CMPE 545, Artificial Neural Networks

# Term Project Report

Çağatay Yıldız - 2014700129

May 25, 2015

# 1 Time Series Clustering using Self-Organized Maps

In [1], SOM's are used to cluster 49 stocks' prices and gold commodity price for portfolio selection. Their method is based on **component planes**, that they define as *a representation that visualizes relative component values in the weight vectors of the SOM*. This is the same as considering a SOM whose each unit is a vector in $N$ dimension as $N$ different grids(=component planes) in $2D$. Then, clustering these $N$ grids corresponds to clustering components(features) in weight vectors.

## 1.1 Training Self-Organized Maps

I applied the same procedure as in [1] for clustering outlets so that outlets whose sale numbers behave similarly in time can be grouped. What's done for clustering $N$ time series, each having length $T$, using a $k \times l$ SOM where each unit is in $N$ dimension is as follows (In my case $N = 340$, $T = 1800$, $k = 20$, $l = 30$)

1. Each value in SOM is initialized to a number between 0 and 1 randomly. Also, each time series is z-normalized. This is because I am interested in how demand to the newspapers in an outlet changes, not the sale numbers themselves.

2. To train the SOM, a vector of dimension $N$ is used. This vector stores the sales of all outlets for a particular day. So, $T$ of such vectors are given to SOM as input in one epoch.

3. At the end of the training, $n$'th grid of the SOM corresponds to a $k \times l$ representation of $n$'th time series.

4. After that, the distance between each grid is calculated using **modified** $R_v$ coefficient, which they define as *a similarity coefficient between positive semi-definite matrices*. The choice of modified $R_v$ coefficient is because traditional $R_v$ coefficient deteriorates as maps get larger. Similarly in [2], the article in which modified $R_v$ coefficient is presented, it is noted that $R_v$ coefficient of completely random matrices increases as the size of matrices increases.
   Given two matrices $C_1$ and $C_2$, modified $R_v$ coefficient is calculated in two steps:

   - First, $T_i = C_i C_i^T - \text{diag}(C_i C_i^T)$ is calculated for $i = 1, 2$.
   - Then, the similarity measure is

   $$\sigma(C_1, C_2) = \frac{\text{trace}\{T_1^T T_2\}}{\sqrt{\text{trace}\{T_1^T T_2\} \times \text{trace}\{T_1^T T_2\}}} \tag{1}$$

## 1.2 Hierarchical Clustering

Similarity matrix of component planes at the hand, hierarchical clustering is then used to complete the procedure. I implemented single linkage clustering and average linkage clustering. In addition to the real data, I generated very simple synthetic data. You can see performances of both approaches in training data in Fig.1.
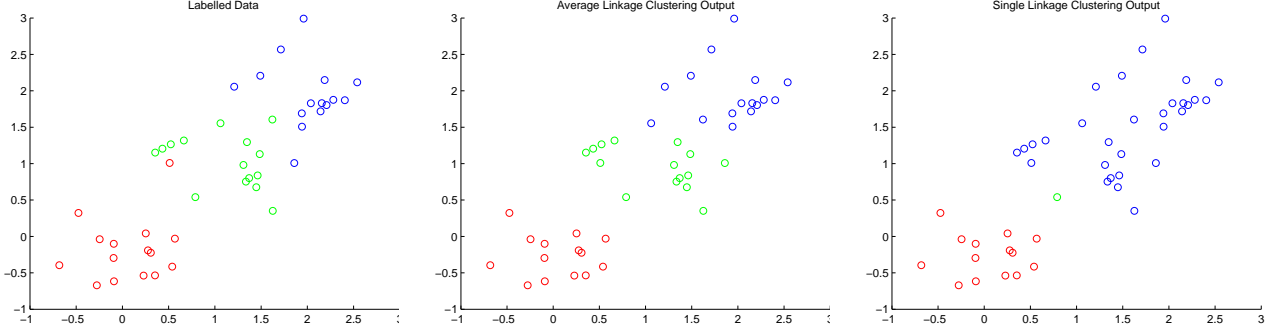


Figure 1: In this simple example, data from 3 different Gaussians centered at (0,0), (1,1) and (2,2) with standard deviation 0.4 is generated. Average linkage approach mislabels some of the data points but this is only because these points are indeed closer to another cluster center than the cluster center from which they are generated. As expected, single linkage approach performs pretty bad in such a data set, which is of course because of the nature of the data.

In Fig.2, you can see the visualization of component planes and sale numbers of randomly chosen 40 outlets. As noted before, we hoped heat maps of two outlets sharing similar trends to be as close as possible, which can roughly be confirmed looking at the figures.

One issue with hierarchical clustering is when to stop merging clusters. Starting with $N$ data points, theoretically it is possible to have any number of clusters from $N$ to 1 at the end of clustering. In this project, I went until merging all clusters to a giant one and recorded cluster labels of outlets at each time. That is, there are $N$ different clustering of the data.

In Fig.3, you can see sale figures of outlets randomly sampled from 4 different clusters. When this figure was drawn, the number of clusters was pulled down to $k = 10$. five of them were single data points, one of them contained 2 data points, three of them had 4-5 data points and the rest belonged to one cluster.

# 2 Time Series Prediction

The other task is to implement ANN's that learn those series and make predictions. For such purposes, recurrent neural networks are used. Two of the most popular RNN's are Elman network [4] and Jordan network [3]. In this project, I have implemented

- Tapped-delay multilayer perceptron, denoted as **TDMLP** in this report.

- Elman RNN, denoted as **EN** in this report.

- Full lag Elman RNN, denoted as **FLEN** in this report.

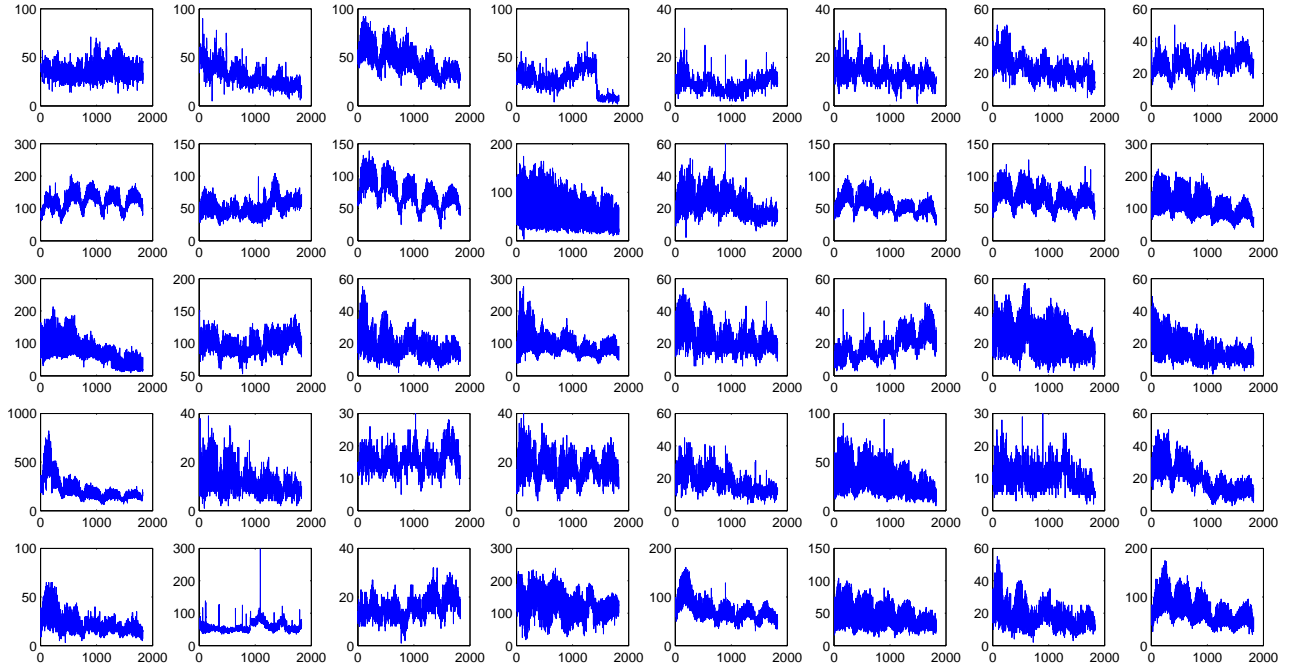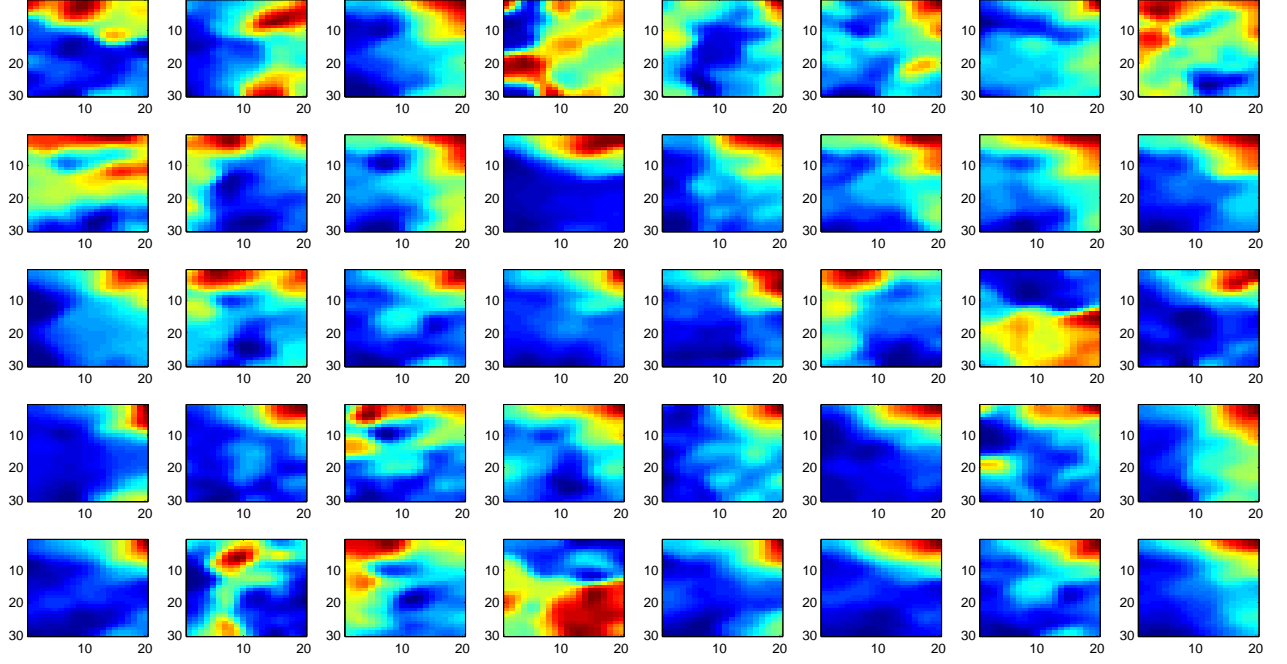- Tapped-delay Jordan RNN, denoted as **TDJN** in this report.

