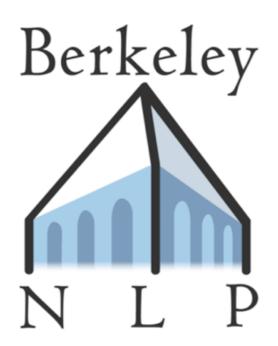
Sequence Modeling



EECS 183/283a: Natural Language Processing

Modeling Language



Modeling the expectations we have over utterances we encounter

- How can we find structure in continuous signals like speech?
- What words are more frequent vs. rare?
- What combinations of words are more likely vs. unlikely?
- What sequences of utterances are plausible vs. implausible?
- How likely is it for certain words (or combinations of words) to appear alongside different contexts vs. others?



- For now: let's assume utterances are sequences of tokens from our vocabulary
- Our vocabulary has a fixed size and consists of discrete wordtypes (we'll get to modeling continuous language signals, like speech, in a few weeks!)
- A sequence is denoted as:

$$\overline{x} = \langle x_1, \dots, x_n \rangle \quad x \in \mathcal{V} \quad x_n = \text{EOS}$$

• We can also consider writing out all possible sequences given our vocabulary (though this set is infinitely large): \mathcal{V}^+



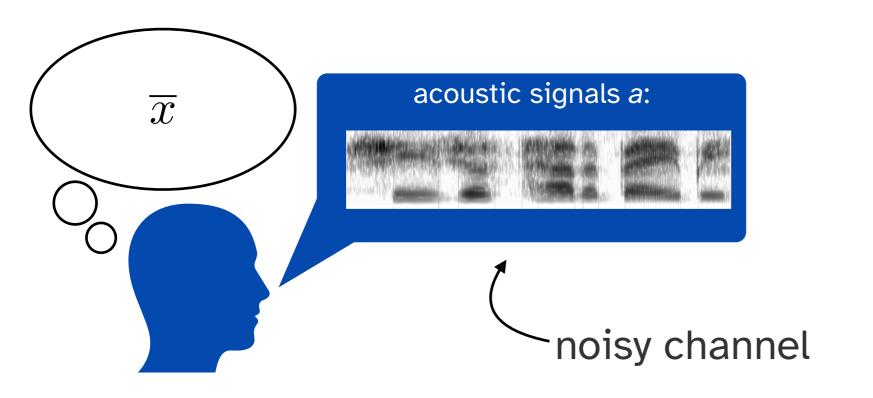
ullet A sequence model imposes a probability distribution over \mathcal{V}^+

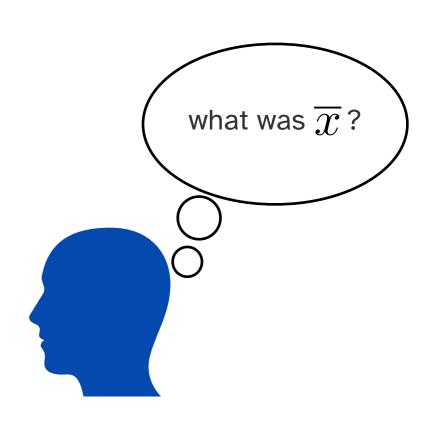
$$p(\overline{X}) \in \Delta^{\mathcal{V}^+} \qquad \overline{x} \in \mathcal{V}^+ \qquad p(\overline{x})$$



ullet A sequence model imposes a probability distribution over \mathcal{V}^+

$$p(\overline{X}) \in \Delta^{\mathcal{V}^+} \qquad \overline{x} \in \mathcal{V}^+ \qquad p(\overline{x})$$





