



Protein sequences classification by means of feature extraction with substitution matrices

Rabie SAIDI
Mondher Maddouri
Engelbert Mephu Nguifo

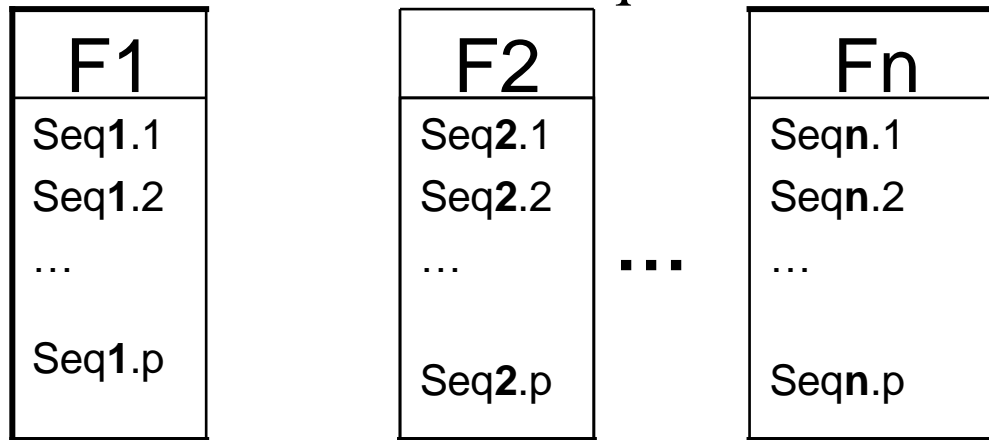


JOBIM'10, Montpellier

Content

- 1 Scope & motivation
- 2 Proposed encoding method : DDSM
- 3 Experiments & results

Families of sequences

New sequence S 

- Alignment
 - BLAST, FAST
- HMM profiles
 - HMMER, SAM
- A query belongs to the subject class with the best hit score
- Fundamentally depend on homology
 - Fail when classes contain distant sequences
- Non-discriminative (a score accompanied by e-value)

Scope & motivation How to involve machine learning

Motif extraction selection



- RAF
- SSH
- FFIT
- MKGD
- FYCG
- VLAA
- LVLQH
- VMVWI
- SYNTV



Encoding

```

>SEQ 1
MEIPAVTEP<b>SYNTV</b>AKNDFMSGFLCFSINV<b>RAF</b>GITVPTPLYSLVFIIGVIGHVLV<b>LVLQH</b>KRLRNMTSIYLFNLAISDLVFLSTLPFWV
DY<b>MKGD</b>WIFGNAMCKFVSGFYLLGLYSDFITLLTIDRYLAVVHVVFALRARTVTFGIISIIHW<b>VLAA</b>LVSIPLCYVFKSQMEFTYH
TCRAILPRKSLIRFLR

>SEQ 2
MAATASQPLATEDADSENSFPYYD<b>RAF</b>YLDEVAFMLCRKDAVVSFGKVFPLPVFYSLIFVLGLSNLLLMVLLRYVPRRRMVEIYLL
NLAISNLLFLVTLPFWGISVAW <b>MKGD</b>HWWVFGSFLCKMHCCFSPILYAF<b>SSH</b>RFRQYLKAFLAAVLGWHLA

>SEQ 3
MPTVASPLPLTTV<b>SYNTV</b>GSENSSIYDYDYLDDMTILVCRKDEVLSFGRVFLPVVYSLIFVLGLAGNLLLVLLHSAPRRRTMELYLLN
LAVSNLLFVVTMPFWAISVAWHVWVFGSFLCKVISTLYSINFYCGI<b>FFIT</b>CMSLDKYLEIVHAQPL<b>FYCG</b>HRPKAQFRNLLI<b>VMVW</b>TSL
AISVPEMLTLFLHSLDLHVF <b>MKGD</b>GNCEISHRLDYT<b>LVLQH</b>LQVTESLAFSHCCFT

>SEQ 4
MPTIASPLPLATTGPENGSSIIYDYDYLDDVTLVCS<b>RAF</b>EVLVSFGRVFLPVVYSLIFVL<b>VLAA</b>GLAGNLLLVLLHSVPQRRRMIELYLL
NLAVSNLLFVVTMPFWAISVAWHVWVFGSFLCKVSTLYSINFYCGI<b>FFIT</b>CMSLDKYLEIVHAQPL<b>VMVWI</b>LHRPKTRFRNLLVWIT
    
```

| | RAF | SSH | FFIT | MKGD | FYCG | VLAA | LVLQH | VMVWI | SYNTV |
|-------|-----|-----|------|------|------|------|-------|-------|-------|
| SEQ 1 | 1 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 1 |
| SEQ 2 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| SEQ 3 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 1 |
| SEQ 4 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 |

- Data: protein sequences
 - Alphabet: 20 amino acids
 - Format inadequate for ML ou DM
 - → Preprocessing is needed
- Preprocessing:
 - Looking for descriptors / motifs
 - Using them as attributes
 - structural and functional importance of preserved regions
 - These regions can be used as descriptors for the bio-sequences [Nevill-Manning et al, 1998]
- Loss of information ?
 - Due to format change
- Efficient preprocessing → better results
 - Exp: accuracy in classification

- N-Grams [Shannon, 1951]
- Active Motifs [Wang et al, 1999]
- Discriminative Descriptors [Maddouri et al, 2004]
- Amino Acid Composition [Zhang et al, 2003]
- Functional Domain Composition [Yu et al, 2006]

- The mentioned methods neglect the substitution
- Some amino acids have similar proprieties
 - Some substitution can be without effect on the function nor the structure of the protein
 - Same thing can be deduced for features
- Idea : use substitution matrices to define similarity between features
 - Only one of the substitutable features is kept
 - Number of features will be reduced
 - Any impact on classification?

- Let M and M' 2 motifs composed of k amino acids
- For each amino acid of M with index i , its probability P_i to not mutate is :

$$P_i = S(M[i], M[i]) / \sum_{j=1}^{20} S^+(M[i], AA_j)$$

- AA_j is the amino acid with index j among the 20 amino acids

- The probability P_m that M mutate to other motifs, is :

$$P_m(M) = 1 - \prod_{i=1}^k P_i$$

- $P_m(\text{LLK}) = 1 - (4/9 * 4/9 * 5/9) = 0.89$

- We note by $S_m(M, M')$ the score of substitution of a motif M' by a motif M :

$$S_m(M, M') = \sum_{i=1}^k S(M[i], M'[i])$$

- $S_m(\text{LLK}, \text{VMK}) = 8$

- $S_m(\text{LLK}, \text{LLK}) = 13$

- We note by $PS(M, M')$ the probability of substitution of a motif M' by a motif M

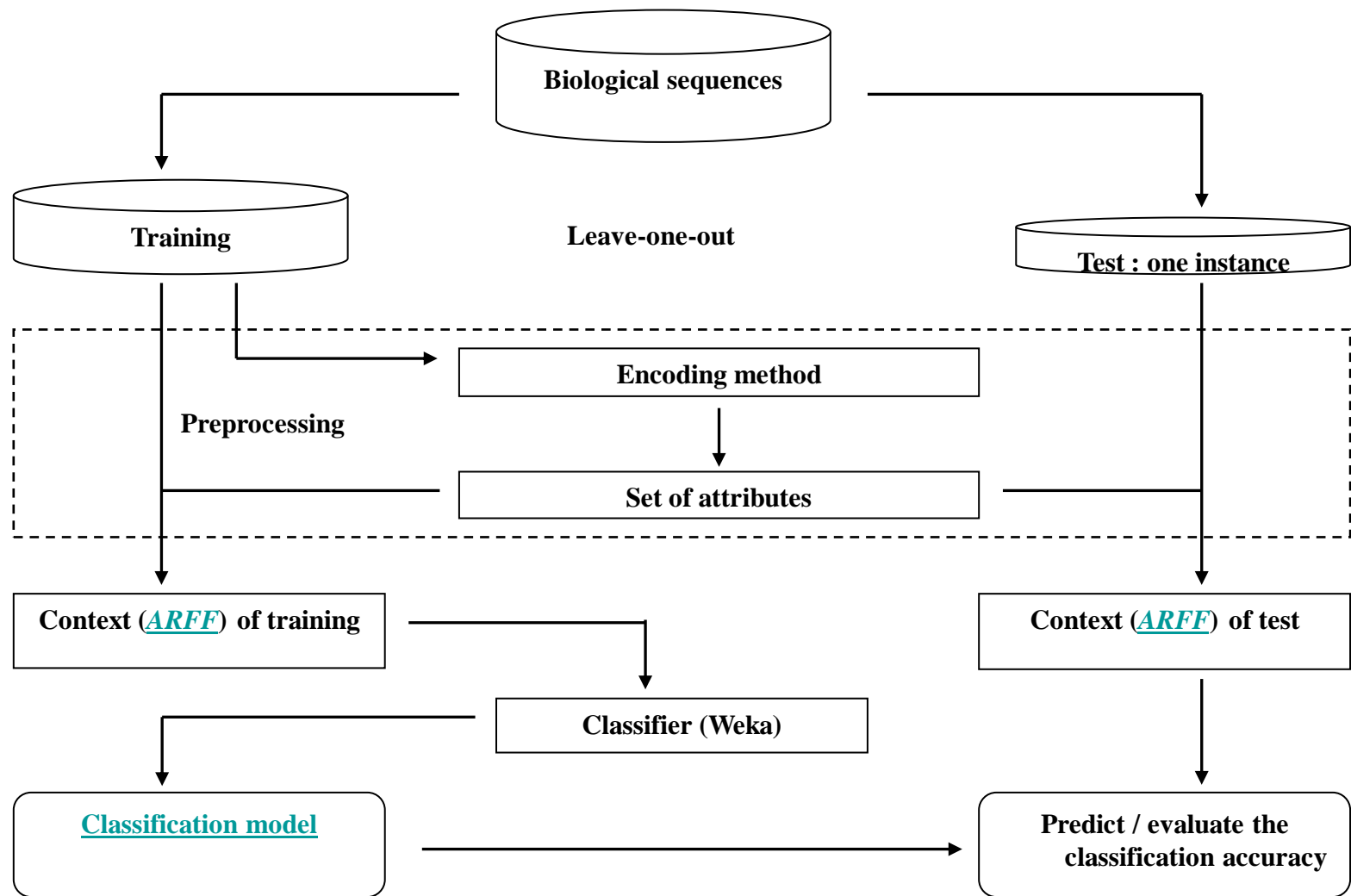
$$PS(M, M') = S_m(M, M') / S_m(M, M)$$

- $PS(\text{LLK}, \text{VMK}) = 8 / 13 = 0.61$

- A motif M substitutes a motif M' if :
 - M and M' have the same length k ,
 - $S(M[i], M'[i]) \geq 0, i = 1.. k$,
 - $PS(M, M') \geq T$, where T is a user-specified threshold : $0 \leq T \leq 1$.

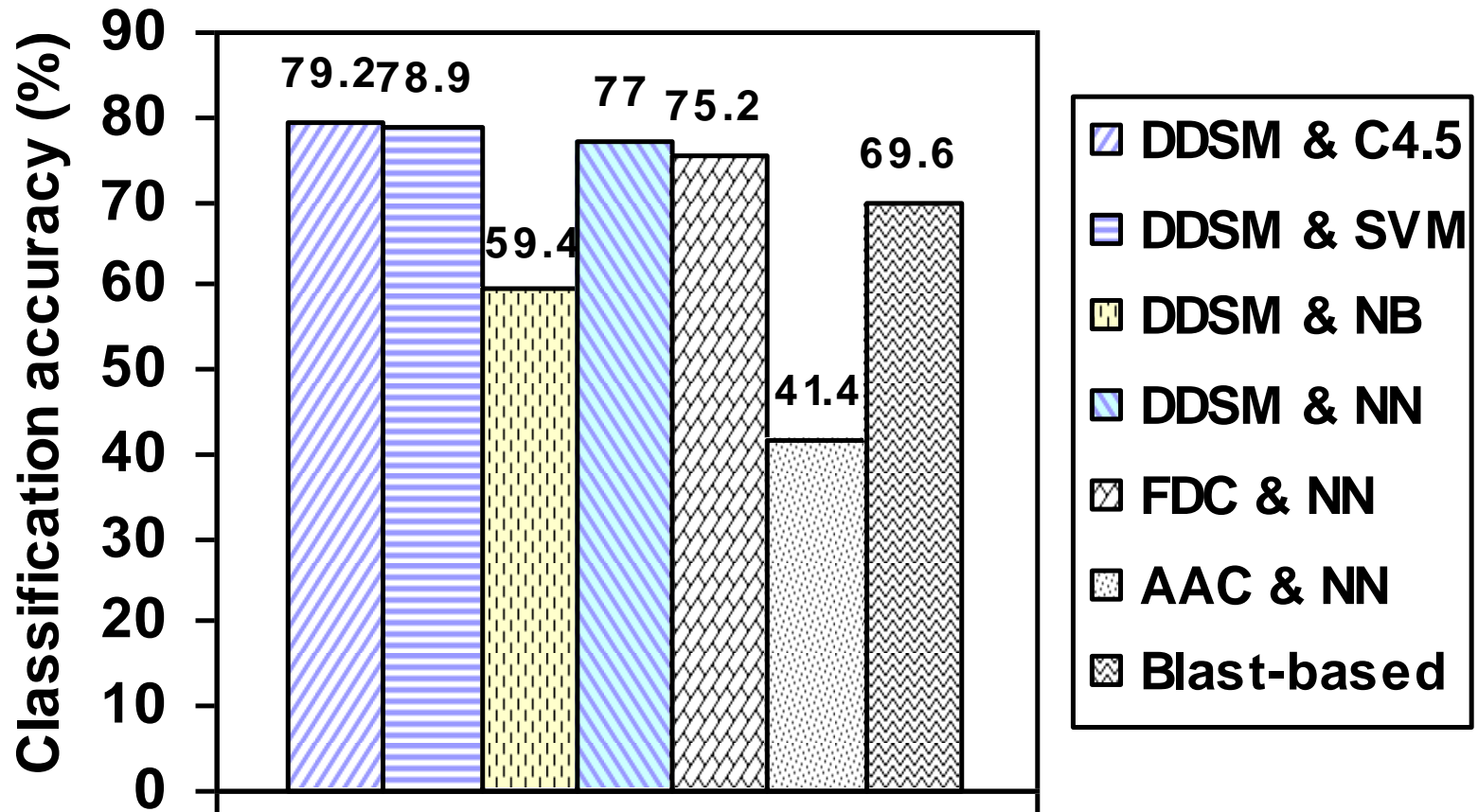
- Research of repeated words (adaptation of the KMR algorithm [Karp et al, 1972])
 - Equivalence notion
- Discrimination : \mathbf{x} is discriminative of a family \mathbf{fi}
 - Occurrence rate of \mathbf{x} in $\mathbf{fi} \geq \alpha$
 - Occurrence rate de \mathbf{x} in others $\mathbf{f} \leq \beta$
- Minimality
 - A sub-sequence is called minimal if it does not contain other discriminative sub-sequences.
- Example

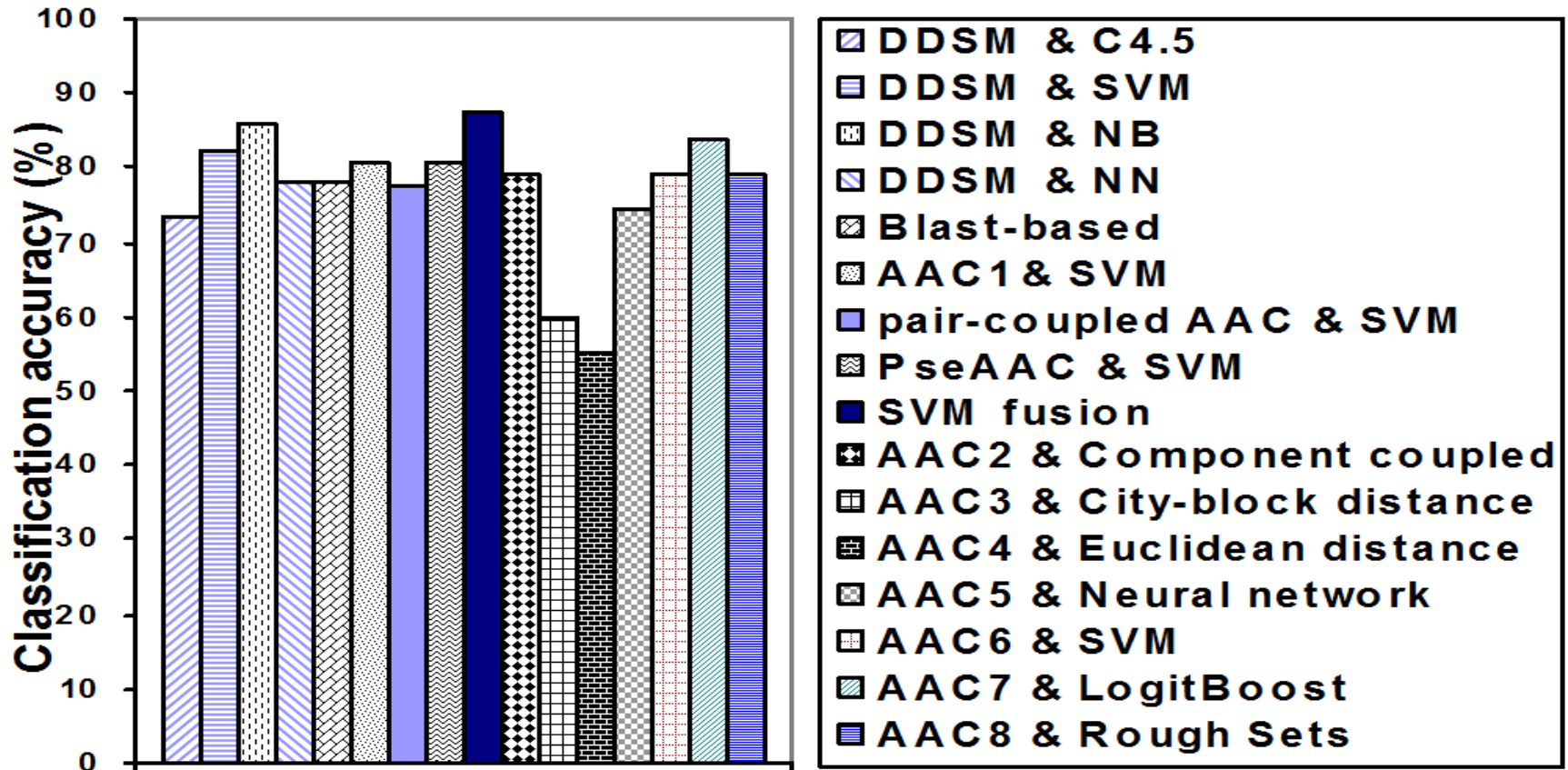
- Identification of main motifs and filtering
 - The motifs are clustered in groups
 - Group : a **main motif** *MM* + other motifs
 - *MM* substitutes all the motifs in its group
 - *MM*: the most likely motif in its cluster to mutate to other motifs i.e having the highest *Pm* in its group
 - Filtering : keep only *MM*
 - [Example](#)
- Construction of the context (binary vectors)
 - Compare each motif with the k-grams (k=length of motif) of each sequence
 - Until find a substitute and note 1 in the binary table
 - Or cross the whole sequence without finding a substitute and note 0
 - → Richer context



Experiments & results data

| Dataset (source) | Identity percentage | Family/class | Size | Total |
|--|---------------------|---------------------------|------|-------|
| DS1 (Yu et al, 2006) <i>BMC Bioinformatics</i> | 25 % | Monomer | 208 | 717 |
| | | Homodimer | 335 | |
| | | Homotrimer | 40 | |
| | | Homotetramer | 95 | |
| | | Homopentamer | 11 | |
| | | Homoheptamer | 23 | |
| | | Homooctamer | 5 | |
| DS2 (Zhou et al, 2006) <i>Anal Biochim</i> | 84 % | All- α domain | 70 | 277 |
| | | All- β domain | 61 | |
| | | α / β domain | 81 | |
| | | α + β domain | 65 | |






| Substitution matrix | Attributes | Accuracy (%) | | | |
|---------------------|------------|--------------|------|------|------|
| | | C4.5 | SVM | NB | NN |
| Blosum45 | 377 | 78.5 | 79.2 | 59.4 | 77.7 |
| Blosum62 | 508 | 79.2 | 78.9 | 59.4 | 77 |
| Blosum80 | 532 | 77.6 | 80.5 | 60 | 77.6 |
| Pam30 | 2873 | 77.8 | 82 | 60.3 | 76.7 |
| Pam70 | 802 | 78.1 | 80.5 | 60.5 | 77 |
| Pam250 | 1123 | 77.3 | 79.4 | 59.6 | 78.7 |

| Substitution matrix | Attributes | Accuracy (%) | | | |
|---------------------|------------|--------------|------|------|----|
| | | C4.5 | SVM | NB | NN |
| Blosum45 | 2603 | 69.3 | 82.3 | 85.9 | 78 |
| Blosum62 | 3083 | 73.3 | 82.3 | 85.9 | 78 |
| Blosum80 | 3146 | 70.1 | 82.3 | 84.1 | 78 |
| Pam30 | 3830 | 69.3 | 82.3 | 84.5 | 78 |
| Pam70 | 3822 | 70.4 | 82.3 | 84.5 | 78 |
| Pam250 | 969 | 66.1 | 85.2 | 79.4 | 78 |

- Motif-based encoding
 - May allow reliable description of protein
 - Allow the injection of external information (pH, temperature,...)
- DDSM (discriminative descriptors with substitution matrix)
 - A discriminative encoding method taking the substitution into account
 - Low number of features
 - Helps with classification task
 - Coupled with SVM : efficient protein classifier
- Blosum matrices with higher numbers and Pam matrices with low numbers allow the building of fewer features
- Variances of accuracies are slight when varying the substitution matrices with the same classifier

- More can be found in [Saidi et al., BMC Bioinformatics 2010]

| | | | |
|------------------------|--|-----------------------------|--|
| Top | Research article | Open Access | BMC Bioinformatics Volume 11 |
| Abstract | Protein sequences classification by means of feature extraction with substitution matrices | | |
| Background | | | |
| Methods and R... | Rabie Saidi ^{1,2,3,4} ✉, Mondher Maddouri ^{4,5} ✉ and Engelbert Mephu Nguifo ^{1,2} ✉ | | |
| Discussion an... | <ol style="list-style-type: none"> 1 LIMOS - Blaise Pascal University - Clermont University, BP 10448, Clermont-Ferrand 63000, France 2 LIMOS - CNRS UMR 6158, Aubière 63173, France 3 Department of Computer Science - FSJ - University of Jendouba, UMA Street, Jendouba 8100, Tunisia 4 URPAH - FST - University of Tunis El Manar, Academic Campus, Tunis 2092, Tunisia 5 Department of Computer Science - FSG - University of Gafsa, Campus of Sidi Ahmed Zarroug, Gafsa 2112, Tunisia | | |
| Competing interests | ✉ author email ✉ corresponding author email | | |
| Authors' contributions | <i>BMC Bioinformatics</i> 2010, 11 :175 doi:10.1186/1471-2105-11-175 | | |
| Acknowledgements | The electronic version of this article is the complete one and can be found online at: http://www.biomedcentral.com/1471-2105/11/175 | | |
| References | <p>Received: 4 September 2009 Accepted: 8 April 2010 Published: 8 April 2010</p> <p>© 2010 Saidi et al; licensee BioMed Central Ltd.</p> | | |
| | | | <p>Viewing options:</p> <ul style="list-style-type: none"> ▪ Abstract ▪ Full text ▪ PDF (947KB) ▪ Additional files <p>Associated material:</p> <ul style="list-style-type: none"> ▪ Readers' comments  ▪ PubMed record <p>Related literature:</p> <ul style="list-style-type: none"> ▪ Articles citing this article <ul style="list-style-type: none"> on Google Scholar on PubMed Central ▪ Other articles by authors <ul style="list-style-type: none"> ⊕ on Google Scholar ⊕ on PubMed ▪ Related articles/pages <ul style="list-style-type: none"> on Google on Google Scholar on PubMed |

THANKS

This work was partially supported by the
French-Tunisian project CMCU-Utique
05G1412 and the LifeGrid
PREFON_META project

JOBIM'10, Montpellier