



IV SEICOM

Conceitos de Data Mining e o uso da linguagem Python para a mineração de dados

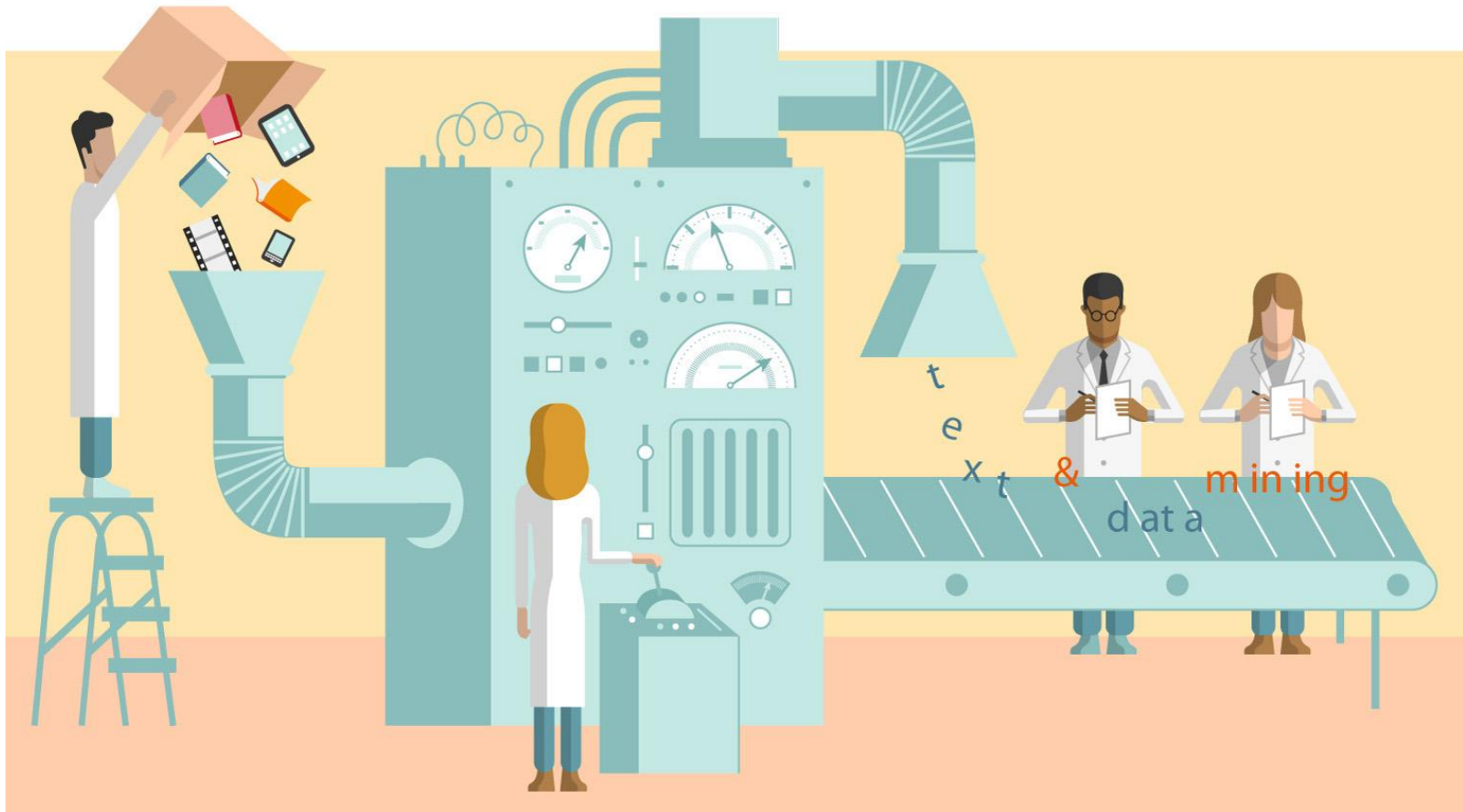
Thiago Medeiros de Souza
Valmei Abreu Junior
João Paulo de Lima Barbosa



