



Comparação de Performance entre PostgreSQL e Cassandra NoSQL utilizando um Dataset de Vagas de Emprego do LinkedIn

Amanda de Medeiros Zepechouka, Juan C. F. R. R. de Moraes, Leandro de Lima Minchuel, Leonardo J. Zanotti, Victoria K. Vieira, Wellington H. Kania

Introdução

- LinkedIn job posts (Kaggle)

<https://www.kaggle.com/datasets/arshkon/linkedin-job-postings>

postings.csv (515.2 MB)



Detail Compact Column

10 of 28 columns

About this file

Add Suggestion

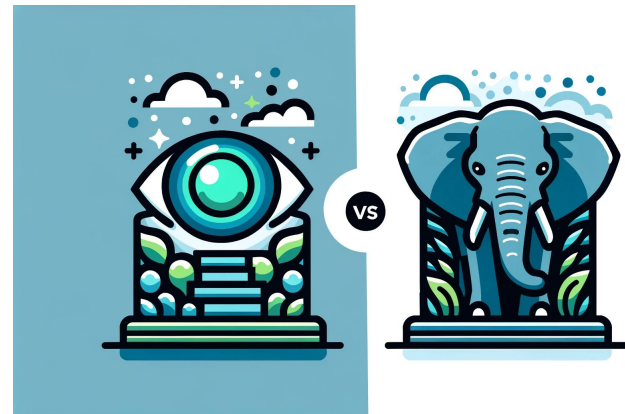
Main table. Supplemental files in folders (companies, job skill/industry tags, etc)

job_id	company_name	title	description	max_salary	pay_per
	<div><div>[null] 1%</div><div>Liberty Healthcare ... 1%</div><div>Other (121022) 98%</div></div>	<div>72521 unique values</div>	<div>107828 unique values</div>	<div></div>	<div>[null] YEARLY</div> <div>Other (15...</div>
922k	3.91b			1	120m
921716	Corcoran Sawyer Smith	Marketing Coordinator	Job descriptionA leading real estate firm in New Jersey is seeking an administrative Marketing Coord...	28.0	HOURLY
1829192		Mental Health Therapist/Counselor	At Aspen Therapy and Wellness , we are committed to serving clients with best practices to help them...	58.0	HOURLY



Tecnologias e comparação

- 01 | Velocidade
- 02 | Simplicidade (de instalação, configuração e escrita de código)
- 03 | Modelagem de dados
- 04 | Error report detalhado e facilidade de correção
- 05 | Abrangência/Especificação





Estudo de caso

1

Consulta Envolvendo Máximo

(Agrupamento): Qual é o maior salário presente no banco de dados?

2

Leitura simples por chave

primária: Qual a vaga com job_id igual a "2147609816"?

3

Atualização de dados por chave

primária: Atualize o título da vaga com job_id igual a "2974397965"

4

Filtragem por índice secundário:

Selecione todas as vagas que são de tempo integral



Resultados

- PostgreSQL com performance superior geral
- Cassandra se destaca nas consultas por chave primária

```
Lendo 123849 linhas do arquivo datasets/postings.csv
Dados importados para o PostgreSQL com sucesso. Tempo decorrido: 16.9958s
Keyspace e tabela criados com sucesso no Cassandra. Tempo decorrido: 0.8872s
Dados importados para o Cassandra com sucesso. Tempo decorrido: 62.0091s
```

```
Função agregada - Salário máximo
PostgreSQL - Tempo: 0.0255s, Resultados: 1
Cassandra - Tempo: 1.3536s, Resultados: 1
```

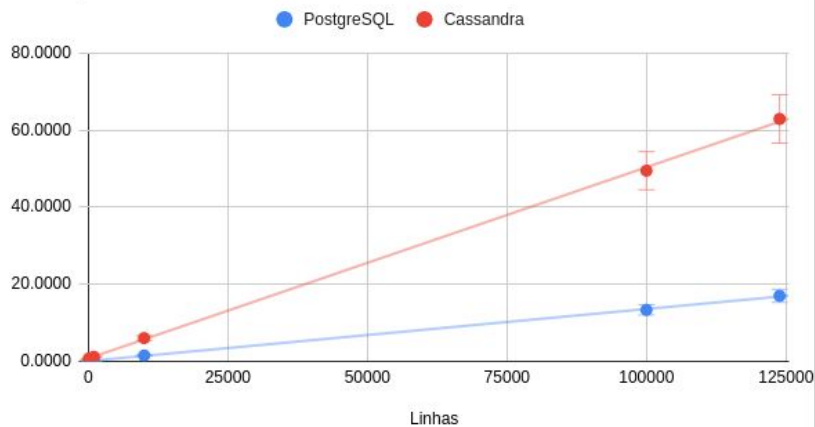
```
Leitura simples por chave primária
PostgreSQL - Tempo: 0.0229s, Resultados: 1
Cassandra - Tempo: 0.0201s, Resultados: 1
```

```
Atualização por chave primária
PostgreSQL - Tempo: 0.0214s, Resultados: 2
Cassandra - Tempo: 0.0168s, Resultados: 0
```

