

Personalizing search using multimodal latent behavioral embeddings



Trey Grainger

Searchkernel

Founder / CTO Author, AI-Powered Search





Building the next generation of Search.



<u>Career</u>

Founder / CTO



Chief Technology Officer

Lucidworks

Chief Algorithms Officer SVP Engineering

Director of Engineering, Search & Recommendations

<u>Books</u>



Education



Trey Grainger Doug Turnbull Max Irwin

Al Powere

AI-Powered Search

- RAG (Retrieval Augmented Generation)
- Generative Search & Summarization
- Learning to Rank
- Semantic Search

 \bullet

- Dense Vector Search
- Fine-tuning LLMs for Search
- Personalized Search & Recommendations
- Knowledge Graph Learning
- User signals boosting & click models
 - Crowdsourced Relevance



(45% Discount Code: **ctwhaystack45**) https://aiPoweredSearch.com

Agenda

- The personalization spectrum between search and recommendations
- Leveraging a user's signals to implement collaborative filtering and personalization from latent features
- Using embedding vectors to generate personalization profiles
- Mixing user signals and content-based attributes to generate multimodal personalization
- Clustering products by embeddings to create personalization guardrails
- Avoiding the pitfalls of personalized search

What is an embedding?

Embeddings ("though vectors")

Word/Phrase Embeddings: [5, 1, 3, 4, 2, 1, 5, 3] [4, 1, 3, 0, 1, 1, 4, 2]

Sentence Embeddings: [2, 3, 2, 4, 2, 1, 5, 3] [5, 3, 2, 3, 4, 0, 3, 4]

Paragraph Embeddings: [5, 1, 4, 1, 0, 2, 4, 0] [1, 1, 4, 2, 1, 0, 0, 0]

Document Embedding: ____ [4, 1, 4, 2, 1, 2, 4, 3]

Definition of an embedding?



An embedding is a numerical representation of a piece of information, for example, text, documents, images, audio, etc.

https://huggingface.co/blog/getting-started-with-embeddings



embeddings are *vectors* created by machine learning models for the purpose of capturing meaningful data about each object

https://www.cloudflare.com/learning/ai/what-are-embeddings/

Google

An **embedding** is a relatively low-dimensional space into which you can translate high-dimensional vectors. Embeddings make it easier to do machine learning on large inputs like sparse vectors representing words. Ideally, an embedding captures some of the semantics of the input by placing semantically similar inputs close together in the embedding space.



Embedding is a means of representing objects like text, images and audio as points in a continuous vector space where the locations of those points in space are semantically meaningful

https://developers.google.com/machine-learning/crash-course/embeddings/video-lecture

An embedding is a set of coordinates in vector space into which we map a concept.

Embedding concepts into a vector space



An inverted index creates embeddings, with one dimension per term

	apple	caffeine	cheese	coffee	drink	donut	food	juice	pizza	tea	water	term N
latte	0	0	0	0	0	0	0	0	0	0	0	
cappuccino	0	0	0	0	0	0	0	0	0	0	0	
apple juice	1	0	0	0	0	0	0	1	0	0	0	
cheese pizza	0	0	1	0	0	0	0	0	1	0	0	
donut	0	0	0	0	0	1	0	0	0	0	0	
soda	0	0	0	0	0	0	0	0	0	0	0	
green tea	0	0	0	0	0	0	0	0	0	1	0	
water	0	0	0	0	0	0	0	0	0	0	1	
cheese bread sticks	0	0	1	0	0	0	0	0	0	0	0	
cinnamon sticks	0	0	0	0	0	0	0	0	0	0	0	

