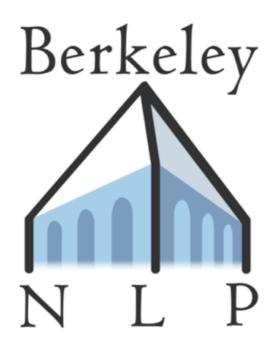
Sequence Modeling: Generation



EECS 183/283a: Natural Language Processing

Generation



- Given a language model $\;p(\overline{X})\in\Delta$ how do we get a plausible sequence out of it?
- Autoregressive language model:

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

- We can sample directly from this approximation
- We can adjust this distribution, then sample from it
- We can try to find the most plausible sequence in \mathcal{V}^+

want some plausible text.

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



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

• As we generate, we build our output sequence \overline{x} , which starts as an empty sequence $\overline{x}=\langle\ \rangle$



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

- As we generate, we build our output sequence \overline{x} , which starts as an empty sequence $\overline{x}=\langle\ \rangle$
- At each step i, we choose an item from the vocabulary $\mathcal V$ by performing some operation on the local probability distribution $p(X_i \mid x_1, \dots, x_{i-1}) \in \Delta^{\mathcal V}$



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

- As we generate, we build our output sequence \overline{x} , which starts as an empty sequence $\overline{x}=\langle\ \rangle$
- At each step i, we choose an item from the vocabulary $\mathcal V$ by performing some operation on the local probability distribution $p(X_i \mid x_1, \dots, x_{i-1}) \in \Delta^{\mathcal V}$
- Then, we append this to our running sequence $\overline{x} \leftarrow \overline{x} + \langle x_i \rangle$



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

- As we generate, we build our output sequence \overline{x} , which starts as an empty sequence $\overline{x}=\langle\ \rangle$
- At each step i, we choose an item from the vocabulary $\mathcal V$ by performing some operation on the local probability distribution $p(X_i \mid x_1, \dots, x_{i-1}) \in \Delta^{\mathcal V}$
- Then, we append this to our running sequence $\overline{x} \leftarrow \overline{x} + \langle x_i \rangle$
- If we ever choose EOS, we stop generation



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

- As we generate, we build our output sequence \overline{x} , which starts as an empty sequence $\overline{x}=\langle\ \rangle$
- At each step i, we choose an item from the vocabulary $\mathcal V$ by performing some **operation** on the local probability distribution $p(X_i \mid x_1, \dots, x_{i-1}) \in \Delta^{\mathcal V}$
- Then, we append this to our running sequence $\overline{x} \leftarrow \overline{x} + \langle x_i \rangle$
- If we ever choose EOS, we stop generation
- Main point: generation methods differ wrt the operation they perform on $p(X_i \mid x_1, \dots, x_{i-1})$



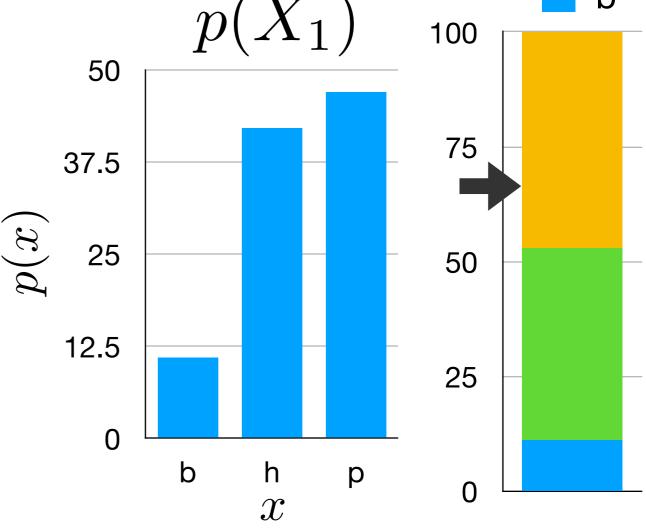
$$p(\overline{x}) = \prod_{i=1}^{n} p(x_i \mid x_1, \dots, x_{i-1})$$

Operation: sample directly from
$$p(X_i \mid x_1, \dots, x_{i-1})$$

$$\overline{x} = \langle \rho \rangle$$

number = random (0, 100)

number is 65





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

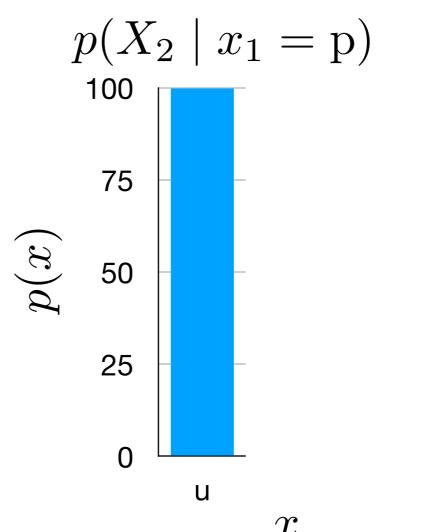
Operation: sample directly from $p(X_i \mid x_1, \dots, x_{i-1})$

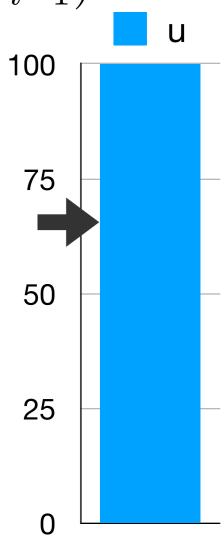
$$p(X_i \mid x_1, \dots, x_{i-1})$$

$$\overline{x} = \langle pu \rangle$$

number = random (0, 100)

