



PREDICTIVE ANALYSIS FOR DECISION MAKING

WEEK 1 AND 2
EXTENDING LINEAR REGRESSION

Dr Issam Malki
School of Finance and Accounting
Westminster Business School

September 2023

1

1



CLASSICAL LINEAR REGRESSION
RECAP

2

2



LINEAR REGRESSION

- The multiple linear regression model can be expressed as:

$$y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3} + \dots + \beta_k x_{ik} + u_i$$

i : subscript refers to the units and takes values from 1 to n .

x_{ij} : explanatory variables where $j=1, 2, \dots, k$ with $x_{i1}=1$.

β_j : unknown parameters to be estimated.

u_i : random error component. Unobserved with variance σ^2 .

We aim to estimate β_j and σ^2 .

3

3

LINEAR REGRESSION



- The multiple linear in a matrix notation:

$$Y = X\beta + \mathbf{u}$$

where Y is $(n \times 1)$, X is $(n \times k)$, \mathbf{u} is $(n \times 1)$ and β is $(k \times 1)$.

- The OLS estimator is define as follows:

$$\hat{\beta} = (X'X)^{-1}X'Y$$

Why? See notes.

4

4

LINEAR REGRESSION



- Assumptions

- X is fixed (non stochastic) with rank k (full column rank, meaning?)

- \mathbf{u} is random vector with $E(\mathbf{u}) = \mathbf{0}$ and $\text{var}(\mathbf{u}) = E(\mathbf{u}'\mathbf{u}) = \sigma^2\mathbf{I}$. Or:

$$\text{var}[u_i] = \text{var}[u] = \begin{bmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \ddots & \dots & \vdots \\ \vdots & & \ddots & \\ 0 & 0 & \dots & \sigma^2 \end{bmatrix}$$

5

5

LINEAR REGRESSION



- Assumptions – Implications

- Assumption 1:

- Explanatory variables are strictly exogenous (i.e. $E(X\mathbf{u})=0$).
- The X 's are linearly independent (i.e. no multicollinearity).

- Assumption 2:

- The fitted line is indeed in the middle.
- The errors are homoscedastic (no heteroskedasticity).
- The errors are serially uncorrelated (no autocorrelation)

- The validity of these assumptions is a must for the OLS to be a reliable estimator.

6

6



LINEAR REGRESSION

- Key results I: unbiasedness

$$E(\hat{\beta}) = \beta$$

See full proof in the notes

- Key results II:

$$\text{var}(\hat{\beta}) = \sigma^2 (X'X)^{-1}$$

where σ^2 is estimated using the sum of squared residuals, $e'e = \sum_{i=1}^n e_i^2$, as follows:

$$\hat{\sigma}^2 = s^2 = \frac{\sum_{i=1}^n e_i^2}{n - k}$$

- Key results III: Gauss-Markov Theorem:

OLS estimator is the Best Linear Unbiased Estimator (BLUE)

7



LINEAR REGRESSION IN PYTHON

- Data: nls80.xls. See problem set for details on data and variables definitions.
- You need the following libraries for the basic linear regression
 - Pandas
 - Numpy
 - Statsmodels.api
- For today, we need one more: matplotlib.pyplot
- Use the file: computer_seminar 1
- The model we are about to estimate is:

$$\ln(\text{wage}_i) = \beta_1 + \beta_2 \text{educ}_{i2} + \beta_3 \text{hours}_{i3} + u_i$$

8



LINEAR REGRESSION IN PYTHON

```

=====
               OLS Regression Results
=====
Dep. Variable:    logwage    R-squared:    0.103
Model:            OLS        Adj. R-squared: 0.101
Method:            Least Squares    F-statistic: 53.62
Date:            Wed, 03 Feb 2021    Prob (F-statistic): 9.12e-23
Time:            10:54:46    Log-Likelihood: -466.72
No. Observations: 935    AIC: 939.4
Df Residuals:      932    BIC: 954.0
Covariance Type:    nonrobust
=====
               coef    std err          t      P>|t|    [0.025    0.975]
-----
Intercept      6.1504      0.109     56.534      0.000      5.937      6.364
educ           0.0612      0.006     10.243      0.000      0.049      0.073
hours          -0.0044      0.002     -2.448      0.015     -0.008     -0.001
=====
Omnibus:            28.032    Durbin-Watson:      1.765
Prob(Omnibus):      0.000    Jarque-Bera (JB):    34.519
Skew:               -0.340    Prob(JB):            3.19e-08
Kurtosis:           3.651    Cond. No.             388.
=====

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
```