Figure 2: Heat maps and sale figures of randomly chosen 40 outlets. Maps and plots that are at the same coordinates in each subfigure belongs to the same outlet. For instance, sale plots at position (4,1) looks very much like that at (5,5) (all in MATLAB's matrix notation). As expected, corresponding heat maps are very alike, meaning that the distance between these component planes is small.

## 2.1 Details of Models

In this section, structures of models are explained. You can see the drawings of models in Fig.4
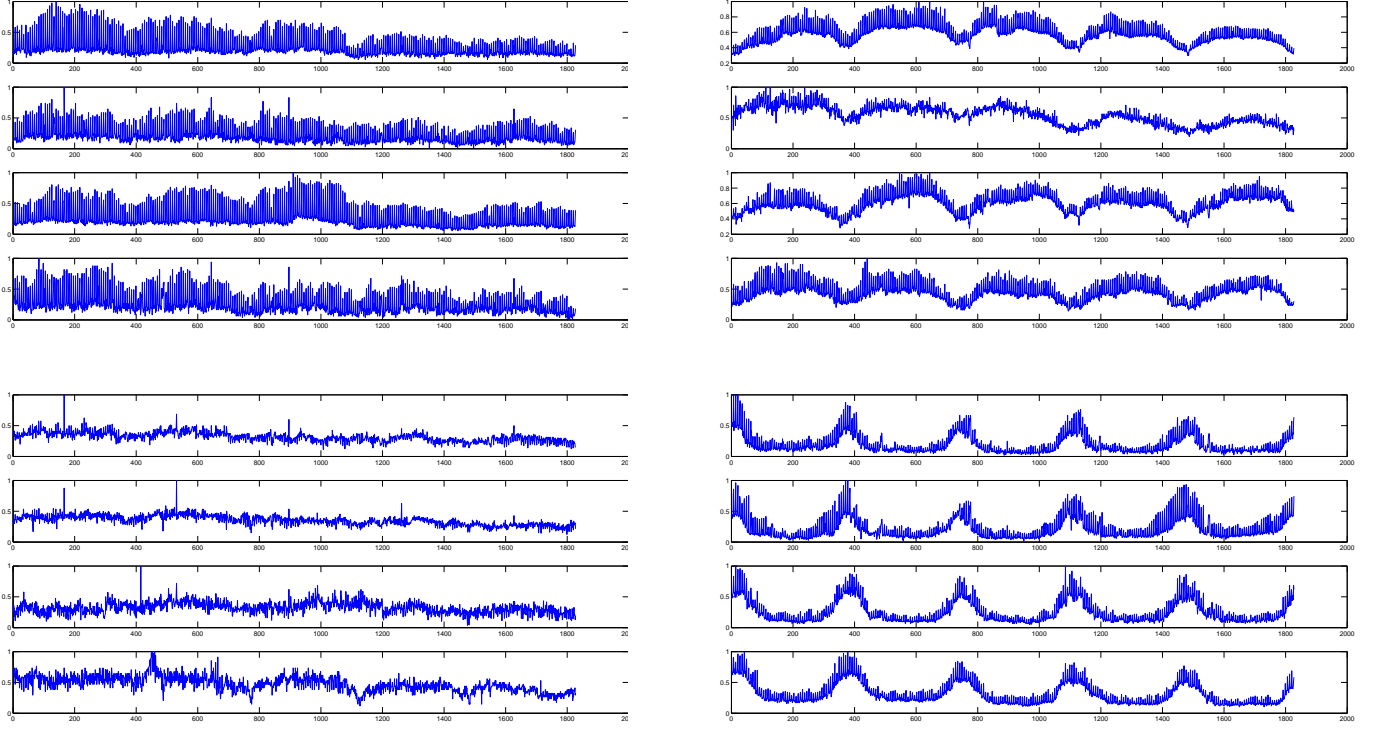Here is the list of all parameters used in networks:

3

Figure 3: Sale plots of outlets in 4 different clusters (out of $k = 10$ total clusters). Single linkage clustering was used to merge clusters.

- $t$: Time index

- $H$: Number of hidden units

- $D$: Dimensionality of the input

- $K$: Time lag (for instance $K = 7$ in newspaper sale problem)

- $z_i^t$: Value of the $i$'th hidden unit at time $t$

- $c_i^t$: Value of the $i$'th context unit at time $t$

- $s_i^t$: Value of the $i$'th state unit at time $t$

- $W^{oh}$: Weights between output and hidden unit

- $W^{hi}$: Weights between hidden unit and input unit

- $W^{hc}$: Weights between hidden unit and context unit

- $W^{hs}$: Weights between hidden unit and state unit

### 2.1.1 Tapped-Delay Multilayer Perceptron

In tapped-delay MLP, for each time $t$, data in a window of length $K$ is given as the input. At time $t + 1$, window is slid one unit on the data. This way, input is hoped to contain all the information that are

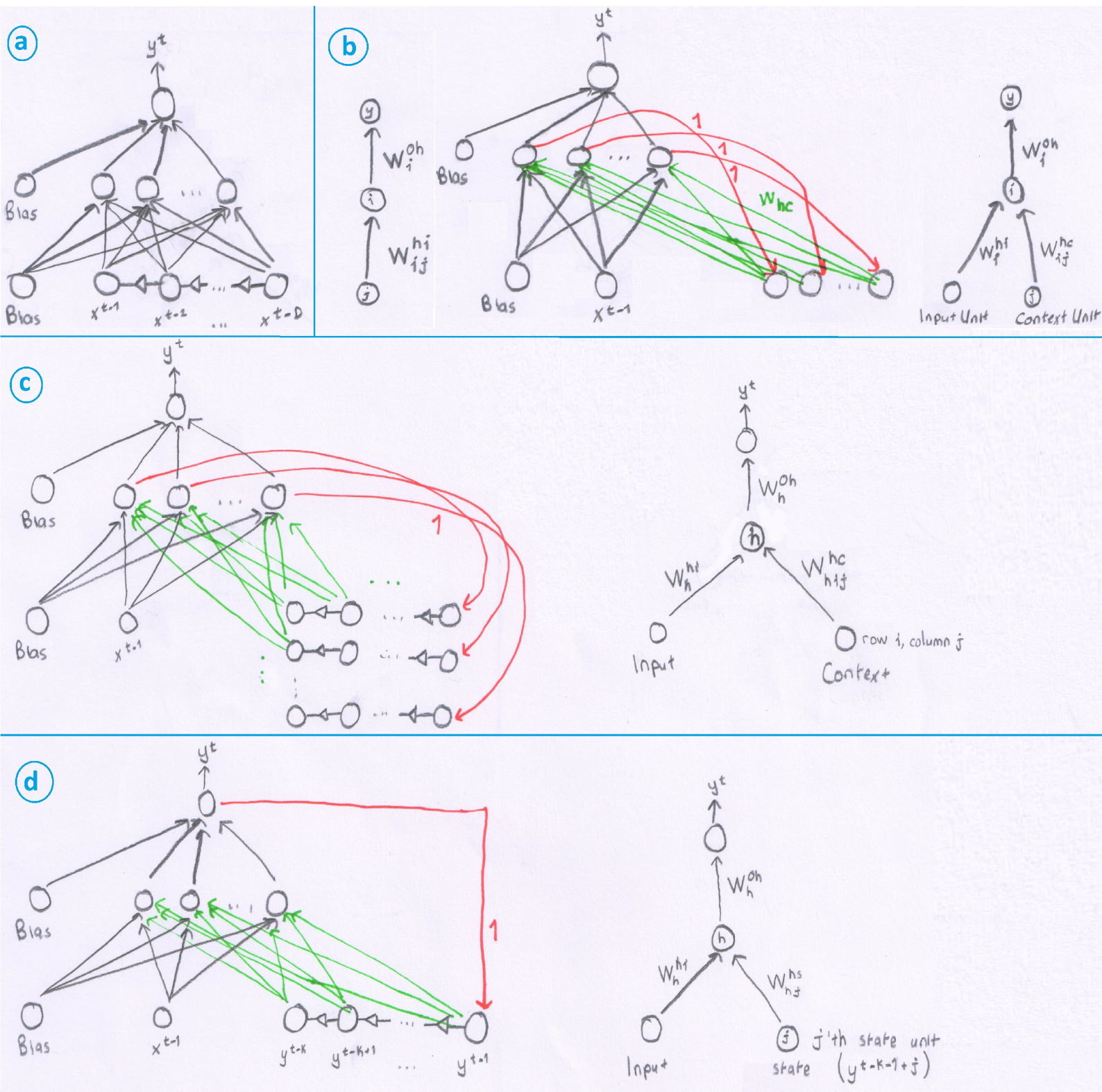4

Figure 4: Structures of all network types that are implemented. (a)Tapped-delay multilayer perceptron, (b)Elman neural network, (c) Full lag Elman RNN, (d)Tapped-delay Jordan RNN

enough to predict current value in time series.