number is 63



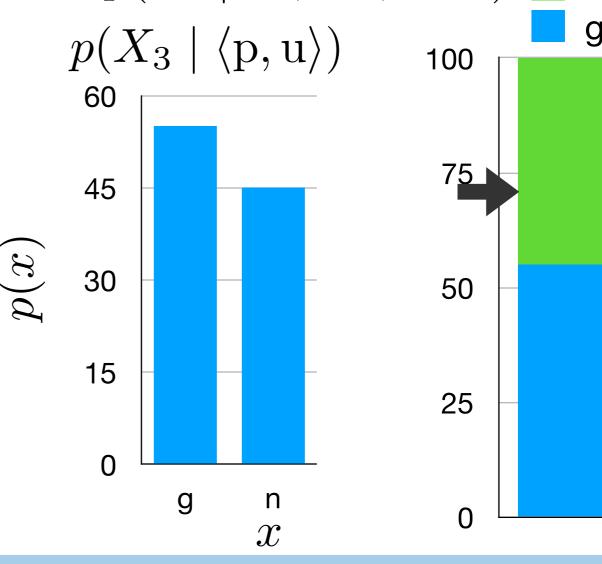




$$p(\overline{x}) = \prod_{i=1}^{n} p(x_i \mid x_1, \dots, x_{i-1})$$

Operation: sample directly from $p(X_i \mid x_1, \dots, x_{i-1})$

$$\overline{x} = \langle pun \rangle$$



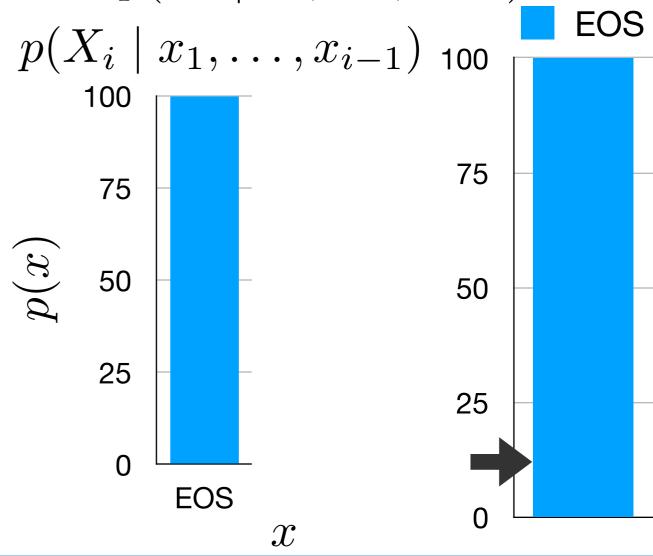


$$p(\overline{x}) = \prod_{i=1}^{n} p(x_i \mid x_1, \dots, x_{i-1})$$

Operation: sample directly from $p(X_i \mid x_1, \dots, x_{i-1})$

$$\overline{x} = \langle pun Eos \rangle$$

pun





Operation: modify the logits before computing probabilities

Sneak peek: computing probabilities over wordtypes using pretty much any modern language model

Score each wordtype independently

$$s(w) = f(w \mid x_1, \dots, x_{i-1}; \theta)$$
 \bullet logits

Renormalize using softmax

$$p(X_i = w \mid x_1, \dots, x_{i-1}) = \frac{\exp(s(w))}{\sum_{w' \in \mathcal{V}} \exp(s(w'))}$$



Operation: modify the logits before computing probabilities

$$s(w) = f(w \mid x_1, \dots, x_{i-1}; \theta) \leftarrow logits$$

$$p(X_i = w \mid x_1, \dots, x_{i-1}) = \frac{\exp(s(w))}{\sum_{w' \in \mathcal{V}} \exp(s(w'))}$$

Temperature parameter controls the "smoothness" of this distribution:

$$p(X_i = w \mid x_1, \dots, x_{i-1}) = \frac{\exp(s(w)/\tau)}{\sum_{w' \in \mathcal{V}} \exp(s(w')/\tau)}$$

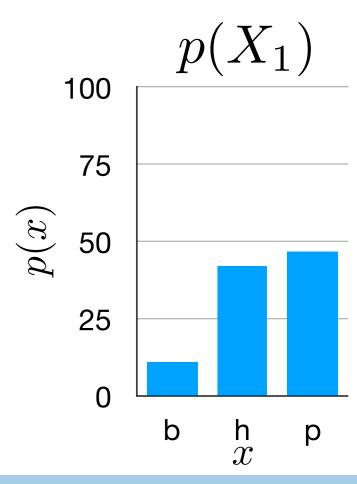


Operation: modify the logits before computing probabilities

$$s(w) = f(w \mid x_1, \dots, x_{i-1}; \theta)$$
 logits

$$p(X_i = w \mid x_1, \dots, x_{i-1}) = \frac{\exp(s(w)/\tau)}{\sum_{w' \in \mathcal{V}} \exp(s(w')/\tau)}$$

• $\tau = 1$: no changes to the probability distribution

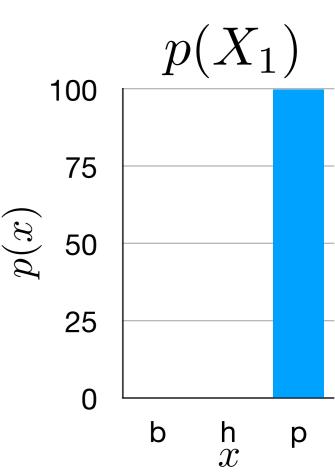




$$s(w) = f(w \mid x_1, \dots, x_{i-1}; \theta)$$
 logits

$$p(X_i = w \mid x_1, \dots, x_{i-1}) = \frac{\exp(s(w)/\tau)}{\sum_{w' \in \mathcal{V}} \exp(s(w')/\tau)}$$

- $\tau = 1$: no changes to the probability distribution
- τ → 0: relative probability assigned to highestprobability item in distribution increases
 - in practice, setting a temperature of 0
 recovers "argmax", putting all of the mass on
 the highest-probability item



