### Learning in Gated Neural Networks

### Pramod Viswanath

University of Illinois at Urbana-Champaign

### Gated Recurrent Neural Networks

- Well-known examples are LSTMs and GRUs
- Achieve state-of-the-art results in many challenging ML tasks



Figure: Google Duplex

# Siri, Alexa and more...

Google Translate

amazon alexa

Break through language barriers

• Language translation

• Speech recognition

Phrase completion









### NNs and RNNs

• Feed-forward neural networks



• Recurrent neural networks (Vanilla)



# Gated RNNs



Figure: Gated Recurrent Unit (GRU)

Key features:

- Gating mechanism
- Non-linear 'switching' dynamical systems
- Provide 'long term memory'





Gates: z<sub>t</sub>, r<sub>t</sub> ∈ [0,1]<sup>d</sup> depend on the input x<sub>t</sub> and the past h<sub>t-1</sub>
States: h<sub>t</sub>, h̃<sub>t</sub> ∈ ℝ<sup>d</sup>

Update equations for each *t*:

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t$$
$$\tilde{h}_t = f(Ax_t + r_t \odot Bh_{t-1})$$

### Building blocks of GRU

 $h_t = (1 - z_t) \odot h_{t-1} + z_t \odot r_t \odot f(Ax_t + Bh_{t-1}) + z_t \odot (1 - r_t) \odot f(Ax_t)$ 



### Building blocks of GRU

 $h_t = (1 - z_t) \odot h_{t-1} + z_t \odot r_t \odot f(Ax_t + Bh_{t-1}) + z_t \odot (1 - r_t) \odot f(Ax_t)$ 



### Mixture-of-Experts: Building blocks of GRU

• Introduced by Robert Jacobs, Michael Jordan, Steven Nowlan and Geoffrey Hinton in 1991



f =sigmoid, g =linear, tanh, ReLU

### MoE as gated feed-forward network



### MoE: Modern relevance

• Outrageously large neural networks



Figure 1: A Mixture of Experts (MoE) layer embedded within a recurrent language model. In this case, the sparse gating function selects two experts to perform computations. Their outputs are modulated by the outputs of the gating network.

### What is known about MoE?

Adaptive mixtures of local experts RA Jacobs, MI Jordan, SJ Nowlan, GE Hinton Neural computation 3 (1), 79-87	3663	1991
Sharing clusters among related groups: Hierarchical Dirichlet processes YW Teh, MI Jordan, MJ Beal, DM Blei Advances in neural information processing systems, 1385-1392	3273	2005
Hierarchical mixtures of experts and the EM algorithm MI Jordan, RA Jacobs Neural computation 6 (2), 181-214	3090	1994

### • No provable learning algorithms for parameters $^1$ $\odot$

<sup>&</sup>lt;sup>1</sup>20 years of MoE, MoE: a literature survey

### Open problem for 25+ years



$$\Leftrightarrow P_{y|\mathbf{x}} = f(\mathbf{w}^{\mathsf{T}}\mathbf{x}) \cdot \mathcal{N}(y|g(\mathbf{a}_{1}^{\mathsf{T}}\mathbf{x}), \sigma^{2}) + (1 - f(\mathbf{w}^{\mathsf{T}}\mathbf{x})) \cdot \mathcal{N}(y|g(\mathbf{a}_{2}^{\mathsf{T}}\mathbf{x}), \sigma^{2})$$

### Open question

Given *n* i.i.d. samples  $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$ , does there exist an efficient learning algorithm with provable theoretical guarantees to learn the regressors  $\mathbf{a}_1, \mathbf{a}_2$  and the gating parameter  $\mathbf{w}$ ?

