

## Laboratory 8

DOVERRE ESSERE UN'ALTERNATIVA ALLA  
DISTRIBUZIONE GAUSSIANA

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.

DOBBIAMO  
USARE  
L'OTTIMIZZAZIONE  
NUMERICA PER  
MINIMIZZARE LA  
PROB. DI  
CLASS. 0  
MINIMIZZARE  
LA MEDIA DELLA  
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.

OPZIONALMENTE GRADIENTE E VINCOLI SU VAR

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

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

**L-BFGS** È LIMITED CIOÈ UTILIZZA UNA QUANTITÀ LIMITATA DI DATI. INVECE DI MEMORIZZARE L'INTERA **MATRICE HESSIANA** (MATRICE DI DERIVATE SECONDE) QUESTO MEMORIZZA SOLO UN CERTO NUMERO DI VETTORI CHE RAPPRESENTA SOLO UN CERTO NUMERO DI VETTORI CHE RAPPRESENTANO APPROSSIMAZIONE DI MATRICE HESSIANA. QUESTO LO RENDE PIÙ EFFICIENTE.

MATRICE DI DERIVATE SECONDE

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

VECTOR  $[0, 1, 0, \dots]$

SMALL VALUE  $\epsilon = 10^{-7}$

$\begin{cases} e_1 = [1, 0, 0, \dots] \\ e_2 = [0, 1, 0, \dots] \end{cases}$

where  $e_i$  is a vector of zeros, except the element in position  $i$ , which is one:  $e_1 = [1, 0, 0 \dots 0]$ ,  $e_2 = [0, 1, 0 \dots 0]$  ...  $e_D = [0, 0, 0 \dots 1]$ , and  $\epsilon$  is a small value, e.g.  $\epsilon = 10^{-7}$ . This requires computing  $f$  a number of times equal to  $2D$ .

*esempio*

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`.  $f$  should accept a 1-D numpy array  $x$  of shape  $(2,)$ . The first component corresponds to variable  $y$ , while the second corresponds to variable  $z$ . The function  $f$  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  $f$ , 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  $x$ ,  $f$ ,  $d$ :

- $x$  is the estimated position of the minimum
- $f$  is the objective value at the minimum
- $d$  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  $f$  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  $x = [-2.57747137, -0.99999927]$ . In this case the numerical approximation was good enough. However, check the values of the third returned value  $d$  in the two cases. 'funcalls' provides the number of times  $f$  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.

*SENZA APPROX-GRAD*

*PER QUESTO MOTIVO NON METTI APPROX-GRAD A TWE*

## 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[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), (DVAL, LVAL) = split_db_2to1(D, L)
```

LA LOGISTIC REGRESSION PUÒ SOFFRIRE DI OVERFITTING QUINDI LA REGULARIZATION AGGIUNGE UN TERMINE DI PENALTA ALLA FUNZIONE DI COSTO CHE IL MODELLO CERCA DI MINIMIZZARE IMPONENDO RESTRIZIONI SUI COEFF

$$J(\mathbf{w}, b) = \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-z_i(\mathbf{w}^T \mathbf{x}_i + b)}) ,$$

$\mathbf{w} \rightarrow$   
peso  $\rightarrow$  è un vettore di dimensione  $n$  (dove  $n$  è il numero delle caratteristiche) che moltiplica ciascuna delle caratteristiche per determinare la prob di appartenenza di una classe

$b \rightarrow$   
bias  $\rightarrow$  è un termine di somma serve per spostare la funzione di attivazione anche quando tutte le caratteristiche sono zero

Function `split_db_2to1` 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):

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

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

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

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

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

and

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

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). This requires minimizing expression (2). For this, you need to write a function `logreg_obj` that, given `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:  $\mathbf{w}$  È ASSOCIATO AD OGNI FEATURES

- Function `logreg_obj` should receive a single numpy array `v` with shape  $(D+1,)$ , where  $D$  is the dimensionality of the feature space (e.g.  $D = 4$  for IRIS). `v` should pack all model parameters, i.e. `v = [w, b]`. Inside the function you can then unpack the array e.g. `w, b = v[0:-1], v[-1]`
- The function `logreg_obj` needs to access also `DTR`, `LTR` and `l`, which are required to compute the objective. You can address this in different, alternative ways:

– Embed the objective function and its optimization inside a function `trainLogReg` that receives both `DTR`, `LTR` and a value for `l`. The objective function is defined inside `trainLogReg`, thus it can access the variables `DTR` and `LTR` in the outer function scope (i.e., those passed to `trainLogReg`. The code looks like:

```
def trainLogReg(DTR, LTR, l):

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

    # Here we find the minimizer of logreg_obj
    xf = scipy.optimize.fmin_l_bfgs_b(func = logreg_obj, x0 = numpy.
        zeros(DTR.shape[0]+1), approx_grad=True)[0]
    return xf
```

- Write a function `logreg_obj` that accepts additional arguments `logreg_obj(v, DTR, LTR, l)`. Optimize the function by calling `scipy.optimize.fmin_l_bfgs_b`. You can pass the additional arguments to the objective function through the `args` parameter of `fmin_l_bfgs_b` (check the documentation of `scipy.optimize.fmin_l_bfgs_b`):

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

# in the main portion, after loading the data:
scipy.optimize.fmin_l_bfgs_b(func = logreg_obj, args=(DTR, LTR, l),
    ...)
```

- 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)
scipy.optimize.fmin_l_bfgs_b(func = logreg_obj, ...)[0]
```

- `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(func = 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, - $z_i(\mathbf{w}^T \mathbf{x}_i + b)$ )`.

- Broadcasting can significantly speed-up the computations. You can compute a vector of “scores”  $\mathbf{S}$

$$\mathbf{S} = [(\mathbf{w}^T \mathbf{x}_1 + b) \dots (\mathbf{w}^T \mathbf{x}_n + b)]$$

using simple matrix-vector multiplication: `S = (vcol(w).T @ DTR + b).ravel()`. Remember to reshape the result to a 1-D array (`.ravel()`). You can then multiply each element of the resulting vector by the corresponding label: `-ZTR * S`, where `ZTR = 2 * LTR - 1`. `numpy.logaddexp`

supports broadcasting so you can compute in a single shot all terms

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

with `numpy.logaddexp(0, -ZTR * S)`

- To speed-up computations, it's useful to provide also the gradient of the function. We can express in vector form the derivatives with respect to the components of  $\mathbf{w}$  and  $b$  as (notice that there is no minus in the exponential term  $z_i(\mathbf{w}^T \mathbf{x}_i + b)$ ):

$$\nabla_{\mathbf{w}} J = \left[ \frac{\partial J}{\partial w_1} \dots \frac{\partial J}{\partial w_d} \right] = \lambda \mathbf{w} + \frac{1}{n} \sum_{i=1}^n \frac{-z_i}{1 + e^{z_i(\mathbf{w}^T \mathbf{x}_i + b)}} \mathbf{x}_i$$

$$\frac{\partial J}{\partial b} = \frac{1}{n} \sum_{i=1}^n \frac{-z_i}{1 + e^{z_i(\mathbf{w}^T \mathbf{x}_i + b)}}$$

where  $d$  is the dimensionality of the feature vectors (`DTR.shape[0]`). These terms can be efficiently computed by computing the vector  $\mathbf{G}$ :

$$\mathbf{G} = \left[ \frac{-z_1}{1 + e^{z_1(\mathbf{w}^T \mathbf{x}_1 + b)}} \dots \frac{-z_n}{1 + e^{z_n(\mathbf{w}^T \mathbf{x}_n + b)}} \right]$$

COME HAI FATTO  $\mathbf{S}$

so that

$$\nabla_{\mathbf{w}} J = \lambda \mathbf{w} + \frac{1}{n} \sum_{i=1}^n \mathbf{G}_i \mathbf{x}_i$$

$$\frac{\partial J}{\partial b} = \frac{1}{n} \sum_{i=1}^n \mathbf{G}_i$$

which can be computed, using broadcasting, from the vector of scores  $\mathbf{S}$  as  $\mathbf{G} = -\text{ZTR} / (1.0 + \text{numpy.exp}(\text{ZTR} * \mathbf{S}))$  (note: you may obtain overflow from the exponentiation due to well-classified samples, however you can safely ignore the issue since it will simply cause the corresponding  $\mathbf{G}_i$  becoming exactly 0 rather than almost 0, and this will not cause numerical issues). You can then compute all terms  $\mathbf{G}_i \mathbf{x}_i$  with broadcasting: `(vrow(G) * DTR)`.

Remember that `scipy.optimize.fmin_l_bfgs_b` requires that you pack  $\nabla_{\mathbf{w}} J$  and  $\frac{\partial J}{\partial b}$  in a single vector  $\mathbf{v}_{grad} = [\nabla_{\mathbf{w}} J, \frac{\partial J}{\partial b}]$ . You have to modify the function `logreg_obj` so that it returns both the objective and the gradient  $\mathbf{v}_{grad}$  (and set `approx_grad = False` when calling the optimizer function).

- $\lambda$  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 ~~2000~~
- 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-probability 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  $\mathbf{S} = \text{vcol}(\mathbf{w}).T @ \text{DVAL} + \text{b}$ . You can then compute class assignments by thresholding the scores with 0 (i.e.  $\mathbf{S}[\mathbf{i}] > 0 \implies \text{LP}[\mathbf{i}] = 1$ , where  $\text{LP}$  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  $\lambda$ .

The logistic regression model provides scores that can be interpreted as log-posterior-ratio for the two classes, with a prior corresponding to the empirical training set prior. To adapt the scores to different applications, however, we need to transform the scores so that they can be interpreted as log-likelihood ratios, otherwise the application-dependent thresholds that we employ for LLRs cannot be used anymore.

To obtain a score that behaves like a LLR we can simply remove the empirical prior log-odds from the score:

$$s_{llr}(\mathbf{x}_t) = \mathbf{w}^T \mathbf{x}_t + b - \log \frac{\pi_{emp}}{1 - \pi_{emp}} = s(\mathbf{x}_t) - \log \frac{\pi_{emp}}{1 - \pi_{emp}} \rightarrow \text{DEVI TOGIERE QUESTO}$$

where  $\pi_{emp}$  is the empirical training set prior (fraction of samples of class 1).

Compute minimum DCF, optimal Bayes predictions and actual DFCs for different applications. Optimal Bayes predictions require comparing  $s_{llr}$  with the standard application-dependent threshold. Below you can find the results for  $\pi_T = 0.5$ .

	$J(\mathbf{w}^*, b^*)$	Error rate ( $\mathbf{w}^T \mathbf{x}_t + b \leq 0$ )	minDCF ( $\pi_T = 0.5$ )	actDCF ( $\pi_T = 0.5$ )
$\lambda = 10^{-3}$	1.100009e-01 ✓	8.8% ✓	0.0625 ✓	0.1181 ✓
$\lambda = 10^{-1}$	4.539407e-01 ✓	11.8% ✓	0.0556 ✓	0.1111 ✓
$\lambda = 1.0$	6.316436e-01 ✓	14.7% ✓	0.1111 ✓	0.1667 ✓

## Prior-weighted logistic regression and calibration

The prior-weighted logistic regression model allows us to simulate different priors for class 1. The objective function becomes

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

where  $\pi_T$  is the target prior, and  $n_T$  and  $n_F$  are the number of samples of class 1 and 0, respectively.

Implement the prior-weighted version of the model, and test it for an application with  $\pi_T = 0.8$  (you can check the results below). In this case, the gradients can be computed as

$$\nabla_{\mathbf{w}} J = \lambda \mathbf{w} + \sum_{i=1}^n \xi_i \frac{-z_i}{1 + e^{-z_i(\mathbf{w}^T \mathbf{x}_i + b)}} \mathbf{x}_i = \lambda \mathbf{w} + \sum_{i=1}^n \xi_i \mathbf{G}_i \mathbf{x}_i$$

$$\frac{\partial J}{\partial b} = \sum_{i=1}^n \xi_i \frac{-z_i}{1 + e^{-z_i(\mathbf{w}^T \mathbf{x}_i + b)}} = \sum_{i=1}^n \xi_i \mathbf{G}_i$$

where  $\mathbf{G}$  is defined as in the standard model

To verify your implementation, you can also compute the solution when  $\pi_T = \pi_{emp}$ . You should obtain the same loss and same output  $(\mathbf{w}, b)$  as for the original, non-weighted model.

Finally, remember that, to obtain scores that behave like LLRs, when using the prior-weighted model we need to remove the prior log-odds of the application (not the empirical prior log-odds), i.e.

$$s_{llr}(\mathbf{x}_t) = \mathbf{w}^T \mathbf{x}_t + b - \log \frac{\pi_T}{1 - \pi_T} = s(\mathbf{x}_t) - \log \frac{\pi_T}{1 - \pi_T}$$

For an application  $\pi_T = 0.8$  you should obtain:

	$J(\mathbf{w}^*, b^*)$	minDCF ( $\pi_T = 0.8$ )	actDCF ( $\pi_T = 0.8$ )
$\lambda = 10^{-3}$	9.401035e-02 ✓	0.1667 ✓	0.2222 ✓
$\lambda = 10^{-1}$	3.606261e-01 ✓	0.0556 ✓	0.7222 ✓
$\lambda = 1.0$	4.724715e-01 ✓	0.1111 ✓	1.0000 ✓



## Multiclass logistic regression (optional)

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

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

where

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

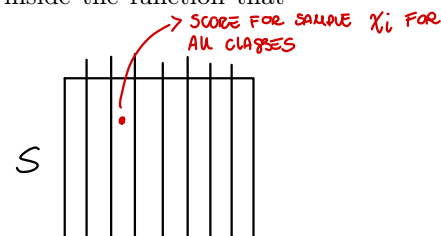
and

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

**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  $\log y_{ik}$ :

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



- 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  $\mathbf{Y}_{ki}^{log} = \log y_{ik}$  (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  $\log \sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}_i + b_j}$ . 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  $z_{ik}$ , 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), (DVAL, LVAL) = 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
$\lambda = 10^{-3}$	9.69097E-2	4.0%
$\lambda = 10^{-1}$	0.500591	6.0%
$\lambda = 1.0$	0.821155	18.0%

# Project

We analyze the binary logistic regression model on the project data. We start considering the standard, non-weighted version of the model, without any pre-processing.

1 Train the model using different values for  $\lambda$ . You can build logarithmic-spaced values for  $\lambda$  using `numpy.logspace`. To obtain good coverage, you can use `numpy.logspace(-4, 2, 13)` (check the documentation). Train the model with each value of  $\lambda$ , score the validation samples and compute the corresponding actual DCF and minimum DCF for the primary application ( $\pi_T = 0.1$ ). To compute actual DCF remember to remove the log-odds of the training set empirical prior. Plot the two metrics as a function of  $\lambda$  (suggestion: use a logarithmic scale for the x-axis of the plot - to change the scale of the x-axis you can use `matplotlib.pyplot.xscale('log', base=10)`). What do you observe? Can you see significant differences for the different values of  $\lambda$ ? How does the regularization coefficient affects the two metrics?

$$s_{LLR}(x_t) = w^T x_t + b - \log \frac{\pi_{emp}}{1 - \pi_{emp}} = s(x_t) - \log \frac{\pi_{emp}}{1 - \pi_{emp}}$$

2 Since we have a large number of samples, regularization seems ineffective, and actually degrades actual DCF since the regularized models tend to lose the probabilistic interpretation of the scores. To better understand the role of regularization, we analyze the results that we would obtain if we had fewer training samples. Repeat the previous analysis, but keep only 1 out of 50 model training samples, e.g. using data matrices `DTR[:, :50]`, `LTR[:, :50]` (apply the filter only on the model training samples, not on the validation samples, i.e., after splitting the dataset in model training and validation sets). What do you observe? Can you explain the results in this case? Remember that lower values of the regularizer imply larger risk of overfitting, while higher values of the regularizer reduce overfitting, but may lead to underfitting and to scores that lose their probabilistic interpretation.

3 In the following we will again consider only the full dataset. Repeat the analysis with the prior-weighted version of the model (remember that, in this case, to transform the scores to LLRs you need to remove the log-odds of the prior that you chose when training the model). Are there significant differences for this task? Are there advantages using the prior-weighted model for our application (remember that the prior-weighted model requires that we know the target prior when we build the model)?

4 Repeat the analysis with the quadratic logistic regression model (again, full dataset only). Expand the features, train and evaluate the models (you can focus on the standard, non prior-weighted model only, as the results you would obtain are similar for the two models), again considering different values for  $\lambda$ . What do you observe? In this case is regularization effective? How does it affect the two metrics?

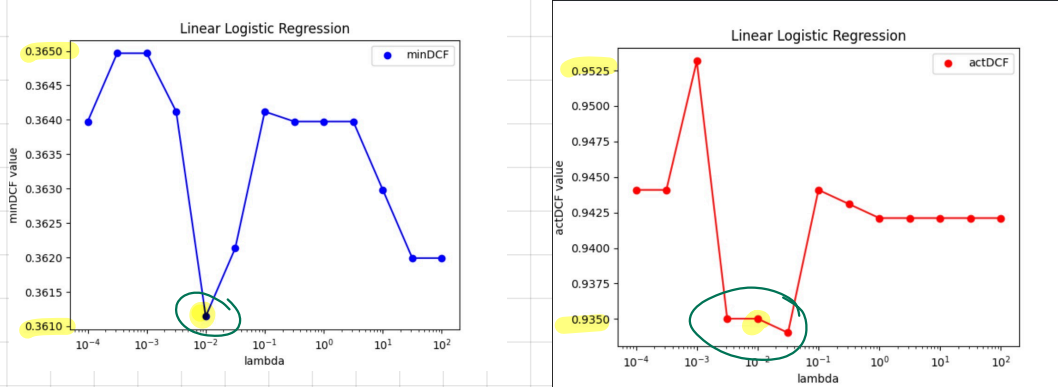
5 The non-regularized model is invariant to affine transformations of the data. However, once we introduce a regularization term affine transformations of the data can lead to different results. Analyze the effects of centering (optionally, you can also try different strategies, including Z-normalization and whitening, as well as PCA) on the model results. You can restrict the analysis to the linear model. Remember that you have to center both datasets with respect to the model training dataset mean, i.e., you must not use the validation data to estimate the pre-processing transformation. For this task, you should observe only minor variations, as the original features were already almost standardized.

6 As you should have observed, the best models in terms of minimum DCF are not necessarily those that provide the best actual DCFs, i.e., they may present significant mis-calibration. We will deal with score calibration at the end of the course. For the moment, we focus on selecting the models that optimize the minimum DCF on our validation set. Compare all models that you have trained up to now, including Gaussian models, in terms of minDCF for the target application ( $\pi_T = 0.1$ ). Which model(s) achieve(s) the best results? What kind of separation rules or distribution assumptions characterize this / these model(s)? How are the results related to the characteristics of the dataset features?

Suggestion for upcoming laboratories: the last laboratories will cover score calibration, and will require to evaluate the results of the models that you tested on an held-out evaluation set. We suggest that you save the models, and the corresponding validation scores as well, since these will be required for score calibration (you can skip the models trained with the reduced dataset, as they won't be needed).

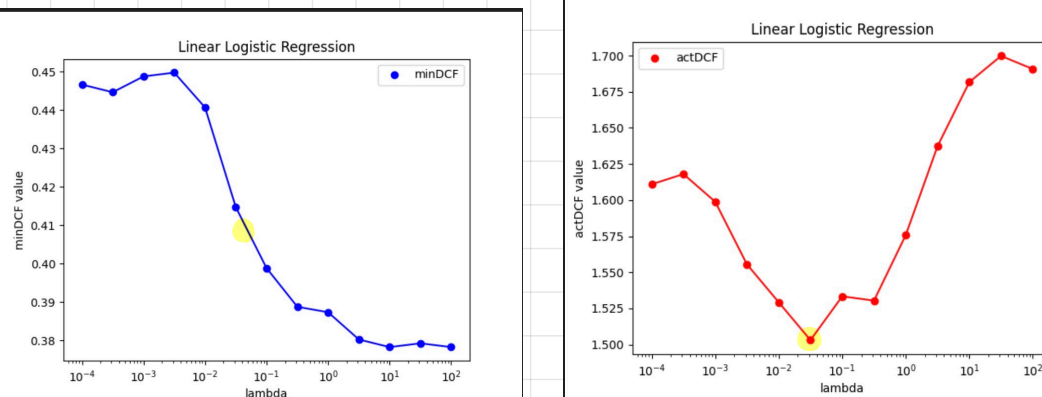
- 1) È normale che i grafici di minDCF e actDCF non siano stabili e che mostrino un comportamento di crescita e diminuzione. Questo perché stai variando il parametro di regolarizzazione lambda (ldb) e osservando come cambiano le prestazioni del modello. Quando il valore di lambda (ldb) è molto piccolo, il modello può adattarsi troppo ai dati di addestramento, portando a un overfitting. Questo può risultare in un errore di validazione più alto, quindi vedrai un aumento nei valori di minDCF e actDCF. Al contrario, quando il valore di lambda (ldb) è molto grande, il modello può diventare troppo semplice e non riuscire a catturare tutte le informazioni nei dati, portando a un underfitting. Anche in questo caso, vedrai un aumento nei valori di minDCF e actDCF. Tra questi due estremi, dovrebbe esserci un valore ottimale di lambda (ldb) che bilancia l'overfitting e l'underfitting, minimizzando sia minDCF che actDCF. Questo è il motivo per cui vedrai i valori di minDCF e actDCF diminuire e poi aumentare di nuovo mentre vari il valore di lambda (ldb).

NO PRIOR WEIGHTED



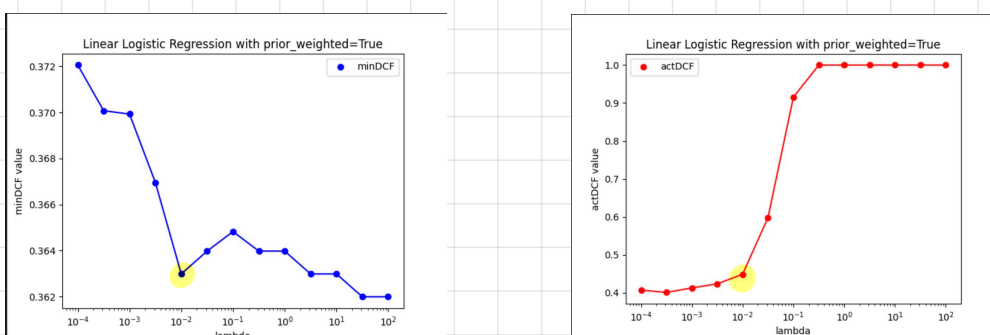
- USANDO SOLO 50 SAMPLE NO HAI VALORI IN GENERALE PIÙ ALTI

2)

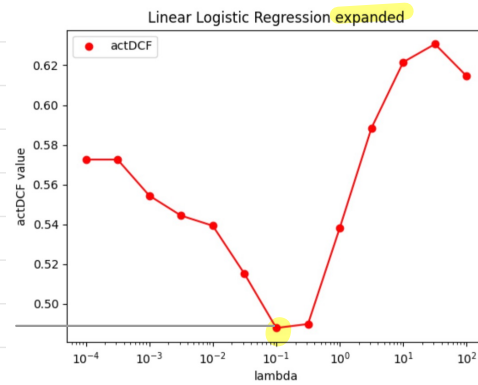
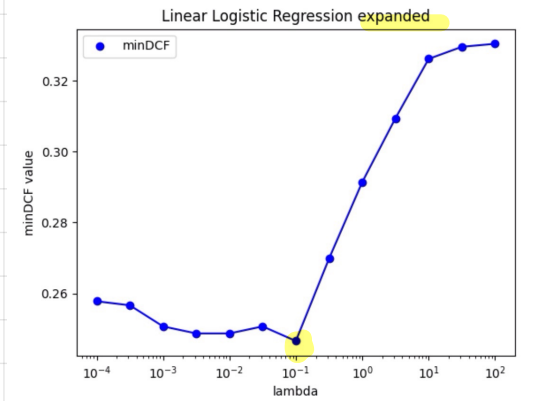


Quando si utilizzano tutti i dati di addestramento, si ha un set di dati più grande e quindi più informazioni per addestrare il modello. In questo caso, sembra che un valore di lambda (ldb) di  $10^{-2}$  sia un buon compromesso tra overfitting e underfitting, minimizzando entrambi i DCF. D'altra parte, quando si utilizza solo una parte dei dati di addestramento (in questo caso, solo 50 dei dati), si ha meno informazioni per addestrare il modello. Questo può rendere più difficile trovare un valore di lambda (ldb) che bilancia efficacemente l'overfitting e l'underfitting. In questo caso, potrebbe non esserci un valore di lambda (ldb) che minimizza entrambi i DCF. Questo suggerisce che la quantità di dati di addestramento può avere un impatto significativo sulla capacità del modello di generalizzare bene ai dati non visti. In generale, più dati di addestramento si hanno, più accurato sarà il modello. Tuttavia, è anche importante bilanciare questo con la complessità del modello (controllata dal parametro lambda ldb) per evitare l'overfitting.

3) WITH PRIOR WEIGHTED



4) PER ESEGUIRE L'ANALISI CON IL MODELLO DI REGRESSIONE LOGISTICA QUADRATA DOVRAI ESPANDERE LE CARATTERISTICHE INCLUDENDO ANCHE I QUADRATI DELLE CARATTERISTICHE. IN QUESTO MODO PUOI CATTURARE ANCHE LE RELAZIONI NON LINEARI TRA CARATTERISTICHE E ETICHETTA DI OUTPUT  
DEVI FARE L'APPEND DI  $\text{np.square}(\text{X})$  e  $\text{DVAL}$    
 ↳ PER COERENZA ANCHE IL DATASET DI VALUTAZIONE DEVE ESSERE ESPANSO



IN QUESTO CASO L'ERRORE EFFETTIVAMENTE È PIÙ BASSO A PARITÀ DI LAMBDA CIO' VUOL DIRE CHE STA PERFORMANDO MEGLIO. SE CAPISCI CHE IN UN MODELLO SEMBRA MEGLIO ANZICHÉ QUELLO POTREBBE ESSERE QUELLO GIUSTO

TUTTE QUESTE VARIANTI CHE SI CONSIDERANO POSSONO ESSERE CONSIDERATI MODELLI.

POTREBBE NON SERVIRE LA 2-NORM POICHÉ I DATI SONO GIÀ NORMALIZZATI

5) GRAFICA PCA È PCA + 2-NORM → NON CAMBIA NULLA

- 6)
- **Modelli migliori:** devi trovare quello con il minDCF più basso, quel modello starà ottenendo i migliori risultati
  - **Analizza regole di separazione e le ipotesi di distribuzione:** Una volta identificati i modelli migliori, dovresti esaminare le regole di separazione o le ipotesi di distribuzione che caratterizzano questi modelli. Ad esempio, se stai utilizzando un modello di regressione logistica, potresti esaminare i coefficienti del modello per capire quali caratteristiche stanno avendo il maggiore impatto sulla previsione.
  - **Correlazioni con le caratteristiche del dataset:** Infine, dovresti cercare di capire come i risultati sono correlati alle caratteristiche del tuo dataset. Ad esempio, se una particolare caratteristica ha un forte impatto sulle previsioni del tuo modello, potrebbe essere perché quella caratteristica è particolarmente rilevante per la variabile target che stai cercando di prevedere. Allo stesso modo, se una caratteristica non ha quasi nessun impatto, potrebbe essere perché quella caratteristica non è rilevante o perché i tuoi dati non contengono abbastanza informazioni su quella caratteristica per renderla utile.