

## MODEL 1 :

### HMM AND VITERBI ALGORITHM WITHOUT MORPH ANALYSIS

#### Experiment -

The utf file was converted to wx notation using the convertor.

The output file was used to train the model.

The train data contains 1000 sentences.

The test data contains 100 sentences.

#### Observations:

Correction of training data POS

- XC tag was not used correctly.
- Proper Nouns were tagged as Nouns
- VAUX were tagged as VM.

#### Emission Matrix -

Contains probability of a word(X) being tagged as Y.

$$P(Y | X) = \text{Count}(X, Y) / \text{Count}(X).$$

These values are computed using the training data. It is a 2-D matrix.

Rows represent unique words and columns represent Tags.

#### Transition Matrix -

Contains probability of Tag1(T1) being followed by Tag2(T2).

$$P(T2 | T1) = \text{Count}(T1, T2) / \text{Count}(T1).$$

These values are computed using the training data. It is a 2-D matrix.

Rows represent T1 tags or previous tags and columns represent T2 tags or current tags.

#### Maximum Likelihood of a Tag sequence -

Probability of tag sequence = prob of the word being tagged as T2 \* probability of T2 following T1

This is computed for all words and the sequence maximising the probability is the best tag sequence.

#### Viterbi -

implemented using dynamic programming

#### HMM -

assumed to be dependent on the previous state only

#### Assumptions:

Handling of unseen words – unseen words are assumed to be nouns.

#### Results -

The unseen words were initially tagged as Nouns but the in best tag sequence result get converted to other tags appropriately

## MODEL 2

### HMM AND VITERBI ALGORITHM WITH NOUN MORPH ANALYSIS

#### Experiment -

The utf file was converted to wx notation using the convertor.

The output file was used to train the model.

The train data contains 1000 sentences.

The test data contains 100 sentences.

#### Observations:

Correction of training data POS

- XC tag was not used correctly.
- Proper Nouns were tagged as Nouns
- VAUX were tagged as VM.

Creation of paradigm tables - 6 categories ( Masculine – 3 , feminine – 3) were observed.

Within each category, Direct and oblique cases were observed for singular and plural forms.

#### Emission Matrix -

Contains probability of a word(X) being tagged as Y.

$P(Y | X) = \text{Count}(X,Y) / \text{Count}(X)$ . observed

These values are computed using the training data. It is a 2-D matrix.

Rows represent unique words and columns represent Tags.

#### Transition Matrix -

Contains probability of Tag1(T1) being followed by Tag2(T2).

$P(T2 | T1) = \text{Count}(T1,T2) / \text{Count}(T1)$ .

These values are computed using the training data. It is a 2-D matrix.

Rows represent T1 tags or previous tags and columns represent T2 tags or current tags.

#### Maximum Likelihood of a Tag sequence -

Probability of tag sequence = prob of the word being tagged as T2 \* probability of T2 following T1

This is computed for all words and the sequence maximising the probability is the best tag sequence.

#### Viterbi -

implemented using dynamic programming

#### HMM -

assumed to be dependent on the previous state only

#### SFST -

6 .lex and .fst files were created on the basis of hindi paradigm rules.

### Feminine Nouns -

	Type 1		Type 2		Type 3	
	Direct	Oblique	Direct	Oblique	Direct	Oblique
Singular	<i>null</i>	<i>null</i>	<i>null</i>	<i>null</i>	<i>null</i>	<i>null</i>
Examples	<i>āg, pyās</i>	<i>āg, pyās</i>	<i>nədī, mənī</i>	<i>nədī, mənī</i>	<i>lātā, bāt</i>	<i>lātā, bāt</i>
Plural	<i>null</i>	<i>null</i>	<i>-yā</i>	<i>-yō</i>	<i>-ē</i>	<i>-ō</i>
Examples	<i>āg, pyās</i>	<i>āg, pyās</i>	<i>nədī-yā, mənī-yā</i>	<i>nədī-yō, mənī-yō</i>	<i>lātā-ē, bāt-ē</i>	<i>lātā-ō, bāt-ō</i>

### Masculine Nouns -

	Type 1		Type 2		Type 3	
	Direct	Oblique	Direct	Oblique	Direct	Oblique
Singular	<i>null</i>	<i>null</i>	<i>null</i>	<i>-e</i>	<i>null</i>	<i>null</i>
Examples	<i>krodh, pyār</i>	<i>krodh, pyār</i>	<i>larkā, sāyā</i>	<i>lark-e, sāy-e</i>	<i>ādmī, ghār</i>	<i>ādmī, ghār</i>
Plural	<i>null</i>	<i>null</i>	<i>-e</i>	<i>-ō</i>	<i>null</i>	<i>-ō/-yō</i>
Examples	<i>krodh, pyār</i>	<i>krodh, pyār</i>	<i>lark-e, sāy-e</i>	<i>lark-ō, sāy-ō</i>	<i>ādmī, ghār</i>	<i>ādmī-yō, ghār-ō</i>

Manually nouns were categorised into these 6 categories.

In the .fst file -

The states in the Finite State Transducers were identified on the basis of the suffixes of the input word.

Alphabet were defined according to occurring inflections.

Delete rules were formed according to each sub-cases within all the 6 categories

Also, the modifications required for inflections were taken care of – separately for all the unique occurrences.

Assumptions:

Handling of unseen words – unseen words are assumed to be nouns.

Morph analysis is done only for nouns.

Results -

The unseen words were initially tagged as Nouns but the in best tag sequence result get covered to other tags appropriately