

Montando a Dimensão subscriptions

Tabelas: bi.subscriptions

Table name: bi.subscriptions

Alimentação diária por ano e mês

In [102].

```
#bi.subscriptions (premium)
df.subscriptions = df.subscriptions\
.withColumn("pt_processamento", current_timestamp())\
.withColumn("pt_mes_referencia", year("payment_date"))\
.withColumn("pt_mes_referencia", month("payment_date"))
```

In [103].

```
df.subscriptions.printSchema()

root
 |-- PaymentDate: string (nullable = true)
 |-- PaymentYear: string (nullable = true)
 |-- StudentId: string (nullable = false)
 |-- pt_processamento: timestamp (nullable = false)
 |-- pt_mes_referencia: integer (nullable = true)
 |-- pt_mes_referencia: integer (nullable = true)
```

In [104].

```
df.subscriptions.show()

-----+-----+-----+
|      PaymentDate|PlantYear|      StudentId|      ts_processamento|pt_mes_referencia|
|-----+-----+-----+
|2017-11-14 19:52:...|      Mensal|29037b0a52c5b576d...|2021-03-28 19:36:...|
|2017-11-08 11:52:...|      Mensal|b2bace77d15c3d3daf...|2021-03-28 19:36:...|
|2017-11-05 21:27:...|      Mensal|f42362de2f8964db6...|2021-03-28 19:36:...|
|2017-11-15 14:36:...|      Mensal|55cc0b5818a2eddbd5...|2021-03-28 19:36:...|
|2017-11-12 22:19:...|      Mensal|bb0cf63fe3e4e20cb...|2021-03-28 19:36:...|
|2017-11-24 19:03:...|      Mensal|69b7bee32821cf76b...|2021-03-28 19:36:...|
|2017-11-11 21:01:...|      Mensal|6553923125fe6364e...|2021-03-28 19:36:...|
|2017-11-12 16:41:...|      Mensal|3903a334d1af8ce83...|2021-03-28 19:36:...|
|2017-11-21 11:52:...|      Mensal|4487f81a4ea9b3c3c...|2021-03-28 19:36:...|
|2017-11-06 22:14:...|      Mensal|bd8436a92ab53ce2...|2021-03-28 19:36:...|
|2017-11-17 19:30:...|      Mensal|3ae0942281b9324e...|2021-03-28 19:36:...|
|2017-11-10 22:33:...|      Mensal|5a0d3781719b95d3c...|2021-03-28 19:36:...|
|2017-11-08 22:06:...|      Mensal|fed3eb368756e019f...|2021-03-28 19:36:...|
|2017-11-01 18:14:...|      Mensal|68bec3e289e331fe...|2021-03-28 19:36:...|
|2017-11-10 15:22:...|      Mensal|ecb4e2cd9585b080...|2021-03-28 19:36:...|
|2017-11-09 01:10:...|      Mensal|2e34dfda3c3ca387a...|2021-03-28 19:36:...|
|2017-11-08 00:41:...|      Mensal|85529f64c1428b69f...|2021-03-28 19:36:...|
|2017-11-13 12:20:...|      Mensal|ife9f61ee85396fbb...|2021-03-28 19:36:...|
-----+-----+-----+
only showing top 20 rows
```

Table name: bi.qtd_disc_premium

Quais disciplinas são mais acessadas por premium dentro do mês?

Alimentação diária por ano e mês

In [105].

```
#Table name: bi.qtd_disc_premium
df.qtdDiscPremium = df.subscriptions.alias("sub")\
.join(df.subjectsUser.alias("subUser"),
      trim(col("sub.StudentId")) == trim(col("subUser.student_id")),
      how = "right")\
.groupBy(\
  col("sub.StudentId"),
  col("subUser.subjects_name")\
)\
.withColumnRenamed("count", "qtd")\
.orderBy(desc("qtd"))\
.withColumn("ts_processamento", current_timestamp())\
.select(\
  col("subUser.subjects_name"),
  col("qtd"),
  col("ts_processamento"),
  col("payment_year_month")\
)
```

