# Sistema de Alerta Processual e Avaliação de Unidades Judiciais

## Tecnologias e ferramentas utilizadas

A linguagem utilizada para a parte do sistema que realiza o processamento dos dados, clusterização, avaliação e alertas foi o Python. Tudo que compõe esse sistema, das bibliotecas de processamento de dados ao modelo de cluster, estão identificadas abaixo:

```
import numpy as np
import pandas as pd
import seaborn as sns
import networkx as nx
import matplotlib.pyplot as plt
import json
import datetime
import pickle
import operator

from google.colab import drive
from unidecode import unidecode
from matplotlib.cm import ScalarMappable
from networkx.drawing.nx_agraph import graphviz_layout
import matplotlib.image as mpimg
from sklearn.preprocessing import normalize
from sklearn.cluster import Birch
```

## Informações sobre os clusters

O principal ponto do desafio, e também alicerce do sistema desenvolvido, é a clusterização. Para isso, optamos por utilizar o modelo <u>Birch</u>, na implementação do scikit-learn.

O motivo da escolha é a flexibilidade do modelo Birch. Sendo implementado a partir de uma estrutura de árvore, o modelo tem suporte ao *online training* (treino em etapas, utilizando a função birch.fit(X)) e é ideal para situações onde o

conjunto de dados é grande e há o risco de outliers. Tais informações partem do próprio <u>scikit-learn</u>.

#### Métricas do Modelo

O modelo BIRCH faz parte da classe de algoritmos não-supervisionados, de modo que não é uma tarefa trivial (e nem mesmo comum) estabelecer métricas de avaliação do modelo. O fluxo padrão em algoritmos de clustering é criar os clustering e realizar análises, o que foi feito.

Como dito, o BIRCH foi escolhido por ser extremamente eficiente ao lidar com grandes quantidades de dados, como dito, inclusive, no seu <u>artigo de origem</u>. Essa eficiência é alcançada ao se construir um sumário das informações e distribuições dos dados e executando a clusterização sobre esses dados.

O funcionamento do Birch se torna ainda mais interessante ao fornecer a opção de especificar um valor de *threshold*, que determina o raio dos subclusters que o algoritmo monta, de modo que há uma certa flexibilidade na avaliação dos clusters.

### Áreas do Conhecimento e Técnicas Envolvidas

A solução envolve, principalmente, conhecimentos e técnicas relacionadas ao tratamento de dados. Como a base de dados do DATAJUD tinha bastante problemas de padronização, padronizar e tratar os dados foram tarefas que consumiram bastante tempo.

Além disso, envolveu-se conhecimentos em inteligência artificial, para o uso e aplicação do modelo de aprendizado não-supervisionado.

Por fim, há conhecimentos em webdevelopment, necessários para a criação e configuração da interface que exibe os dados gerados pela solução.

## Formato dos dados

#### **Processos**

Os processos prontos, cujos dataframes estão no repositório sob o formato .pkl, possuem as seguintes colunas:

- *millisInsercao*: milissegundo de inserção (data-hora em formato UNIX)
- grau
- siglaTribunal
- valorCausa
- nivelSigilo
- numero

- codigoOrgaoJulgador
- procEl: Indicador se o processo tramita em meio eletrônico ou físico
- classeProcessual
- tamanhoProcesso
- dscSistema: Descritor do sistema onde o processo tramita
- competencia: identificador da competência a que pertence o processo, ou da competência a que ele se destina caso se trate de processo inicial
- codigoLocalidade
- totalAssuntos
- dataAjuizamento
- movimentoNacional: Lista de movimentos Nacionais, cada elemento é composto por [codigoMovimento,dataHora]
- movimentoLocal: Lista de movimentos Locais, cada elemento é composto por [codigoMovimento,codigoPaiNacional,dataHora]
- assuntoNacional: Lista de assuntos Nacionais, cada elemento é composto por [codigo,dataHora]
- assuntoLocal: Lista de assuntos Locais, cada elemento é composto por [codigo,codigoPaiNacional,dataHora]
- tempo\_proc\_dias: tempo do processo, em dias, desde a dataAjuizamento ao último movimento
- data\_ultimo\_mov: data do último movimento do processo.
- total\_movs\_loc: total de movimentos locais
- total movs nac: total de movimentos nacionais
- total movs: total de movimentos (soma de locais e nacionais)
- 1\_mov\_nac\_freq: movimento nacional mais frequente
- 2\_mov\_nac\_freq: segundo movimento nacional mais frequente
- desc mov1: descrição do movimento nacional mais frequente
- desc mov2: descrição do segundo movimento nacional mais frequente
- *tempoMedioEntreMovimentos*: tempo médio entre os movimentos, calculado dividindo o tempo total do processo pela quantidade total de movimentos
- cluster processo: id do cluster do processo
- pontuacao\_processo: pontuação do processo, com cálculo explicado abaixo

Imagem do formato dos dados (algumas colunas foram retiradas para melhorar a visualização):

grau	siglaTribunal	numero	totalAssuntos	dataAjuizamento	tempo_proc_dias	data_ultimo_mov	total_movs_loc	total_movs_nac	total_movs	1_mov_nac_freq	2_mov_nac_freq	desc_mov1	desc_mov2	tempoMedioEnt
		00717005820055170003		18446681937950751616	10916	20190826145155								
	TRT17	00401001420085170003		18446681937950751616	10396	20190823153928		0.0						
		00186008620085170003		18446681937950751616	10396	20190823175437								
	TRT17	02041017819925170002		18446742238930339616	7218	20130118163227	50.0	0.0	50.0					
		01317019519945170002		18446742238930339616		20160330151926								
	TRT17	01863014219895170002		623646000000	5383	20090817130708		0.0						
		00352017219915170001		668314800000	6346	20130308182207								
	TRT17	00353012719915170001		668314800000	5782	20120815124131		0.0						
		01102011019935170001		734842800000		20110222094627								
	TRT17	01708019119935170002		740977200000	6958	20180502154156		0.0						
10 1		01533011319935170131		785383200000		19960425001700								
11 1	TRT17	00059015419955170121		814413600000	5431	20101206135101		0.0						
12 1		00760013419955170121		845431200000	8652	20191010160650								
13 1	TRT17	01073008419985170101		895633200000	7824	20170822145126	45.0	0.0	45.0					
14 1		00611014619955170121		924058800000		20040618114315								
15 1		01778009420125170131		937278000000	8267	20191111154940						Expedição de documento		
16 1		03277003020015170005		1006740000000		20170831131454								
17 1	TRT17	00847000320035170131		1052449200000	8233	20200114105546		24.0	26.0	11010		Mero expediente	Mero expediente	
18 1		00834006920045170131		1085367600000	7945	20191127100600								
19 1	TRT17	00825004920055170132		1096340400000	7860	20191114144903		0.0						
talAssunt	os dataA	juizamento tempo_proc_	dias data_ultin	no mov total movs loc	total_movs_nac	total_movs 1_mov	/_nac_freq 2_mo	/_nac_freq desc	mov1 desc m	ov2 tempoMedio	EntreMovimentos	cluster_pr	ocesso pon	tuacao_proce

totalAssuntos	dataAjuizamento	tempo_proc_dias	data_ultimo_mov	total_movs_loc	total_movs_nac	total_movs	1_mov_nac_freq	2_mov_nac_freq	desc_mov1	desc_mov2	tempoMedioEntreMovimentos	cluster_processo	pontuacao_processo
0	18446681937950751616	10916	20190826145155								10916.000000		368.050921
0	18446681937950751616	10396	20190823153928								5198.000000	3	257.317193
0	18446681937950751616		20190823175437								10396.000000		349.589383
0	18446742238930339616	7218	20130118163227	50.0		50.0					144.360000	3	111.193406
0	18446742238930339616										4157.000000		201.879324
0	623646000000	5383	20090817130708	15.0							358.866667	3	82.427253
0	668314800000	6346	20130308182207										130.701809
0	668314800000	5782	20120815124131								2891.000000		134.459205
0	734842800000		20110222094627								2330.500000		104.610093
0	740977200000	6958	20180502154156								2319.333333		145.187016
0	785383200000		19960425001700										53.146057
0	814413600000	5431	20101206135101								1357.750000		101.010980
0	845431200000	8652	20191010160650								1236.000000		156.027253
0	895633200000	7824	20170822145126								173.866667		122.474590
0	924058800000		20040618114315								28.968421		29.866929
0	937278000000	8267	20191111154940						Expedição de documento		1181.000000		148.216602
0	1006740000000		20170831131454										114.148174
0	1052449200000		20200114105546				11010		Mero expediente	Mero expediente	316.653846		274.924323
0	1085367600000	7945	20191127100600										121.911223
0	1096340400000	7860	20191114144903								7860.000000	3	259.553880

Este formato foi escolhido como padrão para lidar com as tabelas massivas de dados, devido a se tratar de uma <u>notação de objetos de Python</u>, assim como o JSON é para o JavaScript. A vantagem é que, enquanto arquivos JSON são lidos em Python como strings ou dicionários, ocupando muito espaço tanto em disco quanto em memória RAM, o Pickle é uma tradução binária de objetos.

Desta forma, nosso fluxo de trabalho envolveu ler os arquivos JSON, converter as colunas para objetos de tamanhos menores (uint, em sua maioria), armazená-los em uma estrutura de dados de consulta bem mais rápida (Pandas DataFrames), e persistir essas estruturas em Pickles, que posteriormente foram ainda compactados (em .zip, individualmente) para possibilitar o upload no GitHub para a submissão.

Isso nos permitiu sair de 673 MB dos arquivos originais zipados, para 13.5 GB combinados dos arquivos JSON, para um tamanho final, após todas as análises e operações nos dados, de aproximadamente 1.3 GB.

Isso seria simplesmente impossível caso tivéssemos permanecido com os arquivos JSON, além de dificultar muito a análise, como pode ser visto nas células que tratam o TRT 5, cujos arquivos eram grandes demais para caber em um só DataFrame em memória, e ele permaneceu dividido em 2 por grande parte do desafio, sendo combinado novamente apenas no início da etapa de Clustering,

quando ele já havia sido bem reduzido através dos métodos e operações mencionados acima.

### Pontuação de Processos

Após realização de processamento nos dados (limpeza e formatação) e também engenharia de recursos (criação de novas colunas), a clusterização de processos foi feita. Para isso, levou-se em conta as seguintes colunas:

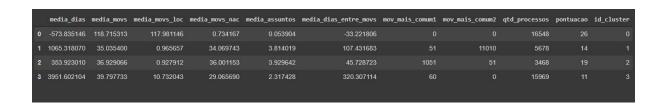
- grau
- totalAssuntos
- tempo\_proc\_dias
- total\_movs\_loc
- total\_movs\_nac
- total\_movs
- 1\_mov\_nac\_freq
- 2 mov nac freq
- tempoMedioEntreMovimentos

A pontuação do processo leva em conta:

- Pontuação do cluster: os processos foram divididos em 4 clusters. Após essa divisão feita pelo modelo Birch, nós analisamos os clusters e montamos estatísticas sobre seu comportamento, que estão disponíveis no repositório.
- Desvio de tempo entre movimentos: compara o tempo entre movimentos do processo com o tempo médio entre movimentos daquele cluster
- Desvio de tempo total: compara o tempo total do processo com o tempo médio total de processos daquele cluster.

A imagem abaixo mostra esse arquivo, contendo estatísticas do TRT1, que são usadas para computar as pontuações. Duas observações:

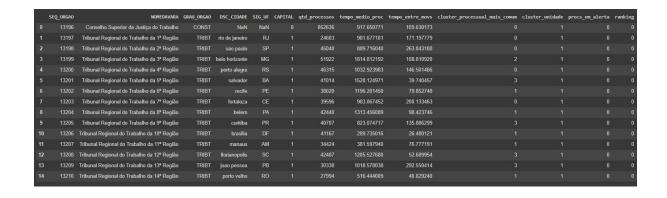
- O BIRCH foi treinado com dados de **todos** os TRTs e o TST, de modo que os clusters são similares em todos os DataFrames, por mais que os dados variem um pouco;
- As médias negativas nos campos relacionados a dias são um problema que enfrentamos devido às datas em formato UNIX e ISO que se encontram nos dados. Pela lógica do nosso código, há uma inversão de sinal quando o formato é ISO. Isso não influencia o código, os clusters ou as conclusões, pois:
  - a) o BIRCH separou de modo perfeito todas as médias negativas e
  - b) o cálculo de média é feito por cluster.



### Unidades Judiciais

Já as unidades judiciais, também presentes no repositório sob o formato .csv, possuem as seguintes colunas:

- ranking: ranking da unidade com base na pontuação.
- SEQ ORGAO
- NOMEDAVARA
- SEQ ORGAO PAI
- GRAU\_ORGAO
- SEQ CIDADE
- DSC CIDADE
- SIG UF
- COD IBGE
- DSC\_TIP\_ORGAO
- TIP ESFERA JUSTICA
- INT\_ORDEM\_ORGAO
- LATITUDE
- LONGITUDE
- CAPITAL: informação se a unidade está localizada em uma capital ou não.
- qtd processos: quantidade de processos atualmente naquela unidade
- tempo medio proc: tempo, em dias, médio dos processos
- tempo\_entre\_movs: tempo médio, em dias, entre movimentações dos processos
- cluster\_processual\_mais\_comum: id do cluster mais comum dos processos daquela unidade (calculado por moda)
- cluster unidade: id do cluster da unidade
- procs em alerta: número de processos em alerta
- pontuacao\_unidade: pontuação da unidade



## Pontuação das Unidades

O cluster de unidades foi feito após o cluter de processos, sendo que o modelo foi alimentado com os seguintes dados:

- CAPITAL
- qtd\_processos
- tempo entre movs
- cluster\_processual\_mais\_comum

### A pontuação das unidades:

- Leva em conta a pontuação do cluster mais comum dos processos;
- Leva em conta a média de dias da unidade em relação à média de dias do cluster:
- Leva em conta a média de dias entre movimentos da unidade em relação à média de dias do cluster;
- Leva em conta a quantidade de processos na unidade.

A imagem do arquivo contendo os detalhes dos clusters.

	media_dias	media_dias_entre_movs	qtd_media_procs	pontuacao_cluster_proc	pontuacao	id_cluster
0	-948.314154	-147.343819	79.875000	25	48	0
1	1387.773044	99.756838	1555.184802	16	27	1
2	6.382353	0.053317	1.176471	21	40	2
3	2284.378721	484.818272	159.532775	8	15	3

Assim, a metodologia de avaliação é a mesma: os aspectos avaliados são os mesmos para todas as unidades e todos os processos. Entretanto, tem-se a ponderação de cluster: os processos e unidades só são comparados com

**similares**. Desse modo, as informações e os rankings são mais valiosos e refletem de maneira mais fidedigna a realidade.

# Licenças

Utilizamos apenas bibliotecas de licença aberta. Todas, sem exceção, permitem uso livre e, inclusive, comercial.

Não foram utilizados programas e ferramentas que possuem licenças pagas, assim como não foram utilizadas soluções de *cloud*.