

Методы аппроксимации обратной ковариационной матрицы для эффективной оптимизации правдоподобия гауссовского процесса

Даниил Кононенко

1. Институт Проблем Передачи Информации
127994, г. Москва, ГСП-4, Большой Каретный Переулок, 19, стр.1

2. DATADVANCE

105064, г. Москва, Садово-Черногорская улица, 13/3

3. МФТИ

141700, г. Долгопрудный, Институтский переулок, 9
daniil.kononenko@datadvance.net

Аннотация

Восстановление регрессии на основе гауссовских процессов — байесовский непараметрический метод, показывающий хорошие результаты во многих приложениях. Параметры модели настраиваются с помощью метода максимума правдоподобия. Каждый раз при подсчете правдоподобия и его производных необходимо выполнять подсчет обратной матрицы ковариации и ее детерминанта, что занимает порядка $O(N^3)$ операций, где N — размер обучающей выборки. В работе предлагается метод аппроксимации обратной матрицы ковариации и ее детерминанта за $O(N^2)$ операций. Проведен ряд вычислительных экспериментов, которые показывают значительное ускорение предложенного метода по сравнению со стандартными подходами.

1. Введение

При построении моделей на основе данных часто необходимо решать задачу аппроксимации неизвестной зависимости. Одним из наиболее популярных методов решения этой задачи является регрессия на основе гауссовских процессов. Наиболее распространенным способом выбора параметров ковариационной функции гауссовского процесса по выборке является метод максимума правдоподобия [1],[7].

Одной из проблем при практической реализации данного метода является значительное время обучения при больших размерах выборок. В стандартных реализациях одна оценка правдоподобия и его производных требует порядка $O(N^3)$ операций, где N — размер обучающей выборки [7]. При использовании большинства классических методов оптимизации для нахождения максимума правдоподобия, на-

пример, градиентных методов, необходимо несколько сотен оценок целевой функции, что делает время обучения достаточно большим уже при $N \approx 500$ [4].

Цель данной работы — предложить алгоритм, уменьшающий время обучения без значительной потери качества аппроксимации. В разделе 2 приводится постановка задачи аппроксимации на основе гауссовских процессов. В разделе 3 делается обзор существующих методов ускорения вычисления правдоподобия и его производных. В разделе 4 описывается предлагаемый алгоритм обучения, и в разделах 5, 6 обсуждаются результаты вычислительного эксперимента.

2. Гауссовские процессы

Пусть $f(x)$ — некоторая неизвестная функциональная зависимость, $f(x) \in \mathbb{R}$, $x \in \mathbb{X} \subset \mathbb{R}^n$, $f(x)$ — гладкая функция. Пусть зависимость $y(x)$ порождена моделью: $y(x) = f(x) + \varepsilon(x)$, $\varepsilon(x) \sim \mathcal{N}(0, \sigma^2)$ — гауссовский белый шум. Пусть также известна обучающая выборка $D = (\mathbf{x}, \mathbf{y}) = \{(x_i, y_i = y(x_i))\}_{i=1}^N$. Задача состоит в построении аппроксимации неизвестной зависимости $y = \hat{f}(x)$. Оценка качества полученной аппроксимации — средняя квадратичная ошибка на независимой тестовой выборке $D' = \{(x'_j, y'_j = f(x'_j))\}_{j=1}^{N'}$:

$$Q(\hat{f}) = \frac{1}{N'} \sum_{j=1}^{N'} (y'_j - \hat{f}(x'_j))^2. \quad (1)$$

Построение аппроксимации с помощью гауссовских процессов основано на следующих вероятностных предположениях. Пусть задано случайное поле $f(x, \omega)$, $\omega \in \Omega$ — случайное событие в Ω . Предполагается, что для произвольного $x \in \mathbb{X}$ существуют

первый и второй моменты:

$$\begin{aligned} M(x) &= \mathbb{E}f(x), \\ K(x_1, x_2) &= \mathbb{E}(f(x_1) - \mathbb{E}f(x_1))(f(x_2) - \mathbb{E}f(x_2)), \end{aligned}$$

