# Minseok Oh,
## Data Scientist
## Portfolio

**MSIS at SCU (Expected June 2025)**
**Ex-ORACLE, LG**

Data Scientist with 10+ years of experience specializing in ML/DL solutions that drive business value. Track record of delivering high-impact projects through experimentation and data-driven decision making.
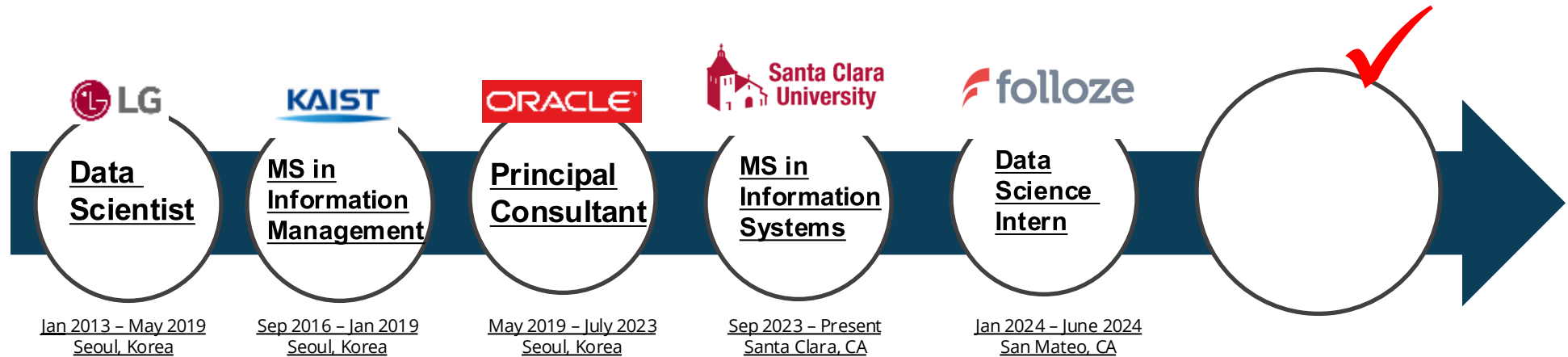
**Cupertino CA, 95014**
**(408) 334 – 5898**
**ohmseok0524@gmail.com**
**https://www.linkedin.com/in/scottmsoh/**

✓ **Driven to push boundaries in data science across industries, from manufacturing to enterprise solutions Now ready to create transformative impact through innovative data-driven solutions in Silicon Valley**
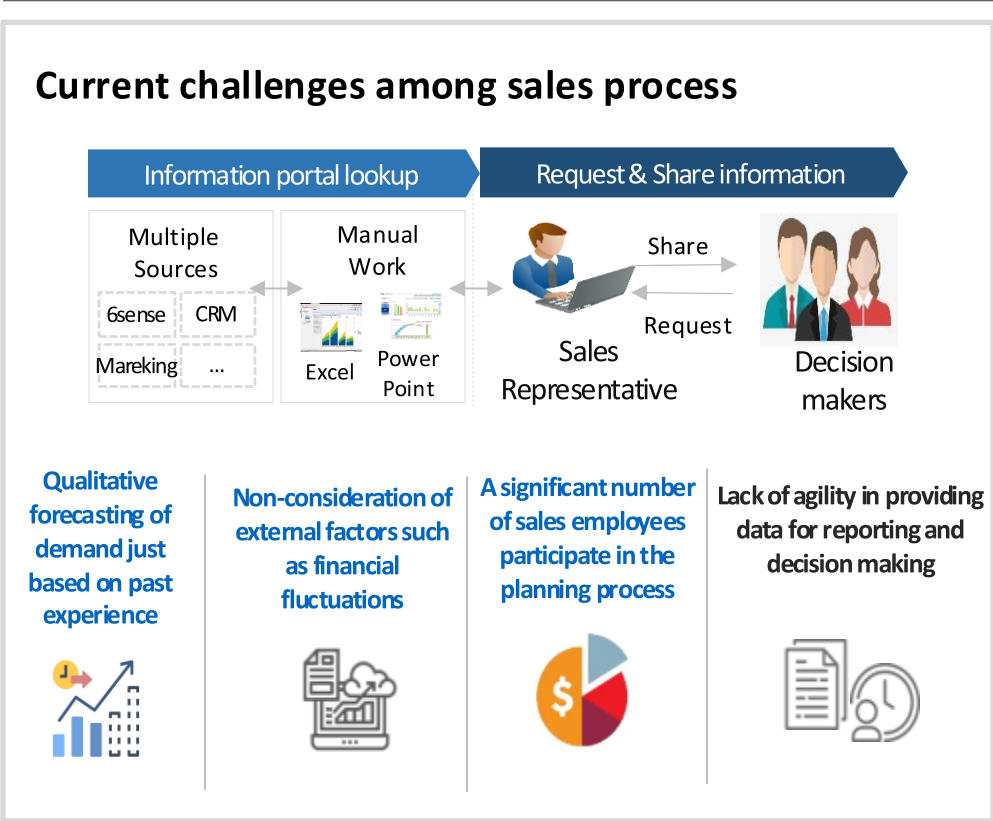


| Data Scientist | MS in Information Management | Principal Consultant | MS in Information Systems | Data Science Intern | |
| --- | --- | --- | --- | --- | --- |
| Jan 2013 – May 2019 Seoul, Korea | Sep 2016 – Jan 2019 Seoul, Korea | May 2019 – July 2023 Seoul, Korea | Sep 2023 – Present Santa Clara, CA | Jan 2024 – June 2024 San Mateo, CA | |

✓ Started career as a **Data Scientist at LG**, working on **process optimization and sentiment analysis using NLP** in the manufacturing sector. To broaden my experience beyond manufacturing, I pursued an **MS in Information Management at KAIST**, where I worked on **machine learning and data analytics projects** across **education, e-commerce, and healthcare**, gaining hands-on experience in **A/B testing and experimental design**.

✓ After graduation, I joined **Oracle**, expanding my expertise with **enterprise-scale projects** in **finance, insurance, retail, and pharmaceuticals**, focusing on the predictive modeling such as **recommendation systems, demand forecasting, and churn prediction**. Seeking **further growth and diverse challenges**, I moved to **Silicon Valley** to advance my career in data science.
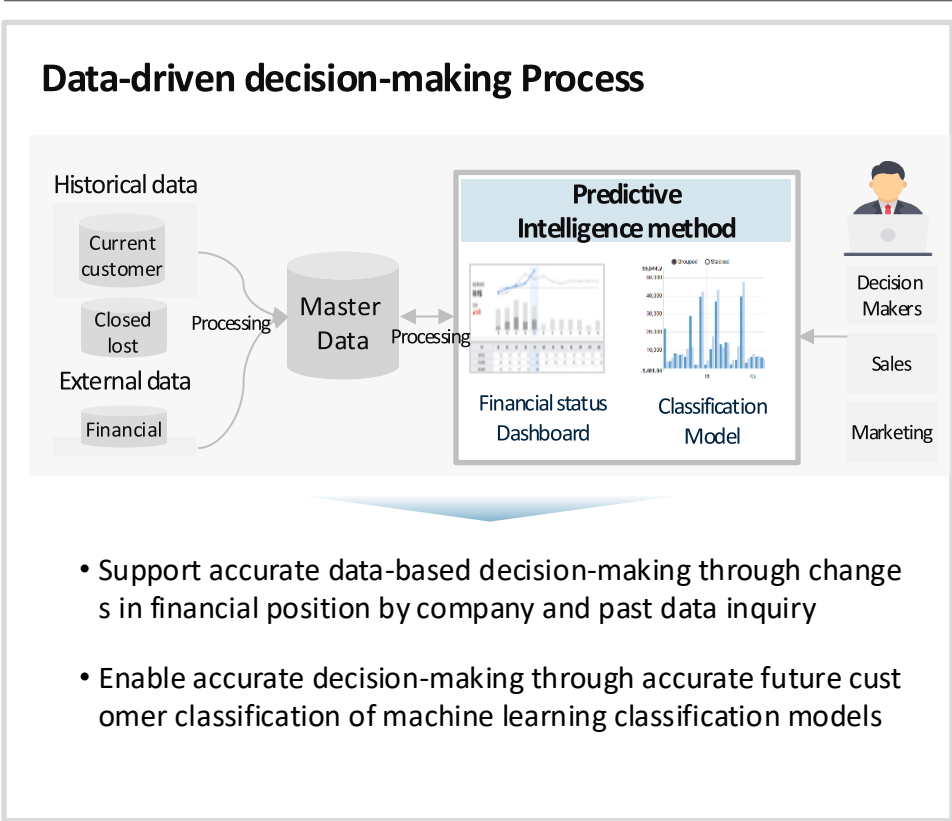
# CONTENTS ····

