

```
# Loading and visualizing the dataset
import pandas as pd
df = pd.read_csv("C:/Users/Ashtin/Downloads/glass.csv")
df.head()
```

```
In [1]:
```

	Ri	Na	Mg	Al	Si	K	Ca	Ba	Fe	Type
0	152101	1354	4.49	1.10	71.78	0.06	8.75	0.00	0.0	1
1	151761	1790	3.60	1.96	72.73	0.48	7.83	0.00	0.0	1
2	151618	1353	3.55	1.54	72.99	0.39	7.78	0.00	0.0	1
3	151766	1321	3.69	1.29	72.61	0.57	8.22	0.00	0.0	1
4	151742	1327	3.62	1.24	73.08	0.55	8.07	0.00	0.0	1

```
In [3]:
```

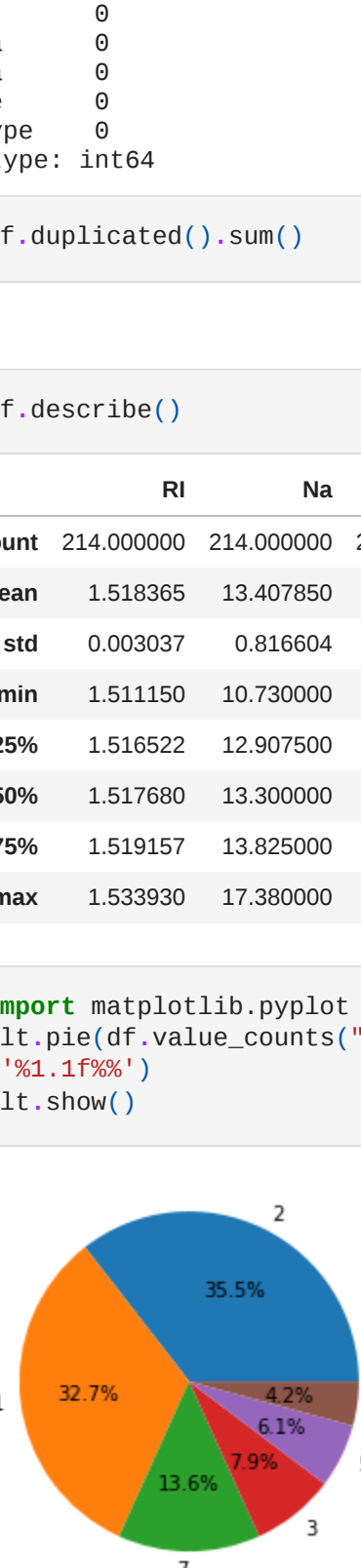
```
df.isnull().sum()
```

```
Out[3]:
```

	Ri	Na	Mg	Al	Si	K	Ca	Ba	Fe	Type
count	214.000000	214.000000	214.000000	214.000000	214.000000	214.000000	214.000000	214.000000	214.000000	214.000000
mean	1513665	13407850	2.684533	1.444907	72.650935	0.497056	8.956963	0.175047	0.057009	2.780374
std	0.000307	0.816604	1.442408	0.499200	0.774546	0.652192	1.423153	0.497219	0.097439	2.103739
min	1511150	10.730000	0.000000	0.290000	69.810000	0.000000	5.430000	0.000000	0.000000	1.000000
25%	1511622	12.907500	2.115000	1.180000	72.280000	0.122500	8.240000	0.000000	0.000000	1.000000
50%	1517680	13.300000	3.480000	1.360000	72.790000	0.555000	8.600000	0.000000	0.000000	2.000000
75%	1519157	13.825000	3.600000	1.630000	73.087500	0.610000	9.172500	0.000000	0.100000	3.000000
max	1533930	17.380000	4.450000	3.500000	75.410000	6.210000	16.390000	3.150000	0.510000	7.000000

```
In [6]:
```

```
import matplotlib.pyplot as plt
plt.pie(df.value_counts("Type"), labels=df.value_counts("Type").index, autopct='%1.1f%%')
plt.show()
```



```
In [7]:
```

```
df_2 = df.copy()
df_2 = df_2.drop_duplicates()
df_2
```

```
Out[7]:
```

