

Personalization and Recommendation

TECHNOLOGY

Google Knows You Better Than You Know Yourself

Predictive analysis combs through calendars and search histories—and gets in the way of routine self-deception.

JAMES CARMICHAEL AUGUST 19, 2014

Facebook Knows You Better than You Know Yourself



Erman Misirlisoy, PhD Oct 18, 2018 · 7 min read ★



The Internet Knows You Better Than You Know Yourself

When Amazon or eBay recommend us something we like but were not looking for, they effectively know us better than we know ourselves.



Netflix: How did it know I was bi before I did?

After BBC reporter Ellie House came out as bisexual, she realised that Netflix already seemed to know. How did that happen?

How Target Figured Out A Teen Girl Was Pregnant Before Her Father Did

Kashmir Hill Former Staff

Welcome to The Not-So Private Parts where technology & privacy collide

What types of data do firms collect?

Data that firms collect

Search history: Almost every search engine collects data on your search histories. Whenever you search on Google, it knows what you are looking for, the items you viewed on Google, the webpages you clicked...

Data that firms collect

Geolocation and device information: A GPS and Wi-Fi chip are installed in every smartphone. The dozens of apps on our phones, most of them free, aren't just serving up information and entertainment. They are collecting and selling your data to digital marketers who will then offer you personalized ads.

Data that firms collect

Purchase histories: Online sellers frequently collect data on your purchase histories. They know what you have purchased and what you have not purchased. As we will figure out later, the purchase history data is the most important data for online sellers.

Data that firms collect

IP addresses: You have to get an IP address to search the internet. You can find your IP address [here](#). Then, even if you are not using a mobile device, firms can still know your geolocation (e.g., your country and even neighborhood).

Google also collects...

Number of email exchanges you've had in Gmail; number of files in Drive; number of photos Google stores for you.

Your location or searches or browsing history: Google Maps keeps track of everywhere you go and when, alongside the photos taken that day and travel times down to the minute.

Your Google Account: your photo and birthdate.

How about traditional firms?

Even the traditional brick-and-mortar (offline) shops are also collecting your data.

- **Your payment method** (Credit? Mobile pay? Cash?)
- **Loyalty program information** (Are you using Yuu?)
- **Personal profile** (If you ever registered there...)

How about traditional firms?

With new technologies, brick-and-mortar stores can also get much more information than what they had before.

- As described in the video, if you use the free Wi-Fi they provide you, they will be able to collect data from your smartphone!
- Facial recognition and mobile payments help collect data from you.

Suppose that you are an Internet company, and you have access to all this consumer data, what would you do?

Personalized Pricing

With personalized pricing, a seller offers each consumer an individualized price, and two persons can receive two different prices at the same time.

Note that personalized pricing is different from dynamic pricing. With dynamic pricing, the price is changing over time. For personalized pricing, the price is changing over consumers.

Example of dynamic pricing: **Uber adjusts prices timely.**

Price Discrimination

Broadly speaking, personalized pricing is a form of price discrimination. Let's review types of price discrimination:

- **1st degree:** The firm sells a product at the maximum price that every consumer is willing to pay.
- **2nd degree:** price varies according to quantity demanded.
- **3rd degree:** charging a different price to different consumer groups.

Price Discrimination

Personalized pricing is close to first-degree price discrimination.

Firms can learn about your income (e.g., from your bank account), your geo-location (e.g., in the US or India), your neighborhood (a high-end one?), your device (iOS or Android), your purchase habits (bargain hunter?), your gender,...

Based on this information, firms can infer how much you are willing to pay for the product and offer you a personalized price.

THE WALL STREET JOURNAL.

[English Edition ▼](#) | [Print Edition](#) | [Video](#) | [Podcasts](#) | [Latest Headlines](#)

[Home](#) [World](#) [U.S.](#) [Politics](#) [Economy](#) [Business](#) **[Tech](#)** [Markets](#) [Opinion](#) [Life & Arts](#) [Real Estate](#) [WSJ.M](#)

On Orbitz, Mac Users Steered to Pricier Hotels

Are you using a Mac or PC?

On Orbitz, Mac users spend as much as 30% more a night on hotels that PC users do.

Websites Vary Prices, Deals Based on Users' Information

Getting Different Deals Online

A Journal examination found online retailers adjusted prices by a shopper's location, among other factors

Staples.com

SnapSafe Titan safe

HIGHER PRICE
\$1,199.99

DISCOUNT PRICE
\$1,099.99

DIFFERENCE:
9.1%



Homedepot.com

A 250-foot spool of electrical wiring



Six pricing groups, including:

\$70.80 in Ashtabula, Ohio

\$72.45 in Erie, Pa.