Conteúdo

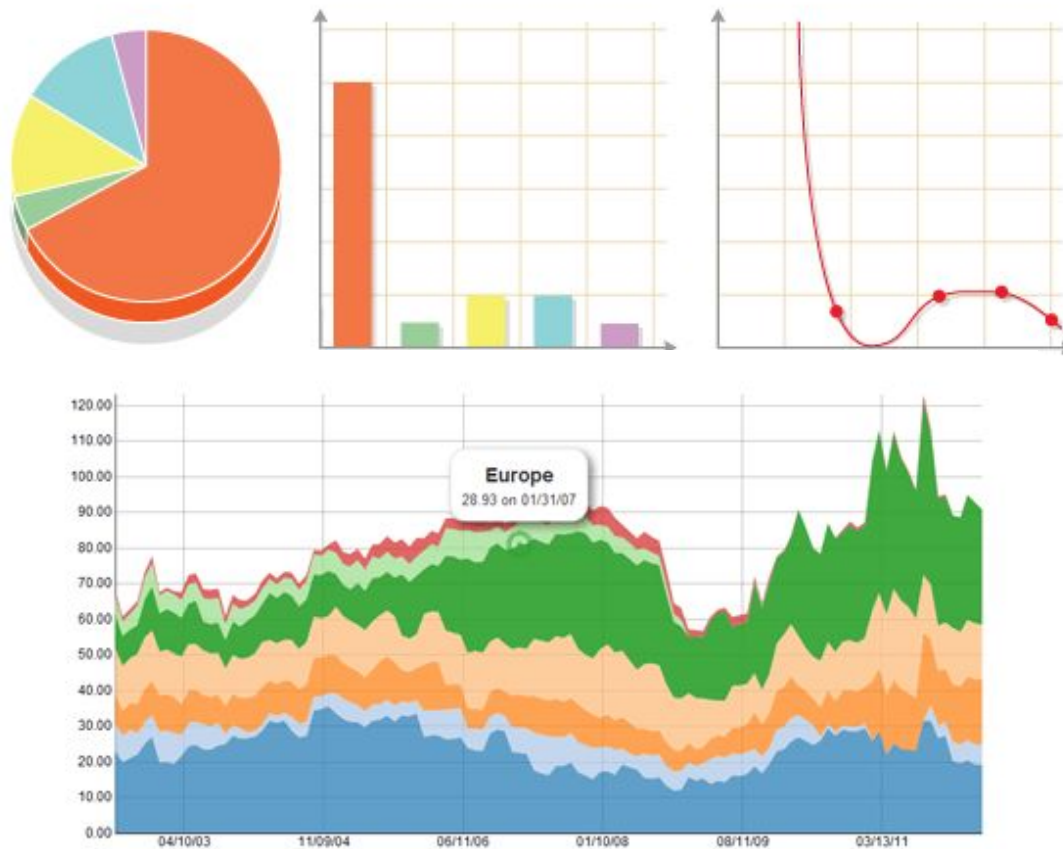
- Introdução
- KDD e Data Mining
- KDD
- Data Mining
- Python
- Casos de Uso
- Considerações Finais

