

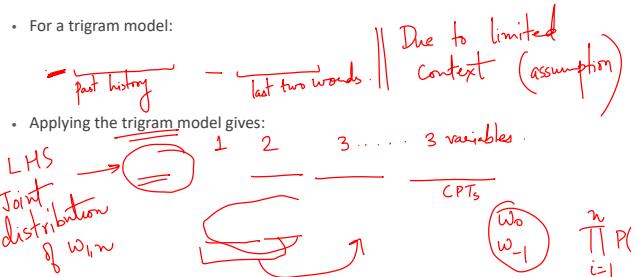
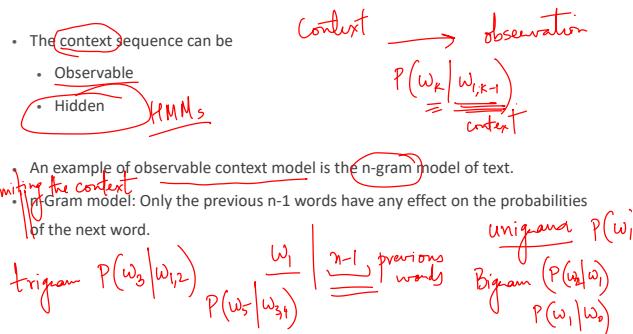
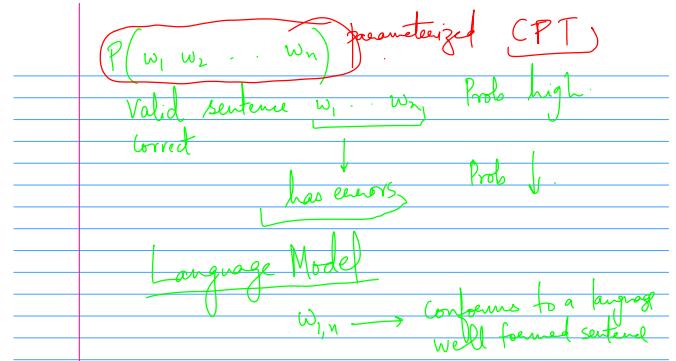
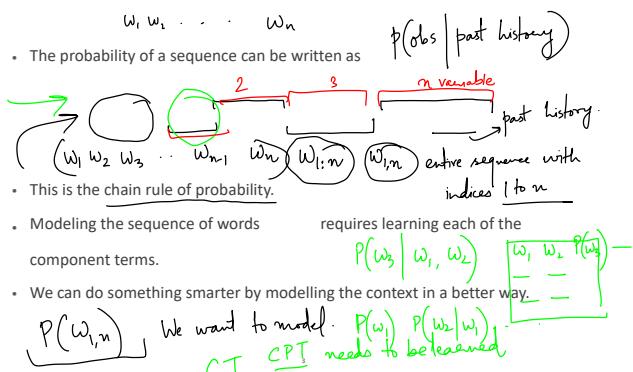
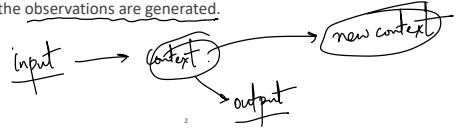
06 Jan 2023.

Hidden Markov Models

Markov Model

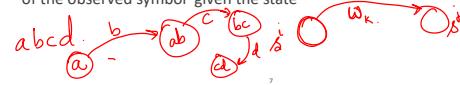
- Let's say there is a process which generates a sequence:
- Sequence of written words text
- Sequence of spoken words (speech)
- DNA, RNA protein's amino acid sequence, etc.

- The observed sequence is governed by some context which itself keeps updating as the observations are generated.



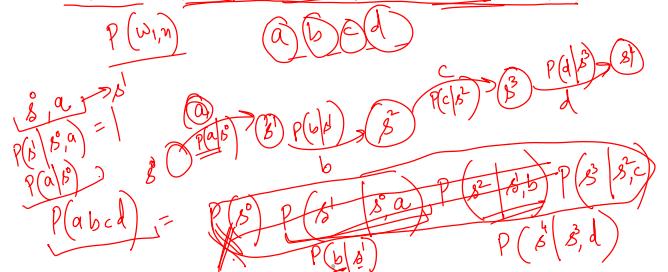
- To create such a trigram model, we require the probability of every possible valid trigram. $P(w_1, w_2, w_3)$
 - The estimate of the conditional probability $P(w_3 | w_1, w_2)$
- Training dataset*
- e. estimate* \rightarrow *Count* \rightarrow *Corpus*
- "To create such a model we simply go through some training text"
- $$P_e(\text{model} | \text{'such', 'a'}) = \frac{C(\text{'such', 'a', 'model'})}{C(\text{'such', 'a'})} = 1$$

- After all the hard work, we now have a model which can generate the probability of any string. $P(w_1, n)$
- This model is called a Markov chain. *context is observable*.
- It can be represented as a graph of nodes and arcs.
- Nodes are the states. *context*
- Arches represent transitions from one state to another.
- An arc is labelled with the observation (emitted symbol) and the probability of the observed symbol given the state.



- Why is the state obvious here?

- Given a Markov chain, the probability of an observed sequence can be computed as the product of the probabilities of each transition in the path.



What if the states are hidden?

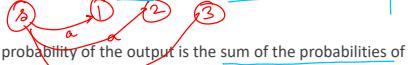
- Now consider the scenario when the states are non-obvious.
- That means, we are uncertain about the state sequence that generated the given observation sequence.
- That means,
 - there can be multiple state sequences that can lead to a given observation sequence.
 - a given state sequence can generate multiple observation sequences
 - for a given state, we can make transition to multiple possible states, even while generating the same output.

Hidden Markov Models

- So now, transition to multiple states is possible from a given state, even while emitting the same observation (symbol). *Difference from a Markov chain*.
- HMMs are a generalisation of Markov Chains in which a given state may have several transitions out of it, all with the same label (symbol).
- Recall that in a Markov chain, given a state, the next state is certain (fixed) once you observe a given symbol.



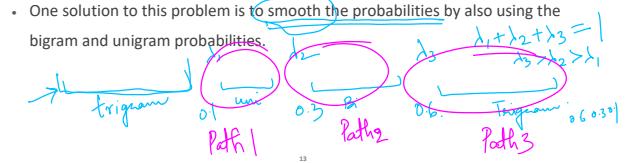
- Since more than one successor state may have the same output, in general there may be several paths through an HMM that produce the same output sequence.
- In such cases, the probability of the output is the sum of the probabilities of all possible paths.
- The Markov chain cannot be used now.
 - What we need is a state transition diagram to represent an HMM.



12

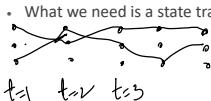
Moving from a Markov Chain to an HMM

- Consider that we are learning a trigram model.
- If the training corpus is sparse, it may happen that there is a trigram in the next text that never appeared in the training corpus.
- In that case, a 0 probability will be assigned to the trigram, i.e. probability of the third word, given the first two words.



13

- Since more than one successor state may have the same output, in general there may be several paths through an HMM that produce the same output.
- In such cases, the probability of the output is the sum of the probabilities of all possible paths.
- The Markov chain cannot be used now.

