

# Churn Analysis

Daniel, López Coto

---



## Main Goals



Customer retention



Getting insights from data analysis





# Univariate Variable exploration

**Number of features:** 99

Numeric: 78

Categorical: 21

**Target variable:** Churn

1: Leave

0: Stays

Company stop billing about  
\$ 2,9 million/month due to churn.

**Missing values:** 43 features.

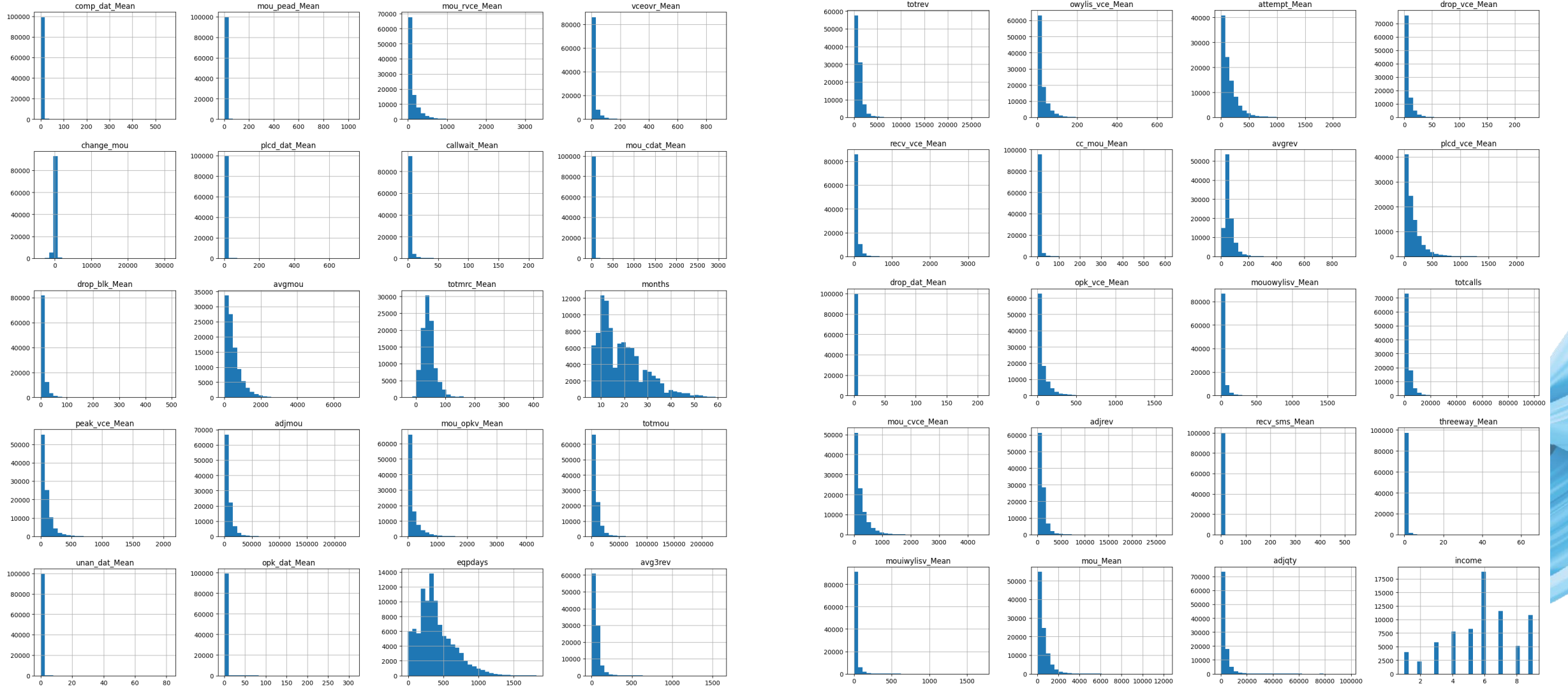
10 features with > 10% of missing values

numbcars	0.49366
dwlsize	0.38308
HHstatin	0.37923
ownrent	0.33706
dwltype	0.31909
lor	0.30190
income	0.25436
adults	0.23019
infobase	0.22079
hnd_webcap	0.10189



