

# Математические основы искусственного интеллекта

## Уравнение регрессии

Солодушкин Святослав Игоревич

Кафедра вычислительной математики и компьютерных наук,  
УрФУ имени первого Президента России Б.Н. Ельцина

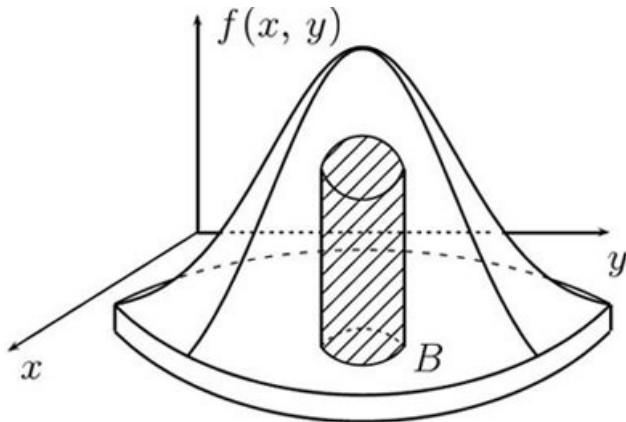
Ноябрь 2021

## Определение

Говорят, что двумерная случайная величина  $\xi, \eta$  имеет двумерное нормальное (гауссовское) распределение с параметрами, если  $\xi, \eta$  имеет следующую плотность распределения:

$$f_{\xi\eta}(x, y) = \frac{1}{2\pi\sigma_\xi\sigma_\eta\sqrt{1-r_{\xi\eta}^2}} \times \\ \times e^{-\frac{1}{2(1-r_{\xi\eta}^2)} \left[ \frac{(x-m_\xi)^2}{\sigma_\xi^2} - 2r_{\xi\eta} \frac{x-m_\xi}{\sigma_\xi} \frac{y-m_\eta}{\sigma_\eta} + \frac{(y-m_\eta)^2}{\sigma_\eta^2} \right]}, \quad (x, y) \in \mathbb{R}^2.$$

Можно доказать, что  $m_\xi, m_\eta$  — математические ожидания,  $\sigma_\xi^2, \sigma_\eta^2$  — дисперсии,  $r_{\xi\eta}$  — коэффициент корреляции случайных величин  $\xi, \eta$ .



# Некоррелированность влечет независимость

Если составляющие двумерной случайной величины некоррелированы, то они и независимы.

$$\begin{aligned} f_{\xi\eta}(x, y) &= \frac{1}{2\pi\sigma_\xi\sigma_\eta} \times e^{-\frac{1}{2}\left[\frac{(x-m_\xi)^2}{\sigma_\xi^2} + \frac{(y-m_\eta)^2}{\sigma_\eta^2}\right]} = \\ &= \frac{1}{\sqrt{2\pi}\sigma_\xi} e^{-\frac{(x-m_\xi)^2}{2\sigma_\xi^2}} \times \frac{1}{\sqrt{2\pi}\sigma_\eta} e^{-\frac{(y-m_\eta)^2}{2\sigma_\eta^2}} = f_\xi(x) f_\eta(y). \end{aligned}$$

Для нормально распределенных составляющих двумерной случайной величины понятие некоррелированности и независимости равносильны.

## Теорема

Если двумерная случайная величина  $\xi, \eta$  распределена нормально, то  $\xi, \eta$  связаны линейной связью.

Доказательство. Сделаем замену  $u = \frac{x - m_\xi}{\sigma_\xi}, \quad v = \frac{y - m_\eta}{\sigma_\eta}$ .

