

Data Cleaning

#imports that may be of use

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

#open file

```
df = pd.read_csv('/Users/jamesmaikara/Downloads/customer_churn.csv')
```

```
df.head()
```

	state	account length	area code	phone number	international	plan	\
0	KS	128	415	382-4657		no	
1	OH	107	415	371-7191		no	
2	NJ	137	415	358-1921		no	
3	OH	84	408	375-9999		yes	
4	OK	75	415	330-6626		yes	

	voice mail plan	number vmail messages	total day minutes	total day calls	\
0	yes	25	265.1	110	
1	yes	26	161.6	123	
2	no	0	243.4	114	
3	no	0	299.4	71	
4	no	0	166.7	113	

	total day charge	...	total eve calls	total eve charge	\
0	45.07	...	99	16.78	
1	27.47	...	103	16.62	
2	41.38	...	110	10.30	
3	50.90	...	88	5.26	
4	28.34	...	122	12.61	

	total night minutes	total night calls	total night charge	\
0	244.7	91	11.01	
1	254.4	103	11.45	
2	162.6	104	7.32	
3	196.9	89	8.86	
4	186.9	121	8.41	

	total intl minutes	total intl calls	total intl charge	\
0	10.0	3	2.70	
1	13.7	3	3.70	

2	12.2	5	3.29
3	6.6	7	1.78
4	10.1	3	2.73

	customer service calls	churn
0	1	False
1	1	False
2	0	False
3	2	False
4	3	False

[5 rows x 21 columns]

#see general outlook of the data
df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3333 entries, 0 to 3332
Data columns (total 21 columns):
```

#	Column	Non-Null Count	Dtype
0	state	3333 non-null	object
1	account length	3333 non-null	int64
2	area code	3333 non-null	int64
3	phone number	3333 non-null	object
4	international plan	3333 non-null	object
5	voice mail plan	3333 non-null	object
6	number vmail messages	3333 non-null	int64
7	total day minutes	3333 non-null	float64
8	total day calls	3333 non-null	int64
9	total day charge	3333 non-null	float64
10	total eve minutes	3333 non-null	float64
11	total eve calls	3333 non-null	int64
12	total eve charge	3333 non-null	float64
13	total night minutes	3333 non-null	float64
14	total night calls	3333 non-null	int64
15	total night charge	3333 non-null	float64
16	total intl minutes	3333 non-null	float64
17	total intl calls	3333 non-null	int64
18	total intl charge	3333 non-null	float64
19	customer service calls	3333 non-null	int64
20	churn	3333 non-null	bool

dtypes: bool(1), float64(8), int64(8), object(4)
memory usage: 524.2+ KB

#check for null values
df.isna().sum()

state	0
account length	0

```

area code          0
phone number       0
international plan  0
voice mail plan    0
number vmail messages 0
total day minutes  0
total day calls     0
total day charge    0
total eve minutes   0
total eve calls     0
total eve charge    0
total night minutes 0
total night calls   0
total night charge  0
total intl minutes  0
total intl calls    0
total intl charge   0
customer service calls 0
churn              0
dtype: int64

```

#Check for placeholder values

```

for col in df.columns:
    unique_values = df[col].unique()
    print(f"{col}: {unique_values}")
    print('\n-----\n')

```

```

state: ['KS' 'OH' 'NJ' 'OK' 'AL' 'MA' 'MO' 'LA' 'WV' 'IN' 'RI' 'IA'
'MT' 'NY'
'ID' 'VT' 'VA' 'TX' 'FL' 'CO' 'AZ' 'SC' 'NE' 'WY' 'HI' 'IL' 'NH' 'GA'
'AK' 'MD' 'AR' 'WI' 'OR' 'MI' 'DE' 'UT' 'CA' 'MN' 'SD' 'NC' 'WA' 'NM'
'NV' 'DC' 'KY' 'ME' 'MS' 'TN' 'PA' 'CT' 'ND']

```

