

R01 - Simple linear regression

STAT 401 (Engineering) - Iowa State University

March 29, 2018

Telomere length

<http://www.pnas.org/content/101/49/17312>

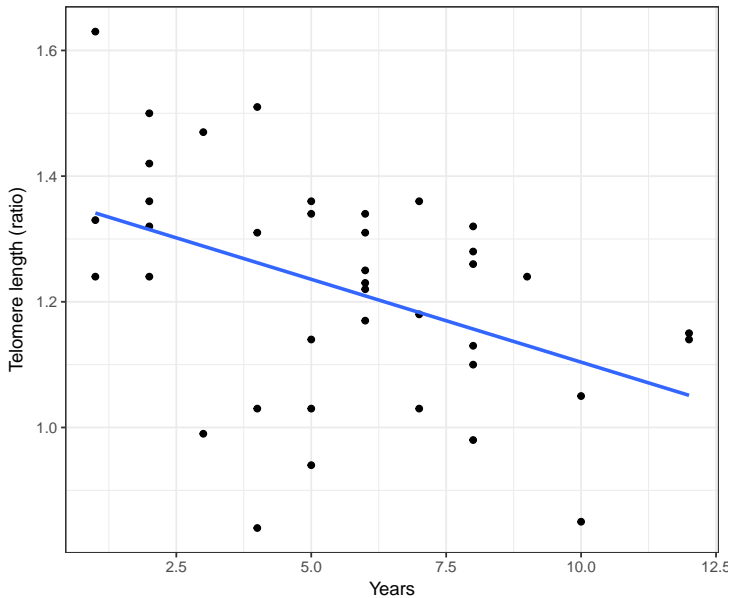
People who are stressed over long periods tend to look haggard, and it is commonly thought that psychological stress leads to premature aging and the earlier onset of diseases of aging.

...

This design allowed us to examine the importance of perceived stress and measures of objective stress (caregiving status and chronicity of caregiving stress based on the number of years since a child's diagnosis).

...

Telomere length values were measured from DNA by a quantitative PCR assay that determines the relative ratio of telomere repeat copy number to single-copy gene copy number (T/S ratio) in experimental samples as compared with a reference DNA sample.



Simple Linear Regression

The **simple linear regression** model is

$$Y_i \stackrel{ind}{\sim} N(\beta_0 + \beta_1 X_i, \sigma^2)$$

where Y_i and X_i are the response and explanatory variable, respectively, for individual i .

Terminology (all of these are equivalent):

response	explanatory
outcome	covariate
dependent	independent
endogenous	exogenous

Parameter interpretation

Recall:

$$E[Y_i|X_i = x] = \beta_0 + \beta_1 x \quad \text{Var}[Y_i|X_i = x] = \sigma^2$$

- If $X_i = 0$, then $E[Y_i|X_i = 0] = \beta_0$.
 β_0 is the **expected** response when the explanatory variable is zero.

- If X_i increases from x to $x + 1$, then

$$\begin{array}{rcl} E[Y_i|X_i = x + 1] & = & \beta_0 + \beta_1 x + \beta_1 \\ - E[Y_i|X_i = x] & = & \beta_0 + \beta_1 x \\ \hline & = & \beta_1 \end{array}$$

β_1 is the **expected** increase in the response for each unit increase in the explanatory variable.

- σ is the standard deviation of the response for a fixed value of the explanatory variable.

Remove the mean:

$$Y_i = \beta_0 + \beta_1 X_i + e_i \quad e_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

So the error is

$$e_i = Y_i - (\beta_0 + \beta_1 X_i)$$

which we approximate by the **residual**

$$r_i = \hat{e}_i = Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_i)$$

The least squares (minimize $\sum_{i=1}^n r_i^2$), maximum likelihood, and Bayesian estimators are

$$\begin{aligned} \hat{\beta}_1 &= SXY/SXX \\ \hat{\beta}_0 &= \bar{Y} - \hat{\beta}_1 \bar{X} \\ \hat{\sigma}^2 &= SSE/(n-2) \quad \text{df} = n-2 \end{aligned}$$

$$\begin{aligned} \bar{X} &= \frac{1}{n} \sum_{i=1}^n X_i \\ \bar{Y} &= \frac{1}{n} \sum_{i=1}^n Y_i \end{aligned}$$

$$\begin{aligned} SXY &= \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) \\ SXX &= \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X}) = \sum_{i=1}^n (X_i - \bar{X})^2 \\ SSE &= \sum_{i=1}^n r_i^2 \end{aligned}$$

How certain are we about $\hat{\beta}_0$ and $\hat{\beta}_1$ being equal to β_0 and β_1 ?

