

线性回归算法

样本空间 $X \subseteq \mathbb{R}^n$

输入： m 个训练数据 $S = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$

输出： 线性模型 $h_{\mathbf{w}^*, b^*}(\mathbf{x}) = \langle \mathbf{w}^*, \mathbf{x} \rangle + b^*$ ，使得 \mathbf{w}^*, b^* 为如下优化问题的最优解

$$\min_{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}} \frac{1}{m} \sum_{i=1}^m (\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle + b - y^{(i)})^2$$

图 3.1 线性回归算法描述

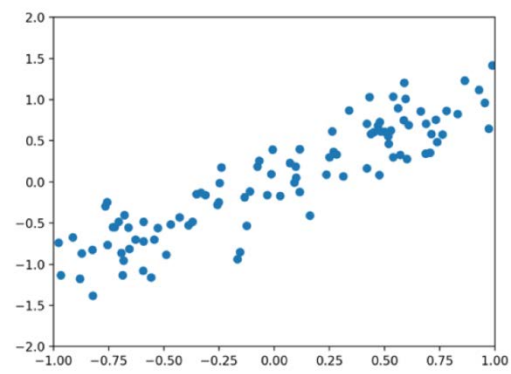


图 3.2 100 个训练数据的散点图

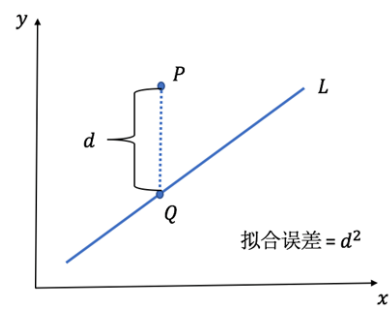


图 3.3 点 P 的拟合误差

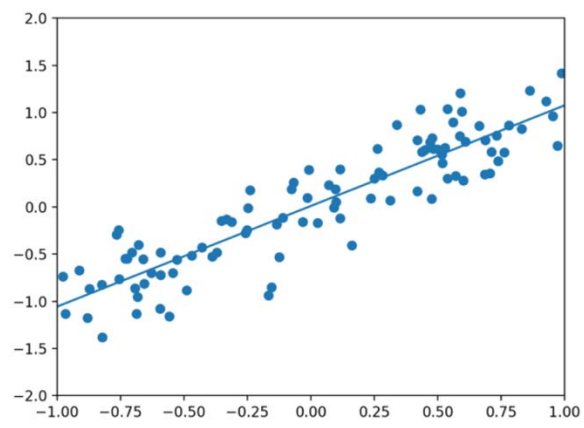


图 3.4 线性回归拟合训练数据

线性回归算法（简化记号）

样本空间 $X \subseteq \mathbb{R}^n$ ，每个样本 $\mathbf{x} \in X$ 首位是 1

输入： m 个训练数据 $S = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$

输出： 线性模型 $h(\mathbf{x}) = \langle \mathbf{w}^*, \mathbf{x} \rangle$ ，使得 \mathbf{w}^* 为如下优化问题的最优解

$$\min_{\mathbf{w} \in \mathbb{R}^n} \frac{1}{m} \sum_{i=1}^m (\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - y^{(i)})^2$$

图 3.5 简化记号的线性回归算法

平均值模型

输入: m 个训练数据 $S = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$

输出: 计算 $\bar{y} = \frac{1}{m} \sum_{i=1}^m y^{(i)}$, 输出常数模型 h_{avg} , 其中 $h_{\text{avg}}(\mathbf{x}) = \bar{y}, \forall \mathbf{x} \in X$

图 3.6 平均值模型

```
machine_learning.lib.linear_regression

1  import numpy as np
2
3  class LinearRegression:
4      def fit(self, X, y):
5          self.w = np.linalg.inv(X.T.dot(X)).dot(X.T).dot(y)
6
7      def predict(self, X):
8          return X.dot(self.w)
9
10 def mean_squared_error(y_true, y_pred):
11     return np.average((y_true - y_pred)**2, axis=0)
12
13 def r2_score(y_true, y_pred):
14     numerator = (y_true - y_pred)**2
15     denominator = (y_true - np.average(y_true, axis=0))**2
16     return 1- numerator.sum(axis=0) / denominator.sum(axis=0)
```