```

-----

account length: [128 107 137  84  75 118 121 147 117 141  65  74 168
95  62 161  85  93
 76  73  77 130 111 132 174  57  54  20  49 142 172  12  72  36  78
136
149  98 135  34 160  64  59 119  97  52  60  10  96  87  81  68 125
116
 38  40  43 113 126 150 138 162  90  50  82 144  46  70  55 106  94
155
 80 104  99 120 108 122 157 103  63 112  41 193  61  92 131 163  91
127
110 140  83 145  56 151 139   6 115 146 185 148  32  25 179  67  19
170
164  51 208  53 105  66  86  35  88 123  45 100 215  22  33 114  24
101
143  48  71 167  89 199 166 158 196 209  16  39 173 129  44  79  31

```

124
37 159 194 154 21 133 224 58 11 109 102 165 18 30 176 47 190
152
26 69 186 171 28 153 169 13 27 3 42 189 156 134 243 23 1
205
200 5 9 178 181 182 217 177 210 29 180 2 17 7 212 232 192
195
197 225 184 191 201 15 183 202 8 175 4 188 204 221]

area code: [415 408 510]

phone number: ['382-4657' '371-7191' '358-1921' ... '328-8230' '364-6381' '400-4344']

international plan: ['no' 'yes']

voice mail plan: ['yes' 'no']

number vmail messages: [25 26 0 24 37 27 33 39 30 41 28 34 46 29 35
21 32 42 36 22 23 43 31 38
40 48 18 17 45 16 20 14 19 51 15 11 12 47 8 44 49 4 10 13 50 9]

total day minutes: [265.1 161.6 243.4 ... 321.1 231.1 180.8]

total day calls: [110 123 114 71 113 98 88 79 97 84 137 127 96
70 67 139 66 90
117 89 112 103 86 76 115 73 109 95 105 121 118 94 80 128 64
106
102 85 82 77 120 133 135 108 57 83 129 91 92 74 93 101 146
72
99 104 125 61 100 87 131 65 124 119 52 68 107 47 116 151 126
122
111 145 78 136 140 148 81 55 69 158 134 130 63 53 75 141 163
59
132 138 54 58 62 144 143 147 36 40 150 56 51 165 30 48 60
42
0 45 160 149 152 142 156 35 49 157 44]

total day charge: [45.07 27.47 41.38 ... 54.59 39.29 30.74]

total eve minutes: [197.4 195.5 121.2 ... 153.4 288.8 265.9]

total eve calls: [99 103 110 88 122 101 108 94 80 111 83 148 71
75 76 97 90 65
93 121 102 72 112 100 84 109 63 107 115 119 116 92 85 98 118
74
117 58 96 66 67 62 77 164 126 142 64 104 79 95 86 105 81
113
106 59 48 82 87 123 114 140 128 60 78 125 91 46 138 129 89
133
136 57 135 139 51 70 151 137 134 73 152 168 68 120 69 127 132
