

GENERALIZED

LINEAR

MODELS

IN

COLLABORATIVE FILTERING

HAO WU

STANFORD UNIVERSITY

Generalized Linear Models

Generalization of linear regression from normal distribution to exponential family

Model components

- distribution from the **exponential family**
- linear parameter
- link function

$$g(y) = \eta = x^T \beta$$

- **sufficient statistic**

$$(y_1, y_2, \dots, y_n \mid x) \equiv (t, n \mid x)$$

Example: linear regression, **logistic regression**

COLLABORATIVE FILTERING

	Alt. Linear Regression	Alt. Logistic Regression
distribution	normal	binomial
link function	$\mu = \mathbf{U}\mathbf{M}^T$	$\log(\mathbf{P}/(1 - \mathbf{P})) = \mathbf{U}\mathbf{M}^T$
loss function	square error	logistic loss
sufficient stat.	sum	# of 1s
application	direct feedbacks (rating)	indirect feedbacks (click)

$$\begin{aligned}
 \min_{\mathbf{U}, \mathbf{M}} \quad & \sum_{i=0}^N L(y_i, \mathbf{u}_{u_i}^T \mathbf{m}_{m_i}) \\
 & + \lambda \left[\alpha \left(\sum_i^{n_U} \|\mathbf{u}_i\|_1 + \sum_i^{n_M} \|\mathbf{u}_i\|_1 \right) + (1 - \alpha) (\|\mathbf{U}\|_F + \|\mathbf{M}\|_F) \right] \\
 \text{s.t.} \quad & \mathbf{U} \in \mathbb{R}^{n_U \times k}, \mathbf{M} \in \mathbb{R}^{n_M \times k}
 \end{aligned}$$

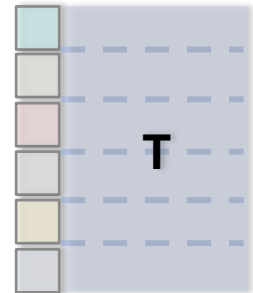
DISTRIBUTED ALG.

(assuming $n_m \times k$ fits in a single machine)

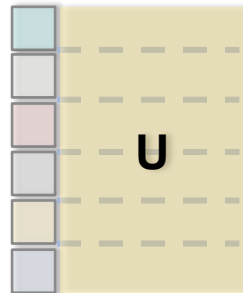
$\text{RDD}[u, m, y]$

 $\xrightarrow[\text{aggregateByKey}]{\text{reduceByKey}}$
 $\text{RDD}[u, \text{map}(m \rightarrow (t, n))]$

Input entries
 Sufficient Stat.



$\text{RDD}[u, \text{vector}]$
 User features



Array[vector]
 Movie features



for each iteration:

- Join **U** with **T** to form **D** (co-partitioned join)
- Update **M**
- Broadcast **M** (communication: $\log(p)(n_M k)$)
- Update **U**

Update M

for each movie:

- prepare dataframe by `filter()` and `map()`
- distributed logistic regression
 - `LogisticRegression()`
 - all-to-one and one-to-all of size k

Update U

Map local logistic regression to users

- added local training method to `LogisticRegression()`
- no communication of data

Summary

- Sparsity is preserved
- Scales in n_U , but not n_M or k
- Communication cost: $\log(p)(n_M k)$
- Computational depth: $\log(n_U)(n_M k)$

