



Section 10

Lavender Jiang

See references for sources of images

PART 01

Prompt Engineering

Temperature

$$\sum p(x_i) = \frac{\Sigma \cdot \square \dots}{\Sigma} = 1$$

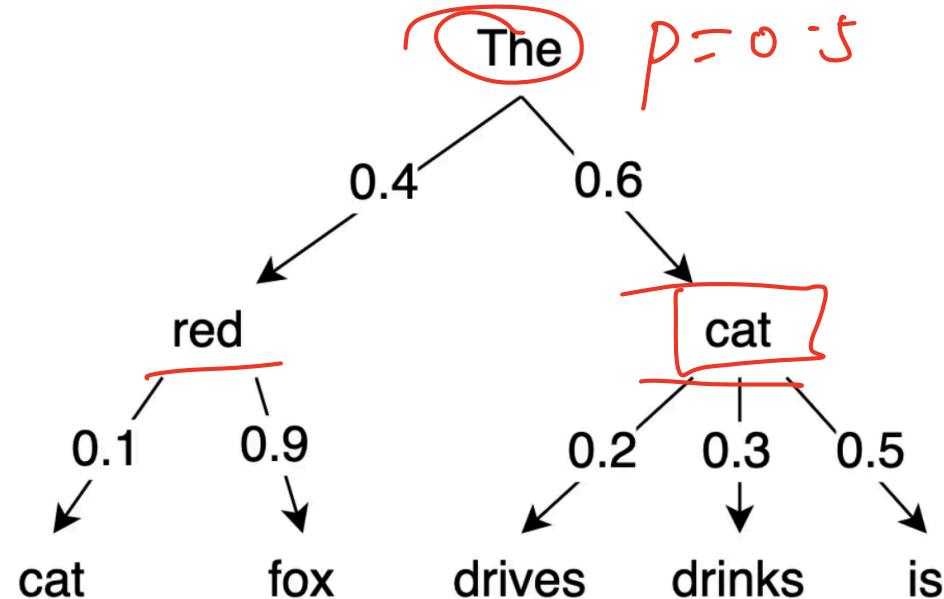
$\prod n^i$ $\rightarrow e^{\frac{x_i}{T}}$

$$p(x_i) = \frac{e^{\frac{x_i}{T}}}{\sum_{j=1}^V e^{\frac{x_j}{T}}}$$

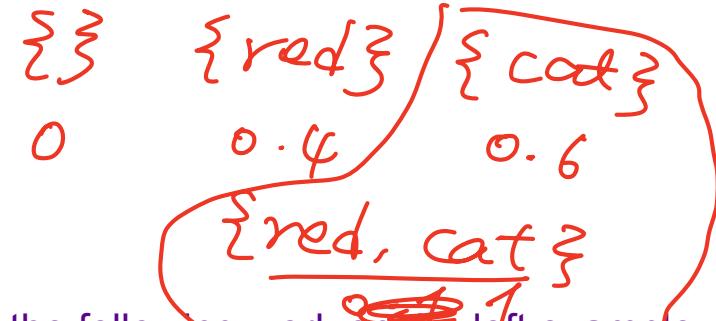
What happens when T approach 0? How about infinity?

$p(x_i) \sim \text{uniform}$

Top-p



$\{\text{red, cat}\}$



How does the following work on the left example:

1. Greedy search
2. beam search with $n_{beam}=2$
3. Top-k sampling
4. Top-p sampling

~~top k possibility~~

0.7

Max Length, Stop Sequences

Predict common sense results of the following actions.
==

- 1 Action: I didn't water the plant for 3 weeks.
Result: The plant died.

==

- 2 Action: I went to school.
Result: I got a diploma.