143
61 124 42 54 131 52 149 56 37 130 49 146 147 55 12 50 157
155
45 144 36 156 53 141 44 153 154 150 43 0 145 159 170]

total eve charge: [16.78 16.62 10.3 ... 13.04 24.55 22.6]

total night minutes: [244.7 254.4 162.6 ... 280.9 120.1 279.1]

total night calls: [91 103 104 89 121 118 96 90 97 111 94 128
115 99 75 108 74 133
64 78 105 68 102 148 98 116 71 109 107 135 92 86 127 79 87
129
57 77 95 54 106 53 67 139 60 100 61 73 113 76 119 88 84
62
137 72 142 114 126 122 81 123 117 82 80 120 130 134 59 112 132
110
101 150 69 131 83 93 124 136 125 66 143 58 55 85 56 70 46
42
152 44 145 50 153 49 175 63 138 154 140 141 146 65 51 151 158
155
157 147 144 149 166 52 33 156 38 36 48 164]

total night charge: [11.01 11.45 7.32 8.86 8.41 9.18 9.57 9.53
9.71 14.69 9.4 8.82
6.35 8.65 9.14 7.23 4.02 5.83 7.46 8.68 9.43 8.18 8.53
10.67
11.28 8.22 4.59 8.17 8.04 11.27 11.08 13.2 12.61 9.61 6.88
5.82
10.25 4.58 8.47 8.45 5.5 14.02 8.03 11.94 7.34 6.06 10.9
6.44
3.18 10.66 11.21 12.73 10.28 12.16 6.34 8.15 5.84 8.52 7.5
7.48
6.21 11.95 7.15 9.63 7.1 6.91 6.69 13.29 11.46 7.76 6.86
8.16
12.15 7.79 7.99 10.29 10.08 12.53 7.91 10.02 8.61 14.54 8.21
9.09
4.93 11.39 11.88 5.75 7.83 8.59 7.52 12.38 7.21 5.81 8.1
11.04
11.19 8.55 8.42 9.76 9.87 10.86 5.36 10.03 11.15 9.51 6.22
2.59
7.65 6.45 9. 6.4 9.94 5.08 10.23 11.36 6.97 10.16 7.88
11.91
6.61 11.55 11.76 9.27 9.29 11.12 10.69 8.8 11.85 7.14 8.71
11.42
4.94 9.02 11.22 4.97 9.15 5.45 7.27 12.91 7.75 13.46 6.32
12.13
11.97 6.93 11.66 7.42 6.19 11.41 10.33 10.65 11.92 4.77 4.38
7.41
12.1 7.69 8.78 9.36 9.05 12.7 6.16 6.05 10.85 8.93 3.48
10.4
5.05 10.71 9.37 6.75 8.12 11.77 11.49 11.06 11.25 11.03 10.82
8.91
8.57 8.09 10.05 11.7 10.17 8.74 5.51 11.11 3.29 10.13 6.8
8.49
9.55 11.02 9.91 7.84 10.62 9.97 3.44 7.35 9.79 8.89 8.14
6.94
10.49 10.57 10.2 6.29 8.79 10.04 12.41 15.97 9.1 11.78 12.75
11.07
12.56 8.63 8.02 10.42 8.7 9.98 7.62 8.33 6.59 13.12 10.46
6.63
8.32 9.04 9.28 10.76 9.64 11.44 6.48 10.81 12.66 11.34 8.75
13.05
11.48 14.04 13.47 5.63 6.6 9.72 11.68 6.41 9.32 12.95 13.37
9.62
6.03 8.25 8.26 11.96 9.9 9.23 5.58 7.22 6.64 12.29 12.93
