

Recent Generative Models

What is a generative model?

↳ In ML, we see a lot of assigning class probability to outcomes
conditional $p(y|x) \rightarrow$ given image, the probability it is a cat, dog, etc

Generative model says: $p^{joint}(x, y) \rightarrow$ what's the probability of the image and the label?

↳ what is the distribution that generates the data? \rightarrow unsupervised, just $p(\tilde{x})$

• more powerful learned representation. Can use it to estimate $p(y|x)$

↑
no relational
mapping

$$p(x, y) = p(x|y)p(y) \propto p(y|x)$$

• can use to generate synthetic data \rightarrow can check in advance if model can explain any data.

• The issue is that it is generally difficult to estimate the most general of parametric distributions. Estimating likelihood $p(\tilde{x})$ is generally intractable. \star New models have clever ways of circumventing this issue.

\star when you assume hidden representation. See P. 121/122 for counter ex.

Circumstance

• Say we are given some $X = \{\tilde{x}_i\}_{i=1}^N$ for some discrete or continuous.

• Say X are actually generated from some hidden variable $\tilde{z} \sim p(\tilde{z})$

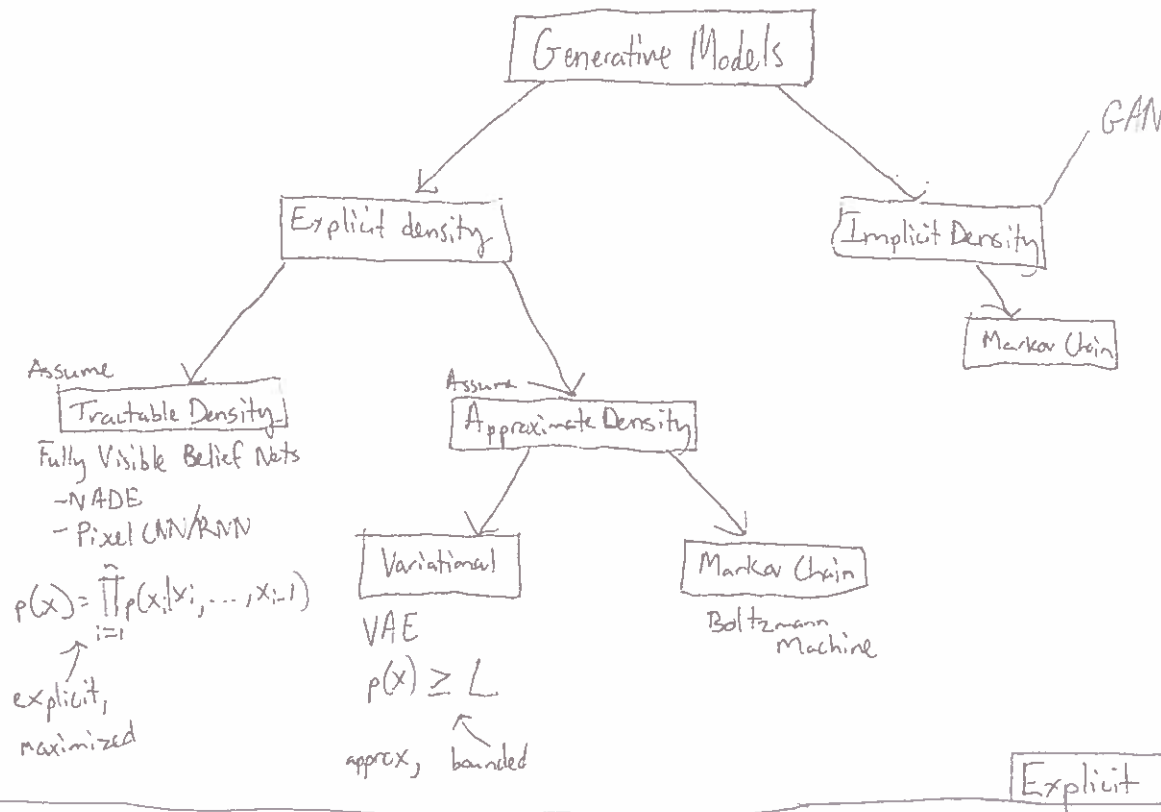
• We want to make inference on the data and its generating process, but all the components to do this seem intractable to compute, namely:

