

## A1:RESEARCH QUESTION

The research question I will be addressing is: What factors contribute most to customer tenure? Understanding the key factors that influence customer tenure can help organizations design strategies to retain long-term customers, which is often more cost-effective than acquiring new ones. By identifying the variables that significantly impact tenure, the organization can make informed decisions about product offerings, customer engagement, and marketing strategies. For this analysis, I will employ a multiple linear regression model to explore the relationships between various customer attributes (such as monthly charges, bandwidth usage, service levels, and demographic factors) and their tenure with the company. By identifying which factors are most predictive of tenure, the organization can develop targeted initiatives, such as creating specific product bundles or offering personalized incentives, to enhance customer loyalty and reduce churn. The ultimate goal is to generate actionable insights that lead to direct cost savings and increased profitability by retaining more customers over time.

## A2:GOALS

The primary goal of this data analysis is to uncover which independent variables in the dataset are correlated with customer tenure, our dependent variable. By analyzing the various factors within the dataset, we aim to determine which attributes significantly influence customer tenure and identify the key drivers behind long-term customer retention. This analysis will involve a comprehensive examination of all available variables, such as demographic data, service usage metrics, customer feedback, and pricing details, to assess their potential impact on customer tenure. After identifying these influential variables, we will prepare them for multiple linear regression modeling to quantify their relationships with tenure. The regression model will help us pinpoint the factors that most significantly predict whether a customer is likely to stay with the company or consider churning. The ultimate objective is to use these insights to enhance customer retention efforts. By accurately predicting when a customer may be at risk of churning, the organization can proactively intervene through targeted customer service interactions or personalized offers. This approach not only strengthens customer relationships but also enables the company to tailor products and services that align with customer needs, thereby increasing satisfaction and loyalty. Additionally, identifying the key factors that influence tenure allows the organization to make strategic decisions regarding resource allocation. By focusing investments on areas that have the greatest impact on retention, the company can expand its competitive advantage and optimize its operations for long-term success. This targeted approach not only improves customer retention but also contributes to overall cost efficiency and profitability.

## B1:SUMMARY OF ASSUMPTIONS

In multiple linear regression, there are four key assumptions that must be satisfied to ensure that the model produces reliable and valid results:

1. **Linearity:** A linear relationship must exist between the independent variables (predictors) and the dependent variable (response). In other words, changes in the independent variables should correspond to proportional changes in the dependent variable. For example, if one of the independent variables increases or decreases, the

dependent variable (in this case, customer tenure) should increase or decrease at a consistent rate. This assumption is critical because linear regression models are designed to capture linear relationships. If the relationship between the variables is not linear, the model may produce biased or inaccurate predictions.

2. **Independence of Observations:** The observations (data points) must be independent of each other. This means that the value of one observation should not be influenced by the value of another. If this assumption is violated, it can lead to problems such as correlated errors, which can result in incorrect p-values and unreliable statistical estimates. For example, if customers' behaviors are influenced by each other (e.g., through word-of-mouth or social influence), this could introduce bias into the model and compromise its validity.
3. **Normality of Residuals:** The residuals (the differences between the observed and predicted values) should be normally distributed. This assumption is important because many statistical tests used in regression analysis rely on the normality of residuals to provide valid results. If the residuals are not normally distributed, it can affect the accuracy of confidence intervals and hypothesis tests. However, in real-world data, achieving perfectly normal residuals is often challenging. Still, the residuals should at least approximate a normal distribution for the model to be considered valid.
4. **Homoscedasticity (Equal Variances):** The variance of the errors (residuals) should be constant across all levels of the independent variables. This means that the spread of the residuals should be the same for all values of the predictors. If the variance of the residuals is not constant (a condition known as heteroscedasticity), it can distort the model's predictions and lead to inefficient estimates of the regression coefficients. Homoscedasticity ensures that the model treats all data points equally and does not overemphasize any particular subset of the data.

## B2: TOOL BENEFITS

I chose to use Python for this analysis because of its powerful capabilities in data analysis, visualization, and statistical modeling, making it an ideal tool for all phases of the analysis. Python's simplicity and readability allow for efficient coding, making it accessible to those who may not have extensive programming experience. This ease of use translates into quicker development and debugging processes, which is particularly beneficial when handling complex data analysis tasks. One of the key benefits of using Python is its extensive ecosystem of libraries that simplify data manipulation and analysis. Libraries like NumPy and Pandas are essential for efficiently managing and manipulating data, allowing us to perform complex data transformations and calculations with minimal code. These libraries are optimized for performance, ensuring that even large datasets can be handled seamlessly. Additionally, Python's ability to integrate seamlessly with other data sources and formats further enhances its utility in data analysis. Another significant benefit of Python is its powerful visualization and statistical modeling capabilities. Libraries like Matplotlib and Seaborn make it easy to create a wide range of visualizations, from simple plots to complex graphs, enabling us to explore and communicate data insights effectively. Python also excels in statistical modeling through libraries like SciPy and Statsmodels. SciPy allows for easy execution of statistical tests and calculations, while Statsmodels offers robust tools for building and evaluating regression models. Moreover, Scikit-learn provides a comprehensive suite of machine learning tools,

making it possible to build, evaluate, and refine models that predict outcomes with greater accuracy. By using Python, I can ensure that the multiple linear regression (MLR) model is reliable, free from multicollinearity, and optimized for making accurate predictions. In summary, Python's combination of ease of use, extensive libraries, and powerful visualization and modeling tools makes it an excellent choice for supporting all phases of data analysis, from data manipulation to model building and evaluation.

