

# F tests

Math 463, Spring 2017, University of Oregon

David A. Levin

University of Oregon

May 8, 2017

## Some housekeeping

- Next HW will be posted by Friday (due following Fri).
- After analysis of variance, will discuss path models (linked regression equations).  
Read Chapter 6 in Freedman.

- Model: random vectors  $\mathbf{Y}, \mathbf{x}_1, \dots, \mathbf{x}_k, \boldsymbol{\varepsilon}$ , with

$$\mathbf{Y} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_k] \boldsymbol{\beta} + \boldsymbol{\varepsilon}.$$

Usual distribution of tests, etc. is *conditional* on the observed values of  $\mathbf{x}_1, \dots, \mathbf{x}_k$ .

$$\boldsymbol{\theta} := \mathbb{E}[\mathbf{Y} \mid \mathbf{X}] = \mathbf{X} \boldsymbol{\beta} \in \mathcal{L}(\mathbf{x}_1, \dots, \mathbf{x}_k)$$

- Suppose we have a sequence of subspaces  $V_0 \subset V_1 \subset \dots \subset V_r = V$ . Let

$$W_k = V_k \cap V_{k-1}^\perp,$$

so that

$$\mathbb{R}^n = V_0 \oplus \underbrace{W_1 \oplus \dots \oplus W_r}_{V_0^\perp} \oplus V^\perp$$

- Thus the projection onto  $V_0^\perp$  can be resolved into orthogonal pieces on each of  $W_1, \dots, W_r$  and  $V^\perp$ : Letting  $\hat{\mathbf{Y}}_j = \Pi_{V_j} \mathbf{Y}$  and recall  $\Pi_{W_j} = \Pi_{V_j} - \Pi_{V_{j-1}}$ ,

$$\begin{aligned} \|\mathbf{Y} - \hat{\mathbf{Y}}_0\|^2 &= \|\hat{\mathbf{Y}}_1 - \hat{\mathbf{Y}}_0\|^2 + \|\hat{\mathbf{Y}}_2 - \hat{\mathbf{Y}}_1\|^2 + \dots + \|\hat{\mathbf{Y}}_r - \hat{\mathbf{Y}}_{r-1}\|^2 + \|\mathbf{Y} - \hat{\mathbf{Y}}_r\|^2 \\ &= \text{SS}_{\text{Reg}}(V_1 \mid V_0) + \dots + \text{SS}_{\text{Reg}}(V_r \mid V_{r-1}) + \text{SS}_{\text{Res}} \\ &= (\|\mathbf{Y} - \hat{\mathbf{Y}}_0\|^2 - \|\mathbf{Y} - \hat{\mathbf{Y}}_1\|^2) + \dots + (\|\mathbf{Y} - \hat{\mathbf{Y}}_{r-1}\|^2 - \|\mathbf{Y} - \hat{\mathbf{Y}}_r\|^2) \\ &\quad + \|\mathbf{Y} - \hat{\mathbf{Y}}_r\|^2 \\ &= \Delta \text{SS}_{\text{Res}}(V_1 \mid V_0) + \dots + \Delta \text{SS}_{\text{Res}}(V_r \mid V_{r-1}) + \text{SS}_{\text{Res}} \end{aligned}$$

- Each of  $\text{SS}_{\text{Reg}}(V_i \mid V_{i-1})$  is chi-squared with noncentrality parameter  $\|\Pi_{W_i}\boldsymbol{\theta}\|^2/\sigma^2$ , independent of  $\text{SS}_{\text{Res}}$ . Thus the hypothesis test of

$$H_0 : \boldsymbol{\theta} \in V_{i-1} \quad \text{vs.} \quad H_1 : \boldsymbol{\theta} \in V_i \setminus V_{i-1}$$

