

8.Git

October 26, 2019

Para entrar no modo apresentação, execute a seguinte célula e pressione -

```
[1]: %reload_ext slide
```

1 Git

Este notebook apresenta os seguintes tópicos:

- Section 1 - Git
- Section 1.1 - Exercício 9
- Section 1.2 - Pandas
- Section 1.3 - Exercício 10

Outra fonte de informações de um repositório de software é o repositório do sistema de controle de versões.

Pelo controle de versões, conseguimos ter acesso a todos os arquivos de todas as versões, todas as mensagens de commit, branches, e colaboradores.

Nesta parte do minicurso, faremos a mineração dessas informações.

No caso do Git, ao clonar um repositório, ficamos com uma cópia local do que está lá. Portanto, começamos a mineração com um clone e não precisamos de nenhum proxy.

```
[2]: !git clone https://github.com/gems-uff/sapos
```

```
Cloning into 'sapos'...
remote: Enumerating objects: 22, done.
remote: Counting objects: 100% (22/22), done.
remote: Compressing objects: 100% (19/19), done.
remote: Total 12954 (delta 4), reused 11 (delta 3), pack-reused 12932
Receiving objects: 100% (12954/12954), 10.41 MiB | 6.43 MiB/s, done.
Resolving deltas: 100% (8011/8011), done.
```

Com o repositório clonado, podemos usar comandos do git para extrair informações.

```
[3]: %cd sapos
```

```
/home/joao/projects/tutorial/sapos
```

```
[4]: !git branch -r
```

```
origin/HEAD -> origin/master  
origin/bugfixes  
origin/hotfixes  
origin/master  
origin/reports
```

Essas informações também podem ser obtidas para tratarmos usando variáveis do Python.

A seguir tentamos descobrir qual é o commit de cada um desses branches.

```
[5]: git_branch_output = !git branch -r  
branches = [  
    branch.strip().split(' ')[0].split('/')[1]  
    for branch in git_branch_output  
]  
branches
```

```
[5]: ['HEAD', 'bugfixes', 'hotfixes', 'master', 'reports']
```

```
[6]: branch_commit = {}  
for branch in branches:  
    __ = !git checkout $branch  
    commit = !git show --pretty=format:"%h" --no-patch  
    branch_commit[branch] = commit  
__ = !git checkout master  
branch_commit
```

```
[6]: {'HEAD': ['d3e4307'],  
      'bugfixes': ['4d22176'],  
      'hotfixes': ['d3e4307'],  
      'master': ['d3e4307'],  
      'reports': ['b294935']}
```

Usamos `__ = !...` para evitar a exibição do output do comando de sistema. O IPython imprime o output quando bang expressions são usadas isoladas e retorna o output quando elas são usadas em atribuições.

Note que apenas o branch `reports` está em um commit diferente.

1.1 Exercício 9

Faça a mesma operação para obter o código dos commits de tags e salve na variável `tag_commit`.

```
[7]: tags = !git tag
```

```
...
```

```
[7]: {'0.1.0': ['cfd4fc2'],
      '0.2.0': ['9cf8be6'],
      '0.3.0': ['336dcc3'],
      '0.3.1': ['336dcc3'],
      '1.0.0': ['8e2de97'],
      '1.0.1': ['2638fc1'],
      '1.1.0': ['bd99c54'],
      '1.1.1': ['359f61a'],
      '1.1.2': ['e702a88'],
      '1.2.0': ['f78d9b0'],
      '1.3.0': ['41edd16'],
      '1.3.1': ['7327e04'],
      '1.3.2': ['c668246'],
      '1.3.3': ['c908b61'],
      '1.3.4': ['f74e3be'],
      '1.4.0': ['77b9d24'],
      '1.4.1': ['6c7e5cf'],
      '1.4.2': ['64c40f2'],
      '1.4.3': ['3be5566'],
      '1.5.0': ['492834d'],
      '1.5.1': ['1f0bdba'],
      '1.5.2': ['1f0bdba'],
      '1.5.3': ['b6504b1'],
      '1.5.4': ['a765764'],
      '1.5.5': ['f4e006b'],
      '1.5.6': ['72522e1'],
      '1.5.7': ['2f11de1'],
      '1.5.8': ['84065cd'],
      '1.5.9': ['53f3623'],
      '1.6.0': ['576aead'],
      '1.6.1': ['6b8a35e'],
      '1.6.2': ['8a22abe'],
      '1.6.3': ['1f3c0b6'],
      '1.7.0': ['1facb9b'],
      '1.7.1': ['25c728e'],
      '1.8.0': ['a691029'],
      '1.8.1': ['f2a374d'],
      '1.8.2': ['d5b4c55'],
      '1.8.3': ['83805d1'],
      '1.8.4': ['ed1a888'],
      '1.8.5': ['df4ad8e'],
      '1.9.0': ['4db772d'],
      '1.9.1': ['ed20262'],
      '1.9.2': ['e04feea'],
      '1.9.3': ['43dee99'],
      '1.9.4': ['197fb76'],
      '2.0.0': ['9a6cf14'],
```

