Identifying potential repeat buyers

THE TEAM | Stakeholders Meeting



REPEAT BUYERS MATTER

Stable revenue

Higher lifetime value (LTV)

Long-term growth

Loyal buyers refer others and boost organic growth

Better ROI for Olist

Repeat buyers already trust Olist, so they need less advertising to buy again





Datasets with duplicates: geolocation

Datasets with null values: reviews, products, orders

cleaning **Key changes:** Issues Removed invalid Lat/Lng outside Brazil coordinates Geolocation noisy zip-level spread Dropped zip prefixes with large spatial variance Filtered to 2016–2018 **Orders** Orders skewed right Imputed missing delivery

dates



Datasets with duplicates: geolocation

Datasets with null values: reviews, products, orders

cleaning **Key changes:** Issues • Filtered to 2016 - 2018 Capped outliers at 99.9 **Price and freight outliers** Order_items percentile

payments

different payment preferences

- Changed not_defined /null to others payment type
- Capped payment_installments at 99 percentile
- Merged multi-method payments per order_id to simplify and preserve totals

FEATURE ENGINEERING OVERVIEW

Four feature buckets:



- Delivery cost
- Discounts
- Total spend
- Distance to seller
- Type of payment

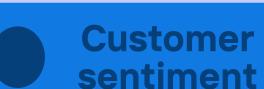


- Photos
- Descriptions
- Product size/weight



Perceived value

- Free shipping
- Good deals



• Review scores







GEOGRAPHICAL PROXIMITY

Seller Performance

Trusted sellers = More repeat buyers

Buyer-Seller Distance

Nearby sellers = stronger buyer loyalty

High Customer Area

Urban areas = better logistics, more repeat buyers

Delivery Speed vs. Expectation

Faster than expected = trust boost

Payment to **Delivery Lag**

Shorter wait = happier customer

Shipping Cost

Lower cost leads to higher satisfaction

Discount Usage

Promotions attract deal-seekers



Photos

Represents the number of images

More images → Higher buyer trust & confidence.

Product size/weight

Important for predicting freight value and order fulfillment success.

Larger or heavier items → higher shipping costs.

Product descriptions

Longer descriptions often reflect product clarity.

Reduces uncertainty → Increased satisfaction & reduced return rates.







REVIEW SENTIMENT & VERIFICATION

Behind the Scenes

Sentiment analysis of comments

comparison with review score

review score verification

FEATURE DASHBOARD

Focus

Drivers of repeat-buyer behaviour.

Data Shown

Key metrics like reviews, repeat rates, spend, and delivery.

Why is it useful?

Highlights where to improve listings and operations for loyalty.



MACHINE LEARNING PIPELINE

Ingestion

review_features.parquet

product_features.parquet

distance-seller-stats.parquet

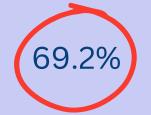
transaction.parquet

Scale numerical features and encode categorical features Model training Model_inputs. parquet data Stratified train-test split **Train/test sets**

