

# ADDRESSING MULTI-CLASS PROBLEMS BY BINARIZATION. NOVEL APPROACHES

Alberto Fernández  
Mikel Galar  
Francisco Herrera

State-of-the-art on One-vs-One, One-vs-All. Novel approaches.

# Outline

2/81

1. Introduction
2. Binarization
  - ▣ Decomposition strategies (One-vs-One, One-vs-All and Others)
  - ▣ State-of-the-art on Aggregations
    - One-vs-One
    - One-vs-All
3. Experimental Study
  - ▣ Experimental Framework
  - ▣ Results and Statistical Analysis
4. Discussion: Lessons Learned and Future Work
5. Conclusions for OVO vs OVA
6. Novel Approaches for the One-vs-One Learning Scheme
  - ▣ Dynamic OVO: Avoiding Non-competence
  - ▣ Distance-based Relative Competence Weighting Approach (DRCW-OVO)
  - ▣ Difficult Classes Problem in OVO Strategy

# Outline

3/81

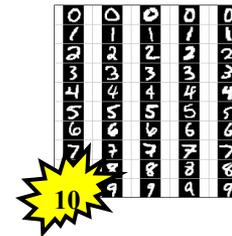
1. Introduction
2. Binarization
  - ▣ Decomposition strategies (One-vs-One, One-vs-All and Others)
  - ▣ State-of-the-art on Aggregations
    - One-vs-One
    - One-vs-All
3. Experimental Study
  - ▣ Experimental Framework
  - ▣ Results and Statistical Analysis
4. Discussion: Lessons Learned and Future Work
5. Conclusions for OVO vs OVA
6. Novel Approaches for the One-vs-One Learning Scheme
  - ▣ Dynamic OVO: Avoiding Non-competence
  - ▣ Distance-based Relative Competence Weighting Approach (DRCW-OVO)

# Introduction

## □ Classification

- 2 class of classification problems:
  - Binary: medical diagnosis (yes / no)
  - Multicategory: Letter recognition (A, B, C...)
- Binary problems are usually easier
- Some classifiers do not support multiple classes
  - SVM, PDFC...

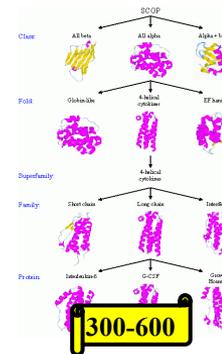
Digit recognition



Phoneme recognition

i:	ɪ	ʊ	u:	ɪə	eɪ
READ	SET	BOOK	TOO	HIRE	SWAY
e	ə	ɜ:	ɔ:	oʊ	ɔɪ
MEN	AMERICA	WRONG	SECRET	TOUR	BOY
æ	ʌ	ɑ:	ɒ	eə	aɪ
CAT	BUT	FIGHT	BOY	WARR	ME
p	b	t	d	f	ç
TOP	BATH	TIME	DO	CRUNCH	SHOES
θ	ð	s	z	ʃ	ʒ
THINK	THE	SEE	SO	SECRET	CASUAL
η	h	l	r	w	j
TRAY	BELLO	LIVE	READ	WINDOW	JES

Automated protein classification



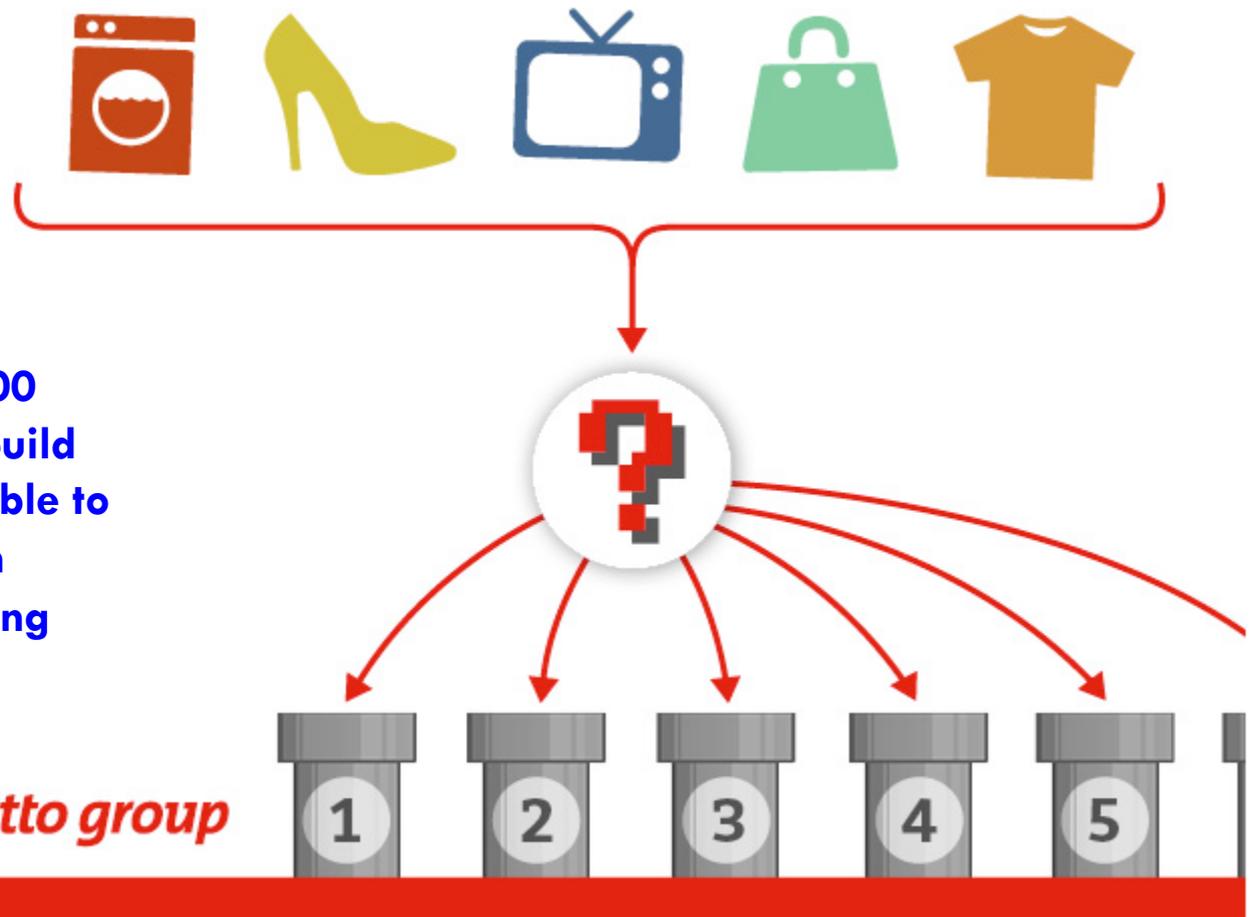
Object recognition



# Introduction. Ejemplo

## Una aplicación real en KAGGLE de Problema Multiclase

For this competition, we have provided a dataset with 93 features for more than 200,000 products. The objective is to build a predictive model which is able to distinguish between our main product categories. The winning models will be open sourced.



# Introduction. Ejemplo

## Una aplicación real en KAGGLE de Problema Multiclase

### Submission Format

You must submit a csv file with the product id, all candidate class names, and a probability for each class. The order of the rows does not matter. The file must have a header and should look like the following:

```
id,Class_1,Class_2,Class_3,Class_4,Class_5,Class_6,Class_7,Class_8,Class_9
1,0.0,0.0,0.0,0.0,1.0,0.0,0.0,0.0,0.0
2,0.0,0.2,0.3,0.3,0.0,0.0,0.1,0.1,0.0
...
etc.
```

1 9 8 7 teams

2 1 1 7 players

1 5 5 0 2 entries

**Started:** 3:56 pm, Tuesday 17 March 2015 UTC

**Ends:** 11:59 pm, Monday 18 May 2015 UTC (62 total days)

**Points:** this competition awards standard [ranking points](#)

**Tiers:** this competition counts towards [tiers](#)

# Introduction. Ejemplo

## Una aplicación real en KAGGLE de Problema Multiclase

### Evaluation

Submissions are evaluated using the multi-class logarithmic loss. Each product has been labeled with one true category. For each product, you must submit a set of predicted probabilities (one for every category). The formula is then,

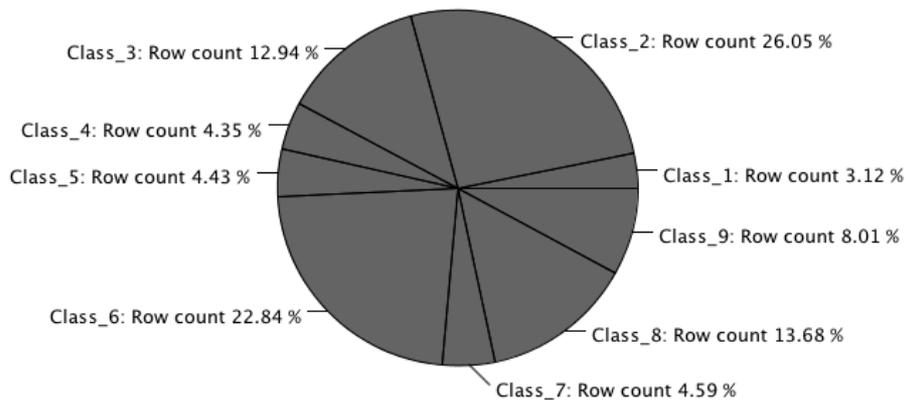
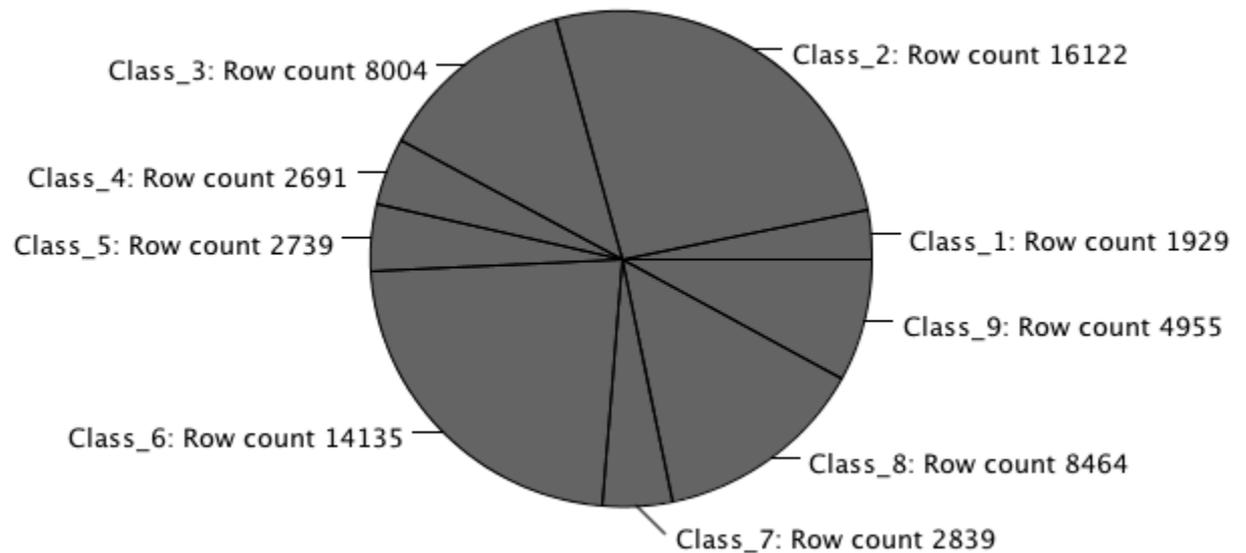
$$\text{logloss} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ij} \log(p_{ij}),$$

where  $N$  is the number of products in the test set,  $M$  is the number of class labels,  $\log$  is the natural logarithm,  $y_{ij}$  is 1 if observation  $i$  is in class  $j$  and 0 otherwise, and  $p_{ij}$  is the predicted probability that observation  $i$  belongs to class  $j$ .

