

```
In [1]: # Hello Folks Hirshay Here a data scientist working for gurgan based logistics(E-commerce) org named HishuPost,
# Today i will be sharing very basic linear regression proj that will help you understand how we deal with real world regression problems using ML algorithms.

# So i am picking here a very basic Linear Regression Dataset where we will find that in respect to Year of exp of a person how salary changes in real world
# Remember that before picking any algorithm in data science project, we have to always figureout what type of problem statement we are going to deal with & which algorithm gonna suit it.
# Linear Regression applies where we have to predict continous value lets say stock price, house price, salary of person, age of a person

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')

In [2]: # Now let me quickly Load the basic dataset which has 2 attributes YearofExp & Salary in it

df = pd.read_csv("Salary_Data.csv")

In [3]: df.head(10)

In [10]: # Now lets check the no of rows and no of columns our dataset contains
# Remember that left side it always shows rows count & and right side cols count

df.shape

Out[10]:
(30, 2)

In [11]: # Now Will check the Dataset Information using

df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30 entries, 0 to 29
Data columns (total 2 columns):
 #   Column                Non-Null Count  Dtype  
---  --
 0   YearsExperience        30 non-null      float64
 1   Salary                 30 non-null      int64   
dtypes: float64(1), int64(1)
memory usage: 668.0 bytes

In [13]: # We can clearly see now, that our dataset do not contains any null values/missing values in any columns as well as, Yearexp column has floating decimal places values in it
# Whereas Salary column has integer values in it

In [14]: # Now i will be using describe function that will help me understand the statistical summary of my dataset

df.describe(include='all')

Out[14]:
      YearsExperience      Salary
count      30.000000      30.000000
mean         5.133333      76033.000000
std          2.837888      7414.429785
min          1.000000      3731.000000
25%          3.000000      49690.000000
50%          4.700000      65237.000000
75%          7.700000      105544.750000
max          10.000000      122391.000000

In [15]: # As we can see from the summary its giving us very meaningful info like:-
# Our data contains 30 different observations(records)/Rows Count
# Avg years of experience of our 30 observations(records)/Employee is 5, means most of employees are having around 5 yrs of exp
# Avg salary of people having 5 yrs of exp is 76085 INR per month
# Stddev is telling us spread or dispersion of yrs of exp & salary around the avg yrs of exp and avg salary
# Or to other words Stddev values telling us that (around 68% of the salary of this dataset falls around 76085+27414)/ Uncase of Stddev of yrs of exp as well
# Next the summary telling us that minimum exp of a person in this 30 observations is 1.1 yrs and minimum salary of a person is 3731 INR / Month Salary
Now 25% is nothing but a quartiles or quartiles, its saying that 25% of people from the dataset are peoples who have 3.1 yrs or less of exp and 49620 INR / Monthly Income
# 50% of people in this dataset having less than or equal to 4.7 yrs of exp and less than or equal to 65237 INR / Monthly Income
75% can be 25% & 50% Explanation
# Last but not least am telling us that maximum experience of our dataset is 10.5 you can say we have a candidate record who has 10.5 yrs of exp and is the highest over all and
# He's salary is also the highest amongst all other candidate (records/observation) we have in our dataset

In [16]: # Now we check if we have any missing values

df.isnull().sum()

# so what we doing here is we simply saying first give us all values that is null in our dataset then sum it all to give us final count of null values

Out[16]:
YearsExperience      0
Salary              0
dtype: int64

In [17]: # No Null values found in our data that good

In [20]: # Now the next very imp step in exploratory data analysis is finding if our data contains any outliers or not
# What is outliers suppose it as that is very dissimilar value in your dataset from majority of datapoints
# If i take yrs of exp example then lets say if we find candidate with 50 yrs of exp that can be considered as outliers

df.skew()

Out[20]:
YearsExperience      0.37956
Salary              0.35412
dtype: float64

In [21]: # You Folks must be thinking what this skew function does in finding outliers
