

# Time Series Analysis - Introduction

ECON 722

---

Elena Pesavento — Emory University

Fall 2024

# Introduction

- History - popular in early 90 s, making comeback now.
- What is a Time Series?  $Y_t$ 
  - The main difference between time series econometrics and cross-section is in dependence structure. Cross section econometrics mainly deals with i.i.d. observations,  $Y_i$ , while in time series each new arriving observation is stochastically depending on the previously observed,  $Y_t$ .
- The dependence is our best friend and a great enemy.
  - On one side, the dependence screw up your inferences (CLT you learned was for i.i.d. data). On the other side, the dependence allow us to do more by exploiting it. For example, we can make forecasts (which are almost non-sense in cross-section).

# Introduction

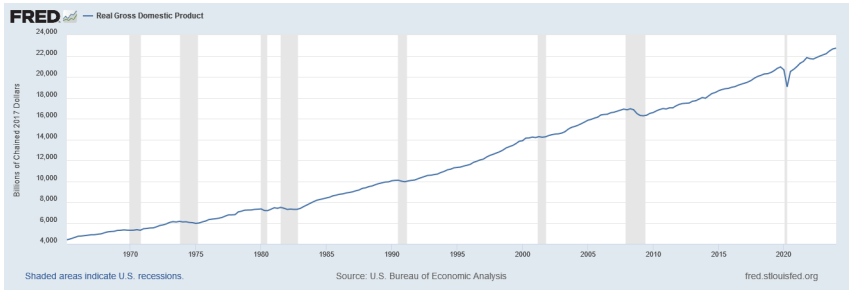
- We are going to ask questions like:
  - What are the characteristics of time series data?
  - Which type of questions can we ask/answer?
  - Which type of questions we cannot answer? Or can we?
- We can roughly divide time series into macro and finance related stuff.
  - Macro Time series mostly focuses on means. Often limited by small number of observations available over long horizon (e.g. 20 years monthly is  $T=300$ ).
  - Financial data usually high-frequency over short period of time. This allows us to model volatility and higher moments.

# Introduction

- Which type of questions can we ask/answer?
  - In the first part of the class we will be using Regression Models for *Forecasting*.
  - Forecasting and estimation of causal effects are quite different objectives (more on this later).
  - In the second half of the class we will discuss if and when we can talk about “causality”.
- Forecasting:
  - It is very important in macro and finance.
  - Omitted variable bias isn't a problem.
  - External validity is paramount: the model estimated using historical data must hold into the (near) future
  - Interpretation of the coefficients is not the goal.

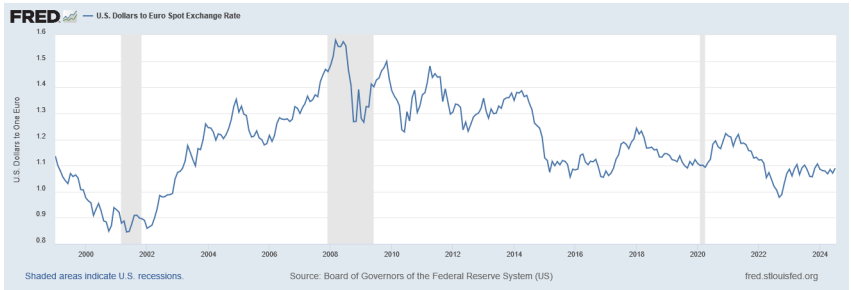
# Some Examples of Time Series Data

Real GDP, 1965:1-2024:1, Quarterly (GDP is only quarterly), Seasonally adjusted, billions of chained 2017 dollars. From FRED.



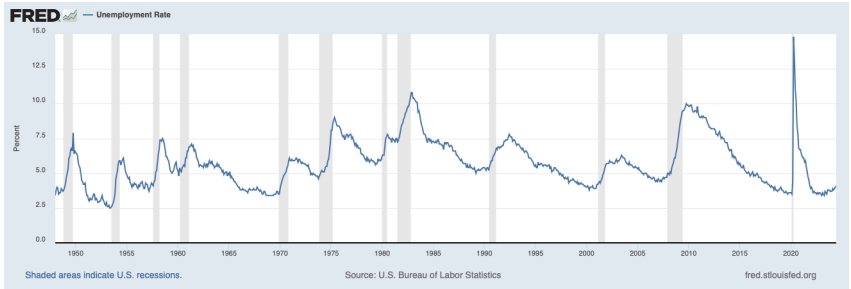
# Some Examples of Time Series Data

US Dollars to Euro Spot Exchange rate. Not Seasonally Adjusted.  
1999:01:04–2024:07:12. From FRED.



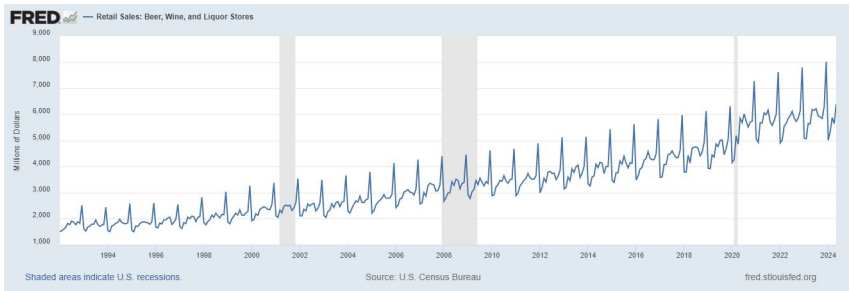
# Some Examples of Time Series Data

Unemployment Rate (in %). Seasonally Adjusted. 1948:01-01 - 2024:06-01. From FRED.



# Some Examples of Time Series Data

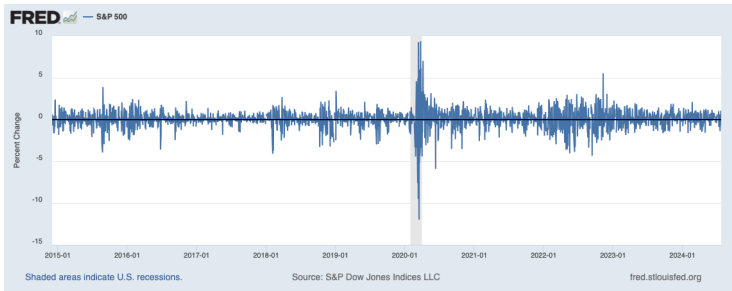
- U.S. Liquor Retail Sales: Beer, Wine, and Liquor Stores
- Millions of Dollars, Not Seasonally Adjusted
- 1992-01-01 to 2024-05-01, Monthly





