



# Deep Learning: Introduction

Adam Moss

School of Physics and Astronomy

[adam.moss@nottingham.ac.uk](mailto:adam.moss@nottingham.ac.uk)



The University of  
Nottingham

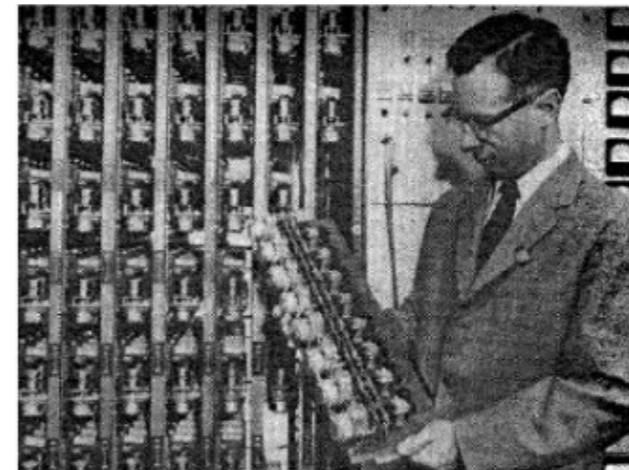
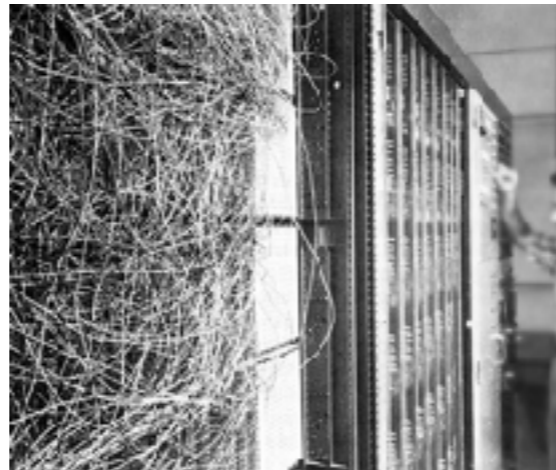
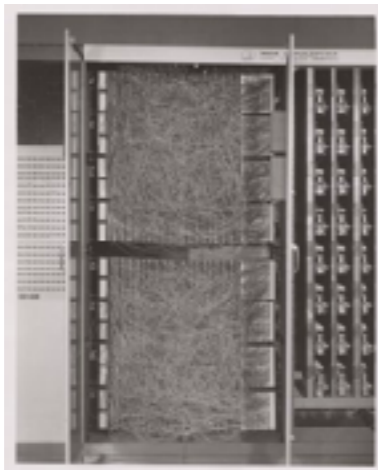
UNITED KINGDOM • CHINA • MALAYSIA

# Early Learning

- ▶ Feature extraction



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

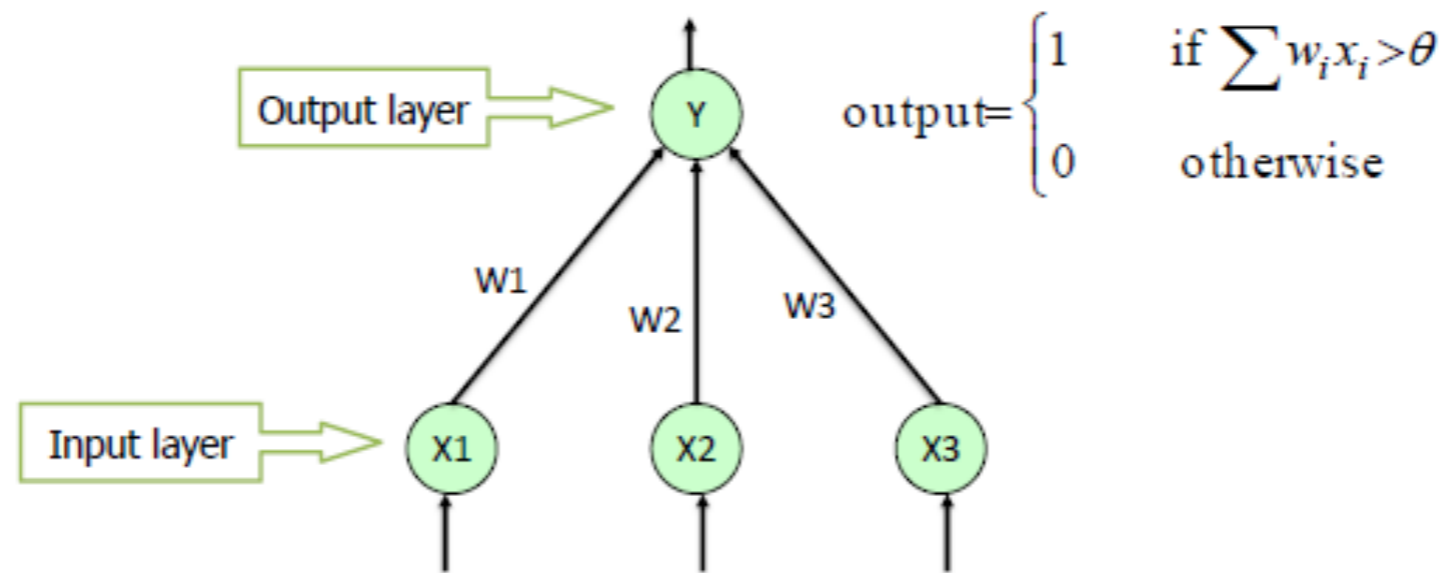




# Perceptron

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

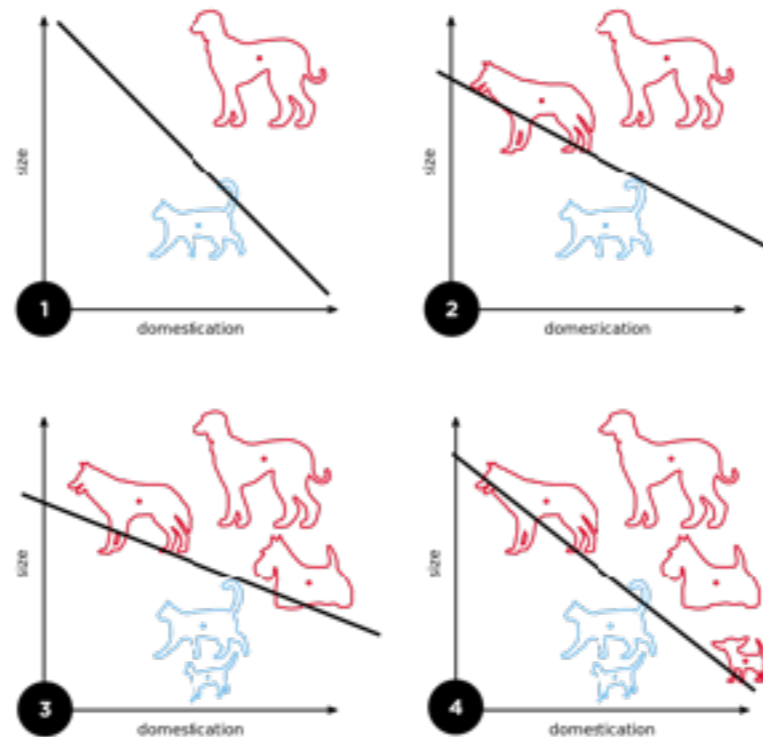
Single Layer Perceptron



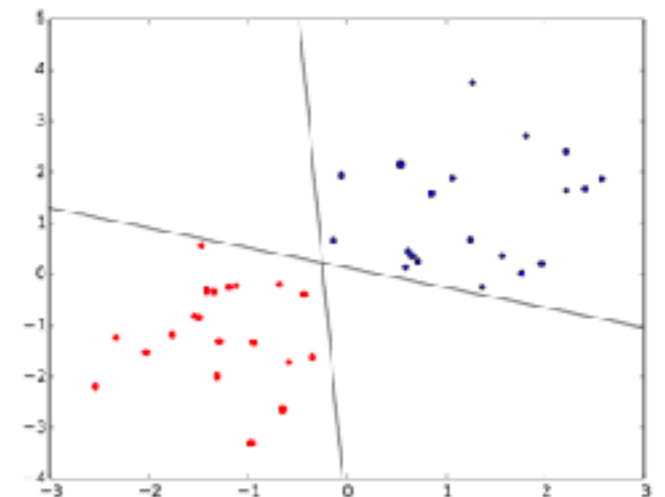
- ▶ Each input has a weight  $w_i$
- ▶ Weights are trained using **supervised learning**
- ▶ Training sets of  $D = \{x_{i,j}, y_j\}$  where  $j$  is the sample number and  $y_j$  the desired output for that sample
- ▶ Prediction is  $\hat{y}$  and weights are updated to minimise loss  $\sum_j (y_j -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

