

QUALITATIVE VARIABLES

A. Qualitative Regressors

(Dummy Variables)

— rhs "exogenous" variables for which numerical values are arbitrary.

↖ * Labor economists use these a lot.

Example 1: $w_t = \alpha + \beta_1 E_t + \beta_2 A_t + \beta_3 S_t + \epsilon_t$
(educ) (age) (sex)

where $S_t = \begin{cases} 1 & \text{if } t^{\text{th}} \text{ obs. is male} \\ 0 & \text{if } t^{\text{th}} \text{ obs. is female} \end{cases}$

Thus S_t is a "qualitative" variable; expresses a quality (not quantity) that has some impact on the dependent variable (the wage rate).

This model presumes that the intercept is different for males and females.

for females: $w_t = \alpha + \beta_1 E_t + \beta_2 A_t + \epsilon_t$

for males: $w_t = (\alpha + \beta_3) + \beta_1 E_t + \beta_2 A_t + \epsilon_t$

$\frac{\partial w_t}{\partial S_t} = \beta_3 \rightarrow$ amount that mean wage rate differs for males & females [wage differential].

This model specification further presumes that the difference this variable makes (S_t) is "constant." \exists no interaction with other variables.

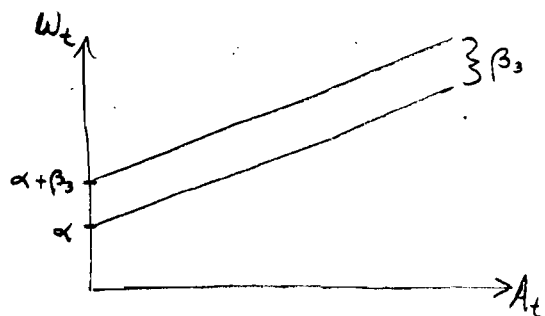
To incorporate interaction, add the variable, $\beta_4 (S_t E_t)$. See example 2.

* Note: Numerical values of dummies are irrelevant; results will be the same.

* Effect: Residuals will sum to zero for each classification (males & females).

Proof: Recall the Normal Equations,

$$\sum_t e_t x_{it} = 0 ; \quad \sum_t e_t = 0 \quad \left(\text{since } \exists \text{ a constant term} \right)$$



$$\text{and } \sum_t e_t S_t = 0$$

[This splits residuals into the two groups; each of which must sum to zero.]

Implication

* Very similar to running 2 separate regressions
 - one for males, one for females;
 - except β & β_3 reflect impact wrt entire sample.

Example 2: Interaction terms

$$w_t = \alpha + \beta_1 E_t + \beta_2 A_t + \beta_3 S_t + \beta_4 (S_t E_t) + \epsilon_t$$

$$\left. \begin{array}{l} \text{for females, } \frac{\partial w_t}{\partial E_t} = \beta_1 \\ \text{for males, } \frac{\partial w_t}{\partial E_t} = \beta_1 + \beta_4 \end{array} \right\} \begin{array}{l} \text{Education may} \\ \uparrow w_t \text{ for males} \\ \text{more than females.} \end{array}$$

This interaction term allows measurement of the impact of education on the two separate groups in the sample

Interpretation; $\beta_3 = \overset{\text{avg.}}{\text{male-female wage differential}}$
 $\beta_4 = \overset{\text{(differential)}}{\text{the effect of Education}}$
 on the wage rate of males & females.

If we think Age affects the wage rate of males and females differently, we can add another interaction term;

$$w_t = \alpha + \beta_1 E_t + \beta_2 A_t + \beta_3 S_t + \beta_4 (S_t E_t) + \beta_5 (S_t A_t) + \epsilon_t$$

→ $\beta_5 = \text{the differential effect of Age on the wage rate of males and females.}$

Note: Running this last Regression, with all possible interaction terms (re: S_t), is (almost) the same as running two separate regressions
 — one for males & one for females.

(i) The coefficients will be identical.

$$\text{i.e. } \frac{\partial W_t}{\partial E_t(\text{females})} = \hat{\beta}_1 ; \frac{\partial W_t}{\partial E_t(\text{males})} = \hat{\beta}_1 + \hat{\beta}_4,$$

etc. whether run with all possible interaction terms or not (separate regressions)

(ii) When running two separate regressions, the estimates of the variance of the residuals ($S_{\text{male}}^2 \neq S_{\text{female}}^2$) need not be identical.

Running one regression with all relevant interaction terms presumes that they are identical ($S_{\text{male}}^2 = S_{\text{female}}^2$).

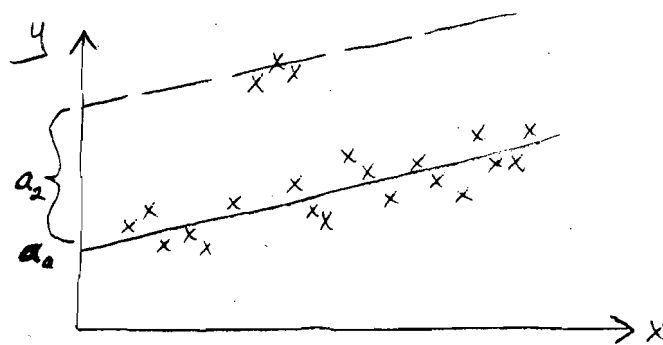
Thus, even though the corresponding estimates will be identical with these two alternatives, the t-ratios may not be.
 (since the standard errors depend on S^2 !)

These procedures actually represent another method of statistical analysis called Analysis of Variance.

Example 3: Outlier Dummies

e.g. | unemployment rates
 | during a year of
 | a strike;
 | commodity prices
 | during a period
 | of an embargo;

(production, ..., during years of a war; etc.



Could include a Dummy with a value of 1 for the outliers.

$$y = a_0 + a_1 x + a_2 D + \epsilon_t$$

This would result in a constant term that shifts the entire regression line up thru the outliers, when they are effective.

