



Universidade Federal do Piauí
Centro de Ciências da Natureza
Programa de Pós-Graduação em Ciência da Computação

Rotulação com Algoritmos Supervisionados

Tarcísio Franco Jaime

Número de Ordem PPGCC: M001

Teresina-PI, Janeiro de 2017

Tarcísio Franco Jaime

Rotulação com Algoritmos Supervisionados

Qualificação de Mestrado apresentada ao Programa de Pós-Graduação em Ciência da Computação da UFPI (área de concentração: Sistemas de Computação), como parte dos requisitos necessários para a obtenção do Título de Mestre em Ciência da Computação.

Universidade Federal do Piauí – UFPI

Centro de Ciências da Natureza

Programa de Pós-Graduação em Ciência da Computação

Orientador: Vinicius Ponte Machado

Teresina-PI

Janeiro de 2017

Tarcísio Franco Jaime

Rotulação com Algoritmos Supervisionados/ Tarcísio Franco Jaime. – Teresina-PI, Janeiro de 2017-

31 p. : il. (algumas color.) ; 30 cm.

Orientador: Vinicius Ponte Machado

Qualificação (Mestrado) – Universidade Federal do Piauí – UFPI

Centro de Ciências da Natureza

Programa de Pós-Graduação em Ciência da Computação, Janeiro de 2017.

1. Rotulação. 2. Algoritmos Supervisionados. 3. CART. 4. Naive Bayes. I. Vinicius Ponte Machado. II. Universidade Federal do Piauí. III. Rotulação com Algoritmos Supervisionados.

CDU 02:141:005.7

Tarcísio Franco Jaime

Rotulação com Algoritmos Supervisionados

Qualificação de Mestrado apresentada ao Programa de Pós-Graduação em Ciência da Computação da UFPI (área de concentração: Sistemas de Computação), como parte dos requisitos necessários para a obtenção do Título de Mestre em Ciência da Computação.

Trabalho aprovado. Teresina-PI, 01 de janeiro de 2018:

Vinicius Ponte Machado
Orientador

Co-Orientador

Professor
Convidado 1

Professor
Convidado 2

Professor
Convidado 3

Teresina-PI
Janeiro de 2017

*Aos meus pais XXXXXXXX e YYYYYYY,
por sempre estarem comigo em todos os momentos.*

Agradecimentos

Agradeço a Deus.

Agradeço aos meus pais, XXXXX e YYYYY, por ...

Aos meus irmãos, por.....

Agradeço ao meu orientador, XXXXXXXXX, por todos os conselhos, pela paciência e ajuda nesse período.

Aos meus amigos ...

Aos professores ...

À XXXXXX pelo apoio financeiro para realização deste trabalho de pesquisa.

*“Não sei o que,
não sei o que,
não sei o que lá.”
(Autor Desconhecido)*

Resumo

Segundo a ABNT, o resumo deve ressaltar o objetivo, o método, os resultados e as conclusões do documento. A ordem e a extensão destes itens dependem do tipo de resumo (informativo ou indicativo) e do tratamento que cada item recebe no documento original. O resumo deve ser precedido da referência do documento, com exceção do resumo inserido no próprio documento. (...) As palavras-chave devem figurar logo abaixo do resumo, antecidas da expressão Palavras-chave:, separadas entre si por ponto e finalizadas também por ponto.

Palavras-chaves: latex. abntex. editoração de texto.

Abstract

This is the english abstract.

Keywords: latex. abntex. text editoration.

Lista de ilustrações

Figura 1 – Breve explicação sobre a figura. Deve vir abaixo da mesma.	1
Figura 2 – Hipóteses ajustadas	5
Figura 3 – Ponto de Corte (R-1)	9
Figura 4 – Discretização EWD	10
Figura 5 – Discretização EFD	11
Figura 6 – Modelo (LOPES; MACHADO; RABELO,)	11
Figura 7 – Modelo (FILHO, 2015)	12

Lista de tabelas

Tabela 1 – Breve explicação sobre a tabela. Deve vir acima da mesma.	2
--	---

Lista de abreviaturas e siglas

ABNT	Associação Brasileira de Normas Técnicas
abnTeX	ABsurdas Normas para TeX

Lista de símbolos

Γ	Letra grega Gama
Λ	Lambda
ζ	Letra grega minúscula zeta
\in	Pertence

Sumário

Introdução	1
figuras	1
tabelas	1
Motivação	2
Objetivos	2
1 REFERENCIAL TEÓRICO	3
1.1 Aprendizado de Máquina	3
1.1.1 Aprendizado Supervisionado	4
1.1.1.1 Algoritmo Classification and Regression Trees - CART	5
1.1.1.2 Algoritmo Naive Bayes	6
1.1.2 Aprendizado Não Supervisionado	7
1.2 Discretização	8
1.2.1 Discretização por Larguras Iguais - EWD	8
1.2.2 Discretização por Frequência Iguais - EFD	9
1.3 Trabalhos Correlatos	10
2 METODOLOGIA / MATERIAIS E MÉTODOS	13
2.1 Considerações do Problema	13
3 RESULTADOS E DISCUSSÃO	15
3.1 Base de Dados	15
3.2 Considerações Finais	15
Conclusão e Trabalhos Futuros	17
REFERÊNCIAS	19
APÊNDICES	21
APÊNDICE A – PRIMEIRO APÊNDICE	23
APÊNDICE B – PERCEBA QUE O TEXTO DO TÍTULO DESSE SEGUNDO APÊNDICE É BEM GRANDE	25

ANEXOS	27
ANEXO A – NOME DO PRIMEIRO ANEXO	29
ANEXO B – NOME DE OUTRO ANEXO	31

Introdução

Este documento segue as normas estabelecidas pela ??, 3.1-3.2).

A proposta deste mestrado bem como outros trabalhos relacionados, onde áreas envolvidas tem como tema principal, Rotulação de Dados, estão alterando a maneira de como Aprendizagem de Máquina define este termo. Em pesquisas realizadas neste área sob supervisão do orientador desta proposta, vários trabalhos estão definindo Rotulação sendo algo diferente da Classificação dos dados.

Apesar de várias literaturas (BARBER, 2011; MITCHELL, 1997) entre outras citarem o termo rotulação como um sinônimo de classificação, neste departamento, esse termo esta ficando obsoleto. Muitos trabalhos feitos aqui neste laboratório estão redefinindo o termo rotulação como algo mais completo e que possui propriedade diferente a apresentada na classificação.

