# Python大数据分析

## 六、数据的分组和聚合运算

| key | data | | Split | | | Apply | | Combine | |
|-----|------|---|-------|---|---|-------|---|---------|---|
| A | 0 | | A | 0 | | | | | |
| B | 5 | | A | 5 | | sum | | | |
| C | 10 | | A | 10 | | | | A | 15 |
| A | 5 | | B | 5 | | | | B | 30 |
| B | 10 | | B | 10 | | sum | | C | 45 |
| C | 15 | | B | 15 | | | | | |
| A | 10 | | C | 10 | | | | | |
| B | 15 | | C | 15 | | sum | | | |
| C | 20 | | C | 20 | | | | | |

| key | data |
|-----|------|
| A | 0 |
| B | 5 |
| C | 10 |
| A | 5 |
| B | 10 |
| C | 15 |
| A | 10 |
| B | 15 |
| C | 20 |

**Split**

| A | 0 |
|---|---|
| A | 5 |
| A | 10 |

| B | 5 |
|---|---|
| B | 10 |
| B | 15 |

| C | 10 |
|---|---|
| C | 15 |
| C | 20 |

**Apply**

sum

sum

sum

**Combine**

| A | 15 |
|---|----|
| B | 30 |
| C | 45 |

mp2 ⟩ Exec.py                                                                                      Exec ▾  ▶ 🐞 🔧 ■  🔍

```python
data = {'ID': ['000001', '000002', '000003', '000004', '000005', '000006', '000007'],
        'name': ['黎明', '赵怡春', '张富平', '白丽', '牛玉德', '姚华', '李南'],
        'gender': [True, False, True, False, True, False, True],
        'age': [16, 20, 18, 18, 17, 18, 16],
        'height': [1.88, 1.78, 1.81, 1.86, 1.74, 1.75, 1.76]
        }
frame = pd.DataFrame(data)
groups = frame.groupby(frame['gender'])
print(groups.count())
```
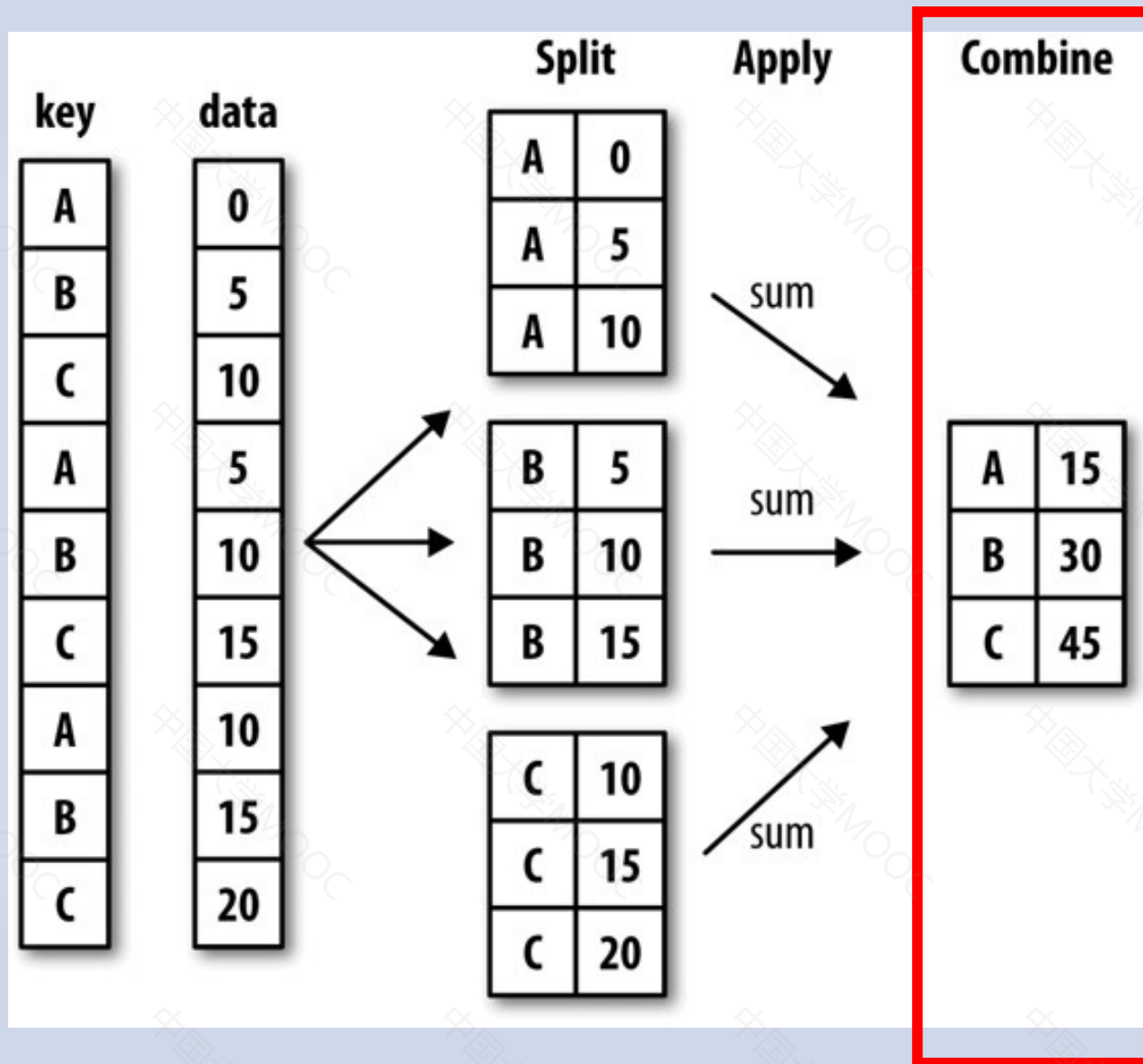
Run:  Exec

```
C:\Users\leeshuqing\PycharmProjects\mp2\venv\Scripts\python.exe C:/Users/leeshuqing/PycharmProjects/mp2/Exec.py
        ID   name   age   height
gender
False   3    3      3     3
True    4    4      4     4

Process finished with exit code 0
```

▶ 4: Run    🐞 5: Debug    ☰ 6: TODO    ⬛ Terminal    🐍 Python Console                                      ① Event Log

PyCharm 2019.3.4 available: // Update... (yesterday 19:05)                          13:40  CRLF  UTF-8  4 spaces  Python 3.7 (mp2)

File  Edit  View  Navigate  Code  Refactor  Run  Tools  VCS  Window  Help    mp2 [C:\Users\leeshuqing\PycharmProjects\mp2] - ...\Exec.py - PyCharm

mp2 ⟩ Exec.py                                                                  Exec ▾  ▶ 🐞 🔧 ■  🔍

```python
 5
 6     data = {'ID': ['000001', '000002', '000003', '000004', '000005', '000006', '000007'],
 7             'name': ['黎明', '赵怡春', '张富平', '白丽', '牛玉德', '姚华', '李南'],
 8             'gender': [True, False, True, False, True, False, True],
 9             'age': [16, 20, 18, 18, 17, 18, 16],
10             'height': [1.88, 1.78, 1.81, 1.86, 1.74, 1.75, 1.76]
11            }
12     frame = pd.DataFrame(data)
13     groups = frame.groupby('gender')
14     print(groups.count())
15
```

Run:  Exec

```
C:\Users\leeshuqing\PycharmProjects\mp2\venv\Scripts\python.exe C:/Users/leeshuqing/PycharmProjects/mp2/Exec.py
        ID   name   age   height
gender
False    3      3     3        3
True     4      4     4        4

Process finished with exit code 0
```
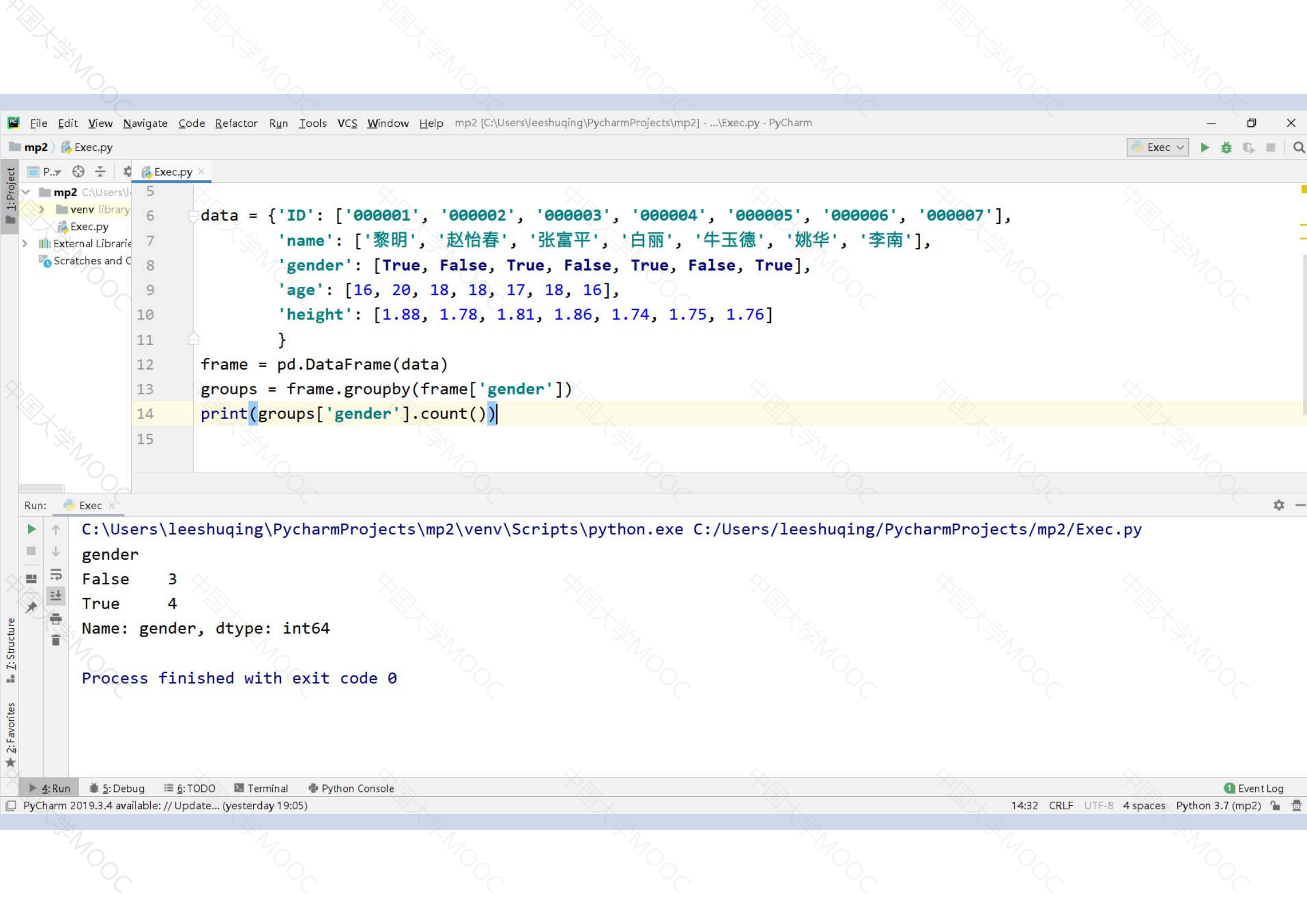
File Edit View Navigate Code Refactor Run Tools VCS Window Help    mp2 [C:\Users\leeshuqing\PycharmProjects\mp2] - ...\Exec.py - PyCharm    —  □  ✕