	Ri	Na	Mg	Al	Si	K	Ca	Ba	Fe	Type
0	152101	13.64	4.49	1.10	71.78	0.06	8.75	0.00	0.0	1
1	151761	13.89	3.60	1.96	72.73	0.48	7.83	0.00	0.0	1
2	151618	13.53	3.55	1.54	72.99	0.39	7.78	0.00	0.0	1
3	151766	13.21	3.69	1.29	72.61	0.57	8.22	0.00	0.0	1
4	151742	13.27	3.62	1.24	73.08	0.55	8.07	0.00	0.0	1
...	...	...	...	...	...	...	...	...	...	...
209	151623	14.14	0.00	2.88	72.61	0.08	9.18	1.06	0.0	7
210	151685	14.82	0.00	1.95	71.06	0.00	8.40	1.99	0.0	7
211	152065	14.36	0.00	2.02	73.42	0.00	8.44	1.64	0.0	7
212	151851	14.38	0.00	1.94	73.61	0.00	8.48	1.57	0.0	7
213	151711	14.23	0.00	2.08	73.86	0.00	8.62	1.67	0.0	7

213 rows x 10 columns

```
In [8]:
```

```
df_2.groupby(["Type"]).mean()
```

```
Out[8]:
```

Type	Ri	Na	Mg	Al	Si	K	Ca	Ba	Fe
1	1518669	13228261	3.548951	1.179413	72.631449	0.452319	8.786007	0.013999	0.057928
2	1518619	1311071	3.002105	1.408156	72.520026	0.521053	9.079594	0.050263	0.079737