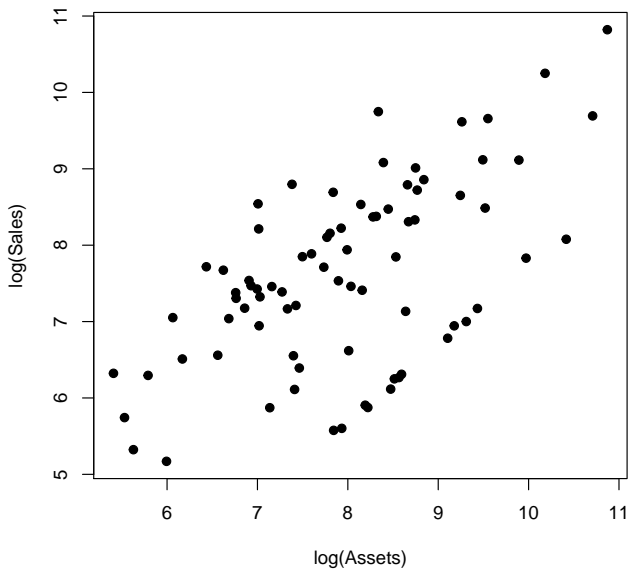
is based on

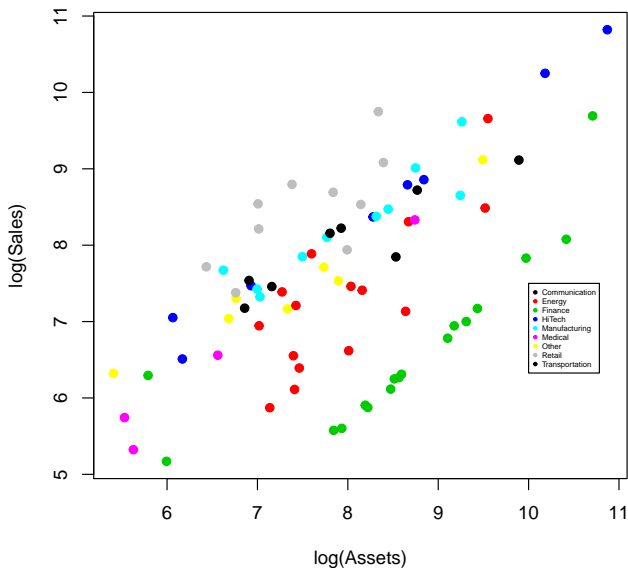
$$F = \frac{\text{SS}_{\text{Reg}}(V_i \mid V_{i-1})/\dim(W_i)}{\text{SS}_{\text{Res}}/(n - \dim(V))}.$$

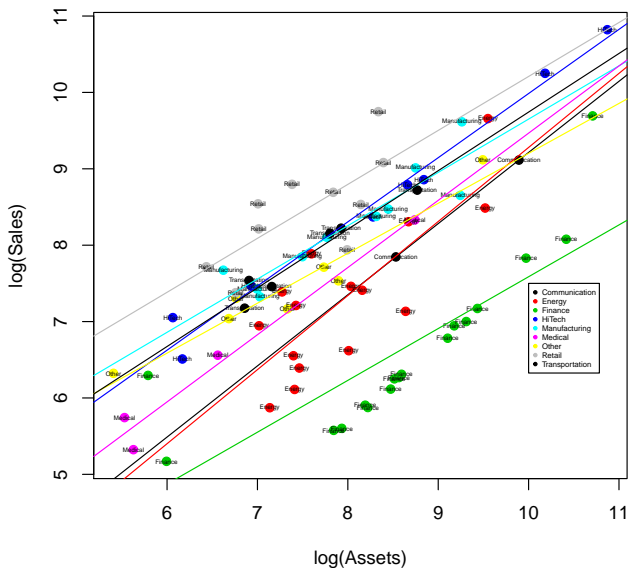
- An *analysis of variance* table gives each of these tests for a sequence of nested models.

