

# Домашнее задание №4 по курсу «Математическая Статистика в Машинном Обучении»

Федор Ерин

## Задача 1 [3 балла]

Пусть данные порождаются моделью

$$t = y(\mathbf{x}, \mathbf{w}) + \varepsilon,$$

где  $\varepsilon \sim \mathcal{N}(\varepsilon|0, \beta^{-1})$ ,  $y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \varphi(\mathbf{x})$ ,  $\varphi(\mathbf{x}) = (\varphi_1(\mathbf{x}), \dots, \varphi_M(\mathbf{x}))$  — некоторый набор функций.

Также обозначим через  $\mathbf{t} = (t_1, \dots, t_N)$  — известный набор целевых значений, через  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$  — соответствующие значения векторов признаков,  $\Phi = (\varphi_j(\mathbf{x}_i)), i = 1, \dots, N, j = 1, \dots, M$ .

Предположим, что априорное распределение  $p(\mathbf{w}|\alpha) = N(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})$ . В таком случае легко подсчитать, что апостериорное распределение  $p(\mathbf{w}|\mathbf{t}, \alpha) = N(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N)$ , где

$$\mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{t},$$

$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi.$$

Апостериорное распределение выходного значения  $t$  в новой точке  $\mathbf{x}$  равно

$$p(t|\mathbf{x}, \mathbf{t}, \alpha, \beta) = \int p(t|\mathbf{w}, \beta) p(\mathbf{w}|\mathbf{t}, \alpha, \beta) d\mathbf{w} = N(t|\mathbf{m}_N^T \varphi(\mathbf{x}), \sigma_N^2(\mathbf{x})),$$

где

$$\sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \varphi(\mathbf{x})^T \mathbf{S}_N \varphi(\mathbf{x}).$$

Используя равенство

$$(\mathbf{M} + \mathbf{v}\mathbf{v}^T)^{-1} = \mathbf{M}^{-1} - \frac{(\mathbf{M}^{-1}\mathbf{v})(\mathbf{v}^T \mathbf{M}^{-1})}{1 + \mathbf{v}^T \mathbf{M}^{-1} \mathbf{v}}, \quad (1)$$

показать, что

$$\sigma_{N+1}^2(\mathbf{x}) \leq \sigma_N^2(\mathbf{x}).$$

## Решение

- Найдем  $\mathbf{S}_{N+1}^{-1}$ , для этого рассмотрим новое наблюдение  $(\mathbf{x}_{N+1}, t_{N+1})$  и выпишем функцию правдоподобия:

$$p(t_{N+1}|\mathbf{x}_{N+1}, \mathbf{w}) = N(t_{N+1}|y(\mathbf{x}_{N+1}, \mathbf{w}), \beta^{-1})$$

Распишем показатель экспоненты в апостериорном распределении:

$$\begin{aligned} & (\mathbf{w} - \mathbf{m}_N)^T \mathbf{S}_N^{-1} (\mathbf{w} - \mathbf{m}_N) + \beta(t_{N+1} - \mathbf{w}^T \varphi(\mathbf{x}_{N+1}))^2 = \\ & = \mathbf{w}^T (\mathbf{S}_N^{-1} + \beta \varphi(\mathbf{x}_{N+1}) \varphi(\mathbf{x}_{N+1})^T) \mathbf{w} - 2\mathbf{w}^T (\mathbf{S}_N^{-1} \mathbf{m}_N + \beta \varphi(\mathbf{x}_{N+1}) t_{N+1}) + const \end{aligned}$$

Отсюда получили:

$$\mathbf{S}_{N+1}^{-1} = \mathbf{S}_N^{-1} + \beta \varphi(\mathbf{x}_{N+1}) \varphi(\mathbf{x}_{N+1})^T$$

Воспользуемся равенством (1) из условия:

$$\begin{aligned} \mathbf{S}_{N+1} &= (\mathbf{S}_N^{-1} + \beta \varphi(\mathbf{x}_{N+1}) \varphi(\mathbf{x}_{N+1})^T)^{-1} = \\ &= \mathbf{S}_N - \frac{(\mathbf{S}_N \sqrt{\beta} \varphi(\mathbf{x}_{N+1})) (\sqrt{\beta} \varphi(\mathbf{x}_{N+1})^T \mathbf{S}_N)}{1 + \sqrt{\beta} \varphi(\mathbf{x}_{N+1})^T \mathbf{S}_N \sqrt{\beta} \varphi(\mathbf{x}_{N+1})} = \mathbf{S}_N - \frac{\beta \mathbf{S}_N \varphi(\mathbf{x}_{N+1}) \varphi(\mathbf{x}_{N+1})^T \mathbf{S}_N}{1 + \beta \varphi(\mathbf{x}_{N+1})^T \mathbf{S}_N \varphi(\mathbf{x}_{N+1})} \end{aligned}$$