==

- 3 Action: I left the AC on all day.
Result: I got a high utility bill.
- ==
- 4 Action: I helped my neighbors when their car broke down.
Result: My neighbors were grateful.
- ==

- 5 Action: I put the ice cream outside for an hour.
Result: **The ice cream melted.**
- ==

<bos>
[eos]
<unk>

Frequency Penalty ↘ ↓ repetitive ↑ unique

“Write a poem where every word starts with Z”

Frequency Penalty = 0	Frequency Penalty = 2
<p>Zebras zigzagging zealously, Zephyrs zipping, zesty, zestfully. Zodiac's zenith, zeal's zodiac, Zinnias zigzag, zircon's zodiac. Zeppelin zooming, zigzag</p>	<p>Zealous zephyrs zoom, zigzagging zestily, Zinnia zones, zenith zeppelins' zone precisely. Zenith's ziggurats zealously zap, Zirconium zebras zip-zap on Zanz</p>

Basic Prompt

Prompt

The sky is

Output:

blue

The sky is blue on a clear day. On a cloudy day, the sky may be gray or white.

Prompt:

instruction

Complete the sentence:

The sky is

Output:

so beautiful today.

Prompt Formatting

This is awesome! // Positive

This is bad! // Negative

Wow that movie was rad! // Positive

What a horrible show! // Negative/positive

Q₁
<Question>?
A₁,
<Question>?
A₂
<Question>?
<Answer>
<Question>?

Elements of Prompts

Instruction - a specific task or instruction you want the model to perform

Add 2 7-digit numbers.

Context - external information or additional context that can steer the model to better responses

examples of addition
rules of addition

Input Data - the input or question that we are interested to find a response for

$$\begin{array}{r} "a+b" \\ \hline "100 + 200" \\ 1,0,0 + 2,0,0 \end{array}$$

Output Indicator - the type or format of the output.

"Answer: 300" || take 1st line
split by : int(s)



Homework: 7-digit addition

Instruction - control with your_prompt

Context - control with your_prompt and your_post_processing

Input Data - given by autograder, can format with your_pre_processing

Output Indicator - controlled by your choice of prompt. Can be further processed with your_post_processing

Other knobs: your_config

100 10M

Metrics: accuracy, mean absolute error, prompt length

>0.1

< 10⁶

Autograder is slow. Recommend local testing with test_prompts.py

(Exact same code as autograder, but different random seeds)

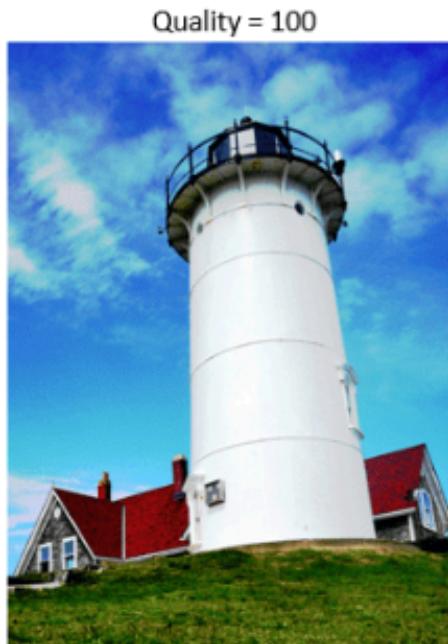
Note: for submission **do not read prompts from files**,
autograder cannot find it

