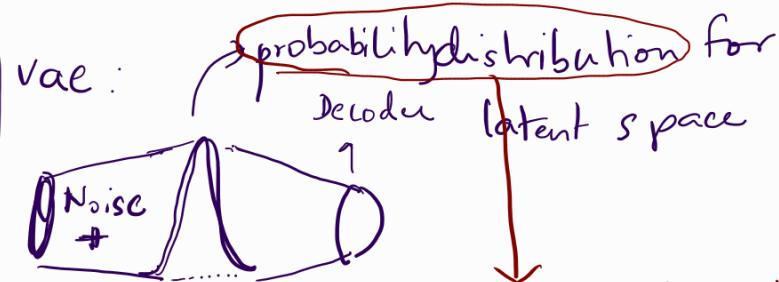
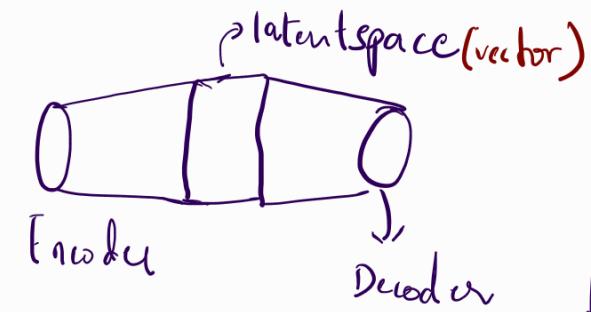


Normal Autoencoder



main goal is to learn these distributions that can approximate the data.

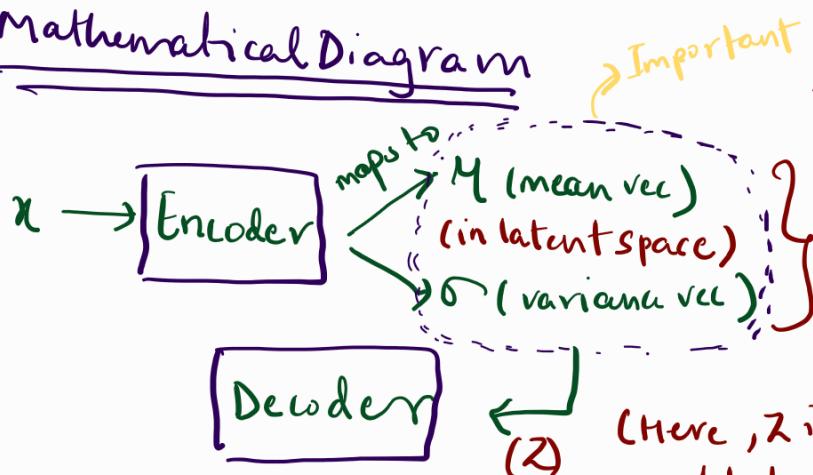
Mathematical Objective : Maximum Likelihood

intuition: we want our data closely related to the actuals. Hence, if our distribution is close to what it should be \rightarrow eureka

How do we achieve : By learning these params from ^{our} distributions.

Q. Why distribution approach is better than direct mapping into latent space and the providing output

Mathematical Diagram



Together define the Gaussian Distribution of each data point.

(Here, z is picked from the distribution of latent space, as mentioned above)

How do we get the best sampling ?

? From deterministic to stochastic
From point value to distribution $\frac{1}{N(\mu, \sigma)}$

Why cannot we just use the mean to sample z ? (variance ?)

Answer: Random variation, rather than randomness causing deterministic nature to a certain degree.

Intuition: sample can have slight differences while still capturing the core features of original data.

Objective Function : (To estimate the a good distribution)

Reconstruction loss = $\dots + \dots$ KL divergence

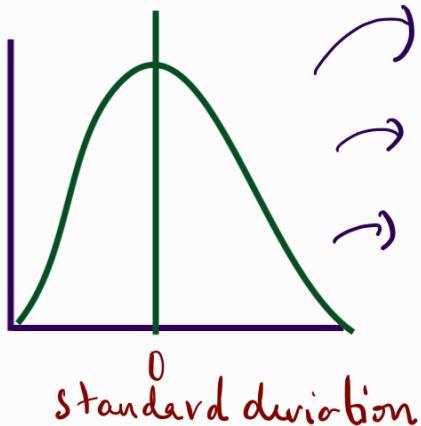
Motive: Decoded output is similar to actual data.

Motive: Makes the learned latent space distribution resemble a standard normal distribution $\rightarrow N(0, 1)$ (as closely as possible).

$$\mathcal{L}_{VAE} = \text{reconstruction loss} + \text{KL Divergence}(N(0, 1))$$

Q) Why do we need the latent space distribution to approximate a standard normal distribution?

Standard normal Distribution



→ Symmetric about mean

→ $\mu = 0, \sigma = 1$

→ (All normal distributions are symmetrical but vice versa \rightarrow not true)

Note 1

(ND) : Normal distribution is key to central limit theorem.

Note 2 : One has to study ND to interpret KL.

Intuition: we introduce stability and regularity into space.

\Rightarrow while sampling points, one is less likely to encounter extreme or erratic values } Hence, coherent outputs.

⇒ Regularization proposed by KL, provides continuity in latent space, making it easier to transit from one feature to another

Why KL divergence? → measures difference between two probability distributions

in our case $\xrightarrow{\textcircled{1}}$ Distribution learned by encoder
 $\xrightarrow{\textcircled{2}} N(0, 1)$ (specific to each data point)

→ Learned distribution for z given x

$$D_{KL}(q(z|x) \parallel N(0,1)) \quad \left\{ \begin{array}{l} \text{intuition: to keep} \\ A \text{ close to } B \end{array} \right.$$

↓
A ↓
B

\Rightarrow So that the learned representation stay close to actual data.

Note: If the distribution stays close to $N(0, 1)$ no large shifts in latent space will be learnt

The Reparameterization trick: (Why, What, How)

Since the beginning of learning process we have been generating and approximating distribution

we need gradients to back propagate | Due to random sampling

Relu is non-smooth gradients
→ continuous functions

makes it impossible

for the gradients to flow through encoder → Model doesn't learn properly

The trick: Reparameterize sampling step } This will make it differentiable

⑧ Instead of sampling z directly from $N(0, 1)$

we break it into 2 parts.

First, sample a standard normal \rightarrow random variable ϵ from $N(0, 1)$

Then shift and scale ϵ using μ and σ

$$z = \mu + \sigma \cdot \epsilon$$

• Our latent vector would look like

since, the randomness parameter is separated from μ and σ , the gradients can flow through μ and σ properly.

Generates randomness but is independent of μ, σ

Why would the above trick work?

→ Since μ and σ are now deterministically used, we have a differentiable path for the latent variable z back to the encoder's parameters.

→ Providing stochastic nature (ϵ) to VAE and also differentiable at the same time.

Reparameterization \longrightarrow back propagation.

Training objective of the VAE: $\{$ Reconstruction loss
+
KL divergence.

(A)

(B)

$$L_{VAE} = E_{q(z|x)} [\log p(x|z)] - D_{KL}(q(z|x) || p(z))$$

Q) Why might balancing these two terms be essential for training an effective VAE? (output sensible)

(A)

Accurate reconstruction from latent space. Sufficient reconstruction accuracy. Capturing key features.

(B)

Prevents latent space from overfitting. Helping the system to not copy and be realistic at the same time.

The balance allows realistic reconstructions while maintaining generalizable and structured latent space. Hence, allowing VAE to generate coherent new samples and interpolate smoothly between different points in latent space.

Until now, we have built these amazing latent spaces that provide the representational learning in vector format. But, does one really know what's within that latent space, what core feature representations has the model learnt \rightarrow Ans \rightarrow INTERPRETABILITY

* Visualizing Latent Space Interpretability.

- Latent space traversals → systematically altering values in the latent space to understand how changes in latent variables affect the output.
compressed, lower space dimensions
(Interpolating) between data points
Eg: $z \rightarrow z + t(z)$

Examples : ① stroke thickness in MNIST ② presence of smile in CelebA

Note: The following traversals involve varying one specific dimension, keeping the remaining constant.

Mathematical Intuition (Important)

Latent space z is modeled as a probability distribution $N(\mu, \sigma)$

Disentangled representation
a way of encoding data such that each dimension in the latent space corresponds to a distinct, interpretable factor of variation of data

During Latent Space Traversal

- Fix all dimensions except one (Eg: z_i)
- Linearly vary z_i across range of values (Eg: $-3, 3$)
- Pass the modified z through decoder $p(x|z)$ to generate

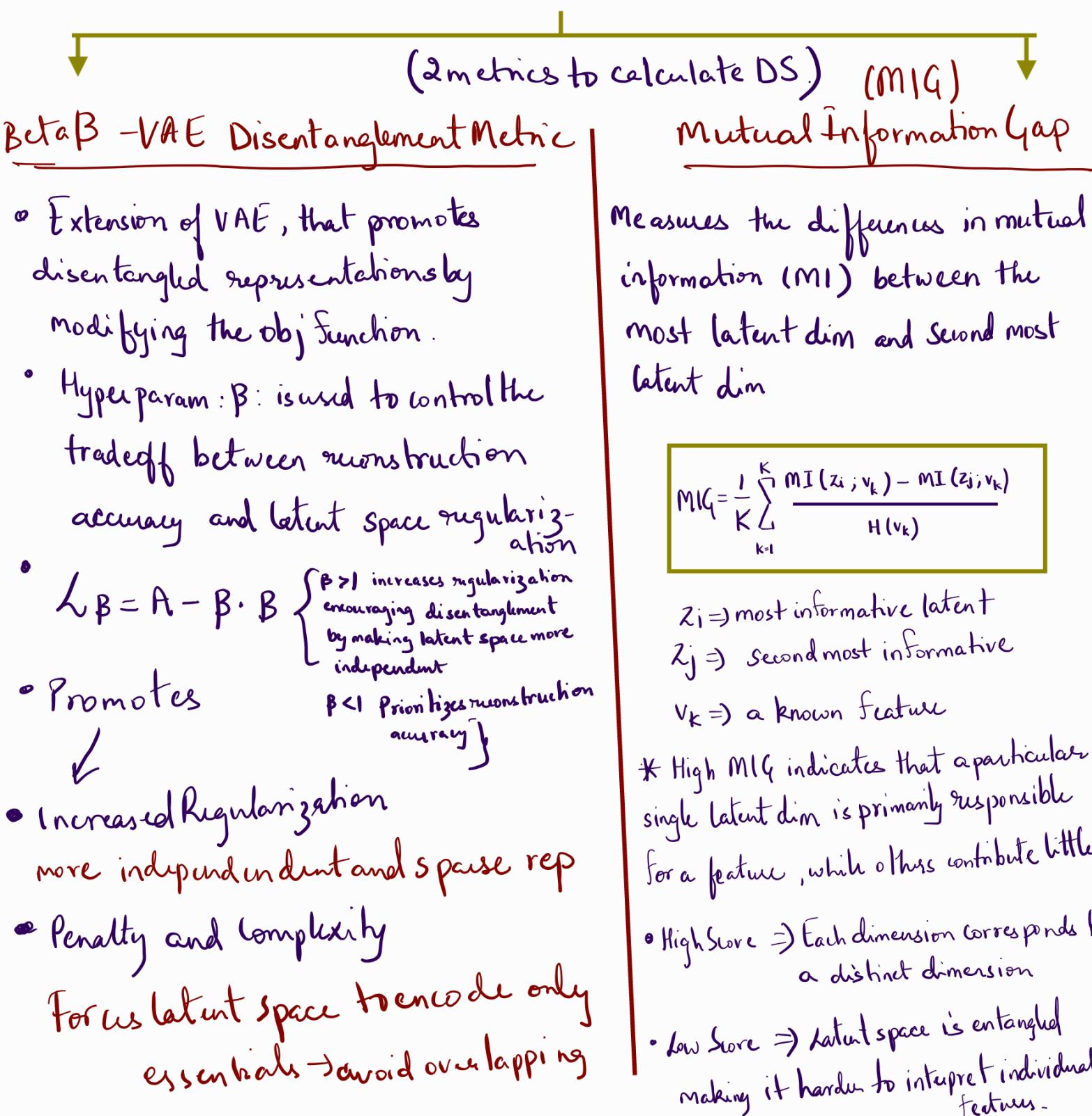
$$z = \mu + \sigma \cdot \epsilon$$

- Why might the ability to vary only one dimension in the latent space be useful for understanding and controlling the behavior of generative model?

By altering a single dimension in latent space we get a clear view of how specific features are represented and controlled. This helps us understand the structure of the latent space and whether the VAE has learned disentangled representations where different features are encoded independently.

- **Disentanglement Score**: The DS corresponds to how compact and meaningful compression (DS) why do we measure this score?
Well each latent dimension corresponds to a single interpretable feature. (altering one dimension shouldn't affect the other features)
Thus shouldn't overlap ↗

Eg: Dim1 : controls presence of smile
Dim2 : controls presence of glasses.



Robustness and Generalization: examining how the latent space and reconstructed outputs behave when the input data is perturbed {intentionally modifying or } altering