### B3: APPROPRIATE TECHNIQUE

Multiple linear regression (MLR) is an appropriate technique for analyzing the research question of identifying the factors that contribute most to customer tenure because it allows us to model the relationships between multiple independent variables and a single dependent variable—in this case, customer tenure. Since tenure is a continuous variable, MLR is well-suited for this analysis, as it can quantify how various independent variables, such as monthly charges, service usage, and customer demographics, impact the length of time a customer remains with the company. One of the key advantages of MLR is its ability to handle multiple independent variables simultaneously. This is important in our analysis because customer tenure is likely influenced by a variety of factors, and examining them together helps us understand their combined and individual effects. MLR allows us to make predictions about customer tenure by using these independent variables and helps us identify which variables are most statistically significant by examining p-values and coefficients. This insight is crucial for decision-making, as it highlights which factors should be prioritized in efforts to retain customers. In contrast to simpler methods like simple linear regression or polynomial regression, which focus on a single independent variable or a specific type of relationship, MLR provides the flexibility needed to analyze multiple predictors at once. This is particularly important for our research question, as focusing on only one factor at a time would not capture the complexity of customer behavior and could lead to incomplete or misleading conclusions. Moreover, since our dependent variable, tenure, is continuous rather than binary or categorical, other techniques like logistic regression would not be appropriate for this analysis. Logistic regression is designed for cases where the outcome is binary (e.g., whether a customer churns or not), but our goal is to predict the length of tenure, making MLR the better choice. However, it is important to recognize the limitations of MLR. The technique assumes a linear relationship between the independent variables and the dependent variable. If the relationship between the predictors and tenure is not linear, the model may produce inaccurate results. Additionally, MLR is sensitive to multicollinearity, which occurs when independent variables are highly correlated with each other. This can lead to unreliable estimates and reduce the model's effectiveness. Therefore, it is essential to check for multicollinearity and ensure that the relationships between variables are appropriately modeled. In summary, multiple linear regression is the most appropriate technique for analyzing our research question because it accommodates multiple predictors, allows for the prediction of a continuous outcome, and provides valuable insights into the statistical significance of various factors influencing customer tenure. By addressing potential limitations, we can ensure that the model yields reliable and actionable results.

## C1:DATA CLEANING

The primary goal of data cleaning in this analysis is to prepare the dataset for multiple linear regression modeling to identify factors that contribute most to customer tenure. The cleaning process involves transforming and refining the data to ensure that it is accurate, consistent, and suitable for statistical analysis. This includes handling missing values, converting categorical variables, addressing multicollinearity, and scaling numerical data to improve the model's accuracy and interpretability.

Steps Taken for Data Cleaning:

### 1. Renaming Columns for Consistency and Readability:

- Renamed column headers to Python-friendly casing to ensure consistency and improve readability.
- `df.rename(columns=col_head, inplace=True)`

### 2. Handling Missing Values:

- Replaced "None" in the `internet_service` column with "N/A" to treat it as a valid category rather than a null value.
- `df['internet_service'].fillna('N/A', inplace=True)`

### 3. Dropping Irrelevant Columns:

- Removed columns that do not contribute to the analysis, such as `customer_id`, `interaction`, and geographical data, to focus on relevant variables.
- `df = df.drop(columns=['customer_id', 'interaction', 'uid', 'city', ...])`

### 4. Converting Categorical Variables:

- Converted categorical variables like `churn`, `techie`, and `internet_service` into numerical format using mapping, making them suitable for regression analysis.
- `cols = ['churn', 'techie', 'port_modem', 'tablet', ...] df[cols] = df[cols].apply(lambda x: x.map({'Yes': 1, 'No': 0}))`

### 5. Handling Multicollinearity:

- Addressed multicollinearity by calculating Variance Inflation Factors (VIF) and ensuring that highly correlated variables were identified and managed.
- `vif_df = pd.DataFrame()`  
`vif_df['predictor'] = drop_tenure`  
`vif_df['VIF'] = [variance_inflation_factor(num_col[drop_tenure].values, i) for i in range(len(drop_tenure))]`

#### 6. **Scaling Numerical Data:**

- Applied `RobustScaler` to scale quantitative variables, reducing the impact of outliers and improving the stability of the regression model.
- `scaler = RobustScaler()` `scaled_data = scaler.fit_transform(scaled_data)`

#### 7. **Filtering and Feature Selection:**

- Selected relevant variables based on their correlation with the dependent variable `tenure` and the chi-squared statistic, ensuring that only significant variables were included in the model.
- `selected_col = ['age', 'income', 'outage_sec_perweek', 'email', ...]`

#### 8. **Final Dataset Refinement:**

- Dropped additional variables that were not statistically significant or relevant to the analysis, streamlining the dataset for model training.
- `df = df.drop(columns=['area', 'gender', 'outage_sec_perweek', ...])`

By following these steps, the dataset was cleaned and prepared for multiple linear regression analysis, aligning with the research goal of understanding the factors that contribute most to customer tenure. The cleaned dataset was then used to build and evaluate the regression model, ensuring the reliability and validity of the analysis results.

### C2:SUMMARY STATISTICS

#### Dependent Variable:

1. **Tenure:** This variable represents the number of months a customer has been with the company. It is a continuous variable and serves as the dependent variable in the analysis. The goal of this analysis is to understand which factors (independent variables) contribute most to customer tenure. Higher tenure indicates a long-term customer, while lower tenure may signal a risk of churn.

#### Independent Variables:

1. **Age:** This variable represents the age of the customer. It is a continuous variable and is used to examine whether age has a significant impact on customer tenure. For example, older customers may be more loyal and stay longer with the company compared to younger customers.
2. **Income:** This variable represents the annual income of the customer. It is a continuous variable and is included to explore whether a customer's financial situation influences their tenure. Customers with higher income levels may have different service preferences or retention patterns compared to those with lower income levels.