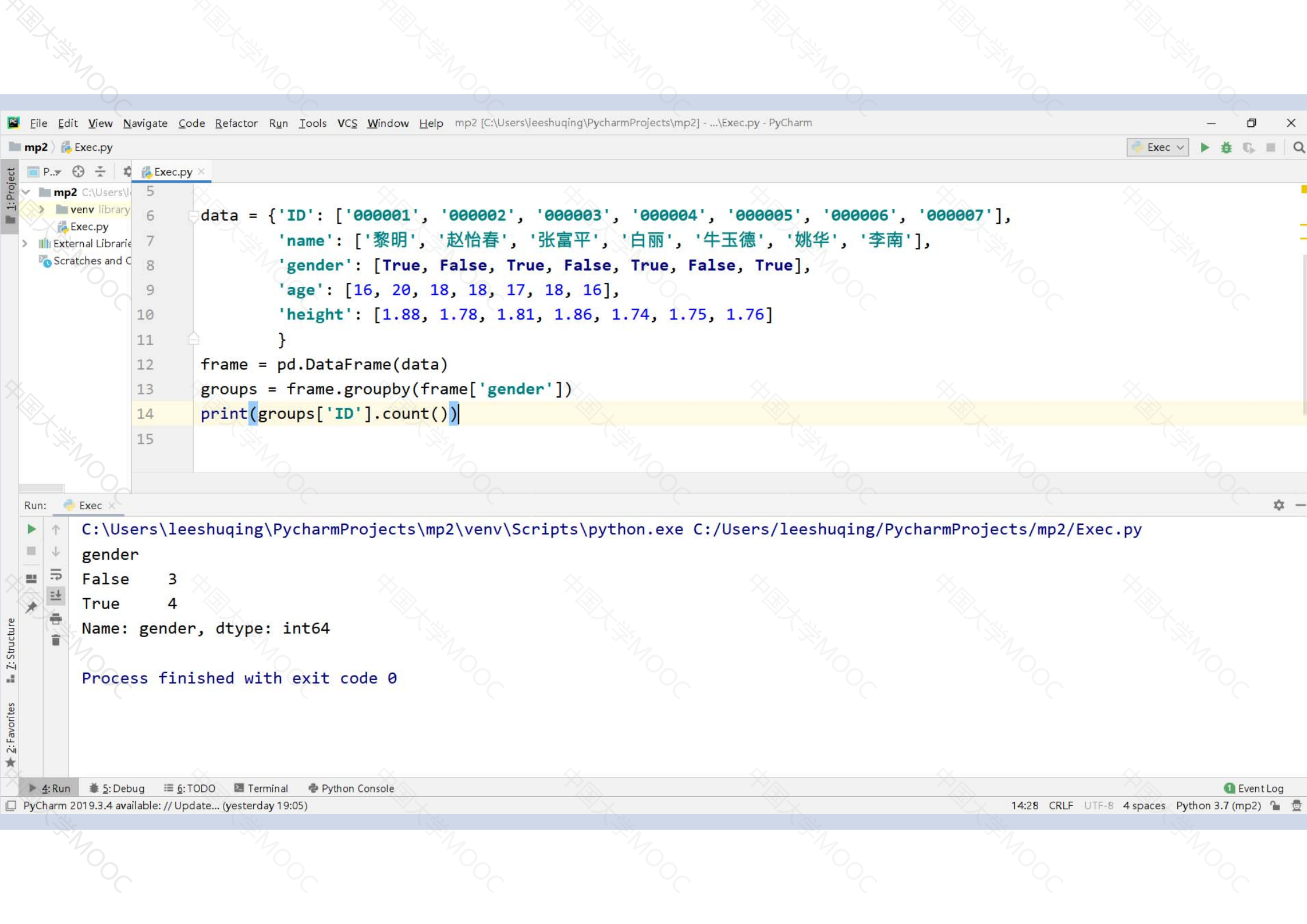mp2 ⟩ Exec.py                                                                                Exec ∨  ► 🐞 🐞 ■    🔍

```python
data = {'ID': ['000001', '000002', '000003', '000004', '000005', '000006', '000007'],
        'name': ['黎明', '赵怡春', '张富平', '白丽', '牛玉德', '姚华', '李南'],
        'gender': [True, False, True, False, True, False, True],
        'age': [16, 20, 18, 18, 17, 18, 16],
        'height': [1.88, 1.78, 1.81, 1.86, 1.74, 1.75, 1.76]
        }
frame = pd.DataFrame(data)
groups = frame.groupby(frame['gender'])
print(groups['gender'].count())
```

Run: Exec

```
C:\Users\leeshuqing\PycharmProjects\mp2\venv\Scripts\python.exe C:/Users/leeshuqing/PycharmProjects/mp2/Exec.py
gender
False    3
True     4
Name: gender, dtype: int64

Process finished with exit code 0
```

File  Edit  View  Navigate  Code  Refactor  Run  Tools  VCS  Window  Help    mp2 [C:\Users\leeshuqing\PycharmProjects\mp2] - ...\Exec.py - PyCharm    —  ☐  ✕

mp2 ⟩ 🐍 Exec.py                                                                        🐍 Exec ∨  ▶ 🐞 ⏭ ▣ Q

```python
data = {'ID': ['000001', '000002', '000003', '000004', '000005', '000006', '000007'],
        'name': ['黎明', '赵怡春', '张富平', '白丽', '牛玉德', '姚华', '李南'],
        'gender': [True, False, True, False, True, False, True],
        'age': [16, 20, 18, 18, 17, 18, 16],
        'height': [1.88, 1.78, 1.81, 1.86, 1.74, 1.75, 1.76]
        }
frame = pd.DataFrame(data)
groups = frame.groupby(frame['gender'])
print(groups['ID'].count())
```

```
Run:    🐍 Exec

C:\Users\leeshuqing\PycharmProjects\mp2\venv\Scripts\python.exe C:/Users/leeshuqing/PycharmProjects/mp2/Exec.py
gender
False    3
True     4
Name: gender, dtype: int64


Process finished with exit code 0
```
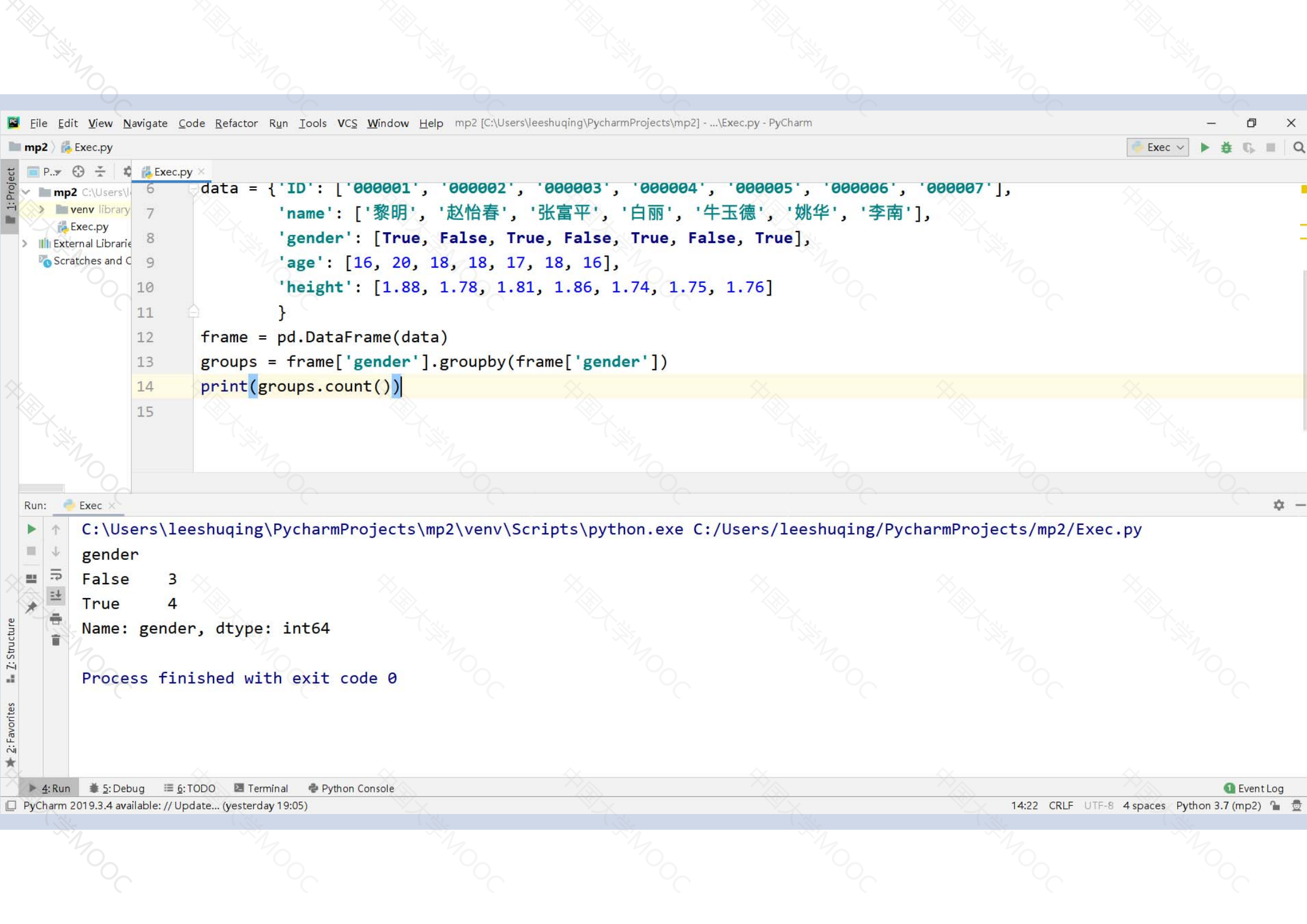
```python
data = {'ID': ['000001', '000002', '000003', '000004', '000005', '000006', '000007'],
        'name': ['黎明', '赵怡春', '张富平', '白丽', '牛玉德', '姚华', '李南'],
        'gender': [True, False, True, False, True, False, True],
        'age': [16, 20, 18, 18, 17, 18, 16],
        'height': [1.88, 1.78, 1.81, 1.86, 1.74, 1.75, 1.76]
        }
frame = pd.DataFrame(data)
groups = frame['gender'].groupby(frame['gender'])
print(groups.count())
```

```
C:\Users\leeshuqing\PycharmProjects\mp2\venv\Scripts\python.exe C:/Users/leeshuqing/PycharmProjects/mp2/Exec.py
gender
False    3
True     4
Name: gender, dtype: int64

Process finished with exit code 0
```
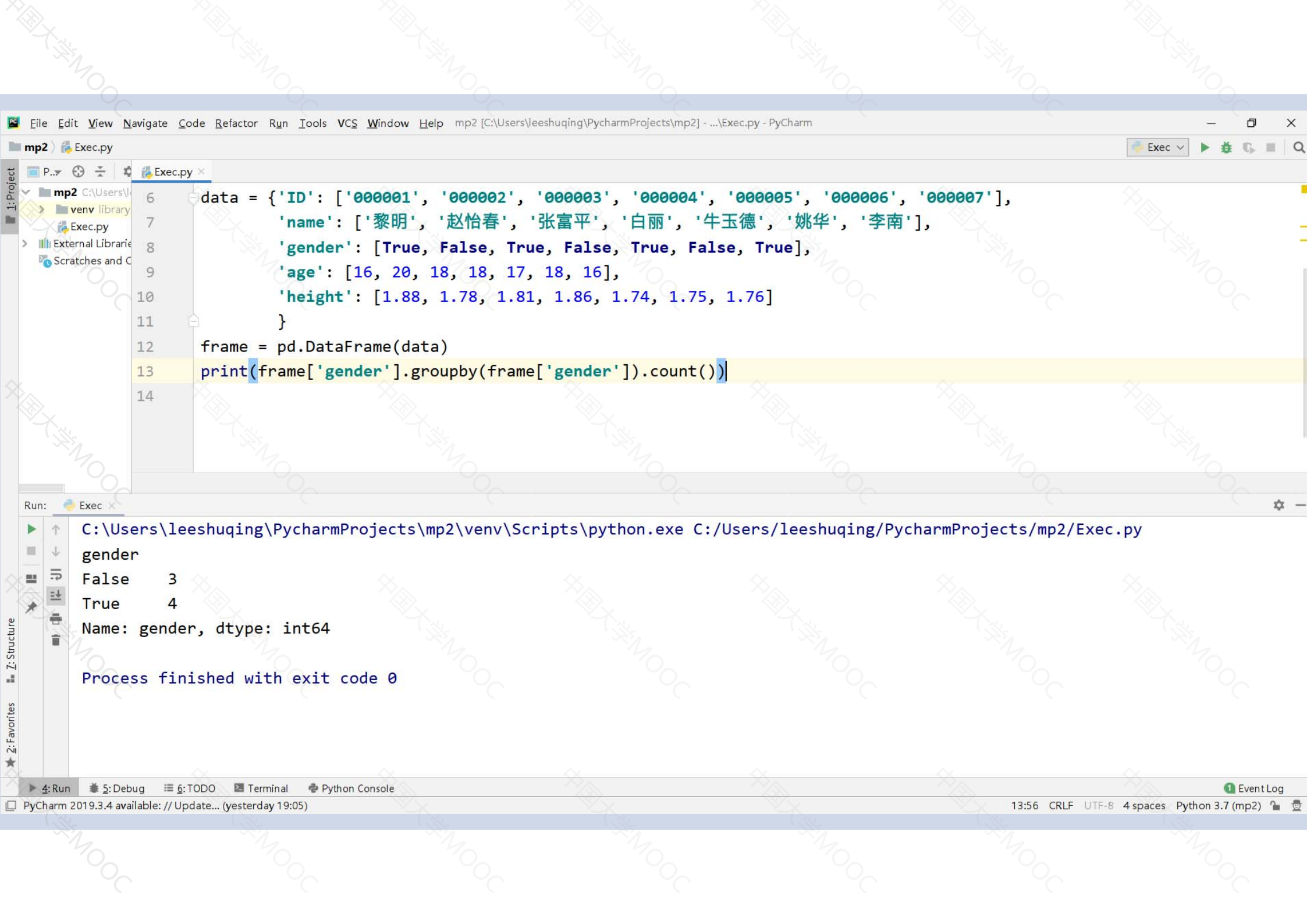
PC  File  Edit  View  Navigate  Code  Refactor  Run  Tools  VCS  Window  Help     mp2 [C:\Users\leeshuqing\PycharmProjects\mp2] - ...\Exec.py - PyCharm     —  ☐  ✕

