

# The Translation Model

## IBM Model 1

$$P(F|E) = \sum_A P(F, A|E)$$

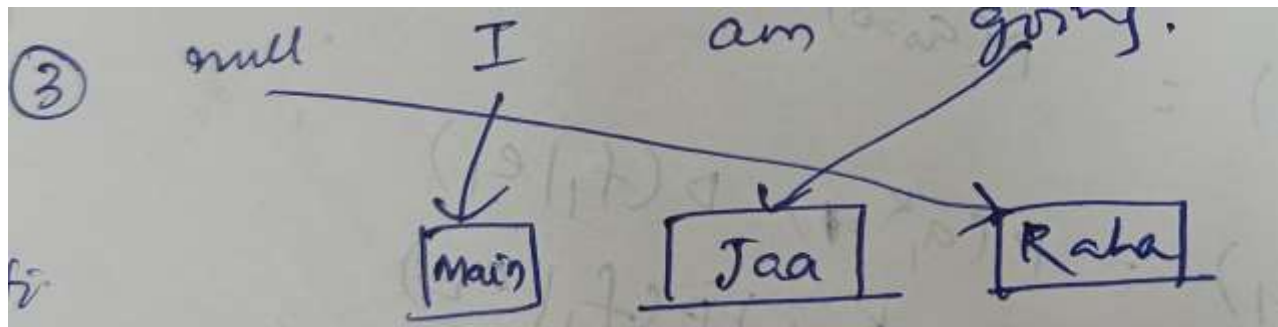
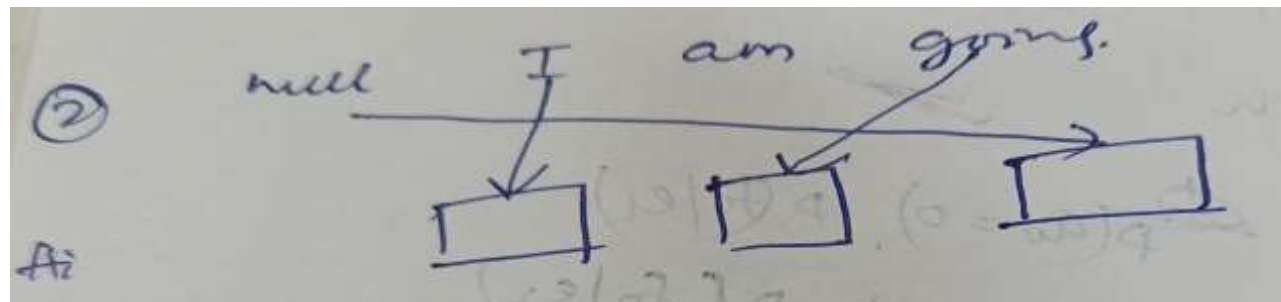
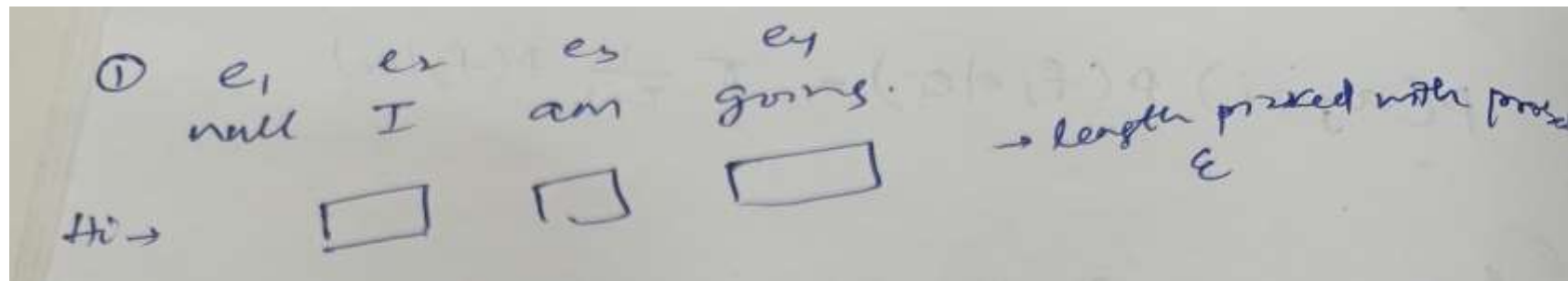
We start with IBM Model 1, so called because it is the first and simplest of five models proposed by IBM researchers in a seminal paper (Brown et al., 1993).

