# Methods of Applied Statistics I

## STA2101H F LEC9101

Week 5

October 12 2022



If you believe they're biased...
US Senate races, Democratic vote margin in pre-election polls, percentage points
September 14th 2022

● Average  ● Average accounting for polling bias*

*Subtracting the overestimation of Democratic vote margins in all statewide federal elections since 2016    Source: The Economist

The Economist

1. Upcoming events
2. Recap
3. Steps in analysis; types of studies
4. In the News

5. HW 4 3rd hour

6. Sections for Project
   - a description of the scientific problem of interest
   - how (and why) the data being analyzed was collected
   - preliminary description of the data (plots and tables)
   - models and analysis
   - summary for a statistician of the analysis and conclusions
   - non-technical summary for a non-statistician of the analysis and conclusions

- October 13 3.30-4.30 : DoSS Seminar Room 9014 (Hydro Building)
- Brenda Betancourt, U Chicago
- Microclustering for record linkage applications



**Brenda Betancourt**
**NORC – University of Chicago**

Brenda is currently a Senior Statistician at NORC at the University of Chicago. Before joining NORC, she was an Assistant Professor in the Department of Statistics at the University of Florida. She obtained her PhD in Statistics at the University of California, Santa Cruz and was a postdoctoral fellow at Duke University working on Bayesian models and algorithms for record linkage and Network analysis.

**Microclustering for record linkage applications**

# Upcoming

- October 15 12.00-1.00 Toronto Data Workshop; Room BL520, Bissell Building and online link    data_4_lyf

- April Wang, U Michigan

  "Reimagining Tools for Collaborative Data Science"

Box-Cox '64    Bickel Doksum '81    Box-Cox revisited
'82    " " "
, rebutted

- Model Selection: hierarchical principle, testing procedures ("$p < 0.05$")
- criterion-based procedures ($AIC$, $BIC$, $C_p$, $R_a^2$)
- regularization/penalization methods: Lasso

BIC $p \log n$    corr'd lm
AIC $2p$    $AIC_c$

- Model Building: plots, partial plots
- consideration of units and types of variation
- potential transformation of variables

$\lambda = 0$

$\lambda \neq 0 \left( \dfrac{y_i^\lambda - 1}{\lambda} = x_i^T \beta + \varepsilon_i \right)$

$\lambda \to 0 \quad \log y \to x_i$

HW3

- clarification of objectives: prediction and/or explanation
- criterion-based methods may be helpful for prediction
- automated methods rarely useful for explanation
- there may be several models consistent with the data

a) $\lambda = 1$

b) $\lambda = 0$

*) d) $\hat{\lambda} \approx 0.22$

$(25, 125)$

$(5, 3000)$

.2 [ ( 10 , 700 ) ]

- common objectives

- to avoid systematic error, that is distortion in the conclusions arising from sources that do not cancel out in the long run

- common objectives

- to avoid systematic error, that is distortion in the conclusions arising from sources that do not cancel out in the long run

- to reduce the non-systematic (random) error to a reasonable level by replication and other techniques

- common objectives

- to avoid systematic error, that is distortion in the conclusions arising from sources that do not cancel out in the long run

- to reduce the non-systematic (random) error to a reasonable level by replication and other techniques

- to estimate realistically the likely uncertainty in the final conclusions

- common objectives

  *bias*

- <u>to avoid systematic error,</u> that is distortion in the conclusions arising from sources that do not cancel out in the long run

$$\bar{y} \qquad \hat{\beta}_7$$

- to reduce the non-systematic (random) error to a reasonable level by replication and other techniques

  *variance*

- to estimate realistically the likely uncertainty in the final conclusions

- to ensure that the <u>scale of effort is appropriate</u>

- we concentrate largely on the careful analysis of individual studies
- in most situations synthesis of information from different investigations is needed

- we concentrate largely on the careful analysis of individual studies
- in most situations synthesis of information from different investigations is needed

- but even there the quality of individual studies remains important
- examples include overviews (such as the Cochrane reviews)

- we concentrate largely on the careful analysis of individual studies
- in most situations synthesis of information from different investigations is needed

- but even there the quality of individual studies remains important
- examples include overviews (such as the Cochrane reviews)

- in some areas new investigations can be set up and completed relatively quickly; design of individual studies may then be less important

*A/B testing in tech*

- formulation of a plan of analysis
- establish and document that proposed data are capable of addressing the research questions of concern