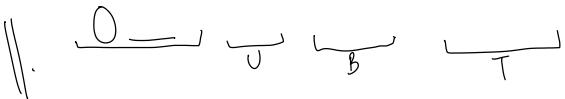


12

$$P(\text{obs}) = \sum_{\text{paths}} \text{prob}(\text{path})$$

Moving from a Markov Chain to an HMM Hidden state

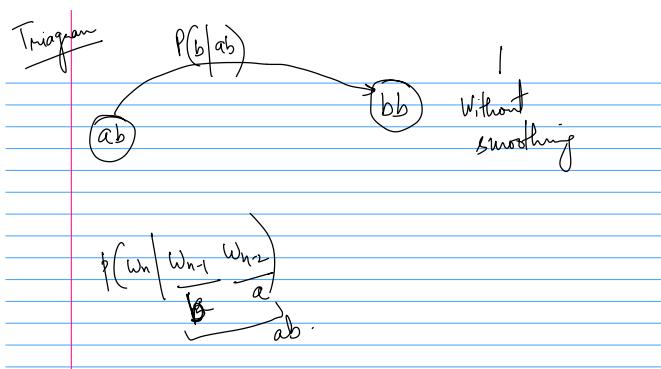
- Consider that we are learning a trigram model.
- If the training corpus is sparse, it may happen that there is a trigram in the next text that never appeared in the training corpus.
- In that case, a 0 probability will be assigned to the trigram, i.e. probability of the third word, given the first two words.
- One solution to this problem is to smooth the probabilities by also using the bigram and unigram probabilities.

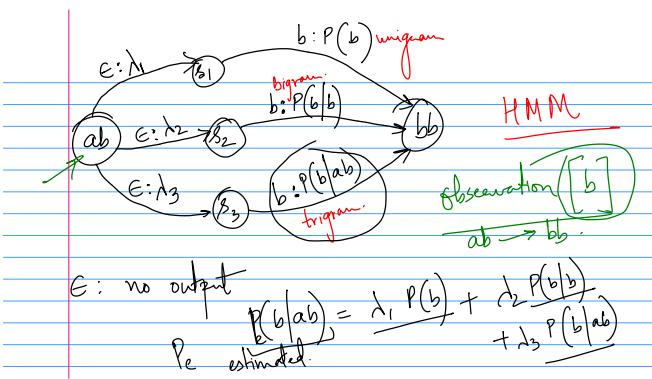


13

- If we have missing trigram then the bigram and unigram probabilities take over.
- If we have missing trigram and bigram, then we fall back on the unigram probability.
- The probability of a transition has become an addition of multiple components (arcs) between the two states.
- An HMM will allow multiple paths between two states.

14



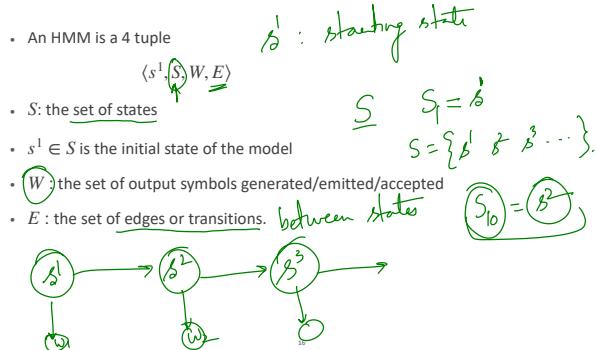


HMM state transition diagram for a trigram model

- Consider that we have seen the symbol ab.
- The next symbol in the sequence is a.
- So the next state is ba
- The trigram $P(a|ab)$ itself can be represented as an HMM

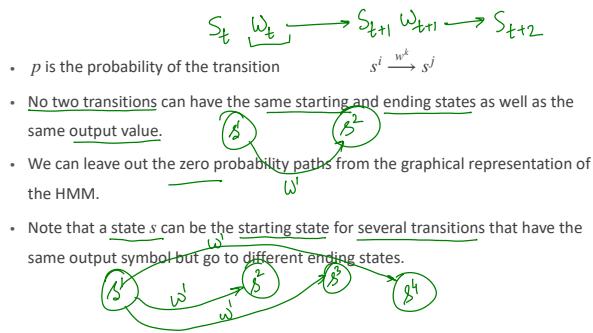
$$\text{Smoothed Trigram } P(a|ab) = \lambda_1 P_e(a) + \lambda_2 P_e(a|b) + \lambda_3 P_e(a|ab) \quad \text{from the corpus}$$

15

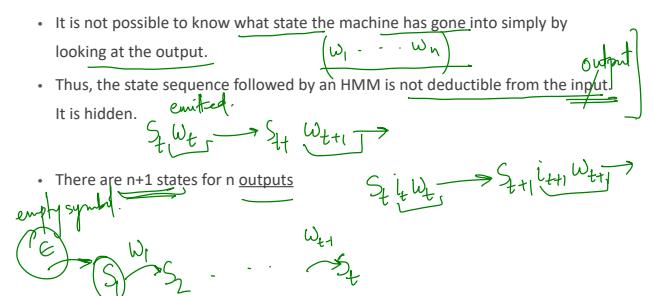


- For each of the sets S, W, E we assume a canonical ordering of the elements
- $S = \langle s^1, s^2, \dots, s^n \rangle$
- $W = \langle w^1, w^2, \dots, w^m \rangle$
- $E = \langle e^1, e^2, \dots, e^r \rangle$
- The starting state of the HMM is the first element in the ordering of states.
- A transition is a 4-tuple $\langle s^i, s^j, w^k, p \rangle$ where s^i is the state from which the transition starts, s^j is the state where the transition ends, w^k is the output symbol generated by the model, p is the probability of taking the transition

17



18



19

- The probability p associated with a transition $s^i \xrightarrow{w_k} s^j$ is $P(s^i \xrightarrow{w_k} s^j) = P(S_{t+1} = s^j | S_t = s^i)$

$$= P(s^j, w^k | s^i)$$

- When writing $P(s^j, w^k | s^i)$, it is understood that s^i is the prior state and s^j is the next state.

$$P(s^i \xrightarrow{w_k} s^j) \quad \leftarrow \text{unknown}$$

$P(s^j | s^i) \quad \leftarrow \text{Learned during training}$

20

- Markov models assume that the only information affecting the probability of an output, or of the next state, is the prior state.

- As per this Markov Assumption

$$P(w_n, s_{n+1} | s_1, \dots, s_{n-1}, s_n) = P(w_n, s_{n+1} | s_n) = P(s^i \xrightarrow{w^k} s^j)$$

- The probability of a sequence $w_{1,n}$ is the probability of all possible paths through the HMM that can produce this sequence.