mp2 ⟩ Exec.py                                                                    Exec ∨  ▶ ✿ ℚ ▣   ℚ

```
6    data = {'ID': ['000001', '000002', '000003', '000004', '000005', '000006', '000007'],
7            'name': ['黎明', '赵怡春', '张富平', '白丽', '牛玉德', '姚华', '李南'],
8            'gender': [True, False, True, False, True, False, True],
9            'age': [16, 20, 18, 18, 17, 18, 16],
10           'height': [1.88, 1.78, 1.81, 1.86, 1.74, 1.75, 1.76]
11           }
12   frame = pd.DataFrame(data)
13   print(frame['gender'].groupby(frame['gender']).count())
14
```

Run:  Exec ✕

```
C:\Users\leeshuqing\PycharmProjects\mp2\venv\Scripts\python.exe C:/Users/leeshuqing/PycharmProjects/mp2/Exec.py
gender
False    3
True     4
Name: gender, dtype: int64

Process finished with exit code 0
```
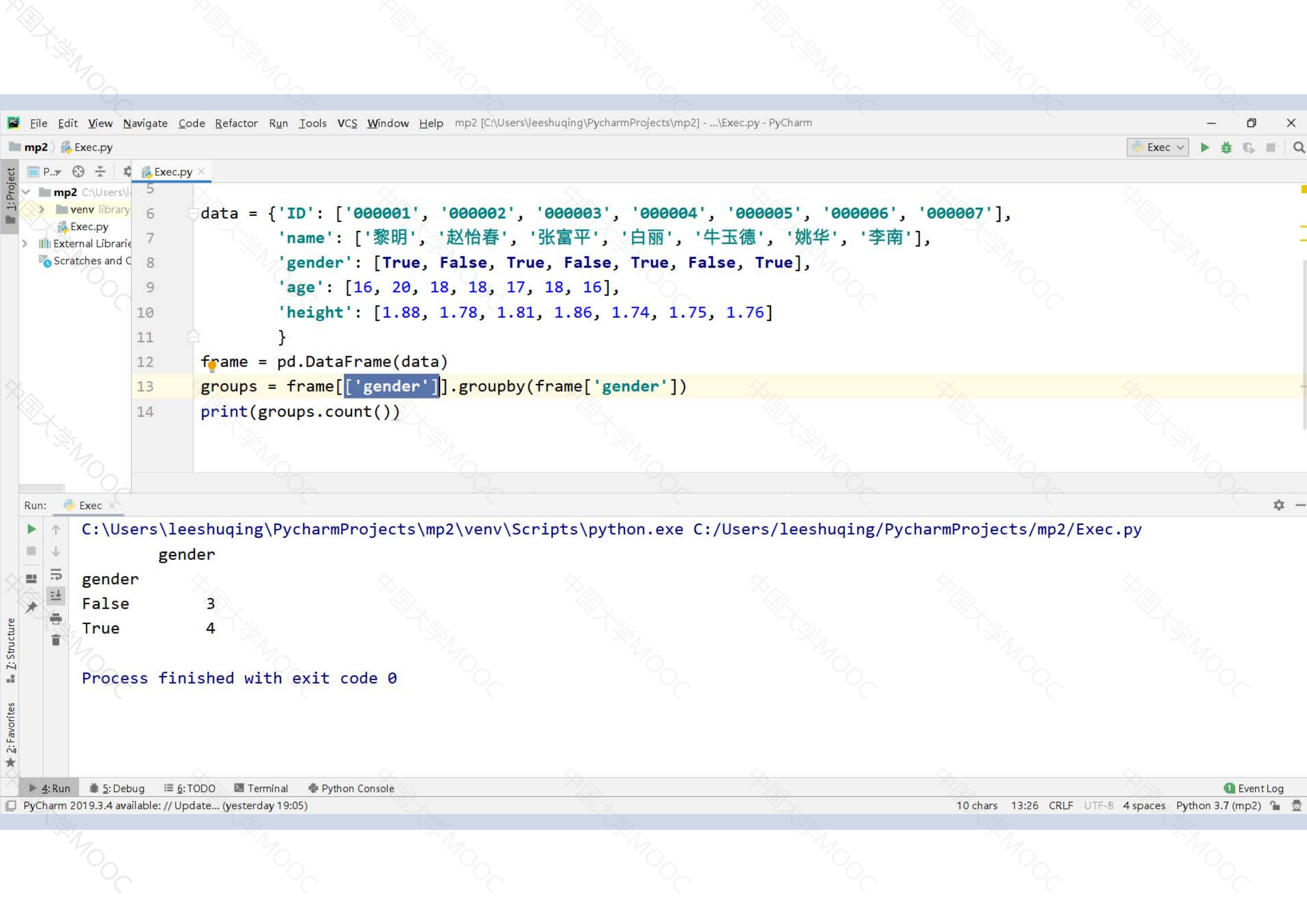
File  Edit  View  Navigate  Code  Refactor  Run  Tools  VCS  Window  Help    mp2 [C:\Users\leeshuqing\PycharmProjects\mp2] - ...\Exec.py - PyCharm    —  ☐  ✕

mp2 › Exec.py                                                                                              Exec ▾  ▶ ✿ ⚐ ■  Q

```python
data = {'ID': ['000001', '000002', '000003', '000004', '000005', '000006', '000007'],
        'name': ['黎明', '赵怡春', '张富平', '白丽', '牛玉德', '姚华', '李南'],
        'gender': [True, False, True, False, True, False, True],
        'age': [16, 20, 18, 18, 17, 18, 16],
        'height': [1.88, 1.78, 1.81, 1.86, 1.74, 1.75, 1.76]
        }
frame = pd.DataFrame(data)
groups = frame[['gender']].groupby(frame['gender'])
print(groups.count())
```

Run:  Exec

```
C:\Users\leeshuqing\PycharmProjects\mp2\venv\Scripts\python.exe C:/Users/leeshuqing/PycharmProjects/mp2/Exec.py
        gender
gender
False        3
True         4


Process finished with exit code 0
```

mp2 > Exec.py

Exec.py

```python
5
6    data = {'ID': ['000001', '000002', '000003', '000004', '000005', '000006', '000007'],
7            'name': ['黎明', '赵怡春', '张富平', '白丽', '牛玉德', '姚华', '李南'],
8            'gender': [True, False, True, False, True, False, True],
9            'age': [16, 20, 18, 18, 17, 18, 16],
10           'height': [1.88, 1.78, 1.81, 1.86, 1.74, 1.75, 1.76]
11           }
12   frame = pd.DataFrame(data)
13   groups = frame[['gender']].groupby(frame['gender'])
14   print(groups.count().rename(columns={'gender': 'genderCount'}))
15
```

Run:    Exec

```
C:\Users\leeshuqing\PycharmProjects\mp2\venv\Scripts\python.exe C:/Users/leeshuqing/PycharmProjects/mp2/Exec.py
        genderCount
gender
False            3
True             4


Process finished with exit code 0
```

File  Edit  View  Navigate  Code  Refactor  Run  Tools  VCS  Window  Help    mp2 [C:\Users\leeshuqing\PycharmProjects\mp2] - ...\Exec.py - PyCharm

mp2   Exec.py                                                      Exec

```python
data = {'ID': ['000001', '000002', '000003', '000004', '000005', '000006', '000007'],
        'name': ['黎明', '赵怡春', '张富平', '白丽', '牛玉德', '姚华', '李南'],
        'gender': [True, False, True, False, True, False, True],
        'age': [16, 20, 18, 18, 17, 18, 16],
        'height': [1.88, 1.78, 1.81, 1.86, 1.74, 1.75, 1.76]
        }
frame = pd.DataFrame(data)
print(frame['gender'].value_counts())
```