✓ **Transitioning from manual, experience-based sales decision-making to a data-driven approach using machine learning techniques for improved accuracy and efficiency.**
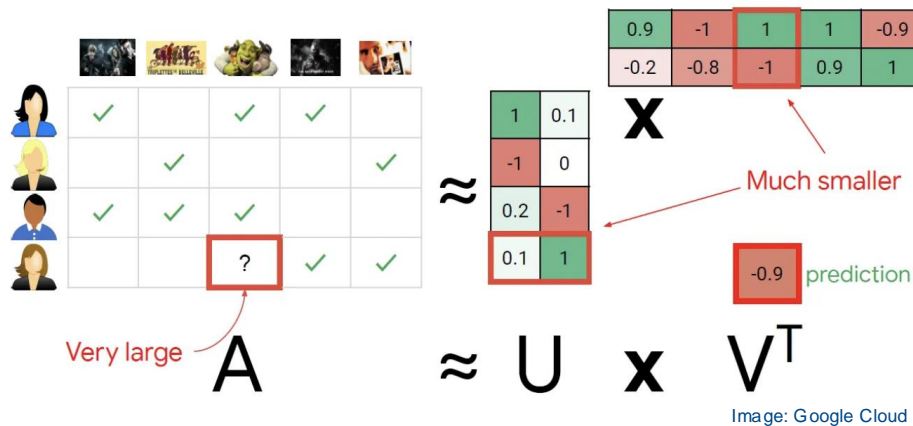
## As-Is

### Current challenges among sales process

Information portal lookup | Request & Share information

Multiple Sources
- 6sense
- CRM
- Mareking
- ...

Manual Work
- Excel
- Power Point

Share

Request

Sales Representative

Decision makers

**Qualitative forecasting of demand just based on past experience**

**Non-consideration of external factors such as financial fluctuations**

**A significant number of sales employees participate in the planning process**

**Lack of agility in providing data for reporting and decision making**

## To-Be

### Data-driven decision-making Process

Historical data
- Current customer
- Closed lost

External data
- Financial

Processing

Master Data

Processing

**Predictive Intelligence method**

Financial status Dashboard | Classification Model

Decision Makers

Sales

Marketing

- Support accurate data-based decision-making through changes in financial position by company and past data inquiry

- Enable accurate decision-making through accurate future customer classification of machine learning classification models

| | Customer Propensity Analysis for Sales Forecasting and Targeting | |
|---|---|---|
| 1 | **Business pain-point** | Limited sales resources make it difficult to manage all customer lists.<br>    ➢ **Aim to focus on high-potential customers by predicting potential vs. non-potential customers.** |
| 2 | **Project Overview** | This project aims to develop a machine learning classification model to identify potential future customers.<br>    ➢ **Optimizes lead prioritization, enhances customer targeting, improves sales efficiency, and drives business growth.** |
| 3 | **Data Collection** | • Customer & Company Data:<br>    ➢ **Demographic/transactional information of clients and non-clients.**<br>• External Financial Data:<br>    ➢ **MRDS dataset containing company profiles and financial indicators.**<br>• Feature Set:<br>    ➢ **130+ features, including company metadata (e.g. industry, contract details) and financial metrics (e.g., revenue, credit rating).** |
| 4 | **Data Preparation** | • Imbalanced Data:<br>    ➢ **Class imbalance between customers (less than 10%) and non-customers.**<br>• Missing Value Imputation: Applied median, mean, and mode where necessary<br>    ➢ **Tree-based model used for imputing crucial missing values (SPC SRC: D~A+)**<br>• Scaling & Transformation:<br>    ➢ **Used Robust Scaler and Log Transformation for skewed data.**<br>• Outlier Detection & Treatment:<br>    ➢ **Identified and handled outliers using IQR** |
| 4 | **Exploratory Data Analysis** | • Examined key financial metrics such as revenue, liabilities, and stockholder equity.<br>• Identified correlations between features impacting customer classification<br>    ➢ **Remove multicollinearity using VIF, PCA, Regularization(L1) techniques**<br>• Visualized distribution trends across different customer segments<br>    ➢ **Small, medium, and large customers, each capturing unique characteristics specific to their segment.**<br>    ➢ **Ensemble models (Develop 3 different models)** |

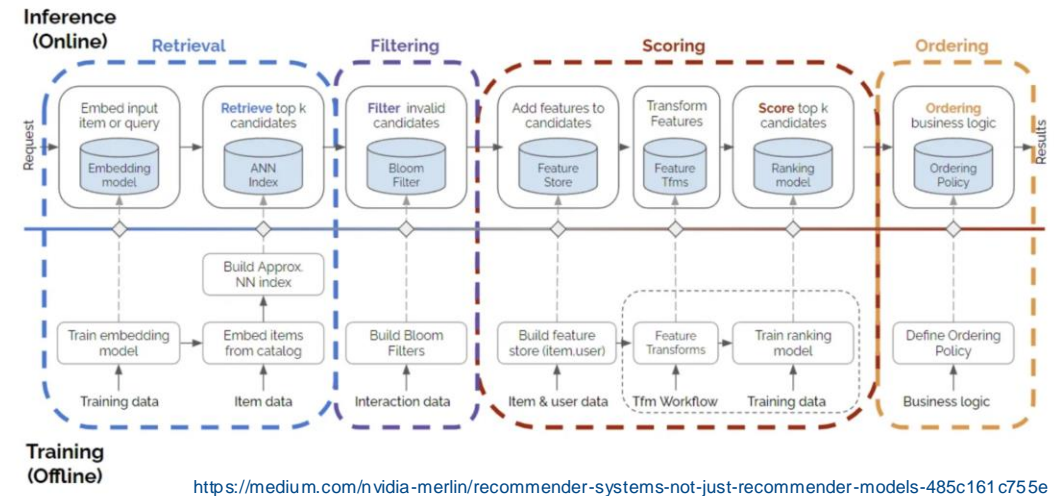| | Customer Propensity Analysis for Sales Forecasting and Targeting | |
|---|---|---|
| 5 | **Feature Engineering** | • Create new features: lifetime value, Report filing delay, Gap preliminary official |
| 6 | **Manage Imbalance data** | • Consider and experiment with different methods to address data class imbalance:<br>      **1) SMOTE (Synthetic Minority Over-sampling Technique)**<br>      2) Random oversample<br>      3) Random undersample<br>      4) Class weight (Cost-sensitive modeling) |
| 7 | **Model development & evaluation** | • Experimented with and trained various models, including SVM, Decision Tree, Random Forest, XGBoost, LightGBM, AdaBoost with Decision Tree, Logistic Regression, and an Ensemble Voting model<br>    ➢ **3 XGBoost models with SMOTE (for small, medium, and large customers) ensembled using weighted majority voting (Tried stacking as well, but it led to overfitting and was too complex to maintain)**<br>    ➢ **Recall is crucial to avoid missing potential customers, ensuring maximum lead capture. A recall of 0.82 minimizes lost opportunities, while an F-beta score of 0.83 balances precision and recall**<br><br>• To Reduce Overfitting: **PCA, Regularization (L1, L2), VIF (Remove features bigger than 8.0)** |
| 8 | **Hyper parameter tuning** | • Fine-tuned hyperparameters to mitigate overfitting and enhance performances<br>    ➢ **max_depth = 5**<br>    ➢ **Eta = 0.2**<br>    ➢ **Colsample_bytree = 0.9**<br>    ➢ **Colsample_bylevel = 0.7** |
| 9 | **Reflections and Future Improvements** | • There were many private clients, but the lack of public financial data required exploring alternative solutions<br>• Aimed to apply the model to the business through A/B testing, but due to time constraints, we handed it over to the responsible team |

