

Домашнее задание №3 по курсу «Математическая Статистика в Машинном Обучении»

Школа Анализа Данных

Задачи

Теоретический блок

Задача 1 [1 балл]

Рассмотрим задачу оптимизации, решаемую в Elastic Net:

$$\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda_1 \|\mathbf{w}\|_1 + \lambda_2 \|\mathbf{w}\|_2^2 \rightarrow \min_{\mathbf{w}}$$

Покажем, что данную задачу оптимизации можно преобразовать к виду, содержащему составляющую только для ℓ_1 -регуляризации. Для этого определим искусственные данные

$$\mathbf{X}^* = \frac{1}{\sqrt{1+\lambda_2}} \begin{pmatrix} \mathbf{X} \\ \alpha \mathbf{I} \end{pmatrix} \in \mathbb{R}^{(n+d) \times d}, \quad \mathbf{y}^* = \begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix} \in \mathbb{R}^{n+d},$$

Покажите, что при таком преобразовании задачу выше можно свести к задаче

$$\|\mathbf{y}^* - \mathbf{X}^* \mathbf{w}^*\|_2^2 + \gamma \|\mathbf{w}^*\|_1 \rightarrow \min_{\mathbf{w}^*}.$$

Найдите требуемые для этого значения α и γ , выразив последние через λ_1 и λ_2 . Найдите связь между решениями оптимизационных задач $\hat{\mathbf{w}}$ и $\hat{\mathbf{w}}^*$.

Решение:

$$(\mathbf{y}^*)^T \mathbf{y}^* = \mathbf{y}^T \mathbf{y}$$

$$(\mathbf{X}^*)^T \mathbf{X}^* = \frac{1}{1+\lambda_2} (\mathbf{X}^T \mathbf{X} + \alpha^2 \mathbf{I}) \Rightarrow \mathbf{X}^T \mathbf{X} = (1+\lambda_2)(\mathbf{X}^*)^T \mathbf{X}^* - \alpha^2 \mathbf{I}$$

$$(\mathbf{y}^*)^T \mathbf{X}^* = \mathbf{y}^T \mathbf{X} \frac{1}{\sqrt{1+\lambda_2}} \Rightarrow \mathbf{y}^T \mathbf{X} = \sqrt{1+\lambda_2} (\mathbf{y}^*)^T \mathbf{X}^*$$

$$\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda_1 \|\mathbf{w}\|_1 + \lambda_2 \|\mathbf{w}\|_2^2 = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda_1 \|\mathbf{w}\|_1 + \lambda_2 \mathbf{w}^T \mathbf{w} =$$

$$= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X}\mathbf{w} + \mathbf{w}^T (\mathbf{X}^T \mathbf{X} + \lambda_2 \mathbf{I}) \mathbf{w} + \lambda_1 \|\mathbf{w}\|_1 =$$

$$= (\mathbf{y}^*)^T \mathbf{y}^* - 2(\mathbf{y}^*)^T \mathbf{X}^* \sqrt{1+\lambda_2} \mathbf{w} + \mathbf{w}^T ((\mathbf{X}^*)^T \mathbf{X}^* (1+\lambda_2) - \alpha^2 \mathbf{I} + \lambda_2 \mathbf{I}) \mathbf{w} + \lambda_1 \|\mathbf{w}\|_1$$

Обозначив $\mathbf{w}^* = \mathbf{w} \sqrt{1+\lambda_2}$ и приравняв $\alpha^2 = \lambda_2$ (чтобы избавиться от ℓ_2 регуляризации), получим нужную нам форму:

$$(\mathbf{y}^*)^T \mathbf{y}^* - 2(\mathbf{y}^*)^T \mathbf{X}^* \mathbf{w}^* + (\mathbf{X}^* \mathbf{w}^*)^T \mathbf{X}^* \mathbf{w}^* + \frac{\lambda_1}{\sqrt{1+\lambda_2}} \|\mathbf{w}^*\|_1 = \|\mathbf{y}^* - \mathbf{X}^* \mathbf{w}^*\|_2^2 + \frac{\lambda_1}{\sqrt{1+\lambda_2}} \|\mathbf{w}^*\|_1$$

Ответ: $\alpha = \sqrt{\lambda_2}, \gamma = \frac{\lambda_1}{\sqrt{1+\lambda_2}}$

Задача 2 [1 балл]

Пусть дана обучающая выборка $\{(X, y): X \in \mathbb{R}^{n \times d}, y \in \mathbb{R}^n\}$. Предположим, что справедлива следующая модель линейной регрессии:

$$y = x^T w + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

Покажите, что модель с наибольшим значением AIC является моделью с наименьшим значением статистики Mallows C_p .

Решение:

- Рассмотрим AIC: $\ell_S - |S| \rightarrow \max$

$$h = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(y_i - x_i^T w)^2} = \frac{1}{(2\pi\sigma^2)^{n/2}} e^{-\frac{1}{2\sigma^2} \|y - Xw\|_2^2}$$

$$\ell_S = \log h = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|y - Xw\|_2^2$$

Максимизация ℓ_S эквивалентна максимизации $-\frac{1}{2\sigma^2} \|y - Xw\|_2^2$, т.к. первое слагаемое не зависит от w

$$\text{AIC: } \boxed{-\frac{1}{2\sigma^2} \|y - Xw\|_2^2 - d} \quad (*)$$

- Рассмотрим C_p -Mallows:

$$\frac{\hat{R}_n(S)}{2\sigma^2} + |S| = \boxed{\frac{\|y - Xw\|_2^2}{2\sigma^2} + d} \quad (**)$$

Получаем, что модель с наибольшим значением AIC (*) является моделью с наименьшим значением статистики Mallows C_p (**)

Ф.И.Г.

