	CS 677 Final Project  Prediction of Data Science Job Salaries  Olivia Lee (U76797374)  Boston University Computer Science Department  Fall 2022
In [42]:	import pandas as pd import numpy as np import matplotlib.pyplot as plt from sklearn.model_selection import train_test_split from sklearn.model_selection matrix_accuracy_score_roc_auc_score_ from sklearn.motrics_import_confusion_matrix_accuracy_score_roc_auc_score_ from sklearn.motrics_import_confusion_matrix_accuracy_score_roc_auc_score_roc_auc_score_roc_auc_score_roc_auc_score_roc_auc_score_roc_auc_score_roc_auc_score_roc_auc_score_roc_auc_score_roc_auc_score_roc_au
	from sklearn.metrics import confusion_matrix, accuracy_score, roc_auc_score import seaborn as sns from sklearn.linear_model import LogisticRegression from sklearn.ensemble import RandomForestClassifier from sklearn.neighbors import KNeighborsClassifier from sklearn.feature_selection import chi2 from matplotlib.pyplot import axhline from sklearn.preprocessing import PolynomialFeatures from sklearn import tree from sklearn naive_bayes import GaussianNB from sklearn preprocessing import StandardScaler MinMaxScaler
	from sklearn.preprocessing import StandardScaler, MinMaxScaler from sklearn.decomposition import PCA  THE DATA SET  The data set I am using contains information on data scientists and similar professions, such as Machine Learning Engineer, Data Analyst, Data Engineer, etc.  Attributes:
	work_year: The year the salary was paid.  experience_level: The experience level in the job during the year with the following possible values: EN Entry-level / Junior MI Mid-level / Intermediate SE Senior-level / Expert EX Executive-level / Director employment_type: The type of employment for the role: PT Part-time FT Full-time CT Contract FL Freelance job_title: The role worked in during the year.  salary: The total gross salary amount paid.  salary_currency: The currency of the salary paid as an ISO 4217 currency code.  salary_in_usd: The salary in USD (FX rate divided by avg. USD rate for the respective year via fxdata.foorilla.com).
In [2]:	employee_residence: Employee's primary country of residence in during the work year as an ISO 3166 country code.  remote_ratio: The overall amount of work done remotely, possible values are as follows: 0 No remote work (less than 20%) 50 Partially remote 100 Fully remote (more than 80%) company_location: The country of the employer's main office or contracting branch as an ISO 3166 country code.  company_size: The average number of people that worked for the company during the year: S less than 50 employees (small) M 50 to 250 employees (medium) L more than 250 employees (large)  df = pd.read_csv("ds_salaries.csv")  df
Out[2]:	Unnamed: 0 work_year experience_level employment_type job_title salary salary_currency salary_in_usd employee_residence remote_ratio company_location company_size  1 0 0 2020 MI FT Data Scientist 70000 EUR 79833 DE 0 DE L  1 1 2020 SE FT Machine Learning Scientist 260000 USD 260000 JP 0 JP S  2 2 2020 SE FT Big Data Engineer 85000 GBP 109024 GB 50 GB M  3 3 2020 MI FT Product Data Analyst 20000 USD 20000 HN 0 HN S  4 2020 SE FT Machine Learning Engineer 150000 USD 150000 USD 50 US 50 US L
	III       I
	DATA CLEANING  1. Filtering and selecting columns
In [3]:	<pre>df = pd.read_csv("ds_salaries.csv")  #dropping index column that was included in csv file df = df.drop(columns = [df.columns[0]])  #dropping all rows that contain NAs. In this dataset, there were none. df = df.dropna() df</pre>
Out[3]:	work_yearexperience_levelemployment_typejob_titlesalarysalary_currencysalary_in_usdemployee_residenceremote_ratiocompany_locationcompany_location02020MIFTData Scientist7000EUR79833DE0DEL12020SEFTMachine Learning Scientist260000USD260000JP0JPS22020SEFTBig Data Engineer85000GBP109024GB50GBM32020MIFTProduct Data Analyst20000USD20000HN0HNS42020SEFTMachine Learning Engineer150000USD150000US50USL
	602         2022         SE         FT         Data Engineer         154000         USD         154000         US         100         US         M           603         2022         SE         FT         Data Engineer         126000         USD         126000         US         100         US         M           604         2022         SE         FT         Data Analyst         129000         USD         129000         US         0         US         M           605         2022         SE         FT         Data Analyst         150000         USD         150000         US         100         US         M           606         2022         MI         FT         AI Scientist         200000         USD         200000         IN         100         US         L
	607 rows × 11 columns  2. Identifying and handling outliers  a. salary_in_usd  Although that there are many minimum outliers with salary prices that seem unrealistic in the US, I have decided to keep them in the dataset as they are from other countries that could have different wages.