Alternatively explained, this brings the outliers down about the regression line, so that they are not "outliers" anymore.

R^2 may = .5 w/o Dummies
 and may = .9 with Dummies.

Big Deal! Implies, w/o Dummies can "explain" 50% of variation; with Dummies can "explain" 90%, BUT 40% of that is due to (weird) attributes of outliers!

This could be useful, if you know why the outliers should be outliers!
e.g. strike, embargo, war, ...

Note: R^2 is improved immensely with Dummies.

If you "throw out" the obs. that don't fit, then the fit of everything else will improve!

[R^2 bigger; SSE smaller; std. errors smaller, ...]

Note: If I just one outlying observation that is Dummied out,

the regression will bring the outlier right down onto the regression line.

This will not affect the slope or position of the fitted line;

\therefore effectively throwing out that obs.!

→ Not "explaining" anything new with Dummy

POINT: Use Caution!

Know what you're doing!

Don't "throw out" outliers without knowing why; you had better be able to tell a good story as to why those obs. are outliers ($\epsilon \therefore$ should be "thrown out").

Don't use Dummies and then brag about your R^2 !



Example 4: Seasonal Dummies

$$C_t = \alpha + \beta y_t + \gamma_1 S_{1t} + \gamma_2 S_{2t} + \gamma_3 S_{3t} + \epsilon_t$$

where $S_{1t} = \begin{cases} 1 & \text{during first quarter} \\ 0 & \text{otherwise} \end{cases}$

$S_{2t} = \begin{cases} 1 & \text{during second quarter} \\ 0 & \text{otherwise} \end{cases}$

$S_{3t} = \begin{cases} 1 & \text{during third quarter} \\ 0 & \text{otherwise} \end{cases}$

This holds the slope, β , constant for all 4 quarters, but allows a different intercept for different quarters.

Computer Problem 4: Seasonality Study

Given below is monthly data on Visits to National Parks (V), Personal Income (Y), U. S. Population (N), and the Consumer Price Index (P), for the years 1971 through 1976.

1. From this data you will need to construct the following time series:

$$\text{Log } (V/N); \quad \text{Log } (Y/(N*P));$$

$$\text{CMA}_t = (.5*V_{t-6} + V_{t-5} + V_{t-4} + V_{t-3} + V_{t-2} + V_{t-1} + V_t + V_{t+1} + V_{t+2} + V_{t+3} + V_{t+4} + V_{t+5} + .5*V_{t+6})/12$$

(Note that this series will run from July, 1971 through June, 1976).

$$\text{RCMA}_t = \frac{V_t}{\text{CMA}_t} \quad (\text{over the same time period})$$

From RCMA, calculate by hand the seasonal adjustment factors, $\overline{\text{RCMA}}_i$;

$$i = \text{Jan.}, \text{Feb.}, \dots, \text{Dec.}$$

Then construct the Seasonally Adjusted Visits series;

$$\text{SAV}_t = V_t / \overline{\text{RCMA}}_t \quad (\text{from July, 1971 through June, 1976})$$

2. Next run the following regressions:

$$(1) \text{Log } (V/N) = a_0 + a_1 \text{Log } (Y/(P*N)) + \epsilon_1$$

$$(2) \text{Log } (V/N) = a_0 + a_1 \text{Log } (Y/(P*N)) + \sum_{i=1}^{12} b_i S_i + \epsilon_2$$

$$\text{where } S_i = \begin{cases} 1 & \text{in month } i; \\ 0 & \text{otherwise} \end{cases}$$

$$(3) \text{Log } (\text{SAV}/N) = a_0 + a_1 \text{Log } (Y/(P*N)) + \epsilon_3$$

Note: Run equations (1) and (2) twice; once over the entire sample, and once over the abridged sample, so that they may be compared with the third regression.

Compare and contrast the results (SSE, R^2 , \hat{a}_1 , t-ratios, etc.) for these three equations, and explain why the results vary.

3. In all cases plot the Actual values $(\frac{V}{N})$ over the Fitted values $(\frac{\hat{V}}{N})$, and the residuals against time.

4. Compare the seasonal adjustment factors computed by hand ($\overline{\text{RCMA}}_i$, $i = \text{Jan.}, \text{Feb.}, \dots, \text{Dec.}$) with those implied by equation (2), the dummy regression. Note that you must calculate the latter ($e^{\hat{b}_i}$ = seasonal adjustment factor for month i ; $i = \text{Jan.}, \text{Feb.}, \dots, \text{Dec.}$).

Problem #5

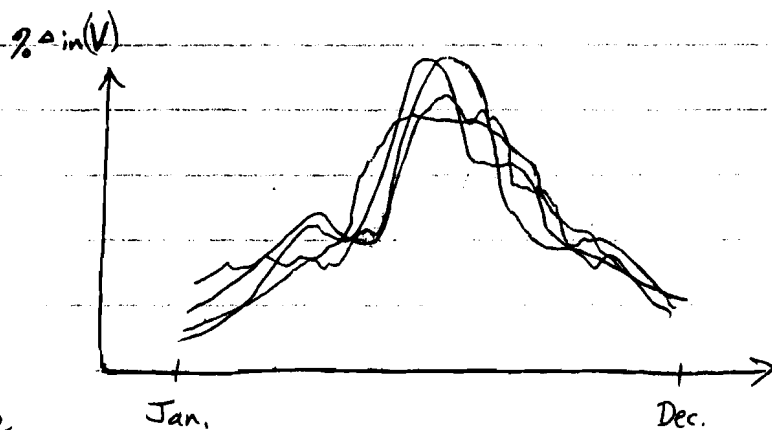
SEASONAL ADJUSTMENT

Consider the series, Visits to National Parks.

— monthly data, 6 yrs.; BUS. STATISTICS (S.C.BUS.)

If you plot variations during each year,
a seasonal pattern appears.