The submitted probabilities for a given product are not required to sum to one because they are rescaled prior to being scored (each row is divided by the row sum). In order to avoid the extremes of the log function, predicted probabilities are replaced with  $\max(\min(p, 1 - 10^{-15}), 10^{-15})$ .

# Introduction. Ejemplo

## Una aplicación real en KAGGLE de Problema Multiclase



# Introduction. Ejemplo

## Una aplicación real en KAGGLE de Problema Multiclase

*otto group*

\$10,000 • 2,296 teams

### Otto Group Product Classification Challenge

Tue 17 Mar 2015

Enter/Merge by  
Mon 18 May 2015 (32 days to go)

Dashboard

## Public Leaderboard - Otto Group Product Classification Challenge

This leaderboard is calculated on approximately 70% of the test data.  
The final results will be based on the other 30%, so the final standings may be different.

See someone using multiple accounts  
[Let us know](#)

#	Δ1w	Team Name <small>* in the money</small>	Score <small>?</small>	Entries	Last Submission UTC (Best - Last Submission)
1	—	i dont know <small>👤 *</small>	0.39067	67	Thu, 16 Apr 2015 21:37:26
2	—	team <small>👤 *</small>	0.40017	20	Thu, 16 Apr 2015 14:45:41
3	new	tk <small>*</small>	0.40110	1	Thu, 16 Apr 2015 16:54:01
4	↓1	IzuiT	0.40311	43	Thu, 16 Apr 2015 05:29:26
5	↓1	Hoang Duong	0.40382	32	Thu, 16 Apr 2015 05:44:21 (-8.8d)
6	↓1	Nicholas Guttenberg	0.40857	55	Thu, 16 Apr 2015 15:46:43 (-2.6d)

# Outline

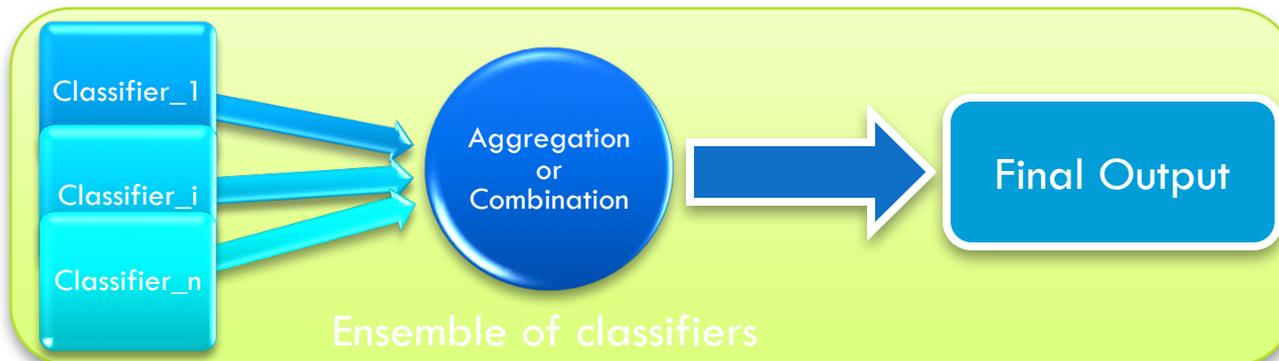
10/81

1. Introduction
2. Binarization
  - ▣ Decomposition strategies (One-vs-One, One-vs-All and Others)
  - ▣ State-of-the-art on Aggregations
    - One-vs-One
    - One-vs-All
3. Experimental Study
  - ▣ Experimental Framework
  - ▣ Results and Statistical Analysis
4. Discussion: Lessons Learned and Future Work
5. Conclusions for OVO vs OVA
6. Novel Approaches for the One-vs-One Learning Scheme
  - ▣ Dynamic OVO: Avoiding Non-competence
  - ▣ Distance-based Relative Competence Weighting Approach (DRCW-OVO)

# Binarization

11/81

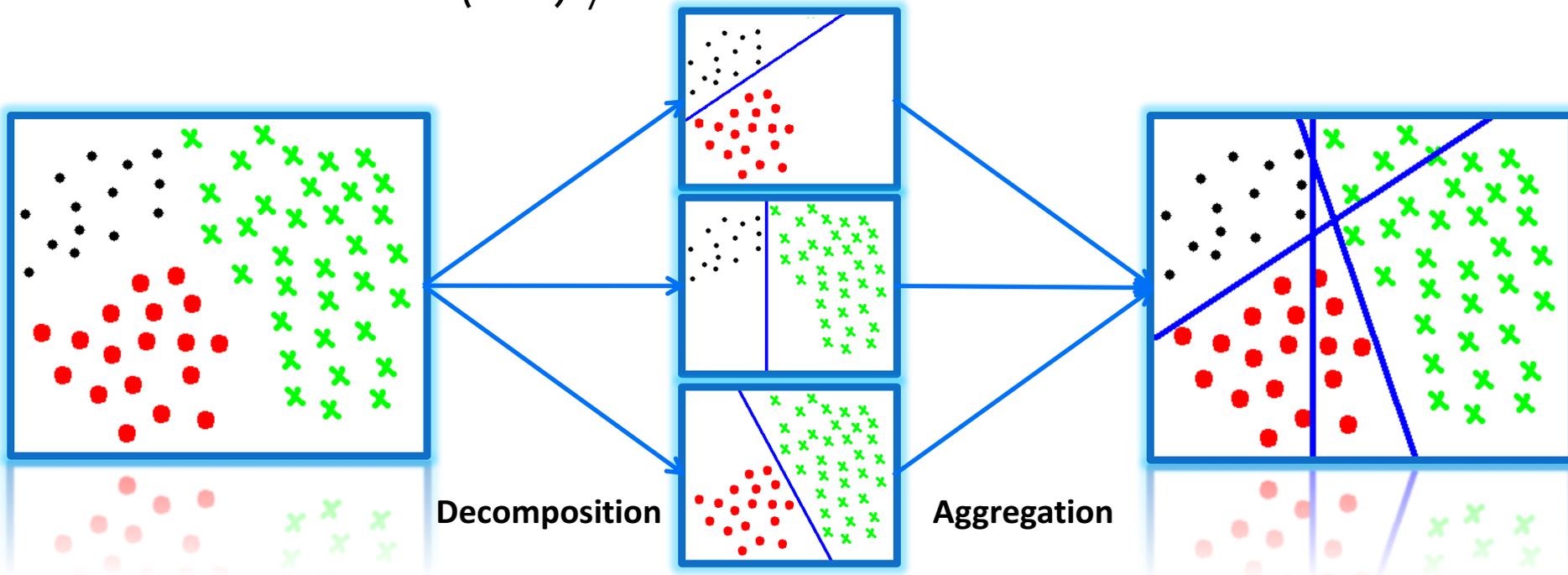
- Decomposition of the multi-class problem
  - Divide and conquer strategy
  - Multi-class → Multiple easier to solve binary problems
    - For each binary problem
      - 1 binary classifier = base classifier
    - Problem
      - How we should make the decomposition?
      - How we should aggregate the outputs?



# Decomposition Strategies

12/81

- “One-vs-One” (OVO)
  - ▣ 1 binary problem for each pair of classes
    - Pairwise Learning, Round Robin, All-vs-All...
    - *Total =  $m(m-1) / 2$  classifiers*



# One-vs-One

13/81

- Advantages
  - Smaller (number of instances)
  - Simpler decision boundaries
    - Digit recognition problem by pairwise learning
      - linearly separable [Knerr90] (first proposal)
  - Parallelizable
  - ...

[Knerr90] S. Knerr, L. Personnaz, G. Dreyfus, Single-layer learning revisited: A stepwise procedure for building and training a neural network, in: F. Fogelman Soulie, J. Hérault (eds.), *Neurocomputing: Algorithms, Architectures and Applications*, vol. F68 of NATO ASI Series, Springer-Verlag, 1990, pp. 41–50.

# One-vs-One

14/81

- Disadvantages
  - Higher testing times (more classifiers)
  - Non-competent examples [Fürnkranz06]
- Many different aggregation proposals
  - Simplest: Voting strategy
    - Each classifier votes for the predicted class
    - Predicted = class with the largest n° of votes

$$R = \begin{pmatrix} - & r_{12} & \cdots & r_{1m} \\ r_{21} & - & \cdots & r_{2m} \\ \vdots & & & \vdots \\ r_{m1} & r_{m2} & \cdots & - \end{pmatrix}$$

[Fürnkranz06] J. Fürnkranz, E. Hüllermeier, S. Vanderlooy, Binary decomposition methods for multipartite ranking, in: W. L. Buntine, M. Grobelnik, D. Mladenic, J. Shawe-Taylor (eds.), Machine Learning and Knowledge Discovery in Databases, vol. 5781(1) of LNCS, Springer, 2006, pp. 359–374.

# One-vs-One

15/81

## □ Related works

- Round Robin Ripper (R3) [Fürnkranz02]
- Fuzzy R3 (FR3) [Huhn09]
- Probability estimates by Pairwise Coupling [Wu04]
- Comparison between OVO, Boosting and Bagging
- Many aggregation proposals
  - There is not a proper comparison between them

[Fürnkranz03]

[Fürnkranz02] J. Fürnkranz, Round robin classification, *Journal of Machine Learning Research* 2 (2002) 721–747.

[Huhn09] J. C. Huhn, E. Hüllermeier, FR3: A fuzzy rule learner for inducing reliable classifiers, *IEEE Transactions on Fuzzy Systems* 17 (1) (2009) 138–149.

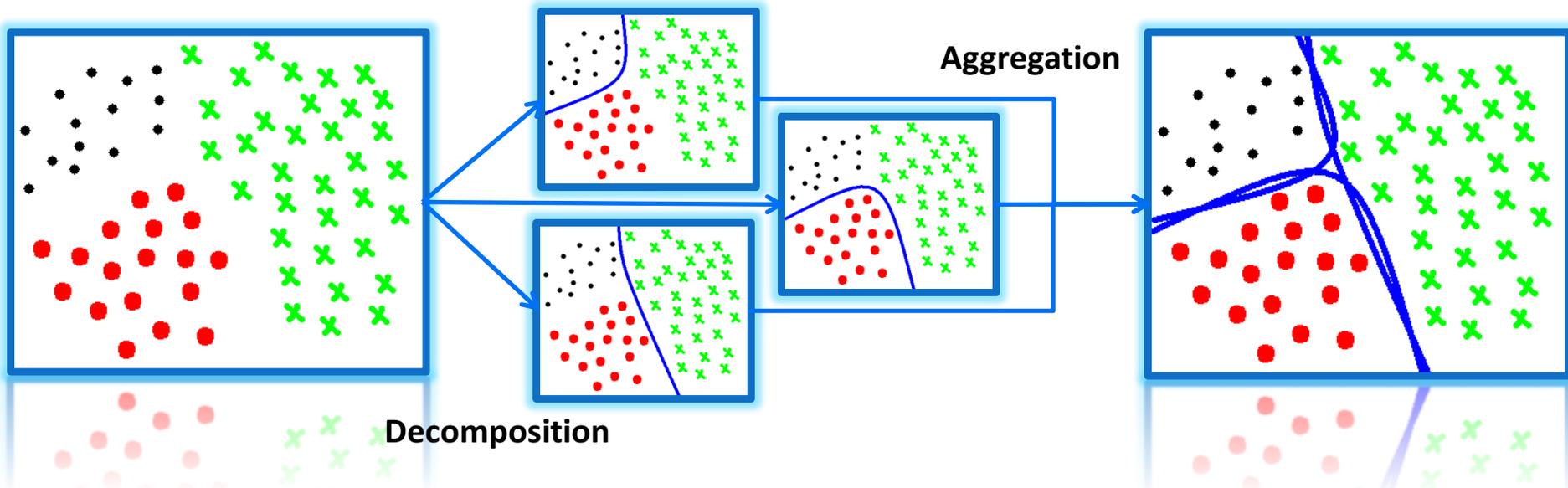
[Wu04] T. F. Wu, C. J. Lin, R. C. Weng, Probability estimates for multi-class classification by pairwise coupling, *Journal of Machine Learning Research* 5 (2004) 975–1005.

