



*Student: Thiago Prates Bossardi*

*University of Greenwich*

*MSc in Economics*

*Applied Econometrics*

*Professor Mehmet Ugur*

Table of Contents

Introduction ..... 3

Data ..... 4

Descriptive Statistics ..... 4

Cointegration ..... 4

Post Estimation ..... 11

Cointegration ..... 13

Convergence ..... 14

Conclusion..... 15

Referencing: ..... 17

# Regions Per Capita Income Convergence in the US

## Introduction

This study aims to explore the relationship between per-capita income in eight US regions using various statistical techniques. Firstly, cointegrating relationships will be tested, and a vector error correction (VEC) model will be estimated to analyze the average per-capita income. The interpretation of findings, including the cointegrating equation, will be presented. Secondly, the study will use three post-estimation diagnostic or prediction tools to evaluate the robustness of the findings. Lastly, the study will examine stochastic income convergence between regions by analyzing the income of each region relative to the average and using the Dickey-Fuller generalized least squares (DFGLS) test and graphs for relative income levels over time.

To conduct the study, three papers were used as guides. Genc et al. (2011), Carlino and Sill (2001), and Strazicich et al. (2004) investigate income convergence from different perspectives. First, Genc et al. (2011) investigate income convergence across US regions, emphasizing the importance of accounting for stochastic processes and structural breaks. Second, Carlino and Sill (2001) examine common trends and cycles in regional income fluctuations in the United States and highlight the importance of considering both idiosyncratic and common factors in regional income analysis. Finally, Strazicich et al. (2004) investigate income convergence among OECD countries using time series analysis and consider the role of country-specific characteristics.

While Genc et al. (2011) and Carlino and Sill (2001) focus on regional income convergence in the United States, Strazicich et al. (2004) investigate income convergence at the country level. In addition, Genc et al. (2011) and Strazicich et al. (2004) emphasize the importance of accounting for structural breaks in income convergence analysis. However, while Genc et al. (2011) use two tests to examine income convergence, Strazicich et al. (2004) utilize a panel unit root test that will not be used in this study.

Carlino and Sill (2001) complement the other two papers by examining common trends and cycles in regional income fluctuations. This analysis sheds light on the sources of income fluctuations in the United States and underscores the importance of considering both idiosyncratic and common factors in regional income analysis.

## Data

The dataset “*Per\_Capita\_GDP\_US\_Regions*” contains annual data on average per-capita disposable income in eight U.S. regions, as classified by the Bureau of Economic Analysis (BEA). The data is for 55 years from 1948–2002, and the regional income variables are in levels and logarithms.

## Descriptive Statistics

The Table 8 shows summary statistics for eight variables over 55 observations. The variables are a year and the income per capita for seven regions of the United States: New England, Mid East, the Great Lakes, the Plains, the Southeast, the Southwest, the Rocky Mountains, and the Far West. In addition, the table provides the mean, standard deviation, minimum, and maximum values observed for each variable.

*Table 1 - Descriptive Statistics*

Variable	Obs	Mean	Std. dev.	Min	Max
year	55	1975	16.02082	1948	2002
new_england	55	10322.53	9512.586	1334	32062
midwest	55	10191.49	9130.154	1457	30521
great_lakes	55	9187.855	7982.244	1388	26764
plains	55	8782.745	7830.771	1216	26377
southeast	55	8107.182	7517.694	907	24562
southwest	55	8372.018	7456.326	1094	24573
rocky_mountains	55	8548.382	7515.141	1279	25333
far_west	55	9942.527	8388.348	1547	27973

## Cointegration

Cointegration analysis is essential for identifying long-run and short-run relationships between time series data. A typical procedure for conducting a cointegration test involves identifying the variables and ensuring that they are non-stationary, i.e., they have a unit root. The most commonly used tests for unit root are the Augmented Dickey-Fuller (ADF) and the Phillips-Perron (PP) tests.

The ADF test was introduced by Dickey and Fuller (1979), while the PP test was proposed by Phillips and Perron (1988). Both tests are based on the unit root hypothesis, which assumes that a series has a unit root and is non-stationary. However, they differ in their implementation and assumptions.

