# Evaluation of Generative Models: *Practice*

Hao Dong

Peking University

# Evaluation of Generative Models: Practice

- DCGAN Evaluation
  - Classification accuracy
  - LPIPS
- VAE Evaluation
  - NLL
  - Beta-VAE metric
  - MIG
  - Clustering
- Others
  - Model Size
  - Tensorlayer Model.weights

- DCGAN Evaluation
  - **Classification accuracy**
  - LPIPS
- VAE Evaluation
  - NLL
  - Beta-VAE metric
  - MIG
  - Clustering
- Others
  - Model Size
  - Tensorlayer Model.weights

# DCGAN Evaluation: Classification Accuracy

- Given pretrained DCGAN for MNIST, how to evaluate it?

- Classification accuracy
  - Use the discriminator's convolutional features from all layers
  - Maxpooling each layers representation to produce a 4 × 4 spatial grid
  - Flatten and concatenate these features to form a 28672 dimensional vector
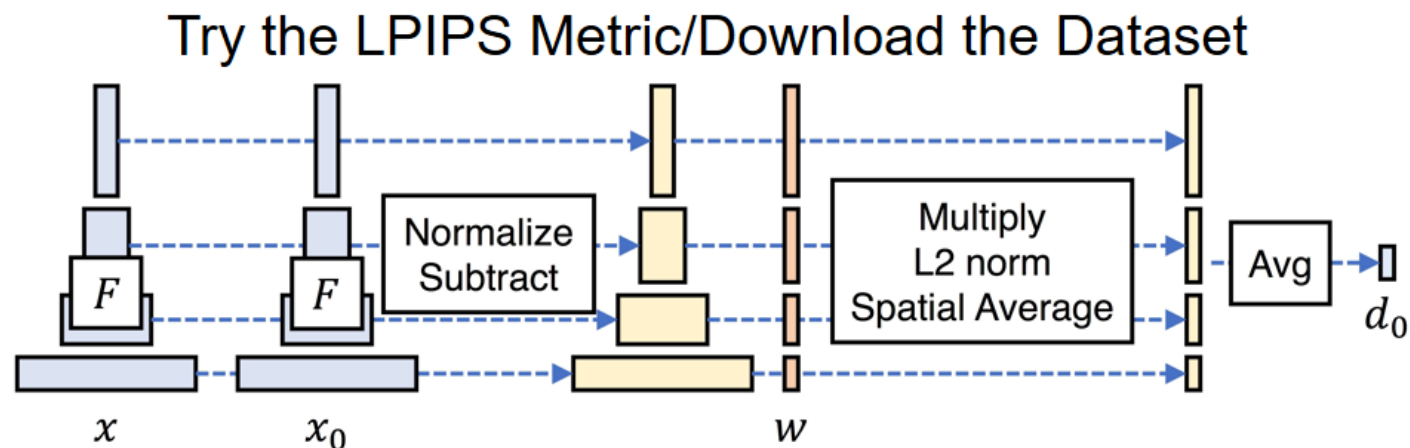  - A regularized linear L2-SVM classifier is trained on top of them

Table 1: CIFAR-10 classification results using our pre-trained model. Our DCGAN is not pre-trained on CIFAR-10, but on Imagenet-1k, and the features are used to classify CIFAR-10 images.

| Model | Accuracy | Accuracy (400 per class) | max # of features units |
|---|---|---|---|
| 1 Layer K-means | 80.6% | 63.7% ($\pm$0.7%) | 4800 |
| 3 Layer K-means Learned RF | 82.0% | 70.7% ($\pm$0.7%) | 3200 |
| View Invariant K-means | 81.9% | 72.6% ($\pm$0.7%) | 6400 |
| Exemplar CNN | 84.3% | 77.4% ($\pm$0.2%) | 1024 |
| DCGAN (ours) + L2-SVM | 82.8% | 73.8% ($\pm$0.4%) | 512 |

- DCGAN Evaluation
  - Classification accuracy
  - **LPIPS**
- VAE Evaluation
  - NLL
  - Beta-VAE metric
  - MIG
  - Clustering
- Others
  - Model Size
  - Tensorlayer Model.weights

# DCGAN Evaluation: LPIPS

- Given pretrained DCGAN for MNIST, how to evaluate it?

- **Learned Perceptual Image Patch Similarity (LPIPS)**
  - To evaluate the diversity of the generation
  - Perceptual similarity is an emergent property shared across deep visual representations.



- **Implementation:** https://github.com/richzhang/PerceptualSimilarity

- DCGAN Evaluation
  - Classification accuracy
  - LPIPS
- VAE Evaluation
  - **NLL**
  - Beta-VAE metric
  - MIG
  - Clustering
- Others
  - Model Size
  - Tensorlayer Model.weights

# VAE Evaluation

- Given pretrained VAE models for MNIST, how to evaluate it?

