# Machine Learning for Recommender Systems

**Laurent CHARLIN, Ph.D.**
Assistant Professor, Department of Decision Sciences, HEC Montreal
Member of Mila – Quebec Artificial Intelligence Institute

**IVADO**

# Introduction

Recommendation task:

- Suggest items of interest to users
- Items: movies, books, articles, humans
- Users: humans

# Is It Worth Our Attention?

Recommendation is the next search

- Search finds items (given a query)
- Recommendation finds items of interest
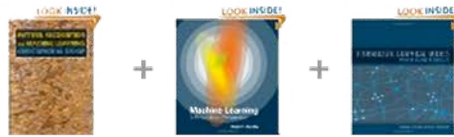
# Is It Worth Our Attention?

new

Recommendation is the ~~next~~ search

- Search finds items (given a query)
- Recommendation finds items of interest

## Frequently Bought Together

**Price for all three:** $230.47

[Add all three to Cart] [Add all three to Wish List]

Show availability and shipping details

- ☑ **This item:** Pattern Recognition and Machine Learning (Information Science and Statistics) by Christopher M. Bishop
  Hardcover $64.01
- ☑ Machine Learning: A Probabilistic Perspective (Adaptive Computation and Machine Learning series) by Kevin P. Murphy
  Hardcover $78.26
- ☑ Probabilistic Graphical Models: Principles and Techniques (Adaptive Computation and Machine Learning ... by Daphne Koller
  Hardcover $88.20

## Customers Who Bought This Item Also Bought

| Learning From Data | Machine Learning: A Probabilistic ... | The Elements of Statistical Learning: ... | Probabilistic Graphical Models: Principles ... | Machine Learning |
|---|---|---|---|---|
| › Yaser S. Abu-Mostafa | › Kevin P. Murphy | Trevor Hastie | › Daphne Koller | › Tom M. Mitchell |
| ★★★★½ (56) | ★★★★ (29) | ★★★★ (31) | ★★★★ (24) | ★★★★½ (47) |
| Hardcover | Hardcover | Hardcover | Hardcover | Hardcover |
| | | | | $195.71 ✓Prime |

# NETFLIX

**Top Picks for Me**

- They are responsible for 4% of US marriages (from 2005 to 2012)

- And lower divorce rates

# Machine Learning for Recommender Systems

- Task: Suggest items of interest to users

- From data how do you determine what denotes interest?

- Item-specific signal (supervised learning)
    1. Score: rating, bid
    2. Consumption: click, buy, watch, bookmark

Imagine
- The data are user ratings
- Task: Recommend items the user will like

How do we set it up as a machine learning problem?

**Data** → | ML | → **Predictions**

Imagine
- The data are user ratings
- Task: Recommend items the user will like

How do we set it up as a machine learning problem?

**Data** → **ML** → **Predictions**

- Task: What do we learn? What do we predict? What is the model?
- Performance measure: How do we evaluate the results?
- Experience: How does our model interact with data?

# Questions?

# Framework for Recommendation Problems

| Data | $\rightarrow$ | Representation of user preferences | $\rightarrow$ | Recommendations |
|------|------|------|------|------|

**User preferences**

★ ★ ★ ★ ☆

**User/Item features**

**Model**



**E.g. Top-N recommendations**

```
┌──────────────────┐         ┌──────────────────┐
│                  │         │                  │
│ Representation of │ ── · ▶ │ Recommendations  │
│ user preferences │         │                  │
│                  │         │                  │
└──────────────────┘         └──────────────────┘
```

Task:

- How we we set it up?
1. Regression (Classification)
2. Ranking

# Ranking vs. Regression

A.  Ranking models

- Computationally more expensive
- E.g., Have to consider a group of items (listwise)

$$\underbrace{f\colon (u, i_1, i_2, ..., i_m)}_{\text{user u's unseen items}} \rightarrow \underbrace{(r_1, r_2, ..., r_m)}_{\text{rank of each item}}$$

B.  Score models

- For each user:
1.  Predict scores of all unseen items  $\mathbf{f}\colon (\mathbf{u}, \mathbf{i}) \rightarrow \mathbb{R}$
2.  Rank items (show top-K)

# Framework for Recommendation Problems

**Data**

**Representation of user preferences**

**Predict Missing Ratings**

User preferences

★ ★ ★ ★ ☆

User/Item features

**Model**

# Score Prediction as Regression

$$\begin{bmatrix} 3 & - & \cdots & 0 \\ - & 0 & \cdots & - \\ \vdots & & \ddots & \cdots \\ 2 & - & \cdots & - \end{bmatrix}$$ users × items $\xrightarrow{\;f\;}$ $$\begin{bmatrix} 3 & 2 & \cdots & 0 \\ 1 & 0 & \cdots & 3 \\ \vdots & & \ddots & \cdots \\ 0 & 2 & \cdots & 2 \end{bmatrix}$$
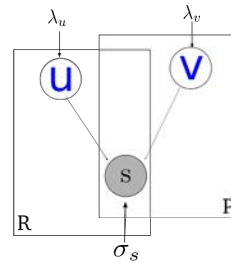
Train : Black  $S^o$
Test  : Red    $S^u$

# Score Prediction as Regression

$$\begin{bmatrix} 3 & - & \cdots & 0 \\ - & 0 & \cdots & - \\ \vdots & & \ddots & \cdots \\ 2 & - & \cdots & - \end{bmatrix}$$ users × items

$\xrightarrow{f}$

$$\begin{bmatrix} 3 & 2 & \cdots & 0 \\ 1 & 0 & \cdots & 3 \\ \vdots & & \ddots & \cdots \\ 0 & 2 & \cdots & 2 \end{bmatrix}$$

**How do we set this up as a learning problem?**
$s^u = f(S^o)$

# Collaborative Filtering (CF)

- Assumption:
  - Users with past similar preferences will have similar future preferences

- Work horse used in many recommender systems



**Customers Who Bought This Item Also Bought**

Learning From Data
› Yaser S. Abu-Mostafa
★★★★½ (56)
Hardcover

Machine Learning: A Probabilistic …
› Kevin P. Murphy
★★★★ (29)
Hardcover

The Elements of Statistical Learning: …
Trevor Hastie
★★★★ (31)
Hardcover

Probabilistic Graphical Models: Principles …
› Daphne Koller
★★★★ (24)
Hardcover

Machine Learning
› Tom M. Mitchell
★★★★½ (47)
Hardcover
$195.71 ✓Prime

# CF - Neighbourhood Approaches

1. For each user, find other users with similar past preferences

2. Predict that user's missing preferences as the weighted combination of its neighbours' preferences

$$
\begin{bmatrix}
3 & - & \cdots & 0 \\
- & 0 & \cdots & - \\
\vdots & & \ddots & \cdots \\
2 & - & \cdots & -
\end{bmatrix}
$$
**users** $\times$ **items**

# CF - Neighbourhood

- Find similarity between every pair of users (or items)

$$Sim(u, u') = \frac{(S_u^o)^\top S_{u'}^o}{\|S_u^o\|\|S_{u'}^o\|}$$

$$\begin{bmatrix} 3 & - & \cdots & 0 \\ - & 0 & \cdots & - \\ \vdots & & \ddots & \cdots \\ 2 & - & \cdots & - \end{bmatrix}$$

$S_{u'}^o$

$S_u^o$

- Predict missing scores using a user's neighbours

$$\widehat{S_{uj}} = \frac{\sum_{u'} Sim(u, u') S_{u'j}^o}{\sum_{u'} Sim(u, u')} \quad \forall u' \text{ that have rated } j$$