MACHINE LEARNING PIPELINE

Model selection *Accuracy

Logistic regression



XGBOOST

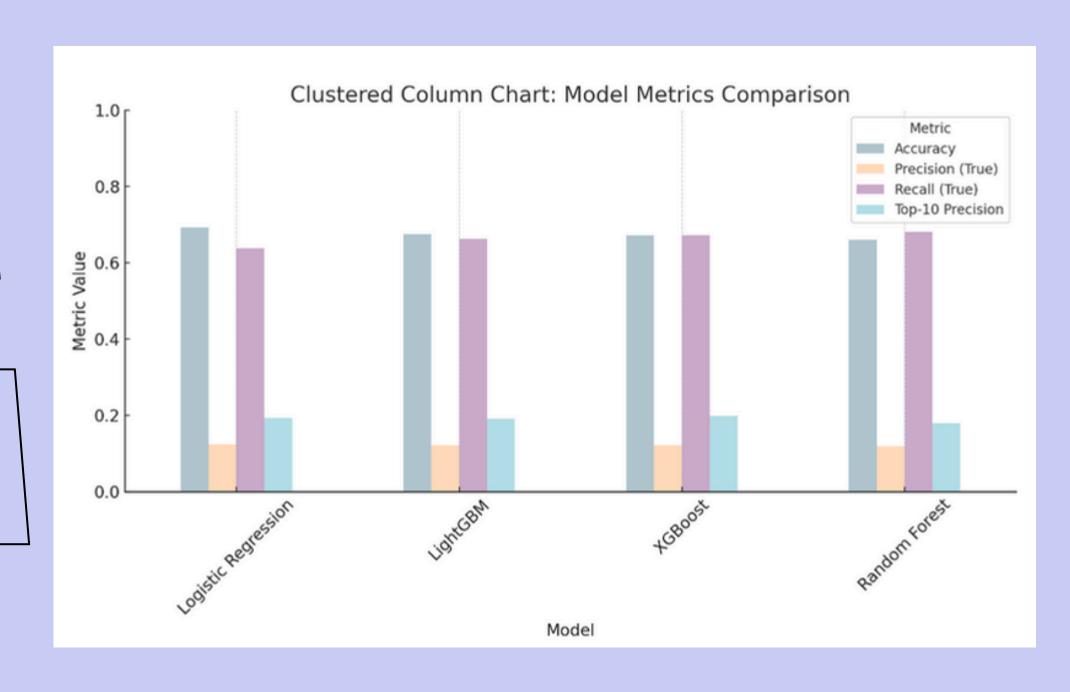
67.1%

LightGBM

67.5%

Random forest

66%



MACHINE LEARNING Model selection *Weighted avg Precision

Logistic regression

91.3%

XGBOOST

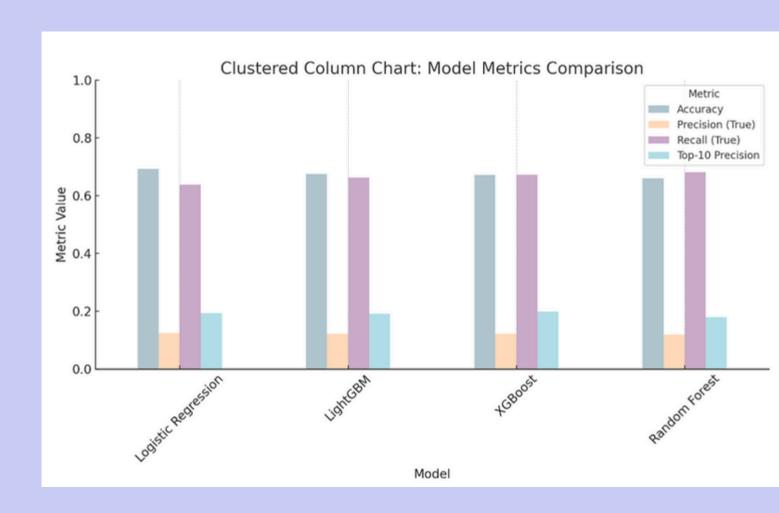