PART 02

Huffman Coding

Compression

JPEG



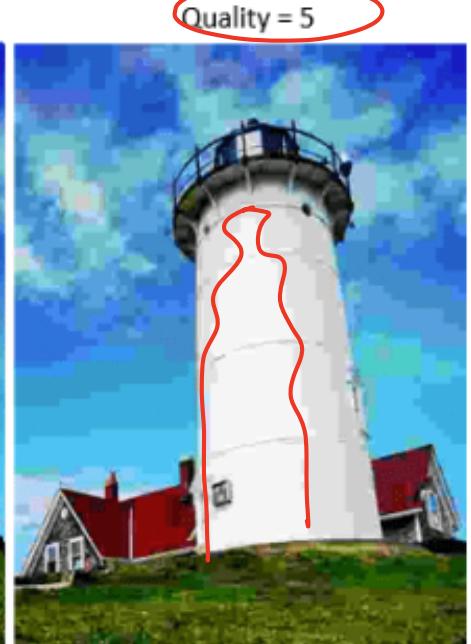
Quality = 100



Quality = 50



Quality = 10



Quality = 5

Less Compression



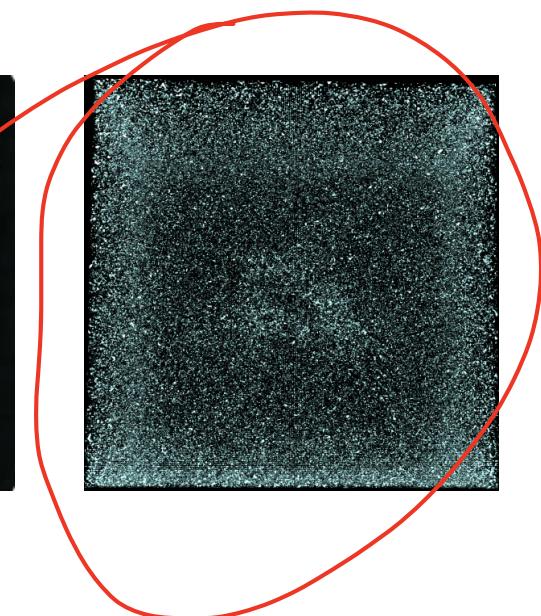
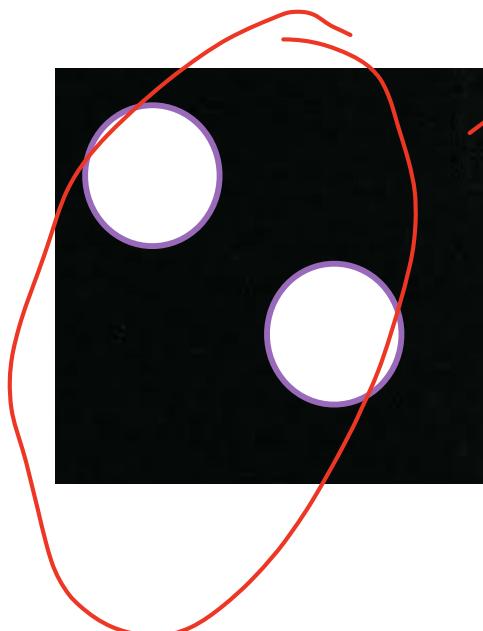
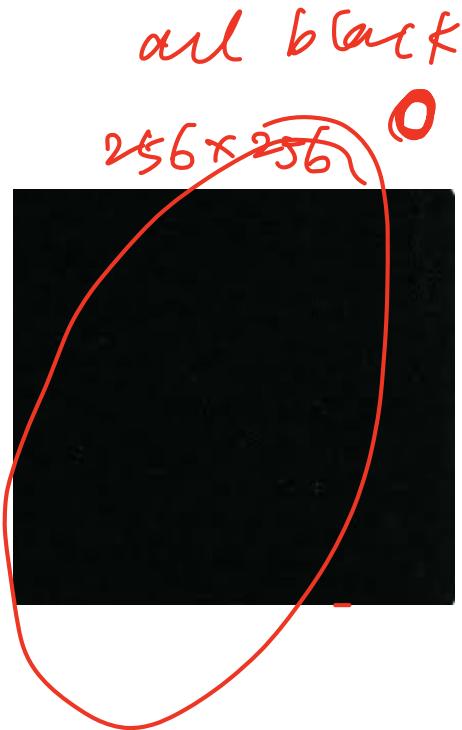
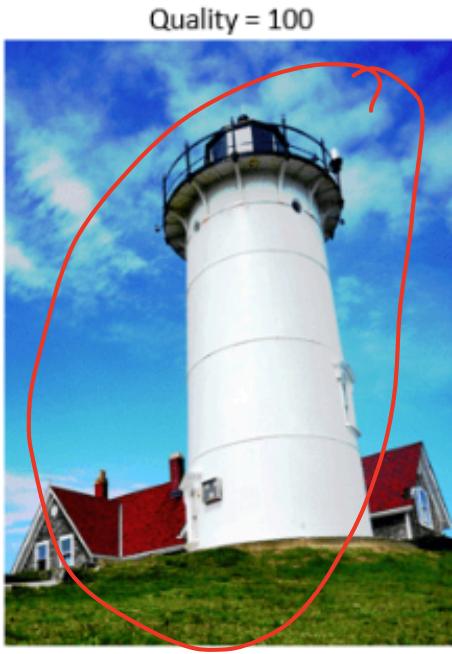
More Compression



