# Breaking the Softmax Bottleneck via Monotonic Functions

Octavian Ganea, Sylvain Gelly, Gary Bécigneul, Aliaksei Severyn

# Softmax Layer (for Language Models)
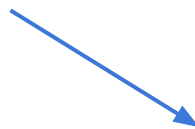
- Natural language as conditional distributions

- Parametric distributions & softmax:

next word        context

$$P_\theta(x|c) = \frac{\exp \mathbf{h}_c^\top \mathbf{w}_x}{\sum_{x'} \exp \mathbf{h}_c^\top \mathbf{w}_{x'}} \approx P^*(x|c)$$

# Softmax Layer (for Language Models)

- Natural language as conditional distributions

next word          context

- Parametric distributions & softmax:

$$P_\theta(x|c) = \frac{\exp \mathbf{h}_c^\top \mathbf{w}_x}{\sum_{x'} \exp \mathbf{h}_c^\top \mathbf{w}_{x'}} \approx P^*(x|c)$$
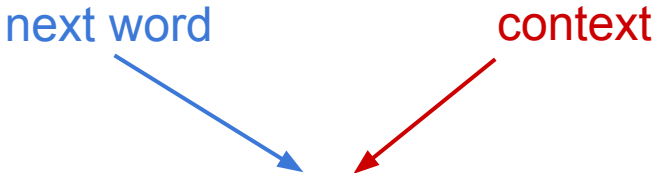
- Challenge: Can we always find $\theta$ s.t. for all $c$ : $\quad P_\theta(X|c) = P^*(X|c) \quad$ ?

# Softmax Layer (for Language Models)

- Natural language as conditional distributions

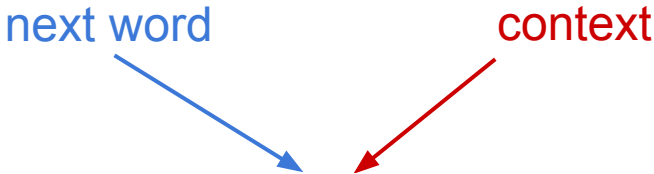next word          context

- Parametric distributions & softmax:

$$P_\theta(x|c) = \frac{\exp \mathbf{h}_c^\top \mathbf{w}_x}{\sum_{x'} \exp \mathbf{h}_c^\top \mathbf{w}_{x'}} \approx P^*(x|c)$$

- Challenge: Can we always find $\theta$ s.t. for all $c$ : $P_\theta(X|c) = P^*(X|c)$ ?

**No**, when embedding size < label cardinality (vocab size) !

# What is the Softmax Bottleneck (Yang et al, '18) ?

- **log-P matrix**:
$$\mathbf{A}_P = \begin{bmatrix} \log P(x_1|c_1) & \log P(x_2|c_1) & \dots & \log P(x_M|c_1) \\ \log P(x_1|c_2) & \log P(x_2|c_2) & \dots & \log P(x_M|c_2) \\ \vdots & \vdots & \ddots & \vdots \\ \log P(x_1|c_N) & \log P(x_2|c_N) & \dots & \log P(x_M|c_N) \end{bmatrix} \in \mathbb{R}^{N \times M}$$

**Label cardinality = Vocabulary size**

Breaking the Softmax Bottleneck: A High-Rank RNN Language Model, Yang et al., ICLR 2018

# What is the Softmax Bottleneck (Yang et al, '18) ?

- **log-P matrix**: $\mathbf{A}_P = \begin{bmatrix} \log P(x_1|c_1) & \log P(x_2|c_1) & \dots & \log P(x_M|c_1) \\ \log P(x_1|c_2) & \log P(x_2|c_2) & \dots & \log P(x_M|c_2) \\ \vdots & \vdots & \ddots & \vdots \\ \log P(x_1|c_N) & \log P(x_2|c_N) & \dots & \log P(x_M|c_N) \end{bmatrix} \in \mathbb{R}^{N \times M}$

**Number of labels = Vocabulary size**

- Then: $\text{rank}(A_{P_\Theta}) \leq d + 1$

Breaking the Softmax Bottleneck: A High-Rank RNN Language Model, Yang et al., ICLR 2018