а также условное математическое ожидание

$$\mathbb{E}(f(x)|f(x_1), f(x_2), \dots, f(x_l)).$$

Далее будем предполагать без ограничения общности, что среднее значение нулевое $M(x) = 0$. Предположим также, что случайное поле — гауссовское. Для такого поля совместное распределение $f(x_1), f(x_2), \dots, f(x_l)$ — нормальное и, следовательно, определяется математическим ожиданием и ковариационной функцией. Пусть ковариационная функция $K_0(x, x')$ гауссовского поля $f(x)$ принадлежит некоторому параметрическому семейству

$$K_0(x, x') = K_0(x, x'|\Theta), \quad (2)$$

где Θ — некоторый набор параметров. Тогда ковариационная функция процесса $y(x)$ имеет вид:

$$K(x, x'|\Theta, \sigma^2) = K_0(x, x'|\Theta) + \sigma^2 \delta(x, x'). \quad (3)$$

При таких предположениях условное распределение $y(x)$ в произвольной точке x не из обучающей выборки имеет вид [7]:

$$\text{Law}(y(x)|y(x_1), y(x_2), \dots, y(x_n)) \sim \mathcal{N}(\hat{f}(x), \hat{\sigma}^2(x)).$$

Математическое ожидание используются для построения аппроксимации [7]:

$$\hat{f}(x) = \mathbf{k}(x)\mathbf{K}^{-1}\mathbf{y},$$

оценка дисперсии имеет вид [7]:

$$\hat{\sigma}^2(x) = K_0(x, x) + \sigma^2 - \mathbf{k}(x)\mathbf{K}^{-1}\mathbf{k}(x)^T,$$

где $\mathbf{k}(x) = \{K(x, x_i)\}_{i=1}^N$, $\mathbf{K} = \{K(x_i, x_j)\}_{i,j=1}^N$.

Значения неизвестных параметров Θ, σ^2 восстанавливаются по обучающей выборке D с помощью метода максимума правдоподобия. Логарифм правдоподобия имеет вид [7]:

$$l(\mathbf{a}) = \log p(\mathbf{y}|\mathbf{x}, \mathbf{a}) = -\frac{1}{2}\mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}| - \frac{N}{2} \log 2\pi, \quad (4)$$

где $\mathbf{a} = \{\Theta, \sigma^2\}$

3. Практическая реализация

Оценки параметров

$$\hat{\mathbf{a}} = \arg \max_{\mathbf{a}} l(\mathbf{a}) \quad (5)$$

предлагается искать с помощью метода сопряженных градиентов с мультистартом. Производная правдоподобия по параметру a_i [7]:

$$\frac{\partial l}{\partial a_i} = -\text{tr} \left(\mathbf{K}^{-1} \frac{\partial \mathbf{K}}{\partial a_i} \right) + \mathbf{y}^T \mathbf{K}^{-1} \frac{\partial \mathbf{K}}{\partial a_i} \mathbf{K}^{-1} \mathbf{y}. \quad (6)$$

Чтобы вычислить правдоподобие и его производные в некоторой точке, необходимо знать значения следующих величин:

$$\mathbf{K}^{-1} \mathbf{y}, \text{tr} \left(\mathbf{K}^{-1} \frac{\partial \mathbf{K}}{\partial a_i} \right), \log(\det(\mathbf{K})). \quad (7)$$

Стандартным подходом является использование разложения Холецкого $\mathbf{K} = \mathbf{U}^T \mathbf{U}$, здесь \mathbf{U} — верхняя треугольная матрица. Матрица \mathbf{K} имеет размерность $N \times N$, поэтому вычисление разложения Холецкого требует порядка $O(N^3)$ операций [7].

Существует несколько подходов для приближенного вычисления величин (7).

Одним из возможных методов является переход к разреженным гауссовским процессам [7]. Их основной идеей является аппроксимация матрицы \mathbf{K} матрицей неполного ранга с помощью введения некоторого (меньшего, чем число прецедентов) числа новых входов, например, из числа прецедентов. Если число новых прецедентов $M < N$, то сложность вычисления величин (7) становится равной $O(M^2N)$. Однако, чтобы добиться значительного ускорения, необходимо в несколько раз уменьшать число входов (например $M \sim \frac{N}{3}$), что значительно уменьшает точность аппроксимации. Кроме того, такой метод может быть нестабилен в силу итеративного характера процедуры оптимизации правдоподобия.

