

1. 什么是回归?

2. 回归模型

2.1 线性回归

2.1.1 普通线性回归

2.1.2 岭回归

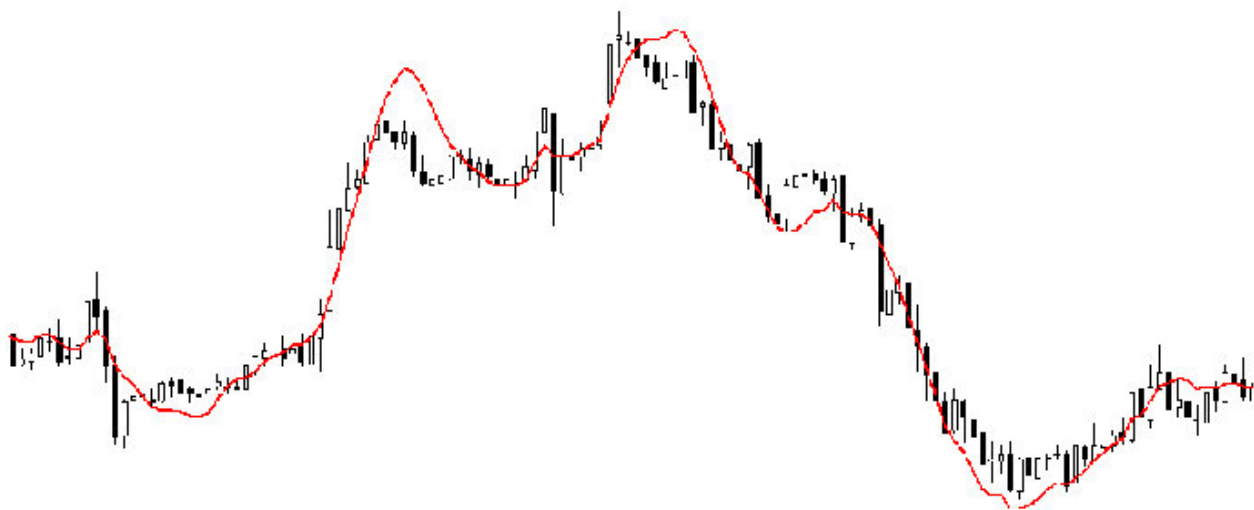
2.2 决策树回归

2.3 SVM回归

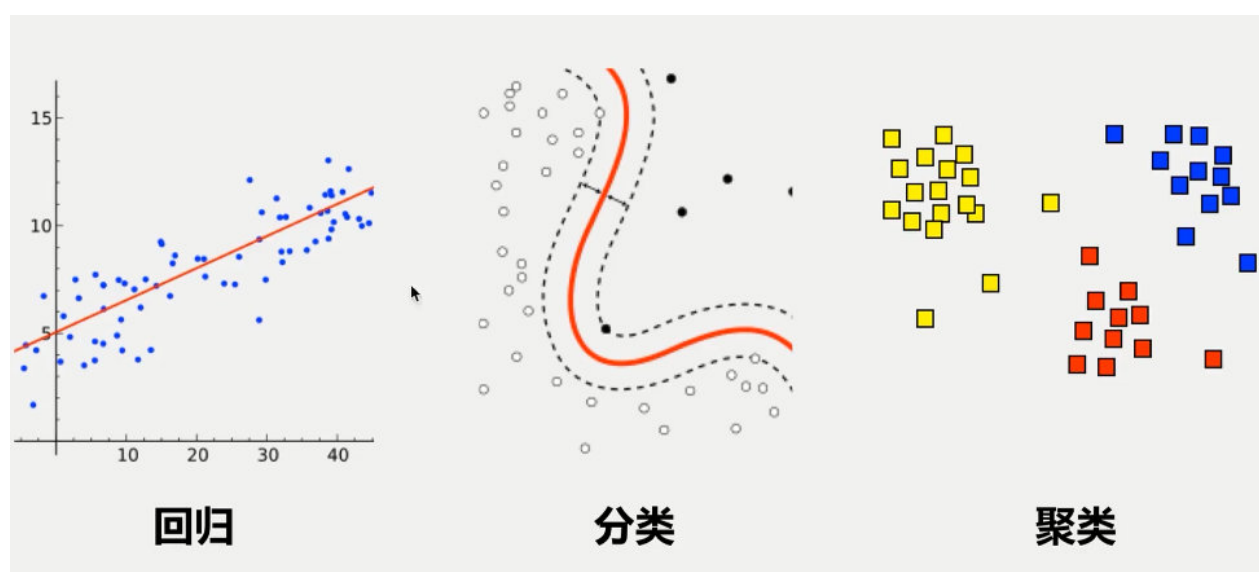
参考

1. 什么是回归?