3	1517654	13437059	3.543329	1.201176	72.424706	0.406471	8.782941	0.008024	0.057059
4	1518938	12826921	0.973986	2.038466	72.963154	1.470000	10.125846	0.367692	0.066769
5	1517466	14.046667	1.305556	1.366667	72.206667	0.200000	9.266667	0.000000	0.000000
7	1517116	14.442069	0.58276	2.122759	72.959882	0.325172	8.451379	1.040000	0.013448

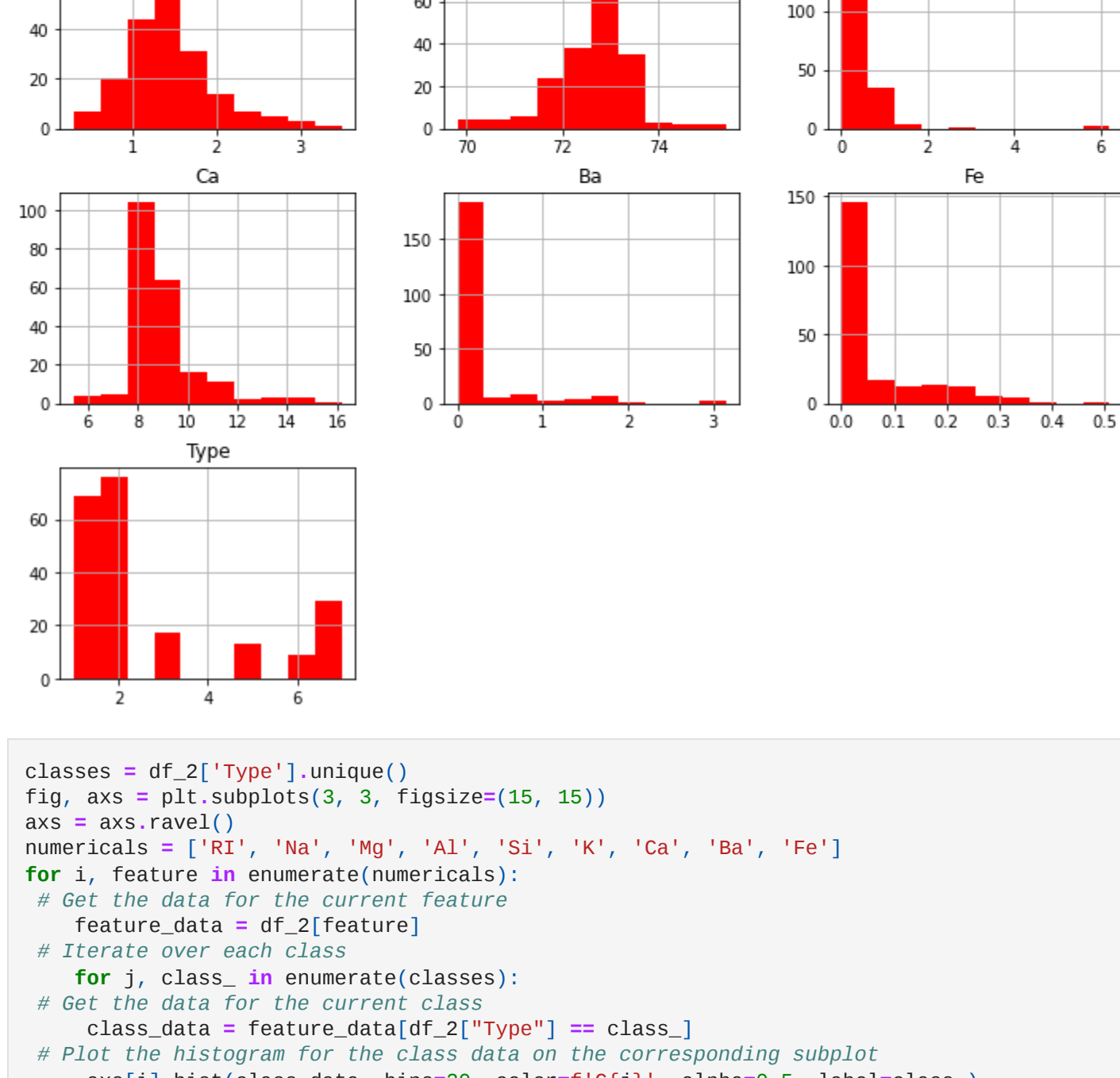
```
In [10]:
```

```
grouped = df_2.groupby(["Type"])
mean_group = grouped.mean()
# Create subplots
colors = {"1": "red", "2": "green", "3": "blue", "5": "pink", "6": "yellow", "7": "orange"}
fig, axes = plt.subplots(nrows=3, ncols=3, figsize=(15, 15))
# Iterate through the columns of mean_group and make a bar plot for each column
for i, column in enumerate(mean_group):
    mean_group[column].plot(ax=axes[i//3, i%3], kind='bar')
    axes[i//3, i%3].set_title(column)
plt.subplots_adjust(hspace=0.5)
plt.show()
```