# What is the Softmax Bottleneck (Yang et al, '18) ?

- **log-P matrix**:
$$\mathbf{A}_P = \begin{bmatrix} \log P(x_1|c_1) & \log P(x_2|c_1) & \ldots & \log P(x \\ \log P(x_1|c_2) & \log P(x_2|c_2) & & \\ \vdots & & & \\ \log P( & & & \end{bmatrix}$$

- Then:

But $\mathbf{A}_{P^*}$ is likely full-rank, so $\mathbf{A}_{P^*} \neq \mathbf{A}_{P\Theta}$ when $d << min(M, N)$
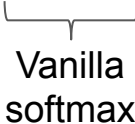
Breaking the Softmax Bottleneck: A High-Rank RNN Language Model, Yang et al., ICLR 2018

# Breaking the Softmax Bottleneck [1]

- MoS [1] : Mixture of **K** Softmaxes

[1] Breaking the Softmax Bottleneck: A High-Rank RNN Language Model, Yang et al., ICLR 2018

# Breaking the Softmax Bottleneck [1]

- MoS [1] : Mixture of **K** Softmaxes
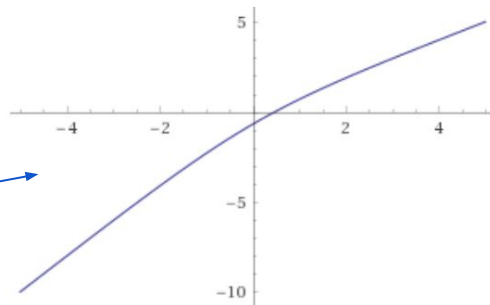
- Improves perplexity

[1] Breaking the Softmax Bottleneck: A High-Rank RNN Language Model, Yang et al., ICLR 2018

# Breaking the Softmax Bottleneck [1]

- MoS [1] : Mixture of **K** Softmaxes

- Improves perplexity

- Slower than vanilla softmax: 2 - 6.4x
- GPU Memory: M x N x **K** tensor

  Vanilla
  softmax

[1] Breaking the Softmax Bottleneck: A High-Rank RNN Language Model, Yang et al., ICLR 2018

# Breaking the Softmax Bottleneck [2]

- Sig-Softmax [2] :

$$\text{softmax}(2\mathbf{y} - \log(1 + \exp(\mathbf{y})))$$

[2] Sigsoftmax: Reanalysis of the Softmax Bottleneck, S. Kanai et al., NIPS 2018
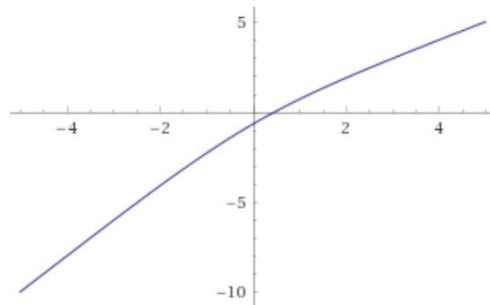
# Breaking the Softmax Bottleneck [2]



- Sig-Softmax [2] :

$$\mathrm{softmax}(2\mathbf{y} - \log(1 + \exp(\mathbf{y})))$$
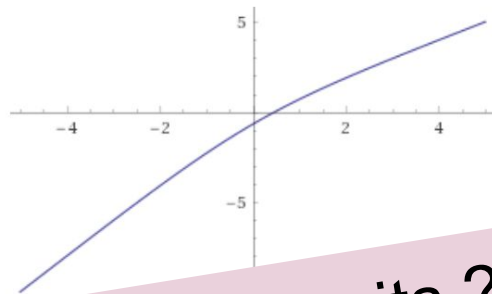
- Small improvement over vanilla Softmax

[2] Sigsoftmax: Reanalysis of the Softmax Bottleneck, S. Kanai et al., NIPS 2018

# Breaking the Softmax Bottleneck [2]

- Sig-Softmax [2] :

$$\text{softmax}(2\mathbf{y} - \log(1 + e^{...}$$

Can we learn the best non-linearity to deform the logits ?