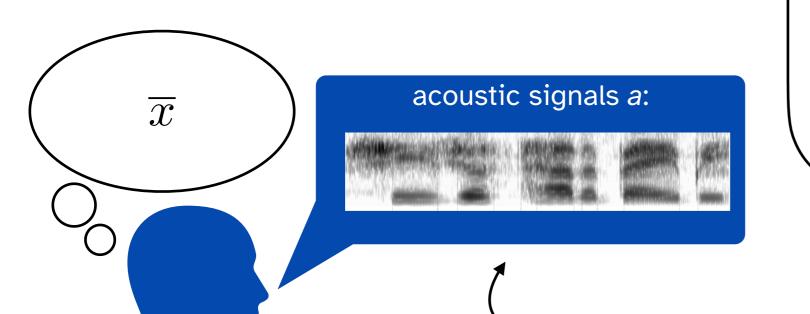


ullet A sequence model imposes a probability distribution over \mathcal{V}^+

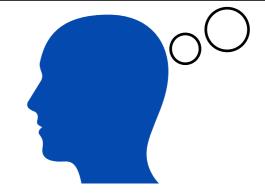
$$p(\overline{X}) \in \Delta^{\mathcal{V}^+} \qquad \overline{x} \in \mathcal{V}^+ \qquad p(\overline{x})$$

noisy channel

 Why is this useful? Bayes' rule



 $\overline{x}^* = \arg\max_{\overline{x}} p(\overline{x} \mid a)$ $\Rightarrow = \arg\max_{\overline{x}} \frac{p(a \mid \overline{x})p(\overline{x})}{p(a)}$ $= \arg\max_{\overline{x}} p(a \mid \overline{x})p(\overline{x})$ acoustic language



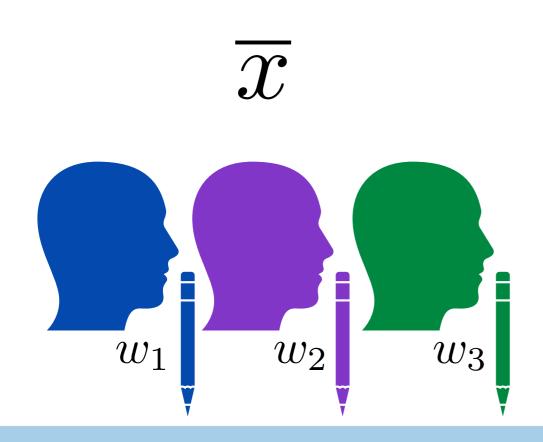
model

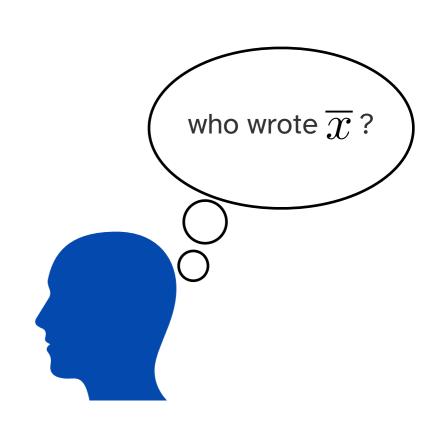
model



ullet A sequence model imposes a probability distribution over \mathcal{V}^+

$$p(\overline{X}) \in \Delta^{\mathcal{V}^+} \qquad \overline{x} \in \mathcal{V}^+ \qquad p(\overline{x})$$

