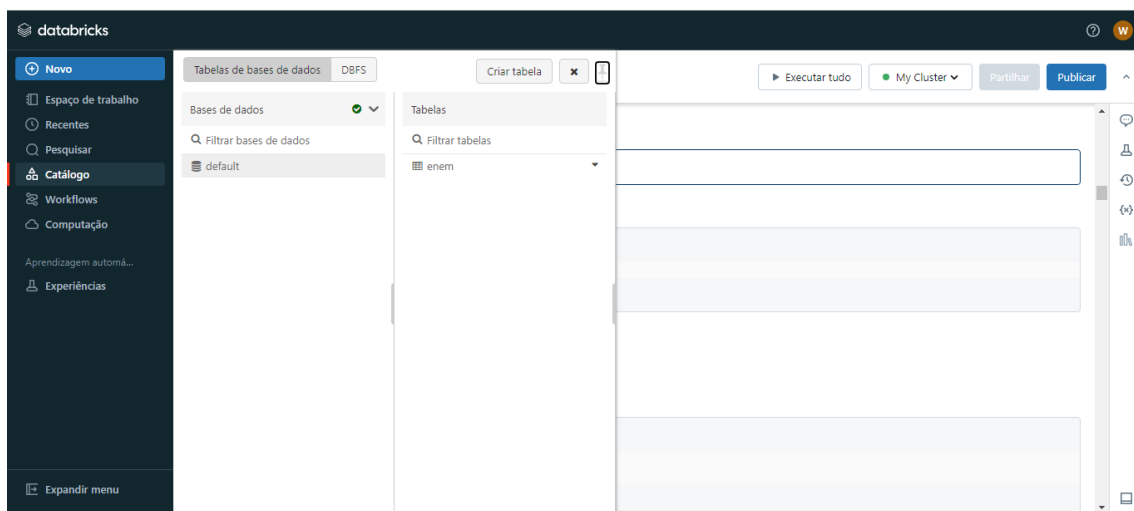


Criação da tabela física:



MVP Engenharia de Dados Python

Arquivo Editar Ver Executar Ajuda A última edição foi há 26 minutos

▶ Executar tudo My Cluster Partilhar Publicar

Abaixo é possível visualizar a tabela, objeto de pesquisa

```
%sql
SELECT * FROM enem LIMIT 10
```

▶ (1) trabalhos Spark

_sqldf: pyspark.sql.dataframe.DataFrame = [NU_ANO: long, CO_UF_ESCOLA: long ... mais 25 campos]

	NU_ANO	CO_UF_ESCOLA	SG_UF_ESCOLA	CO_MUNICIPIO_ESCOLA	NO_MUNICIPIO_ESCOLA	CO_ESCOLA_EDUCACENSO
2	2006	11	RO	1100205	Porto Velho	11000058
3	2005	11	RO	1100205	Porto Velho	11000058
4	2008	11	RO	1100205	Porto Velho	11000058
5	2007	11	RO	1100205	Porto Velho	11000171
6	2008	11	RO	1100205	Porto Velho	11000171
7	2005	11	RO	1100205	Porto Velho	11000171

MVP Engenharia de Dados Python

Arquivo Editar Ver Executar Ajuda A última edição foi há 26 minutos

▶ Executar tudo My Cluster Partilhar Publicar

10 linhas | 1,50 segundos de tempo de execução Atualizada há 2 horas

Este resultado é armazenado como um dataframe PySpark _sqldf e na cache de saída IPython como Out[4] . Saiba mais

Abaixo foi criado a tabela temporária, por onde executei as pesquisas.

```
# Criando uma tabela temporária
df.createOrReplaceTempView("tbl_enem")
```