Другие подходы используют аппроксимацию непосредственно величин (7). Основным методом аппроксимации величины $\mathbf{K}^{-1} \mathbf{y}$ является численная минимизация квадратичной функции [3]

$$g(\mathbf{u}) = \frac{1}{2} \mathbf{u}^T \mathbf{K} \mathbf{u} - \mathbf{u}^T \mathbf{y}. \quad (8)$$

Точный минимум $\mathbf{u}^* = \mathbf{K}^{-1} \mathbf{y}$. Используя градиентный метод оптимизации, например, метод сопряженных градиентов, приближенное решение может быть найдено за $O(N^2)$ операций [3], [4].

Аппроксимация $\text{tr} \left(\mathbf{K}^{-1} \frac{\partial \mathbf{K}}{\partial a_i} \right)$ может быть осуществлена с помощью метода Монте-Карло [3]:

$$\text{tr} \left(\mathbf{K}^{-1} \frac{\partial \mathbf{K}}{\partial a_i} \right) = E \left[d^T \mathbf{K}^{-1} \frac{\partial \mathbf{K}}{\partial a_i} d \right], \quad (9)$$

где d — случайный вектор, имеющий стандартное нормальное распределение.

В работах [5], [6], [8], [9] предложен альтернативный метод аппроксимации величин (7). Итеративный характер любой численной процедуры оптимизации правдоподобия (5) ведет, в частности, к тому, что правдоподобие (4) и его производные (6) на последовательных итерациях j и $(j+1)$ вычисляются в достаточно близких друг к другу точках \mathbf{a}_j и \mathbf{a}_{j+1} . Предлагается итеративный метод аппроксимации обратной ковариационной матрицы \mathbf{K}_{j+1}^{-1} по

известной матрице \mathbf{K}_j^{-1} . Для этого квадратичная форма (8) оптимизируется квазиньютоновским методом [4]. Его преимущество в том, что обратная ковариационная матрица обновляется итеративно вместе с оптимизацией по u :

$$\begin{aligned} u_{j+1} &= u_j - \eta_j \mathbf{K}_j^{-1} \nabla g(u_j), \\ \mathbf{K}_{j+1}^{-1} &= \mathbf{K}_j^{-1} + \Delta \mathbf{K}_j^{-1}. \end{aligned} \quad (10)$$

Здесь j и $(j+1)$ — последовательные итерации оптимизации квадратичной функции (8), u_j и \mathbf{K}_j^{-1} — приближения целевой переменной и обратной матрицы на этих итерациях, $\nabla g(u) = \mathbf{K}u - \mathbf{y}$, η_j — длина шага линейного поиска. Предлагаемый метод обновления обратной матрицы $\Delta \mathbf{K}_j^{-1}$ — формула BFGS [4]:

$$\Delta \mathbf{K}_j^{-1} = \left(1 + \frac{q_j^T \mathbf{K}_j^{-1} q_j}{q_j^T p_j} \right) \frac{p_j p_j^T}{p_j^T q_j} - \frac{p_j q_j^T \mathbf{K}_j^{-1} + \mathbf{K}_j^{-1} q_j p_j^T}{q_j^T p_j}, \quad (11)$$

где $p_j = u_{j+1} - u_j$, $q_j = \nabla g(u_{j+1}) - \nabla g(u_j)$.

При вычислении правдоподобия в точке a_{j+1} в качестве начальных приближений u_0 и \mathbf{K}_0^{-1} берутся значения, полученные в предыдущей точке a_j . Таким образом, величины $\mathbf{K}^{-1} \mathbf{y}$ и \mathbf{K}^{-1} могут быть найдены итеративно. Каждое обновление (10) требует $O(N^2)$ операций. Теоретически метод сходится за N итераций [4], и тогда сложность подсчета правдоподобия равна $O(N^3)$. Однако достаточная точность подсчета величин u_j и \mathbf{K}_j^{-1} часто достигается раньше. Авторы предлагают два критерия остановки оптимизации (8):

$$\|\nabla g(u_j)\|_\infty \leq \frac{\varepsilon}{N} \text{ или } j > J,$$

где ε , J — фиксированные константы. Тогда сложность метода порядка $JN^2 = O(N^2)$ операций.

