

# Trabalho Prático 1

2º semestre de 2019

## I. INTRODUÇÃO

Nos últimos anos, a alta popularidade de plataformas de *microblogging* e redes sociais, como Twitter e Facebook, tem se tornado um meio-chave de comunicação tanto na *World Wide Web* (WWW). Por dia, usuários de serviços como o Twitter, são capazes de gerar centenas de milhões de tweets. Aliado a isso, a era da Internet das Coisas (IoT), nos permite ter acesso a grande quantidade de dispositivos móveis capazes de produzir dados referenciados em localização e tempo.

A análise desses conteúdos é interessante em um grande número de aplicações, a soma de todos os textos de redes como o Twitter (tweets) pode, coletivamente, ser considerada uma fonte de informação sobre opiniões e sentimentos sobre produtos, política, sociedade e eventos. Como o Twitter é atualmente a plataforma de *microblogging* com o maior número de usuários ativos, diversos trabalhos vem sendo produzidos sobre a plataforma. Exemplos incluem previsão de receitas de bilheteria para filmes, flutuações dos mercados bolsistas, surtos de doenças como gripe e dengue, e mesmo eleições políticas.

Dados geo-referenciados têm especial importância pela oportunidade para vincular texto a locais e datas. A ideia de extrair padrões geo-temporais significativos a partir de textos de microblog relacionando a um contexto, permite a uma série de conhecimentos úteis para cientistas de dados de diversas áreas. Por exemplo, departamentos públicos de uma cidade, podem ter interesse em acompanhar um evento/show na cidade ou mapear surtos de alguma doença.

## II. OBJETIVOS

Neste trabalho prático, você está encarregado de analisar assuntos populares no Twitter. Você buscará analisar como eles se comportam espaço-temporalmente, por exemplo: Há determinados assuntos que ficaram restritos a uma determinada região (p.ex., delimitado por um raio)? Como é a popularidade desses assuntos pelo tempo? É possível encontrar algum evento que ocorreu na região? Use sua criatividade.

## III. BASE DE DADOS

A base de dados consiste em extrações de tweets geo-referenciados do Twitter nos meses de março à junho de 2016, cada mês com aproximadamente 600 mil tweets e totalizando aproximadamente 2.4 milhões de tweets. O alvo foi a cidade de Curitiba e precisão da localização para esses tweets depende se usou-se GPS ou outra forma de localização. No pior caso, a precisão a ser considerada será de cidade.

O formato dos dados é o formato tradicional de um tweet<sup>1</sup>, armazenado como um JSON por linha, contendo campos como: nome do usuário, id do usuário, data/hora de coleta do tweet, coordenadas (caso exista) de onde foi feito o tweet, o nome da cidade do tweet, o texto e as *hashtags* utilizadas no tweet. Devido a constante atualização da API do Twitter, pode haver

---

<sup>1</sup>Para mais informações: <https://dev.twitter.com/overview/api/tweets>

algumas diferenas entre os campos da listas na documentao e a base coletada. Por exemplo, na base coletada, a data/hora de envio de um determinado tweet podem ser verificados nos campos `created_at` e `timestamp_ms`, ambos representados em *Unix Time*. Em Python, o usurio pode converter-los em um formato convencional utilizando o comando `fromtimestamp` do modulo de `datetime`.

A base encontra-se disponvel no cluster da disciplina em `hdfs:///datasets/geo_curitiba`, nessa pasta, todos os alunos tem acesso leitura. No entanto, cada aluno tem sua pasta individual em `hdfs://user/login`, onde poder salvar seus resultados.

#### IV. ORIENTAES

A identificao de tpicos no Twitter  rea bastante estudada. Para essa disciplina, voc poder escolher qual estratgia utilizar poder ser, por exemplo, a partir de um processamento de texto (no campo `text`) ou por *hashtags* (no campo `entities`).

Dependendo da preciso da localizao, o geo-referenciamento de um tweet pode ser um ponto no formato (longitude, latitude) no campo `coordinates` ou um polgono, representando um permetro, no campo `place`. Voc  livre para decidir qual abordagem tomar, uma delas, seria converter o polgono pelo seu ponto central.

As coordenadas de latitude e longitude esto representadas em formato de graus decimais, popular em coordenadas de GPS.  interessante utilizar a formula de Haversine<sup>2</sup> para se calcular distncias entre coordenadas, alm do resultado ser em padro de metros, leva-se em considerao a curvatura da Terra.

#### V. DOCUMENTAO E PARMETROS DE AVALIAO

Dever ser escrito um relatrio em que  explicado as anlises, os resultados produzidos, como cada anlise foi construda em Spark e quais suas premissas (caso exista). No precisam colocar o cdigo completo no relatrio, apenas trechos para ajudar na sua ilustrao. O aluno dever enviar um nico arquivo compactado contendo o(s) cdigo(s)-fonte(s) e o relatrio produzido. O aluno  livre para escolher qual linguagem utilizar (Scala, Python ou Java). Caso escolham Java, dever ser enviado tambm as instrues para compilao.

Aproveite essa oportunidade para aprender a extrair informaes teis em grandes volumes de dados. Ser avaliado a capacidade do aluno de interagir com o ambiente Spark bem como a qualidade de suas anlises.

**Obs.:** Como o cluster  compartilhado para todos os alunos, cuidado para extrapolar o armazenamento com mltiplos resultados intermedirios. Para trabalhos como esse, as etapas de preprocessamento dos dados, ajudam na reduo do tamanho final.

**Boa sorte!**

*“The Answer to the Great Question... Of Life, the Universe and Everything... Is... Forty-two.”*  
(The Hitchhiker’s Guide to the Galaxy)

---

<sup>2</sup><https://bit.ly/2XWTUGz>