

Title: PSI 2018/2019 - Machine Learning - Lecture 15

Speakers: Michael Albergo

Collection: PSI 2018/2019 - Machine Learning (Hayward Sierens)

Date: April 12, 2019 - 9:00 AM

URL: <http://pirsa.org/19040010>

Yesterday

Explicit Likelihood

Today

2. Explicit but Approx.

3. Implicit Likelihood

2. Latent Variable

$$p(x) = \int_z p(x, z) = \int p(x|z) p(z) dz$$

$$\text{Approx. } p(x|z) \approx p_\theta(x|z)$$

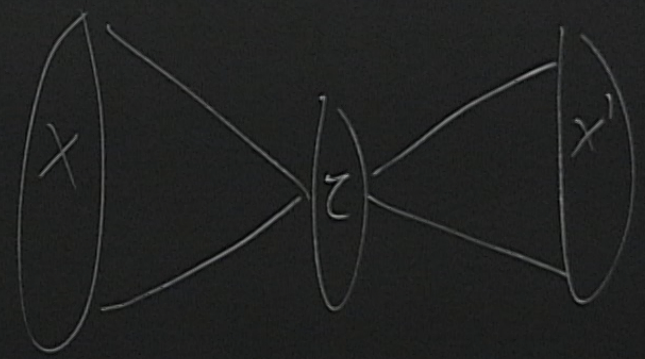
side

~~$p(x) p(z|x)$~~

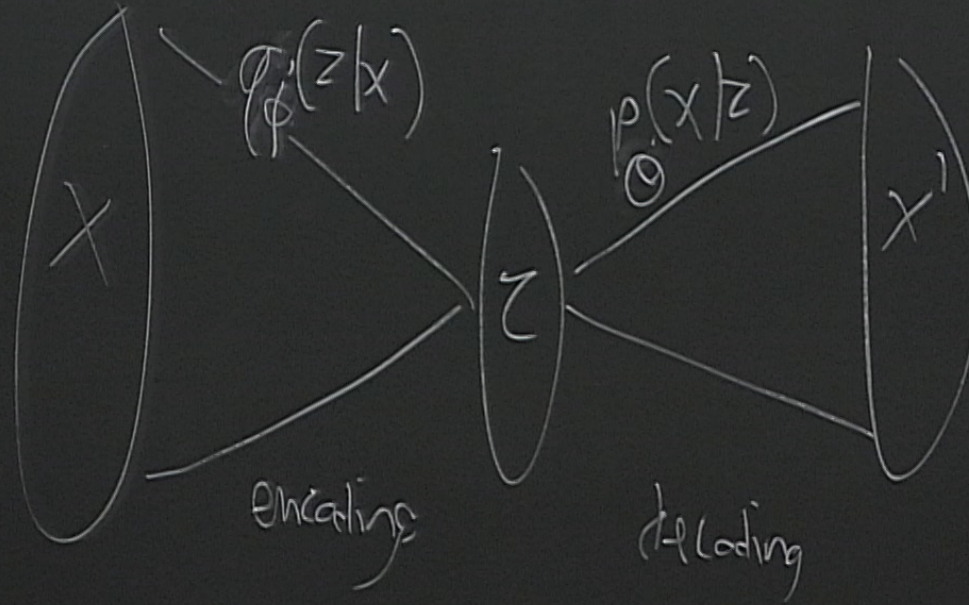
$$\int_z p(x,z) = \int p(x|z)p(z) dz$$

I need $p(z|x) p(z) p(x|z)$

$$p(z) \approx p_{\theta}(x|z)$$



I need $p(z|x)$ $p(z)$ $p(x|z)$



$$\begin{aligned}
\ln p(x) &= E_{z \sim q(z|x)} [\ln p(x)] \\
&= E_{z \sim q(z|x)} \left[\ln \frac{p_0(x|z)p(z)}{p_0(z|x)q_p(z|x)} \right] \\
&= E_{z \sim q} [\ln p(x|z)] - E_{z \sim q} \left[\ln \frac{q_p(z|x)}{p_0(z)} \right] + E_{z \sim q} \left[\ln \frac{q_p(z|x)}{p_0(z|x)} \right]
\end{aligned}$$