✓ **Designed a 3-stage re-ranking architecture incorporating SGD-based matrix factorization for collaborative filtering, optimizing item recommendations based on store-specific sales volume**

✓ **SGD matrix factorization**



Image: Google Cloud

✓ **3-stage recommendation system architecture**



https://medium.com/nvidia-merlin/recommender-systems-not-just-recommender-models-485c161c755e

- Factorizes the store-item interaction matrix into lower-dimensional latent factors.
- SGD (Stochastic Gradient Descent) is used to optimize the factorization.
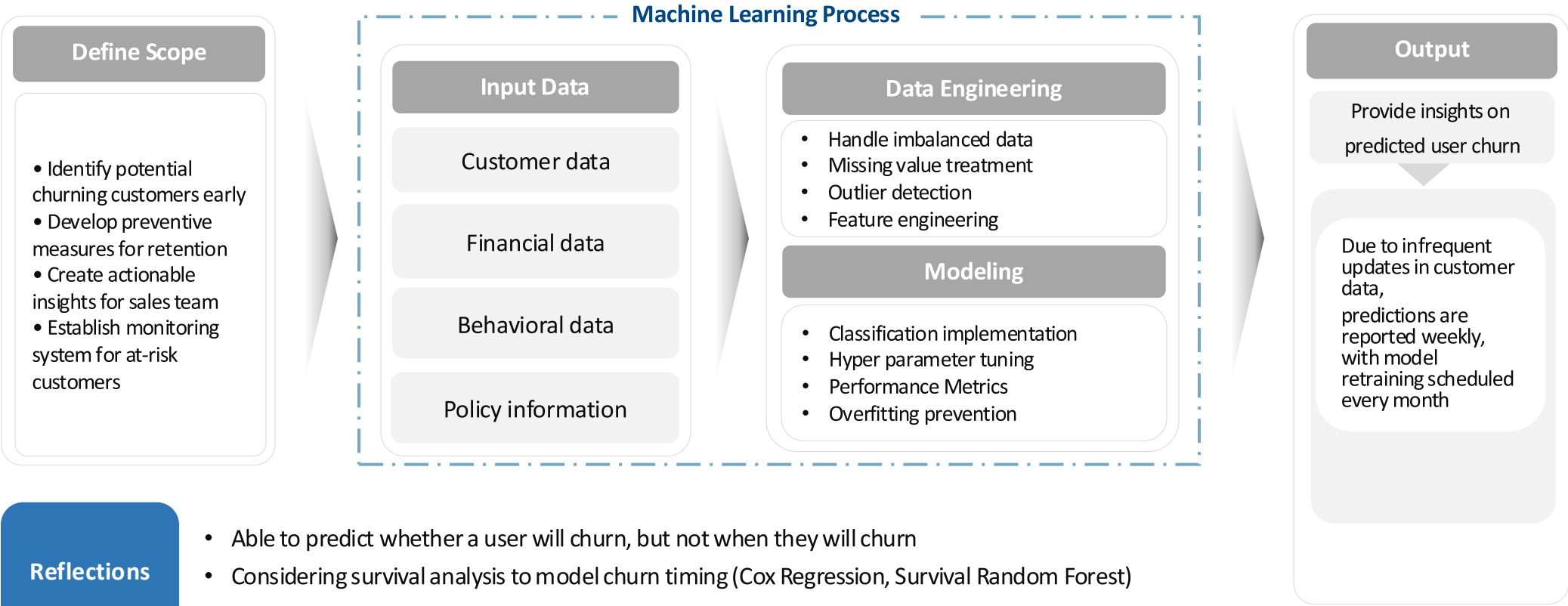- Predicts missing interactions (e.g., ratings, purchases) based on learned embeddings.

- Retrieval: Fetches top-N relevant candidates using an ANN index.
- Filtering: Applies Bloom Filters to remove invalid recommendations.
- Scoring & Ordering:
- Enhances rankings with additional feature transformations.
- Uses a ranking model (e.g., ML-based scoring).

7

| | Optimizing Store-Specific Item Recommendations with a 3-Stage Re-Ranking Architecture | |
|---|---|---|
| 1 | **Business pain-point** | Faced stagnant sales growth and sought new revenue opportunities by leveraging machine learning and a data-driven recommendation system to optimize inventory management and enhance product sales. |
| 2 | **Project Overview** | • Tried to propose market basket analysis using POS sales data<br>• Ultimately designed a 3-stage re-ranking architecture with SGD matrix factorization-based Collaborative Filtering (CF) recommendation model to suggest items based on store-specific sales volume. |
| 3 | **Architecture** | Designed a 3-stage Re-Ranking Architecture<br>    ➢ Stage 1: SGD-based Matrix Factorization – Generates initial item recommendations.<br>    ➢ FAISS Indexer: Enhances efficiency by enabling fast similarity search.<br>    ➢ Stage 2: Bloom Filter: Filters out items already sold at the store.<br>    ➢ Stage 3: XGBoost Ranker – Refines rankings based on additional features. |
| 4 | **Data Collection** | • 27 data tables (Supplier, customer information, POS data and so on)<br>    ➢ **Utilized store-specific item sales quantity for recommendations.**<br>    **(Since they operate only offline stores, they lack customer interaction data)** |
| 5 | **Data Preparation & Feature engineering** | • Sparse Matrix Construction – Converting raw transactional data into a store-item interaction matrix.<br>• Feature Engineering for Stage 2 Re-Ranking Model : Used around 40 features<br>    ➢ Store Average Sales Amount<br>        The mean sales given by a store across all items on sale.<br>    ➢ Item Average Sales Amount<br>        The mean sales an item received from all stores.<br>    ➢ Store-Store Cosine Similarity<br>        Measures how similar a store is to other stores based on sales patterns.<br>    ➢ Item-Item Cosine Similarity<br>        Measures item similarity based on shared store interactions.<br>    ➢ Rating Count Features<br>        Number of sales per store and sales per item, which help assess data density and confidence levels. |

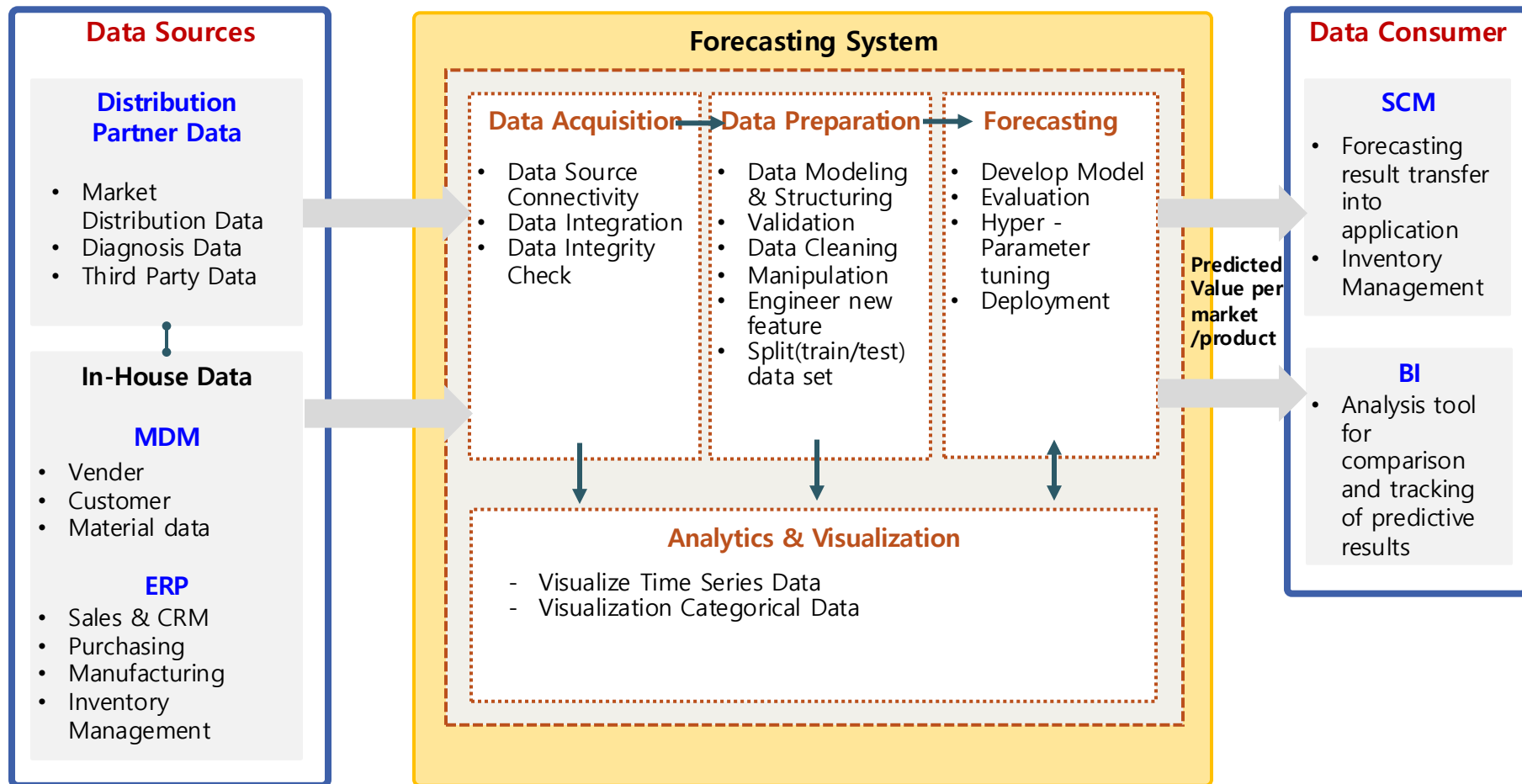| | | Optimizing Store-Specific Item Recommendations with a 3-Stage Re-Ranking Architecture |
|---|---|---|
| 5 | **Model selection & development** | Consider of four recommendation models:<br>• Apriori with Association Rules → Excluded<br>    ➢ Difficult to set optimal lift and support thresholds, automation challenges.<br>• FP-Growth → Excluded<br>    ➢ Interpretation of tree-based patterns is challenging<br>• **Collaborative Filtering → Adopted (Benchmarked from research papers on store-item sales-based recommendation models)**<br>    ➢ KNN-based CF: Too simple; performance was suboptimal<br>    ➢ **SGD-based Matrix Factorization: Suitable for capturing high-dimensional interaction patterns.**<br>• Deep Learning-Based Models → Excluded<br>    ➢ Lack of organizational expertise in machine learning/deep learning<br>    ➢ No available infrastructure to support deep learning deployment. |
| 6 | **Model Evaluation** | Optimizes score calibration for fine-tuning the ranking, where metrics NDCG: NDCG@10: 0.79 → Top 10 item selection |
| 7 | **A/B test & Result** | • Base Metric: Daily Revenue<br>• Expected Metric: 5% Increase in Daily Revenue<br>• Controlled Conditions<br>    ➢ Same Store Locations<br>    ➢ Same Day of the Week (to control demand fluctuations)<br>    ➢ Same Time Window (to eliminate hourly sales variations)<br>• Statistical Parameters:<br>    ➢ Power (1-β): 85%<br>    ➢ Significance Level (α): 5%<br>    ➢ Experiment period (Based on minimum Sample Size): 31.36 days (calculated using t-test formula).<br>• A/B Test Results:<br>    ➢ Comparison by Time Zone & Item Category<br>    ➢ Lift & p-value analysis to determine statistical significance.<br>    ➢ **Revenue (Monitoring Metric) shows a Significant Positive Impact (+4.5%)**<br>    ➢ Indicates Quantifiable Growth due to product recommendations<br>    ➢ Confirms Positive Business Impact from optimization in recommendation strategy |

**Outline**   Early Warning System for Customer Churn

**Define Scope**

• Identify potential churning customers early
• Develop preventive measures for retention
• Create actionable insights for sales team
• Establish monitoring system for at-risk customers

**Machine Learning Process**

**Input Data**

Customer data

Financial data

Behavioral data

Policy information

**Data Engineering**

• Handle imbalanced data
• Missing value treatment
• Outlier detection
• Feature engineering

**Modeling**

• Classification implementation
• Hyper parameter tuning
• Performance Metrics
• Overfitting prevention

**Output**

Provide insights on predicted user churn

Due to infrequent updates in customer data,
predictions are reported weekly, with model retraining scheduled every month

**Reflections**

• Able to predict whether a user will churn, but not when they will churn
• Considering survival analysis to model churn timing (Cox Regression, Survival Random Forest)

| | Early Warning System for Customer Churn | |
|---|---|---|
| **1** | **Business pain-point** | Problem: The company lacks proper churn management, leading to potential revenue loss.<br>➢ **Develop a churn prevention model to identify at-risk customers and provide actionable insights for retention.** |
| **2** | **Project Overview** | Utilize existing customer data to predict churn probability.<br>➢ **Generate reports on at-risk customers for marketing and sales teams to take proactive retention actions.** |
| **3** | **Data Collection** | Customer Data<br>• Feature Types:<br>  ➢ Demographics: Gender, age, marital status, and so on<br>  ➢ Financials: Income, retirement status<br>  ➢ Other Behavioral & Policy Data (Details confidential)<br>  ➢ Target: Churn or Not (Binary Classification) - Churned Customer → Has an End-date in their policy |
| 4 | Data Preparation & Exploratory Data Analysis | • **Highly Imbalanced Data: Churned customers are significantly fewer than non-churned customers.**<br>  ➢ Applied SMOTE, random oversampling/undersampling, or cost-sensitive learning.<br>• Handling Missing Values<br>  ➢ Used median for skewed data, mean imputation.<br>  ➢ For highly correlated features, built a regression model to estimate missing values.<br>• Skewed Data: Price-related features showed skewness.<br>  ➢ Applied log transformation for normalization.<br>• Outlier Detection & Treatment<br>  ➢ Used IQR (Interquartile Range) method for outlier handling. |

| | Early Warning System for Customer Churn | |
|---|---|---|
| 5 | **Feature Engineering** | • Create derived features: cancel_customer, dormant, lifetime_value, and so on.<br>• Create interaction Variables: 20 variables<br>    ➤ Arithmetic Interaction: e.g. income * age<br>    ➤ Conditional Relationship:  e.g. income/household_size<br>    ➤ Statistical Transformation: e.g. log(income) * education level<br>• Partial PCA: Apply PCA only to highly correlated subgroups related to price, preserving raw features for interpretation.<br>• One-hot encoding: Convert values into numeric values.<br>    ➤ Applied Lasso (L1) regression, which forces some coefficients to exactly 0, removing less important features. |
| 6 | **Manage Imbalance data** | Consider and experiment with different methods to address data class imbalance:<br>    **1) SMOTE (Synthetic Minority Over-sampling Technique)**<br>    2) Random oversample<br>    3) Random undersample<br>    4) Class weight (Cost-sensitive modeling) |
| 7 | **Model development & evaluation** | • Prioritizing interpretability over model performance (Understanding churn reasons and root causes is crucial)<br>• Prioritized **churn probability** over estimating user lifetime. (**Classification** vs Cox Regression, Survival Random Forest)<br>• Experimented with and trained various models, including Logistic regression, SVM, Decision Tree, Random Forest, XGBoost, LightGBM<br>    ➤ Model Used: XGBoost<br>    ➤ Primary Metric: Recall (Focus on minimizing false negatives to avoid missing churned customers).<br>    ➤ **Recall: 0.81, Precision: 0.75, F-2(beta) score: 0.80 (More weight on recall than precision)**<br>• To Reduce Overfitting: **PCA, Regularization (L1), VIF (Remove features bigger than 8.0)**<br>• **Due to infrequent updates in customer data, predictions are reported weekly, with model retraining scheduled every month** |
| 8 | **Hyper parameter tuning** | Fine-tuned hyperparameters to mitigate overfitting and enhance performances<br>    ➤ **max_depth = 6**<br>    ➤ **Eta = 0.3**<br>    ➤ **Colsample_bytree = 0.7**<br>    ➤ **Colsample_bylevel = 0.8** |
| 9 | **Reflections and Future Improvements** | • Able to predict whether a user will churn, but not when they will churn<br>• Considering survival analysis to model churn timing (Cox Regression, Survival Random Forest) |

