

# MVP Engenharia de Dados – Felipe Teixeira Marques

## Introdução

Este relatório (MVP) apresenta os resultados de um projeto de Engenharia de Dados que visa analisar a equidade salarial entre homens e mulheres além da checagem do reconhecimento financeiro, através de méritos e promoções para os maiores talentos dentro da Bemobi. O objetivo principal é determinar se as posições nas faixas salariais ou PIR (Position in Range) entre os gêneros (masculino e feminino) estão sendo feitos de forma justa e identificar quaisquer disparidades salariais que possam existir.

## Objetivo

O MVP, visa duas perguntas que são levantadas pela diretoria:

- A distribuição de salários está realizada de forma justa entre gêneros? Em outras palavras, homens e mulheres possuem a mesma distribuição salarial quando comparamos seus salários e suas respectivas medianas?
- Nossos talentos, chamados de Key Talent, foram reconhecidos no último ciclo de mérito?

Para concluir esse objetivo, serão utilizados 3 arquivos em CSV, extraídos de controles da empresa onde temos informações relacionadas à:

- Base de colaboradores;
- Ciclo de mérito;
- Histórico de méritos e promoções anteriores à 2023;

Com isso será possível realizar os próximos passos da análise, que serão:

1. Avaliar e comparar os salários entre funcionários de diferentes gêneros na organização.
2. Identificar possíveis disparidades salariais que possam indicar falta de reconhecimento dos talentos com base no gênero.
3. Realizar análises estatísticas para determinar se as diferenças salariais são estatisticamente significativas.
4. Apresentar resultados claros e conclusões sobre a equidade salarial, focando especialmente no reconhecimento dos talentos dentro da empresa.
5. Propor recomendações, se necessário, para melhorar a equidade salarial e garantir que os talentos sejam reconhecidos de maneira justa.

## Coleta

A coleta de dados foi realizada na empresa Bemobi, e se deu, através de busca por arquivos disponíveis em sistemas internos da empresa.

Esses arquivos são gerados em formato XLSX, como existiam dados sensíveis em 2 das 3 tabelas utilizadas foi necessário realizar a remoção, garantindo assim a segurança e privacidade dos dados e por conseguinte dos funcionários da empresa.

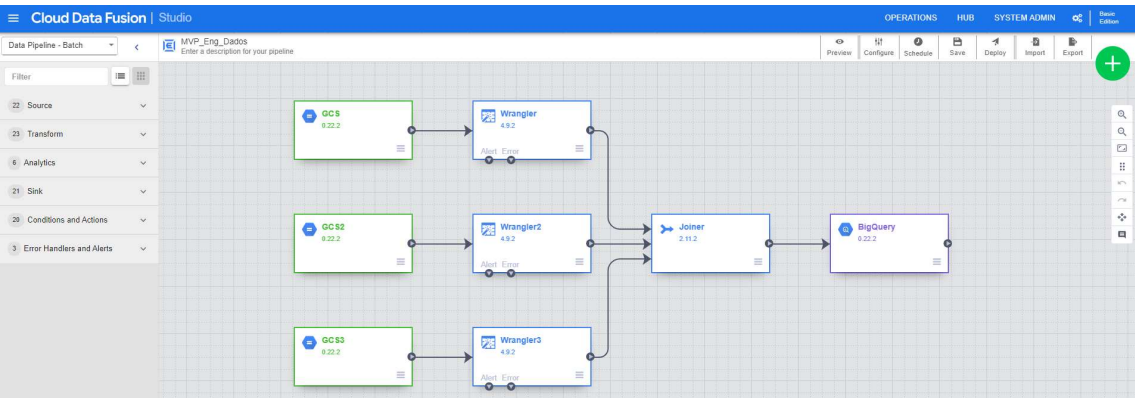
Após essa transformação inicial de dados, houve a conversão para o formato CSV, para que só então seguissem para o Google Cloud, plataforma em nuvem selecionada para o desenvolvimento das análises.

