

---

*Universidade Federal de São Carlos - UFSCar*  
*Departamento de Computação - DC*  
*Programa de Pós-Graduação em Ciência da Computação - PPGCC*

# Experimentos Mover x Weka

---



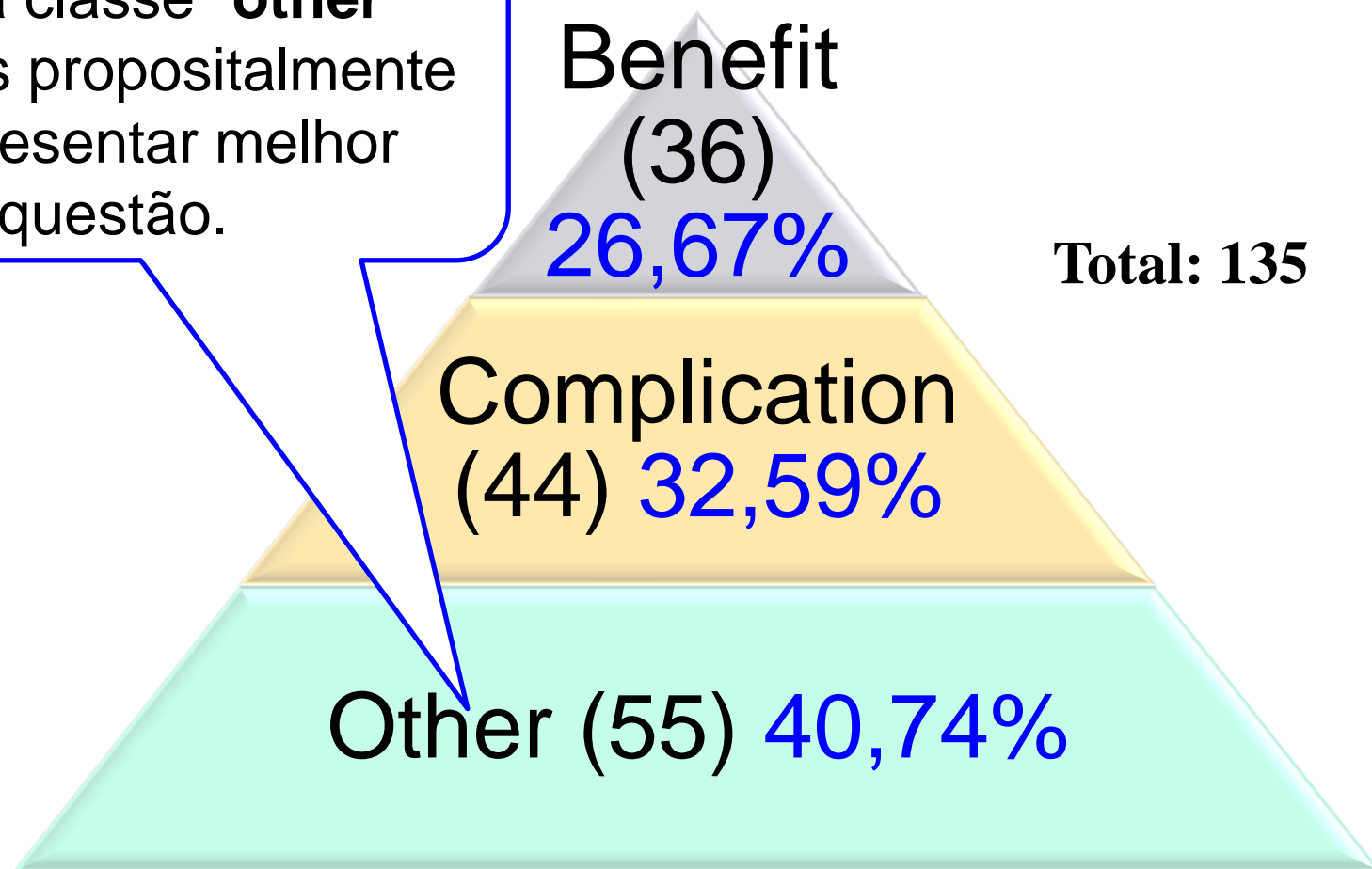
Aluno: Pablo Freire Matos  
Orientador: Dr. Ricardo Rodrigues Ciferri  
Coorientador: Dr. Thiago Alexandre S. Pardo  
Área: Banco de Dados

