

Chaitanya Kohli (S25090)

Q Derive KL Div of 2 gaussian distributions.

$$KL(P||Q) = \int P(x) \log \frac{P(x)}{Q(x)} dx.$$

$$P(x) \sim N(\mu_1, \Sigma_1)$$

$$Q(x) \sim N(\mu_2, \Sigma_2)$$

$$\int N(\mu_1, \Sigma_1) \log \left(\frac{N(\mu_1, \Sigma_1)}{N(\mu_1, \Sigma_2)} \right)$$

$$\int N(\mu_1, \Sigma_1) \log \left(\frac{\cancel{(2\pi)^{d/2}} |\Sigma_2|^{1/2} \exp \left\{ - \left(\frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right) \right\}}{\cancel{(2\pi)^{d/2}} |\Sigma_1|^{1/2} \exp \left\{ - \left(\frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \right) \right\}} \right)$$

$$\int N(\mu_1, \Sigma_1) \left\{ \log \left(\frac{|\Sigma_2|^{1/2}}{|\Sigma_1|^{1/2}} \right) + \left[\begin{array}{l} -\frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \\ + \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \end{array} \right] \right\} dx$$

$$E_{P \sim N(\mu_1, \Sigma_1)} \left\{ \frac{1}{2} \log \left(\frac{|\Sigma_2|}{|\Sigma_1|} \right) + \left[\begin{array}{l} -\frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \\ + \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \end{array} \right] \right\}$$

$$x^T A x = \text{Tr}(x^T A x) = \text{Tr}(A x x^T)$$

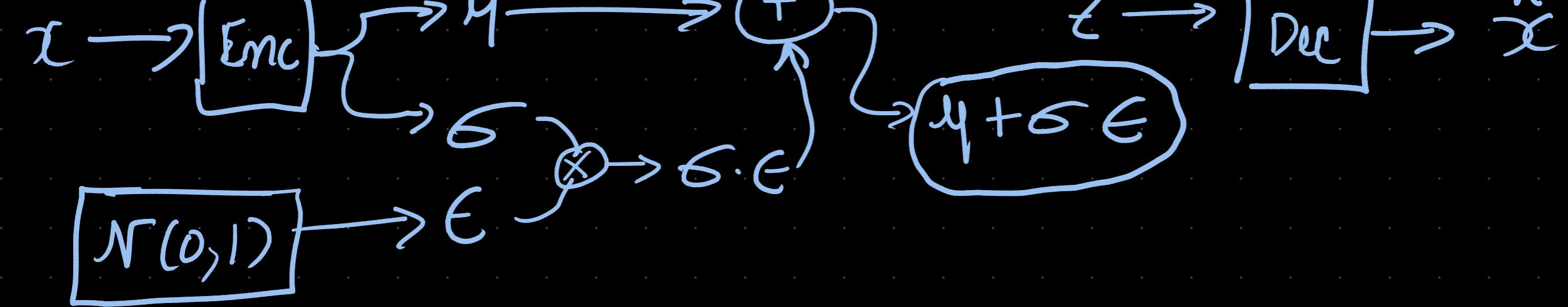
$$E_{P \sim N(\mu_1, \Sigma_1)} \left\{ \frac{1}{2} \log \left(\frac{|\Sigma_2|}{|\Sigma_1|} \right) + \left[\begin{array}{l} -\frac{1}{2} \text{Tr}(\Sigma_1^{-1} (x - \mu_1)^T (x - \mu_1)) \\ + \frac{1}{2} (x - \mu_1 + (\mu_1 - \mu_2))^T \Sigma_2^{-1} (x - \mu_1 + \mu_1 - \mu_2) \end{array} \right] \right\}$$

