

Assignment 2: k -Fold Cross-Validation

In MATLAB implement a function, which compares two learning algorithms using k -fold cross validation.

Suppose that learning algorithms are functions of the form

$$\mathbf{LPar} = \text{name}(\mathbf{Tr}, \mathbf{DTr}, \mathbf{Par}),$$

where

- name is the name of the learning function,
- \mathbf{Tr} is a training set (training vectors are the columns of \mathbf{Tr}),
- \mathbf{DTr} is a row vector of desired outputs, and
- \mathbf{Par} is a cell array containing parameters of the learning algorithm.

Such function returns learned parameters \mathbf{LPar} (also in a cell array).

Further we will need a function which will execute the learned function

$$\mathbf{Out} = \text{nameL}(\mathbf{LPar}, \mathbf{In}),$$

which computes the learned function with parameters \mathbf{LPar} on input vectors from a matrix \mathbf{In} (each column of the matrix is an input vector). \mathbf{Out} is a row vector of the results.

Then it is possible to implement the function

$$\mathbf{E} = \text{Err}(\text{Name}, \text{NameL}, \mathbf{Par}, \mathbf{Tr}, \mathbf{DTr}, \mathbf{Ts}, \mathbf{DTs}),$$

which computes the error of the function NameL learned by the algorithm Name with parameters \mathbf{Par} on a training set \mathbf{Tr} and desired outputs \mathbf{DTr} . The error is computed on a test set \mathbf{Ts} with the desired outputs \mathbf{DTs} . Note that the result of the function Err is always a real value from the closed interval $\langle 0; 1 \rangle$, *not the number of errors* made by the learned algorithm on the test set. Name and NameL are strings containing names of the respective functions. The respective functions can be evaluated using the function `feval`.

Then it is easy to implement the following function

$$[\mathbf{delta}, \mathbf{s}] = \text{CrossVal}(\text{Name1}, \text{Name1L}, \mathbf{Par1}, \text{Name2}, \text{Name2L}, \mathbf{Par2}, \mathbf{Pat}, \mathbf{DOut}, \mathbf{k}, \text{NoShuffle}),$$

which estimates the difference between the errors of the hypothesis Name1L learned by a learning algorithm Name1 with parameters $\mathbf{Par1}$ and the error of the hypotheses Name2L learned by a learning algorithm Name2 with parameters $\mathbf{Par2}$ using k -fold cross-validation on patterns \mathbf{Pat} with the desired outputs \mathbf{DOut} . The function returns the difference of errors \mathbf{delta} and the standard deviation \mathbf{s} of this estimator.

The last parameter `NoShuffle` can be omitted. If it is present, then the order of samples must not be changed before partitioning into folds for k -fold cross-validation and all folds should be a continuous parts of \mathbf{Pat} . If the parameter is omitted, the patterns from \mathbf{Pat} should be randomly assigned into k folds.

For example, we can compare the errors of the perceptron learning algorithm limited to at most 10 epochs and the perceptron learning algorithm limited to at most 100 epochs. As the perceptron learning algorithm we will use the function `PLearn` obtained from the function `perc_learn` from the seminary and similarly for simulating perceptrons we will use `perc_recall` from the seminary:

```

function LPar = PLearn(x,c,Par)
    p = perc_learn(Par{1},x,c,Par{2},Par{3});
    LPar = {p}
end

function Out = PRecall(Par,x)
    Out = perc_recall(Par{1},x);
end

```

Another learning algorithm we will implement is **Memorizer**, which memorizes all training samples together with the desired outputs for them. Learned **Memorizer** answers correctly on the inputs from the training set and randomly otherwise:

```

function LPar = Memorizer(Tr,DTr,Par)
    % The learning algorithm which only remembers all training
    % samples and the desired answers for them
    %
    %     LPar = Memorizer(Tr,Dtr,Par)
    %
    % inputs:
    %   Tr   matrix with training samples in columns
    %   DTr  row vector of the desired outputs
    %   Par  is not used here, it is present here for
    %         compatibility only
    %
    % output:
    %   LPar a cell array containing the training samples and
    %         the desired outputs for them

    LPar = {Tr,DTr};
end

function Out = MemorizerRecall(LPar,In)
    % A function simulating Memorizer
    %
    %     Out = MemorizerRecall(LPar,In)
    %
    % inputs:
    %   Lpar  a cell array with training samples and
    %         the respective desired outputs
    %   In    a matrix of input vectors (as columns)

    % output:
    %   Out   a row vector; whenever the i-th input vector is
    %         contained within the memorized input samples,
    %         Out(i) equals the i-th remembered desired output,
    %         otherwise it will be randomly 0 or 1

    Known = LPar{1};
    KnownOut = LPar{2};

```

```

j=1;
Out = zeros(1,size(In,2));
for i = 1:size(In,2)
    j=1;
    while j<=size(Known,2)
        if In(:,i)==Known(:,j)
            break
        else
            j=j+1;
        end
    end
    if j<=size(Known,2)
        Out(i) = KnownOut(j);
    else
        Out(i) = rand(1)>0.5;
    end
end
end
end

```

The above algorithms can be compared using the above function `CrossVal`.

```

% at first we read training patterns;
% this should a matrix of size 2x600
In1 = csvread('In1.csv')

% then we read desired outputs for the training patterns
% (0/1 row vector of size 600)
c1 = csvread('c1.csv')

% a column In1(:,i) is the i-th training vector with
% the desired output c(i)

% we prepare two sets (cell arrays) of the learning parameters
% for the perceptron learning algorithm consisting of:
%     an extended weight vector,
%     a learning rate, and
%     a maximal number of epochs
Par1 = {[1 1 -1], 1, 10}
Par2 = {[1 1 -1], 1, 100}

#run 5-fold cross-validation
[d,s] = CrossVal('PLearn','PRecall',Par1,'PLearn','PRecall',
                Par2,In1,c1,5,'NoShuffle')

```

Tasks:

- a) Implement the function `CrossVal` and by running the above script estimate the difference in error rates of the perceptron learning algorithm with at most 10 epochs and of the same perceptron learning algorithm with at most 100 epochs. From the resulting error difference and the standard deviation of the estimate compute interval which contains the true error

difference with the probability 95%. Is the error difference statistically significant?

- b) Then modify the above script to compare error rates of the algorithm **Memorizer** and perceptron with at most 50 epochs, learning rate 1.0 and initial extended weight vector $[1 \ 1 \ 1 \ 1 \ -1]$ using 6-fold cross-validation on the following 300 samples:

```
In2 = csvread('In2.csv');  
  
c = csvread('c2.csv');  
  
% in order to obtain the same answers  
% from MemorizerRecall, we reset random  
% numbers generator  
stream = RandStream.getGlobalStream; reset(stream)
```

You should submit:

1. A Zip-file with commented source code of the function **CrossVal** and all the functions you have used for solving the above tasks.
2. A text file (PDF is the preferred format) describing your solution and containing the results of your experiments. You should analyze the obtained results and explicitly write a *recommendation* which algorithm is better to use. Eventually you can give an advice which experiments would be suitable for more thorough comparison of the considered learning algorithms.