

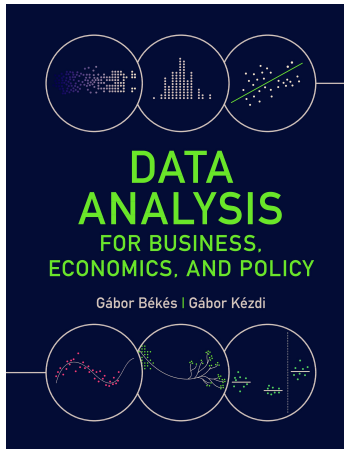
4 Multiple regression

Alice Kügler

Data Analysis 2 – **MS Business Analytics**: Regression Analysis

2023

Slides for the Békés-Kézdi Data Analysis textbook



- ▶ Cambridge University Press, 2021
- ▶ gabors-data-analysis.com
 - ▶ Download all data and code:
gabors-data-analysis.com/data-and-code/
- ▶ These slides are for **Chapter 10**

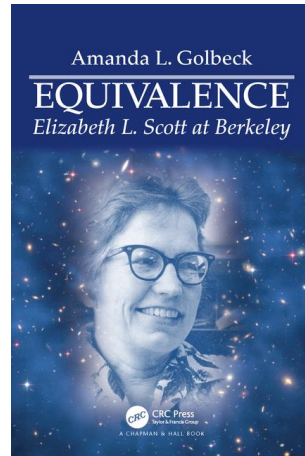
Regression concepts and mechanics

Motivation

- *Interested in finding evidence for or against labor market discrimination of women. Compare wages for men and women who share similarities in wage relevant factors such as experience and education.*
- *Find a good deal on a hotel to spend a night in a European city- analyzed the pattern of hotel price and distance and many other features to find hotels that are underpriced not only for their location but also those other features.*

Motivation II

- ▶ Elizabeth Scott, a Berkeley statistics professor spent two decades analysing inequalities in academic salaries and advocating for change.
- ▶ "How one woman used regression to influence the salaries of many" by Amanda Golbeck (in *Significance*, Dec. 2017)
 - ▶ <http://onlinelibrary.wiley.com/doi/10.1111/j.1740-9713.2017.01092.x/full>



Topics to cover

- ▶ Multiple regression mechanics
- ▶ Estimation and interpreting coefficients
- ▶ Non-linear terms, interactions
- ▶ Variable selection, small sample problems
- ▶ Multiple regression and causality
- ▶ Multiple regression and prediction

Multiple regression analysis

- ▶ Multiple regression analysis uncovers average y as a function of more than one x variable: $y^E = f(x_1, x_2, \dots)$.
- ▶ It can lead to better predictions \hat{y} by considering more explanatory variables.
- ▶ It may improve the interpretation of slope coefficients by comparing observations that are different in terms of one of the x variables but similar in terms of other x variables.
- ▶ Multiple linear regression specifies a linear function of the explanatory variables for the average y .

$$y^E = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots \beta_k x_k$$

Multiple regression: two x variables

$$y^E = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \quad (1)$$

- ▶ β_1 -the slope coefficient on x_1 shows difference in average y across observations with different values of x_1 , *but the same value of x_2* .
 - ▶ β_2 shows difference in average y across observations with different values of x_2 , *but the same value of x_1* .
- ▶ Can compare observations that are similar in one explanatory variable to see the differences related to the other explanatory variable.

Multiple regression - mechanics

Compare slope coefficient in simple (β) and multiple regression (β_1):

$$y^E = \alpha + \beta x_1 \quad (2)$$

$$y^E = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \quad (3)$$

Intermediary step: regression of x_2 on x_1 ("x - x regression") - δ is slope parameter:

$$x_2^E = \gamma + \delta x_1 \quad (4)$$

Multiple regression - mechanics

Compare slope coefficient in simple (β) and multiple regression (β_1):

$$y^E = \alpha + \beta x_1 \quad (2)$$

$$y^E = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \quad (3)$$

Intermediary step: regression of x_2 on x_1 ("x - x regression") - δ is slope parameter:

$$x_2^E = \gamma + \delta x_1 \quad (4)$$

Plug this back

$$y^E = \beta_0 + \beta_1 x_1 + \beta_2 (\gamma + \delta x_1) = \underbrace{\beta_0 + \beta_2 \gamma}_{\text{constant}} + \underbrace{(\beta_1 + \beta_2 \delta)}_{\text{slope}} x_1 . \quad (5)$$

It turns out that

$$\beta - \beta_1 = \delta \beta_2 \quad (6)$$

Omitted variable bias

$$\beta - \beta_1 = \delta\beta_2 \quad (7)$$

- ▶ This $\beta - \beta_1$ difference is called the omitted variable bias (OVB).
- ▶ If we are interested in the coefficient on x_1 , with having x_2 in the regression ($y^E = \beta_0 + \beta_1 x_1 + \beta_2 x_2$), the regression $y^E = \alpha + \beta x_1$ is an incorrect regression as it omits x_2 .
- ▶ Thus, the results from that first regression are biased
 - ▶ Bias is caused by omitting x_2
 - ▶ More: in Chapter 21, Section 21.3. (DA4)

