NLP - NLTK

Ch 6: Learning to classify text
  ↳ Supervised

# gender identific⁀ from name - unethical and wrong

# feature extrac⁀:

Indian parliament → finance doc        } binary classific⁀
              ↳ non finance doc

    feature vector → | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
                        ↳ if the word is present

    ↳ Bag of words model

    ↳ Although, position, structure, grammar have meaning
      But, many tasks can be done w/o this → only Bag of Words enough

# Document classific⁀

# Exploiting content
      ↳ Include position content

# Sequence classific⁀
      Named entity tagger

                                    The
                                    famous
                                    building
                                    Gateway  → start
                                    of
                                    India  → stop

HMM → 0  0 0 0

RNN → 0   0 0

         prob.

LSTM → long sent, influenced by things long behind it

coreference resolution

$$\underset{1}{\text{The famous poet}} \text{ from } \underset{2}{\text{India}},$$

wrote $\underset{3}{\underline{XYZ}}$  In that $\underset{2}{\underline{country}}$, $\underline{he}$,

$\underset{1}{\underline{Kalidasa}}$ was born in year ABC

3 clusters:

# Evaluation
  ⤷ for binary
     classifier
       Φ



predicted

|        |     | C   | N    |
|--------|-----|-----|------|
| actual | C 5 | 0   | 95 → false negative ②|
|        | N 95| 5   | 95   |

① false positive

true positive

true negative

Precision - $\dfrac{TP}{TP+FP} = \dfrac{0}{5} = \dfrac{\text{correct D}}{\text{total D}}$

Recall - $\dfrac{TP}{TP+FN} = \dfrac{\text{correct D}}{\text{total that should be detected}}$

  ⤷ how many relevant were identified

F score = $\dfrac{2 \times \text{prec} \times \text{recall}}{\text{prec} + \text{recall}}$

MNIST data

confusion matrix

                 0   1   2   3   4   5   6   7   8   9

0
1
2
3
4
5
6
7
8
9

diagonal → heavy.

---

Ch 8 : Analyzing Sentence structure

Grammar parser

↳ content free ~~parser~~ assumpⁿ - no subtree in a grammar tree is connected to other subtrees

\# words are like cloud - topic modelling
    5 topics ↝ topicwise segregation
    ↳ complete statistic model

Mallet