


Matrix Factorization



Matrix Factorization

Ratings Matrix ($N \times M$)

N = users M = items

- Factors:

$$10 = 2.5$$

$$15 = 3.5$$

Split matrix into product of 2 matrices

$$\hat{R} = W U^T \quad \hat{R} = \text{approx.}$$

$R_{N \times M}$ is sparse \Rightarrow stored as dict

$$N = 130k \quad M = 26k$$

$$N \times M = 3.38B$$

$$\# \text{ ratings} = 20M$$

W/U should be very skinny

$W_{N \times k}$ $U_{M \times k}$
users matrix movie matrix

$$k \sim 10-15$$

$$\begin{bmatrix} \hat{R} \end{bmatrix} = \begin{bmatrix} W \end{bmatrix} \begin{bmatrix} U^T \end{bmatrix}$$

$$\text{If } k=10, N=130k, M=26k$$

$$Nk + Mk = 1.56M \rightarrow \text{even less space than ratings}$$

One rating

$$\hat{r}_{ij} = w_i^T u_j$$

$$w_i = W[i, :]$$

$$\hat{r}_{ij} = \hat{R}[i, j]$$

$$u_j = U[:, j]$$

Interpretation

k elements in w_i and u_j

$k=5$ and action, comedy, romance, horror, animation

$w_i(1)$ = how much user i likes action

$w_i(2)$ = how much user i likes comedy

$u_j(1)$ = how much movie j contains action

$u_j(2)$ = how much movie j contains comedy

$$w_i^T u_j = \|w_i\| \|u_j\| \cos \theta \propto \text{sim}(i, j)$$

user i correlate with attributes of movie j

Example

$$w_i = \begin{pmatrix} 1 \\ 0.8 \\ -1 \\ 0.1 \\ 1 \end{pmatrix} \quad u_j = \begin{pmatrix} 1 \\ 1.5 \\ -1.3 \\ 0 \\ 1.2 \end{pmatrix} \quad w_i^T \cdot u_j = 4.7$$

+ve. +ve \rightarrow +ve
 -ve. -ve \rightarrow +ve
 \ominus means feature is not present or not liked

$$w_i = \begin{pmatrix} 1 \\ 0.8 \\ -1 \\ 0.1 \\ 1 \end{pmatrix} \quad u_j = \begin{pmatrix} -1 \\ -1 \\ 1 \\ 0 \\ -1 \end{pmatrix} = -3.8$$

- Features

each feature is latent, k is the latent dimensionality
 what features mean is not known before examining them (hidden cause)

Dimensionality reduction

if $k = M$ and $N > M$ and $\text{rank}(X) = M$

$U/S/V$ are approximations (shrunked)

Truncated SVD yields best rank k approx of X

MF reduces the dimensionality of R

Training

$$R \approx \hat{R} = w u^T$$

- Squared Error (for regression)

$$J = \sum_{i,j \in \Omega} (r_{ij} - \hat{r}_{ij})^2 = \sum_{i,j \in \Omega} (r_{ij} - w_i^T u_j)^2$$

Gradient:

$$\frac{\partial J}{\partial w_i} = 2 \sum_{j \in \Psi_i} (r_{ij} - w_i^T u_j) (-u_j) = 0 \rightarrow \text{min}$$

$$\sum_{j \in \Psi_i} (w_i^T u_j) u_j = \sum_{j \in \Psi_i} r_{ij} u_j$$

- Solving for w

$$w_i = \left(\sum_{j \in \Psi_i} u_j u_j^T \right)^{-1} \sum_{j \in \Psi_i} r_{ij} u_j$$

- Solving for u

$$\frac{\partial J}{\partial u_j} = 2 \sum_{i \in \Omega_j} (r_{ij} - w_i^T u_j) (-w_i) = 0 \rightarrow u_j = \left(\sum_{i \in \Omega_j} w_i w_i^T \right)^{-1} \sum_{i \in \Omega_j} r_{ij} w_i$$

- 2-way dependency
solution for w depends on u and reverse
start with random u/w and apply equations (alternating least squares)
- FAQ
loss is symmetric, so w loss = u loss

Expanding Model

- Regression
 $\hat{y} = mx + b$

- Bias terms

$$\hat{r}_{ij} = \underbrace{w_i^T u_j}_{\substack{\text{user vector} \\ \text{movie vector}}} + \underbrace{b_i}_{\text{user bias}} + \underbrace{c_j}_{\text{movie bias}} + \underbrace{\mu}_{\text{global avg.}}$$

predicted rating

μ = centering the dataset with mean 0
 b_i / c_j = Biases

- Movie Bias
great movie \rightarrow everyone likes it; attributes are not enough
Avatar vs SciFi shit

