

## 4F10: Latent Variable and Sequence Models

Mark Gales

Michaelmas 2021

- Previous lecture examined relationships in data samples
  - conditional-independence important
- It is possible to introduce **latent variables** into the model
  - do not have to have any “meaning”
  - these variables are never observed in test (possibly in training)
  - marginalised over to get probabilities
  - **discrete** (mixture models, HMMs), **continuous** (factor-analysis)
- Revise/extend material from 3F8 & previous 4F10 lecture
  - simple language models
  - latent variable models
  - Expectation-Maximisation (EM)
  - hidden Markov models
  - conditional random fields

# Language Modelling

- Consider the word sequence

The cat sat on the ???

- language modelling can be viewed as answering

What is the missing word ???

- Writing this in terms of probabilities

$$P(w_i | \hat{w}_1, \dots, \hat{w}_{i-1}) = P(w_i | \hat{\mathbf{w}}_{1:i-1})$$

- $\hat{\mathbf{w}}_{1:i-1} = \hat{w}_1, \dots, \hat{w}_{i-1}$  observed words at positions 1 to  $i - 1$
- $w_i$  unknown word to be estimated (guessed) at position  $i$

- Predicted word can then be used as input for data generation
  - simple **auto-regressive** data generation process

$$\hat{w}_i \sim P(w_i | \hat{w}_1, \dots, \hat{w}_{i-1}); \quad P(w_{i+1} | \hat{w}_1, \dots, \hat{w}_i)$$

- Very simple process for generating data
  - initialise process with  $w_0 = \langle s \rangle$ , sentence start symbol

$\langle s \rangle$ ???	$P(w_1   \langle s \rangle)$
$\langle s \rangle$ The ???	$P(w_2   \langle s \rangle, \text{The})$
$\langle s \rangle$ The cat ???	$P(w_3   \langle s \rangle, \text{The}, \text{cat})$
$\langle s \rangle$ The cat sat ???	$P(w_4   \langle s \rangle, \text{The}, \text{cat}, \text{sat})$
$\langle s \rangle$ The cat sat on ???	$P(w_5   \langle s \rangle, \text{The}, \text{cat}, \text{sat}, \text{on})$

# Language Modelling - Simple?

- How large would a language model be - consider
  - **vocabulary size**  $|V|$ : (number of possible words)  $65536 = 2^{16}$
  - **sentence length**  $L$ : (length of sentence to generate)  $8 = 2^3$
- How large does the model need to be to predict the last word:
  - number of possible histories (any sentence):  $(|V|)^{L-1} = 2^{112}$
  - number of possible words (any sentence):  $|V| = 2^{16}$
  - total number of parameters:  $(|V|)^L = 2^{128} = 3.4 \times 10^{38}$
- Directly modelling all sentences not possible!
  - size of model dominated by number of histories
  - how to get a compact history representation?

# N-Gram Language Models (Markov-Chain 3M1)

- Simplest approach to make history compact - **truncate**

$$P(w_i | \mathbf{w}_{1:i-1}) \approx P(w_i | \mathbf{w}_{i-N+1:i-1})$$

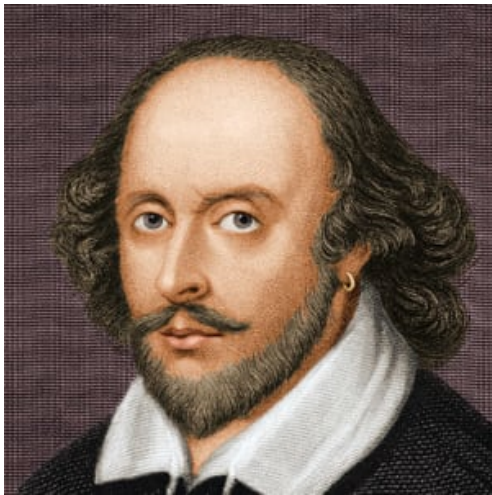
notation:  $\mathbf{w}_{1:i-1}$  (rather than  $\hat{\mathbf{w}}_{1:i-1}$ ) as consider all “histories”

- $N$  determines the length of the history to consider
  - $N = 1$  unigram,  $N = 2$  bigram,  $N = 3$  trigram ...
- 
- A **discrete** model estimated from data

$$P(w_i | \mathbf{w}_{i-N+1:i-1}) \approx \frac{\mathcal{C}(w_{i-N+1}, \dots, w_i)}{\mathcal{C}(w_{i-N+1}, \dots, w_{i-1})}$$

- $\mathcal{C}()$  counts number of occurrences in the training data
- “separate” probability computed for every word
- general probability mass function (PMF) smoothing approaches

# William Shakespeare (1564-1616)





- **Unigram** - No prior information
  - Every enter now severally so, let
  - Will rash been and by I the me loves gentle me not slavish page, the and hour; ill let
- **Bigram** - Current word
  - What means, sir. I confess she? then all sorts, he is trim, captain.
  - The world shall- my lord!
- **Trigram** - Current and previous word
  - Indeed the duke; and had a very good friend.
  - Sweet prince, Falstaff shall die. Harry of Monmouth's grave.
- **4-gram** - Current and previous two words
  - It cannot be but so.
  - Enter Leonato's brother Antonio, and the rest, but seek the weary beds of people sick.

