

# Um Sistema Distribuído para Análise de Recurso de Conteúdo para Prever Informações de Usuários em Mídias Sociais

Pedro Garcia

14 de setembro de 2013

# Objetivo (problema)

- ▶ Prever informações de usuários.
- ▶ Objeto de estudo: estimar idade do usuário com base em informações secundárias (formação, empregos, etc).

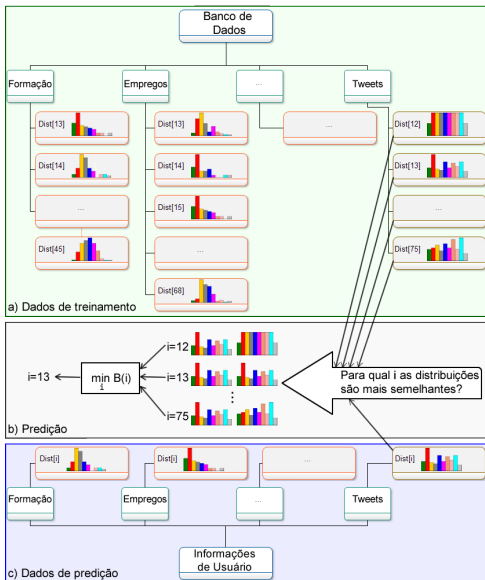
# Solução Proposta

- ▶ Utilizamos um estimador da informação desejada utilizando dois níveis.
- ▶ O primeiro utiliza uma estimação inicial, buscando em todo domínio de dados, minimizando-se a distância de Bhattacharyya.
- ▶ Após a distância inicial, busca-se refinar a estimativa, minimizando-se a divergência de Kullback–Leibler.

# Solução Proposta

- ▶ Inicialmente, computamos a frequência de cada palavra que o usuário de uma rede social produza.
- ▶ Em seguida, agrupamos cada frequência de acordo com o tipo da informação do usuário (emprego, educação, etc), gerando um conjunto de distribuições distintas para cada informação que desejamos estimar.
- ▶ Então, computamos a distância de Bhattacharyya para cada uma das distribuições computadas.
- ▶ Finalmente, teremos um conjunto de estimativas para cada um dos tipos de informações estimadas, buscando qual distribuição corresponda com a distribuição do usuário, cuja idade deseja-se estimar.

# Solução Proposta



Um Sistema Distribuído para Análise de Recurso de Conteúdo para Prever Informações de Usuários em Mídias Sociais

Pedro Garcia

Introduction  
Objetivo

Figura : Correspondência de Distribuições.

# Solução Proposta

- ▶ Dentro desse conjunto de Informações de Usuários estimadas, buscamos quais dessas informações são mais “confiáveis”.
- ▶ O critério que usamos para isso é a mínima Divergência de Kullback–Leibler.

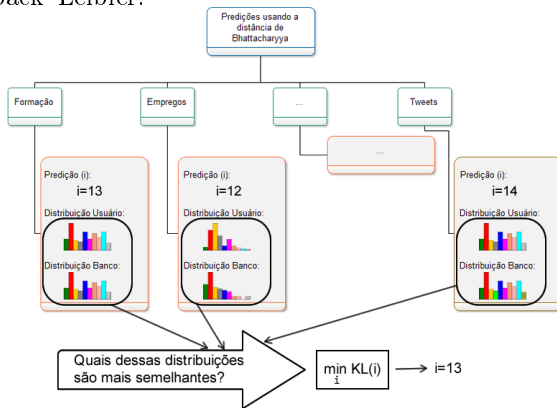


Figura : Seleção da Predição.

# Distribuição da Solução Proposta

- ▶ Redes sociais envolvem, geralmente, um grande volume de dados.
- ▶ Isso complica em uma grande quantidade de dados de análise.
- ▶ Logo, uma abordagem distribuída para esse sistema é desejável.

# Distribuição da Solução Proposta

Um Sistema Distribuído para Análise de Recurso de Conteúdo para Prever Informações de Usuários em Mídias Sociais

Pedro Garcia

Introduction  
Objetivo

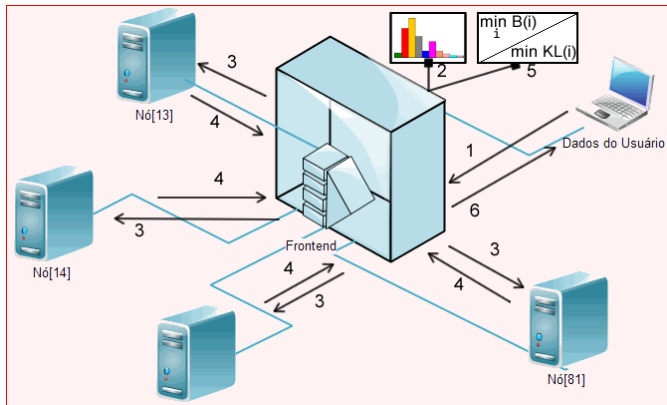


Figura : Distribuição das Tarefas ao Longo dos Nós.



# Resultados

- ▶ Erro quadrático médio: 1.4545
- ▶ Recall: 0.3644,
- ▶ Precisão: 0.5,
- ▶ F-Measure: 0.4215

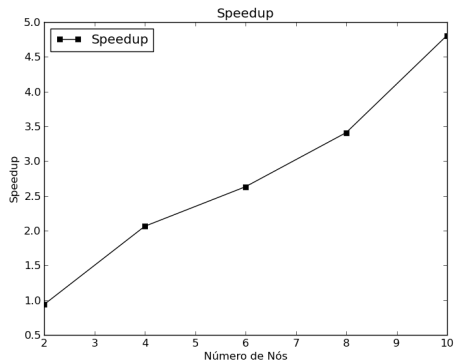


Figura : Speedup.

# Conclusão

- ▶ Os valores indicam que não houve uma relevante taxa de valores preditos que não corresponderam à idade exata do usuário.
- ▶ Pela Figura 4 é possível notar como o speedup cresce conforme o aumento do número dos nós, mostrando que o tempo total de solução do problema diminui, conforme o esperado.
- ▶ O artigo detalhado desse trabalho foi publicado no WSCAD 2013.