# Neural networks with recurrent generative feedback

https://arxiv.org/abs/2007.09200

Yujia Huang, Caltech

yjhuang@caltech.edu





Yujia Huang, Caltech



James Gornet, Caltech



Sihui Dai, Caltech



Zhiding Yu, NVIDIA



Tan Nguyen, Rice University



Doris Y. Tsao, Caltech



Anima Anandkumar, Caltech/NVIDIA



### **Self-Consistency**

Given a joint distribution  $p(h, y, z; \theta)$  parameterized by  $\theta$ ,  $(\hat{h}, \hat{y}, \hat{z})$  are self-consistent if they satisfy the following constraints:



# **Self-Consistency**





h: encoded features z: latent variables

# **Self-Consistency**



x "panda" 57.7% confidence



"nematode" 8.2% confidence



 $m{x} + \epsilon \operatorname{sign}(
abla_{m{x}} J(m{ heta}, m{x}, y))$ "gibbon" 99.3 % confidence



x: imagesh: encoded featuresz: latent variablesy: labels





### **Generative Classifier**

Logistic Regression



Gaussian Naïve Classifier



A. Ng, and M. I. Jordan. On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes. In Neurips 2002.

### Deconvolutional generative model (DGM)



$$\begin{split} y &\sim p(y) \\ z_P^{(i)} &\sim \mathrm{Ber}(\frac{e^{b \cdot g^{(i)}}}{e^{b \cdot g^{(i)}} + 1}) \\ z_R^{(i)} &\sim \mathrm{Ber}(\frac{e^{b \cdot g^{(i)}}}{e^{b \cdot g^{(i)}} + 1}) \\ x &\sim \mathcal{N}(g(0), \mathrm{diag}(\sigma^2)) \end{split}$$

T. Nguyen, N. Ho, A. Patel, A. Anandkumar, M. I. Jordan, and R. G. Baraniuk. A bayesian perspective of convolutional neural networks through a deconvolutional generative model. arXiv:1811.02657, 2018.

### **Inference in the DGM**



- MAP estimate of y:  $\hat{y} = \text{CNN}(h)$
- MAP estimate of h:  $\hat{h} = g(0)$
- MAP estimate of z (informal):  $\hat{z_R} = \mathbb{1}\{\sigma_{\text{AdaReLU}} \neq 0\}$  $\hat{z_P} = \mathbb{1}\{\sigma_{\text{AdaPool}} \neq 0\}$

$$\sigma_{\text{AdaReLU}}(f) = \begin{cases} \sigma_{\text{ReLU}}(f), & \text{if } g \ge 0\\ \sigma_{\text{ReLU}}(-f), & \text{if } g < 0 \end{cases}$$

$$\sigma_{\text{AdaPool}}(f) = \begin{cases} \sigma_{\text{MaxPool}}(f), & \text{if } g \ge 0\\ -\sigma_{\text{MaxPool}}(-f), & \text{if } g < 0 \end{cases}$$

### **Iterative inference**



(y)  $(f) \quad z \quad g$   $(f) \quad z \quad g$  (h)

- → Feedforward
- → Feedback

y

h

 $\longrightarrow$  Inference of z

Label Latent variables

- Image features
- Feedforward layer

Feedback layer

**CNN-F** 



# **Training of CNN-F**

12

 $\mathcal{L}_{\text{Xentropy}}(y_0, \text{target}) \mathcal{L}_{\text{Xentropy}}(y_1, \text{target}) \mathcal{L}_{\text{Xentropy}}(y_2, \text{target})$ 



# **CNN-F** with adversarial training



Reconstruction loss is always between adversarial and natural features.

# **CNN-F on all CNN architectures**





**IN: Instance Normalization** 

VGG/Allconv/...

ResNet

# **CNN-F repairs distorted images**

#### Corrupted

#### Ground-truth



Shot Noise



#### Gaussian Noise





#### Dotted Line





## **CNN-F** improves adversarial robustness



- Standard training on Fashion-MNIST.
- Attack with PGD-40.
- CNN-F has higher adversarial robustness than CNN.

## **CNN-F** improves adversarial robustness



CNN-F trained with different iterations.

CNN-F tested with different iterations.

More iterations are needed for *harder* images.

# **CNN-F combined with adversarial training**



- Adversarial training on Fashion-MNIST.
- Trained with PGD-40 (eps=0.3). Attack with PGD-40.
- CNN-F augmented with adversarial images achieves high accuracy for both clean and adversarial data.

## **CNN-F generalizes better to different attacks**



Feedback helps when there is distribution shift between training and testing data.

# **Train on CIFAR-10**

- CNN-F (on Wide ResNet) combined with adversarial training.
- Clean accuracy decreases over iterations.
- Adversarial accuracy increases over iterations.



# **Neuronal predictivity**



- Used the fifth block and logits in VGG-16 to predict V4 and IT neuronal activities.
- CNN-F predicts V4 and IT neuronal activities better than CNN.
- Call for temporal neuronal data in the neuroscience community.

# **Conclusions and future works**

#### **Biological inspirations**

- Recurrent feedback
  - Generative models (Bayesian brain)
  - Attention
- Lateral connections
- Sparsity

#### Inspirations from other fields

- Signal processing (Kalman filters ...)
- Control (Feedback control, dynamical system)

#### **Down-streaming tasks**

- Robustness
- Few shot learning
- Uncertainty quantification



**Thank You!**