### **Opinion** Artificial intelligence (AI)

A robot wrote this entire article. Are you scared yet, human?

*GPT-3*

- Deep-learning has revolutionised language modelling
  - GPT-3 is a **large** language model from OpenAI
  - about 175 billion model parameters trained on 500 billion tokens
- Architecture of model will be discussed later in course

# Latent Variable Models

# "Static" Latent Variable Generative Models

- N-gram language model use word sequence as context
  - causes exponential growth in size of model with history length
- **Alternative:** introduce a variable that represents the history
  - this variable may be discrete or continuous
  - just needs to represent the history in an appropriate form
- This section introduces **latent variables**
  - initially for "static" data - for example mixture models
  - the next section will discuss application to sequence models - for example HMMs

# Gaussian Mixture Models

- Gaussian mixture models (GMMs) are based on Gaussians
  - form of the Gaussian distribution:

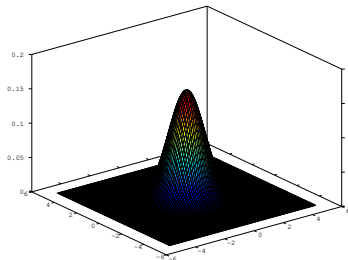
$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)$$

- Each **component** modelled using a Gaussian distribution

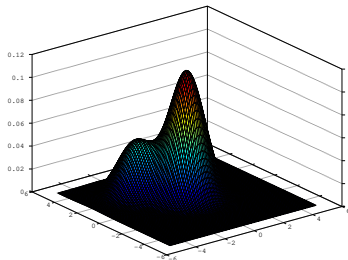
$$p(\mathbf{x}) = \sum_{m=1}^M P(c_m) p(\mathbf{x}|c_m) = \sum_{m=1}^M P(c_m) \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)$$

- **component prior**:  $P(c_m)$
  - **component distribution**:  $p(\mathbf{x}|c_m) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)$
- Flexible model, able to model range of distributions

# Gaussian Mixture Model Example



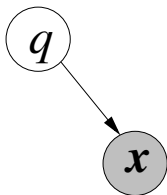
$$\mathcal{N}(\mu_1, \Sigma)$$



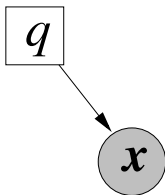
$$0.7\mathcal{N}(\mu_1, \Sigma) + 0.3\mathcal{N}(\mu_2, \Sigma)$$

$$\mu_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad \mu_2 = \begin{bmatrix} -2 \\ 2 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

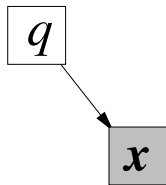
# Latent Variable Generative Models



Factor Analysis



Gaussian Mixture Model



Discrete Mixture Model

- Three Bayesian Networks (BNs) for an observation  $x$ 
  - indicator variable  $q$ : continuous  $z$  or discrete  $c_m$  space
    - factor analysis  $\int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$
    - Gaussian mixture models  $\sum_{m=1}^M P(c_m)p(\mathbf{x}|c_m)$
    - discrete mixture model  $\sum_{m=1}^M P(c_m)P(\mathbf{x}|c_m)$
- Extensively used in many machine learning applications

# Factor Analysis

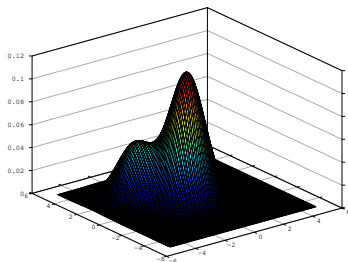
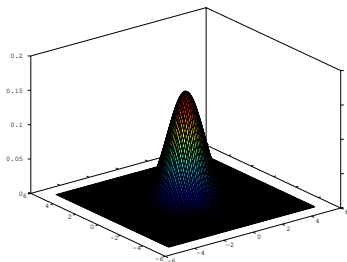
- Two views of Factor Analysis (FA)
  - low-dimensional manifold representation
  - compact covariance matrix for multi-variate Gaussians
- Form of model dimensionality of  $\mathbf{z}$  less than  $\mathbf{x}$ ,  $p < d$ 
  - $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$ : low-dimensional subspace representation
  - $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; \mathbf{C}\mathbf{z}, \mathbf{\Sigma}_{\text{diag}})$ :  
 $\mathbf{C}$  loading matrix,  $\mathbf{\Sigma}_{\text{diag}}$  diagonal covariance matrix
  - if  $\mathbf{\Sigma}_{\text{diag}} = \sigma^2 \mathbf{I}$  model is probabilistic PCA
- As all elements Gaussian closed-form solution

$$p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} = \mathcal{N}(\mathbf{x}; \mathbf{0}, \mathbf{C}\mathbf{C}^T + \mathbf{\Sigma}_{\text{diag}})$$



# Expectation Maximisation

# Gaussian Mixture Model Example



- Would like to estimate latent variable model parameters,  $\theta$ 
  - training data  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
  - generative model: maximum (log-)likelihood

$$\hat{\theta} = \arg \max_{\theta} \{\mathcal{L}(\theta)\} = \arg \max_{\theta} \left\{ \sum_{i=1}^N \log(p(\mathbf{x}_i; \theta)) \right\}$$

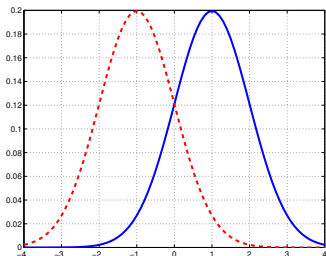
- Gradient descent could be used
  - need to determine learning rate ...

