

In [3]:

```
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
sns.set(style="ticks")
```

Вариант 3. Дата сет - 1

In [12]:

```
data = pd.read_csv('data\data.csv', sep=",")
data.head()
```

Out[12]:

	Unnamed: 0	ID	Name	Age	Photo	Nationality
0	0	158023	L. Messi	31	https://cdn.sofifa.org/players/4/19/158023.png	Argentina
1	1	20801	Cristiano Ronaldo	33	https://cdn.sofifa.org/players/4/19/20801.png	Portugal
2	2	190871	Neymar Jr	26	https://cdn.sofifa.org/players/4/19/190871.png	Brazil
3	3	193080	De Gea	27	https://cdn.sofifa.org/players/4/19/193080.png	Spain
4	4	192985	K. De Bruyne	27	https://cdn.sofifa.org/players/4/19/192985.png	Belgium

5 rows × 89 columns

Масштабирование признака

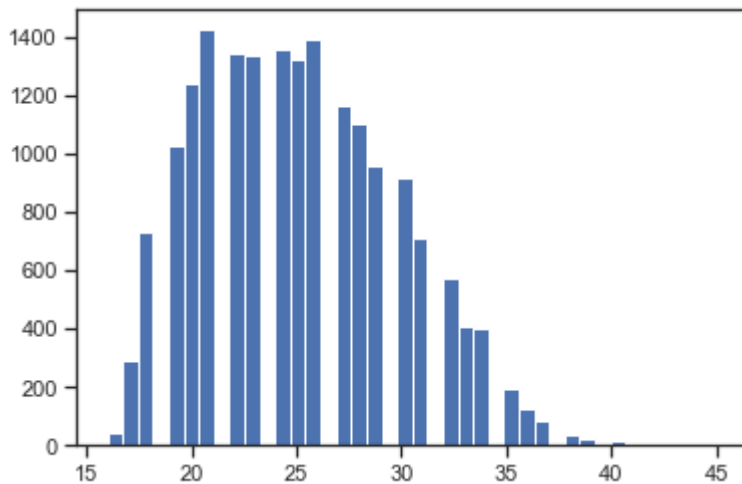
In [6]:

```
from sklearn.preprocessing import MinMaxScaler
```

In [22]:

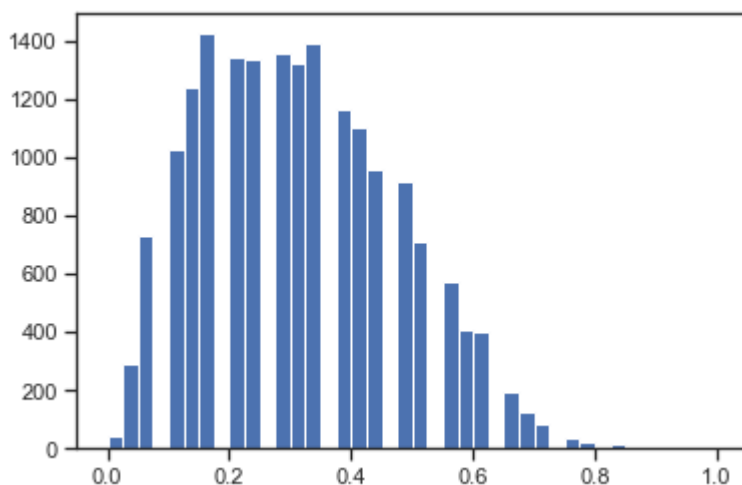
```
sc1 = MinMaxScaler()  
sc1_data = sc1.fit_transform(data[['Age']])  
plt.hist(data['Age'], 40)  
plt.show()
```

D:\Program\Anaconda3\lib\site-packages\sklearn\preprocessing\data.py:323:
DataConversionWarning: Data with input dtype int64 were all converted to float64 by MinMaxScaler.
return self.partial_fit(X, y)



