# INTRO TO DATA SCIENCE
## RECOMMENDATION ENGINES

# I. DATA TYPES
# II. GENERAL DESIGN
# III. CONTENT-BASED FILTERING
# IV. COLLABORATIVE FILTERING
# V. THE NETFLIX PRIZE

A recommendation system aims to match users to products/items/brands/etc that they likely haven't experienced yet and/or predict a user's preference based on past observations.

A recommendation system aims to match users to products/items/brand/etc that they likely haven't experienced yet and/or predict a users preference based on past observations.

A **ranking** or **prediction** is produced by analysing other user/item ratings (and sometimes item characteristics) to provide personalised recommendations to users.

# I. TYPES OF DATA

# THE KIND OF RECOMMENDATIONS YOU CAN GIVE, ARE DEPENDENT ON THE DATA YOU HAVE.

# WE NEED DATA TO RECOMMEND.

- Preferences
- Ratings
- Item meta-data
- User Behavior

**Ratings**
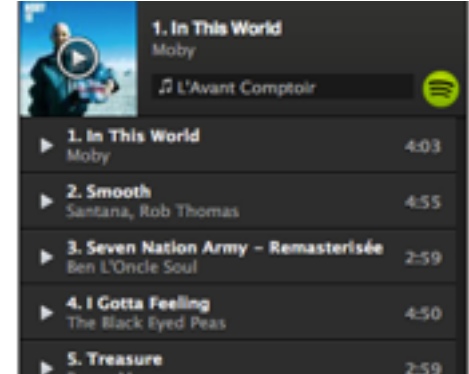**Upvotes / Downvotes**
**Weighted Scale**
**Grades**
**Relevance Feedback**

**Access Logs**
**Session Lengths**
**Time spent on a page**
**Clicks / Non-Clicks**
**Purchase History**
**Product Descriptions**

**Listening History**
**Playlist Creates**
**Follows / Unfriend**
**Impressions**
**Email Reads / Impressions**

*Recommenders need feedback to be useful.*

*Recommenders need feedback to be useful.*

## Explicit

- Explicitly given
- Pro-actively acquired
- Expensive to collect

*Recommenders need feedback to be useful.*

## Explicit

- Explicitly given
- Pro-actively acquired
- Expensive to collect

## Implicit

- Indirectly given
- Larger quantity
- Latent qualities

# Explicit or Implicit?

# Explicit or Implicit?

# Explicit or Implicit?

# Ratings: *Explicit*

# Explicit or Implicit?

# Explicit or Implicit?

# Swipes: *Explicit*

# Explicit or Implicit?

# Explicit or Implicit?

# Both!

# Explicit or Implicit?

# Explicit or Implicit?

# Wifi logs: *Implicit!*

## Explicit

- Ie: Ratings, surveys, reviews
- Easy to interpret
- Expensive

## Implicit

- Ie: Activity logs, clicks, impressions
- Hard to interpret
- Cheap



**NOTE**

*Implicit data collection can involve some privacy issues; any system that would make recommendations must avoid overstepping its bounds.*

# IA. EXPLICIT AND IMPLICIT FEEDBACK

Uber



Yelp



Reddit

Ratings, Votes, Reviews



Ebay

# *Explicit Feedback*

- *Frequently in the form of ratings*
- *Granularly represents preferences*
- *Requires extra effort from the user*

# *Explicit Feedback Questions*

- What does a rating mean?
- Do user preferences change?
- Is what is known about the data accurate?
  - Does what is collected reflect a preference at all?
  - Is it representative of the goal or only reflective of a singular characteristic?

# *Explicit Feedback – Considerations*

- Consistent scale for all ratings
- Can ratings be skewed by self/selection-bias
- Consider the ephemeral nature of preferences
- When the data was collected
  - Before or after experience
- Context of presentation

# Implicit Examples



## Order History



## Session Length

# Implicit Examples



Engagement Metrics



Session Length

# Implicit Feedback

It's still possible to make recommendations when no rating data is explicitly collected from a user.

The goal is to convert user behaviour into user preferences, but it entails one challenge:   How exactly does one infer preference based on actions in a system? This can be a difficult question to answer.