A classificação é dada com um identificador do registro informante qual classe ele pertence....

Figuras

As normas da ??, 3.1-3.2) especificam que o caption da figura deve vir abaixo da mesma.

A Figura 1 ilustra...



Figura 1 – Breve explicação sobre a figura. Deve vir abaixo da mesma.

Tabelas

A Tabela 1 apresenta os resultados...

Tabela 1 – Breve explicação sobre a tabela. Deve vir acima da mesma.

XX	FF	PP	YY	Yr	xY	Yx	ZZ
615	18	2558	0,9930	0,9930	0,9930	0,9930	0,9930
615	18	2558	0,9930	0,9930	0,9930	0,9930	0,9930
615	18	2558	0,9930	0,9930	0,9930	0,9930	0,9930
615	18	2558	0,9930	0,9930	0,9930	0,9930	0,9930
615	18	2558	0,9930	0,9930	0,9930	0,9930	0,9930

Motivação

Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Aenean nonummy turpis id odio. Integer euismod imperdiet turpis. Ut nec leo nec diam imperdiet lacinia. Etiam eget lacus eget mi ultricies posuere. In placerat tristique tortor. Sed porta vestibulum metus. Nulla iaculis sollicitudin pede. Fusce luctus tellus in dolor. Curabitur auctor velit a sem. Morbi sapien. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Donec adipiscing urna vehicula nunc. Sed ornare leo in leo. In rhoncus leo ut dui. Aenean dolor quam, volutpat nec, fringilla id, consectetur vel, pede.

Objetivos

Nulla malesuada risus ut urna. Aenean pretium velit sit amet metus. Duis iaculis. In hac habitasse platea dictumst. Nullam molestie turpis eget nisl. Duis a massa id pede dapibus ultricies. Sed eu leo. In at mauris sit amet tortor bibendum varius. Phasellus justo risus, posuere in, sagittis ac, varius vel, tortor. Quisque id enim. Phasellus consequat, libero pretium nonummy fringilla, tortor lacus vestibulum nunc, ut rhoncus ligula neque id justo. Nullam accumsan euismod nunc. Proin vitae ipsum ac metus dictum tempus. Nam ut wisi. Quisque tortor felis, interdum ac, sodales a, semper a, sem. Curabitur in velit sit amet dui tristique sodales. Vivamus mauris pede, lacinia eget, pellentesque quis, scelerisque eu, est. Aliquam risus. Quisque bibendum pede eu dolor.

1 Referencial Teórico

Será abordado neste capítulo o conteúdo base na compreensão deste trabalho dividido em 3 sessões: Aprendizado de Máquina, Discretização e Trabalhos Correlatos.

A primeira sessão contempla os principais tipos de aprendizados indutivos, não incluindo aqui o aprendizado semi-supervisionado e sim dando ênfase a aprendizagem supervisionada, foco da proposta deste mestrado. O aprendizado indutivo utiliza uma amostra do todo para tirar uma conclusão. Caso os exemplos retirados de uma base de dados não forem suficientes, talvez o conhecimento derivado destes exemplos não mostrem a verdade.

O segundo item dissertará sobre a técnica de discretização adotada nesta pesquisa. Possuindo grande contribuição para os resultados gerados, e ganhando assim uma sessão própria para explanação de como funciona essa técnica. E na terceira sessão serão abordados trabalhos com mesmas características particulares para melhor elucidar o motivo da elaboração dessa proposta de mestrado.

1.1 Aprendizado de Máquina

Aprendizagem de máquina é a capacidade do aprendizado automático com utilização de algoritmos atuando em cima de uma base de dados. Diz-se que o computador está aprendendo quando existe uma melhora de desempenho de tarefas que ele utilizou como exemplo (MITCHELL, 1997). Um exemplo seria a realização do reconhecimento facial de uma pessoa utilizando aprendizado de máquina. Não seria necessário a implementação de várias linhas de código informando que a cor dos olhos são azuis com orelhas e cabelos grandes, seriam de uma certa pessoa. Ao invés disso é observada várias fotos tituladas de uma certa pessoa, e após vários exemplos o computador seria capaz de prever uma foto nova, se é, ou não, da determinada pessoa através de aprendizado anterior.

Existem alguns motivos, onde justificam, que não é possível simplesmente exigir que o projetista implemente melhorias no sistema de forma que ele esteja robusto bastante para lidar com todas as situações (RUSSEL; NORVIG, 2013). Um desses motivos seria a incapacidade da antecipação de todas as situações possíveis de implementação por parte do programador. Fazendo um resumo, aprendizado de máquina seriam algoritmos capazes de aprender automaticamente através de determinados exemplos, ou comportamentos.

A partir desta síntese, tem-se uma observação. A classificação de dados no contexto de aprendizado de máquina, são compostos por dois pilares. Um, seriam os **dados** a serem classificados, e outro, o **algoritmo** que irá atuar nessa base de dados. Existem vários

algoritmos como exemplo: redes neurais, árvores de decisão, Suport Vector Machine – SVM, etc. Qualquer um destes algoritmos são utilizados para solucionar essa classificação. E a escolha apropriada, desse algoritmo, se dará através de métricas que avaliarão o desempenho de cada um, e a melhor métrica, será o algoritmo apropriado para aquele problema de classificação de dados.

Uma analogia referente do que foi dito acima seria um “problema”, comparado a um “motor”, e os algoritmos disponíveis seriam as "ferramentas" para concertar esse motor. A partir daí a ferramenta que fosse mais eficaz, considerando métricas de desempenho, para fazer o motor funcionar, seria a ferramenta(algoritmo) escolhida. Tendo assim a escolha certa para um determinado problema.

1.1.1 Aprendizado Supervisionado

Nesta sessão será abordado um método que através de uma banco de dados já classificado por especialistas, será feita uma predição de novos registros com base em vários desses exemplos já classificados. Os responsáveis por essas predições de novos registros são algoritmos de aprendizado supervisionados projetados para determinados fins.

O termo "Supervisionado" indica que existe um supervisor para cada registro de entrada especificando uma saída para esse registro. Considerando uma base de dados de imagens de rostos, onde cada imagen possui uma saída representado por uma classe: masculino ou feminino. A tarefa seria criar um preditor capaz de acertar a cada novo registro se a imagem é masculina ou feminina. Seria difícil implementar de maneira tradicional, uma vez que são inúmeras as diferenças que difere as faces masculinas e femininas. Mas uma alternativa seria dar exemplos de rostos com suas classificações de fazer que automaticamente a máquina "aprenda" uma regra para predizer se é masculino ou feminino (BARBER, 2011).