Multiple regression - mechanics

- ▶ The slope of x_1 in a simple regression is different from its slope in the multiple regression, the difference being the product of its slope in the regression of x_2 on x_1 and the slope of x_2 in the multiple regression.
- ▶ The slope coefficient on x_1 in the two regressions is different

Multiple regression - mechanics

- ▶ The slope of x_1 in a simple regression is different from its slope in the multiple regression, the difference being the product of its slope in the regression of x_2 on x_1 and the slope of x_2 in the multiple regression.
- ▶ The slope coefficient on x_1 in the two regressions is different
 - ▶ unless x_1 and x_2 are uncorrelated ($\delta = 0$) OR
 - ▶ the coefficient on x_2 is zero in the multiple regression ($\beta_2 = 0$).
- ▶ The slope in the simple regression is larger if x_2 and x_1 are positively correlated and β_2 is positive
 - ▶ or x_2 and x_1 are negatively correlated and β_2 is negative

Quick example

- ▶ Time series regression, sales and prices of beer, regress month-to-month change in log quantity sold (y) on month-to-month change in log price (x_1).
- ▶ $\beta = -0.5$: sales tend to decrease by 0.5% when our price increases by 1%.
 - ▶ Time series regression hence the increase/decrease
- ▶ Extended regression, x_2 : change in ln average price charged by our competitors:
 $\hat{\beta}_1 = -3$ and $\hat{\beta}_2 = 3$

Quick example

- ▶ Time series regression, sales and prices of beer, regress month-to-month change in log quantity sold (y) on month-to-month change in log price (x_1).
- ▶ $\beta = -0.5$: sales tend to decrease by 0.5% when our price increases by 1%.
 - ▶ Time series regression hence the increase/decrease
- ▶ Extended regression, x_2 : change in ln average price charged by our competitors:
 $\hat{\beta}_1 = -3$ and $\hat{\beta}_2 = 3$
- ▶ So: $\hat{\beta} - \hat{\beta}_1 = -0.5 - (-3) = +2.5$. = simple regression gives a biased estimate of the slope coefficient
- ▶ We have $\hat{\delta}_1 = 0.83$ so $\hat{\delta}_1 \times \hat{\beta}_2 = 0.83 * 3 = 2.5$
- ▶ Positive bias = result of two things:
 - ▶ a positive association between the two price changes (δ) and
 - ▶ a positive association between competitor price and our own sales (β_2).

Multiple regression - mechanics

- ▶ If x_1 and x_2 are correlated, comparing observations with or without the same x_2 value makes a difference.
- ▶ If they are positively correlated, observations with higher x_2 tend to have higher x_1 .
- ▶ In the simple regression we ignore differences in x_2 and compare observations with different values of x_1 .
- ▶ But higher x_1 values mean higher x_2 values, too.
- ▶ Corresponding differences in y may be due to differences in x_1 but also differences in x_2 .

Multiple regression - some language

- ▶ Multiple regression with two explanatory variables (x_1 and x_2),
- ▶ We measure differences in expected y across observations that differ in x_1 but are similar in terms of x_2 .
- ▶ Difference in y by x_1 , **conditional on** x_2 . OR **controlling for** x_2 .
- ▶ We condition on x_2 , or control for x_2 , when we include it in a multiple regression that focuses on average differences in y by x_1 .

Multiple regression - some language

- ▶ Multiple regression with two explanatory variables (x_1 and x_2),
- ▶ When we are interested in a regression with x_2 , but we have one without: x_2 is an **omitted variable** in the simple regression.
- ▶ The slope on x_1 in the sample is confounded by omitting the x_2 variable, and thus x_2 is a **confounder**.

Multiple regression - mechanics – SE

- Inference, confidence intervals in multiple regressions is analogous to those in simple regressions.

$$SE(\hat{\beta}_1) = \frac{Std[e]}{\sqrt{n}Std(x_1)\sqrt{1 - R_1^2}} \quad (8)$$

- Same: the SE is small - small Std of the residuals (the better the fit of the regression); large sample, large the Std of x_1 .
- New: $\sqrt{1 - R_1^2}$ term in the denominator - the R-squared of the regression of x_1 on x_2 - correlation between x_1 and x_2 .
- The stronger the correlation between x_1 and x_2 the larger the SE of $\hat{\beta}_1$.
- Note the symmetry: the same would apply to the SE of $\hat{\beta}_2$.
- Also: in practice, use robust SE

Multiple regression - mechanics – collinearity

- **Perfect collinearity** is when x_1 is a linear function of x_2 .