During any year,
∫ systematic variation
in the series
which is not due to
any explanatory variable
such as income or population,
but is due, rather, to seasonality.



→ Want to consider the relationship between
(Visits/N) and Real Personal Income Per Capita

If we regress:

$$\ln\left(\frac{V}{N}\right) = a_0 + a_1 \ln\left(\frac{Y}{CPI * N}\right)$$

we will observe the relationship between

variations in $\frac{Y}{CPI * N}$ and $\frac{V}{N}$;

$\frac{V}{N}$ will display much noise during each year, not due
to variations in $\frac{Y}{CPI * N}$ but due to seasonality.

We are really interested in the relationship between the number of Visits in the typical, avg. month and income, without regard to seasonality (or, holding seasonality constant).

We can capture this desired result by:

- i) Seasonally Adjusting $\frac{V}{N}$ by hand, or
- ii) Fitting $\ln \frac{V}{N} = a_0 + a_1 \ln \frac{Y}{CPI \times N} + b_1 D_1 + b_2 D_2 + \dots + b_{11} D_{11} + b_{12} I$

i) By Hand.

We want the Index; $I_i = \frac{\text{Visits in the } i^{\text{th}} \text{ month}}{\text{Index of avg month}}$

Step 1: Construct a 12-month moving avg.

e.g. observation for June;

Dec Jan. Feb. Mar. Apr. May Je. July Aug. Sept. Oct. Nov. Dec.

$$CMA_{Je.} = \frac{1}{2} \text{Dec.} + \text{Jan.} + \text{Feb.} + \dots + \text{Oct.} + \text{Nov.} + \frac{1}{2} \text{Dec.}$$

SAMPL 7,66 \$ → [72 obs., 1st six & last six used up in MA]

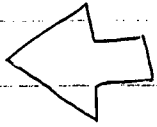
$$\Rightarrow \text{GENR CMA} = [0.5 * V(-6) + V(-5) + \dots + V(4) + V(5) + 0.5 * V(6)] / 12. \$$$

↑
(centered moving avg)

Step 2:
Then

GENR RCMA = $\frac{V}{CMA}$ ‡

↳ Ratio of Visits in any month, to the CMA.
RCMA should vary around 1;
RCMA_{Je} should be > RCMA_{Jan}!



We have 5 years of observations. [lose 6 mo. at front & 6 mo. at end]
 ↳ ∴ 5 obs. of RCMA_{Jan}
 5 " " RCMA_{Feb.}
 etc.

* We are interested in the Average # of Monthly Visits in each month, as a percentage of the avg # of monthly visits in the years centered by that month.

Year	1	2	3	4	5	6	Averages	← Step 3
Jan.		~				~	→ RCMA RCMA _{Jan}	<div style="border: 1px solid black; padding: 5px; width: fit-content;"> We want these Avg's ↓ These tell us the relative seasonal expectation of visits for each month </div>
Feb.		~				~	→ RCMA _{Feb}	
Mar.		~				~	→ RCMA _{Mar}	
Apr.		~				~	→ RCMA _{Apr}	
May		~				~	→ RCMA _{May}	
Je.		~				~	→ RCMA _{Je.}	
July	~	~				~	→ RCMA _{July}	
Aug.	~	~				~	→ RCMA _{Aug.}	
Sept.	~	~				~	→ RCMA _{Sept.}	
Oct.	~	~				~	→ RCMA _{Oct.}	
Nov.	~	~				~	→ RCMA _{Nov.}	
Dec.	~	~				~	→ RCMA _{Dec.}	

These \overline{RCMA}_i should again fluctuate around 1.0

Furthermore, $\sum_{i=1}^{12} \overline{RCMA}_i$ should ≈ 12 . (or 1200 if $\times 100$)

Note: $\sum_{i=1}^{12} \overline{RCMA}_i$ will not exactly = 12, because of the small sample, and the gaps at the front and back.

Suppose $\sum_{i=1}^{12} \overline{RCMA}_i = 15$;

Then adjust each \overline{RCMA}_i by $\times \left(\frac{1200}{15.00} \right)$ so that they do sum to 12.

These numbers so constructed are the seasonal adjustment factors.

Step 4:

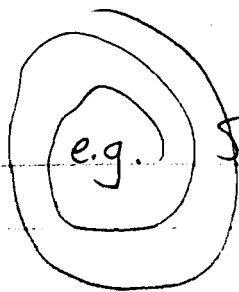
We now construct the seasonally adjusted series, I_t , as follows:

$$\boxed{\text{S.A.V.}} = \frac{V_t}{\overline{RCMA}_t} \quad t = \text{Jan, Feb, } \dots, \text{ Dec.}$$

I_t will now display movements in visits to parks, adjusted for seasonality!

Shows movements in Visits over time (in each month) relative to the # of visits that are expected during that month of the year.

Seasonal adj. factors:



Suppose

$$\overline{RCMA}_{Jan} = .5, \quad \overline{RCMA}_{Je} = 1.2$$

$$V_{Jan} = 50,000 \text{ visits};$$

$$V_{Je} = 120,000 \text{ visits};$$

Seasonally Adjusted:

$$I_{Jan} = \frac{50,000}{.5} = 100,000;$$

$$I_{Je} = \frac{120,000}{1.2} = 100,000$$

This suggests that January and June experienced the same number of monthly visits, after adjusting for seasonality, i.e. w/rt the average number of monthly visits expected in these ~~months~~ months.

Alternatively, suppose

$$V_{Jan} = 60,000;$$

$$V_{Je} = 100,000;$$

$$\text{Then;} \quad I_{Jan} = \frac{60,000}{.5} = 120,000;$$

$$I_{Je} = \frac{100,000}{1.2} = 85,000$$

This suggests that relative to seasonal expectations, June is getting fewer visits than January. [Even though the total # of visits is 40,000 more !!]