Теперь можем рассмотреть разность:

$$\begin{aligned} \sigma_N^2(\mathbf{x}) - \sigma_{N+1}^2(\mathbf{x}) &= \varphi(\mathbf{x}) (\mathbf{S}_N - \mathbf{S}_{N+1}) \varphi(\mathbf{x}) = \\ &= \varphi(\mathbf{x}) \frac{\beta \mathbf{S}_N \varphi(\mathbf{x}_{N+1}) \varphi(\mathbf{x}_{N+1})^T \mathbf{S}_N}{1 + \beta \varphi(\mathbf{x}_{N+1})^T \mathbf{S}_N \varphi(\mathbf{x}_{N+1})} \varphi(\mathbf{x}) = \frac{\beta (\varphi(\mathbf{x})^T \mathbf{S}_N \varphi(\mathbf{x}_{N+1}))^2}{1 + \beta \varphi(\mathbf{x}_{N+1})^T \mathbf{S}_N \varphi(\mathbf{x}_{N+1})} \geq 0, \end{aligned}$$

так как  $\mathbf{S}_N$  положительно определена.

## Задача 2 [5 баллов]

Рассмотрим модель регрессии

$$t = \mathbf{x}^T \mathbf{w} + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \beta^{-1}),$$

где  $p(\mathbf{w}|\beta) = \mathcal{N}(\mathbf{w}|\mathbf{w}_0, \beta^{-1}\mathbf{S}_0)$ ,  $p(\beta) = \text{Gamma}(\beta|a_0, b_0)$ , т.е. совместное априорное распределение  $(\mathbf{w}, \beta)$  имеет вид:

$$p(\mathbf{w}, \beta) = \mathcal{N}(\mathbf{w}|\mathbf{w}_0, \beta^{-1}\mathbf{S}_0) \text{Gamma}(\beta|a_0, b_0).$$

Покажите, что апостериорное распределение  $(\mathbf{w}, \beta)$  после наблюдения выборки  $(\mathbf{X}, \mathbf{t}) = \{(\mathbf{x}_i, t_i)\}_{i=1}^n$  имеет вид

$$p(\mathbf{w}, \beta|\mathbf{X}, \mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{w}_n, \beta^{-1}\mathbf{S}_n) \text{Gamma}(\beta|a_n, b_n).$$

Найдите параметры  $\mathbf{w}_n$ ,  $\mathbf{S}_n$ ,  $a_n$ ,  $b_n$ .

*Примечание. Плотность гамма-распределения  $\text{Gamma}(a, b)$  имеет вид*

$$\text{Gamma}(x|a, b) = \frac{x^{a-1}}{\Gamma(a)b^a} e^{-\frac{x}{b}}.$$

### Решение

- Распишем априорную плотность и функцию правдоподобия:

$$p(\mathbf{w}, \beta) = \mathcal{N}(\mathbf{w}|\mathbf{w}_0, \beta^{-1}\mathbf{S}_0) \text{Gamma}(\beta|a_0, b_0) \sim \frac{\beta^2}{\mathbf{S}_0^2} \exp\left(-\frac{1}{2}(\mathbf{w} - \mathbf{w}_0)^T \beta \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{w}_0)\right) b_0^{a_0} \beta^{a_0-1} \exp(-b_0 \beta)$$

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{i=1}^n \mathcal{N}(t_i|\mathbf{w}^T \varphi(\mathbf{x}_i), \beta^{-1}) \sim \prod_{i=1}^n \sqrt{\beta} \cdot \exp\left(-\frac{\beta}{2}(t_i - \mathbf{w}^T \varphi(\mathbf{x}_i))^2\right)$$

Известно, что

$$p(\mathbf{w}, \beta|\mathbf{t}) \sim p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) \cdot p(\mathbf{w}, \beta)$$