Here's the general IBM Model 1 generative story for how we generate a Spanish sentence from an English sentence  $E = e_1, e_2, \dots, e_I$  of length  $I$ :

1. Choose a length  $J$  for the Spanish sentence, henceforth  $F = f_1, f_2, \dots, f_J$ .
2. Now choose an alignment  $A = a_1, a_2, \dots, a_J$  between the English and Spanish sentences

models proposed by IBM researchers in a seminal paper (2002).  
 Here's the general IBM Model 1 generative story for how we generate a Spanish sentence from an English sentence  $E = e_1, e_2, \dots, e_I$  of length  $I$ :

1. Choose a length  $J$  for the Spanish sentence, henceforth  $F = f_1, f_2, \dots, f_J$ .
2. Now choose an alignment  $A = a_1, a_2, \dots, a_J$  between the English and Spanish sentences.
3. Now for each position  $j$  in the Spanish sentence, choose a Spanish word  $f_j$  by translating the English word that is aligned to it.



$$A = \{ \text{main} \rightarrow I, \text{Jaa} \rightarrow \text{going}, \text{Raha} \rightarrow \text{null} \}$$

- $e_{a_j}$  is the English word that is aligned to the Spanish word  $f_j$

$e_{a_1} \rightarrow I$  ,  $e_{a_2} \rightarrow \text{going}$   $e_{a_3} \rightarrow \text{null}$

- $t(f_x|e_y)$  is the probability of translating  $e_y$  by  $f_x$  (i.e.,  $P(f_x|e_y)$ )

$$t(\text{man}|\text{null}) = 0.1 \quad t(\text{man}|I) = 0.3 \quad t(\text{man}|\text{am}) = 0.1$$

$$t(\text{man}|\text{going}) = 0.2$$

$$t(\text{Jaa}|\text{null}) = 0.1 \quad t(\text{Jaa}|I) = 0.1 \quad t(\text{Jaa}|\text{am}) = 0.1$$

$$t(\text{Jaa}|\text{going}) = 0.5$$

$$t(\text{Raha}|\text{null}) = 0.1 \quad t(\text{Raha}|I) = 0.1 \quad t(\text{Raha}|\text{am}) = 0.1$$

$$t(\text{Raha}|\text{going}) = 0.3$$



We'll work our way backwards from step 3. So suppose we already knew the length  $J$  and the alignment  $A$ , as well as the English source  $E$ . The probability of the Spanish sentence would be

$$P(F|E, A) = \prod_{i=1}^J t(f_j | e_{a_j}) \quad (25.17)$$

$$P(\text{main Taa Raha} | \text{null I am going}, A) = t(\text{main} | \text{I}) * \\ t(\text{Taa} | \text{going}) * t(\text{Raha} | \text{null})$$

$$= 0.1 * 0.5 * 0.1$$

Now let's formalize steps 1 and 2 of the generative story. This is the probability  $P(A|E)$  of an alignment  $A$  (of length  $J$ ) given the English sentence  $E$ . IBM Model 1 makes the (very) simplifying assumption that each alignment is equally likely. How many possible alignments are there between an English sentence of length  $I$  and a Spanish sentence of length  $J$ ? Again assuming that each Spanish word must come from one of the  $I$  English words (or the 1 NULL word), there are  $(I+1)^J$  possible alignments. Model 1 also assumes that the probability of choosing length  $J$  is some small constant  $\epsilon$ . The combined probability of choosing a length  $J$  and then choosing any particular one of the  $(I+1)^J$  possible alignments is

$$P(A|E) = \frac{\epsilon}{(I+1)^J} \quad (25.18)$$

$$P(A | \text{null I am going}) = \frac{\epsilon}{4^3}$$

$$(I+1)^J$$

(25.18)

We can combine these probabilities as follows:

$$P(F, A|E) = P(F|E, A) \times P(A|E)$$

$$= \frac{\epsilon}{(I+1)^J} \prod_{j=1}^J t(f_j|e_{a_j})$$

(25.19)

$$P(\text{main jaa raha, } A \mid \text{null I angong}) =$$

$$\frac{\epsilon}{43} * 0.005$$



$(I+1) \prod_{j=1}^J$  (25.19)