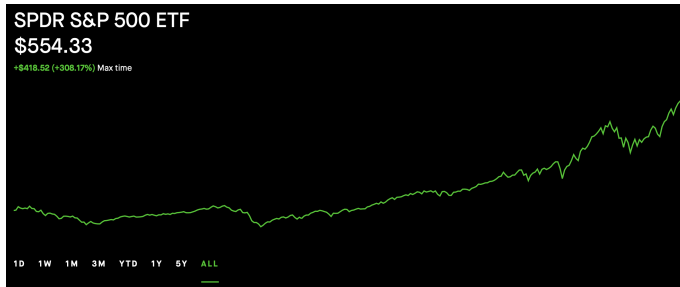
# Some Examples of Time Series Data

- SP-500 returns
- Percentage Change, Not Seasonally Adjusted
- 2014-12-01 to 2024-07-23, Daily



# Some Examples of Time Series Data

SPDR S&P 500 ETF (Ticker: SPY). Not Seasonally Adjusted.  
2000:01-30 - 2024:07-21. From Robinhood



# Components of any Time Series

There are **three** main components in any Time Series:

## **Trend**

Part of a series' movement that corresponds to long-term, slow evolution

## **Seasonality**

Part of a series' movement that repeats each year

## **Cycles**

A catch-all phrase for various forms of dynamic behavior that link the present to the past and the future to the present

# Trends and Breaks

---

- Trend is a slow, long run evolution in the variable that we want to model

- Trend is a slow, long run evolution in the variable that we want to model
- There are two kind of trends: *Deterministic* and *Stochastic*

- Trend is a slow, long run evolution in the variable that we want to model
- There are two kind of trends: *Deterministic* and *Stochastic*
- It is also called **deterministic trend** because evolves in predictable way. How do we model a determinisitc trend?

$$y_t = c + \beta t + u_t$$

- Trend is a slow, long run evolution in the variable that we want to model
- There are two kind of trends: *Deterministic* and *Stochastic*
- It is also called **deterministic trend** because evolves in predictable way. How do we model a determinisitc trend?

$$y_t = c + \beta t + u_t$$

- **Stochastic trends** are non predictable (e.g unit roots, random walk)



# Deterministic Trend

There are various kinds of deterministic trends:

- Linear Trend

$$y_t = c + \beta t$$

# Deterministic Trend

There are various kinds of deterministic trends:

- Linear Trend

$$y_t = c + \beta t$$

- Quadratic Trend

$$y_t = c + \beta_1 t + \beta_2 t^2$$

# Deterministic Trend

There are various kinds of deterministic trends:

- Linear Trend

$$y_t = c + \beta t$$

- Quadratic Trend

$$y_t = c + \beta_1 t + \beta_2 t^2$$

- Exponential Trend

$$y_t = ce^{\beta_1 t} \text{ or } \log(y_t) = \log(c) + \beta_1 t$$

# Deterministic Trend

There are various kinds of deterministic trends:

- Linear Trend

$$y_t = c + \beta t$$

- Quadratic Trend

$$y_t = c + \beta_1 t + \beta_2 t^2$$

- Exponential Trend

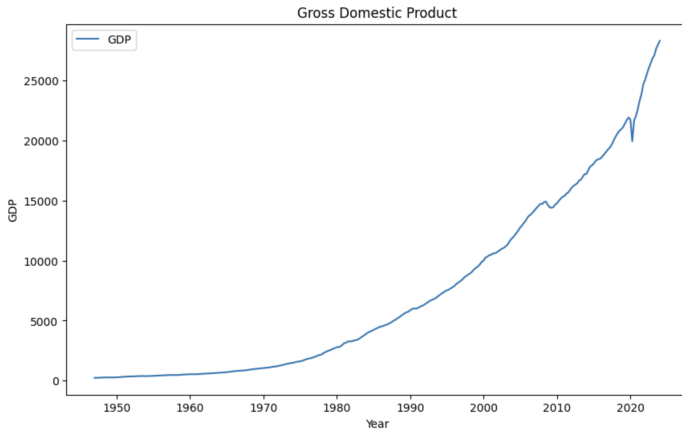
$$y_t = ce^{\beta_1 t} \text{ or } \log(y_t) = \log(c) + \beta_1 t$$

**There are other kinds of trends but those three approximate any trend fairly well!**

Adding an explanatory variable that looks like a trend will explain the same property.

# Looking Closer at GDP

Follow the Module 1 Python code file to graph this using the data 'GDP.csv'



# How do we estimate a Deterministic Trend?

The model we have in mind is

$$GDP_t = c + \beta t + u_t$$

We can run this using any statistical software using some OLS command and we get the following

OLS Regression Results						
=====						
Dep. Variable:	GDP	R-squared:	0.865			
Model:	OLS	Adj. R-squared:	0.865			
Method:	Least Squares	F-statistic:	1973.			
Date:	Tue, 23 Jul 2024	Prob (F-statistic):	1.03e-135			
Time:	19:45:10	Log-Likelihood:	-2884.3			
No. Observations:	309	AIC:	5773.			
Df Residuals:	307	BIC:	5780.			
Df Model:	1					
Covariance Type:	nonrobust					
=====						
	coef	std err	t	P> t	[0.025	0.975]
-----						
Constant	-4763.0379	311.942	-15.269	0.000	-5376.853	-4149.223
Time	77.8620	1.753	44.422	0.000	74.413	81.311
=====						
Omnibus:	34.120	Durbin-Watson:	0.005			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	42.573			
Skew:	0.902	Prob(JB):	5.69e-10			
Kurtosis:	3.236	Cond. No.	355.			
=====						

# How do we estimate a Deterministic Trend?

```
=====
                        OLS Regression Results
=====
Dep. Variable:          GDP      R-squared:                0.865
Model:                  OLS      Adj. R-squared:             0.865
Method:                 Least Squares      F-statistic:          1973.
Date:                   Tue, 23 Jul 2024    Prob (F-statistic):      1.03e-135
Time:                   19:45:10           Log-Likelihood:         -2884.3
No. Observations:       309            AIC:                   5773.
Df Residuals:           307            BIC:                   5780.
Df Model:                1
Covariance Type:        nonrobust
=====
                        coef      std err          t      P>|t|      [0.025      0.975]
=====
Constant    -4763.0379      311.942     -15.269      0.000     -5376.853    -4149.223
Time         77.8620         1.753       44.422      0.000       74.413      81.311
=====
Omnibus:                 34.120    Durbin-Watson:           0.005
Prob(Omnibus):            0.000    Jarque-Bera (JB):        42.573
Skew:                     0.902    Prob(JB):                5.69e-10
Kurtosis:                 3.236    Cond. No.                 355.
=====
```

- Trend explains most of GDP variation. Typical of most time series'.

# How do we estimate a Deterministic Trend?

```
=====
                        OLS Regression Results
=====
Dep. Variable:          GDP      R-squared:                0.865
Model:                  OLS      Adj. R-squared:             0.865
Method:                 Least Squares      F-statistic:          1973.
Date:                   Tue, 23 Jul 2024    Prob (F-statistic):      1.03e-135
Time:                   19:45:10           Log-Likelihood:         -2884.3
No. Observations:       309      AIC:                     5773.
Df Residuals:           307      BIC:                     5780.
Df Model:                1
Covariance Type:        nonrobust
=====
                        coef      std err          t      P>|t|      [0.025      0.975]
-----
Constant    -4763.0379      311.942     -15.269      0.000     -5376.853    -4149.223
Time         77.8620         1.753       44.422      0.000       74.413      81.311
=====
Omnibus:                 34.120    Durbin-Watson:           0.005
Prob(Omnibus):            0.000    Jarque-Bera (JB):        42.573
Skew:                     0.902    Prob(JB):                 5.69e-10
Kurtosis:                 3.236    Cond. No.                 355.
=====
```

- Trend explains most of GDP variation. Typical of most time series'.
- $R^2$  is really high. Is this a good regression? Why or why not?

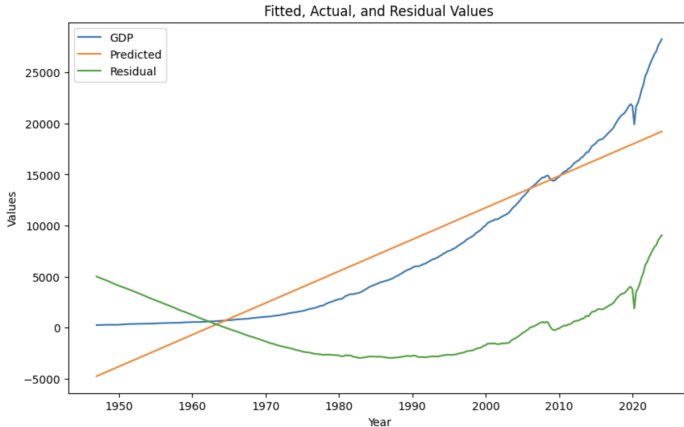


# How do we estimate a Deterministic Trend?

```
=====
                        OLS Regression Results
=====
Dep. Variable:          GDP    R-squared:                0.865
Model:                  OLS    Adj. R-squared:             0.865
Method:                 Least Squares    F-statistic:          1973.
Date:                   Tue, 23 Jul 2024    Prob (F-statistic):    1.03e-135
Time:                   19:45:10    Log-Likelihood:       -2884.3
No. Observations:      309    AIC:                   5773.
Df Residuals:          307    BIC:                   5780.
Df Model:               1
Covariance Type:       nonrobust
=====
                        coef    std err          t      P>|t|      [0.025    0.975]
=====
Constant    -4763.0379    311.942    -15.269    0.000    -5376.853    -4149.223
Time         77.8620     1.753     44.422    0.000     74.413     81.311
=====
Omnibus:                 34.120    Durbin-Watson:           0.005
Prob(Omnibus):            0.000    Jarque-Bera (JB):        42.573
Skew:                     0.902    Prob(JB):                5.69e-10
Kurtosis:                 3.236    Cond. No.                 355.
=====
```

- Trend explains most of GDP variation. Typical of most time series'.
- $R^2$  is really high. Is this a good regression? Why or why not?
- But how do we know if this is truly deterministic trend? We will spent A LOT of time talking about this and the difference between deterministic and stochastic trends

# How do we estimate a Deterministic Trend?



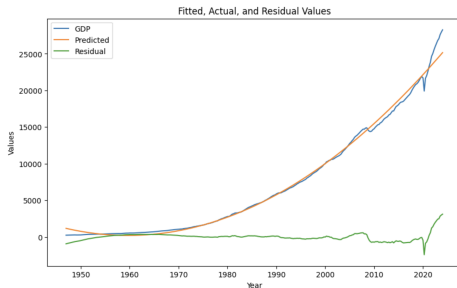
- Look at fitted values! (What do we learn?)
- Which pattern do we want the residuals to have?

# How do we estimate a Deterministic Trend?

Let us now try a quadratic deterministic trend model!

$$GDP_t = c + \beta_1 t + \beta_2 t^2 + u_t$$

OLS Regression Results						
=====						
Dep. Variable:	GDP	R-squared:		0.994		
Model:	OLS	Adj. R-squared:		0.994		
Method:	Least Squares	F-statistic:		2.683e+04		
Date:	Tue, 23 Jul 2024	Prob (F-statistic):		0.00		
Time:	20:19:05	Log-Likelihood:		-2395.0		
No. Observations:	309	AIC:		4796.		
Df Residuals:	306	BIC:		4807.		
Df Model:	2					
Covariance Type:	nonrobust					
=====						
	coef	std err	t	P> t	[0.025	0.975]
-----						
Constant	1174.2928	95.795	12.258	0.000	985.793	1362.792
Time	-38.1771	1.437	-26.570	0.000	-41.004	-35.350
Time_sq	0.3768	0.005	83.425	0.000	0.368	0.386
=====						
Omnibus:	164.549	Durbin-Watson:		0.090		
Prob(Omnibus):	0.000	Jarque-Bera (JB):		1611.955		
Skew:	1.969	Prob(JB):		0.00		
Kurtosis:	13.474	Cond. No.		1.27e+05		
=====						

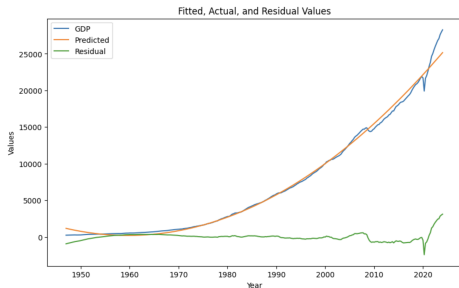


# How do we estimate a Deterministic Trend?

Let us now try a quadratic deterministic trend model!

$$GDP_t = c + \beta_1 t + \beta_2 t^2 + u_t$$

OLS Regression Results						
=====						
Dep. Variable:	GDP	R-squared:		0.994		
Model:	OLS	Adj. R-squared:		0.994		
Method:	Least Squares	F-statistic:		2.683e+04		
Date:	Tue, 23 Jul 2024	Prob (F-statistic):		0.00		
Time:	20:19:05	Log-Likelihood:		-2395.0		
No. Observations:	309	AIC:		4796.		
Df Residuals:	306	BIC:		4807.		
Df Model:	2					
Covariance Type:	nonrobust					
=====						
	coef	std err	t	P> t	[0.025	0.975]
-----						
Constant	1174.2928	95.795	12.258	0.000	985.793	1362.792
Time	-38.1771	1.437	-26.570	0.000	-41.004	-35.350
Time_sq	0.3768	0.005	83.425	0.000	0.368	0.386
=====						
Omnibus:	164.549	Durbin-Watson:		0.090		
Prob(Omnibus):	0.000	Jarque-Bera (JB):		1611.955		
Skew:	1.969	Prob(JB):		0.00		
Kurtosis:	13.474	Cond. No.		1.27e+05		
=====						



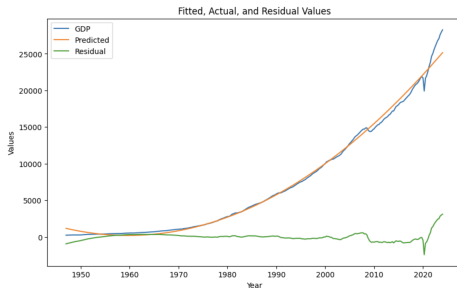
- What can you say from the output? What about from the graph of the residuals?

# How do we estimate a Deterministic Trend?

Let us now try a quadratic deterministic trend model!

$$GDP_t = c + \beta_1 t + \beta_2 t^2 + u_t$$

OLS Regression Results						
=====						
Dep. Variable:	GDP	R-squared:		0.994		
Model:	OLS	Adj. R-squared:		0.994		
Method:	Least Squares	F-statistic:		2.683e+04		
Date:	Tue, 23 Jul 2024	Prob (F-statistic):		0.00		
Time:	20:19:05	Log-Likelihood:		-2395.0		
No. Observations:	309	AIC:		4796.		
Df Residuals:	306	BIC:		4807.		
Df Model:	2					
Covariance Type:	nonrobust					
=====						
	coef	std err	t	P> t	[0.025	0.975]
-----						
Constant	1174.2928	95.795	12.258	0.000	985.793	1362.792
Time	-38.1771	1.437	-26.570	0.000	-41.004	-35.350
Time_sq	0.3768	0.005	83.425	0.000	0.368	0.386
-----						
Omnibus:	164.549	Durbin-Watson:		0.090		
Prob(Omnibus):	0.000	Jarque-Bera (JB):		1611.955		
Skew:	1.969	Prob(JB):		0.000		
Kurtosis:	13.474	Cond. No.		1.27e+05		
=====						



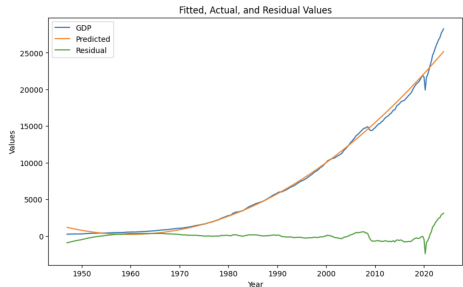
- What can you say from the output? What about from the graph of the residuals?
- In general we would like the residuals to be inside the interval in the graph.

# How do we estimate a Deterministic Trend?

Let us now try a quadratic deterministic trend model!

$$GDP_t = c + \beta_1 t + \beta_2 t^2 + u_t$$

OLS Regression Results						
Dep. Variable:	GDP	R-squared:	0.994			
Model:	OLS	Adj. R-squared:	0.994			
Method:	Least Squares	F-statistic:	2.683e+04			
Date:	Tue, 23 Jul 2024	Prob (F-statistic):	0.00			
Time:	20:19:05	Log-Likelihood:	-2395.0			
No. Observations:	309	AIC:	4796.			
Df Residuals:	306	BIC:	4807.			
Df Model:	2					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
Constant	1174.2928	95.795	12.258	0.000	985.793	1362.792
Time	-38.1771	1.437	-26.570	0.000	-41.004	-35.350
Time_sq	0.3768	0.005	83.425	0.000	0.368	0.386
Omnibus:	164.549	Durbin-Watson:		0.090		
Prob(Omnibus):	0.000	Jarque-Bera (JB):		1611.955		
Skew:	1.969	Prob(JB):		0.00		
Kurtosis:	13.474	Cond. No.		1.27e+05		



Deterministic Trends do not have an economics interpretation, yet it is incredibly powerful for forecasting.

# Information Criteria

---

How do we select between different models?

- The model with the highest  $R^2$  is not always the best model for out of sample forecast.



How do we select between different models?

- The model with the highest  $R^2$  is not always the best model for out of sample forecast.
- Most ICs attempt to find the model with the smallest out-of-sample 1-step-ahead mean squared error.

$$MSE = \frac{\sum_{t=1}^T e_t^2}{T} \quad \text{where } e_t = y_t - \hat{y}_t = \hat{u}_t$$

Note that

$$R^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{\sum e_t^2}{\sum (y_t - \bar{y})^2}$$

Where TTS depends only on the data.

Smallest MSE  $\Leftrightarrow$  highest  $R^2$

We could look at the model with smallest MSE. We said that  $R^2$  increases even if we include irrelevant variables so that we needed to correct for the degrees of freedom (d.f.). The same is true for MSE.

Different Information Criteria correct MSE for the degrees of freedom with different weights.

Different Information Criteria correct MSE for the degrees of freedom with different weights.

- Akaike Information Criteria (AIC)

$$AIC = e^{\frac{2k}{T}} MSE$$

- Schwartz Information Criteria (SIC)

$$SIC = T^{\frac{k}{T}} MSE$$

Different Information Criteria correct MSE for the degrees of freedom with different weights.

- Akaike Information Criteria (AIC)

$$AIC = e^{\frac{2k}{T}} MSE$$

- Schwartz Information Criteria (SIC)

$$SIC = T^{\frac{k}{T}} MSE$$

*Between two models, pick the one with the smallest AIC and SIC.*

# Breaks

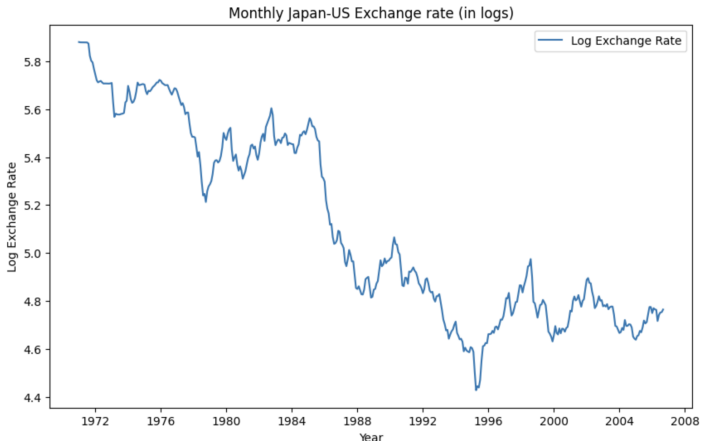
---

Often we may want to model a break in the mean or a break in trend.

- We can use dummy variables to model that.
- Break in the intercept or break in the slope of the trend? (think of an example)
- If we know the date of the break (new law, big event..) then we can impose that break. **We should test for significance of that break.**
- If we do not know the date of the break we can test for a break at unknown time.

# Broken Trend

Some time we may be interested in estimating a “broken trend”. Often there are reason to expect that there is a structural break at a specific date.



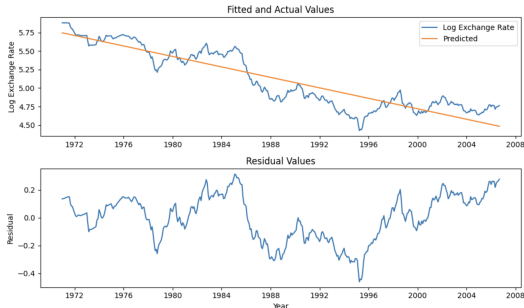


# Broken Trend

As before, we can of course estimate a linear trend

$$yen_t = c + \beta t + u_t$$

OLS Regression Results						
Dep. Variable:	Log_ER	R-squared:	0.831			
Model:	OLS	Adj. R-squared:	0.830			
Method:	Least Squares	F-statistic:	2096.			
Date:	Tue, 23 Jul 2024	Prob (F-statistic):	7.93e-167			
Time:	20:54:44	Log-Likelihood:	165.73			
No. Observations:	429	AIC:	-327.5			
Df Residuals:	427	BIC:	-319.3			
Df Model:	1					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
Constant	5.7447	0.016	361.595	0.000	5.713	5.776
Time	-0.0029	6.43e-05	-45.786	0.000	-0.003	-0.003
Omnibus:	27.346	Durbin-Watson:	0.027			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	17.741			
Skew:	-0.368	Prob(JB):	0.000140			
Kurtosis:	2.328	Cond. No.	494.			



It looks like there is a break in the trend around 1986. Before mid 1986 the trend is flatter and it is decreasing after that!

To model this we use Dummy variables!

To model this we use Dummy variables!

A **dummy variable** is a variables that only take the values of 0 and 1.

Define  $D_t$  as follows

$$D_t = \begin{cases} 0 & \text{if } t < \text{December 1985} \\ 1 & \text{if } t > \text{December 1985} \end{cases}$$

Let's see what model we are estimating if we include this variable in this regression.

$$yen_t = c + \beta_1 D_t + \beta_2 t + u_t$$

# Broken Trend

Clearly, before 1986,  $D_t = 0$

$$yen_t = c + \beta_2 t + u_t$$

However, after 1986,  $D_t = 1$

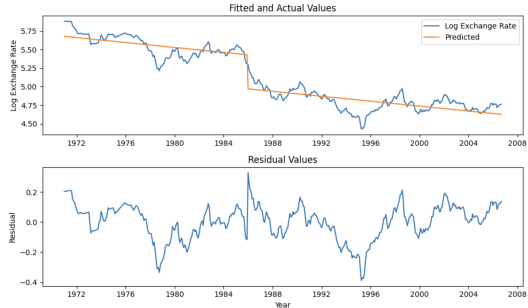
$$yen_t = (c + \beta_1) + \beta_2 t + u_t$$

Ergo, the coefficient on the dummy determines how the intercept has changed before and after 1986!

# Broken Trend

## OLS Regression Results

Dep. Variable:	Log_ER	R-squared:	0.917			
Model:	OLS	Adj. R-squared:	0.917			
Method:	Least Squares	F-statistic:	2354.			
Date:	Tue, 23 Jul 2024	Prob (F-statistic):	5.67e-231			
Time:	20:59:32	Log-Likelihood:	318.55			
No. Observations:	429	AIC:	-631.1			
Df Residuals:	426	BIC:	-618.9			
Df Model:	2					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
Constant	5.6766	0.012	489.406	0.000	5.654	5.699
Post_85	-0.4583	0.022	-21.038	0.000	-0.501	-0.415
Time	-0.0014	8.68e-05	-15.912	0.000	-0.002	-0.001
Omnibus:	30.970	Durbin-Watson:	0.091			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	35.794			
Skew:	-0.661	Prob(JB):	1.69e-08			
Kurtosis:	3.507	Cond. No.	979.			



- How many data points are used to estimate the mean in each time period?
- What does this tell us about how many data points we need in each period?

# Broken Trend

We can use the same idea to allow for the possibility of a change in the slope. Use an "interaction" term ( $D_t \cdot t$ )

$$yen_t = c + \beta_1 D_t + \beta_2 t + \beta_3 (D_t \cdot t) + u_t$$

Before 1986,  $D_t = 0$

$$yen_t = c + \beta_2 t + u_t$$

After 1986,  $D_t = 1$

$$yen_t = (c + \beta_1) + (\beta_2 + \beta_3) t + u_t$$

The coefficient on the dummy determines how the intercept has changed before and after 1986, the coefficient on the interactive term determines how the slope has changed before and after 1986.

# Broken Trend

## OLS Regression Results

```

=====
Dep. Variable:      Log_ER      R-squared:      0.922
Model:              OLS         Adj. R-squared:    0.922
Method:             Least Squares      F-statistic:    1679.
Date:               Tue, 23 Jul 2024    Prob (F-statistic): 3.64e-235
Time:               21:03:39           Log-Likelihood: 332.34
No. Observations:   429              AIC:           -656.7
Df Residuals:       425              BIC:           -640.4
Df Model:           3
Covariance Type:    nonrobust
=====

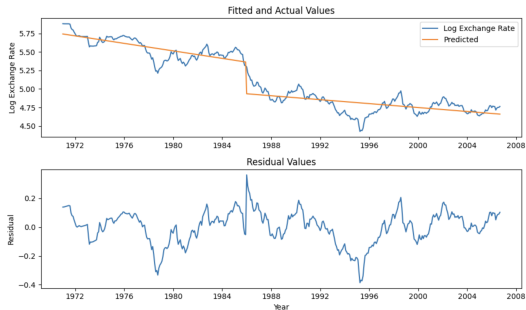
```

	coef	std err	t	P> t	[0.025	0.975]
Constant	5.7417	0.017	345.225	0.000	5.709	5.774
Post_85	-0.6069	0.035	-17.313	0.000	-0.676	-0.538
Time	-0.0021	0.000	-13.120	0.000	-0.002	-0.002
Time_Post_85	0.0010	0.000	5.312	0.000	0.001	0.001

```

=====
Omnibus:            35.302      Durbin-Watson:      0.092
Prob(Omnibus):      0.000      Jarque-Bera (JB):    43.676
Skew:               -0.671     Prob(JB):            3.28e-10
Kurtosis:           3.800      Cond. No.            2.27e+03
=====

```



Why is important to detect breaks?

- If a break occurs in the population regression model in the sample, then the OLS will estimate a relationship that holds 'on average'. The average combines the two periods!



Why is important to detect breaks?

- If a break occurs in the population regression model in the sample, then the OLS will estimate a relationship that holds 'on average'. The average combines the two periods!
  - Depending on the location and the size of the break the estimated 'average' relationship can be very different from the true population line.

Why is important to detect breaks?

- If a break occurs in the population regression model in the sample, then the OLS will estimate a relationship that holds 'on average'. The average combines the two periods!
  - Depending on the location and the size of the break the estimated 'average' relationship can be very different from the true population line.
  - For this reason it is important that we test the presence of breaks in time series.



- Sometime we may suspect that there is a break at a certain date (e.g. Bretton Woods (1973)).

- Sometime we may suspect that there is a break at a certain date (e.g. Bretton Woods (1973)).
- We showed how to use dummy variables to estimate possible breaks in the intercept and in the slope.

- Sometime we may suspect that there is a break at a certain date (e.g. Bretton Woods (1973)).
- We showed how to use dummy variables to estimate possible breaks in the intercept and in the slope.
- If there is no break, then the regression line should be there same in both periods.

$$y_t = c + \beta_1 D_t + \beta_2 t + \beta_3 (D_t \cdot t) + u_t$$

- Sometime we may suspect that there is a break at a certain date (e.g. Bretton Woods (1973)).
- We showed how to use dummy variables to estimate possible breaks in the intercept and in the slope.
- If there is no break, then the regression line should be there same in both periods.

$$y_t = c + \beta_1 D_t + \beta_2 t + \beta_3 (D_t \cdot t) + u_t$$

- We want to test the null hypothesis of no break. What is the null? Which test can we use? **Answer: Let's try the Chow's Test!**

# Operationalizing the Chow's Test

Consider the model

$$yen_t = c + \beta_1 D_t + \beta_2 t + \beta_3 D_t \cdot t + u_t$$

We want to test the hypothesis

$$H_0 : \beta_2 = \beta_3 = 0$$

$$H_a : \text{otherwise}$$

The null hypothesis of joint insignificance of  $D$  can be run as an F-test with  $k$  and  $N_1 + N_2 - 2k$  degrees of freedom where  $N_1$  is  $\text{card}(t_1, \dots, t^*)$  and  $N_2$  is  $\text{card}(t^* + 1, \dots, T)$  and  $k$  parameters.

$$F = \frac{(RSS_C - (RSS_1 + RSS_2))/k}{(RSS_1 + RSS_2)/(N_1 + N_2 - 2k)}$$

Note: The same result can be achieved with dummy variables!



# Operationalizing the Chow's Test

```
# Run a Chow test to test for structural break

# Define the two subsamples
er_pre = er[er['DATE'] <= '1985-12-31']
er_post = er[er['DATE'] > '1985-12-31']

# Run the regression for the two subsamples
model_pre = sm.OLS(er_pre['Log_ER'], er_pre[['Constant', 'Time']])
results_pre = model_pre.fit()

model_post = sm.OLS(er_post['Log_ER'], er_post[['Constant', 'Time']])
results_post = model_post.fit()

# Compute the sum of squared residuals for the two subsamples
SSR_pre = np.sum(results_pre.resid**2)
SSR_post = np.sum(results_post.resid**2)

# Compute the total sum of squared residuals
SSR_total = np.sum((er['Log_ER'] - np.mean(er['Log_ER']))**2)

# Compute the Chow test statistic
Chow = ((SSR_total - (SSR_pre + SSR_post)) / 2) / ((SSR_pre + SSR_post) / (len(er) - 4))
Chow
```

[17] ✓ 0.0s Python

... 2518.029081609119

Clearly, we find that we reject the null hypothesis! The same conclusion was reached through the use of dummy variables.

- More often we don't know when the break occurred but we suspect that it was sometime between date  $t_0$  and  $t_1$ . This is a much more interesting question!
- The Chow test can be modified to handle this by testing for breaks at all possible dates between  $t_0$  and  $t_1$  and then taking the largest of the resulting F-tests.
- This is called the **Quandt test** (or sup-Wald test)
- Because the Quandt test is the *largest* of an individual F-test, its distribution is not the same as the individual tests but it will have its own distribution.

## Digression: Use of Dummy Variables

Dummy variables are very useful. Once you understand the role of the dummy variables and the interactive terms, we can apply the same principle to many questions:

- Test if the demand is on average different at different times of the years (i.e. seasonality, more examples later).
- Test if the elasticity of demand is different at different times of the year (car sales example etc).
- Isolate specific times of the year

Remember the question before. How many data points do we use in each of the periods we isolate?

# Seasonality

---

- Seasonality is a pattern that repeats every year.
- From a micro point of view, seasonality comes from links of technology, preferences and institutions to the calendar. (examples?).
- We will only look at deterministic seasonality = repetition is exact and predictable.
- Seasonality is a very typical component of Time Series.

# How do we deal with Seasonality?

1. If we are interested in forecasting non seasonal fluctuations we may want to remove seasonality and work with seasonally adjusted series.

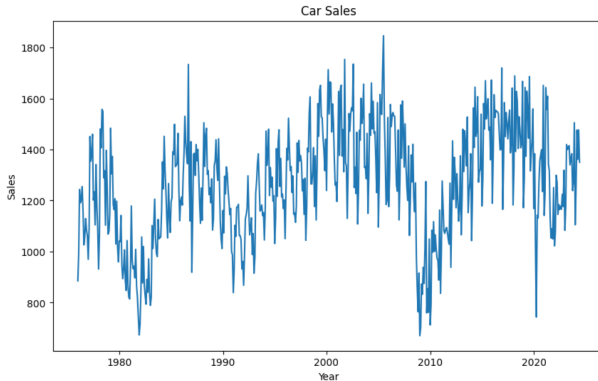
Do we really want to do this? In general we want to forecast all variations in the series.

2. We can take seasonality into account in our forecast and model seasonality.

Often data are already seasonally adjusted (SA), be careful when you download data which kind of data you really want. Data that is SA is passed through a complicated filter.

# How to Model Seasonality

- To model seasonality, we use dummy variables!
- To motivate, let us use the cars.csv dataset which spans from 1967:1 to 2024:6
- It looks like the sales of cars are very different.



## Reminder: Dummy Variable Trap

- If you run your regression with the constant, you need to define one less dummy than categories (in the example, 4 seasons  $\implies$  3 dummy variables).
- Alternatively, you can define the same number of dummies as categories and then run the regression without the constant.



# Intercept and Slope in a Dummy Model

We formulate a pure seasonal dummy model w/ one explanatory variable (for now)

$$cars_t = c + \beta_1 Q_1 + \beta_2 Q_2 + \beta_3 Q_3 + \beta_4 price_t + u_t$$

Clearly from here, we know that

$$\mathbb{E}(cars_t | winter) = c + \beta_1 + \beta_4 price$$

$$\mathbb{E}(cars_t | spring) = c + \beta_2 + \beta_4 price$$

$$\mathbb{E}(cars_t | summer) = c + \beta_3 + \beta_4 price$$

$$\mathbb{E}(cars_t | fall) = c + \beta_4 price$$

The intercept is different for each season. What does this mean in economic terms?

# Some Results

OLS Regression Results						
Dep. Variable:	Car Sales	R-squared:		0.191		
Model:	OLS	Adj. R-squared:		0.185		
Method:	Least Squares	F-statistic:		34.00		
Date:	Tue, 20 Aug 2024	Prob (F-statistic):		1.69e-25		
Time:	15:28:13	Log-Likelihood:		-3906.6		
No. Observations:	582	AIC:		7823.		
Df Residuals:	577	BIC:		7845.		
Df Model:	4					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
Constant	826.2518	46.051	17.942	0.000	735.804	916.700
Q1	-11.4374	23.436	-0.488	0.626	-57.467	34.592
Q2	140.0223	23.435	5.975	0.000	93.995	186.050
Q3	69.5385	23.556	2.952	0.003	23.272	115.805
Car Price	2.9879	0.330	9.048	0.000	2.339	3.637
Omnibus:		2.671	Durbin-Watson:			0.858
Prob(Omnibus):		0.263	Jarque-Bera (JB):			2.658
Skew:		-0.165	Prob(JB):			0.265
Kurtosis:		2.973	Cond. No.			768.