11.32
6.85 8.88 7.03 8.48 3.59 5.86 6.23 7.61 7.66 13.63 7.9
11.82
7.47 6.08 8.4 5.74 10.94 10.35 10.68 4.34 8.73 5.14 8.24
9.99
13.93 8.64 11.43 5.79 9.2 10.14 12.11 7.53 12.46 8.46 8.95

[illegible]

[illegible]

5.9 7.97 5. 10.97 5.88 12.34 12.03 14.97 15.06 12.85 6.54
11.24
12.64 7.06 5.38 13.14 3.99 3.32 4.51 4.12 3.93 2.4 11.75
4.03
15.85 6.81 14.25 14.09 16.42 6.7 12.74 2.76 12.12 6.99 6.68
11.81
7.96 5.06 13.16 2.13 13.17 5.12 5.65 12.37 10.53]

total intl minutes: [10. 13.7 12.2 6.6 10.1 6.3 7.5 7.1 8.7 11.2
12.7 9.1 12.3 13.1
5.4 13.8 8.1 13. 10.6 5.7 9.5 7.7 10.3 15.5 14.7 11.1 14.2 12.6
11.8 8.3 14.5 10.5 9.4 14.6 9.2 3.5 8.5 13.2 7.4 8.8 11. 7.8
6.8 11.4 9.3 9.7 10.2 8. 5.8 12.1 12. 11.6 8.2 6.2 7.3 6.1
11.7 15. 9.8 12.4 8.6 10.9 13.9 8.9 7.9 5.3 4.4 12.5 11.3 9.
9.6 13.3 20. 7.2 6.4 14.1 14.3 6.9 11.5 15.8 12.8 16.2 0. 11.9
9.9 8.4 10.8 13.4 10.7 17.6 4.7 2.7 13.5 12.9 14.4 10.4 6.7 15.4
4.5 6.5 15.6 5.9 18.9 7.6 5. 7. 14. 18. 16. 14.8 3.7 2.
4.8 15.3 6. 13.6 17.2 17.5 5.6 18.2 3.6 16.5 4.6 5.1 4.1 16.3
14.9 16.4 16.7 1.3 15.2 15.1 15.9 5.5 16.1 4. 16.9 5.2 4.2 15.7
17. 3.9 3.8 2.2 17.1 4.9 17.9 17.3 18.4 17.8 4.3 2.9 3.1 3.3
2.6 3.4 1.1 18.3 16.6 2.1 2.4 2.5]

total intl calls: [3 5 7 6 4 2 9 19 1 10 15 8 11 0 12 13 18
14 16 20 17]

total intl charge: [2.7 3.7 3.29 1.78 2.73 1.7 2.03 1.92 2.35 3.02
3.43 2.46 3.32 3.54
1.46 3.73 2.19 3.51 2.86 1.54 2.57 2.08 2.78 4.19 3.97 3. 3.83 3.4
3.19 2.24 3.92 2.84 2.54 3.94 2.48 0.95 2.3 3.56 2. 2.38 2.97 2.11
1.84 3.08 2.51 2.62 2.75 2.16 1.57 3.27 3.24 3.13 2.21 1.67 1.97 1.65
3.16 4.05 2.65 3.35 2.32 2.94 3.75 2.4 2.13 1.43 1.19 3.38 3.05 2.43
2.59 3.59 5.4 1.94 1.73 3.81 3.86 1.86 3.11 4.27 3.46 4.37 0. 3.21
2.67 2.27 2.92 3.62 2.89 4.75 1.27 0.73 3.65 3.48 3.89 2.81 1.81 4.16
1.22 1.76 4.21 1.59 5.1 2.05 1.35 1.89 3.78 4.86 4.32 4. 1. 0.54
1.3 4.13 1.62 3.67 4.64 4.73 1.51 4.91 0.97 4.46 1.24 1.38 1.11 4.4
4.02 4.43 4.51 0.35 4.1 4.08 4.29 1.49 4.35 1.08 4.56 1.4 1.13 4.24
4.59 1.05 1.03 0.59 4.62 1.32 4.83 4.67 4.97 4.81 1.16 0.78 0.84 0.89
0.7 0.92 0.3 4.94 4.48 0.57 0.65 0.68]

customer service calls: [1 0 2 3 4 5 7 9 6 8]

