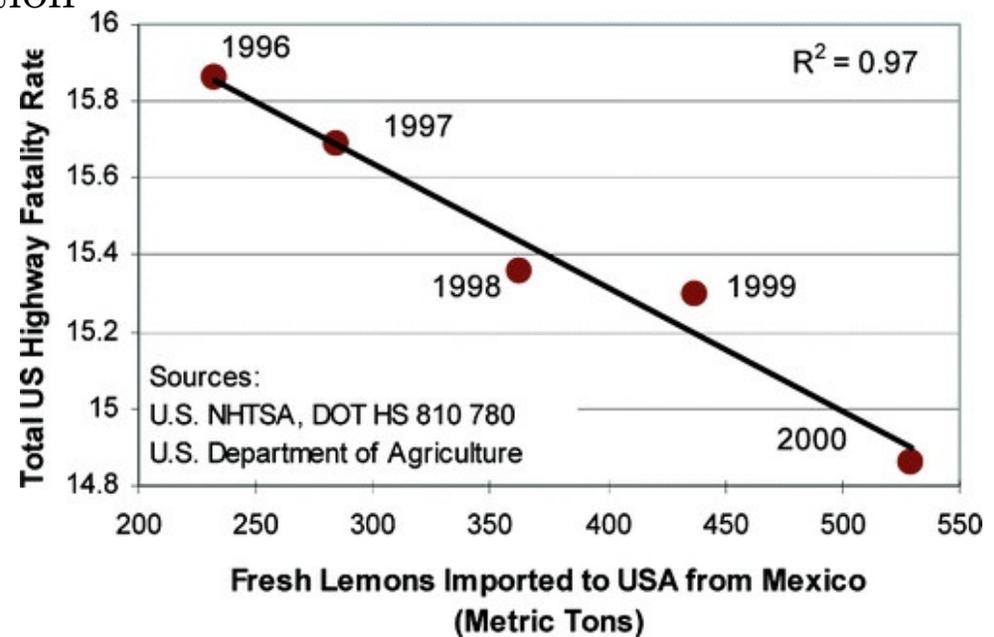


Slides 16: Simple regression model

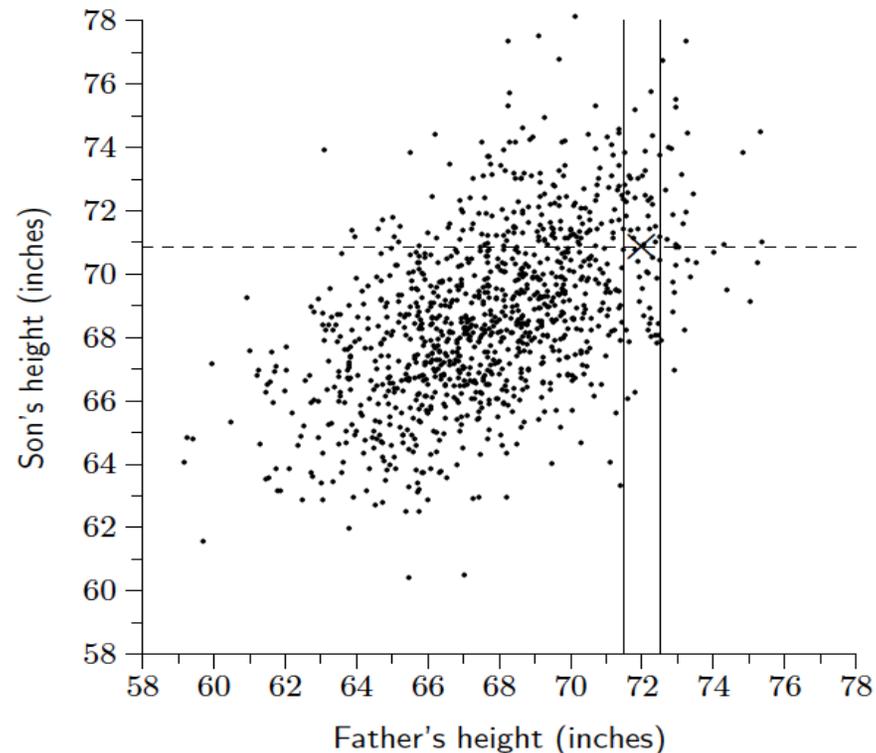
- Regression to mediocrity
- Least squares estimates
- Coefficient of determination
- Residuals
- Confidence intervals
- Model utility test
- Prediction interval



Correlation does not imply causation

Regression to mediocrity

Pearson's father-son data: 1,078 pairs of heights (England, 1900).



