



Deep Learning: Introduction

Adam Moss
School of Physics and Astronomy
adam.moss@nottingham.ac.uk

Early Learning

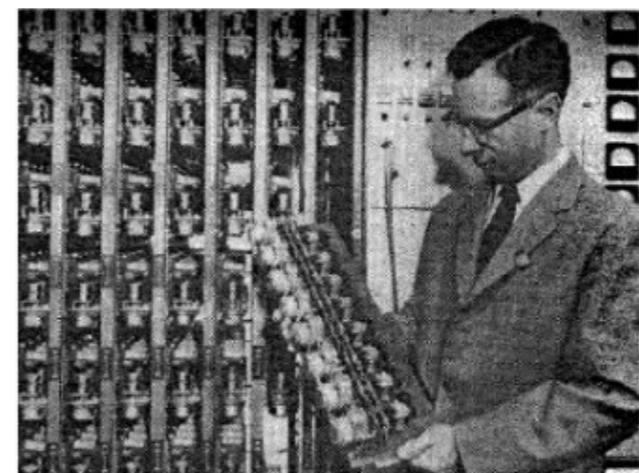
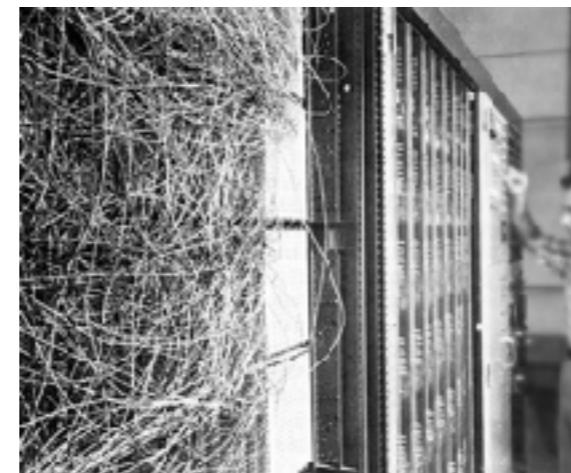
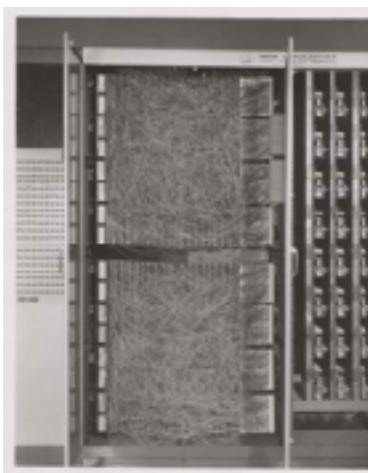
- ▶ Feature extraction



Hand crafted
feature
extraction

Trainable
classifier

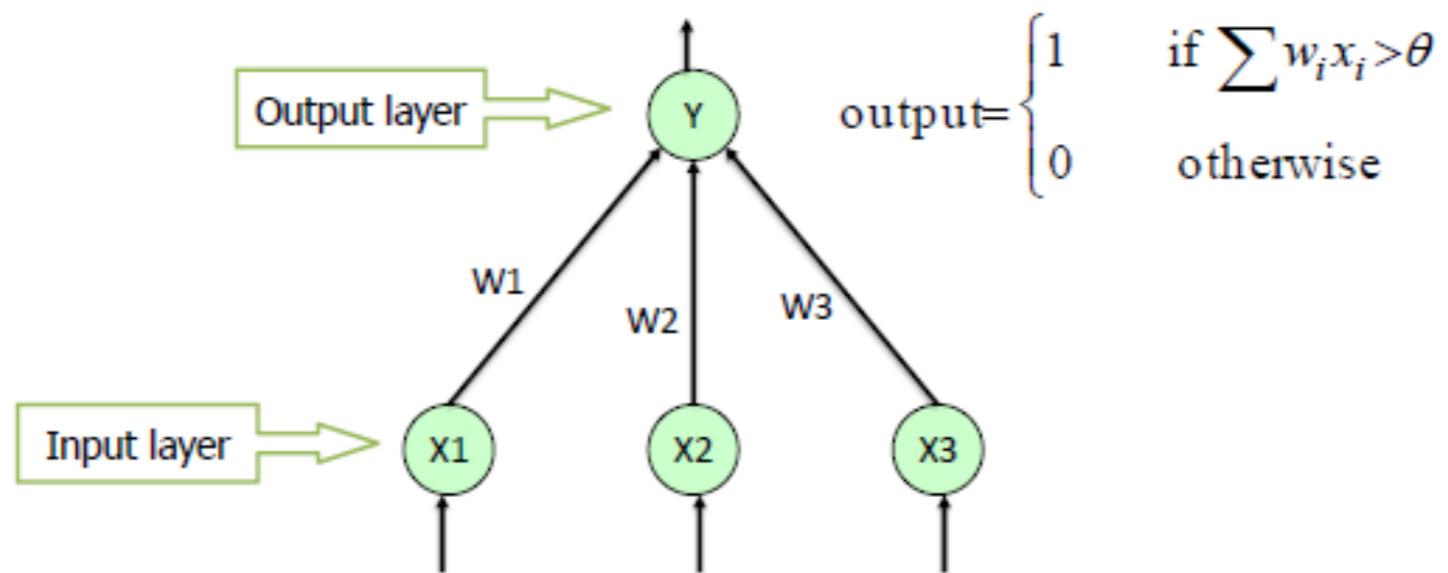
- ▶ Perceptron (binary classifier). Mark 1 perceptron machine (1957) used motors, potentiometers!



Perceptron

- ▶ Simplest perceptron: set of inputs x_i mapped to output \hat{y}

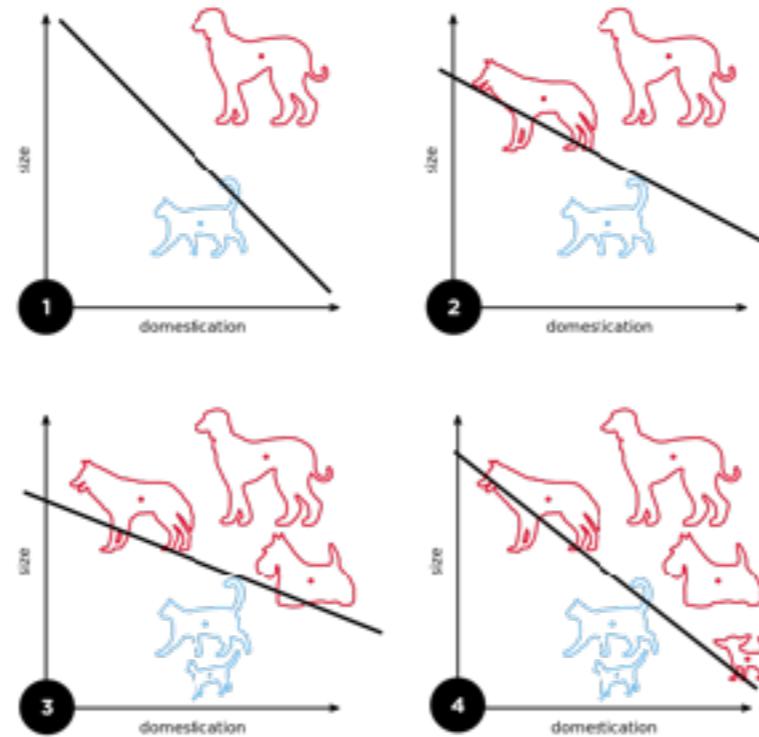
Single Layer Perceptron



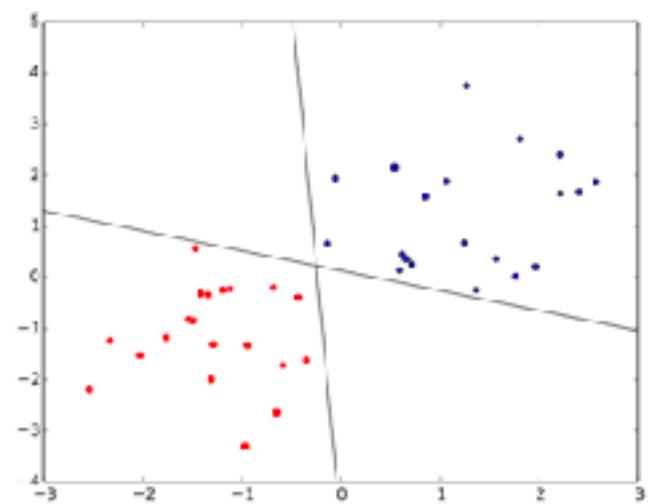
- ▶ Each input has a weight w_i
- ▶ Weights are trained using **supervised learning**
- ▶ Training sets of $D = \{x_{i,j}, y_j\}$ where j is the sample number and y_j the desired output for that sample
- ▶ Prediction is \hat{y} and weights are updated to minimise loss $\sum_j (y_j - \hat{y}_j)^2$

Perceptron

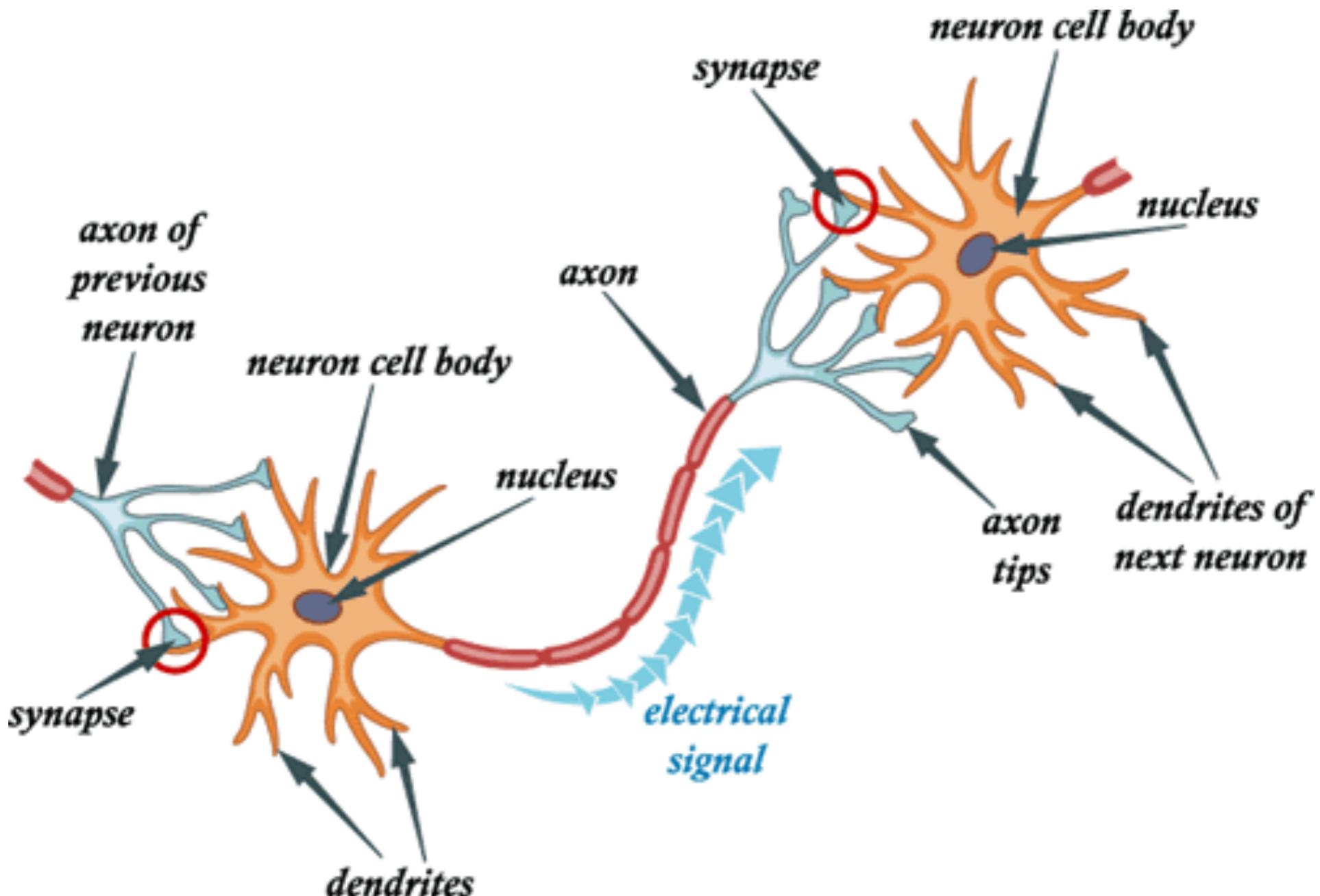
- ▶ The perceptron is only able to classify linearly separable training sets
- ▶ E.g. if two features are size and domestication



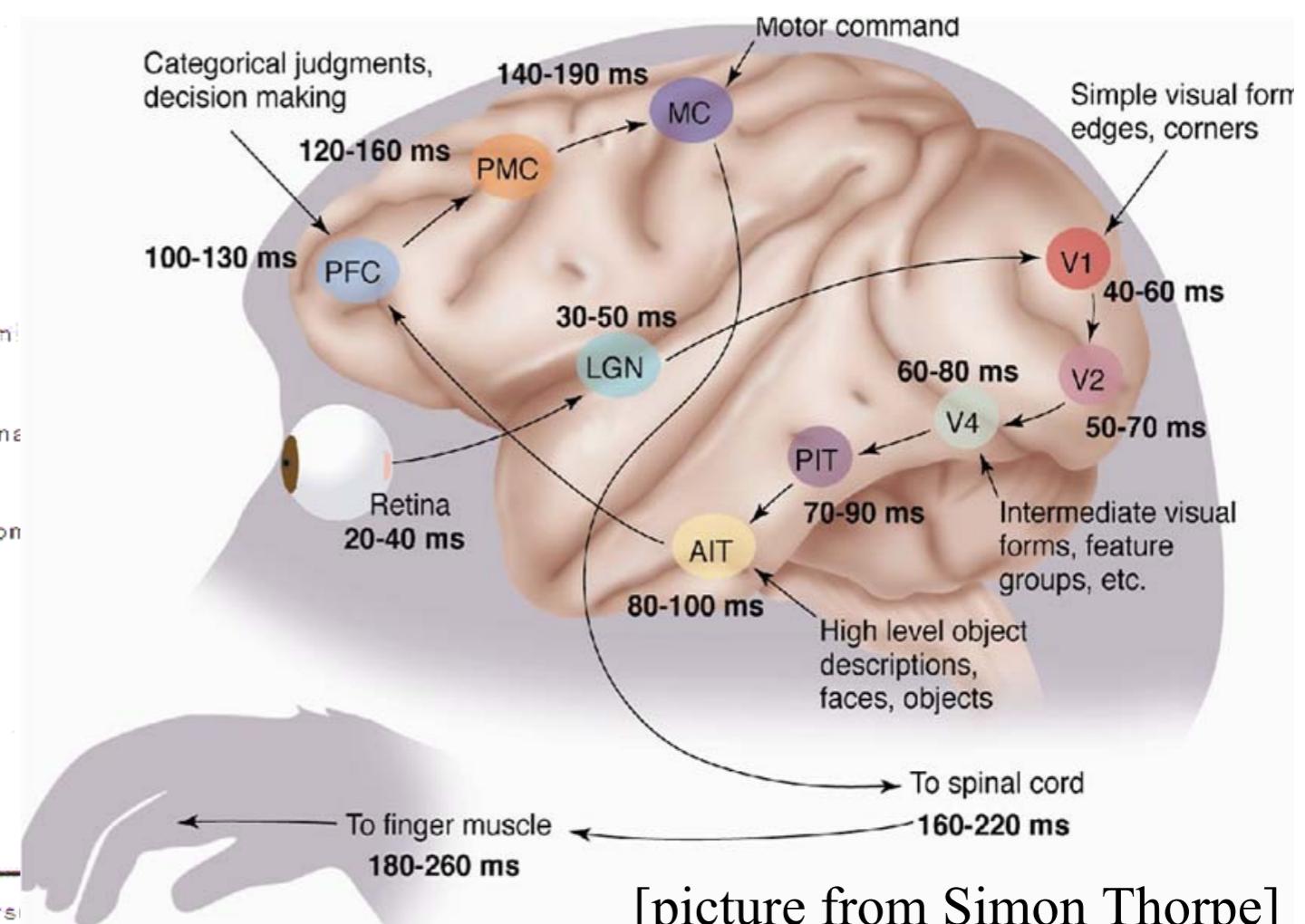
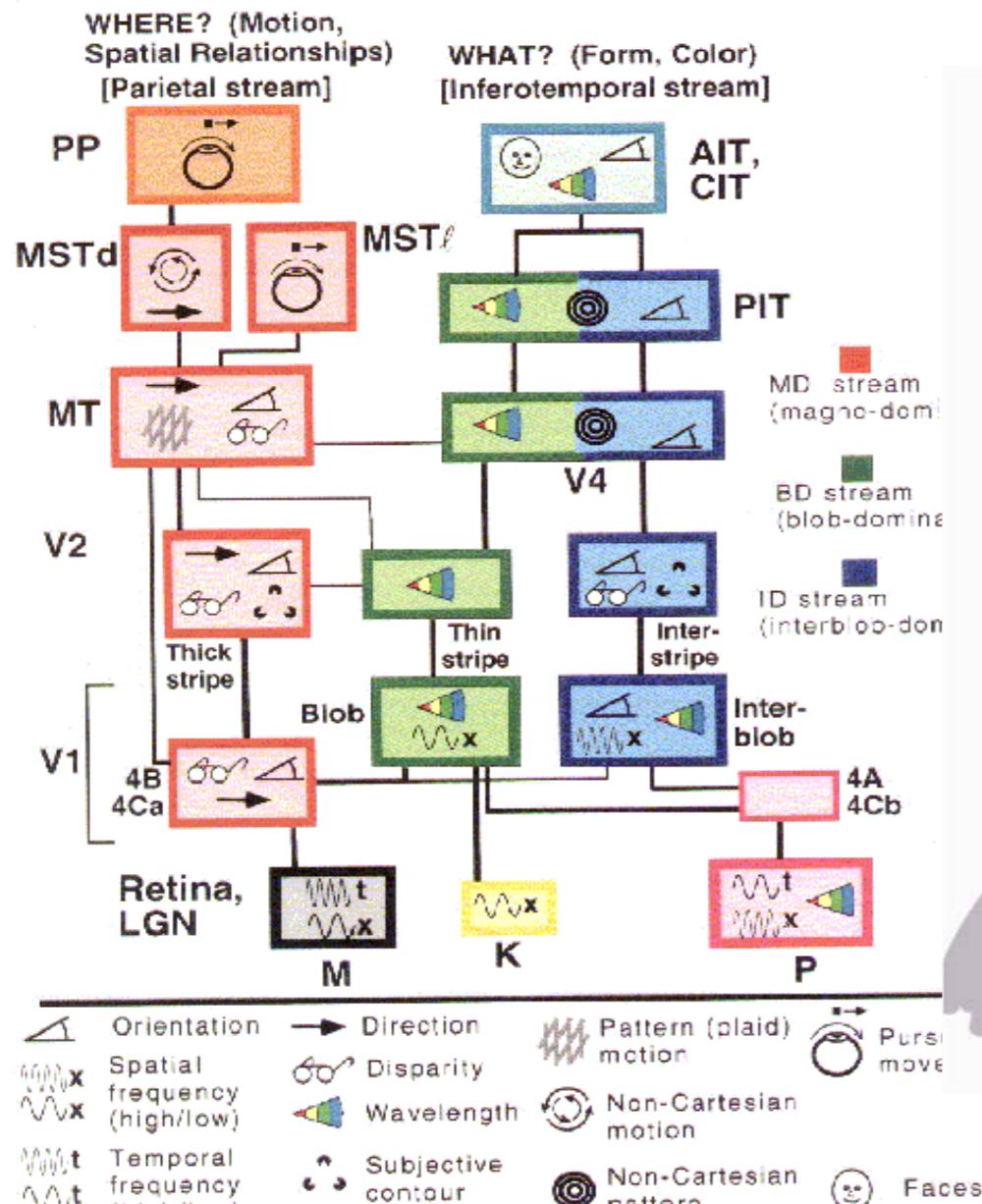
- ▶ May admit solutions of different quality (general problem for machine learning if training data is not representative)



Biological Neuron



Visual Cortex



[Gallant & Van Essen]

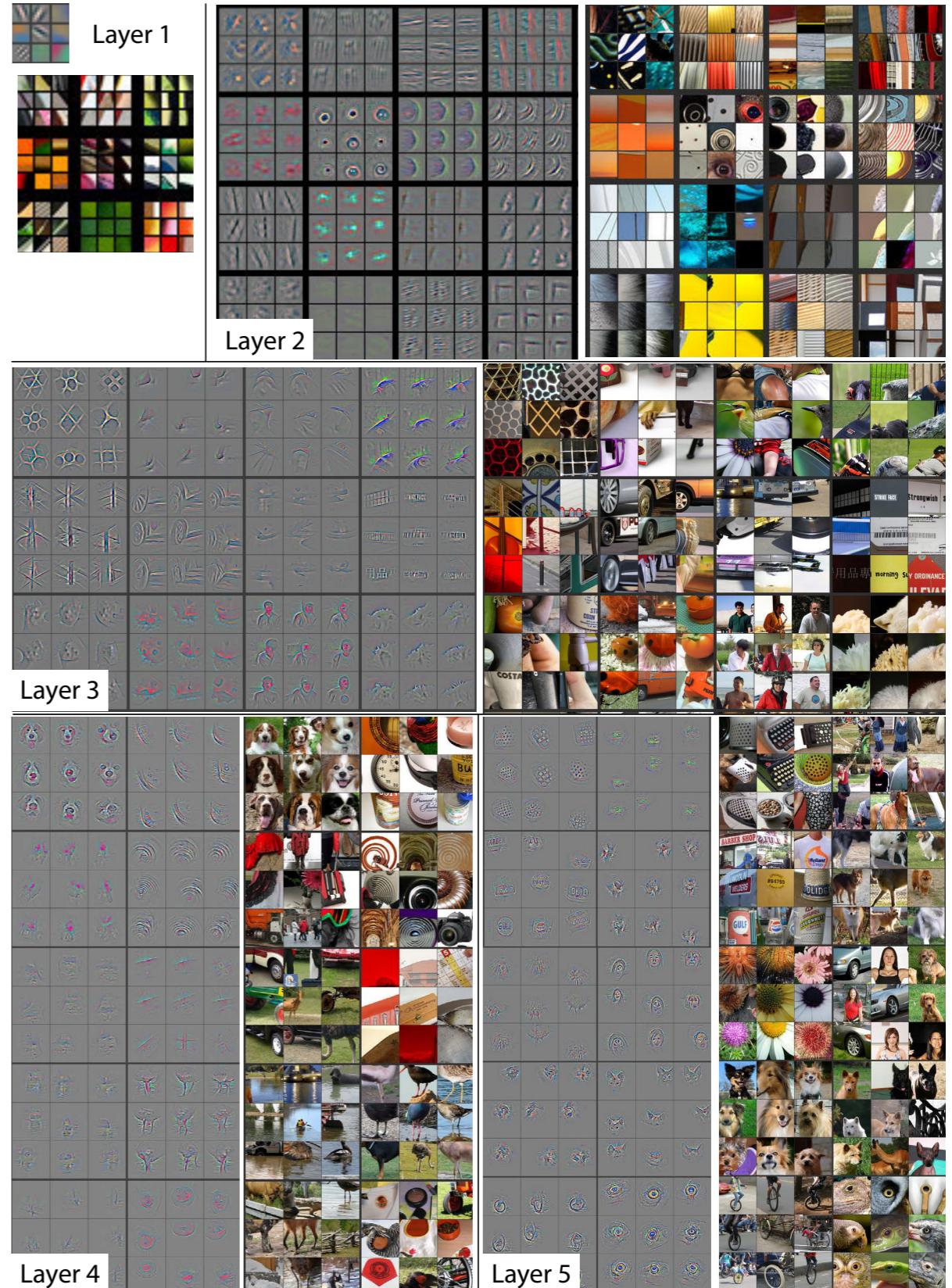
AIT=anterior inferior temporal cortex

PIT=posterior inferior temporal cortex

LGN=lateral geniculate nucleus

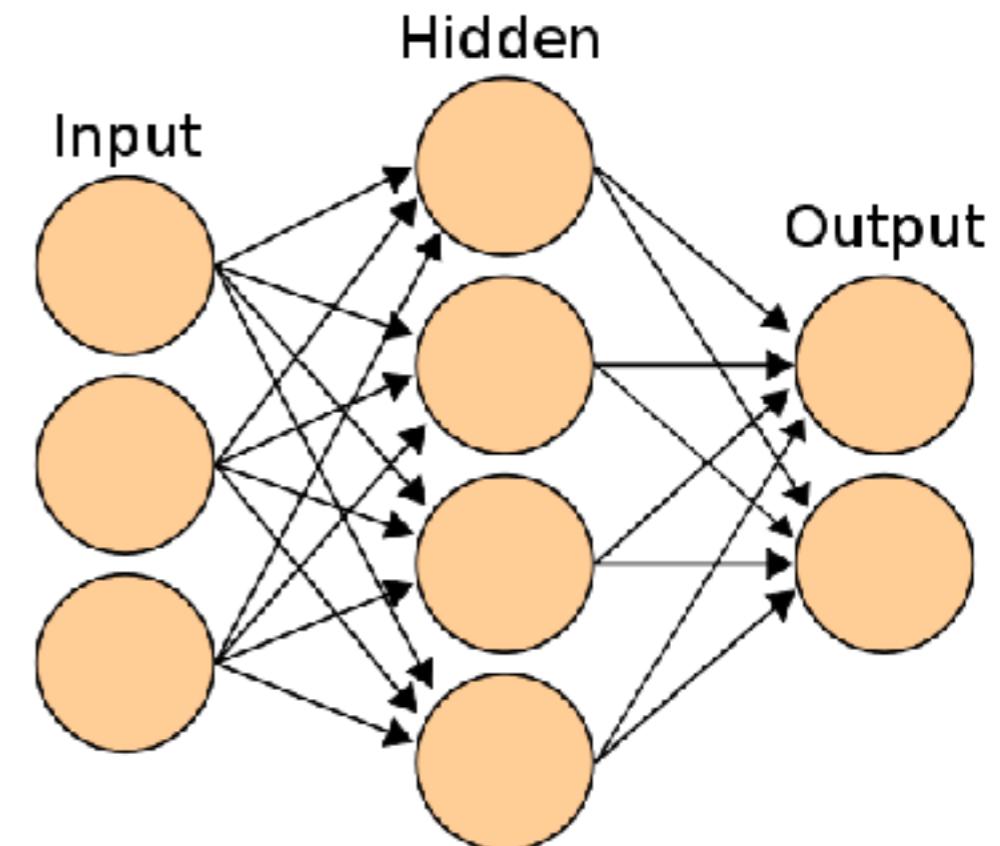
Shallow vs Deep

- ▶ A network is **deep** if it has more than one layer of non-linear feature abstraction
- ▶ Hierarchy of representations with increasing levels of abstraction (e.g. pixel -> edge -> eye -> face)
- ▶ Deep networks can store more memory than equivalent number of units in a single layer



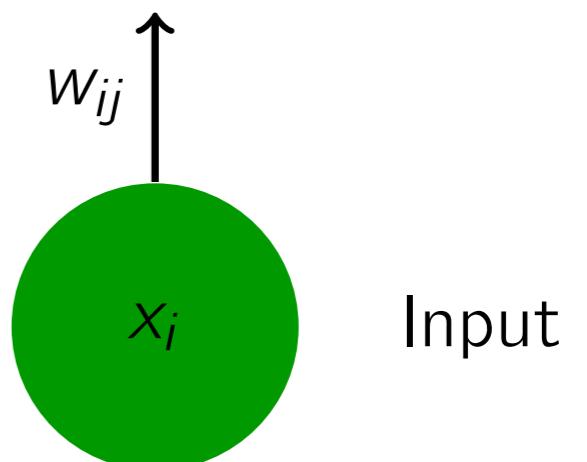
MultiLayer Perceptron

- ▶ Type of **feedforward** artificial neural network
- ▶ Can distinguish data which is not linearly separable
- ▶ Some neurons use non-linear **activation** functions (functions which map the weighted input to the output of a neuron)
- ▶ The brain is thought to work in a similar way when biological neurons are fired
- ▶ MLPs use supervised learning to update network weights using **back-propagation**



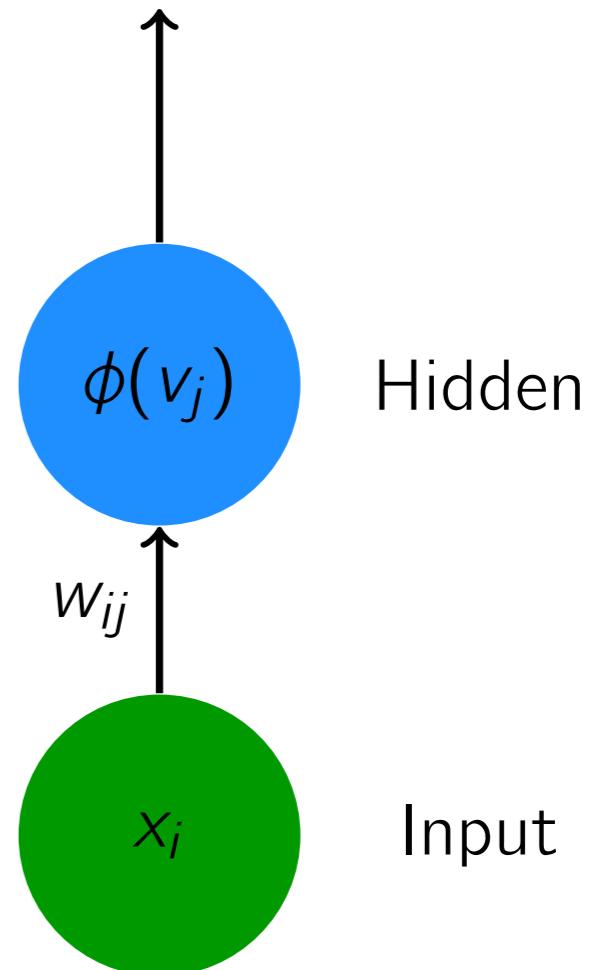
Back-propagation

- Weighted inputs $v_j = w_{ij}x_i$



Back-propagation

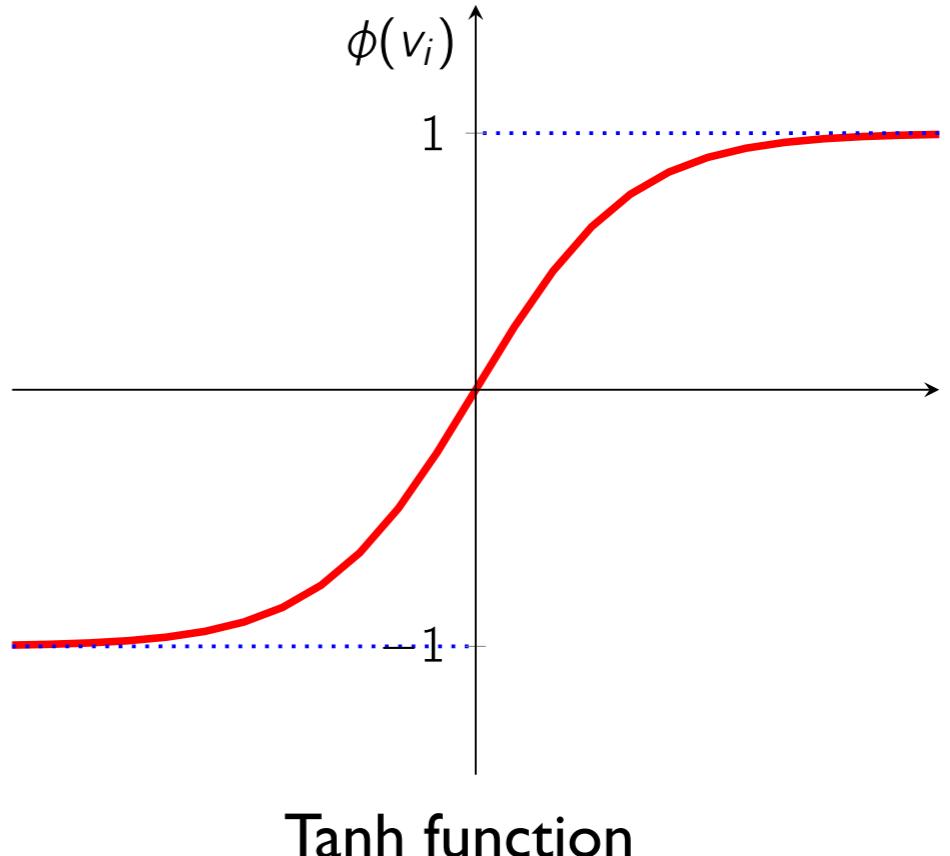
- ▶ Weighted inputs $v_j = w_{ij}x_i$
- ▶ Activation function $\phi(v_j)$



Back-propagation

- ▶ Weighted inputs $v_j = w_{ij}x_i$
- ▶ Activation function $\phi(v_j)$

$$\phi(v_i) = \tanh(v_i)$$

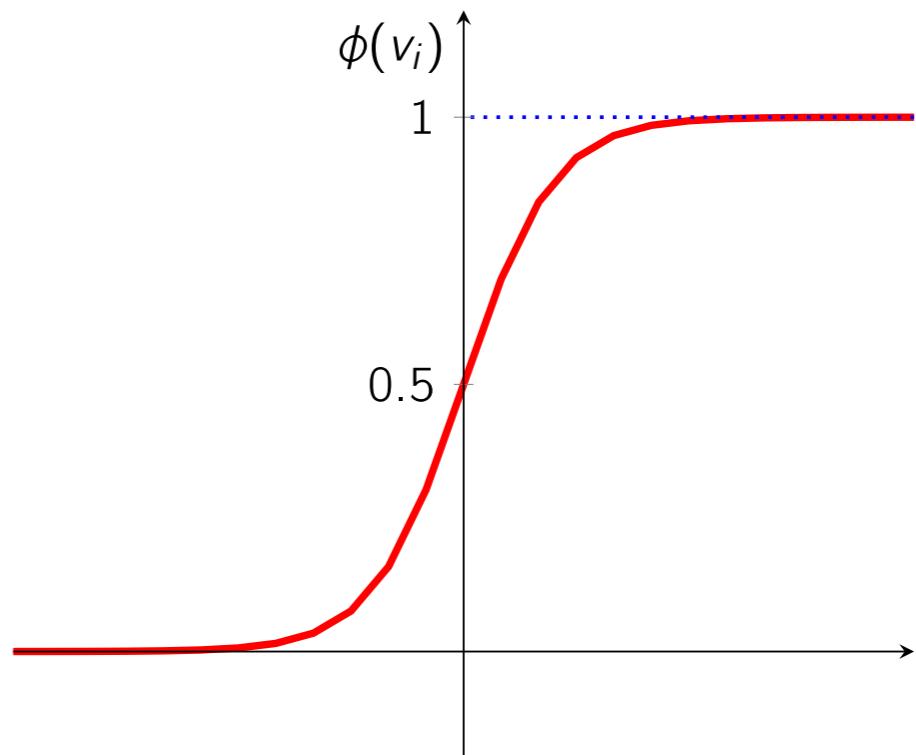


Back-propagation

- ▶ Weighted inputs $v_j = w_{ij}x_i$
- ▶ Activation function $\phi(v_j)$

$$\phi(v_i) = \tanh(v_i)$$

$$\phi(v_i) = (1 + e^{-v_i})^{-1}$$



Sigmoid function

Back-propagation

- ▶ Weighted inputs $v_j = w_{ij}x_i$
- ▶ Activation function $\phi(v_j)$

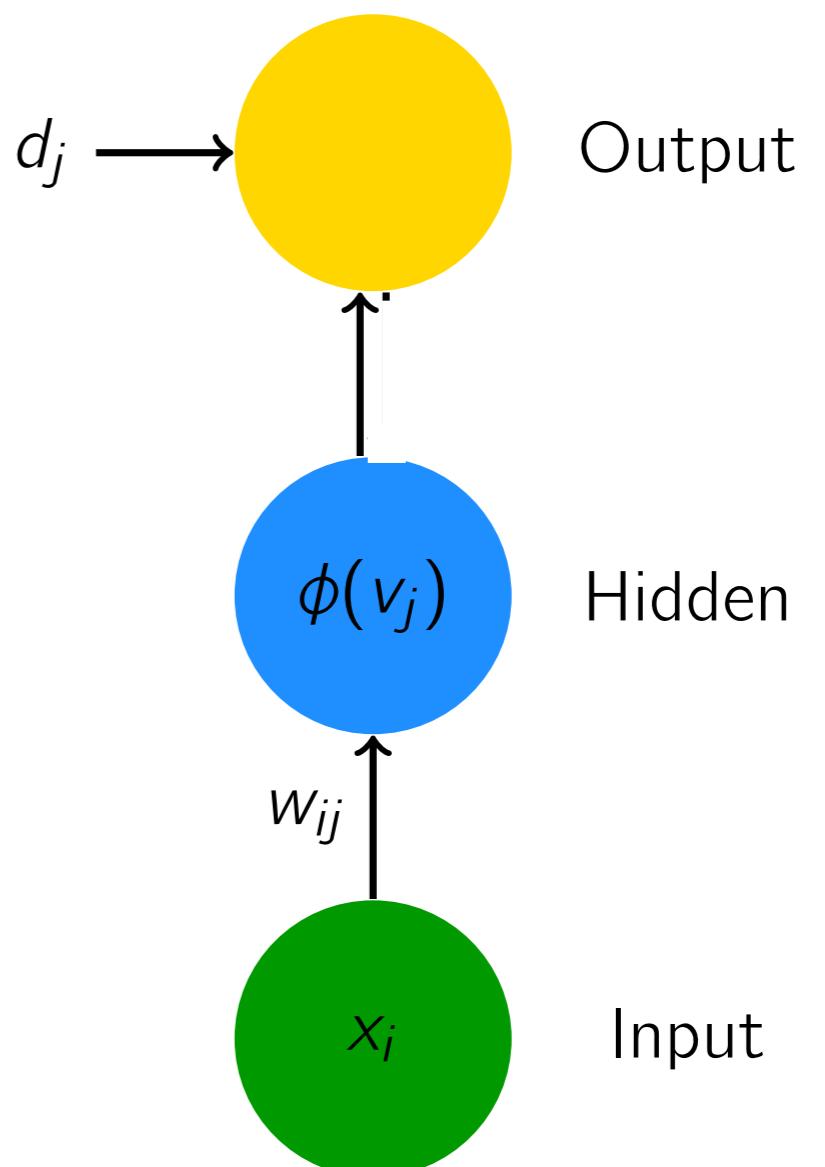
$$\phi(v_i) = \tanh(v_i)$$

$$\phi(v_i) = (1 + e^{-v_i})^{-1}$$

- ▶ Loss function

$$e_j = y_j - \phi(v_j)$$

$$\mathcal{E} = \sum_j e_j^2$$



Back-propagation

- ▶ Weighted inputs $v_j = w_{ij}x_i$
- ▶ Activation function $\phi(v_j)$

$$\phi(v_i) = \tanh(v_i)$$

$$\phi(v_i) = (1 + e^{-v_i})^{-1}$$

- ▶ Loss function

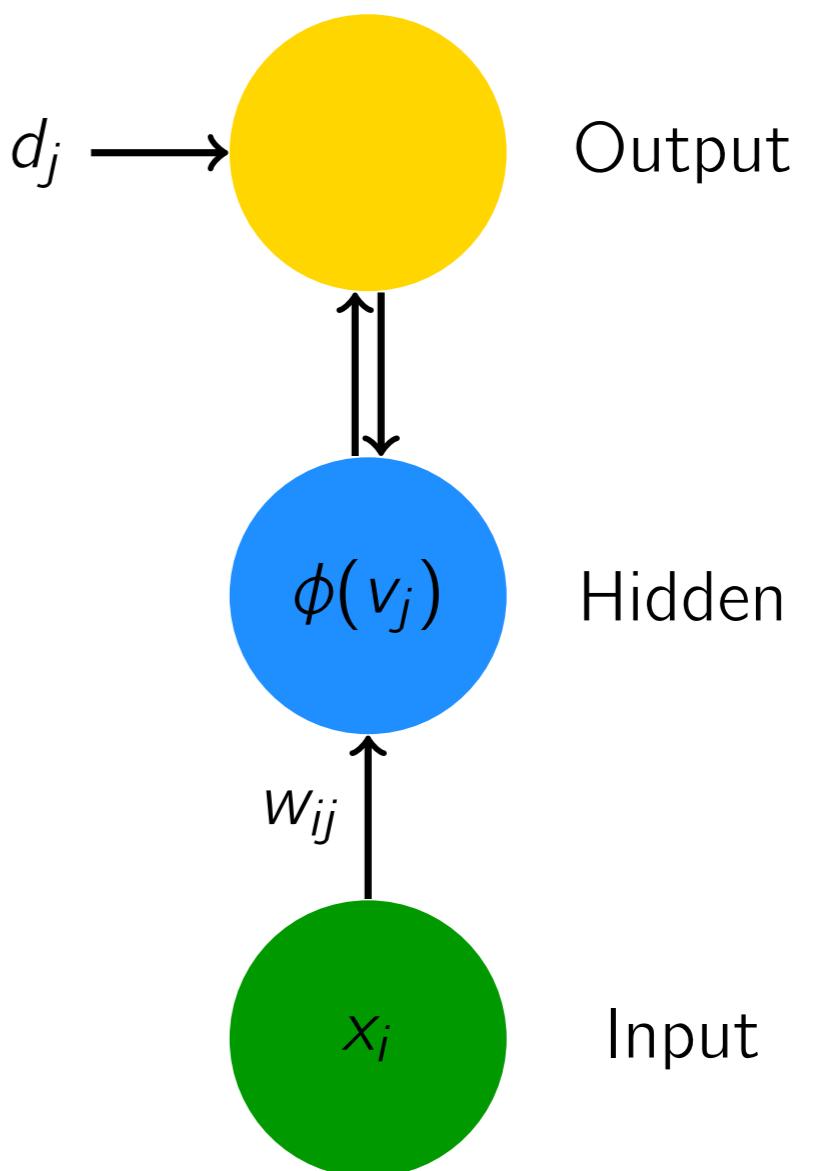
$$e_j = y_j - \phi(v_j)$$

$$\mathcal{E} = \sum_j e_j^2$$

- ▶ Back-propagation

$$\Delta w_{ij} = -\eta \frac{\partial \mathcal{E}}{\partial w_{ij}}$$

$$\Delta w_{ij} = -\eta \frac{\partial \mathcal{E}}{\partial \phi(v_j)} \frac{\partial \phi(v_j)}{\partial v_j} \frac{\partial v_j}{\partial w_{ij}}$$



Back-propagation

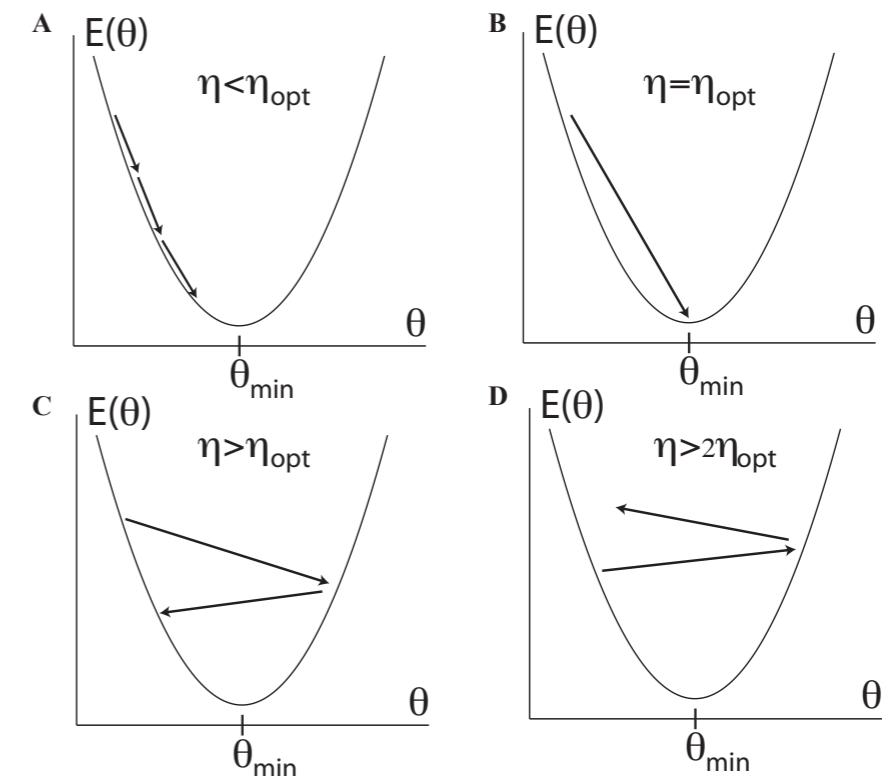
- ▶ Can the cortex do back-propagation?
- ▶ Maybe or maybe not!
- ▶ Many neuroscientists think the brain **can't** back-propagate
 - Source of supervision signal?
 - Neurons send all-or-nothing spikes
 - Neurons must be able to send different signals forward and backward
- ▶ But see work of Geoffrey Hinton (e.g. Lillicrap et. al., Nature Communications **7**, 2016) who presents arguments how the brain can back-propagate
- ▶ How much should we be led by the function of the brain when developing deep learning algorithms?

Gradient Descent (GD)

- Most supervised machine learning algorithms involve finding a set of parameters Θ that minimise a cost/loss function
- Loss functions can be complicated non-convex functions with many local minima
- Simplest algorithm that attempts to minimise the loss $E(\Theta)$ is gradient descent

$$\mathbf{v}_t = \eta_t \nabla_{\theta} E(\boldsymbol{\theta}_t), \quad \eta_t = \text{learning rate}$$
$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \mathbf{v}_t$$

- Requires careful choice of learning rate
- Can become stuck in local minima
- Sensitive to initial conditions
- Learning rate is same for all gradients
- Long time to escape saddle points



Stochastic Gradient Descent (SGD)

- ▶ Stochasticity is added by approximating the gradient on a subset of data called a **mini batch**
- ▶ Size of mini batch is smaller than dataset - typically 10's to 100's of data points
- ▶ Full iteration using all mini batches of a dataset is called an **epoch**
- ▶ Update rule is now just

$$\mathbf{v}_t = \eta_t \nabla_{\theta} E^{MB}(\theta),$$
$$\theta_{t+1} = \theta_t - \mathbf{v}_t.$$

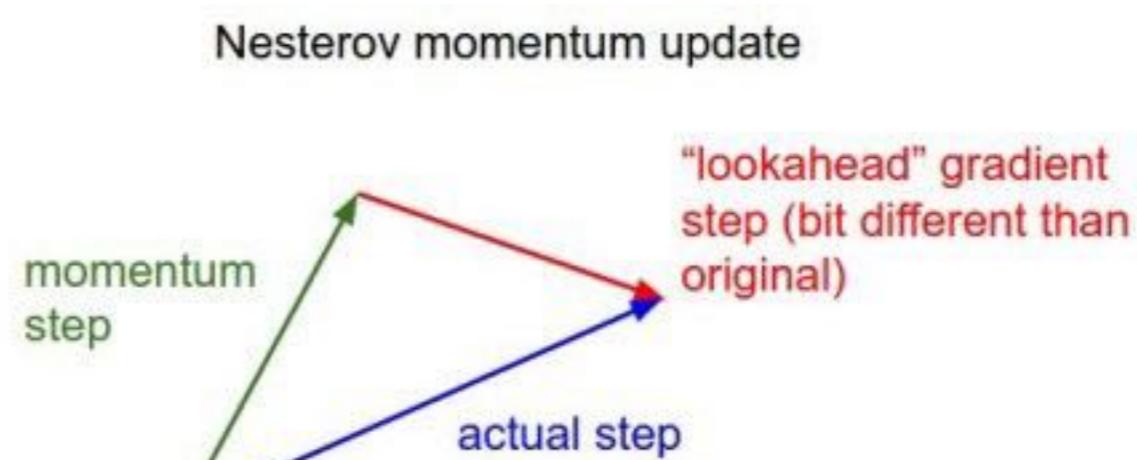
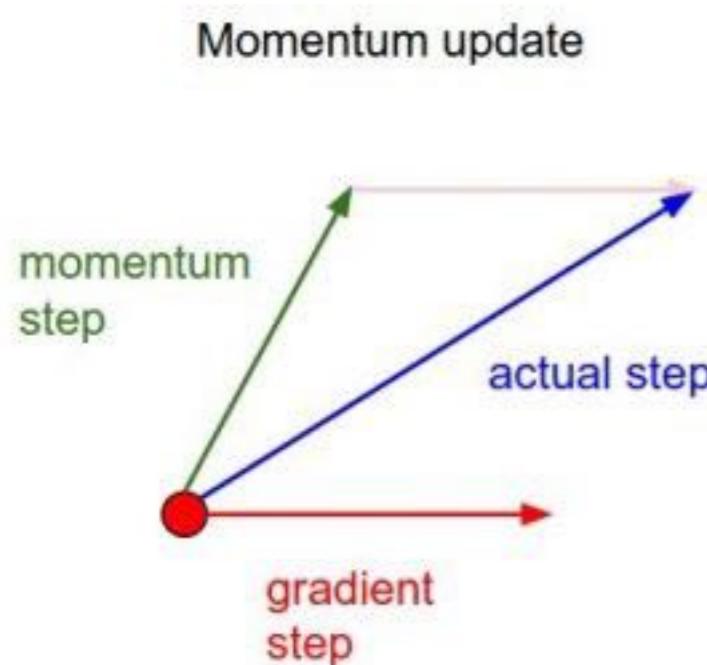
- ▶ Decreases chance of becoming stuck in local minima
- ▶ Also has been shown to help alleviate **over-fitting**

SGD with momentum

- ▶ SGD with momentum adds an inertia term that retains some memory of the direction being moved in

$$\mathbf{v}_t = \gamma \mathbf{v}_{t-1} + \eta_t \nabla_{\theta} E(\theta_t) \quad \gamma = \text{momentum parameter}$$
$$\theta_{t+1} = \theta_t - \mathbf{v}_t,$$

- ▶ SGD helps parameter updates gain ‘speed’ in persistent smaller gradients while suppressing oscillatory high gradients



Second order moments

- ▶ All methods so far require a schedule for the learning rate as a function of time
- ▶ Optimal learning rate is actually inverse of the Hessian $\eta_{\text{opt}} = [\partial_\theta^2 E(\theta)]^{-1}$.
- ▶ Expensive to compute. Second order moment methods keep track of the **squared gradient**
- ▶ Take large steps in shallow directions and small steps in steep directions
- ▶ Algorithms include RMSprop (Teileman and Hinton 2012), AdaDelta (Zeiler, 2012) and ADAM (Kingma and Ba, 2014)
- ▶ E.g. RMSprop update rule

$$\mathbf{g}_t = \nabla_\theta E(\theta)$$

$$\mathbf{s}_t = \beta \mathbf{s}_{t-1} + (1 - \beta) \mathbf{g}_t^2$$

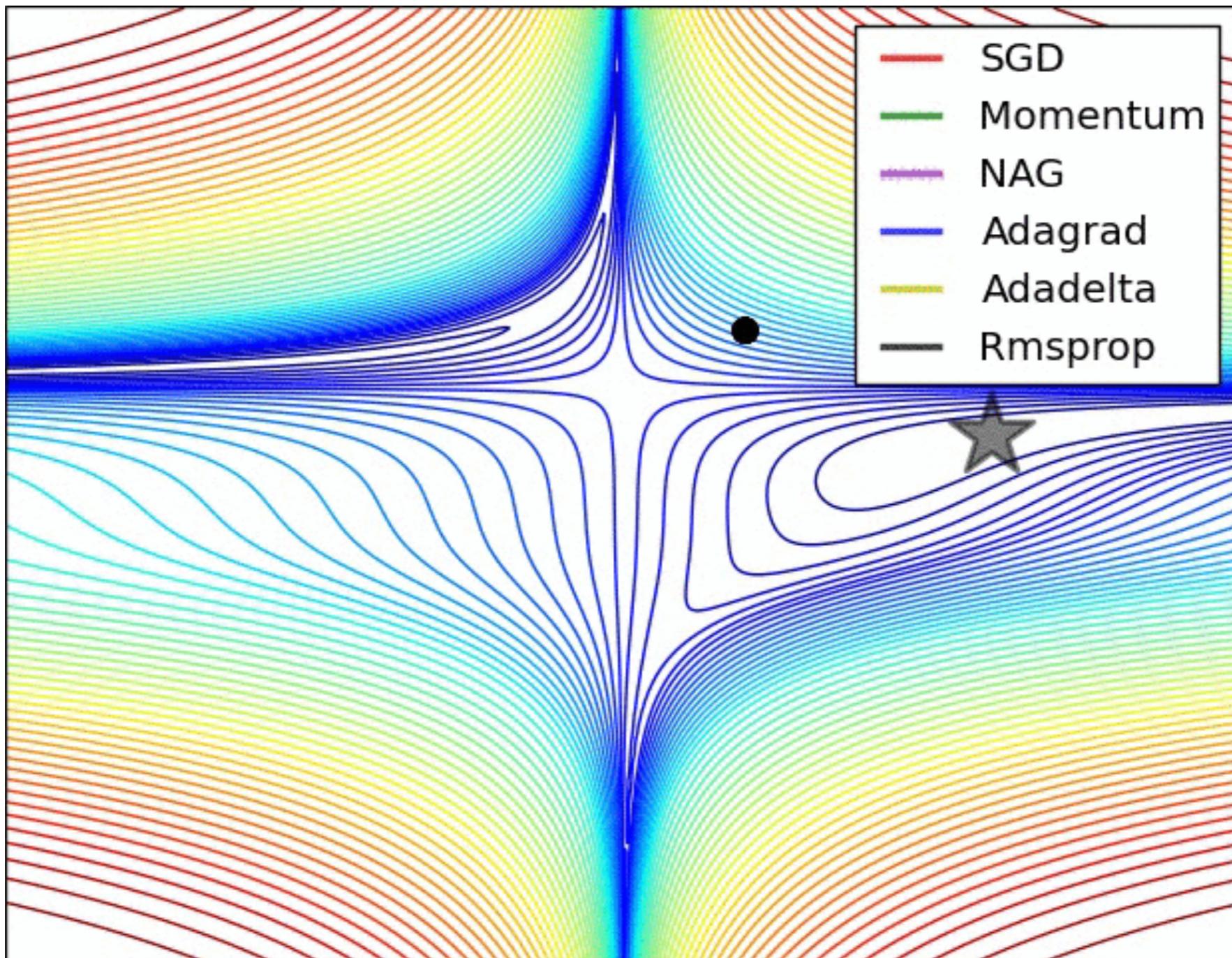
$$\theta_{t+1} = \theta_t - \eta_t \frac{\mathbf{g}_t}{\sqrt{\mathbf{s}_t + \epsilon}},$$

$$\mathbf{s}_t = \mathbb{E}[\mathbf{g}_t^2]$$

β = averaging time

ϵ = regularisation constant

Second order methods



Deep learning in practice

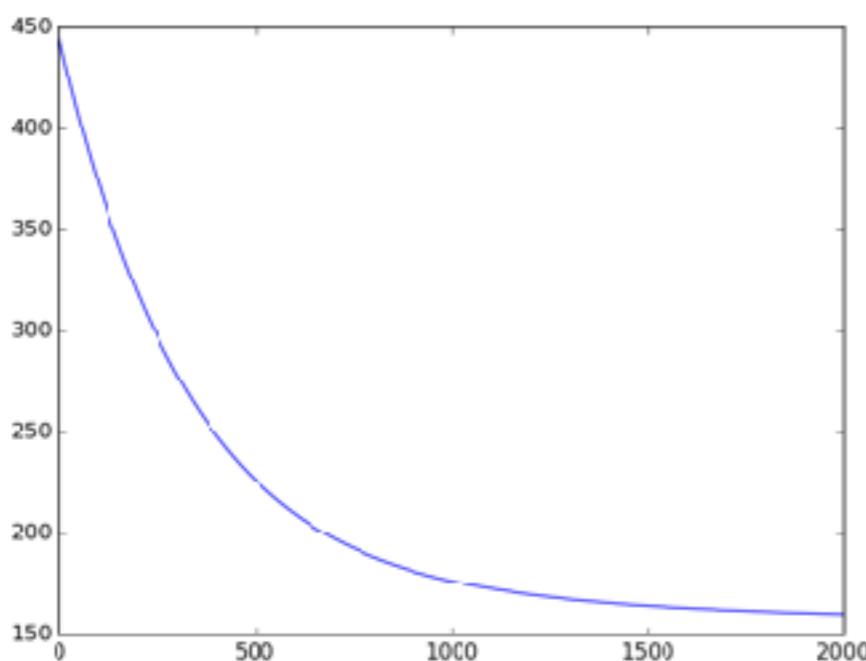
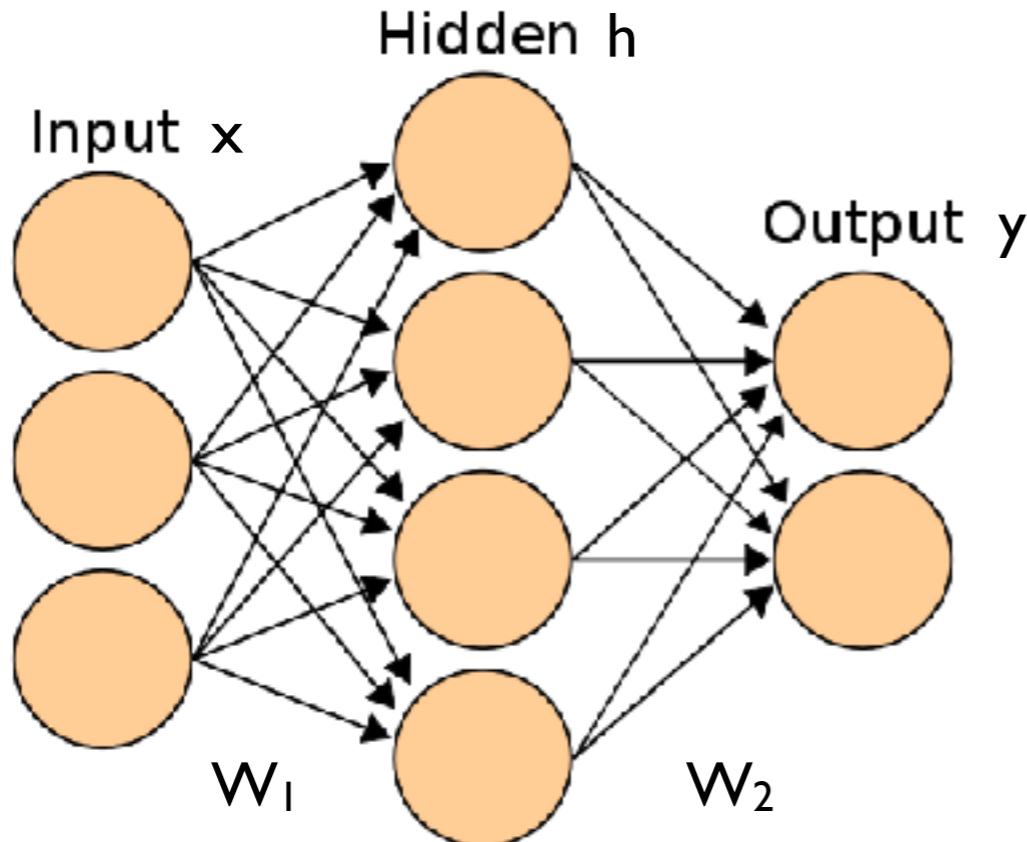


- ▶ Created by Google
- ▶ Math library used for machine learning and neural networks
- ▶ Supports Python and C++
- ▶ Pros
 - ✓ Documentation
 - ✓ Backed by large community
 - ✓ In built monitoring for training processes (Tensorboard)
- ▶ Cons
 - ✗ Static computational graphs
 - ✗ Higher learning curve than other libraries (low level, debugging harder)
 - ✗ Some performance issues

PyTorch

- ▶ Created by Facebook
- ▶ Tensor computation (like numpy) with GPU acceleration
- ▶ Supports Python
- ▶ Pros
 - ✓ Dynamic computational graphs (useful for e.g. RNNs)
 - ✓ Lower learning curve (more pythonic, easier to debug)
 - ✓ Easy to write own layer types
- ▶ Cons
 - ✗ Lacks in built monitoring
 - ✗ Not yet production ready (at v0.4 but less of an issue for research)
 - ✗ Documentation not as detailed

PyTorch Example



```
autograd_example.py <input>
1 import torch
2 import matplotlib.pyplot as plt
3
4 device = torch.device('cpu')
5
6 N, D_in, H, D_out = 64, 3, 4, 2
7
8 x = torch.randn(N, D_in, device=device)
9 y = torch.randn(N, D_out, device=device)
10
11 w1 = torch.randn(D_in, H, device=device, requires_grad=True)
12 w2 = torch.randn(H, D_out, device=device, requires_grad=True)
13
14 learning_rate = 1e-5
15 loss_data = []
16
17 for t in range(2000):
18     h = x.mm(w1)
19     phi = h.tanh()
20     y_pred = phi.mm(w2)
21
22     loss = (y_pred - y).pow(2).sum()
23     loss_data.append(loss.item())
24
25     loss.backward()
26
27     with torch.no_grad():
28         w1 -= learning_rate * w1.grad
29         w2 -= learning_rate * w2.grad
30         w1.grad.zero_()
31         w2.grad.zero_()
32
33 plt.plot(loss_data)
34 plt.show()
```

Deep learning in practice



- ▶ Keras is built on top of TensorFlow/Theano
- ▶ Supports Python
- ▶ Pros
 - ✓ Easiest learning curve
 - ✓ Very intuitive interface for building neural networks
 - ✓ Easy to write own layer types
- ▶ Cons
 - ✗ High level and not always as customisable
 - ✗ Not as many functionalities, less control

Loss Function

- The first thing to do to train a neural network is define a loss function
- For **continuous** outputs these include the mean squared error and mean absolute error

$$E(\mathbf{w}) = \frac{1}{N} \sum_i (y_i - \hat{y}_i(\mathbf{w}))^2 \quad E(\mathbf{w}) = \frac{1}{N} \sum_i |y_i - \hat{y}_i(\mathbf{w})|$$

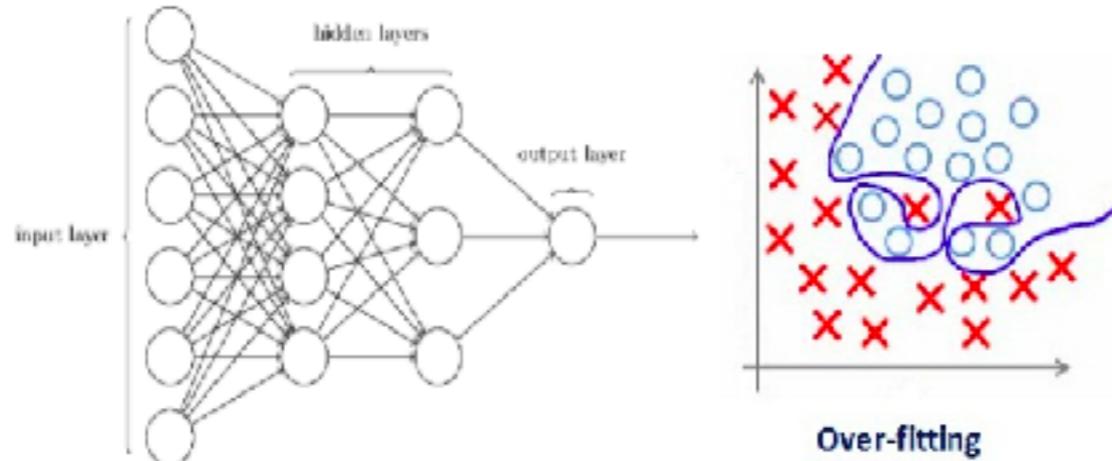
- Full loss function can include additional regularization terms
- For **categorical** outputs loss function is usually the categorical cross-entropy
- Last layer typically has a **soft-max** activation (turns M outputs into normalized probabilities)

$$E((\mathbf{w})) = - \sum_{i=1}^n \sum_{m=0}^{M-1} y_{im} \log \hat{y}_{im}(\mathbf{w}) + (1 - y_{im}) \log [1 - \hat{y}_{im}(\mathbf{w})] \quad y_{im} = \begin{cases} 1, & \text{if } y_i = m \\ 0, & \text{otherwise.} \end{cases}$$

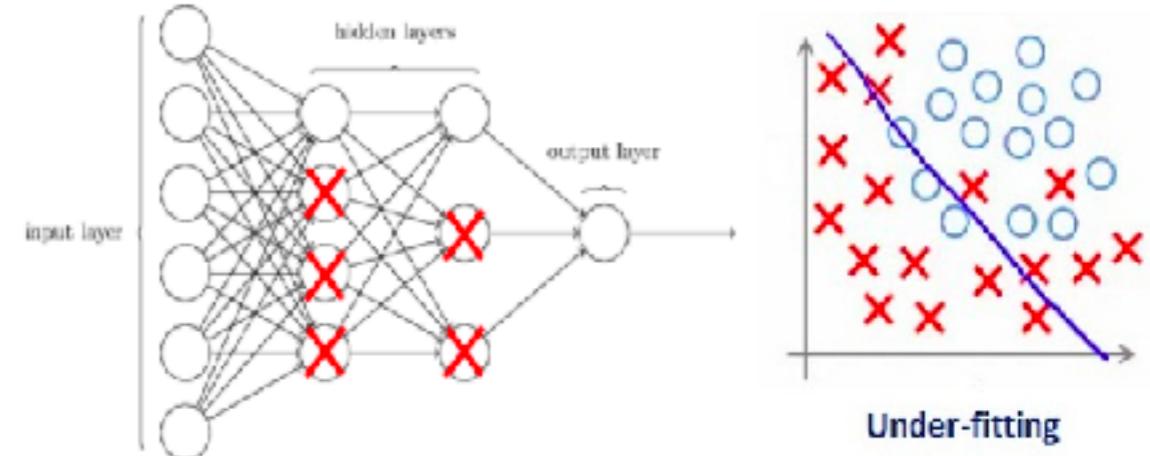
- Loss increased when predicted probability further from actual label

Regularization

- ▶ Regularization helps ensure neural networks do not **over-fit** and generalise well to unseen data
- ▶ L1 and L2 regularization apply penalties on a per layer basis during optimization
- ▶ **Dropout** turns off random neurons with probability p for each mini batch during training



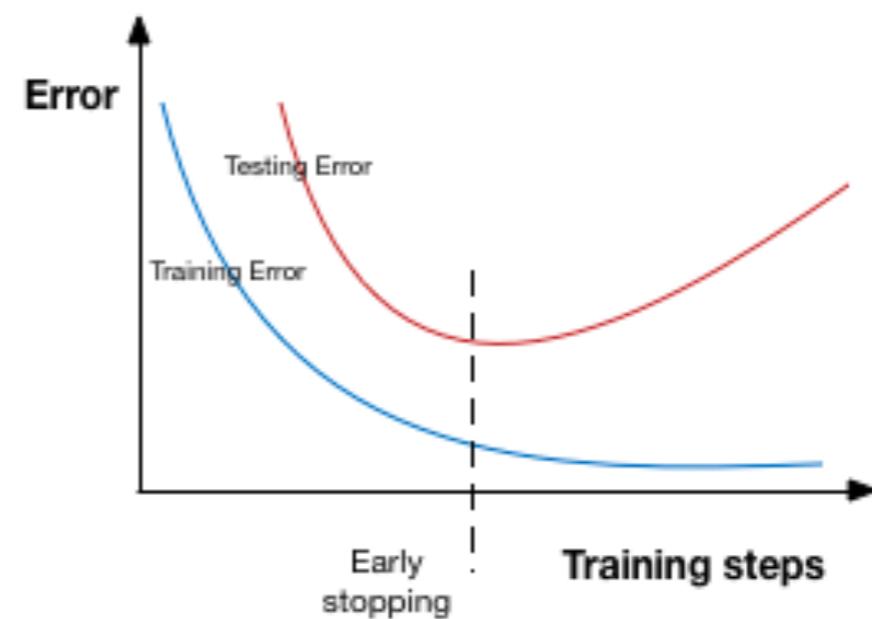
Model has too many free parameters
and data is over-fitted



Too many neurons either dropped out
or regularisation so high that weight
matrices close to zero - under-fitting!

Training

- ▶ Data should first be separated into training, test and holdout datasets
 - Training data is used to fit the parameters of the network
 - Test data to evaluate the performance of the trained model on unseen data. It can be used to tune the various hyper parameters of the network (e.g. number of layers, hidden units etc)
 - Holdout data is used to assess the final performance of the tuned model
- ▶ During training test loss can be monitoring, and training should be stopped when this increases
- ▶ Early stopping can also be seen as a type of regularization and avoids overfitting

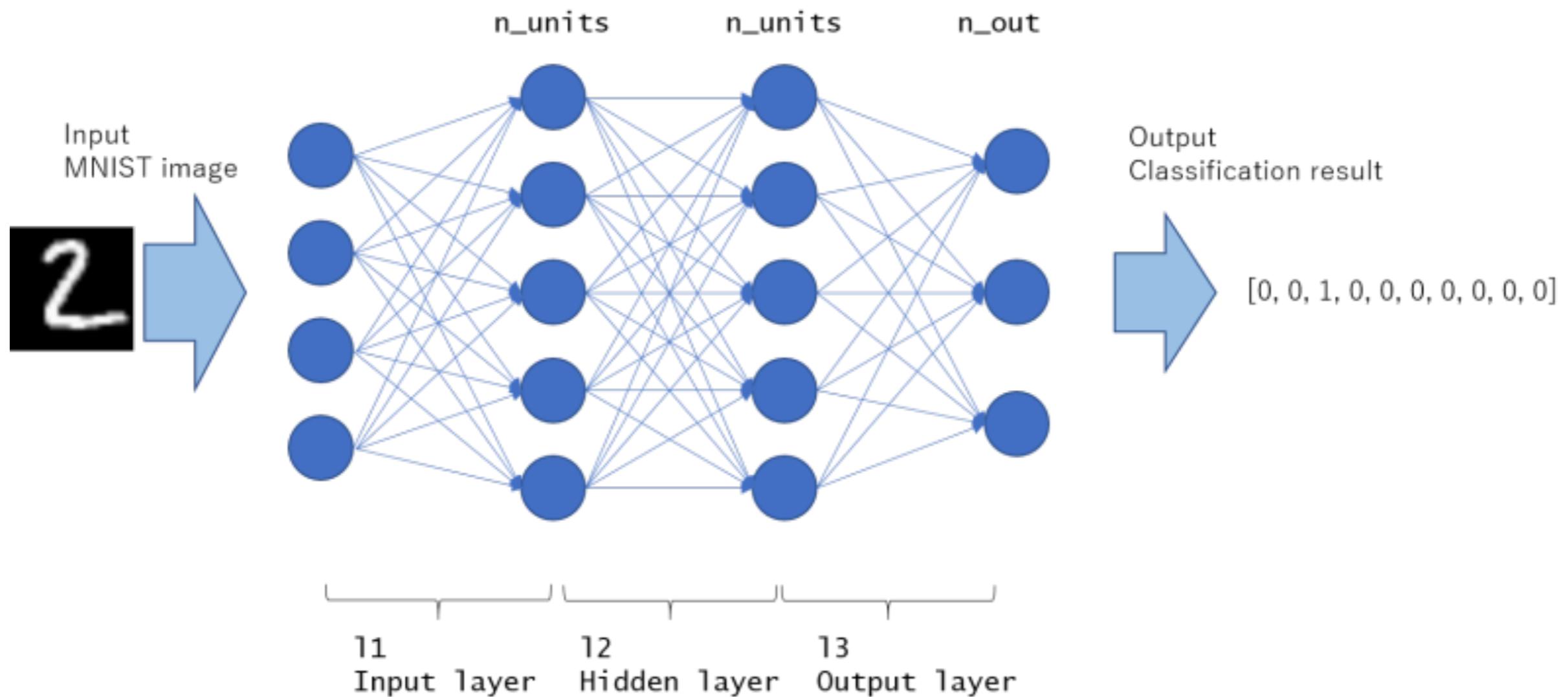


Hello World

- MNIST is ‘Hello Word’ of deep learning
 - Black and white images of integers from 0 to 9
 - 28 x 28 pixel images
 - 60,000 training images and 10,000 test images



MNIST MLP



MNIST MLP (Keras)

- ▶ Very simple API!
- ▶ 2 hidden layers each with 512 units
- ▶ Relu activation functions
- ▶ Dropout with p=0.2
- ▶ Final layer has softmax activation
- ▶ Gets to 98.40% test accuracy after 20 epochs
- ▶ State of the art has 99.8% accuracy



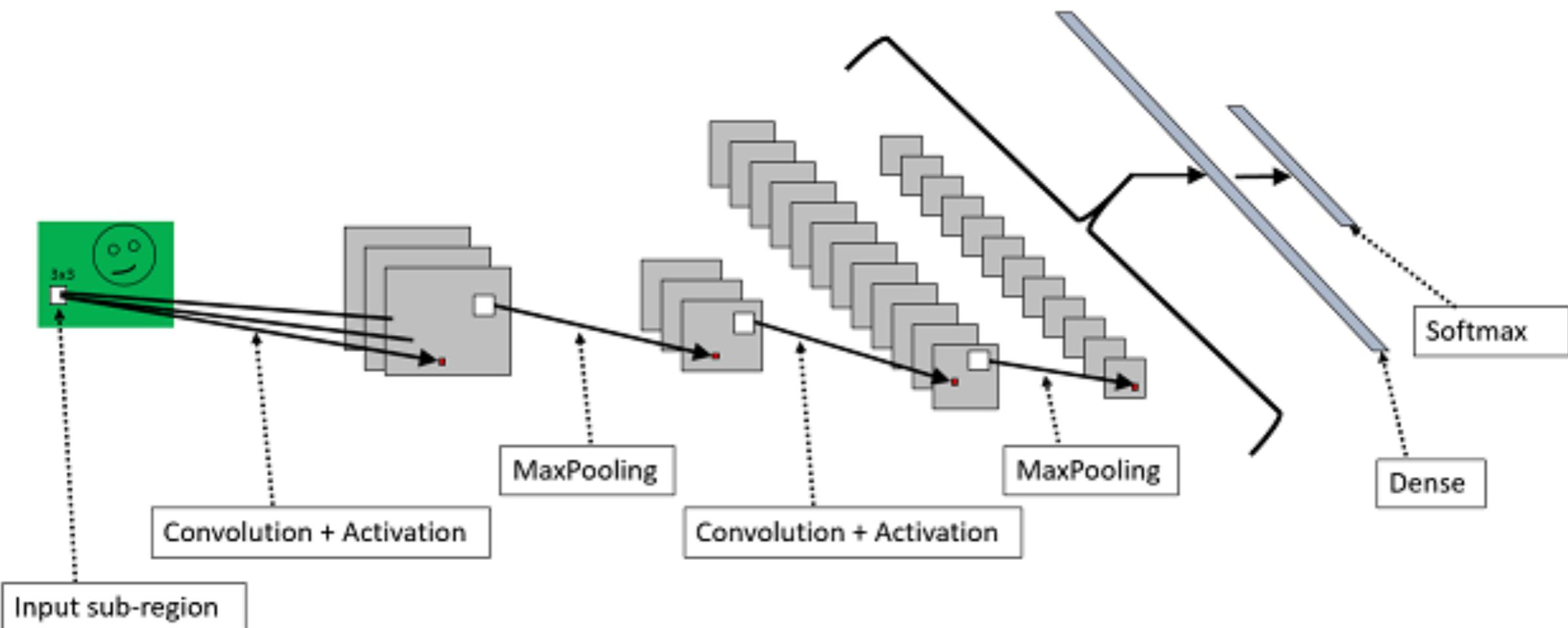
```
autograd_example.py  x  keras_mnist_mlp.py  x
1  from __future__ import print_function
2
3  import keras
4  from keras.datasets import mnist
5  from keras.models import Sequential
6  from keras.layers import Dense, Dropout
7  from keras.optimizers import RMSprop
8
9  batch_size = 128
10 num_classes = 10
11 epochs = 20
12
13 # the data, split between train and test sets
14 (x_train, y_train), (x_test, y_test) = mnist.load_data()
15
16 x_train = x_train.reshape(60000, 784)
17 x_test = x_test.reshape(10000, 784)
18 x_train = x_train.astype('float32')
19 x_test = x_test.astype('float32')
20 x_train /= 255
21 x_test /= 255
22 print(x_train.shape[0], 'train samples')
23 print(x_test.shape[0], 'test samples')
24
25 # convert class vectors to binary class matrices
26 y_train = keras.utils.to_categorical(y_train, num_classes)
27 y_test = keras.utils.to_categorical(y_test, num_classes)
28
29 model = Sequential()
30 model.add(Dense(512, activation='relu', input_shape=(784,)))
31 model.add(Dropout(0.2))
32 model.add(Dense(512, activation='relu'))
33 model.add(Dropout(0.2))
34 model.add(Dense(num_classes, activation='softmax'))
35
36 model.summary()
37
38 model.compile(loss='categorical_crossentropy',
39                 optimizer=RMSprop(),
40                 metrics=['accuracy'])
41
42 history = model.fit(x_train, y_train,
43                      batch_size=batch_size,
44                      epochs=epochs,
45                      verbose=1,
46                      validation_data=(x_test, y_test))
47 score = model.evaluate(x_test, y_test, verbose=0)
48 print('Test loss:', score[0])
49 print('Test accuracy:', score[1])
```

MNIST MLP (pyTorch)

```
4 IP: autograd_exmples.py * kernal_mnist_mlp.py * pytorch_mnist_mlp.py *  
1 from __future__ import print_function  
2 import argparse  
3 import torch  
4 import torch.nn as nn  
5 import torch.nn.functional as F  
6 import torch.optim as optim  
7 from torchvision import datasets, transforms  
8  
9  
10 class Net(nn.Module):  
11     def __init__(self):  
12         super(Net, self).__init__()  
13         self.fc1 = nn.Linear(28*28, 512)  
14         self.fc1_drop = nn.Dropout(0.2)  
15         self.fc2 = nn.Linear(512, 512)  
16         self.fc2_drop = nn.Dropout(0.2)  
17         self.fc3 = nn.Linear(512, 512)  
18  
19     def forward(self, x):  
20         x = x.view(-1, 28*28)  
21         x = F.relu(self.fc1(x))  
22         x = self.fc1_drop(x)  
23         x = F.relu(self.fc2(x))  
24         x = self.fc2_drop(x)  
25         x = F.log_softmax(self.fc3(x))  
26  
27  
28 def train(args, model, device, train_loader, optimizer, epoch):  
29     model.train()  
30     for batch_idx, (data, target) in enumerate(train_loader):  
31         data, target = data.to(device), target.to(device)  
32         optimizer.zero_grad()  
33         output = model(data)  
34         loss = F.nll_loss(output, target)  
35         loss.backward()  
36         optimizer.step()  
37         if batch_idx % args.log_interval == 0:  
38             print('Train Epoch: {} [{}/{} ({:.0f}%)]\tLoss: {:.6f}'.format(  
39                 epoch, batch_idx * len(data), len(train_loader.dataset),  
40                 100. * batch_idx / len(train_loader), loss.item()))  
41  
42 def test(args, model, device, test_loader):  
43     model.eval()  
44     test_loss = 0  
45     correct = 0  
46     with torch.no_grad():  
47         for data, target in test_loader:  
48             data, target = data.to(device), target.to(device)  
49             output = model(data)  
50             test_loss += F.nll_loss(output, target, reduction='sum').item() # sum up batch loss  
51             pred = output.max(1, keepdim=True)[1] # get the index of the max log-probability  
52             correct += pred.eq(target.view_as(pred)).sum().item()  
53  
54     test_loss /= len(test_loader.dataset)  
55     print('\nTest set: Average loss: {:.4f}, Accuracy: {}/{} ({:.0f}%)\n'.format(  
56         test_loss, correct, len(test_loader.dataset),  
57         100. * correct / len(test_loader.dataset)))  
58  
59  
60
```

```
def main():  
    # Training settings  
    parser = argparse.ArgumentParser(description='PyTorch MNIST Example')  
    parser.add_argument('--batch-size', type=int, default=128, metavar='N',  
                        help='input batch size for training (default: 64)')  
    parser.add_argument('--test-batch-size', type=int, default=1000, metavar='N',  
                        help='input batch size for testing (default: 1000)')  
    parser.add_argument('--epochs', type=int, default=20, metavar='N',  
                        help='number of epochs to train (default: 10)')  
    parser.add_argument('--lr', type=float, default=0.01, metavar='LR',  
                        help='learning rate (default: 0.01)')  
    parser.add_argument('--momentum', type=float, default=0.5, metavar='M',  
                        help='SGD momentum (default: 0.5)')  
    parser.add_argument('--no-cuda', action='store_true', default=False,  
                        help='disables CUDA training')  
    parser.add_argument('--seed', type=int, default=1, metavar='S',  
                        help='random seed (default: 1)')  
    parser.add_argument('--log-interval', type=int, default=10, metavar='N',  
                        help='how many batches to wait before logging training status')  
    args = parser.parse_args()  
    use_cuda = not args.no_cuda and torch.cuda.is_available()  
  
    torch.manual_seed(args.seed)  
  
    device = torch.device("cuda" if use_cuda else "cpu")  
  
    kwargs = {'num_workers': 1, 'pin_memory': True} if use_cuda else {}  
    train_loader = torch.utils.data.DataLoader(  
        datasets.MNIST('../data', train=True, download=True,  
                      transform=transforms.Compose([  
                          transforms.ToTensor(),  
                          transforms.Normalize((0.1307,), (0.3081,))  
                      ]),  
                      batch_size=args.batch_size, shuffle=True, **kwargs)  
    test_loader = torch.utils.data.DataLoader(  
        datasets.MNIST('../data', train=False, transform=transforms.Compose([  
            transforms.ToTensor(),  
            transforms.Normalize((0.1307,), (0.3081,))  
        ]),  
        batch_size=args.test_batch_size, shuffle=True, **kwargs)  
  
    model = Net().to(device)  
    optimizer = optim.SGD(model.parameters(), lr=args.lr, momentum=args.momentum)  
  
    for epoch in range(1, args.epochs + 1):  
        train(args, model, device, train_loader, optimizer, epoch)  
        test(args, model, device, test_loader)  
  
    if __name__ == '__main__':  
        main()
```

MNIST CNN



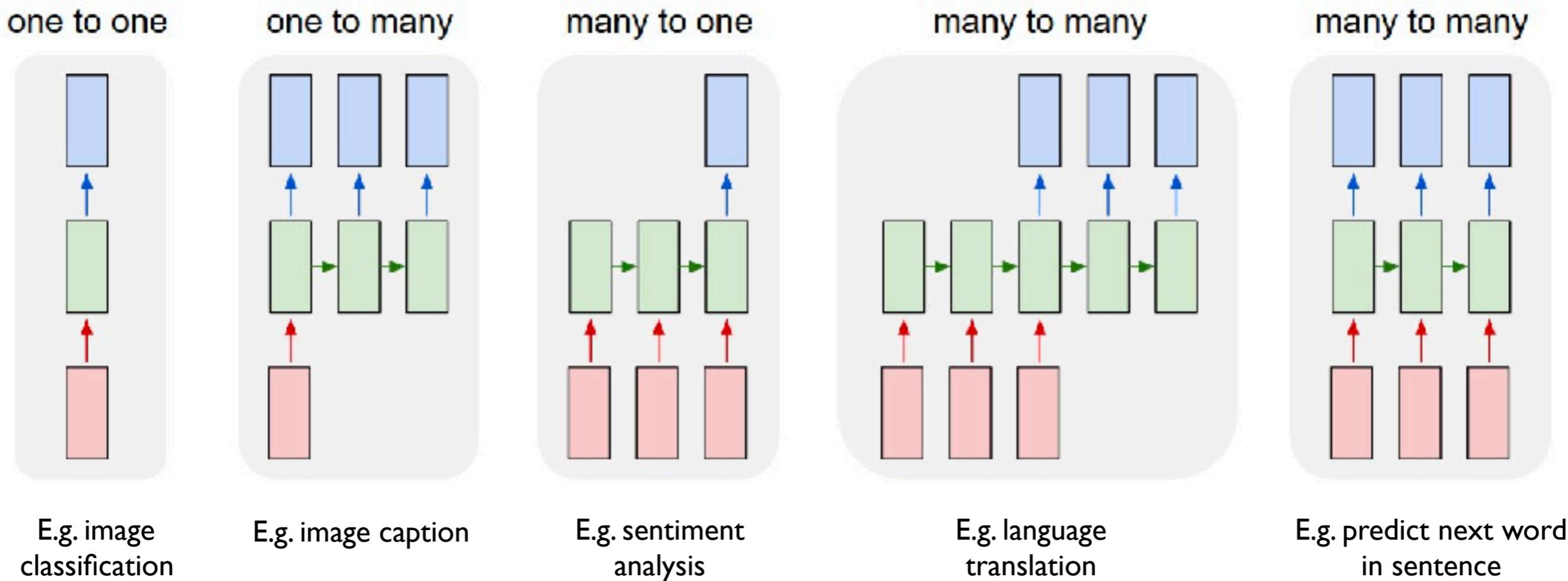
MNIST CNN

- ▶ First uses 32 convolutional units with 3x3 kernel
- ▶ Next 64 convolutional units with 3x3 kernel
- ▶ Max pooling down-samples units
- ▶ Dropout
- ▶ Dense fully connected layer
- ▶ Additional dropout and softmax activation
- ▶ Gets to 99.25% test accuracy after 12 epochs
- ▶ Each epoch takes ~16 seconds on K520 GPU

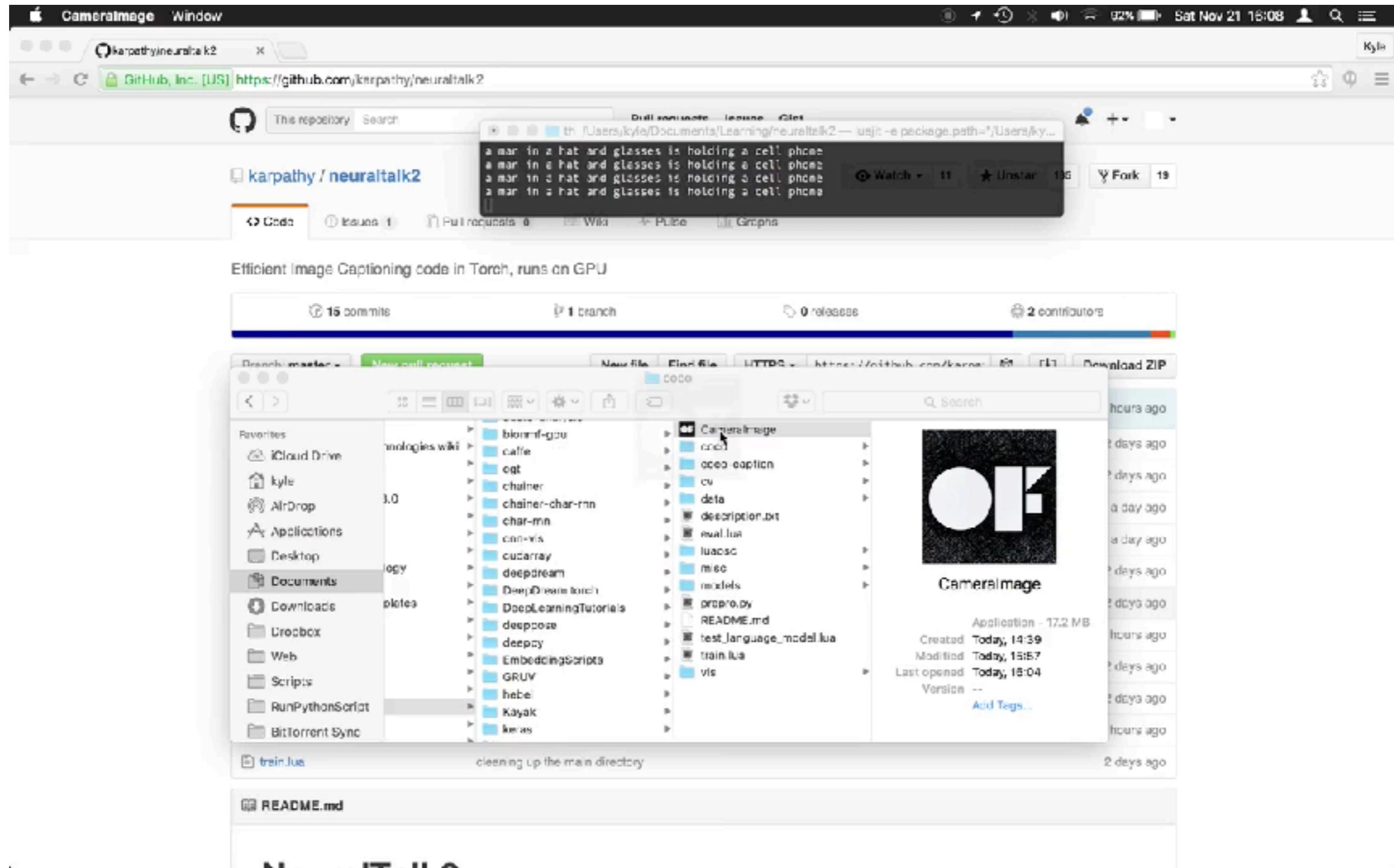
```
14
15 batch_size = 128
16 num_classes = 10
17 epochs = 12
18
19 # input image dimensions
20 img_rows, img_cols = 28, 28
21
22 # the data, split between train and test sets
23 (x_train, y_train), (x_test, y_test) = mnist.load_data()
24
25 if K.image_data_format() == 'channels_first':
26     x_train = x_train.reshape(x_train.shape[0], 1, img_rows, img_cols)
27     x_test = x_test.reshape(x_test.shape[0], 1, img_rows, img_cols)
28     input_shape = (1, img_rows, img_cols)
29 else:
30     x_train = x_train.reshape(x_train.shape[0], img_rows, img_cols, 1)
31     x_test = x_test.reshape(x_test.shape[0], img_rows, img_cols, 1)
32     input_shape = (img_rows, img_cols, 1)
33
34 x_train = x_train.astype('float32')
35 x_test = x_test.astype('float32')
36 x_train /= 255
37 x_test /= 255
38 print('x_train shape:', x_train.shape)
39 print(x_train.shape[0], 'train samples')
40 print(x_test.shape[0], 'test samples')
41
42 # convert class vectors to binary class matrices
43 y_train = keras.utils.to_categorical(y_train, num_classes)
44 y_test = keras.utils.to_categorical(y_test, num_classes)
45
46 model = Sequential()
47 model.add(Conv2D(32, kernel_size=(3, 3),
48                  activation='relu',
49                  input_shape=input_shape))
50 model.add(Conv2D(64, (3, 3), activation='relu'))
51 model.add(MaxPooling2D(pool_size=(2, 2)))
52 model.add(Dropout(0.25))
53 model.add(Flatten())
54 model.add(Dense(128, activation='relu'))
55 model.add(Dropout(0.5))
56 model.add(Dense(num_classes, activation='softmax'))
57
58 model.compile(loss=keras.losses.categorical_crossentropy,
59                 optimizer=keras.optimizers.Adadelta(),
60                 metrics=['accuracy'])
61
62 model.fit(x_train, y_train,
63            batch_size=batch_size,
64            epochs=epochs,
65            verbose=1,
66            validation_data=(x_test, y_test))
67 score = model.evaluate(x_test, y_test, verbose=0)
68 print('Test loss:', score[0])
69 print('Test accuracy:', score[1])
```

Recurrent Network

- Recurrent neural networks (RNNs) are a class of neural network that can learn about sequential data (e.g. time series, natural language)



Recurrent Network



Recurrent Network

Proof. Omitted. □

Lemma 0.1. Let \mathcal{C} be a set of the construction.

Let \mathcal{C} be a gerber covering. Let \mathcal{F} be a quasi-coherent sheaves of \mathcal{O} -modules. We have to show that

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

Proof. This is an algebraic space with the composition of sheaves \mathcal{F} on $X_{\text{étale}}$ we have

$$\mathcal{O}_X(\mathcal{F}) = \{\text{morph}_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

where \mathcal{G} defines an isomorphism $\mathcal{F} \rightarrow \mathcal{G}$ of \mathcal{O} -modules. □

Lemma 0.2. This is an integer Z is injective.

Proof. See Spaces, Lemma ??.

Lemma 0.3. Let S be a scheme. Let X be a scheme and X is an affine open covering. Let $\mathcal{U} \subset \mathcal{X}$ be a canonical and locally of finite type. Let X be a scheme. Let X be a scheme which is equal to the formal complex.

The following to the construction of the lemma follows.

Let X be a scheme. Let X be a scheme covering. Let

$$b : X \rightarrow Y' \rightarrow Y \rightarrow Y \rightarrow Y' \times_X Y \rightarrow X.$$

be a morphism of algebraic spaces over S and Y .

Proof. Let X be a nonzero scheme of X . Let X be an algebraic space. Let \mathcal{F} be a quasi-coherent sheaf of \mathcal{O}_X -modules. The following are equivalent

- (1) \mathcal{F} is an algebraic space over S .
- (2) If X is an affine open covering.

Consider a common structure on X and X the functor $\mathcal{O}_X(U)$ which is locally of finite type. □

This since $\mathcal{F} \in \mathcal{F}$ and $x \in \mathcal{G}$ the diagram

$$\begin{array}{ccc}
 S & \xrightarrow{\quad} & \\
 \downarrow & & \\
 \xi & \longrightarrow & \mathcal{O}_{X'} \\
 & \uparrow & \searrow \\
 & \text{gor}_x & \\
 \\
 = a' & \longrightarrow & \\
 \uparrow & & \\
 = a' & \longrightarrow & \alpha \\
 & & \\
 \text{Spec}(R_S) & & \text{Mor}_{\text{Site}} \quad d(\mathcal{O}_{X_{\mathcal{F}/S}}, \mathcal{G}) \\
 & & X \\
 & & \downarrow \\
 & & \text{d}(\mathcal{O}_{X_{\mathcal{F}/S}}, \mathcal{G})
 \end{array}$$

is a limit. Then \mathcal{G} is a finite type and assume S is a flat and \mathcal{F} and \mathcal{G} is a finite type f_* . This is of finite type diagrams, and

- the composition of \mathcal{G} is a regular sequence,
- $\mathcal{O}_{X'}$ is a sheaf of rings.

Proof. We have see that $X = \text{Spec}(R)$ and \mathcal{F} is a finite type representable by algebraic space. The property \mathcal{F} is a finite morphism of algebraic stacks. Then the cohomology of X is an open neighbourhood of U . □

Proof. This is clear that \mathcal{G} is a finite presentation, see Lemmas ??.
A reduced above we conclude that U is an open covering of \mathcal{C} . The functor \mathcal{F} is a “field”

$$\mathcal{O}_{X,S} \rightarrow \mathcal{F}_S \dashv (\mathcal{O}_{X_{\text{étale}}}) \rightarrow \mathcal{O}_{X'_S}^{-1} \mathcal{O}_{X'_S}(\mathcal{O}_{X'_S}^{\vee})$$

is an isomorphism of covering of $\mathcal{O}_{X'_S}$. If \mathcal{F} is the unique element of \mathcal{F} such that X is an isomorphism.

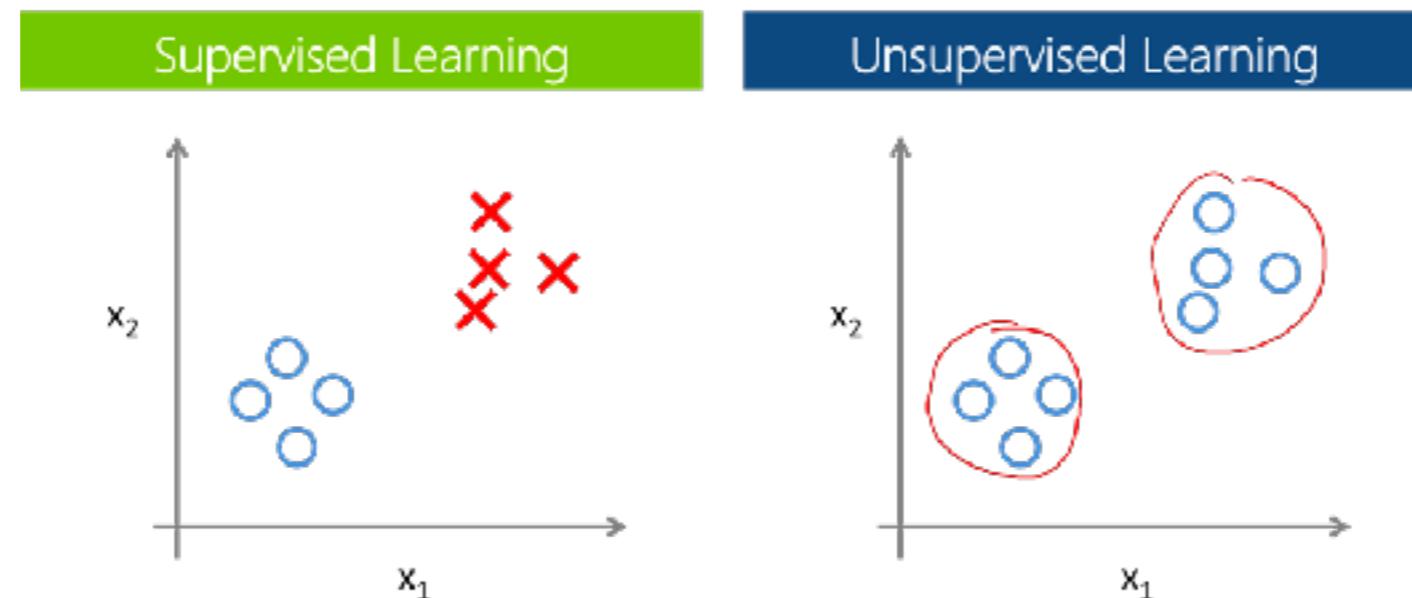
The property \mathcal{F} is a disjoint union of Proposition ?? and we can filtered set of presentations of a scheme \mathcal{O}_X -algebra with \mathcal{F} are opens of finite type over S .

If \mathcal{F} is a scheme theoretic image points. □

If \mathcal{F} is a finite direct sum $\mathcal{O}_{X'_S}$ is a closed immersion, see Lemma ??, This is a sequence of \mathcal{F} is a similar morphism.

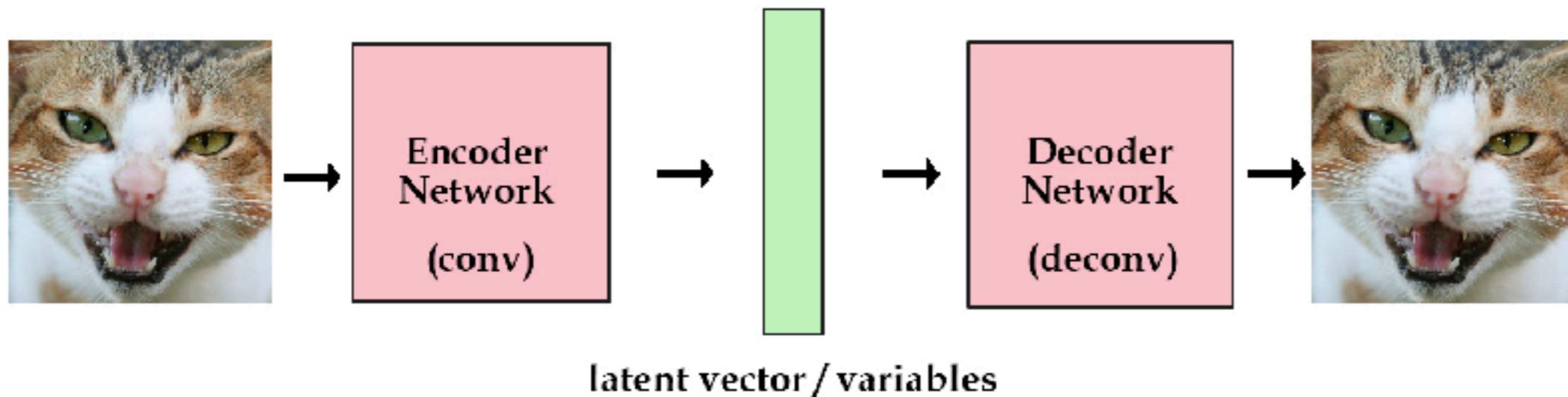
Unsupervised Learning

- ▶ **Much** more unlabelled than labelled data
- ▶ Labelling can be subjective
- ▶ Unsupervised learning finds the most important features in the dataset



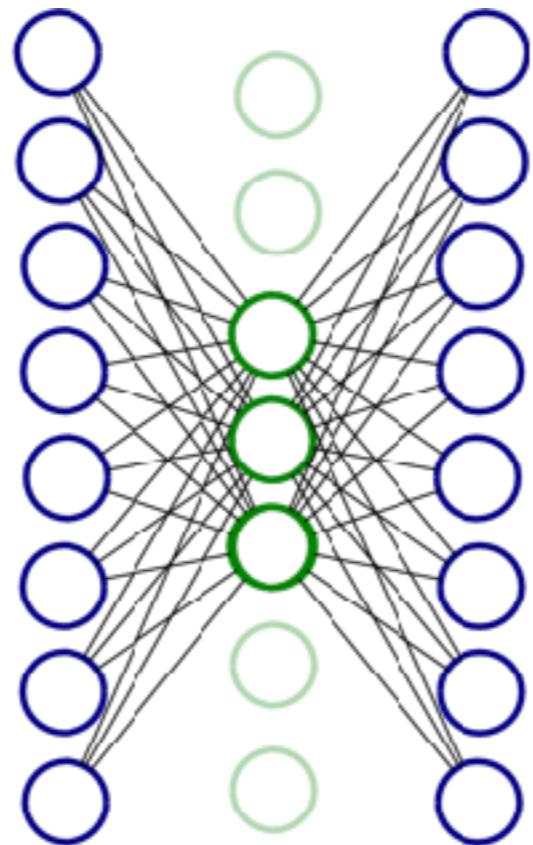
- ▶ Examples of unsupervised learning
 - Auto-encoders and variational auto-encoders (generative)
 - Adversarial networks (generative)
 - Boltzmann machines (generative)
 - Reinforcement learning (self-supervised)

Auto-encoders



- ▶ Data compression algorithm
 - Data specific (will only be able to compress data similar to what they are trained on)
 - Lossy (decompressed output degraded compared to input)
 - No labelling required

Auto-encoders



- ▶ Latent space should have fewer dimensions
 - Under-complete representation
 - Bottleneck architecture
- ▶ Otherwise over-complete representation might just learn the identity function

MNIST

- ▶ Used Keras Python library
- ▶ Discard labels
- ▶ Normalised values between 0 and 1 and reshape into 784 vector as before
- ▶ Encoding dimension 32, corresponds to compression factor of 24.5
- ▶ Can visualise reconstructed inputs from the test set

```
from keras.layers import Input, Dense
from keras.models import Model
from keras.datasets import mnist
import numpy as np

(x_train, _), (x_test, _) = mnist.load_data()

x_train = x_train.astype('float32') / 255.
x_test = x_test.astype('float32') / 255.
x_train = x_train.reshape((len(x_train), np.prod(x_train.shape[1:])))
x_test = x_test.reshape((len(x_test), np.prod(x_test.shape[1:])))
print x_train.shape
print x_test.shape

encoding_dim = 32

input_img = Input(shape=(784,))
encoded = Dense(encoding_dim, activation='relu')(input_img)
decoded = Dense(784, activation='sigmoid')(encoded)

# this model maps an input to its reconstruction
autoencoder = Model(input_img, decoded)

autoencoder.compile(optimizer='adadelta', loss='binary_crossentropy')

autoencoder.fit(x_train, x_train,
                epochs=50,
                batch_size=256,
                shuffle=True,
                validation_data=(x_test, x_test))

decoded_imgs = autoencoder.predict(x_test)
```



- ▶ Losing detail with this approach - how do we improve it?

CNN Auto-encoder

- ▶ Since inputs and outputs are images it makes sense to use convolutional network as encoder and decoder
- ▶ We also do not have to limit ourselves to a single latent layer
- ▶ Both encoder and decoder are stacked CNNs
- ▶ Loss function significantly improved (0.094 vs 0.11)

```
from keras.layers import Input, Dense
from keras.models import Model
from keras.layers import Conv2D, MaxPooling2D, UpSampling2D
from keras import backend as K
from keras.datasets import mnist
import numpy as np

(x_train, _), (x_test, _) = mnist.load_data()

x_train = x_train.astype('float32') / 255.
x_test = x_test.astype('float32') / 255.
x_train = np.reshape(x_train, (len(x_train), 28, 28, 1))
x_test = np.reshape(x_test, (len(x_test), 28, 28, 1))

input_img = Input(shape=(28, 28, 1))

x = Conv2D(16, (3, 3), activation='relu', padding='same')(input_img)
x = MaxPooling2D((2, 2), padding='same')(x)
x = Conv2D(8, (3, 3), activation='relu', padding='same')(x)
x = MaxPooling2D((2, 2), padding='same')(x)
x = Conv2D(8, (3, 3), activation='relu', padding='same')(x)
encoded = MaxPooling2D((2, 2), padding='same')(x)

# at this point the representation is (4, 4, 8) i.e. 128-dimensional
x = Conv2D(8, (3, 3), activation='relu', padding='same')(encoded)
x = UpSampling2D((2, 2))(x)
x = Conv2D(8, (3, 3), activation='relu', padding='same')(x)
x = UpSampling2D((2, 2))(x)
x = Conv2D(16, (3, 3), activation='relu')(x)
x = UpSampling2D((2, 2))(x)
decoded = Conv2D(1, (3, 3), activation='sigmoid', padding='same')(x)

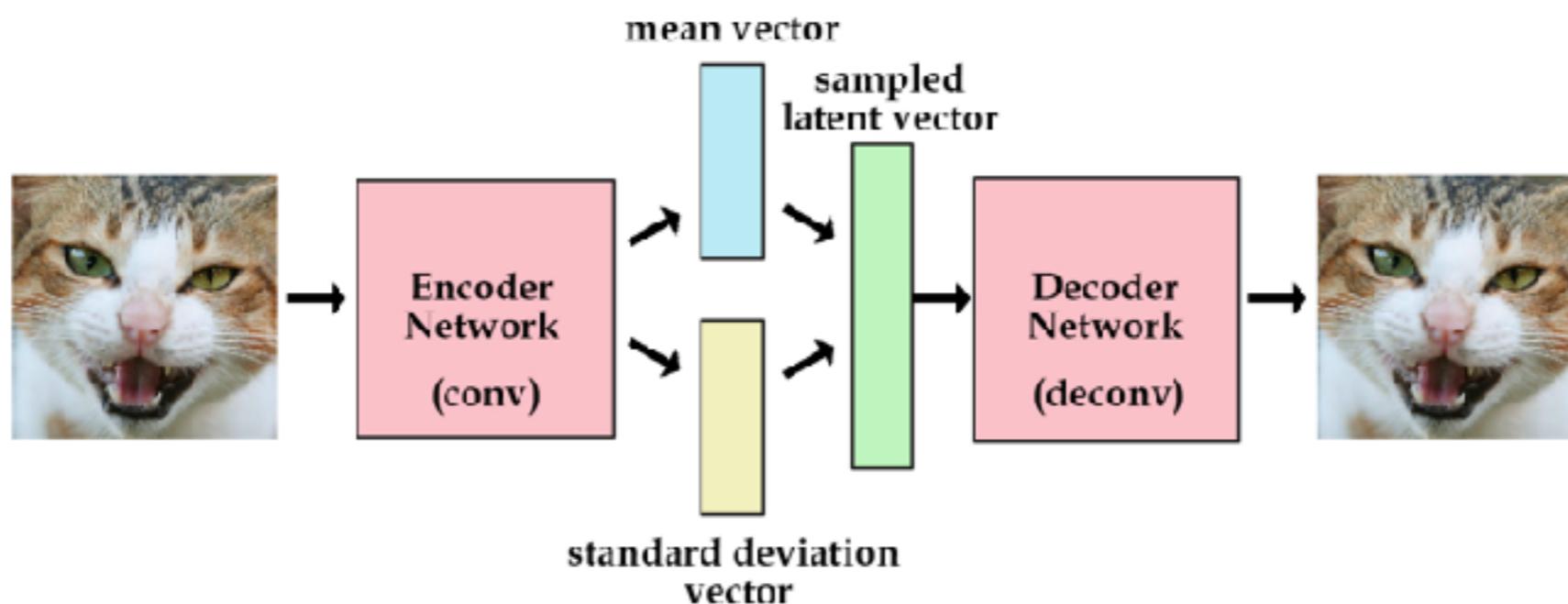
autoencoder = Model(input_img, decoded)
autoencoder.compile(optimizer='adadelta', loss='binary_crossentropy')

autoencoder.fit(x_train, x_train,
                 epochs=50,
                 batch_size=128,
                 shuffle=True,
                 validation_data=(x_test, x_test))
```



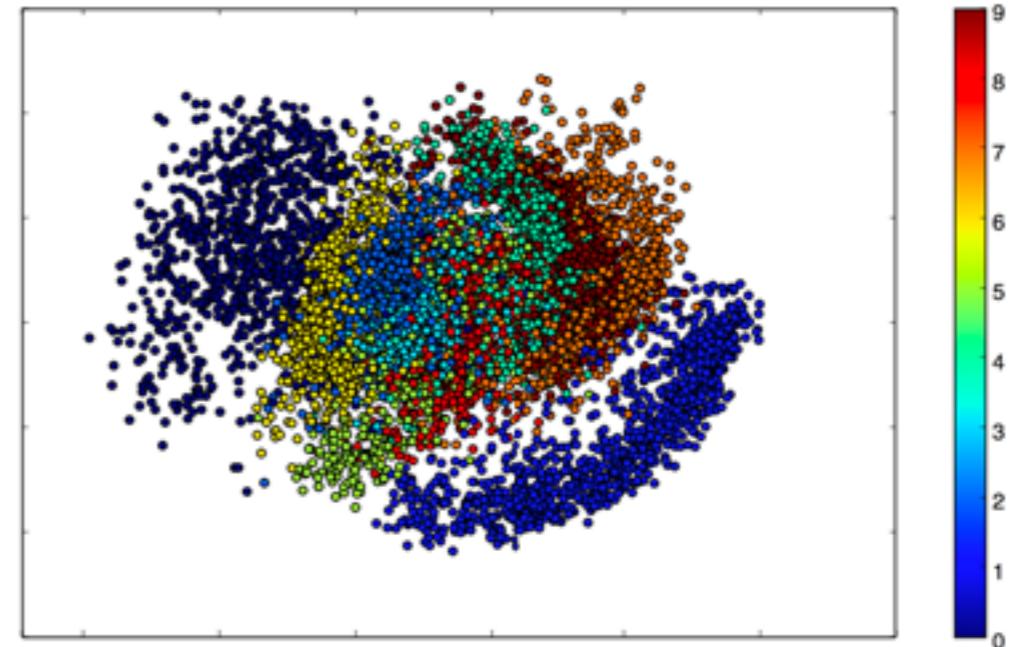
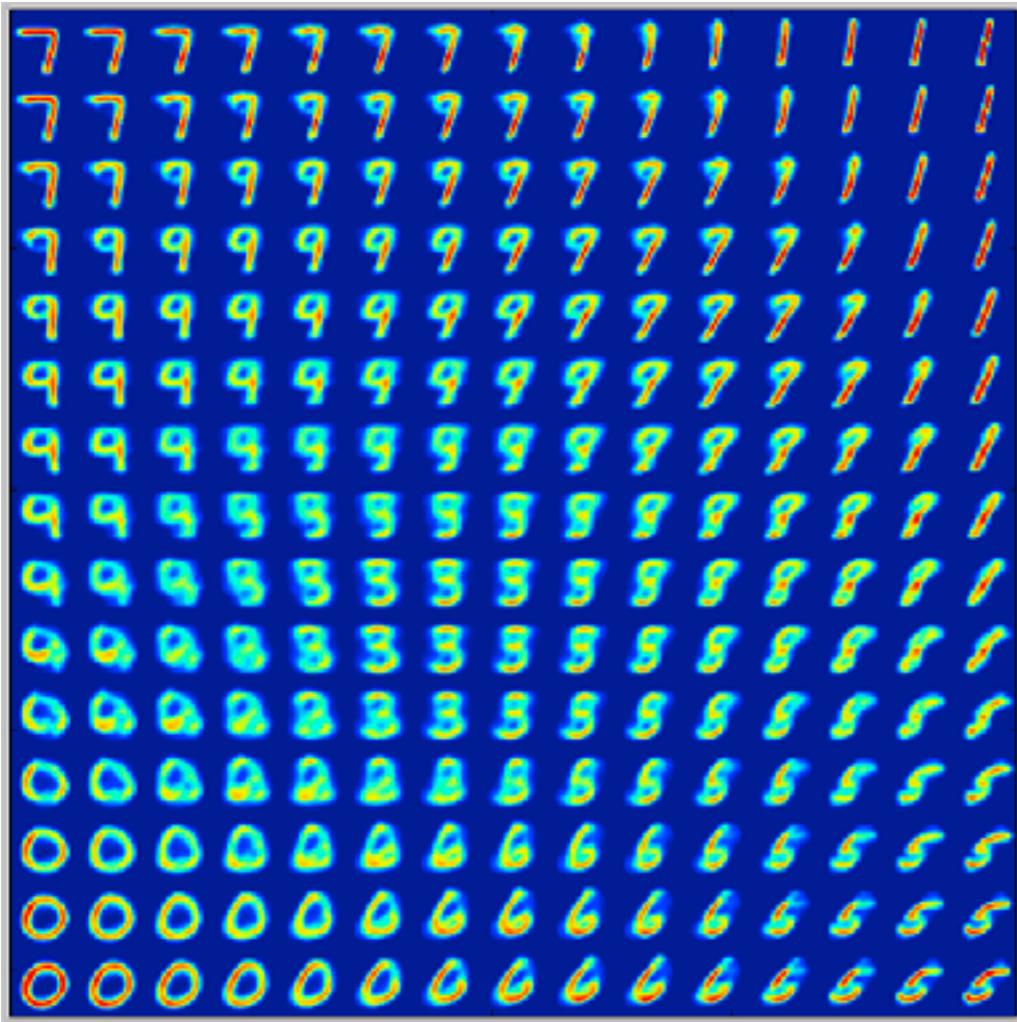
Variational AE

- ▶ An auto-encoder represents data in latent space
- ▶ Can't generate anything yet, since we do not know how to create latent vectors
- ▶ Variational AE adds constraints on latent space, i.e. it learns a latent model such as a Gaussian probability density
- ▶ If you then sample points from this probability density, you can generate new samples



Variational AE

- ▶ Since latent space is 2D can visualise it (each coloured cluster is digit)



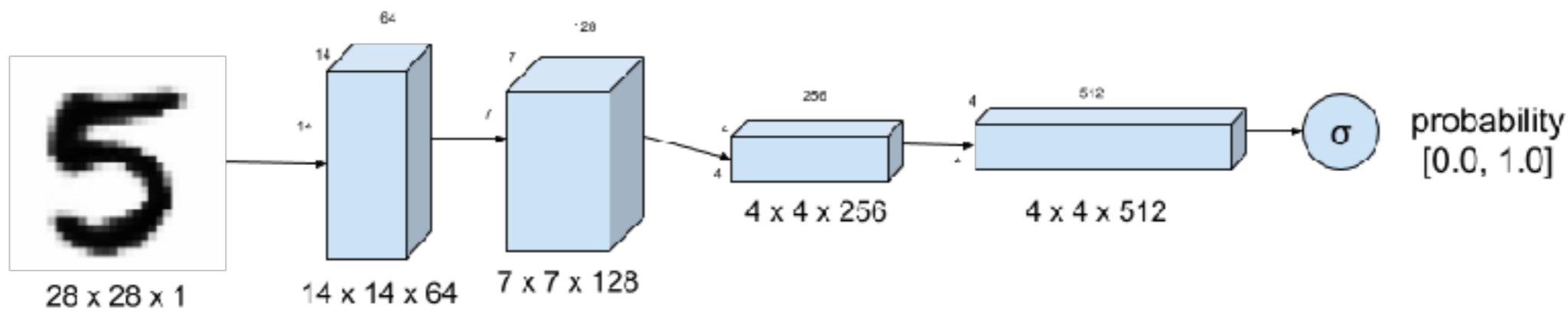
- ▶ Generated digits

Adversarial Networks

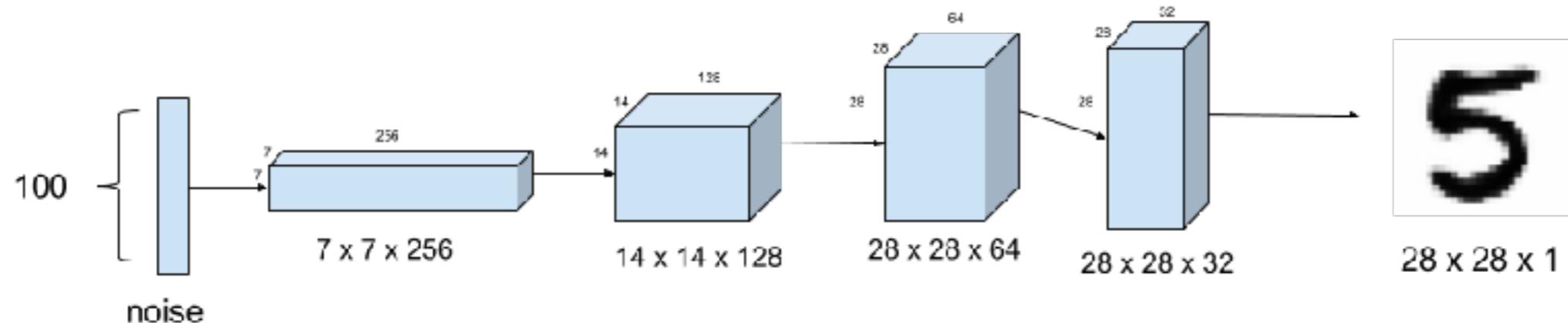
- ▶ VAEs are easy to train and have a tractable likelihood, but can produce images with lower fidelity
- ▶ GANs (Generative Adversarial Network) produce richer images but are much more difficult to train
- ▶ A GAN consists of 2 networks - a generator and discriminator, that compete and co-operate
- ▶ Real world example - counterfeiter (generator) and police (discriminator)
- ▶ Counterfeiter produces fakes, police give feedback why money is fake
- ▶ Counterfeiter improves fakes to trick police based on feedback
- ▶ At same time police are improving how to better tell apart real and fake money

GAN

- Discriminator tells if something is generative or real, e.g money, MNIST digits. Example is CNN

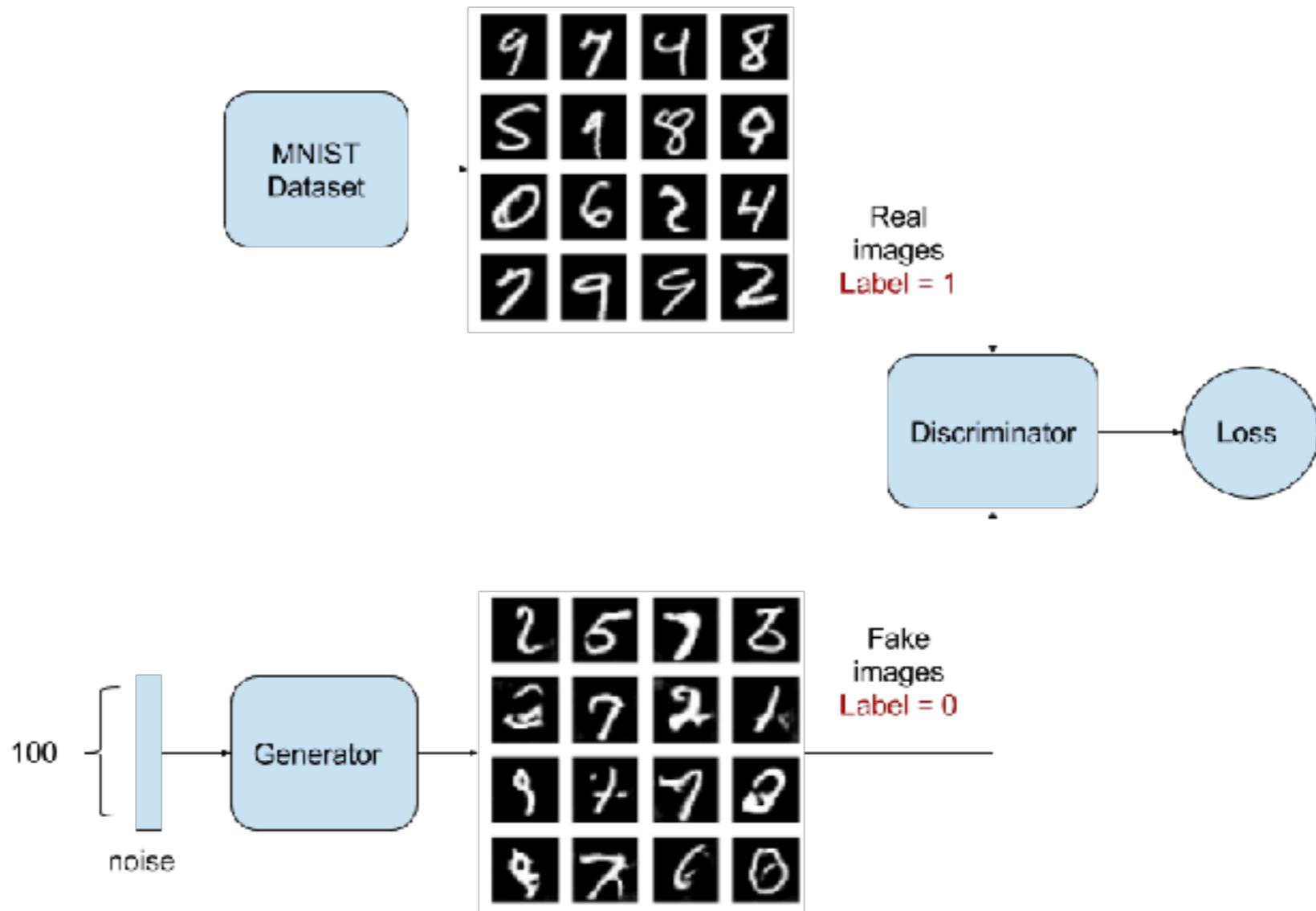


- Generator synthesises fakes. Input is generated from noise

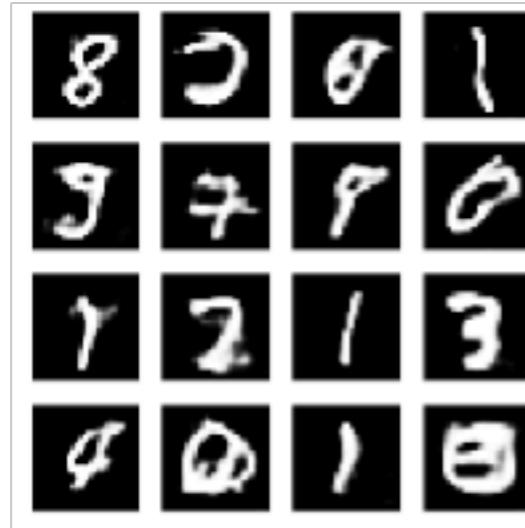


GAN

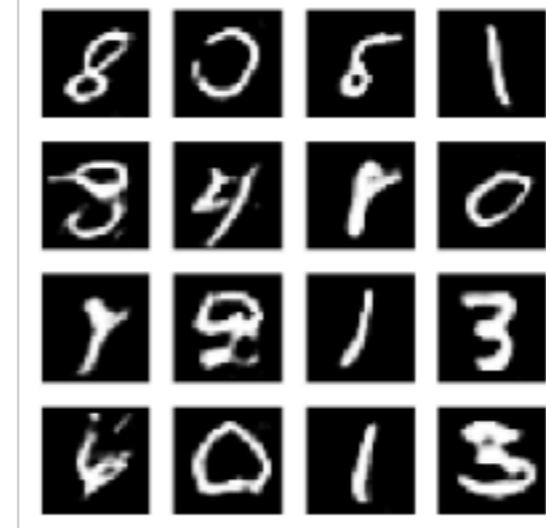
- ▶ Full discriminator model is generator + discriminator with both real and fake data correctly labelled



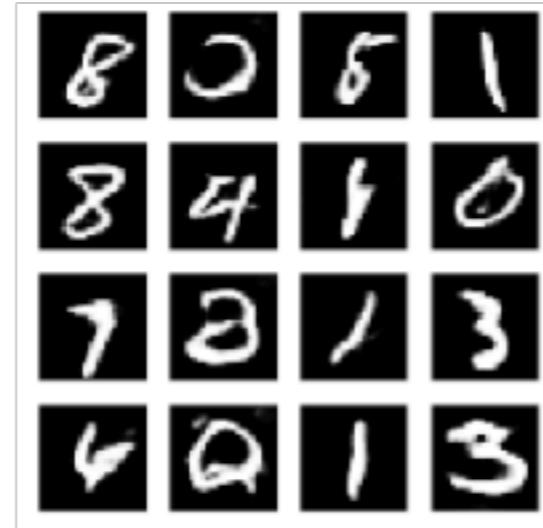
Samples



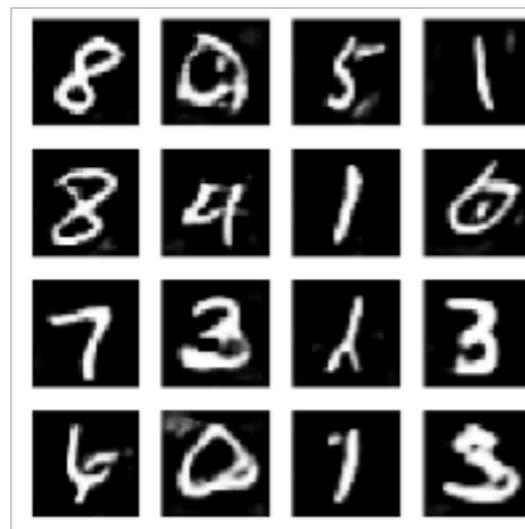
1000 steps



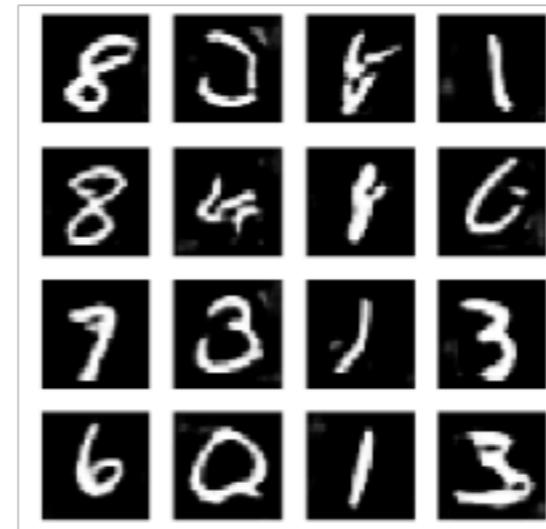
3000 steps



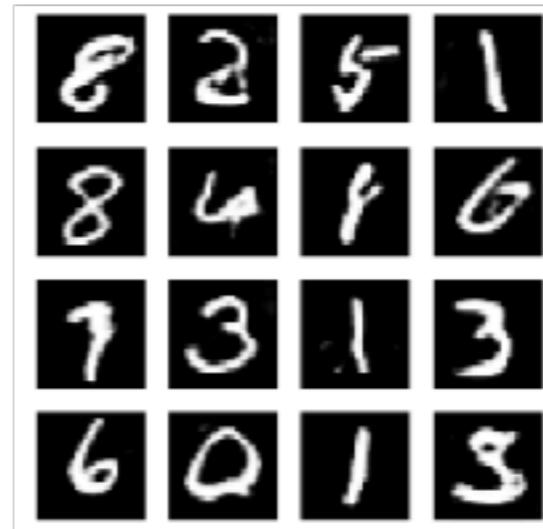
5000 steps



7000 steps



9000 steps



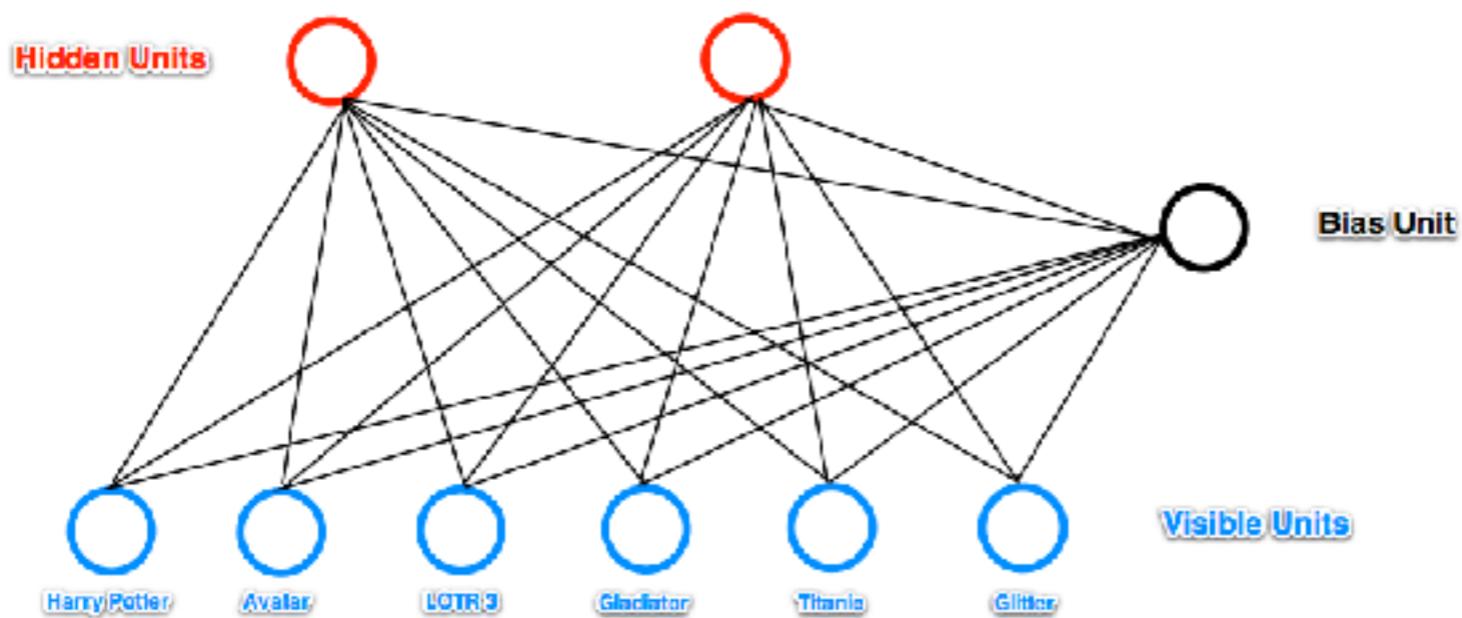
10000 steps

Other Examples



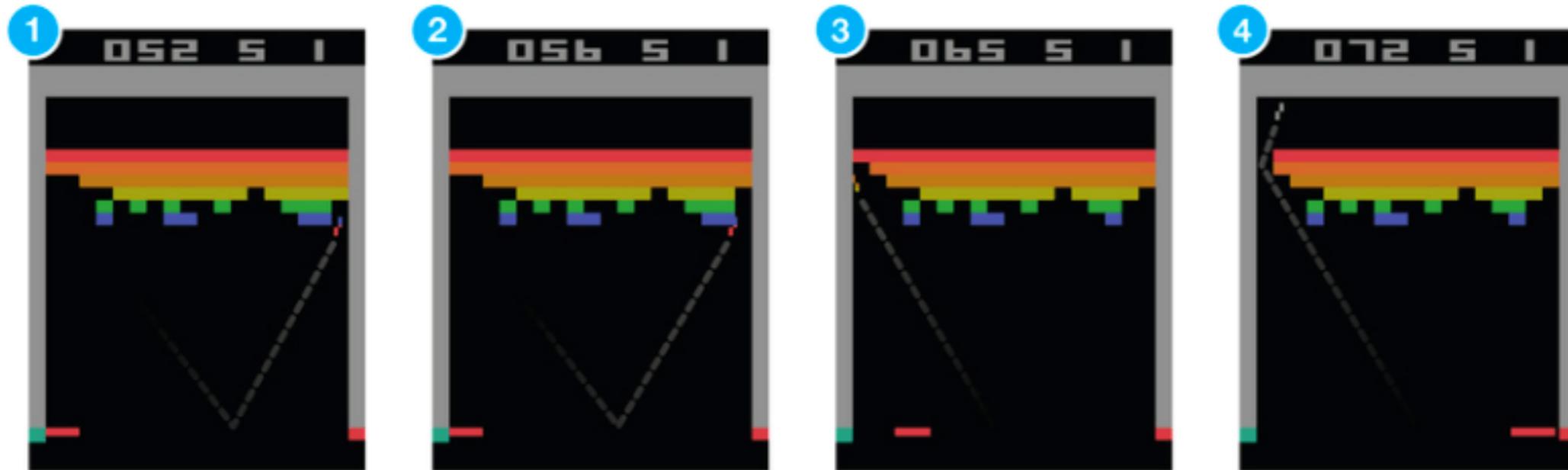
Boltzmann Machines

- ▶ Another type of generative network
- ▶ Important difference is that it is **stochastic**, i.e. activations in network have a probabilistic element
- ▶ Restricted Boltzmann Machine (RBM) consists of one visible layer and one hidden layer (contained latent factors)
- ▶ E.g. we have 6 movies and ask people which ones they want to watch. We will learn 2 latent variables



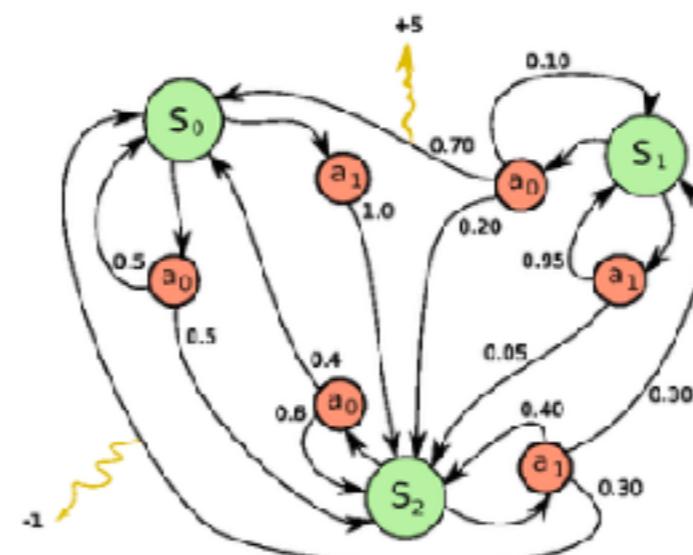
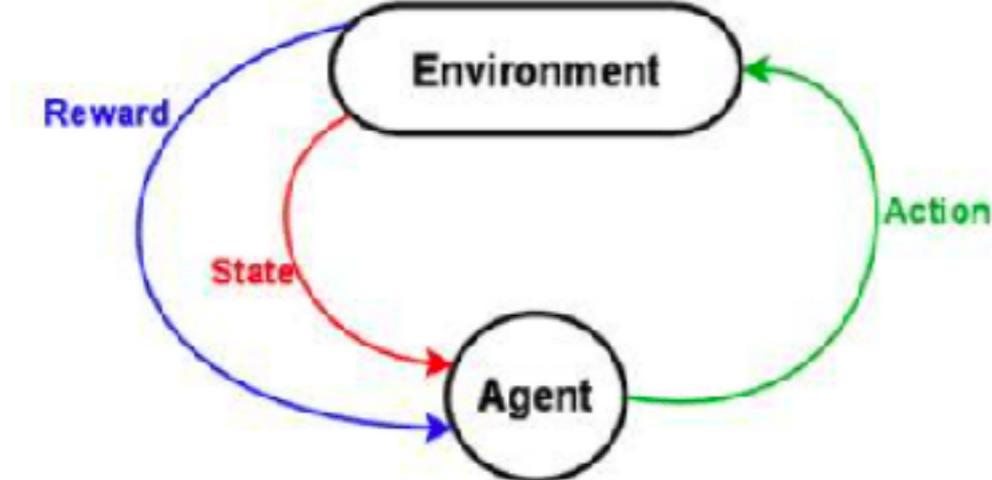
Reinforcement Learning

- ▶ Self-supervised learning, reward based feedback
- ▶ E.g. breakout game. Goal is to get ball though the blocks
- ▶ Can teach neural network how to play by giving the image of the screen as an input to the network
- ▶ Network will then decide which action to take (move paddle left or right) to maximise future possible reward



Reinforcement Learning

- ▶ In general have an agent in an environment
- ▶ Environment is in certain state
- ▶ Agent can perform actions. These sometimes result in a reward
- ▶ Actions transform environment and lead to new state
- ▶ Sequence of states, actions, rewards $s_0, a_0, r_0, s_1, a_1, r_1, \dots$
- ▶ A full sequence before terminal state is called an episode
- ▶ Rules for how to choose actions is called a policy



Q Learning

- ▶ Total discounted future reward in episode from time t onwards

$$R_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots$$

- ▶ Define Q function as maximum possible discounted future reward when we perform an action a in state s at time t

$$Q(s_t, a_t) = \max R_t$$

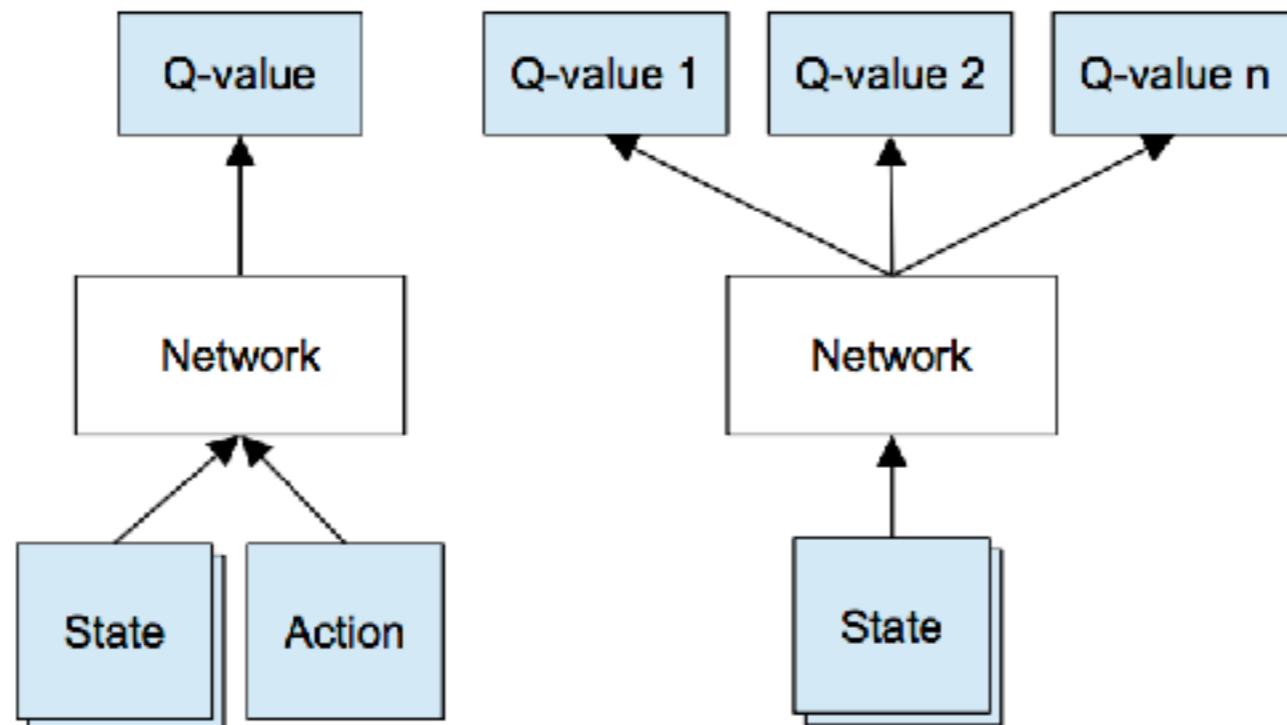
- ▶ Can be iteratively estimated using the **Bellman equation**

$$Q(s, a) = r + \gamma \max_{a'} Q(s', a')$$

- ▶ No neural network yet - could be done with a table of states vs actions

Deep Q Network

- ▶ Can train a neural network to estimate the Q function by playing out many episodes of a game



- ▶ Several subtleties - e.g. training using experience replay

Deep Q Network