Multiple regression - mechanics – collinearity

- ▶ **Perfect collinearity** is when x_1 is a linear function of x_2 .
- ▶ Consequence: cannot calculate coefficients.
 - ▶ One will be dropped by the software
- ▶ Strong but imperfect correlation between explanatory is sometimes called **multicollinearity**.
- ▶ Consequence: We can get the slope coefficients and their standard errors,
 - ▶ The standard errors may be large.

Multicollinearity and SE of beta

- ▶ As a consequence of multicollinearity the standard errors may be large.
 - ▶ Concept: Few variables that are different in x_1 but not in x_2 . Not enough observations for comparing average y when x_1 is different but x_2 remains the same.
 - ▶ Math: R_1^2 is high (x_2 is a good predictor of x_1), thus $\sqrt{1 - R_1^2}$ is (really) small, which makes $SE(\beta_1)$ (very) large.

- ▶ This is a typically a small sample problem (hundreds to few thousand) of observations
 - ▶ May look at pair-wise correlations when start working with data.
 - ▶ Drop one or the other, or combine them (use z-score/average).

Multiple regression - joint significance

- ▶ *Testing joint hypotheses*: null hypotheses that contain statements about more than one regression coefficient.
- ▶ We aim at testing whether a subset of the coefficients (such as all geographical variables) are all zero.
- ▶ F-test answers this.
 - ▶ Individually they are not all statistically different from zero, but together they may be.
- ▶ We may ask if *all slope coefficients are zero* in the regression.
- ▶ "Global F-test", and its results are often shown by statistical software by default. Don't use it, R-squared is fine.

Multiple regression - many explanatory variables

- Having more explanatory variables is no problem.

$$y^E = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_k x_k \quad (9)$$

- Interpreting the slope of x_1 : on average, y is β_1 units larger in the data for observations with one unit larger x_1 but the same value for all other x variables.
- SE formula - small when R_k^2 is small - R^2 of regression of x_k on all *other* x variables.

$$SE(\hat{\beta}_k) = \frac{Std[e]}{\sqrt{n} Std[x_k] \sqrt{1 - R_k^2}} \quad (10)$$

Multiple regression collinearity 2

- ▶ Multicollinearity when many variables
- ▶ Can tabulate pairwise correlations
 - ▶ Do a "heatmap"
- ▶ Drop (combine) strongly correlated variables

Multiple regression collinearity 2

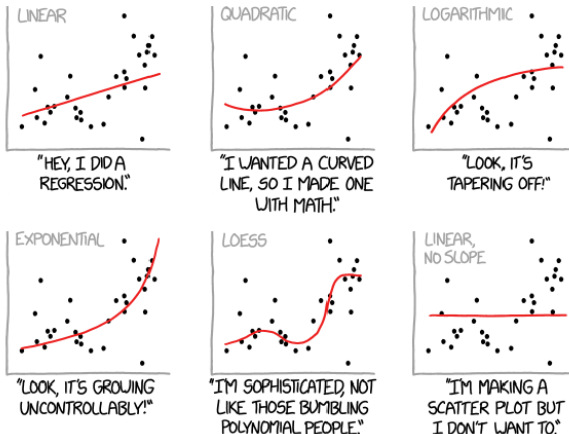
- ▶ Multicollinearity when many variables
- ▶ Can tabulate pairwise correlations
 - ▶ Do a "heatmap"
- ▶ Drop (combine) strongly correlated variables
 - ▶ Strongly: depends on sample size
- ▶ It may be a bit more complicated: some variables may be a function of several others.
- ▶ Check R^2 of regressions like $x_1^E = \alpha + \beta_1 x_2 + \beta_2 x_3 \cdots + \beta_k x_k$
- ▶ Drop variables if very high R^2

Multiple regression - non-linear patterns

- ▶ Uses splines, polynomials - actually like multiple regression - we have multiple coefficient estimates
- ▶ No fear of multicollinearity - not *linear* combinations
- ▶ Non-linear function of various x_i variables may be combined.

Multiple regression - non-linear patterns

CURVE-FITTING METHODS AND THE MESSAGES THEY SEND



Understanding the gender difference in earnings

- ▶ In the USA (2014), women tend to earn about 20% less than men
- ▶ Aim 1: Find patterns to better understand the gender gap. Our focus is the interaction with age.
- ▶ Aim 2: Think about if there is a causal link from being female to getting paid less.

The gender difference in earnings

- ▶ 2014 census data
 - ▶ Age between 15 to 65
 - ▶ Exclude self-employed (earnings is difficult to measure)
 - ▶ Include those who reported 20 hours more as their usual weekly time worked
- ▶ Employees with a graduate degree (higher than 4-year college)
- ▶ Use log hourly earnings ($\ln w$) as dependent variable
- ▶ Use gender and add age as explanatory variables

The gender difference in earnings: regression

We are quite familiar with the relation between earnings and gender:

$$\ln w^E = \alpha + \beta \text{female}, \quad \beta < 0$$

Let's include age as well:

$$\ln w^E = \beta_0 + \beta_1 \text{female} + \beta_2 \text{age}$$

We can calculate the correlation between female and age, which is in fact negative.

