

Predictive models based on SVM with structured output

Jovana Kovačević, Gordana Pavlović-Lažetić

Faculty of Mathematics, Belgrade University

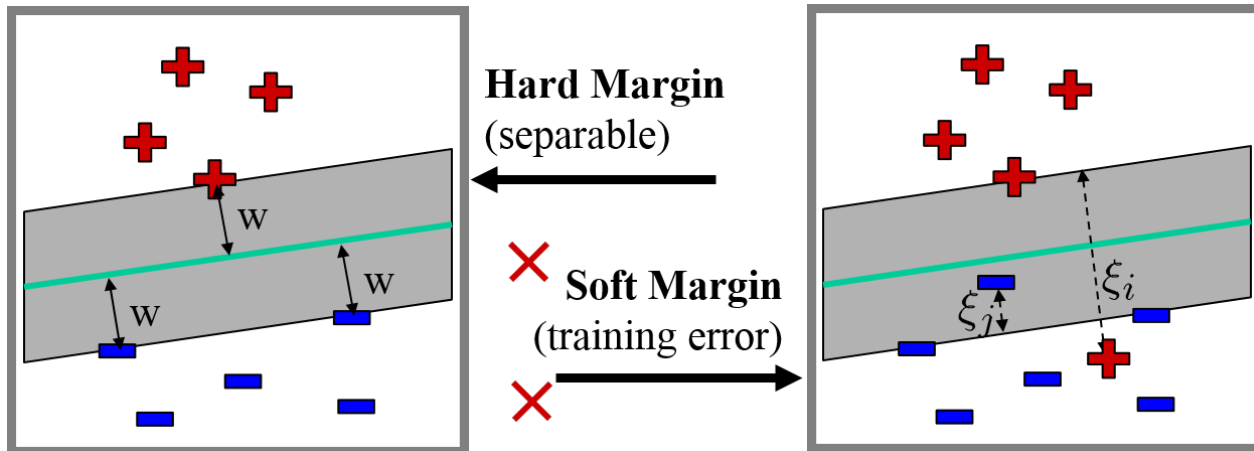
`{jovana,gordana}@matf.bg.ac.rs`

Mathematical Data Science Workshop
Mathematical institute, Belgrade, June 21st 2015

Structural classification vs binary classification

- binary classification: $y \in 0, 1$
- structural classification: $y \in Y$, Y is structured (sequence, tree, graph, ...)
- SSVM - generalization of SVM method on structural output (*T. Joachims et al, 2004*)

SVM



$$f(x) = \langle x, w \rangle + b$$
$$y = \text{sgn}(f(x))$$

Minimize

$$\frac{\|w\|^2}{2} + \sum_{i=1}^n \xi_i$$

such that $y_i \cdot f(x_i) \geq 1 - \xi_i$ for each training example (x_i, y_i)

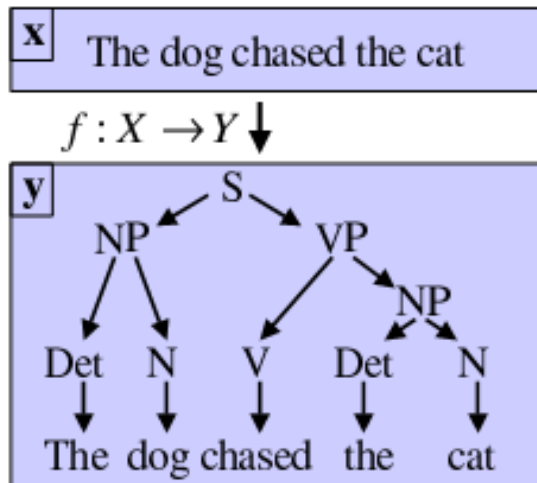
Generalization of SVM on structural output: SSVM

- Goal: for training set $(x_1, y_1), \dots, (x_n, y_n)$, find a function $f : X \rightarrow Y$
- Idea: $F : X \times Y \rightarrow \mathbb{R}$

$$f(x) = \operatorname{argmax}_{y \in Y} (F(x, y))$$

SSVM

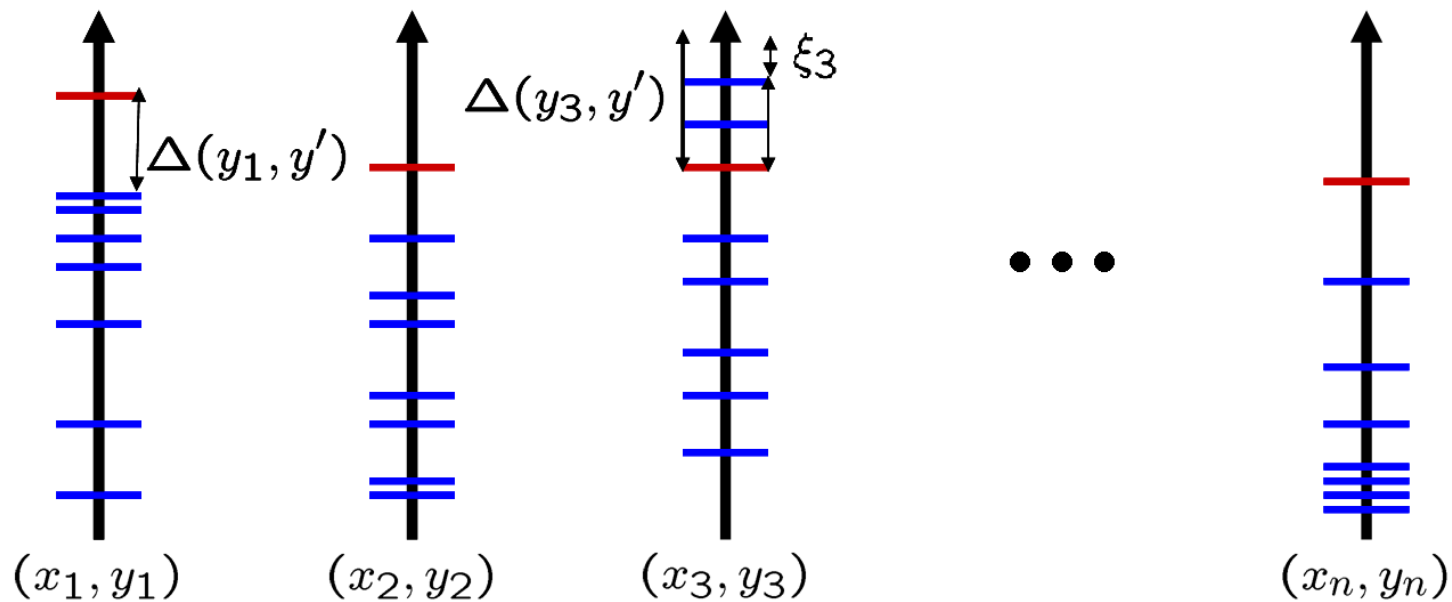
- let us assume: $F(x, y) = \langle \Psi(x, y), w \rangle$
- function Ψ is joint representation of input data x and output data y



$$\Psi(\mathbf{x}, \mathbf{y}) = \begin{pmatrix} 1 \\ 0 \\ 2 \\ 1 \\ \vdots \\ 0 \\ 2 \\ 1 \\ 1 \\ 1 \end{pmatrix} \begin{matrix} S \rightarrow NP VP \\ S \rightarrow NP \\ NP \rightarrow Det N \\ VP \rightarrow V NP \\ \\ Det \rightarrow dog \\ Det \rightarrow the \\ N \rightarrow dog \\ V \rightarrow chased \\ N \rightarrow cat \end{matrix}$$

SSVM - Margin in structured setting

- We would like coefficients w to be such that $\langle \Psi(x, y), w \rangle$ is highest for correct y



SSVM formulation of optimization problem

$$\min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{n} \sum_{i=1}^n \xi_i$$

subject to $\forall i \in \{1, \dots, n\}$

$$\forall \mathbf{y} \in Y : \langle \mathbf{w}, \Psi(x_1, \mathbf{y}_1) \rangle - \langle \mathbf{w}, \Psi(x_1, \mathbf{y}) \rangle \geq \Delta(\mathbf{y}_1, \mathbf{y}) - \xi_1$$

...

$$\forall \mathbf{y} \in Y : \langle \mathbf{w}, \Psi(x_n, \mathbf{y}_n) \rangle - \langle \mathbf{w}, \Psi(x_n, \mathbf{y}) \rangle \geq \Delta(\mathbf{y}_n, \mathbf{y}) - \xi_n$$

- quadratic optimization with linear constraints
- $n|Y|$ constraints!

Algorithm 1 Struct SVM learning algorithm (Cutting plane)

- 1: $W_i = \emptyset$
- 2: **repeat**
- 3: **for** $i = 1, \dots, n$ **do**
- 4: out of all $\hat{y}_i \in Y$, find \hat{y}_i that violates constraint i
most (argmax problem: $\arg \max_{\hat{y}_i \in Y} (\Psi(x_i, \hat{y}_i) + \Delta(y_i, \hat{y}_i))$)
- 5: **if** \hat{y}_i violates constraint for more than ϵ **then**
- 6: update set W_i : $W_i = W_i \cup \hat{y}$
- 7: solve optimization problem with constraints:

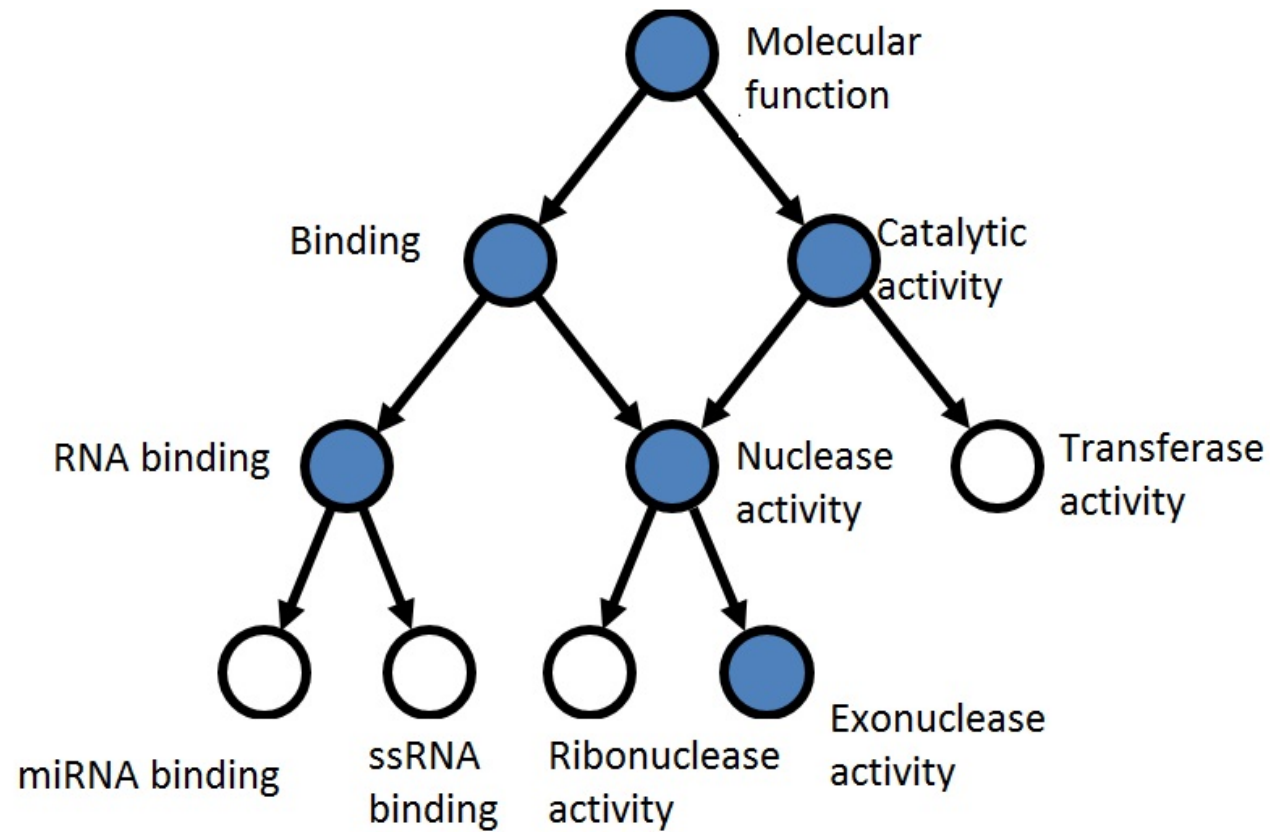
$$\forall y \in W_1 : \langle w, \Psi(x_1, y_1) \rangle - \langle w, \Psi(x_1, y) \rangle \geq \Delta(y_1, y) - \xi_1$$