Em (RUSSEL; NORVIG, 2013) os autores fazem uma apresentação formal do funcionamento da aprendizagem supervisionada. Dado um conjunto de treinamento

$$(x_1, y_1), (x_2, y_2), \dots (x_n, y_n), \quad (1.1)$$

onde cada y_j foi gerado por $y = f(x)$ desconhecida. Encontrar uma função h que se aproxime da função f real.

A função h é uma hipótese onde prevê um melhor desempenho entre as hipóteses possíveis através dos conjuntos de exemplos, que são diferentes do conjunto de treinamento 1.1.

Na figura 2a existe um sobre ajuste da função com o conjunto de dados de treinamento. Esse exemplo acabou exibindo uma função mais complexa para se molda de acordo com os sete pontos do gráfico, especificando para esse conjunto de dados.

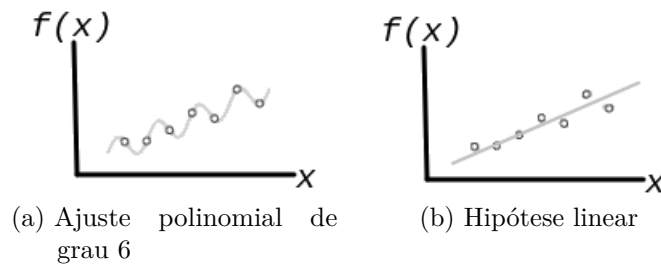


Figura 2 – Hipóteses ajustadas

Ja na figura 2b o ajuste da função se torna mais simples e mesmo não passando por todos os pontos, acabou generalizando melhor o conjunto de treinamento, tornando talvez, um melhor resultado da predição de novos valores.

A figura 2 mostra duas hipóteses que tentam se aproximar ao máximo da função verdadeira, que é desconhecida. Mesmo parecendo que na figura 2a obteve-se melhor resultado, pois todos os pontos são contemplados pela função, mas esta função h acabou ficando muito específica e isso não retrata os dados em um mundo real. Então quanto mais generalizado for h , melhor será para prevê os valores de y para novos conjuntos de dados.

Antes de falar dos algoritmos utilizados nesse texto a aprendizagem supervisionada detem dois tipos de caso: regressão e classificação. A classificação, contém variáveis com valores discretos, onde as amostras destas variáveis de saída estão na forma de categorias. Como exemplo poderia ser masculino e feminino. Já no tipo regressão, possuem valores contínuos: quantidade de água em ml, velocidade de um carro, altura de uma pessoa.

1.1.1.1 Algoritmo Classification and Regression Trees - CART

Esse algoritmo constroi modelos de previsão a partir de dados de treinamento onde seus resultados podem ser representados em uma árvore de decisão. No caso de não ser probabilístico o grau de confiança em seu modelo de predição será embasada em respostas semelhantes em outras circunstâncias antes analisadas.

Inicialmente todas as amostras se concentram no nó raiz, e a partir daí é apresentado uma questão, onde a intenção é separar o nó raiz em dois grupos mais homeogêneos. Dependendo da questão as amostras iram para a folha esquerda ou direita do nó raiz.

O CART faz essa divisão em função da regra Gini¹??, parecida com a regra da entropia usada no algoritmo ID3². O índice Gini varia de 0 a 1, definindo o grau de pureza do nó.

$$Gini(S) = 1 - \sum p^2(j/t) \quad (1.2)$$

¹ O CART pode utilizar outros critérios de divisão de dados como: entropia e critério de Twoing

² Algoritmo abordado por (??)

Onde: $p(j/t)$ é probabilidade a priori da classe j se formar no nó t . E S é um conjunto de dados que contém exemplos de n classes

Para construção de uma árvore existem três componente importantes (YOHANNES; WEBB, 1999):

- Um conjunto de perguntas que servirá de base para fazer uma divisão;
- Regras de divisão para julgar o quanto é boa esta divisão;
- Regras para atribuir uma classe a cada nó;

Abaixo segue um algoritmo de como o critério Gini é aplicado nas variáveis (RAIMUNDO; MATTOS; WALESKA, 2008):

Algorithm 1: Rotina de funcionamento do CART

```

1 melhorGini; /* cria a variável */
2 divisaoCorrente  $\leftarrow$  4.9; /* Ex. recebe o 1º valor do atributo */
3 direita  $\leftarrow$  0;
4 esquerda  $\leftarrow$  6; /* Ex. recebe o total de dados existentes para o
   atributo */
5 while existirem dados do
6   if 1ª Dado Lista do Atributo MAIOR divisaoCorrente then
7      $\lfloor$  valorGini  $\leftarrow$  calculaGini(divisaoCorrente);
8   else
9      $\lfloor$  valorGini  $\leftarrow$  calculaGini(1ª DadoLista);
10  if Primeiro Gini encontrado then
11     $\lfloor$  melhorGini  $\leftarrow$  valorGini;
12  else
13    if valorGini > melhorGini then
14       $\lfloor$  melhorGini  $\leftarrow$  valorGini
15  divisaoCorrente  $\leftarrow$  5.4; /* recebe o próximo dado do atributo */
16  direita recebe o que possui +1 e esquerda o -1;
17  (valorGini + divisaoCorrente)/2; /* encontrar ponto de divisão */

```

1.1.1.2 Algoritmo Naive Bayes

É um algoritmo considerado rápido, em relação a outros algoritmos de classificação, mesmo com grandes volumes de dados em seu conjunto de treinamentos. Utiliza modelo probabilístico, Teorema de Bayes e possui a característica de independência dos atributos, onde as classes não dependem de recursos de outras. Essa independência condicionada entre os atributos, os quais nem sempre ocorrem nos problemas reais, acabou sendo conhecida por Bayes ingênuo, ou Naive Bayes.