Для нахождения  $\mathbf{S}_n$  распишем только квадратичные слагаемые под экспонентой:

$$-\frac{\beta}{2} \mathbf{w}^T \mathbf{S}_0^{-1} \mathbf{w} + \sum_{i=1}^n -\frac{\beta}{2} \mathbf{w}^T \varphi(\mathbf{x}_i) \varphi(\mathbf{x}_i)^T \mathbf{w} = -\frac{\beta}{2} \mathbf{w}^T (\mathbf{S}_0^{-1} + \sum_{i=1}^n \varphi(\mathbf{x}_i) \varphi(\mathbf{x}_i)^T) \mathbf{w}$$

Отсюда получаем:

$$\mathbf{S}_n^{-1} = \mathbf{S}_0^{-1} + \sum_{n=1}^N \varphi(\mathbf{x}_n) \varphi(\mathbf{x}_n)^T \Rightarrow \mathbf{S}_n = (\mathbf{S}_0^{-1} + \sum_{n=1}^N \varphi(\mathbf{x}_n) \varphi(\mathbf{x}_n)^T)^{-1}$$

Для нахождения  $\mathbf{w}_n$  распишем линейные слагаемые:

$$\beta \mathbf{w}_0^T \mathbf{S}_0^{-1} \mathbf{w} + \sum_{i=1}^n \beta t_i \varphi(\mathbf{x}_i)^T \mathbf{w} = \beta (\mathbf{w}_0^T \mathbf{S}_0^{-1} + \sum_{i=1}^n t_i \varphi(\mathbf{x}_i)^T) \mathbf{w}$$

С другой стороны:

$$\mathbf{w}_n^T \mathbf{S}_n^{-1} = \mathbf{w}_0^T \mathbf{S}_0^{-1} + \sum_{i=1}^n t_i \varphi(\mathbf{x}_i)^T$$

Откуда получаем:

$$\mathbf{w}_n = \mathbf{S}_n (\mathbf{S}_0^{-1} \mathbf{w}_0 + \sum_{i=1}^n t_i \varphi(\mathbf{x}_i))$$

Для нахождения  $b_n$  распишем константные слагаемые:

$$\left(-\frac{\beta}{2} \mathbf{w}_0^T \mathbf{S}_0^{-1} \mathbf{w}_0 - b_0 \beta\right) - \frac{\beta}{2} \sum_{i=1}^n t_i^2 = -\beta \left(\frac{1}{2} \mathbf{w}_0^T \mathbf{S}_0^{-1} \mathbf{w}_0 + b_0 + \frac{1}{2} \sum_{i=1}^n t_i^2\right)$$

При этом:

$$\frac{1}{2} \mathbf{w}_n^T \mathbf{S}_n^{-1} \mathbf{w}_n + b_n = \frac{1}{2} \mathbf{w}_0^T \mathbf{S}_0^{-1} \mathbf{w}_0 + b_0 + \frac{1}{2} \sum_{i=1}^n t_i^2$$

Откуда получаем:

$$b_n = b_0 + \frac{1}{2} (\mathbf{w}_0^T \mathbf{S}_0^{-1} \mathbf{w}_0 - \mathbf{w}_n^T \mathbf{S}_n^{-1} \mathbf{w}_n) + \sum_{i=1}^n t_i^2$$

Наконец, найдем  $a_n$ :

$$2 + a_n - 1 = (2 + a_0 - 1) + \frac{N}{2}$$

$$a_n = a_0 + \frac{N}{2}$$

**Ответ:**

- $w_n = S_n(S_0^{-1}w_0 + \sum_{i=1}^n t_i \varphi(x_i))$ ,
- $S_n = (S_0^{-1} + \sum_{i=1}^n \varphi(x_n)\varphi(x_i)^T)^{-1}$ ,
- $a_n = a_0 + N/2$ ,
- $b_n = b_0 + \frac{1}{2}(w_0^T S_0^{-1}w_0 - w_n^T S_n^{-1}w_n + \sum_{i=1}^n t_i^2)$ .