In [23]:

```
plt.hist(sc1_data, 40)  
plt.show()
```



Преобразование категориальных признаков

one hot encoding

In [24]:

```
from sklearn.preprocessing import LabelEncoder, OneHotEncoder
```

In [41]:

```
data['Nationality'].unique()
```

Out[41]:

```
array(['Argentina', 'Portugal', 'Brazil', 'Spain', 'Belgium', 'Croatia',  
      'Uruguay', 'Slovenia', 'Poland', 'Germany', 'France', 'England',  
      'Italy', 'Egypt', 'Colombia', 'Denmark', 'Gabon', 'Wales',  
      'Senegal', 'Costa Rica', 'Slovakia', 'Netherlands',  
      'Bosnia Herzegovina', 'Morocco', 'Serbia', 'Algeria', 'Austria',  
      'Greece', 'Chile', 'Sweden', 'Korea Republic', 'Finland', 'Guinea',  
      'Montenegro', 'Armenia', 'Switzerland', 'Norway', 'Czech Republic',  
      'Scotland', 'Ghana', 'Central African Rep.', 'DR Congo',  
      'Ivory Coast', 'Russia', 'Ukraine', 'Iceland', 'Mexico', 'Jamaica',  
      'Albania', 'Venezuela', 'Japan', 'Turkey', 'Ecuador', 'Paraguay',  
      'Mali', 'Nigeria', 'Cameroon', 'Dominican Republic', 'Israel',  
      'Kenya', 'Hungary', 'Republic of Ireland', 'Romania',  
      'United States', 'Cape Verde', 'Australia', 'Peru', 'Togo',  
      'Syria', 'Zimbabwe', 'Angola', 'Burkina Faso', 'Iran', 'Estonia',  
      'Tunisia', 'Equatorial Guinea', 'New Zealand', 'FYR Macedonia',  
      'United Arab Emirates', 'China PR', 'Guinea Bissau', 'Bulgaria',  
      'Kosovo', 'South Africa', 'Madagascar', 'Georgia', 'Tanzania',  
      'Gambia', 'Cuba', 'Belarus', 'Uzbekistan', 'Benin', 'Congo',  
      'Mozambique', 'Honduras', 'Canada', 'Northern Ireland', 'Cyprus',  
      'Saudi Arabia', 'Curacao', 'Moldova', 'Bolivia',  
      'Trinidad & Tobago', 'Sierra Leone', 'Zambia', 'Chad',  
      'Philippines', 'Haiti', 'Comoros', 'Libya', 'Panama',  
      'São Tomé & Príncipe', 'Eritrea', 'Oman', 'Iraq', 'Burundi',  
      'Fiji', 'New Caledonia', 'Lithuania', 'Luxembourg', 'Korea DPR',  
      'Liechtenstein', 'St Kitts Nevis', 'Latvia', 'Suriname', 'Uganda',  
      'El Salvador', 'Bermuda', 'Kuwait', 'Antigua & Barbuda',  
      'Thailand', 'Mauritius', 'Guatemala', 'Liberia', 'Kazakhstan',  
      'Niger', 'Mauritania', 'Montserrat', 'Namibia', 'Azerbaijan',  
      'Guam', 'Faroe Islands', 'India', 'Nicaragua', 'Barbados',  
      'Lebanon', 'Palestine', 'Guyana', 'Sudan', 'St Lucia', 'Ethiopia',  
      'Puerto Rico', 'Grenada', 'Jordan', 'Rwanda', 'Qatar',  
      'Afghanistan', 'Hong Kong', 'Andorra', 'Malta', 'Belize',  
      'South Sudan', 'Indonesia', 'Botswana'], dtype=object)
```

In [49]:

```
ohe = OneHotEncoder()  
data_one = ohe.fit_transform(data[['Nationality']])  
data_one.todense()
```

Out[49]:

```
matrix([[0., 0., 0., ..., 0., 0., 0.],  
        [0., 0., 0., ..., 0., 0., 0.],  
        [0., 0., 0., ..., 0., 0., 0.],  
        ...,  
        [0., 0., 0., ..., 0., 0., 0.],  
        [0., 0., 0., ..., 0., 0., 0.],  
        [0., 0., 0., ..., 0., 0., 0.]])
```