NYU

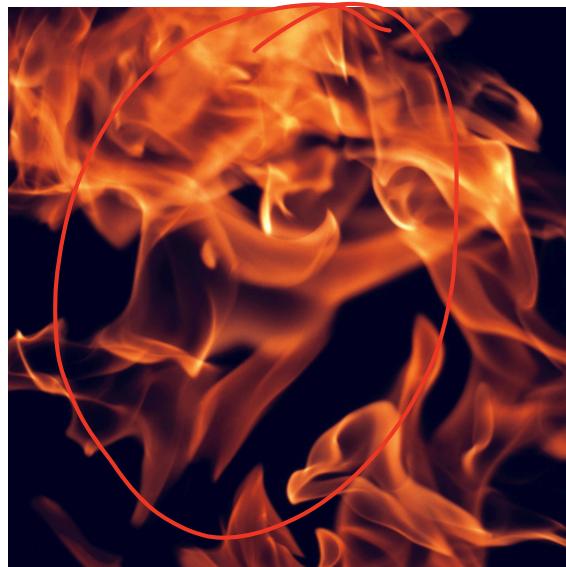
Which is harder to compress?

random ↑ harder ↑

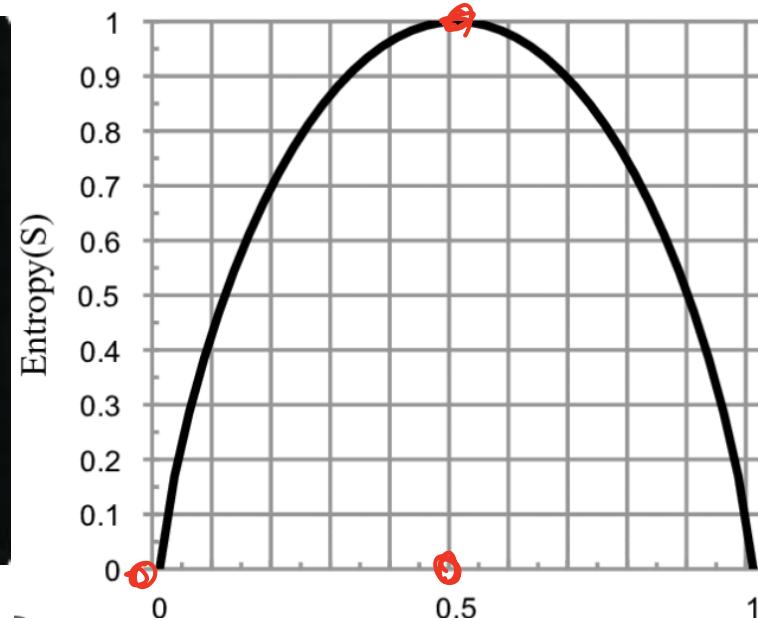


Entropy

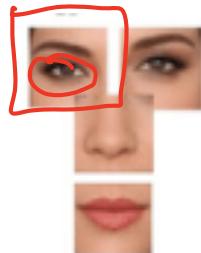
$$H = - \sum p(x) \log p(x)$$



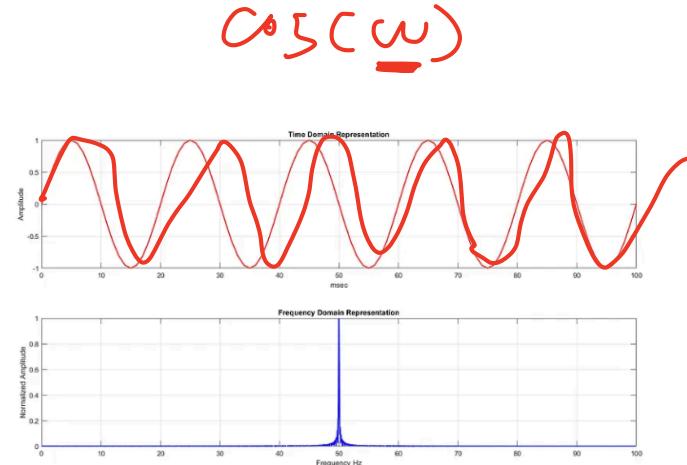
$p(x) = \text{black}$



Redundancy



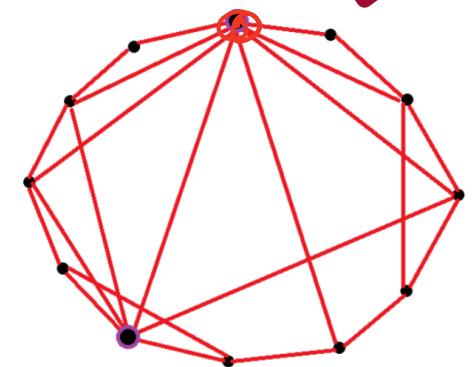
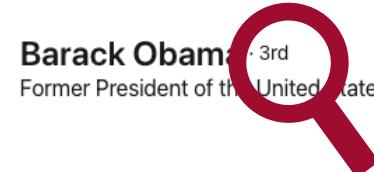
Smoothness/locality prior for image



Periodic prior for sound



Barack Obama, 44th President of the United States of America



Small-world / six-degree prior for social networks

Redundancy in English Text & Entropy of English

Example Rules:

1. i before e except after c.

Hippie, Fries, Field, cake

2. q must always be followed by a u

Quick, quiche, question, quarrel

