

# CS 446/ECE 449: Machine Learning

A. G. Schwing

University of Illinois at Urbana-Champaign, 2020

## Scribe & Exercises

## L8: Deep Neural Networks

## Goals of this lecture

- Getting to know deep nets
- Understanding forward and backward pass
- Learning about deep net components

## Reading material

- I. Goodfellow et al.; Deep Learning; Chapters 6-9

Our current framework:

$$\min_{\mathbf{w}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{i \in \mathcal{D}} \epsilon \ln \sum_{\hat{y}} \exp \frac{L(y^{(i)}, \hat{y}) + \mathbf{w}^T \psi(x^{(i)}, \hat{y})}{\epsilon} - \mathbf{w}^T \psi(x^{(i)}, y^{(i)})$$

What is a possible issue/limitation?

Linearity in the feature space  $\psi(x, y)$ . Fix: use kernels. But still learning a model **linear** in the parameters  $\mathbf{w}$

How to fix this?

Replace  $\mathbf{w}^T \psi(x, y)$  with a general function  $F(\mathbf{w}, x, y)$

$$\min_{\mathbf{w}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{i \in \mathcal{D}} \epsilon \ln \sum_{\hat{y}} \exp \frac{L(y^{(i)}, \hat{y}) + F(\mathbf{w}, x^{(i)}, \hat{y})}{\epsilon} - F(\mathbf{w}, x^{(i)}, y^{(i)})$$

General framework:

$$\min_{\mathbf{w}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{i \in \mathcal{D}} \epsilon \ln \sum_{\hat{y}} \exp \frac{L(y^{(i)}, \hat{y}) + F(\mathbf{w}, x^{(i)}, \hat{y})}{\epsilon} - F(\mathbf{w}, x^{(i)}, y^{(i)})$$

How to get to

- **Logistic regression**
- **Binary SVM**
- **Multiclass regression**
- **Multiclass SVM**
- **Deep Learning**

## Deep Learning:

What function  $F(\mathbf{w}, x, y) \in \mathbb{R}$  to choose?

- Choose any differentiable composite function

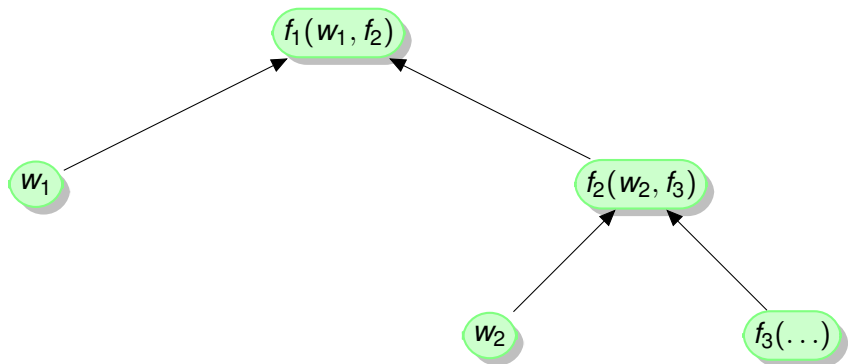
$$F(\mathbf{w}, x, y) = f_1(\mathbf{w}_1, y, f_2(\mathbf{w}_2, f_3(\dots f_n(\mathbf{w}_n, x) \dots)))$$

- More generally: functions can be represented by an acyclic graph (computation graph)

Example:

$$F(\mathbf{w}, x, y) = f_1(w_1, f_2(w_2, f_3(\dots)))$$

Nodes are weights, data, and functions:



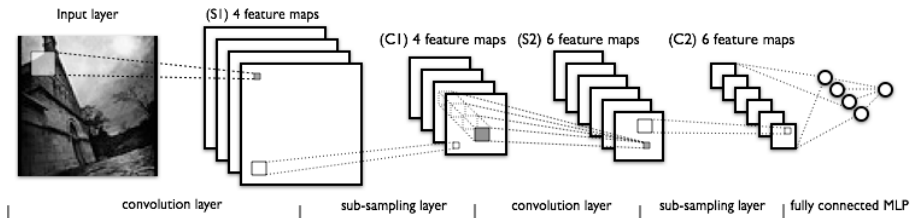
Internal representation used by deep net packages.