## Factors, dummy variables, ANOVA

```
> comp = read.table("http://pages.uoregon.edu/dlevin/DATA/companies.txt",
+                   header=T, sep="\t")
> plot(log(Sales)~log(Assets), data=comp, col=sector, pch=19)
> legend(10,7.5,pch=19,col=1:9,legend=levels(comp$sector),cex=0.4)
> comp.lm = lm(log(Sales)~log(Assets)+sector+log(Assets):sector,
+             data=comp)
> text(log(comp$Sales)~log(comp$Assets),labels=comp$sector,cex=0.3)
> betahat = comp.lm$coef
> basebeta = comp.lm$coef[1:2]
> abline(basebeta,col=1)
> for(i in 0:8){
+   abline(basebeta+betahat[c(2+i,10+i)], col=1+i)
+ }
```









Let for  $j = 2, \dots, 9$

$$\delta_{i,j} = \begin{cases} 1 & \text{if } i \text{ data point in sector } j \\ 0 & \text{else} \end{cases}$$

Consider

$$\mathbb{E}(Y_i \mid x_i, \delta_{i,\cdot}) = \beta_0 + \beta_1 x_i + \sum_{j=2}^9 \beta_j \delta_{i,j} + \sum_{j=2}^9 \beta_{8+j} \delta_{i,j} x_i$$

If the  $i$ -th data point belongs to sector  $k = 2, \dots, 8$ , then

$$\mathbb{E}(Y_i \mid x_i, \delta_{i,\cdot}) = (\beta_0 + \beta_k) + (\beta_1 + \beta_{8+k})x_i$$

In other words,  $\beta_0 + \beta_k$  is the intercept for the sector- $k$  data, and  $(\beta_1 + \beta_{8+k})$  is the slope for the sector- $k$  data, where  $k = 2, \dots, 8$ .

```
> model.matrix(comp.lm)[22:24,]
```

	(Intercept)	log(Assets)	sectorEnergy	sectorFinance	sectorHiTech
22	1	8.393442	0	0	0
23	1	8.841304	0	0	1
24	1	6.759255	0	0	0

	sectorManufacturing	sectorMedical	sectorOther	sectorRetail
22	0	0	0	1
23	0	0	0	0
24	0	0	0	1

	sectorTransportation	log(Assets):sectorEnergy	log(Assets):sectorFinance
22	0	0	0
23	0	0	0
24	0	0	0

	log(Assets):sectorHiTech	log(Assets):sectorManufacturing
22	0.000000	0
23	8.841304	0
24	0.000000	0

	log(Assets):sectorMedical	log(Assets):sectorOther	log(Assets):sectorRetail
22	0	0	8.393442
23	0	0	0.000000
24	0	0	6.759255

	log(Assets):sectorTransportation
22	0
23	0
24	0

## Confidence intervals

Suppose we want to make a confidence interval for

$$E[Y \mid \log(\text{Assets}) = 7, \text{sector} = \text{"Energy"}]$$

Need to create a data frame with new values.

```
> newcomp = data.frame(Assets=1096, sector = "Energy")  
> predict(comp.lm, newdata=newcomp, interval="confidence")
```

	fit	lwr	upr
1	6.370161	5.958412	6.781909

Want to test

$$H_0 : \beta_2 = \cdots = \beta_9 = 0$$

and (separately)