- formulation of a plan of analysis
- establish and document that proposed data are capable of addressing the research questions of concern

- main configurations of answers likely to be obtained should be set out
- level of detail depends on the context

- formulation of a plan of analysis
- establish and document that proposed data are capable of addressing the research questions of concern

- main configurations of answers likely to be obtained should be set out
- level of detail depends on the context

- even if pre-specified methods must be used, it is crucial not to limit analysis
- planned analysis may be technically inappropriate
- more controversially, data may suggest new research questions or replacement of objectives

$$\frac{\bar{y}_a - \bar{y}_b}{\sim N} \qquad \bar{y}_{max} - \bar{y}_{min} \qquad \text{"p-hacking"}$$

- formulation of a plan of analysis
- establish and document that proposed data are capable of addressing the research questions of concern

- main configurations of answers likely to be obtained should be set out
- level of detail depends on the context

- even if pre-specified methods must be used, it is crucial not to limit analysis
- planned analysis may be technically inappropriate
- more controversially, data may suggest new research questions or replacement of objectives

- latter will require confirmatory studies

# Unit of study and analysis

- smallest subdivision of experimental material that may be assigned to a treatment

  context: Expt

- Example: RCT – unit may be a patient, or a patient-month (in crossover trial)
- Example: public health intervention – unit is often a community/school/...

- split plot experiments have two classes of units of study and analysis
- in investigations that are not randomized, it may be helpful to consider what the primary unit of analysis would have been, had a randomized experiment been feasible

$$E(\varepsilon y) = 0 \quad \text{but } \hat{\varepsilon} \neq \varepsilon$$

$$E[\hat{\varepsilon}\hat{y}] = 0 \qquad \hat{\varepsilon} \perp \hat{y} \qquad \text{not} \perp y \qquad E(\hat{\varepsilon}y)$$

$$\text{uncorr'd}$$

$$y - \hat{y} \quad \cancel{\text{ind't}} \text{ of } \hat{y} \qquad \overset{?}{=} (I - H)\sigma^2$$

$i = 1, \ldots, n$

- smallest subdivision of experimental material that may be assigned to a treatment
  context: Expt

- Example: RCT – unit may be a patient, or a patient-month (in crossover trial)
- Example: public health intervention – unit is often a community/school/...

$y$

$x$

- split plot experiments have two classes of units of study and analysis
- in investigations that are not randomized, it may be helpful to consider what the primary unit of analysis would have been, had a randomized experiment been feasible

- the unit of analysis may not be the unit of interpretation – ecological bias
  systematic difference between impact of $x$ at different levels of aggregation

in est $\hat{\beta}$

- CD: Illustration – "For country- or region-based mortality data, countries of regions respectively may reasonably constitute the units of analysis with which to assess the relationship of the data to dietary and other features

- "Yet the objective is interpretation at an individual person level

- "The situation may be eased if supplementary data on explanatory variables are available at the individual level, because this may clarify the connection of between-unit and within-unit variation"

- CD: Illustration – "For country- or region-based mortality data, countries of regions respectively may reasonably constitute the units of analysis with which to assess the relationship of the data to dietary and other features

- "Yet the objective is interpretation at an individual person level

- "The situation may be eased if supplementary data on explanatory variables are available at the individual level, because this may clarify the connection of between-unit and within-unit variation"

- feelings of well-being are associated with socio-economic status link

- the strength of the association is larger in developed nations
  than in developing nations

GDP

- conventional explanation: in nations with a high level of economic development
  perhaps higher SES carries some intrinsic value

- this paper: in nations with a high level of religiosity, the strength of the association
  between SES and well-being is weaker

- religiosity could attenuate the link between well being and SES   Economist, Sep 25 2021

Distribution of national economic development, national religiosity, and estimated means of the cross-level interactions (model 3).