✓ **Developed a biosimilar sales forecasting model using SARIMA to optimize production planning, enabling data-driven inventory management.**
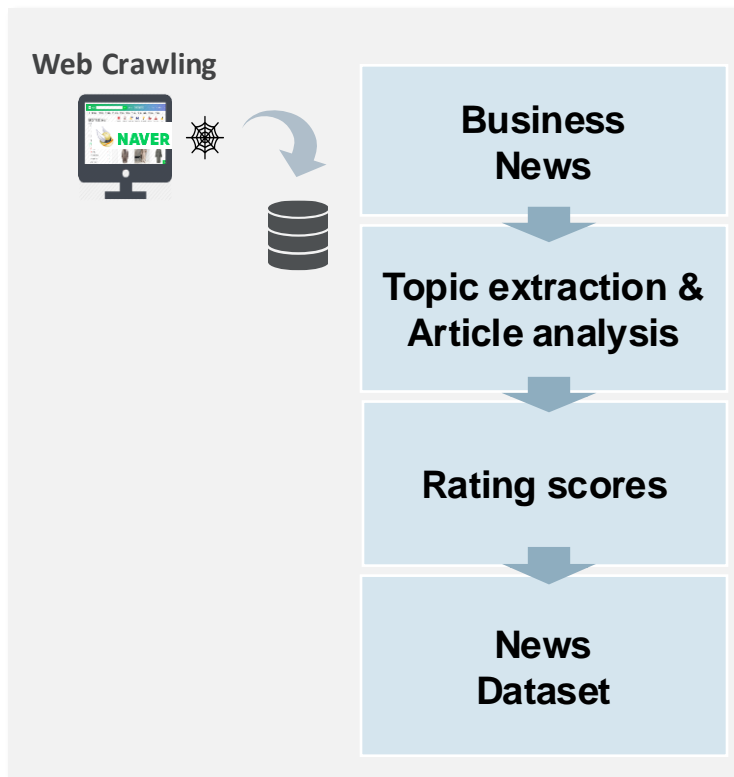
## Data Sources

### Distribution Partner Data

- Market Distribution Data
- Diagnosis Data
- Third Party Data

### In-House Data

#### MDM
- Vender
- Customer
- Material data

#### ERP
- Sales & CRM
- Purchasing
- Manufacturing
- Inventory Management

## Forecasting System

### Data Acquisition → Data Preparation → Forecasting

**Data Acquisition**
- Data Source Connectivity
- Data Integration
- Data Integrity Check

**Data Preparation**
- Data Modeling & Structuring
- Validation
- Data Cleaning
- Manipulation
- Engineer new feature
- Split(train/test) data set

**Forecasting**
- Develop Model
- Evaluation
- Hyper - Parameter tuning
- Deployment

### Analytics & Visualization
- Visualize Time Series Data
- Visualization Categorical Data

**Predicted Value per market /product**

## Data Consumer

### SCM
- Forecasting result transfer into application
- Inventory Management

### BI
- Analysis tool for comparison and tracking of predictive results

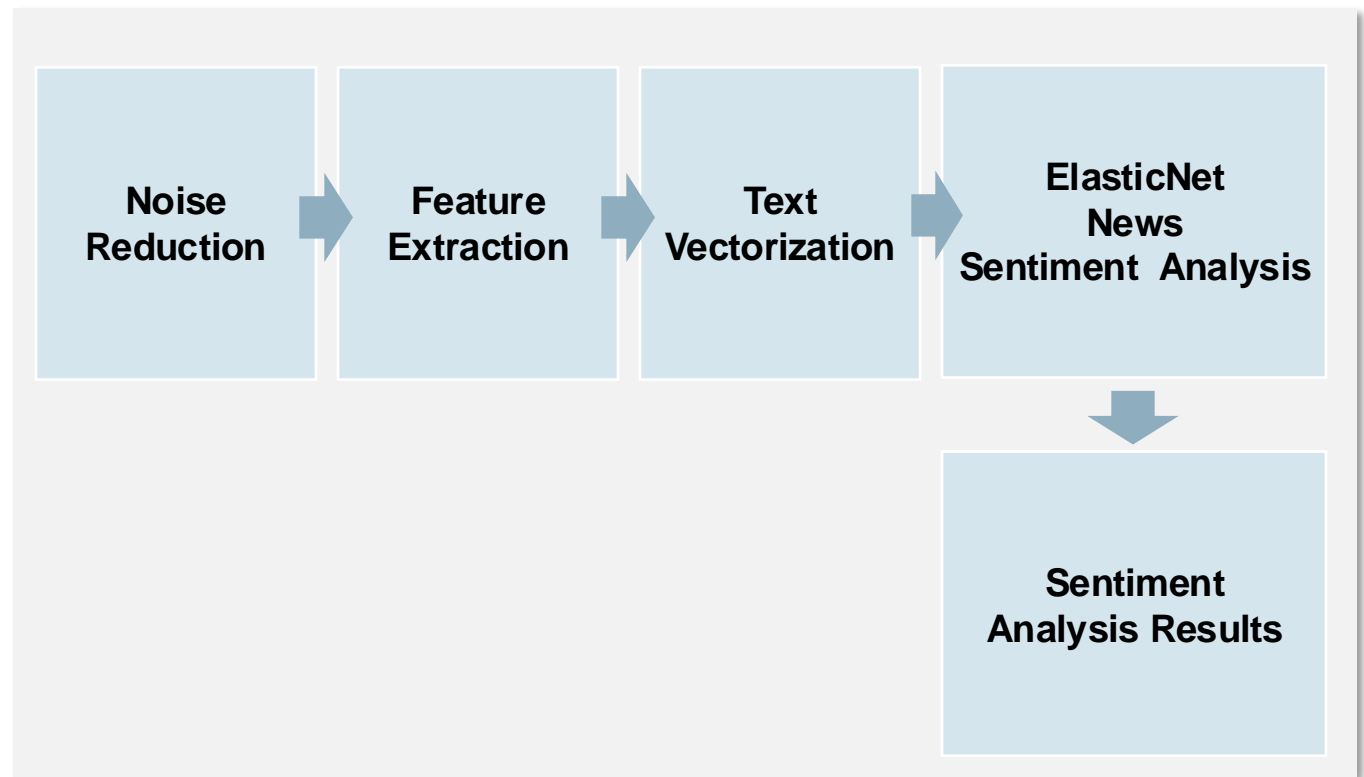| | Biosimilar Sales Forecasting Using ARIMA Time Series Analysis | |
|---|---|---|
| **1** | **Business pain-point** | • Lack of demand forecasting for current products prevents optimization of production lines by product type |
| **2** | **Project Overview** | • Aim to establish a foundation by developing a baseline demand forecasting model |
| **3** | **Data Collection** | • Provided Data<br>Time Range: total 10 years (Monthly herbal medicine data)<br>Regions: Includes US<br>Product Categories:<br>    ➢ Dosage Forms<br>    ➢ Packaging Types<br>Target Variables:<br>    ➢ **Revenue_unit**<br>    ➢ Revenue_value<br>    ➢ Price |
| **4** | **Data Preparation** | • Designed ETL processing to convert wide-format time-series data into long format using pd.melt() for easier analysis.<br>• Baseline modeling: Started with one product as the baseline model.<br>    ➢ Trend Analysis: Observed significant differences in sales unit trends across different products.<br>    ➢ Conclusion: Requires individualized modeling per product rather than a one-size-fits-all approach. |
| **5** | **Exploratory Data Analysis** | • Decomposition Analysis<br>    ➢ Compared Original, Seasonal, Trend, and Residual plots to analyze time-series components.<br>• Autocorrelation & Partial Autocorrelation Analysis<br>    ➢ Checked ACF/PACF plots to identify dependencies over time.<br>    ➢ Observed non-stationarity, indicating the need for transformation.<br>• Differencing for Stationarity<br>    ➢ Applied first and second-order differencing to remove trends and stabilize the series. |