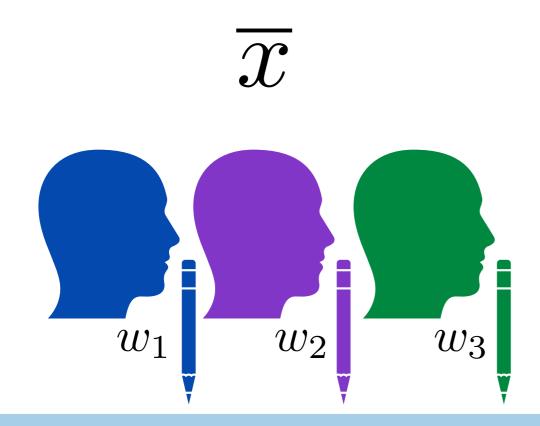


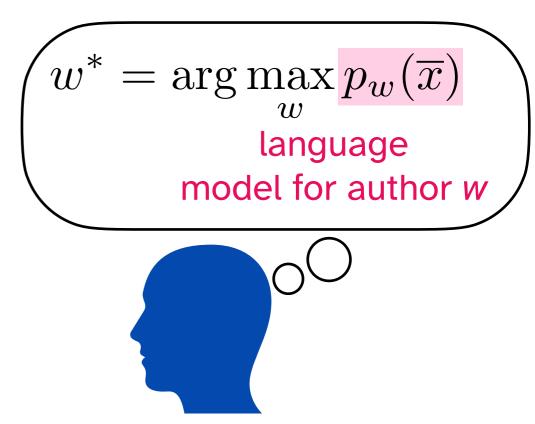




ullet A sequence model imposes a probability distribution over \mathcal{V}^+

$$p(\overline{X}) \in \Delta^{\mathcal{V}^+} \qquad \overline{x} \in \mathcal{V}^+ \qquad p(\overline{x})$$

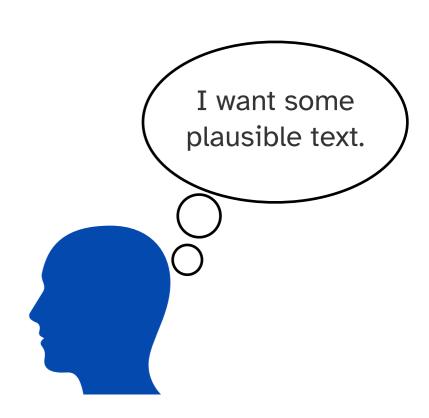






ullet A sequence model imposes a probability distribution over \mathcal{V}^+

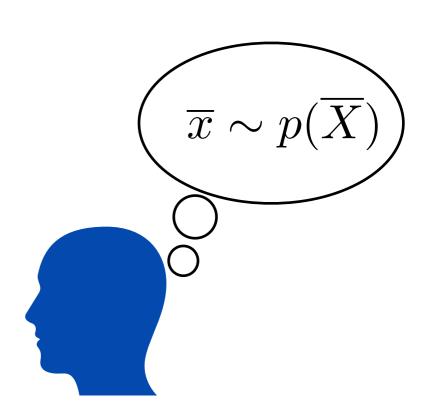
$$p(\overline{X}) \in \Delta^{\mathcal{V}^+} \qquad \overline{x} \in \mathcal{V}^+ \qquad p(\overline{x})$$





ullet A sequence model imposes a probability distribution over \mathcal{V}^+

$$p(\overline{X}) \in \Delta^{\mathcal{V}^+} \qquad \overline{x} \in \mathcal{V}^+ \qquad p(\overline{x})$$





Documents + frequencies: $\mathcal{D} = (\text{'hug', 10), ('pug', 5),} \\ \mathcal{D} = (\text{'pun', 12), ('bun', 4),} \\ (\text{'hugs', 5)} \\ \mathcal{V} = \{b, g, h, n, p, s, u\} \\ \mathcal{V}^+ = \{b, g, h, \dots, bb, bg, \dots, bug, bun, \dots, sssssss, \dots\}$



Documents + frequencies:

Documents + frequencies:
$$\mathcal{D}=\mbox{('hug', 10), ('pug', 5),}\\ \mathcal{D}=\mbox{('pun', 12), ('bun', 4),}\\ \mbox{('hugs', 5)} \label{eq:pun', 5}$$

$$\mathcal{V}^+ = \{b, g, h, \dots, bb, bg, \dots, bug, bun, \dots, sssssss, \dots\}$$

 Learning problem: we want a to estimate the probability distribution $p(\overline{X}) \in \Delta^{\mathcal{V}^+}$ that generated our observations \mathcal{D}



Documents + frequencies:

$$\mathcal{D} = (\text{'hug', 10}), (\text{'pug', 5}), \\ (\text{'pun', 12}), (\text{'bun', 4}), \\ (\text{'hugs', 5})$$

$$\mathcal{V} = \{b, g, h, n, p, s, u\}$$

$$\mathcal{V}^+ = \{b, g, h, \dots, bb, bg, \dots, bug, bun, \dots, sssssss, \dots\}$$