$$H_0 : \beta_{10} = \cdots = \beta_{17} = 0$$

We have nested models  $V_0 \subset V_1 \subset V_2 \subset V_3$

- $V_0 : \theta \in \mathcal{L}(\mathbf{1})$ .  $[H_0 : \beta_1 = \cdots = \beta_{17} = 0, H_1 : \beta_1 \neq 0, \beta_2 = \cdots = \beta_{17} = 0]$
- $V_1 : \theta \in \mathcal{L}(\mathbf{1}, \mathbf{x})$ ,  $\Delta \dim = 1$ .  $[H_0 : \beta_2 = \cdots = \beta_{17} = 0, H_1 : \exists 1 \leq i \leq 9 \text{ s.t. } \beta_i \neq 0, \beta_{10} = \cdots = \beta_{17} = 0]$
- $V_2 : \theta \in \mathcal{L}(\mathbf{1}, \mathbf{x}, \boldsymbol{\delta}_2, \dots, \boldsymbol{\delta}_8)$ ,  $\Delta \dim = 8$ .  $[H_0 : \beta_{10} = \cdots = \beta_{17} = 0]$
- $V_3 : \theta \in \mathcal{L}(\mathbf{1}, \mathbf{x}, \boldsymbol{\delta}_2, \dots, \boldsymbol{\delta}_8, \mathbf{x}\boldsymbol{\delta}_1, \dots, \mathbf{x}\boldsymbol{\delta}_8)$ ,  $\Delta \dim = 8$ .

```
> library(xtable)
> xtable(anova(comp.lm))
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
log(Assets)	1	38.32	38.32	148.84	0.0000
sector	8	58.00	7.25	28.16	0.0000
log(Assets):sector	8	1.00	0.13	0.49	0.8611
Residuals	61	15.70	0.26		

Each line reports the distance between fitted models.

The statistic is

$$\frac{\text{SS}_{\text{Reg}}(V_i | V_{i-1})/\Delta\text{df}}{\text{SS}_{\text{Res}}/\text{df}} \sim F(\Delta\text{df}, \text{df}, \gamma),$$

where  $\gamma_i = \|\Pi_{V_i \cap V_{i-1}^\perp} \theta\|^2 / \sigma^2$ , recalling that  $\theta = \mathbf{X}\beta$ . If  $\theta \in V_{i-1}$ , then  $\gamma_i = 0$ . If  $\theta \in V_i$ , then  $\gamma_i$  depends only on the coefficients and covariates which are in  $V_i$ .

```
> comp2.lm=lm(terms(log(Sales)~log(Assets)+log(Assets):sector+sector,keep.o
> xtable(anova(comp2.lm))
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
log(Assets)	1	38.32	38.32	148.84	0.0000
log(Assets):sector	8	57.54	7.19	27.94	0.0000
sector	8	1.46	0.18	0.71	0.6812
Residuals	61	15.70	0.26		

Note that changes the order of the subspaces does change the  $SS_{\text{Reg}}$ , which is *conditional* on the subspace below. That is, the line marked “sector” in the above table corresponds to

$$SS_{\text{Reg}}(\delta_2, \dots, \delta_8 \mid 1, x, \delta_2 x, \dots, \delta_8 x)$$

In the first table, the line labels sector corresponds to

$$SS_{\text{Reg}}(\delta_2, \dots, \delta_8 \mid 1, x)$$

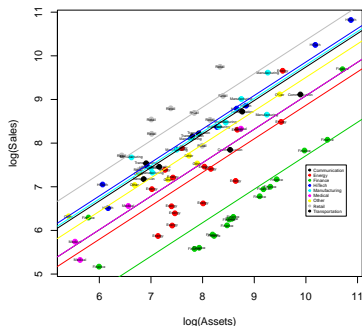
## How to test a linear constant $T\beta = 0$

