

## Матричные вычисления и нормальное распределение

Дата: 19 октября 2011

### Дивергенция Кульбака-Лейблера

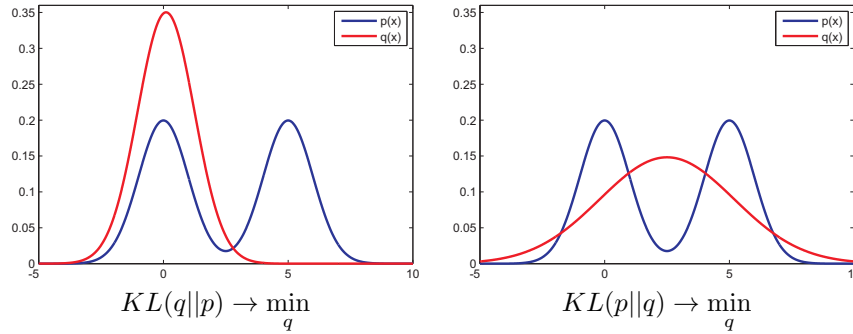


Рис. 1: Приближение двумодального распределения  $p(x)$  с помощью одномодального распределения  $q(x)$ .

Дивергенция Кульбака-Лейблера  $KL(q||p)$  является мерой расстояния между двумя вероятностными распределениями  $q(x)$  и  $p(x)$ :

$$KL(q||p) = - \int q(x) \log \frac{p(x)}{q(x)} dx.$$

КЛ-дивергенция обладает следующими свойствами:

- $KL(q||p) \geq 0$  и  $= 0 \Leftrightarrow q(x) \equiv p(x)$ ,
- $KL(q||p) \neq KL(p||q)$ .

Докажем первое свойство. Для этого рассмотрим произвольную строго вогнутую функцию  $f(y)$ . Тогда по определению

$$f(\alpha y_1 + (1 - \alpha)y_2) \geq \alpha f(y_1) + (1 - \alpha)f(y_2), \quad \forall y_1, y_2, \alpha \in [0, 1],$$

причем равенство достигается, если  $y_1 = y_2$  или  $\alpha = 0, 1$ . Это неравенство может быть обобщено на случай большего числа точек  $y$ :

$$f(\alpha_1 y_1 + \dots + \alpha_N y_N) \geq \alpha_1 f(y_1) + \dots + \alpha_N f(y_N), \quad \sum_{n=1}^N \alpha_n = 1, \alpha_n \geq 0.$$

Последнее неравенство известно как неравенство Йенсена и может быть легко доказано по индукции. Переходя в этом неравенстве от сумм к интегралам, получаем следующее обобщение:

$$f\left(\int \alpha(x)y(x)dx\right) \geq \int \alpha(x)f(y(x))dx, \quad \int \alpha(x)dx = 1, \alpha(x) \geq 0 \quad \forall x.$$

Подставляя в последнее неравенство для функции  $\log$  значения  $\alpha(x) = q(x)$ ,  $y(x) = \frac{p(x)}{q(x)}$ , получаем:

$$0 = \log\left(\int q(x)\frac{p(x)}{q(x)}dx\right) \geq \int q(x)\log\frac{p(x)}{q(x)}dx = -KL(q||p).$$

Так как коэффициенты  $q(x)$ , вообще говоря, отличны от нуля, то равенство в последнем неравенстве достигается только при  $y(x) = \text{const}$ , т.е.  $q(x) \equiv p(x)$ .

#### Влияние несимметричности КЛ-дивергенции.

Рассмотрим задачу аппроксимации распределения  $p(x)$  с помощью распределения  $q(x)$ . Будем искать приближение двумя способами:

1.  $\text{KL}(q||p) \rightarrow \min_q$ ,
2.  $\text{KL}(p||q) \rightarrow \min_q$ .

В первом случае аппроксимация ищется в области высоких значений  $q$ . Поэтому итоговое распределение, как правило, хорошо приближает распределение  $p(\mathbf{x})$  только на подмножестве носителя  $p(\mathbf{x})$  (см. рис. 1, слева). Во втором случае аппроксимация ищется сразу для всего носителя распределения  $p(\mathbf{x})$  (см. рис. 1, справа).