- Learning problem: we want a to estimate the probability distribution $p(\overline{X}) \in \Delta^{\mathcal{V}^+}$ that generated our observations \mathcal{D}
- One simple option: just count based on document occurrence

$$p(\overline{x}) = \frac{C(\overline{x})}{|\mathcal{D}|}$$

| $\overline{\mathcal{X}}$ | $p(\overline{x}) = \frac{C(\overline{x})}{ \mathcal{D} }$ |
|--|---|
| hug | 10/36 = 0.28 |
| pug | 5/36 = 0.14 |
| pun | 12/36 = 0.33 |
| bun | 4/36 = 0.11 |
| hugs | 5/36 = 0.14 |
| $\overline{x} \in \mathcal{V}^+ \setminus \mathcal{D}$ | 0/36 = 0.00 |

$$p(\overline{X})$$



Documents + frequencies:

$$\mathcal{D} = (\text{'hug', 10}), (\text{'pug', 5}), \\ \mathcal{D} = (\text{'pun', 12}), (\text{'bun', 4}), \\ (\text{'hugs', 5})$$

$$\mathcal{V} = \{b, g, h, n, p, s, u\}$$

$$\mathcal{V}^+ = \{b, g, h, \dots, bb, bg, \dots, bug, bun, \dots, sssssss, \dots\}$$

• Let's sample from $p(\overline{X})!$

| \overline{x} | $p(\overline{x}) = \frac{C(\overline{x})}{ \mathcal{D} }$ |
|--|---|
| hug | 10/36 = 0.28 |
| pug | 5/36 = 0.14 |
| pun | 12/36 = 0.33 |
| bun | 4/36 = 0.11 |
| hugs | 5/36 = 0.14 |
| $\overline{x} \in \mathcal{V}^+ \setminus \mathcal{D}$ | 0/36 = 0.00 |

$$p(\overline{X})$$



Documents + frequencies:

$$\mathcal{D} = ('hug', 10), ('pug', 5),$$
 $('pun', 12), ('bun', 4),$
 $('hugs', 5)$

$$\mathcal{V} = \{b, g, h, n, p, s, u\}$$

$$\mathcal{V}^+ = \{b, g, h, \dots, bb, bg, \dots, bug, bun, \dots, sssssss, \dots\}$$

• Let's sample from $p(\overline{X})!$ 0.98 \rightarrow hugs

| | \overline{x} | $p(\overline{x}) = \frac{C(\overline{x})}{ \mathcal{D} }$ |
|-------------|--|---|
| | hug | 10/36 = 0.28 |
| | pug | 5/36 = 0.14 |
| | pun | 12/36 = 0.33 |
| | bun | 4/36 = 0.11 |
| > | hugs | 5/36 = 0.14 |
| • | $\overline{x} \in \mathcal{V}^+ \setminus \mathcal{D}$ | 0/36 = 0.00 |

$$p(\overline{X})$$



Documents + frequencies:

$$\mathcal{D} = (\text{'hug', 10}), (\text{'pug', 5}), \\ (\text{'pun', 12}), (\text{'bun', 4}), \\ (\text{'hugs', 5})$$

$$\mathcal{V} = \{b, g, h, n, p, s, u\}$$

$$\mathcal{V}^+ = \{b, g, h, \dots, bb, bg, \dots, bug, bun, \dots, sssssss, \dots\}$$

• Let's sample from p(X)!

$$0.98 \rightarrow hugs$$

$$0.84 \rightarrow bun$$

| $oldsymbol{\mathcal{X}}$ | $p(\overline{x}) = \frac{C(x)}{ \mathcal{D} }$ |
|--|--|
| hug | 10/36 = 0.28 |
| pug | 5/36 = 0.14 |
| pun | 12/36 = 0.33 |
| bun | 4/36 = 0.11 |
| hugs | 5/36 = 0.14 |
| $\overline{x} \in \mathcal{V}^+ \setminus \mathcal{D}$ | 0/36 = 0.00 |

