

Supplementary

MiRenSVM: Towards Better Prediction of MicroRNA Precursors

Using an Ensemble SVM Classifier with Multi-loop Features

Jiandong Ding¹, Shuigeng Zhou^{1§} and Jihong Guan^{2§}

¹Shanghai Key Lab of Intelligent Information Processing, and School of Computer Science, Fudan University, Shanghai 200433, China.

²Department of Computer Science and Technology, Tongji University, Shanghai 201804, China.

[§]Corresponding author

Email addresses:

Jiandong Ding: jdding@fudan.edu.cn

Shuigeng Zhou: sgzhou@fudan.edu.cn

Jihong Guan: jhguan@tongji.edu.cn

Table S1: Homo sapiens pre-miRNAs with multi-loops secondary structure [1]:

Database	Num_loops	Num_total	Proportion(%)
miRBase 5.0	14	207	6.76
miRBase 8.2	36	462	7.79
miRBase 13.0	34	715	4.76
Total:	84	1384	6.07

*Num_loops: number of *hsa* pre-miRNA with multi-loop secondary structure

*Num_total: number of *hsa* pre-miRNA in the database

Table S2: Homo sapiens pre-miRNAs in miRBase13.0 [2] whose MFE is higher than -16kal/mol (predicted by RNAfold [3] under the default parameters)

>hsa-mir-1973 MI0009983
UAUGUUAACACGGCCAUGGUAUCCUGACCGUGCAAAGGUAGCAUA
(((((.....(((.....)))))).....)))) (-10.80)
>hsa-mir-1978 MI0009988
UAGACGGGCUCACAUCACCCCAUAAACAAGUUUGGUCCUAGCCUUUCUA
(((((.....(((.....)))))).....)))) (-13.70)
>hsa-mir-2052 MI0010486
CUGUUUUGAUAAACAGUAAUGUCCUUUAGUUCAAAGUUAACCAGCUAUCAAAACAA
.(((((((.....)))))).....)) (-10.69)
>hsa-mir-2054 MI0010488
CUGUAAUUAUUUUUUUUUUUUCUCUAUCAUUAAAAAUGUAUUACAG
(((((.....(((.....)))))).....)))) (-8.20)
>hsa-mir-198 MI0000240
UCAUUGGUCCAGAGGGGAGAUAGGUUCCUGUGAUUUUUCCUUCUCUAUAGAAUAAAUGA
(((((.....(((.....)))))).....)))) (-15.44)
>hsa-mir-384 MI0001145
UGUUAUAAUCAGGAAUUUUAACAUAUCCUAGACAUAUUGUAUAAUGUUAUAGUCAUAGUCAUCCUAGAAUUGUUAUAAUGCCUGUAACA
.(((((((.....)))))).....)) (-15.00)
>hsa-mir-923 MI0005715
UAUUUGUCAGCGGAGGAAAAGAAACUAACCAGGAUCCUCAGUAAUGGCAGUG
(((((.....(((.....)))))).....)))) (-13.00)
>hsa-mir-924 MI0005716
AAUAGAGUCUUGUGAUGUCUUGCUUAAGGGCCAACCUAGAGUCUACAAC
..(((((((.....)))))).....)) (-12.30)
>hsa-mir-1279 MI0006426
AUAUUCACAAAAUUCUAUUGUCUUCUUAUUGCCAGAAAGAAGAGUAUAGAACUUC
.....(((((((.....)))))).....)) (-12.10)
>hsa-mir-1308 MI0006441
CCCCGCAUGGGUGGUUCAGUGGCAGAAUUCUCAAUUGUAUCCCAUAAUCCC
.....(((((((.....)))))).....)) (-12.60)
>hsa-mir-1321 MI0006652
ACAUUAUGAAGCAAGUAUUUUUUAUCCUGUUUUACAUAUAGGAAUAAACUCAGGGAGGUGAAUGUGAUCAAAGAUAG
.....(((((((.....)))))).....)) (-11.60)
>hsa-mir-1322 MI0006653
AGUAUCAUGAAUUAGAAACCUACUUUUACAUAUAGUUUACAUAAGAAGCGUGAUGAUGCUGCUGAUGCUGUA
(((((.....(((.....)))))).....)))) (-11.30)

Table S3: Sequences of 14 *hsa* and 13 *aga* pre-miRNA whose identity is lower than 90% (predicted by CD-HIT [4])

