# Latent Factor Model I

STAT3009 Recommender Systems

by Ben Dai    (CUHK-STAT)
on August 13, 2022

* **Training dataset:** [**userID**, **itemID**, **rating**]
* **Testing dataset:** [**userID**, **itemID**, **?**]
* **Evaluation:** Given a testing index set $\Omega^{\text{te}}$ (set of user-item pairs we want to predict),

$$RMSE = \Big( \frac{1}{|\Omega^{\text{te}}|} \sum_{(u,i)\in\Omega^{\text{te}}} \big( \hat{r}_{ui} - r_{ui} \big)^2 \Big)^{1/2}.$$

* **Goal:** Find predicted ratings $(\hat{r}_{ui})_{(u,i)\in\Omega^{\text{te}}}$ such that **minimizes RMSE**
* **Baseline methods:** Global-average, user-average, item-average, user-item average
* **Correlation-based RS:** User-correlation-based RS, item-correlation-based RS, and correlation-based + baseline

Using ML methods to build RS:

Step 1. Introduce a model with some **parameters**

Step 2. Estimate the **parameters** by minimizing (maximizing) the **Evaluation Loss** in **Training Set**

Step 3. Use the estimated model to **predict**

Idea Learning from Data: A model works well in **Training Set**, tend to work well in **Testing Set**

⚠ MATH

* Training dataset $(\mathrm{feat}_i, \mathrm{label}_i)_{i=1}^n$
* Testing dataset $(\mathrm{feat}_j)_{j=1}^m$:

Step 1. Introduce a model with some **parameters**: $f_\theta$

Step 2. Estimate the **parameters** by minimizing (maximizing) the **Evaluation Loss** in **Training Set**

$$\widehat{f_\theta} = \operatorname*{argmin}_{f_\theta} \frac{1}{n} \sum_{i=1}^n L(\mathrm{label}_i, f_\theta(\mathrm{feat}_i)) + \lambda \operatorname{Reg}(f_\theta).$$

Step 3. Use the estimated model to **predict**

$$\widehat{\mathrm{label}}_j = \widehat{f_\theta}(\mathrm{feat}_j).$$

Data (feat, label) is a pair of input features and its outcome

Model $f_\theta$: a parameterized function to map features to label

Loss $L(\cdot, \cdot)$: The measure of how good the predicted outcome compared with the true outcome

Reg $\text{Reg}(f_\theta)$: regularization term in ERM to prevent overfitting

Opt The algorithm for solving the problem

Can we use this learning paradigm to find the best parameters for baseline models?

* **Step 1.** Introduce a method with some **parameters**

| method | math formula | **parameters** |
|---|---|---|
| Global Average | $\hat{r}_{ui} = \mu_0$ | $\mu_0$ |
| User Average | $\hat{r}_{ui} = a_u$ | $\boldsymbol{a} = (a_1, \cdots, a_n)^\intercal$ |
| Item Average | $\hat{r}_{ui} = b_i$ | $\boldsymbol{b} = (b_1, \cdots, b_m)^\intercal$ |
| User-Item Average | $\hat{r}_{ui} = \mu_0 + a_u + b_i$ | $\mu_0, \boldsymbol{a}, \boldsymbol{b}$ |

* **Step 2.** Estimate the **parameters** by minimizing **RMSE**
  * Global Average.

$$\min_{\mu_0 \in \mathbb{R}} \Big( \frac{1}{|\Omega|} \sum_{(u,i) \in \Omega} (\mu_0 - r_{ui})^2 \Big)^{1/2} \iff \min_{\mu_0 \in \mathbb{R}} \sum_{(u,i) \in \Omega} (\mu_0 - r_{ui})^2$$

* **Step 2.** Estimate the **parameters** by minimizing **RMSE**
  * Global Average.

$$\min_{\mu_0 \in \mathbb{R}} \Big( \frac{1}{|\Omega|} \sum_{(u,i) \in \Omega} (\mu_0 - r_{ui})^2 \Big)^{1/2} \iff \min_{\mu_0 \in \mathbb{R}} \sum_{(u,i) \in \Omega} (\mu_0 - r_{ui})^2$$

Taking the derivative to $\mu_0$:

$$2 \sum_{(u,i) \in \Omega} (\mu_0 - r_{ui}) = 0, \implies \mu_0 = \frac{1}{|\Omega|} \sum_{(u,i) \in \Omega} r_{ui} = \bar{r}$$

* **Step 2.** Estimate the **parameters** by minimizing **RMSE**
  * User Average.

$$\min_{\boldsymbol{a} \in \mathbb{R}^n} \sum_{(u,i) \in \Omega} (a_u - r_{ui})^2 \iff \min_{\boldsymbol{a} \in \mathbb{R}^n} \sum_{u=1}^{n} \sum_{i \in \mathcal{I}_u} (a_u - r_{ui})^2$$

