

Resources:

**Demo:** [http://dpkingma.com/sgvb\\_mnist\\_demo/demo.html](http://dpkingma.com/sgvb_mnist_demo/demo.html)

<https://www.youtube.com/watch?v=9zKuYvjFFS8&t=26s>

<https://www.youtube.com/watch?v=vy8q-WnHa9A>

<https://www.youtube.com/watch?v=nKM9875PVtU&t=8s>

[https://www.youtube.com/watch?v=YgSWrafXI8U&list=PLTKMiZHVd\\_2KJtIXOW0zFhFfBaJJilH51&index=142](https://www.youtube.com/watch?v=YgSWrafXI8U&list=PLTKMiZHVd_2KJtIXOW0zFhFfBaJJilH51&index=142)

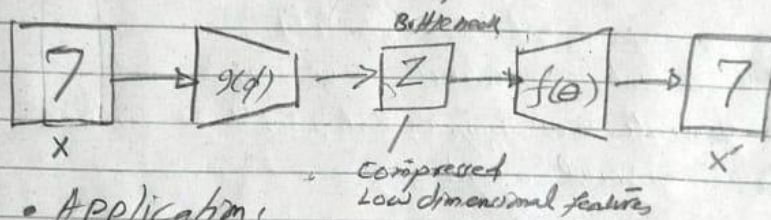
[https://www.youtube.com/watch?v=w8F7\\_rQZxXk&list=PLdxQ7SoCLQANQ9fQcJ0wnnTzkFsJHIWEj&index=19](https://www.youtube.com/watch?v=w8F7_rQZxXk&list=PLdxQ7SoCLQANQ9fQcJ0wnnTzkFsJHIWEj&index=19)

# Variational Autoencoder

- × Review of <sup>stacked</sup> Autoencoder
- × Basic of probability
- × KL Divergence & its significance
- × Derivation of loss function for Variational autoencoder

## Basic Autoencoder

- Architecture: Same image reconstruction



- Applications

- ① Generate compressed image (Dimensional Reduction like PCA)
- ② Denoising (noise removal)
- ③ Image inpainting
- ④ Image Segmentation

$$L(\theta, \phi) = \frac{1}{m} \sum_{i=1}^m [\hat{x}^{(i)} - f_{\theta}(g_{\phi}(\hat{x}^{(i)}))]^2$$

- If no Activation Function is used in autoencoder with 1 hidden layer, this model is sort of identical to PCA.

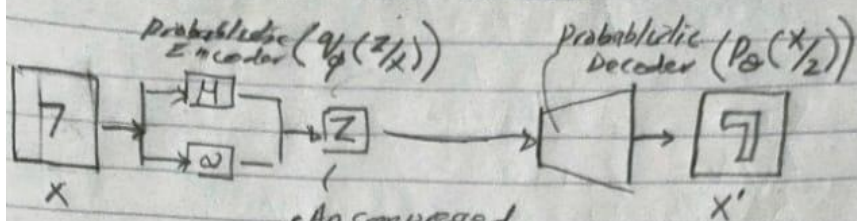
→ In PCA

$$\begin{aligned}
 X &\in \mathbb{R}^n \\
 Z &= U^T X \\
 &\quad \begin{matrix} p \times 1 & p \times n & n \times 1 \end{matrix} \\
 Z &\in \mathbb{R}^p \\
 \hat{X} &= U Z \\
 &\quad \begin{matrix} n \times p & p \times 1 \end{matrix} \\
 \hat{X} &\in \mathbb{R}^n \\
 \min \|X - \underline{U U^T X}\|
 \end{aligned}$$

↳ [due to activation function etc.]

②

## Variational Auto encoder.



- An compressed low dimensional representation of the input
- Sampled latent vector

$$z = \mu + \sigma * \epsilon \quad (\text{Sample } z \text{ from Gaussian distribution})$$

$$\epsilon \sim \mathcal{N}(0, I)$$

- VAE is not reconstruct same sample, but it generate new sample that should be similar to the original sample.

Encoder part

Encoder part

$$L(\theta, \phi) = E_{z \sim q_{\phi}(z|x)} [p_{\theta}(x'|z)] + D_{KL}(q_{\phi}(z|x) || p_{\theta}(z))$$

(learning distribution)      (prior distribution)

- See difference b/w autoencoder (AE) and variational auto encoder (VAE)



•  $P(x)$ : Define the probability of random variable  $x$ .

Example: When a die is tossed once, what is the probability of getting 3.

$$X = \{ \underbrace{1, 2, 3, 4, 5, 6}_{\text{sample space of experiment}} \}$$

random variable

R.V. is a function that assigns real no. to the sample space of experiments

$$P(3) = 1/6$$

•  $P(x/y)$ : Define as the prob. of random variable  $x$  provided  $y$  has happened. Also called an conditional probability.

Example: In tossing a fair die, what is the prob. the 3 has occurred conditioned on the toss being odd.

$$X = \{ \underbrace{1, 2, 3, 4, 5, 6}_{\text{sample space of experiment}} \} \quad \{ \underbrace{1, 3, 5}_{\text{only odd value}} \}$$

sample space of experiment.  
(only odd value)

$$P(3/x) = P(3/(1,3,5)) = 1/3$$

$\neq P(3)$  is increased in case of  $P(3/x \text{ being odd})$

①

Bay's Theorem

Posterior  
prob

$$P(y/x)$$

Attributed value

$$P(x/y) P(y)$$

→ prior prob

(prior information  
available in the  
system)

$$P(x)$$

$$P(y/x) =$$

$$P(x, y)$$

→ joint prob.