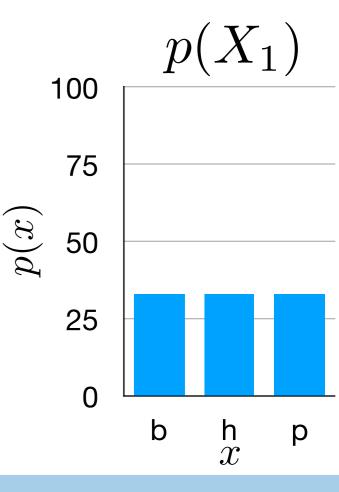


$$s(w) = f(w \mid x_1, \dots, x_{i-1}; \theta) \leftarrow logits$$



$$p(X_i = w \mid x_1, \dots, x_{i-1}) = \frac{\exp(s(w)/\tau)}{\sum_{w' \in \mathcal{V}} \exp(s(w')/\tau)}$$

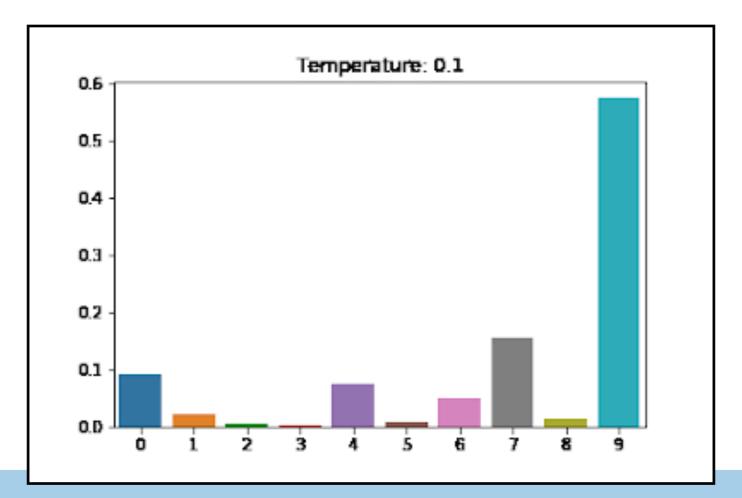
- $\tau = 1$: no changes to the probability distribution
- $\tau \rightarrow 0$: relative probability assigned to highestprobability item in distribution increases
- τ → ∞: distribution becomes more and more uniform





$$s(w) = f(w \mid x_1, \dots, x_{i-1}; \theta)$$
 logits

$$p(X_i = w \mid x_1, \dots, x_{i-1}) = \frac{\exp(s(w)/\tau)}{\sum_{w' \in \mathcal{V}} \exp(s(w')/\tau)}$$





$$s(w) = f(w \mid x_1, \dots, x_{i-1}; \theta)$$
 logits

$$p(X_i = w \mid x_1, \dots, x_{i-1}) = \frac{\exp(s(w)/\tau)}{\sum_{w' \in \mathcal{V}} \exp(s(w')/\tau)}$$

Ten	np = 0	Ter	np = 5
Once upon a time, there was a little		Once upon a time	, there was a young
little	100%	little	
$\begin{array}{c} \text{beautiful} \\ \text{young} \end{array}$	0%	Account to the second s	6 %
girl			5%
	0%	small	
king man	0% 0%		5% 5%
kingdom		kingdom	_
very	0%		4%
boy	0%		4%
princess	0%	princess	4%
great	0%	great	4%
		:	



- Temperature allows us to control the entropy of the output distribution without changing its relative ranking of items
- Higher temperature: closer to a uniform distribution
- Lower temperature: "peakier" distribution (in the limit, gives all probability mass to the most probable item)



Operation: find
$$\arg\max_{\overline{x}\in\mathcal{V}^+}p(\overline{x})$$

• Why is this hard?



Operation: find
$$\underset{\overline{x} \in \mathcal{V}^+}{\arg\max} p(\overline{x})$$

- Why is this hard?
- An approximation: greedy "sampling"

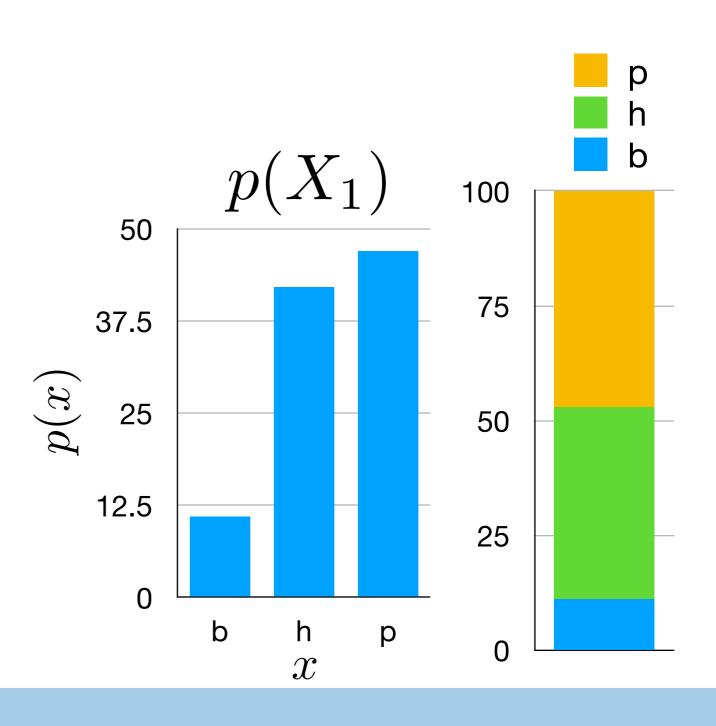
$$x_i \leftarrow \arg\max_{x \in \mathcal{V}} p(X_i \mid x_1, \dots, x_{i-1})$$

 Just choose the most probable wordtype at each generation step (no random sampling needed)



Operation: choose $x_i \leftarrow \arg\max_{x \in \mathcal{V}} p(X_i \mid x_1, \dots, x_{i-1})$

