

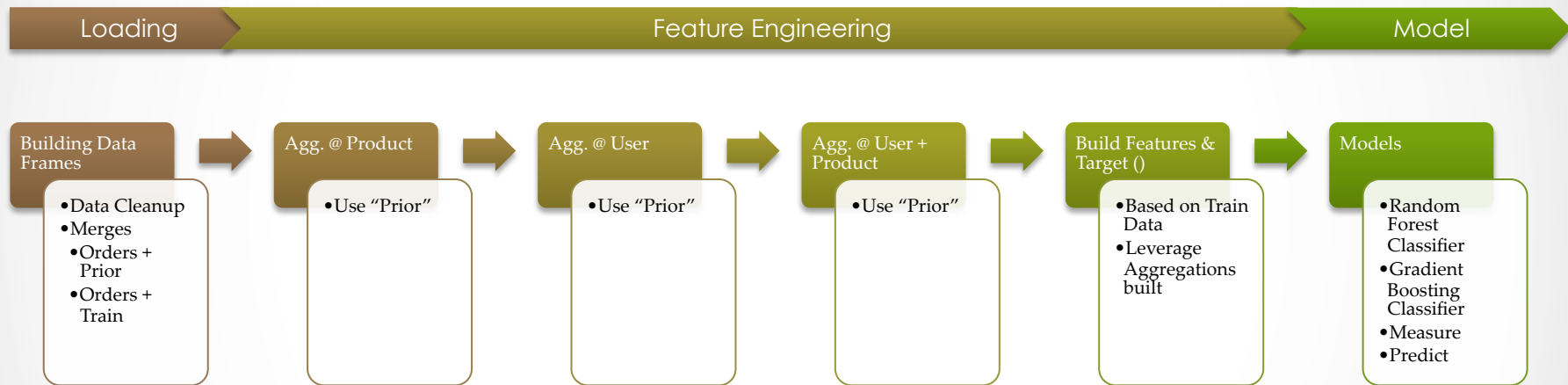
Kaggle Challenge Instacart

Anwar Habeeb

Problem

- Instacart, a grocery ordering and delivery app, aims to make it easy to fill your refrigerator and pantry with your personal favorites and staples when you need them. After selecting products through the Instacart app, personal shoppers review your order and do the in-store shopping and delivery for you.
- Use transactional data to develop models that predict which products a user will buy again, try for the first time, or add to their cart next during a session.
 - **3 Million** Instacart Orders, Open Sourced.
- Using this anonymized data on customer orders to predict which previously purchased products will be in a user's next order.

