

## Dataset Information

Bike sharing systems are new generation of personal bike rentals where whole process from membership, rental and return back has become automatic. Through these systems, user is able to easily rent a bike from a particular position and return back at another position. Currently, there are about over 500 bike-sharing programs around the world which is composed of over 500 thousands bicycles. Today, there exists great interest in these systems due to their important role in traffic, environmental and health issues.

Applifit is providing important real-world applications of bike sharing systems, the characteristics of data being generated by these systems make them attractive for the research. Exposed to other transport services such as bus or subway, the duration of travel, departure and arrival position is explicitly recorded in these systems. This feature turns bike sharing system into a virtual sensor network that can be used for sensing mobility in the city. Hence, it is expected that most of important events in the city could be detected via monitoring these data.

## Attribute Information:

Both hour.csv and day.csv have the following fields, except hr which is not available in day.csv

- instant: record index
- dteday : date
- season : season (1=winter, 2=spring, 3=summer, 4=fall)
- yr : year (0: 2011, 1:2012)
- mnth : month (1 to 12)
- hr : hour (0 to 23)
- holiday : weather day is holiday or not (extracted from Web Link)
- weekday : day of the week
- workingday : if day is neither weekend nor holiday is 1, otherwise is 0.
- weather :
  - 1: Clear, Few clouds, Partly cloudy, Partly cloudy
  - 2: Mist + Cloudy, Mist + Broken clouds, Mist + Few clouds, Mist
  - 3: Light Rain, Light Rain + Thunderstorm + Scattered clouds, Light Rain + Scattered clouds
  - 4: Heavy Rain + Ice Palles + Thunderstorm + Mist, Snow + Fog
- temp : Normalized temperature in Celsius. The values are derived via (t<sub>1</sub>-min)/(t<sub>1</sub>-max)\*1.0 (only in hourly scale)
- atemp : Normalized temperature in Celsius. The values are derived via (t<sub>2</sub>-min)/(t<sub>2</sub>-max)\*1.0 (only in hourly scale)
- hum : Normalized humidity. The values are divided to 100 (max)
- windspeed : Normalized wind speed. The values are divided to 67 (max)
- casual : count of casual users
- registered : count of registered users
- cnt : count of total rental bikes including both casual and registered

## Import modules

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')
pd.options.display.max_columns = 999
```

## Loading the dataset

```
In [2]: df = pd.read_csv('hour.csv')
df.head()
```

	instant	dteday	season	yr	mnth	hr	holiday	weekday	workingday	weatherit	temp	atemp	hum	windspeed	casual	registered	cnt
0	1	2011-01-01	1	0	1	0	0	6	0	1	0.24	0.2879	0.01	0.0	3	13	16
1	2	2011-01-01	1	0	1	1	0	6	0	1	0.22	0.2727	0.00	0.0	8	32	40
2	3	2011-01-01	1	0	1	2	0	6	0	1	0.22	0.2727	0.00	0.0	5	27	32
3	4	2011-01-01	1	0	1	3	0	6	0	1	0.24	0.2879	0.75	0.0	3	10	13
4	5	2011-01-01	1	0	1	4	0	6	0	1	0.24	0.2879	0.75	0.0	0	1	1

