



# **PREDICTING A USER'S NEXT INSTACART ORDER**

Elliot Wilens  
Metis Data Scientist

## Objective and Problem Setup



## Algorithms

- XGBoost
- RandomForest
- Grid Search & K-Fold Cross Validation
- Precision/Recall Curves
- F-1 & F-beta (F-2) scores

## Tools

- PostgreSQL
- Tableau
- Google Cloud Deep Learning VMs
- Jupyter Notebook

## Python Libraries

- Scikit-learn, StatsModels
- multiprocessing
- pickle
- pandas/numpy



“

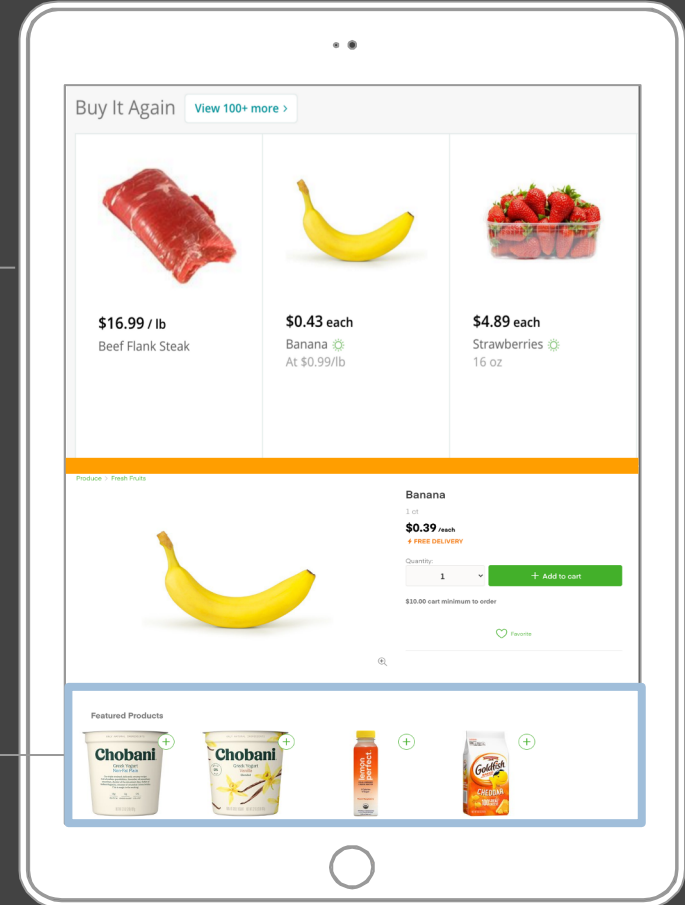
*Objective: To predict all of the reorders in a user's next cart*

*But first... why?*

## Use Cases:

1) Buy-it-again recommendations

2) Frequently bought with...



## Dataset & Feature Engineering





**33,819,106**

Total rows (1 per product per order)

**3,346,083**

Total orders

**337,418**

Total users

## THE DATASET

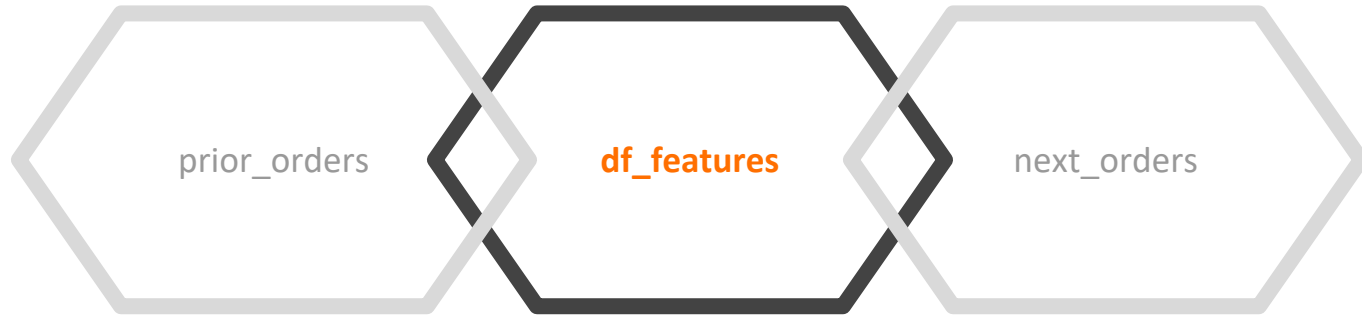


Contains prior order  
details for all users in  
next\_orders

Contains order details for  
each user's 'next' order



## THE DATASET



Contains prior order  
details for all users in  
next\_orders

Contains order details for  
each user's 'next' order

### **df\_features**

- Contains user & product statistics from prior\_orders
- Contains next\_order details
- Modeling done on these inputs

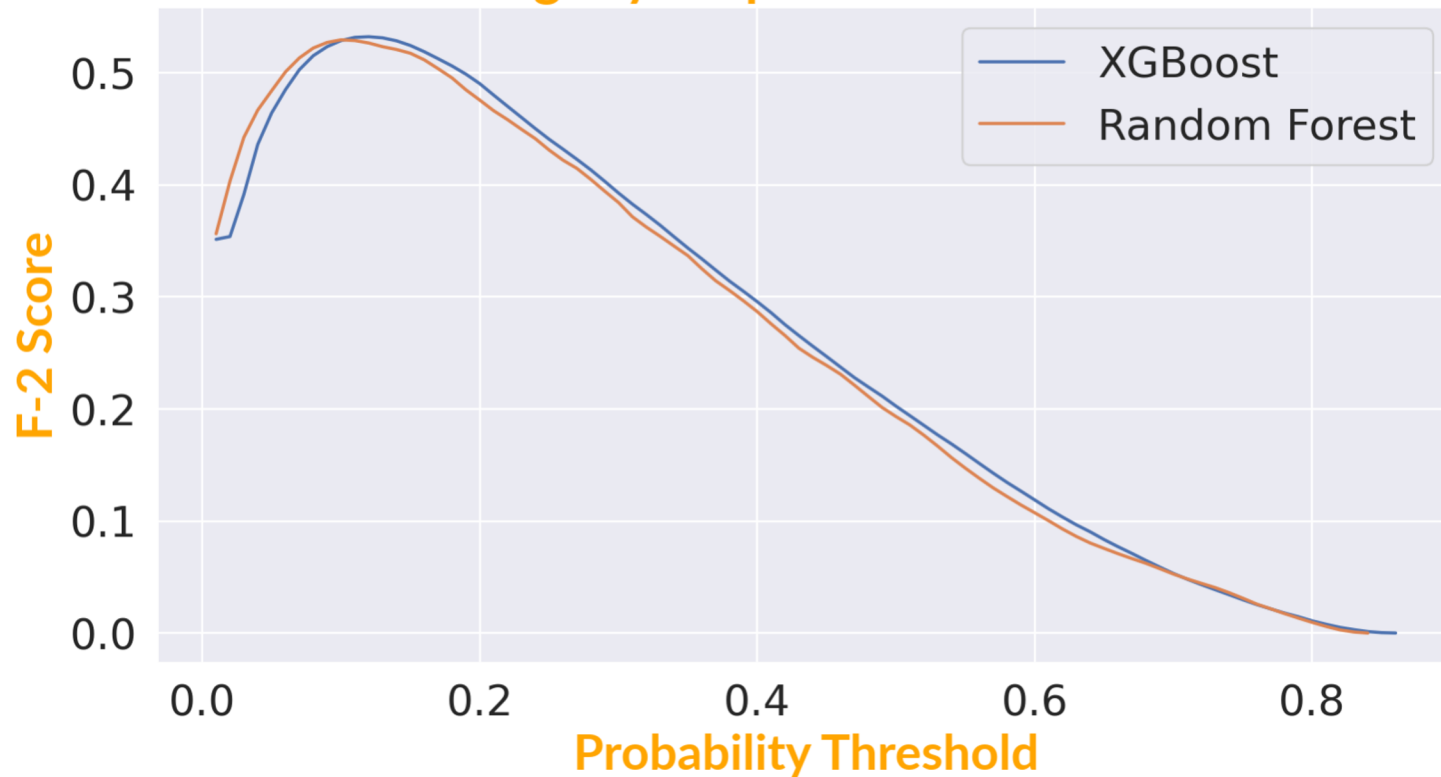
Some key features (32 total):

User Features	Product Features	User/Product Features
avg_cart_size	percent_reorders	<b>order_streak</b>
days_since_prior_order	qty_sold	last_five_buys
avg_time_between_orders	qty_reordered	ln_last_cart (0/1)

## Model Selection & Results



## XGBoost slightly outperformed RandomForest



Used grid search to define the optimal parameters:

learning_rate	<b>0.009</b>
n_estimators	<b>400</b>
max_depth	<b>7</b>
colsample_bytree	<b>0.8</b>
min_child_weight	<b>9</b>



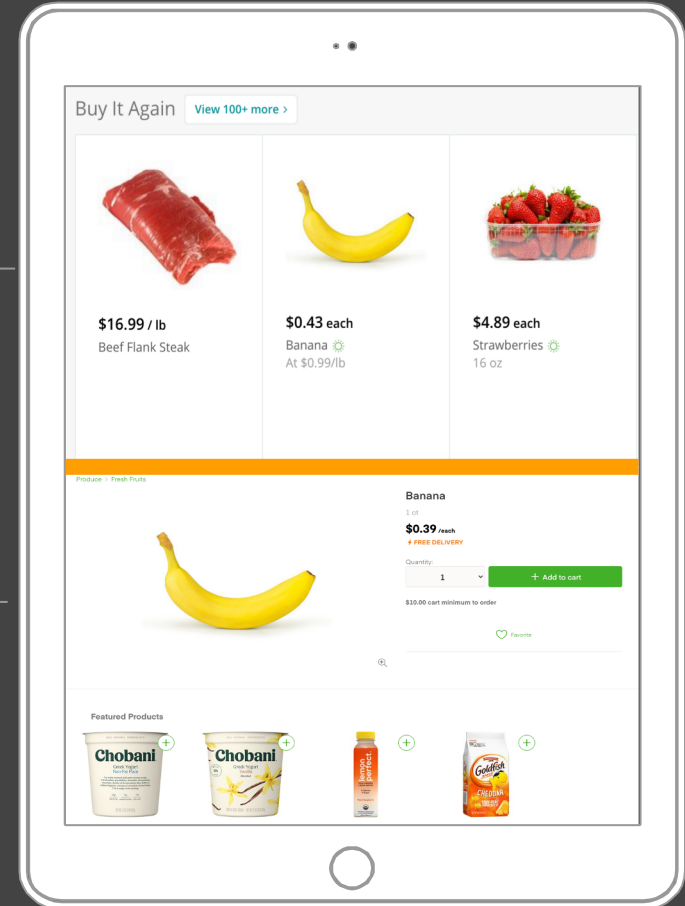
# Scoring

2018 Kaggle competition crowned winner based on resulting  
F-1 scores

*Do you remember our use case?*

## Use Cases:

- 1) Buy-it-again recommendations
- 2) Frequently bought with...



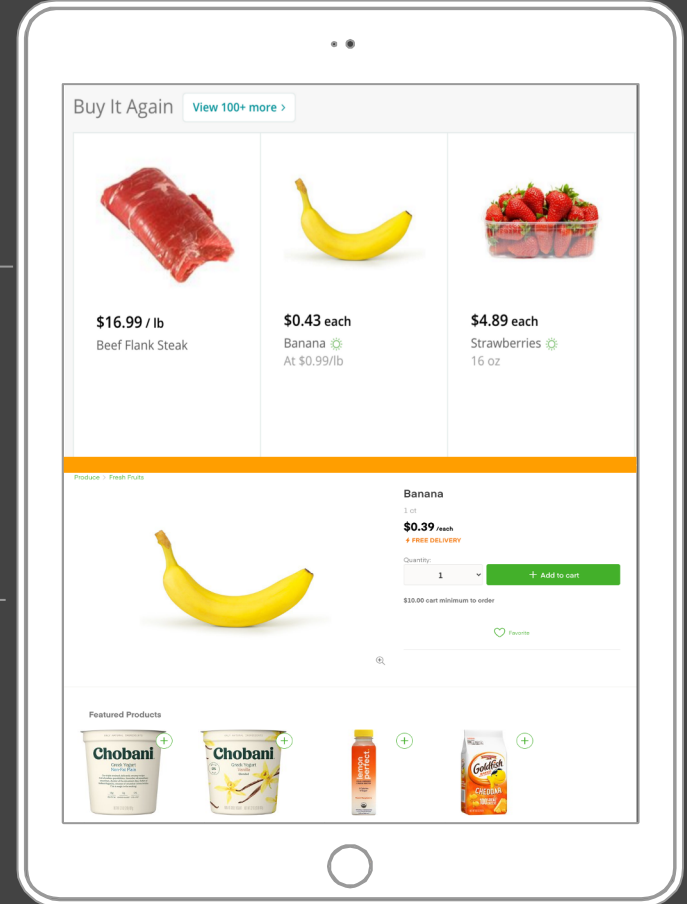
## Use Cases:

- 1) Buy-it-again recommendations
- 2) Frequently bought with...

*But how does this help Instacart?*

1. User ease of use

**1. Increase product conversion rates**





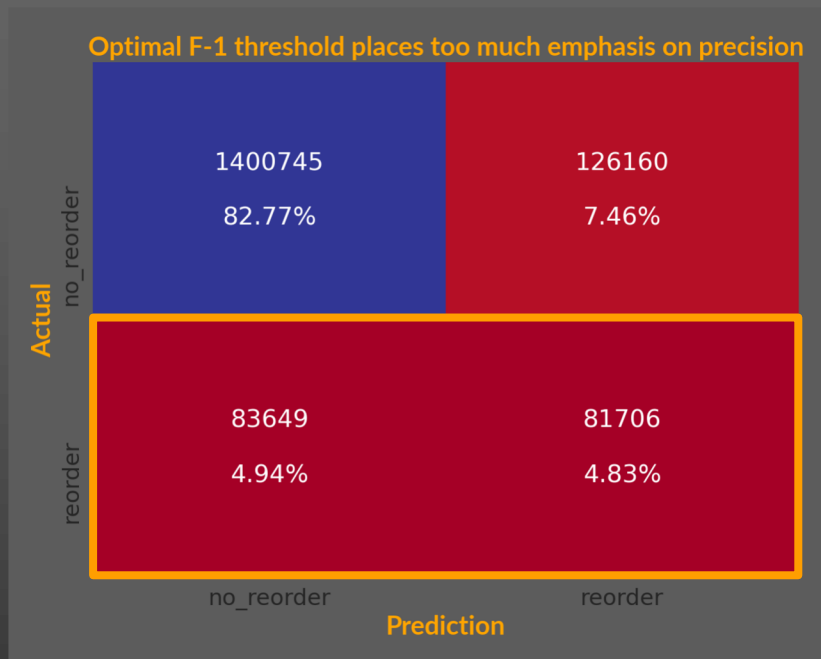
What's the risk of incorrectly classifying an input as positive?

Not Much.

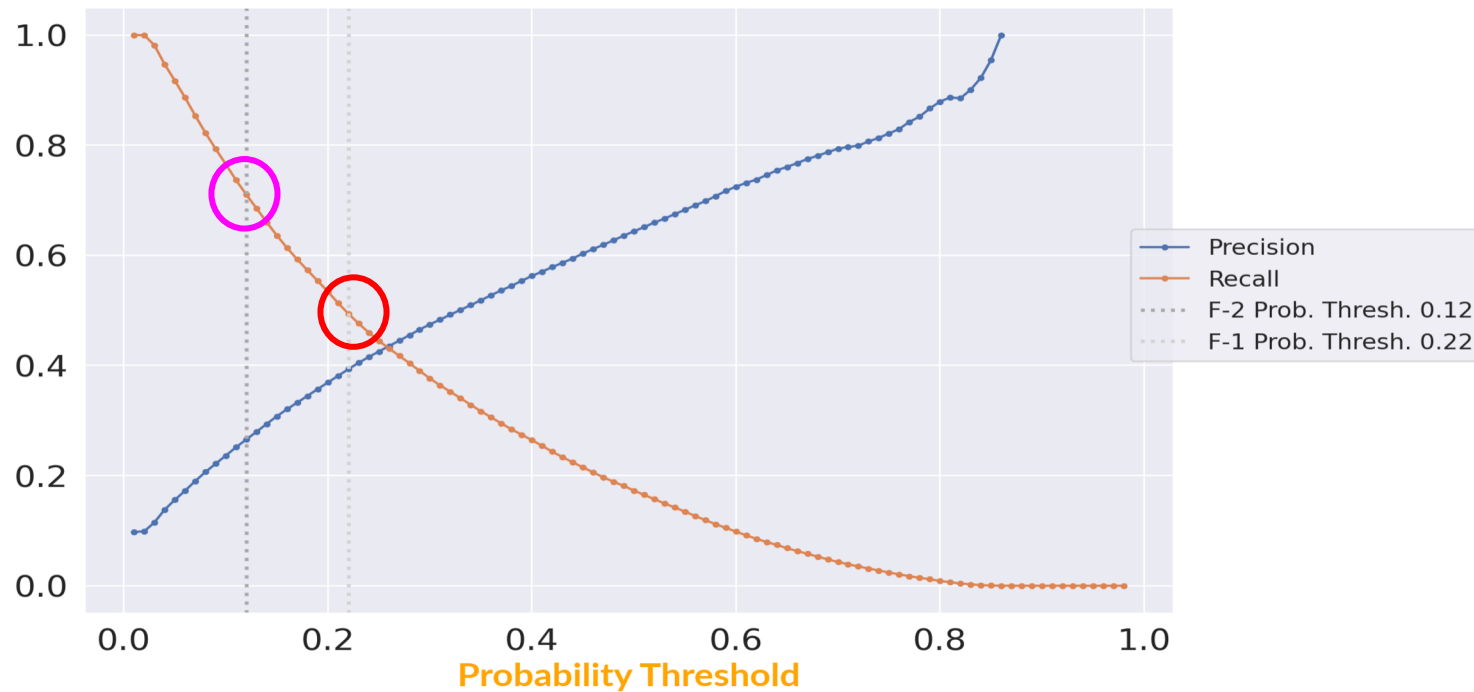
In fact, we may be better off including items that the user is *less likely to buy* based on their prior orders. This will help Instacart to increase conversion!

Therefore, we ought to prioritize **recall**!

$\text{Recall} = \text{True Positives} / (\text{True Positives} + \text{False Negatives})$



## Choosing a threshold with an ideal recall/precision balance using F-2 Scores

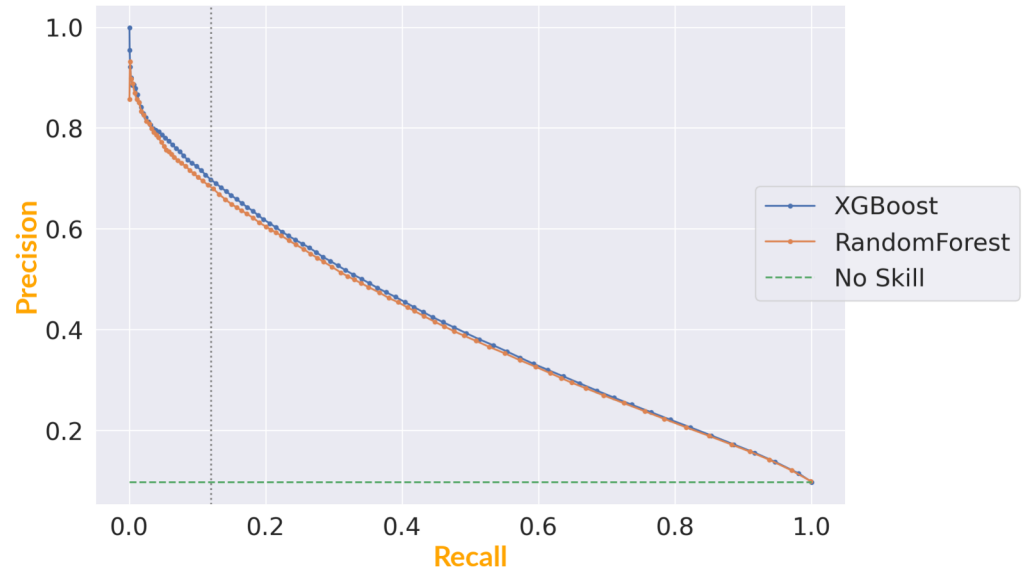


## MODEL RESULTS

### Adjusted F-2 Score

0.53

XGBoost slightly outperformed RF at our chosen recall threshold



Optimized F-2...much better!

Actual	Prediction	
	no_reorder	reorder
no_reorder	1201741 71.01%	325164 19.21%
reorder	47886 2.83%	117469 6.94%

### Probability Threshold

0.12

THANK  
YOU!

**Any questions?**

You can find me...

On LinkedIn: [linkedin.com/elliottwilens](https://www.linkedin.com/elliottwilens)

On GitHub: [github.com/edubu2](https://github.com/edubu2)

[wilensel@gmail.com](mailto:wilensel@gmail.com)