We quantify this uncertainty using their standard errors and posterior standard deviations:

$$SE(\beta_0) = \hat{\sigma} \sqrt{\frac{1}{n} + \frac{\bar{X}^2}{(n-1)s_X^2}} \quad df = n - 2$$

$$SE(\beta_1) = \hat{\sigma} \sqrt{\frac{1}{(n-1)s_X^2}} \quad df = n - 2$$

$$s_X^2 = SXX/(n-1)$$

$$s_Y^2 = SY/(n-1)$$

$$SYY = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

$$r_{XY} = \frac{SXY/(n-1)}{s_X s_Y}$$

$$R^2 = r_{XY}^2$$

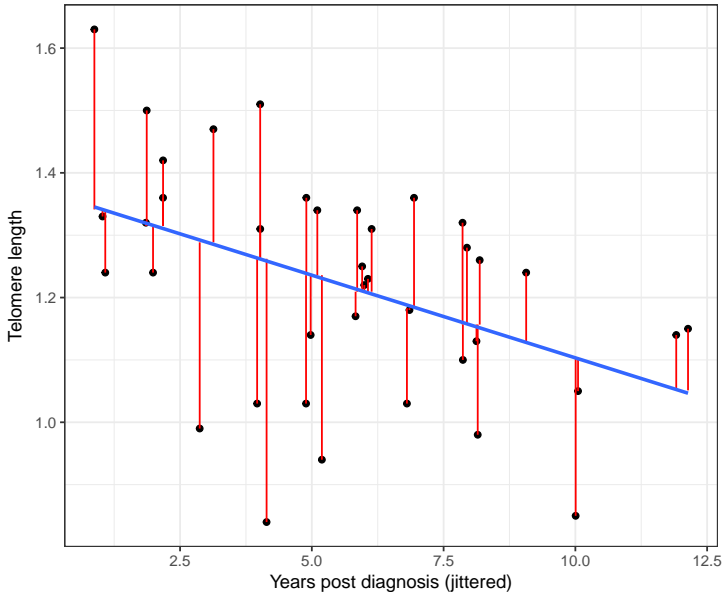
$$SST = SYY = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

$$= \frac{SST - SSE}{SST}$$

correlation coefficient

coefficient of determination

The coefficient of determination (R^2) is the proportion of the total response variation explained by the explanatory variable(s).



Pvalues and confidence interval

We can compute two-sided pvalues, e.g. $H_0 : \beta_0 = 0$ and $H_0 : \beta_1 = 0$, via

$$2P\left(T_{n-2} < -\left|\frac{\hat{\beta}_0}{SE(\beta_0)}\right|\right) \quad \text{and} \quad 2P\left(T_{n-2} < -\left|\frac{\hat{\beta}_1}{SE(\beta_1)}\right|\right)$$

These test the null hypothesis that the corresponding parameter is zero.

We can construct $100(1 - \alpha)\%$ two-sided confidence/credible intervals via

$$\hat{\beta}_0 \pm t_{n-2, 1-\alpha/2} SE(\beta_0) \quad \text{and} \quad \hat{\beta}_1 \pm t_{n-2, 1-\alpha/2} SE(\beta_1)$$

Calculations “by hand” in R

```
sm <- Telomeres %>%
  summarize(n      = n(),
            Xbar    = mean(years),
            Ybar    = mean(telomere.length),
            s_X     = sd(years),
            s_Y     = sd(telomere.length),
            r_XY    = cor(telomere.length, years))

sm
```

	n	Xbar	Ybar	s_X	s_Y	r_XY
1	39	5.589744	1.220256	2.935427	0.1797731	-0.4306534

Calculations “by hand” in R (continued)

$$\begin{aligned}
 SXX &= (n-1)s_x^2 = (39-1) \times 2.9354274^2 = 327.4358974 \\
 SY Y &= (n-1)s_y^2 = (39-1) \times 0.1797731^2 = 1.2280974 \\
 SXY &= (n-1)s_x s_y r_{XY} = (39-1) \times 2.9354274 \times 0.1797731 \times -0.4306534 = -8.6358974 \\
 \hat{\beta}_1 &= SXY/SXX = -8.6358974/327.4358974 = -0.0263743 \\
 \hat{\beta}_0 &= \bar{Y} - \hat{\beta}_1 \bar{X} = 1.2202564 - (-0.0263743) \times 5.5897436 = 1.3676821 \\
 R^2 &= r_{XY}^2 = (-0.4306534)^2 = 0.1854624 \\
 SSE &= SY Y (1 - R^2) = 1.2280974(1 - 0.1854624) = 1.0003316 \\
 \hat{\sigma}^2 &= SSE/(n-2) = 1.0003316/(39-2) = 0.027036 \\
 \hat{\sigma} &= \sqrt{\hat{\sigma}^2} = \sqrt{0.027036} = 0.1644262 \\
 SE(\hat{\beta}_0) &= \hat{\sigma} \sqrt{\frac{1}{n} + \frac{\bar{X}^2}{(n-1)s_x^2}} = 0.1644262 \sqrt{\frac{1}{39} + \frac{5.5897436^2}{(39-1) \times 2.9354274^2}} = 0.0572111 \\
 SE(\hat{\beta}_1) &= \hat{\sigma} \sqrt{\frac{1}{(n-1)s_x^2}} = 0.1644262 \sqrt{\frac{1}{(39-1) \times 2.9354274^2}} = 0.0090867 \\
 p_{H_0: \beta_0=0} &= 2P\left(t_{n-2} < -\left|\frac{\hat{\beta}_0}{SE(\hat{\beta}_0)}\right|\right) = 2P(t_{37} < -23.9058799) = 4.2740348 \times 10^{-24} \\
 p_{H_0: \beta_1=0} &= 2P\left(t_{n-2} < -\left|\frac{\hat{\beta}_1}{SE(\hat{\beta}_1)}\right|\right) = 2P(t_{37} < -2.9025065) = 0.0062047 \\
 CI_{95\% \beta_0} &= \hat{\beta}_0 \pm t_{n-2, 1-\alpha/2} SE(\hat{\beta}_0) \\
 &= 1.3676821 \pm 2.0261925 \times 0.0572111 = (1.2517613, 1.4836028) \\
 CI_{95\% \beta_1} &= \hat{\beta}_1 \pm t_{n-2, 1-\alpha/2} SE(\hat{\beta}_1) \\
 &= -0.0263743 \pm 2.0261925 \times 0.0090867 = (-0.0447858, -0.0079628)
 \end{aligned}$$

Regression in R

```
m = lm(telomere.length~years, Telomeres)
summary(m)
```

```
Call:
lm(formula = telomere.length ~ years, data = Telomeres)
```

```
Residuals:
    Min       1Q   Median       3Q      Max
-0.42218 -0.08537  0.02056  0.10738  0.28869
```

```
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  1.367682   0.057211   23.906  <2e-16 ***
years        -0.026374   0.009087   -2.903   0.0062 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 0.1644 on 37 degrees of freedom
Multiple R-squared:  0.1855, Adjusted R-squared:  0.1634
F-statistic: 8.425 on 1 and 37 DF,  p-value: 0.006205
```

```
confint(m)
```

```
                2.5 %      97.5 %
(Intercept)  1.25176134  1.483602799
years        -0.04478579 -0.007962836
```

Conclusion

Telomere length at the time of diagnosis of a child's chronic illness is estimated to be 1.37 with a 95% confidence interval of (1.25, 1.48). For each year since diagnosis, the telomere length decreases by 0.026 with a 95% confidence interval of (0.008, 0.045) **on average**. The proportional of variability in telomere length described by years since diagnosis is 18.5%.

<http://www.pnas.org/content/101/49/17312>

The zero-order correlation between chronicity of caregiving [years] and mean telomere length, r , is -0.445 ($P < 0.01$). [$R^2 = 0.198$ was shown in the plot.]

Remark I'm guessing our analysis and that reported in the paper don't match exactly due to a discrepancy in the data.

Summary

- The **simple linear regression** model is

$$Y_i \overset{ind}{\sim} N(\beta_0 + \beta_1 X_i, \sigma^2)$$

where Y_i and X_i are the response and explanatory variable, respectively, for individual i .

- Know how to use R to obtain $\hat{\beta}_0$, $\hat{\beta}_1$, $\hat{\sigma}^2$, R^2 , pvalues, CIs, etc.
- Interpret R output
 - At a value of zero for the explanatory variable ($X_i = 0$), β_0 is the expected value for the response (Y_i).
 - For each unit increase in the explanatory variable value, β_1 is the expected increase in the response.
 - At a constant value of the explanatory variable, σ^2 is the variance of the responses.
 - The coefficient of determination (R^2) is the percentage of the total response variation explained by the explanatory variable(s).

What is $E[Y|X = x]$?

We know $\beta_0 = E[Y|X = 0]$, but what about $X = x$?

$$E[Y|X = x] = \beta_0 + \beta_1 x$$

which we can estimate via

$$E[\widehat{Y|X = x}] = \hat{\beta}_0 + \hat{\beta}_1 x$$

but there is uncertainty in both β_0 and β_1 . So the standard error of $E[Y|X = x]$ is

$$SE(E[Y|X = x]) = \hat{\sigma} \sqrt{\frac{1}{n} + \frac{(\bar{X} - x)^2}{(n-1)s_X^2}}$$

and a $100(1 - \alpha)\%$ confidence interval is

$$\hat{\beta}_0 + \hat{\beta}_1 x \pm t_{n-2, 1-\alpha/2} SE(E[Y|X = x])$$

What do we predict about Y at $X = x$?

On the last slide, we calculated $E[Y|X = x]$ and it's uncertainty, but if we are trying to predict a new observation, we need to account for the sampling variability σ^2 . Thus a prediction about Y at a new $X = x$ is still

$$Pred\{Y|X = x\} = \hat{\beta}_0 + \hat{\beta}_1 x$$

but the uncertainty includes the variability due to σ^2 . So the standard error of $Pred\{Y|X = x\}$ is

$$SE(Pred\{Y|X = x\}) = \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(\bar{X} - x)^2}{(n-1)s_X^2}}$$

and a $100(1 - \alpha)\%$ prediction interval is

$$\hat{\beta}_0 + \hat{\beta}_1 x \pm t_{n-2, 1-\alpha/2} SE(Pred\{Y|X = x\}).$$

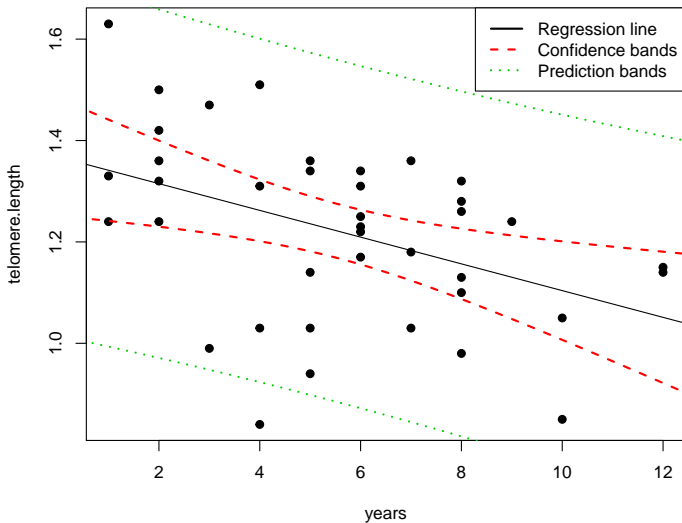
Confidence and prediction intervals for different values of X

```
m = lm(telomere.length~years, Telomeres)
new = data.frame(years=4:6)
new %>% bind_cols(predict(m, new, interval="confidence") %>% as.data.frame)
```

	years	fit	lwr	upr
1	4	1.262185	1.201335	1.323035
2	5	1.235810	1.181369	1.290252
3	6	1.209436	1.155556	1.263316

```
new %>% bind_cols(predict(m, new, interval="prediction") %>% as.data.frame)
```

	years	fit	lwr	upr
1	4	1.262185	0.9235142	1.600855
2	5	1.235810	0.8982324	1.573389
3	6	1.209436	0.8719482	1.546924



Shifting the intercept

The intercept (β_0) is the expected response when the explanatory variable is zero.

So, if we change our explanatory variable, we change the interpretation of our intercept, e.g. if, instead of using number of years since diagnosis, we use “number of years since diagnosis **minus 4**”, then our intercept is the expected response at 4 years since diagnosis.

Let x be number of years since diagnosis, then

$$E[Y|X = x] = \tilde{\beta}_0 + \tilde{\beta}_1(x - 4) = (\beta_0 - 4\beta_1) + \beta_1 x$$

so our new parameters for the mean are

- intercept $\tilde{\beta}_0 = (\beta_0 - 4\beta_1)$ and
- slope $\tilde{\beta}_1 = \beta_1$ (unchanged).

Shifting the intercept (continued)

```
m0 = lm(telomere.length ~ years, Telomeres)
m4 = lm(telomere.length ~ I(years-4), Telomeres)
```

```
coef(m0)
```

```
(Intercept)      years
 1.36768207 -0.02637431
```

```
coef(m4)
```

```
(Intercept) I(years - 4)
 1.26218481 -0.02637431
```

```
confint(m0)
```

```
              2.5 %      97.5 %
(Intercept) 1.25176134 1.483602799
years       -0.04478579 -0.007962836
```

```
confint(m4)
```

```
              2.5 %      97.5 %
(Intercept) 1.20133473 1.323034890
I(years - 4) -0.04478579 -0.007962836
```