Naive Bayes como classificador estatístico possui um modelo de simples construção, e ficou conhecido por ter bons resultados em relação a algoritmos mais sofisticados, mesmo trabalhando com grandes quantidades de dados. Ele agrupa objetos de uma certa classe em razão da probabilidade do objeto pertencer a esta classe.

$$P(c/x) = \frac{P(x/c)P(c)}{P(x)} \quad (1.3)$$

$$P(c/x) = P(x_1|c) * P(x_2|c) * \dots * P(x_n|c) * P(c) \quad (1.4)$$

- $P(c/x)$ probabilidade posterior da classe c , alvo dada preditor x , atributos.
- $P(c)$ é a probabilidade original da classe.
- $P(x|c)$ é a probabilidade que representa a probabilidade de preditor dada a classe.
- $P(x)$ é a probabilidade original do preditor.

A utilização do algoritmo Naive Bayes já é bem difundida, e está presente em vários trabalhos, como classificação de textos, filtro de SPAM, analisador de sentimentos, entre outros ([MADUREIRA, 2017](#); [LUCCA et al., 2013](#); [WU et al., 2008](#); [MCCALLUM; NIGAM, 1997](#)). Mas mesmo atingido popularidade existem pontos negativos. A suposição de ter preditores independentes não acontece muito na vida real, pois acaba sendo difícil ter uma amostra de dados que sejam inteiramente independentes.

Outra situação é caso de existir uma variável categórica que não foi observada na amostra tirada para o conjunto de treinamento, então poderá o modelo atribuir probabilidade 0(zero), não sendo capaz de fazer uma previsão. Quando isso acontecer uma técnica de alisamento é aplicada, chamada estimativa de Laplace, utilizadas em probabilidades condicionadas.

1.1.2 Aprendizado Não Supervisionado

No Aprendizado Não Supervisionado, não existe uma tentativa de se encontrar uma função que se aproxime da real. Logo porque os registros não são classificados, então o conjunto de treinamento não possui informação da saída sobre determinada entrada. Desta forma os algoritmos procuram algum grau de similaridade entre os registros e tenta agrupá-los de forma a ter algum sentido deles estarem juntos.

Quando o algoritmo encontram dados com mesma similaridade ele os agrupa formando clusters. Os números de clusters encontrados irão depender de como os algoritmos funcionam, junto com o grau de dissimilaridade entre elementos de grupos diferentes. Como não existe uma variável classe no Aprendizado Não Supervisionado, então ([BARBER,](#)

2011) diz que o maior interesse seria em uma perspectiva probabilística de distribuição $p(x)$ de um determinado conjunto de dados.

$$D = \{x_n, n = 1, \dots, N\} \quad (1.5)$$

Uma vez que no conjunto 1.5 não existe classe y , encontrado em um conjunto de treinamento 1.1 o algoritmo precisa encontrar padrões nos atributos para fazer os agrupamentos.

1.2 Discretização

A discretização faz parte em duas etapas no modelo defendido nesse trabalho, por isso a preocupação na explanação de seu funcionamento aqui nesta sessão. O método de discretização faz a conversão de valores contínuos em valores discretos. A partir de um atributo com valores contínuos, a discretização irá forçar um ponto inicial e final definindo um intervalo e designando uma faixa para cada intervalo. Assim, ao invés de valores contínuos em cada atributo, será relacionado a faixa que aquele atributo pertence, definindo assim seu novo valor. O melhor método de discretização seria encontrar o conjunto de valores contínuos por faixa de intervalos pequenos (KOTSIANTIS; KANELLOPOULOS, 2006)

A partir de alguns autores (CATLETT, 2006; HWANG; LI, 2002) a discretização melhora a precisão e deixa um modelo classificador mais rápido em seu conjunto de treinamento. Aqui nesse trabalho é utilizado a técnica de discretização antes da execução dos algoritmos e as faixas selecionadas são usadas para identificar o rótulo. Após o conhecimento do rótulo o valor da faixa é trocado pelo início e fim do intervalo.

Os métodos de discretização mais comumente utilizados no âmbito dos métodos não-supervisionados de acordo com (KOTSIANTIS; KANELLOPOULOS, 2006; DOUGHERTY; KOHAVI; SAHAMI, 1995) são os métodos de Discretização por Larguras Iguais(EWD) e Discretização por Frequências Iguais (EFD).

1.2.1 Discretização por Larguras Iguais - EWD

O método de Discretização por Larguras Iguais (EWD) faz a discretização de um intervalo, entre valores contínuos, dividindo em faixas de tamanhos iguais. Logo se existir um intervalo com valores contínuos $[a,b]$, e deseja particionar em R faixas de tamanhos iguais serão necessários $R - 1$ pontos de corte figura 3.

Para haver o ponto de corte antes tem que ser realizado a ordenação dos dados. A largura de cada faixa r_1, \dots, r_R na equação 1.6 é representada por w que é calculada pela

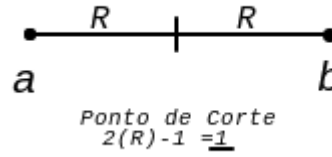


Figura 3 – Ponto de Corte (R-1)

diferença entre os limites superior e inferior do intervalo, dividido pela quantidade R de valores a serem gerados.

$$w = \frac{b - a}{R} \quad (1.6)$$

A variável w determina os pontos de corte (c_1, \dots, c_{R-1}) que irão delimitar o tamanho das faixas de valores. O primeiro ponto de corte, c_1 , é obtido através da soma do limite inferior a com a tamanho de w . E os pontos de corte seguintes são calculados pela soma do ponto de corte anterior com w .

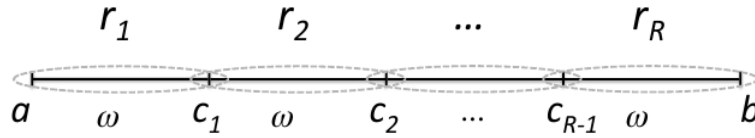
O valor de cada faixa será representado por i , onde i é o índice indicando a faixa. De acordo com a figura 4 para dividir o intervalo $[a, b]$ em R faixas será necessário de $R - 1$ pontos de corte.

$$c_i = \begin{cases} a + w, & \text{se } i = 1 \\ c_{i-1} + w, & \text{caso contrário} \end{cases} \quad (1.7)$$

O valor da faixa do intervalo $[a, c_1]$ será o valor discreto igual ao índice de sua faixa r_1 . Então, um valor na faixa r_1 terá o valor representado por 1(*um*), pois $i = 1$ é limite inferior mais largura da faixa, equação 1.7. E seguindo o mesmo raciocínio o valor da faixa $r_2 =]c_1, c_2]$ é representado por 2(*dois*), e conseqüentemente o valor que se encontra em uma faixa qualquer r_i será representado por i .