Focussing on 6 feet tall fathers, we see that their sons on average are shorter than their fathers. F. Galton called this *regression to mediocrity*.

Question. Which of the previous statistical tools could be applied?

Simple linear regression model

A simple linear regression model is based on the linear relation

$$Y(x) = \beta_0 + \beta_1 x + \epsilon, \quad \epsilon \sim N(0, \sigma),$$

where ϵ is the noisy part of the response, that is not explained by the value x of the main explanatory variable. The assumption of *homoscedasticity* requires that σ is independent of the x -value.

For a given collection of x -values (x_1, \dots, x_n) , and a vector (e_1, \dots, e_n) of independent realisations of ϵ , we get a sample of response values

$$y_i = \beta_0 + \beta_1 x_i + e_i, \quad i = 1, \dots, n.$$

The likelihood is a function of the 3D parameter $\theta = (\beta_0, \beta_1, \sigma^2)$

$$L(\theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma}} \exp \left\{ -\frac{(y_i - \beta_0 - \beta_1 x_i)^2}{2\sigma^2} \right\} = C \sigma^{-n} e^{-\frac{S(\beta_0, \beta_1)}{2\sigma^2}},$$

where

$$C = (2\pi)^{-n/2}, \quad S(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2.$$

Maximum likelihood estimates

Log-likelihood function

$$l(\theta) = \ln L(\theta) = \ln C - n \ln \sigma - \frac{S(\beta_0, \beta_1)}{2\sigma^2}$$

Maximisation of $l(\theta)$ over (β_0, β_1) is equivalent to minimisation of the sum of squares $S(\beta_0, \beta_1)$.

Therefore, the MLEs of (β_0, β_1) are called the least squares estimates.

Observe that

$$\begin{aligned} \frac{S(\beta_0, \beta_1)}{n} &= \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \\ &= \beta_0^2 + 2\beta_0\beta_1\bar{x} - 2\beta_0\bar{y} - 2\beta_1\overline{xy} + \beta_1^2\overline{x^2} + \bar{y}^2 \end{aligned}$$

with the following set of five sufficient statistics:

$$\begin{aligned} \bar{x} &= \frac{x_1 + \dots + x_n}{n}, & \bar{y} &= \frac{y_1 + \dots + y_n}{n} \\ \overline{x^2} &= \frac{x_1^2 + \dots + x_n^2}{n}, & \overline{y^2} &= \frac{y_1^2 + \dots + y_n^2}{n}, & \overline{xy} &= \frac{x_1 y_1 + \dots + x_n y_n}{n} \end{aligned}$$

Normal equations

To obtain MLEs of $\theta = (\beta_0, \beta_1, \sigma^2)$ compute the derivatives

$$\begin{aligned}\frac{\partial l}{\partial \beta_0} &= -\frac{1}{2\sigma^2} \frac{\partial S}{\partial \beta_0}, \\ \frac{\partial l}{\partial \beta_1} &= -\frac{1}{2\sigma^2} \frac{\partial S}{\partial \beta_1}, \\ \frac{\partial l}{\partial \sigma^2} &= -\frac{n}{2\sigma^2} + \frac{S(\beta_0, \beta_1)}{2\sigma^4},\end{aligned}$$

and set them equal to zeros.

Putting $\frac{\partial S}{\partial \beta_0} = 0$ and $\frac{\partial S}{\partial \beta_1} = 0$, we get the so-called normal equations:

$$b_0 + b_1 \bar{x} = \bar{y}, \quad b_0 \bar{x} + b_1 \overline{x^2} = \overline{xy}.$$

Solving this system of linear equations we get

$$b_1 = \frac{\overline{xy} - \bar{x}\bar{y}}{\overline{x^2} - \bar{x}^2} = \frac{rs_y}{s_x}, \quad b_0 = \bar{y} - b_1 \bar{x},$$

where r is the sample correlation coefficient and

$$s_x^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2, \quad s_y^2 = \frac{1}{n-1} \sum (y_i - \bar{y})^2.$$

The sample correlation coefficient

The sample correlation coefficient r is an unbiased estimate of ρ

$$r = \frac{s_{xy}}{s_x s_y}, \quad \rho = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y},$$

where the sample covariance s_{xy} is an unbiased estimate of $\text{Cov}(X, Y)$

$$s_{xy} = \frac{1}{n-1} \sum (x_i - \bar{x})(y_i - \bar{y}), \quad \text{Cov}(X, Y) = \text{E}(X - \mu_x)(Y - \mu_y),$$

provided that (x_1, \dots, x_n) is a random sample from X -distribution.

