Master's Thesis Defense

# MetaRec
# Meta–Learning Meets Recommendation Systems

Presented by **James Le**
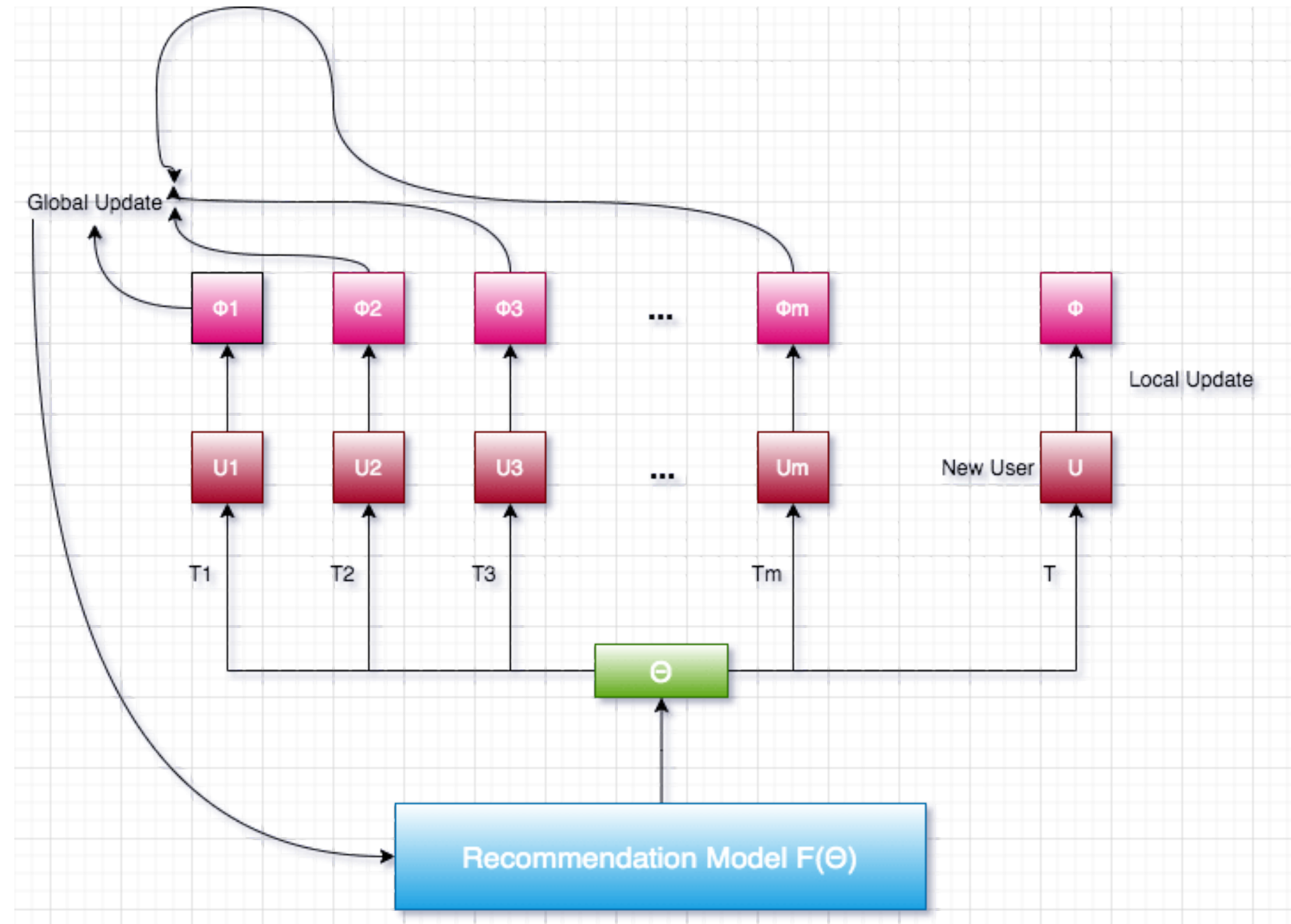Advised