PREDICTION TASK

What is the type of task?

- Classification | Random Forest Classifier (Predict whether a customer will buy Product 6850 or not).

Which entity are predictions made on?

- Predictions are made on individual customers based on their past purchasing behavior.

What are the possible outcomes to predict?

- 1 (Yes) The customer will purchase Product 6850.
- 0 (No) The customer will not purchase Product 6850.

When are outcomes observed?

After each transaction period, particularly in January 2019, based on customer behavior in the previous three months (October-December 2018)

DECISIONS

How are predictions turned into actionable recommendations or decisions for the end-user? (Mention parameters of the process / application for this.)

- Predicted buyers can be targeted with personalized marketing campaigns (e.g., promotions, discounts).
- If a customer is predicted not likely to buy, additional strategies like bundling or cross-selling with other gadgets can be implemented.

Parameters

- Customer demographics (age, gender, location)
- Past purchases (Gadget category spending, frequency of transactions)
- Weather conditions influencing sales behavior

VALUE PROPOSITION

Who is the end beneficiary, and what specific pain points are addressed?

- Retail company's marketing team -> Can better target potential buyers, reducing wasted marketing efforts.
- **Retail customers** → Receive personalized product recommendations based on their
- Business management -> Increases sales and revenue for underperforming products like Product 6850.

How will the ML solution integrate with their workflow, and through which user interfaces?

- Predictions will be fed into the company's marketing platform to create targeted ad
- Customer recommendations can be integrated into email marketing, app notifications, or personalized web banners.

DATA COLLECTION

How is the initial set of entities and outcomes sourced (e.g., database extracts, API pulls, manual labeling)?

- Data is sourced from the company's retail transaction database (df_2018_demo_view table in Spark SQL).

What strategies are in place to update data continuously while controlling cost and maintaining freshness?

- Automated data ingestion pipelines that pull new transactions every month.
- Periodic model retraining (every 3–6 months) to ensure predictive accuracy.

DATA SOURCES

Where can we get data on entities and observed outcomes? (Mention internal and external database tables or API methods.)

Internal sources:

- Retail transaction database (df_2018_demo_view)
- Customer demographics (age, gender, membership status)
- Purchase history (last 3 months)
- Product category spending
- Weather conditions on transaction days

External sources:

- Weather APIs (to get weather conditions per store location).
- Marketing data (response rates to previous promotions).

IMPACT SIMULATION

What are the cost/gain values for (in)correct decisions?

True Positive (TP): Predicting a customer will buy and they actually do → Increased revenue, better targeted marketing.

False Positive (FP): Predicting a customer will buy but they don't → Wasted marketing

True Negative (TN): Predicting a customer will not buy and they don't \rightarrow No unnecessary marketing spend.

False Negative (FN): Predicting a customer won't buy but they actually do \rightarrow Lost revenue opportunity.

Which data is used to simulate pre-deployment impact?

Historical data from October-December 2018 is used to train the model and test predictions for January 2019.

MAKING PREDICTIONS



Are predictions made in batch or

Predictions will be made in batch mode at the start of each month

How frequently?

Every 3 months (quarterly model retraining).

How much time is available for this (including featurization and decisions)?

1-2 hours per batch run, depending on data size and model complexity.

Which computational resources are used?

Google Colab / Spark clusters for training

Retail company's data warehouse for deployment.

BUILDING MODELS



How many models are needed in production?

One main model (binary classification model for purchase prediction).

When should they be updated?

Every 3-6 months, depending on how quickly customer behavior changes.

How much time is available for this (including featurization and analysis)?

1-2 days for retraining and validation.

Which computation resources are used?

Google Colab, Spark

FEATURES

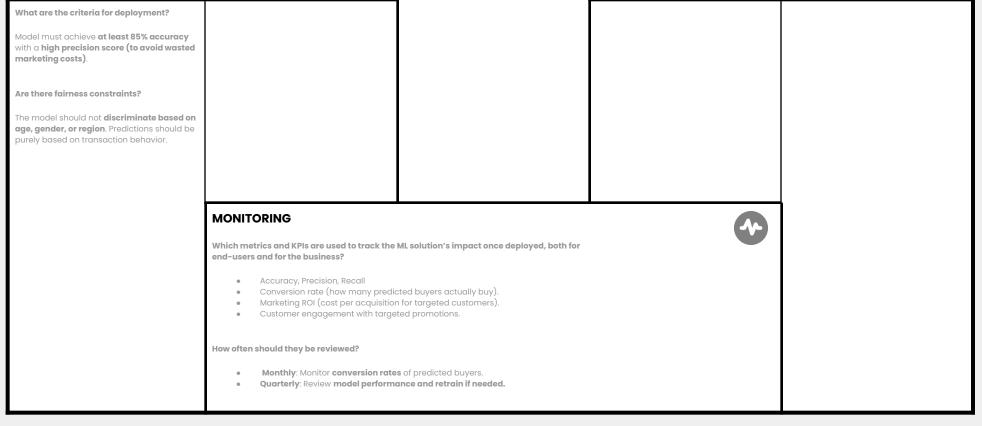


What representations are used for entities at prediction time?

- Customer demographic data (age, gender, region, membership type).
- Spending behavior in different product categories (last 3 months).
- Transaction frequency (number of purchases per category).
- Weather conditions on past purchase dates.

What aggregations or transformations are applied to raw data sources?

- Sum, mean, min, max transactions per category.
- Past 3-month transaction history per customer.
- Categorical encoding (One-Hot Encoding for membership type, region)
- Scaling (Min-Max Normalization for numerical features).











Version 1.2. Created by Louis Dorard, Ph.D. Licensed under a Creative Commons Attribution-ShareAlike 4.0 International License. Please keep this mention and the link to ownml.co when sharing.

