

# CS 280 Computer Vision Midterm Exam

Thursday, April 7<sup>th</sup>, 2016

Trevor Darrell, Alexei A. Efros  
Richard Zhang, Deepak Pathak, Jeff Donahue

<b>First</b>	
<b>Last</b>	
<b>SID</b>	

**Name of neighbor to your:**

<b>Left</b>	
<b>Right</b>	

**For staff use only:**

Q1. True/False	___/14
Q2. Short answers	___/15
Q3. CNN parameters	___/18
Q4. Binocular stereopsis	___/10
Q5. Comparing heights	___/12
Q6. Image processing	___/12
Q7. Pyramids	___/5
<b>Total</b>	<b>___/86</b>

# Q1 [14 pts] True/False

For each of the following, indicate whether the statement is TRUE or FALSE.  
Each True/False question is worth 1 point.

---

T F Parallel lines in an image will remain parallel after applying a homography

T F Perspective projection from a 3D point in the world to a 2D point in the image is linear in Cartesian coordinates

T F Vanishing points of all lines on a plane lie on a line.

T F In a Camera Obscura, the smaller the diameter of the pinhole, the sharper the image becomes.

T F Rotation, Affine and Projective transforms in Euclidean space can be expressed simply as a linear map in homogeneous coordinates.

---

T F Filters approximating the oriented cells in V1 discovered by Hubel and Wiesel emerge naturally if you learn a sparse code on a dataset of natural images

T F Filters approximating the oriented cells in V1 discovered by Hubel and Wiesel emerge naturally in the first layer of a CNN trained for ImageNet classification

T F Given only a brightness histogram [256 bins] of an image, one can estimate an orientation histogram [6 orientation bins] for that image.

---

T F A CNN implements a linear function

T F Optimizing parameters for a CNN is a non-convex problem

T F During a feed-forward pass, *any* convolutional layer can be implemented as a fully connected layer

T F During a feed-forward pass, *any* fully connected layer can be implemented as a convolutional layer

T F It is possible to model varying size sequences with varying spatial size inputs (i.e. each element of the sequence is of different size) using a recurrent LSTM unit.

---

T F RANSAC is used to get rid of outliers in situations when there are not enough point correspondences for least-squares estimation to work well.

## Q2 [15 pts] Short Answer

(a) [2pts] Express the point where 2-D parallel lines intersect in homogeneous coordinates (assume any notation).

$(x, y, 0)$

(b) [2pts] How many degrees of freedom are there in perspective projection transform between two images?

There are 8 degrees of freedom in a perspective projection transform between two images. Given point  $p = (x, y, 1)$  and  $p' = (x', y', 1)$  in a second image, we would like to solve for  $P' = H * p$ , where  $p'$  is  $P'$  normalized over the last coordinate. Matrix  $H$  has 9 values, but one must be fixed to prevent arbitrary scaling.

(c) [2pts]  $E \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = M \begin{pmatrix} X \\ Y \\ Z \\ 1 \end{pmatrix}$  What is matrix  $M$  for orthographic projection?

$$M = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

**(d) [2pts]** Derive the rank of the Essential Matrix  $E$ .

$E = [t]R$  --- Here  $t$  is of rank 2 (because the direction along  $t$  forms a rank 1 null space of  $t$ ; and other two orthogonal directions constitute rank 2) and  $R$  is full rank. Thus  $E$  is rank 2.

**(e) [2pts]** A projective transformation between two images  $X$  and  $X'$  can be modeled by transforming each point  $p$  in one image to a point  $p'$  in a new image using  $p' = Hp$ , where  $H$  is the homography. Briefly explain, why we can't do the same thing with the Fundamental Matrix  $F$ , to help transform each point  $p$  in the Left image to a point  $p'$  in the Right image:  $p' = Fp$ ?

There is an ambiguity in this case, since a point in one image maps to the *epipolar line* in the other image.