As a result, the fitted regression line $y = b_0 + b_1 x$ takes the form

$$y = \bar{y} + r \frac{s_y}{s_x} (x - \bar{x}),$$

Notice that

$$S(b_0, b_1) = \sum_{i=1}^n (y_i - \hat{y}_i)^2,$$

where \hat{y}_i are the predicted responses:

$$\hat{y}_i = b_0 + b_1 x_i, \quad i = 1, \dots, n.$$

Estimating the size of the noise

Putting $\frac{\partial l}{\partial \sigma^2} = 0$, we get

$$0 = -\frac{n}{2\sigma^2} + \frac{S(\beta_0, \beta_1)}{2\sigma^4}$$

and replacing (β_0, β_1) with (b_0, b_1) , we find the MLE of σ^2 to be

$$\hat{\sigma}^2 = \frac{S(b_0, b_1)}{n},$$

The maximum likelihood estimate of $\hat{\sigma}^2$ is only asymptotically unbiased estimate of σ^2 . An unbiased estimate of σ^2 is given by

$$s^2 = \frac{S(b_0, b_1)}{n-2}.$$

The error sum of squares

$$SS_E = S(b_0, b_1) = \sum (y_i - \hat{y}_i)^2 = (n-1)s_y^2(1-r^2)$$

divided by $n-2$ gives a very useful expression

$$s^2 = \frac{n-1}{n-2} s_y^2(1-r^2).$$

Question. If $r = 0.5$, what proportion of s_y^2 is explained by s^2 ?

Residuals

Residuals are differences between observed and predicted responses

$$\hat{e}_i = y_i - \hat{y}_i = (y_i - \bar{y}) - r \frac{s_y}{s_x} (x_i - \bar{x})$$

The residuals $(\hat{e}_1, \dots, \hat{e}_n)$ are linearly connected via

$$\hat{e}_1 + \dots + \hat{e}_n = 0, \quad x_1 \hat{e}_1 + \dots + x_n \hat{e}_n = 0, \quad \hat{y}_1 \hat{e}_1 + \dots + \hat{y}_n \hat{e}_n = 0,$$

so we can say that \hat{e}_i are uncorrelated with x_i and \hat{e}_i are uncorrelated with \hat{y}_i . The residuals \hat{e}_i are realisations of random variables \hat{E}_i having normal distributions with zero means and

$$\text{Var}(\hat{E}_i) = \sigma^2 \left(1 - \frac{\sum_k (x_k - x_i)^2}{n(n-1)s_x^2} \right), \quad \text{Cov}(\hat{E}_i, \hat{E}_j) = -\sigma^2 \cdot \frac{\sum_k (x_k - x_i)(x_k - x_j)}{n(n-1)s_x^2}.$$

Test normality using normal QQ-plot for the standardised residuals

$$\tilde{e}_i = \frac{\hat{e}_i}{s_i}, \quad s_i = s \sqrt{1 - \frac{\sum_k (x_k - x_i)^2}{n(n-1)s_x^2}}, \quad i = 1, \dots, n,$$

where s_i are the estimated standard deviations of \hat{E}_i . In some cases, the non-linearity problem can be fixed by a log-log transformation of the data.

Coefficient of determination r^2

Using $y_i - \bar{y} = \hat{y}_i - \bar{y} + \hat{e}_i$, we obtain a decomposition

$$SS_T = SS_R + SS_E,$$

where

$$SS_T = \sum_i (y_i - \bar{y})^2 = (n - 1)s_y^2$$

is the total sum of squares, and

$$SS_R = \sum_i (\hat{y}_i - \bar{y})^2 = (n - 1)b_1^2 s_x^2 = (n - 1)r^2 s_y^2$$

is the regression sum of squares. Combining these relations, we find that

$$r^2 = \frac{SS_R}{SS_T} = 1 - \frac{SS_E}{SS_T}.$$

Coefficient of determination r^2 is the proportion of variation in the response variable explained by the variation of the predictor.