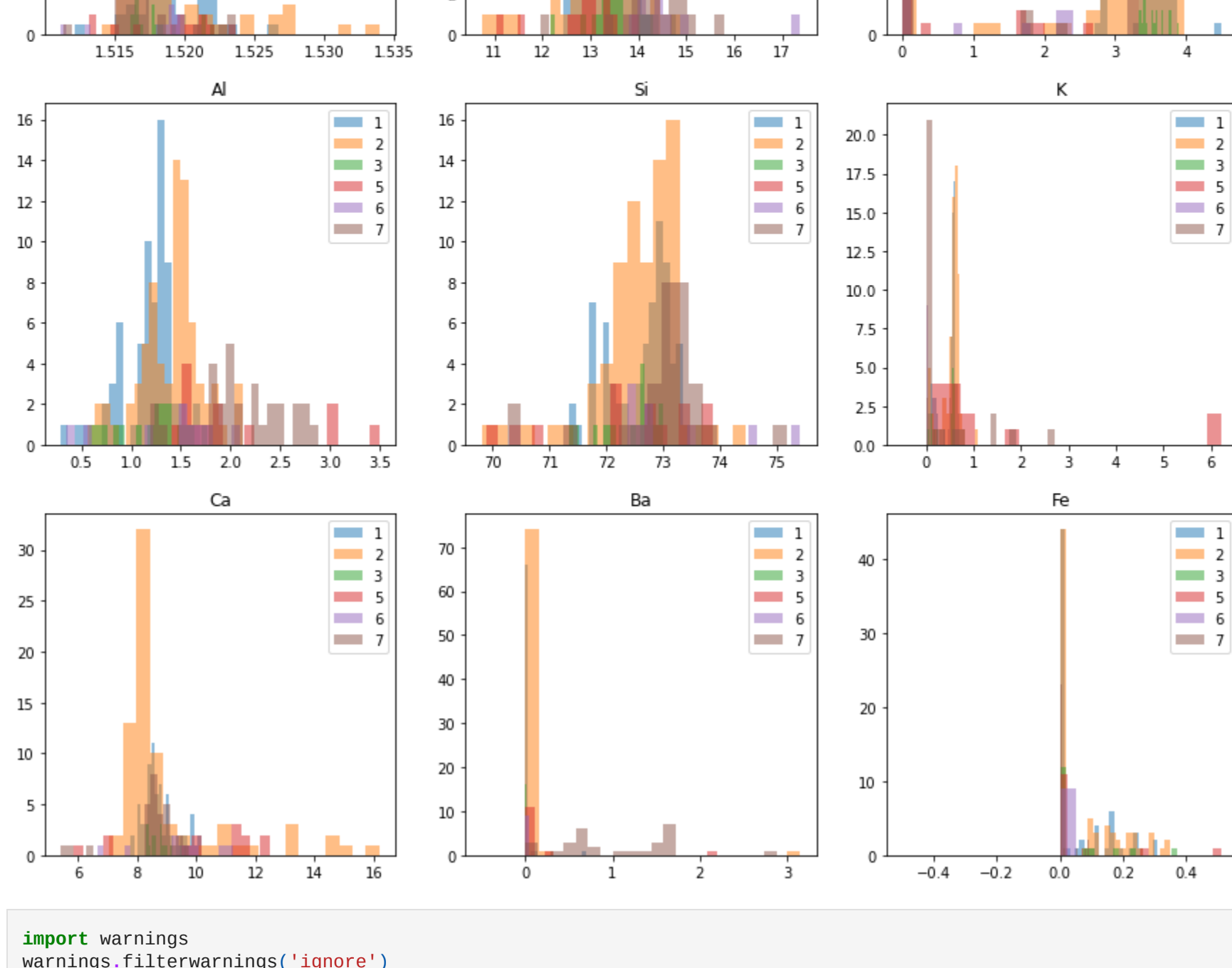
```
In [11]:
```

```
df_2.hist(figsize=(12,12), color='red')
plt.show()
```



```
In [15]:
```

