4F12 Computer Vision (2017)  $Q|_{(a)}$ (i)-Low-pass filter to remove high - frequency additive ruise which is amplified in gradient operations. large of, small low-par filter cut - of fre - J determine, amount of blurr (2)(ü)  $O \quad g_{\sigma}(x) * g_{\sigma}(y) = g_{\sigma}(x) \quad g_{\sigma}(y)$ = 40-(x,y) 202 where  $G_{\sigma}(x,y) = \frac{1}{2\pi\sigma^2}$ o Instead of N<sup>2</sup> multiplater, (N=2n+1) ve perform 2(N). 4 (iii) q= (x) Sample go (x) until tails are <1000 peak volue 90=1(2)= - C e.2 0.0003 · < 1000 ·. n>3 Sizo & filter 2n+1 = n= 3 6 i-3 i-2 i-- ( 6- 0 i-1 i-2 i-3 X 0.004 0.054 0.242 0.399 0.242 0.054 0.004

\_\_\_()|(ь\_) (i) Create a descriptor from vector of pixels in 16×16 patch, 250 Dector X<sub>NC</sub> 256×1 I(K)S Normalize by subtraction of mean , in and of (normalize long th) S.S.D =  $\sum_{i,j} \left( I_{i}(i,j) - I_{i}(i,j) \right) \longrightarrow$ Compute X for comparison. X NC 1 . X NC 2 (lii.) SIFIdescriptor - sample 16x 16 pixel S(x, y, o; ) at dominat orientation - compute VS d-pixels (gradient magnitude + on entation) weight gradient by gamerian weighting ( = 8 pixels at scale of lox 16) - produce lexter cells (16) - produce HOG at 45° bins (8 dir") adding grad may to bins concatorate to 128 Prector nonralize length - renave outliers by truncations and value >0.2 to 0.2. jiij N.C - very simple to compute, 2500, poor of scale + unortation incorrect SIFT - Invariant to scale | orientation + lighting - robust to small geometric/perspective distortion. Poor at boundaries

(a) (i)\_ pin-hole camera planar perupective anto plane No mon-line or distortion. (es no radid les distortion) Xi Xu = = <u>X1</u> = X3 X1 homosarow representation (ü) <u>S</u> represents non-lineor inverse-depth scale <u>S</u> is homogeous scaling (denominator) (geonety) ( alsobro) (iv) Let X4 -> O for ptr. at 00 internation and contraction of matrix, 3 dof) 3×4 (k)R = (iv) (3 dd). Ophica (ontro) 3×1 3×3 orhound fku O L. fky Vo scale pixel rive , whore kn Kv are laternal posseters Uo Vo J-principal point f - food lenot

(3(a) (i) (1) Rotation about optical axis with no translation = K R Kt = ( 3x.3 (B) Viewing a-plane eg Z= O u V  $r_{11}$   $r_{12}$   $t_1$  X where  $R_1 =$ Fil-Fiz-Fiz = K (8) K Where R2= = K 121 121 t2 [1] 22 (23 ٧I ົາ  $= K \left[ \frac{r_{21}}{r_{22}} + \frac{r_{22}}{r_{2}} + \frac{r_{22}}{r_{2}} + \frac{r_{11}}{r_{12}} + \frac{r_{12}}{r_{11}} + \frac{r_{12}}{r_{12}} + \frac{r_{12}}{r_{11}} + \frac{r_{12}}{r_{12}} +$ \* •u

3(b) 8 d d f. rotation, (Idef) scale, s (1 do) tronslation (2 d d) (u.,vo) shear ( 2 ddf; axist mag VP2 - hunizon (2 dof) anning <u>(ä)</u> Under weak peopechis (+1, +1, +1, +2, +2, +1, 0 0  $\begin{bmatrix} t_{11} & t_{12} & t_{13} \\ t_{21} & t_{22} & t_{23} \\ 0 & 0 & t_{33} \end{bmatrix}$ 2 becomes on Affine tran and no non-linear perspecture effects (ie no faming // the stay //

() ()3(6) Rigid-budy motion - epipolar geometry XI = RX+I where  $T_x = \begin{pmatrix} 0 - t_z & t_y \\ t_z & 0 - t_x \\ -t_x & t_y \end{pmatrix}$ E = Ix R  $F = K^{-1} T_{\times} R K^{-1}$ (a) Find n 28 correspondences and volve Af = 0 9×1 FirdF Project F to F with rank 2 is. det F=0 Solve for E = KTFK using known K Now E= TAR Decompose into anti-symphic and orthonord matrix by SVDqE E=UNV where N= 10  $T_{x} = U \begin{bmatrix} 0 & 0 \\ -1 & 0 \\ 0 & 0 \end{bmatrix} U^{T}$  $R = U \begin{bmatrix} 0 + 0 \\ 1 & 0 \end{bmatrix} V^{T}$ ()4 solution Resolve ambis with by check depthr are positive.

