

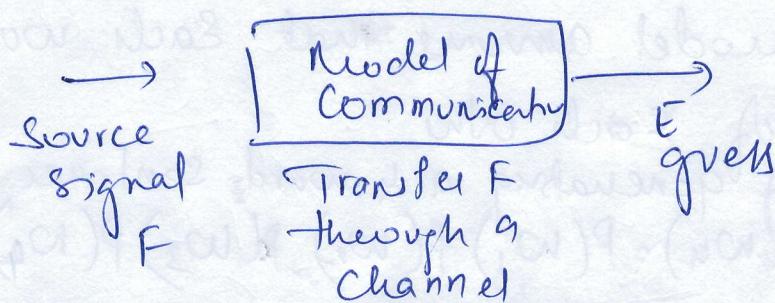
Statistical Machine Translation (SMT)

Given a foreign sentence F , we seek the English sentence E that maximizes $P(E|F)$. The most likely translation is written as

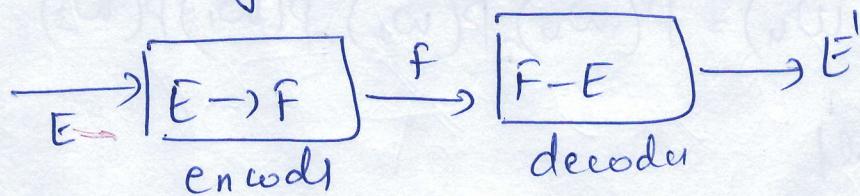
$$\operatorname{argmax}_E P(E|F)$$

The argmax can be read as "The English sentence E , out of all such sentences which yields the highest value for $P(E|F)$ ".

This problem is modelled using noisy channel model



Given F guess E



This model produces E' not E which is an approximation for E . The idea is to produce E' as close as we can be to E .

$$E' = \operatorname{argmax}_E P(E|F) \approx \operatorname{argmax}_E (P(F|E) P(E)) \quad \text{①}$$

$\underbrace{E}_{\substack{\text{lexical} \\ \text{translation} \\ \text{model}}}$ \uparrow Language model

The language model must predict the probability of generating the sentence $P(E)$.

A sentence generation can be modelled as a sequence of words using a chain rule

$$P(w) = P(w_1) P(w_2|w_1) P(w_3|w_2, w_1) \dots P(w_n|w_{n-1}, w_{n-2})$$

↑
This is too lengthy.

To get rid of this lengthy subsequence we apply markov assumption

$$P(w_i|w_{i-1}, w_{i-2} \dots w_1) = P(w_i|w_{i-1} \dots w_{i-n+1})$$

i.e Unigram model assumes that each word is independent of each other

The probability of generating a 4 word sentence is:

$$P(w_1, w_2, w_3, w_4) = P(w_1) P(w_2) P(w_3) P(w_4)$$

drawback of this is word order does not matter.

$$P(w_1, \dots, w_4) = P(w_4) P(w_3) P(w_2) P(w_1)$$

Bigram model

$$P(w_1, \dots, w_4) = P(w_1 | \text{begin of sentence}) P(w_2|w_1) P(w_3|w_2) P(w_4|w_3)$$

The probability of each word is conditioned on the previous word only.

Trigram model

$$P(w_1, \dots, w_4) = P(w_1 \text{ (<begin of sentence)}) P(w_2 | w_1) P(w_3 | w_2, w_1) \\ P(w_4 | w_3, w_2)$$

The longer the better syntactically correct sentences will be generated. But too longer subsequences can be generated with nearly zero probability.

How to compute $P(w_3 | w_2, w_1)$?

$$P(w_3 | w_2, w_1) = \frac{\text{Count}(w_1, w_2, w_3)}{\text{Count}(w_1, w_2)}$$

All parameters can be computed in this form for language models.

why $\arg\max P(E|F)$ has to be modeled as a product of $\arg\max_E P(F|E) P(E)$?

The idea is to estimate $P(E) \propto P(F|E)$ and take the product & then rank all candidates English sentences based on this product and pick the one that is highest.