Is there any alternative?

- Use Gaussian Mixture Models as the example
  - discrete latent variable (for continuous see examples paper)

# “Simple” Expectation Maximisation

- If you knew the component for each observation - simple
  - standard (multi-variate) Gaussian distribution
  - but you don't - make an initial guess of the parameters  $\theta^{(0)}$

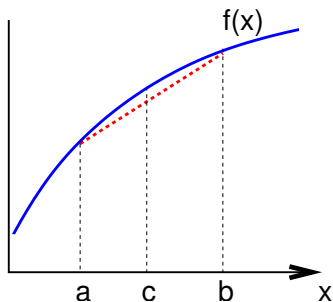


- set  $k = 0$ , using parameters  $\theta^{(k)}$ 
  - “assign” each observation to a component
  - use Bayes' decision rule for assignment
- yields component id for each observation
- simple to update model  $\rightarrow \theta^{(k+1)}$
- repeat process,  $k = k + 1$

Each iteration guaranteed not to decrease likelihood

- “Hard” assignment not always reasonable (overlapping data)
  - derive a more general iterative approach

# Jensen's Inequality



- Jensen's inequality states:

$$f\left(\sum_{m=1}^M \lambda_m x_m\right) \geq \sum_{m=1}^M \lambda_m f(x_m)$$

$f()$  is any **concave function** and

$$\sum_{m=1}^M \lambda_m = 1, \quad \lambda_m \geq 0 \quad m = 1, \dots, M$$

- From the diagram above (left)

$$f(c) = f((1 - \lambda)a + \lambda b) \geq (1 - \lambda)f(a) + \lambda f(b)$$

- Interested in applying to GMMs -  $\lambda_m$  component prior  $P(c_m)$

$$\log\left(\sum_{m=1}^M \lambda_m p(\mathbf{x}|c_m)\right) \geq \sum_{m=1}^M \lambda_m \log(p(\mathbf{x}|c_m)), \quad \sum_{m=1}^M \lambda_m = 1, \quad \lambda_m \geq 0$$

- Changing parameters from  $\theta^{(k)}$  to  $\theta^{(k+1)}$  - desire

$$\mathcal{L}(\theta^{(k+1)}) - \mathcal{L}(\theta^{(k)}) = \sum_{i=1}^N \log \left( \frac{p(\mathbf{x}_i; \theta^{(k+1)})}{p(\mathbf{x}_i; \theta^{(k)})} \right) \geq 0$$

- This can be written as

$$\begin{aligned} \mathcal{L}(\theta^{(k+1)}) - \mathcal{L}(\theta^{(k)}) &= \sum_{i=1}^N \log \left( \frac{1}{p(\mathbf{x}_i; \theta^{(k)})} \sum_{m=1}^M (p(\mathbf{x}_i, c_m; \theta^{(k+1)})) \right) \\ &= \sum_{i=1}^N \log \left( \frac{1}{p(\mathbf{x}_i; \theta^{(k)})} \sum_{m=1}^M \left( \frac{P(c_m | \mathbf{x}_i; \theta^{(k)}) p(\mathbf{x}_i, c_m; \theta^{(k+1)})}{P(c_m | \mathbf{x}_i; \theta^{(k)})} \right) \right) \end{aligned}$$

- Jensen's inequality - use  $P(c_m | \mathbf{x}_i; \theta^{(k)})$  as  $\lambda_m$

$$\mathcal{L}(\theta^{(k+1)}) - \mathcal{L}(\theta^{(k)}) \geq \sum_{i=1}^N \sum_{m=1}^M P(c_m | \mathbf{x}_i; \theta^{(k)}) \log \left( \frac{p(\mathbf{x}_i, c_m; \theta^{(k+1)})}{p(\mathbf{x}_i; \theta^{(k)}) P(c_m | \mathbf{x}_i; \theta^{(k)})} \right)$$

- This inequality can be written as

$$\mathcal{L}(\boldsymbol{\theta}^{(k+1)}) - \mathcal{L}(\boldsymbol{\theta}^{(k)}) \geq \mathcal{Q}(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta}^{(k+1)}) - \mathcal{Q}(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta}^{(k)})$$

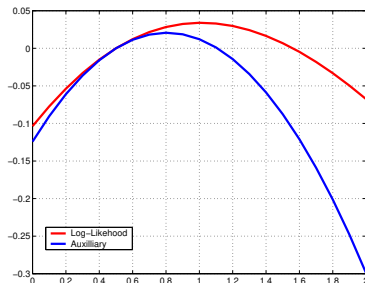
- where

$$\mathcal{Q}(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta}^{(k+1)}) = \sum_{i=1}^N \sum_{m=1}^M P(c_m | \mathbf{x}_i; \boldsymbol{\theta}^{(k)}) \log \left( p(\mathbf{x}_i, c_m; \boldsymbol{\theta}^{(k+1)}) \right)$$

- This is known as the **auxiliary function**
  - optimising the log-likelihood requires

$$\mathcal{Q}(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta}^{(k+1)}) - \mathcal{Q}(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta}^{(k)}) \geq 0$$

- increase is lower bound on the increase in the log likelihood.



