

ECON2228 Notes 6

Christopher F Baum

Boston College Economics

2014–2015

Chapter 7: Multiple regression analysis with qualitative information: Binary (or dummy) variables

We often consider relationships between observed outcomes and qualitative factors: models in which a continuous dependent variable is related to a number of explanatory factors, some of which are quantitative and some of which are qualitative.

In econometrics, we also consider models of qualitative *dependent* variables, but we will not explore those models in this course due to time constraints. But we can readily evaluate the use of qualitative information in standard regression models with continuous dependent variables.

Qualitative information often arises in terms of some coding, or index, which takes on a number of values: for instance, we may know in which one of the six New England states each of the individuals in our sample resides. The data themselves may be coded with the bilateral “MA”, “RI”, “ME”, etc.

How can we use this qualitative factor in a regression equation? In the data, `state` takes on six distinct values. We must create six *binary variables*, or *dummy variables*, each of which will refer to one state: that is, that variable will be 1 if the individual comes from that state, and 0 otherwise. We can generate this set of 6 variables easily in Stata with the command `tab state, gen(st)`, which will create 6 new variables in our dataset: `st1`, `st2`, ... `st6`. Each of these variables are dummies: that is, they only contain 0 or 1 values.

These variables are known as a set of *mutually exclusive and exhaustive* (MEE) measures. They are exclusive, because each individual has only one primary state of residence. They are exhaustive, in that every individual in the sample lives in one of the states.

If we add up these variables, we must get a vector of 1's, suggesting that we will never want to use all 6 variables in a regression (as by knowing the values of any 5...) We may also find the proportions of each state's citizens in our sample very easily: $\sum st^*$ will give the descriptive statistics of all 6 variables, and the mean of each st dummy is the sample proportion living in that state.

In Stata 11+, we actually do not have to create these variables explicitly; we can make use of *factor variables*, which will automatically create the dummies “on the fly” and make them accessible.

How can we use these dummy variables? Say that we wanted to know whether incomes differed significantly across the 6-state region. What if we regressed `income` on *any five* of these `st` dummies? We could do this with explicit dummy variables as

```
regress income st2-st6
```

or with factor variables as

```
regress income i.state
```

In either case, we are estimating the equation

$$income = \beta_0 + \beta_2 st_2 + \beta_3 st_3 + \beta_4 st_4 + \beta_5 st_5 + \beta_6 st_6 + u \quad (1)$$

where I have suppressed the observation subscripts.

What are the regression coefficients in this case? β_0 is the average income in the 1st state: the dummy for which is excluded from the regression. β_2 is the difference between the income in state 2 and the income in state 1. β_3 is the difference between the income in state 3 and the income in state 1, and so on.

What is the ordinary “ANOVA F” in this context—the test that all the slopes are equal to zero? Precisely the test of the null hypothesis:

$$H_0 : \mu_1 = \mu_2 = \mu_3 = \mu_4 = \mu_5 = \mu_6 \quad (2)$$

versus the alternative that not all six of the state means are the same value.

It turns out that we can test this same hypothesis by excluding any one of the dummies, and including the remaining five in the regression. The coefficients will differ, but the p -value of the ANOVA F will be identical for any of these regressions. In fact, this regression is an example of “classical one-way ANOVA”: testing whether a qualitative factor (in this case, state of residence) explains a significant fraction of the variation in income.

What if we wanted to generate point and interval estimates of the state means of income? We could reformulate the model to include all 6 dummies and exclude the constant term, or more usefully, we could just use the `margins` command:

```
regress income i.state  
margins state
```

which will give us the point and interval estimates for each state.

What if we fail to reject the ANOVA F null? Then it appears that the qualitative factor “state” does not explain a significant fraction of the variation in income. Perhaps the relevant classification is between northern, more rural New England states (NEN) and southern, more populated New England states (NES).