 $p(\overline{X})$



Documents + frequencies:

$$\mathcal{D} = (\text{'hug', 10}), (\text{'pug', 5}), \\ (\text{'pun', 12}), (\text{'bun', 4}), \\ (\text{'hugs', 5})$$

$$\mathcal{V} = \{b, g, h, n, p, s, u\}$$

$$\mathcal{V}^+ = \{b, g, h, \dots, bb, bg, \dots, bug, bun, \dots, sssssss, \dots\}$$

• Let's sample from $p(\overline{X})!$

$$0.98 \rightarrow hugs$$

$$0.84 \rightarrow bun$$

$$0.55 \rightarrow pun$$

| | \overline{x} | $p(\overline{x}) = \frac{C(\overline{x})}{ \mathcal{D} }$ |
|----------|--|---|
| | hug | 10/36 = 0.28 |
| | pug | 5/36 = 0.14 |
| + | pun | 12/36 = 0.33 |
| | bun | 4/36 = 0.11 |
| | hugs | 5/36 = 0.14 |
| | $\overline{x} \in \mathcal{V}^+ \setminus \mathcal{D}$ | 0/36 = 0.00 |

$$p(\overline{X})$$



Documents + frequencies:

$$\mathcal{D} = ('hug', 10), ('pug', 5),$$
 $('pun', 12), ('bun', 4),$
 $('hugs', 5)$

$$\mathcal{V} = \{b, g, h, n, p, s, u\}$$

$$\mathcal{V}^+ = \{b, g, h, \dots, bb, bg, \dots, bug, bun, \dots, sssssss, \dots\}$$

- Let's sample from $p(\overline{X})!$
 - $0.98 \rightarrow hugs$
 - $0.84 \rightarrow bun$
 - $0.55 \rightarrow pun$
- This directly generates our observation data! But nobody ever does this
- Why not?

| \overline{x} | $p(\overline{x}) = \frac{C(\overline{x})}{ \mathcal{D} }$ |
|--|---|
| hug | 10/36 = 0.28 |
| pug | 5/36 = 0.14 |
| pun | 12/36 = 0.33 |
| bun | 4/36 = 0.11 |
| hugs | 5/36 = 0.14 |
| $\overline{x} \in \mathcal{V}^+ \setminus \mathcal{D}$ | 0/36 = 0.00 |

$$p(\overline{X})$$

Language Modeling is Hard



$$p(\overline{X}) \in \Delta^{\mathcal{V}^+}$$

$$\overline{x} \in \mathcal{V}^+$$

$$p(\overline{x})$$

- How might we go about assigning a probability to any possible sequence, even ones we've never seen before?
- One intuition: sentences have internal consistencies!

você fala inglês? do you speak English?

2nd person 3rd person singular singular pronoun, present formal

3rd person singular pronoun, masculine

అతను

atanu he accusative case of "dog"