Introdução





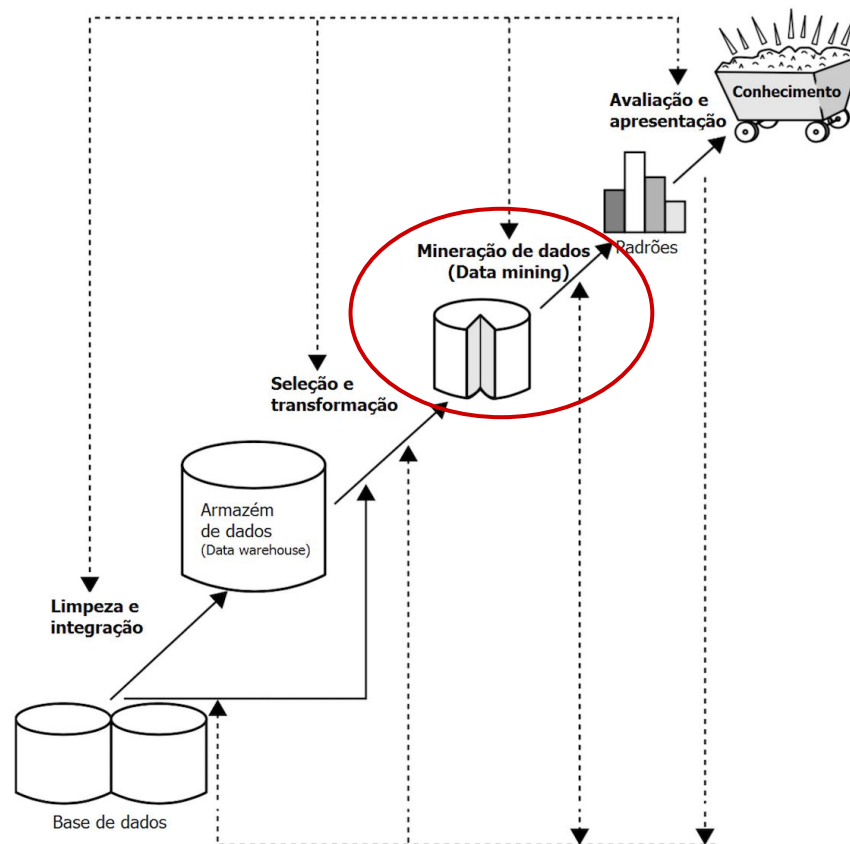
Introdução



KDD e Data Mining



KDD - Knowledge Discovery from Data



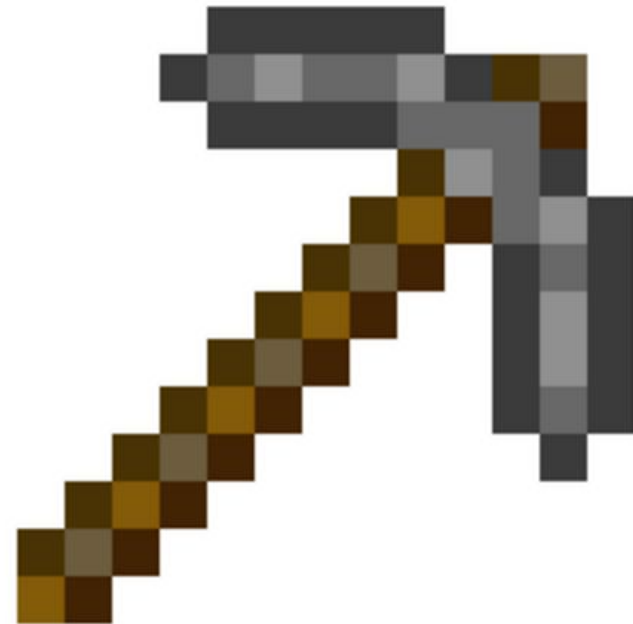
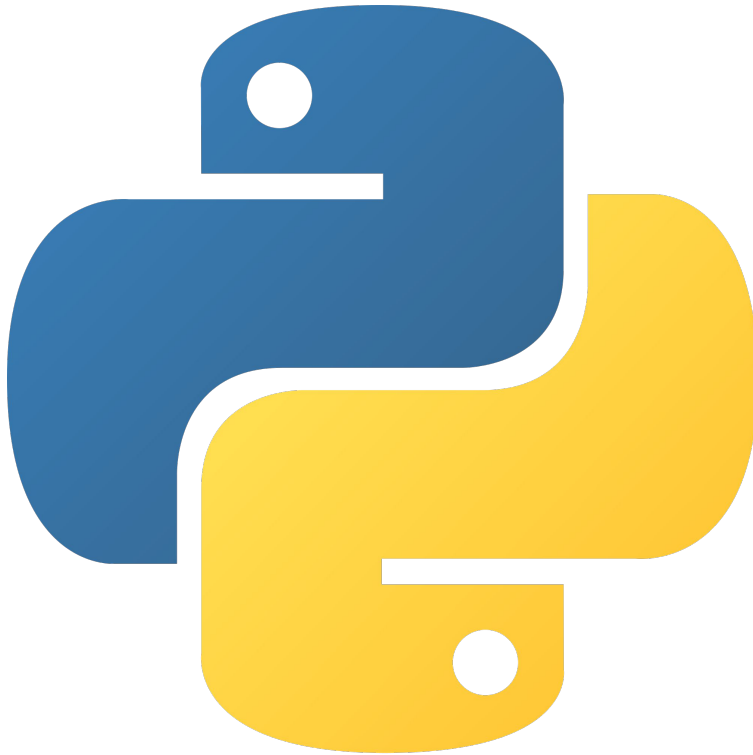
Data Mining