Given the nature of dummy variables, we may generate these dummies two ways. We can express the Boolean condition in terms of the `state` variable: `gen nen = (state=="VT" | state=="NH" | state=="ME")`. This expression, with parens on the right hand side of the `generate` statement, evaluates that expression and returns true (1) or false (0). The vertical bar (|) is Stata's OR operator; since every person in the sample lives in one and only one state, we must use OR to phrase the condition that they live in northern New England.

But there is another way to generate this nen dummy, given that we have $st1 \dots st6$ defined for the regression above. Let's say that Vermont, New Hampshire and Maine have been coded as $st6$, $st4$ and $st3$, respectively. We may just $gen\ nen = st3 + st4 + st6$, since the sum of mutually exclusive and exhaustive dummies must be another dummy.

To check, the resulting nen will have a mean equal to the percentage of the sample that live in northern New England; the equivalent nes dummy will have a mean for southern New England residents; and the sum of those two means must be 1.

We can then run a simplified form of our model as `regress inc nen`. That regression's ANOVA F statistic for that regression tests the null hypothesis that incomes in northern and southern New England do not differ significantly. Since we have excluded `nes`, the coefficient on `nen` measures the amount by which northern New England income differs from southern New England income. The mean income for southern New England is the constant term.

If we want point and interval estimates for those means, we should

```
regress income i.nen
margins nen
```

Regression with continuous and dummy variables

In the above examples, we have estimated “pure ANOVA” models: regression models in which *all* of the explanatory variables are dummies. In econometric research, we often want to combine quantitative and qualitative information, including some regressors that are measurable and others that are dummies.

Consider the simplest example: we have data on individuals’ wages, years of education, and their gender. We could create two gender dummies, male and female, but we will only need one in the analysis: say, female. We create this variable as

```
gen female = (gender=="F"),  
or use the factor variable i.female.
```

We can then estimate the model:

$$wage = \beta_0 + \beta_1 educ + \beta_2 female + u \quad (3)$$

The constant term in this model now becomes the wage for a male with zero years of education. Male wages are predicted as $b_0 + b_1 educ$, while female wages are predicted as $b_0 + b_1 educ + b_2$. The gender differential is thus b_2 .

What is this model saying about wage structure? Wages are a linear function of the years of education. If b_2 is significantly different than zero, then there are two “wage profiles”: parallel lines in $educ, wage$ space, each with a slope of b_1 , with their intercepts differing by b_2 .

How would we test for the existence of “statistical discrimination”: e.g., that females with the same qualifications are paid a lower wage? This would be $H_0 : \beta_2 \geq 0$. The t -statistic for b_2 will provide us with this hypothesis test, which we might conduct as a one-tailed test.

If we have priors about the sign of the coefficient, a one-tailed test will allow us to test this hypothesis more effectively, as the reported p-value will be halved if the estimated coefficient is in the rejection region (in this case, if $b_2 < 0$).

We might question the parallel lines assumption inherent in this model. If there is gender-based discrimination in the labor market, it could take the form of a different intercept for men and women, or a different return to education for men and women, or both. We can allow for this by creating an *interaction term* between gender and education:

$$wage = \beta_0 + \beta_1 educ + \beta_2 female + \beta_3 female \times educ + u \quad (4)$$

Although you could generate this interaction term yourself, using either arithmetic or Boolean logic, it is best to let Stata generate it using the *interaction operator* (#) and the `c.` prefix on the continuous variable education. The model to be estimated then becomes

```
regress wage c.educ i.female c.educ#i.female
```

What if we wanted to expand the original model to consider the possibility that wages differ by both gender and race? Say that each worker is classified as $\text{race}=1$ (white) or $\text{race}=2$ (black). Then we could just add the factor variable $i.\text{race}$ to the specification:

$$\text{wage} = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{female} + \beta_3 \text{black} + u$$

What, now, is the constant term? The wage for a white male with zero years of education. Is there a significant race differential in wages? If so, the coefficient b_3 , which measures the difference between white and black wages, *cet. par.*, will be significantly different from zero.

In educ , wage space, the model can be represented as four parallel lines, with each intercept labelled by a combination of gender and race.