- Training

objective: $J = \sum_{i,j \in \Omega} (r_{ij} - \hat{r}_{ij})^2$

$$\hat{r}_{ij} = w_i^T u_j + b_i + c_j + \mu$$

$$\frac{\partial J}{\partial w_i} \stackrel{!}{=} 0 \leadsto w_i = \left(\sum_{j \in \Psi_i} u_j u_j^T \right)^{-1} \sum_{j \in \Psi_i} (r_{ij} - b_i - c_j - \mu) u_j$$

\uparrow symmetric

$$\frac{\partial J}{\partial u_j} \stackrel{!}{=} 0 \leadsto u_j = \left(\sum_{i \in \Omega_j} w_i w_i^T \right)^{-1} \sum_{i \in \Omega_j} (r_{ij} - b_i - c_j - \mu) w_i$$

$$\frac{\partial J}{\partial b_i} \stackrel{!}{=} 0 \leadsto b_i = \frac{1}{|\Psi_i|} \sum_{j \in \Psi_i} (r_{ij} - w_i^T u_j - c_j - \mu)$$

$$\frac{\partial J}{\partial c_j} \stackrel{!}{=} 0 \leadsto c_j = \frac{1}{|\Omega_j|} \sum_{i \in \Omega_j} (r_{ij} - w_i^T u_j - b_i - \mu)$$

Regularization

prevents overfitting

— For linear regression

Model: $\hat{y} = w^T x$

Objective: $J = \sum_{i=1}^N (y_i - \hat{y}_i)^2 + \lambda \|w\|_2^2$

square magnitude
large weights are penalised

Solution: $w = (\lambda I + X^T X)^{-1} X^T y$

$\|*\|_F$ = Frobenius norm

$$\|w\|_F^2 = \sum_{i=1}^N \sum_{k=1}^K |w_{ik}|^2 = \sum_{i=1}^N \|w_i\|_2^2 = \sum_{i=1}^N w_i^T w_i$$

$$J = \sum_{i,j \in \Omega} (r_{ij} - \hat{r}_{ij})^2 + \lambda (\|w\|_F^2 + \|u\|_F^2 + \|\theta\|_2^2 + \|c\|_2^2)$$

$$\frac{\partial J}{\partial w_i} \stackrel{!}{=} 0 \leadsto w_i = \left(\sum_{j \in \Psi_i} u_j u_j^T + \lambda I \right)^{-1} \sum_{j \in \Psi_i} (r_{ij} - \theta_i - c_j - \mu) u_j$$

$$u_j = \left(\sum_{i \in \Omega_j} w_i w_i^T + \lambda I \right)^{-1} \sum_{i \in \Omega_j} (r_{ij} - \theta_i - c_j - \mu) w_i$$

$$\theta_i = \frac{1}{|\Psi_i| (1+\lambda)} \sum_{j \in \Psi_i} (r_{ij} - w_i^T u_j - c_j - \mu)$$

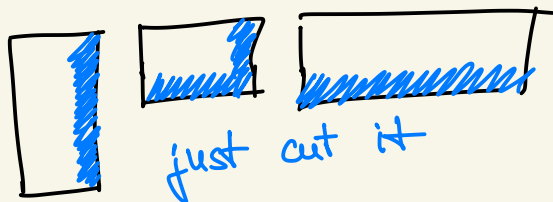
$$c_j = \frac{1}{|\Omega_j| (1+\lambda)} \sum_{i \in \Omega_j} (r_{ij} - w_i^T u_j - \theta_i - \mu)$$

SVD

X can be decomposed into U, S, V

$$X_{N \times M} = U_{N \times N} \cdot S_{N \times M} \cdot V_{M \times M}^T$$

truncated SVD $\rightarrow U, S, V$ ordered by importance



$$U_{N \times K} \cdot S_{K \times K} \cdot V_{M \times K}^T \approx X \quad (\text{approximation})$$

$$\hat{x}_{ij} = \sum_k u_{ik} s_{rk} v_{kj} = \sum_k (u_{ik} \cdot s_{rk}) v_{kj} = \sum_k u'_{ik} v_{kj} = u_i'^T v_j$$

Probabilistic Matrix Factorization

$$\underline{R} \sim N(\overbrace{w u^T}^{\text{mean}}, \sigma^2)$$

$$r_{ij} \sim N(w_i^T u_j, \sigma^2)$$

- Maximum Likelihood estimation

$$L = \prod_{i,j \in \Omega} \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left(-\frac{1}{2\sigma^2} (r_{ij} - w_i^T u_j)^2\right)$$

$$w, u = \arg\max_{w, u} L$$

$$\log(L) = C - \sum_{i,j \in \Omega} \frac{1}{2\sigma^2} (r_{ij} - w_i^T u_j)^2$$