## Modelagem

O ETL foi realizado dentro do Google Cloud, inicialmente foi criado um Bucket, denominado de mvp-eng-dados, e dentro dele inseridos os 3 arquivos:

OBJETOS							
Intervalos > mvp-eng-de-dados							
FAZER UPLOAD DE ARQUIVOS   CARREGAR PASTA   CRIAR PASTA   TRANSFERIR DADOS   GERENCIAR RETENÇÕES   FAZER O DOWNLOAD   EXCLUIR							
Filtrar apenas pelo prefixo do nome   Filtro   Filtrar objetos e pastas   Mostrar dados excluídos							
<input type="checkbox"/>	Nome	Tamanho	Tipo	Criado	Classe de armazenamento	Última modificação	Acesso público
<input type="checkbox"/>	Colaboradores.csv	45,7 KB	text/csv	13 de set. de 2023 16:41:21	Standard	13 de set. de 2023 16:41:21	Não público
<input type="checkbox"/>	Historico_merito_promocao.csv	5,3 KB	text/csv	13 de set. de 2023 16:41:21	Standard	13 de set. de 2023 16:41:21	Não público
<input type="checkbox"/>	Movimentacao.csv	23,4 KB	text/csv	10 de set. de 2023 19:42:45	Standard	10 de set. de 2023 19:42:45	Não público

Em seguida os a realização da extração (extract) das informações na etapa 1 (Google Cloud Storage, ou GCS), dos dados da fonte que estavam no bucket.



As extrações foram realizadas em 3 GCS distintos:

1.

GCS Properties 0.22.2

Reads objects from a path in a Google Cloud Storage bucket.

Properties

Documentation

Label \*

GCS

Connection

Use Connection

no

Project ID

auto-detect

Service Account Type

File Path

JSON

Service Account File Path

auto-detect

Basic

Reference Name \*

Colaboradores

Path \*

gs://mvp-eng-de-dados/Colaboradores.csv

Output Schema

Actions

Metricula	int	+	-
Cargo	string	+	-
Seniority	string	+	-
Categoria	string	+	-
Grade	int	+	-
mediana	int	+	-
Salario	int	+	-
PIR	string	+	-
Diretoria	string	+	-
Area	string	+	-
Manager	string	+	-
PLR	double	+	-
Categoria_para_Promocao	string	+	-
Grade_Promocao	int	+	-
Nova_Mediana	string	+	-

2.

GCS Properties 0.22.2  
Reads objects from a path in a Google Cloud Storage bucket.

Properties Documentation

Label \*  
GCS2

Connection

Use Connection  
☐ NO

Project ID  
auto-detect

Service Account Type  
☒ File Path ☐ JSON

Service Account File Path  
auto-detect

Basic

Reference Name \*  
Historico

Path \*  
gs://mrp-eng-de-dados/historico\_merito\_promocao.csv

Output Schema

Field	Type	Actions
Matricula	int	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Data_Ultimo_reconhecimento	string	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Ultima_movimentacao_realizada	string	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Percentual_ajuste	double	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>

3.

GCS Properties 0.22.2  
Reads objects from a path in a Google Cloud Storage bucket.

Properties Documentation

Label \*  
GCS3

Connection

Use Connection  
☐ NO

Project ID  
auto-detect

Service Account Type  
☒ File Path ☐ JSON

Service Account File Path  
auto-detect

Basic

Reference Name \*  
Movimentacao

Path \*  
gs://mrp-eng-de-dados/Movimentacao.csv

Output Schema

Field	Type	Actions
COMPANY	string	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Diretoria	string	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Categoria	string	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
POSITION_GE_	string	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Cargo_Proposito_BP_	string	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
SENIORITY	string	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Manager	string	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Aprovador_Finial	string	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
REGISTRY	int	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
_da_Faixa	double	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Mediana	string	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
HIRING_AS_EMPLOYEE	string	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Apos	string	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
COMPENSATION	string	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Target_de_PLR	double	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
PLR_Estimada	string	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Encargos	string	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Custo_2023	string	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Data_Ultimo_reconhecimento	string	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>

Após essa etapa, o Wrangler é utilizado para realizar as transformações (transform), convertendo Medianas para Int, realizando substituições de "" e "FALSE" e "FALSO" por "vazio", no primeiro Wrangler.

