

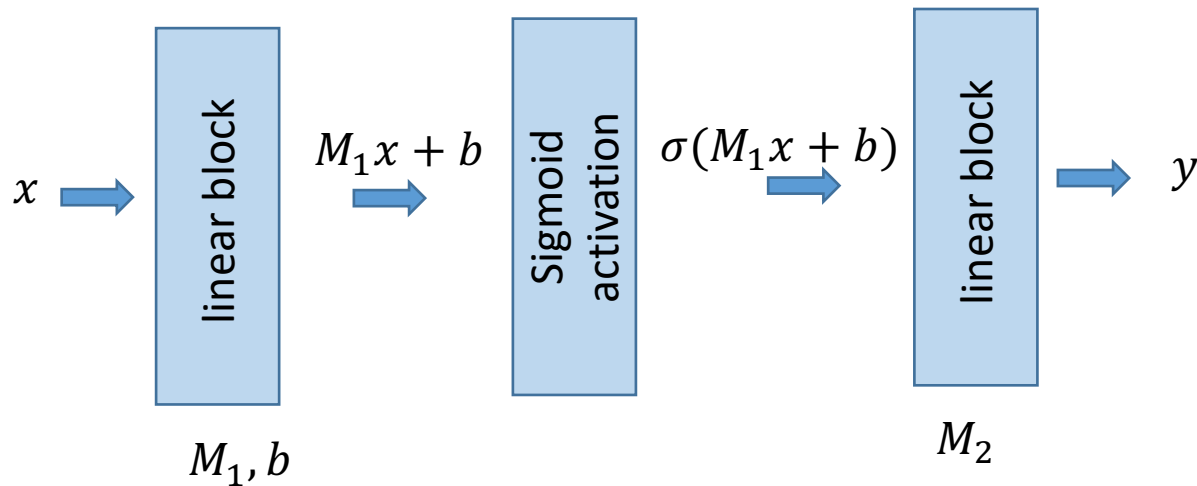


Neural Network Basics II

Today

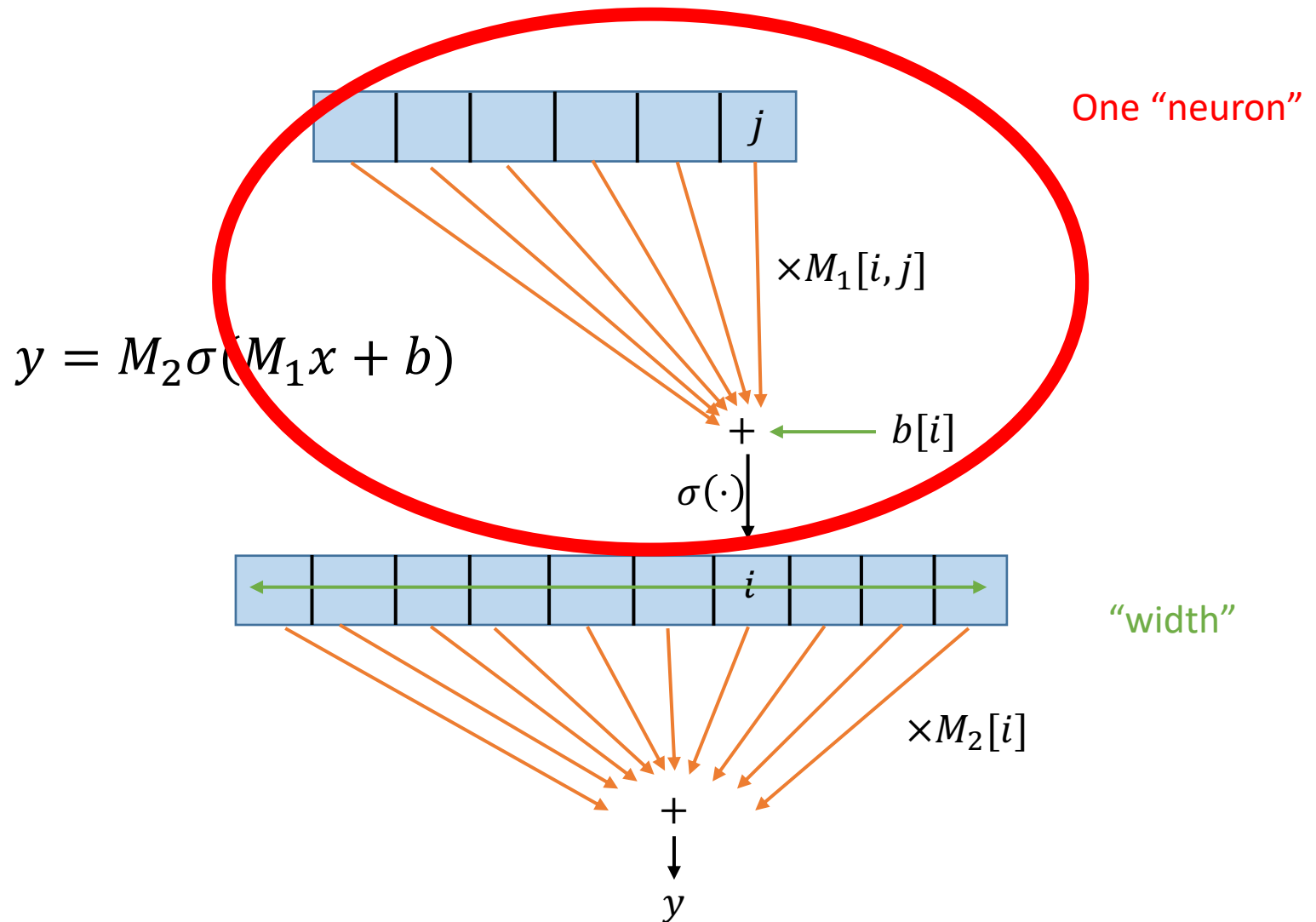
- Multilayer perceptron
- Activation functions
- Surrogate loss functions
- A closer look at softmax/cross entropy loss
- Short intro to PyTorch

Last time: Block Diagrams



- Activation function $\sigma(x)$ is applied per-coordinate.

Last time: Another Visualization



Last Time: Universal Approximation

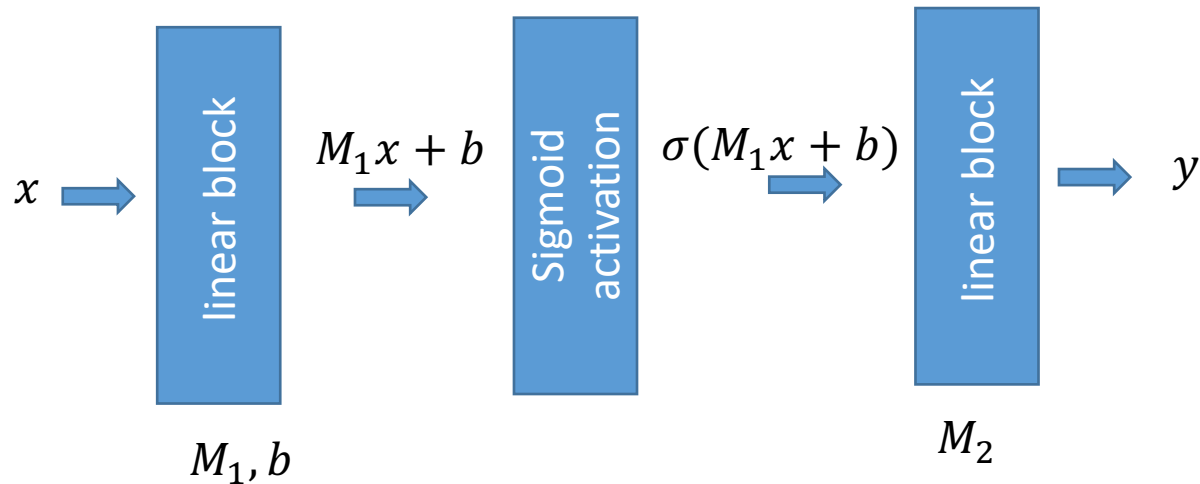
- Given any function $f: [-1,1]^d \rightarrow \mathbb{R}$, any $\epsilon \geq 0$, and any non-polynomial activation function σ , there is some sufficiently large h such that there is a single-layer neural network $W_2 \sigma(W_1 x + b)$ with activation σ , hidden size h such that for all $x \in [-1,1]^d$:

$$|f(x) - W_2 \sigma(W_1 x + b)| \leq \epsilon$$

- In practice, h might have to be absurdly large in order for this theorem to be useful.

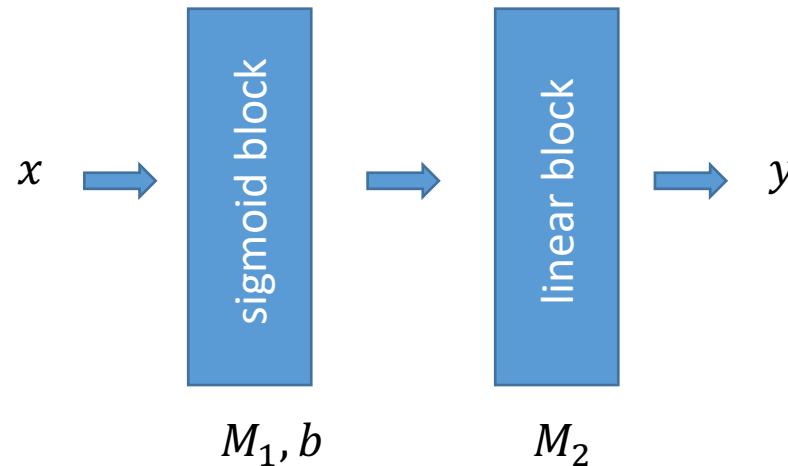
Multi-Layer Perceptron (MLP)

- Last time we talked about a *single layer perceptron*:



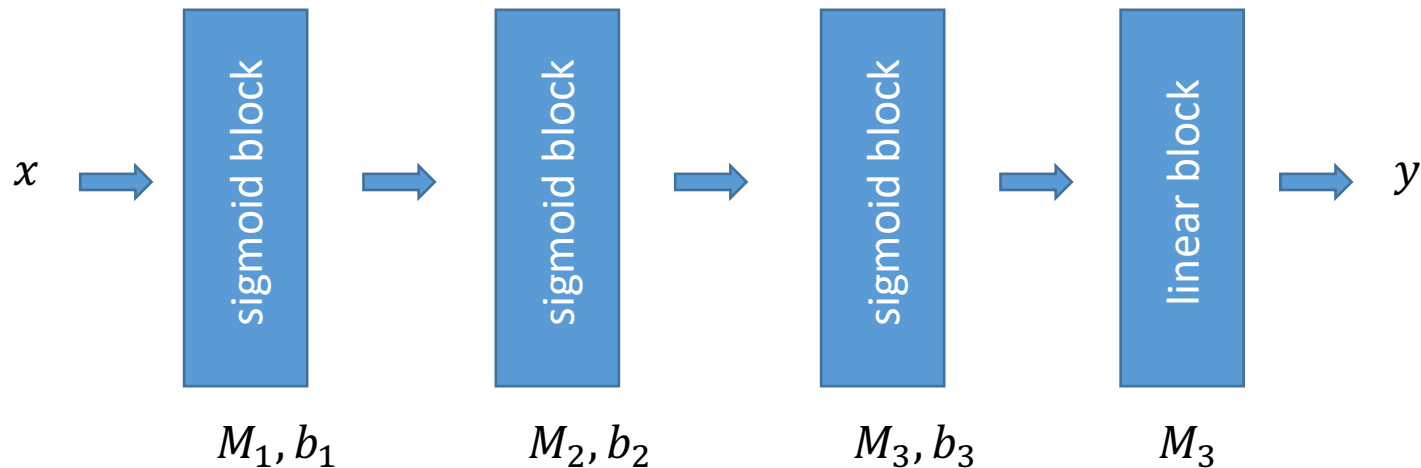
Multi-Layer Perceptron (MLP)

- A more concise diagram of a *single layer perceptron*:



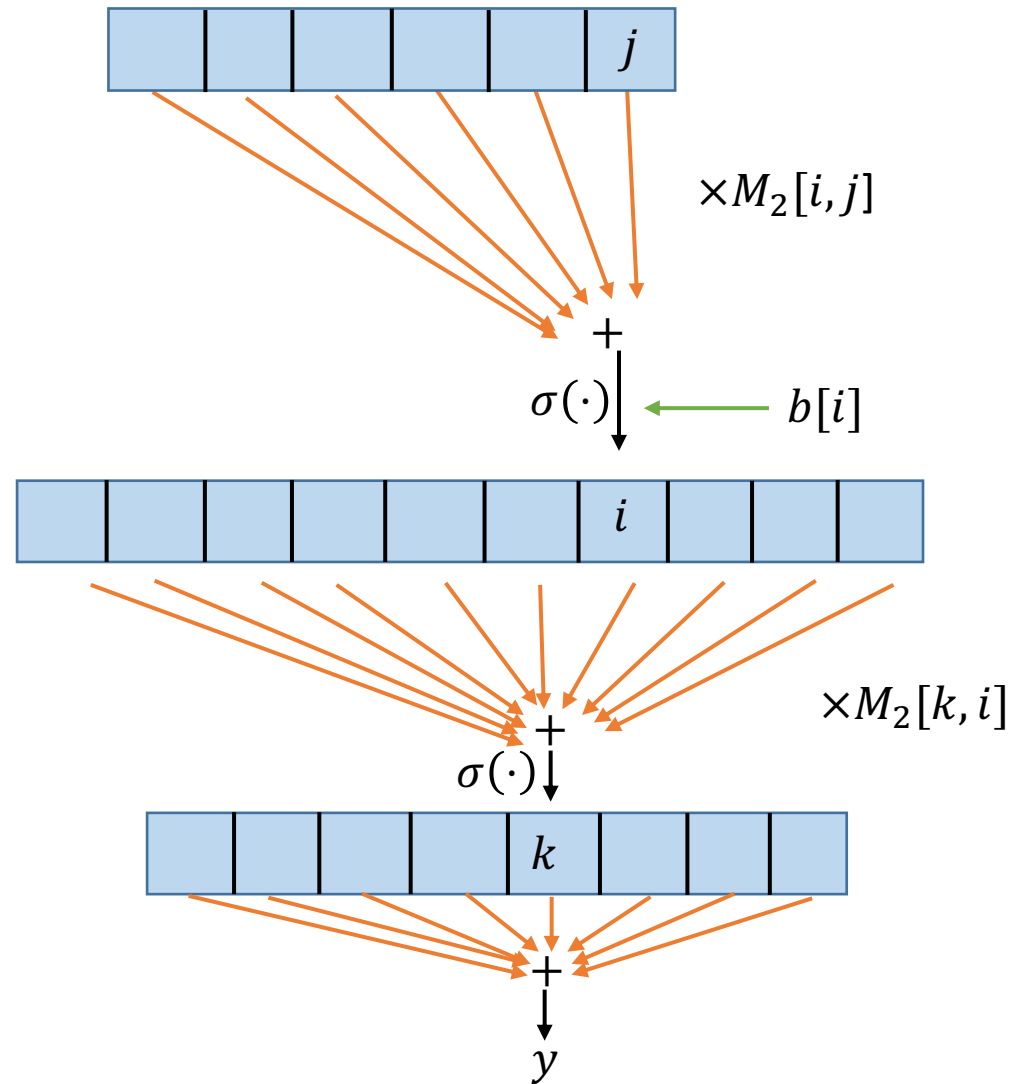
Multi-Layer Perceptron (MLP)

- The MLP just does this many times in a row:



- This is called **depth**. “Deep learning” means using deep architectures.

Multi-Layer Perceptron





Activation Functions

Activation Functions

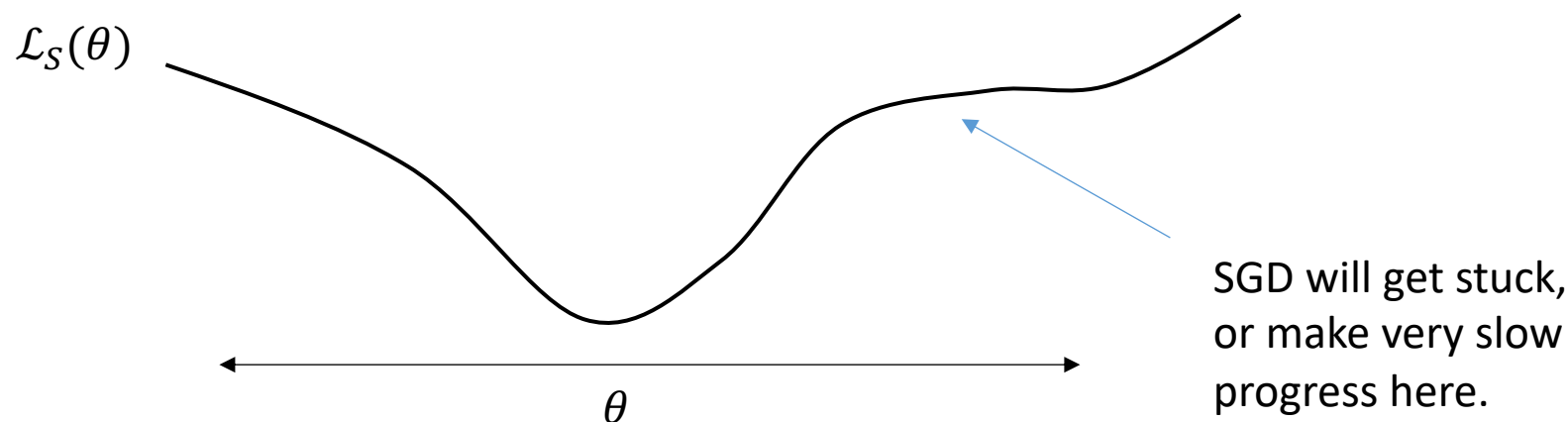
- Last class we covered the *sigmoid activation*:

$$\sigma(x) = \frac{\exp(x)}{1 + \exp(x)}$$

- This activation is no longer so popular because of *vanishing gradients*.

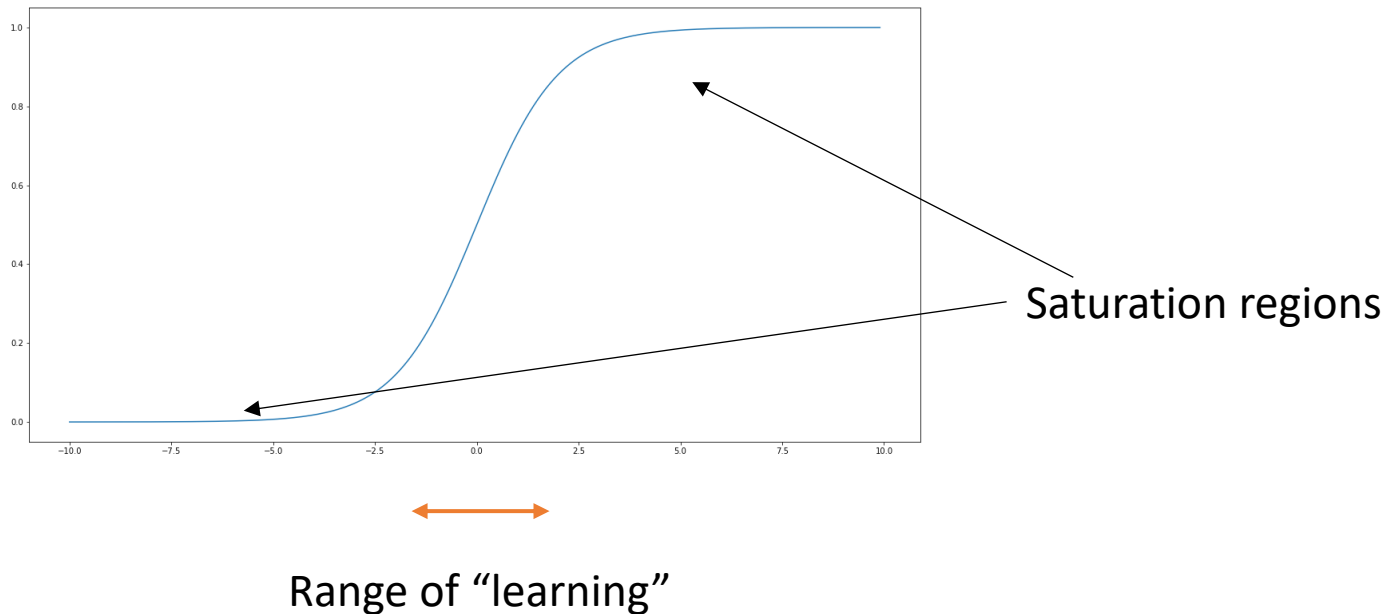
Vanishing Gradients

- Gradient Descent is a “local search” algorithm.



- Small gradient = changing θ by a little does not change the output much.

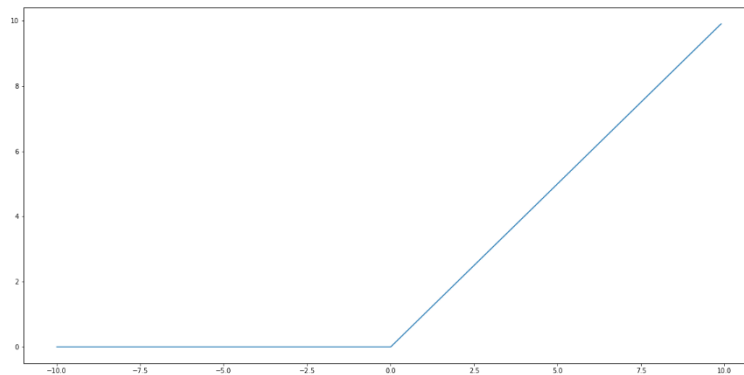
Vanishing Gradients and Sigmoid



- If the network parameters put the sigmoid into the saturated region, then the gradient will be very small.
- The deeper the network, the more likely one is to hit the saturated region.
- This makes it harder to train deep networks with sigmoid activation.