[Fürnkranz03] J. Fürnkranz, Round robin ensembles, *Intelligent Data Analysis* 7 (5) (2003) 385–403.

# Decomposition Strategies

16/81

- “One-vs-All” (OVA)
  - ▣ 1 binary problem for each class
    - All instances in each problem
      - Positive class: instances from the class considered
      - Negative class: instances from all other classes
    - *Total = m classifiers*



# One-vs-All

17/81

- Advantages
  - Less n° of classifiers
  - All examples are “competent”
- Disadvantages
  - Less studied in the literature
    - low n° of aggregations
      - Simplest: Maximum confidence rule ( $\max(r_i)$ )  
$$R = (r_1, r_2, \dots, r_i, \dots, r_m)$$
  - More complex problems
  - Imbalance training sets

# One-vs-All

18/81

- Related Works
  - Rifkin and Klatau [Rifkin04]
    - Critical with all previous literature about OVO
    - OVA classifiers are as accurate as OVO when the base classifier are fine-tuned (about SVM)
- In general
  - Previous works proved goodness of OVO
    - Ripper and C4.5, cannot be tuned

[Rifkin04] R. Rifkin, A. Klatau, In defense of one-vs-all classification, *Journal of Machine Learning Research* 5 (2004) 101–141.

# Decomposition Strategies

19/81

## □ Other approaches

### ▣ ECOC (Error Correcting Output Code) [Allwein00]

- Unify (generalize) OVO and OVA approach
- Code-Matrix representing the decomposition
  - The outputs forms a code-word
  - An ECOC is used to decode the code-word
    - The class is given by the decodification

Class	Classifier								
	C1	C2	C3	C4	C5	C5	C7	C8	C9
Class1	1	1	1	0	0	0	1	1	1
Class2	0	0	-1	1	1	0	1	-1	-1
Class3	-1	0	0	-1	0	1	-1	1	-1
Class4	0	-1	0	0	-1	-1	-1	-1	1

[Allwein00] E. L. Allwein, R. E. Schapire, Y. Singer, Reducing multiclass to binary: A unifying approach for margin classifiers, *Journal of Machine Learning Research* 1 (2000) 113–141.

# Decomposition Strategies

20/81

- Other approaches
  - Hierarchical approaches
    - Distinguish groups of classes in each nodes
  - Detailed review of decomposition strategies in [Lorena09]
    - Only an enumeration of methods
    - Low importance to the aggregation step

[Lorena09] A. C. Lorena, A. C. Carvalho, J. M. Gama, A review on the combination of binary classifiers in multiclass problems, *Artificial Intelligence Review* 30 (1-4) (2008) 19–37.

# Combination of the outputs

21/81

- Aggregation phase
  - *The way in which the outputs of the base classifiers are combined to obtain the final output.*
  - Key-factor in OVO and OVA ensembles
  - Ideally, voting and max confidence works
    - In real problems
      - Contradictions between base classifiers
      - Ties
      - Base classifiers are not 100% accurate
      - ...

# Outline

22/81

1. Introduction
2. Binarization
  - ▣ Decomposition strategies (One-vs-One, One-vs-All and Others)
  - ▣ State-of-the-art on Aggregations
    - **One-vs-One**
    - One-vs-All
3. Experimental Study
  - ▣ Experimental Framework
  - ▣ Results and Statistical Analysis
4. Discussion: Lessons Learned and Future Work
5. Conclusions for OVO vs OVA
6. Novel Approaches for the One-vs-One Learning Scheme
  - ▣ Dynamic OVO: Avoiding Non-competence
  - ▣ Distance-based Relative Competence Weighting Approach (DRCW-OVO)

# State-of-the-art on aggregation for OVO

23/81

## □ Starting from the score-matrix

$$R = \begin{pmatrix} - & r_{12} & \cdots & r_{1m} \\ r_{21} & - & \cdots & r_{2m} \\ \vdots & & & \vdots \\ r_{m1} & r_{m2} & \cdots & - \end{pmatrix}$$

- $r_{ij} =$  confidence of classifier in favor of class  $i$
- $r_{ji} =$  confidence of classifier in favor of class  $j$ 
  - Usually:  $r_{ji} = 1 - r_{ij}$  (required for probability estimates)

# State-of-the-art on aggregation for OVO

24/81

- Voting strategy (VOTE) [Friedman96]
  - Each classifier gives a vote for the predicted class
  - The class with the largest number of votes is predicted

$$Class = arg \max_{i=1, \dots, m} \sum_{1 \leq j \neq i \leq m} s_{ij}$$

- where  $s_{ij}$  is 1 if  $r_{ji} > r_{ji}$  and 0 otherwise.

- Weighted voting strategy (WV)
  - WV = VOTE but weight = confidence

$$Class = arg \max_{i=1, \dots, m} \sum_{1 \leq j \neq i \leq m} r_{ij}$$

[Friedman96] J. H. Friedman, Another approach to polychotomous classification, Tech. rep., Department of Statistics, Stanford University (1996).

# State-of-the-art on aggregation for OVO

25/81

## □ Classification by Pairwise Coupling (PC)[Hastie98]

### ■ Estimates the joint probability for all classes

#### ■ Starting from the pairwise class probabilities

$$■ r_{ij} = \text{Prob}(\text{Class}_i \mid \text{Class}_i \text{ or } \text{Class}_j)$$

#### ■ Find the best approximation $\hat{\mathbf{p}} = (\hat{p}_1, \dots, \hat{p}_m)$

#### ■ Predicts: $\text{Class} = \arg \max_{i=1, \dots, m} \hat{p}_i$

### ■ Algorithm: Minimization of Kullack-Leibler (KL) distance

$$l(\mathbf{p}) = \sum_{1 \leq j \neq i \leq m} n_{ij} r_{ij} \log \frac{r_{ij}}{\mu_{ij}} = \sum_{i < j} n_{ij} \left( r_{ij} \log \frac{r_{ij}}{\mu_{ij}} + (1 - r_{ij}) \log \frac{1 - r_{ij}}{1 - \mu_{ij}} \right)$$

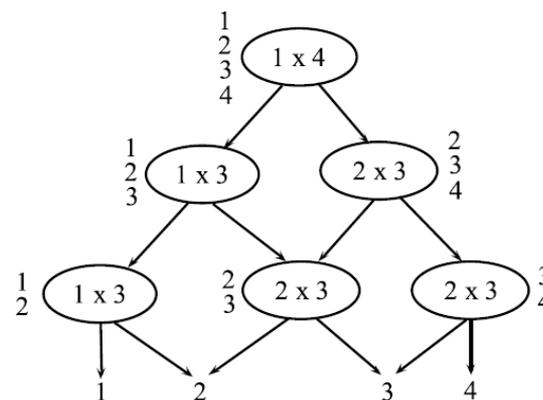
#### ■ where $\mu_{ij} = p_i / (p_i + p_j)$ , $r_{ji} = 1 - r_{ij}$ and $n_{ij}$ is the number of examples of classes $i$ and $j$

# State-of-the-art on aggregation for OVO

26/81

- Decision Directed Acyclic Graph (DDAG) [Platt00]
  - Constructs a rooted binary acyclic graph
    - Each node is associated to a list of classes and a binary classifier
    - In each level a classifier discriminates between two classes
      - The class which is not predicted is removed
    - The last class remaining on the list is the final output class.

[Platt00] J. C. Platt, N. Cristianini and J. Shawe-Taylor, Large Margin DAGs for Multiclass Classification, Proc. Neural Information Processing Systems (NIPS'99), S.A. Solla, T.K. Leen and K.-R. Müller (eds.), (2000) 547-553.



# State-of-the-art on aggregation for OVO

27/81

## □ Learning Valued Preference for Classification (LVPC)

### ■ Score-matrix = fuzzy preference relation

[Hüllermeier08,Huhn09]

### ■ Decomposition in 3 different relations

■ Strict preference

$$P_{ij} = r_{ij} - \min\{r_{ij}, r_{ji}\}$$
$$P_{ji} = r_{ji} - \min\{r_{ij}, r_{ji}\}$$

■ Conflict

$$C_{ij} = \min\{r_{ij}, r_{ji}\}$$

■ Ignorance

$$I_{ij} = 1 - \max\{r_{ij}, r_{ji}\}$$

### ■ Decision rule based on voting from the three relations

$$Class = \arg \max_{i=1, \dots, m} \sum_{1 \leq j \neq i \leq m} P_{ij} + \frac{1}{2} C_{ij} + \frac{N_i}{N_i + N_j} I_{ij}$$

- where  $N_i$  is the number of examples of class  $i$  in training

# State-of-the-art on aggregation for OVO

28/81

## □ Non-Dominance Criterion (ND) [Fernandez09]

### ■ Decision making and preference modeling [Orlovsky78]

### ■ Score-Matrix = preference relation

■  $r_{ji} = 1 - r_{ij}$ , if not  $\rightarrow$  normalize  $\bar{r}_{ij} = \frac{r_{ij}}{r_{ij} + r_{ji}}$

### ■ Compute the maximal non-dominated elements

■ Construct the strict preference relation  $r'_{ij} = \begin{cases} \bar{r}_{ij} - \bar{r}_{ji}, & \text{when } \bar{r}_{ij} > \bar{r}_{ji} \\ 0, & \text{otherwise.} \end{cases}$

■ Compute the non-dominance degree  $ND_i = 1 - \sup_{j \in C} [r'_{ji}]$

■ *the degree to which the class  $i$  is dominated by no one of the remaining classes*

■ Output  $Class = \arg \max_{i=1, \dots, m} \{ND_i\}$

[Fernandez10] A. Fernández, M. Calderón, E. Barrenechea, H. Bustince, F. Herrera, Solving multi-class problems with linguistic fuzzy rule based classification systems based on pairwise learning and preference relations, *Fuzzy Sets and System* 161:23 (2010) 3064-3080,

[Orlovsky78] S. A. Orlovsky, Decision-making with a fuzzy preference relation, *Fuzzy Sets and Systems* 1 (3) (1978) 155–167.

# State-of-the-art on aggregation for OVO

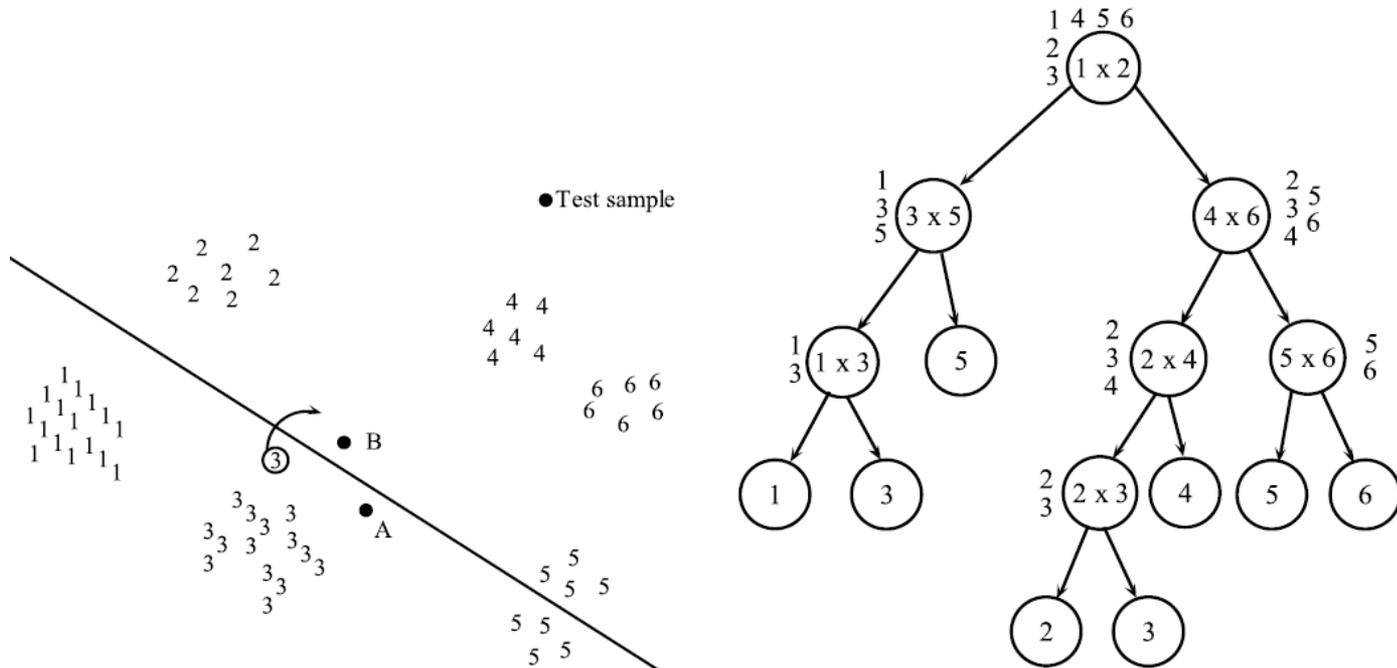
29/81

- Binary Tree of Classifiers (BTC)
  - From Binary Tree of SVM [Fei06]
  - Reduce the number of classifiers
  - Idea: Some of the binary classifiers which discriminate between two classes
    - Also can distinguish other classes at the same time
  - Tree constructed recursively
    - Similar to DDAG
      - Each node: class list + classifier
      - More than 1 class can be deleted in each node
      - To avoid false assumptions: probability threshold for examples from other classes near the decision boundary

