

Loss function. The standard linear model for regression can be written as the dot product between the independent variables and their respective coefficients: $y_i = \mathbf{x}_i^T \boldsymbol{\beta}$. Finding the probability of an input sample being positively classified is performed by adding the logit link function to generate a general linear model. As such, the probability of a positive classification conditioned on the input sample can be written as:

$$P(Y = 1 | \mathbf{x}_i \boldsymbol{\beta}) = \frac{e^{\mathbf{x}_i^T \boldsymbol{\beta}}}{1 + e^{\mathbf{x}_i^T \boldsymbol{\beta}}} \quad (1)$$

However, this transforms the output to a probability between 0 and 1. For ML purposes, I transform the link function's range from $[0, 1]$ to $[-1, 1]$.

$$P(Y = 1 | \mathbf{x}_i \boldsymbol{\beta}) = \frac{2e^{\mathbf{x}_i^T \boldsymbol{\beta}}}{1 + e^{\mathbf{x}_i^T \boldsymbol{\beta}}} - 1 \quad (2)$$

Classifying any input sample has a binary output so can be modeled via a Bernoulli distribution. We can therefore write likelihood of a positive classification as:

$$\begin{aligned} \mathcal{L}(\boldsymbol{\beta}) &= \prod_{i=1}^n P(y_i = 1 | \mathbf{x}_i \boldsymbol{\beta})^{y_i} P(y_i = -1 | \mathbf{x}_i \boldsymbol{\beta})^{1-y_i} \\ &= \prod_{i=1}^n (1 - P(y_i = 1 | \mathbf{x}_i \boldsymbol{\beta}))^{1-y_i} P(y_i = 1 | \mathbf{x}_i \boldsymbol{\beta})^{y_i} \end{aligned}$$

Replacing with equation 2 we find the likelihood and log-likelihood to be:

$$\begin{aligned} (1 - P(y_i = 1 | \mathbf{x}_i \boldsymbol{\beta}))^{1-y_i} P(y_i = 1 | \mathbf{x}_i \boldsymbol{\beta})^{y_i} &= \left(\frac{2e^{\mathbf{x}_i \boldsymbol{\beta}}}{1 + e^{\mathbf{x}_i \boldsymbol{\beta}}} - 1 \right)^{1-y_i} 2^{y_i} \left(1 - \frac{e^{\mathbf{x}_i \boldsymbol{\beta}}}{1 + e^{\mathbf{x}_i \boldsymbol{\beta}}} \right)^{y_i} \\ \mathcal{L}(\boldsymbol{\beta}) &= \prod_{i=1}^n \left(\frac{2e^{\mathbf{x}_i \boldsymbol{\beta}} - (1 + e^{\mathbf{x}_i \boldsymbol{\beta}})}{1 + e^{\mathbf{x}_i \boldsymbol{\beta}}} \right)^{1-y_i} 2^{y_i} \left(\frac{1}{1 + e^{\mathbf{x}_i \boldsymbol{\beta}}} \right)^{y_i} \\ &= \prod_{i=1}^n \left(\frac{e^{\mathbf{x}_i \boldsymbol{\beta}} - 1}{1 + e^{\mathbf{x}_i \boldsymbol{\beta}}} \right)^{1-y_i} 2^{y_i} \left(\frac{1}{1 + e^{\mathbf{x}_i \boldsymbol{\beta}}} \right)^{y_i} \\ \log(\mathcal{L}(\boldsymbol{\beta})) &= \sum_{i=1}^n (1 - y_i) \log \left(\frac{e^{\mathbf{x}_i \boldsymbol{\beta}} - 1}{1 + e^{\mathbf{x}_i \boldsymbol{\beta}}} \right) + y_i \log(1 + e^{-y_i \mathbf{x}_i \boldsymbol{\beta}}) \\ &= \sum_{i=1}^n \log(1 + e^{-y_i \mathbf{x}_i \boldsymbol{\beta}}) \end{aligned}$$

Next we normalize, add in an L_2 norm, and minimize the negative loglikelihood to obtain our objective function.

$$F(\boldsymbol{\beta}) = \min_{\boldsymbol{\beta} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i \mathbf{x}_i \boldsymbol{\beta}}) + \lambda \|\boldsymbol{\beta}\|^2 \quad (3)$$

The objective function implemented in Python is fairly straight forward.

```

def __objective(self, beta):
    """calculate objective value

    Args:
        beta: dX1 (ndarray) weight coefficients

    Returns:
        float
    """
    x, y, n, l = self.__x, self.__y, self.__n, self._lamda

    loss = np.sum([np.log(1 + np.exp(-yi*xi.T@beta)) for xi, yi in zip(x, y)])
    return loss/n + l*np.linalg.norm(beta)**2

```

Gradient descent. In order to optimize the loss function, I choose to use gradient descent; reversing the derivative to descend toward the global minimum. To do this, I take advantage of the linearity of differentiation and begin by assuming both n and d are equal to one. Then, the loss function becomes:

$$\begin{aligned}
 F(\beta) &= \frac{1}{1} \sum_{i=1}^1 \log(1 + e^{-y_1 \mathbf{x}_1 \beta}) + \lambda \|\beta\|^2 \\
 &= \log(1 + e^{-y_1 \mathbf{x}_1 \beta}) + \lambda \|\beta\|^2 \\
 \frac{d}{d\beta} F(\beta) &= \frac{1}{(1 + e^{-y_1 \mathbf{x}_1 \beta})} \log(1 + e^{-y_1 \mathbf{x}_1 \beta}) + 2\lambda\beta \\
 &= \frac{-y_1 \mathbf{x}_1 e^{-y_1 \mathbf{x}_1 \beta}}{1 + e^{-y_1 \mathbf{x}_1 \beta}} + 2\lambda\beta
 \end{aligned}$$

