# Train/Dev/Test set
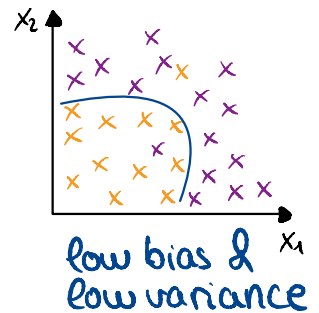
Data

| training set (train) | development set (dev) | test set |
| --- | --- | --- |

Previous era → 60/20/20% ~ #data ~100,000 → 60,000/20,000/20,000

Big data → 98/1/1% ~ #data ~1,000,000 → 980,000/10,000/10,000

# Basic recipe ~ bias/variance tradeoff

High bias (training data problem) —— YES →  bigger network < ↑#hidden layers / ↑#hidden units
train longer
NN architecture search*

↓ NO

High variance (dev set problem) —— YES → more data
regularization
NN architecture search*

↓ NO

Done

*Not allways work



| high bias | high variance | high bias & high variance | low bias & low variance |

|  | | | |
| --- | --- | --- | --- |
| Train set error : | 1% | 15% | 15% | 0.5% |
| Dev set error : | 11% | 16% | 30% | 1% |
|  | high bias | high variance | high bias & high variance | low bias & low variance |

# Deep L-layer Neural Network

## 4-layer Neural Network
layer:     0    1    2    3    4



layer 0 ~ input layer
layer 1-3 ~ hidden layer
layer 4 ~ output layer

shallow NN ~ 1 hidden layer
deep NN ~ > 1 hidden layer

elements:  5    7    7    7    1

$L = 4$ ~ # layers
$n^{[\ell]}$ ~ # nodes/elements in layer $\ell$    $n^{[0]} = 5$   $n^{[1]} = n^{[2]} = n^{[3]} = 7$   $n^{[4]} = 1$
$a^{[\ell]}$ ~ activations in layer $\ell$   $a^{[0]} = X$   $a^{[L]} = a^{[4]} = y$

## Parameters

$W^{[\ell]} \in \mathbb{R}^{(n^{[\ell]}, n^{[\ell-1]})}$ → initialize to: random * 0.01
$b^{[\ell]} \in \mathbb{R}^{(n^{[\ell]}, 1)}$ → initialize to: zeros

*break symmetry (no zeros!!!)*
*lower value → faster convergence*

## Hyperparameters
activation function
$\boxed{\alpha}$ ~ learning rate    *0.9*
$\boxed{\beta}$ ~ optimization param.
$\boxed{\beta_1, \beta_2, \varepsilon}$ ~ Adam optimization params
*0.9  0.999  $10^{-8}$*

# iterations
$\boxed{\text{\# layers}}$
$\boxed{\text{\# hidden units}}$
$\boxed{\text{learning rate decay}}$
$\boxed{\text{mini-batch size}}$

## Hyperparameter tuning *- increase convergence*
$\square$ Most important
$\square$ Second in importance
$\square$ Third in importance
$\square$ Never tuned

# Neural Networks method
## Repeat until convergence

## 1. Forward propagation

$$Z^{[\ell]} = W^{[\ell]} A^{[\ell-1]} + b^{[\ell]}$$ ← broadcasting → $(m,n) \overset{\pm}{*} \begin{matrix}(m,1)\\ \text{or}\\ (1,n)\end{matrix} = (m,n)$

$$A^{[\ell]} = g^{[\ell]}(Z^{[\ell]})$$

## 2. When $A^{[L]} \sim \hat{y}$ obtained:

$$dZ^{[L]} = L(A^{[L]}, y)$$ ~ RMSE, MAE ...
$$dW^{[L]} = 1/m \; dZ^{[L]} A^{[L-1]T}$$
$$db^{[L]} = 1/m \sum dZ^{[L]}$$
↖ columns

## 3. Backward propagation ~ chain rule

$$dZ^{[\ell]} = W^{[\ell+1]} dZ^{[\ell+1]} \overset{\text{element-wise}}{*} g^{[\ell]\prime}(Z^{[\ell]})$$
$$dW^{[\ell]} = 1/m \; dZ^{[\ell]} A^{[\ell-1]}$$
$$db^{[\ell]} = 1/m \sum dZ^{[\ell]}$$

## 4. Parameters update

$$W^{[\ell]} = W^{[\ell]} - \alpha \, dW^{[\ell]}$$
$$b^{[\ell]} = b^{[\ell]} - \alpha \, db^{[\ell]}$$

# Deep Neural Networks in blocks



layer $\ell$

forward

backward

$A^{[\ell-1]}$ → $W^{[\ell]}, b^{[\ell]}$ → $A^{[\ell]}$

cache $Z^{[\ell]}$

$dA^{[\ell-1]}$ ← $\begin{matrix} W^{[\ell]}, b^{[\ell]} \\ dZ^{[\ell]} \end{matrix}$ ← $dA^{[\ell]}$

