

Centro Universitário Carlos Drummond de Andrade

Marcos Dias Barbosa: RA_ 192001090

Alicya Danielle Holanda Silva: RA_ 999001025

Antony Machado Ferreira da Silva: RA_ 192000856

Gustavo Rodrigo Cruz Nunes: RA_ 192001161

Estruturação de dados: Adult-Sensus-Income

São Paulo - SP

2021

Centro Universitário Carlos Drummond de Andrade

Estruturação de dados: Adult-Sensus-Income

Relatório Técnico-Científico apresentado na disciplina de Projeto Integrador para o curso de (Análise e Desenvolvimento de Sistemas) do Centro Universitário Carlos Drummond de Andrade (UNI DRUMMOND).

São Paulo - SP

2021

DIAS, Marcos; DANIELLE, Alicya; MACHADO, Antony; RODRIGO, Gustavo. **Estruturação de dados: Adult-Sensus-Income**. 00f. Relatório Técnico-Científico. Análise e Desenvolvimento de Sistemas – **Centro Universitário Carlos Drummond de Andrade**. Tutor: (ME. Eduardo Palhares Júnior). Polo:(Ponte Rasa), 2021.

RESUMO

Nesta pesquisa, analisaremos uma estruturação de banco de dados: adult-sensus-income, no qual usamos como principal executor dos códigos, o Google Collaboraty (google colab), para que, através dele, possamos entender o funcionamento deste conjunto de dados. Usando como base a linguagem de programação Python, essa estruturação de banco de dados foi destrinchada para testarmos seus códigos em estruturas individuais, mas que irão se interligar como um todo, tais como: suas bibliotecas, discretizações, gêneros, dados de país, remoção de dados etc. Ao final, iremos ter uma visão clara de sua execução e seu resultado de um modo objetivo de acordo com seu funcionamento.

Palavras-Chave: Google Colab; Python; Adult-sensus-income.

LISTA DE ILUSTRAÇÕES

Figura 3.2.1: Bibliotecas.....	15
Figura 3.2.2: Discretização income	16
Figura 3.2.3: Análise das variáveis.....	17
Figura 3.2.4: Discretização das variáveis.....	18
Figura 3.2.5: Discretização das variáveis.	18
Figura 3.2.6: Discretização das variáveis..	19
Figura 3.2.7: Variável native country	20
Figura 3.2.8: Variável marital status.....	21
Figura 3.2.9: Variável marital status.....	21
Figura 3.2.10: Variável marital status	22
Figura 3.2.11: Variável relationship	23
Figura 3.2.12: Variável race	24
Figura 3.2.13: Variável work class	25
Figura 3.2.14: Variável work class	25
Figura 3.2.15: Variável capital Gain	26
Figura 3.2.16: Variável capital loss	27
Figura 3.2.17: Removendo variáveis não explicativas	28
Figura 3.2.18: Remoção	29
Figura 3.2.19: Variáveis não discretizados	30
Figura 3.2.20: Variáveis education num	31
Figura 3.2.21: Variável age	32
Figura 3.2.22: Analisando a correlação	33

Figura 3.2.23: Analisando a correlação	33
Figura 3.2.24: Importando a data split	34
Figura 3.2.25: Importando a data split	34
Figura 3.2.26: Importando os classificadores	35
Figura 3.2.27: Importando os classificadores	35
Figura 3.2.28: Código com erro	36

SUMÁRIO

1. INTRODUÇÃO	1
2. DESENVOLVIMENTO	2
2.1 Objetivos	2
2.2. Justificativa e delimitação do problema.....	2
2. 3. Fundamentação teórica	2
2.4. Aplicação das disciplinas estudadas no Projeto Integrador.....	2
2.5. Metodologia	2
3. Desenvolvimento e avaliação do trabalho.....	3
3.1. Apresentação do cenário.....	3
3.2. Bibliotecas.....	3
3.3. Discretização income.....	3
3.4. Análise das variáveis.....	3
3.5. Discretização das variáveis	3
3.6. Variável native country.....	3
3.7. Variável Marital-status.....	3
3.8. Variável Relationship.....	3
3.9. Variável Race.....	3
3.10. Variável Work Class.....	3
3.11. Variável Capital Gain.....	3
3.12. Variável Capital Loss.....	3
3.13. Removendo Variáveis não Explicativas.....	3
3.14. Remoção.....	3

3.15. Variáveis não Discretizadas.....	3
3.16. Variável Education Núm.....	3
3.17. Variável Age.....	3
3.18. Variável Analisando a correlação.....	3
3.19. Importando o Data Split.....	3
3.20. Importando os classificados.....	3
3.21. Código com erro.....	3
3.22. Solução Final.....	3
4. CONSIDERAÇÕES FINAIS	4
REFERÊNCIAS	5

1. INTRODUÇÃO

Neste trabalho iremos analisar uma base de dados, demonstrando os códigos com os quais usaremos para fazer buscas, discretizações de dados, remoções etc. Posteriormente tendo um resultado final desta análise, trabalhando para que o código não contenha linhas de erro e nos permita ter um código funcional.

Collab:

O termo **collab** é o diminutivo da palavra inglesa **collaboration**, que em português significa colaboração. **Collab** é uma prática utilizada por diversos segmentos com a intenção de unir nomes, personalidades, marcas e outros, para fornecer diferentes serviços ou produtos para um certo público.

O tema proposto foi introduzir uma base de dados, reorganizar esta base e fazer um estudo sobre ela e aplicá-la em um programa viável, no caso como já citado, o Google Collaboration. Ao aplicarmos essa base de dados em um programa, vamos agrupá-los em vários blocos de sequências de códigos de acordo com seus grupos para poder deixar organizado, de uma maneira que se possa entender como a coleta de dados foi realizada. A base de dados é baixada de um site chamado kaggle, em arquivo zipado, onde é extraído para uma pasta e assim poder ter acesso ao banco. A princípio, os dados foram colocados em arquivo Excel, que estava totalmente desorganizado, no qual fizemos uma edição para ser interpretado e, assim, ser submetido ao google collab. Nas informações básicas, podemos citar que, foram selecionadas diversas pessoas americanas com diferentes características sociais, econômicas, étnicas, rentáveis, estudantis etc.

Feito isso, importamos esta base de dados ao google drive, onde foi designado à plataforma pela qual usamos para codificar, usando a linguagem Python como referência e termos um retorno com isso às informações desejadas. Nesta plataforma, importamos as bibliotecas que seriam usadas para expressar os códigos. Nota-se aqui, que é de extrema importância que essas bibliotecas sejam anexadas no início do bloco de códigos para que se possam subir todas as devidas informações contidas no banco de dados. Feito isso, se tornou mais acessível a estrutura do código que, por final,

nos apresentou o resultado que era esperado. A ideia básica deste documento é reestruturar esses documentos em blocos nos quais se encaixam em seus devidos locais de relação de códigos. Sendo assim, cada informação deste dataset, irá ser reescrita nesses blocos que são divididos em parte. A pesquisa em si, se trata de um conjunto de pessoas que foram selecionadas para que fosse possível fazer levantamento de suas características socioeconômicas e demográficas sob os quais se obtém uma estimativa específica da população. Nesta pesquisa, os motivos foram particularmente baseados numa pesquisa feita por nosso professor, com a qual fizemos uma análise do conteúdo e tiramos algumas conclusões deste para dar origem ao documento aqui apresentado. O objetivo do trabalho, como já dito é coletar as informações propostas pelo autor e esboçar uma breve apresentação desses dados, que são as características de um determinado número de pessoas americanas selecionadas, para que obtivesse um resultado final de estimativa de quantas chegam a um valor próximo ou equivalente de 50k por ano. Então a coleta desses dados foi fundamental para que a base de dados e os cálculos usados em linhas de código em python pudessem nos dar todos os resultados mostrados em cada bloco de dados importados.