\$77.87 in Monticello, NY

RosettaStone.com



...for buying multiple levels of German lessons, when test-shopping from the U.S. or Canada. But not from the U.K. or Argentina.

Photos: (l to r) SnapSafe; Home Depot; Rosetta Stone

Source: WSJ testing

The Wall Street Journal

MOST POPULAR NEWS

1. What You Can and Can't Do if You've Been Vaccinated: Travel, Risk Factors, What You Need to Know
2. Europe Confronts Covid Rebound as Vaccine Hopes Recede
3. Biden's \$1,400 Stimulus Checks Hit Bank Accounts Starting Today
4. Schumer and Gillibrand Call for Cuomo to Resign

The US retailer Office Depots uses customers' browsing history and location data to vary prices

These Brands Have Some of the Best Abandoned Cart Email Strategies

Aug 28, 2019 5:03:58 PM

When you abandon an item from your online shopping cart, e-tailers may issue you a discount to lure you to make a purchase.

Behavior-Based Pricing

The more common approach is pricing with consumers' purchase history, a practice known as “behavior-based pricing.”

The idea is very simple: The price you receive depends on whether or not you have purchased the products before. In other words, we offer new and existing consumers different prices.

As consumers, do you like behavior-based pricing? Why?

Recommendation is everywhere!



Recommended for You

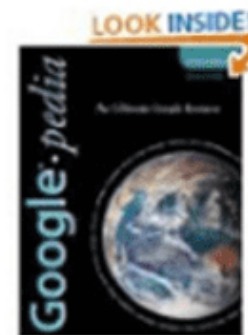
Amazon.com has new recommendations for you based on [items](#) you purchased or told us you own.



[Google Apps
Deciphered: Compute in
the Cloud to Streamline
Your Desktop](#)



[Google Apps
Administrator Guide: A
Private-Label Web
Workspace](#)



[Googlepedia: The
Ultimate Google
Resource \(3rd Edition\)](#)



Arizona Border Ranchers Torn in Support for Trump's Wall

172,275 views

683

249

SHARE

SAVE

...

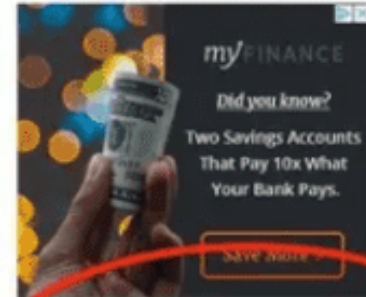


Wall Street Journal
Published on Mar 16, 2017

Despite enthusiastic backing for President Donald Trump and pleas for a stronger border, Arizona ranchers are conflicted in their support for Trump's promise to build a wall along the border with Mexico. Photo/Video: Jake Nicol/The Wall Street Journal

SHOW MORE

SUBSCRIBE 1.2M



Up next

AUTOPLAY



(Part II) A Day in the Life of Arizona Rancher: Fences, II
Center for Immigration Studies
43K views



