## RECOMMANDER SYSTEMS

Recommander Systems are algorithms aimed at suggesting relevant items to users (items = movies to watch, songs to listen, text to read, products to buy, ...)

Ex: Youtube, Spotify, Amazon, Netflix, ...

There are two major paradigms for recommander systems: (i) collaborative filtering and (ii) content based methods.
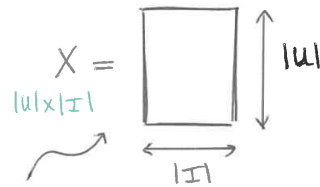
### • Collaborative Filtering (CF) methods.

These are based solely on past interactions recorded between users and items. These interactions are stored in a "user-item interaction" matrix $X$.

user $u \in U$ = set of users
item $i \in I$ = set of items

$|U|$ = # of users
$|I|$ = # of items

$$X = \qquad |U|$$
$$|U| \times |I|$$
$$|I|$$

A sparse matrix
The main idea is that past-interactions are sufficient to detect similar users and items and to make recommendations (fill the missing values in $X$)

CF methods are divided into two sub-categories:

(a) Memory based approaches : no model ; based on nearest neighbour searches

(b) Model based : assumes the existence of an underlying generative

model that explains the user-item interactions (matrix factorization techniques).

CF suffer from the "cold-start" problem : impossible to recommend anything to new users or to recommend any new item to any user.

### • Content - Based approaches.

These approaches use additional information about users and/or items.

Ex: (i) Movies : actors, directors, genre, duration, ...
(ii) User : age, sex, job, ...
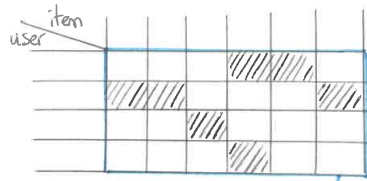↖ These are called user features & item features.

CB methods are divided into (i) classification problems (predict whether a user likes or not an item) and (ii) regression problems (predict the rating of an item). They suffer less from the cold-start problem.

## I. COLLABORATIVE FILTERING

### I.1. Memory - based approaches.

This category corresponds to neighborhood-based collaborative filtering methods. For a certain user, user-based k-nearest neighbour methods first identify a set of similar users, and then recommend items based on what items those similar users have purchased. Similarly, item-based k-nearest neighbour methods first identify a set of similar items for each of the items a user has purchased, and then recommend items based on those similar items.
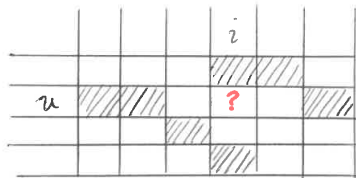
Goal: Fill missing ratings



$X =$ user-item interaction matrix

▨ = known rating

First-Attempt :

<u>User-User Collaborative Filtering</u> vs <u>Item-Item Collaborative Filt.</u>



Estimate $\hat{r}_{ui}$ of rating $r_{ui}$ ?

$$\hat{r}_{ui} = \frac{1}{|U_i^+|} \sum_{\substack{v \in U_i^+ \\ v \neq u}} r_{vi}$$

vs

$$\hat{r}_{ui} = \frac{1}{|I_u^+|} \sum_{\substack{j \in I_u^+ \\ j \neq i}} r_{uj}$$

$U_i^+ = \{ u \in U : (u,i) \in S \}$
$S$ = set of known ratings

$I_u^+ = \{ i \in I : (u,i) \in S \}$
= set of items that user $u$ has rated.

x Better : consider a neighborhood of $(u,i)$.

$\mathcal{N}(u;i)$ = neighborhood of $u$
[users similar to $u$ that have also rated item $i$]

$\mathcal{N}(i;u)$ = neighborhood of $i$
[rated by $u$ and similar to $i$]

$$\hat{r}_{ui} = \frac{1}{k} \sum_{v \in \mathcal{N}(u;i)} r_{vi}$$

$$\hat{r}_{ui} = \frac{1}{k} \sum_{j \in \mathcal{N}(i;u)} r_{uj}$$

Consider the k nearest neighbours only ($\to$ hyperparameter)

---

$\Rightarrow$ We need to define what it means for users to be similar, and for items to be similar.

x <u>Option 1</u> : Jaccard similarity

$$sim(u,v) = \frac{|r_u \cap r_v|}{|r_u \cup r_v|} = \frac{\text{\# items that both } u \& v \text{ have rated}}{\text{total \# of items rated by } u \& v \text{ together}}$$

similarity between two users $u$ & $v$
(similarly define $sim(i,j)$ between two items $i$ and $j$)

= "Intersection over Union".

