eCommerce Purchase Decisions

Alex Qaddourah, Alex McLaughlin, Teddy Li, Junji Wiener

AGENDA



Business Understanding

What are the underlying factors that drives an eCommerce decision.



Data Understanding & Limitations

Collection, Familiarity, Problems, Discovery, Planning



Data Preparation

- Raw data transformation
- Preprocessing
- Subsetting
- Feature Engineering
- Label Encoding



Modeling

What models did we choose?

- Decision Tree
- Random Forest
- Linear Regression



Evaluation

What do our classifiers tell us about the data?



Deployment

What can be used for model deployment?

1 Business Understanding

VIEW --- CART --- PURCHASE

Addressing the Business Problem

Data-Driven Question

What actions within, and circumstances surrounding, a user session on an eCommerce website can help identify the eventual purchase decision?

- Product View
- Cart/Wishlist
- Day or Time

Problem Definition

The average cart abandonment rate for online retailers is 67.91%. To optimize the checkout process, marketers & analysts must be familiar with actions and timelines of user sessions.

- Remarketing
- Sale/Promotion Timelines
- Price Analysis

Motivation & Potential Benefit

A data analytics solution can illuminate characteristics of consumer groups, opportunities to change behavior, and ultimately optimize conversion.

- Motivate (would-be) non-Purchasers
- Retain Purchasers

2

Data Understanding & Limitations

Collection - Familiarity - Problems - Discovery - Planning

Data Understanding



Collection

Where did the data set come from?



Familiarity

What did our team do to familiarize ourselves with the data?



Problems

What problems did we run into during this exploratory phase?

- 42 million rows
- Data integrity gap



Discovery

- User_session
- User_id
- Product id
- Brand
- Various Events



Planning

Is there opportunity to help answer our business problem?

Data Limitations

Limited features

Imbalanced

- Large (42M rows)
- Highly correlated features (view -> cart)

Data structure

3

Data Preparation

Making the data useable

Preprocessing

Null Value Handling

Subsetting

Feature Engineering

One-Hot Encoding

Multicollinearity

Scaling

Cleaned Data Structure



```
> str(phonesSat3)
'data.frame':
               1519759 obs. of 29 variables:
$ product_id.x : int 1000978 1000978 1000978 1000978 1000978 1000978
1000978 1000978 1000978 ...
 $ user_id
              : int 512366656 512387437 512452484 512493424 512494755 512531443
512557960 512574644 512601648 512657606 ...
$ DateTime4
            : int 6 3 14 7 11 17 4 14 14 16 ...
             : Factor w/ 39 levels "apple", "asus", ...: 30 30 30 30 30 30 30 30
 $ brand
30 30 ...
$ price
              : num 333 333 333 333 ...
$ user_session : Factor w/ 405112 levels "00000aaa-
d774-49bc-9c31-0c9f6e1c2f0a",..: 321765 219953 384716 389105 293360 113176 208247
56302 226158 277761 ...
$ colPurchase : int 0000000000...
 $ EarlyMorning : int 1100001000...
 $ Mornina
              : int 0001100000...
 $ Afternoon
              : int 0010010111...
 $ Evenina
              : int 0000000000...
 $ USA
              : int 0000000000...
 $ CHN
              : int 0000000000...
 $ SK
              : int 111111111...
 $ JPN
              : int 0000000000...
 $ UK
              : int 0000000000...
 $ NED
              : int 0000000000...
 $ FIN
 $ SPN
$ RUS
 $ CYP
                  0000000000...
 $ TAI
 $ CAN
                   0000000000...
 $ IND
 $ MAL
 $ BRZ
              : int 0000000000...
 $ product_id.y : int 11 11 6 41 17 2 38 13 38 12 ...
$ countforviews: int 121131111...
 $ countforcarts: int 0000000000...
```

4 Modeling

Classifiers

- Logistic Regression
- Random Forest
- Neural Network
- XGBoost

Performance Metrics

Confusion Matrix

- Accuracy
- TPR
- FPR
- F1 score
- Precision
- ROC AUC

Logistic Regression

- Score 56%
- Accuracy 46%
- TPR 52%
- FPR 44%
- F1 score 54%
- Precision 57%
- ROC AUC 56%