1.2.2 Discretização por Frequência Iguais - EFD

Esse outro método de discretização já possui uma abordagem diferente a do EWD, pois a idéia é manter a quantidade de elementos distintos, entre os pontos de corte, com o mesmo número. Dado um intervalo $[a, b]$ o número de faixas R e a quantidade de valores distintos ξ , onde $\xi \geq R$ o método EFD irá segmentar em R faixas de valores que possuem a mesma quantidade de elementos distintos λ . Então serão realizados $R - 1$ pontos de corte gerando R faixas de valores, (r_1, \dots, r_R) , com a mesma quantidade de elementos distintos λ . Para encontrar λ calcula-se o valor inteiro da divisão entre a quantidade de elementos

Figura 4 – Discretização EWD ³

distintos ξ pela quantidade de faixas de valores R , obtendo o número de elementos da faixa 1.8.

$$\lambda = \frac{\xi}{R} \quad (1.8)$$

Uma observação nesse método é quando ocorrer nos casos de uma amostragem possuir uma má distribuição de valores de um dado atributo, como um número significativo de repetições, isso, irá causar um desequilíbrio nas distribuições dos elementos.

Uma vez no intervalo $[a, b]$ de elementos ordenado e calculado λ contendo R elementos ($v_{[R]}$) pode-se determinar os pontos de corte (c_1, \dots, c_{R-1}) que são os delimitadores das faixas. Cada ponto de corte c_i pode ser calculado por $v_{i\lambda}$ 1.9.

$$\lambda = \frac{\xi}{R} \quad (1.9)$$

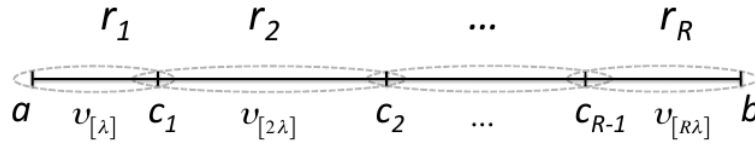
Como na sessão anterior do método EWD o valor que estiver no intervalo $[a, c_1]$ terá seu valor associado a um valor discreto igual ao índice i de sua faixa r_i conforme figura 5. Então, caso o valor esteja na faixa r_2 ele passará a ter o valor de seu índice i igual a 2(*dois*). De maneira consecutiva os valores que estiverem na faixa $r_3 =]c_2, c_3]$ terão valor 3(*três*). Uma outra observação desse método é que diferente do EWD, as faixas podem assumir faixas com tamanhos diferentes.

1.3 Trabalhos Correlatos

Esta sessão propõe relacionar outros trabalhos servindo de complemento teórico, como também leitura imprescindível, para entender a variedade de aplicações referente ao assunto de rotulação de dados. Mas ao longo da escrita desta proposta de mestrado verificou-se uma carência de pesquisas no âmbito de rotulação de dados, referente ao tema aqui proposto neste trabalho, pois acaba sendo redefinido o termo de rotulação.

O trabalho escrito por (LOPES; MACHADO; RABELO,) fez um estudo abordando o tema de rotulação de dados bastante significativo. Foi apresentado nesse trabalho o Problema de Rotulação, que representa também o problema proposto por esse trabalho, mas com abrangência e execução diferente do modelo (LOPES; MACHADO; RABELO,)

³ Figura extraída de (LOPES; MACHADO; RABELO,)

Figura 5 – Discretização EFD⁴

na figura 6 . Na pesquisa de (LOPES; MACHADO; RABELO,) é utilizado como entrada um conjunto dados onde é feito um agrupamento automático formando os clusters, e apresenta como saída um rótulo específico que melhor define o grupo formado. Esses rótulos são formados pela faixa de valor em conjunto com os atributos mais relevantes.

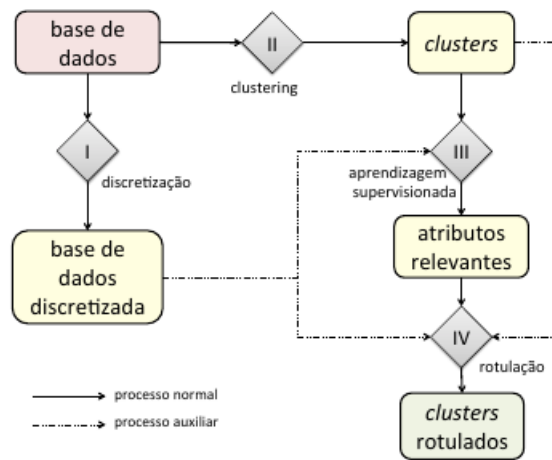


Figura 6 – Modelo (LOPES; MACHADO; RABELO,)

Outra pesquisa aplicada em rotulação está em (FILHO, 2015) onde aborda o mesmo Problema de Rotulação. Mas a atuação é diferenciada, pois o modelo, figura 7 procura diferenças existentes em cada grupo através da seleção dos elementos que representam o grupo, e depois é construído a faixa de valores. Os grupos são formados pelo algoritmo Fuzzy C-Means e após isso que é selecionado os atributos.

Em (LIMA, 2015) o problema em questão é fazer classificação e rotulação em uma base que possuem poucos elementos classificados. O método inicia com uma base dividida em elementos classificados(L) e não classificados(U). Após cada iteração o grupo L vai crescendo e automaticamente diminuindo o grupo U até que não tenha mais nenhum elemento em U. Após isso é realizado uma etapa de agrupamento, sem levar em consideração os dados classificados anteriormente. Terminada essa etapa é feito uma validação para saber quais os rótulos foram considerados corretos.