by **Alexander Ororbia II**
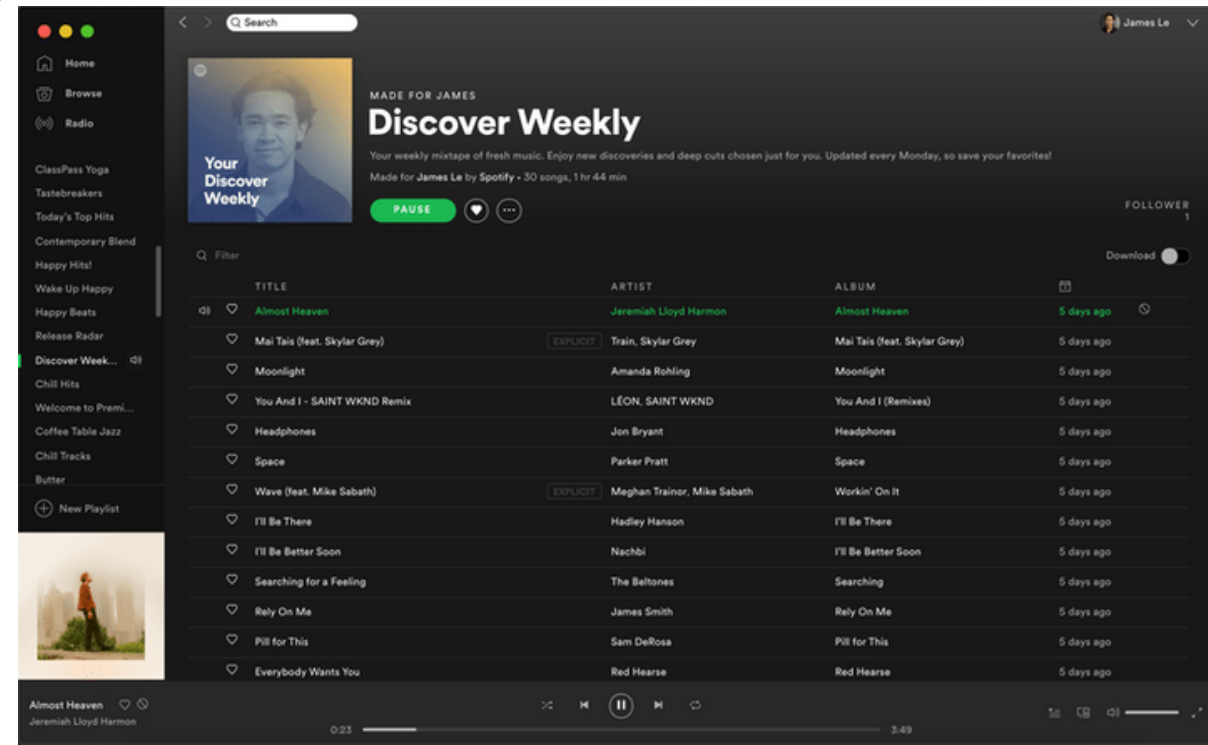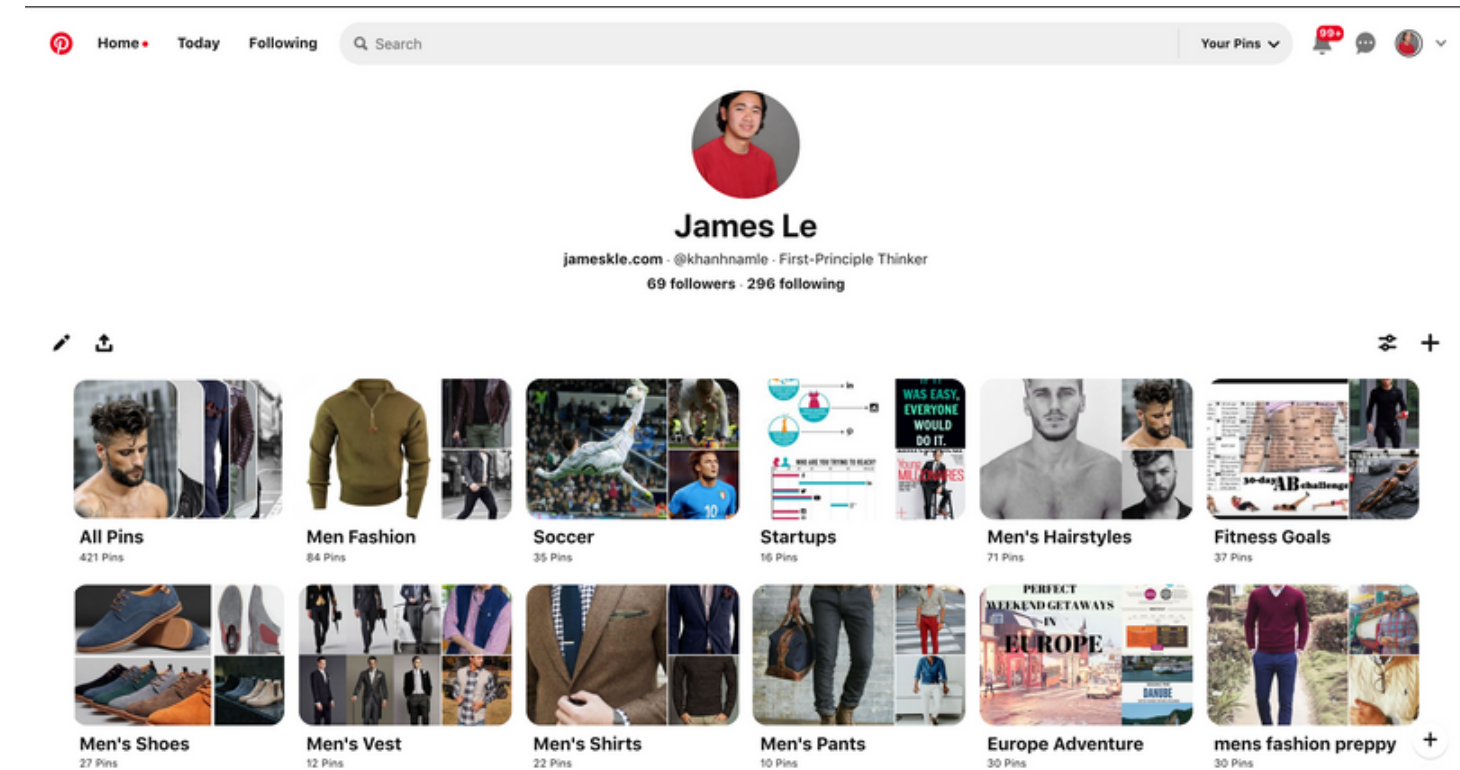As Member of **Neural Adaptive Computing Lab**

