

Q.1

Initially there are 50% matches i.e. $w=0.5$
Homography has 8 degree of freedoms, so
number of points to consider at a time $n=4$

at K iterations, probability of bad model

$$= (1-w^n)^K$$

To get 95% chance of exact homography
means < 0.05 probability of bad model

$$(1-w^n)^K < 0.05$$

$$(1-(0.5)^4)^K < 0.05$$

$$K > \frac{\log(0.05)}{\log(1-(0.5)^4)}$$

$$K > 46.46$$

$$K=47$$

Q.3

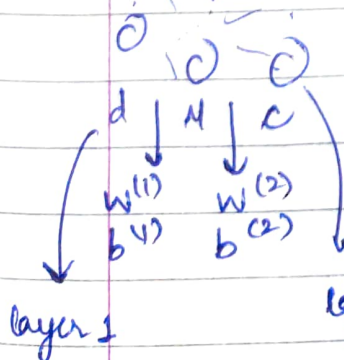
$$\Delta_{ij}^{(2)} = \Delta_{ij}^{(2)} + \delta_i^{(3)} + (a^{(2)})_j$$

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

Q.4



Total weights = $M \times d + C \times M$
Total biases = $M + C$.



$$\frac{\partial E}{\partial W^{(2)}} = \delta^{(3)} (a^{(2)})^T; \quad \frac{\partial E}{\partial W^{(1)}} = \delta^{(2)} (a^{(1)})^T$$

we need to calculate only these.

so, Total independent derivatives = $M + C$

Q-2

$$f(x) = \langle W^{(3)}, h^{(2)} \rangle = (W^{(3)})^T h^{(2)} \quad \left\{ \begin{array}{l} \bullet \text{ means} \\ \text{element wise} \\ \text{product} \end{array} \right.$$

$$\frac{\partial f(x)}{\partial W_{ij}^{(1)}} = \frac{\partial (W^{(3)})^T h^{(2)}}{\partial W_{ij}^{(1)}} = (W^{(3)})^T \frac{\partial h^{(2)}}{\partial W_{ij}^{(1)}} \quad \left\{ \begin{array}{l} h^{(2)} \\ = \sigma(W^{(2)} h^{(1)}) \end{array} \right.$$

$$\frac{\partial f(x)}{\partial W_{ij}^{(1)}} = (W^{(3)})^T \left[h^{(2)} (1 - h^{(2)}) \cdot \frac{\partial W^{(2)} h^{(1)}}{\partial W_{ij}^{(1)}} \right] \quad \left\{ \begin{array}{l} \frac{d\sigma(x)}{dx} \\ = \sigma(x)(1 - \sigma(x)) \end{array} \right.$$

$$\frac{\partial f(x)}{\partial W_{ij}^{(1)}} = (W^{(3)})^T \left[h^{(2)} (1 - h^{(2)}) \left[W^{(2)} \frac{\partial h^{(1)}}{\partial W_{ij}^{(1)}} \right] \right] \quad \left\{ \begin{array}{l} h^{(1)} \\ = \sigma(W^{(1)} x) \end{array} \right.$$

$$\frac{\partial f(x)}{\partial W_{ij}^{(1)}} = (W^{(3)})^T \left[h^{(2)} (1 - h^{(2)}) \cdot \left[W^{(2)} \left[h^{(1)} (1 - h^{(1)}) \cdot \frac{\partial W_{ij}^{(1)} x}{\partial W_{ij}^{(1)}} \right] \right] \right]$$

$$\frac{\partial f(x)}{\partial W_{11}^{(1)}} = (W^{(3)})^T \left[h^{(2)} (1 - h^{(2)}) \cdot \left[W^{(2)} \left[h_1^{(1)} (1 - h_1^{(1)}) x_1 \right] \right] \right]$$

$$\frac{\partial f(x)}{\partial W_{12}^{(1)}} = (W^{(3)})^T \left[h^{(2)} (1 - h^{(2)}) \cdot \left[W^{(2)} \left[\frac{h_1^{(1)} (1 - h_1^{(1)}) x_2}{0} \right] \right] \right]$$

$$\frac{\partial f(x)}{\partial W_{21}^{(1)}} = (W^{(3)})^T \left[h^{(2)} (1 - h^{(2)}) \cdot \left[W^{(2)} \left[\begin{array}{c} 0 \\ h_2^{(1)} (1 - h_2^{(1)}) x_1 \end{array} \right] \right] \right]$$

$$\frac{\partial f(x)}{\partial W_{22}^{(1)}} = (W^{(3)})^T \left[h^{(2)} (1 - h^{(2)}) \cdot \left[W^{(2)} \left[\begin{array}{c} 0 \\ h_2^{(1)} (1 - h_2^{(1)}) x_2 \end{array} \right] \right] \right]$$