Задача 3 [2 балла]

Пусть дана обучающая выборка $\{(x, y): x \in \mathbb{R}^n, y \in \mathbb{R}^n\}$. Предположим, что справедлива модель линейной регрессии:

$$y = w_0 + w_1 x + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2).$$

Сконструируйте тест Вальда для проверки гипотезы $H_0: w_1 = \alpha w_0$.

Решение:

$$W = \frac{w_1 - \alpha w_0}{\hat{se}(w_1 - \alpha w_0)}$$

Найдем оценку для стандартного отклонения. Известно (с семинара), что

$$V(\hat{w}) = \sigma^2 (X^T X)^{-1} = \frac{\sigma^2}{n(\langle X^2 \rangle - \langle X \rangle^2)} \begin{pmatrix} \langle X^2 \rangle & -\langle X \rangle \\ -\langle X \rangle & 1 \end{pmatrix}, \quad \text{где } \langle X^k \rangle = \frac{1}{n} \sum_{i=1}^n X_i^k$$

Значение дисперсии известно, поэтому найдем стандартное отклонение:

$$\begin{aligned} \hat{se}(w_1 - \alpha w_0) &= \sqrt{\mathbb{V}(w_1 - \alpha w_0)} = \sqrt{\mathbb{V}(w_1) + \mathbb{V}(-\alpha w_0) + 2\text{cov}(w_1, -\alpha w_0)} = \sqrt{\mathbb{V}(w_1) + \alpha^2 \mathbb{V}(w_0) - 2\alpha \text{cov}(w_1, w_0)} = \\ &= \frac{\sigma}{\sqrt{n(\langle X^2 \rangle - \langle X \rangle^2)}} \sqrt{1 + \alpha^2 \langle X^2 \rangle + 2\langle X \rangle} \end{aligned}$$

Критерия Вальда размера β отклоняет нулевую гипотезу в пользу альтернативной тогда и только тогда, когда $|W| > z_{\beta/2}$, т.е. когда $|w_1 - \alpha w_0| > z_{\beta/2} \frac{\sigma}{\sqrt{n(\langle X^2 \rangle - \langle X \rangle^2)}} \sqrt{1 + \alpha^2 \langle X^2 \rangle + 2\langle X \rangle}$

Задача 4 [2 балла]

Пусть дана обучающая выборка $\{(\mathbf{X}, \mathbf{y}): \mathbf{X} \in \mathbb{R}^{n \times d}, \mathbf{y} \in \mathbb{R}^n\}$, причем данные соответствуют модели линейной регрессии:

$$\mathbf{y} = \mathbf{x}^T \mathbf{w} + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2),$$

где \mathbf{w} — истинный, но неизвестный нам вектор весов. Пусть $\hat{\mathbf{w}}$ — MLE-оценка вектора весов \mathbf{w} .

Предположим, к нам поступили тестовые данные $\mathbf{X}^* \in \mathbb{R}^{m \times d}$, для которых с помощью оценки $\hat{\mathbf{w}}$ предсказываем вектор $\mathbf{y}^* \in \mathbb{R}^m$. Найдите математическое ожидание и матрицу ковариаций для вектора \mathbf{y}^* (при условии фиксированной матрицы дизайна \mathbf{X}).

Решение:

Из семинара известно, что $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$.

$$\mathbb{E} \mathbf{y}^* = \mathbb{E}(\mathbf{X}^* \hat{\mathbf{w}}) = \mathbf{X}^* \mathbb{E} \hat{\mathbf{w}} = \mathbf{X}^* \mathbb{E}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}) = \mathbf{X}^* (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}(\mathbf{X} \mathbf{w} + \varepsilon) = \mathbf{X}^* (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^* \mathbf{w}$$

$$\begin{aligned} \text{cov}(\mathbf{y}^*) &= \text{cov}(\mathbf{y}^*, \mathbf{y}^*) = \mathbb{E}[(\mathbf{y}^* - \mathbb{E} \mathbf{y}^*)(\mathbf{y}^* - \mathbb{E} \mathbf{y}^*)^T] = \mathbb{E}[\mathbf{y}^* (\mathbf{y}^*)^T] - \mathbb{E} \mathbf{y}^* \mathbb{E} (\mathbf{y}^*)^T = \\ &= \mathbb{E}[\mathbf{X}^* (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \mathbf{y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^*)^T] - \mathbf{X}^* \mathbf{w} \mathbf{w}^T (\mathbf{X}^*)^T = \\ &= \mathbf{X}^* \mathbf{w} \mathbf{w}^T (\mathbf{X}^*)^T + \sigma^2 \mathbf{X}^* (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^*)^T - \mathbf{X}^* \mathbf{w} \mathbf{w}^T (\mathbf{X}^*)^T = \sigma^2 \mathbf{X}^* (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^*)^T \end{aligned}$$

Ответ: $\mathbb{E} \mathbf{y}^* = \mathbf{X}^* \mathbf{w}$, $\text{cov}(\mathbf{y}^*) = \sigma^2 \mathbf{X}^* (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^*)^T$

Задача 5 [2 балла]

Пусть дана выборка $(\mathbf{X}, \mathbf{t}) = \{(\mathbf{x}_i, t_i): \mathbf{x}_i \in \mathbb{R}^d, t_i \in \mathbb{R}\}_{i=1}^n$, $(\mathbf{X} \in \mathbb{R}^{n \times d}, \mathbf{t} \in \mathbb{R}^n)$. Предположим справедливость следующей модели данных

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

где $\varepsilon(\mathbf{x}) \sim \mathcal{N}(0, \sigma(\mathbf{x})^2)$. Найдите MLE-оценку на вектор весов \mathbf{w} в данном случае.