In [4]: Out[4]:	<pre>df.sort_values(by = ['salary_in_usd'])</pre>
	77         2021         MI         PT         3D Computer Vision Researcher         400000         INR         5409         IN         50         IN         M           179         2021         MI         FT         Data Scientist         420000         INR         5679         IN         100         US         S <th< td=""></th<>
	33         2020         MI         FT         Research Scientist         450000         USD         450000         US         0         US         M           252         2021         EX         FT         Principal Data Engineer         600000         USD         600000         US         100         US         L           607 rows × 11 columns         b. work_year         There are no outliers or incorrect entries.         The principal Data Engineer         450000         USD         600000         US         100         US         L
In [5]: Out[5]:	<pre>df['work_year'].unique() array([2020_2021_2022])</pre>
In [6]: Out[6]:	orrow/[IMT] ICEL IENI IEVI] dtymo-object)
In [7]: Out[7]:	<pre>df['employment_type'].unique() array(['FT', 'CT', 'PT', 'FL'], dtype=object)  e. salary_currency There are no outliers of incorrect entries.</pre>
In [8]:	<pre>df['salary_currency'].unique() array(['EUR', 'USD', 'GBP', 'HUF', 'INR', 'JPY', 'CNY', 'MXN', 'CAD',</pre>
In [9]: Out[9]:	2552V/[IDE]   ID]   ICD]   IM]   IUC]   IM]   ID]   ID]   ID]
In [10]: Out[10]:	g. remote_ratio  There are no outliers or inncorrect entries.  df['remote_ratio'].unique()  array([ 0, 50, 100])  h. company_size
	There are no outliers or inncorrect entries.  df['company_size'].unique() array(['L', 'S', 'M'], dtype=object)  i. job_title
In [12]: Out[12]:	'Big Data Engineer', 'Product Data Analyst', 'Machine Learning Engineer', 'Data Analyst', 'Lead Data Scientist', 'Business Data Analyst', 'Lead Data Engineer', 'Lead Data Analyst', 'Data Engineer', 'Data Science Consultant', 'BI Data Analyst', 'Director of Data Science', 'Research Scientist', 'Machine Learning Manager', 'Data Engineering Manager',
	'Machine Learning Infrastructure Engineer', 'ML Engineer', 'AI Scientist', 'Computer Vision Engineer', 'Principal Data Scientist', 'Data Science Manager', 'Head of Data', '3D Computer Vision Researcher', 'Data Analytics Engineer', 'Applied Data Scientist', 'Marketing Data Analyst', 'Cloud Data Engineer', 'Financial Data Analyst', 'Computer Vision Software Engineer', 'Director of Data Engineering', 'Data Science Engineer', 'Principal Data Engineer', 'Machine Learning Developer', 'Applied Machine Learning Scientist', 'Data Analytics Manager',
	'Head of Data Science', 'Data Specialist', 'Data Architect', 'Finance Data Analyst', 'Principal Data Analyst', 'Big Data Architect', 'Staff Data Scientist', 'Analytics Engineer', 'ETL Developer', 'Head of Machine Learning', 'NLP Engineer', 'Lead Machine Learning Engineer', 'Data Analytics Lead'], dtype=object)  DATA PREPROCESSING
In [13]:	