**(g) [3pts]** What is the typical loss function used to train Siamese Networks? Define each term in the expression and the importance of each component.

$$L(x_1, x_2, y) = y * \max(\|f(x_1) - f(x_2)\| - m_1, 0) + (1-y) * \max(m_2 - \|f(x_1) - f(x_2)\|, 0)$$

**y** is whether both examples have the same label

**f** is the network

**m1** is the margin if they have the same label

**m2** is the margin if they have different labels

if they have the same label, we want them to be closer together

if they have different labels, we want them to be further apart

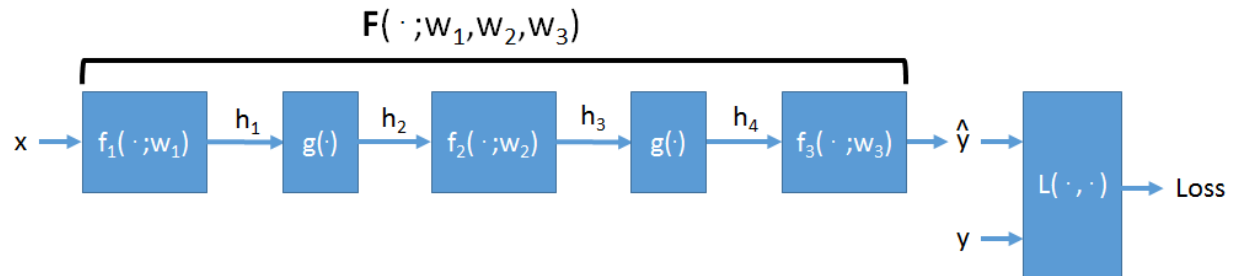
**(h) [1 pt]** I would like to take a derivative of discrete signal  $x[t]$ . What convolutional kernel can I use?

**[1 -1] or [1 0 -1] or flipped**

**(i) [1 pt]** I would like to take the second derivative of discrete signal  $x[t]$ . What convolutional kernel can I use?

**[1 -2 1] or whatever is above convolved with itself**

### Q3 [18 pts] CNN Computations



Given input  $x$  and ground truth  $y$  we have a network  $F$  with parameters  $w_1, w_2, w_3$ , drawn above. During training, we penalize the network with loss function  $L(F(x), y)$ .

In terms of individual layer Jacobians (e.g.,  $\frac{\partial h_3}{\partial h_2}, \frac{\partial h_2}{\partial h_1}, \frac{\partial h_3}{\partial w_2}$ ), write formulas for the following:

(a) [3pts] Write  $\frac{\partial L}{\partial x}$

$$dL/dx = dL/dy\_hat * dy\_hat/dh4 * dh4/dh3 * dh3/dh2 * dh2/dh1 * dh1/dx$$

(b) [3pts] Write  $\frac{\partial L}{\partial w_2}$

$$dL/dw2 = dL/dy\_hat * dy\_hat/dh4 * dh4/dh3 * dh3/dw2$$

---

We have an input image of spatial resolution 224x224 and channel dimension of 3. The first layer of the CNN has 64 convolutional filters of spatial size 3x3. The filters are evaluated with stride of 2. The input image is zero-padded with one (1) extra pixel around all borders (resulting in a padded image of size 226x226). (There are no groups or dilations).

**(c) [3pts]** What are the output dimensions for a single image after applying the first convolutional layer?

**112x112x64**

**Spatial size is  $(X - X_{\text{kern}} + 2 \cdot \text{PAD} + 1) / S = (224 - 3 + 2 + 1) / 2$**

**(d) [3pts]** How many parameters are there in the convolutional layer, including a bias term applied to the convolution output? You may leave your answer as an algebraic expression.

**$X_{\text{kern}} \cdot Y_{\text{kern}} \cdot \text{Chnout} \cdot \text{Chnin} + \text{Chnout} = 3 \cdot 3 \cdot 64 \cdot 3 + 64$**

---

We have an input layer of spatial resolution 32x32 and channel dimension of 3. The first layer of our network is fully connected with output size 1024. The second layer of our network is a ReLU nonlinearity. The third layer of our network is a fully connected layer with output size 1024.

**(e) [3 pts]** How many parameters are there in the first fully connected layer, including a bias term applied to the fully connected output? You may leave your answer as an algebraic expression.

