

Logistic Supplement

Tomoyo

November 21, 2019

Instruction: Derivation for two-class classification.

1 Linear regression:

- 1.1 basis
- 1.2 derivation

1.1 definition

$$\text{probability} : P(x_i) = \text{Pr}(y = 1|x_i) = \frac{1}{1 + e^{-(\beta_0 + \sum_{j=1}^m \beta_j x_{ji})}}$$

$$\text{likelihood} : L = \prod_{i=0}^n p(x_i)^{y_i} * [1 - p(x_i)]^{1-y_i}$$

$$P_\alpha(\beta) = \sum_{j=1}^m \left[\frac{1}{2} (1 - \alpha) \beta_j^2 + \alpha |\beta_j| \right]$$

$$NLL = - \sum_{i=1}^n [y_i \log(p(x_i)) + (1 - y_i) \log(1 - p(x_i))]$$

1.2 derivation

NLL can be written as follow.

$$L = - \sum_{i=1}^n [y_i (\beta_0 + \sum_{j=1}^m \beta_j x_{ij}) - \log(1 + e^{\beta_0 + \sum_{j=1}^m \beta_j x_{ij}})]$$

calculation of grad, hessian

$$g_j = \frac{\partial L}{\partial \beta_j} = - \sum_{i=1}^n [y_i - p(x_i)] x_{ij}$$
$$h_j = \frac{\partial g}{\partial \beta_j} = \sum_{i=1}^n x_{ij}^2 p(x_i) (1 - p(x_i))$$

2 gradient descent

Using g as direction, optimize function based on g , when β_j in k iteration:

$$\beta_j^{k+1} = \beta_j^k - \eta g_j$$

for intercept:

$$\text{intercept}^{(k+1)} = \text{intercept}^{(k)} - \eta[-(y_i - p_i)]$$

3 Newton's method

Taylor's expansion for L in $\beta_j^{(k)}$

$$\begin{aligned} L(\beta_j) &= L(\beta_j^{(k)}) + (\beta_j - \beta_j^{(k)})g_j + \frac{1}{2}(\beta_j - \beta_j^{(k)})^2 h_j + o(n^3) \\ &= \beta_j g_j + \frac{1}{2}\beta_j^2 h_j - \beta_j \beta_j^{(k)} h_j + \frac{(\beta_j^{(k)})^2}{2} h_j + L - \beta_j^{(k)} g_j \\ &= \left(\frac{1}{2}h_j\right)\beta_j^2 + (g_j - \beta_j^{(k)} h_j)\beta_j + L - \beta_j^{(k)} g_j + \frac{(\beta_j^{(k)})^2}{2} h_j \\ \text{donation : } M &= \left(\frac{1}{2}h_j, N = g_j - \beta_j^{(k)} h_j, C = L - \beta_j^{(k)} g_j + \frac{(\beta_j^{(k)})^2}{2} h_j\right. \\ &= M\beta_j^2 + N\beta_j + C, \text{ so that} \\ L &= M\beta_j^2 + N\beta_j + C \end{aligned}$$

The solution is:

$$\beta_j = -\frac{N}{2M} = -\frac{g_j - \beta_j^{(k)} h_j}{h_j} = \beta_j^{(k)} - \frac{g_j}{h_j}$$

4 IRLS(Iterative Reweighted Least Squares)

The Newton's solution can be rewritten as follow.

$$\begin{aligned}
\beta_j &= \frac{h_j \beta_j^{(k)} - g_j}{h_j} \\
&= \frac{\sum_{i=1}^n x_{ij}^2 p(x_i)(1 - p(x_i)) \beta_j^{(k)} + \sum_{i=1}^n [y_i - p(x_i)] x_{ij}}{h_j} \\
&= \frac{\sum_{i=1}^n \{x_{ij}^2 p(x_i)(1 - p(x_i)) \beta_j^{(k)} + [y_i - p(x_i)] x_{ij}\}}{h_j} \\
&= \frac{\sum_{i=1}^n x_{ij} p(x_i)(1 - p(x_i)) \{x_{ij} \beta_j^{(k)} + \frac{[y_i - p(x_i)]}{p(x_i)(1 - p(x_i))}\}}{h_j} \\
\text{donation : } w_i &= p(x_i)(1 - p(x_i)), z_i = \frac{[y_i - p(x_i)]}{p(x_i)(1 - p(x_i))} + x_{ij} \beta_j^{(k)} \\
&= \frac{\sum_{i=1}^n x_{ij} w_i z_i}{h_j}
\end{aligned}$$

the solution is equal to

$$\begin{aligned}
\sum_{i=1}^n [x_{ij}^2 p(x_i)(1 - p(x_i))] \beta_j &= \sum_{i=1}^n x_{ij} w_i z_i \\
\sum_{i=1}^n x_{ij} w_i (x_{ij} \beta_j - z_i) &= 0
\end{aligned}$$

equals to the minimize the problem:

$$\arg \min_{\beta_j} \sum_{i=1}^n w_i (x_{ij} \beta_j - z_i)^2$$

equals to:

$$\arg \min_{\beta_j} \sum_{i=1}^n w_i \left(\sum_{j=1}^m x_{ij} \beta_j + \beta_0 - z_i \right)^2$$

$$\text{where } z_i = \frac{[y_i - p(x_i)]}{p(x_i)(1 - p(x_i))} + \sum_{k=1}^m x_{ik} \beta_k^{(k)}$$

