

# Вычислительно-эффективное сэмплирование из гауссовского процесса в задаче активного обучения

Вайсер Кирилл Олегович

Московский физико-технический институт

Группа 774

Научный руководитель: к.ф.-м.н. Панов М.Е.

## Проблема

- Активное обучение предполагает недостаток исходных обучающих данных и постепенное пополнение обучающей выборки.
- Для пополнения выборки используются функции неопределенности, показывающие уверенность модели в своем ответе. Для их вычисления необходимо сэмплирование.
- Для сэмплирования необходимо обращать ковариационную матрицу и вычислять разложение Холецкого. Сложность этих операций — кубическая по длине выборки.

## Решение

Использовать декомпозицию апостериорного распределения по правилу Маферона. Это позволит сэмплировать за линейное время по длине выборки.

# Постановка задачи: выбор модели классификации

Задана выборка

$$\mathcal{D} = \{\mathbf{x}_i\} \quad i = 1, \dots, n, \quad \mathbf{x}_i \in \mathbb{R}^d$$

и оракул  $g : \mathbb{R}^d \mapsto \{-1, 1\}$  — функция, которая возвращает правильные ответы.

Требуется найти модель, которая аппроксимирует вероятность принадлежности объекта  $\mathbf{x}$  к классу  $g(\mathbf{x})$ :

$$\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) : \mathbb{R}^d \times \mathbb{R}^h \mapsto [0, 1],$$

где  $\boldsymbol{\theta} \in \mathbb{R}^h$  — параметры модели.

Функция стоимости  $c : \{f, g\} \mapsto \mathbb{R}$  стоимости определяет затраты на вычисление функции в точке  $\mathbf{x}$ . В моделях активного обучения  $c(g) \gg c(f)$

Точки в выборку добавляются на основе значения функций неопределенности:

- 1) средняя энтропия в точке,

$$EH(p(\mathbf{x})) = E \sum_{y \in \{0,1\}} p_y(\mathbf{x}) \log p_y(\mathbf{x}),$$

где  $p_y(\mathbf{x})$  — вероятность принадлежности объекта  $\mathbf{x}$  к классу  $y$

- 2) эпистемическая неопределенность в точке,

$$EH(p(\mathbf{x})) - H(Ep(\mathbf{x})),$$

- 3) максимальная вероятность класса

$$\max_y p_y(\mathbf{x}).$$

Обозначим  $\mathbf{f}_m | \mathbf{y}$  — апостериорное распределение гауссовского процесса на тестовой выборке  $\mathbf{X}_m$  после наблюдения обучающей выборки  $(\mathbf{X}_n, \mathbf{y})$ . Апостериорные моменты записываются как

$$\begin{aligned}\boldsymbol{\mu}_{m|n} &= \mathbf{K}_{m,n} (\mathbf{K}_{n,n} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}, \\ \mathbf{K}_{m,m|n} &= \mathbf{K}_{m,m} - \mathbf{K}_{m,n} (\mathbf{K}_{n,n} + \sigma^2 \mathbf{I})^{-1} \mathbf{K}_{n,m},\end{aligned}$$

где  $\boldsymbol{\mu}_{m|n}$  — матожидание полученного вектора, а  $\mathbf{K}$  — ковариационная матрица.

Сэмплирование производится с помощью схемы:

$$\mathbf{f}_m | \mathbf{y} = \boldsymbol{\mu}_{m|n} + \mathbf{K}_{m,m|n}^{1/2} \boldsymbol{\zeta}, \quad \boldsymbol{\zeta} \sim \mathcal{N}(0, \mathbf{I}).$$

Такой подход требует  $O(n^3)$  операций для вычисления  $(\mathbf{K}_{n,n} + \sigma^2 \mathbf{I})^{-1}$  и  $O(m^3)$  операций для вычисления  $\mathbf{K}_{m,m|n}^{1/2}$ .

Байесовская линейная модель

$$f(\mathbf{x}) = \phi(\mathbf{x})^\top \mathbf{w}, \quad y = f(\mathbf{x}) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma_n^2),$$

где  $\phi: \mathbb{R}^d \mapsto \mathbb{R}^\ell$ , а  $\mathbf{w}$  — параметры модели с априорным распределением  $\mathcal{N}(0, \Sigma_p)$ . Апостериорное распределение вектора ответов:

$$\mathbf{f}_m | \mathbf{x}_m, \mathbf{X}, \mathbf{y} \sim \mathcal{N}(\phi_m^\top \Sigma_p \Phi (K + \sigma_n^2 I)^{-1} \mathbf{y} \\ \phi_m^\top \Sigma_p \phi_m - \phi_m^\top \Sigma_p \Phi (K + \sigma_n^2 I)^{-1} \Phi^\top \Sigma_p \phi_m),$$

где  $\Phi = \phi(\mathbf{X})$ ,  $K = \Phi^\top \Sigma_p \Phi$ . Это распределение соответствует гауссовскому процессу с ковариационной функцией

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \Sigma_p \phi(\mathbf{x}') = \psi(\mathbf{x})^\top \psi(\mathbf{x}'),$$

где  $\psi(\mathbf{x}) = \Sigma_p^{1/2} \phi(\mathbf{x})$ .