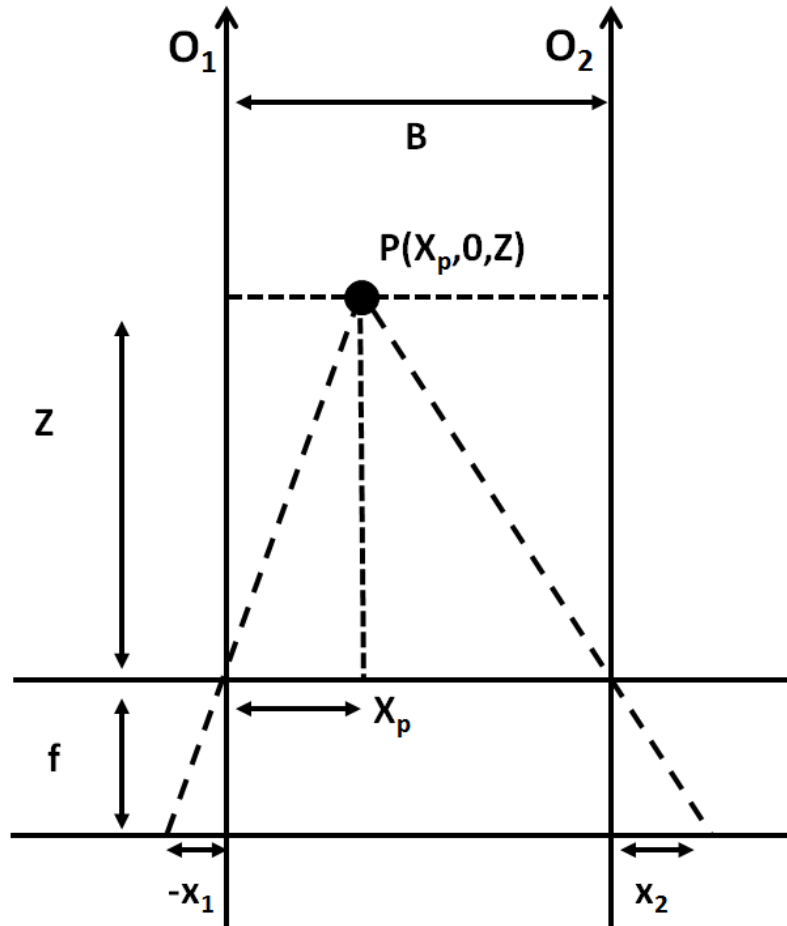
**$X_{\text{kern}} \cdot Y_{\text{kern}} \cdot \text{Chnout} \cdot \text{Chnin} + \text{Chnout} = 32 \cdot 32 \cdot 1024 \cdot 3 + 1024$**

**(f) [3 pts]** How many parameters are there in the last fully connected layer, including a bias term applied to the fully connected output? You may leave your answer as an algebraic expression.

**$X_{\text{kern}} \cdot Y_{\text{kern}} \cdot \text{Chnout} \cdot \text{Chnin} + \text{Chnout} = 1 \cdot 1 \cdot 1024 \cdot 1024 + 1024$**

## Q4 [10 pts] Binocular Stereopsis

We have two optical axes  $O_1$ ,  $O_2$  offset by baseline distance  $B$ . Point  $P(X_p, 0, Z)$  projects onto points  $(x_1, 0)$  and  $(x_2, 0)$  on the image planes.



(a) [5pts] Write values  $x_1, x_2$  in terms of the given constants.

$$x_1 = -f * X_p / Z$$

$$x_2 = f * (B - X_p) / Z$$



**(b) [5pts]** What is the estimated depth  $Z$ , given focal length  $f$ , baseline  $B$ , and difference in project points  $\Delta x = x_2 - x_1$  ?

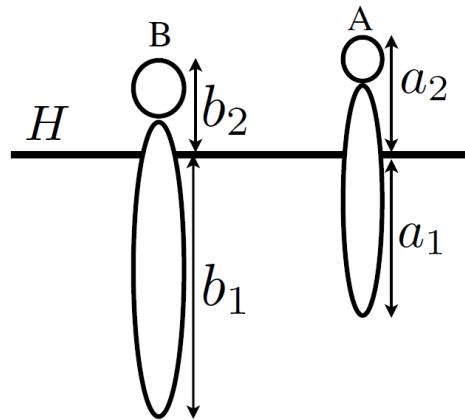
$$\text{delta\_x} = x_2 - x_1$$

$$= f * (B - X_p) / Z - f * X_p / Z$$

$$= f * B / \text{delta\_x}$$

$$Z = f * B / \text{delta\_x}$$

## Q5 [12 pts] Comparing Heights



Consider the image of person A and person B standing on the ground plane, as taken by a perspective camera of focal length  $f$ .  $H$  is the horizon line ( $y = 0$ ).  $a_1$  is the distance between the A's toes and the horizon, and  $a_2$  is the distance between A's head and the horizon in units of pixels. Similarly for person B. Distances  $a_1, a_2, b_1, b_2$ , are expressed in the camera plane. Suppose A's height is  $h$  feet.