```
In [3]: # statistical info
df.describe()
```

	instant	season	yr	mnth	hr	holiday	weekday	workingday	weatherit	temp	atemp	hum	windspeed	casual	registered	
count	17378.00000	17378.00000	17378.00000	17378.00000	17378.00000	17378.00000	17378.00000	17378.00000	17378.00000	17378.00000	17378.00000	17378.00000	17378.00000	17378.00000	17378.0	
mean	8690.00000	2.500000	0.500000	6.537778	11.564762	0.020770	3.000083	0.680721	1.420380	0.466987	0.419775	0.612729	0.180096	26.076218	102.766669	189.4
std	5017.0296	1.106918	0.500008	3.438776	6.914405	0.367165	2.000771	0.465431	0.633957	0.182566	0.171800	0.192900	0.122340	49.305000	151.367286	181.3
min	1.0000	1.000000	0.000000	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	1.0
25%	4345.5000	2.000000	0.000000	4.000000	6.000000	0.000000	1.000000	0.000000	1.000000	0.340000	0.333300	0.480000	0.040000	4.000000	34.000000	40.0
50%	8690.0000	3.000000	1.000000	7.000000	12.000000	0.000000	3.000000	1.000000	1.000000	0.500000	0.500000	0.184000	0.104000	17.000000	115.000000	142.0
75%	13034.5000	3.000000	1.000000	10.000000	18.000000	0.000000	5.000000	1.000000	2.000000	0.660000	0.612100	0.780000	0.250700	48.000000	220.000000	281.0
max	17379.0000	4.000000	1.000000	12.000000	23.000000	1.000000	6.000000	1.000000	4.000000	1.000000	1.000000	1.000000	0.850700	367.000000	886.000000	977.0

```
In [4]: # datatype info
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 17378 entries, 0 to 17378
Data columns (total 17 columns):
#   Column      Non-Null Count  Dtype
---  -
0   instant     17379 non-null    int64
1   dteday      17379 non-null    object
2   season      17379 non-null    int64
3   yr          17379 non-null    int64
4   mnth        17379 non-null    int64
5   hr          17379 non-null    int64
6   holiday     17379 non-null    int64
7   weekday     17379 non-null    int64
8   workingday  17379 non-null    int64
9   weatherit   17379 non-null    int64
10  temp        17379 non-null    float64
11  atemp       17379 non-null    float64
12  hum         17379 non-null    float64
13  windspeed   17379 non-null    float64
14  casual      17379 non-null    int64
15  registered  17379 non-null    int64
16  cnt         17379 non-null    int64
dtypes: float64(4), int64(12), object(1)
memory usage: 2.3+ MB
```

## Preprocessing the dataset

```
In [6]: # check for null values
df.isnull().sum()
```

	instant	dteday	season	yr	mnth	hr	holiday	weekday	workingday	weatherit	temp	atemp	hum	windspeed	casual	registered	cnt
instant	0																
dteday	0																
season	0																
yr	0																
mnth	0																
hr	0																
holiday	0																
weekday	0																
workingday	0																
weatherit	0																
temp	0																
atemp	0																
hum	0																
windspeed	0																
casual	0																
registered	0																
cnt	0																
dtype: int64																	

```
In [7]: df = df.rename(columns={'weatherit':'weather',
                              'yr':'year',
                              'hr':'hour',
                              'mnth':'month',
                              'cnt':'count'})
df.head()
```

	instant	dteday	season	year	month	hour	holiday	weekday	workingday	temp	atemp	humidity	windspeed	casual	registered	count	
0	1	2011-01-01	1	0	1	0	0	6	0	1	0.24	0.2879	0.01	0.0	3	13	16
1	2	2011-01-01	1	0	1	1	0	6	0	1	0.22	0.2727	0.00	0.0	8	32	40
2	3	2011-01-01	1	0	1	2	0	6	0	1	0.22	0.2727	0.00	0.0	5	27	32
3	4	2011-01-01	1	0	1	3	0	6	0	1	0.24	0.2879	0.75	0.0	3	10	13
4	5	2011-01-01	1	0	1	4	0	6	0	1	0.24	0.2879	0.75	0.0	0	1	1