'2.0.1': ['718efaf'],
'2.0.2': ['9d18e2d'],
'2.0.3': ['2967241'],
'2.0.4': ['c9ab787'],
'2.0.5': ['740b9b2'],
'2.0.6': ['de02470'],
'2.0.7': ['4ea98f3'],
'2.1.0': ['391075a'],
'2.1.1': ['f6796cf'],
'2.1.2': ['8e83886'],
'2.1.3': ['861fbae'],
'2.2.0': ['5eb9b88'],
'2.2.1': ['e10a3a0'],
'2.2.2': ['35930ac'],
'2.2.3': ['d413076'],
'2.2.4': ['1306c8a'],
'2.3.0': ['f13421f'],
'2.3.1': ['be0dcf1'],
'2.3.2': ['8687004'],
'2.3.3': ['d1659f1'],
'2.3.4': ['841ca4e'],
'2.3.5': ['7ce001e'],
'2.3.6': ['a50d45f'],
'2.3.7': ['8ee3a43'],
'2.4.0': ['62b985b'],
'2.4.1': ['913c628'],
'2.4.2': ['89ea337'],
'2.4.3': ['8de4bf7'],
'2.4.4': ['db0eef2'],
'2.4.5': ['1825b61'],
'2.4.6': ['2c8c06f'],
'2.4.7': ['8022870'],
'2.4.8': ['aca8f7d'],
'3.0.0': ['9590421'],
'3.1.0': ['8e2d96a'],
'3.2.0': ['cc65c65'],
'3.2.1': ['b4675d1'],
'3.3.0': ['fa5d903'],
'3.3.1': ['797d47a'],
'3.3.2': ['870d838'],
'3.3.3': ['48a5d6a'],
'3.3.4': ['443b11c'],
'3.3.5': ['76edcde'],
'3.3.6': ['7f90171'],
'3.3.7': ['2e61bcf'],
'4.0.0': ['e82315c'],
'4.0.0-migration': ['3b81223'],

'4.0.1': ['7c5512c'],
'4.0.2': ['7275245'],
'4.0.3': ['935bfbb9'],
'4.0.4': ['3ac21d3'],
'4.1.0': ['0d7b0e3'],
'4.1.1': ['3d26153'],
'4.2.0': ['9a7413d'],
'4.3.0': ['bff61af'],
'4.3.1': ['edd6415'],
'4.3.10': ['3aed64e'],
'4.3.11': ['1724358'],
'4.3.12': ['279f6fd'],
'4.3.13': ['50117e2'],
'4.3.14': ['6b24512'],
'4.3.2': ['806a5b5'],
'4.3.3': ['055de18'],
'4.3.4': ['fab0b1a'],
'4.3.5': ['4468cac'],
'4.3.6': ['835a203'],
'4.3.7': ['b9b9565'],
'4.3.8': ['90954dd'],
'4.3.9': ['bcae6f3'],
'4.4.0': ['8082cf7'],
'4.4.1': ['4b02fda'],
'4.4.10': ['5fe03f3'],
'4.4.11': ['befda65'],
'4.4.12': ['5156154'],
'4.4.13': ['b3c4db9'],
'4.4.14': ['356702d'],
'4.4.15': ['ac1ead0'],
'4.4.16': ['97f9a9d'],
'4.4.17': ['fba21c1'],
'4.4.18': ['435ab36'],
'4.4.19': ['840e648'],
'4.4.2': ['3777fa7'],
'4.4.20': ['32a5933'],
'4.4.21': ['499f9ae'],
'4.4.22': ['cd4ef73'],
'4.4.23': ['6f8b402'],
'4.4.24': ['bccafee'],
'4.4.25': ['ae8746e'],
'4.4.26': ['0f5588e'],
'4.4.27': ['d3e4307'],
'4.4.3': ['53b2eb0'],
'4.4.4': ['a1135ef'],
'4.4.5': ['f6b496b'],
'4.4.6': ['813440d'],

```
'4.4.7': ['1eed61d'],  
'4.4.8': ['2dcf377'],  
'4.4.9': ['1ddba59']}
```