...

$$\forall y \in W_n : \langle w, \Psi(x_n, y_n) \rangle - \langle w, \Psi(x_n, y) \rangle \geq \Delta(y_n, y) - \xi_n$$

- 8: **until** each set W_i is unchanged
-

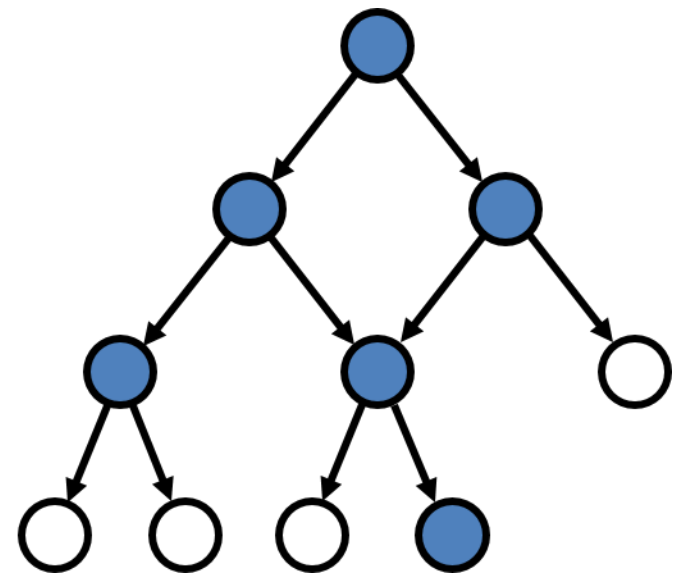
Protein function prediction



Protein function prediction

>sp|P04637|P53_HUMAN

```
MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIE  
QWFTEDPGPDEAPRMPEAAPPVAPAPAAPTPAAPAPAPSWPLSSSVPSQKT  
YQGSYGFRLLGFLHSGTAKSVTCTYSPALNKMFCQLAKTQPVQLWVDSTPPP  
GTRVRAMAIYKQSQHMTDEVVRRCPHHERCSDSDGLAPPQHLLIRVEGNLRVE  
YLDDRNTFRHSVVVPYEPPEVGSDCTTIHYNMCMNSSCMGGMNRRPILTI  
TLEDSSGNLLGRNSFEVRVCACPGRRRTEENLRKKGEPHHELPPGSTKR  
ALPNNTSSSPQPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGK  
EPGGSRAHSSHLKSKKGQSTSRHKKLMFKTEGPDSD
```



Implementation of Ψ function

- input: histogram of tetragrams
- output: sparse 0/1 vector, $y_i(j) = 1$ if protein i contains node j in its function DAG
- dimension of vector Ψ : total number of tetragrams \times total number of nodes that appear in dataset

$$\Psi(x, y) = [\underbrace{0}_{y_1}, \dots, \underbrace{0}_{y_{a-1}}, \underbrace{x}_{y_a}, \underbrace{0}_{y_{a+1}}, \dots, \underbrace{0}_{y_{b-1}}, \underbrace{x}_{y_b}, \underbrace{0}_{y_{b+1}}, \dots, \underbrace{0}_{y_{c-1}}, \underbrace{x}_{y_c}, \underbrace{0}_{y_{c+1}}, \dots, \underbrace{0}_{y_k}]$$

Algorithm 2 Augmented inference algorithm

```
1: Input: training instance  $(x_i, y_i)$ 
2: Output:  $y_{best}$  that maximizes  $H(x_i, y)$  over  $y \in Y$ 
3: Initialization:  $L = \{y_{root}\}, y_{best} = \emptyset, H_{best} = -\infty$ 
4: repeat
5:    $y_{head} :=$  first element from  $L$ 
6:    $Y_{ext} :=$  all extensions of  $y_{head}$  by one node
7:   for each  $y_{ext} \in Y_{ext}$  do
8:     if  $i(y_{ext}) \geq imax$  then
9:       continue
10:    insert  $y_{ext}$  in sorted linked list  $L$ 
11:    if  $H(y_{ext}) > H_{best}$  then
12:      update  $y_{best}, H_{best}$ 
13:  remove  $y_{head}$  from  $L$ 
14:  increment  $step$ 
15: until  $step > smax$  or  $L$  is empty
```

Parameter imax

- model ontology with Bayesian network:

$$Pr(T) = \prod_{v \in T} Pr(v | \mathcal{P}(v))$$

- information content of a DAG T

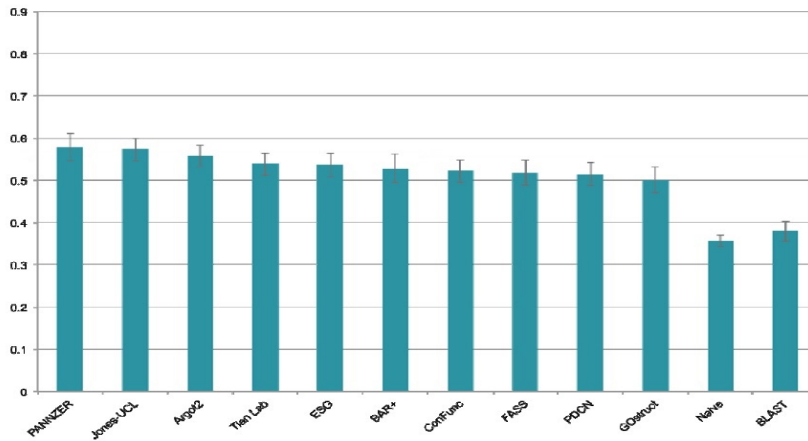
$$\begin{aligned} i(T) &= \log \frac{1}{Pr(T)} \\ &= \sum_{v \in T} \log \frac{1}{Pr(v | \mathcal{P}(v))} = \sum_{v \in T} ia(v) \end{aligned}$$

- information assertion of a node v - rare nodes have high $ia(v)$

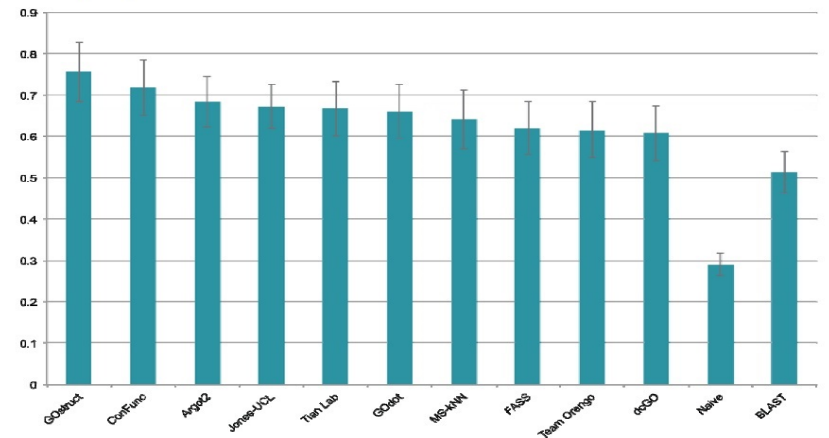
Radivojac, Clark, 2013

Results - comparison with CAFA

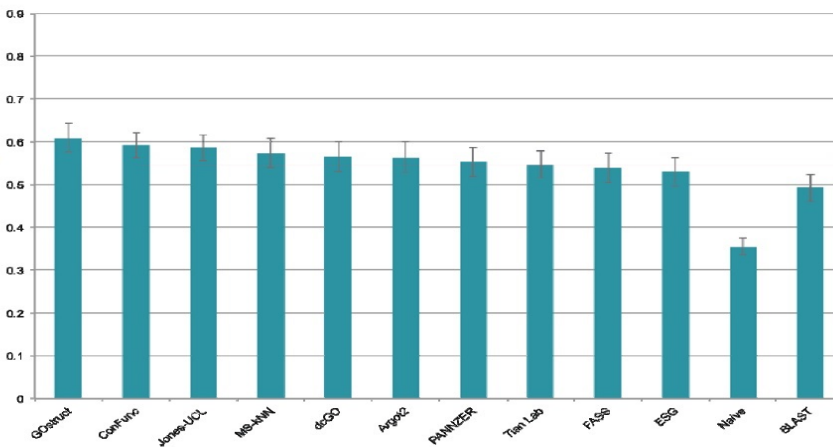
Supplementary Figure 7C: Molecular Function – *H. sapiens*.



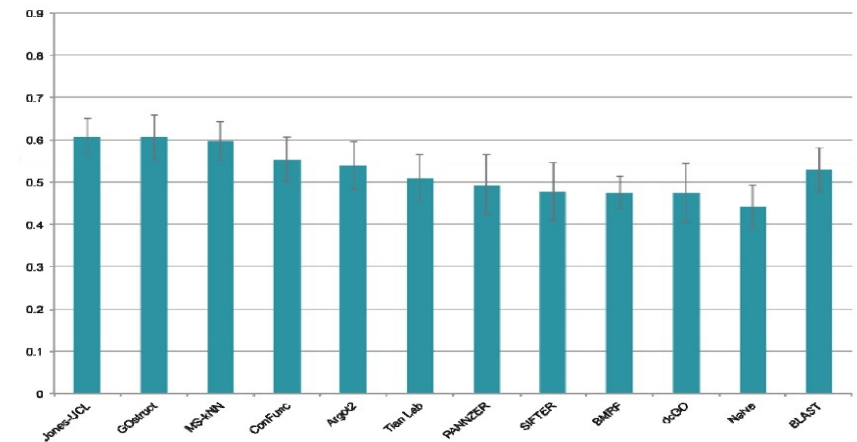
Supplementary Figure 7A: Molecular Function – *A. thaliana*.



Supplementary Figure 7D: Molecular Function – *M. musculus*.

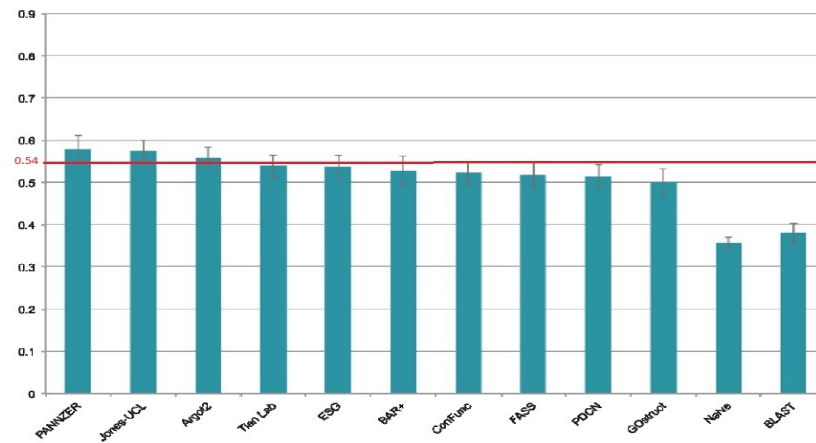


Supplementary Figure 7E: Molecular Function – *R. norvegicus*.

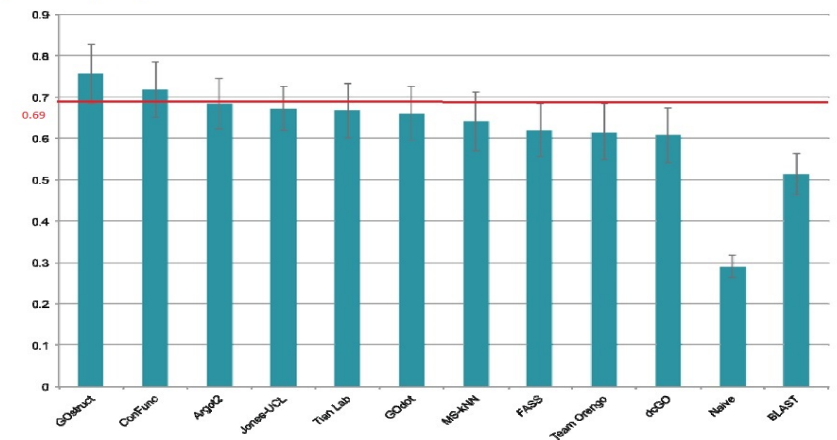


Results - comparison with CAFA

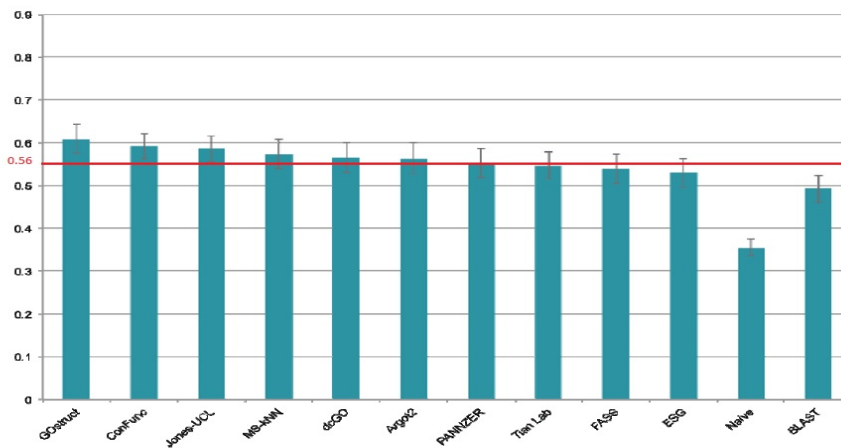
Supplementary Figure 7C: Molecular Function – *H. sapiens*.



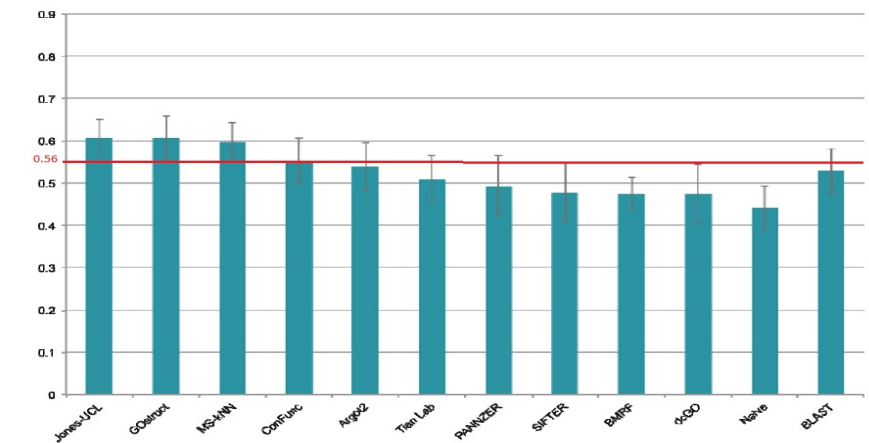
Supplementary Figure 7A: Molecular Function – *A. thaliana*.



Supplementary Figure 7D: Molecular Function – *M. musculus*.

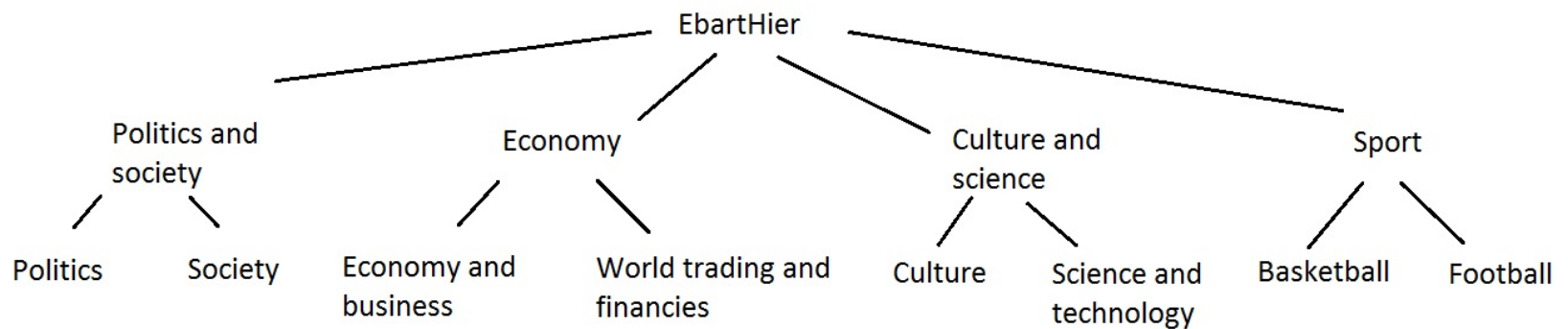


Supplementary Figure 7E: Molecular Function – *R. norvegicus*.



Text classification

- hierarchical classification and multiclassification of text corpora in different languages



Text classification

- each text is represented by ngrams of bytes, $n \in \{2, 3, 4, \dots\}$
- input: each position of vector x corresponds to its tf-idf statistics in the training set
- output in multiclassification: a class from the set $\{1, 2, \dots, N\}$, where N is total number of classes
- output in hierarchical classification: a path from the root to the corresponding leaf

Results

	size of dataset	number of classes	language	our method	result	best known result
Ebart	60637	8	Serbian	hier flat	90.17 90.43	Graovac, 2012, knn: 88.5
Tancorp	14150	60	Chinese	hier flat	87.38 87.19	Li, 2010, knn: 79.37
mesleh-10	7842	20	English	flat	93.44	Mesleh, 2011, SVM: 91.41
20 newsgroups	20000	8	Arabic	flat	83.75	Lan et al, 2009, SVM: 80.81

Thank you for your attention! Questions?