In [64]:

```
data_one.todense()[0:2]
```

Out[64]:

```
matrix([[0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
         0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
         0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
         0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
         0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
         0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
         0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
         0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
         0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
         0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
         0., 0., 0., 0.],
        [0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
         0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
         0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
         0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
         0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
         0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
         0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
         0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
         0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
         0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
         0., 0., 0., 0.])
```

label encoding

In [58]:

```
data['Nationality'].unique()
```

Out[58]:

```
array(['Argentina', 'Portugal', 'Brazil', 'Spain', 'Belgium', 'Croatia',  
      'Uruguay', 'Slovenia', 'Poland', 'Germany', 'France', 'England',  
      'Italy', 'Egypt', 'Colombia', 'Denmark', 'Gabon', 'Wales',  
      'Senegal', 'Costa Rica', 'Slovakia', 'Netherlands',  
      'Bosnia Herzegovina', 'Morocco', 'Serbia', 'Algeria', 'Austria',  
      'Greece', 'Chile', 'Sweden', 'Korea Republic', 'Finland', 'Guinea',  
      'Montenegro', 'Armenia', 'Switzerland', 'Norway', 'Czech Republic',  
      'Scotland', 'Ghana', 'Central African Rep.', 'DR Congo',  
      'Ivory Coast', 'Russia', 'Ukraine', 'Iceland', 'Mexico', 'Jamaica',  
      'Albania', 'Venezuela', 'Japan', 'Turkey', 'Ecuador', 'Paraguay',  
      'Mali', 'Nigeria', 'Cameroon', 'Dominican Republic', 'Israel',  
      'Kenya', 'Hungary', 'Republic of Ireland', 'Romania',  
      'United States', 'Cape Verde', 'Australia', 'Peru', 'Togo',  
      'Syria', 'Zimbabwe', 'Angola', 'Burkina Faso', 'Iran', 'Estonia',  
      'Tunisia', 'Equatorial Guinea', 'New Zealand', 'FYR Macedonia',  
      'United Arab Emirates', 'China PR', 'Guinea Bissau', 'Bulgaria',  
      'Kosovo', 'South Africa', 'Madagascar', 'Georgia', 'Tanzania',  
      'Gambia', 'Cuba', 'Belarus', 'Uzbekistan', 'Benin', 'Congo',  
      'Mozambique', 'Honduras', 'Canada', 'Northern Ireland', 'Cyprus',  
      'Saudi Arabia', 'Curacao', 'Moldova', 'Bolivia',  
      'Trinidad & Tobago', 'Sierra Leone', 'Zambia', 'Chad',  
      'Philippines', 'Haiti', 'Comoros', 'Libya', 'Panama',  
      'São Tomé & Príncipe', 'Eritrea', 'Oman', 'Iraq', 'Burundi',  
      'Fiji', 'New Caledonia', 'Lithuania', 'Luxembourg', 'Korea DPR',  
      'Liechtenstein', 'St Kitts Nevis', 'Latvia', 'Suriname', 'Uganda',  
      'El Salvador', 'Bermuda', 'Kuwait', 'Antigua & Barbuda',  
      'Thailand', 'Mauritius', 'Guatemala', 'Liberia', 'Kazakhstan',  
      'Niger', 'Mauritania', 'Montserrat', 'Namibia', 'Azerbaijan',  
      'Guam', 'Faroe Islands', 'India', 'Nicaragua', 'Barbados',  
      'Lebanon', 'Palestine', 'Guyana', 'Sudan', 'St Lucia', 'Ethiopia',  
      'Puerto Rico', 'Grenada', 'Jordan', 'Rwanda', 'Qatar',  
      'Afghanistan', 'Hong Kong', 'Andorra', 'Malta', 'Belize',  
      'South Sudan', 'Indonesia', 'Botswana'], dtype=object)
```

In [59]:

```
le = LabelEncoder()  
data_new_1 = le.fit_transform(data['Nationality'])  
np.unique(data_new_1)
```

Out[59]:

```
array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12,  
       13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25,  
       26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38,  
       39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51,  
       52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64,  
       65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77,  
       78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90,  
       91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 102, 103,  
      104, 105, 106, 107, 108, 109, 110, 111, 112, 113, 114, 115, 116,  
      117, 118, 119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129,  
      130, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142,  
      143, 144, 145, 146, 147, 148, 149, 150, 151, 152, 153, 154, 155,  
      156, 157, 158, 159, 160, 161, 162, 163])
```

In [60]:

```
le.inverse_transform([0, 1, 2, 3])
```

Out[60]:

```
array(['Afghanistan', 'Albania', 'Algeria', 'Andorra'], dtype=object)
```

In []: