

Regression analyses

Martin Bland

Professor of Health Statistics

University of York

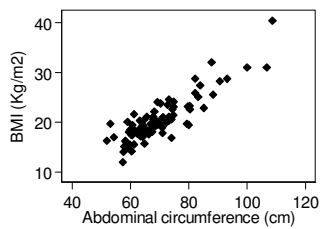
<http://www-users.york.ac.uk/~mb55/msc/>

Regression analyses

- Simple linear regression
- Multiple linear regression
- Curvilinear regression
- Dichotomous predictor variables
- Regression in clinical trials
- Dichotomous outcome variables and logistic regression
- Interactions
- Factors with more than two levels
- Sample size

Simple Linear Regression

Example: Body Mass Index (BMI) and abdominal circumference in 86 women



(Data of Malcom Savage)

What is the relationship?

Regression: predict BMI from observed abdominal circumference.

Simple Linear Regression

Example: Body Mass Index (BMI) and abdominal circumference in 86 women.

What is the relationship?

Regression: predict BMI from observed abdominal circumference.

What is the mean BMI for women with any given observed abdominal circumference?

Simple Linear Regression

Example: Body Mass Index (BMI) and abdominal circumference in 86 women.

What is the relationship?

Regression: predict BMI from observed abdominal circumference.

What is the mean BMI for women with any given observed abdominal circumference?

BMI is the **outcome, dependent, y, or left hand side** variable.

Abdominal circumference is the **predictor, explanatory, independent, x, or right hand side** variable.

Simple Linear Regression

Example: Body Mass Index (BMI) and abdominal circumference in 86 women.

What is the relationship?

Regression: predict BMI from observed abdominal circumference.

What is the mean BMI for women with any given observed abdominal circumference (AC)?

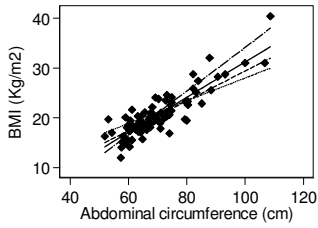
Linear relationship:

$$\text{BMI} = \text{intercept} + \text{slope} \times \text{AC}$$

Equation of a straight line.

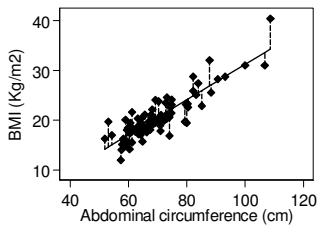
Simple Linear Regression

Which straight line should we choose?



Simple Linear Regression

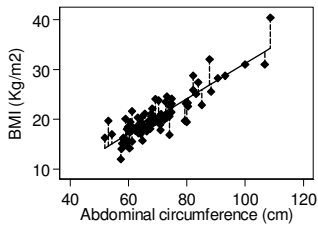
Which straight line should we choose?



Choose the line which makes the distance from the points to the line **in the y direction** a minimum.
Differences between the observed strength and the predicted strength.

Simple Linear Regression

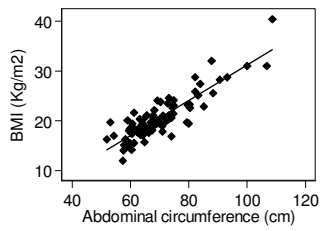
Which straight line should we choose?



Minimise the sum of the squares of these differences.
Principle of least squares, least squares line or equation.

Simple Linear Regression

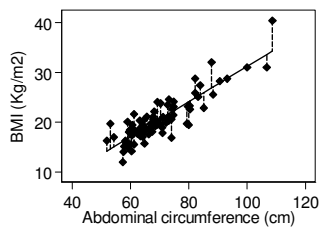
$$\text{BMI} = -4.15 + 0.35 \times \text{AC}$$



We can find confidence intervals and P values for the coefficients subject to assumptions.

Simple Linear Regression

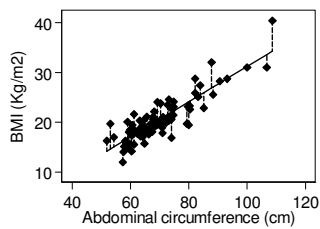
We can find confidence intervals and P values for the coefficients subject to assumptions.



Deviations from line should have a Normal distribution with uniform variance.

Simple Linear Regression

Can find confidence intervals and P values for the coefficients subject to assumptions.



Slope = 0.35 Kg/m²/cm, 95% CI = 0.31 to 0.40 Kg/m²/cm, P<0.001 against zero.

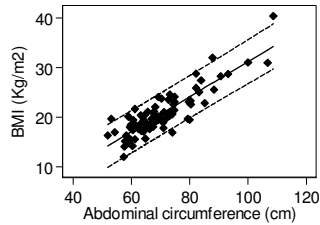
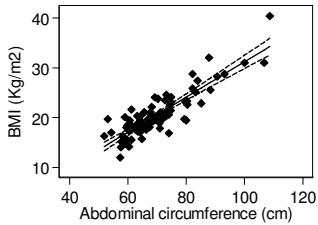
Intercept = -4.15 Kg/m², 95% CI = -7.11 to -1.18 Kg/m².

Simple Linear Regression

We can also find confidence intervals for regression estimates and predicted value for a new subject.

95% confidence intervals for regression estimates for BMI and abdominal circumference

Prediction intervals or 95% confidence intervals for prediction of BMI from abdominal circumference



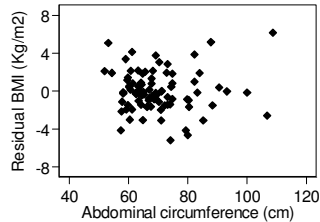
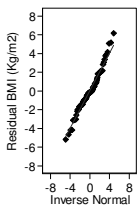
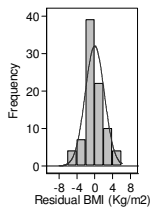
Simple Linear Regression

Assumptions: deviations from line should have a Normal distribution with uniform variance.

Calculate the deviations or residuals, observed minus predicted.

Check Normal distribution:

Check uniform variance:



Dichotomous predictor variable

24 hour energy expenditure (MJ) in two groups of women

Lean			Obese	
6.13	7.53	8.09	8.79	9.69
7.05	7.58	8.11	9.19	9.97
7.48	7.90	8.40	9.21	11.51
7.48	8.08	10.15	9.68	11.85
		10.88		12.79

Can carry out linear regression.

Define variable: obese = 0 if woman lean,
obese = 1 if woman obese.

Regression equation:

$$\text{energy} = 5.83 + 2.23 \times \text{obese}$$

slope: 95% CI = 1.05 to 3.42 MJ, P=0.0008.

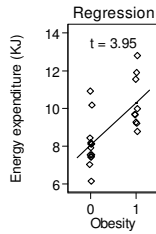
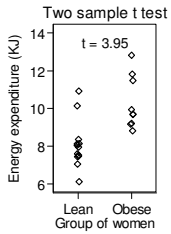
Regression and the two sample t method

Regression:

$$\text{energy} = 5.83 + 2.23 \times \text{obese}$$

slope: 95% CI = 1.05 to 3.42 MJ, P=0.0008.

The two methods are identical.



Difference (obese – lean) = 10.298 – 8.066 = 2.232.

Two sample t method:

95% CI = 1.05 to 3.42 MJ, P=0.0008.

Regression and the two sample t method

Assumptions of two sample t method

1. Energy expenditure follows a Normal distribution in each population.
2. Variances are the same in each population.

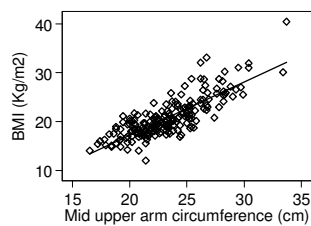
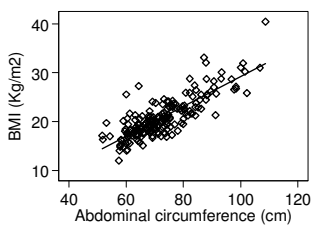
Assumptions of regression

1. Differences between observed and predicted energy expenditure follow a Normal distribution.
2. Variances of differences are the same in whatever the value of the predictor.

These are the same.

Multiple Linear Regression

More than one predictor:



$$\text{BMI} = -1.35 + 0.31 \times \text{AC}$$

$$\text{BMI} = -4.59 + 1.09 \times \text{MUAC}$$

$$\text{BMI} = -5.94 + 0.18 \times \text{AC} + 0.59 \times \text{MUAC}$$

Multiple Linear Regression

More than one predictor:

$$\text{BMI} = -1.35 + 0.31 \times \text{AC} \quad \text{BMI} = -4.59 + 1.09 \times \text{MUAC}$$
$$\text{BMI} = -5.94 + 0.18 \times \text{AC} + 0.59 \times \text{MUAC}$$

We find the coefficients which make the sum of the squared differences between the observed BMI and that predicted by the regression a minimum.

This is called **ordinary least squares** regression or **OLS** regression.

Multiple Linear Regression

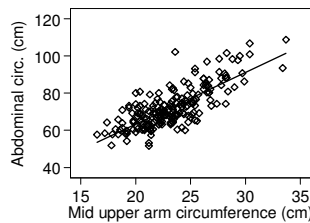
More than one predictor:

$$\text{BMI} = -1.35 + 0.31 \times \text{AC} \quad \text{BMI} = -4.59 + 1.09 \times \text{MUAC}$$
$$\text{BMI} = -5.94 + 0.18 \times \text{AC} + 0.59 \times \text{MUAC}$$

Both coefficients are pulled towards zero because abdominal circumference and arm circumference are related:

$$\text{MUAC} = 7.52 + 2.79 \times \text{AC},$$
$$r = 0.77, P < 0.001$$

AC and MUAC each explains some of the relationship between BMI and the other.



Multiple Linear Regression

More than one predictor:

We can find confidence intervals for the coefficients and test the null hypotheses that coefficients are zero in the population.

$\text{BMI} = -5.94$	$+$	$0.18 \times \text{AC}$	$+$	$0.59 \times \text{MUAC}$
95% CI -8.10 to -3.77		0.14 to 0.22		0.45 to 0.74
		P<0.001		P<0.001

Each predictor reduces the significance of the other because they are related to one another as well as to BMI.

They can both become not significant, even though the regression as a whole is highly significant.

Multiple Linear Regression

Assumptions:

Just as for simple linear regression, for our confidence intervals and P values to be valid, the data must conform to the assumptions that

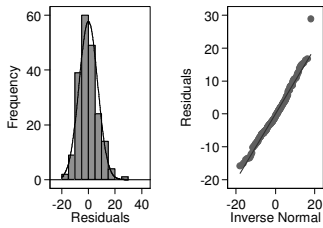
- deviations from line should have a Normal distribution,
- with uniform variance,
- observations must be independent.

Finally, our model of the data is that the relationship with each of our predictors is adequately represented by a straight line rather than a curve.

Multiple Linear Regression

Assumptions:

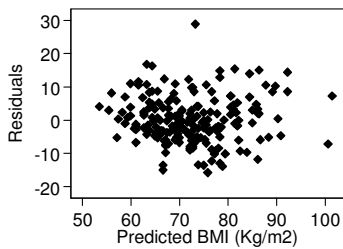
Check by histogram and Normal plot of residuals:



Multiple Linear Regression

Assumptions:

and by plot of residuals against regression estimate:



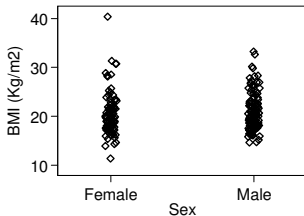
Multiple Linear Regression

Dichotomous predictor: sex.

Variable male = 0 for a female, = 1 for a male.

$$\text{BMI} = 20.51 + 0.40 \times \text{male}$$

95% CI 19.64 to 21.38 -0.75 to 1.55
 P = 0.5



Sex is not a significant predictor alone.

Multiple Linear Regression

Dichotomous predictor: sex.

Variable male = 0 for a female, = 1 for a male.

$$\text{BMI} = 20.51 + 0.40 \times \text{male}$$

95% CI 19.64 to 21.38 -0.75 to 1.55
 P = 0.5

$$\text{BMI} = -6.44 + 0.18 \times \text{AC} + 0.64 \times \text{MUAC} - 1.39 \times \text{male}$$

-8.49 to -4.39 0.14 to 0.22 0.50 to 0.78 -1.94 to -0.84
 P<0.001 P<0.001 P<0.001

Male has become a significant predictor because abdominal circumference and arm circumference have removed a lot of variability.

Mean BMI is lower for men than women **of the same abdominal and arm circumference** by 1.39 units.

Multiple Linear Regression

Dichotomous predictor: sex.

Variable male = 0 for a female, = 1 for a male.

$$\text{BMI} = -6.44 + 0.18 \times \text{AC} + 0.64 \times \text{MUAC} - 1.39 \times \text{male}$$

-8.49 to -4.39 0.14 to 0.22 0.50 to 0.78 -1.94 to -0.84
 P<0.001 P<0.001 P<0.001

When we have continuous and categorical predictor variables, regression is also called **analysis of covariance** or **ancova**.

The continuous variables (here height and age) are called **covariates**.

The categorical variables (here cirrhosis) are called **factors**.

Regression in clinical trials

Used to adjust for prognostic variables and baseline measurements.

An example: specialist nurse education for acute asthma

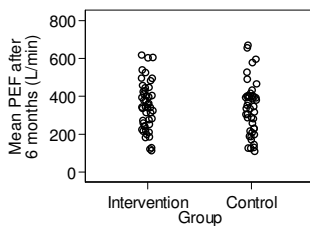
Measurements: peak expiratory flow and symptom diaries made before treatment and after 6 months.

Outcome variables: mean and SD of PEFR, mean symptom score.

Levy ML, Robb M, Allen J, Doherty C, Bland JM, Winter RJD. (2000) A randomized controlled evaluation of specialist nurse education following accident and emergency department attendance for acute asthma. *Respiratory Medicine* 94, 900-908.

Regression in clinical trials

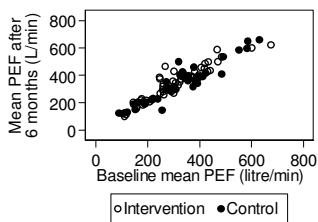
An example: specialist nurse education for acute asthma



Means: 342 338 litre/min
95% CI (intervention – control) –48 to 63 litre/min, P=0.8.

Regression in clinical trials

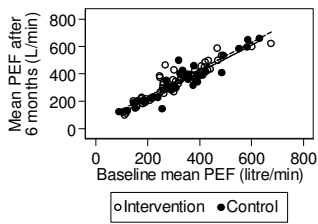
An example: specialist nurse education for acute asthma



If we control for the baseline PEF, we might get a better estimate of the treatment effect because we will remove a lot of variation between people.

Regression in clinical trials

An example: specialist nurse education for acute asthma

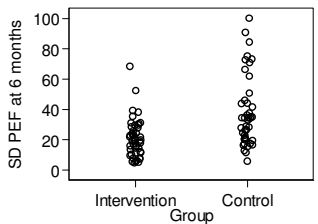


$$\text{PEF@6m} = 58.4 + 0.986 \times \text{PEF@base} + 20.1 \times \text{intervene}$$

$P < 0.001$ $P = 0.046$
95% CI 0.907 to 1.064 0.4 to 39.7

Regression in clinical trials

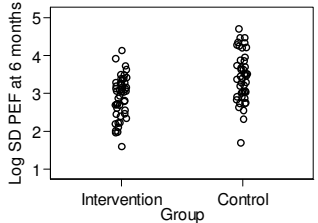
In asthma, large fluctuations in PEF are a bad thing. Use SD.



Clearly we have a skew distribution. Try log transformation.

Regression in clinical trials

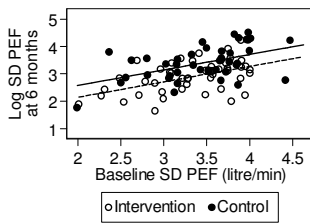
Clearly we have a skew distribution. Try log transformation.



The log scale suggests that this will work.
Log transformed SD of PEF diary, base e.

Regression in clinical trials

Log transformed SD of PEF diary, base e.

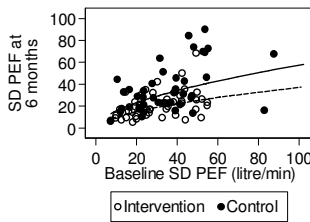


$$\log\text{SD}@6\text{m} = 0.583 + 0.564 \times \log\text{SD}@base - 0.435 \times \text{intervene}$$

$P < 0.001$ $P < 0.001$
95% CI 0.356 to 0.771 -0.651 to -0.218

Regression in clinical trials

SD of PEF diary with fitted lines transformed to natural scale.



Estimated treatment effect = -0.435, 95% CI = -0.651 to -0.218.
Back transform, estimated ratio = 0.65, 95% CI = 0.52 to 0.80.

Regression in clinical trials

Advantages

Reduces variability between subjects and so increase power, narrows confidence intervals.

Removes effects of chance imbalances in predicting variables.

Is adjustment cheating?

It can be if we keep adjusting by more and more variables until we have a significant difference.

We should state before we collect the data what we wish to adjust for and stick to it.

Should include any stratification or minimisation variables, centre in multi-centre trials, any baseline measurements of the outcome variable, known important predictors of prognosis.

Dichotomous outcome variables and logistic regression

Factorial clinical trial: Antidepressant drug counselling and information leaflets to improve adherence to drug treatment.

Patients reporting continuing treatment at 12 weeks

Leaflet	Drug counselling		Total
	Yes	No	
Yes	34/52 (65%)	22/53 (42%)	56/105 (53%)
No	32/53 (60%)	20/55 (36%)	52/108 (48%)
Total	66/105 (63%)	42/108 (39%)	

Peveler R, George C, Kinmonth A-L, Campbell M, Thompson C. Effect of antidepressant drug counselling and information leaflets on adherence to drug treatment in primary care: randomised controlled trial. *BMJ* 1999; **319**: 612-615.

Logistic regression

Patients reporting continuing treatment at 12 weeks

Leaflet	Drug counselling		Total
	Yes	No	
Yes	34/52 (65%)	22/53 (42%)	56/105 (53%)
No	32/53 (60%)	20/55 (36%)	52/108 (48%)
Total	66/105 (63%)	42/108 (39%)	

Counselling: P=0.001 Leaflet: P=0.4

Done by logistic regression.

Logistic regression

Patients reporting continuing treatment at 12 weeks

Leaflet	Drug counselling		Total
	Yes	No	
Yes	34/52 (65%)	22/53 (42%)	56/105 (53%)
No	32/53 (60%)	20/55 (36%)	52/108 (48%)
Total	66/105 (63%)	42/108 (39%)	

Our outcome variable is dichotomous, continue treatment yes or no.

We want to predict the proportion who continue treatment.

We would like a regression equation.

Logistic regression

We want to predict the proportion who continue treatment.

We would like a regression equation:

$$\text{proportion} = \text{intercept} + \text{slope} \times \text{counselling} + \text{slope} \times \text{leaflet}$$

Problem: proportions cannot be less than zero or greater than one. How can we stop our equation predicting impossible proportions?

Find a scale for the outcome which is not constrained.

Odds has no upper limit, but must be greater than or equal to zero.

Log odds can take any value.

Use log odds, called the logit or logistic transformation.

Logistic regression

Predict the log odds of continuing treatment.

$$\text{log odds} = \text{intercept} + \text{slope} \times \text{counselling} + \text{slope} \times \text{leaflet}$$

The slope for counselling will be the increase in the log odds when counselling is used from when counselling is not used.

It will be the log of the odds ratio for counselling, with both the estimate and its standard error adjusted for the presence or absence of the leaflet.

If we antilog, we get the adjusted odds ratio.

Logistic regression

Predict the log odds of continuing treatment.

$$\text{log odds} = \text{intercept} + \text{slope} \times \text{counselling} + \text{slope} \times \text{leaflet}$$

$$\text{log odds} = -0.559 + 0.980 \times \text{counselling} + 0.216 \times \text{leaflet}$$

95% CI	0.426 to 1.53	-0.339 to 0.770
	P=0.001	P=0.4

Antilog:

$$\text{odds} = 0.57 \times 2.66^{\text{counselling}} \times 1.24^{\text{leaflet}}$$

95% CI	1.53 to 4.64	0.71 to 2.16
--------	--------------	--------------

N.B. counselling = 0 or 1, $2.66^0 = 1$, $2.66^1 = 2.66$.

The odds ratio for counselling is 2.66, 95% CI 1.53 to 4.64,
P=0.001.

The odds ratio for the leaflet is 1.24, 95% CI 0.71 to 2.16,
P=0.4.

Interactions

Does the presence of the leaflet change the effect of counseling?

Define an interaction variable = 1 if we have both counseling and leaflet, zero otherwise.

The counseling and leaflet variables are both 0 or 1.

Multiply the counseling and leaflet variables together.

$$\text{Interaction} = \text{counseling} \times \text{leaflet}.$$

$$\text{log odds} = \text{intercept} + \text{slope} \times \text{counseling} + \text{slope} \times \text{leaflet} + \text{slope} \times \text{interaction}$$

$$\text{log odds} = -0.560 + 0.981 \times \text{counseling} + 0.217 \times \text{leaflet} - 0.002 \times \text{interaction}$$

95% CI	0.203 to 1.78	-0.558 to 0.991	-1.111 to 1.107
	P=0.01	P=0.6	P=1.0

Interactions

interaction = counseling × leaflet.

$$\text{log odds} = -0.560 + 0.981 \times \text{counseling} + 0.217 \times \text{leaflet} - 0.002 \times \text{interaction}$$

95% CI	0.203 to 1.78	-0.558 to 0.991	-1.111 to 1.107
	P=0.01	P=0.6	P=1.0

Compare the model without the interaction:

$$\text{log odds} = -0.559 + 0.980 \times \text{counseling} + 0.216 \times \text{leaflet}$$

95% CI	0.426 to 1.53	-0.339 to 0.770
	P=0.001	P=0.4

The estimates of the treatment effects are unchanged by adding this non-significant interaction but the confidence intervals are wider and P values bigger.

We do not need the interaction in this trial and should omit it.

Interactions

BMI data: interaction between AC and MUAC.

$$\text{interaction} = \text{AC} \times \text{MUAC}$$

$$\text{BMI} = -6.44 + 0.18 \times \text{AC} + 0.64 \times \text{MUAC} - 1.39 \times \text{male}$$

P<0.001	P<0.001	P<0.001
---------	---------	---------

Adding the interaction term:

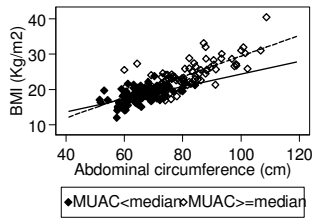
$$\text{BMI} = 8.45 - 0.02 \times \text{AC} + 0.03 \times \text{MUAC} - 1.22 \times \text{male} + 0.0081 \times \text{AC} \times \text{MUAC}$$

P<0.8	P<0.9	P<0.001	P=0.01
-------	-------	---------	--------

If the interaction is significant, both main variables must have a significant effect, so ignore the other P values.

Interactions

BMI data: interaction between AC and MUAC.



Interactions

BMI data: interaction between AC and MUAC.

$$\text{interaction} = AC \times MUAC$$

Adding the interaction term:

$$\text{BMI} = 8.45 - 0.02 \times AC + 0.03 \times MUAC - 1.22 \times \text{male} + 0.0081 \times AC \times MUAC$$

P<0.8 P<0.9 P<0.001 P=0.01

The coefficient for AC now depends on MUAC:

$$\text{slope} = -0.02 + 0.0081 \times MUAC$$

The slope for MUAC depends on AC:

$$\text{slope} = 0.03 + 0.0081 \times AC$$

We cannot interpret the main effects on their own.

Factors with more than two levels

We can use factors with more than two levels, i.e. categorical variables with more than two categories as predictors.

Example:

Serial measurements of tumour size by CT scan and portal vein blood flow (transit time in sec) for six patients (data of Oliver Byass)

Factors with more than two levels

Tumour size by CT scan and portal vein blood flow

Subject 1			Subject 2			Subject 3		
Time	CT	PV	Time	CT	PV	Time	CT	PV
0	10	*	1	12.1	9	1	14.7	13
1	8	*	2	10.4	25	2	14.8	13
2	7.8	13	3	9.4	*	3	14.5	*
3	6.5	11	4	7.2	*	4	14.5	16
4	5.5	18	5	8	18	5	14.1	*
5	4.8	18	6	8.2	*	6	13.6	13.5
6	5	18						

Subject 4			Subject 5			Subject 6		
Time	CT	PV	Time	CT	PV	Time	CT	PV
1	5	19	1	3.6	9	1	7.8	7
2	3.9	15	2	2.6	10	2	6.6	10
3	4.8	17	3	2.6	8	3	5.5	*
4	2.4	18	4	3.2	9	4	4.5	*
			5	3.5	9	5	3.8	10

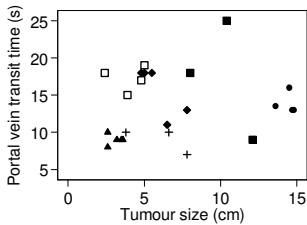
CT = tumour size (cm) by CT scan, PV = portal vein transit time (sec), * = missing data

Factors with more than two levels

We can use factors with more than two levels, i.e. categorical variables with more than two categories as predictors.

Example:

Serial measurements of tumour size by CT scan and portal vein blood flow (transit time in sec) for six patients (data of Oliver Byass)



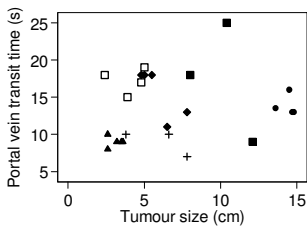
We are interested in whether reduced tumour size is associated with reduced blood flow, not whether people with larger tumour have greater blood flow.

Factors with more than two levels

We can use factors with more than two levels, i.e. categorical variables with more than two categories as predictors.

Example:

Serial measurements of tumour size by CT scan and portal vein blood flow (transit time in sec) for six patients (data of Oliver Byass)



We would like to look at the relationship between tumour size and blood flow within the same subject.

Factors with more than two levels

We would like to look at the relationship between tumour size and blood flow within the same subject.

We can do this by multiple regression or analysis of covariance.

We fit a model which has parallel lines relating transit time and tumour size for each subject separately.

To do this, we need to fit subject as a predictor.

We cannot just put the variable subject number into a regression equation as if it were an interval variable, because there is no sense in which subject 2 is greater than subject 1.

We do not want to assume that our categories are related in this way.

Factors with more than two levels

Instead, we define **dummy variables** or **indicator variables** which enable us to estimate a different mean for each category.

- sub1 = 1 if Subject 1, 0 otherwise,
- sub2 = 1 if Subject 2, 0 otherwise,
- sub3 = 1 if Subject 3, 0 otherwise,
- sub4 = 1 if Subject 4, 0 otherwise,
- sub5 = 1 if Subject 5, 0 otherwise.

If all of these variables are zero, then we have Subject 6. We need five dummy variables to represent a categorical variable with six categories.

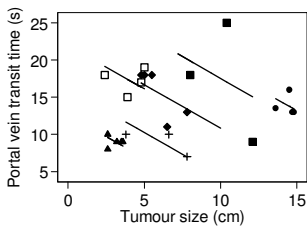
Subject 6 is called the **reference category**.

Factors with more than two levels

We then do regression on our continuous predictor variable and all the dummy variables:

$$PV = 22.5 - 1.17 \times CT + 6.7 \times \text{sub1} + 8.2 \times \text{sub2} - 0.6 \times \text{sub3} - 9.9 \times \text{sub4} - 6.4 \times \text{sub5}$$

P=0.05 P=0.06 P=0.1 P=0.8
P=0.001 P=0.01



We have a significant effect for CT.

Factors with more than two levels

We then do regression on our continuous predictor variable and all the dummy variables:

$$PV = 22.5 - 1.17 \times CT + 6.7 \times \text{sub1} + 8.2 \times \text{sub2} - 0.6 \times \text{sub3} - 9.9 \times \text{sub4} - 6.4 \times \text{sub5}$$

P=0.05 P=0.06 P=0.1 P=0.8
P=0.001 P=0.01

We have a significant effect for CT.

We should ignore the individual tests for the subject coefficients, because they do not mean much.

What we want to know is whether there is any evidence that subject as a whole has an effect.

We get a combined F test for the factor: $F = 6.83$ with 5 and 17 d.f., $P = 0.001$.

Sample size

We should always have more observations than variables.

Rules of thumb:

Multiple regression: at least 10 observations per variable.

Logistic regression: at least 10 observations with a 'yes' outcome and 10 observations with a 'no' outcome per variable.

Otherwise, things get very unstable.

Types of regression

Multiple regression and logistic regression are the types of regression most often seen in clinical trials and in the medical literature in general.

There are many other types for different kinds of outcome variable:

- Cox regression (survival analysis)
- Ordered logistic regression (ordered categories)
- Multinomial regression (unordered categories)
- Poisson regression (counts)
- Negative binomial regression (counts with extra variability)