కుక్కను

kukkanu the dog 3rd person singular present, masculine

చూస్తాడు

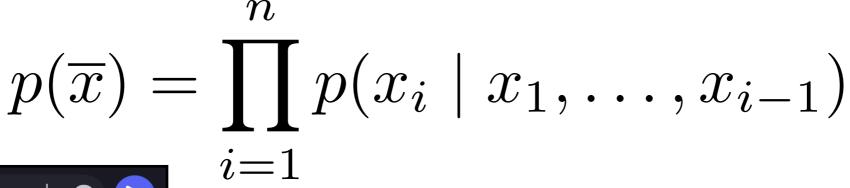
cūstāḍu sees

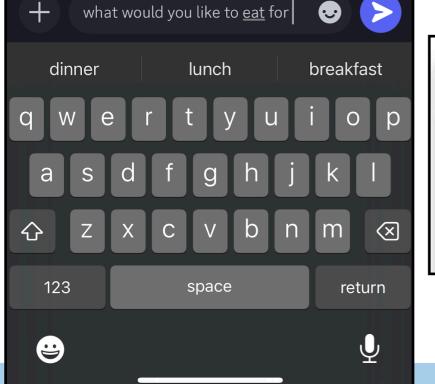
One Approximation

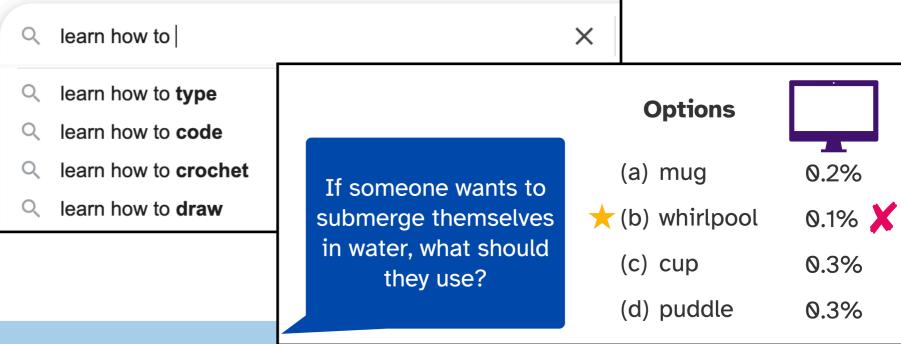


Autoregressive language modeling:

- The probability of a sequence is a product of local token probabilities
- The probability of a token depends on the ones that came before it







Another Approximation



Masked language modeling:

- The probability of a sequence is a product of local token probabilities
- The probability of a token depends on the ones that came before and after it

$$p(\overline{x}) = \prod_{i=1}^{n} p(x_i \mid x_1, \dots, x_{i-1}, x_{i+1} \dots x_n)$$

The name *anteater* refers to the species', which consists mainly of ants and termites.

This approximation will come up later in the class! For now, we'll focus on autoregressive language modeling.

(There are other approximations too!)



$$p(\overline{x}) = \prod_{i=1}^n p(x_i \mid x_1, \dots, x_{i-1})$$
 Chain rule decomposition $= p(x_1)p(x_2 \mid x_1)\dots p(x_n \mid x_1, \dots, x_{n-1})$

Let's sample a sequence from this approximation:

• Sample the first word $x_1 \sim p(X_1) \in \Delta^{\mathcal{V}}$

 $p(X_1)$

$$\overline{x} = \langle \mathit{the} \rangle$$



$$p(\overline{x}) = \prod_{i=1}^{n} p(x_i \mid x_1, \dots, x_{i-1})$$

= $p(x_1)p(x_2 \mid x_1) \dots p(x_n \mid x_1, \dots, x_{n-1})$

- Sample the first word $x_1 \sim p(X_1) \in \Delta^{\mathcal{V}}$ $p(X_2 \mid x_1 = \text{the})$
- Sample the second word $x_2 \sim p(X_2 \mid x_1) \in \Delta^{\mathcal{V}}$

| | -1r | 0.03 |
|----------------------------------|--------|------|
| | same | 0.02 |
| | impact | 0.02 |
| $\overline{x}=\langle$ the, same | line | 0.02 |
| Julie, Saille | ••• | |



$$p(\overline{x}) = \prod_{i=1}^{n} p(x_i \mid x_1, \dots, x_{i-1})$$

= $p(x_1)p(x_2 \mid x_1) \dots p(x_n \mid x_1, \dots, x_{n-1})$

- Sample the first word $x_1 \sim p(X_1) \in \Delta^{\mathcal{V}}$ $p(X_3 \mid \underline{\langle \text{the}, \text{same} \rangle})$
- Sample the second word
- Sample the third word

$$\overline{x}=\langle$$
 the, same, thing

| ${\mathcal X}$ | p(x) |
|----------------|------|
| as | 0.08 |
| way | 0.04 |
| thing | 0.04 |
| time | 0.04 |
| ••• | |



$$p(\overline{x}) = \prod_{i=1}^{n} p(x_i \mid x_1, \dots, x_{i-1})$$

= $p(x_1)p(x_2 \mid x_1) \dots p(x_n \mid x_1, \dots, x_{n-1})$

- Sample the first word $x_1 \sim p(X_1) \in \Delta^{\mathcal{V}} p(X_i \mid \langle x_1, \dots, x_{i-1} \rangle)$
- Sample the second word
- Sample the third word
- Keep sampling until we hit EOS

