

07.09.21
Thursday

Statistical NLP

(Text \rightarrow numbers)

1) TF-IDF (Term freq & Inverse Document-freq)

$$TF = \frac{\text{No. of terms in the document}}{\text{Total no. of terms in the document}} = \left(\frac{n}{N} \right)$$

$$IDF = \log \left(\frac{\text{number of documents}}{\text{no. of documents where term has appeared}} \right) = \left(\frac{D}{F} \right)$$

$$\therefore \text{now } [TF-IDF = TF * IDF]$$

D_1 : I'm reading a newspaper. I'm also reading a book.

D_2 : I'm reading a magazine. It has events

$D=2$, term = reading, $N=$

TFIDF (reading)?

TF(reading, D_1) =

$$TF(\text{reading}, D_1) = \frac{2}{10}$$

$$IDF(\text{reading}) = \log\left(\frac{2}{2}\right) = 0$$

$$TF(\text{newspaper}, D_1) = \frac{1}{10}$$

$$IDF(\text{newspaper}) = \log\left(\frac{2}{1}\right) = TF \cdot IDF = \left(\frac{1}{10}\right) 10.3$$

note:

if a term present in all doc then it's not relevant or important in one content.

Feature Extraction

Tw1: I'm happy because I'm learning NLP

Tw2: I'm happy

Positive

Tw2: I'm sad, I'm not learning NLP

Tw2: I'm sad

negative

row 3
row 3
⋮

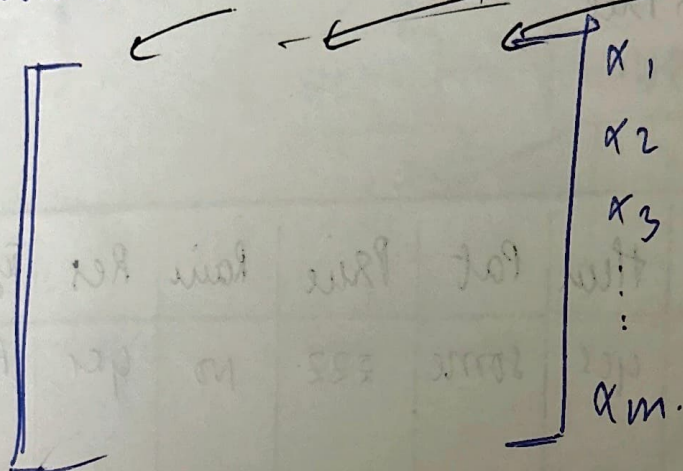
→ x_3

⋮

row m
row m

→ x_m

* First column is Bias, Positive, Negative.



$(x^m \times 3)$ is the dimension of the matrix

- 1) Bag of words
- 2) word to vector
 - Bow.
 - skipgram.
- 3) Glove
- 4) BERT
- 5) LLM based