$$\begin{aligned}
& E_p \left[\frac{1}{2} \log \left(\frac{|\varepsilon_1|}{|\varepsilon_2|} \right) \right] - \frac{1}{2} E_p \left[\text{tr}(\varepsilon_1^{-1} (x - \mu_1)^T (x - \mu_1)) \right] \\
& + \frac{1}{2} E_p \left[(x - \mu_1)^T \varepsilon_2^{-1} (x - \mu_1) + 2(x - \mu_1)^T \varepsilon_2^{-1} (\mu_1 - \mu_2) \right. \\
& \quad \left. + (\mu_1 - \mu_2)^T \varepsilon_2^{-1} (\mu_1 - \mu_2) \right] \\
& E_p \left[\frac{1}{2} \log \left(\frac{|\varepsilon_1|}{|\varepsilon_2|} \right) \right] - \frac{1}{2} \text{tr} \left(\varepsilon_1^{-1} E_p (x - \mu_1)^T (x - \mu_1) \right) \\
& + \frac{1}{2} \text{tr} \left[\varepsilon_2^{-1} E_p (x - \mu_1)^T (x - \mu_1) \right] + 2 \text{tr} \left(\varepsilon_2^{-1} E_p (x - \mu_1)^T (\mu_1 - \mu_2) \right) \\
& + \frac{1}{2} E_p \left[(\mu_1 - \mu_2)^T \varepsilon_2^{-1} (\mu_1 - \mu_2) \right] \\
& = \frac{1}{2} \log \left(\frac{|\varepsilon_1|}{|\varepsilon_2|} \right) - \frac{1}{2} \text{tr}(\varepsilon_1^{-1} \varepsilon_1) + \frac{1}{2} \text{tr}(\varepsilon_2^{-1} \varepsilon_1) \\
& + \frac{1}{2} (\mu_1 - \mu_2)^T \varepsilon_2^{-1} (\mu_1 - \mu_2)
\end{aligned}$$

Q Why is Sampling not differentiable

Sampling involves bit wise AND, OR etc operations to keep the sampling cryptographically secured. The bit wise operations are non-differentiable making Sampling a non-differentiable operation.

Q How reparameterization helps in this case?



$$z = y + \sigma \epsilon$$

$$\frac{\partial z}{\partial y} = 1, \quad \frac{\partial z}{\partial \sigma} = \epsilon$$

Q Spectral Normalization

$$w = \frac{w}{\sigma(w)}$$

$\sigma(w) \rightarrow$ Largest singular value

It represents maximum amount that the matrix W can stretch a vector.

Spectral Normalization helps enforce Lipschitz condition

$$\|f(x) - f(y)\| \leq K \|x - y\|$$

This ensures gradient stability.

Q Explain

a) Deconvolution layer

It is used to upsample the feature maps to higher spatial dimension.

Also known as Transposed Convolution.

Used in VAE, AE, U-nets, in Generator of GAN etc.

i) Strided Convolutional Layer:

Used for downsampling & reducing spatial dimension.

ii) Fractional Convolutional Layers

Another name for Deconv Layer / Transposed Convolutional layer.

QPGR GAN

Progressive Growing GANs for improved stability, and variation.

The paper generated the first high resolution (1024×1024) image of human faces, solving instability seen in previous GAN.

i) Progressive Growing

Initially small Generator & Discriminator are used to train (4×4). Then add layers progressively

Doubling the output. ($4 \rightarrow 8 \rightarrow 16 \rightarrow \dots \rightarrow 1024$)

It's a fully autoencoder based GAN (The autoencoder part is not shown)

This enables network to learn coarse structures, then finer details as output grows. This makes training stable & faster.

2) Fade-in mechanism:

Snapping on an entire layer which increases spatial dimension shocks the pre-trained model & ruins the parameters learned.

- parameter α is used to fade-in the new layers.
It begins from $\alpha=0$, completely ignore the new layers outputs.
- linearly increases over iterations and NN outputs a weighted average of old lower resolution pathways & “new” resolution output.

3) Mode Collapse. (Generator produces only few variations of same image)

PGGAN authors introduced a layer at end of the Dib. that computes standard deviation of features across the minibatch.

Mode collapse occurs if generator realizes it can easily cheat Discriminator with only one image.

calculate std. deviation of each feature map of batch
take the average of std. deviation across this batch.
And append a feature map of avg. std deviation to
the original feature map.

If std. deviation is zero the Discriminator can easily
identify fake batches. &

If std. deviation is high the Discriminator treats it as
a normal feature of real data.

4) Pixel wise Normalization.

Batch Normalization causes GRANS to be unstable.

Pixel Norm \rightarrow It normalizes the feature vector in each
pixel location to unit length after entry conv layer.

$$\theta_{x,y} = \alpha_{x,y} / \sqrt{\frac{1}{N} \sum \alpha_{xy}^2 + \epsilon}$$

This prevents signal escalation without destroying
the semantic info in magnitude.