# Advanced Statistics DS2003 (BDS-4A) Lecture 20

Instructor: Dr. Syed Mohammad Irteza

Assistant Professor, Department of Computer Science, FAST

28 April, 2022

# Previous Lecture

- Use of Python for linear regression analysis:
  - Example with real sample data
  - Example with randomly generated values

- Multiple Linear Regression

# Today

- Revision of Linear Regression
  - Sum of squares for x and y
  - Sum of products for x, y

- Least squares regression model
  - Residuals sum to zero, and the line always passes through $(\bar{x}, \bar{y})$

| | weight (g) | volume (cm$^3$) | cover |
|---|---|---|---|
| 1 | 800 | 885 | hc |
| 2 | 950 | 1016 | hc |
| 3 | 1050 | 1125 | hc |
| 4 | 350 | 239 | hc |
| 5 | 750 | 701 | hc |
| 6 | 600 | 641 | hc |
| 7 | 1075 | 1228 | hc |
| 8 | 250 | 412 | pb |
| 9 | 700 | 953 | pb |
| 10 | 650 | 929 | pb |
| 11 | 975 | 1492 | pb |
| 12 | 350 | 419 | pb |
| 13 | 950 | 1010 | pb |
| 14 | 425 | 595 | pb |
| 15 | 725 | 1034 | pb |



3

*somewhat abbreviated output...*

```
Coefficients:

              Estimate    Std. Error t value   Pr(>|t|)
(Intercept) 107.67931       88.37758   1.218      0.245
Volume        0.70864        0.09746   7.271   6.26e-06


Residual standard error: 123.9 on 13 degrees of freedom
Multiple R-squared:  0.8026,Adjusted  R-squared: 0.7875
F-statistic: 52.87 on 1 and 13  DF,  p-value: 6.262e-06
```

# Modeling weights of books using volume <u>and </u>cover type

Coefficients:

|  | Estimate | Std. Error | t value | Pr(>\|t\|) | |
|---|---|---|---|---|---|
| (Intercept) | 197.96284 | 59.19274 | 3.344 | 0.005841 | ** |
| volume | 0.71795 | 0.06153 | 11.669 | 6.6e-08 | *** |
| cover:pb | -184.04727 | 40.49420 | -4.545 | 0.000672 | *** |

Residual standard error: 78.2 on 12 degrees of freedom
Multiple R-squared: 0.9275,Adjusted R-squared: 0.9154  F-statistic: 76.73  on  2  and   12  DF,    p-value: 1.455e-07

$\hat{\beta}_0$ & $\hat{\beta}_1$ are chosen to minimize

$$\sum e_i^2 = \sum (Y_i - \hat{Y}_i)^2 = \sum (Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_i))^2 \quad \} \text{ sum of squared residuals ..}$$

This is called the method of least squares.

$SS_{xx} = \text{"sum of squares for X"} = \sum (x_i - \bar{x})^2$

$s_x^2 = \dfrac{SS_{xx}}{n-1} \to \text{variance, sample}$

$SS_{yy} = \text{"sum of squares for y"} = \sum (y_i - \bar{y})^2$

$s_y^2 = \dfrac{SS_{yy}}{n-1} \to \text{variance of } y, \text{ sample}$

$SP_{xy} = \text{"sum of products"} = \sum (x_i - \bar{x})(y_i - \bar{y})$
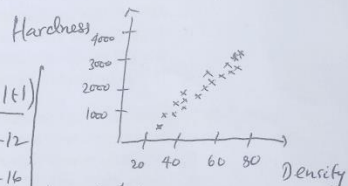
sample covariance $\to Cov(x,y) = \dfrac{SP_{xy}}{n-1}$

$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$

$\hat{\beta}_1 = \dfrac{SP_{xy}}{SS_{xx}} = \dfrac{Cov(x,y)}{Var(x)}$

Hardness vs Density for 36 Australian Tree Species:



| Coefficients | Estimate | Std. Error | t value | Pr(>|t|) |
|---|---|---|---|---|
| (intercept) | -1160.50 | 108.58 | -10.69 | 2.07e-12 |
| Density | 57.507 | 2.279 | 25.24 | <2e-16 |

use "Density" to predict "Hardness"

$\hat{y} = \underset{\hat{\beta}_0}{-1160.50} + \underset{\hat{\beta}_1}{57.507 X}$

For least squares regression :
→ The residuals sum to 0 $(\sum e_i = 0)$.
→ The line passes through $(\bar{X}, \bar{Y})$

