



INSTITUTO DE INFORMÁTICA  
DEPARTAMENTO DE INFORMÁTICA APLICADA

Gabriel Fernandes e Henrique Peixoto

Universidade Federal do Rio Grande do Sul - UFRGS

Trabalho Final - INF01124 - Classificação e Pesquisa de Dados

Porto Alegre  
2021

# Sumário

1. Visão Geral
  - 1.1. Apresentação
  - 1.2. Explicação
  - 1.3. Justificativa
  - 1.4. Importância do problema
2. Implementação
  - 2.1. Estrutura do código
  - 2.2. Estruturas de dados
  - 2.3. Arquivos utilizados
  - 2.4. Algoritmos desenvolvidos
3. Guia de Uso
  - 3.1. Características
  - 3.2. Exemplos de uso
4. Contribuição
5. Elementos de terceiros
6. Considerações finais
7. Referências

# Visão Geral

## Apresentação

Este trabalho destina-se a demonstrar a aplicação, de organizada e prática, dos algoritmos aprendidos durante a disciplina de Classificação e Pesquisa de Dados. Para tanto, foram utilizadas estruturas que são condizentes com sistemas de busca de dados em grande escala, tais como árvores, tabelas hash e arquivos invertidos.

Este trabalho foi desenvolvido utilizando a linguagem de programação Python. Essa escolha foi feita por conta desta linguagem ser bastante versátil e oferecer uma grande facilidade para a implementação e utilização de estruturas de dados mais complexas. Tal feito seria mais complexo se tivéssemos utilizado a linguagem C, por exemplo.

Além disso, vale ressaltar que foram utilizados Git juntamente com GitHub para que fosse mantida a organização do código e para que se pudesse acompanhar o desenvolvimento do mesmo, de modo a preservar e deixar explícito as contribuições de cada membro da equipe envolvida no projeto.

## Explicação

O conjunto de dados utilizado refere-se a mais de 80.000 jogos presentes na plataforma online de venda de jogos Steam. Confira o conjunto de dados através deste [link](#). As informações de cada jogo contém dados tais como nome, preço, categorias, produtora(s), desenvolvedora(s), data de lançamento e descrição. Esses dados são um compilado de informações verdadeiras que podem ser encontradas na Steam, apesar de que alguns dados estejam defasados por conta de mudanças de preço nos jogos assim como a possível retirada de alguns títulos da plataforma.

O aplicativo desenvolvido tem por objetivo organizar esses dados nas estruturas de dados que foram selecionadas, pois chegou-se à conclusão de que essas estruturas seriam as mais adequadas para o tratamentos desses dados, assim como para sua busca, atualização, deleção e ordenação.

## Justificativa

Somente ter acesso a todos esses dados não implica, necessariamente, que sua utilização seja simples e cômoda, haja vista que muitas vezes esses dados estão desorganizados. E por conta disso foi desenvolvida esta aplicação: para que esses dados sejam organizados, através dos algoritmos utilizados, e para que esses dados possam ser utilizados de forma útil.

## Importância do Problema

Em consonância com o que foi exposto acima, organizar dados de forma coerente e utilizável vem tornando-se cada vez mais importante, note, por exemplo, o conceito de Big Data: *“Big data são dados com maior variedade que chegam em volumes crescentes e com velocidade cada vez maior”* (Gardner, 2001).

Além disso, estamos rodeados por informações e dados de toda e qualquer natureza que nos são úteis, tais como datas e horários, previsão do tempo, GPS, notícias, dentre outros, e todos esses dados chegam até nós de forma organizada de maneira que seja fácil compreender o que está sendo exibido. O que este trabalho faz é entender a importância dos dados e a maneira como eles devem ser tratados e apresentados.

# Implementação

## Arquitetura do Código

O código foi organizado através da utilização do padrão MVC (Models, Views, Controllers), é um padrão de arquitetura de software focado no reuso de código e a separação de conceitos em três camadas interconectadas, onde a apresentação dos dados e interação dos usuários (front-end) são separados dos métodos que interagem com o banco de dados (back-end).

