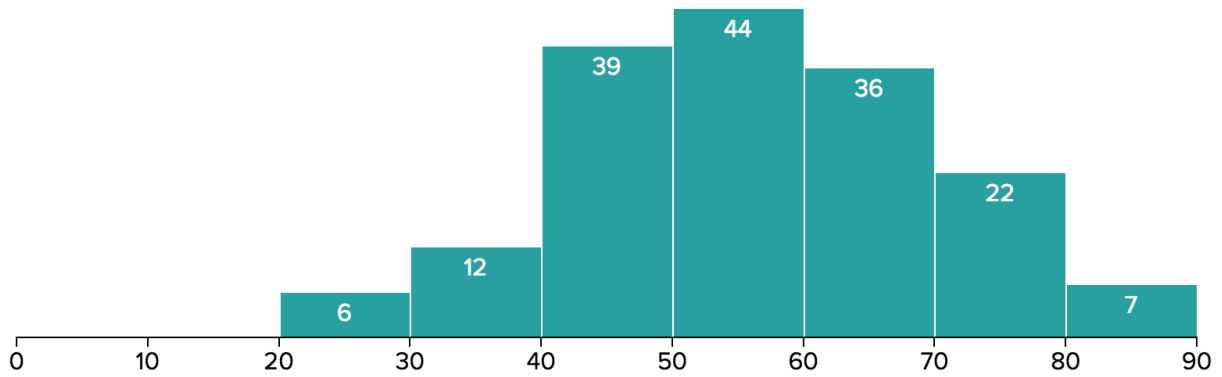


Midterm statistics (graded out of 85, with 5 bonus points):



**Overall statistics (across whole class; 165 exams)**

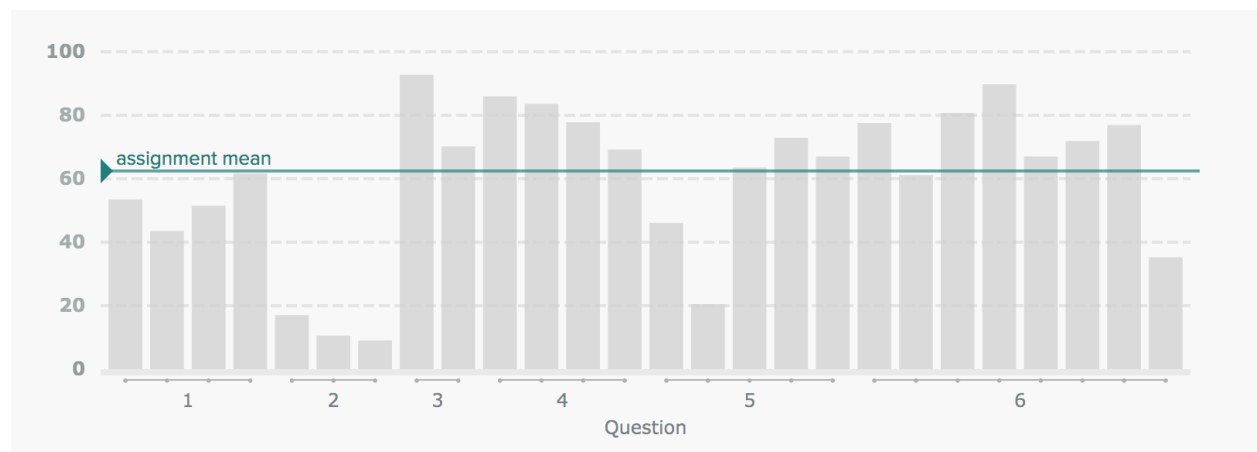
High score: 89 (104.7%)  
Mean score: 56.25 (66.2%)  
Median score: 56 (65.9%)  
Standard deviation: 13.45 (15.8%)

**Graduate student statistics (149 exams)**

High score: 89  
Mean score: 56.1  
Median score: 56  
Standard deviation: 13.4

**Undergraduate student statistics (16 exams)**

High score: 82  
Mean score: 57.7  
Median score: 57.5  
Standard deviation: 14.1



**A comment on the exam:** This was our first time doing an exam for this material. As such, there was the potential that we wouldn't get the difficulty right. This exam, based on the statistics, was on the harder side.