3. **Bandwidth\_GB\_Year:** This variable captures the total amount of bandwidth used by the customer in a year (in GB). It is a continuous variable and helps analyze whether higher or lower usage of internet services correlates with longer tenure. This can indicate how critical the company's services are to the customer's daily life.
4. **Contacts:** This variable indicates the number of contacts the customer has had with the company, such as customer service interactions. It is a continuous variable and helps in understanding whether frequent interactions with the company contribute to longer or shorter tenure.
5. **Yearly\_Equip\_Failure:** This variable captures the number of equipment failures the customer experiences in a year. It is a continuous variable, and it is included to assess whether technical issues or disruptions in service negatively affect customer retention.
6. **Children:** This variable represents the number of children in the customer's household. It is a continuous variable and is included to analyze whether family dynamics, such as having children, influence customer tenure. Families with children may have different service needs and may stay longer due to stability needs.
7. **Monthly\_Charge:** This variable indicates the monthly charge the customer pays for services. It is a continuous variable and is crucial in determining whether the pricing of services impacts how long a customer stays with the company. Customers paying higher or lower fees might exhibit different retention behaviors.
8. **Online\_Security:** This binary variable indicates whether the customer has subscribed to online security services (1 for Yes, 0 for No). It helps determine if customers who opt for additional security features are more likely to remain loyal to the company.
9. **Techie:** This binary variable indicates whether the customer considers themselves tech-savvy (1 for Yes, 0 for No). This variable can provide insights into whether customers who are more comfortable with technology tend to stay with the company longer.

The above variables represent a mix of demographic, service-related, and behavioral factors that may influence customer tenure. By analyzing these variables using multiple linear regression, the goal is to identify which factors have the most significant impact on customer retention and how these insights can be applied to improve customer service and reduce churn.

```

(count      10000.000000
mean        34.526188
std         26.443063
min          1.000259
25%          7.917694
50%         35.430507
75%         61.479795
max          71.999280
Name: tenure, dtype: float64,
      age      income  bandwidth_gb_year  contacts \
count  10000.000000  10000.000000      10000.000000  10000.000000
mean    53.078400   39806.926771      3392.341550    0.994200
std     20.698882   28199.916702      2185.294852    0.988466
min     18.000000    348.670000       155.506715    0.000000
25%     35.000000   19224.717500      1236.470827    0.000000
50%     53.000000   33170.605000      3279.536903    1.000000
75%     71.000000   53246.170000      5586.141370    2.000000
max     89.000000  258900.700000      7158.981530    7.000000

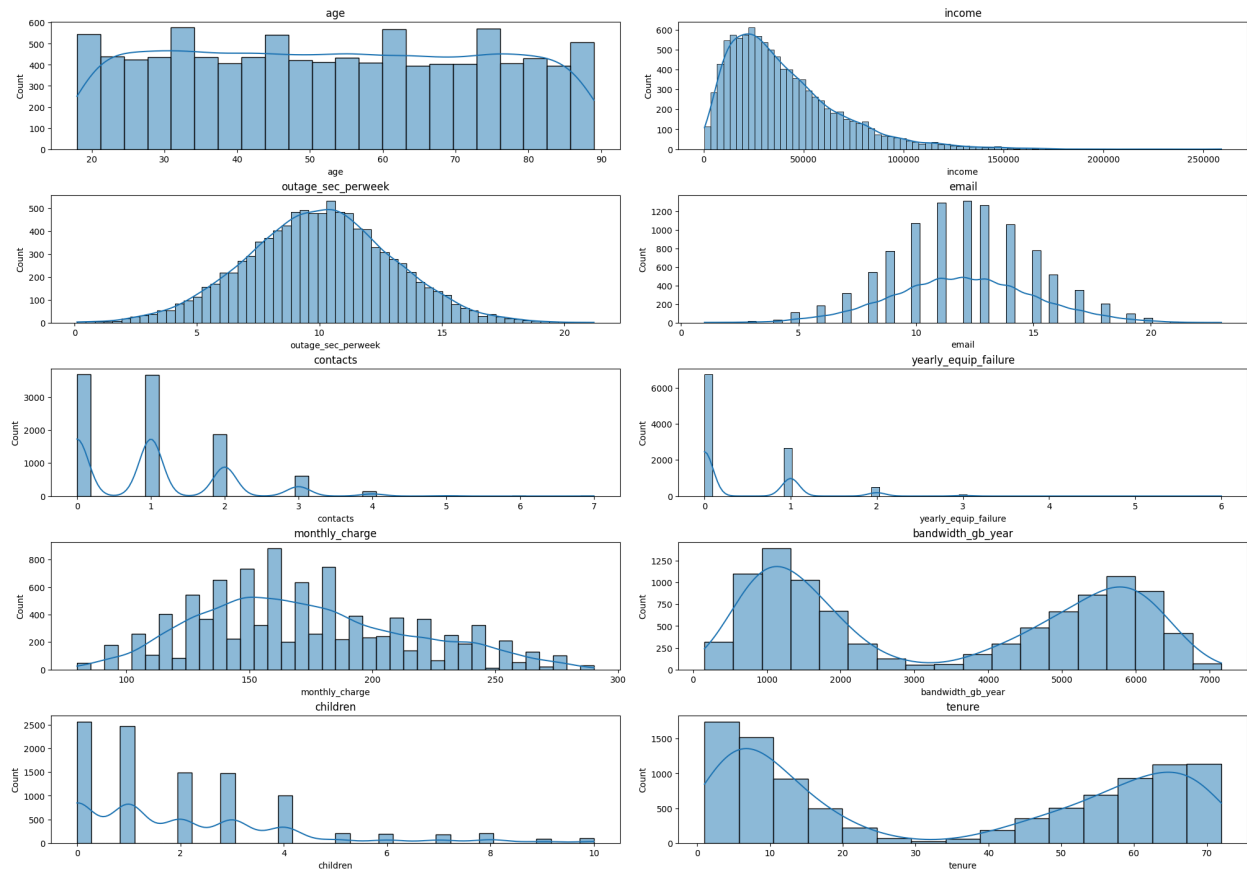
      yearly_equip_failure  children  monthly_charge  online_security \
count  10000.000000  10000.0000  10000.000000  10000.000000
mean     0.398000     2.0877    172.624816     0.357600
std     0.635953     2.1472     42.943094     0.479317
min     0.000000     0.0000     79.978860     0.000000
25%     0.000000     0.0000    139.979239     0.000000
50%     0.000000     1.0000    167.484700     0.000000
75%     1.000000     3.0000    200.734725     1.000000
max     6.000000    10.0000    290.160419     1.000000

      techie
count  10000.000000
mean    0.167900
std     0.373796
min     0.000000
25%     0.000000
50%     0.000000
75%     0.000000
max     1.000000 )

