

Использование регрессионного анализа для создания эффективных малых РНК, играющих важную роль в регуляции работы генов у животных и растений

Матвеева Ольга, Назипова Нафиса, Огурцов Алексей, Шабалина Светлана

Институт математических проблем биологии РАН, г. Пущино, Россия
НЦБИ, Национальная медицинская библиотека, Национальные институты здравоохранения, г. Бетесда, США
Университет штата Юта, г. Солт-Лейк-Сити, США

Frontiers in Genetics, 2012, Vol. 3, Article No. 163

РНК-интерференция (англ. *RNA interference*, *RNAi*) — процесс подавления экспрессии гена или деградаци мРНК при помощи малых молекул РНК. Процессы РНК-интерференции обнаружены в клетках многих эукариот: у животных, растений и грибов. Система РНК-интерференции играет важную роль в защите клеток от вирусов, паразитирующих генов (транспозонов), а также в регуляции развития, дифференцировки и экспрессии генов организма.

Открытие феномена РНК-интерференции (RNAi)



...

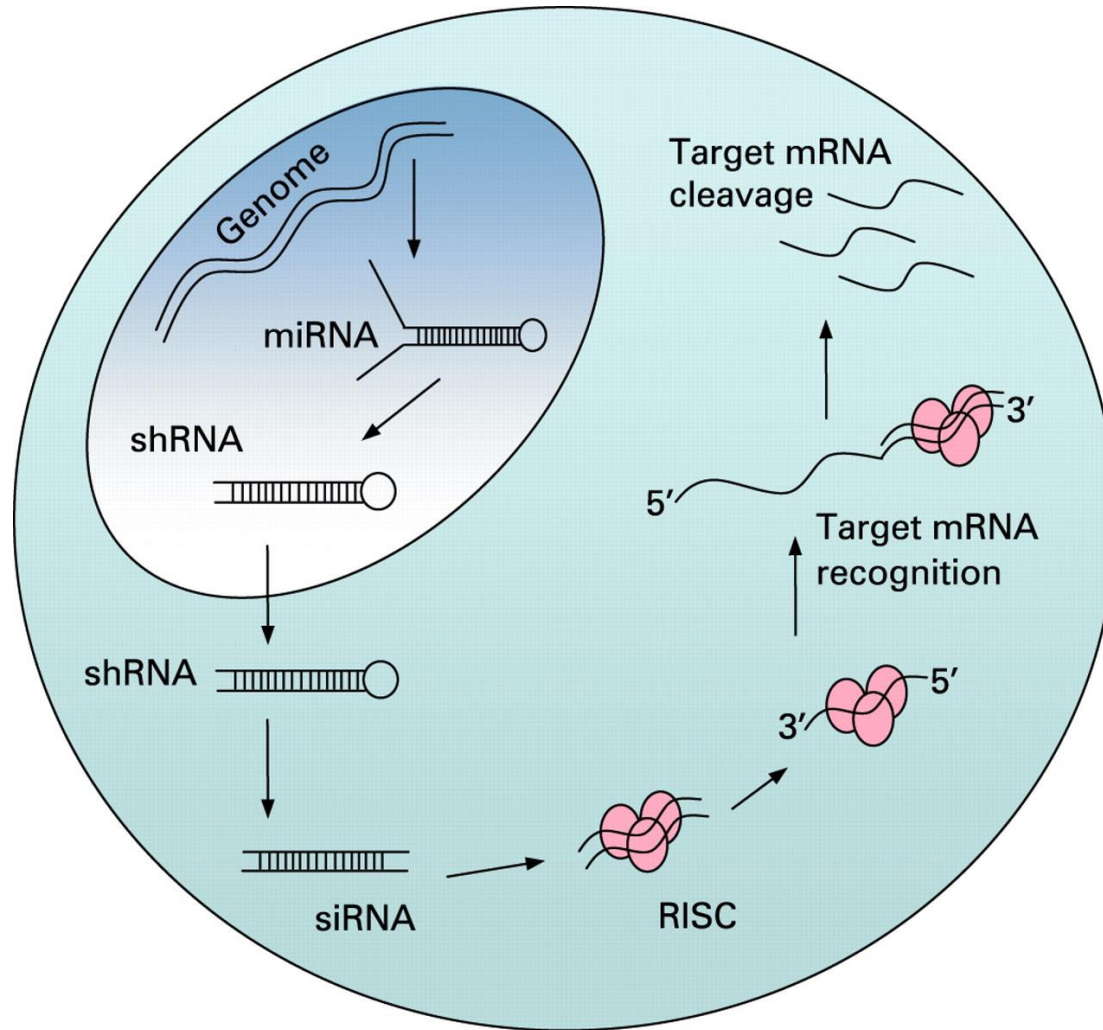
1990

2000

2010




Механизм РНК-интерференции



Lord C J et al. J Clin Pathol 2009;62:195-200

>hsa-mir-30a MI0000088

GCGACUGUAAACAUC CUCGACUGGAAGCUGUGAAGCCACAGAUGGGCUUUCAGUCGGAUGUUUGC
AGCUGC

Stem-loop	<pre> a uc ----- a gcg cuguaaacauc c gacuggaagcu gug a cgu gacguuuuguagg cugacuuucgg cac g c -- guaga c </pre> <p>Get sequence</p>
Deep sequencing	<p>853211 reads, 1.07e+04 reads per million, 79 experiments</p> 