| | Biosimilar Sales Forecasting Using ARIMA Time Series Analysis | |
|---|---|---|
| 6 | **Model selection** | • Decision: Time-Series Analysis Over Regression<br>   ➤ Initially considered both regression-based modeling and time-series forecasting<br>   ➤ The challenge was determining which features to include (e.g., competitor sales, prescriptions, stock levels, clinical trials, claims data)<br>   ➤ Due to the complexity and sequential nature of the data, we chose time-series analysis as the preferred approach<br>• Model Selection for Time-Series Analysis<br>   • Exponential Smoothing is simple but fails to capture seasonality, as observed in the decomposition analysis.<br>   • Since the trend exhibits seasonal patterns, a more advanced model like ARIMA (or SARIMA for seasonality) is required. |
| 7 | **Model development & evaluation** | • ACF & Partial ACF Comparison:<br>   ➤ Compared MA(q), AR(q), and ARMA(p,q) models.<br>   ➤ After second-order differencing, the data exhibited AR(q) characteristics.<br>• SARIMA outperformed ARIMA, effectively capturing seasonal trends with a higher $R^2$ (0.84) and lower AIC (429.78).<br>• Effectiveness validation was handed over to the client MLOps team for A/B testing and the deployment. |
| 8 | **Hyper parameter tuning** | • Hyper parameter tunning & performance comparison:<br>   ➤ ARIMA (p:2, d:1, q:2) → AIC: 670, $R^2$: 0.45<br>   ➤ Seasonal ARIMA (SARIMAX) (p:2, d:1, q:2) with Seasonal (P:2, D:1, Q:0, s:12) → AIC: 429.78, $R^2$: 0.84<br>   ➤ Significantly better performance with seasonal adjustment and hyper parameter tuning. |
| 9 | **Reflections and Future Improvements** | • Extend the SARIMA model to cover all regions and new pharmaceutical products.<br>• Adapt the model for different market conditions and product-specific trends.<br>• Fine-tune parameters for region-specific and product-level forecasting accuracy. |

✓ **Analyze daily news sentiment to provide early notifications, helping stakeholders protect the company's reputation**

✓ **Web crawling**

✓ **Data engineering & model development**

**Web Crawling**

Business News

↓

Topic extraction & Article analysis

↓

Rating scores

↓

News Dataset

Noise Reduction → Feature Extraction → Text Vectorization → ElasticNet News Sentiment Analysis

↓

Sentiment Analysis Results

| | News Sentiment Analysis System for IR Strategy | |
|---|---|---|
| 1 | **Business pain-point** | • Due to high stock price volatility, the company faced increased risk exposure for stakeholders. This volatility created challenges in maintaining financial stability and required proactive measures to mitigate potential risks. |
| 2 | **Project Overview** | • To strengthen risk management and enable proactive decision-making, we identified the need for a real-time news sentiment analysis system. |
| 3 | **Data Collection** | • Scraped 15,000 news articles from NAVER News<br>• Collected 30 days of data (March 2018), averaging 500 articles daily.<br>• Selected only business-related keywords (e.g., Samsung, display, and so on). |
| 4 | **Data Preparation** | • Sentiment Labeling<br>   ➢ Employed domain experts for manual sentiment scoring (1-5 scale) to ensure sentiment classification aligns with LG's perspective, accurately assessing whether news articles have a positive or negative impact on the company.<br>   ➢ Used scores as training data for a regression model to predict sentiment<br>• Quality Check<br>   ➢ Verified sentiment score accuracy before model training. |
| 5 | **Data engineering** | • Morphological Analysis<br>   • Korean text requires morphological analysis to break words into meaningful units.<br>• Part-of-Speech (POS) Filtering<br>   • Select only nouns, verbs, and adjectives<br>• Stopword Removal & Lemmatization<br>   • Remove commonly used words<br>   • Convert words to their base forms<br>• Deduplication & Special Character Removal<br>   • Remove duplicate words to avoid bias in frequency-based models.<br>   • Filter out special characters<br>• TF-IDF Transformation<br>   • Convert text into a Document-Term Matrix<br>   • min_df=3 → Exclude rare words appearing in less than 2–3 documents.<br>   • max_df=0.90 → Remove extremely common words appearing in over 95% of documents. |

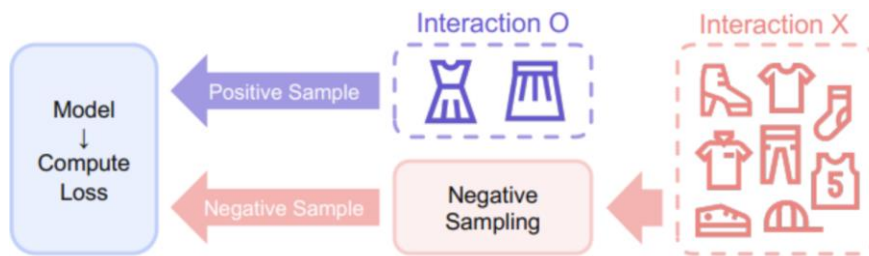| | News Sentiment Analysis System for IR Strategy | |
|---|---|---|
| 6 | **Model selection** | • Model Selection:<br>     • **News titles effectively capture the sentiment and key takeaways of the article.**<br>     • Efficiency—processing full articles requires more computational resources.<br>     • Deep Learning (LSTM, RNN): Not chosen because news titles are short (length <20), making complex models unnecessary.<br>     • Explored reinforcement learning models as well but insufficient data for meaningful reinforcement learning.<br> • **Decision: ElasticNetCV (Regularized Regression)**<br>     • **Applies penalty terms (L1 & L2 regularization) to handle infrequent words.**<br>     ➤ Balances feature sparsity (L1) and weight shrinkage (L2) for better generalization |
| 7 | **Model development & evaluation** | • Regression Performance:<br>     ➤ $R^2$: 0.79<br>     ➤ MAE: 0.35, MSE: 0.2<br> • Adjusting Sentiment Classification Threshold (Conservative Approach)<br>     • Reason: Reduces false positives, making the model more risk-averse in sentiment classification.<br>     ➤ Predicted sentiment score < 3.8 → Negative<br>     ➤ Predicted sentiment score ≥ 3.8 → Positive<br> • **Automated batch process runs daily at 4 AM to collect, score, and report the previous day's news sentiment to the IR team for proactive risk management** |
| 8 | **Hyper parameter tuning** | • l1_ratio = 0.8<br>     • l1_ratio = 1.0 → Pure Lasso Regression (stronger feature selection)<br>     • l1_ratio = 0.5 → Balanced L1 & L2 (default)<br>     • l1_ratio = 0.0 → Pure Ridge Regression (no feature elimination, only weight shrinkage) |
| 9 | **Reflections and Future Improvements** | • Expanding Data Sources:<br>     ➤ News titles captured sentiment well, but lacked context compared to full articles.<br>     ➤ Limited data scope (30 days) might not fully capture long-term sentiment trends.<br> • Real-time Monitoring & Deployment<br>     ➤ Develop a real-time dashboard for monitoring sentiment trends |

✓ **Design and run time zone/age-based A/B tests for female-targeted Facebook ads| 12% improvement in CTR through target segment optimization (Lift 1.12, p-value < 0.05, 30k samples)**