Embedding (verb) maps concepts into another vector space (usually a lower-dimensional space)



	food	drink	dairy	bread	caffeine	sweet	calories	healthy
apple juice	0	5	0	0	0	4	4	3
cappuccino	0	5	3	0	4	1	2	3
cheese bread sticks	5	0	4	5	0	1	4	2
cheese pizza	5	0	4	4	0	1	5	2
cinnamon bread sticks	5	0	1	5	0	3	4	2
donut	5	0	1	5	0	4	5	1
green tea	0	5	0	0	2	1	1	5
latte	0	5	4	0	4	1	3	3
soda	0	5	0	0	3	5	5	0
water	0	5	0	0	0	0	0	5

We then leverage the vector space to explore similarity

R

Phrase:

apple juice:	[0,5,0,0,0,4,4
cappuccino:	[0, 5, 3, 0, 4, 1, 2
cheese bread sticks:	[5, 0, 4, 5, 0, 1, 4
cheese pizza:	[5, 0, 4, 4, 0, 1, 5
cinnamon bread sticks:	[5, 0, 4, 5, 0, 1, 4
donut:	[5, 0, 1, 5, 0, 4, 5
green tea:	[0, 5, 0, 0, 2, 1, 1
latte:	[0, 5, 4, 0, 4, 1, 3
soda:	[0, 5, 0, 0, 3, 5, 5
water:	[0, 5, 0, 0, 0, 0, 0

Vector:

I, 3] 3 2 2 2 5 3] 0]), 5]

Vector Similarity (a, b): $\cos(\theta)$ = a·b a × b

Vector Similarity Scores:

anked	d Results: Cheese Pizza	Ranked Results: Green Tea			
99	cheese bread sticks	0.94	water		
91	cinnamon bread sticks	0.85	cappuccino		
39	donut	0.80	latte		
47	latte	0.78	apple juice		
46	apple juice	0.60	soda		
19	water	0.19	donut		

With LLMs and other Foundation Models, the dimensions learned are Latent Features^{*}.



Dimensions of User Intent



Personalization Spectrum

Traditional Keyword Search (Completely User-specified)

User-guided Recommendations

(Mostly driven by user profile, partially user-specified)

Personalized Search (Mostly user-specified,

partially driven by user profile)

Traditional Recommendations

(Completely driven by user profile)

Personalized Search

Keyword Search

(Completely User-specified)

User-guided Recommendations (Mostly driven by user profile, partially user-specified)

Personalized

Queries (Mostly user-specified, partially driven by user profile)

Collaborative Recommendations

(Completely driven by user behavior)

Recommendations Approaches

Content-based Matching (Content Filtering)



Recommendations Approaches

Behavior-based Matching (Collaborative Filtering)



Collaborative Recommendations

Recommendations Approaches

Content-based Matching (Content Filtering) + I

Behavior-based Matching (Collaborative Filtering)



Collaborative Filtering

•User 1:

•Avengers: Endgame •Black Panther •Black Widow

•User 2:

Black Widow
Captain Marvel
Black Panther

•User 3:

•Black Widow •The Dark Knight •The Batman •User 4: •The Little Mermaid •The Lion King •Toy Story

•User 5: •Frozen

•Toy Story •The Lion King

Users 1-3:

All of these are movies about superheroes
Most of them were made by Marvel Studios, though some were made by Warner Brothers (DC Comics)
They are all action movies
They are not suitable for small children due to violence and/or language

Users 4-5:

All of them are animated movies
All of them are suitable for small children
All of them are made by Disney/Pixar

Collaborative Filtering (Alternating Least Squares)

鱼 🕒 🔹 😓 Alternating Least Squares (Matrix Factorization based Collaborative Filtering)

from pyspark.ml.evaluation import RegressionEvaluator
from pyspark.ml.recommendation import ALS
from pyspark.sql import Row

(training, test) = indexed_prefs.randomSplit([0.8, 0.2], 0)

indexed_user_recs = model.recommendForAllUsers(10)
indexed_user_recs.show(5)

Results:

Root-mean-square error = 0.9589

use:	rIndex	recommendations
+	+	+
I	⊍ [{11,	0.014824464
1	1 [{5,	0.010805087}
1	2 [{7,	0.060263287}
1	3 [{1,	0.02967207},
1	4 [{14,	0.007906279
+	+	+
onlv	showing top	5 rows

