

Uma versão distribuída do algoritmo Apriori em PySpark

Pedro Henrique Parreira
Santo André, São Paulo, Brasil

Abstract

O algoritmo apriori tem um custo computacional bastante elevado devido as inúmeras combinações de elementos que ele gera. Quando aplicado a um banco de dados bastante volumoso pode se tornar um grande desafio por necessitar de um grande volume de memória principal para armazenar e processar todas as combinações. Uma das possíveis saídas para tal problema é o uso da computação distribuída. Neste artigo iremos analisar uma versão distribuída do algoritmo apriori aplicada a uma grande base de dados e compará-la com a sua versão não distribuída.

Keywords: Apriori, Mineração de dados, PySpark,

1. Algoritmo Apriori

Um banco de dados de transações é formado por duas colunas, a primeira contendo o identificador da transação e a segunda os itens desta transação. Um exemplo clássico de um banco de dados de transações é o do histórico de vendas de um supermercado onde cada venda tem um identificador único e os itens da venda.

Sendo $I = \{i_1, i_2, \dots, i_m\}$ o conjunto dos itens de um banco de dados de transações, definimos um *itemset* S como $S \subseteq I$, ou seja, um itemset é um subconjunto de itens do banco de dados. Para todo itemset associamos um *suporte* que é definido como a porcentagem das transações em que o itemset é encontrado.

O objetivo do algoritmo apriori é gerar regras de associação entre os itemsets, sob a forma $A \rightarrow B$, onde tal regra indica uma relação interessante entre tais itemsets. E para toda regra de associação medimos a *confiança* desta como:

$$conf(A \rightarrow B) = \frac{suporte(A \cup B)}{suporte(A)}$$

Ou seja, a confiança de uma regra de associação é dada através da divisão entre a porcentagem de transações em que aparece o itemset $A \cup B$ e a porcentagem de transações em que aparece o itemset A .

O usuário fornece ao algoritmo o banco de dados, o *suporte mínimo* e *confiança mínima* e obtém como saída tais regras de associação.

O algoritmo, primeiramente, gera itemsets de tamanho 1, ou seja, formado apenas por um item que tenha o suporte maior que o suporte mínimo e a cada iteração vai aumentando de forma unitária o tamanho dos itemsets gerados, que tenham o suporte maior que o suporte mínimo, até não ser mais possível gerar itemsets para um determinado tamanho.

A cada iteração é feito o processo de poda que consiste em descartar das transações itens que não aparecem em nenhum itemset, pois tais itens possuem um suporte menor que o suporte mínimo.

Ao final, os itemsets selecionados são usados para gerar as regras de associações, fazendo a combinação de tais itemsets, verificando quais associações possuem confiança maior que a confiança mínima definida pelo usuário usando os suportes que foram gerados na fase anterior.

2. Versão distribuída do algoritmo apriori

A implementação da versão distribuída do algoritmo consiste em duas etapas: a geração dos itemsets e a geração das regras de associação.

2.1. Geração dos itemsets

Cada transação é representada como uma tupla onde cada elemento da tupla é um item da transação. Como por exemplo ('1','2','3'), esta tupla representa uma transação com três itens.

A geração dos itemsets começa com os de tamanho um e a cada nova iteração o tamanho é aumentado de forma unitária. Para um dado momento, seja k o tamanho dos itemsets que vão ser gerados, então para cada transação são realizadas $\binom{N}{k}$, onde N é o número de itens da transação em questão, combinações simples de tamanho k , onde cada combinação é um itemset.

Por exemplo, suponha que $k = 1$ e tomemos a transação ('1','2','3'), portanto são gerados três itemsets: ('1','2'), ('1','3'), ('2','3'). Cada itemset gerado é utilizado para gerar uma nova tupla, onde o itemset é utilizado como chave seguido de um segundo elemento que será o número

50 um, portanto considerando o exemplo acima o resultado do mesmo seria:
51 $((\text{'1'}, \text{'2'}), 1), ((\text{'1'}, \text{'3'}), 1), ((\text{'2'}, \text{'3'}), 1)$.

52 A implementação da geração das tuplas contendo os itemsets e o número
53 um é feita com o uso da transformação *flatMap*. Ao final desta etapa temos,
54 para cada transação, os candidatos que puderam ser utilizados na geração
55 das regras de associação, restando apenas verificar quais tem o suporte maior
56 que o suporte mínimo.

57 Para calcular o valor do suporte para todos os candidatos aplicamos a
58 transformação *reduceByKey*, com o operador de soma, e ao final obtemos
59 tuplas contendo como primeiro elemento o itemset e o segundo a quantidade
60 de transações em que tal itemset aparece.