Решение:

$$\mathcal{L} = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2(\mathbf{x}_i)}} e^{-\frac{1}{2} \cdot \frac{(t_i - \mathbf{x}_i^T \mathbf{w})^2}{\sigma^2(\mathbf{x}_i)}} = \underbrace{\prod_{i=1}^n (2\pi\sigma^2(\mathbf{x}_i))^{-1/2}}_{\text{не зависит от } \mathbf{w}} \cdot e^{-\frac{1}{2} \sum_{i=1}^n \left(\frac{t_i - \mathbf{x}_i^T \mathbf{w}}{\sigma(\mathbf{x}_i)} \right)^2} \rightarrow \max$$

$$\mathcal{L} \rightarrow \max \Leftrightarrow +\frac{1}{2} \sum_{i=1}^n \left(\frac{t_i - \mathbf{x}_i^T \mathbf{w}}{\sigma(\mathbf{x}_i)} \right)^2 \rightarrow \min$$

$$\sum_{i=1}^n \left(\frac{t_i}{\sigma(\mathbf{x}_i)} - \frac{\mathbf{x}_i^T \mathbf{w}}{\sigma(\mathbf{x}_i)} \right)^2 = \sum_{i=1}^n (t_i^* - (\mathbf{x}_i^*)^T \mathbf{w})^2 = \|\mathbf{T}^* - \mathbf{X}^{*T} \mathbf{w}\|_2^2, \text{ где } t_i^* = \frac{t_i}{\sigma(\mathbf{x}_i)} \text{ и } \mathbf{x}_i^* = \frac{\mathbf{x}_i}{\sigma(\mathbf{x}_i)}$$

$$\mathbf{T}^* = \begin{pmatrix} t_1^* \\ t_2^* \\ \vdots \\ t_n^* \end{pmatrix} \quad \mathbf{X}^* = \begin{pmatrix} \mathbf{x}_1^* \\ \mathbf{x}_2^* \\ \vdots \\ \mathbf{x}_n^* \end{pmatrix}$$

$$\text{Тогда } \mathbf{T}^* = \text{diag}\left(\frac{1}{\sigma(\mathbf{x}_1)}, \dots, \frac{1}{\sigma(\mathbf{x}_n)}\right) \cdot \mathbf{T}, \quad \mathbf{X}^* = \text{diag}\left(\frac{1}{\sigma(\mathbf{x}_1)}, \dots, \frac{1}{\sigma(\mathbf{x}_n)}\right) \cdot \mathbf{X}$$

$$\|\mathbf{T}^* - \mathbf{X}^{*T} \mathbf{w}\|_2^2 \rightarrow \min_{\mathbf{w}} \Rightarrow \hat{\mathbf{w}} = ((\mathbf{X}^*)^T \mathbf{X}^*)^{-1} (\mathbf{X}^*)^T \mathbf{T}^*$$

Ответ: $\hat{w} = (DX)^T DX)^{-1} (DX)^T DT$, где $D = \text{diag}(\frac{1}{\sigma(x_1)}, \dots, \frac{1}{\sigma(x_n)})$

Задача 6 [4 балла]

Пусть дана выборка $(\mathbf{x}, \mathbf{y}) = \{(x_i, y_i) : x_i, y_i \in \mathbb{R}\}_{i=1}^n$. Пусть данные соответствуют модели

$$y_i = \beta x_i + \varepsilon_i,$$

где $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$. При этом значения \mathbf{x} наблюдаются с ошибкой, т.е. представлена не выборка (\mathbf{x}, \mathbf{y}) , а выборка $(\mathbf{z}, \mathbf{y}) = \{(z_i, y_i) : z_i, y_i \in \mathbb{R}\}_{i=1}^n$, где $z_i = x_i + \delta_i$, $\delta_i \sim \mathcal{N}(0, \tau^2)$. Шумы ε_i и δ_i независимы. Оценим величину β , используя стандартный метод наименьших квадратов согласно формуле

$$\hat{\beta} = \frac{\sum_{i=1}^n z_i y_i}{\sum_{i=1}^n z_i^2}.$$

Докажите, что оценка $\hat{\beta}$ не является состоятельной. Для этого покажите, что $\hat{\beta} \xrightarrow{P} a\beta$ при $n \rightarrow \infty$. Найдите явное выражение для a в предположении, что точки $\{x_i\}_{i=1}^n$ поступают из некоторого распределения $F(x)$ с конечными первыми и вторыми моментами $\mathbb{E}(X)$ и $\mathbb{E}(X^2)$.

