

CEFET/RJ
Programa de Pós-graduação
em Ciência da Computação
Aprendizado de Máquina - Trabalho 02

Prof. Eduardo Bezerra (ebezerra@cefet-rj.br)

Abril/2019

Conteúdo

1	Introdução	3
2	Sistemas de Recomendação	3
2.1	Conjunto de dados de classificações de filme	3
2.2	Algoritmo de aprendizagem de filtragem colaborativa	3
2.2.1	Função de custo da filtragem colaborativa	4
2.2.2	Gradiente de filtragem colaborativa	4
2.3	Aprendizado de Recomendações para Filmes	4
2.4	Filtragem Colaborativa - método baseado em memória	5
3	Aprendizado de Comitês	6
4	O que deve ser entregue	6
5	Créditos	7

1 Introdução

Neste trabalho, você usará filtragem colaborativa para criar um sistema de recomendação para filmes (Seção 2). Você também irá realizar experimentos no contexto do aprendizado de comitês (Seção 3).

2 Sistemas de Recomendação

Nesta parte, você implementará o algoritmo de aprendizagem de filtragem colaborativa e aplicá-lo-á a um conjunto de dados de avaliações de filmes¹. Este conjunto de dados consiste em classificações em uma escala de 1 a 5. O conjunto de dados tem $n_u = 943$ usuários e $n_m = 1682$ filmes.

2.1 Conjunto de dados de classificações de filme

O arquivo que contém o conjunto de dados tem nome `ex8_movies.mat`, e contém as variáveis Y e R . A matriz Y (de ordem *número de filmes* x *número de usuários*) armazena as classificações $y^{(i,j)}$ (de 1 a 5). A matriz R é uma matriz de indicadores de valor binário, onde $R(i, j) = 1$ se o usuário j forneceu uma classificação para o filme i e $R(i, j) = 0$ em caso contrário. O objetivo da filtragem colaborativa é prever as classificações de filmes para os filmes que os usuários ainda não classificaram, ou seja, as entradas com $R(i, j) = 0$. Isso permitirá recomendar os filmes com classificações mais altas previstas para outros usuários.

Ao longo desta parte, você também irá trabalhar com as matrizes X e Θ . Essas variáveis se encontram no arquivo `ex8_movieParams.mat`. A i -ésima linha de X corresponde ao vetor de característica $x^{(i)}$ para o i -ésimo filme. A j -ésima linha de Θ corresponde a um vetor de parâmetros $\theta^{(j)}$, para o j -ésimo usuário. Tanto $x^{(i)}$ quanto $\theta^{(j)}$ são vetores n -dimensionais. Para os fins deste exercício, você usará $n = 100$ e, portanto, $x^{(i)} \in \mathbb{R}^{100}$ e $\theta^{(j)} \in \mathbb{R}^{100}$. Correspondentemente, X é uma matriz $n_m \times 100$ e Θ é uma matriz $n_u \times 100$.

2.2 Algoritmo de aprendizagem de filtragem colaborativa

Agora, você vai começar a implementar o algoritmo de filtragem colaborativa. Você iniciará implementando a função de custo (sem regularização). O algoritmo de filtragem colaborativa no contexto das recomendações de filmes considera um conjunto de vetores de parâmetros n -dimensionais $x^{(1)}, \dots, x^{(n_m)}$ e $\theta^{(1)}, \dots, \theta^{(n_u)}$, onde o modelo prediz a avaliação para o filme i pelo usuário j como $y^{(i,j)} = (\theta^{(j)})^T x^{(i)}$. Dado um conjunto de dados que consiste em um conjunto de avaliações produzidas por alguns usuários para alguns filmes, o algoritmo deve aprender os vetores de parâmetros $x^{(1)}, \dots, x^{(n_m)}$ e $\theta^{(1)}, \dots, \theta^{(n_u)}$ que produzem o melhor ajuste (i.e., minimizam o erro quadrático).

Você deverá implementar código em um arquivo de nome `cofi_cost_func.py` para calcular a função de custo e o gradiente para a filtragem colaborativa. Os

¹<https://grouplens.org/datasets/movielens/>

parâmetros para a função (ou seja, os valores que você está tentando aprender) devem ser X e $Theta$.