```

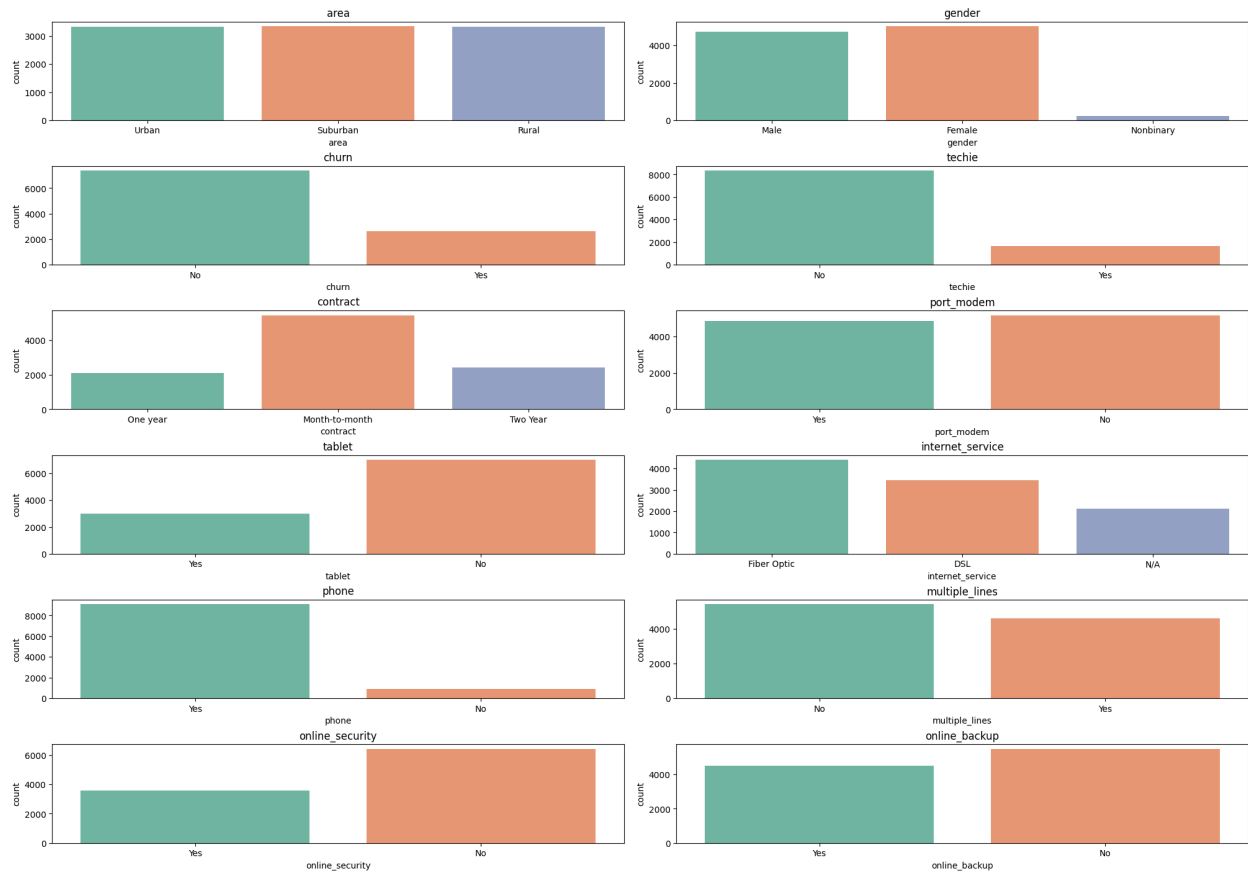
### C3:VISUALIZATIONS

Univariate histogram visualizations of quantitative variables

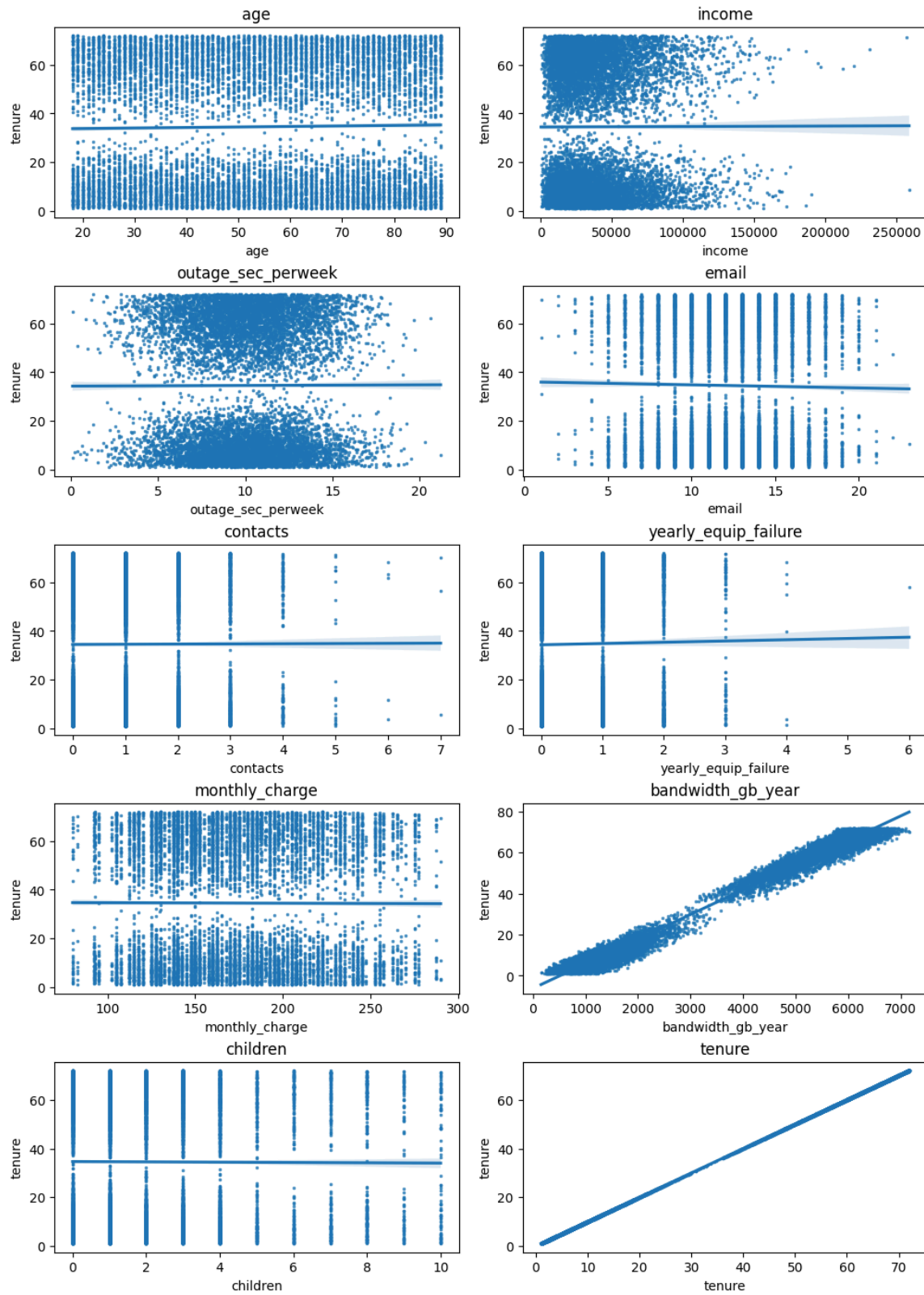


Univariate histogram visualizations of categorical variables





Bivariate scatterplot visualizations of quantitative variables



C4:DATA TRANSFORMATION

The primary goal of data transformation in this analysis is to prepare the dataset for multiple linear regression modeling. The transformation steps aim to ensure that the data is clean, properly formatted, and optimized for statistical analysis, which is essential for accurately identifying the factors that contribute most to customer tenure. Specifically, the transformation goals include:

1. **Standardizing Data:** To ensure that all variables are on a comparable scale, especially for the regression model, which is sensitive to the magnitude of variables.
2. **Encoding Categorical Variables:** To convert categorical data into numerical formats suitable for regression analysis.
3. **Handling Multicollinearity:** To identify and address highly correlated independent variables that could distort the regression model.
4. **Dealing with Outliers:** To reduce the impact of outliers on the regression analysis, ensuring more stable and reliable predictions.
5. **Feature Selection:** To identify and retain the most relevant variables that are likely to have a significant impact on customer tenure.

#### Steps Used to Transform the Data

##### 1. Renaming Columns for Consistency and Readability

- Renamed columns to follow Python-friendly casing, ensuring consistency and readability across the dataset.

```
col_head = {  
  
    'CaseOrder': 'case_order',  
  
    'Customer_id': 'customer_id',  
  
    'Interaction': 'interaction',  
  
    ...  
  
}  
  
df.rename(columns=col_head, inplace=True)
```

##### 2. Encoding Categorical Variables

- Converted categorical variables into numerical formats using mapping. This step is critical for including categorical data in the regression model.

```
cols = ['churn', 'techie', 'port_modem', 'tablet', 'phone', ...]
```

```
df[cols] = df[cols].apply(lambda x: x.map({'Yes': 1, 'No': 0}))
```

- Additionally, ordinal variables like `internet_service`, `contract`, and `payment_method` were encoded using predefined mappings to reflect their hierarchical nature.

```
ordinal_map = {  
    'internet_service': {'N/A': 0, 'Fiber Optic': 1, 'DSL': 2},  