# Experimentos Mover x Weka

- Três Classes
  - Benefit
  - Complication
  - Other

# Quantidade de sentença por classe com 4 artigos (**Seleção Tendenciosa**)

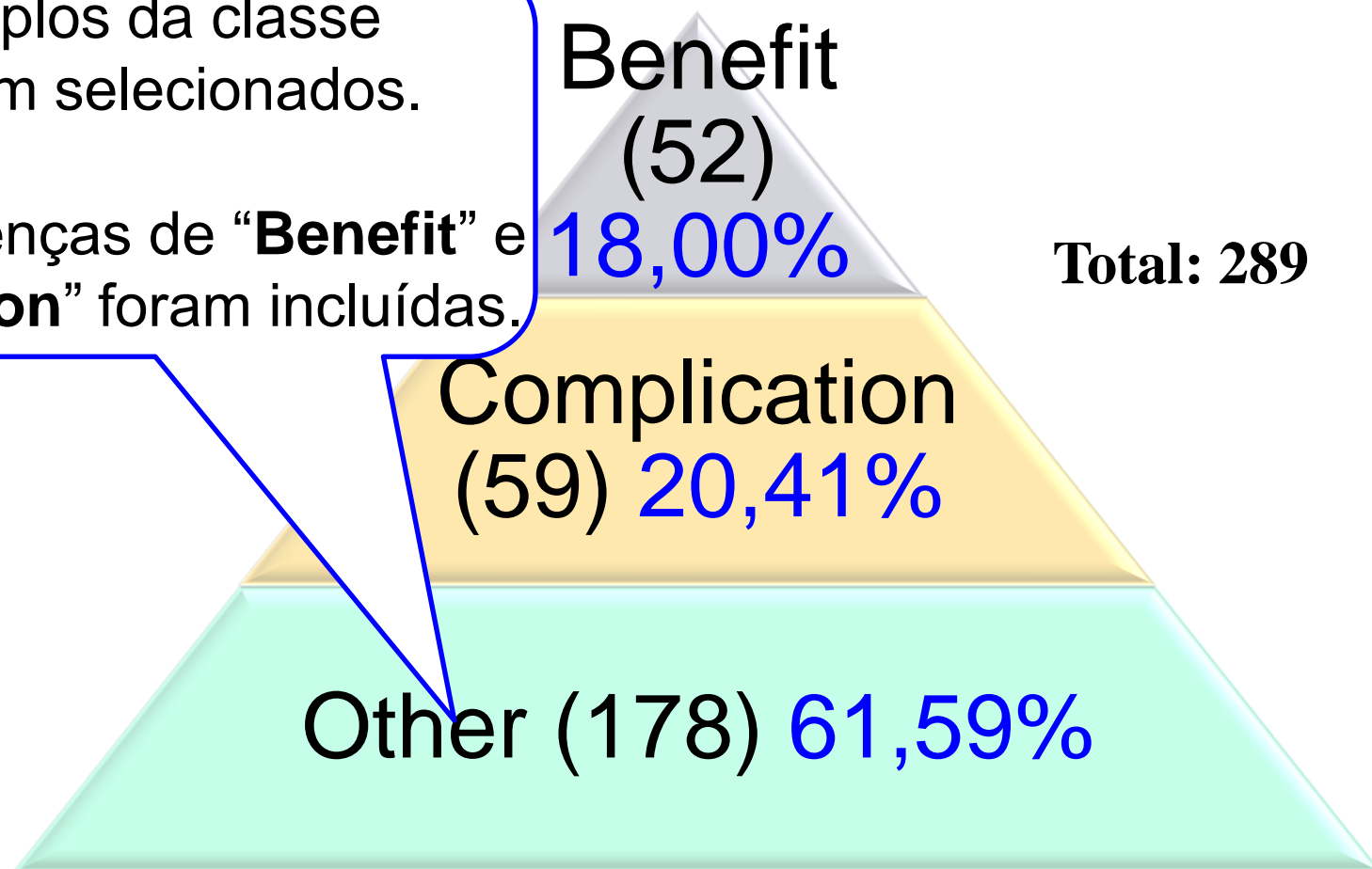
Exemplos da classe “**other**” selecionados propositalmente a fim de representar melhor a classe em questão.



# Quantidade de sentença por classe com 4 artigos (**Seleção Mais Sentenças**)

**Todos** exemplos da classe “**other**” foram selecionados.

Outras sentenças de “**Benefit**” e “**Complication**” foram incluídas.



# Passos para a Classificação no WEKA

1. Carregar os dados
2. Aplicar pré-processamento
  1. Gerar matriz atributo-valor
    1. Frequência, n-gramas, stopwords, stemmer ...
  2. Discretização
  3. RemoveMisclassified (Ruído)
  4. Randomize
  5. Resample (over sampling)
  6. Seleção de atributo
3. Selecionar Classificador

# Resultado Weka

- Matriz atributo-valor
  - Frequência mínima = 2
  - 1 a 3 gramas
- Classificação
  - Ganho de Informação (com ou sem)
  - 10-Fold Cross-Validation