# State-of-the-art on aggregation for OVO

30/81

- BTC for a six class problem
  - ▣ Classes 3 and 5 are assigned to two leaf nodes
    - Class 3 by reassignment (probability threshold)
    - Class 5 by the decision function between class 1 and 2



# State-of-the-art on aggregation for OVO

31/81

- Nesting One-vs-One (NEST) [Liu07,Liu08]
  - Tries to tackle the unclassifiable produced by VOTE
  - Use VOTE
    - But if there are examples within the unclassifiable region
    - Build a new OVO system only with the examples in the region in order to make them classifiable
    - Repeat until no examples remain in the unclassifiable region
  - The convergence is proved
    - No maximum nested OVOs parameter

[Liu07] Z. Liu, B. Hao and X. Yang. Nesting algorithm for multi-classification problems. *Soft Computing*, 11(4):383–389, 2007.

[Liu08] Z. Liu, B. Hao and E.C.C. Tsang. Nesting one-against-one algorithm based on SVMs for pattern classification. *IEEE Transactions on Neural Networks*, 19(12):2044–2052, 2008.

# State-of-the-art on aggregation for OVO

32/81

- Wu, Lin and Weng Probability Estimates by Pairwise Coupling approach (PE)[Wu04]
  - Obtains the posterior probabilities
    - Starting from pairwise probabilities
  - Predicts  $Class = arg \max_{i=1, \dots, m} \hat{p}_i$
  - Similar to PC
    - But solving a different optimization

$$\min_{\mathbf{p}} \sum_{i=1}^m \sum_{1 \leq j \neq i \leq m} (r_{ji}p_i - r_{ij}p_j)^2 \quad \text{subject to} \quad \sum_{i=1}^k p_i = 1, p_i \geq 0, \forall i.$$

# Outline

33/81

1. Introduction
2. Binarization
  - ▣ Decomposition strategies (One-vs-One, One-vs-All and Others)
  - ▣ State-of-the-art on Aggregations
    - One-vs-One
    - **One-vs-All**
3. Experimental Study
  - ▣ Experimental Framework
  - ▣ Results and Statistical Analysis
4. Discussion: Lessons Learned and Future Work
5. Conclusions for OVO vs OVA
6. Novel Approaches for the One-vs-One Learning Scheme
  - ▣ Dynamic OVO: Avoiding Non-competence
  - ▣ Distance-based Relative Competence Weighting Approach (DRCW-OVO)

# State-of-the-art on aggregation for OVA

34/81

- Starting from the score-vector

$$R = (r_1, r_2, \dots, r_i, \dots, r_m)$$

- $r_i$  = confidence of classifier in favor of class  $i$ 
  - Respect to all other classes
- Usually more than 1 classifier predicts the positive class
  - Tie-breaking techniques

# State-of-the-art on aggregation for OVA

35/81

- Maximum confidence strategy (MAX)
  - Predicts the class with the largest confidence
$$\text{Class} = \arg \max_{i=1, \dots, m} r_i$$
- Dynamically Ordered One-vs-All (DOO) [Hong08]
  - It is not based on confidences
  - Train a Naïve Bayes classifier
    - Use its predictions to Dynamically execute each OVA
      - Predict the first class giving a positive answer
  - Ties avoided a priori by a Naïve Bayes classifier

[Hong08] J.-H. Hong, J.-K. Min, U.-K. Cho, and S.-B. Cho. Fingerprint classification using one-vs-all support vector machines dynamically ordered with naïve bayes classifiers. *Pattern Recognition*, 41(2):662–671, 2008.

# Binarization strategies

36/81

- But...
  - ▣ Should we do binarization?
    - When it is not needed? (Ripper, C4.5, kNN...)
      - There exist previous works showing their goodness [Fürnkranz02,Fürnkranz03,Rifkin04]
  - ▣ Given that we want or have to use binarization...
    - How we should do it?
      - Some comparisons between OVO and OVA
        - Only for SVM [Hsu02]
      - No comparison for aggregation strategies

# Outline

37/81

1. Introduction
2. Binarization
  - ▣ Decomposition strategies (One-vs-One, One-vs-All and Others)
  - ▣ State-of-the-art on Aggregations
    - One-vs-One
    - One-vs-All
3. Experimental Study
  - ▣ Experimental Framework
  - ▣ Results and Statistical Analysis
4. Discussion: Lessons Learned and Future Work
5. Conclusions for OVO vs OVA
6. Novel Approaches for the One-vs-One Learning Scheme
  - ▣ Dynamic OVO: Avoiding Non-competence
  - ▣ Distance-based Relative Competence Weighting Approach (DRCW-OVO)

M. Galar, A. Fernández, E. Barrenechea, H. Bustince, F. Herrera, **An Overview of Ensemble Methods for Binary Classifiers in Multi-class Problems: Experimental Study on One-vs-One and One-vs-All Schemes**. *Pattern Recognition* 44:8 (2011) 1761-1776, [doi: 10.1016/j.patcog.2011.01.017](https://doi.org/10.1016/j.patcog.2011.01.017)

# Experimental Framework

38/81

- Different base learners
  - ▣ Support Vector Machines (SVM)
  - ▣ C4.5 Decision Tree
  - ▣ Ripper Decision List
  - ▣ k-Nearest Neighbors (kNN)
  - ▣ Positive Definite Fuzzy Classifier (PDFC)

# Experimental Framework

39/81

## □ Performance measures

### ▣ Accuracy rate

- Can be confusing evaluating multi-class problems

### ▣ Cohen's kappa

- Takes into account random hits due to number of instances

Correct Class	Predicted Class				Total
	$C_1$	$C_2$	...	$C_m$	
$C_1$	$h_{11}$	$h_{12}$	...	$h_{1m}$	$T_{r1}$
$C_2$	$h_{21}$	$h_{22}$	...	$h_{2m}$	$T_{r1}$
⋮			⋮		⋮
$C_m$	$h_{m1}$	$h_{m2}$	...	$h_{mm}$	$T_{rm}$
Total	$T_{c1}$	$T_{c2}$	...	$T_{cm}$	$T_{r1}$

$$\text{kappa} = \frac{n \sum_{i=1}^m h_{ii} - \sum_{i=1}^m T_{ri} T_{ci}}{n^2 - \sum_{i=1}^m T_{ri} T_{ci}}$$

# Experimental Framework

40/81

- 19 real-world Data-sets
- 5 fold-cross validation

Data-set	#Ex.	#Atts.	#Num.	#Nom.	#Cl.
Car	1728	6	6	0	4
Lymphography	148	18	3	15	4
Vehicle	846	18	18	0	4
Cleveland	297	13	13	0	5
Nursery	1296	8	0	8	5
Page-blocks	548	10	10	0	5
Autos	159	25	15	10	6
Dermatology	366	33	1	32	6
Flare	1389	10	0	10	6
Glass	214	9	9	0	6
Satimage	643	36	36	0	7
Segment	2310	19	19	0	7
Shuttle	2175	9	9	0	7
Zoo	101	16	0	16	7
Ecoli	336	7	7	0	8
Led7digit	500	7	0	7	10
Penbased	1099	16	16	0	10
Yeast	1484	8	8	0	10
Vowel	990	13	13	0	11

# Experimental Framework

41/81

- Algorithms parameters
  - ▣ Default configuration

<i>Algorithm</i>	<i>Parameters</i>
SVM	C = 1.0 Tolerance Parameter = 0.001 Epsilon = 1.0E-12 Kernel Type = Polynomial Polynomial Degree = 1 Fit Logistic Models = True
C4.5	Prune = True Confidence level = 0.25 Minimum number of item-sets per leaf = 2
1NN	$k = 1$ Distance metric = Heterogeneous Value Difference Metric (HVDM)
3NN	$k = 3$ Distance metric = Heterogeneous Value Difference Metric (HVDM)
Ripper	Size of growing subset = 66% Repetitions of the optimization stage = 2
PDFC	C = 100.0 Tolerance Parameter = 0.001 Epsilon = 1.0E-12 Kernel Type = Polynomial Polynomial Degree = 1 PDRF Type = Gaussian

# Experimental Framework

42/81

- Confidence estimations
  - SVM: Logistic model
    - SVM for probability estimates
  - C4.5: Purity of the predictor leaf
    - N° of instances correctly classified by the leaf / Total n° of instances in the leaf
  - kNN:  $Confidence = \frac{\sum_{l=1}^k \frac{e_l}{d_l}}{\sum_{l=1}^k \frac{1}{d_l}}$ 
    - where  $d_l$  = distance between the input pattern and the  $l^{th}$  neighbor
    - $e_l = 1$  if the neighbor  $l$  is from the class and 0 otherwise
  - Ripper: Purity of the rule
    - N° of instances correctly classified by the rule / Total n° of instances in the rule
  - PDFC: confidence = 1 is given for the predicted class

# Outline

43/81

1. Introduction
2. Binarization
  - ▣ Decomposition strategies (One-vs-One, One-vs-All and Others)
  - ▣ State-of-the-art on Aggregations
    - One-vs-One
    - One-vs-All
3. Experimental Study
  - ▣ Experimental Framework
  - ▣ Results and Statistical Analysis
4. Discussion: Lessons Learned and Future Work
5. Conclusions for OVO vs OVA
6. Novel Approaches for the One-vs-One Learning Scheme
  - ▣ Dynamic OVO: Avoiding Non-competence
  - ▣ Distance-based Relative Competence Weighting Approach (DRCW-OVO)
  - ▣ Difficult Classes Problem in OVO Strategy

M. Galar, A. Fernández, E. Barrenechea, H. Bustince, F. Herrera, **An Overview of Ensemble Methods for Binary Classifiers in Multi-class Problems: Experimental Study on One-vs-One and One-vs-All Schemes.** *Pattern Recognition* 44:8 (2011) 1761-1776, [doi: 10.1016/j.patcog.2011.01.017](https://doi.org/10.1016/j.patcog.2011.01.017)

