

THIAGO A. LIRA
ADRIANO . DENNANNI
RICARDO Y. NAGANO

NETMAP

Texto apresentado à Escola Politécnica da Universidade de São Paulo como requisito para a conclusão do curso de graduação em Engenharia de Computação, junto ao Departamento de Engenharia de Computação e Sistemas Digitais (PCS).

São Paulo
2016

THIAGO A. LIRA
ADRIANO . DENNANNI
RICARDO Y. NAGANO

NETMAP

Texto apresentado à Escola Politécnica da Universidade de São Paulo como requisito para a conclusão do curso de graduação em Engenharia de Computação, junto ao Departamento de Engenharia de Computação e Sistemas Digitais (PCS).

Área de Concentração:

Engenharia de Computação

Orientador:

Reginaldo Arakaki

FICHA CATALOGRÁFICA

Lira et al.

netmap/ T. A. Lira, A. . Dennanni, R. Y. Nagano. São Paulo, 2016.
18 p.

Monografia (Graduação em Engenharia de Computação) — Escola Politécnica da Universidade de São Paulo. Departamento de Engenharia de Computação e Sistemas Digitais (PCS).

1. Assunto #1. 2. Assunto #2. 3. Assunto #3. I. Universidade de São Paulo. Escola Politécnica. Departamento de Engenharia de Computação e Sistemas Digitais (PCS). II. t.

AGRADECIMENTOS

RESUMO

Mapas físicos tornam-se cada vez menos utilizados com o desenvolvimento progressivo de sistemas de posicionamento cada vez melhores. O sistema americano GPS é possivelmente o mais utilizado, sendo que ele possibilita qualquer um ter informações sobre sua localização, dando apoio, por exemplo, à praticantes de trilhas e acampamentos, principalmente em casos de emergência. Porém, em ambientes fechados, as ondas eletromagnéticas utilizadas pelos satélites sofrem atenuações e interferências devidos aos materiais de construção, e assim o sistema perde precisão e não funciona com toda a precisão esperada. Como uma alternativa para esta dificuldade, procurou-se desenvolver um sistema, que consegue obter a posição do usuário em um ambiente fechado com precisão, sendo usado para isso técnicas de machine learning, aliadas com dados obtidos de redes em fio já instaladas no local. O sistema consistirá de um servidor central, onde serão enviados os dados e os mesmos serão processados. Os dados serão coletados por meio de um aplicativo de Android, este possuirá duas versões. A versão usuário usará os dados do servidor para localizar o usuário, a versão administrador irá coletar dados novos para serem usados em futuras medições. PALAVRAS-CHAVE: Indoor, Localização, Wi-fi, Machine Learning

ABSTRACT

SUMÁRIO

Lista de Ilustrações

Lista de Tabelas

1	Introdução	8
1.1	Apresentação	8
2	Machine Learning	10
2.1	Tratamento dos dados	10
2.1.1	Dados obtidos da UCI	10
2.1.2	Dados do Servidor	10
2.2	Formato dos dados tratados	10
2.3	<i>Cross-Validation</i> para os Modelos Usados	11
2.3.1	KNN : K-Nearest-Neighbors	11
2.3.2	Rede Neural	12
2.4	Métodos de Votação	13
2.4.1	Votação Simples	13
	Referências	17
	Apêndice A - Demonstração do Lema da Bifurcação	18

LISTA DE ILUSTRAÇÕES

1	Erros para uso da distância Euclidiana	11
2	Erros para uso da distância Manhattan	12
3	Erro para diversas quantidades de neurônios na rede neural.	13
4	Erro de votação simples para uma quantidade crescente de pontos de treino.	14
5	Erro do algoritmo KNN para uma quantidade crescente de pontos de treino.	15
6	Erro da Rede Neural para uma quantidade crescente de pontos de treino.	15
7	Erro do algoritmo SVM para uma quantidade crescente de pontos de treino.	16