Métodos utilizados em Data Mining:

- Classificação
- Modelos de Relacionamento
- Análise de Agrupamento (Cluster)
- Sumarização
- Modelo de Dependência
- Regras de Associação
- Análise de Séries Temporais



Python



Python

- pandas

```
[{'created_at': '2015-11-06T13:23:50Z', 'trends': [{'url': 'http://twitter.com/search?q=%23PSYPagyanig', 'query': '23PSYPagyanig', 'name': '#PSYPagyanig', 'promoted_content': None}, {'url': 'http://twitter.com/search?q=%237DaysUntilMITAM', 'query': '237DaysUntilMITAM', 'name': '#7DaysUntilMITAM', 'promoted_content': None}, {'url': 'http://twitter.com/search?q=%23BuenViernes', 'query': '23BuenViernes', 'name': '#BuenViernes', 'promoted_content': None}, {'url': 'http://twitter.com/search?q=%23EB%8B%89%EB%84%A4%EC%9E%84%EC%9D%84_00%EC%9D%B4%EB%9D%BC_%EC%A7%80%EC%9D%80_%EC%9D%B4%EC%9C%A0', 'query': '23EB%8B%89%EB%84%A4%EC%9E%84%EC%9D%84_00%EC%9D%B4%EB%9D%BC_%EC%A7%80%EC%9D%80_%EC%9D%B4%EC%9C%A0', 'name': '#\ub2c9\ubl24\uc784\uc744_00\uc774\ub77c_\uc9c0\uc740_\uc774\uc720', 'promoted_content': None}, {'url': 'http://twitter.com/search?q=%23ManOnTheMoon', 'query': '23ManOnTheMoon', 'name': '#ManOnTheMoon', 'promoted_content': None}, {'url': 'http://twitter.com/search?q=%E3%82%8B%E3%82%8D%E5%89%A3', 'query': 'E3%82%8B%E3%82%8D%E5%89%A3', 'name': 'u\308b\u308d\u5263', 'promoted_content': None}, {'url': 'http://twitter.com/search?q=Coldplay', 'query': 'Coldplay', 'name': 'Coldplay', 'promoted_content': None}, {'url': 'http://twitter.com/search?q=%22DIA+13+%C3%89+DIA+DE+JUSTIN+BIEBER%22', 'query': '22DIA+13+%C3%89+DIA+DE+JUSTIN+BIEBER%22', 'name': 'DIA 13 \xc9 DIA DE JUSTIN BIEBER', 'promoted_content': None}, {'url': 'http://twitter.com/search?q=%E6%9D%BE%E3%81%95%E3%82%93%E8%A8%BA%E6%96%AD', 'query': 'E6%9D%BE%E3%81%95%E3%82%93%E8%A8%BA%E6%96%AD', 'name': 'u\u677e\u3055\u3093\u8a3a\u65ad', 'promoted_content': None}, {'url': 'http://twitter.com/search?q=%E3%82%A2%E3%83%8B%E3%83%A1%E3%82%AD%E3%83%A3%E3%83%A9%E5%8C%96', 'query': 'E3%82%A2%E3%83%8B%E3%83%A1%E3%82%AD%E3%83%A3%E3%83%A9%E5%8C%96', 'name': 'u\u30a2\u30cb\u30e1\u30ad\u30e3\u30e9\u5316', 'promoted_content': None}], 'as_of': '2015-11-06T13:30:05Z', 'locations': [{'woeid': 1, 'name': 'Worldwide'}]}}
```

```
[{'created_at': '2015-11-06T13:28:41Z', 'trends': [{'url': 'http://twitter.com/search?q=%23CarsonOnCNN', 'query': '23CarsonOnCNN', 'name': '#CarsonOnCNN', 'promoted_content': None}, {'url': 'http://twitter.com/search?q=%237DaysUntilMITAM', 'query': '237DaysUntilMITAM', 'name': '#7DaysUntilMITAM', 'promoted_content': None}, {'url': 'http://twitter.com/search?q=%23GetWeirdOutNow', 'query': '23GetWeirdOutNow', 'name': '#GetWeirdOutNow', 'promoted_content': None}, {'url': 'http://twitter.com/search?q=%23FlavorAFilm', 'query': '23FlavorAFilm', 'name': '#FlavorAFilm', 'promoted_content': None}, {'url': 'http://twitter.com/search?q=%23FridayFeeling', 'query': '23FridayFeeling', 'name': '#FridayFeeling', 'promoted_content': None}, {'url': 'http://twitter.com/search?q=%22Sagamore+Bridge%22', 'query': '22Sagamore+Bridge%22', 'name': 'Sagamore Bridge', 'promoted_content': None}, {'url': 'http://twitter.com/search?q=%22Finally+Friday%22', 'query': '22Finally+Friday%22', 'name': 'Finally Friday', 'promoted_content': None}, {'url': 'http://twitter.com/search?q=Coldplay', 'query': 'Coldplay', 'name': 'Coldplay', 'promoted_content': None}, {'url': 'http://twitter.com/search?q=Chaffee', 'query': 'Chaffee', 'name': 'Chaffee', 'promoted_content': None}, {'url': 'http://twitter.com/search?q=%22Art+Modell%22', 'query': '22Art+Modell%22', 'name': 'Art Modell', 'promoted_content': None}], 'as_of': '2015-11-06T13:30:05Z', 'locations': [{'woeid': 23424977, 'name': 'United States'}]}}
```