Exemplos	Mover (Stratified CV) NB	Weka (CV) NB	Weka (CV) SVM
135 Algumas sentenças de "outros"	70,77%	78,52% (Ganho de 7,75%) 75,55% (com <b>stopword</b> ) 74,81% (com stemmer <b>teratedLovins</b> ) 76,29% (com stemmer <b>Lovins</b> ) 75,55% (com stemmer <b>Snowball</b> ) 75,55% (stopword e Snowball)	71,85%
289 Todas sentenças de "outros"	68,21%	70,59% (Ganho de 2,38%)	72,66%

# Configuração dos Experimentos

- Pré-processamento
  - Matriz atributo-valor
    - Frequência mínima = 2
    - 1 a 3 gramas
    - Sem *stopword* e *stemmer*
- Seleção de Atributo
  - Ganho de Informação
- Seis classificadores escolhidos
  - **Valores dos parâmetros padrão**
  - 10-Fold Cross-Validation

**Parâmetros definidos nos experimentos testados**

# Filtro: Randomize

Classificador	135 exemplos		289 exemplos	
	Sem Filtro	Com Filtro	Sem Filtro	Com Filtro
Naïve Bayes	78,51%	 75,55%	70,59%	71,63%
OneR	41,48%	 40,74%	63,67%	 62,63%
Prism	45,92%	 40,00%	48,79%	 46,71%
J48	50,37%	55,55%	64,71%	65,05%
ID3	54,07%	56,29%	62,28%	 59,86%
SVM	71,85%	71,85%	72,66%	 67,82%



# Eliminar Ruído

Classifier  
Choose **NaiveBayes**

Test options  
 Use training set  
 Supplied test set Set...  
 Cross-validation Folds 10  
 Percentage split % 66  
More options...

(Nom) class  
Start Stop

Result list (right-click for options)  
21:26:59 - bayes.NaiveBayes

Classifier output

Correctly Classified Instances	127	94.0741 %
Incorrectly Classified Instances	8	5.9259 %
Kappa statistic	0.9089	
Mean absolute error	0.0467	
Root mean squared error	0.1895	
Relative absolute error	10.666 %	
Root relative squared error	40.5048 %	
Total Number of Instances	135	

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area	Class
	0.886	0	1	0.886	0.94	0.998	1_C
	0.917	0	1	0.917	0.957	0.999	2_F
	1	0.1	0.873	1	0.932	1	5_C
Weighted Avg.	0.941	0.041	0.948	0.941	0.941	0.999	

=== Confusion Matrix ===

```

a  b  c  <-- classified as
39  0  5 | a = 1_Complication_44
 0 33  3 | b = 2_Benefit_36
 0  0 55 | c = 5_Other_55

```

# Filtro: RemoveMisclassified do NB

Classificador	135 exemplos			289 exemplos		
	Sem Filtro		Com Filtro	Sem Filtro		Com Filtro
Naïve Bayes	78,51%		↓ 76,14%	70,59%		-----
OneR	41,48%		↑ 52,29%	63,67%		-----
Prism	45,92%		51,37%	48,79%		-----
J48	50,37%		66,05%	64,71%		-----
ID3	54,07%		75,22%	62,28%		-----
SVM	71,85%		78,89%	72,66%		-----
Distribuição das classes	135 exemplos C[0] C[1] C[2] 44 36 55		109 exemplos C[0] C[1] C[2] 38 16 55	289 exemplos C[1] C[2] C[3] 59 52 178		176 exemplos C[0] C[1] C[2] 0 0 176




# Filtro: RemoveMisclassified (RM) de cada classificador

135 exemplos		
Classificador	Sem Filtro	Com Filtro
Naïve Bayes	78,51%	76,14% (109 exemplos)
OneR	41,48%	100% (64 ex.)
Prism	45,92%	45,92% (135 ex.)
J48	50,37%	65,28 (121 ex.)
ID3	54,07%	54,07% (135 ex.)
SVM	71,85%	71,85% (135 ex.)

289 exemplos															
Sem filtro				<del>RM do NB</del>				<del>RM do OneR</del>				RM do J48			
Dist	C[1]	C[2]	C[3]	Dist	C[1]	C[2]	C[3]	Dist	C[1]	C[2]	C[3]	Dist	C[1]	C[2]	C[3]
3	59	52	178	3	0	0	176	3	13	0	174	3	44	31	176

# Filtro: RemoveMisclassified do J48

<b>Mover (Stratified CV)</b> 135 exemplos	<b>Mover (Stratified CV)</b> 289 exemplos
<b>70,77% NB</b>	<b>68,21% NB</b>

	135 exemplos - ganho de <b>5,26%</b>		289 exemplos - Ganho de <b>11,07%</b>	
<b>Classificador</b>	<b>Sem Filtro</b>	<b>Com Filtro</b>	<b>Sem Filtro</b>	<b>Com Filtro</b>
Naïve Bayes	78,51%	 <b>76,03%</b>	70,59%	 <b>74,90%</b>
OneR	41,48%	 50,41%	63,67%	70,92%
Prism	45,92%	45,45%	48,79%	62,95%
J48	50,37%	<b>65,29%</b>	64,71%	<b>79,28%</b>
ID3	54,07%	63,64%	62,28%	68,92%
SVM	71,85%	<b>72,73%</b>	<b>72,66%</b>	<b>78,88%</b>
Distribuição das Classes	135 exemplos C[0] C[1] C[2] 44 36 55	109 exemplos C[0] C[1] C[2] 38 16 55	289 exemplos C[0] C[1] C[2] 59 52 178	251 exemplos C[0] C[1] C[2] 44 31 176

Preprocess Classify Cluster Associate Select attributes Visualize

Open file...