$$\ln p(x) = \mathbb{E}_{z \sim q(z|x)} [\ln p(x)]$$

$$= \mathbb{E}_{z \sim q(z|x)} \left[\ln \frac{p_0(x|z)p(z)}{p_0(z|x)} \frac{q_\phi(z|x)}{q_\psi(z|x)} \right]$$

$$= \mathbb{E}_{z \sim q} [\ln p(x|z)] - \mathbb{E}_{z \sim q} \left[\ln \frac{q_\phi(z|x)}{p_0(z)} \right] + \mathbb{E}_{z \sim q} [\ln p_0(z)]$$

$$\ln p(x) \geq -D_{KL}(q_\phi(z|x) \parallel p_0(z)) + D_{KL}(p_0(z) \parallel q_\psi(z|x))$$

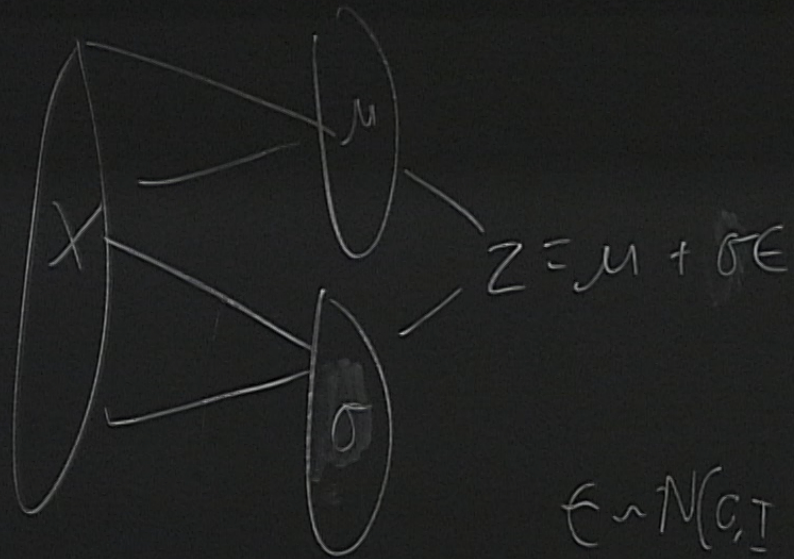
$$\frac{q_f(z|x)}{q_p(z|x)}$$

$$p_0(z) \sim N(0, I)$$

$$-E_{z \sim q} \left[\ln \frac{q_f(z|x)}{p_0(z)} \right] + E_{z \sim q} \left[\ln \frac{q_f(z|x)}{p_0(z|x)} \right]$$

$$-D_{KL}(q_f(z|x) \parallel p_0(z)) + \cancel{D_{KL}(q_f(z|x) \parallel p_0(z|x))} \geq 0$$

