

For the purpose of this project, I will utilize a proprietary database retrieved from our Human Resources Information System (HRIS) platform, ADP. This comprehensive database includes complete records of payroll transactions from the period of January 2021 through May 2023. The structure of this database is organized with rows representing individual employees, and columns denoting various categories such as payroll dates, amounts, and the respective companies the employees were affiliated with at the time of payroll issuance. As noted in the data dictionary below.

My primary objective is to analyze this data by employing statistical methods and data analysis techniques. I will thoroughly cleanse and preprocess the data to ensure its optimal state for the subsequent predictive modeling.

Based on the insights drawn from the initial analysis, I will proceed to apply the most suitable machine learning model that aligns best with my objective - forecasting future payroll expenses. The intent is to develop a robust model that is capable of predicting payroll outflows, thereby assuring the organization in making more informed financial planning decisions.

In [24]:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import sklearn
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_absolute_error, r2_score
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.cluster import KMeans
from sklearn.decomposition import PCA
from sklearn.preprocessing import LabelEncoder
from sklearn.compose import ColumnTransformer
from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
from xgboost import XGBRegressor
import warnings
warnings.filterwarnings("ignore")
```

Data Dictionary

Column Name	Description
Company	The name or identifier of the company that the employee belongs to.
Employee ID	A unique identifier for each employee in the company.
Unique ID	Another unique identifier, used for specific purposes within the dataset.
Title	The job title or position of the employee within the company.
Job Class	The job classification or category of the employee's position.
Address State	The state where the employee's address is located.
Hire Date	The date when the employee was hired by the company.
Tenure	The length of time (in years) the employee has been working in the company.
Gender	The gender of the employee.
Department At Payroll Process	The department to which the employee belongs at the time of payroll processing.
Payroll Code	A code related to the employee's payroll transactions.
G L Account Number	The account number used in the company's General Ledger for financial payroll transactions.
G L Department	The department in the company associated with financial transactions.
Post Date	The date when a transaction or entry was posted.
Description	A description of the check date.
Dollars	The amount of money involved in the transaction.
Section Type	The type or category of the payroll code.
Month	The month in which a particular event or transaction occurred.
Year	The year in which a particular event or transaction occurred.

In [3]:

```
# Loading the data set and viewing dataframe.
df = pd.read_excel(r"C:\Users\Myoung\OneDrive - Bankers Financial Corporation\Desktop\Bellevue\DC8630_T301_221_231_232_233_234_235_236_237_238_239_240_241_242_243_244_245_246_247_248_249_250_251_252_253_254_255_256_257_258_259_260_261_262_263_264_265_266_267_268_269_270_271_272_273_274_275_276_277_278_279_280_281_282_283_284_285_286_287_288_289_290_291_292_293_294_295_296_297_298_299_300_301_302_303_304_305_306_307_308_309_310_311_312_313_314_315_316_317_318_319_320_321_322_323_324_325_326_327_328_329_330_331_332_333_334_335_336_337_338_339_340_341_342_343_344_345_346_347_348_349_350_351_352_353_354_355_356_357_358_359_360_361_362_363_364_365_366_367_368_369_370_371_372_373_374_375_376_377_378_379_380_381_382_383_384_385_386_387_388_389_390_391_392_393_394_395_396_397_398_399_400_401_402_403_404_405_406_407_408_409_410_411_412_413_414_415_416_417_418_419_420_421_422_423_424_425_426_427_428_429_430_431_432_433_434_435_436_437_438_439_440_441_442_443_444_445_446_447_448_449_450_451_452_453_454_455_456_457_458_459_460_461_462_463_464_465_466_467_468_469_470_471_472_473_474_475_476_477_478_479_480_481_482_483_484_485_486_487_488_489_490_491_492_493_494_495_496_497_498_499_500_501_502_503_504_505_506_507_508_509_510_511_512_513_514_515_516_517_518_519_520_521_522_523_524_525_526_527_528_529_530_531_532_533_534_535_536_537_538_539_540_541_542_543_544_545_546_547_548_549_550_551_552_553_554_555_556_557_558_559_560_561_562_563_564_565_566_567_568_569_570_571_572_573_574_575_576_577_578_579_580_581_582_583_584_585_586_587_588_589_590_591_592_593_594_595_596_597_598_599_600_601_602_603_604_605_606_607_608_609_610_611_612_613_614_615_616_617_618_619_620_621_622_623_624_625_626_627_628_629_630_631_632_633_634_635_636_637_638_639_640_641_642_643_644_645_646_647_648_649_650_651_652_653_654_655_656_657_658_659_660_661_662_663_664_665_666_667_668_669_670_671_672_673_674_675_676_677_678_679_680_681_682_683_684_685_686_687_688_689_690_691_692_693_694_695_696_697_698_699_700_701_702_703_704_705_706_707_708_709_710_711_712_713_714_715_716_717_718_719_720_721_722_723_724_725_726_727_728_729_730_731_732_733_734_735_736_737_738_739_740_741_742_743_744_745_746_747_748_749_750_751_752_753_754_755_756_757_758_759_760_761_762_763_764_765_766_767_768_769_770_771_772_773_774_775_776_777_778_779_780_781_782_783_784_785_786_787_788_789_790_791_792_793_794_795_796_797_798_799_800_801_802_803_804_805_806_807_808_809_810_811_812_813_814_815_816_817_818_819_820_821_822_823_824_825_826_827_828_829_830_831_832_833_834_835_836_837_838_839_840_841_842_843_844_845_846_847_848_849_850_851_852_853_854_855_856_857_858_859_860_861_862_863_864_865_866_867_868_869_870_871_872_873_874_875_876_877_878_879_880_881_882_883_884_885_886_887_888_889_890_891_892_893_894_895_896_897_898_899_900_901_902_903_904_905_906_907_908_909_910_911_912_913_914_915_916_917_918_919_920_921_922_923_924_925_926_927_928_929_930_931_932_933_934_935_936_937_938_939_940_941_942_943_944_945_946_947_948_949_950_951_952_953_954_955_956_957_958_959_960_961_962_963_964_965_966_967_968_969_970_971_972_973_974_975_976_977_978_979_980_981_982_983_984_985_986_987_988_989_990_991_992_993_994_995_996_997_998_999_1000_1001_1002_1003_1004_1005_1006_1007_1008_1009_1010_1011_1012_1013_1014_1015_1016_1017_1018_1019_1020_1021_1022_1023_1024_1025_1026_1027_1028_1029_1030_1031_1032_1033_1034_1035_1036_1037_1038_1039_1040_1041_1042_1043_1044_1045_1046_1047_1048_1049_1050_1051_1052_1053_1054_1055_1056_1057_1058_1059_1060_1061_1062_1063_1064_1065_1066_1067_1068_1069_1070_1071_1072_1073_1074_1075_1076_1077_1078_1079_1080_1081_1082_1083_1084_1085_1086_1087_1088_1089_1090_1091_1092_1093_1094_1095_1096_1097_1098_1099_1100_1101_1102_1103_1104_1105_1106_1107_1108_1109_1110_1111_1112_1113_1114_1115_1116_1117_1118_1119_