# Adversarial Autoencoders (AAE)
## Adapted from Makhzani et al., 2015
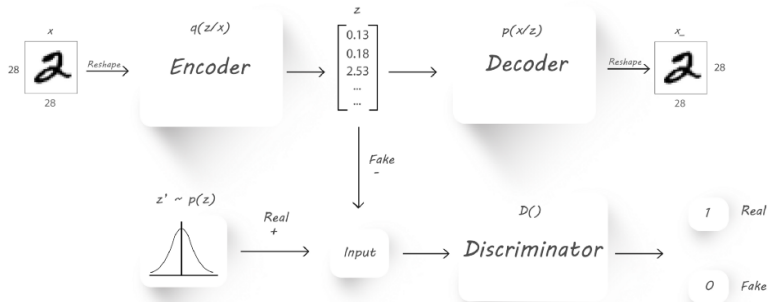
Giovanni Papini

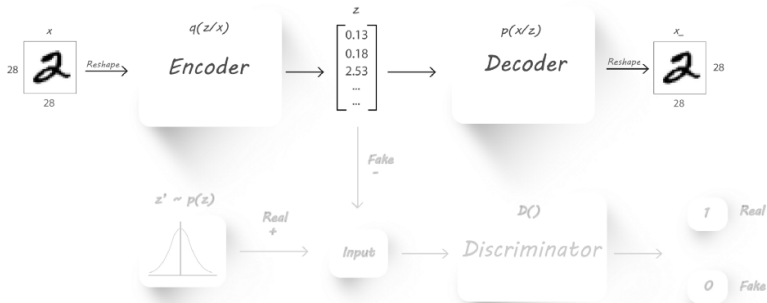Universitá degli Studi di Firenze

February 12, 2019

# Context

- Scalable generative models to capture rich distributions
- Graphical models such as RBMs , DBNs, DBMs are based on MCMC algorithms for doing inference
- VAE [Kingma and Welling, 2013], GAN [Goodfellow et al., 2014], GMMN [Li, Swersky, and Zemel, 2015] are trained via direct back-propagation
- The AAE is trained via back-propagation with dual objectives:
  - **minimizing the reconstruction error** — $||\mathbf{x} - \hat{\mathbf{x}}||^2$
  - **adversarial training criterion** — matching the aggregated posterior distribution of the latent representation to an arbitrary prior
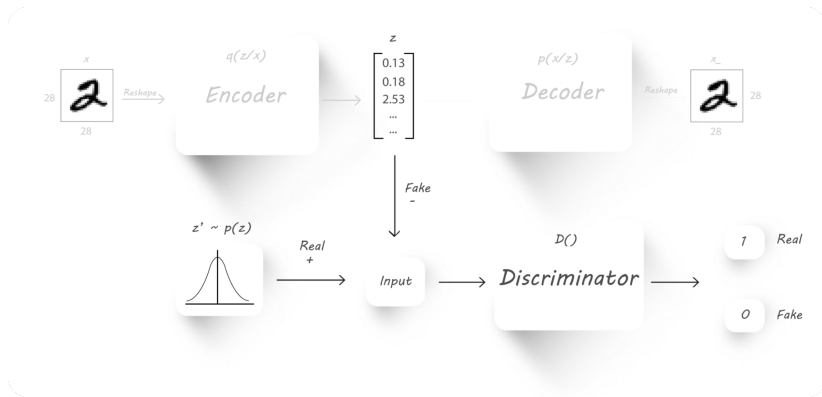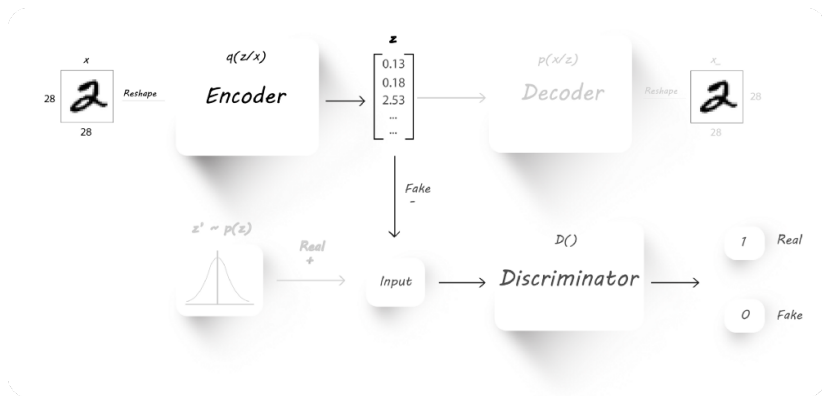
# The AAE architecture