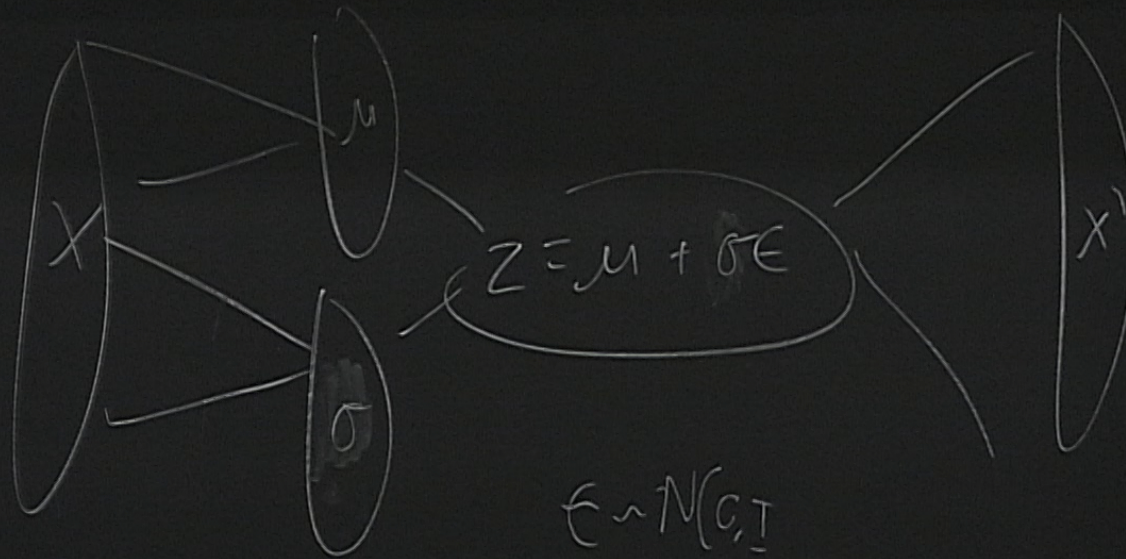
$N(0, I)$



$$Z = \mu + \sigma \epsilon$$

$$\epsilon \sim N(0, I)$$

$N(0, I)$



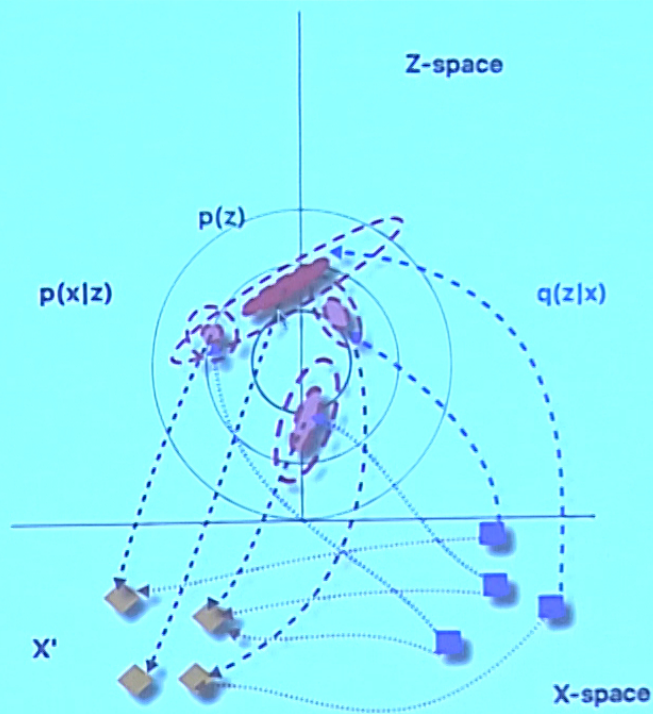


Figure 5: Visualizing how datapoints x are mapped into latent space, and potential redundancy in how they are decoded.

As the KL term is more strongly enforced, data will be better spread around the Gaussian prior, but the generated image quality will be worse. If the KL term is not enforced as strongly (as might be seen above), posterior Gaussians can overlap, and multiple inputs might map to the

