

Формулировка и решение задачи оптимизации, сочетающей классификацию и регрессию, для оценки энергии связывания белка и маленьких молекул

Илья Игашов, Сергей Грудинин, Мария Кадукова, и В.В. Стрижов

В работе рассматривается задача оценки свободной энергии связывания белка с лигандом посредством оптимизации скоринговой функции Convex-PL. Оптимизация сочетает в себе классификацию, основанную на методе опорных векторов, и регрессию, использующую различные функции потерь. Использование одной лишь задачи классификации для предсказания энергии связывания приводит к недостаточно высокой корреляции предсказаний с экспериментальными значениями, в то время как использование только регрессии приводит к переобучению. В этой работе предлагается посторить алгоритм, который объединяет классификацию и регрессию, решает данные проблемы и демонстрирует высокое качество оценки энергии связывания. Для проверки работы алгоритма будут использоваться данные, состоящие из комплексов белков и лигандов, для которых необходимо найти наилучшую позу лиганда или предсказать свободную энергию связывания. Также ниже будут представлены результаты работы алгоритма на датасете, состоящем из белков, для которых нужно найти наиболее сильно связывающий лиганд. Результаты будут получены как на выборках пониженной размерности, так и на полных выборках с большой размерностью.

I. ВВЕДЕНИЕ

Развитие вычислительных методов и появление новых подходов в молекулярном моделировании дает широкие возможности в области изучения химических соединений. В частности, популярный ныне метод виртуального скрининга применяется при разработке новых лекарственных препаратов для поиска и анализа химических соединений, обладающих необходимыми биологическими свойствами¹.

Среди методов молекулярного моделирования широко распространен молекулярный докинг, позволяющий предсказать взаимную ориентацию молекул, наиболее выгодную для образования устойчивого комплекса². Поскольку для решения данной задачи требуется анализ и обработка огромного количества данных, содержащих в себе информацию о химических комплексах и их характеристиках, она является вычислительно трудоемкой и требует решений, обладающих высокой производительностью.

Образование комплекса «белок-лиганд» можно рассматривать как термодинамическое событие, описываемое постоянной степени сродства (аффинности связывания) соединения, которая напрямую сопряжена со свободной энергией связывания белка с лигандом. Нативные конформации отвечают минимуму энергии связывания. Задача молекулярного докинга — предсказать эти взаимные расположения и отвечающие им энергетические значения. Свободная энергия связывания зависит от множества факторов, включающих в себя не только взаимодействие белка с лигандом, но также сольватацию и энтропийные факторы. Строгий подсчет значений энергии связывания потребовал бы семплирования всего конфигурационного пространства, что является, с точки зрения вычислений, крайне затратной задачей ввиду высокой размерности пространства³. В последние годы для ре-

шения данной задачи был предложен ряд аппроксимирующих алгоритмов, оценивающих значение энергии связывания на основе скоринговых функций⁴.

Одной из таких функций является Convex-PL³. Данная функция зависит от взаимного расположения белков и лигандов и достигает минимума на нативной конформации бинарной системы. Основная идея построения Convex-PL состоит в декомпозиции молекул белков и лигандов на отдельные атомы из некоторого конечного набора и последующем подсчете значения функционала на всевозможных комбинациях различных пар атомов, учитывая априорные распределения плотностей расстояний между ними. Как показано³, данную функцию можно приблизить полиномом конечной степени с так называемыми «структурными» коэффициентами, характеризующими взаимное расположение белков и лигандов, и найти эти коэффициенты разложения методами выпуклой оптимизации⁵.

Для решения описанной выше оптимизационной задачи были предложены два подхода. В работе³ описывается алгоритм классификации, основанный на методе опорных векторов⁶. Данный алгоритм был протестирован с помощью D3R Grand Challenge⁷ и CASF 2013. Полученные на CASF 2013 результаты по предсказанию нативных и около-нативных конформаций превзошли показатели других 20 методов, протестированных на тех же данных ранее. В работе⁸ был предложен регрессионный алгоритм, использующий linear ridge regression (LRR) и kernel ridge regression (KRR)^{9–11}. Данный метод был протестирован на датасетах MAP4K4 и HSP90 и продемонстрировал хорошие результаты⁸. Однако, как показали исследования^{3,8}, оба подхода также имеют и слабые стороны: использование одной лишь задачи классификации для предсказания энергии связывания приводит к недостаточно высокой корреляции предсказаний с экспериментальными значениями, в то время как использование только регрессии приводит к пере-

