

Implementação do algoritmo Simultaneous Co-Clustering and Learning utilizando Spark Python API

Fabrcio Costa Antoniasse

1 Introduo

Induzir modelos preditivos nicos que caracterizem a base de dados em estudo nem sempre uma tarefa fcil. Nestes casos, frequentemente, as bases de dados s3o particionadas em grupos a fim de gerar modelos de predies sobre elas, conhecido na literatura como mltiplos modelos localizados. Neste trabalho ser3a apresentado uma possvel implementao, em Apache Spark, do algoritmo SCOAL (Simultaneous Co-Clustering and Learning), o qual realiza de forma iterativa o agrupamento dos dados e o treinamento dos modelos preditivos para cada grupo criado, simultaneamente.

2 SCOAL

Proposto por Deodhar e Ghosh, o algoritmo SCOAL intercala as tarefas de agrupamento e treinamento dos modelos de predio, melhorando assim, tanto a funao de agrupar objetos semelhantes entre si como tambm de gerar modelos preditivos.

Em um contexto de sistema de recomendao, dada uma matriz de avaliao \mathbf{Z} com dimens3o $m \times n$, onde m representa o nmero total de usu3rios e n a quantidade de produtos, cada clula \mathbf{z}_{ij} corresponde a avaliao dada pelo usu3rio i ao produto j , no qual $\mathbf{z}_{ij} \in \mathbf{R}$. Todo par usu3rio-produto descrito por um vetor covariante \mathbf{x}_{ij} , composto pela concatenao dos atributos dos usu3rios \mathbf{u}_i , como por exemplo idade, profisso, sexo, com os atributos dos produtos \mathbf{p}_j , os quais poderiam ser data de lanamento, gnero do filme, diretor etc. A matriz \mathbf{Z} particionada em l clusters de linha e c clusters de coluna, resultando em $l \times c$ co-clusters. Para cada co-cluster um modelo preditivo ser3a induzido. Assume-se que o modelo de predio gerado para cada co-cluster seja construido a partir de uma regress3o linear, dado por:

$$\mathbf{z}_{ij} = \beta^T \mathbf{x}_{ij} + \mathcal{N}(0, \sigma^2) \quad (1)$$

onde $\mathbf{x}_{ij}^T = [1, \mathbf{u}_i^T, \mathbf{p}_j^T]$ e $\beta^T = [\beta_0, \beta_u^T, \beta_p^T]$.

Uma vez estimado β , a partir dos dados hist3ricos representados pela matriz \mathbf{Z} , possvel fazer predies para qualquer par usu3rio-produto, atrav3s de:

$$\hat{\mathbf{z}}_{ij} = \beta^T \mathbf{x}_{ij} \quad (2)$$

A funao objetivo a ser minimizada 3:

$$\sum \mathbf{w}_{ij} (\mathbf{z}_{ij} - \hat{\mathbf{z}}_{ij})^2 \quad (3)$$

onde \mathbf{z}_{ij} o valor observado, $\hat{\mathbf{z}}_{ij}$ o valor estimado e \mathbf{w}_{ij} assume o valor 1 caso \mathbf{z}_{ij} seja conhecido e 0 caso contr3rio. Neste trabalho, a fim de estimar os valores de β , ser3a utilizada a t3cnica do gradiente descendente. Os valores de β ser3o ajustados atrav3s da equao 4.

$$\beta_{new} = \beta_{old} + \alpha \sum (\mathbf{z} - \hat{\mathbf{z}}) \mathbf{x} \quad (4)$$

onde α a taxa de aprendizado, \mathbf{z} o valor observado, $\hat{\mathbf{z}}$ o valor estimado e \mathbf{x} o vetor contendo os atributos.

2.1 Algoritmo SCOAL

O SCOAL recebe uma matriz de avaliao \mathbf{Z} a qual particionada em l grupos de usu3rios e c grupos de produtos, onde cada usu3rio i mapeado por uma funao $\rho(i)$ a um cluster de linha e de forma an3loga, cada produto ser3a mapeado por uma funao $\gamma(j)$ a um cluster de coluna. Em seguida efetuada de forma iterativa, at3 a converg3ncia, os seguintes passos:

1º passo: ajusta-se os modelos para cada co-cluster, i.e., encontra-se os valores de β_{lc} .

2º passo: efetua-se o agrupamento dos usuários mais similares entre si de forma indireta, ou seja, o agrupamento é realizado através dos valores preditos dos modelos induzidos da seguinte maneira:

⇒ deseja-se encontrar, para cada usuário, o cluster de linha l que minimize o erro global dada pela equação 5:

$$\rho(i) = \underset{g}{\operatorname{argmin}} \sum_{g=1}^l \sum_{j=1}^n \mathbf{w}_{ij} (\mathbf{z}_{ij} - \hat{\mathbf{z}}_{ij})^2 \quad (5)$$

onde $\hat{\mathbf{z}}_{ij} = \beta_{g\gamma(j)}^T \mathbf{x}_{ij}$. Efetua-se o mesmo processo em relação aos produtos, como mostra a equação 6:

$$\gamma(j) = \underset{h}{\operatorname{argmin}} \sum_{h=1}^c \sum_{i=1}^m \mathbf{w}_{ij} (\mathbf{z}_{ij} - \hat{\mathbf{z}}_{ij})^2 \quad (6)$$

onde $\hat{\mathbf{z}}_{ij} = \beta_{\rho(i)h}^T \mathbf{x}_{ij}$.

O algoritmo SCOAL, resumido, está representado por meio do Algoritmo 1.

Algoritmo 1: Algoritmo SCOAL (adaptado de (DEODHAR; GHOSH))

Entrada: \mathbf{Z} , \mathbf{x}_{ij}

Saída: Partição dos co-clusters com seus respectivos modelos, β_{lc}

1 Inicializar randomicamente os co-clusters.

2 **repita**

3 Gerar, para cada co-cluster, os modelos preditivos.

4 Atualize $\rho(i)$ - declarar cada linha a um cluster de linha que minimize o erro.

5 **para** $i=1$ até m **faça**

6 $\rho(i) = \underset{g}{\operatorname{argmin}} \sum_{g=1}^l \sum_{j=1}^n \mathbf{w}_{ij} (\mathbf{z}_{ij} - \hat{\mathbf{z}}_{ij})^2$

7 **fim**

8 Atualize $\gamma(j)$ - declarar cada coluna a um cluster de coluna que minimize o erro.

9 **para** $j=1$ até n **faça**

10 $\gamma(j) = \underset{h}{\operatorname{argmin}} \sum_{h=1}^c \sum_{i=1}^m \mathbf{w}_{ij} (\mathbf{z}_{ij} - \hat{\mathbf{z}}_{ij})^2$

11 **fim**

12 **até** o critério de convergência ser atingido;

13 **retorna** $(\rho(i), \gamma(j))$ e β_{lc}

2.2 Implementação em Apache Spark

Neste trabalho será apresentada a implementação da geração dos modelos dos co-clusters e o rearranjo dos usuários aos clusters de linha. Vale salientar que o processo de atualização dos produtos aos clusters de coluna é realizado de forma semelhante ao dos ajustes dos clusters de linha. Foi utilizado o Spark Python API (PySpark) para a implementação do algoritmo SCOAL. O RDD "dataCluster" é composto por:

`dataCluster.take(1) = [idUser, idItem, rating, [1, atbuser1, atbuser2, ..., atbuserp, atbitem1, atbitem2, ..., atbitemq], (($\rho(i), \gamma(j)$), $\rho(i), \gamma(j)$)]`

2.2.1 Geração dos modelos preditivos

A estratégia consiste em dividir a equação 4 em 2 partes:

1 parte - Calcula-se o somatório $\sum (\mathbf{z} - \hat{\mathbf{z}}) \mathbf{x}$

2 parte - Adiciona-se o β com o resultado de $\alpha \sum (\mathbf{z} - \hat{\mathbf{z}}) \mathbf{x}$ a fim de atualizar os valores de β .