### Traditional loss functions

### Loss functions:

Log-likelihood loss

$$L = \log\left(f(\boldsymbol{w}^{\top}\boldsymbol{x}) \cdot e^{-\frac{\|\boldsymbol{y} - \boldsymbol{g}(\boldsymbol{a}_{1}^{\top}\boldsymbol{x})\|^{2}}{2\sigma^{2}}} + (1 - f(\boldsymbol{w}^{\top}\boldsymbol{x})) \cdot e^{-\frac{\|\boldsymbol{y} - \boldsymbol{g}(\boldsymbol{a}_{2}^{\top}\boldsymbol{x})\|^{2}}{2\sigma^{2}}}\right)$$

$$L = \left(y - \left(f(\boldsymbol{w}^{\mathsf{T}}\boldsymbol{x})g(\boldsymbol{a}_{1}^{\mathsf{T}}\boldsymbol{x}) + (1 - f(\boldsymbol{w}^{\mathsf{T}}\boldsymbol{x}))g(\boldsymbol{a}_{2}^{\mathsf{T}}\boldsymbol{x})\right)\right)^{2}$$



# Traditional algorithms

Algorithms: EM, Gradient descent, and their variants

- Practical: Often get stuck in local optima
- Theoretical: Loss surface is hard to analyze because of coupling of **w** and (**a**<sub>1</sub>, **a**<sub>2</sub>). Barely understood for far simpler problem of Gaussian mixtures

### Modular structure

Mixture of classification ( $\boldsymbol{w}$ ) and regression ( $\boldsymbol{a}_1, \boldsymbol{a}_2$ ) problems



### Key observation



### Key observation

If we know the regressors, learning the gating parameter is easy and vice-versa. How to break the gridlock?

# Focus of this talk: Breaking the gridlock

### • First learning guarantees for MoE

### Method 1: Algorithms

We propose an algorithm with first recoverable guarantees

### Method 2: Optimization framework

We design 'unusual' loss function on which traditional algorithms like GD converge to true parameters

- Both approaches work with global initializations
  - x is Gaussian

# Generalizability

*k*-MoE





### Generalizability

Hierarchical mixture of experts (HME)



Figure 2: A two-level hierarchical mixture of experts

# Method 1: Design of algorithms

### Algorithmic approach: An overview

Recall the model for MoE:

$$P_{y|\mathbf{x}} = f(\mathbf{w}^{\mathsf{T}}\mathbf{x}) \cdot \mathcal{N}(y|g(\mathbf{a}_{1}^{\mathsf{T}}\mathbf{x}), \sigma^{2}) + (1 - f(\mathbf{w}^{\mathsf{T}}\mathbf{x})) \cdot \mathcal{N}(y|g(\mathbf{a}_{2}^{\mathsf{T}}\mathbf{x}), \sigma^{2})$$

- We learn  $(a_1, a_2)$  and w separately
- First recover  $(a_1, a_2)$  without knowing w at all
- Later learn w using traditional methods like EM
- Global consistency guarantees (population setting)

### Learning regressors without gating

Recall the model for MoE:

$$P_{y|\mathbf{x}} = f(\mathbf{w}^{\mathsf{T}}\mathbf{x}) \cdot \mathcal{N}(y|g(\mathbf{a}_{1}^{\mathsf{T}}\mathbf{x}), \sigma^{2}) + (1 - f(\mathbf{w}^{\mathsf{T}}\mathbf{x})) \cdot \mathcal{N}(y|g(\mathbf{a}_{2}^{\mathsf{T}}\mathbf{x}), \sigma^{2})$$

In the absence of gating parameters,

$$P_{y|\mathbf{x}} = p \cdot \mathcal{N}(y|g(\mathbf{a}_{1}^{\mathsf{T}}\mathbf{x}), \sigma^{2}) + (1-p) \cdot \mathcal{N}(y|g(\mathbf{a}_{2}^{\mathsf{T}}\mathbf{x}), \sigma^{2})$$

• Mixture of generalized linear models (GLMs)!

- ▶ How do we learn **a**<sub>1</sub> and **a**<sub>2</sub> without knowing *p*?
- Method of moments

### Method of moments in GLMs

• Basic idea: Construct a **third-order super-symmetric** tensor from data such that

$$\mathbb{E}(\psi(X,Y)) = \sum_{i} \boldsymbol{a}_{i} \otimes \boldsymbol{a}_{i} \otimes \boldsymbol{a}_{i} \Rightarrow \boldsymbol{a}_{i}$$
 can be recovered



- How do we construct  $\psi$ ?
  - Stein's lemma

### Stein's lemma 101

# Stein's lemma For $f : \mathbb{R}^d \to \mathbb{R}$ and $\mathbf{x} \sim \mathcal{N}(0, I_d)$ , $\mathbb{E}[f(\mathbf{x}) \cdot \mathbf{x}] = \mathbb{E}[\nabla_{\mathbf{x}} f(\mathbf{x})] \in \mathbb{R}^d$ .

Non-linear regression using Stein's lemma: If  $y = g(a_1^T x) + N$ , then

$$\mathbb{E}[y \cdot \mathbf{x}] = \mathbb{E}[g(\mathbf{a}_{1}^{\mathsf{T}}\mathbf{x}) \cdot \mathbf{x}] + \mathbb{E}[N \cdot \mathbf{x}]$$
  
Estimated from samples  
$$= \mathbb{E}[\nabla_{\mathbf{x}}g(\mathbf{a}_{1}^{\mathsf{T}}\mathbf{x})]$$
$$\propto \mathbf{a}_{1}$$

### Mixture of GLMs: Stein's lemma 101

• Recall, for mixture of GLMs:

$$P_{y|\mathbf{x}} = p \cdot \mathcal{N}(y|g(\mathbf{a}_1^{\mathsf{T}}\mathbf{x}), \sigma^2) + (1-p) \cdot \mathcal{N}(y|g(\mathbf{a}_2^{\mathsf{T}}\mathbf{x}), \sigma^2)$$

• From Stein's lemma,

$$\mathbb{E}[\boldsymbol{y}\cdot\boldsymbol{x}] \propto \boldsymbol{p}\cdot\boldsymbol{a}_1 + (1-\boldsymbol{p})\cdot\boldsymbol{a}_2.$$

- Not unique in **a**<sub>1</sub> and **a**<sub>2</sub>
- How can we ensure uniqueness?

### Stein's lemma 102

### 2nd order Stein's lemma

$$\mathbb{E}[f(\mathbf{x}) \cdot \underbrace{(\mathbf{x}\mathbf{x}^{\top} - I)}_{S_2(\mathbf{x})}] = \mathbb{E}[\nabla_{\mathbf{x}}^{(2)}f(\mathbf{x})] \in \mathbb{R}^{d \times d}.$$

• Mixture of GLMs:

$$P_{y|\mathbf{x}} = p \cdot \mathcal{N}(y|g(\mathbf{a}_{1}^{\mathsf{T}}\mathbf{x}), \sigma^{2}) + (1-p) \cdot \mathcal{N}(y|g(\mathbf{a}_{2}^{\mathsf{T}}\mathbf{x}), \sigma^{2})$$
  
$$\Rightarrow \mathbb{E}[y \cdot (\mathbf{x}\mathbf{x}^{\mathsf{T}} - I)] \propto 2p \cdot \mathbf{a}_{1}\mathbf{a}_{1}^{\mathsf{T}} + 2(1-p) \cdot \mathbf{a}_{2}\mathbf{a}_{2}^{\mathsf{T}}.$$

- Not unique!
- How can we ensure uniqueness?

### Stein's lemma 103

### 3rd order Stein's lemma

$$\mathbb{E}[f(\mathbf{x}) \cdot \mathcal{S}_3(\mathbf{x})] = \mathbb{E}[\nabla_{\mathbf{x}}^{(3)} f(\mathbf{x})] \in \mathbb{R}^{d \times d \times d}$$

• Score transformation  $S_3(\mathbf{x}) = \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x} - \sum_{i \in [d]} \operatorname{sym}(\mathbf{x} \otimes \mathbf{e}_i \otimes \mathbf{e}_i)$ 

#### • Mixture of GLMs:

$$P_{y|\mathbf{x}} = p \cdot \mathcal{N}(y|g(\mathbf{a}_{1}^{\mathsf{T}}\mathbf{x}), \sigma^{2}) + (1-p) \cdot \mathcal{N}(y|g(\mathbf{a}_{2}^{\mathsf{T}}\mathbf{x}), \sigma^{2})$$
  
$$\Rightarrow \mathbb{E}[y \cdot \mathcal{S}_{3}(\mathbf{x})] \propto p \cdot \mathbf{a}_{1} \otimes \mathbf{a}_{1} \otimes \mathbf{a}_{1} + (1-p) \cdot \mathbf{a}_{2} \otimes \mathbf{a}_{2} \otimes \mathbf{a}_{2}.$$

- Unique! (by Kruskal's theorem)
- Can we extend this to MoE?