Python

- pandas

```
{
  "created_at": "2015-11-06T13:23:50Z",
  "trends": [
    {
      "url": "http://twitter.com/search?q=%23PSYPagyanig",
      "query": "%23PSYPagyanig",
      "name": "#PSYPagyanig",
      "promoted_content": null
    },
    {
      "url": "http://twitter.com/search?q=%237DaysUntilMITAM",
      "query": "%237DaysUntilMITAM",
      "name": "#7DaysUntilMITAM",
      "promoted_content": null
    },
    {
      "url": "http://twitter.com/search?q=%23BuenViernes",
      "query": "%23BuenViernes",
      "name": "#BuenViernes",
      "promoted_content": null
    },
    {
      "url": "http://twitter.com/search?q=%23%EB%8B%89%EB%84%A4%EC%9E%84%EC%9D%84_00%EC%9D%B4%EB%9D%BC_%EC%A7%80%EC%9D%80_%EC%9D%B4%EC%9C%A0",
      "query": "%23%EB%8B%89%EB%84%A4%EC%9E%84%EC%9D%84_00%EC%9D%B4%EB%9D%BC_%EC%A7%80%EC%9D%80_%EC%9D%B4%EC%9C%A0",
      "name": "#\ub2c9\ub124\uc784\uc744_00\uc774\ub77c_\uc9c0\uc740_\uc774\uc720",
      "promoted_content": null
    },
    {
      "url": "http://twitter.com/search?q=%23ManOnTheMoon",
      "query": "%23ManOnTheMoon",
      "name": "#ManOnTheMoon",
      "promoted_content": null
    }
  ]
}
```