What do you expect about β, β_1, δ ?

Reminder:

$$\text{age}^E = \gamma + \delta \text{female}$$

The gender gap regression - baseline

Variables	(1) ln wage	(2) ln wage	(3) age
female	-0.195** (0.008)	-0.185** (0.008)	-1.484** (0.159)
age		0.007** (0.000)	
Constant	3.514** (0.006)	3.198** (0.018)	44.630** (0.116)
Observations	18,241	18,241	18,241
R-squared	0.028	0.046	0.005

Note: All employees with a graduate degree. Robust standard errors in parentheses *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$

The gender gap regression - interpretation

$$\beta - \beta_1 = \delta\beta_2$$

which can be calculated easily:

- ▶ $\beta - \beta_1 = -0.195 - (-0.185) = -0.01$
- ▶ $\delta\beta_2 = -1.48 \times 0.007 \approx -0.01$

Interpretation:

- ▶ Age is a confounder, it is different from zero and the beta coefficient changes.
- ▶ But a weak one.

The gender gap regression - interpretation 2

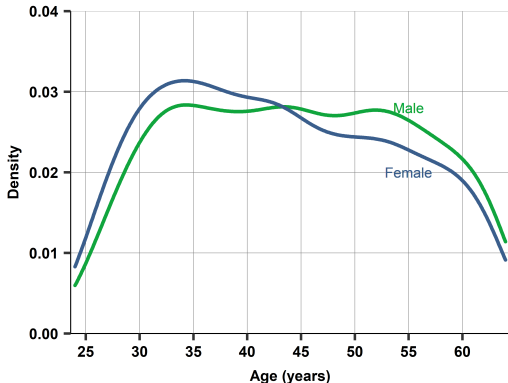
- ▶ Women of the same age have a slightly smaller earnings disadvantage in this data because they are somewhat younger, on average
- ▶ Employees that are younger tend to earn less
- ▶ Part of the earnings disadvantage of women is thus due to the fact that they are younger.
 - ▶ This is a small part: around 1 percentage points of the 20% difference,
 - ▶ = a 5% share of the entire difference.

The gender gap regression - back to modeling

- ▶ A single linear interaction may not be enough.
- ▶ Next: drill down the impact of age

Age distribution

Age distribution of male and female employees with degrees higher than college



- ▶ Relatively few below age 30
- ▶ Above 30
 - ▶ close to uniform for men
 - ▶ for women, the proportion of female employees with graduate degrees drops above age 45, and again, above age 55
- ▶ Two possible things
 - ▶ fewer women with graduate degrees among the 45+ old than among the younger ones
 - ▶ fewer of them are employed

Rerun regression with age in non-linear way

Variables	(1) ln wage	(2) ln wage	(3) ln wage	(4) ln wage
female	-0.195** (0.008)	-0.185** (0.008)	-0.183** (0.008)	-0.183** (0.008)
age		0.007** (0.000)	0.063** (0.003)	0.572** (0.116)
age ²			-0.001** (0.000)	-0.017** (0.004)
age ³				0.000** (0.000)
age ⁴				-0.000** (0.000)
Constant	3.514** (0.006)	3.198** (0.018)	2.027** (0.073)	-3.606** (1.178)
Observations	18,241	18,241	18,241	18,241
R-squared	0.028	0.046	0.060	0.062

Note: *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$

Gender difference in earnings - conditioning on non-linear age

- ▶ Maybe age as confounder is non-linear.
- ▶ So we extended our analysis by including higher orders of age.

Gender difference in earnings - conditioning on non-linear age

- ▶ Maybe age as confounder is non-linear.
- ▶ So we extended our analysis by including higher orders of age.
- ▶ Not much difference re female variable.
- ▶ However R^2 increases as we include higher orders.
- ▶ Regardless how we model age, the gender gap remains.

Qualitative variables and interactions

Multiple regression - qualitative variables

- Can have binary variables as well as other qualitative variables (factors)
- Consider a qualitative variable like continents. How to add it to the regression model?