### Задача 3 [4 балла]

Пусть дана выборка  $(X, t)$ ,  $X = \{x_1, \dots, x_n\}$ . Предположим, что наблюдаемые данные представляют собой зашумленные значения гауссовского случайного процесса  $f(x)$ , т.е. имеет место следующая модель:

$$t(x) = f(x) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2),$$

где  $f(x)$  — стационарный гауссовский процесс с нулевым средним и функцией ковариации  $K(x', x'')$ . Через  $\hat{K}(x', x'')$  обозначим функцию ковариации зашумленного гауссовского процесса  $t(x)$ :

$$\hat{K}(x', x'') = K(x', x'') + \sigma^2 \delta_{x', x''}.$$

Оказалось, что в выборке  $X$  все индексующие параметры идентичны:  $x_1 = \dots = x_n = x$ . Найдите мат. ожидание и дисперсию прогноза в произвольной точке  $y$ . Чему равен прогноз в той же самой точке  $x$ , при условии, что  $f(x) = t$  — истинное значение рассматриваемой реализации в точке  $x$ ? Чему равна дисперсия прогноза? Является ли оценка на  $f(x)$  смещенной?

**Замечание.** К объяснению понятия “переизмерения значения в конкретной точке  $x$ ” можно подходить несколькими способами. Формальный способ предполагает обобщение функции ковариации: рассмотрим  $i$ -ое измерение в точке  $x'$  и  $j$ -ое измерение в точке  $x''$ . Тогда

$$\text{Cov}(t(x', i), t(x'', j)) = \hat{K}(x', x'', i, j) = K(x', x'') + \sigma^2 \delta_{x', x''} \delta_{i, j}.$$

Неформальный способ состоит в том, что измерения происходят в точках  $x_1, \dots, x_n$ , очень близких к  $x$ , но все же не совпадающих с ней. В результате значения шумов в этих точках независимы, а истинные значения целевого гауссовского процесса почти совпадают (в предположении непрерывности функции ковариации  $K(x', x'')$ ).

**Подсказка.** Для обращения матрицы вида  $\alpha I + E$ , где  $I \in \mathbb{R}^{n \times n}$  — единичная матрица, и  $E \in \mathbb{R}^{n \times n}$  — матрица из единиц, можно воспользоваться тождеством Шермана-Моррисона-Вудберри, положив  $E = \mathbf{1} \cdot \mathbf{1}^T$ ,  $\mathbf{1} = (1, 1, \dots, 1)^T \in \mathbb{R}^n$ .

### Задача 4 [5 баллов]

матрим регрессию на основе гауссовских процессов. Пусть нам дана выборка  $(x, t) = \{(x_i, t_i) : x_i, t_i \in \mathbb{R}\}_{i=1}^n$ , где  $x_i = x_1 + (i-1)h$ ,  $h > 0$ , т.е.  $n$  точек расположены равномерно на вещественной оси. Предположим, что выборка  $(x, t)$  представляет собой реализацию некоторого гауссовского случайного процесса  $f(\cdot)$ , т.е.

$$t_i = f(x_i).$$

Будем считать, что шума в наблюдениях нет. Пусть ковариационная функция  $K(\cdot, \cdot)$  имеет вид:

$$K(x', x'') = \alpha \exp\left(-\frac{|x' - x''|}{\gamma}\right),$$

где  $\alpha$  и  $\gamma$  — известные параметры. Найдите апостериорное среднее и дисперсию гауссовского случайного процесса в точке  $x_* \geq x_n$ .

### Задача 5 [3 балла]

Рассмотрим задачу построения адаптивного дизайна на основе гауссовских процессов. Пусть  $(\mathbf{X}, \mathbf{t}) = \{(\mathbf{x}_i, t_i) : \mathbf{x}_i \in \mathbb{X}, t_i \in \mathbb{R}\}$ , где  $\{\mathbf{x}_i\}_{i=1}^n$  — известные точки дизайна из рассматриваемого множества  $\mathbb{X} \subset \mathbb{R}^d$ , и  $\{t_i\}_{i=1}^n$  — измеренные значения аппроксимируемой зависимости. Одним из критериев выбора следующей точки дизайна является *критерий максимальной дисперсии*:

$$\mathcal{I}_{\text{MV}}(\mathbf{x}) \triangleq \hat{\sigma}^2(\mathbf{x}|\mathbf{X}),$$

$$\mathbf{x}_* = \arg \max_{\mathbf{x} \in \mathbb{X}} \mathcal{I}_{\text{MV}}(\mathbf{x}) = \arg \max_{\mathbf{x} \in \mathbb{X}} \hat{\sigma}^2(\mathbf{x}|\mathbf{X}),$$