# Univariate Variable exploration

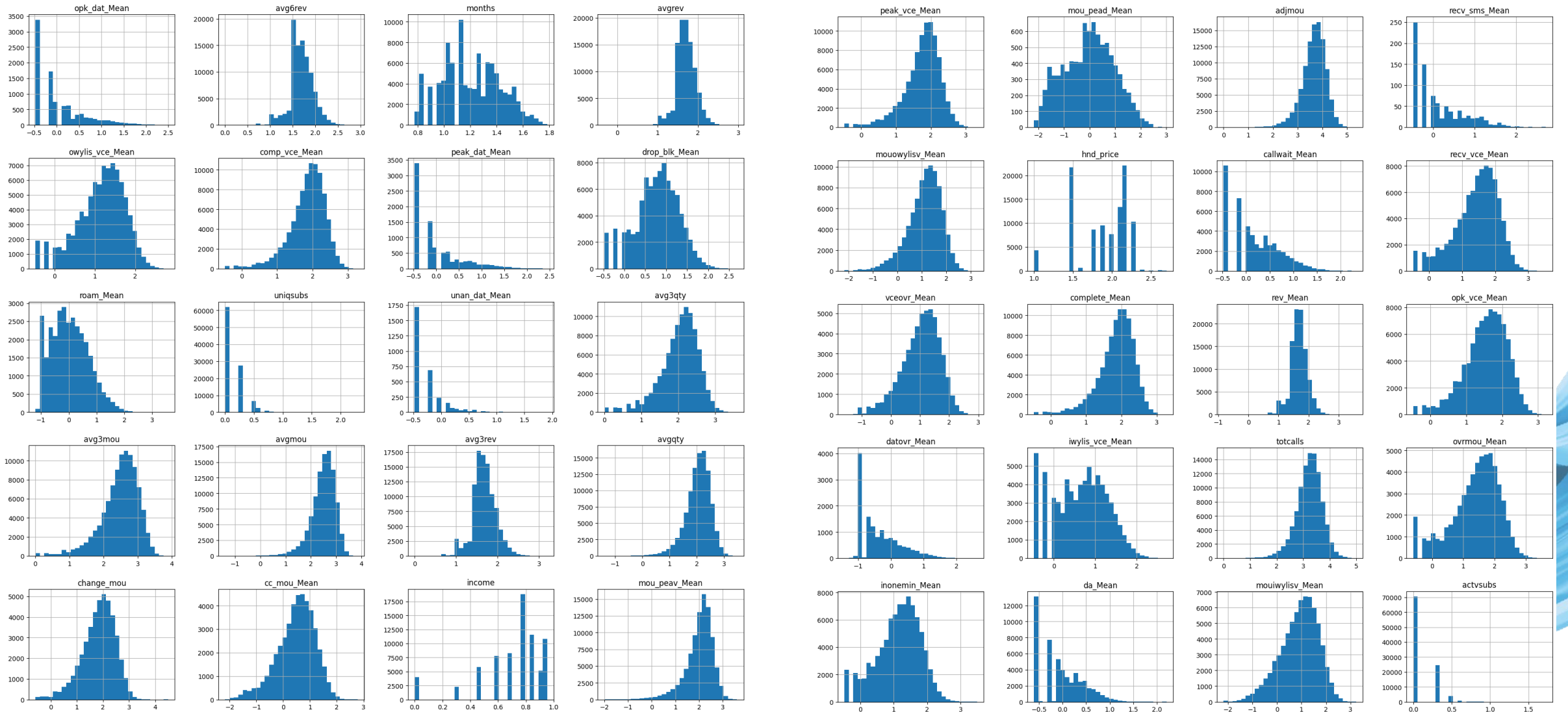
## Distribution shapes





# Univariate Variable exploration

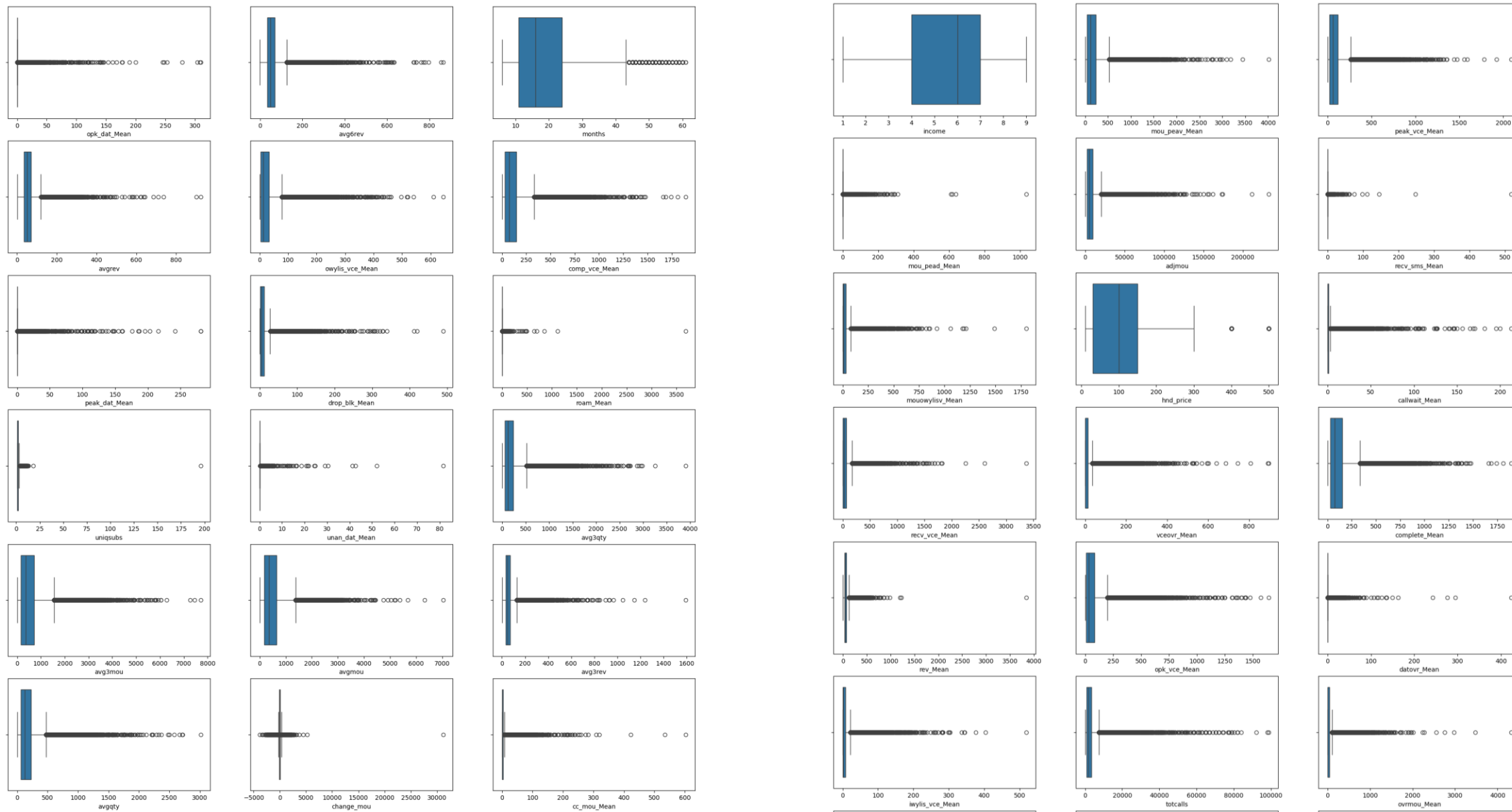
Distribution shapes (log transformed)





# Univariate Variable exploration

## Outlier inspection

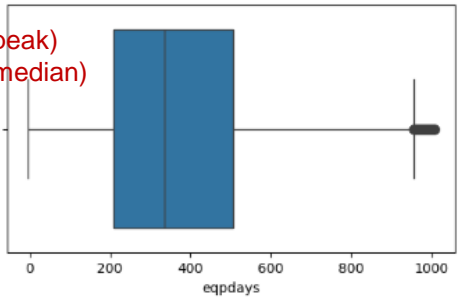
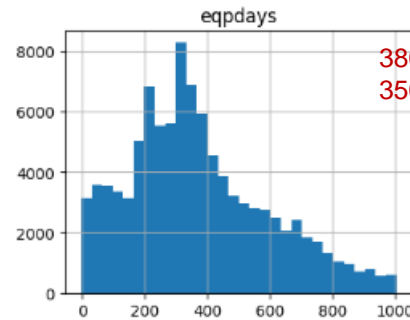
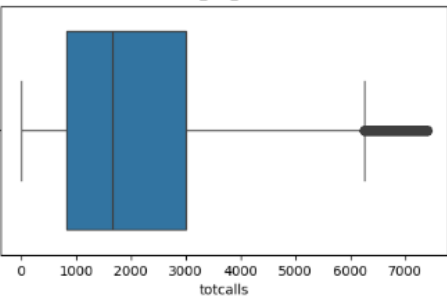
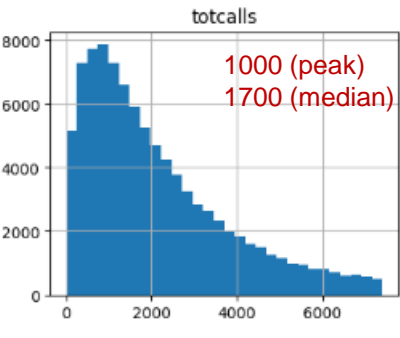
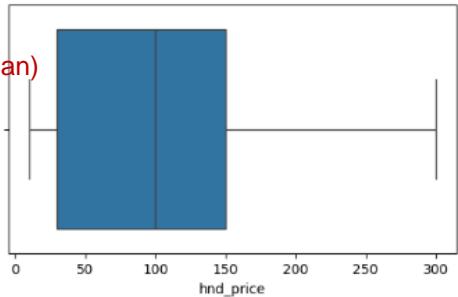
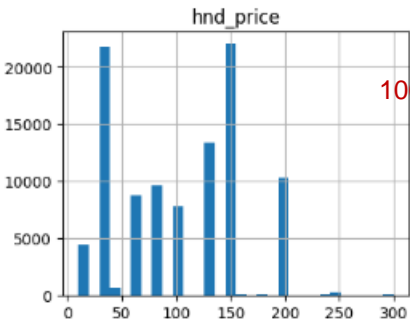
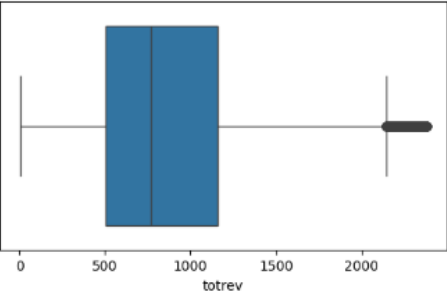
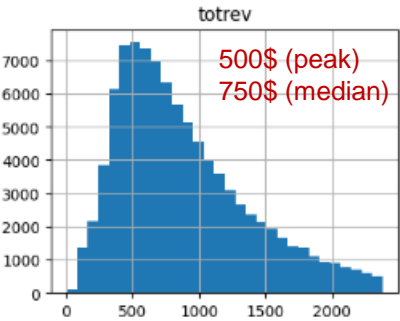
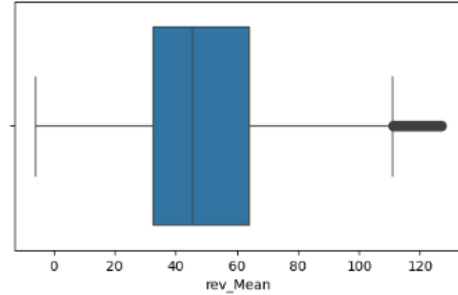
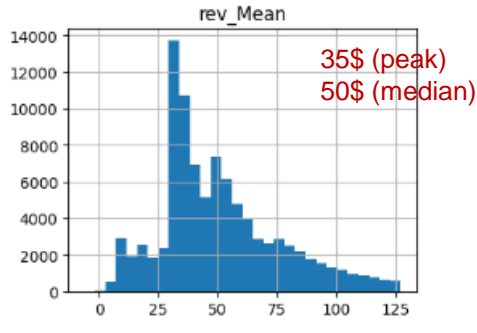
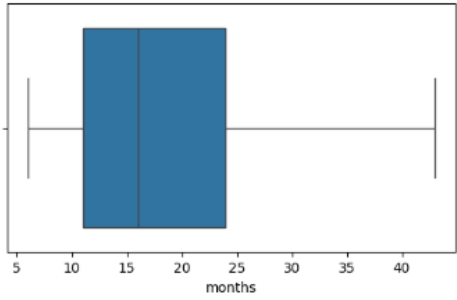
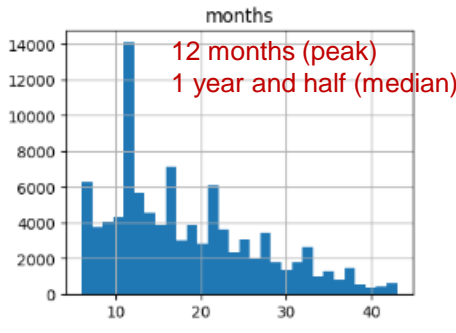