Run:  Exec

```
C:\Users\leeshuqing\PycharmProjects\mp2\venv\Scripts\python.exe C:/Users/leeshuqing/PycharmProjects/mp2/Exec.py
True     4
False    3
Name: gender, dtype: int64


Process finished with exit code 0
```
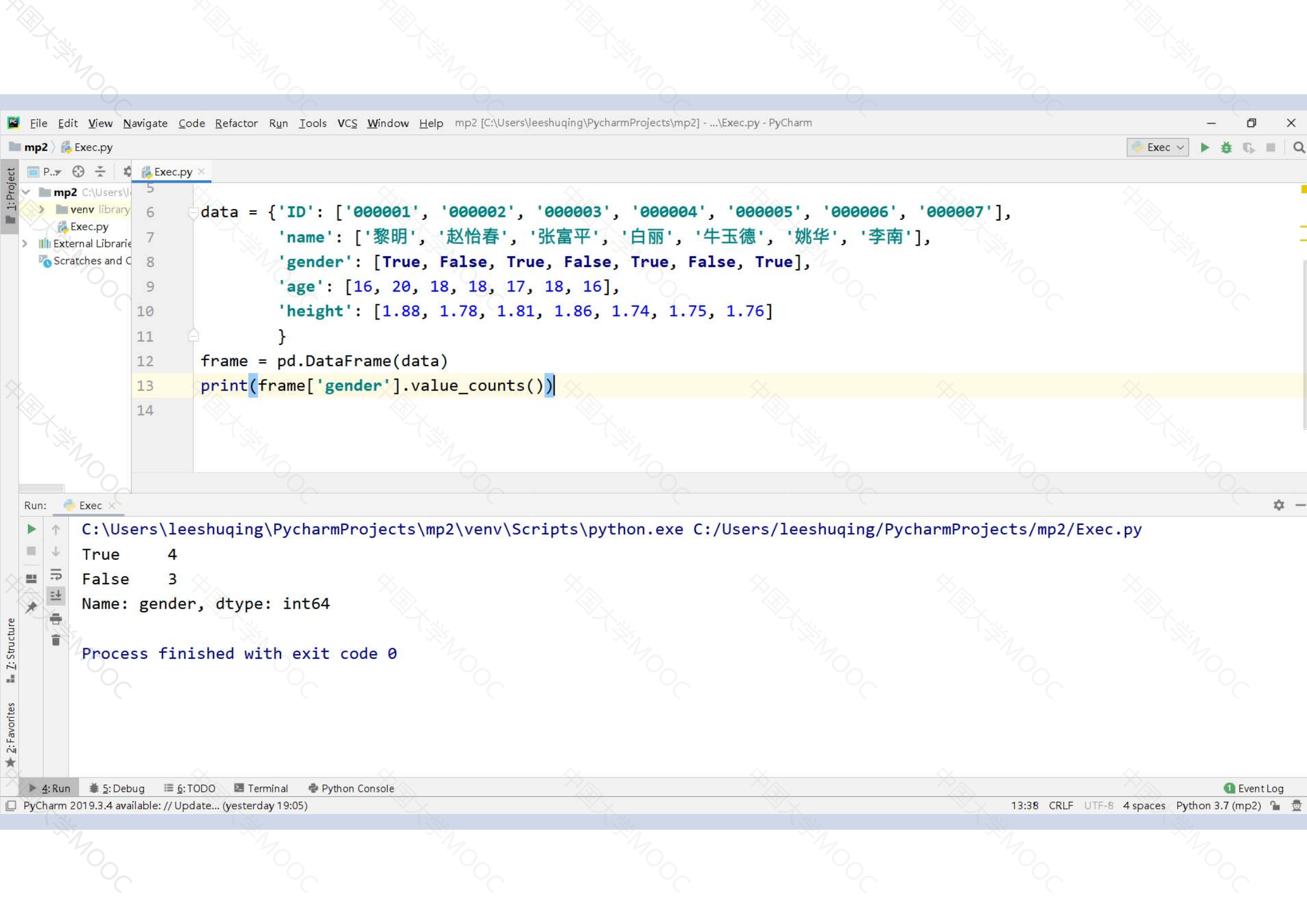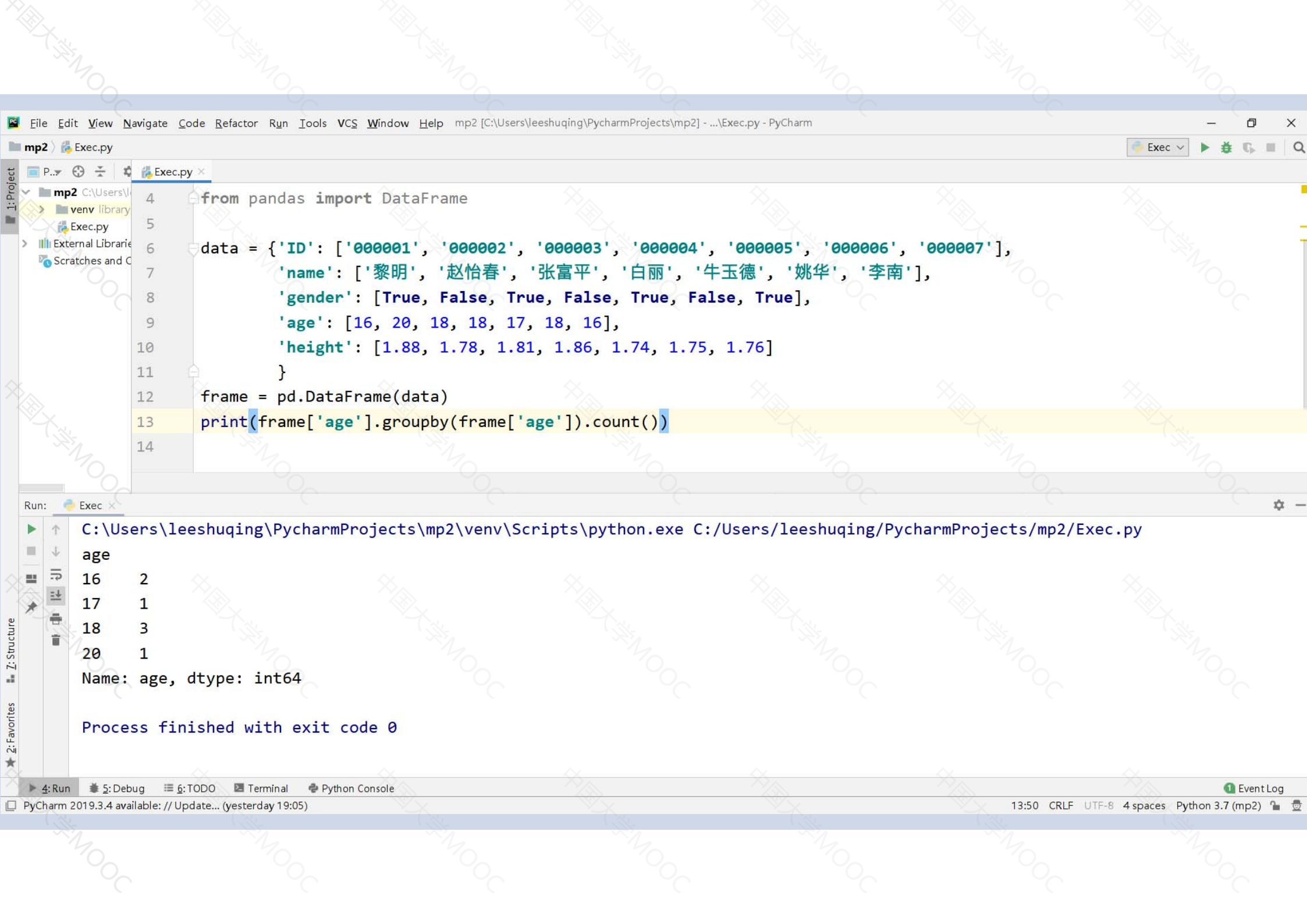
PC  File  Edit  View  Navigate  Code  Refactor  Run  Tools  VCS  Window  Help  mp2 [C:\Users\leeshuqing\PycharmProjects\mp2] - ...\Exec.py - PyCharm  —  □  ✕

mp2 ⟩ 🐍 Exec.py                                                                    🐍 Exec ▾  ▶  🐞  🔗  ■  🔍
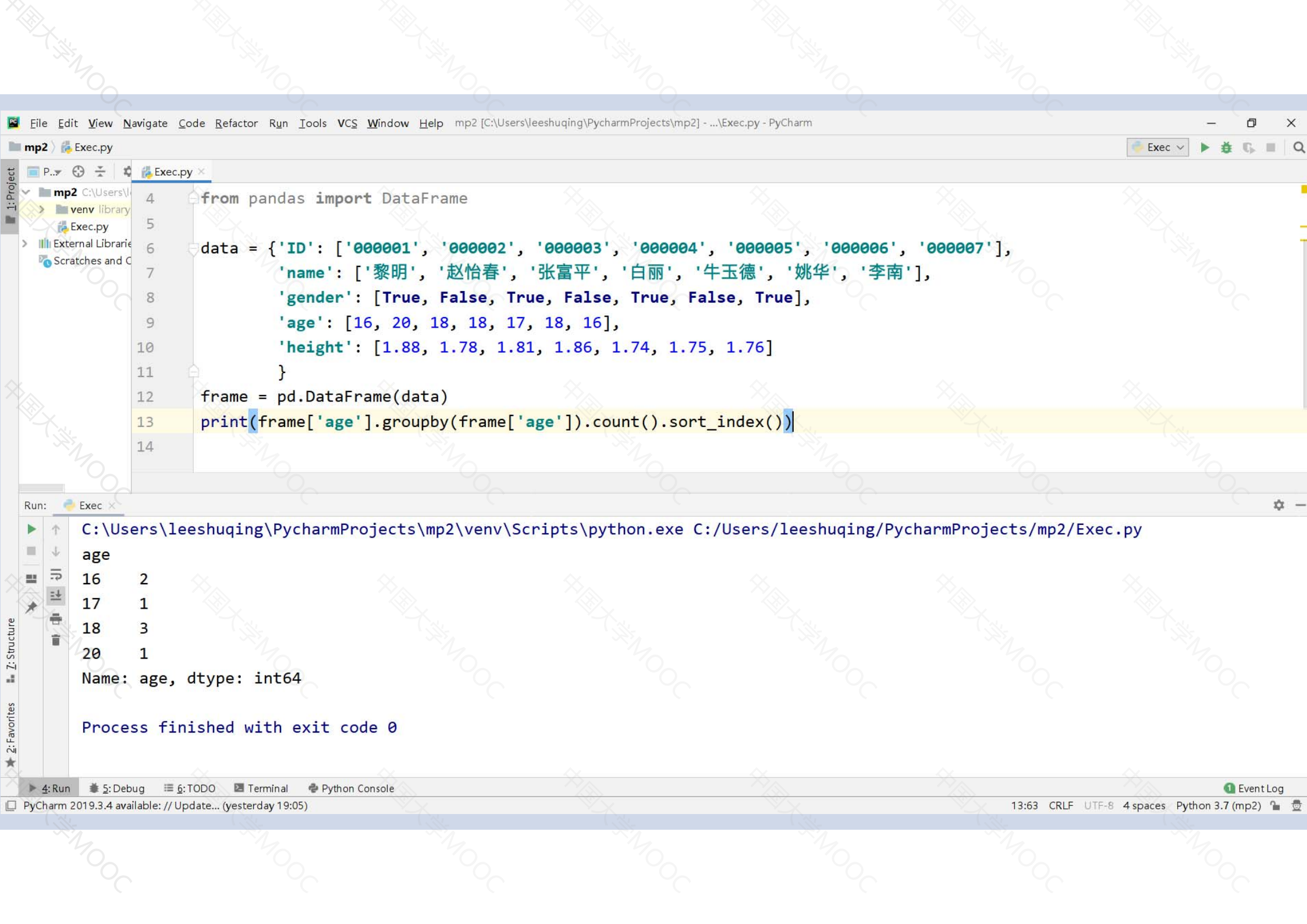
```python
from pandas import DataFrame

data = {'ID': ['000001', '000002', '000003', '000004', '000005', '000006', '000007'],
        'name': ['黎明', '赵怡春', '张富平', '白丽', '牛玉德', '姚华', '李南'],
        'gender': [True, False, True, False, True, False, True],
        'age': [16, 20, 18, 18, 17, 18, 16],
        'height': [1.88, 1.78, 1.81, 1.86, 1.74, 1.75, 1.76]
        }
frame = pd.DataFrame(data)
print(frame['age'].groupby(frame['age']).count())
```

Run:  🐍 Exec  ✕

```
C:\Users\leeshuqing\PycharmProjects\mp2\venv\Scripts\python.exe C:/Users/leeshuqing/PycharmProjects/mp2/Exec.py
age
16    2
17    1
18    3
20    1
Name: age, dtype: int64


Process finished with exit code 0
```