# Outline

# Introduction

## PART 1

ENTERTAINMENT



CONTENT



E-COMMERCE



SERVICES

# Challenges

## WHAT I WANT TO TACKLE

**Scalability**

**Sparsity**

**Accuracy**

# Challenges

## WHAT I WANT TO TACKLE

**Scalability**

**Sparsity**

**Accuracy**

"RECOMMENDATION SYSTEMS THAT META-LEARN INFORMATION FROM RELATED TASKS OUTPERFORM SYSTEMS THAT DO NOT USE THIS INFORMATION ACROSS A WIDE RANGE OF TASKS IN ACCURACY METRICS."

# Recommendation Systems

## PART 2

BY JAMES LE

# Problem Formulation

$$\forall u \in U, i_s = \arg\max_{i \in I} F(u, i)$$

- **U** - set of all users
- **I** - set of all items
- **F(u, i)** - utility function that measures the relevance of item i to user u
- **F** is usually represented as a rating in numeric scale

# Collaborative Filtering

# Content Filtering



Liked By

User A

Item A

Similar Customers

Recommendation System

User B

Recommend Items Based on Similar Customers

Liked By

Item A

User A

Recommendation System

Similar Items

Recommend Items Based on Similar Items

Item B

## Collaborative Filtering

Liked By

Item A

User A

Similar Customers

Recommendation System

Recommend Items Based on Similar Customers

User B

## Content Filtering

Liked By

Item A

User A

Recommendation System

Similar Items

Recommend Items Based on Similar Items

Item B

# Matrix Factorization



**MATRIX FACTORIZATION**

$$u_i^T v_j = \hat{r}_{ij}$$

- The de-facto standard for collaborative filtering
- Represent user preferences as an user-item matrix
- Find two smaller matrices (latent embeddings) that approximate the full matrix
- Minimize the rating prediction errors
- Can be optimized via Gradient Descent
- Prediction is simple

# Neural Networks



**MULTI-LAYER PERCEPTRON**



**AUTOENCODER**

# Neural Networks For Recommendations

## SOLUTIONS FOR EXISTING CHALLENGES

**SCALABILITY** → **SPARSITY** → **ACCURACY**

- Extract low-dimensional factors of high-dimensional user preferences for the items

- Engineer content-based features and integrate them into the algorithm
- Extract high-level representations of user preferences for the items

- Extract user and item's latent factors
- Jointly combine information from different data sources and modalities

# Meta Learning

## PART 3

# Motivation

**Long-Tail**



**Limited Data**

**Fast Inference**

# Model–Agnostic Meta–Learning (MAML, Finn et. al, 2017)

$$\min_{\theta} \sum_{\text{tasks } i} \mathcal{L}_{\text{val}}^{(i)}(\theta - \alpha \nabla_\theta \mathcal{L}_{\text{train}}^{(i)}(\theta))$$

$\theta$   parameter vector being meta-learned

$\theta_i^*$   optimal parameter vector for task i

— meta-learning
---- learning/adaptation

$\theta$

$\nabla \mathcal{L}_3$

$\nabla \mathcal{L}_2$

$\nabla \mathcal{L}_1$

$\theta_3^*$

$\theta_1^*$

$\theta_2^*$

# MetaRec

## PART 4

Binary Classification: $P(r_{ui} = 1 | p_u, q_i) = F_\theta(p_u, q_i, \theta)$

Rating Regression: $\hat{r}_{ui} = F_\theta(p_u, q_i, \theta)$

$$\phi_u = \theta - \alpha \nabla_\theta L_{T_u}(F_\theta)$$

$$\theta = \theta - \beta \nabla_\theta \sum_{T_u \sim B} L_{T_u}(F_{\phi_u})$$

## RECOMMENDATION TASK

- Task T = Recommend items to one user given examples of items already rated by that user
- Binary Classification (Implicit Feedback)
- Rating Regression (Explicit Feedback)
- θ denotes model parameters