```
churn: [False True]
```

```
-----
```

```
#Get the unique number of states
```

```
len(df['state'].unique())
```

```
51
```

```
# Check balance within the data
```

```
df['churn'].value_counts()
```

```
churn
```

```
False    2850
```

```
True      483
```

```
Name: count, dtype: int64
```

```
# look for outliers and spread of data
```

```
df.describe()
```

	account length	area code	number vmail messages	total day
minutes \				
count	3333.000000	3333.000000	3333.000000	
3333.000000				
mean	101.064806	437.182418	8.099010	
179.775098				
std	39.822106	42.371290	13.688365	
54.467389				
min	1.000000	408.000000	0.000000	
0.000000				
25%	74.000000	408.000000	0.000000	
143.700000				
50%	101.000000	415.000000	0.000000	
179.400000				
75%	127.000000	510.000000	20.000000	
216.400000				
max	243.000000	510.000000	51.000000	
350.800000				

	total day calls	total day charge	total eve minutes	total eve
calls \				
count	3333.000000	3333.000000	3333.000000	
3333.000000				
mean	100.435644	30.562307	200.980348	
100.114311				
std	20.069084	9.259435	50.713844	
19.922625				
min	0.000000	0.000000	0.000000	
0.000000				

25%	87.000000	24.430000	166.600000
87.000000			
50%	101.000000	30.500000	201.400000
100.000000			
75%	114.000000	36.790000	235.300000
114.000000			
max	165.000000	59.640000	363.700000
170.000000			

	total eve charge	total night minutes	total night calls \
count	3333.000000	3333.000000	3333.000000
mean	17.083540	200.872037	100.107711
std	4.310668	50.573847	19.568609
min	0.000000	23.200000	33.000000
25%	14.160000	167.000000	87.000000
50%	17.120000	201.200000	100.000000
75%	20.000000	235.300000	113.000000
max	30.910000	395.000000	175.000000

	total night charge	total intl minutes	total intl calls \
count	3333.000000	3333.000000	3333.000000
mean	9.039325	10.237294	4.479448
std	2.275873	2.791840	2.461214
min	1.040000	0.000000	0.000000
25%	7.520000	8.500000	3.000000
50%	9.050000	10.300000	4.000000
75%	10.590000	12.100000	6.000000
max	17.770000	20.000000	20.000000

	total intl charge	customer service calls
count	3333.000000	3333.000000
mean	2.764581	1.562856
std	0.753773	1.315491
min	0.000000	0.000000
25%	2.300000	1.000000
50%	2.780000	1.000000
75%	3.270000	2.000000
max	5.400000	9.000000

Check spread of data for major outliers

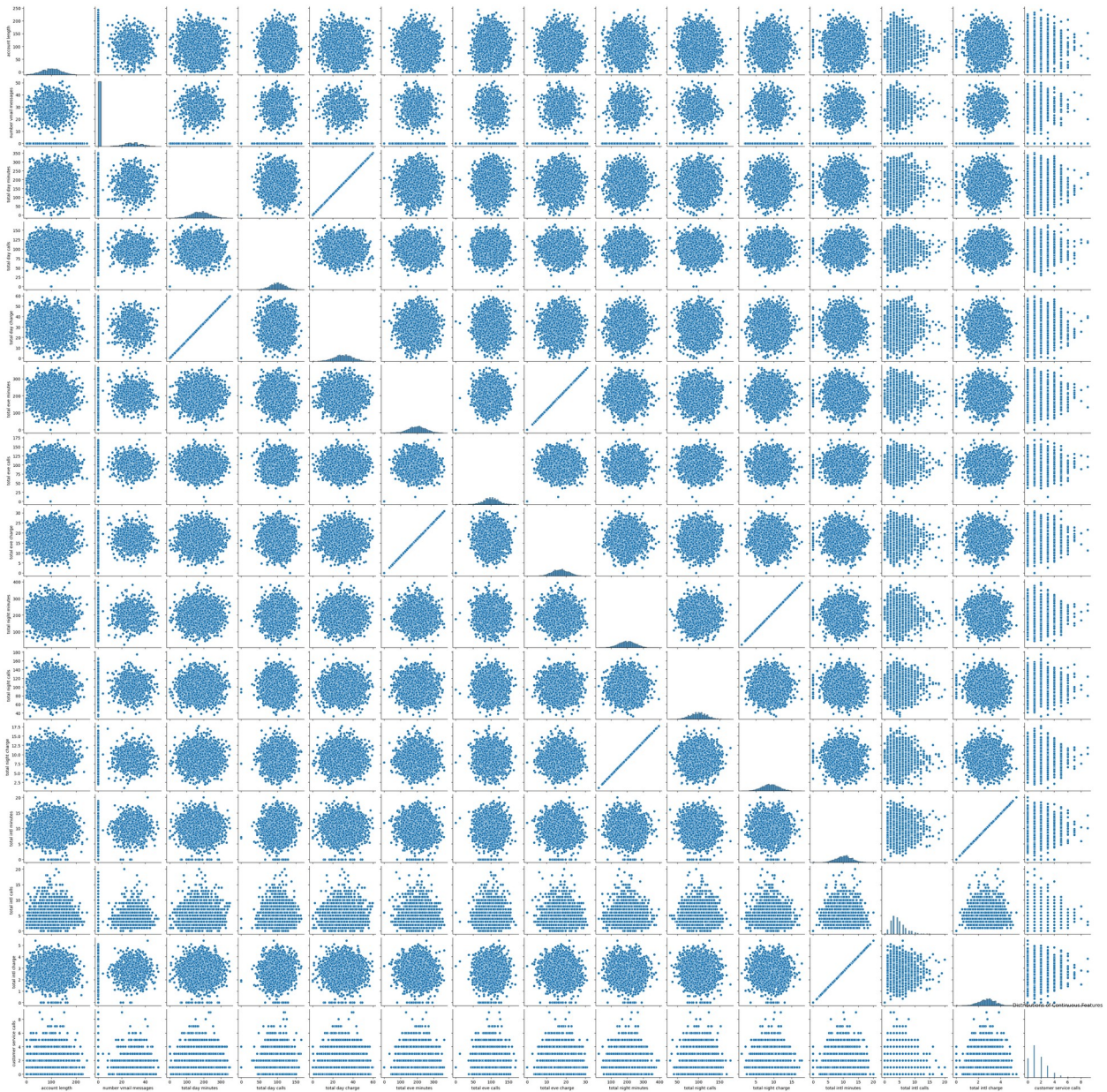
```
cont_df = df.drop(columns=['state', 'phone number', 'international plan', 'voice mail plan', 'churn', 'area code'])
```