- Data generated from the following GMM:

$$x \sim 0.4 \times \mathcal{N}(1, 1) + 0.6 \times \mathcal{N}(-1, 1)$$

- find first component mean: initial estimate is 0.5:

$$x^{(0)} \sim 0.4 \times \mathcal{N}(0.5, 1) + 0.6 \times \mathcal{N}(-1, 1)$$

- log-likelihood and auxiliary function differences to  $\mathcal{L}(\theta^{(0)})$  and  $Q(\theta^{(0)}, \theta^{(0)})$  plotted above



# General form and Continuous Auxiliary Functions

- For GMMs the latent variable of interest is  $c_m$ , generalise
  - consider a set of latent variable  $\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_N\}$
  - associated with the set of observations  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
- The general form of auxiliary function becomes

$$\mathcal{Q}(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta}^{(k+1)}) = \sum_{\forall \mathbf{Z}} P(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}^{(k)}) \log \left( p(\mathbf{X}, \mathbf{Z}; \boldsymbol{\theta}^{(k+1)}) \right)$$

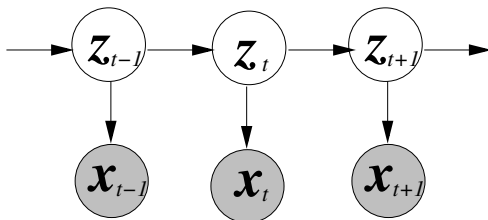
- require posteriors over the latent variables  $\mathbf{Z}$
- The continuous latent variable case version is

$$\mathcal{Q}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k+1)}) = \int p(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}^{(k)}) \log \left( p(\mathbf{X}, \mathbf{Z}; \boldsymbol{\theta}^{(k+1)}) \right) d\mathbf{Z}$$

- applicable to HMMs, FA models, Kalman Smoothers ...

# Hidden Markov Models

# Discrete Kalman Filters (3F8)



- Continuous latent variable (similar to FA)
  - linear relationship between variables

$$\mathbf{z}_t = \mathbf{A}\mathbf{z}_{t-1} + \boldsymbol{\nu}_t; \quad \mathbf{x}_t = \mathbf{C}\mathbf{z}_t + \boldsymbol{\epsilon}_t$$

- $\boldsymbol{\nu}_t$  are Gaussian distributed (and independent)  $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_\nu)$
- $\boldsymbol{\epsilon}_t$  are Gaussian distributed (and independent)  $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_\epsilon)$

# Discrete Kalman Filters (3F8)

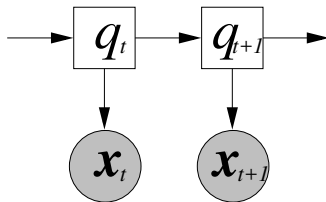
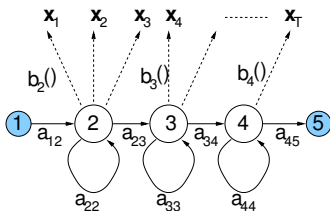
- This yields the following equations

$$p(\mathbf{z}_t | \mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{z}_t; \mathbf{A}\mathbf{z}_{t-1}, \mathbf{\Sigma}_\nu); \quad p(\mathbf{x}_t | \mathbf{z}_t) = \mathcal{N}(\mathbf{x}_t; \mathbf{C}\mathbf{z}_t, \mathbf{\Sigma}_\epsilon)$$

- Everything is Gaussian and linear closed-form solutions
  - don't remember these - for illustration purposes

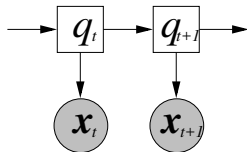
$$\begin{aligned} p(\mathbf{x}_t | \mathbf{x}_1, \dots, \mathbf{x}_{t-1}) &= \int p(\mathbf{x}_t | \mathbf{z}_t) p(\mathbf{z}_t | \mathbf{x}_1, \dots, \mathbf{x}_{t-1}) d\mathbf{z}_t \\ &= \int p(\mathbf{x}_t | \mathbf{z}_t) \int p(\mathbf{z}_t | \mathbf{z}_{t-1}) p(\mathbf{z}_{t-1} | \mathbf{x}_1, \dots, \mathbf{x}_{t-1}) d\mathbf{z}_{t-1} d\mathbf{z}_t \\ &= \int p(\mathbf{x}_t | \mathbf{z}_t) \int p(\mathbf{z}_t | \mathbf{z}_{t-1}) \frac{p(\mathbf{x}_{t-1} | \mathbf{z}_{t-1}) p(\mathbf{z}_{t-1} | \mathbf{x}_1, \dots, \mathbf{x}_{t-2})}{p(\mathbf{x}_{t-1} | \mathbf{x}_1, \dots, \mathbf{x}_{t-2})} d\mathbf{z}_{t-1} d\mathbf{z}_t \end{aligned}$$

# Hidden Markov Model



# Hidden Markov Models

- Important sequence data is the **hidden Markov model** (HMM)
  - an example of a **dynamic Bayesian network** (DBN)
  - consider a sequence of observations  $\mathbf{x}_1, \dots, \mathbf{x}_T$



- add discrete **latent variables**
  - $q_t$  describes discrete **state-space**
  - conditional independence assumptions

$$P(q_t | q_0, \dots, q_{t-1}) = P(q_t | q_{t-1})$$

$$p(\mathbf{x}_t | \mathbf{x}_1, \dots, \mathbf{x}_{t-1}, q_0, \dots, q_t) = p(\mathbf{x}_t | q_t)$$

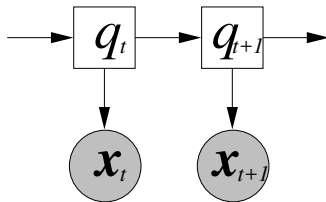
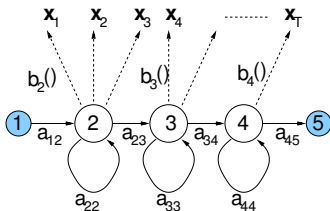
- The likelihood of the data is

$$p(\mathbf{x}_1, \dots, \mathbf{x}_T) = \sum_{\mathbf{q} \in \mathbf{Q}_T} P(\mathbf{q}) p(\mathbf{x}_1, \dots, \mathbf{x}_T | \mathbf{q}) = \sum_{\mathbf{q} \in \mathbf{Q}_T} P(q_0) \prod_{t=1}^T P(q_t | q_{t-1}) p(\mathbf{x}_t | q_t)$$

$\mathbf{Q}_T$  is all possible state sequences for  $T$  observations

- Two types of states often defined for HMMs (total  $N$  states)
  - **emitting** states: produce the observation sequence
  - **non-emitting** states: used to define valid state and end states
- The parameters are normally split into two
  - states are denoted  $s_i$ , assume  $s_1$  and  $s_N$  are non-emitting
  - **transition matrix  $\mathbf{A}$** :  
 $a_{ij} = P(q_t = s_j | q_{t-1} = s_i)$  is the probability of transitioning from state  $s_i$  to state  $s_j$
  - **state output probability  $\{b_2(\mathbf{x}_t), \dots, b_{N-1}(\mathbf{x}_t)\}$** :  
 $b_j(\mathbf{x}_t) = p(\mathbf{x}_t | q_t = s_j)$  is the output distribution for state  $s_j$
- Need to estimate  $\theta = \{\mathbf{A}, b_2(\mathbf{x}_t), \dots, b_{N-1}(\mathbf{x}_t)\}$ 
  - usually trained using **Expectation-Maximisation** (EM)

# Hidden Markov Model



- To design a classifier need to determine:
  - **transition matrix**: discrete state-space and allowed transitions
  - **state output distribution**: form of distribution  $p(\mathbf{x}_t|q_t)$
- Can be used as a generative classifier (HMM for each class  $\omega$ )

$$\hat{\omega} = \arg \max_{\omega} \{P(\omega|\mathbf{x}_1, \dots, \mathbf{x}_T)\} = \arg \max_{\omega} \{P(\omega)p(\mathbf{x}_1, \dots, \mathbf{x}_T|\omega)\}$$

need to be able to compute  $p(\mathbf{x}_1, \dots, \mathbf{x}_T|\omega)$  efficiently



- Viterbi Algorithm important technique (inc. for HMMs)
  - approximate likelihood as

$$p(\mathbf{x}_1, \dots, \mathbf{x}_T) = \sum_{\mathbf{q} \in Q_T} p(\mathbf{x}_1, \dots, \mathbf{x}_T, \mathbf{q}) \approx p(\mathbf{x}_1, \dots, \mathbf{x}_T, \hat{\mathbf{q}})$$

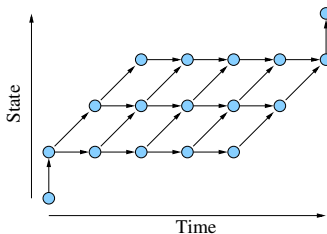
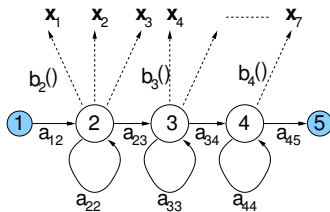
where

$$\hat{\mathbf{q}} = \{\hat{q}_0, \dots, \hat{q}_{T+1}\} = \arg \max_{\mathbf{q} \in Q_T} \{p(\mathbf{x}_1, \dots, \mathbf{x}_T, \mathbf{q})\}$$

- This yields:
  - an approximate likelihood (lower bound) for the model
  - the best state-sequence through the discrete-state space

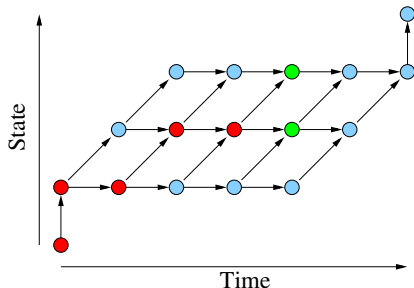
# Viterbi Algorithm

- Need to find the best state-sequence,  $\hat{q}$ ,
  - searching through all possible state-sequences impractical ...



- Consider generating the observation sequence  $x_1, \dots, x_7$ 
  - topology - 3 emitting states with strict **left-to-right**
  - representation of all possible state sequences on the right

# Extending Partial Paths with Time

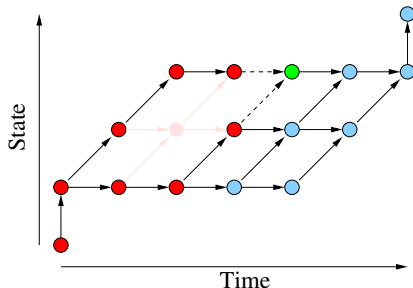


- Red partial path to time 4
- Green possible extensions

- Partial path state sequence  $\{1, 2, 2, 3, 3\}$  ( $\phi_3(4)$ ) extend:
  - stay in state  $s_3$  and generate  $\mathbf{x}_5$ :  $\log(a_{33}b_3(\mathbf{x}_5))$
  - move to state  $s_4$  and generate  $\mathbf{x}_5$ :  $\log(a_{34}b_4(\mathbf{x}_5))$
- Hence:

$$\phi_3(5) = \phi_3(4) + \log(a_{33}b_3(\mathbf{x}_5)), \quad \phi_4(5) = \phi_3(4) + \log(a_{34}b_4(\mathbf{x}_5))$$

# Best Partial Path to a State/Time



- Red possible partial paths
- Green state of interest

- Consider best partial path to state  $s_4$  at time 5 ( $\phi_4(5)$ )
  - move from state  $s_3$  and generate  $\mathbf{x}_5$ :  $\log(a_{34}b_4(\mathbf{x}_5))$
  - stay in state  $s_4$  and generate  $\mathbf{x}_5$ :  $\log(a_{44}b_4(\mathbf{x}_5))$
- Select “best”:

$$\phi_4(5) = \max \{ \phi_3(4) + \log(a_{34}), \phi_4(4) + \log(a_{44}) \} + \log(b_4(\mathbf{x}_5))$$

# Viterbi Algorithm for HMMs

- The Viterbi algorithm for HMMs can then be expressed as:

- Initialisation:** ( $\text{LZERO} = \log(0)$ )

$$\phi_1(0) = 0.0, \quad \phi_j(0) = \text{LZERO}, 1 < j < N, \quad \phi_1(t) = \text{LZERO}, 1 \leq t \leq T$$

- Recursion:**

for  $t = 1, \dots, T$   
for  $j = 2, \dots, N - 1$

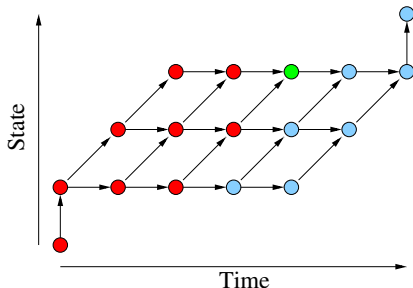
$$\phi_j(t) = \max_{1 \leq k < N} \{ \phi_k(t-1) + \log(a_{kj}) \} + \log(b_j(\mathbf{x}_t))$$

- Termination:**

$$\log(p(\mathbf{x}_1, \dots, \mathbf{x}_T, \hat{\mathbf{q}})) = \max_{1 < k < N} \{ \phi_k(T) + \log(a_{kN}) \}$$

- Can also store best previous state to yield best sequence  $\hat{\mathbf{q}}$ .

# All Paths to a State/Time



- Red possible partial paths
- Green state of interest

$$\text{LAdd}(a, b) = \log(\exp(a) + \exp(b))$$
$$\exp(\text{LAdd}(a, b)) = \exp(a) + \exp(b)$$

- Total path to state  $s_i$  at time  $t$  is  $\alpha_i(t)$ 
  - total path to state  $s_4$  at time 5 given by (see Viterbi)

$$\alpha_4(5) = \text{LAdd}(\alpha_3(4) + \log(a_{34}), \alpha_4(4) + \log(a_{44})) + \log(b_4(\mathbf{x}_5))$$

# Forward-Backward Algorithm

- $\alpha_j(t)$  is related to the **forward-probability** used to train HMMs
  - recursion for this form of model can be expressed as

$$\begin{aligned}\alpha_j(t) &= \log(p(\mathbf{x}_1, \dots, \mathbf{x}_t, q_t = s_j)) \\ &= \log\left(\sum_{k=1}^N \exp(\alpha_k(t-1) + \log(a_{kj}))\right) + \log(b_j(\mathbf{x}_t))\end{aligned}$$

- There's also a term related to the **backward-probability**
  - consider observation at time  $t$  **given** state  $s_j$ ,  $\beta_j(t)$

$$\begin{aligned}\beta_j(t) &= \log(p(\mathbf{x}_{t+1}, \dots, \mathbf{x}_T | q_t = s_j)) \\ &= \log\left(\sum_{k=1}^N \exp(\beta_k(t+1) + \log(a_{jk}) + \log(b_k(\mathbf{x}_{t+1})))\right)\end{aligned}$$

- The posterior required for EM can be expressed as

$$P(q_t = s_j | \mathbf{x}_1, \dots, \mathbf{x}_T) = \exp(\alpha_j(t) + \beta_j(t)) / Z, \quad Z = \sum_{i=1}^N \exp(\alpha_i(t) + \beta_i(t))$$

# Conditional Random Fields



# Discriminative Sequence Models

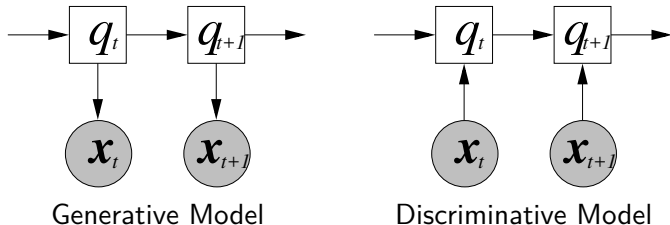
- HMMs compute the joint distribution of words and tags
  - the posterior probability of the tag sequence  $\mathbf{t}_{1:L}$

$$P(\mathbf{t}_{1:L}|\mathbf{w}_{1:L}) = \frac{P(\mathbf{t}_{1:L}, \mathbf{w}_{1:L})}{P(\mathbf{w}_{1:L})} = \frac{P(\mathbf{t}_{1:L}, \mathbf{w}_{1:L})}{\sum_{\tilde{\mathbf{t}}_{1:L}} P(\tilde{\mathbf{t}}_{1:L}, \mathbf{w}_{1:L})}$$

summation is over all possible  $L$ -length tag sequences

- a generative sequence model
- Discriminative sequence models directly model  $P(\mathbf{t}_{1:L}|\mathbf{w}_{1:L})$ 
  - required to be a valid probability mass function (PMF)
  - but only sequence probability needs to be a valid PMF
  - not the individual tag predictions, or word probabilities
- Here tags will be states ( $\mathbf{q}_{1:T}$ ) and words features ( $\mathbf{x}_{1:T}$ )
  - similar to discussion of HMMs

# Simple Discriminative Sequence Models



- Generative model (left) and discriminative model (right)
  - right BN a [maximum entropy Markov model](#)

$$P(q_0, \dots, q_T | \mathbf{x}_1, \dots, \mathbf{x}_T) = \prod_{t=1}^T P(q_t | q_{t-1}, \mathbf{x}_t)$$

state posterior probability ( $Z_t$  normalisation term time  $t$ ):

$$P(q_t | q_{t-1}, \mathbf{x}_t) = \frac{1}{Z_t} \exp \left( \sum_{i=1}^D \lambda_i f_i(q_t, q_{t-1}, \mathbf{x}_t) \right)$$

- State posteriors modelled in Maximum Entropy Markov model
  - can extend to the complete sequence

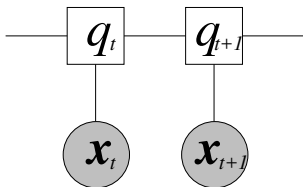
$$P(q_0, \dots, q_T | \mathbf{x}_1, \dots, \mathbf{x}_T) = \frac{1}{Z} \exp \left( \sum_{i=1}^D \lambda_i f_i(q_0, \dots, q_T, \mathbf{x}_1, \dots, \mathbf{x}_T) \right)$$

- Problem is that there are a vast number of possible features:

What features to extract from state/observation sequence?

- need to be able to handle variations in length of the sequence
- keep the number of model parameters  $\lambda$  reasonable

# (Simple) Linear Chain Conditional Random Fields



- Extract features based on undirected graph
  - conditional independence assumptions similar to HMM (though undirected)

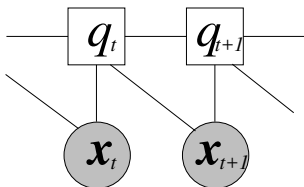
- Posterior model becomes

$$P(q_0, \dots, q_T | \mathbf{x}_1, \dots, \mathbf{x}_T) = \frac{1}{Z} \exp \left( \sum_{t=1}^T \left( \sum_{i=1}^{D_t} \lambda_i^t f_i(q_t, q_{t-1}) + \sum_{i=1}^{D_a} \lambda_i^a f_i(q_t, \mathbf{x}_t) \right) \right)$$

- $D_t$  number of transition style features with parameters  $\lambda^t$
- $D_a$  number of acoustic style features with parameters  $\lambda^a$
- Directly related to (unnormalised) HMM parameters

$$\sum_{i=1}^{D_t} \lambda_i^t f_i(s_i, s_j) \log(a_{ij}) \quad \text{and} \quad \sum_{i=1}^{D_a} \lambda_i^a f_i(s_i, \mathbf{x}_t) \log(b_i(\mathbf{x}_t))$$

# Linear Chain Conditional Random Fields



- Extract features based on undirected graph
  - conditional independence assumptions extended to previous state

- Posterior model becomes

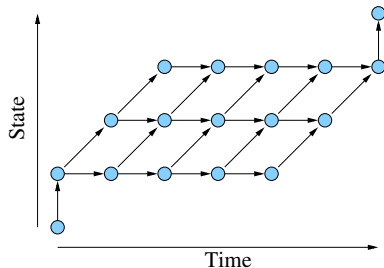
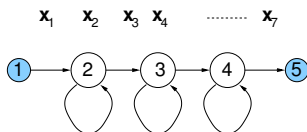
$$P(q_0, \dots, q_T | \mathbf{x}_1, \dots, \mathbf{x}_T) = \frac{1}{Z} \exp \left( \sum_{t=1}^T \left( \sum_{i=1}^D \lambda_i f_i(q_t, q_{t-1}, \mathbf{x}_t) \right) \right)$$

- More interesting than HMM-like features
  - features the similar to (general) MaxEnt Markov model
  - BUT normalised **globally** not **locally**

# Normalisation term

- Need to be able to compute the normalisation term efficiently
  - initially consider the **simple linear chain case**

$$Z = \sum_{\mathbf{q} \in \mathbf{Q}_T} \exp \left( \sum_{t=1}^T \left( \sum_{i=1}^{D_t} \lambda_i^t f_i(q_t, q_{t-1}) + \sum_{i=1}^{D_a} \lambda_i^a f_i(q_t, \mathbf{x}_t) \right) \right)$$



- Left-to-right topology and observation sequence  $\mathbf{x}_1, \dots, \mathbf{x}_7$ 
  - use equivalent of forward-backward algorithm

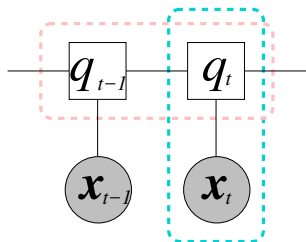
# General Sequence CRFs (reference)

- General CRF use undirected graphical model to define features
  - need to be able to handle sequence data - **dynamic CRF**
  - undirected graph **repeated** each time instance - set of cliques  $\mathcal{C}$
- The posterior probability for this form of model is

$$P(q_0, \dots, q_T | \mathbf{x}_1, \dots, \mathbf{x}_T) = \frac{1}{Z} \exp \left( \sum_{t=1}^T \sum_{\mathcal{C} \in \mathcal{C}} \lambda_{\mathcal{C}}^T \mathbf{f}(\mathbf{q}_{\mathcal{C}t}, \mathbf{x}_1, \dots, \mathbf{x}_T, t) \right)$$

- $\lambda_{\mathcal{C}}^T$  **time-independent** parameters associated with clique  $\mathcal{C}$
- $\mathbf{f}(\mathbf{q}_{\mathcal{C}t}, \mathbf{x}_1, \dots, \mathbf{x}_T, t)$  **time-dependent** features extracted from clique  $\mathcal{C}$  with **time-dependent** label sequence  $\mathbf{q}_{\mathcal{C}t}$

## Example of a Sequence CRF (reference)



- Cliques associated with linear CRF

$$\mathcal{C} = \{\mathcal{C}_1, \mathcal{C}_2\}$$

1. transitions:  $\mathcal{C}_1 = \{q_t, q_{t-1}\}$

2. acoustics:  $\mathcal{C}_2 = \{q_t, \mathbf{x}_t\}$

- Posterior model for the simple linear chain CRF

$$\begin{aligned} P(q_0, \dots, q_T | \mathbf{x}_1, \dots, \mathbf{x}_T) &= \frac{1}{Z} \exp \left( \sum_{t=1}^T \sum_{\mathcal{C} \in \mathcal{C}} \lambda_{\mathcal{C}}^{\top} \mathbf{f}(\mathbf{q}_{\mathcal{C}_t}, \mathbf{x}_1, \dots, \mathbf{x}_T, t) \right) \\ &= \frac{1}{Z} \exp \left( \sum_{t=1}^T \left( \lambda^{\text{t}\top} \mathbf{f}(q_t, q_{t-1}) + \lambda^{\text{a}\top} \mathbf{f}(q_t, \mathbf{x}_t) \right) \right) \end{aligned}$$



- Training for CRFs is normally fully observed

training observation sequence     $\mathbf{x}_1, \dots, \mathbf{x}_T$

training label sequence             $y_1, \dots, y_T$

- where  $y_\tau \in \{\omega_1, \dots, \omega_K\}$
- No need to use EM (or related approaches)
- Need to find the model parameters  $\lambda$  so that

$$\begin{aligned}\hat{\lambda} &= \arg \max_{\lambda} \{P(y_1, \dots, y_T | \mathbf{x}_1, \dots, \mathbf{x}_T; \lambda)\} \\ &= \arg \max_{\lambda} \left\{ \frac{1}{Z} \exp \left( \sum_{i=1}^D \lambda_i f_i(\mathbf{x}_1, \dots, \mathbf{x}_T, y_1, \dots, y_T) \right) \right\}\end{aligned}$$