# Skew function measures the asymmetry of the distribution of values in a variable.
# If the skewness is close to 0, it indicates that the distribution is approximately symmetric.
# If the skewness is positive, it indicates that the distribution is right-skewed (longer tail on the right side).
# If the skewness is negative, it indicates that the distribution is left-skewed (longer tail on the left side).
# Now lets plot the skewness of dataset to properly see and understand

skewness = df.skew()

# Plot the skewness values
plt.figure(figsize=(10, 6))
skewness.plot(kind='bar', color='skyblue')
plt.title('Skewness of Numerical Variables')
plt.xlabel('Variables')
plt.ylabel('Skewness')
plt.xticks(rotation=45)
plt.grid(axis='y', linestyle='--', alpha=0.7)
plt.show()

In [31]: from scipy.stats import norm, skew

# Lets Calculate & Plot skewness of the 'YearsExperience' and 'Salary' columns which shows if our data is normally distributed or right skewed or left skewed

# Remember that if the dataset do not follow normal data distribution it will impact the model performance & accuracy

years_exp_skewness = skew(df['YearsExperience'])
salary_skewness = skew(df['Salary'])

# Create subplots
fig, axs = plt.subplots(2, 2, figsize=(14, 6))

# Plotting histogram for 'YearsExperience'
axs[0].hist(df['YearsExperience'], bins=10, density=True, alpha=0.6, color='b')
axs[0].set_title('Years of Experience (Skewness: {years_exp_skewness:.2f})')
axs[0].set_xlabel('Years of Experience')
axs[0].set_ylabel('Density')
axs[0].grid(True)

# Plot bell curve (normal distribution curve) for 'YearsExperience'
xmin, xmax = axs[0].get_xlim()
x = np.linspace(xmin, xmax, 100)
p = norm.pdf(x, df['YearsExperience'].mean(), df['YearsExperience'].std())
axs[0].plot(x, p, 'k', linewidth=2)

# Plot histogram for 'Salary'
axs[1].hist(df['Salary'], bins=10, density=True, alpha=0.6, color='r')
axs[1].set_title('Salary (Skewness: {salary_skewness:.2f})')
axs[1].set_xlabel('Salary')
axs[1].set_ylabel('Density')
axs[1].grid(True)

# Plot bell curve (normal distribution curve) for 'Salary'
xmin, xmax = axs[1].get_xlim()
x = np.linspace(xmin, xmax, 100)
p = norm.pdf(x, df['Salary'].mean(), df['Salary'].std())
axs[1].plot(x, p, 'k', linewidth=2)

plt.tight_layout()
plt.show()

In [32]: # Now we will check the covariance in our dataset
# What is Covariance ?
# Covariance is a measure of the degree to which two variables change together. It indicates the direction of the linear relationship between two variables. Specifically:
# Positive covariance: Indicates that as one variable increases, the other variable also tends to increase. Similarly, as one variable decreases, the other variable tends to decrease.
# This suggests a positive linear relationship between the variables.
# Negative covariance: Indicates that as one variable increases, the other variable tends to decrease, and vice versa.
# This suggests a negative linear relationship between the variables.

df[['YearsExperience', 'Salary']].cov()

Out[32]:
      YearsExperience      Salary
Salary      6.92305e+06  7.01000e+06
Salary      76106.303440  7515102e+06

In [33]: # Now we will check the correlation in our dataset
# What is correlation ?
# We find correlation to understand the strength and direction of the linear relationship between two variables. Correlation is a standardized measure that ranges from -1 to 1:
# A correlation of 1 indicates a perfect positive linear relationship: as one variable increases, the other variable also increases linearly.
# A correlation of -1 indicates a perfect negative linear relationship: as one variable increases, the other variable decreases linearly.
# A correlation of 0 indicates no linear relationship between the variables.

df.corr()

Out[33]:
      YearsExperience      Salary
YearsExperience      1.000000  0.976242
Salary              0.976242  1.000000

In [34]: # We can clearly see that our data has a positive relationship between them which states that
# If the year of exp increase the salary of a person also increase linearly

In [35]: # Now lets plot a heatmap to visualize the relationship between input and output variable or independent or dependent variable or feature & target variable

sns.heatmap(df.corr(), annot=True)
plt.show()

# We use sns.heatmap(df.corr(), annot=True) to create a heatmap visualization of the correlation matrix for the variables in the dataset df using the Seaborn Library (sns).
# This visualization technique is particularly useful for quickly identifying the strength and direction of the linear relationships between variables.

In [36]: sns.pairplot(df)
plt.show()

# We use sns.pairplot(df) along with plt.show() to create a grid of scatterplots for pairwise relationships in the dataset df using the Seaborn Library (sns).
# This visualization technique is particularly useful for exploring the relationships between multiple variables in the dataset.