# Perceptron

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

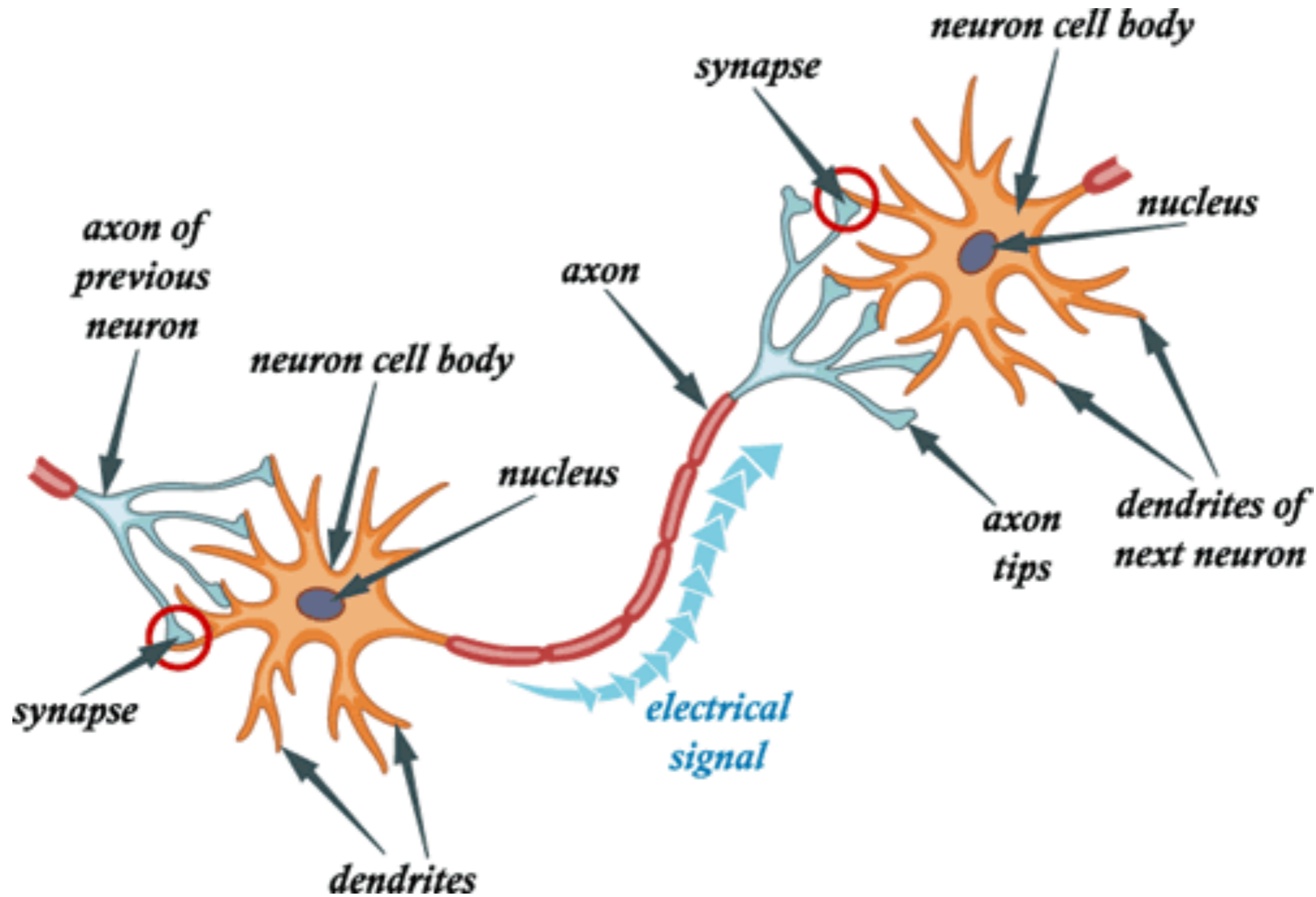


- ▶ May admit solutions of different quality (general problem for machine learning if training data is not representative)

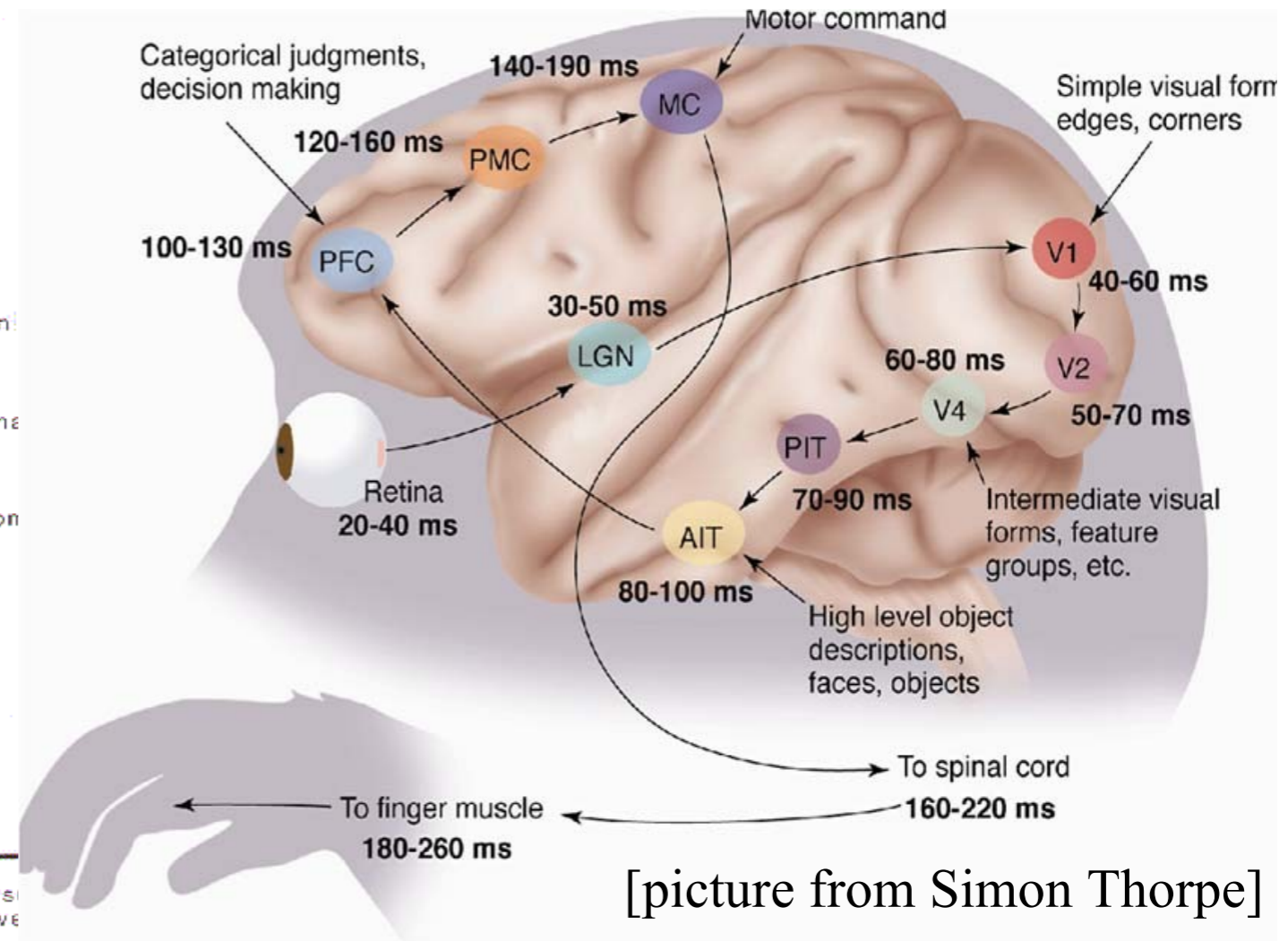
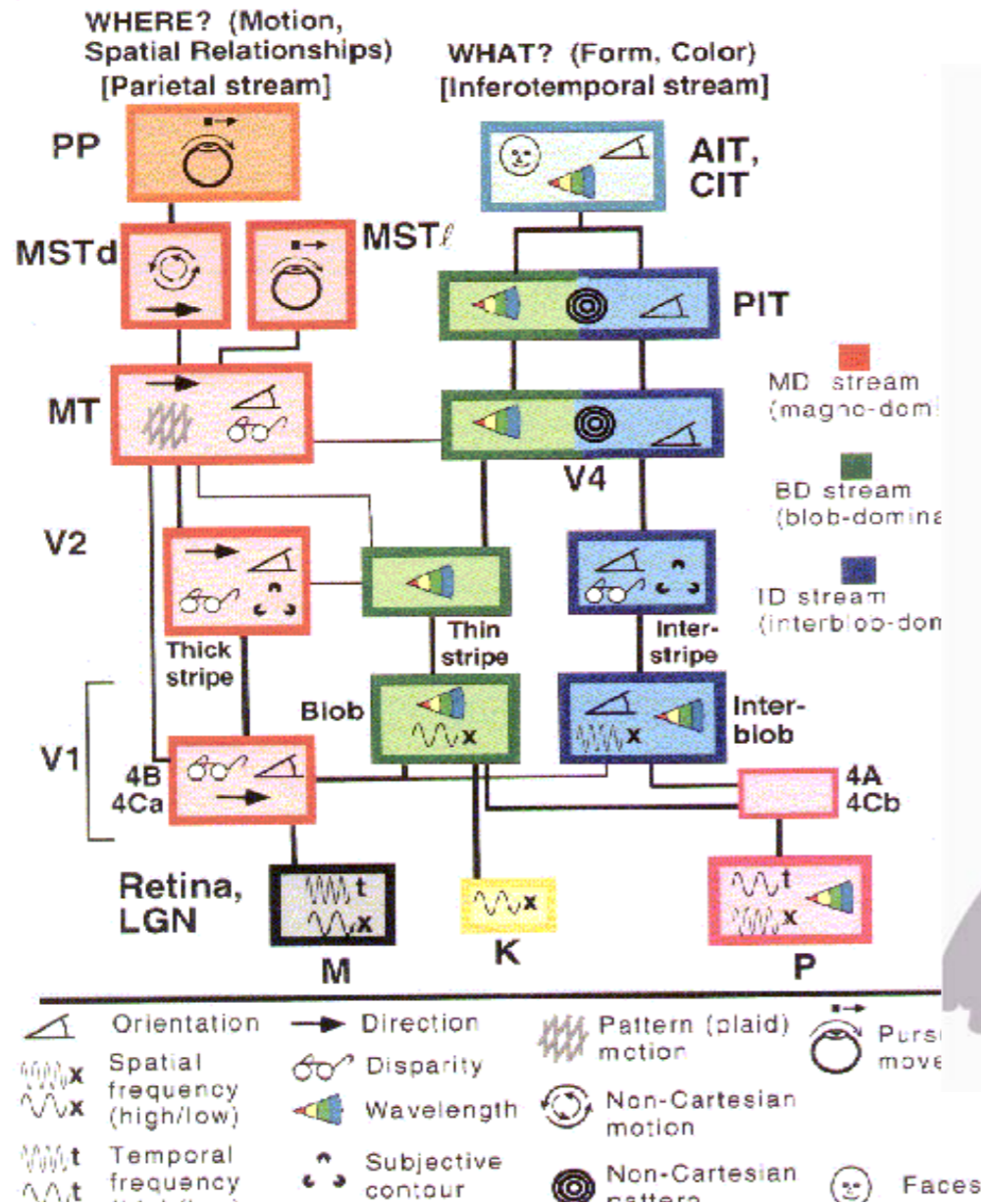




