

Foundations of AI
Yunyu Guo
Assignment 1

- Describe how the choice of training dataset affected your model. What happens if we train a model on one type of data (like music lyrics) and then ask it to work with a different type of data (such as medical reports)? Which datasets do popular language models use, and how might this affect their output? (1 paragraph)

I use medical record as a test of the model, when I use prompts already in the medical record such as “patient” and “physical”, the model seems to generate the rest context more relevant to the medical record. When I use prompt that is not in the medical record, such as “for information on our return policy”, the generated context seems just concatenate random characters. So the choice of training dataset affects the performance and behaviour of a language model.

Models trained on specific types of data learn patterns, vocabulary, and unique to that dataset. Common Datasets: Wikipedia, OpenWebText: A dataset derived from web pages. Effect on Output: Models trained on diverse datasets can handle a wide range of topics. The content and quality of the training data can lead to biases affecting the model's output. (see screenshots of $n=3$, $n=5$, $n=20$ at the end)

- How long did this assignment take you? (1 sentence)
1 day
- Whom did you work with, and how? (1 sentence each)
 - Discussing the assignment with others is encouraged, as long as you don't share the code.
I discuss the outcomes with other student to verify if using relevant and irrelevant prompts would affect the model output.
- Which resources did you use? (1 sentence each)
 - For each, please list the URL and a brief description of how it was useful.
Statistical Language Model: N-gram to calculate the Probability of word sequence using Python: <https://medium.com/codex/statistical-language-model-n-gram-to-calculate-the-probability-of-word-sequence-using-python-2e54a1084250>
- A few sentences about:
 - What was the most difficult part of the assignment?
How to store the n-gram and its probability respectively. What data structure should I use.

- What was the most rewarding part of the assignment?
To be able to generate content after enter the prompt and to design a model that can generate the next character based on probability
- What did you learn doing the assignment?
How to build a nested dictionary.
- Constructive and actionable suggestions for improving assignments, office hours, and class time are always welcome.

Could you provide relevant reading materials to every assignment? For example, in this assignment, provide options of what data structure we should use, the links to materials of how to construct nested dictionary, how a general n-gram model is created with some coding samples.

```

mac@Emmas-Laptop project1 % /usr/local/bin/python3 /Users/mac/Desktop/project1/hw1.py
Testing with n = 3
Enter a prompt (at least 3 characters): patient
Generated text: patient scan confirmed white blood cell confirmed with follower right quadrant. CT scheduled without course was unwever abdominatented on results shower right
quadrant. Lab reveer abdominative complical pain. Physical examinientful. Physical examinientated elevated formed tendectomy. Surgery was unevealed with sevent presults showe
d follow-up appendicitis. Physical pain. Pations. Physicatient was perative course was dscheduled followed on day 3 with seventful. Pations. Pation day 3 without course was ad
mitted appendicitis. Patient was dischargery was unevealed appenderness in the lower right quadrant. CT scan count. Lab revealed.
N-gram probabilities:
Pat: {'l': 1.0}
atl: {'e': 0.5, 'o': 0.3333333333333333, 'v': 0.16666666666666666}
tie: {'n': 1.0}
ien: {'t': 1.0}
ent: {' ': 0.6666666666666666, 'e': 0.16666666666666666, 'f': 0.16666666666666666}
nt: {'p': 0.25, 'w': 0.5, 's': 0.25}
t p: {'r': 1.0}
pr: {'e': 1.0}
pre: {'s': 1.0}
res: {'e': 0.5, 'u': 0.5}
ese: {'n': 1.0}
sen: {'t': 1.0}
nte: {'d': 1.0}
ted: {' ': 1.0}
ed: {'w': 0.375, 't': 0.125, 'e': 0.125, 'a': 0.125, 'f': 0.125, 'o': 0.125}
d w: {'l': 0.6666666666666666, 'h': 0.3333333333333333}
wi: {'t': 1.0}
wit: {'h': 1.0}
ith: {' ': 0.6666666666666666, 'o': 0.3333333333333333}
th: {'s': 0.5, 'f': 0.5}
h s: {'e': 1.0}
se: {'v': 1.0}
sev: {'e': 1.0}
eve: {'r': 0.3333333333333333, 'a': 0.3333333333333333, 'n': 0.3333333333333333}
ver: {'e': 1.0}
ere: {' ': 1.0}
re: {'a': 1.0}
e a: {'b': 1.0}
ab: {'d': 1.0}
abd: {'o': 1.0}
bdo: {'m': 1.0}
dom: {'n': 1.0}
oml: {'n': 1.0}
min: {'a': 1.0}
ina: {'l': 0.5, 't': 0.5}
nat: {' ': 1.0}
al: {'p': 0.5, 'e': 0.5}
l p: {'a': 1.0}
pa: {'l': 1.0}
pai: {'n': 1.0}
ain: {' ': 1.0}
in: {' ': 1.0}
n. : {'P': 1.0}
. P: {'h': 0.25, 'a': 0.5, 'o': 0.25}
Ph: {'y': 1.0}
Phy: {'s': 1.0}
hys: {'l': 1.0}
ysi: {'C': 1.0}