Решение:

$$\hat{\beta} = \frac{\sum_{i=1}^n (x_i + \delta_i)(\beta x_i + \varepsilon_i)}{\sum_{i=1}^n (x_i + \delta_i)^2} = \frac{\frac{\beta}{n} \sum_{i=1}^n x_i^2 + \frac{1}{n} \sum_{i=1}^n x_i \varepsilon_i + \frac{\beta}{n} \sum_{i=1}^n x_i \delta_i + \frac{1}{n} \sum_{i=1}^n \delta_i \varepsilon_i}{\frac{1}{n} \sum_{i=1}^n x_i^2 + \frac{2}{n} \sum_{i=1}^n x_i \delta_i + \frac{1}{n} \sum_{i=1}^n \delta_i^2} \xrightarrow{P} \frac{\beta \mathbb{E}x_i^2}{\mathbb{E}x_i^2 + \tau^2}$$

Покажем сходимость по вероятности слагаемых:

1) $\frac{1}{n} \sum_{i=1}^n x_i^2 \xrightarrow{P} \mathbb{E}x_i^2$ по ЗБЧ

2) $\mathbb{E}\delta_i^2 = \mathbb{V}\delta_i + (\mathbb{E}\delta_i)^2 = \tau^2$, поэтому по ЗБЧ $\frac{1}{n} \sum_{i=1}^n \delta_i^2 \xrightarrow{P} \tau^2$

3) $x_i, \delta_i, \varepsilon_i$ попарно независимы, поэтому $\mathbb{E}(x_i \varepsilon_i) = \mathbb{E}x_i \mathbb{E}\varepsilon_i = 0$, $\mathbb{E}(x_i \delta_i) = \mathbb{E}x_i \mathbb{E}\delta_i = 0$, $\mathbb{E}(\delta_i \varepsilon_i) = \mathbb{E}\delta_i \mathbb{E}\varepsilon_i = 0$

По ЗБЧ соответствующие средние (те все оставшиеся) сходятся к 0 по вероятности.

Ответ: $\alpha = \frac{\mathbb{E}x_i^2}{\mathbb{E}x_i^2 + \tau^2}$

Задача 7 [3 балла]

Пусть дана выборка $(\mathbf{X}, \mathbf{T}) = \{(\mathbf{x}_i, \mathbf{t}_i) : \mathbf{x}_i \in \mathbb{R}^d, \mathbf{t}_i \in \mathbb{R}^m\}_{i=1}^n$, $\mathbf{X} \in \mathbb{R}^{n \times d}$, $\mathbf{T} \in \mathbb{R}^{n \times m}$. Рассмотрим модель *многомерной линейной регрессии*, т.е. регрессии, в которой независимая переменная является вектором:

$$\mathbf{t} = \mathbf{W}^T \mathbf{x} + \varepsilon,$$

где $\mathbf{x} \in \mathbb{R}^d$, $\mathbf{t} \in \mathbb{R}^m$, $\mathbf{W} \in \mathbb{R}^{d \times m}$, $\varepsilon \in \mathbb{R}^m$. Рассмотрим модель, в рамках которой плотность распределения вектора \mathbf{t} при заданном векторе \mathbf{x} имеет вид $p(\mathbf{t}|\mathbf{x}) = \mathcal{N}(\mathbf{t}|\mathbf{W}^T \mathbf{x}, \Sigma)$, т.е. нормальное распределение со средним $\mathbf{W}^T \mathbf{x} \in \mathbb{R}^m$ и матрицей ковариаций $\Sigma \in \mathbb{R}^{m \times m}$. Найдите ML-оценки для матриц \mathbf{W} и Σ .

Подсказка. Вам могут потребоваться следующие формулы матричного дифференцирования:

$$\frac{\partial \mathbf{a}^T \mathbf{X}^{-1} \mathbf{b}}{\partial \mathbf{X}} = -\mathbf{X}^{-T} \mathbf{a} \mathbf{b}^T \mathbf{X}^{-T}, \quad \frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{b}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{b}^T, \quad \frac{\partial \mathbf{a}^T \mathbf{X}^T \mathbf{b}}{\partial \mathbf{X}} = \mathbf{b} \mathbf{a}^T.$$

Внимание. Во возможности ответ следует полностью записать в матричном виде, выразив всё через \mathbf{X} и \mathbf{T} .

Решение:

$$\log L = C - \frac{n}{2} \log |\Sigma| - \frac{1}{2} \sum_{i=1}^n (\mathbf{t}_i - \mathbf{W}^T \mathbf{x}_i)^T \Sigma^{-1} (\mathbf{t}_i - \mathbf{W}^T \mathbf{x}_i),$$

где C - некоторая константа, не зависящая от \mathbf{W} и Σ .

Для удобства домножим выражение на (-2) и будем искать минимум (для удобства вычисления, чтобы не тащить за собой минусы). Также заметим, что Σ^{-1} - симметричная матрица, поэтому она равна транспонированной.

$$\begin{aligned} \frac{\partial (-2 \log L)}{\partial \Sigma} &= \frac{n}{|\Sigma|} |\Sigma| \Sigma^{-T} - \Sigma^{-T} \sum_{i=1}^n (\mathbf{t}_i - \mathbf{W}^T \mathbf{x}_i) (\mathbf{t}_i - \mathbf{W}^T \mathbf{x}_i)^T \Sigma^{-T} = n \Sigma^{-1} - \Sigma^{-1} (\mathbf{T} - \mathbf{X} \mathbf{W})^T (\mathbf{T} - \mathbf{X} \mathbf{W}) \Sigma^{-1} = 0 \\ &\Rightarrow \Sigma = \frac{1}{n} (\mathbf{T} - \mathbf{X} \mathbf{W})^T (\mathbf{T} - \mathbf{X} \mathbf{W})^T \end{aligned}$$

$$\frac{\partial(-2 \log L)}{\partial W} = \sum_{i=1}^n [-2x_i t_i^T \Sigma^{-1} + 2x_i x_i^T w \Sigma^{-1}] = 0$$

$$\Rightarrow W = (X^T X)^{-1} X^T T$$

Ответ: $\Sigma = \frac{1}{n}(T - XW)^T(T - XW)^T$, $W = (X^T X)^{-1} X^T T$

Задача 8 [2 балла]

Рассмотрим задачу восстановления регрессии. Модель регрессии имеет вид

$$t = \mathbf{x}^T \mathbf{w} + \varepsilon,$$

где $\varepsilon \sim \mathcal{N}(0, \beta^{-1})$, и на веса \mathbf{w} наложено априорное распределение вида $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{w}_0, \mathbf{S}_0)$. Пусть дана выборка $(\mathbf{X}, \mathbf{t}) = \{(\mathbf{x}_i, t_i) : \mathbf{x}_i \in \mathbb{R}^d, t_i \in \mathbb{R}\}_{i=1}^n$. Найдите апостериорное распределение $p(\mathbf{w} | \mathbf{X}, \mathbf{t})$.

