

## Autoencoders

**Comparison to PCA:** Autoencoders have an encodec scheme where encoder  $\rightarrow$  lower latent, decoder  $\rightarrow$  recon input. Like PCA where lower dim projection (principal coordinates) by using eigenvectors to capture the directions of most variance, recon the original input through a linear combo of the principal coordinates and their respective principal coordinates. Deep Autoencoders = pseudo-invertible nonlinear dimensionality reduction

**What distribution should the latent follow?** We can enforce a prior distribution by adding a probabilistic distance error (e.g. KLD or JSD). **Kullback-Leibler (KL) Divergence:**

$D_{KL}(p||q) = \mathbb{E}_x \left[ \log \frac{p(x)}{q(x)} \right]$  But this is not defined in non-overlapping support, i.e. when  $q(x) = 0$  Not symmetric. The Jensen-Shannon (JSD) Divergence fixes this by essentially repeating the KL divergence twice to make it symmetric. Formula:

$$D_{JSD}(p||q) = \mathbb{E}_x \left[ -\log \frac{p(x)}{q(x)} \right] + \mathbb{E}_x \left[ -\log \frac{q(x)}{p(x)} \right]$$

How do we calculate a Divergence metric when we don't have access to the data distribution? We're only given access to a batch of samples at training time.

## Adversarial Networks

One way we can enforce the latent distribution to follow the prior distribution is through adversarial networks by using a discriminator (adversary) to judge whether a generated output is real or fake. This training method actually allows us to provably measure JSD at an individual sample level.

Loss:  $\min_{\theta_g, \theta_\phi} \max_{\eta} \mathbb{E}_{x \sim p_X} [\|x - \phi(\theta(x))\|^2] + \lambda \log(1 - \eta(\phi(x))) + \mathbb{E}_{z \sim q_z} [\log(\eta(z))]$  where first term is MSE and other terms are JSD.

Proof deriving JSD:

$$\begin{aligned} \int_z q_z(z) \log(\eta(z)) dz + \int_z p_z(z) \log(1 - \eta(z)) dz &\xrightarrow{L} \frac{\partial L}{\partial \eta} = \frac{q_z(z)}{\eta(z)} - \frac{p_z(z)}{1 - \eta(z)} dz = \\ 0 \Rightarrow \eta^* &= \frac{q_z}{q_z + p_z}, \int_z q_z(z) \log\left(\frac{q_z}{q_z + p_z}\right) dz + \\ \int_z p_z(z) \log\left(1 - \frac{q_z}{q_z + p_z}\right) dz &= \\ \int_z q_z(z) \log\left(\frac{q_z}{q_z + p_z}\right) dz + \int_z p_z(z) \log\left(\frac{p_z}{q_z + p_z}\right) dz &= \\ \int_z q_z(z) \log\left(\frac{q_z}{q_z + p_z}\right) dz + \int_z p_z(z) \log\left(\frac{p_z}{q_z + p_z}\right) dz - &2 \log(2) = 2JS(q_z, p_z) - 2 \log(2) \end{aligned}$$

Extra term  $-2 \log(2)$  acts as a lower bound on the loss if  $p(z) = q(z)$  so if discriminator gives  $p = 0.5$ , meaning that  $\eta(z) = 0.5$  from the initial distance formula, then we get  $-\log(2) + -\log(2) = -2 \log(2)$ .

**Mode Collapse:** However, we may get mode collapse when the generator starts producing a limited variety of outputs, failing to capture the full diversity of the target distribution.

**Wasserstein GAN** Also known as Earth Mover's distance. Quantifies the amount of "work" moving a probability distribution to another. This better quantifies probabilistic distance than other metric, but I won't get into the details here. For the dual formulation for 1-Wasserstein distance:  $W_1(p, q) = \max_{\|\eta\|_L \leq 1} (\int_X \eta(x)p(x)dx - \int_Y \eta(y)q(y)dy)$

$X, Y \subseteq \mathbb{R}^d$  For the Wasserstein Adversarial Net-

works:  $D(p_z, q_z) = \max_{\|\eta\|_L \leq 1} \mathbb{E}_{z \sim q_z} [\ln(\eta(z))] - \mathbb{E}_{z' \sim p_z} [\ln(\eta(z'))]$  For the Loss function:  $\min_{\theta_g, \theta_\phi} \max_{\eta} \mathbb{E}_{x \sim p_X} [\|x - \psi(\phi(x))\|^2] + \lambda (2\eta(\phi(x)) - \mathbb{E}_{z \sim q_z} [\ln(\eta(z))])$

