

4 F12 Computer Vision (2016)

Q1(a) (i) Low pass filter used $g_{\sigma}(x,y) = \frac{1}{2\pi\sigma^2} e^{-\frac{(x^2+y^2)}{2\sigma^2}}$ 1.

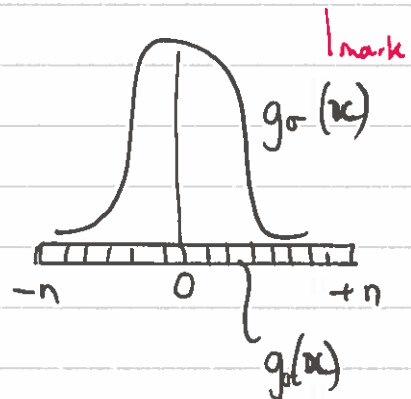
(1) marks

$$= g_{\sigma}(x) g_{\sigma}(y)$$

(2) $S(x,y) = \sum_{-n}^n \sum_{-n}^n I(x-u, y-v) g_{\sigma}(v) g_{\sigma}(u)$ 2 marks

where kernel size is $(2n+1)$

$2n+1$ samples from a 1D gaussian



(1)

(ii). $\nabla S = \nabla (g_{\sigma}(x,y) * I(x,y)) = (S_x, S_y)$ where $S_x = \frac{\partial S}{\partial x}$
 $S_y = \frac{\partial S}{\partial y}$
 2D vector

(1) $\frac{\partial S(x,y)}{\partial x} \approx \frac{S(x+1, y) - S(x-1, y)}{2}$ 1D kernel $\begin{bmatrix} \frac{1}{2} & 0 & -\frac{1}{2} \end{bmatrix}$

(1) $\frac{\partial S(x,y)}{\partial y} \approx \frac{S(x, y+1) - S(x, y-1)}{2}$ $\begin{bmatrix} \frac{1}{2} \\ 0 \\ -\frac{1}{2} \end{bmatrix}$ $\begin{matrix} 0 \rightarrow x \\ \downarrow y \end{matrix}$

1. (a)(iii)

Generate discrete set of smoothed images $S(x, y, \sigma_i)$ logarithmically spaced with s images per octave

$$\sigma_i = 2^{i/s} \sigma_0 \quad \text{and} \quad \sigma_{i+1} = \sigma_i 2^{1/s} \quad (1)$$

— implement all gaussian blurs by 2 x 1D convolutions with 'small' filter kernels

— apply incremental blur

$$g(\sigma_{i+1}) = g(\sigma_i) * g(\sigma_k)$$

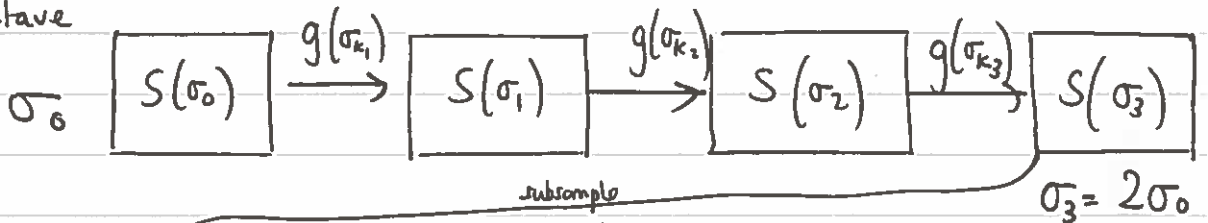
$$\text{where } \sigma_k = \sigma_i \sqrt{2^{2/s} - 1} \quad (1)$$

— sub-sample (reduce image size by $\frac{1}{4}$) after scale doubles (i.e. $\sigma_i = 2\sigma_0$) to generate next octave

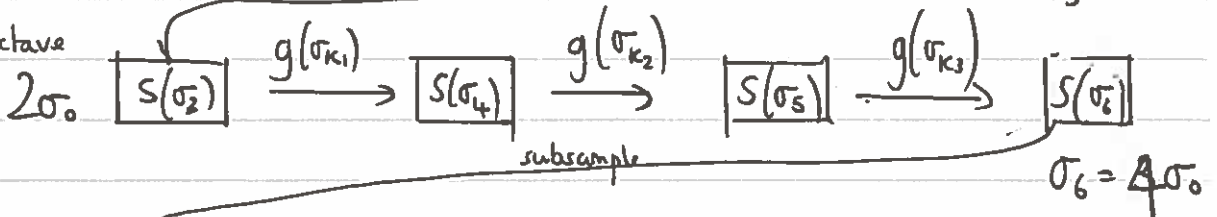
— apply same $g(\sigma_k)$ kernel to next octave image

See example with $s=3$ ($s=3$ used by LOWE in SIFT implementation).

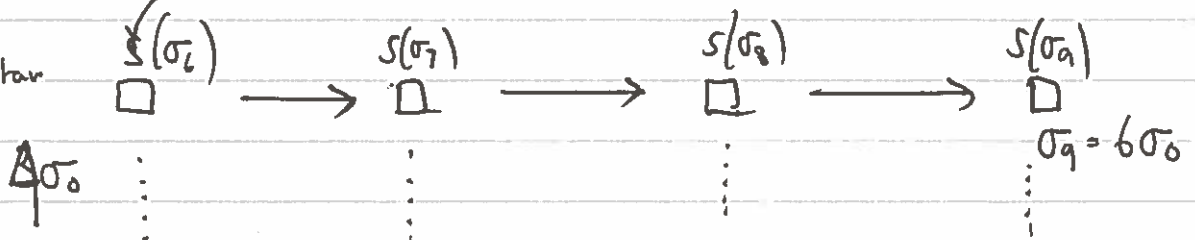
1st octave



2nd octave



3rd octave



top of pyramid 16x16

11x11

11x11

11x11

(b) (i). Band-pass filtering with the Laplacian of a Gaussian

$$\nabla^2 \left(g(x, y, \sigma_i) * I(x, y) \right) = \nabla^2 g(\sigma_i) * I = \nabla^2 S_i$$

(2) $\approx S(\sigma_{i+1}) - S(\sigma_i)$

i.e. Approximated by ^{difference} neighbouring images in some octave for appropriate s

(1) DOG - difference of gaussian image. if $\sigma_{i+1} \approx 1.2 \rightarrow 1.5 \sigma_i$

This is a band-pass filter since difference of gaussian in space and frequency

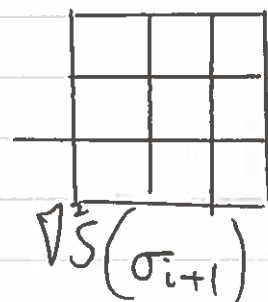
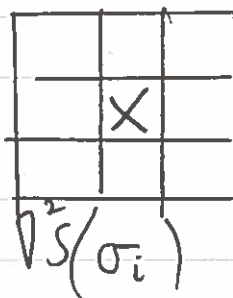
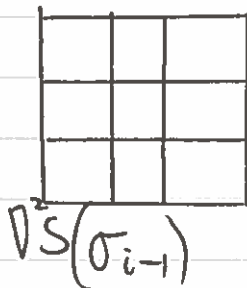


Find blob-like centre by looking for max/min of $\nabla^2 g(\sigma_i) * I$ resp

Search over scale-space representation - use DOG images for efficiency

(1) Need to check (at maximum) 26 neighbors in image beside interest pt.

use
DOG
images



(1) Position (x, y) and scale (σ_i) correspond to local max/min in pyramid of difference of gaussian images.

(b)(ii) Dominant orientation of feature in vicinity of blob centre:

- (1) - Sample 16×16 pixels from $S(x, y, \sigma_i)$ at ^{interest point} centre (x_i, y_i)
 - (1) - Compute $\nabla S(x, y, \sigma_i)$ - i.e. gradient magnitude and orientation
 - Produce a histogram of gradient magnitude against orientation - bin size 10°
 - (1) Bin gradient magnitudes $\times g_{\text{weighted}}(1.5\sigma_i)$ in 10° bins from 0° to 360°
 - (1) - Smooth histogram and find dominant "peak".
This is orientation of feature point and used to give orientation invariance
 - Can produce more accurate orientation by fitting a parabola to nearby bins and find θ
-

Q2.(a)(i). Assumptions: - pin-hole camera, planar perspective projection
 - no non-linear distortion (eg radial lens)