Confidence	<p>Annotation confidence: high Feedback: Do you believe this miRNA is real? <input type="button" value="Yes (+12)"/> <input type="button" value="No (-0)"/> <input type="button" value="Leave comment"/></p>
Comments	<p>The mature sequences miR-30 [1] and miR-97 [2] appear to originate from the same precursor. Subsequent data confirm that both arms of the precursor appear to give rise to mature miRNA sequences (Pfeffer S, pers. comm.). Landgraf et al. later showed that the 5' product is the predominant one [5]. Related miRNAs are processed from the 5' arms of other precursor loci (mir-30b, MI0000441; mir-30c-1, MI0000736; mir-30c-2, MI0000254; mir-30d, MI0000255; mir-30e, MI0000749).</p>
Genome context	<p><i>Coordinates (GRCh38)</i> chr6: 71403551-71403621 [-]</p>
Database links	<p>EMBL: AF480535; AF480535 EMBL: AF480569; AF480569 EMBL: AJ421752; HSA421752 ENTREZGENE: 407029; MIR30A HGNC: 31624; MIR30A</p>

Mature sequence hsa-miR-30a-5p

Accession	MIMAT0000087
Previous IDs	hsa-miR-30a-5p;hsa-miR-30a
Sequence	<p>6 - uguaaacauc c ucgacuggaag - 27</p> <p>Get sequence</p>
Deep sequencing	3520 reads, 61 experiments

Подготовка исходных данных

Исходные данные для анализа - 18719 последовательностей длиной 22 н.п. с известными значениями эффективности сайленсинга (log). Эти данные делились на 2 части: 14972 последовательности для обучения и 3747 для тестирования.