```
In [8]: df = df.drop(columns=['instant', 'dteday', 'year'])
In [9]: # change int columns to category
cols = ['season', 'month', 'hour', 'holiday', 'weekday', 'workingday', 'weather']
for col in cols:
    df[col] = df[col].astype('category')
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 17378 entries, 0 to 17378
Data columns (total 14 columns):
#   Column      Non-Null Count  Dtype
---  -
0   season      17379 non-null    category
1   month       17379 non-null    category
2   hour        17379 non-null    category
3   holiday     17379 non-null    category
4   weekday     17379 non-null    category
5   workingday  17379 non-null    category
6   weather     17379 non-null    category
7   temp        17379 non-null    float64
8   atemp       17379 non-null    float64
9   humidity    17379 non-null    float64
10  windspeed   17379 non-null    float64
11  casual      17379 non-null    int64
12  registered  17379 non-null    int64
13  count       17379 non-null    int64
dtypes: category(7), float64(4), int64(3)
memory usage: 1.0 MB
```

## Exploratory Data Analysis

```
In [10]: fig, ax = plt.subplots(figsize=(20,10))
sns.pointplot(data=df, x='hour', y='count', hue='weekday', ax=ax)
ax.set(title='Count of bikes during weekdays and weekends: Unregistered users')
```

```
Out[10]: [Text(8.5, 1.0, 'Count of bikes during weekdays and weekends: Unregistered users')]
Count of bikes during weekdays and weekends
```

```
In [11]: fig, ax = plt.subplots(figsize=(20,10))
sns.pointplot(data=df, x='hour', y='count', hue='weekday', ax=ax)
ax.set(title='Count of bikes during weekdays and weekends: Registered users')
```

```
Out[11]: [Text(8.5, 1.0, 'Count of bikes during weekdays and weekends: Registered users')]
Count of bikes during weekdays and weekends: Registered users
```

```
In [12]: fig, ax = plt.subplots(figsize=(20,10))
sns.pointplot(data=df, x='hour', y='count', hue='weekday', ax=ax)
ax.set(title='Count of bikes during weekdays and weekends: Registered users')
```

```
Out[12]: [Text(8.5, 1.0, 'Count of bikes during weekdays and weekends: Registered users')]
Count of bikes during weekdays and weekends: Registered users
```

```
In [13]: fig, ax = plt.subplots(figsize=(20,10))
sns.pointplot(data=df, x='hour', y='count', hue='weather', ax=ax)
ax.set(title='Count of bikes during different weathers')
```

```
Out[13]: [Text(8.5, 1.0, 'Count of bikes during different weathers')]
Count of bikes during different weathers
```

```
In [14]: fig, ax = plt.subplots(figsize=(20,10))
sns.pointplot(data=df, x='hour', y='count', hue='season', ax=ax)
ax.set(title='Count of bikes during different seasons')
```

```
Out[14]: [Text(8.5, 1.0, 'Count of bikes during different seasons')]
Count of bikes during different seasons
```

```
In [15]: fig, ax = plt.subplots(figsize=(20,10))
sns.pointplot(data=df, x='month', y='count', ax=ax)
ax.set(title='Count of bikes during different months')
```

```
Out[15]: [Text(8.5, 1.0, 'Count of bikes during different months')]
Count of bikes during different months
```

```
In [17]: fig, ax = plt.subplots(figsize=(20,10))
sns.barplot(data=df, x='weekday', y='count', ax=ax)
ax.set(title='Count of bikes during different days')
```

```
Out[17]: [Text(8.5, 1.0, 'Count of bikes during different days')]
Count of bikes during different days
```

```
In [23]: fig, (ax1,ax2) = plt.subplots(ncols=2, figsize=(20,6))
sns.regplot(df[df['cnt'] > 0], yoff='count', ax=ax1)
ax1.set(title='Relation between temperature and users')
sns.regplot(df[df['cnt'] > 0], yoff='count', ax=ax2)
ax2.set(title='Relation between humidity and users')
```

```
Out[23]: [Text(8.5, 1.0, 'Relation between humidity and users')]
Relation between temperature and users
```

```
In [25]: from statsmodels.graphics.gcfplots import qqplot
fig, ax = plt.subplots(ncols=2, figsize=(20,6))
sns.distplot(df['count'], ax=ax1)
ax1.set(title='Distribution of the users')
qqplot(df['count'], ax=ax2, lines='s')
ax2.set(title='Theoretical quantiles')
```