Open URL...

Open DB...

# Balanceamento dos exemplos

Filter

- weka
  - filters
    - AllFilter
    - MultiFilter
    - supervised
      - attribute
      - instance
        - Resample**
        - SMOTE**
        - Spreadsubsample
        - StratifiedRemoveFolds
    - unsupervised

2 métodos de  
over-sampling

Apply

amentos\_UNID...

Selected attribute

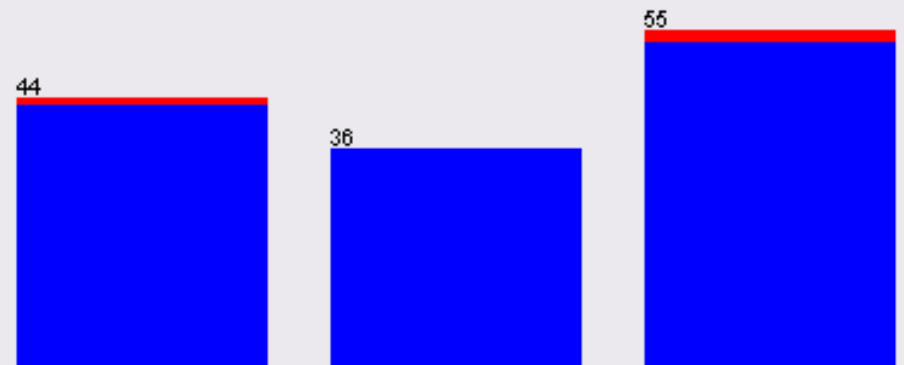
Name: class  
Missing: 0 (0%)      Distinct: 3      Type: Nominal  
Unique: 0 (0%)

No.	Label	Count	Weight
1	1_Complication_44	44	44.0
2	2_Benefit_36	36	36.0
3	5_Other_55	55	55.0

Antes

Class: younger than 2 (Nom)

Visualize All



Filter...

Remove filter

Close

Status

OK

Log

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter

Choose Resample 8 0.0 -5 1 -Z 100.0 Apply

Depois

Current relation

Relation: C:\Documents and Settings\... Instances: 135 Sum of weights: 135

Attributes

- All
  - 1 da
  - 2 1
  - 3 1 of
  - 4 1 recurrent
  - 5 1 recurrent stroke
  - 6 13
  - 7 2
  - 8 2 the
  - 9 2 treatment
  - 10 2 treatment was
  - 11 2 years
  - 12 2 years of
  - 13 3
- Remove

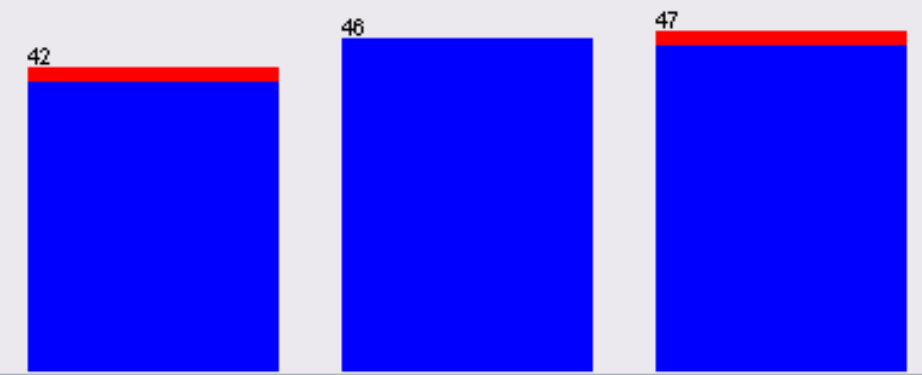
Mantém a quant. original de exemplos

Selected attribute

Name: class Type: Nominal Missing: 0 (0%) Distinct: 3 Unique: 0 (0%)

No.	Label	Count	Weight
1	1_Complication_44	42	42.0
2	2_Benefit_36	46	46.0
3	5_Other_55	47	47.0

Class: younger than 2 (Nom) Visualize All



Status

OK

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter Choose Resample -B 1.0 -S 1 -Z 100.0 Apply

Depois

Current relation Relation: C:\Documents and Settings\... Instances: 135 Sum of weights: 135