где  $\hat{\sigma}^2(\mathbf{x}|\mathbf{X})$  — апостериорная дисперсия в точке  $\mathbf{x} \in \mathbb{X}$  при заданном дизайне  $\mathbf{X}$  (заметим, что в случае гауссовских процессов  $\hat{\sigma}^2(\mathbf{x}|\mathbf{X})$  не зависит от  $\mathbf{t}$ ). Такой критерий легко вычислим и включает в себя информацию о поведении истинной зависимости, однако учитывает только локальное поведение и склонен выдавать точки близкие к границе множества  $\mathbb{X}$ . Поэтому более разумным является *критерий для минимизации ожидаемой среднеквадратичной ошибки аппроксимации на следующей итерации*:

$$\mathcal{I}_{\rho_2}(\mathbf{x}) \triangleq \frac{1}{|\mathbb{X}|} \int_{\mathbb{X}} (\hat{\sigma}^2(\mathbf{v}|\mathbf{X}) - \hat{\sigma}^2(\mathbf{v}|\mathbf{X} \cup \mathbf{x})) d\mathbf{v},$$

$$\mathbf{x}_* = \arg \min_{\mathbf{x} \in \mathbb{X}} \mathcal{I}_{\rho_2}(\mathbf{x}),$$

где  $\hat{\sigma}^2(\mathbf{v}|\mathbf{X})$  и  $\hat{\sigma}^2(\mathbf{v}|\mathbf{X} \cup \mathbf{x})$  — апостериорные дисперсии в точке  $\mathbf{v}$  на дизайнах  $\mathbf{X}$  и  $\mathbf{X} \cup \mathbf{x}$  соответственно. Покажите, что критерий  $\mathcal{I}_{\rho_2}(\mathbf{x})$  может быть записан в виде:

$$\mathcal{I}_{\rho_2}(\mathbf{x}) = \frac{1}{|\mathbb{X}|} \int_{\mathbb{X}} \frac{\hat{K}^2(\mathbf{v}, \mathbf{x})}{\hat{\sigma}^2(\mathbf{x}|\mathbf{X})} d\mathbf{v},$$

где  $\hat{K}(\mathbf{v}, \mathbf{x}) = K(\mathbf{v}, \mathbf{x}) - \mathbf{k}^T(\mathbf{v})\Sigma_{\mathbf{X}}^{-1}\mathbf{k}(\mathbf{x})$ , и

$$\Sigma_{\mathbf{X}} = \begin{pmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & \dots & K(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ K(\mathbf{x}_n, \mathbf{x}_1) & \dots & K(\mathbf{x}_n, \mathbf{x}_n) \end{pmatrix}.$$

### Задача 6 [5 баллов]

Допустим, что в нашем распоряжении имеется выборка  $(\mathbf{X}, \mathbf{Y})$ , где  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{n \times d}$  и  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_n]^T \in \mathbb{R}^{n \times k}$ , т.е. каждому объекту  $\mathbf{x} \in \mathbb{R}^d$  соответствует вектор  $\mathbf{y} \in \mathbb{R}^k$ . Требуется по новому вектору  $\mathbf{x}^*$  предсказать вектор  $\mathbf{y}^*$ . Ниже рассматривается алгоритм, основанный на использовании PCA.

**Этап обучения.** Для объединенной матрицы  $\mathbf{Z} = [\mathbf{X}, \mathbf{Y}] \in \mathbb{R}^{n \times (d+k)}$  находим первые  $r$  (*параметр алгоритма*) главных компонент  $\mathbf{u}_1, \dots, \mathbf{u}_r$ , где  $\mathbf{u}_i \in \mathbb{R}^{d+k}$ . Расположим эти векторы в качестве строк матрицы  $\mathbf{U}$ :

$$\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_r] \in \mathbb{R}^{(d+k) \times r}$$

Каждый из векторов  $\{\mathbf{u}_i\}_{i=1}^r$  представляется как объединение двух подвекторов:  $\mathbf{u}_i^T = (\mathbf{u}_{X,i}, \mathbf{u}_{Y,i})$ , где  $\mathbf{u}_{X,i} \in \mathbb{R}^d$  соответствует пространству признаков, а  $\mathbf{u}_{Y,i} \in \mathbb{R}^k$  — пространству целевых переменных.

