

## Neural Network (Basic) Cheat Sheet

(source: Deep Learning - Goodfellow, Ian)

V2021.01.01

(Dr Yan Xu)

### Core Assumption

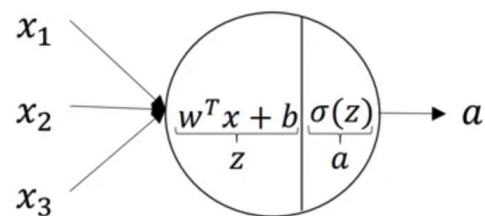
The data (or the target problem) was generated by the composition of factors (or features), potentially at multiple levels in a hierarchy.

### Modelling Mindset

- Models with high representation capability
- Appropriate regularization
- Big training data

### Basic Concepts

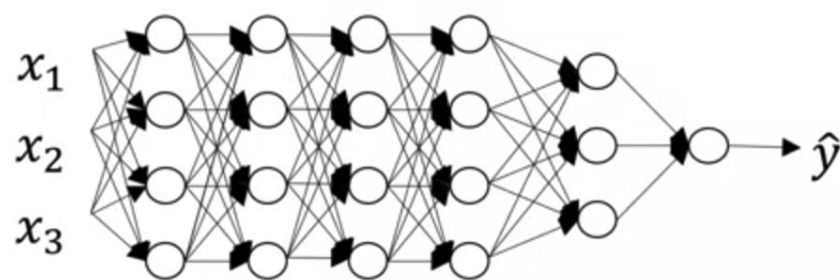
→ Neural:



→ Activation Function:

- Sigmoid
- Tanh
- ReLU (most popular)
- Leaky ReLU

→ Neural Network:



→ Output Units:

- Linear Units (Gaussian Distribution)
- Sigmoid Units (Bernoulli Distribution)
- Softmax Units (Multinoulli Distribution)

### Forward and Backward Propagation

→ Forward Propagation:

```
Require: Network depth,  $l$ 
Require:  $W^{(i)}, i \in \{1, \dots, l\}$ , the weight matrices of the model
Require:  $b^{(i)}, i \in \{1, \dots, l\}$ , the bias parameters of the model
Require:  $x$ , the input to process
Require:  $y$ , the target output
 $h^{(0)} = x$ 
for  $k = 1, \dots, l$  do
   $a^{(k)} = b^{(k)} + W^{(k)} h^{(k-1)}$ 
   $h^{(k)} = f(a^{(k)})$ 
end for
 $\hat{y} = h^{(l)}$ 
 $J = L(\hat{y}, y) + \lambda \Omega(\theta)$ 
```

(note: the symbols  $(a, h)$  are not consistent with other figures)

→ Backward Propagation:

```
After the forward computation, compute the gradient on the output layer:
 $g \leftarrow \nabla_{\hat{y}} J = \nabla_{\hat{y}} L(\hat{y}, y)$ 
for  $k = l, l-1, \dots, 1$  do
  Convert the gradient on the layer's output into a gradient into the pre-
  nonlinearity activation (element-wise multiplication if  $f$  is element-wise):
   $g \leftarrow \nabla_{a^{(k)}} J = g \odot f'(a^{(k)})$ 
  Compute gradients on weights and biases (including the regularization term,
  where needed):
   $\nabla_{b^{(k)}} J = g + \lambda \nabla_{b^{(k)}} \Omega(\theta)$ 
   $\nabla_{W^{(k)}} J = g h^{(k-1)\top} + \lambda \nabla_{W^{(k)}} \Omega(\theta)$ 
  Propagate the gradients w.r.t. the next lower-level hidden layer's activations:
   $g \leftarrow \nabla_{h^{(k-1)}} J = W^{(k)\top} g$ 
end for
```

(source: Deep Learning - Goodfellow, Ian)

### Regularization

→ Parameter Norm

- L1 and L2 norm
- Typically, biases are not regularized

→ Dataset Augmentation

- random rotation, shift, reflection
- random noise

→ EarlyStopping

- must have validation set

→ Dropout (How it works?)

- Motivation: "**Bagging + Parameter Sharing**"
- used after the activation function
- *Training*: random selection & weight adjustment
- *Prediction*: standard forward propagation

→ Dropout (some notes)

- could be used to estimate prediction uncertainty
- could be used in the input layer and hidden layers
- typically performs better than Norm
- typically used on fully connected layers (e.g.  $p = 0.5$ )

### Optimisation

→ Local Minima

- the cost function of neural networks is non-convex
- however, local minima seems not a major problem

→ Second-order (or Newton Method)

- remain difficult to scale to large neural networks
- gradient-based method is still the mainstream

→ Mini-batch Optimisation

- typical values: 32 - 256

→ Stochastic Gradient Descent

→ Batch Normalization

- Training: mini-batch normalization in hidden layers
- Prediction: using learned mean & variance

→ Note: Batch Normalization is not a regularization method

→ Optimisation Method with Adaptive Learning Rates

- RMSProp

**Algorithm 8.5** The RMSProp algorithm

**Require:** Global learning rate  $\epsilon$ , decay rate  $\rho$ .

**Require:** Initial parameter  $\theta$

**Require:** Small constant  $\delta$ , usually  $10^{-6}$ , used to stabilize division by small numbers.

Initialize accumulation variables  $r = 0$

**while** stopping criterion not met **do**

Sample a minibatch of  $m$  examples from the training set  $\{x^{(1)}, \dots, x^{(m)}\}$  with corresponding targets  $y^{(i)}$ .

Compute gradient:  $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$

Accumulate squared gradient:  $r \leftarrow \rho r + (1 - \rho) g \odot g$

Compute parameter update:  $\Delta \theta = -\frac{\epsilon}{\sqrt{\delta + r}} \odot g$ . ( $\frac{1}{\sqrt{\delta + r}}$  applied element-wise)

Apply update:  $\theta \leftarrow \theta + \Delta \theta$

**end while**

- Adam

**Algorithm 8.7** The Adam algorithm

**Require:** Step size  $\epsilon$  (Suggested default: 0.001)

**Require:** Exponential decay rates for moment estimates,  $\rho_1$  and  $\rho_2$  in  $[0, 1]$ . (Suggested defaults: 0.9 and 0.999 respectively)

**Require:** Small constant  $\delta$  used for numerical stabilization. (Suggested default:  $10^{-8}$ )

**Require:** Initial parameters  $\theta$

Initialize 1st and 2nd moment variables  $s = 0, r = 0$

Initialize time step  $t = 0$

**while** stopping criterion not met **do**

Sample a minibatch of  $m$  examples from the training set  $\{x^{(1)}, \dots, x^{(m)}\}$  with corresponding targets  $y^{(i)}$ .

Compute gradient:  $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$

$t \leftarrow t + 1$

Update biased first moment estimate:  $s \leftarrow \rho_1 s + (1 - \rho_1) g$

Update biased second moment estimate:  $r \leftarrow \rho_2 r + (1 - \rho_2) g \odot g$

Correct bias in first moment:  $\hat{s} \leftarrow \frac{s}{1 - \rho_1^t}$

Correct bias in second moment:  $\hat{r} \leftarrow \frac{r}{1 - \rho_2^t}$

Compute update:  $\Delta \theta = -\epsilon \frac{\hat{s}}{\sqrt{\hat{r} + \delta}}$  (operations applied element-wise)

Apply update:  $\theta \leftarrow \theta + \Delta \theta$

**end while**

### Model Tuning

→ Network, Regularization, Optimisation parameters

→ Cross Validation could be very slow