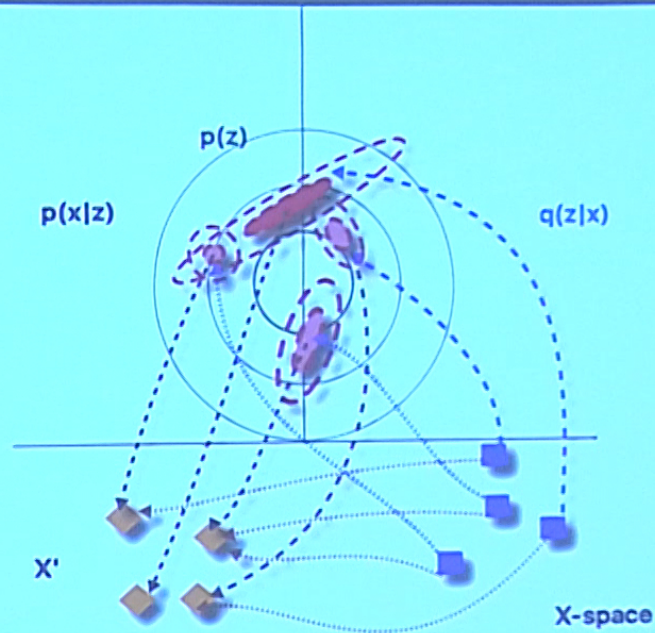


Figure 5: Visualizing how datapoints \mathbf{x} are mapped into latent space, and potential redundancy in how they are decoded.

As the KL term is more strongly enforced, data will be better spread around the Gaussian prior, but the generated image quality will be worse. If the KL term is not enforced as strongly (as might be seen above), posterior Gaussians can overlap, and multiple inputs might map to the same generated \mathbf{x}' . This can be improved with the Wasserstein Autoencoder by Tolstikhin et al in which $q(\mathbf{z}) = \int q(\mathbf{z}|\mathbf{x})dp(\mathbf{x})$ is forced to match $p(\mathbf{z})$ rather than each $q_\phi(\mathbf{z}|\mathbf{x})$ on their own.

Figure 5: Visualizing how datapoints \mathbf{x} are mapped into latent space, and potential redundancy in how they are decoded.

As the KL term is more strongly enforced, data will be better spread around the Gaussian prior, but the generated image quality will be worse. If the KL term is not enforced as strongly (as might be seen above), posterior Gaussians can overlap, and multiple inputs might map to the same generated \mathbf{x}' . This can be improved with the Wasserstein Autoencoder by Tolstikhin et al in which $q(\mathbf{z}) = \int q(\mathbf{z}|\mathbf{x})dp(\mathbf{x})$ is forced to match $p(\mathbf{z})$ rather than each $q_\phi(\mathbf{z}|\mathbf{x})$ on their own. **Note:** there are many many advances and perspectives on these models, and more keep coming. I just reference the Wasserstein Autoencoder case because it elucidates some of the behavior of the latent space visually well.

For example of some other advances in the understanding, one can take a look at ways of disentangling the latent space of these models so that each latent dimension has some more interpretable understanding of its representation. This can be done by enforcing independence between the latent dimensions, which becomes apparent if you break down the KL term into its newly discovered 3 subparts, one of which is the Total Correlation. This can be done adversarially, with d-Hilbert-Schmidt Independence Criteria, or with MMD and variations of adversarial training. One can also tailor their latent space to be a sampling on a manifold that corresponds to the symmetries/transformations that govern the true data.

Advantages:

z) $p(x|z)$

Adv.

Stable training

• Easy to sample $p(x|z)$

Disadv

• Monte Carlo to get $p(x)$

• MLE can limit resolution of images ★

it's definitely better

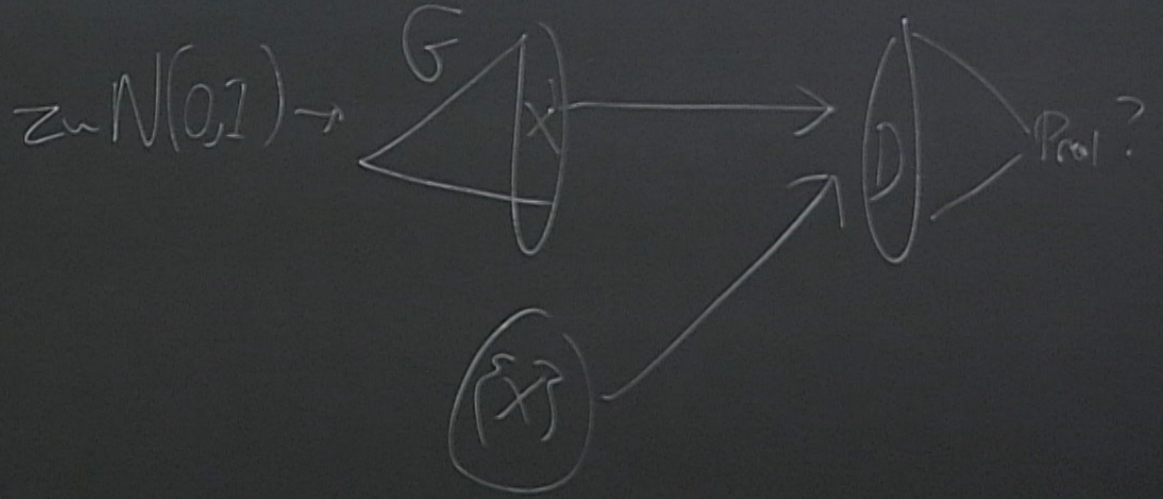
Implicit GM

GAN

$P(x) \rightarrow$ sample

$G_{\theta} \xrightarrow{D_{\phi}} P_{\theta}(x|z)$

$\{X\}$



$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim P_{\text{real}}} [\log D(x)] + \mathbb{E}_{z \sim P_g} [\log |1 - D(G(z))|]$$

$$V(D, G) = \int_x P_{\text{real}} \log D \, d\nu + \int_z P_g(z) \log (1 - D(G(z))) \, dz$$

$$D_{\phi} = \frac{P_{\text{real}}(x)}{P_{\text{real}}(x) + P_g(x)}$$

$$= \int_x P_{\text{real}} \log D(x) + P_g \log (1 - D(x)) \, dx$$

$$\begin{aligned}
 V^* &= \min_G V(D_{\phi}^*, G) = E_X \left[\log D_{\phi}^*(X) \right] + E_{X \sim P_g} \left[\log (1 - D_{\phi}^*(X)) \right] \\
 &= E_X \left[\log \frac{P_{real}}{P_{real} + P_g} \right] + E_{X \sim P_g} \left[\log \frac{P_g}{P_{real} + P_g} \right] = -\log 4 \\
 V^A &= -\log 4 + E_X \left[\log \frac{P_{real}}{(P_{real} + P_g)/2} \right] + E_{X \sim P_g} \left[\log \frac{P_g}{(P_{real} + P_g)/2} \right]
 \end{aligned}$$

