

Understanding GMM, HMM, DNN and LSTM

Pradeep R
12th April 2019

Outline

- Introduction
- K-Means Clustering
- Gaussian Model
- Need for GMM
- Need for HMM
 - Implementation
- Conclusion

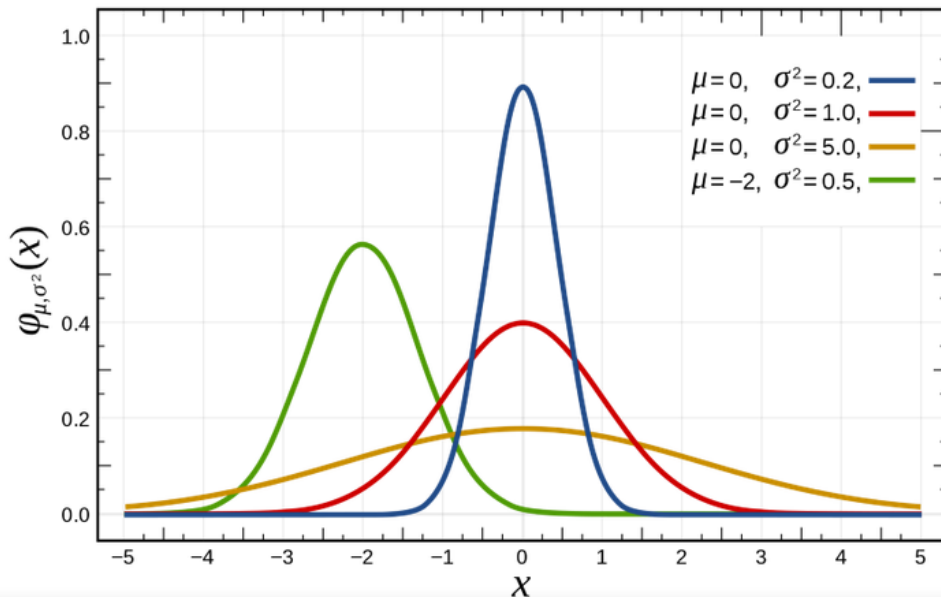
Introduction

- Applications:
 - Hand Character Recognition
 - Spoken Digit Recognition
 - Weather Prediction
 - Human activity detection
- Central limit theorem states that “when we add large number of independent random variables, irrespective of the original distribution of these variables, their normalized sum tends towards a Gaussian distribution.”
- For example, the distribution of total distance covered in an random walk tends towards a Gaussian probability distribution.

Why Gaussian?

- Gaussian remain as a Gaussian even after transformation
 - Product of two Gaussian is a Gaussian
 - Sum of two independent Gaussian random variables is a Gaussian
 - Convolution of Gaussian with another Gaussian is a Gaussian
 - Fourier transform of Gaussian is a Gaussian
- Simpler
 - The entire distribution can be specified using just two parameters- mean and variance

Variants of Gaussian



$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(\mu - x)^2}{2\sigma^2}}$$

[UniVariate Gaussian](#)

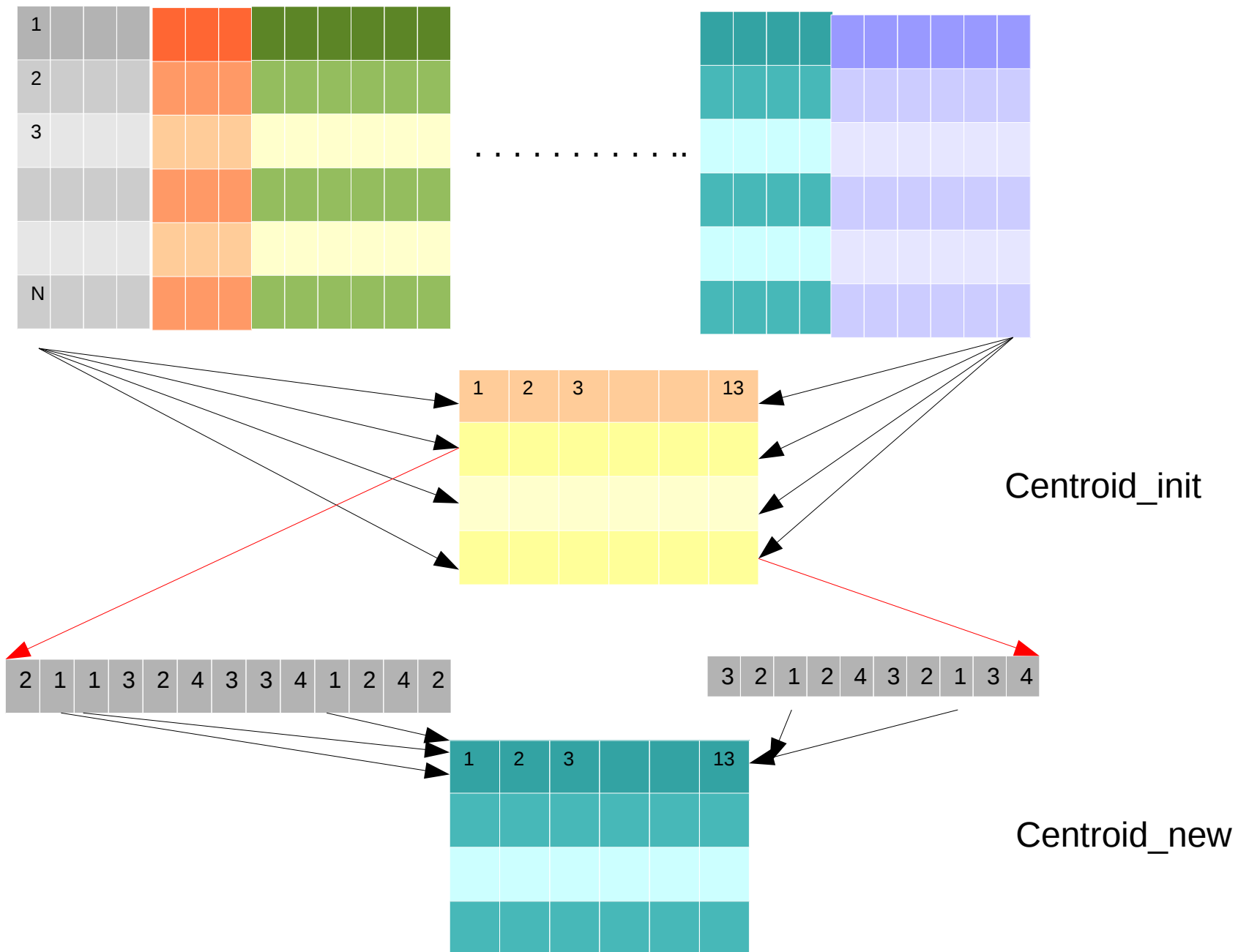
$$p(\mathbf{x}; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right)$$

[MultiVariate Gaussian](#)

$$f(\mathbf{x}) = \sum_{i=1}^N p_i \cdot \frac{1}{(2\pi)^{n/2} \sqrt{\det(\Sigma_i)}} \exp\left(-\frac{1}{2} (\mathbf{x} - \mathbf{m}_i)^T \Sigma_i^{-1} (\mathbf{x} - \mathbf{m}_i)\right)$$

[Gaussian Mixture Model](#)

K-Means Clustering - Idea

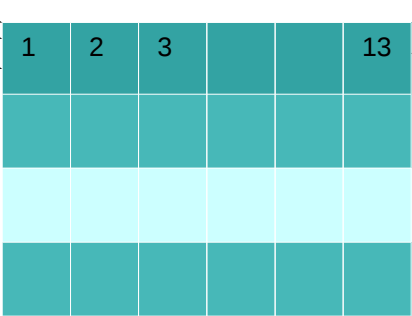


Error
10e5
10e4
0

$$\text{Error} = \text{Sum}(\text{Sum}(\text{centroid_new} - \text{centroid_old})^2)$$

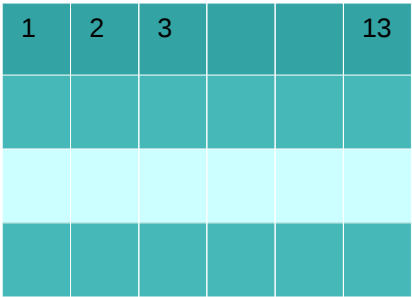


Mean vector



4*13

Variance vector



4*13

Weight vector



4*1

GMM Training: Expectation Maximization Algorithm

- Initialization: k-means clustering

- E-step

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

- M-step

$$N_k = \sum_{n=1}^N \gamma(z_{nk})$$

$$\begin{aligned} \boldsymbol{\mu}_k^{\text{new}} &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n & \pi_k^{\text{new}} &= \frac{N_k}{N} \\ \boldsymbol{\Sigma}_k^{\text{new}} &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}})^T \end{aligned}$$

- Evaluate log likelihood and check for convergence

$$\ln p(\mathbf{X} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

Implementation-Initial Setup



	1	2	3	4	5	6	7	8	9	10	11	12	13	
f1														☆
f2														#
fn														✓
f1														#
f2														✓
fp														
f1														#
f2														☆
fq														✓

Data 1

Data 2

⋮
⋮
⋮
⋮

Data 40

	1	2	3	4	5	6	7	8	9	10	11	12	13
✓													
☆													
#													

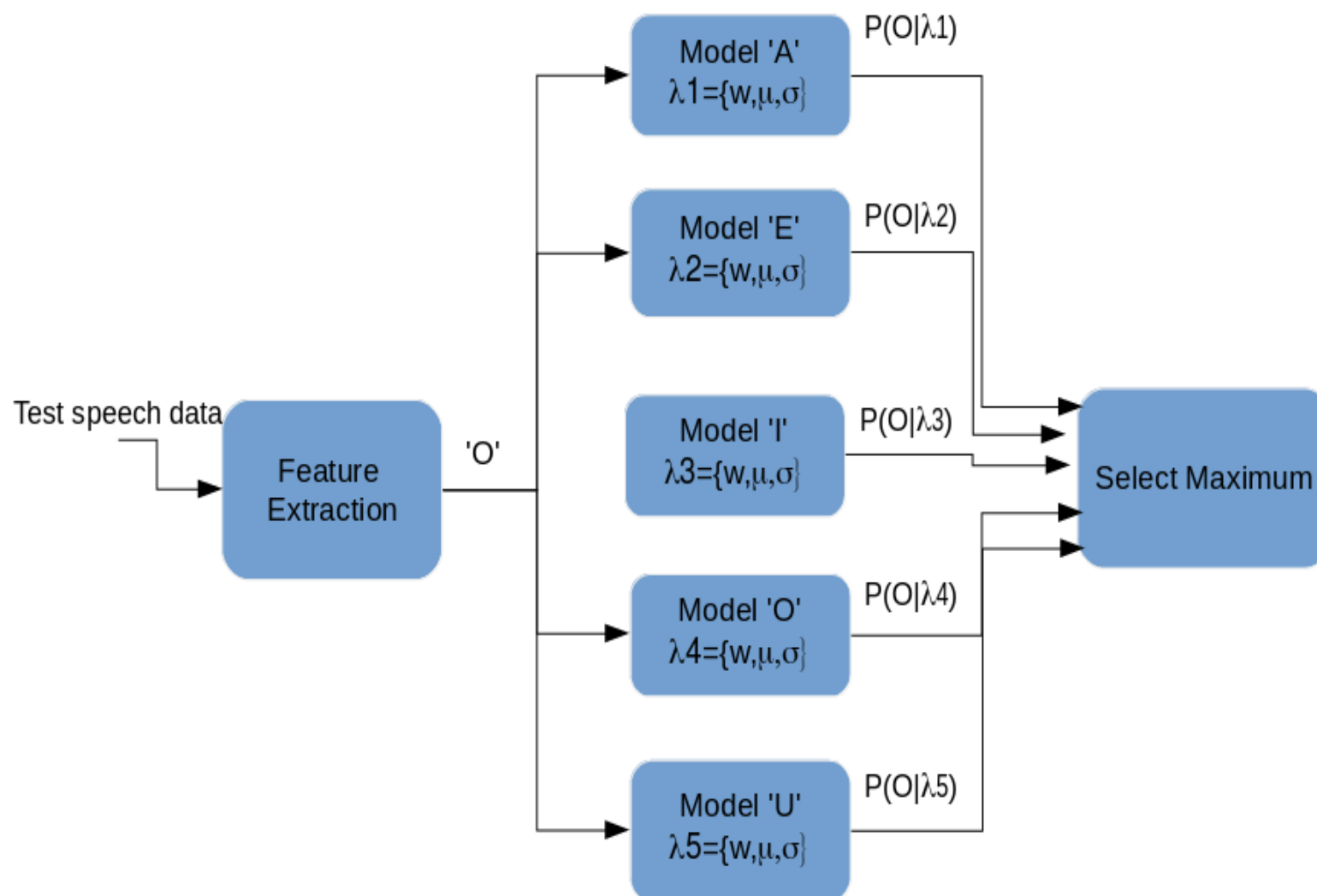
Number of mixtures * coeff per frame

	1	2	3	4	5	6	7	8	9	10	11	12	13
✓													
☆													
#													

Priori weights
✓ ✓ .. ✓
☆ .. ☆
#..#

Number of mixtures * 1

Testing



Drawbacks of GMM

- Cannot alone perform time series prediction
- It is necessary to evaluate probabilities along with time in many practical applications.
 - Ex 01: GMM build for word 'Kamal' and 'Kalam' gathers no sequential information
 - GMM - 'b' and GMM - 'd' may have false substitutions
- Hence it is necessary to capture temporal/sequential information along with spatial information.

Hidden Markov Models

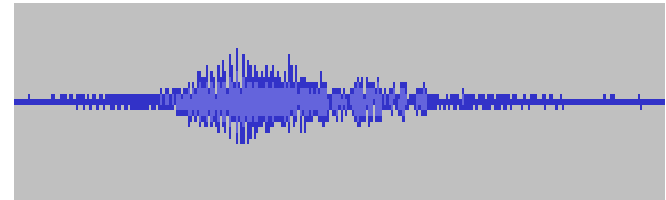
- Terminologies:
 - State
 - State Transition
 - Emission Probability

$$b_j(o_t) = \frac{1}{\sqrt{(2\pi)^n |\Sigma_j|}} e^{-\frac{1}{2}(\mathbf{o}_t - \boldsymbol{\mu}_j)' \boldsymbol{\Sigma}_j^{-1} (\mathbf{o}_t - \boldsymbol{\mu}_j)}$$

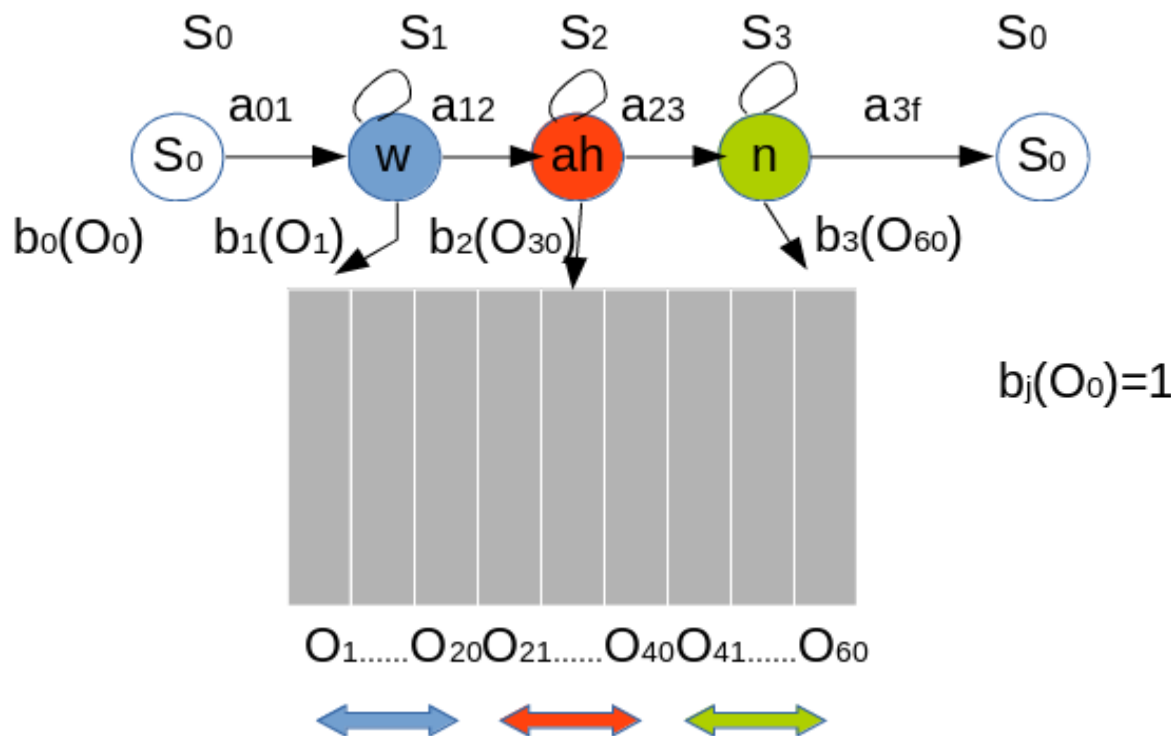
$$\hat{\boldsymbol{\mu}}_j = \frac{1}{T} \sum_{t=1}^T \mathbf{o}_t$$

$$\hat{\boldsymbol{\Sigma}}_j = \frac{1}{T} \sum_{t=1}^T (\mathbf{o}_t - \boldsymbol{\mu}_j)(\mathbf{o}_t - \boldsymbol{\mu}_j)'$$

Illustration of phone level modeling



- Train utterance.wav
 - Transcription: One
 - Phonetic equivalence: One -- w ah n

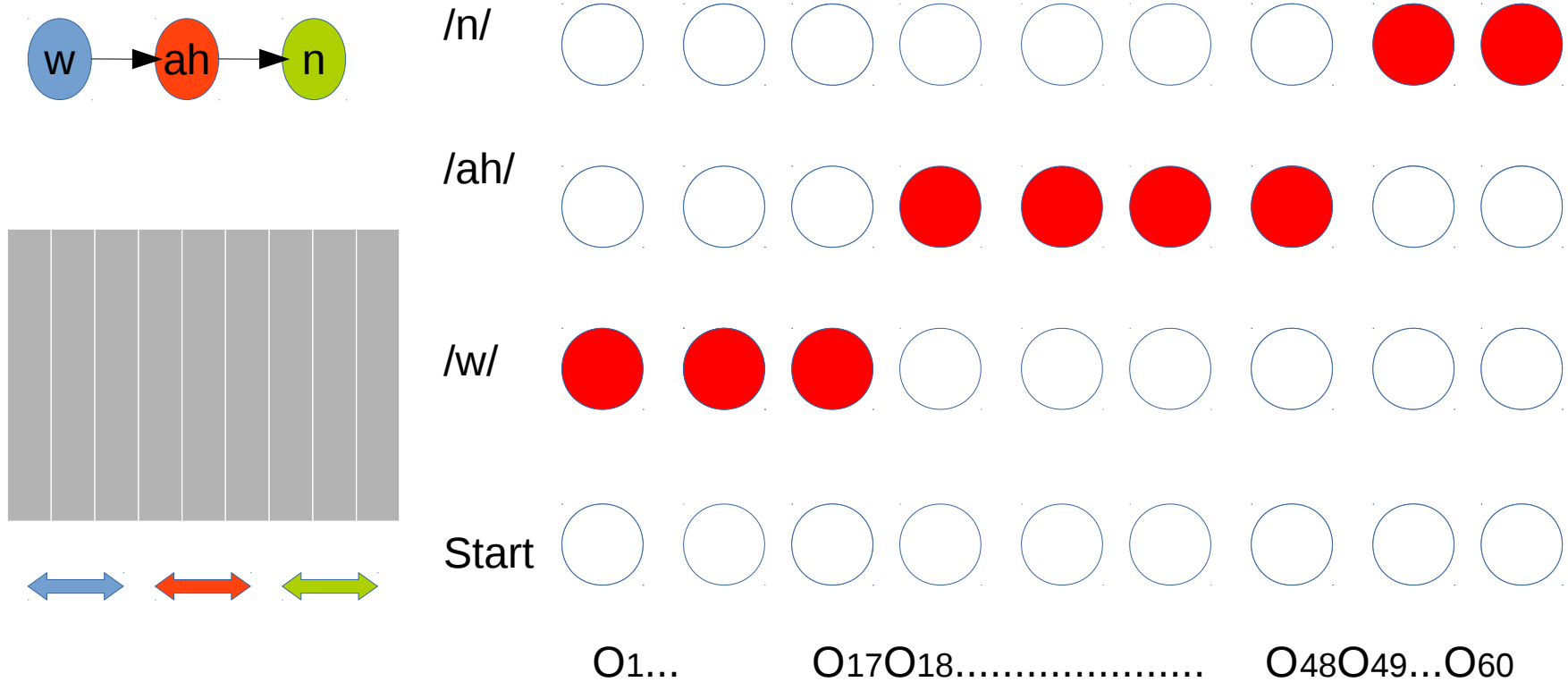


$$b_j(o_t) = \frac{1}{\sqrt{(2\pi)^n |\Sigma_j|}} e^{-\frac{1}{2}(o_t - \mu_j)' \Sigma_j^{-1} (o_t - \mu_j)}$$

$$\hat{\mu}_j = \frac{1}{T} \sum_{t=1}^T o_t$$

$$\hat{\Sigma}_j = \frac{1}{T} \sum_{t=1}^T (o_t - \mu_j)(o_t - \mu_j)'$$

Expected Outcome



0 ; t=0 & j!=initial state

$\alpha_j(t) = 1$; t=0 & j=initial state

$\sum(\alpha_i(t-1) * a_{ij}) * b_j(O_t)$

- $O^4 = \{O_1, O_{25}, O_{46}, O_0\}$
t: 1 2 3 4



$\alpha_j(0)$ $\alpha_j(1)$ $\alpha_j(2)$ $\alpha_j(3)$ $\alpha_j(T=4)$

S0	S1	S2	S3
1	0	0	0
0.2	0.3	0.1	0.4
0.2	0.5	0.2	0.1
0.7	0.1	0.1	0.1

[A]