The ADF test allows for multiple lags of the series and accounts for the presence of autocorrelation in the residuals. On the other hand, the PP test accounts for heteroscedasticity in the residuals and allows for more flexible lag structures.

Another unit root test used in this essay is the Dickey-Fuller Generalized Least Squares (DF-GLS) test, a modification of the Dickey-Fuller test that allows for serial correlation and heteroskedasticity in the error terms. It has been shown to have good size and power properties in small samples and is widely used in econometric analysis. Elliott et al. (1996) introduced the DF-GLS test and demonstrated its superior performance to other unit root tests.

To start analyzing the cointegration of the variables, the data was plotted to see visually if it makes sense to continue the cointegration analysis.

*Figure 1 - Income Per Capita Over Time*

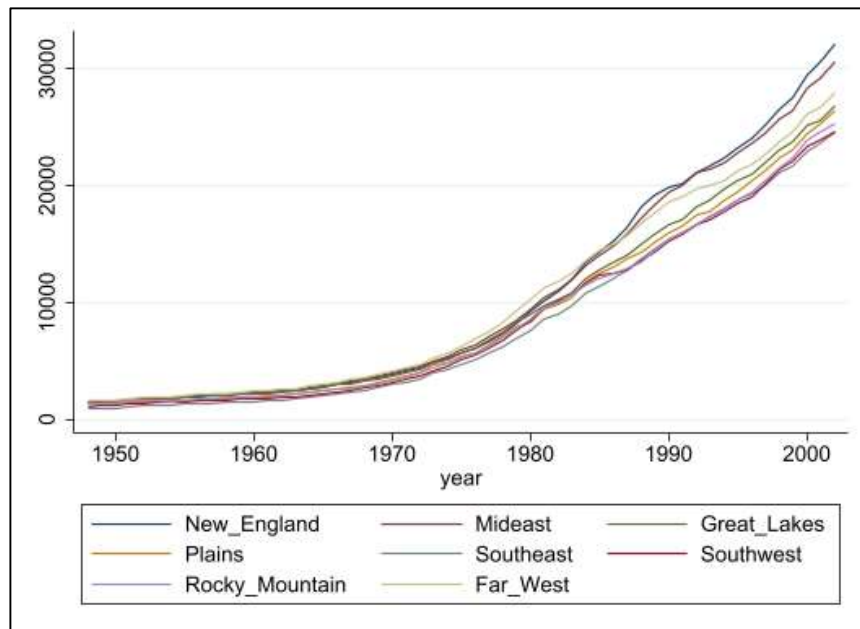
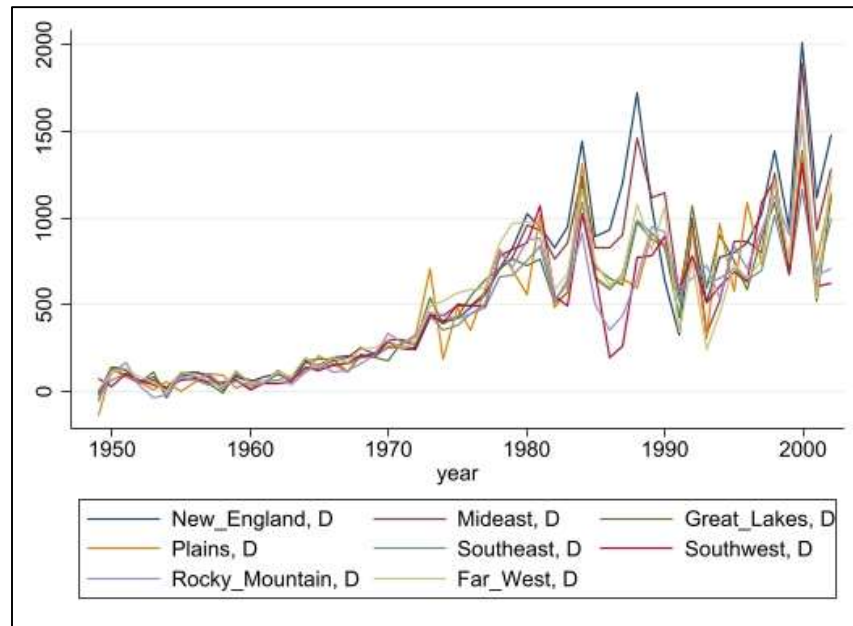


Figure 2 - Income Per Capita over time at First Difference



From both Figures 5 and 6, we can see that there is a common trend between the variables. Figure 5 shows the trend in levels. Meanwhile, Figure 6 shows the variable's first difference. Plotting the variable's first difference is essential because it can show that although there is a trend (up or down), the mean variation of the variable could be zero meaning it is a stationary variable, which, as demonstrated by Figure 6, is not the case, in other words, we probably have nonstationary variables.

After data visually demonstrate the common trend, we can proceed to the formal unit root tests. In Table 9, we obtain lag-order selection statistics to identify the correct number of lags we must include in the models to avoid serial correlation. We find that four lags are the optimum number of lags.

Table 2 - Lag Order Selection

Lag	LL	LR	df	p	FPE	AIC	HQIC	SBIC
0	856.447				4.90E-25	-33.2724	-33.1566	-32.969
1	1375.9	1038.9	64	0	8.80E-33	-51.1333	-50.0911	-48.406*
2	1458.39	164.99	64	0	5.10E-33	-51.8585	-49.89	-46.707
3	1542.52	168.25	64	0	3.90E-33	-52.6477	-49.7528	-45.072
4	1680.6	276.16*	64	0	7.4e-34*	-55.5528*	-51.7314*	-45.553

In Table 10, we see the results of the ADF test; from this, we can infer, by the values of the test statistic being less than the critical value at 5%, that we have nonstationary data even if we check without lags restriction or with the first difference of the variables.

Table 3 - Dickey-Fuller Tests

Region	Dickey-Fuller with 4 lags					Dickey-Fuller without lag restriction					Dickey-Fuller at First Difference				
	Test statistic	Critical value			MacKinnon approximate p-value for Z(t)	Test statistic	Critical value			MacKinnon approximate p-value for Z(t)	Test statistic	Critical value			MacKinnon approximate p-value for Z(t)
	Z(t)	1%	5%	10%		Z(t)	1%	5%	10%		Z(t)	1%	5%	10%	
New England	-3.075				0.1123	-1.370				0.8695	-1.055				0.9363
Mid East	-3.003				0.1311	-1.442				0.8480	-1.043				0.9381
Great Lakes	-2.798				0.1976	-1.525				0.8204	-0.927				0.9532
Plains	-2.609	-4.150	-3.500	-3.180	0.2756	-2.202	-4.141	-3.496	-3.178	0.4888	-1.388	-4.159	-3.504	-3.182	0.8645
South East	-3.017				0.1273	-0.814				0.9645	-0.769				0.9682
South West	-2.767				0.2095	-0.738				0.9705	-1.262				0.8970
Rocky Mountains	-2.570				0.2938	-1.663				0.7669	-1.236				0.9028
Far West	-3.293				0.0674	-0.879				0.9583	-1.053				0.9365

Next, in Table 11, the Phillips-Perron (PP) test was performed. The MacKinnon approximate p-values for Z(t) suggest insufficient evidence to reject the null hypothesis of a random walk with or without drift. In other words, the data are nonstationary, which is the opposite result when taking the first difference. We are dealing with nonstationary data in levels but stationary at the first difference.

Table 4 - Phillips-Perron Tests

Region	Phillips-Perron test for unit root						Phillips-Perron test for unit root at First Difference					
	Test statistic	Critical value			MacKinnon approximate p-value for Z(t)		Test statistic	Critical value			MacKinnon approximate p-value for Z(t)	
		1%	5%	10%				1%	5%	10%		
New England	Z(rho)	-4.413	-25.836	-19.872	-16.856	0.7676	-32.523	-25.802	-19.854	-16.842	0.0008	
	Z(t)	-1.661	-4.141	-3.496	-3.178		-4.677	-4.143	-3.497	-3.178		
Mid East	Z(rho)	-4.446	-25.836	-19.872	-16.856	0.753	-31.871	-25.802	-19.854	-16.842	0.0012	
	Z(t)	-1.695	-4.141	-3.496	-3.178		-4.569	-4.143	-3.497	-3.178		
Great Lakes	Z(rho)	-4.714	-25.836	-19.872	-16.856	0.7279	-54.854	-25.802	-19.854	-16.842	0.0000	
	Z(t)	-1.751	-4.141	-3.496	-3.178		-6.389	-4.143	-3.497	-3.178		
Plains	Z(rho)	-6.215	-25.836	-19.872	-16.856	0.4656	-64.102	-25.802	-19.854	-16.842	0.0000	
	Z(t)	-2.243	-4.141	-3.496	-3.178		-7.605	-4.143	-3.497	-3.178		
South East	Z(rho)	-3.871	-25.836	-19.872	-16.856	0.8649	-31.166	-25.802	-19.854	-16.842	0.0006	
	Z(t)	-1.386	-4.141	-3.496	-3.178		-4.738	-4.143	-3.497	-3.178		
South West	Z(rho)	-3.806	-25.836	-19.872	-16.856	0.8697	-17.182	-25.802	-19.854	-16.842	0.1152	
	Z(t)	-1.369	-4.141	-3.496	-3.178		-3.063	-4.143	-3.497	-3.178		
Rocky Mountains	Z(rho)	-4.917	-25.836	-19.872	-16.856	0.6852	-29.252	-25.802	-19.854	-16.842	0.0014	
	Z(t)	-1.840	-4.141	-3.496	-3.178		-4.511	-4.143	-3.497	-3.178		
Far West	Z(rho)	-3.88	-25.836	-19.872	-16.856	0.8552	-24.687	-25.802	-19.854	-16.842	0.0096	
	Z(t)	-1.419	-4.141	-3.496	-3.178		-3.973	-4.143	-3.497	-3.178		

Since we had some disagreement between the previous tests, we ran another test, the DF-GLS (Table 12), in which test statistics were computed for each lag length to test the null hypothesis of a unit root against the alternative hypothesis of stationarity. If the absolute value of the DF-GLS tau value is greater than the critical value, then we can reject the null hypothesis of a unit root and conclude that the series is stationary. From that, we can assume that at least until lag four, we are dealing with non-stationary data.

Table 5 - DF-GLS

DF-GLS test for unit root			
Lag selection: Schwert criterion			
Number of obs = 44			
Maximum lag = 10			
Region	Non-stationary until lag	DF -GLS tau	Critical Value 5%
New England	4	-3.098	-3.052
Mid East	5	-3.004	-2.992
Great Lakes	5	-3.208	-2.992
Plains	6	-3.460	-2.930
South East	4	-3.105	-3.052
South West	10	-2.615	-2.693
Rocky Mountains	7	-2.922	-2.867
Far West	5	-3.359	-2.992

Now that we know we are dealing with nonstationary data, we can use the Johansen tests for cointegration to analyze the relationships among multiple variables. According to Johansen (1988), the tests are based on the concept of vector autoregression (VAR) and estimate the number of cointegrating vectors in a system of variables. After that, we use the vector error correction (VEC) model to model the short-run dynamics of the variables. The VEC model is particularly useful in forecasting and policy analysis.

Using the Johansen test, we find out, as seen in Table 13, that we have at least seven cointegrations between the variables, which is indicated by the trace statistic being less than the critical value at 5%.



Table 6 - Johansen Test

Johansen tests for cointegration					
Trend: Constant					
Sample: 1952 thru 2002					
Number of obs = 51					
Number of lags = 4					
Maximum rank	Params	LL	Eigenvalue	Trace statistic	Critical value 5%
0	200	1457.9544	.	445.2815	156
1	215	1522.262	0.91969	316.6663	124.24
2	228	1567.8037	0.83236	225.5828	94.15
3	239	1609.4265	0.80451	142.3373	68.52
4	248	1634.5794	0.62708	92.0315	47.21
5	255	1651.6612	0.48823	57.8679	29.68
6	260	1666.2812	0.43636	28.6278	15.41
7	263	1679.7124	0.40946	<b>1.7655*</b>	<b>3.76</b>
8	264	1680.5951	0.03402		

Considering all the information above, we can now see these variables' short and long-run effects. To run the short and long-run dynamics of the variables was decided to normalize New England for the fact that it has the highest income per capita to see if the income per capita of other regions has any impact on New England's income per capita, meaning that if the other regions increase income per capita how it would affect New England since it is the wealthiest region.