# Experimental Study

44/81

## □ Average accuracy and kappa results

Method	Aggregation	SVM		C4.5		1NN	
		Acc <sub>tst</sub>	Avg. Rank	Acc <sub>tst</sub>	Avg. Rank	Acc <sub>tst</sub>	Avg. Rank
Base	-	-	-	80.51 ± 3.85	-	81.24 ± 2.98	-
OVO	VOTE	81.14 ± 3.22	<b>4.37 (1)</b>	81.57 ± 3.29	4.63 (4)	82.06 ± 3.38	3.82 (3)
	WV	81.05 ± 2.92	5.08 (6)	<b>81.59 ± 3.28</b>	3.97 (2)	-	-
	DDAG	81.01 ± 3.28	5.39 (8)	81.02 ± 3.56	6.21 (9)	81.86 ± 3.31	4.32 (5)
	PC	81.08 ± 2.89	5.29 (7)	81.49 ± 3.32	4.34 (3)	82.26 ± 3.33	3.21 (2)
	LVPC	81.14 ± 3.11	4.50 (3)	81.57 ± 3.28	<b>3.87 (1)</b>	-	-
	ND	81.01 ± 3.15	4.92 (5)	81.12 ± 3.24	5.58 (6)	81.48 ± 3.51	4.97 (7)
	BTC	80.82 ± 3.24	6.18 (9)	81.22 ± 2.87	5.61 (7)	82.21 ± 3.12	3.89 (4)
	NEST	<b>81.14 ± 3.32</b>	4.47 (2)	81.20 ± 3.47	5.74 (8)	81.68 ± 3.47	4.68 (6)
	PE	81.03 ± 3.35	4.79 (4)	81.42 ± 3.22	5.05 (5)	<b>82.30 ± 3.11</b>	<b>3.11 (1)</b>
OVA	MAX	78.66 ± 3.00	1.53 (2)	78.01 ± 4.19	1.84 (2)	81.18 ± 4.51	1.63 (2)
	DOO	<b>78.75 ± 3.15</b>	<b>1.47 (1)</b>	<b>78.78 ± 4.36</b>	<b>1.16 (1)</b>	<b>81.77 ± 4.45</b>	<b>1.37 (1)</b>

Method	Aggregation	SVM		C4.5		1NN	
		Kappa <sub>tst</sub>	Avg. Rank	Kappa <sub>tst</sub>	Avg. Rank	Kappa <sub>tst</sub>	Avg. Rank
Base	-	-	-	.7203 ± .0554	-	.7369 ± .0475	-
OVO	VOTE	.7233 ± .0548	4.82 (2)	.7331 ± .0490	5.16 (5)	.7419 ± .0535	3.84 (3)
	WV	.7229 ± .0506	5.05 (6)	<b>.7348 ± .0485</b>	<b>3.76 (1)</b>	-	-
	DDAG	.7230 ± .0555	5.11 (7)	.7304 ± .0535	5.92 (8)	.7402 ± .0522	3.89 (4)
	PC	.7234 ± .0520	5.18 (8)	.7341 ± .0493	4.13 (3)	.7449 ± .0525	3.00 (2)
	LVPC	.7211 ± .0531	5.03 (5)	.7341 ± .0488	4.03 (2)	-	-
	ND	.7225 ± .0533	4.82 (2)	.7286 ± .0489	5.53 (7)	.7340 ± .0556	5.37 (7)
	BTC	.7204 ± .0551	6.05 (9)	.7297 ± .0428	5.42 (6)	.7438 ± .0498	4.29 (5)
	NEST	<b>.7243 ± .0559</b>	<b>4.03 (1)</b>	.7291 ± .0514	6.34 (9)	.7366 ± .0547	4.79 (6)
	PE	.7228 ± .0537	4.92 (4)	.7330 ± .0480	4.71 (4)	<b>.7453 ± .0497</b>	<b>2.82 (1)</b>
OVA	MAX	.6868 ± .0553	1.55 (2)	.6826 ± .0629	1.89 (2)	.7298 ± .0705	1.63 (2)
	DOO	<b>.6868 ± .0565</b>	<b>1.45 (1)</b>	<b>.6938 ± .0649</b>	<b>1.11 (1)</b>	<b>.7368 ± .0701</b>	<b>1.37 (1)</b>

# Experimental Study

45/81

## □ Average accuracy and kappa results

Method	Aggregation	3NN		Ripper		PDFC	
		Acc <sub>tst</sub>	Avg. Rank	Acc <sub>tst</sub>	Avg. Rank	Acc <sub>tst</sub>	Avg. Rank
Base	-	81.54 ± 2.65	-	76.52 ± 4.00	-	-	-
OVO	VOTE	83.00 ± 2.92	5.05 (6)	80.57 ± 3.17	3.89 (3)	<b>84.33 ± 3.10</b>	3.37 (2)
	WV	<b>83.11 ± 2.87</b>	4.47 (3)	80.54 ± 3.03	3.87 (2)	-	-
	DDAG	82.73 ± 2.83	5.87 (8)	77.62 ± 3.61	7.08 (9)	84.05 ± 3.00	3.71 (3)
	PC	83.00 ± 2.96	5.11 (7)	80.33 ± 3.30	4.87 (5)	84.12 ± 3.05	<b>3.29 (1)</b>
	LVPC	83.07 ± 2.79	4.61 (4)	<b>80.58 ± 3.16</b>	<b>3.68 (1)</b>	-	-
	ND	83.07 ± 2.93	<b>4.29 (1)</b>	79.38 ± 3.27	5.29 (7)	84.05 ± 2.96	4.68 (6)
	BTC	82.99 ± 2.98	5.00 (5)	79.19 ± 3.07	6.39 (8)	84.24 ± 3.01	4.29 (5)
	NEST	82.67 ± 2.94	6.16 (9)	80.01 ± 3.50	5.08 (6)	83.88 ± 3.02	4.89 (7)
OVA	PE	83.11 ± 2.94	4.45 (2)	80.07 ± 3.08	4.84 (4)	84.06 ± 3.04	3.76 (4)
	MAX	82.75 ± 4.29	1.58 (2)	78.30 ± 4.94	1.71 (2)	<b>83.59 ± 3.12</b>	<b>1.39 (1)</b>
	DOO	<b>82.76 ± 4.38</b>	<b>1.42 (1)</b>	<b>79.12 ± 4.67</b>	<b>1.29 (1)</b>	83.01 ± 3.10	1.61 (2)

Method	Aggregation	3NN		Ripper		PDFC	
		Kappa <sub>tst</sub>	Avg. Rank	Kappa <sub>tst</sub>	Avg. Rank	Kappa <sub>tst</sub>	Avg. Rank
Base	-	.7335 ± .0452	-	.6799 ± .0554	-	-	-
OVO	VOTE	.7507 ± .0500	5.03 (6)	<b>.7250 ± .0475</b>	4.26 (3)	<b>.7677 ± .0538</b>	3.63 (2)
	WV	.7519 ± .0487	4.71 (3)	.7249 ± .0455	<b>3.68 (1)</b>	-	-
	DDAG	.7479 ± .0487	5.87 (8)	.6957 ± .0489	6.42 (8)	.7659 ± .0518	3.97 (5)
	PC	.7505 ± .0505	4.89 (5)	.7227 ± .0483	4.61 (5)	.7670 ± .0529	<b>3.11 (1)</b>
	LVPC	.7496 ± .0475	5.18 (7)	.7246 ± .0469	4.00 (2)	-	-
	ND	.7524 ± .0500	<b>4.03 (1)</b>	.7098 ± .0479	5.92 (7)	.7625 ± .0524	5.32 (7)
	BTC	.7519 ± .0514	4.87 (4)	.7087 ± .0476	6.58 (9)	.7668 ± .0527	3.79 (3)
	NEST	.7461 ± .0505	6.24 (9)	.7195 ± .0496	4.97 (6)	.7641 ± .0514	4.37 (6)
OVA	PE	<b>.7526 ± .0499</b>	4.18 (2)	.7193 ± .0457	4.55 (4)	.7653 ± .0524	3.82 (4)
	MAX	<b>.7481 ± .0695</b>	1.58 (2)	.6896 ± .0743	1.79 (2)	<b>.7556 ± .0589</b>	<b>1.37 (1)</b>
	DOO	.7473 ± .0710	<b>1.42 (1)</b>	<b>.7004 ± .0716</b>	<b>1.21 (1)</b>	.7478 ± .0587	1.63 (2)

# Which is the most appropriate aggregation?

46/81

- OVO aggregations Analysis
  - SVM: NEST and VOTE, but no statistical differences
  - C4.5: Statistical differences
    - WV, LVPC and PC the most robust
    - NEST and DDAG the weakest
  - 1NN: Statistical differences
    - PC and PE the best → confidences in  $\{0,1\}$ 
      - In PDFC they also excel
    - ND the worst → poor confidences, excessive ties

# Which is the most appropriate aggregation?

47/81

- OVO aggregations Analysis
  - ▣ 3NN: No significant differences
    - ND stands out
  - ▣ Ripper: Statistical differences
    - WV and LVPC vs. BTC and DDAG
  - ▣ PDFC: No significant differences (low p-value in kappa)
    - VOTE, PC and PE overall better performance

# Which is the most appropriate aggregation?

48/81

- OVA aggregations Analysis
  - DOO performs better when the base classifiers accuracy is not better than the Naïve Bayes ones.
  - It helps selecting the most appropriate classifier to use dynamically
  - In other cases, it can distort the results

Base Classifier	Measure	$R^+$	$R^-$	Hypothesis ( $\alpha = 0.05$ )	p-value
SVM	Accuracy	82	108	Not Rejected	0.53213
	Kappa	86	104	Not Rejected	0.75637
C4.5	Accuracy	14	176	Rejected for DOO	0.00179
	Kappa	11.5	178.5	Rejected for DOO	0.00118
1NN	Accuracy	55	135	Not Rejected	0.09097
	Kappa	55	135	Not Rejected	0.09097
3NN	Accuracy	75	115	Not Rejected	0.73532
	Kappa	76	114	Not Rejected	0.86577
Ripper	Accuracy	44.5	145.5	Rejected for DOO	0.04286
	Kappa	42	148	Rejected for DOO	0.03294
PDFC	Accuracy	130.5	59.5	Rejected for MAX	0.0464
	Kappa	138	52	Rejected for MAX	0.02799

# Should we do binarization?

## How should we do it?

49/81

### □ Representatives of OVO and OVA

#### ■ By the previous analysis

	SVM	C4.5	1NN	3NN	Ripper	PDFC
OVO	NEST <sub>ovo</sub>	WV <sub>ovo</sub>	PE <sub>ovo</sub>	ND <sub>ovo</sub>	WV <sub>ovo</sub>	PC <sub>ovo</sub>
OVA	DOO <sub>ova</sub>	DOO <sub>ova</sub>	DOO <sub>ova</sub>	DOO <sub>ova</sub>	DOO <sub>ova</sub>	MAX <sub>ova</sub>

#### ■ Average results

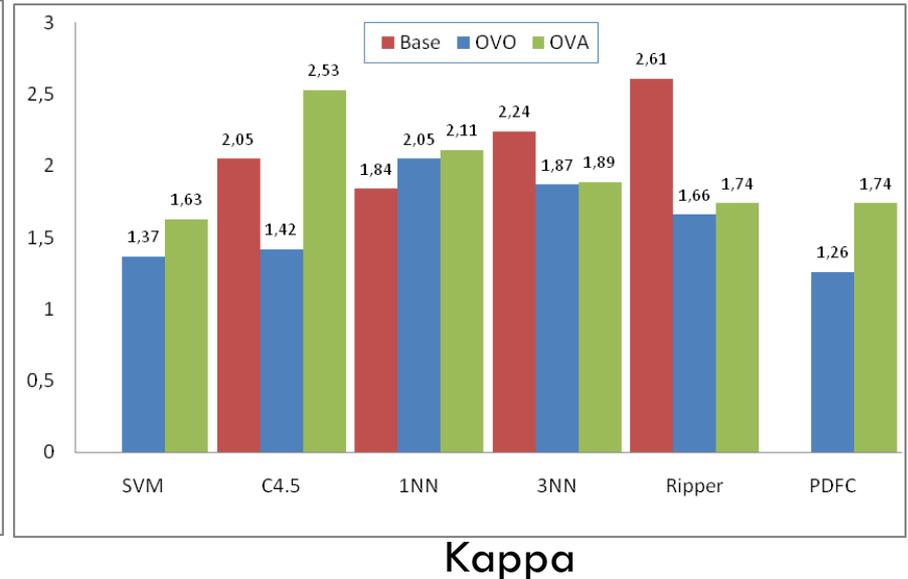
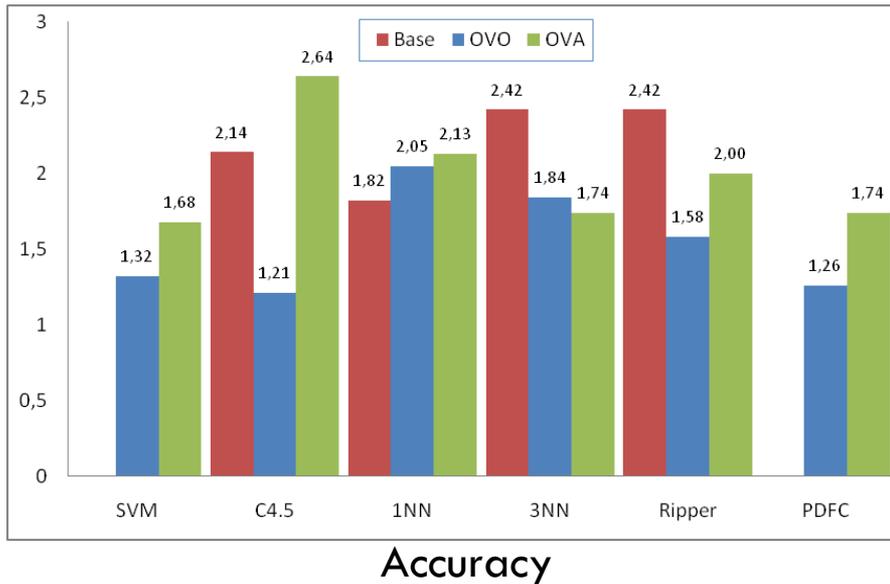
Base Classifier	Aggregation	Accuracy		Kappa	
		Test	Avg. Rank	Test	Avg. Rank
SVM	NEST <sub>ovo</sub>	<b>81.14 ± 3.32</b>	<b>1.37 (1)</b>	<b>.7243 ± .0559</b>	<b>1.32 (1)</b>
	DOO <sub>ova</sub>	78.75 ± 3.15	1.63 (2)	.6868 ± .0565	1.68 (2)
C4.5	C45	80.51 ± 3.85	2.05 (2)	.7203 ± .0554	2.14 (2)
	WV <sub>ovo</sub>	<b>81.59 ± 3.28</b>	<b>1.42 (1)</b>	<b>.7348 ± .0485</b>	<b>1.21 (1)</b>
	DOO <sub>ova</sub>	78.78 ± 4.36	2.53 (3)	.6938 ± .0649	2.64 (3)
1NN	1NN	81.24 ± 2.98	<b>1.84 (1)</b>	.7369 ± .0475	<b>1.82 (1)</b>
	PE <sub>ovo</sub>	<b>82.30 ± 3.11</b>	2.05 (2)	<b>.7453 ± .0497</b>	2.05 (2)
	DOO <sub>ova</sub>	81.77 ± 4.45	2.11 (3)	.7368 ± .0701	2.13 (3)
3NN	3NN	81.54 ± 2.65	2.24 (3)	.7335 ± .0452	2.42 (3)
	ND <sub>ovo</sub>	<b>83.07 ± 2.93</b>	<b>1.87 (1)</b>	<b>.7524 ± .0500</b>	1.84 (2)
	DOO <sub>ova</sub>	82.76 ± 4.38	1.89 (2)	.7473 ± .0710	<b>1.74 (1)</b>
Ripper	Ripper	76.52 ± 4.00	2.61 (3)	.6799 ± .0554	2.42 (3)
	WV <sub>ovo</sub>	<b>80.54 ± 3.03</b>	<b>1.66 (1)</b>	<b>.7249 ± .0455</b>	<b>1.58 (1)</b>
	DOO <sub>ova</sub>	79.12 ± 4.67	1.74 (2)	.7004 ± .0716	2.00 (2)
PDFC	PC <sub>ovo</sub>	<b>84.12 ± 3.05</b>	<b>1.26 (1)</b>	<b>.7670 ± .0529</b>	<b>1.26 (1)</b>
	MAX <sub>ova</sub>	83.59 ± 3.12	1.74 (2)	.7556 ± .0589	1.74 (2)

# Should we do binarization?

## How should we do it?

50/81

- Rankings within each classifier
  - ▣ In general, OVO is the most competitive

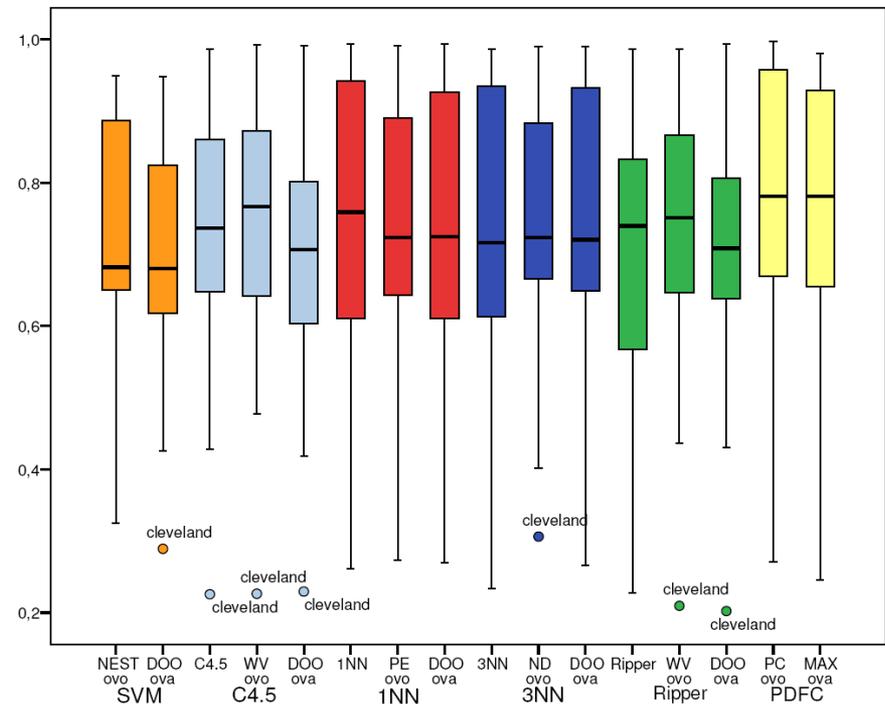
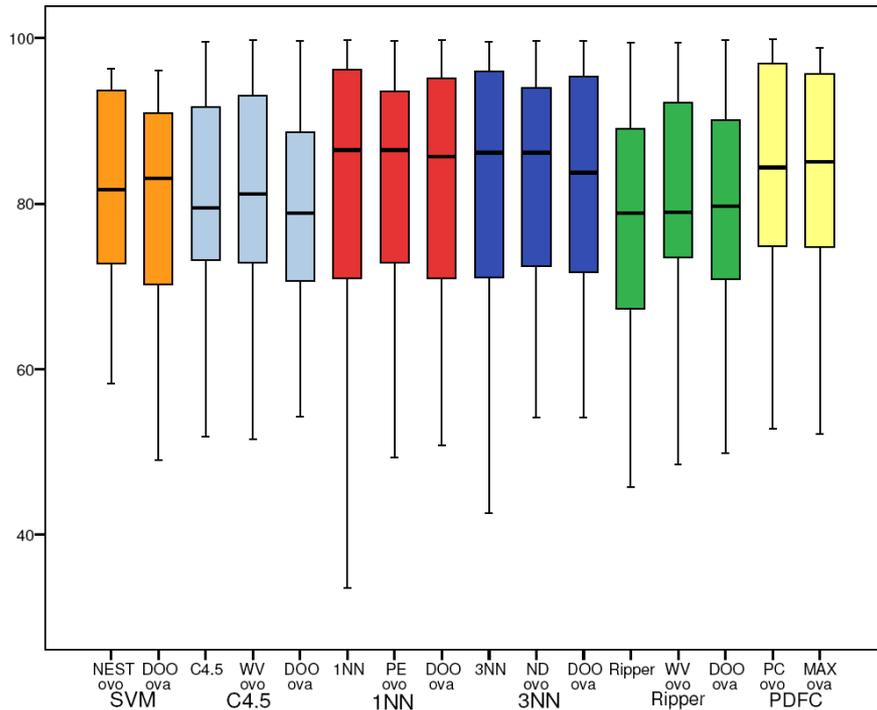


# Should we do binarization?

## How should we do it?

51/81

- Box plots for test results
  - ▣ OVA reduce performance in kappa
  - ▣ OVO is more compact (hence, robust)



# Should we do binarization?

## How should we do it?

52/81

### □ Statistical analysis

#### ■ SVM and PDFC

- OVO outperforms OVA with significant differences

Base Classifier	Comparison	Measure	$R^+$	$R^-$	Hypothesis ( $\alpha = 0.05$ )	p-value
SVM	NEST <sub>ovo</sub> vs. DOO <sub>ova</sub>	Accuracy	153	37	Rejected for NEST <sub>ovo</sub>	0.01959
		Kappa	156	34	Rejected for NEST <sub>ovo</sub>	0.0141
PDFC	PC <sub>ovo</sub> vs. MAX <sub>ova</sub>	Accuracy	146	44	Rejected for PC <sub>ovo</sub>	0.04014
		Kappa	147	43	Rejected for PC <sub>ovo</sub>	0.03639

#### ■ C4.5, 1NN, 3NN and Ripper

- P-values returned by Iman-Davenport tests (\* if rejected)

	C4.5	1NN	3NN	Ripper
Accuracy	0.00134 *	0.70296	0.45982	0.00296 *
Kappa	0.00026 *	0.61089	0.07585	0.02982 *

- Post-hoc test for C4.5 and Ripper

- kNN, no statistical differences, but also not worse results

# Should we do binarization?

## How should we do it?

53/81

### □ Statistical analysis

#### ■ C4.5

##### ■ WV for OVO outperforms the rest

(a) Accuracy			(b) Kappa		
i	Hypothesis	p-value	i	Hypothesis	p-value
1	WV <sub>ovo</sub> vs. DOO <sub>ova</sub>	+(0.00197)	1	WV <sub>ovo</sub> vs. DOO <sub>ova</sub>	+(0.00057)
2	C4.5 vs. WV <sub>ovo</sub>	=(0.05158)	2	C4.5 vs. WV <sub>ovo</sub>	+(0.03496)
3	C4.5 vs. DOO <sub>ova</sub>	=(0.14429)	3	C4.5 vs. DOO <sub>ova</sub>	=(0.10476)

#### ■ Ripper

##### ■ WV for OVO is the best

##### ■ No statistical differences with OVA

##### ■ But OVO differs statistically from Ripper while OVA do not

(a) Accuracy			(b) Kappa		
i	Hypothesis	p-value	i	Hypothesis	p-value
1	Ripper vs. WV <sub>ovo</sub>	+(0.01050)	1	Ripper vs. WV <sub>ovo</sub>	+(0.02833)
2	Ripper vs. DOO <sub>ova</sub>	+(0.01050)	2	Ripper vs. DOO <sub>ova</sub>	=(0.19437)
3	WV <sub>ovo</sub> vs. DOO <sub>ova</sub>	=(0.80775)	3	WV <sub>ovo</sub> vs. DOO <sub>ova</sub>	=(0.19437)

# Outline

54/81

1. Introduction
2. Binarization
  - ▣ Decomposition strategies (One-vs-One, One-vs-All and Others)
  - ▣ State-of-the-art on Aggregations
    - One-vs-One
    - One-vs-All
3. Experimental Study
  - ▣ Experimental Framework
  - ▣ Results and Statistical Analysis
4. Discussion: Lessons Learned and Future Work
5. Conclusions for OVO vs OVA
6. Novel Approaches for the One-vs-One Learning Scheme
  - ▣ Dynamic OVO: Avoiding Non-competence
  - ▣ Distance-based Relative Competence Weighting Approach (DRCW-OVO)

M. Galar, A. Fernández, E. Barrenechea, H. Bustince, F. Herrera, **An Overview of Ensemble Methods for Binary Classifiers in Multi-class Problems: Experimental Study on One-vs-One and One-vs-All Schemes**. *Pattern Recognition* 44:8 (2011) 1761-1776, [doi: 10.1016/j.patcog.2011.01.017](https://doi.org/10.1016/j.patcog.2011.01.017)