$s^o$: observed scores (training data)
$s^u$: unobserved scores (test data)

# CF - Neighbourhood Approaches

- Non-parametric approach

  - A user is represented by a weighted combination of its neighbours
  - New users can change one's recommendations

- Different distance functions to capture different effects

  - Ratings vs. clicks
  - Could consider additional information

- Works well empirically

- Building similarity matrix can be slow (offline)

- Not probabilistic

# Questions?

# CF - Matrix Factorization

$$\begin{array}{c}\text{users}\end{array}\begin{bmatrix} 3 & - & \cdots & 0 \\ - & 0 & \cdots & - \\ \vdots & & \ddots & \cdots \\ 2 & - & \cdots & - \end{bmatrix} \approx \begin{array}{c}\text{users}\end{array}\overset{k}{\begin{bmatrix} \theta \end{bmatrix}} \quad k\overset{\text{items}}{\begin{bmatrix} \beta \end{bmatrix}}$$

- Assumption: the observation matrix is low-rank

- Estimates user and item representations

- k is a hyperparameter

- **k << min(|Users|, |Items|)**

**Model.** $S_{ui} := \theta_u^T \beta_i$

**Parameters.** $\theta_u \ \forall u, \beta_i \ \forall i$

**Objective.** $\Sigma_u \Sigma_i \left( S^o_{ui} - \widehat{S^o}_{ui} \right)^2$

[Salakhutdinov, Mnih, '08]

# CF - Matrix Factorization: Alternative View

Model : $S_{ui} := \theta_u^\top \beta i$

Imagine that $\boldsymbol{\theta}_U$'s are features of users

**The model is then a linear regression for each item:** $\begin{aligned} S_{ui} &= \theta_u^\top \beta i \\ &= \sum_k \theta_{uk} \beta_{ik} \\ &= \sum_k \theta_{u1} \beta_{i1} + \theta_{u2} \beta_{i2} + \cdots + \theta_{up} \beta_{ip} \end{aligned}$

Since the model is symmetric in $\boldsymbol{\theta}$ and $\boldsymbol{\beta}$,
$\boldsymbol{\beta_i}$'s can be seen as features of items

# Model Fitting

Objective $\Sigma_u \Sigma_i \left( S_{ui} - \widehat{S}_{ui} \right)^2$

# Model Fitting

Objective $\quad \Sigma_u \Sigma_i \left( S_{ui} - \widehat{S}_{ui} \right)^2$

- Joint parameter optimization
  - Gradient descent: $(\nabla\theta, \nabla\beta)$

# Model Fitting

$$\text{Objective } \Sigma_u \Sigma_i \left( S_{ui} - \widehat{S}_{ui} \right)^2$$

- Joint parameter optimization
  - Gradient descent: $(\nabla\theta, \nabla\beta)$

- Alternate optimization
  - Fix $\theta$, update $\beta$
  - Fix $\beta$, update $\theta$

- Each step is a (regularized) least-squares problem
- This procedure is known as alternating least squares (ALS)
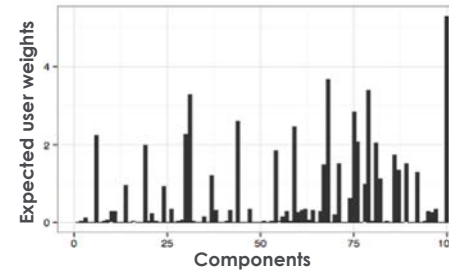
$S^o$

# Matrix Factorization

$S^u$

**User's highly rated movies**

E.T. the Extra-Terrestrial (Children's, Drama)
Full Metal Jacket (Action, Drama, War)
Three Colors: Red (Drama)
Breaker Morant (Drama, War)
Shakespeare in Love (Comedy, Romance)
Shadowlands (Drama, Romance)
Rob Roy (Drama, Romance, War)
The Verdict (Drama)
A Little Princess (Children's, Drama)
Leaving Las Vegas (Drama , Romance)

**User's weights for 100 components**



**Top movies recommended for the user**

Casablanca (Drama, Romance , War)
Breakfast at Tiffany's (Drama, Romance)
Amadeus (Drama)
When Harry Met Sally... (Comedy, Romance)
American Beauty (Comedy, Drama)
Fargo (Crime, Drama, Thriller)
The Right Stuff (Drama)  Gandhi (Drama)
Apocalypse Now (Drama, War)
Toy Story (Children's, Comedy, Animation)

[Gopalan et al.'15]

# Model
# Exporation

$$\text{argsort}_i \, \beta_{ik}$$

[Gopalan et al.'15]

**Movielens**

| "Sci-Fi" | "Drama, Romance" | "Action" |
|---|---|---|
| Day the Earth Stood Still<br>Metropolis<br>Forbidden Planet<br>Them!<br>Invasion of the Body Snatchers<br>The War of the Worlds<br>Godzilla<br>Village of the Damned<br>Night of the Living Dead<br>The Thing From Another World | Strictly Ballroom<br>Like Water for Chocolate<br>The Postman<br>Sense and Sensibility<br>Much Ado About Nothing<br>The Remains of the Day<br>Howards End<br>An Ideal Husband<br>Henry V<br>Shawdowlands | Die Hard 2<br>Die Hard: With a Vengeance<br>Independence Day<br>Air Force One<br>The Rock<br>Con Air<br>Enemy of the State<br>Conspiracy Theory<br>The Matrix<br>Broken Arrow |

6K Users
4K Movies
1M Ratings

**Netflix**

| "Supernatural thriller" | "Literary films" | "Friends sitcom" |
|---|---|---|
| Stir of Echoes<br>The Exorcist<br>The Ring<br>Final Destination<br>Misery<br>What Lies Beneath<br>Poltergeist<br>The Shining<br>Carrie<br>Gothika | Pride and Prejudice<br>Sense and Sensibility<br>Elizabeth<br>Emma<br>Sense and Sensibility<br>Mansfield Park<br>Much Ado About Nothing<br>The Importance of Being Earnest  Anne of Green Gables<br>Shakespeare in Love | Friends: Season 1<br>Friends: Season 2<br>Friends: Season 4<br>The Best of Friends: Vol. 1<br>Friends: Season 3<br>Friends: Season 5<br>The Best of Friends: Season 1<br>The Best of Friends: Season 2<br>The Best of Friends: Season 3<br>Friends: Season 6 |

480K Users
17.7K Movies
100M Ratings

**Mendeley**

| 'Sociology' | "Wireless sensor networks" | "Distributed behavior" |
|---|---|---|
| Social Capital: Its Origins, Institutions and Economic…<br>Institutions and Economic…<br>Increasing Returns and Path Dependence…<br>Diplomacy & Domestic Politics…<br>Comparative Politics and the Comparative…<br>Ethnicity, Insurgency, and Civil  War…<br>Historical Institutionalism in Comparative…<br>Case studies and theory development in social…<br>The Politics, Power, Pathologies…<br>End of the Transition Paradigm… | Wireless sensor networks : a survey…<br>Wireless sensor network survey  An energy-efficient MAC protocol…<br>A survey of routing protocols for…<br>Wireless sensor networks for habitat…<br>Cognitive radio: brain-empowered  wireless…<br>A survey on wireless multimedia  sensor networks<br>NeXt generation/dynamic spectrum…<br>Routing techniques in wireless sensor…<br>Social network analysis… | Flocks , herds and schools<br>Flocking for multi-agent…<br>Market- Based multirobot…<br>Coordination of groups of mobile autonomous…<br>Behavior-based formation control  for multi robot teams…<br>A formal analysis and taxonomy of task allocation…<br>A survey of consensus problem in  multi-agent coordination…<br>Modeling swarm robotic systems:…<br>Cooperative mobile robotics: A case study…<br>The e-puck, a robot designed for education in engineering… |