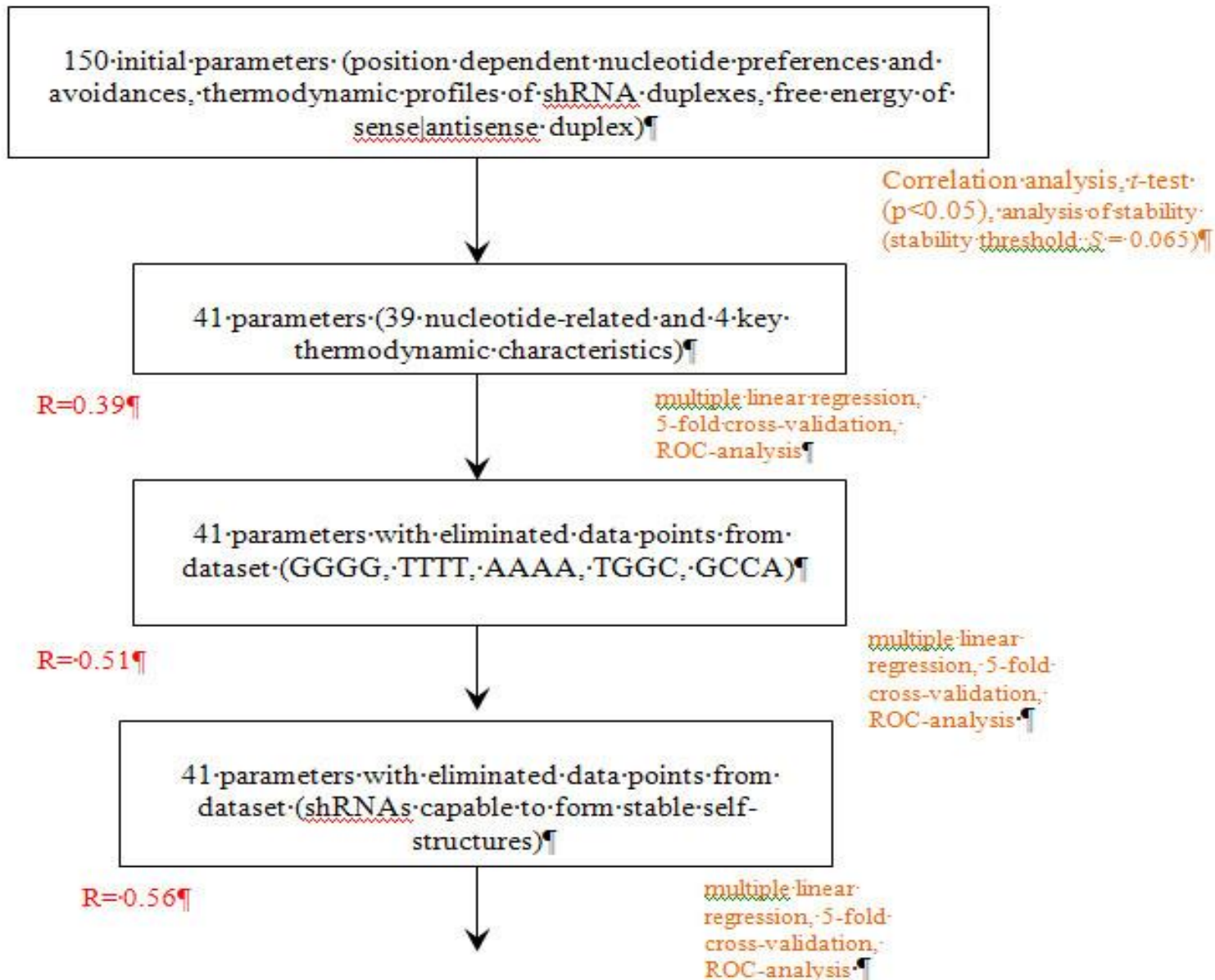
A \Rightarrow 1,0,0,0; T \Rightarrow 0,1,0,0; C \Rightarrow 0,0,1,0; G \Rightarrow 0,0,0,1.

Около 150 гетерогенных параметров из литературы, нуклеотидного контекста, термодинамических расчетов (OligoHybrid, OligoTherm, Afold). Для выбора наиболее универсальных параметров использовалось 2 критерия:

- значимая корреляция с эффективностью для каждого параметра,
- стабильность корреляции для каждого параметра.

Оба критерия оценивались только для обучающей выборки.