3. grammar

Cannot do "subject subject" as a sentence

4. dictionary

Hufamomina is not a word

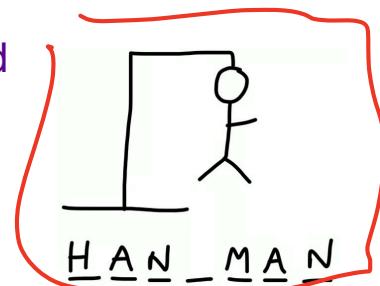
Shannon tried calculating English's entropy using n-gram $p(x_{\{n\}} | x_{\{<n\}})$. Is this a good approach? O complexity?

Another approach: take a set of English words (8000), calculate the word-level entropy based on the subset. Then divide by average number of characters to get Character-level entropy.

Maximum entropy: ~4.7 bit/letter

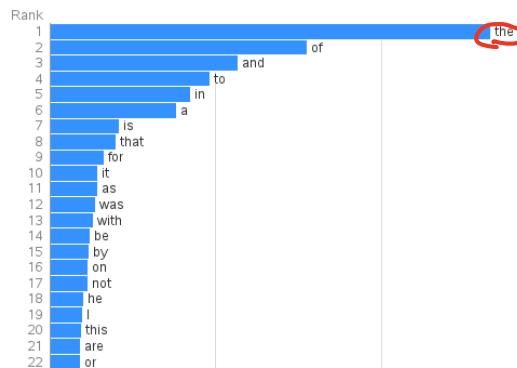
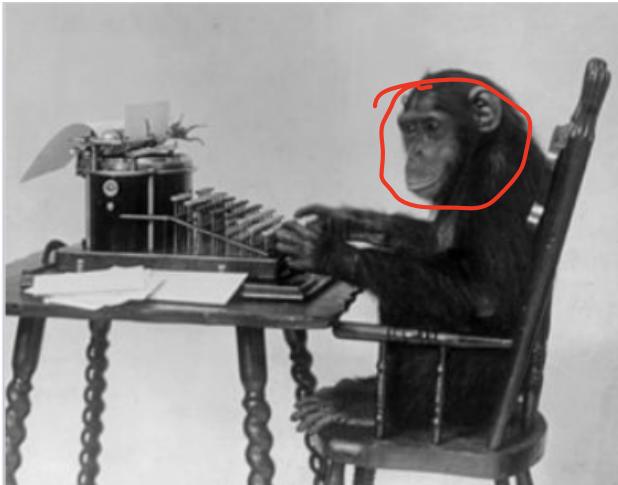
Approximated entropy of English: 2.63 bit/letter (why?)

Is English redundant? Is it good or bad?



Infinite Monkey Theorem

Almost surely, he would type up Shakespeare.