- **Negative Log Likelihood (NLL)**

  - NLL represents the probability of generating real data

  - Less NLL indicated better generation of VAE

- DCGAN Evaluation
  - Classification accuracy
  - LPIPS
- VAE Evaluation
  - NLL
  - **Beta-VAE metric**
  - **MIG**
  - Clustering
- Others
  - Model Size
  - Tensorlayer Model.weights

# VAE Evaluation

- Given pretrained VAE models for MNIST, how to evaluate it?

- **Beta-VAE metric and Mutual Information Gap (MIG)**
  - To evaluate the disentanglement of VAE.
  - Beta-VAE metric is the accuracy of a linear classifier that predicts a fixed factor of variation
  - MIG is the gap between the largest and second largest mutual information
  - Review lecture 20 for more details

| | | |
|---|---|---|
| 📄 beta_vae.py | internal change | 5 months ago |
| 📄 beta_vae_test.py | internal change | 5 months ago |
| 📄 mig.py | internal change | 5 months ago |
| 📄 mig_test.py | internal change | 5 months ago |

- **ICML 2019 Best Paper**
- Implementation:
https://github.com/google-research/disentanglement_lib/tree/master/disentanglement_lib/evaluation/metrics

- DCGAN Evaluation
  - Classification accuracy
  - LPIPS
- VAE Evaluation
  - NLL
  - Beta-VAE metric
  - MIG
- **Clustering**
- Others
  - Model Size
  - Tensorlayer Model.weights

# VAE Evaluation: Clustering

- Given pretrained VAE for MNIST, how to evaluate it?

- Clustering
  - **Completeness score** (between [0, 1])
  - **Homogeneity score** (between [0, 1])
  - **V measure score** (also called normalized mutual information, between [0, 1])

$$c = 1 - \frac{H(K|C)}{H(K)} \qquad h = 1 - \frac{H(C|K)}{H(C)} \qquad v = 2 \cdot \frac{h \cdot c}{h + c}$$

- Review lecture 20's slides for the implementation

- DCGAN Evaluation
  - Classification accuracy
  - LPIPS
- VAE Evaluation
  - NLL
  - Beta-VAE metric
  - MIG
  - Clustering
- Others
  - **Model Size**
  - Tensorlayer Model.weights

# Model Size

- The size of model is also an important metric of generative models
  - The size is the number of parameters of the model
  - It indicates the scalability of the model
  - Less parameters required, stronger scalability of the model

- Example: StarGAN Evaluation

| Method | Classification error | # of parameters |
|---|---|---|
| DIAT | 4.10 | 52.6M × 7 |
| CycleGAN | 5.99 | 52.6M × 14 |
| IcGAN | 8.07 | 67.8M × 1 |
| StarGAN | **2.12** | **53.2M × 1** |
| Real images | 0.45 | - |

Table 3. Classification errors [%] and the number of parameters on the RaFD dataset.

- The smallest size of StarGAN indicated its advantage in multi-domain translation

- DCGAN Evaluation
  - Classification accuracy
  - LPIPS
- VAE Evaluation
  - NLL
  - Beta-VAE metric
  - MIG
  - Clustering
- Others
  - Model Size
  - **Tensorlayer Model.weights**

# Implementation: tensorlayer Model.weight

- Calculate the size of a model using tensorlayer is convenient
- NLL is a term in the loss of classical VAE

```
1194    def count_weights(model):
1195        n_weights = 0
1196        for i, w in enumerate(model.all_weights):
1197            n = 1
1198            # for s in p.eval().shape:
1199            for s in w.get_shape():
1200                try:
1201                    s = int(s)
1202                except:
1203                    s = 1
1204                if s:
1205                    n = n * s
1206            n_weights = n_weights + n
1207        print("num of weights (parameters) %d" % n_weights)
1208        return n_weights
```

```
1211    if __name__ == '__main__':
1212        Ea = get_Ea()
1213        Ea.eval()
1214        print("Ea:")
1215        count_weights(Ea)
1216
1217        Ec = get_Ec()
1218        Ec.eval()
1219        print("Ec:")
1220        count_weights(Ec)
1221
1222        D = get_D()
1223        D.eval()
1224        print("D:")
1225        count_weights(D)
1226
1227        G = get_G()
1228        G.eval()
1229        print("G:")
1230        count_weights(G)
```

- Try to evaluate the size of DCGAN and VAE by yourself!

# Summary

- DCGAN Evaluation
  - Classification accuracy
  - LPIPS
- VAE Evaluation
  - NLL
  - Beta-VAE metric
  - MIG
  - Clustering
- Others
  - Model Size
  - Tensorlayer Model.weights

# Thanks