$$P(x)$$

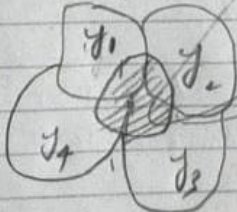
①

Theorem of Total Probability

let  $y_1, y_2, \dots, y_n$  be a set of mutual exclusive events (i.e.  $y_i \cap y_j = \emptyset$ ) and event  $x$  is the union of  $n$  mutually exclusive events, then

$$P(x) = \sum_{i=1}^n P(x/y_i) P(y_i) \quad \text{--- (2)}$$

• let  $y_1, y_2, y_3, y_4$  are mutual exclusive event



Union of 4 ~~into~~ ME

$$P(x) = \sum_{i=1}^4 P(x, y_i)$$

$$P(x) = P(x, y_1) + P(x, y_2) + P(x, y_3) + P(x, y_4)$$

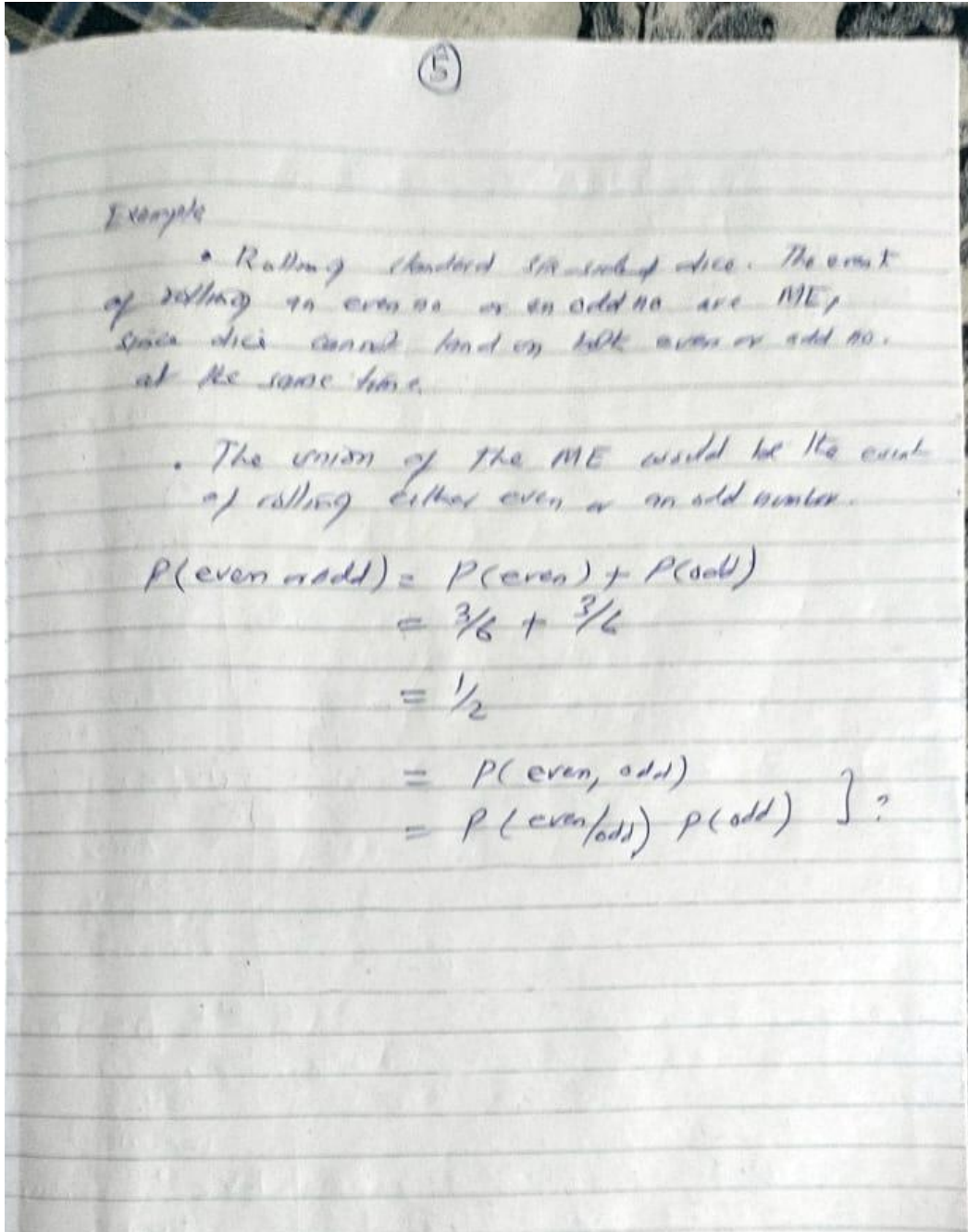
$$\text{or } P(x) = P(x/y_1) P(y_1) + P(x/y_2) P(y_2) + P(x/y_3) P(y_3) + P(x/y_4) P(y_4)$$

By substituting Eq (2) in Eq (1)

$$P(y/x) = \frac{P(x, y)}{\sum_{i=1}^n P(x/y_i) P(y_i)} = \frac{P(x/y) P(y)}{\sum_{i=1}^n P(x/y_i) P(y_i)}$$

- **Example:**

- Rolling a standard six sided die. The events of rolling an **even** number ( $\{2, 4, 6\}$ ) and an **odd** number ( $\{1, 3, 5\}$ ) are mutually exclusive. Since die can not land on both even and odd number at the same time.





6

Weighted avg

Expectation of random variable  $X$ .  $E(X)$

Expectation value of random variable is a weighted average of the possible values of  $X$  that it can take, each value being weighted according to the probability of that event defined as:

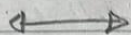
$$E(X) = E(X) = \sum_{i=1}^K x_i \underbrace{P(X=x_i)}_{\substack{\text{state} \\ \text{weight}}} \quad \left\{ \begin{array}{l} \text{// for discrete RV} \\ \text{continuous RV} \\ E(X) = \int x_i P(X=x_i) dx \end{array} \right.$$

$x$  has the Prob. density function  $P$

Example. Let  $X$  represent the outcome of a fair die. What is the  $E(X)$ ?