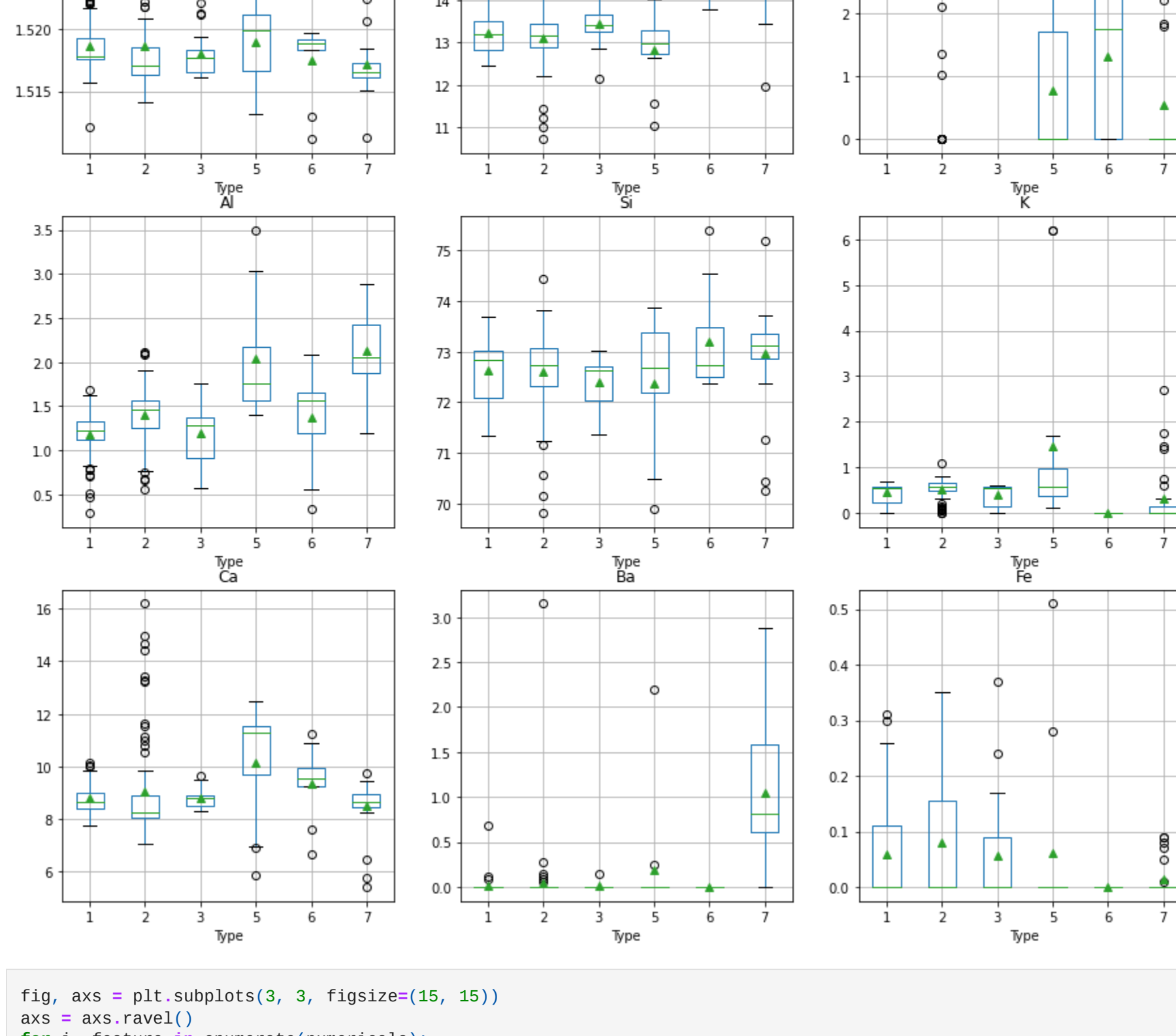
```
classes = df_2["Type"].unique()
fig, axs = plt.subplots(3, 3, figsize=(15, 15))
axs = axs.ravel()
numericals = ["Ri", "Na", "Mg", "Al", "Si", "K", "Ca", "Ba", "Fe"]
for i, feature in enumerate(numericals):
    # Get the data for the current feature
    feature_data = df_2[feature]
    # Iterate over each class
    for j, class_ in enumerate(classes):
        # Get the data for the current class
        class_data = feature_data[feature_data["Type"] == class_]
        # Plot the histogram for the class data on the corresponding subplot
        axs[i].hist(class_data, bins=20, color=f'c[{j}]', alpha=0.5, label=class_)
        axs[i].set_title(feature)
        axs[i].legend()
plt.show()
```



```
In [16]:
```

```
import warnings
warnings.filterwarnings('ignore')
fig, axs = plt.subplots(3, 3, figsize=(15, 15))
axs = axs.ravel()
numericals = ["Ri", "Na", "Mg", "Al", "Si", "K", "Ca", "Ba", "Fe"]
for i, feature in enumerate(numericals):
    # Get the data for the current feature
    feature_data = df_2[feature]
    # Plot the boxplot for the feature data on the corresponding subplot
    feature_data.boxplot(by="Type", ax=axs[i], positions=[j+1 for j in range(len(classes))], widths=0.5, showmeans=True)
    axs[i].set_title(feature)
    axs[i].set_xlabel("Type")
plt.show()
```