Для снижения сложности обучения гауссовского процесса возможен выбор подмножества точек обучающей выборки для вычисления апостериорных моментов. Пусть выбрано множество точек  $\mathbf{Z} = \{z_j\}_{j=1}^v$  и  $u = f(\mathbf{Z}) \sim \mathcal{N}(\boldsymbol{\mu}_u, \Sigma_u)$ .

Тогда апостериорное распределение на тестовой выборке  $f_m \mid \mathbf{u}$  имеет моменты

$$\begin{aligned}\boldsymbol{\mu}_{m|v} &= \mathbf{K}_{m,v} \mathbf{K}_{v,v}^{-1} \boldsymbol{\mu}_v, \\ \mathbf{K}_{m,m|v} &= \mathbf{K}_{m,m} + \mathbf{K}_{m,v} \mathbf{K}_{v,v}^{-1} (\Sigma_u - \mathbf{K}_{v,v}) \mathbf{K}_{v,v}^{-1} \mathbf{K}_{v,m}\end{aligned}$$

Сложность вычисления ковариационной матрицы снижается с  $O(n^3)$  до  $O(v^3)$ ,  $n \gg v$ .

Теорема (Вайсер, 2021, на основе работ Маферона, 1960)

*Пусть  $\mathbf{x}$  и  $\mathbf{z}$  имеют совместное нормальное многомерное распределение. Тогда условное распределение на  $\mathbf{x}$  при  $\mathbf{z} = \beta$  вычисляется как*

$$(\mathbf{a} \mid \mathbf{b} = \beta) \stackrel{d}{=} \mathbf{a} + \text{Cov}(\mathbf{a}, \mathbf{b}) \text{Cov}(\mathbf{b}, \mathbf{b})^{-1} (\beta - \mathbf{b}).$$

Апостериорное распределение может быть представлено суммой априорного распределения и поправки на наблюдаемые данные.



Записав апостериорные распределения весов  $\mathbf{w}$  из байесовской модели и вектора ответов  $\mathbf{f}_m \mid \mathbf{u}$  разреженного процесса в виде разложения по правилу Маферона, получим

$$\begin{aligned}\mathbf{w} \mid \mathbf{y} &\stackrel{d}{=} \mathbf{w} + \Phi^\top (\Phi\Phi^\top + \sigma^2 \mathbf{I})^{-1} (\mathbf{y} - \Phi\mathbf{w} - \boldsymbol{\varepsilon}), \\ \mathbf{f}_m \mid \mathbf{u} &\stackrel{d}{=} \mathbf{f}_m + \mathbf{K}_{m,v} \mathbf{K}_{v,v}^{-1} (\mathbf{u} - \mathbf{f}_v).\end{aligned}$$

Сэмплирование из априорного распределения  $\mathbf{w}$  имеет сложность  $O(\ell)$  и не зависит от размера выборки  $m$ . В то же время вычисление обновления существенно сложнее. В разреженном подходе напротив, сэмплирование из априорного распределения по-прежнему имеет сложность  $O(m^3)$ , однако поправка вычисляется за  $O(m)$ .

Итоговая модель, совмещающая в себе априорное распределение байесовской модели и обновление разреженного процесса

$$(f | \mathbf{u})(\cdot) \stackrel{d}{\approx} \mathbf{w}^\top \phi(\cdot) + \mathbf{h}^\top \mathbf{k}(\cdot, \mathbf{Z}),$$

где  $\mathbf{h} = \mathbf{K}_{v,v}^{-1}(\mathbf{u} - \Phi \mathbf{w})$ .

Сложность сэмплирования из апостериорного распределения  $O(v^3) + O(m)$ .

## 1 AUROC score

$$\frac{\sum_{\mathbf{x}_0 \in \mathcal{Y}^0} \sum_{\mathbf{x}_1 \in \mathcal{Y}^1} 1[f(\mathbf{x}_0) < f(\mathbf{x}_1)]}{|\mathcal{Y}^0| \cdot |\mathcal{Y}^1|},$$

где  $\mathcal{Y}^0, \mathcal{Y}^1$  — множества объектов с негативной и позитивной меткой соответственно.

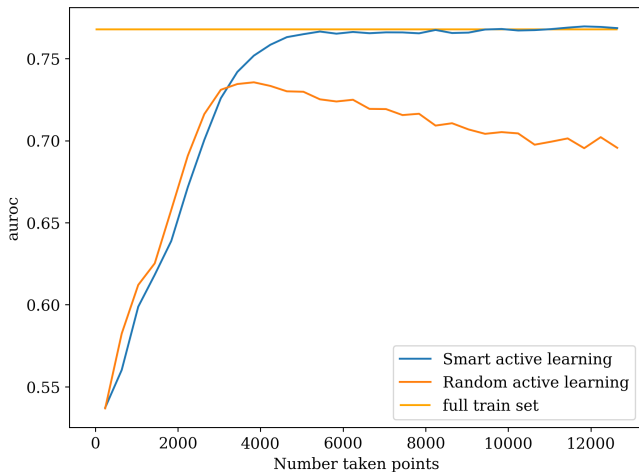
## 2 Общая стоимость вызовов оракула

$$\sum_{i=1}^E n_i c(g),$$

где  $E$  — число шагов активного обучения и  $n_i$  — число добавляемых в выборку точек на  $i$ -ом шаге.

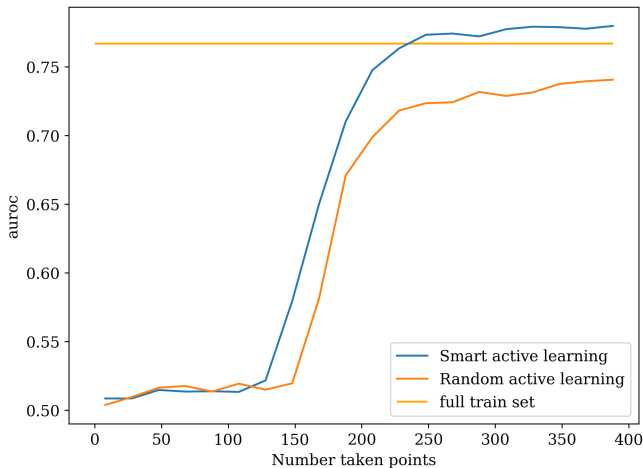
- 1 Обучение модели на всей обучающей выборке и определение метрики качества (AUROC).
- 2 Последовательное обучение на малой подвыборке с добавлением в выборку новых точек, отобранных по значению функции неопределенности.
- 3 Последовательное обучение на малой подвыборке с добавлением в выборку новых точек, отобранных случайно.
- 4 Сравнение результатов.
- 5 Тестирование сложности сэмплинга.

# Зависимость качества AUROC от числа точек в выборке



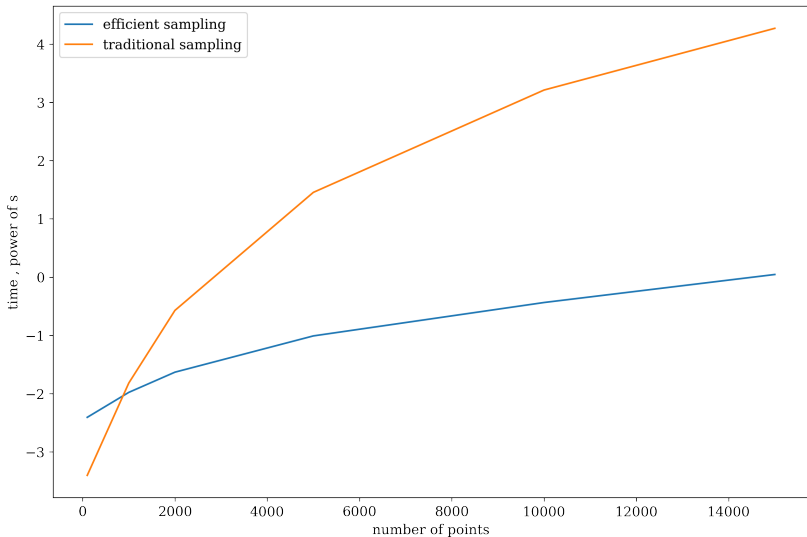
Добавление точек в выборку на основе функций неопределенности позволяет избежать переобучения и достичь уровня качества модели, обученной на всех доступных точках.

# Зависимость качества AUROC от числа точек в выборке



Добавление точек в выборку на основе функций неопределенности позволяет превзойти качество, полученное на всей доступной выборке.

# Сложность сэмплирования (логарифмическая шкала)



- 1 Реализация метода эффективного сэмплирования через разложение апостериорного распределения в сумму априорного и поправки. Использование быстрого сэмплирования из априорного распределения байесовской модели и линейной сложности расчета поправки разреженного гауссовского процесса.
- 2 Применение метода эффективного сэмплирования к задачам активного обучения. Получение превосходящего результата в терминах заданного критерия качества при использовании меньшей обучающей выборки.
- 3 Разработка методики выбора точек для включения в обучающую выборку.



1. Deep learning model structure optimization. Potanin M.S., Vayser K.O., Zholobov V.A., Strijov V.V, Informatics and Applications, 2020, 14(2): 58-65
2. Regularization schedule for neural architecture search Potanin M.S., Vayser K.O., Strijov V.V. Manuscript submitted for publication.