Решение:

$$p(\mathbf{w} | \mathbf{x}, \mathbf{t}) = \frac{p(\mathbf{w}, \mathbf{x}, \mathbf{t})}{p(\mathbf{x}, \mathbf{t})} = \frac{p(\mathbf{t} | \mathbf{x}, \mathbf{w}) \cdot p(\mathbf{x}, \mathbf{w})}{p(\mathbf{x}, \mathbf{t})} = \frac{p(\mathbf{t} | \mathbf{x}, \mathbf{w}) p(\mathbf{x}) p(\mathbf{w})}{p(\mathbf{x}) p(\mathbf{t})} =$$

$$= \frac{\prod_{i=1}^n p(t_i | x_i, \mathbf{w}) \cdot p(\mathbf{w})}{p(\mathbf{t})}$$

$$p(\mathbf{w} | \mathbf{x}, \mathbf{t}) \propto \frac{1}{2} \exp \left\{ -\frac{\beta \|\mathbf{x}\mathbf{w} - \mathbf{t}\|_2^2}{2} - \frac{(\mathbf{w} - \mathbf{w}_0)^T \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{w}_0)}{2} \right\}$$

Покажем, что $\exp \{ \dots \}$ можно записать в виде $\exp \left\{ -\frac{1}{2} (\mathbf{w} - \mathbf{w}_n)^T \Sigma^{-1} (\mathbf{w} - \mathbf{w}_n) \right\}$, тогда получим нормальное распределение $\mathcal{N}(\mathbf{w}_n, \Sigma)$.
($\frac{1}{2}$ - некоторая нормировочная константа)

$$\begin{aligned} \beta \|\mathbf{x}\mathbf{w} - \mathbf{t}\|_2^2 + (\mathbf{w} - \mathbf{w}_0)^T \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{w}_0) &= \beta (\mathbf{x}\mathbf{w} - \mathbf{t})^T (\mathbf{x}\mathbf{w} - \mathbf{t}) + (\mathbf{w} - \mathbf{w}_0)^T \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{w}_0) = \\ &= \beta (\mathbf{w}^T \mathbf{x}^T \mathbf{x} \mathbf{w} - \mathbf{w}^T \mathbf{x}^T \mathbf{t} - \mathbf{t}^T \mathbf{x} \mathbf{w} + \mathbf{t}^T \mathbf{t}) + (\mathbf{w} - \mathbf{w}_0)^T \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{w}_0) = \\ &= \underline{\mathbf{w}^T (\beta \mathbf{x}^T \mathbf{x} + \mathbf{S}_0^{-1}) \mathbf{w}} - \beta \underline{\mathbf{w}^T \mathbf{x}^T \mathbf{t}} - \beta \mathbf{t}^T \mathbf{x} \mathbf{w} + \beta \mathbf{t}^T \mathbf{t} - \underline{\mathbf{w}^T \mathbf{S}_0^{-1} \mathbf{w}_0} - \mathbf{w}_0^T \mathbf{S}_0^{-1} \mathbf{w} + \mathbf{w}_0^T \mathbf{S}_0^{-1} \mathbf{w}_0 = \\ &= \underline{\mathbf{w}^T \Sigma^{-1} \mathbf{w}} - \mathbf{w}^T \Sigma^{-1} \mathbf{w}_n - \mathbf{w}_n^T \Sigma^{-1} \mathbf{w} + \mathbf{w}_n^T \Sigma^{-1} \mathbf{w}_n \end{aligned}$$

$$\Sigma^{-1} = \beta \mathbf{x}^T \mathbf{x} + \mathbf{S}_0^{-1}$$

$$-\mathbf{w}^T \Sigma^{-1} \mathbf{w}_n = -\mathbf{w}^T (\beta \mathbf{x}^T \mathbf{t} + \mathbf{S}_0^{-1} \mathbf{w}_0) \Rightarrow \Sigma^{-1} \mathbf{w}_n = \beta \mathbf{x}^T \mathbf{t} + \mathbf{S}_0^{-1} \mathbf{w}_0 \Rightarrow$$

$$\Rightarrow \mathbf{w}_n = \Sigma (\beta \mathbf{x}^T \mathbf{t} + \mathbf{S}_0^{-1} \mathbf{w}_0) = (\beta \mathbf{x}^T \mathbf{x} + \mathbf{S}_0^{-1})^{-1} (\beta \mathbf{x}^T \mathbf{t} + \mathbf{S}_0^{-1} \mathbf{w}_0)$$

$$p(\mathbf{w} | \mathbf{x}, \mathbf{t}) = \mathcal{N}((\beta \mathbf{x}^T \mathbf{x} + \mathbf{S}_0^{-1})^{-1} (\beta \mathbf{x}^T \mathbf{t} + \mathbf{S}_0^{-1} \mathbf{w}_0), (\beta \mathbf{x}^T \mathbf{x} + \mathbf{S}_0^{-1})^{-1})$$