Т.к. предыдущее приближение используется в качестве инициализации для нового, то ошибка аппроксимации может накапливаться. Также предыдущее значение \mathbf{a}_j может достаточно сильно отличаться от \mathbf{a}_{j+1} . Тогда достаточная точность аппроксимации обратной матрицы может быть не достигнута. Предлагается следующий критерий качества:

$$\frac{|\text{tr}(\tilde{\mathbf{K}}^{-1} \mathbf{K}) - N|}{N} < \varepsilon_1, \quad (12)$$

где $\tilde{\mathbf{K}}^{-1}$ — построенная аппроксимация, ε_1 — фиксированная константа. Если этот критерий не выполнен, то производится рестарт: для данного значения параметров \mathbf{a}_j обратная матрица вычисляется точно с помощью преобразования Холецкого.

Таким образом, в ходе оптимизации (5) при подсчете правдоподобия и его производных может быть 2 варианта:

1. Критерий (12) выполнен, аппроксимация достаточно точна. Тогда сложность равна $O(N^2)$ операций.

2. Критерий не выполнен, тогда производится рестарт и сложность становится равной $O(N^3)$ операций.

В работе [9] предлагается дополнение к методу обновления (11): на каждой итерации выбирается случайная размерность $i \in \{1, \dots, N\}$ матрицы \mathbf{K} , вычисляется величина $s = \mathbf{K}_j^{-1}(i, :) \times \mathbf{K}(:, i)$ и на нее

делится текущее приближение: $\mathbf{K}_j^{-1} := \frac{\mathbf{K}_j^{-1}}{s}$. Таким образом предлагается точнее аппроксимировать информацию об амплитуде матрицы \mathbf{K} .

Для аппроксимации члена $\log(\det(\mathbf{K}))$ в работах [6], [8] предлагается следующий подход.

$$\begin{aligned} \log(\det(\mathbf{K})) &= \text{tr}(\log(\mathbf{K})), \\ \log(\mathbf{K}) &= - \sum_{i=1}^{\infty} \left(\frac{\mathbf{K}^i}{i} \right) - \text{матричный логарифм}. \end{aligned} \quad (13)$$

Используется двухуровневая схема:

1. Берется некоторое конечное число членов в сумме (13): $\log(\mathbf{K}) = - \sum_{i=1}^L \left(\frac{\mathbf{K}^i}{i} \right)$
2. Каждый элемент вида $\text{tr}(\mathbf{K}^i)$ оценивается с помощью метода Монте-Карло (9).

Каждый из двух этапов вносит ошибку в итоговую аппроксимацию, и авторами предлагается ряд компенсирующих схем, уменьшающих эти ошибки.

4. Предлагаемый метод аппроксимации

В данной работе развивается подход, предложенный в [5], [6], [8], [9]. При использовании квазиньютоновских методов для аппроксимации \mathbf{K}^{-1} с помощью (10) используются обновления низкого ранга: $\mathbf{K}_{j+1}^{-1} = \mathbf{K}_j^{-1} + uv^T$, где u , v — некоторые вектора. Тогда, зная детерминант матрицы \mathbf{K}_j^{-1} , можно точно вычислить детерминант матрицы \mathbf{K}_{j+1}^{-1} , используя лемму об определителе матрицы:

$$\det(A + uv^T) = \det(A)(1 + v^T A^{-1} u). \quad (14)$$

Если использовать формулу BFGS (11) для обновления обратной матрицы (10), то из (14), получаем

$$\det(\mathbf{K}_{j+1}^{-1}) = \det(\mathbf{K}_j^{-1}) \frac{p_j^T \mathbf{K}_j p_j}{q_j^T p_j}, \quad (15)$$

где p_j , q_j имеют тот же смысл, что и в (11).

Видно, что для того, чтобы итеративно обновлять детерминант матрицы, необходимо знать текущее приближение для необращенной матрицы \mathbf{K}_j . Для этого на каждой итерации необходимо обновлять еще и значение \mathbf{K}_j , что влечет дополнительные

временные затраты. Поэтому предлагается использовать формулу DFP [4]:

$$\Delta \mathbf{K}_j^{-1} = \frac{p_j p_j^T}{q_j^T p_j} - \frac{\mathbf{K}_j^{-1} q_j q_j^T (\mathbf{K}_j^{-1})^T}{q_j^T \mathbf{K}_j^{-1} q_j}. \quad (16)$$

Тогда формула для обновления значения детерминанта имеет вид

$$\det(\mathbf{K}_{j+1}^{-1}) = \det(\mathbf{K}_j^{-1}) \frac{q_j^T p_j}{q_j^T \mathbf{K}_j^{-1} q_j}, \quad (17)$$

в который уже не входит необращенная матрица \mathbf{K}_j .