Thus, we're going to make a grading scheme *that can't hurt you* and does not penalize students who did well on the midterm. In this scheme, we will calculate grades based on the 40% HW, 40% MT, and 20% project (really 19% project, 1% participation) scale. We will **set the mean of this calculation** to achieve a B+ in the class, and grade based on a scale around this number.

After this, we will calculate everyone's final grades based off of the 40/40/20 scale, as well as a 40/30/30 scale (i.e., if the midterm had 30% of the weight, and the project had 29% of the weight). In this calculation, the statistics of the mean may be adjusted to make it meaningful. Using the set scale from the 40/40/20 calculation, we will assign your grade based on the higher value of your final grade from the 40/40/20 scale vs the 40/30/30 scale.

If you did excellently on the midterm, it will still count 40%. If you didn't do so well, you can earn back points on the project. Again, if you did do well on the exam, this grading scheme does not hurt you; it merely helps your peers.

**Logistics moving forward:** We will allow the submission of regrade requests up until this Friday. If you have questions about the exam, you should talk to the person who graded that exam question. In particular, we all grade in different ways, and e.g., I can't comment on how Tianwei or Cheng graded. The key is that we all did our best to grade consistently, and in assessing regrades, it is based on *consistency*. If I graded all students in a certain way, even if you don't think the point assignments were judicious, what matters is that all students were graded in this way. Thus, regrade requests typically fall into the category of the rubric being used incorrectly to grade your work.

- If you have inquiries about questions #1 and #2, refer to Tianwei.
- If you have inquiries about questions #3 and #4, refer to Cheng.
- If you have inquiries about questions #5 and #6, refer to Prof. Kao.

**Undergraduate scales:** As the undergraduates scored slightly higher than the graduate students, this shouldn't hurt you. If indeed in assigning final grades, the undergraduates perform better than graduate students, we will grade all students together.

**Concerns about your grade:** If you have concerns about your grade, please don't hesitate to talk directly to Prof. Kao.

**ECE 239AS, Winter 2018**  
Department of Electrical Engineering  
University of California, Los Angeles

**Midterm**  
Prof. J.C. Kao  
TAs T. Xing and C. Zheng

UCLA True Bruin academic integrity principles apply.

Open: book, notes, handout, computer.

Closed: internet, except for downloading / viewing from ECE239AS website / CCLE.

6:00 - 7:50pm.

Wednesday, 21 Feb 2017.

State your assumptions and reasoning.

No credit without reasoning.

Show all work on these pages.

Name: \_\_\_\_\_

Signature: \_\_\_\_\_

ID#: \_\_\_\_\_

Problem 1 \_\_\_\_\_ / 18

Problem 2 \_\_\_\_\_ / 5 (BONUS)

Problem 3 \_\_\_\_\_ / 15

Problem 4 \_\_\_\_\_ / 17

Problem 5 \_\_\_\_\_ / 15

Problem 6 \_\_\_\_\_ / 20

Total \_\_\_\_\_ / 85 + 5 BONUS

1. (X points) **Short answer on machine learning basics.**

- (a) (X points) **SVM basics.** In writing the SVM cost function, the hinge loss has a margin of 1. Notice that we set the margin always to 1, as opposed to having it be a hyperparameter  $\Delta$ . Why can we do this? Justify your answer (*at most three sentences*).
- (b) (X points) **Comparing a kNN and linear classifier.** Compare bias and variance of a linear classifier to kNN. Keep your answer to *at most six sentences*.
- (c) (X points) **Bias and variance.** Let  $X$  be a random variable representing the current time (not distinguishing between a.m and p.m). Suppose  $X$  is sampled from a uniform distribution  $X \sim U[0, 12]$ . Let  $Z$  be a random variable representing the difference between the estimated time and the current time, wrapped to  $+/- 6$  hours, so that  $Z \sim U[-6, 6]$ . If  $z = 0$ , that means the estimated time is the current time. Consider two estimators of the time differences.