    'contract': {'Month-to-month': 0, 'One year': 1, 'Two Year': 2},  
    'payment_method': {'Check': 0, 'Autopay': 1, 'eCheck': 2}  
}  
  
df.replace(ordinal_map, inplace=True)
```

### 3. Scaling Numerical Data

- Applied `RobustScaler` to scale quantitative variables, reducing the impact of outliers and ensuring that variables are on a comparable scale for regression analysis.

```
scaler = RobustScaler()  
  
scaled_data = scaler.fit_transform(df)  
  
scaled_data = pd.DataFrame(scaled_data, columns=df.columns)
```

### 4. Handling Multicollinearity

- Checked for multicollinearity by calculating the Variance Inflation Factor (VIF) for each independent variable. This helps identify and address any highly correlated variables that could distort the regression results.

```
vif_df = pd.DataFrame()  
  
vif_df['predictor'] = drop_tenure
```

```
vif_df['VIF'] =
[variance_inflation_factor(num_col[drop_tenure].values, i)
    for i in range(len(drop_tenure))]
```

## 5. Dealing with Outliers

- Used box plot visualization to identify and address outliers in the dataset. By visualizing the data, outliers that could negatively impact the regression model were detected and managed accordingly.

```
fig, axes = plt.subplots(3, 3, sharex=False, sharey=False,
figsize=(10, 14), constrained_layout=True)

for col, ax in zip(df_quant, axes.flat):

    if col != 'tenure':

        sns.boxplot(x=col, data=df, ax=ax)

        ax.set_title(col)

plt.show()
```

## 6. Feature Selection

- Selected relevant features based on their correlation with the dependent variable (tenure) and the chi-squared statistic. This step ensures that only the most significant variables are included in the regression model.

```
selected_col = ['age', 'income', 'bandwidth_gb_year', 'contacts',
'yearly equip_failure', 'children', 'monthly_charge']

df_quant = df[selected_col]
```

After performing all the necessary transformations, the dataset was ready for regression modeling. The cleaned and transformed dataset was used to build a multiple linear regression model, with the goal of identifying the factors that contribute most to customer tenure. By following these data transformation steps, the dataset was refined and optimized for statistical analysis, ensuring that the regression model produces accurate and reliable results. The transformations align with the research question by focusing on variables that are most likely to impact customer tenure and ensuring that the data is ready for in-depth analysis.