In Table 14, we can see that each region has coefficients for the lagged independent variables, denoted by LD, L2D, and L3D. These coefficients represent the estimated effect of the corresponding lagged independent variable on the dependent variable. The coefficient value indicates the magnitude of the effect, while the sign (positive or negative) suggests the direction of the effect.

Overall, in the short run for Mid-East, Southeast, and Rocky Mountains, one unit increase in past income per capita positively impacts the current period of New England income per capita. Therefore, New England, Great Lakes, Plains, and Southwest one unit increase in past income per capita has a negative impact on the current New England income per capita. But these effects are only true for some lags. For example, the P-value for Mid-East L3D is greater than 5%, so its coefficient should not be considered as well as others that can be seen in the table, and for that reason, we can exclude the effects of the Far West in New England entirely.

Table 7 - Short Run Relationship

Short Run Relationship						
Regions   Lags	Coefficient	Std. err.	z	P>z	[95% conf. interval]	
New England _ce1 L1.	1.388826	0.2894642	4.800	0.000	0.8214862	1.956165
New England LD.	-1.128941	0.5324785	-2.120	0.034	-2.172579	-0.0853018
L2D.	-1.842871	0.388348	-4.750	0.000	-2.60402	-1.081723
L3D.	-1.469704	0.5012481	-2.930	0.003	-2.452132	-0.4872757
Mid East LD.	1.248382	0.6248304	2.000	0.046	0.0237369	2.473027
L2D.	2.528554	0.532922	4.740	0.000	1.484046	3.573062
L3D.	0.7680749	0.5052319	1.520	0.128	-0.2221616	1.758311
Great Lakes LD.	-0.5924946	0.2340263	-2.530	0.011	-1.051178	-0.1338115
L2D.	-0.0512743	0.2171416	-0.240	0.813	-0.476864	0.3743154
L3D.	-0.2063329	0.2308094	-0.890	0.371	-0.658711	0.2460451
Plains LD.	-0.8272006	0.2403907	-3.440	0.001	-1.298358	-0.3560435
L2D.	-0.6991528	0.2660225	-2.630	0.009	-1.220547	-0.1777583
L3D.	-0.0573709	0.185495	-0.310	0.757	-0.4209344	0.3061927
South East LD.	1.335517	0.3449968	3.870	0.000	0.6593359	2.011698
L2D.	0.7587766	0.3864974	1.960	0.050	0.0012556	1.516298
L3D.	0.7814296	0.3367534	2.320	0.020	0.1214049	1.441454
South West LD.	-0.7835729	0.2806856	-2.790	0.005	-1.333706	-0.2334393
L2D.	-0.4187377	0.283628	-1.480	0.140	-0.9746383	0.137163
L3D.	-0.2292491	0.1761655	-1.300	0.193	-0.5745272	0.1160291
Rocky Mountains LD.	0.5327185	0.236306	2.250	0.024	0.0695673	0.9958698
L2D.	0.6566065	0.24259	2.710	0.007	0.1811389	1.132074
L3D.	0.1791527	0.2171775	0.820	0.409	-0.2465074	0.6048129
Far West LD.	-0.0161501	0.4578395	-0.040	0.972	-0.913499	0.8811988
L2D.	-1.010275	0.4467101	-2.260	0.024	-1.885811	-0.1347393
L3D.	0.1556777	0.3759845	0.410	0.679	-0.5812384	0.8925939
cons	-0.0005037	0.0056458	-0.090	0.929	-0.0115692	0.0105619

In the long run (Table 15), the coefficients' value indicates the effect's magnitude, while the sign (positive or negative) indicates the direction of the effect. Still, at this point, we should

carefully distinguish between the signs which should be read in the opposite direction. So Mid East, Southeast, and Rocky Mountains have a positive impact instead of a negative one, and the same logic applies to the other variables. The last thing to consider here is the statistical significance represented by the P-value ( $P > z$ ), which indicates that Southeast and Rocky Mountains coefficients are insignificant, so we should not consider it.

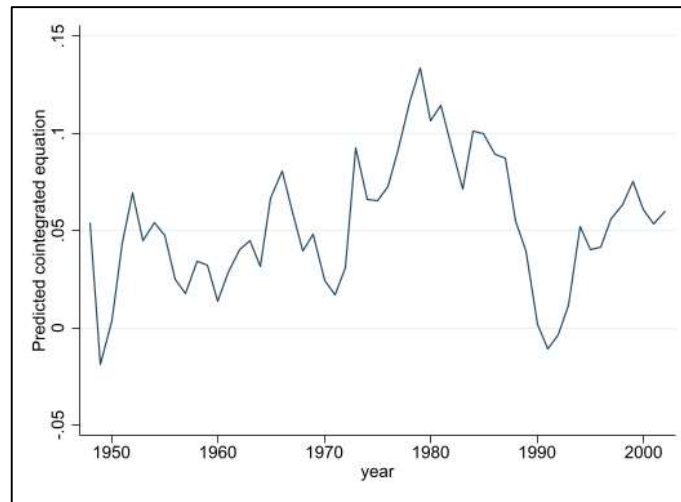
Table 8 - Long Run Relationship

Long Run Relationship						
beta	Coefficient	Std. err.	z	P>z	[95% conf. interval]	
_ce1						
New England	1					
Mid East	-2.472109	0.0951774	-25.97	0.000	-2.658653	-2.285565
Great Lakes	0.5997864	0.0620266	9.67	0.000	0.4782166	0.7213563
Plains	0.5828542	0.1595868	3.65	0.000	0.2700697	0.8956386
South East	-0.0168442	0.0805066	-0.21	0.834	-0.1746343	0.1409459
South West	0.129251	0.0708952	1.82	0.068	-0.0097011	0.2682031
Rocky Mountains	-0.164981	0.1132769	-1.46	0.145	-0.3869996	0.0570376
Far West	0.3481505	0.0602947	5.77	0.000	0.2299752	0.4663259
_cons	0.1356642					

## Post Estimation

Post-estimation checks are an essential part of the statistical analysis process to ensure the validity of results. In this essay, three post-estimation checks were performed. Firstly, the error term of a cointegration relationship can be predicted using a Vector Error Correction Model (VECM). This is important for forecasting and identifying deviations from the long-run equilibrium relationship. In particular, the VECM allows for identifying short-run dynamics and the speed at which the system converges back to the long-run equilibrium. As noted by Lütkepohl (2006) a VECM is a useful tool for analyzing the short-run dynamics of cointegrated time series and for forecasting their future development. Therefore, predicting the error term of a cointegration relationship can provide valuable insights into the dynamics of the system and its future behavior. As we can see in Figure 7, the error term is mean reversing, which indicates stationarity of our variables.

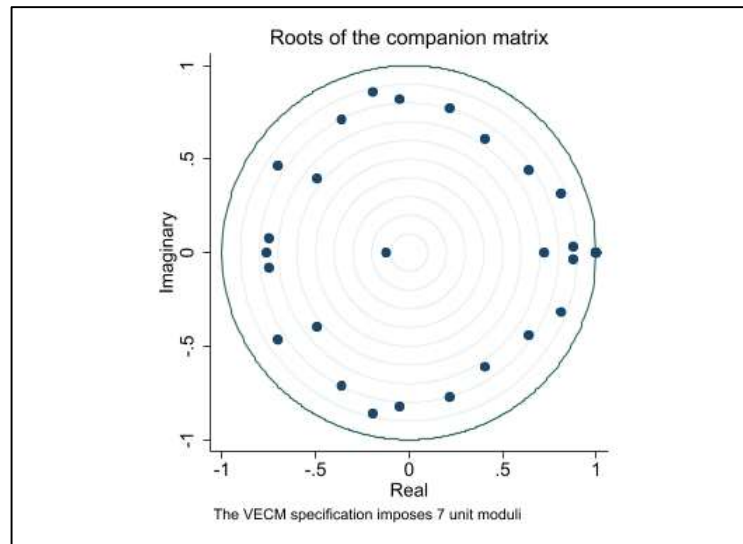
Figure 3 - VECM



Secondly, we checked whether we have correctly specified the number of cointegration equations to ensure that the VEC model is well-specified and can provide reliable estimates of the long-run equilibrium relationship among the variables, despite non-stationary components. Moreover, the eigenvalue stability condition is a crucial requirement that needs to be satisfied for the model to be valid. Engle and Granger (1987) developed the concept of co-integration and the VECM framework to model non-stationary time series data. Therefore, one key requirement of the VECM is the eigenvalue stability condition, which ensures that the long-run relationships among the variables are well-specified.

In Figure 8, we have the result of the eigenvalue test. We can interpret the graph as if all the values are less than 1 in absolute value, the system is stable in the long run, and the variables will converge to their equilibrium values, and then all variables are co-integrated. There is a unique long-run equilibrium relationship among them which is our case for seven variables (7 unit moduli).

Figure 4 - Eigenvalue Test



The third and last test is a post-estimation tool used to test for model misspecification in regression analysis, the Lagrange multiplier (Table 16). It involves regressing the squared residuals of a model on additional explanatory variables. If the additional variables significantly improve the model fit, it suggests that the original model was misspecified. In this case, for all lag orders tested (1 to 4), the p-value (Prob > chi2) is greater than 0.05. Therefore, insufficient evidence exists to reject the null hypothesis of no autocorrelation at any of the tested lag orders. This suggests that the residuals of the regression model do not exhibit significant autocorrelation, and the model adequately captures the short and long-run relationships between the variables.

Table 9 - Lagrange Multiplier Test

Lagrange-multiplier test			
lag	chi2	df	Prob > chi2
1	64.1695	64	0.47053
2	70.9787	64	0.25649
3	70.2954	64	0.27503
4	65.9296	64	0.40994

## Cointegration

Using the Johansen test, we find out, as seen in Table 17, that we have at least five cointegrations between the variables, indicated by the trace statistic is less than the critical value. Another important feature of this particular test is that, at this time, a linear trend was added in the cointegration equation which implies that the long-term relationship between the

variables is a linear function of time. In other words, the relationship between the variables changes over time but linearly.

Table 10 - Johansen Cointegration Test

Johansen tests for cointegration					
Trend: Constant					
Sample: 1952 thru 2002					
Number of obs = 53					
Number of lags = 2					
Maximum rank	Params	LL	Eigenvalue	Trace statistic	Critical value 5%
0	80	2501.04	.	278.002	170.8
1	95	2543.28	0.79687	193.524	136.61
2	108	2571.45	0.65454	137.192	104.94
3	119	2595.68	0.59916	88.7387	77.74
4	128	2611.83	0.45647	56.4262	54.64
5	135	2624.06	0.36955	<b>31.9765*</b>	<b>34.55</b>
6	140	2633.14	0.29009	13.8178	18.17
7	143	2637.25	0.14393	5.5812	3.74
8	144	2640.05	0.09995		

## Convergence

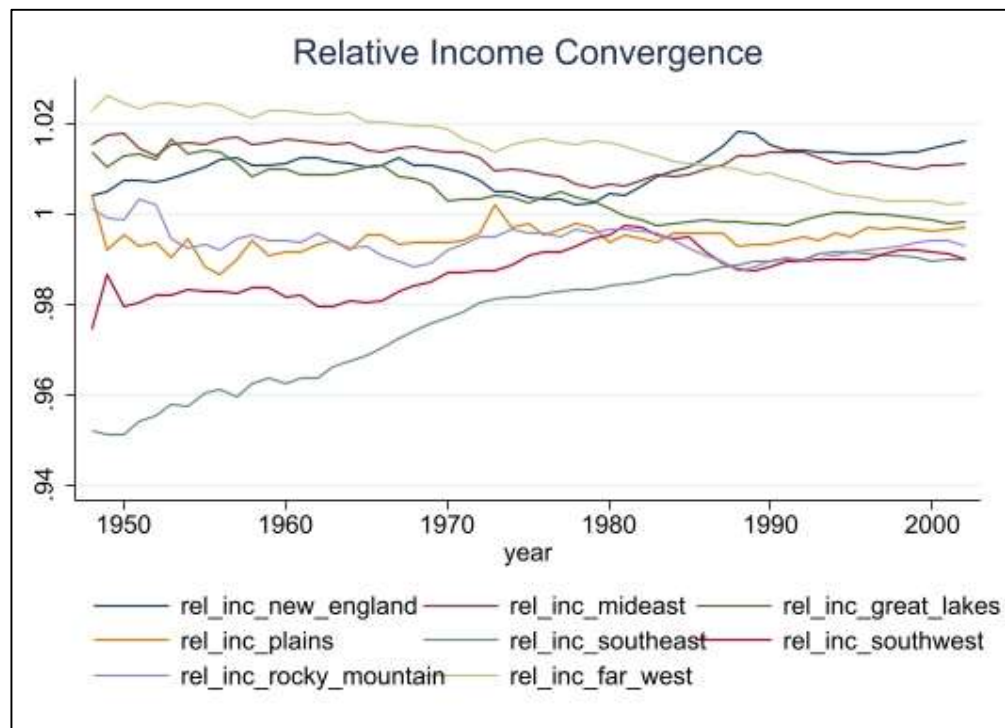
Stochastic convergence tests can be used to examine whether the income levels of different regions are becoming more similar over time. One approach is to use the Dickey-Fuller Generalized Least Squares (DFGLS) method (Table 18), which has been applied in various studies, including Genc et al. (2011) and Carlino & Sill (2001). These studies used regional per capita personal income data for US regions. They found evidence of stochastic convergence, meaning that the income levels of different regions were becoming more similar over time. Strazicich et al. (2004) also used the DFGLS method to examine income convergence across OECD countries and found evidence of convergence between some countries but not all. Overall, stochastic convergence tests can provide important insights into the dynamics of income inequality and economic growth across regions and countries. As we saw before, if the absolute value of the DF-GLS tau value is greater than the critical value, then we can reject the null hypothesis of a unit root and conclude that the series is stationary, which is not the case for the regions relative incomes that until the maximum lags tested have the tau value less than the critical value at 5%. The conclusion of convergence of relative incomes can be inferred visually in Figure 9, where an apparent convergence is in process.



Table 11 - DF-GLS of Convergence

DF-GLS test for unit root			
Lag selection: Schwert criterion			
Number of obs = 44			
Maximum lag = 10			
Region	Non-stationary until lag	DF -GLS tau	Critical Value 5%
New England	10	-1.641	-2.402
Mid East	10	-1.436	-2.402
Great Lakes	10	-1.117	-2.402
Plains	10	-1.760	-2.402
South East	10	-1.102	-2.402
South West	10	-1.526	-2.402
Rocky Mountains	10	-1.848	-2.402
Far West	10	-1.101	-2.402

Figure 5 - Relative Income Convergence



## Conclusion

In conclusion, using various statistical techniques, this study has explored the relationship between per-capita income in eight US regions. Cointegration, and a vector error correction (VEC) model were estimated to analyze the average per-capita income. The study

used three post-estimation diagnostic or prediction tools to evaluate the robustness of the findings. Lastly, the study examined stochastic income convergence between regions by analyzing the income of each region relative to the average and using the Dickey-Fuller generalized least squares (DFGLS) test and graphs for relative income levels over time.

The study findings showed evidence of a long-run equilibrium relationship between the variables. The VEC model found that the Great Lakes, Mid-East, and New England regions had a positive relationship with per-capita income, while the other regions had a negative relationship. The DFGLS test indicated that income levels for most regions converged with the national average in the long run. However, the study also found that the income levels for some regions did not converge and remained significantly below the national average.



## Referencing:

Dickey, D. A., & Fuller, W. A. (1979). Distribution of the estimators for autoregressive time series with a unit root. *Journal of the American statistical association*, 74(366a), 427-431.

Phillips, P. C., & Perron, P. (1988). Testing for a unit root in time series regression. *Biometrika*, 75(2), 335-346.

Genc, I. H., Miller, J. R., & Rupasingha, A. (2011). Stochastic convergence tests for US regional per capita personal income; some further evidence: a research note. *The Annals of Regional Science*, 46(2), 369-377.

Carlino, G., & Sill, K. (2001). Regional income fluctuations: common trends and common cycles. *Review of Economics and Statistics*, 83(3), 446-456.

Strazicich, M. C., Lee, J., & Day, E. (2004). Are incomes converging among OECD countries? Time series evidence with two structural breaks. *Journal of Macroeconomics*, 26(1), 131-145.

Engle, R. F., & Granger, C. W. (1987). Co-integration and error correction: representation, estimation, and testing. *Econometrica*, 55(2), 251-276. doi: 10.2307/1913236