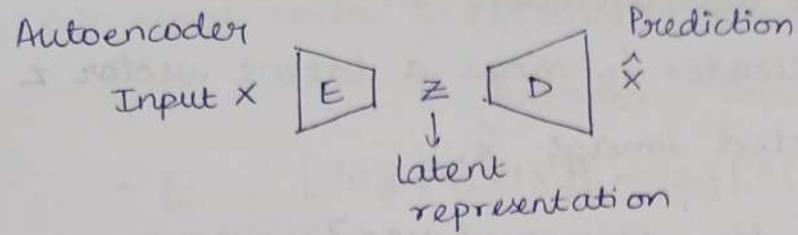


13 Jan 25
Monday.

Computer Vision.



Eg: Denoising autoencoders.

Project noisy image to latent space which is again reconstructed back to give denoised image.

Encoder E : $f_E: X \times \Theta_E \rightarrow Z \rightarrow f_E(x, \theta_E) \rightarrow Z$.

Decoder D : $f_D: Z \times \Theta_D \rightarrow X \hookrightarrow f_D(z, \theta_D) \rightarrow \hat{x}$.

X : input space

Θ_E : Encoder parameter space.

Z : Latent space

Θ_D : Decoder parameter space.

$$(\theta_E, \theta_D)^* = \arg \min_{\theta_E, \theta_D \in \Theta_E \times \Theta_D} L(\theta_E, \theta_D) \rightarrow \text{minimize the reconstruction loss.}$$

$$= \arg \min_{\theta_E, \theta_D} \sum_{i=1}^{N_{tr}} (x^i - \hat{x}^i)^2$$

A potential candidate for f_E is a CNN.

↓
composition of non-linear functions.

A potential candidate for f_D is Deconv Net.

Q. Is it supervised or unsupervised?

Ans:- We generally assume it is unsupervised.

We are trying to predict the input itself.

$p(x)$ is ^{not known} unknown to us.

Note: The encoder f_E maps an input image to a latent representation z .

The decoder f_D maps a latent vector z to a predicted image \hat{x} .

Q. Difference b/w this and VAE?

A. VAE would map to distribution.

We assume smoothness in the latent representation.

Variational Auto Encoder (VAE)

- The encoder $Q(z|x)$ maps an image to a latent representation $N(\mu(x), \Sigma(x))$.
- The decoder maps a latent sample z drawn from $N(0, I)$ to $N(f(z, \theta_D), \sigma^2 I)$
- $N(\mu(x), \Sigma(x))$ need not be a standard normal distribution.

The encoder tries to make the KL-divergence of $N(\mu(x), \Sigma(x))$ and standard normal dist. as close as possible.

- Encoder E : $f_E : X \times \Theta_E \rightarrow \mathbb{Z}$
 $D : \mathbb{Z} \times \Theta_D \rightarrow X$.

Specifically: $Q(z|x)$ is an encoder that maps x to $N(\mu(x), \Sigma(x))$.

$P(x|z)$ is a decoder that maps z to $N(f(z), \sigma^2 I)$

Note that conditioning is important here.

Recall: We start with $D(Q(z|x) \parallel P(z|x))$

Encoder unknown

$$= \mathbb{E}_{z \sim Q(z|x)} [\log [Q(z|x)] - \log [P(z|x)]].$$

$$= (P(z|x) \cdot P(x) = P(x|z) P(z)]$$

$$= D(Q(z|x) \parallel P(z)) - E[\log P(x|z)] + \log P(x).$$

We assume that
this is standard
dist. That is what
we assume for prior.

Here there is
no expectation
because it
doesn't depend
on z .

$$\underbrace{\log P(x) - D(Q(z|x) \parallel P(z|x))}_{\text{unknown}} = \underbrace{E[\log P(x|z)]}_{\text{Decoder}} - \underbrace{D(Q(z|x) \parallel P(z))}_{\text{Encoder.}}$$

$P(x|z)$ is assumed to be gaussian.

$P(z)$ is assumed to be standard normal.

Recall $P(x|z) \sim N(f(z, \theta_b), \sigma^2 I)$

$Q(z|x) \sim N(\mu(x), \Sigma(x))$

$P(z) \sim N(0, I)$.

20/01/2025
Monday

Computer Vision.

Q. How does encoder help in training?

Takes image space, maps to distribution in latent space. Ensures homogeneity / smoothness in manifold. Designed in such a way that conditional distribution is similar to ~~eggs~~ standard gaussian.

- We want to work with tractable solution and make the optimization tractable.
- The RHS terms: the KL divergence b/w encoder output and $P(z)$ is tractable.
- KL divergence b/w two standard normal dist is 0.
- We want to bring the loss to be similar to KL divergence.

Q. Why small time steps?

Non-equilibrium thermodynamics.

Brownian motion in small time steps.

Processes are identical in both forward time and reverse time.

Both distributions have same functional form which is gaussian and loss function of both fwd and bwd would be KL divergence.

Forward process — diffusion till time step T
when it would be normal dist.

Note that this process is controlled.

Controlled amount of noise is added from
time step 0 to time step t to generate x_t
which we sample at time t .

x_t can be expressed in terms of x_0 and some
noise.

β 's can either be learned or determined.

$q(x_0)$ would be distribution of data (original).

→ We need to learn parameters of backward
process such that likelihood is maximized. i.e,
when we draw samples from $p(x_0)$, they should
be similar to those picked from $q(x_0)$.

These are parameters of multivariate gaussian
distribution.

→ We learn how much noise has to be added
at each time step.

→ Note that functional form is same for fwd and bwd
but not identical.

• Markov chain : Both forward and backward
trajectories are Markov chains, specifically
first order.

$q(x^{(0 \dots T)})$ is the forward distribution.

$p_0(x^{(0 \dots T)})$ is the reverse distribution.

Ideally $P_\theta(x^{(0)}) \sim q(x^{(0)})$
 to be learned $\underbrace{q(x^{(0)})}_{\text{data dist.}}$

→ q is deterministic

P has to be parameter θ , which we are going to learn.

$$\text{goal: } \max_{\theta \in \Theta} L(\theta) = E_{x^{(0)} \sim q} [\log P_\theta(x^{(0)})] \quad \text{--- (1)}$$

i.e. we want to maximize the likelihood of generating samples $P_\theta(x^{(0)})$ coming from $q(x^{(0)})$.

At time T , both forward and reverse dist. would have gaussian dist.

$$\begin{aligned} P_\theta(x^{(0)}) &= \underbrace{\int P_\theta(x^{(0 \dots T)}) dx^{(1 \dots T)}}_{\rightarrow \equiv P_\theta(x^{(0)}, x^{(1)}, \dots, x^{(T)})} \\ &= \int P_\theta(x^{(T)}) \prod_{t=1}^T P_\theta(x^{(t-1)} | x^{(t)}) dx^{(1 \dots T)} \quad \text{--- (2)} \end{aligned}$$

(standard result for first order Markov chains)

$$\begin{aligned} [P_\theta(x^{(0 \dots T)})] &= P_\theta(x^{(0)} | x^{(1 \dots T)}) P_\theta(x^{(1 \dots T)}) \\ &\quad \underbrace{\qquad\qquad\qquad}_{\text{(due to Markov property)}} \\ &= P_\theta(x^{(0)} | x^{(1)}) \cdot P_\theta(x^{(1 \dots T)}) \\ &= P_\theta(x^{(T)}) \prod_{t=1}^T P_\theta(x^{(t-1)} | x^{(t)}) \end{aligned}$$

Also we assume $P_\theta(x^{(t-1)}|x^{(t)}) \sim \mathcal{N}(x^{(t-1)}, f_\mu(x^{(t)}, \theta), f_\Sigma(x^{(t)}, \theta))$ — (3).

$$\begin{aligned}
 P_\theta(x^{(0)}) &= \int P_\theta(x^{(T)}) \prod_{t=1}^T P_\theta(x^{(t-1)}|x^{(t)}). dx^{(1\dots T)} \\
 &= \int P_\theta(x^{(T)}) \prod_{t=1}^T P_\theta(x^{(t-1)}|x^{(t)}). \frac{q(x^{(1\dots T)}|x^{(0)})}{q(x^{(1\dots T)}|x^{(0)})} dx^{(1\dots T)} \\
 &= \int P_\theta(x^{(T)}) q(x^{(1\dots T)}|x^{(0)}) \underbrace{\prod_{t=1}^T P_\theta(x^{(t-1)}|x^{(t)})}_{\text{Since forward trajectory is a markov chain.}} dx^{(1\dots T)}
 \end{aligned}$$

21/01/2025
Tuesday

Computer Vision

Recap:

- Diffusion Process → fwd and bwd trajectories have the same dist. function form.
 - Maximizing log likelihood.
 - Markov chain assumption
 - Conditional dist. Gaussian → closed form expression for bound.
- Lower bound on the expected log likelihood.

Recap:

$$L(\theta) = \int q(x^{(0)}) \cdot \log P_\theta(x^{(0)}) \cdot dx^{(0)} \quad \text{--- (1)}$$

$$\theta^* = \underset{\theta \in \Theta}{\operatorname{argmax}} L(\theta) \quad \text{--- (2)}$$

$$P_\theta(x^{(0)}) = \int P_\theta(x^{(0 \dots T)}) \cdot dx^{(1 \dots T)}$$

$$= \int P_\theta(x^{(T)}) \cdot \prod_{t=1}^T P_\theta(x^{(t-1)} | x^{(t)}) \cdot dx^{(1 \dots T)}$$

$$P_\theta(x^{(0)}) = \int P_\theta(x^{(T)}) \cdot q(x^{(1 \dots T)} | x^{(0)}) \cdot \frac{\prod_{t=1}^T P_\theta(x^{(t-1)} | x^{(t)})}{\prod_{t=1}^T q(x^{(t)} | x^{(t-1)})} \cdot dx^{(1 \dots T)}$$

Plug (3) into (1) --- (3)

$$L(\theta) = \int q(x^{(0)}) \cdot \log \left[\int P_\theta(x^{(T)}) \cdot q(x^{(1 \dots T)} | x^{(0)}) \cdot \frac{\prod_{t=1}^T P_\theta(x^{(t-1)} | x^{(t)})}{\prod_{t=1}^T q(x^{(t)} | x^{(t-1)})} \cdot dx^{(1 \dots T)} \right] dx^{(0)}$$

We will apply the Jensen's inequality.
for a concave function, $Q(x)$,

$$Q(E(x)) \geq E(Q(x))$$

We pull expectation outside and log inside.

$L(\theta) \geq k$, where

$$k = \int q(x^{(0)}) \cdot q(x^{(1\dots T)} | x^{(0)}) \log P_\theta(x^{(T)}) \frac{\prod P_\theta(x^{(t-1)} | x^{(t)})}{\prod q(x^{(t)} | x^{(t-1)})} dx^{(0\dots T)}$$

$$= \int q(x^{(0\dots T)}) \log P_\theta(x^{(T)}) \underbrace{\frac{\prod P_\theta(x^{(t-1)} | x^{(t)})}{\prod q(x^{(t)} | x^{(t-1)})}}_{\text{---(4)}}$$

Since the forward trajectory is a forward chain,

$$q(x^{(t)} | x^{(t-1)}) = q(x^{(t)} | x^{(t-1)}, x^{(0)})$$

↳ Since This gives the flexibility of working with image.

Since $x^{(0)}$ is the image R.V.

$$q(x^{(t)} | x^{(t-1)}) = q(x^{(t-1)} | x^{(t)}, x^{(0)}) \cdots q(x^{(t)} | x^{(0)}) \quad \text{---(5)}$$

Plug ⑤ into ④.

$$k = \int q(x^{(0\dots T)}) \log P_\theta(x^{(T)}) \frac{\prod P_\theta(x^{(t-1)} | x^{(t)})}{\prod q(x^{(t-1)} | x^{(t)}, x^{(0)}) \cdot q(x^{(t)} | x^{(0)})} q(x^{(t)} | x^{(0)}) dx^{(0\dots T)} \quad \text{---(6)}$$

Check: We want to express ⑥ as a sum of KL divergence terms and entropy terms.

$$K = \int q(x^{(0\dots T)}) \log p_\theta(x^{(T)}) dx^{(0\dots T)} + \\ \int q(x^{(0\dots T)}) \log \left[\frac{\prod p_\theta(x^{(t-1)}|x^{(t)})}{\prod q(x^{(t-1)}|x^{(t)}, x^{(0)})} \cdot \frac{q(x^{(t-1)}|x^{(0)})}{q(x^{(t)}|x^{(0)})} \right] dx^{(0\dots T)}$$

24/10/2025

Friday.

Computer Vision.

Recap.

- Diffusion Process - forward and backward trajectories.

- Markov chains

$$q(x^{(0\dots T)}) = q(x^{(0)}). \prod_{t=1}^T q(x^{(t)}|x^{(t-1)}).$$

$$p_\theta(x^{(0\dots T)}) = p_\theta(x^{(T)}) \prod_{t=1}^T p_\theta(x^{(t-1)}|x^{(t)}).$$

- Marginal:

$$p_\theta(x^{(0)}) = \int p_\theta(x^{(0\dots T)}). dx^{(1\dots T)}.$$

- Rewriting

$$p_\theta(x^{(0)}) = \int p_\theta(x^{(T)}. q(x^{(1\dots T)}|x^{(0)}). \underbrace{\prod_{t=1}^T p_\theta(x^{(t-1)}|x^{(t)})}_{\overbrace{\prod_{t=1}^T q(x^{(t)}|x^{(t-1)})}^{\text{Reverse}}} dx^{(1\dots T)}$$

forward

• Problem :

$$\max \mathcal{L}(\theta) = \max \int q(x^{(0)}). \log p_{\theta}(x^{(0)}) . dx^{(0)}$$

• Lower bound :

$$\mathcal{L}(\theta) \geq K(\theta) \quad (\text{Jensen's inequality}).$$

$$K(\theta) = \int \cancel{q(x^{(0)})} \cdot q(x^{(0 \dots T)}) . \log \left[p_{\theta}(x^{(T)}) \prod_{t=1}^T p_{\theta}(x^{(t-1)} | x^{(t)}) \right] \\ \frac{\prod_{t=1}^T q(x^{(t)} | x^{(t-1)})}{dx^{(0 \dots T)}} \quad \text{--- (1)}$$

• Using Markov property :

$$q(x^{(t)} | x^{(t-1)}) = q(x^{(t)} | x^{(t-1)}, x^{(0)}) \\ = q(x^{(t-1)} | x^{(t)}, x^{(0)}) \cdot \frac{q(x^{(t)} | x^{(0)})}{q(x^{(t-1)} | x^{(0)})} \quad \text{--- (2)}$$

Plugging (2) into (1), we get

$$K(\theta) = \int q(x^{(0 \dots T)}) \left[\log \left[p_{\theta}(x^{(T)}) + \log \frac{\prod_{t=1}^T p_{\theta}(x^{(t-1)} | x^{(t)})}{\prod_{t=1}^T q(x^{(t-1)} | x^{(t)}, x^{(0)})} \cdot \frac{q(x^{(t-1)} | x^{(0)})}{q(x^{(t)} | x^{(0)})} \right] \right] dx^{(0 \dots T)} \quad \text{--- (3)}$$

Recall :

$$KL(N(x; \mu_0, \Sigma_0) || N(x; \mu_1, \Sigma_1)), \quad (x \in \mathbb{R}^d)$$

$$= \frac{1}{2} (\text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^T \Sigma_1^{-1} (\mu_1 - \mu_0) - d + \underbrace{\text{this term is}}_{\log \left(\frac{\det \Sigma_1}{\det \Sigma_0} \right)} + \text{key to optimization.}) \quad \text{--- (*)}$$

$$\mathcal{N}(x, \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} (\det \Sigma)^{1/2}} \exp \left[-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right]$$

- DDPM makes use of KL divergence.
We can use constants for covariance matrix.
 - ① Parameter of the statistical model.
 - ② Parameter of the learning model.
- $K(\theta)$ contains average w.r.t $q(x^{(0)})$.

$$K(\theta) = \int q(x^{(0..T)}) \left[\log p_\theta(x^{(T)}) + \log \frac{\pi p_\theta(x^{(t-1)} | x^{(t)}) \cdot q(x^{(t)} | x^{(0)})}{\pi q(x^{(t-1)} | x^{(t)}, x^{(0)}) q(x^{(t)} | x^{(0)})} \right] dx^{(0..T)}$$

$$K(\theta) = - \sum_{t=2}^T \int q(x^{(0)}, x^{(t)}) \cdot KL \left(q(x^{(t-1)} | x^{(t)}, x^{(0)}) \parallel p_\theta(x^{(t-1)} | x^{(t)}) \right) dx^{(0)} dx^{(t)}$$

$$+ H_q(x^{(T)} | x^{(0)}) - H_q(x^{(1)} | x^{(0)}) - H_p(x^{(T)})$$

— ④

Note: all terms in ④ can be evaluated exactly under the following conditions :

$$\bullet q(x^{(t)} | x^{(t-1)}) \sim \mathcal{N}(x^{(t)}, \sqrt{1-\beta_t} x^{(t-1)}, \beta_t I) \quad \text{— ⑤}$$

β_t is deterministic schedule for the forward diffusion process.

β_t : deterministic hyperparameter.

$$\cdot p_\theta(x^{(t-1)} | x^{(t)}) \sim \mathcal{N}(x^{(t-1)}, f_\mu(x^{(t)}, t, \theta_\mu), \\ f_\Sigma(x^{(t)}, t, \theta_\Sigma)) \quad (6)$$

These functions can be CNN

$$\cdot q(x^{(t-1)} | x^{(t)}, x^{(0)}) \sim \mathcal{N}(x^{(t-1)}, \tilde{\mu}(x^{(t)}, x^{(0)}), \tilde{\Sigma}_t I) \quad (7)$$

Claim is that this is also gaussian.

$$\text{Key term: } E_q[\|\tilde{\mu}(\cdot) - f_0(\cdot)\|^2]$$

Empirical observation: Covariance should be deterministic.

This comes from the KL divergence between two gaussians.

22/01/2025
Monday

Computer Vision

Recap:

- Maximizing log likelihood: $\mathcal{L}(\theta) = \mathbb{E}_q [\log P_\theta(x^{(0)})]$ — ①

Backward process:

- $P_\theta(x^{(0)})$ evaluated as a marginal:

$$P_\theta(x^{(0)}) = \int P_\theta(x^{(0..T)}) dx^{(1..T)} \quad \text{— ②}$$

- Introduce forward process:

$$P_\theta(x^{(0)}) = \int P_\theta(x^{(T)}) \cdot q(x^{(1..T)} | x^{(0)}) \cdot \frac{\prod P_\theta(x^{(t-1)} | x^{(t)})}{\prod q(x^{(t)} | x^{(t-1)})} dx^{(1..T)}$$

(due to Markov Chain Assumption). — ③

Plug ③ into ①:

$$\mathcal{L}(\theta) = \int q(x^{(0)}) \left[\log \int q(x^{(1..T)} | x^{(0)}) \cdot P_\theta(x^{(T)}) \cdot \frac{\prod P_\theta(x^{(t-1)} | x^{(t)})}{\prod q(x^{(t)} | x^{(t-1)})} \cdot dx^{(1..T)} \right] dx^{(0)} \quad \text{— ④}$$

Lower bound $\mathcal{L}(\theta)$ in terms of $K(\theta)$:

$$\mathcal{L}(\theta) \geq K(\theta).$$

$$K(\theta) = \int q(x^{(0..T)}) \log \left[P_\theta(x^{(T)}) \cdot \frac{\prod P_\theta(x^{(t-1)} | x^{(t)})}{\prod q(x^{(t-1)} | x^{(t)})} \cdot \frac{q(x^{(t-1)} | x^{(0)})}{q(x^{(t)} | x^{(0)})} \right] dx^{(0..T)} \quad \text{— ⑤}$$

After a bit of work we get

$$K(\theta) = - \sum_{t=2}^T \int KL(q(x^{(t-1)} | x^{(t)}, x^{(0)}) || P_\theta(x^{(t-1)} | x^{(t)})) \\ q(x^{(0)}, x^{(t)}) dx^{(0)} dx^{(t)}$$

$$+ H_q(x^{(T)} | x^{(0)}) - H_q(x^{(1)} | x^{(0)}) + H_p(x^{(T)})$$

— ⑥

- Specifics : $q(x^{(t)}|x^{(t-1)}) = \mathcal{N}(x^{(t)}, \sqrt{1-\beta_t} \cdot x^{(t-1)}, \beta_t I)$ — (7)

 $P_0(x^{(t-1)}|x^{(t)}) = \mathcal{N}(x^{(t-1)}, f_u(x^{(t)}, t; \theta_\mu), f_\Sigma(x^{(t)}, t; \theta_\Sigma)).$

(This is a consequence of the diffusion process).

- β_t is assumed to be deterministic.
Covariance of the backward process is deterministic and depends on the covariance of the forward process.
- Further, the variance of the backward process is assumed to be deterministic as well. i.e.,

$$f_\Sigma(x^{(t)}, t; \theta_\Sigma) = \sigma_t^2 I$$

\Rightarrow We end up learning θ_μ .

- The lower bound $K(\theta) = E_q \left[\sum_{t=2}^T L_{t-1} \right] + E_q[L_T] - E_q[L_0]$

$$L_{t-1} = KL(q(x^{(t-1)}|x^{(t)}, x^{(0)})) \parallel P_0(x^{(t-1)}|x^{(t)})$$

$$L_{t-1} = \frac{1}{2\sigma_t^2} \| \tilde{\mu}_t(x^{(t)}, x^{(0)}) - f_u(x^{(t)}, t; \theta_\mu) \|^2 + c$$

constant — (8a)

Since the forward and backward trajectories follow (4) and (8) and

$$q(x^{(t-1)}|x^{(t)}, x^{(0)}) \approx \mathcal{N}(x^{(t-1)}, \tilde{\mu}_t(x^{(t)}, x^{(0)}), \tilde{\beta}_t I) — (9)$$

$$P_0(x^{(t-1)}|x^{(t)}) \approx \mathcal{N}(x^{(t-1)}, f_u(x^{(t)}, t; \theta_\mu), \sigma_t^2 I) — (10)$$

σ_t^2 is related to β_t .

$\tilde{\mu}_t$ estimating \rightarrow this is the mean of $x^{(t-1)} | x^{(t)}, f_t$

$\beta \sigma_t^2$ is related to β_t .

\rightarrow From estimating mean vector to estimating noise.

28 Jan 2025

Tuesday

Computer Vision

Recap:

$$\min_{\theta \in \Theta} E_q [\|\tilde{\mu}(x^{(0)}, x^{(t)}) - f_\mu(x^{(t)}, t; \theta)\|_2^2] \quad \text{--- (1)}$$

$$\hookrightarrow q(x^{(t-1)} | x^{(t)}, x^{(0)}) \sim N(x^{(t-1)}, \tilde{\mu}(x^{(t)}, x^{(0)}), \tilde{\beta}_t I)$$

$$\bullet \tilde{\mu}(x^{(0)}, x^{(t)}) = \frac{\sqrt{\bar{\alpha}_{t-1}}}{(1-\bar{\alpha}_t)} \cdot \beta_t x^{(0)} + \frac{\sqrt{\alpha_t} (1-\bar{\alpha}_{t-1})}{(1-\bar{\alpha}_t)} x^{(t)} \quad \text{--- (2)}$$

$$\bullet P_\theta(x^{(t-1)} | x^{(t)}) \sim N(x^{(t-1)}, f_\mu(x^{(t)}, t; \theta), \sigma_t^2 I) \quad \text{--- (3)}$$

$$q(x^{(t)} | x^{(t-1)}) \sim N(x^{(t)}, \sqrt{1-\beta_t} x^{(t-1)}, \beta_t I) \quad \text{--- (4)}$$

Models like UNet can be used to solve this.

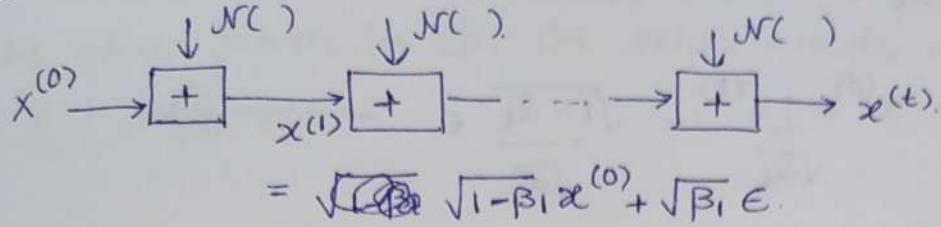
Denoising model:

$$q(x^{(t)} | x^{(0)}) \sim N(x^{(t)}; \sqrt{\bar{\alpha}_t} x^{(0)}, (1-\bar{\alpha}_t) I) \quad \text{--- (5)}$$

$$\text{or } x^{(t)} = \sqrt{1-\bar{\alpha}_t} \sqrt{\bar{\alpha}_t} x^{(0)} + \sqrt{1-\bar{\alpha}_t} \in \epsilon \sim N(0, I) \quad \text{--- (6)}$$

$$\text{where } \alpha_t = 1 - \beta_t \quad \bar{\alpha}_t = \prod_{s=1}^t \alpha_s \quad \text{--- (7)}$$

Can we derive at ⑤ from ④?



Q. Show how $x^{(t)}$ relates to $x^{(0)}$.

$x^{(0)}$ can be expressed in terms of $x^{(t)}$ and ϵ .

$$\begin{aligned} x^{(1)} &= \cancel{\text{IF } x^{(0)}} \sqrt{\alpha_1} x^{(0)} + \sqrt{\beta_1} \epsilon. \\ x^{(2)} &= \cancel{\text{IF } x^{(1)}} \sqrt{\alpha_2} x^{(1)} + \sqrt{\beta_2} \epsilon. \\ &= \sqrt{\alpha_2} (\sqrt{\alpha_1} x^{(0)} + \sqrt{\beta_1} \epsilon) + \sqrt{\beta_2} \epsilon. \\ &= \sqrt{\alpha_1 \alpha_2} x^{(0)} + (\sqrt{1-\beta_2}) \sqrt{\beta_1} \epsilon + \sqrt{\beta_2} \epsilon \\ &= \sqrt{\alpha_2} x^{(0)} + (\sqrt{\beta_1 - \beta_1 \beta_2} + \sqrt{\beta_2}) \epsilon. \end{aligned}$$

$$\begin{aligned} 1 + \beta_1 + \beta_2 - \beta_1 \beta_2 - 1 &\quad 1 - \alpha_1 \alpha_2 \quad 1 - \sqrt{1-\beta_1} \sqrt{1-\beta_2} \\ - 2 \sqrt{\beta_1 + \beta_2 - \beta_1 \beta_2 - 1} &\quad \cancel{1 - \sqrt{1-\beta_1} \sqrt{1-\beta_2}} \\ \sqrt{\alpha_2} \sqrt{\beta_1} \epsilon + \sqrt{1-\alpha_2} \epsilon &\quad (\sqrt{\alpha_2} \sqrt{1-\alpha_1} + \sqrt{1-\alpha_2}) \epsilon \\ = (\sqrt{\alpha_2 - \alpha_1 \alpha_2} + \sqrt{1-\alpha_2}) \epsilon & \\ \sqrt{\alpha_2 - \alpha_1 \alpha_2 + 1 - \alpha_2 + 2 \sqrt{\alpha_2 - \alpha_2^2 - \alpha_1 \alpha_2 + \alpha_1 \alpha_2^2}} \epsilon. & \end{aligned}$$

$$\sqrt{1 - \alpha_1 \alpha_2 + 2 \sqrt{\alpha_2 - \alpha_1 \alpha_2 (\alpha_2 - 1) - \alpha_2 (\alpha_2 - 1)}}.$$

$$\sqrt{1 - \alpha_1 \alpha_2 + 2 \sqrt{(\alpha_2 - 1) \alpha_2 (\alpha_2 - 1)}}.$$

$$\sqrt{1 - \alpha_1 \alpha_2 + 2 \sqrt{\alpha_2 (1 - \alpha_2)}}.$$

$$= \sqrt{1 + 2\sqrt{\alpha_2} - \alpha_2 (\alpha_1 + 2\sqrt{\alpha_2})}$$

$$x^{(t)} = \sqrt{\bar{\alpha}_t} x^{(0)} + \sqrt{1-\bar{\alpha}_t} e \quad e \sim \mathcal{N}(0,1)$$

From ⑥, $x^{(0)} = \frac{1}{\sqrt{\bar{\alpha}_t}} x^{(t)} - \frac{\sqrt{1-\bar{\alpha}_t}}{\sqrt{\bar{\alpha}_t}} e \quad \text{--- ⑧}$

Now we will estimate e .

e is standard gaussian.

Plug ⑧ into ②

We want everything in terms of $x^{(t)}$.

$$\sqrt{\bar{\alpha}_t} = \sqrt{\alpha_1 \dots \alpha_t} ; \quad \alpha_t = 1 - \beta_t$$

$$\sqrt{\bar{\alpha}_t} = \sqrt{\alpha_t \bar{\alpha}_{t-1}}$$

Writing $\tilde{\mu}(x^{(t)}, x^{(0)})$ again,

$$\tilde{\mu}(x^{(t)}, x^{(0)}) = \frac{\sqrt{\bar{\alpha}_{t-1}}}{(1-\bar{\alpha}_t)} \cancel{\beta_t} \cdot x^{(0)} + \frac{\sqrt{\bar{\alpha}_t} (1-\bar{\alpha}_{t-1})}{(1-\bar{\alpha}_t)} x^{(t)}$$

From ②,

$$\begin{aligned} \tilde{\mu}(x^{(t)}, x^{(0)}) &= \frac{\sqrt{\bar{\alpha}_{t-1}}}{(1-\bar{\alpha}_t)} \beta_t \left(\frac{x^{(t)}}{\sqrt{\bar{\alpha}_t}} - \frac{\sqrt{1-\bar{\alpha}_t} e}{\sqrt{\bar{\alpha}_t}} \right) + \\ &= \cancel{\beta_t} \cdot \frac{\sqrt{\bar{\alpha}_t} (1-\bar{\alpha}_{t-1})}{(1-\bar{\alpha}_t)} x^{(t)} \end{aligned}$$

$$= \frac{1}{(1-\bar{\alpha}_t)\sqrt{\bar{\alpha}_t}} (\beta_t \cancel{x^{(t)}} + \alpha_t (1-\bar{\alpha}_{t-1})) x^{(t)} - \frac{\beta_t e}{\sqrt{\bar{\alpha}_t} \sqrt{1-\bar{\alpha}_t}}$$

$$= \frac{1}{(1-\bar{\alpha}_t)\sqrt{\bar{\alpha}_t}} (1 + \alpha_t \bar{\alpha}_{t-1}) x^{(t)} - \frac{\beta_t e}{\sqrt{\bar{\alpha}_t} \sqrt{1-\bar{\alpha}_t}}$$

$$= \frac{1}{\sqrt{\bar{\alpha}_t}} \left(x^{(t)} - \frac{\beta_t e}{\sqrt{1-\bar{\alpha}_t}} \right) \quad \text{--- ⑨}$$

From ⑨, our $f_{\mu}(\cdot)$ effectively learns to predict the noise term in ⑨. In other words we have $f_{\epsilon}(x^{(t)}, \epsilon, t)$.

If $f_{\mu}(x^{(t)}, x^{(0)}) = \frac{1}{\sqrt{\alpha_t}} (x^{(t)} - \frac{\beta_t}{\sqrt{1-\alpha_t}} f_{\epsilon}(\cdot))$, then

⑩ becomes,

$$\mathbb{E}_{\epsilon} \left[\frac{\beta_t^2}{\alpha_t(1-\alpha_t)} \| \epsilon - f_{\epsilon}(x^{(t)}, \epsilon, t) \|^2 \right] \longrightarrow ⑩$$

3rd Feb 2025
Monday

Computer Vision

Introduction to GANs

- Difference is the maximum likelihood approach that we have used in other models.
- Learn parameters of Discriminator and Generator.
- Image classifier and Image generator.
- Discriminator accepts samples that are outputs of other models and raw data ($P(x)$).
- Gradient ascent and Gradient descent
 - discriminator
 - generator.

- Implicit model
 - Game Theoretic Formulation (unlike ML formulation)
 - Generator vs Discriminator.
- Discriminator function:

$$D(x; \theta_D) = \begin{cases} 1, & \text{if } x \in \text{Real Sample Set } (P_{\text{data}}(x)) \\ 0, & \text{if } x \in \text{Generated Sample set } (\cancel{\text{or class}}} \quad (P_{\text{model}}(x)) \end{cases}$$

- The goal of the discriminator is to discriminate (or classify) real and generated samples as accurately as possible.

$$J(\theta_D, \theta_E) = E_{x \sim P_{\text{data}}} [\log D(x; \theta_D)] +$$

$$E_{x \sim P_{\text{model}}} [\log (1 - D(x; \theta_D))] \quad \text{--- ①}$$

- Find θ_D^* that maximizes the above cost in ①.
- Let us now bring the generator into play.
- $x \sim P_{\text{model}}(x)$

$\Rightarrow x = G(z; \theta_G)$ where $G(z; \theta_G)$ is the generator network with parameters θ_G .

Let's plug ② into ①.

$$J_D(\theta_G, \theta_D) = \mathbb{E}_{\substack{x \sim P_{\text{data}} \\ z \sim P_z}} [\log [D(x; \theta_D)]] + \mathbb{E}_{\substack{x \sim P_{\text{data}} \\ z \sim P_z}} [\log (1 - D(G(z; \theta_G); \theta_D))] \quad \text{--- ③}$$

• Let's now solve the problem.

• Let's define

$$J_G(\theta_D, \theta_G) = -J_D(\theta_G, \theta_D) \quad (\text{generator's loss}) \quad \text{--- ④}$$

• The discriminator and the generator play a game which is defined as:

$$(\theta_G^*, \theta_D^*) = \arg \min_G \max_D V(\theta_D, \theta_G).$$

$$= \min_G \max_D \mathbb{E}_{\substack{x \sim P_{\text{data}}}} [\log [D(x; \theta_D)]] + \cancel{\mathbb{E}_{\substack{z \sim P_z}} [\log]} \\ + \mathbb{E}_{\substack{z \sim P_z}} [\log (1 - D(G(z; \theta_G); \theta_D))]$$

Initialize both θ_G, θ_D to random values and then let the discriminator iteratively update the values.

4th Feb'25
Tuesday

Computer Vision

Recap: $J_D(\theta_D, \theta_G) = \mathbb{E}_{\substack{x \sim P_{\text{data}}}} [\log D(x; \theta_D)] + \mathbb{E}_{\substack{x \sim P_{\text{model}}}} [\log (1 - D(x; \theta_D))]$

Minimization wrt generator parameters.

Maximization wrt discriminator parameters.

Q. What model could be better for generator?

A. DeConv Net.

Q. Pass Noise Distribution through a Deconv Net.

Q. what could be better for a discriminator?

A. Any image classifier.

Q. Do we have a guarantee that the generator
~~and~~ will ultimately be able to learn P_{data} ?

A. Yes. But under certain conditions.

Q. When is the generator minimize the best discriminator model?

A. When $P_{\text{model}} = P_{\text{data}}$. Equal to $\frac{1}{2}$.

Q. What are the optimal solutions for $D(x; \theta_D)$ and $G(z, \theta_G)$?

Assumption: The model has high capacity.

Discriminator: D^* the optimal discriminator for a given generator G is

$$D^*(x) = \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_{\text{model}}(x)} \quad \text{--- (1)}$$

Generator: G^* The optimal generator gives

$$P_{\text{model}} = P_{\text{data}} \quad \text{--- (2)}$$

$$\begin{aligned} J(\theta_D) &= \int P_{\text{data}}(x) \log [D(x)] dx + \\ &\quad \int P_z(z) \log [1 - D(G(z))] dz. \\ &= \int P_{\text{data}}(x) \log [D(x)] dx + \int P_{\text{model}}(x) \log [1 - D(x)] dx. \\ &= \int P_{\text{data}}(x) \log D(x) + P_{\text{model}}(x) \log (1 - D(x)) dx. \end{aligned}$$

This is of the form $a \log y + b \log (1-y)$.

What can you tell about 'a' and 'b'?

Non-negative and bounded.

When $y = \frac{a}{a+b}$, it will be optimum.

$$D_G^*(x) = \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_{\text{model}}(x)} \quad \text{--- (1)}$$

Let's now find the optimal G for D^*

$$\begin{aligned} V(G, D^*) &= C(G) = \int P_{\text{data}}(x) \cdot \log \left[\frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_{\text{model}}(x)} \right] \\ &\quad + P_{\text{model}}(x) \cdot \log \left[\frac{P_{\text{model}}(x)}{P_{\text{data}}(x) + P_{\text{model}}(x)} \right] dx. \end{aligned}$$

$$\text{Recall } JS(p \parallel q) = \frac{1}{2} [KL(p \parallel m) + KL(q \parallel m)]$$

$$m = \frac{p+q}{2}$$

$$C(A) = KL \left[P_{\text{data}}(x) \parallel \left(\frac{P_{\text{data}}(x) + P_{\text{model}}(x)}{2} \right) \right] - \log_2 \\ + KL \left[P_{\text{model}}(x) \parallel \left(\frac{P_{\text{data}}(x) + P_{\text{model}}(x)}{2} \right) \right] - \log_2$$

$$\boxed{C(G) = 2 JS(P_{\text{data}}(x) \parallel P_{\text{model}}(x)) - \log 4} \quad \text{--- (2)}$$

$\Rightarrow C(A)$ takes a min-value of $-\log 4$ when
 $P_{\text{model}}(x) = P_{\text{data}}(x)$, i.e., the optimal generator
 G^* learns the data distribution exactly.

Equations ① and ② give answers to our questions.

10/2/2025
Monday

Computer Vision

Metrics for generative models

→ We want to know the quality of images generated by models such as GAN.

Inception Score :

How to calculate?

There is a model Inception V3, takes as input the images we get from GANs (sample them).

Calculate score based on this output.

$$IS = \exp \left[E_{x \sim P_{\text{model}}} \left[KL(P(y|x) \| P(y)) \right] \right]$$

sharpness diversity
or faithfulness

$P(y)$ → R.V y must have high entropy.

$\boxed{\text{diversity}}$ (We want outputs of model to be as diverse as possible, means distributed over all labels).

$P(y|x)$ → R.V $P(y|x)$ must have low entropy.

$\boxed{\text{sharpness or faithfulness}}$ (We want the model to output single output label for given x).

It tells how sharp the image is.

→ We want KL divergence high, means distance between them is high, so that we get images that are diverse across all labels and each image has unique label.

In general,

- ① How realistic is the image? (depends on classifier performance).
- ② Generate samples from variety of classes.

- Ideally $p(y|x)$ should be far from a uniform distribution.
- $P(y)$ must be close to uniform distribution.
- $P(y|x)$ is found from the output of the classifier (Inception v3).

$\textcircled{B} \quad P(y) = \int P(y|x) \cdot P(x) dx \quad (\text{marginal})$

We can approximate it by samples.

$$P(y) \approx \frac{1}{N} \sum_{i=1}^N P(y|x_i) \quad (\text{Approximation})$$

Drawbacks:

We can't ~~train~~ expect perfect accurate results when data like medical data is used.

Why? ~~Biased~~ Becoz biased towards Image-Net like data.

Can't be reliable for non-Image Net data.

Frechet Inception Distance (FID):

- You don't take single image. Sample multiple images.
- Compare dist of generated samples with dist of real samples.
 - Compare mean and covariance of real samples with that of generated ones.

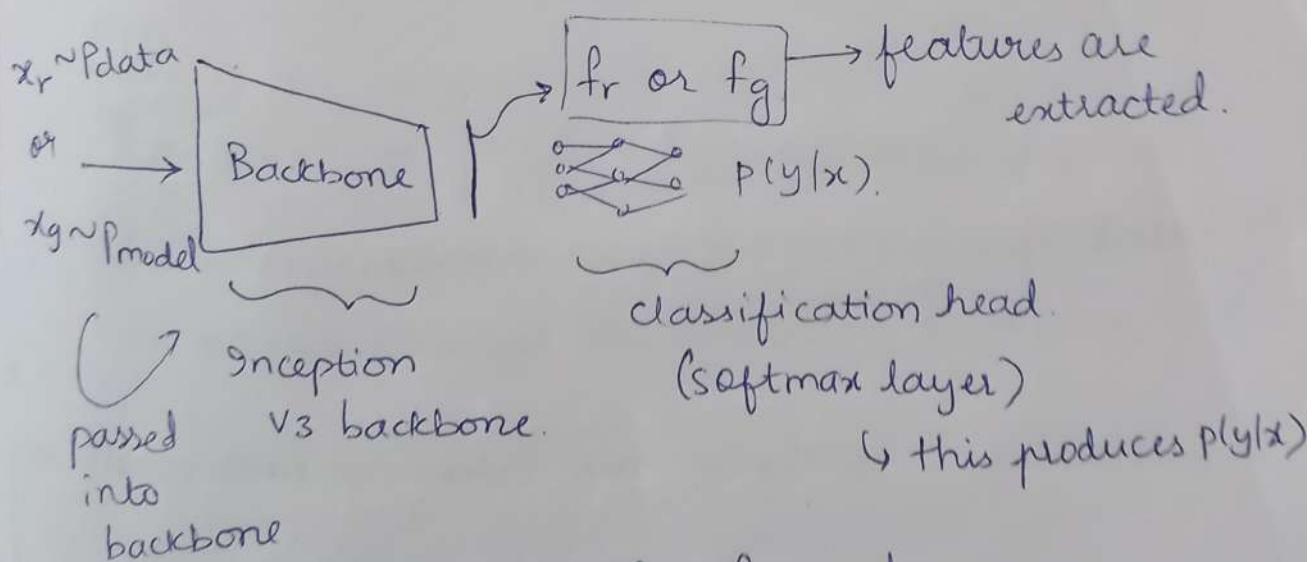
$$FID = \|\mu_r - \mu_g\|_2^2 + \text{tr}(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{-1})$$

→ We want this value to be low.

μ_r : Mean vector of the Inception vector of real samples.

μ_g : Mean vector of the Inception vector of generated samples.

Σ_r, Σ_g : covariance matrices of real and generated samples respectively.



We want f_r, f_g and compare distributions.

The more the samples, the better.

Ideally, a generative model evaluation metric must:

- ① Measure sharpness
- ② Measure diversity
- ③ Measure "realness" / "naturalness".

Naturalness Image Quality Evaluator.

11/02/2025
Tuesday

Computer Vision

Recap:

- ① How real the image is?
- ② How diverse labels are?

Drawbacks of Inception Score:

- 1) Non-ImageNet Type data has problems.

Why this could be issue?

See cond. dist $p(y|x)$ of model.

If y and x come from data it hasn't learnt, then it could be problem.

- 2) Doesn't look for mode collapse.

Mode collapse means model generates not so diverse images.

This model can fool classifier and give high IS. How? Models have few variations in same image.

Here $p(y|x)$ has high entropy and $p(x)$ has less entropy

FID

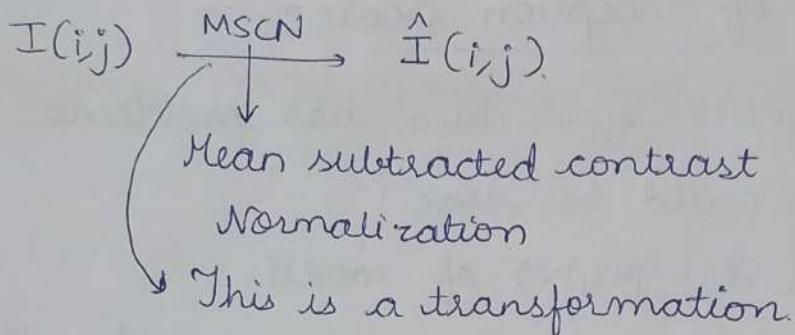
$$\|\mu_r - \mu_g\|_2^2 + \text{tr}(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{1/2})$$

↓
Mean representation we get when real data is passed through the inception backbone.

→ Downstream model we considered is Inception v3.

Naturalness Image Quality Estimator (NIQE)

$I(i,j)$ → Natural image.
 grayscale (captured from camera)
 pixel intensity at spatial location (i,j) .



$$\hat{I}(i,j) = \frac{I(i,j) - \mu(i,j)}{\sigma(i,j) + c}$$

local mean around that square.

local S.D.



why c ?

To prevent division with '0'.

$$\mu(i,j) = \frac{1}{N} \sum_{m,n \in N(i,j)} I(i-m, j-n)$$

Neighbourhood.

$$\sigma^2(i,j) = \frac{1}{N} \sum_{m,n \in N(i,j)} (I(i-m, j-n) - \mu(i,j))^2$$

This is called contrast.

Difference between $I(i,j)$ and $\mu(i,j)$ locally.

→ In natural images, the histogram (~~dist~~
(distribution of \hat{I}) turns out to be gaussian.
It is unimodal (single peak).

So, for natural images, we have more handle
on the statistical patterns of the normalized
 \hat{I} .

This can be used as a signature to characterize
natural images.

→ If we take synthetic images (generated from
models), then the histogram is multi-modal
(many peaks).

Unimodal is true only in case of natural
images.

→ \hat{I} is unimodal, and can be modelled using
either a Gaussian or a Generalized Gaussian
Distribution (GED).

Q. Can we characterize all natural images?

FID score for real images vs synthetic images.

MSCN Histogram

→ When blurring (remove many details), so
pixel values become more similar to the
local ~~average~~ mean

So histogram becomes narrower because
there is less variation.

→ When we add noise, pixel values become random.
Std deviation increases. So histogram becomes wider. Spread in variance of image.

→ If we characterize the statistics of real undistorted images (Gaussian), we can build a model so that when new image comes we can detect whether it is fake or not based on these statistics.

→ (α, β) uniquely identify gaussian.

↳ shape of GGD (Generalized Gauss. Dist.).

Natural images have specific range of α, β .
So we can train a model to check if new image's (α, β) fall into that range.

→ γ is shape parameter (~~"small~~
^{tells} (sharp or flat) & distribution
 B_L, B_R (left and right tails of dist)).

↳ tells how extreme pixel values behave.

14/02/2025
Friday

Computer Vision

Image Net has fixed no of classes.

where Consequence: Generalization.

Observation: Language guidance will help.

→ ResNet, ViT: extract ~~for~~ image features.

Transformers: Extract language (text).
(BERT, GPT)

→ Model should work like, take as input (image, text) pair and give latent representation.
Maps both to common latent space.

Contrastive language Image Pretraining (CLIP)

→ No finetuning is required. It learns generalizable representation.

→ Zero-shot learning: No need of training data for new things.
→ DIDN'T UNDERSTAND.

→ Fixed set of categories is limiting.

→ Learning from raw text about images is promising.

→ Pretraining with (text, image) pairs.

→ Info Noise Contrastive Estimation (info NCE) loss.

Pulls similar pairs closer. in latent spaces
and pushes unrelated ones far apart.

→ Enables zero-shot learning transfer to downstream

tasks.

- Natural language based.
- Cosine similarity + softmax.
 - closeness
 - predict most relevant one

Info NCE loss:

$$L = - \sum_{i=1}^N \log \left(\frac{\exp(s(I_i, T_i)/\tau)}{\sum_{j=1}^N \exp(s(I_i, T_j)/\tau)} \right)$$

temperature parameter.

softmax

controls sharpness of dist.

where I_i is the i^{th} image representation

and T_i is the i^{th} caption representation

$$s(x, y) = \frac{\langle x, y \rangle}{\|x\| \cdot \|y\|}$$

17 Feb '25
Monday

Computer Vision

→ Contrastive pretraining plays major role in unsupervised learning.

CLIP helps in :

- ① Zero-shot learning : No need for fine-tuning
- ② Better generalization.
- ③ Learned features can be used in downstream tasks.

Q. Why same dimension for image and text ?

- Cosine similarity.
- Uniform Representation space.

• Info NCE loss.

- Representation learning in unsupervised learning.
- The role of mutual Information (MI).
 - In contrastive learning, positive pairs \Rightarrow higher MI.
 - Info NCE optimizes MI (minimize negative pair similarity).
- Info NCE Loss Formulation and Negative Sampling.
- Relation to MI.

Info NCE loss:

$$L = - \sum_{i=1}^N \log \left[\frac{\exp(S(I_i, T_i)/\tau)}{\sum_{j=1}^N \exp(S(I_i, T_j)/\tau)} \right]$$

Here we took I_{image} as reference.

Similarly we can have L_{Text} as reference.

What is Info NCE?

↓
Noise → Contrastive

Contrasts positive and negative samples.

Q. Can we do prediction in latent space?

A. Yes. Zero-shot transfer allowed to new tasks
without additional training. →

→ Autoencoder: Deterministic setting.

Maps input data to latent rep. deterministically.

→ Variational Autoencoder: Stochastic setting.

→ CLIP solves a semisupervised setting to learn representations.

Doesn't require labelled data. Can learn from pairs.

→ Goal: Learn latent representations.

Advantage: ① Can work with lower dimension.
② Better generalization.
③ Useful for downstream tasks.

→ Optimizing Info NCE loss is equivalent to maximizing a lower bound for M.I.

→ ~~We want~~ We want Info NCE loss to be low.

Predictive Coding in Latent Space.

→ Helps to model dependencies between past and future representations.

Goal: Maximize Mutual information (MI) between the source (input) and the context (history).

→ If MI is high, context helps in predicting future.

→ We use two-model in Predictive Coding:

① Learn latent representation from input sample (Eg: Speech).

② Auto-regressive model for History.

Remembers the past context and predicts future representations by maximizing MI b/w past and future latent.

Q. What is Mutual Information (M.I)?

If X, Y are R.V, M.I(X, Y) tells how dependent is one R.V on another.

We quantify this using entropy.

$$MI(X, Y) = H(X) - \underbrace{H(X|Y)}$$

Entropy of X : No. of bits to represent X without knowing Y .

If X, Y are independent $MI = 0$.

Mutual Information (M.I):

$$I(X:Y) = KL(P(X,Y) \parallel P(X)P(Y)) \rightarrow p(x,y) \log \frac{p(x,y)}{p(x)p(y)}$$
$$= E_{(x,y) \sim P(X,Y)} \left[\log \left(\frac{P(X|Y)}{P(X)} \right) \right] \quad p(x,y) \log \frac{p(x|y)p(y)}{p(x)p(y)}$$

$$= \int p(x,y) \cdot \log \left[\frac{P(X|Y)}{P(X)} \right] dx \cdot dy.$$

→ We want Z_{t+k} (latent representation) at time ' $t+k$ ' such that M.I between the original speech input X_{t+k} and context representation c_t is maximized.

↓
Summarizes past information.

→ In CPC, we want to learn the latent representation Z_{t+k} at time ' $t+k$ ' such that $I(X_{t+k}; c_t)$ is maximized.

- Find a function $f_k(x_{t+k}, c_t)$ such that $I(x_{t+k}, c_t)$ is maximized.

let $f_k(x_{t+k}; c_t) = \exp(Z_{t+k}^T \cdot W_k \cdot c_t)$ be one candidate function.

18/02/2025
Tuesday

Computer Vision

Auto regressive model helps in looking at long back history and not just the local spatial window.

Applications with the encoder we have learnt:
helps in classifying.

Applications from CLIP: Zero-shot learning.

$$M.I.: I(X;Y) = KL(P(x,y) \parallel p(x).p(y))$$

$$I(X;Y) = \int p(x,y) \cdot \log \frac{p(x|y)}{p(x)} dx dy$$

We want to learn Z_{t+k} such that $I(Z_{t+k}; c_t)$ is maximized.
 \Downarrow
d.r at future time instant.

$$I(Z_{t+k}; c_t) = \sum_{x_{t+k}, c_t} p(x_{t+k}, c_t) \underbrace{\log \frac{p(x_{t+k} | c_t)}{p(x_{t+k})}}$$

we are interested in this term.

- Specifically can we learn a function $f_k(x_{t+k}, c_t)$ that is proportional to $\frac{p(x_{t+k} | c_t)}{p(x_{t+k})}$

\Downarrow
what is that range of this?
[0, ∞).

- Note that $\frac{p(x_{t+k} | c_t)}{p(x_{t+k})}$ is non-negative.

→ Encoder helps with latent rep.

CNN can be used for images and

Transformer can be used for text/speech.

- A candidate function for approximating $\frac{P(x_{t+k} | c_t)}{P(x_{t+k})}$

is $f_k(x_{t+k}, c_t) = \exp(Z_{t+k}^T \cdot W_k \cdot c_t)$

↑
This func. is non-negative.

- Info NCE loss func. is

$$L = -\frac{1}{N} \sum_{k=1}^N \log \left(\frac{f_k}{\sum_{j=1}^N f_j} \right)$$

19th Feb '25
Wednesday

Computer Vision

Recap:

$$L = -E_x \log \left(\frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_j(x_j, c_t)} \right) \quad \textcircled{1}$$

↓
average

- We want to maximally preserve $I(x_{t+k}; c_t)$

$$= \sum_{x_{t+k}, c_t} p(x_{t+k}, c_t) \log \left(\frac{p(x_{t+k}|c_t)}{p(x_{t+k})} \right)$$

↓ raw input ↓ context

- $f_k(x_{t+k}, c_t) = \exp(z_{t+k}^T w_k(c_t))$ → a candidate model we use.
 encoder ↓
 auto regressive model.

to predict

$$\frac{p(x_{t+k}|c_t)}{p(x_{t+k})}$$

This comes from MI preservation condition.

→ We want to learn the encoder and autoregressive model together.

- For us to jointly learn the encoder and the autoregressive model, we use the info NCE loss in $\textcircled{1}$.

$$L = -E_x \left[\log \left[\frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_j(x_j, c_t)} \right] \right]$$

positive pair
 ↓ all pairs of raw inputs and c_t

X is the set of data samples, i.e.,

$$X = \{x_1, x_2, \dots, x_N\}$$

In this set there is one positive sample (correct one) and $N-1$ negative samples.

- Let's assume that we find our optimal encoder and the auto regressor.

$$f_k(x_{t+k}, c_t) = \frac{p(x_{t+k} | c_t)}{p(x_{t+k})} \quad \text{--- (4)}$$

Recap:

$$I(x_{t+k}; c_t) = \sum_{(x_{t+k}, c_t)} p(x_{t+k}, c_t) \cdot \log \frac{p(x_{t+k} | c_t)}{p(x_{t+k})}$$

Plug ④ into ③

$$L_{\text{opt}} = -E_x \log \left[\frac{\frac{p(x_{t+k} | c_t)}{p(x_{t+k})}}{\frac{p(x_{t+k} | c_t)}{p(x_{t+k})} + \sum_{j \neq t+k} \frac{p(x_j | c_t)}{p(x_j)}} \right]$$

$$\star \approx \textcircled{2} E_x \log \left[1 + \frac{p(x_{t+k})}{p(x_{t+k} | c_t)} \cdot (N-1) E_{x_j} \frac{p(x_j | c_t)}{p(x_j)} \right]$$

$$\left(E_{x_j} \frac{p(x_j | c_t)}{p(x_j)} = \int p(x_j) \cdot \frac{p(x_j | c_t)}{p(x_j)} dx_j \right) \\ = 1$$

$$\textcircled{2} = E_x \log \left[1 + \frac{p(x_{t+k})}{p(x_{t+k} | c_t)} \cdot (N-1) \right]$$

$$\leq E_x \log \left[\frac{N \cdot p(x_{t+k})}{p(x_{t+k} | c_t)} \right] \quad \overbrace{\quad}^{I(x_{t+k}; c_t)} \\ = \log N - E_x \log \left[\frac{p(x_{t+k} | c_t)}{p(x_{t+k})} \right]$$

$$L_{\text{opt}} \geq \log N - I(x_{t+k}; c_t)$$

$$\boxed{I(x_{t+k}; c_t) \geq \log N - L_{\text{opt}}} \quad \text{--- (6)}$$

24th Feb '2025
Monday

Computer Vision

Recap:

Downstream task of CLIP (Pretraining) is zero-shot transfer.

- Contrastive learning helps to leverage large scale dataset without explicit annotation or labelling.
- CLIP is an example model.
- Contrastive learning : pulls +ve pairs together.
pushes -ve pairs apart.
- Info NCE loss is a "good" function that encourages contrastive learning.

Noise-Contrastive Estimation (NCE)

We have $X = \{X_1, X_2, \dots, X_N\}$, drawn from unknown data distribution. $P_{\text{data}}(\cdot)$ A dataset of N points.

→ We assume $P_{\text{data}}(\cdot) \in \mathbb{P}\{P_{\text{model}}(\cdot; \alpha)\}_{\alpha}$.
(family of priors).

→ For some optimal $\hat{\alpha}$, we have,

$\exists \hat{\alpha}$, such that $P_{\text{data}}(\cdot) = P_{\text{model}}(\cdot; \hat{\alpha})$.

and $\int \underbrace{P(x; \hat{\alpha}) dx}_\text{constraint} = 1$

→ Any model,

$$P_{\text{model}}(\cdot; \alpha) = \frac{1}{Z(\alpha)} P^{\circ}_{\text{model}}(\cdot; \alpha)$$

where P°_{model} is the $P^{\circ}_{\text{model}}(\cdot; \alpha)$ is an unnormalized function.

$Z(\alpha)$ is the normalization constant.

$$Z(\alpha) = \int P^{\circ}_{\text{model}}(x; \alpha) \cdot dx. \quad \text{--- (1)}$$

→ Finding $Z(\alpha)$ in practice is usually intractable analytically.

becoz it is integral over all possible values of x , computationally expensive or intractable.

→ NCE solves the problem by:

NCE looks at this as a classification problem.

Compares it against a known noise distribution.

It learns $Z(\alpha)$ as a parameter.

It basically differentiates b/w real samples $\sim P_{\text{data}}(x)$ and noise samples $\sim P_{\text{noise}}(x)$ and in the process of doing so, learns finds $P^{\circ}_{\text{model}}(\cdot; \alpha)$.

using

— ~~freezing~~ $P^{\circ}_{\text{model}}(\cdot; \alpha)$

— making $Z(\alpha)$ a learnable parameter.
 $\approx C$

• NCE learns $\theta = \{\alpha, C\}$ in a "noise contrastive" approach.

• In other words, we learn the data dist in a relative or contrastive fashion w.r.t known noise dist. $P_x(\cdot)$.

25/02/25

Computer Vision

- Recaps NCE
- NCE loss
- Connection to supervised learning.
- Recap:
 - $\mathbf{x} = (x_1, x_2, \dots, x_N)$ drawn from $p_{\text{data}}(\cdot)$.
 - $p_{\text{data}}(\cdot) \in \{p_{\text{model}}(\cdot; \alpha)\}_{\alpha}$.
 - Find α that "best fits" the data points
 - We are interested in solving the problem of estimating unnormalized functions $p_{\text{model}}^0(\cdot; \alpha)$ due to tractability issues
 - Specifically finding $Z(\alpha) = \int p_{\text{model}}^0(x; \alpha) dx$ is difficult.
 - If $p_{\text{model}}^0(x; \alpha) \sim \mathcal{N}(x_0, \sigma^2)$, what is $Z(\alpha)$? $Z(\alpha) = \sqrt{2\pi}\sigma$
 - In general, find $Z(\alpha)$ is non-trivial.
 - NCE explores a way to estimate both α and $Z(\alpha)$ by solving a single optimization.
 - NCE accomplishes this by estimating $\Theta = \{\alpha, c\}$ by comparing or contrasting it with an artificial noise distribution from which we draw N samples $\mathbf{y} = (y_1, y_2, \dots, y_N)$.
 - Noise-contrastive estimation works by solving the following optimization. \rightarrow

$$\Theta = \underset{\Theta}{\operatorname{argmax}} J(\Theta), \text{ where}$$

$$J(\Theta) = \frac{1}{2N} \sum_{i=1}^N \log [h(x_i; \Theta)] + \log [1 - h(y_i; \Theta)] \quad \text{---(2)}$$

$$h(x; \Theta) = \frac{1}{1 + \exp(-f(x; \Theta))}; f(x; \Theta) = \ln(p_{\text{model}}^0(x; \Theta)) - \ln(p_n(x)).$$

$$\Rightarrow h(x; \theta) = \frac{1}{1 + \exp \left(\ln \frac{p_n(x)}{p_{\text{model}}^0(x; \theta)} \right)}$$

$$h(x; \theta) = \frac{p_{\text{model}}^0(x; \theta)}{p_{\text{model}}^0(x; \theta) + p_n(x)}$$

Note: $\ln(p_{\text{model}}(\cdot; \alpha)) = \ln c + \ln(p_{\text{model}}^0(\cdot; \alpha))$.

Connection to supervised learning

Let $U = (u_1, u_2, \dots, u_{2N})$ be the union of X and Y .

If $u_i \in X, c_i = 1$ Note $p(c=1) = p(c=0) = \frac{1}{2}$.

$u_i \in Y, c_i = 0$.

- $p(u|c=1; \theta) = p_{\text{data}}(u; \theta)$ or $p_{\text{model}}^0(u; \theta)$

- $p(u|c=0) = p_n(u)$.

What is $p(c=1|u; \theta) = \frac{p(u|c=1; \theta) \cdot p(c=1)}{p(u)}$

$$= \frac{\frac{1}{2} \cdot p_{\text{model}}^0(u; \theta)}{\frac{1}{2}(p_{\text{model}}^0(u; \theta) + p_n(u))}$$

$$= h(u; \theta)$$

$$p(c=0|u; \theta) = 1 - h(u; \theta)$$

Now apply the log loss to these estimates.

$$l(\theta) = \sum_{i=1}^{2N} (-\log(h(u_i; \theta))) + (-\log(1 - h(u_i; \theta))) \quad \text{--- (3)}$$

$$(-\log(1 - h(u_i; \theta)))$$

- o Recap: NCE loss formulation
- o Wrap up contrastive learning discussion
 - SimCLR

Reminders:

- o PPK due tonight
- o Midterm exam coming np.

Setup

- o $X = (x_1, x_2, \dots, x_N)$ drawn from $p_{\text{data}}(x)$.

- o $p_{\text{data}}(x) \in \{p_{\text{model}}(x; \alpha)\}_{\alpha}$, a prior family.

$$\circ p_{\text{model}}(x; \alpha) = \frac{1}{Z(\alpha)} \cdot p^{\circ}_{\text{model}}(x; \alpha)$$

- o We want to learn $Z(\alpha), \alpha$. ($\because Z(\alpha) = \int p^{\circ}_{\text{model}}(x; \alpha) dx$ is intractable)

- o We learn $\Theta = \{\alpha, c\}$ in a unsupervised fashion.

- o $Y = (y_1, y_2, \dots, y_N)$ drawn from $p_m(y)$.

$$\circ J(\Theta) = -\frac{1}{2N} \sum_i \log [h(x_i; \Theta)] + \log [1 - h(y_i; \Theta)]$$

$$h(x_i; \Theta) = \frac{p_{\text{model}}(x_i; \Theta)}{p_{\text{model}}(x_i; \Theta) + p_m(x_i)}$$

$$h(x; \Theta) = \frac{1}{1 + \exp(-f(x; \Theta))},$$

$$f(x; \Theta) = \ln p_{\text{model}}(x; \Theta) - \ln p_m(x). \quad [-\otimes]$$

- o Recall connection with supervised learning.

$$l(\theta) = \sum_i c_i \log \left(p \underbrace{\left(c_i=1 \mid u_i; \theta \right)}_{h(x_i; \theta)} \right) + (1-c_i) \log \left(p \left(c_i=0 \mid u_i; \theta \right) \right)$$

| . Class labelling
 $c_i=1$ if $u_i \in X$
 $c_i=0$ if $u_i \in Y$

$$l(\theta) = \sum_i \log h(u_i; \theta) + \log (1 - h(u_i; \theta))$$

$$U = X \cup Y$$

- Please read the Info NCE and NCE papers.

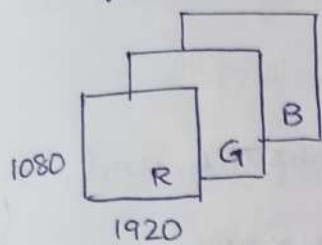
- ^{Hw:} Go back to the InfoNCE formulation and find

$p([d=i] \mid x, c_i)$ where $[d=i]$ is the event corresponding to the positive pair.

3/03/2025
Monday

Computer Vision

How a HD image is represented in memory?
(After encoding)



$1920 \times 1080 \times 3$,
 $w \quad h \quad c$

- There are 256 intensity level values. (0 to 255)
1 byte of data per pixel.
- $I(w, h, c)$ is a pixel intensity corresponding to spatial location (w, h) and channel c .
- Typically, 8 bits are used to represent a pixel's intensity for each colour channel.
 \Rightarrow Since 3 channels, 24 bits per pixel.
- Raw RGB frames:
 $1920 \times 1080 \times 3 \text{ bytes} \approx 6.22 \text{ MB per frame}$.
- Images Further, images (and videos) are compressed before storage.

4/03/2025

Tuesday

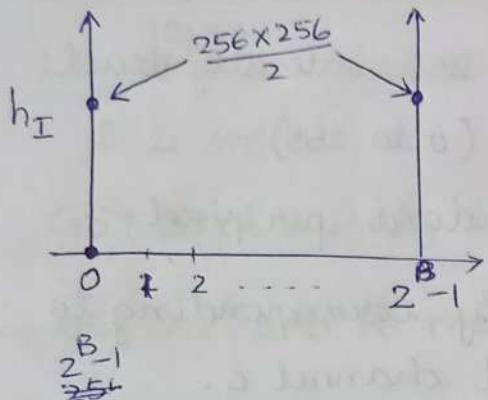
Computer Vision

Recap: Image formation

- Image histogram
- Point operations
 - linear
 - Non linear.
- Image histogram:
Given a gray scale image of size $M \times N$ that uses 8 bits per pixel I , the histogram of I

denoted h_I is the count of the intensity levels in I , i.e.

$h_I(k)$ denotes the count of intensity value k occurring in the image.



If I were 256×256 , and

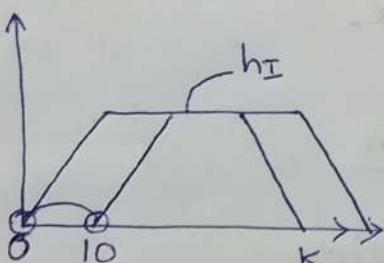
$$h_I(0) = \frac{256 \times 256}{2}$$

$$h_I(255) = \frac{256 \times 256}{2}$$

$$\sum_{k=0}^{2^B-1} h_I(k) = 256 \times 256 \text{ or generally } M \times N.$$

- $P_I(k) = \frac{h_I(k)}{M \times N}$ is the normalized histogram of I . It is a valid probability mass function. (PMF).

- Point operations:



$$\text{Ex: } J = I + k_0 = 10.$$

$$h_J(k) = h_I(k - k_0).$$

In general $J = P_I + Q$ (linear)

Transformed pixel intensity \rightarrow offset (brightness or shift)

contrast (scaling factor)

If $A = \min(I) > 0$

~~B~~ $C = \max(I) < 2^B - 1$

$B = \text{bits per channel}$.

Now we want to stretch the range to $[0, 2^B - 1]$. This is contrast stretching.

We solve for P and L:

$$O = P \cdot A + L$$

$$2^B - 1 = P \cdot C + L$$

- Non-linear operations :-
 - $J = (I)^r$ \log compression
 - $J = \log [1 + I]$
 - $J = \exp [f(I)]$

10/3/2025
Friday

Computer Vision

Recap : Histogram, point operations

- LTI Systems (Linear Time Invariant)
 - Convolution sum.
 - Low pass filter
 - High pass filter
 - Image denoising
- b bit image means 2^b intensity levels can be supported.

Linear filters:

- Operations where output pixel is a linear combination of its neighbouring input pixels.
- We apply the same operation at every pixel.
- We used this in CNN and attention.

- A system $T\{\cdot\}$ is a map from an input sequence $x[n]$ to a desired output sequence $y[n]$.

$$y[n] = T\{x[n]\}.$$

- Linear system:

A system $T\{\cdot\}$ is said to be linear if

$$y_1[n] = T\{x_1[n]\}, \quad y_2[n] = T\{x_2[n]\}$$

$$\Rightarrow T\{ax_1[n] + bx_2[n]\} = ay_1[n] + by_2[n].$$

Desired output sequence would be to smooth out the higher frequencies.

- Time invariant:

A system $T\{\cdot\}$ is said to be time invariant if $y[n] = T\{x[n]\} \Rightarrow T\{x[n-n_0]\} = y[n-n_0]$

(Basically, time-shifted input produces time-shifted output)

- Linear & Time Invariant (LTI) System:

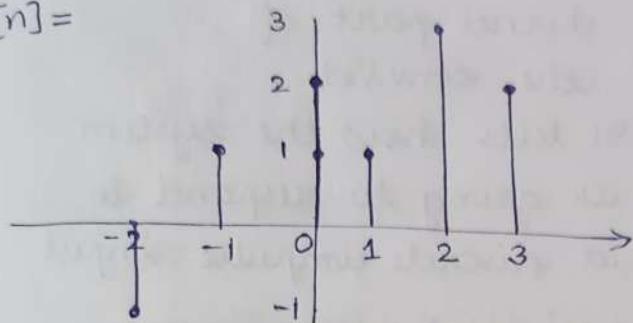
A system $\boxed{T\{\cdot\}}$ that satisfies both the linearity and time invariance properties simultaneously.

→ A consequence of LTI is the convolution sum.

→ Any discrete (digital) signal or sequence $x[n]$ can be expressed as a sum of scaled and shifted discrete delta functions.

$$\delta[n] = 1 \text{ if } n=0 \\ 0, \text{ if } n \neq 0.$$

$$x[n] =$$



$$x[n] = \begin{cases} -1 & n = -2 \\ 1 & n = -1 \\ 2 & n = 0 \\ 1 & n = 1 \\ 3 & n = 2 \\ 2 & n = 3 \end{cases}$$

$x[n]$ can be written as,

$$x[n] = -1 \cdot \delta[n-(-2)] + 1 \cdot \delta[n-(-1)] + 2 \cdot \delta[n-0] + \\ 1 \cdot \delta[n-1] + 3 \cdot \delta[n-2] + 2 \cdot \delta[n-3].$$

- In general, any discrete signal $x[n]$ can be written as

$$x[n] = \sum_{k=-\infty}^{\infty} x[k] \cdot \delta[n-k].$$

→ Let's pass this signal through an LTI system.
 $T\{y\}$.

$$y[n] = T\{x[n]\} \\ = T\left\{ \sum_{k=-\infty}^{\infty} x[k] \delta[n-k] \right\} \\ = \sum_{k=-\infty}^{\infty} x[k] T\{\delta[n-k]\} \quad (\text{due to linearity})$$

$$\text{If } h[n] = T\{\delta[n]\}$$

$$\Rightarrow T\{\delta[n-k]\} = h[n-k] \quad (\text{due to time invariance})$$

$h[n]$ is called the impulse response of the system $T\{y\}$.

$$\therefore y[n] = \sum_{k=-\infty}^{\infty} x[k] \cdot h[n-k] \rightarrow (\text{Convolution sum})$$

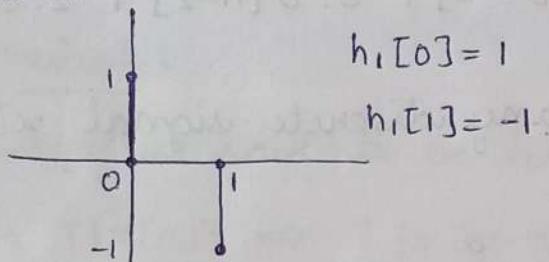
↓
kernel part of
the convNet.

It tells how the system
is going to respond to
a discrete impulse input.

Example :

$$\text{Let } h_1[n] = \frac{\delta[n] - \delta[n-1]}{1}$$

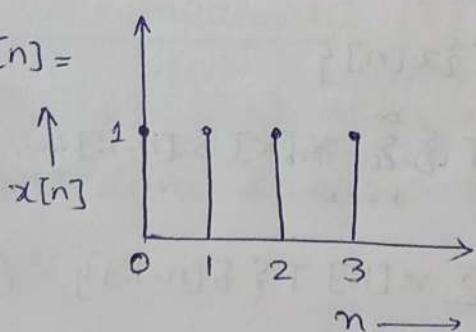
This is called as discrete differentiation operation.



$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

2D version of this is used to detect edges.

$$\text{Let } x[n] =$$



$$\begin{matrix} \delta[n] \\ \delta[n] \end{matrix} \xrightarrow{h_1[n]} \boxed{h_1[n]} \xrightarrow{} \delta[n] - \delta[n-1]$$

$$x[n] = \sum_{k=0}^{3} 1 \cdot \delta[n-k]$$

$$\begin{aligned}
 y[0] &= \sum_{k=0}^3 x[k] h[-k] \\
 &= x[0]h[0] + x[1]h[-1] + x[2]h[-2] + x[3]h[-3] \\
 &= 1 \cdot 1 \rightarrow \\
 &= \underline{\underline{1}}
 \end{aligned}$$

$$\begin{aligned}
 y[2] &= \sum_{k=0}^3 x[k] h[2-k] \\
 &= x[0]h[2] + x[1]h[1] + x[2]h[0] + x[3]h[-1] \\
 &= 1(-1) + 1(1) \\
 &= \underline{\underline{0}}. \\
 y[4] &= \sum_{k=0}^3 x[k] \cdot h[4-k] = \underline{\underline{-1}}.
 \end{aligned}$$

17th March '25
Monday

Computer Vision

- Recap: LTI Systems, Convolution sum.
- Linear Image Filtering → Used to emphasize or suppress certain features.
 - Low Pass
 - High Pass
- LTI systems:

Uniquely specified by its ~~or~~ impulse response.

$$h[n] = T\{\delta[n]\}$$

- For any input $x[n] = \sum_{k=-\infty}^{\infty} x[k] \cdot \delta[n-k]$

$$y[n] = T\{x[n]\} = \sum_{k=-\infty}^{\infty} x[k] h[n-k] \quad \text{---} \circledast$$

Low Pass filters

- Smooth the image
- Remove noise
- Blur sharp details like edges.

High Pass filters

- Edges it can detect.
- Sharp changes features can be extracted.

- Example: $h[n] = \delta[n] - \delta[n-1]$

This is high-pass filter. Why?

A: Becoz if sharp changes are there,
~~if~~ $x[n] \neq x[n-1]$ very different, $h[n]$ will be large, if ~~smooth~~ smooth changes then $x[n] \approx x[n-1]$, so 0.

$$x[n] = \begin{cases} 0 & \text{else} \\ 1 & 0 \leq n \leq 3 \end{cases}$$

We have get

$$y[n] = \begin{cases} -1 & n=4 \\ 0 & 1 \leq n \leq 3 \\ 1 & n=0 \end{cases}$$

We get non-zero response at 2 places.

This high-pass filter is used for edge-detection which is used in a number of downstream tasks.

- Example: $h[n] = \frac{\delta[n] + \delta[n-1]}{2}$

This is a low-pass filter.

$$x[n] = \begin{cases} 1 & 0 \leq n \leq 3 \\ 0 & \text{else.} \end{cases}$$

$$y[n] = \begin{cases} \frac{1}{2} & n=0 \\ 1 & 1 \leq n \leq 3 \\ \frac{1}{2} & n=4 \end{cases}$$

in low-pass, locations of change are smoothed.

$y[n]$ is a "blurred" version of $x[n]$.

→ Why do we need LTI systems?

- They are excellent first order solutions to many problems.
- Applications:
 - Denoising.
 - Edge detection.
 - Approximating visual systems.

Image filtering example:

$$h[i,j] = \frac{1}{9} \quad 0 \leq i \leq 2, 0 \leq j \leq 2. \quad h = \frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

$$I[i,j] = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix}$$

Q. Can you compute $J(0,0)$?

$$J = I * h. \quad J[0,0] = \frac{1}{9} [1+2+5+6] = \frac{14}{9}$$

Q. Find $J[3,3]$.

$$\frac{1}{9} [11+12+15+16] = \frac{54}{9} = \underline{\underline{6}}$$

→ For edge detection, we use high pass filter.

We want to detect the change in intensity.

18/03/2025
Tuesday

Computer Vision

- Recap: Linear Image Filters.

- Optical Flow.

- Convolution sum: $y[n] = \sum_{k=-\infty}^{\infty} x[k] h[n-k]$

∴ System is LTI and any signal $x[n]$ can be written as $x[n] = \sum_{k=-\infty}^{\infty} x[k] s[n-k]$

$$T\{x[n]\} = \sum_{k=-\infty}^{\infty} x[k] T\{s[n-k]\}$$

$\underbrace{\hspace{1cm}}$
 $h[n-k]$

- $s[m,n] = \begin{cases} 1 & \text{if } m=0 \& n=0 \\ 0 & \text{o/w.} \end{cases}$

$$I[i,j] = \sum_u \sum_v I[u,v] s[i-u, j-v]$$

$$J = I * h$$

$$J[i,j] = \sum_k \sum_m (I[k,m]) \cdot h[i-k, j-m]$$

2D convolution sum.

Low-pass filter: (Eg:- Mean, Gaussian)

- Smoothes the image.
- Reduces high frequency components, which means:
 - Reduces noise
 - Also blurs edges
- Think of it like averaging neighbourhood pixel values.

High-pass filter: (Eg: Laplacian).

- Highlights high-frequency components like edges or sudden changes in intensity.

• Suppresses smooth

Q. What happens ~~to~~ when you apply high-pass to noisy image?

A. Noisy image = high frequency everywhere.

Noise causes sudden small changes in pixel values across the image. That means lots of high-frequency content, even in places with no real edge.

Generally, when you apply high-pass filter, these sudden changes are amplified.

~~So, when you app~~ For noisy image, when high-pass filter is applied you end up detecting noise as false edges.

Real edge detection would be hampered.

So you get a black image.

→ Therefore, to denoise, we apply low-pass filter, because it removes unwanted high-frequency components. (Smooths out random noise).

Now you apply high pass filter. (like Laplacian). This will help in responding to significant edges.

Laplacian :-
Filter
$$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 0 \end{bmatrix}$$
 → It is a second-order derivative func filter.
↓

Note that center gets subtracted by its neighbours.

→ When you reduce noise from 100 to 50 to 10, you can see small edges in the averaged image.

Optical Flow: → used in motion estimation

Image acquisition:

When a camera captures a photo:

Light (photons) hits the sensor. The light is converted to electric signals — intensities.

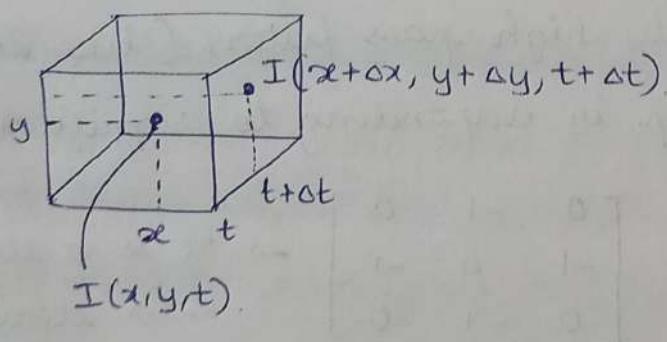
This creates a digital image.

→ $I(x, y, t)$: Intensity of a pixel at time 't'.

→ Here we ~~are not~~ not working with static image but we are working with a video or image sequence over time.

→ Optical flow means "how pixels are moving between one frame and the next?".

→ We want to track how the intensity at point (x, y, t) moves to a new location $(x + \Delta x, y + \Delta y)$ ~~at time $t + \Delta t$~~ .



Let intensity $I(x, y, t)$ at spatial location (x, y) and temporal location 't' move to spatial location $(x + \Delta x, y + \Delta y)$ at time $(t + \Delta t)$.

(Note : Intensity values is the same. Just changing in location.)

$$I(x+\Delta x, y+\Delta y, t+\Delta t) = I(x, y, z) \quad \text{--- (1)}$$

Let's apply the Taylor series expansion to

$$I(x+\Delta x, y+\Delta y, t+\Delta t) \text{ at } (x, y, t)$$

$$\begin{aligned} I(x+\Delta x, y+\Delta y, t+\Delta t) &= I(x, y, z) + \Delta x \cdot \frac{\partial I}{\partial x} + \Delta y \cdot \frac{\partial I}{\partial y} \\ &\quad + \Delta t \cdot \frac{\partial I}{\partial t} + \underbrace{\text{higher order terms}}_{\downarrow} \end{aligned}$$

ignored for natural videos at reasonable rates (15 fmps).

(Bcz they are highly correlated, the changes move slowly and smoothly, so we can ignore higher order terms).

→ Using (1), we get,

$$I(x, y, t) + \Delta x \cdot \frac{\partial I}{\partial x} + \Delta y \cdot \frac{\partial I}{\partial y} + \Delta t \frac{\partial I}{\partial t} = I(x, y, z)$$

We are interested in how quickly the intensity is changing i.e. $\frac{\Delta x}{\Delta t}$ and $\frac{\Delta y}{\Delta t}$.

From (2),

$$\frac{\Delta x}{\Delta t} \cdot \frac{\partial I}{\partial x} + \frac{\Delta y}{\Delta t} \cdot \frac{\partial I}{\partial y} + \frac{\partial I}{\partial t} = 0$$

$$\Rightarrow [I_x \cdot u + I_y \cdot v + I_t = 0] \quad \begin{matrix} \leftarrow & \text{Optical flow} \\ & \text{constraint} \end{matrix}$$

where $u = \lim_{\Delta t \rightarrow 0} \frac{\Delta x}{\Delta t}$, $v = \lim_{\Delta t \rightarrow 0} \frac{\Delta y}{\Delta t}$ $\begin{matrix} \leftarrow & \text{Velocity} \\ & \text{of} \\ & I(x, y, t) \end{matrix}$

$$I_x = \frac{\partial I}{\partial x}; \quad I_y = \frac{\partial I}{\partial y}; \quad I_t = \frac{\partial I}{\partial t}$$

→ We want to know u and v .
 Since we have two variables, we need at least one more equation other than this constraint to solve the problem.

21/03/2025
Friday

Computer Vision

- Recap : Optical Flow (OF)
- OFCE (Optical flow constraint equation)
- Smoothness constant.

Optical Flow:

- $I(x, y, t) = I(x + \Delta x, y + \Delta y, t + \Delta t)$
- $$\boxed{I_x u + I_y v + I_t = 0} \rightarrow \text{OFCE (Optical flow constraint equation)}$$

$$\begin{array}{ccc} \downarrow & \downarrow & \downarrow \\ \frac{\partial I}{\partial x} & \frac{\partial I}{\partial y} & \frac{\partial I}{\partial t} \end{array}$$

$$u = \lim_{\Delta t \rightarrow 0} \frac{\Delta x}{\Delta t} \quad v = \lim_{\Delta t \rightarrow 0} \frac{\Delta y}{\Delta t}$$

The optical flow vector is $\underline{f} = (u, v)$.

u_x, u_y, v_x, v_y — flow gradients can be computed.

Velocities are not going to rapidly change spatially.

- Since OFCE is not sufficient to find (u, v) , we introduce the smoothness condition on the flow.

→ We use Horn-Schunck method.

We assume that neighbouring pixels tend to move in a similar way.

- This method minimizes the energy function.

$$E_0 = \iint_{xy} (I_x u + I_y v + I_t)^2 dx dy.$$

In ideal situation, this is going to be zero. Because higher order terms are going to be zero in taylor expansion series. But this is not the case in real world.

$$E_s = \iint_{xy} (u_x^2 + u_y^2 + v_x^2 + v_y^2) dx dy$$

$$u_x = \frac{\partial u}{\partial x} \quad u_y = \frac{\partial u}{\partial y} \quad v_x = \frac{\partial v}{\partial x} \quad v_y = \frac{\partial v}{\partial y}$$

- Overall loss $E = E_s + \lambda E_0$
- ↓ →
 smoothness flow
 constraint

E_0 enforces brightness constancy.

E_s encourages neighbouring motion vectors to be similar.

λ controls balance between smoothness and flow constraint.

- How to solve?

We can set $\frac{\partial E}{\partial u} = 0$; $\frac{\partial E}{\partial v} = 0$ and solve for u, v .

→ In practice, we solve the discrete version of the problem.

$$E_0 = \sum_{i=0}^{N-1} \sum_{j=0}^{N-1} (I_x u[i,j] + I_y v[i,j] + I_t)^2 \quad \text{--- (1)}$$

$$E_s = \sum_{i=0}^{N-1} \sum_{j=0}^{N-1} ((u_x[i,j])^2 + (u_y[i,j])^2 + (v_x[i,j])^2 + (v_y[i,j])^2) \quad \text{--- (2)}$$

$$\left. \begin{array}{l} u_x = \cancel{E_u(i+1,j)} (u[i+1,j] - u[i,j]); \\ u_y = (u[i,j+1] - u[i,j]); \\ v_x = [v[i+1,j] - v[i,j]]; \\ v_y = [v[i,j+1] - v[i,j]] \end{array} \right\} \quad \text{--- (3)}$$

Find u, v s.t. $E = E_s + \lambda E_0$ is minimized.

$$\frac{\partial E}{\partial u} = 0, \quad \frac{\partial E}{\partial v} = 0. \quad \text{--- (4)}$$

Substitute (3) in (4).

$$\frac{\partial E_0}{\partial u[i,j]} = 2 [I_x u + I_y v + I_t] \cdot I_x$$

$$\frac{\partial E_0}{\partial v[i,j]} = 2 [I_x u + I_y v + I_t] \cdot I_y$$

$$\frac{\partial E_s}{\partial u[i,j]} = \cancel{2E_u} 2 (u[i,j] - u_{avg}(i,j))$$

$$\frac{\partial E_s}{\partial v[i,j]} = 2 (v[i,j] - v_{avg}(i,j))$$

24th march '25
Monday

Computer Vision

Recap: $u = \lim_{t \rightarrow 0} \frac{\Delta x}{\Delta t}; v = \lim_{t \rightarrow 0} \frac{\Delta y}{\Delta t}$.

$$I_x = \frac{\partial I}{\partial x} \quad I_y = \frac{\partial I}{\partial y} \quad I_t = \frac{\partial I}{\partial t} \quad f = (u, v).$$

$$I_x u + I_y v + I_t = 0 \xrightarrow{\text{smoothness condition}}$$

Ideally, the equation won't be zero, because we made some approximations in the Taylor series expansion.

$$E_o = \iint (I_x u + I_y v + I_t)^2 dx dy$$

OFCE → leads to aperture problem.

$$E_s = \iint (u_x^2 + u_y^2 + v_x^2 + v_y^2) dx dy.$$

Minimize $E = E_s + \lambda E_o$; i.e., ensure smoothness and accuracy.

$$E_s = \sum_{ij} (u_x^2 + u_y^2 + v_x^2 + v_y^2)$$

$$u_x(i,j) = u(i+1,j) - u(i,j)$$

$$u_y(i,j) = u(i,j+1) - u(i,j)$$

$$v_x(i,j) = v(i+1,j) - v(i,j)$$

$$v_y(i,j) = v(i,j+1) - v(i,j).$$

Problem: Solve for $u(i,j)$ and $v(i,j)$.

Solution: Set $\frac{\partial E}{\partial u(i,j)} = 0; \frac{\partial E}{\partial v(i,j)} = 0$.

$$E = E_s + \lambda E_0$$

ensures smooth flow

Ensures that we satisfy our intensity constancy.

$$\frac{\partial E}{\partial u} - \bar{u}_{ij} = \frac{1}{4} [u[i,j+1] + u[i+1,j] + u[i-1,j] + u[i,j-1]]$$

$$\frac{\partial E}{\partial u} = 2(u - \bar{u}) + 2\lambda I_x (I_x \cdot u + I_y \cdot v + I_t)$$

$$\frac{\partial E}{\partial v} = 2(v - \bar{v}) + 2\lambda I_y (I_x \cdot u + I_y \cdot v + I_t)$$

$$\text{Now, } \frac{\partial E}{\partial u} = 0, \quad \frac{\partial E}{\partial v} = 0.$$

$$u(1 + \lambda I_x^2) + v(\lambda I_x I_y) + \lambda I_x I_t - \bar{u} = 0$$

$$v(1 + \lambda I_y^2) + u(\lambda I_x I_y) + \lambda I_y I_t - \bar{v} = 0$$

$$u(1 + \cancel{\lambda I_x^2})(1 + \cancel{\lambda I_y^2}) + v(\cancel{\lambda I_x I_y})(1 + \cancel{\lambda I_y^2}) + (\cancel{\lambda I_x I_t} - \bar{u})(1 + \cancel{\lambda I_x^2}) = 0$$

$$u(1 + \lambda I_x^2)(\lambda I_x I_y) + v(\lambda^2 I_x^2 I_y^2) + (\lambda I_x I_t - \bar{u})(\lambda I_x I_y) = 0$$

$$u(1 + \lambda I_x^2)(\lambda I_x I_y) + v(1 + \lambda I_y^2)(1 + \lambda I_x^2) + (\lambda I_y I_t - \bar{v})(1 + \lambda I_x^2) = 0$$

$$v[1 + \lambda I_x^2 + \lambda I_y^2] + \lambda I_y I_t - \bar{v} + \lambda^2 I_x^2 I_y I_t - \lambda I_x^2 \bar{v} - \lambda^2 I_x^2 I_y I_t + \lambda I_x I_y \bar{u} = 0.$$

$$v = \bar{v} - \frac{\lambda(I_x \bar{u} + I_y \bar{v} + I_t) \cdot I_y}{(1 + \lambda I_x^2 + \lambda I_y^2)} \cdot I_y$$

$$u = \bar{u} - \frac{\lambda(I_x \bar{u} + I_y \bar{v} + I_t) \cdot I_x}{(1 + \lambda I_x^2 + \lambda I_y^2)} \cdot I_x$$

Q. Does this Tom-Schunk method converge?

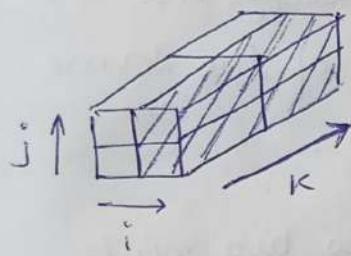
A. Yes, because it's based on minimizing a convex quadratic energy function.

Q. We don't know explicitly I_x, I_y values. How do we get them?

A. Empirical Values.

OR by this method:

We can approximate I_x using intensity values across two frames (at time t and $t+1$).



The slices we consider are the shaded one and unshaded one.

$$I_x = \frac{1}{4} [I(i+1, j, k) + I(i+1, j+1, k) + I(i+1, j, k+1) + I(i+1, j+1, k+1)] - \frac{1}{4} [I(i, j, k) + I(i, j+1, k) + I(i, j, k+1) + I(i, j+1, k+1)]$$

Recall: OFCE desired under the intensity constancy assumption.

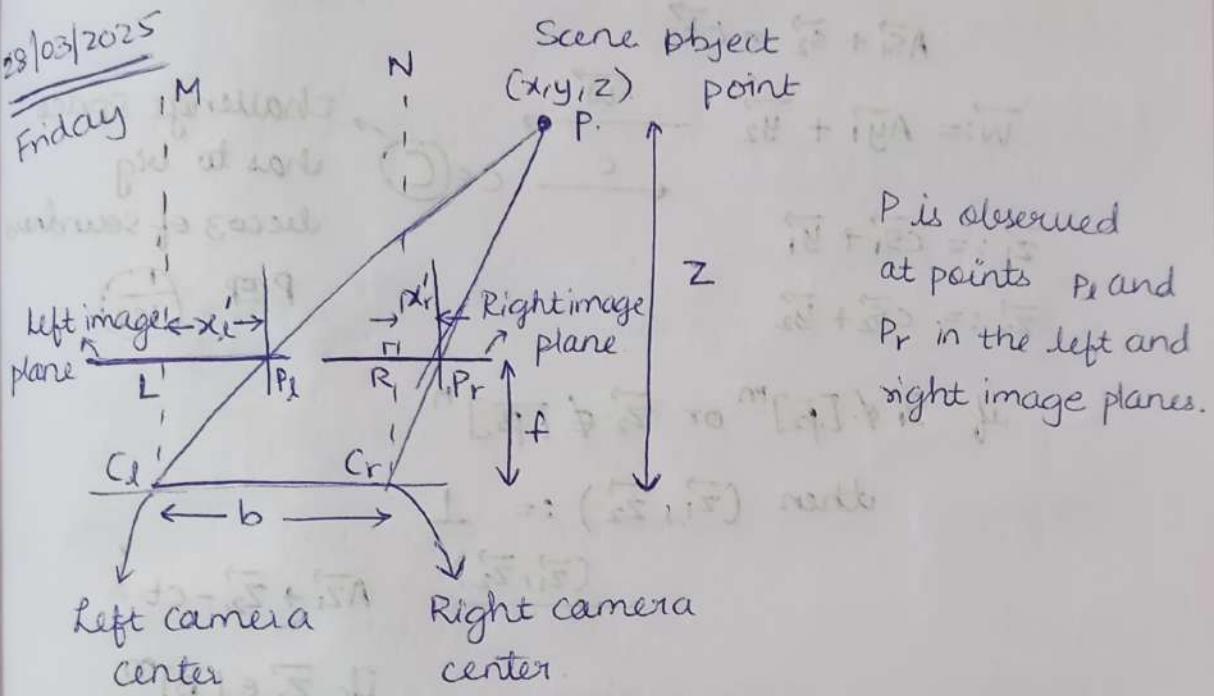
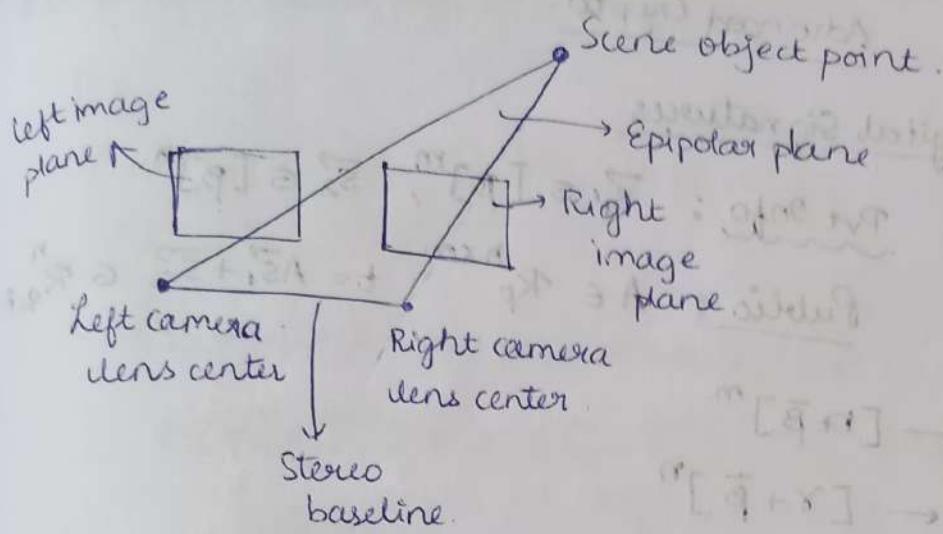
25/03/2025
Tuesday

Computer Vision

Depth Estimation

- Binocular stereo — it tells how machines perceive depth by mimicking human vision.
We use two cameras to estimate how far away things are.
- Two cameras separated by a distance ' b '.
(Distance b/w centers). — baseline ' b '.
- Image planes are aligned coplanar.
- Conjugate pair : Two points in different images that are the projections of the same point.
- Disparity : Horizontal distance b/w two points in a conjugate pair, when the two images are superimposed.
- Epipolar plane : Plane formed by left camera center, right camera center and 3D point.
- Epipolar line : Line where epipolar plane intersects with the image plane.
- Feature in the left image at a particular y location will be present at the same y level in right image, just shifted horizontally, not up or down.

So, no vertical disparity.



PNC_L and $P_L C_L$ are similar.

$$\frac{x}{z} = \frac{x'_i}{f} \rightarrow P_L = \frac{P_L}{f}$$

PNC_R and $P_R C_R$ are similar.

$$\frac{x-b}{z} = \frac{x'_i}{f} \rightarrow P_R = \frac{P_R}{f}$$

$$\frac{x}{z} - \left(\frac{x-b}{z} \right) = \frac{P_L - P_R}{f} \Rightarrow \boxed{z = \frac{bf}{P_L - P_R}}$$

Depth is z . So we can see that depth can be recovered by knowing disparity of corresponding points.

Relating baseline depth and disparity.

$b \rightarrow$ baseline separation.

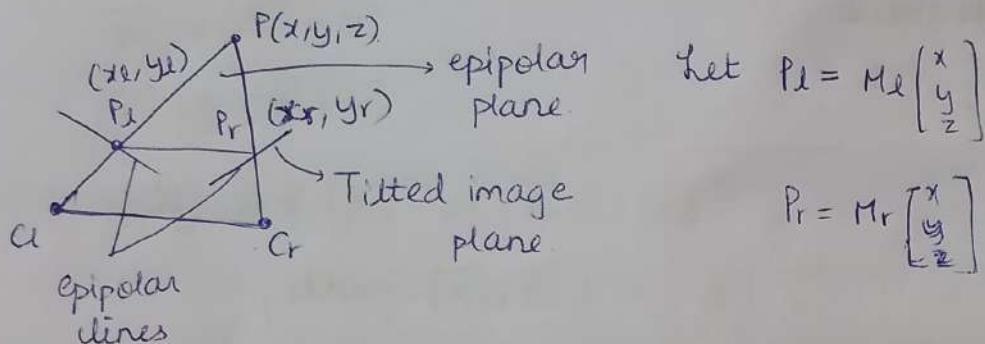
$f \rightarrow$ focal length of the camera.

$P_L - P_r \rightarrow$ disparity

→ We started with the assumption that this setup is ideal, i.e., the image planes are coplanar.

→ The image planes need not be coplanar in a non-ideal scenario.

epipolar plane: plane joining centers of the camera and the scene object.



Correspondence problem: For each point in the left image, find the corresponding point in the right image.

Challenges in depth estimation

→ Image planes need not be coplanar.

→ One image plane could be rotated w.r.t. the other.

→ Correspondence problem.

① Correspondence Problem.

— find key points in left and right images.

— classical methods like SIFT, SURF,
aSIFT etc.

— DL methods.

- Find corresponding key points using matching algs like RANSAC.
- Improve accuracy by our knowledge of 3D-geometry of the scene-specific epipolar geometry.
 - Require Intrinsic Camera Parameter.
- Scale Invariant Feature Transform (SIFT)

10/4/2025
Tuesday

Computer Vision.

Recap: Challenges in depth estimation;
correspondence problem.

- Scale Invariant Feature Transform (SIFT):

- Scale space construction — $\text{DoG} \rightarrow$ Difference of Gaussians.
- Key points from extrema.
- Refining key points.
- Orientation assignment.
- Keypoint description.

Depth and disparity have inverse relation.

→ SIFT algo is used in stereo vision or matching problems where we want to find corresponding points — called matching pairs,
 { basically, left and right camera view of the same scene. }

How to find keypoints in both images?

Goal is :

To detect key points in an image that are:

- Distinctive - easy to match across images
- Invariant to scale : Means whether ^{image} zoomed in or out, same keypoint can be detected.
- Invariant to rotation: Even if image is rotated, same keypoint . . .

• First, we define keyspace

Keypoints are defined in 3D space : (x, y, σ) .

x, y = coord σ = scale at which
keypoint is detected.

• Create scale space:

Take image, apply low-pass filter.

Blurred images are created. Each set of these blurred images is called octave.

→ In octave, image is downsampled by a factor of 2 in both width and height.

• (DoG) : Difference of Gaussians : Bandpass Filtering:

Bandpass filters are obtained by subtracting two Gaussian blurred images.

$$\text{DoG} = G(x, y, \sigma_2) - G(x, y, \sigma_1)$$

* Band pass regions contain keypoints.

• Now we check for ^{local} extrema points in band pass regions.

How do check for extrema?

We compare each pixel to its 26 neighbours.

26 ~~new~~ neighbours?

- (1) 8 in the current scale.
- (2) 9 in the scale above.
- (3) 9 in the scale below.

The extrema we get are potential keypoints.

- Now we got keypoints extrema, but some might not be useful. Like they might be lying on edges. So we need to do refining.
- We do pruning. This will prune spurious points. It discards keypoints with low contrast and also those which lie on edges.
- We make keypoints rotation invariant and scale invariant.
- How to make keypoints rotation invariant?
 - * Around each keypoint, we calculate gradient orientation in neighbourhood.
 - * We construct histogram of orientations.
 - The peak of histogram is assigned as the dominant orientation.
(Multiple orientations are also possible).
- Till now, we have for each keypoint
 $(x, y, \text{scale}, \text{orientation})$.

Q) Why are we making keypoints invariant to scale and orientation?

Ans: Small details get lost when we shrink an image. We want to detect keypoints at multiple scales and ensure that they are not lost during downsampling or zooming.

→ Also, when an image is rotated, the orientation of keypoints changes.

Our algo already knows each keypoint's orientation, so it subtracts this angle.

So makes it rotation invariant.

→ Now we look at Histogram of Gradients (HoG).

Take a keypoint, around it take a 16×16 pixel patch.

Divide this into 4×4 sub-regions $\rightarrow [16]$.

Around each subregion, we compute a 8-bin histogram (like for 8 directions).

$$\text{So } 16 \times 8 = 128.$$

Feature Vector is 128-dimensional. (normalized).

This ~~is~~ is illumination invariant.

(Intensity scaling won't effect).

→ ~~**~~ Keypoint now has:

$(x, y) \rightarrow$ location.
Scale + 81 degrees
Orientation.
HoG Descriptor \rightarrow (128-length)

→ Now do matching, take keypoints from two images and compare feature vectors.

Eg: Suppose we want to detect a clock.

We have a database of objects with precomputed SIFT features.

Now for our image (query), compute SIFT features.

Match these with that of database images.

If good no. of matches, object is present.

Summary:

→ Keypoint candidates are the maxima of DoG-filtered responses.

DoG response to image is difference of Gaussians.

→ Why do we need band-pass filtered regions?

Because keypoints are found there.

→ A DoG-filtered sample is an extrema if largest or smallest of 26 surrounding pixels, in scale-space.

→ The gradient magnitudes and orientations of patch around the keypoint is computed.

→ Then we normalize the histogram.

Why? To be invariant to illumination changes.

Matching:

Now we have a set of keypoints and their representations. We compare the values.

Goal: We want keypoints that are invariant to resolution.

April 4th 2025
Friday

Computer Vision.

- Recap: SIFT algorithm.
- Camera matrix — Notation: Inhomogenous, Homogenous, Augmented.
- Extrinsic — R, t : Rotation and Translation.
- Intrinsic — K ; camera properties — focal length, scaling, center location.
- Projection — $M = K[R, t]$.

→ After detecting keypoint, we now want representation.
We use HoG.

Take a patch (neighbourhood). At each pixel compute gradient magnitudes and orientation.

Quantize orientation into 8 bins. 45° range ($\frac{360}{45} = 8$)

Finally — feature vector (HoG).

Camera Parameters

Goal: We want to project a point from 3D space (real world) into 2D pixel space of image.
(How object is seen on camera?)

Points to note:

- It accommodates rotation as well as translation in our camera matrix.
- Extrinsic Parameters: Defined outside of the camera.
Rotation, Translation.
- Intrinsic Parameters: Built into the camera.
Focal length, Resolution of the sensor, center of

image. ($\frac{\text{img width}}{2}, \frac{\text{height}}{2}$) Eg: 640×480 image
= $(320, 240)$.

algorithm
→ RANSAC Problem helps in like when we have keypoints (feature vectors HOG), we need to match across images. RANSAC filters out bad matches.

→ How does epipolar line help?

Given a key point in the left image, the corresponding point in the right image must lie on the epipolar line.

So instead of searching entire image, we only search along epipolar line (complexity is reduced.)

• 3D to 2D projection.

— Expression

— Example: $\xrightarrow{\text{intrinsic}}$

$[f \ s]$ extrinsic. $\begin{bmatrix} f & s \\ 0 & 1 \end{bmatrix}$

$$K = \begin{bmatrix} 800 & 0 & 320 \\ 0 & 800 & 240 \\ 0 & 0 & 1 \end{bmatrix} \quad R = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad t = \begin{bmatrix} 0 \\ 0 \\ -5 \end{bmatrix}$$

Where does the point $X = [2, 3, 10, 1]^T$ get projected on the image?

• let's operate in \mathbb{R}^3 .

$P_c \in [x \ y \ z]^T$ → Inhomogeneous Representation.

$\tilde{P} \in \mathbb{P}^3 = \mathbb{R}^4 \rightarrow (0, 0, 0, 0); [x, y, z, w]^T$

— Homogenous Representation.

Here w is scaling factor.

$\bar{P} \in \mathbb{R}^4 = [x, y, z, 1]^T$; Augmented Representation.

- K : Intrinsic Parameters:

$$K = \begin{bmatrix} f_x & s & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix}$$

skew: non-zero for non-rectangular pixels.
 location of the center of the image. (in pixels).
 focal length of the camera along the x and y dimensions.

- Extrinsic parameter:

$[R \ t]$ → Translation matrix (3×1)
 R → Rotation matrix (3×3)

- Overall projection matrix

$$\boxed{M = K[R \ t]}$$

$\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$: K : Intrinsic camera parameters
 $[R \ t]$: Extrinsic.

$$\boxed{i = M \cdot \bar{P}}$$

$$M \bar{P} = \begin{bmatrix} 800 & 0 & 320 \\ 0 & 800 & 240 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & -5 \end{bmatrix} \begin{bmatrix} 2 \\ 3 \\ 10 \\ 1 \end{bmatrix} = \begin{bmatrix} 3200 \\ 3600 \\ 5 \end{bmatrix}$$

(with intrinsic parameters)

$$\bar{P} = \begin{bmatrix} 640 \\ 720 \end{bmatrix}$$

Scaling factor

April 7th 2025
Monday

Computer Vision

Recap: $M = K[R \ t]$ project a 3D location to a 2D location.

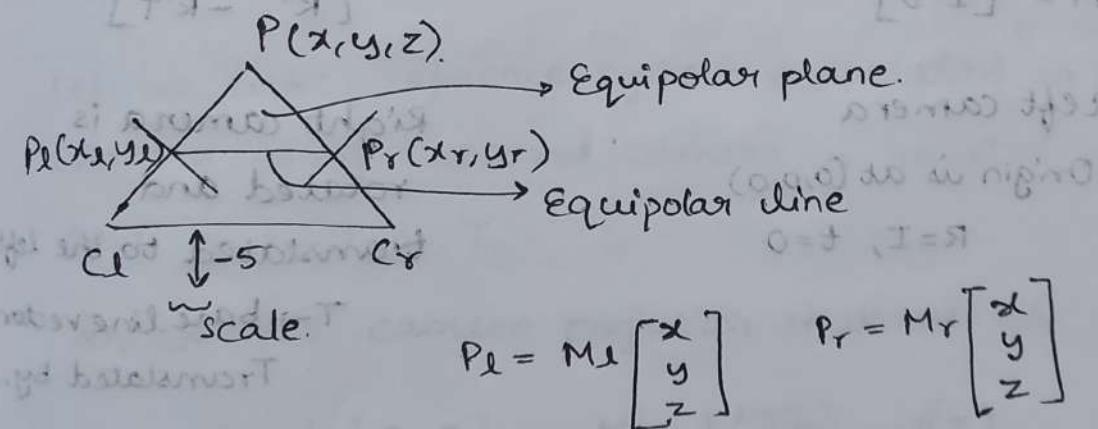
How is K defined? $K = \begin{bmatrix} f_x & s & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix}$

Today:

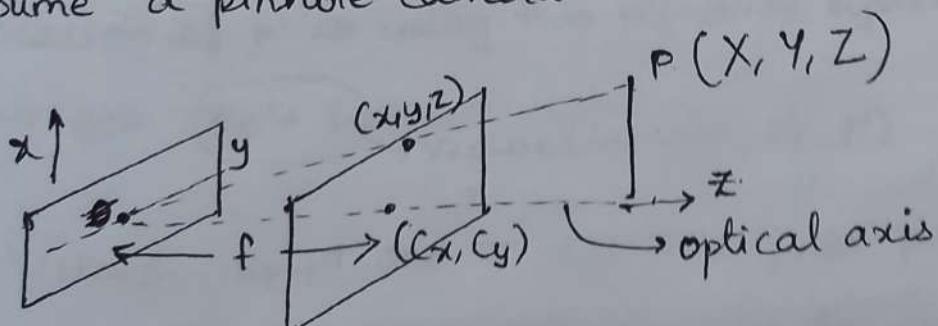
- Pinhole camera model: $x = \frac{fx}{z} + c_x$
- Putting it together:
 - Essential matrix.
 - Epipolar line.
 - Corresponding points.

Depth and disparity are inversely related.

Suppose there are pair of images. We can use DL models (monocular) to circumvent the region where the second can be found.



Assume a pinhole camera model.



$$\frac{x'}{f} = \frac{x}{z} \Rightarrow x' = \frac{xf}{z}$$

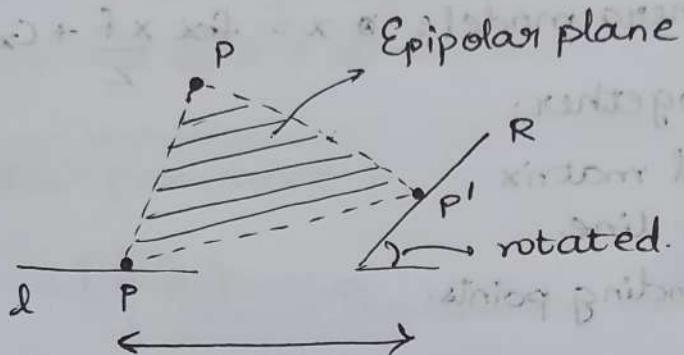
Relative to the image origin,

$$\bullet \quad x' = x - C_x$$

$$x - C_x = \frac{xf}{z}$$

$$\Rightarrow x = \frac{xf}{z} + C_x$$

- Tying it together



Translated by T

$$K^T [R^T \quad -R^T T]$$

$$M = K[R \quad t] \quad ; \quad M' = K^T [R^T \quad R^T t]$$

$$M = K[I \quad 0] \quad ; \quad M' = [R^T \quad R^T T]$$

$$M = [I \quad 0] \quad ; \quad [R^T \quad -R^T T]$$

left camera

Origin is at $(0,0,0)$

$$R = I, \quad t = 0$$

right camera is

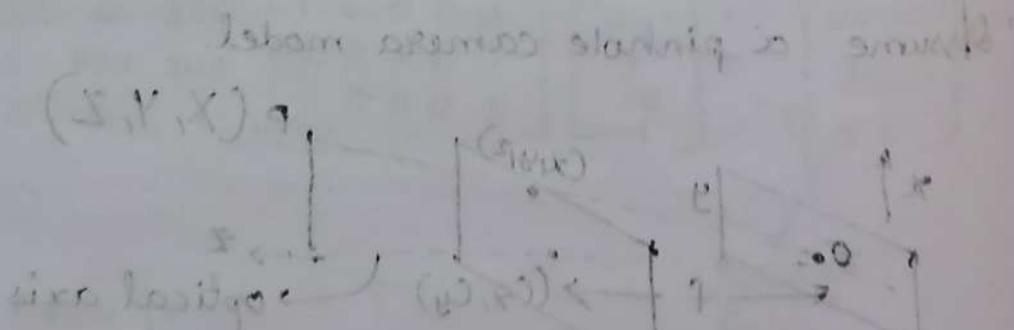
rotated and

translated to the left

$T \rightarrow$ baseline vector

Translated by.

$$\begin{bmatrix} x \\ y \\ z \end{bmatrix} \cdot M = 1 \quad \begin{bmatrix} x \\ y \\ z \end{bmatrix} \cdot M' = 1$$



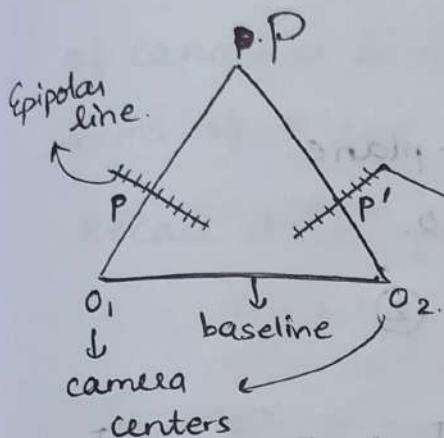
$$\frac{X}{S} = \frac{X'}{S'} \quad \leftarrow \quad \frac{X}{S} = \frac{X'}{\frac{X}{S}}$$

8th April '2025
Tuesday

Computer Vision

Recap: Convex, projection matrix.

- Essential matrix.
- Fundamental matrix.



Two cameras p and p' are observing the same 3D point P .

- We don't know the exact location of P .

Using O_1, O_2 and image point p , we can define epipolar plane. Using epipolar planes, we get the epipolar lines.

$$M = K[R, T]$$

Let us take reference system associated to left camera. and second camera is rotated by R and translation T .

So, we get camera projection matrices as

$$M = K[I \ 0] \quad M' = K' [R^T \ -R^T T]$$

→ location of p' is using this reference system,

we get $(Rp' + T)$ vector (location of p')

→ Vectors $Rp' + T$ and T lie in epipolar plane.

So $\nabla_{Rp' + T} (T \times (Rp' + T)) \rightarrow$ normal to epipolar plane.

We assume the cameras are canonical —
focal lengths are 1.

$$T_x(RP' + T) = T_xRP' \quad [\text{normal to the epipolar plane}]$$

↓ ↓
extrinsic params
of right camera.

①

→ We know p lies on the epipolar plane.
So dot product with normal is 0.

$$P^T \cdot (T_xRP') = 0 \quad \text{--- } ②$$

Generally, if $a = [a_1 \ a_2 \ a_3]^T$, $b = [b_1 \ b_2 \ b_3]^T$
 $c = a \times b = [a_x] \cdot b$ can be written as

$$c = a \times b = \begin{bmatrix} 0 & -a_3 & a_2 \\ a_3 & 0 & -a_1 \\ -a_2 & a_1 & 0 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

Standard property of cross product.

T_xRP' can be written as $\cancel{[T_x][R][P']}$ $[T_x]RP'$ where

$$T_x = \begin{bmatrix} 0 & -a_z & a_y \\ a_z & 0 & -a_x \\ -a_y & a_x & 0 \end{bmatrix}$$

Eqn ② can be written as

$$P^T \cdot [T_x]RP' = 0 \quad \text{--- } ③$$

$$\boxed{P^T \cdot E \cdot P' = 0} \quad \text{where } E = [T_x]R$$

\Leftrightarrow Essential matrix

$$\boxed{\text{Essential Matrix } E = [T_x]R}$$

→ Eqn ③ constraint is used to search for corresponding points.

The above analysis assumed canonical cameras.
i.e. $K = K' = I$.

more general case is non-canonical.

Here also, we will have P and P' (irrespective of canonical or not). Now, to get the canonical forms of P and P' when we have K and K' .

Recall that projection matrices:

$$M = K [I \ 0] \quad M' = K' [R^T \ -R^T T]$$

So, P becomes $K^{-1}P$, P' becomes $K'^{-1}P'$.

→ If we don't have canonical cameras, we find the effective canonical representations of P and P' as follows:

$$P_c = K^{-1} \cdot P \quad \text{--- (4)}$$

$$P'_c = K'^{-1} \cdot P' \quad \text{--- (5)}$$

→ let's plug (4) and (5) into (3).

$$P_c^T \cdot E \cdot P'_c = 0$$

$$\underbrace{P^T (K')^T E \cdot (K'^{-1} P')}_{P^T F P' = 0} = 0$$

$$\boxed{P^T F P' = 0} \quad \text{--- (6)}$$

$$\text{where } F = (K')^T E K'^{-1}$$

Fundamental matrix $F = (K')^T E K'^{-1}$

- F can be estimated directly from "known" points.
- If we know F , then simply knowing a point in an image gives us a constraint (epipolar line) of the corresponding point in other image.
- Difference between this setup and the original setup we discussed few classes before is that:
Before we considered image planes were coplanar with no rotation.

YOLO → You only look once.
 ↳ Globally used for object detection.

11th Apr '2025
 Friday

Computer Vision

Object Detection

IoU Score : (Intersection Over Union)

It basically tells how much overlap does our box have with the desired object.

Q. How can you use the classifier to detect an object?

Take image patches and pass it through the classifier. We have patch level labels. We can even pass overlapping patches. We pick all the patches where the overlapping probability is high.

YOLO

- A popular deep learning based object detection algorithm.
- Tries to identify the potential regions of finding the object.
- Older models like RCNN, Fast RCNN, Faster RCNN lacked in speed although errors were less.
- YOLO formulates object detection as a regression problem. It regresses:
 - (x, y) → object center location
 - width (w) and height (h) of bounding box
 - Confidence Score (class probabilities) of the
 - Class Probabilities bounding box around the Object (class).
- Each bounding box outputs a .
- Divide the input image into $S \times S$ grid.
 - Each grid cell is further divided into B blocks.
 - Each block is assigned a class conditional probability.
- So, each grid cell predicts:
 - B bounding boxes
 - Each box: $x, y, w, h, \text{confidence}$ — 5 parameters.
 - C class probabilities ($C \rightarrow \text{no. of classes}$)

So output tensor has shape :

$$S \times S \times (B \times 5 + C)$$

→ Bounding box outputs confidence score as:

$$\text{Confidence} = \Pr(\text{Object}) \times \text{IOU}_{\text{pred}}^{\text{truth}}$$

prob. that object is present. Intersection over Union b/w predicted box and ground truth

→ Issue: Resolution limitation.

The model is limited, like no. of objects it can be able to detect and where.

If two objects have centers in the same grid cell, one might get missed.

→ We use rectangles for bounding boxes because they have few parameters (x, y, w, h).

→ Combined regression + Classification problem for each box.

→ Q. What can be a candidate model to solve the problem?

Classifier - architecture model (like CNN).

→ Unified detection of all objects in an image.

↓
Does not work on individual proposals.

→ Performs global reasoning about the image for prediction.

(Like all objects simultaneously across the entire image). $[0.1 \times 1 \times 1] \times [2 \times 2]$

- If an object center falls into a grid cell, that grid cell is assigned the responsibility for detecting that object.
- each grid cell predicts B bounding boxes and the confidence for the bounding boxes.
- Bounding box confidence depends on the IOU between ground truth prediction and the probability of an object being present.
- During training, bounding boxes with the highest confidence are retained.
- All other boxes are suppressed.
- During inference, multiple bounding boxes may be predicted.
We use non-maximal suppression to eliminate overlapping boxes with lower confidence.
It basically keeps only the box with highest confidence for overlapping predictions.
- NMS is used to resolve multiple detections.

Non-maximal suppression

It removes duplicate predictions for the same object.

First, it sorts boxes by confidence. It picks the box with the highest confidence, and suppresses all boxes with high IOU overlap.

Cost function:

$$\lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} [(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2] + \text{prior}$$

$$\lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} [(\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2] +$$

$$\sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} (c_i - \hat{c}_i)^2 + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{i,j}^{\text{noobj}} (c_i - \hat{c}_i)^2$$

$$+ \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{\text{classes}} (p_i(c) - \hat{p}_i(c))^2.$$

This term tells what type of object it is.

This term tells us the confidence, or what area of the box is covered with object.

(Intersection of particular box with ground truth.)

This is geometrical.

λ_{noobj} = weighting factor, set to a small value

so that, no-object boxes don't

dominate.

15 April 2025
Tuesday

Computer Vision

$$S \times S = (7 \times 7)$$

$$B = 2 \rightarrow \text{PASCAL VOC dataset}$$

$$C = 20 \text{ (no. of object classes)}$$

$(x, y), (w, h) \rightarrow$ normalized w.r.t dimensions of image.

$$c_i \text{ of } b_i = \Pr(\text{object}) \times \text{IOU}_{\text{truth}}^{\text{pred}}$$

$$P(c) = P(\text{class} = c | \text{object}) \rightarrow \text{class conditional probability.}$$

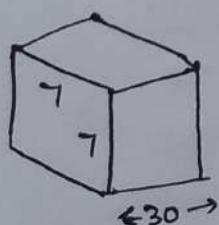
$1_{i,j}^{\text{obj}}$ = j^{th} bounding box predictor for the i^{th} grid cell responsible for the object detection.

→ In the loss function, we take square-root over w_i and h_i because we want to suppress the dominance of width and height.

→ Bounding box need not lie within the cell. There is no such restriction.

→ Why multiple boxes?
more the no. of candidates for prediction; the better would be the overlap with the object.

→ If one box is large and another is smaller but fits the object tighter, the smaller one will have a higher IOU and will be favoured.



We know output tensor has shape $S \times S \times (B \times 5 + C)$.
 $7 \times 7 \times (2 \times 5 + 20)$.
 $\Rightarrow 7 \times 7 \times 30$.

21st April '25
Monday

Computer Vision

Image Segmentation.

Q. How is image segmentation different from object detection?

A. In image segmentation, each pixel in an image is assigned a label.

In object detection, instead of assigning probability to every image, we try to assign probability to every bounding box in object detection.

Semantic segmentation:

Partitions image into semantically meaningful partitions. Assigns each pixel a class label.

Does not distinguish b/w different instances of the same object.

Instance segmentation:

Assigns IDs to each pixel. Distinguishes objects of the same class.

Panoptic Segmentation:

Combines semantic and instance segmentation.

Assigns both class and IDs.

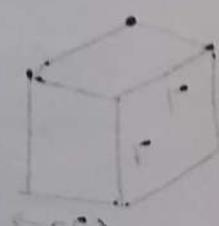
U-Net is used for semantic segmentation

Image segmentation.

(64x64) x 2x2 square grid

(128x128) x 1x1

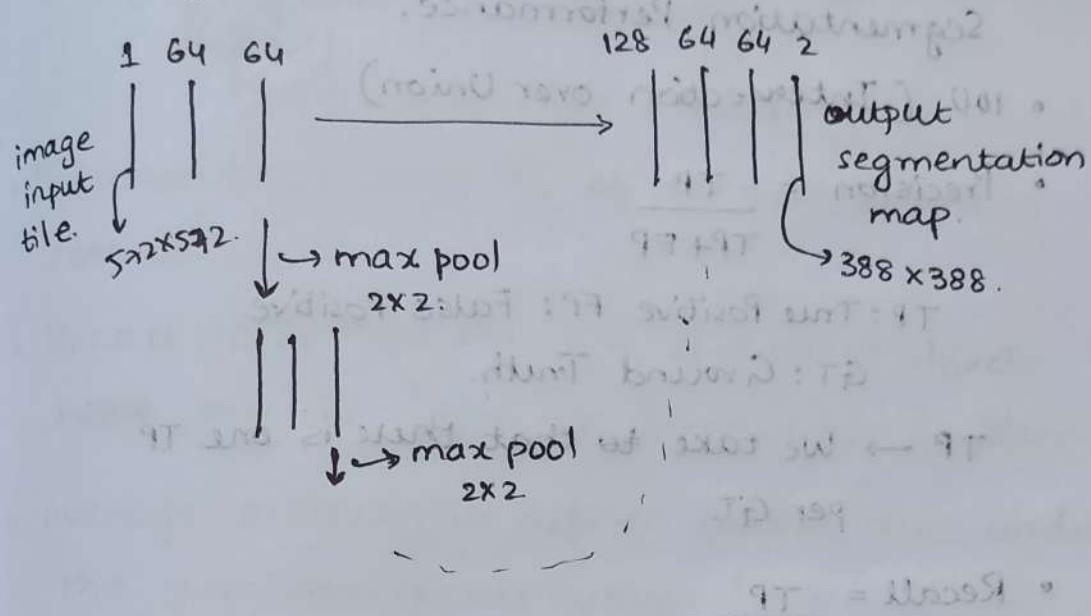
64 x 1x1 =



22nd April '25
Tuesday

Computer Vision

→ UNet architecture is used for semantic segmentation.



→ Input is image of size 572×572 .

Output is a segmentation map of size 388×388 .

→ Left side of this UNet Arch : Encoder (Contracting Path).
This basically learns what is present.

→ Right side : Decoder (Expansive path).

It learns the location.

You start with 1 input channel. Finally you

have 2 output channels. (binary classification).

background or foreground,
and having binary probability.

<u>71/71</u>	<u>Vol</u>	<u>Probability</u>	<u>label</u>
71	P.0	0.0	1
71	E.0	1.0	0
71	30.0	0.0	E
71	F.0	1.0	P
71	30.0	0.0	0

28th APR' 2025
Monday

Computer Vision

- UNet Model
- Detecting Measuring Detection and Segmentation Performance.
- IoU (Intersection over Union)

$$\bullet \text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

TP: True Positive FP: False Positive

GT: Ground Truth.

TP → We take to that there is one TP per GT.

$$\bullet \text{Recall} = \frac{\text{TP}}{\text{Total no. of objects}}$$

Total no. of objects. To spend a day

Average Precision (AP) : Area under the precision-recall curve.

mAP (Mean Average Precision) : Average Average AP values over multiple classes.

Example: 'Dog' class,

4 GroundTruth objects, 5 predictions with 'Dog' label.

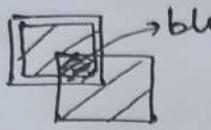
IoU tells about the overlapping b/w the predicted bounding box and ground truth.

Suppose we take ($\text{IoU} \geq 0.5$) to say it as TP.

Pred	Confidence	IoU	TP/FP
1	0.95	0.9	TP
2	0.9	0.3	FP
3	0.8	0.65	TP
4	0.7	0.7	TP
5	0.6	0.25	FP

→ we should have do overlap tiling for smooth segmentation.

• IoU



$$IoU = \frac{\text{blue}}{\text{Total Area}}$$

• Precision tells what % of our detections are correct.

• Recall tells how much % of actual objects have been identified.

→ Average precision is defined as the area under the precision-recall curve.

- Mean AP = $\frac{1}{N_c} \sum_{d=1}^{N_c} AP_d$

- Confidence : Actual probability assigned to obj class

- $\frac{TP}{FP}$ ($IoU \geq 0.5$)

- These are computed w.r.t the closest GT available.

Eg:- 4 GT objects, 5 Preds with 'Dog' label.

Pred	Conf	IoU	TP/FP	Prec	REC
1	0.95	0.9	TP	$\frac{1}{1+0} = 1$	$\frac{1}{4}$
2	0.9	0.3	FP	$\frac{1}{1+1} = 0.5$	$\frac{1}{4}$
3	0.8	0.65	TP	$\frac{2}{2+1} = \frac{2}{3}$	$\frac{2}{4}$
4	0.7	0.7	TP	$\frac{3}{3+1} = \frac{3}{4}$	$\frac{3}{4}$
5	0.6	0.25	FP	$\frac{3}{3+2} = \frac{3}{5}$	$\frac{3}{4}$

* 'Recall' can never be greater than 1.
 We can never flag all the objects as TP.
 (Each GT cannot have more than 1 TP)

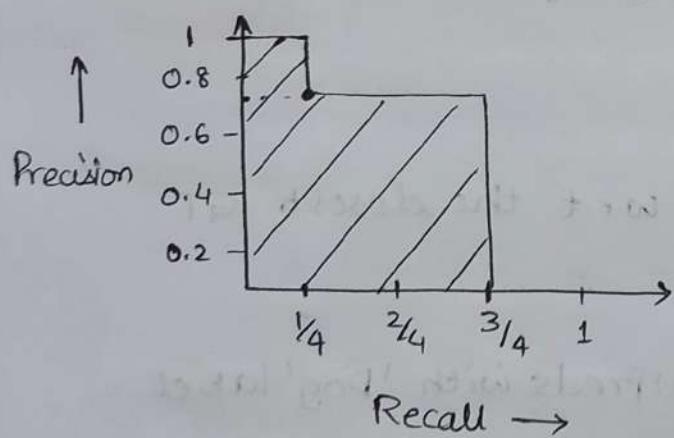
Recall Precision.

0.25	1
0.25	0.5
0.5	0.67
0.75	0.75
0.75	0.6

$$\text{Prec}(r=0.25) = \max(1, 0.5, 0.6, 0.67, 0.75) \\ = 1$$

$$\text{Prec}(r=0.5) = \max(0.67, 0.75, 0.6) \\ = 0.75$$

$$\text{Prec}(r=0.75) = \max(0.75, 0.6) \\ = 0.75$$



Avg precision =
 Area under the
 Precision-recall
 curve.

$$\text{AP} = (0.25 - 0) \times 1 + (0.75 - 0.25) \times 0.75 \\ = 0.25 + 0.5 \times 0.75 \\ = 0.625$$

Q. What happens if we consider $\text{Prec}(r=1)$?
 It does not have a value in the table.

$$\text{P}(r=1) = 0 //$$

$$\text{AP}_{\text{dog}}^{0.5} = 0.625$$

Q. Using this metric, when does the model perform ideally?

- A. When it detects all true objects (high recall) and when it makes few FP (high precision).

Area under curve is close to 1.

In ideal case, model has AP of 1 for each class.

- Average precision (AP) takes care of both precision and recall in one-shot.
- If model is really good, then the fall of mAP values shouldn't be very high.

It should have consistently high AP across all classes. ~~If~~ And also, even if you change the IoU threshold, good model should maintain decent AP.

29 Apr '25
Tuesday

Computer Vision

- A high no. of false positives lowers precision.
- A detector with many FP is considered poor in terms of precision.
- Non maximum suppression is applied after detection (YOLO uses this).

Goal: Keep one bounding box per actual object.

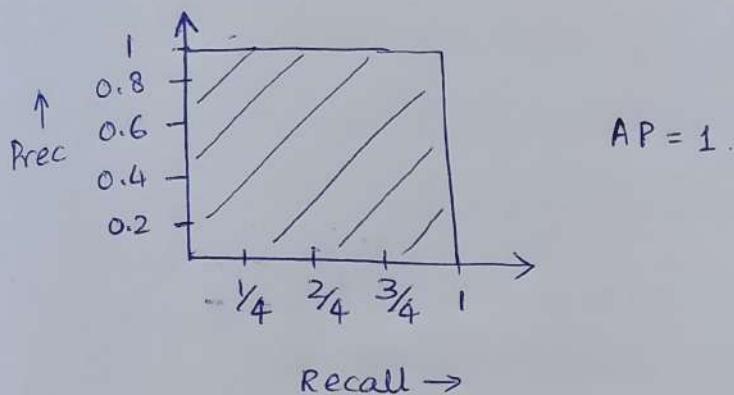
NMS

- Assign confidences to all bounding boxes.
- For each class, sort all detections by confidence score.
- Pick the highest confidence box - best match.
- Now for every other box:
 - Compute IoU with the best box.
 - If $\text{IoU} > \text{threshold}$, discard that box (becoz they are overlapping too much).
- So basically, it reduces duplicate detections for the same object.
- Confidence is used for filtering predictions, ranking detections and plotting precision-recall curves.
- To evaluate whether a prediction is TP, we first match predicted box to a ground truth box using IoU.
If $\text{IoU} \geq \text{threshold}$ its TP. Else it is a FP.
Also, we only want one predicted box per ground truth box. So multiple matches count as extra FPs.
- Can these terms be computed at inference time?
 $\text{Pr}(\text{class}) \rightarrow \text{Yes.}$ But FP, TP need ground truth,
 $\text{NMS} \rightarrow \text{Yes.}$ so they can't.

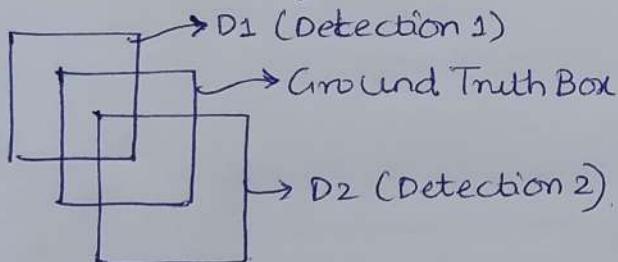
Example: IoU threshold = 0.5, 4 GT, 5 preds.

<u>confidence</u>	<u>IoU</u>	<u>TP/FP</u>	<u>Prec</u>	<u>Recall</u>
0.9	0.95	TP	1	1/4
0.8	0.85	TP	1	2/4
0.6	0.7	TP	1	3/4
0.5	0.6	TP	1	
0.4	0.45	FP	4/5	1

Prec-Recall curve:



→ Now comparing, two detections with Ground truth.



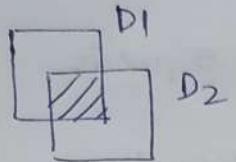
$$\left. \begin{array}{l} P(\text{dog}) = 0.9 \\ \text{IoU} = 0.8 \end{array} \right\} D_1 \quad \left. \begin{array}{l} P(\text{dog}) = 0.85 \\ \text{IoU} = 0.81 \end{array} \right\} D_2$$

Which one do we consider as the best confident box?

$$\text{Conf of } D_1 : 0.9 \times 0.8 = 0.72 \quad \text{Conf of } D_2 = 0.85 \times 0.81 < 0.72.$$

So, we designate D1 as our GT.

Now we check the intersection of D1 and D2.



When is D_2 considered as a potential candidate?

If D_1 and D_2 are overlapping too much, i.e. IoU of D_1 and D_2 is beyond a threshold, say 0.5, we discard it, since D_1 and D_2 in this case would potentially be duplicates.

If $\text{IoU} < 0.5$, then even D_2 is considered as a GT.