

Biostatistics

林 建 甫

C.F. Jeff Lin, MD, PhD.

台北大學統計系助理教授
台北榮民總醫院生物統計顧問
美國密西根大學生物統計博士

2005/10/14

JeffLin, MD, PhD

1

Some Aims of Course

- Learn more about how statistics can help solve real biomedical problems.
 - Including how to figure out what the problem really is!
 - What the statistical result means!
- Learn more about communicating the results of statistical manipulations to the ‘client’,
- Gain and improve specific skills in statistical analysis eg. regression, analysis of variance, multivariate analysis,

2005/10/14

JeffLin, MD, PhD

2

Some Aims of Course

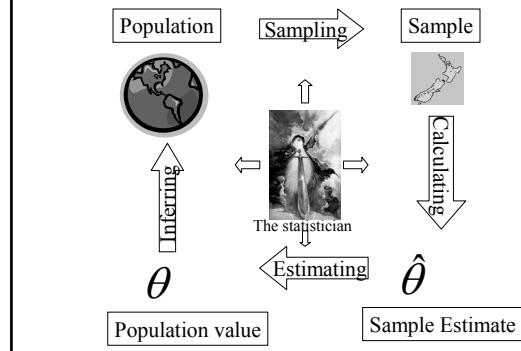
- But **NOT** to learn mathematical theorems which may be behind statistical procedures
- Note the emphasis is on practical statistical inference.
- **R** will be used to perform all statistical calculations. We will also use Excel for data management.
- **R** – <http://lib.stat.cmu.edu/R/CRAN>

2005/10/14

JeffLin, MD, PhD

3

The Statistical Process



2005/10/14

JeffLin, MD, PhD

4

All about R

- Download is about 20+MB
- Quit moaning about it! It won't change my mind ☺
- Keep notes as you work

2005/10/14

JeffLin, MD, PhD

5

Exploration and Presentation of Univariate Data

- Single variable. (there may be multiple samples of the same variable)
- E.g. Reported cancer death for 2003 in Taiwan
 - We could group these into regions.
 - We will look at responses grouped into males and females.

2005/10/14

JeffLin, MD, PhD

6

Statistical Inference About a Relationship

- We have observations on two or more variables for a number of items.
- We believe there is an underlying (linear) relationship between the variables, but we don't know the parameters.
- Assuming a **random** error structure on the observations.
- We make inferences about the parameters of the relationship.
- Whether or not the conclusions extend to the population of items depends how the ones we measured were selected.

2005/10/14

Jeff Lin, MD, PhD

7

Inference on a Relationship

- Observational study--which items go into the treatment group(s) and the control group are not determined using randomisation.
- Randomised Experiment--we use randomisation to determine which items go into treatment groups and which go into control group.
- We can infer the relationship is a **causal** one only when the data comes from a randomised experiment.

2005/10/14

Jeff Lin, MD, PhD

8

Univariate Exploratory Data Analysis (EDA)

- Univariate means the data is comes from measurement of one variable
- Sometimes we have additional information to divide the data into groups
- Analysis and depiction depends on the type of data

2005/10/14

Jeff Lin, MD, PhD

9

Different Types of Data: Statistics

- Discrete data
- Continuous data

2005/10/14

Jeff Lin, MD, PhD

10

Different Types of Data: Statistics

- Discrete data
 - Count data
 - Categorical data
 - Ordinal data
- For example:
 - Number of road accidents due to speed in a year
 - Ethnicity question on the Taiwan Census
 - Socio-economic class (low/med/high)

2005/10/14

Jeff Lin, MD, PhD

11

Different Types of Data: Statistics

- Continuous data
 - Measurements

For example:

- Income in dollars
- Amount of poison (in g/L) required to kill subject
- Age, WT, HT, BP, Cholesterol, Blood sugar

2005/10/14

Jeff Lin, MD, PhD

12

Different Types of Data: Statistics

- Sometimes discrete data can be treated as continuous data, especially when there are many possible outcomes.
- For example – smashing a window and counting the number of glass fragments on 40cmx40cm sheet of paper around the breaker (2-4000)

2005/10/14

Jeff Lin, MD, PhD

13

Plotting Data

- Why plot data?
 - Often the best way to understand the data
- What do we look for when we plot data?
 - Depends on what kind of plot we have.
- How do we select the appropriate plot?
 - Depends on the data

2005/10/14

Jeff Lin, MD, PhD

14

Plotting Discrete Data

- For count data we have 4 options
 1. Dotplots
 2. Stem and leaf plots
 3. Histograms
 4. Boxplots

2005/10/14

Jeff Lin, MD, PhD

15

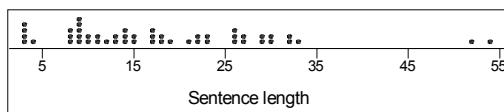
- In the following examples we have data on sentence lengths (in words) from 8 books
 1. Eye of the Dragon
 2. The Tommy- knockers
 3. The Shining
 4. The Stand
 5. The Dark Half
 6. Four Past Midnight
 7. Rising Sun
 8. Disclosure
 - The first six are written by Stephen King and the last two by Michael Crichton

2005/10/14

Jeff Lin, MD, PhD

16

Dotplot



From the dotplot we can see:

- The location of the data
- The spread of the data
- Extremes in the data (outliers)
- Individual values

2005/10/14

Jeff Lin, MD, PhD

17

Stem and Left Plot

Stem-and-leaf of The Stand N = 50
The decimal point is 1 digit(s) to the right of the |

0	33334
0	88899999
1	0011233444
1	55777889
2	12233
2	6667799
3	00223
3	
4	
4	
5	24

We can see:

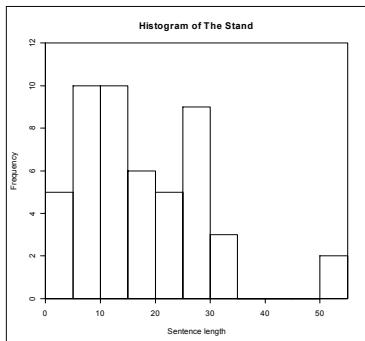
- Location
- Spread
- Order statistics
- Individual values

2005/10/14

Jeff Lin, MD, PhD

18

Histogram



We can see:

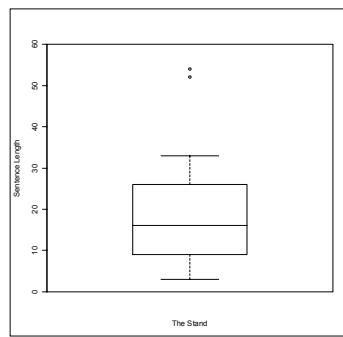
- Location
- Spread
- Approximate mode
- Outliers

2005/10/14

Jeff Lin, MD, PhD

19

Boxplot



We can see:

- Location
- Spread
- Median
- IQR
- Outliers

2005/10/14

Jeff Lin, MD, PhD

20

Which Plot?