## LOCAL UPDATE

- Sample a batch of tasks B
- Uses MAML to compute task-specific parameters φ based on the loss function value for each task T in batch B
- α is the local learning rate

## GLOBAL UPDATE

- Uses MAML to compute model parameters θ on the sum of test losses for all tasks T
- β is the global learning rate

# Experiments

## PART 5

# MovieLens1M

## (HARPER AND KONSTAN, 2015)

| | title | genres | rating |
|---|---|---|---|
| 0 | Toy Story (1995) | Animation\|Children's\|Comedy | 5 |
| 489283 | American Beauty (1999) | Comedy\|Drama | 5 |
| 489259 | Election (1999) | Comedy | 5 |
| 489257 | Matrix, The (1999) | Action\|Sci-Fi\|Thriller | 5 |
| 489256 | Dead Ringers (1988) | Drama\|Thriller | 5 |
| 489237 | Rushmore (1998) | Comedy | 5 |
| 489236 | Simple Plan, A (1998) | Crime\|Thriller | 5 |
| 489226 | Hands on a Hard Body (1996) | Documentary | 5 |
| 489224 | Pleasantville (1998) | Comedy | 5 |
| 489212 | Say Anything... (1989) | Comedy\|Drama\|Romance | 5 |
| 489207 | Beetlejuice (1988) | Comedy\|Fantasy | 5 |
| 489190 | Roger & Me (1989) | Comedy\|Documentary | 5 |
| 489172 | Buffalo 66 (1998) | Action\|Comedy\|Drama | 5 |
| 489171 | Out of Sight (1998) | Action\|Crime\|Romance | 5 |
| 489170 | I Went Down (1997) | Action\|Comedy\|Crime | 5 |
| 489168 | Opposite of Sex, The (1998) | Comedy\|Drama | 5 |
| 489157 | Good Will Hunting (1997) | Drama | 5 |
| 489152 | Fast, Cheap & Out of Control (1997) | Documentary | 5 |
| 489149 | L.A. Confidential (1997) | Crime\|Film-Noir\|Mystery\|Thriller | 5 |
| 489145 | Contact (1997) | Drama\|Sci-Fi | 5 |

| Attributes | Users | Movies | Ratings |
|---|---|---|---|
| Total | 6,040 | 3,883 | 1,000,209 |
| Min Rating | 20 | 1 | 1 |
| Max Rating | 2,314 | 3,428 | 5 |
| Average Rating | 165.6 | 269.9 | 3.58 |

# Matrix Factorization Experiments

## BASELINE MODELS

Vanilla Matrix Factorization
(Koren et al., 2009)

$$r_{ui} = p_u \cdot q_i^T$$

Matrix Factorization with Side Features
(Koren et al., 2009)

$$r_{ui} = b + d_o + \omega_u + \omega_i + (p_u + t_o) \cdot q_i^T$$

Matrix Factorization with Biases
(Koren et al., 2009)

$$r_{ui} = b + \omega_u + \omega_i + p_u \cdot q_i^T$$

Matrix Factorization with Temporal Features
(Koren et al., 2009)

$$r_{ui} = b + \omega_u + \omega_i + p_u \cdot q_i^T + m_u \cdot n_t$$

# Matrix Factorization Experiments

## BASELINE MODELS

Factorization Machine
(Rendle, 2010)

$$r_{ui} = b + \sum_{i=1}^{n} \omega_i \cdot x_i + \sum_{i=1}^{n} \sum_{j=i+1}^{n} <v_i, v_j> x_i \cdot x_j \qquad \mathbf{w} \in \mathbb{R}^n, \quad \mathbf{V} \in \mathbb{R}^{n \times k}$$

Matrix Factorization with Mixture of Taste
(Kula, 2018)

$$r_{ui} = \sigma(A_u \cdot q_i^T) \cdot (U_u \cdot q_i^T) + b + \omega_u + \omega_i$$

$$\sigma(x_i) = \frac{\exp^{x_i}}{\sum_j \exp^{x_j}} \qquad A \in \mathbb{R}^{m \times k}, U \in \mathbb{R}^{m \times k}$$

Variational Matrix Factorization
(Porteus et al., 2010, Kim et al., 2014)

