# Frequent Itemset Mining with tidyclust in R

Advisor: Dr. Bodwin Committee: Dr. Glanz and Dr. Dekhtyar

# **1.** What is Frequent Itemset Mining 2.Explaining tidyclust **3.Clustering Frequent Itemsets 4.Predicting with Frequent Itemsets 5.**Future Directions





# What is Frequent Itemset Mining?

## **Finding Co-occurrences**

- Frequent Itemset Mining (FIM) helps us find items that often appear together in a dataset.
- For example, FIM may discover that customers who buy coffee also frequently buy milk.
- Key Terms
  - Itemset: A <u>collection of items</u>.
  - Frequent Itemset: A collection of items that are <u>commonly found together</u>.
  - Minimum Support: The <u>cutoff point</u> that determines if an itemset is frequent.



beer	milk	bread
0	1	1
1	0	1
1	1	0
1	1	1
0	1	1



# **Uses for Frequent Itemset Mining**

activities.

05/22/2025

### **Market Basket Analysis Understanding customer buying habits** (e.g., products frequently purchased together).

### **Medical Diagnosis** Identifying co-occurring symptoms or conditions.

### **Fraud Detection**

**Detecting unusual combinations of transactions or** 

# A Refresher on tidyclust.

## tidyclust Workflow

- A user can:

```
fi_spec <- <u>freq_itemsets(</u>
 min_support = 0.05,
 mining_method = "apriori",
   |>
  set_engine("arules")
fi_spec_fit <- fi_spec |>
  fit(~., data = groceries)
fi_spec_fit |>
 predict(new_groceries)
```

• tidyclust is a set of tools in R that helps find patterns in data without pre-defined labels (unsupervised learning).

• The goal of tidyclust is to establish a consistent and reproducible workflow.

• Specify a model • Fit a model • Predict with a model all using a standardized syntax.



# Clustering Frequent Itemsets.

## Clustering

- Group similar <u>data points</u> (rows) together.
- Finds natural groupings based on data • characteristics.

		bill		body
sex	island	length	bill depth	mass
male	Torgersen	39.1	18.7	3750
female	Torgersen	39.5	17.4	3800
female	Torgersen	40.3	18	3250
female	Torgersen	36.7	19.3	3450
male	Torgersen	36.8	20.6	3650

- transactions.

beer
0
1
1
1
0

## **Clustering with FIM**

• Groups similar <u>items</u> (<u>columns</u>) together.

• Groups items based on co-occurrence in

milk	bread	diapers	eggs
1	1	0	0
0	1	1	1
1	0	1	0
1	1	1	0
1	1	1	0

## **Finding Frequent Itemsets**

	beer	milk	bread	diapers	eggs	Num
	0	1	1	0	0	<b>Support(X)</b> =
	1	0	1	1	1	
	1	1	0	1	0	<b>Minimum Support</b>
	1	1	1	1	0	
	0	1	1	1	0	
						itemset su
			item	support		{beer, bread}
			beer	0.6		{beer, diapers}
4p te	emset is freq	<u>pie</u> : II an juent,	milk	0.8	$\rightarrow$	{beer, milk}
the State	en so are all bsets	its	bread	0.8		{bread, diapers}
	<b>DSCLS.</b>		diaper	s 0.8		{bread, milk}
			eggs	0.2		{diapers, milk}

### nber of transactions containing itemset X

### **Total number of transactions**

**:t:0.5** 





05/22/2025

### **Handling Outliers**

Items not in any frequent itemset are outliers.

# Predicting with Frequent Itemsets.

beer	milk	br	ead	diapers	eggs			
1	NA		0	1	0			
	t							
itemset		support						
milk		0.8						
bread		0.8						
diamana		0.0		itemset		support	confidence	
mapers		0.8	-	{diapers, mill	k}	0.6	0.75	
beer		0.6			ţ			-
{beer, diapo	ers}	0.6			milk			
{bread, diap	ers}	0.6			0.75		(	201
{bread, mi	lk}	0.6						
{diapers, m	ilk}	0.6						

# Predicting with Frequent Itemsets

- iven a partial set of items, for each missing em:
- . Find Relevant Frequent Itemsets: Find requent itemsets that contain the missing tem and at least one observed item.
- . Calculate Confidence: The likelihood of the nissing item appearing, given the observed tems.

 $idence(X \rightarrow Y) = \frac{Support(X \text{ or } Y)}{Support(observed items in X)}$ 

**Predict:** Average these likelihoods.

## k-means Output

.pred_cluster
Cluster_1
Cluster_1
Cluster_2
Cluster_3
Cluster_2

## FIN Output

.pred_cluster						
	< <b>df [5 x 3]</b> >					
	< <b>df [5 x 3]</b> >					
	< <b>df [5 x 3]</b> >					
item	.obs_item	.pred_item				
beer	0	NA				
milk	NA	1 (0.75)				
bread	1	NA				
liapers	0	NA				
eggs	0	NA				

# Future Additions.

**Cross Validation** Implement item-stratified cross-validation for robust model evaluation.

## **Future Directions**

### **Parameter Tuning** Automate minimum support tuning with datadriven ranges.

## **Enhance Prediction**

Test different confidence weight metrics.

## **Tuning Minimum Support**

- The optimal minimum support value varies • depending on the characteristics of the dataset.
- I propose
  - 1. Calculate the <u>mean support</u> of all 1-itemsets:

$$\mu = \frac{1}{n} \sum_{i=1}^{n} support(\{item_i\})$$

2. Create a confidence-interval-like <u>range around  $\mu$ </u> using the standard deviation  $\sigma$ :

$$\left[\mu - \frac{\sigma}{2}, \mu + \sigma\right]$$

clipping the bounds at [0, 1] to ensure valid support values.

The range is **asymmetric** since transactional datasets tend to have a **<u>right skewed</u>** support distribution.



performance, Martin et. al. 2009

From Fundamental basket size patterns and their relation to retailer

# Putting Everything Together.

## **Functions!**

```
freq_itemsets():
     .freq_itemsets_fit_arules()
<u>fit():</u>
     extract_cluster_assignment.itemsets()
     item_assignment_tibble_w_outliers()
predict():
     itemsets_predict_helper()
     .freq_itemsets_predict_raw_arules()
     .freq_itemsets_predict_arules()
     extract_predictions()
     augment_itemset_predict()
Misc.():
     extract_fit_summary_items()
     tunable.freq_itemsets()
     random_na_with_truth()
```



