

Market Basket Analysis

Provides Answer - How likely customer will buy product A if they already have products B, C, etc.

	A	B	C	D	E	F
P1	✓	✓			✓	✓
P2			✓	✓		
P3	✓	✓			✓	
P4		✓	✓		✓	

↳ Basket = Rows

Idea → People who buy {A,B} likely buy {B,C} as well.

Definitions

Support: Number of baskets (rows) containing the set $\{u, v, \dots\}$
 (can also be a percent or fraction)

Support Threshold "s": Sets $\{u, v, \dots\}$ that appear in $\geq s$ baskets are Frequent Itemsets,

$$\text{Confidence } (I \rightarrow j) = \frac{\text{Support}(I \cup j)}{\text{Support}(I)} \rightarrow \frac{\text{baskets that satisfy } I \rightarrow j}{\text{baskets with } I}$$

Confidence Threshold "c": Sets $\{i_1, i_2, \dots\}$ where $\text{Confidence}(I_j, j) \geq c$
 Association Algebra.

1. Find sets I that satisfy support threshold
2. From above result, calculate confidence and only keep the best

Logical Error : If $\{a, b, c\} \rightarrow \{y\}$ has ST $\geq s$ and CT $\geq c$, then
 Subset of the set also satisfies the conditions

Correct: There is NO SPEEDUP for Apriori ..

~~Proof : Let $s(I) = \text{support}(z)$, $c(I) = \text{confidence}(z)$, $c(I, y) = \text{confidence}(I \cup y)$~~

~~Let $S_{T_0} = s(S_{A,B,C})$, $C_{T_0} = c(S_{A,B,C}) = \frac{s(S_{A,B,C,Y})}{s(S_{A,B,C})}$~~

~~$S_n < S_{\text{anti}}$ because narrowing search space~~

~~Case 1: $S_{A,B,C}$ does not satisfy Support Threshold
 $S_{A,B}$ MAY but... check Case 2.~~

~~Case 2: $S_{A,B,C}$ does not satisfy confidence.~~

~~Proc: $c(S_{A,B,C}) \geq c(S_{A,B})$~~

$\begin{array}{ll} AB & Y \\ AB & Y \\ A\bar{B} & Y \\ A\bar{B} & C \end{array}$

The "proof" is invalid as well.

Truth (rip): There is no real algo optimization for Apriori
2 types of relationships

$S_{a,b} \rightarrow S_{a,b,c}$: Not true bc $s(S_{a,b,c}) \leq s(S_{a,b})$

$S_{a,b,c} \rightarrow S_{a,b}$: Not true bc universal generalization is invalid
Say we had  Then $c(S_{a,b}) \geq c(S_{a,b,c})$

$\text{Not } S_{a,b} \rightarrow \text{Not } S_{a,b,c}$ or $\text{Not } S_{a,b,c} \rightarrow \text{Not } S_{a,b}$ are contrapositives
of the two cases refuted

Any of the $x \rightarrow \text{Not } y$ or $\text{Not } x \rightarrow y$ do not make logical sense

Support Algebra rule!

If $\{a, b\}$ does not satisfy support, then $\{a, b, c\}$ will not either.

Proof: Since $\{a, b\}$ is a subset of $\{a, b, c\}$, then there cannot be more counts of a set than a subset.

Apriori Algo

Key Idea: If $\{a, b, c\}$ satisfies support, then all subsets do as well

Take Frequent Items (Items that appear $\geq s_f$ times) $O(N)$

Take Frequent Pairs (Pairs that appear $\geq s_f$ times) $O(N^2)$

Take F Triples, F Quads, etc $O(N^3)$ why?

Notation: A_n = set of length n

Recurrsively Define Merge: saves "Generating" time

Generating Triples is $O(N^3)$

Can further be optimized by
reproducing the Z

(say $\{m, n\}$ is bad, then $\{a, m\} \cup \{n\}$
should not be included)

BUT! this is worst case. Normally,

We expect the A and B sets to be pruned
So it will be much faster

Merge candidate sets is $O(N_n)$

Merge A_{n-1}, A_1

Say we had candidate sets $Z_n = \{A_1B_1, A_1B_2, \dots\}$

Then run support threshold, generate $A_n \vdots$

Generate all candidate sets $O(2^n)$ but optimized.

Confidence ::

No real rules.

Apriori

Run through support, calculate confidence, choose best sets.

Generally, $k=2$ has most memory (given reasonable support 1%)

PCY \rightarrow Park-Chen-Yu Alg.