$$r_{ui} = b + \omega_u + \omega_i + GS(\mu_u, v_u) \cdot GS(\mu_i, v_i)$$

$$GS(\mu, v) = \mu + \mathcal{N}(0, 1) \cdot \sqrt{v}$$

# MetaRec–MF

## BASE MODEL

Matrix Factorization with Biases

$$\hat{r}_{ui} = F(u, i|p_u, q_i, \theta) = p_u \cdot q_i + b + w_u + w_i$$

## LOSS FUNCTION

Mean-Squared Error with L2 Regularization

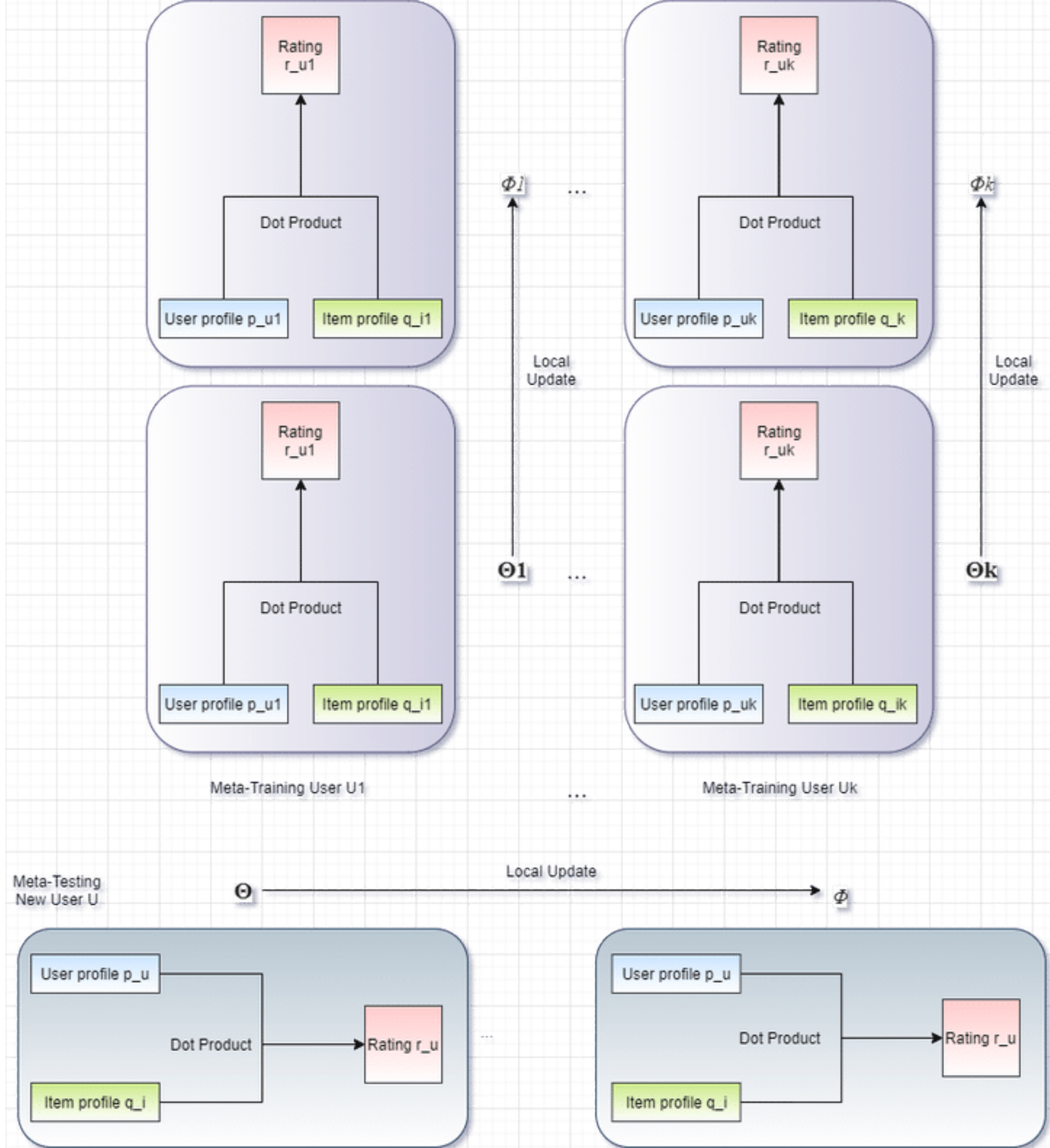$$L(\theta) = \frac{1}{|D_{train}|} \sum_{r_{u,i} \in D_{train}} (r_{u,i} - \hat{r}_{u,i})^2 + \lambda\frac{1}{2}||\theta||_2^2$$

## EVALUATION METRICS

- Rating Regression with Explicit Feedback
- Mean Squared Error
- Mean Absolute Error

$$MAE = \frac{1}{|D_{test}|} \sum_{r_{u,i} \in D_{test}} |r_{u,i} - \hat{r}_{u,i}|$$

$$MSE = \frac{1}{2|D_{test}|} \sum_{r_{u,i} \in D_{test}} (r_{u,i} - \hat{r}_{u,i})^2$$

# Regression Results

| Models | MAE | MSE | Training Time |
|---|---|---|---|
| MF | 0.703 | 0.817 | **6m5s** |
| MF-Bias | 0.694 | 0.790 | 11m38s |
| MF-Side | 0.691 | **0.784** | 13m34s |
| MF-Temporal | 0.689 | 0.793 | 18m51s |
| FM | 0.712 | 0.823 | **3m40s** |
| MF-Mixture | **0.687** | 0.788 | 13m44s |
| Variational-MF | 0.707 | 0.839 | 16m51s |
| MetaRec-MF (Ours) | **0.687** | **0.760** | 12m45s |