Approach

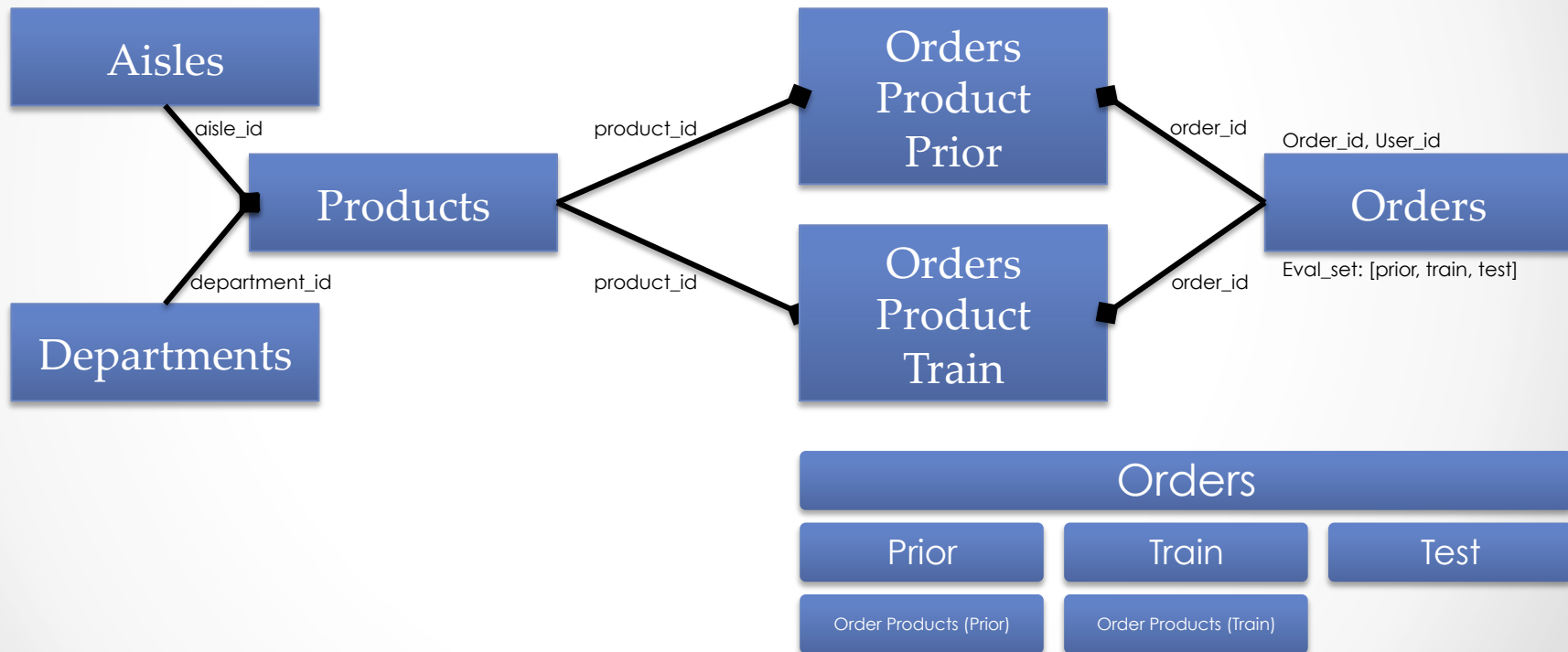


Data Model

...

