

DATA PRODUCT ON REMOTE  
WORK AND MENTAL  
HEALTH

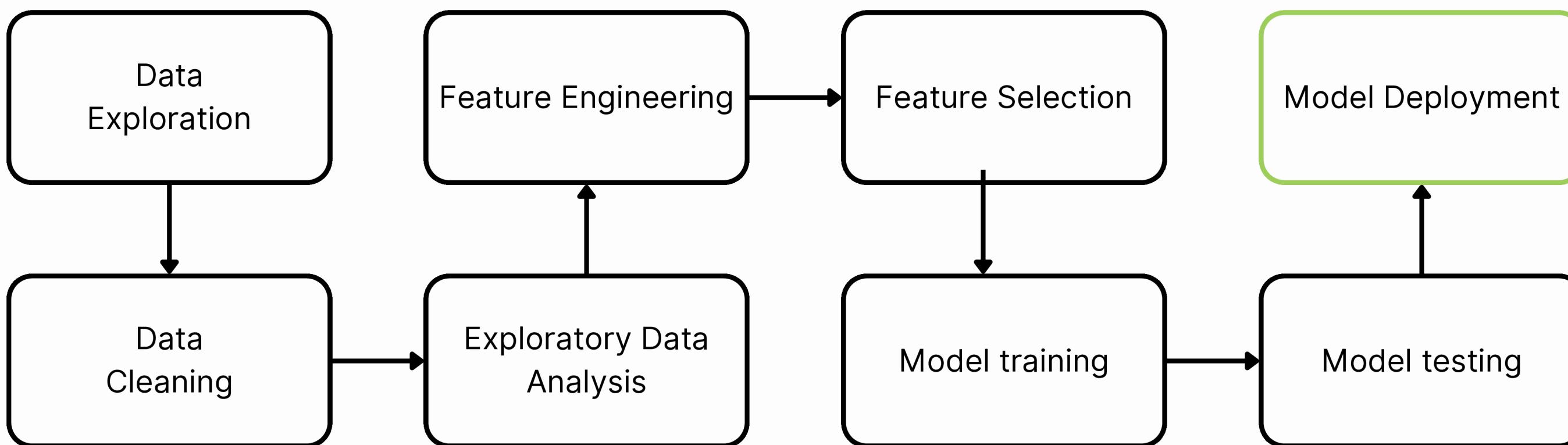
ANALYZING FACTORS  
INFLUENCING EMPLOYEE  
SATISFACTION WITH REMOTE  
WORK

# OBJECTIVE OF THE ASSIGNMENT

- Build a data product around the Kaggle dataset on remote work and mental health.
- Address key questions regarding factors affecting employee satisfaction.
  - What factors influence employee satisfaction in remote work?
  - How to predict employee satisfaction based on identified factors?

**Dataset Source:** Kaggle: Remote Work and Mental Health

# Protocol



# Data Info

Column	Non null	Datatype	Unique
Employee_ID	5000	String	5000
Age	5000	Integer	39
Gender	5000	String	4
Job_Role	5000	String	7
Industry	5000	String	7
Years_of_Experience	5000	Integer	35
Work_Location	5000	String	3
Hours_Worked_Per_Week	5000	integer	41
Number_of_Virtual_Meetings	5000	integer	16
Work_Life_Balance_Rating	5000	integer	5
Stress_Level	5000	String	3
Mental_Health_Condition	3804	String	3
Access_to_Mental_Health_Resources	5000	String	2
Productivity_Change	5000	String	3
Social_Isolation_Rating	5000	integer	5
Satisfaction_with_Remote_Work	5000	String	3
Company_Support_for_Remote_Work	5000	integer	5
Physical_Activity	3371	String	2
Sleep_Quality	5000	String	3
Region	5000	String	6