Multiple regression - qualitative variables

- ▶ Can have binary variables as well as other qualitative variables (factors)
- ▶ Consider a qualitative variable like continents. How to add it to the regression model?
- ▶ Create binary variables (dummy variables) for all options. Add them - all but one.
- ▶ This one will be the base

Multiple regression - qualitative variables

- ▶ x is a categorical variable with three values *low*, *medium* and *high*
- ▶ binary variable m denote if $x = \text{medium}$, h variable denote if $x = \text{high}$.
- ▶ for $x = \text{low}$ is not included. It is called the *reference category* or left-out category.

$$y^E = \beta_0 + \beta_1 x_{\text{medium}} + \beta_2 x_{\text{high}} \quad (11)$$

Multiple regression - qualitative variables

$$y^E = \beta_0 + \beta_1 x_{med} + \beta_2 x_{high} \quad (12)$$

- ▶ Pick $x = low$ as the reference category. Other values compared to this.
 - ▶ This is the omitted variable
- ▶ β_0 shows average y in the reference category. Here, β_0 is average y when both $x_{med} = 0$ and $x_{high} = 0$: this is when $x = low$.
- ▶ β_1 shows the difference of average y between observations with $x = medium$ and $x = low$
- ▶ β_2 shows the difference of average y between observations with $x = high$ and $x = low$.

Multiple regression - qualitative variables

How to pick a reference category?

- ▶ Substantive guide: choose the category to which we want to compare the rest.
 - ▶ Examples include the home country, the capital city, the lowest or highest value group.
- ▶ The statistical guide: chose a category with a large number of observations.
 - ▶ Important when inference is important.
 - ▶ If reference category has few observations - coefficients will have large SE / wide CI.

Qualitative right-hand-side variables summary

- ▶ Include qualitative right-hand-side variables with k categories as a series of $k - 1$ binary (“dummy”) variables
- ▶ Coefficients on each of the $k - 1$ binary variables show average differences in y compared to the reference category

Case study – understanding the gender difference in earnings

- ▶ MA, Professional, and PhD are three categories of graduate degree. Column (2): MA is the reference category. Column (3): the reference category is Professional or PhD.
- ▶ Employees of age 24–65 with a graduate degree and 20 or more work hours per week.

Gender differences in earnings – log earnings, gender and education

Variables	(1) ln wage	(2) ln wage	(3) ln wage
female	-0.195** (0.008)	-0.182** (0.009)	-0.182** (0.009)
ed_Profess		0.134** (0.015)	-0.002 (0.018)
ed_PhD		0.136** (0.013)	
ed_MA			-0.136** (0.013)
Constant	3.514** (0.006)	3.473** (0.007)	3.609** (0.013)
Observations	18,241	18,241	18,241
R-squared	0.028	0.038	0.038

Multiple regression - interactions

- ▶ Many cases, data is made up of important groups: male and female workers or countries in different continents.
- ▶ Some of the patterns we are after may vary across these groups.
- ▶ The strength of a relation may also be altered by a special variable.
 - ▶ In medicine, a *moderator variable* can reduce / amplify the effect of a drug on people.
 - ▶ In business, financial strength can affect how firms may weather a recession.
- ▶ All of these mean different patterns for subsets of observations.

Multiple regression - interactions

- ▶ Regression with two explanatory variables: x_1 is continuous, D is binary denoting two groups in the data (e.g., male or female employees).
- ▶ We wonder if the relationship between average y and x_1 is different for observations with $D = 1$ than for $D = 0$. How?

Multiple regression - qualitative variables: parallel lines

- Option 1: Two *parallel lines* for the $y - x_1$ pattern: one for those with $D = 0$ and one for those with $D = 1$.

$$y^E = \beta_0 + \beta_1 x_1 + \beta_2 D \quad (13)$$

Two groups, $D = 0$ and $D = 1$. Difference is the level

$$y_0^E = \beta_0 + \beta_2 \times 0 + \beta_1 x_1 \quad (14)$$

$$y_1^E = \beta_0 + \beta_2 \times 1 + \beta_1 x_1 \quad (15)$$

Multiple regression - qualitative variables: different slopes

- Option 2: If we want to *allow for different slopes* in the two D groups we have to do something different, add an interaction term.

$$y^E = \beta_0 + \beta_1 x_1 + \beta_2 D + \beta_3 x_1 D \quad (16)$$

Intercepts different by β_2 AND slopes different by β_3 .

$$y_0^E = \beta_0 + \beta_1 x_1 \quad (17)$$

$$y_1^E = \beta_0 + \beta_2 + (\beta_1 + \beta_3) x_1 \quad (18)$$

Multiple regression - separate regressions

- ▶ Separate regressions in the two groups and the regression that pools observations but includes an interaction term yield *exactly the same* coefficient estimates.
- ▶ The coefficients of the separate regressions are easier to interpret.
- ▶ But the pooled regression with interaction allows for a direct test of whether the slopes are the same.

Multiple regression - an extended model

- Extension: D_1 , D_2 are binaries, x continuous:

$$y^E = \beta_0 + \beta_1 D_1 + \beta_2 D_2 + \beta_3 x + \beta_4 D_1 x + \beta_5 D_2 x \quad (19)$$

Interaction with two continuous variable

- ▶ Same model used for two continuous variables, x_1 and x_2 :

$$y^E = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 \quad (20)$$

- ▶ Example: firm level data, 100 industries.
- ▶ y is change in revenue x_1 is change in global demand, x_2 is firm's financial health
- ▶ The interaction can capture that drop in demand can cause financial problems in firms, but less so for firms with better balance sheet.