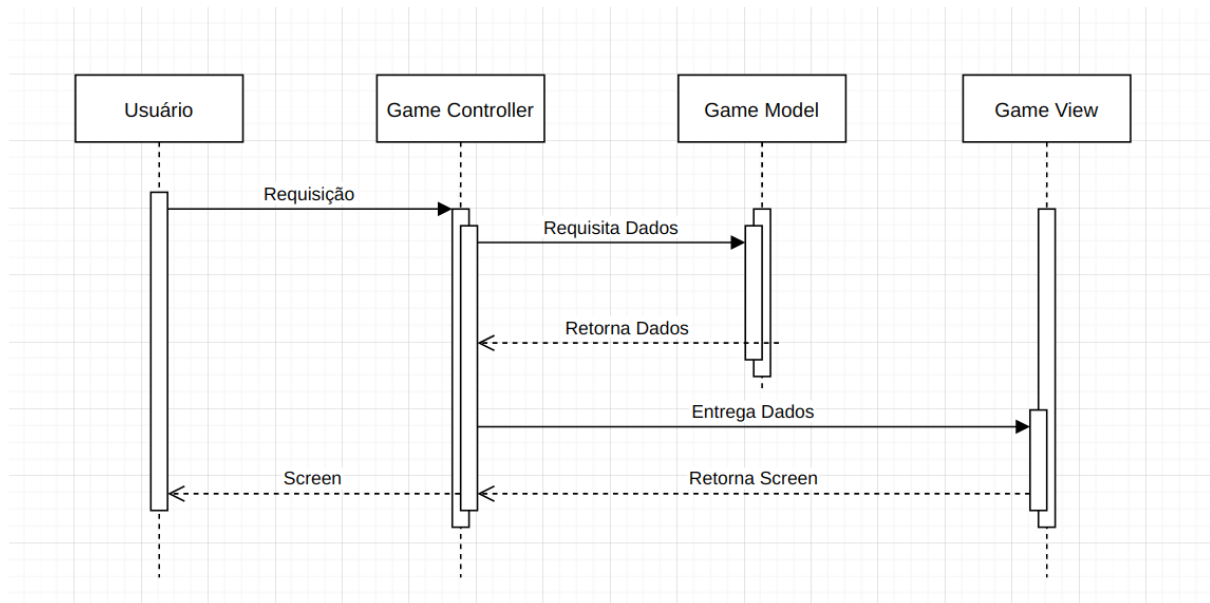


Figura 1: O padrão MVC.

Models, Views e Controllers são tipos de classes com responsabilidades definidas:

- As Classes do tipo Model são responsáveis pela interação com os dados, nestes projeto elas são as estruturas de dados criadas com os dados da Steam;
- As Classes do tipo Views são responsáveis por gerar uma screen dado um conjunto de dados pré-definido. Essa screen pode ser qualquer forma de interação fornecida para o usuário, seja ela gráfica, textual, sonora,...
- As Classes do tipo Controllers são responsáveis pelas ações que o usuário pode executar. São os métodos dos Controllers que, quando executados, usam as Models para pegar os dados, utilizam as Models para pegar os dados, entregam esses dados para uma View, e retorna a screen gerada para o usuário.

Na arquitetura MVC cada entidade da aplicação ganha o seu MVC, e dentro do controller da entidade são implementadas as ações de criação, leitura, atualização e deleção (CRUD) dos dados relativos àquela entidade, cada ação no seu próprio método.

O repositório do GitHub referente a aplicação pode ser encontrado através deste [link](#).

## Estruturas de Dados

As seguintes estruturas de dados foram utilizadas: árvore TRIE, arquivo invertido e conjuntos. O arquivo invertido é utilizado para fazer a consistência dos dados, já que uma vez extraídos devem ser formatados da maneira como se deseja e armazenados em outro arquivo. A árvore TRIE foi utilizada para armazenar atributos como o nome dos jogos, desenvolvedoras, categorias ou tags e colocar nas folhas, referentes a cada jogo, os ID dos registros onde esses nomes aparecem. Já os conjuntos foram utilizados como uma estrutura auxiliar para arquivo invertido. Esses detalhes de implementação serão melhor detalhados na subseção “Algoritmos Desenvolvidos”.