### MoE: Stein's lemma

• For MoE, 
$$p = p(x) = f(\boldsymbol{w}^{\mathsf{T}}\boldsymbol{x})$$
 since

$$P_{y|\mathbf{x}} = f(\mathbf{w}^{\mathsf{T}}\mathbf{x}) \cdot \mathcal{N}(y|g(\mathbf{a}_{1}^{\mathsf{T}}\mathbf{x}), \sigma^{2}) + (1 - f(\mathbf{w}^{\mathsf{T}}\mathbf{x})) \cdot \mathcal{N}(y|g(\mathbf{a}_{2}^{\mathsf{T}}\mathbf{x}), \sigma^{2})$$

- Can we use Stein's lemma to learn **a**<sub>1</sub> and **a**<sub>2</sub>?
- Natural attempt:

$$\mathbb{E}[\mathbf{y} \cdot S_3(\mathbf{x})] = \mathbf{a}_1 \otimes \mathbf{a}_1 \otimes \mathbf{a}_1 + \mathbf{w} \otimes \mathbf{a}_1 \otimes \mathbf{w} + \ldots + \mathbf{a}_1 \otimes \mathbf{a}_1 \otimes \mathbf{w} + \ldots$$

Not a super-symmetric tensor

• Can we construct a super-symmetric tensor for MoE?

### Key insight: Hermite polynomial transformation

Suppose g =linear and  $\sigma$  = 0. Then

$$P_{y|\mathbf{x}} = f(\mathbf{w}^{\mathsf{T}}\mathbf{x}) \cdot \mathbb{1}\{y = \mathbf{a}_{1}^{\mathsf{T}}\mathbf{x}\} + (1 - f(\mathbf{w}^{\mathsf{T}}\mathbf{x}))\mathbb{1}\{y = \mathbf{a}_{1}^{\mathsf{T}}\mathbf{x}\}$$
  
$$\Rightarrow \mathbb{E}[y^{3} - 3y|\mathbf{x}] = \sum_{i \in \{1,2\}} f(\mathbf{w}_{i}^{\mathsf{T}}\mathbf{x})((\mathbf{a}_{i}^{\mathsf{T}}\mathbf{x})^{3} - 3(\mathbf{a}_{i}^{\mathsf{T}}\mathbf{x})), \quad \mathbf{w}_{2} = -\mathbf{w}_{1}$$

Now applying Stein's lemma,

$$\mathbb{E}[(y^3 - 3y) \cdot \mathcal{S}_3(\boldsymbol{x})] = \mathbb{E}[\nabla_{\boldsymbol{x}}^3 \mathbb{E}[y^3 - 3y|\boldsymbol{x}]] = 3\sum_{i \in \{1,2\}^i} \boldsymbol{a}_i \otimes \boldsymbol{a}_i \otimes \boldsymbol{a}_i$$

How do cross terms like  $a_i \otimes a_i \otimes w$  disappear?

- Reason:  $\mathbb{E}[H'_3(Z)] = \mathbb{E}[H''_3(Z)] = \mathbb{E}[H'''_3(Z)] = 0$
- $H_3(z) = z^3 3z$  is third-Hermite polynomial

Does this work for  $\sigma \neq 0$ ?