O0	O1	O25	O46	O0
1	0	0	0	0
0	0.3	0.1	0.4	0.2
0	0.1	0.7	0.1	0.1
0	0.5	0.1	0.2	0.2

[B]

Assumption:

- Let the HMM be at state S1 at t=0
- $\alpha_1(0) = 1$
- Find $P(O^4 | \theta)$

S_3	0	0	0	0	0.0018
S_2	0	0.09	0.0052	0.0052	0
S_1	1	0.01	0.0217	0.0005	0
S_0	0	0.2	0.0057	0.001	0

Use **Viterbi decoding** to find the hidden state sequence that generated the particular observation

$$P(O^4|\theta)=\alpha_0(T)=0.0018$$

$$P(O^4)=\beta_i(0)$$

$t=1$

$$\alpha_0(1)=\{\alpha_0(0)a_{00}+\alpha_1(0)a_{10}+\alpha_2(0)a_{20}+\alpha_3(0)a_{30}\} * b_0(O_1)=0$$

$$\alpha_1(1)=\{\alpha_0(0)a_{01}+\alpha_1(0)a_{11}+\alpha_2(0)a_{21}+\alpha_3(0)a_{31}\} * b_1(O_2)=0.09$$

$$\alpha_2(1)=\{\alpha_0(0)a_{02}+\alpha_1(0)a_{12}+\alpha_2(0)a_{22}+\alpha_3(0)a_{32}\} * b_2(O_3)=0.02$$

$$\alpha_3(1)=\{\alpha_0(0)a_{03}+\alpha_1(0)a_{13}+\alpha_2(0)a_{23}+\alpha_3(0)a_{33}\} * b_3(O_4)=0.05$$

$$\text{In matrix form } [0 \ 1 \ 0 \ 0] * [A] * [B(:,2)] = [0.2 \ 0.01 \ 0.09 \ 0]$$

$$\hat{a}_{ij} = \frac{\sum_{t=1}^{T-1} \gamma_t(i,j)}{\sum_{t=1}^{T-1} \gamma_t(i)}$$

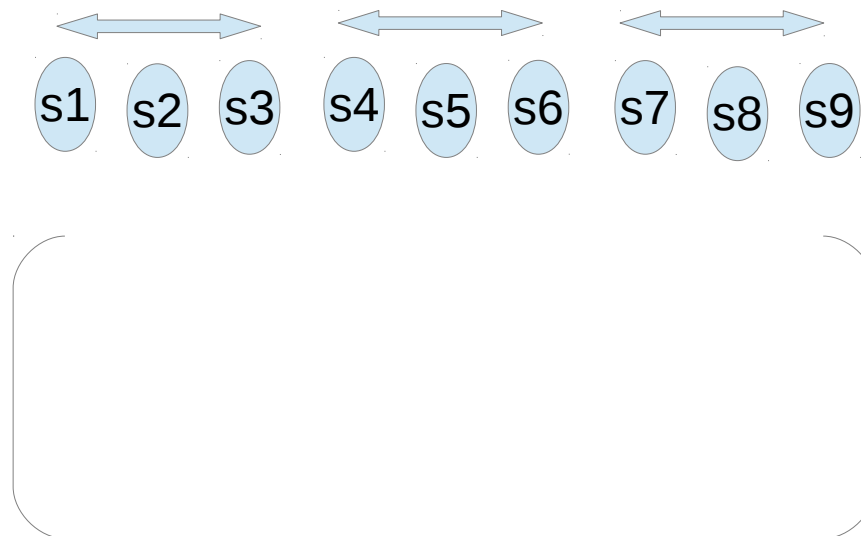
$$\gamma_i(t) = P(X_t = i|Y,\theta) = \frac{\alpha_i(t)\beta_i(t)}{\sum_{j=1}^N \alpha_j(t)\beta_j(t)}$$

$$\hat{b}_j(o_k) = \frac{\sum_{t:O_t=o_k} \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)}$$

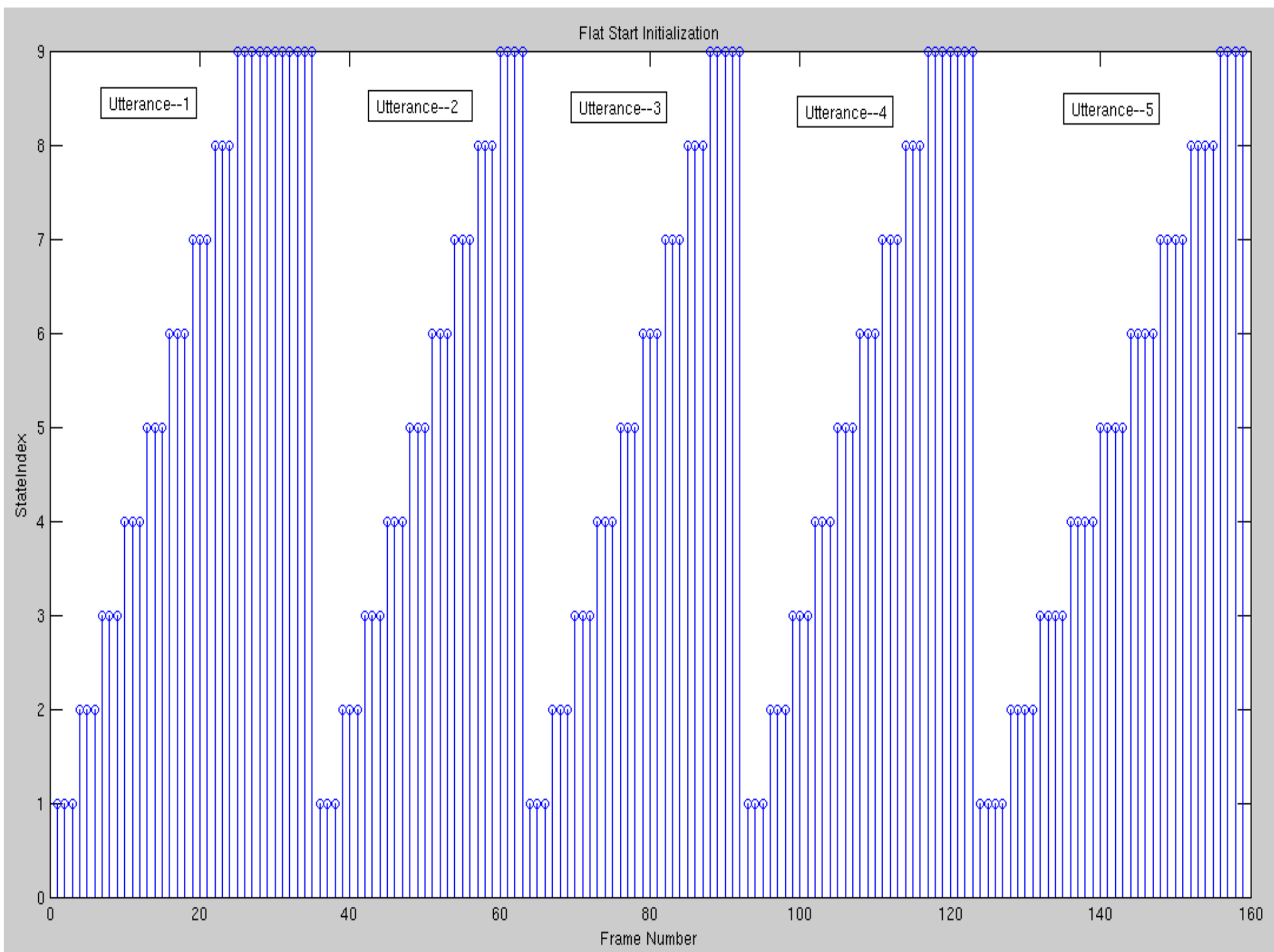
$$\xi_{\mathbf{t}(\mathbf{i},\mathbf{j})} = \frac{\alpha_i(t)a_{ij}\beta_j(t+1)b_j(y_{t+1})}{\sum_{k=1}^N \sum_{l=1}^N \alpha_k(t)a_{kl}\beta_l(t+1)b_l(y_{t+1})} = \frac{\alpha_i(t)a_{ij}\beta_j(t+1)b_j(y_{t+1})}{\sum_{k=1}^N \alpha_k(t)\beta_k(t)}$$

Illustration

- Data:
 - Ten words of the utterance one
 - One : w ah n
- Build 3 state HMM per phone

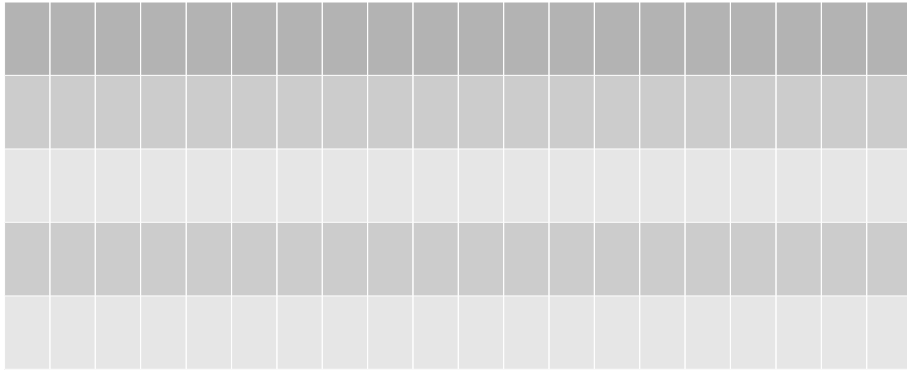


MFCC

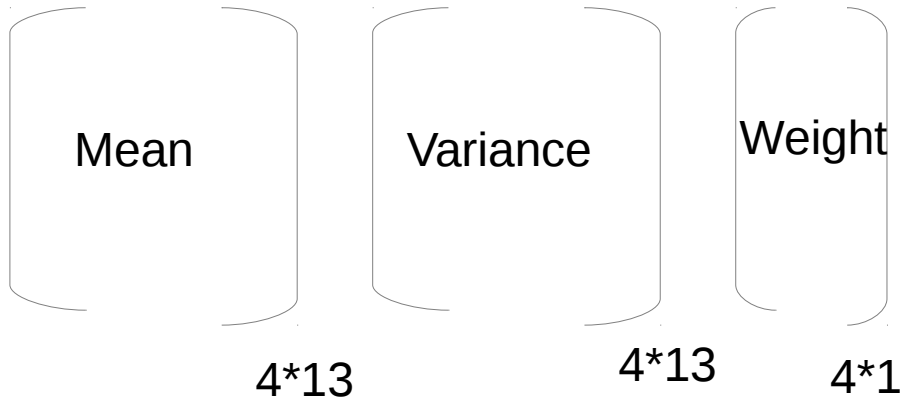


Initial setup

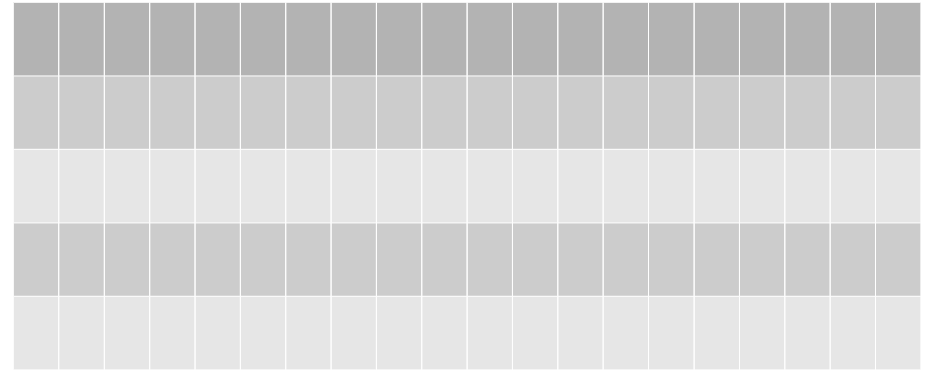
Data_S1



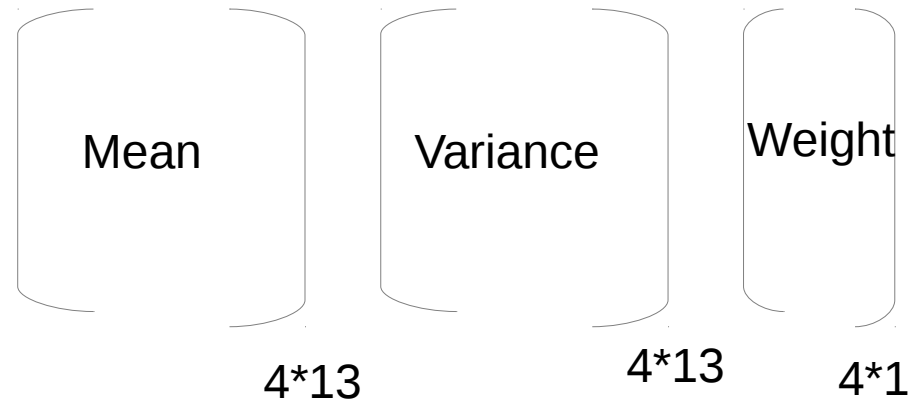
2 4 4 1 1 3 2 1 1 1 4 2 3 1 2 2 4 3 1



Data_S9



3 2 4 2 1 1 1 3 3 2 1 4 3 2 1 3 1 2 4



Forward Algorithm

1) Initialization:

$$\alpha_1(i) = \pi_i b_i(O_1), \quad 1 \leq i \leq N.$$

2) Induction:

$$\alpha_{t+1}(j) = \left[\sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(O_{t+1}), \quad 1 \leq t \leq T-1$$
$$1 \leq j \leq N.$$

3) Termination:

$$P(O|\lambda) = \sum_{i=1}^N \alpha_T(i).$$

Backward Probabilities

1) Initialization:

$$\beta_T(i) = 1, \quad 1 \leq i \leq N.$$

2) Induction:

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(O_{t+1}) \beta_{t+1}(j),$$
$$t = T-1, T-2, \dots, 1, 1 \leq i \leq N.$$

Viterbi Algorithm

1) Initialization:

$$\delta_1(i) = \pi_i b_i(O_1), \quad 1 \leq i \leq N$$

$$\psi_1(i) = 0.$$

2) Recursion:

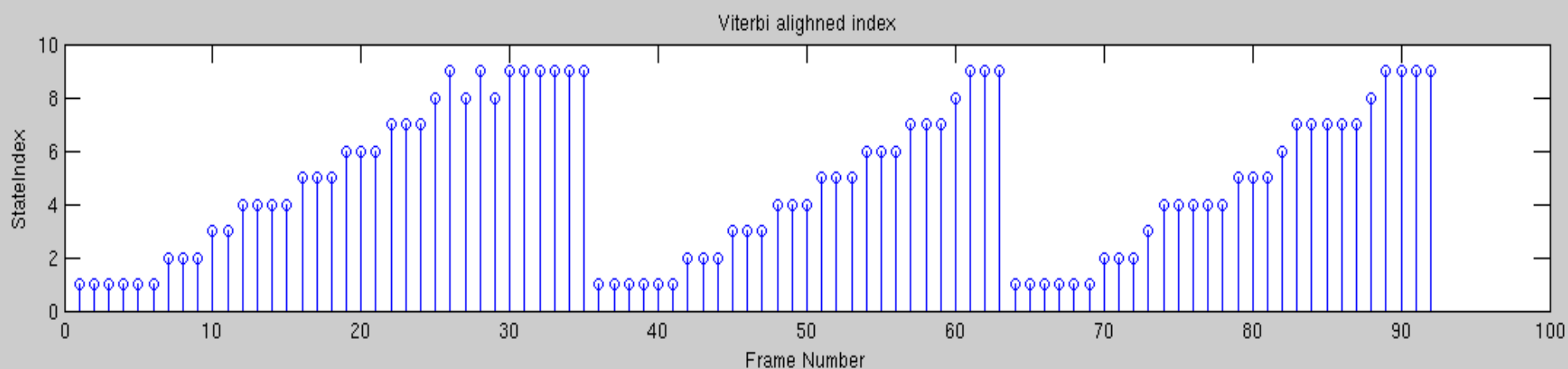
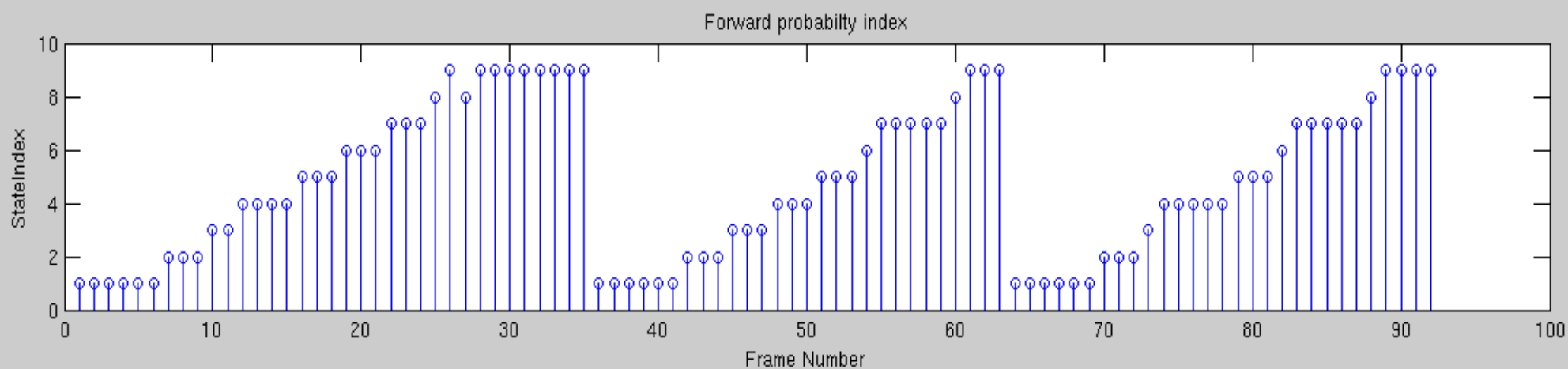
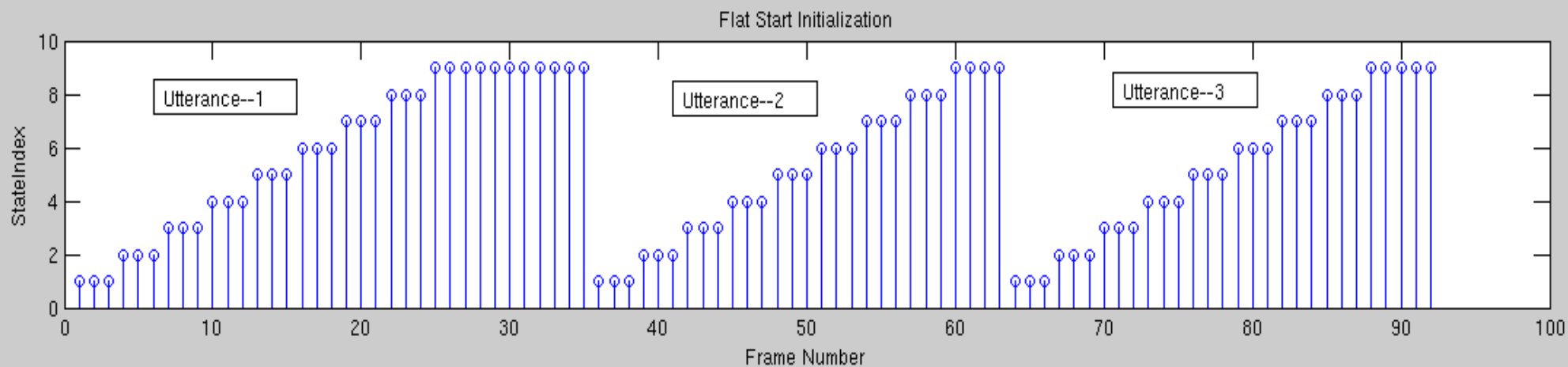
$$\delta_t(j) = \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}] b_j(O_t), \quad 2 \leq t \leq T$$
$$1 \leq j \leq N$$

$$\psi_t(j) = \operatorname{argmax}_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}], \quad 2 \leq t \leq T$$
$$1 \leq j \leq N.$$

