In [2]: import numpy as np import pandas as pd import matplotlib.pyplot as plt import seaborn as sns import warnings warnings.filterwarnings("ignore") In [3]: df = pd.read_csv("ds_salaries.csv") In [4]: **df** $salary_currency_salary_in_usd_employee_residence_remote_ratio_company_location_company_size$ Out[4]: work_year experience_level employment_type job_title ES L 2023 SE FT Principal Data Scientist 80000 EUR 85847 100 ES 2023 MI СТ 30000 USD 30000 US 100 US S ML Engineer 1 ML Engineer 25500 25500 US 100 US 2 2023 MI CT USD S 175000 3 2023 SE FT Data Scientist 175000 USD CA 100 CA M 4 2023 SE FT Data Scientist 120000 USD 120000 CA 100 $\mathsf{C}\mathsf{A}$ M 412000 412000 US US 3750 2020 SE FT Data Scientist USD 100 L 3751 2021 MI 151000 USD 151000 US 100 US L FT Principal Data Scientist ΕN USD 105000 US US S 3752 2020 FT Data Scientist 105000 100 US 3753 2020 ΕN CT Business Data Analyst 100000 USD 100000 US 100 L 50 3754 2021 SE FT Data Science Manager 7000000 INR 94665 IN IN L 3755 rows × 11 columns In [5]: df.head(10) work_year experience_level employment_type job_title salary_salary_currency salary_in_usd employee_residence remote_ratio company_location company_size Out[5]: 0 2023 SE FT Principal Data Scientist 80000 **EUR** 85847 ES 100 ES L US S 1 2023 MI СТ USD 100 US ML Engineer 30000 30000 2 2023 MI CT ML Engineer 25500 USD 25500 US 100 US S 3 FT USD 175000 CA Μ 2023 SE Data Scientist 175000 CA 100 4 2023 SE FT Data Scientist 120000 USD 120000 CA 100 CA M 5 2023 FT Applied Scientist 222200 USD 222200 US 0 US L SE Applied Scientist 136000 6 2023 SE FT USD 136000 US 0 US L Data Scientist 219000 7 SE USD 219000 0 CA Μ 2023 FT CA 8 2023 SE FT Data Scientist 141000 USD 141000 CA 0 CA M SE FT Data Scientist 147100 USD US US 9 2023 147100 0 Μ #Exploratory Data Analysis (EDA): In [6]: df.shape (3755, 11)Out[6]: In [7]: df.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 3755 entries, 0 to 3754 Data columns (total 11 columns): Column Non-Null Count Dtype 0 work_year 3755 non-null int64 experience_level 3755 non-null object 1 employment_type 3755 non-null object 3755 non-null 3 job_title object 3755 non-null 4 salary int64 salary_currency 3755 non-null object salary_in_usd 3755 non-null int64 6 employee_residence 3755 non-null 7 object 8 remote_ratio 3755 non-null int64 company_location 9 3755 non-null object company_size 3755 non-null object 10 dtypes: int64(4), object(7) memory usage: 322.8+ KB In [8]: df.columns Index(['work_year', 'experience_level', 'employment_type', 'job_title', Out[8]: 'salary', 'salary_currency', 'salary_in_usd', 'employee_residence', 'remote_ratio', 'company_location', 'company_size'], dtype='object') df.describe() Out[9]: work_year salary_in_usd remote_ratio salary **count** 3755.000000 3.755000e+03 3755.000000 3755.000000 mean 2022.373635 1.906956e+05 137570.389880 46.271638 63055.625278 48.589050 std 0.691448 6.716765e+05 min 2020.000000 6.000000e+03 5132.000000 0.000000 **25**% 2022.000000 1.000000e+05 95000.000000 0.000000 0.000000 **50%** 2022.000000 1.380000e+05 135000.000000 **75**% 2023.000000 1.800000e+05 175000.000000 100.000000 max 2023.000000 3.040000e+07 450000.000000 100.000000 In [10]: df.dtypes work_year int64 Out[10]: object experience_level employment_type object job_title object salary int64 salary_currency object salary_in_usd int64 employee_residence object remote_ratio int64 company_location object object company_size dtype: object In [11]: df.nunique() work_year 4 Out[11]: experience_level 4 employment_type 4 job_title 93 815 salary salary_currency 20 1035 salary_in_usd employee_residence 78 remote_ratio 3 company_location 72 company_size 3 dtype: int64 There are 4 categorical values in the column "experience_level", such as: EN, which is Entry-level. MI, which is Mid-level. SE, which is Senior-level. EX, which is Executive-level. There are 3 categorical values in the column "remote_ratio", such as: 100, which is Remotely. 