Q1) Degrees of freedom for
homography $= 8$

⟹ we need atleast $n = \lceil \frac{d}{2} \rceil$ points to
compute transformation

⟹ $n = \lceil \frac{8}{2} \rceil = 4$

⟹ given, $w_i = 0.5$

Prob: that the algo never selets $n$ per inlier
points for $k$ iterations is given by

$$(1-w^n)^k = 1 - 0.95$$

⟹ $(1 - (0.5)^4)^k = 0.05$

$$k = \frac{\log(0.05)}{\log(1-(0.5)^4)}$$

$$k = 46.41 \sim 47$$

∴ 47 iterations are required to have 95%
chance of correrted computation

Q2) $\dfrac{\partial f}{\partial w_{ij}}^{(1)} = -?$

Soln.

$w_{ij}^{(1)}$ : Weight of 'link from $j$th to $i$th node after 1st layer

$$\dfrac{\partial f}{\partial w_{ij}}^{(1)} = \dfrac{\partial f}{\partial h^2} \cdot \dfrac{\partial h^2}{\partial h_i^1} \cdot \dfrac{\partial h_i^1}{\partial w_{ij}}^{(1)}$$

\*

$$f = \langle w^3, h^2 \rangle$$

$$\sim \dfrac{\partial f}{\partial h^2} = w^3 \qquad : \begin{bmatrix} w^3_1 \\ w^3_2 \end{bmatrix}$$

\* $h^2 = \sigma(w^2 h^1)$

$\sigma(\cdot) \rightarrow$ sigmoid fn.

$$\boxed{\sigma'(\cdot) = \sigma(\cdot)(1 - \sigma(\cdot))}$$

$\Rightarrow$

$$\dfrac{\partial h^2}{\partial h_i^1} = \sigma'(w^2 h^1) \cdot w_i^{(2)}$$

$$= \sigma(w^2 h^1)(1 - \sigma(w^2 h^1)) \cdot w_i^{(2)}$$

$$\boxed{\dfrac{\partial h^2}{\partial h_i^1} = h^2_{(\cdot)}(1 - h^2) \cdot w_i^{(2)}}$$

Here $W_i^2 = \begin{bmatrix} W_{1i}^{(2)} \\ W_{2i}^{(2)} \end{bmatrix} \Rightarrow$ i th column in $W^2$

$*\quad h^1 = \sigma(W^1 x) \Rightarrow h_i^1 = \sigma\left(\sum_j W_{ij}^1 x_j\right)$

$\Rightarrow \dfrac{\partial h_i^1}{\partial W_{ij}^{(1)}} = \sigma'\left(\sum_j W_{ij}^{(1)} x_j\right) \cdot x_j$

$= \sigma\left(\sum_j W_{ij}^{(1)} x_j\right)\left(1 - \sigma\left(\sum_j W_{ij}^{(1)} x_j\right)\right) \cdot x_j$

$$\boxed{\dfrac{\partial h_i^1}{\partial W_{ij}^{(1)}} = h_i^{(1)}\left(1 - h_i^{(1)}\right) \cdot x_j}\quad \} \text{ scalar}$$

$$\boxed{\dfrac{\partial f}{\partial W_{ij}^{(1)}} = \left\langle W^3, h^2 \odot (1 + h^{(2)}) \odot W_i^{(2)} \right\rangle * h_i^{(1)} * (1 - h_i^{(1)}) \cdot x_j}$$

$$\boxed{= \left[\sum_k W_k^3 \cdot h_k^2 \cdot (1 - h_k^2) \cdot W_{ki}^{(2)}\right] \cdot h_i^{(1)} \cdot (1 - h_i^{(1)}) \cdot x_j} \quad //$$