Блок-схема процедуры подбора параметров модели



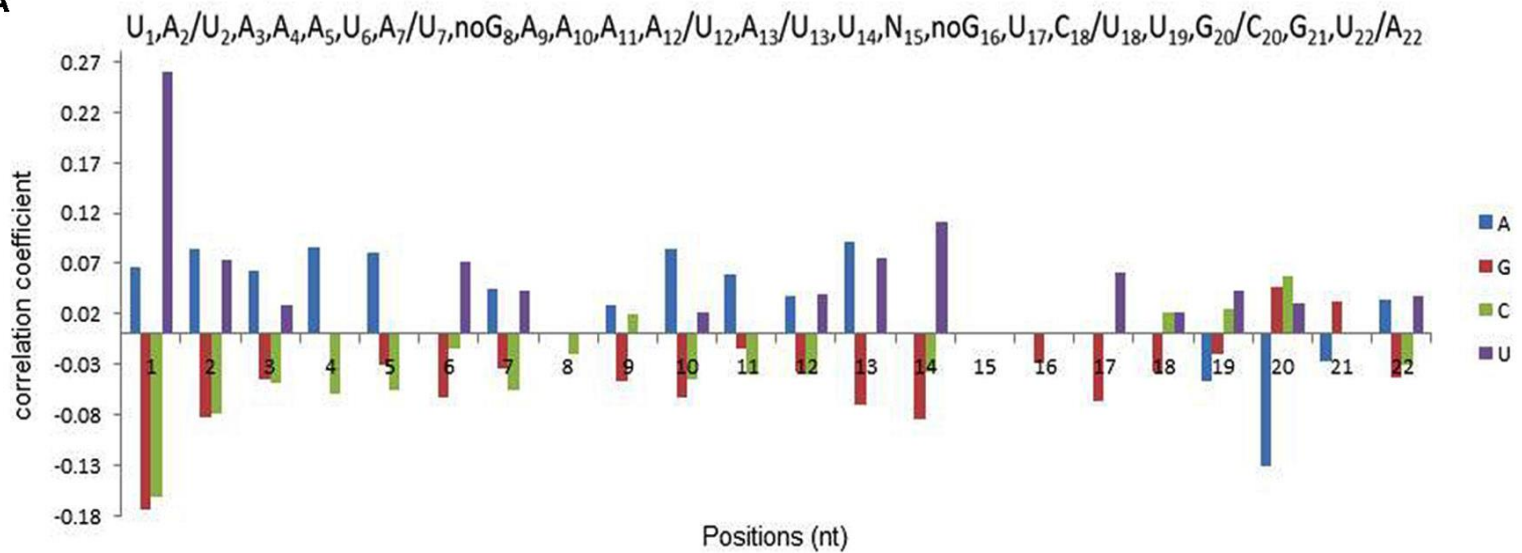
Оценка стабильности корреляции

Оценка стабильности корреляции является индикатором того, насколько предсказательная сила конкретного параметра зависит от выбора подмножества обучающего множества:

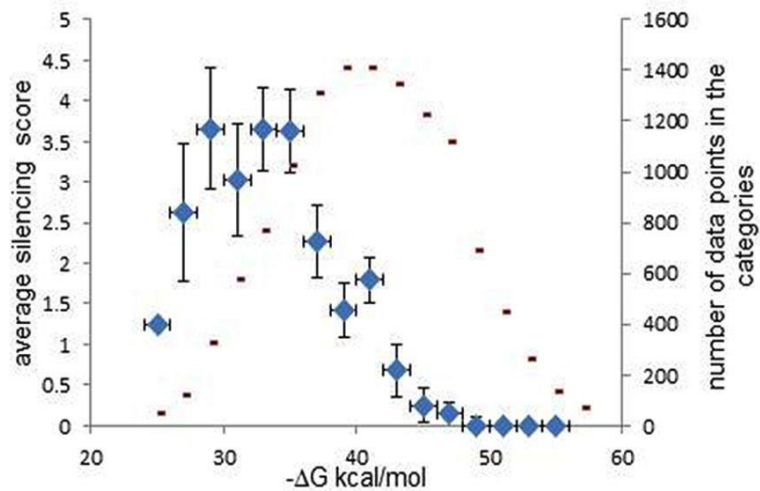
$$S_n = \sqrt{\frac{\sum_{i=1}^m (R_{n_i} - \bar{R}_n)^2}{m-1}}, \quad \bar{R}_n = \frac{\sum_{i=1}^m R_{n_i}}{m}, \quad n = 1, \dots, 1000; \quad m = 10$$

Результаты изучения взаимосвязи эффективности и структурных особенностей shRNA

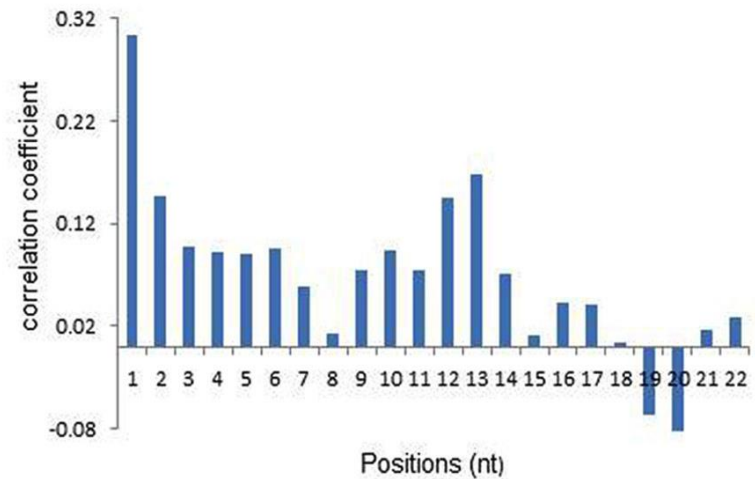
A



B



C



Результат первого
этапа уточнения
модели.
 $\Delta G^2, \Delta G^3$ получены
программой
Afold.