CNN
You promised Mexico would...
2.7M views
New



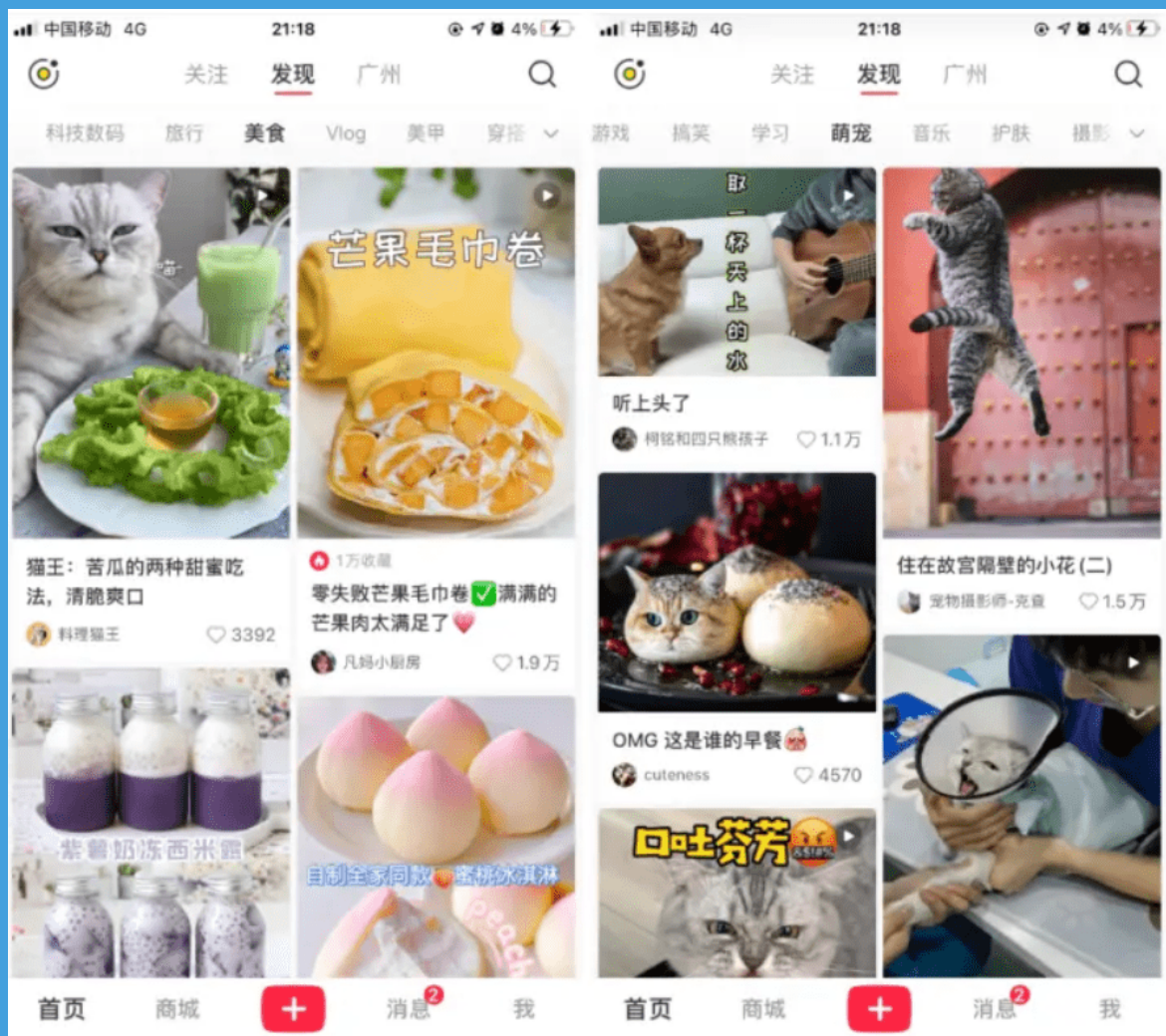
NBC News
People Are Fleeing President Trump's America To This T...
308K views



Sky News
Scrambling onto trucks for better life
2.4M views



BBC Planet Earth | BBC Studios
Polar Bear vs Walrus colony
Recommended for you



NETFLIX

Home TV Shows Movies Latest My List



KIDS



Continue Watching for SmartTV



Trending Now



Korean TV Shows



The Importance of Recommendation

- Netflix: 2 / 3 of the movies watched are recommended.
- Google News: recommendations generate 38% more click-throughs.
- Amazon: 35% sales from recommendations.
- ChoiceStream: 28% of the people would buy more music if they found what they liked.

How to recommend?

A recommendation system must have three inputs:

- **Items** to be recommended: songs, movies, products, restaurants etc. (often many thousands)
- **Users** of the items: watchers, listeners, purchasers, shoppers etc. (often many millions)
- **Feedback** of users on items: 5-star ratings, upvotes/downvotes, clicking “next” or “skipping the ad”, purchases or clicks.

Collaboratives Filtering

Collaborative filtering is not something new. We have done it in many places in the past. Here are a few examples:

- Bestseller list for books
- Top 50 music list
- The “recent returns” shelf at libraries

The intuition behind: People's tastes are correlated.

Collaboratives Filtering

However, in the above examples, recommendations are not personalized, i.e., everybody receives the same recommendation. How to make recommendations personalized?

The intuition: If Alice and Bob both like X and Alice also likes Y , then Bob is more likely to like Y , especially when Alice and Bob know each other.

Suppose that you want to recommend a movie to Emma,
which movie will you recommend?