- It should be obvious that each of the plots we've seen has strengths and weaknesses
- Dotplots
 - Good for small data sets ($n < 50$)
 - and data that is not over dispersed
- Stem and leaf plots
 - Good for very small data sets ($n < 20$)
 - and data that is not over dispersed
 - Useful for ordering data
 - Not useful for any real problem.
 - Sensitive to choice of stems

2005/10/14

Jeff Lin, MD, PhD

21

Which Plot?

- Histograms
 - Good for larger data sets (> 20)
 - Often the best way to summarise data
 - Sensitive to choice of “class intervals”
 - Doesn't highlight individual values
- Boxplots
 - Good for larger data sets (> 15)
 - Often the best way to summarise grouped data
 - Doesn't highlight individual values
 - Does highlight potential outliers (3SD)

2005/10/14

Jeff Lin, MD, PhD

22

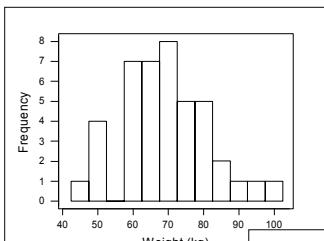
Plotting Continuous (cts.) Data

- The choice of plots is almost the same
 - Dotplots are not much use for cts. data
 - Similarly steam and leaf plots
 - Histograms and boxplots display continuous data well
 - Histograms are even more sensitive for cts. data
 - i.e. no “natural” class breaks

2005/10/14

Jeff Lin, MD, PhD

23



Histogram shows:

- Shape of data well
- Spread and mode

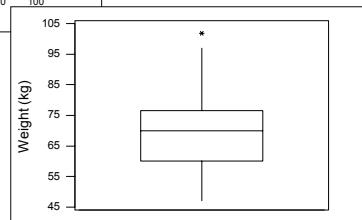
Boxplot shows:

- Spread
- Median (IQR)
- Outlier?

2005/10/14

Jeff Lin, MD, PhD

24



Which Plot?

- Think about the following:
 - What am I trying to show?
 - Which features of the data are important?
 - Shape?
 - Location/Spread?
 - Outliers?
 - Is the plot easy to understand?
 - Chart junk!
 - Clutter!

2005/10/14

Jeff Lin, MD, PhD

25

Grouped Data

- Some ancillary (additional) information on grouping
- We want to compare between groups. Why?
 - To highlight or look for differences between groups
- Points to remember
 - Compare apples with apples!

2005/10/14

Jeff Lin, MD, PhD

26

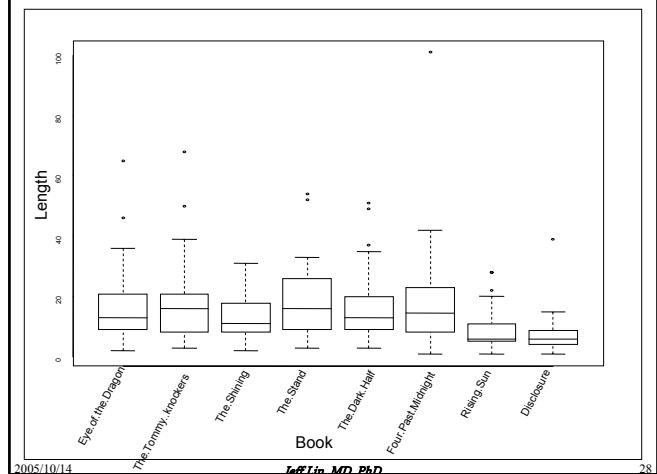
Plotting Grouped Data

- Two groups
 - Back to back stem and leaf plots
 - Parallel or back to back dot plots
- Two or more groups
 - Histogram matrix
 - Boxplots
- Boxplots are often the best way to compare groups

2005/10/14

Jeff Lin, MD, PhD

27



28

Data Description and Transformation

- What do we look for when we describe univariate data?
 - Peaks/Bumps
 - Location
 - Spread
 - Outliers

