

LF12 Computer Vision (2019)

Q1 (a) Smoothing to remove high frequency noise which is amplified
(i) by differentiation:

$$\frac{dI}{dx} \xrightarrow{\text{F.T.}} j\omega I(j\omega)$$

↑ amplification of high spatial frequencies (2)

$$(ii) \quad g_{\sigma}(u) \times g_{\sigma}(v) = g_{\sigma}(u, v) = \frac{1}{2\pi\sigma} e^{-\frac{(u^2+v^2)}{2\sigma^2}}$$

↑
2D gaussian

$$\therefore S(x, y) = \sum_{-n}^n \sum_{-n}^n g_{\sigma}(u, v) I(x-u, y-v)$$

Advantage of 1D scheme vs 2D scheme (efficiency)

$$2N \text{ vs } N^2 \quad \left(\frac{2N}{N^2} \right) \text{ operations; where } N = 2n+1 \quad (\text{size of kernel})$$

(2)

1(b)(i)

$$S(x+\Delta x, y+\Delta y) \approx S(x, y) + \frac{\partial S}{\partial x} \Delta x + \frac{\partial S}{\partial y} \Delta y + \frac{\Delta x^2}{2} \frac{\partial^2 S}{\partial x^2} + \frac{\Delta y^2}{2} \frac{\partial^2 S}{\partial y^2} + \frac{\Delta x \Delta y}{2} \frac{\partial^2 S}{\partial x \partial y}$$

Consider $\Delta x=1, \Delta y=1$

$$\begin{aligned} \therefore S(x+1, y) &= S(x, y) + \frac{\partial S}{\partial x} + \frac{1}{2} \frac{\partial^2 S}{\partial x^2} + \dots \\ S(x-1, y) &= S(x, y) - \frac{\partial S}{\partial x} + \frac{1}{2} \frac{\partial^2 S}{\partial x^2} + \dots \end{aligned}$$

$$\therefore \frac{\partial^2 S}{\partial x^2} \approx \frac{S(x+1, y) - 2S(x, y) + S(x-1, y)}{1} \quad (2)$$

Similarly

$$\frac{\partial^2 S}{\partial y^2} \approx S(x, y+1) - 2S(x, y) + S(x, y-1)$$

$$\therefore \frac{\partial^2 S}{\partial x^2} : \begin{array}{|c|c|c|} \hline +1 & -2 & +1 \\ \hline \end{array} \quad \text{kern} \quad h(u) \quad \frac{\partial^2 S}{\partial x^2} = S''_{xx} = \sum_{-1}^1 h(u) S(x-u)$$

$$\frac{\partial^2 S}{\partial y^2} : \begin{array}{|c|} \hline 1 \\ \hline -2 \\ \hline 1 \\ \hline \end{array} \quad h(v) \quad \frac{\partial^2 S}{\partial y^2} = S''_{yy} = \sum_{-1}^1 h(v) S(x, y-v) \quad (2)$$

(ii) $\nabla^2 S = \frac{\partial^2 S}{\partial x^2} + \frac{\partial^2 S}{\partial y^2}$

0	1	0
1	-4	1
0	1	0

(2)

Q1(c) Image pyramid formed from $S(x, y, \sigma_i)$ where $S(x, y, \sigma_i) = I(x, y) * G_{\sigma_i}$

$$\begin{aligned} \text{(i). } \nabla^2 S(x, y) &= \nabla^2 G_{\sigma}(x, y) * I(x, y) \\ &\approx (G_{\sigma_{i+1}} - G_{\sigma_i}) * I(x, y) \\ &\approx \underline{S(x, y, \sigma_{i+1}) - S(x, y, \sigma_i)} \quad (1) \end{aligned}$$

if $\sigma_{i+1} \approx \frac{1.2-1.6}{\text{range}} \sigma_i$

i.e. Laplacian can be approximated by D.O.G operator
Difference of smoothed images in image pyramid if σ chosen correctly (1)

$$\text{Compute } S(x, y, \sigma_i) = G_{\sigma_i} * I(x, y)$$

for $\sigma_i = \sigma_0 2^{i/s}$ Need $s = 2 \text{ or } 3 \text{ or } 4$ s images per octave

$$\text{— incremental blur } \sigma_k = \sigma_i \sqrt{2^{\frac{2k}{s}} - 1}$$

— subsample after σ doubles $i = s, 2s, \dots$
($\frac{1}{4}$ image size)

$$\text{— } \nabla^2 S(x, y, \sigma_i) \approx S(x, y, \sigma_{i+1}) - S(x, y, \sigma_i) \quad (1)$$