Data Model

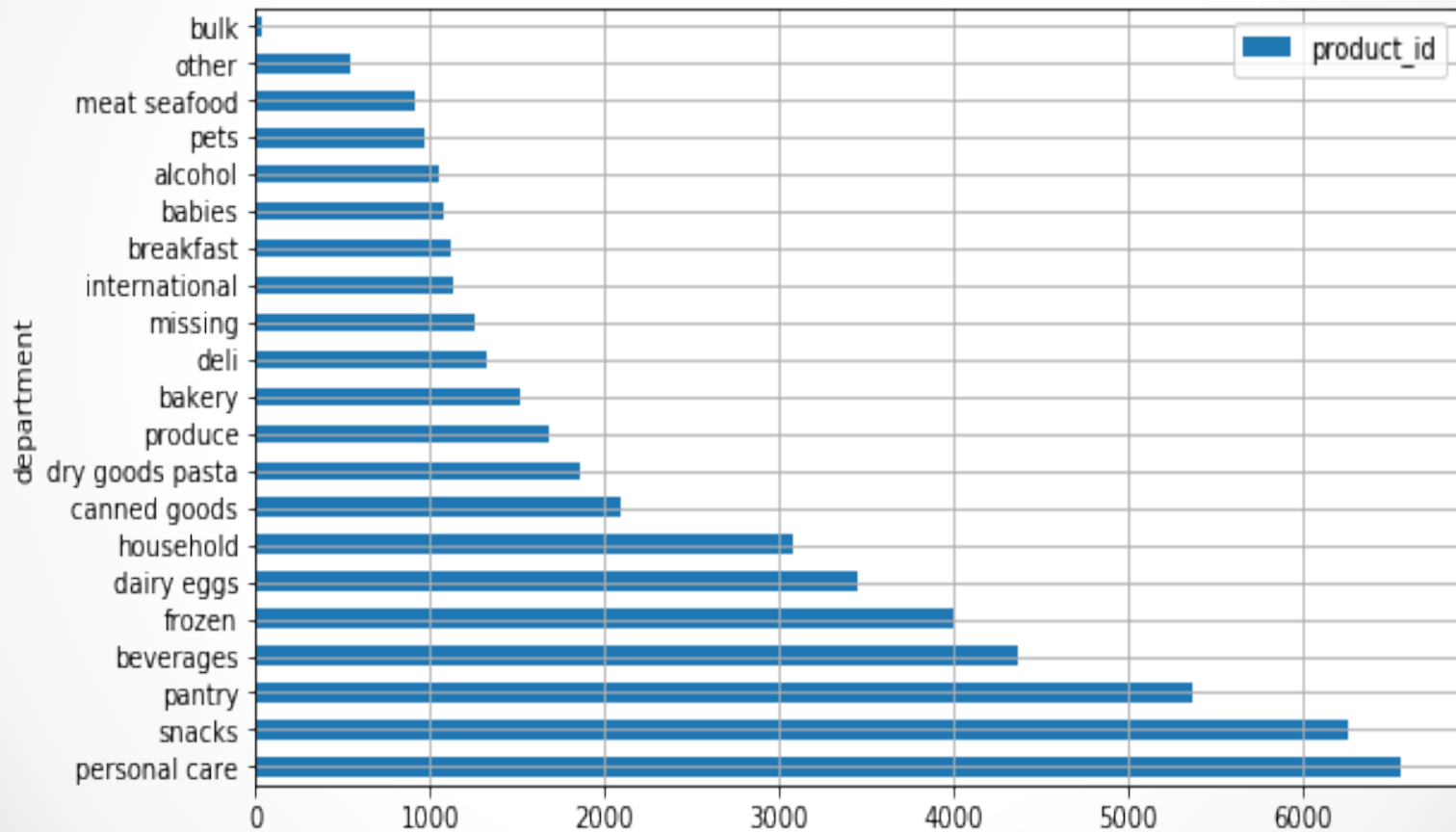
1:M