⁴ Figura extraída de (LOPES; MACHADO; RABELO,)

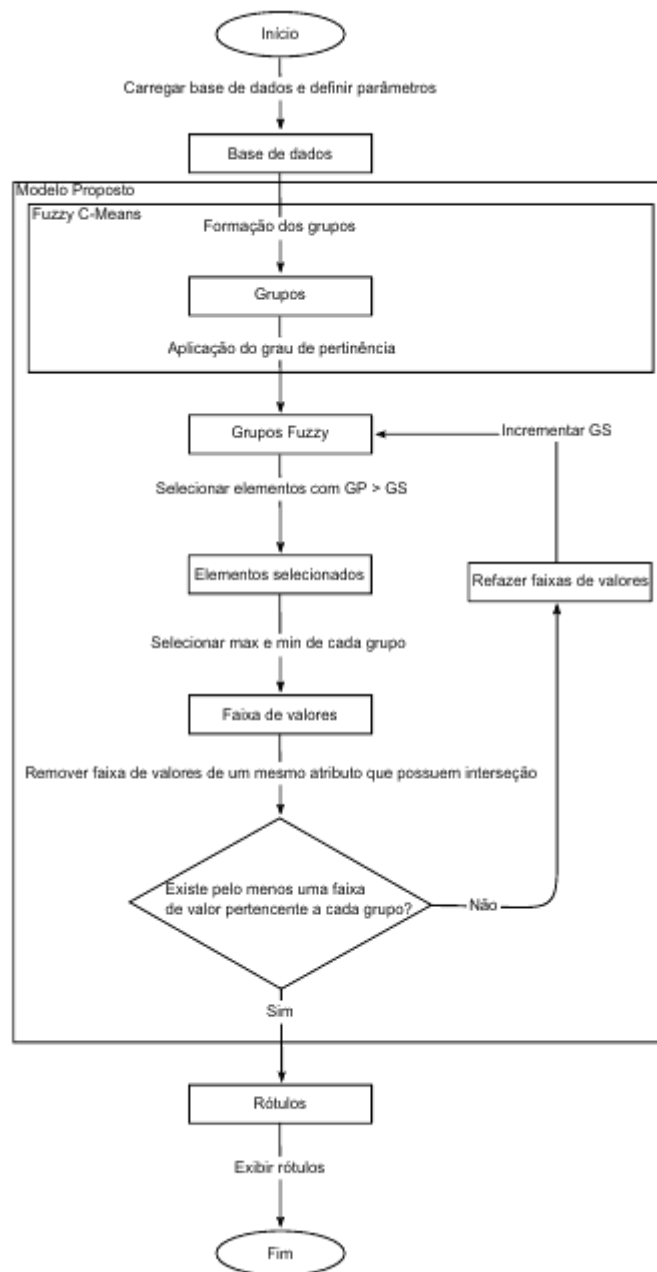


Figura 7 – Modelo (FILHO, 2015)

2 Metodologia / Materiais e Métodos

Esse capítulo abordará em uma sessão o problema proposto por esse trabalho, e logo em seguida, será apresentado um modelo de resolução. O objetivo ao final deste capítulo é poder resolver o problema exibindo seus passos e atribuindo a qualquer outro pesquisador todo o conhecimento necessário para replicar este trabalho através das informações produzidas aqui.

2.1 Considerações do Problema

A abordagem do problema referente a essa proposta de mestrado segue uma linha já pesquisada por (LOPES; MACHADO; RABELO,), que seria o **Problema de Rotulação**. Esse conceito, rotulação de dados, já é estudado na literatura na área de aprendizagem não-supervisionada, sessão 1.1.2, onde é comum os algoritmos lidarem com os agrupamentos dos dados, onde os grupos são criados a partir dos graus similaridade entre os elementos.

Muitas pesquisas realizadas na área de rotulação fazem referencia, de fato, a classificação do dados, e não da rotulação nos termos desse trabalho. Ao agrupar um conjunto de elementos por um determinado critério, esta havendo uma classificação desses elementos escolhidos, mas pouco se sabe, qual é a compreensão desses grupos, já classificados.

Tem-se então o real problema de rotulação, contudo seria necessário ter um rótulo definido para os grupos classificados para melhor compreender o porquê daquele grupo formado. Esse rótulo seria apresentação dos atributo(s) de maior relevância junto com a faixa, onde estaria nessa faixa, seus valores mais frequentes.

O Problema de Rotulação é formalmente definido como segue abaixo:

Dado um conjunto de clusters $C = \{c_1, \dots, c_k | K \geq 1\}$, de modo que cada cluster contém um conjunto de elementos $c_i = \{\vec{e}_1, \dots, \vec{e}_{n(c_i)} | n^{(c_i)} \geq 1\}$ que podem ser representados por um vetor de atributos definidos em \mathbb{R}^m e expresso por $\vec{e}^{c_i} = (a_1, \dots, a_m)$ e ainda que com $c_i \cap c_{i'} = \{0\}$ com $1 \leq i, i' \leq K$ e $i \neq i'$.¹

- K é o número de clusters;
- c_i é o i -ésimo cluster qualquer;
- n^{c_i} é o número de elementos do cluster c_i ;
- $\vec{e}_{n^{(c_i)}}$ se refere ao j -ésimo elemento pertencente ao cluster c_i ;
- m é a dimensão do problema;

¹ Extraída de (LOPES; MACHADO; RABELO,)

O estudo deste trabalho aproveita a perspectiva desse problema e cria um rótulo formado por seus atributos de mais relevância junto com os valores mais frequentes dess atribudo.

Como apresentado na sessão [1.3](#), o autor foca em rotulação automática de grupos utilizando aprendizagem de máquina supervisionada

3 Resultados e Discussão

Integer vel enim sed turpis adipiscing bibendum. Vestibulum pede dolor, laoreet nec, posuere in, nonummy in, sem. Donec imperdiet sapien placerat erat. Donec viverra. Aliquam eros. Nunc consequat massa id leo. Sed ullamcorper, lorem in sodales dapibus, risus metus sagittis lorem, non porttitor purus odio nec odio. Sed tincidunt posuere elit. Quisque eu enim. Donec libero risus, feugiat ac, dapibus eget, posuere a, felis. Quisque vel lectus ut metus tincidunt eleifend. Duis ut pede. Duis velit erat, venenatis vitae, vulputate a, pharetra sit amet, est. Etiam fringilla faucibus augue.

3.1 Base de Dados

Praesent facilisis, augue a adipiscing venenatis, libero risus molestie odio, pulvinar consectetur felis erat ac mauris. Nam vestibulum rhoncus quam. Sed velit urna, pharetra eu, eleifend eu, viverra at, wisi. Maecenas ultrices nibh at turpis. Aenean quam. Nulla ipsum. Aliquam posuere luctus erat. Curabitur magna felis, lacinia et, tristique id, ultrices ut, mauris. Suspendisse feugiat. Cras eleifend wisi vitae tortor. Phasellus leo purus, mattis sit amet, auctor in, rutrum in, magna. In hac habitasse platea dictumst. Phasellus imperdiet metus in sem. Vestibulum ac enim non sem ultricies sagittis. Sed vel diam.