Selected attribute Name: class Type: Nominal Missing: 0 (0%) Distinct: 3 Unique: 0 (0%)

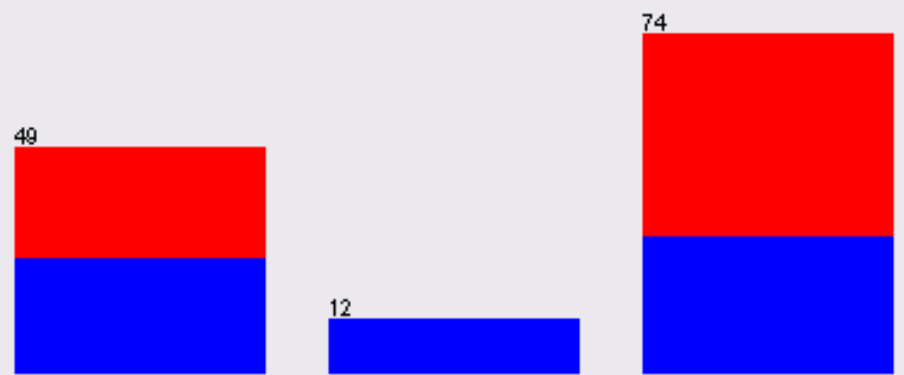
Attributes All

Mantém a quant. original de exemplos

No.	Label	Count	Weight
1	1_Complication_44	49	49.0
2	2_Benefit_36	12	12.0
3	5_Other_55	74	74.0

No.	
1	class
2	1
3	1 of
4	1 recurrent
5	1 recurrent stroke
6	13
7	2
8	2 the
9	2 treatment
10	2 treatment was
11	2 years
12	2 years of
13	3