# Biological Neuron



# Visual Cortex



[picture from Simon Thorpe]

[Gallant & Van Essen]

AIT=anterior inferior temporal cortex

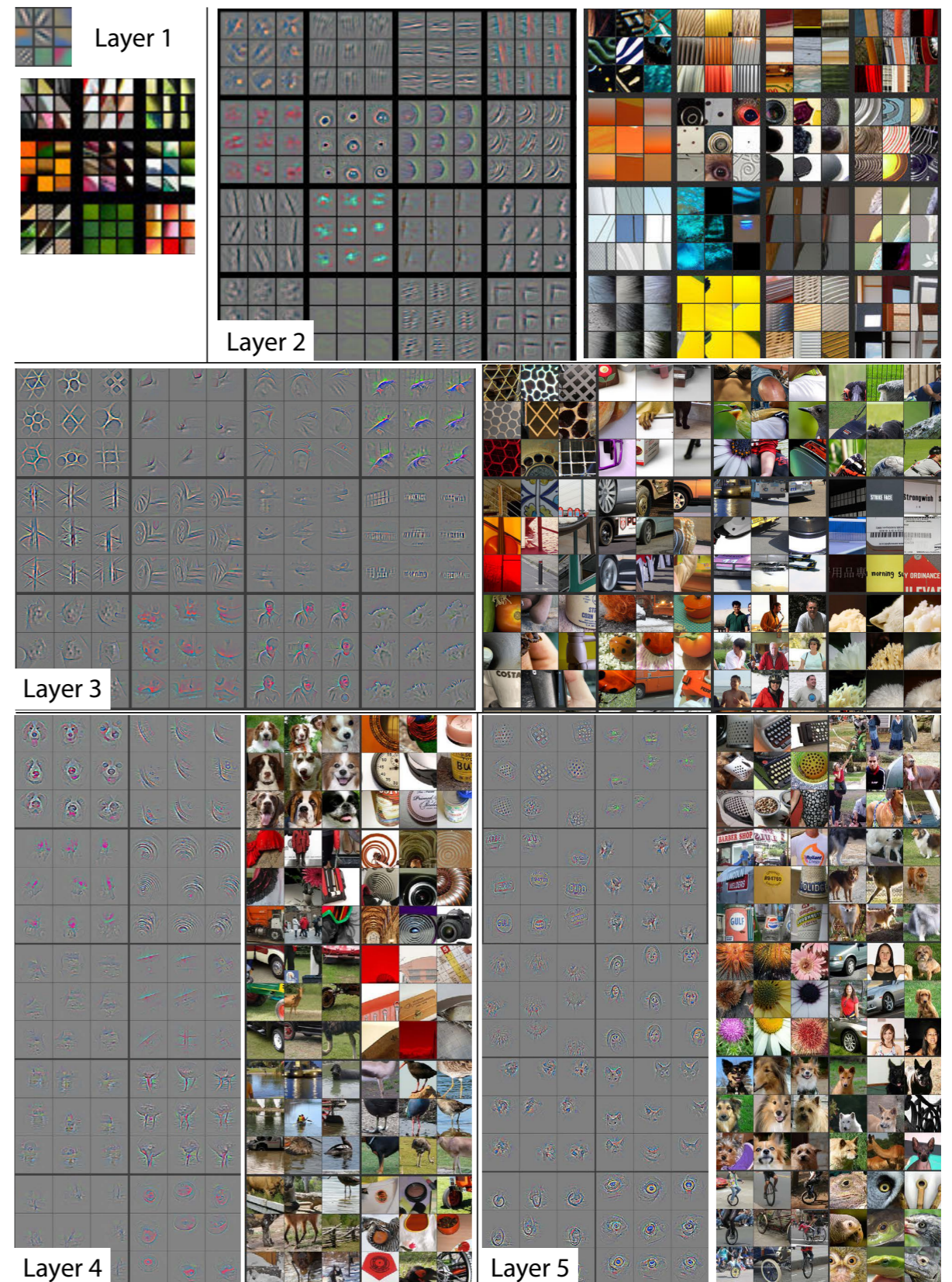
PIT=posterior inferior temporal cortex

LGN=lateral geniculate nucleus



# Shallow vs Deep

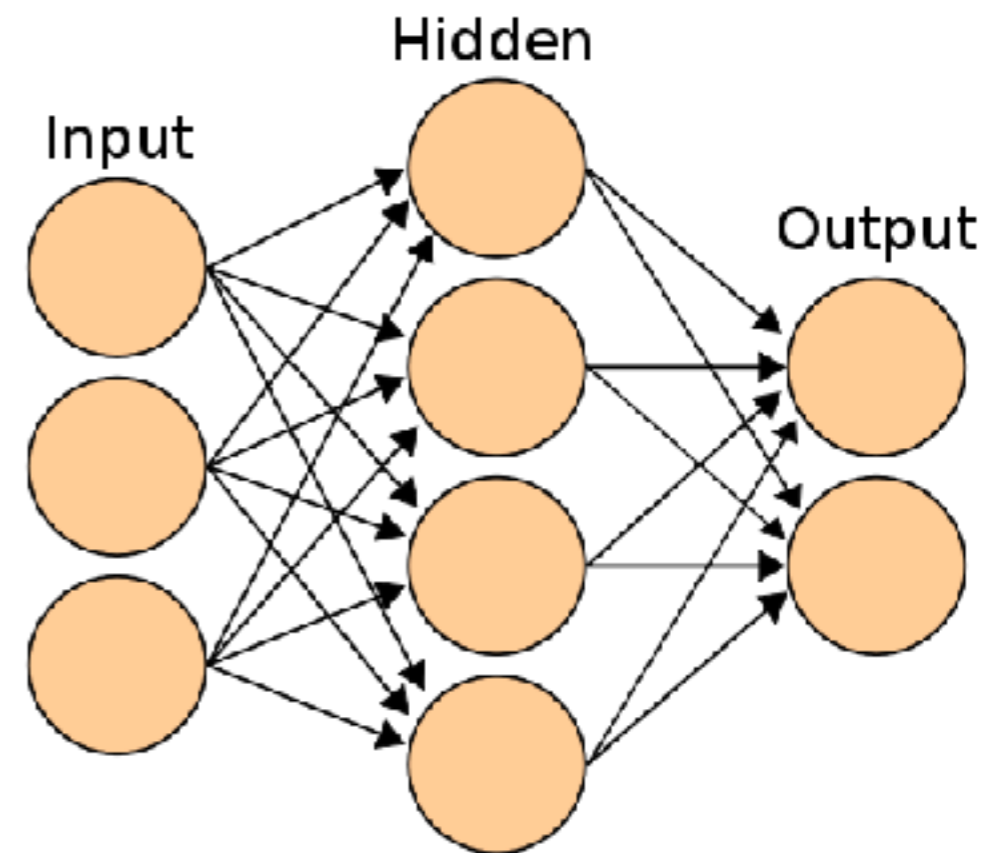
- ▶ A network is **deep** if it has more than one layer of non-linear feature abstraction
- ▶ Hierarchy of representations with increasing levels of abstraction (e.g. pixel -> edge -> eye -> face)
- ▶ Deep networks can store more memory than equivalent number of units in a single layer





# MultiLayer Perceptron

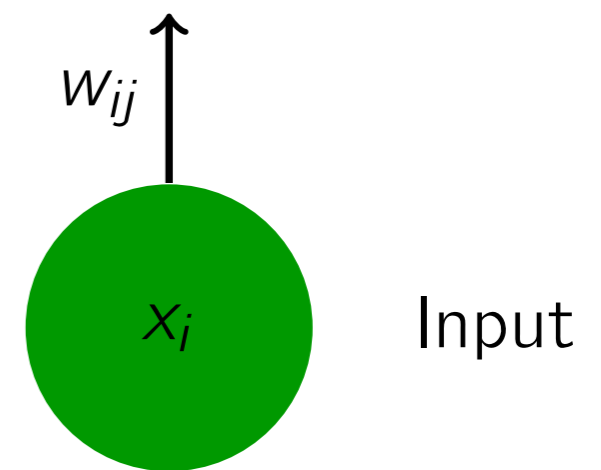
- ▶ Type of **feedforward** artificial neural network
- ▶ Can distinguish data which is not linearly separable
- ▶ Some neurons use non-linear **activation** functions (functions which map the weighted input to the output of a neuron)
- ▶ The brain is thought to work in a similar way when biological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





