

Jefferson Carlos de Mendonça

# **Modelo de Mineração de Dados Análise das Perguntas do Site Stack Overflow**

São Paulo

Agosto 2016, v-1.0.1



Jefferson Carlos de Mendonça

# **Modelo de Mineração de Dados**

## **Análise das Perguntas do Site Stack Overflow**

Modelo canônico de Projeto de Pesquisa em conformidade com as normas ABNT apresentado ao Programa de Pós-Graduação em Engenharia da Computação da Escola Politécnica da Universidade de São Paulo, requisito para a seleção de ingresso ao Curso de Mestrado

---

**Jefferson Carlos de Mendonça**  
Aluno

---

**Edson Satoshi Gomi**  
Tutor/Orientador

São Paulo  
Agosto 2016, v-1.0.1



## Lista de tabelas



# Lista de abreviaturas e siglas

USP	Universidade de São Paulo
API	Application Programming Interface
REST	Representational State Transfer





# Sumário



# Introdução

Não raro a evolução na área computacional se desenvolve em ritmo acelerado, novos algoritmos, técnicas para processamento de imagem, computação distribuída, inteligência artificial etc. Com as *Linguagens de Programação* não é diferente, elas aprimoram suas API's constantemente, construindo e desconstruindo métodos e muitas vezes quebram paradigmas. O que se sabe hoje pode estar obsoleto amanhã, as exigências e tendências mudam com frequência, requerindo capacitação atualizada aos profissionais que atuam neste universo.

O grande desafio das universidades que trabalham na formação de profissionais é prover conhecimento não só para especialistas com perfil para atuar em linhas de pesquisa, onde o conhecimento é mais profundo, mas também para aqueles que, após a formação irão compor o mercado de trabalho, em que o conhecimento pode ser mais raso, porém mais dinâmico.

E para dar vazão a este dinamismo profissionais da área de computação diariamente recorrem a cursos na *internet*, livros, revistas e claro *websites*. Dentre estas mídias, merece destaque o fórum *Stack Overflow*<sup>1</sup>, maior comunidade *online* para programadores aprender, compartilhar conhecimento e progredir na carreira. Também é possível fazer um *tour*<sup>2</sup> pelo site, onde são apresentados as regras e a mecânica desta poderosa ferramenta para troca de conhecimento.

O objeto deste projeto de pesquisa é propor um algoritmo de linguagem de máquina que consiga catalogar as perguntas desta comunidade, nosso objetivo é entender os principais questionamentos dos usuários e detectar padrões nas dúvidas mais frequentes. Estes insumos serão disponibilizados para que as universidades e instituições de ensino possam aperfeiçoar seus cursos e treinamentos, estreitando a distância entre o que é lecionado e os requisitos impostos pelo dia-a-dia nas empresas.

---

<sup>1</sup> <<http://stackoverflow.com/company/about>>

<sup>2</sup> <<http://stackoverflow.com/tour>>



# 1 Objetivos

*Tópicos de interesse, assunto previsto, relevância para a área específica e aplicabilidade do estudo*

Os principais questionamentos dos usuários participantes da rede *Stack Overflow*, objeto de pesquisa deste trabalho, será embasado pelos assuntos das áreas: *Educational Data Mining*, *Learning Analytics*, *Big Data* e *Machine Learning* e versará sobre os tópicos de interesse: *Mineração de dados* para reconhecer padrões nas dúvidas dos usuários e *algoritmos baseados em aprendizagem de máquina* para arranjar as perguntas sintaticamente equivalentes em grupos.

O universo de dados coletados e categorizados, será de grande relevância para instituições de ensino sobre o tema Programação de Computadores - foco das análises, que poderão aprimorar sua grade curricular e com isso reduzir a distância entre as expectativas do mercado de trabalho com o que é lecionado na sala de aula.



## 2 Metodologia

O ponto de partida será a reunião de material literário para a composição da estrutura e argumentação sobre o objeto de estudo proposto.

Na prática será desenvolvido um programa de computador para consumir os dados do *Stack Overflow* utilizando sua API<sup>1</sup>, imputando toda a informação extraída em uma fila de processamento.

Na segunda etapa será proposto um algoritmo de linguagem de máquina (??) para extrair os assuntos principais (??) das dúvidas postadas no *site* e armazenar as informações obtidas. Em adição também serão armazenados a quantidade de visualizações, pontuação e o número de respostas da referida pergunta.

Com os dados processados e armazenados em um banco para grandes massa de dados<sup>2</sup>, faremos uso de técnicas de mineração de dados para detectar padrões nas dúvidas mais frequentes da comunidade e por fim será criado um *ranking* (??) das mais relevantes para o ensino da Linguagem de Programação *Java*.

---

<sup>1</sup> <<https://api.stackexchange.com/docs>>

<sup>2</sup> De acordo com a empresa *Stack Overflow* o *site* recebe cerca de 101 milhões de visitantes únicos mensalmente e conta com 3.7 milhões de perguntas respondidas





## 3 Recursos

Para que os algoritmos propostos, bem como os programas de computadores desenvolvidos para este projeto sejam executados em ambiente de alta performance, serão utilizados servidores na nuvem nas seguintes etapas: consumo dos *endpoints REST* providos pelo *Stack Overflow* e entrada de dados na fila de processamento, *queue*.



## 4 Cronograma

A ?? apresenta o cronograma estimado para a conclusão do projeto proposto.

Tabela 1 – Lista de atividades.

<b>Atividade</b>	<b>Tópico</b>	<b>Tempo</b>
Linguagem de Máquina	Algoritmos para extração de textos	32
Mineração de Dados	Tabulação de dados	32
<i>Big Data</i>	Técnicas de <i>Map Reduce</i>	24
Infra	Criar infraestrutura na nuvem	2
<i>Stack Overflow</i>	Estudo da API da comunidade	4
Programa	Extração dos dados via API	4
Monografia	Revisar e complementar a monografia desenvolvida	3
Apresentação	Material para apresentar a dissertação	3

Observação: Tempo Estimado em semanas