



localhost:8888/notebooks/TNHien/Khai%20thac%20du%20lieu/Lab02/19110315\_TrinhNgocHien\_DM%20Lab02.ipynb

Web Hỗ TrợLatexLibrary GenesisTRANG NHÀ - Tủ S...Moodle HCMUSThư - TRỊNH NGỌC...Drive của tôi - Goo...Zalo WebDuolingo - Cách họ...Sử - web - viethoc.c...

Jupyter19110315\_TrinhNgocHien\_DM Lab02Last Checkpoint: 21 hours ago (autosaved)

Logout

FileEditViewInsertCellKernelWidgetsHelp

TrustedPython 3

Run

Markdown

|   | PassengerId | Survived | Pclass | Name  | Sex    | Age  | SibSp | Parch | Ticket           | Fare    | Cabin | Embarked |
|---|-------------|----------|--------|---|--------|------|-------|-------|------------------|---------|-------|----------|
| 0 | 1           | 0        | 3      | Braund, Mr. Owen Harris                           | male   | 22.0 | 1     | 0     | A/5 21171        | 7.2500  | NaN   | S        |
| 1 | 2           | 1        | 1      | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | 38.0 | 1     | 0     | PC 17599         | 71.2833 | C85   | C        |
| 2 | 3           | 1        | 3      | Heikkinen, Miss. Laina                            | female | 26.0 | 0     | 0     | STON/O2. 3101282 | 7.9250  | NaN   | S        |
| 3 | 4           | 1        | 1      | Futrelle, Mrs. Jacques Heath (Lily May Peel)      | female | 35.0 | 1     | 0     | 113803           | 53.1000 | C123  | S        |
| 4 | 5           | 0        | 3      | Allen, Mr. William Henry                          | male   | 35.0 | 0     | 0     | 373450           | 8.0500  | NaN   | S        |

The shape of data in (nrows,ncols): (891, 12)

### Function find missing percent

```
In [3]: 1 def find_missing_percent(data , showresult = True):
2         miss_df = pd.DataFrame({'ColumnName':[], 'TotalMissingVals':[], 'PercentMissing':[]})
3         for col in data.columns:
4             sum_miss_val = data[col].isnull().sum()
5             percent_miss_val = round((sum_miss_val/data.shape[0])*100,2)
6             missinginfo = {"ColumnName" : col, "TotalMissingVals" : sum_miss_val, "PercentMissing" : percent_miss_val}
7             miss_df = miss_df.append(missinginfo, ignore_index = True)
8
9         miss_df = miss_df[miss_df["PercentMissing"] > 0.0]
10        miss_df = miss_df.reset_index(drop = True)
11        miss_features = miss_df["ColumnName"].values
12        if(showresult):
13            print(data.shape)
14            display(data.head())
15            display(miss_df)
16        return miss_df

In [4]: 1 miss_df = find_missing_percent(data)
```

JupyterLab

Last Checkpoint: 21 hours ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

(891, 12)

| PassengerId | Survived | Pclass | Name | Sex  | Age    | SibSp | Parch | Ticket | Fare             | Cabin   | Embarked |   |
|-------------|----------|--------|------|--|--------|-------|-------|--------|------------------|---------|----------|---|
| 0           | 1        | 0      | 3    | Braund, Mr. Owen Harris                            | male   | 22.0  | 1     | 0      | A/5 21171        | 7.2500  | NaN      | S |
| 1           | 2        | 1      | 1    | Cumings, Mrs. John Bradley (Florence Briggs Th...) | female | 38.0  | 1     | 0      | PC 17599         | 71.2833 | C85      | C |
| 2           | 3        | 1      | 3    | Heikkinen, Miss. Laina                             | female | 26.0  | 0     | 0      | STON/O2. 3101282 | 7.9250  | NaN      | S |
| 3           | 4        | 1      | 1    | Futrelle, Mrs. Jacques Heath (Lily May Peel)       | female | 35.0  | 1     | 0      | 113803           | 53.1000 | C123     | S |
| 4           | 5        | 0      | 3    | Allen, Mr. William Henry                           | male   | 35.0  | 0     | 0      | 373450           | 8.0500  | NaN      | S |

| ColumnName | TotalMissingVals | PercentMissing |
|------------|------------------|----------------|
| 0 Age      | 177.0            | 19.87          |
| 1 Cabin    | 687.0            | 77.10          |
| 2 Embarked | 2.0              | 0.22           |

In [5]:

```
1 drop_cols = list(miss_df[miss_df['PercentMissing'] > 60.0].ColumnName)
2 print(drop_cols)
3 data = data.drop(drop_cols,axis=1)
4 miss_df = find_missing_percent(data)
```

```
['Cabin']
(891, 11)
```

| PassengerId | Survived | Pclass | Name | Sex   | Age    | SibSp | Parch | Ticket | Fare             | Embarked |   |
|-------------|----------|--------|------|---|--------|-------|-------|--------|------------------|----------|---|
| 0           | 1        | 0      | 3    | Braund, Mr. Owen Harris                           | male   | 22.0  | 1     | 0      | A/5 21171        | 7.2500   | S |
| 1           | 2        | 1      | 1    | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | 38.0  | 1     | 0      | PC 17599         | 71.2833  | C |
| 2           | 3        | 1      | 3    | Heikkinen, Miss. Laina                            | female | 26.0  | 0     | 0      | STON/O2. 3101282 | 7.9250   | S |
| 3           | 4        | 1      | 1    | Futrelle, Mrs. Jacques Heath (Lily May Peel)      | female | 35.0  | 1     | 0      | 113803           | 53.1000  | S |