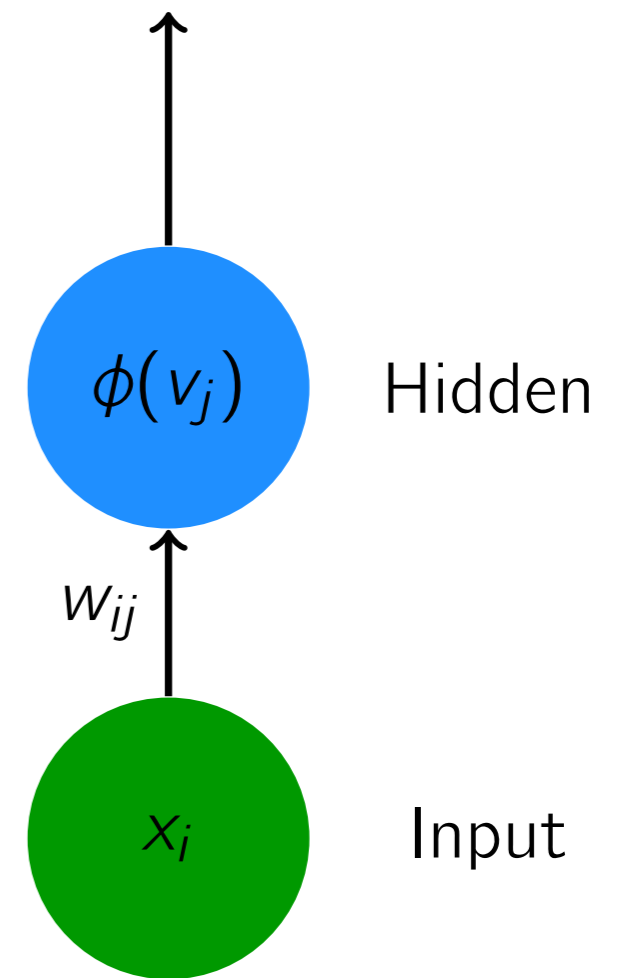
# Back-propagation

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



# Back-propagation

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

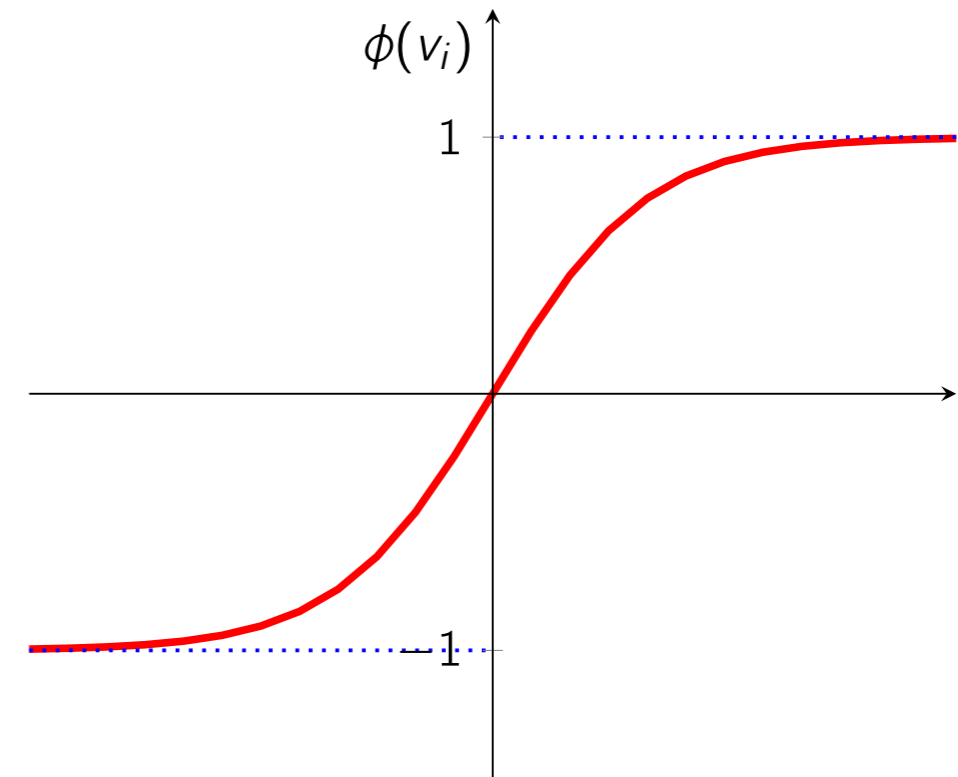




# Back-propagation

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

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



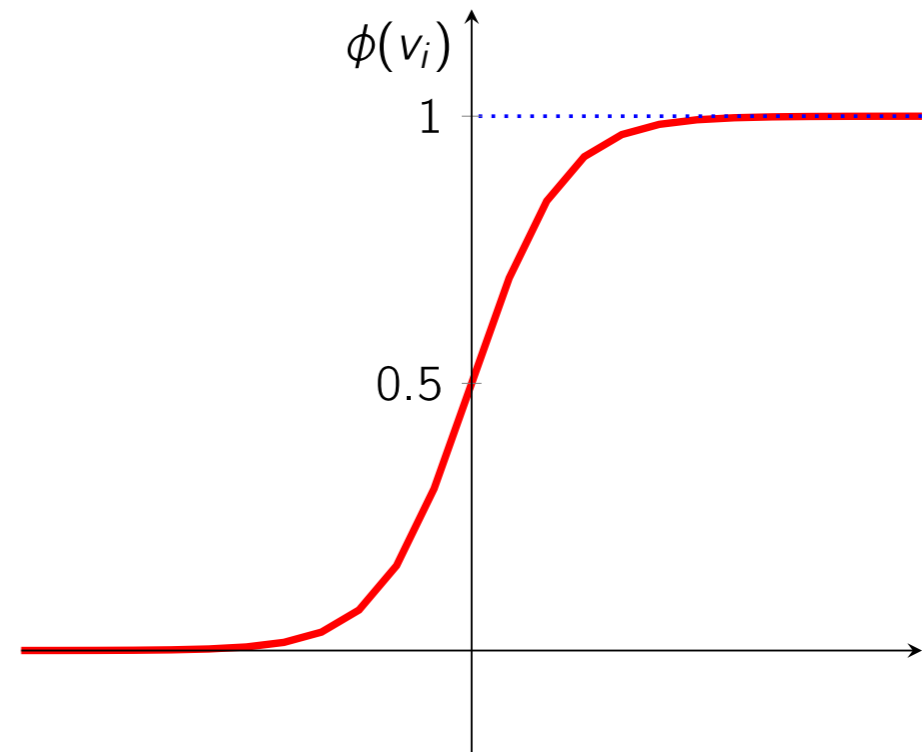
Tanh function

# Back-propagation

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

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