Для более точной аппроксимации информации об амплитуде матрицы \mathbf{K} предлагается следующий метод: на каждой итерации матрица \mathbf{K}_j^{-1} делится на среднее значение диагонали $s = \text{diag}(\mathbf{K}_j^{-1} \times \mathbf{K})$. Тогда значение критерия (12) всегда становится равным нулю. Вместо него предлагается использовать следующие два критерия для рестарта:

1. Если начальная норма градиента $\|\nabla g(u_0)\| > h$, где h — некоторый порог, то выполняется рестарт. Рассматривается l_1 -норма $\|\nabla g(u_0)\| = \max_{i=1, \dots, N} g_i(u_0)$
2. Если в ходе итеративной процедуры норма градиента начинает увеличиваться (что означает, что процесс начинает расходиться в силу численных проблем), то выполняется рестарт.

При вычислении правдоподобия в точке a_{j+1} в качестве начальных приближений для u_0 , \mathbf{K}_0^{-1} и $\det(\mathbf{K}_0^{-1})$ предлагается брать значения в одной из предыдущих точек a_1, \dots, a_j , ближайшей к a_{j+1} в смысле некоторой метрики $\rho(a, a')$. Вид метрики зависит от вида ковариационной функции (2), (3). Например, для экспоненциальной ковариационной функции

$$K_0(x, x') = \sigma_0^2 \exp\left(-\sum_{i=1}^n \theta_i^2 (x_i - x'_i)^2\right),$$

где $a = \{\sigma_0^2, \theta_1^2, \dots, \theta_n^2, \sigma^2\}$, предлагается использовать следующую метрику:

$$\rho(a, a') = \sum_{i=1}^n (\exp(-\theta_i) - \exp(-\theta'_i))^2 + \left(\frac{\sigma^2}{\sigma_0^2} - \frac{\sigma'^2}{\sigma_0'^2}\right)^2.$$

Приведем весь алгоритм аппроксимации величин (7):

1. Инициализируем $\mathbf{K}_0^{-1} = \mathbf{C}^{-1}$, $u_0 = \alpha$, $\det(\mathbf{K}_0^{-1}) = d$ где \mathbf{C}^{-1} , α , d — величины \mathbf{K}^{-1} , $\mathbf{K}^{-1}\mathbf{y}$, $\det(\mathbf{K}^{-1})$, подсчитанные на одной из предыдущих итераций оптимизации (5) в точке, ближайшей к текущей в смысле метрики $\rho(a, a')$.

2. Если $\|\nabla g(u_0)\| > h$, то не пробуем аппроксимировать величины (7), а сразу вычисляем точно с помощью преобразования Холецкого.
3. Иначе для $j := 0, \dots, J-1$
 - (a) Если $\|\nabla g(u_j)\| \leq \frac{\epsilon}{N}$, то аппроксимация достаточно точная, выходим из цикла по j .
 - (b) Оптимальный шаг $\eta_j = -\frac{(\nabla g(u_j))^T \mathbf{K} \nabla g(u_j)}{(\nabla g(u_j))^T \mathbf{K}_j^{-1} \mathbf{K} \mathbf{K}_j^{-1} \nabla g(u_j)}$
 - (c) $u_{j+1} = u_j - \eta_j \mathbf{K}_j^{-1} \nabla g(u_j)$
 - (d) $\mathbf{K}_{j+1}^{-1} = \mathbf{K}_j^{-1} + \Delta \mathbf{K}_j^{-1}$, $\Delta \mathbf{K}_j^{-1}$ определяется по формуле (16).
 - (e) $\det(\mathbf{K}_{j+1}^{-1}) = \det(\mathbf{K}_j^{-1}) \frac{q_j^T p_j}{q_j^T \mathbf{K}_j^{-1} q_j}$, где q_j, p_j определены в (11).
 - (f) $m = \text{mean}(\text{diag}(\mathbf{K}_{j+1}^{-1} \times \mathbf{K}))$, где $\text{mean}(\mathbf{v}) = \frac{1}{n} \sum_{i=1}^n v_i$ — среднее значение вектора.
 - (g) $\mathbf{K}_{j+1}^{-1} := \frac{\mathbf{K}_{j+1}^{-1}}{m}$
 - (h) $\det(\mathbf{K}_{j+1}^{-1}) := \frac{\det(\mathbf{K}_{j+1}^{-1})}{m^N}$
 - (i) Если $\|\nabla g(u_{j+1})\| > \|\nabla g(u_j)\|$, то выходим из цикла по j и вычисляем (7) точно с помощью преобразования Холецкого.

