

# Part 2 Lecture 1 The Statistical Model







# Pascal Tyrrell, PhD

Associate Professor

Department of Medical Imaging, Faculty of Medicine

Institute of Medical Science, Faculty of Medicine

Department of Statistical Sciences, Faculty of Arts and Science







# CREATING A STATISTICAL MODEL

A researcher is studying the relationship between the decrease in a person's blood pressure and personal and environmental variables, to which they are exposed.

Knowledge from observational studies and laboratory findings suggest that the decrease in a person's blood pressure is related to their sex, age, and a drug believed to reduce blood pressure.





# MODEL FOR DECREASE IN BLOOD PRESSURE (DBP)

MODEL  $DBP_j = K + ERROR_j$  j = 1 to n

$$ERROR_j = DBP_j - K$$

Sum of Squares 
$$SS = \sum_{j=1}^{j=n} (DBP_j - K)^2$$

The Sum of Squared Deviations about K is a minimum if K is equal to the sample mean  $\overline{\text{DBP}}$ .

Therefore the sample mean is called the least squares estimate of K.



MODEL DBP; = K + ERROR; j = 1 to



MODEL FOR DECREASE IN BLOOD PRESSURE (DBP)

MODEL  $DBP_j = \mu + ERROR_j$  j = 1 to n

 $ERROR_{j} = DBP_{j} - \mu$  Residual =  $DBP_{j} - DBP$ 

Sample Variance 
$$S^2 = \frac{\sum_{j=1}^{j=n} (DBP_j - \overline{DBP})^2}{n-1}$$

The Sum of Squared Deviations about the sample mean divided by (n-1) is an unbiased estimator of the variance  $\sigma^2$  that is in the formula of the Gaussian probability density function.



Medical Imaging UNIVERSITY OF TORONT

### STRAIGHT LINE MODEL

$$DBP_j = \beta_0 + \beta_1 \times AGE_j + ERROR_j$$

$$ERROR_j = DBP_j - \beta_0 - \beta_1 \times AGE_j$$

$$j = 1 to n$$





### FITTING A STRAIGHT LINE

Least square estimates of the intercept and slope of a straight line model

$$\widehat{\beta_{1}} = \frac{\sum_{j=1}^{j=n} (AGE_{j} - \overline{AGE}) \times DBP_{j}}{\sum_{j=1}^{j=n} (AGE_{j} - \overline{AGE})^{2}} \qquad \widehat{\beta_{0}} = \overline{DBP} - \widehat{\beta_{1}} \times \overline{AGE}$$

 $RESIDUAL_{j} = OBSERVED - PREDICTED = DBP_{j} - \widehat{\beta_{0}} - \widehat{\beta_{1}} \times AGE_{j}$ 

Sample Variance 
$$s^{2} = \frac{\sum_{j=1}^{j=n} (RESIDUAL_{j})^{2}}{n-2}$$

The Sum of Squared Deviations about the sample mean divided by (n - 2) is an unbiased estimator of the variance  $\sigma^2$  that is in the formula of the Gaussian probability density function.



Medical Imaging UNIVERSITY OF TORONTO

### MULTIPLE LINEAR REGRESSION

$$DBP_{j} = \beta_{0} + \beta_{1} \times DRUG_{j} + \beta_{2} \times SEX_{j} + \beta_{3} \times AGE_{j} + ERROR_{j}$$

j = 1 to n DRUG variable can assume values Aspirin and Tynlenol SEX variable can assume values Female and Male

$$Predicted \ DBP_{j} = \widehat{DBP_{j}} = \widehat{\beta_{0}} + \widehat{\beta_{1}} \times DRUG_{j} + \widehat{\beta_{2}} \times SEX_{j} + \widehat{\beta_{3}} \times AGE_{j}$$

$$RESIDUAL_{j} = DBP_{j} - \widehat{DBP}_{j} \qquad Sample Variance \quad s^{2} = \frac{\sum_{j=1}^{j=n} (RESIDUAL_{j})^{2}}{n-4}$$

Sum of Squared Residuals divided by (n - 4) is an unbiased estimator of the variance  $\sigma^2$  that is in the formula of the Gaussian probability density function.





