

# 1 统计模型课程概述

## 1.1 statsmodels vs. sklearn?

## 1.2 准备python环境

In [ ]:

```
# 加载numpy库, 偶尔需要使用相应的函数
import numpy as np
```

In [ ]:

```
# 加载pandas库 (用于数据管理)
import pandas as pd
```

In [ ]:

```
# 加载matplotlib.pyplot库
from matplotlib import pyplot as plt

# 要求图形在notebook中直接显示
% matplotlib inline
```

In [ ]:

```
# 加载seaborn库
import seaborn as sns
# 加载seaborn默认格式设定
sns.set()
```

In [ ]:

```
# 解决中文显示问题
plt.rcParams["font.family"] = "STXIHEI"
```

In [ ]:

```
# 如果需要使用statsmodels, 可以考虑用api接口简化后续调用
import statsmodels.api as sm
```

## 1.3 statsmodels基本操作入门

**模型类的两种调用方式**

In [ ]:

```
# 绝对路径方式引用statsmodels.regression.linear_model.OLS()类
from statsmodels.regression.linear_model import OLS

OLS([1,2,3], [4,5,7])
```

In [ ]:

```
# 使用API方式引用statsmodels.regression.linear_model.OLS()类
import statsmodels.api as sm

sm.OLS([1,2,3], [4,5,7])
```

In [ ]:

```
olsobj = sm.OLS([1,2,3], [4,5,7])
type(olsobj)
```

### 模型类的各种属性

In [ ]:

```
olsobj.endog
```

### 模型拟合后的结果类

In [ ]:

```
olsres = olsobj.fit()
type(olsres)
```

In [ ]:

```
# 取出所需的结果
olsres.params
```

In [ ]:

```
olsres.summary()
```

### 使用模型类的表达式形式

In [ ]:

```
# 可以使用表达式形式的模型类列表
import statsmodels.formula.api as smf

dir(smf)
```

In [ ]:

```
# 模型类的同名小写名称(formula, data, subset=None, *args, **kwargs)
smf.ols('模型表达式', dfname)
```

## 1.4 sklearn的样本数据集

In [ ]:

```
from sklearn import datasets

boston = datasets.load_boston()
```

In [ ]:

```
# 显示数据对象的内容
print(boston.DESCR)
```

In [ ]:

```
# 转换为数据框便于使用
dfboston = pd.DataFrame(boston.data, columns = boston.feature_names)
dfboston.head()
```

## 1.5 sklearn基本操作入门

sklearn中集成了各种数据挖掘所需的变量变换、变量信息处理、统计建模、模型优化、模型评估方法，为便于使用，这些操作基本上都封装成了具有统一API的类，调用时都遵循统一的操作规范。

### 标准的类参数

class sklearn.大类名称.Modelclass(类参数列表)

Modelclass中基本通用的类参数：

```
fit_intercept = True : 模型是否包括常数项
    使用该选项就不需要在数据框中设定cons
n_jobs = 1 : 使用的例程数，为-1时使用全部CPU
max_iter = 200 : int, 模型最大迭代次数
tol = 0.0001 模型收敛标准
warm_start = False : 是否使用上一次的模型拟合结果作为本次初始值
sample_weight = None : 案例权重
random_state = None : int/RandomState instance/None, 随机器的设定
shuffle = True : 是否在拆分前对样本做随机排列
```

)# 大多数类参数都会有默认值

In [ ]:

```
from sklearn import preprocessing

# 完整的类名称为sklearn.preprocessing.StandardScaler()
std = preprocessing.StandardScaler()
std
```

In [ ]:

```
from sklearn import linear_model

# 完整的类名称为sklearn.linear_model.LinearRegression()
reg = linear_model.LinearRegression()
reg
```

### Modelclass中基本通用的类方法

get\_params([deep]) : 获取模型的具体参数设定  
set\_params(\*\*params) : 重新设定模型参数  
fit(X, y[, sample\_weight]) : 使用数据拟合模型/方法

特征处理class: Preprocessing、降维、Feature extraction/selection  
transform(X[, y]) : 使用拟合好的模型对指定数据进行转换  
fit\_transform(X[, y]) : 对数据拟合相应的方法，并且进行转换

建模分析class: Classification、Regression、Clustering  
predict(X) : 使用拟合好的模型对数据计算预测值  
predict\_proba(X) : 模型给出的每个案例（各个类别）的预测概率  
score(X, y[, sample\_weight]) : 返回模型决定系数/模型准确度评价指标

In [ ]:

```
std.get_params()
```

In [ ]:

```
# 使用fit方法，使std类获取数据中相应的信息
std.fit(boston.data)
```

In [ ]:

```
std.mean_
```

In [ ]:

```
std.var_
```

In [ ]:

```
ZX = std.transform(boston.data)
ZX[:2]
```

In [ ]:

```
std.fit_transform(boston.data)[:2]
```

In [ ]:

```
# 使用fit方法, 使reg类基于指定数据估计出回归模型的相应参数
reg.fit(boston.data, boston.target)
```

In [ ]:

```
reg.coef_
```

In [ ]:

```
pred = reg.predict(boston.data)
pred[:10]
```

In [ ]:

```
reg.score(boston.data, boston.target)
```

### **Modelclass中基本通用的类属性**

注意：模型拟合前这些属性可能不存在

coef\_ : array, 多因变量时为二维数组  
intercept\_ : 常数项

classes\_ : 每个输出的类标签  
n\_classes\_ : int or list, 类别数  
n\_features\_ : int, 特征数

loss\_ : 损失函数计算出来的当前损失值  
n\_iter\_ : 迭代次数

In [ ]:

```
std.mean_, std.scale_
```

In [ ]:

```
reg.intercept_, reg.coef_
```

### **简化的调用函数**

特征处理class往往会有简化版本的函数可供调用，功能类似，但使用上更简单。

```
class sklearn.preprocessing.StandardScaler()
sklearn.preprocessing.scale()
```

```
In [ ]:
```

```
preprocessing.scale(boston.data)[:2]
```

### 模型的保存 (持久化)

可以直接使用通过Python的pickle模块将训练好的模型保存为外部文件，但最好使用sklearn中的joblib模块进行操作。

```
In [ ]:
```

```
# 保存为外部文件
from sklearn.externals import joblib

joblib.dump(std, 'f:/std.pkl')
joblib.dump(reg, 'f:/reg.pkl')
```

```
In [ ]:
```

```
# 读入外部保存的模型文件
reg2 = joblib.load('f:/reg.pkl')
reg2.coef_
```

## 1.6 实战练习

尝试使用不同的三种调用方式调用同一个模型。

加载sklearn自带的iris数据集，熟悉该数据集的各种属性，并尝试将其转换为数据框。

尝试在不参考任何帮助文档的情况下，按照sklearn中的标准API操作方式，使用BP神经网络对iris数据进行拟合，并返回各案例的预测类别、预测概率等结果。

BP神经网络对应的类为：`class sklearn.neural_network.MLPClassifier()`  
此处只为API操作演示，不进一步讨论模型拟合前的数据预处理问题

将上题中生成的模型存储为外部文件，并重新读入。

## 2 方差分析模型

### 2.1 一般线性模型概述

### 2.2 方差分析模型的statsmodels实现

statsmodels对方差分析模型的实现，本质上是对通用的一般线性模型框架做了重新的结果解读。

拟合方差分析模型时，使用模型表达式会更为灵活便捷。

patsy包可以提供对模型表达式更好的支持，这里不展开讨论。

#### 2.2.1 方差分析模型的实现