## Матричные вычисления

Здесь и далее вектора будут обозначаться жирным шрифтом  $\mathbf{x}, \mathbf{y}, \dots$ , а матрицы – заглавными буквами  $A, B, \dots$ . При этом под вектором всегда будет пониматься вектор-столбец  $\mathbf{x} = [x_1, \dots, x_d]^T$ , индекс  $T$  обозначает транспонирование. Например, запись  $\mathbf{x}^T A \mathbf{y}$  означает  $\sum_{i,j} a_{ij} x_i y_j$ .

Использование матрично-векторной нотации часто оказывается удобным с точки зрения вывода формул в многомерном пространстве. Кроме того, матричная запись позволяет в некоторых случаях повысить эффективность вычисления соответствующих формул на компьютере. Например, эффективность современных алгоритмов вычисления произведения двух матриц  $AB$  имеет сложность  $n^\alpha$ , где  $n$  – размерность матрицы, а  $\alpha < 2.42$ .

Основные матричные операции представлены в таблице 1. Здесь под производной скалярной функции  $f(\mathbf{x})$  по вектору  $\mathbf{x}$  понимается градиент  $\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = [\frac{\partial f(\mathbf{x})}{\partial x_1}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_d}]^T$ , под производной скалярной функции  $f(A)$  по матрице  $A$  – матрица частных производных  $\left(\frac{\partial f(A)}{\partial A}\right)_{ij} = \frac{\partial f(A)}{\partial a_{ij}}$ , под производной вектор-функции  $\mathbf{f}(x)$  по скалярному аргументу  $x$  – вектор  $\frac{\partial \mathbf{f}(x)}{\partial x} = [\frac{\partial f_1(x)}{\partial x}, \dots, \frac{\partial f_d(x)}{\partial x}]^T$ , под производной матричной функции  $A(x)$  по скалярному аргументу – матрица  $\left(\frac{\partial A(x)}{\partial x}\right)_{ij} = \frac{\partial A_{ij}(x)}{\partial x}$ .

Таблица 1: Основные матричные тождества

### Базовые операции:

$$\begin{aligned} A(B + C) &= AB + AC, \\ (A + B)^T &= A^T + B^T, \\ (AB)^T &= B^T A^T, \\ (AB)^{-1} &= B^{-1} A^{-1}, \\ (A^{-1})^T &= (A^T)^{-1}. \end{aligned}$$

### Производные следа и определителя:

$$\begin{aligned} \frac{\partial}{\partial A} \text{tr} AB &= B^T, \\ \frac{\partial}{\partial A} \det A &= (\det A)(A^{-1})^T, \\ \frac{\partial}{\partial x} \log \det A(x) &= \text{tr} \left( A^{-1} \frac{\partial A}{\partial x} \right). \end{aligned}$$

### След и определитель:

$$\begin{aligned} \det(AB) &= \det A \det B, \\ \det(A^{-1}) &= 1/\det A, \\ \det A &= \prod_j \lambda_j, \\ \text{tr} A &= \sum_j A_{jj} = \sum_j \lambda_j, \\ \text{tr}(ABC) &= \text{tr}(BCA) = \text{tr}(CAB). \end{aligned}$$

### Производные:

$$\begin{aligned} \frac{\partial}{\partial \mathbf{x}} \mathbf{x}^T \mathbf{a} &= \mathbf{a}, \\ \frac{\partial}{\partial \mathbf{x}} \mathbf{x}^T A \mathbf{x} &= (A + A^T) \mathbf{x}, \\ \frac{\partial}{\partial A} \mathbf{x}^T A \mathbf{y} &= \mathbf{x} \mathbf{y}^T, \\ \frac{\partial}{\partial x} A^{-1} &= -A^{-1} \frac{\partial A}{\partial x} A^{-1}. \end{aligned}$$

Используя свойства из таблицы 1, легко показать, например, что

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{t} - A\mathbf{x})^T (\mathbf{t} - A\mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} [\mathbf{t}^T \mathbf{t} - 2\mathbf{x}^T A^T \mathbf{t} + \mathbf{x}^T A^T A \mathbf{x}] = -2A^T \mathbf{t} + 2A^T A \mathbf{x} = 2A^T (A\mathbf{x} - \mathbf{t}).$$