→ Helps us understand ① Robustness ② Generalization.

- Robustness: stable the learned latent representations are when small changes like noise or adversarial perturbations are added to input.

- Generalization: meaning full outputs by VAE.

Mathematical intuition:

→ ① Perturbation Analysis

→ ② Reconstruction stability

① Add Gaussian noise η to the input data x : $\tilde{x} = x + \eta$, where $\eta \sim N(0, \sigma^2)$

• Encode both x, \tilde{x} to obtain $\rightarrow z, \tilde{z} \Rightarrow \text{MSE} = \frac{\|z - \tilde{z}\|^2}{\text{stability}}$

② Decode z and \tilde{z} to obtain $\rightarrow \hat{x}, \tilde{\hat{x}} \Rightarrow \text{MSE} = \frac{\|\hat{x} - \tilde{\hat{x}}\|}{\text{Reconstruction Stability}}$

Generative Quality Evaluation with metrics: Quality, Diversity

• Frobenius Inception distance: compares statistical properties like mean & covariance of features from real & generated extracted from a pretrained model (Eg Inception V3)

$$\text{FID} = \|M_r - M_g\|^2 + \text{Tr}(\Sigma_r + \Sigma_g - 2\sqrt{\Sigma_r \Sigma_g})$$

$M_r, \Sigma_r, M_g, \Sigma_g$: mean and covariance of real & generated outputs.

Inception Score: Measures diversity and quality of generated images based on the predicted class probabilities using pretrained classifiers.

$$IS = \exp(E_{x \sim p_g}[D_{KL}(p(y|x) || p(y))])$$

$p(y|x) \Rightarrow$ predicted label distribution
 $p(y)$ marginal label distribution

SSIM : structural Similarity Index Measure.

↳ Consider → Luminance, contrast, structure as an quality assessment. → used for smoothness measuring continuous meaningful transitions.

INTRODUCTION TO CHAOS THEORY AND ITS IMPLEMENTATION IN VAE

Chaotic Systems : → deterministic systems that exhibit seemingly random behavior

Key characteristics : ① Sensitive dependence to initial conditions
② Deterministic Dynamics ③ Fractal and self similarity.