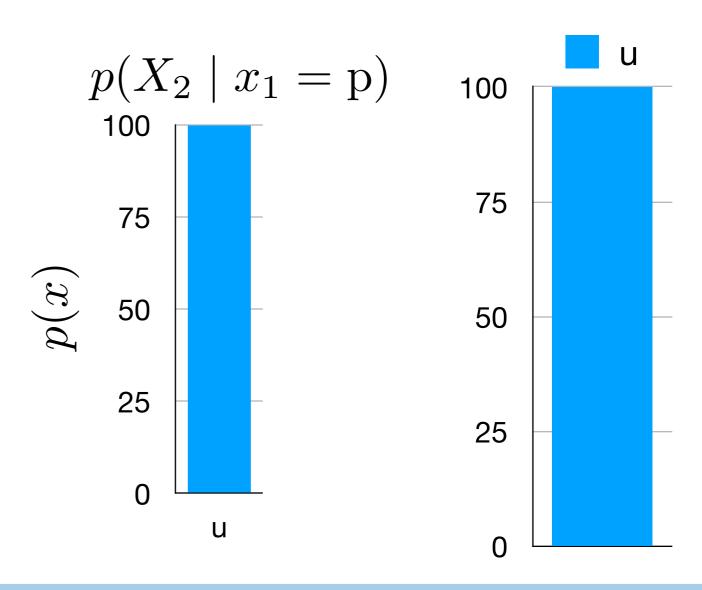
$$\overline{x} = \langle p \rangle$$





Operation: choose $x_i \leftarrow \arg\max_{x \in \mathcal{V}} p(X_i \mid x_1, \dots, x_{i-1})$

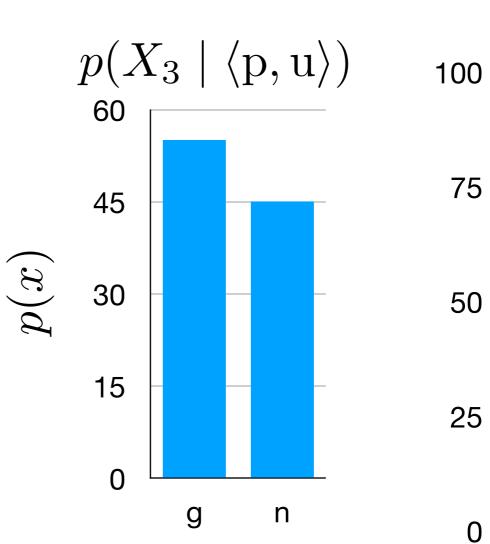
$$\overline{x} = \langle pu \rangle$$

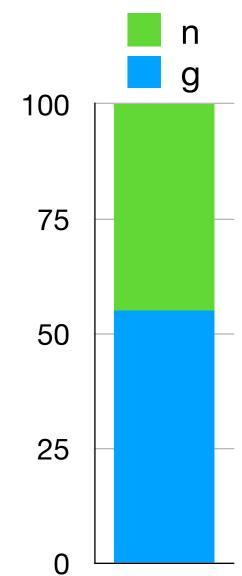




Operation: choose $x_i \leftarrow \arg\max_{x \in \mathcal{V}} p(X_i \mid x_1, \dots, x_{i-1})$

$$\overline{x} = \langle pug \rangle$$



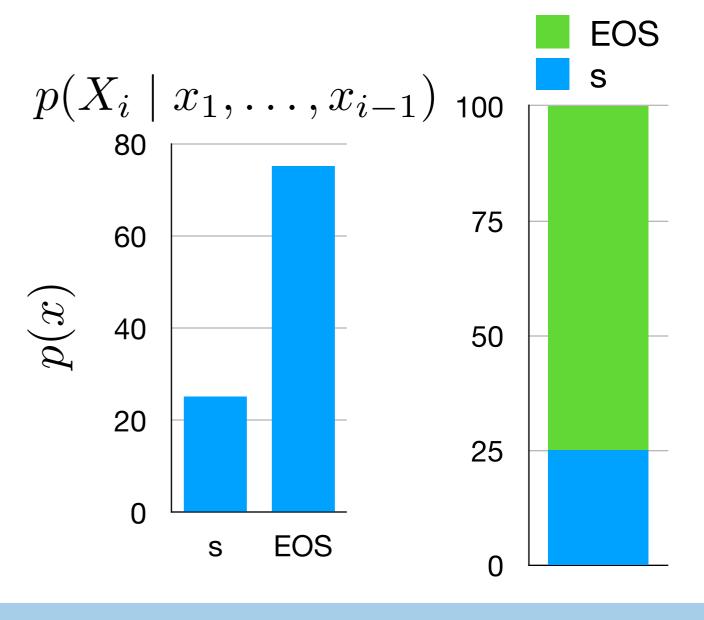




Operation: choose $x_i \leftarrow \arg\max_{x \in \mathcal{V}} p(X_i \mid x_1, \dots, x_{i-1})$

$$\overline{x} = \langle pugeos \rangle$$

pug





Operation: choose
$$x_i \leftarrow \arg\max_{x \in \mathcal{V}} p(X_i \mid x_1, \dots, x_{i-1})$$

- Why isn't this guaranteed to get us the highest-probability sequence?
- A better approximation for global argmax: beam search
 - During generation, we maintain a "beam" of *n* sequences instead of just one
 - At each generation step i,



Operation: choose
$$x_i \leftarrow \arg\max_{x \in \mathcal{V}} p(X_i \mid x_1, \dots, x_{i-1})$$

- Why isn't this guaranteed to get us the highest-probability sequence?
- A better approximation for global argmax: beam search
 - During generation, we maintain a "beam" of *n* sequences instead of just one
 - At each generation step *i*,
 - We select the n most likely next tokens \mathcal{X}_i for each prefix, and create n more sequences



Operation: choose
$$x_i \leftarrow \arg\max_{x \in \mathcal{V}} p(X_i \mid x_1, \dots, x_{i-1})$$

- Why isn't this guaranteed to get us the highest-probability sequence?
- A better approximation for global argmax: beam search
 - During generation, we maintain a "beam" of *n* sequences instead of just one
 - At each generation step *i*,
 - We select the n most likely next tokens \mathcal{X}_i for each prefix, and create n more sequences
 - Then we look at all the n^2 sequences so far, and discard all but the n most likely sequences