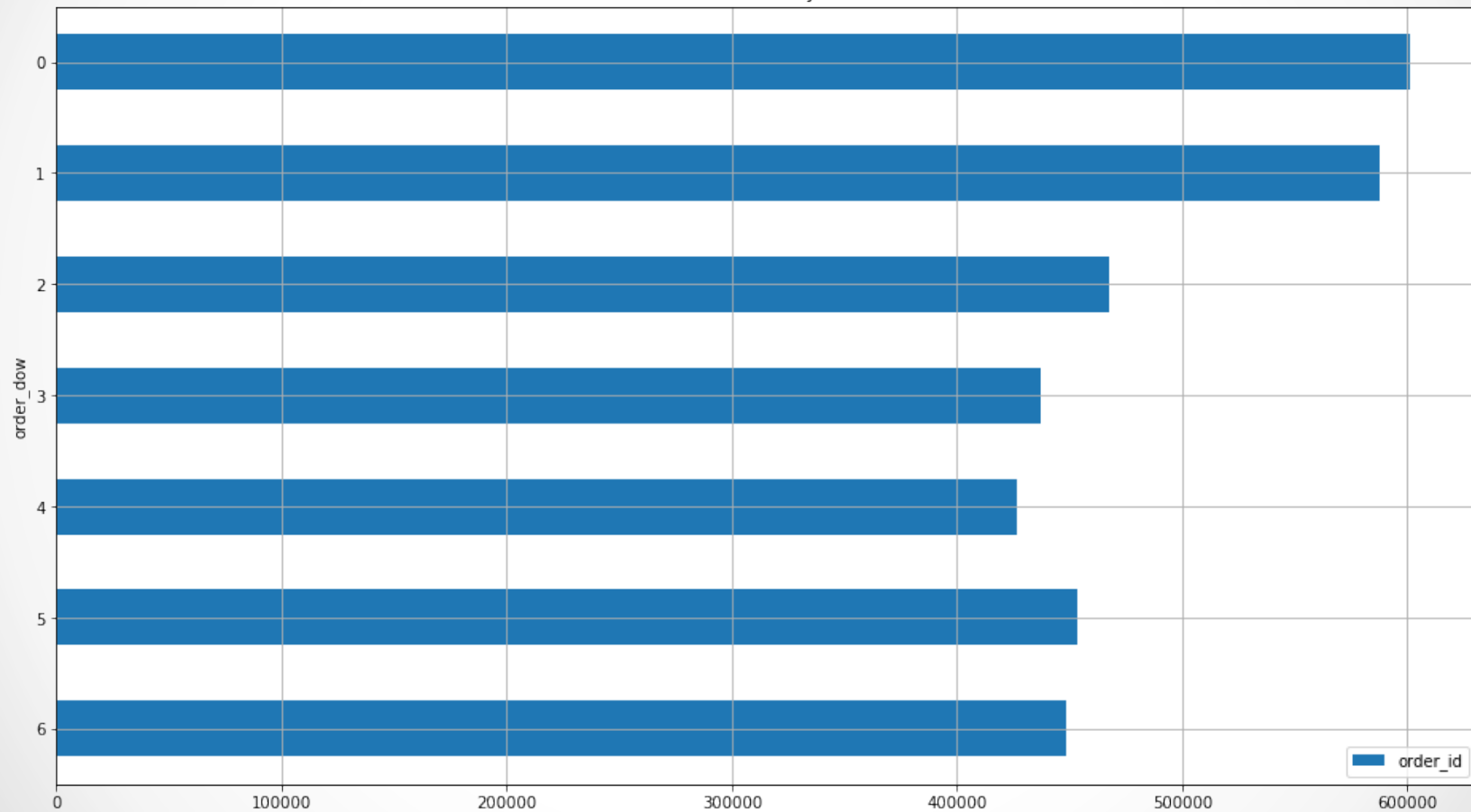
Understanding The Dataset

...