Seasonal Adjustment ii) by Dummy Regression

Case 1: $V = a_0 + a_1 Y + b_1 S_1 + b_2 S_2 + \dots + b_{12} S_{12}$

where $S_i = 1$ in month i ; 0 elsewhere

PROBLEM: \exists perfect multicollinearity ~~multicollinearity~~ here!

a_0 perfectly correlated with $\sum S_i$.

If \exists one solution, \exists an infinite # of solutions!

" " " " " "

I can add any number, k , to $b_1, b_2, \dots, \neq b_{12}$,
and subtract k from a_0 ,
and I'll get the same SSE!

Can solve problem by imposing some constraint.

Which constraint to impose?

$b_{12} = 0$, or $b_i = 0$ $i = 1, \dots, 11$ [or $a_0 = 0$]

— leave one out.

— impose the constraint easiest to use.

★ \Rightarrow It doesn't matter; results all the same.

BUT be careful; since you must explain
the results to someone else!

Story: You do a study with dummies for the 50 states, leaving out Texas.

You present the results to a Senate Subcommittee.

The distinguished senator from Texas says, "What do you mean son - by leaving out Texas!"
the great state of

You respond, "I actually didn't leave Texas out, but rather, I made Texas the basis for the entire study!"

Then the senator from Massachusetts says, "What do you mean, sir!"

It's no use trying to explain that it doesn't matter. All they see is a 1 or a 0.

~~////~~

Or: Use a sex dummy, and they ask, "Why is Male = 1 & Female = 0?!"

It's no use trying to explain it doesn't matter.

* Can ~~avoid~~ ^{avoid} this problem by simply imposing a different constraint. [any ^{relevant} constraint will eliminate multicoll.]
e.g. that the sum of all the coefficients of the dummies must be 0

$$V = a_0 + a_1 Y + b_1 S_1 + b_2 S_2 + \dots + b_{12} S_{12}$$

\rightarrow s.t. $\sum b_i = 0$

Anyone can understand this constraint!
No problems explaining.

Easiest procedure to impose this constraint:

Run regression with $b_{12} = 0$;

get $\hat{a}_0, \hat{a}_1, \hat{b}_1, \hat{b}_2, \dots, \hat{b}_{11}$, and $\hat{b}_{12} = 0$

Can manipulate seasonal dummies;

convert them into deviations from the mean seasonal coeff.

$$\bar{\hat{b}}_0 = \text{Mean Seasonal Coeff.} = \frac{\sum_{i=1}^{12} \hat{b}_i + 0}{12}$$

Then let $\tilde{b}_1 = \hat{b}_1 - \bar{\hat{b}}_0, \tilde{b}_2 = \hat{b}_2 - \bar{\hat{b}}_0, \dots, \tilde{b}_{12} = 0 - \bar{\hat{b}}_0$;

so that $\sum_{i=1}^{12} \tilde{b}_i = 0$.

\Rightarrow Same effective results;
easy to interpret. (For anyone! even a senator!)

Do THIS!!

Converting Results to seasonal coefficients:

[For 2nd regression - seasonality with constant slope]

$$\ln \frac{V}{N} = a_0 + a_1 \ln \frac{Y}{CPI * N} + b_1 S_1 + b_2 S_2 + \dots + b_{12} S_{12}$$

$$\boxed{\text{s.t. } \sum_{i=1}^{12} b_i = 0}$$

take antilog!

$$\Rightarrow \frac{V}{N} = a_0 \left[\frac{Y}{CPI * N} \right]^{a_1} \left\{ e^{b_1 S_1} e^{b_2 S_2} e^{b_3 S_3} \dots e^{b_{12} S_{12}} \right\}$$

Then in Jan., $b_1 S_1 = b_1 \cdot 1$, and $S_2 = S_3 = \dots = 0$;
 and $\frac{V}{N} = a_0 \left[\frac{Y}{CPI * N} \right]^{a_1} e^{b_1}$ → gets you estimate of actual $\frac{V}{N}$

Interpretation; $\hat{a}_0 \left[\frac{Y}{CPI * N} \right]^{\hat{a}_1}$ is the "fitted" average # of visits, over all observations (regardless of seasonality, or holding constant for seasonality) Seasonally adjusted!!

and $[e^{\hat{b}_i}]$ is the seasonal coefficient for Jan.

∴ Seasonal coeff.'s are $(e^{\hat{b}_i})$, $i=1, \dots, 12$
 \hat{b}_i must be constrained to sum to zero here!!

NOTE:

These seasonal coefficients should be (almost) identical to those derived by hand, (\overline{RCMA}_i)

Plot each set of seasonal adjustment factors

B. Binary Dependent Variables

The lhs variable is a Dummy ;
attempting to explain a qualitative variable
with a regression equation.

Models 1.- 3. go from poor models to good models.

1. Linear Model (poor)

$$y_t = X_t \beta + \epsilon_t ; \quad y_t = \begin{cases} 1 & \text{if event occurs} \\ 0 & \text{if not} \end{cases}$$

$T \times K \quad K \times 1$

Interpretation: This model implies that

$$X_t \beta = E(y_t) = P(y_t = 1)$$

= probability that the event occurs

[Note: $E(y_t) = p(y_t=1)(1) + p(y_t=0)(0).$]

Example: Might want to attempt to measure the factors affecting an individual's decision to work.

$$y_t = \begin{cases} 1 & \text{if indiv. supplies labor} \\ 0 & \text{if not} \end{cases}$$

We wish to measure β , the impact of the variables involved (X_t) on the individual's decision to wk

With the Linear Model specification,

$$\beta = \frac{\partial E(y_t)}{\partial x_t} = \frac{\partial P(y_t=1)}{\partial x_t}$$

Problem 1: ϵ_t is strange; it has just 2 values.