19110315\_TrinhNgocHien\_

localhost:8888/notebooks/TNHNien/Khai%20thac%20d%E2%82%A5%20li%E2%82%99/Lab02/19110315\_TrinhNgocHien\_DM%20Lab02.ipynb

Web H%E2%92%81%92 Tr%E2%92%81%92 Latex Library Genesis TRANG NH%E2%92%80 - T%E2%92%81 S... Moodle HCMUS Th%E2%92%81 - TRINH NGOC... Drive c%E2%92%81 t%E2%92%81 Goo... Zalo Web Duolingo - C%E2%92%81 h%E2%92%81 S%E2%92%81 web - viethoc.c...

jupyter 19110315\_TrinhNgocHien\_DM Lab02 Last Checkpoint: 21 hours ago (autosaved)

Logout

File Edit View Insert Cell Kernel Widgets Help

Trusted Python 3

Run

Markdown

['Cabin']

(891, 11)

| PassengerId | Survived | Pclass | Name | Sex   | Age    | SibSp | Parch | Ticket | Fare             | Embarked |   |
|-------------|----------|--------|------|---|--------|-------|-------|--------|------------------|----------|---|
| 0           | 1        | 0      | 3    | Braund, Mr. Owen Harris                           | male   | 22.0  | 1     | 0      | A/5 21171        | 7.2500   | S |
| 1           | 2        | 1      | 1    | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | 38.0  | 1     | 0      | PC 17599         | 71.2833  | C |
| 2           | 3        | 1      | 3    | Heikkinen, Miss. Laina                            | female | 26.0  | 0     | 0      | STON/O2. 3101282 | 7.9250   | S |
| 3           | 4        | 1      | 1    | Futrelle, Mrs. Jacques Heath (Lily May Peel)      | female | 35.0  | 1     | 0      | 113803           | 53.1000  | S |
| 4           | 5        | 0      | 3    | Allen, Mr. William Henry                          | male   | 35.0  | 0     | 0      | 373450           | 8.0500   | S |

| ColumnName | TotalMissingVals | PercentMissing |       |
|------------|------------------|----------------|-------|
| 0          | Age              | 177.0          | 19.87 |
| 1          | Embarked         | 2.0            | 0.22  |

## Missing Handling

### Function List wise delection

In [6]:

```
1 def listwise_deletion(data):
2     for col in data.columns:
3         miss_ind = data[col][data[col].isnull()].index
4         data = data.drop(miss_ind, axis = 0)
5     return data
```

In [7]:

```
1 data_lwd = listwise_deletion(data)
2 miss_df = find_missing_percent(data_lwd)
```

0:21

ENG 8:36 AM

localhost:8888/notebooks/TNHien/Khai%20thac%20du%20liu/Lab02/19110315\_TrinhNgocHien\_DM%20Lab02.ipynb

Web Hỗ TrợLatexLibrary GenesisTRANG NHÀ - Tủ S...Moodle HCMUSThư - TRỊNH NGỌC...Drive của tôi - Goo...Zalo WebDuolingo - Cách họ...Sử - web - viethoc.c...

jupyter19110315\_TrinhNgocHien\_DM Lab02Last Checkpoint: 21 hours ago (autosaved)

Logout

FileEditViewInsertCellKernelWidgetsHelp

TrustedPython 3

Run

Markdown

In [7]:

```
1 data_lwd = listwise_deletion(data)
2 miss_df = find_missing_percent(data_lwd)
```

(712, 11)

|   | PassengerId | Survived | Pclass | Name  | Sex    | Age  | SibSp | Parch | Ticket           | Fare    | Embarked |
|---|-------------|----------|--------|---|--------|------|-------|-------|------------------|---------|----------|
| 0 | 1           | 0        | 3      | Braund, Mr. Owen Harris                           | male   | 22.0 | 1     | 0     | A/5 21171        | 7.2500  | S        |
| 1 | 2           | 1        | 1      | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | 38.0 | 1     | 0     | PC 17599         | 71.2833 | C        |
| 2 | 3           | 1        | 3      | Heikinen, Miss. Laina                             | female | 26.0 | 0     | 0     | STON/O2. 3101282 | 7.9250  | S        |
| 3 | 4           | 1        | 1      | Futrelle, Mrs. Jacques Heath (Lily May Peel)      | female | 35.0 | 1     | 0     | 113803           | 53.1000 | S        |
| 4 | 5           | 0        | 3      | Allen, Mr. William Henry                          | male   | 35.0 | 0     | 0     | 373450           | 8.0500  | S        |

| ColumnName | TotalMissingVals | PercentMissing |
|------------|------------------|----------------|
|------------|------------------|----------------|

In [8]:

```
1 numeric_cols = data.select_dtypes(['float','int']).columns
2 categoric_cols = data.select_dtypes('object').columns
3 print(f"Numeric Columns : {numeric_cols}")
4 print(f"Categoric Columns : {categoric_cols}")
```

Numeric Columns : Index(['Age', 'Fare'], dtype='object')

Categoric Columns : Index(['Name', 'Sex', 'Ticket', 'Embarked'], dtype='object')