In [45]: plt.plot(df, linestyle = '-', linewidth=9.7)

plt.show()

In [52]: df.plot.line(linestyle = '-', linewidth=4)

plt.title('YearExperience-VS-Salary')
plt.show()

In [68]: df = pd.DataFrame(data=df)
df.plot.line(x='YearsExperience',linestyle = '-', linewidth=3)
plt.title('Salary Increase with respect to year of experience Increasing')
plt.show()

In [69]: # The Above Chart Clearly show thas the Linear relationship between dependent and independent variable.

In [70]: # Now Lets select the features(dependent & independent variables) & start Building then Training our model to predict the salary when the yrs of exp given due to model

X = df.drop('Salary', axis=1)
# X is Always be your independent variable or input features for which you want to predict the target variable or values

In [72]: y = df.Salary
# y gonna be always your Dependent/Target variable or output feature that you want to predict with respect to some value of X

In [81]: X.head(), y.head().tolist()

Out[81]:
([ YearsExperience
 0      1.1
 1      1.3
 2      1.5
 3      2.0
 4      2.2
 5      3.0
 6      3.3
 7      3.6
 8      4.0
 9      4.3
10     4.7
11     5.0
12     5.3
13     5.5
14     5.7
15     6.0
16     6.3
17     6.5
18     6.7
19     7.0
20     7.3
21     7.5
22     7.7
23     8.0
24     8.3
25     8.5
26     8.7
27     9.0
28     9.3
29     9.5
30    10.0],
 [3731, 46205, 49690, 6154, 65237, 69885, 76085, 79410, 81613, 84127, 86583, 89089, 91599, 94087, 96583, 99089, 101599, 104087, 106583, 109089, 111599, 114087, 116583, 119089, 121599, 124087, 126583, 129089, 131599, 134087, 136583, 139089, 141599, 144087, 146583, 149089, 151599, 154087, 156583, 159089, 161599, 164087, 166583, 169089, 171599, 174087, 176583, 179089, 181599, 184087, 186583, 189089, 191599, 194087, 196583, 199089, 201599, 204087, 206583, 209089, 211599, 214087, 216583, 219089, 221599, 224087, 226583, 229089, 231599, 234087, 236583, 239089, 241599, 244087, 246583, 249089, 251599, 254087, 256583, 259089, 261599, 264087, 266583, 269089, 271599, 274087, 276583, 279089, 281599, 284087, 286583, 289089, 291599, 294087, 296583, 299089, 301599, 304087, 306583, 309089, 311599, 314087, 316583, 319089, 321599, 324087, 326583, 329089, 331599, 334087, 336583, 339089, 341599, 344087, 346583, 349089, 351599, 354087, 356583, 359089, 361599, 364087, 366583, 369089, 371599, 374087, 376583, 379089, 381599, 384087, 386583, 389089, 391599, 394087, 396583, 399089, 401599, 404087, 406583, 409089, 411599, 414087, 416583, 419089, 421599, 424087, 426583, 429089, 431599, 434087, 436583, 439089, 441599, 444087, 446583, 449089, 451599, 454087, 456583, 459089, 461599, 464087, 466583, 469089, 471599, 474087, 476583, 479089, 481599, 484087, 486583, 489089, 491599, 494087, 496583, 499089, 501599, 504087, 506583, 509089, 511599, 514087, 516583, 519089, 521599, 524087, 526583, 529089, 531599, 534087, 536583, 539089, 541599, 544087, 546583, 549089, 551599, 554087, 556583, 559089, 561599, 564087, 566583, 569089, 571599, 574087, 576583, 579089, 581599, 584087, 586583, 589089, 591599, 594087, 596583, 599089, 601599, 604087, 606583, 609089, 611599, 614087, 616583, 619089, 621599, 624087, 626583, 629089, 631599, 634087, 636583, 639089, 641599, 644087, 646583, 649089, 651599, 654087, 656583, 659089, 661599, 664087, 666583, 669089, 671599, 674087, 676583, 679089, 681599, 684087, 686583, 689089, 691599, 694087, 696583, 699089, 701599, 704087, 706583, 709089, 711599, 714087, 716583, 719089, 721599, 724087, 726583, 729089, 731599, 734087, 