[2] Sigsoftmax: Reanalysis of the Softmax Bottleneck, S. Kanai et al., NIPS 2018

# Can we do better ?

- Our idea - learn a pointwise monotonic function on top of logits:

$$p(y_i) = \frac{\exp(f(y_i))}{\sum_j \exp(f(y_j))}, \ \text{i.e. softmax}(f(\mathbf{y}))$$

# Can we do better ?

- Our idea - learn a pointwise monotonic function on top of logits:

$$p(y_i) = \frac{\exp(f(y_i))}{\sum_j \exp(f(y_j))}, \text{ i.e. } \mathrm{softmax}(f(\mathbf{y}))$$

$f : \mathbb{R} \to \mathbb{R}$ should be:

# Can we do better ?

- Our idea - learn a pointwise monotonic function on top of logits:

$$p(y_i) = \frac{\exp(f(y_i))}{\sum_j \exp(f(y_j))}, \text{ i.e. softmax}(f(\mathbf{y}))$$

$f : \mathbb{R} \to \mathbb{R}$ should be:

1.  With **unbounded image set** -- to model sparse distributions

# Can we do better ?

- Our idea - learn a pointwise monotonic function on top of logits:

$$p(y_i) = \frac{\exp(f(y_i))}{\sum_j \exp(f(y_j))}, \text{ i.e. } \mathrm{softmax}(f(\mathbf{y}))$$

$f : \mathbb{R} \to \mathbb{R}$ should be:

1. With **unbounded image set** -- to model sparse distributions
2. **Continuous** and **(piecewise) differentiable** -- for backprop

# Can we do better ?

- Our idea - learn a pointwise monotonic function on top of logits:

$$p(y_i) = \frac{\exp(f(y_i))}{\sum_j \exp(f(y_j))}, \text{ i.e. softmax}(f(\mathbf{y}))$$

$f : \mathbb{R} \rightarrow \mathbb{R}$ should be:

1. With **unbounded image set** -- to model sparse distributions
2. **Continuous** and **(piecewise) differentiable** -- for backprop
3. **Non-linear** -- to break the softmax bottleneck

# Can we do better ?

- Our idea - learn a pointwise monotonic function on top of logits:

$$p(y_i) = \frac{\exp(f(y_i))}{\sum_j \exp(f(y_j))}, \text{ i.e. softmax}(f(\mathbf{y}))$$

$f : \mathbb{R} \rightarrow \mathbb{R}$ should be:

1. With **unbounded image set** -- to model sparse distributions
2. **Continuous** and **(piecewise) differentiable** -- for backprop
3. **Non-linear** -- to break the softmax bottleneck
4. **Monotonic** -- to preserve the ranking of logits

# Can we do better ?

- Our idea - learn a pointwise monotonic function on top of logits:

$$p(y_i) = \frac{\exp(f(y_i))}{\sum_j \exp(f(y_j))}, \text{ i.e. } \mathrm{softmax}(f(\mathbf{y}))$$

$f : \mathbb{R} \to \mathbb{R}$ should be:

1. With **unbounded image set** -- to model sparse distributions
2. **Continuous** and **(piecewise) differentiable** -- for backprop
3. **Non-linear** -- to break the softmax bottleneck
4. **Monotonic** -- to preserve the ranking of logits
5. **Fast** and **memory efficient** -- comparable w/ vanilla softmax

# Can we do better ?

- Our idea - learn a pointwise monotonic function on top of logits:

$$p(y_i) = \frac{\exp(f(y_i))}{\sum_j \exp(f(y_j))}, \text{ i.e. softmax}(f(\mathbf{y}))$$

$f : \mathbb{R} \to \mathbb{R}$ should be:

1. With _____ -- to model sparse distributions
2. _____ and **(piecewise) differentiable** -- for backprop
3. _____ **linear** -- to break the softmax bottleneck
4. **Monotonic** -- to preserve the ranking of logits
5. **Fast** and **memory efficient** -- comparable w/ vanilla softmax

**Theorem**: these properties are not restrictive in terms of rank deficiency

# Learnable parametric monotonic real functions

- A neural network with 1 hidden layer and positive (constrained) weights [3]

$$f(x) = \sum_{i=1}^{K} v_i \sigma(u_i x + b_i) + b, \ \text{s.t.} \ v_i, u_i \geq 0$$

- Universal approximator for all monotonic functions (when K is large enough !)