**Update Equations:**

$$z_i^t = sigmoid\left(\sum_{j=1}^{D} W_{ij}^{hi} x^{t-j} + W_{i0}^{hi}\right)$$

$$y^t = \sum_{h=1}^{H} W_h^{oh} z_h^t + W_0^{oh}$$

$$\Delta W_{ij}^{hi} = \eta(x^t - y^t)W_i^{oh} z_i^t(1 - z_i^t)x^{t-j}$$

$$\Delta W_{i0}^{hi} = \eta(x^t - y^t)W_i^{oh} z_i^t(1 - z_i^t)$$

$$\Delta W_h^{oh} = \eta(x^t - y^t)z_h^t$$

$$\Delta W_0^{oh} = \eta(x^t - y^t)$$

### 2.1.2 Elman Recurrent Neural Network

In Elman network, there is so-called *context layer* in addition to input, hidden and output layers. Recurrency in this network occurs thanks to the cycles between hidden layer and context layer: Each context unit output is inputted to all hidden units (with different weights) and the output of hidden layer is 1:1 connected to the context layer ($c_i^{t+1} \leftarrow h_i^t, \forall i \in \{1, H\}$ where $c$ and $h$ represents context and hidden units and $t$ is the time index and $H$ is the number of hidden units.).

**Update Equations:**

$$z_i^t = sigmoid\left(W_i^{hi} x^{t-1} + W_0^{hi} + \sum_{j=1}^{H} W_{ij}^{hc} c_j^t\right)$$

$$y^t = \sum_{h=1}^{H} W_h^{oh} z_h^t + W_0^{oh}$$

$$\Delta W_i^{hi} = \eta(x^t - y^t)W_i^{oh} z_i^t(1 - z_i^t)x^{t-1}$$

$$\Delta W_0^{hi} = \eta(x^t - y^t)W_i^{oh} z_i^t(1 - z_i^t)$$

$$\Delta W_{ij}^{hc} = \eta(x^t - y^t)W_i^{oh} z_i^t(1 - z_i^t)c_j^t$$

$$\Delta W_h^{oh} = \eta(x^t - y^t)z_h^t$$

$$\Delta W_0^{oh} = \eta(x^t - y^t)$$

### 2.1.3 Full Lag Elman Recurrent Neural Network

FLEN is a variant of Elman networks. In this architecture, context layer stores not just $h_i^{t-1}, \forall i \in \{1, H\}$ but $h_i^{t-k}, \forall i \in \{1, H\}$ and $\forall k \in \{1, K\}$. As time ticks to $t + 1$, values in context units are shifted in such a way that $h_i^{t-K}, \forall i \in \{1, H\}$ leaves the context layer and $h_i^t, \forall i \in \{1, H\}$ is included.

Consider $c$ as a $H \times K$ matrix whose $j$'th column stores hidden layer values at time $t - K - 1 + j$. For example, for $K = 4$, the first column stores hidden layer values at time $t - 4$, second column stores those at time $t - 3$, etc. In this representation, the following assignments are done:

$$c_{i,j}^{t+1} \leftarrow c_{i,j+1}^t \ \forall i \in \{1, H\} \ \forall j \in \{1, H-1\}$$

$$c_{i,K}^{t+1} \leftarrow h_i^t \ \forall i \in \{1, H\}$$

**Update Equations:**

$$z_i^t = sigmoid\left(W_i^{hi}x^{t-1} + W_0^{hi} + \sum_{k=1}^{K}\sum_{j=1}^{H} W_{ijk}^{hc}c_{jk}^t\right)$$

$$y^t = \sum_{h=1}^{H} W_h^{oh}z_h^t + W_0^{oh}$$

$$\Delta W_i^{hi} = \eta(x^t - y^t)W_i^{oh}z_i^t(1 - z_i^t)x^{t-1}$$

$$\Delta W_0^{hi} = \eta(x^t - y^t)W_i^{oh}z_i^t(1 - z_i^t)$$

$$\Delta W_{ijk}^{hc} = \eta(x^t - y^t)W_i^{oh}z_i^t(1 - z_i^t)c_{jk}^t$$

$$\Delta W_h^{oh} = \eta(x^t - y^t)z_h^t$$

$$\Delta W_0^{oh} = \eta(x^t - y^t)$$

### 2.1.4 Tapped-Delay Jordan Recurrent Neural Network

In Jordan RNN, the additional layer is called *state layer*. Just like in Elman network, recurrency is done by adding cycles between hidden layer and state layer. Each state unit output is inputted to all hidden units. But in contrast to Elman network, state unit values at time $t + 1$ are equal to the outputs at time $t$. Therefore, there are as many units in the state layer as the number of outputs.

Since there is only one output in my problem, Jordan network seems unpromising. To overcome this issue, I set the size of state unit as $K$. At time $t$, state layer stores $y^{t-1:t-K}$. As time ticks, window on outputs is slid one unit forward and values in state layer becomes $y_{t:t-K+1}$.

**Update Equations:**

$$z_i^t = sigmoid\left(W_i^{hi}x^{t-1} + W_0^{hi} + \sum_{j=1}^{H} W_{ij}^{hs}s_j^t\right)$$

$$y^t = \sum_{h=1}^{H} W_h^{oh}z_h^t + W_0^{oh}$$

$$\Delta W_i^{hi} = \eta(x^t - y^t)W_i^{oh}z_i^t(1 - z_i^t)x^{t-1}$$

$$\Delta W_0^{hi} = \eta(x^t - y^t)W_i^{oh}z_i^t(1 - z_i^t)$$

$$\Delta W_{ij}^{hs} = \eta(x^t - y^t)W_i^{oh}z_i^t(1 - z_i^t)s_j^t$$

$$\Delta W_h^{oh} = \eta(x^t - y^t)z_h^t$$

$$\Delta W_0^{oh} = \eta(x^t - y^t)$$

## 2.2 Data Sets Used

I divided the data set (that is mentioned in previous section) into two parts(about 1200 data points in the first set and 600 in the second one), used the first one for training and the second for testing.

In addition, I have generated a number of time series and compared performances of the models. While generating data, I defined various functions where $x^t$ is a nonlinear function of $x^{t-1:t-L}$ for a fixed lag $L$ and a Gaussian noise added.

# 3 Results

For both synthetic and real data, in order to get the highest accuracy, I trained each model with different hyperparameters and then picked the best hyperparameters while comparing models. One interesting finding is that the change in the number of hidden units does not alter the convergence of error in the long run very much. You can see the plots in Fig.6.



Figure 6: Error figures for each method and different number of hidden units. Data set for this plots was generated by the function f = @(x)(x(1)*x(2)*0.3 - abs(x(4)) - sin(pi*x(3)) + 1) but all the functions I have tried yielded similar plots.

## 3.1 Synthetic Data Sets

- I did have hard time to generate time series. Most of my attempts yielded series that converge to the same number. Here, I have incorporated three pairs of error plots, each shows the results gathered from different time series. You can see those in Figures 7, 8 and 9.

- Figure 8 indicates that error rate of TDJN fluctuates so much. So, this learner seems to be problematic.

- Validation set errors converge to the same value in figures 8 and 9. That is somewhat a surprising finding and my best guess is it may be because of the nature of the data set.

- Among all approaches, TDMLP has the lowest overall error rate.

- In contrast to what we had encountered in homework, training error may make sudden increases during the application of a method to a time series. This is something that I cannot explain.



Figure 7: Error plots with function f = @(x)(x(1)*x(2)*0.3 - abs(x(4)) - sin(pi*x(3)) + 1). Number of hidden units are given in the legend.

9

Figure 8: Error plots with function f = @(x) (mod(x(1)*x(2) + 358,547)/211 - sigmoid(x(3)-x(4))*(rand ≥ 0.5)). Number of hidden units are given in the legend.



Figure 9: Error plots with function f = @(x)( sigmoid(x(1)*2-x(2)*5) + sin(pi*x(3)) - abs((x(1))/(x(1)+x(4))) - sin(pi*x(2)*4)). Number of hidden units are given in the legend.

## 3.2 Original Data Sets

- Here, I have incorporated three pairs of plots, each for one outlet. You can see those in Figures 10, 11 and 12.

- Since the raw data are made of integers that may be as big as 500, I needed to somehow scale down the data. Two approaches that I used were to divide each time series to the greatest number in the series (to scale down to 0-1 interval) and taking the derivative. These plots are generated from the data that are preprocessed with the first method described.

- As you see in all figures, FLEN performs really poor. When I first see the plots, I was suspicious of overfitting but decreasing number of hidden units did not change plots much.

- The main goal of a newspaper distributor is to predict demand and deliver newspapers accordingly. But none of these methods was able to do so. In fact, the returns increase substantially in many of the cases that I tested.

- In some of cases, the error rate in validation set does not decrease. So basically learner learns nothing.
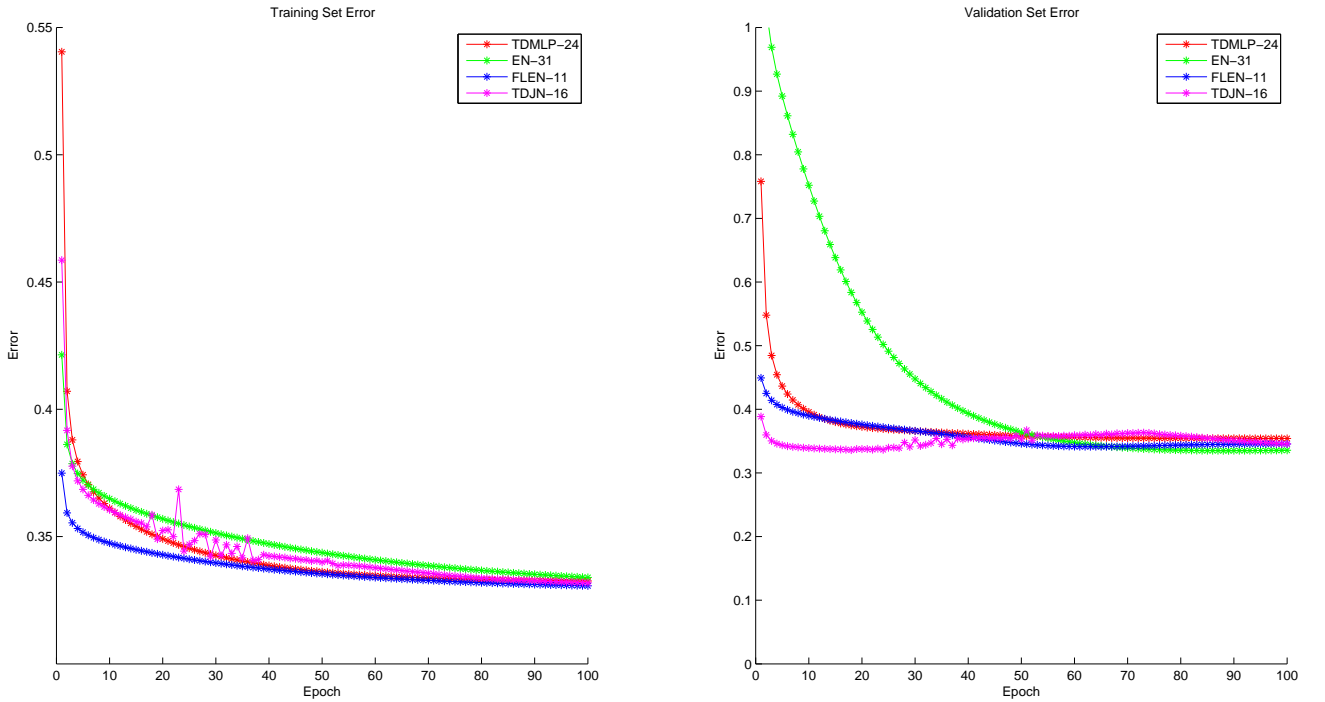
- In addition to those, for each outlet I consider sales in a day-of-week as one time series, which corresponds to 7 different time series for each outlet. The second way of processing data usually results in better predictions because standard deviation is smaller in the latter case, compared to having one long series for one outlet. Standard deviation is higher in the first case because newspaper sales tend to make jumps in weekends whereas numbers are relatively more stable during weekdays. Results are not included in this work because I ran out of time while experimenting.



Figure 10: Result for one of the outlets

# 4 Further Work

- More advanced ways of generating synthetic time series could be discovered.

- Real data can be preprocessed much better.

- It would be nice to compare performances of each method in each (time series) cluster.

11

Figure 11: Result for one of the outlets



Figure 12: Result for one of the outlets

# 5 Source Code

## 5.1 Manual for Source Code

Below is a brief explanation of what each script does. Please note that there are a lot of repetition in the source code this is because I preferred to keep things easy to see and debug.

### 5.1.1 SOM-related Code

- *som_main_script.m*: Just a wrapper script. Prepares data and model parameters, makes calls to other scripts and visualizes the result.

- *feature_clustering.m*: Implements SOM.

- *hierarchical_clustering.m*: Generates synthetic data and make visualization to test implementation of average and single linkage clustering.

- *average_linkage_clustering.m*: Implements average linkage clustering.

- *single_linkage_clustering.m*: Implements single linkage clustering.

### 5.1.2 RNN-related Code

- *generative_model.m*: Creates linear/nonlinear sequences of numbers with some Gaussian noise

- *script_synthetic_data.m*: Runs networks with synthetic data and plots figures

- *script_real_data.m*: Runs networks with one of the outlets and plots figures

- *mlp_updates.m*: Implements tapped-delay multilayer perceptron updates.

- *mlp_real_data_day_by_day.m*: Runs the TDMLP on real data. Here, sales of a particular day for an outlet is assumed to form one time series.

- *elman_updates.m*: Implements updates for Elman network

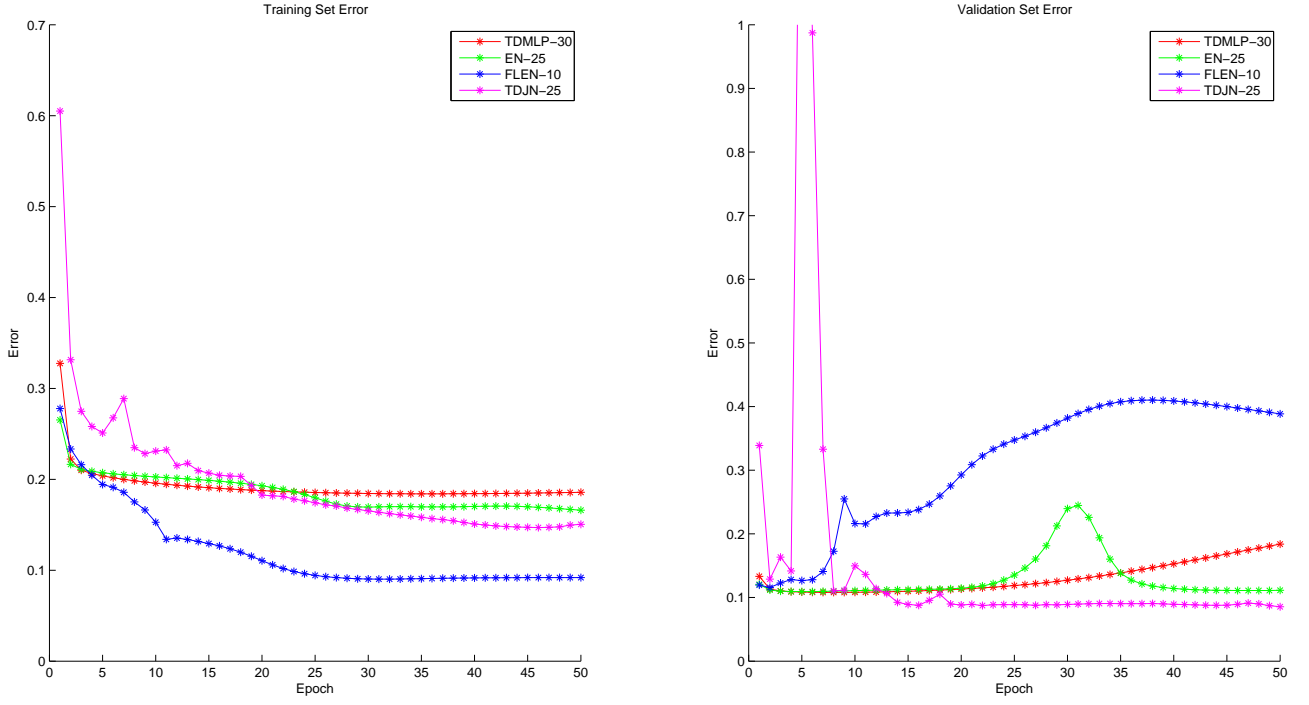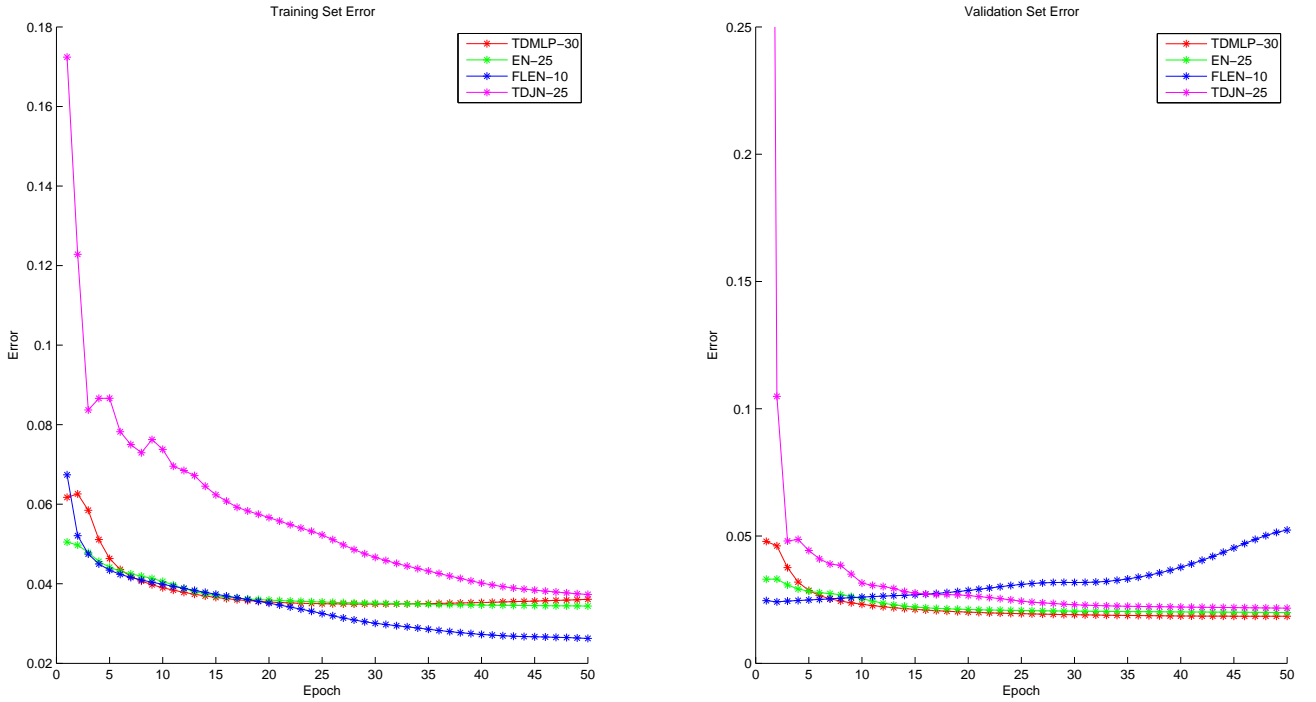- *elman_real_data_day_by_day.m*: Just like in *mlp_real_data_day_by_day.m*, assumes sales of a particular day for an outlet as one time series and runs EN updates.