分类的目标变量是标称型数据，而回归是对连续型数据的预测。回归分析是一种预测建模技术，研究因变量和自变量之间的关系，如销售量预测或制造缺陷预测等，下图中的红线表示的就是回归曲线。



回归不同于分类和聚类，他们的区别可以用下图形象的表达出来。



2. 回归模型

这里使用sklearn进行代码实现，如果想手动实现的话，可以看《机器学习实战》，那本书有部分的算法实现，下面介绍的算法统一使用的函数如下

- 加载数据 数据文件[点这里](#)

```
1 def load_data(file_path):
2     num_feat = len(open(file_path).readline().split("\t")) - 1
3     data_mat = list()
4     lable_mat = list()
5     fr = open(file_path)
6     for line in fr.readlines():
7         line_arr = list()
8         cur_line = line.strip().split("\t")
9         for i in range(num_feat):
10             line_arr.append(float(cur_line[i]))
11         data_mat.append(line_arr)
12         lable_mat.append(float(cur_line[-1]))
13     return data_mat, lable_mat
```

- 绘制回归结果

```
1 def plot_regression(model, x_data, y_data):
2     x_data = np.mat(x_data)
3     y_data = np.mat(y_data).T
4     x_train, y_train = x_data[:150,1:], y_data[:150,:]
5     x_test, y_test = x_data[150:,1:], y_data[150:,:]
6     model.fit(x_train, y_train)
7     score = model.score(x_test, y_test)
8     result = model.predict(x_train)
9     plt.figure()
10    srt_idx = x_train.argsort(0)
11    plt.plot(x_train[srt_idx].reshape(-1,1),
12             y_train[srt_idx].reshape(-1,1), 'go', label = "true value")
13    plt.plot(x_train[srt_idx].reshape(-1,1),
14             result[srt_idx].reshape(-1,1), 'ro-', label = "predict value")
15    plt.title("score:%f" % score)
16    plt.legend()
17    plt.show()
```

2.1 线性回归

2.1.1 普通线性回归

提到回归，首先想到的肯定是线性回归(**linear regression**)，因为它是最容易理解，最简单的回归方法。设待拟合的数据对象为 $X = \{x_1, x_2, \dots, x_m\}$ ，其对应的真实值为 $y = \{y_1, y_2, \dots, y_m\}$ ，线性模型可以写为

$$\hat{y} = Xw \quad (1)$$

其中， w 为回归系数，我们用平方误差来衡量拟合的误差

$$L(X) = \sum_{i=1}^m (y_i - x_i^T w)^2 = (y - Xw)^2 \quad (2)$$

上式对 w 求导等于0可以得到

$$\frac{\partial L(X)}{\partial w} = \frac{\partial(y^T y - w^T X^T y - y^T X w - w^T X^T X w)}{\partial w} = 2X^T (y - Xw) = 0 \quad (3)$$

可以得到

$$\hat{w} = (X^T X)^{-1} X^T y \quad (4)$$

上述方式容易对训练数据欠拟合，一种好的解决方式是局部加权线性回归，为每个误差增加一个权重 w_i （这里的 w 并不是上面的 \hat{w} ），此时误差函数可以写成