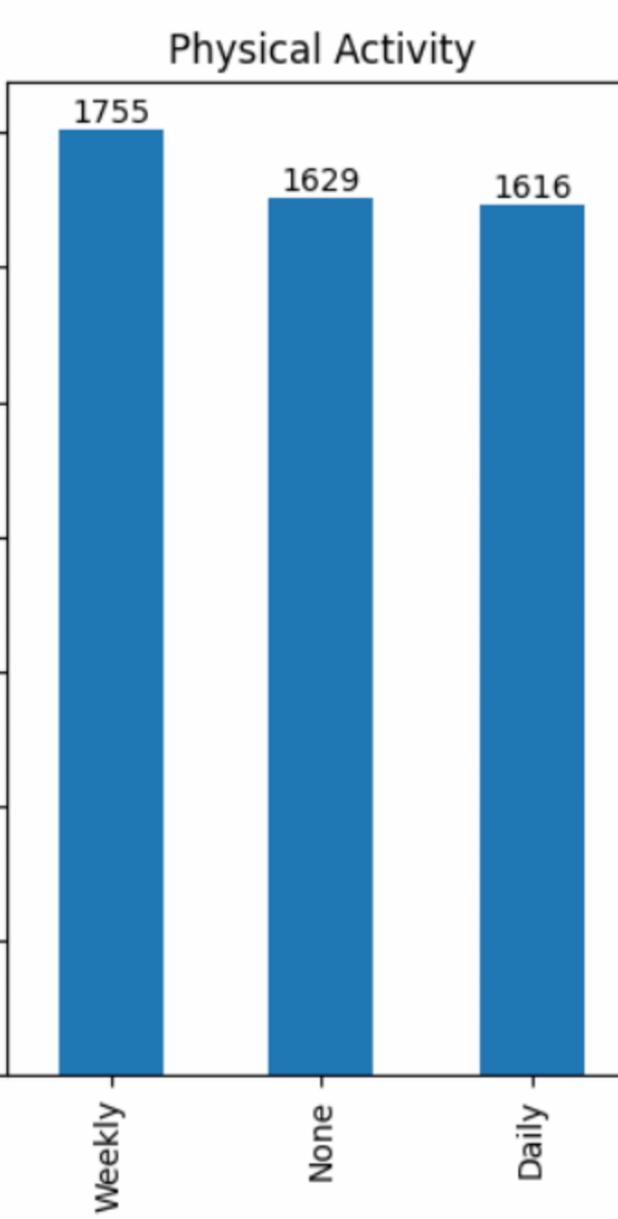
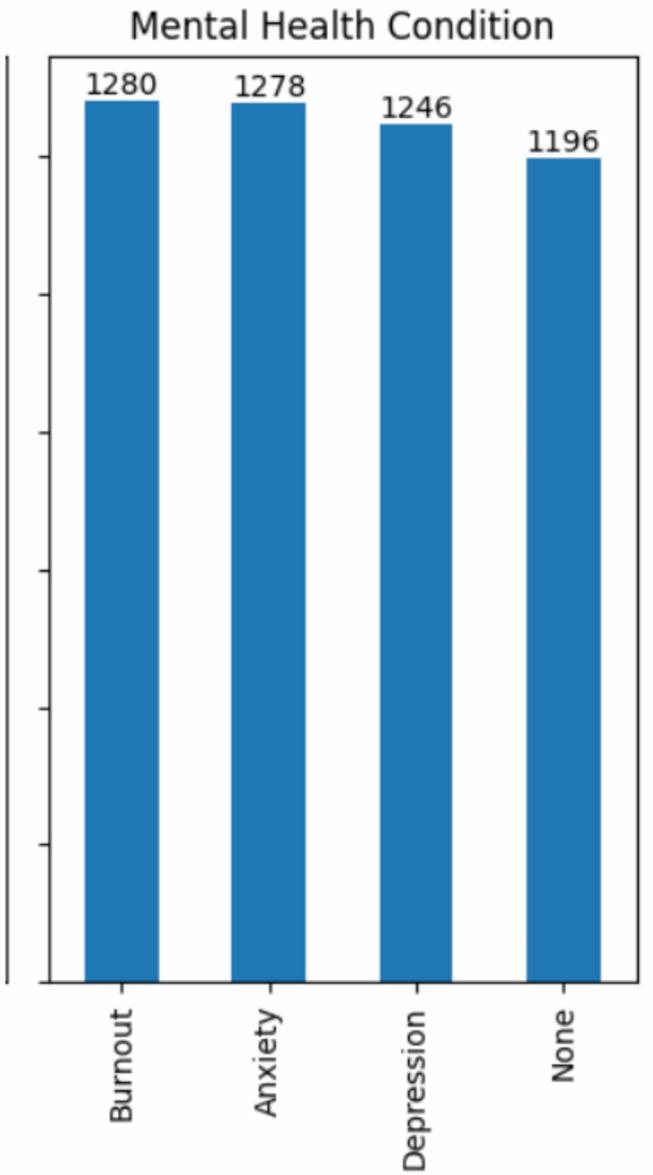