Look for max/min of $\nabla^2 S$ by evaluating 26 neighbours in scale-space differences. Gives position (x, y) and scale, σ_i .

(ii) Look for max/min of $\nabla^2 S$ in LOG images. Position and scale, σ_i (x, y) (1)

Sample 16×16 pixels at $S(x, y, \sigma_i)$. Compute VS for each pixel
Histogram (use interpolation) to 10° bins and smooth with gaussian.

Find max of histogram, $\theta_i = \underline{\text{dominant orientation}}$. (2)

Q1(c)(iii)

SIFT — 16×16 patch is normalized to scale and orientation by blob detection / orientation histogram (2D-viewpoint change)

— use gradients, $\nabla S(x, y, \sigma_i)$ to encode 2D shape (robust to lighting changes)

— normalize 128D vector to unit length and truncate any values > 0.2
(invariance to lighting and specularities)

— Local histograms of gradients give invariance to small viewpoint changes and occlusion.

(3)

Limitations — fail at occluding boundaries (strong edges)

— can not cope with matching under large viewpoint changes (es $> 30^\circ$) from perspective

Q2(a)(i)

Algebraic — homogeneous co-ordinates

$$\begin{matrix} \text{homogeneous} & & \text{image plane} \\ \begin{pmatrix} su \\ sv \\ s \end{pmatrix} & \sim & \begin{pmatrix} u \\ v \end{pmatrix} \\ \text{3D} & & \text{2D} \end{matrix}$$

$$\underline{\tilde{w}} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} \text{ where } s = x_3 \quad \underline{w} = \begin{pmatrix} x_1/x_3 \\ x_2/x_3 \end{pmatrix}$$

Geometric — ^(inverse) depth scaling (distance along optical axis)

$$s = p_{31}X + p_{32}Y + p_{33}Z + p_{34} = Z_{\text{camera}} \quad (2)$$

$$(ii) \quad \begin{matrix} 3 \times 1 & & 3 \times 3 & & 3 \times 4 & & 4 \times 4 \\ \begin{pmatrix} su \\ sv \\ s \\ \underline{\tilde{w}} \end{pmatrix} & = & \begin{bmatrix} k_u & 0 & u_0 \\ 0 & k_v & v_0 \\ 0 & 0 & 1 \end{bmatrix} & \begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} & \begin{bmatrix} R & | & T \\ \hline 0 & 0 & 0 & | & 1 & 1 \end{bmatrix} & \begin{bmatrix} x_i \\ y_i \\ z_i \\ 1 \end{bmatrix} \\ & & \text{4 dof.} & \text{1} & \text{6 dof.} & & \underline{\tilde{X}} \end{matrix}$$

$\tilde{w} = P \tilde{X}$ 11 parameters of 10 d.o.f. since $\alpha_u = f k_u$
 $\alpha_v = f k_v$

$$(iii) \quad \begin{bmatrix} p_{11} & p_{12} & p_{13} \\ p_{21} & p_{22} & p_{23} \\ p_{31} & p_{32} & p_{33} \end{bmatrix}^T = (KR)^T = R^T K^T \quad (4)$$

Decompose by QR to give
 $Q = R^T$ and $R = K^T$ (4)
 orthogonal ~~upper~~ upper triangular

$$\underline{T} = K^{-1} \begin{bmatrix} p_{14} \\ p_{24} \\ p_{34} \end{bmatrix} \quad (1)$$

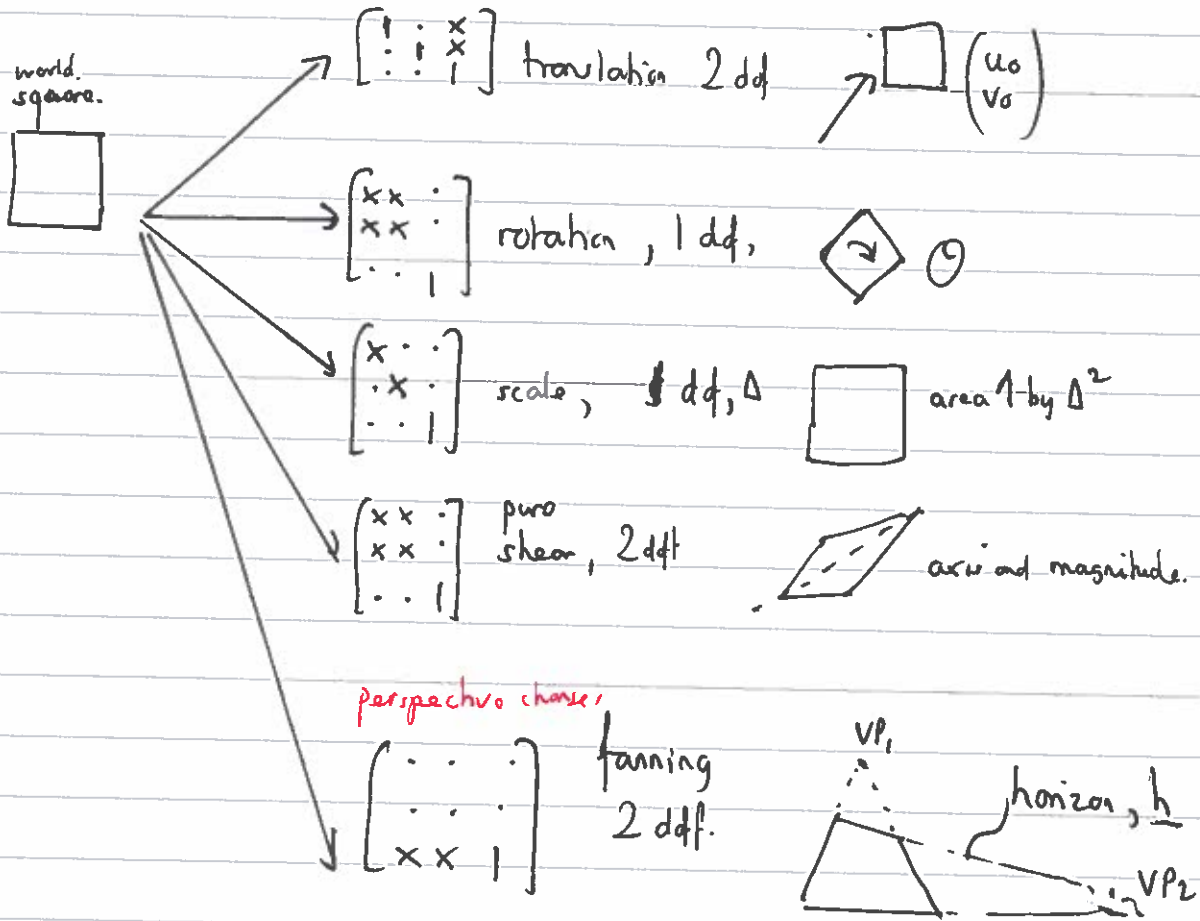
Q2 (b) Look at $Z=0$, X-Y plane.

(i)

$$\begin{matrix} \begin{bmatrix} s_u \\ s_v \\ s \end{bmatrix} \\ 3 \times 1 \end{matrix} = \begin{matrix} \begin{bmatrix} p_{11} & p_{12} & p_{13} \\ p_{21} & p_{22} & p_{23} \\ p_{31} & p_{32} & p_{33} \end{bmatrix} \\ 3 \times 3 \end{matrix} \begin{matrix} \begin{bmatrix} X \\ Y \\ 1 \end{bmatrix} \\ 3 \times 1 \end{matrix}$$

9 parameters up to arbitrary scale \therefore 8 d.o.f. (2)

(i)



(2)

Q2(b)(iii)

Consider lines // to X at ∞

$$V_{p_x} = \begin{bmatrix} p_{11} \\ p_{12} \\ p_{13} \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} p_{11} \\ p_{12} \\ p_{13} \end{bmatrix} \quad X=\infty$$

Consider // lines to Y axis at ∞

$$V_{p_y} = \begin{bmatrix} p_{21} \\ p_{22} \\ p_{23} \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} p_{21} \\ p_{22} \\ p_{23} \end{bmatrix} \quad Y=\infty$$

2

These points lie on horizon, \underline{l} , such that $\underline{l} \cdot \underline{w} = 0$

Horizon \underline{l} must satisfy: $\underline{l} \cdot V_{p_x} = 0$ and $\underline{l} \cdot V_{p_y} = 0$

$$\therefore \underline{l} = \begin{vmatrix} i & j & k \\ p_{11} & p_{12} & p_{13} \\ p_{21} & p_{22} & p_{23} \end{vmatrix}$$



$$\therefore \underline{l} = \begin{pmatrix} p_{12}p_{23} - p_{13}p_{22} \\ p_{13}p_{21} - p_{11}p_{23} \\ p_{11}p_{22} - p_{21}p_{12} \end{pmatrix}$$

line in image: (horizon)

$$\therefore (p_{12}p_{23} - p_{13}p_{22})u + (p_{13}p_{21} - p_{11}p_{23})v + \begin{pmatrix} p_{11}p_{22} \\ -p_{21}p_{12} \end{pmatrix} = 0$$

2

Q2(c) AR need to recover 3D projection matrix in real-time

$$- \underline{w} = \begin{bmatrix} P \\ 3 \times 4 \end{bmatrix} \begin{bmatrix} \tilde{x} \end{bmatrix}$$

- Can calibrate from known 3D points ^(x) and their correspondences (w).
 $N \geq 6$ for 3D or $N \geq 4$ for points on the plane
 $A_p = 0$

- Best to assume known K , internal parameters, which are fixed for mobile phone type

- Can use standard AR tool kit marker for calibration (eg Zappar code)

- Fast solution needed. Find "corners" or FAST features in image and track by normalized correlation. Use plane features.

- Estimate homography $H = K \begin{bmatrix} r_1 & r_2 & t_1 \end{bmatrix}$ and $\underline{\Sigma}_3 = \underline{\Sigma}_1 \underline{\Sigma}_2$
 decompose to get new pose \underline{t} and \underline{R} .
 (approx orientation is given by IMU and gyroscopes)

- Render CG using projection matrix, $K[R|t]$.

(4)

Q3

(a)

(i)

$$\underline{X}' = R \underline{X} + \underline{I}$$

The rays, \underline{X}' , \underline{X} and \underline{I} are co-planar.

\therefore taking cross product with \underline{I} and dot product with $\underline{X}' \Rightarrow$ triple scalar prod = 0

$$\underline{X}' \cdot [\underline{I} \times R] \underline{X} = 0$$

$$\text{E if } \underline{I} \times R = [\underline{T}_*] R$$

$$= \begin{bmatrix} 0 & -T_z & T_y \\ T_z & 0 & T_x \\ -T_y & T_x & 0 \end{bmatrix} R.$$

skew-symmetric matrix

We can express in terms of rays \underline{p}' and \underline{p} since $\underline{p}' \parallel \underline{X}'$ and $\underline{p} \parallel \underline{X}$

$$\therefore \underline{p}'^T E \underline{p} = 0$$

In terms of image coordinates: $\underline{w}' = K' \underline{p}'$ and $\underline{w} = K \underline{p}$

\therefore

$$\underline{w}'^T K'^{-T} E K^{-1} \underline{w} = 0$$

$$\underline{w}'^T \underbrace{[K'^T T_*] R [K^{-1}]}_F \underline{w} = 0 \quad (5)$$

F

4 matrices

Q3(a)(ii)

