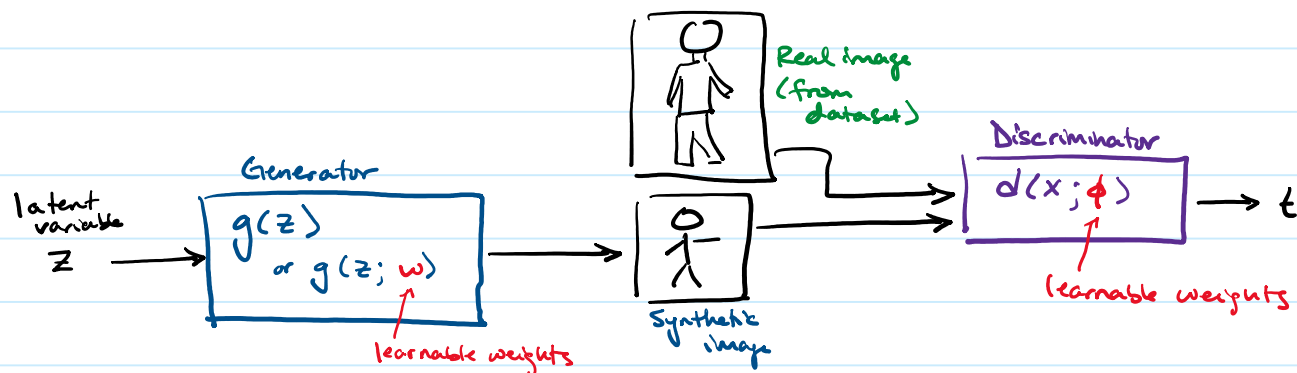


GANS (ch. 17 Bishop)

Sunday, November 17, 2024

4:37 PM

Generative Adversarial Networks (ch. 17 Bishop) Goodfellow et al. '14



$z \sim p(z)$, something simple we can sample from
(almost always $N(0, I)$ Gaussian)

Train generator and discriminator jointly!
 → tries to distinguish synthetic from real
 → tries to fool discriminator
 output: $t = 1$ if it thinks x is real
 $t = 0$ if it thinks x is Synthetic
 or more precisely,

$$P[t = 1] = d(x; \phi)$$

for discriminator, use usual cross-entropy loss: (last layer has sigmoid so output is in $(0, 1)$)

$$E(w, \phi) = -\frac{1}{N} \sum_{n=1}^N t_n \log(d_n) + (1 - t_n) \log(1 - d_n)$$

\uparrow true label \uparrow predicted
 $d_n = d(x_n; \phi)$

... but use a mix of real + synthetic

$$= -\frac{1}{N_{\text{real}}} \sum_{n \in \text{real}} \log(d_n) - \frac{1}{N_{\text{synth}}} \sum_{n \in \text{synth}} \log(1 - d(g(z_n, w), \phi))$$

$t_n = 1$ $t_n = 0$

key idea: minimize E with respect to ϕ (discriminator)

but maximize E with respect to w (generator)

$$\begin{aligned} \text{ie. } \phi &\leftarrow \phi - \eta \cdot \nabla_{\phi} E(w, \phi) && \text{gradient descent} \\ w &\leftarrow w + \eta \cdot \nabla_w E(w, \phi) && \text{gradient ascent} \end{aligned} \quad \left. \vphantom{\begin{aligned} \phi &\leftarrow \phi - \eta \cdot \nabla_{\phi} E(w, \phi) \\ w &\leftarrow w + \eta \cdot \nabla_w E(w, \phi) \end{aligned}} \right\} \text{alternate}$$

"Saddle-point problem"

$$\max_w \min_{\phi} E(w, \phi) \quad \left(\text{sometimes} = \min_{\phi} \max_w E(w, \phi) \right)$$

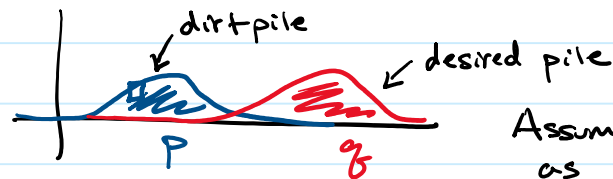
GANS (p. 2)

Tuesday, December 10, 2024

6:33 PM

Details, variants...

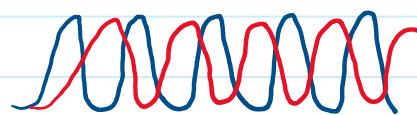
- conditional GANs
- not easy to train due to saddle-point nature
 - no clear metric of progress either
- mode collapse
 - eg. on MNIST, generator learns how to generate realistic "3"s but no other digits
 - a few proposed fixes
- Change metric to... **Wasserstein distance**
aka **Earth mover's distance**, part of **optimal transport**



Assuming you are as efficient as possible, how much work moving dirt is required?

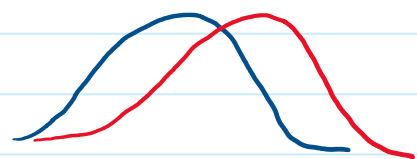
why?

- KL divergence requires shared support
(so $\log(\frac{p(x)}{q(x)})$ well-defined)
If barely so, can have issues



→ huge squared error loss

← similar Wasserstein loss



← small squared error loss

• CycleGANs

Used for things like style transfer

- Representation learning - features in latent space
path from z_1 to z_2 leads to meaningful path x_1 to x_2