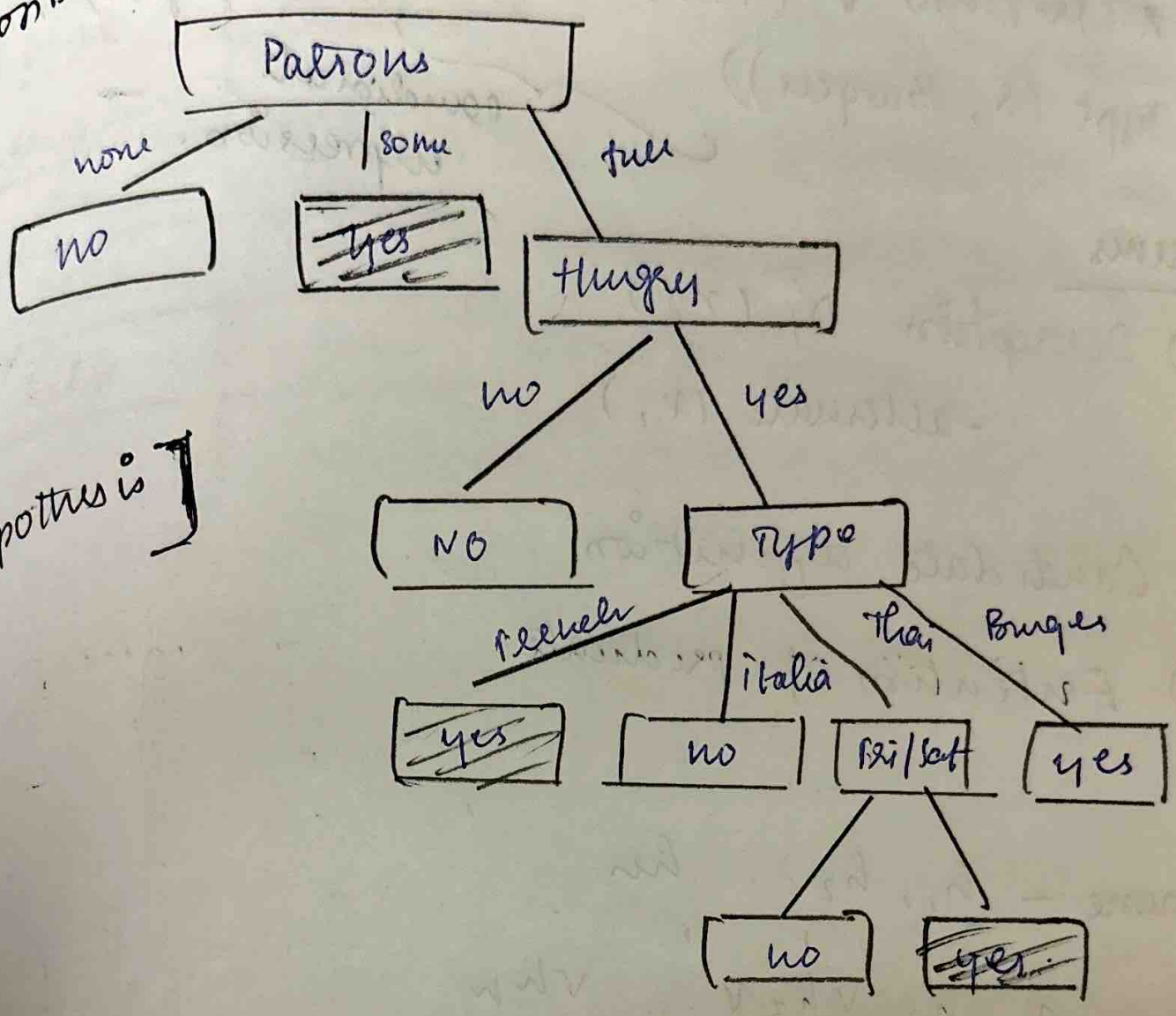
2007
2008

18 Decision Tree

[illegible]

decision tree

[hypothesis is]



Hypothesis as logical sentence

- 1) alternates $(x,)$ \wedge T bar $(x,)$ \wedge T Fri/Pat $(x,)$
- 2) classification - will wait $(x,)$ \vee T will wait $(x,)$

Hypothesis

$$\forall n \text{ goal}(n) \iff G_i(n) \rightarrow \text{candidate expression.}$$

$$\forall x \text{ will wait}(x) \iff \text{Patrons}(x, \text{some}) \vee (\text{Patrons}(x, \text{full}) \wedge \text{Hungry}(x) \wedge \text{Type}(x, \text{french}) \vee (\text{Patrons}(x, \text{full}) \wedge \text{Hungry}(x) \wedge \text{Type}(x, \text{Thai})$$

\wedge (fri/sat) \vee (Patron (x, full) \wedge Hungry (x)
Type (x, Burger))

candidate expression.

Terms

- (1) Description $D_i (n_i)$
- alternate (x_i)
- (2) Candidate definition
- (3) Entention of predicate.

Assume - $h_1, h_2 \dots h_n$

$n_1 \vee h_2 \vee h_3 \vee \dots \vee h_n$

→ logically consistent

→ logically inconsistent

FN

FP

Actual	Predicted	
True	True	→ True positive
False	False	→ True negative
True	False	→ False negative
False	True	→ False positive