## Variational Autoencoders

The prior  $p(z)$  here is a simple Gaussian, which is mathematically convenient to work with, while the conditional output  $p(x|z)$  is complex (image generated). Thus, we would ideally use MLE to estimate the parameters:  $p_\theta(x) = \int p_\theta(x)p_\theta(z)dz$  However, marginalizing across all  $z$  is simply not practically possible, which makes this MLE intractable. Thus, the posterior density  $p_\theta(z|x)$  is also intractable.

**Solution:** we also need to estimate an encoder  $q_\phi(z|x)$  model. We can derive the loss terms based on the Evidence-based Lower Bound (ELBO) that we want to maximize:  $\log p_\theta(x^{(i)}) = \mathbb{E}_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)})]$  (Does

$$\begin{aligned} \text{not depend on } z) &= \mathbb{E}_z \left[ \log \frac{p_\theta(x^{(i)})}{p_\theta(z)} p_\theta(z|x^{(i)}) \right] \\ (\text{Bayes' Rule}) &= \mathbb{E}_z \left[ \log p_\theta(x^{(i)}|z) p_\theta(z) \frac{q_\phi(z|x^{(i)})}{q_\phi(z|x^{(i)})} \right] \\ (\text{Multiply by constant}) &= \mathbb{E}_z [\log p_\theta(x^{(i)}|z)] \\ - \mathbb{E}_z \left[ \log \frac{q_\phi(z|x^{(i)})}{p_\theta(z|x^{(i)})} \right] &+ \mathbb{E}_z \left[ \log \frac{q_\phi(z|x^{(i)})}{p_\theta(z|x^{(i)})} \right] \\ (\text{Logarithms}) &= \mathbb{E}_z [\log p_\theta(x^{(i)}|z)] - \end{aligned}$$

$D_{KL}(q_\phi(z|x^{(i)})||p_\theta(z))$  + positive KL term (intractable, ignore). The first 2 terms are the ELBO, term 1 is decoder distribution estimated through sampling w/ reparam trick (MSE), 2nd term is KLD loss.

The first term is the reconstruction loss and the second term is to make the approximate posterior close to the Gaussian prior.

**Reparam Trick:** To take a forward pass through the VAE, we have to sample from the latent distribution (Gaussian in this case). However, we can't backprop through such a stochastic operation, so the reparameterization trick is used. Since we estimate  $\mu$  and  $\sigma$  through the encoder model, instead of drawing  $z$  from  $N(\mu, \sigma^2)$ , we can instead draw  $\epsilon$  from  $N(0, 1)$  and then get  $z = \mu + \sigma * \epsilon$ , and thus gradient calculation isn't affected.

**Generating Data:** Once the VAE trained, sample data from Gaussian as the latent  $z$  and use the decoder model for output image. Since  $z$  is a multivariate Gaussian, the diagonal prior on  $z$  are independent latent variables that cause different factors of variation (think principal comp in PCA).

**Denoising Autoencoders:** It seems that  $\psi(\phi(x)) - x$  would give us a vector that projects data from outside the manifold onto the manifold: Through the Bengio paper "What Regularized Auto-Encoders Learn from the Data-Generating Distribution", they find that:  $\nabla_x \log(p(x)) = -\nabla E(x)$  In other words, the score function is the gradient of the log probability of the data! If we have the score function, we can sample from the distribution by setting an initial  $x$  drawn from an arbitrary prior distribution and then following the gradient of the data to eventually converge at a sample from  $p(x)$ . This is called Langevin Dynamics:  $x_{i+1} \leftarrow x_i + \epsilon \nabla \log p(x) + \sqrt{2\epsilon} z_i, i \in 0 \dots K$  where the "learning rate"  $\epsilon$  is sufficiently small and

number of steps  $K$  is sufficiently large.

## Diffusion Models

For the training algorithm: Repeat  $x_0 \sim q(x_0)$   $t \sim \text{uniform}\{1, \dots, T\}$   $\epsilon \sim \mathcal{N}(0, I)$   $x_t = \sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}\epsilon$   $\theta \leftarrow \theta + \eta \nabla_\theta \|\epsilon - \epsilon_\theta(x_t, t)\|^2$  Until convergence For the testing algorithm:  $x_T \sim \mathcal{N}(x_0, I)$  For  $t = T, \dots, 1$  do  $\epsilon \sim \mathcal{N}(0, I)$  if  $t > 1$ , else  $\epsilon = 0$   $x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon_\theta(x_t, t) \right) + \sqrt{\alpha_t} \epsilon$  Return  $x_0$

## Adversarial Attacks

White-Box: architecture and weights are known Black-Box: only access to input and output (e.g. API) Gray-Box: architecture known but not weights

We can attain an adversarial generally by adding some perturbation that maximizes the loss:  $\max_{\delta} \text{loss}(\theta, x + \delta, y)$  We want  $\delta$  that's small w.r.t.  $l_p$  norm, Rotation/Translation, VGG feature perturbation or any other perturbation

## White Box Attack Methods

### Fast Gradient Sign Method (FGSM) Attack:

Use pretrained classifier like ResNet50:  $\hat{y} = f(\theta, x)$  Find adversarial example that maximizes the loss:  $\mathcal{L}(x', y) = \mathcal{L}(f(\theta, x'), y)$  Bounded perturbation s.t.:  $\|x' - x\|_\infty \leq \epsilon$ , where  $\epsilon$  is attack strength Optimal adversarial image:  $x' = x + \epsilon \text{sign}(\nabla_x \mathcal{L}(x, y))$

**Iterative FGSM Attack:** Let  $m$  be number of iterations  $x^{(m)} = x^{(m-1)} + \epsilon \text{sign}(\nabla_x \mathcal{L}(x^{(m-1)}, y))$  Both (I)FGSM are fix-perturbation attacks

**(Iterative) Least Likely Attack:** Similar to FGSM but  $y_{LL}$  is the least likely (LL) class predicted by the network on clean image  $x$   $x' = x - \epsilon \text{sign}(\nabla_x \mathcal{L}(x, y_{LL}))$  Strong attack as it emphasizes the least likely class

**Projected Gradient Descent Attack:** We can take a gradient step and then project it back to the feasible set  $\Delta$  since the perturbed input may not lie on the data manifold (similar to denoising AutoEncoder logic)  $\delta := \mathcal{P}_\Delta[\delta + \nabla_\delta \text{Loss}(x + \delta, y; \theta)]$  For example, the projected gradient descent applied to  $l_\infty$  ball, repeat:  $\delta := \text{Clip}_\epsilon[\delta + \alpha \nabla_\delta J(\delta)]$  Slower than FGSM but typically able to find better optima

**Carlini and Wagner (CW) L2 Attack:** zero-confidence attack for all  $t \neq y$ , find the adversarial image that will be classified by  $t$  as solving the problem:  $\min_\delta \|\delta\|_2^2$  subject to  $f(x + \delta) = y, x + \delta \in [0, 1]^n$  Finding the exact solution is difficult so we use relaxed version  $\min_\delta \|\delta\|_2^2 + c \cdot g(x + \delta)$  subject to  $x + \delta \in [0, 1]^n, c \geq 0$  Let  $Z(x)$  be the NN activations before the logit output layer, also called the embeddings  $g(x) = \max(\max_{i \neq t}(Z(x)_i - Z(x)_t), 0)$  Let  $\delta = \frac{1}{2}(\tanh(w) + 1) - x$  With the following constrained optimization problem:  $\min_w \|\frac{1}{2}(\tanh(w) + 1) - x\|_2^2 + c, \text{ReLU}\{\max_{i \neq t} Z(\frac{1}{2}(\tanh(w) + 1))_i - Z(\frac{1}{2}(\tanh(w) + 1))_t\}$  Powerful attack method that resists many defenses

**Universal Adversarial Perturbation Attacks:** A single perturbation on any image with high probability (e.g. 0.8+) Generalize well across different models

**Single Pixel Attack:** self-explanatory. Only modify one pixel

**Poisoning Attacks:** manipulating the training data itself rather than during inference Maintain accuracy but hamper generalization due to outliers through poisoning