- Train-Test Split: 75-25
- Trained for 50 Epochs
- MetaRec-MF outperforms the other models on both MAE and MSE metrics (Lower values are better)

- Slight tradeoff between accuracy performance and compute cost
- The Factorization Machine model took the fastest time to train
- Adding more features to the matrix factorization equation led to longer training time

# Multi–Layer Perceptron Experiments

## BASELINE MODELS

Wide and Deep (Cheng et al., 2016)



Deep Factorization Machine (Guo et al., 2017)



Extreme Deep Factorization Machine (Lian et al., 2018)



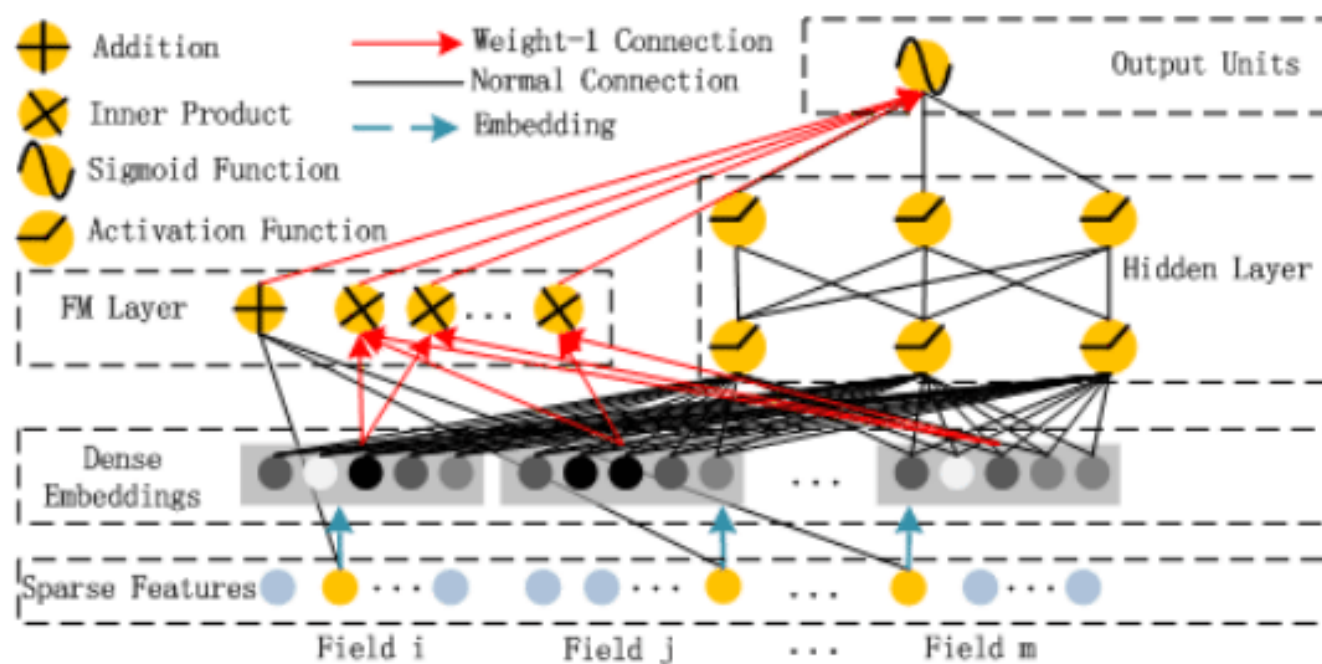BY JAMES LE

# Multi–Layer Perceptron Experiments

## BASELINE MODELS

Neural Factorization Machines (He et al., 2017)

Neural Collaborative Filtering (He et al., 2017)

# MetaRec–MLP

## BASE MODEL

Neural Collaborative Filtering

## LOSS FUNCTION

Binary Cross-Entropy Loss

## EVALUATION METRICS

- Binary Classification with Implicit Feedback
- Area Under the ROC Curve (AUC)

$$e_u = F_{\theta_u}(p_u)$$
$$e_i = F_{\theta_i}(q_i)$$

$$x_0 = [e_u, e_i],$$
$$x_1 = \sigma(W_1 x_0 + b_1),$$
$$\dots$$
$$x_m = \sigma(W_m x_{m-1} + b_m),$$
$$\hat{r}_{ui} = \sigma(W_o x_m + b_o)$$

$$L_u(F_{\theta, \phi_u}) = -\frac{1}{|B|} \sum_{i \in B} (r_{ui} log \hat{r}_{ui}) + (1 - r_{ui}) log(1 - \hat{r}_{ui})$$

# ADAPTED FROM (HE ET AL., 2017)

# Classification Results

| Models | Test AUC | Validation AUC | Training Time |
|---|---|---|---|
| WideDeep | **0.7991** | **0.7995** | 1h12m15s |
| DeepFM | 0.7915 | 0.7918 | 1h10m50s |
| xDeepFM | 0.7429 | 0.7408 | 2h15m17s |
| NeuralFM | 0.7589 | 0.7560 | 1h36m0s |
| NeuralCF | 0.7668 | 0.7673 | **54m15s** |
| MetaRec-MLP (Ours) | **0.8135** | **0.8127** | **1h05m3s** |

- Train-Valid-Test Split: 80-10-10
- Trained for 100 Epochs
- MetaRec-MLP outperforms other methods in AUC performance on both test and validation sets (Higher values are better)

- Tradeoff between model performance and compute cost
- WideDeep and DeepFM perform reasonably well
- xDeepFM overfits

# Autoencoder Experiments

## BASELINE MODELS

Collaborative Denoising Autoencoder
(Wu et al., 2016)

Multinomial Variational Autoencoder
(Liang et al., 2018)

# Autoencoder Experiments

## BASELINE MODELS

Sequential Variational Autoencoder
(Sachdeva et al., 2019)

Embarrassingly Shallow Autoencoder
(Steck, 2019)



$$\hat{x}_{(t+1)} \sim P_\theta\left(x | z_{(t)}\right)$$

$$z_{(t)} \sim Q_\phi\left(z | x_{(1:t)}\right)$$

$$Q_\phi\left(z | x_{(1:t)}\right) = \mathcal{N}\left(\mu_{(t)}, \sigma_{(t)}\right)$$

# MetaRec–AE

## BASE MODEL

Vanilla Autoencoder

$$\hat{r}_{ui} = (recon(r_{ui}, \hat{\theta}))$$

$$recon(r_{ui}; \theta) = f(W \cdot g(V r_{ui} + \mu) + b)$$

## LOSS FUNCTION

Binary Cross-Entropy Loss with L2 Regularization

$$L_u(\theta) = -\frac{1}{m} \sum_{u=1}^{m} ||(r_{ui} log \hat{r}_{ui}) + (1 - r_{ui}) log(1 - \hat{r}_{ui})||_O^2 + \frac{\lambda}{2} \cdot (||W||_2^2 + ||V||_2^2)$$

## EVALUATION METRICS

- Top-K Recommendation with Implicit Feedback
- Precision
- Recall
- Normalized Discounted Cumulative Gain

$$Precision@k = \frac{Hits@k}{k}$$

$$Recall@k = \frac{Hits@k}{|R|}$$

$$NDCG@k = \frac{DCGs@k}{IDCG@k}$$

Meta-Training

# Ranking Results

| Models | Precision@100 | Recall@100 | NDCG@100 | Training Time |
|---|---|---|---|---|
| CDAE | **8.94** | 41.37 | 25.28 | 17m29s |
| MultVAE | **8.86** | 41.15 | 25.08 | **6m31s** |
| SVAE | 8.18 | **58.49** | **38.07** | 6h37m19s |
| ESAE | 7.57 | 41.81 | 25.61 | **10m12s** |
| MetaRec-AE (Ours) | 8.34 | **55.12** | **33.06** | 18m24s |

- Train-Valid-Test Split: 80-10-10
- Trained for 100 epochs with early stopping
- MetaRec-AE had an average performance across three metrics (Higher values are better)

- SVAE was the best performing model in Recall and NDCG metrics
- MultVAE took the shortest time to train
- ESAE performed reasonably well

# Conclusion

## PART 6

# Problem Formulation

## RECOMMENDATION TASK

- Applications
- Challenges
- Collaborative Filtering
- Linear Models
- Neural Network Models

# Proposed Framework

## METAREC

- Model-Agnostic-Meta-Learning
- Leverages Learning From Previous Users To New Users
- Accurate: High Evaluation Metrics
- Rapid: In Small Number of Steps
- Efficient: Using Few Examples

# Experimental Results

## VS STATE-OF-THE-ARTS

- Matrix Factorization Models
- Multi-Layer Perceptron Models
- Auto-Encoder Models
- Explicit and Implicit Feedback
- Improved Accuracy Performance
- Computational Cost Tradeoff

# Areas of Improvement

## NOTES FOR FUTURE RESEARCH

**Make Changes To The Base Model F($\Theta$)**
- Complexify Them
- Simplify Them

**Better MAML Training**
- Training Instability
- Second-Order Derivative Cost
- Fixed Learning Rates During Updates

# Areas of Improvement

## NOTES FOR FUTURE RESEARCH

**Other Meta-Learning Schemes For Recommendations**
- Black-Box Meta-Learning
- Non-Parametric Meta-Learning

**Addressing Scalability and Sparsity Challenges**
- Shorter Training Time ~ Better Scalability
- Warm Users vs Cold Users Setup