обучению. Кроме того, исследования, представленные в^{3,8}, были проведены на данных пониженной размерности. В этой работе представлен алгоритм молекулярного докинга для предсказания нативных конформаций комплексов «белок-лиганд» и значений энергии связывания данных соединений. Описанный ниже алгоритм решает задачу поиска минимума скоринговой функции Convex-PL методами выпуклой оптимизации, сочетая в себе задачи классификации и регрессии.

II. МОДЕЛЬ ВЗАИМОДЕЙСТВИЙ

Пусть имеется P нативных комплексов белков-лигандов $\{C_{i0}\}_{i=1}^P$. Применяя к лигандам изометрические преобразования, сгенерируем для каждой нативной позы D ненативных поз $\{C_{ij}\}_{j=1}^D$. Таким образом, для каждого из P комплексов имеем $(D+1)$ конформаций: одну нативную и D ненативных. Требуется найти скоринговый функционал E , который удовлетворяет следующим неравенствам:

$$E(C_{i0}) < E(C_{ij}) \quad \forall i \in \{1, \dots, P\}, \quad \forall j \in \{1, \dots, D\}. \quad (1)$$

Будем рассматривать свободную энергию связывания как некоторый функционал E , определенный для всех возможных конфигураций белков и лигандов. Для упрощения формы функционала сделаем ряд допущений. Во-первых, будем считать, что E зависит только от взаимодействий белка и лиганда в комплексе. В данном случае под взаимодействиями понимается набор пар атомов, таких, что в каждой паре первый атом является атомом лиганда, а второй – атомом белка. Кроме того, будем рассматривать только те пары, в которых расстояние между атомами не превышает некоторой пороговой величины r_{\max} . В качестве r_{\max} возьмем значение 10\AA , как это было сделано в других работах^{12–17} ранее. Во-вторых, будем рассматривать комплекс "белок-лиганд" как набор атомов, каждый из которых имеет некоторый тип. Тип каждого атома зависит от его химических свойств, таких как номер элемента в периодической таблице, аромат, гибридизация, полярность и т.д. Пусть M_1 – количество типов атомов лиганда, а M_2 – количество типов атомов белка. Тогда получим всего $M_1 \times M_2$ различных атомных взаимодействий. В-третьих, будем считать, что E зависит только от распределения расстояний между взаимодействующими атомами. И, наконец, предположим, что E является линейным функционалом и имеет вид:

$$E(n(r)) = \sum_{k=1}^{M_1} \sum_{l=1}^{M_2} \int_0^{r_{\max}} n^{kl}(r) f^{kl}(r) dr, \quad (2)$$

где $f^{kl}(r)$ – неизвестные функции взаимодействия между атомами типов k и l . Будем называть их *скоринговыми потенциалами*. Функции $n^{kl}(r)$ – численные

плотности распределений пар атомов типов k и l по расстоянию r между ними:

$$n^{kl}(r) = \sum_{i,j} \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{(r - r_{ij})^2}{2\sigma^2} \right], \quad (3)$$

где σ^2 – стандартное отклонение (константа). Сумма берется по всем парам (i, j) атомов с типами k и l соответственно, у которых расстояние между атомами не превышает пороговой величины r_{\max} , атом i принадлежит лиганду, а атом j – белку.

Разложим неизвестные скоринговые потенциалы $f^{kl}(r)$ и плотности $n^{kl}(r)$ по полиномиальному базису:

$$\begin{aligned} f^{kl}(r) &= \sum_q w_q^{kl} \psi_q(r), \\ n^{kl}(r) &= \sum_q x_q^{kl} \psi_q(r), \end{aligned} \quad (4)$$

где $\psi_q(r)$ – ортогональные базисные функции на интервале $[0, r_{\max}]$, а w_q^{kl} и x_q^{kl} – коэффициенты разложения функций $f^{kl}(r)$ и $n^{kl}(r)$ соответственно.