$dW^{[\ell]}, db^{[\ell]}$

# Most used activation functions
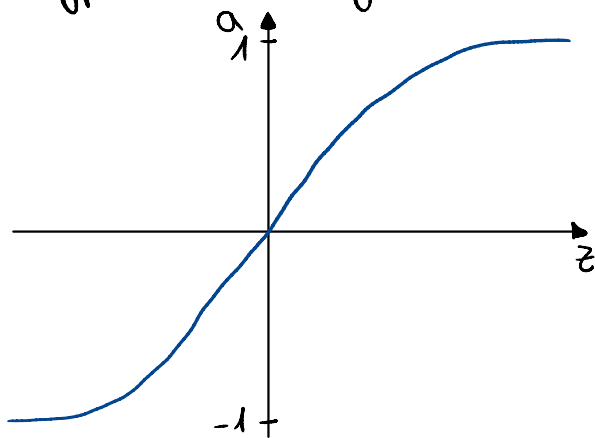## Can be different for different layers

### Sigmoid



$$g(z) = \frac{1}{1+e^{-z}}$$
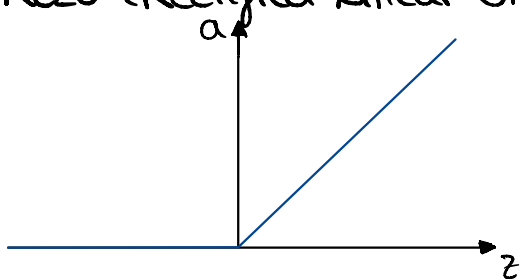
$$g'(z) = a(1-a)$$

### Hyperbolic tangent



$$g(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$
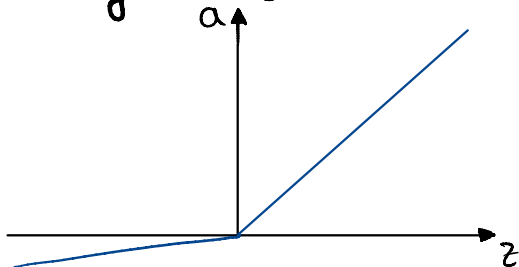
$$g'(z) = 1 - a^2$$

### ReLU (Rectified Linear Unit)



$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

### Leaky ReLU



$$g(z) = \max(0.01z, z)$$

$$g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

# Optimization algorithms
## Mini-batch ~ divisions of the training set

$$X = \left[ x^{(1)}, x^{(2)}, \ldots, x^{(1000)} \mid x^{(1001)}, \ldots, x^{(2000)} \mid \ldots, x^{(m)} \right]$$
$$\underbrace{\hspace{3cm}}_{X^{\{1\}}} \quad \underbrace{\hspace{3cm}}_{X^{\{2\}}} \quad \underbrace{\hspace{2cm}}_{X^{\{5000\}}}$$

$X \in \mathbb{R}^{(n_x, m)}$

$n^{[0]} \sim$ # variables

# training examples

$$y = \left[ y^{(1)}, y^{(2)}, \ldots, y^{(1000)} \mid y^{(1001)}, \ldots, y^{(2000)} \mid \ldots, y^{(m)} \right]$$
$$\underbrace{\hspace{3cm}}_{y^{\{1\}}} \quad \underbrace{\hspace{3cm}}_{y^{\{2\}}} \quad \underbrace{\hspace{2cm}}_{y^{\{5000\}}}$$
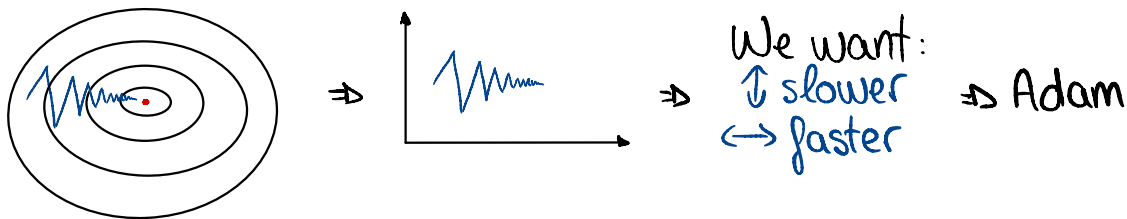
$y \in \mathbb{R}^{(1, m)}$

if 1 output

If $m = 5,000,000$ — mini-batches of 1,000 each

Typical mini-batch sizes: $64, 128, 256, 512 \cong 2^n$

If mini-batch size = m : Batch gradient descent
If mini-batch size = 1 : Stochastic gradient descent

$\Big\}$ In practice somewhere between both

# Adam optimization
## Convergence:



$\Rightarrow$ We want:
$\updownarrow$ slower
$\longleftrightarrow$ faster
$\Rightarrow$ Adam

# Learning rate decay (most used method)

$$\alpha = \frac{1}{1 + \text{decayrate} * \#\text{epochs}} \alpha_0$$

1 epoch = 1 pass through data