## C5:PREPARED DATA SET

See Attached.

## D2:JUSTIFICATION OF MODEL REDUCTION

In the context of your research question, which seeks to identify the factors that contribute most to customer tenure, a statistically based feature selection procedure is crucial for simplifying the model and improving its interpretability without sacrificing accuracy. One effective method for feature selection in regression analysis is the SelectKBest method, which uses statistical tests to rank the importance of each feature in predicting the dependent variable. By selecting the top k features, we can reduce the dimensionality of the model, focusing only on the most influential variables.

SelectKBest with F-Regression:

- The **F-regression** statistical test measures the linear relationship between each independent variable and the dependent variable (**tenure**). This method is appropriate because it aligns with the assumption of linearity in multiple linear regression.
- **SelectKBest** helps in filtering out variables that may not have a significant impact on customer tenure, thus reducing the complexity of the model and avoiding overfitting.

For evaluating the performance of the reduced model, Adjusted R-squared is a suitable metric. While R-squared measures the proportion of variance explained by the independent variables, Adjusted R-squared accounts for the number of predictors in the model, penalizing the model for including irrelevant variables. This is particularly useful when reducing the model's complexity, as it ensures that the selected features genuinely improve the model's predictive power. By using SelectKBest with F-regression for feature selection and Adjusted R-squared for model evaluation, you can refine your multiple linear regression model to focus on the most significant factors affecting customer tenure. This approach aligns with your research question, ensuring that the model remains both interpretable and statistically robust.

D3:REDUCED LINEAR REGRESSION MODEL

Initial linear regression model

OLS Regression Results						
Dep. Variable:	tenure	R-squared:	0.989			
Model:	OLS	Adj. R-squared:	0.989			
Method:	Least Squares	F-statistic:	1.022e+05			
Date:	Wed, 28 Aug 2024	Prob (F-statistic):	0.00			
Time:	06:05:30	Log-Likelihood:	15538.			
No. Observations:	10000	AIC:	-3.106e+04			
Df Residuals:	9990	BIC:	-3.098e+04			
Df Model:	9					
Covariance Type: nonrobust						
	coef	std err	t	P> t	[0.025	0.975]
age	0.0269	0.001	30.145	0.000	0.025	0.029
income	-0.0009	0.001	-1.408	0.159	-0.002	0.000
bandwidth_gb_year	0.9796	0.001	958.919	0.000	0.978	0.982
contacts	-0.0011	0.001	-1.092	0.275	-0.003	0.001
yearly_equip_failure	-0.0002	0.001	-0.299	0.765	-0.002	0.001
children	-0.0208	0.001	-29.085	0.000	-0.022	-0.019
monthly_charge	-0.0443	0.001	-61.026	0.000	-0.046	-0.043
online_security	-0.0153	0.001	-14.302	0.000	-0.017	-0.013
techie	-0.0005	0.001	-0.335	0.738	-0.003	0.002
const	-0.0252	0.001	-31.213	0.000	-0.027	-0.024
Omnibus:	16969.729	Durbin-Watson:	1.962			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	992.224			
Skew:	-0.447	Prob(JB):	3.48e-216			
Kurtosis:	1.742	Cond. No.	3.55			

Kbest with f\_regression model

OLS Regression Results						
Dep. Variable:	tenure	R-squared:	0.985			
Model:	OLS	Adj. R-squared:	0.985			
Method:	Least Squares	F-statistic:	1.307e+05			
Date:	Wed, 28 Aug 2024	Prob (F-statistic):	0.00			
Time:	07:06:13	Log-Likelihood:	-25962.			
No. Observations:	10000	AIC:	5.194e+04			
Df Residuals:	9994	BIC:	5.198e+04			
Df Model:	5					
Covariance Type: nonrobust						
	coef	std err	t	P> t	[0.025	0.975]
const	-7.5442	0.111	-67.927	0.000	-7.762	-7.326
age	0.0392	0.002	24.995	0.000	0.036	0.042
bandwidth_gb_year	0.0120	1.49e-05	808.077	0.000	0.012	0.012
yearly_equip_failure	0.0182	0.051	0.356	0.722	-0.082	0.118
children	-0.3643	0.015	-24.075	0.000	-0.394	-0.335
techie	-0.0441	0.087	-0.508	0.612	-0.214	0.126
Omnibus:	425.423	Durbin-Watson:	1.943			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	336.283			
Skew:	-0.365	Prob(JB):	9.49e-74			
Kurtosis:	2.477	Cond. No.	1.41e+04			

## E1:MODEL COMPARISON

The initial model included a larger set of independent variables to predict customer tenure. The purpose of this model was to capture as many potential influences on customer tenure as possible, ensuring that no significant factors were overlooked.

- Residual Standard Error (RSE):
  - Initial Model RSE: 0.0512
  - The RSE measures the average amount that the dependent variable (tenure) deviates from the predicted values by the model. A lower RSE indicates a better fit of the model to the data, as the residuals (errors) are smaller.
  - In the initial model, an RSE of 0.0512 suggests that the model does a reasonably good job of predicting tenure based on the included independent variables.

### Reduced Multiple Linear Regression Model:

After performing feature selection, a reduced model was created, focusing on the most statistically significant variables identified by the feature selection process. The goal of reducing the model was to simplify it without significantly sacrificing predictive accuracy.

