	Introduction The National Longitudinal Survey of Youth 1997-2011 dataset is one of the most important databases available to social scientists working with US data. It allows scientists to look at the determinants of earnings as well as educational attainment and has incredible relevance for government policy. It can also shed light on politically sensitive issues like how different educational attainment and salaries are for people of different ethnicity, sex, and other factors. When we have a better understanding how these variables affect education and earnings we can also formulate more suitable government policies.
In [1]:	<pre>Upgrade Plotly %pip installupgrade plotly Requirement already satisfied: plotly in c:\users\manda\anaconda3\lib\site-packages (5.11.0) Requirement already satisfied: tenacity>=6.2.0 in c:\users\manda\anaconda3\lib\site-packages (from plotly) (8.0.1)</pre>
In [2]:	Note: you may need to restart the kernel to use updated packages. WARNING: Retrying (Retry(total=4, connect=None, read=None, redirect=None, status=None)) after connection broken by 'ProtocolError('Connection abort ed.', ConnectionResetError(10054, 'An existing connection was forcibly closed by the remote host', None, 10054, None))': /simple/plotly/ Import Statements import pandas as pd import numpy as np
	<pre>import seaborn as sns import plotly.express as px import matplotlib.pyplot as plt from sklearn.linear_model import LinearRegression from sklearn.model_selection import train_test_split Notebook Presentation</pre>
In [3]: In [4]:	Load the Data
	Have a look at the file entitled NLSY97_Variable_Names_and_Descriptions.csv . :Key Variables: 1. S Years of schooling (highest grade completed as of 2011) 2. EXP Total out-of-school work experience (years) as of the 2011 interview. 3. EARNINGS Current hourly earnings in \$ reported at the 2011 interview
	Preliminary Data Exploration Challenge What is the shape of df_data? How many rows and columns does it have?
In [5]:	 What are the column names? Are there any NaN values or duplicates? print(f"df_data shape: {df_data.shape}") print(f"df_data has {df_data.shape[0]} rows and {df_data.shape[1]} columns") print(f"df_data column names: {df_data.columns}") df_data shape: (2000, 96) df_data has 2000 rows and 96 columns df_data column names: Index(['ID', 'EARNINGS', 'S', 'EXP', 'FEMALE', 'MALE', 'BYEAR', 'AGE',
	'AGEMBTH', 'HHINC97', 'POVRAT97', 'HHBMBF', 'HHBMOF', 'HHOMBF', 'HHBMONLY', 'HHBFONLY', 'HHOTHER', 'MSA97NC', 'MSA97NCC', 'MSA97NK', 'ETHBLACK', 'ETHHISP', 'ETHWHITE', 'EDUCPROF', 'EDUCPHD', 'EDUCMAST', 'EDUCBA', 'EDUCAA', 'EDUCHSD', 'EDUCGED', 'EDUCDO', 'PRMONM', 'PRMONF', 'PRMSTYUN', 'PRMSTYPE', 'PRMSTYAN', 'PRMSTYAE', 'PRFSTYUN', 'PRFSTYPE', 'PRFSTYAN', 'PRFSTYAE', 'SINGLE', 'MARRIED', 'COHABIT', 'OTHSING', 'FAITHN', 'FAITHP', 'FAITHC', 'FAITHJ', 'FAITHO', 'FAITHM', 'ASVABAR', 'ASVABWK', 'ASVABPC', 'ASVABMK', 'ASVABNO', 'ASVABCS', 'ASVABC', 'ASVABC4', 'VERBAL', 'ASVABMV', 'HEIGHT', 'WEIGHT04', 'WEIGHT11', 'SF', 'SM', 'SFR', 'SMR', 'SIBLINGS', 'REG97NE',
In [6]: Out[6]:	True
In [7]:	
In [8]: In [9]: In [10]:	df_data.fillna(0, inplace=True) df_data.drop_duplicates(inplace=True) Descriptive Statistics df_data.describe()
Out[10]:	count 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 1,487.00 <t< td=""></t<>
	50% 3,474.00 15.75 15.00 6.63 0.00 1.00 1,982.00 29.00 26.00 41,840.00 1.00 0.00 0.00 0.00 0.00 0.00