2. DESENVOLVIMENTO

2.1 Objetivos

Nosso objetivo é descrever nosso processo de montagem, identificando os erros preexistentes e traçando soluções para a resolução desses erros, descrever as estruturas utilizadas para análise deste código, posteriormente chegando a um resultado bem sucedido no final.

2.2. Justificativa e delimitação do problema

Durante a Produção dessas linhas de códigos, foram encontrados erros inesperados, onde nos impediram de ter um resultado positivo, nossa equipe se uniu para que juntos conseguíssemos encontrar soluções para que esse código pudesse ser efetivamente funcional. Todavia, nossos conhecimentos não foram totalmente eficazes para que solucionasse esses devidos erros. Por este motivo, acabamos por recorrer a outros métodos que fosse capaz de resolver essas incongruências. Um dos métodos mais eficazes que utilizamos foram algumas pesquisas que pudessem nos servir como apoio, e serviu, só que ainda assim nos faltou ferramentas solucionáveis para taparmos essa lacuna que ainda permanece.

2.3. Fundamentação teórica

Com as bases já coletadas, usamos como fontes confiáveis através do site [kaggle.com](https://www.kaggle.com), que é uma plataforma conhecida na área da tecnologia e é muito usada por desenvolvedores de diversos segmentos. Conhecida por ter a oportunidade de fazer competições de Machine Learning, explorar e publicar conjuntos de dados e também com possibilidades de treinamentos. Dito isso, o dataset do qual nos referimos foi retirado desta, podemos dizer uma plataforma conceituada na área tecnológica.

O arquivo baixado é compactado para que não seja tão dificultoso o download do mesmo para o estudo. Além dessa fonte utilizada, podemos contar com alguns documentos e vídeos aulas propostos por nosso professor responsável por este trabalho. Algumas questões que não encontramos no kaggle, os vídeos foram os responsáveis por nos ajudar a suprir dúvidas e questões não tão esclarecedoras. O arquivo, então, pode ser interpretado num compilador de códigos como o vscode, spider e até mesmo o sublime text. Mas por final, usamos mesmo o google colab, que não nos apresentou nenhuma dificuldade em rodar, importar, reestruturar e digitar os códigos.

2.4. Aplicação das disciplinas estudadas no Projeto Integrador

Algumas implementações usadas na construção desse projeto, foi o uso de algumas matérias que nos ajudou a ter uma síntese melhor e compreensível deste tema abordado. Olhando para o dataset, vemos que existe nela uma base desestruturada de dados coletados, fornecidos pelo kaggle. Nas aulas que usamos como acréscimo de informações.

Este item do referencial teórico deve indicar os conteúdos das disciplinas estudadas no curso que foram abordados no projeto. Espera-se que os estudantes relacionem, de forma clara e coerente, o conteúdo estudado à solução desenvolvida durante o projeto. Estruturação de dados baseado em python, programação orientada a objetos usando Java, são as disciplinas que podemos usar para incrementar ainda mais essa dataset. No entanto, como tivemos uma base de dados gravados em uma vídeo aula, foi de grande auxílio, já que o trabalho estava encaminhado. Mesmo tendo isso, observamos que a estruturação desses dados fora também estudada em aula, o que nos fez ter um nível de interpretação melhor.

2.5. Metodologia

Nesta metodologia, usamos ferramentas básicas para a execução de cada pedaço dos códigos. Além de instrumentos simples, mas bem funcionais, contamos com os conceitos e ensinamentos que foram adquiridos no decorrer do curso. Métodos como esse foram objetivos e claros para a realização. Desde o site específico para a coleta dos dados, desde sua construção na plataforma do google colabora. Para esse contexto realizado, tivemos que fazer divisões sobre o que cada participante iria fazer. Como obtemos o esboço de tratativas sobre a ideia do bando de dados, abordamos técnicas para desenvolver pequenos traços relevantes de expressar uma colocação clara do intuito deste trabalho. No mais, contextualizando, a integração de pessoas que tenham estimativas diferentes compõem este dataset.

Com algumas estratégias de diversificar cada espaço de funcionalidade de códigos, fizemos parte por parte, no sentido incremental, de cada bloco. A importação da biblioteca, como dito sendo necessária para a depuração e ligações de cada trecho para não haver erros, a discretização, a análise das variáveis, para se saber se todas estão corretamente declaradas para não ocorrer de algum trecho não rodar, os dados faltantes, os natives country, etc. foram pontos estratégicos para a delimitação correta dessa base de dados. Cada participante, tem um pedaço ideal na construção do projeto, uns mais expostos, outros mais reservados, mas que foram trabalhando em conjunto e mantendo as ligações sempre em dia.

Criar / Prototipar:

A análise de dados foi estudada em vídeo, e após esse estudo foram colocados em prática os conhecimentos de forma qualitativa, no qual há uma necessidade de qualidade na elaboração do projeto pesquisado em si. Em outra parte, na ideia de quantitativa, tivemos um arquivo demasiado grande para executar essa estruturação de dados. No entanto, usando as estratégias, foi possível reverter essa grande quantidade de informações e se enquadrar na organização da pesquisa. A

Descrição das soluções encontradas ou desenvolvidas para o problema investigado.

Com alguns problemas encontrados, sempre estávamos nós recorrendo ao conhecimento individual e coletivo para a resolução do mesmo. Não sendo possível resolver essa questão entre a

equipe, fomos recorrer ao vídeo disponibilizado e alguns websites voltados para a solução do problema.

Implementar / Testar:

A solução do problema que foi encontrado, foi testado na plataforma usada desde a implementação do banco de dados disponibilizado. Ao encontrar a possível solução, já fizemos a aplicação dela direto no código estruturado, pois por lá também já conseguimos executar e ter a certeza se irá ou não funcionar tal mudança no código. As devolutivas foram positivas, apesar de algumas dificuldades que tivemos em alguns códigos. Tendo auxílio de alguns vídeos no youtube, acessando diferentes blogs e páginas relacionadas ao dataset, foi possível converter os erros. Mesmo assim, a gama maior de acertos foi do vídeo aula que também nos foi deixava para ser usado como forma de acesso fácil de correção caso fosse necessário. E foi, então, para se concluir, esta base de dados foi útil na questão de se entender de que forma as pessoas são remuneradas nos EUA.

Que melhorias foram indicadas para as soluções propostas/desenvolvidas?

Como tínhamos já um projeto feito, com a necessidade apenas de digitar esses códigos e fazê-lo rodar de modo correto, viu-se que houve problemas de sintaxe e concatenação com o restante do código. Por aí, tivemos uma noção de que seria necessário aplicar algumas melhorias em certas linhas de código e até mesmo certas mudanças para que todo o conjunto rodasse adequadamente.

