## Chapter 3

# Examples of word Prediction

Google N gram Viewer

_What_ are you? is high et.

# Why probability?

→ Provides methods to predict or make decision to pick the next word in the sequence based on a sampled data.

→ Make the informed depcision when there a certain degree of uncertainity of some observed data.

→ It provides a quantitative description of the chances or likelinood associated with various outcomes.

— Probability of a sentence

→ probability of next word in a sentence

⇒ Probablity :- The probablity is defined as the likelinood that an event will occur

# Discrete sample space

→ Experiment :- Extracting tokens from a document

→ Outcome :- Every token/word $x$ in the document

→ The outcome of the experiment — 52 sample (words). They constitute the sample space, $\Omega$ or the set of all possible outcomes.

Each word in this sample belongs to $\Omega$ represented by $x \in \Omega$

Each sample $x \in \Omega$ is assigned a probability score [0, 1].

A probability function or probability distribution function distributes the probability mass of 1 to the all the samples in the sample space $\Omega$.

$\Rightarrow$ Sample Space - constraints

All the words in the $\Omega$ must satisfy the following constraints:

1. $P(x) \in [0, 1]$ $\forall x \in \Omega$ &

2. $\sum_{x \in \Omega} P(x) = 1$

Example:-

If we are equally likely to pick any word from the BOW, then the probability for any word is

Bag of words Count = 52

$P(x) = 1/52$ $\forall x \in \Omega$ so that

$P(\Omega) = 1$

$P(\text{'weather'}) = 1/52 = 0.0192307692 3$

$\Rightarrow$ EVENTS

Example:-

Total no. of words = 52

The no. of unique words = 37 or there are 37 types of words have frequencies $\geq 1$.

An event is an collection of samples of the same type., $E \subseteq \Omega$

$$P(E) = \sum_{x \in E} P(x)$$

Events can be described as a variable taking a certain value

In the Bow, the word type "the" occurs 6 times. Then

$$E_{the} = 6$$

$$P(E_{the}) = 6 \times \frac{1}{52} = 0.115$$

In the Bow, the word type pack occurs 3 times. Then

$$E_{pack} = 3$$

$$P(E_{pack}) = 3 \times \frac{1}{52} = 0.058$$

# Random variable

A random variable, is a variable whose possible values are numerical outcomes of a random phenomenon

↳ Two types - continuous and discrete for NLP, they are discrete.

To capture the type - token distinction, we random variable w. $W(x)$ maps the sample $x \in \Omega$.

v is the set of types & the value is represented by a variable $v$.

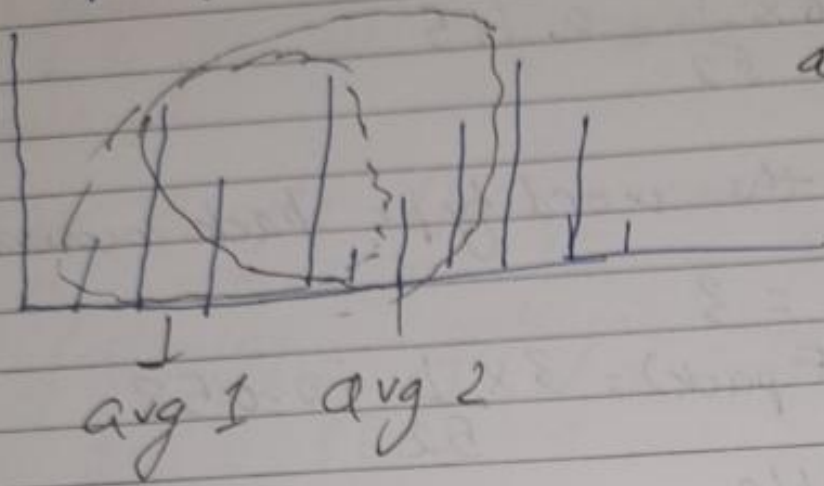Given a random variable $V$ & a value $v$, $P(V = v)$ is the probablity of the event that $V$ takes the value $v$. i.e

$$P(V = v) = P(x \in \Omega : V(x) = v)$$

$$P(V = 'the') = P('the') = 0.115$$

Random Variables are useful in describing, constructing various events.