	1	2	3	4	5	6
Alice	2			4	5	
Bob	5		4			1
Carol			5		2	
Dennis		1		5		4
Emma			4			2
Flora	4	5		1		

User-Based Collaborative Filtering: The Neighbourhood Method

Step 1: Find all the movies rated by Emma before, we get
movies 3 and 6

	1	2	3	4	5	6
Alice	2			4	5	
Bob	5		4			1
Carol			5		2	
Dennis		1		5		4
Emma			4			2
Flora	4	5		1		

Step 2: Identify other users that have rated the same movie,
we get Bob, Carol, and Dennis

	1	2	3	4	5	6
Alice	2			4	5	
Bob	5		4			1
Carol			5		2	
Dennis		1		5		4
Emma			4			2
Flora	4	5		1		

Step 3: Compare the similarity between Emma and her “neighbors” to see who are close to Emma.

	1	2	3	4	5	6
Alice	2			4	5	
Bob	5		4			1
Carol			5		2	
Dennis		1		5		4
Emma			4			2
Flora	4	5		1		

Step 4: Select the top k most similar neighbors and use their average ratings to predict Emma's rating.

	1	2	3	4	5	6
Alice	2			4	5	
Bob	5		4			1
Carol			5		2	
Dennis		1		5		4
Emma			4			2
Flora	4	5		1		

Item-Based Collaborative Filtering

Item-based collaborative filtering

Suppose that we are predicting the who will like movie 5.

Step 1: Who have rated movie 5 before? We get Alice and Carol.

	1	2	3	4	5	6
Alice	2			4	5	
Bob	5		4			1
Carol			5		2	
Dennis		1		5		4
Emma			4			2
Flora	4	5		1		

Step 2: Identify other movies that have rated the same users, we get movies 1 and 3.

	1	2	3	4	5	6
Alice	2			4	5	
Bob	5		4			1
Carol			5		2	
Dennis		1		5		4
Emma			4			2
Flora	4	5		1		

Step 3: Compare the similarity between movie 5 and its “neighbors” to see which movie is close to movie 5.

	1	2	3	4	5	6
Alice	2			4	5	
Bob	5		4			1
Carol			5		2	
Dennis		1		5		4
Emma			4			2
Flora	4	5		1		

Step 4: Select the top k most similar neighbors and use their average ratings to predict movie 5's rating.

	1	2	3	4	5	6
Alice	2			4	5	
Bob	5		4			1
Carol			5		2	
Dennis		1		5		4
Emma			4			2
Flora	4	5		1		

Model-based Collaborative Filtering

What did Netflix do to make recommendations?

In general, how much do you like watching movies from the following genres?						
	Really dislike	Dislike	Neither like nor dislike	Like	Really like	Not sure of genre definition
Action	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Adventure	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Animation	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>
Comedy	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>
Crime/Gangster	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Documentary	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Drama	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Fantasy	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Film-Noir	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Foreign	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Horror	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

The Matrix Factorization Method

Matrix Factorization

Here, we assume that each movie has a number of “latent” or “hidden” factors that affect user preferences. Examples of the factors include the length of the movie, the amount of actions in the movie, the seriousness of the movie, the orientation of the movie for children etc.

The factors are “latent” or “hidden,” implying that we do not know which these factors are, and we do not need to know either.

Matrix Factorization

Each user also has his or her own preference for each factor. For instance, some users prefer long movies over short movies, and some users prefer to have more actions in their movie. If we know the preferences of a user and a movie's attribute values, we can match the movie with the user to see whether the user will like the movie.

Matrix Factorization

For instance, suppose that a movie is long and contains a lot of actions. We also know that

- Alice likes short movies and hates action movies,
- Bob prefers long movies and enjoys action movies.

Then, we can predict that Alice will hate the movie and Bob will like the movie.

Mathematically, our model is as follows:

$$\begin{aligned} \text{Your rating} = & \text{Your preference for length} \times \text{Movie's length} \\ & + \text{Your preference for action} \times \text{Movie's amount of action} \end{aligned}$$

Suppose that the movie's length is 1 and amount of action is 2. Alice's preference for length is 0.5, for action is 0; Bob's preference for length is 1, for action is 1.5. We can predict:

- Alice's rating: $0.5 \times 1 + 0 \times 2 = 0.5$.
- Bob's rating: $1 \times 1 + 1.5 \times 2 = 4$.

Let's generalize the above discussion. Suppose that Alice, Bob, Carol's preferences for length and action are as follows:

	Length	Action
Alice	0.5	0
Bob	1	1.5
Carol	1.5	0.5

There are two movies, whose length and action values are

	Length	Action
1	1	2
2	0	3

We can multiply the two matrix to get user-movie ratings:

$$\begin{bmatrix} 0.5 & 0 \\ 1 & 1.5 \\ 1.5 & 0.5 \end{bmatrix} \times \begin{bmatrix} 1 & 0 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} 0.5 & 0 \\ 4 & 4.5 \\ 2.5 & 1.5 \end{bmatrix}$$

In other words, our prediction is as follows.

