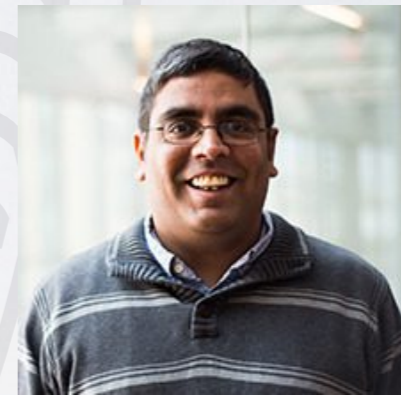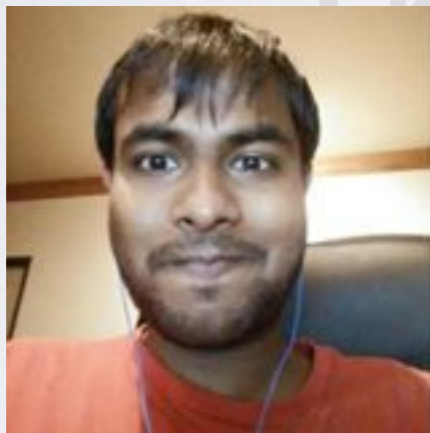AAAI' 19

# EXPLICITLY IMPOSING CONSTRAINTS IN DEEP NETWORKS VIA CONDITIONAL GRADIENTS GIVES IMPROVED GENERALIZATION AND FASTER CONVERGENCE

*Sathya Ravi, Tuan Dinh, Vishnu Lokhande, Vikas Singh*

*Department of Computer Sciences*
*University of Wisconsin–Madison*

11/14/2018

WISCONSIN
UNIVERSITY OF WISCONSIN–MADISON

# DEEP LEARNING

# DEEP LEARNING

**Solve** $\quad \min_{W \in \mathbb{R}^n} L(W)$

# DEEP LEARNING

**Solve**  $\min_{W \in \mathbb{R}^n} L(W)$

$$L(W) = \mathbb{E}_\xi f(W, \xi)$$

$$\xi = (x, y) \sim \mathcal{D}$$

Compute an estimate of gradient

Compute an estimate of gradient

$$W_{t+1} = W_t - \eta_t \nabla \tilde{L}_t(W_t)$$

Compute an estimate of gradient

$$W_{t+1} = W_t - \eta_t \nabla \tilde{L}_t(W_t)$$

$$\mathbb{E}\left[\nabla \tilde{L}_t(W_t)\right] = \nabla L(W_t)$$

$$\mathbb{E}\left[\left\|\nabla \tilde{L}_t(W) - \nabla L(W)\right\|^2\right] \leq \sigma^2$$

Compute an estimate of gradient

$$W_{t+1} = W_t - \eta_t \nabla \tilde{L}_t(W_t)$$

$$\mathbb{E}\left[\nabla \tilde{L}_t(W_t)\right] = \nabla L(W_t)$$

$$\mathbb{E}\left[\left\|\nabla \tilde{L}_t(W) - \nabla L(W)\right\|^2\right] \leq \sigma^2$$

**What about learning?**

# QUALITY/PREDICTIVE PERFORMANCE

# QUALITY/PREDICTIVE PERFORMANCE

$$\mathcal{R}(W) = \mathbb{E}_{(x,y)\sim\mathcal{D}} L(W; (x, y))$$

# QUALITY/PREDICTIVE PERFORMANCE

$$\mathcal{R}(W) = \mathbb{E}_{(x,y)\sim\mathcal{D}} L(W; (x,y))$$

$$\mathcal{R}_S(W) = \frac{1}{n}\sum_{i=1}^{n} L(W; (x_i, y_i))$$

# QUALITY/PREDICTIVE PERFORMANCE

$$\mathcal{R}(W) = \mathbb{E}_{(x,y)\sim\mathcal{D}} L(W; (x, y))$$

$$\mathcal{R}_S(W) = \frac{1}{n} \sum_{i=1}^{n} L(W; (x_i, y_i))$$

## The one true theorem

$$\mathcal{R}(W) = \mathcal{R}_S(W) + \mathcal{R}(W) - \mathcal{R}_S(W)$$

# QUALITY/PREDICTIVE PERFORMANCE

$$\mathcal{R}(W) = \mathbb{E}_{(x,y)\sim\mathcal{D}} L(W; (x,y))$$

$$\mathcal{R}_S(W) = \frac{1}{n} \sum_{i=1}^{n} L(W; (x_i, y_i))$$

## The one true theorem

$$\mathcal{R}(W) = \underbrace{\mathcal{R}_S(W)}_{\textbf{Train error}} + \underbrace{\mathcal{R}(W) - \mathcal{R}_S(W)}_{\Delta_S(W):=\textbf{Test error}}$$

**Train error**          $\Delta_S(W)$:=**Test error**

# QUALITY/PREDICTIVE PERFORMANCE

# QUALITY/PREDICTIVE PERFORMANCE

## **Classical**

$$\Delta_S(W) \propto \#\text{parameters}$$

# QUALITY/PREDICTIVE PERFORMANCE

## Classical

$$\Delta_S(W) \propto \#\text{parameters}$$

## Modern (refined)

$$\Delta_S(W) \propto \|W\|$$

# QUALITY/PREDICTIVE PERFORMANCE

**Classical**

**Modern (refined)**

$$\Delta_S(W) \propto \#\text{parameters}$$

$$\Delta_S(W) \propto \|W\|$$

**Simple idea: Enforce "high" quality**

# HOW TO HAVE YOUR CAKE AND EAT IT (TOO)?

# HOW TO HAVE YOUR CAKE AND EAT IT (TOO)?

Ingredients

# HOW TO HAVE YOUR CAKE AND EAT IT (TOO)?

Ingredients

1. High quality: constraints with nice theoretical properties

# HOW TO HAVE YOUR CAKE AND EAT IT (TOO)?

Ingredients

1. High quality: constraints with nice theoretical properties

2. A "fast" algorithm: Projection free approaches are more parsimonious

# HOW TO HAVE YOUR CAKE AND EAT IT (TOO)?

Ingredients

1. High quality: constraints with nice theoretical properties

2. A "fast" algorithm: Projection free approaches are more parsimonious

3. Resources: GPUs, fellow graduate students, adviser etc.

# CONTRIBUTIONS

# CONTRIBUTIONS

- Enforcing various generic constraints: $R(W) \leq \lambda$

# CONTRIBUTIONS

- Enforcing various generic constraints: $R(W) \leq \lambda$

- Enforcing path norm constraint

# CONTRIBUTIONS

- Enforcing various generic constraints: $R(W) \leq \lambda$

- Enforcing path norm constraint

- Experiments with three different tasks and datasets

# WHY?

# WHY?

# *EXCELLENT* GENERIC CONSTRAINTS

# EXCELLENT GENERIC CONSTRAINTS

$$\|W\|_F, \quad \|W\|_*$$

# *GOOD* GENERIC CONSTRAINTS

# GOOD GENERIC CONSTRAINTS

$$\|W\|_1, \quad \|W\|_\infty$$

# *PESKY* GENERIC CONSTRAINTS

# PESKY GENERIC CONSTRAINTS

$$\|W\|_{\text{TV}}$$

# PESKY GENERIC CONSTRAINTS

$$\|W\|_{\mathbf{TV}}$$

**Lemma: An** $\epsilon$ − **approximate CG step can be computed in** $O\left(1/\epsilon\right)$ **time (independent of dimensions).**
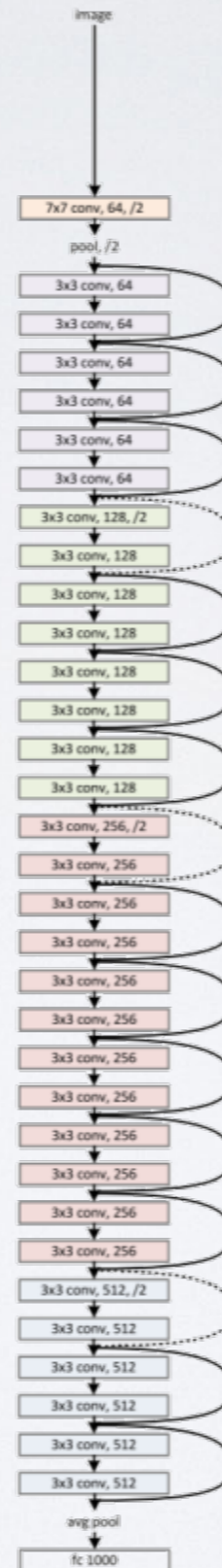
# *PATH* CG

# *PATH* CG

$$\|W\|_\pi^2 = \sum_{v_{in}[i] \xrightarrow{e_1} v_1 \xrightarrow{e_2} \cdots v_{out}[j]} \left| \prod_{k=1}^{l} W_{e_k} \right|^2$$
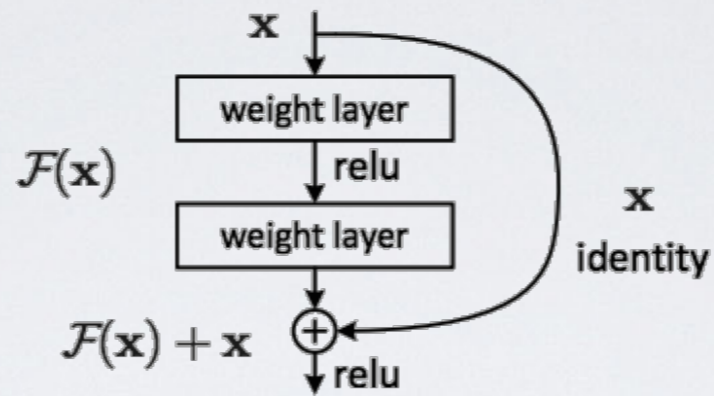
# *PATH* CG

$$\|W\|_\pi^2 = \sum_{v_{in}[i] \xrightarrow{e_1} v_1 \xrightarrow{e_2} \cdots v_{out}[j]} \left| \prod_{k=1}^{l} W_{e_k} \right|^2$$