- $H_0 : T\beta = 0$ .
- Suppose that  $T$  is a  $r \times (k + 1)$  matrix of rank  $r$ .  $H_0$  specifies that  $\beta$  is orthogonal to the row space of  $T$ , so  $\beta$  is constrained to lie in  $k + 1 - r$  dimensional linear space.
- How do we translate this into a hypothesis on  $\theta = \mathbb{E}(Y | X)$ ?
- Two points of view:
  - ▶ Theoretical description of  $F$  test.
  - ▶ Since  $SS_{\text{Reg}}$  is the difference between  $SS_{\text{Res}}$  for large and small models, it is enough to fit both models. How to fit model subject to the constraint in  $H_0$ ?

```

> plot(log(Sales)~log(Assets), data=comp, col=sector, pch=19)
> legend(10,7.5,pch=19,col=1:9,legend=levels(comp$sector),cex=0.4)
> comp.lm.b = lm(log(Sales)~log(Assets)+sector, data=comp)
> text(log(comp$Sales)~log(comp$Assets),labels=comp$sector,cex=0.3)
> betahat = comp.lm.b$coef
> basebeta = comp.lm.b$coef[1:2]
> abline(basebeta,col=1)
> for(i in 0:7){
+ abline(basebeta+c(betahat[3+i],0), col=2+i)
+ }

```





## $F$ test with linear constraint

- Note that  $\beta = (X'X)^{-1}X'\theta$ , so  $T\beta = 0$  iff  $TM^{-1}X'\theta = 0$ .
- This is equivalent to saying that

$$\theta \perp V_1 := \text{row space}(TM^{-1}X') = \text{col space}(\underbrace{XM^{-1}T'}_B)$$

- Since  $SS_{\text{Reg}}$  is the squared length of the projection onto  $V_1$ , we can write the numerator of the  $F$ -statistic as

$$SS_{\text{Reg}} = \|P_{V_1} Y\|^2 = Y'B(B'B)^{-1}B'Y = \hat{\beta}'T'(TM^{-1}T)^{-1}T\hat{\beta}$$

- Note that  $\beta$  is constrained to fall in the null space of  $T$  which has dimension  $k+1-r$ . Thus if  $c_j, j=1, \dots, k+1-r$  is a basis for the null space, then there is  $\gamma \in \mathbb{R}^{k+1-r}$

$$\beta = \gamma_1 c_1 + \dots + \gamma_{k+1-r} c_{k+1-r} = C\gamma$$

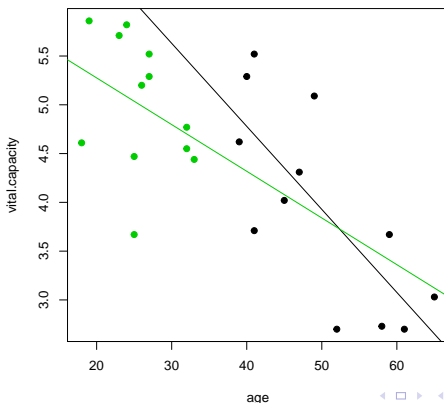
- Fitting the model subject to the constraint is equivalent to the model

$$\mathbb{E}(Y | X) = X\beta = XC\gamma.$$

Thus, fitted value and RSS can be found via OLS using  $XC$  as model matrix.

## Example

```
> library(ISwR)
> plot(vital.capacity~age, data = vitcap, col=group, pch=19)
> vc.lm = lm(vital.capacity~age + age*as.factor(group), data=vitcap)
> vc.coef = vc.lm$coefficients
> abline(vc.coef[1:2])
> abline(vc.coef[1:2]+vc.coef[3:4], col=3)
```



```
> xtable(anova(vc.lm))
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
age	1	12.48	12.48	25.90	0.0001
as.factor(group)	1	0.96	0.96	2.00	0.1732
age:as.factor(group)	1	0.27	0.27	0.57	0.4603
Residuals	20	9.64	0.48		

Let us instead parameterize the model as

$$\mathbb{E}(Y_i | \delta, x) = \beta_0 \delta_i + \eta_0(1 - \delta_i) + \beta_1 \delta_i x_i + \eta_1(1 - \delta_i) x_i.$$

```
> delta = as.numeric(vitcap$group==1)
> vc.lm.2=lm(vital.capacity~delta+I(1-delta)+I(delta*age)
+           +I((1-delta)*age)-1,data=vitcap)
> xtable(summary(vc.lm.2))
```

	Estimate	Std. Error	t value	Pr(> t )
delta	8.1834	1.1608	7.05	0.0000
I(1 - delta)	6.2327	1.1533	5.40	0.0000
I(delta * age)	-0.0851	0.0230	-3.70	0.0014
I((1 - delta) * age)	-0.0479	0.0438	-1.09	0.2878

Note that “a difference in significance is not the same as a significant difference”!

- Want to test  $H_0 : \beta_1 = \eta_1, \beta_0 = \eta_0$ .

- 

$$\begin{bmatrix} 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & -1 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \eta_0 \\ \beta_1 \\ \eta_1 \end{bmatrix} = 0$$

- Want to determine the null space of  $\mathbf{T}$ , although it should be clear that the matrix  $\mathbf{C}$  should add the first two columns of  $\mathbf{X}$  together and add the last two columns of  $\mathbf{X}$  together. The SVD can be used to determine a basis of the null space: if  $\mathbf{T} = \mathbf{V}\mathbf{\Sigma}\mathbf{W}'$ , then the last two columns of  $\mathbf{W}$  are a basis for the null space of  $\mathbf{T}$ .

```
> T=matrix(c(1,-1,0,0,0,0,1,-1),nrow=2,byrow=T)
> print(svd(T,nv=4)$v)
```

```
      [,1]      [,2] [,3] [,4]
[1,]  0.7071068  0.0000000 -0.5  0.5
[2,] -0.7071068  0.0000000 -0.5  0.5
[3,]  0.0000000  0.7071068  0.5  0.5
[4,]  0.0000000 -0.7071068  0.5  0.5
```

The two columns

$$\mathbf{v}_1 = \frac{1}{2} \begin{bmatrix} -1 \\ -1 \\ 1 \\ 1 \end{bmatrix}, \quad \mathbf{v}_2 = \frac{1}{2} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

form a basis for the null space. Simpler to rotate the basis and use  $\mathbf{v}_1 + \mathbf{v}_2$  and  $\mathbf{v}_1 - \mathbf{v}_2$ , i.e.

$$\mathbf{w}_1 = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} \quad \mathbf{w}_2 = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\mathbf{X} \mathbf{C} = \mathbf{X} \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} (\mathbf{x}_1 + \mathbf{x}_2) & (\mathbf{x}_3 + \mathbf{x}_r) \end{bmatrix}$$

- So the end result, to fit the model, regress on  $\mathbf{x}_1 + \mathbf{x}_2$  and  $\mathbf{x}_3 + \mathbf{x}_4$ .
- ```
> vc.lm.3=lm(vital.capacity~I(delta+I(1-delta))+I(I(delta*age)
+
+I((1-delta)*age))-1,data=vitcap)
> xtable(anova(vc.lm.3,vc.lm.2))
```

|   | Res.Df | RSS   | Df | Sum of Sq | F    | Pr(>F) |
|---|--------|-------|----|-----------|------|--------|
| 1 | 22     | 10.87 |    |           |      |        |
| 2 | 20     | 9.64  | 2  | 1.23      | 1.28 | 0.2996 |

>

```
> X = model.matrix(vc.lm.2)
> Y = vitcap$vital.capacity
> Minv = summary(vc.lm.2)$cov.unscaled
> B = X%%Minv%%t(T)
> P = B%%solve(t(B)%%B)%%t(B)
> SS = t(Y)%%P%%Y
> print(SS)
```

[,1]

[1,] 1.234659

Suppose that we have a single quantitative variable  $X_1$  and a qualitative variable  $X_2$  taking on two possible values.

```
> inffile = url(  
+   "http://pages.uoregon.edu/dlevin/DATA/infmort.txt")  
> infmort = read.csv(inffile)  
> infmort[1:10,]
```

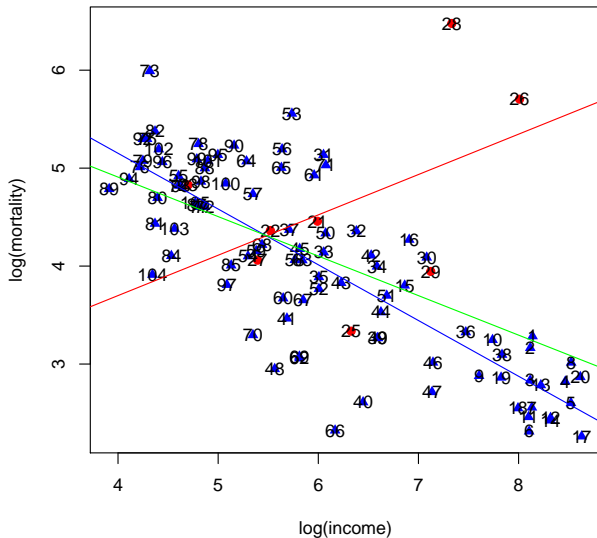
|    | region   | income | mortality | oil            |
|----|----------|--------|-----------|----------------|
| 1  | Asia     | 3426   | 26.7      | no oil exports |
| 2  | Europe   | 3350   | 23.7      | no oil exports |
| 3  | Europe   | 3346   | 17.0      | no oil exports |
| 4  | Americas | 4751   | 16.8      | no oil exports |
| 5  | Europe   | 5029   | 13.5      | no oil exports |
| 6  | Europe   | 3312   | 10.1      | no oil exports |
| 7  | Europe   | 3403   | 12.9      | no oil exports |
| 8  | Europe   | 5040   | 20.4      | no oil exports |
| 9  | Europe   | 2009   | 17.8      | no oil exports |
| 10 | Europe   | 2298   | 25.7      | no oil exports |



```

> oilp = 1:length(infmort$oil)
> oilp[infmort$oil=="no oil exports"]=17
> oilp[infmort$oil!="no oil exports"]=19
> oilc=1:length(infmort$oil)
> oilc[infmort$oil=="no oil exports"]="blue"
> oilc[infmort$oil!="no oil exports"]="red"
> plot(log(mortality)~log(income), data=infmort,pch=oilp,col=oilc)
> attach(infmort)
> text(log(mortality)~log(income),
+       labels=as.character(row.names(infmort)))
> f1 = lm(log(mortality)~log(income),
+         data=subset(infmort,oil=="no oil exports"))
> f2 = lm(log(mortality)~log(income),
+         data=subset(infmort,oil!="no oil exports"))
> infmort2=subset(infmort,
+                 row.names(infmort)!="26"&row.names(infmort)!="28")
> abline(f1,col="blue")
> abline(f2,col="red")
> f3 = lm(log(mortality)~log(income),
+         data=subset(infmort2,oil!="no oil exports"))
> abline(f3,col="green")

```



```

> library(xtable)
> infmort2=na.omit(infmort2)
> f = lm(log(mortality)~log(income)*oil, data=infmort2)
> anovatab = anova(f)
> xtable(anovatab,digits=2,caption="Anova Table",label="anova")

```

|                 | Df    | Sum Sq | Mean Sq | F value | Pr(>F) |
|-----------------|-------|--------|---------|---------|--------|
| log(income)     | 1.00  | 55.77  | 55.77   | 181.76  | 0.00   |
| oil             | 1.00  | 0.02   | 0.02    | 0.07    | 0.79   |
| log(income):oil | 1.00  | 0.09   | 0.09    | 0.31    | 0.58   |
| Residuals       | 95.00 | 29.15  | 0.31    |         |        |

Table: Anova Table