We have $y_t = P_t + \epsilon_t$; where $P_t = P(y_t=1) = x_t\beta$

Thus, since y_t has a binary distribution,
 ϵ_t has a binary distribution.

$$\epsilon_t: \begin{array}{l} \text{possible values:} \\ \text{probability:} \end{array} \left. \begin{array}{cc} * & ** \\ 1 - P_t & -P_t \\ P_t & 1 - P_t \end{array} \right\} \text{or: } \begin{array}{cc} 1 - x_t\beta & -x_t\beta \\ x_t\beta & 1 - x_t\beta \end{array}$$

* If $P(y_t=1) = x_t\beta = P_t$,
then $y_t = 1 = x_t\beta + \epsilon_t$
 $= P_t + (1 - P_t)$.

** Analogously,
if $P(y_t=0) = 1 - P_t$,
then $y_t = 0 = x_t\beta + \epsilon_t$
 $= P_t + (-P_t)$.

Consider the first two moments of ϵ_t .

$$\begin{aligned} E(\epsilon_t) &= (1 - x_t\beta)(x_t\beta) + (-x_t\beta)(1 - x_t\beta) \\ &= 0 \end{aligned}$$

$$\begin{aligned}
 \text{Var}(\epsilon_t) &= E(\epsilon_t^2) \\
 &= (1-P_t)^2(P_t) + (-P_t)^2(1-P_t) \\
 &= (1-X_t\beta)^2(X_t\beta) + (-X_t\beta)^2(1-X_t\beta) \\
 &\quad \text{factoring} \\
 &= (1-X_t\beta) \left[(1-X_t\beta)(X_t\beta) + (-X_t\beta)^2 \right] \\
 &= (1-X_t\beta) X_t\beta
 \end{aligned}$$

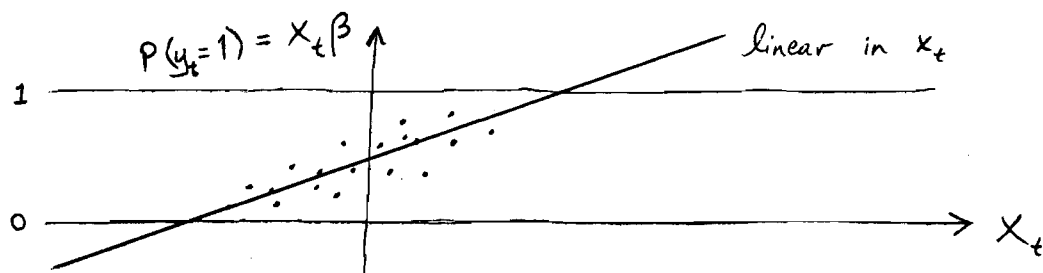
Observe that $\text{Var}(\epsilon_t)$ depends on X_t ; it is not constant; \exists heteroskedasticity.

This is not a very serious problem; can be corrected with GLS-type transformation of the model.

Problem 2: (big problem)

What if $X_t\beta < 0$ or $X_t\beta > 1$?

Can't be! It is a probability!



Since $y_t = X_t\beta + \epsilon_t$ is linear in X_t , for some X_t 's $X_t\beta$ may imply probabilities outside the range, $(0, 1)$.

You may run this and get implied probabilities within the range $(0,1)$. Then you're OK.

But you may get some outliers that are uninterpretable!

Some people do this and just wave their hands and say, "it is obvious that this implies a probability of 1 (or of 0)."

It is not obvious!

This is poor procedure.

Linear Model is poor,
for this reason.

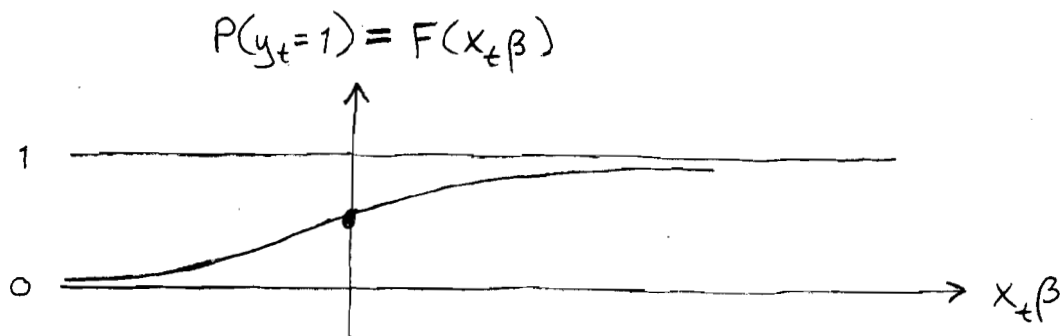
2. Probit Model

How can this second problem be dealt with?

We are attempting to fit a curve to the data, where the data represent the relationship between a probability (that must range between 0 and 1) and factors that influence it, X_t .

It is not a linear relationship that we are trying to fit!

It is actually something like a cumulative distribution function!



The Probit Model simply replaces the linear specification with a nonlinear cumulative distribution function.

Heroic Assumption:

The Probit model uses as F , the Standard Normal c. d. f.
 $[N(0,1); \text{Expected Value} = 0, \text{Variance} = 1]$

The Model:

$$y_t^* = X_t\beta + \epsilon_t \quad \text{with the } \epsilon_t \text{ iid } N(0,1)$$

where y_t^* is unobservable,
 but y_t is observed;

$$y_t = \begin{cases} 1 & \text{if } y_t^* > 0 \\ 0 & \text{if } y_t^* \leq 0. \end{cases}$$

$$\begin{aligned} \text{Then } E(y_t) &= P(y_t=1) = P(y_t^* > 0) = P(X_t\beta + \epsilon_t > 0) \\ &= P(\epsilon_t > -X_t\beta) \\ &= P(\epsilon_t < X_t\beta) \\ &= F(X_t\beta) \end{aligned}$$

(by symmetry of Standard Normal Dist.)

Which dist. to use?