80K Users
260K Sci. articles
100M Ratings

# Probabilistic Matrix Factorization

**Matrix Factorization**

$$\|\beta_i\|_2$$
$$\|\theta_u\|_2$$
$$(S_{ui} - \hat{S}_{ui})^2$$

**Gaussian Matrix Factorization**
[Salakhutdinov et al. '08]

$$\theta_u \sim \mathcal{N}(a, b)$$
$$\beta_i \sim \mathcal{N}(c, d)$$
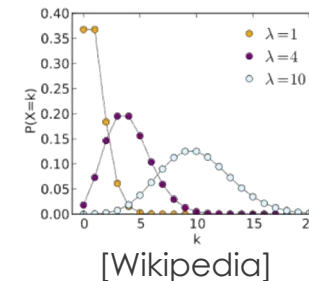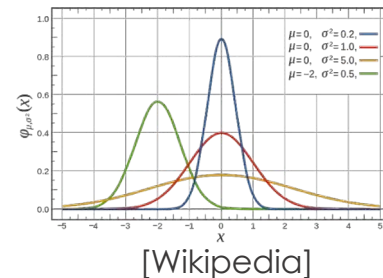$$S_{ui} \sim \mathcal{N}(\theta_u^\top \beta_i, \sigma)$$

**Poisson Matrix Factorization**
[Gopalan et al. '15]

$$\theta_u \sim \mathrm{Gamma}(a, b)$$
$$\beta_i \sim \mathrm{Gamma}(c, d)$$
$$S_{ui} \sim \mathrm{Poisson}(\theta_u^\top \beta_i)$$

- Poisson factorization is correct

- Gaussian factorization is incorrect

- In practice MF typically gives better performance than PF

**Minimizing mean squared error is equivalent to maximizing likelihood under Gaussian noise**


[Wikipedia]


[Wikipedia]

# Questions?

# Towards CF With Deep Learning

$$\begin{bmatrix} 3 & - & \cdots & 0 \\ - & 0 & \cdots & - \\ \vdots & & \ddots & \cdots \\ 2 & - & \cdots & - \end{bmatrix}$$

# Towards CF With Deep Learning

$$\begin{bmatrix} 3 & - & \cdots & 0 \\ - & 0 & \cdots & - \\ \vdots & & \ddots & \cdots \\ 2 & - & \cdots & - \end{bmatrix} \longrightarrow$$

(User 1, Item 1, 3)
(User 1, Item 2, 0)
(User 2 Item 2, 0)
(User n, Item 1, 2)

# Towards CF With Deep Learning

$$\begin{bmatrix} 3 & - & \cdots & 0 \\ - & 0 & \cdots & - \\ \vdots & & \ddots & \cdots \\ 2 & - & \cdots & - \end{bmatrix}$$

(User 1, Item 1, 3)
(User 1, Item 2, 0)
(User 2 Item 2, 0)
(User n, Item 1, 2)

**Encode each categorical variable using a series of indicator variables**

User     Item

(1 0 … 0, 1 0 … 0, 3)
(1 0 … 0, 0 0 … 1, 0)
(0 1 … 0, 0 1 … 1, 0)
(0 0 … 1, 1 0 … 0, 2)

# Towards CF With Deep Learning

$$\begin{bmatrix} 3 & - & \cdots & 0 \\ - & 0 & \cdots & - \\ \vdots & & \ddots & \cdots \\ 2 & - & \cdots & - \end{bmatrix}$$

(User 1, Item 1, 3)
(User 1, Item 2, 0)
(User 2 Item 2, 0)
(User n, Item 1, 2)

**Encode each categorical variable using a series of indicator variables**

User      Item

(1 0 … 0, 1 0 … 0, 3)
(1 0 … 0, 0 0 … 1, 0)
(0 1 … 0, 0 1 … 1, 0)
(0 0 … 1, 1 0 … 0, 2)

(1 0 … 0, 1 0 … 0, 3)

$x$          $y$

# Towards CF With Deep Learning

$$\begin{bmatrix} 3 & - & \cdots & 0 \\ - & 0 & \cdots & - \\ \vdots & & \ddots & \cdots \\ 2 & - & \cdots & - \end{bmatrix}$$

(User 1, Item 1, 3)
(User 1, Item 2, 0)
(User 2 Item 2, 0)
(User n, Item 1, 2)

**Encode each categorical variable using a series of indicator variables**

User      Item

$(1\ 0 \ldots 0, 1\ 0 \ldots 0, 3)$
$(1\ 0 \ldots 0, 0\ 0 \ldots 1, 0)$
$(0\ 1 \ldots 0, 0\ 1 \ldots 1, 0)$
$(0\ 0 \ldots 1, 1\ 0 \ldots 0, 2)$

$(1\ 0 \ldots 0, 1\ 0 \ldots 0, 3)$

$\underbrace{\underbrace{1\ 0 \ldots 0}_{x_1}, \underbrace{1\ 0 \ldots 0}_{x_2}}_{x}, \underbrace{3}_{y}$

# Towards CF With Deep Learning

$$\begin{bmatrix} 3 & - & \cdots & 0 \\ - & 0 & \cdots & - \\ \vdots & & \ddots & \cdots \\ 2 & - & \cdots & - \end{bmatrix}$$

(User 1, Item 1, 3)
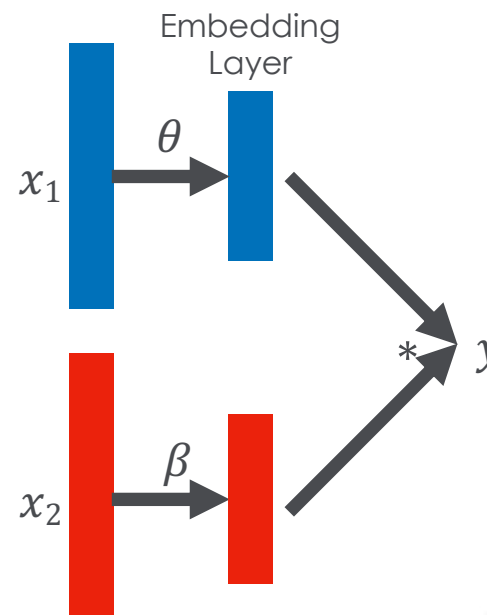(User 1, Item 2, 0)
(User 2 Item 2, 0)
(User n, Item 1, 2)

**Encode each categorical variable using a series of indicator variables**

User          Item

(1 0 … 0, 1 0 … 0, 3)
(1 0 … 0, 0 0 … 1, 0)
(0 1 … 0, 0 1 … 1, 0)
(0 0 … 1, 1 0 … 0, 2)

$$\underbrace{(\underbrace{1\ 0\ \dots\ 0}_{x_1},\ \underbrace{1\ 0\ \dots\ 0}_{x_2}}_{x},\ \underbrace{3}_{y})$$

$$y = (x_1^\mathsf{T}\theta)^\mathsf{T}(x_2^\mathsf{T}\beta)$$

# Towards CF With Deep Learning

$$\begin{bmatrix} 3 & - & \cdots & 0 \\ - & 0 & \cdots & - \\ \vdots & & \ddots & \cdots \\ 2 & - & \cdots & - \end{bmatrix}$$