$\hookrightarrow$ Not taking into account how users liked items ; only how many they have in common.

x <u>Option 2</u> : (Adjusted) cosine similarity (CS)

$$sim(u,v) = \cos(r_u, r_v) = \frac{\langle r_u, r_v \rangle}{\|r_u\| \cdot \|r_v\|} = \frac{\sum\limits_{i \in I_{u,v}} r_{ui} r_{vi}}{\sqrt{\sum\limits_{i \in I_{u,v}} r_{u,i}^2} \sqrt{\sum\limits_{i \in I_{u,v}} r_{v,i}^2}} ,$$

where $I_{u,v}$ = set of items both rated by $u$ and $v$.
The adjusted CS substracts the users' average ratings.

x <u>Option 3</u> : Pearson similarity

$$sim(u,v) = \frac{\sum\limits_{i \in I_{u,v}} (r_{u,i} - \bar{r}_u)(r_{v,i} - \bar{r}_v)}{\sqrt{\sum\limits_{i \in I_{u,v}} (r_{u,i} - \bar{r}_u)^2 \sum\limits_{i \in I_{u,v}} (r_{v,i} - \bar{r}_v)^2}}$$

Rating predictions :

$$\hat{r}_{ui} = \frac{\sum\limits_{v \in \mathcal{N}(u;i)} sim(u,v)\, r_{vi}}{\sum\limits_{v \in \mathcal{N}(u;i)} sim(u,v)}$$

or

$$\hat{r}_{ui} = \frac{\sum\limits_{j \in \mathcal{N}(i;u)} sim(i,j)\, r_{uj}}{\sum\limits_{j \in \mathcal{N}(i;u)} sim(i,j)}$$

user-based CF ↗

item-based CF ↗

## I.2. Explicit Matrix Factorization using SVD.

To gain some insight on how matrix factorization models work, we first review how PCA & SVD work.

Let $X$ be our matrix of observations (the user-interaction matrix in our case), of size $n \times d$ ($n = |U| =$ number of users; $d = |I| =$ number of items).

Let $X = U \Lambda V^t =$ SVD decomposition of $X$ of rank $d < n$.

$(n \times d)$ $(n \times d)$ $(d \times d)$
$(d \times d)$

- $U =$ has orthonormal columns $U^t U = I_d$
  $(n \times d)$
- $V =$ has orthonormal columns $V^t V = I_d$
  $(d \times d)$
- $\Lambda =$ diagonal matrix of singular values
  $(d \times d)$ $= \text{diag}(\lambda_1, \ldots, \lambda_d)$.

↳ Recall that the $j$-th principal component of $X$ is given by the eigenvector of the sample covariance matrix, associated with the $j$-th largest eigenvalue, see  UL: PCA .

Assuming that the colums of $X$ are centered, the sample covariance matrix $S = \frac{1}{n} X^t X = \frac{1}{n} V \Lambda^2 V^t$
$(d \times d)$

$\Rightarrow$ (eigenvalue - eigenvector pairs) $= (\frac{1}{n} \lambda_i^2, v_i)$

$i$-th column of $V = \begin{pmatrix} | & & | \\ v_1 & \cdots & v_d \\ | & & | \end{pmatrix}$
$(d \times d)$

The PCs of $X$ are contained in the colums of $V$, obtained from the SVD decomposition $X = U \Lambda V^t$.

↑ PCs have an interpretation in terms of low rank approximation of $X$. Indeed, it can be shown that the best (with respect

to Frobenius norm) rank $r$ approximation to $X$ is given by $X^* = U \Lambda^* V^t$, where $\Lambda^*$ is obtained from $\Lambda$, setting the smallest $\lambda_{r+1}, \ldots, \lambda_d$ singular values to zero.

Then $\qquad X^* = (X v_1) v_1^t + \cdots + (X v_r) v_r^t = \begin{pmatrix} - \hat{x}_1^t - \\ \vdots \\ - \hat{x}_n^t - \end{pmatrix}$.

The $j$-th row of this matrix is

$$\hat{x}_j = \langle x_j, v_1 \rangle v_1 + \cdots + \langle x_j, v_r \rangle v_r$$

↖ The best rank $r$ approximation is obtained as a linear combination of the first $r$ principal components
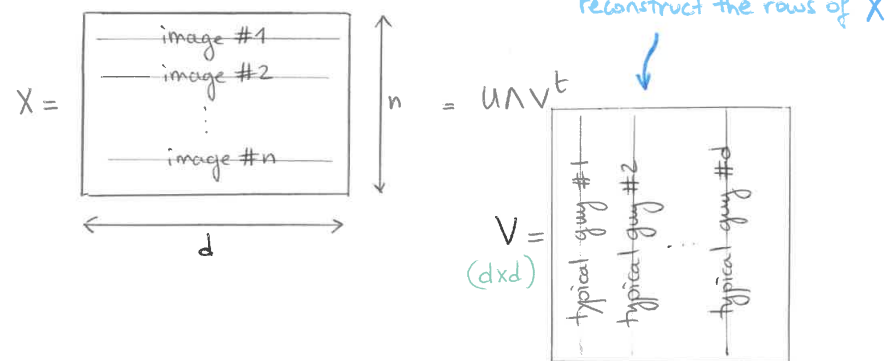
Ex: suppose that $X =$ collection of images.
Then $v_1, \ldots, v_d$ are called  eigenfaces . Each original image in $X$ can be reconstructed from a linear combination of the eigenfaces.

↳ Each vector $v_i$ / eigenface represent one specific aspect underlying the data. In an ideal world,
$v_1 \approx$ a typical elder person
$v_2 \approx$ a typical glasses wearer
$v_3 \approx$ a typical sad looking person, ... / ...

