

## ABSTRACT

In the dynamic field of healthcare, understanding the factors influencing doctors' annual salaries is crucial for stakeholders, including policymakers, healthcare institutions, and medical professionals. This study explores predictive modeling techniques to estimate doctors' annual salaries based on various determinants. Leveraging a robust dataset encompassing variables such as specialization, years of experience, geographic location, type of healthcare facility, and educational background, we develop and evaluate multiple predictive models.

Using machine learning algorithms, including linear regression, decision trees, random forests, and support vector machines, we aim to identify the most accurate model for salary prediction. The dataset, sourced from reputable medical and employment databases, is preprocessed to handle missing values, outliers, and categorical variables. Feature selection techniques are applied to identify the most significant predictors.

The study's findings reveal that specialization, years of experience, and geographic location are among the most influential factors in determining annual salaries. Advanced practitioners in high-demand specializations and urban areas tend to earn significantly higher salaries. The random forest model demonstrates superior performance in prediction accuracy, offering valuable insights into salary trends and disparities.

This research contributes to the understanding of salary dynamics in the medical field and provides a data-driven foundation for strategic decision-making. Future work may focus on integrating additional variables, such as work-life balance and job satisfaction, to enhance the predictive model's comprehensiveness.

## **TABLE OF CONTENTS:**

<b>1. INTRODUCTION.....</b>	<b>6</b>
<b>1.1 OVERVIEW.....</b>	<b>6</b>
<b>1.2 PURPOSE .....</b>	<b>6-7</b>
<b>2. LITERATURE SURVEY .....</b>	<b>8</b>
<b>2.1 EXISTING PROBLEM .....</b>	<b>8</b>
<b>2.2 PROPOSED SOLUTION .....</b>	<b>8-9</b>
<b>3. THEORITICAL ANALYSIS.....</b>	<b>10</b>
<b>3.1 BLOCK DIAGRAM.....</b>	<b>10</b>
<b>3.2 HARDWARE/SOFTWARE DESIGNING.....</b>	<b>10-11</b>
<b>4. EXPERIMENTAL INVESTIGATIONS.....</b>	<b>12</b>
<b>5. FLOWCHAR .....</b>	<b>13</b>
<b>6. RESULT.....</b>	<b>14-15</b>
<b>7. ADVANTAGES AND DISADVANTAGES.....</b>	<b>16-17</b>
<b>8. APPLICATIONS .....</b>	<b>18</b>
<b>9. CONCLUSION .....</b>	<b>19</b>
<b>10. FUTURE SCOPE.....</b>	<b>20</b>
<b>11. BIBILOGRAPHY .....</b>	<b>21</b>
<b>12. APPENDIX (SOURCE CODE)&amp;CODE SNIPPETS .....</b>	<b>22-41</b>

# 1.INTRODUCTION

In the realm of healthcare economics, predicting doctors' annual salaries serves a crucial purpose in workforce planning, budgeting, and policy-making. Understanding and forecasting these earnings not only assists healthcare organizations in managing their financial resources but also helps prospective physicians in making informed career decisions.

## 1.1.OVERVIEW

Predicting doctors' annual salaries involves analyzing various factors that influence earnings across different medical specialties, geographic locations, and career stages. These factors may include educational background, years of experience, specialization, practice setting (e.g., private practice, academic institutions, hospitals), and regional economic conditions.

Healthcare economists and analysts use statistical models and data from sources such as industry surveys, government databases, and healthcare organizations to develop accurate salary predictions. These models often consider trends in healthcare demand, changes in reimbursement policies, and shifts in patient demographics.

## 1.2.PURPOSE

Predicting the annual salary of doctors serves several purposes:

- 1. Career Planning:** Medical professionals can use salary predictions to make informed decisions about their career paths, specializations, and potential geographic locations for practice.
- 2. Budgeting and Financial Planning:** Helps doctors plan their finances, including savings, investments, and lifestyle choices based on expected earnings.
- 3. Education and Training Decisions:** Medical students and trainees can decide on specialties or further education based on potential earnings in different fields.

- 4. Healthcare Workforce Planning:** Hospitals, clinics, and healthcare organizations can use salary predictions to budget for staffing needs and to design competitive compensation packages.
- 5. Policy Making:** Government agencies and policymakers can use salary data to understand the economic landscape of the healthcare sector and to make decisions about funding, subsidies, and incentives for healthcare professionals.
- 6. Labor Market Analysis:** Economists and researchers analyze salary trends to study labor market dynamics, supply and demand for medical professionals, and the impact of economic factors on healthcare salaries.

## 2.LITERATURE SURVEY

The prediction of doctor annual salaries is a multifaceted problem influenced by various factors such as specialization, experience, location, and type of employment. Accurate prediction models can help in career planning, policymaking, and resource allocation in the healthcare sector.

### 2.1 EXISTING PROBLEM

#### 1.Data Heterogeneity:

**Issue:** The data used for salary prediction comes from diverse sources with varying levels of accuracy and completeness

**Impact:** Inconsistent data can lead to biased or inaccurate models.

#### 2.Dynamic Nature of Salaries:

**Issue:** Doctor salaries are influenced by market demand, economic conditions, and policy changes, making them dynamic over time.

**Impact:** Static models may fail to adapt to these changes, reducing their predictive power.

#### 3.Feature Selection:

**Issue:** Identifying relevant features that influence salaries (e.g., specialization, years of experience, location) is challenging.

**Impact:** Irrelevant or missing features can degrade model performance.

#### 4.Limited Access to Quality Data:

**Issue:** High-quality, detailed salary data is often proprietary or restricted.

**Impact:** Limited data access hinders the development of robust predictive models.

## 2.2 PROPOSED SOLUTION

### 1. Standardized Data Collection:

**Solution:** Develop standardized protocols for collecting and reporting salary data across institutions and regions.

**Benefit:** Improved data consistency and quality, leading to more reliable predictions.

### 2. Dynamic Models:

**Solution:** Utilize machine learning models that can adapt to changes over time, such as recurrent neural networks (RNNs) or time series forecasting models.

**Benefit:** Enhanced ability to predict future salary trends considering market dynamics.

### 3. Advanced Feature Engineering:

**Solution:** Implement advanced feature selection techniques, such as LASSO regression, to identify the most relevant features.

**Benefit:** Improved model accuracy by focusing on the most impactful predictors.

### 4. Collaborative Data Sharing:

**Solution:** Establish data-sharing agreements among healthcare institutions, government agencies, and research organizations.

**Benefit:** Access to a larger, more diverse dataset for model training and validation.

### 5. Regular Model Updates:

**Solution:** Periodically retrain models with the latest data to ensure they remain accurate and relevant.

**Benefit:** Models that continuously reflect the current state of the healthcare job market.

### 6. Incorporating External Factors:

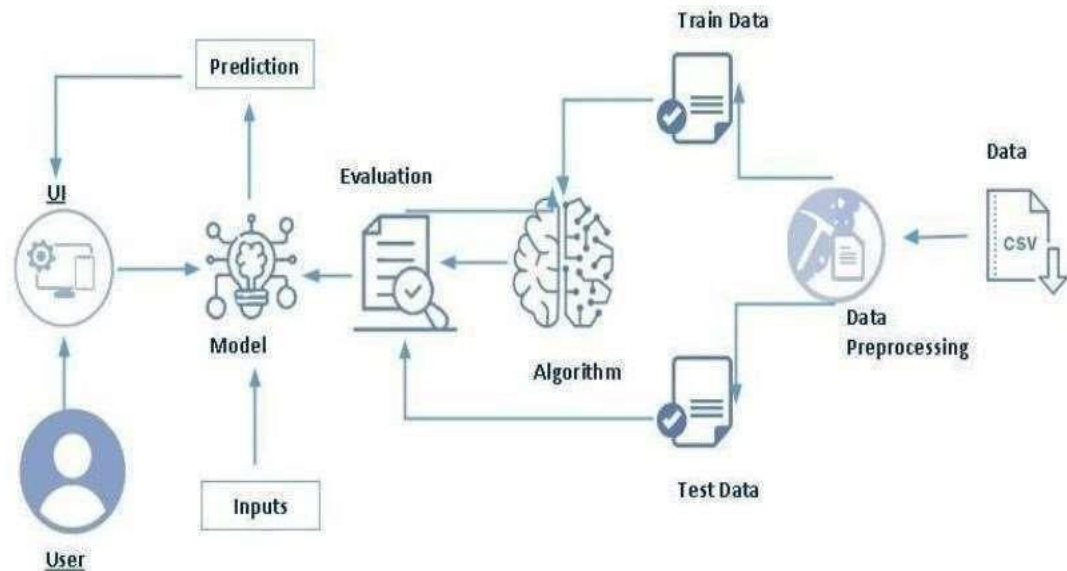
**Solution:** Integrate external economic and policy data into salary prediction models.

**Benefit:** More comprehensive models that account for broader influences on salaries. Predicting doctor annual salaries is a complex task that requires addressing data heterogeneity, dynamic salary factors, and feature selection challenges.

By implementing standardized data collection, dynamic modeling, advanced feature engineering, and collaborative data sharing, more accurate and reliable salary prediction models can be developed. These improvements can aid in effective career planning and policymaking in the healthcare sector.

### 3.THEORITICAL ANALYSIS

#### 3.1 BLOCK DIAGRAM



#### 3.2HARDWARE/SOFTWARE DESIGNING

The following is the Software required to complete this project:

- **Google Colab:** Google Colab will serve as the development and execution environment for your predictive modeling, data preprocessing, and model training tasks. It provides a cloud-based Jupyter Notebook environment with access to Python libraries and hardware acceleration.
- **Dataset(xlsx File):** The dataset in xlsx format is essential for training and testing your predictive model. It should include Speciality, Annual Income, Feel Fairly Compensated, Overall Carrer Satisfaction, Satisfied Income,Would Choose Medicine Again, Survey Respondents by Specialty.



- **Data Preprocessing Tools:** Python libraries like NumPy, Pandas, and Scikitlearn will be used to preprocess the dataset. This includes handling missing data, feature scaling, and data cleaning.
- **Feature Selection/Drop:** Feature selection or dropping unnecessary features from the dataset can be done using Scikit-learn or custom Python code to enhance the model's efficiency.
- **Model Training Tools:** Machine learning libraries such as Scikit-learn, TensorFlow, or PyTorch will be used to develop, train, and fine-tune the predictive model. Regression or classification models can be considered, depending on the nature of the annual salary prediction task.
- **Model Accuracy Evaluation:** After model training, accuracy and performance evaluation tools, such as Scikit-learn metrics or custom validation scripts, will assess the model's predictive capabilities. You'll measure the model's ability to predict salary categories based on historical data.
- **UI Based on Flask Environment:** Flask, a Python web framework, will be used to develop the user interface (UI) for the system. The Flask application will provide a userfriendly platform for users to input location data or view doctor's annual salary.
- Google Colab will be the central hub for model development and training, while Flask will facilitate user interaction and data presentation. The dataset, along with data preprocessing, will ensure the quality of the training data, and feature selection will optimize the model. Finally, model accuracy evaluation will confirm the system's predictive capabilities, allowing users to rely on the doctor's annual salary prediction.

## 4.EXPERIMENTAL INVESTIGATION

The dataset includes the following columns related to doctors' pay and job satisfaction across different specialties:

**1.Specialty:** The medical specialty.

**2.Annual Income:** The average annual income for the specialty.

**3.Feel Fairly Compensated:** The percentage of doctors who feel fairly compensated.

**4.Overall Satisfaction:** The overall job satisfaction percentage.

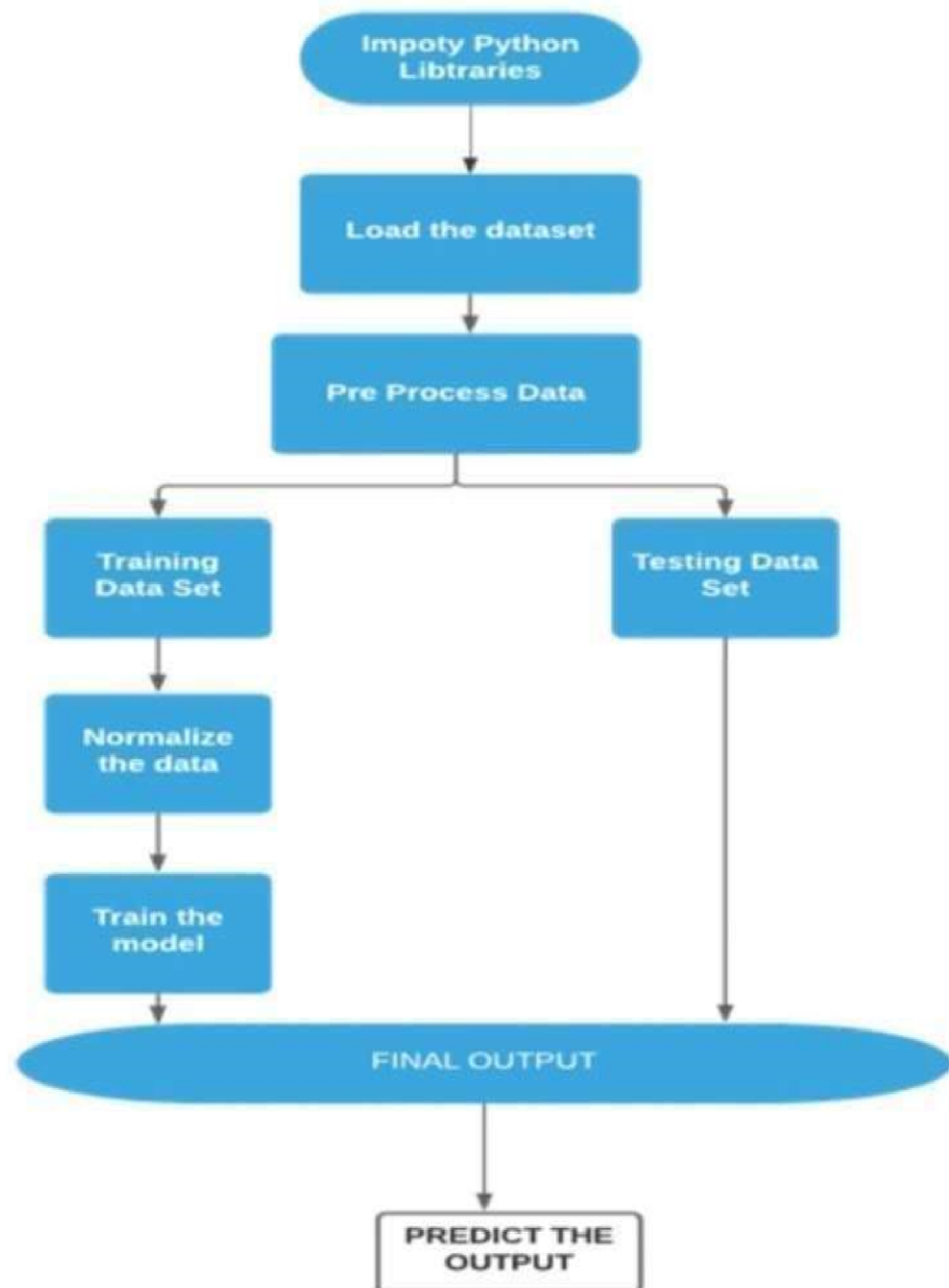
**5.Satisfied Income:** The percentage of doctors satisfied with their income.