What if our racial data classified each worker as white, Black or Asian?
Then we would run the regression:

$$wage = \beta_0 + \beta_1 educ + \beta_2 female + \beta_3 Black + \beta_4 Asian + u \quad (5)$$

Using factor variables, we could specify the model as

```
regress wage educ i.female i.race
```

where the constant term still refers to a white male. In this model, b_3 measures the difference between black and white wages, ceteris paribus, while b_4 measures the difference between Asian and white wages. Each can be examined for significance.

How can we determine whether the qualitative factor `race` affects wages? That is a joint test, that both $\beta_3 = 0$ and $\beta_4 = 0$, and should be conducted as such. If factor variables were used, we could do this with

```
testparm i.race
```

No matter how the equation is estimated, we should not make judgments based on the individual dummies' coefficients, but should rather include both race variables if the null is rejected, or remove them both if it is not.

When we examine a qualitative factor, which may give rise to a number of dummy variables, they should be treated as a group.

For instance, we might want to modify (3) to consider the effect of state of residence:

$$wage = \beta_0 + \beta_1 educ + \beta_2 female + \sum_{j=2}^6 \gamma_j st_j + u \quad (6)$$

where we include any 5 of the 6 `st` variables designating the New England states. The test that wage levels differ significantly due to state of residence is the joint test that $\gamma_j = 0$, $j = 2, \dots, 6$ (or, if factor variables are used, `testparm i.state`). A judgment concerning the relevance of state of residence should be made on the basis of this joint test (an F-test with 5 numerator degrees of freedom).

Note that if the dependent variable was measured in log form, the coefficients on dummies would be interpreted as percentage changes. If (6) was respecified to place $\log(\textit{wage})$ as the dependent variable, the coefficient b_1 would measure the percentage return to education (how many percent does the wage change for each additional year of education), while the coefficient b_2 would measure the (approximate) percentage difference in wage levels between females and males, *ceteris paribus*. The state dummies would, likewise, measure the percentage difference in wage levels between that state and the excluded state (state 1).

We must be careful when working with variables that have an ordinal interpretation, and are thus coded in numeric form, to treat them as ordinal. For instance, if we model the interest rate corporations must pay to borrow (*corprt*) as a function of their credit rating, we consider that Moody's and Standard and Poor's assign credit ratings somewhat like grades: AAA, AA, A, BAA, BA, B, C, et cetera. Those could be coded as 1,2,...,7. Just as we can agree that an "A" grade is better than a "B", a triple-A bond rating results in a lower borrowing cost than a double-A rating.

But while GPAs are measured on a clear four-point scale, the bond ratings are merely ordinal, or ordered: everyone agrees on the rating scale, but the differential between *AA* borrowers' rates and *A* borrowers' rates might be much smaller than that between *B* and *C* borrowers' rates: especially the case if *C* denotes “below investment grade”, which will reduce the market for such bonds. Thus, although we might have a numeric index corresponding to *AAA...C*, we should not assume that $\partial \text{corp}rt / \partial \text{index}$ is constant; we should not treat *index* as a cardinal measure.

Clearly, the appropriate way to proceed is to create dummy variables for each rating class, and include all but one of those variables in a regression of *corprt* on bond rating and other relevant factors. For instance, if we leave out the *AAA* dummy, all of the ratings class dummies' coefficients will then measure the degree to which those borrowers' bonds bear higher rates than those of *AAA* borrowers. But we could just as well leave out the *C* rating class dummy, and measure the effects of ratings classes relative to the worst credits' cost of borrowing.

