

Block A

Training Stability

Block B -

Vision models

Initialization / Optimizers

CNN's Conv, Pooling,
architecture.

Neural network Parameter Initialization.

Training facts \rightarrow vanishing gradients

\rightarrow exploding gradients

\rightarrow Symmetry problems.

Initialization decides how signal flow forward & Gradient
flow backwards.

① Zero Initialization

all weights = 0

all neurons \rightarrow same output

all gradients are identical

neurons never spike ✗

Bias = 0 is fine

weights = 0 is deadly

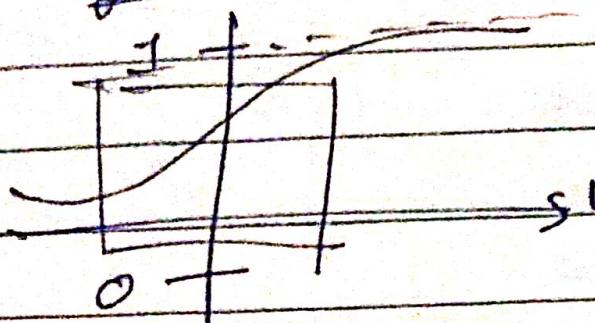
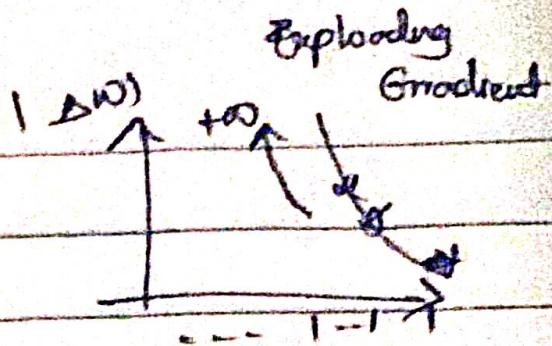
② Random Initialization \rightarrow Small random values.
(normal or uniform)

Problems Too small \rightarrow gradients vanish

too large \rightarrow gradients explode

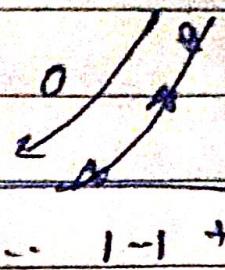
③ Xavier Initialization (Glorot):

Sigmoid



\propto linear combination.

$$f(w)$$



designed for Sigmoid | \tanh

vanishing Gradient

$$\text{Var}(w) = \frac{1}{n_{\text{in}}}$$

→ Output signal should shrink or blow up as depth increases.

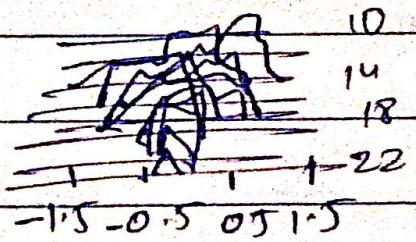
✗ Bad for ReLU.

✓ Great for \tanh / Sigmoid.

④ Kaiming Initialization (He):

$$\sigma = \sqrt{2.0 / n_{\text{in}}}$$

Designed for it



ReLU / Leaky ReLU

$$\text{Var}(w) = \frac{2}{n_{\text{in}}}$$

ReLU → kills half the activation (zeros).

Rule of thumb

ReLU → Kaiming

Tanh → Xavier.

Sigmoid → Xavier

Tanh → Xavier

ReLU → Kaiming

Deep CNN's → Kaiming.

Optimizers → ADAM vs ADAMW.

→ Gradient descent problem.

→ Same learning rates for all parameters

→ Sensitive to Curvature

→ Slow Convergence.

ADAM = momentum +
RMSprop

① First movement → mean of gradients

② Second movement → variance of gradients.

Why works?

Fast convergence

hidden problems

Adaptive learning rate per parameter

Robust to noise gradients

} Adam + L2 regularization

≠

true weight
decay.

AdamW → (modern default)

weight decay is applied directly, not via loss.

→ proper regularization

→ Better Generalization

→ used in Transformers, CNN's, ViTs.

ADAM Adopts Gradient
ADAMW Fixes Regularization

Optimizer Summary

SGD → Theory, Small models.

Adam → fast Convergence

AdamW → production + research default

Block B → CNN: Convolutional Neural Networks

But Connected nets:

→ Ignore Spatial Structure

→ Too many Parameters

CNN's Exploit)

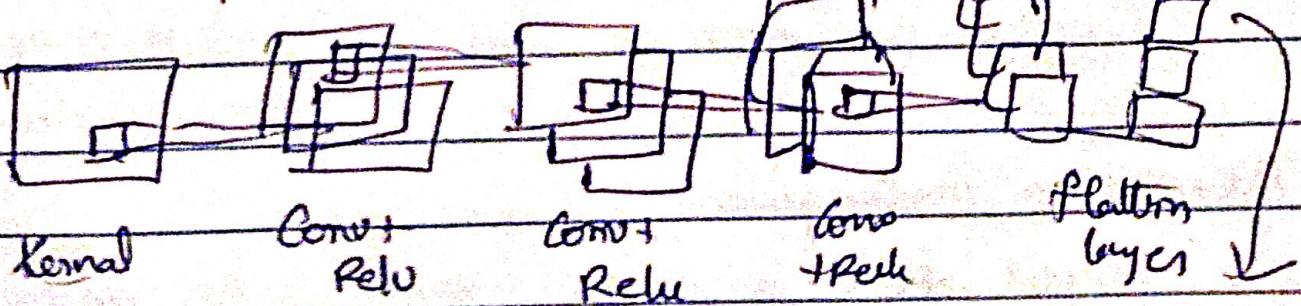
① Locality

② Translation Invariance.

③ Weight Sharing

Pooling.

Convolution operation



feature extraction

output

Fully connected by

Classification

weight sharing = fewer parameters, + translation equivariance +
less over fitting

weight sharing

Small filter slides over Image

→ same filter everywhere

Dot product → feature map

→ learns edges, texture,

patterns -

Parameters

→ kernel size eg (3×3)

→ stride

→ padding

Pooling → → Pooling reduces spatial size

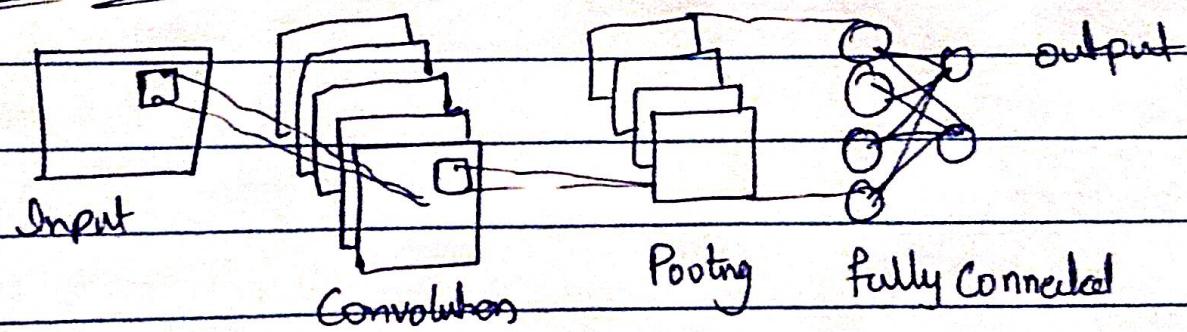
max pooling, average pooling

→ makes model robust to small shifts

Effects → ① DownSampling

② Slight Information loss (Intentional).

CNN Flow



Input



conv → ReLU

Flatten



Pool

Fully Connected.



conv → ReLU

Softmax ||,



Pool

early layer → edges
mid layer → shapes
Deep layer → objects

CNN \rightarrow Feature extraction + classifier

Conv + pool = feature extractor

FC layers = decision maker.