Zero-order approximation	FOML RXKHRJFFJUS ALPWXFWJXYJ FFJEYVJCQSGHYD OPAAMKBZAACIBZLKJQD
First-order approximation	OCRO HLO RGWR NMIELWIS EU LL NBNSEBYA TH EEI ALHENHTTPA OOBTTVA NAH BRL
Second-order approximation	ON IE ANTSOUTINY'S ARE T INCTORE ST BE S DEAMY ACHIN D ILONASIVE TU COOWE AT TEASONARE FUSO TIZIN ANDY TOBE SEACE CTISBE
Third-order approximation	IN NO IST LAT WHEY CRATIC FROURE BIRS GROCID PONDENOME OF DEMONSTURES OF THE REPTAGIN IS REGOACTIONA OF CRE
First-order word approximation	REPRESENTING AND SPEEDILY IS AN GOOD APT OR COME CAN DIFFERENT NATURAL HERE HE THE A IN CAME THE TO OF TO EXPERT GRAY COME TO FURNISHES THE LINE MESSAGE HAD BE THESE
Second-order word approximation	THE HEAD AND IN FRONTAL ATTACK ON AN ENGLISH WRITER THAT THE CHARACTER OF THIS POINT IS THEREFORE ANOTHER METHOD FOR THE LETTERS THAT THE TIME OF WHO EVER TOLD THE PROBLEM FOR AN UNEXPECTED

Compression relies on redundancy

Fixed length v.s. variable length encoding

	a	b	c	d	e	f
Freq in '000s	45	13	12	16	9	5
a fixed-length	000	001	010	011	100	101
* a variable-length	0	101	100	111	1101	1100

Which scheme uses fewer bits to encode the corpus?

How do we encode "bad"?

How do we decode 11000101? (Is the decoding unique? Why?)

{ 1
11
110

1100
fab

unique prefix

Designing Unique Prefixes with Prefix Tree

Idea:

leaves → node

Greedy bottom-up construction of tree

Read encoding based on path from root to leaves.

Why it works:

Each token traces the path of a leave.

A leave has no children.

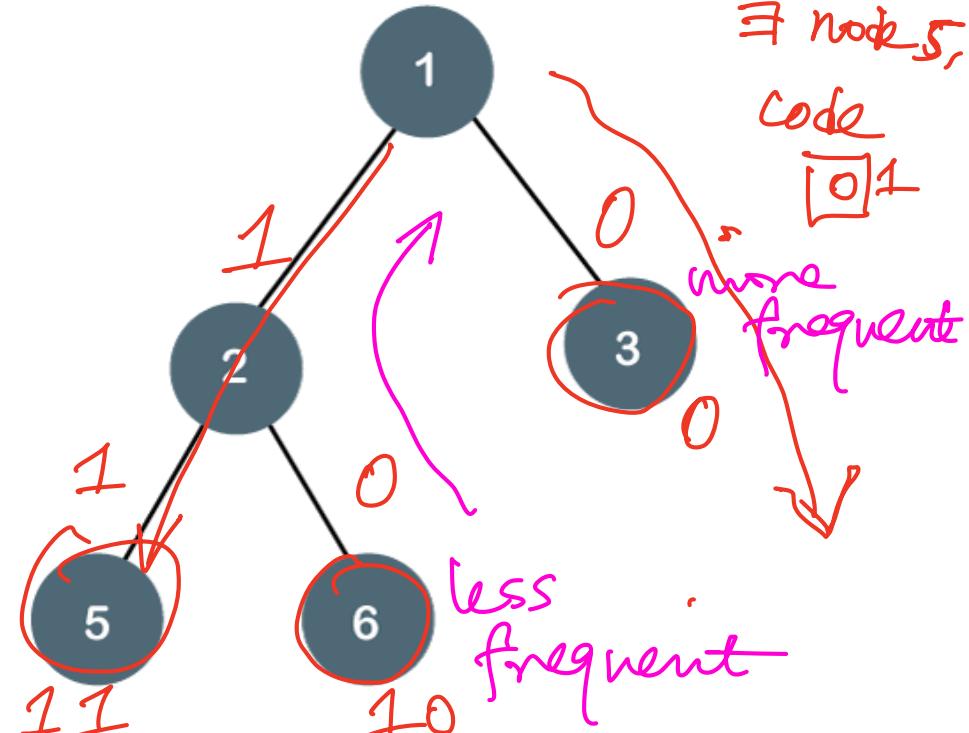
So each token has unique prefix.

