

# Эффективное применение гауссовских процессов к задаче классификации

Вайсер Кирилл Олегович

Московский физико-технический институт

Группа 774, осень 2020

Научный руководитель: Панов М.Е.

## Цель

Предложить модель активного обучения в решении задач классификации, использующую гауссовский процесс в качестве латентной функции.

## Решаемая проблема

Для того, чтобы учесть информацию, содержащуюся в обучающей выборке, необходимо вычислить апостериорное распределение гауссовского процесса. Однако вычисление параметров апостериорного распределения и сэмплирование из него - дорогие по числу вычислений операции.

# Постановка задачи : построение модели

## ❶ выборка

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

## ❷ оракул $g : \mathbb{R}^d \mapsto \{-1, 1\}$ — функция, которая возвращает правильные ответы, но обращение к которой дорогостоящее в некотором смысле, например долго вычисляется или требует много ресурсов. Например командировка бригады бурильщиков для выявления наличия нефти в определенной местности.

## ❸ модель, которая аппроксимирует вероятность объекта $\mathbf{x}$ принадлежать к классу $g(\mathbf{x})$ :

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

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

# Постановка задачи : метрики неопределенности

Для определения точек, в которых необходимо узнать истинную метку для улучшения модели, можно использовать значения метрик неопределенности.

- ①  $\mathbb{E}H(p(x)) = \mathbb{E} \sum_{c \in C} p_c(x) \log p_c(x)$ ,  
где  $p_c(x)$  это вероятность принадлежности объекта  $x$  к классу  $c$ . Эта метрика соответствует средней энтропии в точке.
- ②  $\mathbb{E}H(p(x)) - H(\mathbb{E}p(x))$   
Эпистемическая неопределенность в точке (epistemic uncertainty).
- ③  $\max_c p_c(x)$   
Максимальная вероятность класса.

Гауссовский процесс —случайный процесс, любая конечная совокупность сечений которого имеет совместное нормальное распределение.

Гауссовский процесс характеризуется функцией среднего и ковариационной функцией:

$$m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$$
$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))].$$

Гауссовский процесс задает тем самым распределение функций на своей области определения. Это свойство существенно увеличивает обобщающую способность модели, построенной на гауссовском процессе, а так же делает его хорошим выбором в задачах, требующих расчета метрик неопределенности.

Обозначим  $\mathbf{f}_m \mid \mathbf{y}$  за апостериорное распределение гауссовского процесса в точках  $\{\mathbf{x}'\}_{i=1}^m = X_m$  после наблюдения выборки  $\{\mathbf{x}_i, y_i\}_{i=1}^n$ . Тогда его среднее и матрица ковариации записываются как

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

а сэмплирование производится с помощью стандартной схемы

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

Такой подход требует произвести  $O(n^3) + O(m^3)$  вычислений, что существенно снижает его работоспособность для больших обучающих и тестовых выборок.

Вместо прямого точного вычисления моментов апостериорного распределения и сэмплирования из него факторизуем его на два слагаемых.

Взвешенный подход проистекает из байесовской линейной модели

$$f = \phi(\mathbf{x})^\top \mathbf{w},$$

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

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

где  $K = \Phi^\top \Sigma_p \Phi$ . Тогда мы можем определить ковариационную функцию как

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

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

# Функциональный подход

Гауссовский процесс может быть рассмотрен как распределение над функциями на заданной области определения. Если эта область очень большая, то можно рассмотреть подмножество индуцирующих точек  $\mathbf{Z} = \{z_j\}_{j=1}^V \mathbf{1}^V$ . Пусть  $u = f(\mathbf{Z}) \sim \mathbb{N}(\mu_u, \Sigma_u)$  Тогда  $f_m \mid \mathbf{u}$  имеет следующие моменты

$$\begin{aligned}\mu_{m|v} &= K_{m,v} K_{v,v}^{-1} \mu_v \\ K_{m,m|v} &= K_{m,m} + K_{m,v} K_{v,v}^{-1} (\Sigma_u - K_{v,v}) K_{v,v}^{-1} K_{v,m}\end{aligned}$$

Такой подход позволяет снизить сложность обучения с  $O(n^3)$  до  $O(v^3)$ , что, при выборе подходящего индуцирующего множества, может существенно снизить затраты на вычисления. Поскольку используются не все точки из обучающей выборки, она тем самым *разрезается*, поэтому назовем гауссовский процесс, обученный таким образом, *разреженным*.