- MAP Estimation

$$\text{MLE} : p(R | w, u)$$

$$\text{MLE} : p(w, u | R) \quad \alpha\text{-posteriori}$$

$$p(w, u | R) = \frac{p(R | w, u) p(w) p(u)}{p(R)}$$

$$p(w) = N(0, \lambda^{-1}) \quad p(u) = N(0, \lambda^{-1})$$

↓ normal, centered around zero
 ↗ precision parameter

- MAP objective

$$L = \prod_{i,j \in \Omega} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (r_{ij} - w_i^T u_j)^2\right) \times \frac{\lambda}{\sqrt{2\pi}} e^{-\frac{\lambda}{2} \|w\|_F^2} \cdot \frac{\lambda}{\sqrt{2\pi}} e^{-\frac{\lambda}{2} \|u\|_F^2}$$

Bayesian Matrix Factorization

everything is random variable

- Embeddings

$$\text{cat} \rightarrow \begin{pmatrix} 0.5 \\ -0.2 \\ 1.3 \end{pmatrix}$$

king-men ≈ queen-woman

N word vocabulary $W[i]$ $N \times K$ $K = \text{vector size}$
 $= \text{single word}$

- In Keras

2 Embedding Layers

$\text{emb}_u = \text{Embedding}(N, K)$

$\text{emb}_m = \text{Embedding}(M, K)$

2 inputs \rightarrow user + movie

$\text{dot}(\text{Embed}_{\text{user}}, \text{Embed}_{\text{movie}})$

- Bias Terms

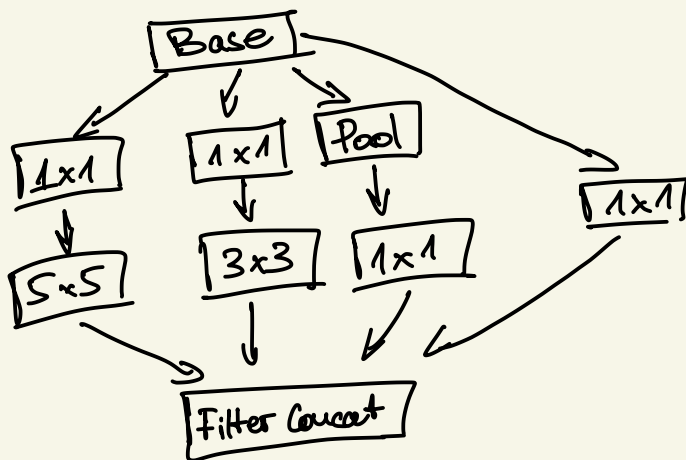
Keras is not low-level library
 (layer level lib)

$\text{Embedding}(N, 1) \rightarrow$ is scalar (user bias)

$(M, 1) \rightarrow$ movie bias

actual rating - μ (full dataset mean)

Residual Learning (Vision Inspired)



Different Branches, summed together
 at the end
 to detect different patterns

- MF performs well, but is linear
- NN can find nonlinear patterns

$\hat{r} = \text{MF}(x) + \text{ANN}(x)$ Residual (can be non-linear)

$\text{ANN}(x) = \hat{r} - \text{MF}(x) = \text{error made by MF model}$



NN will try to predict
 the error made by

MF / non linear pattern may be absent
 \rightarrow not learnable from the residuals

Autoencoders (Auto Dec)

- Feed forward NN that predicts its own input
Error = (output - input)²
general NN: ann. fit(x, y)
autoencoder: ann. fit(x, x)

- Denoising Autoencoders

reconstructing image with blackened parts

- Recommenders → fill missing ratings

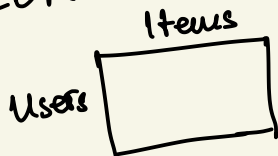
Structure of data in non-recommenders

$N \times D$ (# samples, # features)

N = images

D = pixel values in each image

- Ratings Matrix: 20M cells in the matrix are filled in



N = number of samples (users)

M = number of features (movies)

Complexity → loops over 130k users (not 20M individual ratings)

- New Addition
add noise to input (Dropout layer to the front of NN)
not simply reconstruct the matrix

- Sparse Matrix ($N \times M$ too big)

lil / csr / coo matrices

Keras does not recognise sparse matrix → convert to dense arrays

Missing values (0) have to be skipped

- Test MSE

model.predict(x_test) / y_test → normal

ignore keras val-loss, because it can copy input to output

train ratings should predict test ratings

predictions are entries in $N \times M$, where test values exist

$$J_{\text{test}} = \frac{1}{|\mathcal{I}_{\text{test}}|} \sum_{i=1}^N \sum_{j=1}^M m_{ij}^{(\text{test})} (r_{ij} - \hat{r}_{ij})^2$$

$m_{ij}^{(\text{test})} = 1$ if $(i, j) \in \mathcal{I}_{\text{test}}$ else 0

\hat{r}_{ij} = reconstruction from training data