Epipolar constraint

$$\underline{\omega}'^T F \underline{\omega} = 0$$

- Points on right view must lie on epipolar line $\underline{L}' \cdot \underline{\omega}' = 0$
or $\underline{\omega}'^T \underline{L} = 0$

$$\therefore \underline{L}' = F \underline{\omega}$$

i.e. line corresponding to image point (left) $\underline{\omega}$

- Epipole in right view, \underline{e}' lies on all epipolar lines \underline{L}'

$$\therefore \underline{L}' \cdot \underline{e}' = 0 \text{ for all } \underline{L}' = F \underline{\omega}$$

$$\therefore \underline{\omega}^T F^T \underline{e}' = 0$$

$$\underline{F^T \underline{e}' = 0} \quad \text{null-space of } F^T$$

(3)

Q3(b)

(i) Matching (corresponding) points:

For each image:

- Find interest points such as blobs (scale, position and orientation)
- SIFT descriptor. (128D vector)

- Match in tight view by NN and k-d search tree:
(Compute dot product)

[Alternative is to find Harris corner and match by normalized cross-correlation]

(2)

(ii) Need 8 pt correspondences with a linear method or 7 correspondences if which to include Rank $F = 2$, i.e. $\det F = 2$.

Solve $Af = 0$ by minimizing such that $\|f\| = 1$ $\lambda_1 \leq \frac{f^T A^T A f}{f^T f} < \lambda_2$

(2)

(iii) Least-squares followed by rank $F = 2$ by projection or SVD $\lambda_3 = 0$.
RANSAC to remove incorrect matches.

(2)

Non-linear optimisation of 7 parameters to minimize the algebraic error.

Q3(c)(i) Estimate F
 Comput. $E = K^T F K = T_x R = U \Lambda V^T$ (SVD decom)

$$T_x = U \begin{bmatrix} 0 & 1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} U^T \quad R = U \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} V^T$$

- skew-symmetric
 - unknown scale $|T|$?

- orthonormal matrix

This decomposition gives 4 solutions $\pm T$, R and R^T
 Check by computing depths, must be positive.

Left with unknown scale $|T|$. Need a known length in world.

(4)

$$\underline{P}_L = K \left[\mathbb{I} \mid 0 \right] \quad \text{and} \quad \underline{P}_R = K' \left[R \mid T \right]$$

(c)(ii) We have. $\underline{w}' = K'(R \ T) \underline{X}$ and $\underline{w} = K[\mathbb{I} \mid 0]$
 Find \underline{X} by triangulation
 i.e. by looking at stereo correspondance \underline{w}' and \underline{w}

Each point (visible) gives 2 equations in 3 unknowns (X, Y, Z ,
 i.e. 2 planes defining a ray)

\therefore 4 equations in X, Y, Z . Solve by least-squares

(2)

a) A complete answer will include 4 of the following points

- Stage 1 : 1) Compute high level feature description of both images
 2) both images passed through the same CNN so there are fewer parameters to learn & so ordering of input images does not matter
- Stage 2 : 3) compute similarity of the two images via the weighted squared distance and add a bias.
 4) Again order of inputs does not matter.
 5) Bias required so that when two identical images are passed in, the activation will be large & positive (resulting in $p(y=1|z_1, z_2) \approx 1$)
- Stage 3 : 6) converts real valued activations that lie in $(-\infty, \infty)$ to probabilities that lie between $(0, 1)$.

b) This & the next question or booklet on disguise being identical to logistic regression covered in lectures

$$\begin{aligned} \mathcal{L}(\theta; \mathcal{W}) &= \log \prod_n p(y^{(n)} | z_1^{(n)}, z_2^{(n)}, \mathcal{W}) \\ &= \sum_n \left[y_n \log \sigma(a^{(n)}) + (1-y_n) \log (1-\sigma(a^{(n)})) \right] \end{aligned}$$

$$c) \frac{d\mathcal{L}(\theta; \mathcal{W})}{d\mathcal{W}_d} = \sum_n \frac{d\mathcal{L}}{da^{(n)}} \frac{da^{(n)}}{d\mathcal{W}_d} \quad (\text{chain rule})$$

$$\begin{aligned} \frac{d\mathcal{L}}{da^{(n)}} &= \frac{y_n}{\sigma(a^{(n)})} \frac{d\sigma(a^{(n)})}{da^{(n)}} - \frac{(1-y_n)}{1-\sigma(a^{(n)})} \frac{d\sigma(a^{(n)})}{da^{(n)}} \\ &= \frac{y_n(1-\sigma(a^{(n)})) - (1-y_n)\sigma(a^{(n)})}{\sigma(a^{(n)})(1-\sigma(a^{(n)}))} \frac{d\sigma(a^{(n)})}{da^{(n)}} \end{aligned}$$

$$= \frac{y_n - \sigma(a^{(n)})}{\sigma(a^{(n)})(1-\sigma(a^{(n)}))} \frac{d\sigma(a^{(n)})}{da^{(n)}} \quad \text{cancel}$$

$$\begin{aligned} \frac{d\sigma(a^{(n)})}{da^{(n)}} &= \frac{d}{da^{(n)}} \frac{1}{1+e^{-a^{(n)}}} = -1 \times \left(\frac{1}{1+e^{-a^{(n)}}} \right)^2 \times \left(-e^{-a^{(n)}} \right) \\ &= \sigma(a^{(n)})^2 \left(\frac{1}{\sigma(a^{(n)})} - 1 \right) \\ &= \sigma(a^{(n)})(1-\sigma(a^{(n)})) \end{aligned}$$

$$\frac{da^{(n)}}{d\mathcal{W}_d} = \begin{cases} 1 & \text{if } d=0 \\ h_d - h_{-d} & \text{if } d>0 \end{cases}$$

$$\therefore \frac{d\mathcal{L}}{d\mathcal{W}_d} = - \sum_n \left[\sigma(a^{(n)}) - y^{(n)} \right] \left[\delta_{d,0} + (1-\delta_{d,0}) [h_d - h_{-d}] \right]$$

d) Use gradient ascent of log likelihood:

initialise \mathcal{W} & then use:

$$\underline{\mathcal{W}}^{(\text{new})} = \underline{\mathcal{W}}^{(\text{old})} + \mathcal{N} \frac{\hat{d}\mathcal{L}}{d\mathcal{W}}$$

\mathcal{N} = learning rate

NB. move in direction of gradient here as maximising

$\frac{\hat{d}\mathcal{L}}{d\mathcal{W}}$ = estimator for gradient produced by selecting \mathcal{M} training data points at random @ each iteration

$$= - \frac{\mathcal{N}}{\mathcal{M}} \sum_m \left[\sigma(a^{(m)}) - y^{(m)} \right] \left[\delta_{d,0} + (1-\delta_{d,0}) [h_d - h_{-d}] \right]$$

Can also discuss: initialisation, adapting the learning rate, momentum etc.

Module 4F12 (Computer Vision) Assessor's Comments

1. **Gaussian smoothing, bandpass filtering and SIFT.** Attempted by 69/77 Part IIB candidates, average mark 13.5/20.

The first part was a very straightforward question covering convolution with low pass filters and was well answered by most candidates. The second and third parts were more challenging. There were very few correct attempts at obtaining the laplacian kernel in (b). Candidates also struggled to describe how SIFT achieved invariance to lighting and viewpoint and how it was encoding 2D shape in (c).

2. **Perspective projection and camera calibration.** Attempted by 64/77 candidates, average mark 14.2/20.

A question covering perspective projection. Well answered by most candidates. A few marks were lost on the recovery of the horizon (b). The only challenging component was the application of the theory in (a) and (b) to a practical example of calibration for AR.

3. **Epipolar geometry and stereo vision.** Attempted by 58/77 candidates, average mark 14.3/20.

The derivation in (a) proved to be very easy for most candidates. Marks were lost in (b) for insufficient details in how to find point correspondences using features and descriptors and the estimation of the fundamental matrix in presence of outliers and noise (b). Part (c) was found to be easy for most candidates.

4. **Object detection with convolutional neural networks.** Attempted by 40/77 candidates, average mark 14/20.

Many candidates that attempted this question made excellent progress.