$$X = \{1, 2, 3, 4, 5, 6\}$$

$$E_p(X) = \sum x_i P(X_i) = 1 \times \frac{1}{6} + 2 \times \frac{1}{6} + 3 \times \frac{1}{6} + 4 \times \frac{1}{6} + 5 \times \frac{1}{6} + 6 \times \frac{1}{6} = 3.5$$



also

we can compute average of Function of random variable

$$E[h(X)] = \sum_{i=1}^K h(x_i) P(X=x_i) \quad \left\{ \begin{array}{l} \text{Function} \\ \text{// For discrete RV} \end{array} \right.$$

$$E[h(X)] = \int h(x) \cdot P(X) dx \quad \left\{ \begin{array}{l} \text{// for Function} \\ \text{continuous RV} \end{array} \right.$$

$$\begin{aligned} \text{Max info} \quad E\left[\log\left(\frac{P(x)}{q(x)}\right)\right] &= \sum_{i=1}^N \log\left(\frac{P(x_i)}{q(x_i)}\right) \cdot P(X=x_i) \\ &= P(X) \times \sum_{i=1}^N \log\left(\frac{P(x_i)}{q(x_i)}\right) = \sum_{i=1}^N \log\left(\frac{P(x_i)}{q(x_i)}\right) \cdot P(X=x_i) \end{aligned}$$

⑦

## Information

- We can quantify the information that a statement can carry by the following expression

$$I = -\log P(x) \quad \text{|| log likelihood}$$

$x =$  is a certain event

$P(x) =$  is the prob of  $x$ ; it has value  $0 \rightarrow 1$

if  $P(x)$  is small,  $-\log(P(x))$  is very large

if  $P(x)$  is large,  $-\log(P(x))$  is very small

- unusual event has high information

eg If we have a statement

"In July the weather will be very hot"

$$\text{its } P(x) = 1$$

$$\therefore I = -\log P(x) = -\log 1 = 0$$

• It is very much likely in July

weather will be very hot so  $P(x) \approx 1$

$\therefore I = 0$ , there will be no information

(b) But if we say in July there is a snow, it is very unlikely, but it contains lot of information.

$P(x) \approx 0$  (as it is very unlikely)

$$I = -\log P(x) = -\log(\text{small no})$$

$I = \text{large value}$

In a fair coin, occurrence of head or tail is an event



(8)

∴ unpredictable / unusual events contain lot of information.

## ENTROPY

Average of Information / Expectation of Information is called ENTROPY.

$$H = - \underbrace{\sum P(x) \log P(x)}_{\text{weighted sum}} \rightarrow \text{Information}(I=x)$$

$$E[X] = \sum x \cdot P(x) = \sum P(x) \cdot x$$

KULLBACK - Leibler

## KL DIVERGENCE

- It has close connection with Information or Entropy
- It is also called relative entropy
- It shows dissimilarity b/w 2 different distributions
- It is a measure of how one prob. distribution is different from the 2nd, reference distribution
- For the discrete prob. distributions  $P \neq Q$ , the KL divergence b/w  $P \neq Q$  is defined as

$$D_{KL}(P \parallel Q) = \underbrace{\sum_x P(x=x) \log P(x)}_{\text{approximate}} - \underbrace{\sum_x Q(x=x) \log Q(x)}_{\text{reference}}$$

$$= \underbrace{\sum_x P(x) \log P(x)}_{\text{average info about P}} - \underbrace{\sum_x Q(x) \log Q(x)}_{\text{average info about Q}}$$

$$\rightarrow \text{Entropy } P - \text{Entropy } Q \Rightarrow KL(P \parallel Q)$$

→ it shows how much the average of 2 dist. are dissimilar

⑨  
 • It is not a distance, as distance is always symmetric (distance b/w 2 cities is always same or symmetric)

• But divergence is not symmetric

$$KL(P||Q) = -\sum P(x) \log P(x) + \sum Q(x) \log Q(x)$$

\* KL divergence is always w.r.t other distribution (expectation w.r.t. P)

$$KL(P||Q) = -\left[ \sum P(x) \log Q(x) \right] + \sum P(x) \log P(x)$$

approx reference

$$KL(P||Q) \approx \text{Entropy of } Q - \text{Entropy of } P$$

$$KL(P||Q) \approx -\sum Q(x) \log Q(x) + \sum P(x) \log P(x)$$

KL is always w.r.t to one distribn: KL of Q w.r.t P

$$= -\sum P(x) \log Q(x) + \sum P(x) \log P(x)$$

Compare

Arg h/o Q  
 w.r.t h/o P

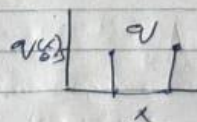
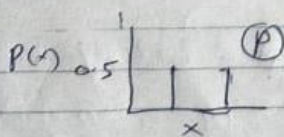
approx ref

$$KL(P||Q) = \sum P(x) \log \frac{P(x)}{Q(x)}$$

Properties of KL divergence:

$$① KL \geq 0$$

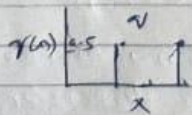
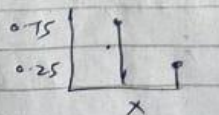
$$② KL(P||Q) \neq KL(Q||P)$$



$$KL(P||Q) = \sum P(x) \log \frac{P(x)}{Q(x)}$$

$$= \sum P(x) \log \frac{0.5}{0.5}$$

$$KL(P||Q) = 0$$



$$KL(P||Q) = 0.75 \log \frac{0.75}{0.5} + 0.25 \log \frac{0.25}{0.5}$$

$$\geq 0$$

$$1) KL(p||q) = 0.5 \log (0.5/0.5) + 0.5 \log (0.5/0.5) = 0$$

$$KL(p||q) = 0.75 \log (0.75/0.5) + 0.25 \log (0.25/0.5) = 0.0568$$

- 2)  $KL(p || q)$  not equal to  $KL(q || p)$ , may find by solving second expression for  
 $KL(q || p) = 0.5 \log (0.5/0.75) + 0.5 \log (0.5/0.25) = 0.0624$



(10)

Suppose we have 2 multivariate normal distributions defined as

$$p(x) = N(x; \mu_1, \Sigma_1)$$

$$q(x) = N(x; \mu_2, \Sigma_2)$$

where  $\mu_1$  &  $\mu_2$  are means and  $\Sigma_1$  &  $\Sigma_2$  are the covariance matrix.

