

Байесовское мультимоделирование: гауссовские процессы

Московский Физико-Технический Институт

2021

Определение (wiki)

- Случайный процесс f_t с непрерывным временем является гауссовским тогда и только тогда, когда для любого конечного множества индексов t_1, \dots, t_k : f_{t_1}, \dots, f_{t_k} — многомерная гауссовская случайная величина.
- Всякая линейная комбинация имеет f_{t_1}, \dots, f_{t_k} одномерное нормальное распределение.

Определение (упрощенное)

Назовем гауссовским процессом $\mathcal{GP}(\mathbf{m}(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$ распределение на множестве функций, такое что для любых \mathbf{x}, \mathbf{x}' : $\mathcal{GP}(\mathbf{m}(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$ — гауссовское распределение.

Пример: регрессия

$$\mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}),$$

где \mathbf{K} — ковариационная матрица для объектов выборки \mathbf{X} .

$$y \sim f + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2).$$

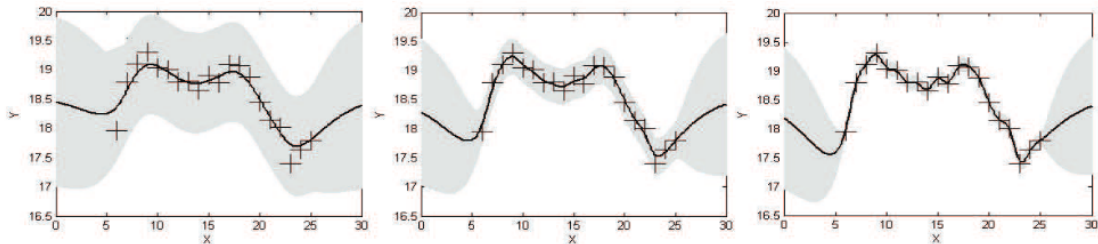
Тогда $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{K} + \sigma^2 \mathbf{I})$, предсказание для новых объектов $\hat{\mathbf{y}}$:

$$\hat{\mathbf{y}} \sim \mathcal{N}(\mathbf{K}'^T (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}, \mathbf{K}'' - \mathbf{K}'^T (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{K}').$$

Отличие гауссовых процессов от других подходов

- Модель непараметрическая
 - ▶ параметризация только ковариационной функции и уровня шума
 - ▶ оптимизация — через ОМП
- Априорное распределение задается для функции, а не для параметров
 - ▶ задание гауссового априорного распределения на параметры: тоже получаем гауссовый процесс с вырожденной матрицей ковариации
- Сложность предсказания: $O(N^3)$.

Влияние σ^2 на предсказательную способность



(McDuff, 2010)

Ковариационная функция

Основные свойства:

- Симметричность $\mathbf{K}(x, x') = \mathbf{K}(x', x)$
- Положительная полуопределенность: $\mathbf{v}^T \mathbf{K} \mathbf{v} \geq 0$
- Стационарность: $K(x, x') = K(x + a, x' + a)$
- Изотропность: зависимость только от $\|x - x'\|$

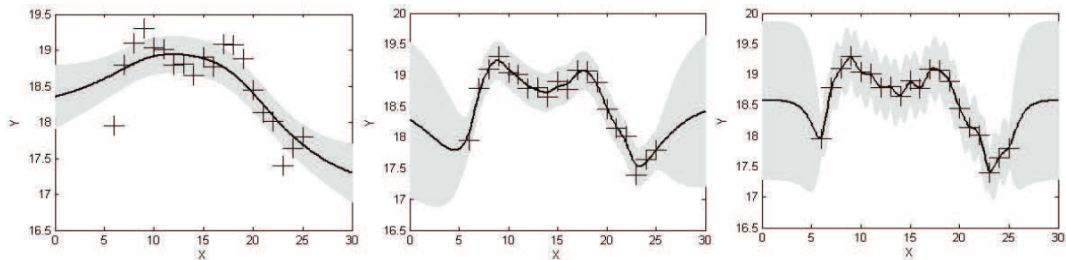
Ковариационная функция: виды

Экспоненциальная:

$$K = \sigma_0^2 \exp\left(\frac{-(x - x')^2}{2\lambda}\right)$$

- Линейная: $K = \sigma_0^2 + xx'$
- Броуновская: $K = \min(x, x')$
- Периодическая: $K = \exp\left(\frac{-2\sin^2(\frac{x-x'}{2})}{\lambda^2}\right)$
- Задаваемая нейросетью

Влияние λ на предсказательную способность



(McDuff, 2010)

Ковариация Матерна

$$K = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}|x - x'|}{\lambda} \right)^\nu K_\nu \left(\frac{\sqrt{2\nu}|x - x'|}{\lambda} \right),$$

K_ν — модифицированная функция Бесселя.

- Стационарная и изотропная функция
- При $\nu \rightarrow \infty$: экспоненциальная функция
- При конечных ν : предсказательная функция менее гладкая

Оптимизация гиперпараметров модели с применением гауссовых процессов

(Snoek et al., 2012)

$$\mathbf{f}(\mathbf{x}, \mathbf{w}, \mathbf{h}) = \mathbf{f}(\mathbf{w}(\mathbf{h}), \mathbf{h}|\mathbf{x}).$$

Модель \mathbf{f} рассматривается как функция от гиперпараметров:

$$\mathbf{f} \sim \mathcal{GP}, \quad \mathbf{y} \sim \mathcal{N}(\mathbf{f}, \sigma^2).$$

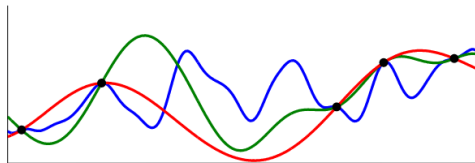
Используется ковариация Матерна, с $\nu = 2.5$.