2.2.1 Função de custo da filtragem colaborativa

A função de custo para a filtragem colaborativa (sem regularização) é dada por

$$J(x^{(i)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}) = \frac{1}{2} \sum_{(i,j):r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2.$$

Você deve implementar código em `cofi_cost_func.py` para retornar esse custo em uma variável de nome J . Note que você deve acumular o custo para o usuário j e o filme i somente se $R(i, j) = 1$.

Depois de implementar a função, você deve testá-la. Você deve ver uma saída de 22,22.

2.2.2 Gradiente de filtragem colaborativa

Agora, você deve implementar o gradiente (sem regularização). Especificamente, você deve implementar código em `cofiCostFunc.py` para retornar as variáveis X_grad e $Theta_grad$. Observe que X_grad deve ser uma matriz do mesmo tamanho que X e, de forma semelhante, $Theta_grad$ é uma matriz do mesmo tamanho que $Theta$. Os gradientes da função de custo são dados por:

$$\begin{aligned} \frac{\partial J}{\partial x_k^{(i)}} &= \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) \theta_k^{(j)} \\ \frac{\partial J}{\partial \theta_k^{(j)}} &= \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)}. \end{aligned}$$

A função deve retornar o gradiente para ambos os conjuntos de variáveis, armazenando-os em um único vetor.

2.3 Aprendizado de Recomendações para Filmes

Depois de concluir a implementação do gradiente e da função de custo de filtragem colaborativa, agora você pode começar a treinar seu algoritmo para fazer recomendações de filmes. Para testar se o treinamento funcionou, você pode definir suas próprias preferências de filmes, para que mais tarde, quando o algoritmo for executado, você possa obter suas próprias recomendações de filmes. A lista de todos os filmes e seus números no conjunto de dados pode ser encontrada no arquivo `movie_ids.txt`.

Depois que você definir suas para alguns filmes e adicioná-las ao conjunto de dados, o script continuará a treinar o modelo de filtragem colaborativa. Isso deve fazer com que os parâmetros X e $Theta$ sejam aprendidos. Para prever a classificação do filme i para o usuário j , você precisa calcular $(\theta^{(j)})^T x^{(i)}$. Você deve implementar um script para calcular as classificações principais para

Principais recomendações:
Previsão de avaliação 9.0 para Titanic (1997)
Previsão de avaliação 8.9 para Star Wars (1977)
Previsão de avaliação 8.8 para Shawshank Redemption, The (1994)
Previsão de avaliação 8.5 para As Good As It Gets (1997)
Previsão de avaliação 8.5 para Good Will Hunting (1997)
Previsão de avaliação 8.5 para Usual Suspects, The (1995)
Previsão de avaliação 8.5 para Schindler's List (1993)
Previsão de avaliação 8.4 para Raiders of the Lost Ark (1981)
Previsão de avaliação 8.4 para Empire Strikes Back, The (1980)
Previsão de avaliação 8.4 para Braveheart (1995)
Avaliações originais fornecidas:
Avaliou 4 para Toy Story (1995)
Avaliou 3 para Twelve Monkeys (1995)
Avaliou 5 para Usual Suspects, The (1995)
Avaliou 4 para Outbreak (1995)
Avaliou 5 para Shawshank Redemption, The (1994)
Avaliou 3 para While You Were Sleeping (1995)
Avaliou 5 para Forrest Gump (1994)
Avaliou 2 para Silence of the Lambs, The (1991)
Avaliou 4 para Alien (1979)
Avaliou 5 para Die Hard 2 (1990)
Avaliou 5 para Sphere (1998)

Tabela 1: Exemplo de saída produzida pelo script de geração de recomendações.

um dados usuário e exibir os filmes recomendados. A saída desse script deve ser parecida com a apresentada na Tabela 2.3, de acordo com as classificações definidas anteriormente. Observe que você pode obter um conjunto diferente de previsões diferentes devido a inicializações aleatórias.

2.4 Filtragem Colaborativa - método baseado em memória

Nessa parte do trabalho, você irá implementar a filtragem colaborativa por meio do método baseado em memória, visto em aula. Inicialmente, leia o artigo intitulado *Empirical Analysis of Predictive Algorithms for Collaborative Filtering*. Para esta parte do trabalho desse artigo, você precisa ler a Seção 2.1, mas é encorajado a ler o artigo todo se tiver tempo.