Ответ: $\mathcal{N}(\mathbf{w} | (\beta \mathbf{X}^T \mathbf{X} + \mathbf{S}_0^{-1})^{-1} (\beta \mathbf{X}^T \mathbf{t} + \mathbf{S}_0^{-1} \mathbf{w}_0), (\beta \mathbf{X}^T \mathbf{X} + \mathbf{S}_0^{-1})^{-1})$

Задача 9 [2 балл]

Пусть $\mathbf{x}^n \sim f(\cdot)$, и пусть $\hat{f}(\cdot) = \hat{f}(\cdot; \mathbf{x}^n)$ обозначает ядерную оценку плотности на основе ядра

$$K(x) = \begin{cases} 1, & x \in (-\frac{1}{2}, \frac{1}{2}); \\ 0, & \text{в противном случае.} \end{cases}$$

Найдите $\mathbb{E}[\hat{f}(x)]$ и $\mathbb{V}[\hat{f}(x)]$. Покажите, что если $h \rightarrow 0$ и $nh \rightarrow \infty$ при $n \rightarrow \infty$, то $\hat{f}(x) \xrightarrow{P} f(x)$ при $n \rightarrow \infty$.

Решение:

$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - X_i}{h}\right)$$

$$\mathbb{E}[\hat{f}(x)] = \frac{1}{h} \mathbb{E}K\left(\frac{x - X_i}{h}\right) = \frac{1}{h} \int K\left(\frac{x - u}{h}\right) f(u) du = \frac{1}{h} \int_{x-h/2}^{x+h/2} f(y) dy$$

$$\mathbb{V}[\hat{f}(x)] = \frac{1}{nh^2} \mathbb{V}K\left(\frac{x - X_i}{h}\right) = \frac{1}{nh^2} \left(\mathbb{E}K^2 - [\mathbb{E}K]^2 \right) = \frac{1}{nh^2} \int_{x-h/2}^{x+h/2} f(y) dy \left(1 - \int_{x-h/2}^{x+h/2} f(y) dy \right)$$

Теперь покажем сходимост по вероятности. Для этого воспользуемся выражениями из лекции:

$$\mathbb{E}[\hat{f}(x)] = \mathbb{E} \frac{1}{h} K\left(\frac{x - X_i}{h}\right) = f(x) + \frac{1}{2} h^2 f''(x) \int u^2 K(u) du + \dots$$

$$\mathbb{V}[\hat{f}(x)] \approx \frac{f(x) \int K^2(x) dx}{nh}$$

Запишем неравенство Чебышева и покажем, что при $h \rightarrow 0$ и $nh \rightarrow \infty$ выполняется сходимост по вероятности $\hat{f}(x) \xrightarrow{P} f(x)$ при $n \rightarrow \infty$:

$$P\left(\left|\hat{f}(x) - \mathbb{E}[\hat{f}(x)]\right| \geq \varepsilon\right) = P\left(\left|\hat{f}(x) - f(x)\right| \geq \varepsilon\right) \leq \frac{\mathbb{V}[\hat{f}(x)]}{\varepsilon^2} \approx \frac{f(x)}{nh\varepsilon^2} \rightarrow 0$$

Что и требовалось показать

Ответ: $\mathbb{E}[\hat{f}(x)] = \frac{1}{h} \int_{x-h/2}^{x+h/2} f(y) dy, \quad \mathbb{V}[\hat{f}(x)] = \frac{1}{nh^2} \int_{x-h/2}^{x+h/2} f(y) dy \left(1 - \int_{x-h/2}^{x+h/2} f(y) dy \right)$

Задача 10 [6 баллов]

Рассмотрим задачу непараметрической оценки плотности распределения $p(x)$ по выборке $\mathbf{x}^{(n)}$. Обозначим через $\hat{p}(x; \mathbf{x}^{(n)})$ оценку плотности, полученную некоторым образом по выборке $\mathbf{x}^{(n)}$. Оценка риска для $\hat{p}(x; \mathbf{x}^{(n)})$ имеет вид:

$$\hat{J}(h) = \int (\hat{p}(x; \mathbf{x}^{(n)}))^2 dx - \frac{2}{n} \sum_{i=1}^n \hat{p}(x_i; \mathbf{x}^{(n \setminus i)}),$$

где $\hat{p}(\cdot; \mathbf{x}^{(n \setminus i)})$ — оценка плотности распределения на основе выборки $\mathbf{x}^{(n \setminus i)}$, т.е. выборки без объекта x_i .

- (Гистограммная оценка) Разобьем диапазон наблюдаемых значений $\mathbf{x}^{(n)}$ на бины ширины h . Пусть в итоге значения $\mathbf{x}^{(n)}$ укладываются в M последовательных бинов B_1, \dots, B_M . Пусть n_m — количество объектов выборки, попавших в B_m . Пусть \hat{p}_m — доля объектов выборки, попавших в бин B_m :

$$n_m = \sum_i I[x_i \in B_m], \quad \hat{p}_m = \frac{n_m}{n}.$$

Покажите, что в случае гистограммной оценки плотности оценка риска имеет вид:

$$\hat{J}(h) = \frac{2}{h(n-1)} - \frac{n+1}{h(n-1)} \sum_{m=1}^M \hat{p}_m^2.$$

Докажите или опровергните равенство

$$\mathbb{E}[\hat{J}(h)] = \mathbb{E}[J(h)].$$

Если равенство не верно, то чему равно $\Delta J(h) = \mathbb{E}[\hat{J}(h)] - \mathbb{E}[J(h)]$?

