

Machine Learning with the Experts: School Budgets

1). Exploring the raw data:

a). Summarizing the Data

Print the summary statistics

```
print(df.describe())
```

Import matplotlib.pyplot as plt

```
import matplotlib.pyplot as plt
```

Create the histogram

```
plt.hist(df['FTE'].dropna())
```

Add title and labels

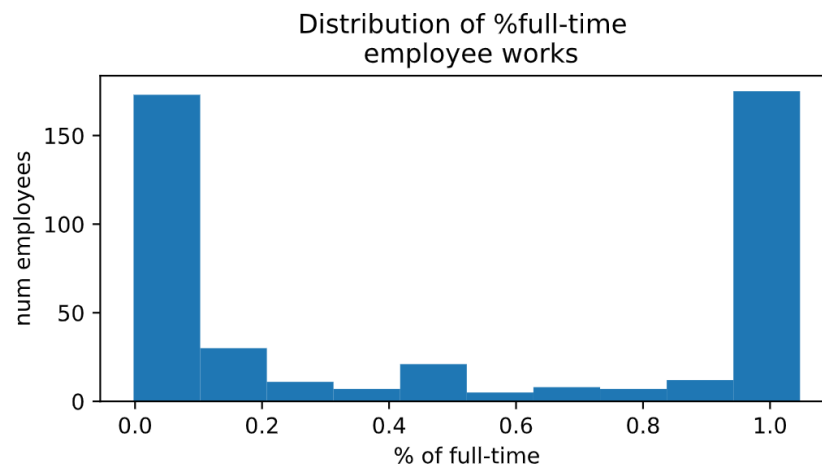
```
plt.title('Distribution of %full-time \n employee works')
```

```
plt.xlabel('% of full-time')
```

```
plt.ylabel('num employees')
```

Display the histogram

```
plt.show()
```



b). Encode the labels as Categorical Values

```
# Define the lambda function: categorize_label
```

```
categorize_label = lambda x: x.astype('category')
```

```
# Convert df[LABELS] to a categorical type
```

```
df[LABELS] = df[LABELS].apply(categorize_label, axis=0)
```

```
# Print the converted dtypes
```

```
print(df[LABELS].dtypes)
```

<script.py> output:

```
Function      category
Use           category
Sharing       category
Reporting     category
Student_Type  category
Position_Type category
Object_Type   category
Pre_K         category
Operating_Status category
dtype: object
```

c). Counting Unique Labels

```
# Import matplotlib.pyplot
```

```
import matplotlib.pyplot as plt
```

```
# Calculate number of unique values for each label: num_unique_labels
```

```
num_unique_labels = df[LABELS].apply(pd.Series.nunique)
```

```
# Plot number of unique values for each label
```

```
num_unique_labels.plot(kind='bar')
```

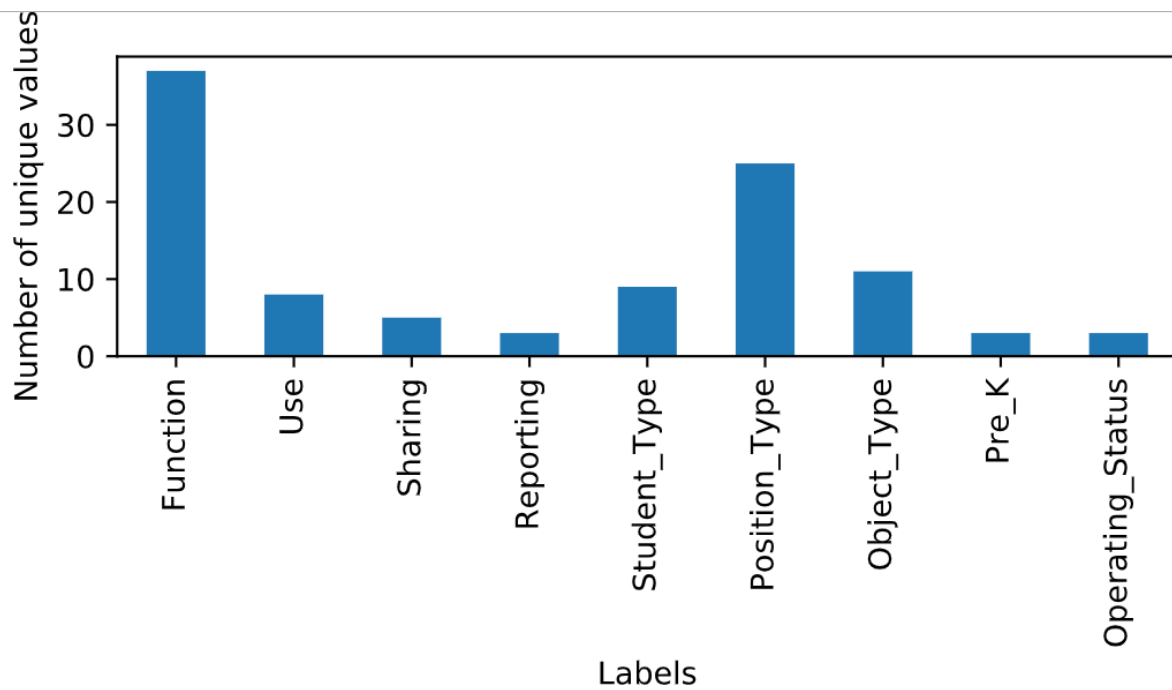
```
# Label the axes
```

```
plt.xlabel('Labels')
```

```
plt.ylabel('Number of unique values')
```

```
# Display the plot
```

```
plt.show()
```



d). Computing logloss with numpy

Compute and print log loss for 1st case

```
correct_confident = compute_log_loss(correct_confident, actual_labels)
```

```
print("Log loss, correct and confident: {}".format(correct_confident))
```

Compute log loss for 2nd case

```
correct_not_confident = compute_log_loss(correct_not_confident, actual_labels)
```

```
print("Log loss, correct and not confident: {}".format(correct_not_confident))
```

Compute and print log loss for 3rd case

```
wrong_not_confident = compute_log_loss(wrong_not_confident, actual_labels)
```

```
print("Log loss, wrong and not confident: {}".format(wrong_not_confident))
```

Compute and print log loss for 4th case

```
wrong_confident = compute_log_loss(wrong_confident, actual_labels)
```

```
print("Log loss, wrong and confident: {}".format(wrong_confident))
```

Compute and print log loss for actual labels

```
actual_labels = compute_log_loss(actual_labels, actual_labels)
```

```
print("Log loss, actual labels: {}".format(actual_labels))
```

<script.py> output:

Log loss, correct and confident: 0.05129329438755058

Log loss, correct and not confident: 0.4307829160924542

Log loss, wrong and not confident: 1.049822124498678

Log loss, wrong and confident: 2.9957322735539904

Log loss, actual labels: 9.99200722162646e-15