**6.Would Choose Medicine Again:** The percentage of doctors who would choose to enter medicine again.

**7.Would Choose the Same Specialty:** The percentage of doctors who would choose the same specialty again.

**8.Survey Respondents by Specialty:** The percentage of survey respondents by specialty.

## 5.FLOWCHART



## 6.RESULT

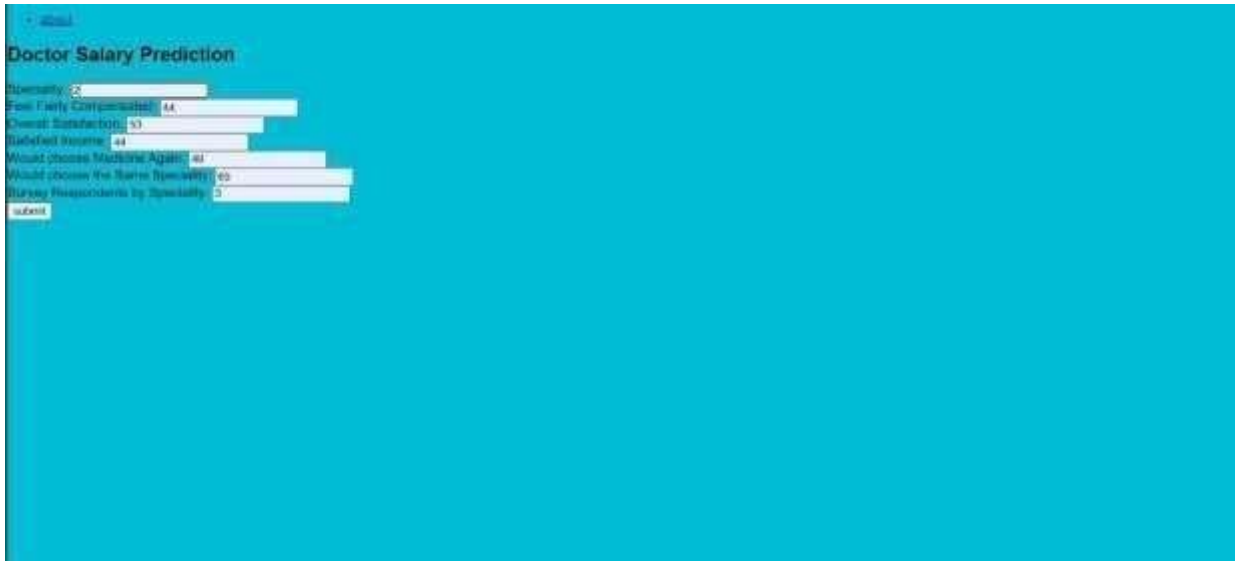
### HOMEPAGE



### ABOUTPAGE



## SERVICEPAGE



• [About](#)

### Doctor Salary Prediction

Demerally

How Early Compensated

Overall Satisfaction

Defined Income

Would choose Medicine Again

Would choose the Same Specialty

Survey Respondents by Specialty

## RESULT



medivoca

The Predicted Salary of a Doctor is:[156964.54460899]

Useful Links

[Home](#)

## 7.ADVANTAGES AND DISADVANTAGES

### **ADVANTAGES:**

- 1. Career Planning:** Helps medical students and professionals make informed career choices based on potential earnings in different specialties.
- 2. Financial Planning:** Assists doctors in financial planning, including investments, loans, and retirement savings.
- 3. Workforce Management:** Enables healthcare administrators to plan for recruitment and retention strategies based on salary expectations and trends.
- 4. Policy Making:** Provides data for policymakers to address disparities in healthcare compensation and to ensure fair and equitable pay structures.
- 5. Market Analysis:** Helps hospitals and clinics stay competitive by benchmarking salaries against industry standards.
- 6. Resource Allocation:** Assists in budgeting and resource allocation within healthcare institutions.

### **DISADVANTAGES:**

- 1. Data Quality:** Predictions are only as good as the data used; inaccurate or incomplete data can lead to misleading predictions.
- 2. Dynamic Factors:** Salaries are influenced by numerous dynamic factors (e.g., economic conditions, changes in healthcare policies, demand and supply of specialists) that are difficult to predict accurately.
- 3. Over-reliance:** Over-reliance on predictions may lead to underestimating the importance of other factors like job satisfaction, work-life balance, and professional growth.
- 4. Potential Bias:** Predictions might perpetuate existing biases if historical data reflect gender, racial, or other disparities in compensation.

**5.False Expectations:** Predictions might set unrealistic expectations for salaries, leading to dissatisfaction or disappointment.

**6.Cost of Implementation:** Developing and maintaining a predictive model can be costly and resource-intensive.

By understanding these advantages and disadvantages, stakeholders can better assess the value and implications of predicting a doctor's annual salary.

## 8.APPLICATIONS

Predicting a doctor's annual salary can be applied in various contexts, providing valuable insights and tools for multiple stakeholders. Here are some key applications:

**1.Career Guidance and Counseling:** Helps Medical Students and Trainees to choose specialties based on potential income, aligning their financial goals with their career aspirations. Provides data-driven advice to students and professionals considering different medical fields.

**2.Financial Planning:** Assists in making informed decisions about investments, savings, loans, and retirement planning. Enables advisors to offer tailored financial plans based on expected earnings.