The multivariate normal density function is defined as

$$p(x) = N(x; \mu_1, \Sigma_1) = \frac{1}{\sqrt{(2\pi)^k |\Sigma_1|}} \exp\left(-\frac{1}{2}(x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1)\right)$$

$$q(x) = N(x; \mu_2, \Sigma_2) = \frac{1}{\sqrt{(2\pi)^k |\Sigma_2|}} \exp\left(-\frac{1}{2}(x-\mu_2)^T \Sigma_2^{-1} (x-\mu_2)\right)$$

• If the ~~2~~ both have some dimension  $k$ .

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \end{bmatrix} \quad \text{each variable follows } N(0,1)$$

Recap:

• Multivariate PD is used to describe the behavior of multiple R.V that are correlated. It is extension of Univariate normal distribution to 2 or more variables where each variable follows a normal distribution & joint distribution can be represented by mean vector & covariance matrix.

• PDF for single variable  $f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$



(11)

Now KL divergence b/w  $P \neq Q$

$$D_{KL}(P(x) \parallel Q(x)) = \frac{1}{2} \left[ \log \frac{|\Sigma_2|}{|\Sigma_1|} - k + \text{tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right]$$

Prove?

Proof

$$KL(P(x) \parallel Q(x)) = \sum_x P(x) \log \left( \frac{P(x)}{Q(x)} \right) \quad \text{--- (1)}$$

We know multivariate normal distribution can be defined as

$$P(x) = \frac{1}{\sqrt{(2\pi)^k |\Sigma_1|}} \exp \left( -\frac{(\mathbf{x} - \mu_1)^T \Sigma_1^{-1} (\mathbf{x} - \mu_1)}{2} \right)$$

determinant of  $\Sigma_1$

$$P(x) = (2\pi)^{-k/2} |\Sigma_1|^{-1/2} \exp \left( -\frac{1}{2} (\mathbf{x} - \mu_1)^T \Sigma_1^{-1} (\mathbf{x} - \mu_1) \right)$$

$$P(x) = (2\pi)^{-k/2} |\Sigma_1|^{-1/2} \exp \left[ -\frac{1}{2} (\mathbf{x} - \mu_1)^T \Sigma_1^{-1} (\mathbf{x} - \mu_1) \right]$$

Taking log on both side

$$\log P(x) = -\frac{k}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_1| - \left[ \frac{1}{2} (\mathbf{x} - \mu_1)^T \Sigma_1^{-1} (\mathbf{x} - \mu_1) \right] \quad \text{--- (2)}$$

Similarly for  $Q(x)$  prob. density function

$$\log Q(x) = -\frac{k}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_2| - \left[ \frac{1}{2} (\mathbf{x} - \mu_2)^T \Sigma_2^{-1} (\mathbf{x} - \mu_2) \right] \quad \text{--- (3)}$$



(12)

Equation ① can be written as

$$KL(P(x) \parallel Q(x)) = \sum_x P(x) [\log P(x) - \log Q(x)]$$

Now put equation ② &amp; ③ in above equation

$$KL(P(x) \parallel Q(x)) = \sum_x P(x) \left[ -\frac{K}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma| \right. \\ \left. - \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + \frac{K}{2} \log(2\pi) + \frac{1}{2} \log |\Sigma_2| \right. \\ \left. + \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \right]$$

$$KL(P(x) \parallel Q(x)) = \sum_x P(x) \left[ \frac{1}{2} \log \frac{|\Sigma_2|}{|\Sigma_1|} + \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) - \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right]$$

How To solve these ② expression  
let consider part by part

$$\sum_x P(x) \cdot \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) = E_p \left[ \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right]$$

As we know

• Expectation of RV

$$E_p[x] = \sum_{i=1}^n x_i \cdot P(x_i)$$

weight

• Expectation of Function of RV

$$E_p[h(x)] = \sum_{i=1}^n h(x_i) \cdot P(x_i)$$

Note

$$\sum_x P(x) \frac{1}{2} \log \frac{|\Sigma_2|}{|\Sigma_1|}$$

$$= E_p \left[ \frac{1}{2} \log \frac{|\Sigma_2|}{|\Sigma_1|} \right]$$

$$= \frac{1}{2} \log \frac{|\Sigma_2|}{|\Sigma_1|}$$

as  $E_p[\text{const}] = \text{const}$



# Recap Linear Algebra.

- If  $X$  is square matrix

$$\begin{bmatrix} d_1 & & \\ & d_2 & \\ & & \ddots \\ & & & d_n \end{bmatrix} \sum_{i=1}^n d_i$$

①

Trace & Expectation trick

NDV

$$E(X^T A X) = E(\text{tr}(X^T A X))$$

→ If  $x$  is scalar, then

→  $x$

$$E(x) = E(\text{tr}(x))$$

$$= E(\text{tr}(A X X^T))$$

• since trace of  $x$  is scalar

$$= \text{tr}(E(A X X^T))$$

$$\rightarrow \text{tr}(CAB) = \text{tr}(BA)$$

$$\rightarrow \text{tr}(ABC) = \text{tr}(BCA) = \text{tr}(CAB)$$

Let's rewrite the expression of previous page

$$\rightarrow \text{tr}(ABC) \neq \text{tr}(ACB)$$

$$\rightarrow E(\text{tr}(x)) = \text{tr}(E(x))$$

$$\frac{1}{2} E_p[(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1)]$$

$$= E_p[\text{tr}(\frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1))]$$

$$= E_p[\text{tr}(\frac{1}{2} (x - \mu_1)(x - \mu_1)^T \Sigma_1^{-1})]$$

$$= \text{tr}(E_p[\frac{1}{2} (x - \mu_1)(x - \mu_1)^T] \Sigma_1^{-1})$$

Covariance matrix

$$= \text{tr}[\Sigma_1 \frac{1}{2} \Sigma_1^{-1}] = \text{tr}[\mathbf{I}_K]$$

$$= K$$

$$K = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & & \\ & & \ddots & \\ 0 & & & 1 \\ & & & & 0 \end{bmatrix}$$