LISTA DE TABELAS

1 INTRODUÇÃO

1.1 Apresentação

Com a modernização das tecnologias de telefonia móvel torna-se cada vez maior o número de pessoas com celulares associados a tecnologias de rede, como WiFi, 3G e 4G, que tem o potencial de fornecer informações sobre seu usuário a todo momento, tais como o conteúdo acessado por seus navegadores ou aplicativos, e também informações sobre posição e deslocamento. Dados de localização por si possuem pouco valor, mas quando aliados a outros conteúdos, é possível fornecer conteúdo personalizado em tempo real, reativo ao ambiente, passando a oferecer um grande retorno por um pouco mais de ocupação na banda[1]. Enquanto sistemas de posicionamento por satélite como GPS (EUA) ou GALILEO (Europa), por um lado, conseguem precisar a posição do usuário em até centímetros em um ambiente outdoor, por outro há uma dificuldade deste sistema em calcular o posicionamento em lugares fechados, devido basicamente à atenuação dos sinais causada pelas paredes e seus materiais. Tendo em vista o crescimento das cidades e consequente aumento no número de construções as pessoas cada vez passam mais tempo em ambientes fechados. A necessidade de serviços de localização indoor tem se tornado cada vez mais crítica[2]. Respondendo a essa necessidade surgiram alternativas para o posicionamento em ambientes fechado, entre eles há o uso da Identificação por Rádio Frequência (RFID ? Radio Frequency Identification), do bluetooth, do Zigbee ou do Wi-fi. A determinação de posicionamento via Wi-fi é uma tecnologia que usa os sinais modulados do Wi-fi para; detectar a presença

de um aparelho, em seguida o sistema é capaz de triangular a posição deste aparelho a partir dos sinais recebidos pelo ponto de acesso [1]. Um dos primeiros exemplos de um sistema de posicionamento utilizando Wi-fi foi RADAR, desenvolvido pela Microsoft, que também criou o RightSPOT que utilizava um ranking das frequências moduladas pelos pontos de acesso, no lugar de usar a potência dos sinais para determinar a posição do aparelho. Tendo em vista este cenário e as condições tecnológicas atuais, nosso projeto procura apresentar uma solução para localização de pessoas em ambientes fechados, como shoppings e eventos em galpões. Para tal, será combinado o uso de Big Data à tecnologia de Wi-fi já citada e ainda implementando um algoritmo de machine learning. O Big Data, que possui como uma de suas principais características a capacidade de armazenar uma grande quantidade de dados, será usado para armazenar todas as leituras de intensidade dos sinais feitas durante uma fase inicial, chamada de treinamento. À partir dessa larga quantidade de dados recebidos serão empregados algoritmos de machine learning processando-os para permitir que o sistema consiga determinar a posição do usuário dentro do ambiente mapeado anteriormente. Esta abordagem se mostra interessante ao ponto de que sua implementação não necessita configuração particular na rede que será usada, uma vez que se baseia em leituras feitas pelo aparelho móvel e no processamento dos dados feitos em um servidor em nuvem. Também se destaca o fato de que em decorrência do machine learning à medida de que o sistema é usado em uma determinada localidade a precisão tende a aumentar, uma vez que cada vez há mais dados para serem consultados para "Aprendizado". Esse approach difere de outro também amplamente empregado em esquemas de localização indoor, que é o de modelar o próprio sinal de Wi-Fi, deduzindo a sua distância a partir da intensidade medida, como em [4] e [5]. Esse método se torna ineficiente porque se torna necessário que sejam conhecidas as coordenadas de todas as APs locais, e também o próprio modelo de atenuação dos sinais RSSI de WiFi sofre todo o tipo de interferência tornando difícil a obtenção de uma boa precisão.

2 MACHINE LEARNING

2.1 Tratamento dos dados

2.1.1 Dados obtidos da UCI

Para grande parte dos testes de *Cross-Validation*, foi usado um dataset obtido por meio do repositório online da UCI, desenvolvido em (TORRES-SOSPEDRA et al., 2014).