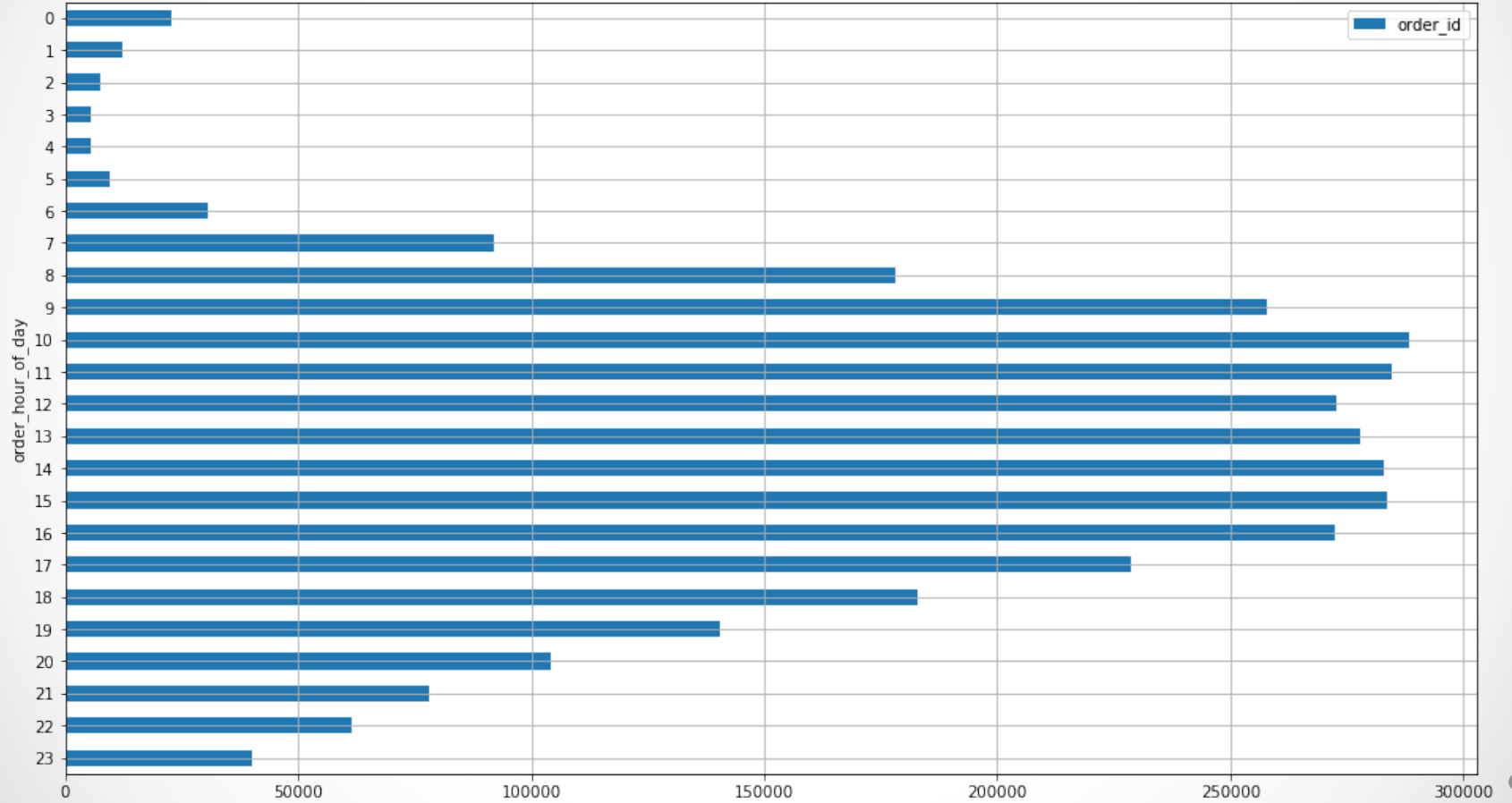
of Products per Department



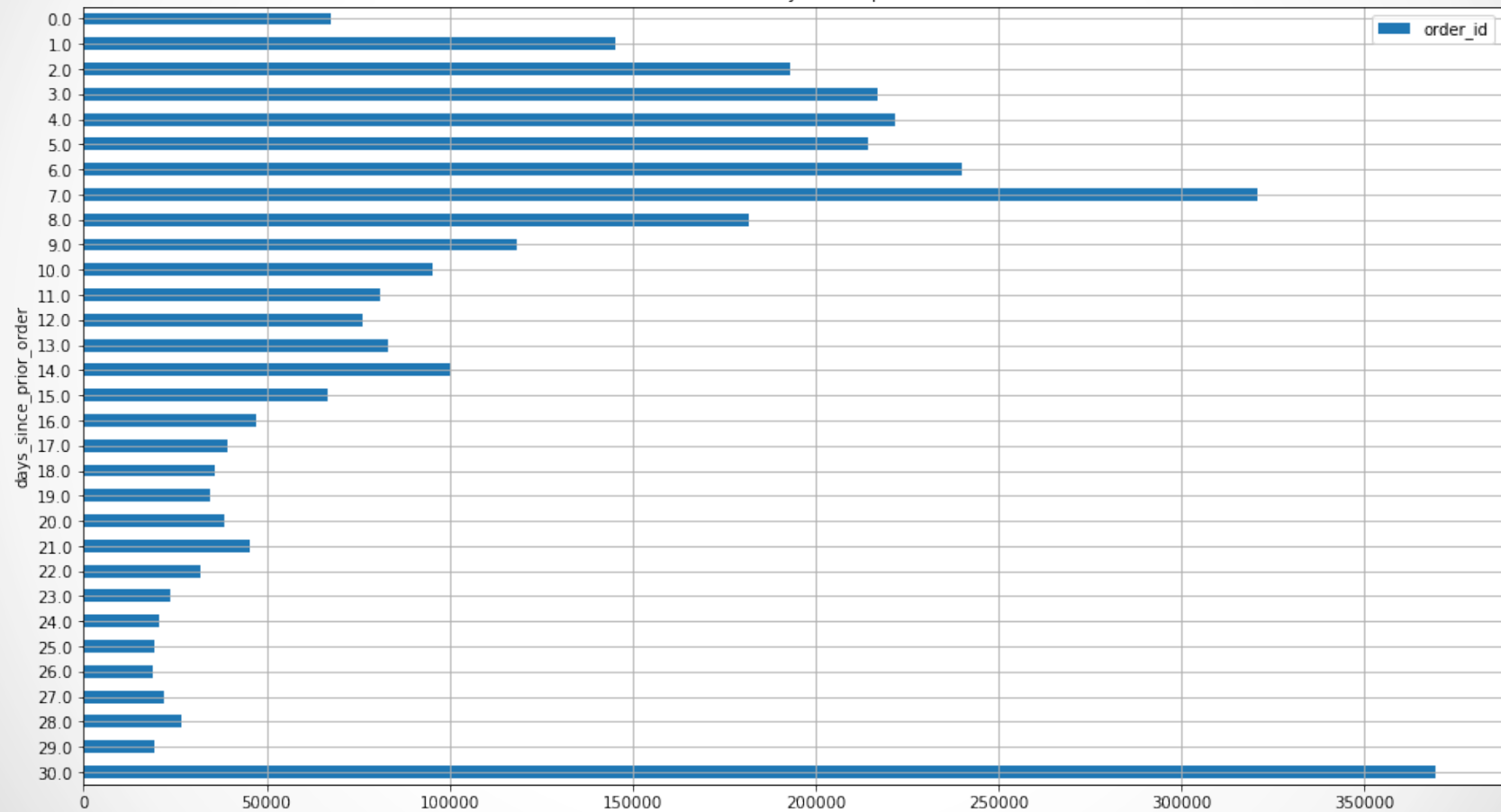
of Orders on Day of Week



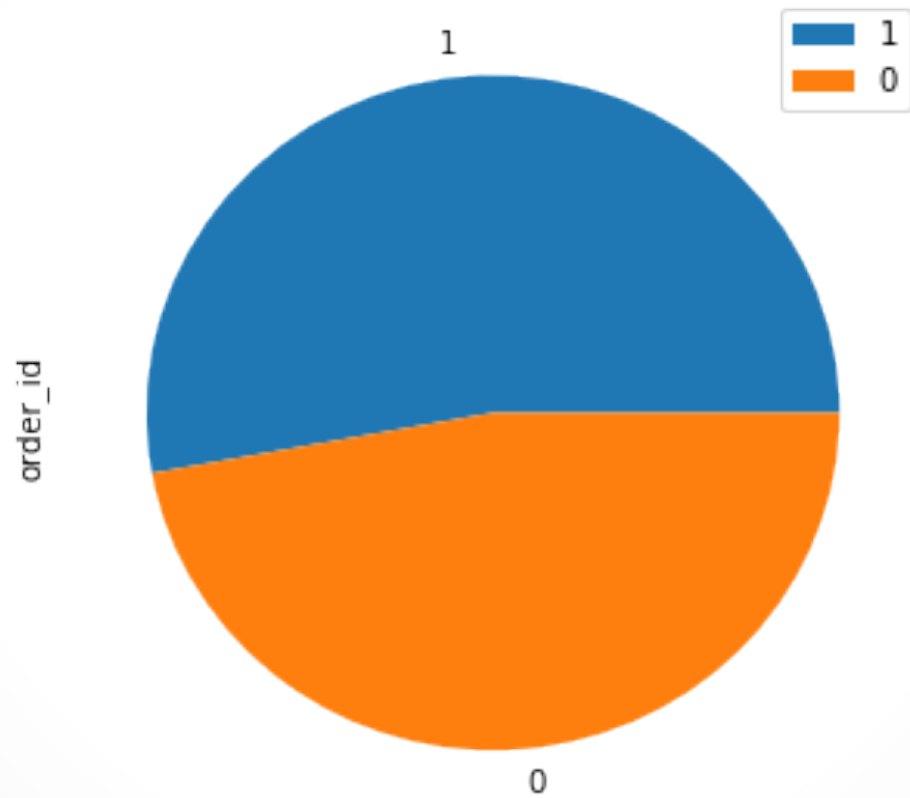
of Orders on Hour of the Day



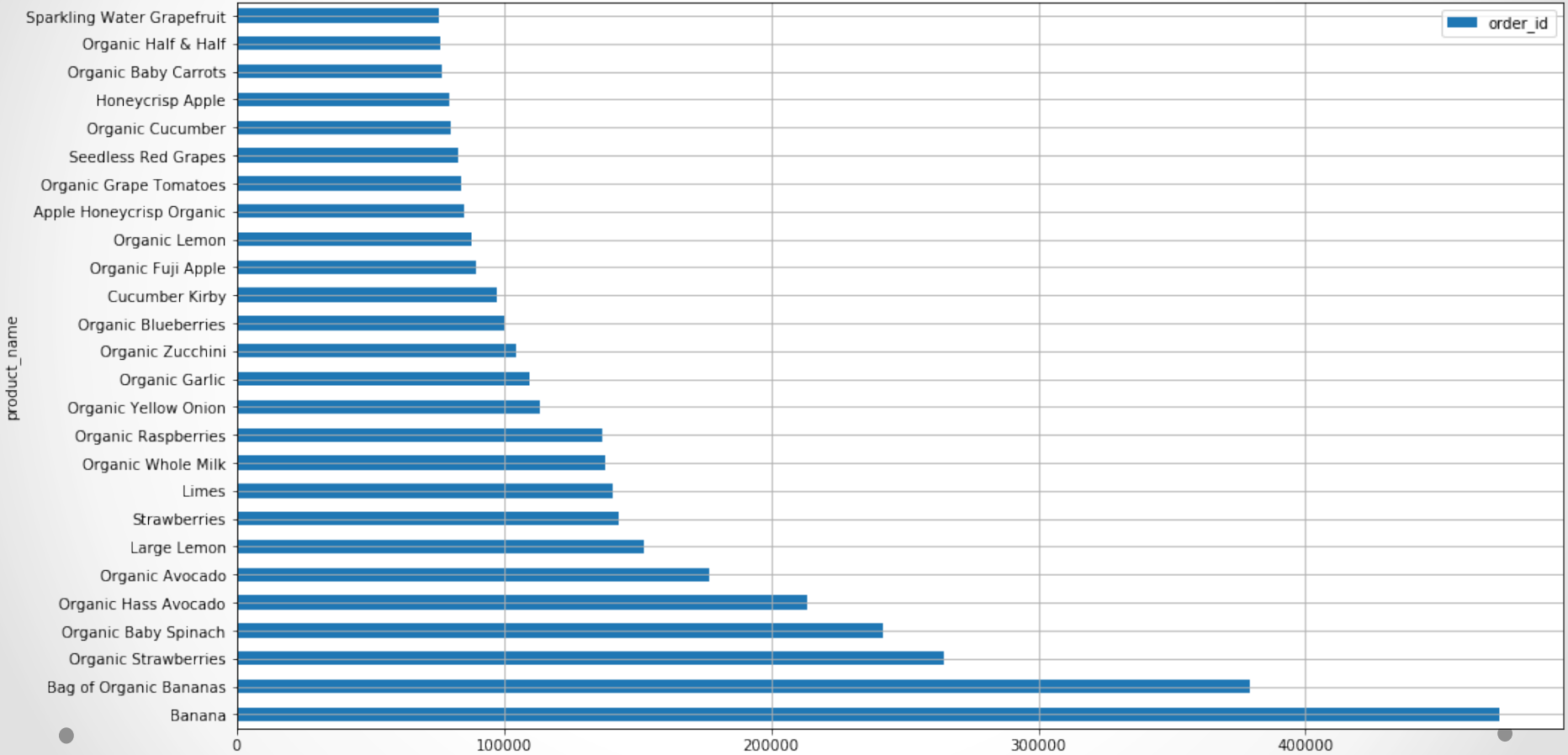
of Orders on Days since prior order



Distribution of Reordered



of Orders per Top 25 Product



Model

...

Choices

Random
Forest
Classifier

Gradient
Boosting
Classifier

Metrics

...

Cross Val Score Comparison

Random Forest

Random Forest :

Score: 0.852013953353

ROC AUC: 0.78906167567

Accuracy: 0.849758501971

Recall: 0.118334225001

Feature importance:	features	importance
12	up_total_orders	0.336196
13	up_max_add_to_cart_order	0.097800