14

Now solve the other expression from (7)

$$\begin{aligned} & \sum_x P(x) \left[ \frac{1}{2} (x - \mu_1)^T \Sigma_2^{-1} (x - \mu_2) \right] \\ & \quad \rightarrow \text{Note add \& subtract } \mu_1 \\ & = \sum_x P(x) \left[ \frac{1}{2} \left[ \underbrace{(x - \mu_1)}_A + \underbrace{(\mu_1 - \mu_2)}_B \right]^T \Sigma_2^{-1} \left[ \underbrace{(x - \mu_1)}_A + \underbrace{(\mu_1 - \mu_2)}_B \right] \right] \\ & = \sum_x P(x) \left[ \frac{1}{2} (x - \mu_1)^T \Sigma_2^{-1} (x - \mu_1) + \frac{1}{2} (x - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right. \\ & \quad \left. + \frac{1}{2} (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right] \end{aligned}$$

$$\begin{aligned} & (A+B)^T \Sigma_2^{-1} (A+B) \\ & = (A^T + B^T) \Sigma_2^{-1} (A+B) \\ & = A^T \Sigma_2^{-1} A + B^T \Sigma_2^{-1} B + \\ & \quad B^T \Sigma_2^{-1} A + A^T \Sigma_2^{-1} B \\ & = A^T \Sigma_2^{-1} A + B^T \Sigma_2^{-1} B + 2A^T \Sigma_2^{-1} B \end{aligned}$$

$$= E_p \left[ \frac{1}{2} (x - \mu_1)^T \Sigma_2^{-1} (x - \mu_1) + (x - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2) + \frac{1}{2} (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right]$$

$$= E_p \left[ \frac{1}{2} (x - \mu_1)^T \Sigma_2^{-1} (x - \mu_1) \right] + E_p \left[ (x - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right] + E_p \left[ \frac{1}{2} (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right]$$

$$\begin{aligned} & = E_p \left[ \text{tr} \left\{ \frac{1}{2} (x - \mu_1)^T \Sigma_2^{-1} (x - \mu_1) \right\} \right] + \text{tr} \left\{ \Sigma_2^{-1} \Sigma_1 \right\} + \text{tr} \left\{ \Sigma_2^{-1} \Sigma_1 \right\} \\ & \quad \quad \quad \downarrow \text{Const} \\ & = \text{tr} \left\{ \frac{\Sigma_2^{-1} \Sigma_1}{2} \right\} + 0 + (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \end{aligned}$$

(as did previously)

(8)

$$E_p \left[ (x - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right] = [E_p(x) - E_p(\mu_1)]^T \Sigma_2^{-1} (E_p(\mu_1) - E_p(\mu_2))$$

$$= (\mu_1 - \mu_1)^T \Sigma_2^{-1} (\mu_2 - \mu_1) = 0 \text{ proved}$$

Substitute (2), (8) in (7) we get

$$KL(P(x) \| Q(x)) = \frac{1}{2} \left[ \log \frac{|\Sigma_2|}{|\Sigma_1|} - K + \text{tr} \left( \Sigma_2^{-1} \Sigma_1 \right) + (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right]$$



(15)

• we proved that the two Gaussian

having  $N(\mu_1, \Sigma_1)$  &  $N(\mu_2, \Sigma_2)$  have

~~the~~ the KL divergence as shown in previous page.

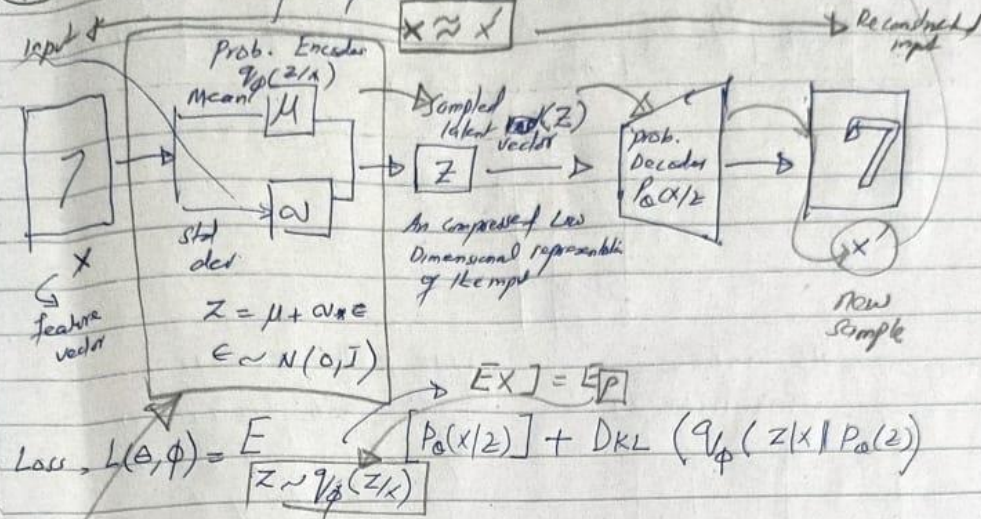
• It is important expression & we will see it



Dataset  
Generate blue sample

(16)

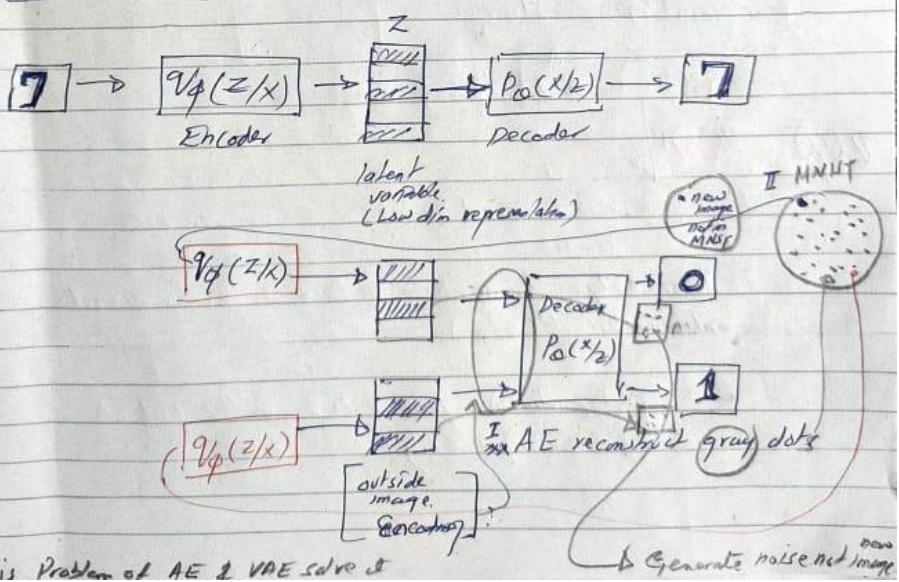
# Working of VAE : Top Level view



Why we need VAE

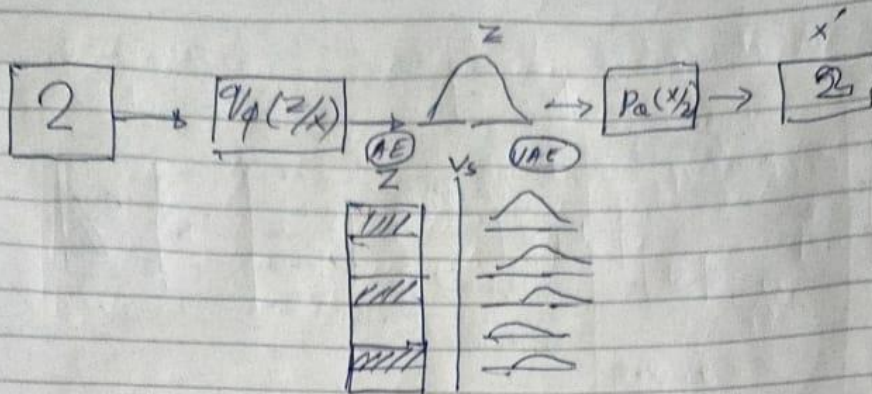
The Goal of VAE: The goal of VAE is to find a distribution  $q_\phi(z/x)$  of some latent variable ( $z$ ), which can sample from  $z \sim q_\phi(z/x)$  to generate new sample  $x' \sim P_\theta(x/z)$

## Typical AE



It is Problem of AE & VAE solve it

# VAE (17)



• AE produce prob. vector ( $z$ )

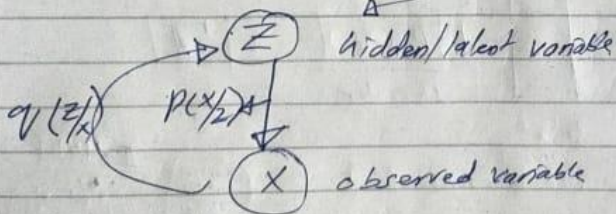
• VAE produce the range of value (prob density function) of each encoding of each encoded value of  $z$ .

• Now if input image to VAE is different from MNIST it still generate reasonable new image.

→ we can't do system

## LATENT VARIABLE | $z \Rightarrow$ hidden variable

We model system as a collection of random variables. Here



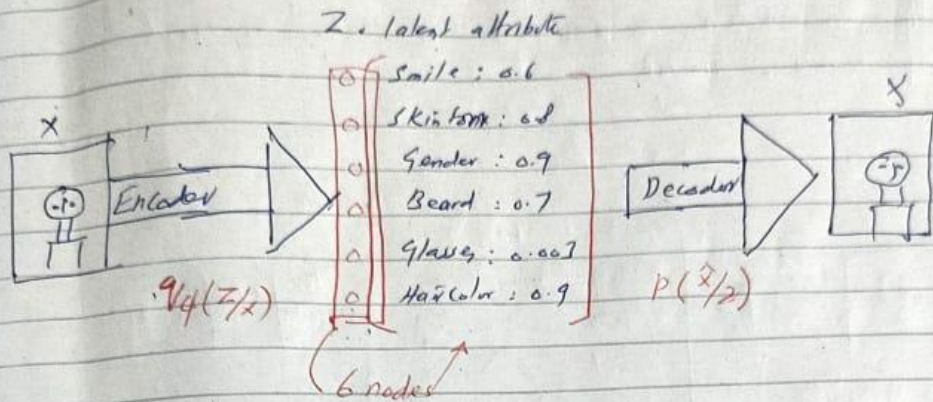
• The edge drawn from  $z$  to  $x$  is the conditional distribution  $P(x/z)$

• The edge drawn from  $x$  to  $z$  is the conditional distribution  $q(z/x)$



(8)

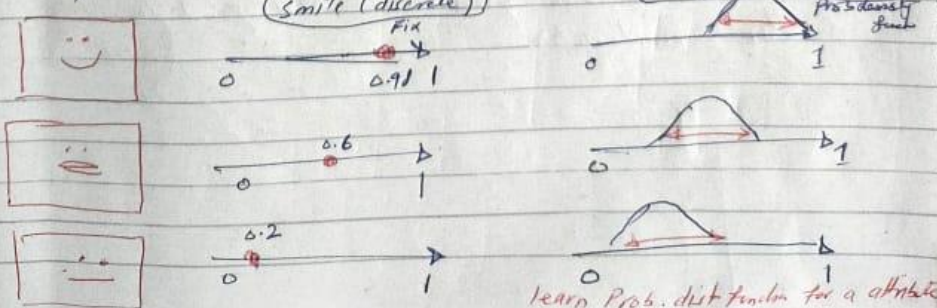
- ① Latent variables corresponding to a real feature of the object that have not been measured (may be technology is not available to do that)



- In above examples, we framed AE on larger dataset of faces with encoding dimension of 6. An ideal AE will learn the descriptive attributes of faces such as smile, skin tone etc. in order to describe observation in some compressed form.

- In above example, we have described the input image in term of latent variables using single value to describe each attributes. For instance, what single value you will assign for photo of monalisa?