  * The loss function is **separable**, thus it suffices to consider user-wise minimization:

$$a_u = \underset{a_u \in \mathbb{R}}{\operatorname{argmin}} \sum_{i \in \mathcal{I}_u} (a_u - r_{ui})^2$$
$$= \frac{1}{|\mathcal{I}_u|} \sum_{i \in \mathcal{I}_u} r_{ui} = \bar{r}_u, \quad \text{for } u = 1, \cdots, n$$

* **Step 2.** Estimate the **parameters** by minimizing **RMSE**
  * Item Average.

$$\min_{\boldsymbol{b} \in \mathbb{R}^m} \sum_{(u,i) \in \Omega} (b_i - r_{ui})^2 \iff \min_{\boldsymbol{b} \in \mathbb{R}^m} \sum_{i=1}^{m} \sum_{u \in \mathcal{U}_i} (b_i - r_{ui})^2$$

  * The loss function is **separable**, thus it suffices to consider item-wise minimization:

$$b_i = \operatorname*{argmin}_{b_i \in \mathbb{R}} \sum_{u \in \mathcal{U}_i} (b_i - r_{ui})^2$$
$$= \frac{1}{|\mathcal{U}_i|} \sum_{u \in \mathcal{U}_i} r_{ui} = \bar{r}_i, \quad \text{for } i = 1, \cdots, m$$

* **Step 2.** Estimate the **parameters** by minimizing **RMSE**
  * User-Item Average.

$$\min_{\boldsymbol{a} \in \mathbb{R}^n, \boldsymbol{b} \in \mathbb{R}^m} \sum_{(u,i) \in \Omega} (\mu_0 + a_u + b_i - r_{ui})^2$$

  * The loss is **non-separable**, taking the gradient w.r.t. $(\mu_0, \boldsymbol{a}, \boldsymbol{b})$ and eqaul to zeros

$$\sum_{(u,i) \in \Omega} (\mu_0 + a_u + b_i - r_{ui}) = 0$$

  * Specify $\mu_0$, $\boldsymbol{a}_u$, and $\boldsymbol{b}_i$ sequentially.

* **Step 3.** Use the **estimated model** to do prediction
  * Global Average: $\hat{r}_{ui} = \hat{\mu}_0$
  * User Average: $\hat{r}_{ui} = \hat{a}_u$
  * Item Average: $\hat{r}_{ui} = \hat{b}_i$
  * User-Item Average: $\hat{r}_{ui} = \hat{\mu}_0 + \hat{a}_u + \hat{b}_i$

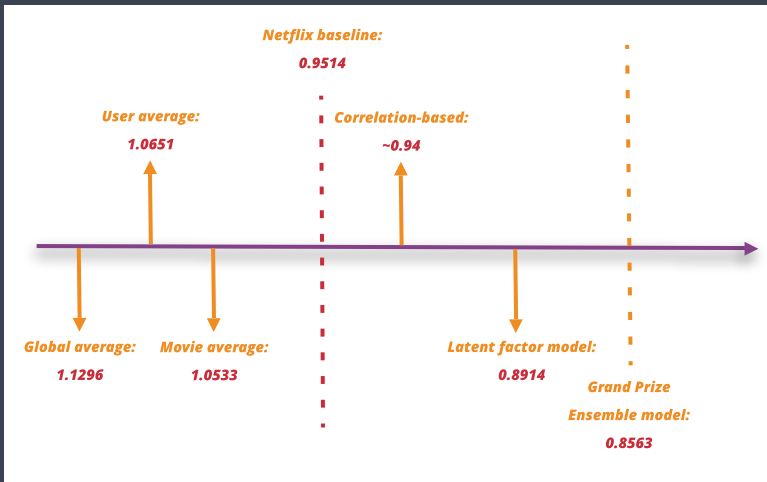*"All models are wrong, but some are useful." — George E. P. Box*

We need to figure out the **assumptions** for each method!

* **Global average** assumes that all users and items are essentially same
* **User average** assumes that a user has equal preference to all items
* **Item average** assumes that all users like "good" items
* **User-item average** assume that additive effects from users and items, **no interaction**

*Example:* I just don't like Action movies, even their ratings is quit high.

We need to model the user-item **interaction**.