Как объединить эти два подхода, чтобы использовать их сильные стороны?



Пусть  $\mathbf{a}$  и  $\mathbf{b}$  имеют совместное нормальное многомерное распределение. Тогда условное распределение на  $\mathbf{a}$  при  $\mathbf{b} = \boldsymbol{\beta}$  вычисляется как

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

Эта теорема позволяет разложить апостериорное распределение в сумму априорного и обновления из-за несоответствия данных априорному распределению.

Применение теоремы Маферона к взвешенному и функциональному подходам позволяет получить следующий вид распределений:

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

Отсюда можно увидеть, что в взвешенном подходе вся сложность приходится на вычисление слагаемого обновления, в то время как в функциональном обновлении линейно по длине входа  $m$ , но имеет кубическую сложность в сэмплировании из априорного распределения. Таким образом, мы можем совместить два подхода, чтобы существенно ускорить сэмплирование.

Итоговая модель, совмещающая в себе сильные стороны обоих подходов, выглядит следующим образом

$$(f | \mathbf{u})(\cdot) \stackrel{d}{\approx} \sum_{i=1}^{\ell} w_i \phi_i(\cdot) + \sum_{j=1}^v h_j k(\cdot, \mathbf{z}_j),$$

где  $\mathbf{h} = \mathbf{K}_{v,v}^{-1}(\mathbf{u} - \Phi \mathbf{w})$ . Таким образом, разложение с помощью теоремы Маферона позволяет существенно ускорить сэмплирование.

## Критерий качества:

① Точность

$$A = \frac{1}{m} \sum_{i=1}^m [\hat{y}_i = y_i]$$

② AUROC score

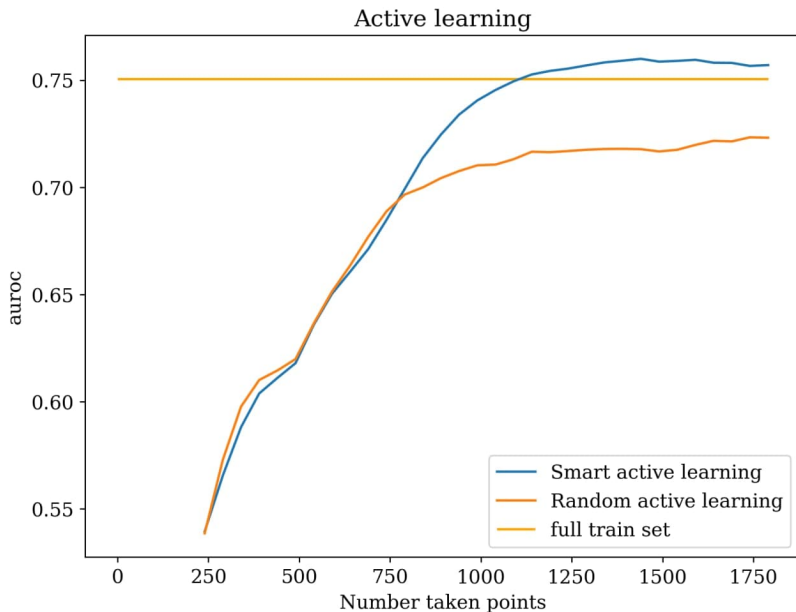
③ Скорость работы

Для эксперимента был рассмотрен датасет default of credit card clients.

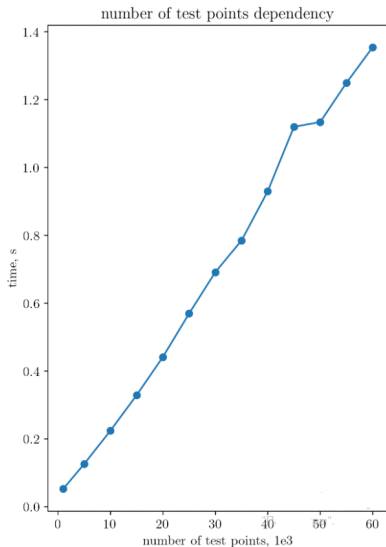
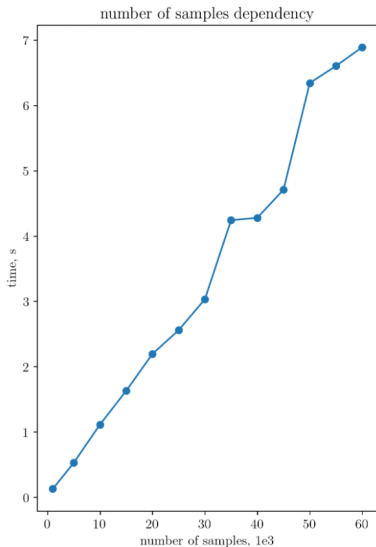
<b>Data Set Characteristics:</b>	Multivariate	<b>Number of Instances:</b>	30000	<b>Area:</b>	Business
<b>Attribute Characteristics:</b>	Integer, Real	<b>Number of Attributes:</b>	24	<b>Date Donated</b>	2016-01-26
<b>Associated Tasks:</b>	Classification	<b>Missing Values?</b>	N/A	<b>Number of Web Hits:</b>	568344

Эксперимент проводился следующим образом:

- 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 метрики неопределенности.
- 3 Последовательное обучение на малой подвыборке с добавлением в выборку новых точек, отобранных случайно.
- 4 Сравнение результатов.
- 5 Тестирование скорости сэмплирования.



# Текущие результаты



Обозначим за  $GP(0, k)$  гауссовский процесс с средним 0 и ковариационной функцией  $k$ . Пусть  $f \sim (0, k)$ . Обозначим за  $f | \mathbf{y}$  точное апостериорное распределение,  $f^{(s)}$  - апостериорное распределение в терминах функционального подхода,  $f^{(d)}$  - апостериорное распределение в терминах эффективного подхода и  $f^{(w)}$  - априорное распределение в терминах взвешенного подхода. Тогда

$$\begin{aligned} W_{2, L^2(\mathcal{X})}(f^{(d)}, f | \mathbf{y}) &\leq \\ &\leq W_{2, L^2(\mathcal{X})}(f^{(s)}, f | \mathbf{y}) + C_1 W_{2, C(\mathcal{X})}(f^{(w)}, f), \end{aligned}$$

где  $W_{2, L^2(\mathcal{X})}$ ,  $W_{2, C(\mathcal{X})}$  - расстояния Вассерштейна в пространствах  $L^2(\mathcal{X})$  и  $C(\mathcal{X})$  и

$$C_1 = \sqrt{2 \left( 1 + \|k\|_{C(\mathcal{X}^2)}^2 \|\mathbf{K}_{mm}^{-1}\|_{L(\ell^\infty; \ell^1)}^2 \right)}.$$



## Ключевые работы

Rasmussen, C. E., Williams, C. K. I. (2005). Gaussian processes for machine learning. MIT Press.

J.T. Wilson, V. Borovitskiy, A. Terenin, P. Mostowsky, M. P. Deisenroth. Efficiently Sampling Functions from Gaussian Process Posteriors. International Conference on Machine Learning, 2020