Parameter	Coefficient	P-value
Intercept	-43.85	2.046×10^{-178}
ΔG	-2.694	4.459×10^{-111}
ΔG^2	-0.06917	1.716×10^{-105}
ΔG^3	-0.0005959	3.017×10^{-106}
ΔG_1	0.3705	0.00048315
A ₁	0.8201	2.61×10^{-06}
A ₂	0.8225	1.4939×10^{-11}
A ₃	0.6747	0
A ₄	0.4644	0
A ₅	0.5027	0
A ₁₁	0.3433	1.0172×10^{-10}
A ₁₂	0.4986	0
A ₁₃	0.9504	0
A ₁₉	-0.4455	2.1798×10^{-17}
A ₂₀	-1.067	2.6228×10^{-98}
A ₂₂	0.3039	2.0603×10^{-08}
G ₆	-0.3063	1.8561×10^{-08}
G ₇	-0.3133	2.2034×10^{-08}
G ₉	-0.2439	2.8698×10^{-06}
G ₁₀	-0.6907	7.4738×10^{-28}
G ₁₄	-0.4776	2.4605×10^{-18}
G ₁₇	-0.2477	4.2719×10^{-06}
G ₁₈	-0.2378	3.9472×10^{-06}
G ₁₉	-0.132	0.0140694
G ₂₁	0.3678	1.1793×10^{-12}
C ₄	-0.239	9.663×10^{-06}
C ₅	-0.3201	3.7375×10^{-09}
C ₇	-0.4681	3.9052×10^{-17}
C ₈	-0.1913	0.00023711
C ₁₀	-0.6128	2.4693×10^{-22}
C ₁₁	-0.2256	3.1081×10^{-05}
C ₁₃	-0.2326	0.00018904
U ₁	1.989	0
U ₂	0.6581	4.3741×10^{-07}
U ₃	0.5757	0
U ₆	0.5343	0
U ₁₀	-0.3281	3.3946×10^{-08}
U ₁₂	0.5828	0
U ₁₃	0.7992	0
U ₁₄	0.6866	0
U ₁₇	0.3141	1.7814×10^{-09}
U ₂₂	0.3012	2.7612×10^{-08}

Множественный коэффициент корреляции

$$R = \sqrt{1 - \frac{\sum_i (x_i - y_i)^2}{\sum_i (x_i - \bar{x})^2}} \quad , \text{ где}$$

x_i - экспериментальные эффективности олигов,

y_i - предсказанные с помощью модели эффективности олигов,

$\bar{x} = \frac{\sum_i x_i}{n}$ - среднее экспериментальных эффективностей,

n - объем выборки,

R^2 - коэффициент детерминации.