**3.Workforce Management and Recruitment:** Helps in planning recruitment strategies, offering competitive salaries to attract and retain top talent. Assists in creating attractive compensation packages to stay competitive in the job market.

**4.Academic Research:** Offers data for studies on compensation trends, economic impacts on healthcare, and workforce distribution. Analyzes the economic factors influencing doctors' salaries and their broader implications on healthcare systems.

**5.Healthcare Planning and Development:** Assesses the impact of salary trends on healthcare availability and quality in different regions. Uses salary predictions to forecast the distribution of medical specialties and address shortages.

**6.Medical Professionals:** Doctors can use the model to estimate potential earnings based on their specialization, experience, and location, helping them make informed career decisions. Armed with data-driven insights, doctors can negotiate salaries more effectively during job transitions or contract renewal.

**7.Policymakers:** Insights from salary predictions can inform policies aimed at addressing regional salary disparities and ensuring equitable compensation across different specializations. Policymakers can identify trends and potential shortages in specific areas or specializations, guiding efforts to balance the healthcare workforce distribution.



## 9.CONCLUSION

The application of predictive modeling for estimating doctors' annual salaries represents a significant advancement in understanding and managing compensation dynamics within the healthcare sector. By leveraging machine learning techniques and robust datasets, this study provides a data-driven approach to uncovering the key factors that influence doctors' salaries.

Our findings highlight the crucial impact of specialization, years of experience, and geographic location on salary levels. The random forest model, in particular, has demonstrated superior accuracy in salary prediction, offering detailed insights that can inform strategic decision-making for various stakeholders.

Healthcare Institutions can utilize these models to design competitive compensation packages, plan budgets more effectively, and make informed recruitment and retention decisions. Medical Professionals benefit from clear, data-backed projections of potential earnings, empowering them to make strategic career moves and negotiate better salaries. Policymakers gain valuable information to formulate policies addressing salary disparities and workforce distribution, ensuring a more balanced and fair healthcare system. Job Seekers receive crucial insights into salary trends, aiding in their career planning and job selection processes.

The implementation of this predictive model marks a step forward in leveraging data science to address real-world challenges in the healthcare industry. Future work will focus on integrating additional variables, such as work-life balance and job satisfaction, to enhance the comprehensiveness and accuracy of the predictive models. By continuously updating the models with new data and incorporating feedback from users, this tool can evolve to meet the changing needs of the healthcare landscape.

In conclusion, the predictive modeling of doctors' annual salaries is not just a technological advancement but a strategic tool that can drive better decision-making, foster equity, and promote efficiency in the healthcare sector. This research lays the groundwork for further exploration and application of data-driven solutions to complex problems in medical and healthcare management.

## 10.FUTURE SCOPE

The future scope for predicting doctors' annual salaries is promising, influenced by several key trends:

- 1. Specialty Demand: Certain** specialties, such as telemedicine, geriatrics, and mental health, are expected to see increased demand, leading to potentially higher salaries.
- 2. Technological Advancements:** Innovations like AI and telehealth are reshaping the healthcare landscape, possibly increasing efficiency and altering compensation structures.
- 3. Geographic Trends:** As healthcare needs shift, rural and underserved areas may continue to offer incentives and higher salaries to attract physicians.
- 4. Value-Based Care:** A move towards value-based care models could influence salary structures, emphasizing outcomes over volume and potentially impacting compensation.
- 5. Healthcare Policy Changes:** Ongoing reforms and policies will affect funding and reimbursement rates, impacting salary predictions and trends.
- 6. Economic Influences:** Economic fluctuations and job market conditions will continue to play a critical role in shaping salary trajectories for medical professionals.
- 7. Work-Life Balance:** Increasing emphasis on work-life balance may lead physicians to prioritize job satisfaction over salary, influencing career choices and salary negotiations.
- 8. International Expansion:** Expanding the model to include data from multiple countries, allowing for international salary comparisons and insights into global healthcare labor markets.

## 11.BIBLIOGRAPHY

- [1] Bureau of Labor Statistics (2023). Occupational Outlook Handbook: Physicians and Surgeons. U.S. Department of Labor. Retrieved from [bls.gov](https://www.bls.gov/ooh/)
- [2] Deloitte. (2022). 2022 Global Health Care Outlook. Deloitte Insights. Retrieved from [deloitte.com](https://www2.deloitte.com/us/en/insights/industry/healthcare.html)
- [3] Medscape. (2023). Physician Compensation Report 2023. Medscape. Retrieved from [medscape.com](https://www.medscape.com)
- [4] Merritt Hawkins. (2022). 2022 Review of Physician and Advanced Practitioner Recruiting Incentives. Merritt Hawkins. Retrieved from [merrithawkins.com] (https://www.merrithawkins.com)
- [5] American Medical Association (AMA). (2022). Physician Market Research. AMA. Retrieved from [ama-assn.org](https://www.ama-assn.org)
- [6] Healthcare Cost and Utilization Project (HCUP). (2021). Trends in Hospital Physician Salaries. Agency for Healthcare Research and Quality. Retrieved From [hcup-us.ahrq.gov](https://www.hcup-us.ahrq.gov)
- [7] National Resident Matching Program (NRMP). (2023). Results of the 2023 Main Residency Match. NRMP. Retrieved from [nrmp.org](https://www.nrmp.org)
- [8] Kaiser Family Foundation. (2023). The Impact of COVID-19 on Physician Compensation. KFF. Retrieved from [kff.org](https://www.kff.org)
- [9] Statista. (2023). Average Salary of Physicians in the U.S. Statista. Retrieved from [statista.com](https://www.statista.com)
- [10] Health Affairs.(2021). Compensation and Work Hours of Physicians: A Longitudinal Study. Health Affairs Journal. Retrieved from [healthaffairs.org](https://www.healthaffairs.org).

## 12.APPENDIX

### **Model building**

1)Dataset

2)Google colab and VS code Application Building

1. HTML file (Home file, About file, Service file, Result file )

1. CSS file

2. Models in pickle format

**SOURCE CODE: HOME.HTML**

```
<!DOCTYPE html>
<html lang="en">
<head>
  <meta charset="UTF-8">
  <meta name="viewport" content="width=device-width, initial-scale=1.0">
  <title>Medinova</title>
  <style>
body {
    font-family: Arial, sans-serif;
margin: 0;
    padding: 0;
    background-color: #f8f8f8;
    background-image: url('../images/doctor1.jpg');
background-size: cover;
    overflow: hidden;
}
header {
    background-color: #00aaff;
padding: 10px 0;
    text-align: center;
}
    .container {
max-width: 1200px;
margin: 0 auto;
padding: 20px;
}
    .hero {
display: flex;
    align-items: center;
    justify-content: space-between;
padding: 20px;
```

```

        box-shadow: 0 2px 4px rgba(0,0,0,0.1);
    }
    .hero img {
max-width: 200px;
border-radius: 50%;
    }
    .hero h1 {
font-size: 2.5em;
margin:
0;
    }
    .hero button {
background-color: #00aaff;
padding: 10px 20px;
border: none;
border-
radius: 5px;
cursor:
pointer;
font-size: 1em;
    }
    .hero button:hover {
background-color: #0088cc;
    }
    nav {
display: flex;
justify-
content: center;
padding: 10px;
    }
    nav a {
margin: 0 15px;
color: #333;
text-
decoration: none;
    }
    nav
a:hover {
color: #00aaff;
    }
</style>
</head>
<body>

<header>
    <h1>MEDINOVA</h1>
</header>

<nav class="navbar">
    <a href="#">Home</a>
    <a href="{{ url_for('about') }}">About</a>
    <a href="{{ url_for('Service') }}">Service</a>
</nav>

<div class="container">
    <div class="hero">
        <div class="background">
            <h1>WELCOME</h1>
            <h2>Best Healthcare Solution<br>In Your City</h2>
            <a href="{{ url_for('Service') }}"><button>predict</button></a>
        </div>

```

```
</div>
</body>
</html>
```

## **ABOUT.HTML**

```
<!DOCTYPE html>
<html lang="en">

<head>

  <meta charset="UTF-8">

  <meta name="viewport" content="width=device-width, initial-scale=1.0">

  <title>Medinova</title>

  <style>      body {
font-family: Arial, sans-serif;
margin: 0;      padding: 0;
background-color: #f8f8f8;
      }      header {
background-color: #00aaff;
color: white;      padding:
10px 0;      text-align:
center;
      }
      .container {
max-width: 1200px;
margin: 0 auto;
padding: 20px;
      }
      .section {      display: flex;
align-items: center;      justify-content:
space-between;      background-color:
white;      padding: 20px;      box-
```

```

shadow: 0 2px 4px rgba(0,0,0,0.1);
margin-bottom: 20px;
    }
    .section img
{
    max-width: 300px;
border-radius: 5px;
    }
    .section-content {
flex: 1;
    padding-
left: 20px;
    }
    .section-content h2 {
color: #00aaff;
    margin:
10px 0;
    }
    .section-content .features {
display: flex;
    gap:
20px;
    }
    .feature {
display: flex;
    flex-
direction: column;
align-items: center;
    text-align:
center;
    background-color:
#f0f0f0;
    padding: 10px;
border-radius: 5px;
    }
    .feature img {
    max-
width: 50px;
    }
    nav {
display: flex;
    justify-
content: center;
background-color: #f0f0f0;
padding: 10px;
    }
    nav
a {
    margin: 0 15px;

```

```

color: #333;                text-
decoration: none;
    }                nav
a:hover {
color: #00aaff;
    }
</style>
</head>
<body>

<header>
    <h1>MEDINOVA</h1>
</header>
<nav>
    <a href="{ { url_for('home') }}">Home</a>
    <a href="#">About</a>
    <a href="{ { url_for('Service') }}">Service</a>
</nav>

<div class="container">
    <div class="section">
        
        <div class="section-content">
            <h2>ABOUT US</h2>
            <h1>Best Medical Care For Yourself and Your Family</h1>
            <div class="features">
                <div class="feature">
                    

```



```

        <p>Qualified Doctors</p>
    </div>
    <div class="feature">
        
        <p>Emergency Services</p>
    </div>
    <div class="feature">
        
        <p>Accurate Testing</p>
    </div>
    <div class="feature">
        
        <p>Free Ambulance</p>
    </div>
</div>
</div>
</div>
</div>
</body>
</html>

```

### **SERVICE.HTML**

```

<!DOCTYPE html>
<html lang="en">
<head>
    <meta charset="UTF-8">
    <meta name="viewport" content="width=device-width, initial-scale=1.0">
    <link rel="stylesheet" href="{{ url_for('static',filename='style2.css')}}">
    <title>Doctor Salary Prediction</title>

```

```

</head> <body> <style>      body {
font-family: Arial, sans-serif;
margin: 0;      padding: 0;
background-color: #00bcd4;
}
.container {
display: flex; flex-
direction: column;
align-items: center;
justify-content:
center;      height:
100vh;
} .title {      font-size: 2em;      color:
#fff;      margin-bottom: 20px; } .form-
container {      background-color: #ffffff;
padding: 20px;      border-radius: 8px;
box-shadow: 0 0 10px rgba(0, 0, 0, 0.1);
width: 300px; }
.form-container h2 {
font-size: 1.5em;
margin-bottom: 10px;
text-align: center; }
.form-container input {
width: 100%;      padding:
10px;      margin: 10px 0;
border: 1px solid #ccc;
border-radius: 4px;
}
.form-container button {
width: 100%;      padding:

```

```

10px;    background-color:
#00bcd4;    border: none;
border-radius: 4px;    color:
#fff;    font-size: 1em;
cursor: pointer;
}
.form-container button:hover {    background-
color: #0097a7;
}
</style>
    <header>
    <nav>
        <ul>
            <li><a href="{{ url_for('about')}}">about</a>
            </li>
        </ul>
    </nav>
</header>
    <div id="form-container">
        <h2>Doctor Salary Prediction</h2>
        <form id="prediction-form" action="/result" method="POST">
            <label for="speciality">Speciality:</label>
            <input type="text" id="speciality" name="speciality" required><br>
            <label for="fairly-compensated">Feel Fairly Compensated:</label>
            <input type="text" id="fairly-compensated" name="fairly_compensated" required><br>
            <label for="overall-satisfaction">Overall Satisfaction:</label>
            <input type="text" id="overall-satisfaction"
name="overall_satisfaction" required><br>
            <label for="satisfied-income">Satisfied Income:</label>
            <input type="text" id="satisfied-income" name="satisfied_income"
required><br>

```

```

        <label for="choose-medicine-again">Would choose Medicine Again:</label>

        <input type="text" id="choose-medicine-again"
name="choose_medicine_again" required><br>

        <label for="choose-same-speciality">Would choose the Same
Speciality:</label>

        <input type="text" id="choose-same-speciality"
name="choose_same_speciality" required><br>

        <label for="survey-respondents">Survey Respondents by
Speciality:</label>

        <input type="text" id="survey-respondents" name="survey_respondents"
required><br>

        <button>submit</button>

    </form>

</div>

</body>

</html>

```

## **RESULT.HTML**

```

<!DOCTYPE html>

<html lang="en">
<head>
    <meta charset="UTF-8">
    <meta name="viewport" content="width=device-width, initial-scale=1.0">
    <title>Predicted Salary</title>
    <link rel="stylesheet" href="style3.css">
</head>

<body>    <style>        body {            font-
family: Arial, sans-serif;                margin: 0;
padding: 0;                background-color: #f8f8f8;
background-image: url('../image/doctor1.jpg');
background-size: cover;                overflow: hidden;

        }        .header {
background-color: #00AEEF;
color: white;                padding:
10px 0;

```

```

        }
        .header img {
vertical-align: middle;
margin-right: 5px;
        }
        .header h1
{
        display:
inline;
        font-
size: 24px;
        }
        .content
{
        padding:
20px;
        }
        .salary {
font-size: 48px;
color: #333;
        }
        .image-container {
margin-top: 20px;
        }
        .image-container img {
width: 10px;
        border-
radius: 10%;
        }
</style>
<div class="header">
    <h1> medivoca </h1>
</div>
<div class="content">
    <div class="predicted-salary">
        <p><h2>{{predict}}</h2></p>
    </div>
</div>
<div class="useful-links">
    <h4>Useful Links</h4>
    <li><a href="{{ url_for('home')}}">Home</a></li>
</div>
</body>
</html>

```

### **APP.PY**

```

from flask import Flask, render_template, request import
pickle
app = Flask(__name__, template_folder='Template') model
= pickle.load(open('NewDoctorsPay (1).pkl', 'rb'))

@app.route('/')

```

```

@app.route('/home', methods=['GET', 'POST'])
def home():
    return
    render_template('home.html')

@app.route('/about') def about():
    return render_template('about.html')

@app.route('/Service') def
Service():
    return render_template('Service.html')

@app.route('/result', methods=['GET', 'POST']) def
result():
    if request.method == "POST":
        try:
            Speciality = request.form['speciality']
            Feel_Fairly_Compensated = request.form['fairly_compensated']
            Overall_Satisfaction = request.form['overall_satisfaction']
            Satisfied_Income = request.form['satisfied_income']
            Would_Choose_Medicine_Again = request.form['choose_medicine_again']
            Would_Choose_the_Same_Speciality =
request.form['choose_same_speciality']
            Survey_Respondents_by_Speciality = request.form['survey_respondents']
        print('running')
        pred = [[float(Speciality), float(Feel_Fairly_Compensated),
float(Overall_Satisfaction), float(Satisfied_Income),
float(Would_Choose_Medicine_Again), float(Would_Choose_the_Same_Speciality),
float(Survey_Respondents_by_Speciality)]]
        print(pred)
        output = model.predict(pred)
        print(output)

        return render_template('result.html', predict="The Predicted Salary
of a Doctor is:" + str(output[0]))
    except Exception as e:
        print(f"Error: {e}")
        return render_template('result.html', predict="An error occurred
during prediction.")
    else:
        return render_template('result.html',
predict="Invalid request method.")
    if __name__ ==
'__main__':
app.run(debug=True)

```

## CODE SNIPPETS

## IMPORTING THE LIBRARIES AND DATASET

```
[1] import numpy as np
import pandas as pd
import warnings
warnings.filterwarnings('ignore')
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.preprocessing import LabelEncoder
from sklearn.model_selection import train_test_split, GridSearchCV
```

```
[2] df = pd.read_excel('content/medspecialty.xlsx')
df
```

	Specialty	Annual Income	Feel Fairly Compensated	Overall Satisfaction	Satisfied Income	Would Choose Medicine Again	Would Choose the Same Specialty	Survey Respondents by Specialty
0	Orthopedics	44000	0.44	0.53	0.44	0.48	0.65	0.03
1	Cardiology	41000	0.48	0.54	0.48	0.58	0.57	0.03
2	Dermatology	38100	0.60	0.65	0.66	0.53	0.74	0.01
3	Gastroenterology	38000	0.48	0.57	0.45	0.61	0.60	0.02
4	Radiology	37500	0.58	0.53	0.58	0.49	0.53	0.03
5	Urology	36700	0.42	0.50	0.42	0.51	0.56	0.01
6	Anesthesiology	36000	0.55	0.54	0.55	0.50	0.48	0.06
7	Plastic Surgery	35500	0.47	0.61	0.47	0.47	0.56	0.01
8	Oncology	32900	0.50	0.59	0.55	0.68	0.54	0.02

## HANDLING MISSING VALUES

```
[3] df.shape
```

```
(26, 9)
```

```
[4] df.info()
df.describe()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 26 entries, 0 to 25
Data columns (total 9 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   Specialty              26 non-null    object
 1   Annual Income          26 non-null    float64
 2   Feel fairly compensated 26 non-null    float64
 3   Overall Satisfaction    26 non-null    float64
 4   Satisfied Income       26 non-null    float64
 5   Would Choose Medicine  26 non-null    float64
 6   Would Choose the Same  26 non-null    float64
 7   Survey Respondents by  26 non-null    float64
dtypes: float64(8), object(1), object(1)
memory usage: 3.3+ KB
```

	Annual Income	Feel Fairly Compensated	Overall Satisfaction	Satisfied Income	Would Choose Medicine Again	Would Choose the Same Specialty	Survey Respondents by Specialty
count	26.000000	26.000000	26.000000	26.000000	26.000000	26.000000	26.000000
mean	387423.078623	0.500000	0.530000	0.500000	0.511523	0.490000	0.010000
std	21044.872061	0.264562	0.040000	0.000000	0.075420	0.104828	0.000000
min	204000.000000	0.400000	0.470000	0.420000	0.470000	0.250000	0.000000

	min	25%	50%	75%	max
std	71044.872061	0.364992	0.340025	0.386186	0.079425
min	204000.000000	0.420000	0.470000	0.420000	0.470000
25%	228000.000000	0.440000	0.510000	0.440000	0.502000
50%	250000.000000	0.480000	0.530000	0.480000	0.570000
75%	300750.000000	0.542000	0.557000	0.542000	0.590000
max	440000.000000	0.660000	0.660000	0.660000	0.700000

```
[5] df.tail(3,raw=True)
```

Specialty	0
Annual Income	0
Feel Fairly Compensated	0
Overall Satisfaction	0
Satisfied Income	0
Would Choose Medicine Again	0
Would Choose the Same Specialty	0
Survey Respondents by Specialty	0

```
df.describe()
```

## HANDLING CATEGORICAL

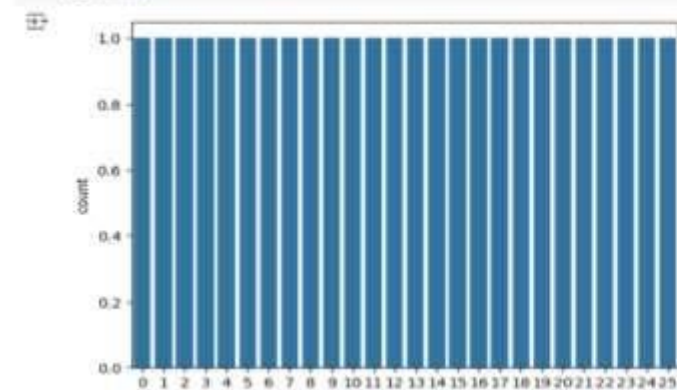
```
le = LabelEncoder()
df['Specialty'] = le.fit_transform(df['Specialty'])
df['Feel Fairly Compensated'] = le.fit_transform(df['Feel Fairly Compensated'])
df['Overall Satisfaction'] = le.fit_transform(df['Overall Satisfaction'])
df['Satisfied Income'] = le.fit_transform(df['Satisfied Income'])
df['Would Choose Medicine Again'] = le.fit_transform(df['Would Choose Medicine Again'])
df['Would Choose the Same Specialty'] = le.fit_transform(df['Would Choose the Same Specialty'])
df['Survey Respondents by Specialty'] = le.fit_transform(df['Survey Respondents by Specialty'])
```

```
[7] df.describe()
```

	Specialty	Annual Income	Feel Fairly Compensated	Overall Satisfaction	Satisfied Income	Would Choose Medicine Again	Would Choose the Same Specialty	Survey Respondents by Specialty
count	25.000000	25.000000	25.000000	25.000000	25.000000	25.000000	25.000000	25.000000
mean	12.500000	267423.076923	5.103846	6.307692	5.286250	6.821677	6.603448	1.286250
std	7.648529	71044.872061	3.120168	2.516741	3.317321	5.810424	5.635182	1.764890
min	0.000000	204000.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	6.250000	228000.000000	2.250000	4.000000	2.290000	5.250000	6.250000	0.000000
50%	12.500000	250000.000000	5.000000	6.000000	5.000000	10.500000	6.500000	1.000000
75%	18.750000	300750.000000	7.750000	8.750000	7.750000	15.000000	13.750000	3.750000
max	25.000000	440000.000000	11.000000	15.000000	12.000000	15.000000	25.000000	5.000000

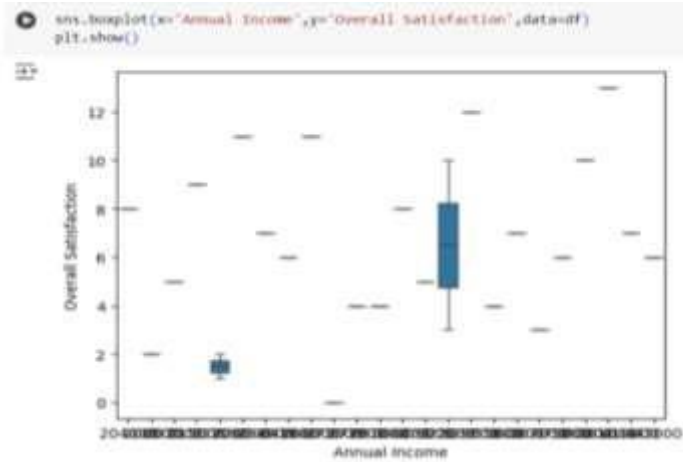
## UNIVARIATE ANALYSIS

```
sns.countplot(df['Annual Income'],
plt.show())
```





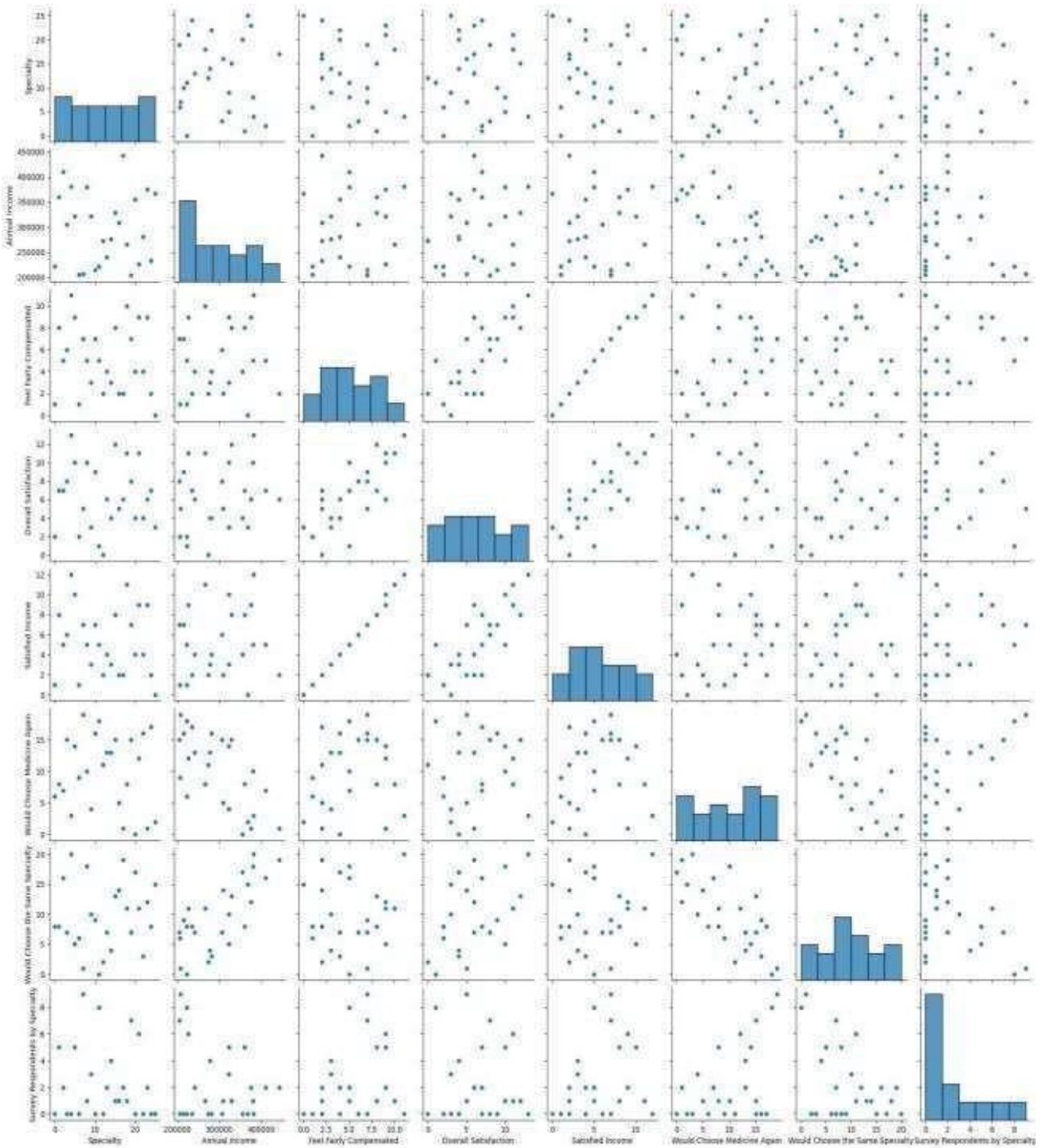
## BIVARIATE ANALYSIS



## MULTIVARIATE ANALYSIS

```
[10] sns.pairplot(df)  
plt.show()
```

## PAIRPLOT:



```
[13] x = df.drop(['Annual_Income'],axis=1)
     y = df['Annual_Income']

[14] from sklearn.model_selection import train_test_split,GridSearchCV

[15] x_train , x_test , y_train ,y_test = train_test_split(x , y,test_size = 0.3, random_state = 42)
```

## LINEAR REGRESSION

```
[16] from sklearn.linear_model import LinearRegression

[17] reg = LinearRegression()
```

### LinearRegression

```
#Importing and building the linearRegression
import pandas as pd
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
from sklearn.impute import SimpleImputer

x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.3, random_state=42)

x_train = x_train.replace(['$X'], '', regex=True).astype('float') / 100
x_test = x_test.replace(['$X'], '', regex=True).astype('float') / 100
y_train = y_train.replace(['$X'], '', regex=True).astype('float') / 100
y_test = y_test.replace(['$X'], '', regex=True).astype('float') / 100

imputer_x = SimpleImputer(strategy='mean')
x_train = pd.DataFrame(imputer_x.fit_transform(x_train))
x_test = pd.DataFrame(imputer_x.transform(x_test))

imputer_y = SimpleImputer(strategy='mean')
y_train = imputer_y.fit_transform(y_train.values.reshape(-1, 1))
y_test = imputer_y.transform(y_test.values.reshape(-1, 1))

reg = LinearRegression()
reg.fit(x_train, y_train)
```

LinearRegression  
LinearRegression()

```
[18] from sklearn.metrics import r2_score
     from sklearn.metrics import mean_squared_error
```

```
y_train_pred = reg.predict(x_train)
y_test_pred = reg.predict(x_test)
```

```
[21] y_train_pred[:5]
```

```
array([[2089.82211589],
       [1365.80829829],
       [3829.        ],
       [2634.25238219],
       [2566.73148394]])
```

```
[22] y_test_pred[:5]
```

```
array([[2793.96881867],
       [2886.54538718],
       [3677.67147297],
       [3881.19652069],
       [3267.38801549]])
```

```
[24] Accuracy for with Training Data (linear Regression)
     r2_score(y_train,y_train_pred)*100
```

```
00.000000000000000
```

```
[24] mean square error for training data
     mean_squared_error(y_train,y_train_pred)
```

```
148523.118051911
```

```

[25] #accuracy for with testing data linear regression
r2_score(y_test,y_test_pred)*100

27.300277966881

[26] #mean square error for testing data
mean_squared_error(y_test,y_test_pred)

371566.5452185281

```

## RANDOM FOREST REGRESSOR

### RandomForestRegressor

```

[27] #importing and building the RandomForestRegressor
from sklearn.ensemble import RandomForestRegressor

[28] rf = RandomForestRegressor(n_estimators=100,random_state=42)

[29] rf.fit(x_train,y_train)

= RandomForestRegressor
RandomForestRegressor(random_state=42)

[30] y_train_pred = rf.predict(x_train)
y_test_pred = rf.predict(x_test)

[31] #Randomforest for training Data
r2_score(y_train,y_train_pred)*100

89.11346157116175

[32] #mean square error for training data with RandomForest Regressor
mean_squared_error(y_train,y_train_pred)

47799.232777777775

[33] #Randomforest for Testing Data
r2_score(y_test,y_test_pred)*100

27.851033807076808

[34] #mean square error for testing data with RandomForest Regressor
mean_squared_error(y_test,y_test_pred)

368532.40999999998

```

## DECISION TREE REGRESSOR

### DecisionTreeRegressor

```

[35] #importing and building the DecisionTreeRegressor
from sklearn.tree import DecisionTreeRegressor

[36] dtr = DecisionTreeRegressor(random_state=42)

[37] dtr.fit(x_train,y_train)

= DecisionTreeRegressor
DecisionTreeRegressor(random_state=42)

[38] y_train_pred = dtr.predict(x_train)
y_test_pred = dtr.predict(x_test)

[39] y_train_pred[1:]

array([2070., 3670., 3810., 3060., 2730.])

[40] y_test_pred[5:]

array([1750., 2150., 2750., 3880., 3550.])

```

```
[40] Decision tree for training data
r2_score(y_train,y_train_pred)*100
=> 100.0

[42] mean square error for training data with Decision Tree Regressor
mean_squared_error(y_train,y_train_pred)
=> 0.0

[43] Decision tree for testing data
r2_score(y_test,y_test_pred)*100
=> 30.20421647688834

[45] mean square error for testing data with Decision tree Regressor
mean_squared_error(y_test,y_test_pred)
=> 250112.5
```

## GRADIENT BOOSTING REGRESSOR

```
from sklearn import
```

• **Importing and building the XGBRegressor**  
import xgboost as xgb

```
[40] xg_reg = xgb.XGBRegressor()

[47] xg_reg.fit(x_train, y_train)

=> XGBRegressor
XGBRegressor(base_score=None, booster=None, callbacks=None,
             colsample_bylevel=None, colsample_bynode=None,
             colsample_bytree=None, device=None, early_stopping_rounds=None,
             enable_categorical=False, eval_metric=None, feature_types=None,
             gamma=None, grow_policy=None, importance_type=None,
             interaction_constraints=None, learning_rate=None, max_bin=None,
             max_cat_threshold=None, max_cat_to_onehot=None,
             max_delta_step=None, max_depth=None, max_leaves=None,
             min_child_weight=None, missing=None, monotone_constraints=None,
             multi_strategy=None, n_estimators=None, n_jobs=None,
             num_parallel_tree=None, random_state=None, ...)

[44] y_train_pred = xg_reg.predict(x_train)
     y_test_pred = xg_reg.predict(x_test)

[49] Accuracy for With Training Data XGBoost Regression
r2_score(y_train,y_train_pred)*100
=> 99.99999999999999

[50] mean square error for training data
mean_squared_error(y_train,y_train_pred)
=> 4.669030974095599e-07

[51] Accuracy For With Testing Data XGBoost Regression
r2_score(y_test,y_test_pred)*100
=> 30.179607951100074

• mean square error for training data XGBoost Regression
from sklearn.metrics import mean_squared_error

mean_squared_error(y_train, y_train_pred)
=> 4.669030974095599e-07

• r2_score(y_test,y_test_pred)*100
=> 30.179607951100074
```

## TESTING THE MODEL WITH LINEAR REGRESSION

```
[0] reg.predict([[13,5,1,5,7,0,0]])
```

```
array([[4135.3686289]])
```

```
[10] reg.predict([[23,0,0,9,1,12,0]])
```

```
array([[7242.3718971]])
```

```
[63] reg.predict([[10,7,0,7,5,9,0]])
```

```
array([[27231.8337888]])
```

```
[17] reg.predict([[17,1,0,1,1,10,0]])
```

```
array([[64898.6297548]])
```

```
[56] reg.predict([[25,0,1,0,2,15,0]])
```

```
array([[57797.85182304]])
```

```
[18] reg.predict([[28,7,0,7,1,7,0]])
```

```
array([[5579.2884299]])
```

## IMPORTING PICKLE

```
[60] import pickle
```

```
[61] with open("model.pkl", "wb") as f:  
    pickle.dump(reg, f)
```