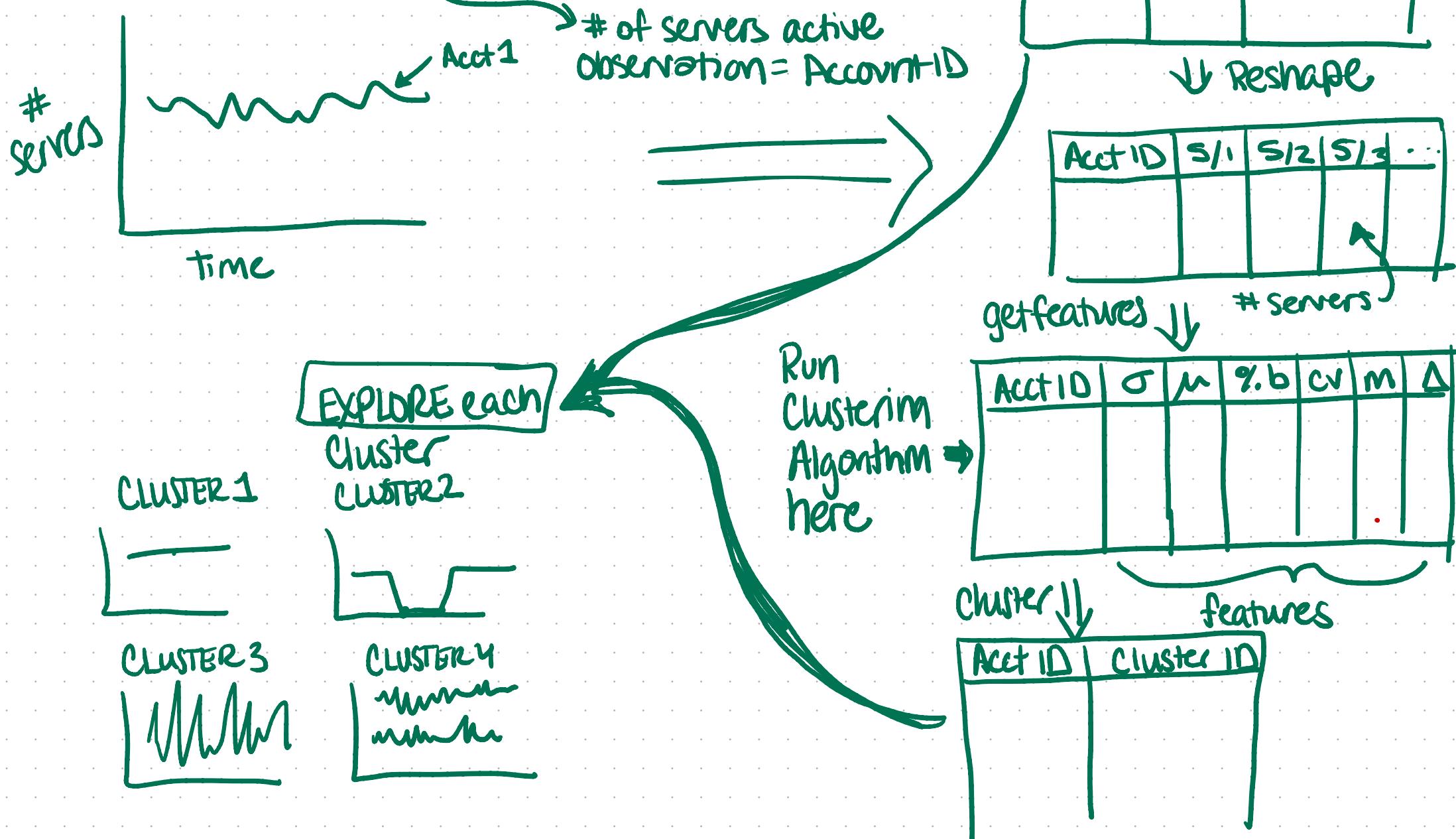
Explore/Analyze

- lots of Viz.
- statistical analysis
- CLUSTERING

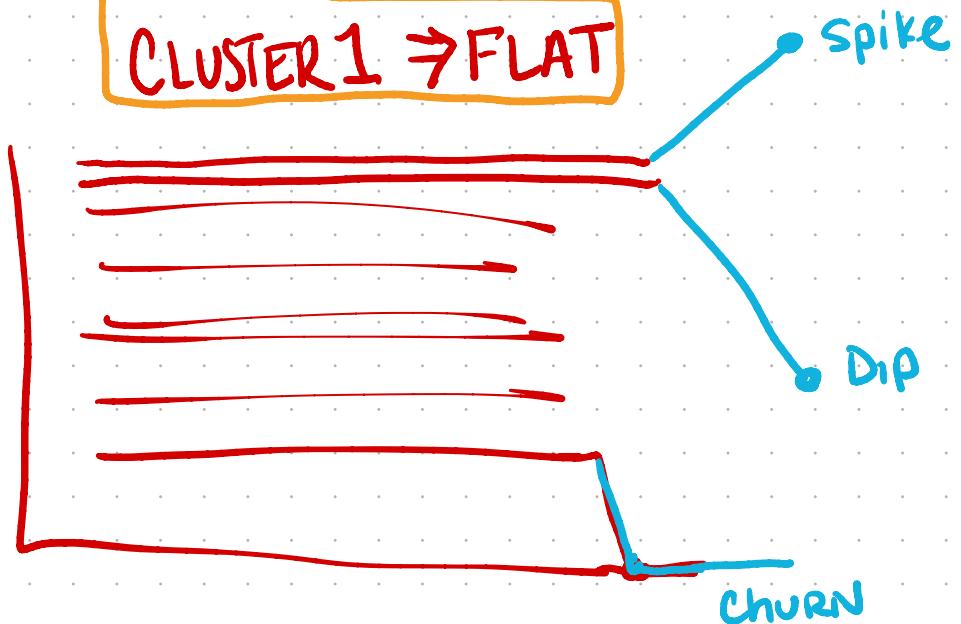
- ↓
- ★ Find customers whose usage patterns are similar
 - define usage

Model
Not yet!

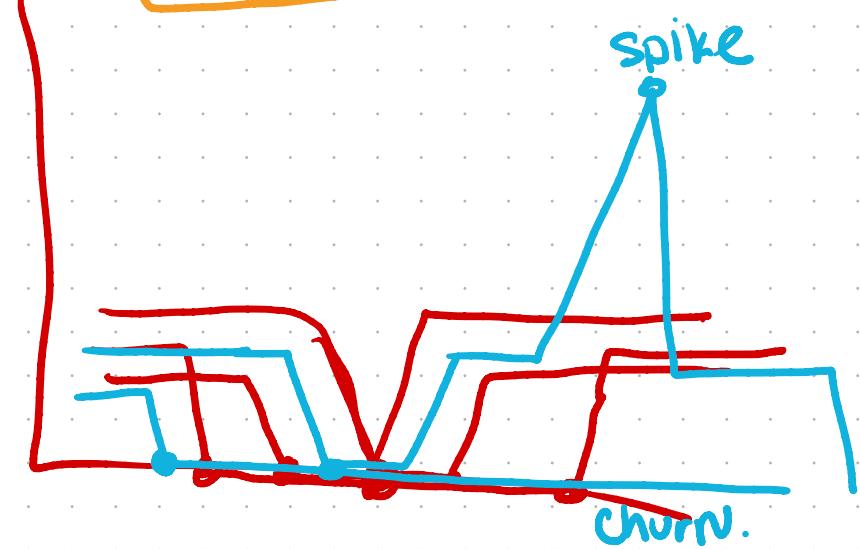
How do customers use Cloud Hosting? What are the common usage patterns?



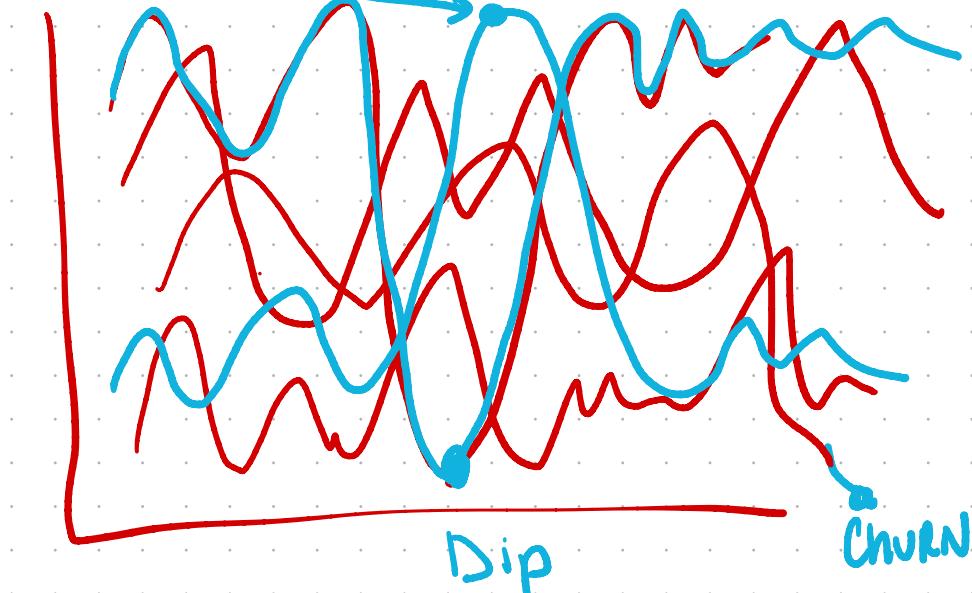
CLUSTER 1 → FLAT



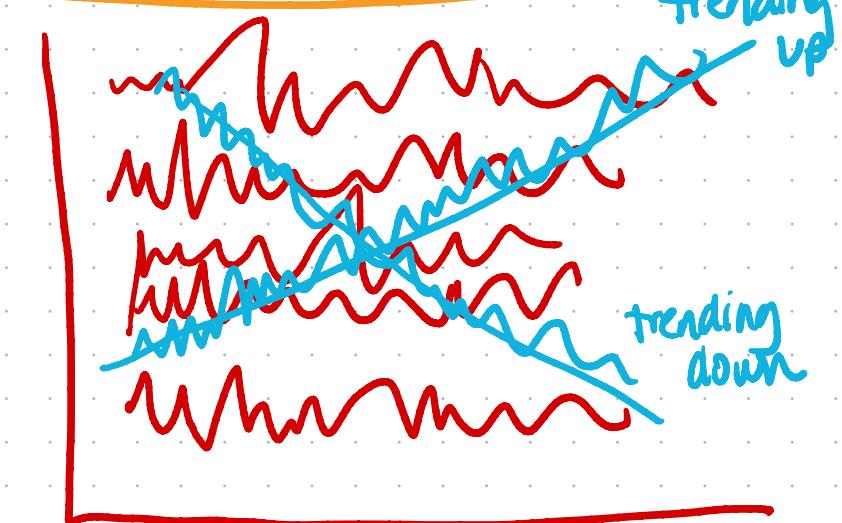
CLUSTER 2 → SPORADIC



CLUSTER 3 → MODERATE / HIGH VOLATILITY



CLUSTER 4 → SMALL VOLATILITY



What do you do with Clusters?

- explore your new clusters, compare & contrast, analyze
- create new features for a future supervised model
- split observations into cluster to build a better model
for each cluster -
makes patterns easier
to detect / extract!
- use clusters as labels for supervised learning - (classification, usually)

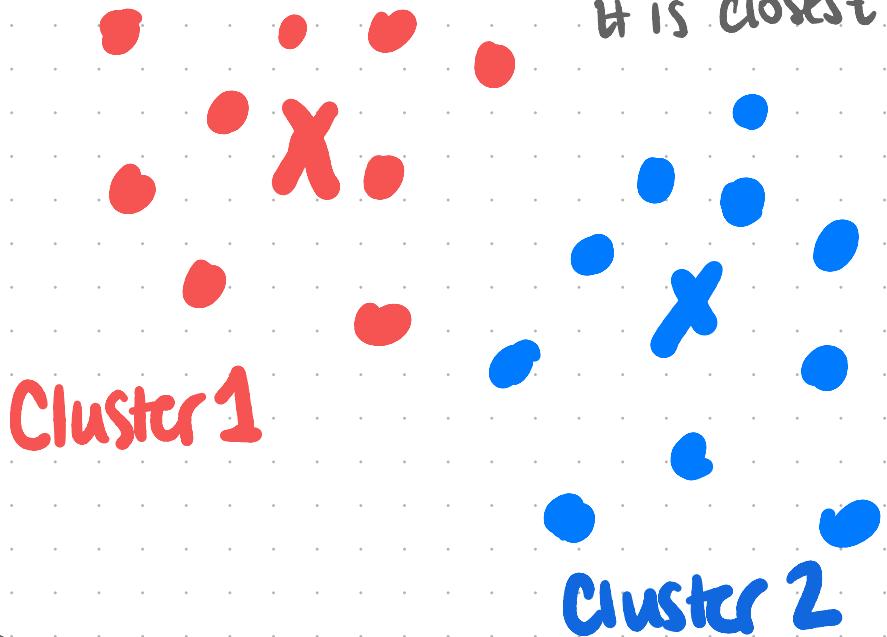
When is Clustering Useful?

- targeted marketing - cust. segmentation, market research
- anomaly detection - account take over, fraud, data breach
- crime zones, housing prices
- text - document classification
- image processing

CLUSTERING

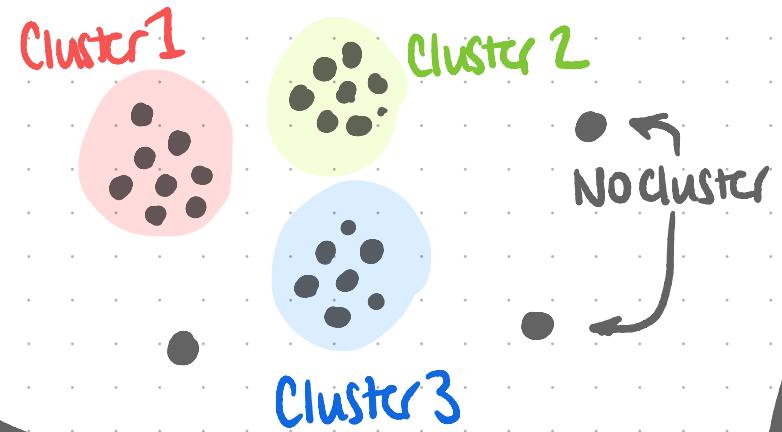
K-Means

$k = \# \text{ of cluster}$
a point is put in the cluster whose centroid (x) it is closest to.

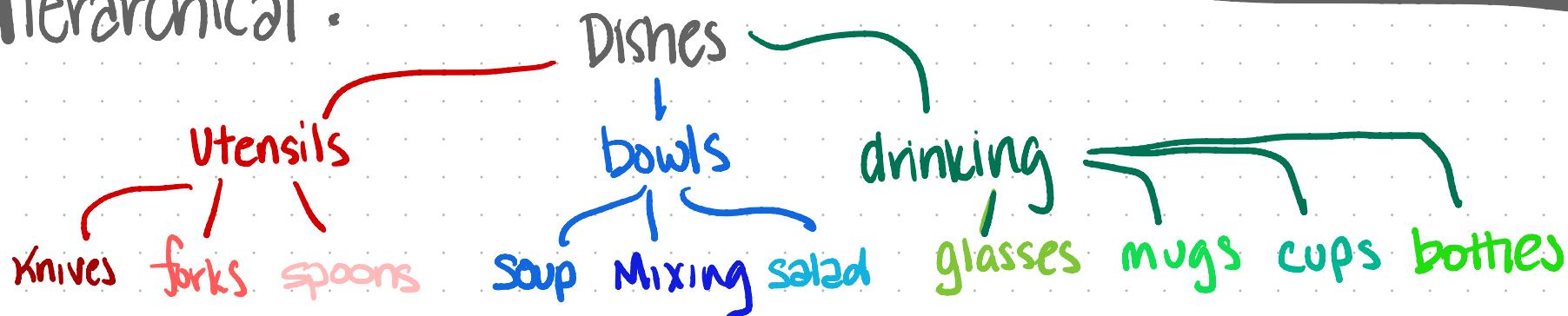


DBScan

A cluster is formed when it has a min. # of pts. within a max distance.

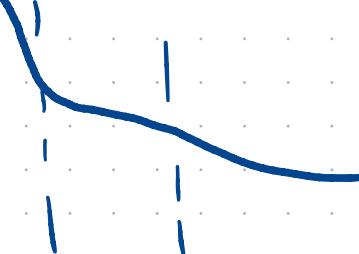


Hierarchical :



① KMeans: SCALED DATA! Which Vars?

② Elbow Method



③ Range of k \Rightarrow Get # of observations / cluster
select K (There IS NOT ONE right answer)



④ Run KMeans(k, random_state)

fit_predict(train_scaled[cluster_vars])
predict(test_scaled[cluster_vars]) \Rightarrow

2 arrays · $\frac{1}{1}$ length of train
 $\frac{1}{1}$ length of test

← add to DF

⑤ Join cluster IDs to Df's (x)
var1 | var2 | ... | var n | cluster_id

⑥ get centroids of each cluster:

cluster_id	centroid_var1	centroid_var2 ...
------------	---------------	-------------------

⑦ join back to Df's var1 | ... | var n | cluster_id | centroid_var1 | centroid_var2 |
(x)

- ① SQL - filter !!!
 - ② NULLS
 - ③ DUMMY Vars of County, & City of LA
 - ④ New features - reduce noise & dependencies
 - ⑤ Remove Outliers for now - Need to Reduce Noise to find Patterns!
 - ⑥ SAMPLE Data as we try out things -
most diff. will be LA, so use LA County
 - ⑦ Split into train / test
 - ⑧ Scale
 - ⑨ Cluster Area - LAT, LON, AGE
 - ⑩ Cluster Size - Sqft, Acres
 - ⑪ once clusters are in df's, group by the 2 cluster ID fields & get summary stats.
 - ⑫ join Median + StdDev back to df's
 - ⑬ use StdDev as feature
 - ⑭ use Median + \$/sqft to compute new feature $\$/\text{sqft} - \text{Median}$
- Create binned vars
- AGE
 Tax Rate
 Acres
 \$/sqft
 \$/acre
 bed:bath Ratio

What have we done?

- Captured area in a useable format (numeric)
- Separated dependencies on var
- extracted a feature that measures variance in assessed value
- decreased dimensions -