### DRUG SEX INTERACTION IN MULTIPLE LINEAR REGRESSION

 $DBP_{j} = \beta_{0} + \beta_{1} \times DRUG_{j} + \beta_{2} \times SEX_{j} + \beta_{3} \times DRUG \times SEX + \beta_{4} \times AGE_{j} + ERROR_{j}$ 

j = 1 to n

 $\begin{array}{l} Predicted \ DBP_{j} \\ = \widehat{DBP_{j}} = \widehat{\beta_{0}} + \widehat{\beta_{1}} \times DRUG_{j} + \widehat{\beta_{2}} \times SEX_{j} + \widehat{\beta_{3}} \times DRUG \times SEX + \widehat{\beta_{4}} \times AGE_{j} \end{array}$ 

j = 1 to n

 $RESIDUAL_{j} = DBP_{j} - \widehat{DBP}_{j} \qquad Sample Variance \ s^{2} = \frac{\sum_{j=1}^{j=n} (RESIDUAL_{j})^{2}}{n-5}$ 

If the estimate  $\widehat{\beta_3}$  is large it may mean that there is DRUG \* SEX interaction which means that the size and possibly also the sign of the drug effect is different among males and females. Sample variance  $s^2$  is unbiased estimator of  $\sigma^2$ .





What values should these betas have to minimize the error? Good question! All statistical programs produce values for these unknown betas so that the variation among the error terms is minimized. How is the variation measured? A statistic called the variance measures the variation among the error terms. Values given to the beta constants will be such that the sum of the error terms is zero and their variance is minimized. Is this variance an estimate of the variance that is in the Gaussian probability model. YES !!!!!





TITLE1 " COMPARING SAME MEANS USING GLM PROCEDURE " ;
DATA STUDY ; INPUT COLOUR \$ NAME \$ ID RTIME ;
DATALINES ;

GREEN	ABEL	1	232.6
RED	ABEL	1	232.0
GREEN	ADAM	2	257.5
RED	ADAM	2	250.5
GREEN	AMOS	3	253.1
RED	AMOS	3	237.1
GREEN	ANDY	4	205.4
RED	ANDY	4	201.5
GREEN	BART	5	226.0
RED	BART	5	211.1

RUN; \*\* NOTE: MOST DATASETS HAVE A LINE OF DATA FOR EACH SUBJECT;





### USING GLM PROCEDURE TO COMPARE TWO OR MORE SAMPLE MEANS

TTTTTTII. ASSUMING A COMPLETELY RANDOMIZED DESTGN ; PROC **GLM** DATA = STUDY ; CLASS COLOUR ; MODEL RTIME = COLOUR ; COLOUR / TDIFF LSMEANS PDTFF STDERR CL ; RUN TTTTE1ASSUMING A RANDOMIZED BLOCK DESIGN ; PROC GLM CLASS COLOUR DATA = STUDY; TD : COLOUR ID ; \*\*Note ID in MODEL statement ; MODEL RTIME =COLOUR / LSMEANS TDIFF PDTFF STDERR CL **RUN** ;

NOTE: ID Variable can be replaced by the NAME variable in the CLASS and MODEL statements.





#### ASSUMING A COMPLETELY RANDOMIZED DESIGN

#### The GLM Procedure

#### Dependent Variable: RTIME

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	1	179.776000	179.776000	0.43	0.5323
Error	8	3377.500000	422.187500		
Corrected Total	9	3557.276000			

R-Square	Coeff Var	Root MSE	<b>RTIME Mean</b>
0.050538	8.907232	20.54720	230.6800

Source	DF	Type I SS	Mean Square	F Value	Pr > F
COLOUR	1	179.7760000	179.7760000	0.43	0.5323

Source	DF	Type III SS	Mean Square	F Value	Pr > F
COLOUR	1	179.7760000	179.7760000	0.43	0.5323

MY NOTE: Square Root 0.426 = 0.65 and before we had t = 0.65

#### < SAME AS BEFORE !





#### The GLM Procedure Least Squares Means

COLOUR		Standard Error	H0:LSMEAN=0	H0:LSMean1=LSMean2		
	RTIME LSMEAN		Pr >  t	t Value	Pr >  t	
GREEN	234.920000	9.188988	<.0001	0.65	0.5323	
RED	226.440000	9.188988	<.0001			

COLOUR	RTIME LSMEAN	95% Confidence Limits		
GREEN	234.920000	213.730156	256.109844	
RED	226.440000	205.250156	247.629844	

	Least Squares Means for Effect COLOUR						
i	j	Difference Between Means	95% Confide	ence Limits for LS	SMean(i)-LSMean(j)		
1	2	8.480000		-21.486965	38.446965		





#### ASSUMING A RANDOMIZED BLOCK DESIGN

#### The GLM Procedure

#### Dependent Variable: RTIME

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	5	3465.762000	693.152400	30.30	0.0028
Error	4	91.514000	22.878500		
Corrected Total	9	3557.276000			

R-Square Coeff Var		Root MSE	<b>RTIME Mean</b>	
0.974274	2.073499	4.783147	230.6800	

Source	DF	Type I SS	Mean Square	F Value	Pr > F
COLOUR	1	179.776000	179.776000	7.86	0.0487
ID	4	3285.986000	821.496500	35.91	0.0022

Source	DF	Type III SS	Mean Square	F Value	Pr > F
COLOUR	1	179.776000	179.776000	7.86	0.0487
ID	4	3285.986000	821.496500	35.91	0.0022

#### Matched Pairs Design

#### << MY NOTE: 3465.76 / 3557.27 = 0.974

#### << Same as before





#### The GLM Procedure Least Squares Means

		Standard Error	H0:LSMEAN=0	H0:LSMean1=LSMean2	
COLOUR	RTIME LSMEAN		Pr >  t	t Value	Pr >  t
GREEN	234.920000	2.139089	<.0001	2.80	0.0487
RED	226.440000	2.139089	<.0001		

COLOUR	RTIME LSMEAN	95% Confid	ence Limits
GREEN	234.920000	228.980938	240.859062
RED	226.440000	220.500938	232.379062

Least Squares Means for Effect COLOUR							
i	j	Difference Between Means	95% Confidence Limits for LSMean(i)-LSMean(j)				
1	2	8.480000	0.080898	16.879102			

MY NOTE: 95% Confidence Interval does not contain zero.









LSMEANS DRUG/TDIFF PDIFF STDERR CL; LSMEANS SEX /TDIFF PDIFF STDERR CL; RUN ;

DRUG SEX DRUG \* SEX AGE WEIGHT / SS3;

MODEL RTIME =

**PROC GLM** DATA=STUDY; CLASS DRUG SEX ;

TITLE1 " COMPARING TWO DRUGS ";

WHAT IS THE ADVANTAGE OF GENERAL LINEAR MODEL (GLM) OVER THE TTEST PROCEDURE ??

# Reaction time variable RTIME is OUTCOME variable.

PRIMARY QUESTION: does the variable DRUG predict OUTCOME ? In other words are the DRUG and OUTCOME variables associated ?

If SEX, AGE and WEIGHT predict OUTCOME including them reduces the residual variance, increases  $R^2$  and reduces their p values. If the effect of drug is greater for males than females that is called interaction (DRUG\*SEX).





If the predictor variables SEX, AGE and WEIGHT are also associated with the DRUG variable then excluding them produces biased estimates of the DRUG effect. They are then called CONFOUNDERS.

Product variable DRUG \* SEX is called INTERACTION. If it is large it means that the size of the DRUG effect is different for males and females.







In the cartoon the loser assumed his friend would let go of the rock and feather at the same time. That was not part of the bet.

In many studies researchers may mistakenly think that the comparison was fair and valid. In a study the proportion of males in the exposed and unexposed groups may be quite

different AND if males are at higher risk of disease the

comparison would be biased.





### PROCEDURES SIMILAR TO GLM USED FOR ALL 4 OUTCOMES

THREECOMPLETELY RANDOMIZEDRANDOMIZEDBLOCKDESIGNSSPLIT PLOT

TREATMENTONE WAY2 BY 2 FACTORIALLAYOUT

OUTCOMECONTINUOUSBlood PressureSerum CholesterolVARIABLEBINARYDeath Yes/NoCure Yes/NoCOUNTNumber of deathsNumber of fallsSURVIVAL TIMETime to DeathTime to Cure







### Next up in Part 2 Lecture 2: Study Design



Medical Imaging

UNIVERSITY OF TORONTO