# Univariate Variable exploration

After removing outliers





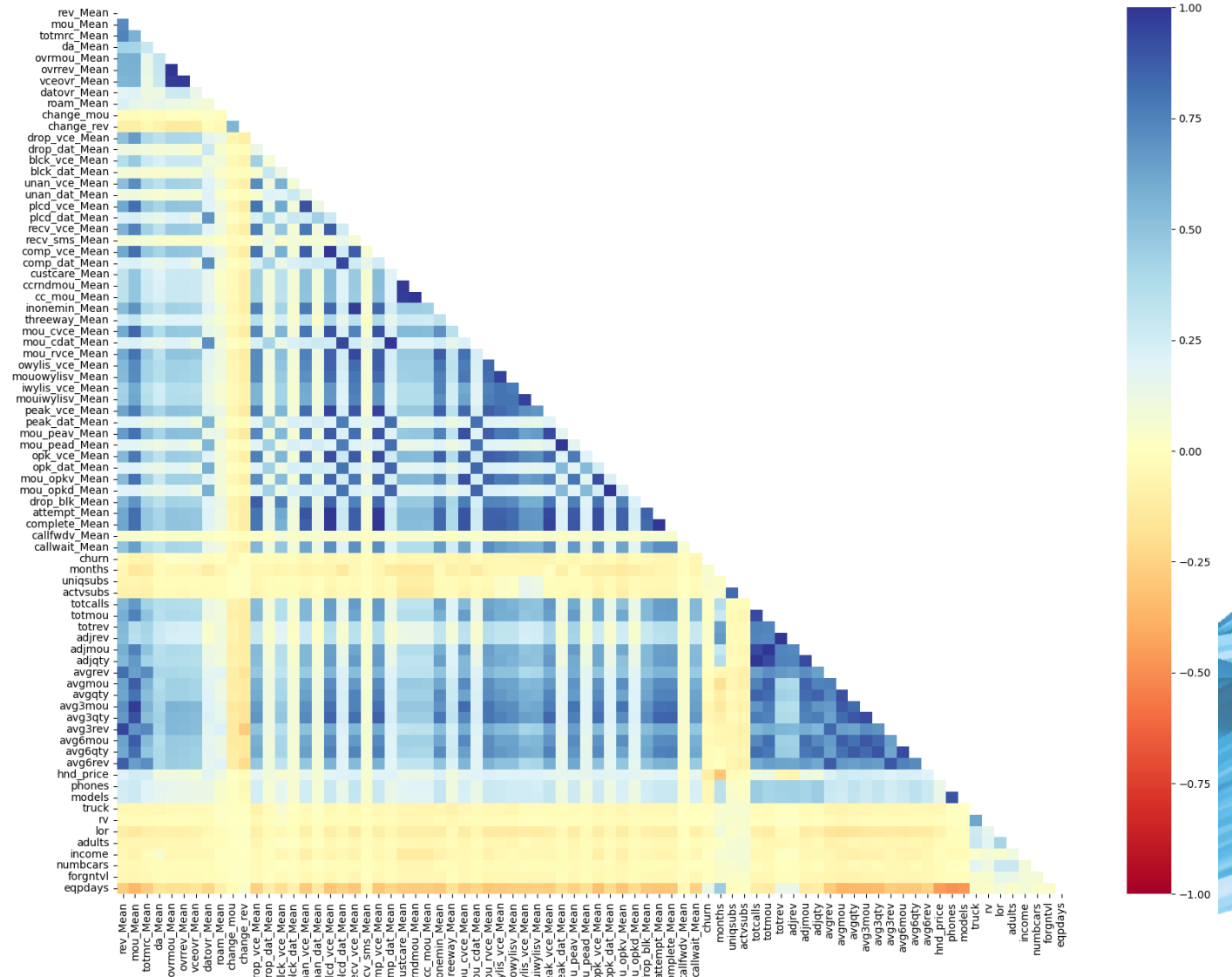
# Univariate Variable exploration

Correlation coefficient

Lot of features highly correlated.

Features related with calls  
(minutes of use, revenue, voice-data calls, etc)  
represent same type of information.

This indicate that these variables have to be  
removed or transformed.