Now, we assume n and d are greater than 1, and using equation 1 (as p) in line 2

$$\begin{aligned}
 &= \frac{1}{n} \left[\frac{-y_1 \mathbf{x}_1 e^{-y_1 \mathbf{x}_1 \beta}}{1 + e^{-y_1 \mathbf{x}_1 \beta}} + \dots + \frac{-y_n \mathbf{x}_n e^{-y_n \mathbf{x}_n \beta}}{1 + e^{-y_n \mathbf{x}_n \beta}} \right] + 2\lambda\beta \\
 &= -\frac{1}{n} \left[-y_1 \mathbf{x}_1 p_1 + \dots + -y_n \mathbf{x}_n p_n \right] + 2\lambda\beta \\
 &= -\frac{1}{n} \sum_{i=1}^n y_i \mathbf{x}_i p_i + 2\lambda\beta \\
 \frac{d}{d\beta} F(\beta) &= -\frac{1}{n} Y^T P X + 2\lambda\beta \quad \text{where} \quad X \in \mathbb{R}^{n \times d}, P \in \mathbb{R}^{n \times n}, Y \in \mathbb{R}^{1 \times n}
 \end{aligned} \tag{4}$$

Again, writing equation in Python is incredibly straight forward.

```

def __gradient(self, b):
    """calculate gradient

    Args:
        b: dX1 (ndarray) weight coefficients

    Returns:
        dX1 ndarray gradient

```

```

"""
x, y, l, n = self.__x, self.__y, self._lamda, self.__n

p = (1 + np.exp(y*(x @ b)))** -1
return 2*l*b - (x.T @ np.diag(p) @ y)/n

```

Optimization is performed iteratively using the fast-gradient descent algorithm and back-tracking line-search (Armijo stopping condition) for finding the optimal learning rate. Proofs of these algorithms are not including in this document but, the Python code is shown below. It can be reused for any optimization technique that allows for gradient descent.

```

def fit(self, x_train, y_train, pos=1, neg=-1, eta=None, queue=None):
    """fit the classifier

    Args:
        x_train: nXd (ndarray) of training samples
        y_train: nX1 (ndarray) of true labels
        pos: (object) positive class label
        neg: (object) optional negative class label
        eta: (float) optional learning rate
        queue: (Queue) optional logging queue
    Returns:
        trained classifier
    """
    self._pos, self._neg = pos, neg

    self.__x, self.__y = x_train, y_train
    self.__n, self.__d = x_train.shape
    self._eta = self.__calc_t_init() if eta is not None else eta

    self.__log_queue = queue
    self.__fgrad()

    self.__log(dict(
        klass=self._pos,
        iter='END',
        obj=self.__objective(self.coef),
        coef=self.coef)
    )
    return self

def __backtracking(self, beta, t_eta=0.5, alpha=0.5):
    """backtracking line search

    Args:
        beta: dX1 (ndarray) weight coefficients
        t_eta: (float) optional [default=0.5] 0 < t_eta < 1 learning rate for eta
        alpha: (float) optional [default=0.5] 0 < alpha < 1 tune stopping condition

    Returns:
        float: optimum learning rate
    """
    l, t = self._lamda, 1

    gb = self.__gradient(beta)

```

```

n_gb = np.linalg.norm(gb)

found_t, i = False, 0
while not found_t and i < 100:
    if self.__objective(beta - t*gb) < self.__objective(beta) - alpha*t*n_gb**2:
        found_t = True
    elif i == self._max_iter-1:
        break
    else:
        t *= t_eta
        i += 1

self.__eta = t
return self.__eta

```

Prediction. We can take advantage of the following proof to simplify the prediction calculations. We begin by showing that $P(Y = 1|\beta) > P(Y = -1|\xi\beta) \iff \xi^T \beta^* > 0$

$$\begin{aligned}
 P(Y = 1|\mathbf{x}_i\beta) &= \frac{e^{\mathbf{x}_i^T \beta}}{1 + e^{\mathbf{x}_i^T \beta}} \\
 &= \frac{1}{1 + e^{\mathbf{x}_i^T \beta}}
 \end{aligned}$$

We need to determine when the ratio of probabilities is > 1 .

$$\frac{P(Y = 1|\mathbf{x}_i\beta)}{P(Y = -1|\mathbf{x}_i\beta)} = \frac{e^{\mathbf{x}_i^T \beta}}{1 + e^{\mathbf{x}_i^T \beta}} (1 + e^{\mathbf{x}_i^T \beta})$$

$$\text{therefore } e^{\mathbf{x}_i^T \beta} > 0 \iff \mathbf{x}_i^T \beta > 0$$

Which gives the following Python snippet for classifying a new sample.

```

def predict(self, x, beta=None):
    """prediction probabilities

    Args:
        x: nXd (ndarray) of input values
        beta: dX1 (ndarray) optional weight coefficients
    Returns:
        nX1 ndarray of class labels
    """
    beta = self.coef if beta is None else beta
    return [self._pos if xi @ beta > 0 else self._neg for xi in x]

```