Interactions involving dummy variables

Just as continuous variables may be interacted in regression equations, so can dummy variables. In the NLSW88 dataset (`sysuse nlsw88`), we have one dummy variable indicating respondents' marital status (*married*) and another indicating whether they belong to a union (*it union*). We could regress their *wage* on these two dummies:

$$wage = b_0 + b_1 union + b_2 married + u$$

```
. reg wage union married
```

Source	SS	df	MS
Model	809.695264	2	404.847632
Residual	31803.7471	1875	16.9619985
Total	32613.4424	1877	17.3753023

```
Number of obs = 1878
F( 2, 1875) = 23.87
Prob > F = 0.0000
R-squared = 0.0248
Adj R-squared = 0.0238
Root MSE = 4.1185
```

wage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
union	1.448355	.2211235	6.55	0.000	1.014681	1.882029
married	-.3705046	.1996102	-1.86	0.064	-.7619862	.0209769
_cons	7.450975	.1719857	43.32	0.000	7.113671	7.788278

This gives rise to the following classification of mean wages, conditional on the two factors, which is thus a classic “two-way ANOVA” setup:

	<i>nonunion</i>	<i>union</i>
<i>unmarried</i>	b_0	$b_0 + b_1$
<i>married</i>	$b_0 + b_2$	$b_0 + b_1 + b_2$

We assume that the two effects, union membership and marital status, have independent effects on the dependent variable. Why? Because this joint distribution is modelled as the product of the marginals. What is the difference between union and nonunion wages? b_1 , irrespective of marital status. What is the difference between unmarried and married wages? b_2 , irrespective of union membership.

If we were to relax the assumption that union membership and marital status had independent effects on wages, we would want to consider their *interaction*. As there are only two categories of each variable, we only need one interaction term, *um*, to capture the possible effects.

That term could be generated as a Boolean (noting that `&` is Stata's AND operator): `gen um=(union==1) & (married==1)`, or we could generate it algebraically, as `gen um=union*married`. In either case, it represents the intersection of the sets.

The additional term added to the estimated equation, corresponding to the interaction, appears as an additive constant in the lower right cell of the table.

If the coefficient on the interaction term is significantly different from zero, the effect of being a union member on the wage differs, depending on marital status, and vice versa. Are the interaction effects important: that is, does the joint distribution meaningfully differ from the product of the marginals? That is easily discerned, as if that is so b_3 will be significantly nonzero.

A much better way to specify this model is to use Stata's factor variables and interaction operators. To interact the `union` and `married` indicators, we can make use of the *factorial interaction* operator:

```
regress wage union married i.union#i.married
```

or, in an even simpler form,

```
regress wage i.union##i.married
```

where the double hash mark indicates the *full factorial* interaction, including both the main effects of each factor and their interaction.

```
. reg wage i.union##i.married
```

Source	SS	df	MS			
Model	811.88412	3	270.62804	Number of obs = 1878		
Residual	31801.5582	1874	16.9698817	F(3, 1874) = 15.95		
Total	32613.4424	1877	17.3753023	Prob > F = 0.0000		
				R-squared = 0.0249		
				Adj R-squared = 0.0233		
				Root MSE = 4.1195		

wage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
union	1.550294	.3598365	4.31	0.000	.8445712	2.256016
married	-.3281956	.2318206	-1.42	0.157	-.7828493	.1264581
union#married	-.1638359	.4561839	-0.36	0.720	-1.058518	.730846
_cons	7.422848	.1890134	39.27	0.000	7.052149	7.793547

With either form of the equation using factor variables, we may then use `margins` to summarize the effects of the two factors, in terms of the predicted means for each combination of factors:

```
margins union#married
```

or indeed for each factor level and their interactions:

```
margins union##married
```



```
. margins union##married
```

```
Predictive margins
```

```
Number of obs = 1878
```

```
Model VCE : OLS
```

```
Expression : Linear prediction, predict()
```

	Delta-method				[95% Conf. Interval]	
	Margin	Std. Err.	t	P> t		
union						
nonunion	7.209294	.1094833	65.85	0.000	6.994572	7.424016
union	8.652981	.1926153	44.92	0.000	8.275218	9.030744
married						
single	7.803405	.1612101	48.41	0.000	7.487235	8.119575
married	7.434992	.1179321	63.04	0.000	7.2037	7.666284
union#married						
nonunion#single	7.422848	.1890134	39.27	0.000	7.052149	7.793547
nonunion #						
married	7.094653	.134219	52.86	0.000	6.831418	7.357887
union#single	8.973142	.3061964	29.31	0.000	8.37262	9.573664
union#married	8.48111	.2461843	34.45	0.000	7.998286	8.963935

An extension of this framework: considering two factors' effects, imagine that instead of marital status we consider *race = white, Black, other*. To run the model without interactions, we would include two of these dummies in the regression: e.g., *Black, other*; the constant term would be the mean wage of a white non-union member (the excluded class).