```

```
Code File Edit Selection View Go Run Terminal Window Help
project1
hw1.py M X requirement stranger
main
class CharNGramLanguageModel:
    def __init__(self, d: dict, ed: dict):
        self.d = d
        self.ed = ed

    def test(self, n: int):
        print(f"Testing with n = {n}")
        prompt = "Enter a prompt (at least 5 characters): physical"
        generated_text = "Generated text: physical examination revealed tenderness in the lower right quadrant. Lab results showed elevated white blood cell count. CT scan confirmed with follow-up appointment scheduled."
        ngram_probs = {}
        for i in range(n):
            ngram_probs[f"{i}": 1.0] = {}
        patie: {'n': 1.0}
        atien: {'t': 1.0}
        tient: {'t': 1.0}
        ient: {'p': 0.3333333333333333, 'w': 0.6666666666666666}
        ent p: {'r': 1.0}
        nt pr: {'e': 1.0}
        t pre: {'s': 1.0}
        pre: {'e': 1.0}
        prese: {'n': 1.0}
        reson: {'t': 1.0}
        esent: {'e': 1.0}
        sente: {'d': 1.0}
        ented: {'t': 1.0}
        nted: {'w': 1.0}
        ted w: {'l': 0.5, 'h': 0.5}
        ed w: {'t': 1.0}
        d wit: {'h': 1.0}
        with: {'s': 0.6666666666666666, 'o': 0.3333333333333333}
        ith s: {'e': 1.0}
        th se: {'v': 1.0}
        h sev: {'e': 1.0}
        sever: {'e': 1.0}
        evere: {'e': 1.0}
        vere: {'a': 1.0}
        ere a: {'b': 1.0}
        re ab: {'d': 1.0}
        e abd: {'o': 1.0}
        abdo: {'m': 1.0}
        abdom: {'t': 1.0}
        bdomi: {'n': 1.0}
        domini: {'a': 1.0}
        omina: {'l': 1.0}
        minal: {'p': 1.0}
        inal: {'p': 1.0}
        nal p: {'a': 1.0}
        al pa: {'i': 1.0}
        l pai: {'n': 1.0}
        pain: {'n': 1.0}
        pain: {'n': 1.0}
        ain: {'P': 1.0}
        in. P: {'h': 1.0}
        n. Ph: {'y': 1.0}

Ln 86, Col 1 Spaces: 4 UTF-8 LF Python 3.12.4 64-bit
```

```
Code File Edit Selection View Go Run Terminal Window Help
project1
hw1.py M X requirement stranger
main
class CharNGramLanguageModel:
    def __init__(self, d: dict, ed: dict):
        self.d = d
        self.ed = ed

    def test(self, n: int):
        print(f"Testing with n = {n}")
        prompt = "Enter a prompt (at least 20 characters): for information on our return policy"
        generated_text = "Generated text: for information on our return policygsiban aweneielsyxw u rt rdsh.ert itnsadidrohtyeroswtqeaufupta mhf o tndw .aey sifveeps s .oectp tt-aqcnqche.noa qnicsg feel n 3psor.readnarse tsrtc loris coaCismreotaclessCrt d t ana y.seo odatui cp. sruhoetwoneatefaweels yptodd wuguonmwa doht cttms lwnebm br rter we c"
        ngram_probs = {}
        for i in range(n):
            ngram_probs[f"{i}": 1.0] = {}
        Patient presented wit: {'t': 1.0}
        atient presented wit: {'h': 1.0}
        tient presented with: {'t': 1.0}
        ient presented with: {'s': 1.0}
        ent presented with s: {'e': 1.0}
        nt presented with se: {'v': 1.0}
        t presented with sev: {'e': 1.0}
        presented with severe: {'e': 1.0}
        presented with severe: {'e': 1.0}
        resented with severe: {'t': 1.0}
        esented with severe: {'a': 1.0}
        sented with severe a: {'b': 1.0}
        ented with severe ab: {'d': 1.0}
        nted with severe abd: {'o': 1.0}
        ted with severe abdo: {'m': 1.0}
        ed with severe abdom: {'t': 1.0}
        d with severe abdomi: {'n': 1.0}
        with severe abdomin: {'a': 1.0}
        with severe abdomina: {'l': 1.0}
        ith severe abdominal: {'p': 1.0}
        th severe abdominal: {'p': 1.0}
        h severe abdominal p: {'a': 1.0}
        severe abdominal pa: {'i': 1.0}
        severe abdominal pai: {'n': 1.0}
        evere abdominal pain: {'n': 1.0}
        vere abdominal pain.: {'n': 1.0}
        ere abdominal pain.: {'P': 1.0}
        re abdominal pain. P: {'h': 1.0}
        e abdominal pain. Ph: {'y': 1.0}
        abdominal pain. Phy: {'s': 1.0}
        abdominal pain. Phys: {'t': 1.0}
        bdominal pain. Physi: {'e': 1.0}
        dominal pain. Physic: {'a': 1.0}
        ominal pain. Physica: {'l': 1.0}
        minal pain. Physical: {'p': 1.0}
        inal pain. Physical: {'e': 1.0}
        nal pain. Physical e: {'x': 1.0}
        al pain. Physical ex: {'a': 1.0}
        l pain. Physical exa: {'m': 1.0}
        pain. Physical exam: {'n': 1.0}
        pain. Physical exami: {'n': 1.0}
        ain. Physical examin: {'a': 1.0}
        in. Physical examina: {'t': 1.0}
        n. Physical examinat: {'t': 1.0}

Ln 86, Col 1 Spaces: 4 UTF-8 LF Python 3.12.4 64-bit
```