In [106].

```
df.qtdDiscPremium.printSchema()

root
 |-- subjects_id: long (nullable = true)
 |-- subjects_name: string (nullable = true)
 |-- qtd: long (nullable = false)
 |-- ts_processamento: timestamp (nullable = false)
 |-- payment_year_month: string (nullable = true)
```

In [107].

```
df.qtdDiscPremium.show()

-----+-----+-----+
|subjects_id|subjects_name|qtd|      ts_processamento|payment_year_month|
|-----+-----+-----+
|669660|Matemática Financeira|70|2021-03-28 19:38:...|2017-11|
|663660|Introdução à Administração|60|2021-03-28 19:38:...|2017-11|
|682100|Administração|48|2021-03-28 19:38:...|2017-11|
|692813|Contabilidade Básica|46|2021-03-28 19:38:...|2017-11|
|659639|Cálculo I|41|2021-03-28 19:38:...|2017-11|
|670592|Cálculo Diferencial|34|2021-03-28 19:38:...|2017-11|
|669404|Administração Financeira|35|2021-03-28 19:38:...|2017-11|
|669744|Cálculo II|35|2021-03-28 19:38:...|2017-11|
|669434|Resistência dos Materiais|34|2021-03-28 19:38:...|2017-11|
|670384|Mecânica Geral|34|2021-03-28 19:38:...|2017-11|
|675399|Metodologia Científica|33|2021-03-28 19:38:...|2017-11|
|670196|Didática|32|2021-03-28 19:38:...|2017-11|
|689481|Direito Civil I|32|2021-03-28 19:38:...|2017-11|
|671126|Direito Constitucional|31|2021-03-28 19:38:...|2017-11|
|671154|Bioquímica|31|2021-03-28 19:38:...|2017-11|
|670661|Cálculo Diferencial|31|2021-03-28 19:38:...|2017-11|
|669796|Física I|31|2021-03-28 19:38:...|2017-11|
|6119915|Planejamento de C...|31|2021-03-28 19:38:...|2017-11|
|672147|Contabilidade / C...|29|2021-03-28 19:38:...|2017-11|
|671761|Álgebra Linear|29|2021-03-28 19:38:...|2017-11|
-----+-----+-----+
only showing top 20 rows
```

Table name: bi.qtd_disc_npremium

#Quais disciplinas são mais acessadas por não premium dentro do mês?

Alimentação diária por ano e mês

In [136].

```
#bi.qtd_disc_npremium
df.qtdDisc = df.subscriptions.alias("sub")\
.join(df.subjectsUser.alias("subUser"),
      trim(col("sub.StudentId")) == trim(col("subUser.student_id")),
      how = "right")\
.filter("sub.StudentId is null")\
.groupBy(\
  col("subUser.subjects_id"),
  col("subUser.subjects_name")\
)\
.count()\
.withColumnRenamed("count", "qtd")\
.orderBy(desc("qtd"))\
.withColumn("ts_processamento", current_timestamp())\
.withColumn("pt_ano_mes_referencia", substr("ts_processamento", 1, 7))
```

In [137].

```
df.DiscGrat.printSchema()

root
 |-- subjects_id: long (nullable = true)
 |-- subjects_name: string (nullable = true)
 |-- qtd: long (nullable = false)
 |-- ts_processamento: timestamp (nullable = false)
 |-- pt_ano_mes_referencia: string (nullable = false)
```

In [138].

```
df.DiscGrat.show()

-----+-----+-----+
|subjects_id|subjects_name|qtd|      ts_processamento|pt_ano_mes_referencia|
|-----+-----+-----+
|671154|Bioquímica|4192|2021-03-28 21:45:...|2021-03|
|669639|Introdução à Administração|14408|2021-03-28 21:45:...|2021-03|
|670780|Anatomia Humana|3930|2021-03-28 21:45:...|2021-03|
|671126|Direito Constitucional|3899|2021-03-28 21:45:...|2021-03|
|669660|Matemática Financeira|3335|2021-03-28 21:45:...|2021-03|
|669796|Física I|2226|2021-03-28 21:45:...|2021-03|
|676833|Direito Penal|12928|2021-03-28 21:45:...|2021-03|
|670401|Fisiologia Humana|2884|2021-03-28 21:45:...|2021-03|
|676780|Direito Constitucional|2829|2021-03-28 21:45:...|2021-03|
|672228|Direito Processual|12556|2021-03-28 21:45:...|2021-03|
|669744|Cálculo II|2502|2021-03-28 21:45:...|2021-03|
|670592|Cálculo Diferencial|2416|2021-03-28 21:45:...|2021-03|
|669454|Resistência dos Materiais|2331|2021-03-28 21:45:...|2021-03|
|671086|Genética|2254|2021-03-28 21:45:...|2021-03|
|672547|Parcologia|2226|2021-03-28 21:45:...|2021-03|
|670384|Mecânica Geral|2113|2021-03-28 21:45:...|2021-03|
|675399|Metodologia Científica|2102|2021-03-28 21:45:...|2021-03|
|682100|Administração|2080|2021-03-28 21:45:...|2021-03|
|670819|Direito Administrativo|12058|2021-03-28 21:45:...|2021-03|
|689481|Direito Civil I|2055|2021-03-28 21:45:...|2021-03|
-----+-----+-----+
only showing top 20 rows
```

Load base B

Table name bi.events

Tabela alimentada diariamente de acordo com o recebimento dos logs

O que acessam na nossa plataforma e como fazem isso.

In [348].

```
#Table name bi.events
#Obs: Criado a coluna student_id. Utilizando a coluna studentId_clienteType, removendo
df.events = spark.read.format("json").load("files/passei/BASE-B/*-json")\
.withColumn("student_id", regexp_replace(col("studentId_clienteType"), "^(.*)@(.*)$",
      col("ts_processamento", current_timestamp())\
      col("pt_mes_referencia", month("at"))\
      col("pt_mes_referencia", month("at"))
```

In [349].

```
df.events.printSchema()

root
 |-- Last Accessed Url: string (nullable = true)
 |-- Page Category: string (nullable = true)
 |-- Page Category 1: string (nullable = true)
 |-- Page Category 2: string (nullable = true)
 |-- Page Category 3: string (nullable = true)
 |-- Page Name: string (nullable = true)
 |-- at: string (nullable = true)
 |-- browser: string (nullable = true)
 |-- carrier: string (nullable = true)
 |-- city: string (nullable = true)
 |-- city_total: long (nullable = true)
 |-- country: string (nullable = true)
 |-- custom_1: string (nullable = true)
 |-- custom_2: string (nullable = true)
 |-- custom_3: string (nullable = true)
 |-- custom_4: string (nullable = true)
 |-- device_new: boolean (nullable = true)
 |-- first-accessed-page: string (nullable = true)
 |-- install: uid: string (nullable = true)
 |-- language: string (nullable = true)
 |-- library_ver: string (nullable = true)
 |-- marketing_campaign: string (nullable = true)
 |-- marketing_medium: string (nullable = true)
 |-- marketing_source: string (nullable = true)
 |-- model: string (nullable = true)
 |-- name: string (nullable = true)
 |-- nch: long (nullable = true)
 |-- os: ver: string (nullable = true)
 |-- platform: string (nullable = true)
 |-- region: string (nullable = true)
 |-- session: uid: string (nullable = true)
 |-- studentId_clienteType: string (nullable = true)
 |-- type: string (nullable = true)
 |-- user_type: string (nullable = true)
 |-- uid: string (nullable = true)
 |-- students_id: string (nullable = true)
 |-- ts_processamento: timestamp (nullable = false)
 |-- pt_mes_referencia: integer (nullable = true)
 |-- pt_mes_referencia: integer (nullable = true)
```

In [350].

```
df.events.head()

Out [350]. Row(Last Accessed Url: '/', Page Category='perfil', Page Category 1='perfil', Page Category 2='Undefined', Page Category 3='Undefined', Page Name='/perfil/22482764/material_s...', at='2017-11-16 02:10:20', browser='Chrome 62', carrier='Telema Norte Leste S.a.', city_name=None, city_total=None, custom_1='core User', custom_2='Pedagogia', custom_3=None, custom_4='Core User', device_new=False, first-accessed-page=None, install: uid='dfdf303505f8a1b17ee40587875f6bb9c8374', language='pt', library_ver='w eb 3.3.3', marketing_campaign=None, marketing_medium=None, marketing_source=None, model='HMDA armv7l', name='Page View', stb='7', os_ver='7', platform='Linux', region=None, session: uid='188031bec37fc43b737c2c493497070Ae89128', studentId_clienteType='34cbef4428b798d94c9d49b43d4e249ce08d52364e09737db5864e4a848Website', type='web', user_type='known', uid='b3ed136094c0ae79f636ed1b03cf245311c8', students_id='34cbef4428b798d94c9d49b43d4e249ce08d52364e09737db5864e4a848', ts_processamento=datetime.datetime(2021, 3, 29, 1, 25, 29, 281000), pt_mes_referencia=2017, pt_mes_referencia=11)
```

Table name bi.qtd_acess_users_all

O que os usuários acessam na plataforma e como acessam?

Agrupado o campo "platform" com diversas versões de Windows para "Windows"

In [331].

```
#Table name: bi.qtd_acess_users_all
df.qtdAcessUser = df.events\
.from_pyspark.sql import functions as f\
df.qtdAcessUserAll = df.events\
.withColumn("platform",
      col("platform") like ("Windows"), "Windows")\
      col("platform")\
)\
.groupBy(\
  to_date("at"), alias("day_event"),
  "Page Name",
  "Page Category",
  "Page Category 1",
  "Page Category 2",
  "Page Category 3",
  "model",
  "platform")\
.count()\
.withColumnRenamed("count", "qtd")\
.orderBy(desc("qtd"))\
.withColumn("ts_processamento", current_timestamp())\
.withColumn("pt_ano_mes_referencia", substr("day_event", 1, 7))
```

In [332].

```
df.qtdAcessUserAll.show()

-----+-----+-----+
|day_event|Page Name|Page Category|Page Category 1|Page Category 2|Page Category 3|model|platform|qtd|      ts_processamento|pt_ano_mes_referencia|
|-----+-----+-----+
|2017-11-16|/|home|Undefined|Undefined|Undefined|Win32|Windows|51411|2021-03-29 01:12:...|2017-11|
|2017-11-16|/cadastro/passo1|cadastro|cadastro|cadastro|Undefined|Win32|Windows|18436|2021-03-29 01:12:...|2017-11|
|2017-11-16|/cadastro/passo2|cadastro|cadastro|cadastro|Undefined|Win32|Windows|1775|2021-03-29 01:12:...|2017-11|
|2017-11-16|/cadastro/premium|cadastro|cadastro|cadastro|Undefined|Win32|Windows|12454|2021-03-29 01:12:...|2017-11|
|2017-11-16|/listas|listas|listas|listas|Undefined|Win32|Windows|9166|2021-03-29 01:12:...|2017-11|
|2017-11-16|/|home|null|null|null|Win32|Windows|8047|2021-03-29 01:12:...|2017-11|
|2017-11-16|/premium|premium|premium|premium|null|Win32|Windows|6003|2021-03-29 01:12:...|2017-11|
|2017-11-16|/configuracoes/conta|configuracoes|configuracoes|configuracoes|Undefined|Win32|Windows|4727|2021-03-29 01:12:...|2017-11|
|2017-11-16|/|home|null|null|null|Linux armv7l|Android|3244|2021-03-29 01:12:...|2017-11|
|2017-11-16|/explorar/materiais/explorar-materiais/explorar-materiais|Undefined|Win32|Windows|3219|2021-03-29 01:12:...|2017-11|
|2017-11-16|/disciplina/seguindo|disciplina|disciplina|disciplina|Undefined|Win32|Windows|3086|2021-03-29 01:12:...|2017-11|
|2017-11-16|/pagamento/premiu...|pagamento|pagamento|pagamento|Undefined|Win32|Windows|2609|2021-03-29 01:12:...|2017-11|
|2017-11-16|/ranking/cursos|ranking|ranking|ranking|Undefined|Win32|Windows|2350|2021-03-29 01:12:...|2017-11|
|2017-11-16|/|home|Undefined|Undefined|Undefined|Win64|Windows|2266|2021-03-29 01:12:...|2017-11|
|2017-11-16|/listas/historico|listas|listas|listas|Undefined|Win32|Windows|2214|2021-03-29 01:12:...|2017-11|
|2017-11-16|/pagamento/premiu...|pagamento|pagamento|pagamento|Undefined|Win32|Windows|2118|2021-03-29 01:12:...|2017-11|
|2017-11-16|/cadastro/passos|cadastro|null|null|null|Linux armv7l|Android|1949|2021-03-29 01:12:...|2017-11|
|2017-11-16|/busca|busca|busca|busca|Undefined|Win32|Windows|1321|2021-03-29 01:12:...|2017-11|
|2017-11-16|/configuracoes/meu...|configuracoes|configuracoes|configuracoes|Undefined|Win32|Windows|1335|2021-03-29 01:12:...|2017-11|
|2017-11-16|/cadastro/passo2|cadastro|null|null|null|Linux armv7l|Android|1231|2021-03-29 01:12:...|2017-11|
-----+-----+-----+
only showing top 20 rows
```

In [255].

```
# Consulta: Após o agrupamento, temos as seguintes plataformas.
df.qtdAcessUserAll\
.groupBy("platform")\
.count()\
.show()

-----+-----+
|platform|count|
|-----+-----+
|IOS|1049|
|Linux|2490|
|Other|34|
|Fedora|38|
|Chrome OS|535|
|BlackBerry OS|2|
|Ubuntu|1051|
|Android|2482|
|Mac OS X|3925|
|Windows|185689|
-----+-----+
```

In [261].

```
df.events\
.groupBy("user_type")\
.count()\
.show()

-----+-----+
|user_type|count|
|-----+-----+
|anonymous|52676|
|known|64266|
-----+-----+
```

Table name bi.events_premium

O que os usuários premium mais acessam? De qual região? Qual Universidade e Curso?

In [366].

```
df.eventsPremium = df.eventsPremium\
.filter(col("user_type") == "known")\
.filter(f.col("students_id").isNotNull())\
.withColumn("platform",
      when(f.col("platform") like ("Windows"), "Windows")\
      otherwise(col("platform"))\
)\
.groupBy(\
  to_date("at"), alias("day_event"),
  "Page Name",
  "Page Category",
  "Page Category 1",
  "Page Category 2",
  "Page Category 3",
  "platform",
  "students_id")\
.count()\
.withColumnRenamed("count", "qtd")\
.orderBy(desc("qtd"))\
.withColumn("ts_processamento", current_timestamp())\
.withColumn("pt_ano_mes_referencia", substr("day_event", 1, 7))
```

In [376].

```
#O que os usuários premium mais acessam? De qual região? Qual Universidade e Curso?
#Table name: bi.events_premium
df.eventsPremium = df.eventsPremium.alias("eventsUsers")\
.join(df.cadastro.alias("cadastro"),
      trim(col("eventsUsers.student_id")) == trim(col("cadastro.student_id")),
      how = "inner")\
      trim(col("cadastro.student_id")) == trim(col("premium.StudentId")),
      how = "inner")\
)\
.select(\
  col("eventsUsers.day_event"),
  col("eventsUsers.Page Name").alias("page_categoria"),
  col("eventsUsers.Page Category").alias("page_categoria_1"),
  col("eventsUsers.Page Category 1").alias("page_categoria_2"),
  col("eventsUsers.Page Category 2").alias("page_categoria_3"),
  col("eventsUsers.Platform").alias("platform"),
  col("eventsUsers.students_id"),
  col("eventsUsers.students_state"),
  col("eventsUsers.universities_id"),
  col("eventsUsers.universities_name"),
  col("eventsUsers.courses_id"),
  col("eventsUsers.courses_name"),
  to_date(col("premium.PaymentDate"), alias("payment_date")),
  col("premium.day_event").alias("day_event"),
  col("eventsUsers.qtd").alias("qtd_acesso")\
).distinct()\
.orderBy(desc("qtd_acesso"))\
.groupBy("ts_processamento", current_timestamp())\
.withColumn("pt_ano_mes_referencia", substr("day_event", 1, 7))
```

In [377].

```
df.eventsPremium.printSchema()

root
 |-- day_event: date (nullable = true)
 |-- page_name: string (nullable = true)
 |-- page_categoria: string (nullable = true)
 |-- page_categoria_1: string (nullable = true)
 |-- page_categoria_2: string (nullable = true)
 |-- page_categoria_3: string (nullable = true)
 |-- platform: string (nullable = true)
 |-- students_id: string (nullable = true)
 |-- students_state: string (nullable = true)
 |-- universities_id: long (nullable = true)
 |-- universities_name: string (nullable = true)
 |-- courses_id: long (nullable = true)
 |-- courses_name: string (nullable = true)
 |-- payment_date: date (nullable = true)
 |-- plan_type: string (nullable = false)
 |-- qtd_acesso: long (nullable = false)
 |-- ts_processamento: timestamp (nullable = true)
 |-- pt_ano_mes_referencia: string (nullable = true)
```

In [378].

```
df.eventsPremium.show()

-----+-----+-----+
|day_event|page_name|page_categoria|page_categoria_1|page_categoria_2|page_categoria_3|platform|students_id|students_state|universities_name|courses_id|courses_name|payment_date|plan_type|qtd|
|-----+-----+-----+
|2017-11-16|disciplina/psico...|disciplina|disciplina|disciplina|disciplina|psicologia-instit...|null|null|disciplina|psicologia-instit...|null|null|2017-11| | | | | | | | | | | | | | | | | | |
|6646231|Undefined|Windows|4d2870e46a505511...|Pernambuco|Pedagogia|2017-11-13|Mensal|93|2021-03-29 10:36:...|
|6646231|ESTACIO|1199704|2017-11|disciplina|metodologia-cient...|São Paulo|null|2017-11|
|6646231|ESTACIO|1199517|2017-11|Direito|2017-11-06|Mensal|26|2021-03-29 10:36:...|
|2017-11-16|disciplina/gesta...|disciplina|disciplina|disciplina|gestao-da-qualidade|Pernambuco|2017-11-13|Mensal|13|2021-03-29 10:36:...|
|6646231|ESTACIO|1199704|2017-11|Pedagogia|2017-11-13|Mensal|22|2021-03-29 10:36:...|
|2017-11-16|disciplina/matem...|disciplina|disciplina|disciplina|matematica-financ...|disciplina|ensino-linico-i...|Ceará|null|2017-11|
|6646231|ESTACIO|1199529|2017-11|Enfermagem|2017-11-14|Mensal|14|2021-03-29 10:36:...|
|2017-11-16|disciplina/avali...|disciplina|disciplina|disciplina|avaliacao-institu...|disciplina|matematica-financ...|São Paulo|null|2017-11|
|6646231|ESTACIO|1199704|2017-11|Pedagogia|2017-11-13|Mensal|93|2021-03-29 10:36:...|
|2017-11-16|disciplina/psico...|disciplina|disciplina|disciplina|psicologia-instit...|Pernambuco|2017-11-13|Mensal|93|2021-03-29 10:36:...|
|2017-11-16|disciplina/direi...|disciplina|disciplina|disciplina|direito-penal-iii|disciplina|logica-de-program...|Goiás|Goiânia|2017-11-14|Mensal|1614530|UNINTER|1199704|2017-11|disciplina|mercado-financeiro|disciplina|matematica-financ...|São Paulo|Embuí|2017-11-08|Mensal|1614530|UNINTER|1199553|Engenharia de Pro...|2017-11-05|Mensal|12|2021-03-29 10:36:...|
|2017-11-16|disciplina/mode...|disciplina|disciplina|disciplina|modelagem-de-dados|disciplina|matematica-financ...|Goiás|null|2017-11-12|Mensal|10|2021-03-29 10:36:...|
|2017-11-16|disciplina/matem...|disciplina|disciplina|disciplina|matematica-para-n...|disciplina|direito-constituc...|São Paulo|null|2017-11-16|Mensal|93|2021-03-29 10:36:...|
|2017-11-16|disciplina/matem...|disciplina|disciplina|disciplina|matematica-financ...|disciplina|matematica-financ...|São Paulo|null|2017-11-14|Mensal|1614530|UNINTER|30724897|Gestão da Tecnolo...|2017-11-14|Mensal|7|2021-03-29 10:36:...|
|6646231|ESTACIO EAD|1199517|2017-11|Direito|2017-11-03|Mensal|93|2021-03-29 10:36:...|
-----+-----+-----+
only showing top 20 rows
```