(a) [3 pts] From the picture, who is taller? Briefly explain why.

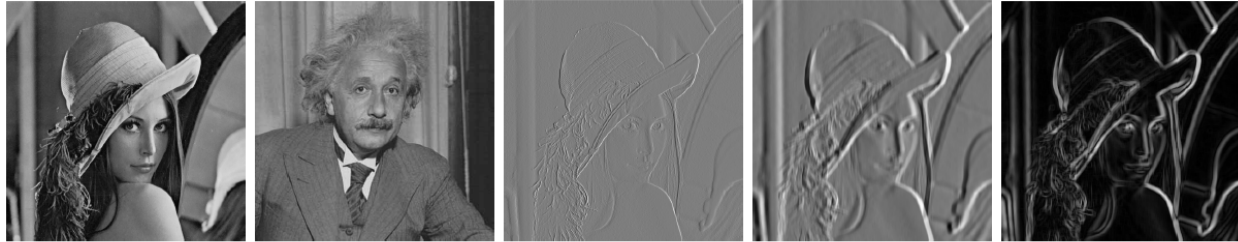
**Person A is taller. Person A is standing behind person B, since you can draw a line from person B's feet to person A's feet and extend to the horizon. If they were the same height, then their feet and their heads would lie on parallel planes. In that case, you could draw a line from person B's head to person A's head to the same point on the horizon. However, if you draw a line from person B's head to person A's head, it actually goes up and doesn't even meet the horizon. Thus, person A is taller.**

(b) [9 pts] Give expressions for the following in terms of  $f, a_1, a_2, b_1, b_2, h$ . Show your work briefly.

How many feet above the ground is the camera?	$h_c = h * a_1 / (a_1 + a_2)$ <p>The horizon line is straight ahead of the camera. The height of the camera is equal to the percentage of person A which is above the horizon times his height.</p>
How tall is person B (in feet)?	$h_b = h_c * (b_1 + b_2) / b_1$ $= h * a_1 / (a_1 + a_2) * (b_1 + b_2) / b_1$ <p>We know portion <math>b_1</math> corresponding to camera height <math>h_c</math>. We then add on the portion of Person B's height which is above the horizon line.</p>
Distance (along the z-axis) from the camera to person B (in feet)?	$b_1 + b_2 = f * h_b / Z$ $Z = f / (b_1 + b_2) * h_b$ $= f * h / b_1 * a_1 / (a_1 + a_2)$ <p>Height of person B projects onto <math>b_1 + b_2</math> in the image plane. Use projection equation and then solve for <math>Z</math>.</p>

## Q6 [12 pts] Image Processing

Consider the input images A and B. Assume that  $g$  is a 2-D Gaussian filter,  $d_1 = \begin{pmatrix} -1 & 1 \end{pmatrix}$  and  $d_2 = \begin{pmatrix} -1 & -1 \end{pmatrix}^T$ .



(a) Image A

(b) Image B

(c) Image C

(d) Image D

(e) Image E

For each of the following desired outputs, write out the image processing operations that you would do with the images A and B using the filters  $g$ ,  $d_1$  and  $d_2$  to get the desired output. The first has been done for you as an example. You can ignore edge artifacts.

Smooth image A	$A * g$
Remove high frequencies from image A	$A * g$
Remove low frequencies from image A	$A + \text{lambda} * (A - A * g)$
Hybrid image which looks like A from far off and looks like B from close by	$A * g + B - B * g$
Generate image C from image A	$A * d_1$
Generate image D from image A	$A * g * d_1$
Generate image E from image A	$\text{Sqrt}((A * g * d_1)^2 + (A * g * d_2)^2)$

## Q7 [5 pts] Pyramids

We are given an image  $I$ . Let  $G_1 = I$ ,  $G_2$ , and  $G_3$  be the resulting images from a 3-levels Gaussian Pyramid constructed from  $I$ . Equivalently, let  $L_1$ ,  $L_2$ ,  $L_3$  be the resulting images from a 3-level Laplacian Pyramid constructed from  $I$ .

Assume that the original image contained all spatial frequencies. Label the resulting images as containing H (high), M (medium), and/or L (low) frequencies. You can use none/one/many labels per image.

$G_1$ :   H  M  L  

$G_2$ :   M L  

$G_3$ :   L  

$L_1$ :   H  

$L_2$ :   M  

$L_3$ :   L