Тогда функция энергии связывания E приблизительно выражается формулой

$$\begin{aligned} E(n(r)) &\approx \sum_{k=1}^{M_1} \sum_{l=1}^{M_2} \sum_{q=0}^Q w_q^{kl} x_q^{kl} = \langle \mathbf{w}, \mathbf{x} \rangle, \\ \mathbf{w}, \mathbf{x} &\in \mathbb{R}^{Q \times M_1 \times M_2}. \end{aligned} \quad (5)$$

Неизвестный вектор \mathbf{w} будем называть *скоринговым вектором*, а вектор \mathbf{x} – *структурным вектором*. Структурный вектор известен для решения задачи: его можно получить из данных. Для оценки энергии связывания в работе⁸ использовался порядок разложения $Q = 10$, и рассматривались $M_1 = 48$ типов атомов.

III. ПОСТАНОВКА ЗАДАЧИ

Пусть имеется P нативных соединений белков и лигандов, где i -я конформация описывается структурным вектором $\mathbf{x}_{i0}^{\text{nat}}$, и $(P \times D)$ ненативных конформаций с соответствующими структурными векторами $\mathbf{x}_{ij}^{\text{nonnat}}$. Учитывая (5) и тот факт, что минимум энергии связывания соответствует нативным конформациям соединений белков и лигандов, справедливы следующие неравенства:

$$\begin{aligned} \langle \mathbf{x}_{i0}^{\text{nat}}, \mathbf{w} \rangle &< \langle \mathbf{x}_{ij}^{\text{nonnat}}, \mathbf{w} \rangle, \\ \langle \mathbf{x}_{ij}^{\text{nonnat}} - \mathbf{x}_{i0}^{\text{nat}}, \mathbf{w} \rangle &> 0, \\ i &= 1, \dots, P, \\ j &= 1, \dots, D. \end{aligned} \quad (6)$$

В такой постановке задачи требуется отыскать неизвестный скоринговый вектор \mathbf{w} , который бы решал

задачу классификации, а именно, определял бы, является ли конформация белка и лиганда со структурным вектором \mathbf{x} нативной или ненативной.

Система неравенств (6) может иметь ноль, одно или бесконечное число решений^{3,6}. Чтобы получить единственное решение задачи (6), переформулируем ее в виде задачи квадратичной оптимизации с мягким зазором⁵:

$$\begin{aligned} \text{minimize:} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{ij} C_{ij} \xi_{ij} \\ \text{subject to:} \quad & y_{ij} [\langle \mathbf{w}, \mathbf{x}_{ij} \rangle - b_i] - 1 + \xi_{ij} \geq 0, \\ & \xi_{ij} \geq 0, \\ & i \in \{1, \dots, P\}, \\ & j \in \{0, \dots, D\}, \end{aligned} \quad (7)$$

где скоринговый вектор \mathbf{w} , вектор смещений b_i и переменные невязки ξ_{ij} – оптимизируемые параметры модели, y_{ij} – класс j -й позы i -го соединения ($y_{i0} = 1$ для нативной позы и $y_{ij} = -1$, $j \in \{1, \dots, D\}$, для ненативной позы), C_{ij} – некоторые коэффициенты регуляризации. Таким образом, можно получить классификатор нативных и ненативных поз соединения, решив оптимизационную задачу (7).

Кроме того, для предсказания свободной энергии связывания белка с лигандом можно решать задачу регрессии⁸:

$$\begin{aligned} \text{minimize:} \quad & \sum_i [\langle \mathbf{w}, \mathbf{x}_{i0} \rangle - s_i]^2 + \alpha \|\mathbf{w}\|^2, \\ & i \in \{1, \dots, P\}, \end{aligned} \quad (8)$$

где s_i – истинное (экспериментально полученное) значение энергии связывания i -го нативного соединения, α – некоторый положительный коэффициент регуляризации для ridge-регрессии.