```
> g = lm(log(mortality)~log(income)+log(income):oil,data=infmort2)
> anovatab2=anova(g)
> xtable(anovatab2,digits=2,caption="Anova Table",label="anova2")
```

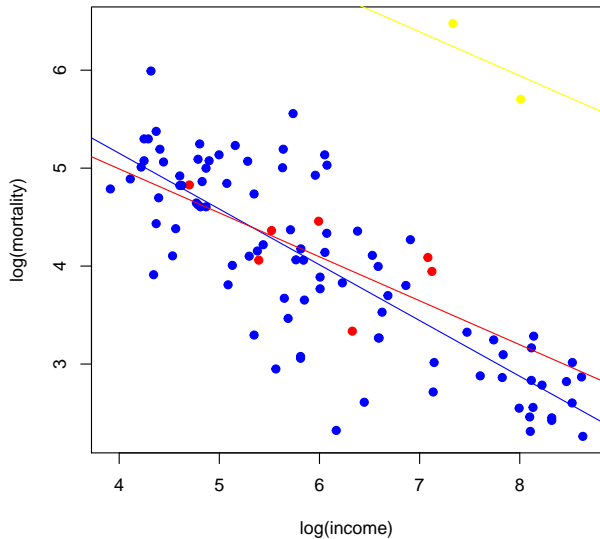
|                 | Df    | Sum Sq | Mean Sq | F value | Pr(>F) |
|-----------------|-------|--------|---------|---------|--------|
| log(income)     | 1.00  | 55.77  | 55.77   | 183.16  | 0.00   |
| log(income):oil | 1.00  | 0.04   | 0.04    | 0.12    | 0.73   |
| Residuals       | 96.00 | 29.23  | 0.30    |         |        |

Table: Anova Table

```
> model.matrix(f)[1:10,]
```

|    | (Intercept) | log(income) | oil | oil exports | log(income):oil | oil exports |
|----|-------------|-------------|-----|-------------|-----------------|-------------|
| 1  | 1           | 8.139149    |     | 0           |                 | 0           |
| 2  | 1           | 8.116716    |     | 0           |                 | 0           |
| 3  | 1           | 8.115521    |     | 0           |                 | 0           |
| 4  | 1           | 8.466110    |     | 0           |                 | 0           |
| 5  | 1           | 8.522976    |     | 0           |                 | 0           |
| 6  | 1           | 8.105308    |     | 0           |                 | 0           |
| 7  | 1           | 8.132413    |     | 0           |                 | 0           |
| 8  | 1           | 8.525161    |     | 0           |                 | 0           |
| 9  | 1           | 7.605392    |     | 0           |                 | 0           |
| 10 | 1           | 7.739794    |     | 0           |                 | 0           |

```
> delt = (row.names(infmort)=="26"|row.names(infmort)=="28")
> infmort$delta = delt
> infmort3 = na.omit(infmort)
> h = lm(log(mortality)~log(income)*oil+delta,data=infmort3)
> hco = 1:length(infmort3$delt)
> hco[oil=="oil exports"]="red"
> hco[oil=="no oil exports"]="blue"
> hco[infmort3$delta]="yellow"
> plot(log(mortality)~log(income),data=infmort3,col=hco,pch=19)
> hc = h$coef
> ab1 = hc[c(1,2)]
> ab2 = hc[c(1,2)]+hc[c(3,5)]
> ab3 = hc[c(1,2)]+hc[c(3,5)]+c(hc[4],0)
> abline(ab1,col="blue")
> abline(ab2,col="red")
> abline(ab3,col="yellow")
```



```
> xtable(anova(h))
```

|                 | Df | Sum Sq | Mean Sq | F value | Pr(>F) |
|-----------------|----|--------|---------|---------|--------|
| log(income)     | 1  | 47.08  | 47.08   | 154.44  | 0.0000 |
| oil             | 1  | 4.58   | 4.58    | 15.03   | 0.0002 |
| delta           | 1  | 12.78  | 12.78   | 41.93   | 0.0000 |
| log(income):oil | 1  | 0.05   | 0.05    | 0.17    | 0.6778 |
| Residuals       | 96 | 29.27  | 0.30    |         |        |