# INTRO ML - M1 ISFA 2024- 2025

## ∨ Import

```python
1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 from io import StringIO
6 import plotly.express as px
7 import plotly.graph_objects as go
8 from sklearn.model_selection import train_test_split
9 from sklearn.preprocessing import OneHotEncoder, StandardScaler
10 from sklearn.compose import ColumnTransformer
11 from sklearn.linear_model import LinearRegression
12 from sklearn.tree import DecisionTreeRegressor, plot_tree
13 from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
14 from sklearn.pipeline import Pipeline
15 from sklearn.metrics import mean_squared_error
16 import statsmodels.api as sm
17 from sklearn.inspection import partial_dependence, PartialDependenceDisplay
18 import time
19
20 # For missing values visualization
21 !pip install missingno
22 import missingno as msno
23
24 # For outlier detection
25 !pip install pyod
26 from pyod.models.iforest import IForest
27
28 # For QQ plot data calculation
29 from scipy.stats import probplot
30
31 # -------------------------
32 # Step 0: Create and Load Dataset
33 # -------------------------
34 csv_data = """Gender,Age,Salary,Height,City
35 Male,25,50000,175,New York
36 Female,30,60000,162,Los Angeles
37 Male,45,,180,Chicago
38 Female,35,70000,158,New York
39 Male,50,120000,185,San Francisco
40 Female,29,40000,,Los Angeles
41 Male,33,65000,172,Chicago
42 Female,44,80000,165,
43 Male,55,130000,188,New York
44 Female,39,90000,163,San Francisco
45 Male,,45000,170,Los Angeles
46 Female,27,48000,155,Chicago
47 Male,60,150000,190,New York
48 Female,22,38000,,San Francisco
49 Male,46,85000,178,Chicago
50 Female,,52000,160,Los Angeles
51 Male,48,92000,,New York
52 Female,41,75000,167,San Francisco
53 Male,34,68000,174,Chicago
54 Female,31,62000,159,New York
55 Male,28,52000,168,Los Angeles
56 Female,,59000,156,Chicago
57 Male,52,100000,182,New York
58 Female,26,43000,153,San Francisco
59 Male,54,,186,Chicago
60 Female,40,85000,161,Los Angeles
61 Male,47,90000,179,New York
62 Female,33,70000,164,San Francisco
63 Male,57,145000,,Chicago
64 Female,29,42000,157,New York
65 Male,42,75000,176,Los Angeles
66 Female,34,64000,160,Chicago
67 Male,,49000,171,New York
68 Female,36,79000,,Los Angeles
69 Male,51,110000,183,San Francisco