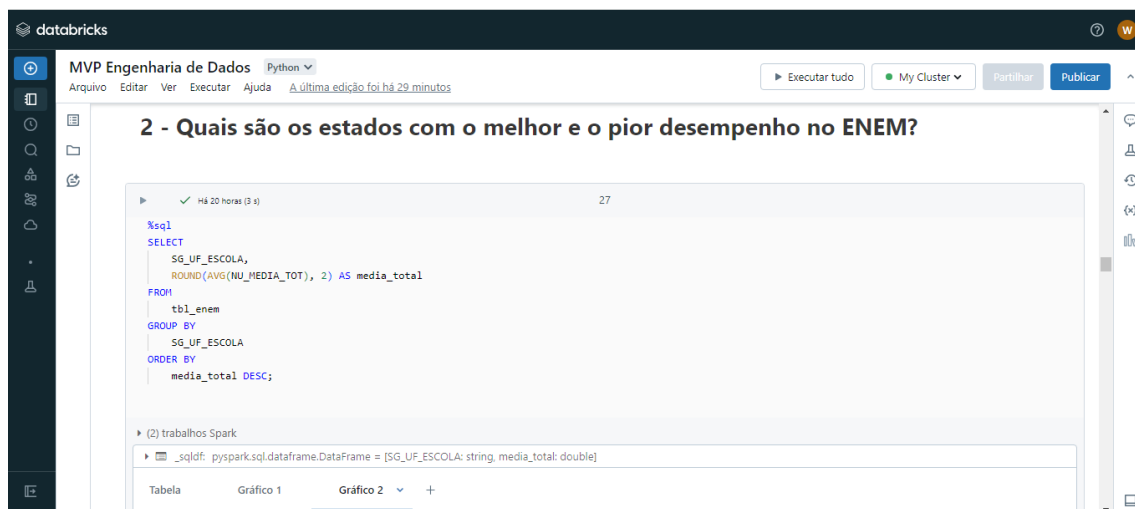
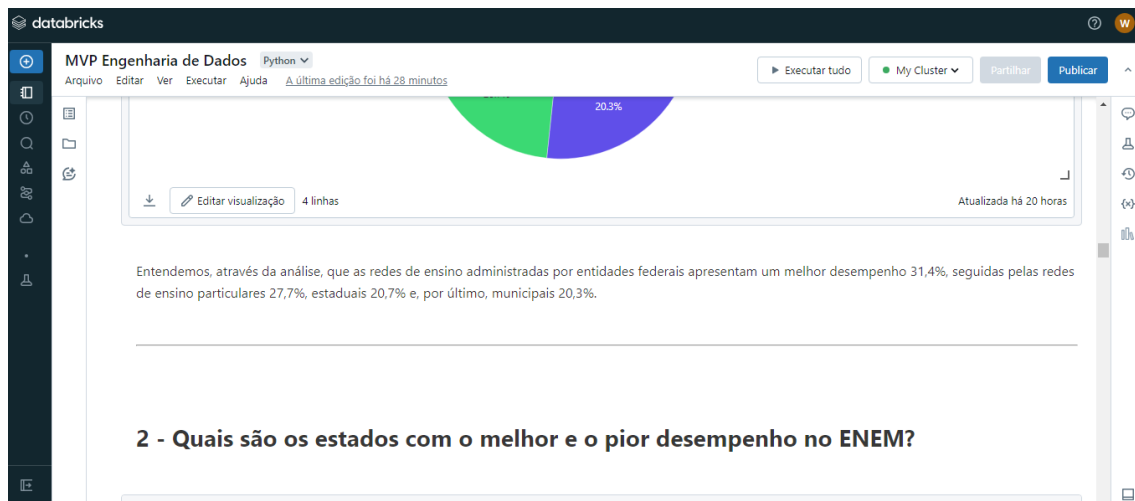
MVP Engenharia de Dados Python