# Project – Stock Price Prediction

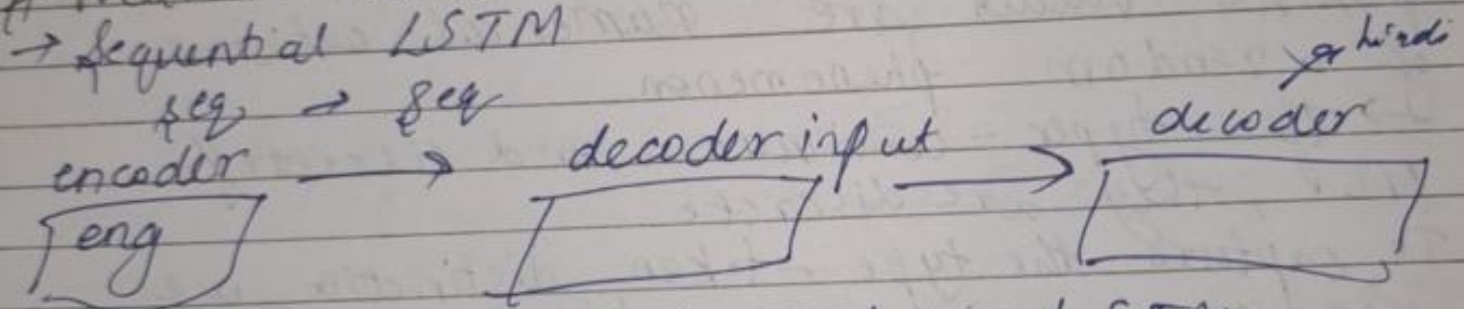→ rolling avg – helpful for traders



avg moves fwd
1 at a
time

avg 1   avg 2

# Neural Machine Translation → Hindi → English

→ Sequential LSTM

seq → seq

encoder ⟶    decoder input ⟶ decoder (hindi)

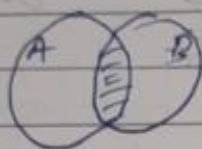[ eng ]       [　　　]       [　　　]

we will use sequential model – LSTM

# Joint Probablity :-

Given any two events $E_1$ & $E_2$, the probablity of their conjuction

$$P(E_1 \cup E_2) = P(E_1 \cap E_2)$$ is called the joint probablity. of $E_1$ & $E_2$ occurs simultaoneously.

Example :- The probablity of the first letter of 't' and the second letter 'h' is $P(F='t', S='h')$. The joint probablity should be as large as the probability of $P$(the')

$P(A) = $ size of $A$ relative to $\Omega$
$P(A,B) = $ size of $A \cap B$ relative to $\Omega$

# Conditional Probablity

When we have partial knowledge influencing the outcome of an expirement, we use it to update the outcome.

The conditional probablity $P(E_2 / E_1)$ is the probablity of event $E_2$ given that event $E_1$ has occured. $P(E_2 / E_1)$ is defined as :

$$P\left(\frac{E_2}{E_1}\right) = \frac{P(E_1 \cap E_2)}{P(E_1)} \quad \text{if } P(E_1) > 0$$

$\Rightarrow$ example - Conditional Probablity - Bigram

let us consider a corpus of Kinematics problems in physics that contains about 280+ problems.

→ Bigram Sample Space - $\{w_1, w_2\} \in \Omega \cong 37A$

→ A $(\{w_1, w_2\}) = \{$ average, * $\}$ - bigram starting with avg

→ B $(w_1, w_2) = \{$ * , speed $\}$ - bigram ending with speed

P(average) = 0.036

P(speed) = 0.114

P(average, speed) = 0.004

P(speed / average) = $\dfrac{0.004}{0.036}$ = 0.111

P(avg / speed) = $\dfrac{0.004}{0.114}$ = 0.035

# Independance

Two events are dependent if the probability of on relies on occurrence of the other; if there is not much interaction; then the events are independent.

Two events $E_1$ & $E_2$ are independent if & only if $P(E_1, E_2) = P(E_1) P(E_2)$

OR

→ $P(E_1) = P(E_1 / E_2)$   $P(E_2) = P(E_2 / E_1)$

Example :-

P(average) = 0.036

P(speed) = 0.114

P(average, speed) = 0.004

The bigram $\{$ average, speed $\}$ did not happen by chance. The words average, speed are NOT independent

# The Language Model
→ Natural language sentences can be described as parse trees which use the morphology of words, syntax & semantic.
→ Probablistic thinking - finding how likely a sentence occurs or formed, given the word sequence
→ In probablistic world, the language model is used to assign a probability $p(w)$ to every possible word sequence $w$.

Application →

Speech Recognition                    Did you hear Recognize
                                      speech or wreck a
                                      nice beach.

Content sensitive                     Once upon a tie.
     spelling                         Their lived asking
Machine Translation                   a+t work is good →
                                      I'oeurve est bonne
Sentence Completion                   Complete a sentence
                                      as the previous word
                                      is given
OCR & Hand writing                    The quick brown
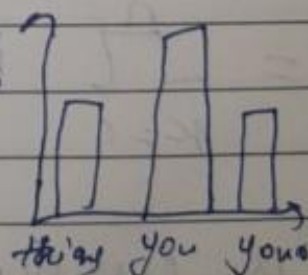Recognition                                        fox

                    predict the next word

How are  →  Language Model  →

   ↑                  ↑

Input Sentence    knowledge about the
                  language - grammar
                    sentence structure,
                      domain etc

                                  the  you  your

# Probablistic Language Model

Goal :- Compute the probablity of a sequence of words

$$P(w) = P(w_1, w_2, w_3 \dots w_n)$$

Task :- To predict the next word using probability. Given the context, find the next word using

$$P(w_n / w_1, w_2, w_3 \dots w_{n-1})$$

A model which computes the probablity for (5) using (6) is called as Probablistic Language Model.