So the minimize the NLL is equal to

$$\arg \min_{\beta} \sum_{i=1}^n w_i \left(\sum_{j=1}^m x_{ij} \beta_j + \beta_0 - z_i \right)^2$$

5 implementation

```
1 from sklearn.datasets import load_breast_cancer
2 from joblib import Parallel, delayed
3 import gc
4 import numpy as np
5 import pandas as pd
6
7 class Logistic_Regression(object):
8     def __init__(self, eta = 0.005, tol = 10e-4):
9         self.eta = eta
10        self.beta = None
11        self.intercept = 0
12        self.tol = tol
13    def process_helper(self, x):
14        x = np.where(x > 1 - 10**-4, 1 - 10**-4, x)
15        x = np.where(x < 10**-4, 10**-4, x)
16        return x
17    def standard_function(self, data):
18        return (data - data.mean(axis = 0))/data.std(axis = 0)
19    def gradient_descent(self, X, y, beta, intercept, parallel = False):
20        linear = intercept + np.sum(beta * X, axis = 1)
21        prob = 1/(1 + np.exp(-linear))
22        new_beta = np.zeros(beta.shape)
23        if parallel:
24            def cal_grad(y, prob, x, beta, eta):
25                return beta - eta * np.sum(-(y - prob) * x)
26            new_beta = Parallel(n_jobs = -1, verbose = 0)\
27                (delayed(cal_grad)(target, prob, X[:,j], beta[j], eta) for j in range(data-
28                new_beta = np.array(new_beta)
29        else:
30            for j in range(X.shape[1]):
31                grad = -np.sum((y - prob) * X[:, j])
32                new_beta[j] = beta[j] - self.eta * grad
33        new_intercept = intercept - self.eta * np.sum(-(y - prob))
34        return new_beta, new_intercept
35    def newton_descent(self, X, y, beta, intercept):
36        linear = intercept + np.sum(beta * X, axis = 1)
37        prob = 1/(1 + np.exp(-linear))
38        new_beta = np.zeros(beta.shape)
39        prob = self.process_helper(prob)
40        for j in range(X.shape[1]):
41            grad = -np.sum((y - prob) * X[:, j])
42            hess = np.sum(prob * (1 - prob) * X[:, j] ** 2)
43            new_beta[j] = beta[j] - self.eta * grad/hess
44        grad_intercept = -np.sum(y - prob)
```

```

45     hess_intercept = np.sum(prob * (1 - prob))
46     new_intercept = intercept - self.eta * grad_intercept/hess_intercept
47     return new_beta, new_intercept
48 def IRLS(self, X, y, beta, intercept):
49     new_beta = np.zeros(beta.shape)
50     linear = intercept + np.sum(beta * X, axis = 1)
51     prob = 1/(1 + np.exp(-linear))
52     prob = self.process_helper(prob)
53     w = prob * (1 - prob)
54     for j in range(X.shape[1]):
55         h = np.sum(w * X[:, j] ** 2)
56         z = (y - prob)/w + beta[j] * X[:, j]
57         beta[j] = np.sum(X[:, j] * w * z)/h
58     new_intercept = np.sum(w * (intercept + (y - prob)/w))/np.sum(w)
59     return beta, new_intercept
60 def fit(self, X, y, method = 'newton'):
61     self.beta = np.zeros(X.shape[1])
62     for i in range(1000):
63         converge = 1
64         beta_p, intercept_p = self.beta, self.intercept
65         if method == 'newton':
66             self.beta, self.intercept = self.newton_descent(data_std, target, self)
67         elif method == 'gradient':
68             self.beta, self.intercept = self.gradient_descent(data_std, target, self)
69         elif method == 'IRLS':
70             self.beta, self.intercept = self.IRLS(data_std, target, self.beta, self)
71         else:
72             raise Exception("only support ['gradient', 'newton', 'IRLS']")
73         converge = np.sum((beta - beta_p) ** 2) + (intercept - intercept_p) ** 2
74         if converge < self.tol:
75             break
76 def predict(self, X):
77     linear = self.intercept + np.sum(self.beta * X, axis = 1)
78     prob = 1/(1 + np.exp(-linear))
79     return np.where(prob > 0.5, 1, 0)

```