(User 1, Item 1, 3)
(User 1, Item 2, 0)
(User 2 Item 2, 0)
(User n, Item 1, 2)

**Encode each categorical variable using a series of indicator variables**

User     Item

(1 0 … 0, 1 0 … 0, 3)
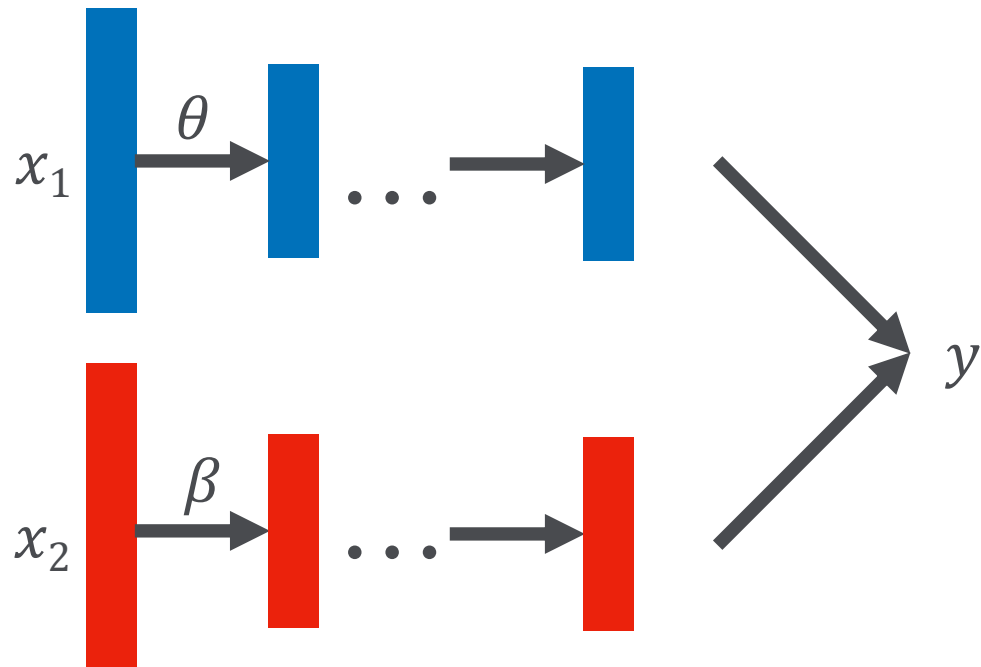(1 0 … 0, 0 0 … 1, 0)
(0 1 … 0, 0 1 … 1, 0)
(0 0 … 1, 1 0 … 0, 2)

(1 0 … 0, 1 0 … 0, 3)

$\underbrace{\phantom{(1 0 … 0,}}_{x_1} \underbrace{\phantom{1 0 … 0)}}_{x_2}$

$x$     $y$

$$y = (x_1^\mathsf{T}\theta)^\mathsf{T}(x_2^\mathsf{T}\beta)$$

Embedding Layer

$x_1$   $\theta$

$x_2$   $\beta$

$*$   $y$

# A Version of Deep Matrix Factorization



Can do more complicated
user and item combinations
(beyond dot product)

[Xue et al.'17]

# Questions?

# Autoencoders

Popular neural-network model often used in unsupervised learning (e.g., dim reduction)

# Autoencoders

Popular neural-network model often used in unsupervised learning (e.g., dim reduction)
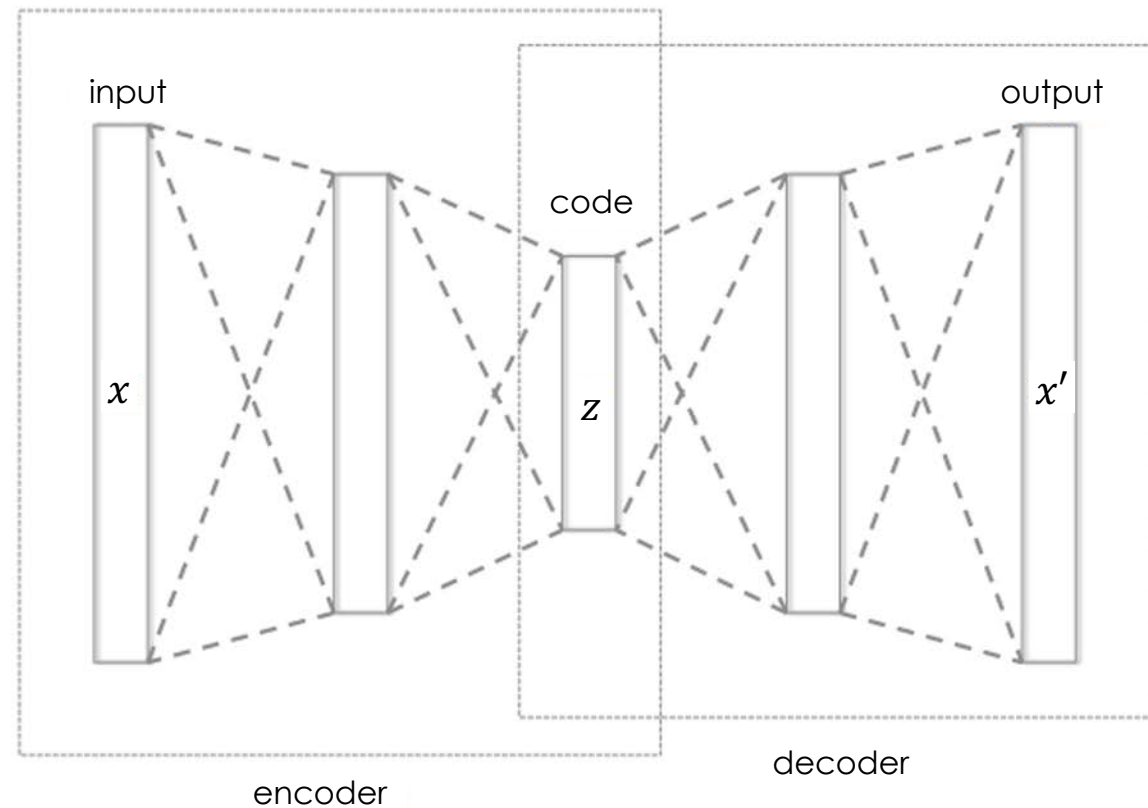
- Non-linear PCA

# Autoencoders

Popular neural-network model often used in unsupervised learning (e.g., dim reduction)

- Non-linear PCA

- Intuition: let's learn to copy the data $x' = f(x)$

# Autoencoders

Popular neural-network model often used in unsupervised learning (e.g., dim reduction)