This probability,  $P(F, A|E)$ , is the probability of generating a Spanish sentence  $F$  through a particular alignment. To compute the total probability  $P(F|E)$  of generating  $F$ , we just sum over all possible alignments:

$$\begin{aligned} P(F|E) &= \sum_A P(F, A|E) \\ &= \sum_A \frac{\epsilon}{(I+1)^J} \prod_{j=1}^J t(f_j|e_{a_j}) \end{aligned} \quad (25.20)$$



# Recap

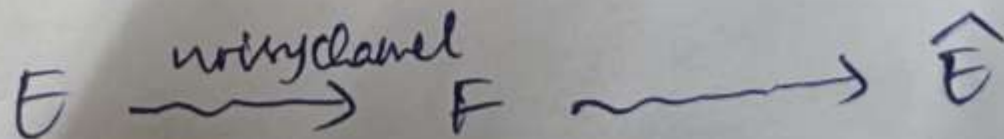
Statistical MT, - 27.03.2021

Given a foreign sentence  $F \rightarrow$   
we want to predict the most likely  
sentence  $E$  with the same meaning.

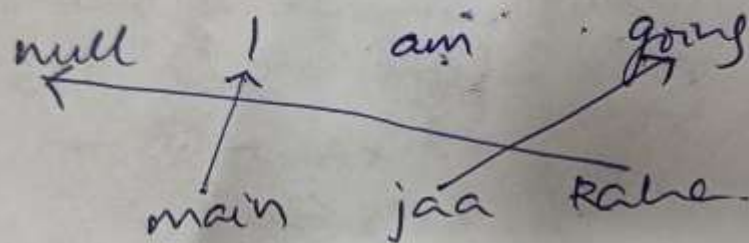
$$\hat{E} = \operatorname{argmax}_E P(E|F)$$

$\rightarrow$  translation model.

$$= \operatorname{argmax}_E P(F|E) * P(E) \rightarrow \text{language model}$$



Alignment:- is a function from output to input of the translation model.



$$P(F, A | E) = \sum_A P(F, A | E)$$

$$= \sum_{j=1}^L \frac{\epsilon}{(I+1)^J} \prod_{j=1}^J t(f_j | e_{a_j})$$

$$\hat{A} = \underset{A}{\operatorname{argmax}} P(F, A | E)$$

~~$$\underset{A}{\operatorname{argmax}} \prod t(f_j | e_{a_j})$$~~

## Training alignment models

sentence level alignment

word level alignment -

for each sentence pair  $(F_s, E_s)$  we need to learn an alignment  $A = a_i^T$  and the corresponding translation probabilities.

$$P(A|E, F) = \frac{P(A, F|E)}{\sum_A P(A, F|E)} \rightarrow \text{if we knew translation probabilities}$$

$$t(\text{main}/I) = \frac{C(\text{main}, I)}{C(I)} \rightarrow \text{if we knew the perfect alignment}$$

$$P(F, A|E) \approx \prod_i t(f_i|e_{a_i}) \rightarrow \text{simplification.}$$



green house the house  
mera ghar EK ghar

①  $V_H = \{ \text{mera, ghar, EK} \}$   $V_E = \{ \text{green, house, the} \}$

$$\begin{aligned} t(\text{mera}|\text{green}) &= \frac{1}{3} & t(\text{ghar}|\text{green}) &= \frac{1}{3} & t(\text{EK}|\text{green}) &= \frac{1}{3} \\ t(\text{mera}|\text{house}) &= \frac{1}{3} & t(\text{ghar}|\text{house}) &= \frac{1}{3} & t(\text{EK}|\text{house}) &= \frac{1}{3} \\ t(\text{mera}|\text{the}) &= \frac{1}{3} & t(\text{ghar}|\text{the}) &= \frac{1}{3} & t(\text{EK}|\text{the}) &= \frac{1}{3} \end{aligned}$$

