

## Math6373 Final (take home exam) Due date May 6th 2020 at midnight

1 *Prediction Task* : Select one major stock on the US stockmarket. On day "t", let  $S(t)$  be the price of this stock at closing time. On each day "t", we want to *predict* the future stock price  $S(t+1)$  given the last 20 observed stock prices  $S(t)$ ,  $S(t-1)$ ,  $S(t-2)$ , ...,  $S(t-19)$

2 *Data Set* : Let  $t=1, 2, \dots, N$  be the days on which the US stock exchange was open during the time period 2014-2015-2016-2017 . Download the time series  $S(t)$  for  $t= 1, 2, \dots, N$

2 *PreProcessing*: Replace *isolated* missing values  $S(t)$  by the mean of two actual values closest to time  $t$ . If there are too many missing values, download another stock . For  $20 \leq t \leq N-1$  , compute the following three moving averages of the time series  $S$  :

$$MA5(t) = [S(t-4) + S(t-3) + S(t-2) + S(t-1) + S(t)] / 5$$

$$MA10(t) = [S(t-9) + S(t-8) + \dots + S(t)] / 10$$

$$MA20(t) = [S(t-19) + S(t-18) + \dots + S(t)] / 20$$

Plot the 4 curves  $S(t)$  ,  $MA5(t)$ ,  $MA10(t)$ ,  $MA20(t)$ , on the same graph

3 *Training and Test sets for an MLP predictor* :On each day  $t \geq 20$  , the *recent past* of the series  $S$  will be defined as the  $1 \times 18$  line vector

$$V_t = [MA5(t), MA10(t), MA20(t), S(t), S(t-1), S(t-2), \dots, S(t-13), S(t-14)]$$

For  $20 \leq t \leq N-1$  the input vector  $V_t$  will be the input of our MLP predictor , which will have a *single* output neuron with state  $Z_t$ . This output  $Z_t$  will be the MLP prediction computed on day  $t$  for the *target*  $TARG_t = S(t+1)$ , which is not known at time  $t$ .

For this prediction task , we have a data set of  $(N-20)$  "cases"  $Case_{20}$   $Case_{21}$   $Case_{22}$  ...  $Case_{N-1}$  , indexed by  $t= 20, 21, \dots, N-1$ . Each  $Case_t$  is described by 18 features = 18 coordinates of vector  $V_t$ . The TRUE output to be predicted at time  $t$  is the yet unknown  $TARG_t = S(t+1)$ .

The data set of  $(N-20)$  cases for MLP prediction learning is denoted

$$PredCases = \{ \text{all pairs } (V_t, TARG_t) \text{ with } t= 20, 21, \dots, N-1 \}$$

Randomly Split the set *PredCases*, with 90% cases in the training set *PredTRAIN*, and 10 % cases in the test set *PredTEST*

4 *MLP predictor* : Our MLP predictor (*MLPpred*) will have the simple 3 layers architecture

$$\text{INPUT} \Rightarrow \text{HiddenLayer } K \Rightarrow \text{OUTPUT} \quad \text{with } \dim(\text{INPUT}) = 18, \dim(\text{OUTPUT}) = 1$$

$\dim(K) = k$  to be selected below. For each training input  $V_t$  we want the MLP output  $Z_t$  to be close to  $TARG_t = S(t+1)$ .

Implement PCA on the set of all input vectors  $V_t$ , with  $t = 20, 21, \dots, N$ . Determine the number  $k$  of principal components which preserves 95% of the variance (see HW3) and fix  $\dim(K) = k$ .

Compute the number  $w$  of weights and thresholds in this MLP, and compare  $w$  to the number of informations provided by the training set.

*5 Training of the MLP predictor:* Implement an automatic training on the training set PredTRAIN, with the options :

RELU response, Loss = "MSE", Stochastic Gradient Descent or ADAM, Batch Learning, Early Stopping

Let RMSE be the root mean squared error  $\sqrt{MSE}$ . Plot the evolution of RMSE versus the number of batches (one curve for the training set and one for the test set). Compare these two curves.

Plot on the same graph the true values  $TARG_t = S(t+1)$  and the predicted values  $Z_t$ . Comments.

Compute the Mean Relative Errors of Prediction MREP on the training set :

$MREP = \text{average} ( | Z_t - TARG_t | / TARG_t )$  over all cases in the Training set

Compute similarly MREP on the test set. Comments .

6. Denote NOD1 NOD2 ... NOD $k$  the hidden neurons . For  $j = 1 \dots k$ , compute and display the mean activity  $Y_j$  of NOD $j$  over all cases in the Training set. Display all the weights  $W_1 \dots W_k$  linking the neurons NOD1 ... NOD $k$  to the output node.

For each hidden NOD $j$  compute  $IMP_j = W_j Y_j = \text{average impact of NOD}_j \text{ on the prediction } Z_t$ . Display these  $k$  impacts and comment. Identify the hidden neuron NOD\* with maximal impact on  $Z_t$

7. Denote INP1 INP2 ... INP18 the 18 input neurons. Compute and display the mean activities  $X_1 \dots X_{18}$  of the 18 input neurons. Display all the weights  $U_1 \dots U_{18}$  linking the input nodes INP1 ... INP18 to the neuron NOD\*. For each input neuron INPs compute  $F_s = U_s X_s$  which is the *average impact* of input feature "s" on the key hidden neuron NOD\*. Identify the 5 input features with the largest impact on NOD\*. Comments.