图 3.7 线性回归的正规方程算法

```
1 import numpy as np
2 import machine_learning.lib.linear_regression as lib
3
4 def generate_samples(m):
5     X = 2 * (np.random.rand(m, 1) - 0.5)
6     y = X + np.random.normal(0, 0.3, (m, 1))
7     return X, y
8
9 def process_features(X):
10     m, n = X.shape
11     X = np.c_[np.ones((m, 1)), X]
12     return X
13
14 np.random.seed(0)
15 X_train, y_train = generate_samples(100)
16 X_train = process_features(X_train)
17 X_test, y_test = generate_samples(100)
18 X_test = process_features(X_test)
19
20 model = lib.LinearRegression()
21 model.fit(X_train, y_train)
22 y_pred = model.predict(X_test)
23 mse = lib.mean_squared_error(y_test, y_pred)
24 r2 = lib.r2_score(y_test, y_pred)
25 print("mse = {} and r2 = {}".format(mse, r2))
```

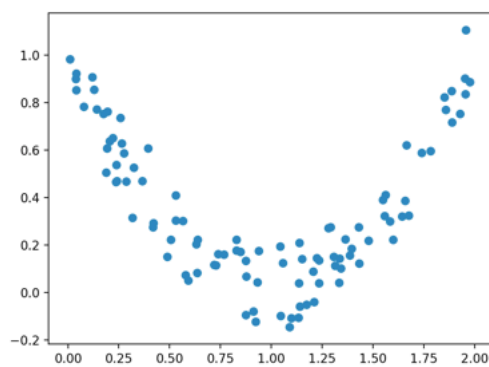
图 3.8 散点的直线拟合


```

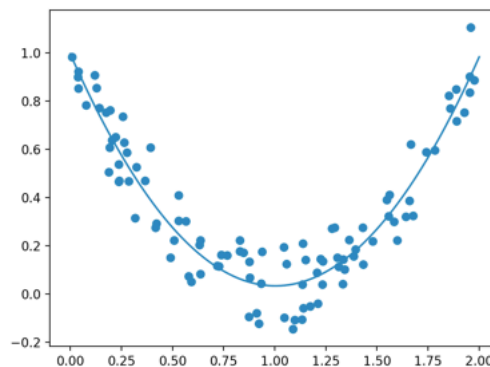
1  import numpy as np
2  from sklearn.model_selection import train_test_split
3  from sklearn.datasets import fetch_california_housing
4  from sklearn.preprocessing import StandardScaler
5  import machine_learning.lib.linear_regression as lib
6
7  def process_features(X):
8      scaler = StandardScaler()
9      X = scaler.fit_transform(X)
10     m,n = X.shape
11     X = np.c_[np.ones((m,1)), X]
12     return X
13
14     housing = fetch_california_housing()
15     X = housing.data
16     y = housing.target.reshape(-1,1)
17     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)
18     X_train = process_features(X_train)
19     X_test = process_features(X_test)
20
21     model = lib.LinearRegression()
22     model.fit(X_train, y_train)
23     y_pred = model.predict(X_test)
24     mse = lib.mean_squared_error(y_test, y_pred)
25     r2 = lib.r2_score(y_test, y_pred)
26     print("mse = {} and r2 = {}".format(mse, r2))

```

图3.9 房价预测问题的线性回归算法



(a)



(b)

图 3.10 多项式模型拟合

```
1 import numpy as np
2 from sklearn.preprocessing import PolynomialFeatures
3 from machine_learning.lib.linear_regression import LinearRegression
4 import matplotlib.pyplot as plt
5
6 def generate_samples(m):
7     X = 2 * np.random.rand(m, 1)
8     y = X**2 - 2 * X + 1 + np.random.normal(0, 0.1, (m, 1))
9     return X, y
10
11 np.random.seed(0)
12 X, y = generate_samples(100)
13 poly = PolynomialFeatures(degree=2)
14 X_poly = poly.fit_transform(X)
15 model = LinearRegression()
16 model.fit(X_poly, y)
17
18 plt.scatter(X, y)
19 W = np.linspace(0, 2, 300).reshape(300, 1)
20 W_poly = poly.fit_transform(W)
21 u = model.predict(W_poly)
22 plt.plot(W, u)
23 plt.show()
```

图 3.11 多项式回归

```
1 import numpy as np
2 from sklearn.preprocessing import PolynomialFeatures
3 from machine_learning.lib.linear_regression import LinearRegression
4 import matplotlib.pyplot as plt
5
6 def generate_samples(m):
7     X = 2 * (np.random.rand(m, 1) - 0.5)
8     y = X + np.random.normal(0, 0.3, (m,1))
9     return X, y
10
11 np.random.seed(100)
12 X, y = generate_samples(10)
13 poly = PolynomialFeatures(degree = 10)
14 X_poly = poly.fit_transform(X)
15 model = LinearRegression()
16 model.fit(X_poly, y)
17
18 plt.axis([-1, 1, -2, 2])
19 plt.scatter(X, y)
20 W = np.linspace(-1, 1, 100).reshape(100, 1)
21 W_poly = poly.fit_transform(W)
22 u = model.predict(W_poly)
23 plt.plot(W, u)
24 plt.show()
```

图 3.12 多项式回归拟合例 2.3 的数据

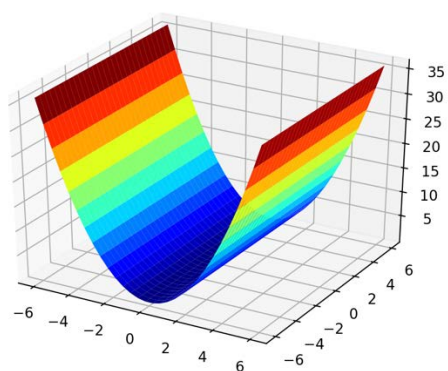


图 3.13 $F(w_1, w_2) = w_1^2$ 的 3 维图像

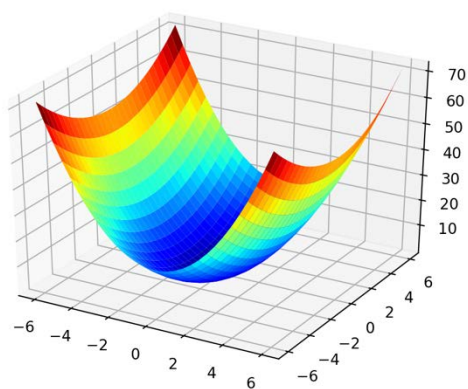


图 3.14 例 3.7 正则化目标函数的 3 维图像

```
machine_learning.lib.ridge_regression

1  import numpy as np
2
3  class RidgeRegression:
4      def __init__(self, Lambda):
5          self.Lambda = Lambda
6
7      def fit(self, X, y):
8          m, n = X.shape
9          r = m * np.diag(self.Lambda * np.ones(n))
10         self.w = np.linalg.inv(X.T.dot(X) + r).dot(X.T).dot(y)
11
12     def predict(self, X):
13         return X.dot(self.w)
```

图 3.15 岭回归算法

```

1  import numpy as np
2  from sklearn.preprocessing import PolynomialFeatures
3  import matplotlib.pyplot as plt
4  from machine_learning.lib.ridge_regression import RidgeRegression
5
6  def generate_samples(m):
7      X = 2 * (np.random.rand(m, 1) - 0.5)
8      y = X + np.random.normal(0, 0.3, (m, 1))
9      return X, y
10
11  np.random.seed(100)
12  X, y = generate_samples(10)
13  poly = PolynomialFeatures(degree = 10)
14  X_poly = poly.fit_transform(X)
15  model = RidgeRegression(Lambda = 0.01)
16  model.fit(X_poly, y)
17
18  plt.axis([-1, 1, -2, 2])
19  plt.scatter(X, y)
20  W = np.linspace(-1, 1, 100).reshape(100, 1)
21  W_poly = poly.fit_transform(W)
22  u = model.predict(W_poly)
23  plt.plot(W, u)
24  plt.show()

```

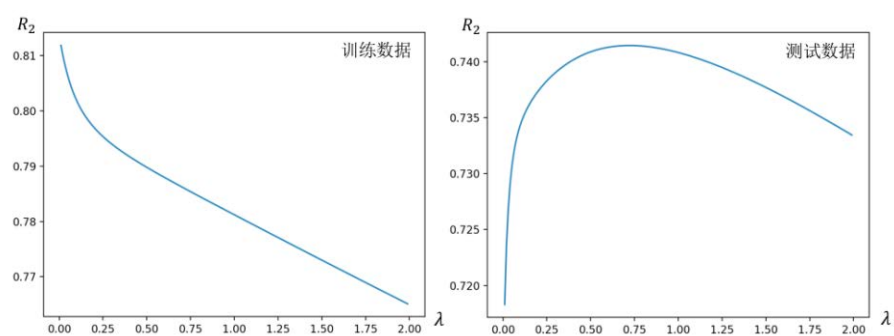
图 3.16 多项式模型和岭回归模型


```

1  import numpy as np
2  from sklearn.preprocessing import PolynomialFeatures
3  import matplotlib.pyplot as plt
4  import machine_learning.lib.linear_regression as lib
5  from machine_learning.lib.ridge_regression import RidgeRegression
6
7  def generate_samples(m):
8      X = 2 * (np.random.rand(m, 1) - 0.5)
9      y = X + np.random.normal(0, 0.3, (m,1))
10     return X, y
11
12     np.random.seed(100)
13     poly = PolynomialFeatures(degree = 10)
14     X_train, y_train = generate_samples(30)
15     X_train = poly.fit_transform(X_train)
16     X_test, y_test = generate_samples(100)
17     X_test = poly.fit_transform(X_test)
18
19     Lambdas, train_r2s, test_r2s = [], [], []
20     for i in range(1, 200):
21         Lambda = 0.01 * i
22         Lambdas.append(Lambda)
23         ridge = RidgeRegression(Lambda)
24         ridge.fit(X_train, y_train)
25         y_train_pred = ridge.predict(X_train)
26         y_test_pred = ridge.predict(X_test)
27         train_r2s.append(lib.r2_score(y_train, y_train_pred))
28         test_r2s.append(lib.r2_score(y_test, y_test_pred))
29
30     plt.figure(0)
31     plt.plot(Lambdas, train_r2s)
32     plt.figure(1)
33     plt.plot(Lambdas, test_r2s)
34     plt.show()

```

图 3.17 拟合岭回归模型



(a)

(b)

图 3.18 正则化系数与决定系数的关系

向前逐步回归算法

StepwiseRegression(\mathbf{X}, \mathbf{y}):

$A = \{1\}, C = \{2, \dots, n\}$

For $i = 2, \dots, n$:

$$mse_A = \min_{\mathbf{w} \in \mathbb{R}^{|A|}} \frac{1}{m} \|\mathbf{X}_A \mathbf{w} - \mathbf{y}\|^2$$

For $j \in C$:

$$mse_{A \cup j} = \min_{\mathbf{w} \in \mathbb{R}^{|A|+1}} \frac{1}{m} \|\mathbf{X}_{A \cup j} \mathbf{w} - \mathbf{y}\|^2$$

$$j^* = \operatorname{argmin}_{j \in C} mse_{A \cup j}$$

If $mse_A / mse_{A \cup j^*}$ pass F test:

$$A \leftarrow A \cup j^*$$

$$C \leftarrow C \setminus j^*$$

Else:

Break

Return A

图 3.19 向前逐步回归算法描述

```

machine_learning.lib.stepwise_regression

1  import numpy as np
2  from scipy.stats import f
3
4  class StepwiseRegression:
5      def fit(self, X, y):
6          return np.linalg.inv(X.T.dot(X)).dot(X.T).dot(y)
7
8      def compute_mse(self, X, y):
9          w = self.fit(X,y)
10         r = y - X.dot(w)
11         return r.T.dot(r)
12
13     def f_test(self, mse_A, mse_min, m):
14         if mse_min > mse_A:
15             return False
16         F = mse_A / mse_min
17         p_value = f.cdf(F, m, m)
18         return p_value > 0.95
19
20     def forward_selection(self, X, y):
21         m,n = X.shape
22         A, C = [0], [i for i in range(1,n)]
23         while len(C) > 0:
24             MSE_A = self.compute_mse(X[:, A], y)
25             MSE_min, j_min = float("inf"), -1
26             j_min = -1
27             for j in C:
28                 MSE_j = self.compute_mse(X[:, A + [j]], y)
29                 if MSE_j < MSE_min:
30                     MSE_min, j_min = MSE_j, j
31                 if self.f_test(MSE_A, MSE_min, m):
32                     A.append(j_min)
33                     C.remove(j_min)
34             else:
35                 break
36         self.w = self.fit(X[:, A], y)
37         self.A = A
38
39     def predict(self, X):
40         return X[:, self.A].dot(self.w)