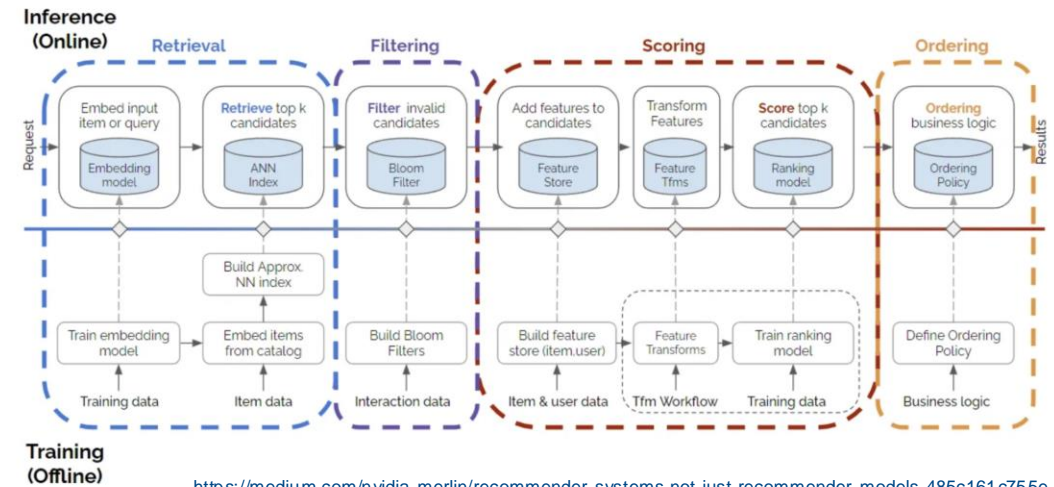| No | Contents | Descriptions |
|---|---|---|
| 1 | Target User | • Potential customers who interact with ads and marketing channels but have not yet converted |
| 2 | Observation | • Users engaging with the service through Facebook (especially women aged 36-54, active between 12-1 PM) and Email channels, where conversion rates are high |
| 3 | Problem Statement | • Reduce ad spend on low-conversion channels and reallocate budget to Facebook<br>• Providing more detailed information (images, details, etc.) on Naver, and other channels, similar to email, can improve conversion rates. |
| 4 | Hypothesis | • $H_0$: Reducing ad spend on low-conversion channels and reallocating it to Facebook does not result in a statistically significant increase in conversion rate<br>• $H_1$: Reducing ad spend on low-conversion channels and reallocating it to Facebook leads to a statistically significant increase in conversion rate |

| No | Contents | Descriptions |
|---|---|---|
| 6 | Experiment Group | • Group A (Control Group): Current ad spend distribution remains unchanged across all channels<br>• Group B (Experimental Group): Increase Facebook ad budget targeting high-converting segments (women aged 36-54, 12-1 PM). |
| 7 | Experiment Period | • Base metric: Facebook avg conversion rate<br>• Expected metric: 0.55%<br>• Alpha (Significance): 95% confidence level ($\alpha = 0.05$)<br>• 1-Beta (Power): 80%<br>   • Since sample collection was not difficult, there was no need to lower the power. Increasing the power was considered, but the cost was too high to justify.<br>• Minimum Sample size (z-test): around 32,000 people<br>• Period calculation: Required days |
| 8 | Metric | • Conversion Rate (CTR) |
| 9 | Trade-off | • Potential loss of customers from other channels, leading to reduced engagement |
| 10 | Andon | • A significant drop in conversion rate and revenue sustained for more than 3 days. |

✓ **Developed a 3-stage recommendation system leveraging negative sampling for implicit feedback and ANN-based retrieval, enhancing ranking with ML-based scoring for optimized item recommendations**

✓ **Negative sampling for implicit data**



✓ **3-stage recommendation system architecture**



https://medium.com/nvidia-merlin/recommender-systems-not-just-recommender-models-485c161c755e

- Limited availability of explicit data, such as purchase history or user ratings.
- Apply negative sampling to handle implicit feedback efficiently.
- Items with no interactions are assumed to be irrelevant or less preferred and are sampled as negative instances to optimize model training while reducing computational cost.

- Retrieval: Fetches top-N relevant candidates using an ANN index.
- Filtering: Applies Bloom Filters to remove invalid recommendations.
- Scoring & Ordering:
- Enhances rankings with additional feature transformations.
- Uses a ranking model (e.g., ML-based scoring).

| | E-commerce Recommendation Engine: Matrix Factorization with Implicit Feedback and Negative Sampling | |
|---|---|---|
| 1 | **Business pain-point** | • Had customer data available but was not effectively utilizing it to drive business impact |
| 2 | **Project Overview** | • Goal: Increase conversion rate (CVR) through a personalized recommendation system.<br>• Initial Approach: Developed an in-house user-based collaborative filtering model themselves.<br>• Issue: Model did not achieve satisfactory performance, requiring further optimization or alternative approaches. |
| 3 | **Architecture** | Existing Challenges & Improvements<br>  ➤ Challenge 1: Insufficient Training Data<br>     a) Explicit data (purchase history, ratings record) was not enough for robust training.<br>     b) Solution: Proposed using click history (implicit feedback) to augment training data and improve model performance.<br>  ➤ Challenge 2: Limitations of User-Based Collaborative Filtering (CF)<br>     a) User-based CF relies only on user-item interactions (ratings, purchases), making it relatively simple.<br>     b) Issue: Struggles to capture complex relationships beyond basic similarities.<br>     c) Solution: Adopted matrix factorization + feature-based ranking for deeper insights.<br>Proposed a 3-stage Re-Ranking Architecture<br>  ➤ Stage 1: SGD-based Matrix Factorization – Generates initial item recommendations.<br>  ➤ FAISS Indexer: Enhances efficiency by enabling fast similarity search.<br>  ➤ Stage 2: Bloom Filter: Filters out items already sold at the store.<br>  ➤ Stage 3: XGBoost Ranker – Refines rankings based on additional features. |
| 4 | **Data Collection** | • Since purchase history is insufficient, click data (implicit feedback) is used, and negative sampling is applied by selecting items the user did not interact with.<br>• Negative Sampling Methods for Implicit Data<br>  ➤ Random Negative Sampling<br>  ➤ **Popularity-Based Negative Sampling: More effective than random sampling in recommendation scenarios** |
| 5 | **Data Preparation** | • Sparse Matrix Construction – Converting raw transactional data into a user-click interaction matrix. |

| | E-commerce Recommendation Engine: Matrix Factorization with Implicit Feedback and Negative Sampling | |
|---|---|---|
| 5 | **Model selection & development** | Consider of four recommendation models:<br>• User-based Collaborative Filtering → SGD Matrix Factorization<br>    ➢ User-based method: pairwise similarity calculations (computationally expensive), struggles users who have few or no interactions<br>    ➢ KNN-based CF: Too simple; performance was suboptimal<br>    ➢ **SGD-based Matrix Factorization: captures latent factors, making it effective even with incomplete data.**<br>• Deep Learning-Based Models → Excluded<br>    ➢ Lack of organizational expertise in machine learning/deep learning<br>    ➢ No available infrastructure to support deep learning deployment. |
| 6 | **Model Evaluation** | Stage 1 focuses on retrieval quality (ranking high-relevance items first) using Recall@1000: 0.68<br>Stage 3 optimizes score calibration for fine-tuning the ranking, where ranking metrics NDCG@10: 0.71 → Top 10 items |
| 7 | **A/B test** | • Base Metric: CTR 3.0%<br>• Expected Metric: CTR 3.5% (Around 20.0% improvement)<br>• Controlled Conditions<br>    ➢ dd<br>• Statistical Parameters:<br>    ➢ Power (1-β): 80%<br>    ➢ Significance Level (α): 5%<br>    ➢ Minimum Sample Size: 5,295 people (calculated using z-test formula).<br>    ➢ Experiment period: daily user: 300 people, required days: 17.7 days *2 group = 35.5days<br>• A/B Test Results:<br>    ➢ Lift & p-value Analysis: Confirmed statistical significance of improvements.<br>    ➢ **CTR improved by 20% (3.0% → 3.6%), with evaluations conducted across OS types and item categories for deeper insights.**<br>    ➢ Business Impact: Demonstrates quantifiable growth driven by optimized product recommendations.<br>    ➢ Validation: Confirms positive business impact from strategic recommendation enhancements. |