3) Termination:

$$P^* = \max_{1 \leq i \leq N} [\delta_T(i)]$$

$$q_T^* = \operatorname{argmax}_{1 \leq i \leq N} [\delta_T(i)].$$



$$\gamma_t(i) = \frac{\alpha_t(i) \beta_t(i)}{P(O|\lambda)} = \frac{\alpha_t(i) \beta_t(i)}{\sum_{i=1}^N \alpha_t(i) \beta_t(i)}$$

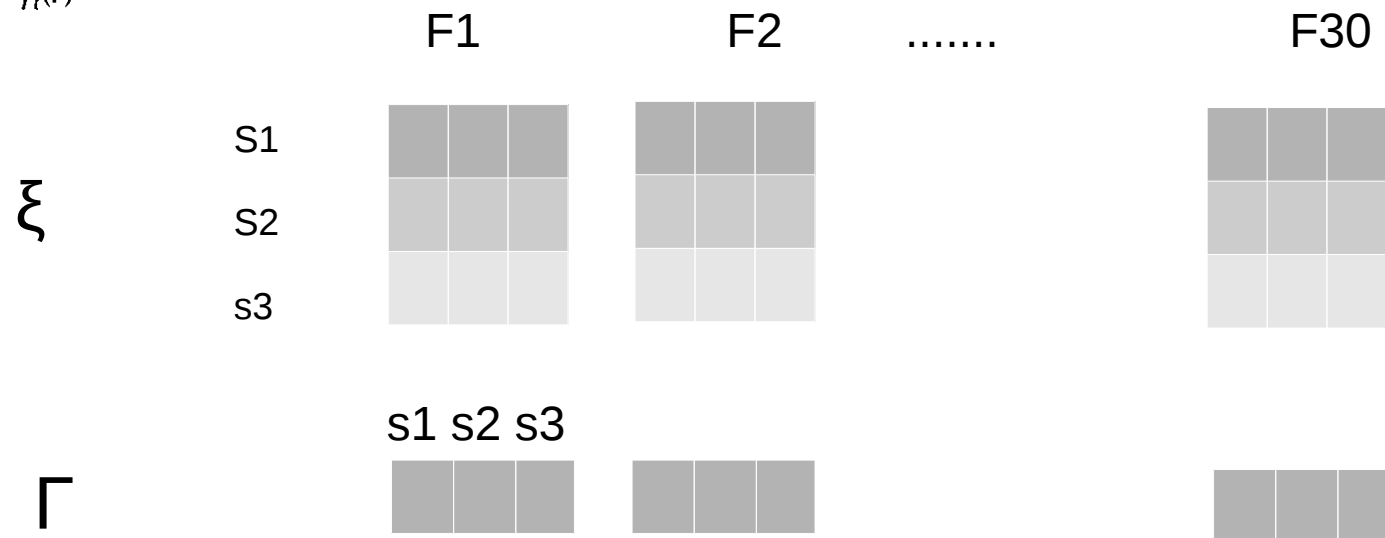
$\bar{\pi}_i$ = expected frequency (number of times) in state S_i at time ($t = 1$)

$$\bar{a}_{ij} = \frac{\text{expected number of transitions from state } S_i \text{ to state } S_j}{\text{expected number of transitions from state } S_i}$$

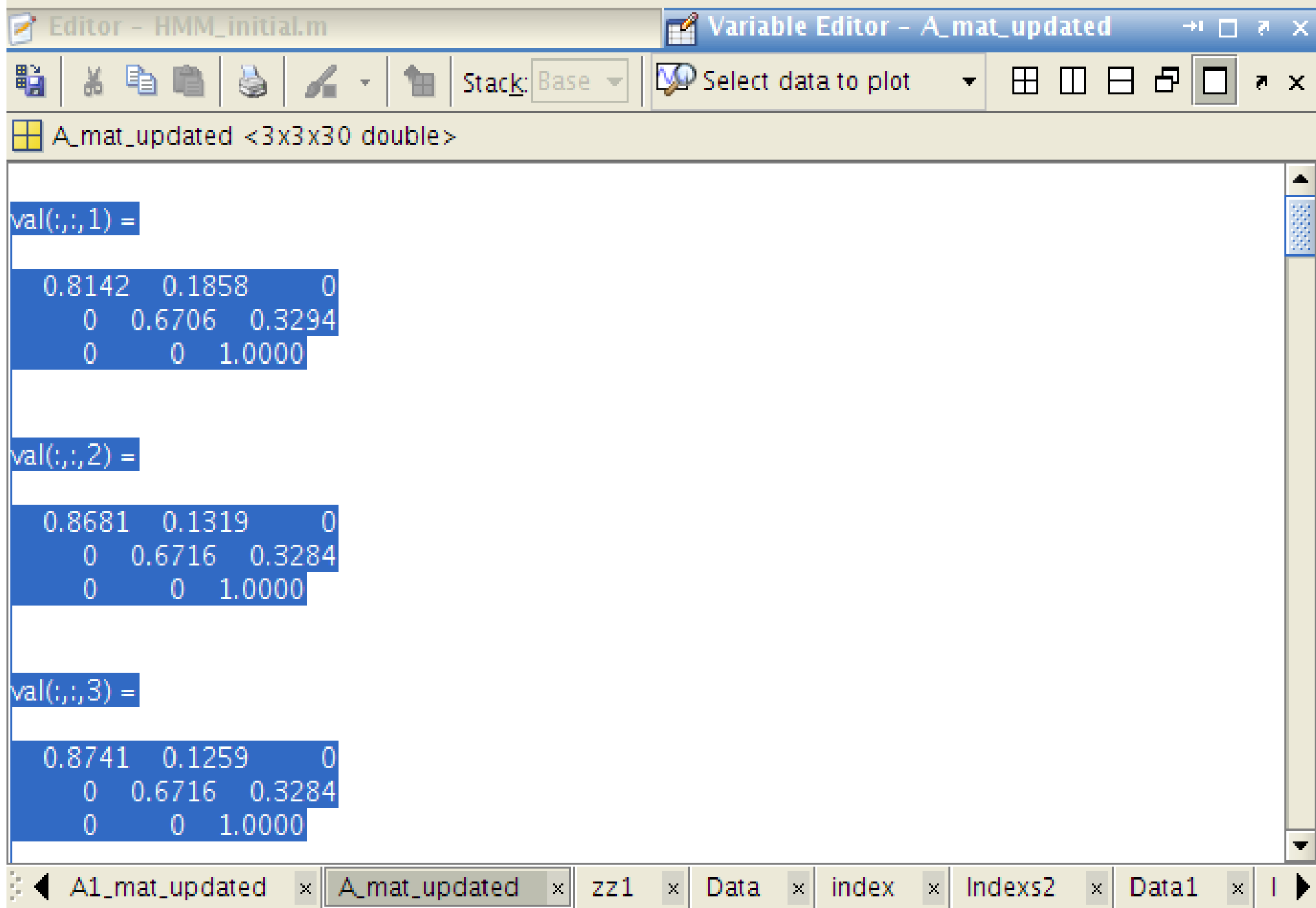
$$= \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}$$

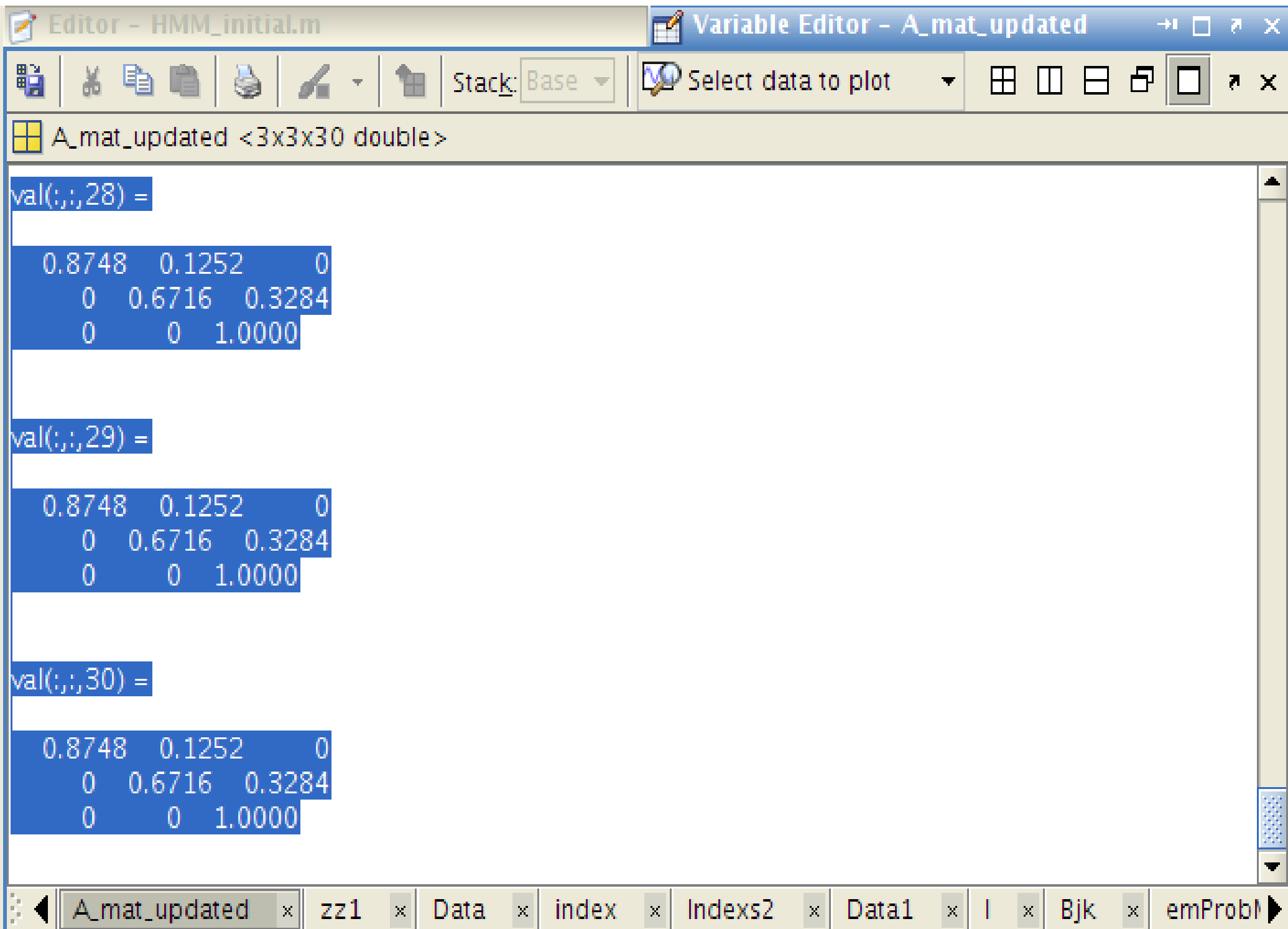
$$\xi_t(i, j) = \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{P(O|\lambda)}$$

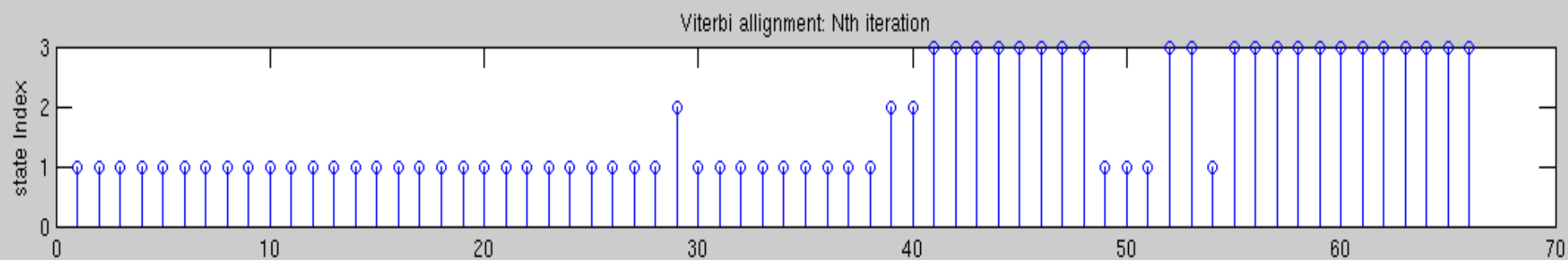
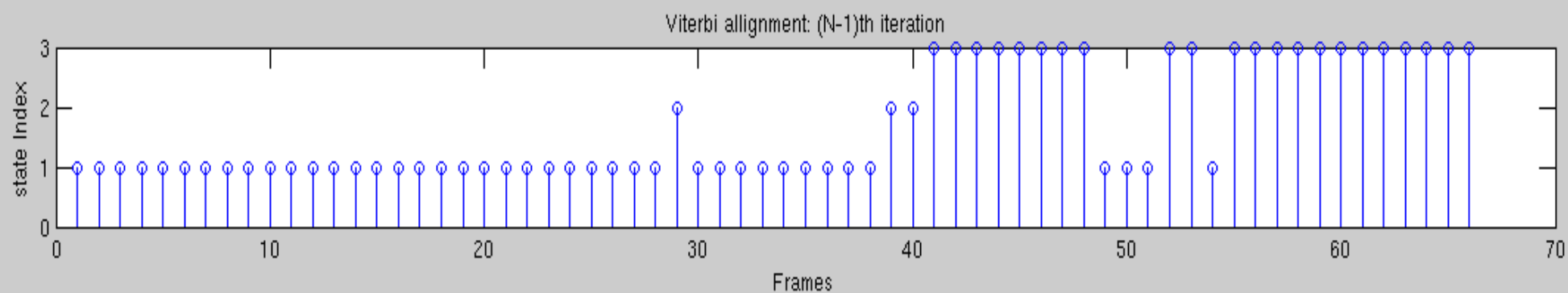
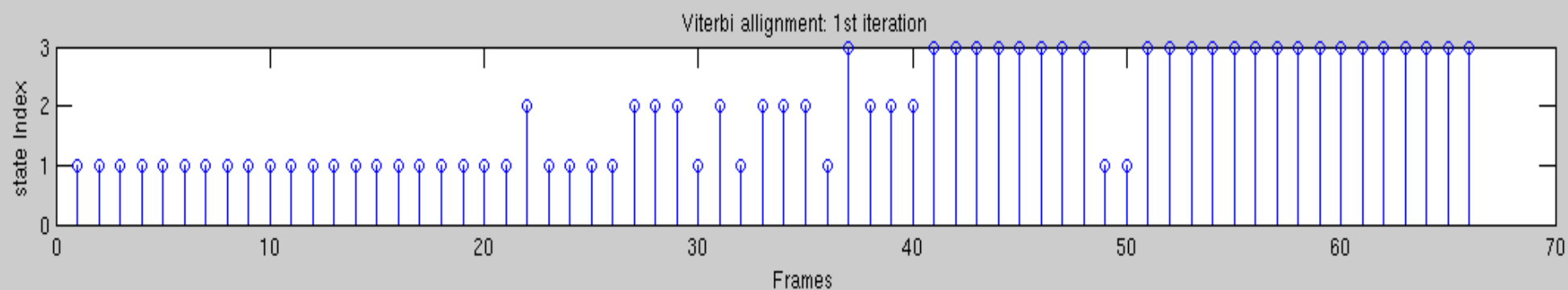
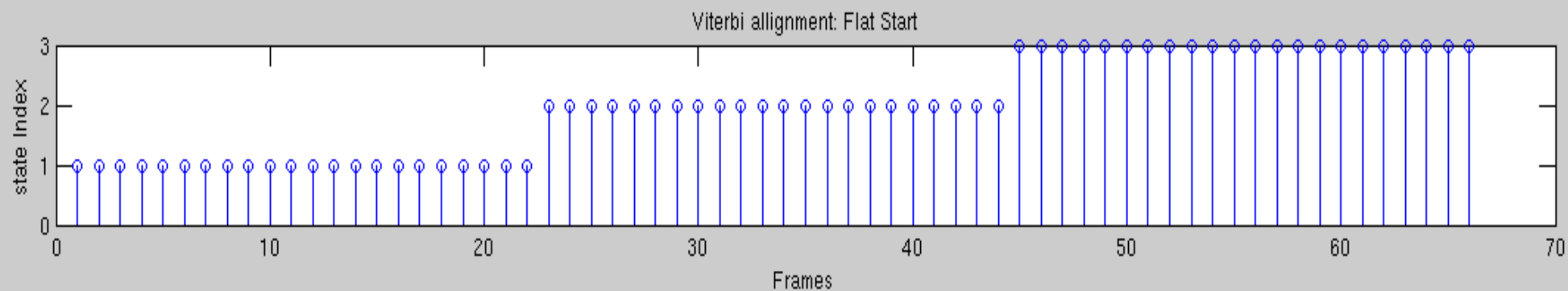
$$= \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{\sum_{i=1}^N \sum_{j=1}^N \alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}$$