- i.  $\hat{z}_1$  is the estimator for a stopped clock (which always shows the same time). Note that it gives the correct estimate twice a day.
- ii.  $\hat{z}_2$  is the estimator for a clock which works with perfect precision but is in the wrong timezone, so it is always one hour late and never gives the correct estimate.

Calculate the bias and the variance of both estimators  $\hat{z}_1$  and  $\hat{z}_2$ .

- (d) (X points) **Cheating a competition.** You are running a machine learning competition where competitors submit a classifier that makes a prediction. You have one private test dataset that you use to evaluate all submitted classifiers. You found out that one of your colleagues allowed competitors to evaluate their classifier on this test dataset an unlimited amount of times (i.e., whenever they wanted to check). How could a competitor take advantage of this to do better in the competition? Keep your answer to *at most three sentences*.

**Solution:**

- (a) Changing the margin parameter does not change the margin or decision boundary. If we change it to 2, the change could be canceled out by increasing  $W$  and  $b$  to twice of their original values.
- (b) kNN can fit a dataset to arbitrarily high accuracy given enough points, so its bias can be arbitrarily small. A linear classifier can never do better than fit a single hyperplane so its bias will be larger. The predictions from kNN depend only on a small number of data points ( $k$ ) and so are subject to high variance they will change with changes in any of these points. The predictions from a linear classifier depend on all the data points. Its variance is typically much smaller for that reason.
- (c) Let  $Z \sim U[-6, 6]$ . Bias:

$$b_z(\hat{z}_1) = E_z(\hat{z}_1) - Z = 0 - 0 = 0$$

$$b_z(\hat{z}_2) = E_z(\hat{z}_2) - Z = -1 - 0 = -1$$

Variance:

$$\text{var}_z(\hat{z}_1) = E_z((\hat{z}_1 - E_z(\hat{z}_1))^2) = E_z((\hat{z}_1 - 0)^2) = 12$$

$$\text{var}_z(\hat{z}_2) = E_z((\hat{z}_2 - E_z(\hat{z}_2))^2) = E_z((\hat{z}_2 + 1)^2) = 0$$

- (d) A competitor could optimize hyperparameters on the testing set.
2. (X points) **Maximum likelihood estimation.** We observe  $n$  i.i.d random samples  $x_1, x_2, \dots, x_n$  that we model as arising from a Poisson distribution. We want to use maximum likelihood estimation method to estimate the parameter,  $\lambda$ , of the Poisson distribution. The probability mass function of a Poisson distribution is given by:

$$\Pr(x = k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

- (a) (X points) Write the log-likelihood of observing the data under the model assumption.
- (b) (X points) Calculate the derivative of the log-likelihood w.r.t. the parameter,  $\lambda$ .
- (c) (X points) Calculate the maximum likelihood estimator of  $\lambda$ .

**Solution:**

- (a) Likelihood function:

$$L(\lambda; x_1, x_2, \dots, x_n) = \prod_{j=1}^n \exp(-\lambda) \frac{1}{x_j!} \lambda^{x_j}$$

Log-likelihood function:

$$l(\lambda; x_1, x_2, \dots, x_n) = -n\lambda - \sum_{j=1}^n \ln(x_j!) + \ln(\lambda) \sum_{j=1}^n x_j$$

- (b) Calculate derivate:

$$\begin{aligned} & \frac{d}{d\lambda} l(\lambda; x_1, x_2, \dots, x_n) \\ &= \frac{d}{d\lambda} (-n\lambda - \sum_{j=1}^n \ln(x_j!) + \ln(\lambda) \sum_{j=1}^n x_j) \\ &= -n + \frac{1}{\lambda} \sum_{j=1}^n x_j \end{aligned}$$

- (c) Set the derivate equal to zero:

$$\lambda = \frac{1}{n} \sum_{j=1}^n x_j$$

The estimator is the sample mean of  $n$  observation samples.

3. (X points). **Backpropagation**

Let  $\mathbf{h} = \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b})$  be a typical layer of a FC network, with  $\mathbf{x} \in \mathbb{R}^n$  as input and  $\mathbf{h} \in \mathbb{R}^m$  as output. The mean square logarithmic error (MSLE) loss function is a variant of mean square error. The MSLE, defined for *one example (the  $i$ th example)* and denoted by  $\mathcal{L}_i$  is given by:

$$\mathcal{L}_i = (\log |\mathbf{h}\mathbf{h}^T + \mathbf{I}| - \log |\mathbf{y}\mathbf{y}^T + \mathbf{I}|)^2$$

where  $\mathbf{y}$  is the target value. Note that  $|\cdot|$  here denotes the determinant of a matrix.

- (a) (5 points) Write the dimension of following variables  $\mathcal{L}_i, \mathbf{W}, \mathbf{b}, \mathbf{I}, \mathbf{y}$ . For example,  $\mathbf{x} \in \mathbb{R}^n$ .

(b) (10 points) Calculate  $\frac{\partial \mathcal{L}_i}{\partial \mathbf{W}}$  and  $\frac{\partial \mathcal{L}_i}{\partial \mathbf{x}}$ . You can use the following formula without proof.

$$\begin{aligned}\frac{\partial \log |\mathbf{X}|}{\partial \mathbf{X}} &= (\mathbf{X}^T)^{-1} \\ \frac{\partial \mathcal{L}_i}{\partial \mathbf{X}} &= \mathbf{W}^T \frac{\partial \mathcal{L}_i}{\partial (\mathbf{W}\mathbf{X})} \\ \frac{\partial \mathcal{L}_i}{\partial \mathbf{W}} &= \frac{\partial \mathcal{L}_i}{\partial (\mathbf{W}\mathbf{x})} \mathbf{x}^T \\ \frac{\partial \mathcal{L}_i}{\partial \mathbf{X}} &= 2 \frac{\partial \mathcal{L}_i}{\partial (\mathbf{X}\mathbf{X}^T)} \mathbf{X}\end{aligned}$$

**Solution:**

(a)  $\mathcal{L}_i \in \mathbb{R}$ ,  $\mathbf{W} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,  $\mathbf{I} \in \mathbb{R}^{m \times m}$ ,  $\mathbf{y} \in \mathbb{R}^m$ .

(b) Let  $a = \log |\mathbf{h}\mathbf{h}^T + \mathbf{I}| - \log |\mathbf{y}\mathbf{y}^T + \mathbf{I}|$  and  $\mathbf{m} = \mathbf{h}\mathbf{h}^T + \mathbf{I}$ . So

$$\frac{\partial \mathcal{L}_i}{\partial a} = 2a, \quad \frac{\partial \mathcal{L}_i}{\partial \mathbf{m}} = (\mathbf{m}^T)^{-1} (2a) \quad (1)$$

$$\frac{\partial \mathcal{L}_i}{\partial \mathbf{h}} = 4a (\mathbf{m}^T)^{-1} \mathbf{h} \quad (2)$$

$$\frac{\partial \mathcal{L}_i}{\partial (\mathbf{W}\mathbf{x} + \mathbf{b})} = \mathbb{I}(\mathbf{W}\mathbf{x} + \mathbf{b} > 0) \odot [4a (\mathbf{m}^T)^{-1} \mathbf{h}] \quad (3)$$

$$\frac{\partial \mathcal{L}_i}{\partial \mathbf{x}} = \mathbf{W}^T \{ \mathbb{I}(\mathbf{W}\mathbf{x} + \mathbf{b} > 0) \odot [4a (\mathbf{m}^T)^{-1} \mathbf{h}] \} \quad (4)$$

$$\frac{\partial \mathcal{L}_i}{\partial \mathbf{W}} = \{ \mathbb{I}(\mathbf{W}\mathbf{x} + \mathbf{b} > 0) \odot [4a (\mathbf{m}^T)^{-1} \mathbf{h}] \} \mathbf{x}^T \quad (5)$$

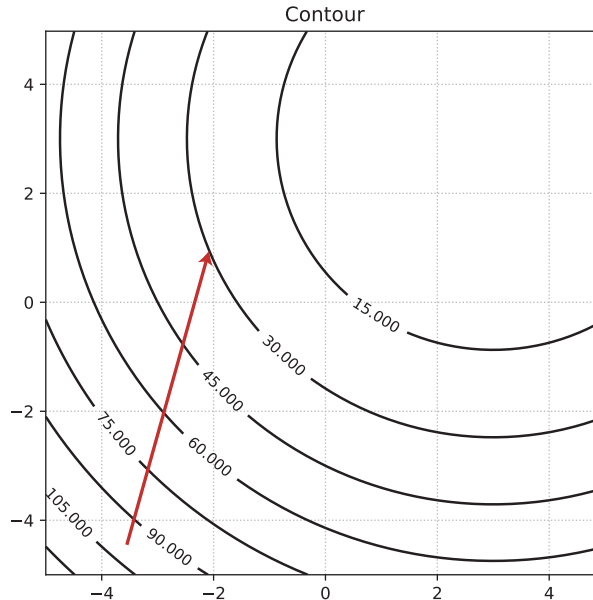
Note that in our solution, we combined the paths of  $\mathbf{h}$  in the 3rd line. You could have backpropagated through  $\mathbf{h}$  and  $\mathbf{h}^T$  separately and then added the contributions of their total derivatives to  $\mathbf{W}$  and  $\mathbf{x}$  to arrive at the same answer.

4. (X points) Optimization (momentum, RMSprop, adam).

(a) (X points) **Momentum**. What parameters does optimizer SGD+momentum maintain that SGD does not? Answer in one or two sentences.

(b) (X points) **Momentum 2**. Imagine separate optimizations, one with SGD and the other SGD with momentum. In both optimizations, the parameters are approaching the same local optima with the same learning rate. Which optimizer is more likely to stay at this local optima, and why? Justify your answer (*at most three sentences*).

(c) (X points) **Adagrad and RMSProp** *Adagrad* and *RMSprop* modify the learning rate in particular directions based off of historical gradients. What weakness of *Adagrad* does *RMSprop* address? Justify your answer (*at most three sentences*).



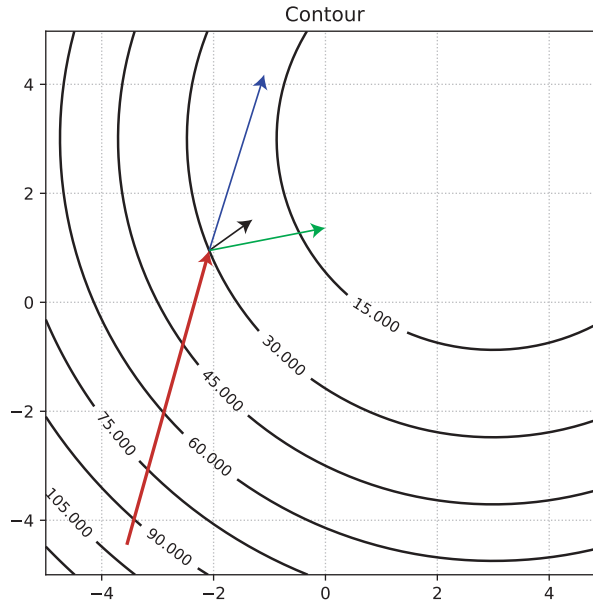
(d) (X points) **Gradient steps** The following figure is the contour plot where the contour lines denote the value of the loss as a function of two weight variables (corresponding to the  $x$  and  $y$ -axis). Imagine we have taken one gradient step in the direction of the red line. Sketch on the same plot the steps taken by SGD, SGD+momentum, and Adagrad. Do not perform any calculations; the sketch should be arrived at through intuition. Give at most a two sentence explanation for each arrow you've drawn.