2005/10/14

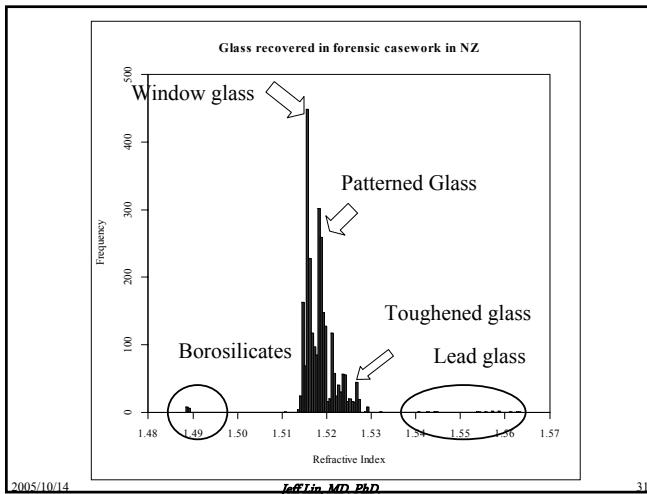
Jeff Lin, MD, PhD

29

2005/10/14

Jeff Lin, MD, PhD

30



Data Description

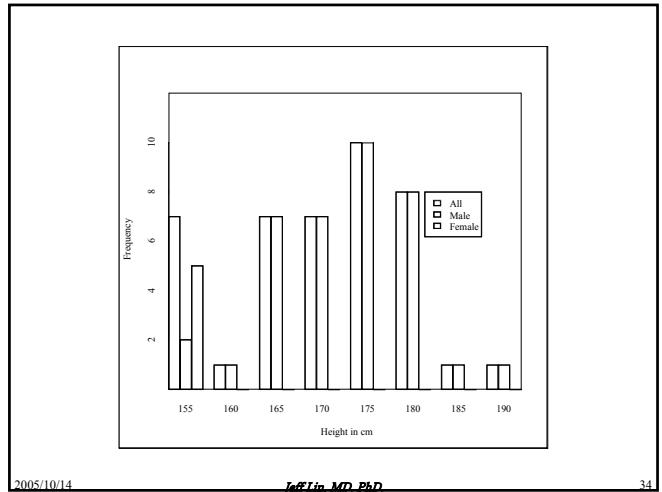
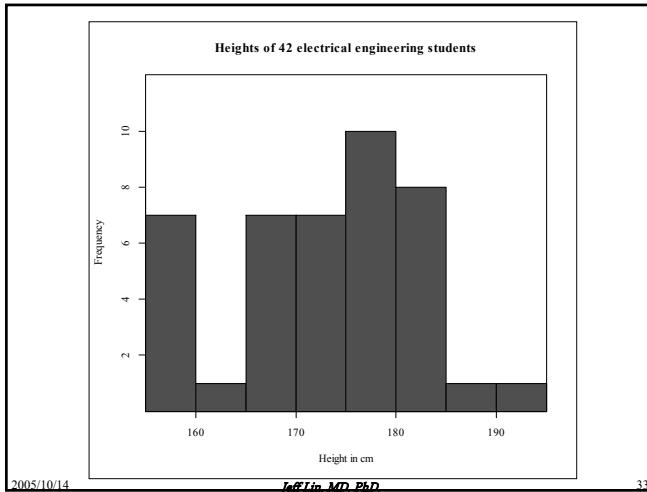
- Peaks/Bumps
 - Unimodal
 - Bimodal
 - Multimodal

Why is there more than one mode?

- Contamination
- Mixing of more than one process?

2005/10/14 JeffLin MD PhD

32



34

Outliers

- Outliers are:
 - Abnormal data points – “different from the bulk of the data”.
- They may be:
 - Data entry errors
 - Measurement errors
 - Just really different pieces of data
- Outliers should always be investigated!

2005/10/14 JeffLin MD PhD 35

Symmetry and Skewness

- Data may also be described in terms of symmetry or skewness
- If the extreme values of the data set lie to the left of the mean or median then the data is said to be left skewed.
- Conversely, if the extreme values of the data set lie to the right of the mean or median then the data is said to be right skewed.

2005/10/14 JeffLin MD PhD

36

- Right Skewed/Positively Skewed
- Long tail to the **right**
- Bulk of the data to the **left** of the mean $\Pr(X > \bar{x}) \approx 0.36$
- Mean = 0.20
- Left skewed/Negatively Skewed
- Long tail to the **left**
- Bulk of the data to the **right** of the mean $\Pr(X < \bar{x}) \approx 0.61$
- Mean = 0.93

2005/10/14 Jeff Lin, MD, PhD 37

Symmetry takes many different forms

2005/10/14 Jeff Lin, MD, PhD 38

Summary Statistics

- Measures of location
 - mean
 - median
- Measures of spread
 - variance/std. deviation
 - range
 - interquartile range (IQR)

2005/10/14 Jeff Lin, MD, PhD 39

Order Statistics and Quantiles

If our data are x_1, x_2, \dots, x_n then the order statistics are simply the data sorted into ascending order and denoted $x_{(1)}, x_{(2)}, \dots, x_{(n)}$ and $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$

If α is of the form $\alpha_i = \frac{i}{n+1}$ for some integer i , define the corresponding quantile to be the i th order statistic, i.e. $q_{\alpha_i} = x_{(i)}$

- The α quantile, q_{α} , is chosen such that approximately $100\alpha\%$ of the data is less than q_{α} .
- For example if $\alpha = 0.1$ then approximately 10% of the data has a value less than the 0.1 quantile $q_{0.1}$.

2005/10/14 Jeff Lin, MD, PhD 40

Example

- Let $x = \{9, 6, 10, 5, 3, 1, 5, 5, 10, 6\}$
- Then the order statistics are $\{1, 3, 5, 5, 5, 6, 6, 9, 10, 10\}$
- The first order statistic $x_{(1)} = 1$, is the $\alpha = \frac{1}{10+1} \approx 0.09$ quantile
- That is, $x_{(1)} = q_{0.09}$
- The last order statistic is $x_{(10)} = 10$, is 0.91 quantile
- Why isn't $x_{(10)}$ the 1.0 quantile?
- Imagine that we're finding quantiles for the normal distribution. What is the point x such that $\Pr(X < x) = 1$?
- The answer is $+\infty$ (positive infinity). This will cause us problems later on.

2005/10/14 Jeff Lin, MD, PhD 41

Quantiles

- The 0.25 quantile is the lower quartile
- The 0.5 quantile is the median
- The 0.75 quantile is the upper quartile

If you want a quantile from the data that is not represented by any α_i (e.g. if $n = 10$, and $\alpha = 0.1$), then it can be interpolated from the data by taking a weighted sum of the two surrounding quantiles

2005/10/14 Jeff Lin, MD, PhD 42

The Symmetry Statistic

- The symmetry statistic provides a numerical summary of symmetry.
- It is defined as:

$$S = \frac{q_{0.5} - q_\alpha}{q_{(1-\alpha)} - q_\alpha}$$

“The ratio of the distance from the median to the lower quantile to the interquantile difference”

α is usually chosen to be 0.1 or 0.05. We will use 0.1

2005/10/14

Jeff Lin, MD, PhD

43

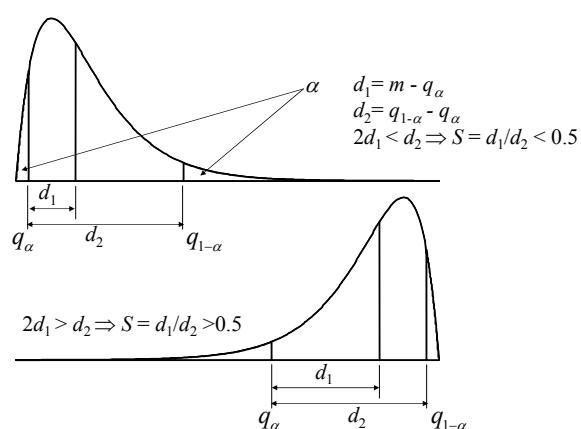
The Symmetry Statistic

- If S is less than 0.5 then the data is right skewed
 - The α quantile is closer to the median than the $1 - \alpha$ quantile
- If S is greater than 0.5 then the data is left skewed
 - The $1 - \alpha$ quantile is closer to the median than the α quantile

2005/10/14

Jeff Lin, MD, PhD

44



2005/10/14

Jeff Lin, MD, PhD

45

Transforming Data

- Sometimes we want to make data more symmetric. Why?
 - Makes it easier to see what is going on
 - Meet the assumptions of some parametric test
 - Makes the data better behaved
- How do we transform the data?
 - Can apply any mathematical function but typically we use power transformations

2005/10/14

Jeff Lin, MD, PhD

46

Power Transformations

- Note these only work for positive data
- Can raise data to an positive integer or fractional power p , for example $p=0.5$ (square root) or $p=3$ (cubed data).
- Or raise data to a negative power – remember $x^p = 1/x^{-p}$.
- When $p=0$, we use the logarithm. Why?