Observation: In pass 1 of Apriori, most memory is idle (we only store item sets)

Can we use idle memory to reduce memory in pass 2? (pairs)

Pass 1 of PCY: Maintain hashtable with as many buckets as fit in memory

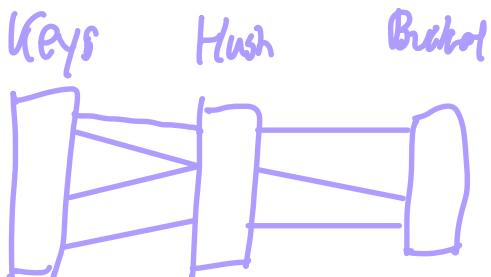
Keep a count for each bucket into which pairs of items are hashed

For each bucket, keep only the count ... not the actual pairs

pseudocode:

```
for (each row)
    for(item in row)
        freq[item] ++
    for (p1 in row)
        for(p2 in row)
            Generate hash
            bucket(hash) ++
```

Hashtable



A faster way to hash/save memory

If no bucket, all keys in memory

KEY OBSERVATION \rightarrow If a bucket contains a frequent pair, the bucket is surely frequent

However, a bucket without a FP is not always bad.

BUT: For a bucket with $\text{count} < 5$, no pairs are frequent.

Even if a pair has 2 frequent items, pairs in this bucket are eliminated as candidate pairs

Proof: Even if $\text{count}(\text{bucket}[h])$ all below 8 (say 6), $\boxed{\text{count} \leq 5}$

Between Passes: Replace buckets by bitset (0/1)

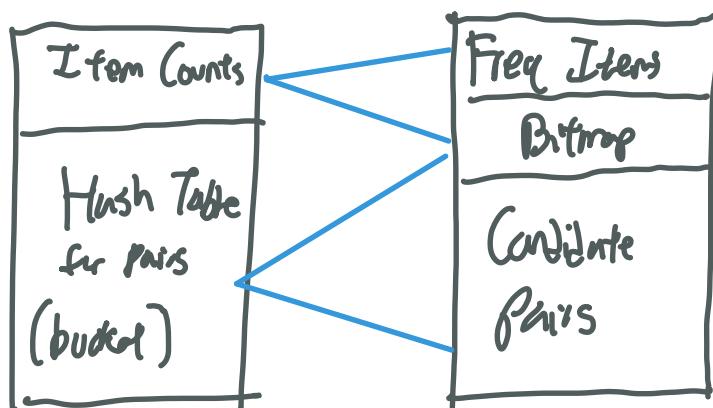
Memory reduction by 1/32 because 4-bytes $\rightarrow 1$ bit

Also list frequent items

Pass 2: Only count pairs in frequent buckets

Candidate Pair Rules $\{i, j\}$

1. i and j are frequent items
2. $\{i, j\}$ hash bucket entry is 1 (frequent bucket)



PCY in Action

Say we had a veggie shop that sells apples, bananas, Carrot, Avocado, blueberries, Celery

Table

Apple Avacado Bananas Blueberry Carrot Celery

	Apple	Avacado	Bananas	Blueberry	Carrot	Celery
1	X		X		X	
2		X		X		
3		X				X
4					X	X
5			X	X	X	
6	X	X	X	X	X	X

Define Hash(a, b) \rightarrow "a[0]" + "a[1]" $\rightarrow h("apple", "banana") = "ab"$

PCY Pass 1

Let Support = 4

Apple	Avacado	Bananas	Blueberry	Carrot	Celery
2	3	4	3	4	2

= F items

Total Table

	AA	AB	AC	BB	BC	CC
1	0	2	0	1	0	0
2	0	1	0	0	0	0
3	0	0	1	0	0	0
4	0	0	0	0	0	1
5	0	0	0	1	2	0
6	1	4	4	1	4	1
Total	1	7	5	3	6	2

Frequent Bits

AB, AC, BC



Exclude AA, BB, CC

< Hash Table

PCY Pass 2

Convert Hash Table Buckets into Bitset

AA	AB	AC	BB	BC	CC
0	1	1	0	1	0

Find Candidate Pairs!

From Frequent Value List (*Bananas, Carrot*)

Generate Pairs (*(Bananas, Carrot)*) and Hash (*BC*)

Check bitset [*BC*] = 1 → Yes

Pruning complete.

Try conditions (*BC*) → only 2 cases, not 4 ;;

Moral of the story:

1. PCY Narrows search space by a lot!

2. Candidate Pairs must be made up of 2 freq. values (rule)

3.