

CTC

- i) Introduce ϵ blank token, can be removed

he $\epsilon \epsilon \epsilon$ ll $\epsilon l \epsilon \epsilon o \epsilon \epsilon$ valid
 \downarrow

he l l o

(h g h) e ϵ ϵ l l l ϵ ϵ o o o invalid \rightarrow bad alignment
 ϵ ll ϵ o invalid missing he

Properties

\rightarrow Monotonic order (left to right) :- (no reordering)

\rightarrow Many To one

\rightarrow $\text{len}(Y) \leq \text{len}(X)$

Problem :-

\rightarrow A sequence X_1, X_2, \dots, X_n goes in

\rightarrow Different sequence Y_1, Y_2, \dots, Y_M comes out

generally $N \neq M$

\rightarrow Speech goes in \rightarrow words come out

English

german

eat an apple

Ich ϵ habe

einen apfel gegessen

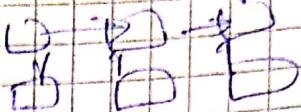
extra

order changed & extra word added

Case 1: With alignment

→ Happens in same order but are asynchronous

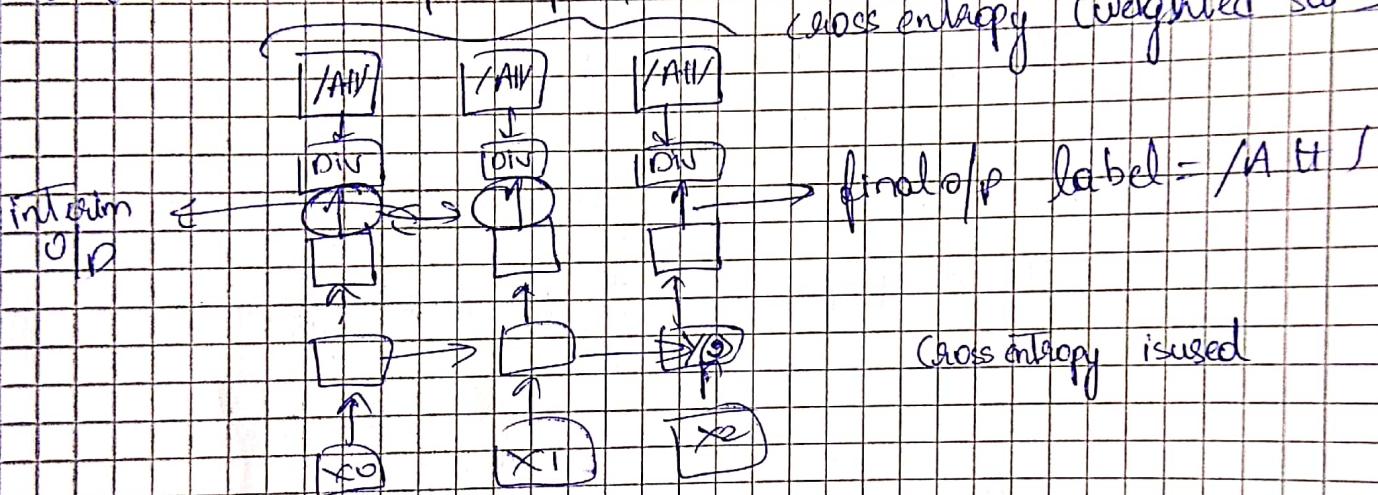
Ex: speech recognition



→ Ignores output at intermediate steps (Ex: at time 0, 1, ..., 2 etc)

→ Why can't we use same label through out each timestep even if we don't have complete input?

cross entropy (weighted sum)



* → It won't work in Question - answer scenario

What is the color of (sky)

→ here until sky comes

in color can be of any thing only after we see sky we should answer blue else we would get blue for other objects

what is the color of (tree)?

→ blue as during training

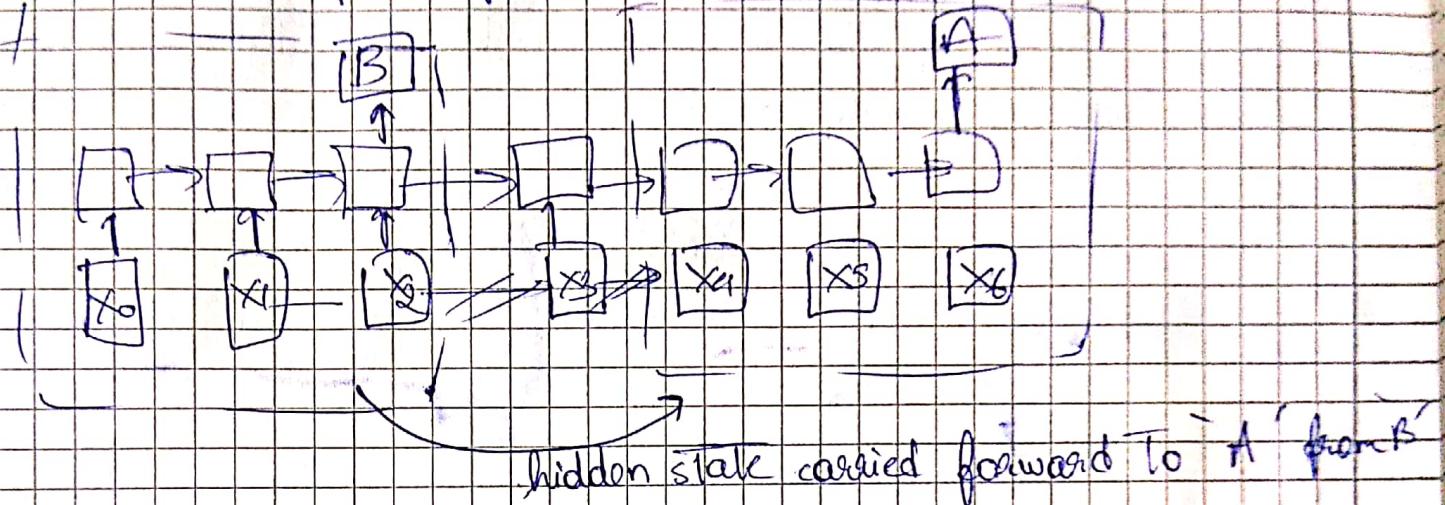
we predicted same answer all over

* → Naive Assumption: Provided we have enough training data above method works

More complex prob.

given a sequence of inputs asynchronously output sequence of symbols

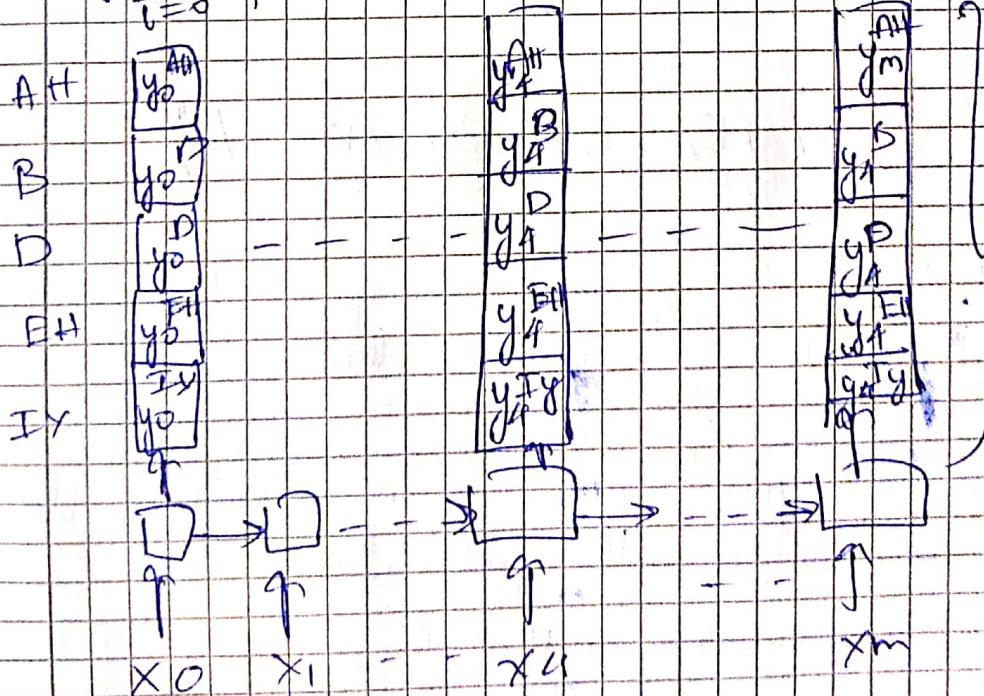
This is just a combination of many copies of simple output extend of input sequence



hidden state carried forward to 'A' from 'B'

* We also don't know when to output symbols during forward prob. Hence make r/w off at each time step even if it's blank
It's not clear when each phoneme ends as every phoneme has some prob at each time step

$$y_A^D = \text{prob}(s_A = D | x_0 \dots x_4)$$



At each step we would receive probability over each letter/phoneme

* *

Don't know 'k'

Find most likely symbol sequence given input

$$s_0 \dots s_{k-1} = \underset{s_0 \dots s_{k-1}}{\operatorname{argmax}} \text{prob}(s_0 \dots s_{k-1} | x_0 \dots x_{n-1})$$

* Option 1:

what happens if i pick only letters with most prob at each time step?

- we would end up having output sequence with same length as input (Normally $\text{len}(\text{inp}) \leq \text{len}(\text{op})$)
- Results in cyclical series diversity of op
- Example: (cannot represent) or distinguish b/w repeated characters v/s extended characters (collapses all)

* Option 2: Impose external constraints like op should be dictionary word it should have some min length. We refer this process as decoding

How To learn?

- i) Within some span of letter, we can assume same symbol repeats

$$\text{Div} = - \sum_{\text{V}} \log Y_i \quad (\text{Y}_i \text{ symbol } i)$$

(Option 1)

g g f f t y d

Actual input: g g f f t y d
Produced output: g f t y d } cannot distinguish b/w extended & repeated char

→ But we can use maximum likelihood and somehow increase the chance of getting correct output with repeated char

→ Generated sequence g f t y d is invalid as there is no word in existence check in dictionary.

a a a → a a a → a a a
split phenom repeat uniqueness maintained

7/6/20

ETC continued (Bhiksha Raj)

gradient w.r.t to the off vector \mathbf{y}_t

$$\nabla_{\mathbf{y}_t} \text{DIV} = [0 \ 0 \ -1 \ 0 \ 0]$$

$\gamma(\cdot, \text{symbol})$

zeros except at component corresponding to target

what if no timing is provided?

