

① For what kind of problems do we use HMMs?

Consider an example

Weather
in Lugano
can be

$\mathcal{S} \in \{\text{Rainy, Cloudy, Sunny}\}$

$\mathcal{S} \in \{R, C, S\}$

① Therefore, weather takes values from a discrete set of States.

① Starting probability of a weather is also a number

$P(Q_1 = R) = 0.1 \rightarrow$ initial probabilities

① Transition ~~probabilities~~ from day to day is also expressed as prob.

$P(Q_{t+1} = \text{Sunny} | Q_t = \text{Sunny}) = 0.6 \rightarrow$ transition probability

① ~~Assume that~~ Assume that there is some guy living in Lugano, and he performs activities based on weather forecast from, again a discrete set of activities (observations)

$\mathcal{A} \in \{\text{~~Study~~, Play}\}$

① Again, activity is ~~a~~ a probabilistic function of the state (weather)

$P(A_t = \text{~~Study~~} | Q_t = \text{Sunny}) = 0.4 \rightarrow$ observation probability

① Assume that Matteo's gf doesn't live in Lugano, so she might have only partial knowledge of the weather & Matteo's activities.

\rightarrow knows weather, wants to guess activities

\rightarrow knows activities, wants to determine most probable weather

① HMM's give us principled way of dealing with following problems

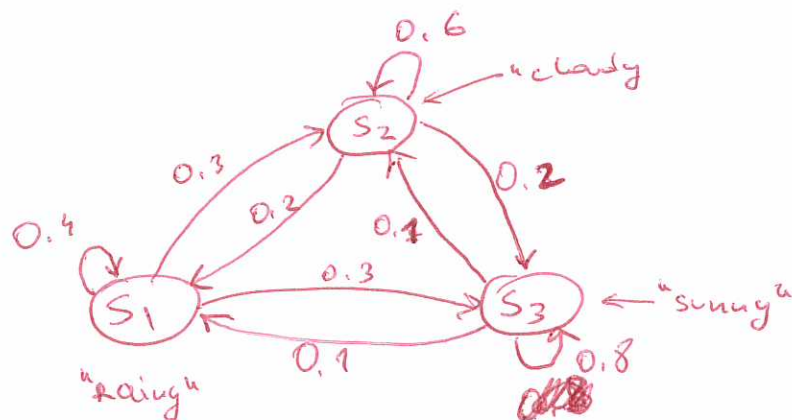
① $P(Q_1 \dots Q_T | \lambda) \rightarrow$ Given model λ , compute most probable obs. seq.
 \rightarrow use simple forward pass

② $P(Q_1 \dots Q_T | A_1 \dots A_T, \lambda) \rightarrow$ Given model & observations, compute most probable weather sequence
 \rightarrow Viterbi algorithm

③ $\text{Argmax}_{\lambda} P(Q_1 \dots Q_T | \lambda) \rightarrow$ Given ~~a given~~ ^{obs. seq.} find best model λ s.t. this is maximized
Not the best example, better \rightarrow EM algorithm \rightarrow speech recognition ①

Markov system

Chain with N (in our case 3) states



⊕ The model is said to be in ONE of the states at each time t
model's state at $t \rightarrow g_t \in \{S_1, \dots, S_N\}$
time t
⚡ note: discrete time steps

⊕ Between each time step state is chosen randomly
therefore the model is stochastic

⊕ Note: probabilities which go out of a state sum to 1.

⊕ What is Markov property:

~~Assume~~ Assume you observed sequence of states

$$S_1, \dots, S_t, S_{t+1}, S_{t+2}, \dots, S_N, \quad p_i, i \in \{1, \dots, N\}$$

$$P(g_{t+1} = S_j \mid g_t = S_i, g_{t-1} = S_k, \dots, g_1 = S_p) = P(g_{t+1} = S_j \mid g_t = S_i)$$

Probability of transitioning to a state depends only on the current state & NOT on past states.

Does not always hold in the real world, but it simplifies the problem & in many cases works quite well.

⊕ Let's write transition probability from S_i to S_j

$$a_{ij} = P(g_{t+1} = S_j \mid g_t = S_i)$$

⇒ defines state transition matrix

$$A = \begin{matrix} & \begin{matrix} \text{from R} & \text{from C} & \text{from S} \end{matrix} \\ \begin{matrix} \text{to R} \\ \text{to C} \\ \text{to S} \end{matrix} & \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.2 & 0.6 & 0.2 \\ 0.1 & 0.1 & 0.8 \end{bmatrix} \end{matrix} \xrightarrow{N \times N} \frac{\sum}{N=3} = 1$$

$$\boxed{\sum_{i=1}^N a_{ij} = 1}$$

⊕ How to compute prob. of a sequence of states

$$\begin{aligned} P(g_1 \dots g_t) &= P(g_t \mid g_{t-1}, \dots, g_1) P(g_{t-1}, \dots, g_1) \quad (\text{chain rule } P(A, B) = P(A|B)P(B)) \\ &= P(g_t \mid g_{t-1}) P(g_{t-1}, \dots, g_1) \quad (\text{Markov property}) \\ &= P(g_t \mid g_{t-1}) P(g_{t-1} \mid g_{t-2}) \dots P(g_2 \mid g_1) P(g_1) \end{aligned}$$

⊕ To proceed need to define

initial state probabilities $\pi_i = P[g_1 = S_i]$, $i \in \{1, \dots, N\}$

$$\pi_i = \{0.1, 0.1, 0.8\}$$

Rainy cloudy sunny

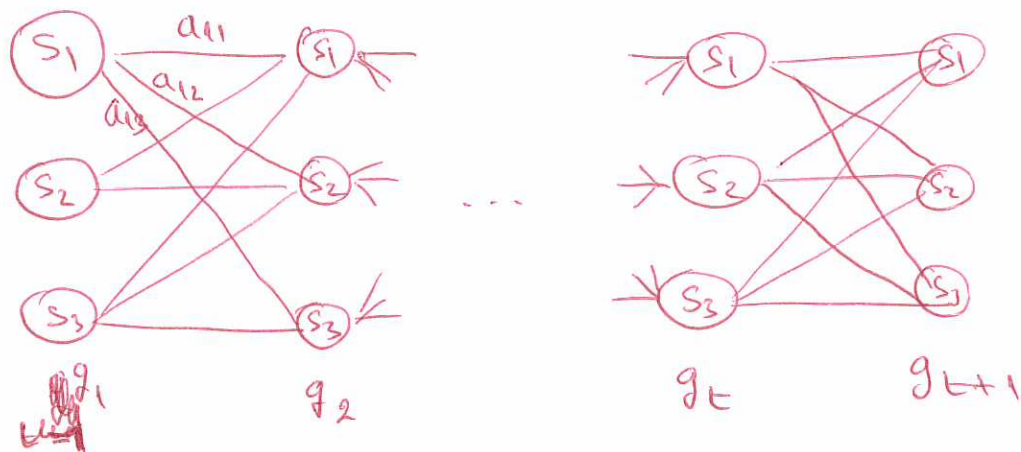
⊕ Example: given state sequence $Q = S, S, S, R, R$

$$\begin{aligned} P(Q \mid \text{model}) &= P(S) \cdot P(S \mid S) P(S \mid S) P(R \mid S) \cdot P(R \mid R) \\ &= \underline{0.8} \cdot 0.8 \cdot 0.8 \cdot 0.1 \cdot 0.4 = \dots \end{aligned}$$

⊕ Okay, this was pretty simple

How about computing $P(q_t = s_i) = ?$

trellis diagram



2 ways: 1) $P(q_t = s_i) = \sum_{Q \in \text{path of length } t, \text{ ending in } s_i} P(Q)$ Slow, stupid $O(TN^2)$

Explain how this looks on trellis diagram

2) Dynamic

Define $p_t(i) = P(q_t = s_i) \leftarrow$ our value

Observe that ~~we can compute~~ if we have $p_t(i), \forall i$

we can compute p_{t+1} ~~recursively~~ with dynamic prog.

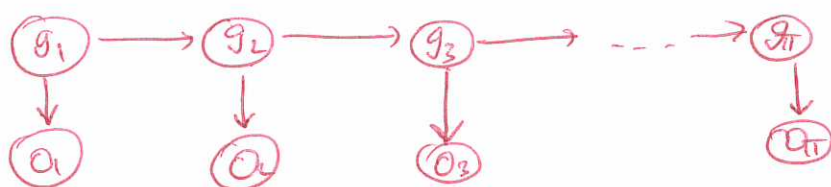
$$\begin{aligned} p_{t+1}(j) &= \sum_{i=1}^N P(q_{t+1} = s_j, q_t = s_i) = \\ &= \sum_{i=1}^N P(q_{t+1} = s_j | q_t = s_i) \cdot P(q_t = s_i) \\ &= \sum_{i=1}^N \underbrace{a_{ij}}_{\text{transition prob}} \cdot \underbrace{p_t(i)}_{\text{previous prob}} \quad O(TN^2) \end{aligned}$$

Trade computational capacity for memory storage.