$$A(S1,:) = \text{zeeta}(S1,:)/ [S1 \ S1 \ S1]$$



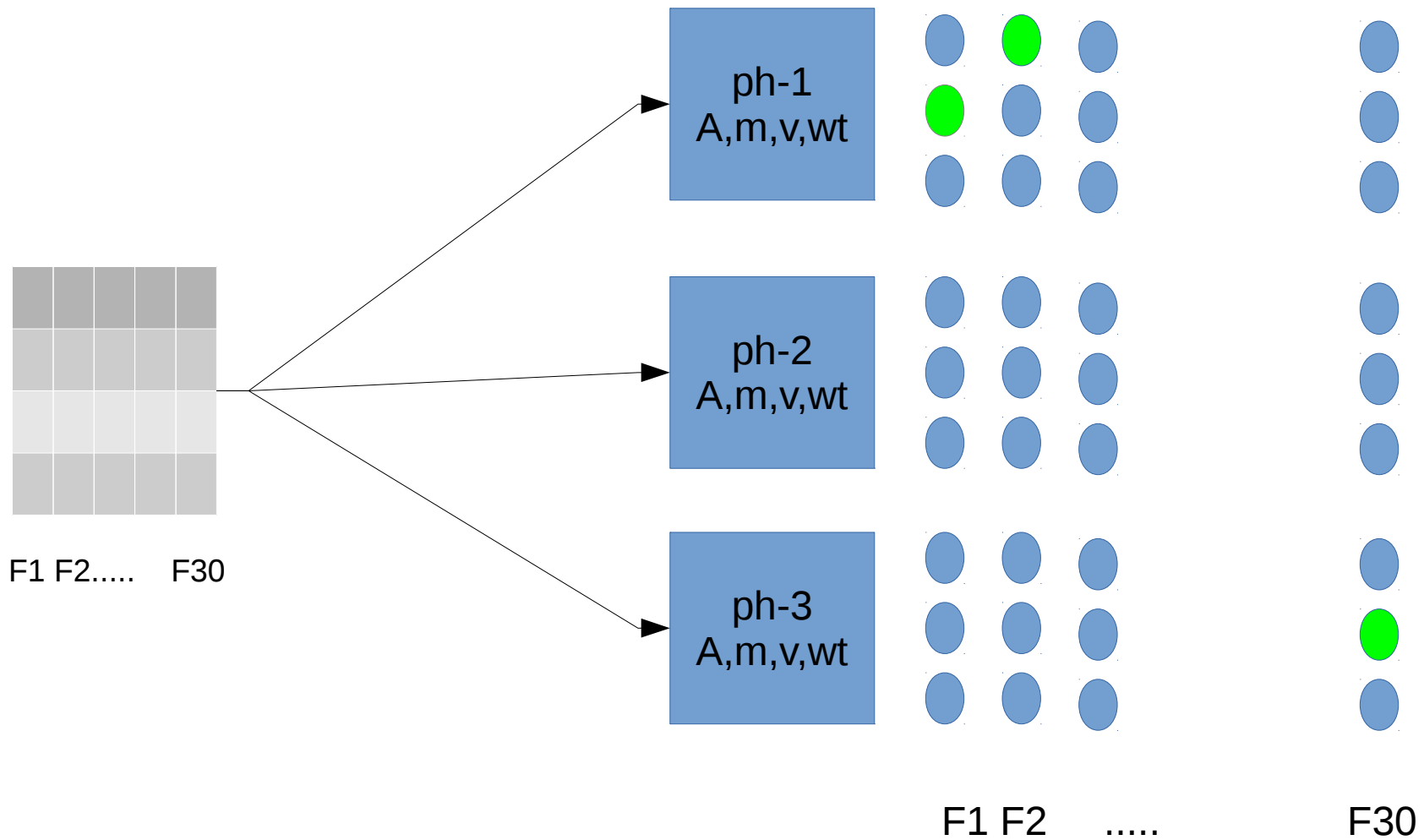




Testing HMM

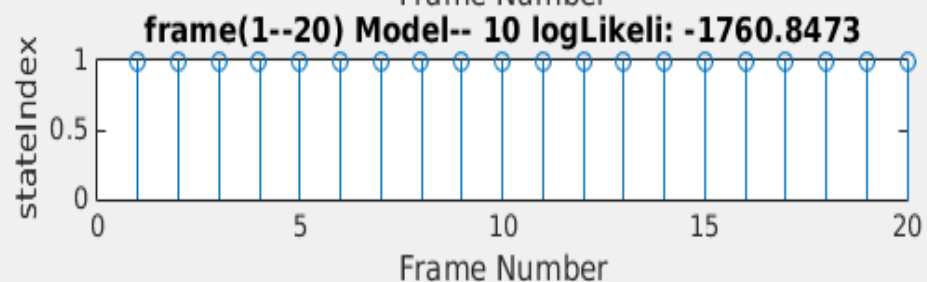
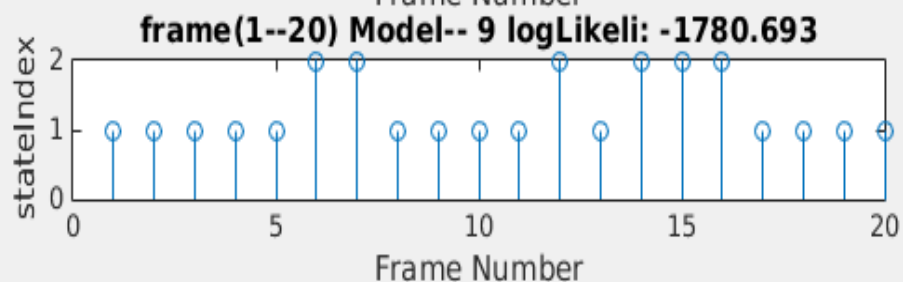
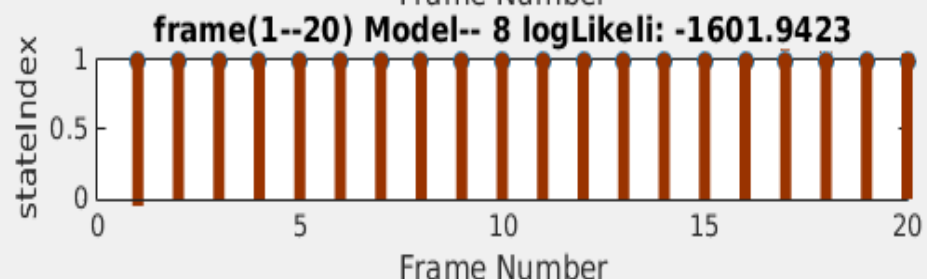
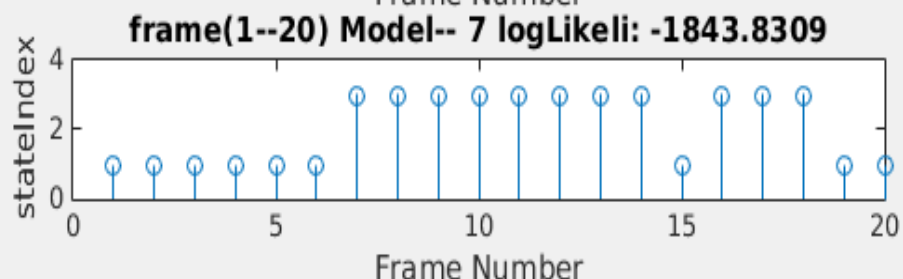
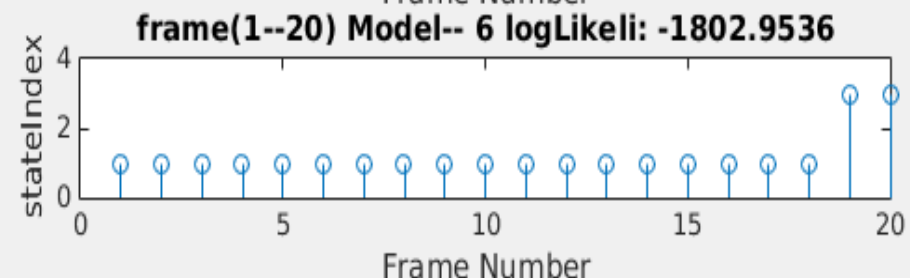
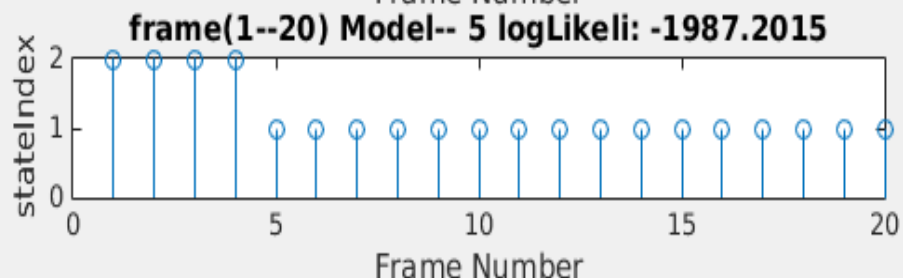
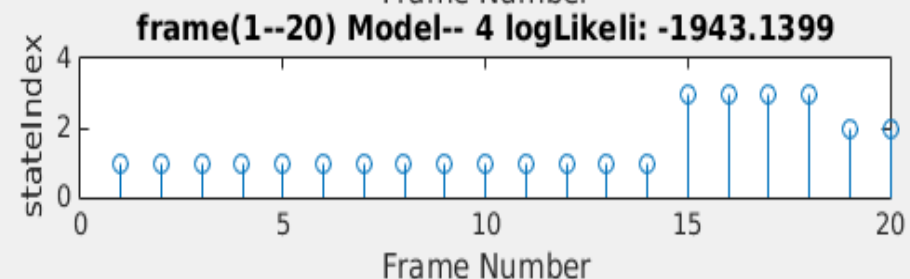
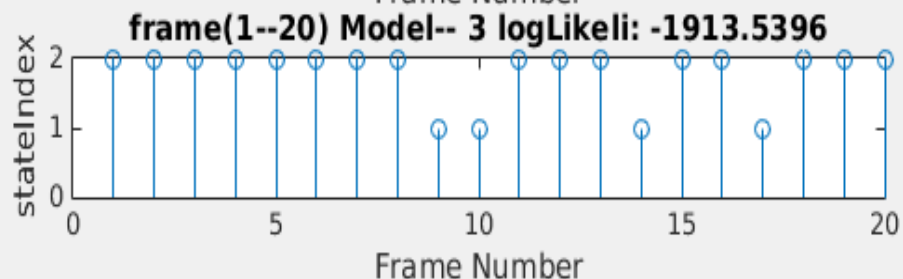
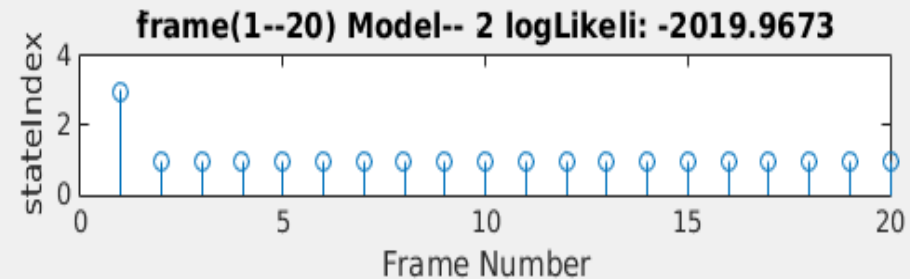
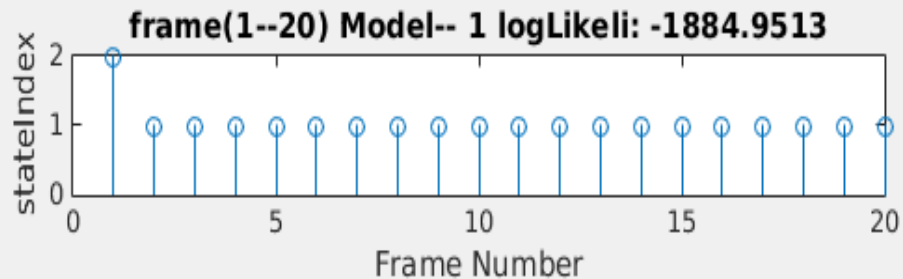
$$\left\{ \sum_{m=1}^{M_s} w_{ism} \mathcal{N}(o_s; \mu_{ism}, \Sigma_{ism}) \right\}$$

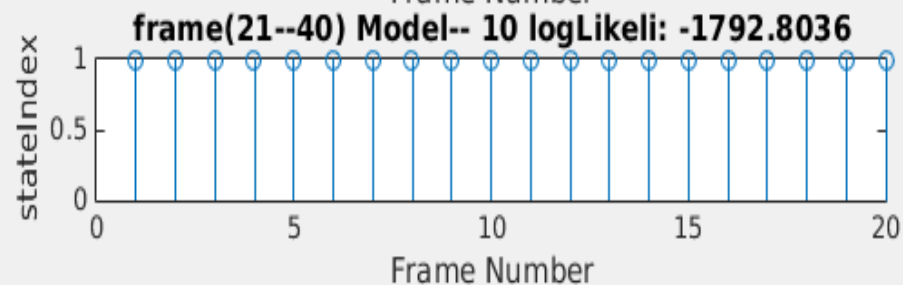
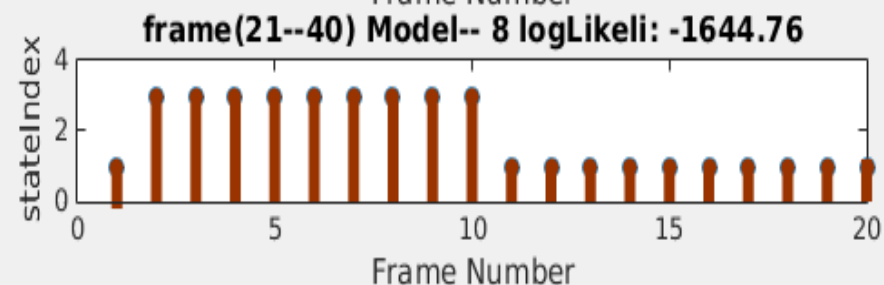
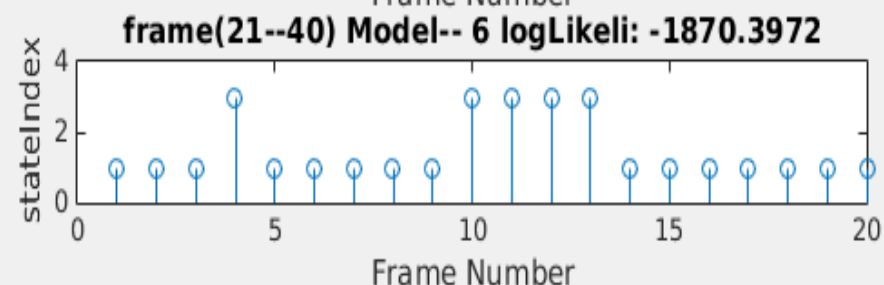
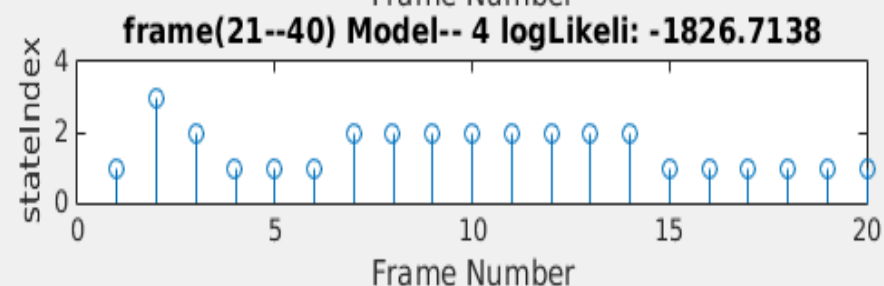
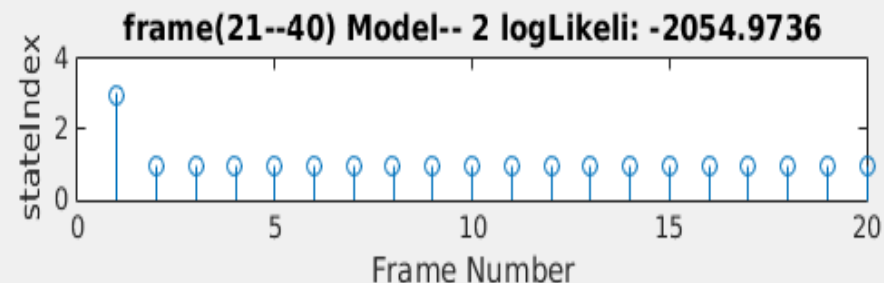
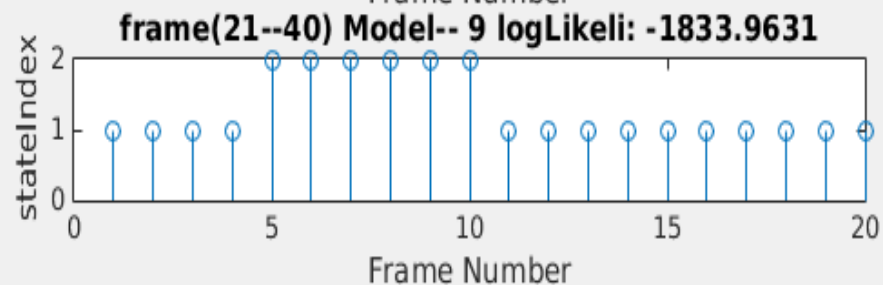
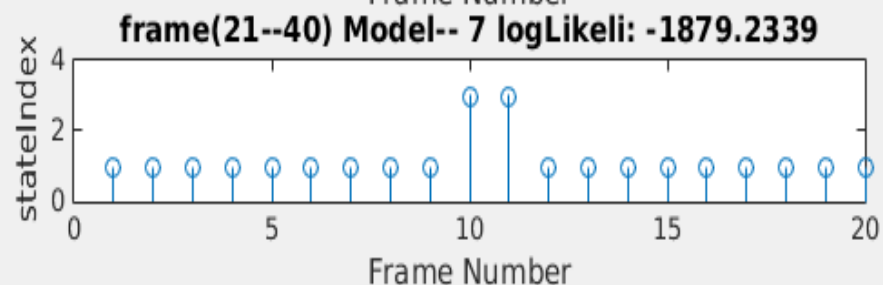
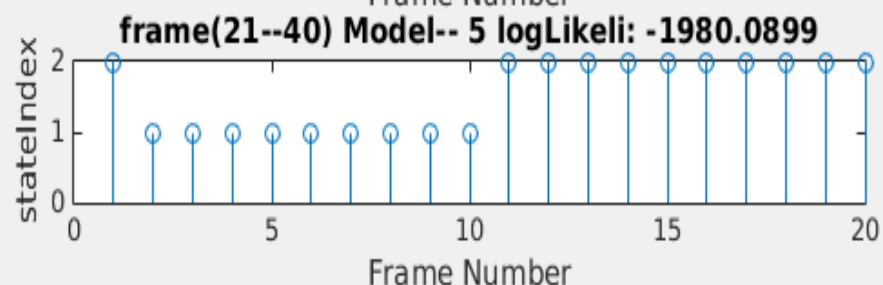
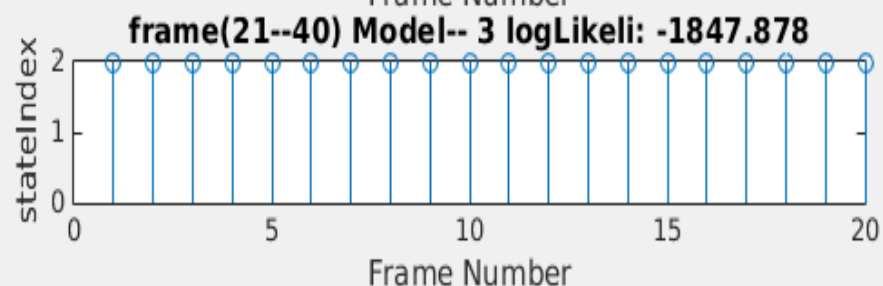
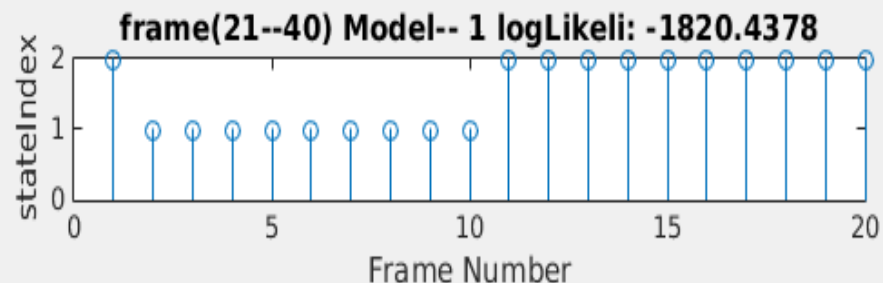
$[A]^*[b1 \ b2 \ b3]$



