| | # Importing the dataset, downloaded file is in the same folder csv_path = "Social_Network_Ads.csv" ads_df = pd.read_csv(csv_path) Run a quick check on the dataset |
|--|--|
| In [3]: Out[3]: In [4]: Out[4]: | ads_df.shape (400, 3) ads_df.head(5) Age |
| | 3 27 57000 0 4 19 76000 0 ads_df.info() <class 'pandas.core.frame.dataframe'=""> RangeIndex: 400 entries, 0 to 399 Data columns (total 3 columns): # Column Non-Null Count Dtype</class> |
| | memory usage: 9.5 KB The independent variables are all numeric and there's no missing data ads_df.describe() Age EstimatedSalary Purchased count 400.000000 400.000000 400.000000 mean 37.655000 69742.500000 0.357500 std 10.482877 34096.960282 0.479864 min 18.00000 15000.000000 0.000000 25% 29.750000 43000.000000 0.000000 50% 37.000000 70000.000000 0.000000 |
| In [7]: | 75% 46.00000 88000.00000 1.000000 max 60.00000 150000.000000 1.000000 #Extracting independent variables: x = ads_df.iloc[:,:-1].values #Extract semua kolom kecuali kolom terakhir print(x) [[19 19000] [35 20000] [26 43000] [27 57000] [19 76000] [27 58000] |
| | [27 8400] [32 150000] [25 33000] [35 65000] [26 80000] [20 86000] [32 18000] [32 18000] [47 25000] [47 25000] [48 29000] [48 29000] [48 29000] [48 29000] [48 41000] |
| | [45 22000] [46 23000] [47 20000] [49 28000] [29 43000] [31 18000] [31 74000] [27 137000] [28 44000] [28 44000] [28 44000] [38 27000] [38 27000] [37 3000] [38 28000] [38 28000] [39 43000] [30 43000] [30 43000] [31 27 31000] |
| | [27 1700] [33 51000] [35 108000] [30 15000] [28 84000] [23 20000] [25 79000] [27 54000] [31 89000] [31 89000] [24 32000] [18 44000] [29 83000] [27 58000] [27 58000] [27 58000] |
| | [28 7900] [22 18000] [32 117000] [27 20000] [28 87000] [23 66000] [32 120000] [59 83000] [24 58000] [22 63000] [22 63000] [31 68000] [22 82000] [24 27000] [31 31 68000] [24 27000] [25 80000] [31 31 68000] [32 18000] |
| | [34 112000] [18 52000] [22 27000] [28 87000] [30 80000] [30 42000] [30 42000] [30 62000] [31 118000] [32 455000] [28 85000] [28 85000] [28 85000] [28 85000] [28 85000] [28 85000] [28 85000] [26 81000] [35 50000] [30 116000] [30 116000] |
| | [29 88000] [35 44000] [35 25000] [38 123000] [38 37000] [28 37000] [27 88000] [32 86000] [33 149000] [19 21000] [27 72000] [26 35000] [27 88000] [38 80000] [39 71000] |
| | [37 71000] [38 61000] [42 80000] [42 80000] [40 57000] [36 52000] [40 59000] [41 59000] [37 72000] [40 75000] [41 51000] [42 65000] [42 65000] [43 69000] [44 65000] |
| | [26 84000] [31 58000] [33 31000] [30 87000] [21 68000] [28 55000] [20 82000] [20 82000] [28 59000] [19 25000] [19 85000] [18 68000] [33 89000] [34 25000] [34 25000] [27 96000] |
| | [41 30000] [29 61000] [20 74000] [26 15000] [41 45000] [31 76000] [40 47000] [31 15000] [46 59000] [29 75000] [29 30000] [25 90000] [32 135000] [32 135000] [37 33000] [37 33000] [38 38000] [37 38000] |
| | [18 86000] [22 55000] [35 71000] [29 148000] [29 47000] [21 88000] [34 115000] [34 3000] [34 72000] [35 47000] [25 22000] [26 16000] [31 34000] [31 71000] |
| | [33 4300] [31 66000] [20 82000] [33 41000] [35 72000] [28 32000] [24 84000] [19 26000] [29 43000] [19 70000] [28 89000] [34 43000] [35 79000] [36 89000] [37 89000] [38 89000] [38 89000] [39 79000] [30 79000] [30 79000] [31 43000] [32 38000] |
| | [49 74000] [39 134000] [41 71000] [58 101000] [47 47000] [55 130000] [40 142000] [48 96000] [59 42000] [35 58000] [47 43000] [48 96000] [49 65000] [49 65000] [40 78000] [40 78000] [40 78000] |
| | [59 143000] [41 80000] [35 91000] [37 144000] [60 102000] [37 53000] [36 126000] [56 133000] [40 72000] [42 80000] [35 147000] [42 80000] [44 80000] [45 147000] [46 79000] [49 86000] [49 86000] |
| | [37 80000] [46 82000] [53 143000] [42 149000] [38 59000] [56 104000] [41 72000] [51 146000] [35 50000] [41 52000] [41 52000] [41 52000] [41 52000] [41 39000] [44 39000] [37 52000] [48 134000] [37 146000] |