(2)

$$\left. \begin{array}{l} \text{Image } u = \frac{x_1}{x_3} \\ v = \frac{x_2}{x_3} \\ \text{Center world 3D } \begin{array}{l} X = \frac{x_1}{x_4} \\ Y = \frac{x_2}{x_4} \\ Z = \frac{x_3}{x_4} \end{array} \end{array} \right\} \text{homogeneous co-ordinates}$$

$$(ii) \quad u_i = \frac{p_{11}X_i + p_{12}Y_i + p_{13}Z_i + p_{14}}{p_{31}X_i + p_{32}Y_i + p_{33}Z_i + p_{34}}$$

$$v_i = \frac{p_{21}X_i + p_{22}Y_i + p_{23}Z_i + p_{24}}{p_{31}X_i + p_{32}Y_i + p_{33}Z_i + p_{34}} \quad (2)$$

We can rearrange to be linear in unknown p_{ij} and assemble into matrix below:

$$\begin{array}{c} 2 \times 12 \\ \left[\begin{array}{cccc|cccc} X_i & Y_i & Z_i & 1 & 0 & 0 & 0 & 0 & -u_i X_i & -u_i Y_i & -u_i Z_i & -u_i \\ 0 & 0 & 0 & 0 & X_i & Y_i & Z_i & 1 & -v_i X_i & -v_i Y_i & -v_i Z_i & -v_i \end{array} \right] \begin{array}{c} p_{11} \\ p_{12} \\ p_{13} \\ p_{14} \\ p_{21} \\ p_{22} \\ p_{23} \\ p_{24} \\ p_{31} \\ p_{32} \\ p_{33} \end{array} = 0 \end{array} \quad (2)$$

Q2(a)(ii)

- Need $N > 6$ 3D points of known (X_i, Y_i, Z_i) which span region/volume of interest at different depths (ie. NOT coplanar). Large FOV
 Measure image (u_i, v_i) correspondences.

(2)

- $N > 6$ stack measurement equations in (a)(i) to give a $2N \times 12$ matrix, A

(2) Solve $A p \approx 0$ where $p = \begin{bmatrix} p_{11} \\ \vdots \\ p_{2N} \end{bmatrix}$

- Solve by least-squares or find singular vector corresponding to smallest singular value (SVD).

$$\text{ie. } \lambda_1 < \frac{p^T A^T A p}{p^T p} < \lambda_{12}$$

This is a linear soln and only approximate.

- Use linear solution to initialize on non-linear optimization (gradient descent)

(2)

$$\min_{p_{ij}} \sum (u_i - \hat{u}_i)^2 + (v_i - \hat{v}_i)^2$$

where \hat{u}_i and \hat{v}_i are projected co-ordinates using p_{ij}
 and (u_i, v_i) is the measured image co-ordinates

Q2a(iv)
$$P = \begin{matrix} 3 \times 4 & 3 \times 3 & 3 \times 4 \\ K & [R | T] \\ & = [KR | KT] \\ & \begin{matrix} 3 \times 3 & 3 \times 1 \end{matrix} \end{matrix}$$

(1)

We can decompose 9 elements of KR by QR decomposition

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

$\begin{matrix} \text{QR decomposition} \\ \swarrow \quad \searrow \\ \text{orthogonal} \quad \text{upper diagonal} \end{matrix}$

Obtain position I

$$I = K^{-1} \begin{bmatrix} p_{14} \\ p_{24} \\ p_{34} \end{bmatrix}$$

(1)

Q2b(i) Weak perspective

$$Z_c = Z_A \quad (\text{orthogonal projection to plane at } Z_A)$$

(1)

$$\therefore x = \frac{f X_c}{Z_A} \quad \text{and} \quad y = \frac{f Y_c}{Z_A}$$

$$\therefore \begin{bmatrix} s_u \\ s_v \\ s \end{bmatrix} = \begin{bmatrix} K \\ \begin{matrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 0 & Z_A \end{matrix} \end{bmatrix} \begin{bmatrix} R & T \\ \text{ooo} & 1 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

weak perspective

$$\begin{bmatrix} s_u \\ s_v \\ s \end{bmatrix} = \begin{matrix} 3 \times 4 \\ \begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ 0 & 0 & 0 & p_{34} \end{bmatrix} \end{matrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} \quad \text{;} \quad \begin{matrix} \text{image} \\ \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} 2 \times 4 \\ 2 \times 4 \end{bmatrix} \end{matrix} \begin{matrix} \text{world} \\ \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} \end{matrix}$$

(3)

$s = \text{constant}$

Weak perspective projection matrix.