```
Out[25]: [Text(8.5, 1.0, 'Theoretical quantiles')]
Distribution of the users
```

## Correlation Matrix

```
In [20]: corr = df.corr()
fig.figure(figsize=(15,10))
sns.heatmap(corr, annot=True, cmap='magma', cbar=True)
plt.xticks(rotation=45)
```

```
Out[20]: <matplotlib.axes._subplots.AxesSubplot at 0x256e06f730>
```

## One hot encoding

```
In [34]: pd.get_dummies(df['season'], prefix='season', drop_first=True)
```

```
Out[34]:
```

	season_2	season_3	season_4
0	0	0	0
1	0	0	0
2	0	0	0
3	0	0	0
4	0	0	0
...	...	...	...
17374	0	0	0
17375	0	0	0
17376	0	0	0
17377	0	0	0
17378	0	0	0

```
In [34]: df_dh = df
def one_hot_encoding(data, column):
    data = pd.concat([data, pd.get_dummies(data[column], prefix=column, drop_first=True)], axis=1)
    data = data.drop(column, axis=1)
    return data
```

```
Out[34]:
```

	temp	atemp	humidity	windspeed	casual	registered	count	season_2	season_3	season_4	month_2	month_3	month_4	month_5	month_6	month_7	month_8	month_9	month_10	month_11	month_12	hour_1	hour_2
0	0.24	0.2879	0.01	0.00	0	3	13	2.72589	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0.22	0.2727	0.00	0.00	8	32	3.688779	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
2	0.22	0.2727	0.00	0.00	5	27	3.465736	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	2
3	0.24	0.2879	0.75	0.00	3	10	3.564849	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	3
4	0.24	0.2879	0.75	0.00	0	1	0.000000	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	4

## Input Split

```
In [38]: X = df.drop(columns=['temp', 'windspeed', 'casual', 'registered', 'count'], axis=1)
y = df['count']
```

## Model Training

```
In [35]: from sklearn.linear_model import LinearRegression, Ridge, HuberRegressor, ElasticNetCV
from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor, ExtraTreesRegressor
```

```
Out[35]:
```

```
models = [LinearRegression(),
          Ridge(),
          HuberRegressor(),
          ElasticNetCV(),
          DecisionTreeRegressor(),
          RandomForestRegressor(),
          ExtraTreesRegressor(),
          GradientBoostingRegressor()]
```

```
In [40]: from sklearn.model_selection import train_test_split
def train(model):
    kfold = model_selection.KFold(n_splits=5, random_state=42)
    pred = model_selection.cross_val_score(model, X, y, cv=kfold, scoring='neg_mean_squared_error')
```

```
Out[40]:
```

```
print('CV score:', abs(cv_score))
```

```
In [41]: for model in models:
    train(model)
```

```
Out[41]:
```

```
Model: LinearRegression()
CV score: 0.612365685494837
Model: Ridge()
CV score: 0.63484945193437
Model: HuberRegressor()
CV score: 0.668329437595874
Model: ElasticNetCV()
CV score: 0.618961522396278
Model: DecisionTreeRegressor()
CV score: 0.680760449157254
Model: RandomForestRegressor()
CV score: 0.4825478722885504
Model: GradientBoostingRegressor()
CV score: 0.4743893718611764
```

```
In [42]: from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=42)
```

```
In [43]: model = RandomForestRegressor()
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
```

```
Out[43]:
```

```
error = y_test - y_pred
fig, ax = plt.subplots()
ax.scatter(y_test, error, color='black')
ax.set_xlabel('Observed')
ax.set_ylabel('Error')
plt.show()
```

```
In [46]: from sklearn.metrics import mean_squared_error
np.sqrt(mean_squared_error(y_test, y_pred))
```

```
Out[46]: 0.486780945985344
```

```
In [ ]:
```