Let's extend this to HMMs

⊕ in Markov system observed directly the states

⊕ in HMMs ~~we don't observe the states~~ set of observations ~~are different~~ is different than ~~states~~ states (which is hidden)



↑ observations noisy
state hidden

in our example

② $P(O_t = \text{study} | g_t = \text{rain}) = 0.8$

① $g_t \in \{\text{raining, cloudy, sunny}\}$
 $O_t \in \{\text{study, play}\}$

⊕ HMM: ~~stochastic~~ model $\lambda = (N, M, \pi, A, B)$

① $N = \{s_1, \dots, s_N\} \leftarrow$ number of states

② $M = \{o_1, \dots, o_M\} \leftarrow$ number of observations note: often $M \neq N$

③ $\pi_i = P(g_1 = s_i) \rightarrow \sum_{i=1}^N \pi_i = 1 \leftarrow$ initial state probabilities

④ ~~$A = \{a_{ij}\}$~~ $a_{ij} = P(g_{t+1} = s_j | g_t = s_i)$

① $A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1N} \\ \vdots & \vdots & \ddots & \vdots \\ a_{N1} & \dots & \dots & a_{NN} \end{bmatrix}_{N \times N}$ ↑ state transition probabilities

⑤ $B = \begin{bmatrix} b_1(1) & b_1(2) & \dots & b_1(M) \\ \vdots & \vdots & \ddots & \vdots \\ b_N(1) & b_N(2) & \dots & b_N(M) \end{bmatrix}_{N \times M}$ ← observation symbol probability distribution

$b_i(k) = P(O_t = k | g_t = s_i)$

~~Sometimes~~ written as $\lambda = (\pi, A, B)$

Often

- 3 problems are:
- ① $P(O|\lambda)$ efficiently compute observation sequence
 - ② $P(Q|O, \lambda)$ optimal state sequence given O
 - ③ $\arg \max_{\lambda} P(O|\lambda)$ adjust model parameters to maximize $P(O|\lambda)$

④ Let's look at the first problem:

Compute probabilities of observation sequence given model λ

$$P(O|\lambda) = P(O_1, O_2, \dots, O_T|\lambda) = \sum_{Q \in \text{all } T} P(O \wedge Q|\lambda)$$

~~enumeration~~

$$= \sum_{Q \in \text{all } T} P(O|Q, \lambda) P(Q|\lambda)$$

$$P(O|Q, \lambda) = \prod_{t=1}^T P(O_t | q_t, \lambda) =$$

$$= b_{q_1}(O_1) \cdot b_{q_2}(O_2) \cdot \dots \cdot b_{q_T}(O_T)$$

$$P(Q|\lambda) = \pi_{q_1} \cdot a_{q_1 q_2} \cdot a_{q_2 q_3} \cdot \dots \cdot a_{q_{T-1} q_T}$$

$$\Rightarrow P(O|\lambda) = \sum_{Q \in \text{all } T} P(O|Q, \lambda) P(Q|\lambda) =$$

$$= \sum_{q_1, q_2, \dots, q_T} \pi_{q_1} \cdot \underline{b_{q_1}(O_1)} \cdot \underline{a_{q_1 q_2}} \cdot \underline{b_{q_2}(O_2)} \cdot \dots \cdot \underline{a_{q_{T-1} q_T}} \cdot \underline{b_{q_T}(O_T)}$$

$O(TN^T)$

Similar to what we had for normal Markov chain,
just that we have one additional component,
Observation probabilities.

Smarter way to calculate this

Forward algorithm

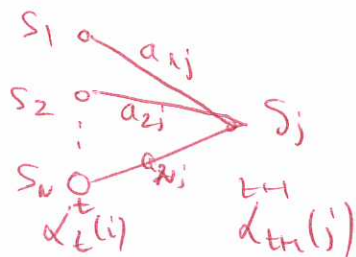
Define $\alpha_t(i) = P(o_1 \dots o_t \wedge q_t = s_i | \lambda)$

"we have observed o_1, \dots, o_t & ended up in state s_i "

① initialization

$$\begin{aligned}\alpha_1(i) &= P(o_1 \wedge q_1 = s_i) = \\ &= P(o_1 | q_1 = s_i) P(q_1 = s_i) \\ &= b_i(o_1) \cdot \pi_i\end{aligned}$$

② update ~~iteration~~ step



$$P(A|B) = P(A)$$

$$P(A, B) = P(A) \cdot P(B)$$

$$P(A, B) = P(A|B) \cdot P(B)$$

$$\alpha_{t+1}(j) = P(o_1 \dots o_t o_{t+1} \wedge q_{t+1} = s_j)$$

$$= \sum_{i=1}^N P(o_1 \dots o_t \wedge q_t = s_i \wedge o_{t+1} \wedge q_{t+1} = s_j)$$

$$P(A|B) = P(A|B) \cdot P(B)$$

$$P(A|B) = P(A) \rightarrow = \sum_{i=1}^N P(o_{t+1}, q_{t+1} = s_j | o_1 \dots o_t \wedge q_t = s_i) P(o_1 \dots o_t \wedge q_t = s_i)$$

$$P(A|B, C) = P(A|C) \rightarrow = \sum_{i=1}^N P(o_{t+1}, q_{t+1} = s_j | q_t = s_i) \cdot \alpha_t(i)$$

if $P(A|B) = P(A)$

$$= \sum_{i=1}^N P(o_{t+1} | q_{t+1} = s_j) P(q_{t+1} = s_j | q_t = s_i) \cdot \alpha_t(i)$$

$$= \sum_{i=1}^N a_{ij} \cdot b_j(o_{t+1}) \cdot \alpha_t(i)$$

③ termination

$$P(o_1 \dots o_T | \lambda) = \sum_{i=1}^N P(o_1 \dots o_T \wedge q_T = s_i | \lambda)$$

$$= \sum_{i=1}^N \alpha_T(i) \sim O(NT^2)$$

④ How to calculate

$$P(q_t = s_i | o_1 \dots o_t) = \frac{P(o_1 \dots o_t \wedge q_t = s_i)}{P(o)} = \frac{\alpha_t(i)}{\sum_{j=1}^N \alpha_t(j)}$$

⑦

Problem 2: Given observations, find the most likely state sequence Q

$$\max_Q P(Q | O, \lambda)$$

Viterbi algorithm

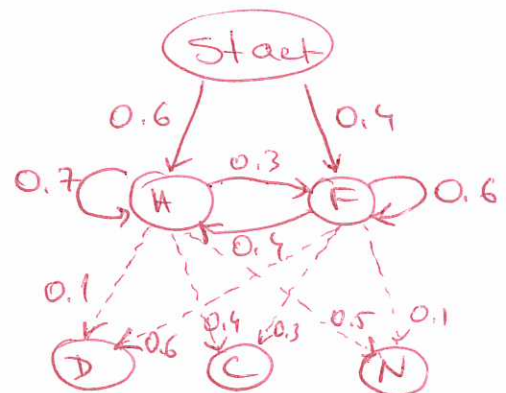
Let's explain it on an example, & later give general definition

Assume 2 states: $S_1 = \{H, F\}$

$O \in \{C, N, D\}$
C N D

$$A = \begin{matrix} & \begin{matrix} \text{Healthy} \\ \text{Fever} \end{matrix} \\ \begin{matrix} \text{Healthy} \\ \text{Fever} \end{matrix} & \begin{pmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \end{pmatrix} \end{matrix}$$

$$\Pi = \begin{matrix} & \begin{matrix} H & F \end{matrix} \\ \begin{matrix} H & F \end{matrix} & \begin{pmatrix} 0.6 & 0.4 \end{pmatrix} \end{matrix}$$



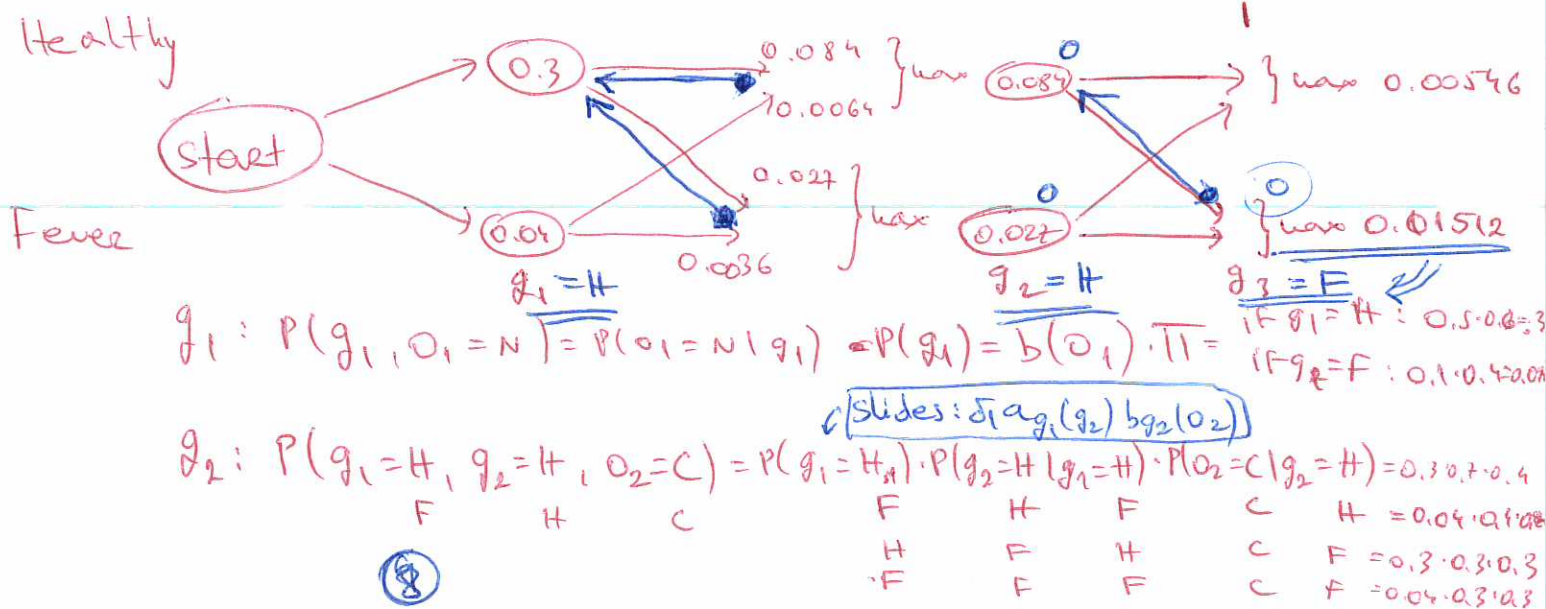
Assume you observe a sequence N, C, D .

Find most probable sequence of states.

$$\max_{g_1, g_2, g_3} P(g_1, g_2, g_3, O_1=N, O_2=C, O_3=D | \lambda)$$

① Find most probable individual states
② or most probable paths / triplets of states

If you look at all possible set of states & find max prob path, again you have many calculations.



In case of a long sequence

=> lots of factors => smaller & smaller numbers

=> convert to log probabilities & sum!

$$\log P(q_1) P(q_2 | q_1) P(o_2 | q_2) = \log P(q_1) + \log P(q_2 | q_1) + \log P(o_2 | q_2)$$

Viterbi formula

$$\delta_t(i) = \max_{q_1 \dots q_{t-1}} P(q_1 \dots q_{t-1} \wedge q_t = S_i \wedge o_1 \dots o_t)$$

= the maximum probability of ending up in state S_i at time t and producing $o_1 \dots o_t$

$$\text{wpp}(i) = \arg \max_{q_1 \dots q_{t-1}} \delta_t(i)$$

$$\delta_1(j) = \pi_j b_j(o_1)$$

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

$$\psi_t(j) = \arg \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}] \quad \begin{matrix} t \in [2, T] \\ j \in [1, N] \end{matrix}$$

↓ keeps track of argument which maximizes $\delta_t(j)$ Q: why don't we have b here?
↳ b/c it wouldn't change anything as we look for arg, not value!

Once we come to the end, we look at maximal probability, take ~~that~~ the state which generated this probability as last, and do backtracking based on ψ_t 's

Q: how is this different from the forward pass?

Difference to forward pass: keeps track of

backtrack variables that lead to maximal probability.