1120_1121_1122_1123_1124_1125_1126_1127_1128_1129_1130_1131_1132_1133_1134_1135_1136_1137_1138_1139_1140_1141_1142_1143_1144_1145_1146_1147_1148_1149_1150_1151_1152_1153_1154_1155_1156_1157_1158_1159_1160_1161_1162_1163_1164_1165_1166_1167_1168_1169_1170_1171_1172_1173_1174_1175_1176_1177_1178_1179_1180_1181_1182_1183_1184_1185_1186_1187_1188_1189_1190_1191_1192_1193_1194_1195_1196_1197_1198_1199_1200_1201_1202_1203_1204_1205_1206_1207_1208_1209_1210_1211_1212_1213_1214_1215_1216_1217_1218_1219_1220_1221_1222_1223_1224_1225_1226_1227_1228_1229_1230_1231_1232_1233_1234_1235_1236_1237_1238_1239_1240_1241_1242_1243_1244_1245_1246_1247_1248_1249_1250_1251_1252_1253_1254_1255_1256_1257_1258_1259_1260_1261_1262_1263_1264_1265_1266_1267_1268_1269_1270_1271_1272_1273_1274_1275_1276_1277_1278_1279_1280_1281_1282_1283_1284_1285_1286_1287_1288_1289_1290_1291_1292_1293_1294_1295_1296_1297_1298_1299_1300_1301_1302_1303_1304_1305_1306_1307_1308_1309_1310_1311_1312_1313_1314_1315_1316_1317_1318_1319_1320_1321_1322_1323_1324_1325_1326_1327_1328_1329_1330_1331_1332_1333_1334_1335_1336_1337_1338_1339_1340_1341_1342_1343_1344_1345_1346_1347_1348_1349_1350_1351_1352_1353_1354_1355_1356_1357_1358_1359_1360_1361_1362_1363_1364_1365_1366_1367_1368_1369_1370_1371_1372_1373_1374_1375_1376_1377_1378_1379_1380_1381_1382_1383_1384_1385_1386_1387_1388_1389_1390_1391_1392_1393_1394_1395_1396_1397_1398_1399_1400_1401_1402_1403_1404_1405_1406_1407_1408_1409_1410_1411_1412_1413_1414_1415_1416_1417_1418_1419_1420_1421_1422_1423_1424_1425_1426_1427_1428_1429_1430_1431_1432_1433_1434_1435_1436_1437_1438_1439_1440_1441_1442_1443_1444_1445_1446_1447_1448_1449_1450_1451_1452_1453_1454_1455_1456_1457_1458_1459_1460_1461_1462_1463_1464_1465_1466_1467_1468_1469_1470_1471_1472_1473_1474_1475_1476_1477_1478_1479_1480_1481_1482_1483_1484_1485_1486_1487_1488_1489_1490_1491_1492_1493_1494_1495_1496_1497_1498_1499_1500_1501_1502_1503_1504_1505_1506_1507_1508_1509_1510_1511_1512_1513_1514_1515_1516_1517_1518_1519_1520_1521_1522_1523_1524_1525_1526_1527_1528_1529_1530_1531_1532_1533_1534_1535_1536_1537_1538_1539_1540_1541_1542_1543_1544_1545_1546_1547_1548_1549_1550_1551_1552_1553_1554_1555_1556_1557_1558_1559_1560_1561_1562_1563_1564_1565_1566_1567_1568_1569_1570_1571_1572_1573_1574_1575_1576_1577_1578_1579_1580_1581_1582_1583_1584_1585_1586_1587_1588_1589_1590_1591_1592_1593_1594_1595_1596_1597_1598_1599_1600_1601_1602_1603_1604_1605_1606_1607_1608_1609_1610_1611_1612_1613_1614_1615_1616_1617_1618_1619_1620_1621_1622_1623_1624_1625_1626_1627_1628_1629_1630_1631_1632_1633_1634_1635_1636_1637_1638_1639_1640_1641_1642_1643_1644_1645_1646_1647_1648_1649_1650_1651_1652_1653_1654_1655_1656_1657_1658_1659_1660_1661_1662_1663_1664_1665_1666_1667_1668_1669_1670_1671_1672_1673_1674_1675_1676_1677_1678_1679_1680_1681_1682_1683_1684_1685_1686_1687_1688_1689_1690_1691_1692_1693_1694_1695_1696_1697_1698_1699_1700_1701_1702_1703_1704_1705_1706_1707_1708_1709_1710_1711_1712_1713_1714_1715_1716_1717_1718_1719_1720_1721_1722_1723_1724_1725_1726_1727_1728_1729_1730_1731_1732_1733_1734_1735_1736_1737_1738_1739_1740_1741_1742_1743_1744_1745_1746_1747_1748_1749_1750_1751_1752_1753_1754_1755_1756_1757_1758_1759_1760_1761_1762_1763_1764_1765_1766_1767_1768_1769_1770_1771_1772_1773_1774_1775_1776_1777_1778_1779_1780_1781_1782_1783_1784_1785_1786_1787_1788_1789_1790_1791_1792_1793_1794_1795_1796_1797_1798_1799_1800_1801_1802_1803_1804_1805_1806_1807_1808_1809_1810_1811_1812_1813_1814_1815_1816_1817_1818_1819_1820_1821_1822_1823_1824_1825_1826_1827_1828_1829_1830_1831_1832_1833_1834_1835_1836_1837_1838_1839_1840_1841_1842_1843_1844_1845_1846_1847_1848_1849_1850_1851_1852_1853_1854_1855_1856_1857_1858_1859_1860_1861_1862_1863_1864_1865_1866_1867_1868_1869_1870_1871_1872_1873_1874_1875_1876_1877_1878_1879_1880_1881_1882_1883_1884_1885_1886_1887_1888_1889_1890_1891_1892_1893_1894_1895_1896_1897_1898_1899_1900_1901_1902_1903_1904_1905_1906_1907_1908_1909_1910_1911_1912_1913_1914_1915_1916_1917_1918_1919_1920_1921_1922_1923_1924_1925_1926_1927_1928_1929_1930_1931_1932_1933_1934_1935_1936_1937_1938_1939_1940_1941_1942_1943_1944_1945_1946_1947_1948_1949_1950_1951_1952_1953_1954_1955_1956_1957_1958_1959_1960_1961_1962_1963_1964_1965_1966_1967_1968_1969_1970_1971_1972_1973_1974_1975_1976_1977_1978_1979_1980_1981_1982_1983_1984_1985_1986_1987_1988_1989_1990_1991_1992_1993_1994_1995_1996_1997_1998_1999_2000_2001_2002_2003_2004_2005_2006_2007_2008_2009_2010_2011_2012_2013_2014_2015_2016_2017_2018_2019_2020_2021_2022_2023_2024_2025_2026_2027_2028_2029_2030_2031_2032_2033_2034_2035_2036_2037_2038_2039_2040_2041_2042_2043_2044_2045_2046_2047_2048_2049_2050_2051_2052_2053_2054_2055_2056_2057_2058_2059_2060_2061_2062_2063_2064_2065_2066_2067_2068_2069_2070_2071_2072_2073_2074_2075_2076_2077_2078_2079_2080_2081_2082_2083_2084_2085_2086_2087_2088_2089_2090_2091_2092_2093_2094_2095_2096_2097_2098_2099_2100_2101_2102_2103_2104_2105_2106_2107_2108_2109_2110_2111_2112_2113_2114_2115_2116_2117_2118_2119_2120_2121_2122_2123_2124_2125_2126_2127_2128_2129_2130_2131_2132_2133_2134_2135_2136_2137_2138_2139_2140_2141_2142_2143_2144_2145_2146_2147_2148_2149_2150_2151_2152_2153_2154_2155_2156_2157_2158_2159_2160_2161_2162_2163_2164_2165_2166_2167_2168_2169_2170_2171_2172_2173_2174_2175_2176_2177_2178_2179_2180_2181_2182_2183_2184_2185_2186_2187_2188_2189_2190_2191_2192_2193_2194_2195_2196_2197_2198_2199_2200_2201_2202_2203_2204_2205_2206_2207_2208_2209_2210_2211_2212_2213_2214_2215_2216_2217_2218_2219_2220_2221_2222_2223_2224_2225_2226_2227_2228_2229_2230_2231_2232_2233_2234_2235_2236_2237_2238_2239_2240_2241_2242_2243_2244_2245_2246_2247_2248_2249_2250_2251_2252_2253_2254_2255_2256_2257_2258_2259_2260_2261_2262_2263_2264_2265_2266_2267_2268_2269_2270_2271_2272_2273_2274_2275_2276_2277_2278_2279_2280_2281_2282_2283_2284_2285_2286_2287_2288_2289_2290_2291_2292_2293_2294_2295_2296_2297_2298_2299_2300_2301_2302_2303_2304_2305_2306_2307_2308_2309_2310_2311_2312_2313_2314_2315_2316_2317_2318_2319_2320_2321_2322_2323_2324_2325_2326_2327_2328_2329_2330_2331_2332_2333_2334_2335_2336_2337_2338_2339_2340_2341_2342_2343_2344_2345_2346_2347_2348_2349_2350_2351_2352_2353_2354_2355_2356_2357_2358_2359_2360_2361_2362_2363_2364_2365_2366_2367_2368_2369_2370_2371_2372_2373_2374_2375_2376_2377_2378_2379_2380_2381_2382_2383_2384_2385_2386_2387_2388_2389_2390_2391_2392_2393_2394_2395_2396_2397_2398_2399_2400_2401_2402_2403_2404_2405_2406_2407_2408_2409_2410_2411_2412_2413_2414_2415_2416_2417_2418_2419_2420_2421_2422_2423_2424_2425_2426_2427_2428_2429_2430_2431_2432_2433_2434_2435_2436_2437_2438_2439_2440_2441_2442_2443_2444_2445_2446_2447_2448_2449_2450_2451_2452_2453_2454_2455_2456_2457_2458_2459_2460_2461_2462_2463_2464_2465_2466_2467_2468_2469_2470_2471_2472_2473_2474_2475_2476_2477_2478_2479_2480_2481_2482_2483_2484_2485_2486_2487_2488_2489_2490_2491_2492_2493_2494_2495_2496_2497_2498_2499_2500_2501_2502_2503_2504_2505_2506_2507_2508_2509_2510_2511_2512_2513_2514_2515_2516_2517_2518_2519_2520_2521_2522_2523_2524_2525_2526_2527_2528_2529_2530_2531_2532_2533_2534_2535_2536_2537_2538_2539_2540_2541_2542_2543_2544_2545_2546_2547_2548_2549_2550_2551_2552_2553_2554_2555_2556_2557_2558_2559_2560_2561_2562_2563_2564_2565_2566_2567_2568_2569_2570_2571_2572_2573_2574_2575_2576_2577_2578_2579_2580_2581_2582_2583_2584_2585_2586_2587_2588_2589_2590_2591_2592_2593_2594_2595_2596_2597_2598_2599_2600_2601_2602_2603_2604_2605_2606_2607_2608_2609_2610_2611_2612_2613_2614_2615_2616_2617_2618_2619_2620_2621_2622_2623_2624_2625_2626_2627_2628_2629_2630_2631_2632_2633_2634_2635_2636_2637_2638_2639_2640_2641_2642_2643_2644_2645_2646_2647_2648_2649_2650_2651_2652_2
```



```
[34]: # 6c. Model Deployment
# Saving the trained model and using it to predict salaries for new data
# Example: Predicting monthly salaries for dataset

from sklearn.preprocessing import LabelEncoder

# Create a copy of the DataFrame to avoid modifying the original data
data_filtered = df.drop.copy()

# Filter the data to include only the rows with Month = 5 and Year = 2023 for most recent actuals
data_filtered = data_filtered[(data_filtered['Month'] == 5) & (data_filtered['Year'] == 2023)]

# Extract the relevant columns for analysis
features = ['Title', 'Job Class', 'Address State', 'Tenure', 'Gender',
            'Department At Payroll Process',
            'G L Account Number', 'Dollars', 'Month', 'Year']
data_filtered = data_filtered[features]

# Encode categorical variables using LabelEncoder
categorical_features = ['Title', 'Job Class', 'Address State', 'Gender',
                        'Department At Payroll Process',
                        'G L Account Number', 'Month', 'Year']

encoder = LabelEncoder()

# Collect all unique values for each categorical feature
unique_values = {}
for feature in categorical_features:
    unique_values[feature] = data_filtered[feature].unique()

# Encode the categorical variables in the filtered data
for feature in categorical_features:
    encoder.classes_ = unique_values[feature]
    data_filtered[feature] = encoder.transform(data_filtered[feature])

# Split the filtered data into features (X) and target variable (y)
X_filtered = data_filtered.drop(['Dollars', 'axis=1'])
y_filtered = data_filtered['Dollars']

# Predict the salaries for the filtered data
predicted_salaries = model.predict(X_filtered)

# Calculate the total monthly salary for those paid within the specified month and year
total_monthly_salary = predicted_salaries.sum()

print("Total Predicted Monthly Salary based upon May 2023:", "${:,2f}".format(total_monthly_salary))
print("Total Predicted Annual Salary based upon May 2023:", "${:,2f}".format(total_monthly_salary*12))

Total Predicted Monthly Salary based upon May 2023: $2,637,237.50
Total Predicted Annual Salary based upon May 2023: $31,646,850.00
```