In [43]:	<pre>labels = ['1', '2', '3', '4'] df['salary_bins'] = pd.cut(df['salary_in_usd'], bins = bins, labels = labels, include_lowest = True) df = df.sort_values(by = ['salary_in_usd'])  plt.hist(df['salary_bins']) plt.title(label = "Distribution of Salaries") plt.xticks(labels) plt.xlabel("Salary Range in USD") plt.ylabel("Frequency") plt.show()</pre>
	Distribution of Salaries
	200 -
	50 - 0 1 2 3 4 Salary Range in USD
In [15]:	<pre>df.loc[df['experience_level'] == 'MI', 'experience_level'] = 0.25</pre>
In [16]:	<pre>df.loc[df['experience_level'] == 'SE', 'experience_level'] = 0.5 df.loc[df['experience_level'] == 'EX', 'experience_level'] = 1  b. company_size  Converting the categorical labels to ordinal, since the ranking matters.  df.loc[df['company_size'] == 'S', 'company_size'] = 0 df.loc[df['company_size'] == 'M', 'company_size'] = 0.5</pre>
	<ul> <li>df.loc[df['company_size'] == 'L', 'company_size'] = 1</li> <li>3. One-hot Encoding</li> <li>a. Creating dummy variables for the following categorical variables:</li> <li>1. employment_type</li> <li>2. salary_currency</li> </ul>
In [18]: In [19]:	<pre>3. employee_residence 4. company_location  df = pd.get_dummies(df, columns = ['employment_type', 'salary_currency', 'employee_residence', 'company_location'], drop_first = False)  b. Creating new dummy variable columns for job_title  job_ml = []</pre>
	<pre>job_analyst = [] job_data_scientist = [] job_data_engineer = [] for i in range(0, len(df)):     #if job_title contains "Machine Learning", "AI", "Computer Vision" or "NLP", we label it as a ML job     if ("Machine Learning" in df['job_title'].iloc[i]) or ("AI" in df['job_title'].iloc[i]) or ("Computer Vision" in df['job_title'].iloc[i]) or ("NLP" in df['job_title'].iloc[i]) or ("NLP" in df['job_title'].iloc[i]) else:         job_ml.append(0)     #if job_title contains "Data Analyst" or "Analytics", we label it as a data analyst job     if ("Data Analyst" in df['job_title'].iloc[i]) or ("Analytics" in df['job_title'].iloc[i]):</pre>
	<pre>job_analyst.append(1) else:     job_analyst.append(0) #if job_title contains "Data Scientist" or "Data Science", we label it as a data scientist job if ("Data Scientist" in df['job_title'].iloc[i]) or ("Data Science" in df['job_title'].iloc[i]):     job_data_scientist.append(1) else:     job_data_scientist.append(0) #if job_title contains "Data Engineer" or "Architect", we label it as a data engineering job if ("Data Engineer" in df['job_title'].iloc[i]) or ("Architect" in df['job_title'].iloc[i]):</pre>
	<pre>job_data_engineer.append(1) else:     job_data_engineer.append(0)  df["job_ml"] = job_ml  df["job_analyst"] = job_analyst  df["job_data_scientist"] = job_data_scientist  df["job_data_engineer"] = job_data_engineer</pre>
In [20]: Out[20]:	df
	238       2021       0       0       0.5       1       0       0       1       0<
	97 2021 0.25 100 1 4 0 1 0 0 0 0 0 0 0 0 3 33 2020 0.25 0 0.5 4 0