70 Female,25,,154,Chicago
71 Male,61,155000,192,New York
72 Female,23,40000,152,San Francisco
73 Male,49,88000,177,Chicago
74 Female,38,82000,164,Los Angeles
75 Male,,70000,173,New York
76 Female,32,61000,158,San Francisco
77 Male,44,87000,,Chicago
78 Female,37,78000,161,New York
79 Male,29,54000,169,Los Angeles
80 Female,28,60000,156,Chicago
81 Male,56,140000,189,New York
82 Female,24,39000,,San Francisco
83 Male,58,148000,191,Chicago
84 Female,21,35000,151,Los Angeles
85 Male,35,72000,175,New York
86 Female,33,68000,164,San Francisco
87 Male,43,82000,176,Chicago
88 Female,29,50000,159,New York
89 Male,41,76000,174,Los Angeles
90 Female,27,62000,,Chicago
91 Male,59,152000,193,New York
92 Female,20,36000,150,San Francisco
93 Male,46,89000,178,Chicago
94 Female,42,80000,163,Los Angeles
95 Male,30,56000,169,New York
```

```
96 Female,35,69000,159,San Francisco
97 """
98
99 df = pd.read_csv(StringIO(csv_data))
```

```
⇥ Requirement already satisfied: missingno in /usr/local/lib/python3.10/dist-packages (0.5.2)
  Requirement already satisfied: numpy in /usr/local/lib/python3.10/dist-packages (from missingno) (1.26.4)
  Requirement already satisfied: matplotlib in /usr/local/lib/python3.10/dist-packages (from missingno) (3.8.0)
  Requirement already satisfied: scipy in /usr/local/lib/python3.10/dist-packages (from missingno) (1.13.1)
  Requirement already satisfied: seaborn in /usr/local/lib/python3.10/dist-packages (from missingno) (0.13.2)
  Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->missingno) (1.3.1)
  Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.10/dist-packages (from matplotlib->missingno) (0.12.1)
  Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib->missingno) (4.55.2)
  Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->missingno) (1.4.7)
  Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib->missingno) (24.2)
  Requirement already satisfied: pillow>=6.2.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib->missingno) (11.0.0)
  Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->missingno) (3.2.0)
  Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.10/dist-packages (from matplotlib->missingno) (2.8.2)
  Requirement already satisfied: pandas>=1.2 in /usr/local/lib/python3.10/dist-packages (from seaborn->missingno) (2.2.2)
  Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas>=1.2->seaborn->missingno) (2024.2)
  Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.10/dist-packages (from pandas>=1.2->seaborn->missingno) (2024.2)
  Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.7->matplotlib->missingno) (1.17.0)
  Requirement already satisfied: pyod in /usr/local/lib/python3.10/dist-packages (2.0.2)
  Requirement already satisfied: joblib in /usr/local/lib/python3.10/dist-packages (from pyod) (1.4.2)
  Requirement already satisfied: matplotlib in /usr/local/lib/python3.10/dist-packages (from pyod) (3.8.0)
  Requirement already satisfied: numpy>=1.19 in /usr/local/lib/python3.10/dist-packages (from pyod) (1.26.4)
  Requirement already satisfied: numba>=0.51 in /usr/local/lib/python3.10/dist-packages (from pyod) (0.60.0)
  Requirement already satisfied: scipy>=1.5.1 in /usr/local/lib/python3.10/dist-packages (from pyod) (1.13.1)
  Requirement already satisfied: scikit-learn>=0.22.0 in /usr/local/lib/python3.10/dist-packages (from pyod) (1.5.2)
  Requirement already satisfied: llvmlite<0.44,>=0.43.0dev0 in /usr/local/lib/python3.10/dist-packages (from numba>=0.51->pyod) (0.43.0)
  Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.10/dist-packages (from scikit-learn>=0.22.0->pyod) (3.5.0)
  Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->pyod) (1.3.1)
  Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.10/dist-packages (from matplotlib->pyod) (0.12.1)
  Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib->pyod) (4.55.2)
  Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->pyod) (1.4.7)
  Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib->pyod) (24.2)
  Requirement already satisfied: pillow>=6.2.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib->pyod) (11.0.0)
  Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->pyod) (3.2.0)
  Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.10/dist-packages (from matplotlib->pyod) (2.8.2)
  Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.7->matplotlib->pyod) (1.17.0)
```
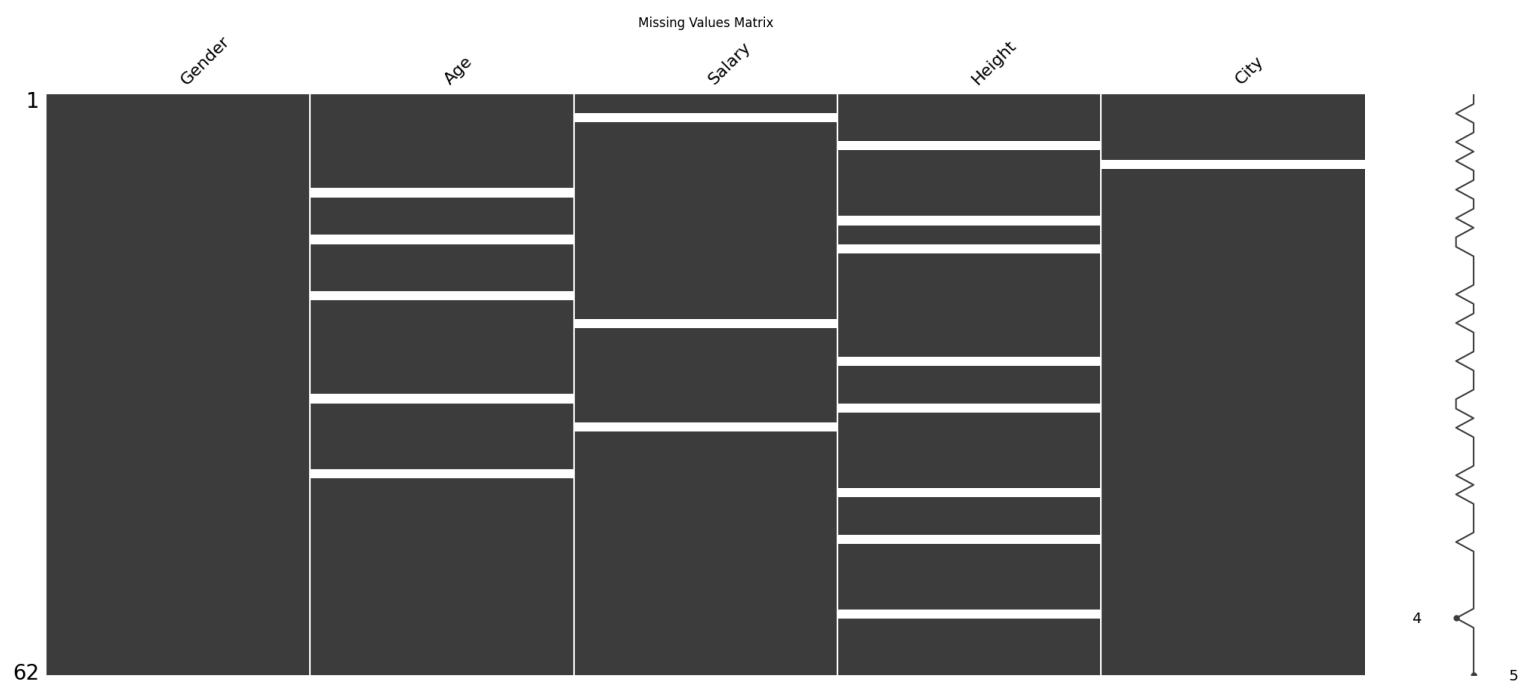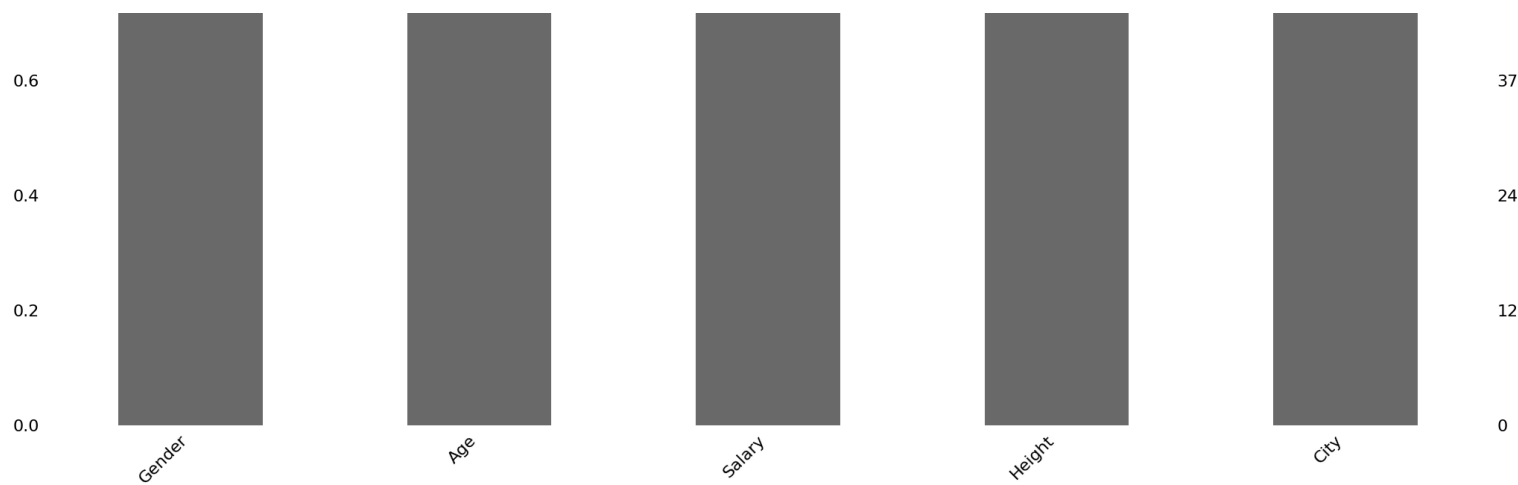
## ⌄ Step 1: Basic Info and Missing Value Evaluation