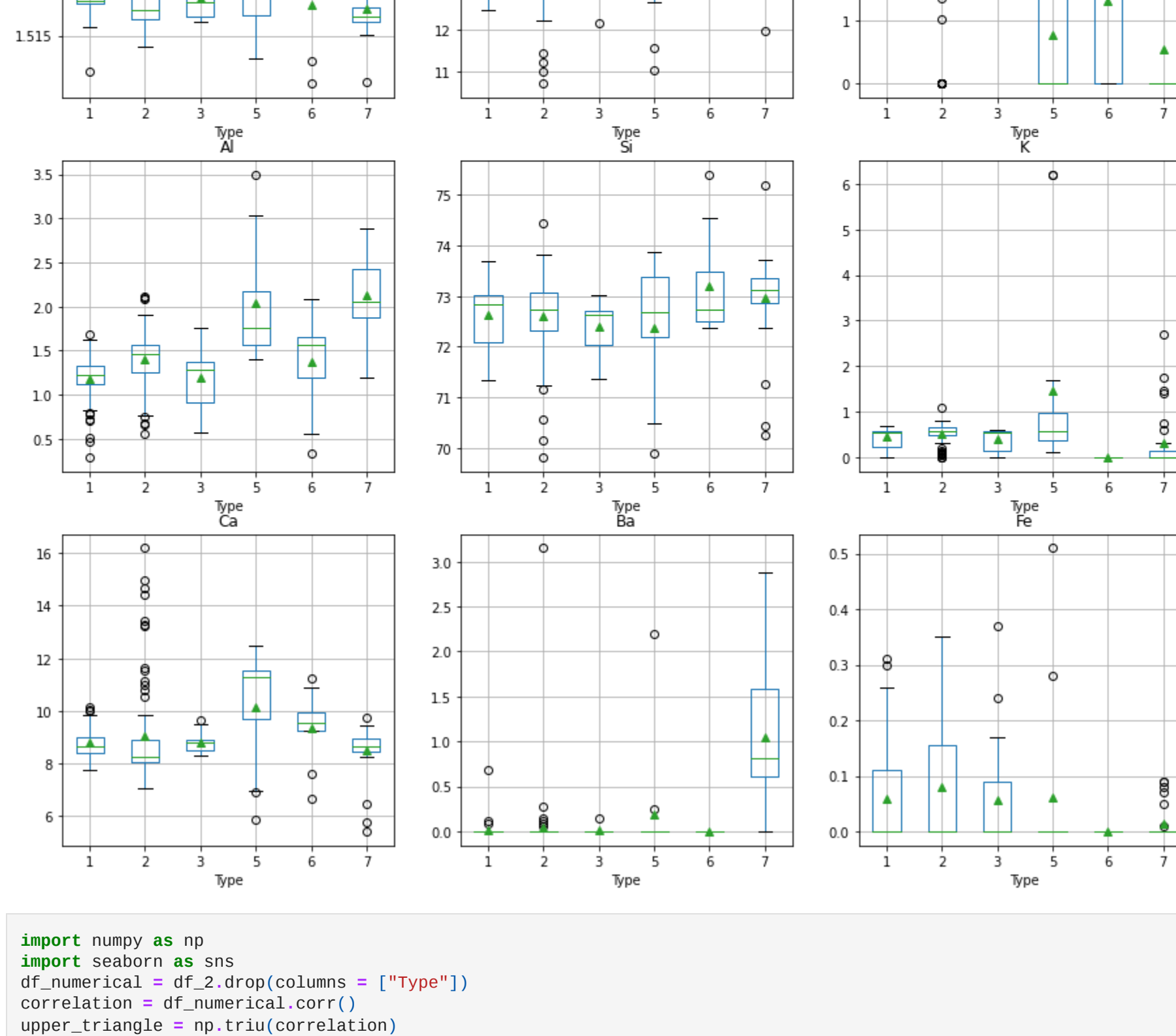
Boxplot grouped by Type



```
In [22]:
```

```
fig, axs = plt.subplots(3, 3, figsize=(15, 15))
axs = axs.ravel()
for i, feature in enumerate(numericals):
    # Get the data for the current feature
    feature_data = df_2[feature]
    # Plot the boxplot for the feature data on the corresponding subplot
    feature_data.boxplot(by="Type", ax=axs[i], positions=[j+1 for j in range(len(classes))], widths=0.5, showmeans=True)
    axs[i].set_title(feature)
    axs[i].set_xlabel("Type")
plt.show()
```