What are the individual functions/layers  $f_1, f_2$  etc.?

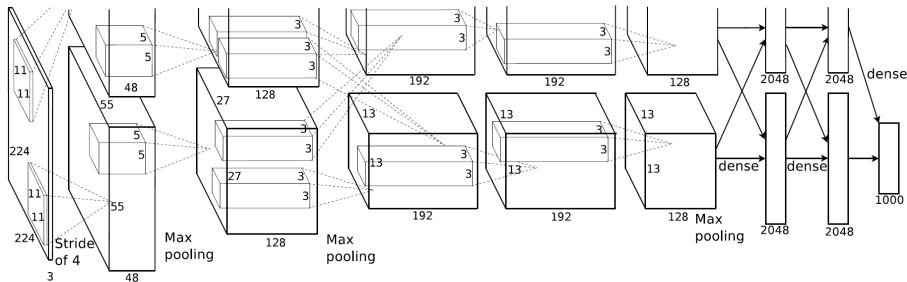
- Fully connected layers
- Convolutions
- Rectified linear units (ReLU):  $\max\{0, x\}$
- Maximum-/Average pooling
- Soft-max layer
- Dropout

## Example function architecture: LeNet



Decreasing spatial resolution and the increasing number of channels

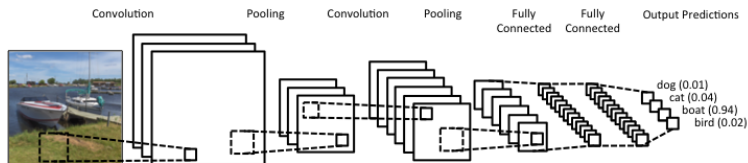
## Example function architecture: AlexNet



Decreasing spatial resolution and the increasing number of channels

Why is the output 1000-dimensional?

## Another deep net:



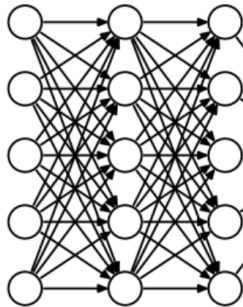
Those nets are structurally simple in that a layer's output is used as input for the next layer. This is not required.

Fully connected layer:

$$\mathbf{W}\mathbf{x} + \mathbf{b}$$

Trainable parameters  $\mathbf{w}$ :

- Matrix
- Bias



What's an issue with fully connected layers?

Issue with fully connected layers:

- Suppose the input is an image of size  $256 \times 256$
- Let the output of this layer have identical size
- How many weights are necessary?

$$2^{32} = 4,294,967,296$$

- How to share weights?

## Convolutions:

120	190	140	150	200
17	21	30	8	27
89	123	150	73	56
10	178	140	150	18
190	14	76	69	87

x

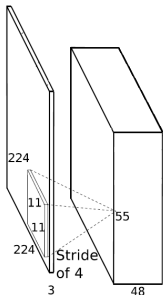
1/9	1/9	1/9
1/9	1/9	1/9
1/9	1/9	1/9

=

	98	98	93	
	84	97	72	
	108	108	91	

Trainable parameters  $\mathbf{w}$ :

- Filters (width, height, depth, number)
- Bias



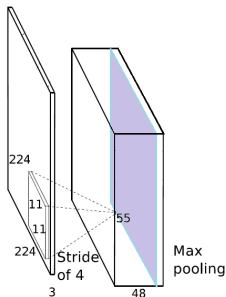
$$\text{width} + 2 * \text{padding} + \dots + 1 \quad 55$$

## Maximum-/Average pooling

Maximum or average over a spatial region

Trainable parameters  $\mathbf{w}$ :

- None





Soft-max layer:

$$x \longrightarrow \frac{\exp x_i}{\sum_j \exp x_j}$$

Trainable parameters  $\mathbf{w}$ :

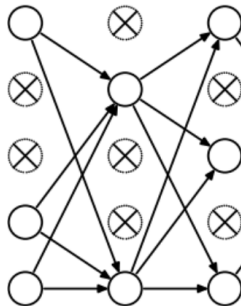
- None

## Dropout layer:

Randomly set activations to zero

Trainable parameters  $\mathbf{w}$ :

- None



Deep net training:

$$\min_{\mathbf{w}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{i \in \mathcal{D}} \ln \sum_{\hat{y}} \exp F(\mathbf{w}, x^{(i)}, \hat{y}) - F(\mathbf{w}, x^{(i)}, y^{(i)})$$

Often also referred to as maximizing the regularized cross entropy: ground truth

$$\max_{\mathbf{w}} -\frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{i \in \mathcal{D}} \sum_{\hat{y}} p_{\text{GT}}^{(i)}(\hat{y}) \ln p(\hat{y} | x^{(i)}) \quad \text{with} \quad \begin{cases} p_{\text{GT}}^{(i)}(\hat{y}) = \delta(\hat{y} = y^{(i)}) \\ p(\hat{y} | x) \propto \exp F(\mathbf{w}, x, \hat{y}) \end{cases}$$

What is  $C$ ? Weight decay (aka regularization constant)

$$\min_{\mathbf{w}} \underbrace{\frac{C}{2} \|\mathbf{w}\|_2^2}_{\text{weight decay}} - \underbrace{\sum_{i \in \mathcal{D}} \sum_{\hat{y}} p_{\text{GT}}^{(i)}(\hat{y}) \ln p(\hat{y} | x^{(i)})}_{\text{torch.nn.CrossEntropyLoss(gt, F)}}$$

Deep learning choices:

- Design a composite function  $F(\mathbf{w}, x, y)$
- Use an appropriate loss function

Know what you are doing, i.e., know all the dimensions.

## Many more loss functions:

- **CrossEntropyLoss**

$$\begin{aligned}\text{loss}(x, \text{class}) &= -\log(\exp(x[\text{class}]) / (\sum_j \exp(x[j]))) \\ &= -x[\text{class}] + \log(\sum_j \exp(x[j]))\end{aligned}$$

- **NLLLoss**

$$\text{loss}(x, \text{class}) = -x[\text{class}]$$

log-likelihood, x is  
log softmax probability

- **MSELoss**

$$\text{loss}(x, y) = 1/n \sum_i |x_i - y_i|^2$$

- **BCELoss**

$$\text{loss}(o, t) = -1/n \sum_i (t[i] * \log(o[i]) + (1-t[i]) * \log(1-o[i]))$$

- **BCEWithLogitsLoss**

$$\begin{aligned}\text{loss}(o, t) &= -1/n \sum_i (t[i] * \log(\text{sigmoid}(o[i])) \\ &\quad + (1-t[i]) * \log(1-\text{sigmoid}(o[i])))\end{aligned}$$

- **L1Loss**

- **KLDivLoss**

Why this form for the NLLLoss?

```
loss(x, class) = -x[class]
```

Intended to be used in combination with ‘LogSoftmax’:

$$f_i(x) = \log \frac{\exp x_i}{\sum_j \exp x_j}$$

Why? Numerical robustness (‘log-sum-exp trick’)

$$\log \sum_j \exp x_j = c + \log \sum_j \exp (x_j - c)$$

Don't try without, it **will** fail!

## Example (PyTorch.py):

```
class Net(nn.Module):
    def __init__(self):
        super(Net, self).__init__()
        self.conv1 = nn.Conv2d(1, 6, 5)
        self.conv2 = nn.Conv2d(6, 16, 5)
        self.fc1 = nn.Linear(16 * 5 * 5, 120)
        self.fc2 = nn.Linear(120, 84)
        self.fc3 = nn.Linear(84, 10)

    def forward(self, x):
        x = F.max_pool2d(F.relu(self.conv1(x)), (2, 2))
        x = F.max_pool2d(F.relu(self.conv2(x)), 2)
        x = x.view(-1, self.num_flat_features(x))
        x = F.relu(self.fc1(x))
        x = F.relu(self.fc2(x))
        x = self.fc3(x)
        return x
```

What are the input dimensions?

## Quiz:

- What are deep nets?
- How do deep nets relate to SVMs and logistic regression?
- What components of deep nets do you know?
- What algorithms are used to train deep nets?



## Important topics of this lecture

- Deep nets
- Relationship to SVMs and logistic regression
- Components of deep nets

## Up next:

- Backpropagation