- Non-linear PCA

- Intuition: let's learn to copy the data $x' = f(x)$

- We force a "bottleneck" $z = f_1(x)$

$$x' = f_2(z)$$

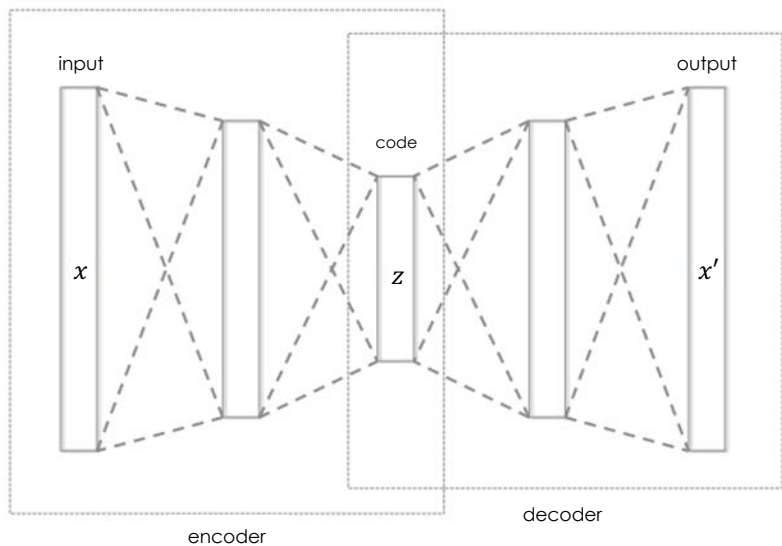$$dim(z) < dim(x)$$

# Autoencoders

$$z = f_1(x)$$
$$x' = f_2(z)$$
$$dim(z) < dim(x)$$



Hyperparameters: size of all neurons (layers) except input one

[From Wikipedia]

# Autoencoders for CF



encoder          decoder

[From Wikipedia]
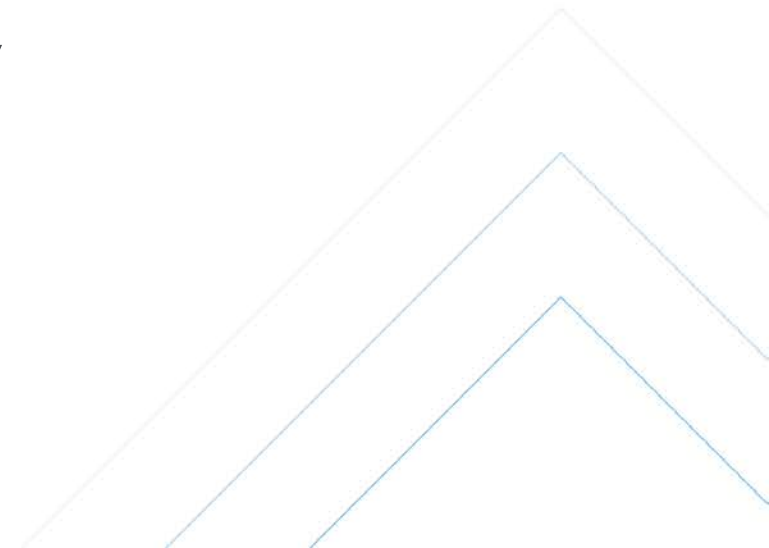
- $x$: Either a row or a column of the ratings matrix

$$x \quad \begin{bmatrix} 3 & - & \cdots & 0 \\ - & 0 & \cdots & - \\ \vdots & & \ddots & \cdots \\ 2 & - & \cdots & - \end{bmatrix}$$ **users $\times$ items**

- Missing entries:
  - Set to 0 in the input
  - Not considered in the output

- Many versions using denoising-autoencoders (DAEs), VAEs, different likelihoods (etc.)

# MF vs. AE for CF

- AE are more "naturally" non-linear

- AE are asymmetric
  - Must choose whether to model users or items

- Versions of AEs are close to the state-of-the-art today

# Questions?

# How to Choose the Right Model?

- Search for papers that compare different models

  - Keep a healthy does of scepticism

    - I.e., results in papers is not necessarily ground truth

- Try it out on your data

  - Compare performance on held-out data

  - Can it handle: your data size, service speed, updating schedule, other desiderata (fairness, uncertainty estimates) …

# Are Deep Models Better?

- 2014: No
- 2018: Yes
- 2019: Maybe not … *

  - "Embarrassingly Shallow Autoencoders for Sparse Data", Steck'19

  - "On the Difficulty of Evaluating Baselines: A Study on Recommender Systems", Rendle et al.'19:

    - "With a careful setup of a vanilla matrix factorization baseline, we are not only able to improve upon the reported results for this baseline but even outperform the reported results of any newly proposed method."

    (* This is not considering possibly available covariates)

# Explicit vs. Implicit Data

- Up to now we assumed ratings

- Ratings are explicit data:

  - "Users explicitly provide their preferences"

  - A high rating means the user liked the item

  - A low rating means the user disliked the item

- In practice implicit data is much more common:

  - click, buy, watch, listen

# Challenge with Implicit Data

- Consuming an item usually implies a positive preference

- Not Consuming an item may either indicate:

  - A negative preference
  - Something else: e.g., lack of exposure or time

# Challenge with Implicit Data

The preference matrix is "full" (as opposed to sparse)

- '1' indicates a consumed item
- '0' indicates an unconsumed item

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 \end{bmatrix}$$

- You must take both 0s and 1s into account
- In practice many models can be adapted to implicit case

# Common Strategy for Implicit Data

- Model the 0s as being "less certain" than the 1s

$$\text{Objective MF: } \frac{1}{|\text{users}|} \sum_u \sum_i (S_{ui} - \hat{S}_{ui})^2$$

$$\text{Objective WMF: } \frac{1}{|\text{users}|} \sum_u \sum_i c_{ui}(S_{ui})(S_{ui} - \hat{S}_{ui})^2$$

$$c_{ui}(0) < c_{ui}(1)$$

- Weighted Matrix Factorization [Hu et al. '08]

- Learn the weight of each zero

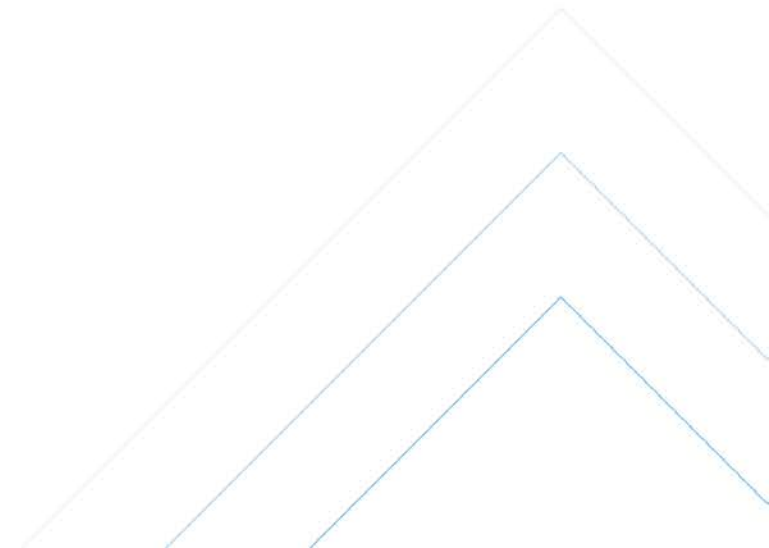- Exposure Matrix Factorization [Liang et al.'15]