Class: younger than 2 (Nom) Visualize All



# Resumo dos resultados: 3 classes

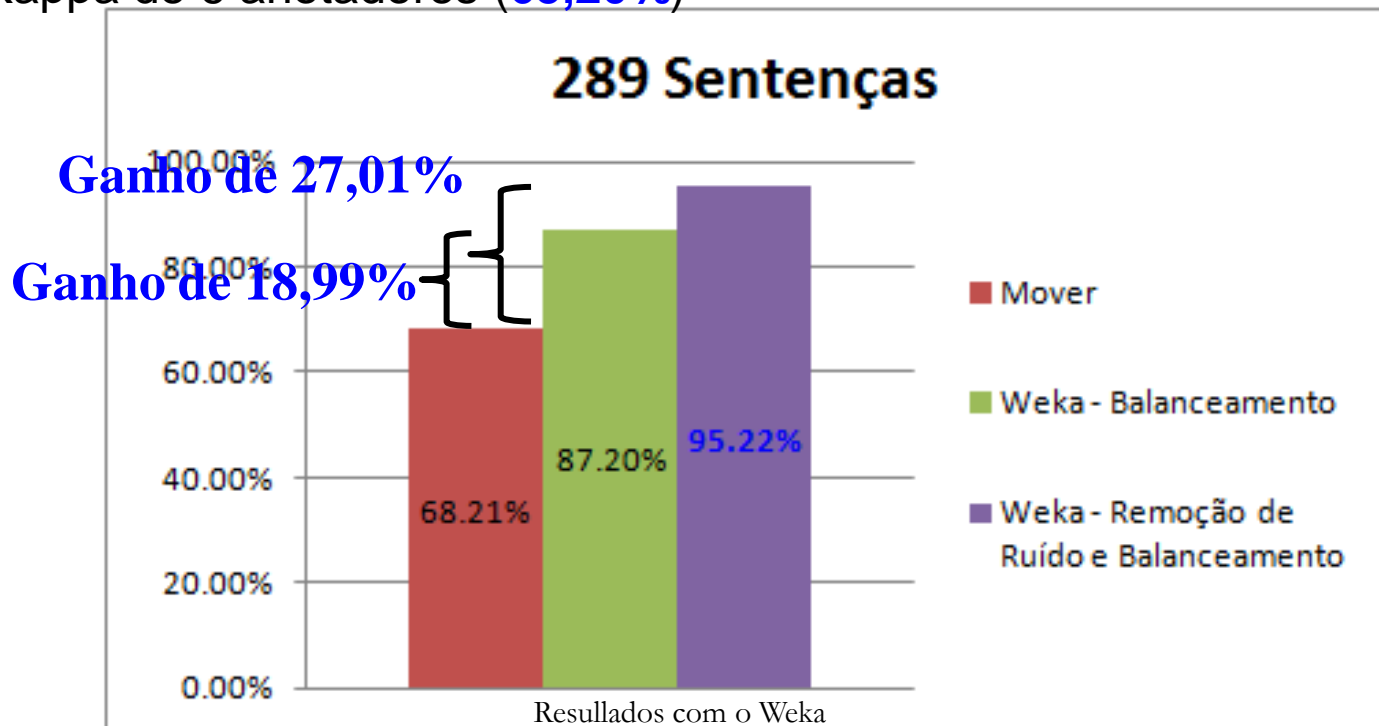
## Seleção de Atributo: Meta

Exemplos	Mover	Sem Filtro	Resample bias 0	RemoveMisclassified (J48) e Resample bias 0
135	70,77% NB	54,07% ID3 71,85% SVM <b>78,52% NB</b> 135 exemplos: C[0] C[1] C[2] 44 36 55	83,70% NB 84,44% SVM <b>85,19% ID3</b> 135 exemplos: C[0] C[1] C[2] 45 47 43 Ganho de <b>6,67%</b>	74,38% ID3 80,17% NB <b>85,12% SVM</b> 121 exemplos: C[0] C[1] C[2] 41 40 40 Ganho de <b>6,6%</b>
289	68,21% NB	64,71% J48 70,59% NB <b>72,66% SVM</b> 289 exemplos: C[0] C[1] C[2] 59 52 178	80,97% ID3 83,04% NB <b>87,20% SVM</b> 289 exemplos: C[0] C[1] C[2] 93 106 90 Ganho de <b>14,54%</b>	88,84% NB 90,04% ID3 <b>95,22% SVM</b> 251 exemplos: C[0] C[1] C[2] 82 91 78 Ganho de <b>22,56%</b>