**Gallup World Poll**

2005-2017

1,567,204 participants in 156 nations

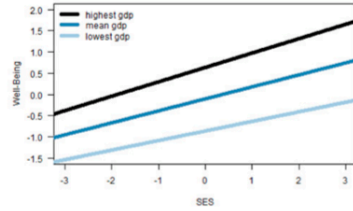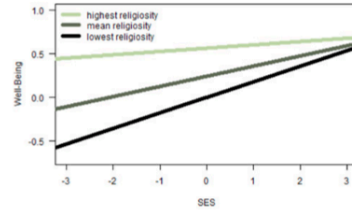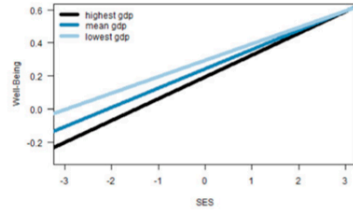**Gosling-Potter Internet Personality Project**

1999-2015

1,493,207 participants in 85 nations

**World Values Survey**

1981-2016
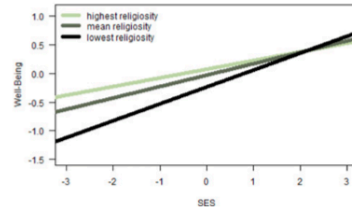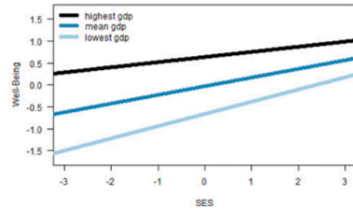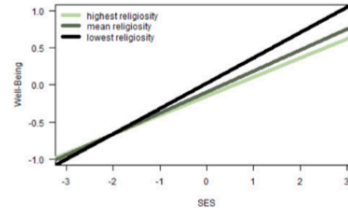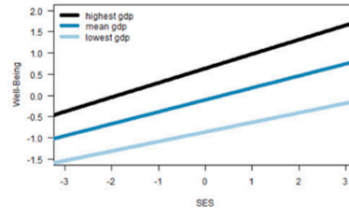
274,393 participants in 92 nations

**Fig. 2.** Distribution of national economic development, national religiosity, and estimated means of the cross-level interactions (model 3). The top row depicts mean national economic development and mean national religiosity worldwide (z-standardized and averaged across data sets). Lighter colors represent lower values, darker colors higher values. The three bottom rows depict estimated marginal means of the cross-level interactions when both national moderators were included in the model (i.e., model 3). Depicted are the moderating effects of national economic development and national religiosity on the association between SES and well-being in all three data sets.

**Fig. 2.** Distribution of national economic development, national religiosity, and estimated means of the cross-level interactions (model 3). The top row depicts mean national economic development and mean national religiosity worldwide (z-standardized and averaged across data sets). Lighter colors represent lower values, darker colors higher values. The three bottom rows depict estimated marginal means of the cross-level interactions when both national moderators were included in the model (i.e., model 3). Depicted are the moderating effects of national economic development and national religiosity on the association between SES and well-being in all three data sets.

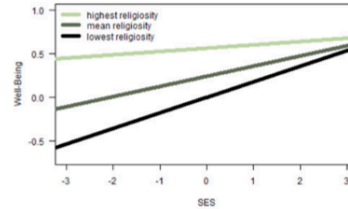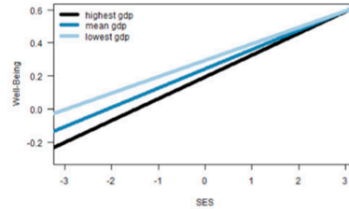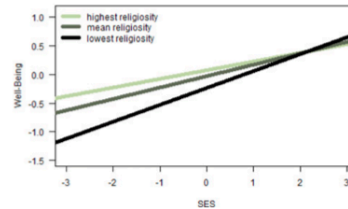Depicted are the moderating effects of national economic development and national religiosity on the association between SES and well-being in all three data sets.

- $\mathrm{E}(y_i \mid x_i, z_i) = \beta_0 + \beta_1 x_i + \beta_2 z_i + \beta_3 x_i z_i$

  `WellB ~ ses + relig + ses:relig`

- $y_i$: Well-being, $x_i$: Socio-economic status, $z_i$: Religiosity – a factor variable

Interaction

- $\mathrm{E}(y_i \mid x_i, z_i) = \beta_0 + \beta_1 x_i + \beta_2 z_i + \beta_3 x_i z_{i1}$    `WellB` $\sim$ `ses` + `relig` + `ses:relig`

  $\pi \cdot \beta_3 z_{i2} + \beta_5 x_i z_{i2}$

- $y_i$: Well-being, $x_i$: Socio-economic status, $z_i$: Religiosity – a factor variable

  $$z_{i1} \quad z_{i2}$$
  $$SES \quad relig\,med \quad relig\,hi$$

- $z_i = ($ "low", "medium", "high" $)$      `model.matrix`    $x_i \quad 0 \quad 0$