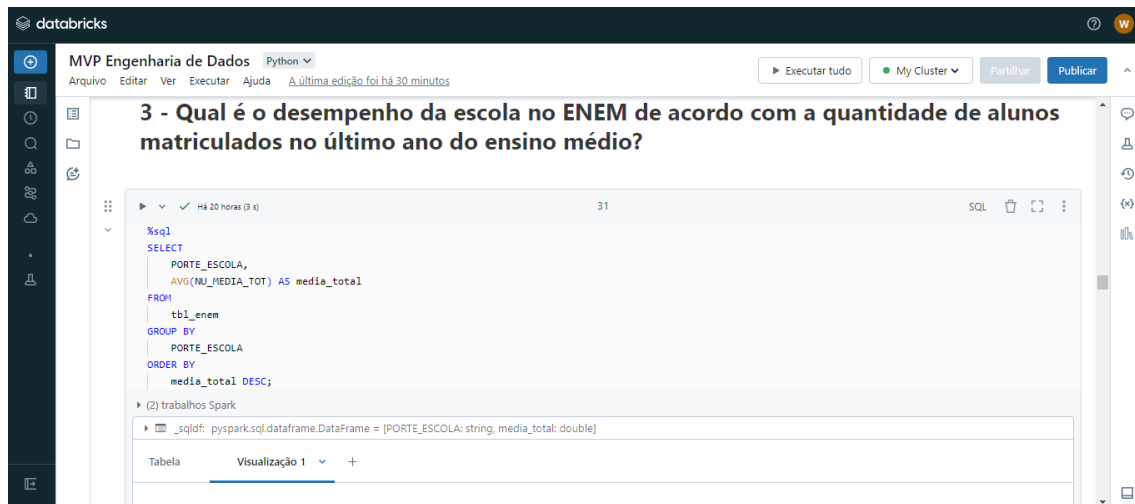
Arquivo Editar Ver Executar Ajuda A última edição foi há 27 minutos

▶ Executar tudo My Cluster Partilhar Publicar

Tentei normalizar os dados no código abaixo retirando as linhas nulas, no entanto alguns dos atributos não possuem dados históricos, como é o caso da média total no Enem que só possui informações nos anos de 2005, 2006 e 2007, então retirar as linhas nulas prejudicaria a conteúdo do dataset de forma geral.

```
%sql
1 DELETE FROM tbl_enem
2 WHERE NU_PARTICIPANTES_NEC_ESP IS NULL
3     OR NU_TAXA_PARTICIPACAO IS NULL
4     OR NU_MEDIA_CN IS NULL
5     OR NU_MEDIA_CH IS NULL
6     OR NU_MEDIA_LP IS NULL
7     OR NU_MEDIA_MT IS NULL
8     OR NU_MEDIA_RED IS NULL
9     OR NU_MEDIA_OB3 IS NULL
10    OR NU_MEDIA_TOT IS NULL
```





MVP Engenharia de Dados Python

Arquivo Editar Ver Executar Ajuda A última edição foi há 31 minutos

▶ Executar tudo My Cluster Partilhar Publicar

4 - Taxa de abandono de acordo com o porte da escola.

```
%sql
SELECT
  PORTE_ESCOLA,
  ROUND(AVG(NU_TAXA_ABANDONO), 2) AS media_taxa_abandono
FROM
  tb1_enem
WHERE
  NU_TAXA_ABANDONO IS NOT NULL
GROUP BY
  PORTE_ESCOLA
ORDER BY
  media_taxa_abandono ASC;
```

▶ (2) trabalhos Spark

MVP Engenharia de Dados Python

Arquivo Editar Ver Executar Ajuda A última edição foi há 31 minutos

▶ Executar tudo My Cluster Partilhar Publicar

▶ (2) trabalhos Spark

_sqldf: pyspark.sql.dataframe.DataFrame = [PORTE_ESCOLA: string, media_taxa_abandono: double]

Tabela	PORTE_ESCOLA	1.2 media_taxa_abandono
1	De 1 a 30 alunos	2.51
2	De 31 a 60 alunos	5.57
3	De 61 a 90 alunos	8.01
4	Maior que 90 alunos	10.5