In other words:

$$P(w_{1,n}) = \sum_{s_{1,n+1}} P(w_{1,n}, s_{1,n+1})$$

$\uparrow \text{hidden/latent}$

$\uparrow \text{observed}$

21

The Markov Assumption helps

$$P(w_{1,n}) = \sum_{s_{1,n+1}} P(w_{1,n}, s_{1,n+1})$$

$\nearrow \text{Joint}$

$\nearrow \text{Summing Out past obs/states}$

$\nearrow \text{Complex}$

$P(w_{1,n}) = \sum_{s_{1,n+1}} P(s_1) P(w_1, s_2 | s_1) P(w_2, s_3 | w_1, s_1) \dots P(w_n, s_{n+1} | w_{1,n-1}, s_{1,n})$

$\nearrow \text{Factorization}$

$P(w_{1,n}) = \sum_{s_{1,n+1}} \prod_{i=1}^n P(w_i, s_{i+1} | s_i)$

$\nearrow \text{past is captured by the current state!}$

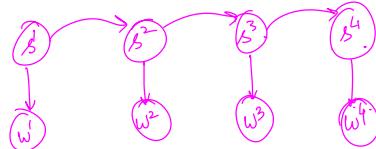
$\nearrow \text{expedited}$

$\nearrow \text{binary}$

$\nearrow \text{Unary terms!}$

22

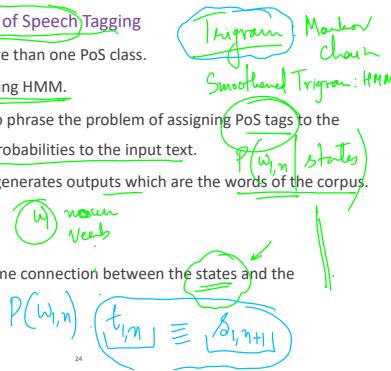
Exercise: Graphically visualize the simplification offered by Markov Assumption



23

An application of HMM: Part of Speech Tagging

- English words can be in more than one PoS class.
- PoS tags can be assigned using HMM.
- But to use HMM we need to phrase the problem of assigning PoS tags to the words as one of assigning probabilities to the input text.
- We assume that the HMM generates outputs which are the words of the corpus.
- We assume that there is some connection between the states and the transitions of the tags.



24

- Earlier we had for Language Model (LM)

$$P(w_{1,n}) = \sum_{s_{1,n+1}} P(w_{1,n}, s_{1,n+1})$$

- For PoS tagging, we assume correspondence between states and the tags $\{\text{sequence}\}$

- So,

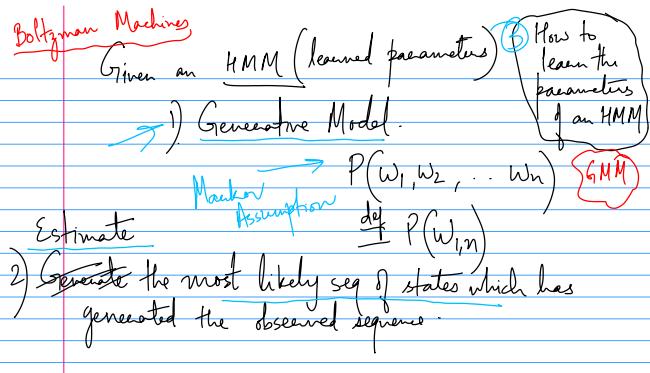
$$P(w_{1,n}) = \sum_{t_{1,n+1}} P(w_{1,n}, t_{1,n+1})$$

- Here, $t_{1,n+1}$ is a sequence of $n+1$ parts of speech or tags.

25

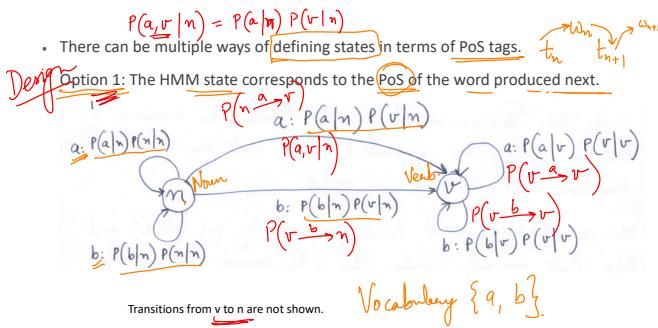
$$\beta^i = (w_{i-1} \ w_{i+1})$$

$$\beta^i = (w_{i-1} \ w_{i+1})$$



- The last tag t_{n+1} is the tag achieved after emitting w_n
 - So t_{n+1} is a tag which corresponds to the non-existent word w_{n+1}
 - We define the problem of PoS tagging as finding
- $$\arg \max_{t_{1:n}} P(t_{1:n} | w_{1:n}) = \arg \max_{t_{1:n}} \frac{P(w_{1:n}, t_{1:n})}{P(w_{1:n})} = \arg \max_{t_{1:n}} P(w_{1:n}, t_{1:n})$$
- Most likely tag sequence \Rightarrow State sequence

26



27

- The joint probability of all the words is

$$P(\text{obs seq}) = \sum_{t_{1:n+1}} P(w_n | t_{1:n+1})$$

- We use Markov assumptions to simplify the joint distribution over the words

- The probability of a word appearing at a particular position given that its PoS occurs at that position is independent of everything else

$$P(w_n | w_{1:n-1}, t_{1:n}) = P(w_n | t_n) \quad \text{Assumption}$$

- The probability of a PoS coming next is dependent only on the previous PoS

$$P(t_n | w_{1:n-1}, t_{1:n-1}) = P(t_n | t_{n-1})$$

28

$$P(w_{1:n}) = \sum_{t_{1:n+1}} P(w_{1:n}, t_{1:n+1})$$

Big team

Markov Assumption

$$= \sum_{t_{1:n+1}} \prod_{i=1}^n P(w_i | t_i) P(t_{i+1} | t_i)$$

- Several state paths can lead to the observed word sequence.
- We can also find the most likely sequence of states that can generate a given word sequence as.

$\arg \max_{t_{1:n+1}} \prod_{i=1}^n P(w_i | t_i) P(t_{i+1} | t_i)$ Tags (as per derrf option 1)

29

Design Option 2

A state can be formulated as representing two parts of speech, that of previous and the current word

$$P(w_{1:n}) = \sum_{t_{1:n+1}} \prod_{i=1}^n P(w_i | t_i) P(t_{i+1} | t_i, t_{i-1})$$

↓ ↓

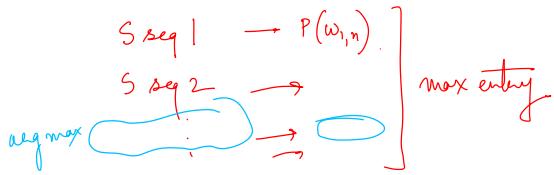
State
 $t_i = (t_i, t_{i-1})$

$t_i t_{i-1} = \text{a state}$

30

Get the POS tags -

Finding the most likely sequence of states



Finding the most likely sequence of states: The Viterbi Algorithm

- We need a way to find the most likely path through an HMM given a particular input $w_{1,n}$

$$s^i = t$$

- The most likely path corresponds to the most likely tags for the sentence.