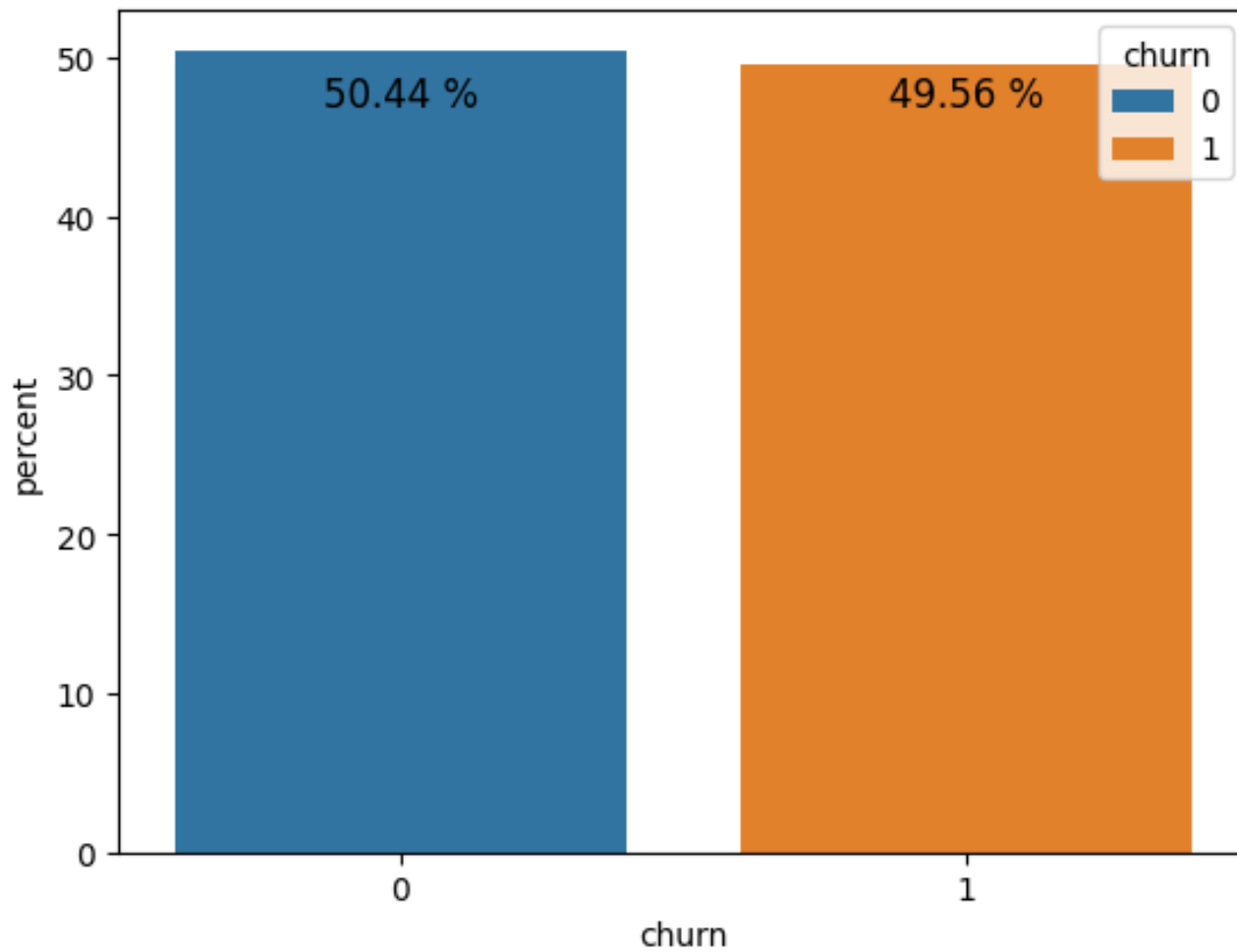




# Univariate Variable exploration

Target variable: Churn

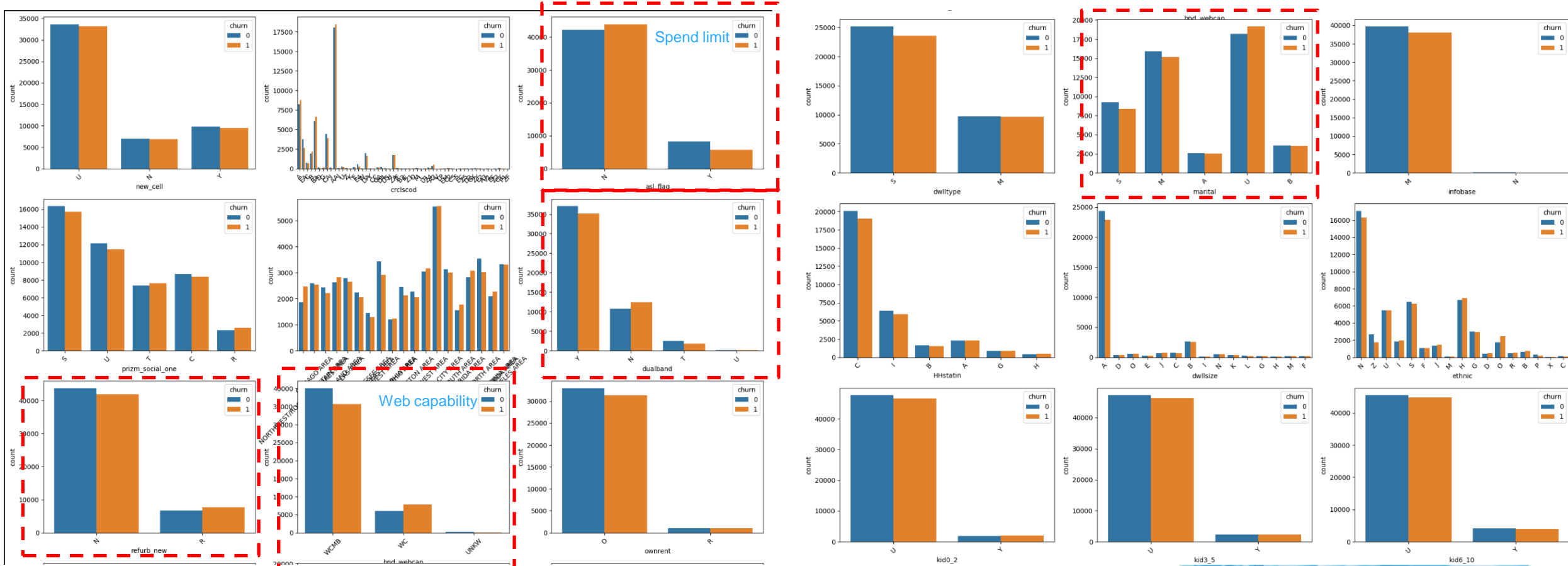
Churn ratio: Not imbalanced target





# Churn-features relationship

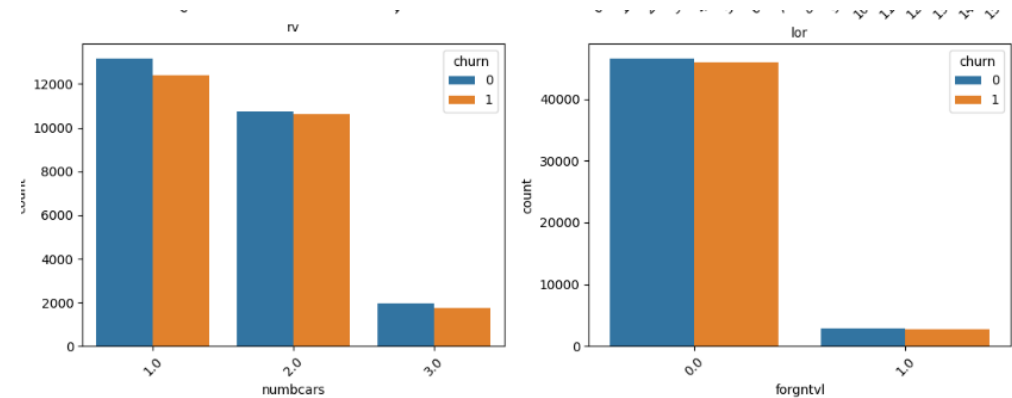
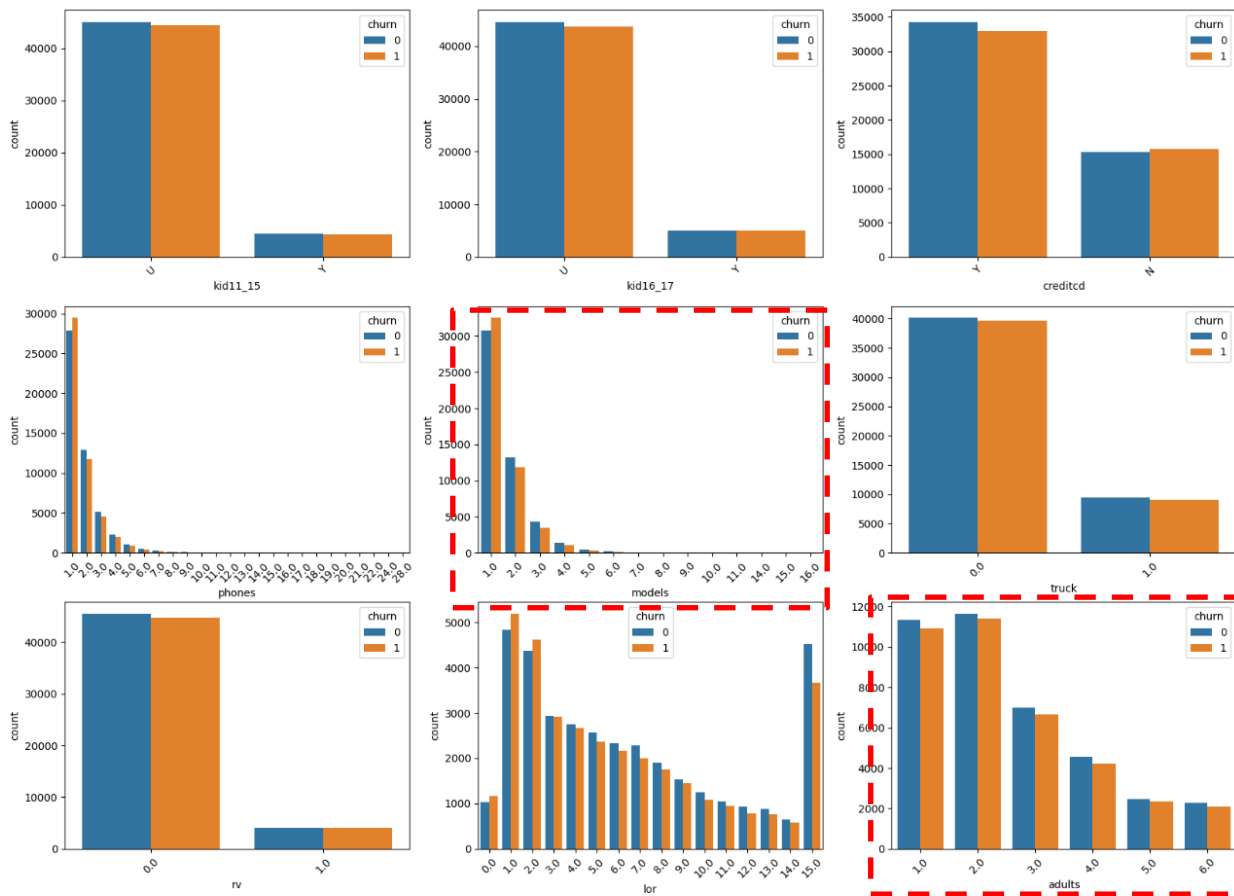
## Categorical variables