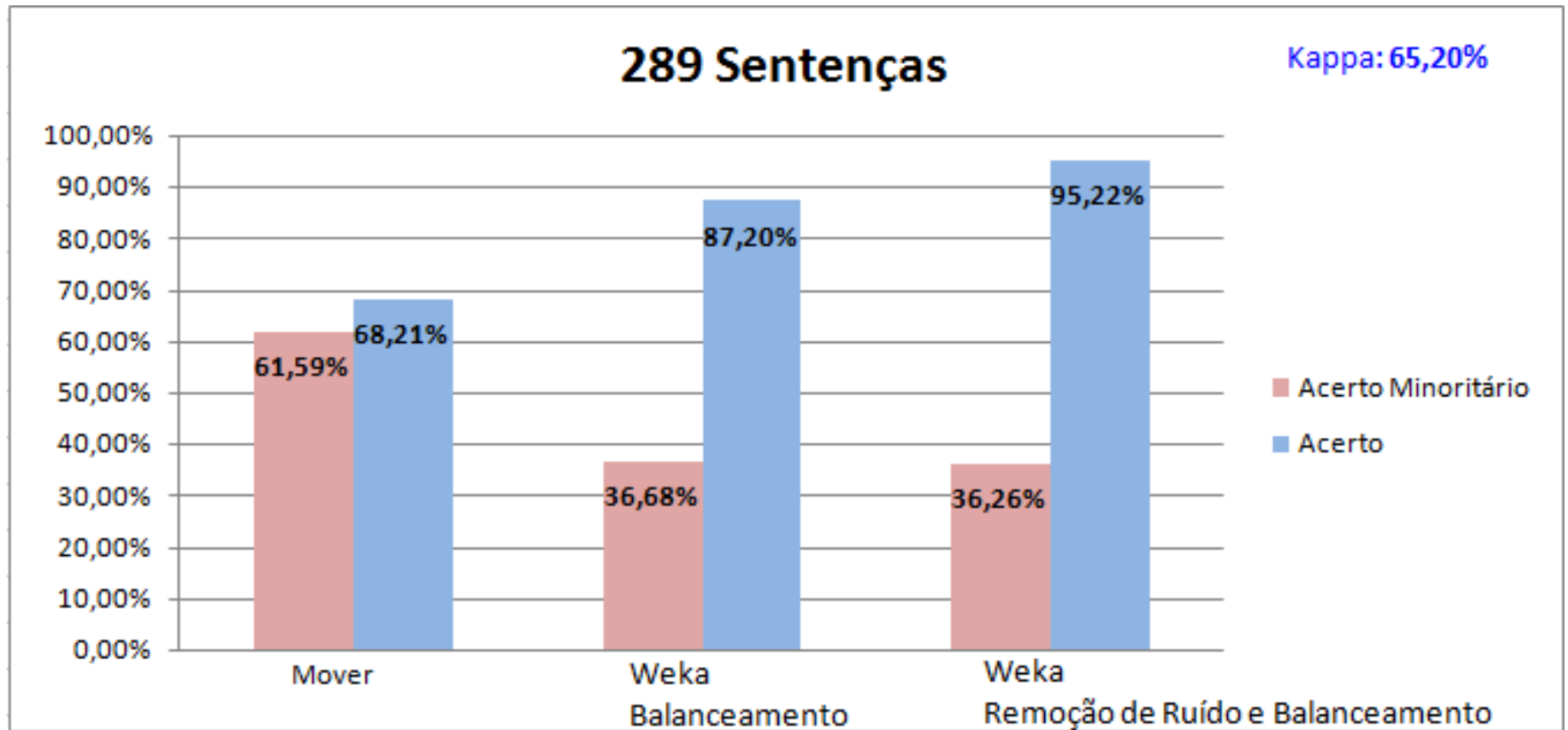


# Análise dos Resultados - Classificação

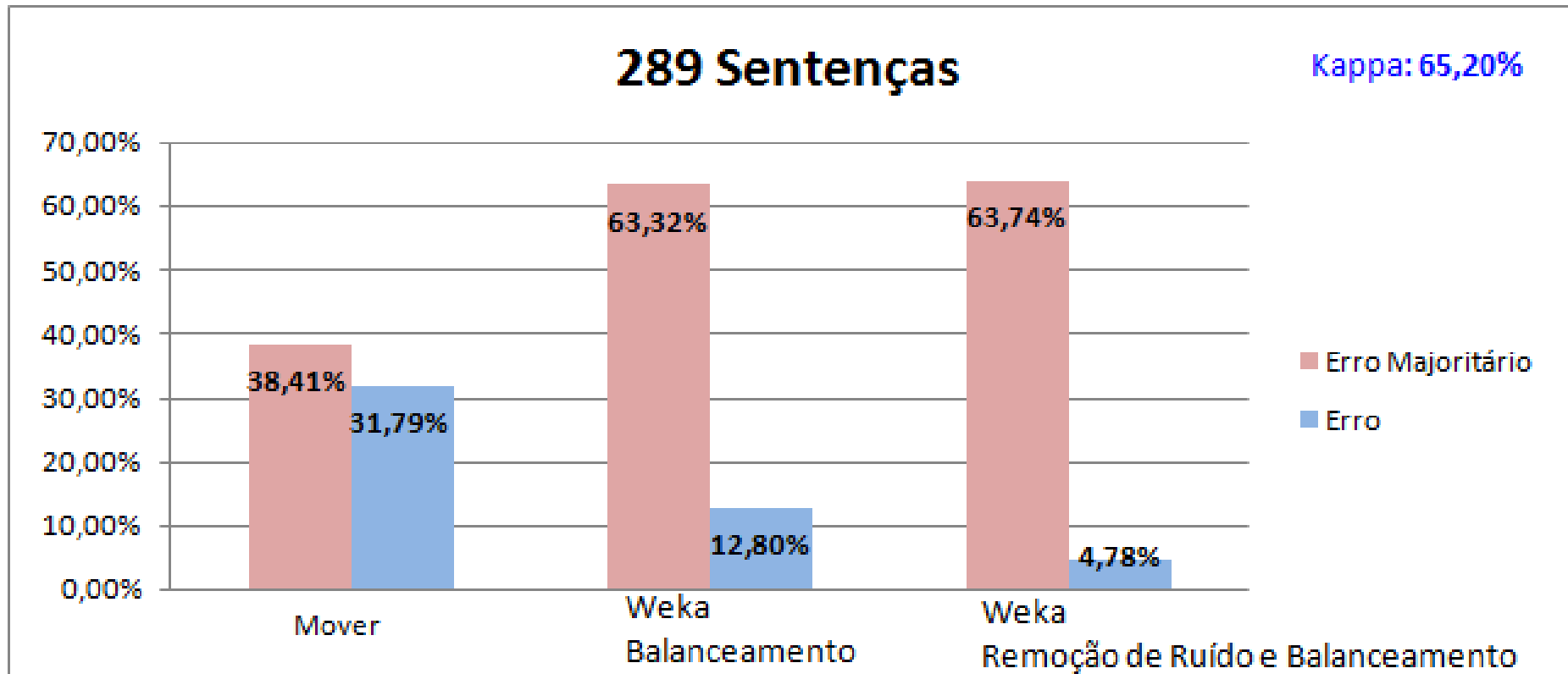
- Taxa de Precisão (P)
  - $Baseline \leq P \leq Topline$
- *Baseline*: Mover
- *Topline*: Taxa de precisão humana
  - 50 sentenças (efeito positivo, efeito negativo e outros)
  - Kappa de 6 anotadores (**65,20%**)



# Análise dos Resultados - Classificação



# Análise dos Resultados - Classificação



# Considerações Sobre a Classificação

- Pré-processamento: 3 classes (*Complicação*, *Benefício* e *Outro*)
  - Matriz atributo-valor
    - Frequência mínima = 2
    - 1 a 3 gramas
    - sem stopword e sem stemmer
  - Eliminar Ruído com **J48**
  - Balanceamento (Bias 1 em Java - Mantém a distribuição das classes)
  - Seleção de Atributo (Meta)
    - Ganho de Informação
- Classificador
  - Naive Bayes
  - Support Vector Machine