Выбор следующей точки для оценки

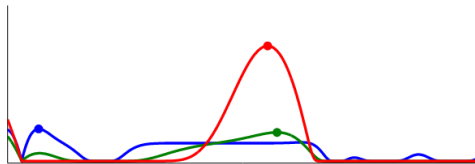
Выбор производится с применением одной из следующих функций (Acquisition function):

- Probability of Improvement: $\Phi(\gamma), \gamma = \frac{f(\mathbf{h}^*) - \mu(\mathbf{h})}{\sigma(\mathbf{h})}$
- Expected Improvement: $\approx E\gamma$
- Confidence Bound: $\mu(\mathbf{h}) - k\sigma(\mathbf{h})$.

Влияние λ на оптимизацию гиперпараметров



(a) Posterior samples under varying hyperparameters



(b) Expected improvement under varying hyperparameters

Результаты

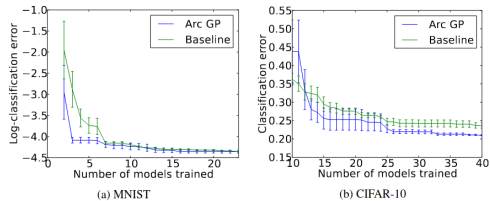
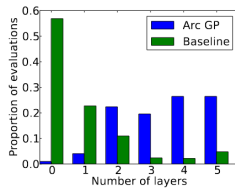


Figure 2: Bayesian optimization results using the arc kernel.



Ковариационная функция для условных параметров

Задача выбора архитектуры нейросети:

- Количество нейронов на скрытом слое: **ок, можем задавать как вещественный гиперпараметр**
- Количество слоев на скрытом слое: **ок, можем задавать как вещественный гиперпараметр**
- Как сравнивать две архитектуры с разным количеством слоев? Как сравнивать архитектуру [100, 100] и [100]?

(Swersky et al., 2014): нужна специальная ковариационная функция!

Идея ковариационной функции

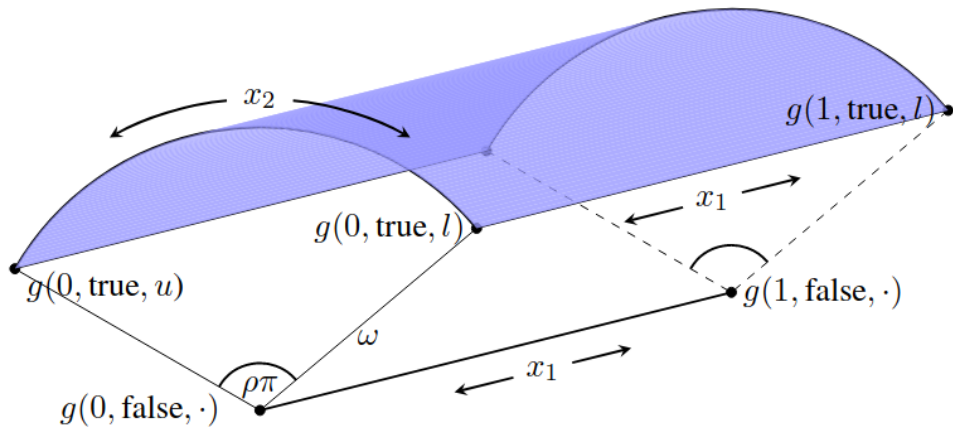
- If we are comparing two points for which the same parameters are relevant, the value of any unused parameters shouldn't matter,

$$k((x_1, \text{false}, x_2), (x'_1, \text{false}, x'_2)) = k((x_1, \text{false}, x''_2), (x'_1, \text{false}, x'''_2)), \forall x_2, x'_2, x''_2, x'''_2; \quad (1)$$

- The covariance between a point using both parameters and a point using only one should again only depend on their shared parameters,

$$k((x_1, \text{false}, x_2), (x'_1, \text{true}, x'_2)) = k((x_1, \text{false}, x''_2), (x'_1, \text{true}, x'''_2)), \forall x_2, x'_2, x''_2, x'''_2. \quad (2)$$

Идея ковариационной функции



Результаты

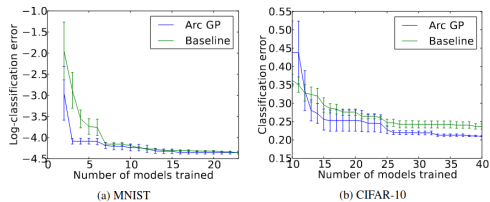
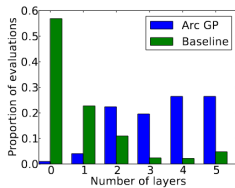


Figure 2: Bayesian optimization results using the arc kernel.



Литература

- Bishop C. M. Pattern recognition //Machine learning. – 2006. – Т. 128. – №. 9.
- Daniel McDuff, Gaussian Processes,
<https://courses.media.mit.edu/2010fall/mas622j/ProblemSets/slidesGP.pdf>
- Chris Williams, Gaussian Processes for Machine Learning,
<https://www.newton.ac.uk/files/seminar/20070809140015001-150844.pdf>
- Ed Snelson, tutorial: Gaussian process models for machine learning,
<https://mlg.eng.cam.ac.uk/tutorials/06/es.pdf>
- Snoek J., Larochelle H., Adams R. P. Practical bayesian optimization of machine learning algorithms //Advances in neural information processing systems. – 2012. – Т. 25.
- Swersky K. et al. Raiders of the lost architecture: Kernels for Bayesian optimization in conditional parameter spaces //arXiv preprint arXiv:1409.4011. – 2014.