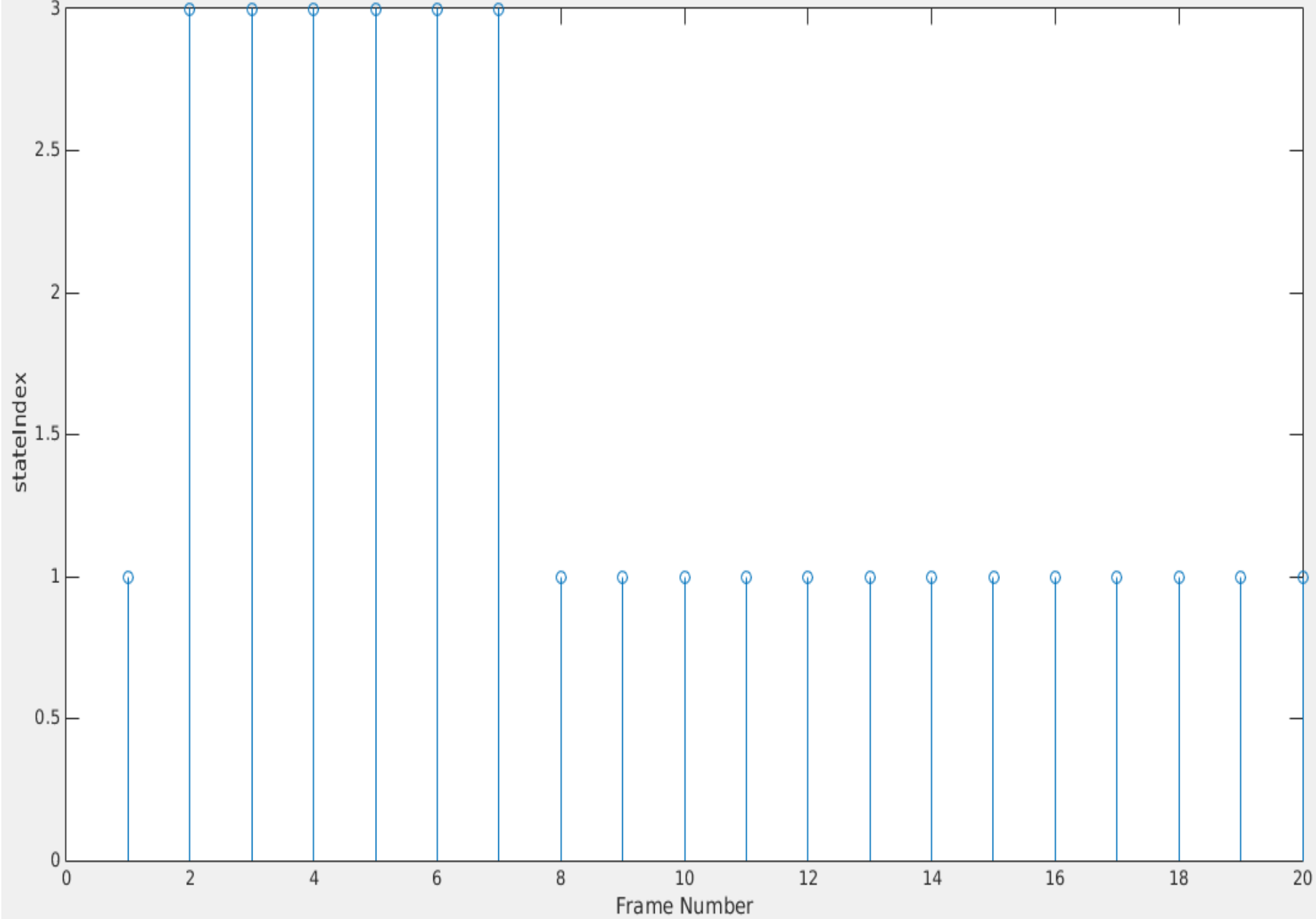
Calculating log-Likelihood from updated parameters

- Trained parameters for each state of a class:
 - Mean vector
 - Diagonal covariance matrix
 - Weighting factors for each Gaussian in a GMM
 - Updated Transition matrix ($N \times N$)
 - priors
- Finding Likelihood:
 - $\alpha(:, t) = B(:, t) .* (A' * \alpha(:, t - 1));$
 - $\text{logLikelihood} = \log(\text{sum}(\alpha(:, t)));$





model---8



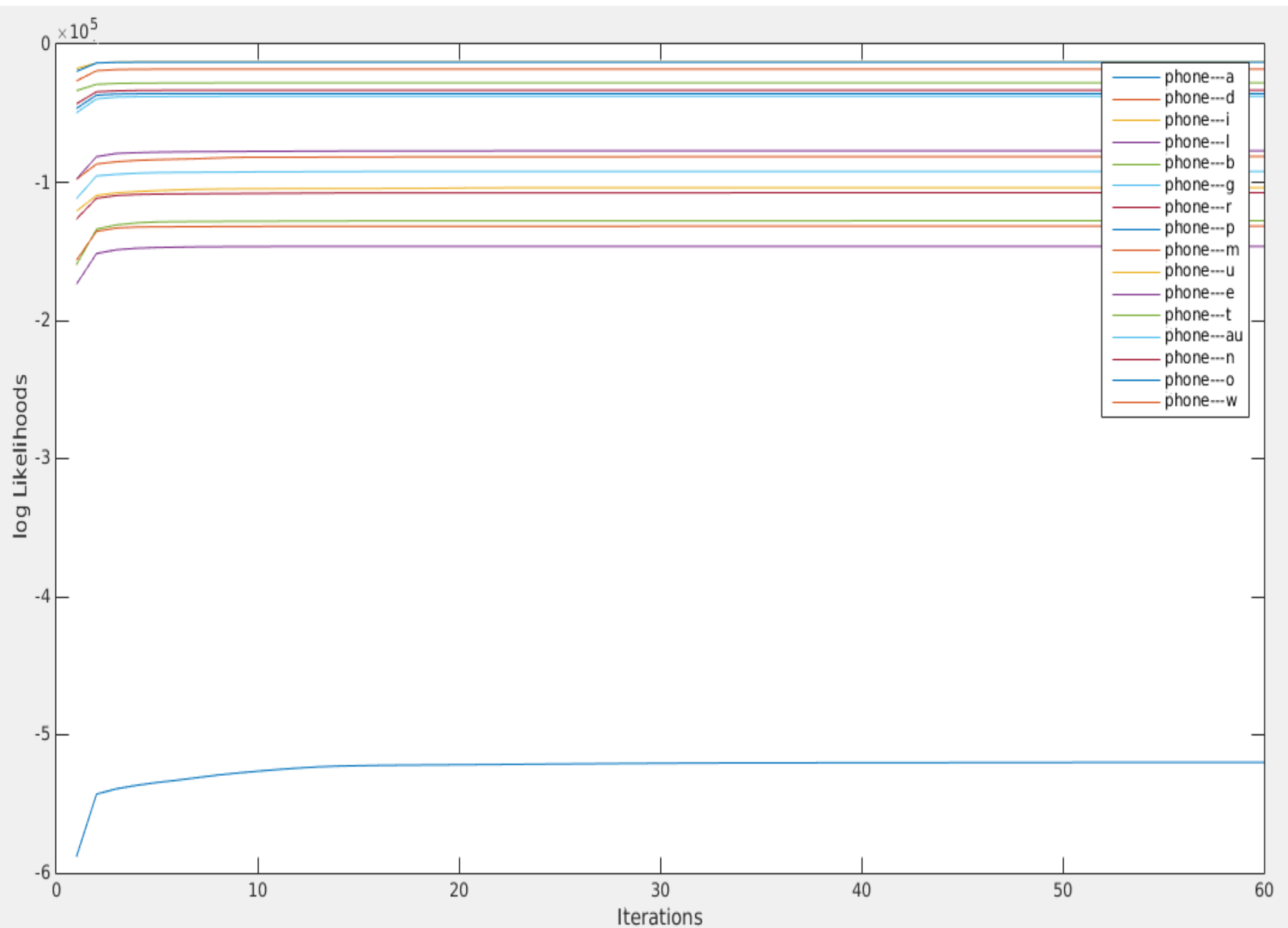
Building a 3-state Phone HMM

a d i l a b a d
a g r a
a l a p i
a m a l a p u r a m
a m b a l a
a m e t i
a u r a n g a b a d
b a n g a l o r e
b e l g a u m
b i l w a r a

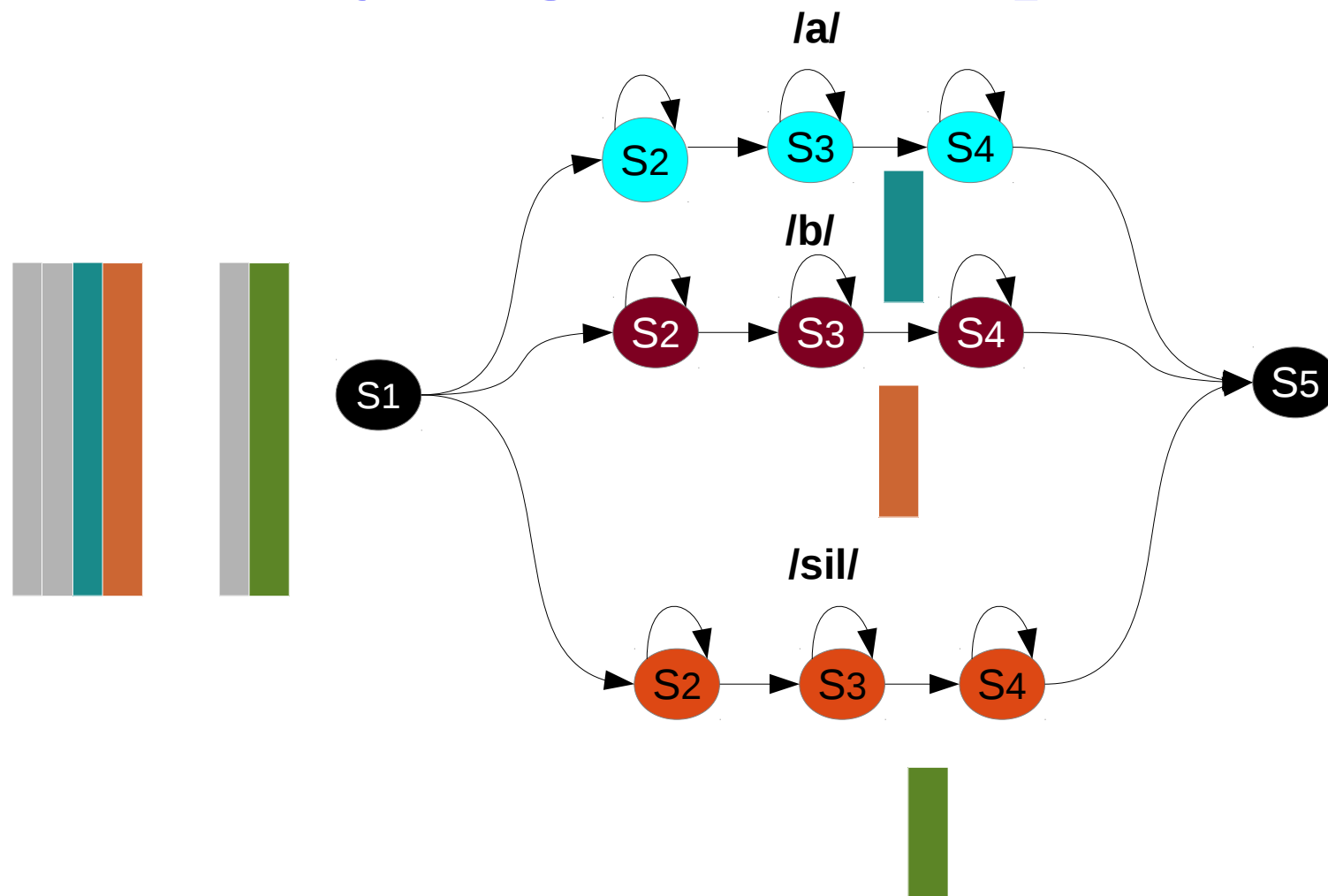
/a/	1
/d/	2
/i/	3
/l/	4
/b/	5
/g/	6
/r/	7
/p/	8
/m/	9
/u/	10
/e/	11
/t/	12
/au/	13
/n/	14
/o/	15
/w/	16

Steps:

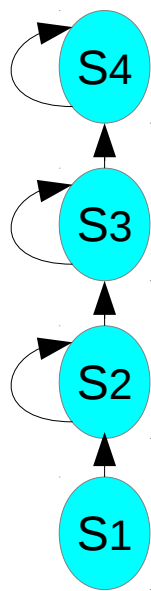
- Flat start initialization----ex: 1 2 3 4 1 5 1 2
- Compute forward and backward probabilities
- Update the HMM parameters
- Stop training if the changes in HMM parameters is negligible



Analysing the break points



- Key idea:
 - Find the hidden state index where the final emitting state starts emitting the MFCC feature vectors
 - Analyze the corresponding viterbi path and its score w.r.t all trained models



x x x

x ✓ x

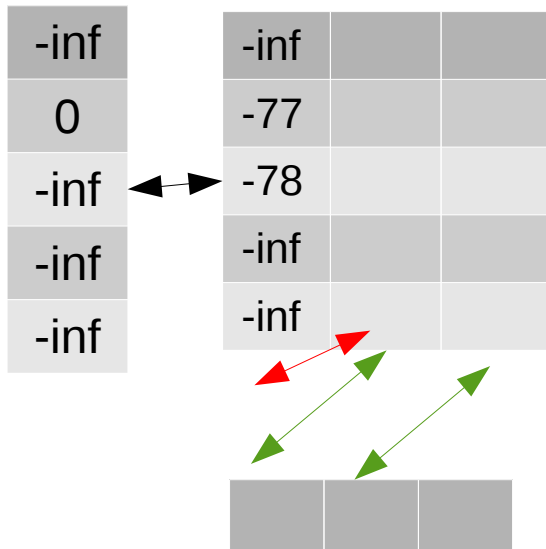
✓ x ✓



x ✓ ✓ ✓

✓ x x x

x x x x



- Identify the path with less significant Viterbi score
- Prune the Viterbi path by considering only the effect of maximum state sequence.

- $t=2$

- $$\begin{array}{|c|} \hline \delta_2(1) \\ \hline \delta_2(2) \\ \hline \delta_2(3) \\ \hline \delta_2(4) \\ \hline \delta_2(5) \\ \hline \end{array} \quad \delta_2(j) = \max \quad \begin{array}{|c|} \hline \delta_1(1).a_{1j} \\ \hline \delta_1(2).a_{2j} \\ \hline \delta_1(3).a_{3j} \\ \hline \delta_1(4).a_{4j} \\ \hline \delta_1(5).a_{5j} \\ \hline \end{array} \quad b_3(F_j)$$

At $t=1$, $p=2$

$$\begin{array}{|c|} \hline \delta_2(1) \\ \hline \delta_2(2) \\ \hline \delta_2(3) \\ \hline \delta_2(4) \\ \hline \delta_2(5) \\ \hline \end{array} = \begin{array}{|c|} \hline \delta_1(1).a_{p1} \\ \hline \delta_1(2).a_{p2} \\ \hline \delta_1(3).a_{p3} \\ \hline \delta_1(4).a_{p4} \\ \hline \delta_1(5).a_{p5} \\ \hline \end{array} \quad \begin{array}{|c|} \hline b_3(F_p) \\ \hline b_3(F_p) \\ \hline b_3(F_p) \\ \hline b_3(F_p) \\ \hline b_3(F_p) \\ \hline \end{array}$$

Token passing

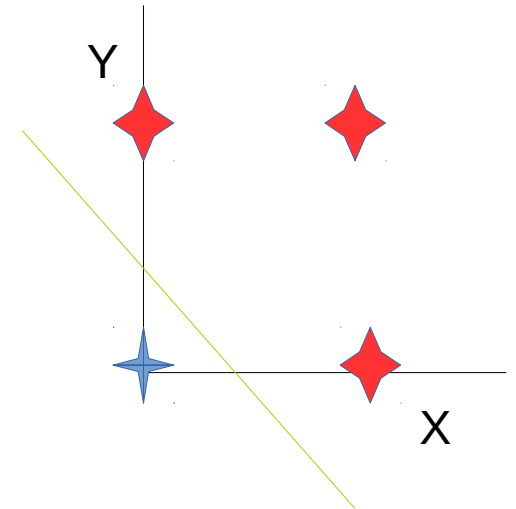
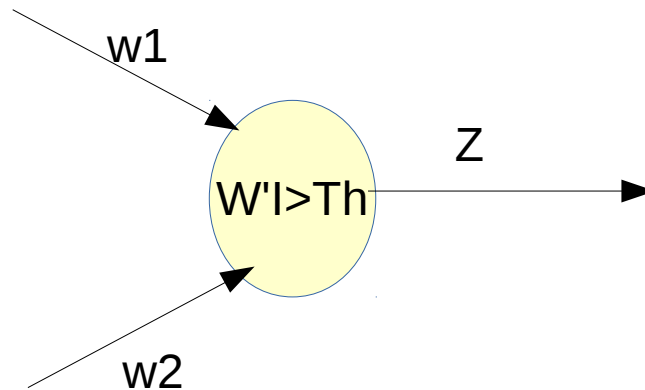
- Keep the best token arriving at each state
- Propagate these tokens to the next states
- Find the probability of the most likely state alignment

Practical Implementation of ANNs for Speech Applications

Introduction

OR-Gate

X	Y	Z
0	0	0
0	1	1
1	0	1
1	1	1



$$\begin{aligned} 0*w1+0*w2 &< Th \\ 0*w1+1*w2 &\geq Th \\ 1*w1+0*w2 &\geq Th \\ 1*w1+1*w2 &\geq Th \end{aligned}$$

$$\begin{aligned} Th &> 0 \\ W2 &\geq Th \\ W1 &\geq Th \\ w1+w2 &\geq Th \end{aligned}$$

$$\begin{aligned} y &= mx + c \\ //m &= w2/w1; \\ //c &= Th; \end{aligned}$$

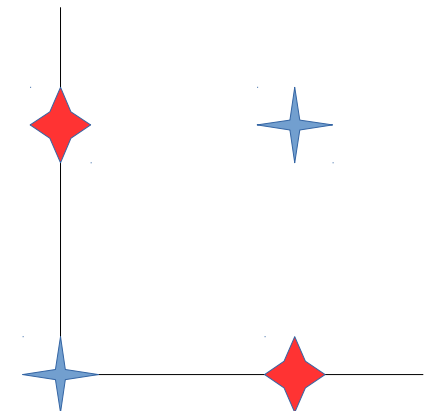
XOR-Gate

X	Y	Z
0	0	0
0	1	1
1	0	1
1	1	0

$$\text{Soln : } Th = 0.1; W = [0.15 \ 0.2]$$

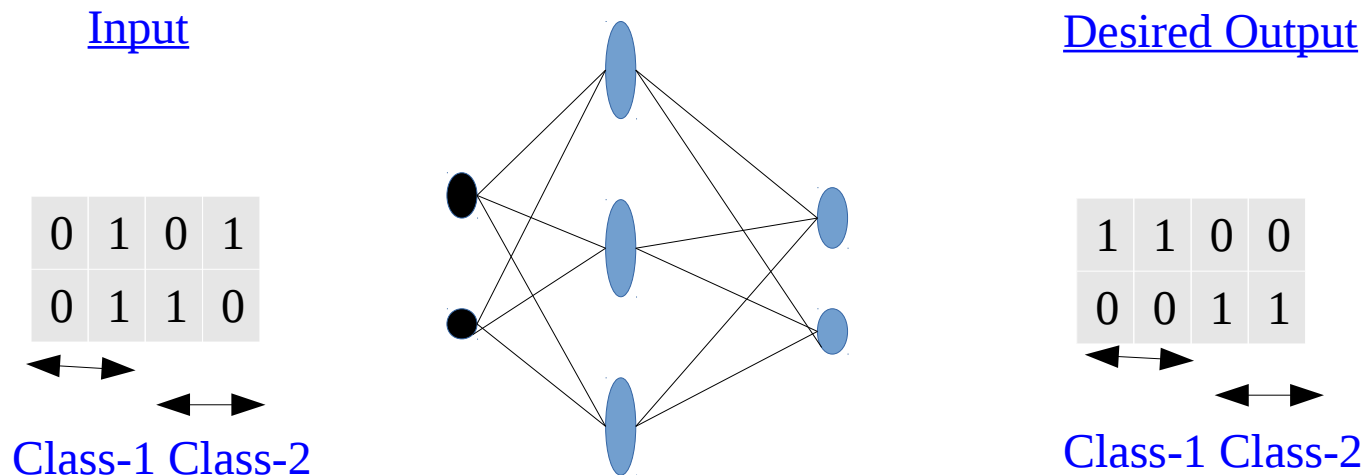
$$\begin{aligned} 0*w1+0*w2 &< Th \\ 0*w1+1*w2 &\geq Th \\ 1*w1+0*w2 &\geq Th \\ 1*w1+1*w2 &\geq Th \end{aligned}$$

$$\begin{aligned} Th &> 0 \\ W2 &\geq Th \\ W1 &\geq Th \\ w1+w2 &< Th \end{aligned}$$



Need of Hidden Layer

- Project the lower dimensional data into higher one if they are not separable.
- These projections are via weight vectors.