# Discussion

55/81

- Lessons learned
  - Binarization is beneficial
    - Also when the problem can be tackled without it
  - The most robust aggregations for OVO
    - WV, LVPC, PC and PE
  - The most robust aggregations for OVA
    - Not clear
    - Need more attention, can be improved
  - Too many approaches to deal with the unclassifiable region in OVO (NEST, BTC, DDAG)

# Discussion

56/81

- Lessons learned
  - ▣ OVA problem
    - Imbalanced data-sets
    - Not against Rifkin's findings
      - But, this means that OVA are less robust
        - Need more fine-tuned base classifiers
  - ▣ Importance of confidence estimates of base classifiers
  - ▣ Scalability
    - Number of classes: OVO seems to work better
    - Number of instances: OVO natures make it more adequate

# Discussion

57/81

- Future work
  - Detection of non-competent examples
  - Techniques for imbalanced data-sets
  - Studies on scalability
  - OVO as a decision making problem
    - Suppose inaccurate or erroneous base classifiers
  - New combinations for OVA
    - Something more than a tie-breaking technique
  - Data-complexity measures
    - A priori knowledge extraction to select the proper mechanism

# Outline

58/81

1. Introduction
2. Binarization
  - ▣ Decomposition strategies (One-vs-One, One-vs-All and Others)
  - ▣ State-of-the-art on Aggregations
    - One-vs-One
    - One-vs-All
3. Experimental Study
  - ▣ Experimental Framework
  - ▣ Results and Statistical Analysis
4. Discussion: Lessons Learned and Future Work
5. Conclusions for OVO vs OVA
6. Novel Approaches for the One-vs-One Learning Scheme
  - ▣ Dynamic OVO: Avoiding Non-competence
  - ▣ Distance-based Relative Competence Weighting Approach (DRCW-OVO)

M. Galar, A.Fernández, E. Barrenechea, H. Bustince, F. Herrera, **An Overview of Ensemble Methods for Binary Classifiers in Multi-class Problems: Experimental Study on One-vs-One and One-vs-All Schemes.** *Pattern Recognition* 44:8 (2011) 1761-1776, [doi: 10.1016/j.patcog.2011.01.017](https://doi.org/10.1016/j.patcog.2011.01.017)

# Conclusions

59/81

- Goodness of using binarization
  - ▣ Concretely, OVO approach
    - WV, LVPC, PC and PE
    - The aggregation is base learner dependant
- Low attention to OVA strategy
  - ▣ Problems with imbalanced data
- Importance of confidence estimates
- Many work remind to be addressed

# Outline

60/81

1. Introduction
2. Binarization
  - ▣ Decomposition strategies (One-vs-One, One-vs-All and Others)
  - ▣ State-of-the-art on Aggregations
    - One-vs-One
    - One-vs-All
3. Experimental Study
  - ▣ Experimental Framework
  - ▣ Results and Statistical Analysis
4. Discussion: Lessons Learned and Future Work
5. Conclusions for OVO vs OVA
6. Novel Approaches for the One-vs-One Learning Scheme
  - ▣ Dynamic OVO: Avoiding Non-competence
  - ▣ Distance-based Relative Competence Weighting Approach (DRCW-OVO)

M. Galar, A. Fernández, E. Barrenechea, H. Bustince, F. Herrera, **Dynamic Classifier Selection for One-vs-One Strategy: Avoiding Non-Competent Classifiers**. *Pattern Recognition* 46:12 (2013) 3412–3424, [doi: j.patcog.2013.04.018](https://doi.org/10.1016/j.patcog.2013.04.018)

# Dynamic OVO: Avoiding Non-competence

61/81

## □ *Non-Competent Classifiers:*

- Those whose output is not relevant for the classification of the query instance
- They have not been trained with instances of the real class of the example to be classified
  - Classify  $x$ , whose real class is  $c_1$

$$\square R(x) = \begin{pmatrix} & c1 & c2 & c3 & c4 & c5 \\ c1 & - & 0,55 & 0,6 & 0,75 & 0,7 \\ c2 & 0,45 & - & 0,4 & 1 & 0,8 \\ c3 & 0,4 & 0,6 & - & 0,5 & 0,4 \\ c4 & 0,25 & 0,0 & 0,5 & - & 0,1 \\ c5 & 0,30 & 0,2 & 0,6 & 0,9 & - \end{pmatrix}$$

# Dynamic OVO: Avoiding Non-competence

62/81

## □ *Non-Competent Classifiers:*

- Consider WV aggregation,  $c_2$  is predicted
- **None** of the classifiers **considering  $c_1$  failed**
- **Non-competent** classifiers **strongly voted for  $c_2$**

$$\square R(x) = \left( \begin{array}{cccccc|c} & c1 & c2 & c3 & c4 & c5 & WV \\ c1 & - & 0,55 & 0,6 & 0,75 & 0,7 & 2,6 \\ c2 & 0,45 & - & 0,4 & 1 & 0,8 & \mathbf{2,65} \\ c3 & 0,4 & 0,6 & - & 0,5 & 0,4 & 1,9 \\ c4 & 0,25 & 0,0 & 0,5 & - & 0,1 & 0,85 \\ c5 & 0,30 & 0,2 & 0,6 & 0,9 & - & 2,1 \end{array} \right)$$

# Dynamic OVO: Avoiding Non-competence

63/81

- Dynamic Classifier Selection:
  - Classifiers specialized in different areas of the **input** space
  - Classifiers **complement** themselves
  - **The most competent one** for the instance is selected:
    - Instead of combining them all
    - Assuming that several misses can be done (they are corrected)

# Dynamic OVO: Avoiding Non-competence

64/81

- **Avoiding non-competence problem**
- Adapting Dynamic Classifier Selection (DCS) to OVO
  - ▣ Baseline classifiers competent over their pair of classes
- **Search for a lower set of classes** than those that probably the instance belongs to.
  - ▣ Remove those (probably) non-competent classifiers
  - ▣ Avoid misclassifications
- **Neighbourhood of the instance** is considered<sup>[Woods97]</sup>
  - ▣ Local precisions cannot be estimated
  - ▣ **Classes in the neighbourhood** → reduced score matrix

[WOODS97] K. Woods, W. Philip Kegelmeyer, K. Bowyer. Combination of multiple classifiers using local accuracy estimates, IEEE Transactions on Pattern Analysis and Machine Intelligence 19(4):405-410, 1997.

# Dynamic OVO: Avoiding Non-competence

65/81

- DCS ALGORITHM FOR OVO STRATEGY
  1. Compute the  $k$  nearest neighbors of the instance ( $k = 3 \cdot m$ )
  2. Select the classes in the neighborhood (if it is unique  $k++$ )
  3. Consider the subset of classes in the reduced-score matrix
- **Any existing OVO aggregation can be used**
- Difficult to misclassify instances
- $k$  value is larger than the usual value for classification

---

**Algorithm 1** Dynamic Classifier Selection for OVO scheme

---

```
1: procedure DYNAMIC OVO( $e, R$ )
2:    $k = 3 \cdot m$       ▷  $m$  is the number of classes
3:   repeat
4:      $Neighbours \leftarrow kNN(e)$ 
5:      $C \leftarrow Classes(Neighbours)$  ▷ We select
        the class labels in the neighbourhood
6:      $k++$ 
7:   until  $\#C > 1$  or  $k == 6 \cdot m$ 
8:   if  $C > 1$  then
9:      $R' \leftarrow [R - rows(i), cols(i)]; i \notin C$ 
10:    return  $R'$  ▷ A subset of the score matrix
11:  else
12:    return  $R$    ▷ Standard OVO approach
13:  end if
14: end procedure
```

---

# Dynamic OVO: Avoiding Non-competence

66/81

- Classify  $x$ , whose real class is  $c_1$

$$\square R(x) = \begin{pmatrix} & c1 & c2 & c3 & c4 & c5 \\ c1 & - & 0,55 & 0,6 & 0,75 & 0,7 \\ c2 & 0,45 & - & 0,4 & 1 & 0,8 \\ c3 & 0,4 & 0,6 & - & 0,5 & 0,4 \\ c4 & 0,25 & 0,0 & 0,5 & - & 0,1 \\ c5 & 0,30 & 0,2 & 0,6 & 0,9 & - \end{pmatrix}$$

# Dynamic OVO: Avoiding Non-competence

67/81

- Consider WV aggregation,  $c_2$  is predicted
- **None of the classifiers considering  $c_1$  failed**
- **Non-competent classifiers strongly voted for  $c_2$**

$$\square R(x) = \left( \begin{array}{cccccc|c} & c1 & c2 & c3 & c4 & c5 & WV \\ c1 & - & 0,55 & 0,6 & 0,75 & 0,7 & 2,6 \\ c2 & 0,45 & - & 0,4 & 1 & 0,8 & \mathbf{2,65} \\ c3 & 0,4 & 0,6 & - & 0,5 & 0,4 & 1,9 \\ c4 & 0,25 & 0,0 & 0,5 & - & 0,1 & 0,85 \\ c5 & 0,30 & 0,2 & 0,6 & 0,9 & - & 2,1 \end{array} \right)$$

# Dynamic OVO: Avoiding Non-competence

68/81

- Applying DynamickNN
  - ▣ Compute the  $k$ NN of  $x$  ( $k = 3 \cdot 5 = 15$ )
  - ▣ Subset of classes =  $\{c_1, c_4, c_5\}$
  - ▣ Remove  $\{c_2, c_3\}$  from the score-matrix
  - ▣ Apply WV to the reduced score-matrix

$$\square R_{dyn}(x) = \left( \begin{array}{cccccc|c} & c1 & c2 & c3 & c4 & c5 & WV \\ c1 & - & 0,55 & 0,6 & 0,75 & 0,7 & 1,45 \\ c2 & 0,45 & - & 0,4 & 1 & 0,8 & - \\ c3 & 0,4 & 0,6 & - & 0,5 & 0,4 & - \\ c4 & 0,25 & 0,0 & 0,5 & - & 0,1 & 0,35 \\ c5 & 0,30 & 0,2 & 0,6 & 0,9 & - & 1,2 \end{array} \right)$$

# Dynamic OVO: Avoiding Non-competence

69/81

## □ Summary:

- We **avoid** some of the **non-competent** classifiers by DCS
- It is **simple, yet powerful**
- Positive synergy between Dynamic OVO and WV
- All the differences are due to the aggregations
  - Tested with **same score-matrices** in all methods
  - Significant differences only changing the aggregation

# Outline

70/81

1. Introduction
2. Binarization
  - ▣ Decomposition strategies (One-vs-One, One-vs-All and Others)
  - ▣ State-of-the-art on Aggregations
    - One-vs-One
    - One-vs-All
3. Experimental Study
  - ▣ Experimental Framework
  - ▣ Results and Statistical Analysis
4. Discussion: Lessons Learned and Future Work
5. Conclusions for OVO vs OVA
6. Novel Approaches for the One-vs-One Learning Scheme
  - ▣ Dynamic OVO: Avoiding Non-competence
  - ▣ Distance-based Relative Competence Weighting Approach (DRCW-OVO)

M. Galar, A. Fernández, E. Barrenechea, F. Herrera, **DRCW-OVO: Distance-based Relative Competence Weighting Combination for One-vs-One Strategy in Multi-class Problems**. *Pattern Recognition* 48 (2015), 28-42, [doi: 10.1016/j.patcog.2014.07.023](https://doi.org/10.1016/j.patcog.2014.07.023).