$$L(X) = \sum_{i=1}^m w_i (y_i - x_i^T w)^2 = [W(y - Xw)]^2 \quad (5)$$

其中， W 是一个对角矩阵，也叫做核，核的类型可以自由选择，最常见的就是高斯核，高斯核对应的权重如下

$$W(j, j) = \exp\left(\frac{\|x^i - x^j\|^2}{-2k^2}\right) \quad (6)$$

同样的，对新的误差函数 $L(X)$ 求导可以得到此时回归系数为

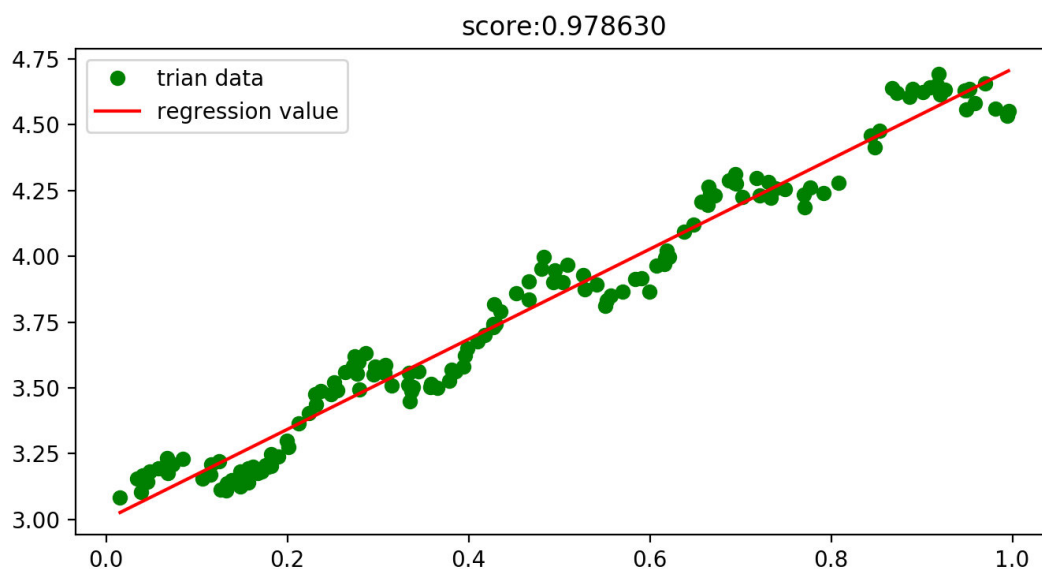
$$\hat{w} = (X^T W X)^{-1} X^T W y \quad (7)$$

这里的 W 其实是 $W^T W$ ，但是使用 W 代替具有同样的意义并且简便。

sklearn调用代码：

```
1 x_data, y_data = load_data("ex0.txt")
2 from sklearn import linear_model
3 # 线性回归
4 model_linear_regression = linear_model.LinearRegression()
5 plot_regression(model_linear_regression, x_data, y_data)
```

绘制出的回归曲线如下图所示



2.1.2 岭回归

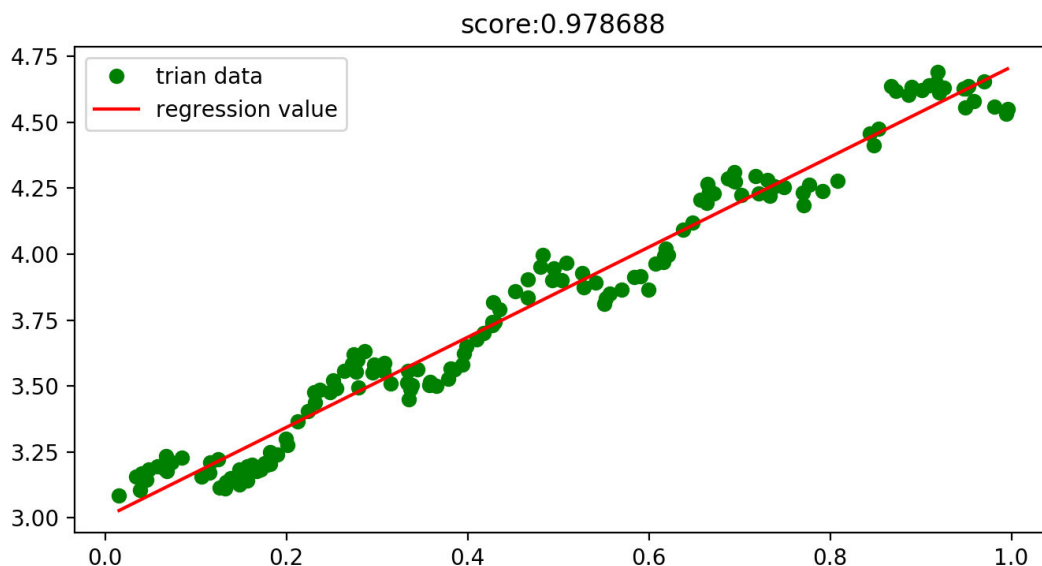
我们看线性回归中的输入集 $X = \{x_1, x_2, \dots, x_m\}$ ，假设其维度为 n ，当 $n > m$ 的时候， X 不是满秩矩阵，无法求解逆矩阵，这时候就需要用到岭回归(ridge regression)了，在矩阵 $X^T X$ 上加上一个 λI 让其成为满秩矩阵，那么这个时候的回归系数为

$$\hat{w} = (X^T X + \lambda I)^{-1} X^T y \quad (8)$$

sklearn调用代码:

```
1 x_data, y_data = load_data("ex0.txt")
2 from sklearn import linear_model
3 # Ridge回归
4 model_ridge = linear_model.Ridge(alpha = 0.01)
5 plot_regression(model_ridge, x_data, y_data)
```