$SE_{b_0} = \dfrac{s}{\sqrt{n}} * \sqrt{1 + \dfrac{(\bar{X})^2}{Var(x)}}$

$SE_{b_1} = \dfrac{s}{\sqrt{n}} * \dfrac{1}{stdev(x)}$

$s (\text{st. error of the regression}) = \sqrt{\dfrac{1}{n-2} \sum_{t=1}^{n} e_t^2} = stdev(errors) \times \sqrt{\dfrac{(n-1)}{(n-2)}}$

$R^2 = 1 - \dfrac{Var(errors)}{Var(Y)}$

---

A study investigated a possible relationship btw eggshell thickness & environmental contaminants in brown pelican eggs.

pesticide, banned

The figure shows a scatterplot of shell thickness vs. DDT level in a sample of 65 brown pelican eggs on Anacapa Island, California.

We use a computer to fit the least-squares regression line.

| | Estimate | Std. Error | t value | Pr (>|t|) |
|---|---|---|---|---|
| Intercept | 4.231e-01 | 2.128e-02 | 19.880 | <2e-16 |
| DDT (ppm) | -8.732e-05 | 1.603e-05 | -5.448 | 9e-07 |

$\hat{Y} = 0.4231 - 0.00008732 x$

egg with DDT = 2000 → prediction of thickness?

$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X = 0.4231 - (0.00008732 * 2000)$
$= 0.24846 \text{ mm}$

What proportion of the variability in eggshell thickness can be explained by the linear relationship w/ DDT?



2 reasons → → decreasing trend
→ random variability about the line..

proportion due to the "decreasing trend" is $r^2$ ($r = $ correlation coefficient)

$r^2$ → proportion of the variability in response variable Y that is attributable to the linear relationship w/ X.

"coefficient of determination"

Multiple R-squared = 0.3202 ⟹ $r^2 = 0.3202$.

"32% of the variability in eggshell thickness can be explained by the linear relationship with DDT"

⟹ Thus, the other 68% of variability ⟹ random variation...

Before doing any statistical inference,
we should check the residual plots.

Simple reg. model

$(\text{Plots of } e_i = Y_i - \hat{Y}_i)$

→ errors follow normal
distribution, ... etc.

$\underset{obs.}{Y} \quad \underset{pred.}{Y}$

→ Constant variance

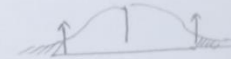Plot looks fine..
Normal? → plot theoretical quantiles

Test the $H_0$ : no linear relationship btw Y & X.

$$H_0: \beta_1 = 0$$

$$t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)} = \frac{-8.732 \times 10^{-5} - 0}{1.603 \times 10^{-5}}$$

But, we don't conclude "causation"!

$$= -5.448$$

p-value = 9e-07 = 0.0000009
"strong evidence against $H_0$".

Weights of books (g)

| | | volume (cm³) |
|---|---|---|
| 1 | 800 | |
| 2 | 950 | 885 |
| 3 | 1050 | 1016 |
| 4 | 250 | 1125 |
| 5 | 750 | 239 |
| 6 | 600 | 701 |
| 7 | 1075 | 641 |
| 8 | 250 | 1228 |
| 9 | 700 | 412 |
| 10 | 650 | 953 |
| 11 | 975 | 929 |
| 12 | 350 | 1492 |
| 13 | 950 | 419 |
| 14 | 425 | 1010 |
| 15 | 725 | 595 |
| | | 1034 |

Coefficients

| | Est | SE | t-val |
|---|---|---|---|
| (Intercept) | 107.679 | 88.37 | 1.218 |
| Volume | 0.70864 | 0.09746 | 7.271 |

$R^2 = 0.8026$

# Jbstatistics (Youtube)

- **Simple Linear Regression: The Least Squares Regression Line**
  - https://www.youtube.com/watch?v=coQAAN4eY5s

- **Simple Linear Regression: An Example**
  - https://www.youtube.com/watch?v=xIDjj6ZyFuw

# Sources

- [openintro.org/os](http://openintro.org/os) (Chapter 9, Section 9.1)

Linear Regression using Excel:

[https://1drv.ms/x/s!Apc0G8okxWJ1zCUKXCGBs8TgfywO?e=I69d5e](https://1drv.ms/x/s!Apc0G8okxWJ1zCUKXCGBs8TgfywO?e=I69d5e)