[3] Monotone and Partially Monotone Neural Networks, Daniels and Velikova, 2010, IEEE TRANSACTIONS ON NEURAL NETWORKS

# Synthetic Experiment

- **Goal**: *separate softmax bottleneck from context embedding bottleneck*
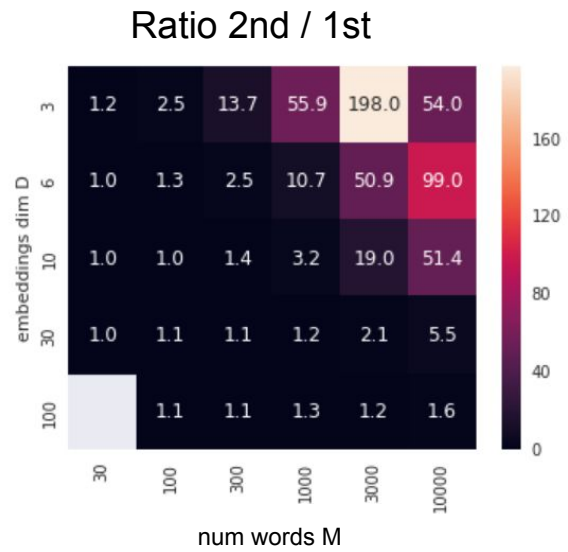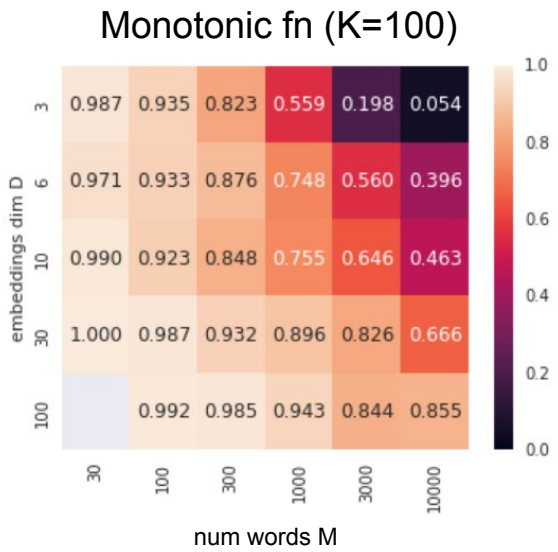
# Synthetic Experiment

- **Goal**: *separate softmax bottleneck from context embedding bottleneck*
- Sample N different categorical distributions with M outcomes:

$$P^*(\cdot|c_j) \sim \mathrm{Dir}(\alpha), \text{ for } j = 1, \ldots, N$$
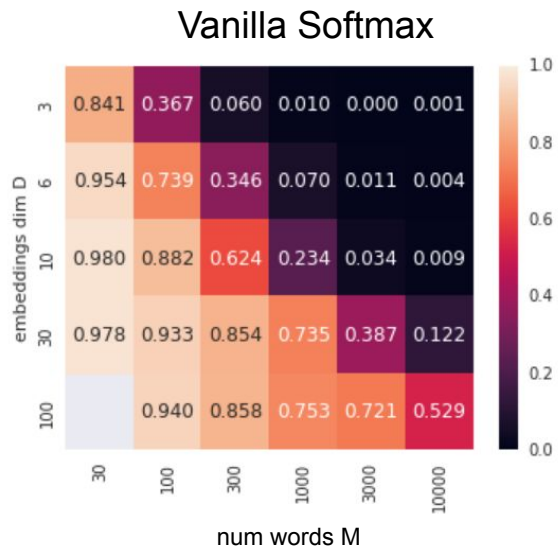
# Synthetic Experiment

- **Goal**: *separate softmax bottleneck from context embedding bottleneck*
- Sample N different categorical distributions with M outcomes:

$$P^*(\cdot|c_j) \sim \text{Dir}(\alpha), \text{ for } j = 1, \ldots, N$$

- Goal: $\quad P_\theta(x|c) = \dfrac{\exp \mathbf{h}_c^\top \mathbf{w}_x}{\sum_{x'} \exp \mathbf{h}_c^\top \mathbf{w}_{x'}} \approx P^*(x|c)$