91.5%

LightGBM

91.4%

Random forest

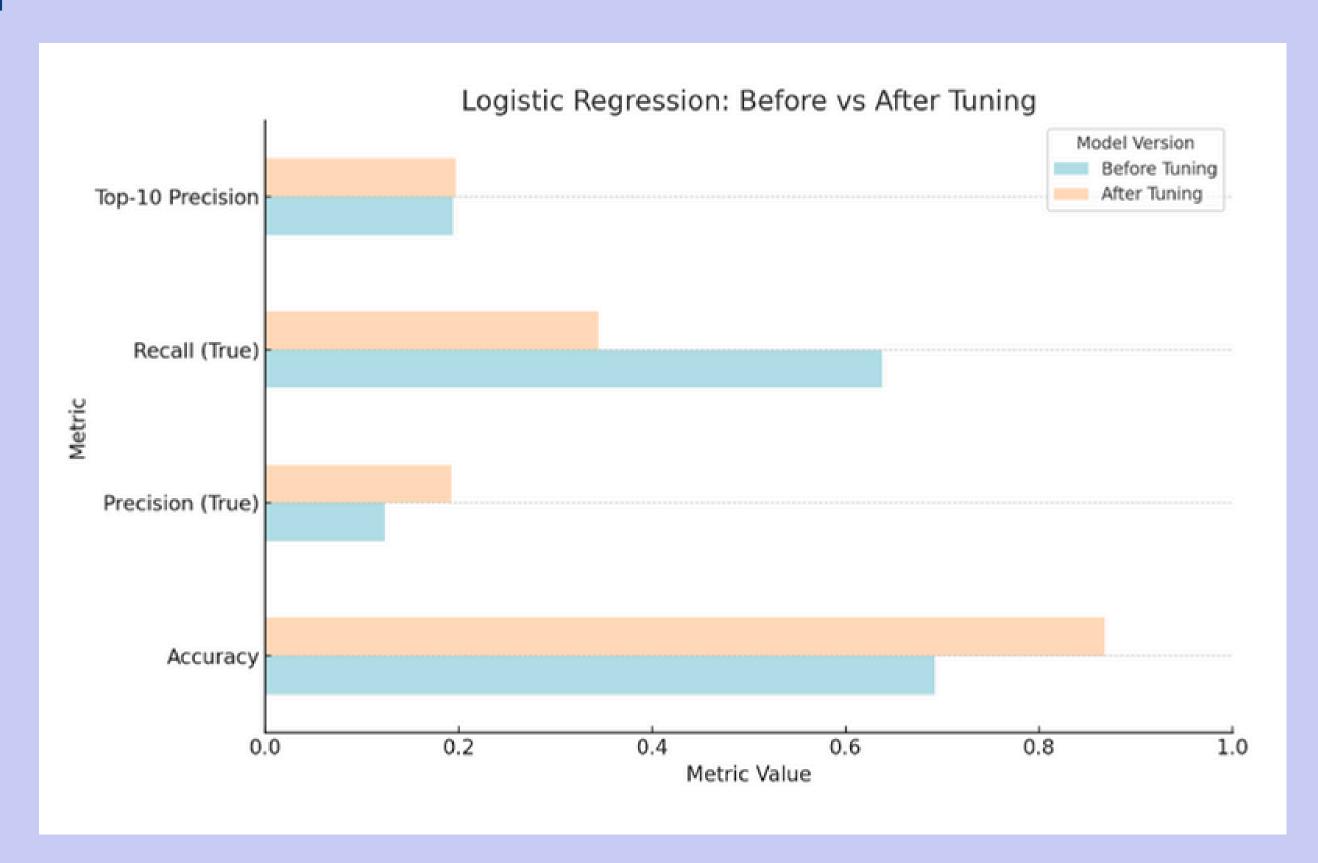
91.5%



KEYMODEL INSIGHTS

Fine-tuned Logistic Regression

Metrics used

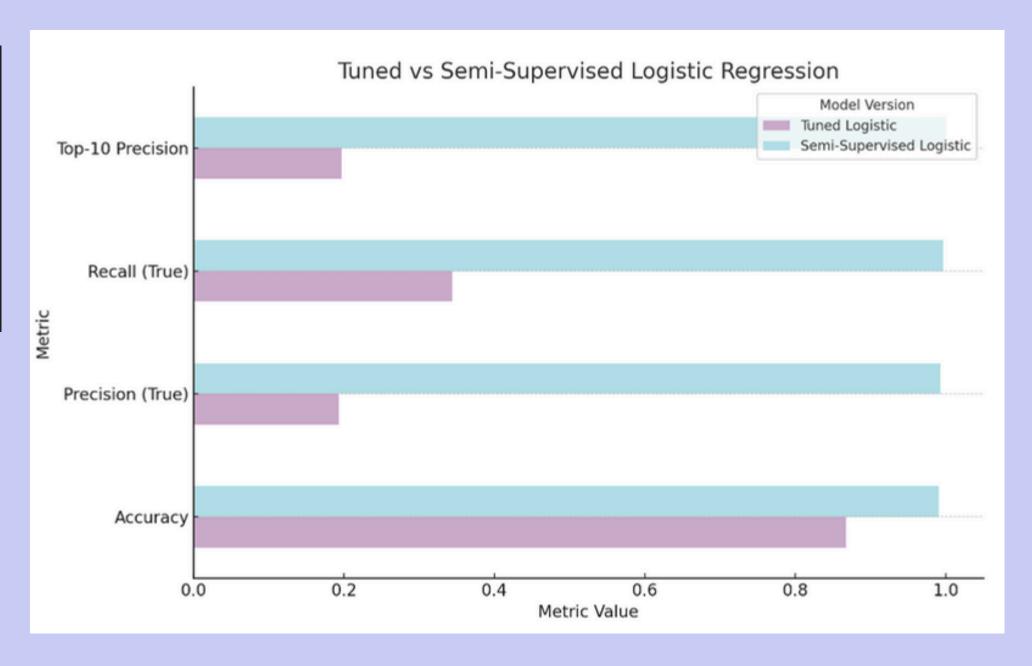


KEY MODEL INSIGHTS

Metrics used

```
mask = (
    (df["review_score"] > 3)
    | (df["deli_duration_exp"] <= -7)
    | (df["voucher"] >= 0.3)
    | (df["total_spent"] >= df["total_spent"].quantile(0.8))
    | (df["product_category_name"].isin(top_categories))
)
return pd.DataFrame({"weak_positive": mask})
```

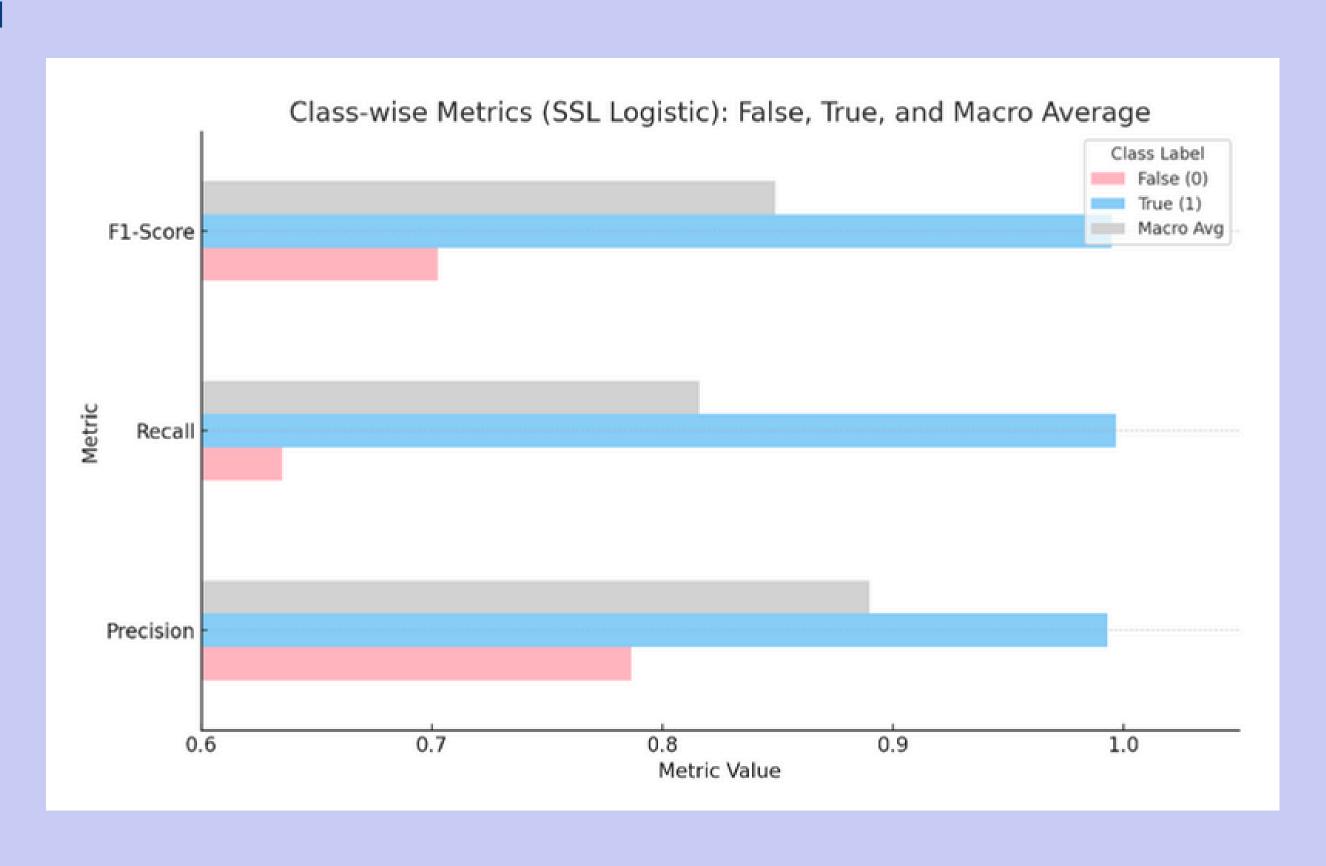
Semi-supervised Fine-tuned Logistic Regression