Решение:

$$\hat{J}(h) = \sum_{m=1}^M \frac{\hat{p}_m^2}{h^2} h - \frac{2}{nh} \sum_{i=1}^n \sum_{m=1}^M \frac{n_m - 1}{n - 1} I[x_i \in B_m] = \frac{1}{h} \sum_{m=1}^M \hat{p}_m^2 - \frac{2}{h} \sum_{m=1}^M \left(\hat{p}_m - \frac{1}{n} \right) \frac{n_m}{n - 1} =$$

$$= \frac{1}{h} \sum_{m=1}^M \hat{p}_m^2 - \frac{2n}{h(n-1)} \sum_{m=1}^M (\hat{p}_m^2 - \frac{1}{n}) \frac{n_m}{n} = \frac{1}{h} \sum_{m=1}^M \hat{p}_m^2 - \frac{2n}{h(n-1)} \sum_{m=1}^M \hat{p}_m^2 + \frac{2}{h(n-1)} = \frac{2}{h(n-1)} - \frac{n+1}{h(n-1)} \sum_{m=1}^M \hat{p}_m^2$$

Оценка смещенная, почему?

- (Ядерная оценка) Покажите, что в случае ядерной оценки плотности оценка риска имеет вид:

$$\hat{J}(h) \approx \frac{1}{hn^2} \sum_{i,j} K^* \left(\frac{X_i - X_j}{h} \right) + \frac{2}{nh} K(0),$$

где $K^*(x) = K^{(2)}(x) - 2K(x)$ и $K^{(2)}(z) = \int K(z-y)K(y)dy$. В частности, если $K(x)$ — это плотность нормального распределения $\mathcal{N}(0, 1)$, т.е. гауссово ядро, то $K^{(2)}(z)$ — плотность распределения $\mathcal{N}(0, 2)$.

Докажите или опровергните равенство

$$\mathbb{E}[\hat{J}(h)] = \mathbb{E}[J(h)].$$

Если равенство не верно, то чему равно $\Delta J(h) = \mathbb{E}[\hat{J}(h)] - \mathbb{E}[J(h)]$?

Решение:

$$\begin{aligned} \hat{J}(h) &= \int \left(\frac{1}{nh} \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) \right)^2 dx - \frac{2}{n} \sum_{i=1}^n \frac{1}{(n-1)h} \sum_{j=1}^n k\left(\frac{x_i-x_j}{h}\right) = \\ &= \frac{1}{n^2 h^2} \int \left[\sum_{i=1}^n k^2\left(\frac{x-x_i}{h}\right) + \sum_{i \neq j} k\left(\frac{x-x_i}{h}\right) k\left(\frac{x-x_j}{h}\right) \right] dx - \\ &= \frac{2}{hn(n-1)} \sum_{i,j} k\left(\frac{x_i-x_j}{h}\right) + \frac{2}{nh} k(0) = \\ &= \frac{h}{n^2 h^2} \sum_{i,j} \int k\left(\frac{x_i-x_j}{h} - y\right) k(y) dy - \frac{2}{n(n-1)h} \sum_{i,j} k\left(\frac{x_i-x_j}{h}\right) + \frac{2}{nh} k(0) \approx \\ &\approx \frac{1}{hn^2} \sum_{i,j} k^*\left(\frac{x_i-x_j}{h}\right) + \frac{2}{nh} k(0) \end{aligned}$$

Покажем, что $\mathbb{E}[\hat{J}(h)] = \mathbb{E}[J(h)]$

$$\hat{J}(h) = \int_{-\infty}^{+\infty} \hat{f}_h^2(x) dx - \frac{2}{h} \sum_{i=1}^n \hat{f}_{(-i)}(x_i)$$

$$J(h) = \int_{-\infty}^{+\infty} f_h^2(x) dx - 2 \int_{-\infty}^{+\infty} f_h(x) f(x) dx$$

Достаточно показать, что

$$\mathbb{E} \int_{-\infty}^{+\infty} \hat{f}_h(x) f(x) dx = \mathbb{E} \left[\frac{1}{h} \sum_{i=1}^n \hat{f}_{(-i)}(x_i) \right]$$

$$\begin{aligned} \mathbb{E} \int_{-\infty}^{+\infty} \hat{f}_h(x) f(x) dx &= \mathbb{E} \int_{-\infty}^{+\infty} \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right) f(x) dx = \\ &= \frac{1}{h} \mathbb{E} \int_{-\infty}^{+\infty} K\left(\frac{x-x_i}{h}\right) f(x) dx = \frac{1}{h} \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} K\left(\frac{x-u}{h}\right) f(x) f(u) dx du \quad (*) \end{aligned}$$

$$\begin{aligned} \mathbb{E} \left[\frac{1}{h} \sum_{i=1}^n \hat{f}_{(-i)}(x_i) \right] &= \mathbb{E} \hat{f}_{(-i)}(x_i) = \mathbb{E} \frac{1}{(n-1)h} \sum_{\substack{j=1 \\ j \neq i}}^n K\left(\frac{x_i-x_j}{h}\right) = \\ &= \frac{1}{h} \mathbb{E} K\left(\frac{x_i-x_j}{h}\right) = \frac{1}{h} \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} K\left(\frac{x-u}{h}\right) f(x) f(u) dx du \quad (**) \end{aligned}$$

(*) = (**) \Rightarrow равенство доказано

Задача 11 [3 балла]

Рассмотрим задачу непараметрической регрессии:

$$Y_i = f(X_i) + \varepsilon_i, \quad i \in \{1, \dots, n\}, \quad X_i \in \mathbb{R}, \quad Y_i \in \mathbb{R}.$$

где ε_i и X_i независимы, $\mathbb{E}\varepsilon_i = 0$, $\mathbb{V}\varepsilon_i = \sigma^2$, выборка $\{X_i\}_{i=1}^n$ одномерная и сэмплируется из отрезка $[0, 1]$. Необходимо по имеющимся данным оценить функцию регрессии $f(x) = \mathbb{E}(Y|X = x)$.