736583, 739089, 741599, 744087, 746583, 749089, 751599, 754087, 756583, 759089, 761599, 764087, 766583, 769089, 771599, 774087, 776583, 779089, 781599, 784087, 786583, 789089, 791599, 794087, 796583, 799089, 801599, 804087, 806583, 809089, 811599, 814087, 816583, 819089, 821599, 824087, 826583, 829089, 831599, 834087, 836583, 839089, 841599, 844087, 846583, 849089, 851599, 854087, 856583, 859089, 861599, 864087, 866583, 869089, 871599, 874087, 876583, 879089, 881599, 884087, 886583, 889089, 891599, 894087, 896583, 899089, 901599, 904087, 906583, 909089, 911599, 914087, 916583, 919089, 921599, 924087, 926583, 929089, 931599, 934087, 936583, 939089, 941599, 944087, 946583, 949089, 951599, 954087, 956583, 959089, 961599, 964087, 966583, 969089, 971599, 974087, 976583, 979089, 981599, 984087, 986583, 989089, 991599, 994087, 996583, 999089, 1001599, 1004087, 1006583, 1009089, 1011599, 1014087, 1016583, 1019089, 1021599, 1024087, 1026583, 1029089, 1031599, 1034087, 1036583, 1039089, 1041599, 1044087, 1046583, 1049089, 1051599, 1054087, 1056583, 1059089, 1061599, 1064087, 1066583, 1069089, 1071599, 1074087, 1076583, 1079089, 1081599, 1084087, 1086583, 1089089, 1091599, 1094087, 1096583, 1099089, 1101599, 1104087, 1106583, 1109089, 1111599, 1114087, 1116583, 1119089, 1121599, 1124087, 1126583, 1129089, 1131599, 1134087, 1136583, 1139089, 1141599, 1144087, 1146583, 1149089, 1151599, 1154087, 1156583, 1159089, 1161599, 1164087, 1166583, 1169089, 1171599, 1174087, 1176583, 1179089, 1181599, 1184087, 1186583, 1189089, 1191599, 1194087, 1196583, 1199089, 1201599, 1204087, 1206583, 1209089, 1211599, 1214087, 1216583, 1219089, 1221599, 1224087, 1226583, 1229089, 1231599, 1234087, 1236583, 1239089, 1241599, 1244087, 1246583, 1249089, 1251599, 1254087, 1256583, 1259089, 1261599, 1264087, 1266583, 1269089, 1271599, 1274087, 1276583, 1279089, 1281599, 1284087, 1286583, 1289089, 1291599, 1294087, 1296583, 1299089, 1301599, 1304087, 1306583, 1309089, 1311599, 1314087, 1316583, 1319089, 1321599, 1324087, 1326583, 1329089, 1331599, 1334087, 1336583, 1339089, 1341599, 1344087, 1346583, 1349089, 1351599, 1354087, 1356583, 1359089, 1361599, 1364087, 1366583, 1369089, 1371599, 1374087, 1376583, 1379089, 1381599, 1384087, 1386583, 1389089, 1391599, 1394087, 1396583, 1399089, 1401599, 1404087, 1406583, 1409089, 1411599, 1414087, 1416583, 1419089, 1421599, 1424087, 1426583, 1429089, 1431599, 1434087, 1436583, 1439089, 1441599, 1444087, 1446583, 1449089, 1451599, 1454087, 1456583, 1459089, 1461599, 1464087, 1466583, 1469089, 1471599, 1474087, 1476583, 1479089, 1481599, 1484087, 1486583, 1489089, 1491599, 1494087, 1496583, 1499089, 1501599, 1504087, 1506583, 1509089, 1511599, 1514087, 1516583, 1519089, 1521599, 1524087, 1526583, 1529089, 1531599, 1534087, 1536583, 1539089, 1541599, 1544087, 1546583, 1549089, 1551599, 1554087, 1556583, 1559089, 1561599, 1564087, 1566583, 1569089, 1571599, 1574087, 1576583, 1579089, 1581599, 1584087, 1586583, 1589089, 1591599, 1594087, 1596583, 1599089, 1601599, 1604087, 1606583, 1609089, 1611599, 1614087, 1616583, 1619089, 1621599, 1624087, 1626583, 1629089, 1631599, 1634087, 1636583, 1639089, 1641599, 1644087, 1646583, 1649089, 1651599, 1654087, 1656583, 1659089, 1661599, 1664087, 1666583, 1669089, 1671599, 1674087, 1676583, 1679089, 1681599, 1684087, 1686583, 1689089, 1691599, 1694087, 1696583, 1699089, 1701599, 1704087, 1706583, 1709089, 1711599, 1714087, 1716583, 1719089, 1721599, 1724087, 1726583, 1729089, 1731599, 1734087, 1736583, 1739089, 