


```
In [22]: cleaned
Out[22]:
```

	0	1	2
0	1.0	6.5	3.0

```
In [23]: data.dropna(how='all')
Out[23]:
```

	0	1	2
0	1.0	6.5	3.0
1	1.0	NaN	NaN
3	NaN	6.5	3.0

```
In [25]: data
Out[25]:
```

	0	1	2	4
0	1.0	6.5	3.0	NaN
1	1.0	NaN	NaN	NaN
2	NaN	NaN	NaN	NaN
3	NaN	6.5	3.0	NaN

```
In [26]: data.dropna(axis=1, how='all')
Out[26]:
```

	0	1	2
0	1.0	6.5	3.0
1	1.0	NaN	NaN
2	NaN	NaN	NaN
3	NaN	6.5	3.0

```
In [30]: df
Out[30]:
```

	0	1	2
0	-0.204708	NaN	NaN
1	-0.555730	NaN	NaN
2	0.092908	NaN	0.769023
3	1.246435	NaN	-1.296221
4	0.274992	0.228913	1.352917
5	0.886429	-2.001637	-0.371843
6	1.669025	-0.438570	-0.539741

```
In [32]: df.dropna(thresh=2)
Out[32]:
```

	0	1	2
2	0.092908	NaN	0.769023
3	1.246435	NaN	-1.296221
4	0.274992	0.228913	1.352917
5	0.886429	-2.001637	-0.371843
6	1.669025	-0.438570	-0.539741

你可能不想滤除缺失数据（有可能会丢弃跟它有关的其他数据），而是希望通过其他方式填补那些“空洞”。对于大多数情况而言，`fillna`方法是最主要的函数。通过一个常数调用`fillna`就会将缺失值替换为那个常数值：

若是通过一个字典调用fillna，就可以实现对不同的列填充不同的值：

fillna默认会返回新对象，但也可以对现有对象进行就地修改：

对reindexing有效的那些插值方法也可用于fillna：

? æ□°æ□®æ, □æ' □å□□å□□åx?.pdf[2020/7/14 18:20:03]

```
5 -1.265934    0.124121 -2.370232
```

In [42]: df.fillna(method='ffill', limit=2)
Out[42]:

	0	1	2
0	0.476985	3.248944	-1.021228
1	-0.577087	0.124121	0.302614
2	0.523772	0.124121	1.343810
3	-0.713544	0.124121	-2.370232
4	-1.860761	NaN	-2.370232
5	-1.265934	NaN	-2.370232

只要有些创新，你就可以利用fillna实现许多别的功能。比如说，你可以传入Series的平均值或中位数：

```
In [43]: data = pd.Series([1., NA, 3.5, NA, 7])
In [44]: data.fillna(data.mean())
Out[44]:
0    1.000000
1    3.833333
2    3.500000
3    3.833333
4    7.000000
dtype: float64
```

表7-2列出了fillna的参考。

参数	说明
value	用于填充缺失值的标量值或字典对象
method	插值方式。如果函数调用时未指定其他参数的话，默认为“ffill”

参数	说明
axis	待填充的轴，默认axis=0
inplace	修改调用者对象而不产生副本
limit	（对于前向和后向填充）可以连续填充的最大数量

fillna函数参数

7.2 数据转换

本章到目前为止介绍的都是数据的重排。另一类重要操作则是过滤、清理以及其他的转换工作。

移除重复数据

DataFrame中出现重复行有多种原因。下面就是一个例子：

```
In [45]: data = pd.DataFrame({'k1': ['one', 'two'] * 3 + ['two'],
.....:                       'k2': [1, 1, 2, 3, 3, 4, 4]})
```

```
In [46]: data
```

Out[46]:

	k1	k2
0	one	1
1	two	1
2	one	2
3	two	3
4	one	3
5	two	4
6	two	4

DataFrame的duplicated方法返回一个布尔型Series，表示各行是否是重复行（前面出现过的行）：

```
In [47]: data.duplicated()
```

Out[47]:

```
0    False
1    False
2    False
3    False
4    False
5    False
6     True
dtype: bool
```

还有一个与此相关的`drop_duplicates`方法，它会返回一个`DataFrame`，重复的数组会标为`False`：

```
In [48]: data.drop_duplicates()
```