Введем обозначения

$$\langle \mathbf{X} \rangle \triangleq \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \in \mathbb{R}^d, \quad \langle \mathbf{Y} \rangle \triangleq \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i \in \mathbb{R}^k, \quad \langle \mathbf{Z} \rangle \triangleq \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \in \mathbb{R}^{d+k}.$$

В таких обозначениях

$$\langle \mathbf{Z} \rangle = \begin{pmatrix} \langle \mathbf{X} \rangle \\ \langle \mathbf{Y} \rangle \end{pmatrix} \in \mathbb{R}^{d+k}.$$

- Преобразование произвольного объединенного вектора  $\mathbf{z} = (\mathbf{x}, \mathbf{y})$  в сжатое описание  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_d)$  происходит согласно следующей формуле:

$$\lambda_i = \mathbf{u}_i^T (\mathbf{z} - \langle \mathbf{Z} \rangle), \quad i = 1, \dots, d,$$

или в векторном виде

$$\boldsymbol{\lambda} = \mathbf{U}^T (\mathbf{z} - \langle \mathbf{Z} \rangle).$$

- Восстановление  $\mathbf{z} = (\mathbf{x}, \mathbf{y})$  по его сжатому описанию  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_r)$  происходит следующим образом:

$$\tilde{\mathbf{x}}(\boldsymbol{\lambda}) = \langle \mathbf{X} \rangle + \sum_{i=1}^r \lambda_i \mathbf{u}_{X,i} = \langle \mathbf{X} \rangle + \mathbf{U}_X \boldsymbol{\lambda},$$

$$\tilde{\mathbf{y}}(\boldsymbol{\lambda}) = \langle \mathbf{Y} \rangle + \sum_{i=1}^r \lambda_i \mathbf{u}_{Y,i} = \langle \mathbf{Y} \rangle + \mathbf{U}_Y \boldsymbol{\lambda}.$$

**Этап предсказания.** Теперь опишем алгоритм восстановления  $\mathbf{y}$  по признакам  $\mathbf{x}$ :

1. По заданному  $\mathbf{x}$  определяется вектор коэффициентов  $\boldsymbol{\lambda}^*$ :

$$\boldsymbol{\lambda}^*(\mathbf{x}) = \arg \min_{\boldsymbol{\lambda} \in \mathbb{R}^r} \|\tilde{\mathbf{x}}(\boldsymbol{\lambda}) - \mathbf{x}\|_2^2.$$

2. По вычисленному вектору  $\boldsymbol{\lambda}^*$  оценивается вектор целевых переменных по формуле

$$\mathbf{y}^*(\mathbf{x}) = \langle \mathbf{Y} \rangle + \mathbf{U}_Y \boldsymbol{\lambda}^*(\mathbf{x}).$$

**Задание.** Предполагается, что известны обучающая выборка  $\mathbf{X}$ ,  $\mathbf{Y}$  и матрица главных компонент  $\mathbf{U}$ .

- Найдите формулы для коэффициентов  $\boldsymbol{\lambda}^*$ , подсчитанных по вектору признаков  $\mathbf{x}$ .
- Формулы для  $\mathbf{y}$ , восстановленному по  $\mathbf{x}$  согласно описанной выше процедуре.

*Внимание. Ответы будут отличаться для случаев когда  $r \leq d$  и  $r > d$ . В решении следует предусмотреть оба случая. Для простоты при решении следует полагать, что  $\text{rg} \mathbf{Z} = d + k$ .*

### Задача 7 [10 баллов]

В данной задаче требуется реализовать подбор параметров алгоритма классификации с помощью байесовской оптимизации на основе гауссовских процессов.