2005/10/14

Jeff Lin, MD, PhD

47

Choosing an Optimal Transformation

- In general:
 - Increasing power (>1) will correct left skew
 - Decreasing power (<1) will correct right skew
 - Always choose a transformation you can explain (rounding is a good start)
 - A log transformation is often all you need
- Symmetry power plot
 - Plot power (p) versus symmetry statistic S

2005/10/14

Jeff Lin, MD, PhD

48

05BioST02 EDA

Optimal Transformation

- Note: for $p < 0$ we must transform the data x to $-x^p$ to preserve the order of the data.
- Box-Cox transformations

$$x^{(p)} = \begin{cases} \frac{x^p - 1}{p}, & p \neq 0 \\ \log x, & p = 0 \end{cases}$$

- Box-Cox transformations try to make the data more normal
- Box-Cox method assumes data raised to some power p , comes from a normal distribution with mean μ and std. dev. σ
- To find p , μ and σ we use maximum likelihood estimation (mle)

2005/10/14

Jeff Lin, MD, PhD

49

Box-Cox Transformations

Minimise:

$$n \log s^{(p)} - (p - 1) \sum_{i=1}^n \log x_i$$

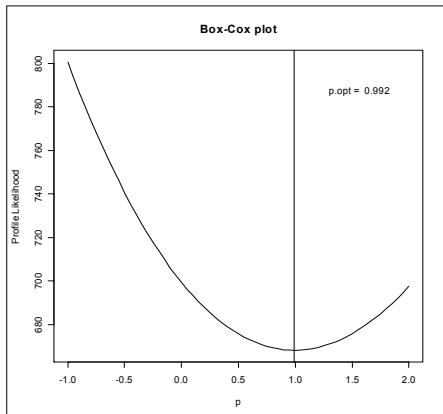
This is a “profile likelihood” – which you might learn about in the third or fourth year. **YOU DO NOT NEED TO REMEMBER THIS OR EVEN KNOW HOW IT WORKS!**

Typically this is minimised using a graphical method or by some sort of numerical optimisation routine. R can do both. The routine I have supplied is available on the web page and from the server

2005/10/14

Jeff Lin, MD, PhD

50



2005/10/14

Jeff Lin, MD, PhD

51

Transformations - Summary

- Only works on positive data**
- Box-Cox usually works well
- Make sure you can explain your transformation
- Taking the log often is the easiest thing to do
- Don't transform if you don't need to!

2005/10/14

Jeff Lin, MD, PhD

52

Symmetry and Skewness Measures

- We've seen the symmetry statistic S .
 - Robust to outliers
 - Crude symmetry statistic can be obtained using the quartiles. These can be gained in R by typing: `quantile(x)`
 - For the standard symmetry statistic, there is a function called s in the file called `symmetry.r` on the web and in the course folder. We type `s(x)`

2005/10/14

Jeff Lin, MD, PhD

53

Skewness Statistic

- Traditional measure of skewness is defined as:

$$\text{Skewness} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{s^3}$$

where s is the sample standard deviation. This is sometimes called the 3rd central moment (the variance is the second)

- If this statistic is positive, then the data are right skewed and
- if it is negative then they are left skewed.

2005/10/14

Jeff Lin, MD, PhD

54

Quantile Plots

- Box-Cox transformations try to “transform to normality”
- Is there any way we can find out whether the data behave as though they came from a normal distribution?
- Yes – a normal quantile plot or a norplot.

2005/10/14

Jeff Lin, MD, PhD

55

Norplots

- Define the quantiles of a standard normal distribution $N(0,1)$

$$z_\alpha = z \text{ s.t. } \Pr(Z < z) = \alpha$$

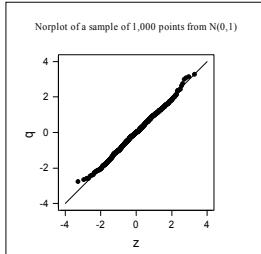
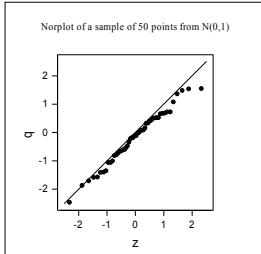
If the data come from a standard normal distribution then the empirical quantiles will be in approximately the same position as the theoretical quantiles

2005/10/14

Jeff Lin, MD, PhD

56

Example of a norplot



2005/10/14

Jeff Lin, MD, PhD

57

Norplots

- What happens if the data don't come from a standard normal?
- Remember that if $Y \sim N(\mu, \sigma)$ and $Z \sim N(0,1)$ then

$$Y = \sigma Z + \mu$$

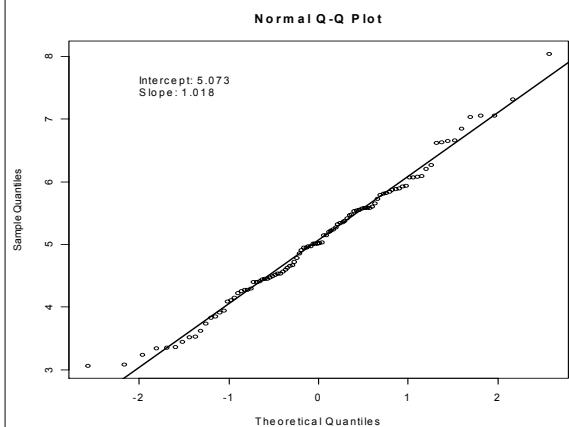
- Therefore if the data are normal, but not standard normal, the data will have slope σ and intercept μ