In [9]:

```
1 def mean_imputation(data_numeric):
2     for col in data_numeric.columns:
3         mean = data_numeric[col].mean()
4         data_numeric[col] = data_numeric[col].fillna(mean)
5     return data_numeric
6 def mode_imputation(data_categoric):
7     for col in data_categoric.columns:
8         mode = data_categoric[col].mode().iloc[0]
```

localhost:8888/notebooks/TNHien/Khai%20thac%20dữ%20liệu/Lab02/19110315\_TrinhNgocHien\_DM%20Lab02.ipynb

Web Hỗ TrợLatexLibrary GenesisTRANG NHÀ - Tủ S...Moodle HCMUSThư - TRINH NGOC...Drive của tôi - Goo...Zalo WebDuolingo - Cách họ...Sử - web - viethoc.c...

jupyter19110315\_TrinhNgocHien\_DM Lab02Last Checkpoint: 21 hours ago (autosaved)

Logout

FileEditViewInsertCellKernelWidgetsHelp

TrustedPython 3

In [9]:

```
1 def mean_imputation(data_numeric):
2     for col in data_numeric.columns:
3         mean = data_numeric[col].mean()
4         data_numeric[col] = data_numeric[col].fillna(mean)
5     return data_numeric
6 def mode_imputation(data_categorical):
7     for col in data_categorical.columns:
8         mode = data_categorical[col].mode().iloc[0]
9         data_categorical[col] = data_categorical[col].fillna(mode)
10    return data_categorical
```

In [10]:

```
1 data_numeric = data[numeric_cols]
2 data_numeric_mean_imp = mean_imputation(data_numeric)
3 data_categorical = data[categorical_cols]
4 data_categorical_mode_imp = mode_imputation(data_categorical)
5
6 data_imputed_value = pd.concat([data_numeric_mean_imp, data_categorical_mode_imp], axis = 1)
7 miss_df = find_missing_percent(data_imputed_value)
```

(891, 6)

|   | Age  | Fare    | Name  | Sex    | Ticket           | Embarked |
|---|------|---------|---|--------|------------------|----------|
| 0 | 22.0 | 7.2500  | Braund, Mr. Owen Harris                           | male   | A/5 21171        | S        |
| 1 | 38.0 | 71.2833 | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | PC 17599         | C        |
| 2 | 26.0 | 7.9250  | Heikkinen, Miss. Laina                            | female | STON/O2. 3101282 | S        |
| 3 | 35.0 | 53.1000 | Futrelle, Mrs. Jacques Heath (Lily May Peel)      | female | 113803           | S        |
| 4 | 35.0 | 8.0500  | Allen, Mr. William Henry                          | male   | 373450           | S        |

ColumnNameTotalMissingValsPercentMissing

In [11]:

```
1 import xgboost
```

0:18

ENG8:37 AM



localhost:8888/notebooks/TNHien/Khai%20thac%20d%E2%82%A5%20li%E2%82%99u/Lab02/19110315\_TrinhNgocHien\_DM%20Lab02.ipynb

Web Hỗ Trợ | LaTeX | Library Genesis | TRANG NHÀ - Tủ S... | Moodle HCMUS | Thư - TRINH NGOC... | Drive của tôi - Goo... | Zalo Web | Duolingo - Cách họ... | Sử - web - viethoc.c...

jupyter 19110315\_TrinhNgocHien\_DM Lab02 Last Checkpoint: 21 hours ago (autosaved)

Logout

File Edit View Insert Cell Kernel Widgets Help

Trusted Python 3

In [11]:

```
1 import xgboost
2 from sklearn.experimental import enable_iterative_imputer
3 from sklearn.impute import IterativeImputer
4 from sklearn.preprocessing import OrdinalEncoder
5 from sklearn.ensemble import (GradientBoostingRegressor, GradientBoostingClassifier)
```

Function find missing index

Function xgboost imputation

In [12]:

```
1 def find_missing_index(data_numeric_xgboost, target_cols):
2     miss_index_dict = {}
3     for tcol in target_cols:
4         index = data_numeric_xgboost[tcol][data_numeric_xgboost[tcol].isnull()].index
5         miss_index_dict[tcol] = index
6     return miss_index_dict
7
8 def xgboost_imputation(data_numeric_xgboost, target_cols, miss_index_dict):
9     predictors = data_numeric_xgboost.drop(target_cols, axis=1)
10    for tcol in target_cols:
11        y = data_numeric_xgboost[tcol]
12        y = y.fillna(y.mean())
13        xgb = xgboost.XGBRegressor(objective="reg:squarederror", random_state=42)
14        xgb.fit(predictors, y)
15        predictions = pd.Series(xgb.predict(predictors), index=y.index)
16        index = miss_index_dict[tcol]
17        data_numeric_xgboost[tcol].loc[index] = predictions.loc[index]
18    return data_numeric_xgboost
```

In [13]:

```
1 miss_df = find_missing_percent(data, showresult = False)
2 miss_features = miss_df["ColumnName"].values
3 target_cols = [feature for feature in miss_features if feature in numeric_cols]
4 print(target_cols)
```

ttvh đã online.

0:18 ENG 8:37 AM





Browser tabs: [KHTN] NMMH HK2 21-22, Machine Learning cơ bản, (238) MIN - ĐỪNG YẾU, overleaf - Tìm trên Google, TNHien/Khai thác dữ liệu/L, 19110315\_TrinhNgocHien\_

Address bar: localhost:8888/notebooks/TNHien/Khai%20thác%20dữ%20liệu/Lab02/19110315\_TrinhNgocHien\_DM%20Lab02.ipynb

Navigation: Web Hỗ Trợ, Latex, Library Genesis, TRANG NHÀ - Tủ S..., Moodle HCMUS, Thư - TRINH NGOC..., Drive của tôi - Goo..., Zalo Web, Duolingo - Cách họ..., Sử - web - viethoc.c...