Q2 (ii) Projection equations become linear.

$$u_i = p_{11} X_i + p_{12} Y_i + p_{13} Z_i + p_{14}$$

$$v_i = p_{21} X_i + p_{22} Y_i + p_{23} Z_i + p_{24}$$

marks
(2)

Solve optimally by ~~linear~~ least squares. Need only 4 pts to calibrate.

Q3(a)(i) - Sample 16×16 pixels from image $S(x, y, \sigma_i)$ and dominant orientation of interest point.

(1) - Compute gradient, ∇S in 16×16 pixels

(1) - weight by g_0 where scale is $\frac{1}{2}$ window of pixels (8), to weight dist from blob centre.

- produce 4×4 cells (16 cells)

(1) - in each cell produce a HOG at 45° bins (8 dir's)
ie. add gradient magnitudes to each bin (8 sums?)



- concatenate to give 128D vector (8×16).

- normalize to unit length

(1) - remove effect of outliers such as highlights by truncating all values > 0.2 (ie max entry = 0.2)

- renormalize to unit vector of 128D (SIFT)

(ii) Invariant to lighting (gradients and normalization) and saturation (truncation)
Encodes 2D shape in vicinity of blob centre
(ie. edges and their orientations and spatial distribution)
(2) Invariant to small translations (histogramming, pooling).

Q3(a)(iii). Each interest point produces a 128D descriptor vector

which is invariant to scale, orientation and lighting

— Finding a correspondence is equivalent to NN matching of these vectors either by correlation of ^{nearest neighbour} euclidean distance.

(1)

— Compute euclidean distance between \underline{x} and nearest neighbour \underline{x}_1
or Accept as a match if $\frac{\underline{x} \cdot \underline{x}_2}{\underline{x} \cdot \underline{x}_1} < 0.7$ (no need for a global threshold)

where \underline{x}_2 is second nearest neighbour.

(2) Efficient NN by using a k-D tree (binary search). Organize keypoint descriptors into a tree and search $\log_2 N$.

Q3(b)

(i) - Assume we have n correspondences or correct matches between the 2 views (u, v) and (u', v') .

- Each correspondence must satisfy epipolar constraint

$$(2) \quad \begin{bmatrix} u' & v' & 1 \end{bmatrix} \begin{bmatrix} F \end{bmatrix} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = 0$$

where using $n \geq 8$ we can solve for F and approximate to nearest 3×3 matrix with $\det F = 0$ by SVD.

$$(1) \quad \text{ie. Solve } A \underline{f} \approx 0$$

$$\begin{bmatrix} 9 \times 1 \end{bmatrix}$$

$$(1) - F = K^{-1} T_x R K^{-1}$$

(1) - Solve for $E = K^T F K$ using known K (internal parameters)

$$- E = T_x R.$$

(1) - Solve for T_x and R matrices by SVD of $E = U \Lambda V^T$

Q3(b)(ii) Left projection matrix

$$P_1 = K[I|0]$$

$$\text{or } \begin{bmatrix} su \\ sv \\ s \end{bmatrix} = \begin{bmatrix} P_1 \\ 3 \times 4 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

right projection matrix

(2) $P_2 = K[R|T]$

$$\text{or } \begin{bmatrix} su' \\ sv' \\ s \end{bmatrix} = \begin{bmatrix} P_2 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

Q3(b)(iii) Each image pt gives equation of a ray (2 planes)

Solve by triangulation of rays.

Equivalent to solving for (X, Y, Z) from 4 linear equations by least squares.

(2)

Optimize both P_{ij} and X_k over many views by tracking features and bundle-adjustment.

(2)

Complete reconstruction requires matching more densely. Use feature matcher to guide dense pixel matching in stereo.

Q 4 a) Provide a probabilistic interpretation for the network's output and use this to justify the form of the objective function. [20%]

Answer

Let the probability that a person's age is t^n given the image of their face $Z^{(n)}$ and the network weights V, W be given by a Gaussian $p(t^{(n)}|Z^{(n)}, V, W) = \mathcal{G}(t^{(n)}; x^{(n)}, \sigma^2)$. In this way, the output of the network $x^{(n)}(Z^{(n)}, V, W)$ is the mean of the Gaussian and σ^2 is its variance. The probability of the dataset labels given the parameters is therefore:

$$p(D|V, W) = \prod_{n=1}^N p(t^{(n)}|Z^{(n)}, V, W) = \frac{1}{\mathcal{Z}(\sigma^2)} \exp\left(-\frac{1}{2\sigma^2} \sum_{n=1}^N (t^{(n)} - x^{(n)})^2\right) \quad (1)$$

Placing independent zero-mean Gaussian priors over the weights $V_{i,j}$ and $W_{i,j}$ with variance $1/\alpha$ and $1/\beta$ respectively, yields:

$$p(V, W|\alpha) = \prod_{i,j} p(w_{i,j})p(v_{i,j}) = \prod_{i,j} \frac{1}{\mathcal{Z}(\alpha, \beta)} \exp\left(-\frac{1}{2} (\alpha W_{i,j}^2 + \beta V_{i,j}^2)\right) \quad (2)$$

In this way we can interpret the objective function as relating to the probability of the weight vector given the training data, $p(W, V|D, \alpha) = \frac{1}{\mathcal{Z}(\alpha, \beta)\mathcal{Z}(\sigma^2)} \exp(-G(V, W))$.

b) Describe how to train the network's convolutional weights W using gradient descent. Compute the derivative required to implement gradient descent. Simplify your expression and interpret the terms. [40%]

Answer

The gradient descent algorithm operates as follows:

- i. initialise the weights (e.g. using Gaussian noise with a small variance)
- ii. compute the derivative of the objective function with respect to the convolutional weights $\frac{dG(V,W)}{dW_{i,j}}$
- iii. step down the gradient $W_{i,j} \leftarrow W_{i,j} - \eta \frac{dG(V,W)}{dW_{i,j}}$ (η is a user defined learning rate)
- iv. loop to step (2) until $\Delta G(V, W) < \text{tol}$

To compute the derivative we use backpropagation (aka the chain rule)

$$\frac{d}{dW_{a,b}} G(V, W) = \sum_{n=1}^N \sum_{i,j} \frac{dG(V, W)}{dx^{(n)}} \frac{dx^{(n)}}{dy_{i,j}^{(n)}} \frac{dy_{i,j}^{(n)}}{da_{i,j}^{(n)}} \frac{da_{i,j}^{(n)}}{dW_{a,b}} + \beta W_{a,b} \quad (3)$$

where each of the terms are,

$$\frac{dG(V, W)}{dx^{(n)}} = \frac{1}{\sigma^2} (x^{(n)} - t^{(n)}), \quad \frac{dx^{(n)}}{dy_{i,j}^{(n)}} = V_{i,j} \quad (4)$$

$$\frac{dy_{i,j}^{(n)}}{da_{i,j}^{(n)}} = f'(a_{i,j}^{(n)}) \quad \frac{da_{i,j}^{(n)}}{dW_{a,b}} = Z_{i-a,j-b}^{(n)} \quad (5)$$

14/14

Q4(cont)

Combining the terms together yields

$$\frac{d}{dW_{a,b}} G(V, W) = -\frac{1}{\sigma^2} \sum_{n=1}^N (t^{(n)} - x^{(n)}) \sum_{i,j} V_{i,j} f'(a_{i,j}^{(n)}) Z_{i-a,j-b}^{(n)} + \beta W_{ik}. \quad (6)$$

So, the derivative is simply the sum over all datapoints of the error between the predicted and true ages $(x^{(n)} - t^{(n)})$ multiplied by the sensitivity of the network's output on the convolutional weights plus a linear weight decay term. The sensitivity is a convolution between the product of the output weights and non-linearity derivative $V_{i,j} f'(a_{i,j}^{(n)})$ and the image flipped in the x and y directions $Z_{-i,-j}^{(n)}$.

- c) Describe enhancements to the architecture of the network that might improve its ability to estimate the age of a person from an image of their face. [40%]

Answer

There are lots of possible ways of improving the architecture of the network.

- i. A first enhancement would use **additional sets of convolutional weights**. Currently the method only uses one set and this means that it is only able to extract a single feature (e.g. a specific oriented edge) to perform regression.
- ii. A second enhancement, would use a **pooling/subsampling** stage after the non-linear stage. This would pool over a local neighbourhood and pick e.g. the max or average value. This will introduce shift invariance and reduce the number of parameters that are required in the layers above.
- iii. A third enhancement would be to use a neural network with **many layers** each of which is structured as above. Together these enhancements lead to deep convolutional neural networks.

The answer should describe these enhancements in detail which is bookwork.