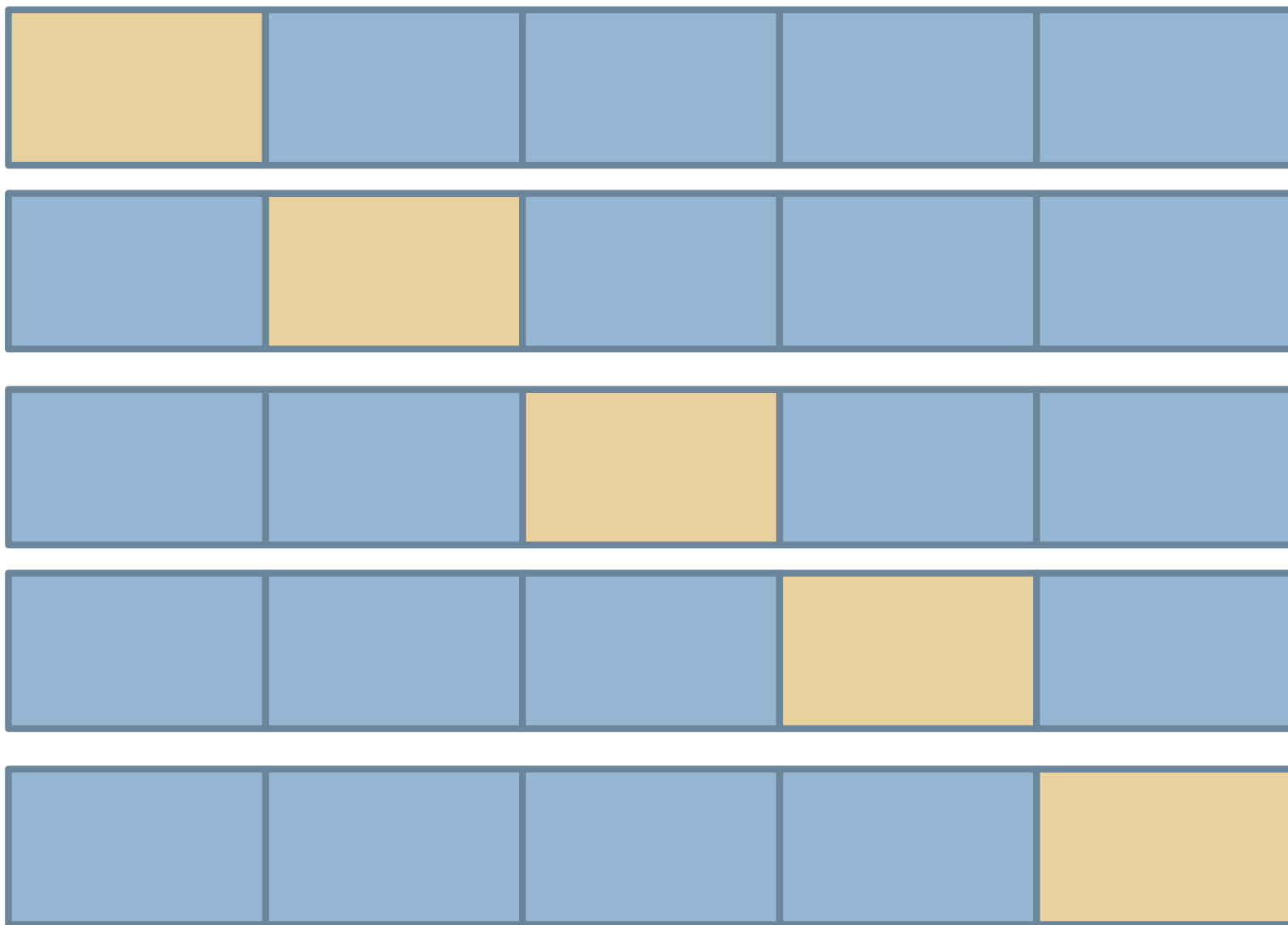
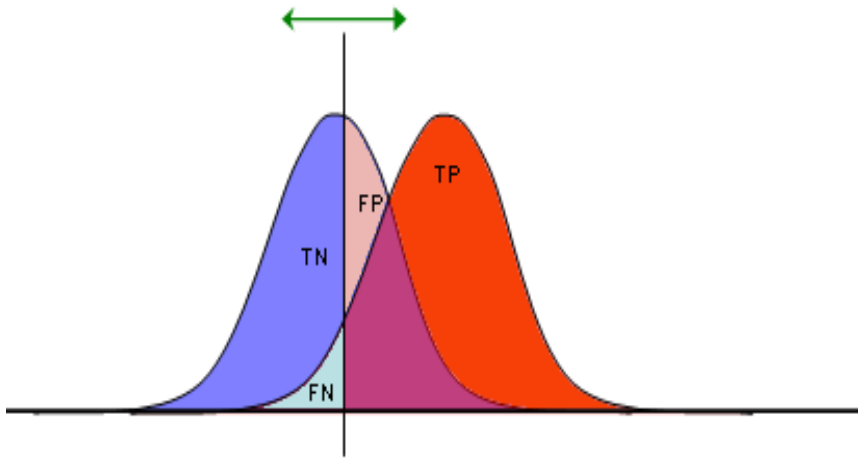
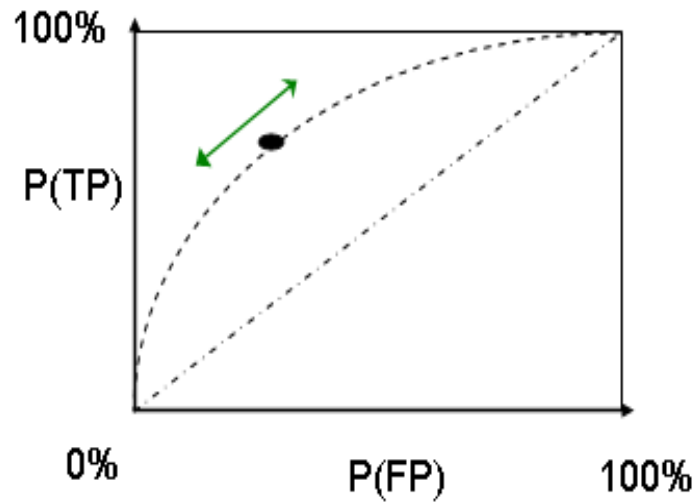


