



Прикладная статистика и анализ данных

Съезд II

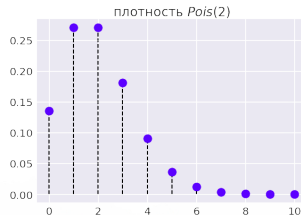


Обобщенная модель линейной регрессии

Пуассоновское распределение

$$\text{Pois}(\lambda) : p(x) = \frac{\lambda^x}{x!} e^{-\lambda}, x \in \mathbb{Z}_+$$

Смысл: число событий,
произошедших за единицу времени



Условия:

1. события происходят с фиксированной интенсивностью λ .
2. независимо друг от друга.

Утверждение: время между двумя событиями имеет распр. $\text{Exp}(\lambda)$
(см. пуассоновские случайные процессы)

Примеры:

1. число клиентов в час
2. число запросов на сервер за минуту

Интенсивность не постоянна, может зависеть от каких-то факторов.



Обобщенная модель линейной регрессии

Гауссовская линейная модель

Ожидаемый отклик:

$$y = \mu_{\theta}(x) = x^T \theta.$$

Наблюдаемый отклик:

$$Y_i = x_i^T \theta + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, \sigma^2).$$

$$\text{или } Y_i \sim \mathcal{N}(\mu_{\theta}(x_i), \sigma^2)$$

Оценка отклика:

$$\hat{y} = x^T \hat{\theta}.$$

Generalized Linear Models (GLM)

Ожидаемый отклик:

$$y = \mu_{\theta}(x), \text{ причем } g(\mu_{\theta}(x)) = x^T \theta,$$

т.е. g — линейризация ожид. отклика

Наблюдаемый отклик:

$$Y_i \sim P_{\mu_{\theta}(x_i)},$$

где $\{P_{\psi} \mid \psi \in \Psi\}$ — семейство распр.

Оценка отклика:

$$\hat{y} = g^{-1} \left(x^T \hat{\theta} \right).$$



Натуральный отклик

$y \in \mathbb{Z}_+$ — значения наблюдаемого отклика

$\mu_\theta(x) = E_x Y$ — ожидаемый отклик

$Y_i \sim \text{Pois}(\mu_\theta(x_i)) = P_{\mu_\theta(x_i)}$ — наблюдаемый отклик

Линеаризация ожид. отклика: $g(z) = \ln z$

т.е. $g(\mu_\theta(x)) = \ln E_x Y = x^T \theta$

Тогда $\mu_\theta(x) = g^{-1}(x^T \theta) = \exp(x^T \theta)$

\Rightarrow это **пуассоновская регрессия**

Смысл: пытаемся приблизить интенсивность с помощью регрессии.

Примечание. Интенсивность имеет неравномерный масштаб. Пусть $\lambda = 1$.

В два раза чаще это $\lambda = 2$, в два раза реже это $\lambda = 1/2$.

Логарифмирование это исправляет.

Свойства GLM

В качестве $\hat{\theta}$ берется ОМП (ищется численно)

$$L_X(\theta) = \prod_{i=1}^n p_{\mu_\theta(x_i)}(Y_i) = \prod_{i=1}^n p_{g^{-1}(x_i^T \theta)}(Y_i) \rightarrow \max_{\theta}$$

Если $\{P_\psi \mid \psi \in \Psi\}$ лежит в экспоненциальном классе, то $\hat{\theta}$:

1. существует и единственна;
2. состоятельна;
3. асимптотически нормальна: $\sqrt{I(\theta)} (\hat{\theta} - \theta) \xrightarrow{d} \mathcal{N}(0, I_d)$,
где $I(\theta) = \left(-E \frac{\partial^2 \log L_X(\theta)}{\partial \theta_j \partial \theta_k} \right)_{jk}$ — информационная матрица Фишера.

Для пуассоновской регрессии $l(\theta) = \sum_{i=1}^n e^{x_i^T \theta} x_i x_i^T = X^T V(\theta) X$,
где $V(\theta) = \text{diag} [e^{x_i^T \theta}]$.

Асимпт. доверительные интервалы в GLM

Для параметров (\implies критерий для гипотезы $H_0: \theta_j = 0$)

$$\theta_j \in \left(\hat{\theta}_j \pm z_{1-\alpha/2} \sqrt{\left(I^{-1}(\hat{\theta}) \right)_{jj}} \right)$$

Для преобразованного ожидаемого отклика

$$x_0^T \theta \in \left(x_0^T \hat{\theta} - \delta, x_0^T \hat{\theta} + \delta \right)$$

Для ожидаемого отклика

$$\mu(x_0) = g^{-1}(x_0^T \theta) \in \left[g^{-1} \left(x_0^T \hat{\theta} - \delta \right), g^{-1} \left(x_0^T \hat{\theta} + \delta \right) \right],$$

$$\delta = z_{1-\alpha/2} \sqrt{x_0^T I^{-1}(\hat{\theta}) x_0}$$