Netflix baseline: 0.9514

User average: 1.0651

Correlation-based: ~0.94

Global average: 1.1296

Movie average: 1.0533

Latent factor model: 0.8914

Grand Prize

Ensemble model: 0.8563

* **Step 1.** Introduce a method with some **parameters** (latent factors)
    * Introduce $K$-length latent factors $\boldsymbol{p}_u$ for the user $u$:

    $$\boldsymbol{p}_u = (p_{u1}, \cdots, p_{uK})^\mathsf{T}$$

    * Introduce $K$-length latent factors $\boldsymbol{q}_i$ for the item $i$:

    $$\boldsymbol{q}_i = (q_{i1}, \cdots, q_{iK})^\mathsf{T}$$

    * Model the user-item **interaction** as **inner production**

    $$\hat{r}_{ui} = \boldsymbol{p}_u^\mathsf{T} \boldsymbol{q}_i \to r_{ui}$$

* $K$ is pre-specified number: **#Latent Factors**
* We tend to learn the **latent factors** from the data

* **Step 2.** Estimate the **parameters** by minimizing **RMSE**

$$\min_{\boldsymbol{P} \in \mathbb{R}^{n \times K}, \boldsymbol{Q} \in \mathbb{R}^{m \times K}} \frac{1}{|\Omega|} \sum_{(u,i) \in \Omega} (r_{ui} - \boldsymbol{p}_u^{\mathsf{T}} \boldsymbol{q}_i)^2 \tag{1}$$

* $K$ is a pre-specified #Latent Factor, can **NOT** be solved by (3). In ML, we call it **tuning parameter** or **hyperparameter**.
* $K$ increases $\implies$ more **parameters** $\implies$ lower **training loss**
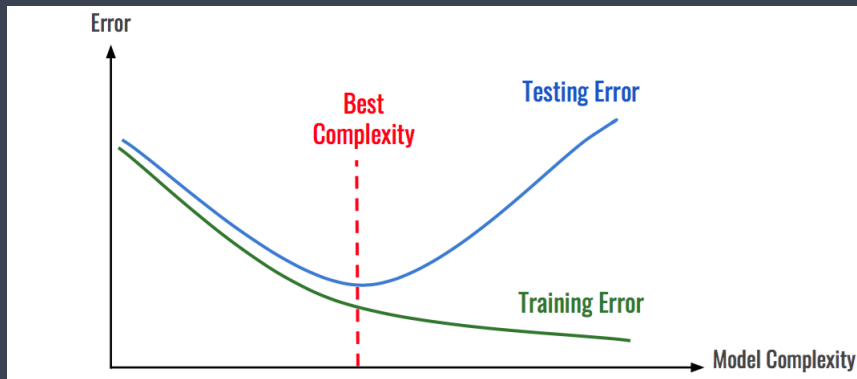* Lower training loss is always better? **NO!**

Source[1]

* **Overfitting**: fit the noise
* Too many **parameters** (**model complexity**) leads to **overfitting**

---

[1]https://hackernoon.com/memorizing-is-not-learning-6-tricks-to-prevent-overfitting-in-machine-learning-820b091dc42

Source[2]

* **Complexity** too large $\implies$ Low Training loss but high Testing loss

[2]https://hackernoon.com/memorizing-is-not-learning-6-tricks-to-prevent-overfitting-in-machine-learning-820b091dc42

|  | Low *Training* Error | High *Training* Error |
|---|---|---|
| **Low *Testing* Error** | The model is learning! | Probably some error in your code. Or you've created a *psychic* AI. |
| **High *Testing* Error** | O V E R F I T T I N G | The model is not learning. |

Source[3]

---
[3] https://hackernoon.com/memorizing-is-not-learning-6-tricks-to-prevent-overfitting-in-machine-learning-820b091dc42

* **Q1:** How to **quantify** the Model **Complexity**:
    * **#Parameters**: $(n + m)K$
    * **Magnitude** of Parameters: $\sum_{u=1}^{n} \|\boldsymbol{p}_u\|_2^2, \sum_{i=1}^{m} \|\boldsymbol{q}_i\|_2^2$
* **A1:** Control (#Parameters by $K$, **Magnitude** by $l_2$-norm).
* **Regularized** Latent Factor Model:

$$\min_{\boldsymbol{P}, \boldsymbol{Q}} \underbrace{\frac{1}{|\Omega|} \sum_{(u,i) \in \Omega} (r_{ui} - \boldsymbol{p}_u^\top \boldsymbol{q}_i)^2}_{\text{Training loss}} + \lambda \big( \underbrace{\sum_{u=1}^{n} \|\boldsymbol{p}_u\|_2^2 + \sum_{i=1}^{m} \|\boldsymbol{q}_i\|_2^2}_{\text{Params magnitude}} \big) \quad (2)$$

