

Shopper Behavior Prediction

Analyzing User Behavior to Predict Revenue Generation

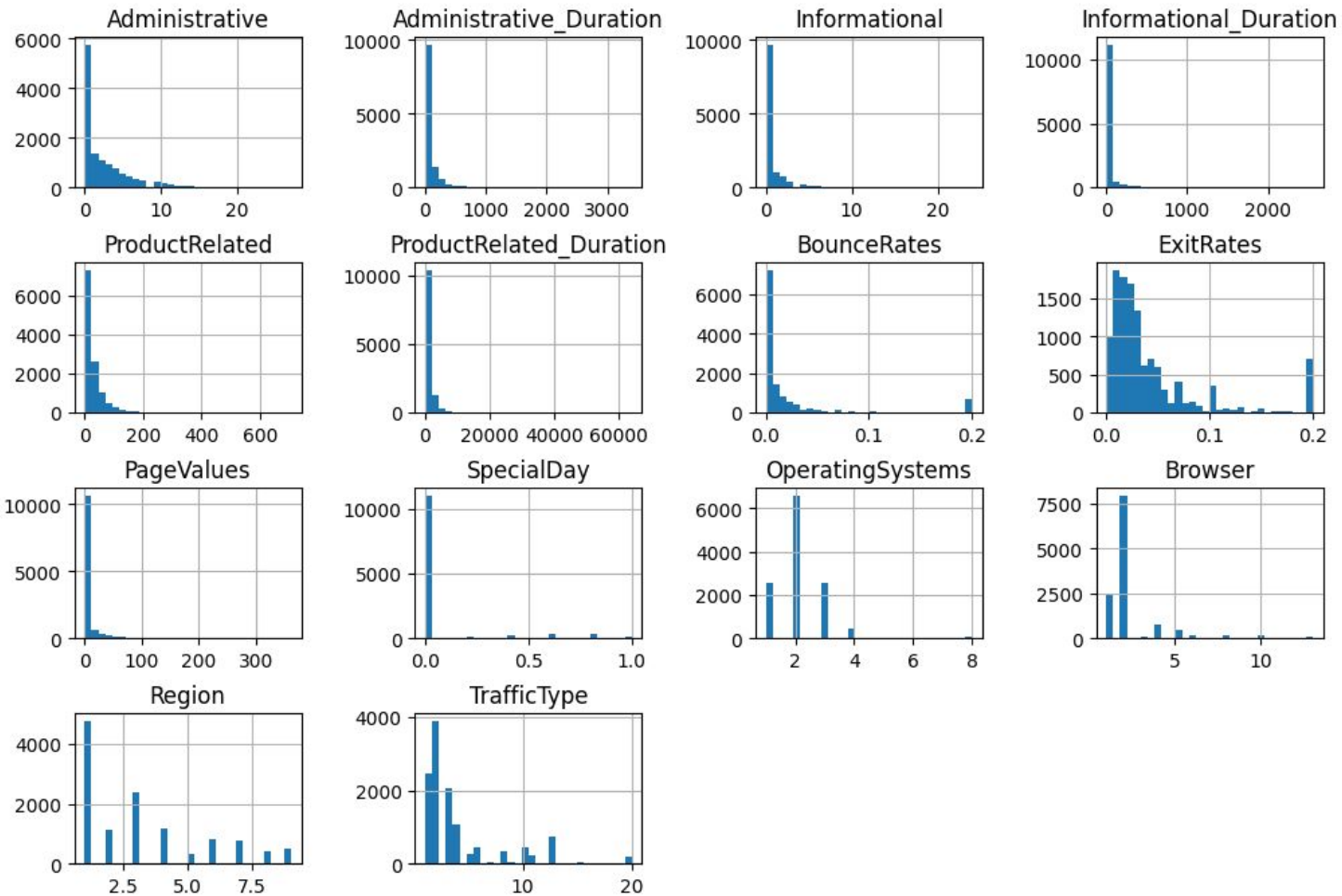


By: Matthew Neba

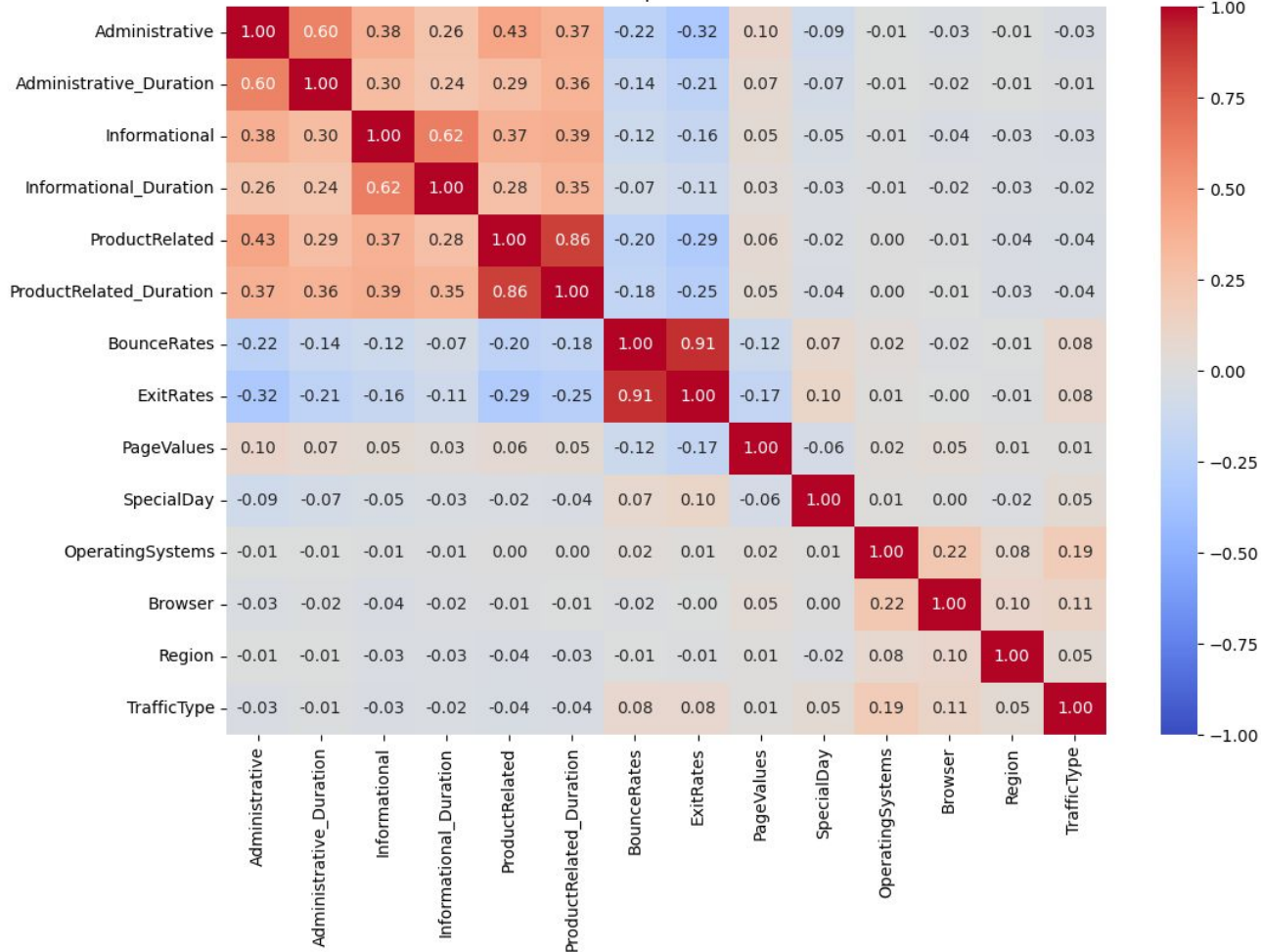
Introduction

- **Dataset:** UCI Online Shopper Intention Dataset
- **Objective:** Analyze user behavior to predict revenue generation
- **Data Composition:**
 - **Samples:** Online shopping sessions
 - **Features:** 17 variables including session duration, visited pages, bounce rates, and visitor type
 - **Response Variable:** Revenue (True/False)