# References

- Chelsea Finn, Pieter Abbeel, and Sergey Levine. **Model-Agnostic Meta-Learning For Fast Adaptation Of Deep Networks**. In *International Conference on Machine Learning*, 2017.
- Maxwell Harper and Joseph Konstan. **The MovieLens Datasets: History and Context.** In *ACM Transactions on Interactive Intelligent Systems*, 2015.
- Yehuda Koren, Robert Bell, and Chris Volinsky. **Matrix Factorization Techniques For Recommender Systems.** In *IEEE Computer Society*, 2009.
- Steffen Rendle. **Factorization Machines.** In *IEEE Conference on Data Mining*, 2010.
- Maciej Kula. **Mixture-Of-Tastes Models For Representing Users With Diverse Interests.** 2017.
- Ian Porteous, Arthur Asuncion, and Max Welling. **Bayesian Matrix Factorization With Side Information and Dirichlet Process Mixtures.** 2010.
- Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, et al. **Wide and Deep Learning For Recommendation Systems.** In *RecSys*, 2016.
- Huifeng Guo, Ruiming Tang, Yunming Ye, et al. **DeepFM: A Factorization-Machine Based Neural Network For CTR Prediction.** 2017.
- Jianxun Lian, Xiaohuan Zhou, Fuzheng Zhang, et al. **xDeepFM: Combining Explicit and Implicit Feature Interactions For Recommender Systems.** 2018.
- Xiangnan He and Tat-Seng Chua. **Neural Factorization Machines For Sparse Predictive Analytics.** 2017.
- Xiangnan He, Lizi Liao, Hangwang Zhang, et al. **Neural Collaborative Filtering.** In *The World Wide Web Conference,* 2017.
- Yao Wu, Christopher DuBois, Alice Zheng, and Martin Ester. **Collaborative Denoising Autoencoders For Top-N Recommender Systems.** In *Web Search and Data Mining,* 2016.
- Dawen Liang, Rahul Krishman, Matthew Hofman, and Tony Jebara. **Variational Autoencoders For Collaborative Filtering.** 2018.
- Noveen Sachdeva, Giuseppe Manco, Ettore Ritacco, and Vikram Pudi. **Sequential Variational Autoencoders For Collaborative Filtering.** In *Web Search and Data Mining,* 2019.
- Harld Steck. **Embarrassingly Shallow Autoencoders For Sparse Data.** In *The World Wide Web Conference,* 2019.

# Thank You!

## jl1165@rit.edu 😊

**Codebase: github.com/khanhnamle1994/MetaRec**

# Appendix

## Extra Details

# MetaRec–MF

---

**Algorithm 1:** MAML Algorithm for MetaRec-MF

---

**Input** : Trainable global parameters $\theta$, user set $U$, local and global

hyper-parameters $\alpha, \beta$

**Output:** Learned global parameters $\theta$

1 Initialize $\theta$ randomly;

2 **while** *not converged* **do**

3      Sample a batch of users $B$ from $U$;

4      **for** *user $u \in B$* **do**

5          Evaluate the gradient $\nabla_\theta L_u(F_\theta)$;

6          Perform the local update $\phi_u \leftarrow \theta - \alpha \nabla_\theta L_u(F_\theta)$;

7      **end**

8      Perform global update $\theta \leftarrow \theta - \beta \sum_u L_u(F_{\phi_u})$

9 **end**

---

# MetaRec–MLP

---

**Algorithm 2:** MAML Algorithm for MetaRec-MLP

---

    **Input**   : Trainable global parameters $\theta$, user set $U$, local and global hyper-parameters $\alpha, \beta$

    **Output:** Learned global parameters $\theta$

1  Initialize $\theta$ randomly;

2  Initialize $\phi$ randomly;

3  **while** *not converged* **do**

4      Sample a batch of users $B$ from $U$;

5      **for** *user $u \in B$* **do**

6          Evaluate the gradient $\nabla_{\phi_u} L_u(F_{\theta,\phi_u})$;

7          Perform the local update $\phi_u \leftarrow \phi_u - \alpha \nabla_{\phi_u} L_u(F_{\theta,\phi_u})$;

8      **end**

9      Perform global update $\theta \leftarrow \theta - \beta \sum_u \nabla_\theta L_u(F_{\theta,\phi_u})$ ;

10     Perform global update $\phi \leftarrow \phi - \beta \sum_u \nabla_\phi L_u(F_{\theta,\phi_u})$

11 **end**

---

# MetaRec–AE

---

**Algorithm 4:** MAML Algorithm for MetaRec-AE

---

**Input** : Trainable global parameters $\theta$, user set $U$, local and global
hyper-parameters $\alpha, \beta$

**Output:** Learned global parameters $\theta$

1 initialize $\theta^{(0)}$;

2 $i = 0$;

3 **while** *not converged* **do**

4      **for** $u \in U$ **do**

5          update autoencoder

6          adapt locally: $\phi^u = \theta^{(0,i)} - \alpha \nabla_\theta L(\theta^{(0,i)})$

7      **end**

8      update globally: $\theta^{(0,i+1)} \leftarrow \theta^{(0,i)} - \beta \nabla_\phi L(\phi^u)$

9      $i \leftarrow i + 1$

10 **end**

---