# The AAE architecture

# The AAE architecture

# The AAE architecture

# The encoder

$$q(\boldsymbol{z}) = \int_{\Omega_x} q(\boldsymbol{z}|\boldsymbol{x}) p_d(\boldsymbol{x}) \mathrm{d}\boldsymbol{x} \rightarrow p(\boldsymbol{z})$$

○ **Deterministic** — (used in the paper)

○ **Gaussian posterior** — $z_i|\boldsymbol{x} \sim \mathcal{N}(\mu_i(\boldsymbol{x}), \sigma_i(\boldsymbol{x}))$

○ **Universal approximator posterior** — $\boldsymbol{z}|\boldsymbol{x}, \boldsymbol{\eta} \sim \delta(\boldsymbol{z} - f(\boldsymbol{x}, \boldsymbol{\eta}))$
   where $\boldsymbol{\eta}$ is random noise with a fixed distribution

$$q(\boldsymbol{z}|\boldsymbol{x}) = \int_{\Omega_\eta} q(\boldsymbol{z}|\boldsymbol{x}, \boldsymbol{\eta}) p_\eta(\boldsymbol{\eta}) \mathrm{d}\boldsymbol{\eta}$$

# Relationship with VAEs

VAE aims to minimize the negative ELBO:

$$-\log p(\mathbf{x}) < \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})}[-\log p(\mathbf{x}|\mathbf{z})] + \mathbb{KL}(q(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$
$$= \mathbb{E}_{\mathbf{z}}[-\log p(\mathbf{x}|\mathbf{z})] - H(q(\mathbf{z}|\mathbf{x})) + \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})}[-\log p(\mathbf{z})]$$
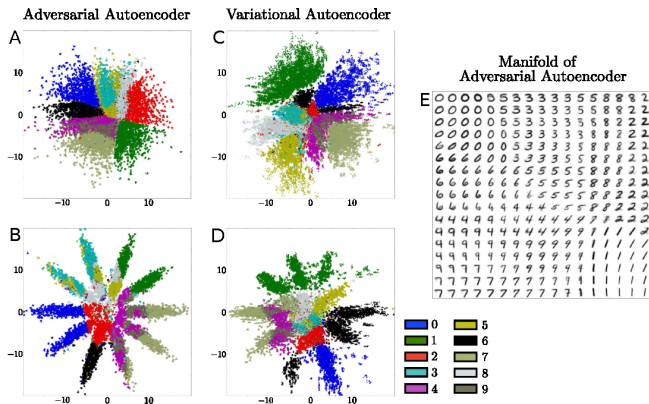$$= \text{Reconstr. Err.} - \text{Entropy} + \text{Crossentropy}$$

| VAE | AAE |
|---|---|
| ○ encourages $q(\mathbf{z})$ to match $p(\mathbf{z})$ modes because of crossentropy penalty | ○ encourages $q(\mathbf{z})$ to match the *whole* $p(\mathbf{z})$ because of the adversarial training |
| ○ needs the exact functional form of the prior distribution in order to back-propagate through KL divergence | ○ can impose even complicated distributions just through the capability of sampling from them |

# Relationship with VAE



Adversarial Autoencoder

Variational Autoencoder

Manifold of
Adversarial Autoencoder

# Relationship with GAN and GMMN

| GAN | AAE |
|---|---|
| ○ imposes the data distribution to the output of a neural network | ○ relies on the autoencoder to capture the data distribution <br><br> ○ shapes a much lower dimensional space into a much simpler distribution |