- For a string of length n , the Viterbi algorithm finds $s^{(n+1)}$, which is defined as

$$\begin{aligned} \max_{s_{1,t}} \text{prob path from } s_1 \text{ to } s_t \\ &= \arg \max_{s_{1,t}} P(s_{1,t} | w_{1,t-1}) \\ &= \arg \max_{s_{1,t}} \frac{P(s_{1,t}, w_{1,t-1})}{P(w_{1,t-1})} \quad (\text{up to time } t). \end{aligned}$$

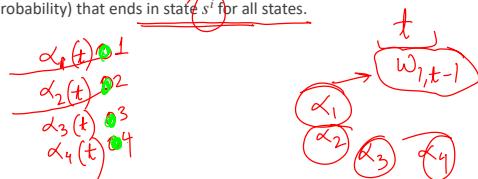
32

The Viterbi Algorithm

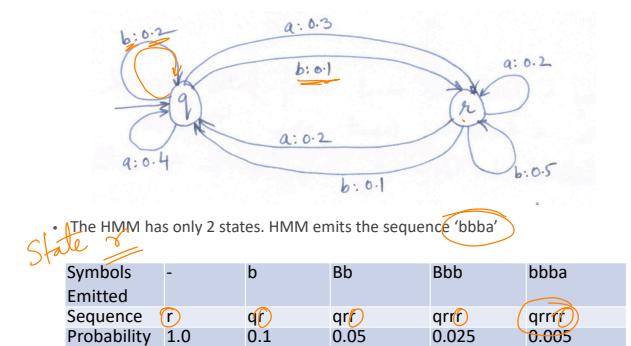
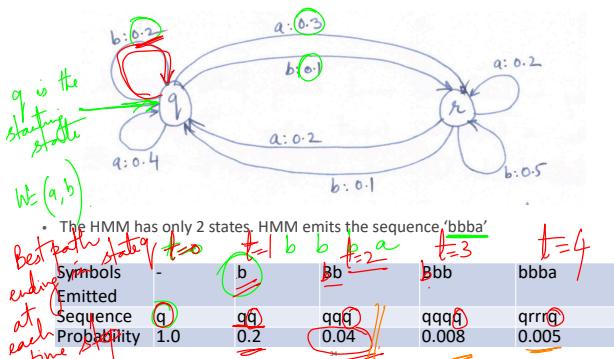
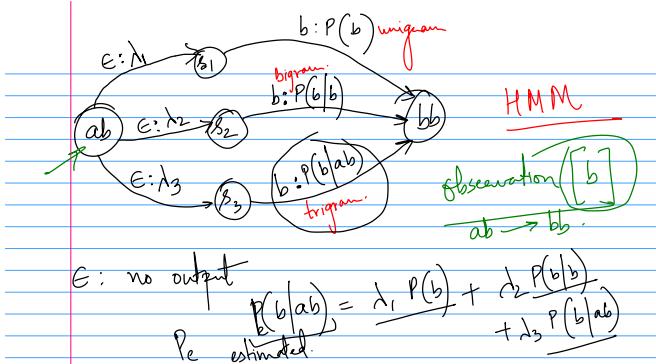
- The Viterbi algorithm computes the most likely state sequence starting with the empty output sequence, working up to the answer one output at a time.

time step

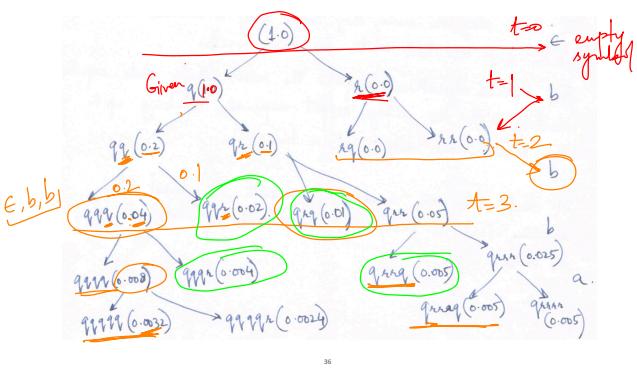
- At each stage one need only compute the most likely sequence (and its probability) that ends in state s^t for all states.



33



35



36

The Viterbi Algorithm

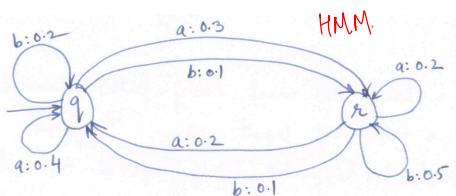
- The Viterbi algorithm computes the most likely state sequence starting with the empty output sequence, working up to the answer one output at a time.

- At each stage one need only compute the most likely sequence (and its probability) that ends in state s^t for all states.

$$\text{Most likely path} = \arg \max_{S_{1:t+1}} P(w_{1:t}, S_{1:t+1})$$

(1) Best path to each state at each step
 (2) Overall best path
 which has generated the observation

34



- The HMM has only 2 states. HMM emits the sequence 'bbba'

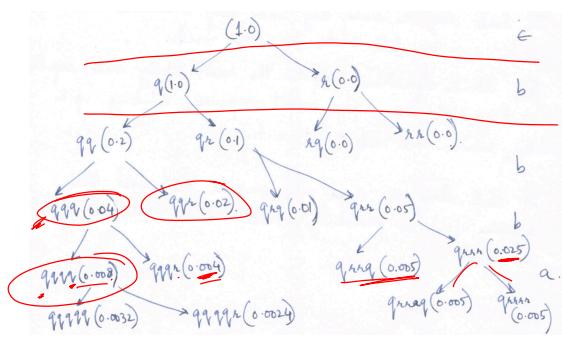
	$t=1$	$t=2$	$t=3$	$t=4$	$t=5$
Symbols Emitted	-	b	Bb	Bbb	bbba
Sequence1	q	qq	qqq	qqqq	qrrr

35

- The HMM has only 2 states. HMM emits the sequence 'bbba'

	-	b	Bb	Bbb	bbba
Sequence2	r	qr	qrr	qrrr	qrrrr
Probability	1.0	0.1	0.05	0.025	0.005

36



37

$$P(S_{t-1} = s^i \xrightarrow{w_t} S_t = s^j)$$

- For each time tick it is necessary to store the best path and its probability for each possible state.

- Storing only the overall best path is not sufficient.

$$S_{t-1} = s^i$$

$$S_t = s^j$$

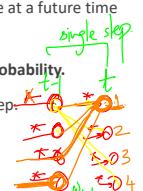
- This allows us to discover paths that become more probable at a future time step.

- So we need to store the best path to each state and the probability.

- We never need to store more than $\sigma = |S|$ paths at any step.

$$P(s^i \xrightarrow{w_t} s^j)$$

38



Step by Step constructing the best forward path from time t to get the overall best path

- Let $\alpha_i(t)$ be the state sequence that has maximum probability of generating

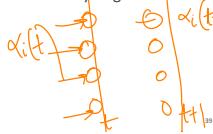
$w_{1,t-1}$ and ends up in state s^i at time t . $s^i = s^t$

Path $\alpha_i(t)$ is the best partial path.

$$\alpha_i(t) = \arg \max P(w_{1,t-1}, s_{1,t-1}, s_t = s^i)$$

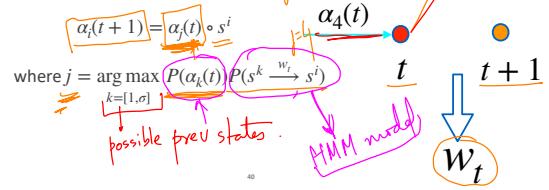
fixed for $\alpha_i(t)$

- Given the best way to get to all states at time t , $\alpha_k(t)$ for $k = 1$ to σ , we can compute the best way to get to a state $t+1$ as a simple extension.



Step by Step constructing the best forward path from time t to get the overall best path

- The best path leading to state i at time $t+1$



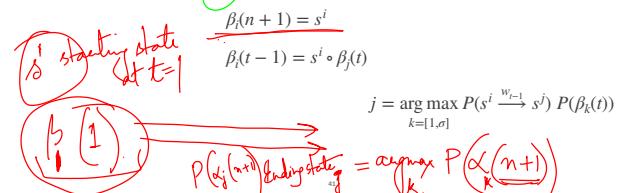
40

41

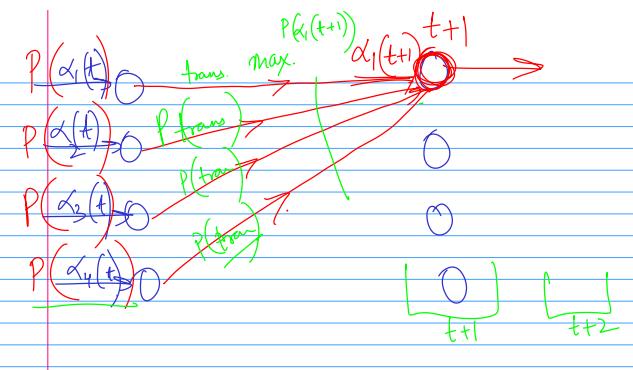
Step by Step constructing the best backward path from time t

- It is also possible to use the same algorithm to work backward from the end of the output state sequence to the beginning.

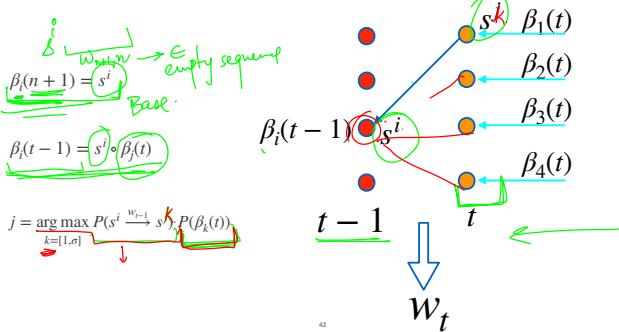
- Let $\beta_i(t)$ be a state sequence with the maximum probability of generating $w_{t,n}$ while starting in state s^i



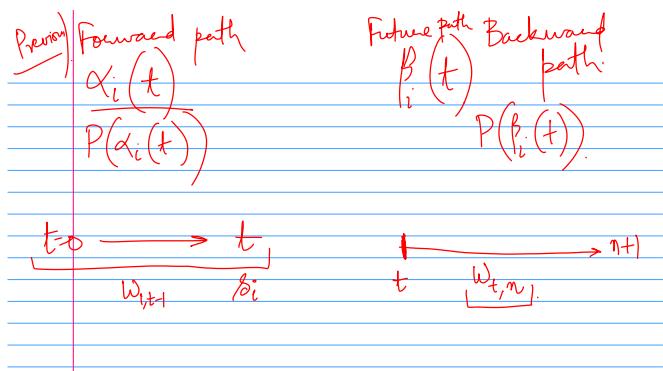
42



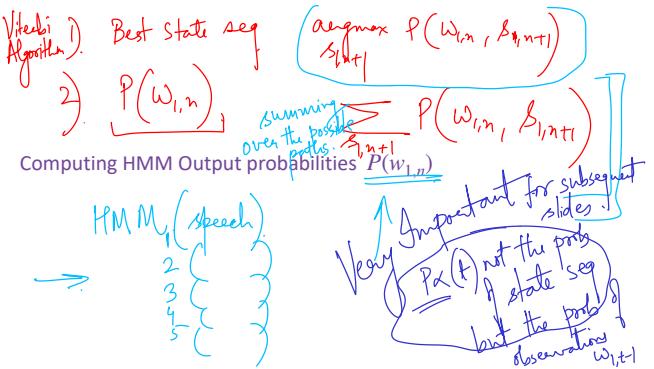
Step by Step constructing the best backward path from time t



42



- We have seen how the Viterbi algorithm is able to construct the best path (state sequence) using forward recursion or backward recursion.



43

Computing HMM output probabilities using forward recursion

- There are two algorithms for computing the probability of any particular HMM output. The algorithms are mirror images of each other.

- Let $P\alpha_i(t)$ be the probability of producing $w_{1:t-1}$ while ending up in state s^i .
 i.e. $P\alpha_i(t) = P(w_{1:t-1}, S_t = s^i)$ for $t > 1$
- For $t = 1$ we can think of $w_{1,0}$ as the empty string
 Given that we start in state s^1 part of the learned HMM model.

$$P\alpha_i(1) = \begin{cases} 1.0 & \text{if } i = 1 \\ 0 & \text{otherwise} \end{cases}$$

45

- The quantity $P\alpha_i(t)$ is called the forward probability.

- That is, the probability of generating the word sequence up to time t and ending in state i

- We can compute $P\alpha_i(t) \forall t$ in linear time by working forward through the entire output sequence, word by word.

$P\alpha_i(t+1) = P(w_{1:t}, S_{t+1} = s^j)$

$$\begin{array}{|c|c|} \hline i=1 & 0 \\ \hline i=2 & 0 \\ \hline \vdots & \vdots \\ \hline i=3 & 0 \\ \hline P\alpha_i(t) & t \\ \hline \end{array}$$

$$\begin{array}{|c|c|} \hline i=1 & 0 \\ \hline i=2 & 0 \\ \hline \vdots & \vdots \\ \hline i=3 & 0 \\ \hline P\alpha_i(t+1) & t+1 \\ \hline \end{array}$$

46

$$\begin{aligned} P\alpha_j(t+1) &= P(w_{1:t}, S_{t+1} = s^j) \\ &= \sum_{i=1}^{\sigma} P(w_{1:t}, S_t = i, S_{t+1} = s^j) \\ &\stackrel{\text{Bayes Theorem}}{=} \sum_{i=1}^{\sigma} P(w_{1:t}, S_{t+1} = s^j | w_{1:t-1}, S_t = s^i) P(w_{1:t-1}, S_t = s^i) \\ &\stackrel{\text{Apply Markov Assumption.}}{=} \sum_{i=1}^{\sigma} P(w_{1:t}, S_{t+1} = s^j | S_t = s^i) P(w_{1:t-1}, S_t = s^i) \\ &\quad \text{transition prob. } P(s^i \xrightarrow{w_t} s^j) \end{aligned}$$

47

$$P\alpha_j(t+1) = \sum_{i=1}^{\sigma} P(s^i \xrightarrow{w_t} s^j) P\alpha_i(t)$$

Forward Recursion

If we compute all the $P\alpha_i(n+1)$'s then it is easy to compute $P(w_{1:n})$

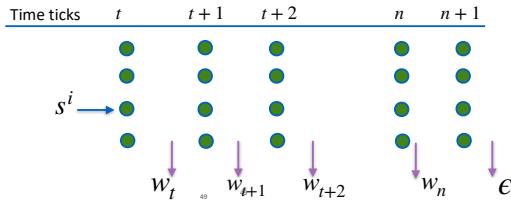
$$\begin{aligned} P(w_{1:n}) &= \sum_{i=1}^{\sigma} P(w_{1:n}, S_{n+1} = s^i) \\ &= \sum_{i=1}^{\sigma} P\alpha_i(n+1) \end{aligned}$$

$$\begin{array}{|c|c|} \hline 0 & n+1 \\ \hline 0 & \\ \hline 0 & \\ \hline 0 & \\ \hline \end{array}$$

48

Computing $P(w_{1,n})$ using a Mirror Algorithm (Backward recursion)

- To calculate the probability of a sequence, we work backward from the end.
- Backward probability $\underline{P\beta}_i(t)$ is the probability of seeing the sequence $w_{t,n}$ if the state of the HMM at time t is s^i



$$\underline{P\beta}_i(t) = P(w_{t,n} | S_t = s^i)$$

$$\begin{aligned}\underline{P\beta}_1(1) &= P(w_{1,n} | S_1 = s^1) \\ &= P(w_{1,n})\end{aligned}$$

$$\begin{aligned}\underline{P\beta}_i(n+1) &= P(\epsilon | S_{n+1} = s^i) \\ &= 1\end{aligned}$$

50

51

52

- We can calculate $\underline{P\beta}_i(t)$ starting from the end of the output sequence and working our way backward, at each step calculating $\beta_i(t)$ for all states s^i .

$$\underline{P\beta}_i(t-1) = \sum_{j=1}^{\sigma} P(w_{t-1}, S_t = s^j | S_{t-1} = s^i) P(w_{t,n} | w_{t-1}, S_t = s^j, S_{t-1} = s^i)$$

$$= \sum_{j=1}^{\sigma} P(w_{t-1}, S_t = s^j | S_{t-1} = s^i) P(w_{t,n} | w_{t-1}, S_t = s^j, S_{t-1} = s^i)$$

$$= \sum_{j=1}^{\sigma} P(w_{t-1}, S_t = s^j | S_{t-1} = s^i) P(w_{t,n} | S_t = s^j)$$

$$= \sum_{j=1}^{\sigma} P(s^i \xrightarrow{w_{t-1}} s^j) P\beta_j(t)$$

51

52

- Computation of $\underline{P\beta}_i(t-1)$ can be formulated as backward recursion.

$$\underline{P\beta}_i(t-1) = P(w_{t-1,n} | S_{t-1} = s^i)$$

$$= \sum_{j=1}^{\sigma} P(w_{t-1,n}, S_t = s^j | S_{t-1} = s^i)$$

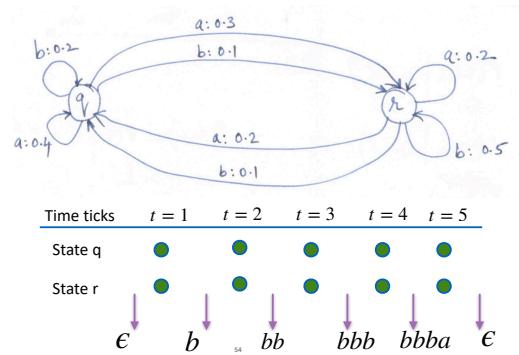
51

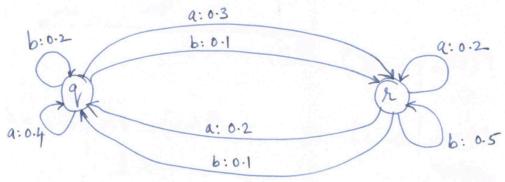
52

- $\underline{P\beta}_i(n+1) = 1 \quad \forall i$

$$\underline{P\beta}_i(t-1) = \sum_{j=1}^{\sigma} P(s^i \xrightarrow{w_{t-1}} s^j) \underline{P\beta}_j(t)$$

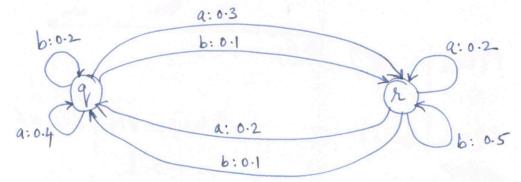
53





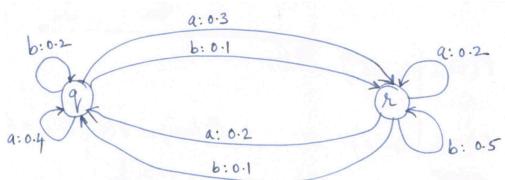
Time ticks	$t = 1$	$t = 2$	$t = 3$	$t = 4$	$t = 5$
$P\alpha_q(t)$	1.0	0.2	0.05	0.017	0.0148
$P\alpha_r(t)$	0.0	0.1	0.07	0.04	0.0131
$P(w_{1,t})$	1.0	0.3	0.12	0.057	0.0279

$= \sum_{i=1}^{\sigma} P\alpha_i(t+1)$



Time ticks	$t = 1$	$t = 2$	$t = 3$	$t = 4$	$t = 5$
State q	●	●	●	●	●
State r	●	●	●	●	●

Below the states, arrows point down to the sequence: bbba at time t=1, bba at time t=2, ba at time t=3, a at time t=4, and ε at time t=5.



Time ticks	$t = 1$	$t = 2$	$t = 3$	$t = 4$	$t = 5$
$P\beta_q(t)$	0.0279	0.063	0.18	0.7	1
$P\beta_r(t)$	0.153	0.27	0.4	1	

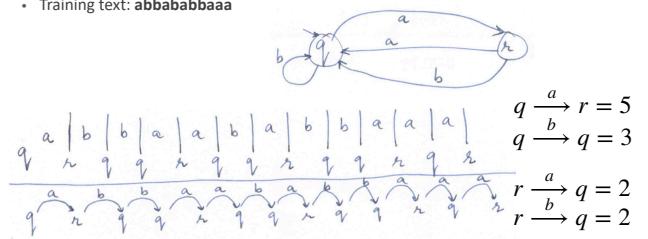
HMM Training: Learning the parameters

HMM Training

- The training algorithm is called as the forward-backward algorithm or the **Baum-Welch algorithm**.
- The algorithm adjusts the probabilities to make the training sequence as likely as possible.
- First, let's consider how to train a Markov chain model

Estimating Probabilities for a Markov Chain

- Given the representative text created by some Markov Chain (with unknown arc probabilities), we want to estimate those probabilities.
- Training text: abbababbaaa



Estimating Probabilities for a Markov Chain

$$P_e(q \xrightarrow{a} r) = \frac{5}{8} \quad P_e(q \xrightarrow{b} q) = \frac{3}{8}$$

- The estimate of the transition probability $P_e(s^i \xrightarrow{w^k} s^j)$ can be written as

$$P_e(s^i \xrightarrow{w^k} s^j) = \frac{C(s^i \xrightarrow{w^k} s^j)}{\sum_{l,m} C(s^i \xrightarrow{w^m} s^l)}$$

61

How to estimate the transition counts for an HMM?

