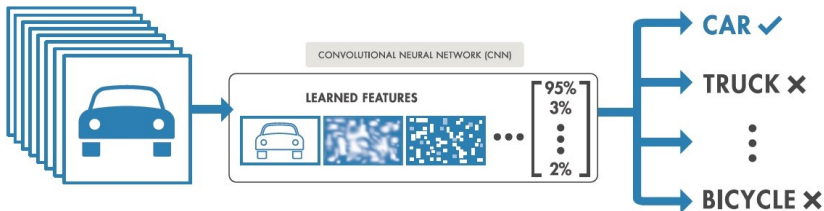


# REVE: REGULARIZING DEEP LEARNING WITH VARIATIONAL ENTROPY BOUND

Tuesday 24<sup>th</sup> September, 2019

Antoine Saporta, Yifu Chen, Michael Blot,  
Matthieu Cord

# Image Classification



Source: MathWorks (<https://goo.gl/zondfq>)

## Notations

- $X$  is the input image from  $\mathcal{X}$
- $C$  is the class label from  $\mathcal{C}$
- $Y$  is an intermediate representation of  $X$  from which is determined the predicted class  $\hat{C}$

# Regularization?

## Problem

- Huge number of parameters compared to the number of training samples
- Deep networks prone to overfitting
- Regularization: way to mitigate this gap and improve generalization

## Common strategies

- Weight decay
- Dropout
- Batch normalization

# Information-based Regularization Criteria

## Information Bottleneck (Tishby et al., 1999)

→ Principle: minimize  $I(Y, X)$  at optimal  $I(Y, C)$

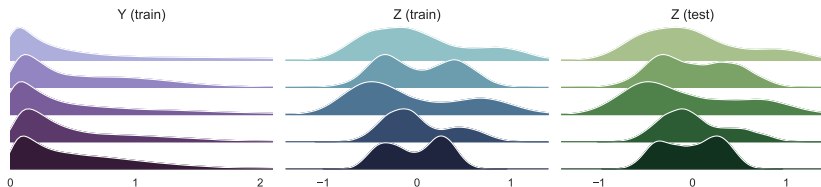
## SHADE (Blot et al., 2018)

→ Principle: minimize  $H(Y | C)$

## REVE Contribution

- Identify a *new* variable  $Z$  better suited for regularization
- Develop a variational bound over the criterion  $H(Z | C)$

# REVE Variable



## Definition

- Linear decoder:  $W_d Y$
- Unique decomposition:  $Y = Z + Y^{\ker}$  where  $Y^{\ker} \in \ker W_d$  and  $Z \in (\ker W_d)^\perp$ . Thus,  $W_d Y = W_d Z$
- REVE Variable:  $Z$ , the part of  $Y$  effectively used for prediction

# REVE Criterion



→ The conditional entropy can be written:

$$H(Z | C) = H(Z) + H(C | Z) - H(C)$$

with  $H(C)$  entirely determined by the problem

→ Objective Function:

$$\mathcal{L}_{REVE} = H(Z) + H(C | Z)$$

→ For any  $q(Z)$  and  $r(C | Z)$  variational approximations of  $p(Z)$  and  $p(C | Z)$ , resp.:

$$\mathcal{L}_{REVE} \leq - \int_{\mathbf{z}} p(\mathbf{z}) \log q(\mathbf{z}) d\mathbf{z} - \iint_{\mathbf{z}, \mathbf{c}} p(\mathbf{z}, \mathbf{c}) \log r(\mathbf{c} | \mathbf{z}) d\mathbf{z} d\mathbf{c}$$

# REVE Instantiation

## Stochastic Encoding

A stochastic encoding is needed for computing entropies:

$$Y = h(\mathbf{W}_e, X) + \epsilon \text{ with } \epsilon \sim \mathcal{N}(0, \sigma^2)$$