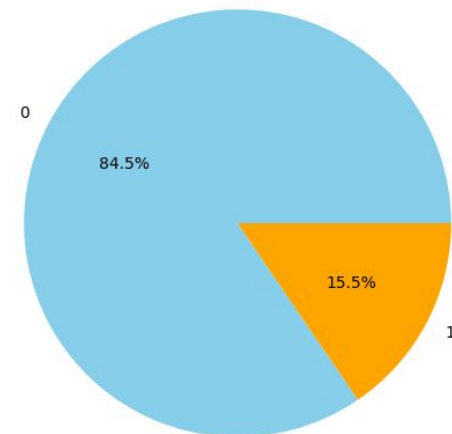
Distribution of Numerical Features



Correlation Heatmap of Numerical Features



Proportion of Revenue Classes



Feature Preprocessing:

- **Numerical Features:** Imputed missing values using `SimpleImputer(strategy="mean")`, standardized using `StandardScaler()`
- **Categorical Features:** Imputed missing values using `SimpleImputer(strategy='constant', fill_value='missing')`, one-hot encoded using `OneHotEncoder(handle_unknown='ignore')`

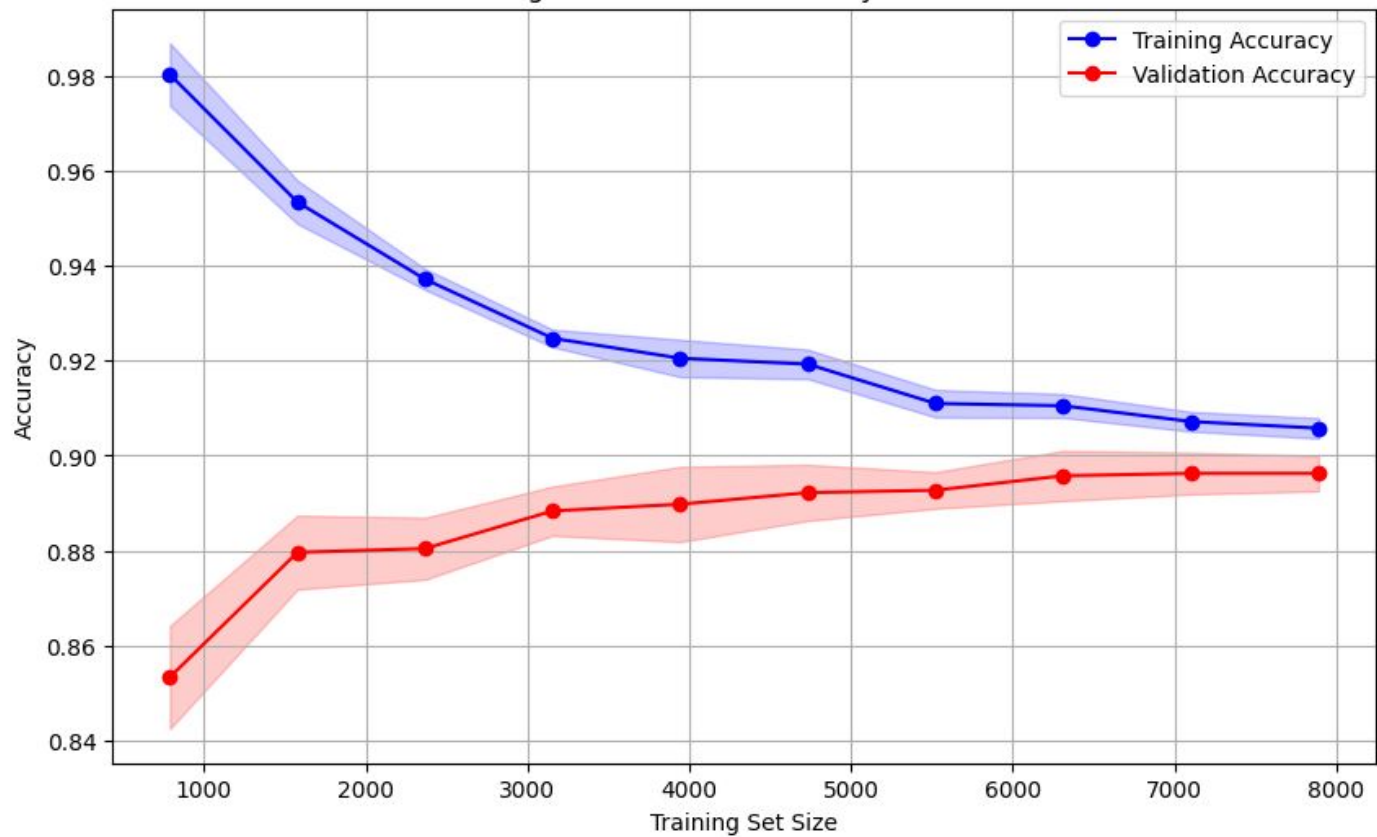
Methodology & Tools

- **Algorithm Used:** Support Vector Machine (SVM) with a polynomial kernel
- **Hyperparameter Tuning:** 5-fold internal cross-validation to determine optimal cost and degree
- **Evaluation Metric:** Accuracy
-
- `pandas` for data manipulation
- `scikit-learn` for model training, preprocessing, and evaluation
- `datasets` to load UCI dataset

Model Training

- **Model Construction:**
