

“Doesn’t This Look Familiar?”

A Multimodal Approach to Finding Duplicate Products

Macy’s Team: Rohan Talati, Albert Egido De Poy, Jay Malepati
Faculty Advisor: Professor Alexandre Jacquillat



Jeremy Michael



Yolanda Wang

Project Background

519 stores across the USA and online E-commerce;

has

Launched a third-party marketplace in 2022, contributing approximately 4.2% to the total revenue of \$24 billion in 2023

Problem Statement



Macy's E-commerce

Macy's Products

Marketplace Products

Macy's expansion into the Marketplace risks potential **revenue loss** from including third-party products that **closely resemble** existing Macy's products.

Our solution **identifies these overlaps**, ensuring a unique and optimized product selection.



Increased Revenue

Improved Customer Experience

Challenges

Product nuances:

- ✓ silhouette, color, pattern, key elements
- brand, price, size, material



VS.



Diverse categories:

- Phase 1: Women's swimwear
 - Phase 2: Women's dresses
 - ⋮
 - Phase N: ALL product types
- Style diversity

Attribute variability:

- Duplicates vary in specifications, titles, and images, making simple matching unreliable

Data Limitation:

- Only five labeled pairs are provided, insufficient for building a supervised classification model

Data Sources

1.8M
Unique Products

20%
from Marketplace

425
Product Categories



Tabular

Category, color, style pattern



Text

Product titles, descriptions



Image

Front view, close-up details



Miscellaneous

Inventory, fabric, occasion

Methodology

Multimodal Data

- Tabular Data**
Category
Style List
- Text Data**
Product Title
Product Description
- Image Data**
Product Image
Swatch Image

Preprocessing

- Feature Selection**
450+ features → Top 10
- Text Cleaning**
Keep Only Meaningful Text
- Semantic Segmentation**
Isolate Clothing from Model and Background
- Color Extraction**
K-Means to Extract Color

Model Development

- Embedding Model**
Fashion-CLIP: OpenAI-based
Fine-Tuned for Fashion
- Dataset Creation**
Build Image Triplet Dataset – Anchor,
Positive, and Negative
- Fine-Tuning on Macy's Data**
Triplet Network Architecture for
Contrastive Learning

Similarity Matching

- Text Embeddings**
- Image Embeddings**
- Pairwise Cosine Similarity**
- Similarity Score (0 – 1)**

Candidate Ranking

- Filter Candidates – Category**
- Filter Candidates – Style**
- Filter Candidates – Color**
- Filter Candidates – Similarity Scores**
- Rank Top 10 Candidates – Weighted Sum**
- Human Elaboration**

Model Results

Precision@10

5% Baseline **↑** 93% Final Model

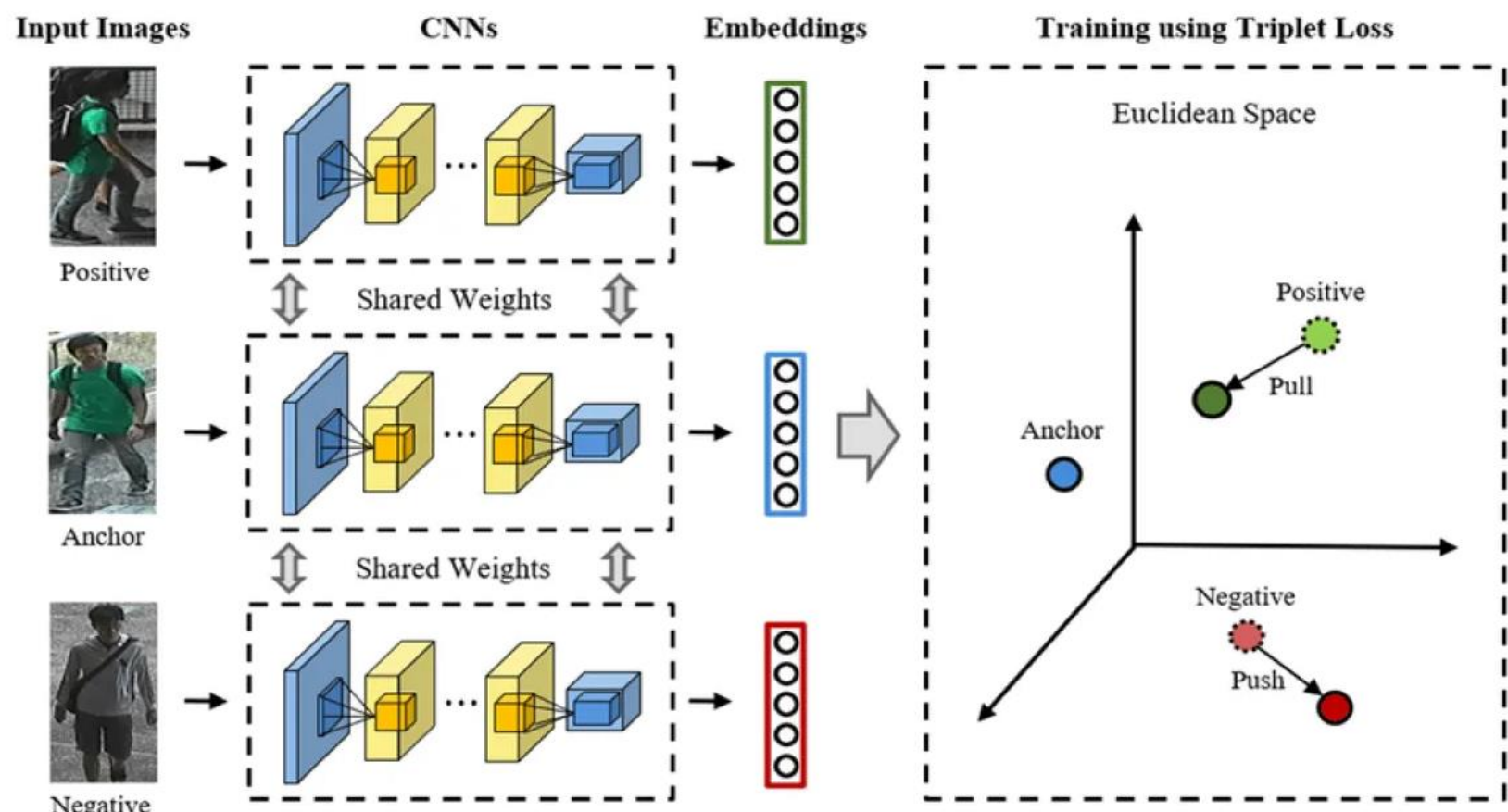
Identified Duplicates

↑ 550 %

Precision@10: Number of top 10 items matching at least 3 of silhouette, color, pattern, or key elements

Baseline: Filter candidates only by tabular features and ranked randomly

Fine-Tuning with Triplet Network



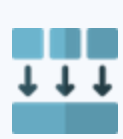
Example Similarity Search



The model returns the true duplicate product with the highest similarity score.

Business Impact

Multiple Use Cases



Group existing similar products together and categorize new products into clusters of similar items.



Store ranked matches with similarity scores in Macy's cloud and build a dashboard for labeling exact duplicates and non-duplicates.



Integrate model into product recommendation systems for enhanced relevance and diversity of suggestions.

By eliminating duplicate products across all categories, the estimated cost savings is up to **\$3.9M.**



Boosted the current number of identified duplicates by **550%.**