Notes:

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

For a given car price, the average sales are higher in the second and third quarters than in the fourth quarter.

## Some Extensions

- What is the economic intuition of our estimates?
- What other variable/s would you include in this model?
- Do you think that differences in the average demand for cars across seasons is the only kind of seasonality we need to take into account?

- It may also be that not only is the sales higher in Q2 for a given price but also that the demand for cars is more sensitive to the interest rate in the different seasons!
- We want to model possible shifts in the slope of the regression line for different seasons
- We can do this by looking at the interaction of the dummy variables with the price variable

# Interactive Dummy Variables

Consider this simple model with intercept and slope dummy

$$cars_t = c + \beta_1 Q1 + \beta_2 Q2 + \beta_3 Q3 + \beta_4 p_t + \beta_5 (Q1 \cdot p_t) + \beta_6 (Q2 \cdot p_t) + \beta_7 (Q3 \cdot p_t) + u_t$$

```
=====
                        OLS Regression Results
=====
Dep. Variable:          Car Sales   R-squared:                0.195
Model:                  OLS        Adj. R-squared:              0.185
Method:                 Least Squares   F-statistic:             19.84
Date:                   Tue, 20 Aug 2024   Prob (F-statistic):      7.26e-24
Time:                   16:13:23         Log-Likelihood:          -3905.1
No. Observations:      582             AIC:                    7826.
Df Residuals:          574             BIC:                    7861.
Df Model:               7
Covariance Type:       nonrobust
=====
                        coef    std err          t      P>|t|      [0.025      0.975]
-----
Constant             748.0431     89.883      8.322     0.000     571.503     924.583
Q1                   142.5440     123.148      1.157     0.248    -99.332     384.420
Q2                   279.7676     124.202      2.253     0.025     35.822     523.713
Q3                   72.6679     126.668      0.574     0.566    -176.122     321.457
Car_Price             3.5895       0.679      5.283     0.000      2.255      4.924
Q1_Price             -1.1878       0.932     -1.274     0.203     -3.019      0.644
Q2_Price             -1.0761       0.939     -1.146     0.252     -2.921      0.769
Q3_Price             -0.0206       0.960     -0.021     0.983     -1.907      1.865
=====
Omnibus:                 2.873   Durbin-Watson:           0.850
Prob(Omnibus):           0.238   Jarque-Bera (JB):         2.879
...
Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
=====
```

What is the underlying model for each season?

# Forecasting with Seasonality

- Given that we are only looking at deterministic seasonality, forecasting with seasonal dummy variables is very easy! (i.e. we know exactly how the dummy variable looks in the future)
- You can predict what the value for the dummy will be next quarter and you can use the estimated values for your parameters to forecast.
- We can re-estimate using only part of our sample and see how we are doing in forecasting.
- We don't have a good  $R^2$  so most likely we will not do too well in forecasting.

# Cycles

---

- Cycles include any sort of dynamics that is not captured by trend and seasonality.
- This includes dynamic, persistence, and any way in which the present is linked to the past or the future.



- We can explain this as thinking that the error terms/data are positively correlated with each other over time.

- We can explain this as thinking that the error terms/data are positively correlated with each other over time.
- One good year will tend to be followed by another, and one bad year will be followed by another bad year.

- We can explain this as thinking that the error terms/data are positively correlated with each other over time.
- One good year will tend to be followed by another, and one bad year will be followed by another bad year.
- Eventually, of course, something unusual will happen and a bad year will be followed by a good year. This represents the end of the recession and the start of the next boom.

- We can explain this as thinking that the error terms/data are positively correlated with each other over time.
- One good year will tend to be followed by another, and one bad year will be followed by another bad year.
- Eventually, of course, something unusual will happen and a bad year will be followed by a good year. This represents the end of the recession and the start of the next boom.
- The economy is subject to *shocks*! These shocks move the economy up or down in the period that occurs, then persist for a period of time.

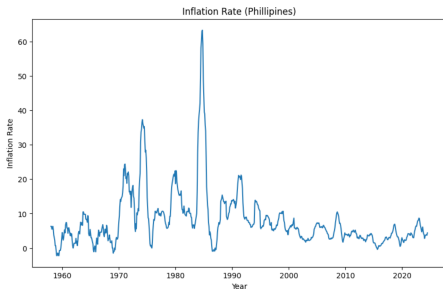
- We can explain this as thinking that the error terms/data are positively correlated with each other over time.
- One good year will tend to be followed by another, and one bad year will be followed by another bad year.
- Eventually, of course, something unusual will happen and a bad year will be followed by a good year. This represents the end of the recession and the start of the next boom.
- The economy is subject to *shocks*! These shocks move the economy up or down in the period that occurs, then persist for a period of time.
- The shocks are represented by the values of the error terms.

- We can explain this as thinking that the error terms/data are positively correlated with each other over time.
- One good year will tend to be followed by another, and one bad year will be followed by another bad year.
- Eventually, of course, something unusual will happen and a bad year will be followed by a good year. This represents the end of the recession and the start of the next boom.
- The economy is subject to *shocks*! These shocks move the economy up or down in the period that occurs, then persist for a period of time.
- The shocks are represented by the values of the error terms.
- Think of technology shocks slowly moving to different sectors.

- Before modeling cycles, we need to define when a time series is *stationary* since we can only deal (for the most part) with data that is stationary.
- Think about a time series. We observe a part of the path of the series. In theory, a time series begins in the infinite past and continues in the infinite future and we only observe a small part of it.
- Since we only observe a finite period, we would like the series to be stable over all periods.
- **Stationarity** is then very important!

# Autocorrelation

When we think of cycles, we often think about *autocorrelation*, the correlation of a series with itself lagged.



	Inflation Rate	Change in the Inflation Rate
Lag 1	0.983451	0.443811
Lag 2	0.952210	0.311938
Lag 3	0.910634	0.233348
Lag 4	0.861315	0.186828

- The inflation rate is **highly serially correlated** since  $\rho(1) = 0.98$
- Last month's inflation rate contains much information about this month's inflation. Moreover, the plot is dominated by multiyear swings but there are still surprise movements!



# How do we model cycles?

- The best way to capture the fact that things move slowly is to include a lagged dependent variable as an explanatory variable (ala  $AR(1)$ ). *Why the name?*
- I could also add two lags  $\implies AR(2)$
- Or we can use  $MA$  or  $ARMA$  models! We will learn to see which model/s are best when
- You can really only talk about  $AR$  terms when you have stationary data. If you don't (unit roots) you have to transform the data so they are stationary (i.e. through differencing)

- Most time series will need **at least** an  $AR(1)$ , often not more than an  $AR(2)$
- Once you add lags and other variables, the coefficients may become insignificant. Intuition?
- $R^2$  will increase significantly!
- $AR$  terms are your best friends when it comes to forecasting.
- They are at the core of time series analysis.
- **Always** look at the residuals or the ACF or PACF of the residuals to make sure there is no serial correlation in the errors. If there is, add one more lag.

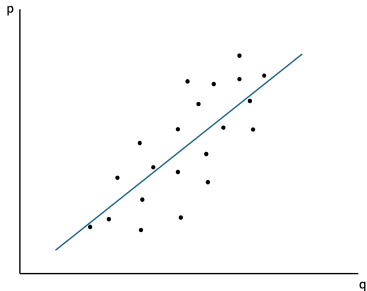
- What if there are breaks? → Testing
- How do we forecast and evaluate the forecasts?
- What if there is autocorrelation in the variance instead of the mean  
→ ARCH and GARCH
- Extension to multivariate → VAR and SVAR
- What if the data is non stationary? → Unit Root, Unit Root Testing, Cointegration, VECM
- Estimation in the Frequency Domain
- Kalman Filter
- ... We will see how much time we have!

# Structural vs Time Series Models

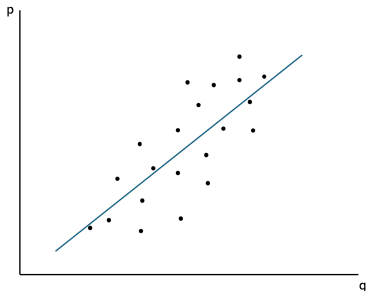
---

# Simple Intuition

- If we regress  $p$  on  $q$ , the estimated regression line will be the blue line. Is this the *demand curve* or the *supply curve*?

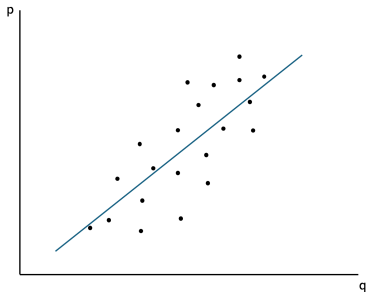


# Simple Intuition



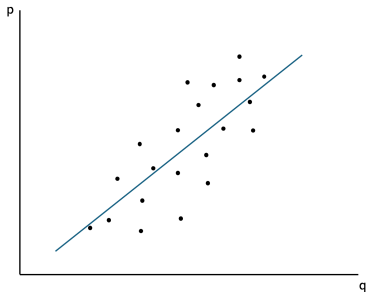
- If we regress  $p$  on  $q$ , the estimated regression line will be the blue line. Is this the *demand curve* or the *supply curve*?
- The estimated slope is not the slope of the demand function but a combination of the true slope of demand and supply

# Simple Intuition



- If we regress  $p$  on  $q$ , the estimated regression line will be the blue line. Is this the *demand curve* or the *supply curve*?
- The estimated slope is not the slope of the demand function but a combination of the true slope of demand and supply
- The estimator is not correct, it is biased! The slope is not estimating the elasticity. There is a simultaneity bias!

# Simple Intuition



- If we regress  $p$  on  $q$ , the estimated regression line will be the blue line. Is this the *demand curve* or the *supply curve*?
- The estimated slope is not the slope of the demand function but a combination of the true slope of demand and supply
- The estimator is not correct, it is biased! The slope is not estimating the elasticity. There is a simultaneity bias!
- The prices are a **combined** effect of supply and demand, because both are price contingent but also influence prices *simultaneously*!



- Suppose you are interested in the *price elasticity of demand*.

# Simple Intuition

- Suppose you are interested in the *price elasticity of demand*.
  - If you set prices based on the expectation that you will have an increase in demand...

- Suppose you are interested in the *price elasticity of demand*.
  - If you set prices based on the expectation that you will have an increase in demand...
  - You will not be estimating the elasticity by simply regressing demand on prices

# Simple Intuition

- Suppose you are interested in the *price elasticity of demand*.
  - If you set prices based on the expectation that you will have an increase in demand...
  - You will not be estimating the elasticity by simply regressing demand on prices
- If we really want to estimate the slope of the demand function, we need to make sure that the RHS variable is "exogenous/predetermined" (e.g. say weather)

- But if we only care about **forecasting** prices and/or quantities, then we could just regress  $q$  on  $p$  or vice-versa and get good forecasts even if we can't exactly estimate the demand function

- But if we only care about **forecasting** prices and/or quantities, then we could just regress  $q$  on  $p$  or vice-versa and get good forecasts even if we can't exactly estimate the demand function
- In pure time series, we do not worry about giving econ interpretations to coefficients because we are only interested in forecasting and not the structural parameters

# Simple Intuition

- But if we only care about **forecasting** prices and/or quantities, then we could just regress  $q$  on  $p$  or vice-versa and get good forecasts even if we can't exactly estimate the demand function
- In pure time series, we do not worry about giving econ interpretations to coefficients because we are only interested in forecasting and not the structural parameters
- What you do is a combination of these two which introduces difficulties in interpreting your coefficients.

The models you currently use somewhere in between: prices are in a way predetermined, gas cost is not. → Your models are a mix of the two



The models you currently use somewhere in between: prices are in a way predetermined, gas cost is not. → Your models are a mix of the two

- If you want to interpret the price elasticities very strictly, you need to make sure  $p$  is exogenous! ( $\mathbb{E}[p|\varepsilon] = 0$ )

The models you currently use somewhere in between: prices are in a way predetermined, gas cost is not. → Your models are a mix of the two

- If you want to interpret the price elasticities very strictly, you need to make sure  $p$  is exogenous! ( $\mathbb{E}[p|\varepsilon] = 0$ )
- If you are not sure, then DO NOT give too much weight on the estimates of price elasticity, as a bias is probably there!

The models you currently use somewhere in between: prices are in a way predetermined, gas cost is not. → Your models are a mix of the two

- If you want to interpret the price elasticities very strictly, you need to make sure  $p$  is exogenous! ( $\mathbb{E}[p|\varepsilon] = 0$ )
- If you are not sure, then DO NOT give too much weight on the estimates of price elasticity, as a bias is probably there!
- We cannot always give a **causal** interpretation unless we are sure that  $p$  is exogenous

The models you currently use somewhere in between: prices are in a way predetermined, gas cost is not. → Your models are a mix of the two

- If you want to interpret the price elasticities very strictly, you need to make sure  $p$  is exogenous! ( $\mathbb{E}[p|\varepsilon] = 0$ )
- If you are not sure, then DO NOT give too much weight on the estimates of price elasticity, as a bias is probably there!
- We cannot always give a **causal** interpretation unless we are sure that  $p$  is exogenous
- This will not affect the quality of your forecasts!