

Linear Models

①

Recall $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$

pop. intercept pop. slope pred. variable error

$$y_1 = \beta_0 + \beta_1 x_1 + \varepsilon_1$$

$$y_2 = \beta_0 + \beta_1 x_2 + \varepsilon_2$$

$$y_3 = \beta_0 + \beta_1 x_3 + \varepsilon_3$$

...

$\varepsilon_1, \varepsilon_2, \varepsilon_3, \dots \equiv$ residuals

Deviation of y_i from expected \hat{y}_i

Recall $\varepsilon \sim N(0, \sigma^2)$

Then $Y \sim \beta_0 + \beta_1 X + \varepsilon$, which can be written as: $Y \sim N(\beta_0 + \beta_1 X, \sigma^2)$

Aside show that ~~$N(10, \sigma)$~~ $N(10, \sigma) = 10 + N(0, \sigma)$

We now use OLS or ML to estimate β_0 and β_1 .

$$\beta_1 : b_1 = \frac{\sum_{i=1}^n [(x_i - \bar{x})(y_i - \bar{y})]}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\beta_0 : b_0 = \bar{y} - b_1 \bar{x}$$

$$\varepsilon_i : e_i = y_i - \hat{y}_i$$

~~What if instead of continuous x_i , we had categorical values?~~

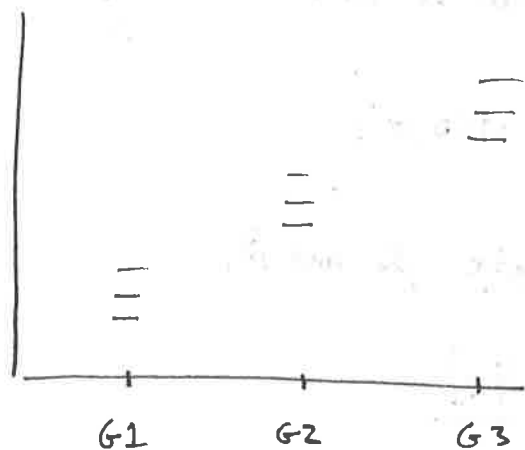
What if we have more than one predictor variable?

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \varepsilon_i$$

What if instead of continuous X_i , we had categorical variables?

- Instead of β_0 , use overall mean μ .
- Use dummy variables, one for each category / group.
- i treatment levels ($1 + p$)
- j replicates within each treatment ($1 + n$)

$$Y_{ij} = \mu + \beta_1 (\text{dummy}_1)_{ij} + \beta_2 (\text{dummy}_2)_{ij} + \dots + \epsilon_{ij}$$



9 observations

$$(\text{dummy}_1) =$$

1
1
0
0
0
0
0
0

~~xxxx~~

$$(\text{dummy}_2) =$$

0
0
0
1
1
0
0
0

More typically this model expressed as an effects model in which treatment levels are represented by a single term that ^{denotes the} effect of each level on the overall mean

$$Y_{ij} = \mu + \alpha_i + \epsilon_{ij}$$

~~xxxx~~ A problem with this is that if you have e.g. three groups, w/ three means, and the overall mean, you have 4 params and 3 groups of data

Correlation

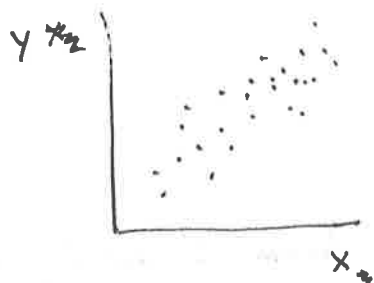
3

Let's take a step back and ignore the potential cause-effect relationship implied by linear models (i.e., change in x , results in a change in y). Perhaps their change is correlated, or we may say that \underline{x} and y covary.

We measure the strength of this relationship as

covariation
or
correlation

How do we visualize co-variation?



How do we calculate co-variation?

Recall variance in x :
$$\frac{\sum (x_i - \bar{x})^2}{n-1}$$

Similarly, variance in y :
$$\frac{\sum (y_i - \bar{y})^2}{n-1}$$

Co-variance then is :
$$\frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

A.K.A.
Sum of cross-products

Go to example

Co-variance (much like variance) is dependent on magnitude of two variables of interest.

[multiply sepal length by 10 and see the effect]

→ We can standardize co-variance by dividing by the standard deviations of both variables, yielding the Pearson's correlation coefficient (product-moment coeff)

$$\rho_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

When using $\hat{\bar{x}}, \hat{\bar{y}}$; $\rho_{xy} = r_{xy}$

$$\rho_{xy} \in [-1, 1]$$

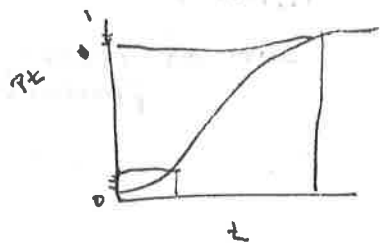
Also, when using samples from X and Y, then

r is the sample correlation-coeff statistic.

It has a dist'n, just like all other sample statistics. We can write its

standard error: $\sqrt{\frac{1-r^2}{(n-2)}}$

And then we can test whether r is significantly different from 0 or not, using a t test: $t = \frac{r}{s_r}$



Linear Regression

5

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

Null Hypotheses:

$$\beta_0 = 0 \quad (\text{the population } y\text{-intercept equals zero})$$

$$\beta_1 = 0 \quad (\text{the population slope equals zero})$$

We are rarely interested in the first null hypothesis.

Both of these are tested using t-statistics.

<u>Full model</u>	$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$	} Q: Which one will likely have greater error (variation)?
<u>Reduced model</u>	$Y_i = \beta_0 + \varepsilon_i$	

The difference in the variation explained by these two models is known as: Explained Variation

Does the full model explain significantly more variation?

For this, we use the F-test.

What's going on here visually?

Go to p. 172 in Logan.