(1) Q: List probabilities HMM, MEMM, and CRF need to estimate
during training. (how these probabilities sum to constant 1).

(point out sample space of each probability distribution).

Ans:

① MEMM:

- to compute maximum entropy we need to define a window first.
we can compute $q_t = c$ using last state $q_{t-1}$ or last two
states $q_{t-1}, q_{t-2}$ and so on. I assume we are using
bigram model that predict current $q_t$ based on previous
state $q_{t-1}$. so we have:

$$Pr(q_t = c | q_{t-1}, O_t ; \theta^c) = \frac{1}{Z(q_{t-1}, O_t)} \exp \left\{ \sum_{i=1}^{d} \theta_i^c f_i(q_{t-1}, O_t, q_t = c) \right\}$$

and we have optimize log-likelihood to find $\theta^*$.

- so with bigram MEMM or window we have to compute
$$\sum_{c = \{all\ tags\}}^{T} Pr(q_t = c | q_{t-1}, O_t ; \theta^c) = 1 \quad \boxed{T\ times} \quad T = number\ of\ sequences)$$

① cont'd

CRF:

For CRF we have $\frac{1}{Z(\theta)}$ of $\left(\begin{smallmatrix} y_1 \\ \end{smallmatrix}\right)$ $\beta t$ $\left(\begin{smallmatrix} y_2 \\ \end{smallmatrix}\right)$ $\times exp(S(y_1, y_2, \theta))$

for three sequences as an example. So we need to compute $TN$ probabilities for $\alpha$, same number for $\beta$, and $(T-1)$ possible triangle cliques for each $\circ$ and since we have $T$ sequences, its $(T-1) \times T$.

total number won't be $T(T-1) \times 2TN \simeq \boxed{T^2 + 2TN}$

probabilities.

① Con'td

Hmm parameters:
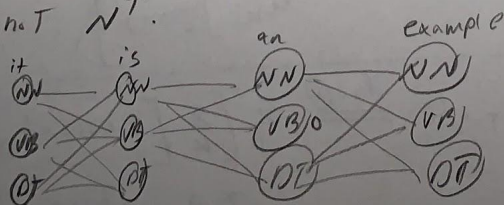
We can use Expectation maximization to learn $A$, $B$, and $\mathcal{M}$.

in E step:

each TxN

computation
$$\begin{cases} \alpha_t(j) \text{ for all tags (js) and all sequences } t. \\ \beta_t(j) \quad \text{''} \quad \text{''} \quad \text{''} \quad \text{''} \quad \text{''} \quad \text{''} \quad \text{''} . \end{cases}$$

M-step:
$$\overbrace{Pr(\ell_t = i, \ell_{t+1} = j \mid 0, \lambda)}^{T^2 \text{ computation}} \longrightarrow \text{sum to 1}$$

$$\underbrace{Pr(\ell_t = i \mid 0, \lambda)}_{T \text{ computation}} \longrightarrow \text{sum to 1}$$

then $T^2 + T + 2TN$ probabilities should be computed.

* if we look at a lattice of Hmm with 3 tags like below we see for $T=4$ we have to compute $3^4$ probabilities, but since lots of these are zero, and we avoid those paths, the real number is not $N^T$.

②

I can think of starting with capital letter for a feature.

So $f_q : c = NNP$ And $O_L = X$ And starts with capital letter.

In this case if term $X$ did not appear in our training set, the model cannot generalize and correctly classify it, because of $O_L = X$ in $f_q$. removing that, however makes the accuracy falls below 100%, due to tagging any word starting with a capital letter.

★ maybe similar features such as capital letter can be added to $f_q$.

③

$$\theta'_{f_i} = \theta \mid i + \eta \frac{\partial}{\partial \theta_c} \left[ -L \cdot g \; Pr \left( q_t = c \mid o_t ; \theta \right) \right]$$

$$= \theta \mid i - \eta \left\{ - \left\{ [q_t = c] - Pr \left( q_t = c \mid o_t ; \theta \right) \right\} f_i \left( o_t, q_t = c \right) \right.$$

$$= \theta \mid i - \eta \left\{ - \left\{ [q_t = N] - pr \left( q_t = N \mid o_t = race ; \theta \right) \right\} f_i \left( o_t = 'race', \; q_t = N \right) \right.$$

$$\underbrace{\qquad\qquad\qquad 0 \swarrow \qquad\qquad\qquad}_{①}$$

②

$$Pr \left( q_t = N \mid o_t = race ; \theta \right) = \frac{1}{z} \; exp \left\{ \sum_{i=1}^{6} \theta^c_i \, f_i \left( o_t = race, \, q_t = N \right) \right\}$$

$$z = \sum_{c' \in \{all \; tys\}} exp \left\{ \sum_{i=1}^{6} \theta^{c'}_i \, f_i \left( o_t = race, \, q_t = c' \right) \right\}$$

$$c' = NGN$$
$$= \left[ exp \left\{ 0.8 - 1.3 \right\} \right] + \left[ exp \left\{ 0.8 + 0.9 + 0.1 \right\} \right] +$$

$$c' = VBG$$
$$\left[ exp \left\{ 0 \right\} \right] = 0.607 + 2.484 + 1.0 = 4.091 \quad ③$$

$$from \; ② \; and \; ③ : pr \left( q_t = N \mid o_t = race ; \theta \right) = \frac{1}{4.091} \times exp \left( -0.5 \right)$$

$$= \frac{0.607}{4.091} = 0.148$$

④

③ cont'd

from ④ and ① :

$$\theta_i = \theta_i - \eta \left( - \left[ \dots 0.148 \right] \right) f_i (0_{E=race}, \, \ell_{E}=N)$$

$$\theta_1 = 0.8 - \eta \, (0.148) \cdot 1 = 0.8 - 0.148 \, \eta$$

$$\theta_2 = 1 - \eta \, (0.148) \times 0 = 1$$

$$\theta_3 = -2 - \eta \, (\eta) \times 0 = -2$$

$$\theta_4 = 3 - \eta \, (\eta) \times 0 = 3$$

$$\theta_5 = 0.1 - \eta \, (\eta) \times 0 = 0.1$$

$$\theta_6 = -1.3 - \eta \, (0.148) \cdot 1 = -1.3 - \eta \, (0.148)$$

\* the parameters/weights $\theta_i$ for these features ($f_1$ and $f_6$)
that were wrong are getting smaller due to the wrong classification
and <u>tagging</u> NN instead of VB.

on the other hand, $f_2$, $f_4$, and $f_5$ for $C=VB$ should increase

for $\theta VB$.

$$a^t = b + W h^{t-1} + U x^t$$

$$h^t = \tanh(a^t)$$

$$o^t = C + V h^t$$

$$\hat{y}^t = \text{softmax}(o^t)$$

$o^t$: len of 2

$h^t$: len of 3

$x^t$: len of 2

$b, W, U, C, V$   size?

- $h^t$ is of length 3, so $a^t$ should be as of length 3 as well.

So, $b \longrightarrow 3\times1$

$W \longrightarrow 3\times3$

$U \longrightarrow 3\times2$

- since $h^t$ is of length 3 and $o^t$ has length of 2, then

$C \longrightarrow 2\times1$  and  $V \longrightarrow 2\times3$

(5)

$$b = [1, 1, 1]$$

$$W = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \quad U = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$c = [1, 1] \qquad V = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

$$h^0 = [1 \ 1 \ 1]^T$$

$$x^1 = [1, 1]^T, \quad x^2 = [2, 2]$$

$$\underline{\quad t = 1 \quad}$$

$$x^1: \quad a^1 = [1, 1, 1] + \begin{bmatrix} 3 & 3 & 3 \end{bmatrix} + \begin{bmatrix} 2 & 2 & 2 \end{bmatrix}$$

$$= [6 \ 6 \ 6]$$

$$h^1 \simeq [1 \ 1 \ 1]^T \ \tanh(a^1)$$

$$o^1 = [1 \ 1] + [3 \ 3] = [4 \ 4]$$

$$\hat{y}^1 = \text{S.ftmax}([4 \ 4]) = [0.5 \ 0.5]$$

+ since $h^1 = h^0$, all the results would be the same for $t = 2$.

$$a^1 = a^2$$
$$o^1 = o^2$$
$$h^1 = h^2$$
$$\hat{y}^1 = \hat{y}^2$$

* since $\tanh(6) \simeq 3.9999\ldots$
values are round.

⑤ cn'td

$$\frac{t=4}{x^2}:$$

$$a^1 = [8 \ 8 \ 8]$$

$$h^1 = \tanh([8 \ 8 \ 8]) \approx [1 \ 1 \ 1]$$

$$o^1 = [4 \ 4]$$

$$\hat{y}^1 = [0.50.5]$$

yaw   Since $h^1 = h^2 \longrightarrow a^1 = a^2 / o^1 = o^2 / \hat{y}^1 = \hat{y}^2$