

## Laboratory 8

In this laboratory we will focus on Logistic Regression models for classification.

### Numerical optimization

The Logistic Regression model is obtained by minimizing the average cross-entropy between the model predictions and the observed labels. As we have seen, this corresponds also to a Maximum Likelihood solution for the observed labels. While for Gaussian models closed form expressions are available for the ML solutions, this is not the case for Logistic Regression. Therefore, we turn to numerical optimization to find the maximizer of the class likelihoods, or, equivalently, the minimizer of the average cross-entropy.

Numerical optimization algorithms look for the minima of a function  $f(\mathbf{x})$  with respect to the argument  $\mathbf{x}$ . A simple, iterative method to find a local minimum of  $f$  is gradient descent (GD). Given a point  $\mathbf{x}_t$ , gradient descent looks for a descent direction of the function. The direction is given by the negative of the gradient of  $f$ . The algorithm then moves a step  $\alpha_t$  from  $\mathbf{x}_t$  along the descent direction:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \alpha_t \nabla_{\mathbf{x}} f(\mathbf{x})$$

Under mild assumptions on  $\alpha_t$  (e.g.  $\alpha_t \rightarrow 0$ ,  $\sum_{t=1}^{\infty} \alpha_t \rightarrow \infty$ ) the algorithm converges to a local minimum of  $f$ .

A drawback of gradient descent is that it can be quite slow. Faster convergence can be obtained by considering second-order information, such as the Hessian of the function. In this laboratory we will use the L-BFGS algorithm. L-BFGS builds an incremental approximation of the Hessian, that is used to identify a search direction  $\mathbf{p}_t$  at each iteration. The algorithm then proceeds at finding an acceptable step size  $\alpha_t$  for the search direction  $\mathbf{p}_t$ , and uses the direction and step size to update the solution.

The algorithm is implemented in `scipy` (requires importing `scipy.optimize`). We will use the `scipy.optimize.fmin_l_bfgs_b` interface to the numerical solver.

`scipy.optimize.fmin_l_bfgs_b` requires at least 2 arguments (check the documentation for more details):

- **func**: the function we want to minimize.
- **x0**: the starting value for the algorithm.

The L-BFGS algorithm requires computing the objective function and its gradient. To pass the gradient we have different options:

- Through **func**: **func** should return a tuple  $(f(\mathbf{x}), \nabla_{\mathbf{x}} f(\mathbf{x}))$
- Through the optional parameter **fprime**: **fprime** is a function computing the gradient. In this case, **func** should only return the objective value  $f(\mathbf{x})$
- Let the implementation compute an approximated gradient: pass **approx\_grad = True**. Also in this case, **func** should only return the objective value  $f(\mathbf{x})$

The last option does not require writing a function that computes the gradient, as an approximation of the gradient is automatically obtained through finite differences. While this has the advantage that we do not need to derive and implement the gradient, it has two drawbacks:

- The gradient computed through finite differences may not be accurate enough
- The computations are much more expensive, since we need to evaluate the objective function a number of times at least  $D$ , where  $D$  is the size of  $\mathbf{x}$ , at each iteration, and if we want a more accurate approximation of the gradient we may need to evaluate  $f$  many more times

For example, a way to compute a numerical approximation of the gradient consists in computing

$$\frac{\partial f(\mathbf{x})}{\partial x_i} \approx \frac{f(\mathbf{x} + \epsilon \mathbf{e}_i) - f(\mathbf{x} - \epsilon \mathbf{e}_i)}{2\epsilon}$$

where  $\mathbf{e}_i$  is a vectors of zeros, except the element in position  $i$ , which is one:  $\mathbf{e}_1 = [1, 0, 0 \dots 0]$ ,  $\mathbf{e}_2 = [0, 1, 0 \dots 0] \dots \mathbf{e}_D = [0, 0, 0 \dots 1]$ , and  $\epsilon$  is a small value, e.g.  $\epsilon = 10^{-7}$ . This requires computing  $f$  a number of times equal to  $2D$ .

Using the numerical solver, find the minimum of

$$f(y, z) = (y + 3)^2 + \sin(y) + (z + 1)^2$$

The function is convex, so it has a unique minimum.

*STEP 1:* Implement  $f$ . The sin function can be computed using `numpy.sin`.  $\mathbf{f}$  should accept a 1-D numpy array  $\mathbf{x}$  of shape  $(2,)$ . The first component corresponds to variable  $y$ , while the second corresponds to variable  $z$ . The function  $\mathbf{f}$  should return the value  $f(y, z)$ .

*STEP 2:* Call the numerical optimization function `scipy.optimize.fmin_l_bfgs_b`. Pass to the function the previously implemented  $\mathbf{f}$ , and `approx_grad = True`. As starting point you can use values  $[0, 0]$  (pass a numpy array, not a list). If you pass the optional argument `iprint = 1` you can visualize the iterations of the algorithm.

`scipy.optimize.fmin_l_bfgs_b` returns a tuple with three values  $\mathbf{x}$ ,  $\mathbf{f}$ ,  $\mathbf{d}$ :

- $\mathbf{x}$  is the estimated position of the minimum
- $\mathbf{f}$  is the objective value at the minimum
- $\mathbf{d}$  contains additional information (check the documentation)

You should find the minimum at  $[-2.57747138, -0.99999927]$ , with value (truncated)  $-0.356143012$ .

We can also try providing an explicit gradient:

$$\frac{\partial f(y, z)}{\partial y} = 2(y + 3) + \cos(y), \quad \frac{\partial f(y, z)}{\partial z} = 2(z + 1)$$

Rewrite function  $\mathbf{f}$  so that it returns  $f(x)$  as well the gradient of  $f$  as a numpy array with shape  $(2,)$ . Call again the solver, but do not pass `approx_grad`. You should obtain  $\mathbf{x} = [-2.57747137, -0.99999927]$ . In this case the numerical approximation was good enough. However, check the values of the third returned value  $\mathbf{d}$  in the two cases. `'funcalls'` provides the number of times  $\mathbf{f}$  was called. The numerical approximation of the gradient is significantly more expensive, and the cost becomes relatively worse when the dimensionality of the domain of  $f$  increases.

## Binary logistic regression

We can now turn our attention to Logistic Regression. In this section we will implement the binary version of the logistic regression to discriminate between iris virginica and iris versicolor. We will ignore iris setosa. We will represent labels with 1 (iris versicolor) and 0 (iris virginica).

You can load the filtered data with

```
def load_iris_binary():
    D, L = sklearn.datasets.load_iris()['data'].T, sklearn.datasets.load_iris(
        ['target'])
    D = D[:, L != 0] # We remove setosa from D
    L = L[L!=0] # We remove setosa from L
    L[L==2] = 0 # We assign label 0 to virginica (was label 2)
    return D, L

D, L = load_iris_binary()
(DTR, LTR), (DTE, LTE) = split_db_2to1(D, L)
```

Function `split_db_2tol` was defined in Laboratory 5.

The regularized Logistic Regression objective can be written in different ways (we adopt the average-risk expression which divides the loss by the number of samples):

$$J(\mathbf{w}, b) = \frac{\lambda}{2} \|\mathbf{w}\|^2 - \frac{1}{n} \sum_{i=1}^n [c_i \log \sigma(\mathbf{w}^T \mathbf{x}_i + b) + (1 - c_i) \log (1 - \sigma(\mathbf{w}^T \mathbf{x}_i + b))] \quad (1)$$

$$J(\mathbf{w}, b) = \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \log \left( 1 + e^{-z_i(\mathbf{w}^T \mathbf{x}_i + b)} \right), \quad z_i = \begin{cases} 1 & \text{if } c_i = 1 \\ -1 & \text{if } c_i = 0 \end{cases} \quad (\text{i.e. } z_i = 2c_i - 1) \quad (2)$$

$$J(\mathbf{w}, b) = \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \left[ c_i \log \left( 1 + e^{-\mathbf{w}^T \mathbf{x}_i - b} \right) + (1 - c_i) \log \left( 1 + e^{\mathbf{w}^T \mathbf{x}_i + b} \right) \right] \quad (3)$$

Note that (3) follows either from (2), observing that  $c_i = 1$  for  $z_i = 1$  and  $c_i = 0$  for  $z_i = -1$ , or from (1), observing that

$$\log \sigma(\mathbf{w}^T \mathbf{x}_i + b) = -\log \left( 1 + e^{-\mathbf{w}^T \mathbf{x}_i - b} \right)$$

and

$$\log (1 - \sigma(\mathbf{w}^T \mathbf{x}_i + b)) = -\log \sigma(-(\mathbf{w}^T \mathbf{x}_i + b)) = -\log \left( 1 + e^{\mathbf{w}^T \mathbf{x}_i + b} \right)$$

The reason for preferring (2) or (3) to (1) is due to numerical issues that may arise when explicitly computing sigmoids followed by natural logarithms (see below).

Implement Logistic regression using expression (2). You need to write a function `logreg_obj` that, given  $\mathbf{w}$  and  $b$ , allows computing  $J(\mathbf{w}, b)$ . You can then provide this function to the numerical solver to obtain the minimizer of  $J$ .