Q3)
$$\Delta_{ij}^{(2)} := \Delta_{ij}^{(2)} + \delta_i^{(3)} * (a^{(2)})_j^o$$

In vector form:

$$\Delta^{(2)} = \Delta^{(2)} + (\delta)^{(3)} \cdot (a^{(2)})^T$$

Q4)  d inputs, M hidden units, c outputs

a) No. of weights:

$$M \times d + M \times c \geq M(d+c)$$

b) No. of biases:

$$M + C$$

c) No. of independent derivatives :-

for 2nd hidden layer

$$\frac{\partial E}{\partial W_2} = \delta^3 (a^2)^T \quad \left\{ \begin{array}{c} M \text{ to } \\ c \text{ indepent} \\ \text{derivatives} \end{array} \right.$$

For 1st hidden layer

$$\frac{\partial E}{\partial W_1} = \delta^2 (a')^T \quad \left\{ \begin{array}{c} M \text{ independent} \\ \text{derivatives} \end{array} \right.$$

Total no. of independent derivative

$$= M + C \text{ } \text{//}$$

Q5) showing minimizing Sum of Square error is equivalent to MLE

$\Rightarrow$ Let $x_n$ be the input
$w:$ weights

$f(x_n, w):$ Estimate of nueral network

$y_n:$ target data

Sum of Squares error is defined as for '$\underline{N}$' datapoints

$$SOS = \sum_{i=1}^{N} ( y_i - \tilde{y}_i )^2$$

where $\tilde{y}_i = f(x_n, w)$

* Considering that the target data is of the form

$$y_n = f(x_n, w) + \varepsilon_n$$

$\varepsilon_\eta \sim N(0, \Sigma)$ : Multivariate gaussian

Let $y_n$ is a 'd' dim data point

$\quad \varepsilon_n \rightarrow$ also is a $(d \times 1)$ vector.

$\Rightarrow$ Since for a given input $x_n$, the estimate is deterministic using 'w'.

$\quad f(x_n, w)$ is not a RV.

$$\Rightarrow \quad \underline{y_n \sim N\left(f(x_n, w), \Sigma\right)} \Rightarrow dx_1$$
$$\quad\quad x_n, w$$

$\Rightarrow$ It can be seen in a way that $y_n$ has been drawn from a normal distribution

of mean : $f(x_n, w)$ $\displaystyle\int_{}$ $x_n, w$ are given

$\quad$ Covariance : $\Sigma$

$$P\left(\frac{y_n}{x_n, w}\right) =$$

$$\frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left\{\frac{-1}{2}(y_n - \mu_n)^T \Sigma^{-1}(y_n - \mu_n)\right\}$$

where $\mu_n = f(x_n, w)$

\* Considering

N independent data points are drawn

Likelihood of collection of N data points

$$\prod_{i=1}^{N} P\left(y_i / x_i, \Sigma\right)$$

$$= \prod_{i=1}^{N} \left( \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left\{ -\frac{1}{2} (y_i - \mu_i)^T \Sigma^{-1} (y_i - \mu_i) \right\} \right)$$

where $\mu_i = f(x_i, w)$

\# Our goal is to find the parameter $(W)$,

so that the likelihood is maximum.

Assuming $\Sigma$ is constant

$$\Rightarrow w^* = \arg\max_{w} \prod_{i=1}^{N} P\left(y_i / x_i, \Sigma\right)$$

## log-likelihood

$$\log\left( \prod_{i=1}^{N} P\left( \sqrt[y_i]{x_i, \Sigma}\right)\right) =$$

$$-\frac{N}{2}\log|\Sigma| \underbrace{-\frac{Nd}{2}\log(2\pi)}_{constant} -\frac{1}{2}\sum_{i=1}^{N}(y_i-\mu_i)\Sigma^{-1}(y_i-\mu_i)$$

The objective function here is

$$\max\left\{\log\left(\prod_{i=1}^{N} P\left(\dot{y_i}|x_i,\Sigma\right)\right)\right\}$$

$$= \max\left\{-\frac{1}{2}\sum_{i=1}^{N}(y_i-\mu_i)\Sigma^{-1}(y_i-\mu_i)\right\}$$

$$= \min\left\{\sum_{i=1}^{N}(y_i-\mu_i)\Sigma^{-1}(y_i-\mu_i)\right\}$$

$$\therefore \; \underset{w}{w^*} = \arg\min_{w}\left(\sum_{i=1}^{N}(y_i-\mu_i)\Sigma^{-1}(y_i-\mu_i)\right)$$

optimum
w

$$\mu_i = f(x_i, w)$$

$$w^* = \arg\min_{w}\left(\sum_{i=1}^{N}(y_i-f(x_i,w))\Sigma^{-1}(y_i-f(x_i,w))\right)$$

MLE

For MLE : error function :

$$e(w) = \sum_{i=1}^{N} (y_i - f(x_i, w))^T \, \Sigma^{-1} \, (y_i - f(x_i, w)) \quad \rightarrow \textcircled{1}$$

For a NN using sum-of-squares

error function is :

$$e(w) = \frac{N}{2} \sum_{i=1}^{N} \| y_i - f(x_i, w) \|^2 \quad \left. \begin{array}{l} \text{Sum of} \\ \text{Squares} \end{array} \right. \quad \rightarrow \textcircled{2}$$

$\therefore$  for  $\Sigma = \sigma^2 I$  :  $I$ : identity matrix

$\rightarrow$ eq $\textcircled{1}$

$$e(w)_1 = \frac{1}{\sigma^2} \sum_{i=1}^{N} (y_i - f(x_i, w))^T \, I \, (y_i - f(x_i, w))$$

$$= \frac{1}{\sigma^2} \sum_{i=1}^{N} \| y_i - f(x_i, w) \|^2$$

$\therefore$ we can see that :

$\qquad$ $\cancel{\text{then}}$  $\underset{w}{\arg \min} (e_1(w)) = \underset{w}{\arg \min} (e_2(w))$

Both are equivalent if $\boxed{\Sigma = \sigma^2 I}$

Q6) Scale-space symmetry:

1)

Problem: Vanishing gradient problem

* But left layers are scaled

Let incoming layers are scaled by $\gamma$

& outgoing by $\frac{1}{\gamma}$

→ During Accumulation of gradients in
back-propagation:

→ Let $(\delta_l, \delta_l')$ are the gradients before
and after scaling

→ for before layer: $(\delta_{l-1}, \delta'_{l-1})$
    After    "    $(\delta_{l+1}, \delta'_{l+1})$

Gradient of weights

$$\Delta_l = \delta'_{l+1} \cdot a_l'$$

$$= \delta_{l+1} * \gamma * a_l$$

$$\boxed{|\Delta_l'| = \gamma \Delta_l}$$

Since, $l$ is changed

$\delta_{l+1} = \delta'_{l+1}$

$\delta_l = \frac{1}{\gamma} \delta_l'$

$\delta_{l-1} = \delta'_{l-1}$

Hence, $\gamma$ is scaled,

af $\gamma$ is very small, while computing
initial layers gradient $\longrightarrow 0$

Q6)

b) The permutation-symmetry in weight space.

Since every node is connected to all the nodes in the previous layer. This makes the weights & nodes / the neural network invariant to permutation of nodes.

⟹ . This leads to $ permutation-symmetry in weight space of each layer.

Consequences:

1) Gives rise to multiple-equal global minima of the loss function.

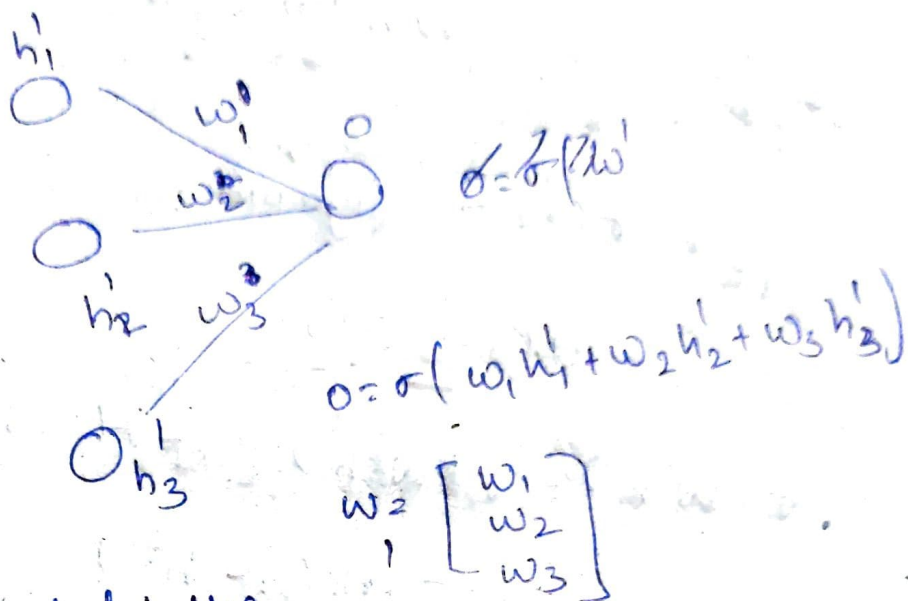2) Also creates saddle points on the path between these minima.

___

* For a multilayer network with $d-1$ hidden layers → each with 'n' neurons.

$n!$ : No. of possible permutation per layer

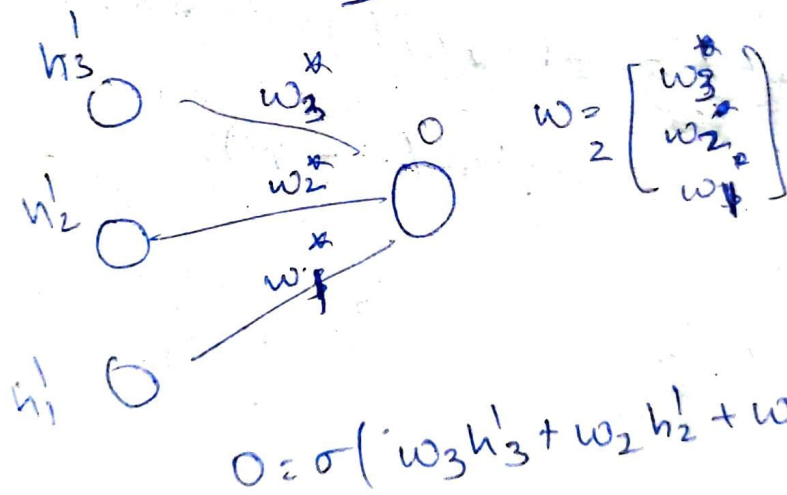∴ Total no. of equivalent configuration which

yield the same ~~los~~ error $= (n!)^{d-1}$

Eg:



$O = \sigma(\omega_1 h'_1 + \omega_2 h'_2 + \omega_3 h'_3)$

$W_2 \begin{bmatrix} \omega_1 \\ \omega_2 \\ \omega_3 \end{bmatrix}$

$\Rightarrow$ Permutated : the
nodes



$W = \begin{bmatrix} \omega_3^* \\ \omega_2^* \\ \omega_1^* \end{bmatrix}$

$O = \sigma(\omega_3 h'_3 + \omega_2 h'_2 + \omega_1 h'_1)$

· For $\bar{\omega}_1 = \begin{bmatrix} \omega_1 \\ \omega_2 \\ \omega_3 \end{bmatrix}$ ; $\bar{\omega}_2 = \begin{bmatrix} \omega_3 \\ \omega_2 \\ \omega_1 \end{bmatrix}$

$\rightarrow$ Yield the same activation

$\rightarrow$ Thus will yield same loss

- Hence they are permutation symmetric.