# Recommendation Systems Overview

Qianchao Liu
Oct. 2017

# Outline

- Introduction
- Recommendation System
  - Collaborative Filtering
  - Feature based Ranking Model
    - Pointwise (CTR Models)
      - Linear Model
        - GLMMs, FFM
      - Tree Model
        - GBDT
    - Pairwise / Listwise
      - RankNet
      - LambdaMART
      - ListNet

# Introduction

- Background
  - Netflix: 80% of the movies watched are recommended
  - Google News: recommendation generate 38% more click-through
  - Amazon: 30% page views from recommendations

- Model Classification
  - Collaborative Filtering
  - Feature based Recommendation(Learning to rank)
    - Pointwise
    - Pairwise/Listwise

Ref: Two Decades of Recommender Systems at Amazon.com(Brent Smith, Greg Linden 2017)

# Collaborative Filtering

# Collaborative Filtering

- *What is CF:* recommendations are based on user's <u>past</u> behaviour
  - User based CF
    - Represent each User as a *N*-dimensional vector
    - Find top k similar Users for each target user
    - select recommendations from the similar users' items (count/score)

  - Item based CF (2003 Amazon.com)
    - Generate Item to Item co-rated matrix and compute the similarity
    - Predict ratings for the target item

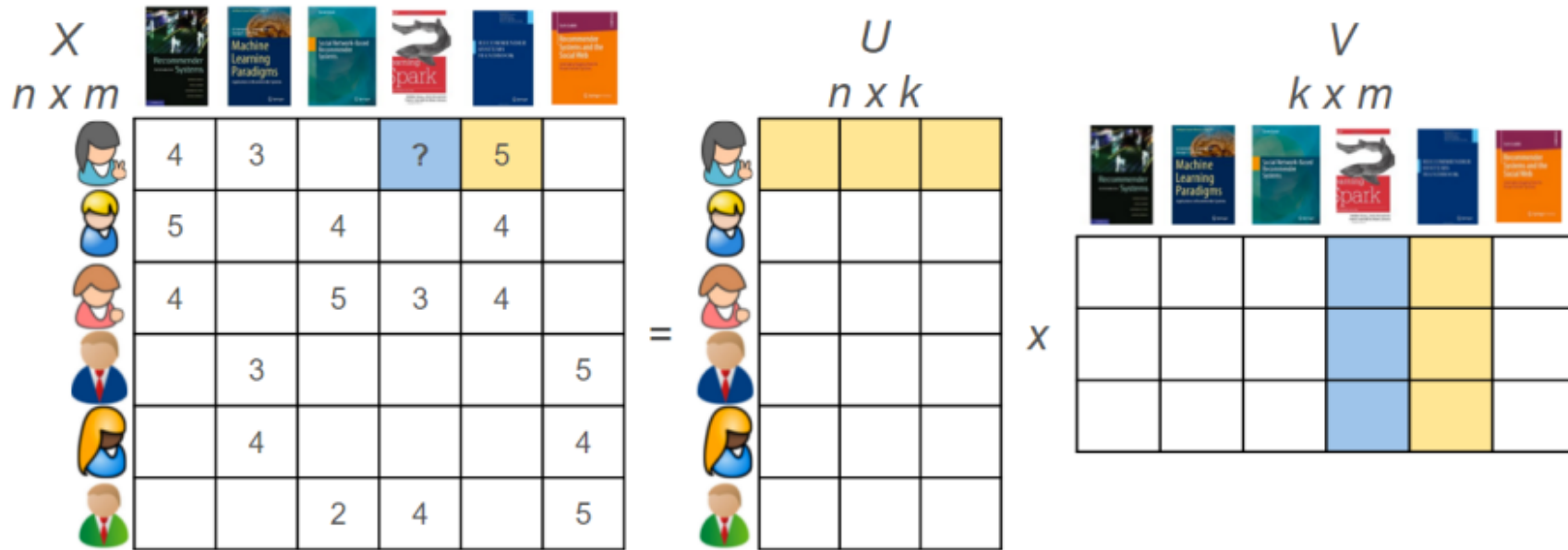$$\cos(\vec{w}_c, \vec{w}_s) = \frac{\vec{w}_c \cdot \vec{w}_s}{\|\vec{w}_c\|_2 \times \|\vec{w}_s\|_2}$$

$$sim(u,v) = \frac{\sum_{i \in I_{uv}} (y_{u,i} - \hat{y}_u)(y_{v,i} - \hat{y}_v)}{\sqrt{\sum_{i \in I_{uv}} (y_{u,i} - \hat{y}_u)^2 \sum_{i \in I_{uv}} (y_{v,i} - \hat{y}_v)^2}}$$

Cosine Similarity                                          Pearson Correlation

# Collaborative Filtering

- Matrix Factorization (2008 Netflix Competition)
  - learning: Gradient Descent with Square Loss



$$\hat{y}_{ij} = U_i^T V_j$$

$$e_{ij} = \frac{1}{2}(y_{ij} - \hat{y}_{ij})^2 + \frac{\lambda}{2}\sum_{p=1}^{P}(u_{pi}^2 + v_{pj}^2)$$

# Pointwise Feature-based Recommendation System

# Generalized Linear Model(GLM)

- ## Starting with GLM

    - $X$: Feature Vector
    - $\beta$: Coefficient Vector
    - $g(\cdot)$: Link Function

$$g(E[y]) = \beta X$$

| Model | Link Function | Mean Function | Application Content |
|---|---|---|---|
| Linear Regression | identity function: $\beta X = u$ | $u = \beta X$ | continuous variable |
| Logistic Regression | Logistic function: $\beta X = ln(\frac{u}{1-u})$ | $u = \frac{1}{1 + e^{-\beta X}}$ | binary varibale |
| Possion Regression | log function: $\beta X = ln(u)$ | $u = e^{\beta X}$ | count variable |

- GLMMs =  Fixed Effect (GLM) + Random Effect
- GLMMs allows the response to include information at the <span style="color:red">group level</span>

$$g(E[y]) = \beta X + uZ + e$$

Fixed     $u \sim \mathcal{N}(0, \sigma^2 I)$

- Multi-task Learning(Parameter based Transfer Learning)
  - For a user with many responses, it's able to accurately estimate personal coefficient
  - For a user doesn't have much past data, $u$ is close to 0 and model will fall back to the global fixed effect

# Generalized Linear Mixed Models(GLMMs): GLMMs

- GLMix: Generalized Linear Mixed Models For Large-Scale Response Prediction 2015
  - Application in *Linkedin* (with 400 million members)

- GLMix: Generalized Linear Mixed Models For Large-Scale Response Prediction 2015
  - Application in *Linkedin* (with 400 million members)

$$g(E[y]) = \underbrace{wX_{ij}}_{1} + \underbrace{\alpha_i X_j}_{2} + \underbrace{\beta_j X_i}_{3}$$

- Global Model: Linear combination between user and job profile
- Per-user Model: $x_{userid}x_{job\_title}$ if a member has applied to a job with title "Data Scientist"
- Per-item Model: $x_{jobid}x_{user\_years}$ if a job get an apply from a user with "4 year experiences"

| Model | LR-Cosine | LR-Full | MF | GLMix-Member | GLMix-Job | GLMix-Full | GLMix-Full+MF |
|-------|-----------|---------|------|--------------|-----------|------------|---------------|
| AUC | 0.664 | 0.723 | 0.772 | 0.780 | 0.758 | 0.799 | **0.801** |

Offline model performance for LinkedIn Job Recommendation data.

- Improving Demand Prediction in Bike Sharing System by Learning Global Features (2016)
  - Regression by individuals level + Global Features(GBDT, NNs)



The station-centric model with global feature transformation framework. The decision tree and neural network are trained on the entire dataset. Then the global features (gray) are transformed by the models from the station-centric features and concatenated to the original features. The regression model will be trained for each station individually. In this way, we can integrate the station-centric features and global features.

# FM & FFM (Factorization Machine & Field-aware FM)

- Introduce Feature Combination to Linear Model

$$y(X) = w_0 + \sum_{i=1}^{n} w_i x_i + \sum_{i=1}^{n-1} \sum_{j=i+1}^{n} w_{ij} x_i x_j$$

- Using latent vector

$$y(X) = w_0 + \sum_{i=1}^{n} w_i x_i + \sum_{i=1}^{n-1} \sum_{j=i+1}^{n} <v_i, v_j> x_i x_j \qquad \text{FM}$$

- Extending feature to field-feature

$$y(X) = w_0 + \sum_{i=1}^{n} w_i x_i + \sum_{i=1}^{n-1} \sum_{j=i+1}^{n} <v_{i,f_j}, v_{j,f_i}> x_i x_j \qquad \text{FFM}$$

- Definition
  - Decision rules same as in decision tree
  - Contains continuous/categorical value in each leaf for regression/ classification problem
- How to build a tree?
  - Branching criteria (decision stump/subset for different feature types)
    - Leaf value
      - classification: Majority (regression: Mean)
    - Evaluated by purifying
      - classification: Gini $1 - \sum_{k=1}^{K} \left( \frac{\sum_{n=1}^{N} [\![y_n = k]\!]}{N} \right)^2$ (regression: SSE)
  - Termination criteria
    - All $y_n$ the same / Height limitation
    - Pruning when growing or after fully-grown by $Error(G) + \lambda \Omega(G)$