O conjunto de dados que você deve usar é um subconjunto dos dados de classificações de filmes do Prêmio Netflix (<https://www.netflixprize.com>). Ele contém um conjunto de treinamento, um conjunto de testes, um arquivo de filmes, um arquivo de descrição de conjunto de dados e um arquivo README. Os conjuntos de treinamento e teste são ambos subconjuntos dos dados de treinamento do Netflix. Você usará as classificações fornecidas no conjunto de treinamento para prever as do conjunto de testes. Você irá comparar suas pre-

visões com as classificações reais fornecidas no conjunto de testes. As métricas de avaliação que você precisa medir são o *erro absoluto médio* (*mean absolute error*) e o Erro Quadrático Médio da Raiz (*root-mean-square error* (RMSE)). O arquivo de descrição do conjunto de dados descreve ainda mais o conjunto de dados e ajudará você a começar. O arquivo README é do conjunto original de arquivos Netflix e foi incluído para cumprir os termos de uso desses dados.

Implemente o algoritmo de filtragem colaborativa descrito na Seção 2.1 do artigo (Equações 1 e 2; ignore a Seção 2.1.2) para fazer as previsões.

Crédito extra: adicione a si mesmo como um novo usuário ao conjunto de treinamento. Para fazer isso, você precisará criar um novo ID de usuário para você mesmo. Selecione alguns filmes que você viu entre os que estão no conjunto de treinamento e adicione suas classificações para eles no conjunto de treinamento. Estenda seu sistema para gerar previsões para os filmes que você não classificou, classificados em ordem decrescente de classificação. Você concorda com as previsões do seu sistema? Confira alguns dos filmes mais bem classificados que você não viu (mas somente depois de terminar o trabalho no projeto).

3 Aprendizado de Comitês

Nesta seção, você irá tomar como ponto de partida o experimento contido no notebook fornecido na aula sobre aprendizado de comitês. Esse notebook treina um comitê de classificadores utilizando o algoritmo AdaBoost, cuja implementação em Python é fornecida. Sua tarefa nessa parte do trabalho é usar a biblioteca scikit-learn para replicar o experimento do notebook, com a diferença que você irá a classe `sklearn.ensemble.AdaBoostClassifier`. Compare os resultados obtidos no notebook e os resultados obtidos com sua implementação. Para isso, apresente em um mesmo gráfico (que plota a taxa de erro contra a quantidade de iterações) as taxas de erro de treinamento e de teste, tanto para a implementação fornecida no notebook, quanto para sua implementação usando a classe `sklearn.ensemble.AdaBoostClassifier`. (de 1 até 400, de 10 em 10).

Agora, usando a classe `sklearn.ensemble.BaggingClassifier`, realize um experimento também similar ao notebook fornecido. Para isso, apresente em um mesmo gráfico (que plota a taxa de erro contra a quantidade de iterações) as taxas de erro de treinamento e de teste obtidas com sua implementação usando a classe `sklearn.ensemble.BaggingClassifier`. Use a mesma faixa de variação de componentes utilizada no notebook (de 1 até 400, de 10 em 10).

4 O que deve ser entregue

Você deve preparar um único relatório para a apresentar sua análise e conclusões sobre as diversas partes desse trabalho. O formato desse relatório deve ser em PDF. Alternativamente à entrega do relatório em PDF, você pode entregar um notebook Jupyter².

²<http://jupyter.org/>

Independente de escolher entregar um relatório em PDF ou na forma de um notebook Jupyter, entregue também todos os arquivos em Python que você criou para cada parte deste trabalho. Todos os arquivos em Python devem estar em uma única pasta.

Crie um arquivo compactado que contém o relatório (ou notebook Jupyter) e os arquivos (*scripts*) em Python. Esse arquivo compactado deve se chamar `SEU_NOME_COMPLETO_T2.zip`. **Não serão aceitos arquivos que usem outro formato de compactação que não seja ZIP.** Esse arquivo compactado deve ser entregue pelo Moodle, até a data acordada.

5 Créditos

Esse trabalho é uma tradução/adaptação dos *programming assignments* encontrados no curso *Machine Learning*³ encontrado no Coursera. O material original é de autoria do prof. Andrew Ng. A parte do trabalho sobre filtragem colaborativa baseada em memória é uma adaptação do projeto de programação da disciplina Data Mining, de autoria do prof. PEdro Domingos (<https://courses.cs.washington.edu/courses/csep546/17au/>)

³<https://www.coursera.org/learn/machine-learning>