4 linhas | 2,96 segundos de tempo de execução

Atualizada há 20 horas

Nesta análise a taxa de abandono segue uma proporcionalidade de acordo com a quantidade alunos matriculados no último ano do ensino médio por escola em suas turmas, quanto menos alunos, menor é a taxa de abandono.

databricks

🏠

📁

🔍

🔗

📊

📄

🔧

👤

MVP Engenharia de Dados Python

Arquivo

Editar

Ver

Executar

Ajuda

A última edição foi há 32 minutos

▶ Executar tudo

● My Cluster

Partilhar

Publicar

▶

✓

Há 20 horas (4 s)

39

```

WHERE
  NU_TAXA_PERMANENCIA IS NOT NULL
  AND NU_TAXA_APROVACAO IS NOT NULL
  AND NU_TAXA_REPROVACAO IS NOT NULL
  AND NU_TAXA_ABANDONO IS NOT NULL

GROUP BY
  PORTE_ESCOLA

ORDER BY
  PORTE_ESCOLA;

```

▶ (2) trabalhos Spark

▶ _sqldf: pyspark.sql.dataframe.DataFrame = [PORTE_ESCOLA: string, media_taxa_permanencia: double ... mais 3 campos]

Tabela

▼

+

🔍

🔗

📄

	A ₀ PORTE_ESCOLA	1.2 media_taxa_permanencia	1.2 media_taxa_aprovacao	1.2 media_taxa_reprovacao	1.2 media_taxa_abandono	
1	De 1 a 30 alunos	75.09	91.79	6.32	1.9	
2	De 31 a 60 alunos	77.74	88.44	8.21	3.36	
3	De 61 a 90 alunos	78.77	85.63	9.8	4.57	
4	Maior que 90 alunos	77.29	82.54	11.44	6.02	

The screenshot shows the Databricks web interface. At the top, there's a header with the Databricks logo and a navigation bar. Below the header, the main workspace area displays a table titled 'Tabela'. The table has 7 columns: 'PORTE_ESCOLA', '1.2 media_taxa_permanencia', '1.2 media_taxa_aprovacao', '1.2 media_taxa_reprovacao', '1.2 media_taxa_abandono', and an empty column. The table contains 4 rows of data. Below the table, there's a status bar indicating '4 linhas' and '4,34 segundos de tempo de execução'. The interface also shows a sidebar with navigation icons and a top bar with buttons like 'Executar tudo', 'My Cluster', 'Partilhar', and 'Publicar'.

	PORTE_ESCOLA	1.2 media_taxa_permanencia	1.2 media_taxa_aprovacao	1.2 media_taxa_reprovacao	1.2 media_taxa_abandono	
1	De 1 a 30 alunos	75.09	91.79	6.32	1.9	
2	De 31 a 60 alunos	77.74	88.44	8.21	3.36	
3	De 61 a 90 alunos	78.77	85.63	9.8	4.57	
4	Maior que 90 alunos	77.29	82.54	11.44	6.02	

4 linhas | 4,34 segundos de tempo de execução

Atualizada há 20 horas

MVP Engenharia de Dados Python

Arquivo Editar Ver Executar Ajuda A última edição foi há 34 minutos

Executar tudo My Cluster Partilhar Publicar

8 - Qual a taxa de aprovação no ensino médio de acordo com o estado?

```
%sql
SELECT
  SG_UF_ESCOLA,
  ROUND(AVG(NU_TAXA_APROVACAO), 2) AS media_taxa_aprovacao
FROM
  tbl_enem
WHERE
  NU_TAXA_APROVACAO IS NOT NULL
GROUP BY
  SG_UF_ESCOLA
ORDER BY
  media_taxa_aprovacao DESC;
```

(2) trabalhos Spark

_sqlidf: pyspark.sql.dataframe.DataFrame = [SG_UF_ESCOLA: string, media_taxa_aprovacao: double]