marginal likelihood $p(\tilde{x}) = \int p(\tilde{x}|\tilde{z})p(\tilde{z})d\tilde{z}$ $\xrightarrow{\text{hard}}$ posterior density $p(\tilde{z}|\tilde{x}) = \frac{p(\tilde{x}|\tilde{z})p(\tilde{z})}{p(\tilde{x})} \leftarrow \text{hard}$

• w/ the nonlinearities of NNs and large datasets now used.

This is the circumstance ^{some} new models try to circumvent.

Generative Model Taxonomy



PixelRNN, for example, by van der Oord et al. 2016

- Dependency on previous pixels modeled using an RNN (LSTM). They work very well
- Drawback → Sequential Generation can be slow for many pixels

PixelCNN, van der Oord et al 2016

- Previous pixel dependency modeled w/ CNN over neighboring region. Could be interesting for 2D grids
- No longer looking at "previous" but rather proximal.
- Softmax at each pixel location, on images for $[0, \dots, 255]$ possible values
- faster training b/c conv, but still slow generation

- explicit $p(x)$ → good evaluation metric, good samples
- slow

- Main idea: if you can't maximize $p(x)$ b/c of this intractability, can you at least put a lower bound on it?

We want to be able to sample from $p(z)$ and use that sample to generate a sample from $p(x|z)$ ← this will help us generate new samples that approximately come from original distribution

We will work with 2 conditional distributions: $p(x|z)$ modeled by $p_\theta(x|z)$

$q(z|x)$ modeled by $q_\phi(z|x)$

where θ, ϕ are the set of parameters for 2 neural networks

We use these to try and maximize our estimate of the evidence of the data $p_\theta(x)$:

$$\ln p_\theta(x) = E_{z \sim q_\phi(z|x)} [\ln p_\theta(x, z)]$$

$$= E_z \left[\ln \frac{p_\theta(x|z) p_\theta(z)}{p_\theta(z|x)} \right]$$

$$= E_z [\ln p_\theta(x|z)] - E_z \left[\ln \frac{q_\phi(z|x)}{p_\theta(z)} \right] + E_z \left[\ln \frac{q_\phi(z|x)}{p_\theta(z|x)} \right]$$

$$= E_z [\ln p_\theta(x|z)] - D_{KL}(q_\phi(z|x) \| p_\theta(z)) + D_{KL}(q_\phi(z|x) \| p_\theta(z|x))$$

① estimate through sampling, reconstruction

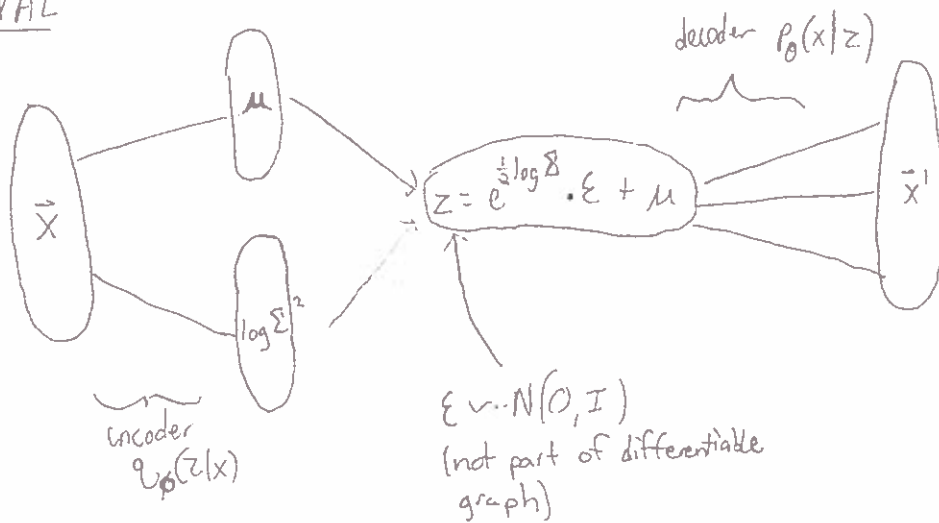
② analytically calculable if we assume $p(z)$ is an easily parametrized dist like $p(z) \sim N(0, I)$