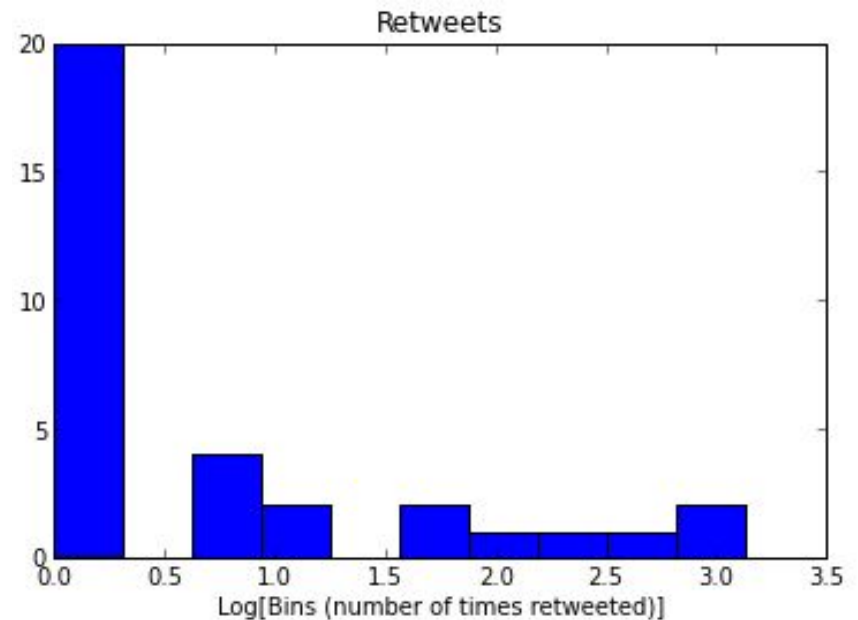
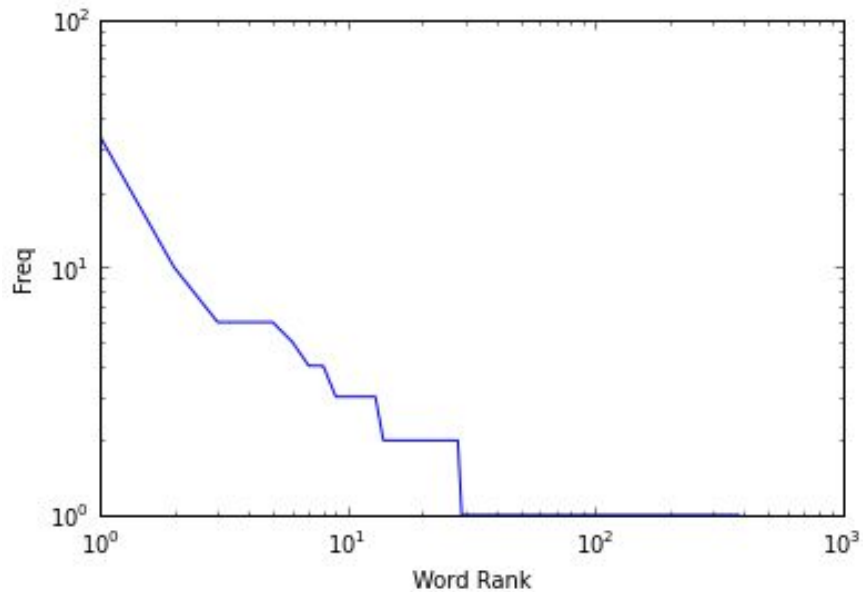
Python

- NumPy

```
>>> a = np.ones((3, 3)) # reminder: (3, 3) is a tuple
>>> a
array([[ 1.,  1.,  1.],
       [ 1.,  1.,  1.],
       [ 1.,  1.,  1.]])
>>> b = np.zeros((2, 2))
>>> b
array([[ 0.,  0.],
       [ 0.,  0.]])
>>> c = np.eye(3)
>>> a = np.random.rand(4) # uniform in [0, 1]
>>> a
array([ 0.95799151,  0.14222247,  0.08777354,  0.51887998])
>>> b = np.random.randn(4) # Gaussian
>>> b
array([ 0.37544699, -0.11425369, -0.47616538,  1.79664113])
>>> np.random.seed(1234) # Setting the random seed
```