Operation: choose
$$x_i \leftarrow \arg\max_{x \in \mathcal{V}} p(X_i \mid x_1, \dots, x_{i-1})$$

- Why isn't this guaranteed to get us the highest-probability sequence?
- A better approximation for global argmax: beam search
 - During generation, we maintain a "beam" of *n* sequences instead of just one
 - At each generation step *i*,
 - We select the n most likely next tokens \mathcal{X}_i for each prefix, and create n more sequences
 - Then we look at all the n^2 sequences so far, and discard all but the n most likely sequences
 - At the end, we select the sequence that has the highest probability among the set



```
p(X_1)
the
0.4
```

0.3

and 0.2

every

way0.1



```
p(X_1)
the
0.4
```

in 0.3

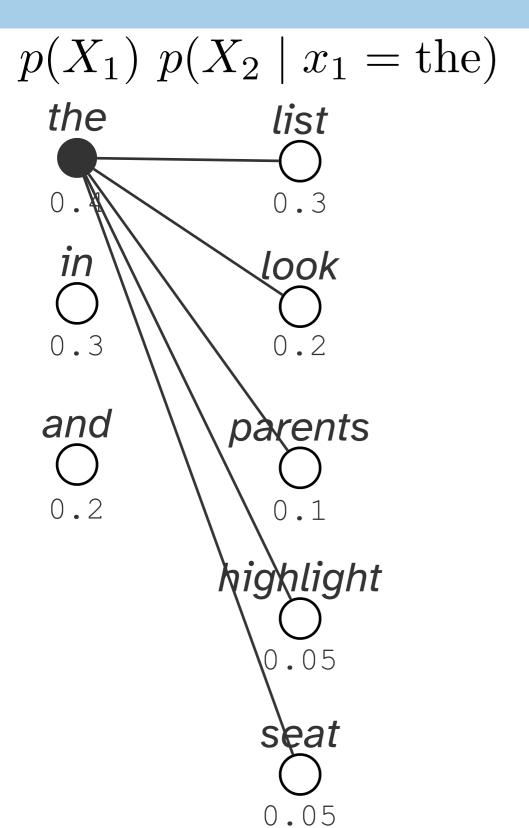
and 0.2

Discard all but the top 3 continuations

E
Q
(1)
M

Prefix	Probability
the	0.4
in	0.3
and	0.2





Continue with beam item #1 as prefix

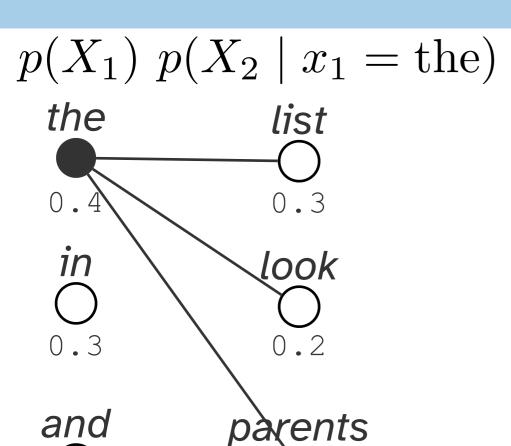
E	
Q	
3)

Prefix	Probability
the	0.04
in	0.02
and	0.02



Discard all but the

top 3 continuations



Current Beam Candidates

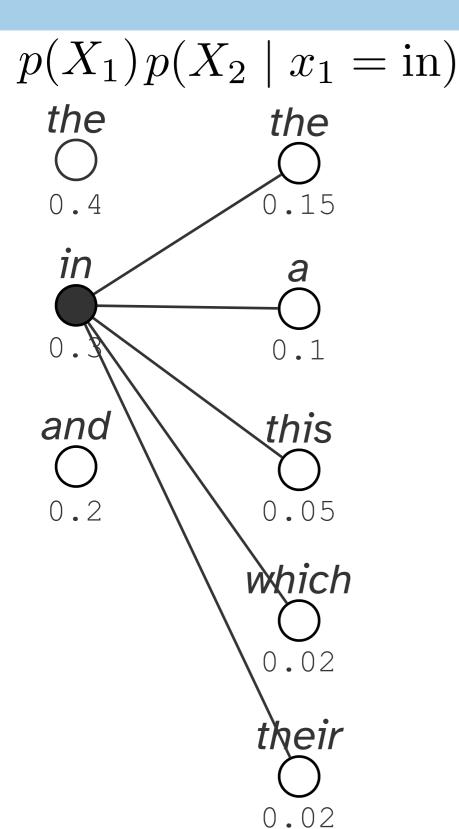
Prefix	Probability
the list	0.4*0.3
the look	0.3*0.2
the parents	0.2*0.1