2.1.2 Dados do Servidor

2.2 Formato dos dados tratados

A matriz de dados tratados usada diretamente pelos algoritmos de ML tem o seguinte formato:

$$\begin{array}{ccccc}
 ZoneID & BSSID_1 & BSSID_2 & \dots & BSSID_n \\
 \left(\begin{array}{ccccc}
 1 & -70 & -92 & \dots & -87 \\
 2 & -89 & -80 & \dots & -63 \\
 3 & -28 & -120 & \dots & -35 \\
 \vdots & \vdots & \vdots & \ddots & \vdots \\
 1 & -48 & -36 & \dots & -29
 \end{array} \right) & \begin{array}{l} Measure_1 \\ Measure_2 \\ Measure_3 \\ \vdots \\ Measure_n \end{array}
 \end{array}$$

2.3 *Cross-Validation* para os Modelos Usados

A chamada *Cross-Validation* dos modelos de ML é o teste para encontrar os parâmetros ótimos para o treinamento. O método escolhido de *Cross-Validation* foi o *K-Fold*. O método *K-Fold* divide o dataset em K subconjuntos de igual tamanho e então um dos conjuntos é usado como validação do treinamento feito pelos $K - 1$ subconjuntos restantes. O processo é repetido K vezes e em cada subdivisão possível podem ser usados valores de parâmetros de treino distintos.

2.3.1 KNN : K-Nearest-Neighbors

Para o algoritmo de KNN, foi usado o método *K-Fold* com $K = 10$. O método foi aplicado duas vezes. Na primeira o modelo de KNN foi treinado com os cálculos de distância usando-se a distância Euclidiana, e variando-se o valor de K do número de vizinhos em cada *fold*. Então, foi feito o mesmo processo, mas com a distância de Manhattan. Os resultados são mostrados a seguir:

Figura 1: Erros para uso da distância Euclidiana

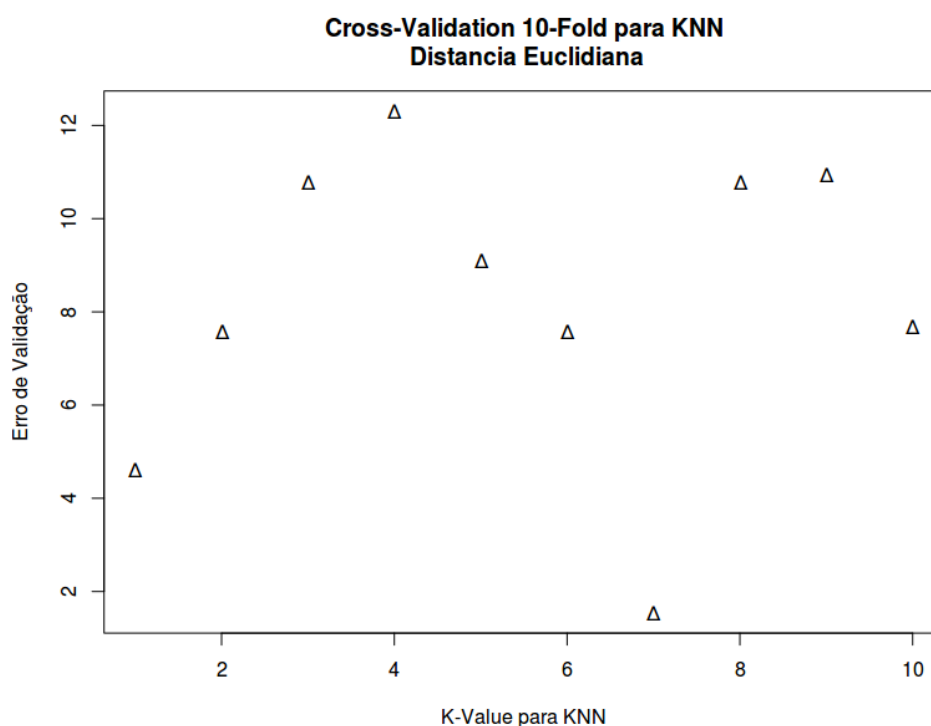
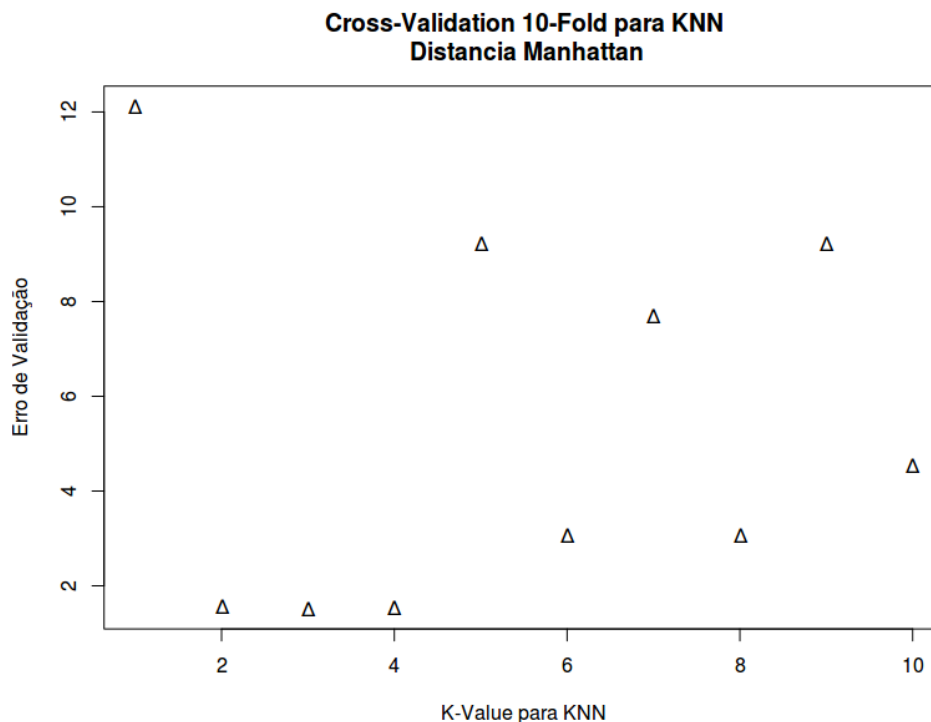


Figura 2: Erros para uso da distância Manhattan



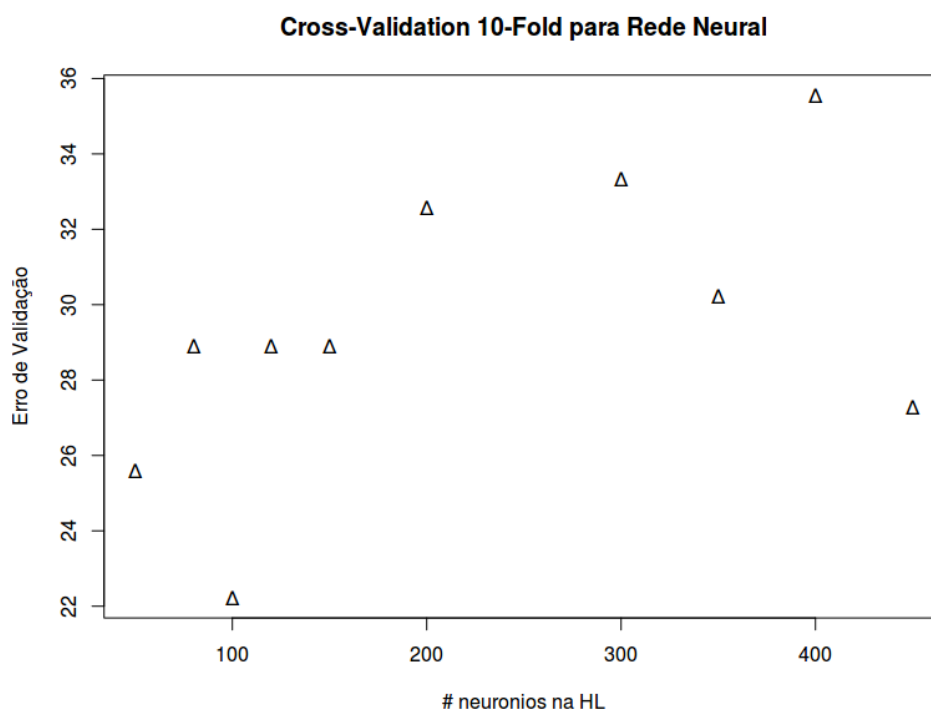
Podemos concluir então que (1) o erro é menor no geral com a distância de Manhattan e (2) o valor de K que minimiza o erro fica entre 2 e 4 vizinhos. No restante do trabalho foi escolhido o valor de $K = 3$ vizinhos.

