

# Ganho de Informação com Processamento Paralelo

Sergio Polimante Souto

*São Paulo, Brasil*

---

## Abstract

Ganho de informação (GI) mede o quanto de informação um atributo nos fornece sobre uma classe. Ou seja, atributos que dividem perfeitamente a classe fornecem informação máxima sobre a classe, e atributos que não se correlacionam com a classe não fornecem informação alguma. A medição de Ganho de Informação utiliza medidas de redução de Entropia. Essa técnica é utilizada para reduzir o conjunto de informação de um determinado banco de dados e minimizar os efeitos da maldição da dimensionalidade [1]. Sabendo o ganho de informação de cada atributo, é possível selecionar quais os melhores atributos para serem analisados sobre uma classe. Entretanto, esses cálculos apesar de simples, podem se tornar computacionalmente complicados dado a grande quantidade de informação de um banco de dados. Por isso, é importante a implementação de algoritmos que possibilitem o processamento em paralelo para otimizar o processamento do Ganho de Informação. Esse trabalho ilustra uma aplicação dessa técnica. Foi desenvolvido um algoritmo para computar o ganho de informação total dos atributos de um banco de dados. O algoritmo foi programado em Python utilizando a biblioteca pyspark. [2] [4]

*Keywords:* Ganho de Informação, Computação Paralela, pyspark

---

## 1. Algoritmos

O algoritmo para cálculo do ganho de informação é dividido em duas funções. A primeira, calcula a Entropia de um conjunto de dados categóricos, a segunda, calcula o Ganho de Informação de um atributo sobre uma classe. A segunda função utiliza a primeira.

### 1.1. Cálculo da Entropia de Shannon

O cálculo de Entropia de Shannon é feito através da equação 1 [2].

$$H = \sum_{i=0}^{N-1} -p_i \times \log_2 p_i \quad (1)$$

O algoritmo que realiza o cálculo está exibido abaixo:

---

```

8
9
10 def Entropia (classe):
11     """Calcula a Entropia de Shannon de uma distribuio de dados.
12     Args:
13         classe (RDD): RDD contendo conjunto de dados a ser calculado
14             a entropia.
15             Valores devem ser categoricos.
16     Returns:
17         float: valor de Entropia de Shannon calculado para o RDD.
18     """
19     #counts calcula paralelamente o contedo do RDD
20     #como tuplas contendo (tipo, quantidade)
21     counts = (classe.map(lambda x: (x, 1))
22               .reduceByKey(lambda a,b: a + b))
23     # n recebe o valor total de itens do RDD
24     n = classe.count()
25
26     # probs calcula a probabilidade de cada um dos estados do RDD
27     probs = counts.map(lambda x: x[1]/float(n))
28
29     # Entropia calcula a entropia do RDD
30     ## a funo map faz o calculo da Entropia de cada um dos estados
31     ## a funo reduce faz o somatrio da entropia de Shannon
32     entropia = (probs.map(lambda p: -p*math.log(p,2))
33                .reduce(lambda a,b: a + b))
34
35     # retorna valor escalar referente a entropia do RDD.
36     return entropia

```

---

### 1.2. Cálculo do Ganho de Informação

O ganho de informação é uma medida que expressa o quanto de informação um atributo fornece em relação a classe. Atributos que dividem perfeitamente a classe fornecem máximo ganho de informação (1), enquanto que

42 atributos que se relacionam aleatoriamente em relação a classe fornecem mí-  
 43 nimo ganho de informação (0). O ganho de informação é calculado utilizando  
 44 a equação 2 [2]:

$$GI = H - \sum_{i=0}^{N-1} \frac{n_i}{n} \times H_i \quad (2)$$

45 Onde H é a entropia calculada para a classe,  $n_i$  é a quantidade de atributos  
 46 de um determinado tipo de valor, N é o número de tipos diferentes de valores  
 47 que o atributo pode assumir, n é o total de dados do atributos e  $H_i$  é a  
 48 entropia calculada da classe para cada tipo específico de atributo. A classe  
 49 é filtrada e dividida em grupos de classe, um para cada tipo de valor do  
 50 atributo.

51 O código abaixo calcula o Ganho de Informação.

---

```

52
53 def infoGain (feature, classe, H):
54     """Calcula o ganho de informao de um atributo em relao a uma
55         classe.
56
57     Args:
58         feature (RDD): RDD contendo os conjuntos de dados do
59             atributo a ser
60                 calculado o Ganho de Informao
61
62         classe (RDD): RDD contendo conjunto de dados da classe
63
64         H (float): Entropia da Classe, previamente calculada.
65
66     Returns:
67         float: valor de ganho de informao ( reduo da Entropia) que
68             o atributo fornece sobre a classe
69     """
70     # calcula paralelamente o contedo do RDD
71     # como tuplas contendo (tipo, quantidade)
72     feat_count = feature.map(lambda x: (x, 1))\
73         .reduceByKey(lambda a,b: a + b)\
74         .collect()
75     # calcula as Entropias de um conjunto da classe dado cada um dos
76         estados do atributo
77     entropiasN = [Entropia(classe.zip(feature).filter(lambda x:
```

```

78         x[1]==v).map(lambda x: x[0])) for v,_ in feat_count]
79
80     # calcula a quantidade de itens no atributo
81     n = classe.count()
82
83     # calcula o ganho de informao do atributo.
84     ig = H - sum([(f[1]/float(n))*p for f,p in zip(feat_count,
85         entropiasN)])
86
87     return ig
88

```

---

## 89 2. Resultados

90 A base de dados escolhida é um banco de dados de atraso em voês co-  
91 merciais [3]. A classe escolhida é Vãos cancelados. O objetivo é encontrar  
92 qual dos atributos gera o maior ganho de informação.

93 A Entropia calculada para a classe é no valor de 0.1150, isso significa que  
94 os dados são fortemente puros. Existem 5729195 entradas com valor '0' (não  
95 cancelados), e 89884 entradas com o valor '1' (cancelados).

Tabela 1: Resultados

| -                                      | Mês       | Dia do<br>Mês | Dia da<br>Semana | cia. Aérea |
|--|-----------|---------------|------------------|------------|
| <b>Ganho de<br/>Informação</b>         | 0.0048013 | 0.0019261     | 0.0008652        | 0.0043092  |
| <b>Tempo com 8<br/>Partições (min)</b> | 37.64     | 94.51         | 21.79            | 43.14      |
| <b>Tempo com 1<br/>Partição (min)</b>  | 36.46     | 90.98         | 21.26            | 41.94      |

96 Verificou-se que é obtido o melhor desempenho com 1 partição para o  
97 hardware utilizado (i7 4500 2 cores). Também verificou-se que os atributos  
98 que fornecem maior ganho de informação para os voos cancelados são o mês  
99 em que o voo acontece a **companhia aérea** responsável. Ou seja, os voos  
100 cancelados possuem maior correlação com o mês em que o voo é realizado e  
101 a companhia aérea responsável pelo voo. Há baixa correlação com o dia do  
102 mês e muito baixa correlação com o dia da semana.

### 103 3. Agradecimentos

104 Agradeço ao Prof. Dr. Fabricio Olivetti por ministrar o curso e me  
105 encorajar a continua-lo até o fim. Agradeço ao Prof. Dr. Ronaldo Prati,  
106 meu coorientador no mestrado, pela ajuda na realização deste projeto.

### 107 4. Referências

- 108 [1] **Dimention Reduction** Acessado em 10/05/2018. Disponível em:  
109 [https://www.knime.com/blog/seven-techniques-for-data-dimensionality-](https://www.knime.com/blog/seven-techniques-for-data-dimensionality-reduction)  
110 [reduction](https://www.knime.com/blog/seven-techniques-for-data-dimensionality-reduction)
- 111 [2] **Information Gain** Acessado em 10/05/2018. Disponível em:  
112 <https://www.kaggle.com/usdot/flight-delays>
- 113 [3] **2015 Flight Delays and Cancellations** Acessado em 10/05/2018.  
114 Disponível em: <https://www.kaggle.com/usdot/flight-delays>
- 115 [4] **Information Gain** Acessado em 10/05/2018. Dis-  
116 ponível em: [https://courses.cs.washington.edu](https://courses.cs.washington.edu/courses/cse455/10au/notes/InfoGain.pdf) /[cour-](https://courses.cs.washington.edu/courses/cse455/10au/notes/InfoGain.pdf)  
117 [ses/cse455/10au/notes/InfoGain.pdf](https://courses.cs.washington.edu/courses/cse455/10au/notes/InfoGain.pdf)