Why it's efficient:

Greedy algorithm always look for least frequent token

Less frequent tokens become leaves earlier

Less frequent tokens → longer path → longer code



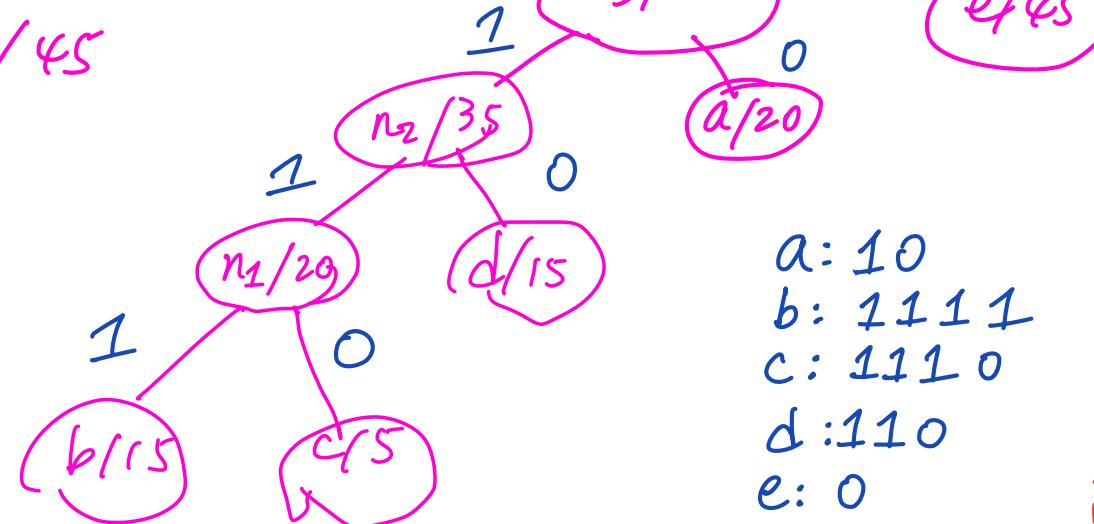
Walkthrough (Unigram)

~~a/20, b/15, c/5, d/15, e/45~~

~~a/20, n₁/20, d/15, e/45~~

~~a/20, n₂/35, e/45~~

~~n₃/55, e/45~~



1. Pick two least frequent words
2. Use them as leaves of a subtree
3. Merge frequency on their common parent
4. Add common parent back to list
5. Repeat

token

$a: 10$
 $b: 1111$
 $c: 1110$
 $d: 110$
 $e: 0$

check:
 encode "bad"
 $1111|10|110$

decode:

$\text{no } 1$
 $\text{no } 11$
 $\text{no } 111$
 $1111 \rightarrow b$
 $\text{no } 1$
 $10 \rightarrow a$
 $110 \rightarrow d$

Walkthrough (Bigram)

$a/20, b/15, c/5, d/15, e/45$

The tokens here (a,b,c,d,e) can be bigrams!

e.g., $\underbrace{a = \text{cat} | \text{the}, b = \text{on} | \text{was}}$

Hierarchical

no code

$$\text{example: } P(X_{t+1} | X_t = \langle \text{bos} \rangle) = \begin{cases} 1, & \text{if } X_{t+1} = a \\ 0, & \text{otherwise} \end{cases}$$

$$\text{Code } \checkmark P(X_{t+1} | X_t = a) = \begin{cases} \frac{1}{2}, & \text{if } X_{t+1} \in \{b, c\} \\ 0, & \text{otherwise} \end{cases}$$

$$X P(X_{t+1} | X_t = b) = \begin{cases} 1, & \text{if } X_{t+1} = \langle \text{eos} \rangle \\ 0, & \text{otherwise} \end{cases}$$

$$X P(X_{t+1} | X_t = c) = \begin{cases} \frac{2}{3}, & \text{if } X_{t+1} = \langle \text{eos} \rangle \\ \frac{1}{3}, & \text{if } X_{t+1} = a \\ 0, & \text{otherwise} \end{cases}$$

- ~~joint~~
1. Pick two least frequent words
 2. Use them as leaves of a subtree
 3. Merge frequency on their common parent
 4. Add common parent back to list
 5. Repeat

Let's rank them by joint!

Only need code for non-deterministic.