## Computing $Z$

- Compact Singular Value Decomposition of  $\mathbf{W}_d$ :  $\mathbf{W}_d = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$  where the  $r$  column vectors of  $\mathbf{U}$  and  $\mathbf{V}$  correspond to the non-zero singular values of  $\mathbf{W}_d$
- Computation of  $Z$ :

$$Z = \mathbf{V}\mathbf{V}^\top Y$$

# REVE Loss Function



Using Bayes' Theorem and DNN Markov-Chain hypothesis

$C \leftrightarrow X \leftrightarrow Z$ :

$$\mathcal{L}_{\text{REVE}} \leq - \iiint_{\mathcal{X} \mathcal{Z} \mathcal{C}} p(\mathbf{x})p(\mathbf{c} | \mathbf{x})p(\mathbf{z} | \mathbf{x}) (\log q(\mathbf{z}) + \log r(\mathbf{c} | \mathbf{z})) d\mathbf{x}d\mathbf{z}d\mathbf{c}$$

Thus, applying Monte-Carlo methods using the empirical distribution of  $(X, C)$  and the sampling of  $Y$  resp., we obtain the upperbound to minimize:

$$\Omega_{\text{REVE}}((\mathbf{x}_n, \mathbf{c}_n); \mathbf{W}_e; \mathbf{W}_d) = -\frac{1}{NS} \sum_{n=1}^N \sum_{s=1}^S (\log q(\mathbf{z}_{n,s}) + \log r(\mathbf{c}_n | \mathbf{z}_{n,s}))$$



Approximation  $r(\mathbf{c}|\mathbf{z})$ 

$$\Omega_{\text{REVE}} = -\frac{1}{NS} \sum_{n=1}^N \sum_{s=1}^S (\log q(\mathbf{z}_{n,s}) + \log r(\mathbf{c}_n | \mathbf{z}_{n,s}))$$

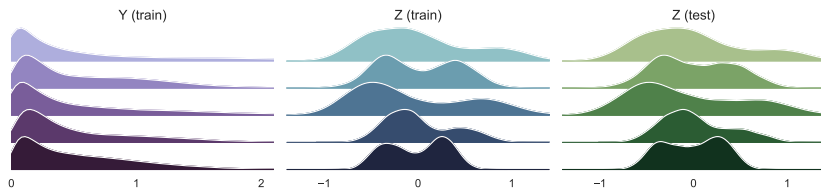
Variational approximation  $r(\mathbf{c} | \mathbf{y})$ : we use the learned classifier:

$$r(\mathbf{c} | \mathbf{z}) = \mathcal{S}(\mathbf{W}_d \mathbf{z} + \mathbf{b})_c$$

Approximation  $q(\mathbf{z})$ 

$$\Omega_{\text{REVE}} = -\frac{1}{NS} \sum_{n=1}^N \sum_{s=1}^S (\log q(\mathbf{z}_{n,s}) + \log r(\mathbf{c}_n | \mathbf{z}_{n,s}))$$

Variational approximation  $q(\mathbf{z})$ : how to model  $Z$ ?



# Bimodal Approximation

## Model

→ Independence between coordinates:

$$q(\mathbf{z}) = \prod_{i=1}^{\dim(\mathcal{Z})} q(z_i)$$

→ Bimodal approximation:

$$q(z_i) = \alpha_i \mathcal{N}(z_i | \mu_{1,i}, \sigma_{1,i}^2) + (1 - \alpha_i) \mathcal{N}(z_i | \mu_{0,i}, \sigma_{0,i}^2).$$

→  $\alpha_i$ : probability of the semantic attribute being present.

→  $\alpha, \mu_1, \sigma_1^2, \mu_0, \sigma_0^2$  to determine.

# Compute the Parameters

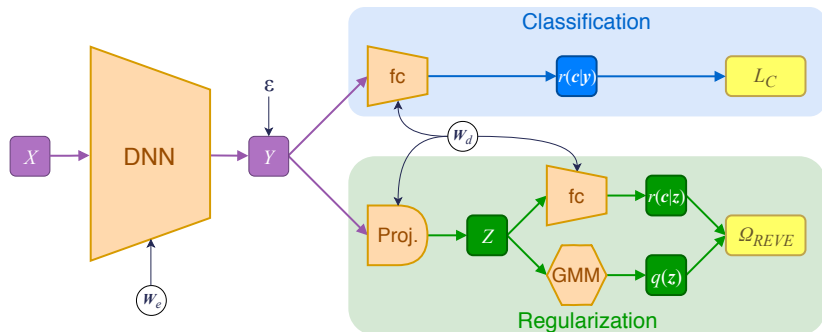
## Expectation Maximization for two Gaussian Mixture Model?