In [21]:	<pre>X = df.loc[:, df.columns != 'salary_bins'] #features y = df['salary_bins'] #target X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.5)</pre> MACHINE LEARNING ALGORITHMS
In [22]:	<pre>1. Logistic Regression  def log_reg(X_train, y_train, X_test, y_test):     scaler = StandardScaler()     X_train_scaled = scaler.fit_transform(X_train)     X_test_scaled = scaler.fit_transform(X_test)     classifier = LogisticRegression(multi_class = "ovr")     classifier.fit(X_train_scaled, y_train)     y_pred = classifier.predict(X_test_scaled)</pre>
In [23]:	<pre>return y_pred  2. Gaussian Naive Bayes  def naive_bayes(X_train, y_train, X_test, y_test):     classifier = GaussianNB()     classifier.fit(X_train, y_train)     y_pred = classifier.predict(X_test)</pre>
In [24]:	<pre>return y_pred  3. Random Forest  def rand_forest(X_train, y_train, X_test, y_test):     classifier = RandomForestClassifier()     classifier.fit(X_train, y_train)     y_pred = classifier.predict(X_test)     return y_pred</pre>
In [25]:	
In [26]:	<pre>scaler = StandardScaler() X_train_scaled = scaler.fit_transform(X_train) X_test_scaled = scaler.fit_transform(X_test) classifier = KNeighborsClassifier(n_neighbors = i) classifier.fit(X_train_scaled, y_train) y_pred = classifier.predict(X_test_scaled) return y_pred  #using number of clusters, k from k = 1 to k = 15 k = np.arange(1, 15, 1)</pre>
	<pre>#getting accuracies for each k number of clusters accuracies = [] for i in k:     y_pred_kNN = kNN(X_train, y_train, X_test, y_test, i)     accuracies.append(accuracy_score(y_pred_kNN, y_test))  plt.plot(k, accuracies) plt.xticks(k) plt.show()</pre>
	0.70 - 0.68 -
	0.64 - 0.62 -
	DIMENSIONALITY REDUCTION
In [27]:	1. Chi-square Test  We will use the Chi-square test for dimensionality reduction. We will remove any attribute that has a p-value of less than 0.1.  #Chi-Square Test chi_scores = chi2(X_train, y_train) results = pd.DataFrame({'Attribute' : X.columns, 'p-value' : chi_scores[1]}) results = results.sort_values(by = ['p-value'], ascending = False)
	<pre>f = plt.figure() f.set_figwidth(18) f.set_figheight(4) plt.bar(results['Attribute'], results['p-value']) plt.xticks(rotation = 90, fontsize = 7) plt.title(label = "Distribution of Salaries") plt.xlabel("Attribute") plt.ylabel("p-value") axhline(0.1, color = "red")</pre>
	<pre>#Reduced Dataset remove = list(results[results[results['p-value'] &gt; 0.1]['Attribute']) X_train_1 = X_train.drop(columns = remove) X_test_1 = X_test.drop(columns = remove)</pre> Distribution of Salaries
	0.8 - 0.6 - 0.4 -
	employee residence Aremployee residence Bremployee residence Aremployee residence Aremployee residence Bremployee
	MODEL EVALUATION  1. Without dimensionality reduction
	<pre>y_pred_lr = log_reg(X_train, y_train, X_test, y_test) y_pred_nb = naive_bayes(X_train, y_train, X_test, y_test) y_pred_rf = rand_forest(X_train, y_train, X_test, y_test) y_pred_kNN = kNN(X_train, y_train, X_test, y_test, 4) y_preds = [y_pred_lr, y_pred_nb, y_pred_rf, y_pred_kNN]  fig, (ax1, ax2, ax3, ax4) = plt.subplots(ncols=4, sharey=True, figsize=(15, 3)) mat_lr = confusion_matrix(y_test, y_pred_lr) mat_nb = confusion_matrix(y_test, y_pred_nb) mat_rf = confusion_matrix(y_test, y_pred_rf)</pre>
	mat_rf = confusion_matrix(y_test, y_pred_rf) mat_kNN = confusion_matrix(y_test, y_pred_kNN) sns.heatmap(mat_lr.T,square=True, annot=True, fmt = 'd', cbar=True, xticklabels=[1,2,3,4], yticklabels=[1,2,3,4], ax = ax1).set(title = "Logistic Regression") sns.heatmap(mat_nb.T,square=True, annot=True, fmt = 'd', cbar=True, xticklabels=[1,2,3,4], yticklabels=[1,2,3,4], ax = ax2).set(title = "Naive Bayes") sns.heatmap(mat_rf.T,square=True, annot=True, fmt = 'd', cbar=True, xticklabels=[1,2,3,4], yticklabels=[1,2,3,4], ax = ax3).set(title = "Random Forest") sns.heatmap(mat_kNN.T,square=True, annot=True, fmt = 'd', cbar=True, xticklabels=[1,2,3,4], yticklabels=[1,2,3,4], ax = ax4).set(title = "k-NN") plt.show()  Logistic Regression  Naive Bayes  Random Forest  - 120  - 100
	H - 125 23 1 0 -100 - 70 8 1 0 -80 - 126 33 2 1 -100 - 108 25 2 0 -80 N - 28 107 16 3 -80 - 44 14 1 0 -60 -24 94 14 2 -80 - 42 105 13 3 -60 M - 1 0 0 0 0 -80 - 40 - 40 - 40 - 4 3 0 0 -40 - 40 - 40
In [37]:	y_pred_lr_1 = log_reg(X_train_1, y_train, X_test_1, y_test) y_pred_nb_1 = naive_bayes(X_train_1, y_train, X_test_1, y_test)
In [38]:	<pre>y_pred_rf_1 = rand_forest(X_train_1, y_train, X_test_1, y_test) y_pred_kNN_1 = kNN(X_train_1, y_train, X_test_1, y_test, 4) y_preds_1 = [y_pred_lr_1, y_pred_nb_1, y_pred_rf_1, y_pred_kNN_1]  fig, (ax1, ax2, ax3, ax4) = plt.subplots(ncols=4, sharey=True, figsize=(15, 3)) mat_lr_1 = confusion_matrix(y_test, y_pred_lr_1) mat_nb_1 = confusion_matrix(y_test, y_pred_nb_1) mat_rf_1 = confusion_matrix(y_test, y_pred_rf) mat_kNN_1 = confusion_matrix(y_test, y_pred_kNN_1) sns.heatmap(mat_lr_1.T,square=True, annot=True, fmt = 'd', cbar=True, xticklabels=[1,2,3,4], yticklabels=[1,2,3,4], ax = ax1).set(title = "Logistic Regression")</pre>
	sns.heatmap(mat_nb_1.T, square=True, annot=True, fmt = 'd', cbar=True, xticklabels=[1,2,3,4], yticklabels=[1,2,3,4], ax = ax2).set(title = "Naive Bayes") sns.heatmap(mat_rf_1.T, square=True, annot=True, fmt = 'd', cbar=True, xticklabels=[1,2,3,4], yticklabels=[1,2,3,4], ax = ax3).set(title = "Random Forest") sns.heatmap(mat_kNN_1.T, square=True, annot=True, fmt = 'd', cbar=True, xticklabels=[1,2,3,4], yticklabels=[1,2,3,4], ax = ax4).set(title = "k-NN") plt.show()    Logistic Regression
	$\begin{array}{cccccccccccccccccccccccccccccccccccc$
In [39]:	<pre>3. Performance  accuracies = [] for i in y_preds:     accuracies.append(accuracy_score(i, y_test))  accuracies_1 = []</pre>
In [40]:	<pre>for i in y_preds_1:     accuracies_1.append(accuracy_score(i, y_test))</pre>
In [41]: Out[41]:	
	table = pd.DataFrame({"Model": models, "Accuracy without Dimensionality Reduction": accuracies, "Accuracy with Dimensionality Reduction": accuracies_1})  table  Model Accuracy without Dimensionality Reduction Accuracy with Dimensionality Reduction