③ posterior intractable as seen earlier, but $D_{KL} \geq 0 !!!$

maximize the lower bound of the evidence

What does this look like in practice? Turn the page

Training a VAE



- Feed \vec{x} to $\vec{\mu}$, $\log \Sigma^2$ layers \rightarrow use to generate some \vec{z} using a random $\epsilon \sim N(0, I)$ that doesn't require differentiation.
- Optimize reconstruction with cross-entropy or MSE
- Optimize $D_{KL}(q_\phi(z|x) || p(z)) \leftarrow$ can do directly with μ, Σ

\hookrightarrow we know $q(z|x)$ goes like $\sim N(\mu, \Sigma)$ from our model, where μ, Σ are functions of x :

$$D_{KL}(q_\phi(z|x) || p(z)) = D_{KL}(N(\mu, \Sigma) || N(0, I))$$

$$= \frac{1}{2} \left(\text{tr}(\Sigma) + \mu^T \mu - k - \log \det(\Sigma) \right) \quad (\Sigma \text{ already diagonal})$$

where k is dimension of Gaussian:

$$= \frac{1}{2} \left(\sum_k \Sigma + \sum_k \mu^2 - \sum_k 1 - \log \prod_k \Sigma \right) \quad \left(\text{treating diagonals of } \Sigma \text{ as a vector} \right)$$

$$= \frac{1}{2} \left(\sum_k \Sigma + \sum_k \mu^2 - \sum_k 1 - \sum_k \log \Sigma \right)$$

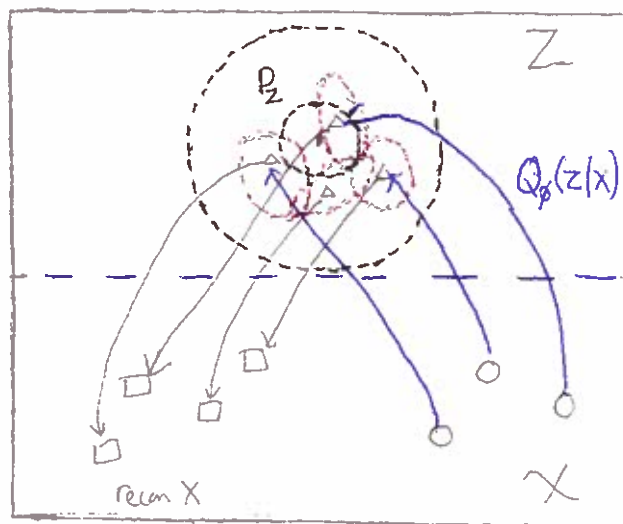
$$= \frac{1}{2} \sum_k \left(\Sigma + \mu^2 - 1 - \log \Sigma \right) \quad \checkmark$$

Overall

- Reconstruction requires that encoded z 's correspond to decoded real samples
 \hookrightarrow enc & dec approximate inverses of each other
- KL latent distribution mapping ensures continuous, navigable latent space, ideally s.t. navigating through latent space corresponds to navigating through space of samples

\nwarrow Wasserstein Autoencoders improve this.

Visualizing VAE Latent Space Correspondence



- When well trained, sample μ, Σ from encoder are very near Gaussian/whatever prior is chosen.

- Including reconstruction error forces them to not all be exact / overlapping

- stable training, but ^{pro} limited resolution because of approx. likelihood max and reconstruction error.

Con

- multiple original x 's run the risk of being most similar the same reconstructed x' b/c of intersection of latent embedd gaussians

↓

Wasserstein Autoencoder improves this by making $q(z|x) = \int q(z|x=x') dP_x$ match $P(z)$, which is a continuous mixture

What else can you do with them?

