

# Вероятностный подход для задачи предсказания биологической активности ядерных рецепторов

Володин Сергей Евгеньевич

Московский физико-технический институт

Курс: Численные методы обучения по прецедентам  
(практика, В. В. Стрижов)/Группа 374, осень 2016

Предсказание взаимодействия двух типов молекул: лиганд с рецепторов. Необходимо оценить вероятность связывания и построить бинарный классификатор

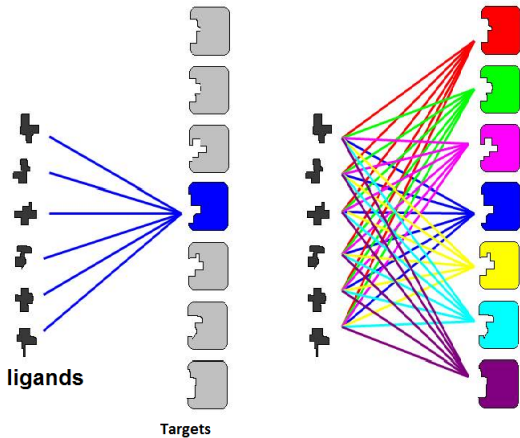
## Проблема

В ответах содержатся пропуски, поэтому независимое рассмотрение рецепторов дает неудовлетворительное качество классификации.

Оказывается, что события реакции лиганда с различными рецепторами не независимы.

## Задача

Необходимо построить вероятностную модель, учитывающую схожесть рецепторов.



- 1 Olexandr Isayev Sherif Farag Stephen J. Capuzzi, Regina Politi and Alexander Tropsha. Qsar modeling of tox21 challenge stress response and nuclear receptor signaling toxicity assays.
- 2 Geoff Holmes Eibe Frank Jesse Read, Bernhard Pfahringer. Classifier chains for multi-label classification.
- 3 Eyke H.0 Krzysztof Dembczynski, Weiwei Cheng. Bayes optimal multilabel classification via probabilistic classifier chains. 2010.

# Постановка задачи

Задана выборка  $\mathfrak{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\} = \mathfrak{L} \sqcup \mathfrak{T}$ .  $\mathbf{x}_i \in \mathbb{R}^n$ .  $\mathbf{y}_i \in \{0, 1, \square\}^I$ ,  
 $\square$  — пропуск в данных.

$\mathbf{X}, \mathbf{Y}$  — случайные величины, между классами есть зависимости.

Модель классификации: функция  $f: \mathbf{W} \times \mathbf{X} \times \mathbf{Y} \rightarrow [0, 1]$ ,

$$f(\mathbf{w}, \mathbf{x}, \mathbf{y}) = P(\mathbf{Y} = \mathbf{y} | \mathbf{X} = \mathbf{x}; \mathbf{w})$$

Функция потерь — логарифм правдоподобия

$$Q(f | \mathbf{w}, \mathcal{Z}) = - \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{Z}} \log f(\mathbf{w}, \mathbf{x}, \mathbf{y}) P(\mathbf{X} = \mathbf{x})$$

Требуется минимизировать  $Q$ :

$$\mathbf{w}^* = \arg \min_{\mathbf{w} \in \mathbf{W}} Q(f | \mathbf{w}, \mathcal{Z})$$

Для оценки конкретной модели используется AUC.



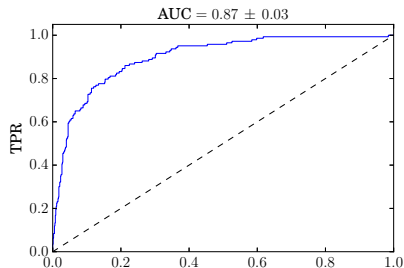
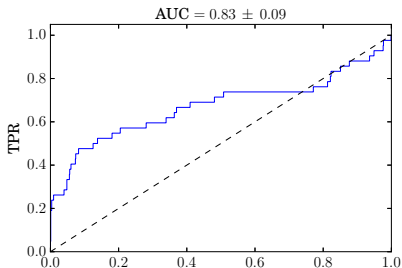
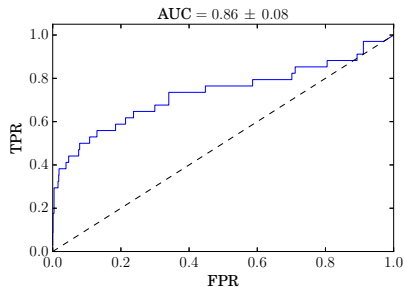
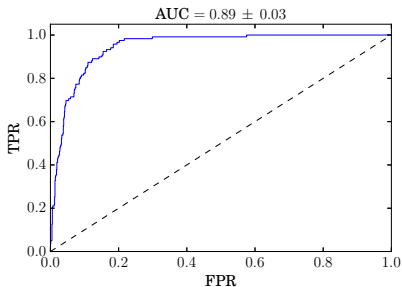
## Цель эксперимента

Сравнение различных моделей по критерию AUC для различных классов. Выбор гиперпараметров исходя из внешних требований к решению задачи.

## Сравниваемые модели

- 1 Binary Relevance
- 2 PCC — предлагаемое решение
- 3 Random Forest

# Вычислительный эксперимент





# Вычислительный эксперимент

## Результаты эксперимента

Рецептор	Binary Relevance	Random Forest	PCC
NR-AhR	<b>0.83</b> $\pm$ 0.03	<b>0.93</b>	
NR-AR-LBD	<b>0.86</b> $\pm$ 0.08	<b>0.88</b>	
NR-AR	<b>0.83</b> $\pm$ 0.09	<b>0.83</b>	
SR-MMP	<b>0.87</b> $\pm$ 0.03	<b>0.95</b>	
NR-ER	<b>0.78</b> $\pm$ 0.04	<b>0.81</b>	
SR-HSE	<b>0.79</b> $\pm$ 0.04	<b>0.86</b>	
SR-p53	<b>0.79</b> $\pm$ 0.07	<b>0.88</b>	
NR-PPAR-gamma	<b>0.79</b> $\pm$ 0.04	<b>0.86</b>	
SR-ARE	<b>0.78</b> $\pm$ 0.02	<b>0.84</b>	
NR-Aromatase	<b>0.81</b> $\pm$ 0.05	<b>0.84</b>	
SR-ATAD5	<b>0.81</b> $\pm$ 0.06	<b>0.83</b>	
NR-ER-LBD	<b>0.80</b> $\pm$ 0.07	<b>0.83</b>	

- 1 Предложена модель для предсказания взаимодействия, учитывающая зависимости между классами
- 2 Проведено сравнение модели с другими по критерию AUC
- 3 Модель BR имеет худшие показатели AUC, чем Random Forest
- 4 PCC лучше BR для классов a, b, c по критерию AUC.