In [11]:	<pre>plt.figure(figsize=(8, 4), dpi=200) sns.displot(df_data["EARNINGS"], bins=50, aspect=2, kde=True) plt.title(f'1997-2011 Earnings. Average: \${(df_data.EARNINGS.mean()):.3} per hour') plt.xlabel("Earnings Per Hour") plt.ylabel("Number of earners") plt.show() </pre> <pre></pre>
	250 - 200 - 8 150 -
In [12]: Out[12]:	0 0 1 1 2 0 3 0 4 0
	1984
In [13]: In [14]:	We can't use all the entries in our dataset to train our model. Keep 20% of the data for later as a testing dataset (out-of-sample data). y = df_data["EARNINGS"] to_drop = ["EARNINGS"] X = df_data.drop(to_drop, axis=1) X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=10) train_pct = 100 * len(X_train) / len(X) print(f"Training data is : {train_pct:.3}% of the total data")
In [15]:	Training data is : 80.0% of the total data
In [16]: Out[16]:	Only use the years of schooling to predict earnings. Use sklearn to run the regression on the training dataset. How high is the r-squared for the regression on the training data? X_axis=pd.DataFrame(X_train["S"]) X_axis.head() S 179 14
In [17]: In [18]:	
In [18]: In [19]: In [20]:	model = lm.fit(X_axis, y_train) print(f"r-squared the years of schooling to predict earnings: {model.score(X_axis, y_train):.2}") r-squared the years of schooling to predict earnings: 0.081 print("Conclusion: Using Years of Schooling is not enough variable to predict earnings.") Conclusion: Using Years of Schooling is not enough variable to predict earnings. Evaluate the Coefficients of the Model
In [21]:	Evaluate the Coefficients of the Model Here we do a sense check on our regression coefficients. The first thing to look for is if the coefficients have the expected sign (positive or negative). Interpret the regression. How many extra dollars can one expect to earn for an additional year of schooling? print(f"model predict : {model.predict(X_axis)}") theta_1 = model.coef_ print(f"model coef: {theta_1}") print(f"model incercept: {model.intercept_}")
In [22]: Out[22]:	<pre>model predict : [18.09739916 20.53575512 13.22068723 18.09739916 20.53575512 20.53575512] model coef: [1.21917798] model incercept: 1.0289074107474185 coef_regr = pd.DataFrame(data=theta_1, index=X_axis.columns, columns=["Coeffisien"]) coef_regr</pre>
In [23]:	S 1.22
In [24]:	How good our regression is also depends on the residuals - the difference between the model's predictions ($\hat{y}i$) and the true values (yi) inside y_train. Do you see any patterns in the distribution of the residuals? predicted_values = model.predict(X_axis) residuals = (y_train - predicted_values) print(f"predicted values \n {predicted_values}") print(f"residuals \n {residuals}") predicted_values
	[18.09739916 20.53575512 13.22068723 18.09739916 20.53575512 20.53575512] residuals 179 -12.34 142 -5.16 1626 6.78 707 2.12 1288 -6.97 1762 23.92
In [25]:	1660 7.50 527 7.90 1263 19.46 1541 5.90 Name: EARNINGS, Length: 1189, dtype: float64 plt.figure(figsize=(8,4), dpi=200) plt.scatter(x=y_train, y=predicted_values, c="indigo", alpha=0.5) plt.plot(y_train, y_train, color="cyan") plt.title("Actual vs Predicted Earning: \$y _i\$ vs \$\hat y_i\$") plt.xlabel("Actual Earning")
	<pre>plt.xlabel("Predicted Earning") plt.show() plt.figure(figsize=(8,4), dpi=200) plt.scatter(x=predicted_values, y=residuals, c="indigo", alpha=0.5) plt.title("Residual vs Predicted Values") plt.xlabel("Predicted Values") plt.ylabel("Residuals") plt.show()</pre>
	Actual vs Predicted Earning: y_i vs $\hat{y_i}$
	100 - 80 - 80 - 40 - 40 - 40 - 40 - 40 -
	9 60 - 20 -
	0 20 40 60 80 100 120 Actual Earning
	Residual vs Predicted Values 100 - 80 -
	Residuals 40 -
	20 - 0 - -20 -
	7.5 10.0 12.5 15.0 17.5 20.0 22.5 25.0 Predicted Values Multivariable Regression
In [27]:	Now use both years of schooling and the years work experience to predict earnings. How high is the r-squared for the regression on the training data? new_X = df_data[["S", "EXP"]]
	regr = LinearRegression() regr.fit(X_train, EXS_y_train) rsquared = regr.score(X_train, EXS_y_train) EXS_predictions = regr.predict(X_train) residuals = (EXS_y_train - EXS_predictions)
In [29]:	new_coef
Out[29]: In [30]: Out[30]:	<pre>df_coef = pd.DataFrame(new_coef, index=X_train.columns, columns=['coef']) df_coef</pre>
In [31]: In [32]:	The coefficients interpret: having additional year of school and more work experience results in higher earning, because S and EXP has a coefficien t greater than zero. Analyse the Estimated Values & Regression Residuals EXP_predicted_values = regr.predict(new_X)
In [33]:	<pre>EXP_residuals = (new_Y - EXP_predicted_values)</pre>
	plt.figure(figsize=(8,4), dpi=200) plt.scatter(x=EXP_predicted_values, y=EXP_residuals, c="indigo", alpha=0.5) plt.title("Residual vs Predicted Values") plt.xlabel("Predicted Values") plt.ylabel("Residuals") plt.show() Actual vs Predicted Earning: y _i vs ŷ _i
	120 -
	Predicted Earning 60 - 08 40 - 40 - 09 1
	20 - 0 -
	0 20 40 60 80 100 120 Actual Earning Residual vs Predicted Values
	100 - 80 - 60 -
	8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8 60 - 8
	0 - -20 -
	5 10 15 20 25 Predicted Values Use Your Model to Make a Prediction How much can someone with a bachelors degree (12 + 4) years of schooling and 5 years work experience expect to earn in 2011?
In [34]:	<pre>experience = 5 # Starting Point: Average Values in the Dataset features = df_data[['S', 'EXP']] average_vals = features.mean().values someone = pd.DataFrame(data=average_vals.reshape(1, len(features.columns)),</pre>
Out[34]: In [35]:	someone
Out[35]: In [36]:	<pre>someone['EXP'] = experience someone</pre>
In [37]:	the earning estimate is: \${estimate:.3} per hour.") Experiment and Investigate Further Which other features could you consider adding to further improve the regression to better predict earnings?
<pre>In [37]: Out[37]:</pre>	
	4 1994 36.06 15 2.94 0 1 1984 27 23.00 44,188.00 1 0 0 0 1 0 0 1 0 0 1 0 0 1 0 0 1 0 1 0 1 0 0 0 0 1 0 0 1 0 1984 2400 9.00 12 10.83 1 0 1982 29 29.00 73,100.00 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 0 0 0 0 1 0 0 0 0 0 1 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 1 0 0 <
In [38]: Out[38]:	1487 rows × 96 columns new_features = df_data[['EXP', 'S', 'MALE']] new_y = df_data['EARNINGS'] new_features
	0 9.71 12 1 1 5.71 17 1 2 9.94 14 1 3 1.54 18 1 4 2.94 15 1 1984 10.83 12 0
	1989 9.37 12 0 1991 6.29 17 0 1994 9.12 12 0 1995 7.87 8 0 1487 rows × 3 columns
In [39]: In [40]:	<pre>new_y, test_size=0.2, random_state=10)</pre> Lregr = LinearRegression() Lregr.fit(LX_train, L_y_train) Lrsquared = Lregr.score(LX_train, L_y_train)
In [41]: In [42]: Out[42]:	residuals = (L_y_train - L_predictions) print(f'Training data r-squared: {Lrsquared:.2}') Training data r-squared: 0.15
Out[42]: In [43]:	
Out[43]:	<pre>coef MALE 4.23 S 1.93 EXP 1.00</pre>