13 <i>aga</i> pre-miRNAs:
>aga-mir-137 MI0010601
AAAACUUGGUUGGCCACGCGUAUUCUUGGGUUACUAACACACUGUUUAUGUUGUUUAUUGCUUGAGAAUACACGUAUUGACUAGUGUUGUA
>aga-mir-190 MI0010594
UGUUUGGUGGGACAGUUUCGGUGAGAUUGUUUGAUUUCUUGGUUGUUAAGAUUUCAAUUUACCCAGGAAUCAACAUAUUUUACCGUG ACUGUCGU
>aga-mir-263b MI0010606
UGACCAAUAUGGACCUUGGCACUGGGAGAAUUCACAGUGGAUCGUACCAUAUCGUUUUCUGUGGAUCUUUCGUGCCAUCGUUCAGAUUUGGU GCC
>aga-mir-286 MI0010605
GUCAUAAUUUGGGCGAUUUGUCGGACUAGUCGCGUCGAGUCAUUCGGUUUACCCUAGUGACUAGACCGAACACUCGCGUCCUAAACGAAC GAC
>aga-mir-309 MI0010607
AAGAUGCACAAACUCCGUCCAGAGGGUGUCCGAUUUCAAGAACUCAUCUGGGCAAAGUUUGUCGCAUAAA
>aga-mir-927 MI0010596
UAGUUAUGGUUUUUUAGAAUUCUACGCUUUACCCGUGAUUAAAGUAGUGCGGCAAAGCGUUUGGAUUCUGAAACGAAACGUUACAACG
>aga-mir-929 MI0010595
CCUACGCGUGGGAUUAAAUUGACUCUAGUAGGGAGUCCUUUUUACGAGCGACUCCCUAACGGAGUCAGAUUGAUUCCGGUAGUGCACC
>aga-mir-957 MI0010600
UGAUCACUGCGUGCGUUAGUUUUGGGCGGUUUUAGUGUAUUUCGAUGAGAAUUCUAUUGAAACCGUCCAAAACUGAGGCUGGCAGAUUGGUUA C
>aga-mir-965-1 MI0010602
CGAAGUAAAUCGACGAAUUCUGAAAGGGAAUGCUGUACACUCUAUGUGCUAAAAUAUCCAUAAGCGUAUAGCUUUUCCAUUUAGCGUUCG UUGCGAGCAA
>aga-mir-970 MI0010599
AUGGCUGACGGAAGAUAGCCAGCGUUUGUUUUUUUGGUAUGUUACUACUGCUGCAAGCUAUUAUCAUAAGACACACGCGCUUUUUCCGCAA CCCGA
>aga-mir-988 MI0010598
CUGUGGAGCAUUGGACAUUGCCGGUGUGUUUUUGUGACAAUGAGAUUUAAUAGUUUAAGUUAUCCCUUGUUGCAAACCUCACGCUGGCGG UGUUUCCGA
>aga-mir-993 MI0010597
UGCUCUCCUGGACCUACCCUGUAGUUCGGGCUUUUGUGGUUGAAUAUUUAUCAUAUAUAUCUCAUAUACGCUUAUCAGAAGCUCGU UUCUAUAGAGGUAUCUCAGGGAGUGAA
>aga-mir-1000 MI0010604
CCUAGCAGUCGAUGAAUUGUCCUGUCACAGCAGUACUUUUUGCCUAGUUUACUGUUGUUUCGGGACAUUCCAUCGACGCUAGGGUUUUCAU
14 <i>hsa</i> pre-miRNAs:
>hsa-mir-1204 MI0006337
ACCUCUGGCGUGUCUCCAUAUUUUGAGAUGAGUUACAUCUUGGAGGUGAGGACGUGCCUCGUGGU
>hsa-mir-1972 MI0009982
UAUAGGCAUGUGCCACCACACUUGGCUUAAAUGUGUCAUUUAAAAUUCAGGCCAGGCACAGUGGCUCAUGCCUGUA

>hsa-mir-1973 MI0009983
UAUGUUAACGGCCAUGGUAUCCUGACCGUGCAAAGGUAGCAUA
>hsa-mir-1974 MI0009984
UGUUCUUGUAGUUGAAAUAACAACGAUGGUUUUUCAUAUCAUUGGUCGUGGUUGUAGUCCGUGCGAGAAUA
>hsa-mir-1975 MI0009985
AGUUGGUCGAGUGUUGGGUUUAUUGUUAAGUUGAUUUAAACAUUGUCUCCCCCACAACCGCGCUUGACUAGCU
>hsa-mir-1976 MI0009986
GCAGCAAGGAAGGCAGGGGUCCUAAGGUGUGUCCUCCUGCCCUCCUUGCUGU
>hsa-mir-1977 MI0009987
UUGAUUAGGGUGCUUAGCUGUUAAACUAAGUUGUUUGGGUUUAAGUCCAUUGGUCUAGUAAGGGCUUAGCUAAUUA
>hsa-mir-1978 MI0009988
UAGACGGGCUCACAUACCCCAUAAACAAAUAGGUUUGGUCCUAGCCUUUCUA
>hsa-mir-1979 MI0009989
UCUUUACUCCACUGCUUCACUUGACUAGCCUUUAAAAAAGAAAGGCUUGGUUGAUGAAUGGGUGAGAGAAAAGG
>hsa-mir-2052 MI0010486
CUGUUUUGAUAAACAGUAAUGUCCUUUAGUUAAGUUAACAGCUAUCAAAACAA
>hsa-mir-2053 MI0010487
CUUGCCAUGUAAAACAGAUUUAAUUAACAUAUUGCAACCUGUGAAGAUGCAAACUUUAA
>hsa-mir-2054 MI0010488
CUGUAAUAAAAUUAAUUUUAUUCUCAUCAUUAAAAAUGUAUUACAG
>hsa-mir-2110 MI0010629
CAGGGUUUGGGAAACGGCCGUGAGUGAGGCGUCGGCUGUUGUUCUACCGCGGUCUUUCCUCCACUCUUG
>hsa-mir-2113 MI0003939
UUUUCAAAGCAAUGUGUGACAGGUACAGGGACAAAUCCCGUUAAUAAGUAAGAGGAUUUGUCUUGGCUCUGUCACAUGCCACUUUGAAAA