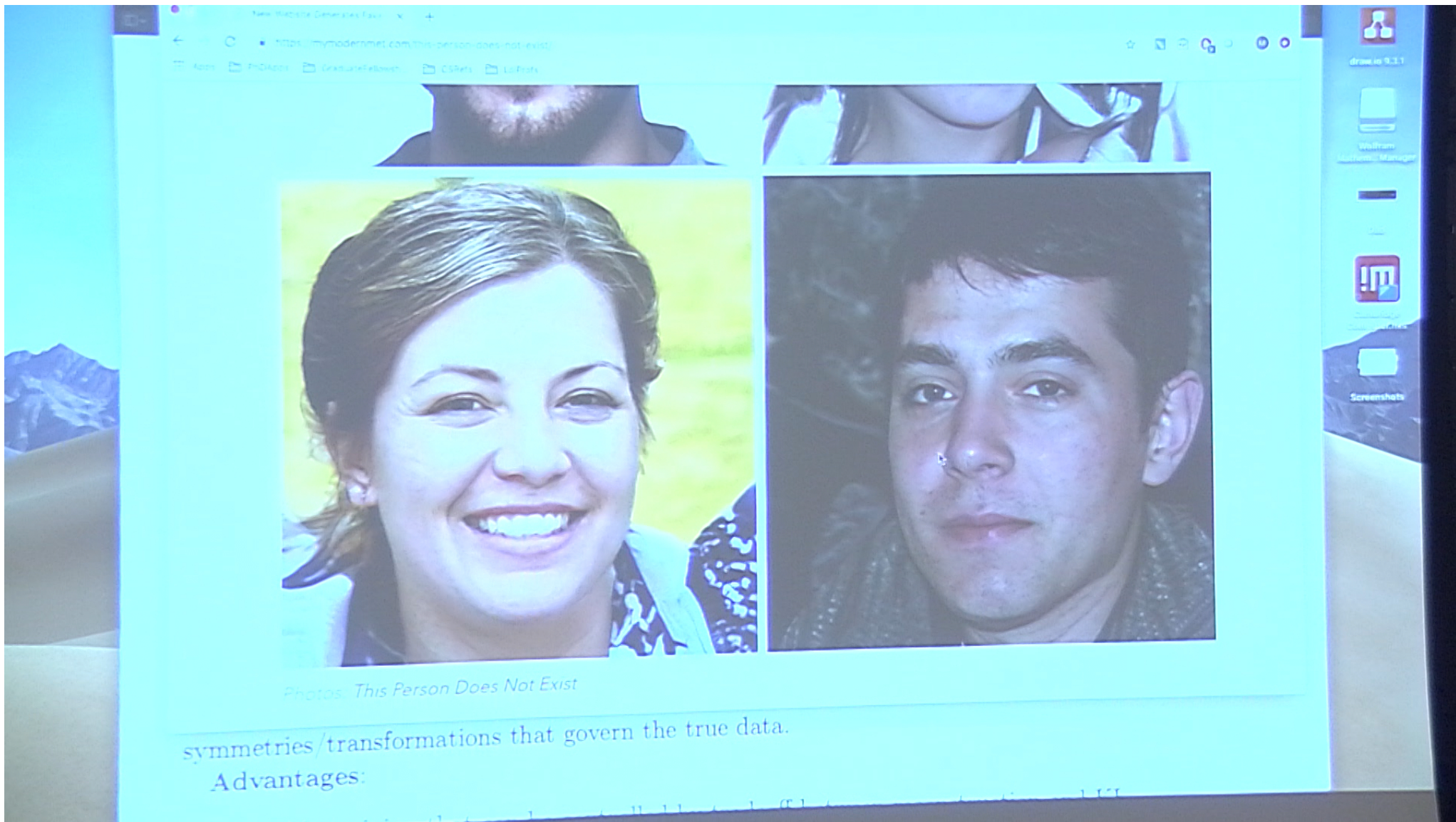
$$\begin{aligned}
G) &= -E_x \left[\log D_f^*(x) \right] + E_{x \sim P_g} \left[\log \left(1 - D_f^*(x) \right) \right] \\
&= E_x \left[\log \frac{P_{real}}{P_{real} + P_g} \right] + E_{x \sim P_g} \left[\log \frac{P_g}{P_{real} + P_g} \right] = -\log 4 \\
&= E_{x \sim P_g} \left[\log \frac{P_g}{(P_{real} + P_g)/2} \right] = -\log 4 + D_{JS} \left(P_{real} \parallel P_g \right)
\end{aligned}$$