Jupyter interface: 19110315\_TrinhNgocHien\_DM Lab02, Last Checkpoint: 21 hours ago (autosaved), Logout

Menu: File, Edit, View, Insert, Cell, Kernel, Widgets, Help

Toolbar: Trusted, Python 3, Run, Stop, Restart, Clear, Help, Markdown

```
In [14]: 1 from sklearn.ensemble import GradientBoostingRegressor
2 from sklearn.preprocessing import OrdinalEncoder
3 from sklearn.ensemble import GradientBoostingClassifier
4 from sklearn.impute import IterativeImputer
5 from sklearn.preprocessing import StandardScaler
6
7 def mice_imputation_numeric(train_numeric):
8     iter_imp_numeric = IterativeImputer(GradientBoostingRegressor())
9     imputed_train = iter_imp_numeric.fit_transform(train_numeric)
10    train_numeric_imp = pd.DataFrame(imputed_train, columns = train_numeric.columns, index= train_numeric.index)
11    return train_numeric_imp
12
13 def mice_imputation_categorical(train_categorical):
14     ordinal_dict={}
15     for col in train_categorical:
16         ordinal_dict[col] = OrdinalEncoder()
17         nn_vals = np.array(train_categorical[col][train_categorical[col].notnull()]).reshape(-1,1)
18         nn_vals_arr = np.array(ordinal_dict[col].fit_transform(nn_vals)).reshape(-1,)
19         train_categorical[col].loc[train_categorical[col].notnull()] = nn_vals_arr
20
21     iter_imp_categorical = IterativeImputer(GradientBoostingClassifier(), max_iter =5, initial_strategy='most_frequent')
22     imputed_train = iter_imp_categorical.fit_transform(train_categorical)
23     train_categorical_imp = pd.DataFrame(imputed_train, columns =train_categorical.columns,index = train_categorical.index).astype
24
25     for col in train_categorical_imp.columns:
26         oe = ordinal_dict[col]
27         train_arr= np.array(train_categorical_imp[col]).reshape(-1,1)
28         train_categorical_imp[col] = oe.inverse_transform(train_arr)
29
30     return train_categorical_imp
```

```
In [15]: 1 data_numeric_imp = mice_imputation_numeric(data_numeric)
2 data_categorical_imp = mice_imputation_categorical(data_categorical)
3
4 data_imputed_mice = pd.concat([data_numeric_imp, data_categorical_imp], axis = 1)
```

Windows taskbar: 0:17, ENG, 8:39 AM

Browser tabs: [KHTN] NMMH HK2 21-22, Machine Learning cơ bản, (238) MIN - ĐỪNG YẾU, overleaf - Tìm trên Google, TNHien/Khai thác dữ liệu/L, 19110315\_TrinhNgocHien\_

Address bar: localhost:8888/notebooks/TNHien/Khai%20thác%20dữ%20liệu/Lab02/19110315\_TrinhNgocHien\_DM%20Lab02.ipynb

Navigation: Web Hỗ Trợ, Latex, Library Genesis, TRANG NHÀ - Tủ S..., Moodle HCMUS, Thư - TRINH NGOC..., Drive của tôi - Goo..., Zalo Web, Duolingo - Cách họ..., Sử - web - viethoc.c...

Jupyter interface: 19110315\_TrinhNgocHien\_DM Lab02, Last Checkpoint: 21 hours ago (autosaved), Logout

Menu: File, Edit, View, Insert, Cell, Kernel, Widgets, Help

Buttons: +, %, Copy, Paste, Up, Down, Run, Stop, Refresh, Markdown

```
In [15]: 1 data_numeric_imp = mice_imputation_numeric(data_numeric)
2 data_categorical_imp = mice_imputation_categorical(data_categorical)
3
4 data_imputed_mice = pd.concat([data_numeric_imp, data_categorical_imp], axis = 1)
5 miss_df = find_missing_percent(data_imputed_mice)

-----
KeyboardInterrupt                                Traceback (most recent call last)
<ipython-input-15-578255edd542> in <module>
      1 data_numeric_imp = mice_imputation_numeric(data_numeric)
----> 2 data_categorical_imp = mice_imputation_categorical(data_categorical)
      3
      4 data_imputed_mice = pd.concat([data_numeric_imp, data_categorical_imp], axis = 1)
      5 miss_df = find_missing_percent(data_imputed_mice)

<ipython-input-14-09a8be3e20e0> in mice_imputation_categorical(train_categorical)
     20
     21 iter_imp_categorical = IterativeImputer(GradientBoostingClassifier(), max_iter =5, initial_strategy='most_frequent')
----> 22 imputed_train = iter_imp_categorical.fit_transform(train_categorical)
     23 train_categorical_imp = pd.DataFrame(imputed_train, columns =train_categorical.columns,index = train_categorical.index).a
     24 stype(int)
     24

~\anaconda3\lib\site-packages\sklearn\impute\_iterative.py in fit_transform(self, X, y)
     655                                     feat_idx,
     656                                     abs_corr_mat)
--> 657         Xt, estimator = self._impute_one_feature(
     658             Xt, mask_missing_values, feat_idx, neighbor_feat_idx,
     659             estimator=None, fit_mode=True)

~\anaconda3\lib\site-packages\sklearn\impute\_iterative.py in _impute_one_feature(self, X_filled, mask_missing_values, feat_id
X, neighbor_feat_idx, estimator, fit_mode)
     307         y_train = _safe_indexing(X_filled[:, feat_idx],
     308                                 ~missing_row_mask)
--> 309         estimator.fit(X_train, y_train)
     310
```

