# Offline evaluation of product search with model-based judgments

Alberto Castelo
Sr. Applied Machine Learning Engineer
MICES 2024

shopify

# Agenda

- Search@Shopify
- Problem statement
- Solution: Model-based judgments
- Building a binary judgment model
- Conclusions

# Search@Shopify

# Shop

# Storefront - MrMaple Example

# Admin

shopify

tweeting-shirts    twe

- Home
- Orders                                    3
- Products
- Customers
- Content
- Finances
- Analytics
- Marketing
- Discounts

Last 30 days    All channels ⌄

Next payout: $0.00

| Online store sessions ✎ | Total sales | Total orders | Conversion rate | ⌃ |
|---|---|---|---|---|
| 0 — | $0.00 — | 0 — | 0% — | |

## No sessions in this date range
Try selecting a different date range or channel.

# App store

# Problem statement

# Problem statement

- We need offline evaluation
  - Release with confidence
  - Faster iteration cycles
- Implicit judgments approach is not good enough:
  - Judgment scarcity
  - Bias & noise
  - Not aligned with desired UX

# Implicit judgments - CTR Scarcity

- Scarcity tends towards reinforcing old good products

| Product Id | Judgment |
|------------|----------|
| 4634 | 0.23 |
| 2156 | 0.25 |
| 1234 | 0 or Mean |
| 7891 | 0 or Mean |
| 12945 | 0.12 |

# Implicit judgments - UX impact

# Solution: Model-based Judgments

# Solution

# Tasks

- Binary relevance

- Ranking

# Building a binary judgment model

# Model flywheel

- Build your golden eval dataset
  - **Manual annotation**
  - Representative of business objectives
- Set a training loop:
  - Train dataset
  - Model
  - Error analysis

# Building a training dataset

- Leverage public datasets
  - [ESCI (Amazon search)](#)
- Synthetic data
  - Distilling from GPT4 labels
  - Real data + synthetic
- Manual annotation
  - Clear guidelines

```
SYSTEM_PROMPT = """
You are an expert on ecommerce working on a product search engine. Your job is to:
1. Understand the product provided.
2. Generate an inexact query that:
 * Could happening during a ecommerce product search session.
 * where the product does not satisfy all the conditions or characteristics of the query.
 * has a similar structure than the exact query provided.
3. Answer in JSON format.

Use these definitions:
* Exact query: a query that matches the product characteristics.
* Inexact query: query that do not match the product (even if the product is close to match).

Examples:
1. Product: an Iphone 14 phone.
 a. Exact query: "iphone 14",
 b. Inexact query: "iphone 15"
2. Product: jordan's sneakers
 a. Exact queries: "nike"
 b. Inexact queries: "adidas"
3. Product: A bed for dogs
 a. Exact queries: "dog bed"
 b. Inexact queries: "cat bed"
"""


USER_PROMPT = """
Take a look at the product:
{product_dict}

In JSON format, generate a close yet inexact query with a similar structure as the exact query "{query_string}":
"""
```
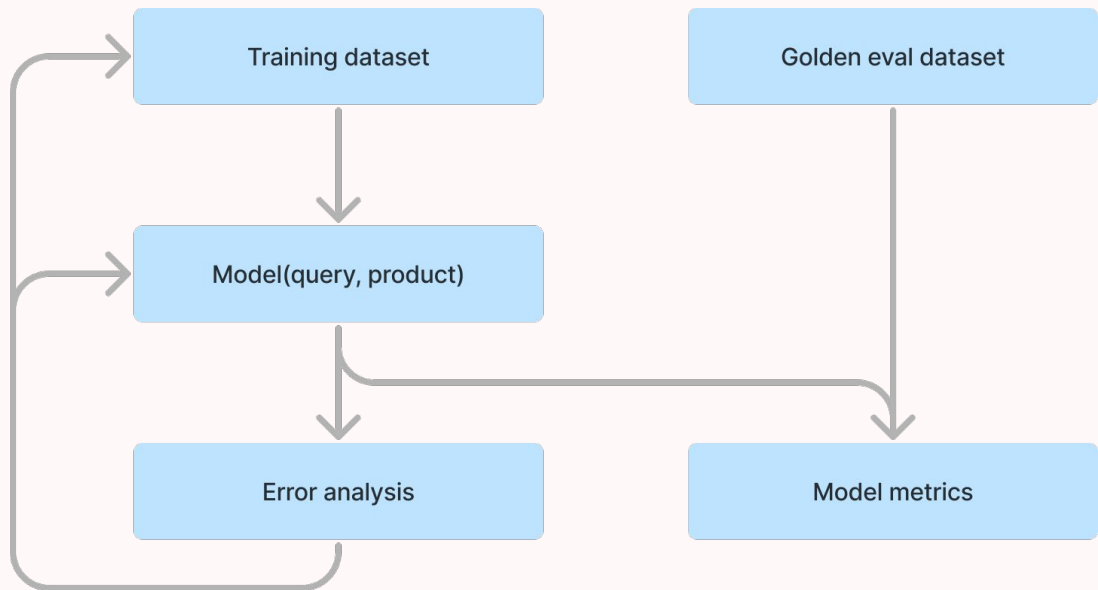
# Binary relevance model-based judgments