▶ 4: Run  🐞 5: Debug  ☰ 6: TODO  🖥 Terminal  🐍 Python Console                                ① Event Log
☐ PyCharm 2019.3.4 available: // Update... (yesterday 19:05)        13:50  CRLF  UTF-8  4 spaces  Python 3.7 (mp2)

File  Edit  View  Navigate  Code  Refactor  Run  Tools  VCS  Window  Help    mp2 [C:\Users\leeshuqing\PycharmProjects\mp2] - ...\Exec.py - PyCharm    —    ☐    ✕

mp2 ⟩ Exec.py                                                                                Exec ▾  ▶  🐞  🔄  ⬛  🔍

```python
from pandas import DataFrame

data = {'ID': ['000001', '000002', '000003', '000004', '000005', '000006', '000007'],
        'name': ['黎明', '赵怡春', '张富平', '白丽', '牛玉德', '姚华', '李南'],
        'gender': [True, False, True, False, True, False, True],
        'age': [16, 20, 18, 18, 17, 18, 16],
        'height': [1.88, 1.78, 1.81, 1.86, 1.74, 1.75, 1.76]
        }
frame = pd.DataFrame(data)
print(frame['age'].groupby(frame['age']).count().sort_index())
```

Run:    Exec

```
C:\Users\leeshuqing\PycharmProjects\mp2\venv\Scripts\python.exe C:/Users/leeshuqing/PycharmProjects/mp2/Exec.py
age
16     2
17     1
18     3
20     1
Name: age, dtype: int64


Process finished with exit code 0
```
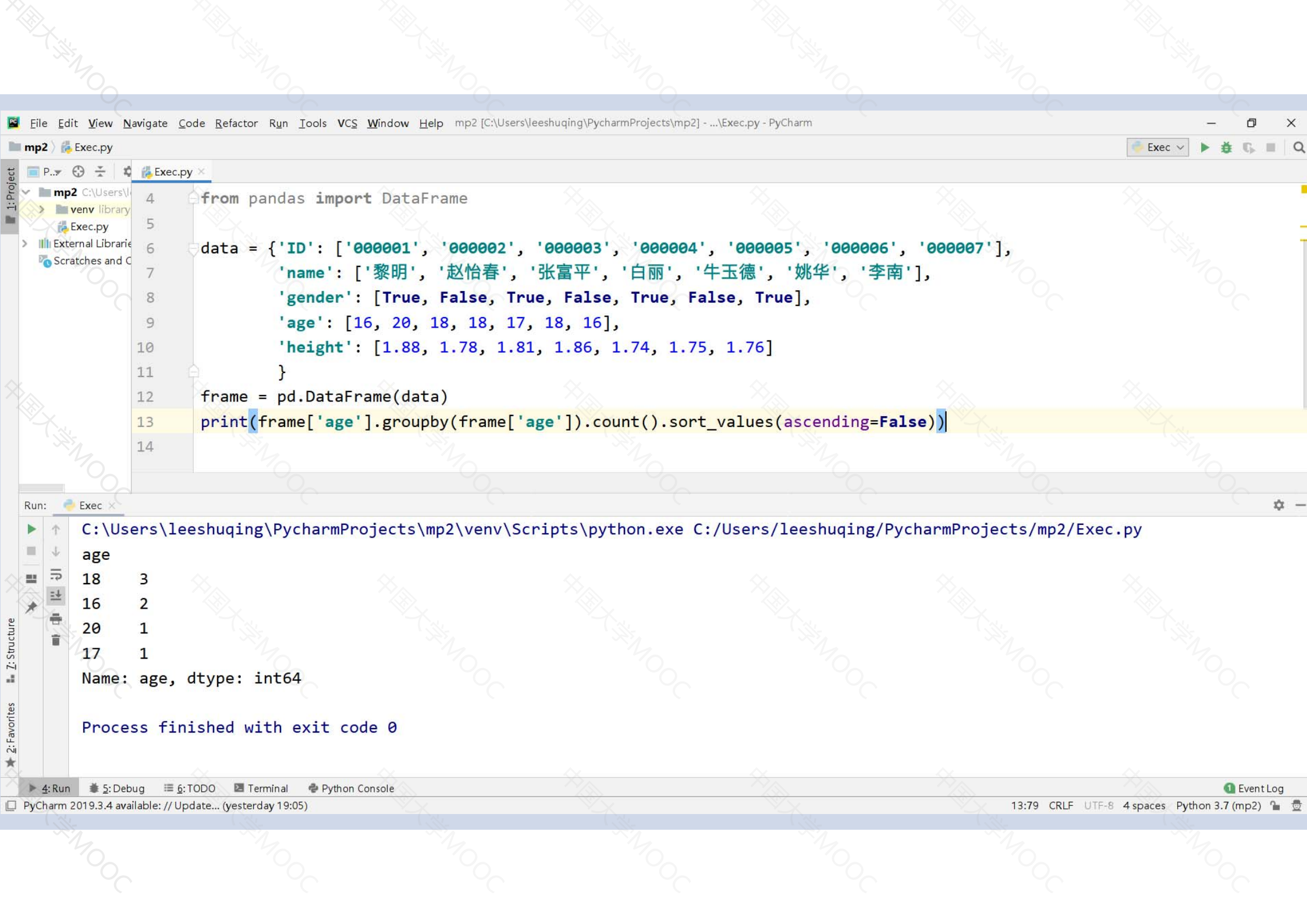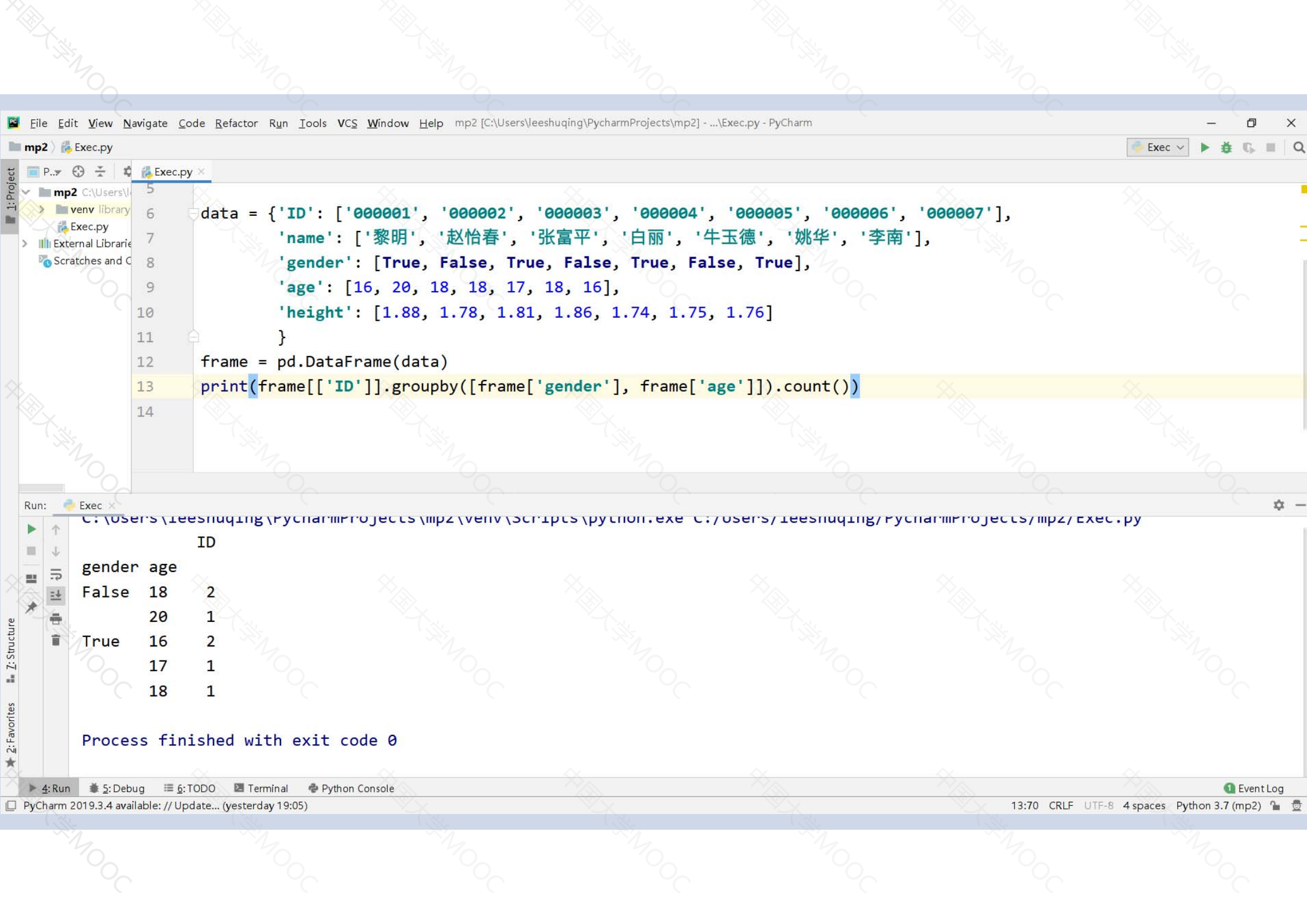
File  Edit  View  Navigate  Code  Refactor  Run  Tools  VCS  Window  Help    mp2 [C:\Users\leeshuqing\PycharmProjects\mp2] - ...\Exec.py - PyCharm    —  □  ✕

mp2 › Exec.py    Exec ∨  ► 🐞 🔂 ■ Q

```python
from pandas import DataFrame

data = {'ID': ['000001', '000002', '000003', '000004', '000005', '000006', '000007'],
        'name': ['黎明', '赵怡春', '张富平', '白丽', '牛玉德', '姚华', '李南'],
        'gender': [True, False, True, False, True, False, True],
        'age': [16, 20, 18, 18, 17, 18, 16],
        'height': [1.88, 1.78, 1.81, 1.86, 1.74, 1.75, 1.76]
        }
frame = pd.DataFrame(data)
print(frame['age'].groupby(frame['age']).count().sort_values(ascending=False))
```

Run:  Exec

```
C:\Users\leeshuqing\PycharmProjects\mp2\venv\Scripts\python.exe C:/Users/leeshuqing/PycharmProjects/mp2/Exec.py
age
18     3
16     2
20     1
17     1
Name: age, dtype: int64


Process finished with exit code 0
```