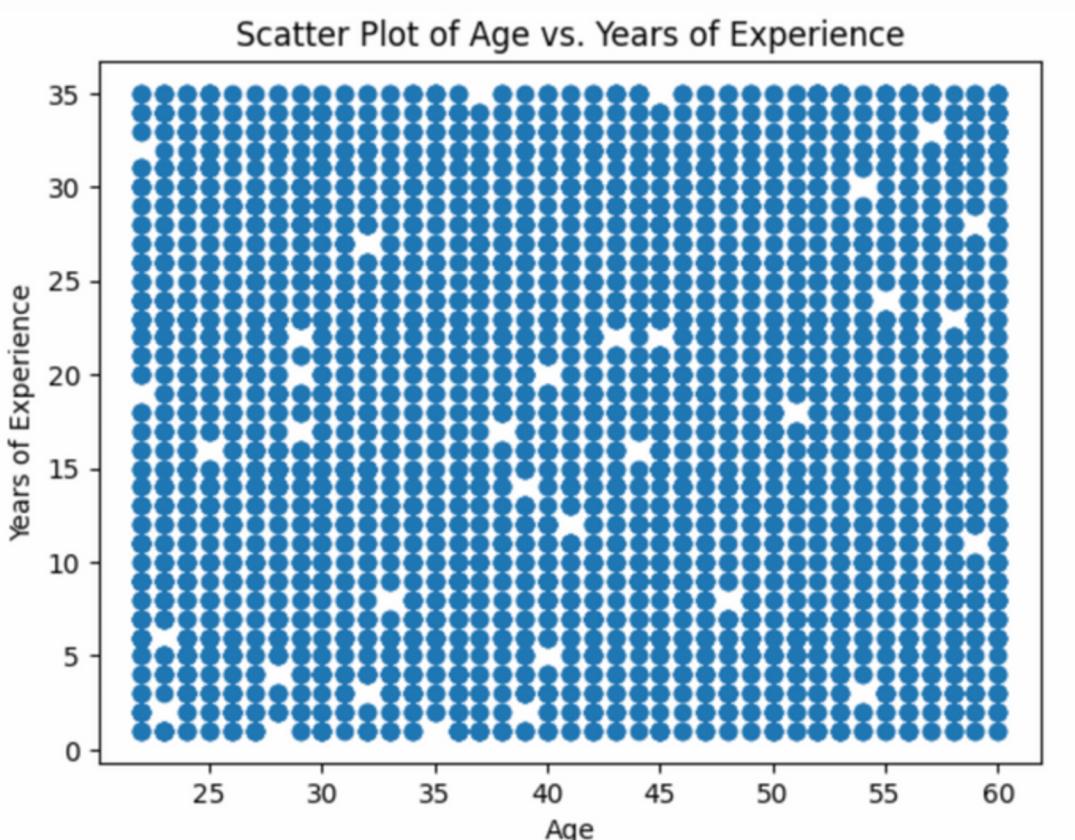
## Key Points