Consider translating a foreign sentence F to E

$$\overline{F} \in \overline{D}_{\text{MT}} \in E$$

Op Sentence	$P(E)$	$P(F E)$	$P(F E) P(E)$
This rose is	low	high	low
is this rose	low	high	low
a rose this	low	high	low
This pink is	low	low	low
pink is good	low	low	low
This is rose	high	high	high

What factors influence $P(F|E)$?

1. Alignment factor a

How to assign values for $P(F|E)$?

$$P(F|E) = \frac{\text{Count}(F, E)}{\text{Count}(E)} \leftarrow \text{Sentence level}$$

Not possible to get true translation of sentences
since the no of such exact sentences is very

less. Hence we model this using a hidden variable a , that represent alignment between the individual words in a sentence as marginalization over all alignment variables.

$$P(F|E) = \sum_a P(f, a | e) \leftarrow \text{word level}$$

$$\sum_a P(a, f | e)$$

According to Marginalization rule

For continuous random variables $P(x) = \int_y P(x, y) dy$

For discrete random variables

$$P(x) = \sum_y P(x, y)$$

Since our alignment variables are discrete $P(F|E)$ is marginalized using alignment variables as $\sum_a P(a, f|e)$ at word level.

E : vector of English words

$i=1$ to 1
 $j=1$ to 2
 a 1 2
This is nose

F : vector of Foreign "

$i=1$ to 4
 $j=1$ to 2
 a 1 2

a : vector of alignment indices

$$a[1] \rightarrow 1$$

$$a[2] \rightarrow 3$$

$$a[3] \rightarrow 2$$

$$a = \langle 1 3 2 \rangle$$

$P(a, f|e)$ is called alignment

Probability.

Interpretation of a :

$$P(1 3 2 | "This is rose") > P(1 2 3 | "This is rose")$$

How probable are the alignment a and the translation f given e ?

2) length factor of $F(l_f)$

$$P(F, a | E) = \sum_{l_f} P(F, a, l_f | E)$$

Again we introduce the marginalization over the length of foreign sentence l_f .

$$P(F, a, l_f | E) = P(l_f | E) P(F, a | E, l_f)$$

Apply chain rule by which joint probabilities are converted into conditional probabilities.

Now if a 3 words sentence are given in $E \# F$

$$P(F, a, 3 | E) = \underbrace{P(3 | E)}_{\text{Probability of generating 3 worded foreign sentence}} P(f_1 a_1, f_2 a_2, f_3 a_3 | e, 3)$$

which can be generalized to a length l_f in foreign sentence as

$$P(l_f | E) \prod_{j=1}^{l_f} P(f_j, a_j | f_1^{j-1}, a_1^{j-1}, e, n)$$

Assumption is that there is one-to-one mapping between the words in E and F

If the English sentence has l_e words

$$E = E^{l_e} = e_1, e_2, \dots, e_{l_e}$$

Foreign sentence has l_f words

$$F = F^{l_f} = f_1, f_2, \dots, f_{l_f}$$

then the alignment a , can be represented by a series $a^{l_f} = a_1, a_2, \dots, a_{l_f}$ of l_f values each between 0 to l_e such that if the word in position j of the foreign sentence is connected to the word in position i of the English sentence then $a_j = i$ & if it is not connected to any English word then $a_j = 0$.

$$P(F, a | E) = P(l_f | E) \prod_{j=1}^{l_f} P(a_j | a_1^{j-1}, f_1^{j-1}, l_f, E) \rightarrow ②$$

$$\prod_{j=1}^{l_f} P(f_j | a_1^j, f_1^{j-1}, l_f, E)$$

where $P(l_f | E)$ is called as length probability

$\prod_{j=1}^{l_f} P(a_j | a_1^{j-1}, f_1^{j-1}, l_f, E)$ is an alignment probability

$\prod_{j=1}^{l_f} P(f_j | a_1^j, f_1^{j-1}, l_f, E)$ is called translation probability

IBM-Model I

In the equation ② the conditional probabilities on the RHS cannot be taken as independent parameters.