#### NOTES:

- Function `logreg_obj` should receive a single numpy array  $\mathbf{v}$  with shape  $(D+1,)$ , where  $D$  is the dimensionality of the feature space (e.g.  $D = 4$  for IRIS).  $\mathbf{v}$  should pack all model parameters, i.e.  $\mathbf{v} = [\mathbf{w}, b]$ . Inside the function you can then unpack the array e.g.  $\mathbf{w}, b = \mathbf{v}[0:-1], \mathbf{v}[-1]$
- The function `logreg_obj` has to access also `DTR`, `LTR` and  $\lambda$ , which are required to compute the objective. You can address this in different ways (choose one, the first is the easiest to implement, the other two options offer more insights on the language):

- Write function `logreg_obj` so that it accepts additional arguments `logreg_obj(v, DTR, LTR, l)`. Pass the additional arguments when calling `scipy.optimize.fmin_l_bfgs_b`. This can be achieved by passing `args=(DTR, LTR, l)` to `fmin_l_bfgs_b` (check the documentation of `scipy.optimize.fmin_l_bfgs_b`)

- Write a function `logreg_obj_wrap` that accepts as input `DTR`, `LTR` and  $\lambda$ . Inside the function, define `logreg_obj` as before. `logreg_obj` has now access to the scope of the enclosing function `logreg_obj_wrap`. Make `logreg_obj_wrap` return the created function.

```
def logreg_obj_wrap(DTR, LTR, l):
    def logreg_obj(v):
        # ...
        # Compute and return the objective function value using DTR,
        # LTR, l
        # ...
    return logreg_obj
```

```
# in the main portion, after loading the data:
logreg_obj = logreg_obj_wrap(DTR, LTR, l)
```

- `logreg_obj_wrap` can also be any callable object that accepts a single parameter. You can use an instance (object) of a class that has a single method `logreg_obj(self, v)`, and store in the object the values you need to access. These can be passed to the method that initializes the object

```

class logRegClass:
    def __init__(self, DTR, LTR, l):
        self.DTR = DTR
        self.LTR = LTR
        self.l = l
    def logreg_obj(self, v):
        # Compute and return the objective function value. You can
        # retrieve all required information from self.DTR, self.LTR,
        # self.l

# in the main portion, after loading the data, instantiate a new object
logRegObj = logRegClass(DTR, LTR, l)
# You can now use logRegObj.logreg_obj as objective function:
scipy.optimize.fmin_l_bfgs_b(logRegObj.logreg_obj, ...)

```

- The computation of  $\log(1 + e^{-z_i(\mathbf{w}^T \mathbf{x}_i + b)})$  can lead to numerical issues when  $z_i(\mathbf{w}^T \mathbf{x}_i + b)$  is large, since the sum will make the contribution of the exponential term disappear. We can avoid the issue by using the `numpy.logaddexp` function, which computes

$$\text{numpy.logaddexp}(a, b) = \log(e^a + e^b) .$$

In our example, we need to compute `numpy.logaddexp(0,  $-z_i(\mathbf{w}^T \mathbf{x}_i + b)$ )`.

- $\lambda$  is a hyper-parameter. As usual, we should employ a validation set to estimate good values of  $\lambda$ . For this laboratory, we can simply try different values and see how this affects the performance
- The starting point does not significantly influence the result, since the objective function is convex (there may be slight differences, but should be very small). You can use as initial value an array of all zeros `x0 = numpy.zeros(DTR.shape[0] + 1)`
- The `scipy` implementation of L-BFGS calls the objective function a maximum of 15000 times, and the algorithm stops when this threshold is reached. You can specify a larger amount for the maximum number of calls through the `maxfun` argument
- You can also control the maximum number of allowed iterations through the argument `maxiter`
- You can control the precision of the L-BFGS solution through the parameter `factr`. The default value is `factr=10000000.0`. Lower values result in more precise solutions (i.e. closer to the optimal solution), but require more iterations. Below the default value was used.

Once you have trained the model, you can compute posterior log-likelihood ratios by simply computing, for each test sample  $\mathbf{x}_t$ , the score

$$s(\mathbf{x}_t) = \mathbf{w}^T \mathbf{x}_t + b$$

Compute the array of scores **S**. You can then compute class assignments by thresholding the scores with 0 (i.e. `S[i] > 0  $\implies$  LP[i] = 1`, where **LP** is the array of predicted labels for the test samples).

To check that you implemented the algorithm correctly, you can find below values of the objective function and error rate (1 - accuracy) for different values of  $\lambda$ .

