

Closed book exam, no books, papers, notes, phones etc. allowed. The exam has 8 coding questions (labs), 8 open questions (lectures) and 14 multiple choice questions (papers). **Answer on the separate answer sheets.** Explain all answers, i.e. explicitly show intermediate steps to clarify if needed for coding / calculations / motivations / etc. Good luck!

## 1 Lab assignments (38pts)

1. **Question (4pts)** Implement the forward pass of a linear layer, followed by a ReLU activation function:

$$y = xW + b$$

$$y = \text{Sigmoid}(y)$$

Here  $x$  has dimensions [batch\_size, input\_channels]. Use the code template provided below. You may use elementary operations from the PyTorch library only, i.e. no predefined layers from `torch.nn`. Sub-questions:

A) Placeholder for layer weight and bias (1pts)

B) Forward pass (3pts)

```
1 import torch
2 class Linear(object):
3     """
4     Fully connected layer.
5     Args:
6         in_feat: number of input features
7         out_feat: number of output features
8     """
9
10    def __init__(self, in_feat, out_feat):
11        super(Linear, self).__init__()
12        """
13        # A) Define placeholder tensors for
14        # layer weight and bias. The placehol-
15        # der tensors should have the correct
16        # dimension according to the in_feat
17        # and out_feat variables.
18        """
19        # A) Your code goes here.
20        """
21        # END OF YOUR CODE
22        """
23
24        self.init_params() # Init. parameters
25
26    def init_params(self):
27        """
28        Initialize layer parameters. Sample
29        weight from Gaussian distribution and
30        bias uniform distribution.
31        """
32        self.weight = torch.randn_like(
33            self.weight)
34        self.bias = torch.rand_like(self.bias)
35
36    def forward(self, x):
37        """
38        Forward pass of Linear layer: multiply
39        input tensor by weights and add bias.
40        Args:
41            x: input tensor
42        Returns:
43            y: output tensor
44        """
45        """
46        # B) Forward pass.
```

```
47        """
48        # B) Your code goes here.
49        """
50        # END OF YOUR CODE
51        """
52
53        return y
```

2. **Question (4pts)** Implement the backward passes for the ReLU and Sigmoid non-linearities. The non-linearities are defined as follows:

$$\text{ReLU}(x) = \max(0, x)$$

$$\text{Sigmoid}(x) = \frac{1}{1 + \exp(-x)}$$

You may use elementary operations from PyTorch library only, i.e. no predefined layers from `torch.nn`. Sub-questions:

A) ReLU backward (2pts)

B) Sigmoid backward (2pts)

```
1 import torch
2 class ReLU(object):
3     """
4     ReLU non-linear activation function.
5     """
6
7    def __init__(self):
8        super(ReLU, self).__init__()
9
10    # Cache forward pass variables for use
11    # during backward pass.
12    self.cache = None
13
14    def forward(self, x):
15        """
16        Forward pass of ReLU non-linear
17        activation function: y=max(0,x). Store
18        input tensor as cache variable.
19        Args:
20            x: input tensor
21        Returns:
22            y: output tensor
23        """
24
25        y = torch.clamp(x, min=0) # forward pass
26        self.cache = y # update cache
27
28        return y
29
30    def backward(self, dupstream):
31        """
32        Backward pass of ReLU non-linear
33        activation function: return downstream
34        gradient dx.
35        Args:
36            dupstream: Gradient of loss with
37            respect to output of this layer.
38        Returns:
39            dx: Gradient of loss with respect to
40            input of this layer.
```

```

41     """
42
43     # Making sure that we don't modify the
44     # incoming upstream gradient
45     dupstream = dupstream.clone()
46
47     #####
48     # A) ReLU Backward pass
49     #####
50     # A) Your code goes here.
51     #####
52     # END OF YOUR CODE
53     #####
54
55     return dx
56
57 class Sigmoid(object):
58     """
59     Sigmoid non-linear activation function.
60     """
61
62     def __init__(self):
63         super(Sigmoid, self).__init__()
64
65         # Cache forward pass variables for use
66         # during backward pass.
67         self.cache = None
68
69     def forward(self, x):
70         """
71         Forward pass of Sigmoid non-linear
72         activation function:  $y=1/(1+\exp(-x))$ .
73         Store input tensor as cache variable.
74         Args:
75             x: input tensor
76         Returns:
77             y: output tensor
78         """
79
80         y = 1.0 / (1.0 + torch.exp(-x))
81         self.cache = y # update cache
82
83         return y
84
85     def backward(self, dupstream):
86         """
87         Backward pass of Sigmoid non-linear
88         activation function: return downstream
89         gradient dx.
90
91         Args:
92             dupstream: Gradient of loss with
93             respect to output of this layer.
94
95         Returns:
96             dx: Gradient of loss with respect to
97             input of this layer.
98         """
99
100         #####
101         # B) Sigmoid Backward pass
102         #####
103         # B) Your code goes here
104         #####
105         # END OF YOUR CODE
106         #####
107
108         return dx

```

3. Question (6pts) Implement the forward pass of a convolutional layer. You may use elementary operations from PyTorch library only, i.e. no pre-

defined layers from `torch.nn`. Sub-questions:

- A) Placeholder for weight and bias (1pt)
- B) The corresponding input window (2pts)
- C) Forward pass for batch, in for loop (3pts)

```

1 import torch
2 class Conv2d(object):
3     """
4     2D convolutional layer.
5     """
6
7     def __init__(self, in_channels, out_channels,
8                 kernel_size, stride=1, padding=0):
9
10        """
11        Initialize the layer with given params
12        Args:
13            in_channels: num. of input channels
14            out_channels: num. of output channels
15            kernel_size: kernel size, single int
16            stride: step size of conv. operation
17            padding: num pixels to zero-pad input
18        """
19
20        self.in_channels = in_channels
21        self.out_channels = out_channels
22        self.kernel_size = kernel_size
23        self.stride = stride
24        self.padding = padding
25
26        #####
27        # A) Create placeholder tensors for
28        # weight and bias with correct dims.
29        #####
30        # A) Your code goes here.
31        #####
32        # END OF YOUR CODE
33        #####
34
35        # Initialize parameters
36        self.init_params()
37
38    def init_params(self):
39        """
40        Initialize layer parameters. Sample
41        weight from Gaussian distribution
42        and initialize bias as zeros.
43        """
44
45        self.weight = torch.randn_like(
46            self.weight)
47        self.bias = torch.rand_like(self.bias)
48
49    def forward(self, x):
50        """
51        Forward pass of convolutional layer
52
53        Args:
54            x: input tensor (N, C, H, W)
55
56        Returns:
57            y: output tensor (N, F, H', W'),
58            where:
59                H' = 1 + (H + 2 * padding -
60                    kernel_size) / stride
61                W' = 1 + (W + 2 * padding -
62                    kernel_size) / stride
63        """
64
65        # Pad the input
66        x_padded = torch.nn.functional.pad(

```

```

66         x, [self.padding] * 4)
67
68     # Unpack the needed dimensions
69     N, _, H, W = x.shape
70
71     # Calculate output height and width
72     Hp = 1 + (H + 2 * self.padding -
73             self.kernel_size) // self.stride
74     Wp = 1 + (W + 2 * self.padding -
75             self.kernel_size) // self.stride
76
77     # Create an empty output to fill in
78     y = torch.empty(
79         (N, self.out_channels, Hp, Wp))
80
81
82     #####
83     # B-C) Compute the output y by looping
84     # over its height and width dimensions
85     # and defining an input window over x
86     # called window, using which you
87     # calculate the output for each input
88     # sample by convolving it with the
89     # weight and adding a bias.
90     #####
91     for i in range(Hp):
92         for j in range(Wp):
93             # B) Make window.
94             # Your code goes here.
95             for k in range(N):
96                 # C) Forward pass.
97                 # Your code goes here.
98             #####
99     # END OF YOUR CODE
100    #####
101
102    # Cache input to use in backward pass
103    self.cache = x_padded
104
105    return y

```

4. Question (4pts) Implement the gradient update of SGD with momentum. The formula for a gradient update with momentum is:

$$v_i = \rho v_{i-1} + (1 - \rho) \nabla_{\theta}$$

$$\theta' = \theta - \epsilon v_i$$

```

1 def momentum(X, rho, learning_rate, prev_value,
2             Grad=Grad.f):
3     """
4     Gradient descent with momentum optimization
5     step.
6     Args:
7     X: Current value of objective function.
8     rho: Optimization hyperparameter -
9         see formula above.
10    learning_rate: Optimization step size.
11    prev_value: Momentum parameter from
12        previous iteration.
13    Grad: Function that returns gradient
14        of objective function.
15    """
16    # Gradient of current value
17    gradient = Grad(*X)
18    # Momentum parameter
19    v = 0
20    # Momentum parameter from previous iteration
21    v_prev = prev_value
22    #####

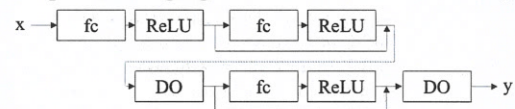
```

```

22     # A) Create gradient descent with
23     # momentum: update v and X.
24     #####
25     # A) Your code goes here
26     #####
27     # END OF YOUR CODE
28     #####
29
30     return X, v

```

5. Question (4pts) Implement the Dropout layers and the forward pass of the network architecture depicted in the figure (fc := fully connected layer; DO:= dropout; merging of arrows := summation):



Sub-questions:

- A) Initialize Dropout layers (1pts)  
B) Implement forward pass (3pts)

```

1 import torch.nn as nn
2 class FCNet_do(nn.Module):
3     """
4     Simple fully connected neural network with
5     residual connections and dropout
6     layers in PyTorch. Layers are defined in
7     __init__ and forward pass
8     implemented in forward.
9     """
10    def __init__(self):
11        super(FCNet_do, self).__init__()
12
13        self.fc1 = nn.Linear(16, 32)
14        self.fc2 = nn.Linear(32, 32)
15        self.fc3 = nn.Linear(32, 64)
16
17        #####
18        # A) Define dropout layers.
19        # A) Your code goes here.
20        # END OF YOUR CODE
21        #####
22
23    def forward(self, x):
24        #####
25        # B) Implement forward pass as in figure
26        # B) Your code goes here
27        # END OF YOUR CODE
28        #####
29        return y

```

6. Question (4pts) Implement the vanilla RNN. [Hint]: think about all the architectural hyperparameters and especially sizes of the input, hidden states, weights and outputs. Sub-questions:

- A) Parameter initialization (2pts)  
B) Forward pass implementation (2pts)

```

1 class VanillaRNN(nn.Module):

```

```

2 """
3 Vanilla recurrent neural network which
4 has the following update rule:
5     ht = tanh(W_xh * xt + b_xh + W_hh *
6         h(t-1) + b_hh)
7 """
8 def __init__(self, input_size, hidden_size):
9     super(VanillaRNN, self).__init__()
10
11     self.hidden_size = hidden_size
12
13     #####
14     # A) Create weight and bias tensors
15     # with correct sizes. NOTE: Don't
16     # forget to encapsulate weights and
17     # biases in nn.Parameter.
18     #####
19     # A) Finish the following lines
20     self.weight_xh = # input-to-hidden w
21     self.weight_hh = # hidden-to-hidden w
22     self.bias_xh = # input-to-hidden b
23     self.bias_hh = # hidden-to-hidden b
24     #####
25     # END OF YOUR CODE
26     #####
27
28 def forward(self, x):
29     """
30     Args:
31         x: input with shape (N, T, D):
32             N: number of samples
33             T: number of time steps
34             D: input size, equal to
35               self.input_size.
36     Returns:
37         y: output with shape (N, T, H):
38             H: hidden size
39     """
40
41     # Transpose input for efficient
42     # vectorized calculation. After
43     # transposing # the input will have
44     # shape (T, N, D).
45     x = x.transpose(0, 1)
46
47     # Unpack dimensions
48     T, N = x.shape[0], x.shape[1]
49
50     # Initialize hidden states to zero.
51     h0 = torch.zeros(N, self.hidden_size)
52
53     # Define a list to store outputs.
54     y = []
55
56     #####
57     # B) Implement the RNN forward pass
58     #####
59     # B) Your code goes here.
60     #####
61     # OF YOUR CODE
62     #####
63
64     # Stack outputs to shape (T, N, H)
65     y = torch.stack(y)
66
67     # Switch time and batch dimension
68     # (T, N, H) -> (N, T, H)
69     y = y.transpose(0, 1)
70
71     return y

```

7. **Question** (6pts) Implement the basic self-attention layer in PyTorch. Recall that the input serve three different purposes in the self-attention operation, namely as the query, key and value:

$$\mathbf{q}_i = \mathbf{W}_q \mathbf{x}_i, \quad \mathbf{k}_i = \mathbf{W}_k \mathbf{x}_i, \quad \mathbf{v}_i = \mathbf{W}_v \mathbf{x}_i$$

And the basic self-attention operation:

$$w'_{ij} = \frac{\mathbf{q}_i^T \mathbf{k}_j}{\sqrt{k}}, \quad w_{ij} = \text{softmax}(w'_{ij}),$$

$$\mathbf{y}_i = \sum_j w_{ij} \mathbf{v}_j$$

Sub-questions:

- A) Transform input tensor to queries, keys and values (2pts)
- B) Compute weights, scale weights, apply softmax (2pts)
- C) Compute output tensor (2pts)

```

1 import torch.nn as nn
2 class SelfAttention(nn.Module):
3     """
4     Self-attention operation with learnable key,
5     query and value embeddings.
6     Args:
7         d: embedding dimension
8     """
9     def __init__(self, k):
10         super(SelfAttention, self).__init__()
11
12         # Linear map to queries, keys, values
13         self.k = nn.Linear(d, d, bias=False)
14         self.q = nn.Linear(d, d, bias=False)
15         self.v = nn.Linear(d, d, bias=False)
16
17     def forward(self, x):
18
19         # Get tensor dimensions: batch size,
20         # sequence length and embedding
21         # dimension.
22         b, t, d = x.size()
23
24         #####
25         # A) Perform self-attention operation
26         # with learnable query, key and value
27         # mappings.
28         # B) Calculate w_prime, apply scaling
29         # and softmax.
30         # C) Compute the output tensor y.
31         #####
32         # Your code goes here.
33         #####
34         # END OF YOUR CODE
35         #####
36
37     return y

```

8. **Question** (6pts) Implement the auto-encoder. The whole auto-encoder is implemented in the code block below, which consists of 3 classes: Encoder, Decoder and Autoencoder. Sub-questions:

- A) Create necessary layers for Encoder (2pts)
- B) Finish forward pass of Encoder (1pts)
- C) Implement forward pass of Decoder (2pts)
- D) Implement forward pass of Autoencoder (1pts)

```

1 import torch.nn as nn
2 class Encoder(nn.Module):
3     """
4     Encoder, projects input to latent space.
5     Args:
6         l_dim: latent space dimensionality (int)
7         s_img: size of square input image (int)
8         h_dim: dim. of hidden layers (list)
9     """
10    def __init__(self, l_dim, s_img, h_dim):
11        super(Encoder, self).__init__()
12
13        #####
14        # A) Create the necessary layers
15        #####
16        # A) Your code goes here.
17        #####
18        # END OF YOUR CODE
19        #####
20
21    def forward(self, x):
22
23        x = torch.flatten(x, start_dim=1)
24        x = self.relu(self.linear1(x))
25        x = self.relu(self.linear2(x))
26
27        #####
28        # B) Apply final layer of the encoder
29        # [Hint]: do you need a non-linearity?
30        #####
31        # B) Your code goes here.
32        #####
33        # END OF YOUR CODE
34        #####
35
36        return x
37
38 class Decoder(nn.Module):
39    """
40    Decoder, projects latent space to image.
41    Args:
42        l_dim: latent space dimensionality (int)
43        s_img: size of square input image (int)
44        h_dim: dim. of hidden layers (list)
45    """
46    def __init__(self, l_dim, s_img, h_dim):
47        super(Decoder, self).__init__()
48
49        self.linear1 = nn.Linear(l_dim,
50                                  h_dim[1])
51        self.linear2 = nn.Linear(h_dim[1],
52                                  h_dim[0])
53        self.linear3 = nn.Linear(h_dim[0],
54                                  s_img*s_img)
55
56        self.relu = nn.ReLU()
57        self.sigmoid = nn.Sigmoid()
58
59    def forward(self, z):

```

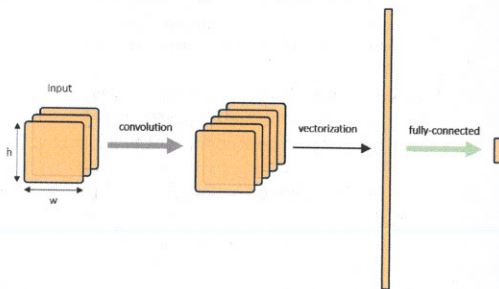
```

60        #####
61        # C) Implement forward pass of decoder
62        # [Hint]: Have a close look at the
63        # forward pass of the encoder.
64        #####
65        # C) Your code goes here.
66        #####
67        # END OF YOUR CODE
68        #####
69
70        return z
71
72 class Autoencoder(nn.Module):
73     """
74     Autoencoder model.
75     Args:
76         l_dim: latent space dimensionality (int)
77         s_img: size of square input image (int)
78         h_dim: dim. of hidden layers (list)
79     """
80    def __init__(self, l_dim, s_img,
81                  h_dim = [256, 128]):
82        super(Autoencoder, self).__init__()
83
84        self.encoder = Encoder(l_dim, s_img,
85                                h_dim)
86        self.decoder = Decoder(l_dim, s_img,
87                                h_dim)
88
89    def forward(self, x):
90
91        #####
92        # D) Autoencoder forward pass.
93        #####
94        # D) Your code goes here.
95        #####
96        # END OF YOUR CODE
97        #####
98
99        return y

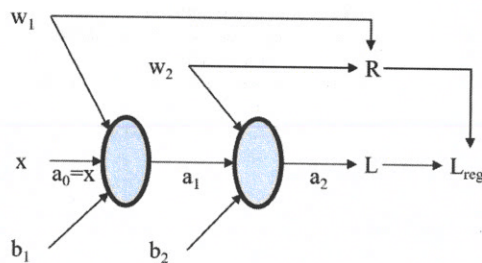
```

## 2 Lectures (38pt)

1. **Question (4pts)** Consider a 2-layered network. The first layer is convolutional with a kernel size of  $3 \times 3$ , with padding = 1, stride = 1. The second layer is a fully connected layer. The number of input channels is 3, the number of hidden channels is 5. The amount of output features is 2. The network uses bias terms. The size of an input image is:  $h \times w = 16 \times 16$ . The image below clarifies the network further. **What is the total amount of learnable parameters?** Motivate your answer with a detailed step-by-step approach. (Don't forget the bias terms).



2. **Question (6pts)** Consider the following scalar feedforward neural network with input  $x$  and target  $y$ .



$L = (a_2 - y)^2$  and  $R = \sum_i w_i^2$  are the loss and regularization term respectively and the regularized loss is computed as  $L_{\text{reg}} = L + \lambda R$ , in which  $\lambda$  is some predefined weighting hyper-parameter. Furthermore the activations  $a_i$  are computed as

$$a_i = \sigma(z_i) = \sigma(w_i a_{i-1} + b_i),$$

in which  $\sigma$  is some arbitrary continuously differentiable non-linear activation function (e.g. sigmoid, tanh). We are interested in computing  $\frac{\partial L_{\text{reg}}}{\partial w_1}$ . We have already computed  $\frac{\partial L_{\text{reg}}}{\partial w_2}$  and  $\frac{\partial L_{\text{reg}}}{\partial b_2}$  for which we used  $\bar{z}_2 = \frac{\partial L_{\text{reg}}}{\partial z_2}$ . Efficient backpropagation uses this intermediate result for the computation of gradients of earlier layers. Compute  $\frac{\partial L_{\text{reg}}}{\partial w_1}$  based on this result. Motivate your answer (do not only write down the formula). If needed, use  $\sigma'()$  for the derivative of the non-linear activation function.

3. **Question (6pts)** We would like to perform a convolution on an  $5 \times 5 \times 2$  image (ie: 2 channels, so 2 values per pixel) with an "average value kernel with horizontal and vertical spatial size 3" (no padding) and a stride of 1. The value of the input image  $X$  is given below:

$$X[:, :, 0] = \begin{bmatrix} 1 & 0 & 2 & 6 & 4 \\ 0 & 2 & 1 & 2 & 5 \\ 2 & 1 & 0 & 4 & 2 \\ 6 & 2 & 4 & -4 & -3 \\ -4 & -2 & 0 & 3 & -2 \end{bmatrix}$$

$$X[:, :, 1] = \begin{bmatrix} 0 & -2 & 1 & 2 & 2 \\ -3 & 0 & -2 & 3 & 2 \\ 1 & -2 & -2 & 2 & 2 \\ 0 & 4 & 4 & -1 & -1 \\ 2 & 2 & 0 & 3 & -2 \end{bmatrix}$$

(2pt) 1. Give the kernel matrix.

(4pt) 2. Compute output of the convolution.

4. **Question (4pts)** About regularization.

(1pt) 1. What is underfitting?

(1pt) 2. What is overfitting?

(1pt) 3. How to detect underfitting?

(1pt) 4. How to detect overfitting?

5. **Question (6pts)** The exponentially weighted moving average (EWMA) is used in many optimization methods. The equation is given by:  $S_t = (\rho S_{t-1}) + (1 - \rho)y_t$ .

Consider a case where the gradient has a constant value of 1.0, over all considered training samples. Let  $\rho = 0.95$ .

(2pt) 1. Compute  $S_2$ .

The bias correction equation is:  $\hat{S}_t = \frac{S_t}{1 - \rho^t}$ .

(2pt) 2. Compute the bias corrected  $\hat{S}_2$ .

(1pt) 3. What is the goal of bias correction?

(1pt) 4. What role does bias correction play after a large number of iterations?

6. **Question (4pts)** Can you explain why vanishing/exploding gradients happen in RNNs? What changes can be made to RNNs to solve these problems?

7. **Question (4pts)** Regarding self-attention.

(1pt) 1. What is the computational advantage of self-attention over recurrent networks?

(1pt) 2. What is the purpose of having multiple attention heads?

(1pt) 3. What is the difference between narrow and wide multi-head self-attention?

(1pt) 4. How does the number of parameters in a self-attention layer change if the input sequence length doubles?

8. **Question** (4pts) Explain the GAN objective function given below:

$$\arg \min_{\theta_G} \max_{\theta_D} \mathbb{E}_{x \sim p_{\text{data}}} \log D_{\theta_D}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_D}(G_{\theta_G}(z)))$$

### 3 Papers (14pts)

Each paper has a multiple-choice question. Select the single best fitting answer per question. Each question counts for 1 point.

- Paper: *A Step Toward Quantifying Independently Reproducible Machine Learning Research*. What statement is mentioned in the paper:
  - Without releasing code it is very difficult to reproduce ML research.
  - Code should be reviewed in the scientific peer-review process.
  - Its OK if approximately 25% of the claims in a paper cannot be reproduced.
  - Reproducibility should become part of university courses.
- Paper: *Troubling Trends in Machine Learning Scholarship*. The paper mentions "Explanation vs. Speculation". What exactly is meant by that?
  - Papers are intentionally made overly complex.
  - There is too much math
  - There is not sufficient explanation for some claims
  - The method is not explained.
- Paper: *Do ImageNet Classifiers Generalize to ImageNet?* So, do they generalize?
  - Yes, although the accuracy is reduced, the ranking between models stays the same
  - Yes, because it showed to generalize to other, slightly changed, domains.
  - No, because the technical term "generalization" only applies to the original data
  - No, because there was insufficient data
- Paper: *Scaling down deep learning*. According to the paper, why scale down?
  - Complex models are too difficult to interpret
  - Results on small-scale datasets often translate to large-scale datasets
  - Scaling down is the only option for smaller labs and universities
  - Scaling down is essential for scientific reproductions
- Paper: *Highway and Residual Networks learn Unrolled Iterative Estimation*. What best describes a residual connection?
  - Add the output of the previous layer to the input of the current layer
  - Add the output of the previous layer to the output of the current layer
  - Add the input of the previous layer to the output of the current layer
  - Add the input of the previous layer to the input of the current layer
- Paper: *ResNet strikes back: An improved training procedure in timm*.
  - An independent reproduction of a standard ResNet with an improved training procedure improves over the original paper
  - An independent reproduction of a modified ResNet with an improved training procedure improves over the original paper
  - An independent reproduction of a standard ResNet with an improved training procedure cannot improve over the original paper
  - An independent reproduction of a modified ResNet with an improved training procedure cannot improve over the original paper
- Paper: *Deep Image Prior*. What statement fits best with the paper:
  - Neural network architectures are a form of prior knowledge
  - ConvNets are a form of prior knowledge for images
  - Optimization is less important as the network architecture
  - ConvNets can learn to turn noise to images
- Paper: *Approximating CNNs with Bag-of-local-Features models works surprisingly well on ImageNet*. Why is it called a BagNet?
  - Its based on unordered bagging and boosting bootstrapping
  - The receptive field is artificially unordered
  - The class evidence used for explainability is bagged across unordered classes
  - It treats images as a unordered bag of patches

9. Paper: *Group Normalization*. The statistics are computed over groups of pixels. The grouping is done by:
- A) Grouping spatially
  - B) Grouping feature maps
  - C) Grouping images
  - D) Grouping weights
10. Paper: *Torch.manual\_seed(3407) is all you need: On the influence of random seeds in deep learning architectures for computer vision*.
- A) Black swans can be found, although the variance is small
  - B) Black swans cannot be found because the variance is small
  - C) Black swans can be found as the variance is large
  - D) Black swans cannot be found, although the variance is large
11. Paper: *Attention Is All You Need*. Select the best answer.
- A) The query, key and value are shared over the inputs, similar to a 1x1 CNN
  - B) The connections between the network layers are gated, similar to a GRU
  - C) In a self-attention layer, information flows sequentially between tokens, similar to an LSTM.
  - D) Without positional embeddings, it would overfit on the word position, similar to a MLP.
12. Paper: *Perceiver IO: A General Architecture for Structured Inputs & Outputs*.
- A) Uses standard self-attention, and it is outperformed by specialized solutions
  - B) Uses a special form of self-attention, and it is outperformed by specialized solutions
  - C) Uses standard self-attention and it outperforms specialized solutions
  - D) Uses a special form of self-attention and it outperforms specialized solutions
13. Paper: *Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks*. What does the "cycle" refer to?
- A) Cyclic activation functions
  - B) Training cycles
  - C) Cycling between the annotated image pairs
  - D) Cycling between domains
14. Paper: *GANORCON: Are Generative Models Useful for Few-shot Segmentation?* So, are GANs useful for few-shot segmentation?
- A) Yes, more useful than contrastive learning.
  - B) Yes, but not as useful as contrastive learning
  - C) No, but contrastive learning is
  - D) No, and contrastive learning also is not