# Implicit feedback is everywhere.

- Email impressions
- Email click-throughs
- Conversions
- Demographic
- Session lengths
- Login attempts
- Track plays
- Money spent

- Ad impressions
- Ad clicks
- Ad click-purchase
- Web "click depth"
- # of swipes
- Profile views
- Message initiations
- Poll Votes

- Friend / unfriend
- Follow / unfollow
- *Like
- Post text
- Image EXIF
- Friends in common
- Message text
- Food purchases

- Geospatial data
- Store cameras
- Wifi logins / MAC
- Time series
- Objects in photos
- Driving record
- Credit history
- Topics most read

# *Implicit Feedback Caveats*

# Implicit Feedback Caveats

**(ie: Users don't tell you what you want to know.)**

- Preferences can be vague
- You may need to process tons of data to get what you want
- Analysis can be complicated / meaning hard to find

- Identities can be indistinguishable
- Users don't tell you what you want to know
- Easy to project bias onto data
- Positive / negative experience hard to assess

# Implicit Feedback General Advice:
# Question Everything.

- Can a preference actually be observed?
- Is the lack of data actually a negative preference?
- Is there enough data to describe feedback or only a portion of it?

- Is the data scaled properly?
- Are there hidden correlations?
- Are there contradictory patterns?
- What's missing?
- Can new features be created?

# Implicit + Explicit Feedback:
# Work together

**If a user rates an item, can you use implicit feedback to validate credibility?**

• **Did they read the article?**
• **Do they own the item?**
• **Did they rate before or after experience?**
• **Do other users mention them?**
• **Does user tend to rate high or low?**
• **How likely was the rating automated?**

**Use implicit data to understand the context and characteristics of a rating.**

• **Does time of day affect rating?**
• **Which kinds of reviews do they typically write?**
• **Are the reviews positive or negative?**
• **Do  other users like their reviews?**

# *Implicit + Explicit Feedback:  Final Caveat*

**Take care when creating explicit data from implicit data.**

- **Does the set of actions reflect a preference?**
- **Does the scale make sense?**
- **Is the outcome prediction (ratings) or recommendation?**

# Explicit

- **Higher value with respect to preferences**
- **Usually collected as a "rating"**
- **Collection is responsibility of user**
- **More direct evaluation of items**

# Implicit

- **Easy to collect in large quantities**
- **More difficult to work with**
- **Assumes nothing about the user (could be anyone!)**
- **Goal is to convert into preferences**

# II. GENERAL DESIGN

*There are two general approaches to the design:*

*There are many approaches to the design, but these are common modelling techniques:*

*In **content-based filtering**, items are mapped into a feature space, and recommendations depend on item characteristics.*

*In contrast, an important assumption underlying all of **collaborative filtering**, is: users who have similar preferences in the past are likely to have similar preferences in the future.*

## Recommendations for You in Books

| Cracking the Coding Interview: 150... | Introduction to Algorithms | Data Mining: Practical Machine... | Elements of Programming Interviews... | The Algorithm Design Manual |
|---|---|---|---|---|
| > Gayle Laakmann McDowell | Thomas H. Cormen, Charles E... | > Ian H. Witten, Eibe Frank, Mark A. Hall | > Amit Prakash, Adnan Aziz, Tsung-Hsien Lee | > Steve Skiena |
| Paperback | Hardcover | Paperback | Paperback | Paperback |
| ★★★★★ (166) | ★★★★☆ (85) | ★★★★☆ (27) | ★★★★½ (25) | ★★★★½ (47) |
| ~~$39.95~~ $23.22 | ~~$92.00~~ $80.00 | ~~$69.95~~ $42.09 | ~~$29.99~~ $26.18 | ~~$89.95~~ $71.84 |
| Why recommended? | Why recommended? | Why recommended? | Why recommended? | Why recommended? |

## Customers Who Bought This Item Also Bought



**Pitch Dark (NYRB Classics)**
› Renata Adler
Paperback
$11.54



**How Literature Saved My Life**
› David Shields
★★★★☆ (60)
Hardcover
$18.08



**Bleeding Edge**
Thomas Pynchon
Hardcover
$18.05



**The Flamethrowers: A Novel**
› Rachel Kushner
★★★★☆ (17)
Hardcover
$15.79

Recommended for you because you watched
Sugar Minott - Oh Mr Dc (Studio One)

**Mikey Dread - Roots and Culture**

by klaxonklaxon · 1,164,133 views

Lyrics:
Now here comes a special request
To each and everyone

6:00

Recommended for you because you watched
Thelonious Monk Quartet - Monk In Denmark

**Bill Evans Portrait in Jazz (Full Album)**

by hansgy1 · 854,086 views

Bill Evans Portrait in Jazz 1960
1. Come Rain or Come Shine - 3.19 (0:00)
2. Autumn Leaves - 5.23 (3:24)