Рассмотрим некоторые случаи использования матричных операций для повышения эффективности вычислений на компьютере. Пусть необходимо решить систему линейных уравнений с невырожденной квадратной положительно-определенной матрицей  $A\mathbf{x} = \mathbf{t}$ . Прямое решение данной системы  $\mathbf{x} = A^{-1}\mathbf{t}$  требует вычисления обратной матрицы  $A^{-1}$ . Вычисление обратной матрицы – численно неустойчивая процедура, особенно для плохо обусловленных матриц. Однако, в данном случае нам не требуется знать полностью значение обратной матрицы, достаточно знать лишь вектор  $A^{-1}\mathbf{t}$ . Для его эффективного вычисления воспользуемся разложением Холецкого  $A = RR^T$ ,

где  $R$  – нижнетреугольная матрица, т.е. матрица, у которой все элементы выше главной диагонали равны нулю. Тогда

$$Ax = t \Leftrightarrow \underbrace{RR^T}_y x = t \Leftrightarrow \begin{cases} Ry = t, \\ R^T x = y. \end{cases}$$

Таким образом, решение системы линейных уравнений  $Ax = t$  эквивалентно решению двух систем линейных уравнений с треугольными матрицами. Эти системы могут быть легко решены путем последовательного исключения неизвестных.

Другой пример связан с вычислением выражения  $\log \det A$  для положительно-определенной матрицы  $A$ . Снова воспользуемся разложением Холецкого:

$$A = RR^T, \log \det A = \log \det RR^T = 2 \log \det R = 2 \sum_i \log R_{ii}.$$

Рассмотрим четыре матрицы  $A \in \mathbb{R}^{n \times n}$ ,  $C \in \mathbb{R}^{m \times m}$ ,  $U \in \mathbb{R}^{n \times m}$ ,  $V \in \mathbb{R}^{m \times n}$ . Тогда справедливы следующие утверждения:

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1} - \text{тождество Вудбери},$$

$$\det(A + UCV) = \det A \det C \det(C^{-1} + VA^{-1}U) - \text{лемма об определителе матрицы}.$$

Предположим, что  $n \gg m$  и величины  $A^{-1}$ ,  $\det A$  известны. Тогда данные утверждения позволяют свести вычисление обратной матрицы и определителя матрицы размера  $n \times n$  к вычислению обратной матрицы и определителя матрицы размера  $m \times m$ .

## Нормальное распределение

Случайная величина  $x \in \mathbb{R}$  имеет нормальное (гауссовское) распределение с параметрами  $\mu$  и  $\sigma^2$ , если ее плотность задается выражением (см. рис. 2,а)

$$p(x|\mu, \sigma^2) = \mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), \quad \sigma > 0.$$

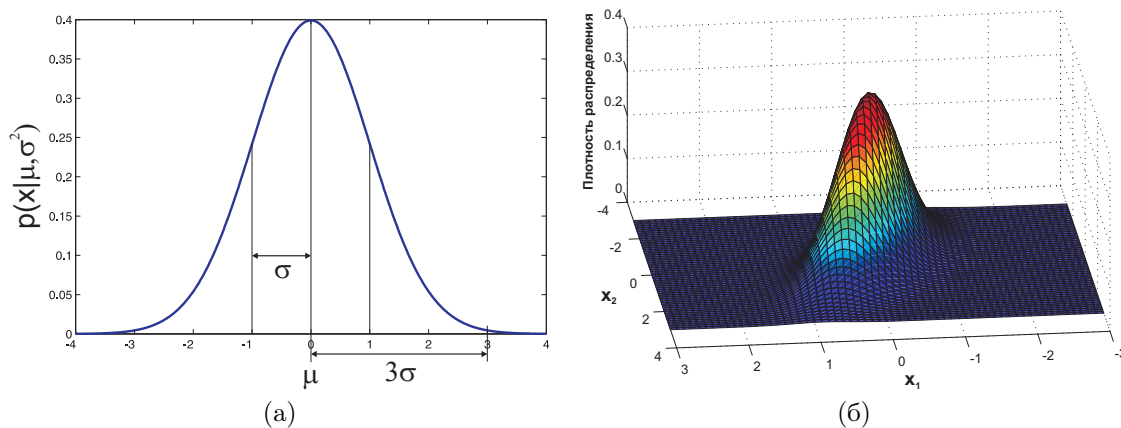


Рис. 2: Плотность одномерного (а) и многомерного (б) нормальных распределений.