- improve interpolation in the latent space by learning independent features → minimizing Total Correlation

↓

- can be done with adversary or independence criteria
- implicitly what β -VAE does.

Avoiding explicit likelihood estimation with GAN (Generative Adversarial Networks) Goodfellow 2014

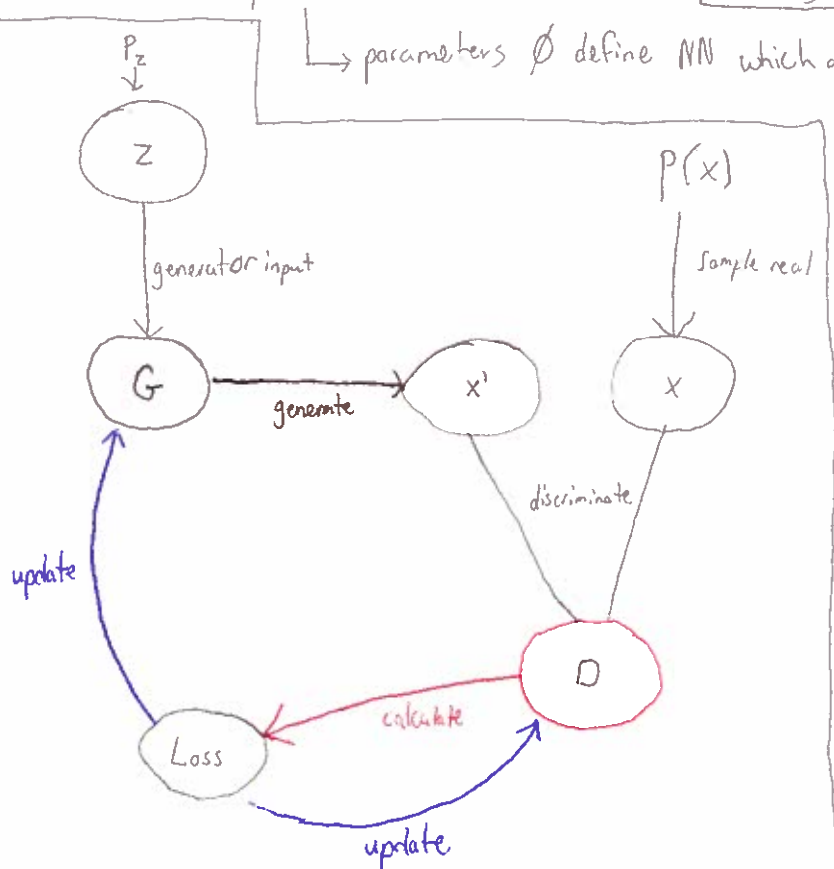
- This time, no explicit estimate of density $p(x)$, which is complex, high-dimensional
- Two networks are jointly optimized in tandem with one another \rightarrow one that learns to generate samples from an approximate distribution $p_\theta(x) \approx p(x)$ and one that tries to discriminate samples from p_θ from those of p . *
- two-player game between generator and discriminator to transform samples from latent $z \sim p(z)$ to $z \rightarrow x \sim p_\theta(x)$

*In reality, we don't have access to the true p , just a set of finite samples from it.

\rightarrow parameters θ define neural network which transforms z in z -space to some x in x -space $\equiv G(z; \theta)$

\rightarrow parameters ϕ define NN which outputs probability that a given x or x' is real $\equiv D(x \text{ or } x'; \phi)$

Model:



The loss function

- Generator and Discriminator participate in min-max game for value function V :

$$\min_G \max_D V(D, G) = E_{x \sim p(x)} [\ln D_\phi(x)] + E_{z \sim p(z)} [\ln (1 - D(G(z)))]$$