Random Forest

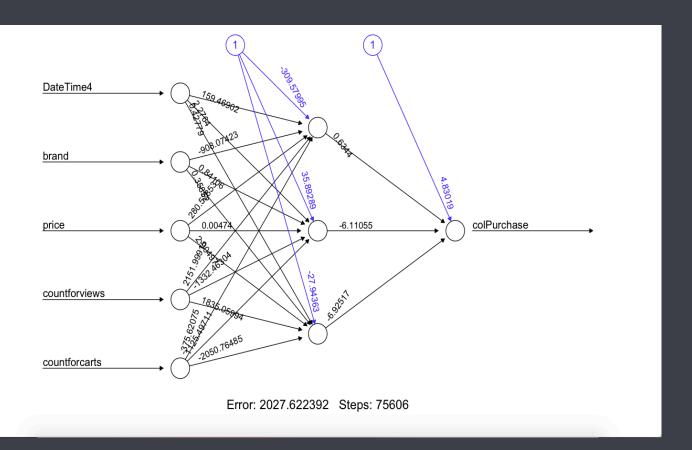
RandomForestClassifier(n_estimators = 100, verbose=3, n_jobs=-1, max_depth=20, random_state=42)

- Score 64%
- Accuracy 38%
- TPR 53%
- FPR 62%
- F1 score 60%
- Precision 66%
- ROC AUC 70%

Neural Network

- Accuracy 91%
- TPR 0.12%
- FPR 0.04%
- F1 score 0.23%
- Precision 21%

Neural Network



XGBoost

- Score 57%
- Accuracy 40%
- TPR 57%
- FPR 41%
- F1 score 58%
- Precision 59%
- ROC AUC 62%

5 Evaluation

What did our classifiers find?

Confusion Matrix

Random Forest Classifier

	1	0
1	148497	128850
0	69668	207211

Logistic Regression (upsampled)

5	1	0
1	144237	133110
0	109328	167551

Neural Network

0 1

0 1384991 572

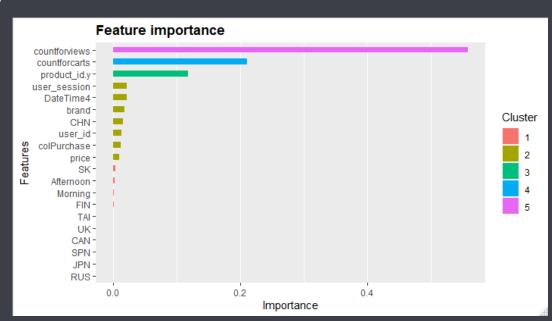
1 134040 156

XGBoost (upsampled)

	1	0
1	164157	113190
0	123354	153525

Feature Ranking

- Product View Counts (total)
- Cart View Counts (total)
- Unique Products Viewed
- User Sessions (total)
- Hour of Day (24hr)
- Brand
 - Samsung
 - Xiaomi
 - Apple
 - Huawei
 - Oppo
 - Meizu
 - Nokia
 - Vivo
 - TP-Link



Interpreting Models

- Which model had the highest Accuracy?
 - NeuralNetwork(unbalanced)

- Which model had the best TPR?
 - XGBoost

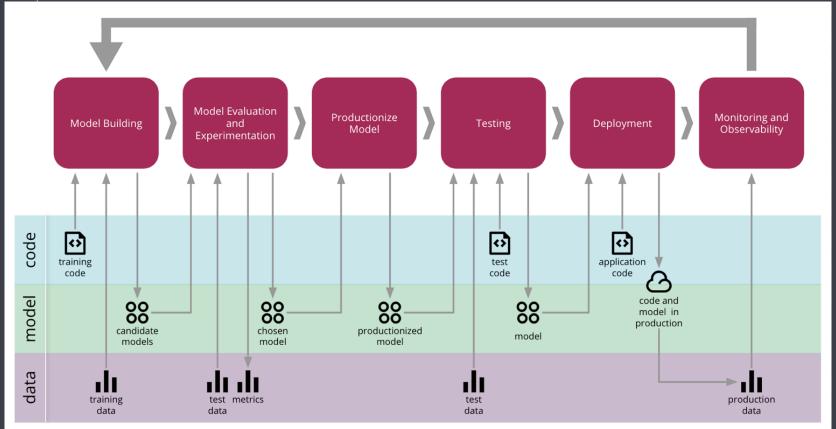
- Which model had the best F1 Score?
 - RandomForest

- Which model had the highest Precision?
 - RandomForest

6 Deployment

Continuous Delivery for Machine Learning (CD4ML)

End to End Continuous Delivery for Machine Learning



Thanks!

ANY QUESTIONS?

GitHub's:

- 8Jun
- alexqaddourah
- almc6742
- TeddyCU

Citations

- Continuous Delivery for Machine Learning
 - https://martinfowler.com/articles/cd4ml.html#initial-ml-process.png
 - Danilo Sato, Arif Wider, Christoph Windheuser
- Kaggle https://www.kaggle.com/mkechinov/ecommerce-behavior-data-from-multi-category-store