The probability of $P(\text{The cat roars})$ is less likely to happen than $P(\text{The cat meows})$

# Chain Rule

It is difficult to compute the probability of the entire sequence $P(w_1, w_2, w_3, \dots w_n)$?

Chain rule is used to decompose the joint probablity of a sequence into a product of conditional probablity.

$$P(w) = P(w_1, w_2, w_3 \dots, w_n) = P(w_1^n)$$
$$= P(w_1) \, P(w_2 / w_1) \, P(w_3 / w_2, w_1) \dots$$
$$= \prod_{k=1}^{n} P(w_k / w_1^{k-1})$$

It is possible to $P(w/h)$ but it doesn't really help in reducing the computational complexity

→ We use innovative ways to string words to form new sentences.

→ Finding the probability for a long sentence may not yield good outcome as the content may never occur in the corpus.

→ Short sequences may provide better results.

# Markov Assumption

The future behaviour of a dynamic system depends on its recent & not on the entire history.

The product of the conditional probablities can be written approximately for a bigram as

$$P\left(\omega_k \mid \omega_1^{k-1}\right) \approx P(\omega_k \mid \omega_{k-1}) \quad - \text{⑩}$$

equation ⑩ can be generalized for an n-gram as

$$P(\omega_k \mid \omega_1^{k-1}) \approx P\left(\omega_k \mid \omega_{k-K+1}^{k-1}\right)$$

Now, the joint probability of a sequence can be re-written as

$$P(\omega) = P(\omega_1, \omega_2, \omega_3 \ldots \omega_n) = P(\omega_1^n)$$
$$= P(\omega_1) P\left(\omega_2/\omega_1\right) P\left(\omega_3 \mid \omega_2 \omega_1\right) \ldots$$
$$= \prod_{k=1}^{n} \left(\omega_k \mid \omega_1^{k-1}\right)$$
$$\approx \prod_{k=1}^{n} P\left(\omega_k \mid \omega_{k-K+1}^{k-1}\right)$$
$$\downarrow$$
$$\text{n-gram}$$

# Generative Models

⇒ Target & Context words

Next words in the sentence depends on its immediate past words, known as context words

$$P(W_{k+1} \mid W_{i-k}, W_{i-k+1} \dots W_k)$$

<u>content words</u>

n-grams
unigram      —    $P(W_{k+1})$
bigram      —    $P(W_{k+1} \mid W_k)$
trigram     —    $P(W_{k+1} \mid W_{k-1}, W_k)$
4-gram     —    $P(W_{k+1} \mid W_{k-2}, W_{k-1}, W_k)$

# Language Modelling using Unigrams

→ A unigram language model all words are generated independently $W_1, W_2, W_3 \dots W_n$ and none of them depend on the other.

→ This is not a good model for language generation

→ It may generate the the the the as a sentence

→ cannot string words according to high probability

→ Generates a document containing N words using n-gram.

→ A good model assigns higher probability to the word that actually occurs

→ $$\sum_{i=1}^{N} P(W_i) = 1 \quad W_i \text{ to be estimated}$$

in this model is $P(w_i)$ & it must satisfy this

$$P(w) = \prod_{i=1}^{N} P(w_i)$$

→ The location of the word is the document is not important
- $P(N)$ is the distribution over $N$ & if same for all documents.

# Maximum Likelihood Estimate
→ One of the methods to find the unknown parameter(s) is the use of Maximum Likelihood Estimate.
→ Estimate the parameter value for which the observed data has the highest probability.
→ Training data may not have all the words in a vocabulary.
→ If a sentence with an unknown word is presented, then the MLE is 0.
→ Add a smoothing parameter to the equation without affecting the overall probability requirements

$$P(w) = \frac{Cw_i + \alpha}{Cw + \alpha/|w|}$$

# Bigram Language Model
A bigram language model generates a sequence one word at a time, starting with the first word & then generating each succeeding word conditioned on the previous one.

→ A bigram model is defined as follows:-

$$P(w) = \prod_{i=1}^{n+1} P(w_i / w_{i-1}),$$

where $w = w_1, w_2 \ldots w_n$

→ Estimate the parameter $P(w_i / w_{i-1})$ for all bigrams.

→ The parameter estimation does not depend on the location of the word.