User: u478462

Previous Searches: --apple --macbook

Previous Product Interactions:

-type: click, name: Apple[®] - iPad[®] 2 with Wi-Fi - 16GB - Black --type: add-to-cart, name: Apple[®] - iPad[®] 2 with Wi-Fi - 16GB - Black --type: purchase, name: Apple[®] - iPad[®] 2 with Wi-Fi - 16GB - Black <u>--type: click, name</u>: Apple® - MacBook® Air - Intel® Core™ i5 Processor ...

.

Alternating Least Squares (Matrix Factorization based Collaborative Filtering)

def get_query_time_boosts(user, boosts_collection): request = {"query": "", "return_fields": ["product", "boost"], "filters": [("user", user)] if user else [], "limit": 10, "order_by": [("boost", "desc")]}

response = boosts_collection.search(request) signals_boosts = response["docs"] return " ".join($f'''{b["product"]}''{b["boost"] * 100}'$ for b in signals_boosts)

def run_main_query(query, signals_boosts):

request = product_search_request(query if query else "") request["query_boosts"] = signals_boosts if signals_boosts else "*" return products_collection.search(request)

recs_collection = engine.get_collection("user_item_recommendations") user = "u478462"boosts = get_query_time_boosts(user, recs_collection)

response = run_main_query(None, boosts)

print(f"Boost Query:\n{boosts}") display_product_search("", response["docs"])

Boost Ouery:

"885909457588"^{75.393623} "097360810042^{*18.904798} "821793013776^{*15.852094} "610839379408^{*10.} 217768000000001 "635753493559"^9.087185 "885909395095"^8.304988 "885909457595"^7.917564000000 0005 "885909431618"^7.375394 "885909459858"^6.592548 "885909436002"^6.1031554

Recommendations:



9.7" widescreen display; 802.11a/b/g/n Wi-Fi; Bluetooth; iBooks support; measures just 0.34" thin and weighs only 1.35 lbs.

Search



Name: HTC - Flyer Tablet with 16GB Internal Memory - White | Manufacturer: HTC

Android 2.3 Gingerbread operating system7" color touch screenWi-Fi16GB memoryHTC notes, Watch and Listen apps



Name: Asus - Eee Pad Transformer Tablet with 16GB Storage Memory - Brown/Black | Manufacturer: Asus

Android 3.0 Honeycomb10.1" WXGA IPS touch screen Wi-Fi16GB hard drive



Name: Samsung - Galaxy Tab 10.1 - 16GB - Metallic Gray | Manufacturer: Samsung

Android 3.1 (Honeycomb) operating system10.1" WXGA touch screenWi-Fi



Name: Apple® - iPod touch® 32GB* MP3 Player (4th Generation - Latest Model) - Black | Manufacturer: Apple®



FaceTime camera, HD video recording, Retina display, Multi-Touch interface; gorgeous 3.5" widescreen display; Wi-Fi Web browsing



Name: Apple® - iPad® 2 with Wi-Fi - 32GB - Black | Manufacturer: Apple®

9.7" widescreen display; 802.11a/b/g/n Wi-Fi; Bluetooth; iBooks support

Output Search Results

query = "tablet"

```
response = run_main_query(
    query, None, 7
```

```
print(f"Non-personalized Query")
display_product_search(
   query, response["docs"]
```

Non-personaliz	zed Query	
tablet		Search
	Name: Init™ - Tablet Sleeve - Olive Manufacturer: Init™	
	Fits most tablets with up to a 10" display; heavy-duty neopr	ene material
	Name: Memorial Tablet - CD Manufacturer: Ltm/cd41	



