

$$1. \hat{f}(x_0) = x_0^T \hat{\beta}$$

$$E[\hat{f}(x_0)] = E[x_0^T (X^T X)^{-1} X^T \hat{y}]$$

$$\text{注意到 } X = \begin{bmatrix} x_0^T \\ \vdots \\ x_n^T \end{bmatrix}$$

$$y(x_0) = x_0^T \beta^*$$

$$\hat{y} = \begin{bmatrix} y(x_0) \\ y(x_1) \\ \vdots \\ y(x_n) \end{bmatrix} = \begin{bmatrix} x_0^T \beta^* \\ x_1^T \beta^* \\ \vdots \\ x_n^T \beta^* \end{bmatrix} = X \beta^*$$

$$E[\hat{f}(x_0)] = E[x_0^T (X^T X)^{-1} X^T X \beta^*] = x_0^T \beta^* = y(x_0)$$

$$\text{即 } \text{bias}(\hat{f}(x_0)) = E[\hat{f}(x_0)] - y(x_0) = 0$$

$$x_i: p \times 1 \quad X: n \times p \quad y: n \times 1 \quad X^T X: p \times p$$

$$(2) \text{法1} \quad \hat{f}(x_i) = x_i^T \hat{\beta} = x_i^T (X^T X)^{-1} X^T y$$

$x$  是自变量 不具有随机性, 只有  $y$  有随机性

$$\text{Var}(\hat{f}(x_i)) = x_i^T (X^T X)^{-1} X^T \text{Var}(y) X (X^T X)^{-1} x_i$$

$p \times p \quad p \times n \quad n \times n \quad n \times p \quad p \times p \quad p \times 1$

$$\text{Var}(y) = \begin{pmatrix} \sigma^2 & & \\ & \sigma^2 & \\ & & \ddots \\ & & & \sigma^2 \end{pmatrix} = \text{diag}_n(\sigma^2, \sigma^2, \dots, \sigma^2)$$

$$\therefore \text{Var}(\hat{f}(x_i)) = \sigma^2 \cdot x_i^T (X^T X)^{-1} X^T X (X^T X)^{-1} x_i$$

$$\begin{aligned} \sum_{i=1}^n \text{Var}(\hat{f}(x_i)) &= \sigma^2 \sum_{i=1}^n x_i^T [(X^T X)^{-1}]^T x_i \\ &= \sigma^2 \sum_{i=1}^n x_i^T (X^T X)^{-1} x_i \end{aligned}$$

$$(X^T X)^{-1} \text{ 又被称为 投影矩阵, 令 } (X^T X)^{-1} = A = Q^T I Q = P^T P \quad P = (X^T X)^{-\frac{1}{2}} \quad P x_i = (X^T X)^{-\frac{1}{2}} x_i$$

$$\sum_{i=1}^n \text{Var}(\hat{f}(x_i)) = \sigma^2 \left[ \sum_{i=1}^n x_i^T P^T P x_i \right]$$

$$= \sigma^2 \text{tr} \sum_{i=1}^n x_i^T P^T P x_i$$

$$= \sigma^2 \cdot \text{tr}(X^T P^T P X) = \sigma^2 P$$

$$\frac{1}{n} \sum_{i=1}^n \text{Var}(\hat{f}(x_i)) = \frac{P}{n} \sigma^2$$

法II:

证明: 若  $X \sim N(u, \Sigma)$  则  $E(X^T A X) = \text{tr}(A \Sigma) + u^T A u$

$$y = x \beta^* + \varepsilon \sim N(x \beta^*, \sigma^2)$$

$$\text{Var}(\hat{f}(x_i)) = E[(\hat{f}(x_i))^2] - (E \hat{f}(x_i))^2$$

①                      ②

$$\textcircled{1} = E \left[ (x_i^T \hat{\beta})^T (x_i^T \hat{\beta}) \right]$$

$1 \times p \quad p \times 1$

$$= E \left[ \underbrace{x_i^T (X^T X)^{-1} X^T y}_{H_i} \underbrace{x_i^T (X^T X)^{-1} X^T y}_{H_i} \right]$$

$$= E y^T H_i^T H_i y$$

$$= \text{tr}(H_i^T H_i \Sigma) + (x_i \beta^*)^T H_i^T H_i (x_i \beta^*) \rightarrow y_i^2$$