Out[48]:

	k1	k2
0	one	1
1	two	1
2	one	2
3	two	3
4	one	3
5	two	4

这两个方法默认会判断全部列，你也可以指定部分列进行重复项判断。假设我们还有一列值，且只希望根据k1列过滤重复项：

```
In [49]: data['v1'] = range(7)
```

```
In [50]: data.drop_duplicates(['k1'])
```

```
Out[50]:
```

	k1	k2	v1
0	one	1	0
1	two	1	1

`duplicated`和`drop_duplicates`默认保留的是第一个出现的值组合。传入`keep='last'`则保留最后一个：

```
In [51]: data.drop_duplicates(['k1', 'k2'], keep='last')
```

Out[51]:

	k1	k2	v1
0	one	1	0
1	two	1	1
2	one	2	2
3	two	3	3
4	one	3	4
6	two	4	6

利用函数或映射进行数据转换


```
4      cow
5      pig
6      cow
7      pig
8  salmon
Name: food, dtype: object
```

使用map是一种实现元素级转换以及其他数据清理工作的便捷方式。

替换值

利用fillna方法填充缺失数据可以看做值替换的一种特殊情况。前面已经看到，map可用于修改对象的数据子集，而replace则提供了一种实现该功能的更简单、更灵活的方式。我们来看看下面这个Series：

```
In [60]: data = pd.Series([1., -999., 2., -999., -1000., 3.])
In [61]: data
Out[61]:
0      1.0
1    -999.0
2      2.0
3    -999.0
4   -1000.0
5      3.0
```

-999这个值可能是一个表示缺失数据的标记值。要将其替换为pandas能够理解的NA值，我们可以利用replace来产生一个新的Series（除非传入inplace=True）：

```
In [62]: data.replace(-999, np.nan)
Out[62]:
0          1.0
1          NaN
2          2.0
3          NaN
4    -1000.0
5          3.0
dtype: float64
```

如果你希望一次性替换多个值，可以传入一个由待替换值组成的列表以及一个替换值：

```
In [63]: data.replace([-999, -1000], np.nan)
Out[63]:
0      1.0
1      NaN
2      2.0
3      NaN
4      NaN
5      3.0
dtype: float64
```

要让每个值有不同的替换值，可以传递一个替换列表：

```
In [64]: data.replace([-999, -1000], [np.nan, 0])
Out[64]:
0      1.0
1      NaN
2      2.0
3      NaN
4      0.0
5      3.0
dtype: float64
```

传入的参数也可以是字典：

```
In [65]: data.replace({-999: np.nan, -1000: 0})
Out[65]:
0      1.0
1      NaN
```



```
2      2.0
3      NaN
4      0.0
5      3.0
dtype: float64
```

笔记：data.replace方法与data.str.replace不同，后者做的是字符串的元素级替换。我们会在后面学习Series的字符串方法。

重命名轴索引

跟Series中的值一样，轴标签也可以通过函数或映射进行转换，从而得到一个新的不同标签的对象。轴还可以被就地修改，而无需新建一个数据结构。接下来看看下面这个简单的例子：

[illegible]

跟Series一样，轴索引也有一个map方法：

```
In [67]: transform = lambda x: x[:4].upper()

In [68]: data.index.map(transform)
Out[68]: Index(['OHIO', 'COLO', 'NEW ', dtype='object')
```

你可以将其赋值给index，这样就可以对DataFrame进行就地修改：

```
In [69]: data.index = data.index.map(transform)
In [70]: data
Out[70]:
```

	one	two	three	four
OHIO	0	1	2	3
COLO	4	5	6	7
NEW	8	9	10	11

如果想要创建数据集的转换版（而不是修改原始数据），比较实用的方法是`rename`：

```
In [71]: data.rename(index=str.title, columns=str.upper)
Out[71]:
```

	ONE	TWO	THREE	FOUR
Ohio	0	1	2	3
Colo	4	5	6	7
New	8	9	10	11

特别说明一下，rename可以结合字典对象实现对部分轴标签的更新：

```
In [72]: data.rename(index={'OHIO': 'INDIANA'},
Out[72]:
      one two  peekaboo  four
INDIANA  0     1         2     3
COLO     4     5         6     7
NEW      8     9        10    11
```

rename可以实现复制DataFrame并对其索引和列标签进行赋值。如果希望就地修改某个数据集，传入inplace=True即可：

```
In [73]: data.rename(index={'OHIO': 'INDIANA'}, inplace=True)
In [74]: data
Out[74]:
```

	one	two	three	four
INDIANA	0	1	2	3
COLO	4	5	6	7
NEW	8	9	10	11

离散化和面元划分

为了便于分析，连续数据常常被离散化或拆分为“面元”（bin）。假设有一组人员数据，而你希望将它们划分为不同的年龄组：

```
In [75]: ages = [20, 22, 25, 27, 21, 23, 37, 31, 61, 45, 41, 32]
```