Prefix	Probability
the	0.04
in	0.02
and	0.02



Continue with beam

item #2 as prefix



Current Beam Candidates

Prefix	Probability
the list	0.4*0.3
the look	0.3*0.2
the parents	0.2*0.1

Prefix	Probability
the	0.04
in	0.02
and	0.02



Discard all but the

top 3 continuations

$$p(X_1)p(X_2 \mid x_1 = \text{in})$$

the the O 0.15

in a 0.1

this

0.05

and

Current Beam Candidates

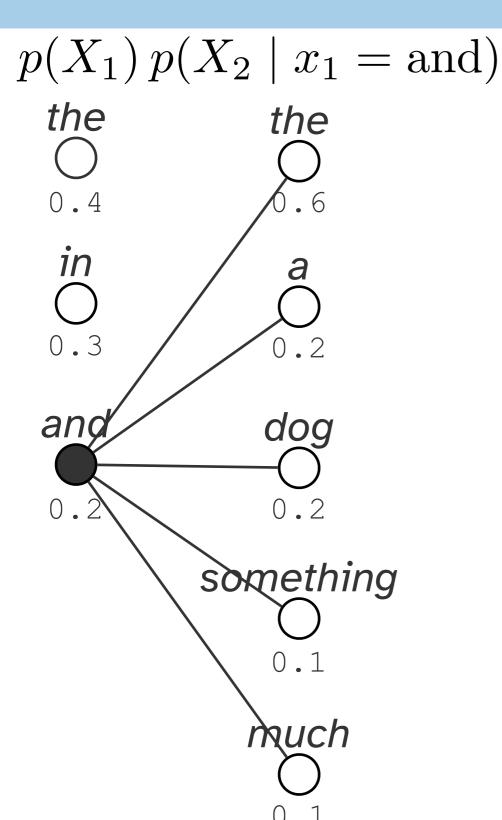
Prefix Probability the list 0.4*0.3 the look 0.3*0.2 the parents 0.2*0.1 in the 0.3*0.15 in a 0.3*0.1 in this 0.3*0.05

Prefix	Probability
the	0.04
in	0.02
and	0.02



Continue with beam

item #3 as prefix



Current Beam Candidates

Prefix	Probability
the list	0.4*0.3
the look	0.3*0.2
the parents	0.2*0.1
in the	0.3*0.15
in a	0.3*0.1
in this	0.3*0.05

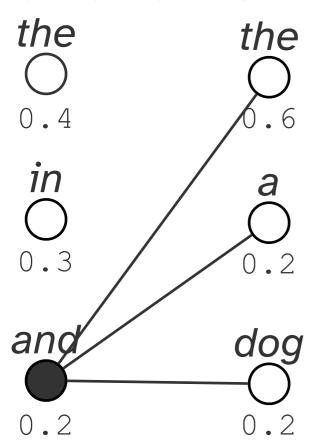
Prefix	Probability
the	0.04
in	0.02
and	0.02



Discard all but the

top 3 continuations

$$p(X_1) p(X_2 \mid x_1 = \text{and})$$



Current Beam Candidates

andidates

Prefix	Probability
the list	0.4*0.3
the look	0.3*0.2
the parents	0.2*0.1
in the	0.3*0.15
in a	0.3*0.1
in this	0.3*0.05
and the	0.2*0.6
and a	0.2*0.2
and dog	0.2*0.2

Prefix	Probability
the	0.04
in	0.02
and	0.02



Compute probabilities of all candidates

Current Beam Candidates

Prefix	Probability
the list	0.12
the look	0.06
the parents	0.02
in the	0.05
in a	0.03
in this	0.02
and the	0.12
and a	0.04
and dog	0.04

Prefix	Probability
the	0.04
in	0.02
and	0.02



Refresh the beam with top *n* candidates

Current Beam Candidates

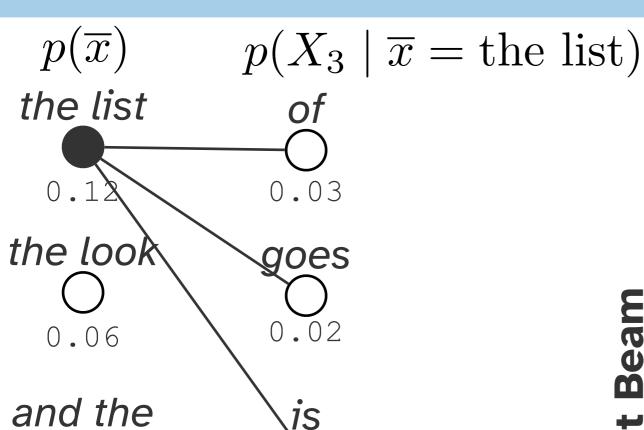
Prefix	Probability
the list	0.12
the look	0.06
the parents	0.02
in the	0.05
in a	0.03
in this	0.02
and the	0.12
and a	0.04
and dog	0.04

Prefix	Probability
the list	0.12
the look	0.06
and the	0.12



Keep generating

with new beam



0.12

ent Beam	ndidates
Curre	Cano

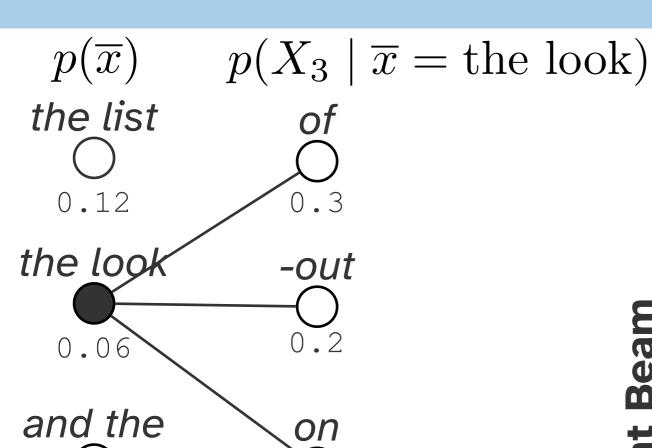
Prefix	Probability
the list of	0.004
the list goes	0.002
the list is	0.002

Prefix	Probability
the list	0.12
the look	0.06
and the	0.12



Keep generating

with new beam



0.12

Surrent Beam Candidates

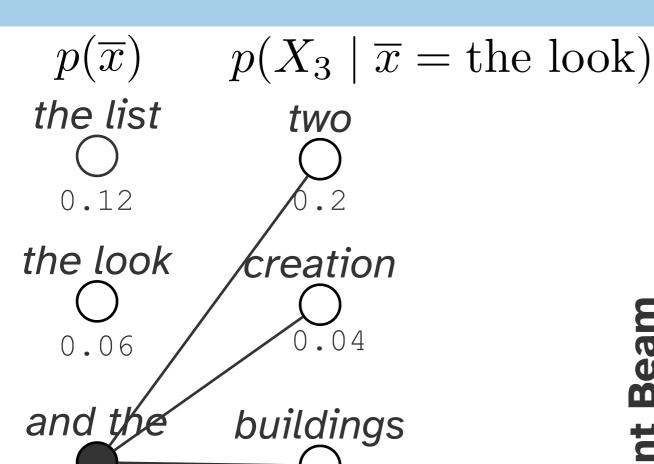
Prefix	Probability
the list of	0.004
the list goes	0.002
the list is	0.002
the look of	0.018
the lookout	0.012
the look on	0.006