**Expected Values on the Continuous Intention to Use IoT Products from the Perspective of Expectation-Confirmation Theory**

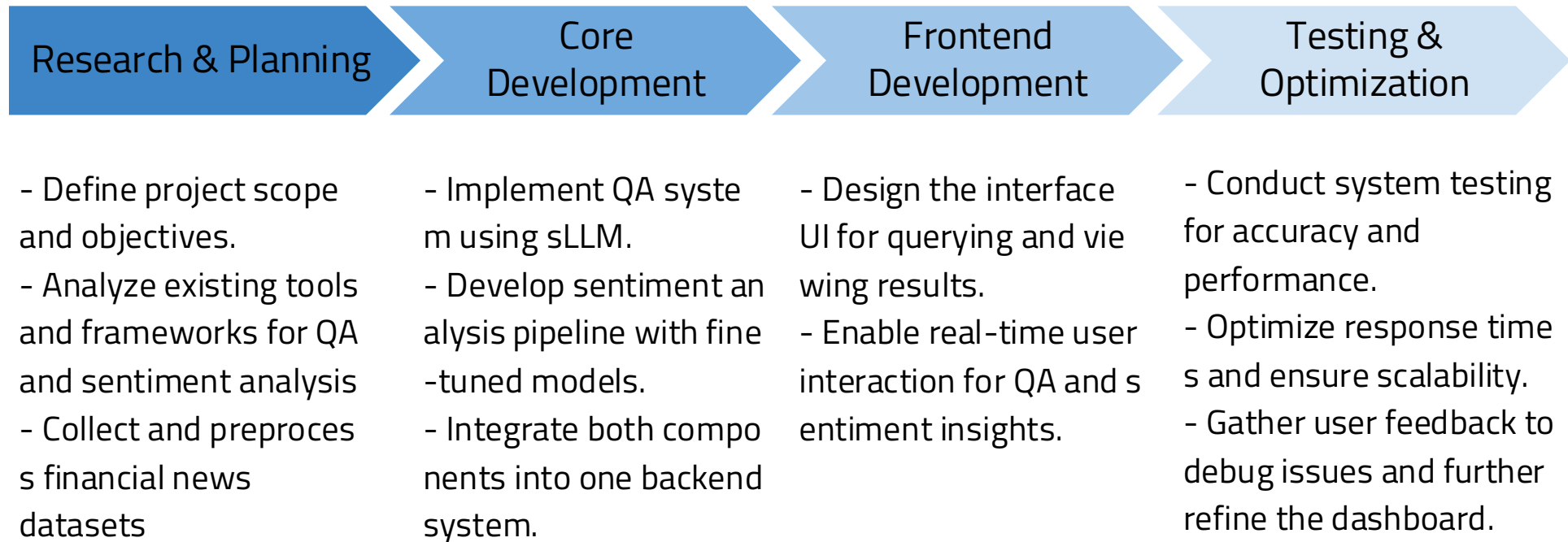**Master of Science in Information Management**
**Minseok Oh**

**KAIST**

https://nbviewer.org/github/whommso/ddaf/blob/main/Publication%20Thesis%20KAIST.pdf

| | Expected Values on the Continuous Intention to Use IoT Products from the Perspective of Expectation-Confirmation Theory | |
|---|---|---|
| 1 | **Objective** | • Identify gaps between user expectations and actual experience with IoT devices.<br>• Detect pain points and improvement areas to develop user-driven product strategies. |
| 2 | **Methodology** | • Structural Equation Modeling (SEM)<br>   ➢ Models latent variables to understand how different factors influence user satisfaction.<br>• Causal Analysis<br>   ➢ Examines independent-dependent variable relationships to establish causality between user expectations and actual usage experience.<br>• Exploratory Research<br>   ➢ Conducts data-driven hypothesis generation to uncover key insights on IoT adoption and usability. |
| 3 | **Data collection** | • User Data (Demographics & Behavior)<br>   ➢ Collected demographic and usage data:<br>     a) Basic Information: Sex, age, income, region, household size, marital status, children.<br>     b) IoT Experience: Prior exposure to IoT products, average daily usage<br>• Survey Data (User Experience & Perception)<br>     a) Designed 50 survey questions based on previous research, with 4 questions per key variable.<br>     b) Key Variables for Analysis:<br>       Perceived Manageability, Scalability, Entertainment Value, Reliability, Compensability, Expectation-Performance Alignment, Perceived Usefulness, Perceived Ease of Use, Social Influence, Satisfaction, Intention to Continue Use |
| 4 | **Exploratory Data Analysis** | • Reference: document page p.12~14<br>   ➢ https://nbviewer.org/github/whommso/ddaf/blob/main/Publication%20Thesis%20KAIST.pdf |
| 5 | **Research model** | • Reference: document page p.8<br>   ➢ https://nbviewer.org/github/whommso/ddaf/blob/main/Publication%20Thesis%20KAIST.pdf |

| | | Expected Values on the Continuous Intention to Use IoT Products from the Perspective of Expectation-Confirmation Theory |
|---|---|---|
| 5 | **Feature engineering & Assessment** | • Reliability Assessment:<br>  ➢ Cronbach's Alpha (CA) ≥ 0.7, CR, and AVE confirm internal consistency & validity.<br>• Correlation Analysis:<br>  ➢ All coefficients < 0.85 ensure construct distinctiveness.<br>  ➢ Diagonal AVE's square root > correlations confirms discriminant validity.<br>• Model Explanatory Power ($R^2$):<br>  ➢ $R^2 > 0.26$ indicates strong predictive power for dependent variables. |
| 6 | **Evaluation** | • Model Construction<br>  ➢ Defined path relationships between independent and dependent variables.<br>  ➢ Mapped survey scores to respective variables.<br>• PLS-SEM Execution<br>  ➢ Path coefficient calculation to quantify relationships.<br>  ➢ Bootstrapping (400 resamples) to test statistical significance.<br>  ➢ P-value analysis:<br>    a.  $p < 0.05$ → Reject $H_0$ (independent variable has a significant impact).<br>    b.  $p > 0.05$ → Fail to reject $H_1$ (no significant relationship found).<br>• Interpretation of Results<br>  ➢ If $H_0$ is rejected, e.g., Manageability significantly influences Agreement of Expectation.<br>  ➢ If $H_0$ is not rejected, the factor does not have a statistically significant effect. |
| 7 | **Result** | • Income Level: ≥50M KRW → Manageability & Compensability had no impact on satisfaction or continued use.<br>• Household Size:<br>  ➢ 1-2 person → Manageability & Entertainingness not significant.<br>  ➢ 3+ person → Significant impact on satisfaction & continued use.<br>• Regional Differences:<br>  ➢ Metro → Manageability was a dissatisfaction factor.<br>  ➢ Non-Metro → Compensability was valued for time & cost savings.<br>• Core Drivers:<br>  ➢ Perceived Usefulness → Significant for ≤50M KRW & 3+ person households.<br>  ➢ Perceived Ease → Strong positive impact on satisfaction & continued use. |

https://github.com/scottmsoh/ref_deep_leraning/tree/main/LLM_chatbot_SCU

| Research & Planning | Core Development | Frontend Development | Testing & Optimization |
|---|---|---|---|
| - Define project scope and objectives.<br>- Analyze existing tools and frameworks for QA and sentiment analysis<br>- Collect and preprocess financial news datasets | - Implement QA system using sLLM.<br>- Develop sentiment analysis pipeline with fine-tuned models.<br>- Integrate both components into one backend system. | - Design the interface UI for querying and viewing results.<br>- Enable real-time user interaction for QA and sentiment insights. | - Conduct system testing for accuracy and performance.<br>- Optimize response times and ensure scalability.<br>- Gather user feedback to debug issues and further refine the dashboard. |

End of Document