<pre>def nabla(vet_x, n_beta, y, pesos): produto = np.dot(vet_x, pesos[n_beta-1][1]) return np.dot((y- produto), vet_x)</pre>	<p>A função "nabla" recebe um vetor covariante \mathbf{x}_{ij}, o id do co-cluster associado ao vetor \mathbf{x}_{ij} e a matriz de pesos β.</p> <p>Tem como retorno o cálculo de $(\mathbf{z} - \hat{\mathbf{z}}) \mathbf{x}$</p>
<pre>def addBeta(n_beta, vet_xB, pesos): return np.array(pesos[n_beta-1][1]) + vet_xB</pre>	<p>A função "addBeta" recebe como parâmetros o id do co-cluster, o produto do resultado da função "nabla" e α e a matriz de pesos β.</p> <p>Tem como retorno o valor dos pesos β atualizados.</p>

```

def modelos(pesos,dataCluster,alpha):
    bNovo = dataCluster.map(lambda x : (x[-1][0],nabla(x[3],x[-1][0],x[2],pesos)))
        .reduceByKey(lambda x,y : np.array(x)+np.array(y))
        .map(lambda x : (x[0],addBeta(x[0],x[1]*(alpha),pesos)))
        .sortByKey()
    return bNovo.collect()

```

2.2.2 Atualização de $\rho(i)$ - rearranjo dos usuários aos clusters de linha

```

def mse(vet_x, n_beta,y,peso):
    produto = np.dot(vet_x,peso[n_beta-1][1])
    return pow((y - produto),2)

```

A função "mse" recebe um vetor covariante \mathbf{x}_{ij} , o id do co-cluster associado ao \mathbf{x}_{ij} e a matriz de pesos β .

Tem como retorno o valor do erro médio quadrático (MSE).

```

def mseClustersLinha (user,rating,n_clus,clusterLinha,cl_L,cl_C,peso):
    quantidadeClusterLinha = len(cl_L)
    quantidadeClusterColuna = len(cl_C)
    aux = n_clus - clusterLinha*(quantidadeClusterColuna)
    erros = []
    for linha in range(quantidadeClusterLinha):
        erros.append(mse(user,aux,rating,peso))
        aux += quantidadeClusterColuna
    return erros

```

A função "mseClustersLinha" recebe como parâmetros o id do usuário, \mathbf{z}_{ij} , id do co-cluster ($\rho(i), \gamma(j)$), $\rho(i)$, listas contendo os limites dos clusters de linha e de coluna e a matriz de pesos β .

Tem como retorno uma lista contendo os erros globais de cada cluster de linha de forma ordenada.

```

def def mudanca(x_erros):
    return x_erros.index(min(x_erros))

```

A função "mudanca" recebe como parâmetro uma lista contendo os erros dos clusters de linha.

Ela retorna a posição do menor valor encontrado na lista.

```

mudancaL = dataCluster.map(lambda x : [x[0], mseClustersLinha(x[3],x[2],x[-1][0],x[-1][1],cl_L,cl_C,beta)])
    .reduceByKey(lambda x ,y : np.array(x)+np.array(y))
    .mapValues(lambda x : mudanca(list(x)))
dsAuxiliar = dataCluster.map(lambda x : [x[0], x[1:]])
    .leftOuterJoin(mudancaL)
    .map(lambda x : [x[0]]+ x[1][0][:-1] +[(x[1][0][-1][0],x[1][-1],x[1][0][-1][2])])

```

Referências

DEODHAR, Meghana; GHOSH, Joydeep. A framework for simultaneous co-clustering and learning from complex data. In: Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining. ACM, 2007. p. 250-259.

KARAU, Holden et al. Learning spark: lightning-fast big data analysis. " O'Reilly Media, Inc.", 2015.

PEREIRA, Andre Luiz Vizine; HRUSCHKA, Eduardo Raul. Simultaneous co-clustering and learning to address the cold start problem in recommender systems. Knowledge-Based Systems, v. 82, p. 11-19, 2015.