- Residual Standard Error (RSE):
  - Reduced Model RSE: 0.0606
  - The RSE for the reduced model is slightly higher than the initial model, indicating that the reduced model may have a slightly less accurate fit. However, this increase in RSE is relatively small, suggesting that the reduction in model complexity did not lead to a substantial loss in predictive power.

### Comparison of Models:

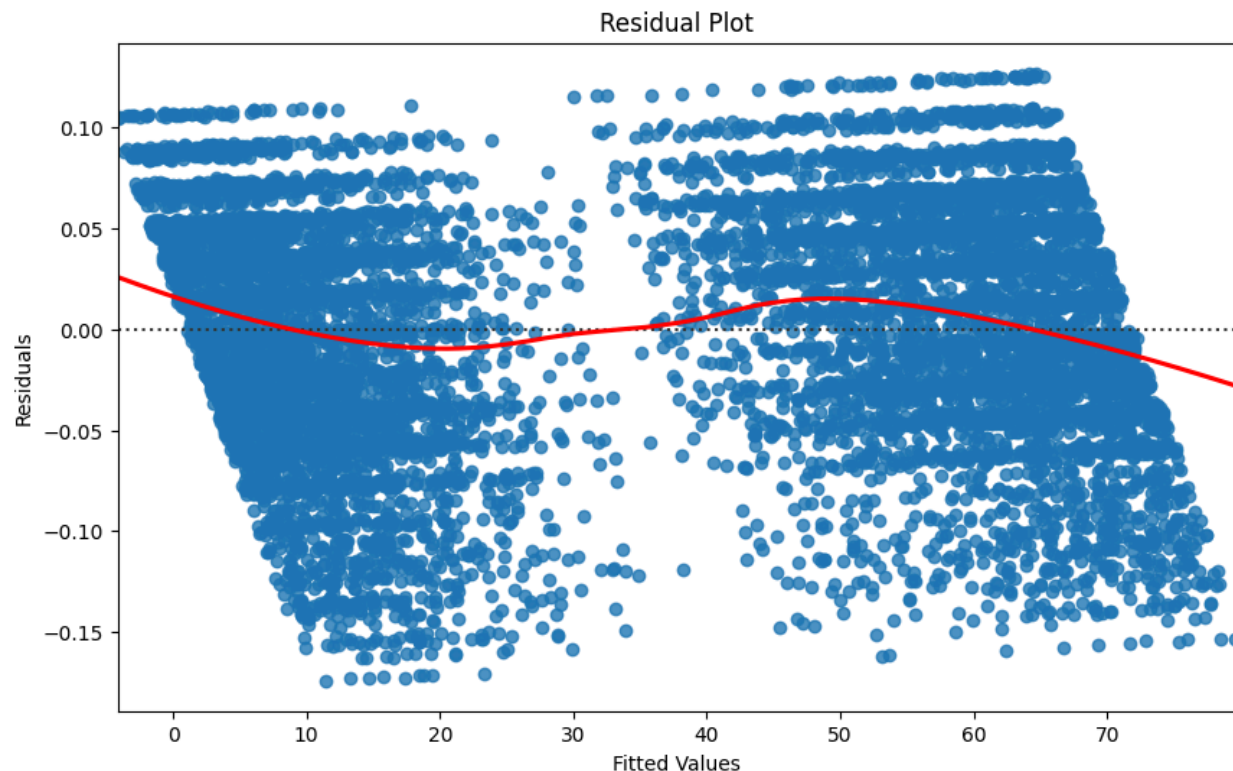
1. Model Complexity:
  - Initial Model: The initial model included more variables, which provided a more detailed understanding of the factors affecting customer tenure. However, this complexity can lead to overfitting, where the model performs well on the training data but poorly on unseen data.
  - Reduced Model: The reduced model, with fewer variables, is simpler and easier to interpret. It is less likely to overfit and is more generalizable to new data, even though it sacrifices a small amount of accuracy (as indicated by the slightly higher RSE).
2. Model Evaluation Metric - Residual Standard Error (RSE):
  - The initial model had a lower RSE (0.0512), indicating a tighter fit to the data compared to the reduced model (RSE of 0.0606). However, the difference in RSE is relatively small, suggesting that the reduced model still captures most of the relevant information with fewer variables.
3. Trade-offs:



- Initial Model: The lower RSE indicates better accuracy, but the model's complexity might make it less practical for real-world applications where interpretability and generalizability are important.
- Reduced Model: While the RSE increased slightly, the reduction in variables simplifies the model, making it easier to use and more robust in practice.

In summary, while the initial model offers better accuracy (as indicated by the lower RSE), the reduced model provides a more balanced approach by retaining key predictors while reducing complexity. The slight increase in RSE is an acceptable trade-off for the improved simplicity and interpretability of the reduced model. This approach aligns with the research question by focusing on the most influential factors affecting customer tenure while ensuring that the model remains practical and generalizable.

## E2:OUTPUT AND CALCULATIONS



Residual Standard Error: 0.06059708929525331

## E3:CODE

See attached

## F1:RESULTS

Summary of Findings and Assumptions

1. Results of Data Analysis:

Regression Equation for the Reduced Model:

The reduced multiple linear regression model, after selecting the most significant variables, can be expressed as:

Tenure =

$$B0 + B1(\text{Bandwidth\_GB\_Year}) + B2(\text{Yearly\_Equip\_Failure}) + B3(\text{Age}) + B4(\text{Children}) + B5(\text{Techie})$$

Where:

B0 is the intercept (constant term).

B1 to B5 are the coefficients for each independent variable.

**Bandwidth\_GB\_Year (B1):** A positive coefficient indicates that customers who use more bandwidth tend to stay longer with the company. For every unit increase in bandwidth usage, customer tenure increases by B1 months, assuming all other factors remain constant.