2005/10/14

Jeff Lin, MD, PhD

58

Normal Q-Q Plot



2005/10/14

Jeff Lin, MD, PhD

59

Quantile Plots

- It is relatively easy to construct quantile plots for any of the standard distributions.
- However, the parameters are often unknown and so most people don't bother for anything other than a uniform distribution
- Sometimes, however, we wish to see if two sets of data are similarly distributed

2005/10/14

Jeff Lin, MD, PhD

60

QQ-Plots

- QQ or quantile-quantile plots compare the quantiles of two sets of data
- If the data sets are similarly distributed then the quantiles will be similar and hence the points will cluster around the 45° line
- If the data sets have the same spread but differing means the relationship will be linear but the intercept will be different from zero
- If the data sets have the differing means and spreads the relationship will be linear but the slope will be different from 1 and the intercept non-zero

2005/10/14

Jeff Lin, MD, PhD

61

Constructing a QQ-plot

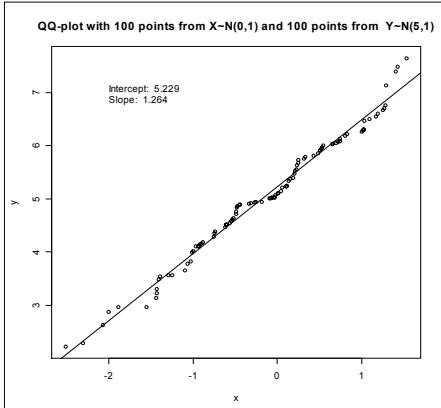
- Let X represent data set 1, and Y data set 2.
- If X and Y are of equal length then we can plot the order statistics of X vs. the order statistics of Y
- If $n_x > n_Y$ then we need to interpolate the quantiles from Y with

$$\alpha_i = \frac{i}{n_X + 1}$$

2005/10/14

Jeff Lin, MD, PhD

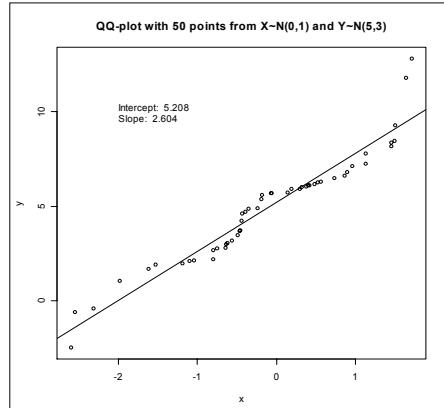
62



2005/10/14

Jeff Lin, MD, PhD

63

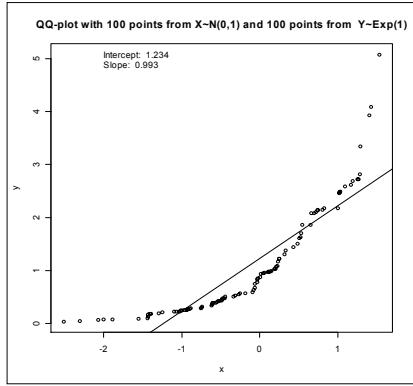


2005/10/14

Jeff Lin, MD, PhD

64

QQ-Plot – Different Distributions



2005/10/14

Jeff Lin, MD, PhD

65

Some Points about QQ-plots

- Easier to interpret if aspect ratio is 1:1
- Good for detect mean and variance/std. dev. shifts
- Not really convincing evidence against normality (need a lot of data)
- Heavy curvature usually means the data sets are from different distributions

2005/10/14

Jeff Lin, MD, PhD

66

A Small Case Study

- Description: Rainfall from Cloud-Seeding. The rainfall in acre-feet from 52 clouds, 26 of which were chosen at random and seeded with silver nitrate.
- Reference: Chambers, Cleveland, Kleiner, and Tukey. (1983). Graphical Methods for Data Analysis. Wadsworth International Group, Belmont, CA, 351.
- Original Source: Simpson, Alsen, and Eden.(1975). A Bayesian analysis of a multiplicative treatment effect in weather modification. Technometrics 17, 161-166.

2005/10/14

Jeff Lin, MD, PhD

67

Number of cases: 26

Variable Names:

- Unseeded:** Amount of rainfall from unseeded clouds (in acre-feet)
- Seeded:** Amount of rainfall from seeded clouds with silver nitrate (in acre-feet)

Unseeded	Seeded	Unseeded	Seeded
1202.6	2745.6	41.1	200.7
830.1	1697.8	36.6	198.6
372.4	1656	29	129.6
345.5	978	28.6	119
321.2	703.4	26.3	118.3
244.3	489.1	26.1	115.3
163	430	24.4	92.4
147.8	334.1	21.7	40.6
95	302.8	17.3	32.7
87	274.7	11.5	31.4
81.2	274.7	4.9	17.5
68.5	255	4.9	7.7
47.3	242.5	1	4.1

2005/10/14

Jeff Lin, MD, PhD

68

Step 1: Visually Inspect the Data

- It is easier to look at the order statistics \Rightarrow sort the data. First let's make the data easier to work with

```
clouds<-data.frame(seed=Seeded, u.seed=Unseeded)
attach(clouds)
```

Now sort the data. Type:

```
seed<-sort(seed)
u.seed<-sort(u.seed)
```

Type `clouds` to display the data

2005/10/14

Jeff Lin, MD, PhD

69

Visual Inspection of the Data

- There are some very large values
- There are some very small values
- The data does not increase linearly
- The data sets look as though they increase at the same rate but there is a shift in mean

2005/10/14

Jeff Lin, MD, PhD

70

Plot the Data

- A histogram of each data set would be helpful
 - Type:

```
par(mfrow=c(2,1))      # this sets the
graph screen    # to 2 rows 1 column
hist(u.seed)
hist(seed)
```

2005/10/14

Jeff Lin, MD, PhD

71

Interpret the Histograms

- Note the histograms don't have the same scales so it makes them difficult to compare
- Each histogram has a long tail to the right so the data sets are right skewed.
- Confirm this by calculating a symmetry statistic for each.
- Type:

```
s(seed)
s(u.seed)
```

Note the 0.1 – this means we're using the 0.1 quantile and the 0.9 quantile to calculate the symmetry statistic

2005/10/14

Jeff Lin, MD, PhD

72

Get Summary Statistics

- Type:
`summary(seed)`
`summary(u.seed)`
- This gives you a five number summary (Min, LQ, Med, UQ, Max) plus the mean
- To get the standard deviations type
`sd(seed)`
`sd(u.seed)`
- To get the skewness statistics type
`skewness(seed)`
`skewness(u.seed)`