Оба описанных выше подхода в конечном итоге решают одну и ту же задачу поиска неизвестного скорингового вектора \mathbf{w} . Объединив эти методы, мы также будем решать задачу оптимизации модели, предсказывающей энергию связывания белков с лигандами:

$$\begin{aligned} \text{minimize:} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{ij} C_{ij} \xi_{ij} + C_r \sum_i f(\mathbf{x}_{i0}, \mathbf{w}, s_i) \\ \text{subject to:} \quad & y_{ij} [\langle \mathbf{w}, \mathbf{x}_{ij} \rangle - b_i] - 1 + \xi_{ij} \geq 0, \\ & \xi_{ij} \geq 0, \\ & i \in \{1, \dots, P\}, \\ & j \in \{0, \dots, D\}, \end{aligned} \quad (9)$$

где $f(\mathbf{x}_{i0}, \mathbf{w}, s_i)$ – функция потерь регрессии (Mean Squared Error), C_r – коэффициент регуляризации для функции потерь регрессии.

Поднимем размерность пространства векторов \mathbf{w} и

\mathbf{x} , избавившись от вектора смещения b_i :

$$\begin{aligned} \hat{\mathbf{w}} &= (\mathbf{w}, b_1, \dots, b_P)^T, \\ \hat{\mathbf{x}}_{1j} &= (\mathbf{w}, -1, 0, \dots, 0)^T, \\ \hat{\mathbf{x}}_{2j} &= (\mathbf{w}, 0, -1, 0, \dots, 0)^T, \\ &\dots \\ \hat{\mathbf{x}}_{Pj} &= (\mathbf{w}, 0, \dots, 0, -1)^T, \\ \hat{\mathbf{w}}, \hat{\mathbf{x}}_{ij} &\in \mathbb{R}^{(Q \times M_1 \times M_2) + P}, \\ j &\in \{0, \dots, D\}. \end{aligned} \quad (10)$$

Тогда задача оптимизации (9) принимает вид:

$$\begin{aligned} \text{minimize:} \quad & \frac{1}{2} \|\hat{\mathbf{w}}\|^2 + \sum_{ij} C_{ij} \xi_{ij} + C_r \sum_i f(\hat{\mathbf{x}}_{i0}, \hat{\mathbf{w}}, s_i) \\ \text{subject to:} \quad & y_{ij} \langle \hat{\mathbf{w}}, \hat{\mathbf{x}}_{ij} \rangle - 1 + \xi_{ij} \geq 0, \\ & \xi_{ij} \geq 0, \\ & i \in \{1, \dots, P\}, \\ & j \in \{0, \dots, D\}. \end{aligned} \quad (11)$$

IV. БАЗОВЫЙ ЭКСПЕРИМЕНТ

A. Уточнение задачи

Сформулированная в общем виде задача, которую предстоит решить, представлена формулой (11). В данном разделе приведен ряд уточнений формулировки задачи (11), которые позволят минимизировать трудности в реализации алгоритма и в запуске базового эксперимента.

Во-первых, заменим в задаче (11) функцию потерь классификации "hinge loss" :

$$V(\hat{\mathbf{w}}, \hat{\mathbf{x}}_{ij}, y_{ij}) = \max\{0, 1 - y_{ij} \hat{\mathbf{w}}^T \hat{\mathbf{x}}_{ij}\}$$

на функцию потерь "logistic loss" :

$$V(\hat{\mathbf{w}}, \hat{\mathbf{x}}_{ij}, y_{ij}) = \frac{1}{\ln 2} \ln [1 + \exp(-y_{ij} \hat{\mathbf{w}}^T \hat{\mathbf{x}}_{ij})].$$

Таким образом, вместо задачи квадратичного программирования (11) решаем задачу выпуклой оптимизации:

$$\begin{aligned} \text{minimize:} \quad & \frac{1}{2} \|\hat{\mathbf{w}}\|^2 + \sum_{ij} C_{ij} \ln [1 + \exp(-y_{ij} \hat{\mathbf{w}}^T \hat{\mathbf{x}}_{ij})] + \\ & + C_r \sum_i f(\hat{\mathbf{x}}_{i0}, \hat{\mathbf{w}}, s_i), \\ & i \in \{1, \dots, P\}, \\ & j \in \{0, \dots, D\}. \end{aligned} \quad (12)$$

Во-вторых, с помощью замены переменных сведем два квадратичных слагаемых в целевой функции из

задачи (12) к одному. Функция потерь регрессии – Mean Squared Error:

$$f(\hat{\mathbf{x}}_{i0}, \hat{\mathbf{w}}, s_i) = (\hat{\mathbf{w}}^T \hat{\mathbf{x}}_{i0} - s_i)^2. \quad (13)$$

Тогда для выборки $\hat{\mathbf{X}} = (\hat{\mathbf{x}}_{10}, \dots, \hat{\mathbf{x}}_{P0})^T$, состоящей из нативных конфигураций, и целевого вектора $\mathbf{s} = (s_1, \dots, s_P)^T$ квадратичные слагаемые целевой функции из задачи (11) принимают вид:

$$\begin{aligned} & \frac{1}{2} \|\hat{\mathbf{w}}\|^2 + C_r \sum_i f(\hat{\mathbf{x}}_{i0}, \hat{\mathbf{w}}, s_i) = \\ & \frac{1}{2} \hat{\mathbf{w}}^T \hat{\mathbf{w}} + C_r (\hat{\mathbf{w}}^T \hat{\mathbf{X}} - \mathbf{s})^2 = \\ & \frac{1}{2} \hat{\mathbf{w}}^T \hat{\mathbf{w}} + C_r \hat{\mathbf{w}}^T \hat{\mathbf{X}} \hat{\mathbf{X}}^T \hat{\mathbf{w}} - 2C_r \hat{\mathbf{w}}^T \hat{\mathbf{X}} \mathbf{s} + C_r \mathbf{s}^T \mathbf{s} = \\ & \hat{\mathbf{w}}^T \left(\frac{1}{2} \mathbf{I} + C_r \hat{\mathbf{X}} \hat{\mathbf{X}}^T \right) \hat{\mathbf{w}} - 2C_r \hat{\mathbf{w}}^T \hat{\mathbf{X}} \mathbf{s} + C_r \mathbf{s}^T \mathbf{s} = \\ & \left(\left[\frac{1}{2} \mathbf{I} + C_r \hat{\mathbf{X}} \hat{\mathbf{X}}^T \right]^{\frac{1}{2}} \hat{\mathbf{w}} \right)^2 - \\ & - 2C_r \left(\hat{\mathbf{w}}^T \left[\frac{1}{2} \mathbf{I} + C_r \hat{\mathbf{X}} \hat{\mathbf{X}}^T \right]^{\frac{1}{2}} \right) \left(\left[\frac{1}{2} \mathbf{I} + C_r \hat{\mathbf{X}} \hat{\mathbf{X}}^T \right]^{-\frac{1}{2}} \hat{\mathbf{X}} \mathbf{s} \right) = \\ & \left(\left[\frac{1}{2} \mathbf{I} + C_r \hat{\mathbf{X}} \hat{\mathbf{X}}^T \right]^{\frac{1}{2}} \hat{\mathbf{w}} - C_r \left[\frac{1}{2} \mathbf{I} + C_r \hat{\mathbf{X}} \hat{\mathbf{X}}^T \right]^{-\frac{1}{2}} \hat{\mathbf{X}} \mathbf{s} \right)^2 + \\ & + C_r \mathbf{s}^T \mathbf{s} - C_r^2 \left[\frac{1}{2} \mathbf{I} + C_r \hat{\mathbf{X}} \hat{\mathbf{X}}^T \right]^{-1} \hat{\mathbf{X}} \mathbf{s}. \end{aligned}$$

Введем замену переменных:

$$\begin{aligned} \mathbf{w}' &= \mathbf{A} \hat{\mathbf{w}} - \mathbf{B}, \text{ где} \\ \mathbf{A} &= \left[\frac{1}{2} \mathbf{I} + C_r \hat{\mathbf{X}} \hat{\mathbf{X}}^T \right]^{\frac{1}{2}}, \\ \mathbf{B} &= C_r \left[\frac{1}{2} \mathbf{I} + C_r \hat{\mathbf{X}} \hat{\mathbf{X}}^T \right]^{-\frac{1}{2}} \hat{\mathbf{X}} \mathbf{s}. \end{aligned} \quad (14)$$

Тогда, учитывая, что

$$C_r \mathbf{s}^T \mathbf{s} - C_r^2 \left[\frac{1}{2} \mathbf{I} + C_r \hat{\mathbf{X}} \hat{\mathbf{X}}^T \right]^{-1} \hat{\mathbf{X}} \mathbf{s} = \text{const},$$

задача оптимизации принимает вид:

$$\begin{aligned} \text{minimize: } & \|\mathbf{w}'\|^2 + \\ & + \sum_{ij} C_{ij} \ln \left(1 + e^{-y_{ij} [\mathbf{A}^{-1}(\mathbf{w}' + \mathbf{B})]^T \hat{\mathbf{x}}_{ij}} \right). \end{aligned} \quad (15)$$

Заменяем параметры регуляризации C_{ij} внутри суммы одним параметром C перед суммой и получим L2-регуляризованную логистическую регрессию:

$$\begin{aligned} \text{minimize: } & \|\mathbf{w}'\|^2 + \\ & + C \sum_{ij} \ln \left(1 + e^{-y_{ij} [\mathbf{A}^{-1}(\mathbf{w}' + \mathbf{B})]^T \hat{\mathbf{x}}_{ij}} \right). \end{aligned} \quad (16)$$

В. Данные

Для обучения модели были выбраны данные из базы PDBBind^{18,19}, которая содержит экспериментально определенные комплексы белков и лигандов и соответствующие им измеренные значения постоянной аффинности связывания. Создатели PDBBind предлагают для исследований два датасета. Собранный вручную "general set" содержит 14620 комплексов белков и лигандов, а также значения аффинности связывания (K_d , K_i и IC_{50}). Кроме того, авторы отдельно отобрали данные лучшего качества и собрали дополнительный "refined set", также состоящий из комплексов белков и лигандов и соответствующих им значений аффинности связывания.

В представленном эксперименте используются данные из "refined set", для которых известно значение константы K_i : всего 39444 конформации, 2076 из которых являются нативными, и 37368 – ненативными. В данном наборе каждому соединению соответствует одна нативная поза с известным значением константы K_i и 18 ненативных поз, для которых значение K_i неизвестно. Таким образом, $P = 2076$ и $D = 18$. В базе PDBBind рассматривается $M_1 = 40$ типов атомов лигандов и $M_2 = 23$ типов атомов белков, а значение Q принято равным 7. Таким образом, каждая конформация характеризуется структурным вектором \mathbf{x}_{ij} размерности $Q \times M_1 \times M_2 = 6440$.

С. Эксперимент

Из исходной выборки \mathbf{X} и соответствующего вектора классов конформаций \mathbf{y} на обучающую выборку ($\mathbf{X}_{\text{train}}$, $\mathbf{y}_{\text{train}}$) было выделено 60% объектов, а на тестовую – все нативные комплексы и соответствующие значения постоянной аффинности связывания (\mathbf{X}_{test} , \mathbf{s}_{test}) из оставшихся 40% выборки \mathbf{X} . Для замены переменных из обучающей выборки были выделены исключительно нативные конформации с известными значениями постоянной K_i , которые образовали выборку $\mathbf{X}_{\text{nat_train}}$ и целевой вектор $\mathbf{s}_{\text{train}}$ соответственно. С помощью этих данных были вычислены матрицы \mathbf{A} и \mathbf{B} .

Чтение и обработка данных, включая замену переменных и вычисление ошибки предсказаний, проводились на языке Python с использованием пакетов NumPy, SciPy и scikit-learn. Поиск минимума функции (16) и оптимизация скорингового вектора \mathbf{w} осуществлялись с помощью библиотеки LIBLINEAR, написанной на языке C++. Поскольку методы LIBLINEAR позволяют работать с формулой логистической регрессии в стандартном виде, в экспоненте в формуле (16) не учитывалось второе слагаемое, не содержащее \mathbf{w}' .

С помощью кросс-валидации было подобрано значение коэффициента регуляризации $C = 1024$, а методом перебора было получено оптимальное значение

коэффициента регуляризации $C_r = 100$.

С помощью полученной модели были предсказаны значения постоянной аффинности для нативных поз тестовой выборки. Корреляция между истинными и предсказанными значениями оказалась равна $\rho = 0.6201$ с p -значением $p = 2.6774 \times 10^{-133}$, коэффициент детерминации $R^2 = -0.2094$, среднеквадратичное отклонение $MSE = 5.4160$.