- Why path norm? Rescale invariance

# PATH CG

$$\|W\|_\pi^2 = \sum_{v_{in}[i] \xrightarrow{e_1} v_1 \xrightarrow{e_2} \cdots v_{out}[j]} \left| \prod_{k=1}^{l} W_{e_k} \right|^2$$

- Why path norm? Rescale invariance

# *PATH* CG

$$\|W\|_\pi^2 = \sum_{v_{in}[i] \xrightarrow{e_1} v_1 \xrightarrow{e_2} \cdots v_{out}[j]} \left| \prod_{k=1}^{l} W_{e_k} \right|^2$$

- Why path norm? Rescale invariance

- Projection at least NP Hard

# *PATH* CG

$$\|W\|_\pi^2 = \sum_{v_{in}[i] \overset{e_1}{\to} v_1 \overset{e_2}{\to} \cdots v_{out}[j]} \left| \prod_{k=1}^{l} W_{e_k} \right|^2$$

• Why path norm? Rescale invariance

•Projection at least NP Hard

# PATH CG

$$\|W\|_\pi^2 = \sum_{v_{in}[i] \overset{e_1}{\to} v_1 \overset{e_2}{\to} \cdots v_{out}[j]} \left| \prod_{k=1}^{l} W_{e_k} \right|^2$$

- Why path norm? Rescale invariance

- Projection at least NP Hard

- Subproblems in CG can be solved in  O(Bl)
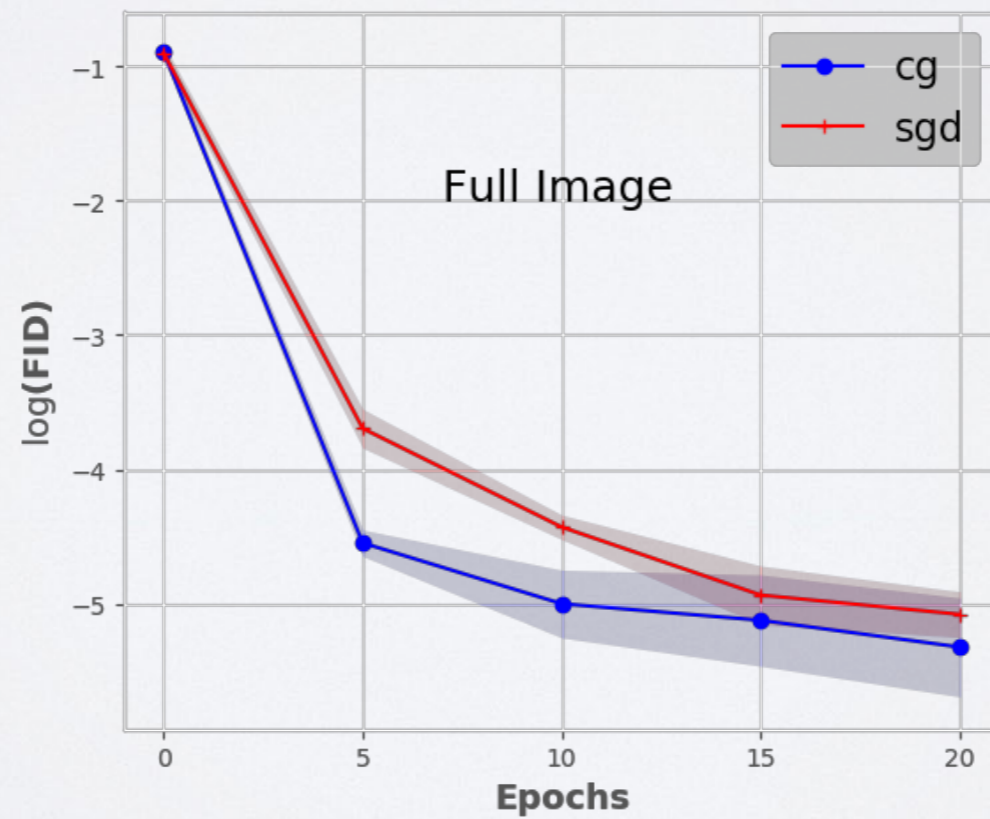
# RESNETS



34-layer residual

# RESNETS

# PATH-CG

# DC-GAN

# DC-GAN

# DC-GAN



| Masked | Ground Truth | CG | SGD |