Получим

$$f_{\xi\eta}(u, v) = \frac{1}{2\pi\sigma_\xi\sigma_\eta\sqrt{1-r^2}} \times e^{-\frac{1}{2(1-r^2)}[u^2 - 2ruv + v^2]}.$$

Найдем маргинальную функцию распределения

$$f_\xi(u) = \int_{-\infty}^{+\infty} f_{\xi\eta}(u, v) dv = \frac{1}{\sqrt{2\pi}\sigma_\xi} e^{-\frac{u^2}{2}}.$$

# Теорема о линейности регрессии

Найдем условную плотность распределения

$$\begin{aligned} f_{\eta}(y|\xi = x) &= \frac{1}{\sqrt{2\pi}\sigma_{\eta}\sqrt{1-r^2}} e^{-\frac{(y-rx)^2}{2\sigma_{\eta}^2(1-r^2)}} = \\ &= \frac{1}{\sqrt{2\pi}\sigma_{\eta}\sqrt{1-r^2}} e^{-\frac{1}{2\sigma_{\eta}^2(1-r^2)} \left( y - \left[ m_{\eta} - r \frac{\sigma_{\eta}}{\sigma_{\xi}} (x - m_{\xi}) \right] \right)^2}. \end{aligned}$$

Условное математическое ожидание

$$E(\eta|\xi = x) = m_{\eta} - r \frac{\sigma_{\eta}}{\sigma_{\xi}} (x - m_{\xi}).$$

так что

$$\mu_{2v} = \mu'_{2v} - (\mu'_{1v})^2 = (1 - \rho^2) \left\{ 2\rho^2 v + \frac{1}{2} n (1 - \rho^2)^2 \right\}. \quad (28.23)$$

Соотношения (28.22) и (28.23) показывают, что регрессии среднего и дисперсии величины  $u$  по  $v$  линейны.

## Критерии линейности регрессии

**28.5** Пусть  $\psi(t_1, t_2) = \log \varphi(t_1, t_2)$  — совместная п. ф. с. величин  $x$  и  $y$ . Сейчас мы докажем следующий факт: если регрессия величины  $y$  по  $x$  линейна, так что

$$\mu'_{1x} = M(y|x) = \beta_0 + \beta_1 x, \quad (28.24)$$

то

$$\left[ \frac{\partial \psi(t_1, t_2)}{\partial t_2} \right]_{t_2=0} = i\beta_0 + \beta_1 \frac{\partial \psi(t_1, 0)}{\partial t_1}; \quad (28.25)$$

и наоборот, если выполнено некоторое условие полноты, то (28.25) не только необходимо, но и достаточно для (28.24).

Используя (28.24), из (28.9) при  $r = 1$  находим

$$\left[ \frac{\partial \psi(t_1, t_2)}{\partial t_2} \right]_{t_2=0} = i \int_{-\infty}^{\infty} \exp(it_1 x) g(x) (\beta_0 + \beta_1 x) dx = \quad (28.26)$$

$$= i\beta_0 \varphi(t_1, 0) + \beta_1 \frac{\partial}{\partial t_1} \varphi(t_1, 0). \quad (28.27)$$

Полагая в (28.27)  $\psi = \log \varphi$  и деля обе части на  $\varphi(t_1, 0)$ , получаем (28.25).

Обратно, если имеет место соотношение (28.25), то, используя (28.9), перепишем его в виде

$$i \int_{-\infty}^{\infty} \exp(it_1 x) (\beta_0 + \beta_1 x - \mu'_{1x}) g(x) dx = 0. \quad (28.28)$$

Теперь видим, что соотношение (28.28) влечет тождественно по  $x$

$$\beta_0 + \beta_1 x - \mu'_{1x} = 0, \quad (28.29)$$

если только семейство  $\exp(it_1 x) g(x)$  полно. Следовательно, мы получили (28.24).

# Прямая среднеквадратичной регрессии

Для упрощения математической стороны изложения, мы наложим ограничения<sup>1</sup> на случайные величины  $\xi$  и  $\eta$ . Пусть условное распределение  $\eta$  относительно своего среднего (которое, как и раньше, является функцией от  $x$ ) одно и то же для любого  $x$ , т. е. только среднее значение  $\eta$  изменяется при изменении  $x$ . Говорят, что  $\eta$  имеет «однородные ошибки».

Таким образом, существует случайная величина  $\varepsilon$  такая, что

$$\eta = M[\eta|\xi = x] + \varepsilon.$$

В частности, когда регрессия линейная, имеем

$$\eta = \beta_1\xi + \beta_0 + \varepsilon.$$

---

<sup>1</sup>Эти ограничения не сильно обременительны и в прикладных задачах обычно выполняются хотя бы приближенно.



Априори известен только вид уравнения, а конкретные значения коэффициентов  $\beta_0$  и  $\beta_1$  неизвестны, естественно выбрать их так, чтобы, зная значение, которое в эксперименте приняла величина  $\xi$ , наиболее точно спрогнозировать значение, которое примет величина  $\eta$ .

Таким образом мы приходим к постановке задачи построения линейной регрессии:

$$M[\eta - \beta_1\xi - \beta_0]^2 \xrightarrow{\beta_0, \beta_1} \min .$$

## Теорема

Линейная среднеквадратическая регрессия  $\eta$  на  $\xi$  имеет вид

$$f(x) = m_{\eta} + r \frac{\sigma_{\eta}}{\sigma_{\xi}} (x - m_{\xi}),$$

где  $m_{\eta}$ ,  $m_{\xi}$ ,  $\sigma_{\eta}$ ,  $\sigma_{\xi}$  — математические ожидания и средние квадратические отклонения случайных величин  $\eta$  и  $\xi$ , соответственно,  $r$  — коэффициент корреляции случайных величин  $\eta$  и  $\xi$ .

Рассмотрим функцию  $F(\beta_0, \beta_1) = M[\eta - \beta_1\xi - \beta_0]^2$ .

Пользуясь формулами

$$M[\xi^2] = M^2[\xi] + D[\xi]$$

$$M[\xi\eta] = M[\xi]M[\eta] + \text{cov}(\xi, \eta) = M[\xi]M[\eta] + r\sigma_\xi\sigma_\eta,$$

получим

$$F(\beta_0, \beta_1) = \sigma_\eta^2 + \beta_1^2\sigma_\xi^2 - 2r\sigma_\xi\sigma_\eta\beta_1 + (m_\eta - \beta_1m_\xi - \beta_0)^2.$$

$$F(\beta_0, \beta_1) = \sigma_\eta^2 + \beta_1^2 \sigma_\xi^2 - 2r\sigma_\xi\sigma_\eta\beta_1 + (m_\eta - \beta_1 m_\xi - \beta_0)^2.$$

Исследуем функцию  $F$  на минимум, приравняв частные производные к нулю

$$\frac{\partial F}{\partial \beta_0} = -2(m_\eta - \beta_1 m_\xi - \beta_0) = 0,$$

$$\frac{\partial F}{\partial \beta_1} = 2\beta_1 \sigma_\xi^2 - 2r\sigma_\xi\sigma_\eta = 0.$$

Откуда получаем

$$\beta_0 = m_\eta - r \frac{\sigma_\eta}{\sigma_\xi} m_\xi, \quad \beta_1 = r \frac{\sigma_\eta}{\sigma_\xi}.$$

## Интеграл Пуассона

$$I = \int_{-\infty}^{\infty} e^{-x^2/2} dx = \sqrt{2\pi}.$$

$$I = \int_{-\infty}^{\infty} e^{-ax^2} dx = \frac{1}{\sqrt{a}} \Gamma\left(\frac{1}{2}\right).$$

Для сведения  $\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}.$

Двумерная случайная величина задана плотностью распределения

$$f_{\xi\eta}(x, y) = \frac{3\sqrt{3}}{\pi} e^{-4x^2 - 6xy - 9y^2}.$$

Найти условные законы распределения составляющих. Найти уравнение регрессии  $\eta$  на  $\xi$ .