绘制出的回归曲线如下图所示



2.2 决策树回归

决策树学习常用的算法有ID3、C4.5、CART(classification and regression tree)，这介绍用于回归的决策树CART，具体的方法理论参考李航的《统计学习方法》。

我们考虑输入的训练数据 $D = \{X, y\} = \{(x_1, y_1), (x_1, y_1), \dots, (x_m, y_m)\}$ ，一个回归树对应着输入空间（即特征空间）的一个划分以及在划分的单元上的输出值，假设已将输入空间划分为 M 个单元 R_1, R_2, \dots, R_M ，并且在每一个单元上都有一个固定的输出值 c_m ，那么回归树模型可以表示为

$$f(x) = \sum_{m=1}^M c_m I(x \in R_m) \quad (9)$$

其中，函数 I 对应着0-1函数。当输入空间的划分确定时，可以用平方误差 $\sum_{x_i \in R_m} (y_i - f(x_i))$ 来表示回归树对于训练数据的预测误差，用平方误差最小的准则求解每个单元上的最优输出值，那么单元 R_m 上的最优值 \hat{c}_m 是 R_m 上的所有输入实例 x_i 对应的输出 y_i 的均值，即

$$\hat{c}_m = \text{ave}(y_i | x_i \in R_m) \quad (10)$$

上面是整个树的输出形式，关键的问题来了，怎么对输入空间进行划分？这里采用启发式的算法，选择第 j 个变量和它取的值 s 作为切分变量(splitting variable)和切分点(splitting point)，并定义两个区域

$$R_1(j, s) = \{x | x_j \leq s\} \quad R_2(j, s) = \{x | x_j > s\} \quad (11)$$

然后寻找最优切分变量和最优切分点，即

$$m(s) = \min_{j,s} \left[\min_{c_1} \sum_{x_i \in R_1(j,s)} (y_i - c_1)^2 + \min_{c_2} \sum_{x_i \in R_2(j,s)} (y_i - c_2)^2 \right] \quad (12)$$

简单的理解，就是在要求切分点 s 两边的区域的均方差都尽量小的同时，保证两个区域的最小均方差和是最小的。

对每一对 (j, s) ，均值表示为

$$\hat{c}_1 = \text{ave}(y_i | x_i \in R_1(j, s)) \quad \hat{c}_2 = \text{ave}(y_i | x_i \in R_2(j, s)) \quad (13)$$

遍历所有输入变量，找到最优的对 (j, s) ，从而将输入空间切分为两个区域，接着对切分的两个区域重复上述划分过程，直到满足停止条件为止，这样一个回归树的生成就完成了。

举个🍎，输入数据 D 如下表所示。

x_i	1	2	3	4	5	6
y_i	5.56	5.70	5.91	6.40	6.90	7.95

对上述连续型变量，只有一个切分变量，那么考虑切分点为1.5, 2.5, 3.5, 4.5, 5.5。对切分点依次求解 $R_1, R_2, c_1, c_2, m(s)$ ，例如当切分点为2.5时， $R_1 = \{1, 2\}, R_2 = \{3, 4, 5, 6\}$ ，其他的计算如下

$$c_1 = \frac{1}{N_1} \sum_{x_i \in R_1(j,s)} y_i = \frac{1}{2}(5.56 + 5.70) = 5.63 \quad (14)$$

$$c_2 = \frac{1}{N_2} \sum_{x_i \in R_2(j,s)} y_i = \frac{1}{4}(5.91 + 6.40 + 6.90 + 7.95) = 6.79 \quad (15)$$

$$s_m = \min_{j,s} \left[\min_{c_1} \sum_{x_i \in R_1(j,s)} (y_i - c_1)^2 + \min_{c_2} \sum_{x_i \in R_2(j,s)} (y_i - c_2)^2 \right] = 2.294 \quad (16)$$