Ex Step 1:-

green house	green house	the house	the house
mera ghar	<del>mera ghar</del>	EK ghar	<del>EK ghar</del>
(a) -	$\frac{1}{9}$	$\frac{1}{9}$	$\frac{1}{9}$

$$P(a, f|e) = \frac{1}{3} \times \frac{1}{3} = \frac{1}{9}$$

(b) normalize  $P(a, e|f) = \frac{P(a, e|f)}{\sum_a P(a, e|f)}$

~~$P(a, e|f)$~~   $\frac{\frac{1}{9}}{\frac{1}{9} + \frac{1}{9} + \frac{1}{9}} = \frac{1}{3}$

(c)  $t(\text{count}(\text{mera}|\text{green})) = \frac{1}{2}$   $t(\text{ghar}|\text{green}) = \frac{1}{2}$   $t(\text{EK}|\text{green}) = 0$  total greens = 1  
 $t(\text{count}(\text{mera}|\text{house})) = \frac{1}{2}$   $t(\text{ghar}|\text{house}) = \frac{1}{2} + \frac{1}{2}$   $t(\text{EK}|\text{house}) = \frac{1}{2}$  wise 2  
 $t(\text{count}(\text{mera}|\text{the})) = 0$   $t(\text{ghar}|\text{the}) = \frac{1}{2}$   $t(\text{EK}|\text{the}) = \frac{1}{2}$  the = 1

$$\begin{aligned} t(\text{mera}|\text{green}) &= \frac{1/2}{1} = \frac{1}{2} & t(\text{ghar}|\text{green}) &= 1/2 / 1 = \frac{1}{2} & t(\text{EK}|\text{green}) &= 0 \\ t(\text{mera}|\text{house}) &= \frac{1/2}{2} = \frac{1}{4} & t(\text{ghar}|\text{house}) &= \frac{1}{2} & t(\text{EK}|\text{house}) &= \frac{1/2}{2} = \frac{1}{4} \\ t(\text{mera}|\text{the}) &= 0 & t(\text{ghar}|\text{the}) &= \frac{1}{2} & t(\text{EK}|\text{the}) &= \frac{1}{2} \end{aligned}$$



## Phrase-based Translation models

- Sequence of words are the units.
- organize english source words into phrases  
 $E = \bar{e}_1, \bar{e}_2, \dots, \bar{e}_I$
- translate each english phrase to corresponding french phrase  $\bar{f}_i$  if  $F = f_1, f_2, \dots, f_I$
- reorder  $f_i$ 's if required.

$\phi(\bar{f}_i | \bar{e}_i)$ : translation probability of generating  $\bar{f}_i$  from  $\bar{e}_i$

$a_i$ : start of foreign phrase generated by  $i$ th english phrase  $\bar{e}_i$

$b_{i-1}$ : end pos<sup>n</sup> of foreign phrase generated by  $i-1$ th english phrase  $\bar{e}_{i-1}$

$d(a_i, b_{i-1})$ : distortion probability  $\propto |a_i - b_{i-1}|$

$$P(F|E) = \prod_{i=1}^I \phi(\bar{f}_i | \bar{e}_i) d(a_i, b_{i-1})$$

$e$  many did not slap the green words  
 $f$  manyne kari bhutni ko nahi mara

$$P(f|e) = P(\text{manyne} | \text{many}) d(0-0)$$

$$P(\text{nahi} | \text{did not}) d(3-1)$$

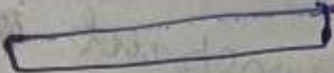
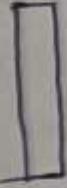
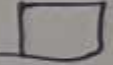



$$P(\text{mara} | \text{slap}) d(4-3)$$

$$P(\text{kari bhutni ko} | \text{the green words}) d(2-4)$$

— parameters  $\phi(\bar{f}_i, \bar{e}_i)$

— to ~~learn~~ find  $\phi$  as well as  $\alpha$ , we need a  
 parallel corpus where <sup>in each pair,</sup> each english phrase is  
 mapped to ~~one~~ foreign language phrase.

$$\phi(\bar{f}_i, \bar{e}_i) = \frac{\phi \text{ count}(\bar{f}_i, \bar{e}_i)}{\sum_{\bar{f}_i} \text{count}(\bar{f}_i, \bar{e}_i)}$$

many did not slap the green witch.  
many ne hari batri ko nahi mara.  
many ne hari batri ko nahi mara.  
many   
did   
not   
slap   
the   
green   
witch.