Using missingno to visualize missing data

```
1
2 print("Initial DataFrame shape:", df.shape)
3 print("\nColumn-wise Missing Values:\n", df.isnull().sum())
4
5 # Visualize missing data using missingno
6 msno.bar(df)
7 plt.title("Missing Values Bar Chart")
8 plt.show()
9
10 msno.matrix(df)
11 plt.title("Missing Values Matrix")
12 plt.show()
13
14 # Impute missing values:
15 num_cols = df.select_dtypes(include=[np.number]).columns
16 for col in num_cols:
17     if df[col].isnull().sum() > 0:
18         df[col].fillna(df[col].mean(), inplace=True)
19
20 cat_cols = df.select_dtypes(include=['object']).columns
21 for col in cat_cols:
22     if df[col].isnull().sum() > 0:
23         df[col].fillna(df[col].mode()[0], inplace=True)
24
25 print("\nDataFrame shape after missing value imputation:", df.shape)
26 print("Missing Values after imputation:\n", df.isnull().sum())
27
```

Missing Values Matrix



```
DataFrame shape after missing value imputation: (62, 5)
Missing Values after imputation:
 Gender    0
Age       0
Salary    0
```

## Step 2: Outlier Detection Using IQR and more

```
 1
 2 Q1 = df[num_cols].quantile(0.25)
 3 Q3 = df[num_cols].quantile(0.75)
 4 IQR = Q3 - Q1
 5 outlier_mask = ((df[num_cols] < (Q1 - 1.5 * IQR)) | (df[num_cols] > (Q3 + 1.5 * IQR)))
 6 outliers_present = outlier_mask.any(axis=1)
 7 print("\nNumber of Outliers Detected:", outliers_present.sum())
 8
 9 df_clean = df[~outliers_present]
10 print("DataFrame shape after outlier removal:", df_clean.shape)
11
```

```
Number of Outliers Detected: 6
DataFrame shape after outlier removal: (56, 5)
```

```
 1
 2 # Using an Isolation Forest for outlier detection on numeric features
 3 X_numeric = df[num_cols]
 4 clf = IForest(contamination=0.1, random_state=42)
 5 clf.fit(X_numeric)
 6 outlier_labels = clf.labels_  # 1 for outliers, 0 for inliers
 7
```
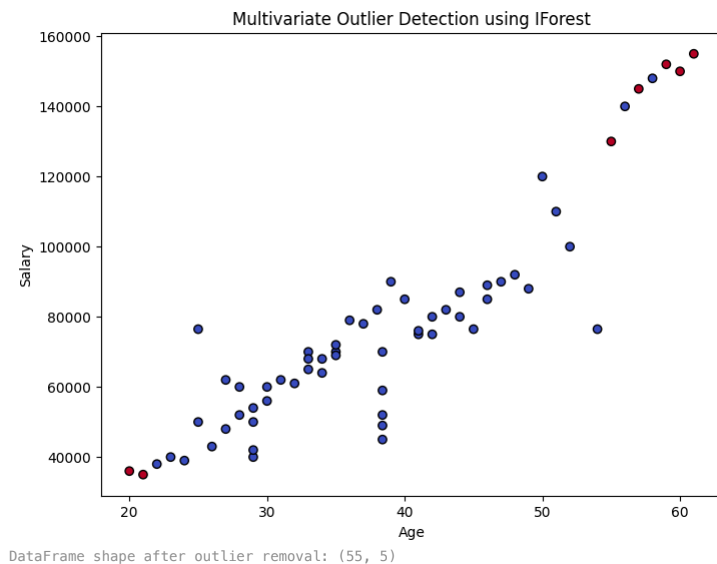
```
 8 # Count outliers
 9 outliers_present = (outlier_labels == 1)
10 print("\nNumber of Outliers Detected by PYOD IForest:", outliers_present.sum())
11
12 # Visualize outliers in a 2D scatter plot (choose two numeric features)
13 if len(num_cols) >= 2:
14     plt.figure(figsize=(8,6))
15     plt.scatter(X_numeric[num_cols[0]], X_numeric[num_cols[1]], c=outlier_labels, cmap='coolwarm', edgecolor='k')
16     plt.xlabel(num_cols[0])
17     plt.ylabel(num_cols[1])
18     plt.title("Multivariate Outlier Detection using IForest")
19     plt.show()
20
21 df_clean = df[~outliers_present]
22 print("DataFrame shape after outlier removal:", df_clean.shape)
23
```
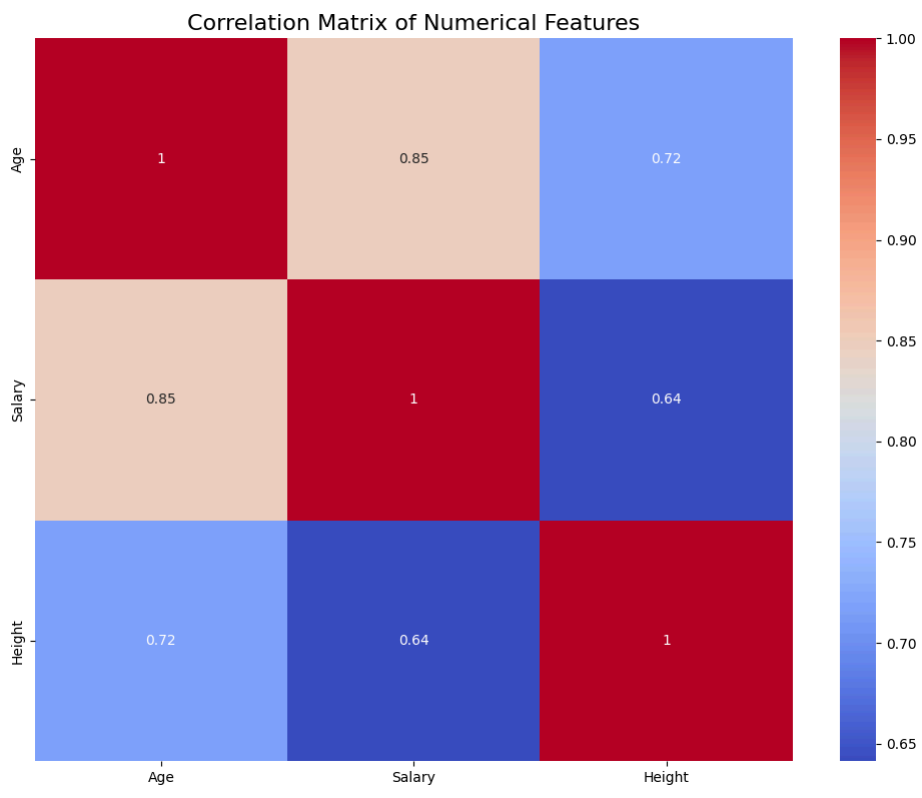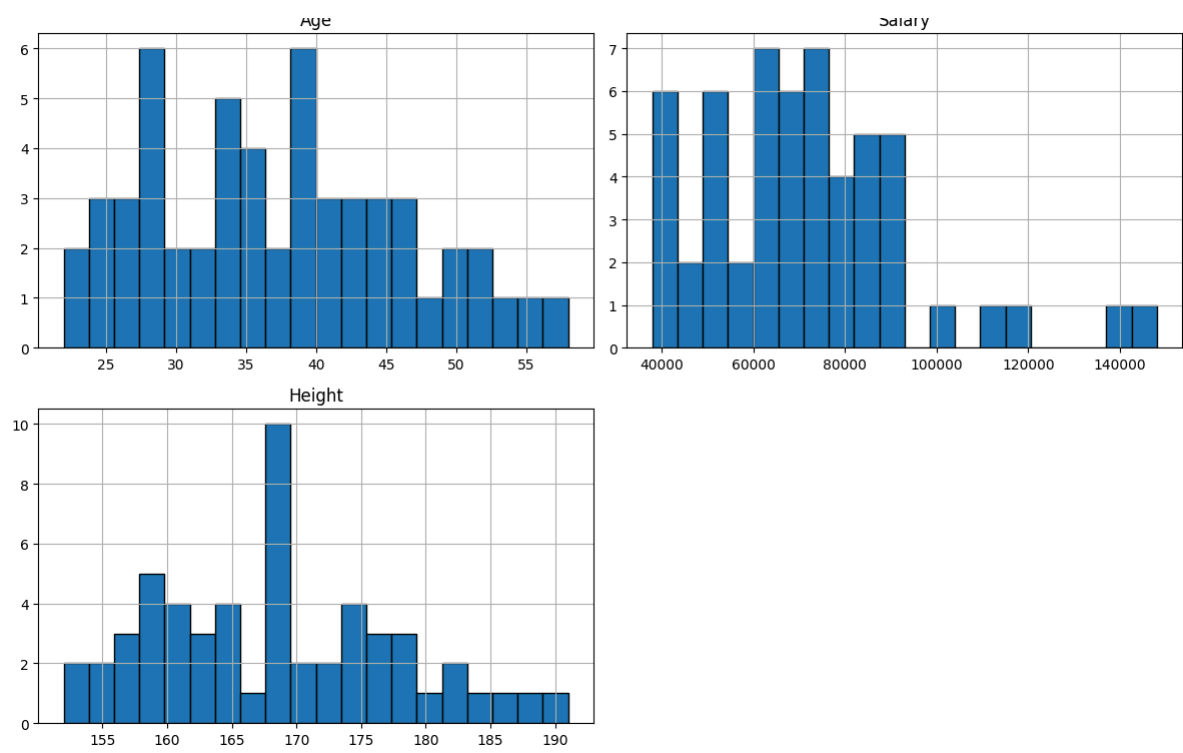
Number of Outliers Detected by PYOD IForest: 7