2.3.2 Rede Neural

Para a *Cross-Validation* das Redes Neurais, também foi usado o método de *K-Fold* com $K = 10$. Para cada uns dos *folds* foi testado um número de neurônios em uma *hidden-layer* única. (Outros testes mostraram não valer a pena para o nosso problema usar mais de uma camada de *hidden layer* ou *deep learning*). Os resultados são mostrados a seguir:

Concluimos que o número ideal de neurônios na *hidden-layer* está próximo de 200. E é esse valor que usaremos daqui para frente.

Figura 3: Erro para diversas quantidades de neurônios na rede neural.



2.4 Métodos de Votação

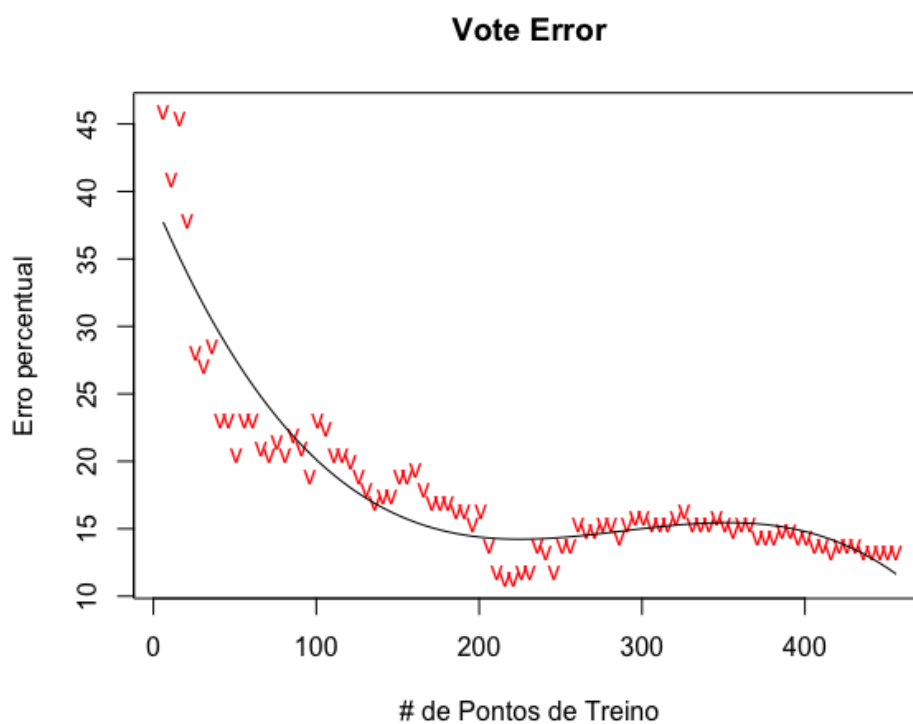
Os chamados sistemas de *Ensemble Learning* são aqueles em que diversos modelos de ML são treinados e então é feito um sistema de votação para que a classificação final seja feita. No nosso sistema experimentamos com 2 modelos de votação diferentes, documentados a seguir.