Agora vamos agrupar as tags por versões minor e ordenar as versões patch.

```
[8]: from itertools import groupby  
groups = groupby(tags, lambda x: x.rsplit(".", 1)[0])  
minor_tags = {}  
for minor, elements in groups:  
    minor_tags[minor] = sorted(  
        elements,  
        key=lambda x: int(x.split('-')[0].split('.')[1])  
    )  
minor_tags['4.3']
```

```
[8]: ['4.3.0',  
      '4.3.1',  
      '4.3.2',  
      '4.3.3',  
      '4.3.4',  
      '4.3.5',  
      '4.3.6',  
      '4.3.7',  
      '4.3.8',  
      '4.3.9',  
      '4.3.10',  
      '4.3.11',  
      '4.3.12',  
      '4.3.13',  
      '4.3.14']
```

Fazendo o mesmo para agrupar versões major.

```
[9]: groups = groupby(minor_tags, lambda x: x.rsplit(".", 1)[0])  
major_tags = {}  
for major, elements in groups:  
    major_tags[major] = sorted(  
        elements,  
        key=lambda x: int(x.split('-')[0].split('.')[1])  
    )  
major_tags['4']
```

```
[9]: ['4.0', '4.1', '4.2', '4.3', '4.4']
```

Com isso, podemos escolher versões major (e.g., 3 e 4) e obter a última versão patch para cada minor delas.

```
[10]: last_patch_for_v3v4 = {
        minor: minor_tags[minor][-1]
        for minor in major_tags['3'] + major_tags['4']
    }
last_patch_for_v3v4
```

```
[10]: {'3.0': '3.0.0',
       '3.1': '3.1.0',
       '3.2': '3.2.1',
       '3.3': '3.3.7',
       '4.0': '4.0.4',
       '4.1': '4.1.1',
       '4.2': '4.2.0',
       '4.3': '4.3.14',
       '4.4': '4.4.27'}
```

Agora queremos ver a evolução de linhas de código para as versões selecionadas. Para isso, vamos percorrer o dicionário fazendo checkout de cada versão, carregar o número de linhas usando `cloc` e parsear o resultado para extrair as colunas para construir linhas de uma tabela.

```
[11]: from collections import defaultdict
columns = {"id"}
rows = []

for minor, tag in last_patch_for_v3v4.items():
    __ = !git checkout $tag
    lines = !cloc .
    filtered_lines = lines[lines.index("-" * 79) + 3:]
    commit_result = defaultdict(int)
    commit_result["id"] = minor
    for line in filtered_lines:
        if not line.startswith("-"):
            split = line.split()
            language = split[0]
            commit_result[language + "_files"] = int(split[1])
            commit_result[language + "_blank"] = int(split[2])
            commit_result[language + "_comment"] = int(split[3])
            commit_result[language + "_code"] = int(split[4])
            columns |= {
                language + "_files", language + "_blank",
                language + "_comment", language + "_code"
            }
    rows.append(commit_result)
```

1.2 Pandas

Podemos usar `pandas` para construir a tabela a partir da lista de dicionários.

```
[12]: import pandas as pd
df = pd.DataFrame(rows)
df
```

```
[12]:      id  Ruby_files  Ruby_blank  Ruby_comment  Ruby_code  HTML_files  \
0  3.0         336        1880         1362        19147         108
1  3.1         336        1937         1363        19359         108
2  3.2         364        2437         1385        21739         108
3  3.3         396        2749         1527        23319         109
4  4.0         447        2965         1653        24719         109
5  4.1         449        2977         1667        24759         109
6  4.2         449        2978         1667        24767         109
7  4.3         505        3221         2100        25564           5
8  4.4         508        3186         2042        25765           4
```

```
      HTML_blank  HTML_comment  HTML_code  Sass_files  ...  XML_comment  \
0         10891           22        15740           3  ...           0.0
1         10891           22        15740           3  ...           0.0
2         10891           22        15740           8  ...           0.0
3         10901           22        15887           8  ...           0.0
4         10903           22        15799           8  ...           0.0
5         10903           22        15799           8  ...           0.0
6         10903           22        15799           8  ...           0.0
7            28           22          227           9  ...           NaN
8            16           22          168           9  ...           NaN
```

```
      XML_code  SUM:_files  SUM:_blank  SUM:_comment  SUM:_code  \
0          9.0         508        13520         1525        38849
1          9.0         508        13582         1526        39079
2          9.0         546        14249         1643        42458
3          9.0         583        14604         1790        44486
4          9.0         645        14903         1966        46478
5          9.0         647        14917         1980        46532
6          9.0         647        14918         1980        46540
7          NaN         594         4237         2418        31506
8          NaN         600         4192         2321        31740
```

```
      CoffeeScript_files  CoffeeScript_blank  CoffeeScript_comment  \
0                   NaN                   NaN                   NaN
1                   NaN                   NaN                   NaN
2                   NaN                   NaN                   NaN
3                   NaN                   NaN                   NaN
4                   NaN                   NaN                   NaN
5                   NaN                   NaN                   NaN
6                   NaN                   NaN                   NaN
7                   1.0                   0.0                   3.0
8                   1.0                   0.0                   3.0
```


	CoffeeScript_code
0	NaN
1	NaN
2	NaN
3	NaN
4	NaN
5	NaN
6	NaN
7	0.0
8	0.0

[9 rows x 45 columns]

O pandas permite descrever a tabela com o método `.describe()`.

```
[13]: df.describe()
```

```
[13]:
```

	Ruby_files	Ruby_blank	Ruby_comment	Ruby_code	HTML_files	\
count	9.000000	9.000000	9.000000	9.000000	9.000000	
mean	421.111111	2703.333333	1640.666667	23237.555556	85.444444	
std	66.325795	507.187589	275.371113	2566.055587	45.894202	
min	336.000000	1880.000000	1362.000000	19147.000000	4.000000	
25%	364.000000	2437.000000	1385.000000	21739.000000	108.000000	
50%	447.000000	2965.000000	1653.000000	24719.000000	108.000000	
75%	449.000000	2978.000000	1667.000000	24767.000000	109.000000	
max	508.000000	3221.000000	2100.000000	25765.000000	109.000000	

	HTML_blank	HTML_comment	HTML_code	Sass_files	Sass_blank	...	\
count	9.000000	9.0	9.000000	9.000000	9.000000	...	
mean	8480.777778	22.0	12322.111111	7.111111	319.666667	...	
std	4795.680160	0.0	6874.179268	2.368778	23.690715	...	
min	16.000000	22.0	168.000000	3.000000	277.000000	...	
25%	10891.000000	22.0	15740.000000	8.000000	315.000000	...	
50%	10891.000000	22.0	15740.000000	8.000000	333.000000	...	
75%	10903.000000	22.0	15799.000000	8.000000	334.000000	...	
max	10903.000000	22.0	15887.000000	9.000000	336.000000	...	