  $z_i = 0$     $\mathrm{E}(y_i \mid x_i, z_i = "low") \;=\; \boxed{\beta_0 + \beta_1 x_i}$    $\vdots \quad \longrightarrow 1 \quad 0$

  $z_i = $     $\mathrm{E}(y_i \mid x_i, z_i = "med") \;=\; \beta_0 + \beta_2 + (\beta_1 + \beta_4)x_i$    $0 \quad 1$

  mediator of $\{$ x–y assoc$^n$    $\mathrm{E}(y_i \mid x_i, z_i = "hi") \;=\; \beta_0 + \beta_3 + (\beta_1 + \beta_5)x_i$   $x_n$

- as usual, the paper's a bit more complicated
- some data collected on people, and some on countries – multi-level model
- "Following standard practice, we averaged person-level religiosity within each nation "
- there's another covariate – GDP
- "Following a standard economic method, we log-transformed the GDP data"

regression lm w̄ 1 cont^s & 1 factor

- a factor variable is treated as categorical
- a non-factor variable is treated as continuous
- it depends on the application which is preferred

- a factor variable is treated as categorical
- a non-factor variable is treated as continuous
- it depends on the application which is preferred

- a linear model with one factor and one continuous variable might be written as, for example:

$$y_{ij} = \mu + \alpha_j + \beta x_{ij} + \epsilon_{ij}, \quad j = 1, \ldots, J; \quad i = 1, \ldots m$$

- linear in $x$, but arbitrary changes in $\mathbb{E}(y)$ by category (here indexed by $j$)

$J$ - level factor

$\beta$

$\alpha_j$
$\alpha_{j-}$
$\alpha_{j''}$

$E(y_{ij}) = \mu + \alpha_j + \beta_j x_{ij} + \epsilon_{ij}$

$\& \alpha_c$

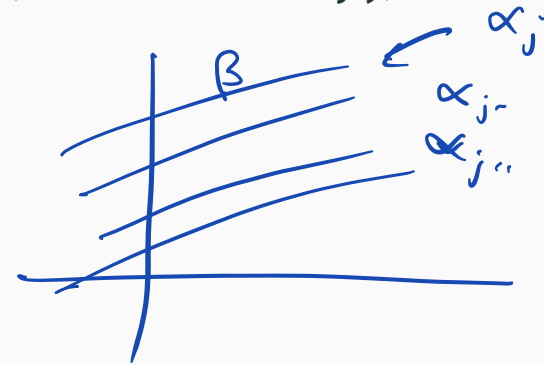$E(y|a) \begin{vmatrix} \alpha_1 \\ \alpha_2 \end{vmatrix}$

$, \alpha_3$

- a factor variable is treated as categorical
- a non-factor variable is treated as continuous
- it depends on the application which is preferred

*age level 1*
*group 2*
*⋮*
*6*

$\alpha_j$

- a linear model with one factor and one continuous variable might be written as, for example:

$$E\ y_{ij} = \mu + \alpha_j + \beta x_{ij} + \epsilon_{ij}, \quad j = 1, \ldots, J; \quad i = 1, \ldots m$$

- linear in *x*, but arbitrary changes in $\mathbb{E}(y)$ by category (here indexed by *j*)
- R doesn't distinguish this at the modelling phase:

```
lm(response ~ variable1 + variable2, data = ...)
```

- but uses metadata in the data frame to accommodate factors
- `is.factor(variable)` and `newvar <- as.factor(oldvar)` are helpful

*Tukey's*
*Least*
*Significc*
*Difference*

*Check   fruitfly.html   &*

12



**Figure S1.** Mixed Effects Mediated Moderation Model (model 4). Statistical model to calculate the portion of the cross-level interaction

- "treatment" is not assigned to units, only observed
- any observed effect of treatment cannot be assumed to be causal

"correlation is not causation"

- "treatment" is not assigned to units, only observed
- any observed effect of treatment cannot be assumed to be causal

"correlation is not causation"

- we can try to assess the effect by controlling for other variables that may also influence the response
- but we cannot control for unmeasured variables

- "treatment" is not assigned to units, only observed
- any observed effect of treatment cannot be assumed to be causal

  "correlation is not causation"