Interaction between gender and age I

- ▶ Why do we assume that age has the same slope regardless of gender? We might want to check, whether they are different.
- ▶ Are the slopes significantly different?
- ▶ Can one get the slope for age for female only from the regression with the interaction?
- ▶ How is the gender dummy's coefficient changed?

Interaction between gender and age II

- Look at men and women separately.
Earning for men rises faster with age.
- Or, look at them pooled with interaction.
 - Observe that pooling with interaction is the SAME as two separate models.
- Female dummy is close to zero. Does this mean no gender gap?

Variables	(1) Women ln wage	(2) Men ln wage	(3) All ln wage
female			-0.036 (0.035)
age	0.006** (0.001)	0.009** (0.001)	0.009** (0.001)
female × age			-0.003** (0.001)
Constant	3.081** (0.023)	3.117** (0.026)	3.117** (0.026)
Observations	9,685	8,556	18,241
R-squared	0.011	0.028	0.047

Interaction between gender and age: interactions and non-linearities

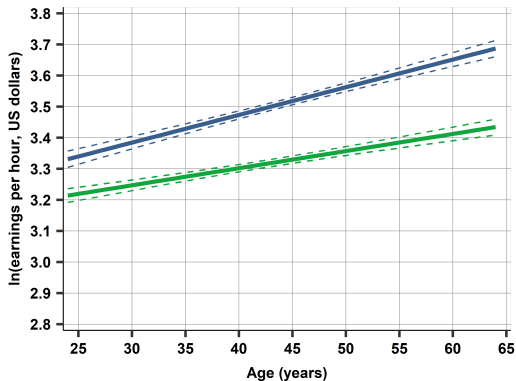
We estimate two models

$$\ln w^E = \beta_0 + \beta_1 \text{age} + \beta_2 \text{female} + \beta_3 \text{female} \times \text{age} \quad (21)$$

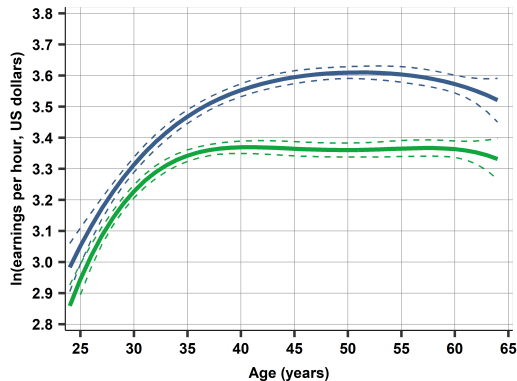
and now add age as non-linear

$$\begin{aligned} \ln w^E = & \beta_0 + \beta_1 \text{age} + \beta_2 \text{age}^2 + \beta_3 \text{age}^3 + \beta_4 \text{age}^4 \\ & + \beta_5 \text{female} + \beta_6 \text{female} \times \text{age} + \beta_7 \text{female} \times \text{age}^2 \\ & + \beta_8 \text{female} \times \text{age}^3 + \beta_9 \text{female} \times \text{age}^4 \end{aligned} \quad (22)$$

Interaction between gender and age: interactions and non-linearities



Log earnings per hour and age by gender: predicted values and confidence intervals from a linear regression interacted with gender.



Log earnings per hour and age by gender: predicted values and confidence intervals from a regression with 4th-order polynomial interacted with gender.

Understanding the gender difference in earnings

- ▶ the average earnings difference is around 10% between ages 25 and 30
- ▶ increases to around 15% by age 40, and reaches 22% by age 50,
- ▶ from where it decreases slightly to age 60 and more by age 65.
- ▶ confidence intervals around the regression curves are rather narrow, except at the two ends.
- ▶ Conclusion?

Conditioning and causality

Multiple regression - causal analysis

- ▶ One main reason to estimate multiple regressions is to get closer to a causal interpretation.
- ▶ By conditioning on other observable variables, we can get closer to comparing similar objects – “apples to apples” – even in observational data.
- ▶ But getting closer is not the same as getting there.
- ▶ In principle, one may help that by conditioning on *every* potential confounder: variables that would affect y and the causal variable x_1 at the same time.
- ▶ Ceteris paribus = conditioning on **every** such relevant variable.

Multiple regression - causal analysis

- ▶ Ceteris paribus = conditioning on **every** such relevant variable.
- ▶ *Ceteris paribus* prescribes what we want to condition on; a multiple regression can condition on **what's in the data** the way it is measured.
- ▶ Importantly, conditioning on everything is impossible in general.
- ▶ Multiple regression is never (hardly ever) ceteris paribus.