$$\textcircled{2} = (y_i)^2$$

$$\sum_{i=1}^n [\textcircled{1} - \textcircled{2}] = \sum_{i=1}^n \text{tr}(H_i^T H_i \Sigma)$$

$$= \sigma^2 \sum_{i=1}^n \text{tr}(H_i^T H_i)$$

$$= \sigma^2 \cdot \sum_{i=1}^n \text{tr}(x_i (X^T X)^{-1} x_i^T x_i (X^T X)^{-1} X^T)$$

$$= \sigma^2 \cdot \text{tr}(X^T X)^{-1} \sum_{i=1}^n x_i x_i^T$$

$$= \sigma^2 \cdot p$$

$$\frac{1}{n} \sum_{i=1}^n \text{Var}(\hat{f}(x_i)) = \frac{\sigma^2 p}{n}$$

$$(3) E(y_i' - \hat{f}(x_i))^2$$

$$= \text{Var}(\varepsilon) + E[\text{bias}(f(x_i))]^2 + \text{Var}(\hat{f}(x_i))$$

$$= \sigma^2 + 0 + \text{Var}(\hat{f}(x_i))$$

$$\frac{1}{n} \sum_{i=1}^n E(y_i' - \hat{f}(x_i))^2$$

$$= \frac{1}{n} \sum_{i=1}^n (\sigma^2 + \text{Var}(\hat{f}(x_i)))$$

$$= \frac{1}{n} \cdot n \sigma^2 + \frac{\sigma^2 p}{n} = (1 + \frac{p}{n}) \sigma^2$$

## 10\_2

(1) 请由仿真结果计算得到x1与x2的皮尔森相关系数

10000次仿真结果的平均值是:

x1, x2皮尔逊相关系数: 0.438843

x1, x2皮尔逊相关系数: 0.248606

平均下来的x1, x2皮尔逊相关系数: 0.386986

$x_1$ 和 $x_2$ 的皮尔逊相关系数是0.386986

## (2) 请多次生成数据，观察正则化系数为 1 情况下三种模型拟合参数的稳定性

生成30次数据。

### 线性回归

线性回归的斜率和截距的均值及方差

```
beta1的平均值:3.016474,beta0的平均值:0.054433
beta1的方差:0.009815,beta1的方差:0.835609
```

线性回归的斜率的均值是3.016 ,方差是0.0098

截距的均值是0.0544, 方差是0.835

### 岭回归

岭回归的斜率和截距的均值及方差

```
beta1的平均值:2.990340,beta0的平均值:0.196614
beta1的方差:0.013563,beta1的方差:0.653476
```

岭回归的斜率的均值是2.990 ,方差是0.0136

截距的均值是0.197, 方差是0.653

### lasso回归

lasso回归的斜率和截距的均值及方差

```
beta1的方差:0.009823,beta1的方差:0.000353
beta1的平均值:2.939963,beta0的平均值:0.003490
beta1的方差:0.009823,beta1的方差:0.000353
```

lasso回归的斜率的均值是2.940 方差是0.0098

截距的均值是0.0035 方差是0.000353

不难发现, lasso的截距和斜率的方差是最小的, 说明lasso对于此题的稳定性最强, 岭回归的斜率的稳定性不如线性回归, 但是截距的稳定性高于线性回归。

### 采用fisher法

```
.512050 1128280554e-14, 572], [4.655885747,
nums_of_features = 1,train_acc:0.880000
nums_of_features = 1,test_acc:0.890000
nums_of_features = 5,train_acc:0.956667
nums_of_features = 5,test_acc:0.960000
nums_of_features = 10,train_acc:0.966667
nums_of_features = 10,test_acc:0.960000
nums_of_features = 20,train_acc:0.970000
nums_of_features = 20,test_acc:0.950000
nums_of_features = 50,train_acc:1.000000
nums_of_features = 50,test_acc:0.950000
nums_of_features = 100,train_acc:1.000000
nums_of_features = 100,test_acc:0.940000
全部特征,test_acc:0.93
```

特征数	训练集上的acc	测试集上的acc
1	0.88	0.89
5	0.957	0.96
10	0.967	0.96
20	0.97	0.95
50	1	0.95
100	1	0.94
全部特征	1	0.93

显然，做特征选择后预测结果有所上升。原因主要是通过特征选择减少了过拟合，提高了模型的泛化能力，减少了不必要的特征对于分类的影响。

## 最大信息系数

```
nums_of_features = 1,train_acc:0.880000
nums_of_features = 1,test_acc:0.890000
nums_of_features = 5,train_acc:0.963333
nums_of_features = 5,test_acc:0.960000
nums_of_features = 10,train_acc:0.973333
nums_of_features = 10,test_acc:0.950000
nums_of_features = 20,train_acc:0.970000
nums_of_features = 20,test_acc:0.960000
nums_of_features = 50,train_acc:1.000000
nums_of_features = 50,test_acc:0.880000
nums_of_features = 100,train_acc:1.000000
nums_of_features = 100,test_acc:0.910000
全部特征,test_acc:0.93
```

类似的表：

特征数	训练集上的acc	测试集上的acc
1	0.88	0.89
5	0.963	0.96
10	0.973	0.95
20	0.97	0.96
50	1	0.88
100	1	0.91
全部特征	1	0.93

与fisher法的类似，都有过拟合的样子。

除此之外，请比较两种方法在这些特征个数时挑选出的特征子集有多少特征是相同的：

先找出选择多少特征的时候效果比较好：可以看出5个特征或者10个特征的时候效果比较好。

在这个特征数下，寻找相同的特征：

	5	10
fisher	[47, 916, 4, 219, 415]	[47, 916, 219, 4, 415, 476, 271, 224, 825, 461]
最大化信息	[47, 916, 4, 219, 415]	[47, 916, 219, 4, 415, 835, 468, 407, 747, 634]

两者在5个特征的时候，选择的特征一模一样，在10个特征的时候，选择的特征只有前5个一样，后面的都不一样。

这说明在选择特征较少时，可能选择的特征是一样的，但是当选择的特征较多时，选择的特征可能会出现差异。

## 前向算法

- 1) 数据有400个，每一个的都有1000个特征
- 2) 初始化一个空列表 $M_0$ ， $M_i$ 表示存放了 $i$ 个特征。
- 3) 已知 $M_k$ 的时候，剩下 $1000-k$ 个特征，这时候把每个特征分别加入 $M_k$ 中，用logistic回归子啊训练集上训练一次，然后在验证集上看看效果，计算RSS等参数
- 4) 选择RSS最好的特征，与之前的 $M_k$ 一起，变成 $M_{k+1}$
- 5) 在 $M_1, M_2, M_3, \dots, M_{1000}$ 中，选择RSS最高的提取出来。

## 实验结果

当特征选到6个的时候， $R^2$ 已经变成了1,继续是实验下去也没有意义了，所以可以直接结束了

特征数	$R^2$	选择的特征
1	0.196	47
2	0.464	47, 4
3	0.732	47, 4, 103
4	0.866	47, 4, 103, 283
5	0.866	47, 4, 103, 283, 0
6	1	47, 4, 103, 283, 0, 134

与上一题的对比，发现差别还是不小的，这前5个选择之中，只有47和4是两者共有的，这说明算法之间还是存在差异的。我认为这里的差异主要是异步算法的问题。无论是fisher还是最大信息系数，都是在同步计算特征的影响，但前向算法是在固定了前面的特征的情况下，观察后续特征的。这可能会造成比较大的选择差异。

## 决策树算法

```
1032 1033
C:\Users\hutter_sadan\anaconda3\envs\environment_test\python.exe D:/大三下/模式识别与机器学习/hw_10/1033.py
[[47, 0.6331313267176789], [4, 0.1891946068699743], [552, 0.04640269719829376], [916, 0.02853904680203866], [311, 0.02690100430416069], [851, 0.017595625456809486], [117, 0.01743583612306709], [217, 0.012777977044476323], [89, 0.010087876614060258], [856, 0.008967001434720229], [943, 0.008967001434720229]]
```

决策树算法会自动选择特征，用训练集的数据训练决策树，最后选择的特征是

```
[[47, 0.6331313267176789], [4, 0.1891946068699743], [552, 0.04640269719829376], [916, 0.02853904680203866], [311, 0.02690100430416069], [851, 0.017595625456809486], [117, 0.01743583612306709], [217, 0.012777977044476323], [89, 0.010087876614060258], [856, 0.008967001434720229], [943, 0.008967001434720229]]
```

与 (1) 相比, 共同选取了47, 4, 916这三个特征。