3. DESENVOLVIMENTO E AVALIAÇÃO DO TRABALHO

Este conjunto de dados é uma coleção de informações relacionadas a uma pessoa. A tarefa de previsão é prever se uma pessoa está ganhando um salário acima ou abaixo de \$50 mil. Exploramos extensivamente este conjunto de dados no intuito de analisar com o máximo de cuidado, as informações de cada habitante do país, checando assim se todos dentro dos parâmetros tem uma participação fundamental. Os dados foram pegos de pessoas com idade de 17 a 90 anos, incluindo a classe de trabalho, governo, peso final, educação, estado civil etc. Com isso esses conjuntos de dados é possível verificar o número de pessoas dentre as diferentes faixas etárias, sexo, ganho de capital, horas trabalhadas, país nativo dentre outras coisas.

```

3         'United-States', 'native.country'] = 'Non-US'
4     dataset.loc[dataset['native.country'] == \
5         'United-States', 'native.country'] = 'US'

[18] 1 native_country_map = {'US':1, 'Non-US':0}
      2 df['native.country'] = df['native.country'].map(native_country_map).astype(int)

[19] 1 print ("native.countryca classified by binary \n\n", df.head(10), "\n")
      native.countryca classified by binary
      age      workclass  fnlwtg  ...  hours.per.week  native.country  income
      1      82      Private  132870  ...      18          1          1
      3      54      Private  140359  ...      40          1          1
      4      41      Private  264663  ...      40          1          1
      5      34      Private  216864  ...      45          1          1
      6      38      Private  150601  ...      40          1          1
      7      74      State-gov  88638  ...      20          1          0
      8      68      Federal-gov  422813  ...      40          1          1
      10     45      Private  172274  ...      35          1          0
      11     38      Self-emp-not-inc  164526  ...      45          1          0
      12     52      Private  129177  ...      20          1          0

[10 rows x 15 columns]

Variável "marital-status" (matrimônio)

[20] 1 print(df[['marital.status', 'income']].groupby(['marital.status']).mean(), "\n")
      marital.status      income
      Divorced          0.892738
      Married-AF-spouse  0.523810
      Married-civ-spouse 0.545841
  
```

Começando pelo capital ganho por pessoa, podemos observar que o ganho dos \$50 mil são mais predominantes do que por pessoas que ganham menos que este valor. Já as pessoas mais velhas são prováveis que ganhem mais de \$50 mil do que os jovens.

```

1 print (df[['sex', 'income']].groupby(['sex']).mean(), "\n")
      income
      sex
      Female  0.890539
      Male    0.694263

[9] 1 sex_map = {'Male':1, 'Female':0}
      2 df['sex'] = df['sex'].map({'Male':1, 'Female':0})

[10] 1 print ("Sex classified by binary \n\n", df.head(10), "\n")
      Sex classified by binary
      age      workclass  fnlwtg  ...  hours.per.week  native.country  income
      1      82      Private  132870  ...      18          1          1
      3      54      Private  140359  ...      40          1          1
      4      41      Private  264663  ...      40          1          1
      5      34      Private  216864  ...      45          1          1
      6      38      Private  150601  ...      40          1          1
      7      74      State-gov  88638  ...      20          1          0
      8      68      Federal-gov  422813  ...      40          1          1
      10     45      Private  172274  ...      35          1          0
      11     38      Self-emp-not-inc  164526  ...      45          1          0
      12     52      Private  129177  ...      20          1          0
  
```

Na educação, podemos ver que as pessoas com diploma ganham mais que as pessoas que não possuem. Ex: apenas uma pequena proporção de pessoas com menos de 12 anos de educação ganha mais de \$50 mil por ano. Essa proporção aumenta quase que linearmente em 12 anos de escolaridade. Portanto, além de um diploma universitário, o que você pode fazer para aumentar suas

chances de ganhar mais de \$50 mil é maior. Também podemos observar que, pessoas que são casadas têm uma renda superior a \$50 mil. Mas somente 47 pessoas tinham um cônjuge. Sobre gênero e raça, os homens têm muito mais probabilidade de ganhar mais de \$50 mil anuais do que as mulheres. Mais trabalho pode ser feito usando esses dados para analisar as razões por trás dessa desigualdade de renda. Como avaliação, podemos dizer que os conjuntos de dados foram favoráveis para se distinguir a questão de idade, raça, gênero, educação e pais, como uma das formas de se avaliar uma pessoa no país.

3.1. APRESENTAÇÃO DO CENÁRIO

O que é o Google Colab?

Google Colab, ou Google Collaboratory, é um serviço de armazenamento em nuvem de notebooks voltados à criação e execução de códigos em Python, diretamente em um navegador, sem a necessidade de nenhum tipo de instalação de software em uma máquina. Em outras palavras, com o Google Colab você é capaz de ler, desenvolver e rodar códigos e rich texts em documentos interativos que agrupam células de códigos — chamados de notebooks —, compartilhá-los com outros programadores, modificá-los a quaisquer momentos e mantê-los salvos de maneira totalmente online. Todo o poder computacional utilizado para executar o software que você escreve é fornecido pela nuvem de computadores da Google. Dando gratuitamente ao usuário a possibilidade de processar uma quantidade bem grande de dados; segundo Ugo Roveda Co-fundador e COO na Kenzie Academy Brasil.

3.2.1 BIBLIOTECAS

The screenshot shows a Jupyter Notebook with the following code cells:

```
[1] 1 import numpy as np
    2 import pandas as pd
    3 import matplotlib.pyplot as plt
    4 import seaborn as sns
    5 from pandas.plotting import scatter_matrix
    6 from google.colab import drive

[2] 1 from google.colab import drive
    2 drive.mount('/content/drive')

Mounted at /content/drive

[3] 1 drive.mount('/content/drive')

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True)

[4] 1 df = pd.read_csv("/content/drive/MyDrive/adult.csv",1,")
    2 data = [df]
    3 print ("Inicial dataset \n\n", df.head(10), "\n")
```

O aumento na quantidade de dados gerados e a necessidade que algumas aplicações têm de analisar esses dados em tempo real estão entre as principais demandas para bancos de dados em memória. Nesta imagem 01 você verá o que leu na página 06 desta pesquisa.

3.2.3. DISCRETIZAÇÃO INCOME

The screenshot shows a Jupyter Notebook with the following code cell and output:

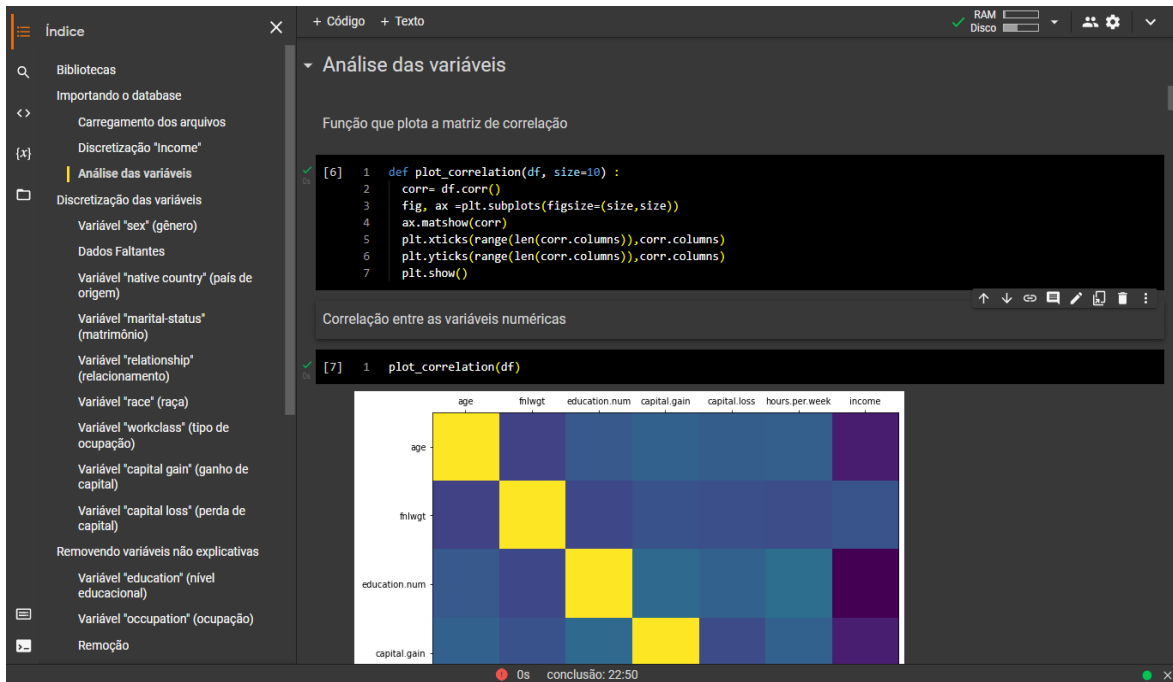
```
1 df['income'] = df['income'].map({'<=50K': 1, '>50K': 0}, {'<=50K.': 1, '>50K.': 0})
2 #income_map = {'<=50K': 1, '>50K': 0}
3 #df['income'] = df['income'].map(income_map).astype(int)
4
5 print ("Income classified by binary \n\n", df.head(10), "\n")
```

Income classified by binary

	age	workclass	fnlwt	...	hours.per.week	native.country	income
0	90	?	77053	...	40	United-States	1
1	82	Private	132870	...	18	United-States	1
2	66	?	186061	...	40	United-States	1
3	54	Private	140359	...	40	United-States	1
4	41	Private	264663	...	40	United-States	1
5	34	Private	216864	...	45	United-States	1
6	38	Private	150601	...	40	United-States	1
7	74	State-gov	88638	...	20	United-States	0
8	68	Federal-gov	422013	...	40	United-States	1
9	41	Private	70037	...	60	?	0

Como, nesta imagem 02 Tornar discreto ou descontínuo; transformar uma distribuição contínua em unidades individuais.

3.2.4. ANÁLISE DAS VARIÁVEIS



As variáveis podem ser classificadas em quantitativas ou qualitativas. Uma variável quantitativa possui como especificidade ter como atribuição um valor numérico. Uma variável qualitativa possui como especificidade ter como atribuição um valor não-numerado.

3.2.5. DISCRETIZAÇÃO DAS VARIÁVEIS

The screenshot shows a Jupyter Notebook interface with a sidebar on the left containing a project index. The main area displays a heatmap for variables 'capital loss', 'hours per week', and 'income'. Below the heatmap, the notebook is titled 'Discretização das variáveis' and shows the following code and output:

```

1 print(df[['sex', 'income']].groupby('sex').mean(), "\n")

```

	income
sex	
Female	0.890539
Male	0.694263

```

[9] 1 sex_map = {'Male':1, 'Female':0}
    2 df['sex'] = df['sex'].map({'Male':1, 'Female':0})

[10] 1 print("Sex classified by binary \n", df.head(10), "\n")

```

Sex classified by binary

The screenshot shows a Jupyter Notebook interface with a sidebar on the left containing a project index. The main area displays a heatmap for variables 'native.country' and 'income'. Below the heatmap, the notebook shows the following code and output:

```

1 print(df[['native.country', 'income']].groupby('native.country').mean())

```

native.country	income
?	0.749571
Cambodia	0.631579
Canada	0.677686
China	0.733333
Columbia	0.966102
Cuba	0.736842
Dominican-Republic	0.971429
Ecuador	0.857143
El-Salvador	0.915894
England	0.666667
France	0.586207
Germany	0.678832
Greece	0.724138
Guatemala	0.953125
Haiti	0.909091
Holand-Netherlands	1.000000
Honduras	0.923077
Hong	0.700000
Hungary	0.769231
India	0.600000
Iran	0.581395
Ireland	0.791667
Italy	0.657534
Jamaica	0.876543
Japan	0.612903
Laos	0.888889
Mexico	0.948678
Nicaragua	0.941176
Outlying-US(Guam-USVI-etc)	1.000000
Peru	0.935484
Philippines	0.691919
Poland	0.800000
Portugal	0.891892
Puerto-Rico	0.894737
Scotland	0.750000
South	0.800000

Poland 0.880000
Portugal 0.891892
Puerto-Rico 0.894737
Scotland 0.750000
South 0.800000
Taiwan 0.807843
Thailand 0.833333
Trinidad&Tobago 0.894737
United-States 0.754165
Vietnam 0.925373
Yugoslavia 0.625000

Precisamos remover as linhas com "?", que são dados faltantes. Se não sabemos de que país essa pessoa veio, esse dado não é informativo.

```
[12] 1 print("\nTotal registers", df.shape)
```

Total registers (32561, 15)

```
[13] 1 df["native.country"] = df["native.country"].replace('?', np.nan)
2 df["workclass"] = df["workclass"].replace('?', np.nan)
3 df["occupation"] = df["occupation"].replace('?', np.nan)
```

```
[14] 1 print("Dados faltantes \n", df.head(10), "\n")
```

Dados faltantes

	age	workclass	fnlwtg	...	hours.per.week	native.country	income
0	90	NaN	77053	...	40	United-States	1
1	82	Private	132870	...	18	United-States	1
2	66	NaN	186061	...	40	United-States	1
3	54	Private	140359	...	40	United-States	1
4	41	Private	264663	...	40	United-States	1
5	34	Private	216864	...	45	United-States	1
6	38	Private	150601	...	40	United-States	1
7	74	State-gov	88638	...	20	United-States	0

A escala ordinal é definida como uma escala de medição de variáveis utilizada para simplesmente representar a ordem das variáveis e não a diferença entre cada uma das variáveis. O que vimos nas imagens 04,05 e 06.

3.2.6. VARIÁVEL NATIVE COUNTRY

```
[15] 1 df.dropna(how='any', inplace=True)
```

```
[16] 1 print("Valid registers", df.shape)
2 print("\nTabel with invalid values removed \n\n", df.head(10), "\n")
```

Valid registers (30162, 15)

Tabel with invalid values removed

	age	workclass	fnlwtg	...	hours.per.week	native.country	income
1	82	Private	132870	...	18	United-States	1
3	54	Private	140359	...	40	United-States	1
4	41	Private	264663	...	40	United-States	1
5	34	Private	216864	...	45	United-States	1
6	38	Private	150601	...	40	United-States	1
7	74	State-gov	88638	...	20	United-States	0
8	68	Federal-gov	422013	...	40	United-States	1
10	45	Private	172274	...	35	United-States	0
11	38	Self-emp-not-inc	164526	...	45	United-States	0
12	52	Private	129177	...	20	United-States	0

Variável "native country" (país de origem)

```
[17] 1 for dataset in data:
2 dataset.loc[dataset['native.country'] != \
```

O conjunto de dados possui variáveis, numéricas e categóricas. `native.country`: país de origem; `gender`: gênero. Nesta imagem 07

