**Group C Member:** Zhenning Li (B00780786) Yixin Zhao (B00780592) Miao Wang (B00761718) Jiaqian Ma (B00782247) In [33]: import pandas as pd import numpy as np import requests import io import matplotlib.pyplot as plt import statsmodels.api as sm import math from statsmodels.stats.outliers influence import summary table from scipy.stats import norm from scipy.stats import t from scipy import stats I. Summary statistics In [3]: url = 'https://raw.githubusercontent.com/adufays/GDP expectancy/main/APM Returns.csv' s = requests.get(url).content df = pd.read\_csv(io.StringIO(s.decode('utf-8'))) df = df.drop('Unnamed: 0', axis = 1) df.iloc[:,:10] Date Mkt-RF SMB HML RMW CMA...6 MMM ABT ATVI ADBE Out[3]: **0** 2000/01 -4.74 4.15 -0.29 -6.05 4.73 -1.20 -7.37 -5.35 -16.43 1 2000/02 2.45 18.32 -9.93 -18.33 -0.51 -6.24 0.91 -19.23 84.81 2 2000/03 5.20 -14.91 7.38 11.68 -1.05 -0.04 6.16 -5.40 8.66 3 2000/04 -6.40 -5.55 7.55 -2.65 8.78 -48.65 8.19 8.61 -3.68 2.56 4 2000/05 -4.42 4.63 0.74 -1.51 5.35 -1.50 -7.43 **205** 2017/02 3.57 -2.12 -1.79 0.78 -1.72 6.56 7.88 12.20 4.33 206 2017/03 0.78 -3.17 0.68 -1.00 2.64 -1.52 10.45 9.93 207 2017/04 1.09 0.49 -1.87 2.00 -1.55 2.30 -1.78 4.74 2.72 **208** 2017/05 1.06 -3.05 -3.78 -1.88 4.35 4.57 12.05 6.01 **209** 2017/06 0.78 2.48 1.35 -2.01 -0.07 1.76 6.40 -1.78 -0.36 210 rows × 10 columns 1.1 Computing monthly average return and standard deviation In [4]: df 1 = df.iloc[:, 6:]In [5]: asset name list = list(df 1.columns) # we used describe() function to generate dataframes # of average return and standard deviation asset average return list = list(df 1.describe().iloc[1]) asset standard deviation list = list(df 1.describe().iloc[2]) df\_1.describe().iloc[:, :7] **ADBE** Out[5]: **MMM ABT ATVI AES AET AMG** count 210.000000 210.000000 210.000000 210.000000 210.000000 210.000000 210.000000 0.572381 0.734286 2.379476 0.488429 1.569095 1.784286 1.396762 mean 14.581512 9.353504 5.661786 5.714191 12.059177 11.942540 10.969253 std -15.500000 -20.870000 -33.860000 -44.100000 -48.650000 -62.330000 -31.210000 min 25% -2.592500 -2.625000 -3.460000 -5.452500 -4.190000 -4.582500 -3.870000 0.950000 **50%** 0.955000 2.265000 2.710000 -0.665000 2.495000 1.930000 3.885000 7.800000 6.417500 **75**% 4.520000 7.600000 6.905000 6.677500 20.150000 74.290000 34.480000 13.820000 59.500000 84.810000 49.710000 max In [6]: plt.rcParams['figure.figsize'] = (10.0, 7.0) plt.scatter(asset\_average\_return\_list, asset\_standard\_deviation\_list, c = 'black', s alpha = 0.7, label = 'samples') plt.grid() plt.legend(fontsize = 14)plt.xlabel('Average return') plt.ylabel('Standard deviation') plt.show() samples 25 Standard deviation 15 10 5 Average return 1.2 Use OLS to estimate the model The model is:  $std_i = \beta_1 + \beta_2 \times ave_i + \epsilon_i$ In [7]: ave = np.array(asset average return list) std = np.array(asset standard deviation list) x model = sm.add constant(ave)  $y \mod el = std$ model = sm.OLS(y model, x model)result = model.fit() result.params Out[7]: array([6.82753557, 2.65556998]) Here,  $\boldsymbol{\beta}_1 = 6.82753557$ , and  $\boldsymbol{\beta}_2 = 2.65556998$  $std_i = 6.83 + 2.66 \times ave_i$ 1.3 Interpret the coefficient  $\beta_2$ The literal interpretation of the estimated coefficient  $\beta_2$  is that a one-unit increase in  $ave_i$  (asset average return) will produce an expected increase in  $std_i$  (asset standard deviation) of  $\beta_2$  = 2.66 units. 1.4 Scatter plot with regression line In [8]: x\_reg= np.asarray(asset\_average\_return\_list) y reg = result.params[0] + result.params[1]\*x reg plt.scatter(asset\_average\_return\_list, asset\_standard\_deviation\_list, c = 'black', s = 5, alpha = 0.7, label = 'sample') plt.xlabel('Average return') plt.ylabel('Standard deviation') plt.plot(x\_reg, y\_reg, label = 'regression') plt.grid() plt.legend(fontsize = 14) plt.rcParams['figure.figsize'] = (10.0, 7.0) plt.show() regression sample 25 Standard deviation 15 10 5 Average return 1.5 Write explicit mathematical formula for the estimators  $\hat{\beta}_2 = \frac{\sum_{t=1}^{T} (x_t - \bar{x})(y_t - \bar{y})}{\sum_{t=1}^{T} (x_t - \bar{x})^2} = \frac{cov(ave, std)}{var(ave)}$  $\hat{\boldsymbol{\beta}}_1 = average(std) - \hat{\boldsymbol{\beta}}_2 * average(ave)$ 1.6 Is Beta2 statistically significant? Perform a hypothesis testing In [9]: result.summary() **OLS Regression Results** Out[9]: OLS Model: Adj. R-squared: Method: **Least Squares** F-statistic: 107.4 Fri, 28 May 2021 **Prob (F-statistic):** 2.47e-22 Log-Likelihood: Time: 18:07:44 -975.32 No. Observations: 384 AIC: 1955. **Df Residuals:** BIC: 382 1963. **Df Model: Covariance Type:** nonrobust t P>|t| [0.025 0.975] 0.293 23.303 0.000 **const** 6.8275 6.251 7.404 **x1** 2.6556 0.256 10.361 0.000 2.152 3.160 **Omnibus:** 108.212 **Durbin-Watson:** 1.901 Prob(Omnibus): 0.000 Jarque-Bera (JB): **Prob(JB):** 2.97e-56 Skew: 1.405 **Kurtosis:** 5.844 Cond. No. 3.48 Notes: [1] Standard Errors assume that the covariance matrix of the errors is correctly specified. From the table above, we can see for  $\beta 2$ , the t-value is 10.361, and P>|t| is nearly 0, thus we can say  $\beta 2$ is statistically significant. We can also calculate t-value by using the formula: In [10]: y = np.array(asset standard deviation list) y\_reg = result.params[0] + result.params[1]\*np.array(asset average return list) residual square = (np.square(y-y reg)).sum()/(len(y)-2)x\_variance = np.var(asset\_average\_return\_list) x\_error = x\_variance\*(len(np.array(asset\_average\_return\_list))-1) t = 2.656/(math.sqrt(residual\_square/x\_error)) Out[10]: 10.349174525454748  $H_0$ :  $\beta_2 = 0$  $H_1$ :  $\beta_2 \neq 0$ Test statistic:  $t_{eta_2}=rac{2.656-0}{SE}=rac{2.6556-0}{0.2563}=10.36$  $P-value = 2[1-p(t_{382} \le t_{eta_2})] = 0 < 5\%$ We should reject the null hypothesis  $\beta_2 = 0$ .  $eta_2=2.656$  is thus statistically significant. II. Linear Regression In [11]: url = 'https://raw.githubusercontent.com/adufays/GDP expectancy/main/life gdp.csv' s = requests.get(url).content df\_2 = pd.read\_csv(io.StringIO(s.decode('utf-8'))) 2.1 Scatter plot x-y model In [12]: plt.scatter(df\_2['GDPpcap'], df\_2['Life exp.'], c = 'black', s = 5, alpha = 0.7, label = 'samples') slope, intercept, r\_value, p\_value, std\_err = stats.linregress(df\_2['GDPpcap'], df 2['Life exp.']) x reg 1 = np.linspace(0, 120, 10)y reg 1 = intercept + slope\*x reg 1 plt.plot(x\_reg\_1, y\_reg\_1, label = 'regression') plt.grid() plt.legend(fontsize = 14) plt.xlabel('GDP per cap') plt.ylabel('Life expectancy') plt.rcParams['figure.figsize'] = (10.0, 7.0) regression samples 90 Life expectancy 70 60 120 GDP per cap 2.2 Scatter plot logx-y model In [13]: life exp list = np.array(df 2['Life exp.']) GDP list = np.array(df 2['GDPpcap']) GDP list log = np.log(GDP list) In [14]: plt.scatter(GDP\_list\_log, life\_exp\_list, c = 'black', s = 5, alpha = 0.7, label = 'samples') slope, intercept, r value, p value, std err = stats.linregress(GDP list log, life\_exp\_list)  $x_reg_2 = np.linspace(-1, 5, 10)$ y\_reg\_2 = intercept + slope\*x\_reg\_2 plt.plot(x\_reg\_2, y\_reg\_2, label = 'regression') plt.grid() plt.legend(fontsize = 14) plt.rcParams['figure.figsize'] = (10.0, 7.0) plt.xlabel('log GDP per cap') plt.ylabel('Life expectancy plt.show() regression 85 samples 80 75 Life expectancy 70 55 log GDP per cap 2.3 Coefficient of determination x-y model In [15]: x = np.array(df 2['GDPpcap']) y = np.array(df\_2['Life exp.'])  $x_{model} = sm.add_{constant}(x)$  $y \mod el = y$  $model_2 = sm.OLS(y_model, x_model)$ result 2 = model 2.fit() In [16]: result 2.params array([67.51367284, 0.24932456]) Out[16]: In [17]: result 2.summary() **OLS Regression Results** Out[17]: Dep. Variable: R-squared: 0.442 У Model: OLS Adj. R-squared: 0.439 Method: **Least Squares** F-statistic: 142.4 Fri, 28 May 2021 Prob (F-statistic): 1.47e-24 Date: Time: 18:07:45 Log-Likelihood: -574.97 No. Observations: 182 AIC: 1154. **Df Residuals:** 180 BIC: 1160. **Df Model: Covariance Type:** nonrobust t P>|t| [0.025 0.975] coef std err 0.585 115.478 0.000 66.360 **const** 67.5137 0.2493 0.021 11.933 0.000 0.208 0.291 х1 **Omnibus:** 22.094 2.006 **Durbin-Watson:** Prob(Omnibus): 0.000 Jarque-Bera (JB): 26.566 Skew: -0.925 **Prob(JB):** 1.70e-06 **Kurtosis:** 3.285 Cond. No. 38.5 Notes: [1] Standard Errors assume that the covariance matrix of the errors is correctly specified. The adjusted coefficient of determination for x-y model  $R_x^2$  is 0.439Inx-y model In [18]: x\_model = sm.add\_constant(GDP\_list\_log) y\_model = life\_exp\_list  $model_3 = sm.OLS(y_model, x_model)$ result 3 = model 3.fit() In [19]: result 3.params Out[19]: array([59.65932948, 5.33364812]) In [20]: result 3.summary() **OLS Regression Results** Out[20]: R-squared: 0.680 Dep. Variable: У OLS Model: Adj. R-squared: 0.678 Least Squares F-statistic: 382.8 Method: Fri, 28 May 2021 **Prob (F-statistic):** 2.00e-46 Date: Time: 18:07:45 Log-Likelihood: -524.28 No. Observations: 182 AIC: 1053. **Df Residuals:** 180 BIC: 1059. **Df Model: Covariance Type:** nonrobust [0.025 0.975] coef std err t P>|t| **const** 59.6593 0.722 82.628 0.000 61.084 58.235 5.3336 0.273 19.565 0.000 4.796 5.872 **Omnibus:** 40.029 **Durbin-Watson:** 2.238 Prob(Omnibus): Jarque-Bera (JB): 0.000 68.168 **Prob(JB):** 1.58e-15 Skew: -1.116 Cond. No. 5.002 6.65 **Kurtosis:** Notes: [1] Standard Errors assume that the covariance matrix of the errors is correctly specified. The adjusted coefficient of determination for Inx-y model  $R^2_{lnx}$  is 0.678Conclusion Here we have  $R_x^2=0.439$  and  $R_{lnx}^2=0.678$ . With x-y model, we improve the fit by 43.9% compared to the sample average, and we improve the fit by 67.8% with the lnx-y model. It's clear that the Inx-y model performs better. Since GDP has a long-tailed distribution (right skewed), the logarithmic transformation on GDP(x variable) reduces or removes the skewness of our original data, creating the desired asymptotic association between input (X variable )and target (Y variable) on the original input scale. Besides, we can infer from the scatter plot that the lnx-y model has smaller error. 2.4 Empirical correlation between the log of GDP and life expectancy The correlation coefficient is: In [21]: fanx = np.array(df 2['GDPpcap']) x 1 = np.log(x)y = np.array(df 2['Life exp.']) cor = np.vstack((x 1, y))print('The empirical correlation is:', (np.corrcoef(cor)[0][1])) print('The squared empirical correlation is:', (np.corrcoef(cor)[0][1])\*\*2) print('The coefficient of determination for lnx-y model is: ', result 3.rsquared) The empirical correlation is: 0.8247155249526711 The squared empirical correlation is: 0.6801556970979599 The coefficient of determination for lnx-y model is: 0.6801556970979601 The coefficient of determination is equal to the square of the empirical correlation. III. Generating spurious correlation In [22]: url = 'https://raw.githubusercontent.com/adufays/GDP\_expectancy/main/Appl\_Russell.csv s = requests.get(url).content df\_3 = pd.read\_csv(io.StringIO(s.decode('utf-8')))  $df_3 = df_3.drop('Unnamed: 0', axis = 1)$ df 3 Out[22]: Date Apple returns Russell returns **0** 2016-04-25 -11.989084 -1.260337 **1** 2016-05-02 -1.094085 -0.561795 2 2016-05-09 -0.581591 -1.794334 **3** 2016-05-16 5.061906 0.377336 **4** 2016-05-23 5.247411 2.334956 **257** 2021-03-29 1.465979 1.245510 **258** 2021-04-05 7.816477 2.372223 **259** 2021-04-12 0.868402 1.428772 **260** 2021-04-19 -0.641943 0.505571 **261** 2021-04-19 0.074139 -0.000118 262 rows × 3 columns In [23]: df 4 = df 3.set index('Date') df\_4.plot() <AxesSubplot:xlabel='Date'</pre> Out[23]: 15 Apple returns Russell returns 10 5 0 -10-15-20 2016-04-25 2017-04-10 2018-03-26 2019-03-11 2020-02-24 2021-02-08 Date 3.1 Generate explanatory variables from Normal Distribution In [24]: df explanVar = pd.DataFrame() for i in range(1000): df explanVar[i] = norm.rvs(size=262, loc=0, scale=1) df explanVar.iloc[:, :7] 1 Out[24]: 0 2 5 6 -1.103979 0 1.297910 0.804206 -0.068985 -0.415247 -0.368813 -1.151460 1 0.949897 -0.141375 -0.871852 1.291526 0.419458 -0.498054 -0.174098 2 1.216040 -0.263023 -0.141120 -1.244212 -1.805208 1.005682 0.454859 -0.749674 0.270932 0.692110 -0.069647 1.760037 -0.773310 1.254317 0.274575 4 1.032070 0.277932 1.272772 -0.416011 1.425553 -1.010408 -0.159323 257 -0.782331 -1.363049 0.971107 0.676408 -0.945307 -0.601152 258 -0.557416 0.606416 1.174797 0.120092 -0.875927 0.722186 -0.542585 0.029605 259 -1.240483 0.434995 0.636910 1.106097 -1.190532 -0.411126 -0.786905 -0.505423 0.352118 -1.304097 260 2.497227 -1.754366 0.417685 0.809549 -1.084726 261 0.639011 1.469021 -0.934225 1.164064 -0.631074 262 rows × 7 columns Now we have 1000 explanatory variables of size 262 3.2 Test statistics In [25]: y\_model\_1 = df\_3['Apple returns'] test statistics = [] x for plot = []beta1 = [] beta2 = []p value lst = [] std err lst = [] In [26]: **for** i **in** range (1000): x model = df explanVar.iloc[:,i] slope, intercept, r\_value, p\_value, std\_err = stats.linregress(x\_model, x for plot.append(i) test statistics.append(slope/std err) beta1.append(intercept) beta2.append(slope) p value lst.append(p value) std err lst.append(std err) In [27]: df lm = pd.DataFrame(list(zip(beta1, beta2, test\_statistics, p\_value\_lst, std\_err\_lst)), columns=["intercept", "slope", "test-statistics", "p-value", "std err"])  $df_lm$ Out[27]: p-value slope test-statistics -0.089326 0.928892 0.242300 **0** 0.647768 -0.021644 **1** 0.646012 -0.192805 -0.740348 0.459757 0.260425 2 0.650820 -0.230201 -0.897561 0.370250 0.256474 -0.017935 0.985705 0.246016 0.650115 -0.004412 0.837105 0.403302 0.243056 0.661927 0.203463 -0.050252 0.959961 0.228480 **995** 0.651106 -0.011481 0.630077 -0.238862 -1.019512 0.308907 0.234291 **997** 0.638431 -0.188627 -0.732820 0.464329 0.257399 -0.272627 0.785356 0.251017 **998** 0.650794 -0.068434 -1.298278 0.195342 0.238115 **999** 0.670594 -0.309140 1000 rows × 5 columns In [28]: plt.rcParams['figure.figsize'] = (10.0, 7.0) plt.scatter(x\_for\_plot, test\_statistics, c = 'black', s = 5, alpha = 0.7, label = 'test statistics values') plt.grid() plt.legend(fontsize = 14) one\_line = np.linspace(0, 1000, 1000)  $y_{value} = np.array([1.96]*1000)$ plt.plot(one\_line, y\_value, label = 'upper bound')  $y_value_2 = np.array([-1.96]*1000)$ plt.plot(one\_line, y\_value\_2, label = 'lower bound') plt.legend(loc= 'upper right') plt.xlabel('1000 random normal distribution values') plt.ylabel('test statistics values') plt.show() upper bound lower bound test statistics values 1 test statistics values  $^{-1}$ -2 -3 800 1000 1000 random normal distribution values Here, the test statistic formula is:  $t = rac{\hat{eta}_2 - 0}{SE(\hat{eta}_2)}$ where  $SE(\hat{eta}_2) = \sqrt{rac{\hat{\sigma}^2}{\sum_{t=1}^T (x_t - ar{x})^2}}$ 3.3 Hypothesis test In [29]: maximum t = np.absolute(np.array(test statistics)).max() print("So, the maximum value of test statistics is:", maximum t) So, the maximum value of test statistics is: 3.378476711517529 Since  $P - value = 2(1 - P[t_{T-2} \le |t|])$ , we have: In [34]: p value 1 = 2\*(1- t.cdf(maximum t,260))print('P-value is: ', p\_value\_1) P-value is: 0.0008406438969521535 With the maximun value of test statistics, we got a p-value far smaller than 0.05, which indicates that we can confidently reject the Null Hypothesis that  $\beta_2=0$  at a significant level of 95%. This experiment shows that, even if we generate some random numbers, it is still possible that the result can reject the Null Hypothesis, so we cannot blindly trust the result we got and write down the conclusion without thinking.