Used also in telecommunications to decode a sequence

starts at

01
10
11

→ noisy channel →



3rd problem of HMMs

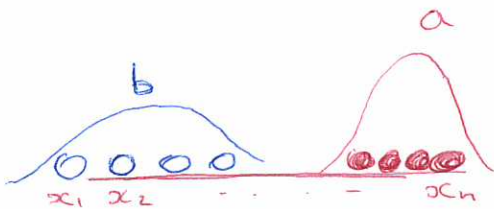
Given obs. O

Find λ s.t. $\max_{\lambda} P(O|\lambda)$

- ↳ no known way to analytically compute such model
- ↳ no optimal way of estimating model parameters

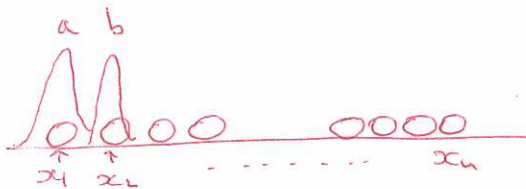
Approximate ways to choose λ s.t. $\max_{\lambda} P(O|\lambda)$:
① iterative procedure
EM (Baum-Welch algo)
② gradient descent

1-D EM algorithm



+ Labeled \rightarrow easy to fit the two distributions
+ ~~we~~ have distributions \rightarrow easy to classify

Not clear what to do if data is unlabeled!



$P(a|x_1) > P(b|x_1)$ "belongs more to a"

$P(a|x_2) < P(b|x_2)$ "belongs more to b"

To solve this issue you can use iterative EM procedure

EM algorithm used in:

① Example here: Gaussian Mixture model

① K-means clustering (special case of GMMs) \leftarrow covariances diagonal equal & such

② don't mix with
K-nearest neighbor!
 \rightarrow supervised!

① HMMs

\rightarrow minimize pairwise distance of points in the same cluster
 \downarrow
initialize $\sum_{x_i \in C_j} \|x_i - \mu_j\|^2$
mean pts \rightarrow minimize this distance to μ_j

1D-EM

① initialize randomly $\mu_a, \mu_b, \sigma_a, \sigma_b, P(a), P(b)$

means and variances of the Gaussians

for each of the "clusters"
 $P(a), P(b)$: priors (numbers)
 how likely is a data point to belong to classes a and b

② E-step: compute the belief for each x_i using our current Gaussian parameters

where is each datapoint expected to belong?

$$a_i = P(a|x_i) = \frac{P(x_i|a) \cdot P(a)}{P(x_i|a) \cdot P(a) + P(x_i|b) \cdot P(b)}$$

As we have assumed we are working with Gaussian dist:

$$P(x_i|a) = \frac{1}{\sqrt{2\pi\sigma_a^2}} e^{-\frac{(x_i - \mu_a)^2}{2\sigma_a^2}}$$

$$b_i = 1 - a_i$$

③ M-step: Re-estimate the Gaussian parameters using all data points

$$\mu_a = \frac{a_1 x_1 + \dots + a_n x_n}{a_1 + \dots + a_n}$$

$$\sigma_a^2 = \frac{\sum_{i=1}^n (x_i - \mu_a)^2 \cdot a_i}{\sum_{i=1}^n a_i}$$

$$P(a) = \frac{a_1 + \dots + a_n}{n}$$

$$P(b) = 1 - P(a)$$

analytical solutions to MLE

$$\text{MLE } \theta_{\text{MLE}} = \arg \max_{\theta} P(x|\theta) = \arg \max_{\theta} \log P(x|\theta) = \arg \max_{\theta} \sum_i \log P(x_i|\theta)$$

↑ maximum likelihood estimate

$$\theta = \{\mu_a, \mu_b, \sigma_a, \sigma_b, P(a), P(b)\}$$

can optimize it also via GD

It is just a special case of MAP (Maximum a posteriori estimate)

We can see this through Bayes theorem

$$\text{Posterior} \rightarrow P(\theta|x) = \frac{P(x|\theta)P(\theta)}{P(x)} \leftarrow \begin{matrix} \text{likelihood} \\ \text{prior} \end{matrix}$$