1741599, 1744087, 1746583, 1749089, 1751599, 1754087, 1756583, 1759089, 1761599, 1764087, 1766583, 1769089, 1771599, 1774087, 1776583, 1779089, 1781599, 1784087, 1786583, 1789089, 1791599, 1794087, 1796583, 1799089, 1801599, 1804087, 1806583, 1809089, 1811599, 1814087, 1816583, 1819089, 1821599, 1824087, 1826583, 1829089, 1831599, 1834087, 1836583, 1839089, 1841599, 1844087, 1846583, 1849089, 1851599, 1854087, 1856583, 1859089, 1861599, 1864087, 1866583, 1869089, 1871599, 1874087, 1876583, 1879089, 1881599, 1884087, 1886583, 1889089, 1891599, 1894087, 1896583, 1899089, 1901599, 1904087, 1906583, 1909089, 1911599, 1914087, 1916583, 1919089, 1921599, 1924087, 1926583, 1929089, 1931599, 1934087, 1936583, 1939089, 1941599, 1944087, 1946583, 1949089, 1951599, 1954087, 1956583, 1959089, 1961599, 1964087, 1966583, 1969089, 1971599, 1974087, 1976583, 1979089, 1981599, 1984087, 1986583, 1989089, 1991599, 1994087, 1996583, 1999089, 2001599, 2004087, 2006583, 2009089, 2011599, 2014087, 2016583, 2019089, 2021599, 2024087, 2026583, 2029089, 2031599, 2034087, 2036583, 2039089, 2041599, 2044087, 2046583, 2049089, 2051599, 2054087, 2056583, 2059089, 2061599, 2064087, 2066583, 2069089, 2071599, 2074087, 2076583, 2079089, 2081599, 2084087, 2086583, 2089089, 2091599, 2094087, 2096583, 2099089, 2101599, 2104087, 2106583, 2109089, 2111599, 2114087, 2116583, 2119089, 2121599, 2124087, 2126583, 2129089, 2131599, 2134087, 2136583, 2139089, 2141599, 2144087, 2146583, 2149089, 2151599, 2154087, 2156583, 2159089, 2161599, 2164087, 2166583, 2169089, 2171599, 2174087, 2176583, 2179089, 2181599, 2184087, 2186583, 2189089, 2191599, 2194087, 2196583, 2199089, 2201599, 2204087, 2206583, 2209089, 2211599, 2214087, 2216583, 2219089, 2221599, 2224087, 2226583, 2229089, 2231599, 2234087, 2236583, 2239089, 2241599, 2244087, 2246583, 2249089, 2251599, 2254087, 2256583, 2259089, 2261599, 2264087, 2266583, 2269089, 2271599, 2274087, 2276583, 2279089, 2281599, 2284087, 2286583, 2289089, 2291599, 2294087, 2296583, 2299089, 2301599, 2304087, 2306583, 2309089, 2311599, 2314087, 2316583, 2319089, 2321599, 2324087, 2326583, 2329089, 2331599, 2334087, 2336583, 2339089, 2341599, 2344087, 2346583, 2349089, 2351599, 2354087, 2356583, 2359089, 2361599, 2364087, 2366583, 2369089, 2371599, 2374087, 2376583, 2379089, 2381599, 2384087, 2386583, 2389089, 2391599, 2394087, 2396583, 2399089, 2401599, 2404087, 2406583, 2409089, 2411599, 2414087, 2416583, 2419089, 2421599, 2424087, 2426583, 2429089, 2431599, 2434087, 2436583, 2439089, 2441599, 2444087, 2446583, 2449089, 2451599, 2454087, 2456583, 2459089, 2461599, 2464087, 2466583, 2469089, 2471599, 2474087, 2476583, 2479089, 2481599, 2484087, 2486583, 2489089, 2491599, 2494087, 2496583, 2499089, 2501599, 2504087, 2506583, 2509089, 2511599, 2514087, 2516583, 2519089, 2521599, 2524087, 2526583, 2529089, 2531599, 2534087, 2536583, 2539089, 2541599, 2544087, 2546583, 2549089, 2551599, 2554087, 2556583, 2559089, 2561599, 2564087, 2566583, 2569089, 2571599, 2574087, 2576583, 2579089, 2581599, 2584087, 2586583, 2589089, 2591599, 2594087, 2596583, 2599089, 2601599, 2604087, 2606583, 2609089, 2611599, 2614087, 2616583, 2619089, 2621599, 2624087, 2626583, 2629089, 2631599, 2634087, 2636583, 2639089, 2641599, 2644087, 2646583, 2649089, 2651599, 2654087, 2656583, 2659089, 2661599, 2664087, 2666583, 2669089, 2671599, 2674087, 2676583, 2679089, 2681599, 2684087, 2686583, 2689089, 2691599, 2694087, 2696583, 2699089, 2701599, 2704087, 2706583, 2709089, 2711599, 2714087, 27165
```