# Questions?

# User/Item Features (I)

Often additional information exists

- Users: demographic information, social networks

# User/Item Features (I)

Often additional information exists

- Users: demographic information, social networks

- Items: content (e.g., movie genre/trailer, book text)
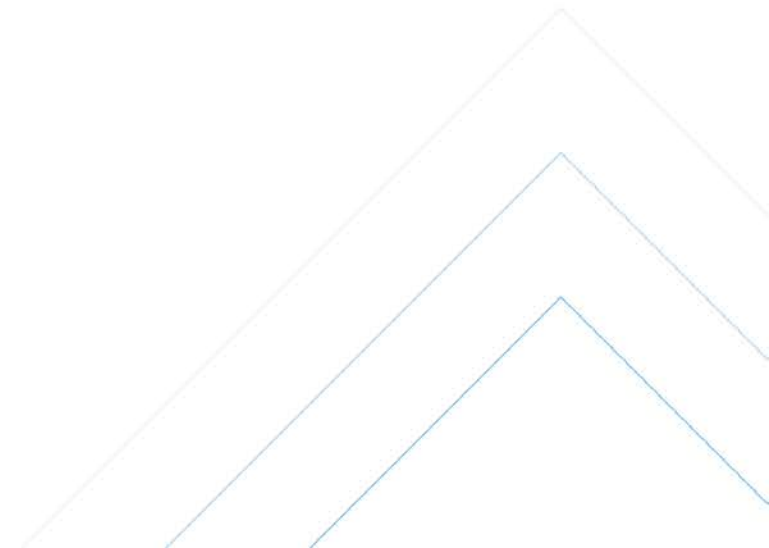
- Users & items:

# User/Item Features (I)

Often additional information exists

- Users: demographic information, social networks

- Items: content (e.g., movie genre/trailer, book text)

- Users & items:
  - timestamps, session information

# User/item Features (II)

- Allow for content-based recommendations
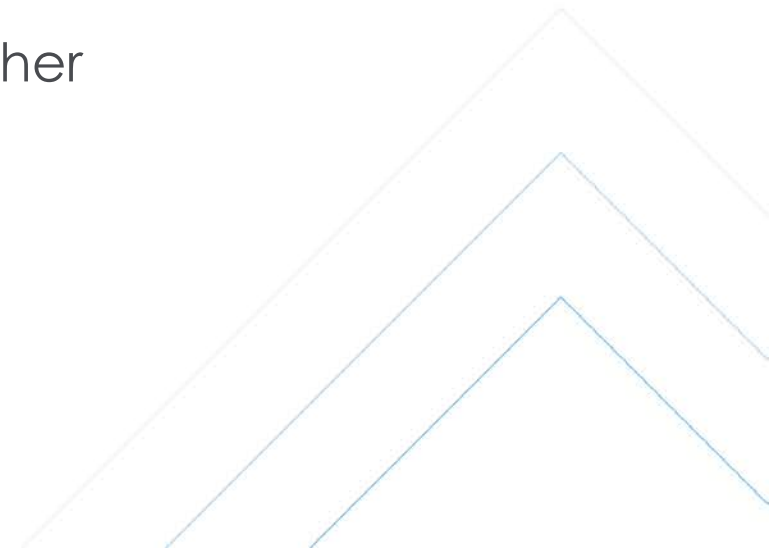  - Good to combat the cold-start problem

# User/item Features (II)

- Allow for content-based recommendations
  - Good to combat the cold-start problem

- Assume that features are predictive of preferences

# User/item Features (II)

- Allow for content-based recommendations
  - Good to combat the cold-start problem

- Assume that features are predictive of preferences
  - More difficult in some domains than others (e.g., movies)

- A practical approach is to bootstrap with content-based to gather preference data and then switch to CF

# User/item Features (II)

- Allow for content-based recommendations
  - Good to combat the cold-start problem

- Assume that features are predictive of preferences
  - More difficult in some domains than others (e.g., movies)

- A practical approach is to bootstrap with content-based to gather preference data and then switch to CF

- In the next slides we explore hybrid models for these data
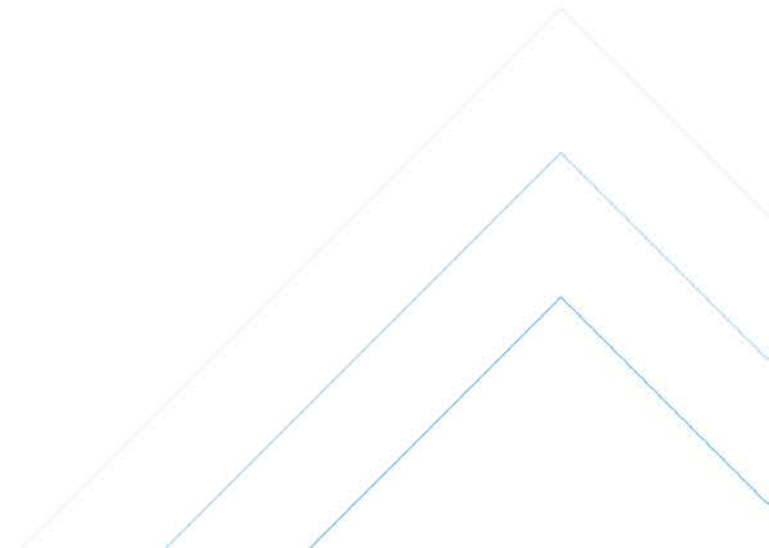
# Modelling Strategy

1. Generic models

   - Easily extend to many different use cases

2. Tailored modelling for specific features

   - This is where neural nets shine (images, text, networks)

# Factorization Machines (FM) [Generic Model]

Model "all" additional information

$$\begin{bmatrix} 3 & - & \cdots & 0 \\ - & 0 & \cdots & - \\ \vdots & & \ddots & \cdots \\ 2 & - & \cdots & - \end{bmatrix}$$

[Rendle'10]

# Factorization Machines (FM) [Generic Model]

## Model "all" additional information

$$\begin{bmatrix} 3 & - & \cdots & 0 \\ - & 0 & \cdots & - \\ \vdots & & \ddots & \cdots \\ 2 & - & \cdots & - \end{bmatrix}$$

**Encode each categorical variable using a series of indicator variables**

$\longrightarrow$

User  Item

$(1\ 0 \ldots 0, 1\ 0 \ldots 0, 3)$
$(1\ 0 \ldots 0, 0\ 0 \ldots 1, 0)$
$(0\ 1 \ldots 0, 0\ 1 \ldots 1, 0)$
$(0\ 0 \ldots 1, 1\ 0 \ldots 0, 2)$

[Rendle'10]

# Factorization Machines (FM) [Generic Model]

Model "all" additional information

$$\begin{bmatrix} 3 & - & \cdots & 0 \\ - & 0 & \cdots & - \\ \vdots & & \ddots & \cdots \\ 2 & - & \cdots & - \end{bmatrix}$$

**Encode each categorical variable using a series of indicator variables** →

User     Item

(1 0 … 0, 1 0 … 0, 3)
(1 0 … 0, 0 0 … 1, 0)
(0 1 … 0, 0 1 … 1, 0)
(0 0 … 1, 1 0 … 0, 2)

**Add features as columns** →

User    Item    Age

(1 0 … 0, 1 0 … 0, 25, 3)
(1 0 … 0, 0 0 … 1, 22, 0)
(0 1 … 0, 0 1 … 1, 55, 0)
(0 0 … 1, 1 0 … 0, 60, 2)

[Rendle'10]

# Factorization Machines (FM) [Generic Model]

Model "all" additional information

$$\begin{bmatrix} 3 & - & \cdots & 0 \\ - & 0 & \cdots & - \\ \vdots & & \ddots & \cdots \\ 2 & - & \cdots & - \end{bmatrix}$$

**Encode each categorical variable using a series of indicator variables**
⟶

User     Item

(1 0 … 0, 1 0 … 0, 3)
(1 0 … 0, 0 0 … 1, 0)
(0 1 … 0, 0 1 … 1, 0)
(0 0 … 1, 1 0 … 0, 2)

**Add features as columns**
⟶

User     Item     Age

(1 0 … 0, 1 0 … 0, 25, 3)
(1 0 … 0, 0 0 … 1, 22, 0)
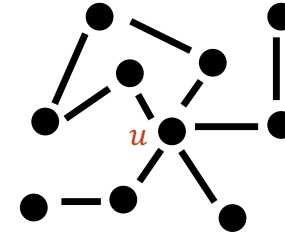(0 1 … 0, 0 1 … 1, 55, 0)
(0 0 … 1, 1 0 … 0, 60, 2)

**Model.** $S_{ui} := W_0 + \underbrace{\sum_i^p w_i x_i}_{\text{per-feature regression}} + + \underbrace{\sum_{j=0}^p \sum_{j'=j+1}^p \theta_j^\top \theta x_j x_{j'}}_{\text{per-pair regression}}$

[Rendle'10]

# Factorization Machines (FM) [Generic Model]

Model "all" additional information

$$\begin{bmatrix} 3 & - & \cdots & 0 \\ - & 0 & \cdots & - \\ \vdots & & \ddots & \cdots \\ 2 & - & \cdots & - \end{bmatrix}$$

**Encode each categorical variable using a series of indicator variables** →

User     Item

$(1\,0\,\ldots\,0,\ 1\,0\,\ldots\,0,\ 3)$
$(1\,0\,\ldots\,0,\ 0\,0\,\ldots\,1,\ 0)$
$(0\,1\,\ldots\,0,\ 0\,1\,\ldots\,1,\ 0)$
$(0\,0\,\ldots\,1,\ 1\,0\,\ldots\,0,\ 2)$

**Add features as columns** →

User    Item    Age

$(1\,0\,\ldots\,0,\ 1\,0\,\ldots\,0,\ 25,\ 3)$
$(1\,0\,\ldots\,0,\ 0\,0\,\ldots\,1,\ 22,\ 0)$
$(0\,1\,\ldots\,0,\ 0\,1\,\ldots\,1,\ 55,\ 0)$
$(0\,0\,\ldots\,1,\ 1\,0\,\ldots\,0,\ 60,\ 2)$

**Model.** $S_{ui} := W_0 + \underbrace{\sum_i^p w_i x_i}_{\text{per-feature regression}} + + \underbrace{\sum_{j=0}^p \sum_{j'=j+1}^p \theta_j^\top \theta x_j x_{j'}}_{\text{per-pair regression}}$

- Features added to the data (extra columns) are "automatically" used in the model

- Modelling extra information implies adding the feature

[Rendle'10]

# A. Social Network



- Data: user ratings and users' friends

- Assume:

  1. Friends influence your preferences

  2. Different levels of trusts for different friends

[Chaney et al. '15]

# A. Social Network

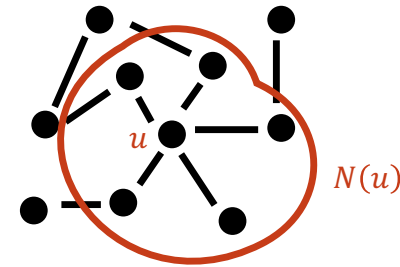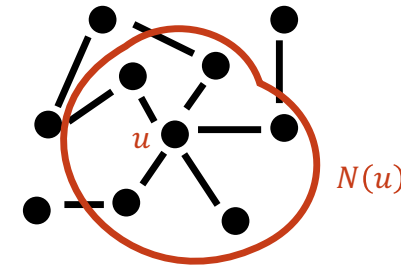- Data: user ratings and users' friends

- Assume:

  1. Friends influence your preferences

  2. Different levels of trusts for different friends



$$N(u)$$

[Chaney et al. '15]

# A. Social Network

- Data: user ratings and users' friends



- Assume:

1. Friends influence your preferences

2. Different levels of trusts for different friends

**Model** $\quad S_{ui} := \theta_u^\top \beta_i + \sum_{u' \in N(u)} \tau_{un} S_{u'i}$

The rating of $u'$ on item $i$

How much $u$ "trusts" $u'$

[Chaney et al. '15]

# A. Social Network

- Recent models use Graph Convolutional Networks (GCNs)

  - Powerful model for graph data

# A. Social Network

- Recent models use Graph Convolutional Networks (GCNs)

  - Powerful model for graph data

# B. Item Content

Data: user ratings and item text/image/…

**Model** $\quad S_{ui} := \theta_u^\top (\beta_i + \gamma_i)$

[Wang et al. '14]

# B. Item Content

Data: user ratings and item text/image/…

**Model** $S_{ui} := \theta_u^\top (\beta_i + \gamma_i)$

**Content features**

[Wang et al. '14]

# Questions?

# C. Dynamic Modelling

- Data: user ratings with timestamps

- Assume:
  - User tastes change over time
  - Item popularity change over time

**Model**
$$S_{ui}^t := \theta_u^{t\top} \beta_i^t$$
$$\theta^t = \theta^{t-1} + \epsilon$$

# C.1 Session-Based Modelling

- Data: user ratings with timestamps

- Assume: Users consume related items over short periods of time
  - Domains: Music playlist, exercises, short videos
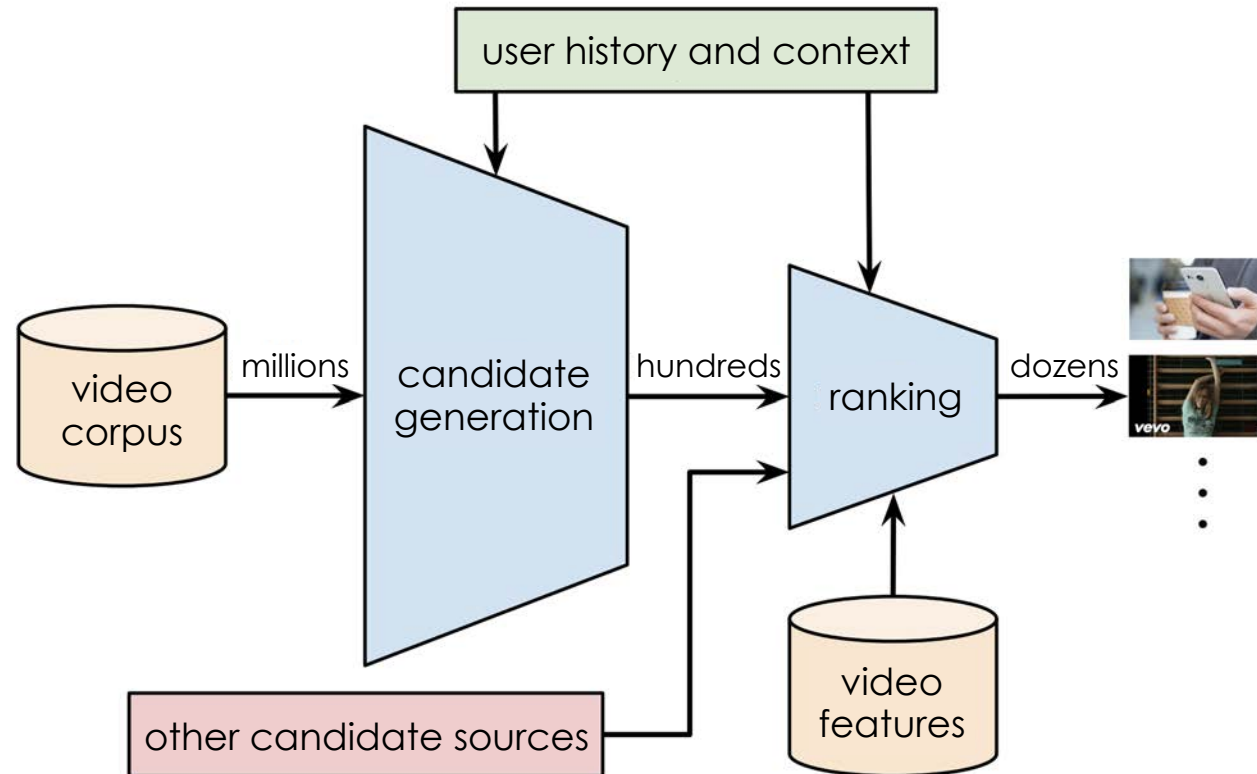
- Model. Sequential models like RNNs.



**Next movie prediction**

# Session-based + Social Networks



Figure 2: A schematic view of our proposed model for dynamic social recommandation

[Song et al. 2019]

# An Example From YouTube



[Covington et al., '16]

# Evaluation

- Evaluate performance on held-out data (standard)

- Splitting data into train/validation/test:

  - Split by user to give equal "weight" to each user

  - Ensure that each user has enough data (no cold-start)

# Evaluation Metrics

1. Score prediction (explicit data only)

   - Mean squared error: $\frac{1}{|\mathbf{users}|} \sum_u \sum_i (S_{ui} - \hat{S}_{ui})^2$

# Evaluation Metrics

1.  Score prediction (explicit data only)

    - Mean squared error: $\frac{1}{|\textbf{users}|}\sum_u \sum_i \left(S_{ui} - \hat{S}_{ui}\right)^2$

2.  Information retrieval

    - Precision, Recall
    - Average rank, Mean average precision
    - Normalized Discounted Cumulative Gain (NDCG)
        - Compares the ranking of your system with the optimal ranking
        - (Exponentially) Discounts items lower ranked items
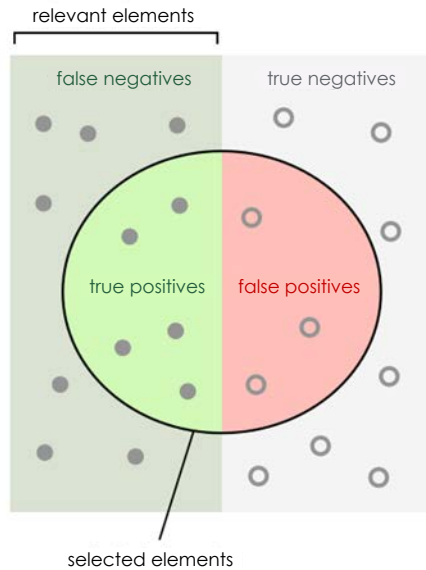
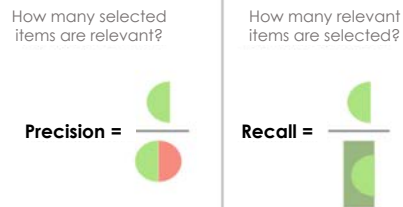# Precision/Recall



[From Wikipedia]

# Precision/Recall


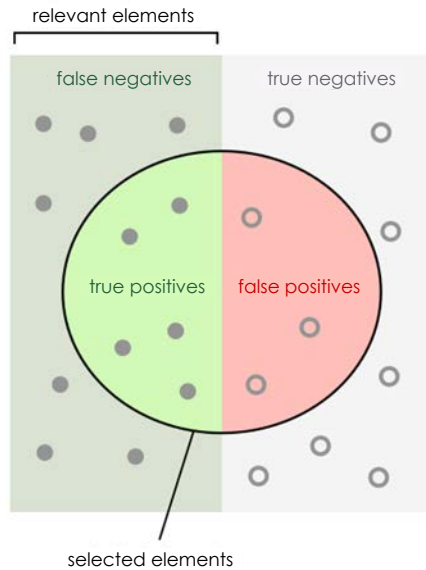
[From Wikipedia]

- For implicit data recall is more appropriate
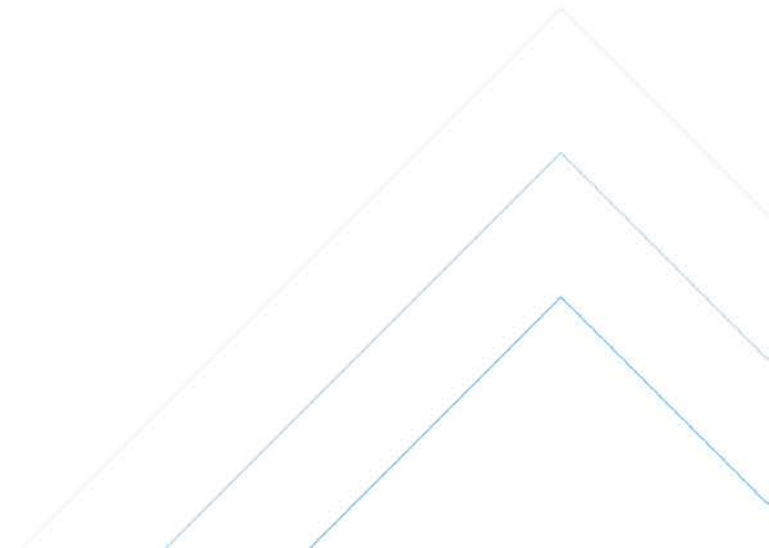
- Recall $:= \dfrac{TP}{TP+FN}$

# Precision/Recall



[From Wikipedia]

- For implicit data recall is more appropriate

- Recall := $\frac{TP}{TP+FN}$

- Consider only the top items (Recall@K)

# Other Topics

- Lots of other possible signals

  - Search queries, engagement (time spent on page)

- Structured recommendations

  - E.g., Recommend a trip, a curriculum of courses

# Concluding Remarks (I)

Type of models we have discussed are useful for:

- Domains with large number of items (and users for CF)

- Subjective preferences over attributes (features)
    - E.g., movies and not plane tickets

- Items can be consumed relatively fast
    - E.g., restaurants/movies and not cars/houses

# Concluding Remarks (II)

- CF models "work well" especially in large-data regimes

  - Commercial systems are reasonably good

  - There is evidence that companies derive value from them

- Much progress remains to be done

  - Modelling preferences is a very active research topic

  - Good preference models gave rise to other questions

# References

- Marital satisfaction and break-ups differ across on-line and off-line meeting venues,
  J. Cacioppo et al., PNAS'13

- A Probabilistic Model for Using Social Networks in Personalized Item Recommendation, A. Chaney et al., Recsys'15

- Deep Neural Networks for YouTube Recommendations, P. Covington et al., Recsys'16

- Content-based recommendations with Poisson factorization, P. Gopalan et al., NIPS'14

- Scalable Recommendation with Hierarchical Poisson Factorization., P. Gopalan et al., UAI'15

- Collaborative Filtering for Implicit Feedback Datasets, Hu et al., ICDM'08

- Exposure Matrix Factorization, D. Liang et al., WWW'16

- Factorization Machines, S. Rendle, ICDM '10

- Probabilistic Matrix Factorization, R. Salakhutdinov and A. Mnih, NIPS'08

- Session-based Social Recommendation via Dynamic Graph Attention Networks, W. Song et al., WSDM'19

- Collaborative Topic Modeling for Recommending Scientific Articles, C. Wang and D. Blei, KDD'11

- Deep matrix factorization models for recommender systems, HJ Xue et al., IJCAI'17