# Churn-features relationship

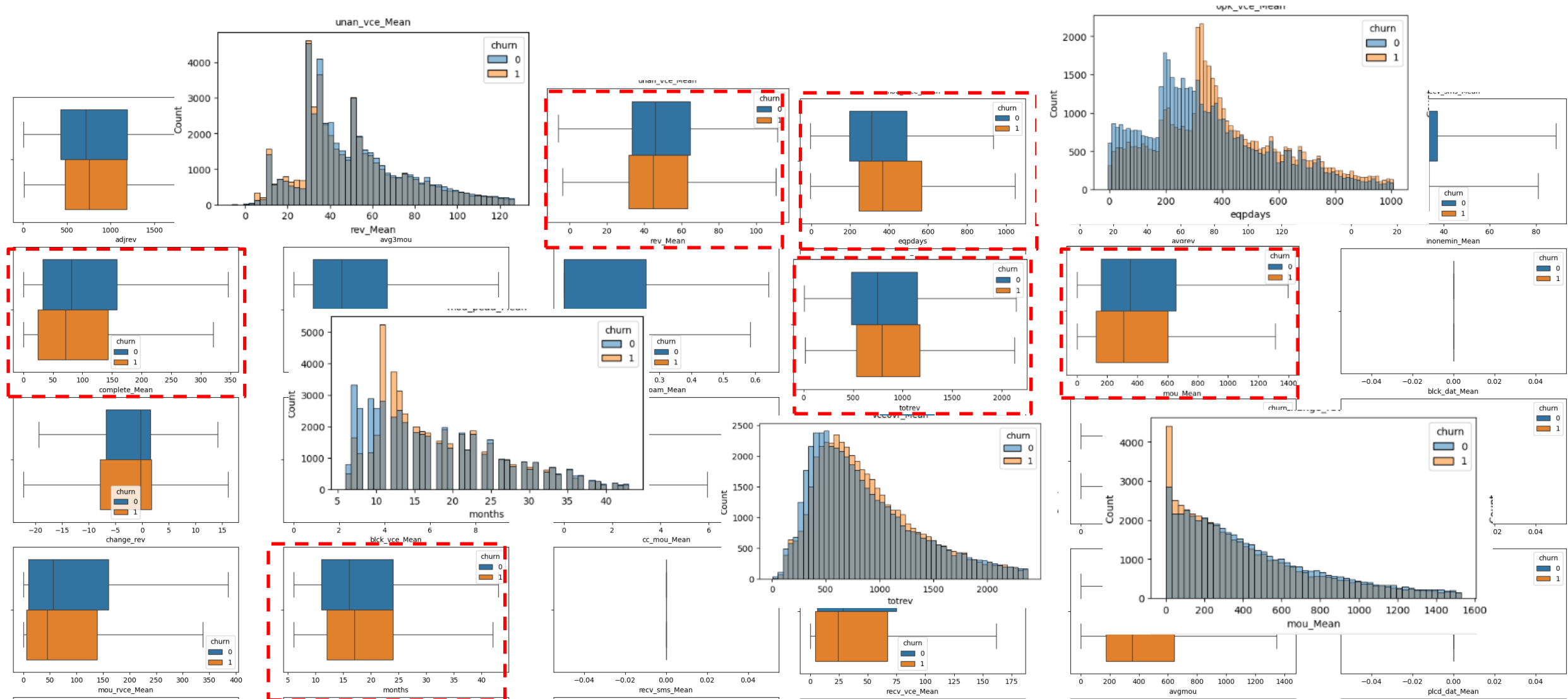
## Categorical variables





# Churn-features relationship

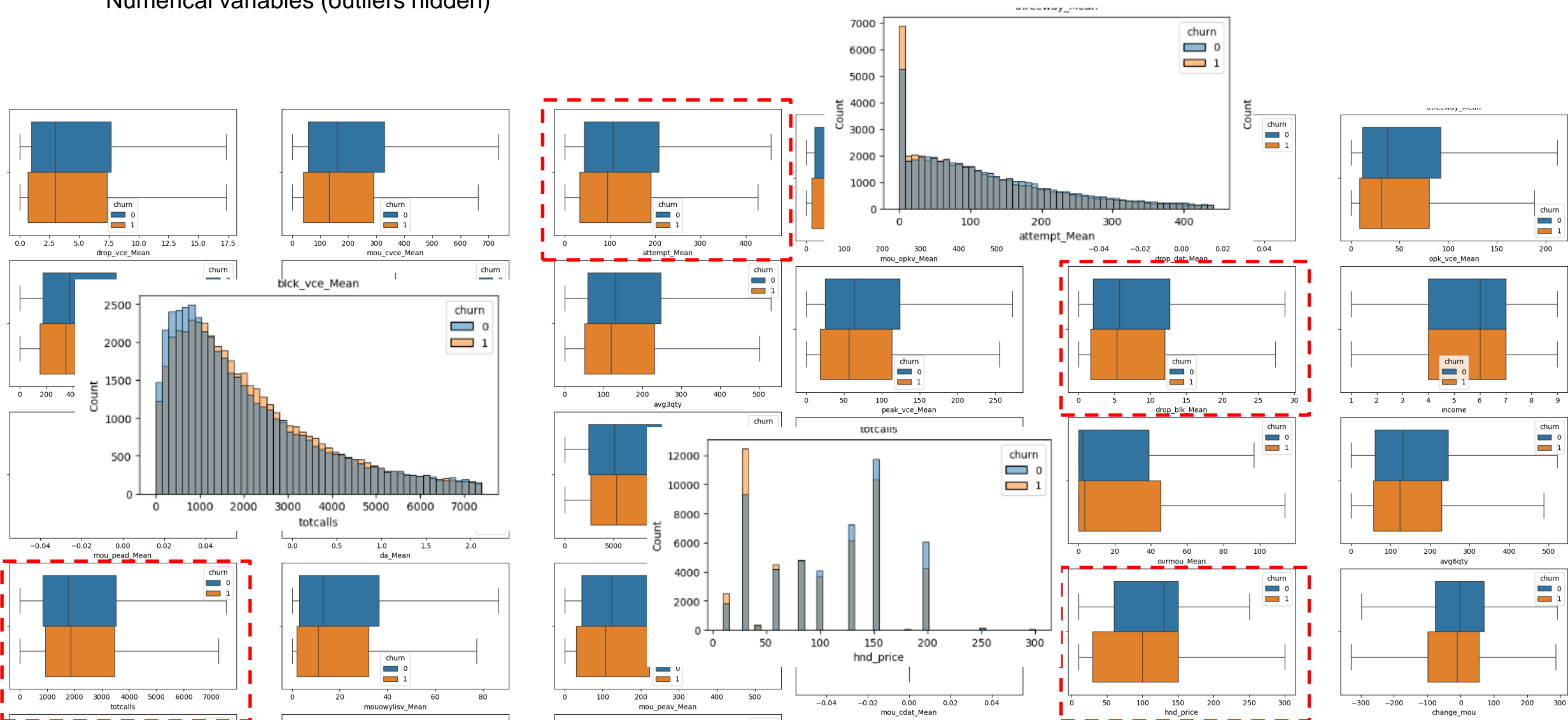
Numerical variables (outliers hidden)





# Churn-features relationship

Numerical variables (outliers hidden)





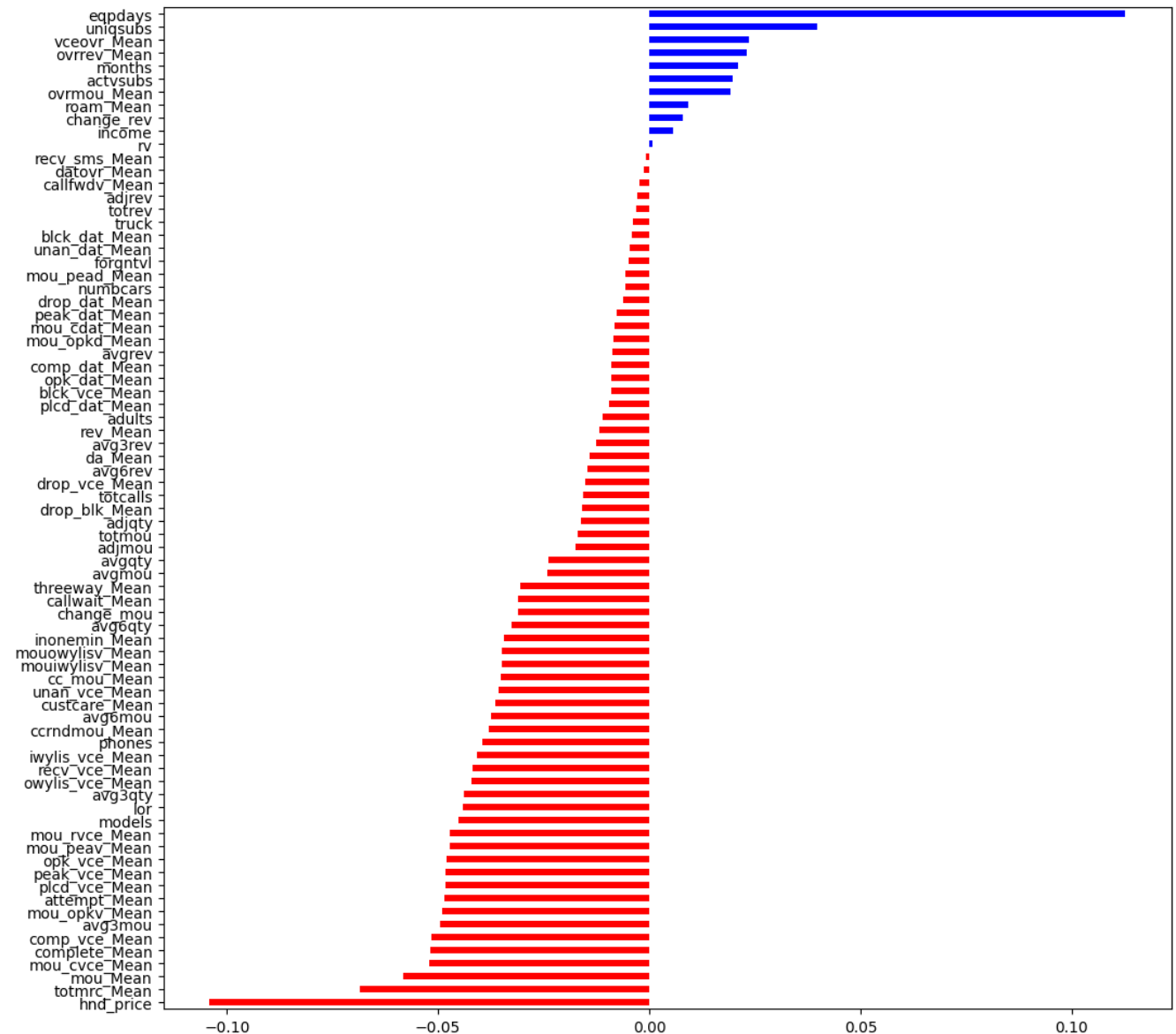
# Churn-features relationship

Point Biserial Correlation factor.

$$r_{pb} = \frac{\overline{Y_1} - \overline{Y_0}}{s_y} \sqrt{\frac{N_0 N_1}{N(N-1)}}$$

Used to check the correlation between the numeric features and the target variable.