# Distance-based Relative Competence Weighting Approach

71/81

## □ *Non-Competent Classifiers:*

- Those whose output is not relevant for the classification of the query instance
- They have not been trained with instances of the real class of the example to be classified

$$R(\mathbf{x}) = \begin{pmatrix} & \mathbf{c}_1 & \mathbf{c}_2 & \mathbf{c}_3 & \mathbf{c}_4 & \mathbf{c}_5 \\ \mathbf{c}_1 & - & 0.55 & 0.45 & 0.80 & 0.90 \\ \mathbf{c}_2 & 0.45 & - & 0.55 & 1.00 & 0.80 \\ \mathbf{c}_3 & 0.55 & 0.45 & - & 0.45 & 0.40 \\ \mathbf{c}_4 & 0.20 & 0.00 & 0.55 & - & 0.10 \\ \mathbf{c}_5 & 0.10 & 0.20 & 0.60 & 0.90 & - \end{pmatrix}$$

# Distance-based Relative Competence Weighting Approach

72/81

- *Non-Competent Classifiers:*
  - Consider WV aggregation,  $c_2$  is predicted
  - **None** of the classifiers **considering  $c_1$  failed**
  - **Non-competent** classifiers **strongly voted for  $c_2$**

$$R(\mathbf{x}) = \left( \begin{array}{cccccc|c|c} & \mathbf{c}_1 & \mathbf{c}_2 & \mathbf{c}_3 & \mathbf{c}_4 & \mathbf{c}_5 & V & WV \\ \mathbf{c}_1 & - & 0.55 & 0.45 & 0.80 & 0.90 & 3 & 2.70 \\ \mathbf{c}_2 & 0.45 & - & 0.55 & 1.00 & 0.80 & 3 & \mathbf{2.80} \\ \mathbf{c}_3 & 0.55 & 0.45 & - & 0.45 & 0.40 & 1 & 1.85 \\ \mathbf{c}_4 & 0.20 & 0.00 & 0.55 & - & 0.10 & 1 & 0.85 \\ \mathbf{c}_5 & 0.10 & 0.20 & 0.60 & 0.90 & - & 2 & 1.80 \end{array} \right)$$

# Distance-based Relative Competence Weighting Approach

73/81

- Designed to address the **non-competence** classifier problem
- It carries out a **dynamic adaptation** of the score-matrix
  - ▣ More competent classifiers should be those whose pair of classes are “**nearer**” to the query instance.
  - ▣ Confidence degrees are **weighted** in accordance to the former distance.
  - ▣ This distance is computed by using the standard **kNN approach**

# Distance-based Relative Competence Weighting Approach

74/81

## DRCW ALGORITHM FOR OVO STRATEGY

1. Compute the  $k$  nearest neighbors of each class for the given instance and store the average distances of the  $k$  neighbors of each class in a vector  $\mathbf{d} = (d_1, \dots, d_m)$ .
2. A new score-matrix  $R^w$  is created where the output  $r_{ij}$  of a classifier distinguishing classes  $i, j$  are weighted as follows,

$$r_{ij}^w = r_{ij} \cdot w_{ij},$$

where  $w_{ij}$  is the relative competence of the classifier on the corresponding output computed as

$$w_{ij} = \frac{d_j^2}{d_i^2 + d_j^2},$$

being  $d_i$  the distance of the instance to the nearest neighbor of class  $i$ .

3. Use weighted voting strategy on the modified score-matrix  $R^w$  to obtain the final class.

$$\text{Class} = \arg \max_{i=1, \dots, m} \sum_{1 \leq j \neq i \leq m} r_{ij} \cdot w_{ij}$$

Distance is computed with respect to all classes:

- ▣  $k \cdot m$  neighbors are used
- ▣  $k = 1$  is not the same as using 1NN classifier

**With  $k = 1$  a neighbor for each class is obtained, therefore it would use the  $m$  neighbours (1 per class). Next experimental example use  $k=5$ .**

# Distance-based Relative Competence Weighting Approach

75/81

- Classify  $x$ , whose real class is  $c_1$

$$R(\mathbf{x}) = \begin{pmatrix} & \mathbf{c}_1 & \mathbf{c}_2 & \mathbf{c}_3 & \mathbf{c}_4 & \mathbf{c}_5 \\ \mathbf{c}_1 & - & 0.55 & 0.45 & 0.80 & 0.90 \\ \mathbf{c}_2 & 0.45 & - & 0.55 & 1.00 & 0.80 \\ \mathbf{c}_3 & 0.55 & 0.45 & - & 0.45 & 0.40 \\ \mathbf{c}_4 & 0.20 & 0.00 & 0.55 & - & 0.10 \\ \mathbf{c}_5 & 0.10 & 0.20 & 0.60 & 0.90 & - \end{pmatrix}$$

# Distance-based Relative Competence Weighting Approach

76/81

- Distances to  $k$  nearest neighbors of each class ( $\mathbf{d}$ ) are computed:  $\mathbf{d} = (0.8, 0.9, 0.6, 1.2, 1.6)$
- A Weight-matrix  $W$  is computed to represent all  $w_{ij}$

$$\square W(x) = \begin{pmatrix} & c1 & c2 & c3 & c4 & c5 \\ c1 & - & 0,56 & 0,36 & 0,69 & 0,80 \\ c2 & 0,44 & - & 0,31 & 0,64 & 0,76 \\ c3 & 0,64 & 0,69 & - & 0,80 & 0,88 \\ c4 & 0,31 & 0,36 & 0,5 & - & 0,64 \\ c5 & 0,20 & 0,24 & 0,6 & 0,36 & - \end{pmatrix}$$

# Distance-based Relative Competence Weighting Approach

77/81

- Apply the weight-matrix  $W$  to the score-matrix  $R$
- $WV$  is applied to obtain the predicted class in DRCW-OVO

$$\square R^w(x) = \left( \begin{array}{c|cccccc} & c1 & c2 & c3 & c4 & c5 & WV \\ \hline c1 & - & 0,31 & 0,16 & 0,55 & 0,72 & \mathbf{1,74} \\ c2 & 0,20 & - & 0,17 & 0,64 & 0,61 & 1,66 \\ c3 & 0,35 & 0,31 & - & 0,36 & 0,35 & 1,37 \\ c4 & 0,06 & 0,00 & 0,11 & - & 0,06 & 0,24 \\ c5 & 0,02 & 0,05 & 0,07 & 0,32 & - & 0,47 \end{array} \right)$$

# Distance-based Relative Competence Weighting Approach

78/81

## □ Experimental Analysis

**Table 8**  
Average accuracy results in test of the representative combinations, DCS method and DRCW-OVO method (with  $k=5$ ) for each base classifier.

Data-set	C45			SVM <sub>Poly</sub>			SVM <sub>Phk</sub>			3NN			PDFC			Ripper		
	WV	DCS	DRCW	PE	DCS	DRCW	PE	DCS	DRCW	ND	DCS	DRCW	PC	DCS	DRCW	WV	DCS	DRCW
Autos	76.24	74.96	<b>80.96</b>	73.75	73.81	<b>79.48</b>	69.02	70.27	<b>71.45</b>	<b>78.88</b>	76.96	75.14	78.82	79.40	<b>80.74</b>	<b>85.09</b>	84.42	84.58
Car	94.68	94.50	<b>96.99</b>	93.58	93.58	<b>97.16</b>	64.99	<b>84.84</b>	81.65	93.57	93.40	<b>96.93</b>	99.77	<b>99.88</b>	99.42	92.59	93.52	<b>96.35</b>
Cleveland	52.55	53.55	<b>55.23</b>	58.97	<b>59.31</b>	58.66	47.53	47.87	<b>48.88</b>	<b>58.31</b>	57.96	56.61	53.92	55.93	<b>56.61</b>	52.18	54.54	<b>56.90</b>
Dermatology	95.24	<b>98.32</b>	98.06	94.71	94.99	<b>95.55</b>	97.20	97.20	<b>97.48</b>	92.14	95.49	<b>96.90</b>	84.66	<b>93.85</b>	91.90	93.32	94.43	<b>95.27</b>
Ecoli	81.06	81.94	<b>85.58</b>	79.37	79.63	<b>82.25</b>	77.11	77.11	<b>81.64</b>	81.66	82.52	<b>84.30</b>	84.07	83.78	<b>84.68</b>	78.47	78.74	<b>82.34</b>
Flare	<b>75.34</b>	73.62	75.27	75.43	75.46	<b>75.86</b>	69.28	<b>73.39</b>	72.04	71.21	71.59	<b>72.43</b>	73.64	<b>73.92</b>	73.69	75.24	74.83	<b>75.60</b>
Glass	72.03	71.63	<b>74.81</b>	62.14	63.14	<b>71.04</b>	73.72	74.15	<b>76.19</b>	73.35	74.27	<b>74.33</b>	68.72	<b>70.12</b>	<b>70.12</b>	68.56	68.12	<b>75.40</b>
Led7digit	64.51	<b>65.35</b>	65.33	67.90	<b>68.09</b>	66.47	61.33	61.57	<b>62.54</b>	66.68	67.88	<b>68.26</b>	62.17	62.60	<b>65.42</b>	63.98	63.86	<b>64.19</b>
Lymphography	74.50	<b>76.44</b>	<b>76.44</b>	82.48	82.48	<b>83.10</b>	81.87	81.87	<b>82.50</b>	68.19	<b>79.55</b>	79.52	<b>83.19</b>	<b>83.19</b>	<b>83.19</b>	75.68	75.68	<b>77.04</b>
Nursery	89.66	89.81	<b>90.90</b>	92.13	92.13	<b>94.53</b>	80.33	89.05	<b>90.83</b>	93.29	93.29	<b>93.68</b>	<b>97.92</b>	<b>97.92</b>	97.84	90.66	90.81	<b>92.44</b>
Pageblocks	95.64	95.46	<b>95.82</b>	94.90	94.53	<b>95.27</b>	94.58	94.76	<b>95.11</b>	<b>95.63</b>	95.46	95.09	<b>95.09</b>	94.91	<b>95.09</b>	95.45	95.11	<b>96.00</b>
Penbased	91.10	91.11	<b>95.64</b>	95.92	96.01	<b>97.01</b>	97.55	97.64	<b>98.00</b>	<b>97.00</b>	96.64	96.91	<b>98.19</b>	98.10	98.10	91.38	91.11	<b>96.01</b>
Satimage	82.15	82.92	<b>85.41</b>	84.48	84.16	<b>86.34</b>	84.77	85.70	<b>87.56</b>	87.58	87.73	<b>88.34</b>	86.79	86.95	<b>87.25</b>	82.61	82.14	<b>86.01</b>
Segment	96.28	96.71	<b>97.97</b>	92.68	92.90	<b>95.58</b>	97.23	97.36	<b>97.40</b>	96.58	96.80	<b>96.84</b>	97.32	<b>97.36</b>	97.27	96.58	96.88	<b>97.84</b>
Shuttle	99.59	99.68	<b>99.72</b>	96.55	97.61	<b>99.50</b>	99.59	99.63	<b>99.63</b>	<b>99.50</b>	99.40	99.40	97.43	98.03	<b>98.76</b>	99.40	<b>99.68</b>	99.54
Vehicle	72.33	72.81	<b>73.88</b>	73.53	74.00	<b>74.48</b>	81.92	81.92	<b>82.04</b>	72.11	72.23	<b>72.23</b>	<b>84.53</b>	84.40	84.41	69.27	70.20	<b>71.29</b>
Vowel	83.43	83.64	<b>94.75</b>	71.41	71.82	<b>95.05</b>	<b>99.70</b>	<b>99.70</b>	99.29	<b>97.78</b>	<b>97.37</b>	<b>97.27</b>	98.28	98.08	<b>98.59</b>	80.20	79.39	<b>94.44</b>
Yeast	59.57	59.84	<b>60.46</b>	60.52	59.98	<b>60.92</b>	59.31	59.51	<b>62.14</b>	56.68	56.54	<b>58.30</b>	60.25	59.98	<b>60.92</b>	58.30	58.10	<b>61.81</b>
Zoo	92.17	92.17	<b>93.22</b>	95.72	95.72	<b>96.77</b>	78.06	84.13	<b>90.80</b>	89.90	91.86	<b>94.64</b>	96.77	96.77	<b>97.82</b>	94.05	94.05	<b>96.10</b>
Average	81.48	81.81	<b>84.02</b>	81.38	81.55	<b>84.48</b>	79.74	81.98	<b>83.01</b>	82.63	83.52	<b>84.06</b>	84.29	85.01	<b>85.36</b>	81.21	81.35	<b>84.17</b>

# ?Questions?

79/81

- Thank you for your attention!