$$W = \begin{bmatrix} w_{11}, w_{12}; \\ w_{21}, w_{23}; \\ w_{31}, w_{32} \end{bmatrix}$$

$$W[3*2] * I[2*4] = I_h [3*4]$$

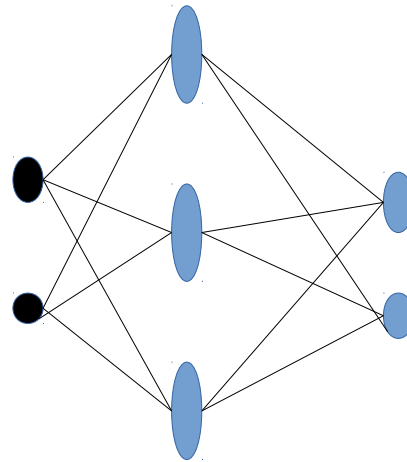
Trained weights

Input

0	1	0	1
0	1	1	0



Class-1 Class-2



Desired Output

1	1	0	0
0	0	1	1



Class-1 Class-2

$$W = \begin{bmatrix} W_{11}, W_{12}; \\ W_{21}, W_{23}; \\ W_{31}, W_{32} \end{bmatrix}$$

$$W[3 \times 2] * I[2 \times 4] = I_h [3 \times 4]$$

$$W_{ih} = [-6.793, 12.214; \\ 8.051, 8.053; \\ 12.208, -6.791];$$

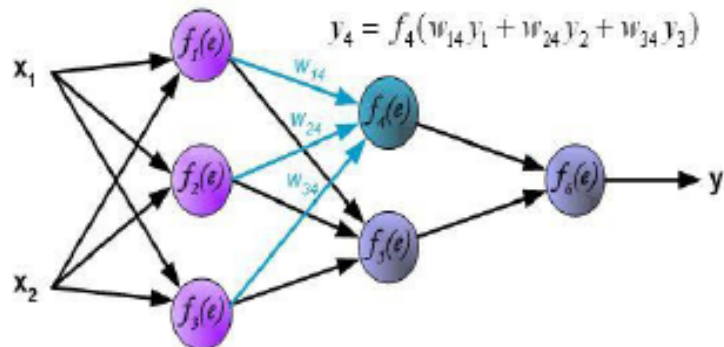
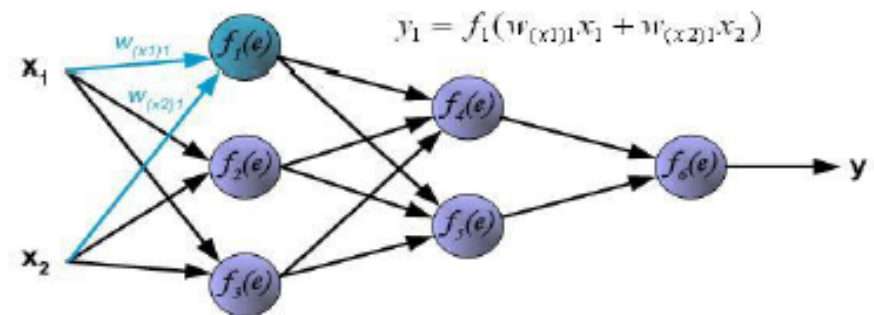
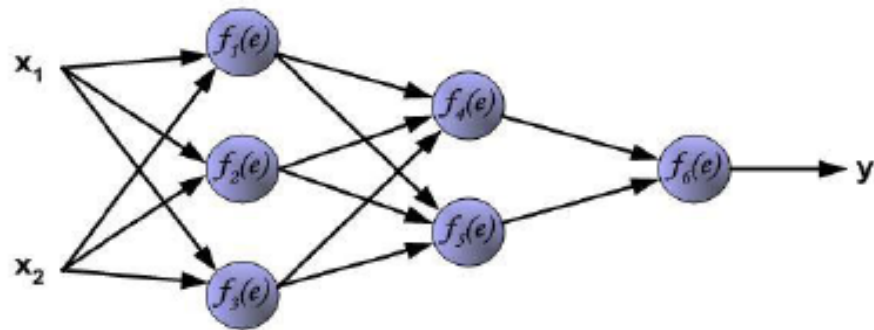
$$W_{ho} = [-11.58, 8.4, 8.4; \\ 11.57, -8.39, -8.39];$$

Actual output
(after training)

0.93	0.99	0.04	0.04
0.06	0.01	0.96	0.96

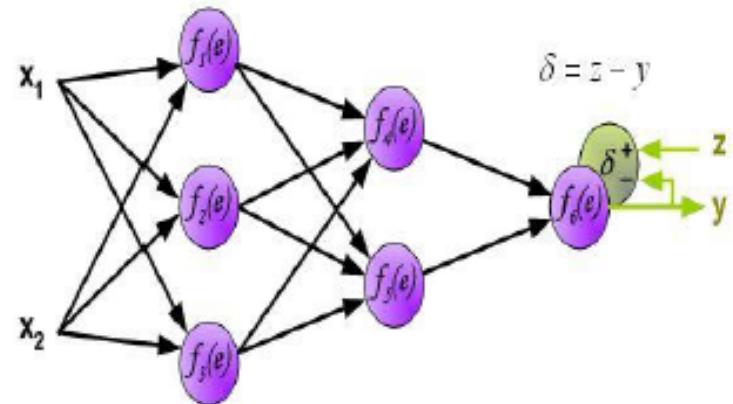
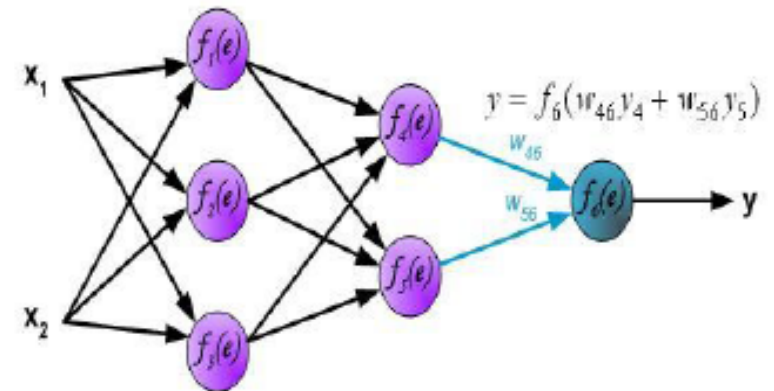
$$MSE = \text{sum}(\text{Desired-Actual}).^2 \\ = 9.9 \times 10^{-4}$$

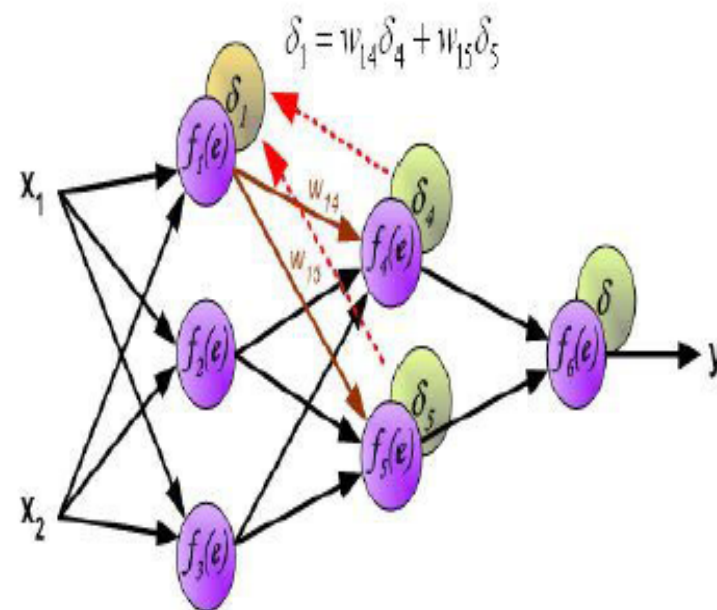
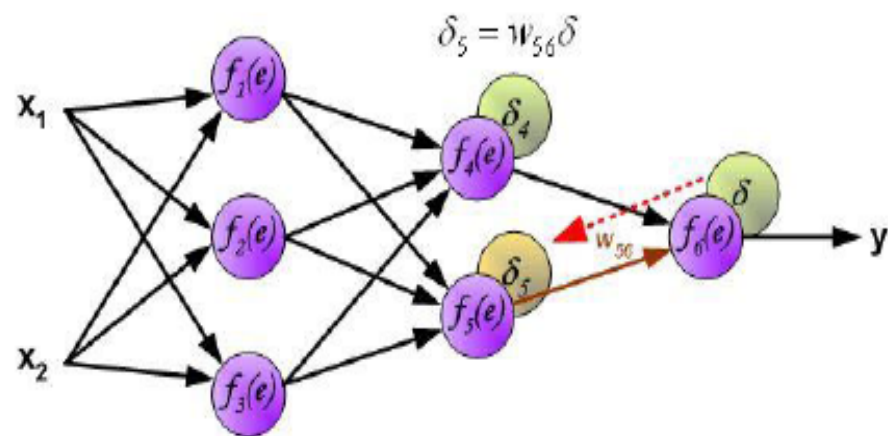
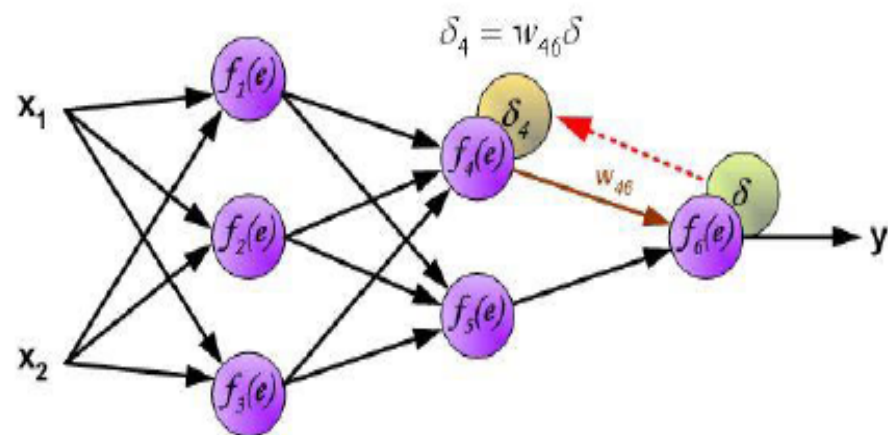
Illustration of Back-Propagation learning algorithm



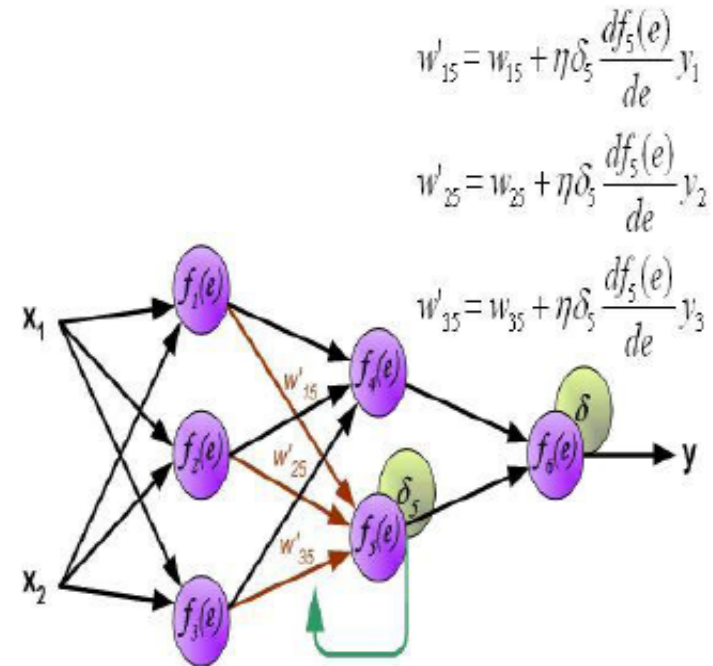
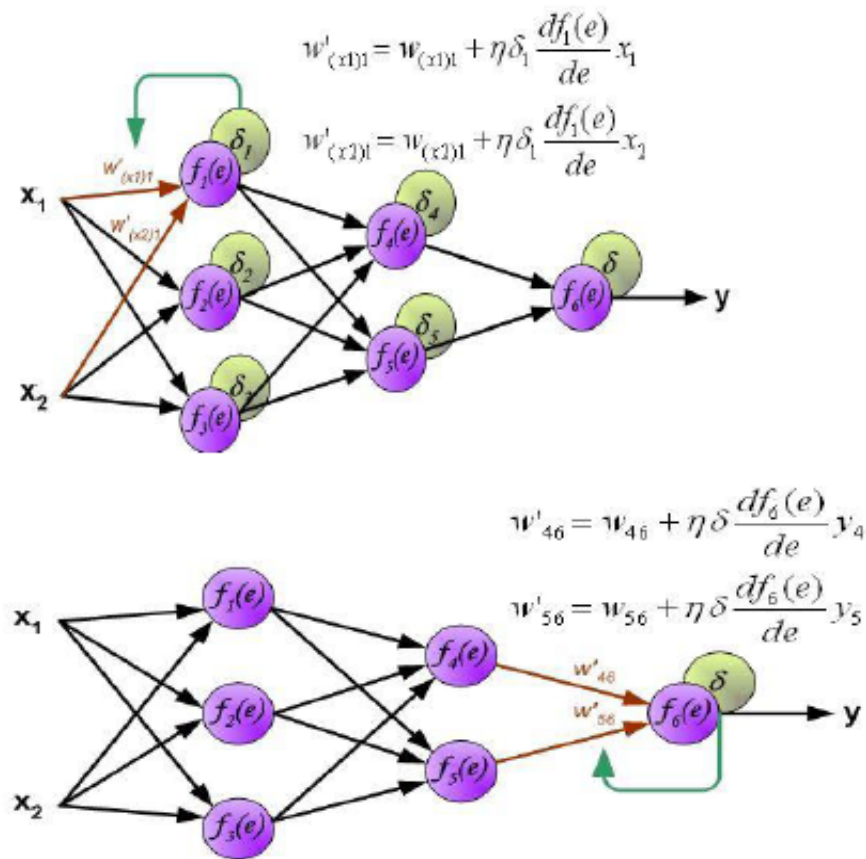
Sigmoid activation fn.

$$f(x) = 1 / (1 + \exp(-x))$$



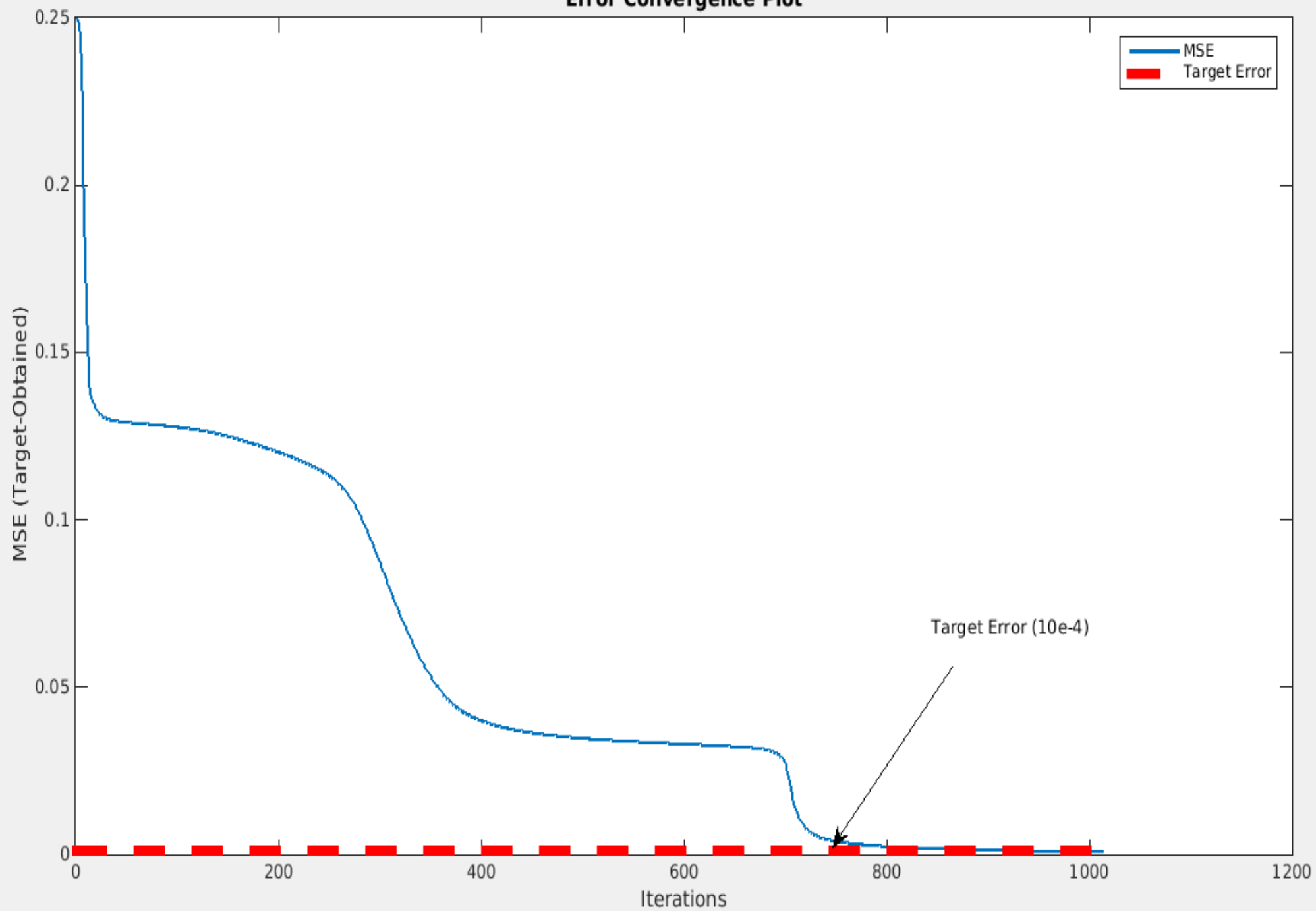


▪ Weight Update: Delta learning rule

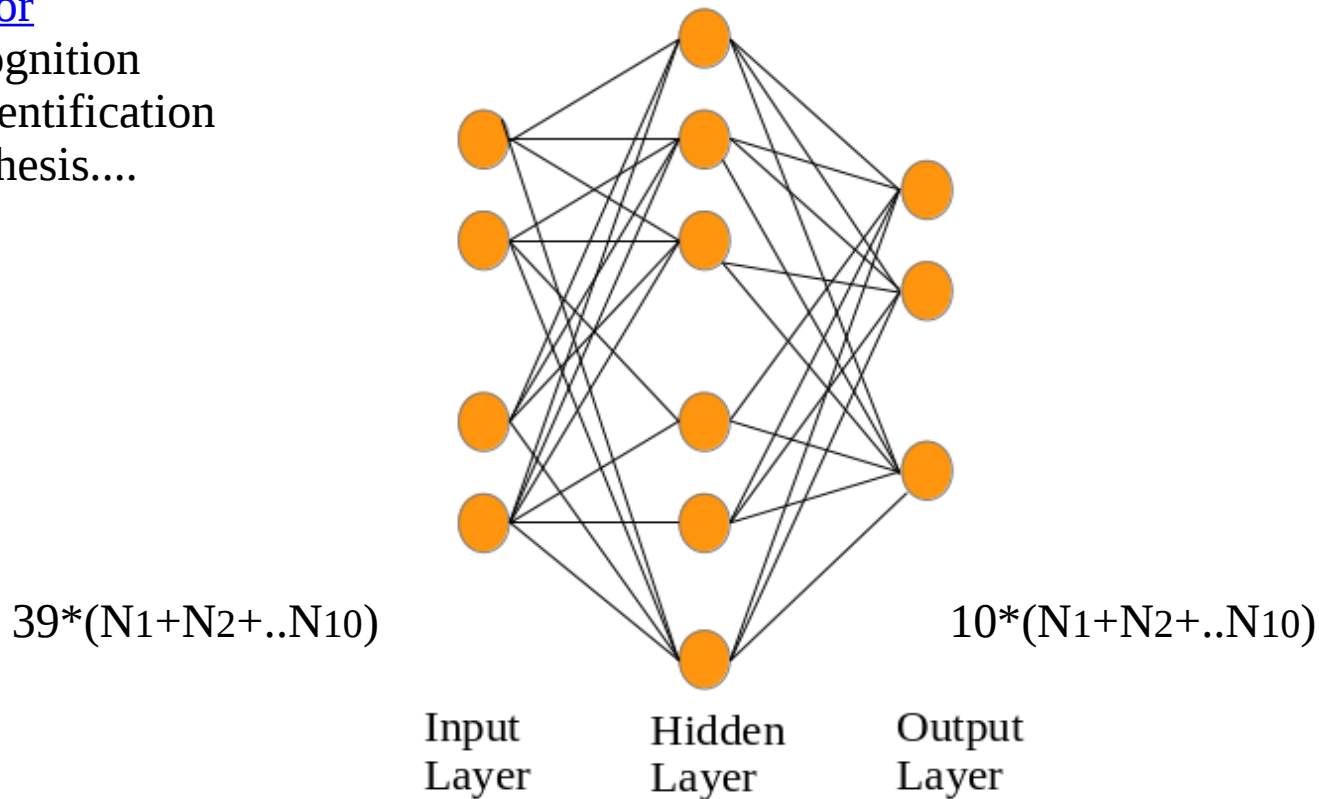


▪ Stopping Criteria: Mean square error \leq Goal error

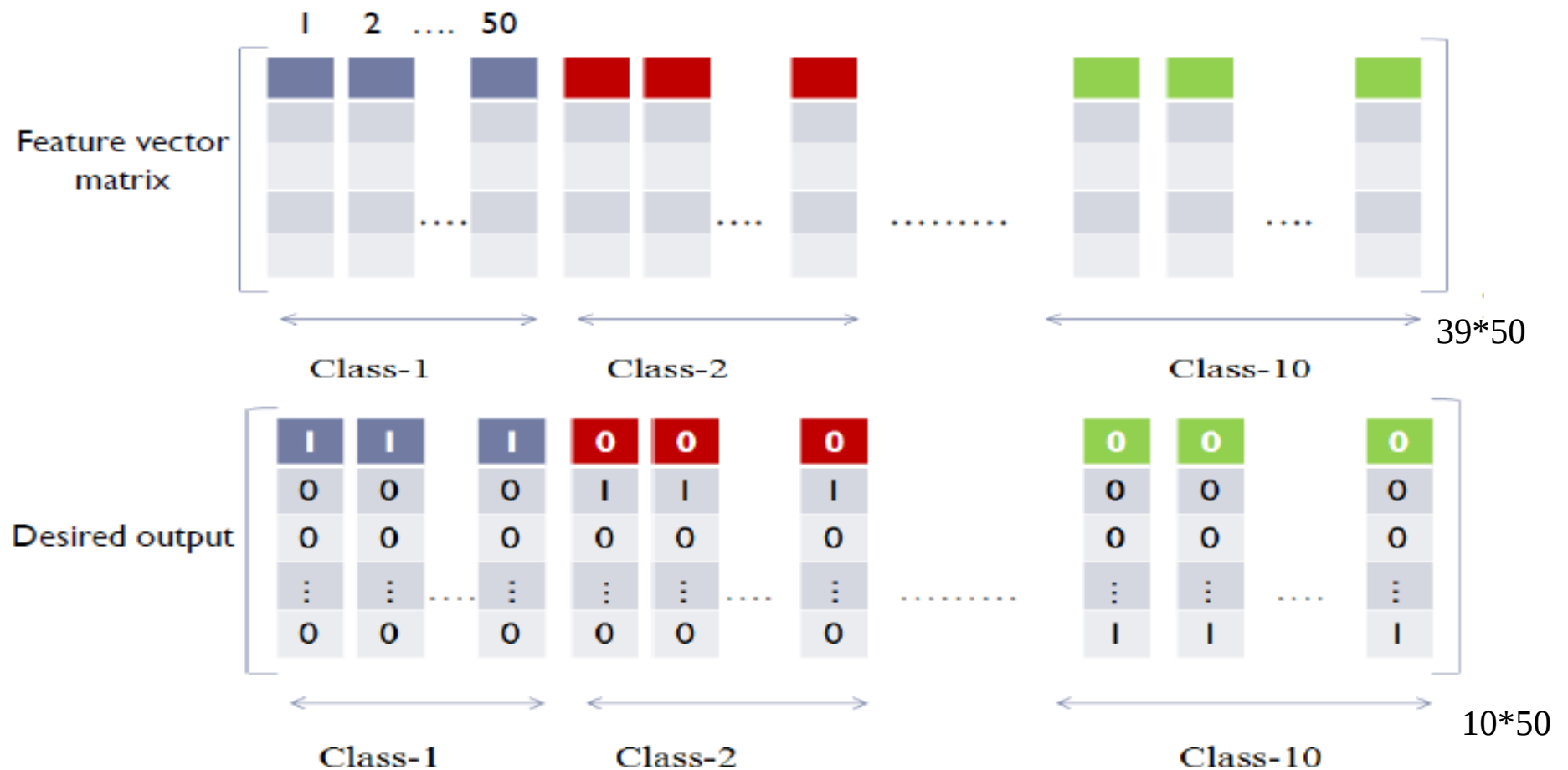
Error Convergence Plot



Applicable for
Speech Recognition
Language Identification
Speech Synthesis....



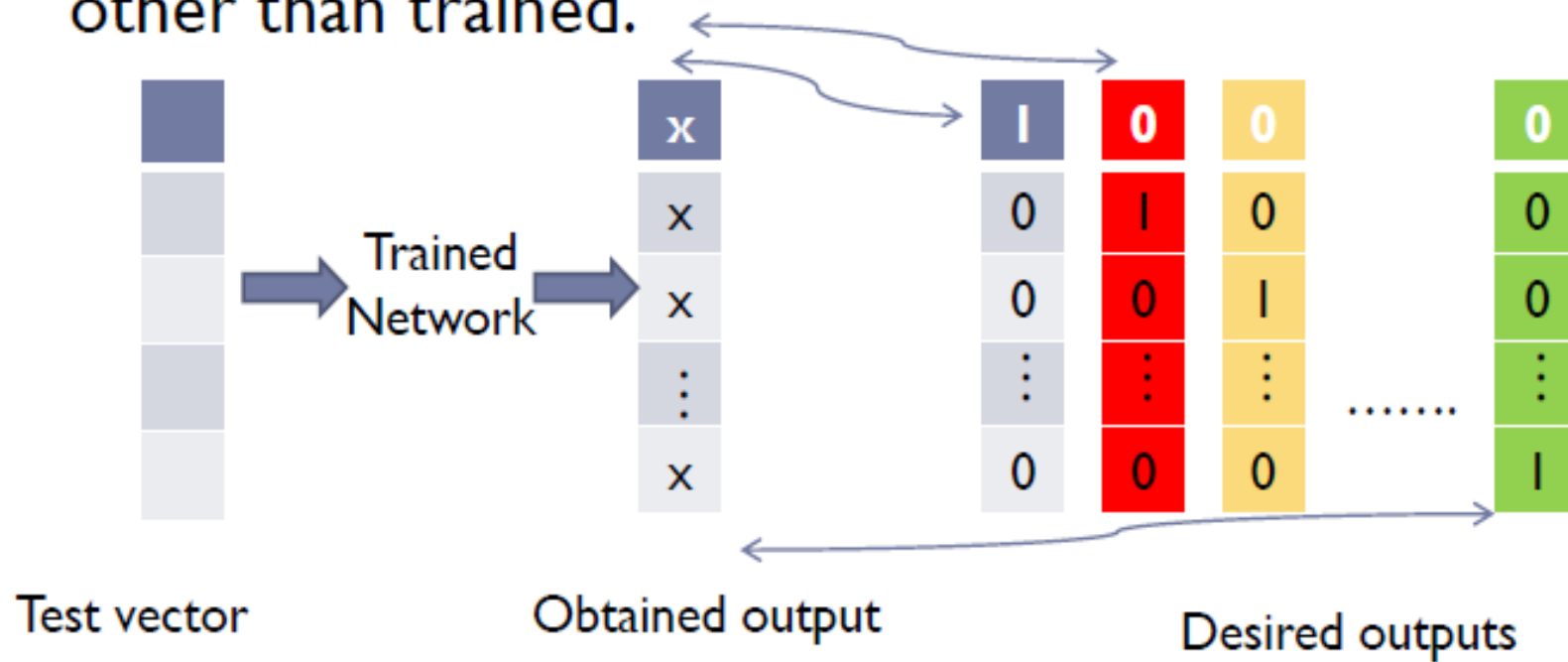
Training Feed forward Neural networks



- Learning rate=0.001
- Min. MSE=0.0001

Testing FFNN

- ▶ Weights are updated implying that the network has learnt.
- ▶ Check the algorithm efficiency by feeding a new input other than trained.

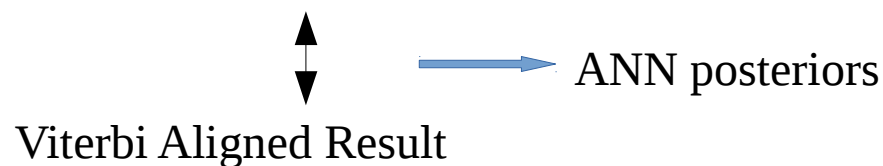


Comparing GMMs and ANNs

GMMs	ANNs
$\hat{w} = \underset{w}{\operatorname{argmax}} p(w \mathbf{x}) = \underset{w}{\operatorname{argmax}} p(\mathbf{x} w) p(w) / p(\mathbf{x})$	$y_j = \operatorname{logistic}(x_j) = \frac{1}{1 + e^{-x_j}}, \quad x_j = b_j + \sum_i y_i w_{ij},$
<ul style="list-style-type: none"> ▪ Maximize the Likelihood score ▪ Needs appropriate choice of Mixture components 	<ul style="list-style-type: none"> ▪ Maximize the posterior probabilities ▪ Choice of nodes and layers plays a vital role

Future Work

- ♦ To create HMM-ANN baseline
- ♦ HMM State Transitions



[HMM-ANN Baseline](#)

Deep Neural Networks(DNNs)

[Hinton et. al, *Neural Computation*, 2006]

- DNNs or Deep Belief Networks (DBNs) are Multi layer perceptrons with many hidden layers but the difference exist in the training phase.

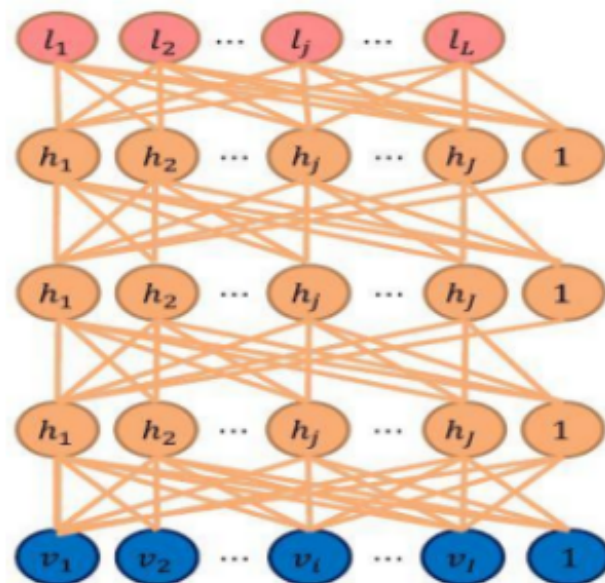


Fig. DBN/DNN Architecture

RBM-4

RBM-3

Restricted Boltzmann Machines

RBM-2

RBM-1

- G.Hinton discovered learning strategy for Multilayer Perceptrons
 - a) Pre-train each layer from bottom to top
 - b) Each pair of layers is an Restricted Boltzmann Machine(RBM)
 - c) Jointly fine-tune all layers using back-propagation

Restricted Boltzmann Machines

[Hinton et. al, *Neural Computation*, 2006]

- DNNs are formed by stacking several RBMs one upon another.

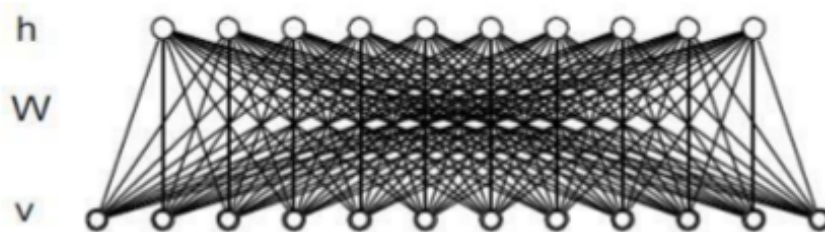


Fig. Restricted Boltzmann Machine (RBM)

Bernoulli-Bernoulli RBM

$$p(h_j = 1|\mathbf{v}; \theta) = \sigma \left(\sum_{i=1}^I w_{ij} v_i + a_j \right),$$
$$p(v_i = 1|\mathbf{h}; \theta) = \sigma \left(\sum_{j=1}^J w_{ij} h_j + b_i \right),$$

Gaussian-Bernoulli RBM

$$p(h_j = 1|\mathbf{v}; \theta) = \sigma \left(\sum_{i=1}^I w_{ij} v_i + a_j \right),$$
$$p(v_i|\mathbf{h}; \theta) = \mathcal{N} \left(\sum_{j=1}^J w_{ij} h_j + b_i, 1 \right),$$

$$\sigma(x) = 1/(1 + \exp(-x))$$

Unsupervised Pre-training

- ❖ $v=x_I$ is the input distribution.
- ❖ Find activations at the hidden layer

$$p(h_j = 1 | v; \theta) = \sigma\left(\sum_{i=1}^V w_{ij} v_i + a_j\right)$$

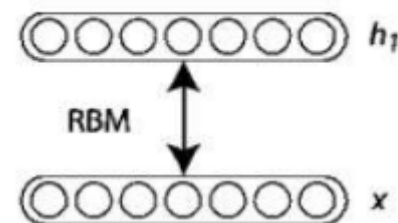
- ❖ Convert activation probabilities to binary representation, h_2
- ❖ Compute the activations at the visible unit given the binary data at the hidden unit.

$$p(v_i = 1 | h; \theta) = N\left(\sum_{j=1}^H w_{ij} h_j + b_i\right)$$

- ❖ Again the activation probabilities at the hidden units given the new input distribution are found, $P(h_2=I|x_2)$.
- ❖ The weight update rule for 1st iteration as described by the above procedure is given as

$$W = W + (h_I x_I^T - P(h_2=I|x_2) x_2^T)$$

Contrastive
Divergence



Stacking RBMs

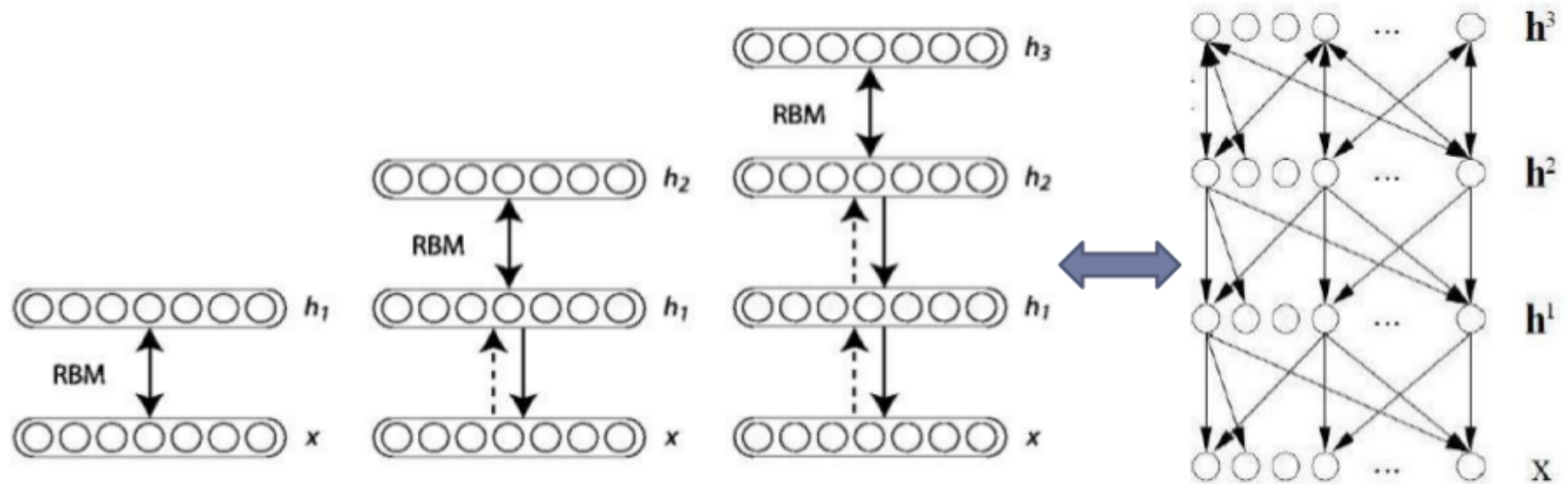
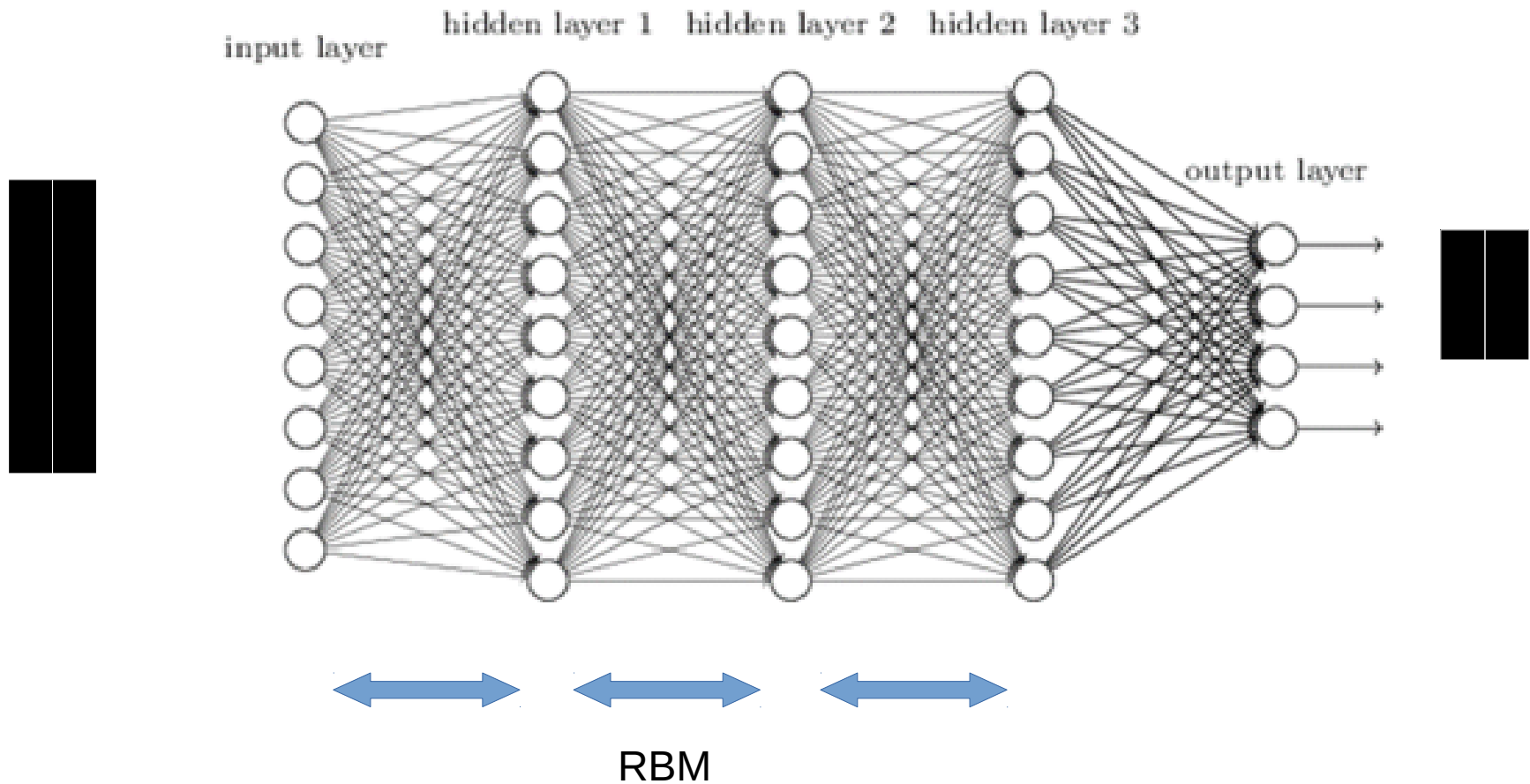


Fig. Stacking RBMs to form a DBN

- The activation probabilities of hidden units are used as the visible data for the next layer.
- This process is repeated to create other layers.
- Now Back-propagation algorithm is applied to adjust the DBN weights (fine-tuning phase)

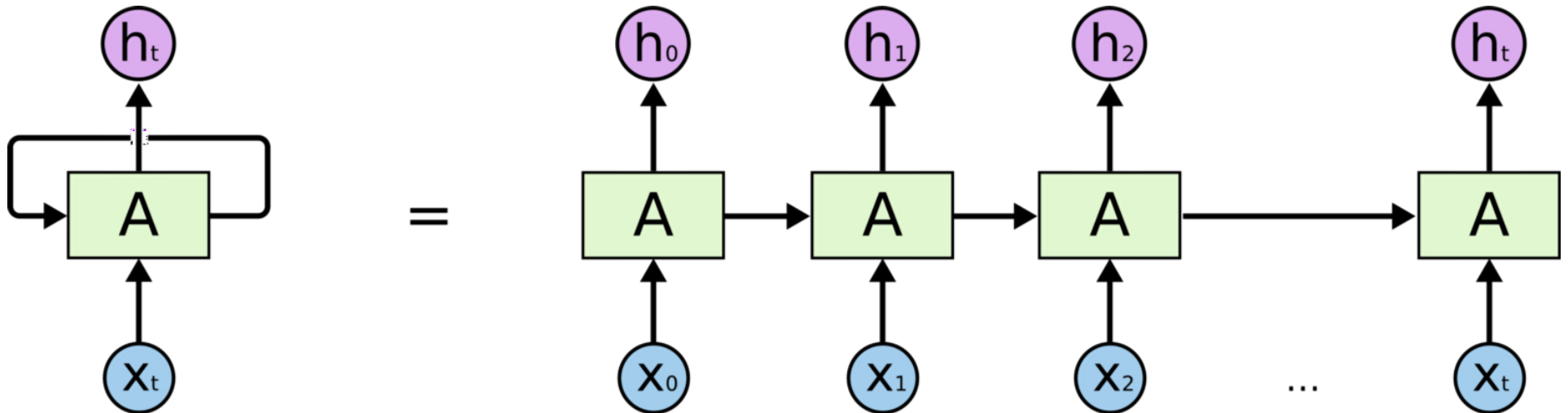
DNN



Training:

- Unsupervised Pre-training
- Supervised Fine tuning

RNN

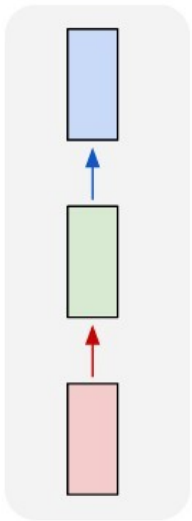


Applications:

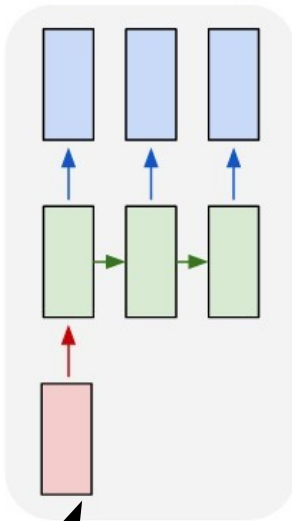
- Time series prediction
- Language modeling
- Text sentiment analysis
- Speech recognition
- Translation

Variations of RNN

one to one

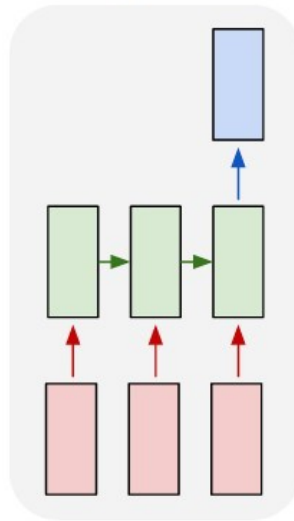


one to many



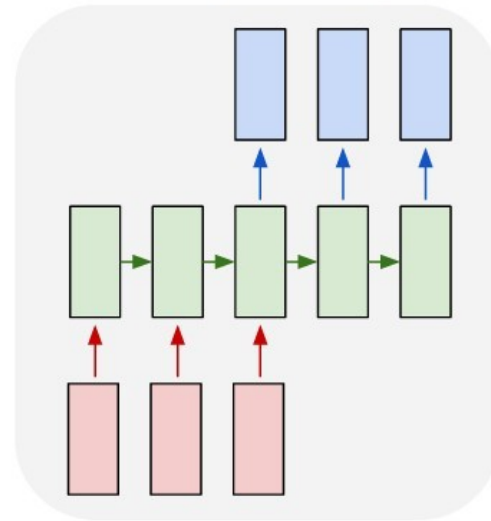
Ex: Image-to-Words

many to one



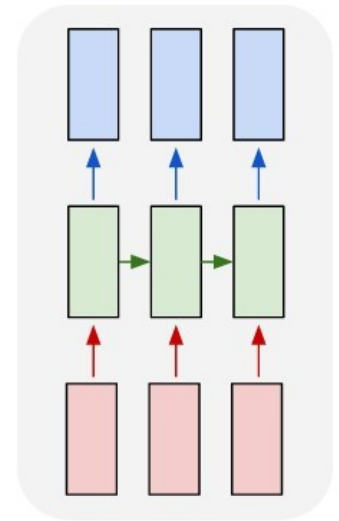
Ex: Sentence-to-Sentiment

many to many

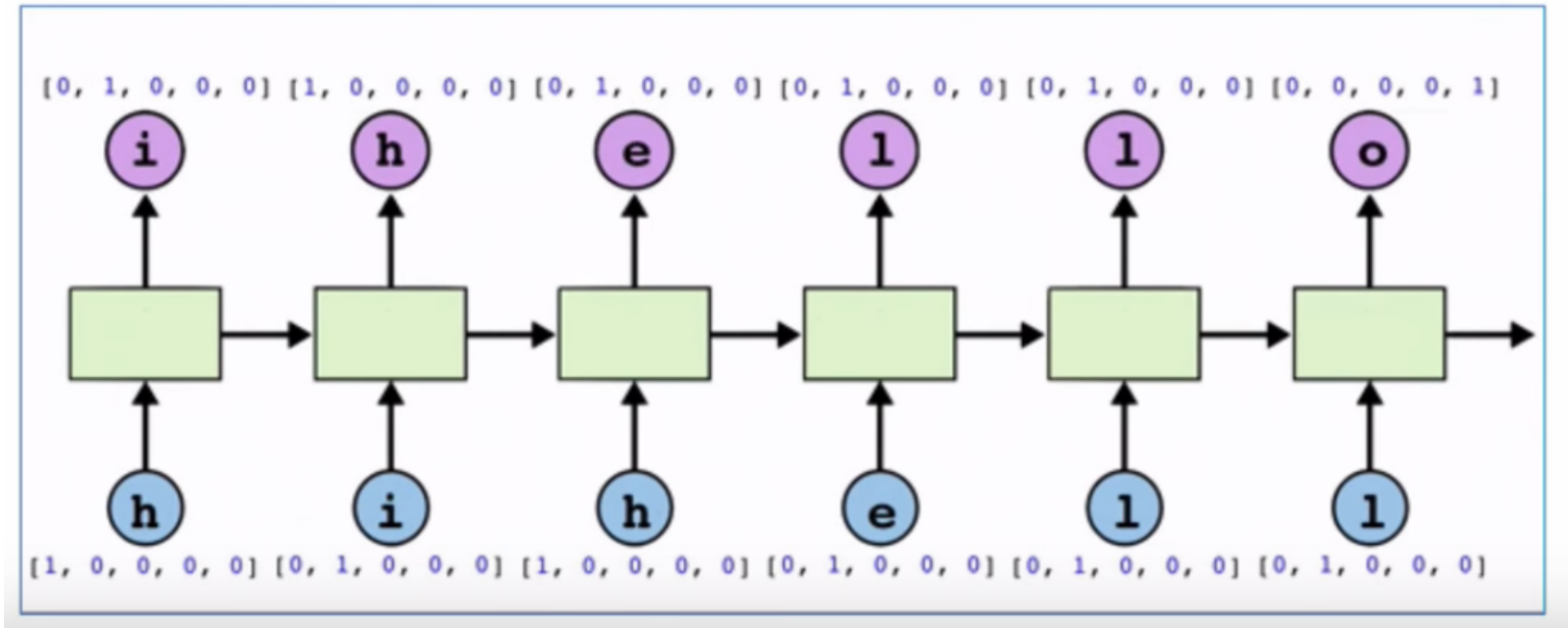


Ex: Translation

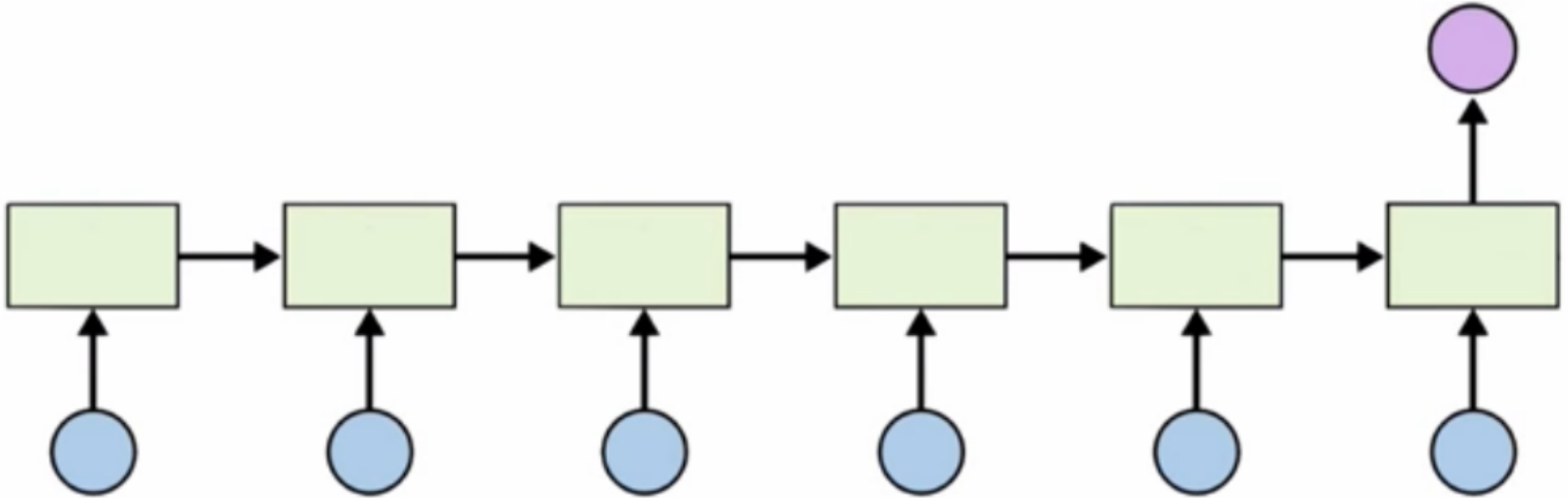
many to many



RNNs- Language modeling



RNNs --Classification

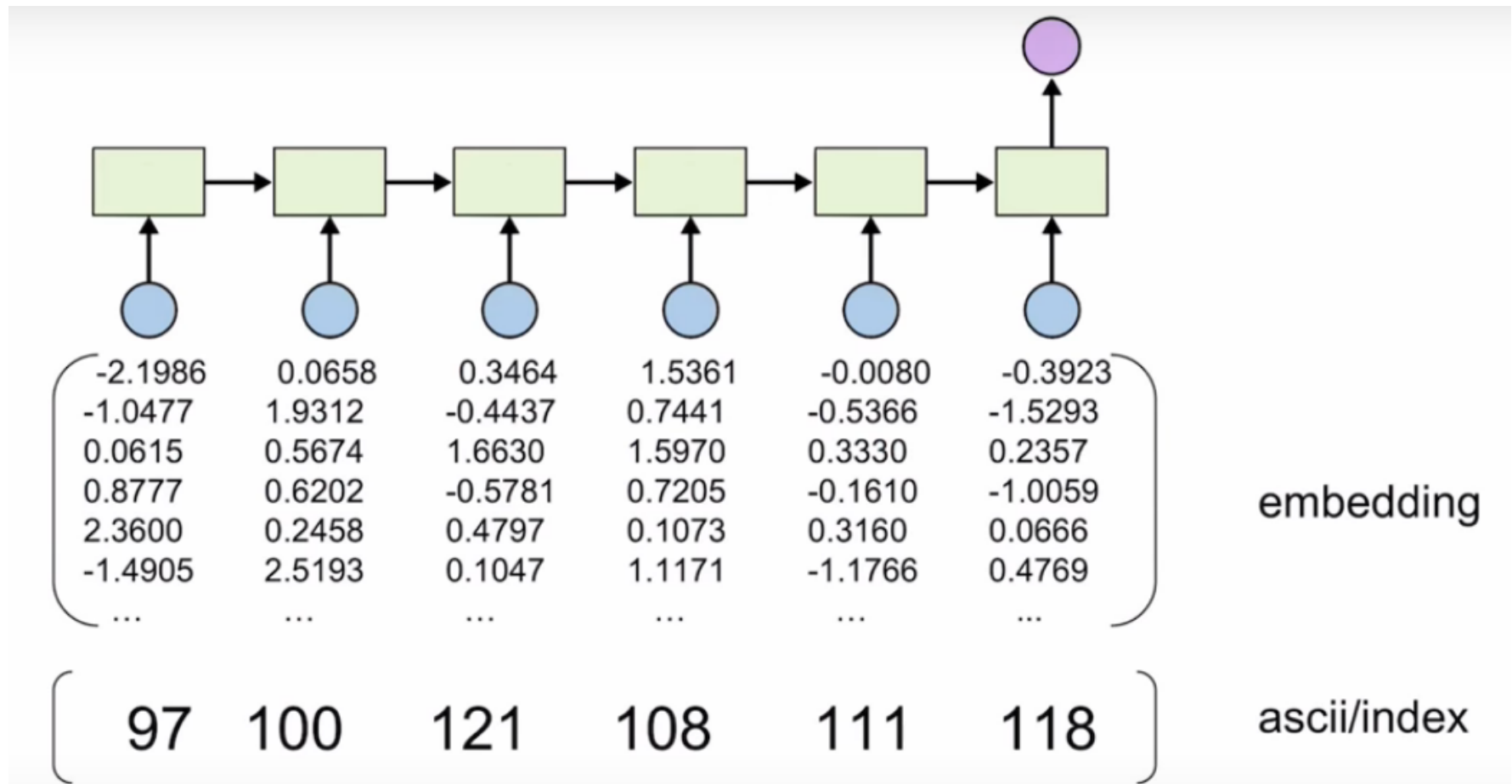


Names Language

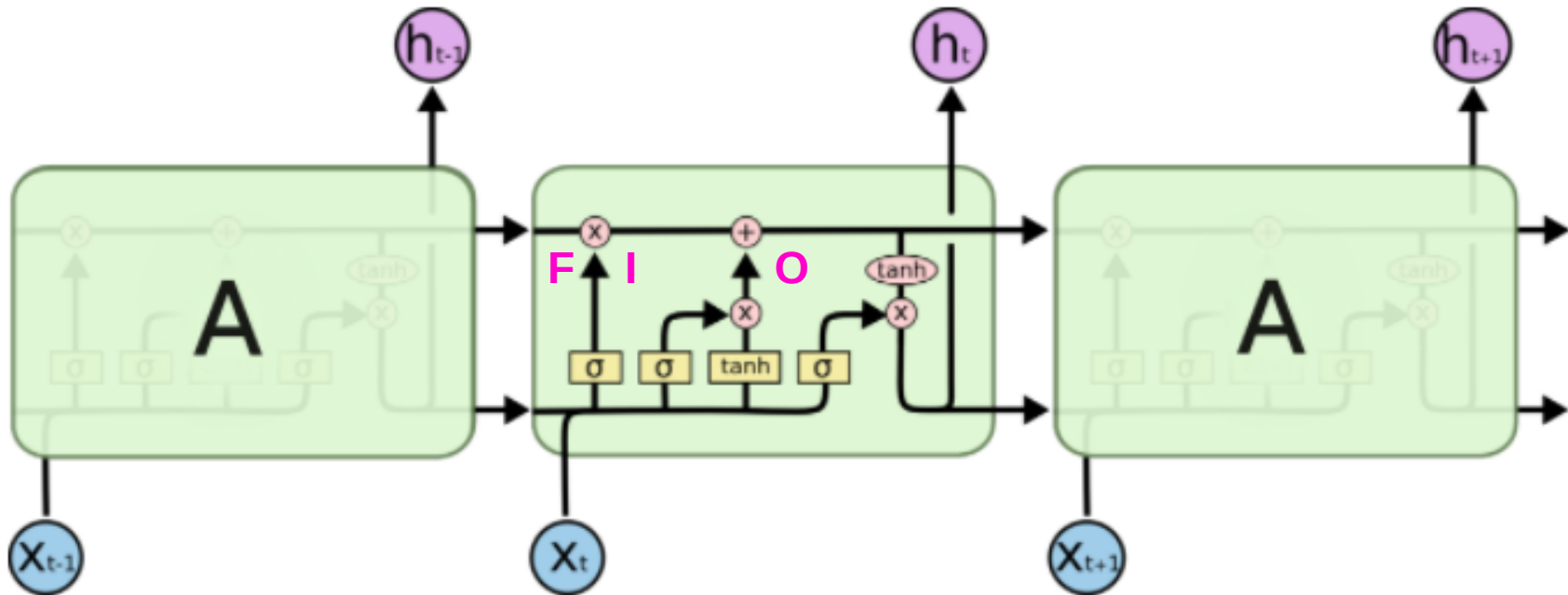
- Adalad German
- Adele German
- Lin Dan Chinese
- Chang Lee Chinese

....

RNNs--Classification



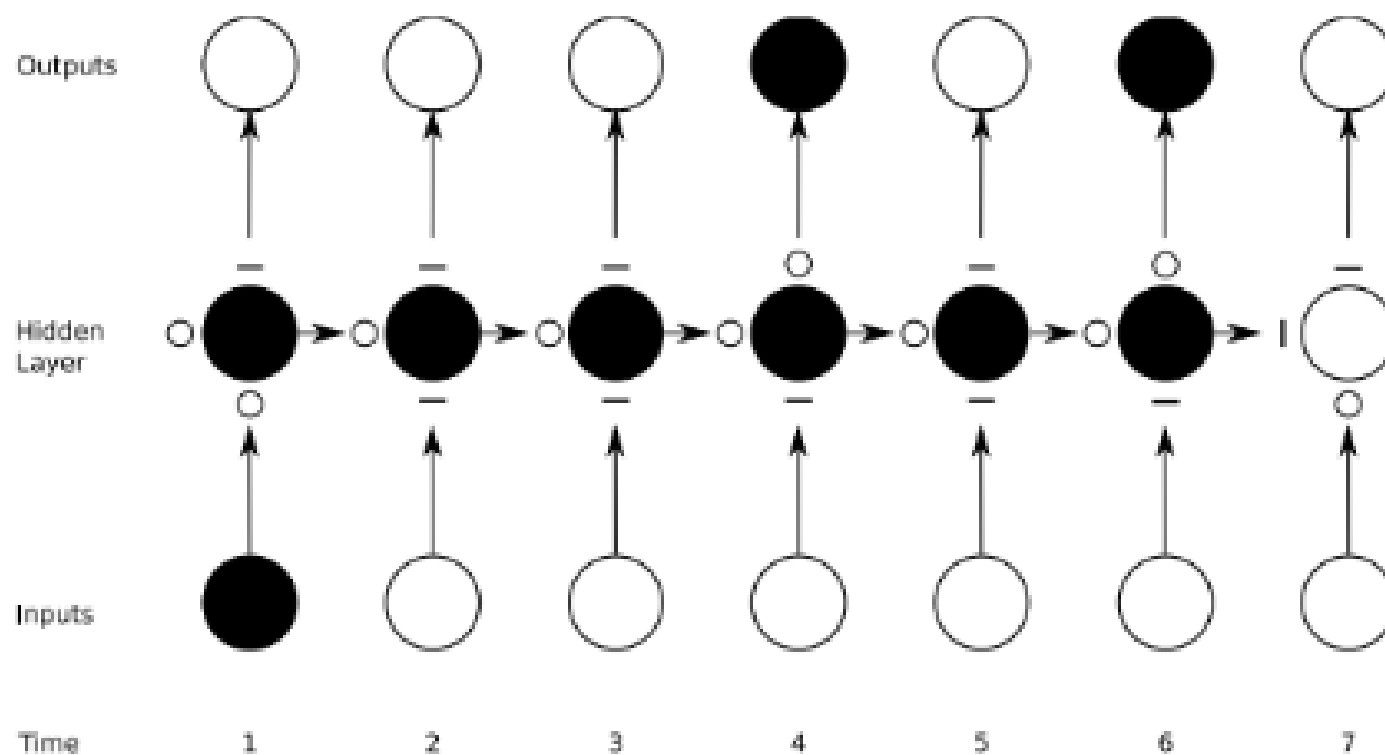
LSTM



Gates:

- Forget gate: Determines whether current contents of memory will be forgotten (erased)
- Input gate: Determines whether the input will be stored in the memory cell
- Output Gate: Determines if current memory contents will be output

LSTM full network



Graves et al 2013

Thank You