```
sns.pairplot(cont_df)
```

```
plt.title('Plots of the continuous features')
```

```
plt.show()
```

```
/Users/jamesmaikara/anaconda3/lib/python3.11/site-packages/seaborn/axisgrid.py:118: UserWarning: The figure layout has changed to tight
  self._figure.tight_layout(*args, **kwargs)
```



```
# A closer look at these 3 weird distributions
for col in ['total intl calls', 'customer service calls', 'number
vmail messages']:
    sns.distplot(df[col], bins=5)
    plt.title(f'Distribution of Feature: {col.title()}')
    plt.show()
```

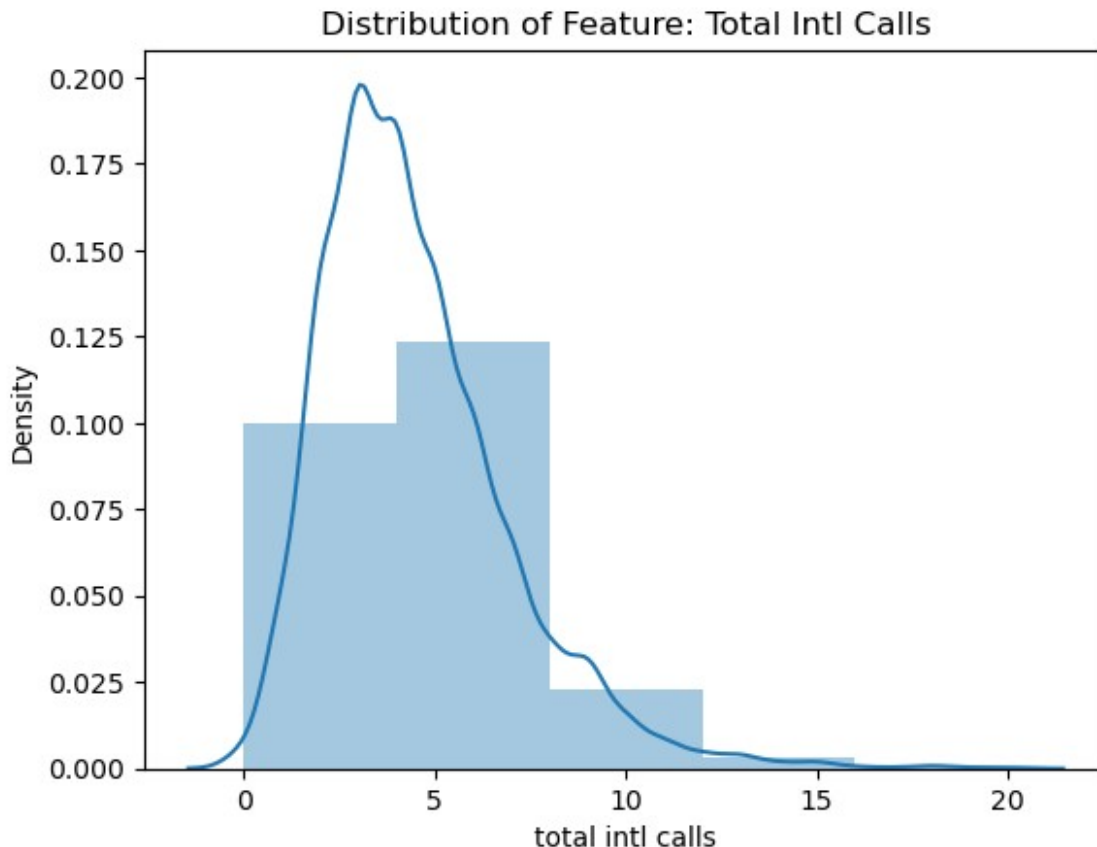
/var/folders/mn/l_cg50ls31vcyix2tfp461z00000gq/T/
ipykernel_8479/3875214485.py:3: UserWarning:

`distplot` is a deprecated function and will be removed in seaborn
v0.14.0.

Please adapt your code to use either ``displot`` (a figure-level function with similar flexibility) or ``histplot`` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>