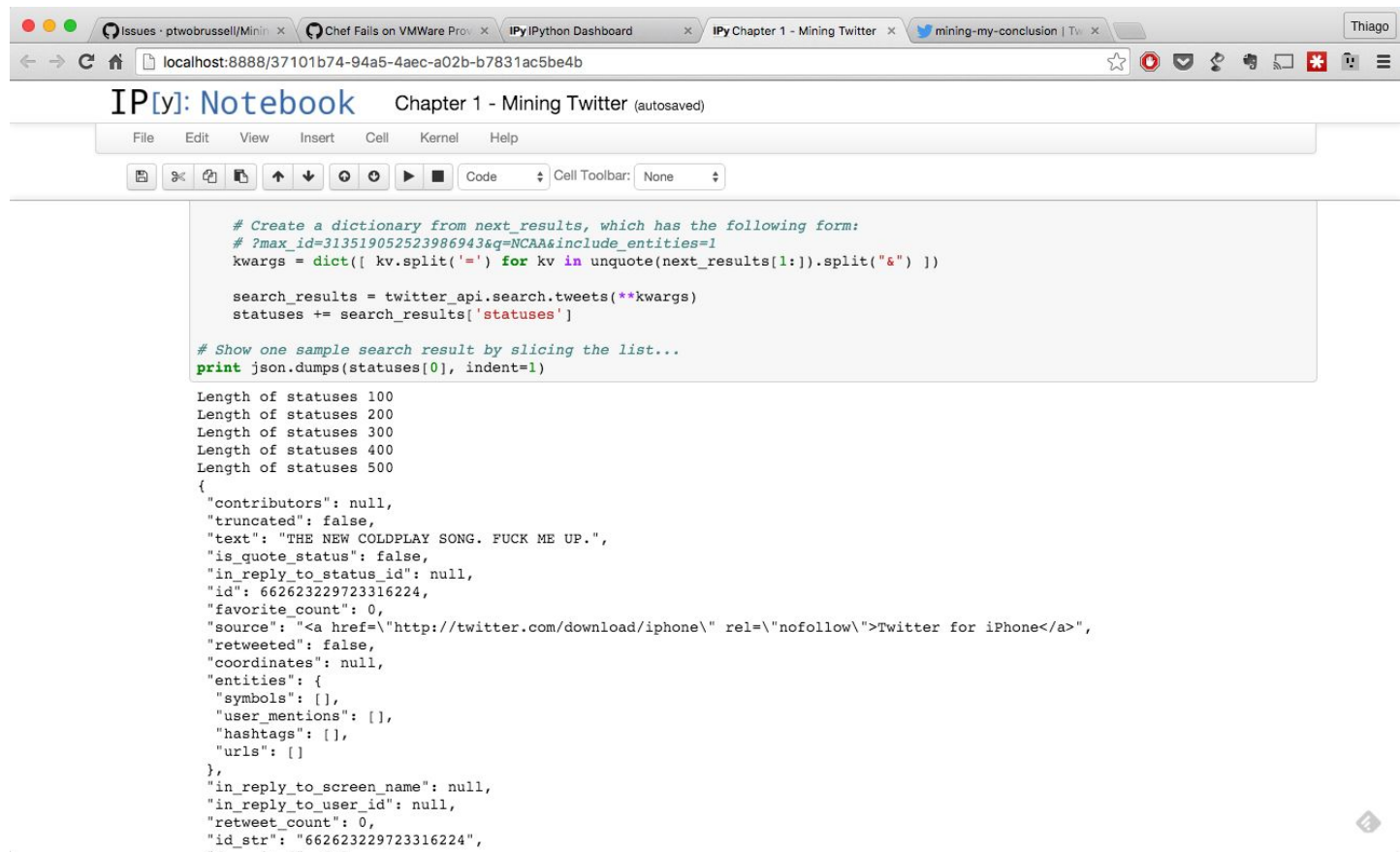

Python

- matplotlib



Python

- IPython



The screenshot shows a web browser window with multiple tabs. The active tab is titled "IPython Notebook Chapter 1 - Mining Twitter (autosaved)". The address bar shows the URL "localhost:8888/37101b74-94a5-4aec-a02b-b7831ac5be4b". The notebook interface includes a menu bar (File, Edit, View, Insert, Cell, Kernel, Help) and a toolbar with icons for file operations and execution. The main area contains a Python code cell with the following text:

```
# Create a dictionary from next_results, which has the following form:
# ?max_id=313519052523986943&q=NCAA&include_entities=1
kwargs = dict([ kv.split('=') for kv in unquote(next_results[1:]).split("&") ])

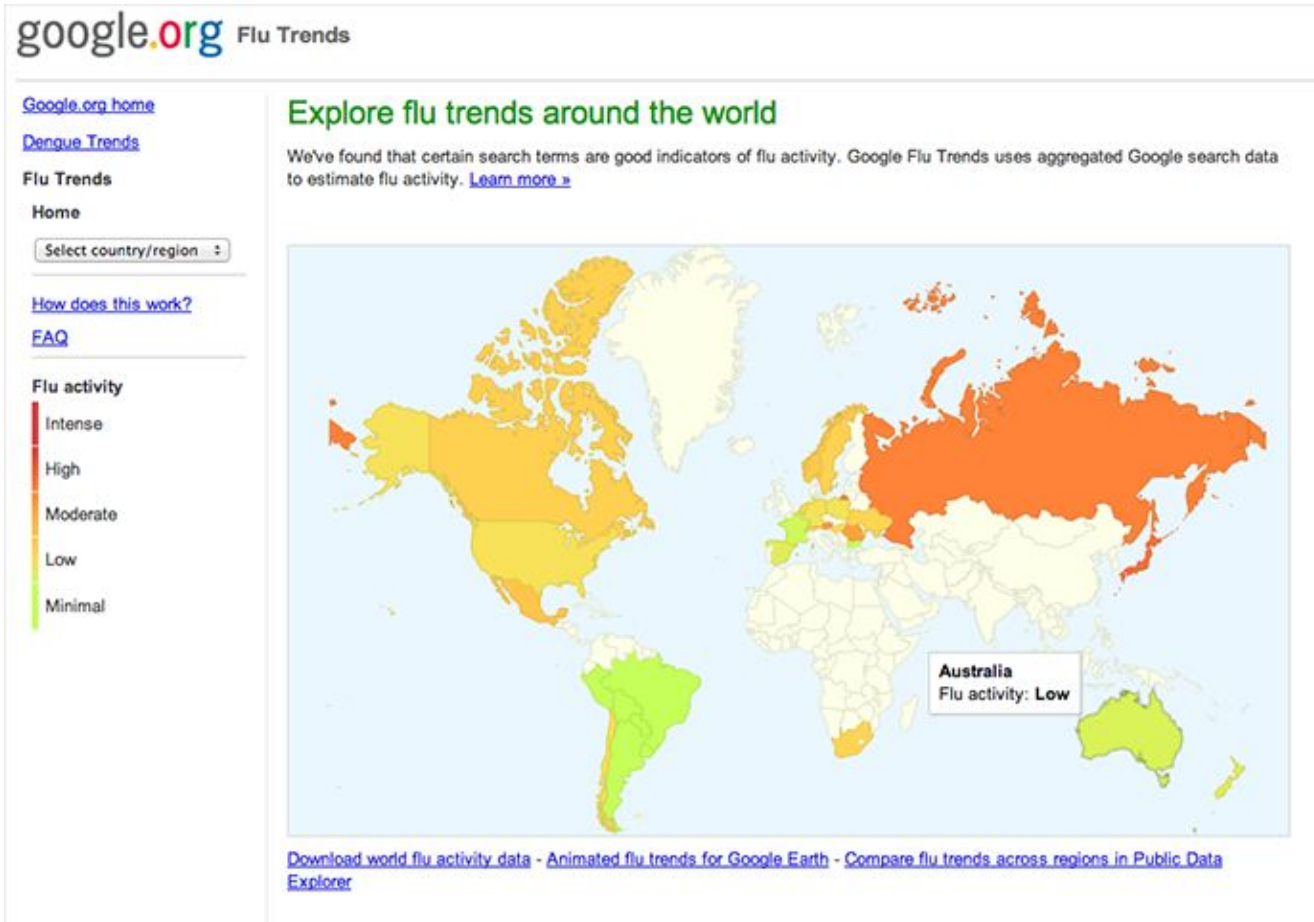
search_results = twitter_api.search.tweets(**kwargs)
statuses += search_results['statuses']

# Show one sample search result by slicing the list...
print json.dumps(statuses[0], indent=1)

Length of statuses 100
Length of statuses 200
Length of statuses 300
Length of statuses 400
Length of statuses 500
{
  "contributors": null,
  "truncated": false,
  "text": "THE NEW COLDPLAY SONG. FUCK ME UP.",
  "is_quote_status": false,
  "in_reply_to_status_id": null,
  "id": 662623229723316224,
  "favorite_count": 0,
  "source": "<a href='\"http://twitter.com/download/iphone\"' rel='\"nofollow\"'>Twitter for iPhone</a>",
  "retweeted": false,
  "coordinates": null,
  "entities": {
    "symbols": [],
    "user_mentions": [],
    "hashtags": [],
    "urls": []
  },
  "in_reply_to_screen_name": null,
  "in_reply_to_user_id": null,
  "retweet_count": 0,
  "id_str": "662623229723316224",
  ...
```



Casos de Uso



Considerações Finais