## Extracting aligned pair of phrases

- use two different aligners - one from F to E, the other from E to F
  - take intersection of the two alignments to get high precision aligned words
  - take the union of the alignments to get low precision aligned words
  - use a classifier to incrementally update union and add back to alignment. (Koehn 2003)
- Select words from minimal intersection



## Decoding for phrase-based MT

$$\hat{E} = \underset{E}{\operatorname{argmax}} P(E|F) R(E)$$

- finding the english sentence that maximises the translation & lang. model probabilities is a search problem — Decoding.

MT decoders are a special case of A\* search.

- stack decoding. implemented using a priority queue.

Function StackDecoding(source sentence) returns (target sentence)

init - stack contains null hypothesis

do.

pop the best hypothesis off the stack

if  $h$  is complete, return  $h$

for each possible expansion  $h'$  of  $h$

assign a score to  $h'$

push  $h'$  onto the stack

## Decoding for phrase-based MT

$$\hat{E} = \underset{E}{\operatorname{argmax}} P(E|F) P(E)$$

- finding the english sentence that maximises the translation & lang. model probabilities is a search problem - Decoding.
- MT decoders are a special case of A\* search.
- stack decoding. implemented using a priority queue.

Function StackDecoding(source sentence) returns (target sentence)

init - stack contains null hypothesis

do.