2005/10/14

Jeff Lin, MD, PhD

73

Descriptive Statistics

- Note that the skewness statistic confirms our histogram interpretation and our symmetry statistic.
- We can plot the mean and median on each histogram:
`hist(seed)`
`abline(v=median(seed), col="red", lwd=2)`
`abline(v=mean(seed), col="blue", lwd=2)`
`hist(u.seed)`
`abline(v=median(u.seed), col="red", lwd=2)`
`abline(v=mean(u.seed), col="blue", lwd=2)`
 Notice how the median lines up with the centre of the data much better than the mean.
- If we want to compare groups it is going to be easier if we transform the data

2005/10/14

Jeff Lin, MD, PhD

74

Transforming the Data

- Choose a sensible transformation either using the `sympowerplot` function, e.g.
`sympowerplot(u.seed)`
`sympowerplot(seed)`
- Or by using the Box-Cox method.
`bcp(u.seed)`
`bcp(seed)`

2005/10/14

Jeff Lin, MD, PhD

75

Transforming the Data

- Box-Cox recommends $p = 0.047$ and $p = 0.118$ – therefore a log transformation might work well. If you're willing to wait awhile try
 - `bcp.bounds(seed)` and
`bcp.bounds(u.seed)`
 - These will give approximate 95% confidence intervals on the powers and they both include 0
- Symmetry power plots recommend $p = -0.1$ and $p = 0.16$ for unseeded and seeded respectively.
- Log is easiest to explain, so try that.
- Look at log transformed histograms
 - Still lumpy, but skewness has been reduced

2005/10/14

Jeff Lin, MD, PhD

76

Compare the Data

- Boxplot the transformed data
 - Join the data into one vector
`rainfall<-c(log.u.seed, log.seed)`
 - Use the `rep` command to make a vector of labels corresponding to the unseeded and seeded data
`gp<-rep(c("Unseeded","Seeded"),c(26,26))`
`boxplot(split(rainfall, gp))`
- QQ-Plot the data
`qqplot(log.u.seed,log.seed)`
`qq.plot(log.u.seed,log.seed)`

2005/10/14

Jeff Lin, MD, PhD

77

Compare the Data

- QQ-plot shows linear trend \Rightarrow logged data are from similar distributions
- Slope is approximately equal to 1 \Rightarrow similar spread
- Is the difference in the means significant ?
- Two-sample t -test

2005/10/14

Jeff Lin, MD, PhD

78

Bivariate Data

2005/10/14

Jeff Lin, MD, PhD

79

Bivariate Data

- Bivariate data is simply data where we have measurements on two variables.
- The variables can be either discrete or continuous, leading to three possible combinations: discrete/discrete, continuous/ discrete, continuous/continuous
- The type of analysis depends on the combination you have.
- Note if one of your variables is in continuous time units (such as years or seconds) then time series analysis is usually the appropriate treatment. We will not cover this subject

2005/10/14

Jeff Lin, MD, PhD

80

Displaying Discrete/Discrete Data

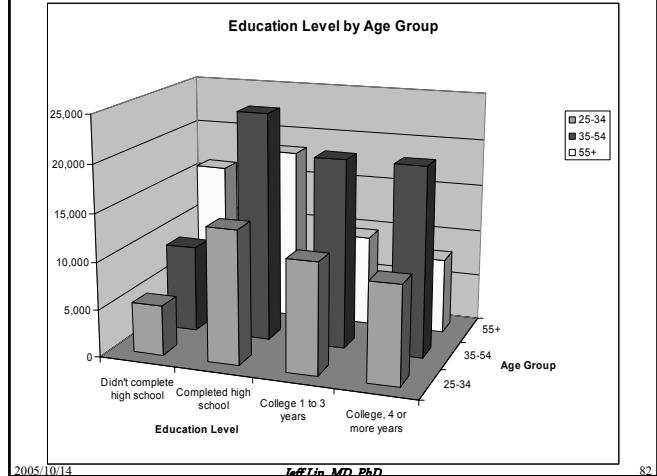
- Some texts recommend a three dimensional barplot

	Age Group			
Education	25-34	35-54	55+	Total
Didn't complete high school	5,325	9,152	16,035	30,512
Completed high school	14,061	24,070	18,320	56,451
College 1 to 3 years	11,659	19,926	9,662	41,247
College, 4 or more years	10,342	19,878	8,005	38,225
Total	41,387	73,026	52,022	166,435

2005/10/14

Jeff Lin, MD, PhD

81



2005/10/14

Jeff Lin, MD, PhD

82

3D Perspective Bar Charts

- Difficult to read
- Increase in categories leads to increase in complexity
- Increase in categories leads to more obscuration of data/features
- Perspective distorts the information
- **Don't do it!**

2005/10/14

Jeff Lin, MD, PhD

83

Displaying Discrete/Discrete Data

- **Two-way tables** are usually the best
- Think about what your aims are:
 - If the data are for someone else to use in an analysis or for reference in an appendix, then apart from organising the data in a legible fashion don't do anything to it.
 - If the data are for presentation or summary then apply Ehrenberg's rules

2005/10/14

Jeff Lin, MD, PhD

84

Displaying Tabular Data – Ehrenberg's Rules

1. Round drastically “2 busy digits”
2. Arrange the numbers so that comparisons are made column-wise and not row-wise
3. Order the columns by size or use ordinal information
4. Use row and column averages or sums as a focus
5. Use white space well
6. Provide verbal summaries

2005/10/14

Jeff Lin, MD, PhD

85

Education	Age Group			Total
	25-34	35-54	55+	
Didn't complete high school	5,325	9,152	16,035	30,512
Completed high school	14,061	24,070	18,320	56,451
College 1 to 3 years	11,659	19,926	9,662	41,247
College, 4 or more years	10,342	19,878	8,005	38,225
Total	41,387	73,026	52,022	166,435

2005/10/14

Jeff Lin, MD, PhD

86

Education	Age Group			Total
	25-34	35-54	55+	
Didn't complete high school	5,300	9,100	16,000	30,400
Completed high school	14,000	24,000	18,300	56,300
College 1 to 3 years	11,700	20,000	9,700	41,400
College, 4 or more years	10,300	20,000	8,000	38,300
Total	41,300	73,100	52,000	166,400

