





Ct After debating whether to bow to the **king** or the **woman** first, the jester decided on the



Architecture Enhancement



Figure 2: Comparison between different architectures. The **#S**, **#I**, and **#P** are the number of softmaxes, input hidden states, and partitions, respectively. The green boxes refer to embeddings/vectors. The vocab means the embeddings of all words in the vocabulary. \oplus refers to concatenation. L^h , L^f , and L^{π} are linear projection layers.

Softmax Bottleneck Makes Language Models **Unable to Represent Multi-mode Word Distributions** Haw-Shiuan Chang and Andrew McCallum

Theoretical Analysis Ideal next word probability **Theorem 1** (simplified): If many word embeddings are linearly dependent, the softmax in a LM cannot rank the words arbitrarily **woman** 0.4 **king** 0.4 lady 0.03 **Example**: If "*woman* - *man* = *queen* - *king*", GPT-2 cannot rank the word woman and king as the top 2 words $= \frac{\exp(\boldsymbol{h}_{c_t}^T \boldsymbol{w}_x)}{\sum_{x'} \exp(\boldsymbol{h}_{c_t}^T \boldsymbol{w}_{x'})}$ **Example**: If "<u>UMass</u> = 0.2 <u>University</u> + 0.2 <u>Massachusetts</u>", GPT-2 cannot rank a rare word <u>UMass</u> on top of the similar popular words Dot product <u>University</u> and <u>Massachusetts</u> (Demeter et al., 2020). **Linear Algebra Intuition**: N+1 words are linear dependent **→** They are in subspace with $d < N \rightarrow$ cannot have arbitrary probabilities • 1 \underline{w}_{king} - 1 \underline{w}_{queen} = 1 \underline{w}_{man} - 1 \underline{w}_{woman} • $1 \underline{w_{king}} + 1 \underline{w_{woman}} = 1 \underline{w_{queen}} + 1 \underline{w_{man}}$ × \underline{h}^{T} (hidden state) on both side Learning analogical word embedding structure Generalization • 1 $\underline{h}^T \underline{w}_{king}$ + 1 $\underline{h}^T \underline{w}_{woman} \neq$ 1 $\underline{h}^T \underline{w}_{queen}$ + 1 $\underline{h}^T \underline{w}_{man}$ • If $\exists \underline{h}$, s.t min($\underline{h}^T \underline{w}_{king}$, $\underline{h}^T \underline{w}_{woman}$) > max($\underline{h}^T \underline{w}_{queen}$, $\underline{h}^T \underline{w}_{man}$) • 1 $\underline{h}^T \underline{w}_{king} + 1 \underline{h}^T \underline{w}_{woman} \ge$ 2 min($\underline{h}^T \underline{w}_{king}, \underline{h}^T \underline{w}_{woman}$) > Contradict $2 \max(\underline{h}^T \underline{w}_{queen}, \underline{h}^T \underline{w}_{man}) >$ $1 \underline{h}^T \underline{w}_{queen} + 1 \underline{h}^T \underline{w}_{man} (\rightarrow \leftarrow)$ Ranking next word • Thus, the logits of LM cannot rank both king and woman arbitrarily on top of <u>queen</u> and <u>man</u> Improvements over Yang et al., (2018) • Serious among which words? • Affect the top words? If yes, when? • Disappears after making D>V? (d) Multi-facet Softmax (Ours) **Theorem 2** (simplified): If many word embeddings are approximately linearly dependent and the magnitude of the hidden state has a upperbound, the softmax in a LM cannot assign very small probabilities Softmax Softmax to some words dot[product **Example**: If "*woman* + *king* = *queen* + *man* + $\underline{\varepsilon}$ ", GPT-2 cannot make the logits of *queen* and *man* much smaller than the logits of *king* and *woman* **Example**: If "*woman* = *man* + $\underline{\varepsilon}$ ", GPT-2 cannot make the logits of <u>man</u> GELU(*L^h*(.)) much smaller than the logits of *woman* put Hidden States (#1 **Intuition**: $\underline{h}^T \underline{king} + \underline{h}^T \underline{woman} = \underline{h}^T \underline{queen} + \underline{h}^T \underline{man} + \underline{h}^T \underline{\varepsilon}$, and we can ignore $\underline{h}^T \underline{\varepsilon}$ if $||\underline{h}||$ and $||\underline{\varepsilon}||$ are both small

- Linearly dependent among $\{\underline{w}_{l_1}, \ldots, \underline{w}_{l_L}, \underline{w}_{r_1}, \ldots, \underline{w}_{r_R}\}$
- $a_{l_1} \underline{W}_{l_1} + \ldots + a_{l_L} \underline{W}_{l_L} = a_{r_1} \underline{W}_{r_1} + \ldots + a_{r_R} \underline{W}_{r_R}$
 - All coefficient $a_{li} > 0, a_{ri} > 0$
 - WLOG $a_{l_1} + \ldots + a_{l_L} \ge a_{r_1} + \ldots + a_{r_R}$
- $a_{l_1} \underline{h}^T \underline{w}_{l_1} + \ldots + a_{l_L} \underline{h}^T \underline{w}_{l_L} = a_{r_1} \underline{h}^T \underline{w}_{r_1} + \ldots + a_{r_R} \underline{h}^T \underline{w}_{r_R}$
- If $\exists \underline{h}$, s.t min_i($\underline{h}^T w_{li}$) > max_i($\underline{h}^T w_{Ri}$)
 - $a_{l_1} \underline{h}^T \underline{w}_{l_1} + \ldots + a_{l_L} \underline{h}^T \underline{w}_{l_L} \ge$ $(a_{l_1} + \ldots + a_{l_k}) \min_i (\underline{h}^T \underline{w}_{l_i}) >$ $(a_{r_1} + \ldots + a_{r_R}) \max_{i} (\underline{h}^T \underline{w}_{R_i}) \geq 0$ $a_{r_1} \underline{h}^T \mathbf{w}_{r_1} + \ldots + a_{r_R} \underline{h}^T \mathbf{w}_{r_R} (\rightarrow \leftarrow)$
- Thus, the logits of LM cannot rank all the left words on top of the right words.
- -> Among words in a small subspace
- -> Yes. When the ideal distribution is multi-mode
- -> No, if some words are in a small subspace

	Μ
Only adding	Softm
nonlinearity is not	→ SigSoftmax (
enough (Parthiban et	Softmax
al., 2021)	Softmax +
	MoS (Yang
	MoS (Yang
	DOC (Tak
	MFS w/o
	MFS w/
	MF
	Table 1: Perple K), input hidde are the test set p smaller than 0.0

Improvement of MFS over Softmax is around 15% between GPT-2 Small and GPT-2 Medium (with 3x parameters)

$\text{Corpus} \rightarrow$	Op
	The Elastic
Input Context	Elastic SIEM s
	this post are not
Softmax (GPT-2)	the 0.087, 1
MFS (Ours)	Elastic 0.220
MFS Softmax 1	end 0.051, the
MFS Softmax 2	Elastic 0.652
MFS Softmax 3	the 0.193
1	1





Theory

Multi-mode distribution must exist if some word embeddings are in a small subspace Stolen Probability (Demeter et al., 2020) Generalizatio Multi-mode Distribution (Ours) Model." In ICLR. 2018.

[2] David Demeter, Gregory Kimmel, and Doug Downey. Stolen probability: A structural weakness of neural language models. In ACL. 2020 3] Sekitoshi Kanai, Yasuhiro Fujiwara, Yuki Yamanaka, and Shuichi Adachi. Sigsoftmax: Reanalysis of the softmax bottleneck. In NeurIPS 2018 [4] Dwarak Govind Parthiban, Yongyi Mao, and Diana Inkpen. On the softmax bottleneck of recurrent language models. In AAAI 2021 [5] Sho Takase, Jun Suzuki, and Masaaki Nagata. 2018. Direct output connection for a high-rank language model. In EMNLP 2018 [6] Sharan Narang, Hyung Won Chung, Yi Tay, Liam Fedus, Thibault Févry, Michael Matena, Karishma Malkan et al. "Do Transformer Modifications Transfer Across Implementations and Applications?." In EMNLP 2021





Empirical Analysis

	Cor	nfigur	ation	GPT-2 Small				GPT-2 Medium				
Models↓	#S	#I	#P	Size	Time	OWT	Wiki	Size	Time	OWT	Wiki	
max (GPT-2)	1	1	1	163.6M	84ms	18.72	24.06	407.3M	212ms	15.89	20.34	
(Kanai et al., 2018)	1	1	1	163.6M	91ms	18.63	24.06	407.3M	221ms	16.07	20.65	
x + Multi-input	1	9	1	169.5M	87ms	18.50	23.89	417.8M	219ms	15.76	20.29	
+ Multi-partition	1	1	4	165.4M	88ms	18.77	24.08	410.5M	218ms	15.89	20.30	
ng et al., 2018) (4)	4	1	1	165.4M	152ms	18.61	23.77	410.5M	299ms	15.75	20.08	
ng et al., 2018) (3)	3	1	1	164.8M	130ms	18.63	23.81	409.4M	270ms	15.79	20.11	
kase et al., 2018)	3	3	1	164.8M	130ms	18.69	24.02	409.4M	270ms	15.88	20.34	
o Multi-partition	3	9	1	171.9M	133ms	18.37	23.56	422.0M	276ms	15.65	20.06	
v/o Multi-input	3	1	4	166.6M	134ms	18.60	23.72	412.6M	275ms	15.71	20.08	
IFS (Ours)	3	9	4	175.4M	138ms	18.29	23.45	428.3M	283ms	15.64	20.02	

Multiple input nidden states heli

Multiple partitions

exity comparison between MFS (Ours) and baselines. #S, #I, #P are the number of softmaxes (i.e. len states, and partitions, respectively. The top four baselines use a single softmax. OWT and Wiki perplexity of OpenWebText and Wikipedia 2021, respectively. The standard errors of all models are .02 perplexity. We also compare the number of parameters and the inference time on one batch.



References

[1] Zhilin Yang, Zihang Dai, Ruslan Salakhutdinov, and William W. Cohen. "Breaking the Softmax Bottleneck: A High-Rank RNN Language