1. Employee\_ID: all data is unique don't add anything to the model
2. Most of the columns are not null
3. Physical\_Activity and Mental\_Health\_Condition are having null values
4. Age , Years\_of\_Experience , Hours\_Worked\_Per\_Week , Number\_of\_Virtual\_Meetings , Work\_Life\_Balance\_Rating , Social\_Isolation\_Rating , Company\_Support\_for\_Remote\_Work are integer, but we can't say which are ordinal
5. Gender , Job\_Role , Industry , Work\_Location , Stress\_Level , Mental\_Health\_Condition , Access\_to\_Mental\_Health\_Resources , Productivity\_Change , Satisfaction\_with\_Remote\_Work , Physical\_Activity , Sleep\_Quality , Region are strings so they need to be encoded.

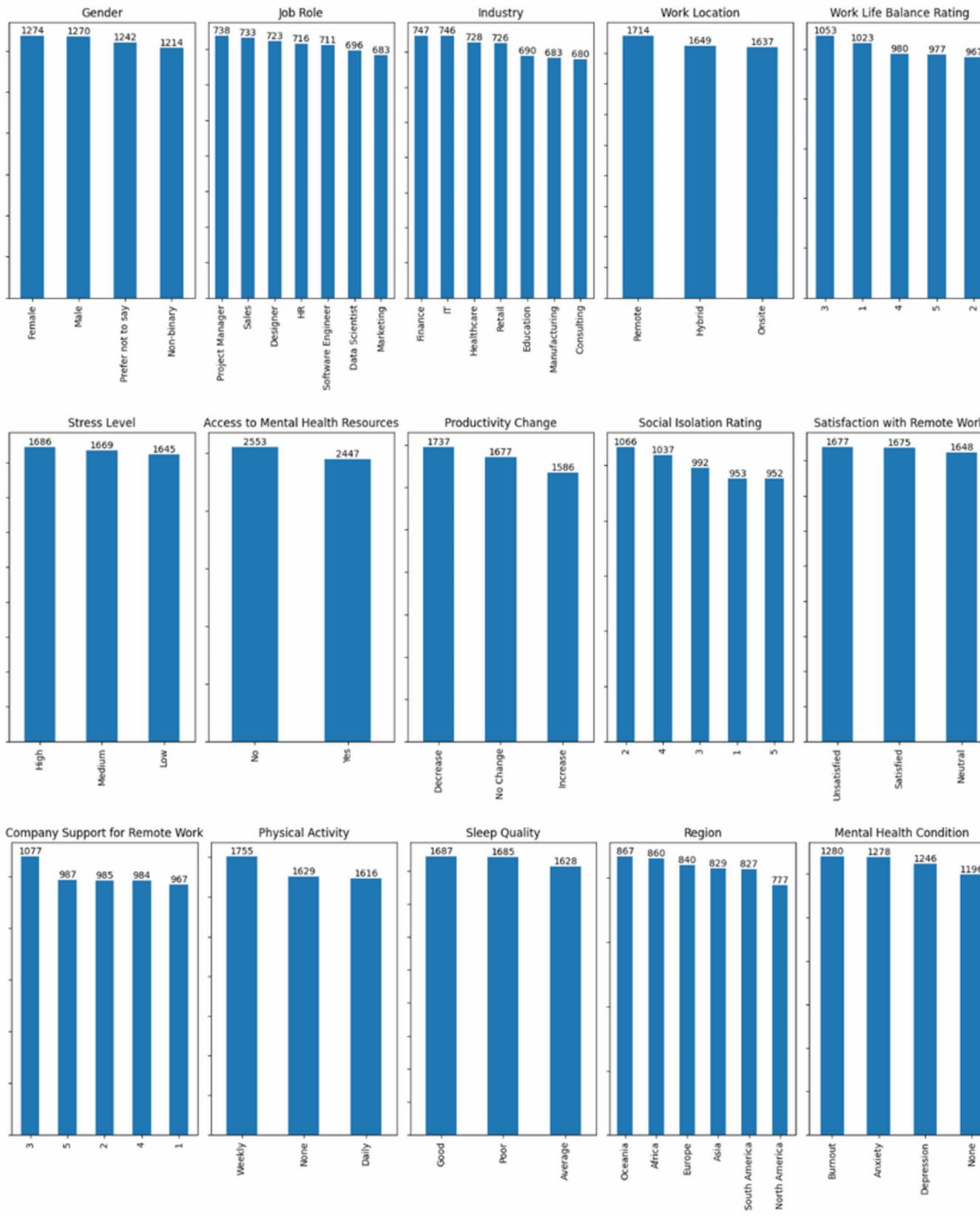
# Data Cleaning

1. On checking data I have found the data is clean according to the syntax but logically it is having issue when we check year of experience and Age which shows we have 25 age person with 35 year of experience which don't make sense. For cleaning this we have removed anyone who is having more than 18 years of between Age and Experience considering an individual can start working from age 18.