## Arquivos utilizados

Foram usados um arquivo, aquele do qual são extraídos os dados sobre os jogos, esse [link](#) dá acesso ao arquivo em questão e que também é utilizado para persistir os dados.

## Algoritmos desenvolvidos

Mostraremos, através de um exemplo, como é feita a extração e preparação dos dados, isso abrange o preenchimento dos conjuntos e da árvore TRIE bem como o salvamento dos dados. Os dados que serão usados neste exemplo estão presentes no conjunto de dados utilizado, por praticidade serão utilizados somente alguns registros do conjunto de dados e de maneira simplificada.

A primeira coisa a ser feita é a abertura do arquivo que contém os dados dos jogos. Esse arquivo não está organizado de modo legível, por isso é feito o parsing do mesmo no formato JSON (JavaScript Object Notation). Esse processo gera um array onde cada jogo é um dicionário (de acordo com a nomenclatura utilizada pela linguagem Python). Originalmente, esses dicionários não tem um campo “id”, então temos que percorrer este array e, para cada jogo, adicionar um campo “id”, que é determinado pela classe Entity. Suponha que o arquivo com os jogos tenha os seguintes dados:

```

1  [{
2      "name": "Counter-Strike: Global Offensive",
3      "price": "free",
4      "data": "Aug 21, 2012",
5      "developer": "Valve, Hidden Path Entertainment",
6      "publisher": "Valve"
7  },
8  {
9      "name": "Counter-Strike",
10     "price": "free",
11     "data": "Jan 28, 2007",
12     "developer": "Valve, Hidden Path Entertainment",
13     "publisher": "Valve"
14 },
15 [
16     "name": "Global Warming Simulator",
17     "price": "155",
18     "data": "Mar 15, 2013",
19     "developer": "Naughty Dog",
20     "publisher": "Naughty Dog"
21 ]]

```

a) Os dados sobre os jogos.

```

9  def create(self, item):
10     item = super().create(item)
11
12     for word in item["name"].strip().split():
13         word = word.lower()
14         if word not in self._names_table:
15             # cria conjunto (set) de itens
16             self._names_table[word] = {item["id"]}
17         else:
18             # add o item em conjunto existente
19             self._names_table[word].add(item["id"])
20     return item

```

b) Algoritmo de criação dos conjuntos (game.py).

Figura 2: Preparação dos dados sobre os jogos.

Então, é selecionado o nome de cada jogo, executa-se o parsing do nome, dividindo-o através de seus espaços em branco, por exemplo, o jogo “Metal Gear Solid: Phantom Pain” após o processo de parsing torna-se o array [“metal”, “gear”, “solid”, “phantom”, “pain”], as palavras também são transformadas para minúsculas, para que se mantenha um padrão, assim como quebras de linha e “dois pontos” também são retirados.

Dessa maneira vamos preenchendo os conjuntos, de modo que cada palavra, palavra essa que pertence ao nome de um jogo, é uma chave e o seu valor associado é um conjunto contendo todos os IDs dos jogos onde essa palavra aparece.

No momento em que inserimos um registro também o inserimos na TRIE, onde a folha, referente ao ID inserido, irá guardar os dados sobre o jogo associado ao ID em questão.

```
1  [{
2    "name": "Counter-Strike: Global Offensive",
3    "price": "free",
4    "data": "Aug 21, 2012",
5    "developer": "Valve, Hidden Path Entertainment",
6    "publisher": "Valve",
7    "id": "0"
8  },
9  {
10   "name": "Counter-Strike",
11   "price": "free",
12   "data": "Jan 28, 2007",
13   "developer": "Valve, Hidden Path Entertainment",
14   "publisher": "Valve",
15   "id": "1"
16 },
17 {
18   "name": "Global Warming Simulator",
19   "price": "155",
20   "data": "Mar 15, 2013",
21   "developer": "Naughty Dog",
22   "publisher": "Naughty Dog",
23   "id": "2"
24 }]
```

a) Dados sobre os jogos depois da adição dos IDs.

```
34   "counter-strike":{"0","1"},
35   "global":{"0","2"},
36   "offensive":{"0"},
37   "warming":{"2"},
38   "simulator":{"2"}
```

b) Criação dos conjuntos a partir dos nomes dos jogos.

Figura 3: Criação dos conjuntos.

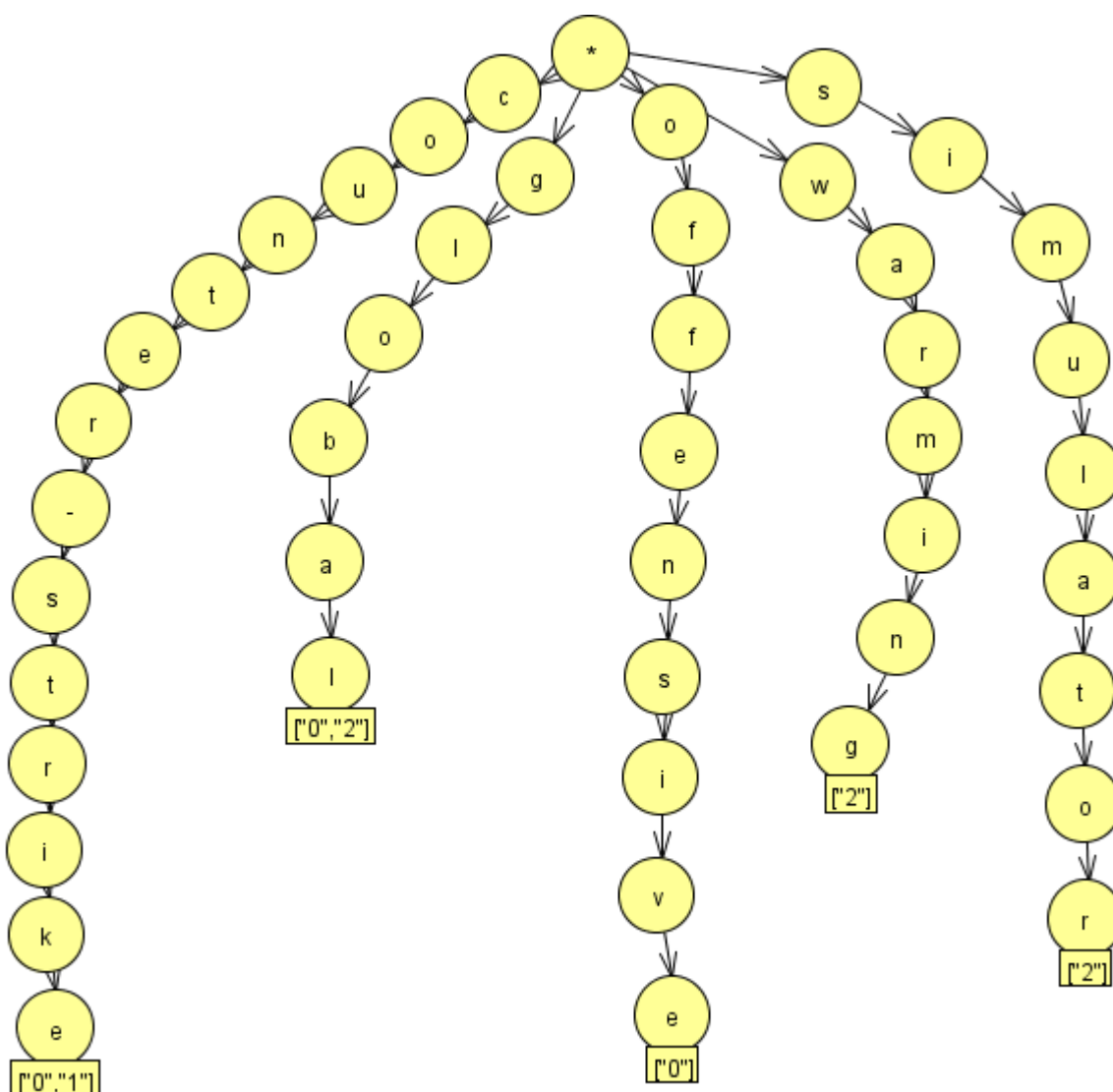
```

30 def __setitem__(self, word: str, id):
31     """
32     Método de Trie: adiciona uma palavra na TRIE.
33
34     Param: add(word, id)
35           word:str - palavra que será inserida
36           id:int - id da palavra
37     """
38     word = self._standardize_string(word)
39     node = self.root
40     for char in word:
41         found_in_child = False
42         # Procura pelo caractere nos filhos do nodo atual
43         for child in node.children:
44             if child.char == char:
45                 # Encontramos o caractere, incrementamos o contador pois