- we can try to assess the effect by controlling for other variables that may also influence the response
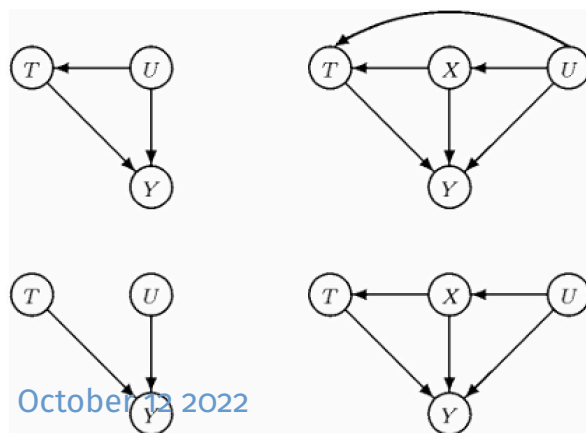- but we cannot control for unmeasured variables



418

*9 · Designed Experiments*

**Figure 9.1** Directed acyclic graphs showing consequences of randomization. An arrow from $T$ to $Y$ indicates dependence of $Y$ on $T$, and so forth. In general both response $Y$ and treatment $T$ may depend on properties $U$ of units (upper left). Randomization (lower left) makes treatments and units independent, so any observed dependence of $Y$ on $T$ cannot be ascribed to joint dependence on $U$. The upper right graph shows the general dependence of $Y$, $T$, and covariates $X$ on $U$.