2.4.1 Votação Simples

Na votação simples os resultados da classificação de todos os modelos são contabilizados e é escolhido aquele com o maior número de ocorrências entre os classificadores. Os testes foram feitos calculando-se o erro de validação no MESMO dataset de treino com uma quantidade cada vez maior de pontos de treino, estes retirados de um dataset de treino também fixado.

A votação simples foi testada com os algoritmos de Rede Neural, KNN e SVM

Figura 4: Erro de votação simples para uma quantidade crescente de pontos de treino.



Para comparação, foram calculados os erros de validação independentes de cada um dos algoritmos usados.

Figura 5: Erro do algoritmo KNN para uma quantidade crescente de pontos de treino.

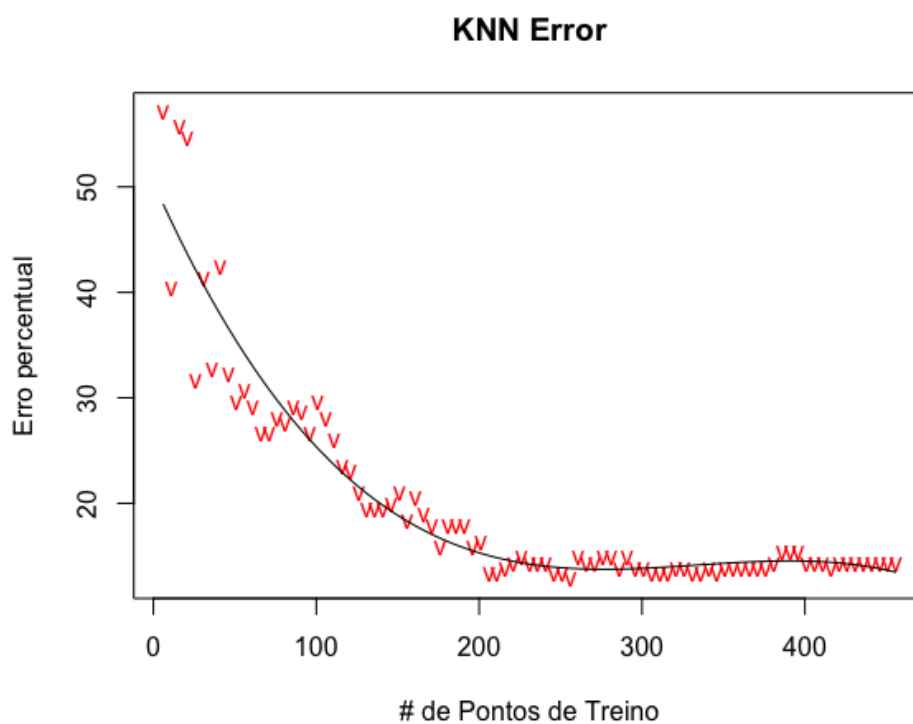


Figura 6: Erro da Rede Neural para uma quantidade crescente de pontos de treino.

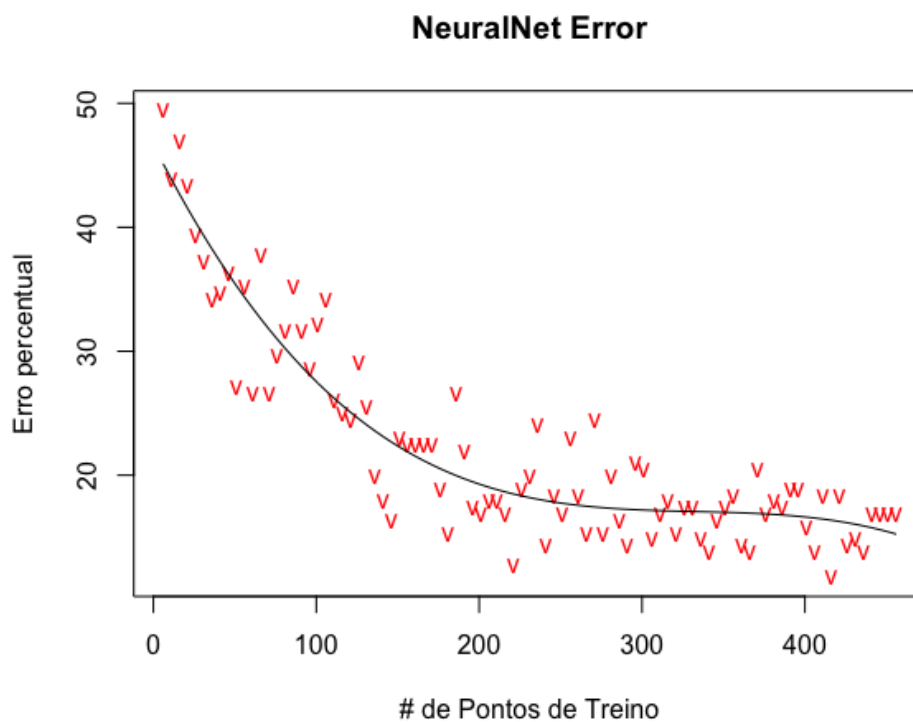
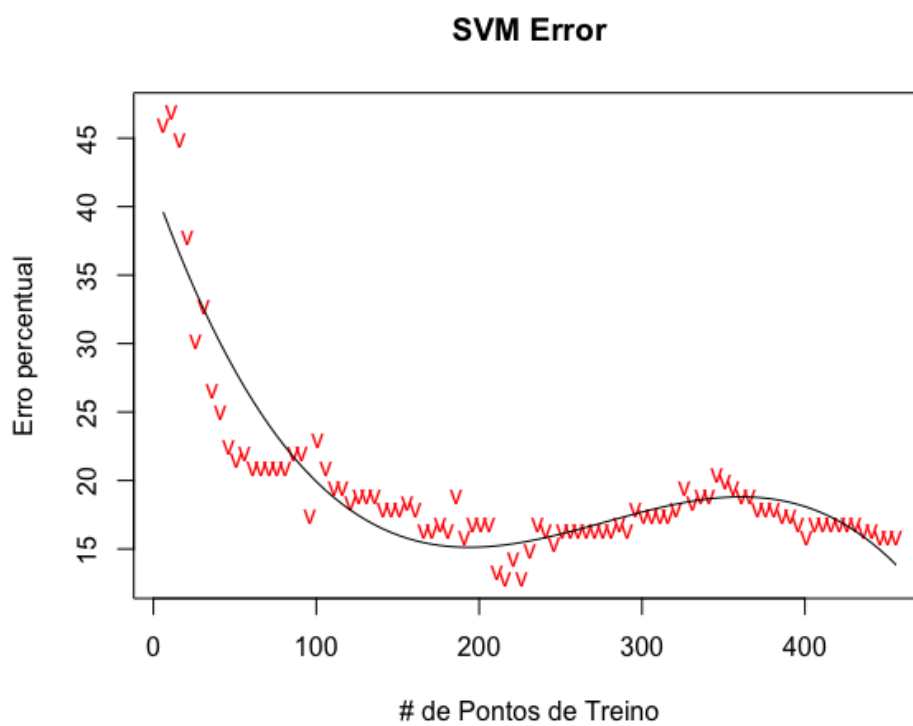


Figura 7: Erro do algoritmo SVM para uma quantidade crescente de pontos de treino.



REFERÊNCIAS

TORRES-SOSPEDRA, J.; MONTOLIU, R.; MARTÍNEZ-USÓ, A.; AVARIENTO, J. P.; ARNAU, T. J.; BENEDITO-BORDONAU, M.; HUERTA, J. Ujiindoorloc: A new multi-building and multi-floor database for wlan fingerprint-based indoor localization problems. In: *Indoor Positioning and Indoor Navigation (IPIN), 2014 International Conference on*. [S.l.: s.n.], 2014. p. 261–270.

APÊNDICE A - DEMONSTRAÇÃO DO LEMA DA BIFURCAÇÃO