Observe that r^2 has a more intuitive meaning than the sample correlation coefficient r .

Confidence intervals and hypothesis testing

The least squares estimators (b_0, b_1) are unbiased and consistent. Due to the normality assumption we have the following exact distributions

$$B_0 \sim N(\beta_0, \sigma_0), \quad \sigma_0^2 = \frac{\sigma^2 \sum x_i^2}{n(n-1)s_x^2}, \quad s_{b_0}^2 = \frac{s^2 \sum x_i^2}{n(n-1)s_x^2}, \quad \frac{B_0 - \beta_0}{S_{B_0}} \sim t_{n-2},$$
$$B_1 \sim N(\beta_1, \sigma_1), \quad \sigma_1^2 = \frac{\sigma^2}{(n-1)s_x^2}, \quad s_{b_1}^2 = \frac{s^2}{(n-1)s_x^2}, \quad \frac{B_1 - \beta_1}{S_{B_1}} \sim t_{n-2}.$$

There is a weak correlation between the two estimators:

$$\text{Cov}(B_0, B_1) = -\frac{\sigma^2 \bar{x}}{(n-1)s_x^2}$$

which is negative, if $\bar{x} > 0$, and positive, if $\bar{x} < 0$.

Exact $100(1 - \alpha)\%$ confidence intervals $I_{\beta_i} = b_i \pm t_{n-2}\left(\frac{\alpha}{2}\right) \cdot s_{b_i}$

For $i = 0$ or $i = 1$ and a given value β^* , one would like to test the null hypothesis $H_0: \beta_i = \beta^*$. Use the test statistic

$$t = \frac{b_i - \beta^*}{s_{b_i}},$$

having the exact null distribution $T \sim t_{n-2}$.

Model utility test

Two important examples of hypothesis testing for the linear regression.

1. Model utility test is built around the null hypothesis

$$H_0 : \beta_1 = 0$$

stating that there is no relationship between the predictor variable x and the response y . The corresponding test statistic, called t-value,

$$t = \frac{b_1}{s_{b_1}} = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}}$$

has an exact t_{n-2} null distribution.

2. Zero-intercept test aims at

$$H_0 : \beta_0 = 0.$$

Compute its t-value

$$t = b_0/s_{b_0},$$

and find whether this value is significant, again using t-distribution with $df = n - 2$.

Intervals for individual observations

Given the earlier sample of size n consider a new value x of the predictor variable. We wish to say something on the unobserved response value

$$Y = \beta_0 + \beta_1 x + \epsilon.$$

Its expected value

$$\mu = \beta_0 + \beta_1 x$$

is estimated by

$$\hat{\mu} = b_0 + b_1 x.$$

The standard error of $\hat{\mu}$ is computed as the square root of

$$\text{Var}(\hat{\mu}) = \frac{\sigma^2}{n} + \frac{\sigma^2}{n-1} \cdot \left(\frac{x-\bar{x}}{s_x}\right)^2.$$

An exact $100(1 - \alpha)\%$ confidence interval

$$I_\mu = b_0 + b_1 x \pm t_{n-2}\left(\frac{\alpha}{2}\right) \cdot s \sqrt{\frac{1}{n} + \frac{1}{n-1} \left(\frac{x-\bar{x}}{s_x}\right)^2}$$

Question. In what sense $\hat{\mu}$ is a random variable? Is it independent of ϵ that defines the random variable Y ?

Prediction interval

This I_μ should be compared to the prediction interval for Y

$$I = b_0 + b_1x \pm t_{n-2}\left(\frac{\alpha}{2}\right) \cdot s \sqrt{1 + \frac{1}{n} + \frac{1}{n-1} \left(\frac{x-\bar{x}}{s_x}\right)^2}$$

obtained from

$$\text{Var}(Y - \hat{\mu}) = \text{Var}(\mu + \epsilon - \hat{\mu}) = \sigma^2 + \text{Var}(\hat{\mu}) = \sigma^2 \left(1 + \frac{1}{n} + \frac{1}{n-1} \cdot \left(\frac{x-\bar{x}}{s_x}\right)^2\right).$$

Prediction interval I

has wider limits than I_μ ,
since it contains uncertainty
due to the noise component ϵ

The further x lies from \bar{x} ,
the wider are the intervals.