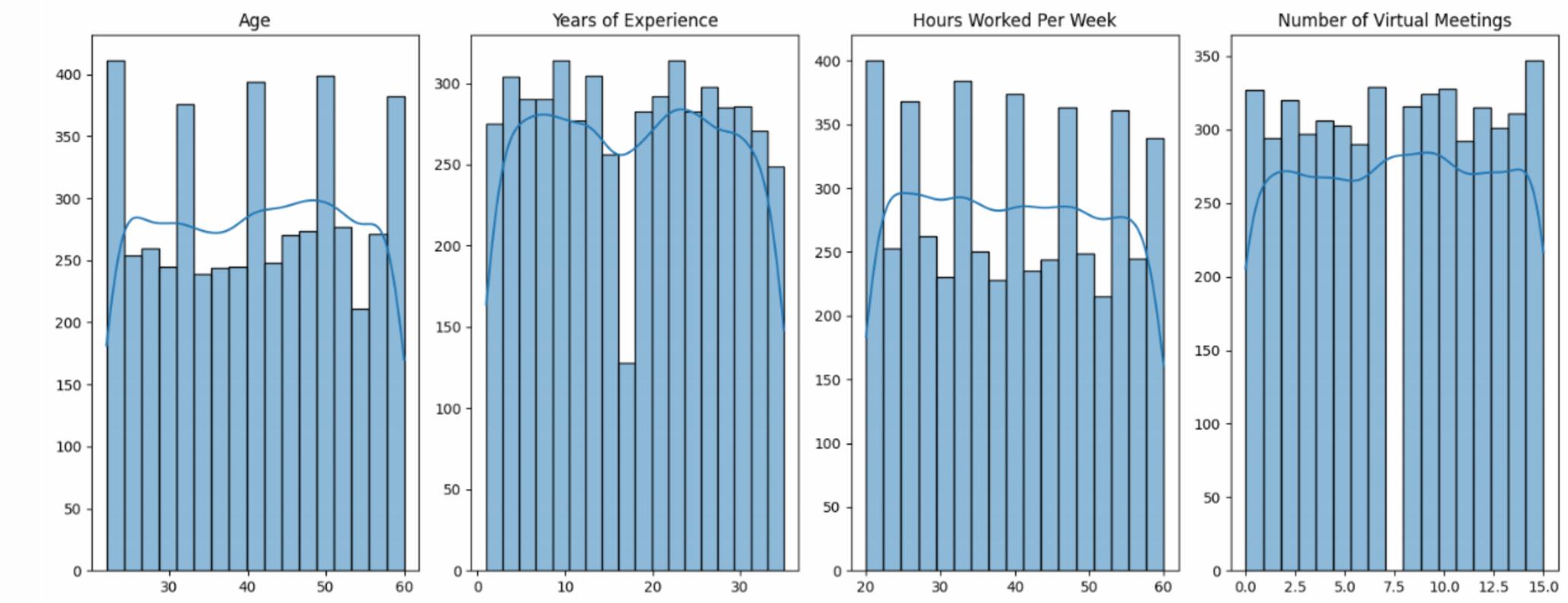
2. Physical\_Activity and Mental\_Health\_Condition are having None in them which make sense and we considered it as employees are not active and not feeling any tirednes, or burnout that's why we got None, so we replaced Nan as None



# Exploratory Data Analysis (EDA)



- 1- On a quick glimpse we can see that data is very balance in both numerical and string columns, so we don't have to take care of class imbalancing .
- 2- This behavior is same when we checked interaction between different features.



# String Feature check

Gender	Job_Role	Industry	Work_Location	Stress_Level	Mental_Health_Condition	Access_to_Mental_Health_Resources	Productivity_Change	Satisfaction_with_Remote_Work	Physical_Activity	Sleep_Quality	Region
Non-binary	HR	Healthcare	Hybrid	Medium	Depression	No	Decrease	Unsatisfied	Weekly	Good	Europe
Female	Data Scientist	IT	Remote	High	Anxiety	Yes	Increase	Satisfied	None	Poor	Asia
Male	Software Engineer	Education	Onsite	Low	None		No Change	Neutral	Daily	Average	North America
Prefer not to say	Sales	Finance			Burnout						South America
	Marketing	Consulting									Oceania
	Designer	Manufacturing									Africa
	Project Manager	Retail									

1. If we see columns like Stress\_Level, Productivity\_Change, Satisfaction\_with\_Remote\_Work and Sleep\_Quality, these are the values can be scaled and this we can scale with Ordinal Encoder
2. Whereas as columns like Gender, we can't say Female is 3 and Male is 2, so this columns we will encode with one hot encoder

# Feature Categories

## CATEGORICAL

Gender  
Job\_Role  
Industry  
Work\_Location  
Mental\_Health\_Condition  
Access\_to\_Mental\_Health\_Resources  
Physical\_Activity  
Region

## ORDINAL

Stress\_Level  
Productivity\_Change  
Satisfaction\_with\_Remote\_Work  
Sleep\_Quality