接下来将这些数据划分为“18到25”、“26到35”、“35到60”以及“60以上”几个面元。要实现该功能，你需要使用pandas的cut函数：

```
In [76]: bins = [18, 25, 35, 60, 100]
In [77]: cats = pd.cut(ages, bins)
In [78]: cats
Out[78]:
[(18, 25], (18, 25], (18, 25], (25, 35], (18, 25], ..., (25, 35], (60, 100], (35, 60], (35, 60], (25, 35]]
Length: 12
Categories (4, interval[int64]): [(18, 25] < (25, 35] < (35, 60] < (60, 100]]
```

pandas返回的是一个特殊的Categorical对象。结果展示了pandas.cut划分的面元。你可以将其看做一组表示面元名称的字符串。它的底层含有一个表示不同分类名称的类型数组，以及一个codes属性中的年龄数据的标签：

```
In [79]: cats.codes
Out[79]: array([0, 0, 0, 1, 0, 0, 2, 1, 3, 2, 2, 1], dtype=int8)

In [80]: cats.categories
Out[80]:
IntervalIndex([(18, 25], (25, 35], (35, 60], (60, 100]]
              closed='right',
              dtype='interval[int64]')

In [81]: pd.value_counts(cats)
Out[81]:
(18, 25]      5
(35, 60]      3
(25, 35]      3
(60, 100]     1
dtype: int64
```

pd.value_counts(cats)是pandas.cut结果的面元计数。

跟“区间”的数学符号一样，圆括号表示开端，而方括号则表示闭端（包括）。哪边是闭端可以通过`right=False`进行修改：

```
In [82]: pd.cut(ages, [18, 26, 36, 61, 100], right=False)
Out[82]:
[[18, 26), [18, 26), [18, 26), [26, 36), [18, 26), ..., [26, 36), [61, 100), [36,
 61), [36, 61), [26, 36)]
Length: 12
Categories (4, interval[int64]): [[18, 26) < [26, 36) < [36, 61) < [61, 100)]
```

你可以通过传递一个列表或数组到labels, 设置自己的面元名称:

```
In [83]: group_names = ['Youth', 'YoungAdult', 'MiddleAged', 'Senior']

In [84]: pd.cut(ages, bins, labels=group_names)
Out[84]:
[Youth, Youth, Youth, YoungAdult, Youth, ..., YoungAdult, Senior, MiddleAged, MiddleAged, YoungAdult]
Length: 12
Categories (4, object): [Youth < YoungAdult < MiddleAged < Senior]
```

如果向cut传入的是面元的数量而不是确切的面元边界，则它会根据数据的最小值和最大值计算等长面元。下面这个例子中，我们将一些均匀分布的数据分成四组：

```
In [85]: data = np.random.rand(20)

In [86]: pd.cut(data, 4, precision=2)
Out[86]:
[(0.34, 0.55], (0.34, 0.55], (0.76, 0.97], (0.76, 0.97], (0.34, 0.55], ..., (0.34
```

选项precision=2, 限定小数只有两位。

```
In [87]: data = np.random.randn(1000) # Normally distributed

In [88]: cats = pd.qcut(data, 4) # Cut into quartiles

In [89]: cats
Out[89]:
[(-0.0265, 0.62], (0.62, 3.928], (-0.68, -0.0265], (0.62, 3.928], (-0.0265, 0.62],
, ..., (-0.68, -0.0265], (-0.68, -0.0265], (-2.95, -0.68], (0.62, 3.928], (-0.68,
-0.0265]]
Length: 1000
Categories (4, interval[float64]): [(-2.95, -0.68] < (-0.68, -0.0265] < (-0.0265,
0.62] <
(0.62, 3.928]]

In [90]: pd.value_counts(cats)
Out[90]:
(0.62, 3.928]      250
(-0.0265, 0.62]    250
(-0.68, -0.0265]    250
(-2.95, -0.68]      250
dtype: int64
```

```
In [91]: pd.qcut(data, [0, 0.1, 0.5, 0.9, 1.])
Out[91]:
[(-0.0265, 1.286], (-0.0265, 1.286], (-1.187, -0.0265], (-0.0265, 1.286], (-0.0265, 1.286], ..., (-1.187, -0.0265], (-1.187, -0.0265], (-2.95, -1.187], (-0.0265, 1.286], (-1.187, -0.0265]]
Length: 1000
Categories (4, interval[float64]): [(-2.95, -1.187] < (-1.187, -0.0265] < (-0.0265, 1.286] < (1.286, 3.928]]
```