Effect: We generate a probability that is nice.

Comments:

- 1) Can add a constant term,
with no change in the implications;

$$y_t = \begin{cases} 1 & \text{if } y_t^* = \alpha + \beta x_t + \epsilon_t > 0 \\ 0 & \text{if } y_t^* = \alpha + \beta x_t + \epsilon_t \leq 0 \end{cases}$$

- 2) We are estimating β ;
the impact of x_t on the probability
of the event occurring, $P(y_t=1)$.
— more explanation below.

- 3) The Standard Normal c.d.f. is used
for good reason. Consider an alternative
assumption that is more general.
let ϵ_t be iid $N(0, \underline{\sigma^2})$.

$$\begin{aligned} \rightarrow E(y_t) &= P(y_t=1) = P(y_t^* > 0) = P(x_t\beta + \epsilon_t > 0) \\ &\stackrel{(\neq 0)}{=} P(\epsilon_t > -x_t\beta) \end{aligned}$$

(Must standardize to obtain
the cdf, F)

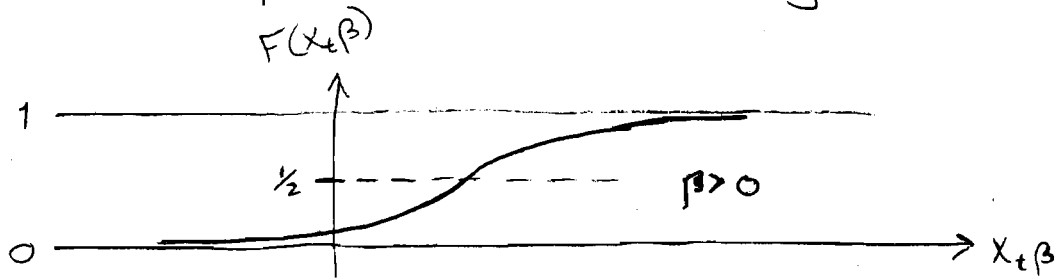
Here we can estimate $\frac{\beta}{\sigma}$,
but we cannot obtain $\beta \neq \sigma$.
 \exists an identification problem.

$$\begin{aligned} &= P(\epsilon_t < x_t\beta) \\ &= P\left(\frac{\epsilon_t}{\sigma} < x_t \frac{\beta}{\sigma}\right) \\ &= F\left(x_t \frac{\beta}{\sigma}\right) \end{aligned}$$

$\sigma^2 = 1$ is a "normalization" assumption.

Comments, cont.

4) A drawback (for some models) is that there is symmetry above and below the probability of $\frac{1}{2}$, due to the assumption of Normality.



Estimation of Probit Models : by MLE

$$\text{Here, } E(y_t) = P(y_t=1) = P_t = F(x_t\beta)$$

$$\text{and } P(y_t=0) = 1 - P_t = 1 - F(x_t\beta)$$

$$\mathcal{L} = P(y_1) P(y_2) P(y_3) \dots P(y_T)$$

$$= \prod_{t \rightarrow y_t=1} P_t \prod_{t \rightarrow y_t=0} (1 - P_t)$$

$$= \prod_{t=1}^T P_t^{y_t} (1 - P_t)^{(1-y_t)}$$

These powers will mechanically include a factor, P_t , if $y_t=1$, and $(1-P_t)$ if $y_t=0$.

$$= \prod_{t=1}^T [F(x_t\beta)]^{y_t} [1 - F(x_t\beta)]^{(1-y_t)}$$

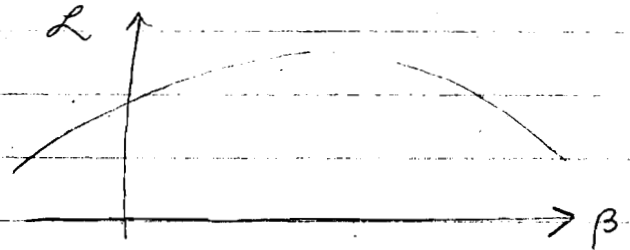
Binary
Dist.

$$\text{Thus, } \log \mathcal{L} = \sum_{t=1}^T \left\{ y_t \log [F(x_t, \beta)] + (1 - y_t) \log [1 - F(x_t, \beta)] \right\}$$

Find the values of the parameters, β , that maximize this; MLE ($\hat{\beta}$).

Note: $\log \mathcal{L}$ has a global maximum with the Probit model specification.

Thus, MLE ($\hat{\beta}$) is "easy" to find with any iterative technique.



— find the top of a hill.

What β means:

$$\frac{\partial P(y_t=1)}{\partial x_{ti}} = \frac{\partial F(x_t, \beta)}{\partial x_{ti}} = f(x_t, \beta) \frac{\partial x_t \beta}{\partial x_{ti}} = f(x_t, \beta) \beta_i$$

where f is the Standard Normal density fn.

Observe that x_{ti} affects $P(y_t=1)$ by an amount, $f(x_t, \beta) \beta_i$.

→ If $\beta_i > 0$, then if x_{ti} ↑, $P(y_t=1)$ ↑.

Attended meetings;

Paper presented, studying the likelihood that loan applicants would default on loans;

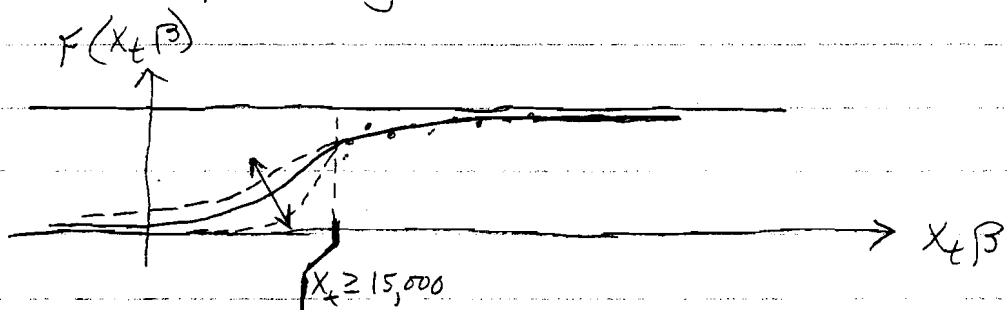
This probability of default might be influenced by factors such as the applicants' income, age, education, ...