→ If we consider the sentence as a sequence in time, they are time-invariant MLE picks up the word that is $\dfrac{n_{w,w'}}{n_{w,0}}$

where $n_{w,w'}$ is the number of times the words $w, w'$ occur together & $n_{w,0}$ is the number of times the word $w$ appears in the bigram sequence.

# Probablistic Languacge Model - Example

1] Peter Piper picked a peck of pickled pepper.
2] A peck of pickled peppers Peter Piper picked.
3] If Peter Piper picked a peck of pickled peppers
4] Where's the peck of pickled peppers Peter Piper picked?

| Bigram | Freq |
|---|---|
| <s> Peter | 1 |
| Peter Piper | 4 |
| Piper picked | 4 |
| picked a | 2 |
| a peck | 2 |
| peck of | 4 |

The joint probability of a sentence formed with n words can be expressed as a product of conditional probabilities - we use immediate context & not the entire history.

$$P(w_1 | <s>) \times P(w_2/w_1) \times \ldots P(<E>/w_n)$$

and $$P(w_{i+1}/w_i) = \frac{||w_i \cdot w_{i+1}||}{||w_i||}$$

# Out of Vocabulary Words
→ In a closed vocabulary language model, there is no unknown words or out of vocabulary words (OOV)
→ In an open vocabulary system, you will find new words that are not present in the trained model.
→ Pick words below certain frequency and replace them as OOV.
→ Treat every OOV as a regular word.
→ During testing the new words would be treated as OOV & the corresponding frequency will be used for computation
→ this eliminates zero probability for sentences containing OOV


# Curse of Dimensionality
→ A fundamental problem that makes language and other learning problems difficult is the curse of dimensionality
→ It is particularly obvious in the case when one wants to model the joint

distribution b/w many discrete random variable
→ If one wants to estimate the joint probability
distribution of 10 words in a language
with a million words as vocabulary,
then we need to estimate $10\,000\,000^{10} - 1 =$
$10^{60} - 1$ free parameters.

# Naive Bayes Classification

⇒ Bayes Theorem

Let us consider two random variables X & Y.
Then joint probability. Then joint probability
$P(X=x, Y=y)$ refers to the probability that
the variable X takes the value x and the
variable Y takes the value y. The conditional
probability $(P\&Y=y \mid X=x)$ refers to the probability
that the variable Y takes the value y
given the observation the variable X takes
the value x.

$$P(X,y) = P(y \mid X) \times P(x) = P(x \mid y) \times P(y)$$

$$P(y \mid X) = \frac{P(x \mid y)\, P(y)}{P(X)}$$

# Bayes theorem for Email Classification.
→ Map Baye's theorem using statistical properties
of data
→ Let $X = \{X_1, X_2, X_3 \ldots X_n\}$ where
X is a set of attributes & Y represents a class
The relationship b/w X & y can be found
using the conditional properties $P(y \mid X)$

→ The conditional probability $P(y/x)$ is known as the posterior probability of $y$

→ $P(y)$ is known as the prior probability

→ In the classification problem, it is important to learn the parameters $P(y/x)$. Given the attributes of the email (TF, TF-IDF) find the class to which the email belongs

→ The parameters are obtained from training data. During the training process, we will learn $P(y/x)$ for every word in the corpus.

## # Supervised Classification.

→ Set of input parameters / attributes $X = X_1, X_2 ... X_m$ and a fixed set of classes $Y = y_1, y_2 ... y_n$

→ Every element of the training set $D = d_1, d_2 ... d_n$ is manually assigned a class
$(d_1, y_1) (d_2, y_2) .. (d_n, y_n)$.

→ Goal is to learn the classifier, so that it can map a new document $d$ to any of the classes $y \in Y$

→ Bayes classifier would assign a probability based on the observation to the new document to aid the class selection.

→ The probability score for each class is computed as given by the the equation
$$P(y/x) = \frac{P(x/y) \, P(y)}{P(x)}$$

→ The class will be found using $argmax \, P(y/x)$

$$\hat{y} = \arg\max_{y \in Y} P(y/x)$$

$$= \arg\max_{y \in Y} P(x/y)\, P(y)$$

$$= \arg\max\; P(y)\, P(x_1/y) \times P(x_2/y) \times \dots \; P(x_m/y)$$

$$= \arg\max_{y \in Y} P(y) \prod_{i=1}^{m} P(x_i/y)$$

## TRAINING

1] Prior Probability $\qquad P(y) = \dfrac{Count(y)}{Count(y)} = \dfrac{1}{2}$

2] Learn $\qquad P(x_1/y) = \dfrac{Count(x_1, y)}{Count(y)}$

given the class,
find the probability of
the word in it

for new word, the probability will become 0
so we add 1 to all numerators
dont change denominator. — Smoothing
every corpus is not complete hence
smoothing is important acc to domain.