MVP Engenharia de Dados Python

Arquivo Editar Ver Executar Ajuda A última edição foi há 34 minutos

Executar tudo My Cluster Partilhar Publicar

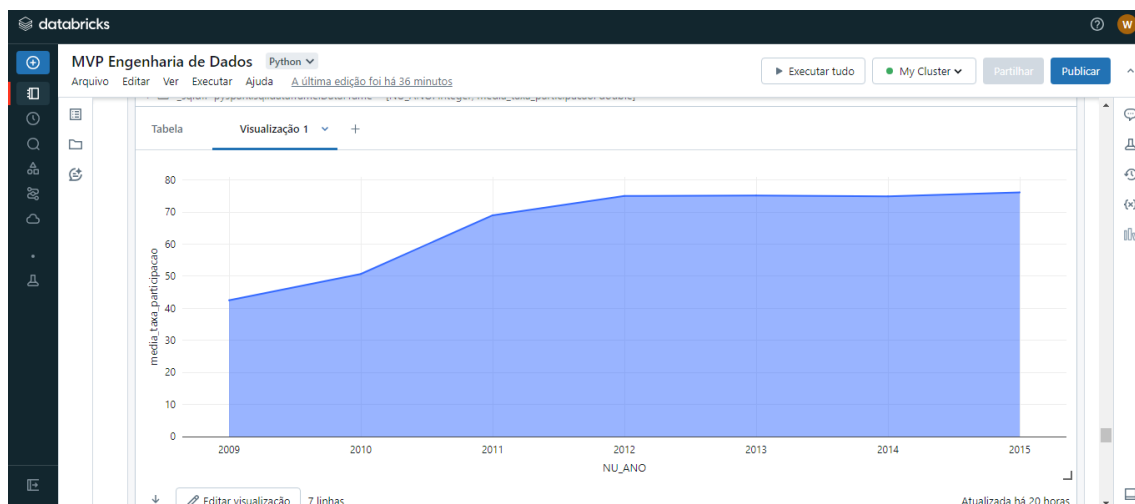
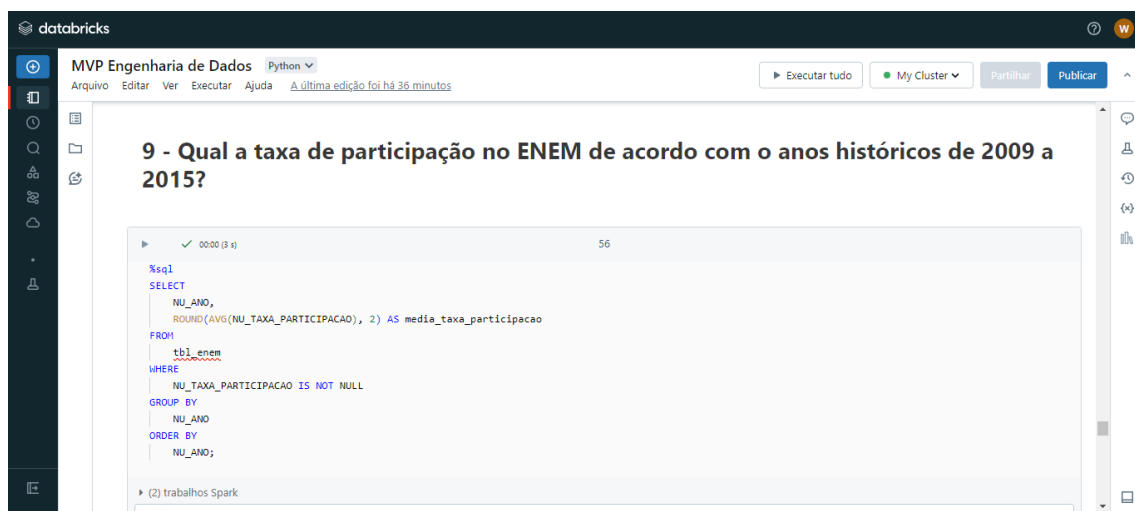
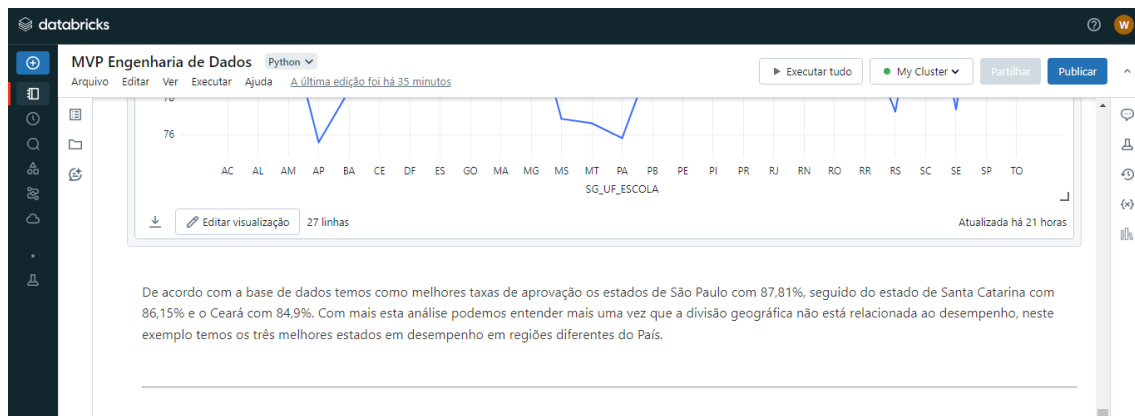
(2) trabalhos Spark