Coefficients > 0 means correlation with binary state 1.  
Coefficients < 0 means correlation with binary state 0.





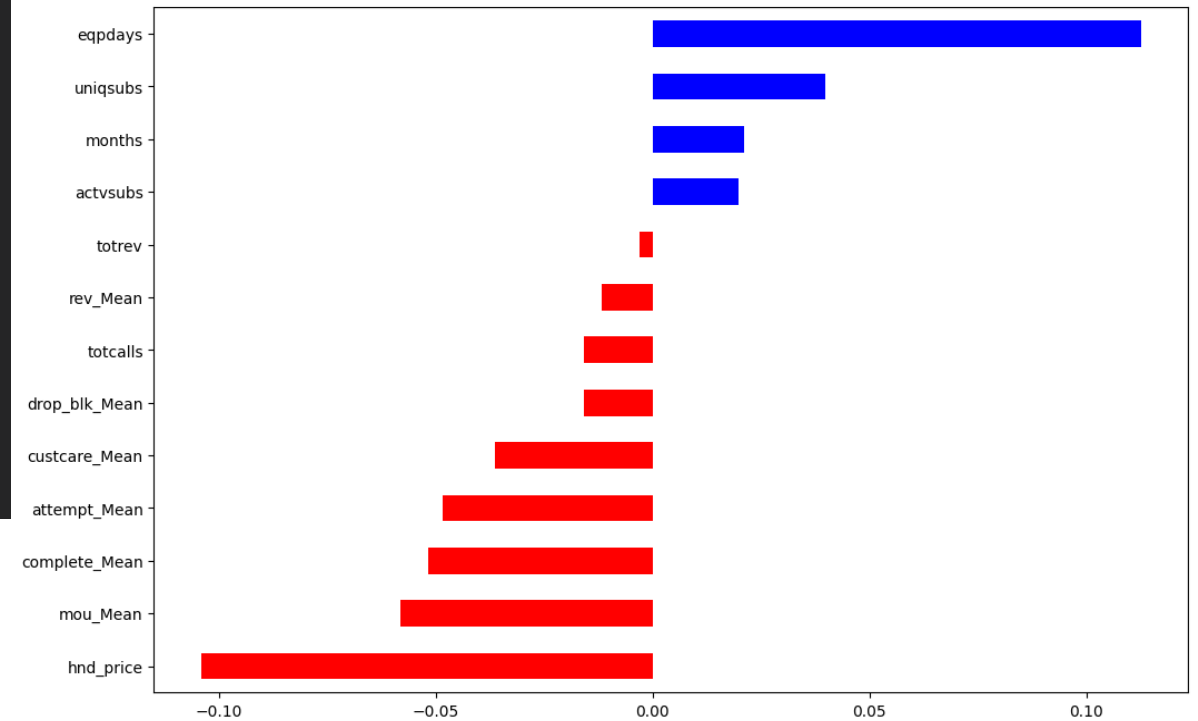


# Features selected

From correlation analysis some features were removed.

```
feat_num_to_keep = ['rev_Mean', 'eqpdays', 'drop_blk_Mean', 'uniqusubs', 'actvsubs',  
                    'custcare_Mean',  
                    'complete_Mean',  
                    'attempt_Mean',  
                    'hnd_price',  
                    'months',  
                    'mou_Mean', 'totcalls', 'totrev']  
  
feat_cat_to_keep = ['asl_flag',  
                    'dualband',  
                    'hnd_webcap',  
                    'refurb_new',  
                    'marital',  
                    'models',  
                    'adults',  
                    'kid0_2',  
                    'kid3_5', 'kid6_10', 'kid11_15', 'kid16_17']
```

Point Biserial Correlation factor





# Machine Learning

Classification problem.

Target variable identified (churn).

Features selected from EDA.

Next steps:

- Split dataset into train-test.
- Clean dataset (outliers and missing values).
- Scaling features.
- Encoding categorical variables.
- Select model and train it.
- Model evaluation.





# Machine Learning

Split dataset into train-test.

```
1 target_column = "churn"
2 df_target = data[target_column]
3 df_features = data[feat_num_to_keep + feat_cat_to_keep]
4
5 X_train, X_test, y_train, y_test = train_test_split(
6     df_features, df_target, stratify=df_target, test_size=0.2, random_state=42)
```

```
cols_to_drop, null_ratio = detect_missing_values(X_train, threshold=0.25)
imputer = drop_and_impute_missing(X_train, cols_to_drop=[], set_type="train")
drop_and_impute_missing(X_test, cols_to_drop=[], set_type="test", imputer=imputer)
```

Clean dataset

```
replace_outliers(X_train)
replace_outliers(X_test)
```

```
1 def replace_outliers(df, whiskers=1.5):
2     """Replace outliers with median value
3
4     Args:
5         df (_type_): Input data frame with numeric values only
6         whiskers (float, optional): Threshold to compute the upper and lower bounds. Defaults to 1.5.
7
8     Returns:
9         Data Frame: With outliers replaced.
10    """
11    # Make sure the selected values are the numeric ones
12    df_numeric = df.select_dtypes('number')
13    # Compute the quantiles for each column
14    quant = df_numeric.quantile(q=[0.75, 0.25])
15    # Compute the IQR and the upper and lower bounds used to consider outliers
16    iqr = quant.iloc[0] - quant.iloc[1]
17
18    up_bound = quant.iloc[0] + (whiskers*iqr)
19    low_bound = quant.iloc[1] - (whiskers*iqr)
20
21    # Replace values above upper bound with median
22    df_numeric = df_numeric.apply(lambda x: x.mask(
23        x > up_bound.loc[x.name], np.nan), axis=0)
24
25    # Replace values below lower bound with median
26    df_numeric = df_numeric.apply(lambda x: x.mask(
27        x < low_bound.loc[x.name], np.nan), axis=0)
28
29    # Replace the numeric columns in the data set
30    df[df_numeric.columns] = df_numeric
31
32    return df
```

```
1 def detect_missing_values(df, threshold=0.1):
2     """Detect missing values and return columns to drop (if pass the threshold) and the ratio of the missing values.
3
4     Args:
5         df (DataFrame): Data Frame
6         threshold (float, optional): Threshold to drop columns. Defaults to 0.1.
7
8     Returns:
9         tuple: Tuple containing columns to drop and their ratio
10    """
11    nrows = len(df)
12    # Retrieves the columns with null values
13    col_w_nulls = df.columns[df.isna().any()]
14    # Ratio of null values
15    null_ratio = df[col_w_nulls].isna().sum(
16        ).sort_values(ascending=False) / nrows
17    # Columns with null values greater than the threshold
18    cols_to_drop = list(null_ratio.index[null_ratio > threshold])
19
20    return cols_to_drop, null_ratio
```

```
1 def drop_and_impute_missing(df, cols_to_drop=[], set_type="train", imputer=None):
2     """Drop columns based on a list of column names, and impute values based on the imputer passed (if any)
3     and depending on if the df is the test or training set.
4
5     Args:
6         df (DataFrame): Data Frame
7         cols_to_drop (list, optional): List of column names to drop. Defaults to [].
8         imputer (object, optional): Imputer object to impute missing values. Defaults to None.
9         set_type (str, optional): check if the df is the test or train set. Options: test, train.
10
11     Returns:
12         DataFrame: DataFrame after dropping columns and imputing missing values
13    """
14    if cols_to_drop:
15        df.drop(columns=cols_to_drop, axis=1, inplace=True)
16
17    if set_type.lower() == "train":
18        if imputer is None:
19            imputer = SimpleImputer(strategy="most_frequent")
20            df[df.columns] = imputer.fit_transform(df)
21            df[df.columns] = df[df.columns].infer_objects()
22            # print(df.info())
23            # print(imputer)
24        elif set_type.lower() == "test":
25            if imputer is None:
26                raise ValueError(
27                    "An imputer object is required for imputing missing values in the test set.")
28            df[df.columns] = imputer.transform(df)
29            df[df.columns] = df[df.columns].infer_objects()
30            # print(df.info())
31    else:
32        raise ValueError(
33            "Invalid set_type. Allowed options are 'train' or 'test'.")
34
35    return df
```