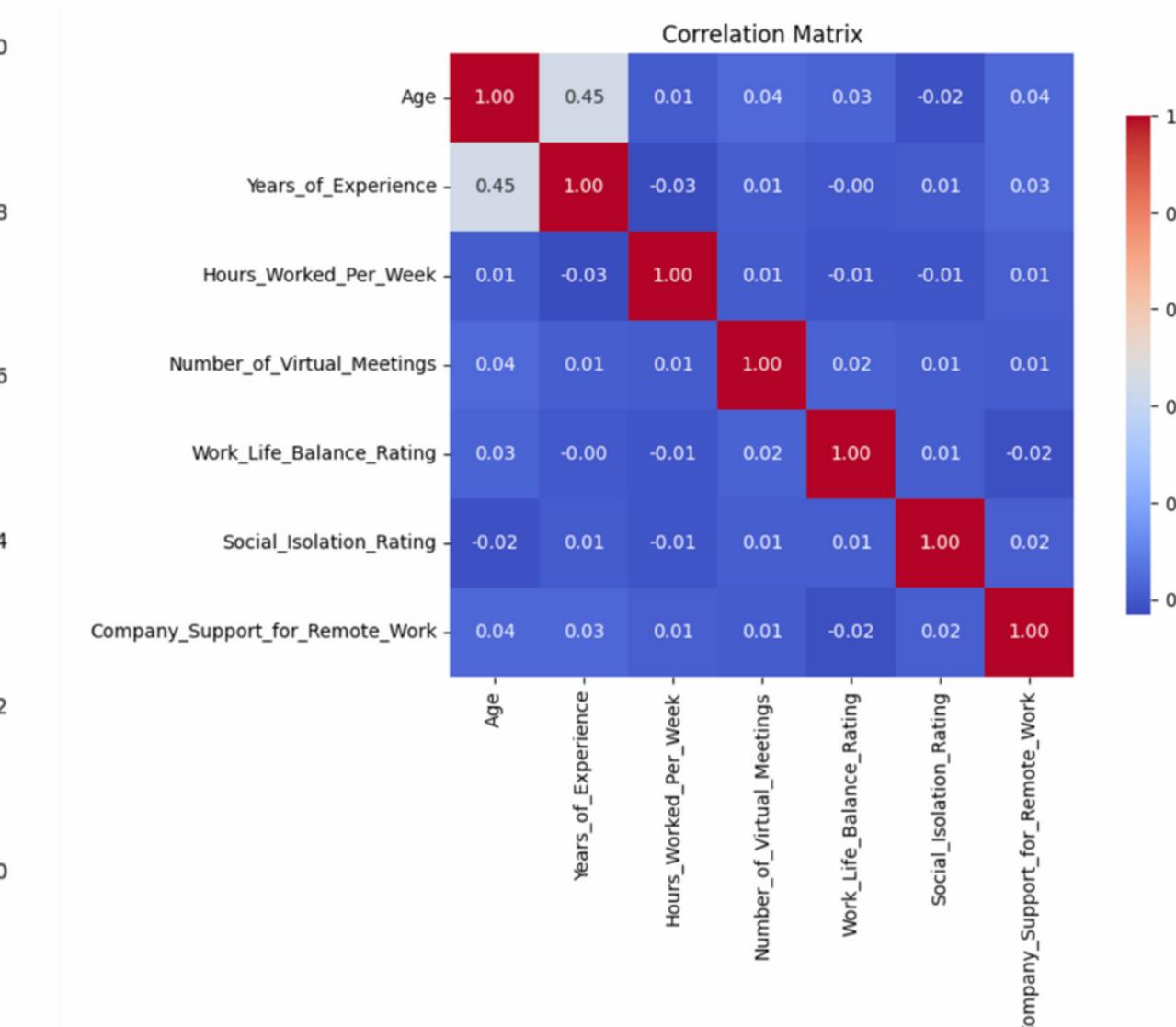
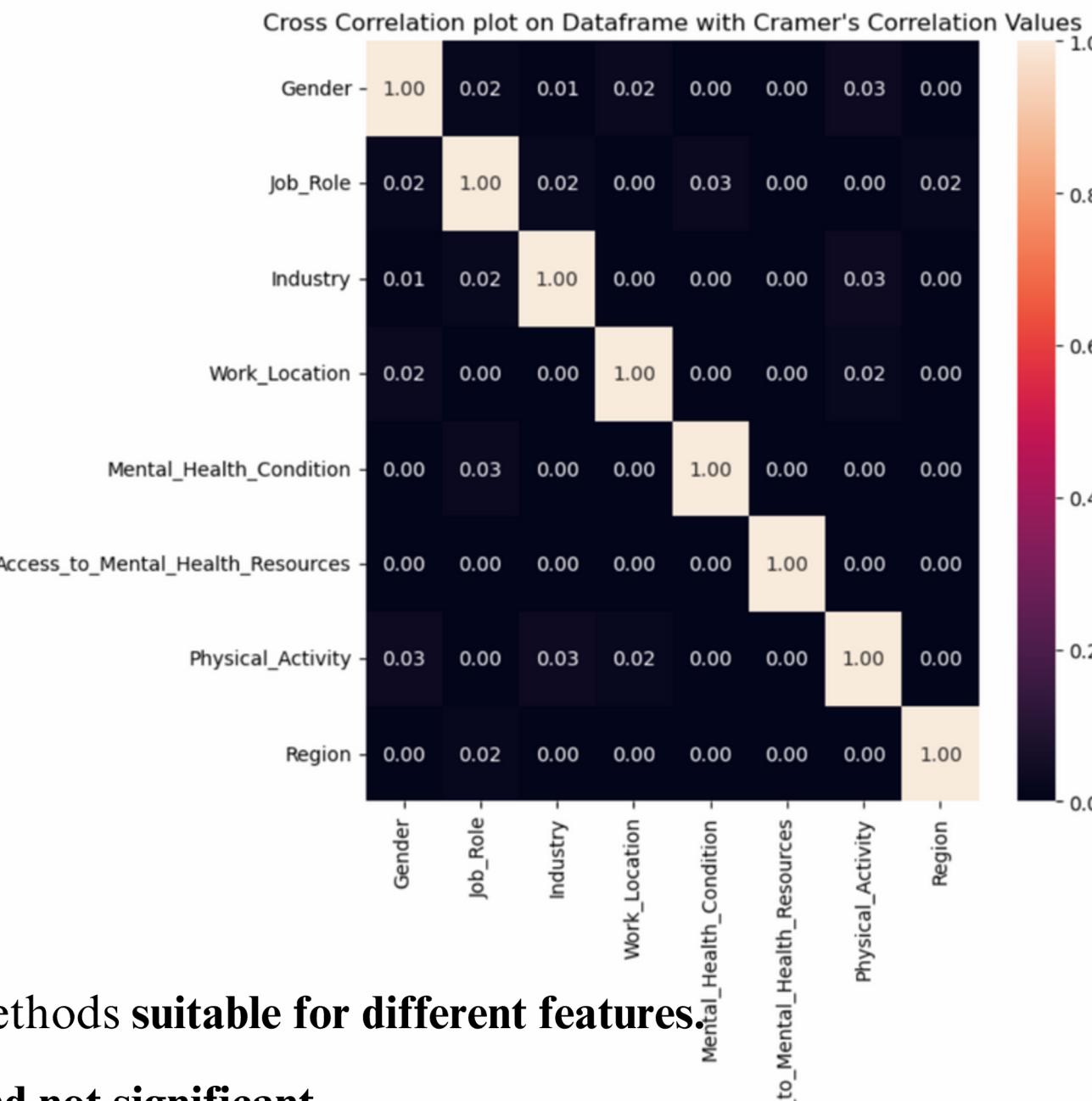
