

matrix factorization

Matrix Factorization takes into account both the similarities between other user's ratings and items.

$$\text{Carmen} \begin{bmatrix} \text{Alien} & \text{Bug's Life} & \text{Cars} & \text{Dark Knight} \\ 5 & - & 1 & 3 \end{bmatrix}$$

$$\text{Carmen}^1 \begin{bmatrix} \text{Scary} & \text{Kiddy} \\ 5 & 1 \end{bmatrix} \rightarrow \text{Dark Knight} \begin{bmatrix} \text{Scary} & \text{Kiddy} \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

Prediction from **Carmen's** rating of **Dark Knight**: $5 \times \frac{1}{2} + 1 \times \frac{1}{2} = 3 \rightarrow$ inner product of the vectors

$$\text{Carmen}^1 \begin{pmatrix} \text{Scary} & \text{Kiddy} \\ 5 & 1 \end{pmatrix} \begin{pmatrix} \text{Dark Knight} \\ 1/2 \\ 1/2 \end{pmatrix} \begin{pmatrix} \text{Scary} \\ \text{Kiddy} \end{pmatrix} = \text{Carmen}^1 \begin{pmatrix} \text{Dark Knight} \\ 3 \end{pmatrix}$$

The above logic implies predictions can be made for **Carmen's** rating of many different **movies** and equally many different **users**:

$$\text{Carmen} \begin{bmatrix} \text{Alien} & \text{Bug's Life} & \text{Cars} & \text{Dark Knight} & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 5 & - & 1 & 3 & 2 & - & 5 & 1 & 3 & 1 & 5 & - \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

Matrix Factorization determines the **latent factors**

Determining **latent space**

- .. How to represent each user in the space
- .. How to represent each movie in the space

If the latent factors are ideal, few are required to recover ratings for accurate approximations of real ratings, regardless of having thousands of different movies to rate. If a new movie exists that a user has not seen prior, all that matters is the latent state of the movie and the latent state of the user; if these two attributes are known, predications can be accurately assigned:

$$\text{Rating for user } i \text{ on movie } j \approx \hat{R}_{ij} = \sum_{\text{latent factors } \ell} \text{user}_{u,\ell} \text{movie}_{\ell,m}$$

Prediction from **Leonore's** rating of **Dark Knight**: $3 \times \frac{1}{2} + 2 \times \frac{1}{2} = 2.5 \rightarrow$ vectors inner product

$$\begin{matrix} \text{Carmen}^1 \\ \text{Leonore}^1 \end{matrix} \begin{pmatrix} \text{Scary} & \text{Kiddy} \\ 5 & 1 \\ 3 & 2 \end{pmatrix} \begin{pmatrix} \text{Dark Knight} \\ 1/2 \\ 1/2 \end{pmatrix} \begin{pmatrix} \text{Scary} \\ \text{Kiddy} \end{pmatrix} = \begin{matrix} \text{Carmen}^1 \\ \text{Leonore}^1 \end{matrix} \begin{pmatrix} \text{Dark Knight} \\ 3 \\ 2.5 \end{pmatrix}$$

Matrix Factorization applied to multiple movies and users:

Customer → Latent

Latent → Movies

$$\begin{array}{c}
 \text{Carmen}^1 \\
 \text{Leonore}^1 \\
 \text{Joeseeph}^1
 \end{array}
 \begin{pmatrix}
 \text{Scary} & \text{Kiddy} & & \text{Gore} \\
 5 & 1 & \dots & 4 \\
 3 & 2 & \dots & 1 \\
 \vdots & \vdots & \ddots & \vdots \\
 2 & 3 & \dots & 1 \\
 4 & 5 & \dots & 2
 \end{pmatrix}
 \begin{pmatrix}
 \text{Dark Knight} & \text{Jurassic Park} & & \text{Zootopia} \\
 1/2 & 2/3 & \dots & 1/4 & 1/9 \\
 1/2 & 1/3 & \dots & 1/3 & 1/3 \\
 \vdots & \vdots & \ddots & \vdots & \vdots \\
 1/8 & 1/4 & \dots & 1/5 & 1/7
 \end{pmatrix}
 \begin{array}{c}
 \text{Scary} \\
 \text{Kiddy} \\
 \\
 \text{Gore}
 \end{array}$$

$$= \begin{array}{c}
 \text{Carmen}^1 \\
 \text{Leonore}^1 \\
 \text{Joeseeph}^1
 \end{array}
 \begin{pmatrix}
 \text{Dark Knight} & \text{Jurassic Park} & & \text{Zootopia} \\
 3 & 3.7 & \dots & \dots \\
 2.5 & 2.7 & \dots & \dots \\
 \vdots & \vdots & \ddots & \vdots \\
 \dots & \dots & \dots & \dots
 \end{pmatrix}$$

The **first matrix** in the above illustration represents the **latent factors of users**. The **second matrix** represents the **latent factors of movies**. The **third matrix** represents the **inner product** of the latter:

Customer → Latent

Latent → Movies

Customer → Movies

$$\begin{pmatrix}
 5 & 1 & \dots & 4 \\
 3 & 2 & \dots & 1 \\
 \vdots & \vdots & \ddots & \vdots \\
 2 & 3 & \dots & 1 \\
 4 & 5 & \dots & 2
 \end{pmatrix}
 \begin{pmatrix}
 1/2 & 2/3 & \dots & 1/4 & 1/9 \\
 1/2 & 1/3 & \dots & 1/3 & 1/3 \\
 \vdots & \vdots & \ddots & \vdots & \vdots \\
 1/8 & 1/4 & \dots & 1/5 & 1/7
 \end{pmatrix}
 = \begin{pmatrix}
 3 & 3.7 & \dots & \dots \\
 2.5 & 2.7 & \dots & \dots \\
 \vdots & \vdots & \ddots & \vdots \\
 \dots & \dots & \dots & \dots
 \end{pmatrix}$$

The above model: $P \times Q^T = \hat{R}$ with the **goal** of choosing P and Q^T such that \hat{R} is accurate.

\hat{R} is accurate when: P and Q^T are chosen to **minimize**:

$$\sum_{\text{users } U, \text{ movies } m} (R_{u,m} - \hat{R}_{u,m})^2 = \sum_{\text{users } U, \text{ movies } m} (\text{error}_{u,m})^2$$

The estimated rating should ideally be as close to the true rating as possible. The squared distance between the true ratings and the estimated ratings are minimized (**Sum of Squares Error**).

The **optimization problem** above is solved through **gradient descent** on the **errors**. The method begins by minimizing P , then minimizing Q , then reiterating until convergence.

Alternate minimization scheme:

$$\text{error}_{u,m}^2 = (R_{u,m} - \hat{R}_{u,m})^2 = \left(r_{ij} - \sum_{\text{latent } \ell=1}^L \mathbf{p}_{u,\ell} \mathbf{q}_{\ell,m} \right)^2$$

$$\frac{d}{d\mathbf{p}_{u,\ell}} \text{error}_{u,m}^2 = -2(R_{u,m} - \hat{R}_{u,m}) \mathbf{q}_{\ell,m} = -2 \text{error}_{u,m} \mathbf{q}_{\ell,m}$$

$$\mathbf{p}_{u,\ell} := \mathbf{p}_{u,\ell} - \alpha \frac{d \text{error}_{u,m}^2}{d\mathbf{p}_{u,\ell}} = \mathbf{p}_{u,\ell} + 2\alpha \text{error}_{u,m} \mathbf{q}_{\ell,m}$$

$$\frac{d}{d\mathbf{q}_{u,\ell}} \text{error}_{u,m}^2 = -2(R_{u,m} - \hat{R}_{u,m}) \mathbf{p}_{\ell,m} = -2 \text{error}_{u,m} \mathbf{p}_{\ell,m}$$

$$\mathbf{q}_{u,\ell} := \mathbf{q}_{u,\ell} - \alpha \frac{d \text{error}_{u,m}^2}{d\mathbf{q}_{u,\ell}} = \mathbf{q}_{u,\ell} + 2\alpha \text{error}_{u,m} \mathbf{p}_{\ell,m}$$

If the learning rate α is set to be too **large**, gradient descent is likely to overshoot the minimum and oscillate back and forth. If the learning rate α is set to be too **small**, gradient descent will be slow to converge and might fail entirely. **Regularization** can also be added to the above illustration of **gradient descent**, adding an additional factor to the expression.

It is important to note that the above **minimization objective** does not necessarily provide the best possible solution at a **global minimum**. The objective is not **convex** and thus, **local minimizers** may result when implementing but the technique works well in practice.

Questioning the above method:

Is **Matrix Factorization** illustrated above synonymous to setting the **unknowns** to **0** and performing least squares regression?

$$\begin{array}{cc}
 \begin{array}{c} \text{Carmen}^1 \\ \text{Leonore}^1 \\ \text{Joeseeph}^1 \end{array} \begin{array}{c} \hat{R} \\ \left(\begin{array}{cccc} 3 & 3.7 & 2.7 & 2.9 \\ 2.5 & 2.7 & 4.1 & 1.3 \\ 2.1 & 0.9 & 0.5 & 2.8 \\ 4.6 & 2.6 & 4.2 & 0.1 \end{array} \right) \end{array} & \begin{array}{c} \text{Carmen}^1 \\ \text{Leonore}^1 \\ \text{Joeseeph}^1 \end{array} \begin{array}{c} R \\ \left(\begin{array}{cccc} 3 & 3.7 & ? & 2.9 \\ 2.5 & ? & 4.1 & ? \\ ? & 0.9 & 0.5 & 2.8 \\ ? & 2.6 & 4.2 & ? \end{array} \right) \end{array} \\
 \\
 \begin{array}{c} \text{Carmen}^1 \\ \text{Leonore}^1 \\ \text{Joeseeph}^1 \end{array} \begin{array}{c} \hat{R} \\ \left(\begin{array}{cccc} 3 & 3.7 & 2.7 & 2.9 \\ 2.5 & 2.7 & 4.1 & 1.3 \\ 2.1 & 0.9 & 0.5 & 2.8 \\ 4.6 & 2.6 & 4.2 & 0.1 \end{array} \right) \end{array} & \begin{array}{c} \text{Carmen}^1 \\ \text{Leonore}^1 \\ \text{Joeseeph}^1 \end{array} \begin{array}{c} R \\ \left(\begin{array}{cccc} 3 & 3.7 & 0 & 2.9 \\ 2.5 & 0 & 4.1 & ? \\ 0 & 0.9 & 0.5 & 2.8 \\ 0 & 2.6 & 4.2 & 0 \end{array} \right) \end{array}
 \end{array}$$

The answer is **no**; if the **unknowns** are replaced with **zero** in **matrix factorization**, any missing votes will suffer a large penalty which becomes problematic when many votes are missing. Only the errors that can be calculated are considered.