- An HMM has many possible paths which can generate the same output symbol sequence.
- However, unlike a Markov Chain, we do not know which transitions were actually used.
- Therefore, we simply pretend that all possible paths were used, and weigh (prob rate) the transitions according to the probability of the path in which the arc appears.

$$C(s^i \xrightarrow{w^k} s^j) = \sum_{s_{1,n+1}} P(s_{1,n+1} | w_{1,n}) \eta(s^i \xrightarrow{w^k} s^j, s_{1,n}, w_{1,n})$$

62

$$\star C(s^i \xrightarrow{w^k} s^j) = \frac{1}{P(w_{1,n})} \sum_{s_{1,n+1}} P(s_{1,n+1}, w_{1,n}) \eta(s^i \xrightarrow{w^k} s^j, s_{1,n}, w_{1,n})$$

- Here $\eta(s^i \xrightarrow{w^k} s^j, s_{1,n}, w_{1,n})$ counts the number of times the transition appears in the state sequence (path) $s_{1,n+1}$ when generating the output $w_{1,n}$

- But note the chicken and the egg problem

Computing $P(w_{1,n})$, $P(s_{1,n+1}, w_{1,n})$ and path probabilities $P(s_{1,n+1} | w_{1,n})$ require the transition probabilities

63

- The iterative process continues till the parameter estimates converge.
OR
we can compute $P(w_{1,n})$ after every parameter update and declare convergence if the probability assigned to the training corpus does not change much.

- Corpus/training sequence 01011
4 possible paths could produce this output
ababaa
abaaaa
aabaaa
aaaaaa

64

- The term $\frac{1}{P(w_{1,n})}$ need not be calculated since it gets cancelled for the numerator and

the denominator when computing $P_e(s^i \xrightarrow{w_k} s^j)$

- A limitation of the training procedure is that it requires iterating through all possible paths for the input sequence.
- Number of paths grow exponentially with the length of the sequence.
- We make certain changes in our computation so that we can train in a more reasonable time.

65

$$\begin{aligned} C(s^i \xrightarrow{w_k} s^j) &= \sum_{s_{1,n+1}} P(s_{1,n+1} | w_{1,n}) \eta(s^i \xrightarrow{w_k} s^j, s_{1,n}, w_{1,n}) \\ &= \frac{1}{P(w_{1,n})} \sum_{s_{1,n+1}} P(s_{1,n+1}, w_{1,n}) \eta(s^i \xrightarrow{w_k} s^j, s_{1,n}, w_{1,n}) \\ &= \frac{1}{P(w_{1,n})} \sum_{i=1}^n \sum_{s_{1,n+1}} P(s_{1,n+1}, w_{1,n}) \delta_i(s^i \xrightarrow{w_k} s^j, s_{1,n+1}, w_{1,n}) \\ &= \frac{1}{P(w_{1,n})} \sum_{i=1}^n \text{sum}_{s_{1,n+1}} P(S_i = s^i, S_{i+1} = s^j, w_i = w^k, s_{1,n+1}, w_{1,n}) \end{aligned}$$

66

- The function $\delta_t(s^i \xrightarrow{s_k} s^j, s_{1,n}, w_{1,n})$ takes value 1 if the transition is used at time t and 0 otherwise.

- We then get rid of the δ_t function by putting the requirement that the path under consideration would go through the transition $s^i \xrightarrow{s_k} s^j$ at time tick t .

- There is no harm in putting this requirement because if the path does not go through the transition $s^i \xrightarrow{s_k} s^j$ at time t , then the probability will be zero.

- Also note that

$$\sum_{s_{1,n+1}} P(S_t = s^i, S_{t+1} = s^j, w_t = w^k, s_{1,n+1}, w_{1,n}) = P(S_t = s^i, S_{t+1} = s^j, w_t = w^k, w_{1,n})$$

$$\therefore C(s^i \xrightarrow{s_k} s^j) = \frac{1}{P(w_{1,n})} \sum_{t=1}^n P(S_t = s^i, S_{t+1} = s^j, w_t = w^k, w_{1,n})$$

$$= \sum_{t=1}^n P(s^i \xrightarrow{s_k} s^j \text{ at time tick } t \mid w_{1,n})$$

$$= \frac{1}{P(w_{1,n})} \sum_{t=1}^n P(S_t = s^i, S_{t+1} = s^j, W_t = w^k, w_{1,n})$$

67

68

$$= \frac{1}{P(w_{1,n})} \sum_{t=1}^n P(w_{1,t-1}, S_t = s^i, W_t = w^k, S_{t+1} = s^j, w_{t+1,n})$$

$$= \frac{1}{P(w_{1,n})} \sum_{t=1}^n P(w_{1,t-1}, S_t = s^i) P(W_t = w^k, S_{t+1} = s^j \mid w_{1,t-1}, S_t = s^i)$$

- We have already seen how to compute the values for $\underline{P\alpha}$ and $\underline{P\beta}$

$$\underline{P\alpha}_j(t) = P(w_{1,t-1}, S_t = s^j)$$

$$\underline{P\beta}_j(t+1) = P(w_{t+1,n} \mid S_{t+1} = s^j)$$

69

71

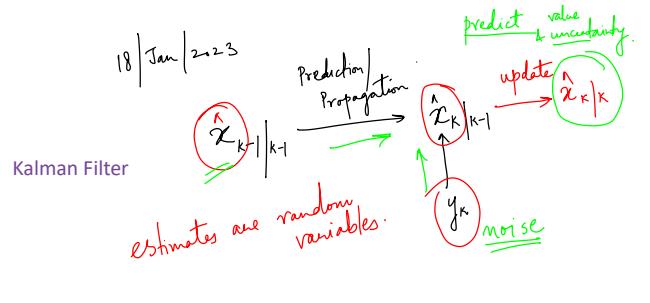
$$\cdot P(w_{t+1,n} \mid w_{1,t}, S_t = s^i, S_{t+1} = s^j)$$

$$= \frac{1}{P(w_{1,n})} \sum_{t=1}^n P(w_{1,t-1}, S_t = s^i)$$

$$P(W_t = w^k, S_{t+1} = s^j \mid S_t = s^i)$$

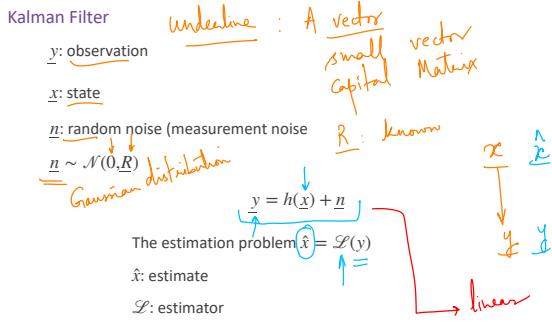
$$P(w_{t+1,n} \mid S_{t+1} = s^j)$$

$$= \frac{1}{P(w_{1,n})} \sum_{t=1}^n \underline{P\alpha}_j(t) P(s^i \xrightarrow{s_k} s^j) \underline{P\beta}_j(t+1)$$



70

72



- If the observation function $h(\cdot)$ is linear, then

$$\underline{y} = \underline{H}\underline{x} + \underline{n}$$

where \underline{y} is a linear function of \underline{h} but with some added noise.

$$\text{Then } \hat{\underline{x}} = \underline{L}\underline{y}$$

where $\underline{L} = \underline{H}^{-1}$

$$\text{Matrix } h(x) = \underline{H}\underline{x}$$

Over constrained System

- If \underline{H} has more rows than the number of columns then there is no inverse of \underline{H} .
- But we can find the best \hat{x} such that the residue $\|\underline{y} - \underline{H}\hat{x}\|$ is as small as possible.
 L_2 norm
- The best estimator for a linear observation function happens to be a linear function.
- We denote the residue as $\underline{n} = \underline{y} - \underline{H}\hat{x}$

$$\begin{array}{c} \hat{x} \\ \downarrow \\ \underline{H}\hat{x} \leftrightarrow \underline{y} \end{array}$$

- Different equations can have noise terms of different variance. So some equations may be more reliable.
- We formulate minimisation of the norm of $\underline{W}\underline{n}$ where \underline{W} is a matrix as

$$\text{Instead of minimizing } \|\underline{n}\|^2$$

$$\|\underline{W}\underline{n}\|^2 = \underline{n}^T \underline{W}^T \underline{W}\underline{n}$$

$$\text{which is equivalent to solving the system}$$

$$\underline{W}\underline{y} = \underline{W}\underline{H}\hat{x}$$

$$\underline{y} - \underline{H}\hat{x} \leftrightarrow \underline{W}\underline{y} = \underline{W}\underline{H}\hat{x}$$

$$\underline{W}^T \underline{W}\underline{y} = \underline{W}^T \underline{W}\underline{H}\hat{x}$$

$$\text{because } \underline{W}\underline{n} = \underline{W}(\underline{y} - \underline{H}\hat{x}) = \underline{W}\underline{y} - \underline{W}\underline{H}\hat{x}$$
- That is we want to solve $\underline{W}\underline{y} = \underline{W}\underline{H}\hat{x}$ subject to minimising $\|\underline{W}\underline{y} - \underline{W}\underline{H}\hat{x}\|^2$

- Consider the standard notation for a system of equations $\underline{A}\underline{x} = \underline{b}$ to be solved for \underline{x} so as to minimise $\|\underline{b} - \underline{A}\underline{x}\|^2$

$$\|\underline{u}\|^2 = \underline{u}^T \underline{u}$$

- The objective to be minimised can be re-written as $\|\underline{A}\underline{x} - \underline{b}\|^2$ or

$$(\underline{A}\underline{x} - \underline{b})^T (\underline{A}\underline{x} - \underline{b}) \quad \text{or} \quad (\underline{x}^T \underline{A}^T - \underline{b}^T)(\underline{A}\underline{x} - \underline{b})$$

- $\frac{df}{d\underline{x}} = 2\underline{A}^T \underline{A} * \underline{x} - 2\underline{A}^T \underline{b} = 0$ for $\underline{A}\underline{x} = \underline{b}$
- Therefore $\underline{A}^T \underline{A}\underline{x} = \underline{A}^T \underline{b}$ or $\hat{\underline{x}} = (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{b}$ pseudo inverse.
- The given system of equations $\underline{W}\underline{H}\hat{x} = \underline{W}\underline{y}$ solve it

$$\text{Substituting } \underline{A} \equiv \underline{W}\underline{H} \text{ and } \underline{b} \equiv \underline{W}\underline{y} \text{ in } \hat{\underline{x}} = (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{b}$$

$$\hat{\underline{x}} = ((\underline{W}\underline{H})^T \underline{W}\underline{H})^{-1} (\underline{W}\underline{H})^T \underline{W}\underline{y} = (\underline{H}^T \underline{W}^T \underline{W}\underline{H})^{-1} \underline{H}^T \underline{W}^T \underline{W}\underline{y}$$

$$= (\underline{H}^T \underline{R}^{-1} \underline{H})^{-1} (\underline{H}^T \underline{R}^{-1} \underline{b})$$

- \underline{R} is the covariance matrix of the measurement noise.
- Entries in \underline{R} which correspond to a large variance or covariance will translate to smaller entries in the \underline{W} matrix.

More noise \Rightarrow uncertainty
Covariance \downarrow
 \downarrow less weightage

$$\underline{R}^{-1} = \underline{W}^T \underline{W}$$

$$y \sim N(0, \underline{R})$$

- Quality of the estimator for the estimate
- The size of the covariance matrix should be small.
 - Covariance matrix $P = \mathbb{E}[(\underline{x} - \hat{\underline{x}})(\underline{x} - \hat{\underline{x}})^T]$ for the estimate $\hat{\underline{x}}$
 - A smaller norm of P implies reduced uncertainty or fluctuations
 - The estimate $\hat{\underline{x}}$ should be unbiased
 - For a linear estimator $\hat{\underline{x}} = \underline{L}\underline{y}$ corresponding to the system $\underline{y} = \underline{H}\underline{x} + \underline{n}$
- $$\mathbb{E}[\underline{x} - \hat{\underline{x}}] = 0 \quad \underline{x} - \mathbb{E}[\hat{\underline{x}}] = 0 \quad \mathbb{E}[\hat{\underline{x}}] = \underline{x}$$
- $$\mathbb{E}[\underline{x} - \hat{\underline{x}}] = \mathbb{E}[\underline{x} - \underline{L}\underline{y}] = \mathbb{E}[\underline{x} - \underline{L}(\underline{H}\underline{x} + \underline{n})] = \mathbb{E}[(\underline{I} - \underline{L}\underline{H})\underline{x}] - \mathbb{E}[\underline{L}\underline{n}] = 0$$
- 3.

- $\mathbb{E}[(\underline{I} - \underline{L}\underline{H})\underline{x}] = 0$
 - or $(\underline{I} - \underline{L}\underline{H})\mathbb{E}[\underline{x}] = 0$
 - $\therefore \underline{I} = \underline{L}\underline{H}$ or $\underline{L}\underline{H} = \underline{I}$
 - $\mathbb{E}[\hat{\underline{x}}] = \underline{x}$
- $\underline{y} = \underline{H}\underline{x} + \underline{n}$
 $\hat{\underline{x}} = \underline{L}\underline{y}$.
 $\hat{\underline{x}} = \underline{H}^{-1}\underline{y}$.

- The covariance matrix of our estimate
- $$\underline{P} = \mathbb{E}[(\underline{x} - \hat{\underline{x}})(\underline{x} - \hat{\underline{x}})^T]$$
- $$= \mathbb{E}[(\underline{x} - \underline{L}\underline{y})(\underline{x} - \underline{L}\underline{y})^T]$$
- Substituting $\underline{y} = \underline{H}\underline{x} + \underline{n}$ we have
- $$\underline{P} = \mathbb{E}[(\underline{x} - \underline{L}\underline{H}\underline{x} - \underline{L}\underline{n})(\underline{x} - \underline{L}\underline{H}\underline{x} - \underline{L}\underline{n})^T]$$
- $$= \mathbb{E}[[\underline{(I - L H)x - L n}] [\underline{(I - L H)x - L n}]^T] = \mathbb{E}[(\underline{L n})(\underline{L n})^T]$$
- L n n^T L

\underline{n} : measurement noise

$$\begin{aligned} &= \mathbb{E}[\underline{L n n^T L}] \\ &= \underline{L} \mathbb{E}[\underline{n n^T}] \underline{L} \\ \underline{P} &= \underline{L} \underline{R} \underline{L}^T \end{aligned}$$

We define $P \equiv \underline{L} \underline{R} \underline{L}^T$

- Here \underline{R} is the covariance matrix of the measurement noise.
- Substituting for the estimator $\underline{L} = (\underline{H}^T \underline{R}^{-1} \underline{H})^{-1} \underline{H}^T \underline{R}^{-1}$ we get