# Note for Regularization Paths for Generalized Linear Models

Tomoyo

November 10, 2019

Instruction: Derivation for two-class classification.

# 1 Linear regression:

- 1.1 basis
- 1.2 derivation

## 1.1 definition

$$probability : Pr(y=1|x) = \frac{1}{1+e^{-(\beta_0+\sum\limits_{i=1}^{m}\beta_i x_i)}}$$

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

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

## 1.2 derivation

Based on definition, penalized log-likelihood can be presented as follow.

$$L = \sum_{i}^{n}\{y_i log(p(x_i)) + (1-y_i)log(1-p(x_i))\} - \sum_{j=1}^{m}[\frac{1}{2}(1-\alpha)\beta_j^2 + \alpha|\beta_j|]$$

Note that the first part can be written as follow.

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

$$= \sum_{i=1}^{n}[y_i(\beta_0 + \sum_{i=1}^{m}\beta_i x_i) + log(\frac{e^{-(\beta_0 + \sum_{i=1}^{m}\beta_i x_i)}}{1 + e^{-(\beta_0 + \sum_{i=1}^{m}\beta_i x_i)}})]$$

$$= \sum_{i=1}^{n}[y_i(\beta_0 + \sum_{i=1}^{m}\beta_i x_i) + log(1 + e^{(\beta_0 + \sum_{i=1}^{m}\beta_i x_i)})]$$

Using quadratic approximation(taylor expansion), The log-likelihood can be written as follow.

$$l_Q = -\frac{1}{2}\sum_{i=1}^{N}w_i(z_i - \beta_0 - \sum_{i=1}^{m}\beta_i x_i)^2 + C(\hat{\beta}_0, \hat{\beta})^2$$

$$z_i = \hat{\beta}_0 + \sum_{i=1}^{m}\hat{\beta}_i x_i + \frac{y_i - \hat{p}(x_i)}{\hat{p}(x_i)(1 - \hat{p}(x_i))}$$

$$w_i = \hat{p}(x_i)(1 - \hat{p}(x_i))$$

Using coordinate descent to solve the penalized weighted least-squares problem, which is

$$\min_{\beta_0, \beta}\{-l_q(\beta_0, \beta) + \lambda P_\beta)\}$$

derive derivative to solve this problem, for $\beta_j, j \neq 0$:

$$part_1 = \frac{\partial L}{\partial \beta_j} = \sum_{i=1}^{N}w_i[z_i - (\beta_0 + \sum_{j=1}^{m}\beta_j x_j)]x_{ij}$$

$$= \sum_{i=1}^{N}w_i[z_i - (\beta_0 + \sum_{k \neq j}^{m}\beta_k x_k)]x_{ij} - \sum_{i=1}^{N}w_i x_{ij}^2 \beta_j$$

$$= \sum_{i=1}^{N}w_i[\hat{\beta}_j x_{ij} + \frac{y_i - \hat{p}(x_i)}{\hat{p}(x_i)(1 - \hat{p}(x_i))}]x_{ij} - \sum_{i=1}^{N}w_i x_{ij}^2 \beta_j$$

$$= \sum_{i=1}^{N}[w_i x_{ij}^2 \hat{\beta}_j + (y_i - \hat{p}(x_i))x_{ij}] - \sum_{i=1}^{N}w_i x_{ij}^2 \beta_j$$

$$denotation, P = \sum_{i=1}^{N}w_i x_{ij}^2, Q = \sum_{i=1}^{N}x_{ij}(y_i - \hat{p}(x_i))$$

$$= P * \hat{\beta}_j + Q - P * \beta_j$$

Now,derive the penalty part.

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

Using Subderivative, the equation can be written as follow.

$$part_2 = \frac{\partial P_\alpha(\beta)}{\partial \beta_j} = [(1-\alpha)\beta_j + \alpha\partial|\beta_j|]$$

$$= \begin{cases} (1-\alpha)\beta_j + \alpha, & \beta_j > 0 \\ [(1-\alpha)\beta_j - \alpha, (1-\alpha)\beta_j + \alpha], & \beta_j = 0 \\ (1-\alpha)\beta_j - \alpha, & \beta_j < 0 \end{cases}$$

to find the minimum point, $-part_1 + \lambda part_2 = 0$, which is

$$0 = \begin{cases} \lambda[(1-\alpha)\beta_j + \alpha] - (P\hat{\beta}_j + Q - P\beta_j), & \beta_j > 0 \\ \{\lambda[(1-\alpha)\beta_j - \alpha] - (P\hat{\beta}_j + Q - P\beta_j), \lambda[(1-\alpha)\beta_j + \alpha] - (P\beta_j\} + \hat{Q} - P\beta_j), & \beta_j = 0 \\ \lambda[(1-\alpha)\beta_j - \alpha] - (P\hat{\beta}_j + Q - P\beta_j), & \beta_j < 0 \end{cases}$$

Take $\beta_j > 0$ as an example.

$$\lambda\beta_j - (P\hat{\beta}_j + Q - P\beta_j) = 0$$

$$\beta_j = \frac{Q - \lambda\alpha + P\hat{\beta}_j}{P + \lambda(1-\alpha)}$$

$$condition : Q + P\hat{\beta}_j > \lambda\alpha$$

Implmentation as follow

```
1   from joblib import Parallel, delayed
2   import gc
3   import numpy as np
4   import pandas as pd
5
6
7   class Lasso_CR(object):
8       def __init__(self, **kwargs):
9           self.beta = 0
10          self.intercept = 0
11          self.lambda_ = 0
12          self.max_iter = 500
13          self.lambda_ = kwargs.get('lambda_')
14
15      def standardlize(self, data):
16          return (data - data.mean(axis = 0))/data.std(axis = 0)
17
18      def soft_thresholding(self, beta_old, P, Q, threshold):
19          if P * beta_old + Q > threshold:
20              beta_new = beta_old + (Q - threshold)/P
21          elif P * beta_old + Q < - threshold:
```

```python
22             beta_new = beta_old + (Q + threshold)/P
23         else:
24             beta_new = 0
25         return beta_new
26     def coordinate_descent(self, x, y, beta, intercept, lambda_):
27         def process_helper(x):
28             x = np.where(x > 1 - 10**-5, 1 - 10**-5, x)
29             x = np.where(x < 10**-5, 10**-5, x)
30             return x
31         for j in range(len(beta)):
32             linear = intercept + np.sum(beta * x, axis = 1)
33             prob = 1/(1 + np.exp(-linear))
34             prob = process_helper(prob)
35             w = prob * (1 - prob)
36             P = np.mean(w * (x[:, j] ** 2))
37             Q = np.mean((y - prob) * x[:, j])
38             beta[j] = self.soft_thresholding(beta[j], P, Q, lambda_)
39
40         linear = intercept + np.sum(beta * x, axis = 1)
41         prob = process_helper(1/(1 + np.exp(-linear)))
42         w_ = prob * (1 - prob)
43         dir_ = y - prob
44         intercept_new = intercept + np.mean(dir_)/np.mean(w)
45         return beta, intercept_new
46
47     def fit(self, X, y):
48         X_std = self.standardlize(X)
49         max_iter = 500
50         self.beta = np.ones(X.shape[1])/X_std.shape[1]
51         self.intercept = 0
52         for iteration in range(max_iter):
53             self.beta, self.intercept = self.coordinate_descent(X_std, y,
54                 self.beta, self.intercept, self.lambda_)
```