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



Sigmoid function



# Back-propagation

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

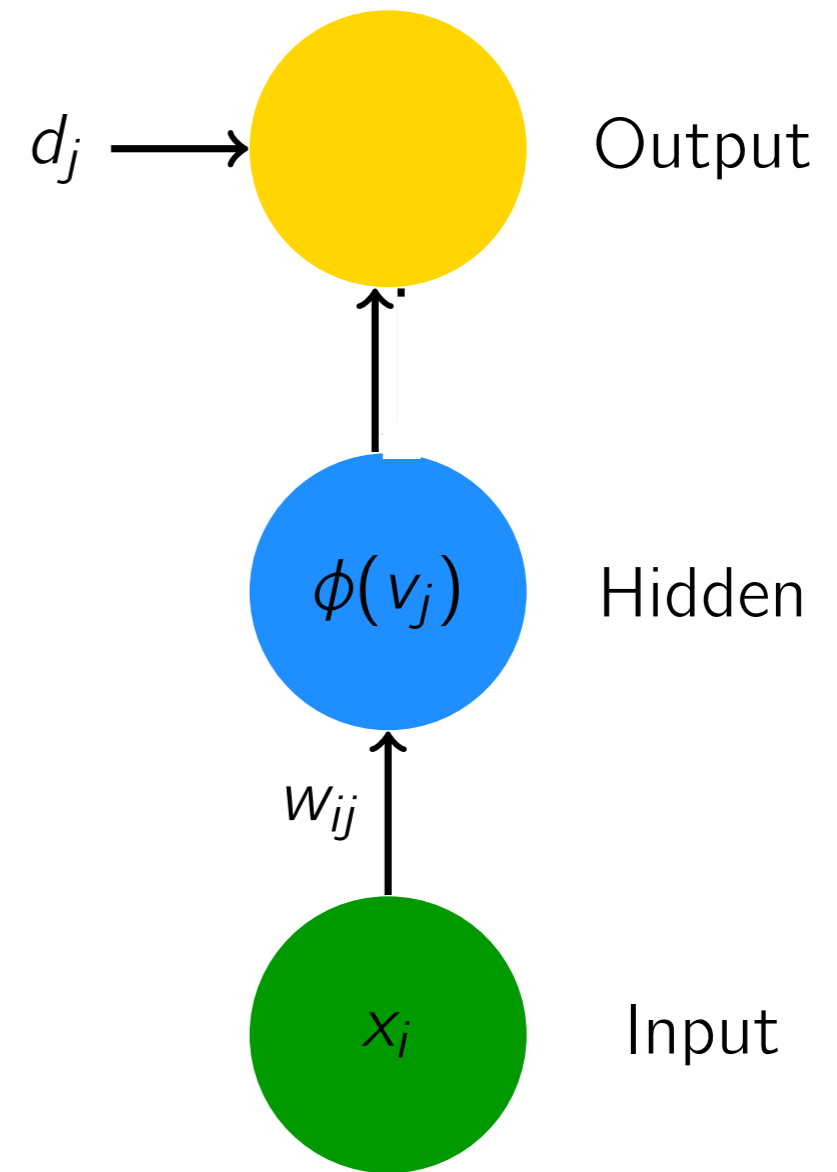
$$\phi(v_i) = \tanh(v_i)$$

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

- ▶ 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

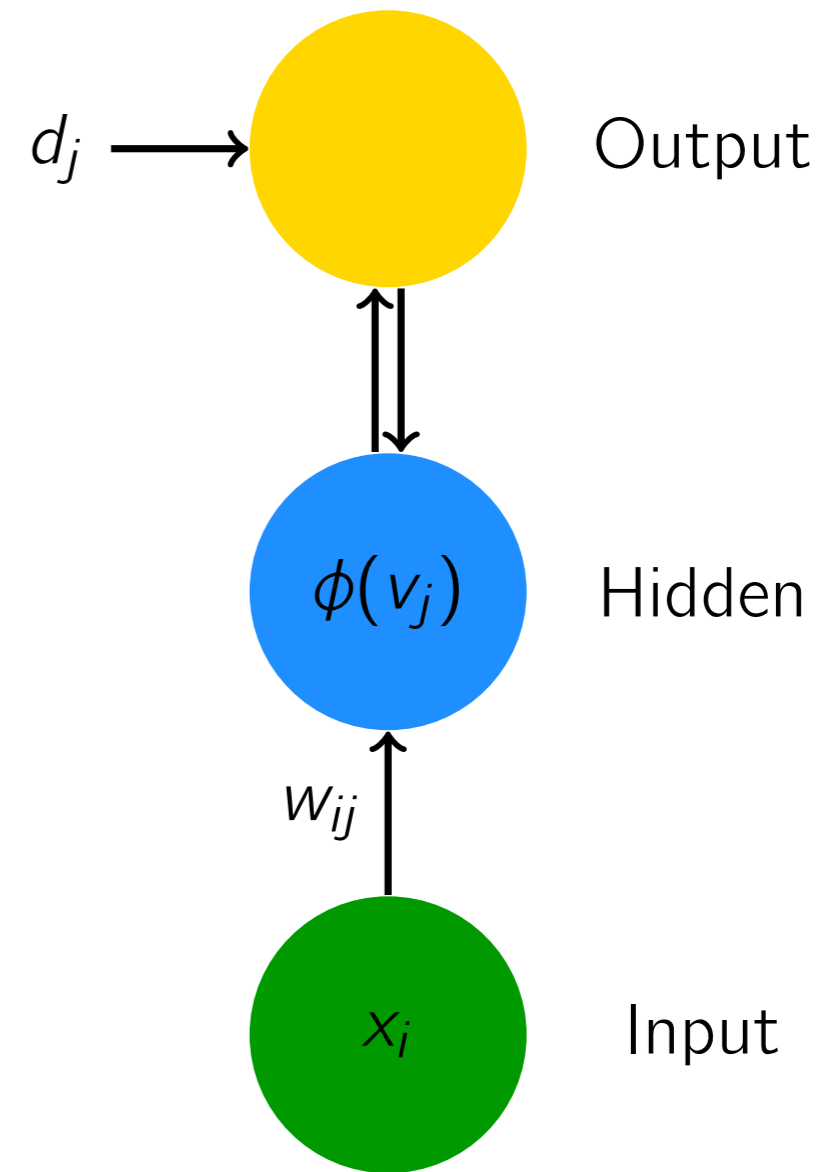
$$e_j = y_j - \phi(v_j)$$

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

- ▶ 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)

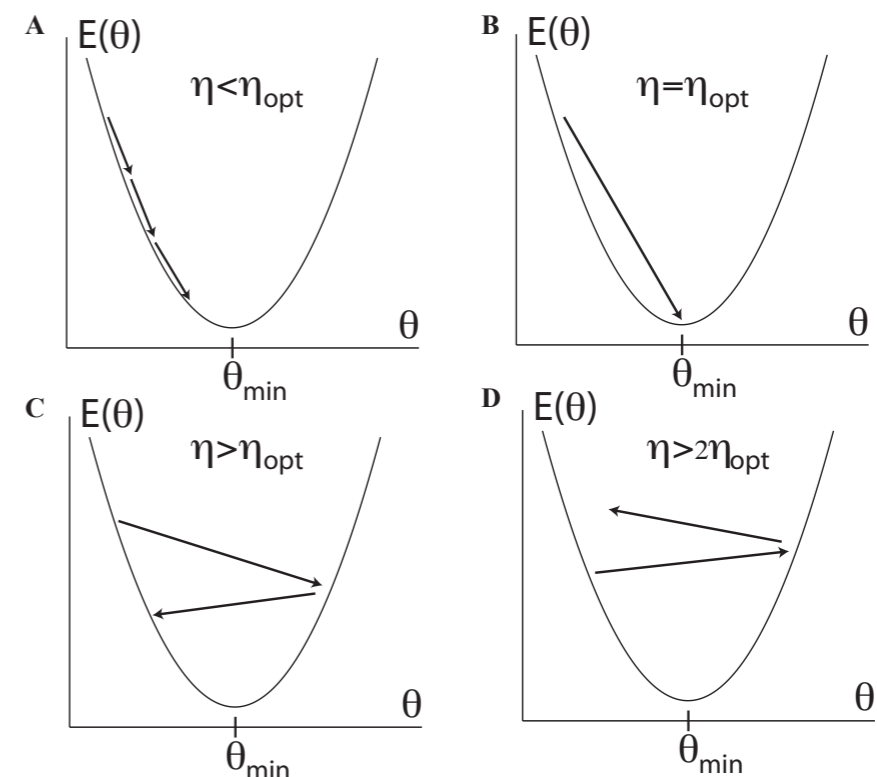
- ▶ Most supervised machine learning algorithms involve finding a set of parameters  $\Theta$  that minimise a cost/loss function
- ▶ Loss functions can be complicated non-convex functions with many local minima
- ▶ Simplest algorithm that attempts to minimise the loss  $E(\Theta)$  is gradient descent

$$\mathbf{v}_t = \eta_t \nabla_{\theta} E(\boldsymbol{\theta}_t),$$

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \mathbf{v}_t$$

$\eta_t$  = learning rate

- ▶ 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

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

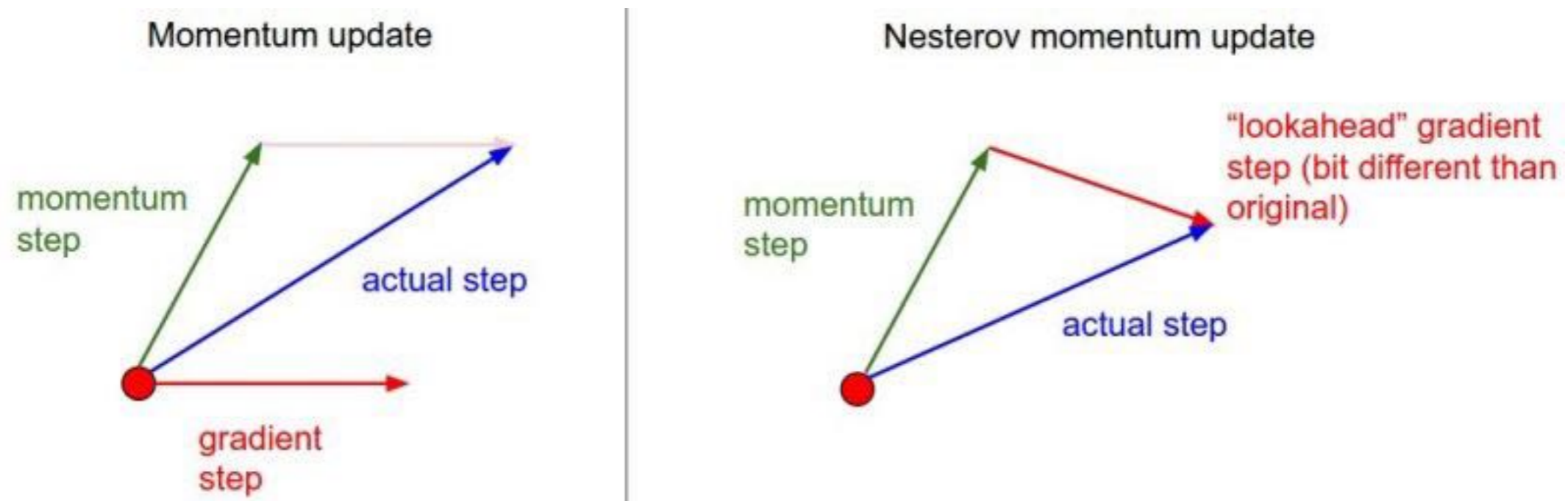
- ▶ Decreases chance of becoming stuck in local minima
- ▶ Also has been shown to help alleviate **over-fitting**

# SGD with momentum

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

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

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





# Second order moments

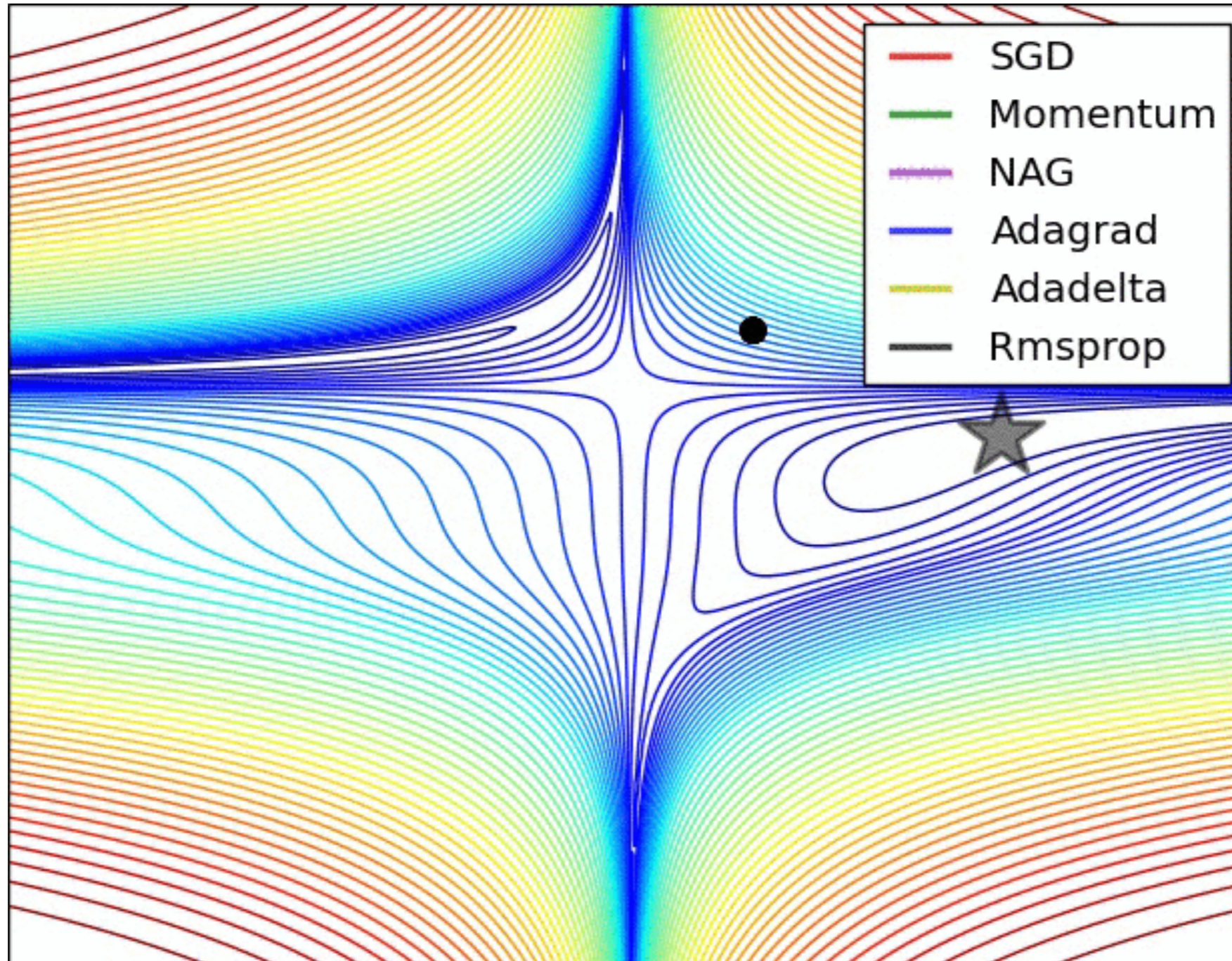
- ▶ All methods so far require a schedule for the learning rate as a function of time
- ▶ Optimal learning rate is actually inverse of the Hessian  $\eta_{\text{opt}} = [\partial_{\theta}^2 E(\theta)]^{-1}$ .
- ▶ Expensive to compute. Second order moment methods keep track of the **squared gradient**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

$$\begin{aligned}\mathbf{g}_t &= \nabla_{\theta} E(\boldsymbol{\theta}) \\ \mathbf{s}_t &= \beta \mathbf{s}_{t-1} + (1 - \beta) \mathbf{g}_t^2 \\ \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t - \eta_t \frac{\mathbf{g}_t}{\sqrt{\mathbf{s}_t + \epsilon}},\end{aligned}$$

$$\begin{aligned}\mathbf{s}_t &= \mathbb{E}[\mathbf{g}_t^2] \\ \beta &= \text{averaging time} \\ \epsilon &= \text{regularisation constant}\end{aligned}$$



# Second order methods





# Deep learning in practice



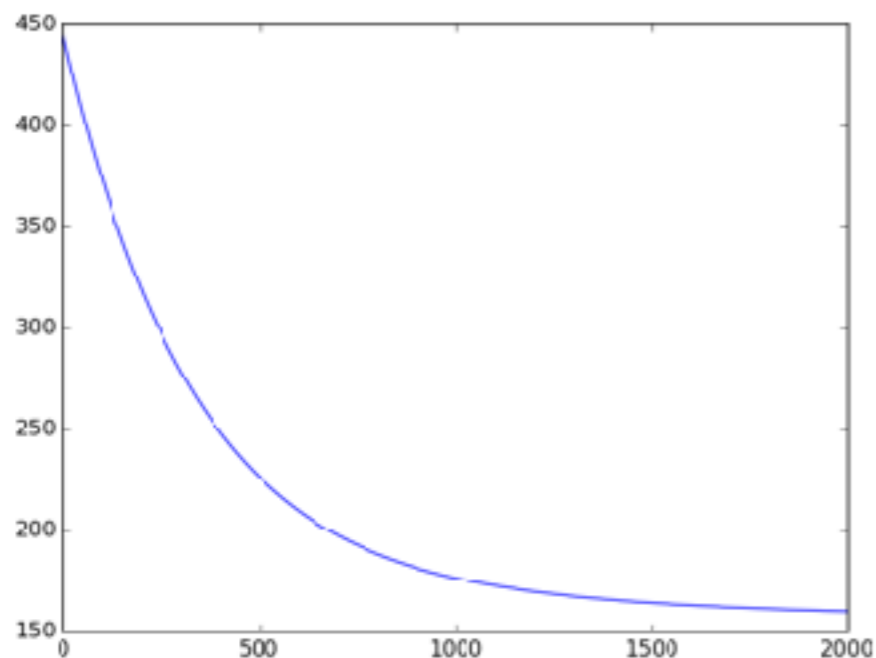
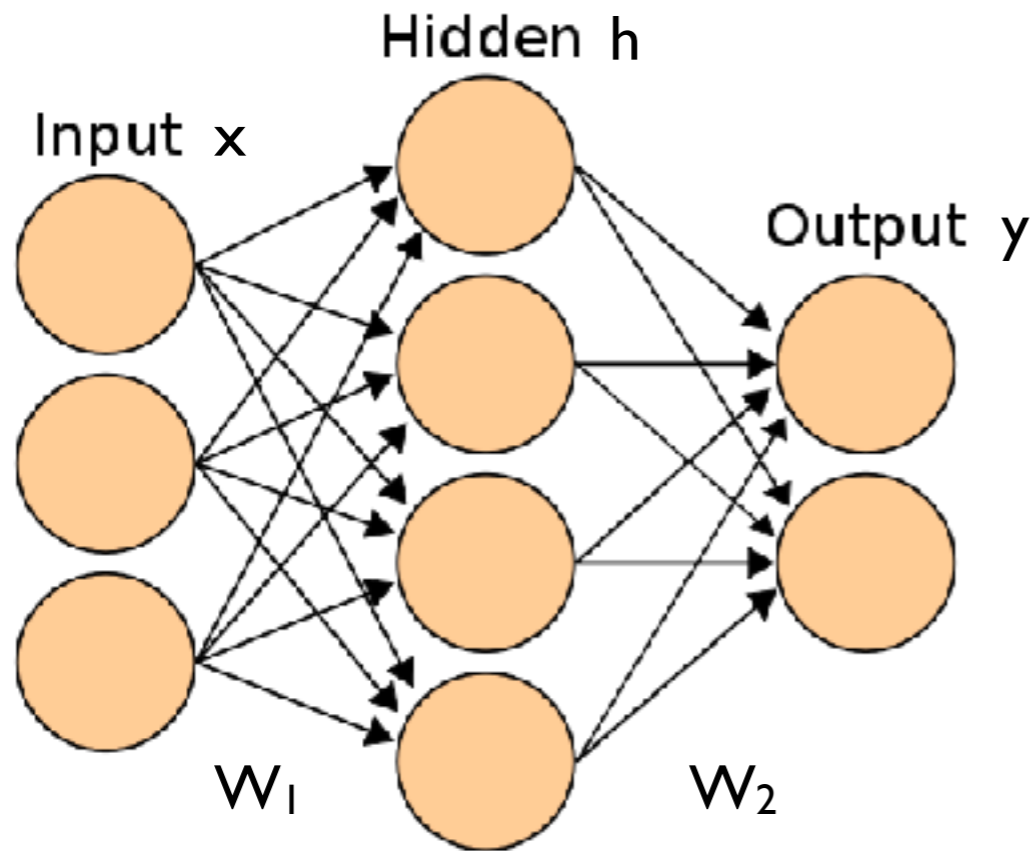
- ▶ Created by Google
- ▶ Math library used for machine learning and neural networks
- ▶ Supports Python and C++
- ▶ Pros
  - ✓ Documentation
  - ✓ Backed by large community
  - ✓ In built monitoring for training processes (Tensorboard)
- ▶ Cons
  - ✗ Static computational graphs
  - ✗ Higher learning 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 x
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()
35
```

# 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

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

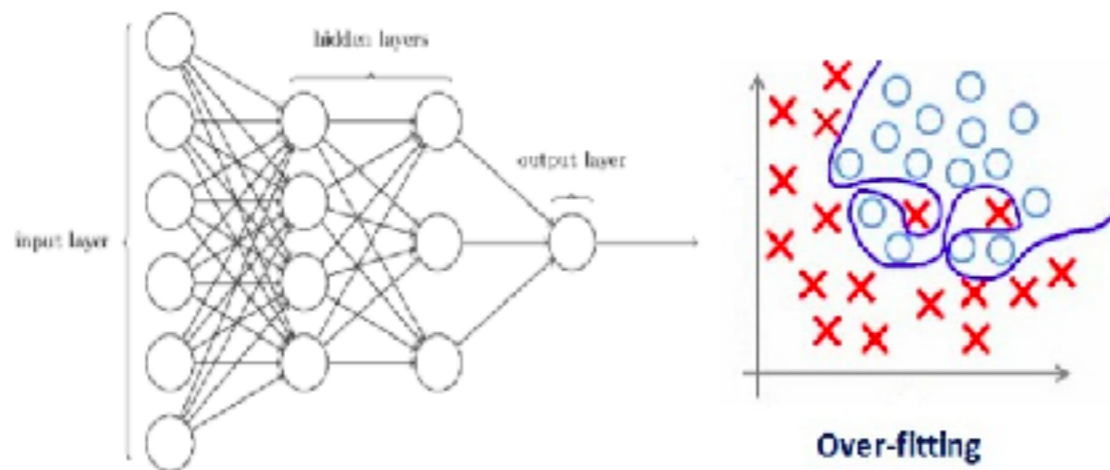
- ▶ Full loss function can include additional regularization terms
- ▶ For **categorical** outputs loss function is usually the categorical cross-entropy
- ▶ Last layer typically has a **soft-max** activation (turns M outputs into normalized probabilities)

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

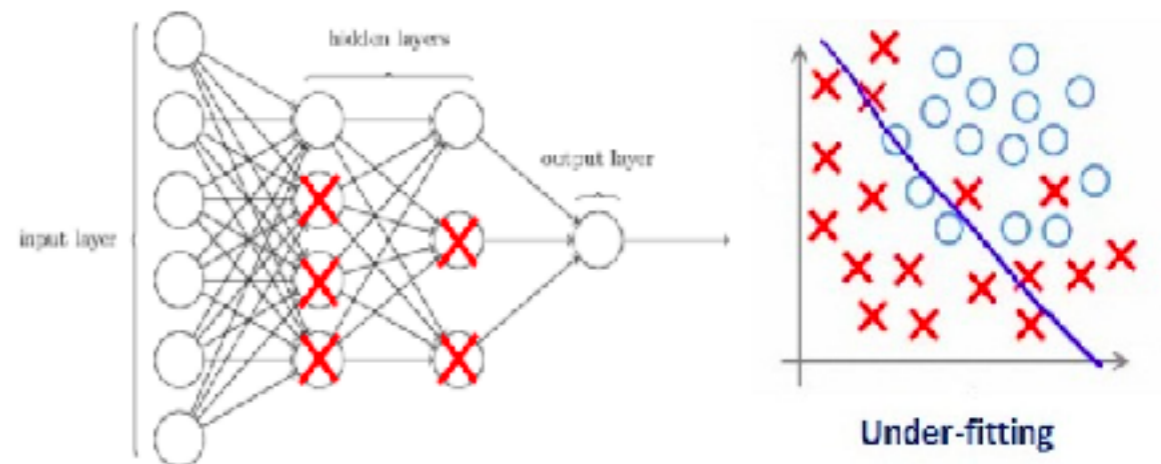
- ▶ Loss increased when predicted probability further from actual label

# 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



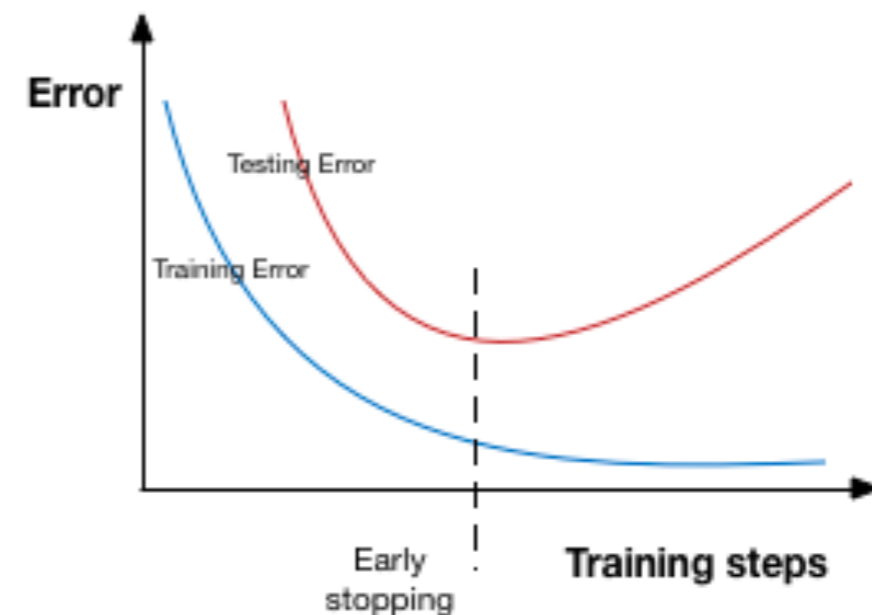
Model has too many free parameters  
and data is over-fitted



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 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 over-fitting



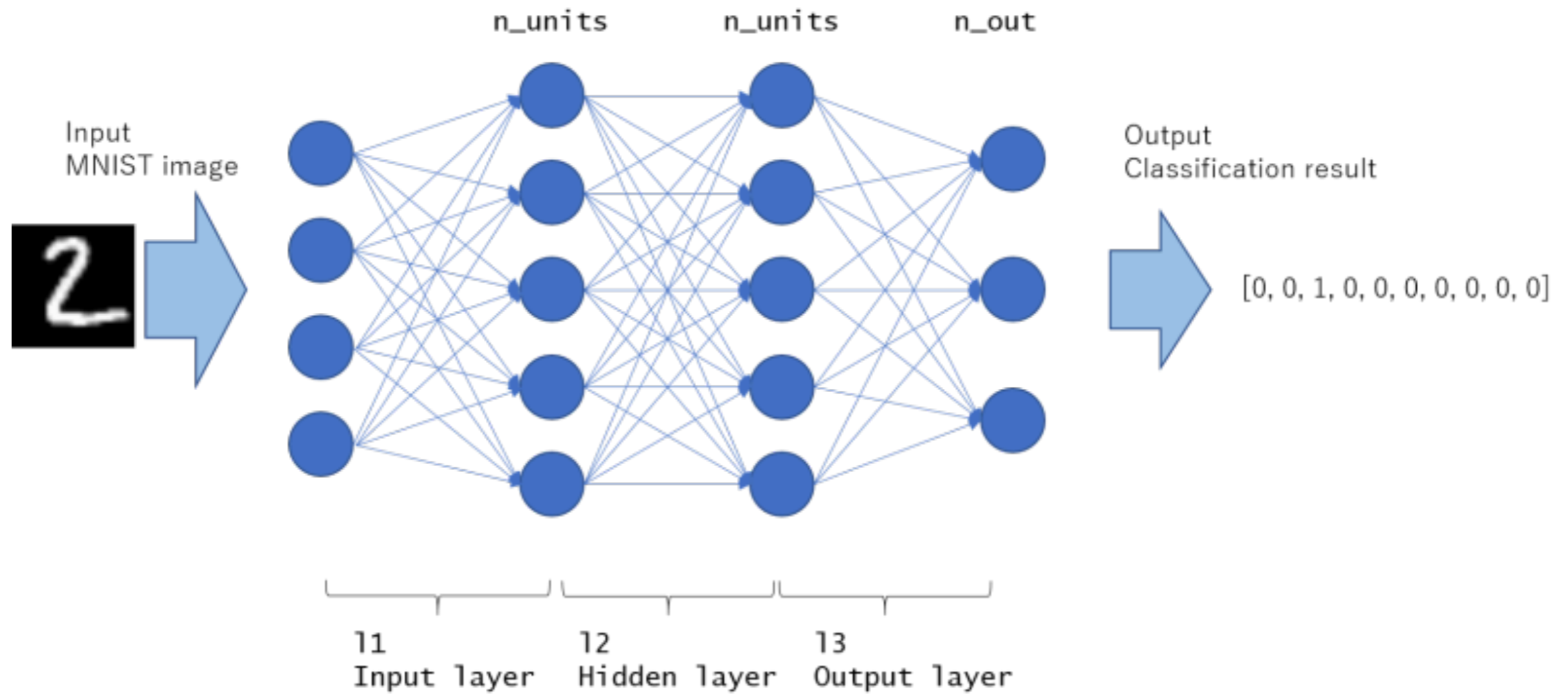


# Hello World

- ▶ 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])
50
```



# MNIST MLP (pyTorch)



```
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, 10)
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         return 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
43 def test(args, model, device, test_loader):
44     model.eval()
45     test_loss = 0
46     correct = 0
47     with torch.no_grad():
48         for data, target in test_loader:
49             data, target = data.to(device), target.to(device)
50             output = model(data)
51             test_loss += F.nll_loss(output, target, reduction='sum').item() # sum up batch loss
52             pred = output.max(1, keepdim=True)[1] # get the index of the max log-probability
53             correct += pred.eq(target.view_as(pred)).sum().item()
54
55     test_loss /= len(test_loader.dataset)
56     print('\nTest set: Average loss: {:.4f}, Accuracy: {}/{} ({:.0f}%)\n'.format(
57         test_loss, correct, len(test_loader.dataset),
58         100. * correct / len(test_loader.dataset)))
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()
```