5. Вычислительный эксперимент

В вычислительных экспериментах сравнивалось качество и время работы алгоритма с использованием аппроксимации, предложенной в данной работе, и алгоритма без аппроксимации, в котором величины (7) всегда вычисляются точно с помощью разложения Холецкого.

Для демонстрации экспериментальных результатов был использован большой набор тестовых функций, которые применяются для тестирования задач оптимизации [10], [11]. Размерности функций — от 2 до 8. Всего тестирование проводилось на 211 различных функциях, для каждой из которых генерировались случайные обучающие выборки размером 320 и 1000 точек. Результаты сравнивались по времени обучения и по среднеквадратичной ошибке (1) на больших контрольных выборках из 10000 точек. Для удобства результаты представлены в виде кривых Долан-Мора [2] на рисунках 1, 2, 3, 4. Чем выше кривая находится на графике, тем выше качество или меньше время работы соответствующего алгоритма. Видно, что точность работы двух алгоритмов практически не отличается, а время обучения алгоритма, использующего предложенную аппроксимацию, значительно меньше: ускорение до 8 раз при размере выборки 320 и до 16 раз при размере выборки 1000.

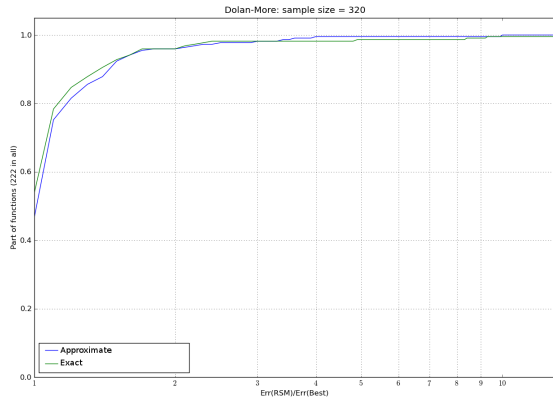


Рис. 1. Размер выборки 320, точность аппроксимации на тестовой выборке

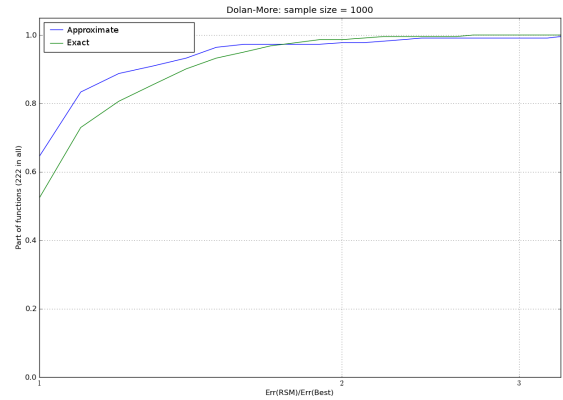


Рис. 3. Размер выборки 1000, точность аппроксимации на тестовой выборке

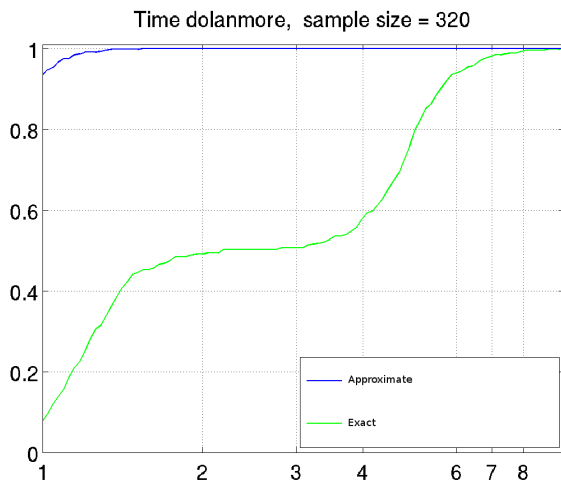


Рис. 2. Размер выборки 320, время обучения

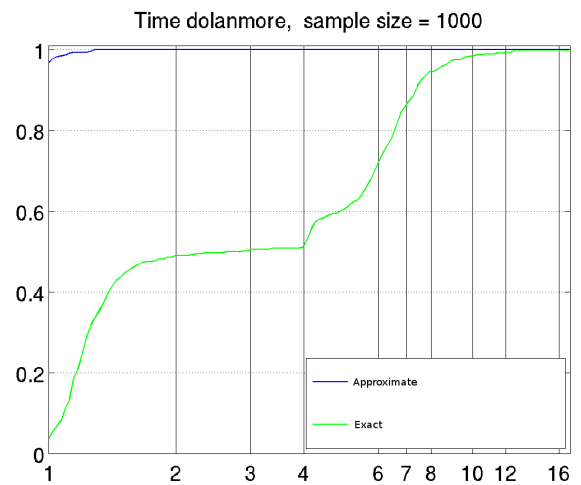


Рис. 4. Размер выборки 1000, время обучения

6. Выводы

Вычислительный эксперимент показал, что предложенный алгоритм дает значительное ускорение по времени обучения без потери точности аппроксимации. Предложенный метод позволяет избежать громоздких схем для аппроксимации $\log(\det(\mathbf{K}))$, предложенных в [6], [8]. Полученное значение аппроксимации детерминанта в предложенном алгоритме в точности совпадает со значением детерминанта аппроксимации обратной матрицы:

$$\tilde{d} = \frac{1}{\det(\tilde{\mathbf{K}}^{-1})},$$