# Synthetic Experiment

- **Goal**: *separate softmax bottleneck from context embedding bottleneck*
- Sample N different categorical distributions with M outcomes:

$$P^*(\cdot|c_j) \sim \mathrm{Dir}(\alpha), \text{ for } j = 1, \ldots, N$$

- Goal:    $P_\theta(x|c) = \dfrac{\exp \mathbf{h}_c^\top \mathbf{w}_x}{\sum_{x'} \exp \mathbf{h}_c^\top \mathbf{w}_{x'}} \approx P^*(x|c)$
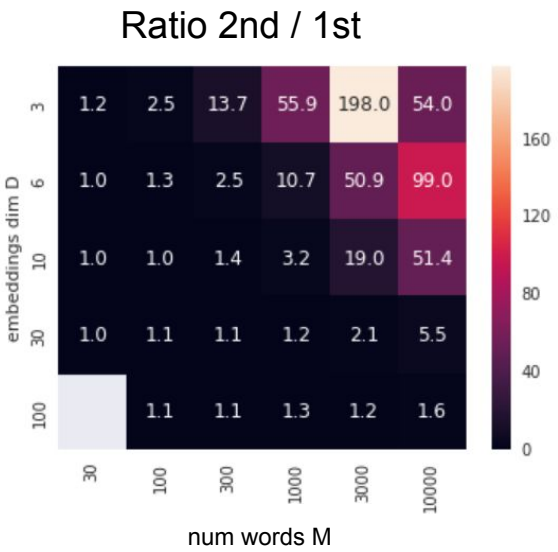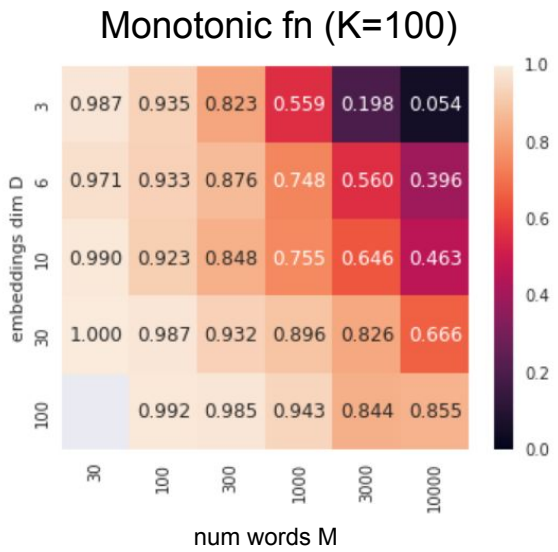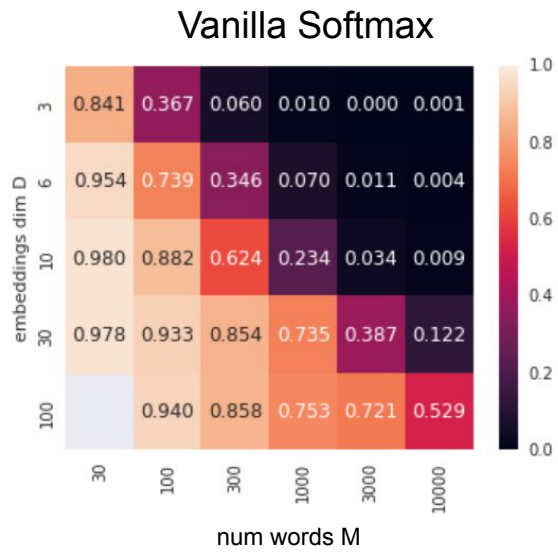
- Independent context embeddings; shared word embeddings

# Synthetic Experiments - Mode Matching $(\alpha=0.01)$

- Percentage of contexts c for which $\mathrm{argmax}_x P^*(x|c) = \mathrm{argmax}_x P_\Theta(x|c)$



Vanilla Softmax · Monotonic fn (K=100) · Ratio 2nd / 1st

# Synthetic Experiments - Mode Matching  ($\alpha$=0.01)

- Percentage of contexts c for which $\mathrm{argmax}_x P^*(x|c) = \mathrm{argmax}_x P_\Theta(x|c)$