Loss function analysis of original GAN framework

- What are optimal D, G ?
- Let's look at case of optimal D .
- Using \int instead of expectation b/c dealing w/ complete distributions

$$\text{So } V(D, G) = \int_x p_{\text{real}}(x) \log D(x) dx + \int_z p_g(z) \log(1 - D(G(z))) dz$$

$$= \int_x \underbrace{p_{\text{real}} \log D(x) + p_g(x) \log(1 - D(x))}_{\text{Optimal discriminator to maximize this is: } D^*(x) = \frac{p_{\text{real}}(x)}{p_{\text{real}}(x) + p_g(x)}} dx \quad \text{where } p_g(x) \text{ is probability that generator produces } x$$

Using this, we rewrite

$$\begin{aligned} V^* &= \min_G V(D^*, G) = E_x[\log D^*(x)] + E_{x \sim p_g}[\log(1 - D^*(x))] \\ &= E_x\left[\log \frac{p_{\text{real}}(x)}{p_{\text{real}}(x) + p_g(x)}\right] + E_{x \sim p_g}\left[\log \frac{p_g}{p_{\text{real}} + p_g}\right] \end{aligned}$$

when $p_g = p_{\text{real}}$, $D^* = \frac{1}{2}$. $p_g = p_{\text{real}}$ is optimal scenario. In this case

$$E_x\left[\log \frac{1}{2}\right] + E_{x \sim p_g}\left[\log \frac{1}{2}\right] = -\log 4 \text{ is lowest value}$$

Factoring this from the RHS of V^* will give us the minimum plus whatever needs to be minimized to 0 to achieve that minimum:

$$\begin{aligned} V^* &= -\log 4 + E_x\left[\log \frac{p_{\text{real}}}{(p_{\text{real}} + p_g)/2}\right] + E_{x \sim p_g}\left[\log \frac{p_g}{(p_{\text{real}} + p_g)/2}\right] \\ &= -\log 4 + D_{\text{KL}}\left(p_{\text{real}} \parallel \frac{p_{\text{real}} + p_g}{2}\right) + D_{\text{KL}}\left(p_g \parallel \frac{p_{\text{real}} + p_g}{2}\right) \\ &= -\log 4 + D_{\text{JS}}(p_{\text{real}} \parallel p_g) \end{aligned}$$

[\therefore Finding optimal value function sol'n is equivalent to minimizing the Jensen-Shannon Divergence.]

GAN perks

- highly parameterized replacement for likelihood means potential for diverse and expressive samples
- Theoretically motivated/sound optimal equilibrium for distribution matching

Original GAN problems

- In practice, hard to minimize this D_{JS}
 - Generator/Discriminator can outperform one another. Training "gets stuck"
- ↳ This is symptomatic of:

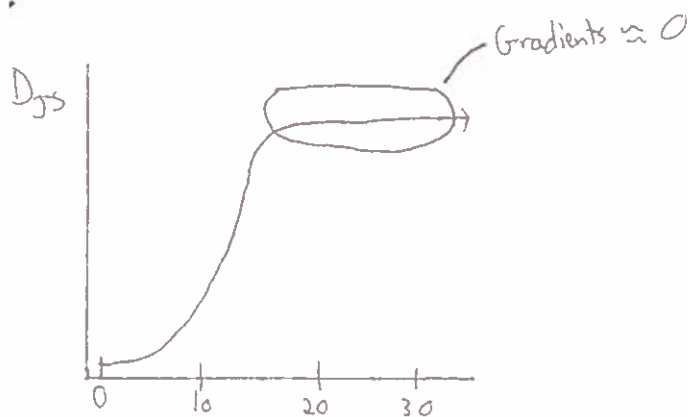
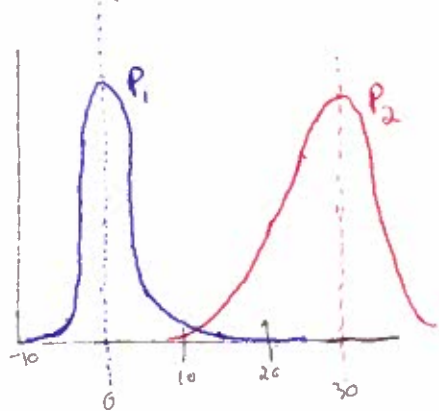
- Overfitting
- Density misspecification
→ no parameters for which p_g and p_{real} are similar enough
- Dimensional misspecification
→ real data might be on a lower dim. manifold