Hence IBM Model-I makes assumptions

① $\Pr(l_f | E)$ is independent of E as it is fixed as a constant parameter ϵ :

② $P(a_j^j | a_1^{j-1}, f_1^{j-1}, l_f, E)$ depends on l_e , hence must be $\frac{1}{(l_e + 1)^{l_f}}$

③ $P(f_j^j | a_1^j, f_1^{j-1}, l_f, E)$ depends on f_j and a_{aj}^j

\therefore Estimation of translation probabilities for model-I is the joint likelihood of a Foreign String F and an alignment given an English String E

$$P(F, a | E) = \frac{\epsilon}{(l_e + 1)^{l_f}} \prod_{j=1}^{l_f} t(f_j^j | e_{aj}^j)$$

$$P(F | E) = \frac{\epsilon}{(l_e + 1)^{l_f}} \sum_{a_1=0}^{l_e} \sum_{a_2=0}^{l_e} \dots \sum_{a_{lf}=0}^{l_e} \prod_{j=1}^{l_f} t(f_j^j | e_{aj}^j)$$

We wish to adjust the translation probabilities so as to maximize $P(F | E)$ subject to the constraint that for each

$$\sum_f t(f | e) = 1$$

There are many ways to define $P(F|E)$.

IBM Model-1 assumes lexical translations (word to word)
i.e Each word in Foreign sentence is a translation
of exactly zero/more of English sentence.

$$F = \langle f_1, \dots, f_{l_f} \rangle \quad E = \langle e_1, \dots, e_{l_e} \rangle$$

Assumption: Each foreign word aligned to exactly
one English word, we can represent an alignment
of the French words by an array a of
length l_f $\langle a_1, \dots, a_{l_f} \rangle$

Each a_i takes a value 0 to be denoted
the index of English word to which the Foreign
word f_i is aligned to.

If $a_i = 0$ means the foreign word f_i is aligned
to Null.

Consider a sentence pair $\langle F, E \rangle$

This 2 3
is ~~is~~ rose
I ~~I~~ ~~E~~
The ~~old~~ ~~E~~

$$a = \langle 1, 3, 2 \rangle \quad l_e = 3 \quad l_f = 3$$

~~A length l_f is chosen according to a distribution~~
 $P(l_f | l_e)$ in this example $P(3 | 3)$

Then each Foreign word portion aligns to an English word (or null) according to uniform distribution

$$P(a_i=j | le) = \frac{1}{le+1} \text{ in this case } \frac{1}{3+1}$$

Finally each foreign word f_i is translated according to a ~~uniform~~ distribution $P(f_i | a_i)$ conditioned on the aligned English word $P(f_i | e_{a_i})$

Hence for this alignment we multiply

$$P(\text{rose} | \text{the}) P(\text{dotted} | \text{rose}) P(\text{is} | \text{is})$$

How to estimate $P(f | e)$

Suppose you have ever seen a English word rose aligned with dotted & dotted
 ↓ true ↓ is true

$$P(\text{dotted} | \text{rose}) = 0.7 \quad P(\text{dotted} | \text{rose}) = 0.3$$

Since these correspond to the proportions in the data you observed.

$$\therefore P(\text{dotted} | \text{rose}) = \frac{\text{Count}(\text{dotted}, \text{rose})}{\text{Count}(\text{dotted}, \text{rose}) + \text{Count}(\text{dotted}, \text{not rose})} = \frac{3}{10} = 0.3$$

$$P(\text{dotted} | \text{not rose}) = \frac{\text{Count}(\text{dotted}, \text{not rose})}{\text{Count}(\text{dotted}, \text{not rose})}$$

$$P(F|E) = \frac{\epsilon}{(k_e+1)^{lf}} \prod_{j=1}^{lf} \sum_{i=0}^{le} t(f_j | e_{ai})$$

IBM Model-1 & EM

Need the expected no of times word e connects to word f in the translation $(f|e)$

This is the (expected) count of f given e for $(f|e)$

$$C(f|e, F, E) = \sum_a P(a|E, F) \sum_{j=1}^{lf} \delta(f, f_j) \delta(e, e_{aj})$$

denotes the no of times
connects to f in a

\therefore we need to compute $P(a|E, F)$

$$P(a|E, F) = \frac{P(F, a|E)}{P(F|E)} \rightarrow \text{we have derive this for IBM model-1}$$

$$= \frac{\epsilon / (k_e+1)^{lf}}{\prod_{j=1}^{lf} \sum_{i=0}^{le} t(f_j | e_{aj})}$$

$$= \frac{\epsilon / (k_e+1)^{lf}}{\prod_{j=1}^{lf} \sum_{i=0}^{le} t(f_j | e_i)}$$

$$= \frac{\prod_{j=1}^{lf} t(f_j | e_{aj})}{\sum_{i=0}^{le} t(f_j | e_i)}$$

Maximization Step

Now we to collect the counts

Evidence from a sentence pair E, F that a word f is a translation of word e .