本章稍后在讲解聚合和分组运算时会再次用到cut和qcut，因为这两个离散化函数对分位和分组分析非常重要。

检测和过滤异常值

```
In [92]: data = pd.DataFrame(np.random.randn(1000, 4))
In [93]: data.describe()
Out[93]:
```

	0	1	2	3
count	1000.000000	1000.000000	1000.000000	1000.000000
mean	0.049091	0.026112	-0.002544	-0.051827
std	0.996947	1.007458	0.995232	0.998311
min	-3.645860	-3.184377	-3.745356	-3.428254
25%	-0.599807	-0.612162	-0.687373	-0.747478
50%	0.047101	-0.013609	-0.022158	-0.088274
75%	0.756646	0.695298	0.699046	0.623331
max	2.653656	3.525865	2.735527	3.366626

```
In [94]: col = data[2]
```

要选出全部含有“超过3或-3的值”的行，你可以在布尔型DataFrame中使用any方法：

根据这些条件，就可以对值进行设置。下面的代码可以将值限制在区间-3到3以内：

根据数据的值是正还是负，`np.sign(data)`可以生成1和-1：

排列和随机采样

然后就可以在基于iloc的索引操作或take函数中使用该数组了：

ç-?07ç«? æ□°æ□@æ, □æ' □å□□å□□å?.pdf[2020/7/14 18:20:03]

```
In [104]: df.take(sampler)
Out[104]:
```

	0	1	2	3
3	12	13	14	15
1	4	5	6	7
4	16	17	18	19
2	8	9	10	11
0	0	1	2	3

```
In [105]: df.sample(n=3)
Out[105]:
```

	0	1	2	3
3	12	13	14	15
4	16	17	18	19
2	8	9	10	11

```
In [106]: choices = pd.Series([5, 7, -1, 6, 4])
In [107]: draws = choices.sample(n=10, replace=True)
In [108]: draws
Out[108]:
4      4
1      7
4      4
2     -1
0      5
3      6
1      7
4      4
0      5
4      4
dtype: int64
```

另一种常用于统计建模或机器学习的转换方式是：将分类变量（categorical variable）转换为“哑变量”或“指标矩阵”。

```
In [109]: df = pd.DataFrame({'key': ['b', 'b', 'a', 'c', 'a', 'b'],
.....:                      'data': range(6)})

In [110]: pd.get_dummies(df['key'])
Out[110]:
```

	a	b	c
0	0	1	0
1	0	1	0
2	1	0	0
3	0	0	1
4	1	0	0
5	0	1	0

```
In [111]: dummies = pd.get_dummies(df['key'], prefix='key')
In [112]: df_with_dummy = df[['data1']].join(dummies)
```

```
In [113]: df_with_dummy
Out[113]:
```

	data1	key_a	key_b	key_c
0	0	0	1	0
1	1	0	1	0
2	2	1	0	0
3	3	0	0	1
4	4	1	0	0
5	5	0	1	0

如果DataFrame中的某行同属于多个分类，则事情就会有点复杂。看一下MovieLens 1M数据集，14章会更深入地研究它：

```
In [114]: mnames = ['movie_id', 'title', 'genres']

In [115]: movies = pd.read_table('datasets/movielens/movies.dat', sep='::',
.....:                             header=None, names=mnames)

In [116]: movies[:10]
Out[116]:
```

	movie_id	title	genres
0	1	Toy Story (1995)	Animation Children's Comedy
1	2	Jumanji (1995)	Adventure Children's Fantasy
2	3	Grumpier Old Men (1995)	Comedy Romance
3	4	Waiting to Exhale (1995)	Comedy Drama
4	5	Father of the Bride Part II (1995)	Comedy
5	6	Heat (1995)	Action Crime Thriller
6	7	Sabrina (1995)	Comedy Romance
7	8	Tom and Huck (1995)	Adventure Children's
8	9	Sudden Death (1995)	
9	10	GoldenEye (1995)	Action Adventure Thriller