Boxplot grouped by Type



```
In [24]:
```

```
import numpy as np
import seaborn as sns
df_numerical = df_2.drop(columns = ["Type"])
correlation = df_numerical.corr()
upper_triangle = np.triu(correlation)
fig=plt.gcf()
fig.set_size_inches(10,10)
sns.heatmap(data=correlation,mask=upper_triangle,square=True,annot=True)
```

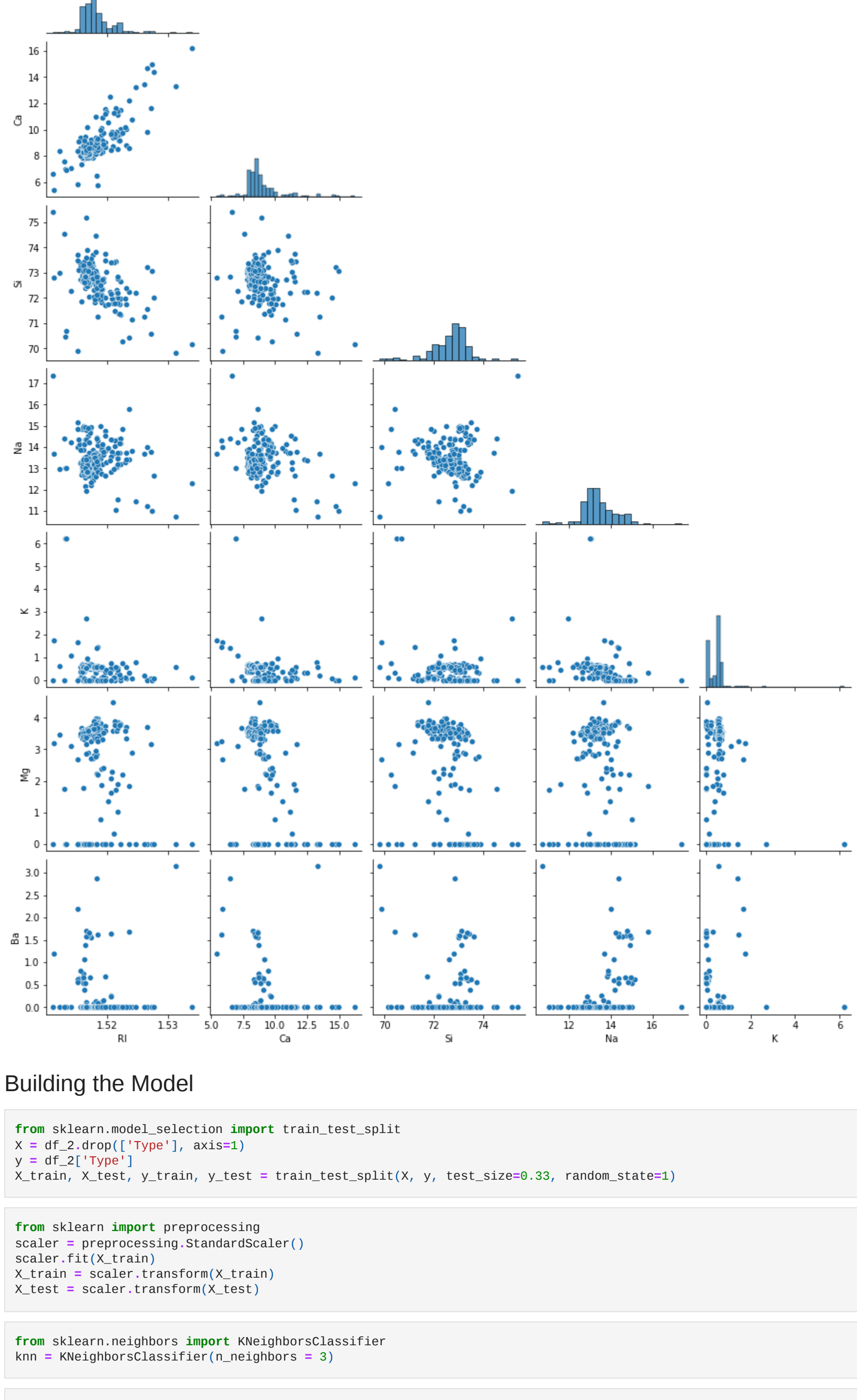


```
In [25]:
```

```
sns.pairplot(df_2[["Ri","Ca","Si","Na","K","Mg","Ba"]], corner=True)
```

```
Out[25]:
```

```
<seaborn.axisgrid.PairGrid at 0x1abbb8c6a0>
```