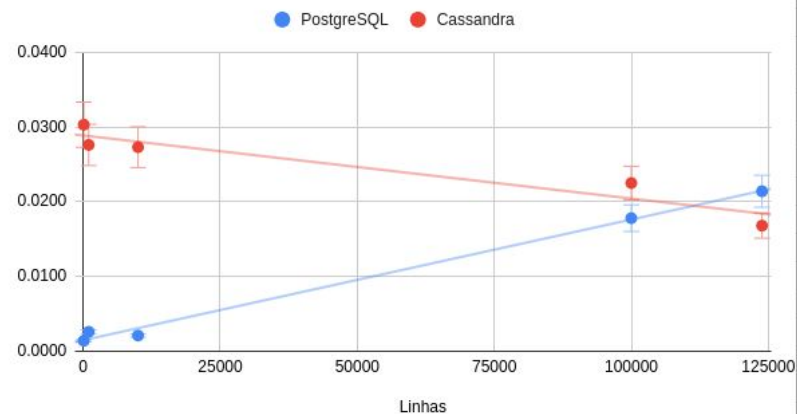
```
Filtragem por índice secundário
PostgreSQL - Tempo: 1.7680s, Resultados: 98814
Cassandra - Tempo: 6.2949s, Resultados: 98814
```

Regressão linear dos resultados por tempo

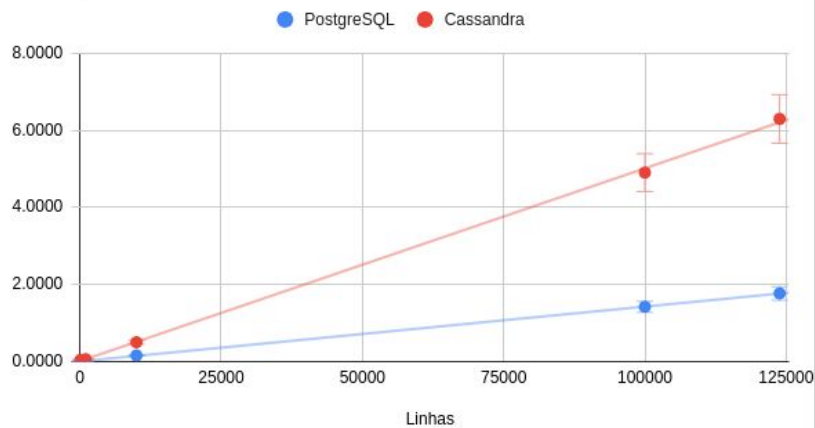
Inserção



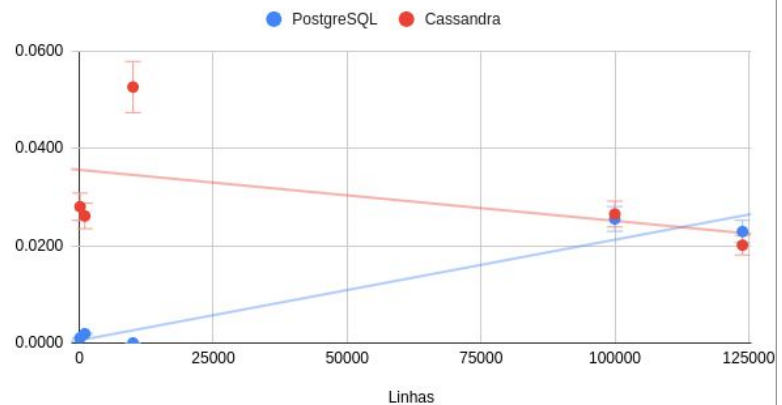
Atualização



Filtragem por índice secundário



Leitura simples





Simplicidade, abrangência e erros

- ❖ Match exatos
- ❖ Group By, JOIN, LIKE %word%
- ❖ Elasticsearch, Apache Spark
- ❖ Pré-agregação dos dados, contadores e sumários
- ❖ Atualizar uma coluna é inviável

2

Linhas de código para inserir
no PostgreSQL

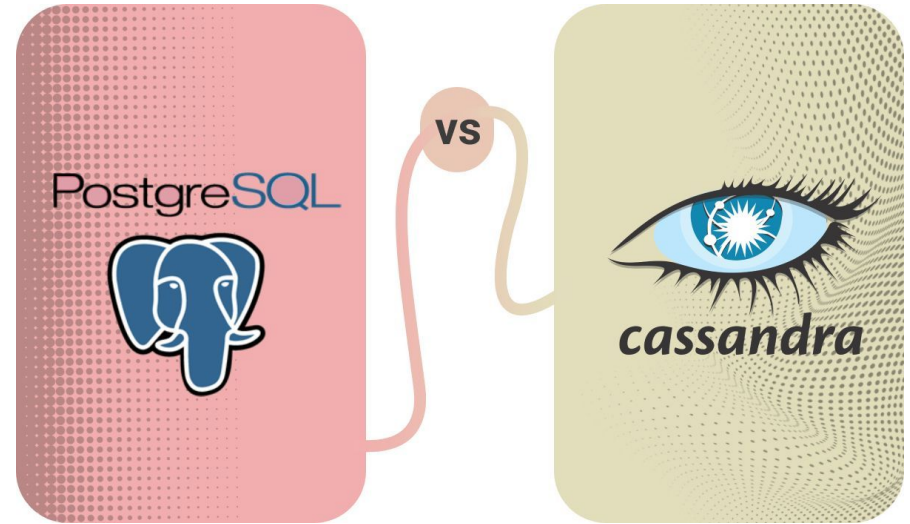
178

Linhas de código para
converter dados, remodelar e
inserir no Cassandra

```
Erro ao importar dados para o Cassandra: <Error  
from server: code=2000 [Syntax error in CQL qu  
ery] message="line 32:28 no viable alternative  
at input ',' (...,  
nsation_type  
                ) VALUES (  
                ["%]s",...)">
```

Conclusão

- Resultados foram os esperados
- Dataset maior (milhões de linhas)
- Usar outra modelagem de dados
- Usar outra ferramenta para fazer a comparação
- Padronizar o ambiente





Thank you.

“Pensaram que iam deixar o pai offline, mas não...” - Neymar Jr.