### Linear experts: Hermite-like-polynomials

Suppose g = linear and  $\sigma \neq 0$ :

$$P_{y|\mathbf{x}} = f(\mathbf{w}^{\mathsf{T}}\mathbf{x}) \cdot \mathcal{N}(y|\mathbf{a}_{1}^{\mathsf{T}}\mathbf{x},\sigma^{2}) + (1 - f(\mathbf{w}^{\mathsf{T}}\mathbf{x})) \cdot \mathcal{N}(y|\mathbf{a}_{2}^{\mathsf{T}}\mathbf{x},\sigma^{2})$$

Super-symmetric tensor

$$\mathcal{T}_{3} = \mathbb{E}[(y^{3} - 3y(1 + \sigma^{2})) \cdot \mathcal{S}_{3}(\boldsymbol{x})] = 3(\boldsymbol{a}_{1} \otimes \boldsymbol{a}_{1} \otimes \boldsymbol{a}_{1} + \boldsymbol{a}_{2} \otimes \boldsymbol{a}_{2} \otimes \boldsymbol{a}_{2})$$

• This very much needs special linear structure. What about other non-linearities for g?

### Generalization: Cubic polynomial transformations

• For a wide class of non-linearities such as *g*=linear, sigmoid, ReLU, etc.

$$\mathcal{T}_3 = \mathbb{E}[(y^3 + \alpha y^2 + \beta y) \cdot \mathcal{S}_3(\mathbf{x})] = c(\mathbf{a}_1 \otimes \mathbf{a}_1 \otimes \mathbf{a}_1 + \mathbf{a}_2 \otimes \mathbf{a}_2 \otimes \mathbf{a}_2)$$

• How do we choose  $\alpha$  and  $\beta$ ?

- Solving a linear system
- Example: For sigmoid,

$$\begin{bmatrix} 0.2067 & 0.2066 \\ 0.0624 & -0.0001 \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \begin{bmatrix} -0.1755 - 0.6199\sigma^2 \\ -0.0936 \end{bmatrix}$$

• Key idea: Acts like a 'Hermite' like polynomial for general g and cancels cross terms

### Learning regressors: Spectral decomposition

#### Algorithm

- Input: Samples  $(\mathbf{x}_i, y_i)$
- Compute  $\hat{\mathcal{T}}_3 = (1/n) \sum_i H_3(y_i) \cdot \mathcal{S}_3(\boldsymbol{x}_i)$
- $\hat{a}_1, \hat{a}_2 = \text{Rank-2}$  decomposition on  $\mathcal{T}_3$

# Learning the gating

Recall

$$P_{y|\mathbf{x}} = f(\mathbf{w}^{\mathsf{T}}\mathbf{x}) \cdot \mathcal{N}(y|\mathbf{a}_{1}^{\mathsf{T}}\mathbf{x},\sigma^{2}) + (1 - f(\mathbf{w}^{\mathsf{T}}\mathbf{x})) \cdot \mathcal{N}(y|\mathbf{a}_{2}^{\mathsf{T}}\mathbf{x},\sigma^{2})$$

- If we know  $a_1$  and  $a_2$ , learning w is a classification problem!
- Traditional methods:
  - EM algorithm
  - Gradient descent on log-likelihood

## Theoretical contributions

- Show global convergence for existing methods
- Provide convergence rate
- Finite sample complexity
- First theoretical guarantees

### Learning the gating parameters

Suppose spectral methods give  $\hat{a}_i$  with  $\|\hat{a}_i - a_i\|_2 \le \sigma^2 \varepsilon$ 

For high SNR, i.e.  $\sigma < \sigma_0$ ,  $\sigma_0$  is a dimension independent constant:

- EM iterates converge geometrically to  $\hat{\boldsymbol{w}}$
- Convergence rate is a dimension-independent constant depending on  $\sigma$  and  $\|{\pmb a}_1 {\pmb a}_2\|$
- $\hat{\boldsymbol{w}}$  is  $\varepsilon$ -close to the ground truth

Method 2: Optimization framework-loss function design

### Regressors: Loss function design

$$P_{y|\boldsymbol{x}} = f(\boldsymbol{w}^{\mathsf{T}}\boldsymbol{x}) \cdot \mathcal{N}(y|g(\boldsymbol{a}_{1}^{\mathsf{T}}\boldsymbol{x}), \sigma^{2}) + (1 - f(\boldsymbol{w}^{\mathsf{T}}\boldsymbol{x})) \cdot \mathcal{N}(y|g(\boldsymbol{a}_{2}^{\mathsf{T}}\boldsymbol{x}), \sigma^{2})$$