- *full_lag_elman_updates.m*: Runs full lag Elman network updates on synthetic data

- *tapped_jordan_updates.m*: Runs tapped-delay Jordan network updates on synthetic data

## 5.2 SOM-related Code

### 5.2.1 som_main_script.m

```matlab
clc;
clear all;
close all;

%% initialization
T = 261*7; % length of time series = number of objects
N = 340; % number of time series = number of features

sRow = 30; % number of rows in SOM
sCol = 20; % number of columns in SOM

eta = 0.1; % learning rate
eta_update = 1 - 1e-3; % update parameter of eta
epC = 10; % number of epochs

```

```matlab
16  sigmaI = 5; % initial width
17  t1 = epC / log(sigmaI); % time constant
18
19  %% data preprocessing
20  load full_outlets_data; % only outlets with no missing info
21  Sat = reshape(Sat,T,N); % each column contains 261 week info of one outlet
22  zSat = bsxfun(@rdivide, Sat, max(Sat))'; % z-normalization
23
24  %% running feature clustering algorithm
25  feature_clustering
26
27  %% visualization of 40 random outlets info
28  nums = randi(N,1,40);
29  figure(1);
30  for i=1:40
31      subplot(5,8,i);
32      plot(Sat(:,nums(i)));
33      %ylim([-5,5]);
34  end
35  figure(2);
36  %som = permute(som,[2 3 1]);
37  for i=1:40
38      subplot(5,8,i);
39      imagesc(som(:,:,nums(i)));
40  end
41
42  %% extracting similarity matrix
43  similarity_matrix
44
45  %% hierarchical clustering
46  average_linkage_clustering
47  % single_linkage_clustering
48
49  %% end
50  beep;
51  pause(1);
52  beep;
```

### 5.2.2   feature_clustering.m

```matlab
1  % each day is another sample
2  % each sample in #_of_outlets=N dimension
3  som = rand(N, sRow, sCol); % som object
4  dis = zeros(sRow, sCol); % euclidian distance of current time series to all
      maps
5
6  for ep=1:epC
7      ep
8      %% updates
9      % sigma update
10     %{
```

```matlab
11      sigma = sigmaI * exp(-ep/t1);
12      variance = sigma^2;
13      %}
14
15      %another sigma update
16      sigma = sigmaI / (1 + 2*ep/epC);
17      var = sigma^2;
18
19      % learning rate update
20      eta = eta*eta_update;
21
22      % stores neighbourhood functions of neurons in the map
23      nhood = zeros(sRow, sCol, T);
24
25      %%
26      for t=randperm(T)
27
28          %% distance calculation bw current data point and all som
29          % objects
30          for r=1:sRow
31              for c=1:sCol
32                  tmp = som(:,r,c) - zSat(:,t);
33                  %dis(r,c) = tmp' * tmp;
34                  dis(r,c) = sqrt(tmp' * tmp);
35              end
36          end
37
38          %% winner som object
39          [minRow, minCol] = find(dis==min(min(dis)));
40
41          %% calculating neighborhood
42          for r=1:sRow
43              for c=1:sCol
44                  tmp = (minRow - r)^2 + (minCol - c)^2;
45                  nhood(r, c, t) = exp(-tmp/(2*var));
46              end
47          end
48
49          %% som update
50          for r=1:sRow
51              for c=1:sCol
52                  delta = eta*nhood(r, c, t)*(zSat(:,t)-som(:,r,c));
53                  som(:,r,c) = delta + som(:,r,c);
54              end
55          end
56      end
57  end
```

### 5.2.3  hierarchical_clustering.m

```matlab
1 %% synthetic data generation
```

```
2  clear all;
3  clc;
4  close all;
5  N = 45;
6  cc = [0,0; 1,1; 2,2]';
7  sigma = [0.4, 0.4, 0.4]';
8  nums = zeros(2,N);
9
10 for c=1:3
11     for i=1:N/3
12         nums(:,(c-1)*N/3+i) = cc(:,c) + sigma(c)*randn(2,1);
13     end
14 end
15
16 figure;
17 colors = {'ro','go','bo'};
18 hold on;
19 for c=1:3
20     for i=1:N/3
21         plot(nums(1,(c-1)*N/3+i),nums(2,(c-1)*N/3+i), colors{c});
22     end
23 end
24 title('Real Case');
25 hold off;
26
27 sigma = zeros(N); % distance matrix
28 for i=1:N
29     tmp = repmat(nums(:,i),1,N)-nums;
30     sigma(i,:) = -1*sum(tmp.^2); % inverting distances
31     sigma(i,1:i) = -Inf;
32 end
33
34 average_linkage_clustering;
35 single_linkage_clustering
```

### 5.2.4    average_linkage_clustering.m

```
1  %% hierarchical clustering (average linkage clustering)
2
3  % load sigma.mat;
4  c = num2cell(1:N); % cell i stores ids of outlets in cluster i
5  for i=1:N-1
6      c_count = length(c); % number of clusters
7      dis = ones(c_count)*-Inf; % distances between clusters
8
9      % calculation of inter-cluster distances
10     for c1 = 1:c_count-1
11         for c2 = c1+1:c_count
12             cluster1 = c{c1};
13             cluster2 = c{c2};
14             tmp = combvec(cluster1, cluster2);
```

```matlab
15              d = 0;
16               for  t=1:size(tmp,2)
17                   %lower  triangle  of  sigma  (similarity)  matrix  is  set  to −Inf
18                   %so  that  max  fnc.  can  be  used  with  one  line.
19                   d = d + max(sigma(tmp(1,t),tmp(2,t)),sigma(tmp(2,t),tmp(1,t)
                         ));
20               end
21               d = d/size(tmp,2); % average  similarity  bw  two  clusters
22               dis(c1,c2) = d;
23          end
24      end
25      [maxRow, maxCol] = find(dis==max(max(dis))); % most  similar  clusters
26
27      c{maxCol} = [c{maxCol}, c{maxRow}]; % merge  clusters
28      c(maxRow) = []; % remove  cluster
29      snapshots{i} = c; % save  current  clusters
30  end
31
32  if  0
33      %% visualization  of  clustering  with  real  data
34      k = 20; % number  of  clusters  at  the  time  we  stopped
35      m = 4; % only  clusters  with  m+ members  are  plotted
36      for  i=1:length(clusters)
37          if  length(clusters{i})>=m
38              figure;
39              for  j=1:m
40                  subplot(m,1,j);
41                  plot(zSat(clusters{i}(j),:));
42              end
43          end
44      end
45  else
46      %% visualization  of  synthetic  data
47      figure;
48      colors = {'ro','go','bo'};
49      hold  on;
50      for  c=1:3
51          for  i=1:length(snapshots{N−3}{c})
52              plot(nums(1,snapshots{N−3}{c}(i)),nums(2,snapshots{N−3}{c}(i)),
                     colors{c});
53          end
54      end
55      title('Average  Linkage  Clustering  Output');
56      hold  off
57  end
```

### 5.2.5 single_linkage_clustering.m

```matlab
1  %% hierarchical  clustering  (single  linkage  clustering)
2  % load  sigma.mat;
3  snapshots = zeros(N,N);
```

```matlab
c_id = (1:N)';
for c=1:N-1
    snapshots(:,c) = c_id;
    while 1
        [maxRow, maxCol] = find(sigma==max(max(sigma)));
        sigma(maxRow, maxCol) = -Inf;
        sigma(maxCol, maxRow) = -Inf;
        if c_id(maxCol) ~= c_id(maxRow)
            id1 = c_id(maxCol);
            id2 = c_id(maxRow);
            c_id(find(c_id==id2)) = id1;
            break;
        end
    end
end

if 0
    %% visualization of real data
    k = 5; % number of clusters
    m = 4; % only clusters with m+ members are plotted
    ss = snapshots(:,N-k+1); % current snapshot = cluster id's
    unqs = unique(ss); % number of clusters
    counts = zeros(1,k);
    for i=1:k
        counts(i) = length(find(ss==unqs(i)));
    end
    clusters = unqs(counts>=m); % clusters with m+ members
    for i=1:length(clusters)
        figure;
        outlet_ids = find(ss==clusters(i));
        for j=1:4
            subplot(4,1,j);
            plot(zSat(outlet_ids(j),:));
        end
    end
else
    %% visualization of synthetic data
    k = 3; % number of clusters
    ss = snapshots(:,N-k+1); % current snapshot = cluster id's
    unqs = unique(ss); % number of clusters
    figure;
    colors = {'ro','go','bo'};
    hold on;
    for c=1:3
        ids = find(ss==unqs(c));
        for j=1:length(ids)
            plot(nums(1,ids(j)),nums(2,ids(j)), colors{c});
        end
    end
    title('Single Linkage Clustering Output');
```

```
54        hold off;
55
56   end
```