Prefix	Probability
the list	0.12
the look	0.06
and the	0.12



Keep generating

with new beam



0.02

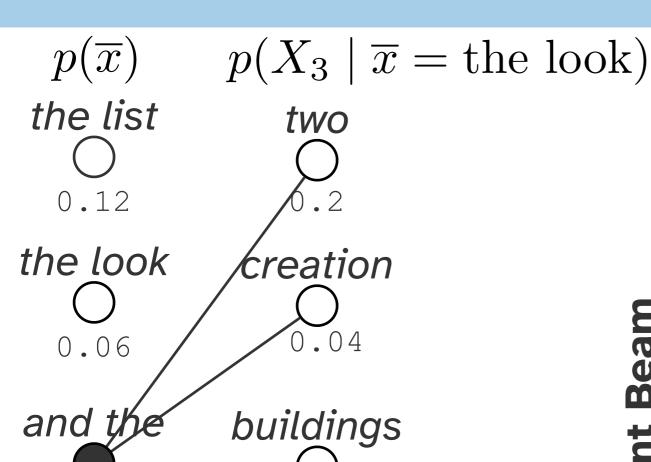
0.12

Current Beam Candidates

Prefix	Probability
the list of	0.004
the list goes	0.002
the list is	0.002
the look of	0.018
the lookout	0.012
the look on	0.006
and the two	0.024
and the creation	0.005
and the buildings	0.002

Prefix	Probability
the list	0.12
the look	0.06
and the	0.12





0.02

0.12

Refresh the beam with top *n* candidates

Current Beam Candidates

Probability
0.004
0.002
0.002
0.018
0.012
0.006
0.024
0.005
0.002

Prefix	Probability
the look of	0.018
the lookout	0.012
and the two	0.024



- How do we know when to stop?
 - When all of the items in the beam have EOS (we don't expand these prefixes, just keep them around for the end)
 - Or, when we've reached a maximum sequence length
- Let's say we're done sampling at this point
- We'll select the sequence with the highest probability in the beam

Prefix	Probability
the look of	0.018
the lookout	0.012
and the two	0.024



- How do we know when to stop?
 - When all of the items in the beam have EOS (we don't expand these prefixes, just keep them around for the end)
 - Or, when we've reached a maximum sequence length
- Let's say we're done sampling at this point
- We'll select the sequence with the highest probability in the beam

Prefix	Probability
the look of	0.018
the lookout	0.012
and the two	0.024



- How do we know when to stop?
 - When all of the items in the beam have EOS (we don't expand these prefixes, just keep them around for the end)
 - Or, when we've reached a maximum sequence length
- Let's say we're done sampling at this point
- We'll select the sequence with the highest probability in the beam
- What if our sequences have different lengths?

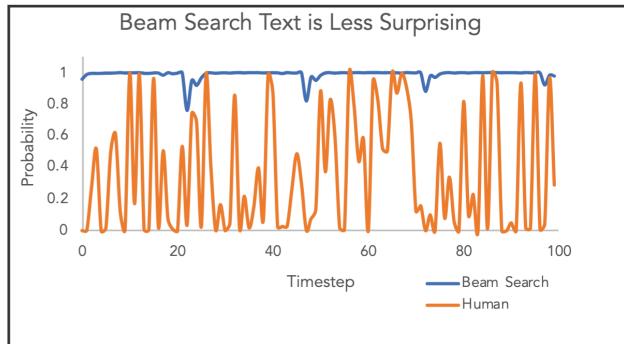
Length normalization:
$$p(\overline{x}) \propto -\frac{1}{|\overline{x}|^{\alpha}} \sum_{i=1}^{|x|} \log p(x_i \mid x_1, \dots, x_{i-1})$$



 Should we always try to approximate the argmax?



- Should we always try to approximate the argmax?
- Maybe not!
- Argmax produces repetitive, less diverse, and overall tooprobable output sequences
- What's missing?



Beam Search

...to provide an overview of the current state-of-the-art in the field of computer vision and machine learning, and to provide an overview of the current state-of-the-art in the field of computer vision and machine learning, and to provide an overview of the current state-of-the-art in the field of computer vision and machine learning, and to provide an overview of the current state-of-the-art in the field of computer vision and machine learning, and...

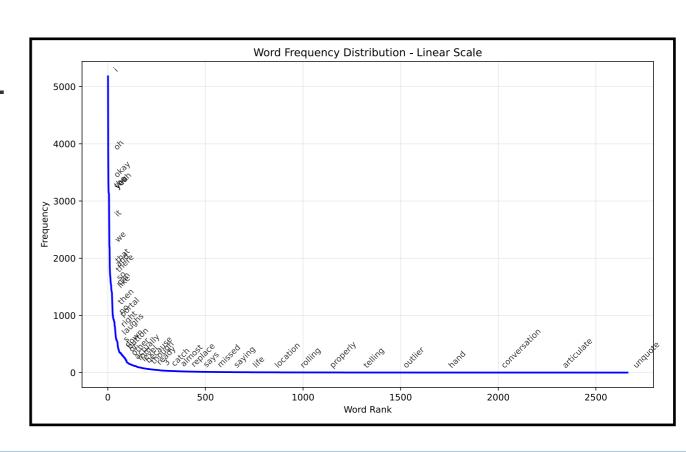
Human

...which grant increased life span and three years warranty. The Antec HCG series consists of five models with capacities spanning from 400W to 900W. Here we should note that we have already tested the HCG-620 in a previous review and were quite satisfied With its performance. In today's review we will rigorously test the Antec HCG-520, which as its model number implies, has 520W capacity and contrary to Antec's strong beliefs in multi-rail PSUs is equipped...