Table S4: Results of miRenSVM on animal pre-miRNAs published in miRBase13.0.

Species (Animal)	Evaluated pre-miRNAs	Correctly Predicted	Accuracy (%)
<i>Apis mellifera</i>	64	60	93.75
<i>Branchiostoma floridae</i>	74	62	83.78
<i>Bombyx mori</i>	61	58	95.08
<i>Bos taurus</i>	359	341	94.99
<i>Capitella sp. I</i>	72	59	81.94
<i>Caenorhabditis briggsae</i>	98	89	90.82
<i>Caenorhabditis elegans</i>	155	148	95.48
<i>Canis familiaris</i>	325	317	97.54
<i>Ciona intestinalis</i>	25	25	100
<i>Ciona savignyi</i>	27	26	96.30
<i>Drosophila melanogaster</i>	152	138	90.79
<i>Drosophila pseudoobscura</i>	73	70	95.89
<i>Danio rerio</i>	337	331	98.22
<i>Fugu rubripes</i>	133	132	99.25
<i>Gallus gallus</i>	471	437	92.78
<i>Lottia gigantea</i>	57	46	80.70
<i>Monodelphis domestica</i>	119	115	96.64
<i>Macaca mulatta</i>	458	417	91.05
<i>Mus musculus</i>	568	521	91.73
<i>Nematostella vectensis</i>	49	43	87.76
<i>Ornithorhynchus anatinus</i>	344	307	89.24
<i>Pan troglodytes</i>	599	528	88.15
<i>Rattus norvegicus</i>	286	271	94.76
<i>Schmidtea mediterranea</i>	79	75	94.94
<i>Sus scrofa</i>	55	54	98.18
<i>Tribolium castaneum</i>	55	53	96.36
<i>Tetraodon nigroviridis</i>	143	140	97.90
Total	5238	4863	92.84

Features:

Triplet element (32):

Triplet structure-sequence element is proposed by Xue *et.al* [5]. The detail of these features has already been well described in the main article.

Base pair feature (15):

11 of these features have been used in previous miRNA gene predicting methods [6, 7]

Four new features relevant to loop number in the predicted secondary structure are introduced:

- dP/n_loops , where n_loops is the number of loops in secondary structure.
- $\%(A-U)/n_loops$, $\%(G-C)/n_loops$, $\%(G-U)/n_loops$, where $\%(X-Y)$ is the ratio of X-Y base pairs in the secondary structure.

Thermodynamics features (18):

6 MFE related features; 8 other global thermodynamics features and 4 statistically significant features are chosen from previous research [6, 7, 8].

The definitions of other features already used by existing pre-miRNA classification methods are available in [7]'s supplementary data.

References:

1. Jiang P, Wu H, Wang W, et al.: **MiPred: classification of real and pseudo microRNA precursors using random forest prediction model with combined features.** *Nucleic acids research* 2007, **35**:W339-44.
2. Griffiths-Jones S, Saini HK, van Dongen S, Enright AJ: **miRBase: tools for microRNA genomics.** *Nucleic acids research* 2008, **36**:D154-8.
3. Hofacker IL: **Vienna RNA secondary structure server.** *Nucleic acids research* 2003, **31**:3429-31.
4. Li W, Godzik A: **Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide sequences.** *Bioinformatics* 2006, **22**:1658-9.
5. Xue C, Li F, He T, et al.: **Classification of real and pseudo microRNA precursors using local structure-sequence features and support vector machine.** *BMC bioinformatics* 2005, **6**:310.
6. Ng KL, Mishra SK: **De novo SVM classification of precursor microRNAs from genomic pseudo hairpins using global and intrinsic folding measures.** *Bioinformatics* 2007, **23**:1321-30.
7. Batuwita R, Palade V: **microPred: effective classification of pre-miRNAs for human miRNA gene prediction.** *Bioinformatics* 2009, **25**:989-95.
8. Freyhult E, Gardner PP, Moulton V: **A comparison of RNA folding measures.** *BMC bioinformatics* 2005, **6**:241.