# Considerações Finais

## ■ Pré-processamento

1. Gerar matriz atributo-valor
  1. Frequência, n-gramas, stopwords, stemmer ...
  2. **IDF e TF transform, lematização**
2. Converter numérico para nominal
3. Randomize
4. RemoveMisclassified
5. Resample (over sampling)
6. **Igualar os exemplos (under sampling)**
7. Seleção de atributo

# Discussão com 135 sentenças

- Randomize
  - influencia **negativamente** o **NB**
  - **NÃO** influencia o **SVM**
- Attribute Selection
  - **NÃO** influencia o **NB**
  - **NÃO** influencia o **SVM**
  - Influencia **positivamente** todos os outros algoritmos testados

# Discussão (cont. 1) - com 135 sentenças

Filtro	Problema	Observações
<b>Seleção de Atributo</b>	<ul style="list-style-type: none"><li>• <b>Filter</b> e <b>Meta</b> geram atributos diferentes</li><li>• O <b>Filter</b> é melhor do que o <b>Meta</b></li></ul> <p>Resultados diferentes</p>	<ul style="list-style-type: none"><li>• O <b>Filter</b> não interfere positivamente nem negativamente no resultado do <b>NB</b> e <b>SVM</b></li><li>• Entretanto, o <b>Filter</b> melhora o resultado dos outros algoritmos (J48, ID3, Prism, OneR)</li></ul>

# Discussão (cont. 2) - 135 sentenças

Filtro	Problema	Definição Bias
<b>Resample</b>	<p>A implementação do Bias 0 no Explorer é similar a implementação do Bias 1 no Java</p> <p>Resultados diferentes - Algo está errado!!!</p>	<p>A value of 0 leaves the class distribution as-is.</p> <p>A value of 1 ensures the class distribution is uniform in the output data.</p>

## Distribuição Inicial das sentenças

C[0]	C[1]	C[2]
44	36	55

Total: 135 sentenças

	Explorer			Java		
<b>Bias 0</b>	C[0] 42	C[1] 46	C[2] 47	C[0] 49	C[1] 26	C[2] 60
<b>Bias 1</b>	C[0] 49	C[1] 12	C[2] 74	C[0] 45	C[1] 47	C[2] 43



# Discussão (cont. 3) – 135 sentenças

## Filtro: Resample

	Explorer			Java		
Bias 0	C[0] 49	C[1] 12	C[2] 74	C[0] 49	C[1] 26	C[2] 60
Bias 1	C[0] 42	C[1] 46	C[2] 47	C[0] 45	C[1] 47	C[2] 43

	Explorer	Java
Bias 0	J48 - 85,18 % - 3º melhor NB - 88,14% - 2º melhor SVM - 92,59%	NB - 78,52% - 2º pior ID3 - 86,67% - 2º melhor SVM - 91,11%
Bias 1	ID3 - 77,03 % - 3º melhor SVM - 82,96% - 2º melhor NB - 83,70%	NB - 83,70% - 3º melhor SVM - 84,44% - 2º melhor ID3 - 85,19%

# Referências

- ANTHONY, L.; LASHKIA, G. V. Mover: a machine learning tool to assist in the reading and writing of technical papers. **IEEE Transactions on Professional Communication**, v. 46, n. 3, p. 185-193, 2003.
- BURSTEIN, J.; MARCU, D.; KNIGHT, K. Finding the WRITE stuff: automatic identification of discourse structure in student essays. **Intelligent Systems**, IEEE, v. 18, n. 1, p. 32-39, 2003.
- BATISTA, G. E. A. P. A.; PRATI, R. C.; MONARD, M. C. A study of the behavior of several methods for balancing machine learning training data. **SIGKDD Explorations Newsletter**, v. 6, n. 1, p. 20-29, 2004.
- HEY, D. F.; FELTRIM, V. D. Uma investigação sobre a aplicação de algoritmos de aprendizado à classificação de papéis retóricos. In: VIII Fórum de Informática e Tecnologia de Maringá, XI Mostra de Trabalhos de Informática, 2008, Maringá. **Anais...** Universidade Estadual de Maringá, 2008. p. 94-104.

# Referências (Cont.)

- WITTEN, I. H.; FRANK, E. **Data mining**: practical machine learning tools and techniques with Java implementations. 2nd ed. San Francisco, CA: Morgan Kaufmann, 2005. 525 p.

---

*Universidade Federal de São Carlos - UFSCar*  
*Departamento de Computação - DC*  
*Programa de Pós-Graduação em Ciência da Computação - PPGCC*

# Experimentos Mover x Weka

---

**Obrigado!!!**

Aluno: Pablo Freire Matos  
Orientador: Dr. Ricardo Rodrigues Ciferri  
Coorientador: Dr. Thiago Alexandre S. Pardo  
Área: Banco de Dados