```

图 3.20 向前逐步回归算法实现

```
1 import numpy as np
2 from sklearn.preprocessing import PolynomialFeatures
3 from machine_learning.lib.stepwise_regression import StepwiseRegression
4
5 def generate_samples(m):
6     X = 2 * (np.random.rand(m, 1) - 0.5)
7     y = X + np.random.normal(0, 0.3, (m,1))
8     return X, y
9
10 np.random.seed(100)
11 X, y = generate_samples(10)
12 poly = PolynomialFeatures(degree = 10)
13 X_poly = poly.fit_transform(X)
14 model = StepwiseRegression()
15 model.forward_selection(X_poly, y)
16 print(model.A, model.w)
```

图 3.21 向前逐步回归与过度拟合

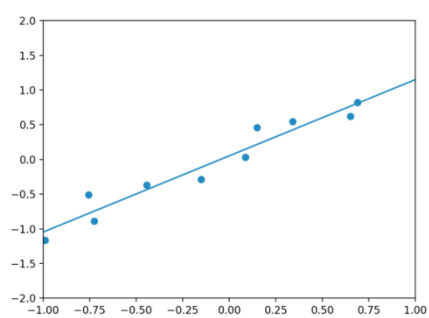


图 3.22 向前逐步回归拟合效果

```
1 import numpy as np
2 from machine_learning.lib.stepwise_regression import StepwiseRegression
3
4 X= np.array(
5     [[ 0.06,  0.34,  0.03]
6      ,[ 0.44,  0.76,  0.28]
7      ,[ 0.86,  0.44,  0.20]
8      ,[ 0.26,  0.09,  0.25]])
9 y = np.array([[ 0.42], [ 1.32 ], [ 0.84], [ 0.61]])
10
11 model = StepwiseRegression()
12 model.forward_selection(X, y)
13 print(model.A, model.w)
```

图 3.23 向前逐步回归得到次优解

分段回归算法

StagewiseRegression(\mathbf{X} , \mathbf{y} , N , η):

$\mathbf{w} = (0, 0, \dots, 0), t = 0$

While $t < N$:

$\mathbf{r} = \mathbf{y} - \mathbf{X}\mathbf{w}$

$j^* = \operatorname{argmax}_{1 \leq j \leq n} |\operatorname{corr}(\mathbf{X}_j, \mathbf{r})|$

$\mathbf{w} \leftarrow \mathbf{w} + \eta \cdot \operatorname{sign}(\operatorname{corr}(\mathbf{X}_{j^*}, \mathbf{r})) \cdot \mathbf{X}_{j^*}$

$t \leftarrow t + 1$

Return \mathbf{w}

图 3.24 分段回归算法描述

```
machine_learning.lib.stagewise_regression

1  import numpy as np
2
3  class StagewiseRegression:
4      def feature_selection(self, X, y, N, eta):
5          m, n = X.shape
6          norms = np.linalg.norm(X, 2, axis=0).reshape(-1, 1)
7          w = np.zeros(n)
8          t = 0
9          r = y
10         while t < N:
11             c = X.T.dot(r) / norms
12             j_max = np.argmax(abs(c))
13             delta = eta * np.sign(c[j_max])
14             w[j_max] = w[j_max] + delta
15             r = r - delta * X[:, j_max].reshape(-1,1)
16             t = t + 1
17         self.w = w
18         return w
19
20     def predict(self, X):
21         return X.dot(self.w)
```

图 3.25 分段回归算法实现

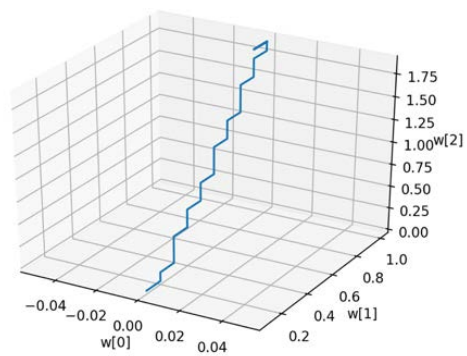


图 3.26 分段回归算法的搜索轨迹