$E(y_i) = \beta_0 + \beta_1 x_i + \beta_2 t_i$

tmt

hidden

measured confounders

obs'l study

$\beta \neq 0$

response

obs'l

unmeas'd

obs'l

unmeasured confounder

randomized expt

obs'd

od exp't

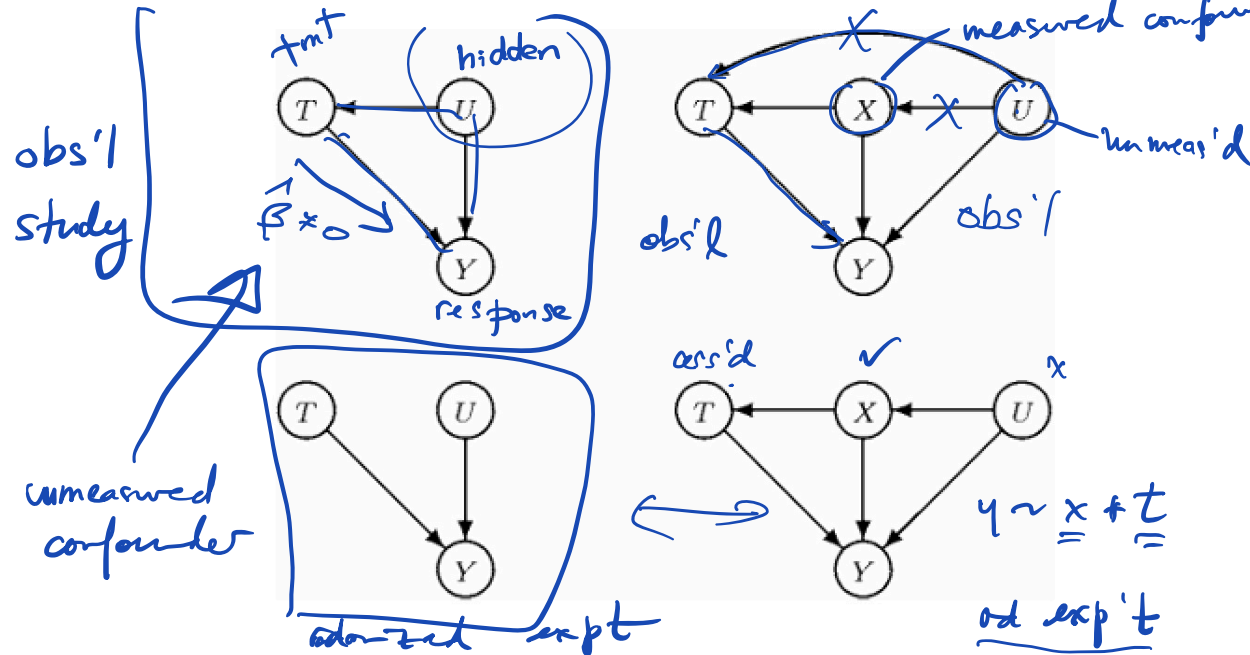$y \sim x + t$

**Figure 9.1**   Directed acyclic graphs showing consequences of randomization. An arrow from $T$ to $Y$ indicates dependence of $Y$ on $T$, and so forth. In general both response $Y$ and treatment $T$ may depend on properties $U$ of units (upper left). Randomization (lower left) makes treatments and units independent, so any observed dependence of $Y$ on $T$ cannot be ascribed to joint dependence on $U$. The upper right graph shows the general dependence of $Y$, $T$, and covariates $X$ on $U$. Randomization makes $T$ and $U$ independent, conditional on $X$ (lower right), so any influence of $U$ on $T$ is mediated through $X$, for which adjustment is possible in principle. Thus having adjusted for $X$, dependence of $Y$ on $T$ cannot be due to $U$.

the control group. The response is to be the blood pressure of an individual measured a fixed time after the drug has first been administered. We calculate the average changes for the treated and control groups, $\bar{y}_1$ and $\bar{y}_0$, observe that $\bar{y}_1 - \bar{y}_0$ is significantly less than zero, and declare that the drug plays an effect in reducing blood pressure. Is this headline news? No!

A key difficulty is that the procedure does not avoid biased allocation of treatments to units. For example, if the control group mostly consisted of those patients with

*from obs'l study*

*Bradford-Hill criteria*

- strength of the association
- consistency of the association  *across many exp'ts*
- specificity of the proposed causal factor
- potential cause occurs before its effect (temporality)
- dose-response relationship
- a subject-matter theory exists *⊛*

- "natural experiments" e.g. minimum wage

# Types of observational studies

- secondary analysis of data collected for another purpose

# Types of observational studies

- secondary analysis of data collected for another purpose

- estimation of some feature of a defined population

  could in principle be found exactly

- tracking across time of such features

# Types of observational studies

- secondary analysis of data collected for another purpose

- estimation of some feature of a defined population

  could in principle be found exactly

- tracking across time of such features

- study of a relationship between features, where individuals may be examined
  - at a single time point
  - at several time points for different individuals
  - at different time points for the same individual

# Types of observational studies

- secondary analysis of data collected for another purpose

- estimation of some feature of a defined population

  could in principle be found exactly

- tracking across time of such features

- study of a relationship between features, where individuals may be examined
    - at a single time point
    - at several time points for different individuals
    - at different time points for the same individual

- census

- meta-analysis: statistical assessment of a collection of studies on the same topic

- $y_i \sim Bin(6, p_i), \quad i = 1, \ldots, 23$

# Modelling numbers/proportions of events

- $y_i \sim Bin(6, p_i), \quad i = 1, \ldots, 23$

- in general: $n_i$ trials, $y_i$ successes, probability of success $p_i$

- for regression: associated covariate vector $x_i$, e.g. temperature

- SM uses $m_i$ and $r_i$ instead of $n_i$ and $y_i$

- each $y_i$ could in principle be the sum of $n_i$ independent Bernoulli trials

- each of the $n_i$ trials having the same probability $p_i$

- with the same covariate vector $x_i$ \qquad FELM 'covariate classes', p.26

```
> library(faraway); data(orings)
> logitmod <- glm(cbind(damage,6-damage) ~ temp, family = binomial, data = orings)
> summary(logitmod)
Call:
glm(formula = cbind(damage, 6 - damage) ~ temp, family = binomial,
    data = orings)
...
Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept) 11.66299    3.29626   3.538 0.000403 ***
temp        -0.21623    0.05318  -4.066 4.78e-05 ***
---

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 38.898  on 22  degrees of freedom
Residual deviance: 16.912  on 21  degrees of freedom
```

```
> library(SMPracticals) # this is for datasets in
                          #Statistical Models by Davison
> data(shuttle) # same example, different name
> shuttle2 <- data.frame(as.matrix(shuttle)) # this is a kludge to avoid
                                              #an error with head(shuttle)

> head(shuttle2)
  m r temperature pressure
1 6 0          66       50
2 6 1          70       50
3 6 0          69       50
4 6 0          68       50
5 6 0          67       50
6 6 0          72       50
> par(mfrow=c(2,2)) # puts 4 plots on a page



> with(orings,plot(temp,damage,main="Faraway",xlim=c(31,80)))
> with(shuttle,plot(temperature,r,main="Davison",xlim=c(31,80),
+ ylim=c(0,5)))
```