**Solution:**

- (a) *Momentum* maintain the running mean of the gradients.
- (b) SGD+momentum is less likely to stay in the local optima because the momentum terms may cause the optimization to leave the local optima.
- (c) *Adagrad* accumulates the squared gradients in the denominator so that learning rates shrink and eventually become infinitesimally small. *RMSProp* addresses this by using a running average that exponentially decays historical gradients.
- (d) SGD is black, SGD+momentum is blue, and Adagrad is green. SGD is the gradient according to the contours. SGD+momentum takes into account our historical step going along the red line. Adagrad will downweight the gradient step going vertically and have a larger gradient step in the horizontal direction.

5. (X points) **Short answer on training neural networks.**

- (a) (X points) **Alien brain.** Imagine you are designing a neural network inspired from an alien brain. In this alien brain, experiments show that their neurons are organized like a feedforward network.
  - i. (X points) Curiously, after several experiments, the researchers find that neurons in the first and last layers capture similar features, indicating that they have similar weights. They want you to build a neural network based off of the principle that the



first and last layers have similar weights. Assume that the number of parameters in the first and last layer are the same. How would you train a neural network that achieves this? Justify your answer (*at most three sentences*).

- ii. (X points) Your colleague returns and says “these alien brains make no sense! Now it turns out that activations in the 10th layer of the network are almost all entirely zero.” You want to incorporate this knowledge into your neural network design. How would you train a neural network that achieves this? Justify your answer (*at most three sentences*).

- (b) (X points) **Activation functions.** For each of these activation functions, state whether it (1) saturates (and if so, for what values does it saturate), (2) has the zig-zagging gradients problem (i.e., in a network using only this activation function, are gradient steps zig-zagging?), (3) is differentiable everywhere.

- i. ELU:  $x$  if  $x > 0$  and  $\alpha(\exp(x) - 1)$  if  $x \leq 0$ .
- ii. Exponential:  $\exp(x)$ .
- iii. Shifted logarithm:  $\max(0, \log(1 - x))$ .

**Solution:**

- (a) i. The goal is to have weights in the first and last layers capture the same features. One way to achieve this is to, by denoting  $\mathbf{w}_1$  and  $\mathbf{w}_k$  to be the first and last layers respectively, regularize the cost function by incorporating the penalty:

$$\lambda \|\mathbf{w}_1 - \mathbf{w}_k\|^2$$

- ii. The activations in layer 10 are sparse, so we should regularize the activations to be close to 0, i.e., incorporate the following penalty to the cost:

$$\lambda \|\mathbf{h}_{10}\|_1$$

where  $\mathbf{h}_{10}$  are the activations of unit  $\mathbf{h}_{10}$ .



- (b)
  - i. Saturates when  $x$  is very negative. Does not have the zig-zagging gradients problem. Is not differentiable at  $x = 0$ .
  - ii. Saturates when  $x$  is very negative. Has the zig-zagging gradients problem. Is differentiable.
  - iii. Saturates when  $x$  is very positive. Has the zig-zagging gradients problem. Is not differentiable at  $x = 0$ .

## 6. Short answer on convolutional neural networks.

- (a) (X points) **AlexNet.** AlexNet is an architecture that we've mentioned several times in class. In this question, you will calculate attributes of some AlexNet layers.
  - i. (X points) The 3rd, 4th, and 5th convolutional layers in AlexNet have the following properties:
    - 3rd Conv Layer: 384 filters that are  $3 \times 3$  applied at stride 1, pad 1.
    - 4th Conv Layer: 384 filters that are  $3 \times 3$  applied at stride 1, pad 1.
    - 5th Conv Layer: 256 filters that are  $3 \times 3$  applied at stride 1, pad 1.

The input to the 3rd conv layer is a tensor that is  $27 \times 27 \times 96$ . After the 5th conv layer, the output is a tensor that is  $27 \times 27 \times 256$ . Calculate the number of parameters in the 3rd, 4th, and 5th convolutional layers, not including biases. (To be clear, do not include biases in your calculation, but only the convolutional filter weights.)
  - ii. (X points) How many input neurons is each output neuron in the 3rd convolutional layer connected to? Do not include biases in this calculation.
  - iii. (X points) After the 5th Conv Layer, there are  $3 \times 3$  max pooling filters applied at stride 2. The output of this operation is  $6 \times 6 \times 256$ . Calculate the number of trainable parameters in the max pooling layer. Do not include biases in this calculation.
  - iv. (X points) After this pooling layer, the data representation is a  $6 \times 6 \times 256$  tensor. This data representation becomes the input to a fully connected layer with 4096 units. Calculate the number of parameters in this fully connected layer. Do not include biases in this calculation.
  - v. (X points) Consider the memory required to store the output of each layer. Assume all numbers are represented as singles (i.e., with 4 bytes). How much memory is required to store the output of the 5th Conv Layer, the pooling layer, and the fully connected layer?
- (b) (X points) **Receptive fields.** We define the receptive field as the spatial extent of the inputs (to the first layer) seen by a given output neuron (in the last layer). For example, if we use one convolutional layer with a  $3 \times 3$  filter applied at stride 1, then the effective receptive field is  $3 \times 3$ , since each output neuron is connected to a  $3 \times 3$  patch of the inputs. We would say that the receptive field is 3 (assume that the height and width of the filters are always matched). Assume the input is an  $m \times m \times d$  tensor. Consider the following questions:
  - i. What is the receptive field of a stack of three  $1 \times 1$  convolutional filters applied at stride 1?
  - ii. What is the receptive field of a fully connected layer?

- iii. What is the receptive field of a stack of three  $3 \times 3$  convolutional filters applied at stride 2? Assume the stride is lawful.

**Solution:**

- (a)
- i. The number of parameters in each layer is given as follows:
    - 3rd Conv Layer:  $3 \times 3 \times 96 \times 384 = 331,776$ .
    - 4th Conv Layer:  $3 \times 3 \times 384 \times 384 = 1,327,104$ .
    - 5th Conv Layer:  $3 \times 3 \times 384 \times 256 = 884,736$ .
  - ii. Each output neuron is connected to  $3 \times 3 \times 96$  input neurons, or 864.
  - iii. The max pooling layers have zero parameters.
  - iv. The FC layer has connections from every input to every output. Thus, it has  $6 \times 6 \times 256 \times 4096 = 37,748,736$  trainable parameters.
  - v.
    - The 5th conv layer output stores:  $256 \times 27 \times 27$  numbers, which is 186,624 numbers. This corresponds to 729 KB.
    - The pooling layer output stores:  $256 \times 6 \times 6$  numbers, which is 9,216 numbers. This corresponds to 36 KB.
    - The FC layer output stores 4096 numbers, which is 16 KB.
- (b)
- i. Each convolutional filter only sees one neuron at the input. So its receptive field is 1.
  - ii. The receptive field is  $m$ , since each output neuron is fully connected to every input.
  - iii. With one convolutional layer, the receptive field is 3. When stacking two, due to the stride, the receptive field is 7. Each output is connected to three inputs, and they overlap with their adjacent neuron only in one input. So the receptive field is  $3 \times 2 + 1 = 7$ . When stacking three, due to stride, the receptive field is  $7 \times 2 + 1 = 15$ .