Параметры  $\mu$  и  $\sigma^2$  определяют, соответственно, мат.ожидание и дисперсию нормальной случайной величины.

По центральной предельной теореме среднее арифметическое независимых случайных величин с ограниченными мат.ожиданием и дисперсией стремится к нормальному распределению. Поэтому это распределение часто используется в качестве модели шума, который определяется суммой большого количества независимых друг от друга случайных факторов.

Из неравенства Чебышева известно, что для произвольной одномерной случайной величины с конечными мат.ожиданием и дисперсией вероятность отклонения от своего мат.ожидания на  $k$

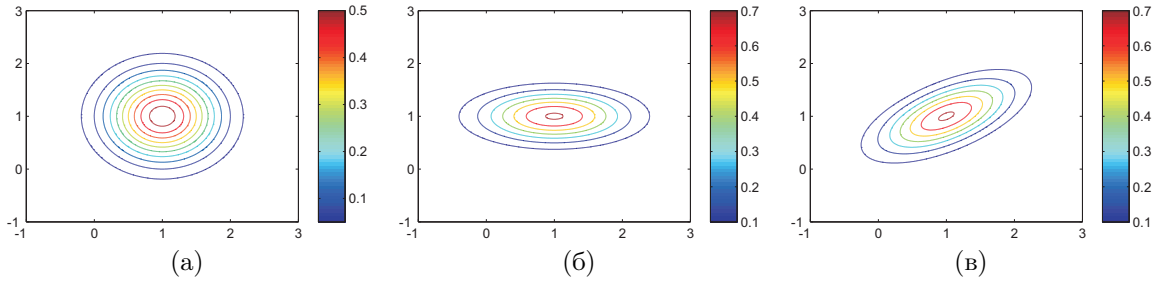


Рис. 3: Линии уровня различных нормальных распределений: (а) — нормальное распределение с матрицей ковариации, пропорциональной единичной  $\Sigma = \alpha I$ , (б) — нормальное распределение с диагональной матрицей ковариации, (в) — нормальное распределение с матрицей ковариации общего вида.

стандартных отклонений не превышает  $1/k^2$ . В частности, для  $k = 3$  эта вероятность отклонения составляет 11.2%. Однако, для нормального распределения вероятность отклонения от своего мат.ожидания на 3 стандартных отклонения составляет всего 0.3%, что намного меньше общего случая. Этот факт известен как «правило трех сигма». Таким образом, для нормального распределения большие отклонения от мат.ожидания практически невозможны. Говорят, что нормальное распределение имеет очень легкие хвосты. Это обстоятельство необходимо учитывать при приближении случайных величин нормальными.

Случайная величина  $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$  имеет многомерное нормальное распределение с параметрами  $\boldsymbol{\mu} \in \mathbb{R}^d$  и  $\Sigma \in \mathbb{R}^{d \times d}$ , если ее плотность задается выражением (см. рис. 2,б)

$$p(\mathbf{x}|\boldsymbol{\mu}, \Sigma) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \Sigma) = \frac{1}{\sqrt{2\pi}^d \sqrt{\det \Sigma}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right).$$

Параметр  $\boldsymbol{\mu}$  является мат.ожиданием, а  $\Sigma$  — матрицей ковариации нормального распределения. Матрица  $\Sigma$  является симметричной и положительно определенной (в дальнейшем положительная определенность будет обозначаться как  $\Sigma \succ 0$ ).

Нормальное распределение **полностью определяется** своими первым и вторым моментом, т.е. мат.ожиданием и матрицей ковариации. В частности, все зависимости между переменными  $x_i$  3-го и более порядков (например, выражения вида  $\mathbb{E}x_{i_1}x_{i_2}\dots x_{i_m}$ ,  $m > 2$ ) являются функциями от зависимостей порядка  $\leq 2$ .

Линии уровня плотности нормального распределения соответствуют линиям уровня квадратичной формы  $(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})$  и представляют собой эллипсы. Можно выделить три основных вида этих линий уровня (см. рис. 3) в зависимости от значения матрицы ковариации. Если матрица ковариации пропорциональна единичной, т.е. имеет вид  $\Sigma = \alpha I$ , то все компоненты нормального распределения  $x_i$  являются независимыми друг от друга и имеют одинаковую дисперсию  $\alpha$ . Линии уровня при этом образуют окружности. Диагональная матрица ковариации соответствует независимым компонентам  $x_i$ , но с различными дисперсиями. Линии уровня в этом случае являются эллипсами, параллельными координатным осям. Наконец, произвольная положительная определенная матрица ковариации соответствует эллипсам общего вида.

Матрица квадратичной формы  $(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})$  является симметричной и положительно определенной. Поэтому с помощью ортогонального преобразования координат (поворота координатных осей) эту квадратичную форму можно привести к каноническому виду с положительными коэффициентами:

$$(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}) = \left\{ \begin{array}{l} \hat{\mathbf{x}} = Q(\mathbf{x} - \boldsymbol{\mu}) \\ Q^{-1} = Q^T \end{array} \right\} = \sum_i \frac{\hat{x}_i^2}{s_i}.$$

Здесь  $s_i > 0$  — собственные значения матрицы  $\Sigma$ , а матрица поворота  $Q$  состоит из собственных векторов матрицы  $\Sigma$ . При этом значения  $\sqrt{s_i}$  определяют длины полуосей эллипса, соответствующего линии уровня нормальной плотности  $\exp(-1/2)/(\sqrt{2\pi}^d \sqrt{\det \Sigma})$ , а базис в пространстве координат  $\hat{\mathbf{x}}$  определяется собственными векторами матрицы  $\Sigma$  (см. рис. 4а).

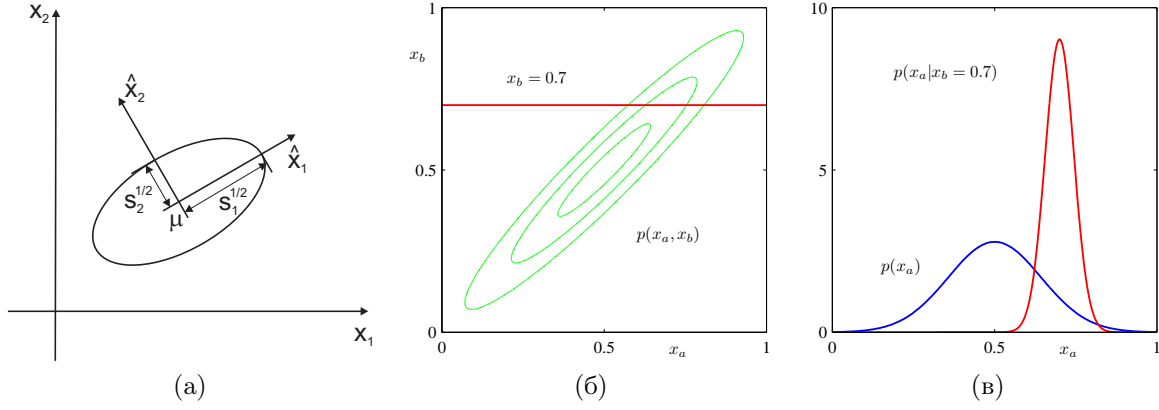


Рис. 4: (а) — линия уровня нормальной плотности, соответствующая значению  $\exp(-1/2)/(\sqrt{2\pi}^d \sqrt{\det \Sigma})$ , (б) — линии уровня нормального распределения общего вида  $p(x_a, x_b)$  в двухмерном пространстве, (в) — маргинальное распределение  $p(x_a)$  (синяя кривая) и условное маргинальное распределение  $p(x_a|x_b = 0.7)$  (красная кривая).