	Movie 1	Movie 2
Alice	0.5	0.0
Bob	4.0	4.5
Carol	2.5	1.5

Overall, if we know the user matrix and the movie matrix, we can multiply the two to get the user-movie rating matrix. The issue is: **We do not yet know the user matrix and the movie matrix.**

But how to get the user matrix and the movie matrix?

The short answer is, we can guess. We guess different user matrices and movie matrices, and see if the predicted rating matrix is close to the actual rating given by users. When the two are close enough, we can use the two matrices to construct the new rating matrix. This is known as matrix factorization.

We can do better than guessing. There are some advanced statistical methods for estimating the user matrix and movie matrix, but this is beyond the scope of our class. If you are interested, you can search for “stochastic gradient descent.”

The matrix factorization algorithm

```
1 matrix_factorization <- function(R, P, Q, K, steps=5000, alpha=0.0002,  
  beta=0.02) {  
2   Q <- t(Q)  
3   for (step in 1:steps) {  
4     for (i in 1:nrow(R)) {  
5       for (j in 1:ncol(R)) {  
6         if (R[i, j] > 0) {  
7           eij <- R[i, j] - sum(P[i,] * Q[,j])  
8           for (k in 1:K) {  
9             P[i, k] <- P[i, k]+alpha*(2*eij*Q[k, j] - beta*P[i, k])  
10            Q[k, j] <- Q[k, j]+alpha*(2*eij*P[i, k] - beta*Q[k, j])  
11          }  
12        eR <- P %*% Q  
13        e <- 0  
14        for (i in 1:nrow(R)) {  
15          for (j in 1:ncol(R)) {  
16            if (R[i, j] > 0) {  
17              e <- e + (R[i, j] - sum(P[i,] * Q[,j]))^2  
18              for (k in 1:K) {  
19                e <- e + (beta/2) * (P[i, k]^2 + Q[k, j]^2)  
20              }  
21            }  
22          if (e < 0.001){break}  
23        return(list(P = P, Q = t(Q)))  
24      }
```

```

1  set.seed(123)
2  R <- matrix(c(5, 3, 0, 1,
3               4, 0, 0, 1,
4               1, 1, 0, 5,
5               1, 0, 0, 4,
6               0, 1, 5, 4,
7               2, 1, 3, 0), nrow = 6, ncol = 4, byrow = TRUE)
8  # N: num of User
9  N <- nrow(R)
10 # M: num of Movie
11 M <- ncol(R)
12 # Num of Features
13 K <- 2
14 P <- matrix(runif(N * K), nrow = N, ncol = K)
15 Q <- matrix(runif(M * K), nrow = M, ncol = K)
16 result <- matrix_factorization(R, P, Q, K)
17 nP <- result$P
18 nQ <- result$Q
19 nR <- nP %*% t(nQ)
20 print(nP)
21 print(nQ)
22 print(nR)

```


	Movie 1	Movie 2	Movie 3	Movie 4
Alice	4	4		1
Bob		2	2	3
Carol	1	5	3	
Dennis	3		4	1
Emma	5	2	1	4
Flora	3	1		5

$$\begin{bmatrix} 4 & 4 & ? & 1 \\ ? & 2 & 2 & 3 \\ 1 & 5 & 3 & ? \\ 3 & ? & 4 & 1 \\ 5 & 2 & 1 & 4 \\ 3 & 1 & ? & 5 \end{bmatrix} \approx \begin{bmatrix} 1.71 & 0.74 \\ 0.84 & 1.35 \\ 1.92 & -0.69 \\ 2.05 & 0.46 \\ 0.70 & 1.95 \\ 0.13 & 1.93 \end{bmatrix} \times \begin{bmatrix} 1.24 & 2.44 & 1.77 & -0.20 \\ 1.83 & 0.14 & 0.10 & 2.32 \end{bmatrix}$$

$$= \begin{bmatrix} 3.47 & 4.28 & 3.09 & 1.37 \\ 3.52 & 2.26 & 1.63 & 2.97 \\ 1.14 & 4.63 & 3.34 & -1.99 \\ 3.41 & 5.09 & 3.67 & 0.66 \\ 4.48 & 1.99 & 1.43 & 4.40 \\ 3.69 & 0.60 & 0.43 & 4.46 \end{bmatrix}$$

<https://www.youtube.com/embed/n3RKsY2H-NE?enablejsapi=1>

The Cold Start Problem