3.2 Considerações Finais

Aenean velit sem, viverra eu, tempus id, rutrum id, mi. Nullam nec nibh. Proin ullamcorper, dolor in cursus tristique, eros augue tempor nibh, at gravida diam wisi at purus. Donec mattis ullamcorper tellus. Phasellus vel nulla. Praesent interdum, eros in sodales sollicitudin, nunc nulla pulvinar justo, a euismod eros sem nec nibh. Nullam sagittis dapibus lectus. Nullam eget ipsum eu tortor lobortis sodales. Etiam purus leo, pretium nec, feugiat non, ullamcorper vel, nibh. Sed vel elit et quam accumsan facilisis. Nunc leo. Suspendisse faucibus lacus.

Conclusões e Trabalhos Futuros

Proin non sem. Donec nec erat. Proin libero. Aliquam viverra arcu. Donec vitae purus. Donec felis mi, semper id, scelerisque porta, sollicitudin sed, turpis. Nulla in urna. Integer varius wisi non elit. Etiam nec sem. Mauris consequat, risus nec congue condimentum, ligula ligula suscipit urna, vitae porta odio erat quis sapien. Proin luctus leo id erat. Etiam massa metus, accumsan pellentesque, sagittis sit amet, venenatis nec, mauris. Praesent urna eros, ornare nec, vulputate eget, cursus sed, justo. Phasellus nec lorem. Nullam ligula ligula, mollis sit amet, faucibus vel, eleifend ac, dui. Aliquam erat volutpat.

Conclusões

Fusce vehicula, tortor et gravida porttitor, metus nibh congue lorem, ut tempus purus mauris a pede. Integer tincidunt orci sit amet turpis. Aenean a metus. Aliquam vestibulum lobortis felis. Donec gravida. Sed sed urna. Mauris et orci. Integer ultrices feugiat ligula. Sed dignissim nibh a massa. Donec orci dui, tempor sed, tincidunt nonummy, viverra sit amet, turpis. Quisque lobortis. Proin venenatis tortor nec wisi. Vestibulum placerat. In hac habitasse platea dictumst. Aliquam porta mi quis risus. Donec sagittis luctus diam. Nam ipsum elit, imperdiet vitae, faucibus nec, fringilla eget, leo. Etiam quis dolor in sapien porttitor imperdiet.

Cras pretium. Nulla malesuada ipsum ut libero. Suspendisse gravida hendrerit tellus. Maecenas quis lacus. Morbi fringilla. Vestibulum odio turpis, tempor vitae, scelerisque a, dictum non, massa. Praesent erat felis, porta sit amet, condimentum sit amet, placerat et, turpis. Praesent placerat lacus a enim. Vestibulum non eros. Ut congue. Donec tristique varius tortor. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Nam dictum dictum urna.

Phasellus vestibulum orci vel mauris. Fusce quam leo, adipiscing ac, pulvinar eget, molestie sit amet, erat. Sed diam. Suspendisse eros leo, tempus eget, dapibus sit amet, tempus eu, arcu. Vestibulum wisi metus, dapibus vel, luctus sit amet, condimentum quis, leo. Suspendisse molestie. Duis in ante. Ut sodales sem sit amet mauris. Suspendisse ornare pretium orci. Fusce tristique enim eget mi. Vestibulum eros elit, gravida ac, pharetra sed, lobortis in, massa. Proin at dolor. Duis accumsan accumsan pede. Nullam blandit elit in magna lacinia hendrerit. Ut nonummy luctus eros. Fusce eget tortor.

Trabalhos Futuros

Ut sit amet magna. Cras a ligula eu urna dignissim viverra. Nullam tempor leo porta ipsum. Praesent purus. Nullam consequat. Mauris dictum sagittis dui. Vestibulum sollicitudin consectetur wisi. In sit amet diam. Nullam malesuada pharetra risus. Proin lacus arcu, eleifend sed, vehicula at, congue sit amet, sem. Sed sagittis pede a nisl. Sed tincidunt odio a pede. Sed dui. Nam eu enim. Aliquam sagittis lacus eget libero. Pellentesque diam sem, sagittis molestie, tristique et, fermentum ornare, nibh. Nulla et tellus non felis imperdiet mattis. Aliquam erat volutpat.

Referências

- BARBER, D. *Bayesian Reasoning and Machine Learning*. [s.n.], 2011. ISSN 9780521518147. ISBN 9780511804779. Disponível em: <<http://ebooks.cambridge.org/ref/id/CBO9780511804779>>. Citado 3 vezes nas páginas 1, 4 e 8.
- CATLETT, J. Into Ordered Discrete Attributes. v. 3, n. 1989, p. 2006, 2006. Citado na página 8.
- DOUGHERTY, J.; KOHAVI, R.; SAHAMI, M. Supervised and Unsupervised Discretization of Continuous Features. *Machine Learning Proceedings 1995*, v. 0, p. 194–202, 1995. ISSN 0717-6163. Disponível em: <<http://linkinghub.elsevier.com/retrieve/pii/B9781558603776500323>>. Citado na página 8.
- FILHO, V. P. R. *Rotulacao de grupos utilizando conjuntos fuzzy*. Tese (Doutorado) — Universidade Federal do Piauí, 2015. Citado 3 vezes nas páginas 15, 11 e 12.
- HWANG, G. J.; LI, F. A Dynamic Method for Discretization of Continuous Attributes. *Lecture Notes in Computer Science - Intelligent Data Engineering and Automated Learning - IDEAL 2002: Third International Conference*, v. 2412/2002, p. 506, 2002. ISSN 16113349. Disponível em: <<http://www.springerlink.com/content/4n05b2n6x0cx4tlk>>. Citado na página 8.
- KOTSIANTIS, S.; KANELLOPOULOS, D. Discretization Techniques : A recent survey. *GESTS International Transactions on Computer Science and Engineering*, v. 32, n. 1, p. 47–58, 2006. Citado na página 8.
- LIMA, B. V. A. Método Semissupervisionado de Rotulação e Classificação Utilizando Agrupamento por Sementes e Classificadores. 2015. Citado na página 11.
- LOPES, L. A.; MACHADO, V. P.; RABELO, R. D. A. L. Automatic Labeling of Groupings through Supervised Machine Learning. Citado 4 vezes nas páginas 15, 10, 11 e 13.
- LUCCA, G. et al. Uma implementação do algoritmo Naïve Bayes para classificação de texto. *Centro de Ciências Computacionais - Universidade Federal do Rio Grande (FURG) Rio Grande - RS - Brasil*, p. 1–4, 2013. Citado na página 7.
- MADUREIRA, D. F. *Análise de sentimento para textos curtos*. Tese (Doutorado) — Fundacao Getulio Vargas, Rio de Janeiro, 2017. Citado na página 7.
- MCCALLUM, A.; NIGAM, K. A Comparison of Event Models for Naive Bayes Text Classification. 1997. Citado na página 7.
- MITCHELL, T. M. *Machine learning*. [S.l.]: McGraw-Hill Science/Engineering/Math, 1997. 432 p. ISSN 10450823. ISBN 9781577354260. Citado 2 vezes nas páginas 1 e 3.
- RAIMUNDO, L. R.; MATTOS, M. C. D.; WALESKA, P. O Algoritmo de Classificação CART em uma Ferramenta de Data Mining. 2008. Citado na página 6.