СПИСОК ЛИТЕРАТУРЫ

- ¹Christoph Sotriffer and Hans Matter. **Virtual Screening: Principles, Challenges, and Practical Guidelines**. Wiley Online Library, 10.1002/9783527633326.ch7 edition, 2011.
- ²Lengauer T, Rarey M (Jun 1996). "Computational methods for biomolecular docking". *Current Opinion in Structural Biology*. 6 (3): 402–6.
- ³Maria Kadukova, Sergei Grudinin. **Convex-PL: a novel knowledge-based potential for protein-ligand interactions deduced from structural databases using convex optimization**. *Journal of Computer-Aided Molecular Design*, October 2017, Volume 31, Issue 10, pp 943–958.
- ⁴Xuan-Yu Meng, Hong-Xing Zhang, Mihaly Mezei, and Meng Cui. **Molecular Docking: A powerful approach for structure-based drug discovery**. *Curr Comput Aided Drug Des*. 2011 Jun 1; 7(2): 146–157.
- ⁵S.P. Boyd and L. Vandenberghe. **Convex optimization**. Cambridge Univ Press, 2004.
- ⁶V. Vapnik. **The nature of statistical learning theory**. Springer, 2000.
- ⁷Maria Kadukova and Sergei Grudinin. **Docking of small molecules to farnesoid X receptors using AutoDock Vina with the Convex-PL potential : lessons learned from D3R Grand Challenge 2**. *J. Comput.-Aided Mol. Des.*, 2017.
- ⁸Sergei Grudinin, Maria Kadukova, Andreas Eisenbarth, Simon Marillet, Frédéric Cazals. **Predicting binding poses and affinities for protein-ligand complexes in the 2015 D3R Grand Challenge using a physical model with a statistical parameter estimation**. *J Comput Aided Mol Des*. 2016 Sep;30(9):791-804. Epub 2016 Oct 7.
- ⁹D. Barber. **Bayesian reasoning and machine learning**. Cambridge University Press, Cambridge, 2012.
- ¹⁰Kevin Vu, John Snyder, Li Li, Matthias Rupp, Brandon F. Chen, Tarek Khelif, Klaus-Robert Müller, Kieron Burke. **Understanding Kernel Ridge Regression: Common behaviors from simple functions to density functionals**.
- ¹¹Press, WH; Teukolsky, SA; Vetterling, WT; Flannery, BP (2007). "Section 19.5. Linear Regularization Methods". **Numerical Recipes: The Art of Scientific Computing (3rd ed.)**. New York: Cambridge University Press. ISBN 978-0-521-88068-8.
- ¹²Sheng-You Huang and Xiaoqin Zou. **An iterative knowledge-based scoring function for protein-protein recognition**. *Proteins: Struct., Funct., Bioinf.*, 72(2):557–579, 2008.
- ¹³Gwo-Yu Chuang, Dima Kozakov, Ryan Brenke, Stephen R Comeau, and Sandor Vajda. **Dars (decoys as the reference state) potentials for protein-protein docking**. *Biophys. J.*, 95(9):4217–4227, 2008.
- ¹⁴Vladimir N Maiorov and Gordon M Grippen. **Contact potential that recognizes the correct folding of globular proteins**. *J. Mol. Biol.*, 227(3):876–888, 1992.
- ¹⁵Jian Qiu and Ron Elber. **Atomically detailed potentials to recognize native and approximate protein structures**. *Proteins: Struct., Funct., Bioinf.*, 61(1):44–55, 2005.
- ¹⁶Dror Tobi and Ivet Bahar. **Optimal design of protein docking potentials: Efficiency and limitations**. *Proteins: Struct., Funct., Bioinf.*, 62(4):970–981, 2006.
- ¹⁷Myong-Ho Chae, Florian Krull, Stephan Lorenzen, and Ernst-Walter Knapp. **Predicting protein complex geometries with a neural network**. *Proteins: Struct., Funct., Bioinf.*, 78(4):1026–1039, 2010.
- ¹⁸Renxiao Wang, Xueliang Fang, Yipin Lu, Chao-Yie Yang, and Shaomeng Wang. **The PDBbind Database: Methodologies And Updates**. *J. Med. Chem.*, 48(12):4111–9, June 2005.
- ¹⁹Renxiao Wang, Xueliang Fang, Yipin Lu, and Shaomeng Wang. **The PDBbind Database: Collection of Binding Affinities for Protein-Ligand Complexes with Known Three-Dimensional Structures**. *J. Med. Chem.*, 47(12):2977–80, June 2004.