где \tilde{d} — аппроксимация величины $\det(\mathbf{K})$, $\tilde{\mathbf{K}}^{-1}$ — аппроксимация величины \mathbf{K}^{-1} . Т.к. качество аппроксимации $\tilde{\mathbf{K}}^{-1}$ контролируется и при необходимости производятся рестарты, то автоматически контролируется качество аппроксимации \tilde{d} , в отличие от

методов из статей других авторов.

Список литературы

- [1] A. S. Alexander I. J. Forrester and A. J. Keane. *Engineering Design via Surrogate Modelling: A Practical Guide*. Wiley, 2008.
- [2] E. Dolan and J. Moré. Benchmarking optimization software with performance profiles. *Mathematical Programming*, 91(2):201–213, 2002.
- [3] M. Gibbs and D. MacKay. Efficient implementation of gaussian processes. 1997.
- [4] S. J. W. Jorge Nocedal. *Numerical Optimization*. Springer, 1999.
- [5] W. Leithead and Y. Zhang. $O(n^2)$ -operation approximation of covariance matrix inverse in gaussian process regression based on quasi-newton bfgs method. *Communications in Statistics—Simulation and Computation*, 36(2):367–380, 2007.
- [6] W. Leithead, Y. Zhang, and D. Leith. Efficient

- gaussian process based on bfgs updating and logdet approximation. In *the 16th IFAC world congress*, 2005.
- [7] C. E. Rasmussen and C. K. I. Williams. *Gaussian processes for machine learning*. The MIT Press, 2006.
- [8] Y. Zhang and W. Leithead. Approximate implementation of the logarithm of the matrix determinant in gaussian process regression. *journal of Statistical Computation and Simulation*, 77(4):329–348, 2007.
- [9] Y. Zhang, W. Leithead, and D. Leith. Random scaling of quasi-newton bfgs method to improve the $O(n^2)$ -operation approximation of covariance-matrix inverse in gaussian process. In *Intelligent Control, 2007. ISIC 2007. IEEE 22nd International Symposium on*, pages 452–457. IEEE, 2007.
- [10] GDR MASCOT-NUM Toy Functions benchmark. <http://gdr-mascotnum.math.cnrs.fr/data2/benchmarks/jm.pdf>
- [11] Lappeenranta University of Technology: evolutionary computation pages - the function testbed. <http://www.it.lut.fi/ip/evo/functions/functions.html>