42:26

Recommended for you because you watched
Bob Marley One Drop

**Bob Marley - She's gone**

by Dionysios29 · 1,058,704 views

This is one of the eleven songs of album Kaya that Bob Marley
and The Wailers creative in 1978.
Lyrics:

2:53

## How can we find good recommendations?

- Manual Curation

  *Songza*

  content-based filtering

- Manually Tag Attributes

  **P** PANDORA

- Audio Content, Metadata, Text Analysis

  the echonest

- Collaborative Filtering

  last.fm

  Spotify

MOST E-MAILED | RECOMMENDED FOR YOU

1. **How Big Data Is Playing Recruiter for Specialized Workers**

2. SLIPSTREAM
**When Your Data Wanders to Places You've Never Been**

3. MOTHERLODE
**The Play Date Gun Debate**

4. **For Indonesian Atheists, a Community of Support Amid Constant Fear**

5. **Justice Breyer Has Shoulder Surgery**

6. BILL KELLER
**Erasing History**

## 8. How do you determine my Most Read Topics?

Back to top ▲

Each NYTimes.com article is assigned topic tags that reflect the content of the article. As you read articles, we use these tags to determine your most-read topics.

To search for additional articles on one of your most-read topics, click that topic on your personalized Recommendations page. To learn more about topic tags, visit Times Topics.

**NOTE**

**Collaborative or Content based?**

## 8. How do you determine my Most Read Topics?

Back to top ▲

Each NYTimes.com article is assigned topic tags that reflect the content of the article. As you read articles, we use these tags to determine your most-read topics.

To search for additional articles on one of your most-read topics, click that topic on your personalized Recommendations page. To learn more about topic tags, visit Times Topics.

**NOTE**

**Collaborative or Content based?**

*CONTENT BASED* ☺

# III. CONTENT-BASED FILTERING

**Content-based filtering** *begins by mapping each item into a feature space. Both users and items are represented by vectors in this space.*

**Content-based filtering** *begins by mapping each item into a feature space. Both users and items are represented by vectors in this space.*

***Item vectors*** *measure the degree to which the item is described by each feature, and **user vectors** measure a user's preferences for each feature.*

**Content–based filtering** *begins by mapping each item into a feature space. Both users and items are represented by vectors in this space.*

***Item vectors*** *measure the degree to which the item is described by each feature, and* ***user vectors*** *measure a user's preferences for each feature.*

*Ratings are generated by taking* ***dot products*** *of user & item vectors.*

*One notable example of content-based filtering is Pandora, which maps songs into a feature space using features (or "genes") designed by the Music Genome Project.*

*Using song vectors that depend on these features, Pandora can create a station with music having similar properties to a song the user selects.*

*Content-based filtering has some difficulties:*

*Content-based filtering has some difficulties:*

- *Must map items into a feature space (usually by hand!)*

- *Recommendations are limited in scope (items must be similar to each other)*

- *Hard to create cross-content recommendations (eg books/music films...this would require comparing elements from different feature spaces!)*

# IV. COLLABORATIVE FILTERING

**Collaborative filtering** *refers to a family of methods for predicting ratings where instead of thinking about users and items in terms of a feature space, we are only interested in the existing user-item ratings themselves.*

**Collaborative filtering** *refers to a family of methods for predicting ratings where instead of thinking about users and items in terms of a feature space, we are only interested in the existing user-item ratings themselves.*

*In this case, our dataset is a ratings matrix whose columns correspond to items, and whose rows correspond to users.*

**Collaborative filtering** *refers to a family of methods for predicting ratings where instead of thinking about users and items in terms of a feature space, we are only interested in the existing user-item ratings themselves.*

*In this case, our dataset is a ratings matrix whose columns cor~~~~ to items, and whose rows correspond to users.*

**NOTE**

**The idea here is that users get value from recommendations based on other users with similar** *tastes*.

| | ← 18,000 movies → | | | | |
|---|---|---|---|---|---|
| x | 1 | 1 | x | ... | x |
| x | x | x | 5 | ... | x |
| x | x | 3 | x | ... | x |
| x | 4 | 3 | x | ... | 2 |
| ... | x | x | x | ... | x |
| x | 5 | x | 1 | ... | x |
| x | x | 3 | 3 | ... | x |
| x | 1 | x | x | ... | 2 |