where $K$ and $\lambda > 0$ are **tuning** parameters to **balance** the model complexity and training loss.
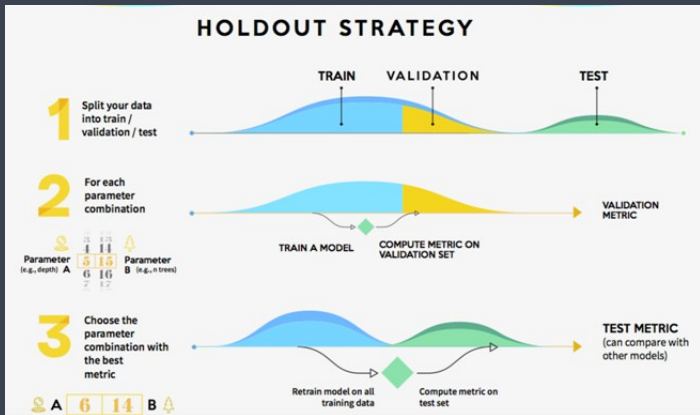* Why the later term can control the magnitude?

[18/23]

* **Q2:** How to find the **best** tuning parameters ($K, \lambda$)
* **A2: Cross-validation** (CV) based on
  **Training/Validation** splitting

* **Q2:** How to find the **best** tuning parameters $(K, \lambda)$
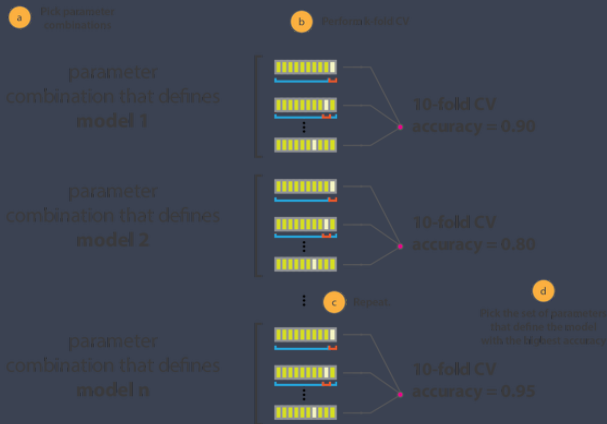* **A2: Cross-validation** (CV) based on
  **Training/Validation** splitting



Source[4]

∗ Typical splitting method: $k$-**fold** CV



Source[5]

---

[5] https://cambridgecoding.wordpress.com/2016/04/03/scanning-hyperspace-how-to-tune-machine-learning-models/

* **Step 2.** Estimate the **parameters** by minimizing **RMSE**

$$\min_{\boldsymbol{P},\boldsymbol{Q}} \frac{1}{|\Omega|} \sum_{(u,i)\in\Omega} (r_{ui} - \boldsymbol{p}_u^{\mathsf{T}}\boldsymbol{q}_i)^2 + \lambda \, \big( \sum_{u=1}^{n} \|\boldsymbol{p}_u\|_2^2 + \sum_{i=1}^{m} \|\boldsymbol{q}_i\|_2^2 \big)$$

* Using **Cross-Validation** to determine the optimal **tuning parameters** $(K, \lambda)$, denote as $K^*$ and $\lambda^*$
* **Refit** the model based on **full training data** with $K^*$ and $\lambda^*$.
* The final estimator is denoted as $(\hat{\boldsymbol{p}}_u)_{u=1}^n$ and $(\hat{\boldsymbol{q}}_i)_{i=1}^m$

* **Step 3.** Using the **estimated model** with the best tuning parameters to do prediction

$$\hat{r}_{ui} = \hat{\boldsymbol{p}}_u^{\mathsf{T}} \hat{\boldsymbol{q}}_i, \text{ for } (u, i) \in \Omega^{\text{te}}$$

* **Step 1.** Introduce a method with some **parameters** (latent factors)
  * Model the user-item **interaction** as **inner production**

  $$\hat{r}_{ui} = \boldsymbol{p}_u^\mathsf{T} \boldsymbol{q}_i \rightarrow r_{ui}$$

* **Step 2.** Estimate the **parameters** by minimizing **RMSE**

$$\min_{\boldsymbol{P}, \boldsymbol{Q}} \frac{1}{|\Omega|} \sum_{(u,i) \in \Omega} (r_{ui} - \boldsymbol{p}_u^\mathsf{T} \boldsymbol{q}_i)^2 + \lambda \left( \sum_{u=1}^{n} \|\boldsymbol{p}_u\|_2^2 + \sum_{i=1}^{m} \|\boldsymbol{q}_i\|_2^2 \right) \quad (3)$$

  * Using **Cross-Validation** to determine the optimal **tuning parameters** $(K, \lambda)$, denote as $K^*$ and $\lambda^*$
  * **Refit** the model based on **full training data** with $K^*$ and $\lambda^*$.

* **Step 3.** Using the **estimated model** with the best tuning parameters to do prediction