Label \*  
Wrangler

Input Selection and Prefilters

Input field name \*  
\*

Precondition Language  
☒ JEXL ☐ SQL

Precondition (JEXL)  
false

Directives

Recipe

```
1 find-and-replace :Area s/"/"/g
2 find-and-replace :Nova_Mediana s/FALSE/g
3 find-and-replace :Nova_Mediana s/FALSO/g
4 set-type :Nova_Mediana long
```

Output Schema

Field	Type	Actions
Matricula	int	<input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Cargo	string	<input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Seniority	string	<input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Categoria	string	<input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Grade	int	<input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
mediana	int	<input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Salario	int	<input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
PIR	string	<input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Diretoria	string	<input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Area	string	<input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Manager	string	<input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
PLR	double	<input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Categoria_para_Promocao	string	<input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Grade_Promocao	int	<input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Nova_Mediana	long	<input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>

Além disso, foram realizadas diversas transformações no 3 Wrangler, segue a descrição abaixo:

Label \*

Wrangler3

Input Selection and Prefilters

Input field name \*

\*

Precondition Language

☒ JEXL ☐ SQL

Precondition (JEXL)

false

Directives

Recipe

1 drop :COMPANY

2 drop :Aprovador\_Final

3 rename \_de\_Faixa PIR

4 drop :Avos

5 drop :PLR\_Estimada\_2

6 drop :Novo\_Custo\_2023

7 drop :Sal\_P\_s\_Diss\_dio\_6

WRANGLE

User Defined Directives(UDD)

Output Schema

Actions

Diretoria	string	+	+
Categoria	string	+	+
POSITION_GE_	string	+	+
Cargo_Proposto_BPs_	string	+	+
SENIORITY	string	+	+
Manager	string	+	+
REGISTRY	int	+	+
PIR	string	+	+
Mediana	string	+	+
HIRING_AS_EMPLOYEE	string	+	+
COMPENSATION	string	+	+
Target_de_PLR	double	+	+
movimentacao_Proposta	string	+	+
Percentual_Proposto	string	+	+
Valor_do_Ajuste	string	+	+
Sal_rio_Proposto	string	+	+
Novo_Target_PLR	double	+	+

Depois de realizar essas transformações foi utilizado o Joiner Properties realizando um Inner Join:

Joiner Properties 2.11.2

Performs join operation on records from each input based on required inputs. If all the inputs are required inputs, Inner join will be performed. Otherwise inner join will be performed on required inputs and records from non-required inputs will only be present if they match join criteria. If there are no required inp...

Validar

Properties Documentation

Input Schema

Wrangler

Wrangler2

Wrangler3

Matrícula	int	+	+
Cargo	string	+	+
Seniority	string	+	+
Categoria	string	+	+
Grade	int	+	+
mediana	int	+	+
Salario	int	+	+
PIR	string	+	+
Diretoria	string	+	+
Area	string	+	+
Manager	string	+	+
PLR	double	+	+
Categoria_para_Promocao	string	+	+
Grade_Promocao	int	+	+
Nova_Mediana	long	+	+

Basic

Fields \*

Wrangler

Wrangler2

Wrangler3

Join Type

Inner

Join Condition Type

Basic

Join Condition

Wrangler

Wrangler2

Wrangler3

Matrícula

Matrícula

REGISTRY

Output Schema

Actions

Matrícula	int	+	+
Cargo	string	+	+
Seniority	string	+	+
Categoria	string	+	+
Grade	int	+	+
mediana	int	+	+
Salario	int	+	+
Diretoria	string	+	+
Area	string	+	+
Manager	string	+	+
PLR	double	+	+
Categoria_para_Promocao	string	+	+
Grade_Promocao	int	+	+
Nova_Mediana	long	+	+
Data_Ultimo_reconhecimento	string	+	+
Ultima_movimentacao_realiz	string	+	+
Percentual_ajuste	double	+	+
POSITION_GE_	string	+	+
Cargo_Proposto_BPs_	string	+	+
SENIORITY	string	+	+

E finalmente após essas etapas, o BigQuery é utilizado, para realizar o carregamento de dados (Load).

BigQuery Properties 0.22.2

This sink writes to a BigQuery table. BigQuery is Google's serverless, highly scalable, enterprise data warehouse. Data is first written to a temporary location on Google Cloud Storage, then loaded into BigQuery from there.

Validates

PropertiesDocumentation

Input Schema

Matricula	int	-	+	🔍
Cargo	string	-	+	🔍
Seniority	string	-	+	🔍
Categoria	string	-	+	🔍
Grade	int	-	+	🔍
mediana	int	-	+	🔍
Salario	int	-	+	🔍
Distortor	string	-	+	🔍
Area	string	-	+	🔍
Manager	string	-	+	🔍
PIR	double	-	+	🔍
Categoria_para_Promocao	string	-	+	🔍
Grade_Promocao	int	-	+	🔍
Nova_Mediana	long	-	+	🔍
Data_Ultimo_reconhecimento	string	-	+	🔍
Ultima_movimentacao_realizada	string	-	+	🔍
Percentual_ajuste	double	-	+	🔍
POSITION_GE_	string	-	+	🔍
Cargo_Proposito_BP's_	string	-	+	🔍

Label \*

BigQuery

Connection

Use connection

NO

Project ID

auto-detect

Dataset Project ID

Project the dataset belongs to, if different from the Project ID.

Service Account Type

☒ File Path ☐ JSON

Service Account File Path

auto-detect

Basic

Reference Name

Name used to identify this sink for lineage

BROWSE

Na próxima etapa o Deploy do fluxo foi realizado para que os dados possam ser analisados e responder os questionamentos iniciais do MVP.

Análise

Considerando os dados foi possível notar que das 73 promoções ou méritos realizados no último ciclo, 25% (18) ocorreram em pessoas do sexo feminino, o que em primeiro momento pode parecer que essas pessoas estão sendo prejudicadas, entretanto por se tratar de uma empresa de tecnologia, precisamos avaliar o quadro inicial, e nesse temos a seguinte distribuição por sexo:

Gênero	Valor Absoluto	Percentual
Masculino	285	74%
Feminino	99	26%
Total	384	100%

Levando essa tabela em consideração, a perspectiva muda pois se olharmos para cada grupo de forma individual, percebemos que 19% dos homens (pessoas do sexo masculino), receberam alguma promoção ou mérito, e 18% das mulheres (pessoas do sexo feminino), foram promovidas ou contempladas, com isso percebemos que os dois sexos foram considerados de forma semelhante no ciclo de mérito, uma vez que uma diferença de 1% não pode ser considerada grande suficiente para que uma ação seja necessária para gerar uma igualdade de gêneros da empresa.

Em relação ao PIR, notamos que as mulheres possuem uma média de **99%** em suas respectivas faixas, mostrando que de forma geral as mulheres estão sendo pagas de forma justa pelo seu trabalho e consideradas adequadas para suas funções, a média de PIR para homens é de **108%**, o que mostra, que de forma geral estão mais acelerados em suas respectivas faixas e para que ocorra um equilíbrio de gêneros, mostrando que em média homens estão sendo considerados como candidatos a promoções, já que para ser elegível a uma promoção a pessoa precisa estar acima dos 100% de PIR.

Embora a distribuição de promoções e méritos sejam justas uma informação que merece atenção é a distribuição dos talentos, conforme tabela abaixo:

Gênero	Talentos	Percentual
Masculino	37	71%
Feminino	15	28%
Total	52	100%

Ao cruzar os dois quadros notamos que 13% dos homens são considerados talentos e 15% das mulheres são consideradas como talento.

Portanto a empresa não possui um problema de reconhecimento de mulheres, o que atualmente ainda pode ser considerado um diferencial, especialmente por ser tratar de um mercado majoritariamente masculino.

### **Autoavaliação**

Ao longo do desenvolvimento foram encontradas algumas dificuldades, em geral, relacionadas a plataforma, e senti que não tive assistência através do Discord, por isso, precisei recorrer a colegas que atuam na área e a longas buscas na internet.

Apesar das dificuldades o trabalho foi concluído com sucesso, e para minha surpresa, o resultado obtido foi muito diferente do esperado.

Como citado na introdução foram utilizados aqui dados da empresa em que atuo, a Bemobi. E existe uma percepção de que mulheres não são reconhecidas dentro da empresa, inclusive com propostas de como resolver esse problema sendo levantadas por outras sub áreas dentro do RH.

Ao concluir a análise de dados, pude perceber que na verdade não existe um problema com o reconhecimento de mulheres na empresa, o que ocorre é que majoritariamente temos homens na empresa, tornando o espaço amostral de mulheres muito menor.

Sendo assim, o que podemos pensar é em como captar mais talentos femininos para a empresa, uma vez que não foi notada nenhuma diferença significativa nos percentuais de reconhecimento ou mesmo posicionamento na faixa entre homens e mulheres.