### 5.2.6 RNN-related Code

### 5.2.7 generative_model.m

```matlab
1  D = 4; % time lag
2  N = 3000; % length of the data
3  mn = 0; % mean of the noise
4  var = 0.05; % variance of the noise
5  dataOriginal = zeros(N,1); % original data
6  dataOriginal(1:D) = rand(D,1);
7  data = dataOriginal; % data with gaussian noise
8  sigmoid = @(x)(1./(1+exp(-x))); % sigmoid function
9
10 %f = @(x)(x(1)*x(2)*0.3 - abs(x(4)) - sin(pi*x(3)) + 1);
11 % f = @(x)(min(sqrt(abs(x(1)*x(2))),0.5) - 0.25 - sin(pi*x(4)));
12 %f = @(x) (mod(x(1)*x(2) + 358,547)/211 - sigmoid(x(3)-x(4))*(rand>0.5) );
13 % f = @(x)( sigmoid(x(1)^2-x(2)) + sigmoid(-x(3)+x(4)/3) - (x(2)<(x(1)+x(3)+
       x(4))/3)*3 - (x(1)<(x(1)+x(3)+x(4))/3)*2   );
14 %f = @(x)( sigmoid(x(1)*2-x(2)*5) + sin(pi*x(3)) - abs((x(1))/(x(1)+x(4))) -
       sin(pi*x(2)*4)  );
15
16 for i = D+1:N
17     dataOriginal(i) = f(dataOriginal(i-D:i-1));
18 end
19 data(D+1:end) = dataOriginal(D+1:end) + mn*ones(N-D,1) + var*randn(N-D,1);
20
21 fprintf('max in data %d\n',max(data))
22 fprintf('min in data %d\n',min(data))
23
24 figure;
25 subplot(211);
26 plot(dataOriginal);
27 title('Original Data');
28 subplot(212);
29 plot(data);
30 title('Perturbated Data');
```

### 5.2.8 script_synthetic_data.m

```matlab
1  addpath MLP; % 1
2  addpath ELMAN; % 2
3  addpath FULL_LAG_ELMAN; % 3
4  addpath TAPPED_JORDAN; % 4
5
6  methods{1}='TDMLP';
7  methods{2}='EN';
8  methods{3}='FLEN';
9  methods{4}='TDJN';
10
11 %% data generation
12 generative_model;
13
```

```matlab
%% preparing  data
N = 2000; % number of data points in the training set
Nval = 999; % number of data points in the validation set
% dimensionality of the input
if config==1, D=4; else D=1; end

x = data(1:N);
r = data(D+1:N+1);
xval = data(N-D+2:end-1);
rval = data(N+2:end);

%% learning
Hs = [24, 31, 11, 16];
H = Hs(1);
mlp_updates
c{1} = error;
cv{1} = errorval;
H = Hs(2);
elman_updates
c{2} = error;
cv{2} = errorval;
H = Hs(3);
full_lag_elman_updates
c{3} = error;
cv{3} = errorval;
H = Hs(4);
tapped_jordan_updates
c{4} = error;
cv{4} = errorval;

%%
figure;
subplot(121);
hold on;
plot(1:length(c{1}), c{1},'-r*');
plot(1:length(c{2}), c{2},'-g*');
plot(1:length(c{3}), c{3},'-b*');
plot(1:length(c{4}), c{4},'-m*');
legend(sprintf('TDMLP-%d',Hs(1)),...
sprintf('EN-%d',Hs(2)),...
sprintf('FLEN-%d',Hs(3)),...
sprintf('TDJN-%d',Hs(4)));
hold off;
title('Training Set Error');
xlabel('Epoch');
ylabel('Error');
subplot(122);
hold on;
plot(1:length(cv{1}), cv{1},'-r*');
plot(1:length(cv{2}), cv{2},'-g*');
```

```matlab
64  plot(1:length(cv{3}), cv{3}, '-b*');
65  plot(1:length(cv{4}), cv{4}, '-m*');
66  legend(sprintf('TDMLP-%d',Hs(1)),...
67  sprintf('EN-%d',Hs(2)),...
68  sprintf('FLEN-%d',Hs(3)),...
69  sprintf('TDJN-%d',Hs(4)));
70  hold off;
71  title('Validation Set Error');
72  xlabel('Epoch');
73  ylabel('Error');
74  ylim([0,1000]);
```

### 5.2.9   script_real_data.m

```matlab
1   close all;
2   load full_outlets_data.mat; % only outlets with no missing info
3
4   methods{1}='TDMLP';
5   methods{2}='EN';
6   methods{3}='FLEN';
7   methods{4}='TDJN';
8
9   T = 261*7;
10  N = 340;
11  Sat = reshape(Sat,T,N); % each column contains 261 week info of one outlet
12  Iade = reshape(Iade,T,N); % each column contains 261 week info of one outlet
13  Sat = log(Sat);
14  %Ncons = max(Sat);
15  %Sat = bsxfun(@rdivide, Sat/5, Ncons); % division by max/5
16
17  %% reading data
18  N = 1200; % number of data points in the training set
19  Nval = 626; % number of data points in the validation set
20  index = randi(340);
21  x = Sat(1:N,index);
22  r = Sat(D+1:N+1,index);
23  xval = Sat(N-D+2:end-1,index);
24  rval = Sat(N+2:end,index);
25
26  Hs = [30, 25, 10, 25];
27
28  %% learning
29  H = Hs(1);
30  mlp_updates
31  c{1} = error;
32  cv{1} = errorval;
33  H = Hs(2);
34  elman_updates
35  c{2} = error;
36  cv{2} = errorval;
37  H = Hs(3);
```

22

```matlab
38  Lag = 7;
39  full_lag_elman_updates
40  c{3} = error;
41  cv{3} = errorval;
42  H = Hs(4);
43  Lag = 7;
44  tapped_jordan_updates
45  c{4} = error;
46  cv{4} = errorval;
47
48  %% plotting error figures
49  if 1
50      figure;
51      subplot(121);
52      hold on;
53      plot(1:length(c{1}), c{1}, '-r*');
54      plot(1:length(c{2}), c{2}, '-g*');
55      plot(1:length(c{3}), c{3}, '-b*');
56      plot(1:length(c{4}), c{4}, '-m*');
57      legend(sprintf('TDMLP-%d',Hs(1)),...
58          sprintf('EN-%d',Hs(2)),...
59          sprintf('FLEN-%d',Hs(3)),...
60          sprintf('TDJN-%d',Hs(4)));
61      hold off;
62      title('Training Set Error');
63      xlabel('Epoch');
64      ylabel('Error');
65      subplot(122);
66      hold on;
67      plot(1:length(cv{1}), cv{1}, '-r*');
68      plot(1:length(cv{2}), cv{2}, '-g*');
69      plot(1:length(cv{3}), cv{3}, '-b*');
70      plot(1:length(cv{4}), cv{4}, '-m*');
71      legend(sprintf('TDMLP-%d',Hs(1)),...
72          sprintf('EN-%d',Hs(2)),...
73          sprintf('FLEN-%d',Hs(3)),...
74          sprintf('TDJN-%d',Hs(4)));
75      hold off;
76      title('Validation Set Error');
77      xlabel('Epoch');
78      ylabel('Error');
79      ylim([0,1]);
80  end
81
82  %% plotting sale figures
83  if 1
84      figure;
85      plot(1:N,y*Ncons(index),'-r*',1:N,r*Ncons(index),'-b*')
86      title('Training plots')
87      legend('prediction','real sale')
```

```matlab
88      figure;
89      plot(1:Nval,yval*Ncons(index),'-r*',1:Nval,rval*Ncons(index),'-b*')
90      title('Test plots')
91      legend('prediction','real sale')
92  end
93
94  %% plotting error figures
95  if 0
96      figure;
97      subplot(1,3,1);
98      plot(1:N,r*Ncons(indexes),'-b')
99      title('Training sale plots')
100     subplot(1,3,2);
101     plot((y-r)*Ncons(indexes))
102     title('Training diff plots')
103     subplot(1,3,3);
104     plot((yval-rval)*Ncons(indexes))
105     sum(abs((yval-rval)*Ncons(indexes)))
106     sum(rval*Ncons(indexes))
107     title('Test diff plots')
108 end
109
110 %% errors
111 diff = yval - rval ;
112 returns = sum(ceil(diff(diff>0)*Ncons(index)));
113 misses = sum(ceil(diff(diff<0)*Ncons(index)));
114 [sum(Sat(N+2:end,index)*Ncons(index)),sum(Iade(N+2:end,index)),sum(returns),
        sum(misses)]
```

### 5.2.10   mlp_updates.m

```matlab
1  %% constants
2  %H = 10; % number of hidden units
3  K = 1; % number of outputs
4  Whi = rand(H,D+1)-0.5; % weights bw real input and hidden unit, weights in
       rows
5  Woh = rand(K,H+1)-0.5; % weights bw hidden unit and output unit
6  epoch_c = 50; % number of epochs
7  sigmoid = @(x)(1./(1+exp(-x))); % sigmoid function
8  eta = 0.2; % eta
9  eta_update = 0.95; % in each epoch, eta is multiplied with this number
10 error = zeros(epoch_c,1); % training error
11 errorval = zeros(epoch_c,1); % validation error
12
13 %% perceptron updates
14 for ep = 1:epoch_c  % for each epoch
15     ep
16     dwoh = zeros(K,H+1);
17     dwhi = zeros(H,D+1);
18     eta = eta*eta_update;
19     y = zeros(N-D+1,1);
```

```
20        for  t = 1:N−D+1
21            tmp = sigmoid ( Whi*[1; x(t:t+D−1)] );
22            Z = [1; tmp];
23            y(t) = Woh * Z;
24            dwoh = eta*(r(t)−y(t))*Z';
25            dwhi = zeros(H,D+1);
26            for h=1:H
27                dwhi(h,:) = eta*(r(t)−y(t))*Woh(h+1)*Z(h+1)*(1−Z(h+1))*[1;x(t:t+
                     D−1)];
28            end
29            Woh = Woh + dwoh;
30            Whi = Whi + dwhi;
31        end
32
33        %error calculation
34        yval = zeros(Nval,1);
35        for t=1:Nval
36            tmp = sigmoid ( Whi*[1; xval(t:t+D−1)] );
37            Z = [1; tmp];
38            yval(t) = Woh * Z;
39        end
40
41        error(ep) = ((y−r)'*(y−r))/N;
42        errorval(ep) = ((yval−rval)'*(yval−rval))/Nval;
43    end
```

### 5.2.11   mlp_real_data.m

```
1  load 'full_outlets_data.mat'; % only outlets with no missing info
2  T = 261*7;
3  N = 340;
4
5  X = 1; % number of outlets to process
6  indexes = randi(340,1,X);
7
8  for i=1:X
9
10      index = indexes(i);
11      index = 10;
12      Ncons = zeros(1,7);
13      returns = zeros(1,7);
14      misses = zeros(1,7);
15
16      for d=1:7
17          Satd = Sat(d,:,index); % each column contains 261 week info of one
                 outlet
18          Ncons(d) = max(Satd);
19          Satd = Satd'/Ncons(d);
20
21          D = 7; % dimensionality of the input
22          N = 200; % number of data points in the training set
```

25

```matlab
            Nval = 60; % number of data points in the validation set

        %% reading data
        x = Satd(1:N);
        r = Satd(D+1:N+1);
        xval = Satd(N-D+2:end-1);
        rval = Satd(N+2:end);

        %mlp_day_by_day_updates;
        mlp_updates;

        diff = yval - rval ;
        returns(d) = sum( ceil( diff( diff >0)*Ncons(d)));
        misses(d) = sum( ceil( diff( diff <0)*Ncons(d)));


        %% plotting error figures
        if 0
            figure;
            subplot(211);
            plot(1:length(error), error,'-r*');
            title('Training Set Error');
            subplot(212);
            plot(1:length(errorval), errorval,'-r*');
            title('Validation Set Error');
        end
    end
    [sum(sum(Sat(:,end-Nval:end,index))),sum(sum(Iade(:,end-Nval:end,index))
        ),sum(returns),sum(misses)]


    %% plotting error figures
    if 0
        figure;
        subplot(211);
        plot(1:length(error), error,'-r*');
        title('Training Set Error');
        subplot(212);
        plot(1:length(errorval), errorval,'-r*');
        title('Validation Set Error');
    end

    %% plotting sale figures
    if 0
        figure;
        plot(1:N,y*Ncons(indexes),'-r',1:N,r*Ncons(indexes),'-b')
        title('Training plots')
        legend('prediction','real sale')
        figure;
```

```
71          plot(1:Nval,yval*Ncons(indexes),'-r',1:Nval,rval*Ncons(indexes),'-b'
                )
72          title('Test plots')
73          legend('prediction','real sale')
74      end
75
76      %% plotting error figures
77      if 0
78          figure;
79          subplot(3,1,1);
80          plot(r*Ncons(index))
81          title('Training sale plots')
82          subplot(3,1,2);
83          plot((y-r)*Ncons(index))
84          title('Training diff plots')
85          subplot(3,1,3);
86          plot((yval-rval)*Ncons(index))
87          sum(abs((yval-rval)*Ncons(index)))
88          sum(rval*Ncons(index))
89          title('Test diff plots')
90      end
91 end
```

### 5.2.12    elman_updates.m

```
1  %% constants
2  D = 1; % dimensionality of the input
3  H = 15; % number of hidden units
4  K = 1; % number of outputs
5  C = ones(H,1)/2; % context unit
6  Whc = rand(H,H)-0.5; % weights bw hidden unit and context unit. first row =
       weights from context units to the first hidden unit
7  Whi = rand(H,D+1)-0.5; % weights bw real input and hidden unit, weights in
       rows
8  Woh = rand(K,H+1)-0.5; % weights bw hidden unit and output unit
9  epoch_c = 50; % number of epochs
10 sigmoid = @(x)(1./(1+exp(-x))); % sigmoid function
11 eta = 0.2; % eta
12 mom = 0; % momentum term
13 eta_update = 0.98; % in each epoch, eta is multiplied with this number
14 error = zeros(epoch_c,1); % training error
15 errorval = zeros(epoch_c,1); % validation error
16
17 %% perceptron updates
18 for ep = 1:epoch_c  % for each epoch
19     ep
20     dwhc = zeros(H,H);
21     dwhi = zeros(H,D+1);
22     dwoh = zeros(K,H+1);
23     eta = eta*eta_update;
24     y = zeros(N-D+1,1);
```

```matlab
25        yval = zeros(Nval,1);
26        for t = 1:N-D+1
27            C_old = C;
28            C = sigmoid(Whi*[1; x(t:t+D-1)] + Whc*C_old); % stores hidden value
                  for this iteration
29            Z = [1; C];
30            y(t) = Woh * Z;
31
32            dwoh_old = dwoh;
33            dwhi_old = dwhi;
34            dwhc_old = dwhc;
35
36            dwoh = eta*(r(t)-y(t))*Z';
37            dwhi = zeros(H,D+1);
38            dwhc = zeros(H,H);
39            for h=1:H
40                dwhi(h,:) = eta*(r(t)-y(t))*Woh(h+1)*Z(h+1)*(1-Z(h+1))*[1;x(t:t+
                      D-1)]';
41                dwhc(h,:) = eta*(r(t)-y(t))*Woh(h+1)*Z(h+1)*(1-Z(h+1))*C_old';
42            end
43            Woh = Woh + (1-mom)*dwoh + mom*dwoh_old;
44            Whi = Whi + (1-mom)*dwhi + mom*dwhi_old;
45            Whc = Whc + (1-mom)*dwhc + mom*dwhc_old;
46        end
47
48        %error calculation
49        for t=1:Nval
50            C = sigmoid (Whi*[1; xval(t:t+D-1)] + Whc*C);
51            Z = [1; C];
52            yval(t) = Woh * Z;
53        end
54        error(ep) = ((y-r)'*(y-r))/N;
55        errorval(ep) = ((yval-rval)'*(yval-rval))/Nval;
56
57    end
```

### 5.2.13 elman_real_data_day_by_day.m

```matlab
1  close all;
2  %clc;
3  clear all;
4
5  load 'full_outlets_data.mat'; % only outlets with no missing info
6  T = 261*7;
7  N = 340;
8
9  X = 1; % number of outlets to process
10 indices = randi(340,1,1);
11 indices = 20;
12
13 results = zeros(4,X);
```

```matlab
14
15  for i=1:X
16
17      index = indices(i);
18      Ncons = zeros(1,7);
19      returns = zeros(1,7);
20      misses = zeros(1,7);
21
22      for d=1:7
23
24          %% preparing data
25          tmp = Sat(d,:,index); % each column contains 261 week info of one
                  outlet
26          Ncons(d) = max(tmp);
27          tmp = tmp'/Ncons(d);
28
29          N = 200; % number of data points in the training set
30          Nval = 60; % number of data points in the validation set
31          x = tmp(1:N);
32          r = tmp(2:N+1);
33          xval = tmp(N-1:end-1);
34          rval = tmp(N+2:end);
35
36          elman_updates;
37
38          diff = yval - rval ;
39          returns(d) = sum(ceil(diff(diff>0)*Ncons(d)));
40          misses(d) = sum(ceil(diff(diff<0)*Ncons(d)));
41
42          %% plotting error figures
43          if 0
44              figure;
45              subplot(211);
46              plot(1:length(error), error,'-r*');
47              title('Training Set Error');
48              subplot(212);
49              plot(1:length(errorval), errorval,'-r*');
50              title('Validation Set Error');
51          end
52
53          %% plotting sale figures
54          if 1
55              figure;
56              plot(1:length(y),ceil(y*Ncons(d)),'-r*',1:length(r),ceil(r*Ncons
                  (d)),'-b*')
57              title('Training plots')
58              legend('prediction','real sale')
59              figure;
60              plot(1:length(yval),ceil(yval*Ncons(d)),'-r*',1:length(rval),
                  ceil(rval*Ncons(d)),'-b*')
```

```
61              title('Test plots')
62              legend('prediction','real sale')
63           end
64
65        %% plotting miss-return figures
66         if 0
67              figure;
68              subplot(3,1,1);
69              plot(r*Ncons(d))
70              title('Training sale plots')
71              subplot(3,1,2);
72              plot((y-r)*Ncons(d))
73              title('Training diff plots')
74              subplot(3,1,3);
75              plot((yval-rval)*Ncons(d))
76              sum(abs((yval-rval)*Ncons(d)))
77              sum(rval*Ncons(d))
78              title('Test diff plots')
79           end
80
81     end
82     results(:,i) = [sum(sum(Sat(:,end-Nval:end,index))),sum(sum(Iade(:,end-
          Nval:end,index))),sum(returns),sum(misses)]';
83
84 end
85
86 results
```

### 5.2.14  full_lag_elman_update.m

```
1  %% constants
2  D = 1; % dimensionality of the input
3  %H = 10; % number of hidden units
4  K = 1; % number of outputs
5  Lag = 4; % number of hidden unit "groups" stored
6  C = ones(H,Lag)/2; % context unit
7  Whc = rand(H,H,Lag)-0.5; % weights bw hidden unit and context unit. first
      row = weights from context units to the first hidden unit
8  Whi = rand(H,D+1)-0.5; % weights bw real input and hidden unit, weights in
      rows
9  Woh = rand(K,H+1)-0.5; % weights bw hidden unit and output unit
10 epoch_c = 50; % number of epochs
11 sigmoid = @(x)(1./(1+exp(-x))); % sigmoid function
12 eta = 0.1; % eta
13 mom = 0.5; % momentum term
14 eta_update = 0.98; % in each epoch, eta is multiplied with this number
15 error = zeros(epoch_c,1); % training error
16 errorval = zeros(epoch_c,1); % validation error
17
18 %% perceptron updates
19 for ep = 1:epoch_c  % for each epoch
```

```matlab
20        ep
21        dwhc = zeros(H,H,Lag);
22        dwhi = zeros(H,D+1);
23        dwoh = zeros(K,H+1);
24
25        eta = eta*eta_update;
26        y = zeros(N-D+1,1);
27        yval = zeros(Nval,1);
28        hdns = zeros(H,Lag+1); % stores hidden unit values
29
30        for t = 1:N-D+1
31            C_old = C;
32            hdns(:,1) = Whi*[1; x(t:t+D-1)]; % input to hidden units
33            for lag=1:Lag
34                hdns(:,lag+1) = Whc(:,:,lag)*C(:,lag); % contribution of each
                    lag
35            end
36            hdns = sigmoid(sum(hdns,2)); % stores hidden values for this
                iteration
37            Z = [1; hdns];
38            y(t) = Woh * Z;
39
40            dwoh_old = dwoh;
41            dwhi_old = dwhi;
42            dwhc_old = dwhc;
43
44            dwoh = eta*(r(t)-y(t))*Z';
45            dwhi = zeros(H,D+1);
46            dwhc = zeros(H,H,Lag);
47            for h=1:H
48                dwhi(h,:) = eta*(r(t)-y(t))*Woh(h+1)*Z(h+1)*(1-Z(h+1))*[1;x(t:t+
                    D-1)]';
49            end
50            for lag=1:Lag
51                for h=1:H
52                    dwhc(h,:,lag) = eta*(r(t)-y(t))*Woh(h+1)*Z(h+1)*(1-Z(h+1))*
                        C_old(:,lag)';
53                end
54            end
55            Woh = Woh + (1-mom)*dwoh + mom*dwoh_old;
56            Whi = Whi + (1-mom)*dwhi + mom*dwhi_old;
57            Whc = Whc + (1-mom)*dwhc + mom*dwhc_old;
58            C = [C(:,2:end), hdns];
59        end
60
61        hdns = zeros(H,Lag+1); % stores hidden unit values
62        %error calculation
63        for t=1:Nval
64            C_old = C;
65            hdns(:,1) = Whi*[1; x(t:t+D-1)]; % input to hidden units
```

```
66        for  lag=1:Lag
67            hdns(:,lag+1) = Whc(:,:,lag)*C(:,lag); % contribution of each
                 lag
68        end
69        hdns = sigmoid(sum(hdns,2)); % stores hidden values for this
             iteration
70
71        Z = [1; hdns];
72        yval(t) = Woh * Z;
73        C = [C(:,2:end), hdns];
74    end
75    error(ep) = ((y-r)'*(y-r))/N;
76    errorval(ep) = ((yval-rval)'*(yval-rval))/Nval;
77
78 end
```

### 5.2.15   tapped_jordan.m

```
1  %% constants
2  D = 1; % dimensionality of the input
3  %H = 20; % number of hidden units
4  K = 1; % number of outputs
5  Lag = 4; % number of outputs stored and given as input to hidden layer
6  S = rand(Lag,1)-0.5; % state layer
7  Whs = rand(H,Lag)-0.5; % weights bw hidden unit and state unit. first row =
       weights from context units to the first hidden unit
8  Whi = rand(H,D+1)-0.5; % weights bw real input and hidden unit, weights in
       rows
9  Woh = rand(K,H+1)-0.5; % weights bw hidden unit and output unit
10 epoch_c = 50; % number of epochs
11 sigmoid = @(x)(1./(1+exp(-x))); % sigmoid function
12 eta = 0.2; % eta
13 mom = 0.4; % momentum term
14 eta_update = 0.98; % in each epoch, eta is multiplied with this number
15 error = zeros(epoch_c,1); % training error
16 errorval = zeros(epoch_c,1); % validation error
17 % changes in weights
18 dwhs = zeros(H,Lag);
19 dwhi = zeros(H,D+1);
20 dwoh = zeros(K,H+1);
21
22 %% perceptron updates
23 for ep = 1:epoch_c  % for each epoch
24     ep
25     eta = eta*eta_update;
26     y = zeros(N-D+1,1);
27     yval = zeros(Nval,1);
28     S = ones(Lag,1)/2;
29
30     for t = 1:N-D+1
31         Hunits = sigmoid (Whi*[1; x(t:t+D-1)] + Whs*S); % hidden unit values
```

```matlab
32
33            Z = [1; Hunits];
34            y(t) = Woh * Z;
35
36            % storing last updates
37            dwoh_old = dwoh;
38            dwhi_old = dwhi;
39            dwhs_old = dwhs;
40
41            % updates
42            dwoh = eta*(r(t)-y(t))*Z';
43            for h=1:H
44                %dwhi(h,:) = dwhi(h,:) +  eta*(r(t)-y(t))*Woh(h+1)*Z(h+1)*(1-Z(h
                      +1))*[1;x(t:t+D-1)]';
45                %dwhc(h,:) = dwhc(h,:) +  eta*(r(t)-y(t))*Whc(h+1)*Z(h+1)*(1-Z(h
                      +1))*C';
46
47                dwhi(h,:) = eta*(r(t)-y(t))*Woh(h+1)*Z(h+1)*(1-Z(h+1))*[1;x(t:t+
                      D-1)]';
48                dwhs(h,:) = eta*(r(t)-y(t))*Whs(h+1)*Z(h+1)*(1-Z(h+1))*S';
49            end
50
51            % update with momentum
52            Woh = Woh + (1-mom)*dwoh + mom*dwoh_old;
53            Whi = Whi + (1-mom)*dwhi + mom*dwhi_old;
54            Whs = Whs + (1-mom)*dwhs + mom*dwhs_old;
55
56            % update state layer
57            S = [S(2:end); y(t)];
58            %pause;
59
60        end
61
62        %error calculation
63        for t=1:Nval
64            tmp = sigmoid (Whi*[1; xval(t:t+D-1)] + Whs*S);
65            Z = [1; tmp];
66            yval(t) = Woh * Z;
67            S = [S(2:end); yval(t)];
68        end
69        error(ep) = ((y-r)'*(y-r))/N;
70        errorval(ep) = ((yval-rval)'*(yval-rval))/Nval;
71
72 end
73
74 %{
75 %error calculation
76 C = ones(H,1)/2; % context unit
77 for t=1:Nval
78
```

```matlab
79      C = sigmoid (Whi*[1; xval(t:t+D-1)] + Whc*C);
80      Z = [1; C];
81      yval(t) = Woh * Z;
82
83      if 0
84          dwoh = eta*(rval(t)-yval(t))*Z';
85          dwhi = zeros(H,D+1);
86          dwhc = zeros(H,H);
87          for h=1:H
88              %dwhi(h,:) = dwhi(h,:) +  eta*(r(t)-y(t))*Woh(h+1)*Z(h+1)*(1-Z(h
                    +1))*[1;x(t:t+D-1)]';
89              %dwhc(h,:) = dwhc(h,:) +  eta*(r(t)-y(t))*Whc(h+1)*Z(h+1)*(1-Z(h
                    +1))*C';
90              dwhi(h,:) = eta*(rval(t)-yval(t))*Woh(h+1)*Z(h+1)^2*exp(-1*Whc(h
                    ,:)*C)*[1;xval(t:t+D-1)]';
91              dwhc(h,:) = eta*(rval(t)-yval(t))*Whc(h+1)*Z(h+1)^2*exp(-1*Whi(h
                    ,:)*[1; xval(t:t+D-1)])*C';
92          end
93          Woh = Woh + (1-mom)*dwoh + mom*dwoh_old;
94          Whi = Whi + (1-mom)*dwhi + mom*dwhi_old;
95          Whc = Whc + (1-mom)*dwhc + mom*dwhc_old;
96      end
97  end
98  error(ep) = ((y-r)'*(y-r))/(N*std(r));
99  errorval(ep) = ((yval-rval)'*(yval-rval))/(Nval*std(rval));
100 %}
```

# References

[1] B. Silva and N. C. Marques, "Feature clustering with self-organizing maps and an application to financial time-series for portfolio selection." in *IJCCI (ICFC-ICNC)*, 2010, pp. 301–309.

[2] A. K. Smilde, H. A. L. Kiers, S. Bijlsma, C. M. Rubingh, and M. J. van Erk, "Matrix correlations for high-dimensional data: the modified rv-coefficient." *Bioinformatics*, vol. 25, no. 3, pp. 401–405, 2009. [Online]. Available: http://dblp.uni-trier.de/db/journals/bioinformatics/bioinformatics25.html

[3] M. I. Jordan, "Serial order: A parallel, distributed processing approach," Institute for Cognitive Science, University of California, San Diego, Tech. Rep. 8604, 1986.

[4] J. L. Elman, "Finding structure in time," *Cognitive science*, vol. 14, no. 2, pp. 179–211, 1990.