3.2.7. VARIÁVEL MARITAL-STATUS

RAM Disco

```

3 |         'United-States', 'native.country' ] = 'Non-US'
4 | dataset.loc[dataset['native.country'] == \
5 |             'United-States', 'native.country' ] = 'US'

[18] 1 | native_country_map = {'US':1, 'Non-US':0}
     2 | df['native.country'] = df['native.country'].map(native_country_map).astype(int)

[19] 1 | print("native.countryca classified by binary \n\n", df.head(10), "\n")

```

native.countryca classified by binary

	age	workclass	fnlwtg	...	hours.per.week	native.country	income
1	82	Private	132870	...	18	1	1
3	54	Private	140359	...	40	1	1
4	41	Private	264663	...	40	1	1
5	34	Private	216864	...	45	1	1
6	38	Private	150601	...	40	1	1
7	74	State-gov	88638	...	20	1	0
8	68	Federal-gov	422013	...	40	1	1
10	45	Private	172274	...	35	1	0
11	38	Self-emp-not-inc	164526	...	45	1	0
12	52	Private	129177	...	20	1	0

[10 rows x 15 columns]

▼ Variável "marital-status" (matrimônio)

```

[20] 1 | print(df[['marital.status', 'income']].groupby(['marital.status']).mean(), "\n")

```

	income
Divorced	0.892738
Married-AF-spouse	0.522810
Married-civ-spouse	0.545841

0s conclusão: 22:50

RAM Disco

```

Married-civ-spouse      0.545841
Married-spouse-absent  0.916216
Never-married           0.951676
Separated               0.929712
Widowed                 0.903265

[22] 1 | df['marital.status'] = df['marital.status'].replace\
     2 |     (('Divorced', 'Married-spouse-absent', 'Never-married', 'Separated', \
     3 |      'Widowed'), 'Single')
     4 |
     5 | df['marital.status'] = df['marital.status'].replace\
     6 |     (('Married-AF-spouse', 'Married-civ-spouse'), 'Couple')

[21] 1 | print("marital.status classified by binary class\n\n", df.head(10), "\n")

```

marital.status classified by binary class

	age	workclass	fnlwtg	...	hours.per.week	native.country	income
1	82	Private	132870	...	18	1	1
3	54	Private	140359	...	40	1	1
4	41	Private	264663	...	40	1	1
5	34	Private	216864	...	45	1	1
6	38	Private	150601	...	40	1	1
7	74	State-gov	88638	...	20	1	0
8	68	Federal-gov	422013	...	40	1	1
10	45	Private	172274	...	35	1	0
11	38	Self-emp-not-inc	164526	...	45	1	0
12	52	Private	129177	...	20	1	0

[10 rows x 15 columns]

```

[23] 1 | print(df[['marital.status', 'income']].groupby(['marital.status']).mean())

```

	income
Couple	0.545841
Single	0.916216

0s conclusão: 22:50

```
[23] 1 print(df[['marital.status', 'income']].groupby(['marital.status']).mean())
```

marital.status	income
Couple	0.545009
Single	0.931637

```
[24] 1 print(df[['marital.status', 'relationship', 'income']].groupby\
2 | (('marital.status', 'relationship')).mean())
```

marital.status	relationship	income
Couple	Husband	0.544331
	Not-in-family	0.714286
	Other-relative	0.857143
	Own-child	0.821429
Single	Wife	0.506401
	Not-in-family	0.893802
	Other-relative	0.976623
	Unmarried	0.988818
	Unmarried	0.933686

```
[25] 1 print(df[['marital.status', 'relationship', 'income']].groupby\
2 | (('relationship', 'marital.status')).mean())
```

relationship	marital.status	income
Husband	Couple	0.544331
	Single	0.714286
Not-in-family	Couple	0.893802
	Single	0.803802
Other-relative	Couple	0.857143
	Single	0.976623
Own-child	Couple	0.821429
	Single	0.988818
Unmarried	Single	0.933686
	Couple	0.506401

```
[26] 1 marital_status_map = {'Single':1, 'Couple':0}
```

Traduções em contexto de "marital status" em inglês-português da Reverso Context: This referente à variável situação conjugal. Nesta imagem 08,09 e 10

3.2.8. VARIÁVEL RELATIONSHIP

```
[26] 1 marital_status_map = {'Single':1, 'Couple':0}
2 df['marital.status'] = df['marital.status'].map(marital_status_map)
```

```
[27] 1 print("\nmarital.status classified by binary values\n\n", df.head(10), "\n")
```

marital.status classified by binary values							
	age	workclass	fnlwt	...	hours.per.week	native.country	income
1	82	Private	132870	...	18	1	1
3	54	Private	140359	...	40	1	1
4	41	Private	264663	...	40	1	1
5	34	Private	216864	...	45	1	1
6	38	Private	150601	...	48	1	1
7	74	State-gov	88638	...	20	1	0
8	68	Federal-gov	422013	...	40	1	1
10	45	Private	172274	...	35	1	0
11	38	Self-emp-not-inc	164526	...	45	1	0
12	52	Private	129177	...	20	1	0

[10 rows x 15 columns]

```
[28] 1 print(df[['relationship', 'income']].groupby(['relationship']).mean(), "\n")
```

relationship	income
Husband	0.544331
Not-in-Family	0.893477
Other-relative	0.960630
Own-child	0.985670
Unmarried	0.933686

Variáveis relevantes para um bom relacionamento entre cooperativa e relationship between customers, serving as a subsidy for the construction of the. Nesta imagem 11

varável: workclass (tipo de ocupação)

Pressione F11 para sair do modo tela cheia

```
[33] 1 def f(x):
2     if x['workclass'] == 'Federal-gov' or x['workclass'] == 'Local-gov' or \
3         x['workclass'] == 'State-gov': return 'govt'
4     elif x['workclass'] == 'Private': return 'private'
5     elif x['workclass'] == 'Self-emp-inc' or x['workclass'] == \
6         'Self-emp-not-inc': return 'self_employed'
7     else: return 'without_pay'
```

```
[34] 1 df['workclass'] = df.apply(f, axis=1)
```

```
[35] 1 print("\nworkclass classified by groups\n\n", df.head(10), "\n")
```

workclass classified by groups

	age	workclass	fnlwgt	...	hours.per.week	native.country	income
1	82	private	132870	...	18	1	1
3	54	private	140359	...	40	1	1
4	41	private	264663	...	40	1	1
5	34	private	216864	...	45	1	1
6	38	private	150601	...	40	1	1
7	74	govt	88638	...	20	1	0
8	68	govt	422013	...	40	1	1
10	45	private	172274	...	35	1	0
11	38	self_employed	164526	...	45	1	0
12	52	private	129177	...	20	1	0

[10 rows x 15 columns]

```
[36] 1 print(df[['workclass', 'income']].groupby(['workclass']).mean())
```

income

workclass	income
-----------	--------

0s conclusão: 22:50

Catégoricas, a análise de classes latentes é uma possível técnica a ser utilizada ... Variável categórica que corresponde ao tipo de ocupação declarado pelo. Nessas imagens 13 e 14

3.2.11. VARIÁVEL CAPITAL GAIN

workclass

workclass	income
govt	0.692702
private	0.781208
self_employed	0.632242
without_pay	1.000000

```
[37] 1 workclass_map = {'govt':0, 'private':1, 'self_employed':2, 'without_pay':3}
2
3 df['workclass'] = df['workclass'].map(workclass_map)
```

```
[38] 1 print("\nworkclass calssified by values\n\n", df.head(10), "\n")
```

workclass calssified by values

	age	workclass	fnlwgt	...	hours.per.week	native.country	income
1	82	1.0	132870	...	18	1	1
3	54	1.0	140359	...	40	1	1
4	41	1.0	264663	...	40	1	1
5	34	1.0	216864	...	45	1	1
6	38	1.0	150601	...	40	1	1
7	74	0.0	88638	...	20	1	0
8	68	0.0	422013	...	40	1	1
10	45	1.0	172274	...	35	1	0
11	38	2.0	164526	...	45	1	0
12	52	1.0	129177	...	20	1	0

[10 rows x 15 columns]

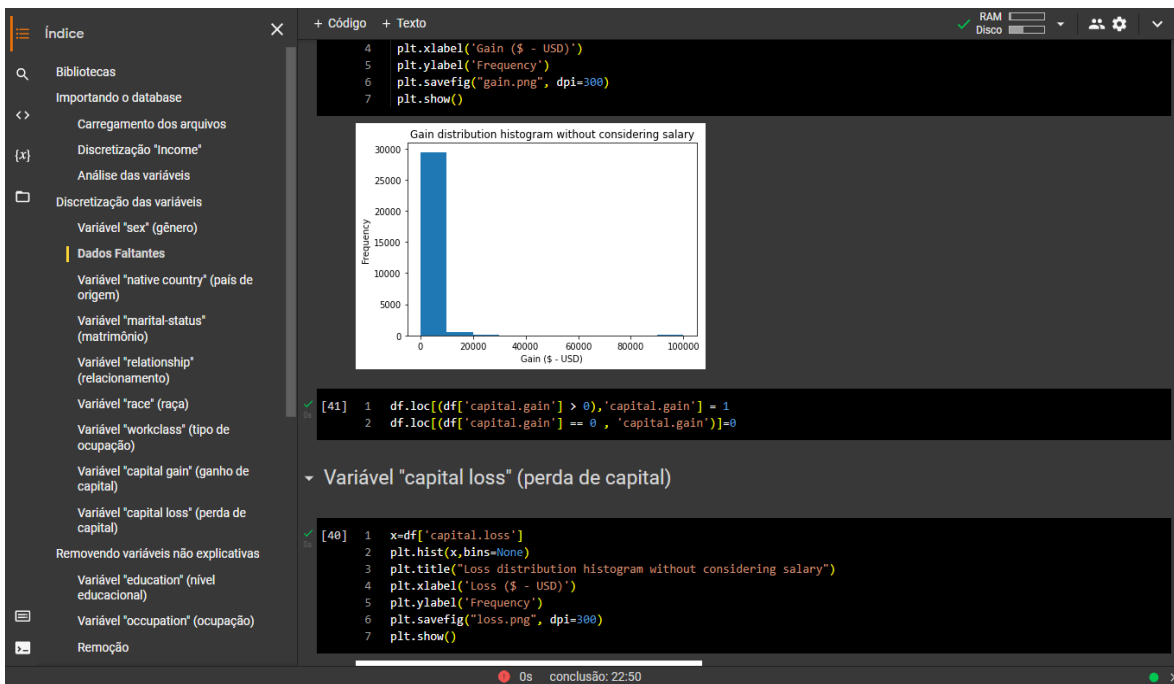
▼ Variável "capital gain" (ganho de capital)

```
[39] 1 x=df['capital.gain']
2 plt.hist(x,bins=None)
3 plt.title("gain distribution histogram without considering salary")
4 plt.xlabel('Gain ($ - US$)')
```

0s conclusão: 22:50

Dentre os resultados obtidos, constatou-se que a variável que ... Social capital, Growth and Poverty: a Survey of Cross-Country Evidence. Nesta imagem 15

3.2.12. VARIÁVEL CAPITAL LOSS



A variável capital próprio indica uma eventual utilização da provisão. Nesta imagem 16

3.2.13. REMOVENDO VARIÁVEIS NÃO EXPLICATIVAS

Índice

Bibliotecas

Importando o database

Carregamento dos arquivos

Discretização 'Income'

Análise das variáveis

Discretização das variáveis

Variável 'sex' (gênero)

Dados Faltantes

Variável 'native country' (país de origem)

Variável 'marital-status' (matrimônio)

Variável 'relationship' (relacionamento)

Variável 'race' (raça)

Variável 'workclass' (tipo de ocupação)

Variável 'capital gain' (ganho de capital)

Variável 'capital loss' (perda de capital)

Removendo variáveis não explicativas

Variável 'education' (nível educacional)

Variável 'occupation' (ocupação)

Remoção

+ Código + Texto

Variável "capital loss" (perda de capital)

```
[40] 1 x=df['capital.loss']
      2 plt.hist(x,bins=None)
      3 plt.title("Loss distribution histogram without considering salary")
      4 plt.xlabel('loss ($ - USD)')
      5 plt.ylabel('Frequency')
      6 plt.savefig("loss.png", dpi=300)
      7 plt.show()
```

```
[42] 1 df.loc[(df['capital.loss'] > 0), 'capital.loss'] = 1
      2 df.loc[(df['capital.loss'] == 0, 'capital.loss')] = 0
```

Removendo variáveis não explicativas

Variável "education" (nível educacional)

0s conclusão: 22:50

3.2.14. REMOÇÃO

Índice

Bibliotecas

Importando o database

Carregamento dos arquivos

Discretização 'Income'

Análise das variáveis

Discretização das variáveis

Variável 'sex' (gênero)

Dados Faltantes

Variável 'native country' (país de origem)

Variável 'marital-status' (matrimônio)

Variável 'relationship' (relacionamento)

Variável 'race' (raça)

Variável 'workclass' (tipo de ocupação)

Variável 'capital gain' (ganho de capital)

Variável 'capital loss' (perda de capital)

Removendo variáveis não explicativas

Variável 'education' (nível educacional)

Variável 'occupation' (ocupação)

Remoção

+ Código + Texto

```
[44] 1 print(df[['occupation', 'income']].groupby(['occupation']).mean())
```

occupation	income
Adm-clerical	0.866165
Armed-Forces	0.888889
Craft-repair	0.774690
Exec-managerial	0.514780
Farming-fishing	0.883721
Handlers-cleaners	0.938519
Machine-op-inspct	0.875381
Other-service	0.958904
Priv-house-serv	0.993007
Prof-specialty	0.551511
Protective-serv	0.673913
Sales	0.729353
Tech-support	0.695175
Transport-moving	0.797074

Remoção

Variável redundante

```
[45] 1 df.drop(labels=['education', 'occupation'], axis=1, inplace=True)
```

```
[46] 1 print("\neducation and occupation are removed because are not significant " \
      2 | | | "explain variables \n\n", df.head(10), "\n")
```

```
education and occupation are removed because are not significant explain variables
   age workclass  fnlwgt  ...  hours.per.week  native.country  income
1   82      1.0  132870  ...             18              1         1
```

0s conclusão: 22:50

3.2.15. VARIÁVEIS NÃO DISCRETIZADOS

education and occupation are removed because are not significant explain variables

	age	workclass	fnlwt	...	hours.per.week	native.country	income
1	82	1.0	132870	...	18	1	1
3	54	1.0	140359	...	40	1	1
4	41	1.0	264663	...	40	1	1
5	34	1.0	216864	...	45	1	1
6	38	1.0	150601	...	40	1	1
7	74	0.0	88638	...	20	1	0
8	68	0.0	422013	...	40	1	1
10	45	1.0	172274	...	35	1	0
11	38	2.0	164526	...	45	1	0
12	52	1.0	129177	...	20	1	0


[10 rows x 13 columns]

- ▼ Variáveis não-discretizáveis
- ▼ Variável "hours per week" (horas trabalhadas por semana)

```

1 x= df["hours.per.week"]
2 plt.hist(x,bins=None,density=True,histtype='bar')
3 plt.title("Hour worked distribution histogram")
4 plt.xlabel("Hours worked per week (Total hours)")
5 plt.ylabel("Frequency")
6 plt.savefig("hours.png", dpi=300)
7 plt.show()

```



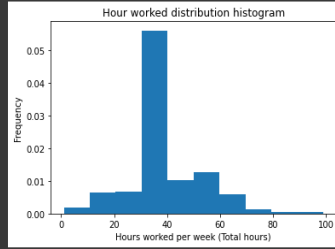
0s conclusão: 22:50

3.2.16. VARIÁVEL EDUCATION NUM

```

7 plt.show()

```

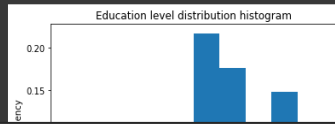


- ▼ Variável "education num" (número de anos na escola)

```

[48] 1 x= df["education.num"]
2 plt.hist(x,bins=None,density=True,histtype='bar')
3 plt.title("Education level distribution histogram")
4 plt.xlabel("Education level")
5 plt.ylabel("Frequency")
6 plt.savefig("education.png", dpi=300)
7 plt.show()

```



0s conclusão: 22:50

3.2.17. VARIÁVEL AGE

The screenshot shows a Jupyter Notebook interface. On the left, a sidebar lists various data processing steps under the heading 'Índice'. The main area displays a histogram titled 'Education level distribution histogram' with 'Education level' on the x-axis (ranging from 2 to 16) and 'Frequency' on the y-axis (ranging from 0.00 to 0.20). Below the histogram, a code cell is shown with the following Python code:

```
[49] 1 x= df['age']
      2 plt.hist(x,bins=None,density=True,histtype='bar')
      3 plt.title("Age distribution histogram")
      4 plt.xlabel("Age (Years)")
      5 plt.ylabel("Frequency")
      6 plt.savefig("age.png", dpi=300)
      7 plt.show()
      8 plt.savefig("age.pdf")
      9
```

Below the code cell, a smaller histogram titled 'Age distribution histogram' is visible, showing the distribution of age with 'Age (Years)' on the x-axis (ranging from 20 to 90) and 'Frequency' on the y-axis (ranging from 0.000 to 0.025).

3.2.18. ANALISANDO A CORRELAÇÃO

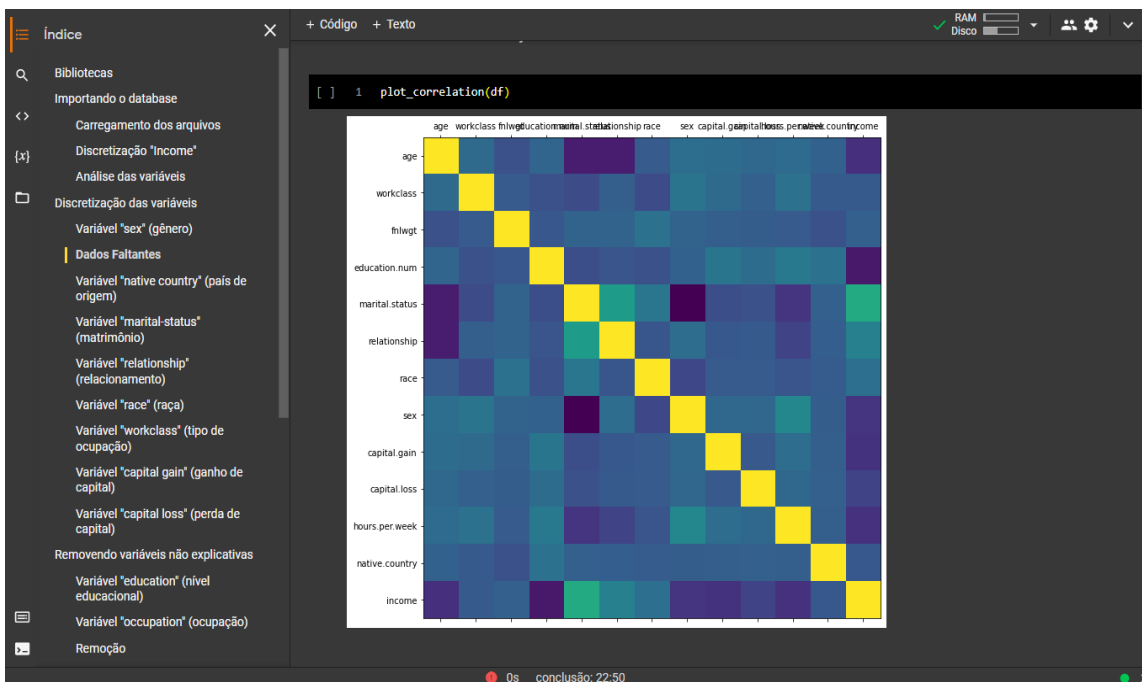
The screenshot shows a Jupyter Notebook interface. On the left, a sidebar lists various data processing steps under the heading 'Índice'. The main area displays a histogram titled 'Age distribution histogram' with 'Age (Years)' on the x-axis (ranging from 20 to 90) and 'Frequency' on the y-axis (ranging from 0.000 to 0.025). Below the histogram, a code cell is shown with the following Python code:

```
[50] 1 df['age'].count()
      30162
```

Below the code cell, a code cell is shown with the following Python code:

```
[ ] 1 plot_correlation(df)
```

Below the code cell, a correlation matrix heatmap is displayed, showing the relationships between variables: age, workclass, fnlwgt, education, marital, statelocation, race, sex, capital, gsp, ita, loss, per, week, country, and ome. The heatmap uses a color scale from blue (low correlation) to yellow (high correlation).



3.2.19. IMPORTANDO O DATA SPLIT

Construção dos classificadores

▼ Importando o dataSplit

```
[51] 1 from sklearn.model_selection import train_test_split
      2 from sklearn.model_selection import Kfold
      3
```

▼ Particionando os dados

O conjunto de dados será dividido da seguinte forma

- 50% para treinamento
- 20% para validação
- 30% para teste

```
[52] 1 x = df.drop(['income'], axis=1)
      2 y = df['income']

[53] 1 split_size=0.3
      2
      3 #Creation of Train and Test dataset
      4 x_train, x_test, y_train, y_test = train_test_split\
      5 | (x,y, test_size=split_size, random_state=22)
      6
      7 #Creation of Train and validation dataset
```

3.2.20. IMPORTANDO OS CLASSIFICADORES

```

17
18
19
[55] 1 print ("Train dataset: {}".format(x_train.shape, y_train.shape))
      2 print ("Validation dataset: {}".format(x_val.shape, y_val.shape))
      3 print ("Test dataset: {}".format(x_test.shape, y_test.shape))

Train dataset: (16890, 12)(16890,)
Validation dataset: (4223, 12)(4223,)
Test dataset: (9049, 12)(9049,)

▼ Importando os classificadores

[56] 1 from sklearn.neighbors import KNeighborsClassifier
      2 from sklearn.naive_bayes import GaussianNB
      3 from sklearn.tree import DecisionTreeClassifier
      4 from sklearn.ensemble import RandomForestClassifier
      5 from sklearn.linear_model import LogisticRegression
      6 from sklearn.svm import SVC

▼ Criando os modelos de treinamento

[57] 1 models = []
      2
      3 names = ["Nearest Neighbors", "Naive Bayes", "Decision Tree", \
      4         "Random Forest", "Logistic Regression", "SVC"]
      5

[58] 1 models.append((KNeighborsClassifier(n_neighbors =50)))
      2 models.append((GaussianNB()))
  
```

```

5
[58] 1 models.append((KNeighborsClassifier(n_neighbors =50)))
      2 models.append((GaussianNB()))
      3 models.append((DecisionTreeClassifier()))
      4 models.append((RandomForestClassifier(n_estimators=100)))
      5 models.append((LogisticRegression()))
      6 models.append((SVC()))

[59] 1 print(models)

[KNeighborsClassifier(n_neighbors=50), GaussianNB(), DecisionTreeClassifier(), RandomForestClassifier(), LogisticRegression()]

▼ Importando os métodos de treinamento

[60] 1 from sklearn import model_selection
      2 from sklearn.model_selection import Kfold
      3 from sklearn.metrics import accuracy_score
      4

[64] 1 KF = model_selection.KFold(n_splits=5, random_state=None, shuffle=True)

▼ Importando as bibliotecas de avaliação

[61] 1

[68] 1
      2
  
```

3.2.21. CÓDIGO COM ERRO

```
[61] 1
[68] 1
2
3
4 for i in range(0, len(models)):
5     kfold = model_selection.KFold(n_splits=5, random_state=None)
6     cv_result = model_selection.cross_val_score
7     (models[i], x_train, y_train, cv=kfold, scoring='accuracy')
8     score=models[i].fit(x_train, y_train)
9     prediction = models[i].predict(x_val)
10    acc_score = accuracy_score(y_val, prediction)
11    print ('-'*40)
12    print ('{0}: {1}'.format(names[i], acc_score))

Importando as bibliotecas de avaliação

[63] 1 from sklearn.metrics import classification_report
2     from sklearn.metrics import confusion_matrix
3
4     from sklearn.model_selection import GridSearchCV

Avaliação de resultados

[63] 1
```

3.2.22. Solução Final

```
[ ] 1
2
3
4 for i in range(0, len(models)):
5     kfold = model_selection.KFold(n_splits=5, random_state=None)
6     cv_result = model_selection.cross_val_score /
7     (models[i], x_train, y_train, cv=kfold, scoring='accuracy')
8     score=models[i].fit(x_train, y_train)
9     prediction = models[i].predict(x_val)
10    acc_score = accuracy_score(y_val, prediction)
11    print ('-'*40)
12    print ('{0}: {1}'.format(names[i], acc_score))

File "c:\python-input-1-ecb14d9eb7ab2", line 6
cv_result = model_selection.cross_val_score /
^
SyntaxError: invalid syntax

SEARCH STACK OVERFLOW

Importando as bibliotecas de avaliação

[ ] 1 from sklearn.metrics import classification_report
2     from sklearn.metrics import confusion_matrix
3
4     from sklearn.model_selection import GridSearchCV

Avaliação de resultados
```

Nesta parte do código, ficamos com uma enorme dor de cabeça sobre este pedaço de bloco de código, pois de nenhuma forma obtemos resultado sobre a correção da string. Com isso percebemos que no prosseguimento das linhas, estão faltando pequenos trechos, e nisto, paramos na edição e digitação deste bloco. O restante do conjunto está funcionando perfeitamente, tendo como finalidade

um acerto significativo de linhas corretas. Sendo assim, concluímos que a solução deste pedaço de código ficará para uma próxima checagem, onde iremos trazer soluções plausíveis e coerentes para que seja aplicado um pouco mais de conhecimento.

4. CONSIDERAÇÕES FINAIS

A pesquisa, como um todo, demonstrou ter um grande número de observações relevantes na construção de estruturas a partir desse banco de dados adulto. Tivemos, a princípio, algumas dificuldades para se estudar tal tema, já que esta pesquisa, que está disponível na kaggle.com, apresentou diferenças em nosso método de formatação.

Começamos, em primeiro lugar, agindo na base de análise, parte por parte, antes mesmo de se fazer qualquer outra coisa. Pensamos que seria de certa forma, difícil chegar a todo aquele amontoado de informações diversas, mas que de mesmo ponto de interesse - que ainda assim, havia partes dos dados, faltantes. Então, a partir daí, tivemos que começar com as breves edições, apenas para dar uma outra forma neste dataset.

Nisto, começamos a parte de edição usando o aplicativo do pacote office, Excel. Já que nele conseguimos ter uma visão mais encadeada de quais colunas usar na plataforma do google colab, onde foi o inspecionado, digitalizado, estruturado e assim, deixando de uma forma com processo fácil de interpretação, que assim já temos uma estrutura feita, compreensível, legível e entendível. A importância que se tem nesses pontos é que o acréscimo de funções baseadas em Python no banco de dados se torna acessível.

O banco, como vimos no início para o momento pós edição, é completamente diferente, já que seus dados foram divididos de acordo com seus blocos. Com isso, a estrutura continua com a sua adequação, ligação e correspondentes. Adult census income, como já dito, é uma base de dados com informações de pessoas moradoras da américa, onde há diversas etnias, classes sociais, matrimônios, patrimônios, educação. Pessoas dos seus 17 anos até seus 90. O objetivo deste dataset, é mostrar em seus blocos a divisão com suas respectivas sequências de informações para que se pudesse extrair suas referências.

Contudo, chega-se à conclusão de que a base de dados voltada para entender como é dado o valor anual de que cada pessoa recebe em seu trabalho, é mostrado em casa sequência de código baseado em python que foi escrito. Mas durante todo esse caminho, houve momentos em que tivemos que recorrer a outros métodos para que fosse entendido e corrigido linhas de código que não estavam em combinação de resultados. Vendo vídeos em youtube, acessando diferentes blogs e páginas relacionadas ao dataset, foi possível converter os erros. Mesmo assim, a gama maior de acertos foi do vídeo aula que também nos foi deixava para ser usado como forma de acesso fácil de correção caso fosse necessário. E foi, então, para se concluir, esta base de dados foi útil na questão de se entender de que forma as pessoas são remuneradas nos EUA. Muito se deve à educação e família, onde se tem muita influência, já que se a família fosse de boa situação, de ensino, de bom emprego, de boa influência, ajudaria suas gerações; mas nesse caso, o impacto forma que tem na educação, tempo de estudo, o trabalho em si, reflete a curto, médio e longo prazo, em relação à faixa etária das pessoas pesquisadas.

REFERÊNCIAS BIBLIOGRÁFICAS

http://www.teses.usp.br/index.php?option=com_content&view=article&id=52&Itemid=67

<https://medium.com/data-warriors/eda-of-adult-census-income-dataset-cc9ac1a3d552>

https://drive.google.com/drive/folders/1UCPP0pXZPXLQnk6O_EETuqgrS9gjAf98?usp=sharing

<https://pt.stackoverflow.com/questions/66777/d%C3%BAvidas-na-utiliza%C3%A7%C3%A3o-de-stratified-k-fold-no-scikit-learn>

<https://ieeexplore.ieee.org/document/8748528>

<https://www.kaggle.com/uciml/adult-census-income>

<https://www.youtube.com/watch?v=kEPYo45n51M>