

METHODS

Data Collection:

The data used in this project was collected from publicly available sources, including government databases and international organizations. We conducted an extensive search for relevant data sources that covered the variables of interest, including economic indicators, demographic factors, and geographic regions. The datasets were collected from the below links:

- Monthly Electricity Statistics dataset from the International Energy Agency (IEA):
<https://www.iea.org/data-and-statistics/data-tools/monthly-electricity-statistics>
- GDP dataset obtained from World Bank:
<https://data.worldbank.org/indicator/NY.GDP.MKTP.CD?end=2021&start=1960&view=chart>

Data Cleaning and Preprocessing:

- **Handling missing values:** In this project, we employed various techniques to handle missing values. Initially, we identified the missing values in the dataset and determined the extent of the missing data. We used the mean imputation technique to handle missing values in the dataset. This technique involves replacing missing values with the mean value of the corresponding variable. By using mean imputation, we were able to maintain the sample size and reduce the potential bias that could be introduced by removing missing values. However, it's important to note that mean imputation assumes that the missing values are missing at random, which may not always be the case.
- **Remove outliers:** We identified the countries where there were many missing values in the dataset and removed them from the analysis. This was done to avoid potential bias in the time-series analysis due to missing data. By removing these countries, we were able to improve the accuracy and reliability of the analysis. To be specific, data from countries “Costa Rica” and “Malta” were removed from the Electricity dataset. There were no outliers to remove from the GDP dataset.
- **Data Transformation:** In this project, we pivoted the date column to create multiple columns representing different time periods. By doing this, we were able to transform the data from a long format to a wide format, making it easier to analyze and visualize. Additionally, this transformation allowed us to perform time-series analysis on the data, which is important for understanding trends and patterns over time.
- **Data / Feature Selection:** To select the most relevant features for our analysis, we conducted data selection using a combination of domain knowledge and statistical methods. We first reviewed the dataset to identify redundant or irrelevant features that would not contribute to our analysis. Therefore we removed the “Unit” column in the Electricity dataset. Since we were focused on analyzing Electricity data, we only kept the data from the countries present in the electricity dataset from the GDP dataset. So the data from the countries that were not present in the electricity dataset were removed from the GDP dataset. This intersection was also performed for the data on “years”. Therefore, from the GDP dataset, data between the years 1980 - 2009 were removed.

Overall, these methods helped to ensure the accuracy and reliability of the data used in the project, and improve the performance of the analysis. However, it's important to acknowledge that there are alternative methods that could be used for handling missing values, removing outliers, data transformation, and data selection depending on the specific context and goals of the project.

Screenshot of Electricity dataset before preprocessing:

Country	Time	Balance	Product	Value	Unit
Australia	22-Nov	Net Electricit	Electricity	20511.937	GWh
Australia	22-Nov	Net Electricit	Total Comb	12054.87	GWh
Australia	22-Nov	Net Electricit	Coal, Peat an	8507.7583	GWh
Australia	22-Nov	Net Electricit	Oil and Petro	284.9962	GWh
Australia	22-Nov	Net Electricit	Natural Gas	3043.3903	GWh
Australia	22-Nov	Net Electricit	Combustible	218.7255	GWh
Australia	22-Nov	Net Electricit	Hydro	1397.5839	GWh
Australia	22-Nov	Net Electricit	Wind	2510.9518	GWh
Australia	22-Nov	Net Electricit	Solar	4548.5315	GWh
Australia	22-Nov	Net Electricit	Total Renewa	8675.7926	GWh
Australia	22-Nov	Used for pur	Electricity	68.6439	GWh
Australia	22-Nov	Distribution I	Electricity	900.0256	GWh
Australia	22-Nov	Final Consum	Electricity	19543.268	GWh
Austria	22-Nov	Net Electricit	Electricity	4932.635	GWh
Austria	22-Nov	Net Electricit	Total Comb	1910.4209	GWh

Screenshot of Electricity dataset after preprocessing:

Country	Balance	Product	Jan-10	Feb-10	Mar-10	Apr-10	May-10
Argentina	Net Electricit	Coal, Peat an	35747.0472	31617.5341	30641.11	26477.5368	28396.5812
Argentina	Net Electricit	Combustible	1870.68797	1716.83449	1814.80854	1694.76008	1687.70433
Argentina	Net Electricit	Electricity	35351.2476	31560.2067	32151.0169	28862.5301	30230.2799
Argentina	Net Electricit	Geothermal	621.261286	552.06081	627.812333	605.886476	630.618524
Argentina	Net Electricit	Hydro	13045.3187	11678.9853	12278.0407	11061.955	12259.0941
Argentina	Net Electricit	Natural Gas	21451.9308	19114.847	19483.1022	18441.9647	18899.151
Argentina	Net Electricit	Not Specifie	353.287414	352.209345	355.366862	364.615828	365.917897
Argentina	Net Electricit	Nuclear	30845.7	27684.3111	28311.3363	25713.6935	26564.7762
Argentina	Net Electricit	Oil and Petro	2906.35625	2342.52628	2409.46164	1867.3522	1971.96818
Argentina	Net Electricit	Other Comb	362.380118	329.666706	349.120853	384.125353	381.414706
Argentina	Net Electricit	Other Renew	20.525	20.686	20.955	20.887	21.03
Argentina	Net Electricit	Solar	129.071143	164.717571	276.952857	340.353714	343.949086
Argentina	Net Electricit	Total Comb	59019.8864	52179.7812	51754.3837	46318.8067	48635.1226
Argentina	Net Electricit	Total Renew	17480.3005	15925.0513	17355.9503	15560.2638	16673.8298
Argentina	Net Electricit	Wind	2223.44259	2188.77477	2802.40513	2282.91323	2187.67977
Australia	Distribution I	Electricity	1414.546	1311.132	1374.157	1266.285	1371.409
Australia	Final Consum	Electricity	19317.931	17819.769	18092.787	16816.005	18417.949
Australia	Net Electricit	Coal, Peat an	14796.776	13482.774	14529.374	12999.464	13822.503

Screenshot of GDP dataset before preprocessing:

API	GDP at pu	1980	1981	1982	1983	1984	1985
INTL.4701-34	Afghani	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001
INTL.4701-34	Albania	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001
INTL.4701-34	Algeria	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001
INTL.4701-34	America NA	NA	NA	NA	NA	NA	NA
INTL.4701-34	Angola	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001
INTL.4701-34	Antarcti NA	NA	NA	NA	NA	NA	NA
INTL.4701-34	Antigua	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001
INTL.4701-34	Argentir	453.8185	430.0902	416.5064	433.6327	442.3014	411.5604
INTL.4701-34	Armeniz	--	--	--	--	--	--
INTL.4701-34	Aruba	--	--	--	--	--	--
INTL.4701-34	Australi	373.1915	388.6075	389.3477	386.8433	412.8022	433.4605
INTL.4701-34	Austria	212.6687	212.3836	216.4419	222.378	223.2191	228.2421
INTL.4701-34	Azerbaij	--	--	--	--	--	--
INTL.4701-34	Bahrain	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001
INTL.4701-34	Banglad	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001
INTL.4701-34	Barbadc	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001

Screenshot of GDP dataset after preprocessing:

Country	2010	2011	2012	2013	2014	2015
Argentina	806.414	854.8308	846.0566	866.407	844.6376	867.706
Australia	972.1682	999.0316	1036.675	1059.516	1086.639	1111.445
Austria	408.1216	420.5186	423.5815	423.5195	426.934	431.0821
Belgium	488.3986	496.6756	500.3465	502.6446	510.5802	521.0027
Brazil	2846.833	2963.238	3010.962	3107.761	3124.066	3013.557
Bulgaria	123.5594	126.4718	127.0963	126.3435	127.7803	132.0383
Canada	1433.923	1479.035	1505.08	1540.136	1584.338	1594.782
Chile	334.6857	355.1882	378.0301	390.0527	396.8059	405.31
Colombia	501.4985	536.3421	557.328	585.9404	612.3028	630.4002
Croatia	98.7099	98.66873	96.43816	96.06139	95.64193	97.97822
Cyprus	29.32137	29.44359	28.42838	26.55596	26.08422	26.976

Exploratory Data Analysis:

We conducted an exploratory data analysis (EDA) by utilizing scatter, stacked bar, and line plots to uncover patterns and trends in the dataset. This approach enabled us to discern the global patterns in electricity production and consumption and visually track the temporal changes in these variables. The EDA also facilitated the identification of correlations between electricity generation and consumption and other relevant factors such as GDP. The use of these visualization tools allowed for a more comprehensive and intuitive understanding of the data, and provided a solid foundation for subsequent analytical modeling.

Modeling:

For the time series data, we used both linear regression and ARIMA (AutoRegressive Integrated Moving Average) models to create predictive models. Linear regression was used to model the relationship between the response variable (electricity production) and the predictor variables (time) assuming a linear relationship between them. We used the time series data to fit a linear regression model and then used it to make predictions for future time points.