- Expensive
- The size of the mini-batches is in general too small for obtaining a coherent model

## Mini-batch Computation

- We assume  $P(M_i = 1 | z_i) = \sigma(z_i)$  (e.g.  $\sigma$  the sigmoid function)
- On the mini-batch,  $\alpha_j = \sum_n \sigma(z_i^{(n)})$ ,  $\mu_{1,j} = \sum_n \frac{\sigma(z_i^{(n)})}{\alpha_j} z_i^{(n)}, \dots$

## REVE Loss



$$\Omega_{REVE}((x_n, c_n); W_e; W_d) = -\frac{1}{NS} \sum_{n=1}^N \sum_{s=1}^S \left[ \log \mathcal{S}(W_d z_{n,s} + b)_{c_n} \right. \\ \left. + \sum_{i=1}^{\dim(Z)} \log \left( \alpha_i \mathcal{N}(z_{n,s,i} | \mu_{1,i}, \sigma_{1,i}^2) + (1 - \alpha_i) \mathcal{N}(z_{n,s,i} | \mu_{0,i}, \sigma_{0,i}^2) \right) \right].$$

# Performance Analysis

Reve Performance Analysis. Classification error (%) results on CIFAR-10 and CIFAR-100 test sets.

	CIFAR-10		CIFAR-100	
	AlexNet	Inception	AlexNet	Inception
Baseline	15.62	6.10	48.29	27.36
SGM Reve	14.24	6.17	-	-
KDE Reve	13.86	6.04	-	-
Reve	13.92	5.92	48.07	26.94
Reve + DO	<b>12.54</b>	<b>5.78</b>	<b>41.13</b>	<b>26.02</b>

## Results

### Classification error (%) results on CIFAR test sets.

	CIFAR-10			CIFAR-100		
	AlexNet	Inception Net	ResNet	AlexNet	Inception Net	ResNet
Baseline	15.62	6.10	4.08	48.29	27.36	20.70
Dropout	12.63	6.04	3.93	41.32	27.26	20.16
Information DO	14.97	6.04	NC	47.97	27.34	NC
Shade + DO	13.93	5.90	4.30	41.25	26.99	20.37
Reve + DO	<b>12.54</b>	<b>5.78</b>	<b>3.88</b>	<b>41.13</b>	<b>26.02</b>	<b>20.05</b>

### Classification error (%) results on SVHN test set.

	SVHN		
	AlexNet	Inception Net	ResNet
Baseline	7.68	3.78	3.40
Reve	<b>6.55</b>	<b>3.29</b>	<b>3.11</b>

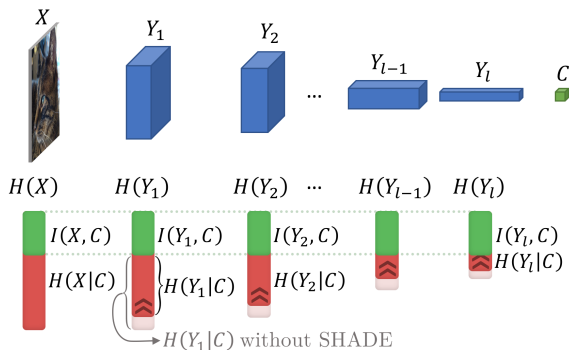
# Conclusion

- REVE is a tractable regularization loss for image classification
- Identifies the part of the representation orthogonal to the kernel of the classifier as the variable to constrain
- Penalizes the conditional entropy of the REVE variable given the class
- REVE shows consistent positive results on multiple architectures and datasets



# Questions?

Thank you for your attention!



$\approx$  Effect of SHADE: Reduction of conditional entropy

- Layer-wise criterion:  $\Omega_{\text{SHADE}} = \sum_{l=1}^L \sum_{i=1}^{D_l} H(Y_{l,i} | C)$
- Uses a latent Bernoulli variable  $B$  as minimal sufficient statistic of  $C$  for  $Y$ :  $I(Y, C) = I(Y, B)$