46                 # agora temos mais uma palavra que tem o mesmo prefixo
47                 child.counter += 1
48                 # E vamos para o nodo filho onde encontramos o caractere de 'word'
49                 node = child
50                 found_in_child = True
51                 break
52         # Se não achamos o caractere de 'word' temos que adicioná-lo
53         if not found_in_child:
54             new_node = TrieNode(char, id)
55             node.children.append(new_node)
56             node = new_node
57
58     # Quando terminarmos a palavra, marcamos o seu final
59     node.word_finished = True

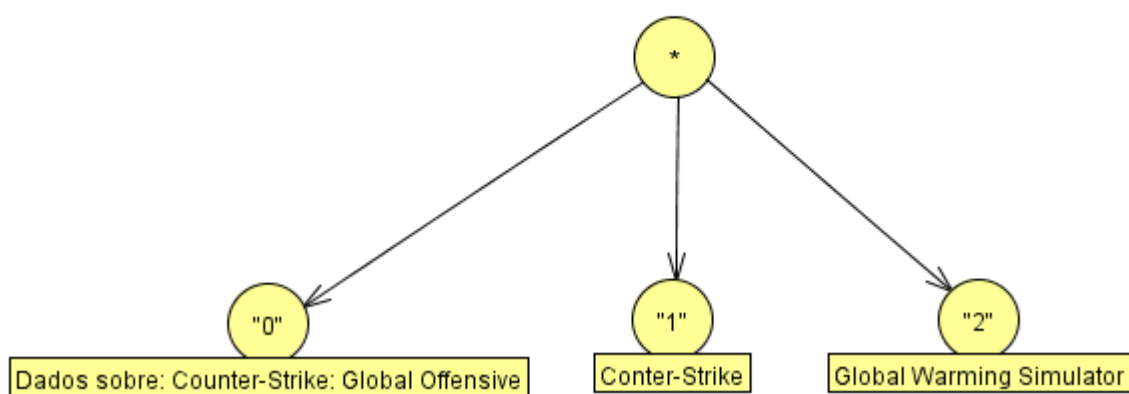
```

a) Algoritmo que implementa a TRIE (trie.py).





b) TRIE criada com base nos nomes dos jogos.



c) TRIE criada com base nos IDs dos jogos.

Figura 4: Criação da TRIE a partir dos conjuntos e algoritmo que implementa a TRIE.

Pela TRIE construída pode-se levar a pensar que não poderiam ter prefixos no meio de palavras ou que não ocorreriam ramificações na árvore. Abaixo são exibidas duas situações que a implementação dessa TRIE é capaz de suportar. Estes casos também podem ocorrer para outros elementos inseridos na TRIE (IDs, desenvolvedores, preços, dentre outros).

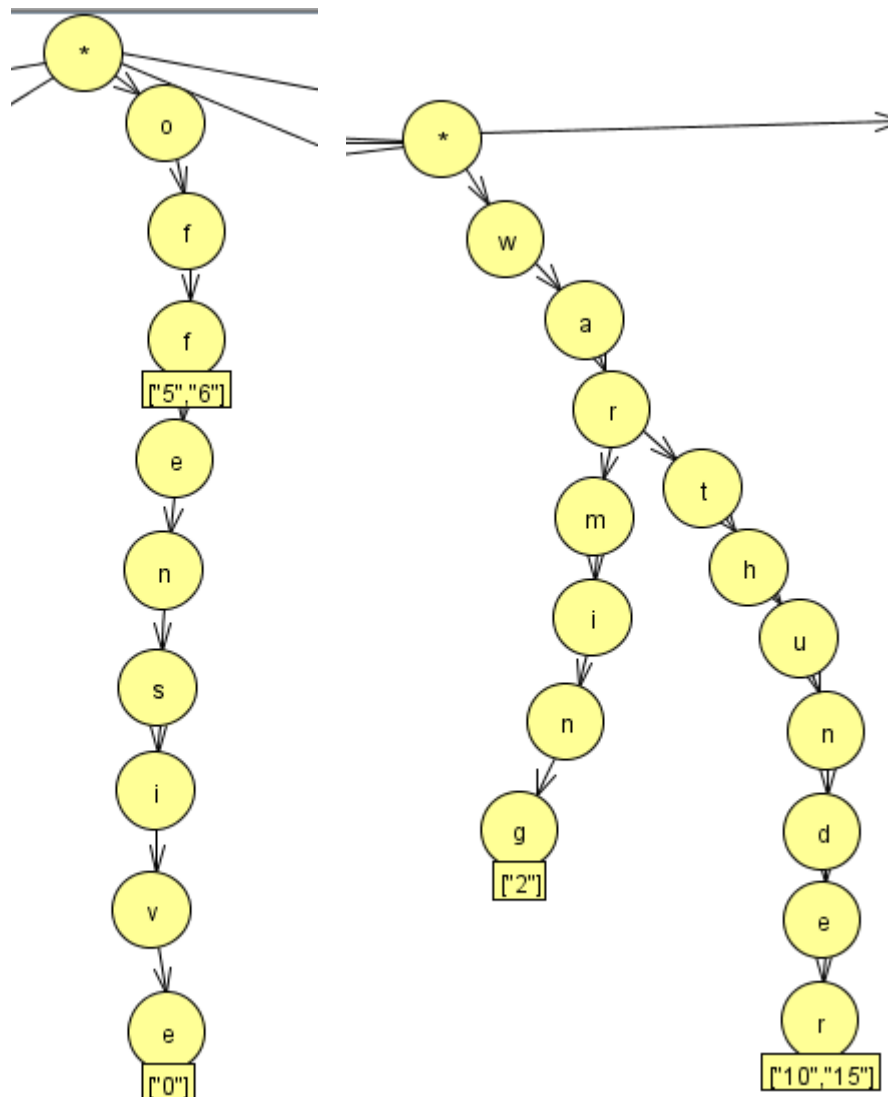


Figura 5: Outras possibilidades para a TRIE.

## Guia de Uso

### Características

### Exemplos de Uso

- Inicialização das estruturas: precisamos importar as classes Game e Entity. Após isso, abrimos o arquivo que contém os jogos e os colocamos nas estruturas utilizadas.

```
from steam.models.game import Game
from steam.models.entity import Entity
```

a) Importando as classes Game e Entity.

```
Games = Game()
```

b) Instanciamos um objeto do tipo Game.

```
file_content = open("archive/steam_data.json")
file = file_content.read()
file = json.loads(file)
file_content.close()

file_content = open("archive/steam_data.json", "w")

for game in file:
    game = Games.create(game)

json.dump(file, file_content)
file_content.close()
```

c) Abrimos o arquivo onde estão os jogos e os inserimos nas estruturas, após isso, salvamos as alterações.

Figura 6: A inicialização das estruturas com os dados dos arquivos.

- Inserção de dados: para inserir um jogo precisamos usar o método *create* da classe Game. Primeiro, precisamos importar a classe Game e a classe Entity, de onde Game herda alguns métodos..

```
from steam.models.game import Game
from steam.models.entity import Entity
```

a) A importação da classe Game e da classe Entity.

Depois, precisamos criar um jogo; para isso devemos escrever pelo menos um nome para o jogo, o ID será escolhido pelo método *create*. Você pode adicionar outros atributos se preferir, tais como preço, desenvolvedora, produtora, descrição do jogo, tags, categorias e data de publicação. Caso algum desses atributos não seja preenchido, eles serão preenchidos pelo método *create*.

```
Games = Game()
```

b) Instanciamos um objeto do tipo Game.

```

28 |     a_game = {
29 |         "name": "War simulator",
30 |         "price": "10",
31 |         "publisher": "Valve",
32 |         "developer": "Sony",
33 |         "full_desc": "Who cares? Just gimme a good grade!"
34 |     }

```

c) O jogo a ser inserido.

```
a_game = Games.create(a_game)
```

d) Adicionando um jogo, o retorno da função é o jogo criado já com o ID.

```

file_content = open("archive/steam_data.json")
file = file_content.read()
file = json.loads(file)
file_content.close()

for game in file:
    game = Games.create(game)

game1 = Game.create(game1)
game2 = Game.create(game2)
...
gamen = Game.create(gamen)

file_content = open("archive/steam_data.json", "w")

json.dump(file, file_content)
file_content.close()

```

e) Salvando o jogo criado no arquivo.

Figura 7: Adição de um jogo.

Note que *file* também contém os dados dos jogos que foram utilizados para inicializar as estruturas, portanto, o que estamos fazendo é sobrescrever os dados sobre os jogos. Podemos inserir tantos jogos quanto forem necessários.

- Ordenamento de dados: após inserir os jogos, podemos ordenar esses dados. No primeiro exemplo, iremos ordenar os jogos através de seus nomes em ordem decrescente.

```
Games.sort_by("name")
games_sorted = Games.paginate(1,80893)

for game in games_sorted:
    print(game["name"])
```

a) Ordenando todos os jogos pelo nome.

Figura 8: Ordenando os jogos pelo nome em ordem crescente.

Agora iremos ordenar os jogos pelo preço em ordem crescente.

```
Games.sort_by("price", True)
games_sorted = Games.paginate(1,80893)

for game in games_sorted:
    print(game["price"])
```

a) Ordenando os jogos pelo preço.

Figura 9: Ordenando todos os jogos pelo preço.

Quando o método `sort_by` não recebe o valor booleano, ele ordena os dados em ordem crescente, mas quando o valor booleano passado é `False`, os dados são ordenados em ordem decrescente.

## Elementos de Terceiros

O algoritmo de ordenamento utilizado pelo Python é o Timsort. Vejamos a sua definição: *“Timsort is a hybrid stable sorting algorithm, derived from merge sort and insertion sort, designed to perform well on many kinds of real-world data. It was implemented by Tim Peters in 2002 for use in the Python programming language. The algorithm finds subsequences of the data that are already ordered (runs) and uses them to sort the remainder more efficiently. This is done by merging runs until certain criteria are fulfilled”*<sup>[1]</sup>.

## Considerações Finais

A realização deste trabalho é a culminação de um processo de aprendizado que ocorreu ao longo de um semestre e que teve o objetivo de explicar algoritmos muito importantes para a computação como um todo. É necessário que se tenha conhecimento desses algoritmos para que se possa aplicá-los nas situações corretas. Não só isso, mas também é preciso que se saiba analisar um processo, no caso deste trabalho, classificar e ordenar dados, e utilizar as estruturas de dados e os algoritmos mais adequados possíveis para a resolução deste processo, seja esta resolução visando um melhor desempenho da aplicação em detrimento da quantidade de memória consumida ou o caso contrário.

Este trabalho também revelou aos seus participantes as maiores e mais prestigiadas habilidades que concernem o pensamento computacional, a saber: abstração, reconhecimento de padrões, decomposição de problemas e a criação de algoritmos.

Além disso, é muito satisfatório entender como as estruturas usadas neste trabalho, árvore TRIE, conjuntos e arquivos, trabalham juntas para compor a aplicação.

A maior dificuldade foi trabalhar com os arquivos, por isso escolhemos um formato um tanto flexível, JSON, que facilitou bastante a manipulação dos dados.

## Referências

[1] Timsort, Wikipedia. <https://en.wikipedia.org/wiki/Timsort>