9



CLASSICAL LINEAR REGRESSION

EXTENSION I: RELAXING ASSUMPTIONS ABOUT THE ERROR TERM

10

10



ASSUMPTION VIOLATIONS: GENERAL VIEW

- We will now study these assumptions further, and in particular look at:
 - How we test for violations
 - Causes
 - Consequences
- In general we could encounter any combination of 3 problems:
 - the coefficient estimates are wrong
 - the associated standard errors are wrong
 - the distribution that we assumed for the test statistics will be inappropriate
- Solutions
 - The assumptions are no longer violated
 - we work around the problem so that we
 - use alternative techniques which are still valid

11



ASSUMPTION VIOLATIONS: $E(\varepsilon_i) = 0$

- Assumption that the mean of the disturbances is zero.
- For all diagnostic tests, we cannot observe the disturbances and so perform the tests of the residuals.
- The mean of the residuals will always be zero provided that there is a constant term in the regression.

12

12

ASSUMPTION VIOLATIONS: $\text{Var}(u_i)$ is not constant

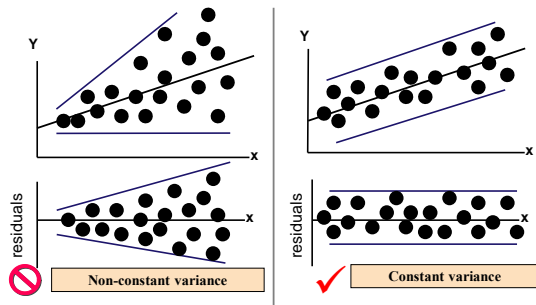


- Assumption that the variance of the error terms is finite and constant.
- Implications of this property
 - Homoscedasticity
 - Deviation from the fitted line is – on average- constant
 - All observed individuals have the same errors on average
- Violation of this assumptions leads to
 - Heteroscedasticity (spread of the variance)
 - Affect the standard errors not the estimated coefficients.
 - Ignoring this produces incorrectly smaller errors and higher t-statistics
 - OLS is no longer efficient and thus no longer BLUE.

13

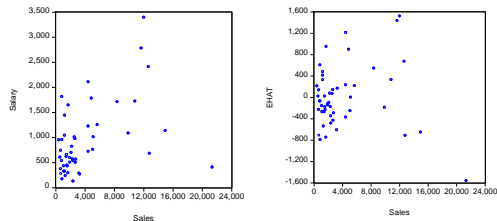
13

ASSUMPTION VIOLATIONS: $\text{Var}(u_i)$ is not constant



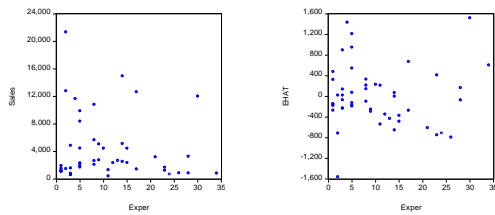
14

EXAMPLES: GRAPHICAL INSPECTION



15

EXAMPLES: GRAPHICAL INSPECTION



16

ASSUMPTION VIOLATIONS: $\text{Var}(\varepsilon_i) = \sigma^2 < \infty$ **White's Test for Heteroscedasticity**

- White's general test for heteroscedasticity is one of the best approaches because it makes few assumptions about the form of the heteroscedasticity.
- The hypotheses

$$H_0: \text{Residuals are homoscedastic}$$

$$H_1: \text{Residuals are heteroscedastic}$$

- The test is carried out as follows:

- Assume that the regression we carried out is as follows

$$y_i = \beta_1 + \beta_2 x_{2i} + \beta_3 x_{3i} + \varepsilon_i$$

And we want to test $\text{Var}(\varepsilon_i) = \sigma^2$. We estimate the model, obtaining the residuals, $\hat{\varepsilon}_i$

- Then run the auxiliary regression

$$\hat{\varepsilon}_i^2 = \alpha_0 + \alpha_1 x_{2i} + \alpha_2 x_{3i} + \alpha_3 x_{2i}^2 + \alpha_4 x_{3i}^2 + \alpha_5 x_{2i} x_{3i} + v_i$$

17

ASSUMPTION VIOLATIONS: $\text{Var}(\varepsilon_i) = \sigma^2 < \infty$ **White's Test for Heteroscedasticity**

- Obtain R^2 from the auxiliary regression and multiply it by the number of observations, T . It can be shown that

$$T R^2 \sim \chi^2(m)$$

where m is the number of regressors in the auxiliary regression excluding the constant term.

- If the χ^2 test statistic from step 3 is greater than the corresponding value from the statistical table then reject the null hypothesis that the disturbances are homoscedastic.

18



ASSUMPTION VIOLATIONS: $\text{Var}(\varepsilon_i) = \sigma^2 < \infty$

White's Test for Heteroscedasticity

- If we confirm the presence of heteroscedasticity:
 - Estimates are still unbiased, but not efficient
 - This implies that the standard errors and therefore the t statistics are incorrect
 - OLS estimator is not BLUE.
- Use White's Heteroscedasticity Consistent Standard Errors

19



ASSUMPTION VIOLATIONS: $\text{Var}(\varepsilon_i) = \sigma^2 < \infty$

Table A.4 χ^2 Distribution: critical values of χ^2 at 5%, 1%, and 0.1% significance levels

Degrees of freedom	5%	1%	0.1%
1	3.8415	6.6349	10.828
2	5.9915	9.2103	13.816
3	7.8797	11.3449	16.266
4	9.4877	13.2767	18.467
5	11.0705	15.0863	20.515
6	12.5916	16.8119	22.458
7	14.0671	18.4753	24.278
8	15.5073	20.0902	26.125
9	16.9190	21.6660	27.879
10	18.3070	23.2093	29.588
11	19.6751	24.7259	31.264
12	21.0261	26.2170	32.909
13	22.3645	27.6882	34.528
14	23.6848	29.1412	36.123
15	24.9958	30.578	37.567
16	26.2962	31.9995	38.932
17	27.5871	33.4097	40.289
18	28.8691	34.8053	41.635
19	30.1435	36.1909	42.979
20	31.4105	37.5662	44.314
21	32.6716	38.9322	45.642
22	33.9244	40.2893	46.963
23	35.1785	41.6384	48.278
24	36.4210	42.9798	49.589
25	37.6525	44.3141	50.892
26	38.8821	45.6422	52.188
27	40.1131	46.9635	53.476
28	41.3371	48.2782	54.755
29	42.5540	49.5891	56.026
30	43.7730	50.8922	57.289
40	55.7585	63.6957	72.402
50	67.5048	76.1539	86.561
60	79.0819	88.3774	99.607
70	90.5312	100.425	112.317
80	101.879	112.329	124.604
90	113.145	124.116	137.566
100	124.342	135.805	149.449

Reprinted from E. S. Pearson and H. O. Hooley (ed.), *Biometrika Tables for Statisticians*, Cambridge: Cambridge University Press, 1970 with kind permission of the Biometrika Trustees.

20



ASSUMPTION VIOLATIONS: $E[\varepsilon_i \varepsilon_j] = 0$

- This assumption implies that the errors are independent.
- If violated, then the errors are autocorrelated.
 - Errors are not independent
 - Errors in one period (for one individual observation) is correlated with another.
- Autocorrelation is a prominent feature of time series data. Could be found in cross section too.

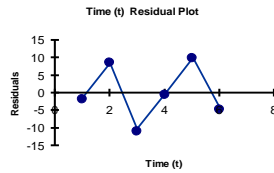
21



ASSUMPTION VIOLATIONS: $E[\varepsilon_i \varepsilon_j] = 0$

- Autocorrelation is correlation of the errors (residuals) over time

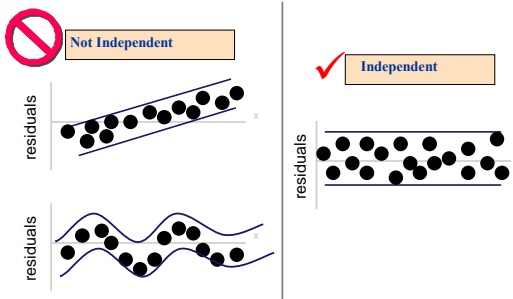
- Here, residuals show a cyclic pattern, not random



- Violates the regression assumption that residuals are random and independent

22

ASSUMPTION VIOLATIONS: $E[\varepsilon_i \varepsilon_j] = 0$



23

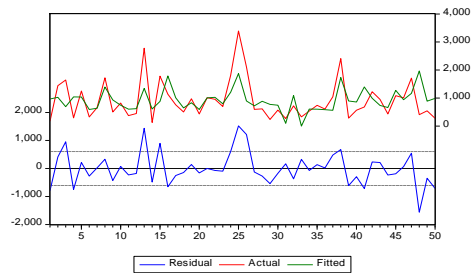
ASSUMPTION VIOLATIONS: $E[\varepsilon_i \varepsilon_j] = 0$



- The consequences of autocorrelation
 - The coefficient estimates derived using OLS are still unbiased, but they are inefficient, i.e. they are not BLUE, even in large sample sizes.
 - Thus, if the standard error estimates are inappropriate, there exists the possibility that we could make the wrong inferences.
 - R^2 is likely to be inflated relative to its “correct” value for positively correlated residuals.

24

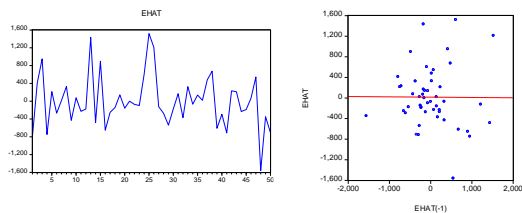
APPLICATION 6 DETECTION OF AUTOCORRELATION



25

25

APPLICATION 6 DETECTION OF AUTOCORRELATION



26

26

ASSUMPTION VIOLATIONS: $E[\varepsilon_i \varepsilon_j] = 0$



The Durbin-Watson Test

- The Durbin-Watson (DW) is a test for first order autocorrelation - i.e. it assumes that the relationship is between an error and the previous one

$$\varepsilon_t = \rho \varepsilon_{t-1} + v_t \quad (1)$$

where $v_t \sim N(0, \sigma_v^2)$.

- The DW test statistic actually tests

$H_0: \rho = 0$ (no autocorrelation) and $H_1: \rho \neq 0$

- The test statistic is calculated by

$$DW = \frac{\sum_{t=2}^T (\varepsilon_t - \varepsilon_{t-1})^2}{\sum_{t=1}^T \varepsilon_t^2}$$

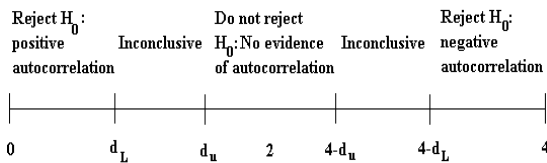
27

27



ASSUMPTION VIOLATIONS: $E[\varepsilon_i \varepsilon_j] = 0$

The Durbin-Watson Test



Conditions which Must be Fulfilled for DW to be a Valid Test

1. Constant term in regression
2. Regressors are non-stochastic
3. No lags of dependent variable

28

28



ASSUMPTION VIOLATIONS: $E[\varepsilon_i \varepsilon_j] = 0$

Table A-9 Durbin-Watson Statistics of autocorrelation ρ_1 and ρ_2 at 5% significance level

n	k = 1				k = 2				k = 3				k = 4				k = 5			
	d_L	d_U	d_L	d_U	d_L	d_U	d_L	d_U	d_L	d_U	d_L	d_U	d_L	d_U	d_L	d_U	d_L	d_U		
15	1.08	1.36	0.93	1.24	0.82	1.13	0.69	1.07	0.56	0.91	0.45	0.77	0.36	0.66	0.27	0.55	0.20	0.48		
20	1.10	1.37	0.95	1.26	0.84	1.15	0.71	1.09	0.58	0.93	0.47	0.80	0.38	0.68	0.29	0.57	0.22	0.50		
25	1.12	1.39	0.97	1.28	0.86	1.17	0.73	1.11	0.60	0.95	0.49	0.82	0.40	0.70	0.31	0.59	0.24	0.52		
30	1.14	1.41	0.99	1.30	0.88	1.19	0.75	1.13	0.62	0.97	0.51	0.84	0.42	0.72	0.33	0.61	0.26	0.54		
35	1.16	1.43	1.01	1.32	0.90	1.21	0.77	1.15	0.64	0.99	0.53	0.86	0.44	0.74	0.35	0.63	0.28	0.56		
40	1.18	1.45	1.03	1.34	0.92	1.23	0.79	1.17	0.66	1.01	0.55	0.88	0.46	0.76	0.37	0.65	0.30	0.58		
45	1.20	1.47	1.05	1.36	0.94	1.25	0.81	1.19	0.68	1.03	0.57	0.90	0.48	0.78	0.39	0.67	0.32	0.60		
50	1.22	1.49	1.07	1.38	0.96	1.27	0.83	1.21	0.70	1.05	0.59	0.92	0.50	0.80	0.41	0.69	0.34	0.62		
55	1.24	1.51	1.09	1.40	0.98	1.29	0.85	1.23	0.72	1.07	0.61	0.94	0.52	0.82	0.43	0.71	0.36	0.64		
60	1.26	1.53	1.11	1.42	1.00	1.31	0.87	1.25	0.74	1.09	0.63	0.96	0.54	0.84	0.45	0.73	0.38	0.66		
65	1.28	1.55	1.13	1.44	1.02	1.33	0.89	1.27	0.76	1.11	0.65	0.98	0.56	0.86	0.47	0.75	0.40	0.68		
70	1.30	1.57	1.15	1.46	1.04	1.35	0.91	1.29	0.78	1.13	0.67	1.00	0.58	0.88	0.49	0.77	0.42	0.70		
75	1.32	1.59	1.17	1.48	1.06	1.37	0.93	1.31	0.80	1.15	0.69	1.02	0.60	0.90	0.51	0.79	0.44	0.72		
80	1.34	1.61	1.19	1.50	1.08	1.39	0.95	1.33	0.82	1.17	0.71	1.04	0.62	0.92	0.53	0.81	0.46	0.74		
85	1.36	1.63	1.21	1.52	1.10	1.41	0.97	1.35	0.84	1.19	0.73	1.06	0.64	0.94	0.55	0.83	0.48	0.76		
90	1.38	1.65	1.23	1.54	1.12	1.43	0.99	1.37	0.86	1.21	0.75	1.08	0.66	0.96	0.57	0.85	0.50	0.78		
95	1.40	1.67	1.25	1.56	1.14	1.45	1.01	1.39	0.88	1.23	0.77	1.10	0.68	0.98	0.59	0.87	0.52	0.80		
100	1.42	1.69	1.27	1.58	1.16	1.47	1.03	1.41	0.90	1.25	0.79	1.12	0.70	1.00	0.61	0.89	0.54	0.82		

29

29



ASSUMPTION VIOLATIONS: $E[\varepsilon_i \varepsilon_j] = 0$

The Breusch-Godfrey Test (LM Test)

- It is a more general test for r^{th} order autocorrelation:

$$\varepsilon_t = \rho_1 \varepsilon_{t-1} + \rho_2 \varepsilon_{t-2} + \dots + \rho_r \varepsilon_{t-r} + v_t, v_t \sim iid(0, \sigma^2)$$

- The null and alternative hypotheses are:

$$H_0: \rho_1 = 0 \text{ and } \rho_2 = 0 \text{ and } \dots \text{ and } \rho_r = 0 \text{ (No autocorrelation)}$$

$$H_1: \rho_1 \neq 0 \text{ or } \rho_2 \neq 0 \text{ or } \dots \text{ or } \rho_r \neq 0$$

- The test is carried out as follows:

1. Estimate the linear regression using OLS and obtain the residuals, $\hat{\varepsilon}_i$.

2. Regress $\hat{\varepsilon}_i$ on all of the regressors from stage 1 (the x 's) plus

Obtain R^2 from this regression.

3. It can be shown that $(T-r)R^2 \sim \chi^2(r)$

- If the test statistic exceeds the critical value from the statistical tables, reject the null hypothesis of no autocorrelation.

30

ASSUMPTION VIOLATIONS: $E[\varepsilon_i, \varepsilon_j]=0$



- There are many remedies to correct for serial correlation
 - Using lagged dependent variable (will be revisited when we deal with time series)
 - Use Methods such as Cochrane-Orcutt if the form is known.
 - Use Autocorrelation Consistent Standard Errors such as Newey-West
- Time Series Data
 - Autocorrelation is a key feature
 - Easily fixed if the dynamic process is stationary – by adding lags.
 - Has crucial implications on the stability of the process and long-run analysis
 - We will revisit this issue when dealing with time series

31

31



THANK YOU

32

32