| Summer/WinterSemester2 | **Tutorial 4** | 2015 |
|---|---|---|

> This tutorial explores bivariate data.
> Most of this week's tutorial class will consist of the Report 1 presentations.
> Complete the rest of the tutorial questions at home.

---

**Bivariate Data**

For paired observations $\{(x_i, y_i)\}$ for $i = 1, \ldots, n$

summary statistics

$$S_{xy} = \sum_{i=1}^{n} x_i y_i - \frac{1}{n}\left(\sum_{i=1}^{n} x_i\right)\left(\sum_{i=1}^{n} y_i\right)$$

$$S_{xx} = \sum_{i=1}^{n} x_i^2 - \frac{1}{n}\left(\sum_{i=1}^{n} x_i\right)^2 = (n-1)s_x^2$$

$$S_{yy} = \sum_{i=1}^{n} y_i^2 - \frac{1}{n}\left(\sum_{i=1}^{n} y_i\right)^2$$

least squares regression line $\quad y = a + bx$

where $a = \bar{y} - b\bar{x}$ and $b = \dfrac{S_{xy}}{S_{xx}}$

correlation coefficient $\quad r = \dfrac{S_{xy}}{\sqrt{S_{xx} S_{yy}}} = b\dfrac{s_x}{s_y}$

---

**1.** Bivariate Data by hand

Of environmental interest is the relation between carbon monoxide concentration and traffic density. The following table gives the traffic density (vehicles per hour to the nearest 500 vehicles) and carbon monoxide concentration (CO) in ppm for a particular street corner in Newtown.

Note the table reads as $(x_1, y_1) = (1.0, 9), (x_2, y_2) = (1.0, 6.8), \ldots (x_{12}, y_{12}) = (3.0, 20.6)$

| $x$: Traffic density (in thousands) | $y$: CO concentration (in ppm) |
|---|---|
| 1.0 | 9.0   6.8   7.7 |
| 1.5 | 9.6   6.8   10.3 |
| 2.0 | 12.3   11.8 |
| 3.0 | 20.7   20.2   21.6   20.6 |

(a) Produce a scatter plot. Does a linear fit seem reasonable?

(b) Calculate the summary statistics from

$$\sum_{i=1}^{12} x_i y_i = 361.05, \sum_{i=1}^{12} x_i^2 = 53.75, \sum_{i=1}^{12} y_i^2 = 2449, \sum_{i=1}^{12} x_i = 23.5, \sum_{i=1}^{12} y_i = 157.4$$

(c) Calculate the mean and variance of $x$ and $y$.

(d) Find the least squares regression line and correlation coefficient. What does $r$ suggest?

(e) What is the predicted CO concentration on a day with traffic density 2,500?

(f) For the reading (2.0,12.3), what is the residual?

(g) On a day with traffic density of 5,000, would it be appropriate to use the least squares regression line to predict the CO concentration? Explain.

## 2. Bivariate Data in R

Using the data in Question 1, produce a scatterplot, least squares regression line, correlation coefficient and residual plot. Is the linear fit appropriate?

```
> x=c(1,1,1,1.5,1.5,1.5,2,2,3,3,3,3)
> y=c(9.0,6.8,7.7,9.6,6.8,10.3,12.3,11.8,20.7,20.2,21.6,20.6)
> cor(x,y)
> fit <- lm(y~x)
> a= fit$coefficients[[1]]
> b= fit$coefficients[[2]]
> plot(x,y,xlab="Traffic density", ylab="CO")
> abline(fit)

> res = y - (a + b * x)
> plot(x, res)                OR > plot(x,fit$residuals)
> abline(h = 0)                   > abline(h = 0)
```

## 3. Bivariate Data

J.B. Haldane is responsible for showing how carbon dioxide levels in the blood influences breathing rates by affecting the acidity of the blood. In one experiment he administered varying doses of sodium bicarbonate with the following results:

| Dose (in grams) | $x$ | 30 | 40 | 50 | 60 | 70 | 80 | 90 | 100 |
|---|---|---|---|---|---|---|---|---|---|
| Breathing rate | $y$ | 16 | 14 | 13 | 13 | 11 | 12 | 9 | 9 |

(a) By hand, find the least squares regression line and correlation coefficient. What breathing rate would you predict for a dose of 85g?

(b) Using R, produce a scatter plot and residual plot. Is the linear fit appropriate?

(c) Refit the model, so that you could predict the dose from the breathing rate.

**4.** Bivariate Data in R

The following data describes the relationship between obesity measured as the percentage over ideal weight ($x$) and individual's response to pain measured on a certain scale ($y$):

| $x$ | 89% | 90% | 75% | 30% | 51% |
|---|---|---|---|---|---|
| $y$ | 2 | 3 | 4 | 4.5 | 5.5 |

(a) Using R, produce a scatter plot, least squares regression line, correlation coefficient and residual plot.

(b) Do the plots suggest that the linear fit is appropriate? If so, use the least squares line to predict the $y$-value for an $x$-value of 60%.