- LLM as a judge

- Classifier as a judge

# LLM-Llama 3

- 8B-Instruct
  - A couple of few-shot examples.

- 8B (un-instruct) LoRA finetune
  - O(100) samples

```
You are a relevance engineer working on ecommerce search.
Your job is to tag a pair of query and product as 'E' (exact match) or 'I'
(inexact match) depending on whether the product fully satisfies/matches the
intent of the query
### Query:
nike
### Product:
# title
jordan's
# vendor
nike
# shop name
shoes retailer
# product category
shoes
# product attributes
black
# image description
A man wearing a Nike shoe
# price
120
# description
Great Jordan's

### Response:
```

# Cross-encoders

## Bi-encoder (2 tower)

Cosine Similarity Score

Encoder model

Encoder model

tokenizer("{query}")

tokenizer("{product_text}")

## Cross-encoder

Label

Classification head

Encoder model

tokenizer("{query} <SEP><SEP> {product_text}")

# CLIP + crossencoder



CLIP + Cross-encoder

Label

Classification head

Text Encoder model

CLIP text encoder

CLIP image encoder

tokenizer("{query} <SEP><SEP> {product_text}")

tokenizer("{query}")

processor(product_image)

# Performance Summary

| Model | Data required | Inference Cost (1M pairs) | Classification Performance |
|---|---|---|---|
| Llama 3 8B-Instruct | O(1) | O($100) | Low |
| Llama 3 8B LoRA FT | O(100) | O($100) | High |
| Crossencoder | O(1k) | O($10) | High |

# SERP-Precision

- **Precision@1**: 0/1
- **Precision@4**: 0/4
- **Precision@8**: 3/8

# Winner/Loser comparison

- Comparing 2 search strategies: Precision@1

Query: ebike

2024_04_10_elasticsearch_multi_match

2024_04_11_elasticsearch_basic



Position: "0"
Title: "Dubbel Ebike"
Shop Name: "Blix Electric Bikes"
Predicted Category: "Sporting Goods > Outdoor Recreation > Cycling > Bicycles"
Relevance score: "0.98"

Position: "0"
Title: "eBike Cranks"
Shop Name: "20TwentyStore"
Predicted Category: "Sporting Goods > Outdoor Recreation > Cycling > Bicycle Parts > Bicycle Drivetrain Parts > Bicycle Cranks"
Relevance score: "0.17"

# Conclusions

# Conclusions

- Model-based judgments for the win
- Specialized models >> LLM
  - For now...
- Data is key

# Thank you!

| | |
|---|---|
| 𝕏 🐦 | **acaste10** |
| in | **alberto-castelo-becerra** |

shopify

# Offline evaluation of product search with model-based judgments

Alberto Castelo
Sr. Applied Machine Learning Engineer
MICES 2024

shopify