# Machine Learning

Split dataset into train-test.

```
X_train_num_scaled, scaler_obj = my_scaler(X_train.select_dtypes('number'), set_type="train", scaler_type="standard")
X_test_num_scaled, _ = my_scaler(X_test.select_dtypes('number'), set_type="test", scaler_chosen=scaler_obj)

X_train[X_train_num_scaled.columns] = X_train_num_scaled
X_test[X_test_num_scaled.columns] = X_test_num_scaled
```

Encode categorical features

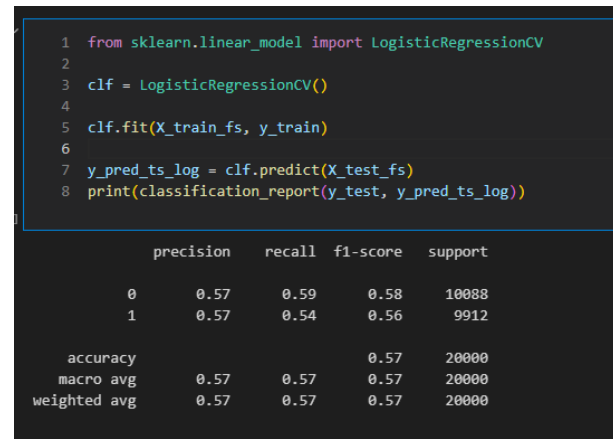
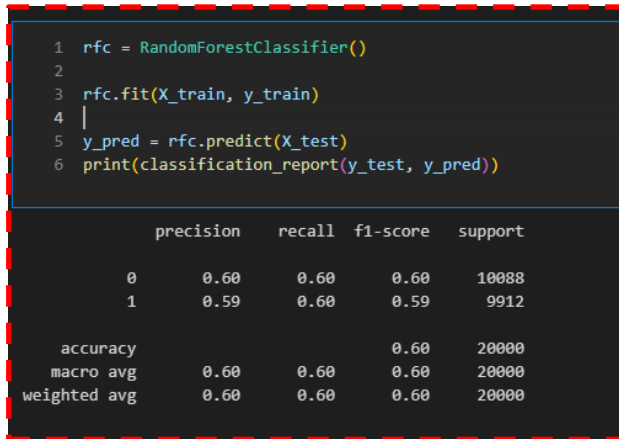
```
1 from sklearn.preprocessing import OneHotEncoder
2
3 cat_var = X_train.select_dtypes('object').columns #feat_cat_to_keep
4
5 cat_encoder = OneHotEncoder(handle_unknown='ignore')
6 feat_encoded = cat_encoder.fit_transform(X_train[cat_var]).toarray()
7
8 feature_labels = cat_encoder.get_feature_names_out()
9 X_train_cat = pd.DataFrame(data=feat_encoded, columns=feature_labels)
10 X_train_cat.index = X_train.index
11 X_train = pd.concat([X_train, X_train_cat], axis=1)
12 X_train.drop(columns=cat_var, axis=1, inplace=True)
13
14 feature_test_encoded = cat_encoder.transform(X_test[cat_var]).toarray()
15 X_test_cat = pd.DataFrame(data=feature_test_encoded, columns=cat_encoder.get_feature_names_out())
16 X_test_cat.index = X_test.index
17 X_test = pd.concat([X_test, X_test_cat], axis=1)
18 X_test.drop(columns=cat_var, axis=1, inplace=True)
```

```
1 def my_scaler(df, set_type="train", scaler_chosen=None, scaler_type="robust"):
2     """Scale a dataset with a passed/selected scaler type.
3
4     Args:
5         df (_type_): Dataset to scale.
6         set_type (str, optional): Dataset type 'train' or 'test'. Defaults to "train".
7         scaler_chosen (_type_, optional): Scaler object to use. Needed if set_type = test. Defaults to None.
8         scaler_type (str, optional): Scaler type to be used. Needed for trainin datasets. Defaults to "robust".
9
10    Returns:
11        _type_: Dataframe scaled.
12    """
13
14    from sklearn.preprocessing import RobustScaler, StandardScaler, PowerTransformer
15
16    scaler_dict = {"robust": RobustScaler(),
17                  "standard": StandardScaler(),
18                  "power": PowerTransformer()}
19
20    if set_type.lower() == "train":
21        if scaler_type is None:
22            return ValueError("A valid scaler type has to be chosen. Valid scalers are: 'robust', 'standard' and 'power'")
23        scaler_chosen = scaler_dict[scaler_type]
24        array_scaled = scaler_chosen.fit_transform(df)
25        df_scaled = pd.DataFrame(
26            data=array_scaled, columns=scaler_chosen.get_feature_names_out())
27
28    elif set_type.lower() == "test":
29        if scaler_chosen is None:
30            return ValueError("An scaler object has to be passed.")
31        array_scaled = scaler_chosen.fit_transform(df)
32        df_scaled = pd.DataFrame(
33            data=array_scaled, columns=scaler_chosen.get_feature_names_out())
34
35    df_scaled.index = df.index
36
37    return df_scaled, scaler_chosen
```



# Machine Learning

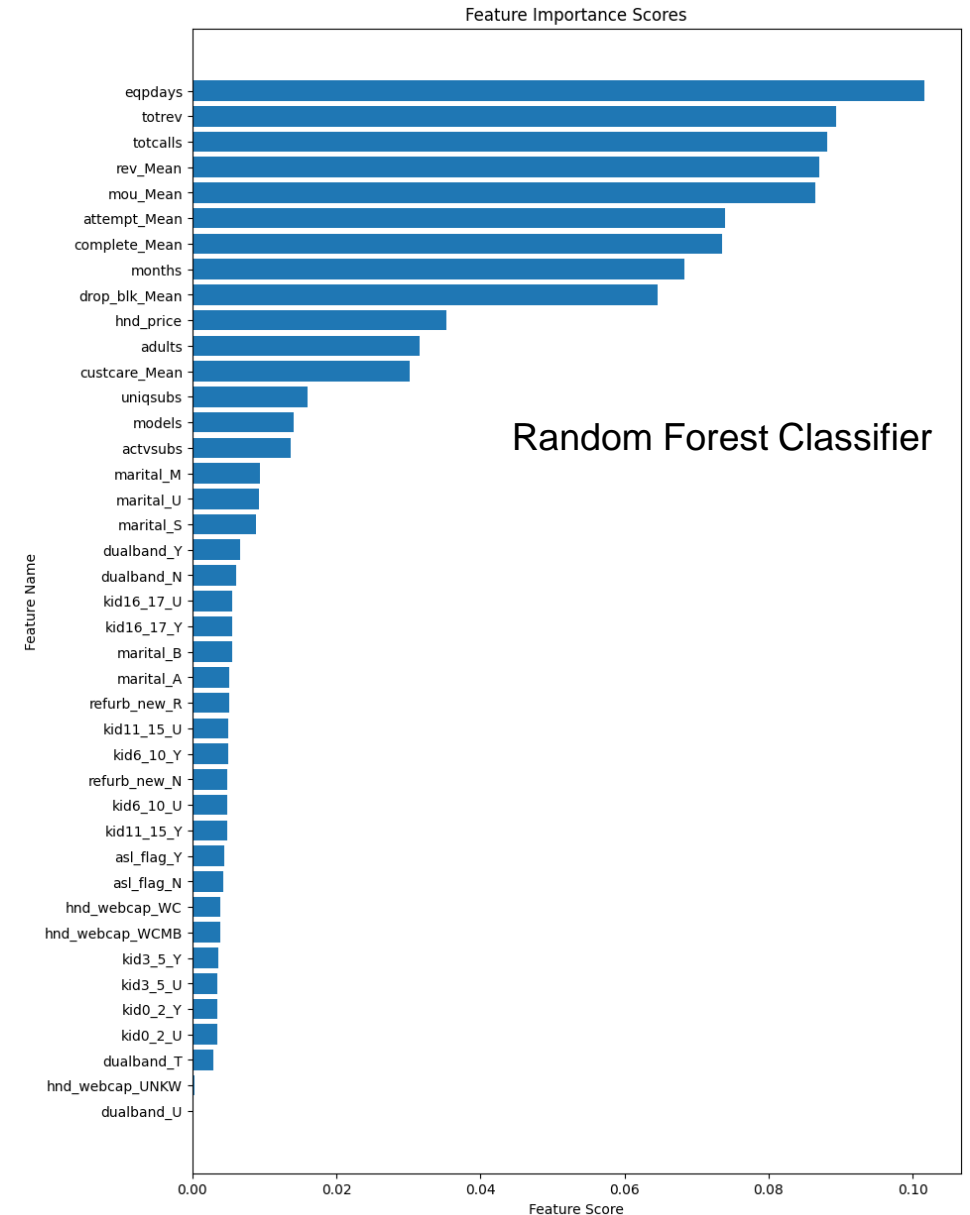
Model selection, evaluation and feature importances.