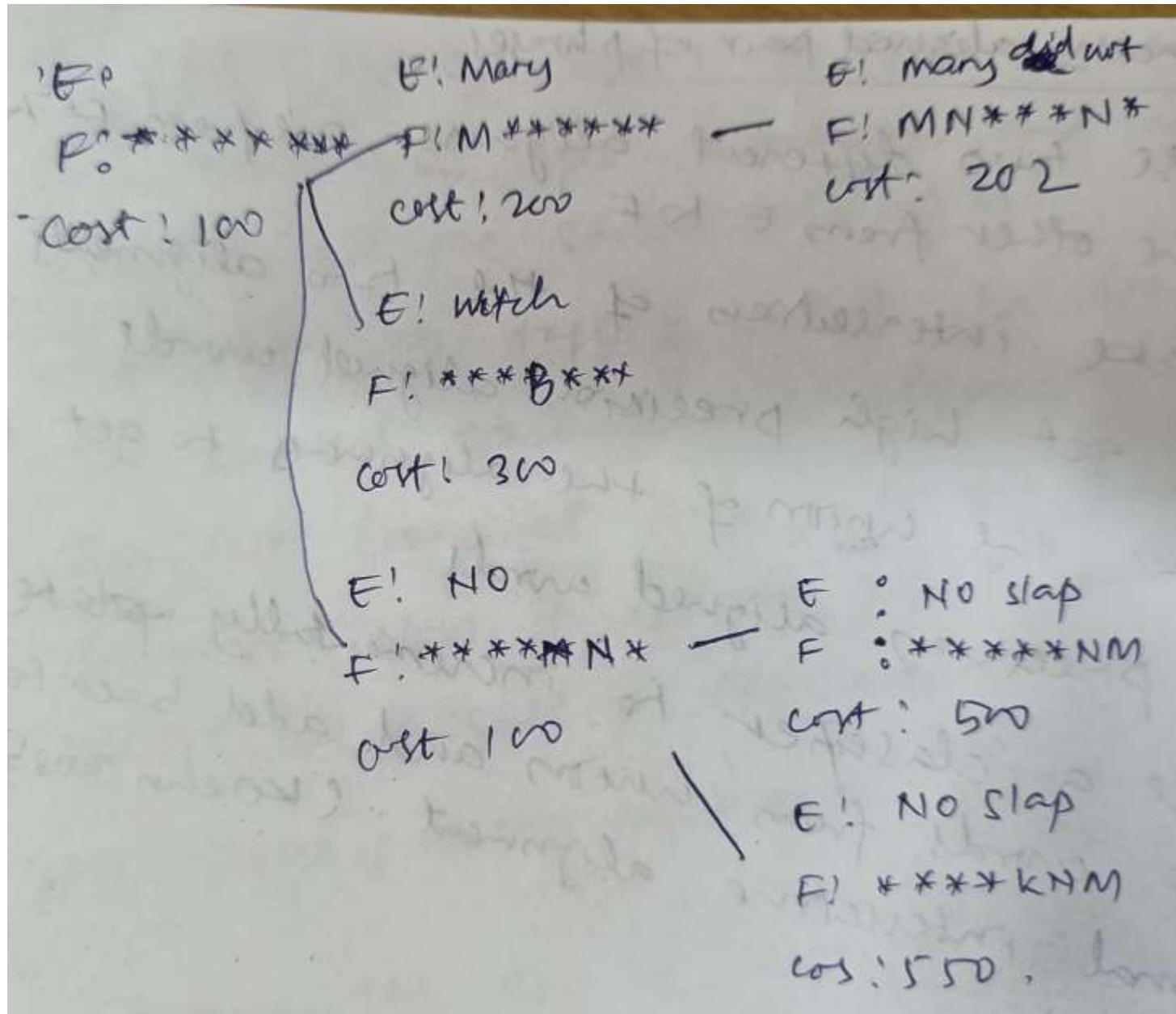
pop the best hypothesis off the stack

if  $h$  is complete, return  $h$

for each possible expansion  $h'$  of  $h$

align a score to  $h'$

push  $h'$  onto the stack





For a set of partially translated phrases

$$S = (E|F)$$

$$\text{cost}(E, F) = \prod_{i \in S} \phi(\bar{f}_i, \bar{e}_i) d(a_i - b_{i-1}) P(F)$$

By combining the current cost with future cost,  
for current F for remaining phrases F

the state cost gives an estimate of total probability of search path for eventual complete translation sentence passing through the current node

for future cost, distribution prob. is ignored for simplicity.

beam search pruning - keep only the  $K$  promising states at each level.



# Evaluation Metrics

# Accuracy

		Predicted	
		T	F
Actual	T	TP	FN
	F	FP	TN

$$Acc = \frac{1}{N} \sum_{i=0}^N \delta(\frac{y_i}{y} = \hat{y})$$

class imbalance

# Precision, Recall, F1 measure

		Predicted	
		T	F
Actual	T	TP	FN
	F	FP	TN

$$\text{Precision } \hat{p} \rightarrow \frac{TP}{TP + FP}$$

$$\text{Recall } \hat{r} = \frac{TP}{TP + FN}$$

$$F_1\text{-measure} \rightarrow \frac{2\hat{r}\hat{p}}{\hat{r} + \hat{p}}$$

# Receiver operating characteristic (ROC) curve

		Predicted	
		T	F
Actual	T	TP	FN
	F	FP	TN

X axis      False positive rate       $\frac{FP}{FP+TN}$

Y axis      True positive rate       $\frac{TP}{TP+FN}$

ROC can be summarised using AUC.



# Receiver operating characteristic (ROC) curve

		Predicted	
		T	F
Actual	T	TP	FN
	F	FP	TN

x axis False positive rate  $\frac{FP}{FP+TN}$   
y axis True positive rate  $\frac{TP}{TP+FN}$

ROC can be summarised using AUC.

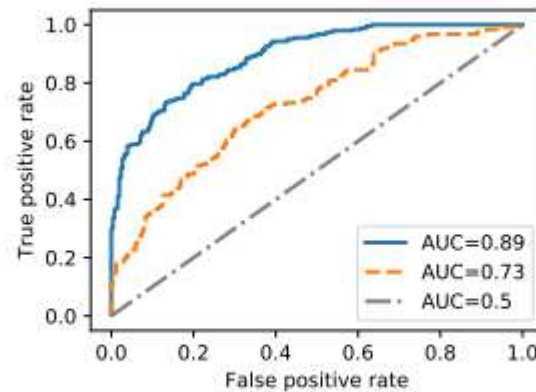


Figure 4.4: ROC curves for three classifiers of varying discriminative power, measured by AUC (area under the curve)

# Evaluation Metrics for MT: BLEU

$$- P_n (\text{modified } n\text{-gram precision}) = \frac{\# \text{ } n\text{-grams in both reference \& hypothesis translation}}{\# \text{ } n\text{-grams in the hypothesis translation}}$$

$$\text{BLEU} = \exp\left(\frac{1}{N} \sum_{n=1}^N \log p_n\right) * \text{BP}, \quad \text{BP} = \begin{cases} 1 & \text{if } c > r \\ e^{1 - \frac{r}{c}} & \text{if } c \leq r \end{cases}$$

$c$  = total length of ~~reference~~ hypothesis

$r$  = effective reference length.