Had data from Credit Unions on people with Incomes $\geq 15,000$; data on Income, age, education, ..., and whether or not they defaulted.

Fit probit model; got interesting results, claimed ability to "predict" the likelihood that an applicant would default, based on these factors.
 — Potentially Valuable for Credit Unions, and all financial institutions!

De Min Wu, U of Kansas, discussant; argued that sample was not random; truncated at income levels $\geq 15,000$.

Thus, the fitted curve is only one possibility; many other curves might be consistent with the given data. Thus, results ($\hat{\beta}$) not necessarily reliable for predicting this likelihood.



3. Logit Model



$$\log \left[\frac{P(y_t=1)}{P(y_t=0)} \right] = X_t \beta$$

X_t affects the (log of) relative probabilities of events, $y_t=1$ and $y_t=0$, by the amount, β .

Note: \exists no error term here.

This is a description of the probability of an event occurring, and this probability is deterministic, as the model is specified.

It is then the actual occurrence or non-occurrence of the event which is random.

Observe: IF $P(y_t=1) = P_t$, then Logit is

$$\log \left(\frac{P_t}{1-P_t} \right) = X_t \beta$$

$$\Rightarrow \frac{P_t}{1-P_t} = e^{X_t \beta}$$

$$\Rightarrow P_t = (1-P_t) e^{X_t \beta}$$

$$\Rightarrow P_t + P_t e^{X_t \beta} = e^{X_t \beta}$$

$$P(y_t=1) \Rightarrow P_t = \frac{e^{X_t \beta}}{1 + e^{X_t \beta}} = \frac{1}{e^{-X_t \beta} (1 + e^{X_t \beta})} = \frac{1}{1 + e^{-X_t \beta}}$$



$$= G(X_t \beta)$$

[cdf for the Logistic Distribution]

Probit uses the Standard Normal cdf;
Logit uses the Logistic cdf.

18

Note: $G(x) = \frac{1}{1+e^{-x}}$ is the Logistic Dist.

Note: $P(y_t = 0) = 1 - P(y_t = 1)$

$$\begin{aligned} &= 1 - \frac{e^{x_t\beta}}{1+e^{x_t\beta}} &&= 1 - \frac{1}{1+e^{-x_t\beta}} \\ &= \frac{1}{1+e^{x_t\beta}} &&= \frac{e^{-x_t\beta}}{1+e^{-x_t\beta}} \end{aligned}$$

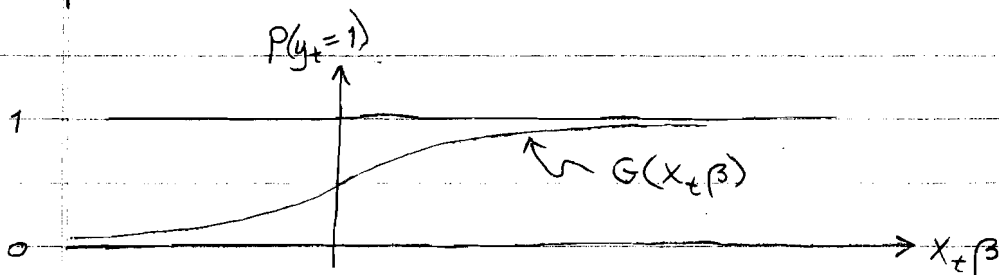
SKIP

Advantage of Logit:

Although $P(y_t = 1)$ and $[1 - P(y_t = 1)]$ must be in the interval $(0, 1)$,

$\left[\frac{P(y_t = 1)}{P(y_t = 0)} \right]$ may range from $[0 \text{ to } \infty]$,

and thus, $\log \left[\frac{P(y_t = 1)}{P(y_t = 0)} \right] = x_t\beta$ [what we're trying to explain w/ x_t] may range from $[-\infty \text{ to } \infty]$.



The Logistic Dist., G , is much like the Standard Normal cdf used in Probit. The difference amounts to slightly thicker tails in one of the two.

Estimation of Logit: by MLE

$$\mathcal{L} = \prod_{t=1}^T G(x_t\beta)^{y_t} [1 - G(x_t\beta)]^{(1-y_t)}$$

$$= \prod_{t: y_t=1} \pi(P_t) \prod_{t: y_t=0} \pi(1-P_t)$$

$$= \prod_{t: y_t=1} \left(\frac{e^{x_t\beta}}{1 + e^{x_t\beta}} \right) \prod_{t: y_t=0} \left(\frac{1}{1 + e^{x_t\beta}} \right)$$

Same as Binary Probit,
with $G(x_t\beta)$ replacing $F(x_t\beta)$.
Like Probit, this \mathcal{L} has a
global maximum.

So $\max \mathcal{L}$ wrt $\beta \rightarrow \hat{\beta}_{MLE}$.

Thus ends discussion of Binary Dependent variables.

- Use Probit or Logit!

- Slightly more expensive than least squares,
but much better.

↑

[on Linear Model]

SAS - PROC PROBIT
PROC LOGIT (Logit)

PROC LOGIST

C. Qualitative Dependent Variables with more than 2 values.

$$y = f(x's)$$

$$\text{where } y = \begin{cases} 1 & \text{if in classific. I} \\ 2 & \text{if in classific. II} \\ \vdots & \\ k & \text{if in classific. k} \end{cases}$$

e.g. attempting to explain how various factors might influence how individuals are classified into:

- income groups
- education levels
- good products - faulty prod.
- ⋮

Any situation in which you hypothesize that certain factors are influential in discriminating among various classifications.

7 many possible models to choose from

1. Multi^{nomial} ~~nominal~~ Logit
2. Conditional Logit
3. Multinomial Probit (with ordered or unordered responses)
4. Discriminant Analysis

D. Limited Dependent Variables

This is situation in which lhs variable is drawn from or described by a truncated distribution.

1. Tobit Model (by Tobin)

The model; $y_t^* = X_t \beta + E_t \quad E_t \sim N(0, \sigma^2)$
(under $X_t \beta$)
 usual rhs variables

y_t^* is the variable from the entire population of this particular distribution.

However, we don't observe y_t^* .

Instead, we have a sample that is drawn from a subsample or truncated portion of the distribution of y_t^* .

e.g. observe $y_t = \begin{cases} y_t^* & \text{if } y_t^* > 0 \\ 0 & \text{if } y_t^* \leq 0 \end{cases}$

Note: $E(y_t^*) = X_t \beta$; this is the mean of the dist. of y_t^* , but not of the truncated distribution of y_t !

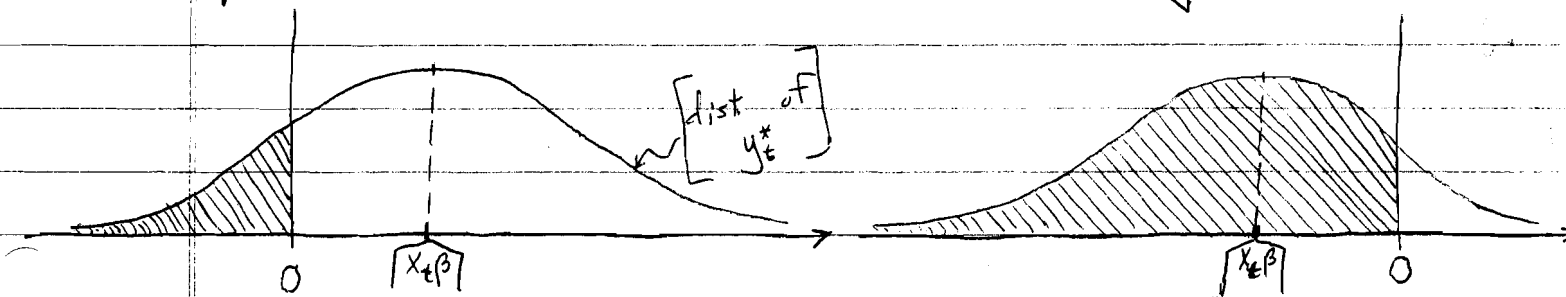
Tobin's example:

The amount of expenditures on cars
in one year, by a household.

This is bounded, or truncated at zero.

$$[\text{Car exp. in 1 yr.}] = \begin{cases} 0 & \text{for 70\% of population} \\ 300-50,000 & \text{for 30\% of population} \end{cases}$$

2 examples:



Over a 5-yr. pd,
most people buy a car

Over a 1-yr. pd,
most do not buy a car
(70% in above ex.)

~~THE KEY~~
The Difference between Tobit & Binary Probit (or Logit):

$$\text{Probit} \quad \text{observed } y_t = \begin{cases} 1 & \text{if } y_t^* > 0 \\ 0 & \text{if } y_t^* \leq 0 \end{cases} \quad \text{where } y_t^* = X_t\beta + \epsilon_t \text{ is unobserved}$$

$$\text{Tobit} \quad \text{observed } y_t = \begin{cases} y_t^* & \text{if } y_t^* > 0 \\ 0 & \text{if } y_t^* \leq 0 \end{cases}$$

Thus, Tobit is for a truly Limited dependent var.

Interpretation of β :

$$\beta = \frac{\partial E(y_t^*)}{\partial x_t} = \frac{\partial x_t \beta}{\partial x_t}$$

Thus, β is the impact on the mean of the entire (unlimited) population, y_t^* , from a change in x_t .

This interpretation is clouded when we observe y_t and not y_t^* .

Tobin shows how this kind of model can be estimated.

— Somewhat gory, but can be done.

E. Sample Selection Bias

||| In Tobit model, our sample consisted of all people in the entire population; but only some bought cars this year.

What if you didn't have a random sample, but instead, you only had a sample of people who bought cars (nonzero obs.).

— e.g. 1; only have car registration data.
 [e.g. 2; financial data on loan applicants with $y \geq 15,000$].

Here, the sample is not random!

— sample selection bias exists.

This is not a Tobit model since we don't have data on people who do not buy cars. We only have data for the truncated distribution, y_t ; not on the entire distribution, y_t^* .

★ See Hausman & Weiss, ECON, 1977.
 Heckman, ECON, 1979.

They show how to deal with this model.

— Get something like Max. Likelihood Estimates, only the density functions involved are conditional densities (conditional on $y_t^* > 0, \dots$).

F. Simultaneous equation models in which one equation involves a Qualitative Dependent Var.

e.g. Relationship between Wage Rates and the Quit Rate (Q) (Union Membership (U_t))

— Simultaneously determined. $W_t \leftrightarrow Q_t$.

Cross-sectional Survey data on individuals' Wage Rate, Income, Education Level, Age, ..., and whether or not they recently "Quit." (belong to a Union)

System:

$$W_t = F(Y, E, A, Q_t)$$

$$Q_t = F(W_t, Y, E, A)$$

$$Q \text{ (or } U_t) \rightarrow W_t$$

$$W_t \rightarrow Q \text{ (or } U_t)$$

$$\text{where } Q_t = \begin{cases} 1 & \text{if (union member) quit recently} \\ 0 & \text{otherwise} \end{cases}$$

Can estimate this kind of model.

1. Simultaneous Probit Model
2. Simultaneous Logit Model
3. Simultaneous Tobit Model

⋮