Suppose $P(a) = \frac{1}{3}$, $P(b) = \frac{1}{4}$

$$b|a / \frac{1}{2} \times \frac{1}{3} = \frac{1}{6} = \frac{2}{12}$$

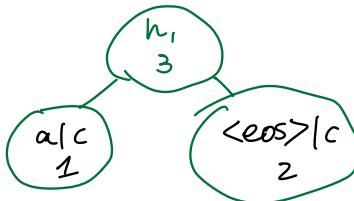
$$c|a / \frac{1}{2} \times \frac{1}{3} = \frac{1}{6} = \frac{2}{12}$$

$$\langle \text{eos} \rangle | c / \frac{2}{3} \times \frac{1}{4} = \frac{1}{6} = \frac{2}{12}$$

$$a|c / \frac{1}{3} \times \frac{1}{4} = \frac{1}{12}$$

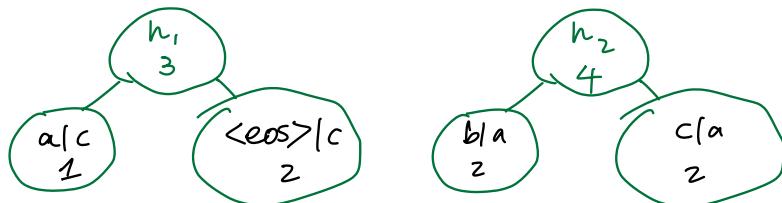
Step 1: $L = \{ \text{bla}/2, \text{cla}/2, \text{<eos>}/2, \text{a/c} / 1 \}$

Pick 2 least frequent, a/c and <eos>/c



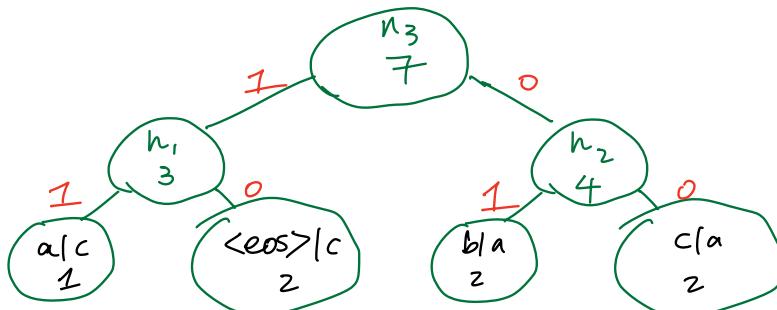
Step 2: $L = \{ \text{bla}/2, \text{cla}/2, \text{n1}/3 \}$

Pick 2 least frequent, bla and cla



Step 3: $L = \{ n1/3, n2/4 \}$

Pick 2 least frequent, n1/3 and n2/4



Step 4: left edge 1, right edge 0 (order is arbitrary)
read code by tracing path from root to leave.

a/c: 11 <eos>/c: 10 bla: 01 cla: 00

Test: encode <bos>a c a b <eos> → 00 11 01

decode: ① <bos>a ② no 0, 00 → cla
100% start with this

③ no 1, 11 → a/c ④ no 0, 01 → bla
⑤ b → 100% <eos> next
out: <bos>a c a b <eos>

References

<https://medium.com/@lazyprogrammerofficial/what-is-temperature-in-nlp-langs-aa2a7212e687>

<https://medium.com/nlplanet/two-minutes-nlp-most-used-decoding-methods-for-language-models-9d44b2375612>

<https://docs.ai21.com/docs/when-the-generation-stops>

<https://towardsdatascience.com/guide-to-chatgpts-advanced-settings-top-p-frequency-penalties-temperature-and-more-b70bae848069>

<https://www.promptingguide.ai/introduction/basics>

<https://www.mathworks.com/help/images/jpeg-image-deblocking-using-deep-learning.html>

https://www.researchgate.net/figure/Top-The-CNN-trained-for-the-task-of-full-face-detection-Bottom-The-CNN-trained-for-the_fig2_308944615

https://cs.stanford.edu/people/eroberts/courses/soco/projects/1999-00/information-theory/entropy_of_english_9.html

<https://home.cse.ust.hk/faculty/golin/COMP271Sp03/Notes/MyL17.pdf>