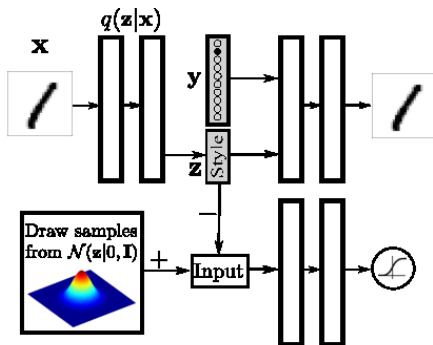| GMMN | AAE |
|---|---|
| ○ (first) trains a dropout autoencoder (then) fits a distribution in the code-space of the pretrained network | ○ uses adversarial training as a regularizer that shapes the code distribution while training the autoencoder from scratch |

# Likelihood analysis

- Benchmarks on: MNIST, <u>Toronto Face Dataset</u> (TFD)
- Not direct likelihood measure, but lower bound approximation (KDE using Gaussian Parzen window, $\sigma$ selected by cross-validation)
- Samples of 10K and 10M units to estimate the test set log-likelihood

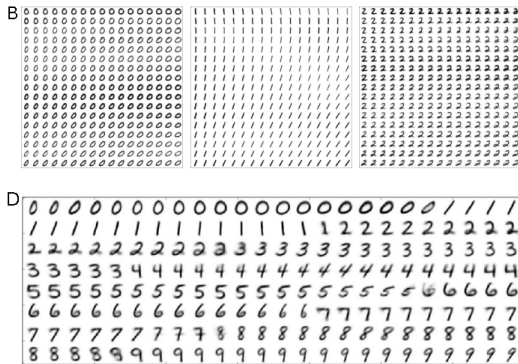|  | MNIST(10K) | MNIST(10M) | TFD(10K) | TFD(10M) |
|---|---|---|---|---|
| DBN | $138 \pm 2.0$ | - | $1909 \pm 66$ | - |
| Stacked CAE | $121 \pm 1.6$ | - | $2110 \pm 50$ | - |
| Deep GSN | $214 \pm 1.1$ | - | $2890 \pm 29$ | - |
| GAN | $225 \pm 2.0$ | 386 | $2057 \pm 26$ | - |
| GMMN + AE | $282 \pm 2.0$ | - | $2204 \pm 20$ | - |
| **AAE** | $340 \pm 2.0$ | 427 | $2252 \pm 16$ | 2522 |

# Semi-supervised approach
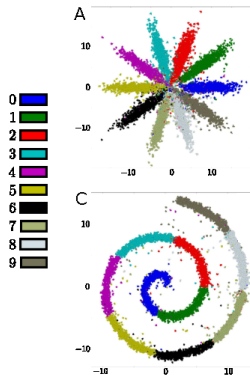
○ Incorporate one-hot vector in the latent representation
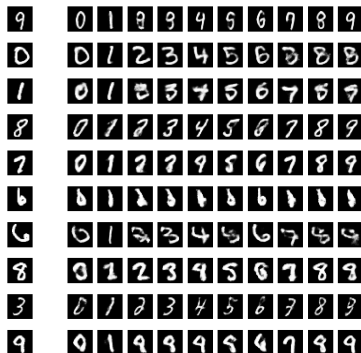
# Semi-supervised approach

Incorporating label information
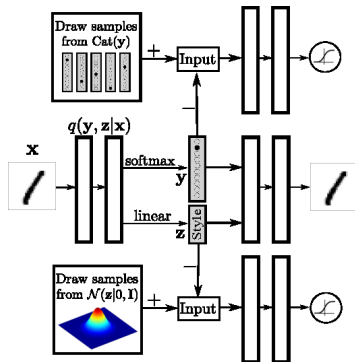
# Semi-supervised approach

Incorporating label information

- ○ The AAE disentangles *style* features from *content/semantic* features
- ○ Experiments on MNIST and <u>Street View House Number</u> dataset (SVHN)

# Semi-supervised classification

Architecture



- Improve classification performance using both labeled and unlabeled data
- Assume the latent space is a mixed Categorical and Gaussian distribution

$$p(\boldsymbol{y}) = \mathrm{Cat}(\boldsymbol{y}) \quad p(\boldsymbol{z}) = \mathcal{N}(\boldsymbol{z}; 0, \mathbf{I})$$

- For each distribution a different adversarial network regularizes the latent representation

# Semi-supervised classification

Training

3-phase training, all with SGD:

1. *reconstruction* — the AE updates the encoder $q(\boldsymbol{z}, \boldsymbol{y}|\boldsymbol{x})$ to minimize the reconstruction error
2. *regularization* — each adversarial network
   - first updates the discriminative network to distinguish true samples from the generated samples
   - then update the encoder to confuse the discriminative networks
3. *semi-supervised classification* — the AE updates $q(\boldsymbol{z}, \boldsymbol{y}|\boldsymbol{x})$ to minimize the cross-entropy cost on a labeled minibatch

# Semi-supervised classification

|  | MNIST(100) | MNIST(1000) | MNIST(All) | SVHN(1000) |
|---|---|---|---|---|
| NN Baseline | 25.80 | 8.73 | 1.25 | 47.5 |
| VAE (M1) + TSVM | $11.82 \pm 0.03$ | $4.24 \pm 0.07$ | - | $55.33 \pm 0.11$ |
| VAE (M2) | $11.97 \pm 1.71$ | $3.60 \pm 0.56$ | - | $36.02 \pm 0.10$ |
| VAE (M1 + M2) | $3.33 \pm 0.14$ | $2.40 \pm 0.02$ | 0.96 | 24.63 |
| CatGAN | $1.91 \pm 0.10$ | $1.73 \pm 0.18$ | 0.91 | - |
| Ladder Networks | $1.06 \pm 0.37$ | $0.84 \pm 0.08$ | $0.57 \pm 0.02$ | - |
| ADGM | $0.96 \pm 0.10$ | - | - | $16.61 \pm 0.24$ |
| **AAE** | $1.90 \pm 0.10$ | $1.60 \pm 0.08$ | $0.85 \pm 0.02$ | $17.70 \pm 0.30$ |

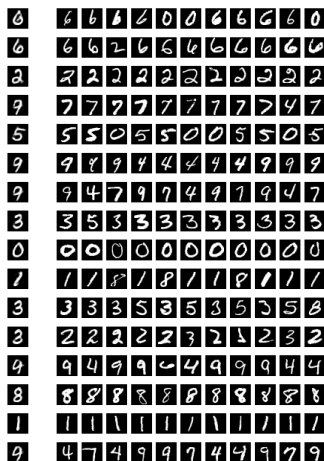Table: Error rate on semi-supervised classification.

# Unsupervised clustering

- Similar to the semi-supervised architecture, but without the semi-supervised training stage
- Not necessarily 10 classes, arbitrary number of clusters
- Benchmark criterion: find $\arg\max_{x_n} q(y_i|x_n)$ and assign label of $x_n$ to all elements of $i$-th cluster, then compute error rate based on the assigned class labels

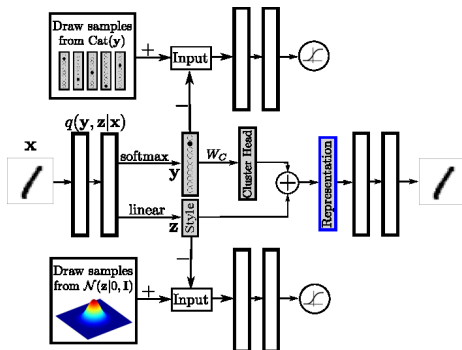|                       | MNIST (Unsupervised) |
|-----------------------|:--------------------:|
| CatGAN (20 clusters)  | 9.70                 |
| AAE (16 clusters)     | $9.55 \pm 2.05$      |
| AAE (30 clusters)     | $4.10 \pm 1.13$      |

Table: Error-rate on unsupervised clustering

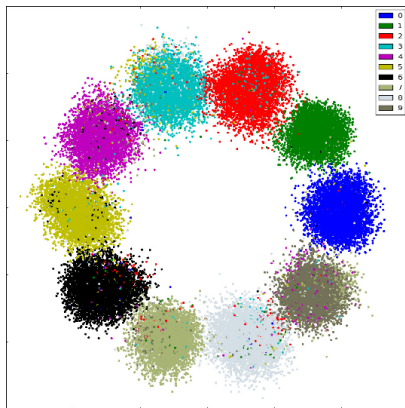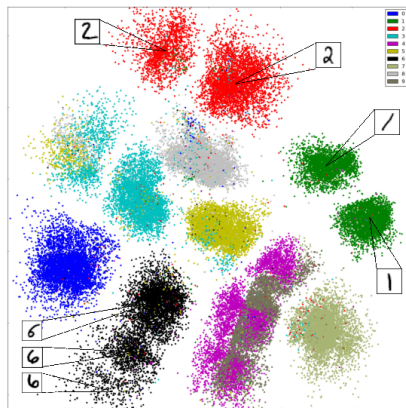# Unsupervised clustering

# Dimensionality reduction

Architecture

- ○ Typically, non-regularized autoencoders fracture the manifolds ⇒ very different codes for similar images
- ○ Little modification label-integrated AAE architecture



- ○ Additional cost to penalize euclidean distance between cluster heads

# Dimensionality reduction

○ Reducing the dimensionality of the hidden space (from 10 to 2) can impact on the predictive power of the whole network

○ Ad hoc tricks can balance this trade-off

# Conclusions

**Pros**

- ○ it is framework capable of modeling complex distributions only requiring to be able to sample from it
- ○ the AAE is highly flexible, could be combined with variational objectives (see Rosca et al., 2017)

**Cons**

- ○ like GANs, it requires much hyper-parameter tuning to perform at the top
- ○ it could suffer from complex adversarial game dynamics
- ○ it could be an overkill just to treat Gaussian/Gaussian-mixture distributions

\* It would need some real-world testing.

# References

Bengio, Yoshua et al. (2013). "Better mixing via deep representations". In: *International Conference on Machine Learning*, pp. 552–560.

Kingma, Diederik P and Max Welling (2013). "Auto-encoding variational bayes". In: *arXiv preprint arXiv:1312.6114*.

Bengio, Yoshua et al. (2014). "Deep generative stochastic networks trainable by backprop". In: *International Conference on Machine Learning*, pp. 226–234.

Goodfellow, Ian et al. (2014). "Generative adversarial nets". In: *Advances in neural information processing systems*, pp. 2672–2680.

Li, Yujia, Kevin Swersky, and Rich Zemel (2015). "Generative moment matching networks". In: *International Conference on Machine Learning*, pp. 1718–1727.

Makhzani, Alireza et al. (2015). "Adversarial autoencoders". In: *arXiv preprint arXiv:1511.05644*.

Maaløe, Lars et al. (2016). "Auxiliary deep generative models". In: *arXiv preprint arXiv:1602.05473*.

Nagabushan, Naresh (2017). *A wizard's guide to Adversarial Autoencoders: Part 2, Exploring latent space with Adversarial Autoencoders*. URL: https://towardsdatascience.com/a-wizards-guide-to-adversarial-autoencoders-part-2-exploring-latent-space-with-adversarial-2d53a6f8a4f9.

Rosca, Mihaela et al. (2017). "Variational approaches for auto-encoding generative adversarial networks". In: *arXiv preprint arXiv:1706.04987*.

Dr. Harang, Richard and Madeline Schiappa (2018). *Adversarial Autoencoders*. Tech. rep. Sophos. URL: https://www.sophos.com/en-us/medialibrary/PDFs/technical-papers/Adversarial-Autoencoders.pdf.