 - SVM Classifier implemented using `SVC(kernel="poly")`
 -
- **Hyperparameter Optimization:**
 - Grid search performed using `GridSearchCV()` with parameter grid:
 - `classifier__degree`: [2, 3, 4, 7]
 - `classifier__C`: [0.1, 10, 100, 1000]
 - Best parameters identified using 5 fold cross-validation, scoring based on accuracy

Learning Curve for SVM with Polynomial Kernel

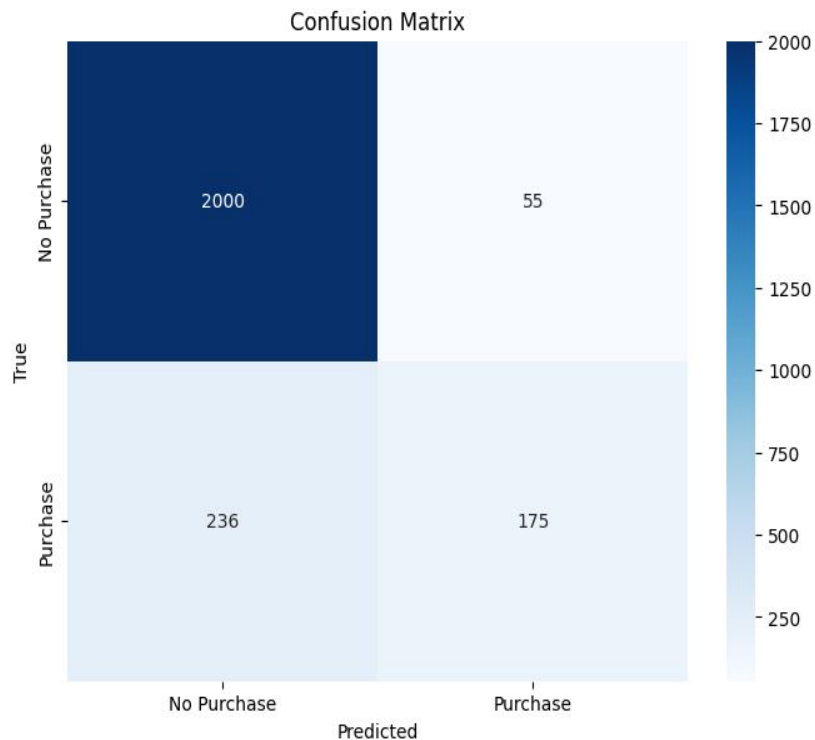


Evaluation

- **Accuracy Score:** Used as the primary metric to assess model performance (`accuracy_score()` from `sklearn.metrics`)
- **Confusion Matrix:** Provides a breakdown of true positives, true negatives, false positives, and false negatives (`confusion_matrix()` from `sklearn.metrics`)
- **Feature Importance (Permutation Importance Analysis):**
 - Computed using `permutation_importance()` from `sklearn.inspection`
 - Evaluates how much model accuracy decreases when a feature is randomly shuffled

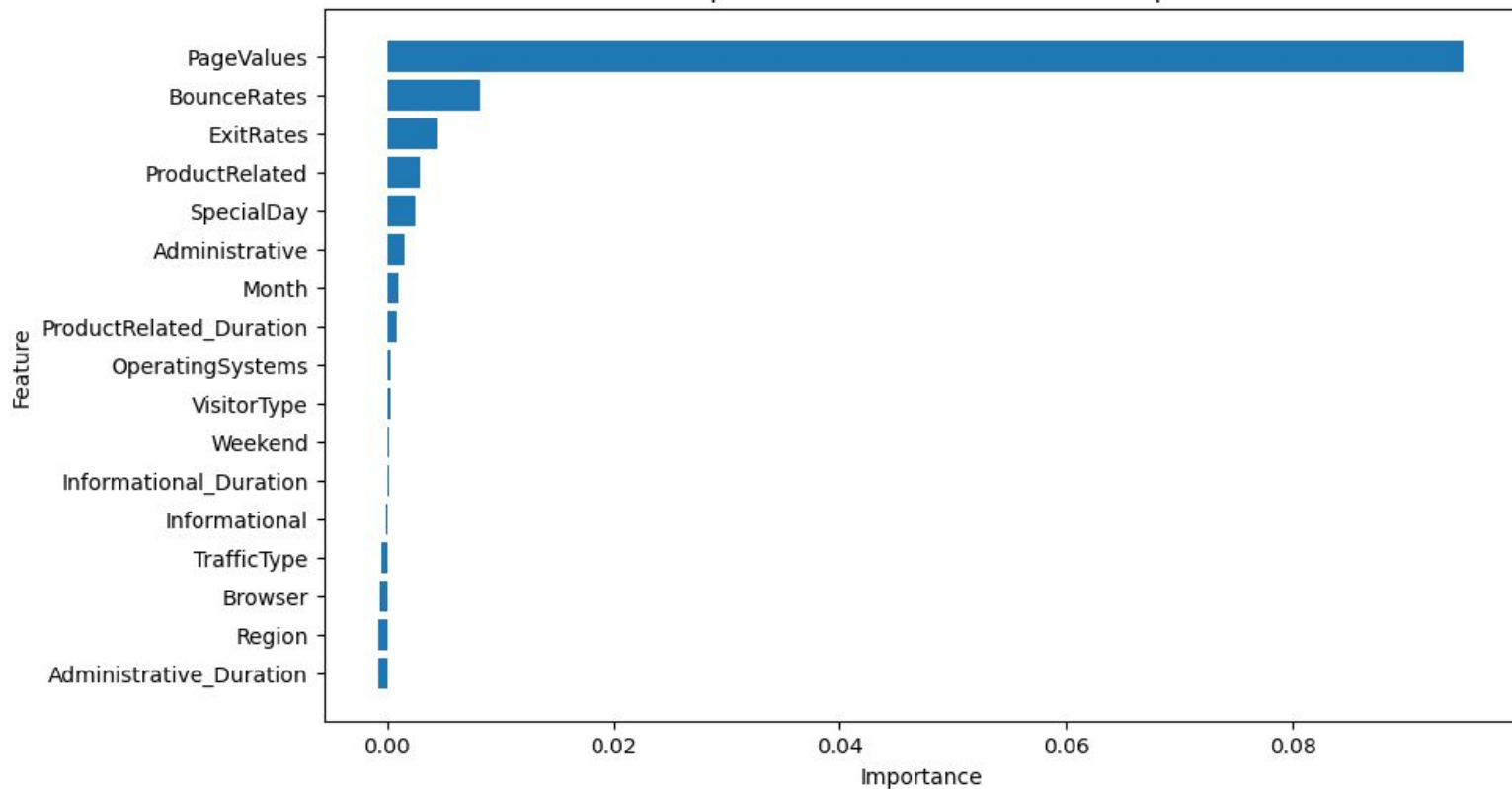
Key Findings & Insights

- **Best SVM Parameters Found:** Cost: 100, Degree of Polynomial Kernel: 2
- **Training Accuracy:** 0.8963
- **Test Accuracy:** 0.8820



Feature Importance (Permutation Importance Analysis)

Feature Importance based on Permutation Importance



Conclusion

Business Implications

- **Customer engagement with high-value pages** is the strongest predictor of purchasing behavior
- **Session quality metrics** (bounce/exit rates) are more valuable than visitor demographics
- Product browsing behavior is moderately predictive of purchase intent
- Temporal factors (Month, Weekend) and regional information have minimal predictive power
- Technical attributes (Browser, OS) show little correlation with purchasing decisions