Random Forest Classifier chosen due to better performance.

Approach for (tentatively) improve the performance:

- Perform feature engineering (transform variables into another variables) to get features with higher importance.
- Drop features with missing values over a threshold.
- Test other models.
- Hyperparameter tuning (select different methods for doing stimations, select different number of trees in a random forest...)





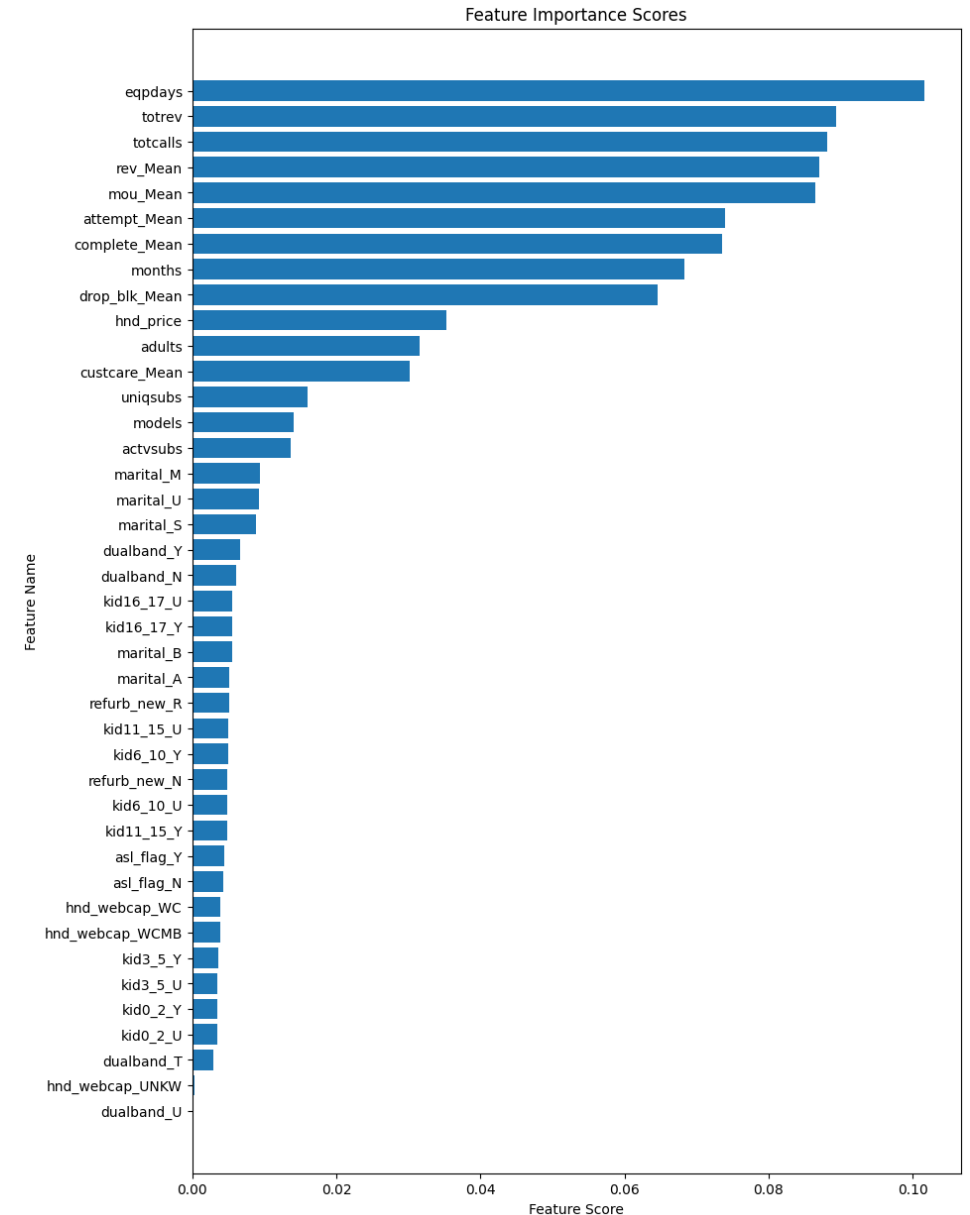
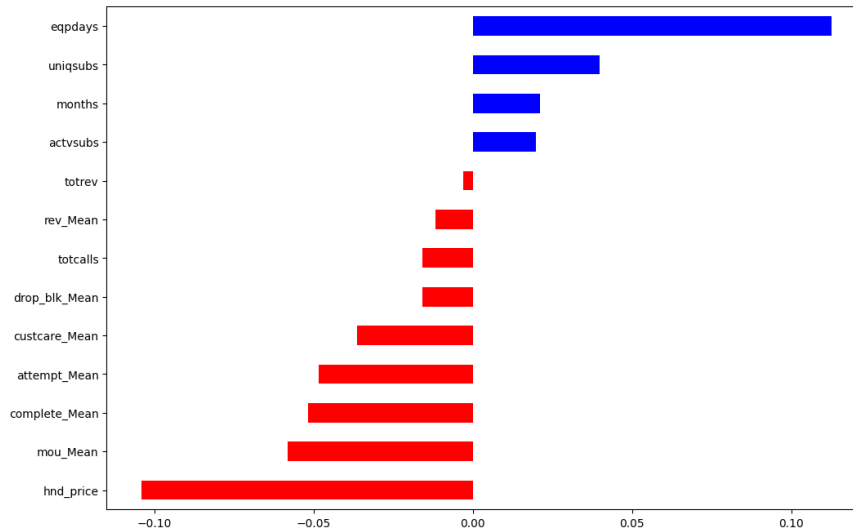
# Conclusion

Based on feature importances and correlation coefficients of features with the target variable:

- Eqpdays: Number of days of current equipment.
- Hnd\_price: Price of the device.
- uniqsubs: Number of unique subscribers on the house.
- totcalls: Total number of calls over the life of the customer.
- Custcare\_Mean: Mean number of customer care calls.

## Suggestions:

- Renew the device of the customers (new promotions).
- Family promotions (reduce the number of unique subscribers in the household, and reduce the cost of services).
- Promotions with reduced rates/fee for usual customers.







# Implementation in Airflow

## Setting up environment

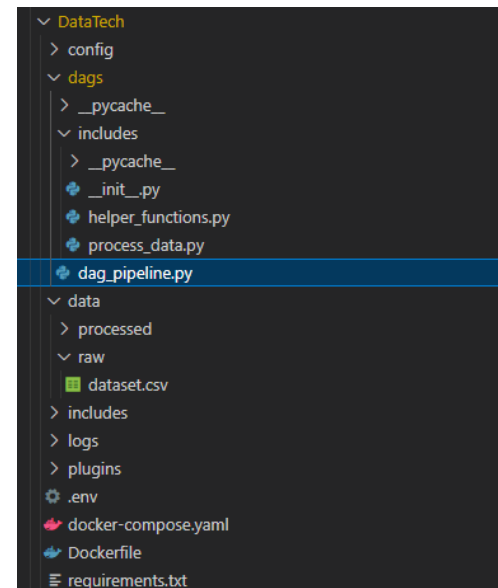
- Install docker.
- Create a folder for the project.
- Download the “*docker-compose.yaml*” on the project folder.
- Run the *docker compose up airflow-init* to initialize the database for the airflow project.

```
1 joblib==1.3.2
2 jsonschema==4.6.0
3 matplotlib==3.7.4
4 numpy==1.21.2
5 openpyxl==3.1.2
6 packaging==23.1
7 pandas==2.0.1
8 pickleshare==0.7.5
9 scikit-learn==1.3.2
10 scipy==1.10.1
11 seaborn==0.13.1
12 statsmodels==0.14.1
```

- List needed dependencies in a *requirements.txt* file.
- Write the instructions in the *Dockerfile*.

```
1 FROM apache/airflow:2.8.1
2 COPY requirements.txt /requirements.txt
3 RUN pip install --user --upgrade pip
4 RUN pip install --no-cache-dir --user -r /requirements.txt
```

- Build the new image with the extended packages.
- Change the image name in the “*docker-compose.yaml*”.
- Run again *docker compose up* to build the image with the extended image of airflow.



```
docker build . --tag extending_airflow:2.8.1
```

```
#
# Feel free to modify this file to suit your needs.
---
x-airflow-common:
  &airflow-common
  # In order to add custom dependencies or upgrade provider packages you can use yo
  # Comment the image line, place your Dockerfile in the directory where you placed
  # and uncomment the "build" line below, Then run `docker-compose build` to build
  image: ${AIRFLOW_IMAGE_NAME:-extending_airflow:2.8.1} # apache/airflow:2.8.1
  # build: .
```



# Implementation in Airflow

Docker and Airflow up and running

Structure of the data pipeline