61 Considerando o exemplo acima, suponha que ao final desta etapa, as tu-
62 plas tenham sido alteradas para $((\text{'1'}, \text{'2'}), 10), ((\text{'1'}, \text{'3'}), 8), ((\text{'2'}, \text{'3'}), 4)$, ou seja,
63 o itemset $(\text{'1'}, \text{'2'})$ apareceu em 10 transações, o itemset $(\text{'1'}, \text{'3'})$ apareceu em
64 8 e o itemset $(\text{'2'}, \text{'3'})$ apareceu em 4.

65 Tendo o número de transações em que cada itemset aparece, basta dividir
66 o mesmo pelo número total de transações e assim obter o valor de suporte
67 para cada itemset. Para isto é aplicado uma transformação *map* onde em cada
68 tupla, o segundo elemento é dividido pelo número total de transações, obtendo
69 assim o valor de suporte, e em seguida é aplicado uma transformação *filter*
70 onde são selecionados os itemsets que contém o valor de suporte maior que
71 o valor de suporte mínimo.

72 A cada iteração é feita a poda, ou seja, é descartado das transações itens
73 que não aparecem nos itemsets, evitando a geração desnecessária de can-
74 didatos.

75 A etapa de geração dos itemsets é encerrada quando não for possível obter
76 nenhum itemset para aquele determinado tamanho, restando apenas a etapa
77 de geração das regras de associação.

78 2.2. Geração das regras de associação

79 Ao final da primeira etapa temos os itemsets e seus respectivos valores de
80 suporte, nesta última etapa são geradas as regras de associação que tenham o
81 valor de confiança e lift maiores que os definidos pelo usuário.

82 É gerado um dicionário de dados com a transformação *collectAsMap*, ou
83 seja, através de um itemset podemos acessar o valor de seu suporte. Através
84 de uma transformação *map* são gerados os produtos cartesianos de todos
85 os itemsets, ou seja, são gerados tuplas de tamanho dois onde ambos os

86 elementos são itemsets distintos. Tais produtos cartesianos são as possíveis
87 regras de associações.

88 É aplicado posteriormente uma transformação filter onde é descartados
89 as tuplas em que existe algum elemento em comum entre os dois itemsets. A
90 tupla (X, Y) , é um exemplo de tupla gerada onde X e Y são os itemsets.

91 Posteriormente é necessário gerar um terceiro elemento $X \cup Y$, para cada
92 tupla que é utilizado para acessar o dicionário e calcular a confiança e o lift da
93 regra de associação $X \rightarrow Y$, assim é obtido no final deste processo uma tupla
94 com três elementos que no caso do exemplo citado acima seria, $(X, Y, X \cup Y)$.

95 Antes de calcular a confiança e o lift da regra de associação, é necessário
96 verificar se existe a chave $X \cup Y$ no dicionário, portanto é aplicado um trans-
97 formação filter onde é descartada todas as tuplas que não tenham o terceiro
98 elemento como chave do dicionário. Isto significa na prática que o itemset
99 $X \cup Y$ não tem um valor de suporte maior que o valor de suporte mínimo.

100 Posteriormente é feito o cálculo da confiança e lift utilizando os valores de
101 suporte do dicionário, que são acessados utilizando os itemset como chave,
102 restando apenas a aplicação de uma transformação filter para selecionar as
103 regras de associações que tiveram valores de confiança e lift maiores que os
104 definidos pelo usuário.

105 3. Resultados experimentais

106 Para os experimentos foi utilizado a base de dados da competição *In-*
107 *stacart Market Basket Analysis* do Kaggle, mais especificamente o arquivo
108 *order_products_prior.csv* [1]. Tal arquivo possui 3.214.875 de transações.

109 Em sua versão não distribuída o algoritmo foi executado totalmente em
110 15 minutos e 30 segundos. A sua versão distribuída para quatro threads foi
111 executada totalmente em 10 minutos, ou seja, a sua versão distribuída levou
112 cerca de 35 % a menos.

113 Como o esperado verificamos que o algoritmo, em sua versão distribuída,
114 obteve um ganho bastante interessante de tempo de processamento do que a
115 sua sua versão não distribuída.

117 **References**

- 118 [1] Competição Instacart Market Basket Analysis. [https://www.kaggle.](https://www.kaggle.com/c/instacart-market-basket-analysis/data)
119 [com/c/instacart-market-basket-analysis/data](https://www.kaggle.com/c/instacart-market-basket-analysis/data), 2018.
- 120 [2] Rakesh Agrawal and Ramakrishnan Srikant. Fast algorithms for mining
121 association rules in large databases. In *Proceedings of the 20th Interna-*
122 *tional Conference on Very Large Data Bases*, VLDB '94, pages 487–499,
123 San Francisco, CA, USA, 1994. Morgan Kaufmann Publishers Inc.
- 124 [3] Wikipedia. Apriori algorithm, [https://en.wikipedia.org/wiki/](https://en.wikipedia.org/wiki/Apriori_algorithm)
125 [Apriori_algorithm](https://en.wikipedia.org/wiki/Apriori_algorithm), 2018.