Multiple regression - causal analysis

- ▶ In randomized experiments, we use causal language, as treated and untreated units are similar - by random grouping.
- ▶ In observational data, comparisons don't uncover causal relations.
 - ▶ Cautious with language. Do not use "effect" or "increase".
 - ▶ Regression, even with multiple x is just a comparison, a conditional mean.

Multiple regression - causal analysis

- ▶ Not all variables should be included as control variables even if correlated both with the causal variable and the dependent variable.
- ▶ *Bad conditioning variables* are variables that are correlated both with the causal variable and the dependent variable but are actually part of the causal mechanism.
- ▶ This is the reason for excluding them.
- ▶ For example, when we want to see how TV advertising affects sales. Should we control for how many people viewed the advertising?

Multiple regression - causal analysis

- ▶ Not all variables should be included as control variables even if correlated both with the causal variable and the dependent variable.
- ▶ *Bad conditioning variables* are variables that are correlated both with the causal variable and the dependent variable but are actually part of the causal mechanism.
- ▶ This is the reason for excluding them.
- ▶ For example, when we want to see how TV advertising affects sales. Should we control for how many people viewed the advertising?
 - ▶ No. Part of why less advertising may hurt sales - fewer heads.
- ▶ Super hard.
- ▶ More in DA4

Multiple regression - causal analysis

- ▶ A multiple regression on observational data is rarely capable of uncovering a causal relationship.
 - ▶ Cannot capture all potential confounders (not ceteris paribus).
 - ▶ Potential bad conditioning variables (bad controls).
 - ▶ We can never really know. BUT:
- ▶ Multiple regression can get us **closer** to uncovering a causal relationship.
 - ▶ Compare units that are the same in many respects - controls.

Understanding the gender difference in earnings - causal analysis

What may cause the difference in wages?

- ▶ Labor discrimination - one group earns less even if they have the same *marginal product*
- ▶ Try control for marginal product (or for variables which matter to marginal product)
 - ▶ Eg.: occupation (as an indicator for inequality in gender roles), or industry, union status, hours worked and other socio-economic characteristics
- ▶ Use variables as controls - does comparing apples to apple change the coefficient of the female variable?

Understanding the gender difference in earnings -regression

- Restricted sample:
employees of age 40
to 60 with a
graduate degree that
work 20 hours per
week or more
- More and more
confounders added

Variables	(1) ln wage	(2) ln wage	(3) ln wage	(4) ln wage
female	-0.224** (0.012)	-0.212** (0.012)	-0.151** (0.012)	-0.141** (0.012)
Age and education		YES	YES	YES
Family circumstances			YES	YES
Demographic background			YES	YES
Job characteristics			YES	YES
Union member			YES	YES
Age in polynomial				YES
Hours in polynomial				YES
Observations	9,816	9,816	9,816	9,816
R-squared	0.036	0.043	0.182	0.195

Understanding the gender difference in earnings - interpretation

- Interpret the last coefficient

Understanding the gender difference in earnings - interpretation

- ▶ Interpret the last coefficient
- ▶ Let us compare two people, with same age, hours, industry, occupation, geography, background (=confounders) but one is a man and the other is a woman.
- ▶ The women tend to earn 14% less, on average.
 - ▶ On average, women are expected to earn 14% less

Regression table detour

- ▶ Regression table with many x variables is hard to present
- ▶ In presentation, suppress unimportant coefficients
- ▶ In paper, you may present more, but mostly if you want to discuss them or for sanity check
 - ▶ Sanity check: do control variable coefficient make sense by and large?
- ▶ Check N of observations: if the same sample, they should be exactly the same.
- ▶ R^2 is enough, no need for other stuff

Discussion

- ▶ Could not safely pin down the role of labor market discrimination and broader gender inequality
- ▶ Multiple regression - *closer* to causality
- ▶ But, broader gender inequality seems to matter: inequality in gender roles likely plays a role in the fact that women earn less per hour than men
- ▶ We cannot prove that that the remaining 14% is due to discrimination - *plenty of remaining heterogeneity*
- ▶ Also: selection and bad conditioning variables?

Prediction

Multiple regression - prediction and benchmarking

- ▶ Reason to estimate a multiple regression is to make a *prediction*
- ▶ Find the best guess for the dependent variable y_j for a particular *target observation* j

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots$$

- ▶ The $\hat{y}-y$ plot is a good way to visualize the fit of a prediction.
- ▶ It is a scatterplot with \hat{y} on the horizontal axis and y on the vertical axis,
- ▶ Together with the 45 degree line, which is the regression line of y on \hat{y}
 - ▶ observations to the right of the 45 degree line show overpredictions ($\hat{y} > y$)
 - ▶ observations to the left of the 45 degree line show underpredictions ($\hat{y} < y$).