2	p_reordered_rate	0.091143
14	order_number	0.088994
5	u_total_orders	0.087582
1	p_sum_reordered	0.069906
0	p_total_orders	0.048231
11	up_max_order_number	0.044468
10	up_max_order_id	0.038993
7	u_total_products	0.034003
6	u_avg_days_since_prior_order	0.018562
17	days_since_prior_order	0.014353
9	u_avg_cart	0.013312
3	aisle_id	0.008460
8	u_total_distinct_items	0.002738
16	order_hour_of_day	0.001784
4	department_id	0.001742
15	order_dow	0.001734

Gradient Boosting

Gradient Boosting :

Score: 0.997289572245

ROC AUC: 0.754068356894

Accuracy: 0.840725429577

Recall: 0.22290055707

Feature importance:	features	importance
10	up_max_order_id	0.123060
13	up_max_add_to_cart_order	0.108933
6	u_avg_days_since_prior_order	0.096258
9	u_avg_cart	0.078599
2	p_reordered_rate	0.076827
8	u_total_distinct_items	0.070864
7	u_total_products	0.061259
12	up_total_orders	0.050924
0	p_total_orders	0.049493
1	p_sum_reordered	0.047270
11	up_max_order_number	0.046568
17	days_since_prior_order	0.042338
16	order_hour_of_day	0.039318
3	aisle_id	0.029388
15	order_dow	0.028233
5	u_total_orders	0.022189
14	order_number	0.021393
4	department_id	0.007087

Confusion Matrix

Random Forest

n = 275776	Predicted NO	Predicted YES	
Actual NO	TN =	FP =	
Actual YES	FN =	TP =	

Accuracy (TP+TN) / total:	
Error Rate 1 - Accuracy:	
Recall or True Positive Rate TP / Actual YES:	
False Positive Rate FP / Actual NO:	
Specificity 1 - False Positive Rate:	
Precision TP / Predicted YES:	
Prevalence Actual YES / total:	

Gradient Boosting

n = 275776	Predicted NO	Predicted YES	
Actual NO	TN =	FP =	
Actual YES	FN =	TP =	

Accuracy (TP+TN) / total:	
Error Rate 1 - Accuracy:	
Recall or True Positive Rate TP / Actual YES:	
False Positive Rate FP / Actual NO:	
Specificity 1 - False Positive Rate:	
Precision TP / Predicted YES:	
Prevalence Actual YES / total:	

Next Steps

- Huge Potentials to enrich the features to improve model's performance.
- Building label could be improved.
- For validation of the model's in lieu of the test results to compare with, Prior data could be sampled for test set.

Thank You!

...