Holtzman et al. 2019



- Should we always try to approximate the argmax?
- Maybe not!
- Argmax produces repetitive, less diverse, and overall tooprobable output sequences
- What's missing?
- Long tail!

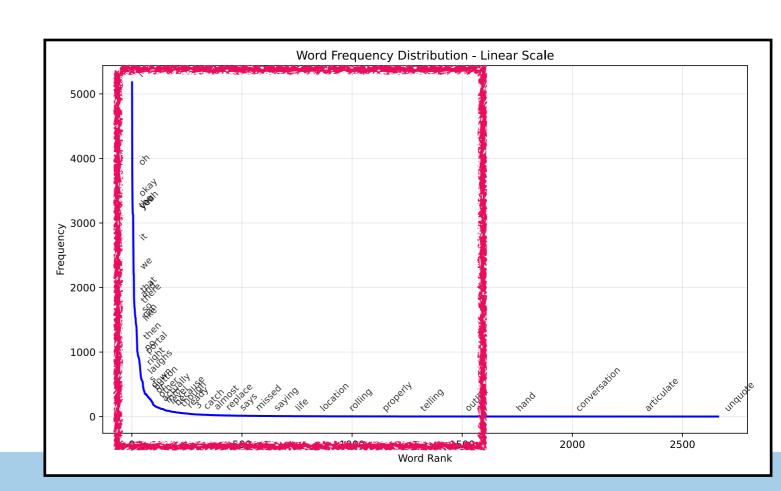




Operation: & sampling

- Identify the set of n tokens \mathcal{E} such that $\forall x \in \mathcal{E}$, $p(x \mid x_1, \dots, x_{i-1}) \geq \epsilon$
- Set the probabilities of all but these tokens to 0
- Renormalize by dividing remaining probabilities by

$$\sum_{x \in \mathcal{E}} p(X_i \mid x_1, \dots, x_{i-1})$$

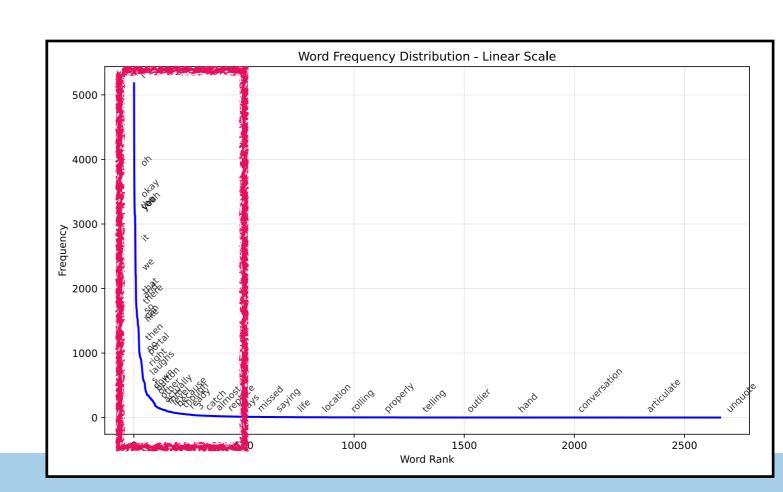




Operation: top-k sampling

- Identify the set of k tokens $\mathcal K$ that have the highest probabilities under $p(X_i \mid x_1, \dots, x_{i-1})$
- Set the probabilities of all but these tokens to 0
- Renormalize by dividing remaining probabilities by

$$\sum_{x \in \mathcal{K}} p(X_i \mid x_1, \dots, x_{i-1})$$

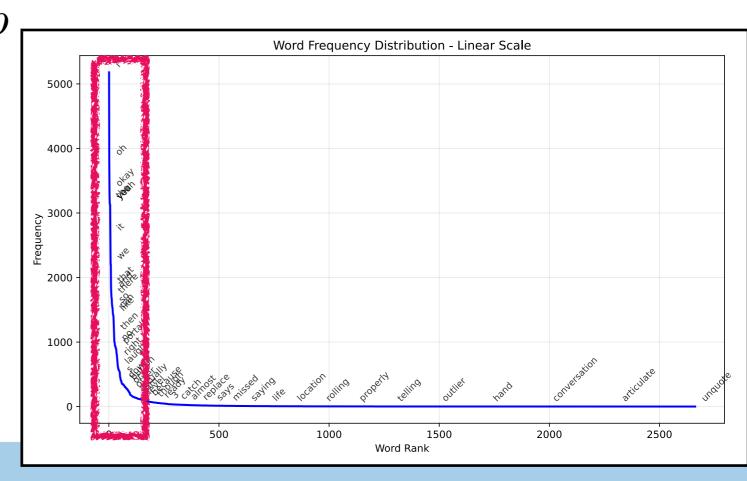




Operation: top-p (nucleus) sampling

- Identify the set of n tokens $\mathcal P$ that have the highest probabilities under $p(X_i \mid x_1, \dots, x_{i-1})$ and their cumulative probability is p
- Set the probabilities of all but these tokens to 0
- Renormalize by dividing remaining probabilities by

$$\sum_{x \in \mathcal{P}} p(X_i \mid x_1, \dots, x_{i-1}) = p$$





Operation: constrained decoding

- For some tasks, we have additional information about what wordtypes can or cannot be next
- E.g., in code generation, I can't generate more) than I have (
- While modern LLMs can learn these patterns from data at scale, it can sometimes still be useful to constrain our output space



Operation: constrained decoding

- For some tasks, we have additional information about what wordtypes can or cannot be next
- E.g., in code generation, I can't generate more) than I have (
- While modern LLMs can learn these patterns from data at scale, it can sometimes still be useful to constrain our output space
- Similar to before: given a set of possible continuations $\mathcal{C} \subseteq \mathcal{V}$ we will set the probabilities of all other tokens to 0, then renormalize using $\sum p(X_i \mid x_1, \dots, x_{i-1})$