Разобьем вектор  $\mathbf{x}$  на две группы переменных  $\mathbf{x}_a, \mathbf{x}_b$  и обозначим

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{bmatrix}, \quad \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{bmatrix}, \quad \Lambda = \Sigma^{-1} = \begin{bmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{bmatrix}.$$

Матрицу  $\Lambda$  называют также **матрицей точности**. Тогда можно показать, что

$$\begin{aligned} p(\mathbf{x}_a) &= \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_a, \Sigma_{aa}), \\ p(\mathbf{x}_a | \mathbf{x}_b) &= \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_a - \Lambda_{aa}^{-1} \Lambda_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b), \Lambda_{aa}^{-1}). \end{aligned} \quad (1)$$

Этот результат означает, что вектор мат.ожиданий  $\boldsymbol{\mu}$  состоит из мат.ожиданий отдельных компонент  $x_i$ , а на диагонали матрицы ковариации  $\Sigma$  стоят дисперсии соответствующих компонент  $x_i$ . Кроме того, у многомерного нормального распределения все маргинальные и маргинальные условные распределения также являются нормальными (см. рис. 4б,в).

Рассмотрим величину  $\mathbf{y} \in \mathbb{R}^D$ , которая с точностью до нормального шума связана линейно с величиной  $\mathbf{x}$ , т.е.

$$\begin{aligned} p(\mathbf{y} | \mathbf{x}) &= \mathcal{N}(\mathbf{y} | A\mathbf{x}, \Gamma), \quad A \in \mathbb{R}^{D \times d}, \Gamma \in \mathbb{R}^{D \times D}, \\ p(\mathbf{x}) &= \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \Sigma), \quad \boldsymbol{\mu} \in \mathbb{R}^d, \Sigma \in \mathbb{R}^{d \times d}. \end{aligned}$$

Тогда можно показать, что

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | A\boldsymbol{\mu}, \Gamma + A\Sigma A^T), \quad (2)$$

$$p(\mathbf{x} | \mathbf{y}) = \mathcal{N}(\mathbf{x} | P(A^T \Gamma^{-1} \mathbf{y} + \Sigma^{-1} \boldsymbol{\mu}), P), \quad P = (\Sigma^{-1} + A^T \Gamma^{-1} A)^{-1}. \quad (3)$$

В частности, если  $\Gamma = 0$ , то результат (2) говорит о том, что любые линейные комбинации компонент нормального распределения также распределены нормально.

Найдем оценки максимального правдоподобия для параметров нормального распределения  $\boldsymbol{\mu}, \Sigma$ . Пусть имеется н.о.р. выборка объема  $N$   $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ , взятая из распределения  $\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \Sigma)$ . Тогда оценки максимального правдоподобия являются решениями следующей задачи:

$$p(X | \boldsymbol{\mu}, \Sigma) = \prod_n \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}, \Sigma) \rightarrow \max_{\boldsymbol{\mu}, \Sigma}.$$

Переходя к логарифму, получаем:

$$\log p(X | \boldsymbol{\mu}, \Sigma) = -\frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) - \frac{1}{2} \log \det \Sigma + \text{Const.}$$

Отсюда

$$\nabla_{\boldsymbol{\mu}} \log p(X|\boldsymbol{\mu}, \Sigma) = -N\Sigma^{-1}\boldsymbol{\mu} - \sum_{n=1}^N \Sigma^{-1}\mathbf{x}_n = -\Sigma^{-1}\left(N\boldsymbol{\mu} - \sum_{n=1}^N \mathbf{x}_n\right) = 0 \quad \Rightarrow \quad \boldsymbol{\mu}_{ML} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n.$$

Для получения оценки  $\Sigma_{ML}$  удобно перейти к матрице точности  $\Lambda$ :

$$\begin{aligned} \nabla_{\Lambda} \log p(X|\boldsymbol{\mu}, \Sigma) &= -\frac{1}{2} \nabla_{\Lambda} \text{trace}\left(\Lambda \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^T\right) + \frac{1}{2 \det \Lambda} \nabla_{\Lambda} \det \Lambda = \\ -\frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^T + \frac{1}{2 \det \Lambda} \det \Lambda \Lambda^{-1} &= 0 \quad \Rightarrow \quad \Sigma_{ML} = \Lambda_{ML}^{-1} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^T. \end{aligned}$$