- using VE, we define latent attributes in probabilistic terms AE

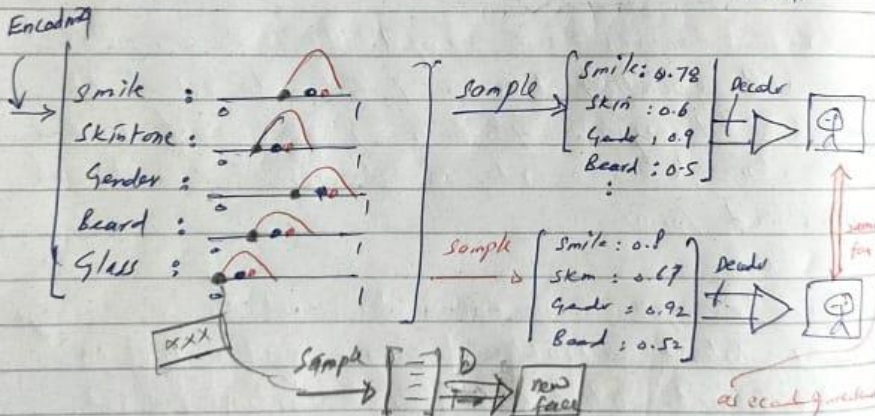
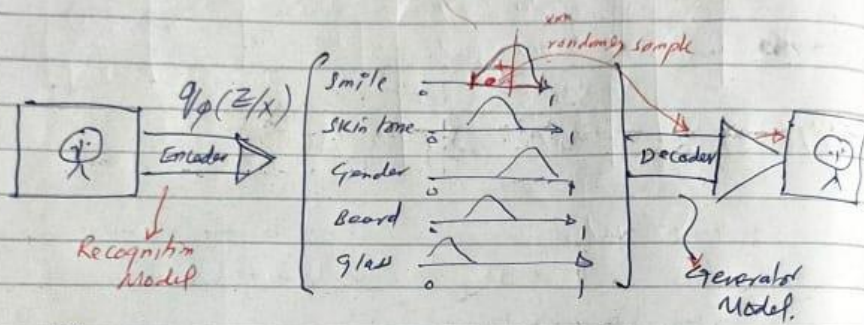




→ during Training we use that range of values to generate faces so it works

(19)

With this approach, we now represent each latent attribute for a given input as a probability distribution. When decoding, we will randomly sample from each latent state distribution to generate a vector as input for decoder model.



• By constructing our encode model to output range of possible values (or statistically distribution) from which we will randomly sample to feed into our decoder model. The values which are near each other in latent space must correspond to similar reconstruction.

if

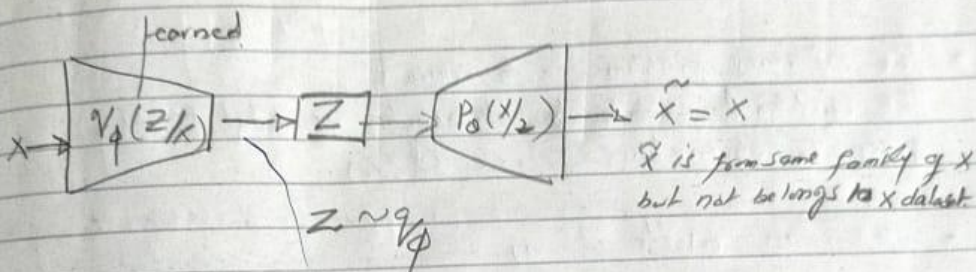
(20)

## Derivation of The Loss Function

Goal

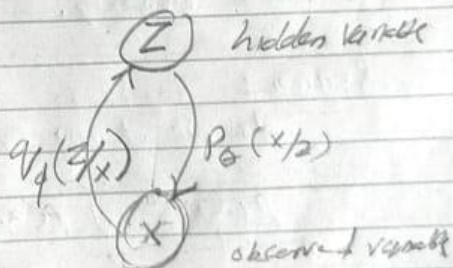
Let's recall the goal of VAE:

The goal of VAE is to find a distribution  $q_\phi(z/x)$  of same latent variables ( $z$ ), which we can sample from  $z \sim q_\phi(z/x)$  to generate new samples  $x'$  from  $p_\theta(x/z)$ .



$\phi, \theta$  are learned through B.P

## Graphical Model



- \*  $x$  is feed to Encoder,
- \* we have to learn  $q_\phi(z/x)$  & then samp

$z \sim q_\phi(z/x)$ . Feed  $z$  to  $p_\theta(x/z)$  to generate the new sample ( $\tilde{x}$ )

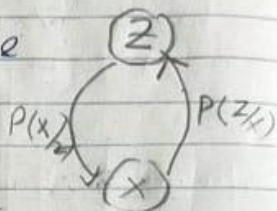


(21)

The Derivation of loss function <sup>of VAE</sup> is solved through the problem approximate inference.

The Problem Of Approximate Inference :

- let  $x$  be a set of observed variables & let  $z$  be the set of latent variables with joint distribution  $P(z, x)$ .



- The inference problem is to compute the conditional probability distribution of latent variables given the observations ( $x$ ) - i.e  $P(z|x)$

we can write it as

$$P(z|x) = \frac{P(x|z)P(z)}{P(x)} \quad \text{--- (A)}$$

Evaluating (A) is difficult because  $p(x)$  cannot be solved properly.

Reason :

As we know by the law of total probability

$$\begin{aligned} P(x) &= \int_z P(x|z) P(z) dz \\ &= \int_z P(x, z) dz \end{aligned}$$

This integral is not available in close form or in tractable (ie requiring exponential time to compute) due to multiple integrals involved for latent variable vector  $z$ .

$z = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$    
 smile   
 skin tone   
 glasses   
 :   
 :



(22)

Variational Inference (VI)

$$P(Z/x) = \frac{P(x/z) P(z)}{P(x)}$$

observed      prior      normalizing constant

• (VI) is used to solve this problem

• Approximating  $P(Z/x)$  with another distribution  $q_\phi(z/x)$

G.T

Prediction

• Goal is to approximate intractable posterior distribution  $P(Z/x)$  with simpler distribution ( $q_\phi(z/x)$ ), such as Gaussian distribution

By choosing tractable distribution (like Gaussian distribution), let say  $Q(Z/x)$  and play with the parameters of  $Q(Z/x)$  distribution it become closer to  $P(Z/x)$

IDEA

we make  $Q(Z/x) \approx P(Z/x)$ , by minimizing  $KL(Q||P)$  [as close as possible]

as  $KL(Q||P)$  minimization make  $Q \approx P$ .

$$D_{KL}(Q_\phi(Z/x) || P_\theta(Z/x)) = \sum_z Q_\phi(Z/x) \log \frac{Q_\phi(Z/x)}{P_\theta(Z/x)}$$

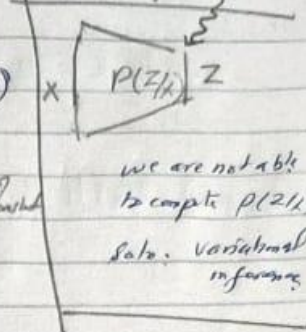
$$= E_{Z \sim Q_\phi(Z/x)} \left[ \log \frac{Q_\phi(Z/x)}{P_\theta(Z/x)} \right]$$

$$= E_{Z \sim Q_\phi(Z/x)} [\log(Q_\phi(Z/x)) - \log(P_\theta(Z/x))]$$

$$= E_{Z \sim Q_\phi(Z/x)} [\log Q_\phi(Z/x)] - E_{Z \sim P_\theta(Z/x)} [\log P_\theta(Z/x)]$$

— (B)

make Gaussian  $Q(Z/x)$  distribution close  $P(Z/x)$



$$(23) \quad Z = z \sim Q_q(z/x)$$

Substituting (A) in (B)

$$\begin{aligned} D_{KL}(Q_q(z/x) \parallel P_0(z/x)) &= E_z [\log Q_q(z/x)] \\ &\quad - E_z \left[ \log \frac{P_0(x/z) P(z)}{P(x)} \right] \\ &= E_z [\log Q_q(z/x)] - E_z \left[ \log \frac{P(z, x)}{P(x)} \right] \quad \text{Bayes' Theorem} \\ &= E_z [\log Q_q(z/x)] - E_z [\log P(z, x)] + E_z [\log P(x)] \\ &= E_z [\log Q_q(z/x)] - E_z [\log P(z, x)] + \int Q_q(z/x) \log P(x) dz \\ &= E_z [\log Q_q(z/x)] - E_z [\log P(z, x)] + \log P(x) \int Q_q(z/x) dz \\ &= E_z [\log Q_q(z/x)] - E_z [\log P(z, x)] + \log P(x) \end{aligned}$$

$$\log P(x) = \underbrace{E_z [\log P(z, x)] - E_z [\log Q_q(z/x)]}_{\text{Component 1}} + \underbrace{D_{KL}(Q_q(z/x) \parallel P_0(z/x))}_{\text{Component 2}}$$

$\downarrow$  Fix no  $\downarrow$   $\downarrow$   $\downarrow$   $\downarrow$   
 Component 1  $\uparrow$  max  $\downarrow$  min Component 2  
 $\downarrow$   $\downarrow$   $\downarrow$   $\downarrow$   
 we want to min. it

- In order to make the above equation equally true, if we minimize component (2) we have to maximize the component 1 (1).
- Component 1 is called Evidence Lower Bound (ELBO).
- Now if we maximize ELBO, it indirectly minimizes the KL divergence.



Q4)

• Why  $\mathcal{L}$  is called lower bound

As  $KL \geq 0$  (it is property of KL)

$$\therefore \log p_a(x) = \mathcal{L} + \underbrace{(\text{+ve no})}_{\geq 0}^{KL}$$

$$\mathcal{L} = \log p_a(x) - KL$$

• If  $KL = 0$

$$\mathcal{L} = \log p_a(x)$$

• If  $KL > 0$

$$\mathcal{L} < \log p_a(x)$$

$$\therefore \mathcal{L} \leq \log p_a(x)$$

So it is lower bound of  $\log p_a(x)$

NOW As

$$\begin{aligned} \mathcal{L} = ELBO &= E_z [\log p(x/z)] - E_z [\log q_\phi(z/x)] \\ &= E_z [\log p(x/z) p_a(z)] - E_z [\log q_\phi(z/x)] \\ &= E_z [\log p(x/z)] + \underbrace{E_z [\log p_a(z)] - E_z [\log q_\phi(z/x)]}_{\geq 0} \\ &= E_z [\log p(x/z)] - E_z [\log \underbrace{q_\phi(z/x)}_{p_a(z)}] \end{aligned}$$

$$\mathcal{L} = E_z [\log p(x/z)] - D_{KL}[q_\phi(z/x) \parallel p_a(z)]$$

↙ expected reconstruction error

↘

if  $z$  is gaussian, then reconstruction is square error  
 $p(x/z) = \frac{1}{\sigma^2} e^{-|x-z|^2/\sigma^2}$   
 $\log p(x/z) = -|x-z|^2/\sigma^2$  if  $z$  is better, then reconstruction error is smaller

learn  $\approx$  prior distribution



# Expectation of random variable $X$ ( $E(X)$ )

Expectation

Random variable  
is a function that  
assigns real no. to  
the sample space  
components

e.g.

Experiment

- Pick Red Ball  
Score is 10
- pick Blue Ball  
Score is 20
- pick Green Ball  
Score is 40

$$U = \{Red, Blue, Green\}$$

$$X = \text{Score of Pick}$$

$$X(Red) = 10$$

$$X(Blue) = 20$$

$$X(Green) = 40$$

$$X = \{10, 20, 40\}$$

Random  
variable

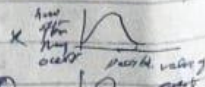
use (PMF)

R.V.  $\rightarrow$  Discrete (countable)

$\rightarrow$  Continuous

$\rightarrow$  (PDF)

prob. distribution is finite  
that shows possible values  
of a variable & how  
often it occurs



$$\Rightarrow$$