```
sns.distplot(df[col], bins=5)
```



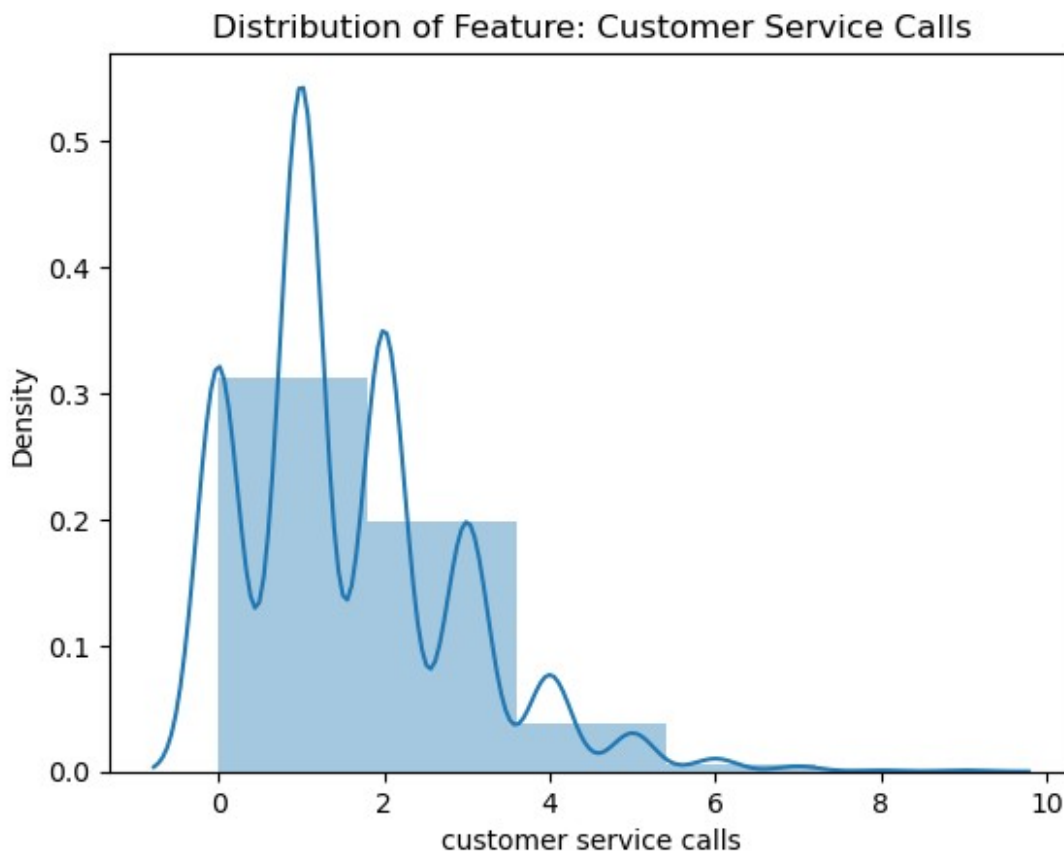
```
/var/folders/mn/l_cg501s3lvcyjx2tftp461z00000gq/T/  
ipykernel_8479/3875214485.py:3: UserWarning:
```

``distplot`` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either ``displot`` (a figure-level function with similar flexibility) or ``histplot`` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>


