

Эффективное применение гауссовских процессов к задаче классификации

Вайсер Кирилл Олегович

Московский физико-технический институт

Отчет о НИР/Группа 774, осень 2020

Научный руководитель: Панов М.Е.

Цель

Предложить быстрый и эффективный в терминах заданных критериев качества подход к использованию гауссовских процессов в задаче классификации.

Решаемая проблема

Вычисление параметров апостериорного распределения - долгая по времени операция. Кроме того, в явном виде оно не обладает свойством сопряженности.

Метод решения

Предлагаемый метод заключается в следующем:

- 1 Использовать аппроксимацию Лапласа для получения свойства сопряженности.
- 2 Использовать нейронную сеть для обучения ковариационной функции
- 3 Использовать методы эффективного сэмплирования для получения предсказаний.

- ❶ выборка

$$\mathcal{D} = \{x_i, y_i\} \quad i = 1, \dots, m, \quad x_i \in \mathbb{R}^m \quad y_i \in \{-1, 1\}$$

- ❷ модель

$$g(x, w) : \mathbb{R}^m \times \mathbb{R}^n \longrightarrow \mathbb{R}^m \times \mathbb{R}^h,$$

где $w \in \mathbb{R}^n$ — пространство параметров модели.

- ❸ априорное распределение гауссовского процесса

$$p(f) \sim \mathcal{GP}(\mu, K)$$

Постановка задачи

1 Факторизация

Ковариационная функция факторизуется как

$$k(x, x') = h(x)^\top h(x') + \sigma^2$$

Модель обучается для получения факторизации h и представления ковариации как

$$K = H^\top H + \sigma^2 I$$

2 Классификация

Гауссовский процесс используется для получения латентных переменных f . После чего эти латентные переменные отображаются в отрезок $[0, 1]$ и используются для оценки вероятностей классов. Рассматривается преобразование

$$p = \Phi\left(\frac{f}{1 + \sigma^2}\right),$$

где Φ - нормальная функция распределения.

Критерий качества:

- 1 Точность

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

- 2 AUC score
- 3 Неопределенность

$$\mathcal{H}_{pred} = - \sum_{c=-1,1} p(y = c|x) \log p(y = c|x)$$

- 4 Скорость работы

1 Аппроксимация Лапласа

$$K_{post} = (K_{pr}^{-1} + W)^{-1},$$

где $W = -\nabla^2 \log p(y|f)$

2 Факторизация ковариационной функции

$$K = H^T H + \sigma^2 I$$

где G - выход сети.

3 Эффективное сэмплирование

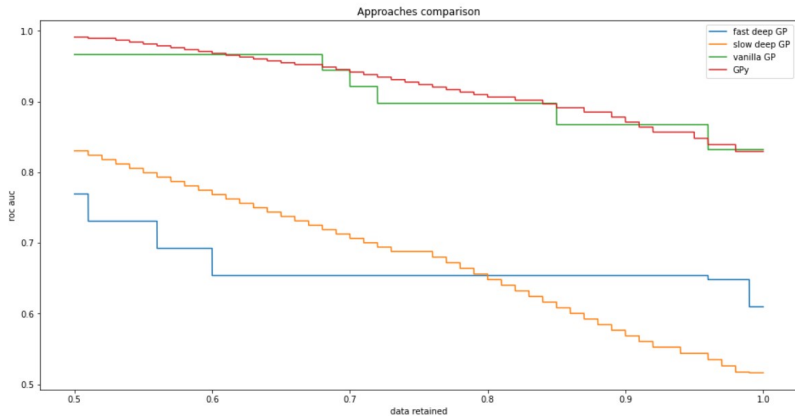
Разреженное представление:

$$p(\mathbf{f}_* | \mathbf{y}) \approx \int_{\mathbb{R}^m} p(\mathbf{f}_* | \mathbf{u}) q(\mathbf{u}) d\mathbf{u}$$

Представление Фурье:

$$f(\cdot) = \sum_{i=1}^l w_i \phi_i(\cdot)$$

Текущие результаты



На данный момент предложенные модели проигрывают уже существующим по критерию AUC. Предлагается выяснить причину такого расхождения и устранить ее, если возможно.

- ❶ Добиться улучшения результатов работы сети.
- ❷ Реализовать эффективное сэмплирование
- ❸ Исследовать подходы к поиску оптимальной подвыборки (BALD)