

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

Overtime, 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

LEARNING GOALS

- Clustering use cases
- General knowledge of common clustering algorithms
- Implement k-Means using SKLearn
- Make use of clusters later down the pipeline.
- Handle outliers using IQR
- More practice scaling
- Strategies for Missing values
- Plotting Clusters

Machine Learning

SUPERVISED

CLASSIFICATION REGRESSION

TARGET:

Discrete
Categorical

Continuous
numeric, ordered

SUB-METHODOLOGIES:

time series
Anomaly detection
NLP

time series
anomaly detection

ALGORITHMS:

Logistic Regression
SVM
Decision Tree
Random Forest
KNN
Neural Network

OLS (linear regression)
GLM (tweedie regressor)
Polynomial Regression
Support Vector Regressor
Decision Tree Regressor
Neural Networks
LASSO, LARS

UNSUPERVISED

CLUSTERING DIMENSIONALITY REDUCTION

SUB-METHODOLOGIES:

Anomaly Detection
NLP

ALGORITHMS:

K-Means
DBSCAN
Hierarchical

PCA
t-SNE
factor analysis

REINFORCEMENT

SUB-METHODOLOGIES

Anomaly Detection
NLP
Recommender Systems

ALGORITHMS:

Collaborative Filtering
Content Filtering
Neural Networks

CLUSTERS & CAREERS

They hate us!!
Oh, they love
us! 😊
OMG these
surveys ...



Maggie, Customer Intel Analyst
Rackspace C. 2010

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



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

Lanham Napier, CEO
Rackspace

WHO ARE OUR CUSTOMERS?

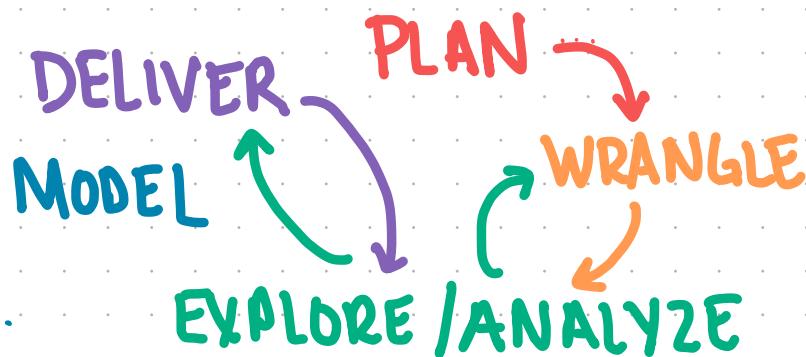
Why? are we meeting customer needs? support strategy for scaling support?
target customer acquisition? fire customers?

Define the Question

- A report w/ visualizations exhibiting the top discovered usage patterns (clusters) w/ useful labels applied
- Use NPS to see if one cluster is significantly more engaged than others.
- Add info about MRR + # of accounts.
- Any specific customer use cases I can discover w/ in each cluster.

Will not do predictive modeling yet.
We must assess the usefulness
of the clusters created before
taking that step.

How do our customers use our cloud hosting product?
Can we identify <10 common patterns of number of active servers daily, overtime,
across accounts such that when we know a customer follows pattern x, we then
know what to expect from customer. So when usage deviates from the expected,
THAT'S when we know to take action, but not everytime a customer, any customer
spins down all their servers.



GOAL

- Identify the most common usage patterns of cloud hosting by clustering accounts based on their patterns of usage measured by number of slices.
- Define, label those patterns discovered.

How?

- Use k-means, No time series.
- find ways to represent the usage patterns + attributes seen overtime through summary stats + other creative summarizing features.

Details to consider:

- Group by account ID.
- Start w/ history over last 3 months
- Exclude accounts whose 1st server was spun up < 4 months ago.
- # of servers in day = # of distinct server IDs 'active' status at any point in the day.

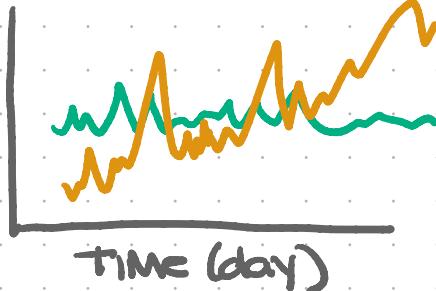
challenges

- hundreds of millions of records
- indexing not ideal
- tens of tables, handful of databases
- permission issues in ability to write
- no unique customer identifier
- accounts may become stale but not closed.

RAW DATA → NAMED CLUSTERS

Acct # 112358, 235813

Slices



Acct.	Date	Count
112358	10/09	5
:	10/09	6
112358	7/11	3
235813	10/09	4
:	10/09	10
235813	7/11	9
235813	7/10	2

90 days / Acct
10k Accounts

Reshape
1 Row /
Account

Account	7/10	7/11	...	10/8	10/9
112358	4	3	...	6	5
235813	1	2	...	9	10

Why?

- * Ensure every account has an entry on every day. i.e. fill missing days w/ 0.
- * 1 Row per account.

Explore, Understand,
Analyze the clusters

- do they make sense?
- do I need to refine features?
- More/Less Clusters?
- Descriptive names of clusters?

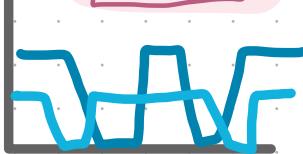
cluster 1

FLAT



CLUSTER 2

SPORADIC



cluster 3

HIGH VOLATILITY /



CLUSTER 4

LOW VOLATILITY



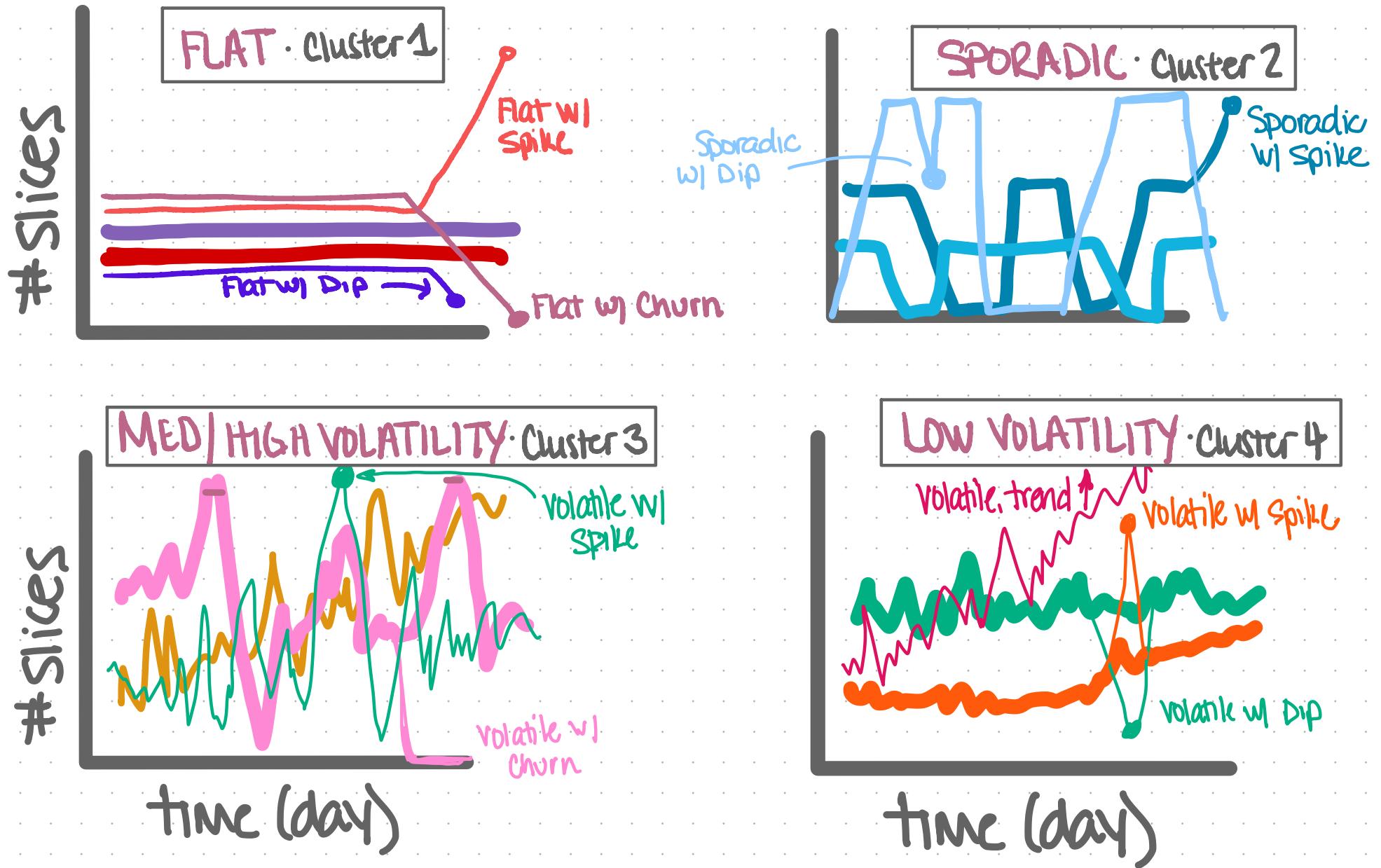
Account	$\bar{\mu}$	$CV(\bar{\mu})$	$\% \Delta$	UBB	LBB	slope	day/day Δ	$\% \Delta$	MA	zero-use days	Cluster
112358											4
235813											3

< data is sampled (train, validate, test)
=> scaled >

Cluster the observations
using k-Means

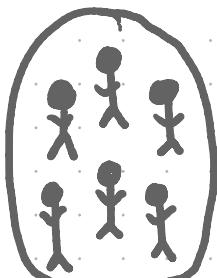
- A cluster ID is returned (an integer)
- That ID means nothing in/of itself
- All it is is a link that ties all the accounts in that cluster together.
- The ID could be "pink" if it wanted!

THE CLUSTERS



WHEN IS CLUSTERING USEFUL?

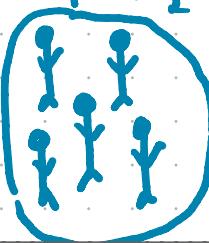
Fresh only



Our own Brand



FROZEN + Pantry



Marketing, Customer Segmentation

Anomaly Detection

file Access



.py files

HR files

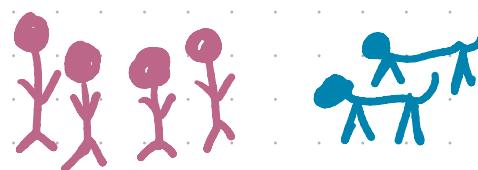
Document Classification

Support
UpTime
Support
Price

SURVEY response..12:00
Survey response 12:01
Survey response.. 1:05
Survey Response .. 1:11
...

THEMES of responses

Image Processing

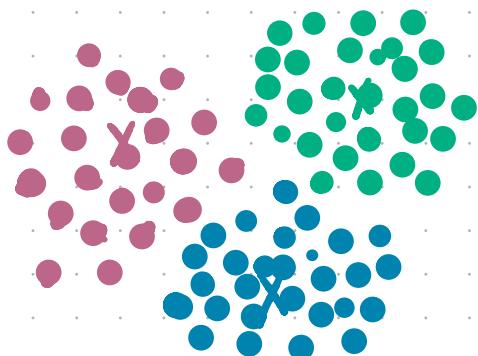


What are the different
subjects of these photos?

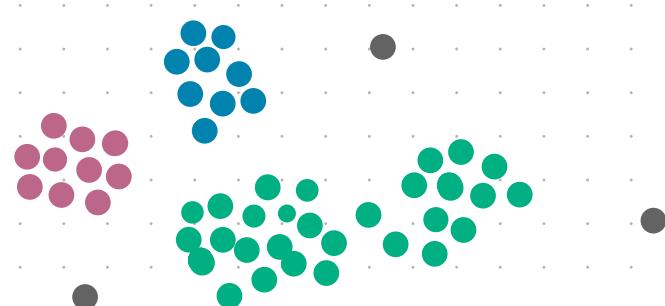
CLUSTERING ALGORITHMS

K-Means

$k = \# \text{ of clusters}$



DBSCAN



Hierarchical

DISHES

utensils

forks knives spoons

dowls

soup Mixing Salad

drinking

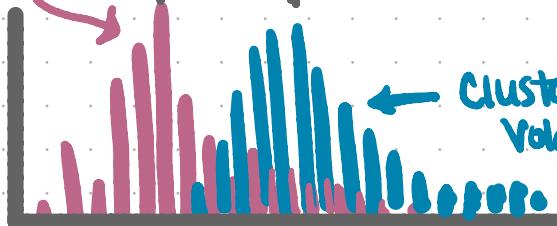
mugs glasses cups bottles

I HAVE CLUSTERS. Now WHAT?

EXPLORE & ANALYZE CLUSTER GROUPS

Cluster 1:
flat usage

Do volatile customers
spend more?



Monthly Spend

Cluster 2:
Volatile Usage

MODEL EACH CLUSTER SEPARATELY

PREDICTING CHURN

Model 1 - flat cluster

```
.fit(X_train_flat, y_train_flat)
```

Model 2 - volatile cluster

```
.fit(X_train_volatile, y_train_volatile)
```

Model 3 - sporadic cluster

```
.fit(X_train_sporadic, y_train_sporadic)
```

TURN CLUSTERS INTO FEATURES

target features

Account	churn	tenure	MRR	is_Flat	is_volatile	is_sporadic
112358	1	12	100	0	1	0
123581	0	7	200	0	1	0
235813	0	9	10	1	0	0
358132	1	15	20	0	0	1

TURN CLUSTERS INTO LABELS

Classify NEW customers' usage
patterns using labels identified
through clustering
target

Account	cluster	σ	μ	$\%b$	$\Delta d/d$
112358	volatile				
123581	volatile				
235813	flat				
358132	sporadic				

Machine Learning

SUPERVISED

CLASSIFICATION REGRESSION

TARGET:

Discrete
Categorical

Continuous
numeric ordered

SUB-METHODOLOGIES:

time series
Anomaly detection
NLP

time series
Anomaly detection

ALGORITHMS:

Logistic Regression
SVM-classifier
Decision Tree
Random Forest
KNN
Neural Network
Naïve Bayes

OLS (linear regression)
GLM (tweedie regressor)
Polynomial Regression
Support Vector Regressor
Decision Tree Regressor
Neural Networks
LASSO, LARS
Holt-Winters
Auto Regression

Anomaly Detection Module

UNSUPERVISED

CLUSTERING DIMENSIONALITY REDUCTION

SUB-METHODOLOGIES:

Anomaly detection
NLP

ALGORITHMS:

K-Means
DBSCAN
Hierarchical

PCA
t-SNE
Factor Analysis

REINFORCEMENT

SUB-METHODOLOGIES

Anomaly Detection
NLP
Recommender Systems

ALGORITHMS:

Collaborative Filtering
Content Filtering
Neural Networks