Q.5

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

$$\Rightarrow y_n \sim N(f(x_n, w), \Sigma)$$

$$P(y_n = y | x_n) \propto e^{-\frac{1}{2}(y - f(x_n, w))^T \Sigma^{-1} (y - f(x_n, w))}$$

$$\text{Now } P(Y|X) = \prod_{i=1}^N P(y_i | x_i)$$

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

Loss is negative log likelihood.

$$L(w) = -\log P(Y|X)$$

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

$$w^* = \underset{w}{\operatorname{argmin}} L(w).$$

Σ is p.s.d so can be written as $S \Lambda S^T$ [diagonalization]

$$\Rightarrow L(w) = \sum_{i=1}^N (y_i - f(x_i, w))^T S^{-T} \Lambda^{-1} S^{-1} (y_i - f(x_i, w))$$

$$= \sum_{i=1}^N (S^{-1} y_i - S^{-1} f(x_i, w))^T \Lambda^{-1} (S^{-1} y_i - S^{-1} f(x_i, w))$$

$$= \sum_{i=1}^N (y_i' - f'(x_i, w))^T \Lambda^{-1} (y_i' - f'(x_i, w))$$

here $y_i' = S^{-1} y_i$ and $f'(x_i, w) = S^{-1} f(x_i, w)$

so the problem becomes simple regression. i.e

$$y_n' = f'(x_n, w) + \varepsilon_n' \quad \varepsilon_n' \sim N(0, \Lambda)$$

$$\Lambda = \operatorname{diag}[\sigma_1^2, \sigma_2^2, \dots, \sigma_d^2]$$

Q.6 a) Let $\mathcal{L}(\mathcal{D}, w_1, w_2)$ be the loss function of neural network on dataset \mathcal{D} with two weights w_1 and w_2 .

Scale-symmetry exists i.e.

$$\mathcal{L}(\mathcal{D}, w_1, w_2) = \mathcal{L}(\mathcal{D}, \gamma w_1, \frac{1}{\gamma} w_2) \quad - \text{Eq. (1)}$$

~~$\frac{\partial \mathcal{L}(\mathcal{D}, w_1, w_2)}{\partial w_1}$~~ Now, taking / calculating gradients for updates.

$$\begin{aligned} \frac{\partial \mathcal{L}(\mathcal{D}, \gamma w_1, \frac{1}{\gamma} w_2)}{\partial w_1} &= \frac{\partial \mathcal{L}(\mathcal{D}, \gamma w_1, \frac{1}{\gamma} w_2)}{\partial (\gamma w_1)} \cdot \frac{\partial (\gamma w_1)}{\partial w_1} \\ &= \frac{\partial \mathcal{L}(\mathcal{D}, w_1, w_2)}{\partial w_1} \cdot \gamma \end{aligned}$$

Similarly,

$$\frac{\partial \mathcal{L}(\mathcal{D}, \gamma w_1, \frac{1}{\gamma} w_2)}{\partial w_2} = \frac{\partial \mathcal{L}(\mathcal{D}, w_1, w_2)}{\partial w_2} \cdot \frac{1}{\gamma}$$

$\left\{ \begin{array}{l} \because \text{using} \\ \text{Eq. (1)} \end{array} \right.$

Hence gradients w.r.t to w_1 and w_2 are also scaled with γ and $\frac{1}{\gamma}$ respectively.

~~∴~~ This will create uneven updates as if γ is high w_1 gradient will explode and w_2 gradient will vanish. and Similarly if γ is low w_1 grad. can vanish and w_2 will explode.

This will increase instability of training as updates will be uneven.

b) Permutation Symmetry:

In any Multi layered network d with d layers and n neurons. There will be $(n!)^d$ permutations of hidden units.

Given any permutation, we can find weight configuration that gives same output. ~~So~~ And it is more-likely these configurations are local-minima. Thus ~~but~~ we could encounter many local-minimas during training. We need to initialize network carefully for better performance.