Q4

 a) Describe the components that make up a standard convolutional neural network that is used for image classification. Briefly explain the rationale for the form of each of the components. [40%]

## Answer

Bookwork that should include the following:

- Convolutional filtering, motivated by the translation invariance of images and the need to reduce the number of parameters to be learned
- (Pointwise) Non-linearities, motivated by the fact that non-linear computation needs to be made, and certain point-wise non-linearities followed by linear weighting are universal approximators.
- Pooling / sub-sampling, motivated by the fact that high-level features are coarser and that subsampling again reduces the number of parameters. It also helps to build in translation invariance.
- b) Compute the derivative required to implement gradient descent learning of the network's convolutional weights W. Simplify your expression and interpret the terms. [40%]

#### Answer

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

$$\frac{\mathrm{d}}{\mathrm{d}W_{a,b}}G(V,W) = \sum_{n=1}^{N} \sum_{i,j} \frac{\mathrm{d}G(V,W)}{\mathrm{d}x^{(n)}} \frac{\mathrm{d}x^{(n)}}{\mathrm{d}y^{(n)}_{i,j}} \frac{\mathrm{d}y^{(n)}_{i,j}}{\mathrm{d}a^{(n)}_{i,j}} \frac{\mathrm{d}a^{(n)}_{i,j}}{\mathrm{d}W_{a,b}} + \beta W_{a,b}.$$
 (1)

where each of the terms are,

$$\frac{\mathrm{d}G(V,W)}{\mathrm{d}x^{(n)}} = \frac{x^{(n)} - \mathbf{t}^{(n)}}{x^{(n)}(1 - x^{(n)})}, \quad \frac{\mathrm{d}x^{(n)}}{\mathrm{d}y^{(n)}_{i,j}} = V_{i,j}x^{(n)}(1 - x^{(n)}) \tag{2}$$

$$\frac{\mathrm{d}y_{i,j}^{(n)}}{\mathrm{d}a_{i,j}^{(n)}} = f'(a_{i,j}^{(n)}) \quad \frac{\mathrm{d}a_{i,j}^{(n)}}{\mathrm{d}W_{a,b}} = Z_{i-a,j-b}^{(n)} \tag{3}$$

Combining the terms together yields

$$\frac{\mathrm{d}}{\mathrm{d}W_{a,b}}G(V,W) = -\sum_{n=1}^{N} \left( \mathbf{t}^{(n)} - x^{(n)} \right) \sum_{i,j} V_{i,j} f'(a_{i,j}^{(n)}) Z_{i-a,j-b}^{(n)} + \alpha W_{ik}.$$
 (4)

So, the derivative is simply the sum over all datapoints of the error between the predicted and true labels  $(x^{(n)} - t^{(n)})$  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) The company would like to extend the network to estimate the average population density of the area in each image. Describe how to extend the architecture of the network to perform this additional task. Explain your design.

#### Answer

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

The first step is to **add a second output to the network**  $p^{(n)}$  that models the average population density  $\rho$ . It makes sense to use the same features as are used for classification to do this as you would expect population density to be higher in urban areas. These common features can be passed through a different set of output weights  $U_{i,i}$  in order to form the scalar population density estimate output of the network,

$$p^{(n)} = \sum_{i,j} U_{i,j} y_{i,j}^{(n)}.$$

The network's regression weights U can then be trained by adding a regression term to the objective function and optimising all parameters together,

$$H(U,W) = \frac{1}{2\sigma^2} \sum_{n=1}^{N} \left(\rho^{(n)} - p^{(n)}\right)^2 + \frac{\alpha}{2} \sum_{i,j} U_{i,j}^2.$$

Strictly the question does not ask for the form of the new objective, but it interacts with the design decisions.

There are also more general enhancements that could be mentioned, but which are not a substitute for the above.

- i. One enhancement would be to 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 classification / regression.
- ii. A second enhancement, would use a **pooling/subsampling stage** after the nonlinear 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.

#### 4F12 Comments:

# Q1

Well answered by most candidates. Only part to cause problems being (b)i- describing a descriptor for normalised intensities for use in cross-correlation.

## Q2

Well answered by most candidates. Only difficulty was in (b)ii in deriving the equations for the perspective projection of a line with the given projection matrix and then computing the vanishing point. Many struggled with a simple derivation of the equation of the horizon.

# Q3

A well-answered question. Marks were lost in not expressing the matrices in terms of the camera motion parameters.

#### Q4

An unpopular question because this was newer material but those that did attempt it tended to provide good answers. When describing the convolutional neural network in part (a) a common error was to describe just a single convolution / non-linear / subsampling component without mentioning how these are used to build a network. Many candidates failed to identify the last part of the question as a regression problem. The few candidates that attempted made good progress.