	XML_comment	XML_code	SUM:_files	SUM:_blank	SUM:_comment	\
count	7.0	7.0	9.000000	9.000000	9.000000	
mean	0.0	9.0	586.444444	12124.666667	1905.444444	
std	0.0	0.0	55.854971	4516.526874	320.371703	
min	0.0	9.0	508.000000	4192.000000	1525.000000	
25%	0.0	9.0	546.000000	13520.000000	1643.000000	
50%	0.0	9.0	594.000000	14249.000000	1966.000000	
75%	0.0	9.0	645.000000	14903.000000	1980.000000	
max	0.0	9.0	647.000000	14918.000000	2418.000000	

	SUM:_code	CoffeeScript_files	CoffeeScript_blank	\
count	9.000000	2.0	2.0	
mean	40852.000000	1.0	0.0	
std	6016.483711	0.0	0.0	
min	31506.000000	1.0	0.0	
25%	38849.000000	1.0	0.0	
50%	42458.000000	1.0	0.0	
75%	46478.000000	1.0	0.0	
max	46540.000000	1.0	0.0	

	CoffeeScript_comment	CoffeeScript_code
count	2.0	2.0
mean	3.0	0.0
std	0.0	0.0
min	3.0	0.0
25%	3.0	0.0
50%	3.0	0.0
75%	3.0	0.0
max	3.0	0.0

[8 rows x 44 columns]

Além disso, é possível fazer seleções nos dados.

```
[14]: df[df["Ruby_code"] > 25000]
```

```
[14]:
```

	id	Ruby_files	Ruby_blank	Ruby_comment	Ruby_code	HTML_files	\
7	4.3	505	3221	2100	25564	5	
8	4.4	508	3186	2042	25765	4	

	HTML_blank	HTML_comment	HTML_code	Sass_files	...	XML_comment	\
7	28	22	227	9	...	NaN	
8	16	22	168	9	...	NaN	

	XML_code	SUM:_files	SUM:_blank	SUM:_comment	SUM:_code	\
7	NaN	594	4237	2418	31506	
8	NaN	600	4192	2321	31740	

	CoffeeScript_files	CoffeeScript_blank	CoffeeScript_comment	\
7	1.0	0.0	3.0	
8	1.0	0.0	3.0	

	CoffeeScript_code
7	0.0
8	0.0

[2 rows x 45 columns]

1.3 Exercício 10

Selecione as versões que usam CoffeeScript e as versões que não usam XML.

```
[15]: with_coffee = ...  
with_coffee
```

```
[15]:      id  Ruby_files  Ruby_blank  Ruby_comment  Ruby_code  HTML_files  \  
7  4.3         505        3221         2100      25564         5  
8  4.4         508        3186         2042      25765         4  
  
      HTML_blank  HTML_comment  HTML_code  Sass_files  ...  XML_comment  \  
7          28          22        227          9  ...         NaN  
8          16          22        168          9  ...         NaN  
  
      XML_code  SUM:_files  SUM:_blank  SUM:_comment  SUM:_code  \  
7          NaN         594        4237         2418      31506  
8          NaN         600        4192         2321      31740  
  
      CoffeeScript_files  CoffeeScript_blank  CoffeeScript_comment  \  
7              1.0              0.0              3.0  
8              1.0              0.0              3.0  
  
      CoffeeScript_code  
7              0.0  
8              0.0
```

[2 rows x 45 columns]

```
[16]: without_xml = ...  
without_xml
```

```
[16]:      id  Ruby_files  Ruby_blank  Ruby_comment  Ruby_code  HTML_files  \  
7  4.3         505        3221         2100      25564         5  
8  4.4         508        3186         2042      25765         4  
  
      HTML_blank  HTML_comment  HTML_code  Sass_files  ...  XML_comment  \  
7          28          22        227          9  ...         NaN  
8          16          22        168          9  ...         NaN  
  
      XML_code  SUM:_files  SUM:_blank  SUM:_comment  SUM:_code  \  
7          NaN         594        4237         2418      31506  
8          NaN         600        4192         2321      31740
```

	CoffeeScript_files	CoffeeScript_blank	CoffeeScript_comment	\
7	1.0	0.0	3.0	
8	1.0	0.0	3.0	

	CoffeeScript_code
7	0.0
8	0.0

[2 rows x 45 columns]