$$c(f|e; E, F) = \sum_a P(a|E, F) \sum_{j=1}^{|f|} \delta(f, f_j) \delta(e, e_{aj})$$

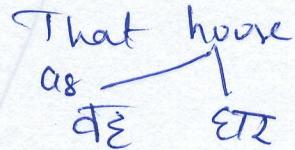
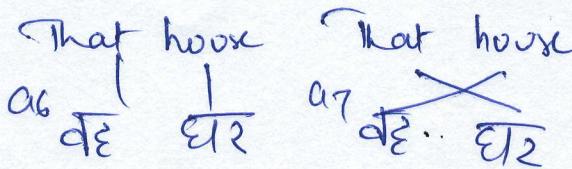
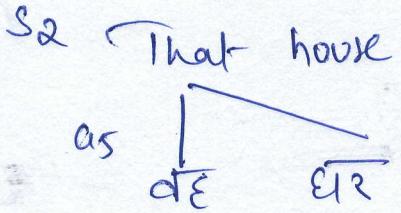
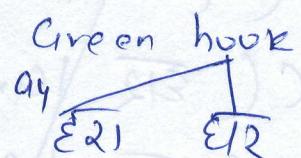
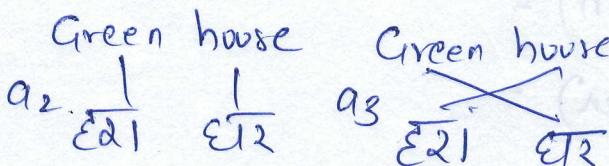
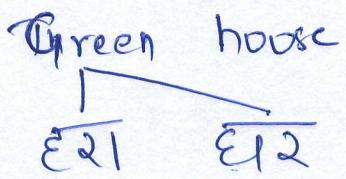
$$= \frac{t(f|e)}{\sum_{i=0}^{|e|} t(f|e_i)} \sum_{j=1}^{|f|} \delta(f, f_j) \delta(e, e_j)$$

After collecting the counts over the corpus, we can estimate the model as

$$t(f|e; E, F) = \frac{\sum_{(E, F)} c(f|e; E, F)}{\sum_f \sum_{(E, F)} c(f|e; E, F)}$$

Learning parameters : IBM Model - I

s_1



Given \rightarrow Initial all $t(\text{file})$ uniformly

Green	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$
house	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$
That	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$

Compute posterior for each alignment

$$P(a_1, \text{file}) = P(\overline{ET1} | \text{Green}) P(\overline{ET2} | \text{Green}) = \frac{1}{3} \times \frac{1}{3} = \frac{1}{9}$$

$$P(a_2, \text{file}) = P(\overline{ET1} | \text{Green}) P(\overline{ET1} | \text{house}) = \frac{1}{3} \times \frac{1}{3} = \frac{1}{9}$$

!

!

-

$$P(a_5, \text{file}) = P(\overline{dE} | \text{That}) P(\overline{ET2} | \text{That}) = \frac{1}{3} \times \frac{1}{3} = \frac{1}{9}$$

!

$$P(a_8) = \frac{1}{9}$$

Normalize alignments

i.e. compute $P(a | F, E) = \frac{P(a, F | E)}{P(F | E)}$

$$P(F | E) = a_1 + a_2 + a_3 + a_4 = \frac{1}{9} + \frac{1}{9} + \frac{1}{9} + \frac{1}{9} = \frac{4}{9}$$

$$\frac{P(a_1, f | e)}{P(F | E)} = \frac{\frac{1}{9}}{\frac{4}{9}} = \frac{1}{4}$$

$$P(a_2) = \frac{1}{4}$$

:

$$P(a_4) = \frac{1}{4}$$

$$P(a_5, f | e) = \frac{1}{4}$$

:

$$P(a_8, f | e) = \frac{1}{4}$$

Compute fractional counts

$$C(\overline{E}\overline{Q}\overline{I}, \text{Green}) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2} \quad (\text{a}_1, \text{a}_2)$$

$$C(\overline{E}\overline{R}, \text{Green}) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2} \quad (\text{a}_1, \text{a}_3)$$

$$C(d\overline{E}, \text{Green}) = 0 \quad (\text{These words never occurred in any alignments})$$

$$C(\overline{E}\overline{Q}\overline{I}, \text{house}) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2} \quad (\text{a}_3, \text{a}_4)$$

$$C(\overline{E}\overline{R}, \text{house}) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2} \quad (\text{a}_2, \text{a}_4)$$

$$C(d\overline{E}, \text{house}) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2} \quad (\text{a}_7, \text{a}_8)$$

$$C(\overline{E}\overline{Q}\overline{I}, \text{that}) = 0$$

$$C(\overline{E}\overline{R}, \text{that}) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2} \quad (\text{a}_5, \text{a}_7)$$

$$C(d\overline{E}, \text{that}) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2} \quad (\text{a}_5, \text{a}_6)$$

$$\text{total count (Green)} = \frac{1}{2} + \frac{1}{2} + 0 = 1$$

$$(\text{house}) = \frac{1}{2} + \frac{1}{2} + \frac{1}{2} = \frac{3}{2}$$

$$(\text{that}) = 0 + \frac{1}{2} + \frac{1}{2} = 1$$

Maximization step

$$t(\overline{ex1} | \text{Green}) = \frac{1}{2}$$

$$t(\overline{ex2} | \text{Green}) = \frac{1}{2} / \frac{3}{2} = \frac{1}{3}$$

$$t(\overline{d\bar{e}} | \text{Green}) = 0$$

$$t(\overline{ex1} | \text{house}) = \frac{1}{2} / \frac{3}{2} = \frac{1}{3}$$

$$t(\overline{ex2} | \text{house}) = \frac{1}{2} / \frac{3}{2} = \frac{1}{3}$$

$$t(\overline{d\bar{e}} | \text{house}) = \frac{1}{2} / \frac{3}{2} = \frac{1}{3}$$

$$t(\overline{ex1} | \text{That}) = 0$$

$$t(\overline{ex2} | \text{That}) = \frac{1}{2} / 1 = \frac{1}{2}$$

$$t(\overline{d\bar{e}} | \text{That}) = \frac{1}{2} / 1 = \frac{1}{2}$$

Keep iterating over Expansion & Maximization steps until convergence of parameters $t(z_j | e_{aj})$

IBM Model-2

In IBM model-1 all alignments are Equally Likely which is wrong. Since the languages having SOV order like most Indian languages the Verbs Predominately move to the end of the Sentence. Once the position of the noun is fixed adjectives qualify these nouns. Hence it is realistic to assume a probability distribution over position of word as

$$a(\underset{\text{English}}{a_j} | \underset{\text{Foreign}}{j}, \text{le}, \text{lf})$$

Where j is the position of the foreign word and

a_j is the position of English word. The function $a(a_j, j, \alpha)$ models the probability of word in the position a_j in source sentence being reordered to j^{th} position in the target sentence.

$P = \text{That house is beautiful}$

$$\alpha = \langle 1, 2, 4, 3 \rangle$$

$j = \overline{dE} \quad \overline{E12} \quad \cancel{\overline{s342}} \quad \overline{i}$