File  Edit  View  Navigate  Code  Refactor  Run  Tools  VCS  Window  Help    mp2 [C:\Users\leeshuqing\PycharmProjects\mp2] - ...\Exec.py - PyCharm  —  □  ✕

mp2  Exec.py                                                                                           Exec ▾  ▶  🐞  🔧  ■  🔍
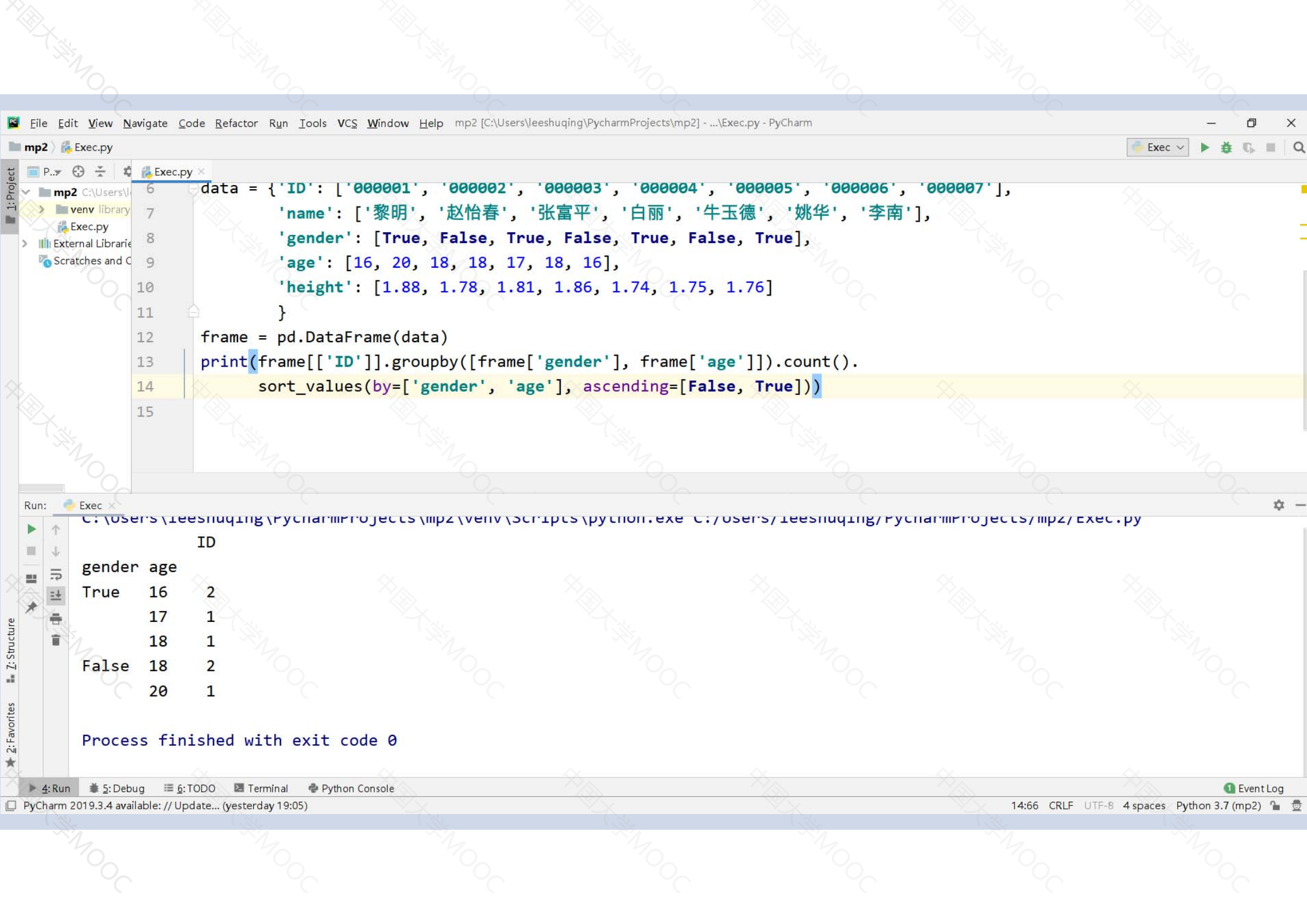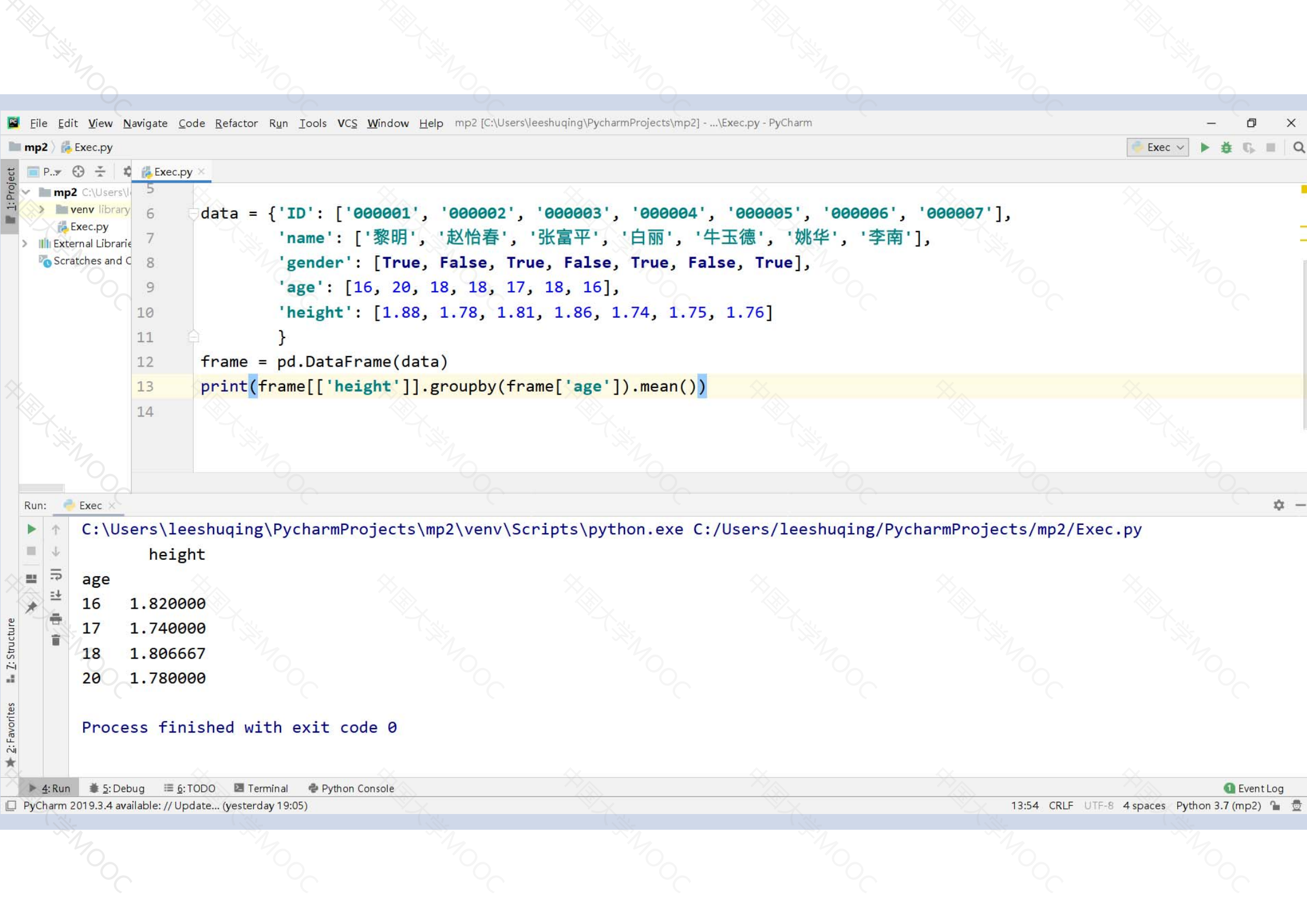
```python
data = {'ID': ['000001', '000002', '000003', '000004', '000005', '000006', '000007'],
        'name': ['黎明', '赵怡春', '张富平', '白丽', '牛玉德', '姚华', '李南'],
        'gender': [True, False, True, False, True, False, True],
        'age': [16, 20, 18, 18, 17, 18, 16],
        'height': [1.88, 1.78, 1.81, 1.86, 1.74, 1.75, 1.76]
        }
frame = pd.DataFrame(data)
print(frame[['ID']].groupby([frame['gender'], frame['age']]).count())
```

Run:    Exec

```
C:\Users\leeshuqing\PycharmProjects\mp2\venv\Scripts\python.exe C:/Users/leeshuqing/PycharmProjects/mp2/Exec.py
            ID
gender age
False  18    2
       20    1
True   16    2
       17    1
       18    1


Process finished with exit code 0
```

PC ☰  File  Edit  View  Navigate  Code  Refactor  Run  Tools  VCS  Window  Help    mp2 [C:\Users\leeshuqing\PycharmProjects\mp2] - ...\Exec.py - PyCharm    ─  ☐  ✕

mp2  ＞  Exec.py                                                                                      Exec ∨  ▶  🐞  ⟳  ▦  Q

```python
 6      data = {'ID': ['000001', '000002', '000003', '000004', '000005', '000006', '000007'],
 7             'name': ['黎明', '赵怡春', '张富平', '白丽', '牛玉德', '姚华', '李南'],
 8             'gender': [True, False, True, False, True, False, True],
 9             'age': [16, 20, 18, 18, 17, 18, 16],
10             'height': [1.88, 1.78, 1.81, 1.86, 1.74, 1.75, 1.76]
11             }
12      frame = pd.DataFrame(data)
13      print(frame[['ID']].groupby([frame['gender'], frame['age']]).count().
14          sort_values(by=['gender', 'age'], ascending=[False, True]))
15
```

Run:  🐍 Exec  ✕

```
C:\Users\leeshuqing\PycharmProjects\mp2\venv\Scripts\python.exe C:/Users/leeshuqing/PycharmProjects/mp2/Exec.py
            ID
gender age
True   16    2
       17    1
       18    1
False  18    2
       20    1


Process finished with exit code 0
```

mp2 > Exec.py                                                                              Exec
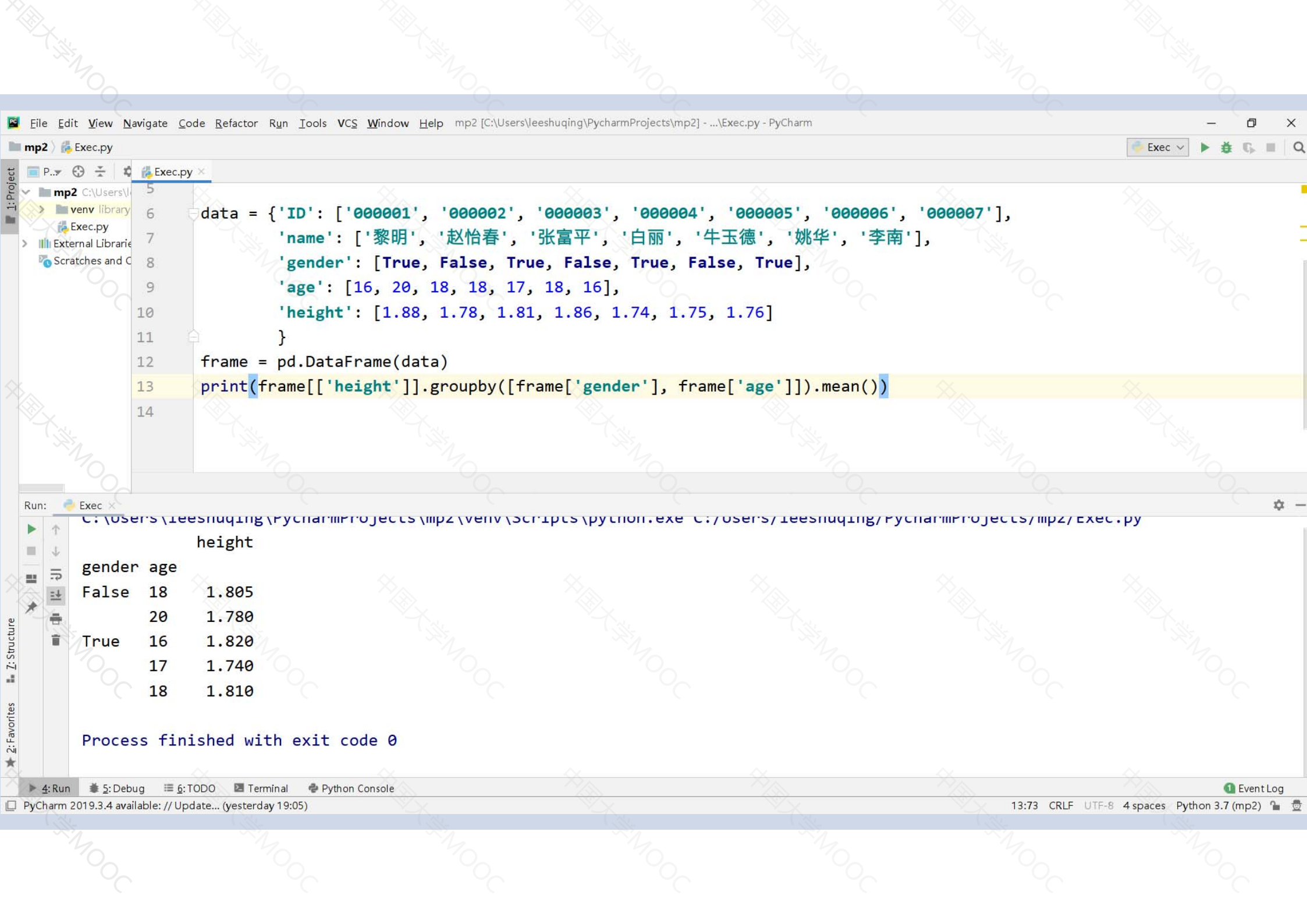
Exec.py ×

```python
5
6    data = {'ID': ['000001', '000002', '000003', '000004', '000005', '000006', '000007'],
7           'name': ['黎明', '赵怡春', '张富平', '白丽', '牛玉德', '姚华', '李南'],
8           'gender': [True, False, True, False, True, False, True],
9           'age': [16, 20, 18, 18, 17, 18, 16],
10          'height': [1.88, 1.78, 1.81, 1.86, 1.74, 1.75, 1.76]
11          }
12   frame = pd.DataFrame(data)
13   print(frame[['height']].groupby(frame['age']).mean())
14
```

Run:  Exec

```
C:\Users\leeshuqing\PycharmProjects\mp2\venv\Scripts\python.exe C:/Users/leeshuqing/PycharmProjects/mp2/Exec.py
        height
age
16    1.820000
17    1.740000
18    1.806667
20    1.780000


Process finished with exit code 0
```

4: Run      5: Debug      6: TODO      Terminal      Python Console                                                    Event Log

PyCharm 2019.3.4 available: // Update... (yesterday 19:05)                          13:54   CRLF   UTF-8   4 spaces   Python 3.7 (mp2)

File  Edit  View  Navigate  Code  Refactor  Run  Tools  VCS  Window  Help    mp2 [C:\Users\leeshuqing\PycharmProjects\mp2] - ...\Exec.py - PyCharm

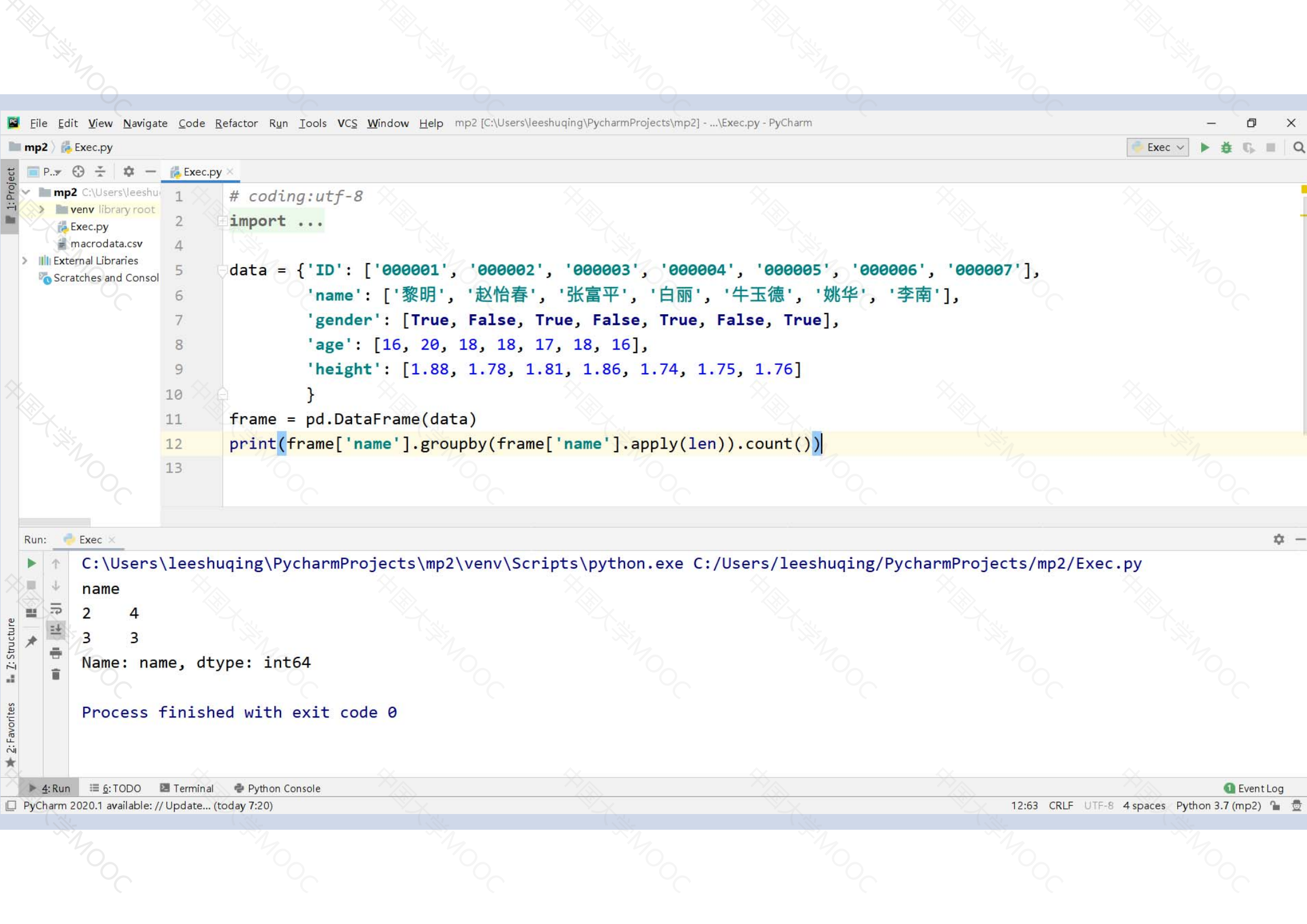mp2  Exec.py                                                                                          Exec

Exec.py

```python
 5
 6    data = {'ID': ['000001', '000002', '000003', '000004', '000005', '000006', '000007'],
 7           'name': ['黎明', '赵怡春', '张富平', '白丽', '牛玉德', '姚华', '李南'],
 8           'gender': [True, False, True, False, True, False, True],
 9           'age': [16, 20, 18, 18, 17, 18, 16],
10           'height': [1.88, 1.78, 1.81, 1.86, 1.74, 1.75, 1.76]
11           }
12    frame = pd.DataFrame(data)
13    print(frame[['height']].groupby([frame['gender'], frame['age']]).mean())
14
```

Run:    Exec

```
C:\Users\leeshuqing\PycharmProjects\mp2\venv\Scripts\python.exe C:/Users/leeshuqing/PycharmProjects/mp2/Exec.py
            height
gender age
False  18   1.805
       20   1.780
True   16   1.820
       17   1.740
       18   1.810


Process finished with exit code 0
```

4: Run     5: Debug     6: TODO     Terminal     Python Console                                    Event Log

PyCharm 2019.3.4 available: // Update... (yesterday 19:05)                    13:73  CRLF  UTF-8  4 spaces  Python 3.7 (mp2)

■ mp2 〉 🐍 Exec.py                                                                                    🐍 Exec ∨  ▶  🐞  🔾  ▣  Q

```
# coding:utf-8

import ...


data = {'ID': ['000001', '000002', '000003', '000004', '000005', '000006', '000007'],
        'name': ['黎明', '赵怡春', '张富平', '白丽', '牛玉德', '姚华', '李南'],
        'gender': [True, False, True, False, True, False, True],
        'age': [16, 20, 18, 18, 17, 18, 16],
        'height': [1.88, 1.78, 1.81, 1.86, 1.74, 1.75, 1.76]
        }
frame = pd.DataFrame(data)
print(frame['name'].groupby(frame['name'].apply(len)).count())
```