Multivariate Outlier Detection using IForest

DataFrame shape after outlier removal: (55, 5)

## ⌄ Step 3 to 5: Visulizations
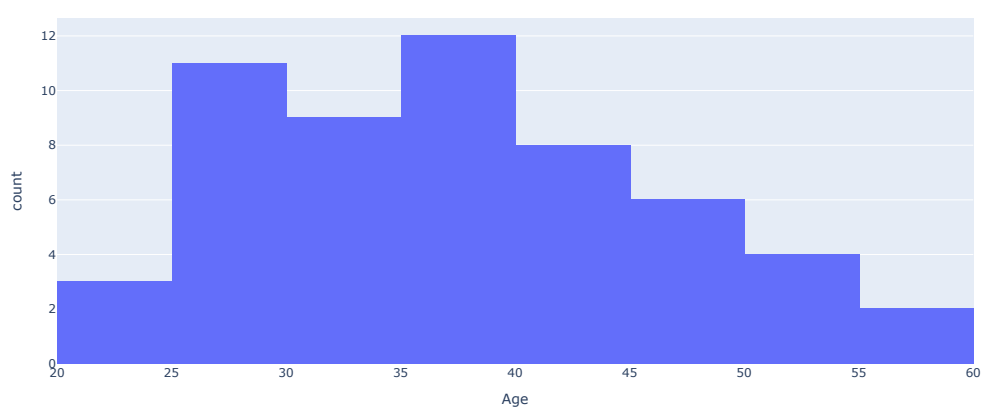
```
 1 # -------------------------
 2 # Step 3: Distribution Plots using Matplotlib/Seaborn (No Change)
 3 # -------------------------
 4 plt.figure(figsize=(12,8))
 5 df_clean[num_cols].hist(bins=20, edgecolor='black', figsize=(12,8))
 6 plt.suptitle("Histograms of Numerical Features", fontsize=16)
 7 plt.tight_layout()
 8 plt.show()
 9
10 # -------------------------
11 # Step 4: Correlation Analysis Using Seaborn (No Change)
12 # -------------------------
13 plt.figure(figsize=(10,8))
14 corr = df_clean[num_cols].corr()
15 sns.heatmap(corr, annot=True, cmap='coolwarm')
16 plt.title("Correlation Matrix of Numerical Features", fontsize=16)
17 plt.tight_layout()
18 plt.show()
19
20 # -------------------------
21 # Step 5: Interactive Visualization using Plotly (No Change)
22 # -------------------------
23 fig_age = px.histogram(df_clean, x="Age", nbins=10, title="Age Distribution (Plotly)")
24 fig_age.show()
25
26 fig_corr = go.Figure(data=go.Heatmap(
27     z=corr.values,
28     x=corr.columns,
29     y=corr.columns,
30     colorscale='RdBu',
31     zmin=-1, zmax=1))
32 fig_corr.update_layout(title="Correlation Matrix (Plotly)", xaxis_nticks=36)
33 fig_corr.show()
34
```
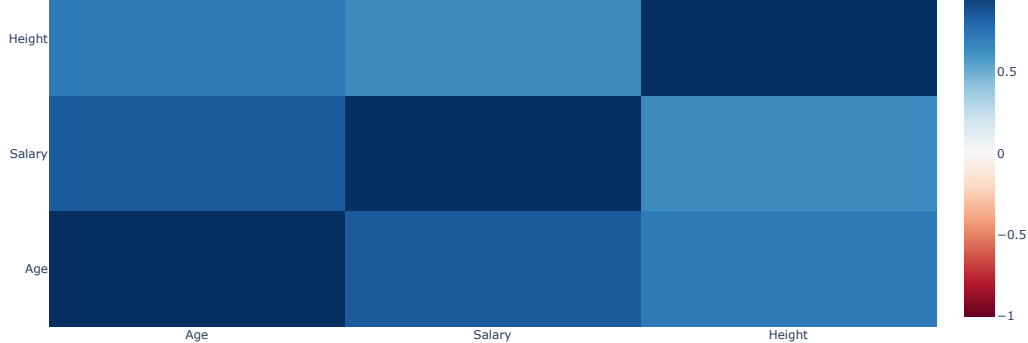
## Age