Soln :- guess the alignment

Initiate :- Randomly or based on some heuristic or any other

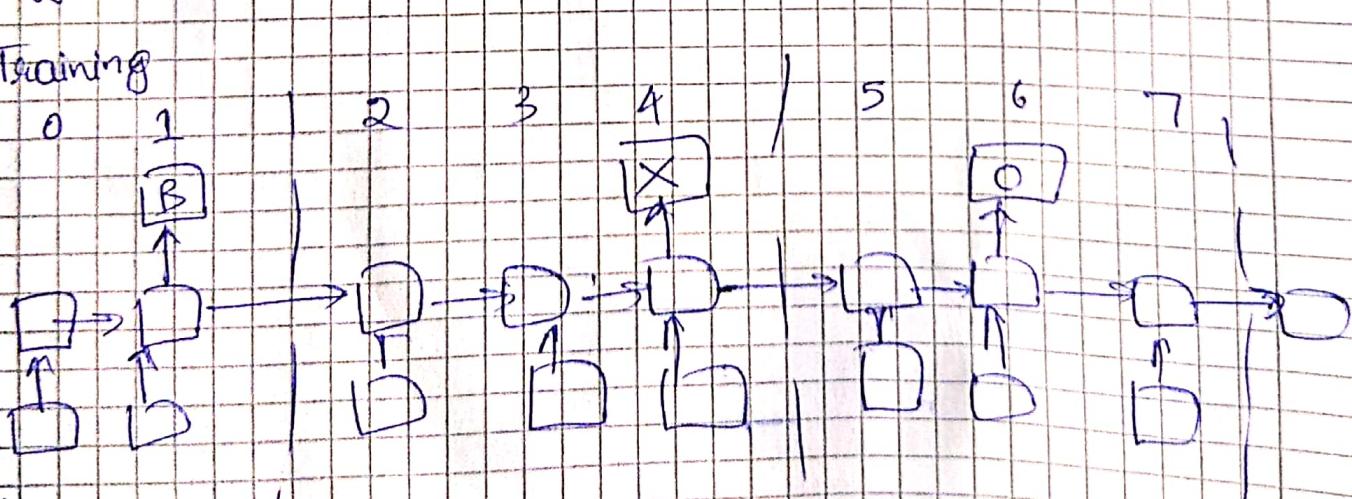
Update :- Train nfw over current alignment

- Recompute alignment for each training instance by using decoding method

Decoding :- Output from RNN obtained, This can be addressed with option 1, 2

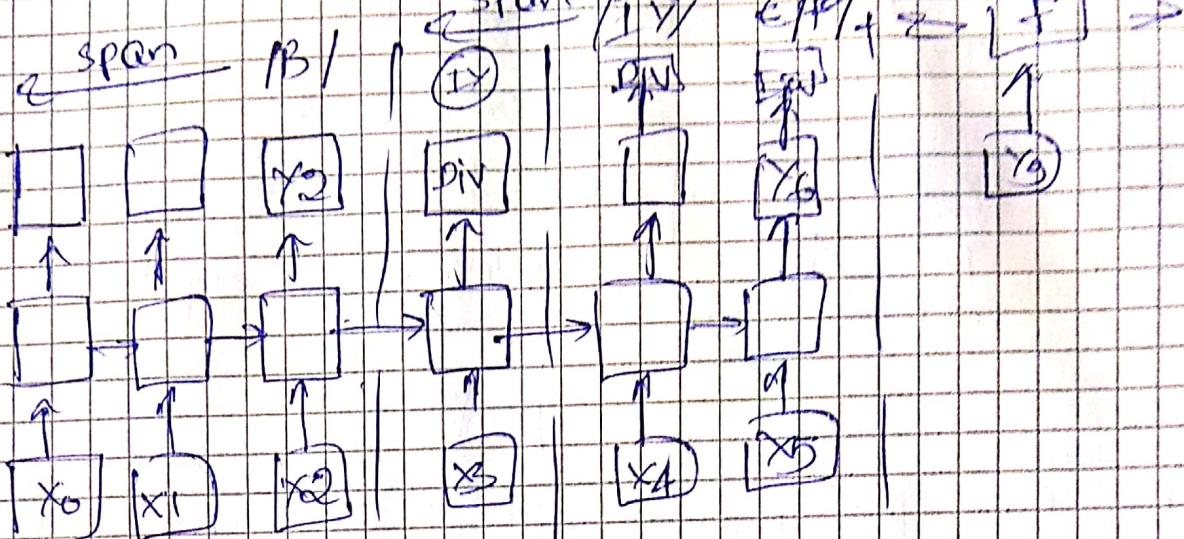
Bigger issue is training

Training



* If I know at which timestep symbol occurs $(1, 4, 6)$
we can compute crossentropy gradients only at those α locations

case - within some span same symbol appear



(case 1)

$$\text{Divergence} = \text{Xent}(Y_2, B) + \text{Xent}(Y_3, IY) + \text{Xent}(Y_4, F) + \text{Xent}(Y_5, IY)$$

at individual step

or repeat symbol during their span

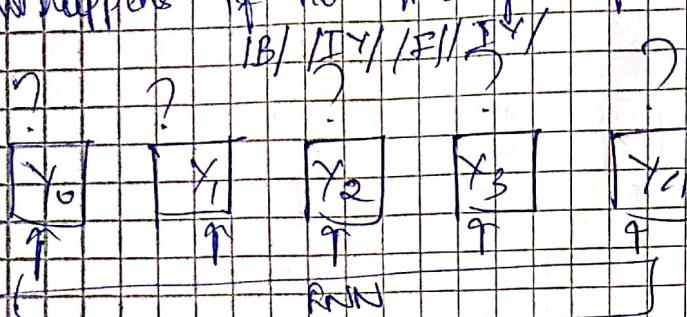
$$\text{se 2 } \nabla = \text{Xent}(Y_1, B) + \text{Xent}(Y_2, B) + \text{Xent}(Y_3, IY) + \text{Xent}(Y_4, F) + \text{Xent}(Y_5, IY) + \text{Xent}(Y_6, F) + \text{Xent}(Y_7, F) + \text{Xent}(Y_8, IY)$$

* final gradient : $\nabla_{\text{DIV}} = \begin{bmatrix} 0 & 0 & -1 & 0 & \dots & 0 \end{bmatrix}^T$

(goes at all location symbols except locations out span of Target)

partial derivative for each symbol results in being treated other sym as constant so after differentiate we get zero for other symbol

What happens if no timing info is provided?



Only we know off labels or sequence
but don't know when it occurs

How do we compute divergence?

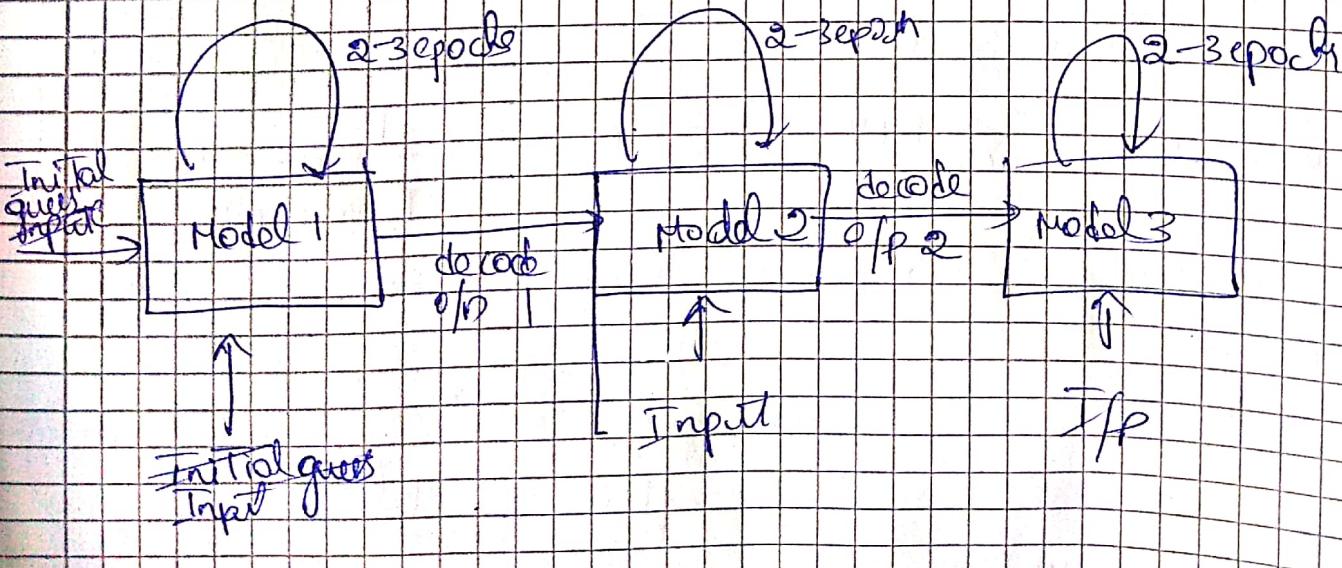
Solution 1: Guess the alignment at start

Initialization: Guess the alignment randomly / manually

Training: Pass the inputs get the off & find gradient

Learning: Use the output train network for 2-3 epochs for current alignment

Iterate: Decode Trained n/w off & provide it as input for next n/w. Again retain for 2-3 epochs & get back



Estimating alignment

- a) Given: \rightarrow Unaligned K -length symbol sequence $S = S_0 \dots S_{K-1}$ (e.g. /B/ /E/ /F/ /E/)
- \rightarrow An N -length input ($N \geq K$)
- \rightarrow Find a trained recurrent n/w

- * b) find (By separating symbols) final sequence length = Input spectrogram frame length
- \rightarrow An N -length expansion $S_0 \dots S_{N-1}$ comprising of symbols in S in
strict order e.g. $S_0 S_1 S_1 S_2 S_3 S_3 S_3 \dots S_{K-1}$
- $\rightarrow s_i \neq S_K \Rightarrow i > K$
- $\rightarrow s_i = S_K, s_j = S_L, i < j \Rightarrow K \leq L$

- c) Outcome: An alignment of Tgt symbol sequence $S_0 \dots S_{K-1}$ To
input $X_0 \dots X_{N-1} \rightarrow K \rightarrow i$
- ~~$S_0 S_1 (S_1) S_2 S_3 S_3 S_3 S_4$~~ $\dots S_{K-1}$
- $S_0 S_1 (S_2) S_3 S_4 S_5 S_6 S_7 \dots S_{N-1}$
- $i \geq K$

\rightarrow We would use viterbi algorithm

Unconstrained decoding

- \rightarrow Actual op of network at each timestep is simply probability distribution over all symbols

\rightarrow We can condition on input to get most likely symbols

\rightarrow Output may not correspond to target sequence

Problem : Symbols in unconstrained opp may not be
of symbol required

Solution : Block out all output symbols not there in target
sequence

→ Even after this we may get different sequence

Input: BIYFIY

Output without constraints :

(g) IYBF(M)

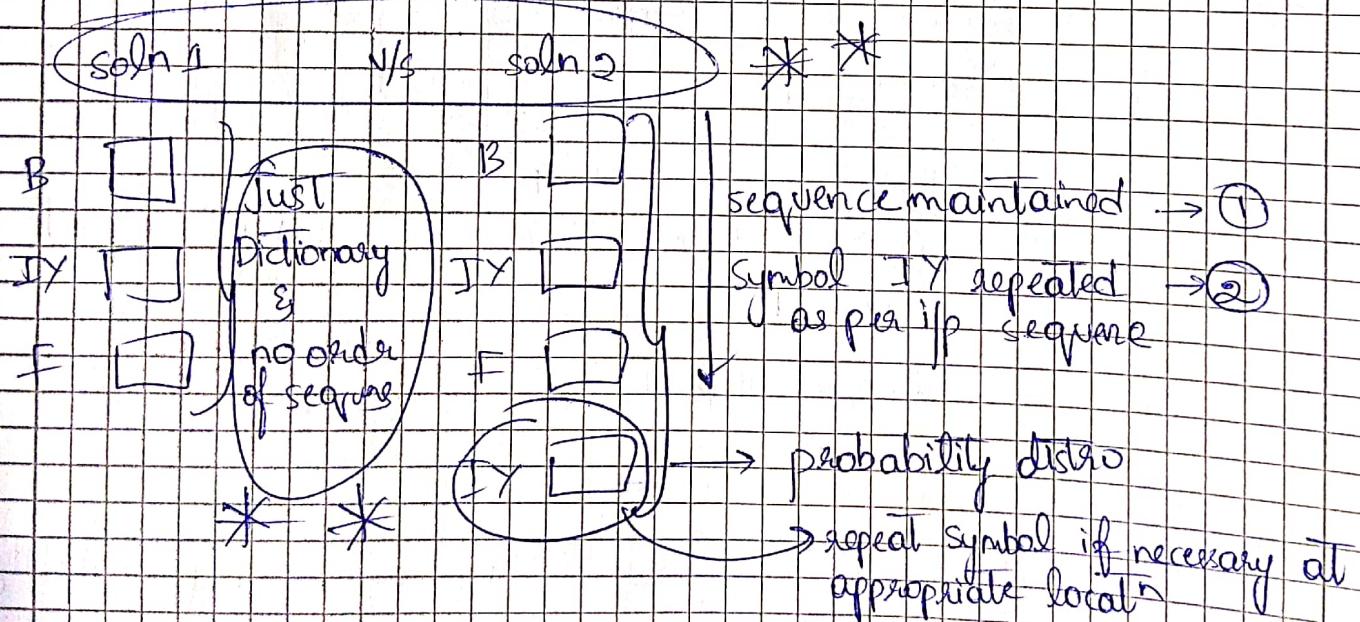
Symbols not in dictionary

Output with dictionary.
Sln 1 constraint:

B F B I Y B

Incorrect sequence

Sln 2 : Symbols should be in dictionary and also sequence must be the



Explicit constraint on alignment

- First symbol must be in top left block when decoded
- Last symbol at bottom right block during decoding

Explicit

Different paths possible

B

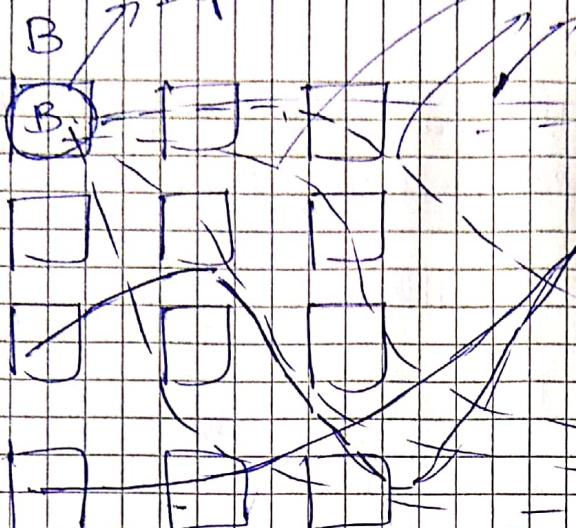
B

IY

F

IY

Explicit



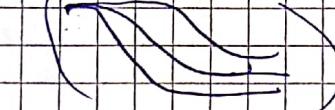
Paths non
monotonic
are not poss.

(IY)

→ explicit

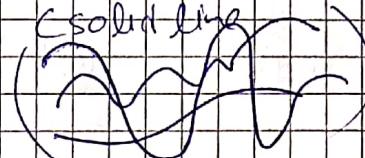
Only

Monotonic Traversal
(dotted line)



non-monotonic
(solid line)

not poss



Monotonicity assures final o/p has some order maintained even if it fails to get some symbols in between

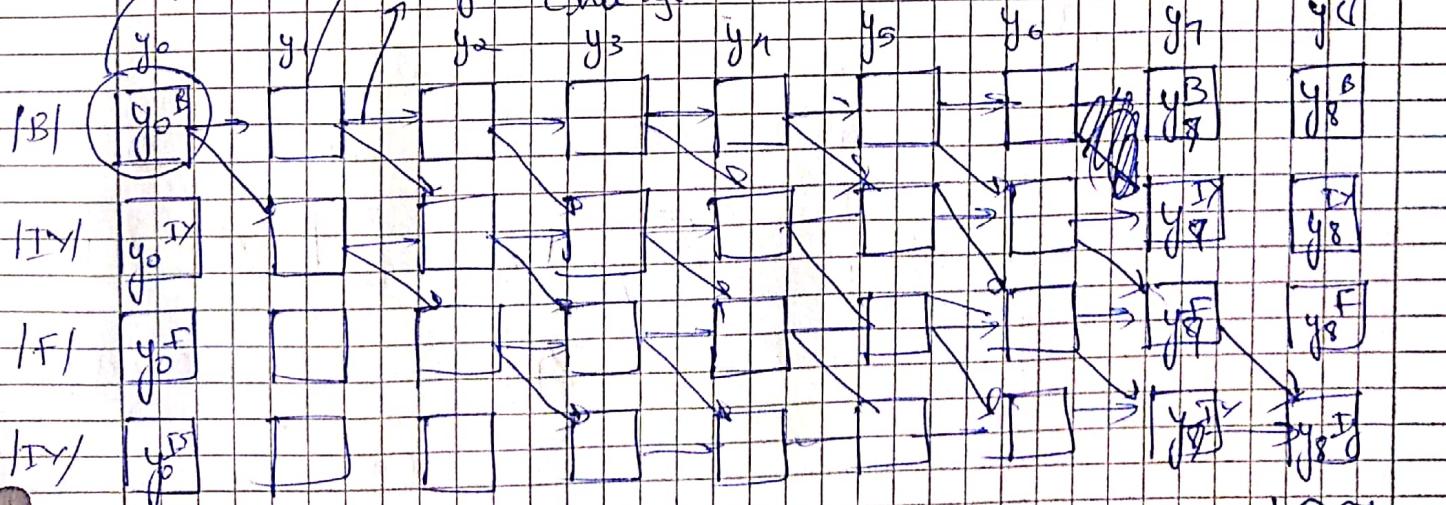
Example :-

B(B IY F IY)
B IY IY IY
B F IY IY

possible sequences

constant

prob = 1
vertices
edges
change

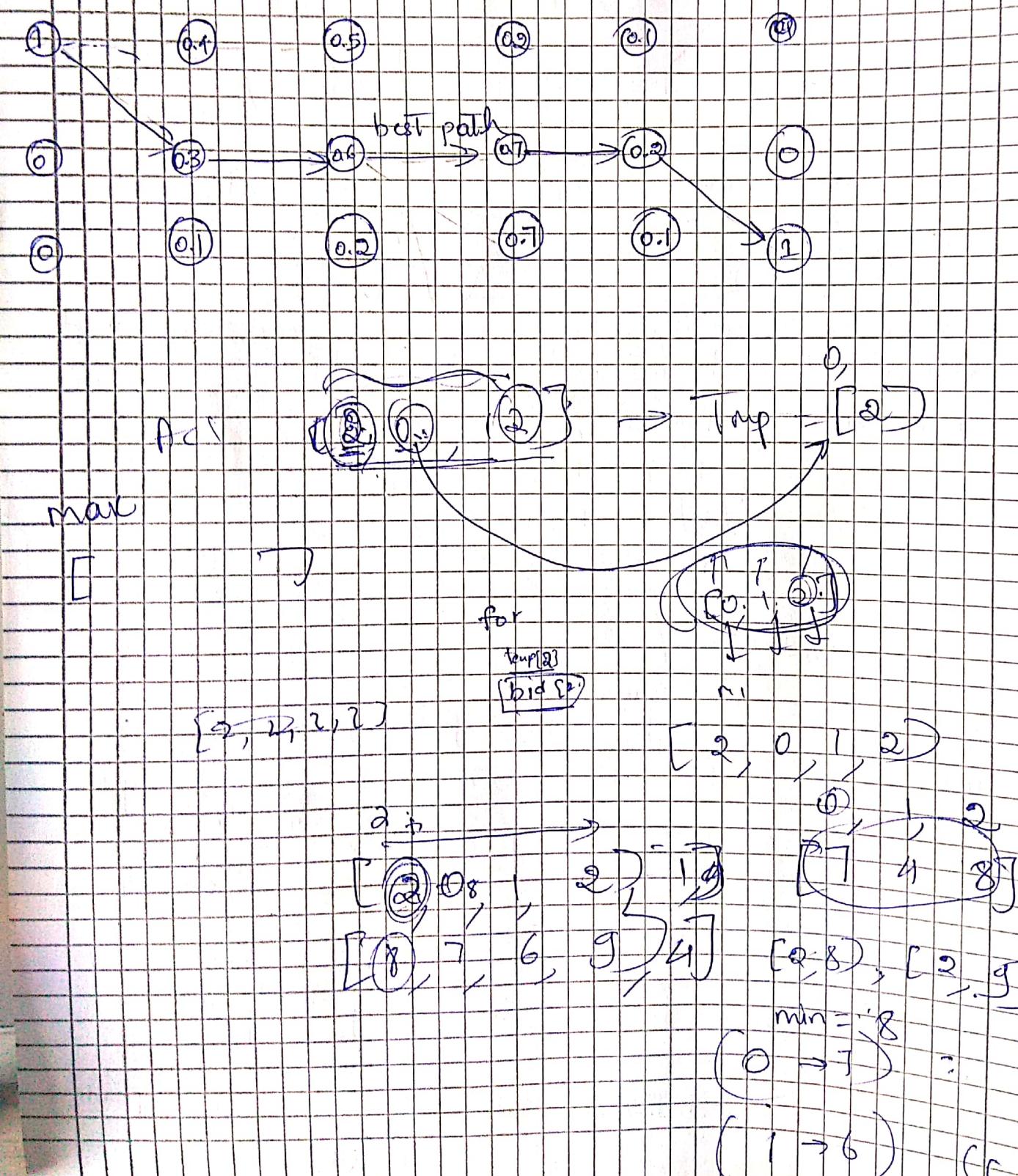


Vertices contain probability, edges value hard coded

prob = 1

To ①

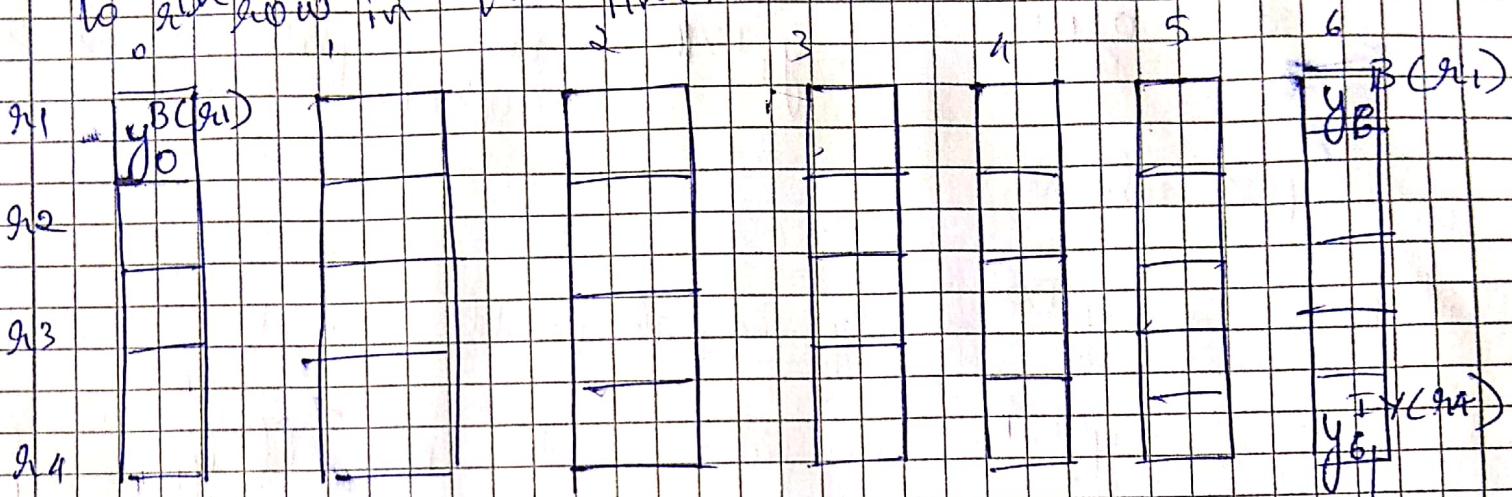
- Most probable sequence is the path having max score where score is product of probabilities from $B \rightarrow Y$
- * Use dynamic program Viterbi algorithm



Viterbi Algorithmen

Notation : Let $y_v^{(sc)}$ is prob of tgt symbol assigned

To catch snow, in T^{th} time.



Initialization :

1st Time Step

No predecessor $\Rightarrow \text{BPC}(0, i) = \text{null}$, $i = 0 - k - 1$

$$B_{SCA}(0,0) = y_0^{(s_0)}, B_{SCA}(0,i) = -\infty, i=1, \dots, k-1$$

for $T=1 \dots T-1$

$$BP(t, 0) = 0 ; BSCN(t, 0) = BSCN(t-1, 0) * y_t^{(s, 0)}$$

for $l=1 \dots k-1$

$$\rightarrow \text{BP}(T, l) = (\text{if } (\text{BSGR}(T-1, l-1) > \text{BSGR}(T-1, l)) \\ l-1, \text{else } l)$$

$$B_{S2}(T, l) = B_{S2} \left(BP(T, l) \right) * y_T^{(cl)}$$

Timestep 0 :- No predecessor
score = y_0^B

Timestep (1) :-

rows. { 0, 1y }

$y_1^B \text{ or } y_1^{IY}$ \rightarrow Only one path possible hence update

Compare whether

Timestep (2) :- $y_2^B \rightarrow$ only one path possible use it

$y_2^{IY} \rightarrow$ Two paths possible now if $y_1^{IY} > y_1^B$
choose y_1^{IY} as predecessor

$y_2^F \rightarrow$ Only one path y_1^{IY}

T3: $y_3^{B0} \rightarrow y_2^{B0}$ \rightarrow predecessors

$y_3^{IY} \rightarrow$ $y_2^B * y_3^{IY}$ if $y_2^B > y_2^{IY}$

$y_3^{IY} \rightarrow y_2^{IY} * y_3^{IY}$ else

\rightarrow Predecessor

$y_3^F \rightarrow$ $y_2^{IY} * y_3^F$ if $y_2^{IY} > y_3^F$
else

$y_3^{IY} \rightarrow y_2^F$

Greedy v/s Beam

- Greedy can pick only one symbol at each time step, helpful if one/two symbols have much higher prob than others
 - Beam can help with generating multiple alignments for a single o/p
- B E I Y C F I Y
 B B E I Y F I Y
 B C E C I Y F I Y
- But not guaranteed to find most likely Y

$$\text{Beam search} \quad \frac{1}{T_y^{\alpha}} \sum_{t=1}^{T_y} \log P(Y^{(t)} | X, Y^{(t-1)})$$

$\alpha = 0.7$

6/06/20

Stober - CTC

$$P(Y|X) = \text{Sum of all valid alignments} \rightarrow \text{dynamic programming}$$

$$\text{RNN output} = \prod_{t=1}^T p(a_t|x) \rightarrow \text{loss each alignment probability}$$

→ should be differentiable

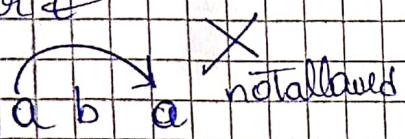
- Since we don't know exact time position where each character occurs, we need to aggregate the each individual valid alignment & treat them all as one

$$\begin{aligned} & B \ I \ Y \ C \ F \ C \ Y \\ & B \ E \ I \ Y \ F \ E \ Y \\ & B \ C \ I \ Y \ C \ F \ C \ Y \\ & B \ B \ I \ Y \ C \ F \ F \ E \ Y \ Y \end{aligned} \quad \left. \right\} \text{all are valid}$$

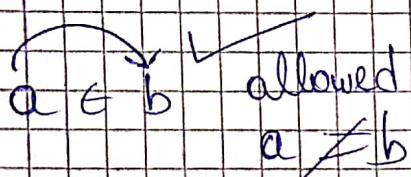
Allowed Transitions v/s not allowed transition

Introduce ϵ To differentiate b/w Two repeated characters

Collapse every thing /symbol repeated To unique symbol before & after ϵ



$a \xrightarrow{\epsilon} a$ not allowed if
a is not repeating by ϵ
belong to single



greedy

→ Doesn't consider multiple alignment

* MAP v/s MLE hypothesis

sum of prob belonging
To one class should be
great

greedy

If we consider greedy Then most
probable alignment is ~~bbb~~
with prob of 0.4 It's like
MLE hypothesis

Example	Prob
a a c	greedy
a a a	batch
b a b	0.05
a a c	0.3
a c a	0.05
b b b	0.4
c a a	0.1

$$\text{But realistically } P(\text{aab} | \text{aa} \epsilon \text{ka}) = 0.1 + 0.05 + 0.3 \\ = 0.48$$

hence aab is required sequence
which is like MAP

Beam Search Heuristic

Applying Language Model

- Just beam search / greedy results in ~~one~~ output at each time step being generated conditionally independent of its previous timestep
- In realworld This is different, Use some language prior while calculating combined prob

$$Y^* = \underset{y}{\operatorname{argmax}} P(Y|X) \cdot \underbrace{P(Y)}_{\text{CTC conditional prob}}^\alpha \cdot \underbrace{L(Y)}_{\substack{\text{language model prob} \\ \text{word insertion bonus}}}^\beta$$

* → Normally used during deployment / decoding & not during RNN Training

** → $P(Y)$ → can be trained independently on Wikipedia Text To know general structure of language (Domain law/medical)

*** → $L(Y)$ → Helps to remove bias on short Target sequence

$$\log \operatorname{argmax} P(Y|X) + \alpha \log P(Y) + \beta \log L(Y)$$

→ $P(Y)$ → can be swapped b/w domains

HMM

Monotonicity → cannot go back, skip over blank state