	$J(\mathbf{w}^*, b^*)$	Error rate
$\lambda = 10^{-6}$	<b>7.54150E-3</b>	11.8%
$\lambda = 10^{-3}$	<b>0.110001</b>	8.8%
$\lambda = 10^{-1}$	<b>0.453941</b>	11.8%
$\lambda = 1.0$	<b>0.631644</b>	14.7%

## Multiclass logistic regression (optional)

Implement the multiclass version of logistic regression. The objective function to minimize is

$$J(\mathbf{W}, \mathbf{b}) = \frac{\lambda}{2} \|\mathbf{W}\|^2 - \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K \mathbf{z}_{ik} \log \mathbf{y}_{ik}$$

where

$$\mathbf{y}_{ik} = \frac{e^{\mathbf{w}_k^T \mathbf{x}_i + b_k}}{\sum_j e^{\mathbf{w}_j^T \mathbf{x}_i + b_j}}$$

and

$$\mathbf{W} = [\mathbf{w}_1 \dots \mathbf{w}_K] \quad , \quad \mathbf{b} = \begin{bmatrix} b_1 \\ \dots \\ b_K \end{bmatrix}$$

**NOTE:**  $\mathbf{W}$  is a  $D \times K$  matrix, where  $D$  is the dimensionality of the features space (i.e. dimensionality of  $\mathbf{x}_i$ ). The python function should accept a 1-D numpy array. Represent  $\mathbf{W}$  as a 1-D numpy array of shape  $(D \cdot K,)$ , and reshape it only when performing computations (i.e. inside the function that computes  $J$  and when computing predictions).

*Suggestion:* to avoid numerical issues work directly with  $\log \mathbf{y}_{ik}$ :

$$\log \mathbf{y}_{ik} = \mathbf{w}_k^T \mathbf{x}_i + b_k - \log \sum_{j=1}^n e^{\mathbf{w}_j^T \mathbf{x}_i + b_j}$$

- 1) Compute the matrix of scores  $\mathbf{S}$ :  $\mathbf{S}_{ki}$  should be equal to  $\mathbf{S}_{ki} = \mathbf{w}_k^T \mathbf{x}_i + b_k$ , i.e. each column of  $\mathbf{S}$  contains the scores for sample  $\mathbf{x}_i$  for all classes. You can compute the score matrix with a single product, exploiting broadcasting:  $\mathbf{S} = \text{numpy.dot}(\mathbf{W}, \mathbf{T}) + \mathbf{b}$
- 2) Compute matrix  $\mathbf{Y}^{log}$  containing  $\mathbf{Y}_{ki}^{log} = \log \mathbf{y}_{ik}$  (note that the indices are swapped: in  $\mathbf{Y}^{log}$  the first index represents the class, the second the sample).  $\mathbf{Y}_{ki}^{log}$  can be computed from  $\mathbf{S}$ . Each row of  $\mathbf{Y}_{log}$  corresponds to the same row of  $\mathbf{S}$  minus the expression  $\log \sum_{j=1}^n e^{\mathbf{w}_j^T \mathbf{x}_i + b_j}$ . The last expression is the log-sum-exp of the rows of  $\mathbf{S}$ .
- 3) Use the 1-of-K encoding of the labels: matrix  $\mathbf{T}$  should contain the labels, encoded as  $\mathbf{T}_{ki} = 1 \iff c_i = k$ . The other elements should be 0 (again, the indices are swapped with respect to  $\mathbf{z}_{ik}$ , the first index of  $\mathbf{T}$  represents the class).
- 4) The summation in  $J(\mathbf{w}, \mathbf{b})$  can be computed by element-wise multiplication of  $\mathbf{T}$  and  $\mathbf{Y}_{log}$ , followed by the summation of all elements of the resulting matrix
- 5) The squared norm of  $\mathbf{W}$  corresponds to  $\sum_i \sum_j (\mathbf{W}_{ij})^2$ , i.e.  $(\mathbf{W}^* \mathbf{W}).\text{sum}()$

Train the model using the data that we used for the Gaussian classifier:

```
D, L = load_iris()
(DTR, LTR), (DTE, LTE) = split_db_2to1(D, L)
```

To test the model, compute class posterior probabilities as  $P(C = c | \mathbf{x}_t) = \mathbf{w}_c^T \mathbf{x}_t + b_c$ . Predict the class with highest posterior probability (we assume uniform mis-classification costs).

The following table contains the training loss and the test error rate for different values of  $\lambda$

	$J(\mathbf{w}^*, \mathbf{b}^*)$	Error rate
$\lambda = 10^{-6}$	3.96791E-2	4.0%
$\lambda = 10^{-3}$	9.69097E-2	4.0%
$\lambda = 10^{-1}$	0.500591	6.0%
$\lambda = 1.0$	0.821155	18.0%