Além de selecionar linhas, podemos selecionar colunas.

```
[17]: columns = ['SUM:_files', 'SUM:_blank', 'SUM:_comment', 'SUM:_code']
      ndf = df[columns]
      ndf
```

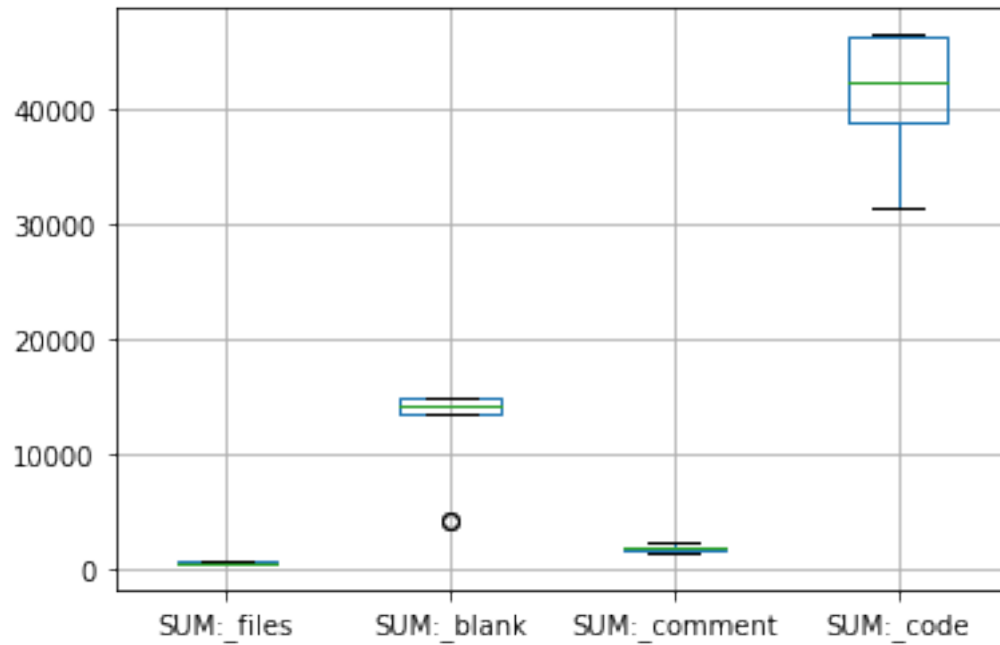
```
[17]:
```

	SUM:_files	SUM:_blank	SUM:_comment	SUM:_code
0	508	13520	1525	38849
1	508	13582	1526	39079
2	546	14249	1643	42458
3	583	14604	1790	44486
4	645	14903	1966	46478
5	647	14917	1980	46532
6	647	14918	1980	46540
7	594	4237	2418	31506
8	600	4192	2321	31740

O pandas também oferece algumas funções que facilitam a geração de gráficos.

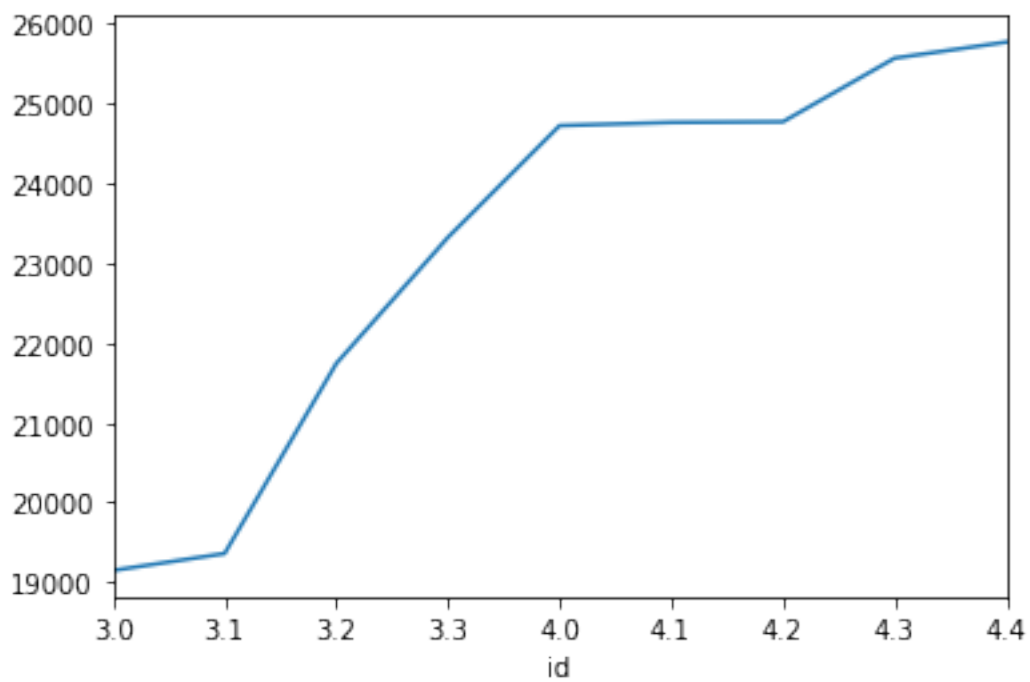
```
[18]: %matplotlib inline
      ndf.boxplot()
```

```
[18]: <matplotlib.axes._subplots.AxesSubplot at 0x7f79bd12df28>
```



```
[19]: df.set_index("id")["Ruby_code"].plot()
```

```
[19]: <matplotlib.axes._subplots.AxesSubplot at 0x7f79bc83bac8>
```



É possível aplicar operações em colunas e criar novas colunas.

```
[20]: df.loc[:, "tag"] = df["id"].apply(lambda minor: last_patch_for_v3v4[minor])
```

```
[21]: df
```

```
[21]:
```

	id	Ruby_files	Ruby_blank	Ruby_comment	Ruby_code	HTML_files	\
0	3.0	336	1880	1362	19147	108	
1	3.1	336	1937	1363	19359	108	
2	3.2	364	2437	1385	21739	108	
3	3.3	396	2749	1527	23319	109	
4	4.0	447	2965	1653	24719	109	
5	4.1	449	2977	1667	24759	109	
6	4.2	449	2978	1667	24767	109	
7	4.3	505	3221	2100	25564	5	
8	4.4	508	3186	2042	25765	4	

	HTML_blank	HTML_comment	HTML_code	Sass_files	...	XML_code	SUM:_files	\
0	10891	22	15740	3	...	9.0	508	
1	10891	22	15740	3	...	9.0	508	
2	10891	22	15740	8	...	9.0	546	
3	10901	22	15887	8	...	9.0	583	
4	10903	22	15799	8	...	9.0	645	
5	10903	22	15799	8	...	9.0	647	
6	10903	22	15799	8	...	9.0	647	
7	28	22	227	9	...	NaN	594	
8	16	22	168	9	...	NaN	600	

	SUM:_blank	SUM:_comment	SUM:_code	CoffeeScript_files	\
0	13520	1525	38849	NaN	
1	13582	1526	39079	NaN	
2	14249	1643	42458	NaN	
3	14604	1790	44486	NaN	
4	14903	1966	46478	NaN	
5	14917	1980	46532	NaN	
6	14918	1980	46540	NaN	
7	4237	2418	31506	1.0	
8	4192	2321	31740	1.0	

	CoffeeScript_blank	CoffeeScript_comment	CoffeeScript_code	tag
0	NaN	NaN	NaN	3.0.0
1	NaN	NaN	NaN	3.1.0
2	NaN	NaN	NaN	3.2.1
3	NaN	NaN	NaN	3.3.7
4	NaN	NaN	NaN	4.0.4
5	NaN	NaN	NaN	4.1.1
6	NaN	NaN	NaN	4.2.0

7	0.0	3.0	0.0	4.3.14
8	0.0	3.0	0.0	4.4.27

[9 rows x 46 columns]

Existem muitas outras operações que podem ser vistas na documentação:
<https://pandas.pydata.org/pandas-docs/stable/>.

Continua: [9.Pygit2.pdf](#)