What if we wanted to include interactions? Then we would define u_{Black} and u_{other} , and include those two regressors as well. The test for the significance of interactions is now a joint test that these two coefficients are jointly zero.

It is much easier to estimate this model using factor variables:

```
regress wage i.union##i.race
```

where the factorial interaction includes all race categories, both in levels and interacted with the union dummy.

```
. regress wage i.union##i.race
```

Source	SS	df	MS			
Model	1531.00192	5	306.200385	Number of obs = 1878		
Residual	31082.4404	1872	16.6038678	F(5, 1872) = 18.44		
Total	32613.4424	1877	17.3753023	Prob > F = 0.0000		
				R-squared = 0.0469		
				Adj R-squared = 0.0444		
				Root MSE = 4.0748		

wage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
union	1.153829	.2660411	4.34	0.000	.6320603	1.675597
race						
black	-1.614053	.2514712	-6.42	0.000	-2.107247	-1.12086
other	1.881194	1.026421	1.83	0.067	-.1318556	3.894244
union#race						
union#black	1.492629	.4776786	3.12	0.002	.5557899	2.429467
union#other	-3.140969	1.784377	-1.76	0.079	-6.640547	.3586095
_cons	7.5821	.1256907	60.32	0.000	7.335591	7.828608

We can also request that `margins` compute the derivatives of the regression function with respect to each of the factors:

```
. margins, dydx(*)
```

```
Average marginal effects
```

```
Number of obs = 1878
```

```
Model VCE : OLS
```

```
Expression : Linear prediction, predict()
```

```
dy/dx w.r.t. : 1.union 2.race 3.race
```

	dy/dx	Delta-method Std. Err.	t	P> t	[95% Conf. Interval]	
union						
union	1.511882	.2201069	6.87	0.000	1.080201	1.943562
race						
black	-1.247652	.2143378	-5.82	0.000	-1.668018	-.8272859
other	1.110168	.8533268	1.30	0.193	-.5634038	2.78374

Note: dy/dx for factor levels is the discrete change from the base level.

These marginal effects take the interaction terms into account as well.

Analysis of covariance models

What if we want to consider a regular regression, on quantitative variables, but want to allow for different slopes (as well as intercepts) for different categories of observations? Then we create interaction effects between the dummies that define those categories and the measured variables. For instance,

$$wage = b_0 + b_1 married + b_2 tenure + b_3 (married \times tenure) + u$$

Here, we are in essence estimating two separate regressions in one: a regression for single women, with an intercept of b_0 and a slope of b_2 , and a regression for married women, with an intercept of $(b_0 + b_1)$ and a slope of $(b_2 + b_3)$.

Why would we want to do this? We could clearly estimate the two separate regressions, but if we did that, we could not answer the questions:

(a) do single and married women have the same intercept?

(b) Do they have the same slope, or return to one more year of experience?

If we use interacted dummies, we can run one regression, and test all of the special cases of this model which are nested within: that the slopes are the same, that the intercepts are the same, and the “pooled” case in which we need not distinguish between single and married women. Since each of these special cases merely involves restrictions on this general form, we can run this equation and then just conduct the appropriate tests.

This can be easily done with factor variables as

```
regress wage i.married##c.tenure
```

where we *must* use the `c.` operator to tell Stata that `tenure` is to be treated as a continuous variable, rather than considering all possible levels of that variable in the dataset.

. regress wage i.married#c.tenure

Source	SS	df	MS
Model	2478.69035	3	826.230118
Residual	71623.1373	2227	32.1612651
Total	74101.8276	2230	33.2295191

Number of obs = 2231
 F(3, 2227) = 25.69
 Prob > F = 0.0000
 R-squared = 0.0334
 Adj R-squared = 0.0321
 Root MSE = 5.6711

wage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
married					
married	-.079158	.369911	-0.21	0.831	-.8045644 .6462484
tenure	.2184467	.0349072	6.26	0.000	.1499926 .2869008
married#c.tenure					
married	-.0556633	.0447069	-1.25	0.213	-.1433349 .0320083
_cons	6.745483	.2965052	22.75	0.000	6.164027 7.326938

. margins, dydx(*)

Average marginal effects

Number of obs = 2231

Model VCE : OLS

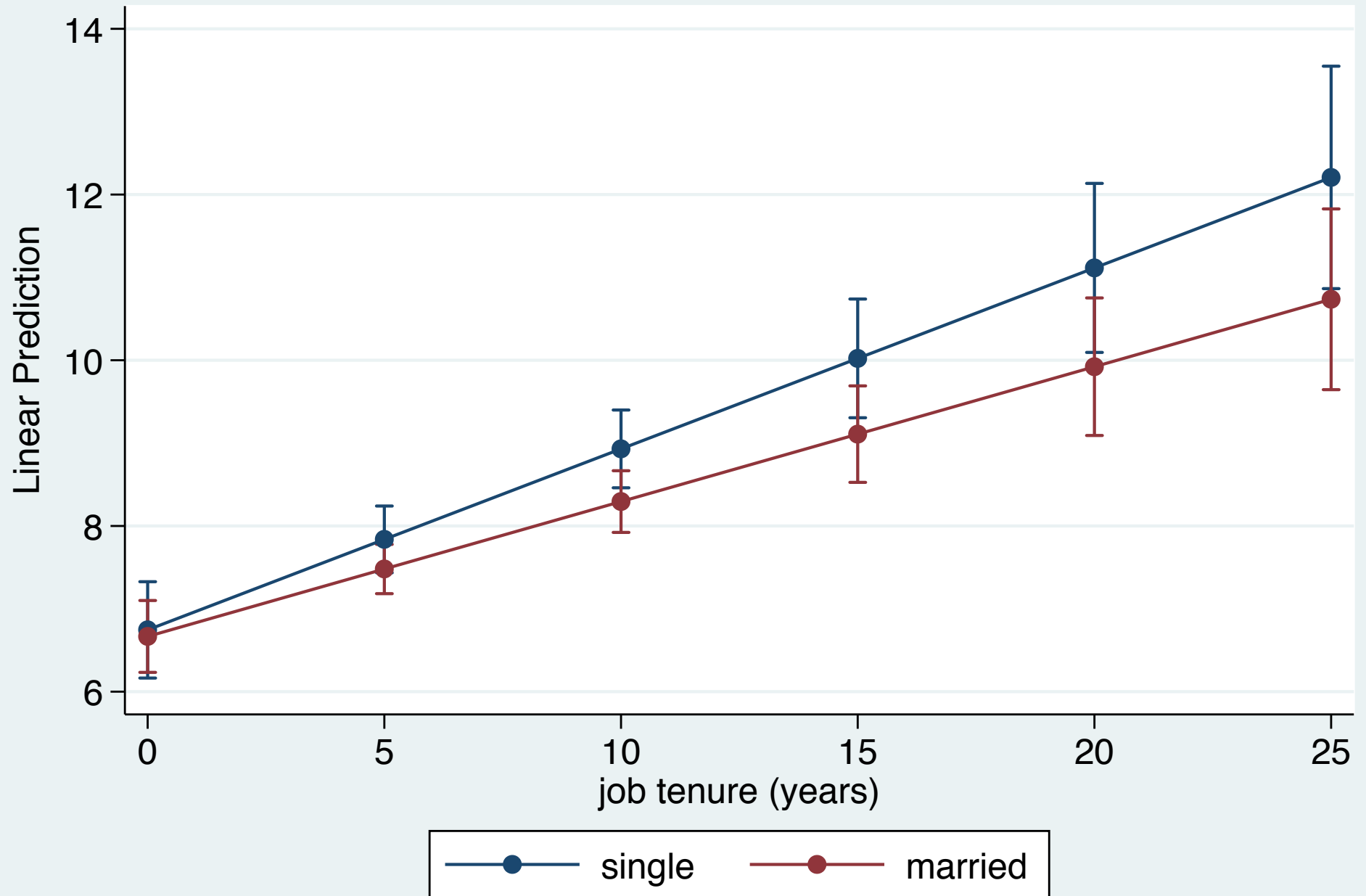
Expression : Linear prediction, predict()

dy/dx w.r.t. : 1.married tenure

	Delta-method					[95% Conf. Interval]	
	dy/dx	Std. Err.	t	P> t			
married							
married	-.4119048	.2506443	-1.64	0.100	-.9034257	.0796161	
tenure	.1827184	.0218568	8.36	0.000	.1398566	.2255802	

Note: dy/dx for factor levels is the discrete change from the base level.

Predictive Margins with 95% CIs



If we extended this logic to include *race*, as defined above, as an additional factor, we would include two of the race dummies (say, *Black* and *other*) and interact each with *tenure*.

This would be a model without interactions, where the effects of marital status and race are considered to be independent, but it would allow us to estimate different regression lines for each combination of marital status and race, and test for the importance of each factor.

```
. margins, dydx(*)
```

```
Average marginal effects
```

```
Number of obs = 2231
```

```
Model VCE : OLS
```

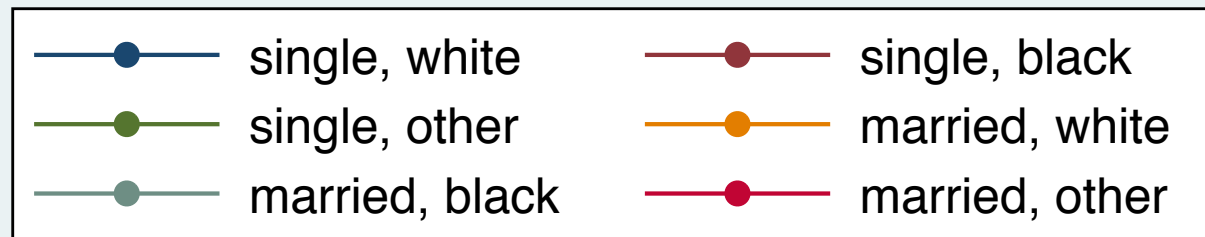
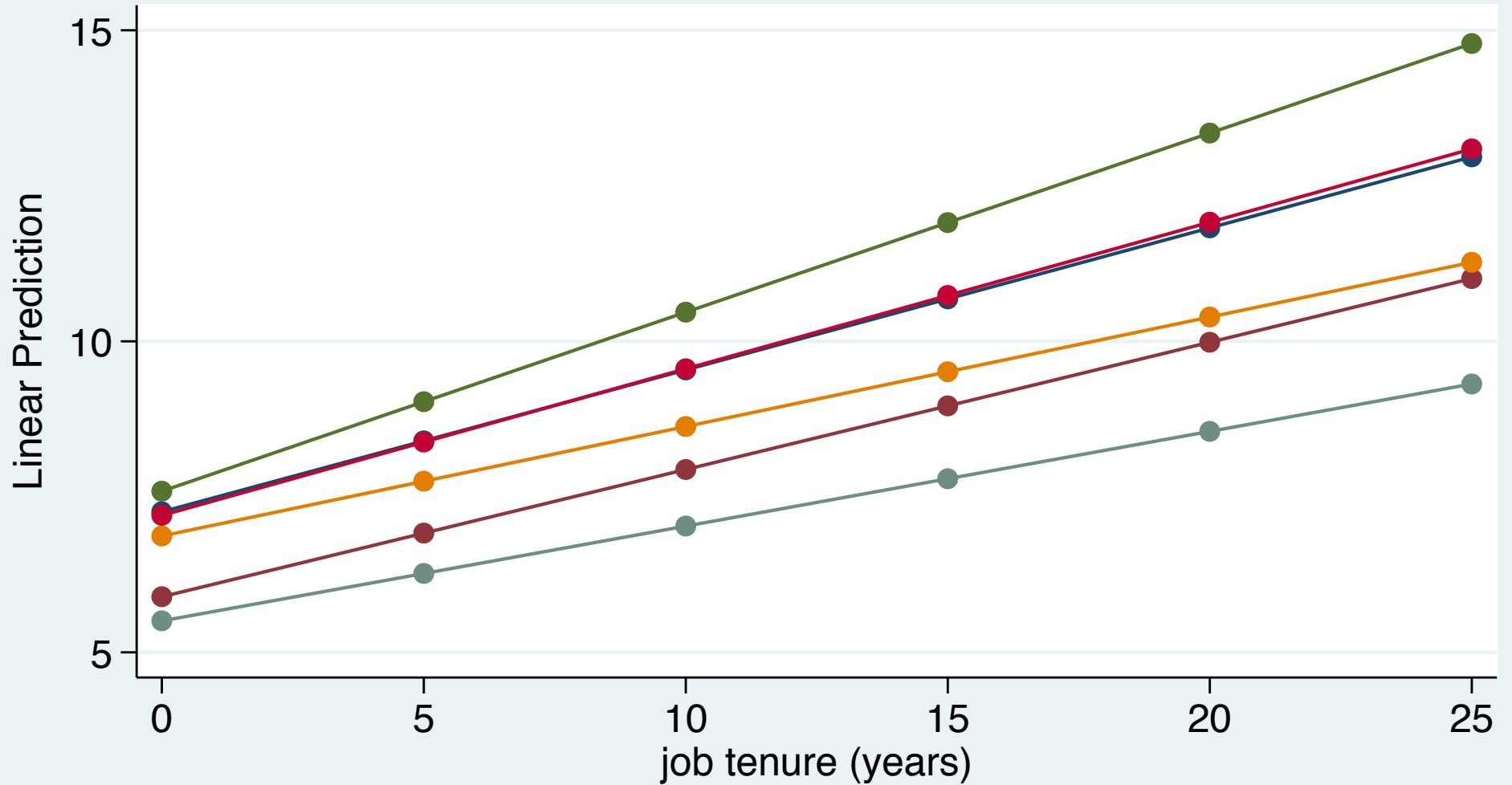
```
Expression : Linear prediction, predict()
```

```
dy/dx w.r.t. : 1.married 2.race 3.race tenure
```

	dy/dx	Delta-method Std. Err.	t	P> t	[95% Conf. Interval]	
married						
married	-.7001233	.2551511	-2.74	0.006	-1.200483	-.1997639
race						
black	-1.506149	.2805171	-5.37	0.000	-2.056252	-.9560463
other	.6878529	1.136602	0.61	0.545	-1.541059	2.916765
tenure	.1892636	.0217633	8.70	0.000	.1465852	.231942

Note: dy/dx for factor levels is the discrete change from the base level.

Predictive Margins



These interaction methods are often used to test hypotheses about the importance of a qualitative factor. For instance, in a sample of companies from which we are estimating their profitability, we may want to distinguish between companies in different industries, or companies that underwent a significant merger, or companies that were formed within the last decade, and evaluate whether their expenditures on R&D or advertising have the same effects across those categories.

All of the necessary tests involving dummy variables and interacted dummy variables may be easily specified and computed, since models without interacted dummies (or without certain dummies in any form) are merely restricted forms of more general models in which they appear.

The standard “subset F” testing strategy that we have discussed for the testing of joint hypotheses on the coefficient vector may be readily applied in this context. The text describes how a “Chow test” may be formulated by running the general regression, running a restricted form in which certain constraints are imposed, and performing a computation using their sums of squared errors; this computation is precisely that done with Stata’s `test` command.

The advantage of setting up the problem for the `test` command is that any number of tests (e.g. above, for the importance of marital status, or for the importance of race) may be conducted after estimating a single regression; it is not necessary to estimate additional regressions to compute any possible “subset F” test statistic, which is what the “Chow test” is doing.