| | [52 90000] [41 72000] [40 57000] [58 95000] [45 131000] [35 77000] [36 144000] [35 125000] [48 90000] [48 90000] [40 75000] [47 144000] [47 144000] [48 61000] [40 61000] [40 61000] [40 642000] |
| | [39 106000] [57 26000] [57 74000] [38 71000] [49 88000] [52 38000] [59 88000] [35 61000] [37 70000] [48 141000] [37 93000] [48 138000] [47 79000] [48 138000] [48 138000] [48 138000] [48 138000] [48 138000] [48 138000] |
| | [49 89000] [37 77000] [35 57000] [36 63000] [42 73000] [43 112000] [45 79000] [46 117000] [58 38000] [48 74000] [37 137000] [37 137000] [40 60000] [42 54000] [42 54000] [43 134000] [47 113000] [36 125000] |
| | [38 50000] [42 70000] [39 96000] [38 50000] [49 141000] [39 75000] [54 104000] [35 55000] [45 32000] [45 32000] [52 138000] [53 82000] [41 52000] [41 52000] [41 52000] [41 52000] [41 52000] [41 52000] |
| | [42 75000] [36 118000] [47 107000] [48 119000] [48 119000] [40 65000] [57 60000] [36 54000] [35 79000] [38 55000] [38 55000] [39 122000] [39 122000] [37 75000] [38 65000] [47 51000] [47 510000] [47 510000] |
| | [53 72000] [54 108000] [39 77000] [38 61000] [38 113000] [42 90000] [42 90000] [36 99000] [60 34000] [60 34000] [41 72000] [40 71000] [42 54000] [42 54000] [42 54000] [43 129000] [47 50000] [47 50000] |
| | [42 104000] [59 29000] [58 47000] [46 88000] [38 71000] [60 46000] [60 83000] [39 73000] [37 80000] [46 32000] [46 74000] [47 74000] [42 53000] [42 64000] [42 64000] [43 83000] |
| | [44 13900] [49 28000] [57 33000] [56 60000] [49 39000] [47 34000] [48 35000] [48 33000] [47 23000] [45 45000] [60 42000] [39 59000] [46 41000] [51 23000] [50 20000] [36 33000] [49 36000]] |
| In [8]: | #Extracting dependent variable: y = ads_df.iloc[:,2].values #Extract kolom terakhir print(y) [0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 1 1 1 1 |
| In [9]: | Splitting training and test data Before we do feature scaling, we must split our data first With 400 data, I decided to split the dataset into 3:1 ratio (Training set contains 300 data, Test set contains 100 data) from sklearn.model_selection import train_test_split x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.25, random_state = 0) print('Train set: ', x_train.shape, y_train.shape) print('Test set: ', x_test.shape, y_test.shape) Train set: (300, 2) (300,) Test set: (100, 2) (100,) Feature Scaling |
| n [12]: | <pre>from sklearn.preprocessing import StandardScaler sc = StandardScaler() x_train = sc.fit_transform(x_train) x_test = sc.fit_transform(x_test) x_train[0:5] array([[0.58164944, -0.88670699],</pre> |
| n [13]: | array([[-0.54748976, 0.5130727], |
| ut[14]: n [15]: | <pre>mean_acc = np.zeros((K-1)) std_acc = np.zeros((K-1)) ConfusionMatrix = [] for n in range(1,K): #Training the model and predict neighb = KNeighborsClassifier(n_neighbors = n, metric = 'minkowski').fit(x_train, y_train) #using Minkowski distance metric yhat = neighb.predict(x_test) mean_acc[n-1] = metrics.accuracy_score(y_test, yhat) std_acc[n-1] = np.std(yhat == y_test) / np.sqrt(yhat.shape[0]) mean_acc array([0.88, 0.88, 0.93, 0.93, 0.93, 0.93, 0.93, 0.93, 0.93])</pre> plt.plot(range(1,K), mean_acc, 'g') plt.fill_between(range(1,K), mean_acc - 1 * std_acc, mean_acc + 1 * std_acc, alpha = 0.1) |
| | plt.legend(('Accuracy', '+/- 3xstd')) plt.ylabel('Accuracy') plt.xlabel('Number of Neighbors (K)') plt.tight_layout() plt.show() |
| : | print("The best accuracy was with", mean_acc.max(), "with k=", mean_acc.argmax()+1) The best accuracy was with 0.93 with k= 3 We found that k = 3 resulted in best accuracy (even though the same thing could be said for k = 4 to k = 9) So, for the classifier, we'll use K = 3 Predict using K-NN Classifier with K = 3 |