Run:    🐍 Exec

```
C:\Users\leeshuqing\PycharmProjects\mp2\venv\Scripts\python.exe C:/Users/leeshuqing/PycharmProjects/mp2/Exec.py
name
2    4
3    3
Name: name, dtype: int64


Process finished with exit code 0
```
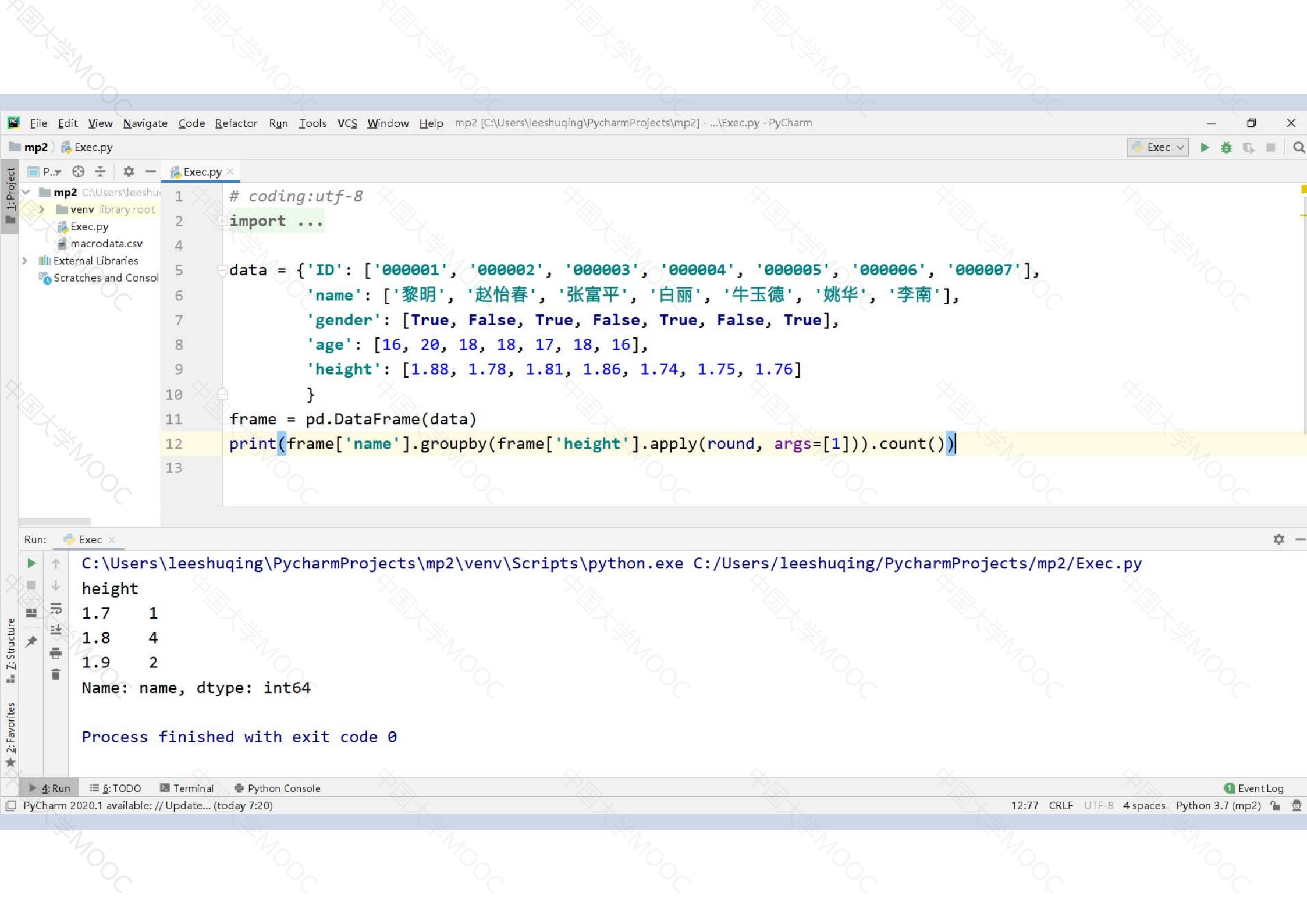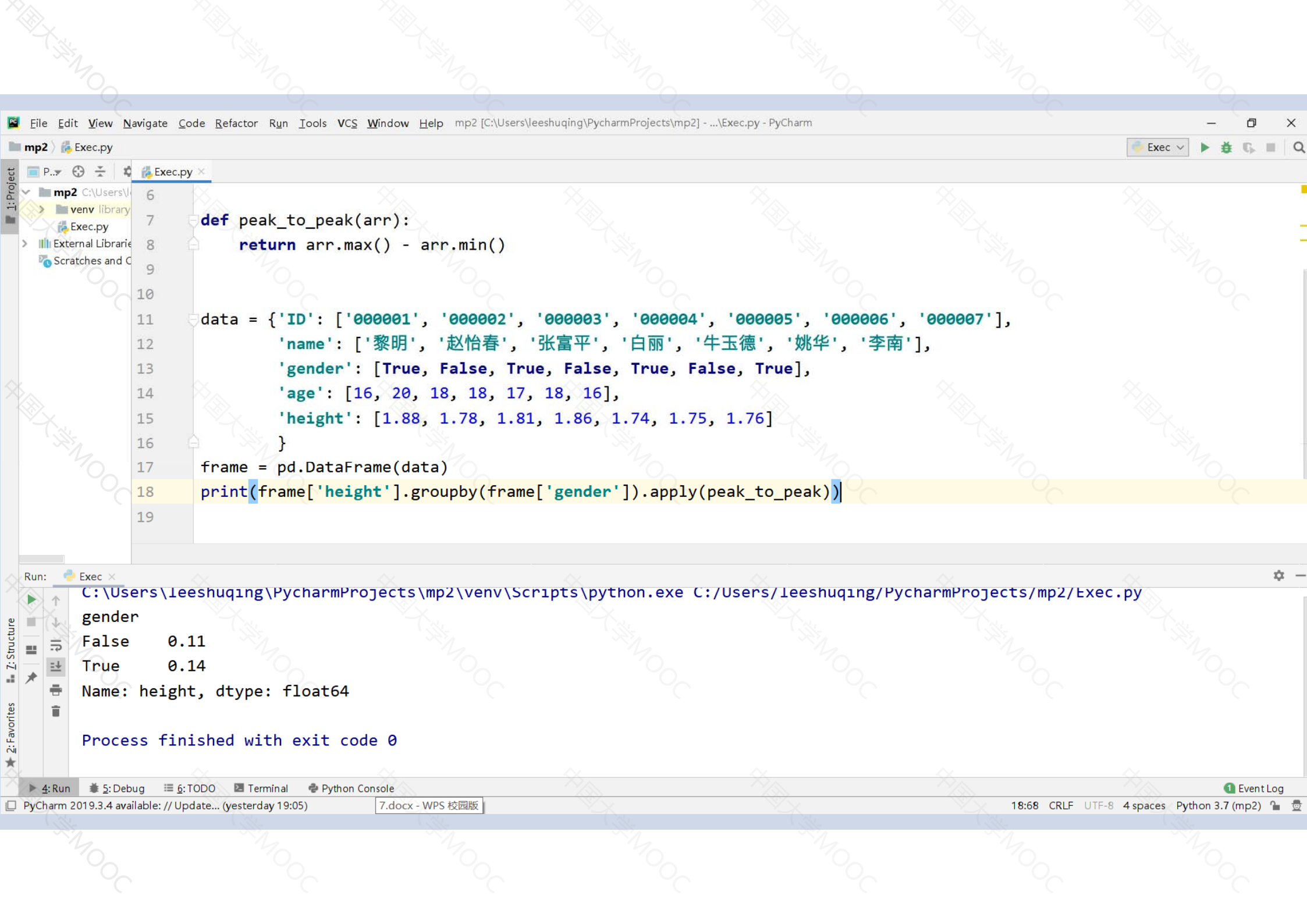
mp2 › Exec.py

Exec.py ×

```python
# coding:utf-8

import ...


data = {'ID': ['000001', '000002', '000003', '000004', '000005', '000006', '000007'],
        'name': ['黎明', '赵怡春', '张富平', '白丽', '牛玉德', '姚华', '李南'],
        'gender': [True, False, True, False, True, False, True],
        'age': [16, 20, 18, 18, 17, 18, 16],
        'height': [1.88, 1.78, 1.81, 1.86, 1.74, 1.75, 1.76]
        }
frame = pd.DataFrame(data)
print(frame['name'].groupby(frame['height'].apply(round, args=[1])).count())
```

Run: Exec ×

```
C:\Users\leeshuqing\PycharmProjects\mp2\venv\Scripts\python.exe C:/Users/leeshuqing/PycharmProjects/mp2/Exec.py
height
1.7    1
1.8    4
1.9    2
Name: name, dtype: int64


Process finished with exit code 0
```

File Edit View Navigate Code Refactor Run Tools VCS Window Help    mp2 [C:\Users\leeshuqing\PycharmProjects\mp2] - ...\Exec.py - PyCharm

mp2 ⟩ Exec.py    Exec ∨ ▶ 🐞 🔲 Q

P.. 🔅 🔧    Exec.py ×

```python
def peak_to_peak(arr):
    return arr.max() - arr.min()


data = {'ID': ['000001', '000002', '000003', '000004', '000005', '000006', '000007'],
        'name': ['黎明', '赵怡春', '张富平', '白丽', '牛玉德', '姚华', '李南'],
        'gender': [True, False, True, False, True, False, True],
        'age': [16, 20, 18, 18, 17, 18, 16],
        'height': [1.88, 1.78, 1.81, 1.86, 1.74, 1.75, 1.76]
        }
frame = pd.DataFrame(data)
print(frame['height'].groupby(frame['gender']).apply(peak_to_peak))
```

```
C:\Users\leeshuqing\PycharmProjects\mp2\venv\Scripts\python.exe C:/Users/leeshuqing/PycharmProjects/mp2/Exec.py
gender
False    0.11
True     0.14
Name: height, dtype: float64


Process finished with exit code 0
```

File  Edit  View  Navigate  Code  Refactor  Run  Tools  VCS  Window  Help   mp2 [C:\Users\leeshuqing\PycharmProjects\mp2] - ...\Exec.py - PyCharm  — □ ✕

mp2 > 📄 Exec.py                                                          🚀 Exec ⌄ ▶ 🐞 🔄 ▣ 🔍

```
2      import ...

4
5      data = {'ID': ['000001', '000002', '000003', '000004', '000005', '000006', '000007'],
6              'name': ['黎明', '赵怡春', '张富平', '白丽', '牛玉德', '姚华', '李南'],
7              'gender': [True, False, True, False, True, False, True],
8              'age': [16, 20, 18, 18, 17, 18, 16],
9              'height': [1.88, 1.78, 1.81, 1.86, 1.74, 1.75, 1.76]
10             }
11     frame = pd.DataFrame(data)
12     print(frame['height'].groupby(frame['gender']).apply(lambda arr: arr.max() - arr.min()))
13
```

Run:  Exec

```
C:\Users\leeshuqing\PycharmProjects\mp2\venv\Scripts\python.exe C:/Users/leeshuqing/PycharmProjects/mp2/Exec.py
gender
False    0.11
True     0.14
Name: height, dtype: float64


Process finished with exit code 0
```
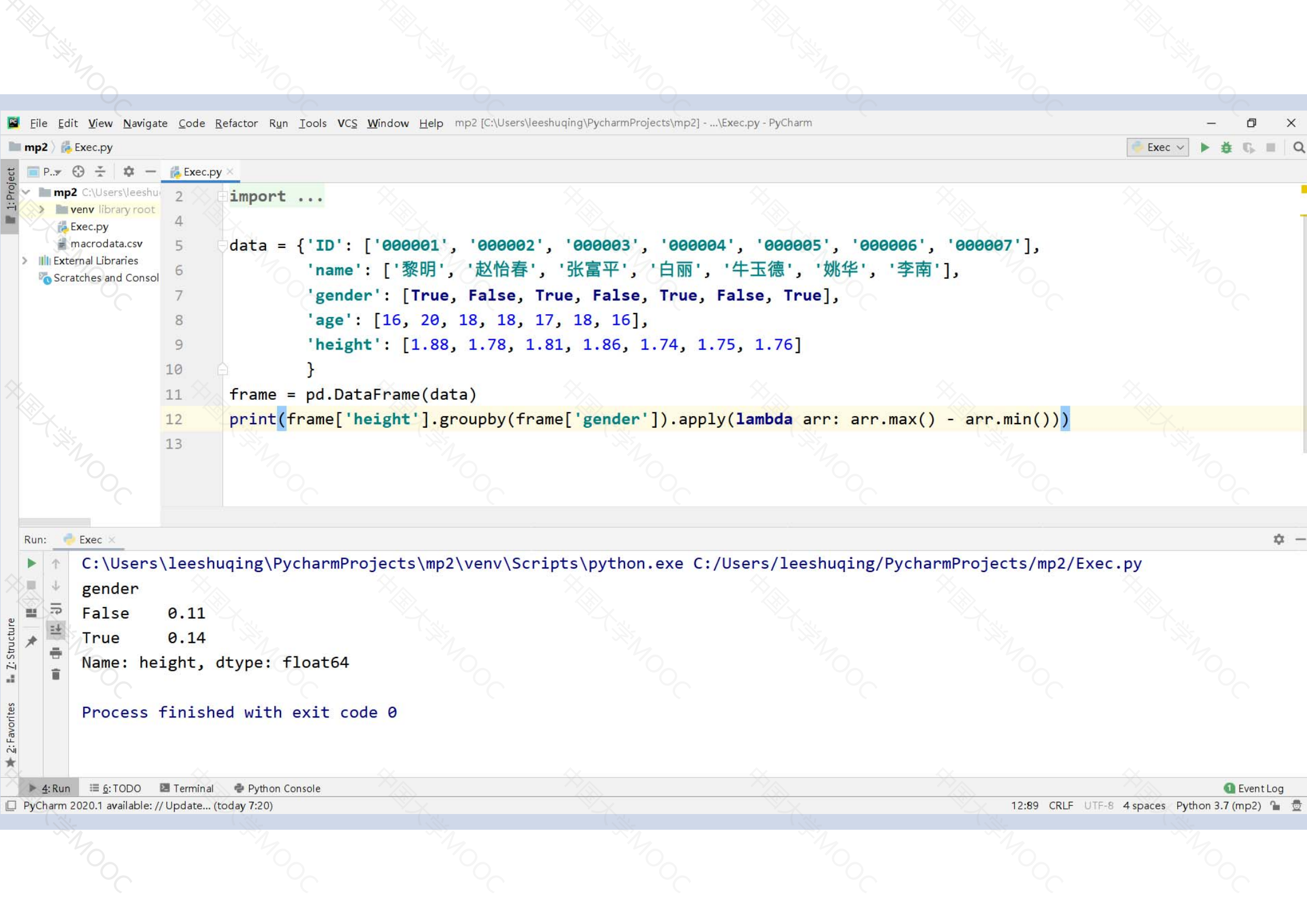
一次不学多，下次再学