

## СЕМ, лекция 15

(2021-01-21)

### Припомняне:

- Имаме случайна величина  $X$ , която искаме „да разберем“, т.е. да извлечем някаква информация за нея, и чиято функция на разпределение зависи от някакъв параметър  $\theta$  (едномерен);
- $\vec{X}$  са някакви наблюдения дадени като  $n$ -мерен вектор;
- $\alpha$  (алфа) е предварително зададена грешка от I<sup>-ви</sup> род;
- Тестваме две прости хипотези:  
 $H_0 : \theta = \theta_0$  (базова)  
 $H_1 : \theta = \theta_1$  (алтернативна)

Означаваме:  $L_0(x) = L(x; \theta_0)$  (функция на правдоподобие за  $\theta = \theta_0$ ) и  $L_1(x) = L(x; \theta_1)$  (функция на правдоподобие за  $\theta = \theta_1$ ).

Търсим  $W^* \subseteq \mathbb{R}^n$  (област на  $\mathbb{R}^n$ ), така, че когато нашия  $n$ -мерен вектор от наблюдения попадне в него, ние да отхвърляме нулевата хипотеза и да приемаме алтернативната.  $\alpha = \mathbb{P}(\vec{X} \in W^* | H_0)$  е грешка от първи род, т.е. да попаднем в  $W^*$  и да отхвърлим нулевата хипотеза, но тя да е била вярна. Тази грешка е презададена (предефинирана) и се контролира от изследователя. Целта на оптималната критична област  $W^*$  е да намери тази област, която минимизира грешката от първи род.

$$\beta = \min_{\alpha = \mathbb{P}(\vec{X} \in W | H_0)} \mathbb{P}(\vec{X} \notin W | H_1).$$

Лемата на Нейман-Пийрсън е „добра“, защото ни характеризира даден критерий, по който да определим дали една област е оптимална критична област и по тази лема знаем, че  $W^*$  е о.к.о. (оптимална критична област), ако съществува някаква константа  $K$  ( $K$  може да зависи от  $n$  и от  $\theta$ , но не може да зависи от  $x$ ), за която

$$W^* \in \{x \in \mathbb{R}^n : L_1(x) > K \times L_0(x)\}$$
$$W^{*c} \in \{x \in \mathbb{R}^n : L_1(x) \leq K \times L_0(x)\}$$

и ако знаем, че е изпълнено равенството:  $\alpha = \mathbb{P}(\vec{X} \in W^* | H_0)$ , то  $W^*$  е о.к.о.

$\oplus X \in \mathcal{N}(\mu, \sigma^2)$ , където  $\sigma^2$  е известно и искаме да построим оптимална критична област за тестване на хипотезата на  $\mu$  (за намиране на средното, знаейки каква е дисперсията)

$$H_0 : \mu = \mu_0$$
$$H_1 : \mu = \mu_1$$

Допускаме за улеснение, че  $\mu_1 > \mu_0$ . При зададено  $\alpha$ .

$$L_0(x) = (\sqrt{2\pi}\sigma)^{-n} e^{-\frac{\sum_{j=1}^n (x_j - \mu_0)^2}{2\sigma^2}}, L_1(x) = (\sqrt{2\pi}\sigma)^{-n} e^{-\frac{\sum_{j=1}^n (x_j - \mu_1)^2}{2\sigma^2}}.$$

От лемата на Нейман-Пиърсън знаем, че оптималните критични области се намират лесно с неравенства от вида:

$$\{X \in \mathbb{R}^n : L_1(x) \geq K \times L_0(x)\} = \left\{ X \in \mathbb{R}^n : \frac{-\sum_{j=1}^n (x_j - \mu_1)^2}{2\sigma^2} \geq \ln K - \frac{\sum_{j=1}^n (x_j - \mu_0)^2}{2\sigma^2} \right\} =$$

$$= \left\{ X \in \mathbb{R}^n : -\frac{\sum_{j=1}^n x_j^2}{2\sigma^2} + \frac{1}{\sigma^2} \mu_1 \sum_{j=1}^n x_j - \frac{n\mu_1^2}{2\sigma^2} \geq \right.$$

не зависи от  
наблюденията  $\vec{X}$

$$\geq -\frac{\sum_{j=1}^n x_j^2}{2\sigma^2} + \frac{1}{\sigma^2} \mu_0 \sum_{j=1}^n x_j - \frac{n\mu_0^2}{2\sigma^2} \left. \right\} =$$

не зависи от  
наблюденията  $\vec{X}$

$$= \left\{ X \in \mathbb{R}^n : \frac{1}{\sigma^2} \underbrace{(\mu_1 - \mu_0)}_{\mu_1 > \mu_0} \sum_{j=1}^n x_j \geq K_1 \right\} = , \text{ където } K_1 = \ln K - \frac{n\mu_0^2}{2\sigma^2} + \frac{n\mu_1^2}{2\sigma^2}.$$

по допускане

$$= \left\{ X \in \mathbb{R}^n : \sum_{j=1}^n x_j \geq K_2 \right\} = , \text{ където } K_2 = \frac{K_1 \sigma^2}{\mu_1 - \mu_0}$$

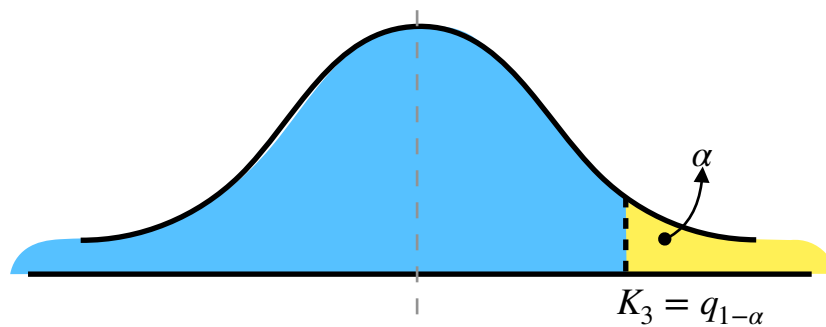
$$= \left\{ X \in \mathbb{R}^n : \frac{\sum_{j=1}^n x_j}{n} \geq \frac{K_2}{n} \right\} = \left\{ X \in \mathbb{R}^n : \bar{X} \geq \frac{K_2}{n} \right\} =$$

$$= \left\{ X \in \mathbb{R}^n : \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \geq \frac{K_2}{\sigma\sqrt{n}} = K_3 \right\}.$$

$$\alpha = \mathbb{P}(\vec{X} \in W | H_0) = \mathbb{P}(\underbrace{L_1(\vec{X}) \geq K_0 L_0(\vec{X})}_{\text{критична област}} | H_0) = \mathbb{P}\left(\frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \geq K_3 | H_0\right),$$

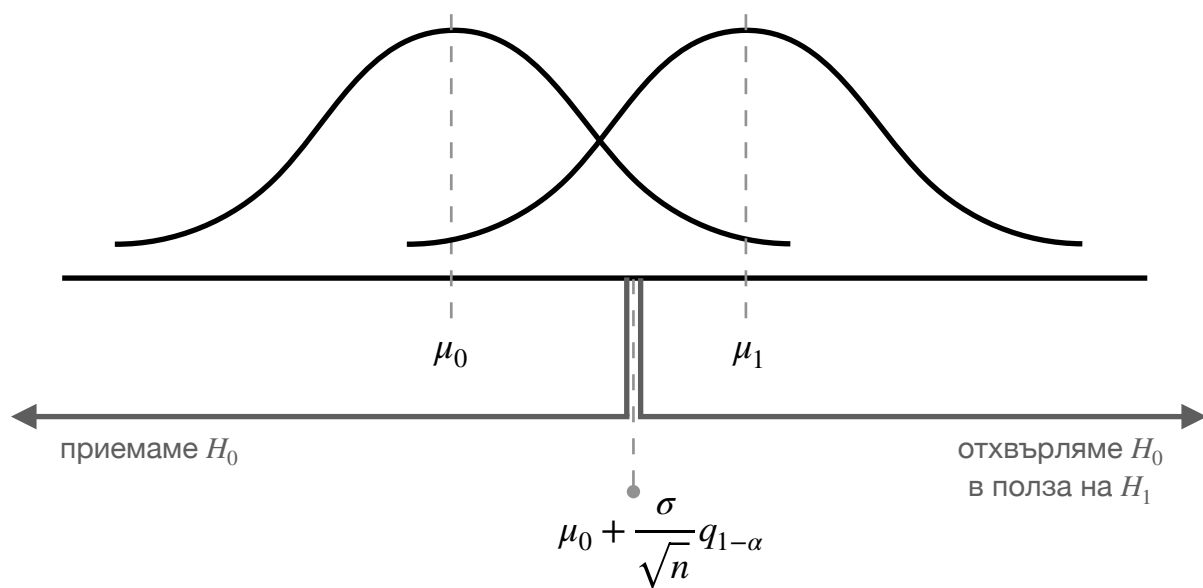
$$\text{където } \bar{X} = \frac{\sum_{j=1}^n x_j}{n}.$$

$$Z = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \in \mathcal{N}(0,1) = \mathbb{P}(Z \geq K_3)$$



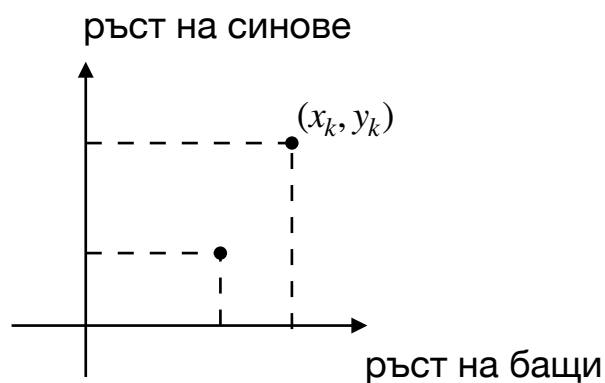
$$\Rightarrow K_3 = q_{1-\alpha}$$

$$\text{o.k.o.: } \left\{ \bar{X} \geq \mu_0 + q_{1-\alpha} \cdot \frac{\sigma}{\sqrt{n}} \right\}$$



### Линейна регресия (Галтон)

Нека  $a$  е средния ръст на мъжете.



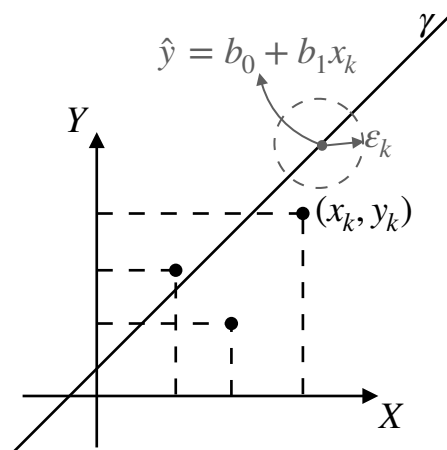
$$y_{\text{син}} = a + \beta(x_{\text{баща}} - a).$$

Галтон е забелязал, че по неговите данни, коефициента бета е  $\beta = 0.6$ . Тоест, ако бащата е 10 см. над средния ръст, то сина му ще е с 6 см. над средния ръст. Тази по-слаба зависимост е влязла в теорията като регрес (завръщане) към средното. Синовите на високите бащи не са чак толкова високи в средно както бащите им, а са на около половината от отклонението на бащата над средния ръст за мъжете.

Модел:  $Y = \beta_0 + \beta_1 X + \varepsilon$

↑                      ↑                      ↗

отклик                      предиктор                      грешка



Допускаме че в множеството от точки  $(x_1, x_2, \dots, x_n)$  и  $(y_1, y_2, \dots, y_n)$  има някакъв линейен модел. Т.е. предполагаме, че има линейен модел  $y_k = b_0 + b_1 x_k + \varepsilon_k$ . Тоест имаме някаква права  $\gamma$  (от чертежа). Искаме да си построим линейен модел, а не някакъв друг, за да не рискуваме да интерполираме, тъй като интерполацията няма добра статистическа стойност. Т.е. не е добре да обхванем всички данни с много сложна крива и в момента, в който добавим данни – кривата ни да е твърде динамична и да няма никаква прогнозна сила. Интерполацията може да мине през  $n$  точки, но при добавянето на  $n + 1$ -вата точка – точността на кривата да рухне.

$$\bar{X} = \frac{1}{n} \sum_{k=1}^n x_k \text{ и } \bar{Y} = \frac{1}{n} \sum_{k=1}^n y_k.$$

$$\text{Търсим: } \min_{b_0, b_1} \sum_{k=1}^n (y_k - \hat{y}_k)^2 = \min_{b_0, b_1} \sum_{k=1}^n (y_k - b_0 - b_1 x_k)^2.$$

(С квадратични грешки се смята по-лесно, а освен това имат и статистическо значение. Въпреки това тук може да имаме най-разнообразни метрики, които искаме да оптимизираме (например абсолютната стойност или максималното отклонение измежду всички възможни отклонения и т.н.))

Искаме да минимизираме функцията по-горе по две променливи. За целта ще си вземем производната по  $b_0$  и тя трябва да бъде нула и аналогично за производната по  $b_1$ :

$$0 = \frac{\partial}{\partial b_0} \sum_{k=1}^n (y_k - b_0 - b_1 x_k)^2 = -2 \sum_{k=1}^n (y_k - b_0 - b_1 x_k) \Rightarrow$$

$$\Rightarrow \sum_{k=1}^n y_k - n b_0 - b_1 \sum_{k=1}^n x_k = 0 \Rightarrow n \bar{Y} - n b_0 - n b_1 \bar{X} = 0 \Rightarrow \bar{Y} = b_0 + b_1 \bar{X}$$

$$0 = \frac{\partial}{\partial b_1} \sum_{k=1}^n (y_k - b_0 - b_1 x_k)^2 = -2 \sum_{k=1}^n (b_0 + b_1 x_k - y_k) x_k \Rightarrow$$

$$\Rightarrow n \bar{X} (\bar{Y} - b_1 \bar{X}) + b_1 \sum_{k=1}^n x_k^2 - \sum_{k=1}^n x_k y_k$$

$$\begin{cases} b_0 = \bar{Y} - b_1 \bar{X} \\ b_1 = \frac{\sum_{k=1}^n x_k y_k - n \bar{X} \bar{Y}}{\sum_{k=1}^n x_k^2 - n (\bar{X})^2} = \frac{\sum_{k=1}^n (y_k - \bar{Y})(x_k - \bar{X})}{\sum_{k=1}^n (x_k - \bar{X})^2} = \frac{\sum_{k=1}^n (x_k - \bar{X}) y_k}{\sum_{k=1}^n (x_k - \bar{X})^2} \end{cases}$$

$$A = \sum_{k=1}^n (X_k - \bar{X})^2$$

Допускаме, че  $Y_k$  като отговор на  $X_k$  е случайна величина, в смисъл, че съществуват неизвестни коефициенти  $\beta_0, \beta_1$ , които при зададено  $X_k$  дават следната линейна зависимост, където  $\varepsilon_k$  е случайна грешка.

Допускаме, че  $Y_k = \beta_0 + \beta_1 X_k + \varepsilon_k$ , където  $k = \overline{1, n}$ .

Правим и следните допускания за епсилон грешките:

$(\varepsilon_i)_{i=1}^n$  са независими еднакво разпределени случайни величини, като  $\varepsilon_i \in \mathcal{N}(0, \sigma^2)$

хомоскедастичност

Т.е. грешките са нормално разпределени и независими една от друга. Т.е. нямаме системна грешка. Хомоскедастичността е малко по-тежко допускане, но тя придава простота на модела.

$$Y_k \in \mathcal{N}(\beta_0 + \beta_1 X_k, \sigma^2)$$

$$\begin{cases} \hat{\beta}_0 = b_0 = \bar{Y} - b_1 \bar{X} \\ \hat{\beta}_1 = b_1 = \frac{\sum_{k=1}^n X_k Y_k - n \bar{X} \bar{Y}}{\sum_{k=1}^n X_k^2 - n (\bar{X})^2} = \frac{\sum_{k=1}^n (Y_k - \bar{Y})(X_k - \bar{X})}{\sum_{k=1}^n (X_k - \bar{X})^2} = \frac{\sum_{k=1}^n (X_k - \bar{X}) Y_k}{\sum_{k=1}^n (X_k - \bar{X})^2} \end{cases}$$

(1)                      (2)                      (3)

$$A = \sum_{k=1}^n (X_k - \bar{X})^2 \text{ си остава същото, тъй като е фиксирано число в знаменателя.}$$

$$\begin{aligned}\mathbb{E}b_1 &\stackrel{(3)}{=} \frac{1}{A} \sum_{k=1}^n (X_k - \bar{X}) \mathbb{E}Y_k = \frac{1}{A} \sum_{k=1}^n (X_k - \bar{X})(\beta_0 - \beta_1 X_k) = \\ &= \underbrace{\frac{\beta_0}{A} \sum_{k=1}^n (X_k - \bar{X})}_{=0} + \underbrace{\frac{\beta_1}{A} \sum_{k=1}^n (X_k - \bar{X})X_k}_{=A} = \beta_1.\end{aligned}$$

Оказва се, че очакването на  $b_1$  е равно на  $\beta_1$ , което ни казва, че  $b_1$  е неизместена оценка на  $\beta_1$ .

$$\begin{aligned}\mathbb{E}b_0 &= \mathbb{E}\bar{Y} - \bar{X}\mathbb{E}b_1 = \frac{1}{n} \sum_{k=1}^n Y_k - \beta_1 \bar{X} = \\ &= \frac{1}{n} \sum_{k=1}^n (\beta_0 + \beta_1 X_k) - \beta_1 \bar{X} = \beta_0.\end{aligned}$$

И  $b_0$  и  $b_1$  са неизместени оценки на неизвестните параметри  $\beta_0$  и  $\beta_1$ . По този начин знаем, че нямаме систематична грешка, когато правим тези оценки.

$$\begin{aligned}Db_1 &\stackrel{(3)}{=} D \frac{\sum_{k=1}^n (X_k - \bar{X})Y_k}{A} = \frac{1}{A^2} \sum_{k=1}^n (X_k - \bar{X})^2 \underbrace{DY_k}_{\sigma^2} = \\ &= \sigma^2 \frac{\sum_{k=1}^n (X_k - \bar{X})^2}{A^2} = \frac{\sigma^2}{A}.\end{aligned}$$

$$\Rightarrow b_1 = \underset{\text{оценка}}{\hat{\beta}_1} \in \mathcal{N}\left(\beta_1, \frac{\sigma^2}{A}\right)$$

Това означава, че вече може да тестваме хипотези за  $b_1$ .

За дисперсията на  $b_0$  по същата логика може да докажем, че:

$$Db_0 = \sigma^2 \left( \frac{1}{n} + \frac{\bar{X}^2}{A} \right) \Rightarrow b_0 = \underset{\text{оценка}}{\hat{\beta}_0} \in \mathcal{N}\left(\beta_0, \sigma^2 \left( \frac{1}{n} + \frac{\bar{X}^2}{A} \right)\right).$$

Двете дисперсии клонят към нула.

Оценка на  $\sigma^2$  (ако не го знаем априорно):  $Y_k \in \mathcal{N}(\beta_0 + \beta_1 X_k, \sigma^2)$ . Проблема е, че не знаем  $\beta_0$  и  $\beta_1$ , тъй като, ако допуснем, че ги знаем, щяхме да имаме  $\frac{Y_k - \beta_0 - \beta_1 X_k}{\sigma} \in \mathcal{N}(0,1)$  и тогава

$$\frac{\sum_{k=1}^n (Y_k - \beta_0 - \beta_1 X_k)^2}{\sigma^2} \in \mathcal{X}^2(n).$$

Но, ако са ни верни допусканията за модела, тогава:

$$\frac{\sum_{k=1}^n (Y_k - b_0 - b_1 X_k)^2}{\sigma^2} \in \mathcal{X}^2(n-2) \text{ („изхабили“ (използвали) сме две степени на}$$

свобода (две данни), за да оценим  $b_0$  и  $b_1$ )

$$\mathbb{E} \frac{\sum_{k=1}^n (Y_k - b_0 - b_1 X_k)^2}{\sigma^2} = n - 2$$

$$\hat{\sigma}^2 = \frac{\sum_{k=1}^n (Y_k - b_0 - b_1 X_k)^2}{n - 2}, \text{ т.е. } \sigma^2 = \mathbb{E} \hat{\sigma}^2.$$

Оттук нататък ние може да тестваме хипотези. Може да си конструираме множество хипотези от следния вид:

$$H_0 : \beta_1 = \tilde{\beta}$$

$$H_1 : \beta_1 = \tilde{\tilde{\beta}}$$

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

$$\frac{b_1 - \tilde{\beta}}{\sqrt{\sigma^2/A}} \in \mathcal{N}(0,1) \text{ при } H_0.$$