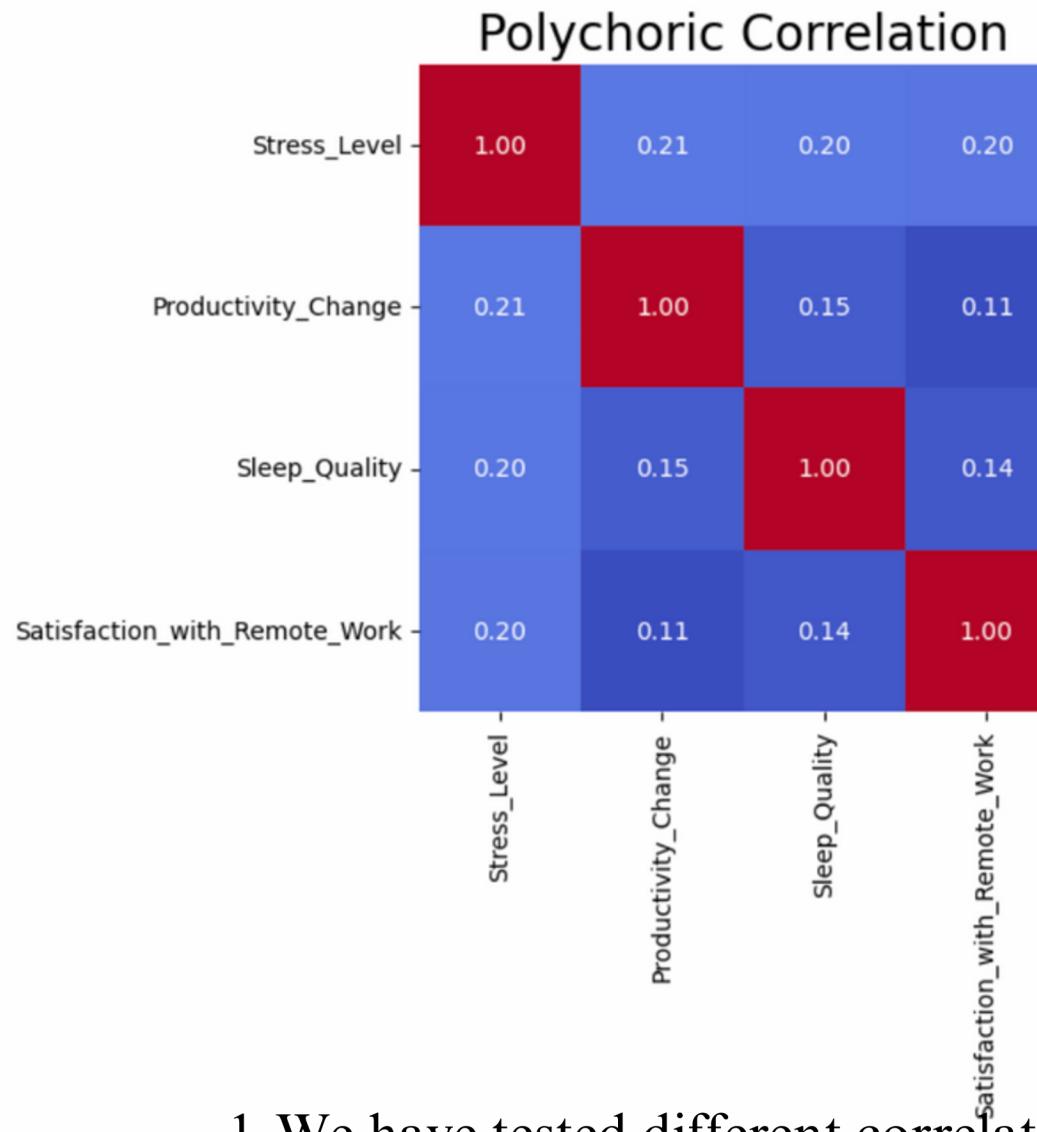
## NUMERICAL