_sqlidf: pyspark.sql.dataframe.DataFrame = [SG_UF_ESCOLA: string, media_taxa_aprovacao: double]

Tabela Visualização 1

	SG_UF_ESCOLA	1.2 media_taxa_aprovacao
13	MA	81.52
14	RR	81.45
15	AL	80.32
16	PB	80.22
17	RJ	80.06
18	RO	79.01
19	PI	78.88
20	RN	78.54
21	BA	78.36
22	SE	77.35
23	RS	77.23
24	MS	76.85
25	MT	76.61







MVP Engenharia de Dados Python

Arquivo Editar Ver Executar Ajuda A última edição foi há 36 minutos

▶ Executar tudo My Cluster Partilhar Publicar

10 - Quais estados ficaram em primeiro lugar com maior taxa de aprovação, de acordo com o histórico dos anos presentes no estudo?

```
%sql
WITH taxa_por_estado AS (
  SELECT
    NU_ANO,
    SG_UF_ESCOLA,
    ROUND(AVG(NU_TAXA_APROVACAO), 2) AS media_taxa_aprovacao
  FROM
    tbl_enem
  WHERE
    NU_TAXA_APROVACAO IS NOT NULL
  GROUP BY
    NU_ANO,
    SG_UF_ESCOLA
),
ranked_taxa AS (
  SELECT
```

MVP Engenharia de Dados Python

Arquivo Editar Ver Executar Ajuda A última edição foi há 37 minutos

▶ Executar tudo My Cluster Partilhar Publicar

```
    NU_TAXA_APROVACAO IS NOT NULL
  GROUP BY
    NU_ANO,
    SG_UF_ESCOLA
),
ranked_taxa AS (
  SELECT
    NU_ANO,
    SG_UF_ESCOLA,
    media_taxa_aprovacao,
    ROW_NUMBER() OVER (PARTITION BY NU_ANO ORDER BY media_taxa_aprovacao DESC) AS rank
  FROM
    taxa_por_estado
)
SELECT
  NU_ANO,
  SG_UF_ESCOLA,
  media_taxa_aprovacao
FROM
  ranked_taxa
WHERE
  rank = 1
```