$$\overline{x}=\langle$$
 the, same, thing, that

| $\overline{\mathcal{X}}$ | p(x) |
|--------------------------|------|
| • | 0.11 |
| as | 0.08 |
| that | 0.07 |
| , | 0.06 |
| ••• | |



$$p(\overline{x}) = \prod_{i=1}^{n} p(x_i \mid x_1, \dots, x_{i-1})$$

= $p(x_1)p(x_2 \mid x_1) \dots p(x_n \mid x_1, \dots, x_{n-1})$

- Sample the first word $x_1 \sim p(X_1) \in \Delta^{\mathcal{V}} p(X_i \mid \langle x_1, \dots, x_{i-1} \rangle)$
- Sample the second word
- Sample the third word
- Keep sampling until we hit EOS

$$\overline{x} = \langle$$
 the, same, thing, that, the

| \mathcal{X} | p(x) |
|---------------|------|
| happened | 0.11 |
| happens | 0.08 |
| makes | 0.07 |
| the | 0.06 |
| ••• | |



$$p(\overline{x}) = \prod_{i=1}^{n} p(x_i \mid x_1, \dots, x_{i-1})$$

= $p(x_1)p(x_2 \mid x_1) \dots p(x_n \mid x_1, \dots, x_{n-1})$

- Sample the first word $x_1 \sim p(X_1) \in \Delta^{\mathcal{V}} p(X_i \mid \langle x_1, \dots, x_{i-1} \rangle)$
- Sample the second word
- Sample the third word
- Keep sampling until we hit EOS

$$\overline{x} = \langle$$
 the, same, thing, that, the, -y

| 1 \ | , <i>o</i> <u> </u> |
|----------------|---------------------|
| ${\mathcal X}$ | p(x) |
| -у | 0.03 |
| -ir | 0.02 |
| person | 0.02 |
| line | 0.02 |
| ••• | |
| | |



$$p(\overline{x}) = \prod_{i=1}^{n} p(x_i \mid x_1, \dots, x_{i-1})$$

= $p(x_1)p(x_2 \mid x_1) \dots p(x_n \mid x_1, \dots, x_{n-1})$

Let's sample a sequence from this approximation:

- Sample the first word $x_1 \sim p(X_1) \in \Delta^{\mathcal{V}} p(X_i \mid \langle x_1, \dots, x_{i-1} \rangle)$
- Sample the second word
- Sample the third word
- Keep sampling until we hit EOS

 $\overline{x} = \langle$ the, same, thing, that, the, -y, had

| |) ** t-1/ |
|----------------|-----------|
| ${\mathcal X}$ | p(x) |
| are | 0.12 |
| were | 0.05 |
| have | 0.04 |
| had | 0.02 |
| ••• | |
| | |



$$p(\overline{x}) = \prod_{i=1}^{n} p(x_i \mid x_1, \dots, x_{i-1})$$

= $p(x_1)p(x_2 \mid x_1) \dots p(x_n \mid x_1, \dots, x_{n-1})$

Let's sample a sequence from this approximation:

- Sample the first word $x_1 \sim p(X_1) \in \Delta^{\mathcal{V}} p(X_i \mid \langle x_1, \dots, x_{i-1} \rangle)$
- Sample the second word
- Sample the third word
- Keep sampling until we hit EOS

 $\overline{x} = \langle \text{the, same, thing, that, the, -y, had, EOS} \rangle$ the same thing that they had

| $oldsymbol{\mathcal{X}}$ | p(x) |
|--------------------------|------|
| done | 0.08 |
| said | 0.04 |
| EOS | 0.04 |
| \ to | 0.02 |
| ? / | |

Autoregressive Language Models



$$p(\overline{x}) = \prod_{i=1}^{n} p(x_i \mid x_1, \dots, x_{i-1})$$

$$= p(x_1) p(x_2 \mid x_1) \dots p(x_n \mid x_1, \dots, x_{n-1})$$

probability that probability that the first word is x_1 the second word is x_2 , given that the first word is x_1

probability that the sentence ends after the sequence $\langle x_1, \dots, x_{n-1} \rangle$

Core modeling challenge:

How do we compute these conditional probabilities?