- Создайте искусственный датасет для классификации. Возьмите 15 информативных признаков и 5 шумовых. Количество сэмплов возьмите от 1000 и выше: 500 сэмплов отведите на обучающую выборку, а оставшиеся на тестовую.
- Используйте Kernel SVM с RBF-ядром для того, чтобы обучиться на данных. Для параметров  $C$  и  $\gamma$  в некоторой сетке (например, в сетке `np.logspace(-4, 3, 71) × np.logspace(-4, 3, 71)`) подсчитайте значения ROC-AUC на тестовой выборке (не забудьте использовать для подсчета ROC AUC метод `decision_function`, а не `predict`). Отрисуйте в виде `heat map`-ы значения ROC AUC (см. для примера код с семинара или рис. 1 ниже). В качестве альтернативного подхода можно провести `grid search` по сетке с помощью класса `GridSearchCV`, который хранит в себе значения качества во всех подсчитанных узлах. Это сэкономит время, так как для ускорения можно передать параметр `jobs = -1`. Но при указанных выше параметрах подсчет в один поток не занимает много времени.
- Реализуйте следующий алгоритм выбора параметров  $(C, \gamma)$  с помощью байесовской оптимизации:
  1. На вход алгоритма подается множество пар значений  $C$  и  $\gamma$ , среди которых будет проводится поиск наилучших параметров. Для простоты пусть данные параметры образуют сетку `np.logspace(-4, 3, 701) × np.logspace(-4, 3, 701)`. Стоит заметить, что эта сетка гораздо плотнее построенной в предыдущем пункте.
  2. На первом шаге просэмплируйте из сетки 5-10 пар значений  $(C, \gamma)$ , на которых явно подсчитайте ROC AUC.
  3. Обучите гауссовскую регрессию на полученных данных. Подсчитайте дисперсии предсказаний на всем множестве значений пар  $(C, \gamma)$ , поданных на вход алгоритма. Выберите параметры  $C_*$  и  $\gamma_*$ , дисперсия прогноза ROC AUC для которых максимальна.
  4. Подсчитайте ROC AUC для  $C_*$  и  $\gamma_*$  и добавьте ее в выборку известных соответствий  $(C, \gamma)$  и ROC AUC.
  5. Далее вновь обучите регрессию на обновленных данных и т.д. Повторяем процесс итеративно до тех пор, пока не будет достигнуто заданное количество итераций `max_iter`, либо дисперсия прогноза ROC AUC в каждой точке не станет меньше некоторого порога `var_threshold` (параметры `max_iter` и `var_threshold` стоит подобрать так, чтобы получить более-менее хорошее качество).
  6. В тот момент, когда подсчеты остановились, в качестве  $(C, \gamma)$  выдаем точку с максимальным ROC AUC.
- [10 баллов] Отрисуйте все полученные в процессе оптимизации точки на `heatmap`-е с ROC AUC. Выделите цветом точку с лучшим качеством. Получилось ли с помощью байесовской оптимизации добиться того же качества, что и с помощью поиска по сетке? Помните, что сетка параметров для байесовской оптимизации плотнее, поэтому хотя чисто теоретически результат и может получаться один в один, но на практике это далеко не факт. К тому же в данной задаче не требуется получить результат, лучший, чем поиск по сетке. Требуется лишь показать, что этот результат приближается к результату, полученному по сетке.

Используйте ядро следующего вида

$$K(\mathbf{x}, \mathbf{x}') = \theta_1 \exp \left\{ -\frac{1}{2} \left( \frac{(x_1 - x'_1)^2}{2\eta_1} + \frac{(x_2 - x'_2)^2}{2\eta_2} \right) \right\} + \theta_2 \delta_{\mathbf{x}, \mathbf{x}'}$$

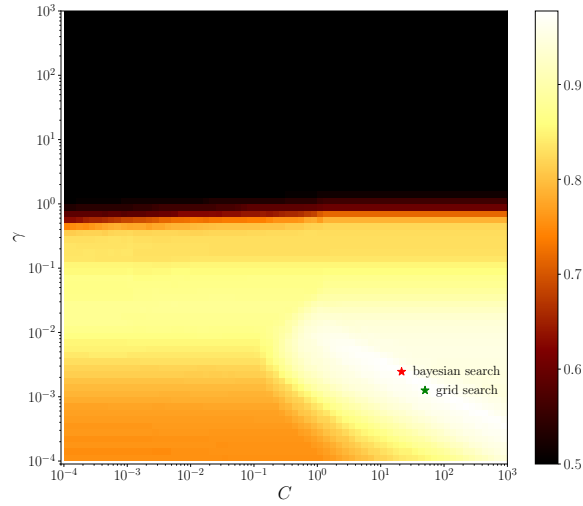


Рис. 1:

Так как на каждой итерации алгоритма выборка увеличивается, то вполне разумно переобучать параметры  $\theta_1$ ,  $\theta_2$ ,  $\eta_1$ ,  $\eta_2$  перед тем, как выбирать очередную точку  $\mathbf{x} = (C, \gamma)$  для подсчета ROC AUC. В данной задаче это допустимо, так как количество точек и размерность невелики. В практических задачах можно просто перестать оптимизировать данные параметры как только их изменение от итерации к итерации становится незначительным.

**Замечание.** Конкретные числа (число сэмплов, размер сетки, порог `var_threshold`) остаются на ваше усмотрение.

**Замечание.** Код должен запускаться от `u` до `do`, чтобы была возможность полностью воспроизвести эксперимент. Ну и чем лучше логическое разделение кода, тем проще его проверять.