⇓

colums of $V$ are used to reconstruct the rows of $X$

$X = \begin{bmatrix} \text{image \#1} \\ \text{image \#2} \\ \vdots \\ \text{image \#n} \end{bmatrix}$ $\updownarrow n$ $= U \Lambda V^t$

$\xleftrightarrow{\quad d \quad}$

$V = \begin{bmatrix} \text{typical guy \#1} & \text{typical guy \#2} & \cdots & \text{typical guy \#d} \end{bmatrix}$
$(d \times d)$

⇒ If instead of a matrix of images we have a user-item interaction matrix, in an ideal world, the interpretation of the columns of $V$ would be:

$v_1$ = typical action movie fan
$v_2$ = typical romance movie fan
$v_3$ = typical sci-fy movie fan, ... /...

context of movie rating



$$X = \underset{(n \times d)}{\begin{bmatrix} \text{—— user 1 ——} \\ \text{—— user 2 ——} \\ \vdots \\ \text{—— user n ——} \end{bmatrix}} = U \wedge V^t \; ; \quad V = \underset{(d \times d)}{\begin{bmatrix} | & | & | & | \\ \text{action movie fan} & \text{romance mov. fan} & \text{sci-fi movie fan} & \cdots \\ | & | & | & | \end{bmatrix}}$$

$v_1 \quad v_2 \quad v_3 \quad \quad v_d$

And similarly, taking the transpose of $X$,

$$X^t = \underset{(d \times n)}{\begin{bmatrix} \text{—— movie 1 ——} \\ \text{—— movie 2 ——} \\ \vdots \\ \text{—— movie d ——} \end{bmatrix}} = V \wedge U^t \; ; \quad \underset{(n \times d)}{U} = \begin{bmatrix} | & | & | \\ \text{typical action movie} & \text{typical romance mov} & \text{typical sci-fi mov} \\ | & | & | \end{bmatrix}$$

The columns of $U$ contain the eigenvalues of $X X^t$; related to the sample covariance matrix of $X^t$.

In this case, PCA does not reveal typical faces or typical users, but typical movies. These typical movies can be used to build back the original movies.

↳ SVD decomposition reveals all typical users and typical movies in one go; a very powerful tool.

---

r latent factors



$$\underset{(n \times d)}{X} = \quad \approx \quad \cdots \quad = U_r \wedge_r V_r^t$$

n users, d items, typical movie, typical user

keep only r PCs
& set $\lambda_{r+1}, ..., \lambda_d$ to zero; or keep only the first r columns of $U$ and $V$ (denoted $U_r$, $V_r$)

⚠ $X$ is usually sparse, with many missing values, and the SVD decomposition is not defined.

A rating $\hat{r}_{ui}$ in $X$ (entry corresponding to user $u$, item $i$) is expressed as a dot product between a vector $p_u \in \mathbb{R}^d$ and $q_i \in \mathbb{R}^d$ (up to a scaling factor):

$$\hat{r}_{ui} = \langle p_u, q_i \rangle = p_u^t q_i$$

≡ "user-item" interaction.

$$= \underset{\substack{f \in \text{latent} \\ \text{factor}}}{\sum} \text{affinity of } u \text{ for } f \times \text{affinity of } i \text{ for } f$$

## I.3. Explicit Matrix Factorization: ALS, SGD, ..

The user-item interaction matrix contains many missing entries, so its SVD decomposition is undefined & cannot be computed. An easy fix would be to replace each missing value with a row mean, column mean, or 0, and compute the SVD decomposition of the filled matrix $X$. A better approach would be

to solve the optimization problem

$$\min_{p_u, q_i} \left\{ \sum_{r_{ui} \in X} (r_{ui} - p_u^t q_i)^2 + \lambda_p \sum_u \|p_u\|^2 + \lambda_q \sum_i \|q_i\|^2 \right\}$$

$$\| \\ \mathcal{L}(p_u, q_i)$$

add regularization terms to prevent overfitting.

The are two popular methods for minimizing this criterion: Alternating Least Squares (ALS) and Stochastic Gradient Descent (SGD).

### x Alternating Least Squares (ALS).

For ALS minimization, we hold one set of latent vectors constant ( $q_i$, say ), and optimize w.r.t. $p_u$ (a simple LS problem). With this updated value of $p_u$, we hold it constant, and optimize w.r.t. $q_i$. And Repeat.

↳ derivation.

$$\frac{\partial \mathcal{L}}{\partial p_u} = -2 \sum_i (r_{ui} - \hat{p}_u^t q_i)^2 q_i^t + 2 \lambda_p \hat{p}_u^t = 0$$

$$\Leftrightarrow -\sum_i (r_{ui} - \hat{p}_u^t q_i)^2 q_i^t + \lambda_p \hat{p}_u^t = 0$$

$$\sum_i r_{ui} q_i^t \qquad p_u^t \sum_i q_i q_i^t = p_u^t Q^t Q$$
$$\|$$
$$(r_{u1} \cdots r_{ud}) \begin{pmatrix} - q_1^t - \\ \vdots \\ - q_d^t - \end{pmatrix} \qquad \Leftrightarrow \hat{p}_u^t (Q^t Q + \lambda_p I_r) = r_u^t Q$$
$$=: r_u^t Q \qquad \Leftrightarrow \hat{p}_u^t = r_u^t Q (Q^t Q + \lambda_p I_r)^{-1}$$
$$(1 \times d)\,(d \times r) \qquad \Leftrightarrow \boxed{\hat{p}_u = (Q^t Q + \lambda_p I_r)^{-1} Q^t r_u}$$