Age  
Years\_of\_Experience  
Hours\_Worked\_Per\_Week  
Number\_of\_Virtual\_Meetings

## Target

Satisfaction\_with\_Remote\_Work

# Correlation



1. We have tested different correlation methods **suitable for different features.**
2. We identified **correlation is very less and not significant.**
3. So for identifying important feature, used random forest after encoding and scaling to identify the importance of each feature, which is itself not significant.
4. We scaled every ordinal values except the target as, we decided to use classifier to find which class they belong to and also, **we need to see that in this classification, we need to see, by how much we are wrong.** For example. if a employee is satisfied with the the remote work, but if we predict unsatisfied it will completely wrong but identifying as Neutral can be better result.

# Feature Importance

Feature	Importance
Hours_Worked_Per_Week	0.089146
Years_of_Experience	0.087274
Age	0.085568
Number_of_Virtual_Meetings	0.075396
Company_Support_for_Remote_Work	0.049098
Social_Isolation_Rating	0.047085
Work_Life_Balance_Rating	0.047075
Stress_Level	0.033506
Productivity_Change	0.03306
Sleep_Quality	0.033035

# Key Points

1. We can see that feature importance is not much but we have chosen top 10 features and trained it using this.
2. After this we have used Weight and bias Sweep method to run different models like **RandomForestClassifier**, **XGBoostClassifier** and **SVC** for classification and identified the best result which we received is little above 33% which supports, that the data is not proper for prediction.
3. Feature\_engineering and Training model scripts will be found in [GITHUB](#)

## Conclusion

1. After overall analysis it can be said that data is balanced artificially and without any dependency
2. With the random data, it is hard to retrieve any valuable insights