## Salary

## Height

## Correlation Matrix of Numerical Features

|        | Age  | Salary | Height |
|--------|------|--------|--------|
| Age    | 1    | 0.85   | 0.72   |
| Salary | 0.85 | 1      | 0.64   |
| Height | 0.72 | 0.64   | 1      |

## Age Distribution (Plotly)

## Correlation Matrix (Plotly)

## Step 6: Prepare Data for Supervised Machine Learning

```python
1
2  # For demonstration, let's predict Salary from the other features.
3  # We'll consider "Salary" as the target variable and the rest as features.
4  X = df_clean.drop("Salary", axis=1)
5  y = df_clean["Salary"]
6
7  numeric_features = X.select_dtypes(include=np.number).columns.tolist()
8  categorical_features = X.select_dtypes(include=['object']).columns.tolist()
9
10 preprocessor = ColumnTransformer(
11     transformers=[
12         ("num", StandardScaler(), numeric_features),
13         ("cat", OneHotEncoder(drop='first'), categorical_features)
14     ]
15 )
16
17 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
18
```

## Step 7: Extend Modeling - Compare Multiple Models using RMSE and QQ plots

Models: LinearRegression, DecisionTreeRegressor, RandomForestRegressor, GradientBoostingRegressor

```python
1
2
3  models = {
4      "Linear Regression": LinearRegression(),
5      "Decision Tree": DecisionTreeRegressor(random_state=42),
6      "Random Forest": RandomForestRegressor(random_state=42, n_estimators=500),
7      "Gradient Boosting": GradientBoostingRegressor(random_state=42, n_estimators=500)
8  }
9
10
11 results = {}
12 inference_times = {}
13
14 for model_name, model_obj in models.items():
15     pipe = Pipeline(steps=[
16         ('preprocessor', preprocessor),
17         ('regressor', model_obj)
18     ])
19     pipe.fit(X_train, y_train)
20
21     # Inference time measurement
22     start_time = time.time()
23     y_pred = pipe.predict(X_test)
24     end_time = time.time()
25     inference_time = end_time - start_time
26     inference_times[model_name] = inference_time
27
28     rmse = np.sqrt(mean_squared_error(y_test, y_pred))
29     results[model_name] = (pipe, rmse, y_pred)
30
31 # Print RMSE results
32 print("\nModel Performance (RMSE):")
33 for model_name, (pipe, rmse, y_pred) in results.items():
34     print(f"{model_name}: RMSE = {rmse:.2f}")
35
```

```
Model Performance (RMSE):
  Linear Regression: RMSE = 16652.39
  Decision Tree: RMSE = 13013.29
  Random Forest: RMSE = 11602.97
  Gradient Boosting: RMSE = 13201.14
```

```python
1
2  # Select best model based on RMSE
3  best_model_name = min(results, key=lambda x: results[x][1])
4  best_pipe, best_rmse, best_pred = results[best_model_name]
5  print(f"\nBest Model: {best_model_name} with RMSE = {best_rmse:.2f}")
6
```

```
Best Model: Random Forest with RMSE = 11602.97
```
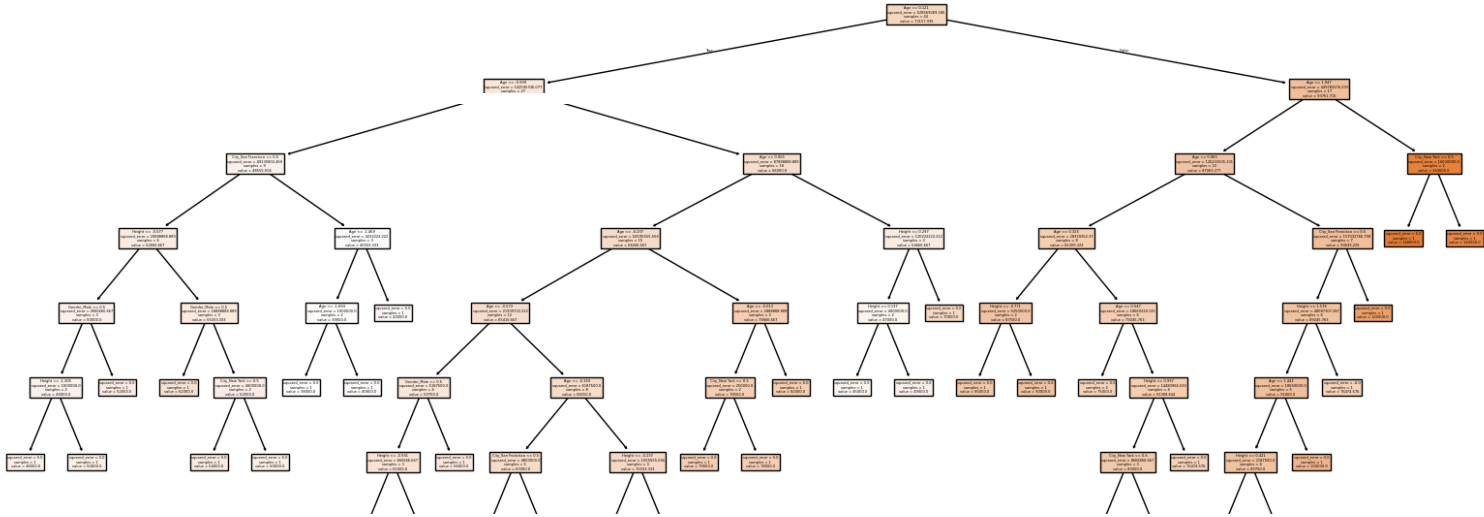
```python
1
2  # 7a: Tree representation of the regression tree
3  # Note: Only works for the Decision Tree model
4  tree_model = results["Decision Tree"][0].named_steps['regressor']
5  plt.figure(figsize=(20,10))
6  plot_tree(tree_model, feature_names=numeric_features+list(results["Decision Tree"][0].named_steps['preprocessor'].transformers_[1][1].get_feature_names_out(ca
7  plt.title("Decision Tree Regression Tree")
```

# Decision Tree Regression Tree



```
1
2 # 7b: QQ plots using plotly
3 def qqplot_plotly(residuals, model_name):
4     # Calculate theoretical quantiles and sample quantiles
5     qq = probplot(residuals, dist="norm")
6     theoretical_q = qq[0][0]
7     sample_q = qq[0][1]
8     slope, intercept = qq[1][0], qq[1][1]
9
10     fig = go.Figure()
11     fig.add_trace(go.Scatter(x=theoretical_q, y=sample_q, mode='markers', name='Data'))
12     # line of best fit
13     line_y = slope*theoretical_q + intercept
14     fig.add_trace(go.Scatter(x=theoretical_q, y=line_y, mode='lines', name='Best Fit Line'))
15     fig.update_layout(title=f"QQ Plot of Residuals – {model_name}",
16                       xaxis_title='Theoretical Quantiles',
17                       yaxis_title='Sample Quantiles')
18     fig.show()
19
20 for model_name, (pipe, rmse, y_pred) in results.items():
21     residuals = y_test – y_pred
22     qqplot_plotly(residuals, model_name)
23
```