$$t(\overline{dE} | \text{That}) a(1 | 1, 4, 4) \times$$

$$t(\overline{E12} | \text{house}) a(2 | 2, 4, 4) \times$$

$$t(\cancel{\overline{s342}} | \text{beautiful}) a(4 | 3, 4, 4) \times$$

$$t(\overline{i} | \text{is}) a(3 | 4, 4, 4)$$

Another

Null Peter advised me to invest in shares

पीटर ने मुझे शेयर में पैसा लगाने का सुझाव दिया।

$a = \langle 1 \ 0 \ 3 \ 7 \ 6 \ 5 \ 5 \ 4 \ 2 \ 2 \rangle$

$t(\text{पीटर} | \text{Peter}) \ a(1|1, 7, 10)$

$t(\text{नों} | \text{Null}) \ a(0|2, 7, 10)$

$t(\text{मुझे} | \text{me}) \ a(3|3, 7, 10)$

$t(\text{शेयर} | \text{shares}) \ a(7|4, 7, 10)$

$t(\text{में} | \text{in}) \ a(6|5, 7, 10)$

$t(\text{पैसा} | \text{invest}) \ a(5|6, 7, 10)$

$t(\text{लगाने} | \text{invest}) \ a(5|7, 7, 10)$

$t(\text{का} | \text{to}) \ a(4|8, 7, 10)$

$t(\text{सुझाव} | \text{invest}) \ a(2|9, 7, 10)$

~~$t(\text{दिया} | \text{invest}) \ a(2|10, 7, 10)$~~

$P(f_1, f_2, \dots, f_{14}, a_1, a_2, \dots, a_{14} | e_1, e_2, \dots, e_{14})$

$\underbrace{\prod_{j=1}^{14} a(a_j | j, 1e, 1f)}_{\text{j}} + t(f_j | e_{aj})$

To estimate $a(a_j | j, le, lf)$ we use EM as we did in Model-1 for translation parameters.

$$a(a_j | j, le, lf) = \frac{c(a_j | j, le, lf)}{c(j, le, lf)}$$

Where $c(a_j | j, le, lf)$ is the no of times we see an English sentence of length le & foreign sentence length lf where word j in the foreign sentence is aligned to word a_j in the English sentence.

$c(j, le, lf)$ the no of times we see English & foreign sentences with length le & lf respectively.

IBM Model-3

The assumption in IBM Models 1 & 2 is that each word in English generates one or many foreign words. but their distribution is uniform whether it generates one word or many words.

This factor is handled as a fertility factor in IBM Model 3.

for each English word e in English sentence E , we choose a fertility ϕ . The choice of ϕ depends only on e which is modeled as a probability distribution parameter $n(\phi|e)$. These parameters for each English word e is computed from the corpus.

The parameter $n(\phi|e)$ is used for non Null words only.

The fertility of a Null token depends on the context of English words & foreign sentence length.

Hence the fertility of a Null word is treated differently. $n(1|\text{NULL})$, $n(2|\text{NULL})$, ..., $n(25|\text{NULL})$, ... needs to be computed which is infeasible.

∴ we represent a parameter P_i which is the probability of having No Null token inserted after each English word e . P_i is the probability of generating one Null token inserted after each English word e .

The number of Foreign words generated from a Null token is represented as ϕ_0 .

$$\text{No of OP words} = \sum_{i=1}^{l_f} \phi_i = l_f - \phi_0$$

The estimation of ϕ_0 is treated as a binomial distribution since each English word e may generate a Null (P_1) or not (P_0).

$P(\phi_0)$ is the probability of generating ϕ_0 foreign words from a null token is estimated as

$$P(\phi_0) = \binom{l_f - \phi_0}{\phi_0} P_1^{\phi_0} P_0^{l_f - \phi_0} \quad (\text{using binomial distribution})$$

2BM Model-3 Parameters:

1. word / lexical translations from 1BM Model 1: $t(f_j | e_i)$
2. Fertility of Non Null words from model 3: $n(\phi_i | e_i)$
3. Fertility of Null words: $P(\phi_0)$
4. Distortion probability distribution which will be in the reverse direction of alignment probability shown in 2BM Model 2. $d(j | a(j), le, lt)$

$$\therefore P(F|E) = \sum_a P(F, a | E)$$

$$= \sum_{a_1=0}^{l_e} \cdots \sum_{a_{l_f}=0}^{l_e} P(F, a | E)$$

$$= \sum_{a(1)=0}^{l_e} \cdots \sum_{a(l_f)=0}^{l_e} \frac{l_f}{\prod_{j=1}^f} \binom{l_f - \Phi_0}{\Phi_0} P_1^{\Phi_0} P_0^{l_f - 2\Phi_0} \times$$

$$\frac{l_e}{\prod_{i=1}^e} \phi_i n(\phi_i | e_i) \times$$

$$\frac{l_f}{\prod_{j=1}^f} t(f_j | e_{a_j}) d(j | a(j), l_e, l_f)$$