

# BigMart Sales Prediction

---

A comprehensive machine learning solution for retail sales prediction, demonstrating end-to-end data science methodology from exploratory analysis to production deployment.

## Project Overview

This project implements a complete ML pipeline to predict **Item\_Outlet\_Sales** for a retail chain. The solution follows industry best practices and achieves a 12.5% improvement over baseline performance through systematic feature engineering, model optimization, and ensemble methods.

## Key Results

- **Final Model Performance:** RMSE 1245.87 (12.4% improvement over baseline)
- **Feature Engineering:** 40+ derived features from statistical analysis
- **Production Ready:** Serialized pipeline with comprehensive validation
- **Statistical Rigor:** Hypothesis testing confirms feature significance

## Technical Highlights

### Data Science Methodology

1. **Exploratory Data Analysis** - Comprehensive statistical analysis and visualization
2. **Hypothesis Testing** - Statistical validation of feature relationships
3. **Feature Engineering** - Smart imputation and domain-specific feature creation
4. **Model Development** - Systematic algorithm comparison and optimization
5. **Ensemble Methods** - Advanced stacking with optimized weights
6. **Production Pipeline** - Deployment-ready artifacts with monitoring

### Machine Learning Architecture

- **Ensemble Model:** LightGBM + CatBoost + Neural Network
- **Validation Strategy:** 5-fold Group Cross-Validation
- **Optimization:** Bayesian hyperparameter tuning with Optuna
- **Feature Count:** 40+ engineered features from 12 original features

## Repository Structure

```
├─ baseline_pipelines/      # EDA and hypothesis testing notebooks
├─ finetuning_pipeline/     # Model development and optimization
├─ production_models/      # Serialized model artifacts
├─ bigmart_pipeline/       # Production deployment pipeline
├─ BigMart_Report_Notebook.ipynb # Complete interactive analysis
├─ REPORT.md               # Detailed technical documentation
└─ Email_Report.md         # Executive summary for stakeholders
```

# Quick Start

## Prerequisites

```
pip install pandas numpy scikit-learn lightgbm catboost matplotlib seaborn  
scipy
```

## Run Analysis

```
# Complete interactive analysis  
jupyter notebook BigMart_Report_Notebook.ipynb  
  
# Individual pipeline components  
jupyter notebook baseline_pipelines/EDA.ipynb  
jupyter notebook finetuning_pipeline/final_version.ipynb
```

## Generate Predictions

```
from finetuning_pipeline.BigMartpreprocessing import BigMartPreprocessor  
import pickle  
  
# Load trained pipeline  
with open('bigmart_pipeline/bigmart_preprocessor.pkl', 'rb') as f:  
    preprocessor = pickle.load(f)  
  
# Process new data and predict  
predictions = model.predict(new_data)
```

# Key Features

## Statistical Analysis

- **Missing Value Analysis:** Smart imputation strategies using hierarchical grouping
- **Hypothesis Testing:** Kruskal-Wallis and Spearman correlation tests
- **Feature Relationships:** Comprehensive correlation and effect size analysis

## Feature Engineering

- **Statistical Features:** Item and outlet-level aggregated statistics
- **Business Logic:** Premium item flags, age calculations, ratio features
- **Categorical Encoding:** One-hot encoding with fallback handling
- **Target Engineering:** Log transformation for improved model stability

## Model Development

- **Algorithm Comparison:** RandomForest, LightGBM, CatBoost, XGBoost
- **Hyperparameter Optimization:** Bayesian search with cross-validation
- **Ensemble Methods:** Stacking with optimized weight combinations
- **Validation:** Robust cross-validation preventing data leakage

Production Features

- **Serialization:** Complete pipeline saved as pickle artifacts
- **Input Validation:** Schema checking and error handling
- **Scalability:** Efficient batch processing capabilities
- **Monitoring:** Prediction confidence intervals and interpretability

Performance Summary

Model	RMSE	Improvement
Baseline RandomForest	1421.89	-
Optimized LightGBM	1287.65	9.4%
Final Ensemble	1245.87	12.4%

Documentation

- [REPORT.md](#) - Complete technical methodology and results
- [Email\\_Report.md](#) - Executive summary for stakeholders
- [BigMart\\_Report\\_Notebook.ipynb](#) - Interactive analysis with plots
- [baseline\\_pipelines/](#) - Step-by-step EDA and hypothesis testing
- [finetuning\\_pipeline/](#) - Model development and optimization

Business Impact

The model provides reliable sales predictions enabling:

- **Inventory Optimization:** Accurate demand forecasting
- **Revenue Planning:** Data-driven sales projections
- **Store Performance:** Outlet-specific insights and recommendations
- **Product Strategy:** Item-level performance analysis

Technical Skills Demonstrated

- **Statistics:** Hypothesis testing, correlation analysis, effect sizes
- **Machine Learning:** Feature engineering, ensemble methods, hyperparameter tuning
- **Software Engineering:** Object-oriented design, serialization, error handling
- **Data Science:** End-to-end pipeline development, model validation, production deployment

Contributing

This project demonstrates professional data science practices including:

- Reproducible analysis with clear documentation
- Production-ready code with proper error handling
- Comprehensive testing and validation
- Clean, maintainable codebase following best practices

## License

This project is for educational and demonstration purposes.

---

**Contact:** [Your Email]

**LinkedIn:** [Your Profile]

**Portfolio:** [Your Website]