ReLU activation

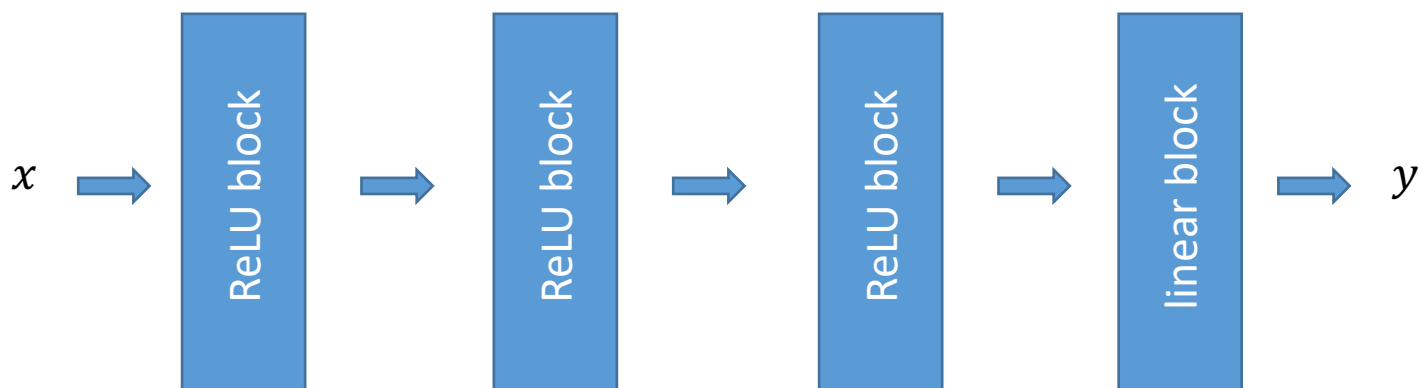
- Modern deep networks usually use some variant of the *rectified linear unit* activation:
$$\sigma(x) = \max(x, 0)$$



- This still shares some inspiration with sigmoid: it is “off” when the input is negative.
- Since it **does not saturate** on the positive end, it has a much **larger range of “learning”**.

ReLU activation MLP

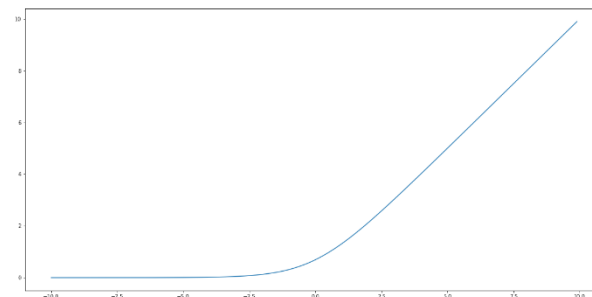
- A more modern MLP looks like:



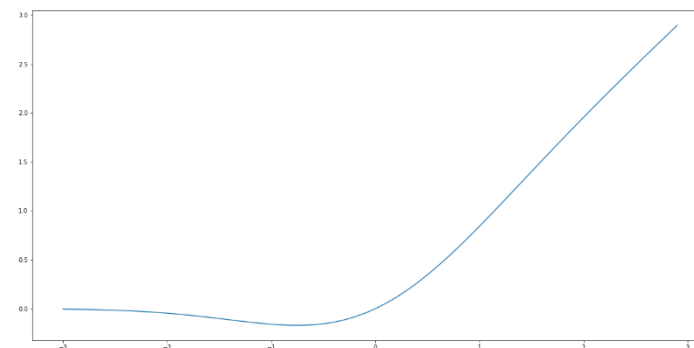
- ReLU networks are piece-wise linear.

Variations on ReLU

- Softplus: $\sigma(x) = \log(1 + \exp(x))$
 - “smooth ReLU”



- GeLU: $\sigma(x) = \frac{x \left(1 + \operatorname{erf}\left(\frac{x}{\sqrt{2}}\right)\right)}{2}$
 - Much more recent.





Loss Functions

Classification Losses

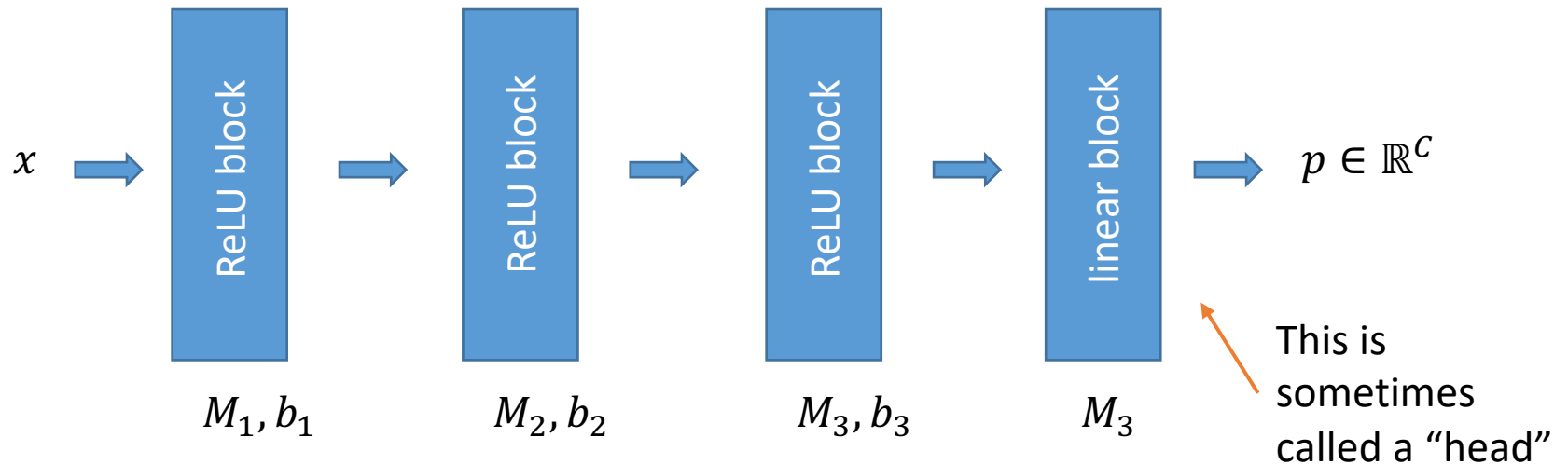


Cat?

Dog?

Horse?

Classification Architectures



- Prediction is the coordinate of p with the highest value:
 $prediction = \operatorname{argmax}_j p[j]$

True Classification Loss

- We want to optimize the true train error:

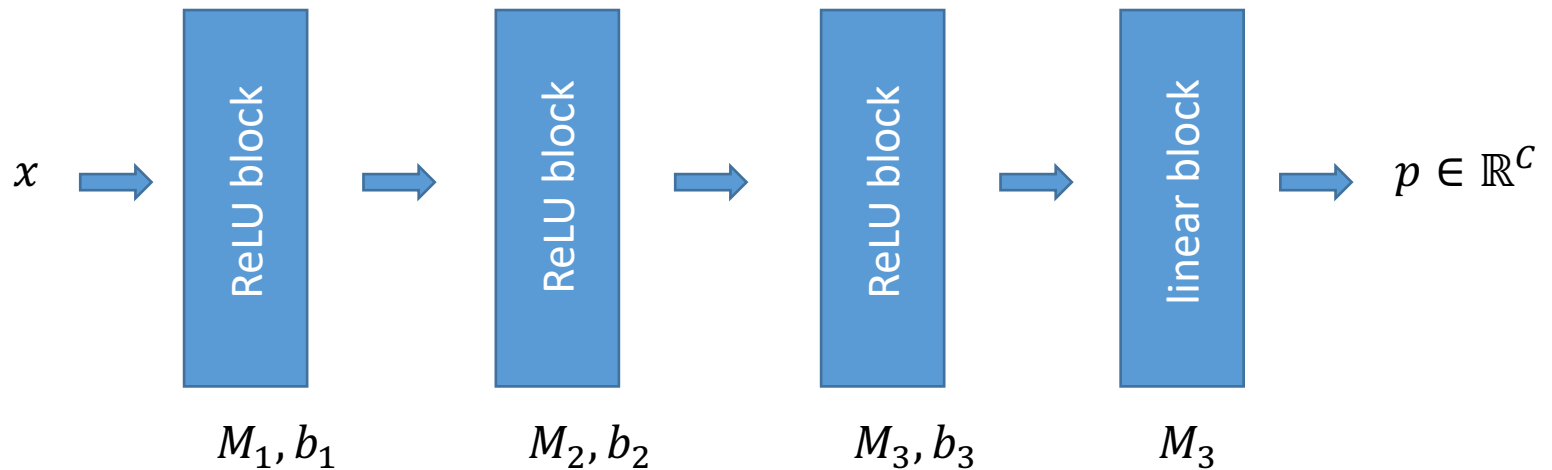
$$\sum_{dataset} 1[prediction \neq label]$$

- Let $S = \{(x_1, y_1), \dots, (x_N, y_N)\}$ where $y_i \in \mathbb{R}^C$ is a 1-hot vector:

$$y_i = (0, \dots, 1, \dots, 0)$$

$y_i[k] = 1$ only if x_i has label k .

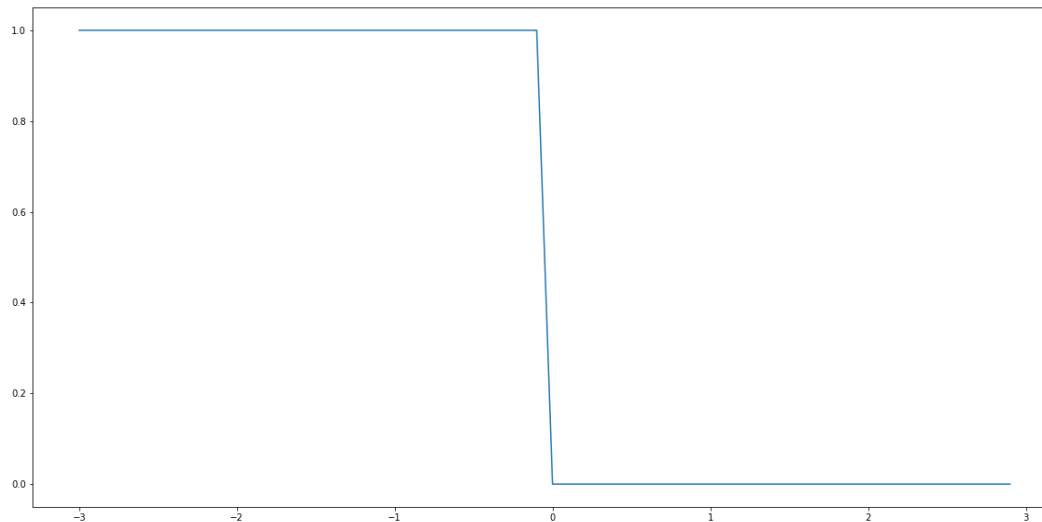
True Classification Loss



- Prediction is the coordinate of p with the highest value:
 $prediction = \operatorname{argmax}_j p[j]$
- Loss is $\ell(y, p) = 1 - y[prediction]$
- Why **can't we use this loss?**
 - **Non-differentiable**
 - Even if it were, there would be **vanishing gradients.**

True Classification Loss

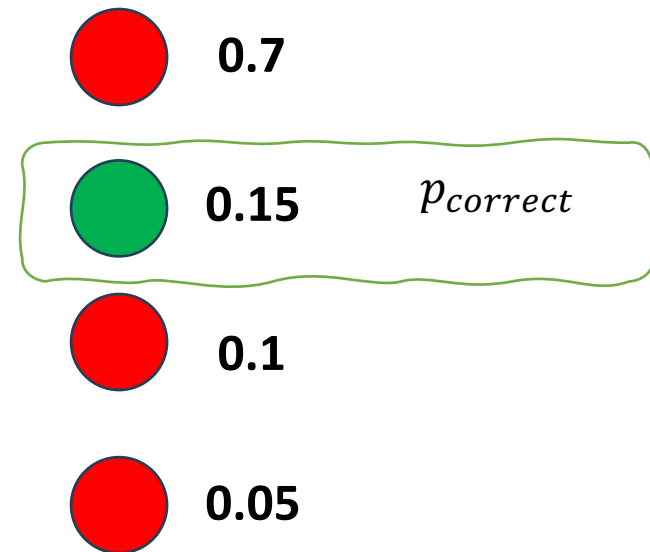
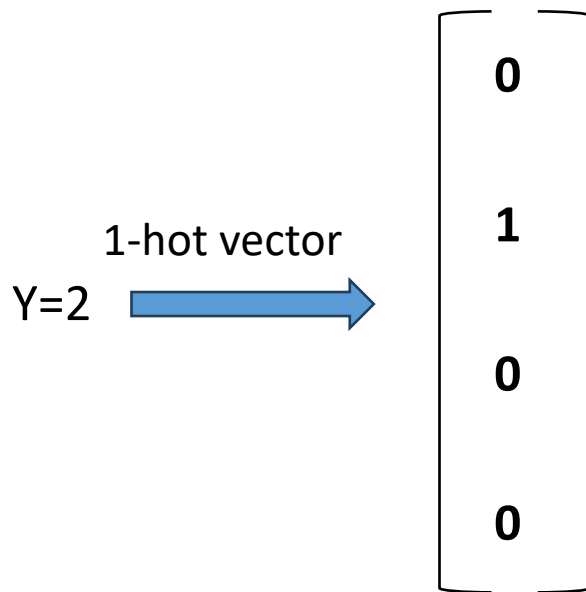
- It's best to visualize this loss for binary classification:



$$p_{correct} - p_{incorrect}$$

Hinge Loss

- Hinge is to true loss as ReLU is to sigmoid:
- $hinge(y, p) = \max[\max_{i \neq correct} p_i - p_{correct} + 1, 0]$

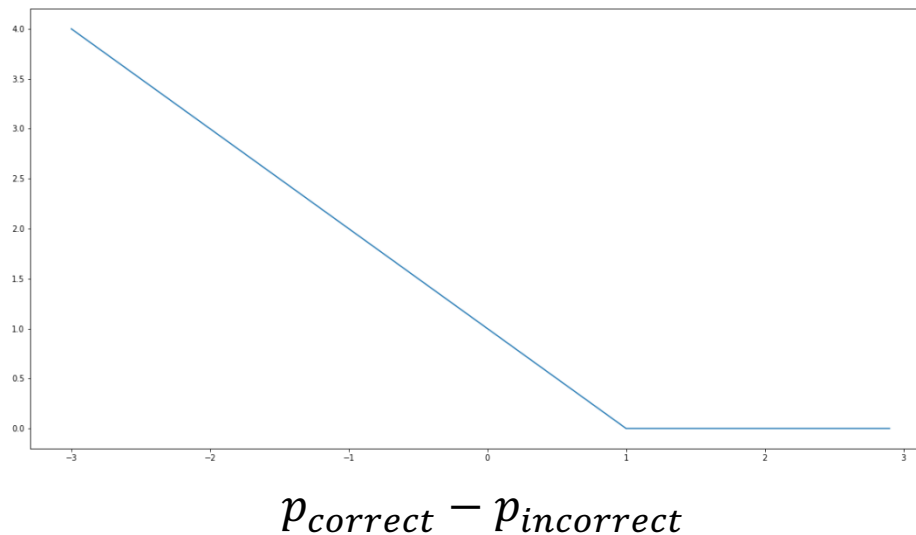


Output of the NN

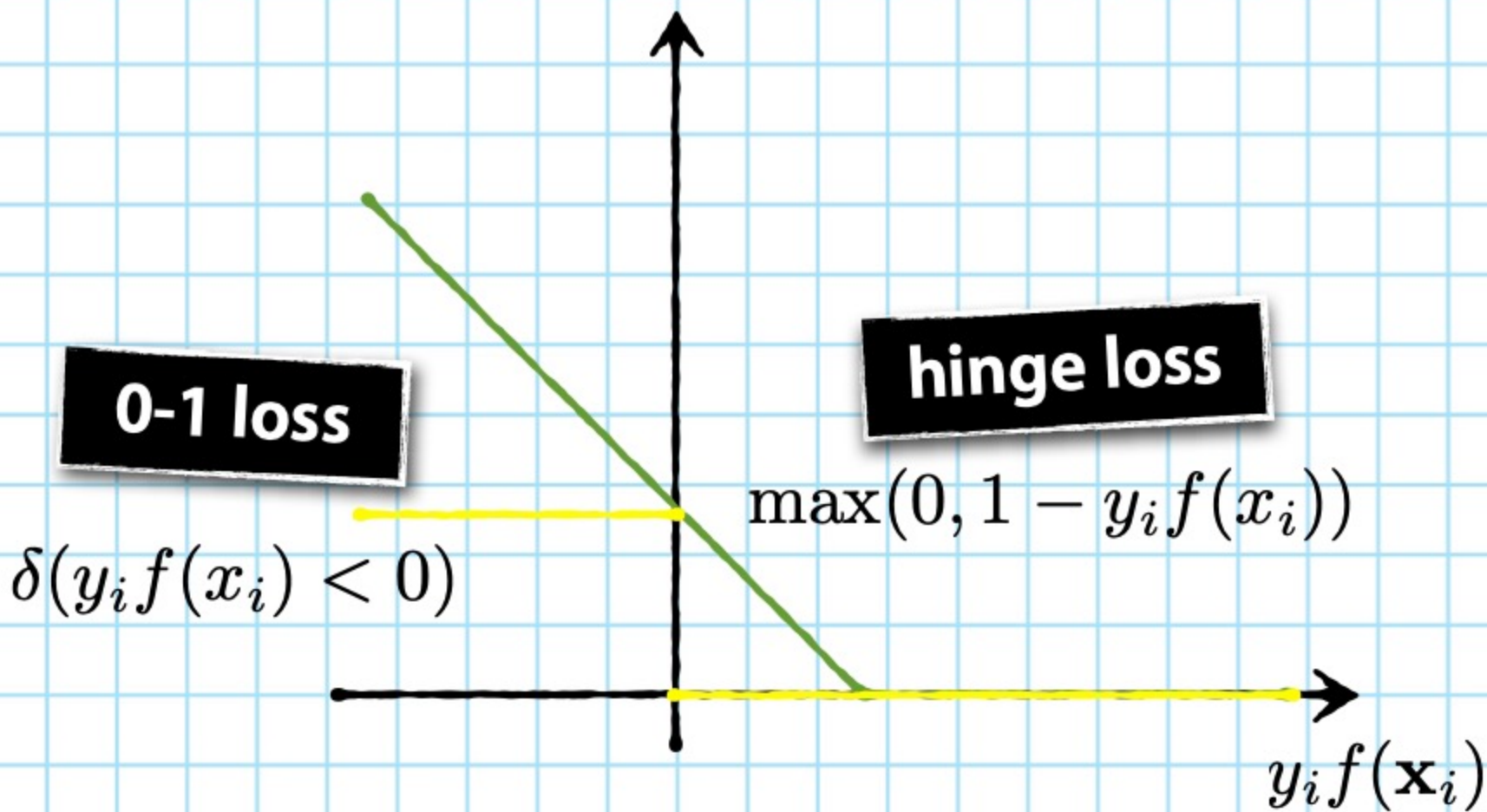
$$hinge(y, p) = 0.7 - 0.15 + 1$$

Hinge Loss

- Hinge is to true loss as ReLU is to sigmoid:
- $hinge(y, p) = \max[\max_{i \neq correct} p_i - p_{correct} + 1, 0]$
- The score for the “true” label needs to be 1 more than all the other scores to have zero loss.



Hinge Loss vs True Loss



Hinge Loss vs True Loss

- If the average hinge loss is 0.1, what can you say about the average number of mistakes?
- If the average hinge loss is 0.99, what can you say about the average number of mistakes?

Cross-Entropy Loss

- Input 1: a *distribution* over the C possible output classes

$$p = (p_1, \dots, p_C), \quad \sum_{i=1}^C p_i = 1$$

- Input 2: the true class value as a 1-hot vector:

$$y = (0, \dots, 1, \dots, 0)$$

- The cross-entropy loss is:

$$\begin{aligned} \ell(p, y) &= \sum_i -1[y = i] \log(p_i) \\ &= -\log(p_{\text{correct class}}) \end{aligned}$$

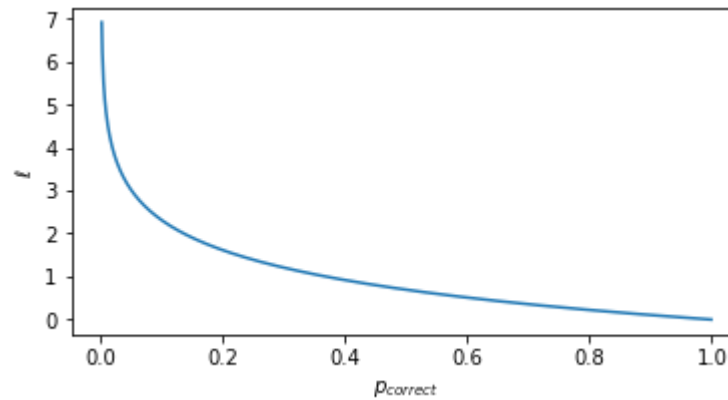
Cross-Entropy Loss

- If the correct class has probability $p_{correct} = 1$,

$$\ell(p, y) = -\log(p_{correct}) = 0$$

- If the correct class has probability $p_{correct} = 0$,

$$\ell(p, y) = -\log(0) = \infty$$



Cross-Entropy Loss

- The hinge loss is the smallest convex upper-bound of the classification loss.
- The cross-entropy is also convex, and positive.
- How is it related to the classification loss?
 - When do we not choose class k ?
 - When $k \neq \operatorname{argmax}_i(p_i)$.
 - If $k \neq \operatorname{argmax}_i(p_i)$, then $p_k \leq \frac{1}{2}$.
- If incorrect,

$$\ell(p, y) = -\log(p_k) \geq -\log\left(\frac{1}{2}\right) = \log(2)$$

Cross-Entropy and Classification Loss

- Since $\ell(p, y) \geq 0$, when the prediction is correct we have $\ell(p, y) \geq 1[\textit{prediction} \neq \textit{true label}]$.

- Otherwise, we have

$$\ell(p, y) = -\log(p_k) \geq -\log\left(\frac{1}{2}\right) = \log(2)$$

- So, in general

$$\ell(p, y) \geq \log(2)1[\textit{prediction} \neq \textit{true label}]$$

Softmax

- Cross-Entropy loss requires the input to be a *distribution*.
- The output of a linear layer is never normalized to 1.
- The *softmax layer* turns arbitrary vectors into probability distributions:

$$\text{softmax}(x)[i] = \frac{\exp x[i]}{\sum_j \exp x[j]}$$

$$\sum_i \text{softmax}(x)[i] = 1$$

UNSTABLE SOFTMAX

$$s_i = \frac{e^{a_i}}{\sum_{j=1}^N e^{a_j}}$$

```
import numpy as np

def softmax(x):
    exps = np.exp(x)
    return exps / np.sum(exps)

x = softmax([1, 2, 3])
y = softmax([1000, 2000, 3000])

print('x = ', x)
print('y = ', y)
```

```
import numpy as np
```

```
def softmax(x):  
    exps = np.exp(x)  
    return exps / np.sum(exps)
```

```
x = softmax([1, 2, 3])  
y = softmax([1000, 2000, 3000])
```

```
print('x = ', x)  
print('y = ', y)
```

```
import numpy as np
```

```
def softmax(x):  
    exps = np.exp(x)  
    return exps / np.sum(exps)
```

```
x = softmax([1, 2, 3])  
y = softmax([1000, 2000, 3000])
```

```
print('x = ', x)  
print('y = ', y)
```

```
import numpy as np
```

```
def softmax(x):  
    exps = np.exp(x)  
    return exps / np.sum(exps)
```

```
x = softmax([1, 2, 3])  
y = softmax([1000, 2000, 3000])
```

```
print('x = ', x)  
print('y = ', y)
```

```
import numpy as np

def softmax(x):
    exps = np.exp(x)
    return exps / np.sum(exps)

x = softmax([1, 2, 3])
y = softmax([1000, 2000, 3000])

print('x = ', x)
print('y = ', y)
```

```
import numpy as np

def softmax(x):
    exps = np.exp(x)
    return exps / np.sum(exps)

x = softmax([1, 2, 3])
y = softmax([1000, 2000, 3000])

print('x = ', x)
print('y = ', y)
```

Console

```
x = [0.09003057 0.24472847 0.66524096]
```



```
import numpy as np

def softmax(x):
    exps = np.exp(x)
    return exps / np.sum(exps)

x = softmax([1, 2, 3])
y = softmax([1000, 2000, 3000])

print('x = ', x)
print('y = ', y)
```

Console

```
x = [0.09003057 0.24472847 0.66524096]
```

```
y = [nan nan nan]
```

$$\begin{aligned}
s_i &= \frac{e^{a_i}}{\sum_{j=1}^N e^{a_j}} \\
&= \frac{C e^{a_i}}{\sum_{j=1}^N C e^{a_j}} \\
&= \frac{e^{a_i + \log C}}{\sum_{j=1}^N e^{a_j + \log C}}
\end{aligned}$$

Cross Entropy and KL-Divergence

- Entropy: For a distribution π :

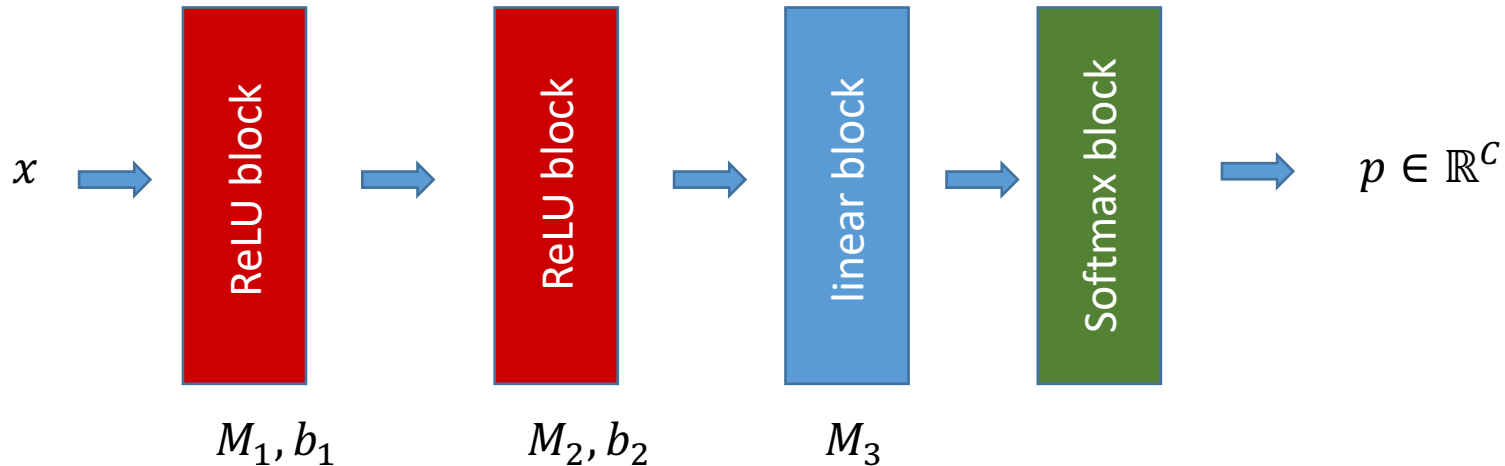
$$H(\pi) = - \sum_i \pi_i \log(\pi_i)$$

- KL-Divergence:

$$D(\pi||p) = \sum_i \pi_i \log\left(\frac{\pi_i}{p_i}\right)$$

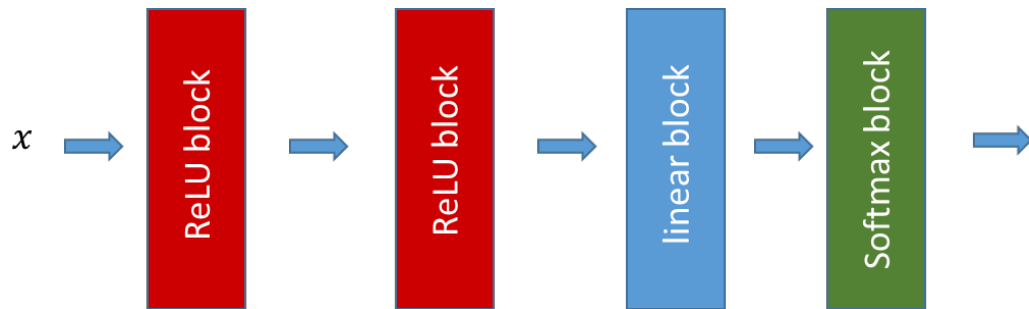
- $D(\pi||p) \geq 0$, and $D(\pi||p) = 0$ only if $\pi = p$.
- Think of $y = \pi$ as a distribution. Then
 $\text{crossentropy}(p, y) = H(y) + D(y||p)$

Summary



- Most common architecture: ReLU blocks, followed by a linear block with softmax (for classification)
- What about regression?
- Other tasks?

Deep Learning is modular

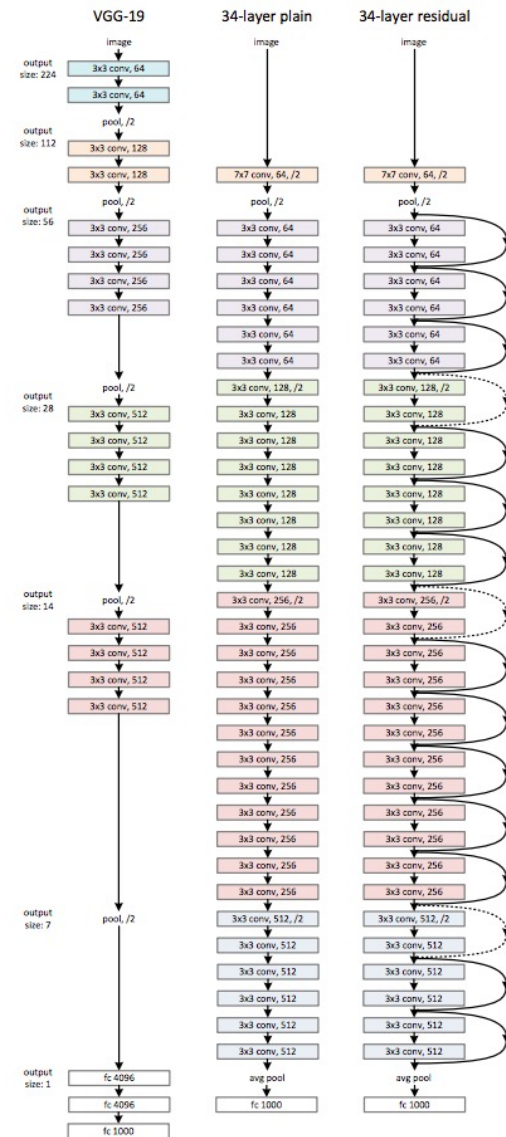
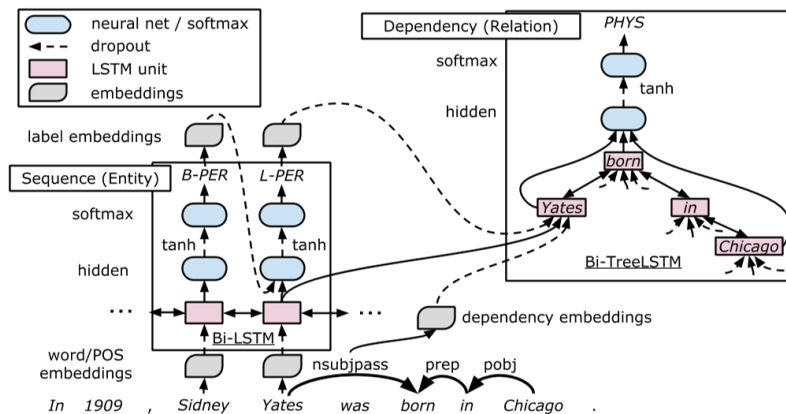




PyTorch

Why do we need deep learning libraries?

- Deep learning architectures can get very complicated!
- Instead of writing specialized code, use “auto-grad” library
- Easy to specify computation graph, find gradients



Popular DL libraries



PYTORCH



dmlc
mxnet



theano



Caffe

Some Pytorch

```
1  class TwoLayerNet(torch.nn.Module):
2      def __init__(self, input_size, layer_size, output_size):
3          super(TwoLayerNet, self).__init__()
4
5          self.Linear1 = torch.nn.Linear(input_size, layer_size)
6          self.activation1 = torch.nn.ReLU()
7
8          self.Linear2 = torch.nn.Linear(layer_size, output_size)
9          self.activation2 = torch.nn.ReLU()
10         self.softmax = torch.nn.Softmax(-1)
11
12
13     def forward(self, x):
14         out = self.Linear1(x)
15         out = self.activation1(out)
16
17         out = self.Linear2(out)
18         out = self.activation2(out)
19         out = self.softmax(out)
20
21     return out
```

Some Pytorch

```
1 x = torch.tensor([[1.0, 1.0],
2                   [1.0, 0.0],
3                   [-1.0, 2.0],
4                   [2.0, 4.0]])
5
6 y = torch.tensor([[0, 0, 1],
7                   [1, 0, 0],
8                   [1, 0, 0],
9                   [0, 1, 0]])
```

“batch size” = 4

- A “tensor” is an n-dimensional matrix.
- By convention, the first dimension of a tensor of datapoints is the “batch dimension”
- In this example, x is a batch of 4 datapoints, each of which is 2-dimensional.

Some Pytorch

```
1 x = torch.tensor([[1.0, 1.0],
2                   [1.0, 0.0],
3                   [-1.0, 2.0],
4                   [2.0, 4.0]])
5
6 y = torch.tensor([[0,0,1],
7                   [1,0,0],
8                   [1,0,0],
9                   [0,1,0]])
```

```
1 net = TwoLayerNet(2, 10, 3)
2 scores = net(x)
3 loss_func = torch.nn.CrossEntropyLoss()
4
5 label_indices = np.argmax(y,axis=1)
6
7 loss = loss_func(scores, label_indices)
8
9 print("scores: ", scores, "\n\nlabels: ", label_indices, "\n\nloss: ", loss)
```

```
scores: tensor([[0.3275, 0.4101, 0.2624],
               [0.3724, 0.3478, 0.2798],
               [0.2729, 0.4804, 0.2467],
               [0.2259, 0.5482, 0.2259]]), grad_fn=<SoftmaxBackward>)
```

```
labels: tensor([2, 0, 0, 1])
```

```
loss: tensor(1.0730, grad_fn=<NllLossBackward>)
```

Some Pytorch

```
1 class MultiLayerNet(torch.nn.Module):
2     def __init__(self, input_size, layer_sizes, output_size):
3         super(MultiLayerNet, self).__init__()
4
5         layer_sizes = [input_size] + layer_sizes
6         self.linear_layers = [torch.nn.Linear(layer_sizes[i], layer_sizes[i+1]) for i in range(len(layer_sizes)-1)]
7         self.activations = [torch.nn.ReLU() for _ in self.linear_layers]
8
9         self.final_linear = torch.nn.Linear(layer_sizes[-1], output_size)
10
11         self.softmax = torch.nn.Softmax(-1)
12
13
14     def forward(self, x):
15         out = x
16         for linear_layer, activation in zip(self.linear_layers, self.activations):
17             out = linear_layer(out)
18             out = activation(out)
19         out = self.final_linear(out)
20
21         out = self.softmax(out)
22
23         return out
```

Next Time

- Backpropagation algorithm
 - How to build an automatic differentiation package.
 - What are these weird “graph” and “leaf” things my tensorflow/pytorch code keeps complaining about?