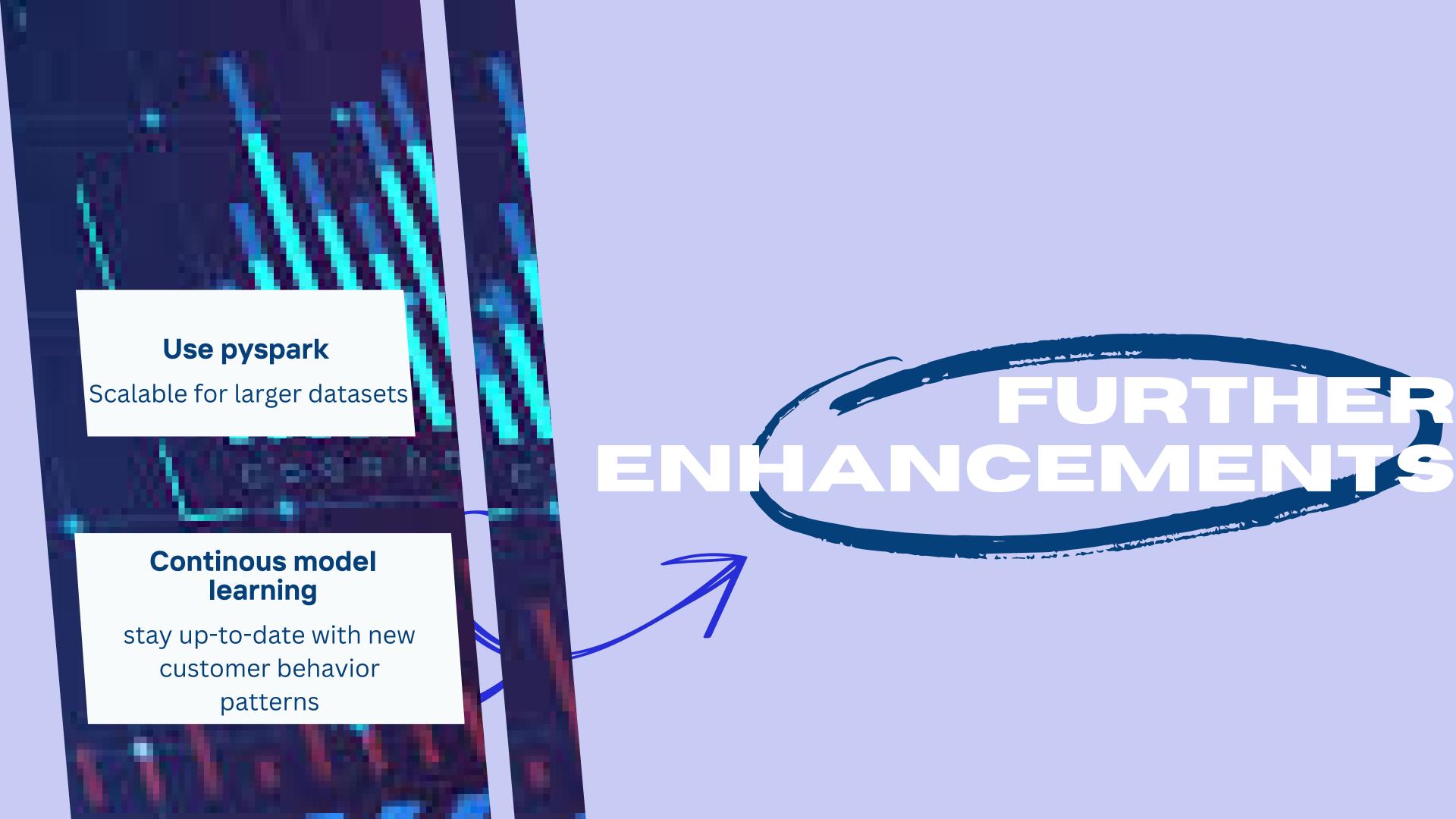


KEYMODEL INSIGHTS

Semi-supervised Fine-tuned Logistic Regression

Metrics used





THANK YOU &