480,000 users

**NOTE**

**This matrix will always be *sparse*!**

*source: http://www.eecs.berkeley.edu/~zhanghao/main/publications/subfolder/netflix.png*

*Main difference between content and collaborative filtering:*

*Content Based:*

> *maps items and users into a feature space*

*Collaborative:*

> *relies on previous user-item ratings*

*We will look at collaborative filtering in a user–user sense.*

*We will look at collaborative filtering in a user-user sense.*

*We will take a given user, and find the K most similar users, and then recommend brands from the similar users!*

*We will look at collaborative filtering in a user-user sense.*

*We will take a given user, and find the K most similar users, and then recommend brands from the similar users!*

**NOTE**

**Sounds familiar? It's similar to KNN!**

## Customers Who Bought This Item Also Bought

**Pitch Dark (NYRB Classics)**
› Renata Adler
Paperback
$11.54

**How Literature Saved My Life**
› David Shields
★★★★☆ (60)
Hardcover
$18.08

**Bleeding Edge**
Thomas Pynchon
Hardcover
$18.05

**The Flamethrowers: A Novel**
› Rachel Kushner
★★★☆☆ (17)
Hardcover
$15.79

*The system cannot draw inferences because it hasn't gathered enough information yet.*

*The cold start problem arises because we've been relying only on ratings data, or on* **explicit feedback** *from users.*

*The cold start problem arises because we've been relying only on ratings data, or on* **explicit feedback** *from users.*

*Until users rate several items, we don't know anything about their preferences!*

*The cold start problem arises because we've been relying only on ratings data, or on* **explicit feedback** *from users.*

*Until users rate several items, we don't know anything about their preferences!*

*We can get around this by enhancing our recommendations using* **implicit feedback***, which may include things like item browsing behaviour, search patterns, purchase history, etc.*

*While explicit feedback (ratings, likes, purchases) leads to high quality ratings, the data is sparse and cold starts are problematic.*

*While explicit feedback (ratings, likes, purchases) leads to high quality ratings, the data is sparse and cold starts are problematic.*

*Meanwhile implicit feedback (browsing behaviour, etc.) leads to less accurate ratings, the data is much more dense (and less invasive to collect).*

# V. THE NETFLIX PRIZE

*The Netflix prize was a competition to see if anyone could make a 10% improvement to Netflix's recommendation system (accuracy measured by RMSE).*

*The Netflix prize was a competition to see if anyone could make a 10% improvement to Netflix's recommendation system (measured by RMSE).*

*The grand prize was $1m dollars.*

*The Netflix prize was a competition to see if anyone could make a 10% improvement to Netflix's recommendation system (accuracy measured by RMSE).*

*The grand prize was $1m dollars.*

*The ratings matrix contained >100mm numerical entries (1-5 stars) from ~500k users across ~17k movies. The data was split into train/quiz/test sets to prevent overfitting on the test data by answer submission (this was a clever idea!).*

*The competition began in 2006, and the grand prize was eventually awarded in 2009. The winning entry was a stacked ensemble of 100's of models (including neighbourhood and matrix factorisation models) that were blended using boosted decision trees.*

*Ultimately, the competition ended in a photo finish. The winning strategy came down to last-minute team mergers and creative blending schemes to shave 3rd & 4th decimals off RMSE (concerns that would not be important in practice).*

*The competition began in 2006, and the grand prize was eventually awarded in 2009. The winning entry was a stacked ensemble of 100's of models (including neighbourhood & matrix factorisation models) that were blended using boosted decision trees.*

*Ultimately, the competition ended in a photo finish. The winning strategy came down to last-minute team mergers and creative blending schemes to shave 3rd & 4th decimals off RMSE (concerns that would not be important in practice).*

*The competition did much to spur interest and research advances in recommender systems technology, and the prize money was donated to charity.*

*Though they adopted some of the modelling techniques that emerged from the competition, Netflix never actually implemented the prizewinning solution.*

*Why do you think that's true?*

# VI. SUMMARY

- *Want to predict how users are going to rate items*
- *Obtain ratings implicitly or explicitly*
- *Try to predict these ratings through*
  - *Content based filtering*
  - *Collaborative filtering*
  - *Need to measure the similarity between user and item pairs*

- *Data Sparsity*
- *Cold Start*
- *Scalability*
- *Accurate but also recommendation of new content*
- *Evaluation*
- *Transparency to users*
- *Temporal changes*
- *Vulnerability to attacks*