Name: Stone Tablet - CD | Manufacturer: Important Records



Name: Sony - AC Power Adapter for Sony Tablet S | Manufacturer: Sony

Compatible with Sony Tablet S; charges your tablet

```
Personalized Search Results
query = "tablet"
recs_collection = engine.get_collection(
  "user_item_recommendations"
user = "u478462"
boosts = get_query_time_boosts(
  user, recs_collection
```

```
response = run_main_query(query, boosts, 7)
print(f"Personalized Query")
display_product_search(query, response["docs"])
```

```
Personalized Query
tablet
                Name: Apple® - iPad® 2 with Wi-Fi - 16GB - Black | Manufacturer: Apple®
```

9.7" widescreen display; 802.11a/b/g/n Wi-Fi; Bluetooth; iBooks support; measures just 0.34" thin and weighs only 1.35 lbs.



Name: HTC - Flyer Tablet with 16GB Internal Memory - White | Manufacturer: HTC

Android 2.3 Gingerbread operating system7" color touch screenWi-Fi16GB memoryHTC notes, Watch and Listen apps

Search



Name: Asus - Eee Pad Transformer Tablet with 16GB Storage Memory - Brown/Black | Manufacturer: Asus

Android 3.0 Honeycomb10.1" WXGA IPS touch screen Wi-Fi16GB hard drive



Name: Samsung - Galaxy Tab 10.1 - 16GB - Metallic Gray | Manufacturer: Samsung

Android 3.1 (Honeycomb) operating system10.1" WXGA touch screenWi-Fi



Name: Apple® - iPad® 2 with Wi-Fi - 32GB - Black | Manufacturer: Apple®

9.7" widescreen display; 802.11a/b/g/n Wi-Fi; Bluetooth; iBooks support

Name: Init[™] - Tablet Sleeve - Olive | Manufacturer: Init[™]

Fits most tablets with up to a 10" display; heavy-duty neoprene material



Name: Memorial Tablet - CD | Manufacturer: Ltm/cd41

Personalizing search using embeddings

Learned embeddings are latent features



User-Item Interaction Matrix



Matrix Factorization... learning latent features from user signals



User-Item Rankings (R)

User Feature Matrix (U)

Item Feature Matrix (I)

Learning embeddings of USERS allows us to personalize search results

Current search:

microwave

Previous Interactions:

VS.

<mark>Hello Kitty</mark> Plush Toy	GE Electric Razor (<mark>Black</mark>)			
GE Bright White	Samsung Stainless-steel			
ight Bulbs	Refrigerator			

No personalization guardrails:

GE - 0.7 Cu. Ft. Compact Microwave - Black a

Cook or reheat your meals in minutes with this compact 700-watt microwave oven that features instant-on LED controls which heat food with the simple touch of a button. It also provides peace of mind with a control lockout feature that locks the microwave when it's not in use.

Dept Appliance Brand GE Color Black

Hello Kitty - 0.7 Cu. Ft. Compact Microwave 0



Dept Photo/Commodities Brand Hello Kitty

GE - 1.1 Cu. Ft. Mid-Size Microwave - White ิด

Prepare hassle-free meals and snacks with this

With categorical personalization guardrails:



Samsung - Mid-Size Microwave - Stainless-steel 1

This mid-size microwave is large enough for almost any meal, but compact enough to fit into almost any space. Take an interactive product tour of this Samsung Microwave! (479KB Flash demo) Dept Appliance Brand Samsung Color Stainless-Steel



Samsung - 1.7 Cu. Ft. Over-the-Range Microwave - Stainless-Steel

0



The capabilities of a toaster oven, convection oven, range hood and microwave combine in this microwave that features a 1.7 cu. ft. capacity for simultaneously cooking multiple dishes and sensor cook options to take the guesswork out of guick meal prep.

Dept Appliance Brand Samsung Color Stainless-

Clustering products by embedding to generate dynamic categories

```
• • • • Clustering products to generate dynamic category contexts
```

```
def get_clusters(data, algorithm, args, kwds):
    return algorithm(*args, **kwds).fit(data)
```

```
def assign_clusters(labels, product_names):
    clusters = defaultdict(lambda:[], {})
    for i in range(0, len(labels)):
        clusters[labels[i]].append(product_names[i])
    return clusters
```

#<1> Generate 100 clusters using a KMeans clustering algorithms #<2> Assign each product name to its corresponding cluster label

Visualizing the clusters

0.4

0.3

0.2

0.0

-0.1

-0.2

-0.3

-0.3

39_Range_Self-Cleaning

6 Cu. Ft.

89_98WEF_PRWEEtos 4_Case_Black

-0.2

-0.1

40_Vacuum_Bagless

95 Razor_Black

....

65 Case Mobile

88 Laptop Case

21 30" Range 76 Microwave Cu.

33 Ft. Cu.

47 Vehicles Select

9 Headphones Earbud

68-Apple0#RhenCase

90 Charger Adapter

0.1 -58_Air_White --

• • •

Visualize the category contexts

import collections, numpy as np, matplotlib.pyplot as plt from adjustText import adjust_text from sklearn.decomposition import PCA

plt.figure(figsize=(15, 15)) pca = PCA(100, svd_solver="full") centers = algo.cluster_centers_ plot_data = pca.fit_transform(centers) #<1>

```
points = []
for i, cluster_name in enumerate(plot_data): #<2>
 plt.scatter(plot_data[i,0], plot_data[i,1], #<3>
              s=30, color="k") (2)
 label = f"{i}_{"_".join(top_words(clusters[i], 2))}"
 points.append(plt.text(plot_data[i,0], #<3>
                         plot_data[i,1], #<3>
                         label, size=12)) #<3>
adjust_text(points, arrowprops=dict(arrowstyle="-", #<4>
                           color="gray", alpha=.3)) #<4>
```

plt.show()

#<1> Perform PCA to reduce to 2 dimensions for visualization *#<4> Display improvement: adjust the text labels to* minimize overlap

80 Lens Cameras 32 Cable Monster 74_Wireless-N_Switchayer_MP3ystem_Camera 23 Windows Mac/Windows 82_Windows_Mac/Windows 99 Card Prepaid 89 Mouse Wireless _36_DVD_Disc 91_USB_Drive 54 Blu-ray_Disc 75 Adobe Windows 50_Drive_Hard 71_Memory_Card 2_Home_Theater 41 Camera Digital 26 TVs Flat-Panel 83_DVD_WWE: 77 Card PCI 46_Memory_Processor 17_Laptop_Memory 42 Camcorder HD •70 Monitor Widescreen 49_Player_3D

0.0

61_Guitar_Electric

27 Stainless-Steel Black

92 Amplifier Class

63 Subwoofer 10

34 M-Edge Accessories

45_Cartridge_Ink

67 Keyboard Wireless

98 Subtitle DVD 66_DVD_Widescreen _ 1_AC3_Subtitle 59_HDTV Class 97 AC3 Widescreen 12 Widescreen DVD

0.1

22 Deck Apple®

60 GPS Garmin

48_Camera_Black

37 Black Stand

28 Black Chair

87 CD [PA]

94 CD You

11 VINYL [LP]

7 CD Various

8 CD Greatest

55_CD_Original

24_DVD_Live_

52 DVD Widescreen

81_DVD_Discs]

31 DVD Discs]

62_DVD ¶ullscreen

13 DVD Discs]

93_Season_DVD

0.4

18 DVD Season

53 DVD Workout

14 CD Best

20 CD [Digipak

3 CD Various

•72_CD_[PA]

51 CD de

16 (Pair) Speakers

35 Watch Monitor

.

79 PlavStation PSP

0.2

38 Microphone USB

78 Light Lighting

19 Nintendo Wii

29 DVD Discs]•

0.3

86 DVD Season

56_Dolby_Widescreen

44_Widescreen_DVD

84_CD [ECD]

30 CD Best

• 73 CD [CD

Mapping queries into clusters (multiple approaches)

•••

🗧 🥏 Different approaches for mapping queries into clusters

import sentence_transformer, heapq

def get_top_labels_centers(query, centers, n=2): #<4>
 q_emb = transformer.encode([query], convert_to_tensor=False)
 similarities = sentence_transformers.util.cos_sim(q_emb, centers)
 sim = similarities.tolist()[0]
 return [sim.index(i) for i in heapq.nlargest(n, sim)]

- def get_query_cluster(query): #<5>
 q_emb = transformer.encode([query], convert_to_tensor=False)
 return algo.predict(q_emb)
- def get_cluster_description(cluser_num):
 return "_".join(top_words(clusters[cluser_num], 5))

query = "microwave"
kmeans_predict = get_query_cluster(query)[0] #<1>
print("KMeans Predicted Cluster:")
print(f" {kmeans_predict} ({get_cluster_description(kmeans_predict)})")

closest_sim = get_top_labels_centers(query, centers, 1)[0] #<2>
print(f"\nCosine Predicted Cluster:")
print(f" {closest_sim} ({get_cluster_description(closest_sim)})")

knn_cosine_similarity = get_top_labels_centers(query, centers, 5) #<3>
print(f"\nKNN Cosine Predicted Clusters: {knn_cosine_similarity}")
for n in knn_cosine_similarity:

print(f" {n} ({get_cluster_description(n)})")

#<1> Option 1: Predict nearest cluster (KMeans)
#<2> Option 2: Find most-similar cluster (Cosine similarity)
#<3> Option 3 (recommended): Find N-most-similar clusters (Cosine similarity)
#<4> Get the top N clusters based on cosine similarity with the cluster centry
#<5> Get the cluster based on the KMeans model's prediction

Results:

KMeans Predicted Cluster:

76 (Microwave_Cu._Ft._Stainless-Steel_Oven)

Cosine Predicted Cluster:

76 Microwave_Cu._Ft._Stainless-Steel_Oven)

KNN Cosine Predicted Clusters: [76, 27, 21, 39, 33]

- 76 (Microwave_Cu._Ft._Stainless-Steel_Oven)
- 27 (Stainless-Steel_Black_KitchenAid_Cuisinart_Maker)
- 21 (30"_Range_Stainless-Steel_Gas_Cooktop)
- 39 (Range_Self-Cleaning_30"_Freestanding_Stainless-Steel)
- **33** (Ft._Cu._Refrigerator_Water_Thru-the-Door)

Approaches for integrating embedding-based personalization into search results

• Perform a weighted average between the query vector (embedding for microwave) and the vectors for the user's previous interactions within the predicted clusters. This would generate a single vector representing a personalized version of the user's query, so all results would be personalized.

• Perform a standard search, but then boost the results based on the average of the embeddings from a user's previous interactions within the predicted clusters. This would be a hybrid keyword and vector-based ranking function, where the keyword search would be the primary driver of the results, but the user's previous interactions would be used to boost related results higher.

• Do one of the above, but then only personalize a few items in the search results instead of all the results. This follows a light-touch mentality so as not to disturb all of the user's search results, while still injecting novelty to enable the user to discover personalized items they may not have otherwise found.

• Perform a standard search (keyword or vector), but then re-rank the results based on the weighted average between the query vector and the vectors for the user's previous interactions within the predicted clusters. This uses the original search to find the candidate results using the default relevance algorithm, but then those results are re-ranked to boost personalized preferences higher.

Generating personalization vectors for a user's query

. Generating a personalization vector for the user's query import pandas, numpy def get_user_embeddings(products=[]): #<1> values = [] for p in products: values.append([product_ids_emb[p], top_clusters_for_embedding(product_ids_emb[p], 1)[0]]) column_names = ["embedding", "cluster"] return pandas.DataFrame(*data*=numpy.array(values), *index*=products, columns=column_names) def get_personalization_vector(query=None, #<2> user_items=[], query_weight=1, #<3> user_items_weights=[]): #<4> query_embedding = transformer.encode(query) if query else None if *len*(user_items) > 0 and *len*(user_items_weights) = 0: #<3> user_items_weights = numpy.full(shape=len(user_items), fill_value=1 / len(user_items)) embeddings = [] embedding_weights = []

```
for weight in user_items_weights: #<4>
    embedding_weights.append(weight) #<4>
for embedding in user_items:
    embeddings.append(embedding)
if query_embedding.any():
    embedding_weights.append(query_weight) #<4>
    embedding.append(query_embedding)
```


#<1> Returns a dataframe with the embedding and guardrail cluster for each product #<2> Returns a vector that combines (weighted average) an embedding for the query # with the embeddings for the passed-in user_items #<3> By default, the weight is split 1:1 (50% each) between the query embedding # and the user_items_weight #<4> You can optionally specify a query_weight and user_item_weights to influence

how much each embedding influences the personalization vector

Generating contextual user personalization vectors

• • • Contextual (filtered) vs. Non-contextual (unfiltered) personalization vectors

user_embeddings = get_user_embeddings(product_interests)
query = "microwave"

unfiltered_personalization_vector = get_personalization_vector(query=query, user_items=user_embeddings['embedding'].to_numpy()) #<1> print("\nPersonalization Vector (No Cluster Guardrails):") print(format_vector(unfiltered_personalization_vector))

clustered = user_embeddings.cluster.isin(query_clusters) #<3>
products_in_cluster = user_embeddings[clustered] #<3>
print("\nProducts Filtered to Query Clusters:\n" + str(products_in_cluster))

filtered_personalization_vector = get_personalization_vector(query=query, user_items=filtered['embedding'].to_numpy()) #<4> print("\nFiltered Personalization Vector (With Cluster Guardrails):") print(format_vector(filtered_personalization_vector))

#<1> Personalization vector with no guardrails
(uses query and all past item interactions)
#<2> Get the top 5 clusters for the query to use as guardrails
#<3> Filter down to only items in the guardrail query clusters
#<4> Generate a personalization vector with guardrails
(uses query and only items related to the query)

Results:

 Products Interactions for Personalization:
 cluster

 product
 embedding
 cluster

 7610465823828
 [0.06417941, 0.04178553, -0.0017139615, -0.020...]
 28

 36725569478
 [0.0055417763, -0.024302201, -0.024139373, -0...]
 39

Personalization Vector (No Cluster Guardrails): [0.016, -0.006, -0.02, -0.032, -0.016, 0.008, -0.0, 0.017, 0.011, 0.007,...]

Query Clusters ('microwave'): [76, 27, 21, 39, 33]

Products Filtered to Query Clusters: product embedding cluster 36725569478 [0.0055417763, -0.024302201, -0.024139373, -0....] 39

Filtered Personalization Vector (With Cluster Guardrails): [0.002, -0.023, -0.026, -0.037, -0.025, 0.002, -0.009, 0.007, 0.033, -0....]

Run different personalization scenarios

.

Run Different Personalization Scenarios

```
def rerank_with_personalization(docs, personalization_vector):
  result_embeddings = numpy.array([product_ids_emb[docs[x]["upc"]]
                                  for x in range(len(docs))]).astype(float)
  similarities = sentence_transformers.util.cos_sim(
      personalization_vector, result_embeddings).tolist()[0]
  reranked = [similarities.index(i)
              for i in heapq.nlargest(len(similarities), similarities)]
 reranked, v = zip(sorted(enumerate(similarities),
                            key=itemgetter(1), reverse=True))
  return [docs[i] for i in reranked]
query = "microwave"
request = {"query": query,
           "query_fields": ["name", "manufacturer"],
           "return_fields": ["upc", "name", "manufacturer", "score"],
           "limit": 100,
           "order_by": [("score", "desc"), ("upc", "asc")]}
response = products_collection.search(*request)
docs = response["docs"]
print("No Personalization:")
display_product_search(query, docs[0:4])
print("Global Personalization (no category guardrails):")
reranked_seach_results_no_guardrails = \
  rerank_with_personalization(docs,
   unfiltered_personalization_vector)
display_product_search(query, reranked_seach_results_no_guardrails[0:4])
print("Contextual Personalization (with category guardrails):")
reranked_seach_results_with_guardrails = \
  rerank_with_personalization(docs,
   filtered_personalization_vector)
display_product_search(query, reranked_seach_results_with_guardrails[0:4])
```

No Personalization

microwave	Search
111-1111-1111-111-111-11-11-11-11-11-11	
Name: LG - LMHM2017SW Micro Manufacturer: LG	owave Oven - White



Name: Panasonic - NNCD989S Microwave Oven - Stainless-Steel Manufacturer: Panasonic

Name: West Bend - Microwave Popcorn Popper - Red Manufacturer: West Bend



Name: Hello Kitty - 0.7 Cu. Ft. Compact Microwave Manufacturer: Hello Kitty

Global Personalization (No Category Guardrails)





Name: Electrolux - 1.5 Cu. Ft. Built-In Microwave - Stainless-Steel | Manufacturer: Electrolux



Name: Panasonic - 1.2 Cu. Ft. Mid-Size Microwave - Stainless-Steel | Manufacturer: Panasonic



Name: Panasonic - 1.2 Cu. Ft. Mid-Size Microwave - White | Manufacturer: Panasonic

Contextual Personalization (with Category Guardrails)

microwave	Search	
Name: Electrolux - 1.5 Cu. Ft. Built-In Manufacturer: Electrolux	Microwave - Stainless-Stee	1
Name: Panasonic - 1.2 Cu. Ft. Mid-Siz	ze Microwave - Stainless-	

Steel | Manufacturer: Panasonic



Name: Panasonic - 2.2 Cu. Ft. Full-Size Microwave - Stainless-Steel | Manufacturer: Panasonic



Name: Panasonic - 2.2 Cu. Ft. Full-Size Microwave - Stainless-Steel | Manufacturer: Panasonic

Multimodal models (language + behavior) can integrate these latent features to discover new insights and enhance relevance



Avengers: Endgame (Movie)	Black Panther (Movie)	Notting Hill (Movie)	The Notebook (Movie)	Minecra [:] (Game)	The ft Bachelo (TV Sho	or w)
0.09	0.71	1.91	1.74	1.80	0.28	
0.75	0.75	2.17	2.54	-1.17	3.02	
3.43	2.21	0.01	0.46	1.89	0.30	



Calculating Preferences from Latent Factors:

Avengers: Endgame (Movie)



The Notebook (Movie)

9.9

3 * 1.74) + (3.18 * 2.54) + (-0.13 * 0.46) =



It's also powerful to use many *query modalities* to optimize search quality



Reflected Intelligence: So many ways to integrate user signals...



Signal Processing & Machine Learning

Trey Grainger Doug Turnbull Max Irwin

Al Powe

MANNING

AI-Powered Search

- RAG (Retrieval Augmented Generation)
- Generative Search & Summarization
- Learning to Rank
- Semantic Search
- Dense Vector Search
- Fine-tuning LLMs for Search
- Personalized Search & Recommendations
- Knowledge Graph Learning
- User signals boosting & click models
- Crowdsourced Relevance

Get a copy @ http://aiPoweredSearch.com

(35% Discount Code: **ctwhays**tack45)

Thank You!

Trey Grainger



trey@searchkernel.com @treygrainger



Other presentations:

http://treygrainger.com

(45% Discount Code: ctwhaystack45)



Books:

http://aiPoweredSearch.com http://solrinaction.com