RUSSEL, S.; NORVIG, P. *Inteligência Artificial*. 3ª. ed. Rio de Janeiro: [s.n.], 2013. ISBN 9780136042594. Citado 2 vezes nas páginas 3 e 4.

WU, X. et al. *Top 10 algorithms in data mining*. [S.l.: s.n.], 2008. v. 14. 1–37 p. ISSN 02191377. ISBN 1011500701. Citado na página 7.

YOHANNES, Y.; WEBB, P. *Classification and Regression Trees, CART: A User Manual for Identifying Indicators of Vulnerability to Famine and Chronic Food Insecurity*. International Food Policy Research Institute, 1999. (Microcomputers in policy research). ISBN 9780896293373. Disponível em: <<https://books.google.com.br/books?id=7iuq4ikyNdoC>>. Citado na página 6.

Apêndices

APÊNDICE A – Primeiro Apêndice

Quisque facilisis auctor sapien. Pellentesque gravida hendrerit lectus. Mauris rutrum sodales sapien. Fusce hendrerit sem vel lorem. Integer pellentesque massa vel augue. Integer elit tortor, feugiat quis, sagittis et, ornare non, lacus. Vestibulum posuere pellentesque eros. Quisque venenatis ipsum dictum nulla. Aliquam quis quam non metus eleifend interdum. Nam eget sapien ac mauris malesuada adipiscing. Etiam eleifend neque sed quam. Nulla facilisi. Proin a ligula. Sed id dui eu nibh egestas tincidunt. Suspendisse arcu.

APÊNDICE B – Perceba que o texto do título desse segundo apêndice é bem grande

Maecenas dui. Aliquam volutpat auctor lorem. Cras placerat est vitae lectus. Curabitur massa lectus, rutrum euismod, dignissim ut, dapibus a, odio. Ut eros erat, vulputate ut, interdum non, porta eu, erat. Cras fermentum, felis in porta congue, velit leo facilisis odio, vitae consectetur lorem quam vitae orci. Sed ultrices, pede eu placerat auctor, ante ligula rutrum tellus, vel posuere nibh lacus nec nibh. Maecenas laoreet dolor at enim. Donec molestie dolor nec metus. Vestibulum libero. Sed quis erat. Sed tristique. Duis pede leo, fermentum quis, consectetur eget, vulputate sit amet, erat.

Donec vitae velit. Suspendisse porta fermentum mauris. Ut vel nunc non mauris pharetra varius. Duis consequat libero quis urna. Maecenas at ante. Vivamus varius, wisi sed egestas tristique, odio wisi luctus nulla, lobortis dictum dolor ligula in lacus. Vivamus aliquam, urna sed interdum porttitor, metus orci interdum odio, sit amet euismod lectus felis et leo. Praesent ac wisi. Nam suscipit vestibulum sem. Praesent eu ipsum vitae pede cursus venenatis. Duis sed odio. Vestibulum eleifend. Nulla ut massa. Proin rutrum mattis sapien. Curabitur dictum gravida ante.

Phasellus placerat vulputate quam. Maecenas at tellus. Pellentesque neque diam, dignissim ac, venenatis vitae, consequat ut, lacus. Nam nibh. Vestibulum fringilla arcu mollis arcu. Sed et turpis. Donec sem tellus, volutpat et, varius eu, commodo sed, lectus. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Quisque enim arcu, suscipit nec, tempus at, imperdiet vel, metus. Morbi volutpat purus at erat. Donec dignissim, sem id semper tempus, nibh massa eleifend turpis, sed pellentesque wisi purus sed libero. Nullam lobortis tortor vel risus. Pellentesque consequat nulla eu tellus. Donec velit. Aliquam fermentum, wisi ac rhoncus iaculis, tellus nunc malesuada orci, quis volutpat dui magna id mi. Nunc vel ante. Duis vitae lacus. Cras nec ipsum.

Anexos

ANEXO A – Nome do Primeiro Anexo

Sed mattis, erat sit amet gravida malesuada, elit augue egestas diam, tempus scelerisque nunc nisl vitae libero. Sed consequat feugiat massa. Nunc porta, eros in eleifend varius, erat leo rutrum dui, non convallis lectus orci ut nibh. Sed lorem massa, nonummy quis, egestas id, condimentum at, nisl. Maecenas at nibh. Aliquam et augue at nunc pellentesque ullamcorper. Duis nisl nibh, laoreet suscipit, convallis ut, rutrum id, enim. Phasellus odio. Nulla nulla elit, molestie non, scelerisque at, vestibulum eu, nulla. Ut odio nisl, facilisis id, mollis et, scelerisque nec, enim. Aenean sem leo, pellentesque sit amet, scelerisque sit amet, vehicula pellentesque, sapien.

ANEXO B – Nome de Outro Anexo

Phasellus id magna. Duis malesuada interdum arcu. Integer metus. Morbi pulvinar pellentesque mi. Suspendisse sed est eu magna molestie egestas. Quisque mi lorem, pulvinar eget, egestas quis, luctus at, ante. Proin auctor vehicula purus. Fusce ac nisl aliquam ante hendrerit pellentesque. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Morbi wisi. Etiam arcu mauris, facilisis sed, eleifend non, nonummy ut, pede. Cras ut lacus tempor metus mollis placerat. Vivamus eu tortor vel metus interdum malesuada.