Windows taskbar: 0:17, ENG, 8:39 AM

localhost:8888/notebooks/TNHien/Khai%20thac%20dữ%20liệu/Lab02/19110315\_TrinhNgocHien\_DM%20Lab02.ipynb

Web Hỗ Trợ | LaTeX | Library Genesis | TRANG NHÀ - Tủ S... | Moodle HCMUS | Thư - TRINH NGOC... | Drive của tôi - Goo... | Zalo Web | Duolingo - Cách họ... | Sử - web - viethoc.c...

jupyter 19110315\_TrinhNgocHien\_DM Lab02 Last Checkpoint: 21 hours ago (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

```
sample_weight=sample_weight)

172
193

~\anaconda3\lib\site-packages\sklearn\ensemble\_gb_losses.py in negative_gradient(self, y, raw_predictions, k, **kwargs)
718     """
719     return y - np.nan_to_num(np.exp(raw_predictions[:, k] -
--> 720                             logsumexp(raw_predictions, axis=1)))
721
722     def _update_terminal_region(self, tree, terminal_regions, leaf, X, y,

~\anaconda3\lib\site-packages\scipy\special\_logsumexp.py in logsumexp(a, axis, b, keepdims, return_sign)
112     # suppress warnings about log of zero
113     with np.errstate(divide='ignore'):
--> 114         s = np.sum(tmp, axis=axis, keepdims=keepdims)
115         if return_sign:
116             sgn = np.sign(s)

<__array_function__ internals> in sum(*args, **kwargs)

~\anaconda3\lib\site-packages\numpy\core\fromnumeric.py in sum(a, axis, dtype, out, keepdims, initial, where)
2245     return res
2246
--> 2247     return _wrapreduction(a, np.add, 'sum', axis, dtype, out, keepdims=keepdims,
2248                           initial=initial, where=where)
2249

~\anaconda3\lib\site-packages\numpy\core\fromnumeric.py in _wrapreduction(obj, ufunc, method, axis, dtype, out, **kwargs)
85     return reduction(axis=axis, out=out, **passkwargs)
86
---> 87     return ufunc.reduce(obj, axis, dtype, out, **passkwargs)
88
89

KeyboardInterrupt:
```

0:16 ENG 8:39 AM



The screenshot displays a Jupyter Notebook environment running locally at localhost:8888. The notebook file is named "19110315\_TrinhNgocHien\_DM Lab02.ipynb". The interface includes a top toolbar with menus like File, Edit, View, Insert, Cell, Kernel, Widgets, and Help. Below the menu bar are icons for saving, adding new files, undo, redo, and running code.

The notebook content shows three input cells:

- In [16]:

```
1 data_modelling = data_lwd.copy()
2
```
- In [17]:

```
1 data_modelling
```

The output of In [17] is displayed as a table with 11 columns: PassengerId, Survived, Pclass, Name, Sex, Age, SibSp, Parch, Ticket, Fare, and Embarked. The first few rows of the table are shown, followed by an ellipsis indicating more rows, and then rows 885 through 890.

|     | PassengerId | Survived | Pclass | Name   | Sex    | Age  | SibSp | Parch | Ticket           | Fare    | Embarked |
|-----|-------------|----------|--------|--|--------|------|-------|-------|------------------|---------|----------|
| 0   | 1           | 0        | 3      | Braund, Mr. Owen Harris                            | male   | 22.0 | 1     | 0     | A/5 21171        | 7.2500  | S        |
| 1   | 2           | 1        | 1      | Cummings, Mrs. John Bradley (Florence Briggs Th... | female | 38.0 | 1     | 0     | PC 17599         | 71.2833 | C        |
| 2   | 3           | 1        | 3      | Heikkinen, Miss. Laina                             | female | 26.0 | 0     | 0     | STON/O2. 3101282 | 7.9250  | S        |
| 3   | 4           | 1        | 1      | Futrelle, Mrs. Jacques Heath (Lily May Peel)       | female | 35.0 | 1     | 0     | 113803           | 53.1000 | S        |
| 4   | 5           | 0        | 3      | Allen, Mr. William Henry                           | male   | 35.0 | 0     | 0     | 373450           | 8.0500  | S        |
| ... | ...         | ...      | ...    | ...  | ...    | ...  | ...   | ...   | ...              | ...     | ...      |
| 885 | 886         | 0        | 3      | Rice, Mrs. William (Margaret Norton)               | female | 39.0 | 0     | 5     | 382652           | 29.1250 | Q        |
| 886 | 887         | 0        | 2      | Montvila, Rev. Juozas                              | male   | 27.0 | 0     | 0     | 211536           | 13.0000 | S        |
| 887 | 888         | 1        | 1      | Graham, Miss. Margaret Edith                       | female | 19.0 | 0     | 0     | 112053           | 30.0000 | S        |
| 889 | 890         | 1        | 1      | Behr, Mr. Karl Howell                              | male   | 26.0 | 0     | 0     | 111369           | 30.0000 | C        |
| 890 | 891         | 0        | 3      | Dooley, Mr. Patrick                                | male   | 32.0 | 0     | 0     | 370376           | 7.7500  | Q        |

Below the table, it states "712 rows x 11 columns".

The third input cell, In [18], contains a function definition:

```
1 def FeatureEngineering(data_modelling):
2     data_modelling['Total_Age'] = data_modelling['Age'].sum()
3     data_modelling['Total_Survived_for_all'] = data_modelling['Survived'].sum()
4     data_modelling['Total_Fare_for_all'] = data_modelling['Fare'].sum()
5     data_modelling['Total_Age_and_Fare'] = data_modelling['Age'] + data_modelling['Fare']
6     return data_modelling
```

localhost:8888/notebooks/TNHien/Khai%20thac%20du%20lieu/Lab02/19110315\_TrinhNgocHien\_DM%20Lab02.ipynb

Web Hỗ TrợLatexLibrary GenesisTRANG NHÀ - Tủ S...Moodle HCMUSThư - TRINH NGOC...Drive của tôi - Goo...Zalo WebDuolingo - Cách họ...Sử - web - viethoc.c...

jupyter19110315\_TrinhNgocHien\_DM Lab02Last Checkpoint: 21 hours ago (autosaved)

Logout

FileEditViewInsertCellKernelWidgetsHelp

TrustedPython 3

In [18]:

```
1 def FeatureEngineering(data_modelling):
2     data_modelling['Total_Age'] = data_modelling['Age'].sum()
3     data_modelling['Total_Survived_for_all'] = data_modelling['Survived'].sum()
4     data_modelling['Total_Fare_for_all'] = data_modelling['Fare'].sum()
5     data_modelling['Total_Age_and_Fare'] = data_modelling['Age'] + data_modelling['Fare']
6     return data_modelling
7
8 data_modelling = FeatureEngineering(data_modelling)
9 display(data_modelling.head())
10 print(data_modelling.shape)
```

|   | PassengerId | Survived | Pclass | Name | Sex   | Age    | SibSp | Parch | Ticket | Fare             | Embarked | Total_Age | Total_Survived_for_all | Total_Fare_for_all | Tr        |
|---|-------------|----------|--------|------|---|--------|-------|-------|--------|------------------|----------|-----------|------------------------|--------------------|-----------|
| 0 |             | 1        | 0      | 3    | Braund, Mr. Owen Harris                           | male   | 22.0  | 1     | 0      | A/5 21171        | 7.2500   | S         | 21105.17               | 288                | 24611.883 |
| 1 |             | 2        | 1      | 1    | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | 38.0  | 1     | 0      | PC 17599         | 71.2833  | C         | 21105.17               | 288                | 24611.883 |
| 2 |             | 3        | 1      | 3    | Heikkinen, Miss. Laina                            | female | 26.0  | 0     | 0      | STON/O2. 3101282 | 7.9250   | S         | 21105.17               | 288                | 24611.883 |
| 3 |             | 4        | 1      | 1    | Futrelle, Mrs. Jacques Heath (Lily May Peel)      | female | 35.0  | 1     | 0      | 113803           | 53.1000  | S         | 21105.17               | 288                | 24611.883 |
| 4 |             | 5        | 0      | 3    | Allen, Mr. William Henry                          | male   | 35.0  | 0     | 0      | 373450           | 8.0500   | S         | 21105.17               | 288                | 24611.883 |

(712, 15)

0:15

ENG8:42 AM

The screenshot displays a Jupyter Notebook environment running locally at localhost:8888. The notebook file is named "19110315\_TrinhNgocHien\_DM Lab02.ipynb". The interface includes a top toolbar with standard file operations and a bottom toolbar with execution controls like Run, Step, and Restart.

The notebook contains two input cells:

In [19]:

```
1 skew_limit = 0.5  
2 skew_vals = data_modelling[numeric_cols].skew()  
3 skew_cols = (skew_vals  
4               .sort_values(ascending=False)  
5               .to_frame()  
6               .rename(columns={0:'Skew'}))  
7               .query('abs(Skew) > {}'.format(skew_limit)))  
8 display(skew_cols.T)
```

This cell outputs a DataFrame titled "Fare":

|      | Fare     |
|------|----------|
| Skew | 4.667009 |

In [20]:

```
1 fig, (ax_positive, ax_target, ax_negative) = plt.subplots(1, 3, figsize=(15, 5))  
2 sns.histplot(data_modelling['Age'], kde=True, stat='density', linewidth=0, color = '#236AB9', ax=ax_positive)  
3 sns.histplot(data_modelling['Fare'], kde=True, stat='density', linewidth=0, color = 'blue', ax=ax_target)  
4 sns.histplot(data_modelling['Parch'], kde=True, stat='density', linewidth=0,color='#B85B14', ax=ax_negative)  
5 plt.show()
```

The output of In [20] consists of three side-by-side density plots:

- The first plot shows the density distribution for "Age", with a peak around 30-35 years, colored blue (#236AB9).
- The second plot shows the density distribution for "Fare", with a sharp peak at low fare values, colored blue.
- The third plot shows the density distribution for "Parch" (number of passengers), with a peak at 0 passengers, colored orange (#B85B14).





Browser tabs: [KHTN] NMMH HK2 21-22, Machine Learning cơ bản, (238) Chi Pu | ANH OI, overleaf - Tìm trên Google, TNHien/Khai thác dữ liệu/L, 19110315\_TrinhNgocHien\_

URL: localhost:8888/notebooks/TNHien/Khai%20thác%20dữ%20liệu/Lab02/19110315\_TrinhNgocHien\_DM%20Lab02.ipynb

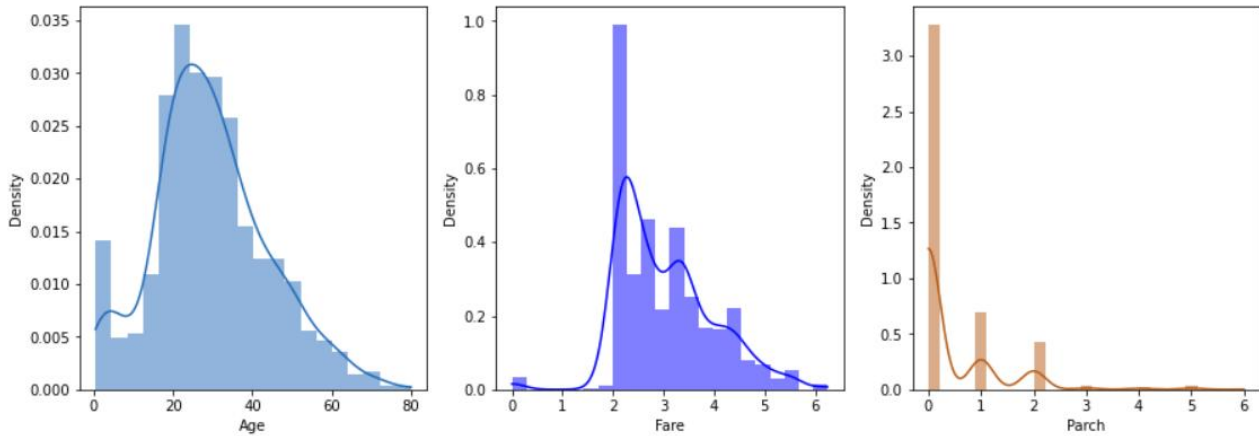
Navigation: Web Hỗ Trợ, Latex, Library Genesis, TRANG NHÀ - Tủ S..., Moodle HCMUS, Thư - TRINH NGOC..., Drive của tôi - Goo..., Zalo Web, Duolingo - Cách họ..., Sử - web - viethoc.c...

Jupyter 19110315\_TrinhNgocHien\_DM Lab02 Last Checkpoint: 21 hours ago (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

In [22]:

```
1 fig, (ax_positive, ax_target, ax_negative) = plt.subplots(1, 3, figsize=(15, 5))
2 sns.histplot(data_modelling['Age'],kde=True, stat='density', linewidth=0, color = '#236AB9', ax=ax_positive)
3 sns.histplot(data_modelling['Fare'],kde=True, stat='density', linewidth=0, color = 'blue', ax=ax_target)
4 sns.histplot(data_modelling['Parch'], kde=True, stat='density', linewidth=0,color='#B85B14', ax=ax_negative)
5 plt.show()
```



In [23]:

```
1 def FeatureEncoding(data_modelling):
2     data_modelling = pd.get_dummies(data_modelling, columns=categoric_cols, drop_first=True)
3     return data_modelling
4
5 data_modelling = FeatureEncoding(data_modelling)
6 display(data_modelling.head())
7 print(data_modelling.shape)
```

Windows taskbar: 0:14, ENG, 8:43 AM

localhost:8888/notebooks/TNHien/Khai%20thac%20dữ%20liệu/Lab02/19110315\_TrinhNgocHien\_DM%20Lab02.ipynb

Web Hỗ TrợLatexLibrary GenesisTRANG NHÀ - Tủ S...Moodle HCMUSThư - TRINH NGOC...Drive của tôi - Goo...Zalo WebDuolingo - Cách họ...Sử - web - viethoc.c...

jupyter19110315\_TrinhNgocHien\_DM Lab02Last Checkpoint: 21 hours ago (autosaved)

Logout

FileEditViewInsertCellKernelWidgetsHelpTrustedPython 3

RunCode

|   | PassengerId | Survived | Pclass | Age  | SibSp | Parch | Fare     | Total_Age | Total_Survived_for_all | Total_Fare_for_all | Total_Age_and_Fare | Name_Abbott, Mr. Rossmore Edward | Name, Mrs. (Ro |
|---|-------------|----------|--------|------|-------|-------|----------|-----------|------------------------|--------------------|--------------------|----------------------------------|----------------|
| 0 | 1           | 0        | 3      | 22.0 | 1     | 0     | 2.110213 | 21105.17  | 288                    | 24611.883          | 29.2500            | 0                                |                |
| 1 | 2           | 1        | 1      | 38.0 | 1     | 0     | 4.280593 | 21105.17  | 288                    | 24611.883          | 109.2833           | 0                                |                |
| 2 | 3           | 1        | 3      | 26.0 | 0     | 0     | 2.188856 | 21105.17  | 288                    | 24611.883          | 33.9250            | 0                                |                |
| 3 | 4           | 1        | 1      | 35.0 | 1     | 0     | 3.990834 | 21105.17  | 288                    | 24611.883          | 88.1000            | 0                                |                |
| 4 | 5           | 0        | 3      | 35.0 | 0     | 0     | 2.202765 | 21105.17  | 288                    | 24611.883          | 43.0500            | 0                                |                |

5 rows × 1265 columns

(712, 1265)

In [24]:

```
1 from sklearn.linear_model import Ridge, RidgeCV, Lasso, LassoCV
2 from sklearn.model_selection import train_test_split
3 from sklearn.metrics import mean_squared_error, r2_score
4
5 def DataSplitTrainTest(data_modelling):
6     train = data_modelling.copy()
7     X = train.drop('Fare', axis=1)
8     y = train['Fare']
9     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=12345)
10    print("Train Data", X_train.shape)
11    print("Test Data", X_test.shape)
12    return X_train, X_test, y_train, y_test
13
14 X_train, X_test, y_train, y_test = DataSplitTrainTest(data_modelling)
```

Train Data (498, 1264)  
Test Data (214, 1264)

0:13ENG8:43 AM



Browser tabs: (3) [KHTN] NMMH HK2 21... Machine Learning cơ bản... (238) Aloha - Cool || Pi... overleaf - Tìm trên Google... TNHien/Khai thác dữ liệu/L... 19110315\_TrinhNgocHien\_...

Address bar: localhost:8888/notebooks/TNHien/Khai%20thác%20dữ%20liệu/Lab02/19110315\_TrinhNgocHien\_DM%20Lab02.ipynb

Navigation: Web Hỗ Trợ, Latex, Library Genesis, TRANG NHÀ - Tủ S..., Moodle HCMUS, Thư - TRINH NGOC..., Drive của tôi - Goo..., Zalo Web, Duolingo - Cách họ..., Sử - web - viethoc.c...

Jupyter 19110315\_TrinhNgocHien\_DM Lab02 Last Checkpoint: 21 hours ago (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

In [25]:

```
1 def BuildLassoModel(X_train, X_test, y_train, y_test):
2     lasso = Lasso(max_iter = 100000, normalize = True)
3     lassocv = LassoCV(alphas = None, cv = 10, max_iter = 100000, normalize = True)
4     lassocv.fit(X_train, y_train)
5
6     lasso.set_params(alpha=lassocv.alpha_)
7     lasso.fit(X_train, y_train)
8
9     print('The Lasso:')
10    print("Alpha =", lassocv.alpha_)
11    print("RMSE =", mean_squared_error(y_test, lasso.predict(X_test), squared=False))
12    print("R2 Score =", r2_score(y_test, lasso.predict(X_test)))
13    return lasso
14
15 lasso = BuildLassoModel(X_train, X_test, y_train, y_test)
```

The Lasso:  
Alpha = 3.335189513517674e-05  
RMSE = 0.2743078224257173  
R2 Score = 0.9089543698393528

In [26]:

```
1 def BuildRidgeModel(X_train, X_test, y_train, y_test):
2     alphas = np.geomspace(1e-9, 5, num=100)
3     ridgecv = RidgeCV(alphas = alphas, scoring = 'neg_mean_squared_error', normalize = True)
4     ridgecv.fit(X_train, y_train)
5
6     ridge = Ridge(alpha = ridgecv.alpha_, normalize = True)
7     ridge.fit(X_train, y_train)
8
9     print('Ridge Regression:')
10    print("Alpha =", ridgecv.alpha_)
11    print("RMSE =", mean_squared_error(y_test, ridge.predict(X_test), squared=False))
12    print("R2 Score =", r2_score(y_test, lasso.predict(X_test)))
13    return ridge
14
15 ridge = BuildRidgeModel(X_train, X_test, y_train, y_test)
```

Windows taskbar: A/C, 9:03 AM

The screenshot displays a web-based Jupyter Notebook environment. The browser's address bar shows the URL: localhost:8888/notebooks/TNHien/Khai%20thac%20du%20lieu/Lab02/19110315\_TrinhNgocHien\_DM%20Lab02.ipynb. The notebook title is "19110315\_TrinhNgocHien\_DM Lab02". The top navigation bar includes tabs for File, Edit, View, Insert, Cell, Kernel, Widgets, and Help. Below this is a toolbar with icons for saving, adding cells, undo, redo, running, and other functions. The main area contains two code cells. The first cell, labeled "In [26]:", defines a function named BuildRidgeModel and executes it, displaying the following output:

The Lasso:  
Alpha = 3.335189513517674e-05  
RMSE = 0.2743078224257173  
R2 Score = 0.9089543698393528

```
def BuildRidgeModel(X_train, X_test, y_train, y_test):  
    alphas = np.geomspace(1e-9, 5, num=100)  
    ridgecv = RidgeCV(alphas = alphas, scoring = 'neg_mean_squared_error', normalize = True)  
    ridgecv.fit(X_train, y_train)  
  
    ridge = Ridge(alpha = ridgecv.alpha_, normalize = True)  
    ridge.fit(X_train, y_train)  
  
    print('Ridge Regression:')  
    print("Alpha =", ridgecv.alpha_)  
    print("RMSE =", mean_squared_error(y_test, ridge.predict(X_test), squared=False))  
    print("R2 Score = ", r2_score(y_test, lasso.predict(X_test)))  
    return ridge  
  
ridge = BuildRidgeModel(X_train, X_test, y_train, y_test)
```

Ridge Regression:  
Alpha = 1e-09  
RMSE = 0.438342413483562  
R2 Score = 0.9089543698393528

The second cell, labeled "In [ ]:", is currently empty and shows a single line number "1". The bottom status bar indicates the system time as 9:03 AM and the language as ENG.