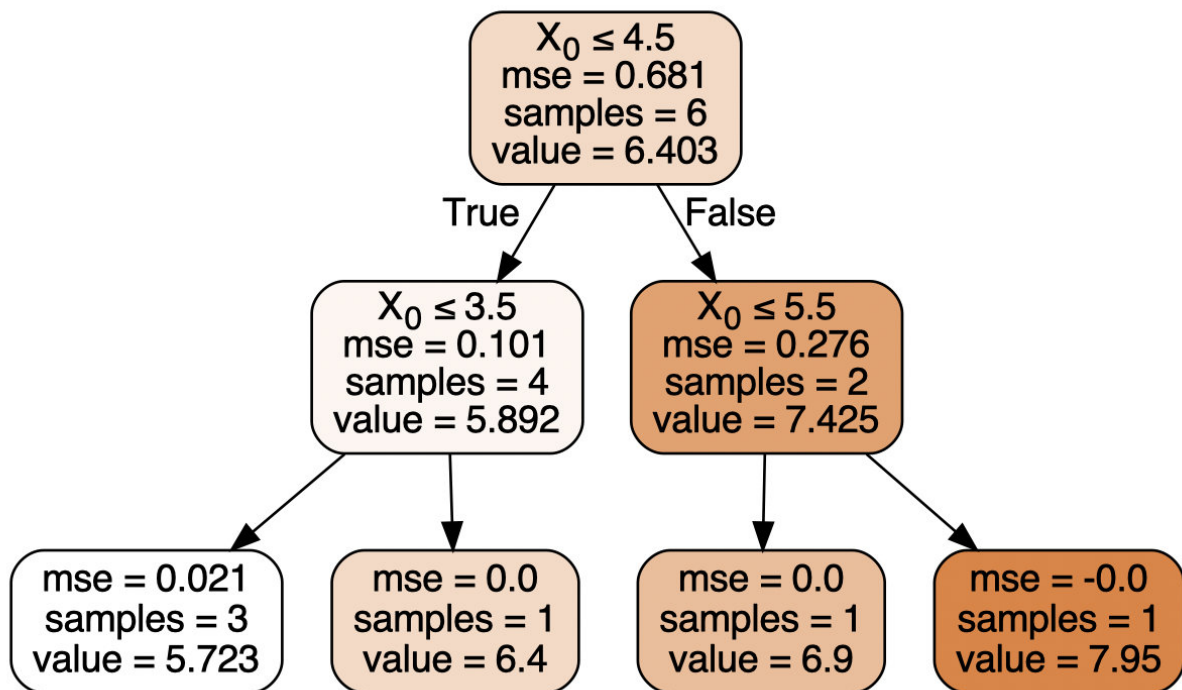
第一次切分时，对象为全体输入，计算出来的 s_m 值如下表所示。

切分点	1.5	2.5	3.5	4.5	5.5
$s(m)$	3.23468	2.294	1.31373333	0.956725	1.21752

可以看到，当 $s = 4.5$ 时，取得最小的 $s(m)$ 值，此时的回归估计值为全体输入的均值6.403，递归求解左子树和右子树的回归估计值，最终求解的回归方程为

$$f(x) = \begin{cases} 5.723 & x \leq 3.5 \\ 6.4 & 3.5 < x \leq 4.5 \\ 6.9 & 4.5 < x \leq 5.5 \\ 7.95 & x > 5.5 \end{cases} \quad (17)$$