Multiple regression - benchmarking

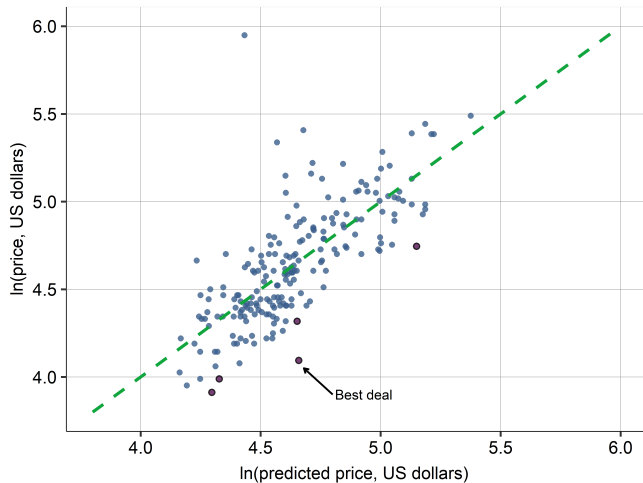
- ▶ When the goal is prediction we want the regression to produce as good a fit as possible.
- ▶ As good a fit as possible in the general pattern that is representative of the target observation j .
- ▶ A common danger is *overfitting* the data: finding patterns in the data that are not true in the general pattern.

Case study – finding a good deal among hotels with multiple regression

- ▶ Regress log price
 - ▶ on distance to the city center (piecewise linear spline with knots at 1 and 4 miles),
 - ▶ average customer rating (piecewise linear spline with knot at 3.5),
 - ▶ binary variables for 3.5 stars and 4 stars (reference category is 3 stars).

- ▶ The process to create a scatterplot and regression line with multiple x
 - ▶ Estimate model for (log) hotel prices, many predictor variables, ($R^2 = 0.56$), and get $\hat{y} - y$
 - ▶ Scatterplot and line: $\hat{y} - y$ plot for log hotel price
 - ▶ Identify good deals?

$\hat{y}-y$ plot for log hotel price



Good deals for hotels: the five hotels with the most negative residuals

List of the five observations with the smallest (most negative) residuals from the multiple regression

Hotel name	Price	Residual in $\ln(\text{price})$	Distance	Stars	Rating
21912	60	-0.565	1.1	4	4.1
21975	115	-0.405	0.1	4	4.3
22344	50	-0.385	3.9	3	3.9
22080	54	-0.338	1.1	3	3.2
22184	75	-0.335	0.7	3	4.1

Multiple regression - variable selection

- ▶ How should one decide which variables to include and how?
- ▶ Depends on the purpose: prediction or causality.
- ▶ Lot of judgment calls to make
 - ▶ Very hard task. No super-duper solution.
- ▶ Non-linear fit - use a non-parametric first and if non-linear, pick a model that is close - quadratic, piecewise spline.
- ▶ If two or many variables strongly correlated, pick one of them. Sample size will help decide.
- ▶ Keep it as simple as possible.
- ▶ Key topic for DA3

Multiple regression - variable selection for causal questions

- ▶ Causal question in mind x impact on y . Having z variables to condition on, to get closer to causality
- ▶ Our aim is to focus on the coefficient of one variable. What matter here are the estimated value of the coefficient and its confidence interval
- ▶ Keep z – keep many variables that help comparing apples to apples
- ▶ Drop z if they not matter
- ▶ Functional form for z matters only for crucial confounders
- ▶ Present the model you judge is best, and then report a few other solutions – robustness.
- ▶ More in DA4

Multiple regression - variable selection – process

- ▶ Select control variables you want to include
- ▶ Select functional form one by one
- ▶ Focus on key variables by domain knowledge, add the rest linearly

- ▶ Key issue is sample size
 - ▶ For 20-40 obs, about 1-2 variables.
 - ▶ For 50-100 obs, about 2-4 variables
 - ▶ Few hundred obs, 5-10 variables could work
 - ▶ Few thousand obs, few dozen variables, including industry/country/profession etc dummies, interactions.
 - ▶ 10-100K obs - many variables, polynomials, interactions

Multiple regression - variable selection for prediction

- ▶ IF prediction - keep whatever works
- ▶ Balance is needed to ensure it works beyond the data at hand
- ▶ Overfitting: building a model that captures some patterns that may fit the data we have but would not generalize to the data we use the prediction for.
- ▶ Focus on functional form, interactions
- ▶ Value simplicity. Easier to explain, more robust.
- ▶ Formal way: BIC (“Bayesian Information Criterion”). Similar to R-squared but takes into account number of variables.
 - ▶ The smaller, the better
 - ▶ Do not use adjusted R-squared
- ▶ More in DA3

Summary take-away

- ▶ Multiple regression are linear models with several x variables.
- ▶ May include binary variables and interactions
- ▶ Multiple regression can take us closer to a causal interpretation and help make better predictions.