	Didn't complete high school	Completed high school	College 1 to 3 years	College, 4 or more years
25-34	17%	25%	28%	27%
35-54	30%	43%	48%	52%
55+	53%	33%	23%	21%

2005/10/14

Jeff Lin, MD, PhD

87

Statistical Analysis of Discrete Data

- There are two possible situations when dealing with bivariate discrete data
 1. Single sample cross-classified by two variables
 - E.g. 221 students classified by age group and birth order
 2. Multiple samples classified by a single variable
 - Genotype of individuals in different ethnic DNA databases
- As you would expect the methods of analysis are quite different

2005/10/14

Jeff Lin, MD, PhD

88

Multiple Samples Classified By a Single Categorical Variable

- If the samples are independent (and they usually are) then the category proportions can be compared across samples
- For example, in the US the FBI has DNA databases for African Americans, Caucasians and South Western Hispanics (Florida keeps its own database for South Eastern Hispanics)

2005/10/14

Jeff Lin, MD, PhD

89

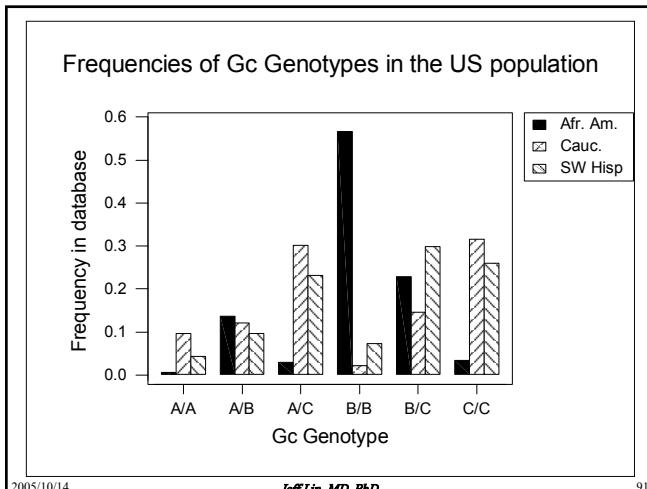
Analysing Multiple Sample Discrete Data

- Say we wished to compare genotype frequencies at a particular DNA locus called Gc
- Gc has 3 alleles A, B, and C. We have one allele from our mother and one from our father to make a genotype, so there are six possible genotype frequencies: A/A, A/B, A/C, B/B, B/C, and C/C
- How do we go about comparing the proportions of each genotype?
- We can plot them – always a good idea

2005/10/14

Jeff Lin, MD, PhD

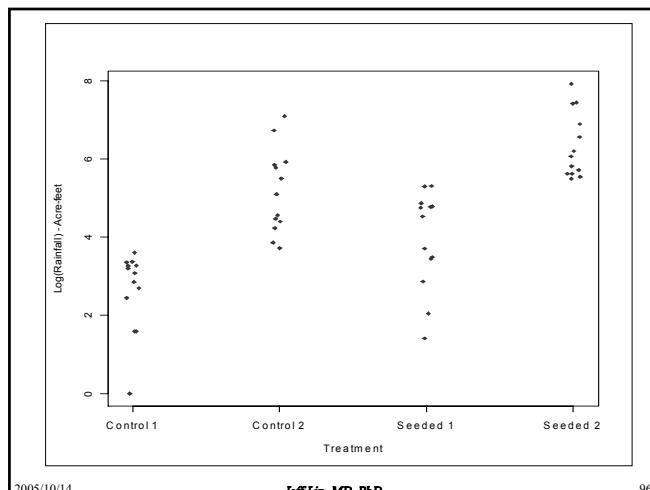
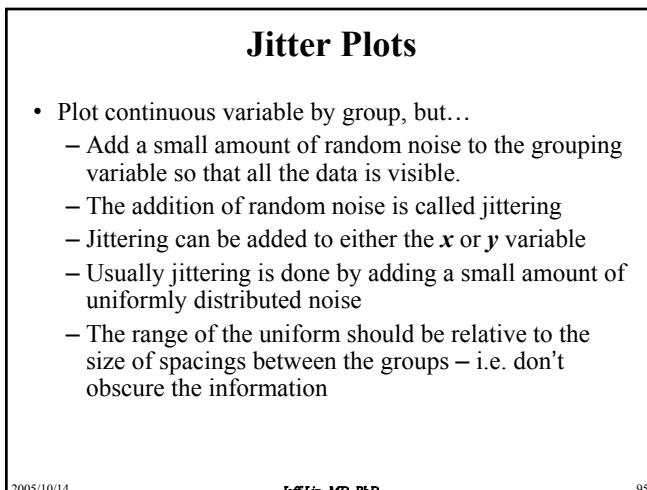
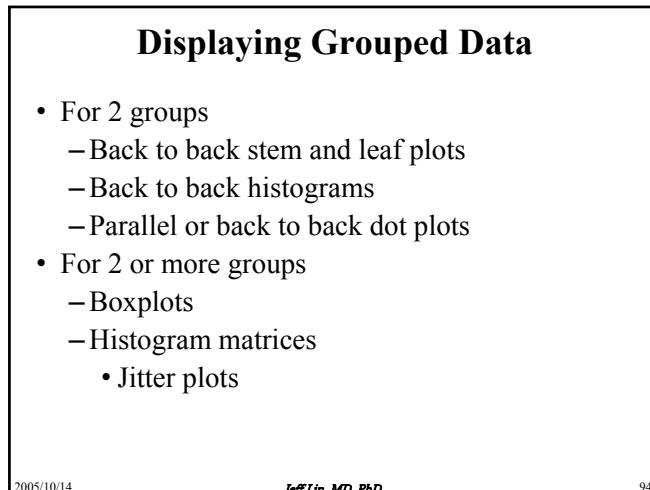
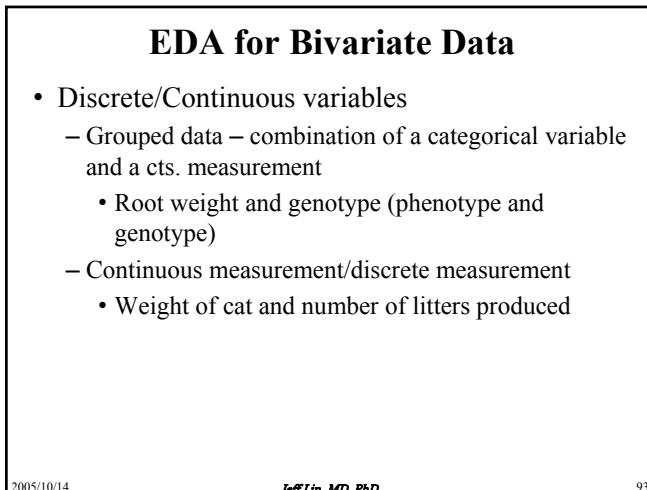
90



Analysing Multiple Sample Discrete Data

- How about a statistical test?
- It boils down to the question we want to answer.
- If that question is: "Is there a difference in genotype probabilities between the races" then we can answer it. How?
- A chi-square analysis.
- Note, it turns out that this is the same analysis that we apply to discrete data from a single sample cross classified by two (or more) categorical variables, so we will assume we are talking about either case from now on.

2005/10/14 Jeff Lin, MD, PhD 92



Analysing Grouped Data

- 2 groups
 - Paired data (before and after data say)
 - Paired *t*-test
 - Sign test (non-parametric)
 - 2 independent groups
 - 2 sample *t*-test
 - Wilcoxon rank sum test (non-parametric)

2005/10/14

Jeff Lin, MD, PhD

97

Analysing Grouped Data

- k independent groups
 - One-way analysis of variance (ANOVA)
 - Kruskal-Wallis test (non-parametric)
- Plotting options for k groups
 - LSD plots
 - Tukey's HSD

2005/10/14

Jeff Lin, MD, PhD

98

Displaying Continuous Bivariate Data

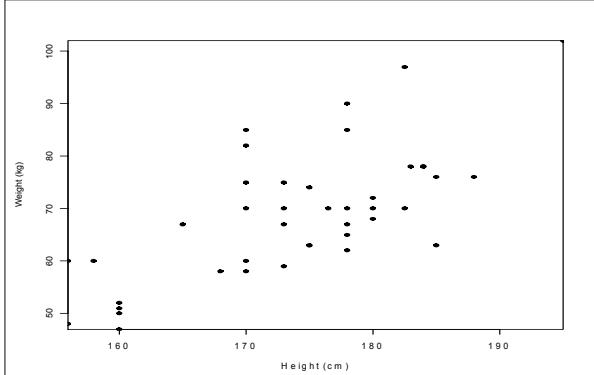
- Scatterplot
 - Looking for a relationship between the two variables
 - This could be a **linear** relationship – the data follows a straight line
 - The relationship could be **non-linear** – the data follows a curve
 - It is rare that any relationship is perfect

2005/10/14

Jeff Lin, MD, PhD

99

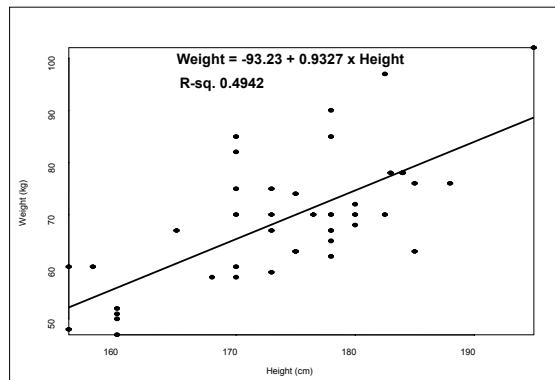
A Linear Relationship?



2005/10/14

Jeff Lin, MD, PhD

100



2005/10/14

Jeff Lin, MD, PhD

101

Non-linear Relationships

- Sometimes non-linear relationships can be described by simple mathematical functions
 - Log/Exponential

$$y = \exp(\beta_0 + \beta_1 x)$$
 - Polynomial regression

$$y = \beta_0 + \beta_1 x^2$$
 - Trigonometric functions (uncommon)

$$y = \beta_0 \sin(x)$$

2005/10/14

Jeff Lin, MD, PhD

102

05BioST02 EDA

More Advanced Problems

- Sometimes no mathematical function will adequately describe the relationship that exists in the data
- Solutions to this problem are harder
 - lowess (locally weighted least squares)
 - basis splines (beyond the scope of this course)

2005/10/14

Jeff Lin, MD, PhD

103

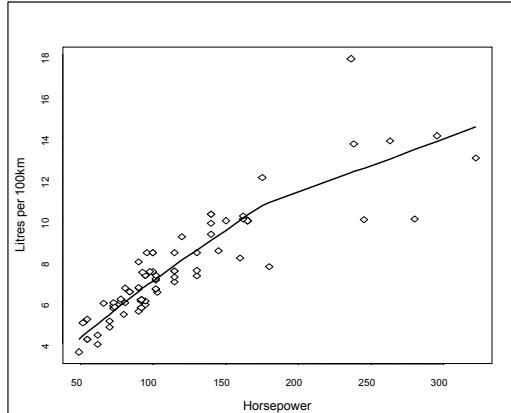
Lowess Scatterplot Smoothing

- We will discuss the methodology behind lowess smoothing when we learn about regression
- Why smooth?
 - Why do we plot data?
 - To describe a relationship
 - If data fail to conform to standard models do we give up?

2005/10/14

Jeff Lin, MD, PhD

104



2005/10/14

Jeff Lin, MD, PhD

105

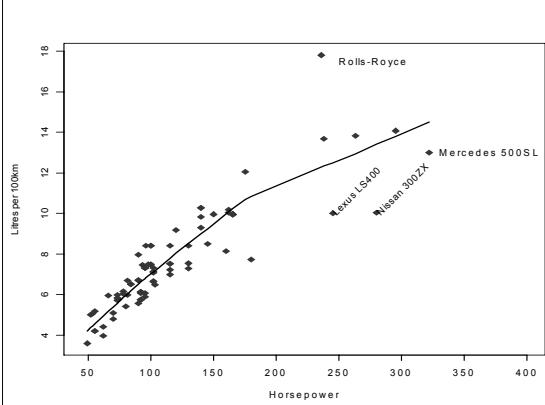
Useful things to add to a scatterplot

- A fitted line
- Identifiers on extreme or outlying points
- Correlation or R² (only good for linear data)
- Identify groups or clusters of data

2005/10/14

Jeff Lin, MD, PhD

106



2005/10/14

Jeff Lin, MD, PhD

107

Thanks !

2005/10/14

Jeff Lin, MD, PhD

108