0, which is On-site. 50, which is Hybrid plt.figure(figsize=(8, 6)) sns.countplot(data=df, x='experience_level') plt.xlabel('Experience Level') plt.ylabel('Count') plt.title('Distribution of Work Experience Levels') plt.show() Distribution of Work Experience Levels 2500 2000 1500 Count 1000 500 0 SE ΜI EΝ EX Experience Level df1 = df.groupby("work_year")["salary"].mean() df1.plot(kind="bar") plt.title("Average Salary by Work Year") plt.xlabel("Work Year") plt.ylabel("Salary") plt.show() Average Salary by Work Year 500000 400000 300000 200000 100000 0 Work Year In [18]: df.hist() plt.show() work_year salary 1500 3000 1000 2000 500 1000 0 -2022 2023 2020 2021 0 1 3 salary in usd remote ratio 1e7 2000 1000 1500 750 1000 500 500 250 -0 0 1000020000B0000400000 0 25 50 75 100 In [19]: print("Total value counts of the roles:\n", df["experience_level"].value_counts()) Total value counts of the roles: experience_level SE 2516 ΜI 805 ΕN 320 EX 114 Name: count, dtype: int64 In [20]: roles = ["SE", "MI", "EN", "EX"] people = [2516, 805, 320, 114] plt.pie(people, labels=roles, autopct='%1.1f%%') plt.title('Distribution of Roles') plt.axis('equal') plt.show() Distribution of Roles SE 67.0% 3.0% ΕX 8.5% ΕN 21.4% MΙ In [21]: print(df["employment_type"].value_counts()) employment_type FT 3718 PΤ 17 CT 10 FL 10 Name: count, dtype: int64 In [22]: company_numbers = [3153, 454, 148] company_size = ["M", "L", "S"] explode = [0.1, 0, 0]plt.pie(company_numbers, labels=company_size, explode=explode, autopct='%1.1f%%') plt.axis('equal') plt.title("Company Size") plt.show() Company Size М 84.0% 3.9% S 12.1% In [23]: df["job_title"].value_counts() job_title Out[23]: Data Engineer 1040 Data Scientist 840 Data Analyst 612 Machine Learning Engineer 289 Analytics Engineer 103 Principal Machine Learning Engineer 1 Azure Data Engineer 1 Manager Data Management 1 Marketing Data Engineer 1 Finance Data Analyst Name: count, Length: 93, dtype: int64 In [24]: sns.catplot(x="experience_level", y="salary_in_usd", kind="bar", data=df) plt.show() 200000 150000 salary_in_usd 125000 100000 75000 50000 25000 0 SE ΕN EX experience_level There are 4 categorical values in the column "employment_type", such as: FT, which is Full-time. PT, which is Pant-time. CT, which is Contractual. FL, which is Fneelancer. In [26]: sns.catplot(x="employment_type", y="salary_in_usd", kind="bar", data=df) plt.show() 200000 175000 150000 125000 psn_in_100000 75000 50000 25000 0 ст ΡT FΤ FL employment_type In [27]: sns.catplot(x="work_year", y="salary_in_usd", kind="bar", data=df) plt.show() 140000 120000 100000 salary_in_usd 80000 60000 40000 20000 0 2020 2021 2022 2023 work_year In [28]: sns.catplot(x="company_size", y="salary_in_usd", kind="box", data=df) plt.title("Box plot grouped by company size") plt.show() Box plot grouped by company size 400000 300000 salary_in_usd 200000 100000 0 S Μ company_size In [29]: sns.catplot(x="experience_level", y="salary_in_usd", hue="company_size", kind="box", data=df) plt.title("Box plot grouped by company size") plt.show() Box plot grouped by company size 400000 300000 salary_in_usd company_size 200000 100000 0 EX SE ΕN experience_level In [30]: cat_list = [i for i in df.select_dtypes("object")] In [31]: cat_list ['experience_level', Out[31]: 'employment_type', 'job_title', 'salary_currency', 'employee_residence', 'company_location', 'company_size'] In [33]: for i in cat_list: df[i] = df[i].factorize()[0] In [34]: **df** Out[34]: $work_year \quad experience_level \quad employment_type \quad job_title$ salary salary_currency salary_in_usd employee_residence remote_ratio company_location company_size 0 2023 0 0 0 80000 0 85847 0 100 0 0 30000 30000 1 2023 1 1 100 1 1 1 1 2 2023 1 1 1 25500 1 25500 1 100 1 1 0 2 3 2023 175000 175000 2 100 2 0 2 120000 2 2 4 2023 0 1 120000 100 2 ••• 3750 2020 0 0 2 412000 1 412000 1 100 1 0 2021 0 0 151000 151000 0 3751 1 100 3752 2020 2 0 2 105000 1 105000 1 100 1 1 0 3753 2020 20 100000 100000 100 1 2 6 6 3754 2021 26 7000000 94665 50 0 3755 rows × 11 columns Correlat ton In []: plt.figure(figsize=(12, 12)) sns.heatmap(df.corr(), annot=True, linewidths=0.7, cmap="viridis", fmt=".2f") plt.show() 1.0 work_year -1.00 -0.17 -0.12 -0.16 -0.09 -0.13 0.23 -0.28 -0.24 -0.26 0.39 -0.17 1.00 0.13 0.13 0.03 0.14 -0.24 0.21 0.04 0.19 -0.21 experience_level -- 0.8 employment_type --0.12 0.13 1.00 0.06 -0.01 0.06 -0.13 0.23 0.06 0.13 -0.09 - 0.6 job_title --0.16 0.13 0.06 1.00 0.04 0.08 -0.07 0.19 0.06 0.18 -0.18 salary --0.09 0.03 -0.01 0.04 1.00 0.32 -0.02 0.16 0.03 0.19 -0.14 - 0.4 salary_currency --0.13 0.14 0.06 0.08 0.32 1.00 -0.20 0.25 0.03 0.28 -0.14 0.23 -0.24 -0.13 -0.07 -0.02 -0.20 1.00 -0.30 -0.06 -0.30 0.17 salary_in_usd -0.2 -0.28 0.21 0.23 0.19 0.16 0.25 -0.30 1.00 0.10 0.82 -0.24 employee_residence -- 0.0 -0.24 0.04 0.06 0.06 0.03 0.03 -0.06 0.10 1.00 0.07 -0.14 remote_ratio --0.26 0.19 0.13 0.18 0.19 0.28 -0.30 0.82 0.07 1.00 -0.23 company_location -- -0.2 0.39 -0.21 -0.09 -0.18 -0.14 -0.14 0.17 -0.24 -0.14 -0.23 1.00 company_size job_title salary experience_level employment_type salary_currency salary_in_usd employee_residence remote_ratio company_location In [36]: X = df.drop(["salary_in_usd"], axis=1) Y = df["salary_in_usd"] In []: Splitting the Dataset into Train & Test: In [37]: from sklearn.model_selection import train_test_split In [38]: X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=42) In [39]: print("X_train shape:", X_train.shape) print("X_test shape:", X_test.shape) print("Y_train shape:", Y_train.shape) print("Y_test shape:", Y_test.shape) X_train shape: (3004, 10) X_test shape: (751, 10) Y_train shape: (3004,) Y_test shape: (751,) In []: Modeling In [40]: **from** sklearn.linear_model **import** Ridge, Lasso, RidgeCV, LassoCV, ElasticNet, ElasticNetCV, LinearRegression from sklearn.tree import DecisionTreeRegressor from sklearn.neighbors import KNeighborsRegressor from sklearn.neural_network import MLPRegressor from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor, AdaBoostRegressor from sklearn import neighbors from sklearn.svm import SVR In [41]: dt=DecisionTreeRegressor() In [42]: dt.fit(X_train,Y_train) Out[42]: ▼ DecisionTreeRegressor DecisionTreeRegressor() In [43]: y_predict = dt.predict(X_test) In [44]: dt.score(X_train,Y_train) Out[44]: 1.0 In [45]: dt.score(X_test,Y_test) 0.9939366761794373 Out[45]: In [46]: **from** sklearn.metrics **import** r2_score, mean_squared_error, mean_absolute_error print(r2_score(Y_test, y_predict) * 100) print(mean_squared_error(Y_test, y_predict)) print(mean_absolute_error(Y_test, y_predict)) 99.39366761794372 23936728.14780293 801.9800266311585 In [48]: knn = KNeighborsRegressor().fit(X_train, Y_train) ada = AdaBoostRegressor().fit(X_train, Y_train) svm = SVR().fit(X_train, Y_train) ridge = Ridge().fit(X_train, Y_train) lasso = Lasso().fit(X_train, Y_train) rf = RandomForestRegressor().fit(X_train, Y_train) gbm = GradientBoostingRegressor().fit(X_train, Y_train) In [49]: models=[ridge, lasso, knn, ada, svm, rf, gbm] In [50]: **def** ML(Y, model): y_pred = model.predict(X_test) mse = mean_squared_error(Y_test, y_pred) rmse = np.sqrt(mse) $r2 = r2_score(Y_test, y_pred) * 100$ return mse, rmse, r2 In [51]: for model in models: print("\n", model, "\n\nDifferent models success rate:", ML("salary_in_usd", model)) Ridge() Different models success rate: (3464466548.979133, 58859.71923972398, 12.24288289418518) Lasso() Different models success rate: (3464396847.135084, 58859.12713534821, 12.244648485742848) KNeighborsRegressor() Different models success rate: (390614700.0352064, 19763.97480354613, 90.10548392093921) AdaBoostRegressor() Different models success rate: (227090924.79969683, 15069.536316678654, 94.24764401688705) SVR() Different models success rate: (3944880838.892905, 62808.28638717112, 0.07368671254132098) RandomForestRegressor() Different models success rate: (21217502.093839817, 4606.245987117906, 99.46254732449627) GradientBoostingRegressor() Different models success rate: (16574508.687009491, 4071.180257248442, 99.58015726829696)