#### • Traditional approaches: *l*<sub>2</sub>-loss, log-likelihood loss

- Get stuck in local minima
- No theoretical analysis
- Single loss function for both  $(a_1, a_2)$  and w
- Formulation of right loss function is critical (Jacobs et. al 1991)

### Theoretical contributions

• Separate loss functions  $L_4$  and  $L_{log}$  to learn  $(a_1, a_2)$  and w



• Gradient descent on both  $L_4$  and  $L_{log}$ . What are they?

### Tensor based loss function for regressors

• For linear experts,

$$P_{y|\mathbf{x}} = f(\mathbf{w}^{\mathsf{T}}\mathbf{x}) \cdot \mathcal{N}(y|\mathbf{a}_{1}^{\mathsf{T}}\mathbf{x},\sigma^{2}) + (1 - f(\mathbf{w}^{\mathsf{T}}\mathbf{x})) \cdot \mathcal{N}(y|\mathbf{a}_{2}^{\mathsf{T}}\mathbf{x},\sigma^{2})$$

• Stein's lemma+ 4-Hermite polynomial implies

$$\mathcal{T}_4 = \mathbb{E}[(y^4 - 6y^2(1 + \sigma^2)) \cdot \mathcal{S}_4(\boldsymbol{x})] = 12(\boldsymbol{a}_1^{\otimes 4} + \boldsymbol{a}_2^{\otimes 4})$$

• If  $\hat{a}_1$  and  $\hat{a}_2$  are parameters,

$$\begin{split} L_4(\hat{\boldsymbol{a}}_1, \hat{\boldsymbol{a}}_2) &\triangleq \sum_{j \neq k} \mathcal{T}_4(\hat{\boldsymbol{a}}_j, \hat{\boldsymbol{a}}_j, \hat{\boldsymbol{a}}_k, \hat{\boldsymbol{a}}_k) - \mu \sum_{j \in \{1, 2\}} \mathcal{T}_4(\hat{\boldsymbol{a}}_j, \hat{\boldsymbol{a}}_j, \hat{\boldsymbol{a}}_j, \hat{\boldsymbol{a}}_j) \\ &+ \lambda \sum_{j \in \{1, 2\}} (\|\hat{\boldsymbol{a}}_j\|^2 - 1)^2 \end{split}$$

# Landscape of $L_4$

### Properties

- No spurious local minima: All local minima are global
- Global minima are ground truth (upto permutation and sign-flip)
- All saddle points have negative curvature
- SGD converges to approximate global minima

Why  $L_4$ ?

# Why $L_4$ ?

- Connection to tensor based losses
- We show that

$$L_4(\hat{\boldsymbol{a}}_1, \hat{\boldsymbol{a}}_2) = 12 \sum_i \sum_{j \neq k} \langle \boldsymbol{a}_i, \hat{\boldsymbol{a}}_j \rangle^2 \langle \boldsymbol{a}_i, \hat{\boldsymbol{a}}_k \rangle^2 - 12\mu \sum_i \sum_j \langle \boldsymbol{a}_i, \hat{\boldsymbol{a}}_j \rangle^4 + \lambda \sum_j (\|\boldsymbol{a}_j\|^2 - 1)^2$$

- 4-order tensor loss
  - Landscape analysis in (Ge et. al 2018)

# Summary

- Algorithmic innovation: First provably consistent algorithms for MoE in 25+ years
- Loss function innovation: First SGD based algorithm on novel loss functions with provably nice landscape properties
- Sample complexity: First sample complexity results for MoE
- Global convergence: Our algorithms work with global initializations

### Conclusion



- 1. Theoretical understanding  $\checkmark$
- 2. Novel algorithms ✓

- 1. Theoretical understanding?
- 2. Algorithms?

# Collaborators





# Ashok Vardhan Makkuva

Sreeram Kannan

# Thank you!