

Prova 1

1. A ideia principal desse tipo de filtragem é de se aproveitar da "wisdom of the crowds". User apóia a rating média nos dá alguns problemas:

a. Podemos ter um item 1 com 5 avaliações máximas: tendo uma média de 5. Ao mesmo tempo, podemos ter um item 2 com 20 mil avaliações cuja média é 4,89. O segundo item é mais popular e, mesmo assim, têm boa avaliação - ele deveria ser o recomendado.

b. Os usuários não avaliam da mesma forma. Alguns são mais gentis (avaliações mais altas), enquanto outros são mais críticos (avaliações mais baixas). Se levarmos em conta a média do item, sem levar em conta esse bias, não vamos ter uma boa representação das preferências.

2. a. $m = 10$ e $n = 1.000.000$

Item-based: a combinação de poucos usuários e muitos itens gera vetores muito esparsos. Posteriormente, temos a vantagem de computar a lista de vizinhos para cada item de forma offline.

b. $m = 1.000.000$ e $n = 10$.

User-based: como há pouca quantidade de itens os vetores de usuários não são esparsos.

c. $m = 1.000.000$ e $n = 1.000.000$

Item-based: podemos aproveitar o fato de que os vetores de item são menos esparsos e mais estáveis. Essa estabilidade permite a pré-computação das similaridades p/ todos os pares de itens. Posteriormente, não é necessário manter todos os vizinhos de modo a aumentar eficiência.

3.

Vantagens:

1. Dados futuros não vazam. Ex. não sabemos que um usuário assistiu a Star Wars Ep. V antes de ver a Ep. IV (pois que ele assiste a trilogia na ordem cronológica)
2. Mix redistrito de cold e não-cold start cases.
3. É possível verificar consumos negativos com o uso de janelas deslizantes.

Desvantagens

1. O tempo em que o corte será feito pode afetar negativamente os resultados
2. Podemos fazer uma partição que sofre de efeitos negativos, fato que nos forçaria a utilizar o método de janelas deslizantes para a correção.

4. Meu usuário: U₀

↳ queremos: item-based / similaridade coseno / $K=3$

- Cálculo da similaridades (U₀ analisei 1, 3, 5)

$$S(0, 1) = 0,19 \quad S(2, 1) = 0,43 \quad S(4, 1) = 0,17$$

$$S(0, 3) = 0,38 \quad S(2, 3) = 0,48 \quad S(4, 3) = 0,16$$

$$S(0, 5) = 0,55 \quad S(2, 5) = 0,19 \quad S(4, 5) = 0,36$$

- Aplicando esses valores em $N_{ui} = N_{oi} = [1, 3, 5]$

$$\hat{\mu}_{00} = \frac{0,19 \cdot 4 + 0,38 \cdot 3 + 0,55 \cdot 3}{0,19 + 0,38 + 0,55} = 3,17$$

$$\hat{\mu}_{02} = \frac{0,43 \cdot 4 + 0,48 \cdot 3 + 0,19 \cdot 3}{0,43 + 0,48 + 0,19} = 3,38$$

$$\hat{\mu}_{04} = \frac{0,17 \cdot 4 + 0,16 \cdot 3 + 0,36 \cdot 3}{0,17 + 0,16 + 0,36} = 3,24$$

Portanto,

$$\hat{\pi}_{00} = 3,17 \quad | \quad \hat{\pi}_{02} = 3,38 \quad | \quad \hat{\pi}_{04} = 3,24$$

5. $RMSE_0 \rightarrow n=3$ tal que o ranking é $[2, 4, 0]$

$$\sqrt{\frac{1}{3} \left[\underbrace{(\pi_{00} - \hat{\pi}_{00})^2}_{\substack{1 \\ 4,7089}} + \underbrace{(\pi_{02} - \hat{\pi}_{02})^2}_{\substack{2 \\ 1,9044}} + \underbrace{(\pi_{04} - \hat{\pi}_{04})^2}_{\substack{4 \\ 0,5776}} \right]}$$

$$RMSE_0 = 1,55$$

$$DCG_0 = \left(\frac{2^{\pi_{02}} - 1}{\log_2 2} + \frac{2^{\pi_{04}} - 1}{\log_2 3} + \frac{2^{\pi_{00}} - 1}{\log_2 4} \right)$$

$$DCG_0 = 12,96$$