这个过程可以使用graphviz模块显示出来。

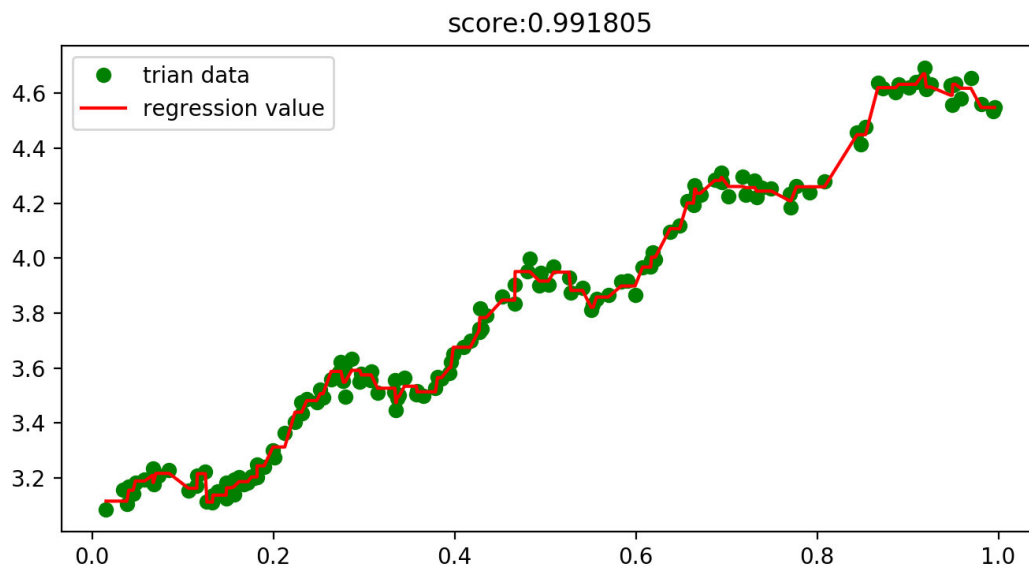


使用本文一开始提到的数据，决策树回归的代码如下

```

1 # 决策树回归
2 from sklearn import tree
3 model_decisiontree_regression =
  tree.DecisionTreeRegressor(min_weight_fraction_leaf=0.01)
4 plot_regression(model_decisiontree_regression, x_data, y_data)

```



2.3 SVM回归

先回顾一下在基本线性可分情况下的SVM模型:

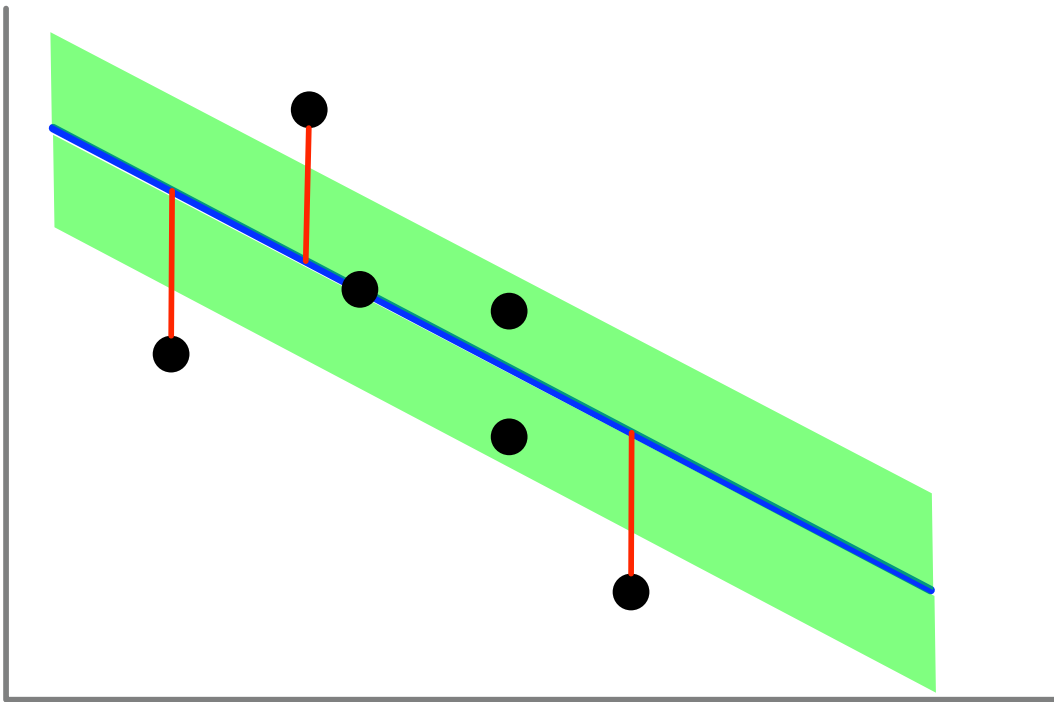
$$\begin{aligned}
\min_{w,b,\xi} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \xi_i \\
\text{s.t.} \quad & -y_i (w \cdot x_i + b) - \xi_i + 1 \leq 0, \quad i = 1, 2, \dots, N \\
& -\xi_i \leq 0, \quad i = 1, 2, \dots, N
\end{aligned} \tag{18}$$

分类SVM模型中要让训练集中的每个样本尽可能远离自己类别一侧的支持向量，回归模型也一样，沿用的是最大建哥分类器的思想。

对于回归模型，优化的目标函数和分类模型保持一致，依然是 $\min_{w,b} \frac{1}{2} \|w\|^2$ ，但是约束条件不一样，回归模型的目标是让训练集中的每个样本点 (x_i, y_i) 尽量拟合到一个线性模型 $y_i = wx_i + b$ 上，对于一般的回归模型使用均方误差MSE作为损失函数的，但是SVM回归不是这样定义的。

SVM需要我们定义一个常量 $\varepsilon > 0$ ，对于某一个点 (x_i, y_i) ，如果 $|y_i - wx_i - b| \leq \varepsilon$ ，则完全没有损失，如果 $|y_i - wx_i - b| > \varepsilon$ ，则对应的损失为 $|y_i - wx_i - b| - \varepsilon$ ，这个和均方差损失不同，对于均方差，只要 $|y_i - wx_i - b| \neq 0$ 就会有损失。

如下图所示，在蓝色条带里面的点是没有损失的，但是在外面的点是有损失的，损失大小为红色线的长度。



总结下，我们的SVM回归模型的损失函数度量度为

$$\text{err}(x_i, y_i) = \begin{cases} 0 & |y_i - w \cdot x_i - b| \leq \varepsilon \\ |y_i - w \cdot x_i - b| - \varepsilon & |y_i - w \cdot x_i - b| > \varepsilon \end{cases} \tag{19}$$

有了损失函数之后，我们就可以定义SVM回归的目标函数为

$$\begin{aligned}
\min \quad & \frac{1}{2} \|w\|_2^2 \\
\text{s.t.} \quad & |y_i - w \cdot x_i - b| \leq \varepsilon (i = 1, 2, \dots, m)
\end{aligned} \tag{20}$$

这个模型的最优解求解过程这里不再赘述，有兴趣的可以看[参考\[2\]或者\[3\]](#)中的论述。

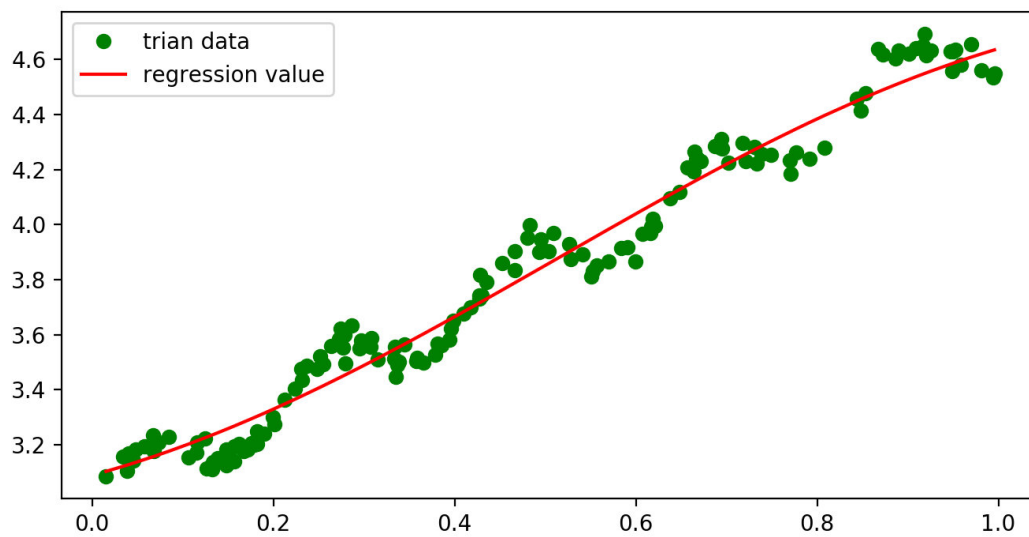
使用本文一开始提到的数据，SVM回归的代码如下

```

1 # SVM回归
2 from sklearn import svm
3 model_svr = svm.SVR()
4 plot_regression(model_svr, x_data, y_data)

```

score:0.980294



参考

[1] 李航. 统计学习方法, 清华大学出版社

[2] [CSDN-SVM回归博客](#)

[3] [cnblog-SVM回归博客](#)