---

And similarly,

$$\frac{\partial \mathcal{L}}{\partial q_i} = -2 \sum_u (r_{ui} - p_u^t \hat{q}_i) p_u^t + 2 \lambda_q \hat{q}_i^t = 0$$

$$q_i^t p_u$$

Introducing $\underset{(n \times r)}{P} := \begin{pmatrix} - p_1^t - \\ \vdots \\ - p_n^t - \end{pmatrix}$ and $\underset{(n \times 1)}{r_i} := \begin{pmatrix} r_{u_1 i} \\ \vdots \\ r_{u_n i} \end{pmatrix}$,

we get

$$\boxed{\hat{q}_i = (P^t P + \lambda_q I_r)^{-1} P^t r_i}$$

↙ The ridge solution, see
SL: RR AND LASSO

### x Stochastic Gradient Descent

Here we take the derivative of the loss function with respect to each variable in the model. Samples at taken one at a time, or in batch.

We consider the more general model:

$$\hat{r}_{ui} = \mu + b_u + b_i + p_u^t q_i$$

interaction term.

global mean

user bias = some users tend to rate higher than others on average

item bias = some items are more popular than others

The loss becomes:

$$\mathcal{L}(p_u, q_i, b_u, b_i, \mu) = \sum_{r_{ui} \in X} (r_{ui} - (\mu + b_u + b_i + p_u^t q_i))^2$$
$$+ \lambda_{pb} \sum_u \|b_u\|^2 + \lambda_{qb} \sum_i \|b_i\|^2$$
$$+ \lambda_{pf} \sum_u \|p_u\|^2 + \lambda_{qf} \sum_i \|q_i\|^2$$

We need the gradient w.r.t. each parameter.

- $\dfrac{\partial \mathcal{L}}{\partial b_u} = -\cancel{2}\Big(\underbrace{r_{ui} - (\rho + b_u + b_i + p_u^t q_i)}_{\substack{\| \\ \varepsilon_{ui}\, =\, \text{current prediction error}}}\Big) + \cancel{2}\,\lambda_{pb}\, b_u = 0$ ⑪

$\dfrac{\partial \mathcal{L}}{\partial b_u} = 0 \iff \lambda_{pb}\, b_u - \varepsilon_{ui} = 0$

- $\dfrac{\partial \mathcal{L}}{\partial b_i} = 0 \iff \lambda_{qb}\, b_i - \varepsilon_{ui} = 0$

- $\dfrac{\partial \mathcal{L}}{\partial p_u} = 0 \iff \lambda_{pf}\, p_u - \varepsilon_{ui}\, q_i = 0$

- $\dfrac{\partial \mathcal{L}}{\partial q_i} = 0 \iff \lambda_{qf}\, q_i - \varepsilon_{ui}\, p_u = 0$

The SGD updates are:

$$\times\ b_u \leftarrow b_u + \eta\,(\varepsilon_{ui} - \lambda_{pb}\, b_u)$$
$$\times\ b_i \leftarrow b_i + \eta\,(\varepsilon_{ui} - \lambda_{qb}\, b_i)$$
$$\times\ p_u \leftarrow p_u + \eta\,(\varepsilon_{ui}\, q_i - \lambda_{pf}\, p_u)$$
$$\times\ q_i \leftarrow q_i + \eta\,(\varepsilon_{ui}\, p_u - \lambda_{qf}\, q_i)$$

### I.4. Implicit Matrix Factorization.

With implicit feedback, we do not have ratings anymore, we have user's preferences for items (e.g. a like on Instagram). The criterion to optimize must be modified accordingly. The following criterion was proposed in <u>Collaborative Filtering for Implicit Feedback Datasets</u>, Hu, Koren & Volinsky.

---

We introduce a set of binary variables $d_{ui}$, which indicates the preference of user $u$ to item $i$. The $d_{ui}$ are obtained by binarizing the $r_{ui}$ values:

$$d_{ui} = \begin{cases} 1 & \text{if } r_{ui} > 0 \\ 0 & \text{otherwise} \end{cases}$$

↑ If user $u$ consumed item $i$ ($r_{ui} > 0$), we have an indication that $u$ likes it ($d_{ui} = 1$). On the other hand, if $u$ never consumed $i$, we believe no preference ($d_{ui} = 0$).

However, values $d_{ui} = 0$ are associated with low confidence: not taking any positive action on an item can have many other reasons than not liking it (unaware, or unable to consume it since too pricy).

In addition $d_{ui} = 1$ can be the result of many $\neq$ reasons: gift to a friend even if they dislike the product, watching a TV show because staying on the same channel as before...

In general, as $r_{ui}$ grows, we have a stronger indication that the user likes it.

$\Rightarrow$ Introduce a set of variables $c_{ui}$ which indicate/measure our confidence in observing $d_{ui}$. The authors take

$$c_{ui} = 1 + \alpha\, r_{ui}$$

↑ hyperparameter

↑ e.g. $r_{ui} = \#$ of clicks for item $i$. The more clicks, the more likely user $u$ likes $i$.

⇓

$$\mathcal{L} = \sum_{\forall (u,i)} c_{ui}\,(d_{ui} - p_u^t q_i)^2 + \lambda_p \sum_u \|p_u\|^2 + \lambda_q \sum_i \|q_i\|^2$$

(use ALS)

# I.5. Sparse Linear Methods (SLIM).

All methods described above correspond to Matrix Factorization methods since they factorize the user-item interaction matrix $X$ into a product of low-rank user factors and item factors. Similarly, SLIM learns a sparse matrix $W$ by solving an $\ell_1 + \ell_2$-norm regularized optimization problem.

↗ MF methods approximate $X$ using $PQ^t$
↗ user & item factors

↘ SLIM = a special case setting $P \equiv X$, and $W \equiv Q^t$

SLIM does not learn user representation
⇒ this will simplify the learning process.
(& information is fully preserved)

to be learned, a $(d \times d)$ matrix

$$\text{minimize}_{W \atop (d \times d)} \quad \frac{1}{2} \| X - X W \|_F^2 + \frac{\beta}{2} \| W \|_F^2 + \lambda \| W \|_1$$
$$\text{subject to} \quad W \geqslant 0$$
$$\text{diag}(W) = 0$$

($X$ is $(n \times d)$)

- $\ell_1$-penalization ensures a sparse solution.
- $\text{diag}(W) = 0$ ensures that $x_{ij}$ is not used to compute $\tilde{x}_{ij} := (\tilde{X})_{ij}$ ; $\tilde{X} := XW$.

↳ the optimization problem can be decoupled into a set of sub-problems $\quad \text{minimize}_{w_j} \frac{1}{2} \| x_j - X w_j \|_2^2 + \frac{\beta}{2} \| w_j \|_2^2 + \lambda \| w_j \|_1$
subject to $\begin{cases} w_j \geqslant 0 \\ w_{jj} = 0 \end{cases}$

$j$-th column of $W$

An elastic net problem (use e.g. CDA)

$\| \cdot \|_F$ = Frobenius norm
$\| W \|_1 = \sum_{i,j} |w_{ij}|$

---

# II. LIGHT FM: LEARNING WITH SIDE INFORMATION

A simple idea to incorporate side information : one-hot-encode each (user and/or item) feature into an attribute space, and assume that each attribute has its own latent vector.

The user vector is then decomposed into two components.

$$p_u + \sum_{a \in \mathcal{N}(u)} s_a$$

↗ All $p_u$ & $s_a$ are learned

↘ set of attributes of user $u$.

× **Light FM model.**

In the above decomposition, the user vector is the sum of two parts. Light FM treats everything as side information, or features. If we want to have a specific vector for each user, one must one-hot-encode that as a single feature for that user.

↳ latent representation of user $u$ is $\quad p_u := \sum_{a \in \mathcal{N}(u)} s_a$

——— " ——— of item $i$ is $\quad q_i := \sum_{\alpha \in \mathcal{M}(i)} u_\alpha$

**Remark** : (i) Without side information, one-hot-encoding users (and/or item) yields

$$\begin{pmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{pmatrix} \begin{pmatrix} - s_1^t - \\ \vdots \\ - s_n^t - \end{pmatrix}.$$

$(n \times n)$      $(n \times r)$
to be learned

aka the "feature matrix"

(ii) With side information, consider



<span style="color:green">↑ n ↓</span> $\begin{pmatrix} 1 & & 0 & \vdots & \\ & \diagdown & & \vdots & * \\ 0 & & 1 & \vdots & \end{pmatrix}$

<span style="color:green">← n →</span>

<span style="color:green">length of the set of user features</span>

This matrix is passed on to LightFM for training in user_features. The default value for user_features is nothing is passed on is the Id matrix.

Each feature is also described by a scalar bias term $\beta_a^u$ & $\beta_\alpha^I$ :

$$b_u = \sum_{a \in W(u)} \beta_a^u \quad ; \quad b_i = \sum_{\alpha \in W(i)} \beta_\alpha^I .$$

LightFM model is $\hat{r}_{ui} = f\left( p_u^t q_i + b_u + b_i \right)$

This model is able to compute recommendation for new users and items.

In the original paper ( Metadata Embeddings for User and Item Cold-start Recommendations, Maciej Kula (2015) ), the author is interested in predicting binary data : positive interactions denoted 1, and negative interactions denoted 0. For this reason, the author takes $f(x) = \dfrac{e^x}{1+e^x}$ = sigmoid function.

An identity function would work well in predicting ratings.

Remark: If the feature sets consist solely of indicator variables for each user and item, LightFM reduces to the standard matrix factorization approach.

---

× Criterion

We discuss three criteria + learning algorithms to fit a LightFM model.

(a) Likelihood .

The author originally proposed to maximize the likelihood of the data:

$$\mathcal{L}\left( p_u, q_i, b_u, b_i \right) = \prod_{(u,i) \in S_+} \hat{r}_{ui} \prod_{(u,i) \in S_-} (1 - \hat{r}_{ui}) .$$

set of positive interactions

—"— negative —"—

$\mathcal{L}$ is maximized using SGD + AdaGrad.

(b) Learning to Rank — WARP / k-os WARP.

This technique and the one presented next are novel alternatives to ALS and SGD, in the context of implicit feedback. WARP stands for Weighted Approximate-Rank Pairwise loss, and was introduced in WSABIE: Scaling Up To Large Vocabulary Image Annotation, J. Weston, S. Bengio & N. Usunier (2011).

It was introduced in the context of image annotation, where each image $x$ is given a set of possible annotations $\{y\}$ to be ranked. Each image is assumed to has only one true annotation.

Specifically, let $\varphi(x) \in \mathbb{R}^d$ denote a scoring function for each annotation $y$ associated with image $x$ :

$$\varphi(x) = \begin{pmatrix} \varphi_1(x) \\ \vdots \\ \varphi_d(x) \end{pmatrix} \in \mathbb{R}^d$$

<span style="color:green">What we called $r_{xy}$ before, and $\hat{r}_{xy}$ for their estimated versions.</span>

Let $\text{rank}_y(\varphi(x))$ denote the rank of annotation $y$ given by $\varphi(x)$:

$$\text{rank}_y(\varphi(x)) := \sum_{j \neq y} \mathbb{1}\left(\varphi_y(x) \leq \varphi_j(x)\right) \quad\text{———(*)}$$

↑ top rank is zero

The rank is transformed into a loss via a function $L$:

$$L\left(\text{rank}_y[\varphi(x)]\right) = \text{loss for ranking true annotation } y \text{ for user } x$$
using $\varphi$. (assuming there is only one true annotation)

where

$$L(k) := \sum_{j=1}^{k} \alpha_j \quad , \quad \alpha_1 \geq \alpha_2 \geq \ldots \geq 0$$

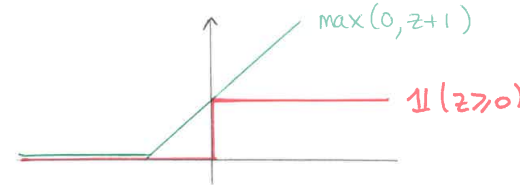→ $\alpha_1 = \ldots = \alpha_d = \frac{1}{d-1}$   minimizes the mean rank

→ $\alpha_1 = 1$ ; $\alpha_{j>1} = 0$   optimizes the proportion of top-ranked correct labels.

→ larger values of $\alpha$ in the top positions optimize the top $k$ in the ranked list (good for optimizing precision at $k$).

We see from (*) that $\sum_{j \neq y} \dfrac{\mathbb{1}(\varphi_y(x) \leq \varphi_j(x))}{\text{rank}_y[\varphi(x)]} = 1,$

so that

$$L\left(\text{rank}_y[\varphi(x)]\right) = \left( \sum_{j \neq y} \frac{\mathbb{1}(\varphi_y(x) \leq \varphi_j(x))}{\text{rank}_y[\varphi(x)]} \right) \times L\left(\text{rank}_y[\varphi(x)]\right)$$

$$= \sum_{j \neq y} L\left(\text{rank}_y[\varphi(x)]\right) \frac{\mathbb{1}(\varphi_y(x) \leq \varphi_j(x))}{\text{rank}_y[\varphi(x)]}.$$

$\max(0, z+1)$
$\mathbb{1}(z \geq 0)$

We consider the convex surrogate given by the hinge loss to upper-bound the term $\mathbb{1}\left(\varphi_j(x) - \varphi_y(x) \geq 0\right)$

$$\text{loss} \leq \sum_{j \neq y} L\left(\text{rank}_y^1[\varphi(x)]\right) \frac{\max\left(0, \varphi_j(x) - \varphi_y(x) + 1\right)}{\text{rank}_y^1[\varphi(x)]}$$

where we replaced the $\text{rank}_y[\varphi(x)]$ functions by the margin-penalized rank of $y$:

$$\text{rank}_y^1[\varphi(x)] = \sum_{j \neq y} \mathbb{1}\left(\underbrace{1 + \varphi_j(x)}_{\text{margin of size } 1} > \varphi_y(x)\right) \quad (\geq \text{rank}_y[\varphi(x)])$$

The expected risk of $\varphi$ is:

$$R(\varphi) := \mathbb{E}_{(X,Y)} \ell(Y, \varphi(X)),$$

where

$$\ell(y, \varphi(x)) := \sum_{j \neq y} L\left(\text{rank}_y^1[\varphi(x)]\right) \frac{\max\left(0, \varphi_j(x) - \varphi_y(x) + 1\right)}{\text{rank}_y^1[\varphi(x)]}.$$

↑ Criterion to minimize.

The authors suggest the following sampling procedure
  (i) Select a pair $(x, y)$ according to $\mathbb{P}$
  (ii) For the chosen pair $(x, y)$ select a violating label $\bar{y}$ such that $1 + \varphi_{\bar{y}}(x) > \varphi_y(x)$.

Indeed, note that

$$R(\varphi) = E_{X,Y}\left[\left\{\sum_{j \neq Y} L(\text{rank}_Y^1[\varphi(X)])\max(0, \varphi_j(X) - \varphi_Y(X)+1)\right\} \times \frac{1}{\text{rank}_Y^1 \varphi(X)}\right]$$

$$= E_{X,Y} E_{\bar{Y}|X,Y}\left\{L(\text{rank}_Y^1[\varphi(X)])\max(0, \varphi_{\bar{Y}}(X) - \varphi_Y(X)+1)\right\}$$

$$= E_{X,Y,\bar{Y}}\left\{L(\text{rank}_Y^1[\varphi(X)]\max(0, \varphi_{\bar{Y}}(X) - \varphi_Y(X)+1)\right\}$$

where

$$\mathbb{P}(\bar{Y}=j \mid X=x, Y=j) = \begin{cases} \dfrac{1}{\text{rank}_y^1[\varphi(x)]} & \text{if } 1+\varphi_j(x) > \varphi_y(x) \\ 0 & \text{otherwise} \end{cases}$$

↑ sums to 1 since

$$\sum_{j \neq y}\mathbb{P}(\bar{Y}=j \mid X=x, Y=y) = \sum_{\substack{j: \\ 1+\varphi_j(x) > \varphi_y(x)}} \frac{1}{\text{rank}_y^1[\varphi(x)]} = 1$$

Goal = minimize a function $R(\varphi)$ expressed as the expected value of another one $\Rightarrow$ we can use the Robbins & Monro (1951) procedure to update the parameters $\beta$ of the model

$$\beta \leftarrow \beta - \gamma_t \frac{\partial}{\partial \beta}\left\{L(\text{rank}_y^1[\varphi(x)]\max(0, \varphi_{\bar{y}}(x) - \varphi_y(x)+1)\right\}$$

↑ Parameters governing the expression of $\varphi$.

expensive to compute since $\varphi_i(x)$ must be calculated $\forall i$

---

Remarks: (i) To select $\bar{y}$ randomly, we must know in advance with labels are violators, i.e. we must compute $\varphi_i(x)$ $\forall i$. To get around this issue, the authors suggest to sample labels $i$ uniformly with replacement until a violating label is found.

(ii) Following the procedure above, assuming there are $k = \text{rank}_y^1[\varphi(x)]$ violating labels, the number of trials $N_k$ in the sampling step thus follows a geometric distribution with probability of success $\frac{k}{d-1}$.

Thus $\mathbb{E}N_k = \frac{d-1}{k}$, and we may approximate $\text{rank}_y^1[\varphi(x)]$ by $\left\lfloor \frac{d-1}{N} \right\rfloor$; $N$ number of trials.

Since there are precisely $k$ points that satisfy $1+\varphi_j(x) > \varphi_y(x)$.

(iii) In the original paper, it is assumed that each image is assigned to a unique label. In the context of collaborative filtering, this amounts in assuming that each user is associated with a single positive item. WARP thus ignores the fact that there are multiple positive items per user, and treat those items independently (sample a positive example $(x, y)$ at random).

not unique

A modification of WARP, called k-os WARP (see Learning to Rank Recommendations with the k-order statistic loss, J. Weston, H. Yee, R. J. Weiss) uses the k-th positive

example for any user as a basis for pairwise
updates. Specifically :

→ Select a user $u$ at random
→ Draw randomly $K$ positive items for that user, and
  compute the score $\varphi_i(u)$ for each of the picked items.
→ Select as the positive example the $k$-th largest
  in the list.

*hyperparameter.*

(c) <u>Learning to Rank - BRP.</u>

An alternative to WARP is **Bayesian** **Personalized**
**Ranking** (BPR), introduced in the context of item
recommendation with implicit feedback. The optimization
criterion is derived from the maximum posterior estimator
of the model parameters. Similar to WARP, we
iteratively sample triplets $(u, i, j)$ of seen $i$ and
unseen $j$ items of a randomly picked user $u$. The procedure
goes as follows =

↳ If an item $i$ has been viewed by user $u$, we
  say that $i$ is a positive item. We write $(u,i) \in S$
↳ If $u$ has not viewed item $j$, we assume that
  $u$ prefers $i$ over $j$, and we write $i >_u j$.
↳ Consider the training set
$$\mathcal{L}_S := \{(u,i,j) \mid i \in I_u^+ , \; j \in I \setminus I_u^+\}$$

*No preference between two seen items, or btw to unseen items*

$$I_u^+ := \{i \in I \mid (u,i) \in S\}$$   *$I$ = set of all items*

---

↳ Likelihood function of the learning sample $\mathcal{L}_S$ :

$$\ell(\Theta \mid \mathcal{L}_S) = \prod_{(u,i,j) \in \mathcal{L}_S} p(i >_u j \mid \Theta)$$

Assuming a prior $p(\Theta)$ on $\Theta$, the posterior distribution
is proportional to $\ell(\Theta \mid \mathcal{L}_S) \, p(\Theta)$.

$$\| \quad \prod_{(u,i,j) \in \mathcal{L}_S} p(i >_u j \mid \Theta) \, p(\Theta)$$

$$:= \sigma(\hat{x}_{uij}(\Theta)), \text{ where}$$
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$\mathcal{N}(0, \Sigma_\Theta)$
with $\Sigma_\Theta = \lambda_\Theta I$
diagonal

$\hat{x}_{uij}(\Theta)$ = arbitrary real-valued
function parametrized by $\Theta$,
capturing the relationship
between user $u$ and items
$i$ and $j$. The authors suggest
to take $\hat{x}_{uij} = \hat{x}_{ui} - \hat{x}_{uj}$,
so that we can apply a standard
collaborative filtering model that
predicts $\hat{x}_{u\ell}$, $\forall \ell \in I$