а) Рассмотрим следующее семейство функций

$$\mathfrak{F}_M = \left\{ f(x) = \sum_{i=1}^M c_i I[x \in B_i], c_i \in \mathbb{R}, i = \overline{1, M} \right\}, \text{ где } B_i = \left[\frac{i-1}{M}, \frac{i}{M} \right).$$

Последний отрезок B_M включает обе граничные точки. Найдите функцию из класса \mathfrak{F}_M , которая минимизирует сумму квадратов ошибок:

$$r(x; \mathbf{X}^n) = \arg \min_{f(x) \in \mathfrak{F}_M} \sum_{i=1}^n (Y_i - f(X_i))^2$$

б) Найдите функцию регрессии поточечно, решив в каждой точке x следующую оптимизационную задачу:

$$r(x; \mathbf{X}^n) = \arg \min_{y \in \mathbb{R}} \sum_{i=1}^n K\left(\frac{x - X_i}{h}\right) (Y_i - y)^2,$$

где $K(x)$ — заданная ядерная функция, h — ширина ядра.

с) Какая оценка получится, если изменить задачу на следующую:

$$r(x; \mathbf{X}^n) = \arg \min_{a, b \in \mathbb{R}} \sum_{i=1}^n K\left(\frac{x - X_i}{h}\right) (Y_i - a - bX_i)^2,$$

где $K(x)$ — заданная ядерная функция, h — ширина ядра?

Решение:

$$a) \argmin_{f(x) \in \mathcal{F}(x)} \sum_{i=1}^n (y_i - f(x_i))^2$$

Чтобы минимиз. сумму квадратов, нужно минимиз. её во каждой бинке B_m , т.е. $\bar{y}_m = \text{mean}(y_i)_{x_i \in B_m}$. Тогда получаем, что

$$f(x) = \sum_{i=1}^M I(x \in B_i) \text{mean}(y_j)_{x_j \in B_i}$$

$$b) \argmin_{y \in \mathbb{R}} \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) (y_i - y)^2$$

$$\frac{\partial r}{\partial y} = -2 \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) (y_i - y) = +2 \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) (y - y_i) = 0$$

$$\Rightarrow y = \frac{\sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) y_i}{\sum_{i=1}^n k\left(\frac{x-x_i}{h}\right)}$$

$$b) \argmin_{a, b \in \mathbb{R}} \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) (y_i - a - bx_i)^2$$

$$\begin{cases} \frac{\partial r}{\partial a} = -2 \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) (y_i - a - bx_i) = 0 \\ \frac{\partial r}{\partial b} = -2 \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) (y_i - a - bx_i) x_i = 0 \end{cases}$$

$$\Rightarrow b = \frac{\sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) x_i \cdot \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) y_i - \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) \cdot \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) x_i y_i}{\left(\sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) x_i\right)^2 - \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) \cdot \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) x_i^2}$$

$$a = \frac{\sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) x_i y_i \cdot \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) x_i - \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) y_i \cdot \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) x_i^2}{\left(\sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) x_i\right)^2 - \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) \cdot \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) x_i^2}$$

Практический блок

Задача 12 [7 баллов]

Скачайте по ссылке данные о связи между оценкой качества вина от различных характеристик вина <https://archive.ics.uci.edu/ml/datasets/wine+quality>. По ссылке представлено два набора данных: для белых и для красных вин. Далее предполагается использование данных для белых вин (`winequality-white.csv`). Разбейте данные на обучающую и тестовую выборку: для тестовой выборки возьмите 25% данных.

- Обучите простую линейную регрессию по обучающей выборке. Примените модель к тестовой выборке и найдите MSE.
- По обучающей выборке оцените наилучший набор признаков, описывающих выходную переменную. Используйте для этого статистику Cp Mallow, AIC-критерий, BIC-критерий, LOO-проверку и 10-кратную кросс-проверку. Выбор подмножества признаков проведите полным перебором. Позволяет ли какой-нибудь набор признаков получить значение MSE на тестовых данных меньше, чем на всех признаках.
- Обучите гребневую регрессию и Lasso. Оптимальные параметры подберите с помощью поиска по сетке и 10-кратной кросс-проверки. Получилось ли обучить модель, имеющую лучшее качество на тестовой выборке, чем простая линейная регрессия.

Задача 13 [5 баллов]

Скачать данные со страницы курса (значения коэффициента преломления для разных типов стекла; первый столбец). Оценить плотность распределения этих значений, используя гистограмму и ядерную оценку. Для подбора ширины ячейки или ширины ядра использовать перекрестную проверку (кросс-проверку). Для выбранных значений ширины ячейки и ширины ядра построить 95%-ые доверительные интервалы для полученной оценки плотности.

Задача 14 [5 баллов]

По данным из предыдущей задачи, используя в качестве выходной переменной y значения преломления для разных типов стекла, а в качестве входной переменной x — данные о содержании алюминия (четвертая переменная в матрице данных), восстановить зависимость между y и x с помощью ядерной непараметрической регрессии. Оценку ядра проводить с помощью перекрестной проверки. Построить 95%-ые доверительные интервалы для полученной оценки функции регрессии.