Схема пятикратной перекрестной проверки модели

ROC KURBUNA



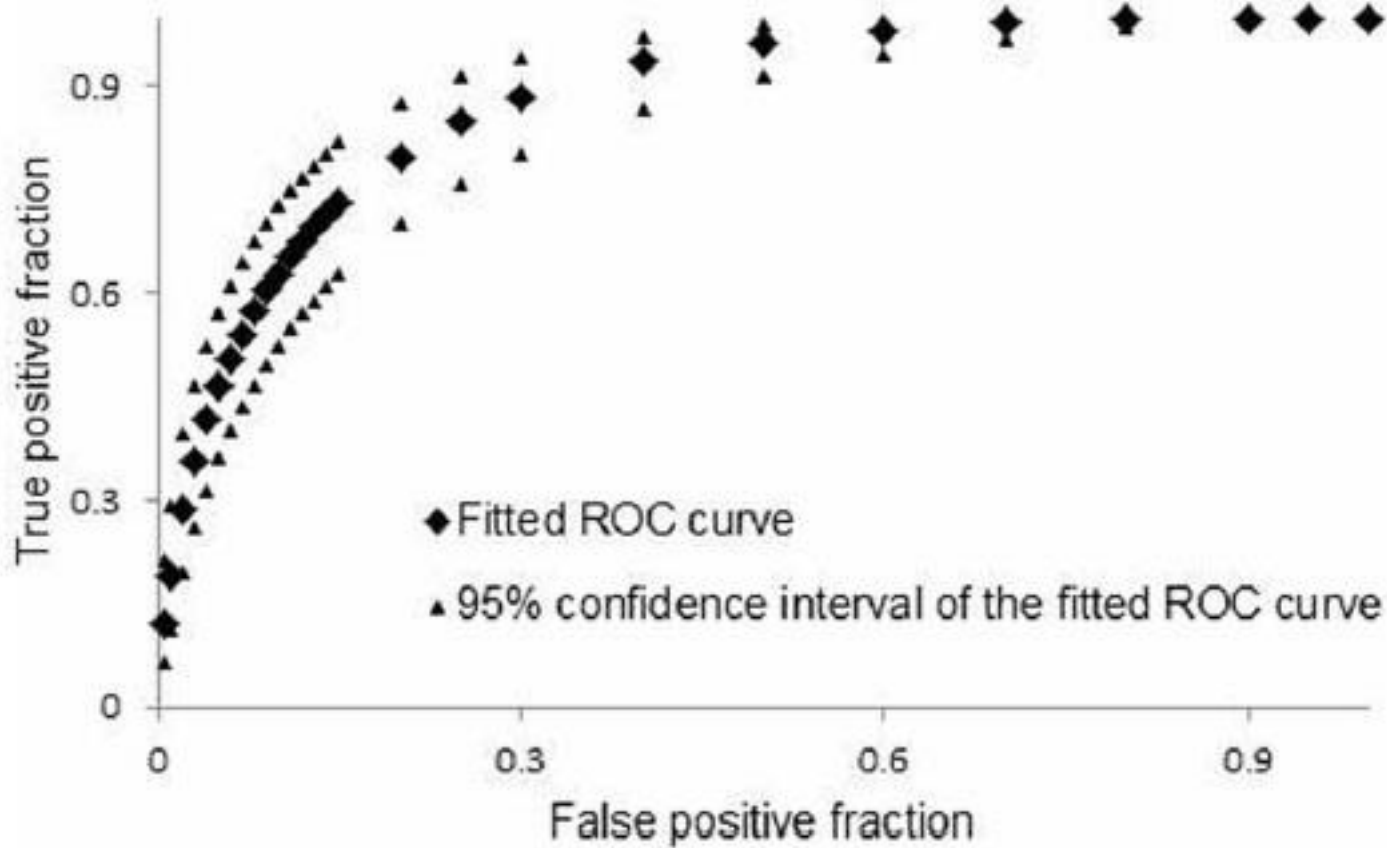
TP	FP
FN	TN
1	1



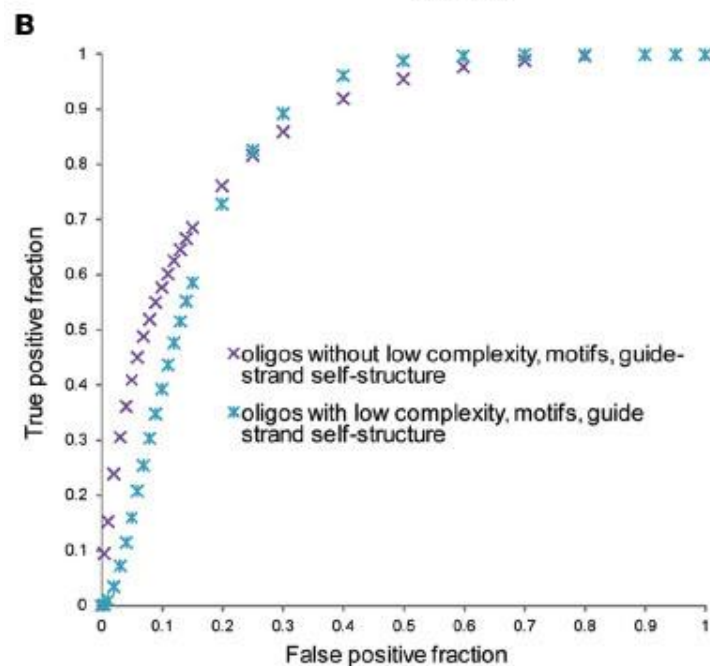
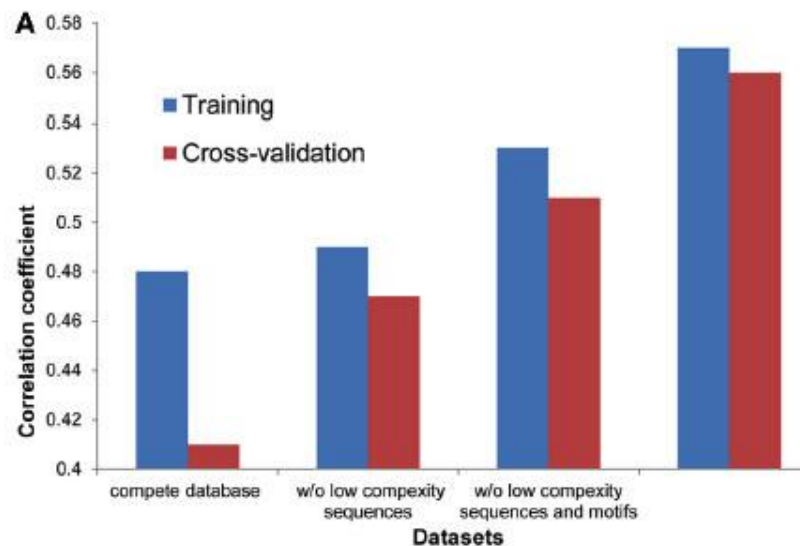
$$Sensitivity = \frac{TP}{TP + FN},$$

$$Specificity = \frac{TN}{TN + FP}$$

ROC-кривая для результатов предсказания эффективности сайленсинга по 41 параметру (0.882)



Результат улучшения предсказательной силы модели после удаления из набора данных простых последовательностей



miR_Scan

www.ncbi.nlm.nih.gov/staff/ogurtsov/projects/mi30/

Работа частично поддержана средствами гранта РФФИ
№12-07-00530а