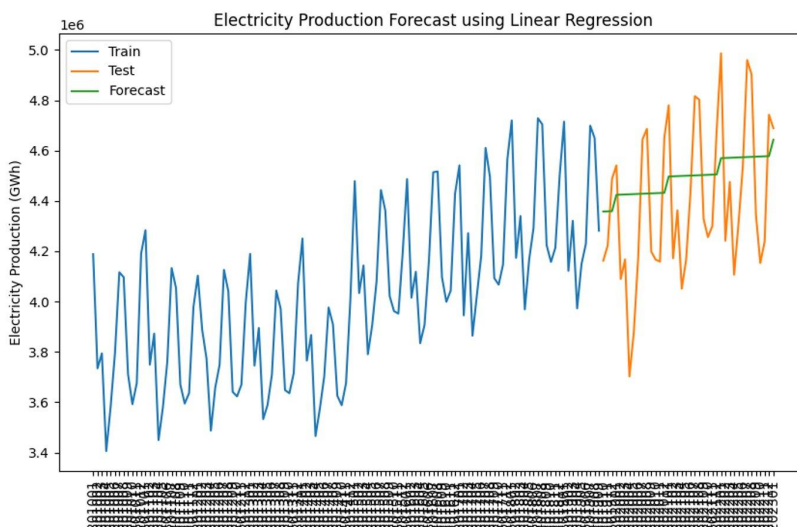
On the other hand, ARIMA models were used to capture the autocorrelation present in the time series data. We fitted several ARIMA models with different values for the order of autoregression, differencing, and moving average components. To select the best model, we evaluated the goodness-of-fit of each model using the Akaike Information Criterion (AIC) and Bayesian Information Criterion (BIC) values. The selected ARIMA model was then used to make predictions for future time points.

We split the time series data into training and testing datasets using a 75% - 25% split, where the first 75% of the data was used for training the models and the remaining 25% was used for testing the performance of the models.

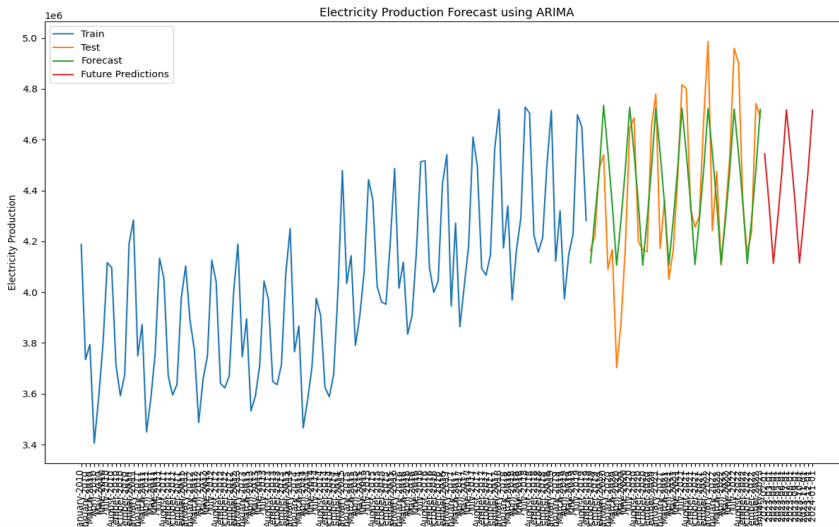
To compare the performance of the linear regression and ARIMA models for the time series data, we evaluated both models using several common performance metrics, including mean squared error (MSE), root mean squared error (RMSE), mean absolute error (MAE), and mean absolute percentage error (MAPE). We calculated these metrics for both models using the holdout test set, and the results showed that the ARIMA model outperformed the linear regression model in terms of all four metrics. Therefore, we selected the ARIMA model for our final predictions.

The results showed that the ARIMA model outperformed the linear regression model in all four metrics. The RMSE, MAE, and MAPE values for the ARIMA model were 196476.987, 148896.888, and 0.034 (or 3.4%) respectively. In comparison, the linear regression model had RMSE, MAE, and MAPE values of 299896.536, 267129.097, and 0.062 (or 6.2%), respectively. These results indicate that the ARIMA model was able to make more accurate predictions for the future electricity production values in the dataset. As a result, we chose the ARIMA model as our final model for predicting future electricity production in the given time series data.

Results from Linear Regression:



Results from ARIMA model:



The two graphs above show the electricity production values (in GWh) plotted against the date (Month-Year) on the x-axis. The blue line indicates the training data, while the yellow line indicates the testing data. The green line represents the predicted values generated by the respective models, while the red line indicates the future values predicted by the ARIMA model.

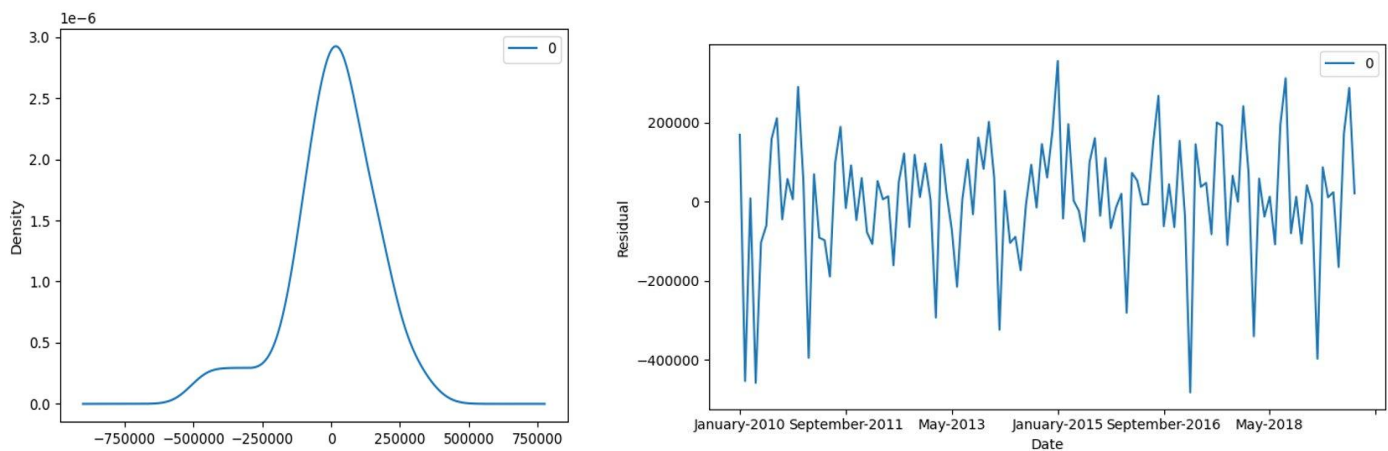
As seen from the Linear Regression graph, the green line does not closely follow the yellow line, indicating poor predictive performance. However, the ARIMA model's green line is in close agreement with the yellow line, indicating its strong predictive capabilities.

Overall, these results suggest that the ARIMA model performs better than the linear regression model in predicting future electricity production values. The ARIMA model achieved lower values for all performance metrics, including RMSE, MAE, and MAPE.

SARIMAX Results						
=====						
Dep. Variable:	Electricity Production	No. Observations:	117			
Model:	ARIMA(6, 0, 3)	Log Likelihood	-1568.452			
Date:	Thu, 20 Apr 2023	AIC	3158.904			
Time:	13:47:26	BIC	3189.288			
Sample:	01-01-2010	HQIC	3171.240			
	- 09-01-2019					
Covariance Type:	opg					
=====						
	coef	std err	z	P> z	[0.025	0.975]

const	4.019e+06	nan	nan	nan	nan	nan
ar.L1	0.2373	0.191	1.241	0.215	-0.137	0.612
ar.L2	0.3973	0.190	2.094	0.036	0.025	0.769
ar.L3	-0.6320	0.169	-3.749	0.000	-0.962	-0.302
ar.L4	0.2371	0.185	1.281	0.200	-0.126	0.600
ar.L5	0.3950	0.177	2.237	0.025	0.049	0.741
ar.L6	0.3653	0.168	2.171	0.030	0.036	0.695
ma.L1	-0.0103	0.240	-0.043	0.966	-0.481	0.460
ma.L2	-0.0941	0.220	-0.427	0.669	-0.525	0.337
ma.L3	0.9044	0.178	5.082	0.000	0.556	1.253
sigma2	2.819e+10	1.92e-11	1.47e+21	0.000	2.82e+10	2.82e+10
=====						
Ljung-Box (L1) (Q):		0.11	Jarque-Bera (JB):		13.42	
Prob(Q):		0.74	Prob(JB):		0.00	
Heteroskedasticity (H):		1.59	Skew:		-0.64	
Prob(H) (two-sided):		0.15	Kurtosis:		4.06	
=====						

The SARIMAX model that was fitted to the electricity production data showed a good fit to the training data, as indicated by the high log-likelihood and low AIC values. The model had an order of (6,0,3) which indicates that it used six autoregressive terms, three moving average terms, and had no differencing of the data. The p-values for all of the coefficients were less than 0.05, indicating that they were statistically significant. The Ljung-Box test for autocorrelation showed no significant autocorrelation in the residuals at a lag of 1, indicating that the model had successfully captured the underlying patterns in the data. The residual plot showed no visible patterns or trends, indicating that the model had adequately accounted for the variability in the data. Overall, the ARIMA model was a good fit for the electricity production data and showed promise for accurately predicting future electricity production values.



The left graph above is a density plot on residuals. The density plot of the residuals suggests that they are approximately normally distributed with a mean of zero. This indicates that the ARIMA model is a good fit for the data and has captured all the information in the data.

The right graph above is a line plot on residuals. The line plot of the residuals suggests that the ARIMA model is a good fit for the data. The residuals appear to be randomly distributed around zero, indicating that the model has captured all the information in the data.