```
sns.distplot(df[col], bins=5)
```



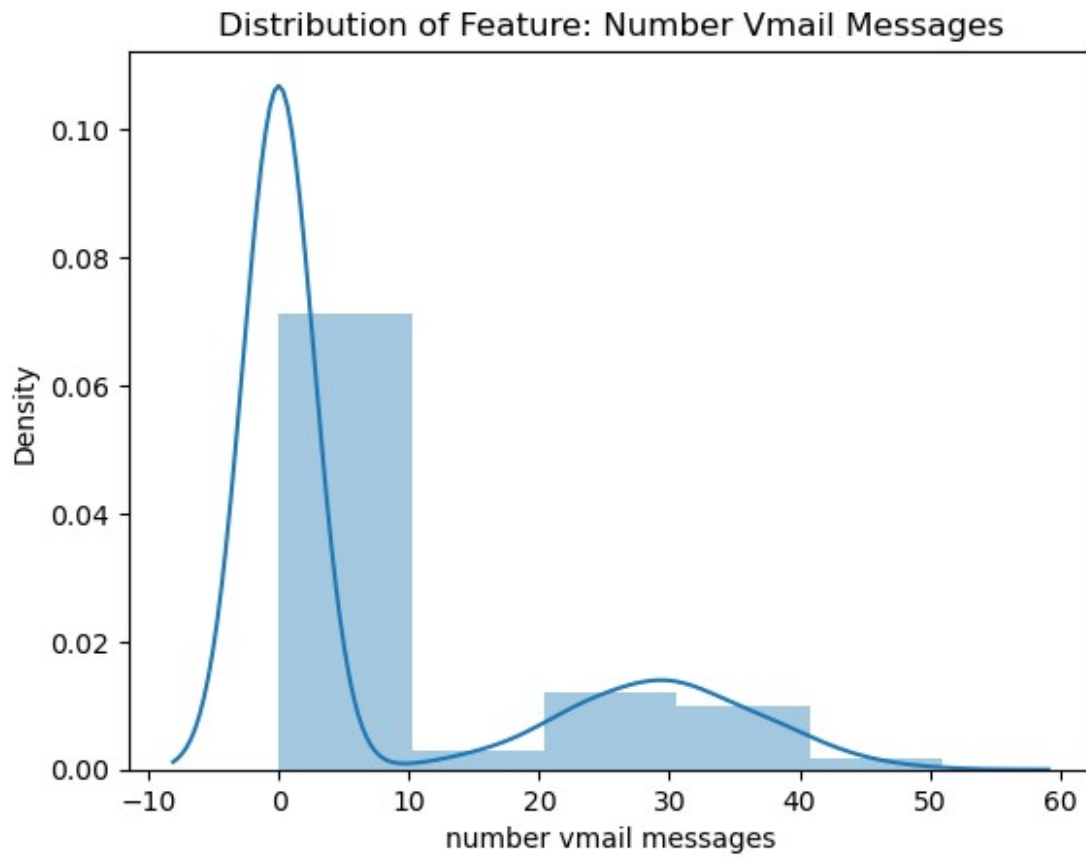
```
/var/folders/mn/l_cg50ls31vcyjsx2tfp461z00000gq/T/  
ipykernel_8479/3875214485.py:3: UserWarning:
```

``distplot`` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either ``displot`` (a figure-level function with similar flexibility) or ``histplot`` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>

```
sns.distplot(df[col], bins=5)
```



Conclusion

The dataset was clean and did not need much alteration.