# GBDT Model

- Functional Gradient Descent + Tree Model = Gradient Boosting Decision Tree

- Gradient Descent

$$f(x + x') \approx f(x) + f'(x)x' \qquad \longrightarrow \qquad \theta_{k+1} = \theta_k - \eta g_k$$

Taylor' formula  Gradient Descent

- Functional Gradient Descent

$$F_0(x) = f_0(x)$$
$$F_1(x) = f_0(x) + f_1(x) = F_0(x) + f_1(x)$$
$$F_2(x) = f_0(x) + f_1(x) + f_2(x) = F_1(x) + f_2(x)$$
$$\dots$$
$$F_m(x) = f_0(x) + \dots + f_m(x) = F_{m-1}(x) + f_m(x)$$

$$\longrightarrow$$

$$L(y, F_m(x)) = L(y, F_{m-1}(x) + f_m(x))$$
$$\approx L(y, F_{m-1}(x)) + \frac{\delta L(y, F_{m-1}(x))}{\delta F_{m-1}(x)} f_m(x)$$
$$\therefore -g_m f_m(x) < 0$$
$$\therefore \eta f_m(x) = g_m(x)$$

# GBDT Model

- Gradient Boosting Algorithm

$$g_m = -\left[\frac{\delta L(y_i,\ F(x))}{\delta F(x)}\right]_{F(x)=F_{m-1}(x)} \qquad \sum_{i=1}^{N}(g_{im}-h_m(x_i))^2 \qquad \gamma_m = argmin_\gamma\ L(y_i, F_{m-1}(x) + \gamma h_m(x))$$

gradient residual                    week learner                         step length

$$A\ F_m(x) = F_{m-1}(x) + \gamma_m h_m(\mathrm{x})$$

- Gradient Boosting Decision Tree
  - Compute gradient residual $g_m$ and use tree model to fit the residual
  - Compute $\gamma_m$ (line search step)to re-estimate the parameters at the leaves of the tree

$$\gamma_{jm} = argmin_\gamma \sum_{x_i \in R_{jm}} L(y_i, F_{m-1}(x_i) + \gamma)$$

$$F_m(x) = F_{m-1}(x) + \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$$

# Relationship between GBDT and XGBoost

- *XGBoost* (Extreme Gradient Boosting) is an implement of gradient boosting using C++

- Model Difference between GBDT and XGBoost
  - Use second order Taylor approximation
  - Regularization of XGBoost: $\gamma T + \frac{1}{2}\lambda \sum_{j=1}^{T} w_j^2$
  - Column subsampling

- Engineering Contribution
  - Support C++, Python, Java, R, etc.
  - Distributed

# Summary of Pointwise Approach

- Pointwise Approach
  - GLM
    - GLMMs: Individual Level
    - FFM: Feature Combination
  - Tree and Ensemble Model
    - GBDT: Feature Combination and Regularization


- However, they doesn't take into account the location of each alternative.

# Pairwise/Listwise Recommendation System

# Pairwise Ranking Model: RankNet

- *RankNet (Burges et al. 2005)* is used by Microsoft Bing Search

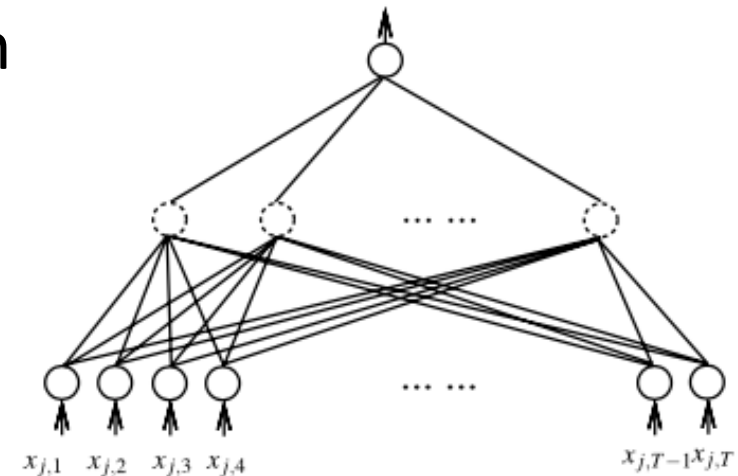- Evaluating target probability through scoring function $f(x; w)$

$$\hat{P}_{i,j}(f) = \frac{1}{1 + e^{-\sigma(f(x_i) - f(x_j))}}$$

- Using Cross Entropy loss function

$$C_{ij} = -P_{i,j}\log\hat{P}_{i,j}(f) - (1 - P_{i,j})\log(1 - \hat{P}_{i,j}(f)) \qquad C_{ij} = -P_{i,j}\log\hat{P}_{i,j}(f) - (1 - P_{i,j})\log(1 - \hat{P}_{i,j}(f))$$