And evidenced by:

- mode collapse
and/or
diminished gradients

- non-convergence:
→ oscillating parameters,
no stability in training

Imagine comparing 2 distributions using D_{JS} :



Two Common Solutions to GAN Problems: Wasserstein GAN and GAN w/ Discriminator Gradient Penalty

Wasserstein GAN: Replace f-divergence with integral probability metric (WGAN-GP) (minimum cost)
(and improvements to it)

Arjovsky 2017
Gulrajani 2017

- question of optimal transport (OT) \rightarrow what is the best measure of how much density of the model distribution needs to be moved to make distributions equal?

Wasserstein Distance: $W(P_r, P_g) = \inf_{\gamma \in \Pi(P_r, P_g)} E_{(x,y) \sim \gamma} [\|x - y\|]$ aka Earth Mover distance

- $\Pi(P_r, P_g)$ is set of all joint distributions $\gamma(x,y)$ whose marginals are P_r, P_g
- infimum is greatest lower bound for any transport plan

* this equation is intractable to estimate, but can be expressed under Lipschitz constraint:

$$W(P_r, P_g) = \sup_{\|f\|_{Lip} \leq 1} E_{x \sim P_r} [f(x)] - E_{x \sim P_g} [f(x)]$$

- 1-Lipschitz functions are just functions that slopes/gradients ≤ 1 for all points
- This simplification is derived via the Kantorovich-Rubinstein duality
- New objective: The discriminator learns this supremal function $f(x)$ and the distance is minimized
- Discriminator no longer a discriminator \rightarrow now a critic
 - used to calculate a metric of distributional distance

$$\frac{|f(x) - f(y)|}{\|x - y\|} \leq 1$$

1 Problem

- hard to maintain Lipschitz constraint when learning $f(x)$
- can penalize model if gradient becomes > 1 :

$$L = E_{\tilde{x} \sim P_g} [D(\tilde{x})] - E_{x \sim P_r} [D(x)] + \lambda E_{\hat{x} \sim P_{\hat{x}}} [\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1]$$

where $\hat{x} = t\tilde{x} + (1-t)x$ w/ t uniformly sampled b/w 0 and 1

GAN w/ Discriminator Gradient Penalty (GAN-DP) Roth 2017

- similar gradient penalty to improving WGAN can be used to improve original GAN training.

- this method is locally convergent in general case, whereas WGAN is not

- $$\mathcal{L} = E_x [\ln D_\theta(x)] + E_z [\ln(1 - D(G(z)))] + \underbrace{\frac{\gamma}{2} E_x [\|\nabla D_\theta(x)\|^2]}_{\text{regularizer that prevents Discriminator from deviating from Nash Equilibrium}}$$

- Stability studies on next page. WGAN-GP only stable in this context at very low learning rate $lr = 10^{-5}$

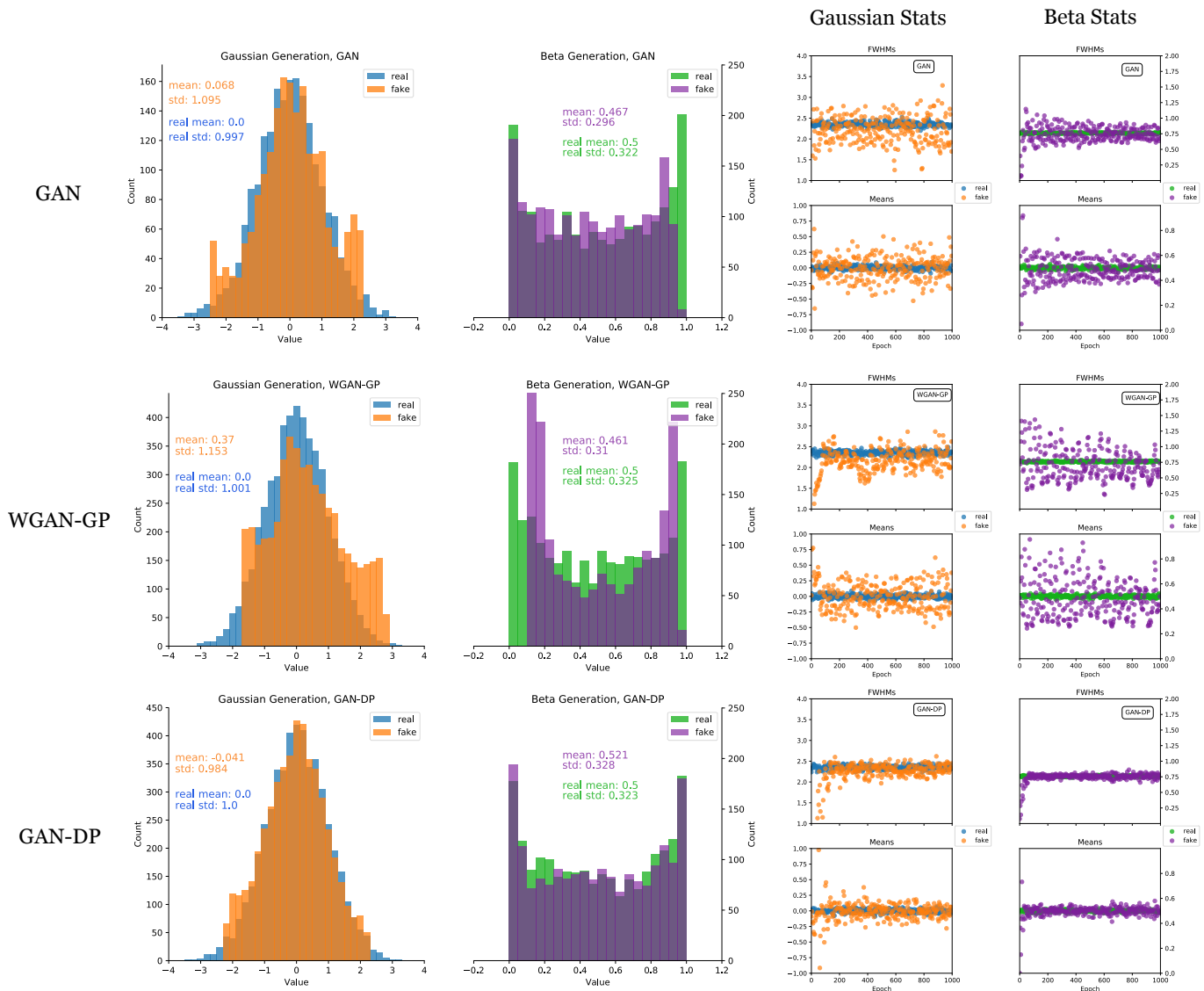


Figure 2.3: Top Row: Generation of Gaussian and Beta distributions using original GAN loss. Middle Row: Generation of Gaussian and Beta distributions using WGAN-GP loss. Bottom Row: Generation of Gaussian and Beta distributions using GAN-DP loss. The right half contains values of the sample mean and standard widths of the two generated distributions across epochs compared to ground truth.

accuracy, both due to my selecting of sample distributions to display and due to the stochasticity of sample generation. Yet, the continuous comparison of some metrics that characterize the Gaussian and Beta distributions across training epochs allows for more formative evaluation between the models. The WGAN-GP and unregularized GAN showed evidence of deviating from expected characteristics even after approximating them better in previous epochs, while the GAN-DP does not show such inconsistencies. Upon closer inspection of the WGAN-GP results, in

which the training in the Beta case seems to be converging but at a much slower rate, another test was done with minute learning rate of 10^{-5} which was outside of the original hyperparameter search. Results show that the WGAN-GP could move down the loss surface more stably if trained with very small gradient updates, as shown in Figure 2.4.³ Moreover, a comparison of the KL-divergences between estimates of the probability distributions from the generated samples and the true samples corroborates these claims. The average KL-divergences over 100 samples of 10000 values each were lowest for the GAN-DP and higher for the unregularized GAN and WGAN-GP models on both the Gaussian and Beta tests, as shown in Table 2.2. The WGAN-GP value improved significantly to 0.024 when training was locally convergent in the slow learning rate case. The KL-divergence was computed by creating normalized histograms of the 10000 values so that there could be a density estimates $p_{GAN}(\mathbf{x})$ and $q_{true}(\mathbf{x})$ at each bin. A Kernel Density Estimation was also tested to compare to the normalized histogram binning technique and it yielded similar results.

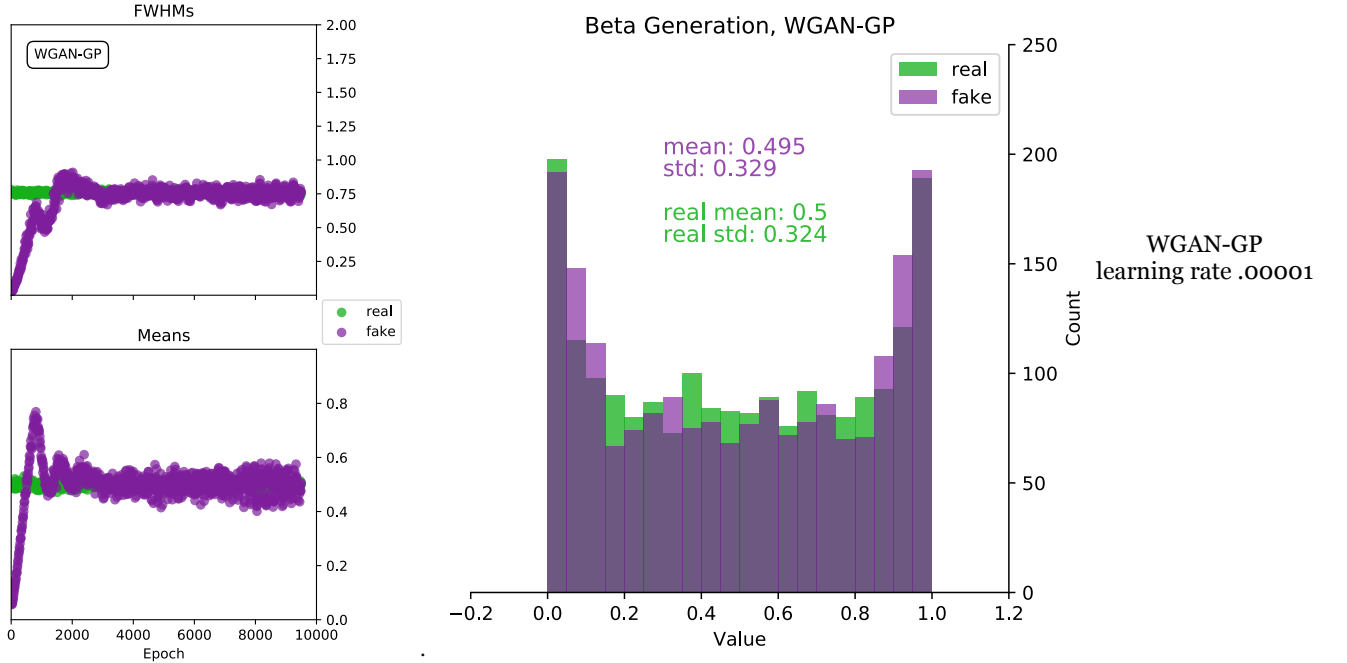


Figure 2.4: WGAN-GP training only achieves accuracy with small learning rate and many epochs.

Overall, the GAN-DP method could accurately model both the Gaussian and Beta distributions, capturing natural parameters of the distributions. It empirically struggles to generate values

³See Appendix B.1 for comparison of local stability across learning rates.