## Building the Model

```
In [26]:
```

```
from sklearn.model_selection import train_test_split
X = df_2.drop(["Type"], axis=2)
y = df_2["Type"]
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=1)
```

```
In [27]:
```

```
from sklearn import preprocessing
scaler = preprocessing.StandardScaler()
scaler.fit(X_train)
X_train = scaler.transform(X_train)
X_test = scaler.transform(X_test)
```

```
In [28]:
```

```
from sklearn.neighbors import KNeighborsClassifier
knn = KNeighborsClassifier(n_neighbors = 3)
```

```
In [29]:
```

```
from sklearn.model_selection import cross_validate
cv_result = cross_validate(knn, X_train, y_train, cv=10)
cv_result
```

```
Out[29]:
```

```
{'fit_time': array([0.03700042, 0.00206367, 0.00199238, 0.00198865, 0.00099993,
0.00209551, 0.00199842, 0.00296665, 0.00201392, 0.00199676]),
'score_time': array([0.89762924, 0.00199795, 0.00200772, 0.00180817, 0.00208462,
0.00209311, 0.00299854, 0.00398676, 0.00285955, 0.00300881]),
'test_score': array([0.53333333, 0.53333333, 0.5       , 0.78571429, 0.71428571,
0.57142857, 0.64285714, 0.57142857, 0.64285714, 0.64285714])}
```

```
In [30]:
```

```
scores = cv_result["test_score"]
print("The mean accuracy using cross validation is: ")
print((scores.mean()::3f) * (scores.std()::3f))

The mean accuracy using cross validation is: 0.614 ± 0.085
```

## Hyperparameter Tuning

```
In [31]:
```

```
from sklearn.model_selection import GridSearchCV
param_grid = {
    "n_neighbors": [3, 5, 7, 9],
    "weights": ["uniform", "distance"],
}
model_grid_search = GridSearchCV(
    knn, param_grid=param_grid, n_jobs=2, cv=10
)
model_grid_search.fit(X_train, y_train)
```

```
Out[31]:
```

```
GridSearchCV(cv=10, estimator=KNeighborsClassifier(n_neighbors=3), n_jobs=2,
              param_grid={'n_neighbors': [3, 5, 7, 9],
                           'weights': ['uniform', 'distance']})
```

```
In [32]:
```

```
cv_results = pd.DataFrame(model_grid_search.cv_results_)
cv_results
```

```
print("Accuracy score: " + str(accuracy_score(y_test, y_pred_final)))
print("Balanced accuracy score: " + str(balanced_accuracy_score(y_test, y_pred_final)))
print("\nConfusion matrix: \n" + str(confusion_matrix(y_test, y_pred_final)))
print("\nClassification report: \n" + str(classification_report(y_test, y_pred_final)))
```

---

Accuracy score: 0.7746478073230436  
Balanced accuracy score: 0.7345695758732737

Confusion matrix:

```
[[ 12  1  0  0  0  0]
 [ 2 20  0  0  1  0]
 [ 3  2  2  0  0  0]
 [ 0  1  2  0  0  0]
 [ 0  0  0  3  0  0]
 [ 1  1  0  0  0  7]]
```

Classification report:

	precision	recall	f1-score	support
1	0.78	0.81	0.79	26
2	0.80	0.87	0.83	23
3	0.23	0.29	0.22	7
5	1.00	0.67	0.80	3
6	0.75	1.00	0.86	3
7	1.00	0.78	0.88	9

accuracy

0.77

71

```
In [33]:
```

```
from sklearn.metrics import classification_report, accuracy_matrix, balanced_accuracy_score
knn_final = KNeighborsClassifier(n_neighbors = 3, weights = "distance")
knn_final.fit(X_train, y_train)
y_pred_final = knn_final.predict(X_test)
print("Accuracy score: " + str(accuracy_score(y_test, y_pred_final)))
print("Balanced accuracy score: " + str(balanced_accuracy_score(y_test, y_pred_final)))
print("Unconfusion matrix: \n" + str(confusion_matrix(y_test, y_pred_final)))
print("\nConfusion report: \n" + str(classification_report(y_test, y_pred_final)))

Accuracy score: 0.774647887329436
Balanced accuracy score: 0.7345693758737237

Confusion matrix:
[[ 21  4  0  0  0]
 [ 2 20  0  1  0]
 [ 2  2  0  0  0]
 [ 0  1  0  2  0]
 [ 0  0  0  3  0]
 [ 1  1  0  0  7]]

Classification report:
      precision    recall  f1-score   support

   1         0.78         0.81         0.79         26
   2         0.80         0.87         0.83         23
   3         0.33         0.29         0.31         7
   4         0.88         1.00         0.94         3
   5         0.75         1.00         0.86         9

 accuracy          0.78         0.73         0.77         71
 macro avg         0.78         0.74         0.77         71
 weighted avg         0.78         0.77         0.77         71
```