Vanilla Softmax          Monotonic fn (K=100)          Ratio 2nd / 1st

- Similar results for cross-entropy and other values of $\alpha$

# Piecewise Linear Increasing Functions (PLIF)

- NN w/ 1 hidden layer ⇒ memory hungry:
  - Tensor of size N x M x K on GPU , where K >= 1000

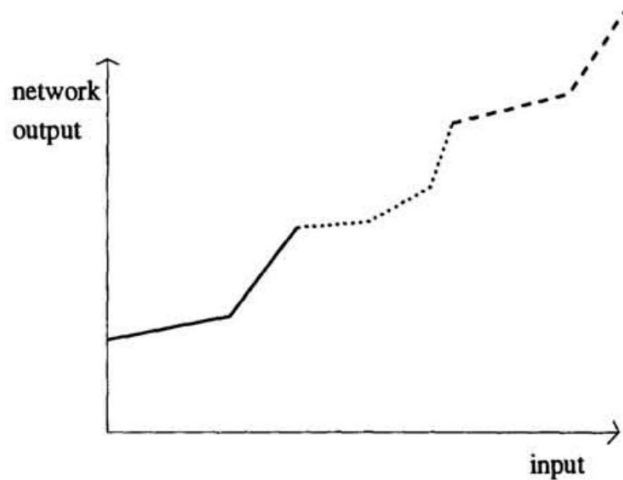# Piecewise Linear Increasing Functions (PLIF)

- NN w/ 1 hidden layer ⇒ memory hungry:
    - Tensor of size N x M x K on GPU , where K >= 1000
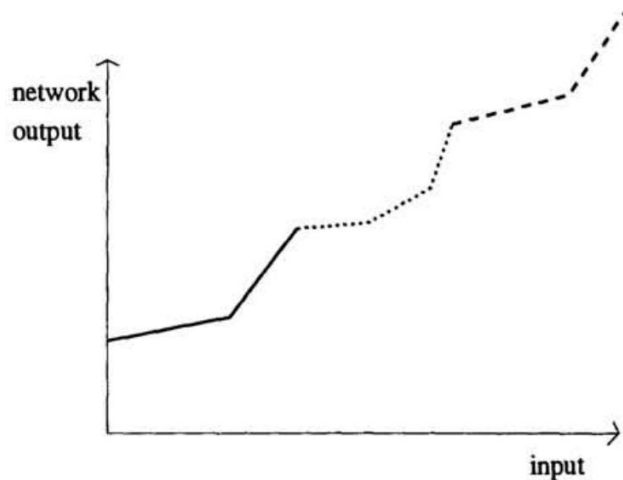
- **PLIF:**

# Piecewise Linear Increasing Functions (PLIF)

- NN w/ 1 hidden layer ⇒ memory hungry:
  - Tensor of size N x M x K on GPU , where K >= 1000

- **PLIF**:



- Forward & backward passes: just a lookup in two K dim vectors
- Memory and running time very efficient (comparable with Vanilla Softmax)

# Language Modeling Results

| | PENN TREEBANK | | | | WIKITEXT-2 | | | |
|---|---|---|---|---|---|---|---|---|
| | #PARAM | VALID PPL | TEST PPL | #SEC/EP | #PARAM | VALID PPL | TEST PPL | #SEC/EP |
| LINEAR-SOFTMAX w/ AWD-LSTM, w/o FINETUNE (MERITY ET AL., 2017) | 24.2M | 60.83 | 58.37 | ~60 | 33M | 68.11 | 65.22 | ~120 |
| OURS LMS-PLIF, $10^5$ KNOTS w/ AWD-LSTM, w/o FINETUNE | 24.4M | 59.45 | 57.25 | ~70 | 33.2M | 67.87 | 64.86 | ~150 |
| MOS, K = 15 w/ AWD-LSTM, w/o FINETUNE (YANG ET AL., 2017) | 26.6M | 58.58 | 56.43 | ~150 | 33M | 66.01 | 63.33 | ~550 |
| MOS(15 COMP) + OUR PLIF ($10^6$ KNOTS) w/ AWD-LSTM, w/o FINETUNE | 28.6M | 58.20 | 56.02 | ~220 | - | - | - | - |

GPU Memory: N x M x K

GPU Memory: N x M

# Thank you!

**Poster #23**