要为每个genre添加指标变量就需要做一些数据规整操作。首先，我们从数据集中抽取出不同的genre值：

```
In [117]: all_genres = []
In [118]: for x in movies.genres:
.....:     all_genres.extend(x.split('|'))
In [119]: genres = pd.unique(all_genres)
```

现在有：

```
In [120]: genres
Out[120]:
array(['Animation', 'Children's', 'Comedy', 'Adventure', 'Fantasy',
      'Romance', 'Drama', 'Action', 'Crime', 'Thriller', 'Horror',
      'Sci-Fi', 'Documentary', 'War', 'Musical', 'Mystery', 'Film-Noir',
      'Western'], dtype=object)
```

构建指标DataFrame的方法之一是从一个全零DataFrame开始：

```
In [121]: zero_matrix = np.zeros((len(movies), len(genres)))
In [122]: dummies = pd.DataFrame(zero_matrix, columns=genres)
```

现在，迭代每一部电影，并将dummies各行的条目设为1。要这么做，我们使用dummies.columns来计算每个类型的列索引：

```
In [123]: gen = movies.genres[0]

In [124]: gen.split('|')
Out[124]: ['Animation', 'Children's', 'Comedy']

In [125]: dummies.columns.get_indexer(gen.split('|'))
Out[125]: array([0, 1, 2])
```

然后，根据索引，使用`.iloc`设定值：

```
In [126]: for i, gen in enumerate(movies.genres):
```


这些运算大部分都能使用正则表达式实现（马上就会看到）。

表7-3: Python内置的字符串方法

方法	说明
count	返回子串在字符串中的出现次数（非重叠）
endswith、startswith	如果字符串以某个后缀结尾（以某个前缀开头），则返回True
join	将字符串用作连接其他字符串序列的分隔符
index	如果在字符串中找到子串，则返回子串第一个字符所在的位置。如果没有找到，则引发ValueError。
find	如果在字符串中找到子串，则返回第一个发现的子串的第一个字符所在的位置。如果没有找到，则返回-1
rfind	如果在字符串中找到子串，则返回最后一个发现的子串的第一个字符所在的位置。如果没有找到，则返回-1
replace	用另一个字符串替换指定子串

如果打算对许多字符串应用同一条正则表达式，强烈建议通过`re.compile`创建`regex`对象。这样将可以节省大量的CPU时间。

```
text = """Dave dave@google.com
Steve steve@gmail.com
Rob rob@gmail.com
Ryan ryan@yahoo.com
"""
pattern = r'[A-Z0-9._%+]+@[A-Z0-9.-]+\.[A-Z]{2,4}'
# re.IGNORECASE makes the regex case-insensitive
regex = re.compile(pattern, flags=re.IGNORECASE)
```

```
In [155]: regex.findall(text)
Out[155]:
['dave@google.com',
 'steve@gmail.com',
 'rob@gmail.com',
 'ryan@yahoo.com']
```

```
In [156]: m = regex.search(text)

In [157]: m
Out[157]: <_sre.SRE_Match object; span=(5, 20), match='dave@google.com'>

In [158]: text[m.start():m.end()]
Out[158]: 'dave@google.com'
```

```
In [159]: print(regex.match(text))
None
```

```
In [160]: print(regex.sub('REDACTED', text))
Dave REDACTED
Steve REDACTED
Rob REDACTED
Ryan REDACTED
```

```
In [161]: pattern = r'([A-Z0-9._%+])@([A-Z0-9.-]+)\.([A-Z]{2,4})'
```

```
In [162]: regex = re.compile(pattern, flags=re.IGNORECASE)
```

```
In [163]: m = regex.match('wesm@bright.net')
In [164]: m.groups()
Out[164]: ('wesm', 'bright', 'net')
```

ç-?07ç«? æ□°æ□®æ, □æ' □å□□å□□å?.pdf[2020/7/14 18:20:03]


```
Out[177]:
Dave      NaN
Rob       NaN
Steve     NaN
Wes       NaN
dtype: float64
```

你可以利用这种方法对字符串进行截取：

```
In [178]: data.str[:5]
Out[178]:
Dave      dave@
Rob       rob@g
Steve     steve
Wes       NaN
dtype: object
```

表7-5介绍了更多的pandas字符串方法。