| n [18]: | <pre>classifier = KNeighborsClassifier(n_neighbors = 3, metric = 'minkowski') classifier.fit(x_train, y_train) classifier KNeighborsClassifier(n_neighbors=3) #Predicting y_pred = classifier.predict(x_test) Evaluating accuracy using confusion matrix from sklearn.metrics import confusion_matrix,accuracy_score conf_mtx = confusion_matrix(y_test, y_pred) accuracy = accuracy_score(y_test, y_pred)</pre> |
| n [21]: ut[21]: | <pre>accuracy 0.93 conf_mtx array([[64, 4], [3, 29]], dtype=int64) #Visualizing confusion matrix sns.heatmap(conf_mtx, annot = True) <axessubplot:></axessubplot:></pre> |
| | - 64 4 - 50 - 40 - 30 - 20 - 10 - 10 |
| n [23]: | <pre>Visualizing the Training and Test set results from matplotlib.colors import ListedColormap x1, x2 = np.meshgrid(np.arange(start = x_train[:, 0].min()-1, stop = x_train[:, 0].max()+1, step = 0.01),</pre> |
| | plt.title('K-NN (Training Set)') plt.xlabel('Age') plt.ylabel('Estimated Salary') plt.legend() plt.show() *c* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length ches with *x* & *y*. Please use the *color* keyword-argument or provide a 2D array with a single row if you intend to specify the same RGB or value for all points. *c* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length ches with *x* & *y*. Please use the *color* keyword-argument or provide a 2D array with a single row if you intend to specify the same RGB or value for all points. K-NN (Training Set) **C** **C* |
| n [24]: | x1, x2 = np.meshgrid(np.arange(start = x_test[:, 0].min()-1, stop = x_test[:, 0].max()+1, step = 0.01), |
| | <pre>alpha = 0.75, cmap = ListedColormap(('red', 'green'))) plt.xlim(x1.min(), x1.max()) plt.ylim(x2.min(), x2.max()) for i, j in enumerate(np.unique(y_test)): plt.scatter(x_test[y_test == j, 0], x_test[y_test == j, 1], c = ListedColormap(('red', 'green'))(i), label = j) plt.title('K-NN (Test Set)') plt.xlabel('Age') plt.ylabel('Estimated Salary') plt.legend() plt.show()</pre> |
| | *c* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length ches with *x* & *y*. Please use the *color* keyword-argument or provide a 2D array with a single row if you intend to specify the same RGB or Foundation value for all points. *c* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length ches with *x* & *y*. Please use the *color* keyword-argument or provide a 2D array with a single row if you intend to specify the same RGB or Foundation value for all points. *K-NN (Test Set) **C** Argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length ches with *x* & *y*. Please use the *color* keyword-argument or provide a 2D array with a single row if you intend to specify the same RGB or Foundation value for all points. **K-NN (Test Set)** **C** Brown and **C** Color** RGBA sequence, which should be avoided as value-mapping will have precedence in case its length ches with *x* & *y*. Please use the *color* keyword-argument or provide a 2D array with a single row if you intend to specify the same RGB or Foundation value for all points. **C** Brown and **C** Color** RGBA sequence, which should be avoided as value-mapping will have precedence in case its length ches with *x* & *y*. Please use the *color* keyword-argument or provide a 2D array with a single row if you intend to specify the same RGB or Foundation value for all points. **C** Brown and **C** Color* Color* RGBA sequence, which should be avoided as value-mapping will have precedence in case its length ches will be avoided as value-mapping will have precedence in case its length ches will be avoided as value-mapping will have precedence in case its length ches will be avoided as value-mapping will have precedence in case its length ches will be avoided as value-mapping will have precedence in case its leng |
| | -2 -1 0 1 2 3 Age |
| | |
| | |

2440016804 - Rio Pramana - LA01 - Assignment 7