$P(x) \propto \text{const}$ (proportional)

$$\Rightarrow \theta_{\text{MAP}} = \arg \max_{\theta} \sum_i \log P(x_i|\theta) \cdot P(\theta)$$

$$\text{if prior is uniform } P(\theta) = \text{const} \Rightarrow \theta_{\text{MAP}} = \arg \max_{\theta} \sum_i \log P(x_i|\theta)$$

$$= \theta_{\text{MLE}}$$

④ repeat 2 & 3 until converged

Notes: ④ can lead to local maxima, no guaranteed convergence in general case

④ in practice you typically make several runs with different initializations

④ Question also how many clusters to use

Baum-Welch algorithm (EM for HMMs)

We need 4 variables: $\alpha_t(i)$, $\beta_t(i)$, $\gamma_t(i)$, $\xi_t(i,j)$

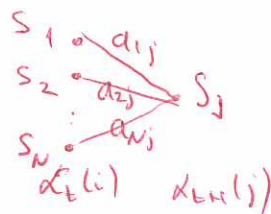
① Previously we have considered

$$\alpha_t(i) = P(O_1, \dots, O_t, q_t = S_i)$$

which can be computed recursively as

$$\alpha_1(i) = \pi_i \cdot b_i(O_1)$$

$$\alpha_{t+1}(i) = \sum_{j=1}^N \alpha_t(j) \cdot a_{ji} \cdot b_i(O_{t+1})$$



② Similarly, we introduce backward variable

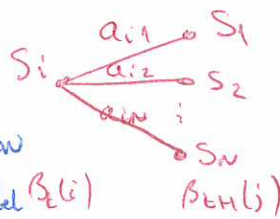
$$\beta_t(i) = P(O_{t+1}, \dots, O_T, q_t = S_i)$$

PROBABILITY OF THE PARTIAL OBSERVATION

SEQUENCE, GIVEN THE STATE S_i @ t & model $\beta_t(i)$

$$\beta_T(i) = 1$$

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



③ Probability of being in state S_i @ time t , given O

$$\gamma_t(i) = P(q_t = S_i | O) = \frac{P(q_t = S_i, O_1, \dots, O_T)}{P(O_1, \dots, O_T)} = \frac{\alpha_t(i) \beta_t(i)}{\sum_{i=1}^N \alpha_t(i) \beta_t(i)}$$

Note $\sum_i \gamma_t(i) = 1$

↑

for sure we are in state S_i

Most likely state @ time t : $q_t = \arg \max_i \gamma_t(i)$

What is the problem with this?

+ Only ~~individualizes~~ individual states, not the path as a whole

+ In the resulting path there could be impossible transitions (if $a_{ij} = 0$)

Solutions: + could observe pairs of consecutive states

+ or triplets, and so on...

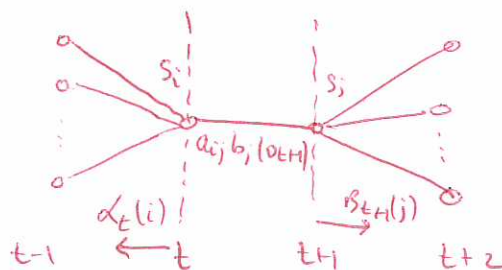
+ ~~Viterbi~~ Viterbi algorithm gives us the most prob. path.

④ Probability of being in state S_i at time t & S_j at $t+1$

$$\xi(i,j) = P(q_t = S_i, q_{t+1} = S_j | O, \lambda)$$

$$= \frac{P(q_t = S_i, q_{t+1} = S_j, O_1, \dots, O_T)}{P(O_1, \dots, O_T)}$$

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



Note: $\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j)$

$\sum_{t=1}^{T-1} \gamma_t(i) = \text{expected number of transitions from } S_i$

$\sum_{t=1}^{T-1} \xi_t(i, j) = \text{expected number of transitions from } S_i \text{ to } S_j$

Baum-Welch algo

① Start with initial $\lambda = (\pi, A, B)$

② E-step: Calculate $\gamma_t(i), \xi_t(i, j)$

③ M-step: Determine new model as

⊕ $\pi_i = \gamma_1(i) \leftarrow \text{expected number of times in state } S_i$

⊕ $a_{ij} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)} = \frac{\text{expected number of transitions from } S_i \text{ to } S_j}{\text{expected number of ~~transitions~~ times in } S_i}$

⊕ $b_j(k) = \frac{\sum_{t=1, s.t. O_t=k}^T \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)} = \frac{\text{expected number of times in } S_j \text{ \& observed symbol } k}{\text{expected number of times in } S_j}$

④ Repeat 2 & 3 until converged