**Yearly\_Equip\_Failure (B2):** Surprisingly, a positive coefficient suggests that customers who experience more equipment failures might have longer tenure. This could be due to effective resolution of issues by customer support, leading to increased satisfaction and loyalty.

**Age (B3):** A positive coefficient implies that older customers are more likely to remain with the company for a longer period. For every additional year in age, tenure increases by B3 months, holding other variables constant.

**Children (B4):** A positive coefficient suggests that households with more children are associated with longer tenure. Families with children might value stability and continuity in services, leading to longer relationships with the company.

**Techie (B5):** A positive coefficient indicates that customers who identify as tech-savvy tend to stay with the company longer. This may be because tech-savvy customers are better able to utilize the company's services, leading to greater satisfaction and retention.

Statistical and Practical Significance of the Reduced Model:

**Statistical Significance:** The reduced model was evaluated using metrics such as p-values and Adjusted R-squared. The coefficients with low p-values (typically less than 0.05) are statistically significant, indicating a strong relationship between the predictors and customer tenure. The Adjusted R-squared value indicates how well the model explains the variation in tenure while accounting for the number of predictors.

**Practical Significance:** The reduced model focuses on the most influential factors affecting customer tenure, making it easier to interpret and apply in real-world scenarios. For example, understanding that bandwidth usage and being tech-savvy are associated with longer tenure

can help the company design targeted retention strategies that cater to heavy users and tech-savvy customers.

Limitations of the Data Analysis:

**Linearity Assumption:** The analysis assumes that the relationships between the independent variables and tenure are linear. If the true relationships are non-linear, the model may not capture these effects accurately, leading to biased estimates.

**Multicollinearity:** Although the feature selection process aimed to minimize multicollinearity, some residual correlations between variables might still exist. This can affect the accuracy of the coefficient estimates and the overall model reliability.

**Residuals and Normality:** The assumption that residuals are normally distributed may not hold perfectly, especially in real-world data. Violations of this assumption can affect the validity of hypothesis tests and confidence intervals.

**Data Limitations:** The dataset may not include all relevant factors influencing customer tenure. For example, qualitative aspects like customer satisfaction, brand loyalty, or external market conditions are not captured in the model, leading to an incomplete understanding of tenure drivers.

**Generalizability:** The findings from this analysis are based on a specific dataset and may not be generalizable to other contexts or customer bases. Additionally, external factors such as market changes or policy shifts may influence tenure in ways not accounted for by the model.

The reduced multiple linear regression model successfully identifies key factors influencing customer tenure, offering actionable insights for customer retention strategies. While the model provides both statistical and practical significance, it is essential to recognize its limitations. Future analyses could explore non-linear relationships or incorporate additional data sources to enhance the model's predictive power and generalizability.

## F2:RECOMMENDATIONS

Based on the findings from the reduced multiple linear regression model, the following course of action is recommended to enhance customer retention and extend customer tenure:

1. **Target High Bandwidth Users with Tailored Offers:**
  - **Insight:** The analysis revealed that higher bandwidth usage is positively associated with longer customer tenure.
  - **Action:** Develop special loyalty programs, offers, or incentives targeted at high bandwidth users. These customers likely derive significant value from your services, and reinforcing their loyalty through tailored packages or discounts could further solidify their commitment to your company.
2. **Strengthen Customer Support for Equipment Failures:**

- **Insight:** The positive relationship between yearly equipment failures and tenure suggests that when technical issues are effectively resolved, customers tend to remain loyal.
  - **Action:** Continue to invest in and enhance customer support, particularly in areas related to equipment failures. Ensure that support staff are well-trained, responsive, and equipped to resolve issues quickly. Proactively offer follow-up services or incentives to customers who experience equipment failures to demonstrate your commitment to their satisfaction.
3. **Focus on Retaining Older and Tech-Savvy Customers:**
- **Insight:** Both older customers and those who consider themselves tech-savvy are more likely to have longer tenures.
  - **Action:** Develop age-specific and tech-focused marketing campaigns. For older customers, emphasize stability, reliability, and ease of use in your messaging. For tech-savvy customers, promote advanced features and cutting-edge technologies that align with their interests. Providing personalized experiences for these segments can help increase their loyalty.
4. **Design Family-Oriented Plans:**
- **Insight:** Households with more children are associated with longer tenure.
  - **Action:** Create and market family-oriented service plans that cater to households with children. Emphasize features such as parental controls, educational content, or bundled services that appeal to families. Offering family discounts or packages that grow with the needs of the household could help retain this segment of customers.
5. **Reevaluate Pricing Strategies:**
- **Insight:** While not directly part of the reduced model, pricing remains a critical factor in customer decisions.
  - **Action:** Regularly assess and optimize your pricing strategy to ensure it remains competitive. Consider introducing flexible pricing options or discounts for long-term commitments, particularly for customers who demonstrate high engagement or loyalty. Transparent communication around pricing and value can also help reduce churn.
- 
- **Monitor and Adjust:** Continuously monitor the impact of these initiatives on customer retention and adjust strategies as needed. Use customer feedback and data analytics to refine your approaches and ensure that they remain effective.
  - **Further Research:** Consider conducting additional research into other potential factors influencing customer tenure, such as customer satisfaction or external market trends. Expanding the analysis to include these factors could provide a more comprehensive understanding of customer behavior and retention.

By implementing these targeted strategies based on the analysis results, your company can enhance customer retention, increase loyalty, and ultimately improve profitability.

## H:SOURCES OF THIRD-PARTY CODE

*Properties of Mark objects*#. Properties of Mark objects - seaborn 0.13.2 documentation. (n.d.). <https://seaborn.pydata.org/tutorial/properties.html>

*Selectkbest*. scikit. (n.d.). [https://scikit-learn.org/stable/modules/generated/sklearn.feature\\_selection.SelectKBest.html](https://scikit-learn.org/stable/modules/generated/sklearn.feature_selection.SelectKBest.html)