*And indeed, user $u$ prefers item $i$ over $j$ iff*
*$\hat{x}_{uij} \geqslant 0$, i.e. $\hat{x}_{ui} \geqslant \hat{x}_{uj}$, i.e. iff $\sigma(\hat{x}_{uij}) \geqslant \frac{1}{2}$*

The MAP estimator is then

$$\hat{\Theta}_{MAP} \in \underset{\Theta}{\arg\max} \sum_{(u,i,j) \in \mathcal{L}_S} \log \sigma(\hat{x}_{uij}(\Theta)) - \lambda_\Theta \| \Theta \|^2$$

Since the criterion is differentiable, we can use a gradient descent algorithm for maximization:

$$\theta \leftarrow \theta + \alpha \left( \frac{e^{-\hat{x}_{uij}}}{1 + e^{-\hat{x}_{uij}}} \frac{\partial}{\partial \theta} \hat{x}_{uij} + \lambda_\theta \theta \right)$$

observation $(u, i, j)$ is chosen randomly.

The expression $\frac{\partial}{\partial \theta} \hat{x}_{uij}$ depends on the model considered.

× Ex: Matrix Factorization

Goal = estimation of a matrix $X$ of dimension $|U| \times |I|$

How = using the approximation $X \approx W H^t$, for low-rank matrices $W$ and $H$.

$|U| \times r$      $|I| \times r$

Then $\hat{x}_{ui} = \langle \omega_u, h_i \rangle = \sum_{\ell=1}^{r} w_{u\ell} h_{i\ell}$

row of $W$    column of $H$

[The model parameters are $\theta = (W, H)$]

We obtain

$$\frac{\partial}{\partial \theta} \hat{x}_{uij} = \begin{cases} h_{i\ell} - h_{j\ell} & \text{if} \quad \theta = w_{u\ell} \\ w_{u\ell} & \text{if} \quad \theta = h_{i\ell} \\ -w_{u\ell} & \text{if} \quad \theta = h_{j\ell} \\ 0 & \text{otherwise} \end{cases}$$

• Remark = Analogies to AUC optimization.

Reminder: Probabilistic interpretation of (ROC) AUC

Notation = $\mathcal{L}_n = \{(X_1, Y_1), ..., (X_n, Y_n)\}$ our learning sample, $Y_i \in \{0, 1\}$, used to construct a learning function $f_n$, which classifies a new observed point $x$ as $1$ if $f_n(x) \geq t$, and as $0$ otherwise, for some threshold $t$.
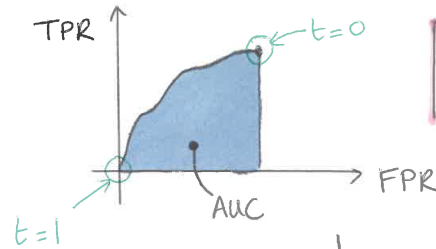
• TPR (aka recall / sensitivity) is

$$\frac{TP}{P} \approx \mathbb{P}(f_n(X) \geq t \mid Y=1, \mathcal{L}_n) = \bar{F}_1(t) = \int_t^1 f_1(u)\, du$$

• FPR (aka 1- specificity) is

$$\frac{FP}{N} \approx \mathbb{P}(f_n(X) \geq t \mid Y=0, \mathcal{L}_n) = \bar{F}_0(t) = \int_t^1 f_0(u)\, du$$



$$AUC = \int_1^0 \bar{F}_1(t)\, d\bar{F}_0(t)$$

With $\bar{F}_1(t) = \int_t^1 f_1(u)\, du$ & $d\bar{F}_0(t) = -f_0(t)\, dt$,

we observe that

$$AUC = \int_0^1 \int_t^1 f_1(u)\, du\, f_0(t)\, dt = \iint f_0(u) f_1(t)\, du\, dt$$

$$= \mathbb{P}(f_n(X_1) \geq f_n(X_0) \mid \mathcal{L}_n, Y_0 = 0, Y_1 = 1),$$

where $(X_0, Y_0 = 0)$ and $(X_1, Y_1 = 1)$ are randomly chosen pairs.

In other words, the AUC corresponds to the probability that the score $f_n(X_1)$ is larger than $f_n(X_0)$, for a randomly generated sample $X_1 \mid Y_1 = 1$ from category $1$, and a randomly generated sample $X_0 \mid Y_0 = 0$ from category $0$.

In our context, $AUC(u) = \dfrac{1}{|I_u^+| |I \setminus I_u^+|} \sum_{i \in I_u^+} \sum_{j \in I \setminus I_u^+} \mathbb{1}(\hat{x}_{ui} \geq \hat{x}_{uj})$

Empirical Probability

The average AUC is $\quad AUC = \dfrac{1}{|U|} \displaystyle\sum_{u \in U} AUC(u)$ $\qquad$ (25)

$$= \sum_{(u,i,j) \in \mathcal{R}_S} \alpha_u \underbrace{\mathbb{1}\left(\hat{x}_{ui} \geqslant \hat{x}_{uj}\right)}_{\parallel};$$

with $\quad \alpha_u = \dfrac{1}{|U|\,|I_u^+|\,|I \smallsetminus I_u^+|}$ .

$$\mathbb{1}\left(\hat{x}_{uij} \geqslant 0\right)$$

Instead of an indicator function, BRP uses the differentiable loss $\log \sigma(x)$.