# Evaluation Metrics for MT: BLEU

$$P_n (\text{modified } n\text{-gram precision}) = \frac{\# \text{ } n\text{-grams in both reference hypothesis translation}}{\# \text{ } n\text{-grams in the hypothesis translation}}$$

$$BLEU = \exp\left(\frac{1}{N} \sum_{n=1}^N \log p_n\right) * BP, \quad BP = \begin{cases} 1 & \text{if } c > r \\ e^{\frac{1-r}{2}} & \text{if } c \leq r \end{cases}$$

$c$  = total length of ~~reference~~ hypothesis

$r$  = effective reference length.

Reference - Vinay likes programming in python

Hypothesis 1 TO vinay it like to program python

2 Vinay likes python

3 vinay likes programming in <sup>his</sup> pyjamas

	$e^{\frac{1-5}{3}} = 0.51$					
	$P_1$	$P_2$	$P_3$	$P_4$	$BP$	$BLEU$
1	$\frac{2}{7}$	0	0	0	1	0.21
2	$\frac{3}{3}$	$\frac{1}{2}$	0	0	0.51	0.33
3	$\frac{4}{6}$	$\frac{3}{5}$	$\frac{2}{4}$	$\frac{1}{3}$	1	0.76

# Evaluation Metrics for MT: chrF

- character F-score

**chrP** percentage of character 1-grams, 2-grams, ..., k-grams in the hypothesis that occur in the reference, averaged.

**chrR** of character 1-grams, 2-grams, ..., k-grams in the reference that occur in the hypothesis, averaged.

The metric then computes an F-score by combining chrP and chrR using a weighting parameter  $\beta$ . It is common to set  $\beta = 2$ , thus weighing recall twice as much as precision:

$$\text{chrF}\beta = (1 + \beta^2) \frac{\text{chrP} \cdot \text{chrR}}{\beta^2 \cdot \text{chrP} + \text{chrR}} \quad (10.24)$$

For  $\beta = 2$ , that would be:

$$\text{chrF2} = \frac{5 \cdot \text{chrP} \cdot \text{chrR}}{4 \cdot \text{chrP} + \text{chrR}}$$



# Evaluation Metrics for MT: chrF

- character F-score

For example, consider two hypotheses that we'd like to score against the reference translation *witness for the past*. Here are the hypotheses along with chrF values computed using parameters  $k = \beta = 2$  (in real examples,  $k$  would be a higher number like 6):

REF: witness for the past,	
HYP1: witness of the past,	chrF <sub>2,2</sub> = .86
HYP2: past witness	chrF <sub>2,2</sub> = .62

Let's see how we computed that chrF value for HYP1 (we'll leave the computation of the chrF value for HYP2 as an exercise for the reader). First, chrF ignores spaces, so we'll remove them from both the reference and hypothesis:

REF: witnessforthepast,	(18 unigrams, 17 bigrams)
HYP1: witnessofthepast,	(17 unigrams, 16 bigrams)

# Evaluation Metrics for MT: chrF

- character F-score

Next let's see how many unigrams and bigrams match between the reference and hypothesis:

unigrams that match: w i t n e s s f o t h e p a s t , (17 unigrams)

bigrams that match: wi it tn ne es ss th he ep pa as st t, (13 bigrams)

We use that to compute the unigram and bigram precisions and recalls:

unigram P:  $17/17 = 1$       unigram R:  $17/18 = .944$

bigram P:  $13/16 = .813$     bigram R:  $13/17 = .765$

Finally we average to get chrP and chrR, and compute the F-score:

$$\text{chrP} = (17/17 + 13/16)/2 = .906$$

$$\text{chrR} = (17/18 + 13/17)/2 = .855$$

$$\text{chrF}_{2,2} = 5 \frac{\text{chrP} * \text{chrR}}{4\text{chrP} + \text{chrR}} = .86$$

# Evaluation Metrics for MT

Translations are evaluated along two dimensions:

1. **adequacy**: how well the translation captures the exact meaning of the source sentence. Sometimes called **faithfulness** or **fidelity**.
2. **fluency**: how fluent the translation is in the target language (is it grammatical, clear, readable, natural).