- Using two layer neural nets as scoring function
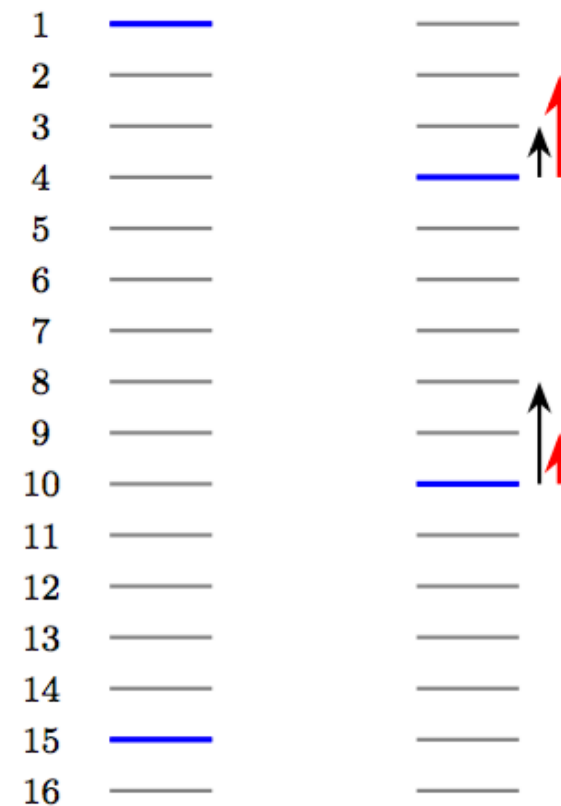
Scoring Function: $s = f(x; w)$

- Dive into RankNet

  - Single pair

    $$\frac{\delta C_{ij}}{\delta w} = \frac{\delta C_{ij}}{\delta s_i}\frac{\delta s_i}{\delta w} + \frac{\delta C_{ij}}{\delta s_j}\frac{\delta s_j}{\delta w} \qquad \text{where} \quad \frac{\delta C_{ij}}{\delta s_i} = -\frac{\delta C_{ij}}{\delta s_j} = -\frac{\sigma}{1+e^{\sigma(s_i-s_j)}} \overset{\text{def}}{=} \lambda_{ij}$$

    $$= \lambda_{ij}\left(\frac{\delta s_i}{\delta w} - \frac{\delta s_j}{\delta w}\right)$$

  - mini-batch

    $$\frac{\delta C}{\delta w} = \sum_{\{i,j\}\in I} \lambda_{ij}\left(\frac{\delta s_i}{\delta w} - \frac{\delta s_j}{\delta w}\right) = \left(\sum_{j:\{i,j\}\in I} \lambda_{ij} - \sum_{j:\{j,i\}\in I} \lambda_{ji}\right)\frac{\delta s_i}{\delta w}$$
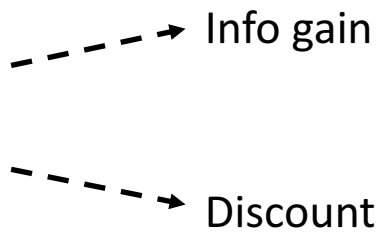
  - LambdaRank

    $$\lambda'_{ij} = \lambda_{ij}|\Delta(NDCG)|$$

# From RankNet to LambdaMART

- Ranking performance evaluation: Normalized Discounted Cumulative Gain(NDCG)

$$N = Z_n \sum_{j=1}^{n} \frac{2^{r(j)} - 1}{\log(j + 1)}$$

Info gain

Discount

- LambdaMART: Combining MART(Multivariate Adaptive Regression Trees) and LambdaRank

- LambdaMART won 1st at 2010 *Yahoo! Learning to Rank Challenge*

# Other Listwise Ranking Model

- *ListRank(2007):* using Plackett-Luce dist to model the uncertainty of a rank

$$P(\pi|s) = \prod_{j=1}^{m} \frac{\exp(s_j)}{\sum_{u=j}^{m} \exp(s_u)}$$

$$-\sum_{i}\sum_{\boldsymbol{\pi}} p(\boldsymbol{\pi}|\mathbf{y}_i) \log p(\boldsymbol{\pi}|\mathbf{s}_i)$$

evaluate the probability                    loss function

- It's intractable since π is over m!. We can consider permutations of top K: $A_m^k$

- In the special case where only one document is relevant——MNL.

$$P(\pi|s) = \frac{\exp(s_j)}{\sum_{u=1}^{m} \exp(s_u)}$$

$$-\sum_{i}\sum_{\boldsymbol{\pi}} p(\boldsymbol{\pi}|\mathbf{y}_i) \log p(\boldsymbol{\pi}|\mathbf{s}_i)$$