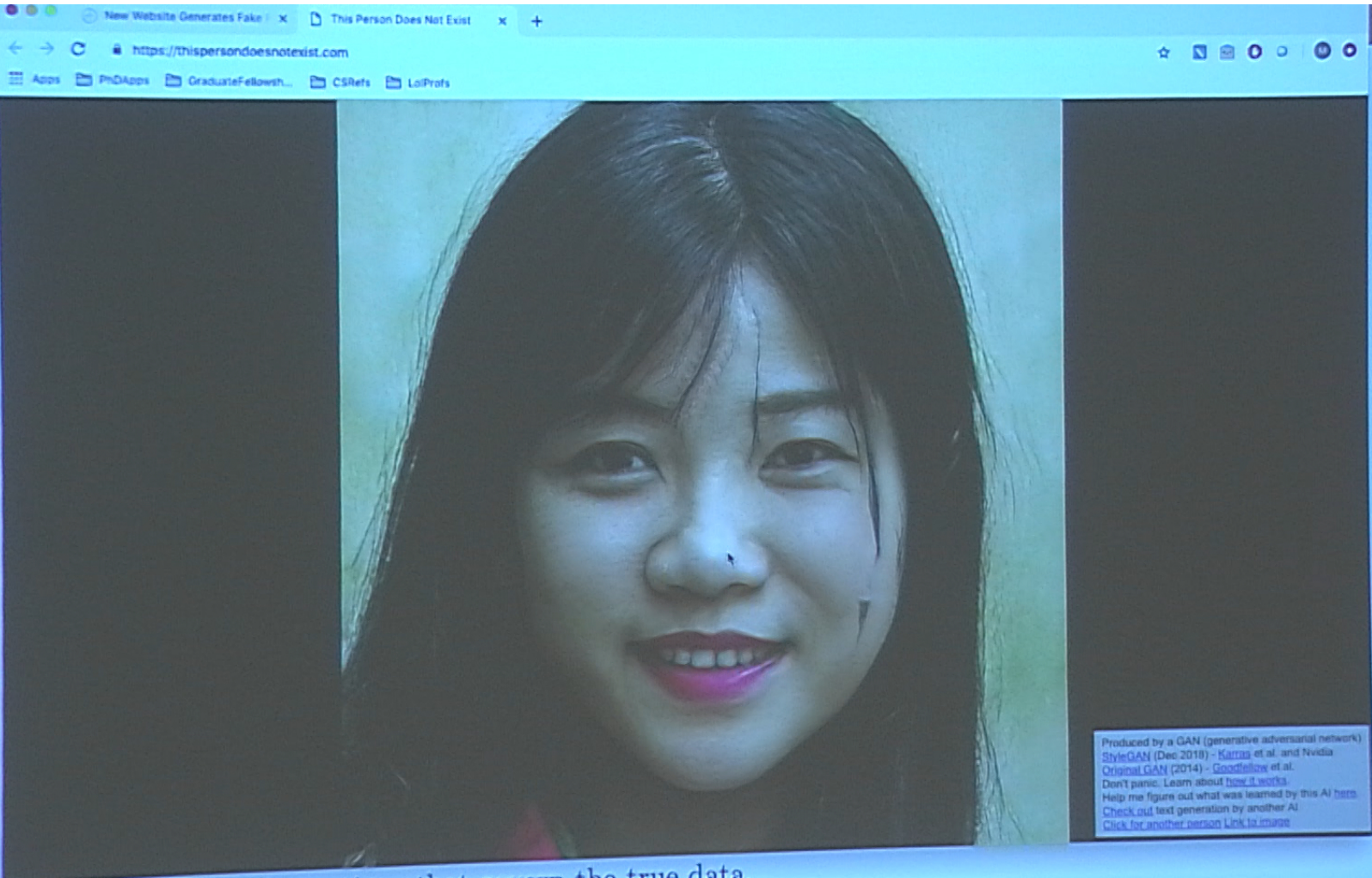
Wasserstein GAN

$$W(P_r, P_g) = \inf_{\gamma \in \Pi(P_r, P_g)} \mathbb{E}_{x, y \sim \gamma} [\|x - y\|]$$

$$W(P_r, P_g) = \sup_{\|f\|_L \leq 1} \left(\mathbb{E}_{x \sim P_r} (f(x)) - \mathbb{E}_{x \sim P_g} (f(x)) \right)$$

$$\begin{aligned}
 & \lambda \|x - y\| \\
 & \mathbb{E}_{x \sim p_S}(f(x)) - \mathbb{E}_{x \sim p_T}(f(x)) + \lambda \left\{ \mathbb{E}_{\lambda \sim p_\lambda} \left[\left(\|\nabla_x f(x)\|_2 - 1 \right)^2 \right] \right\} \\
 & \lambda = \frac{1}{2} + (1 - \frac{1}{2}) \lambda \in [0, 1]
 \end{aligned}$$





symmetries/transformations that govern the true data.

Advantages: