# Neural networks basics

$x_1$

$x_2$ → output
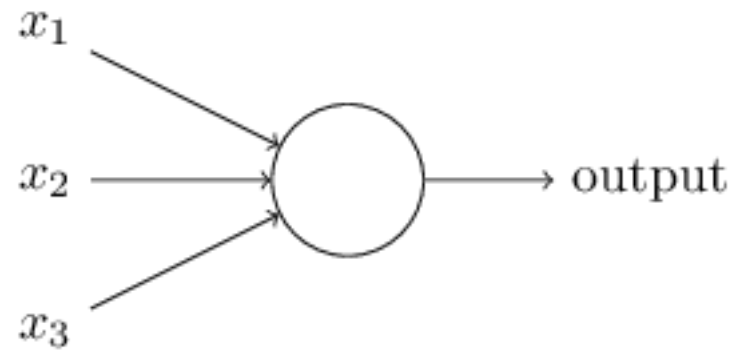
$x_3$

(almost) all images taken from:

**_Neural Networks and Deep Learning_**

Michael Nielsen

http://neuralnetworksanddeeplearning.com/

# Neural networks basics

# Neural networks basics

$x_1$

w1

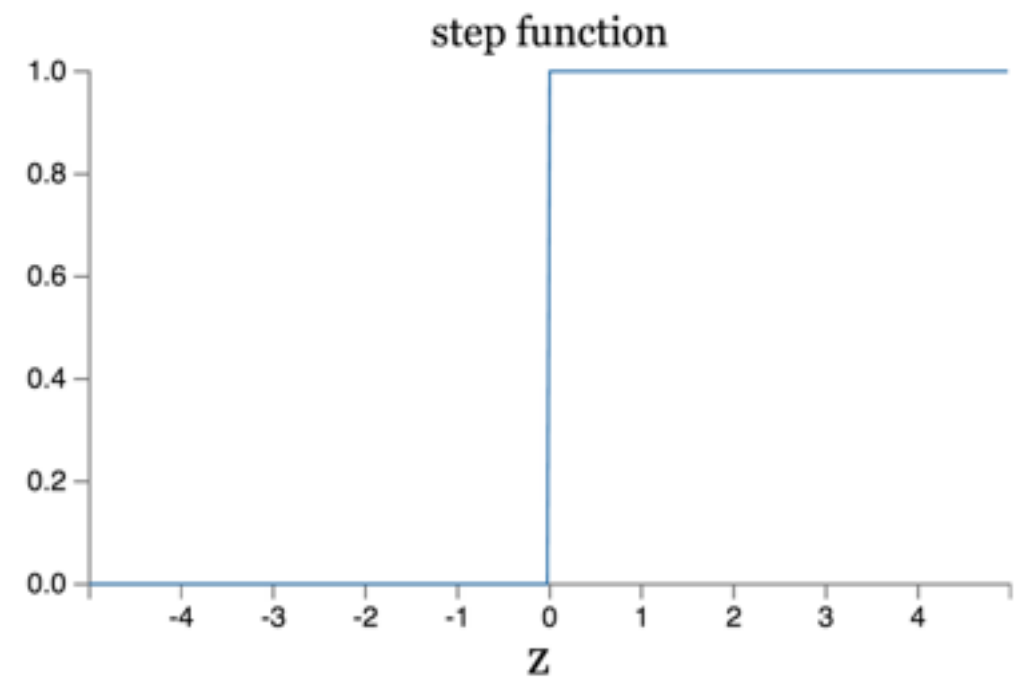$x_2$ —— w2 —→ $\;b\;$ —→ output

$x_3$  w3

# Neural networks basics



Perceptron

$$\text{output} = \begin{cases} 0 & \text{if } w \cdot x + b \leq 0 \\ 1 & \text{if } w \cdot x + b > 0 \end{cases}$$

step function

$x_1$

$w_1$

$x_2$ $\xrightarrow{\phantom{w_2}}$ $w_2$ $\quad$ $b$ $\longrightarrow$ output

$x_3$ $w_3$

output

### Perceptron

$$\text{output} = \begin{cases} 0 & \text{if } w \cdot x + b \leq 0 \\ 1 & \text{if } w \cdot x + b > 0 \end{cases}$$

**step function**

1.0

0.8

0.6

0.4

0.2

0.0

-4   -3   -2   -1   0   1   2   3   4

**z**

# Neural networks basics



$x_1$
$w_1$
$x_2$ $w_2$ $b$ → output
$x_3$ $w_3$

## Sigmoid

$$\sigma(z) \equiv \frac{1}{1 + e^{-z}} \qquad z \equiv w \cdot x + b$$

sigmoid function

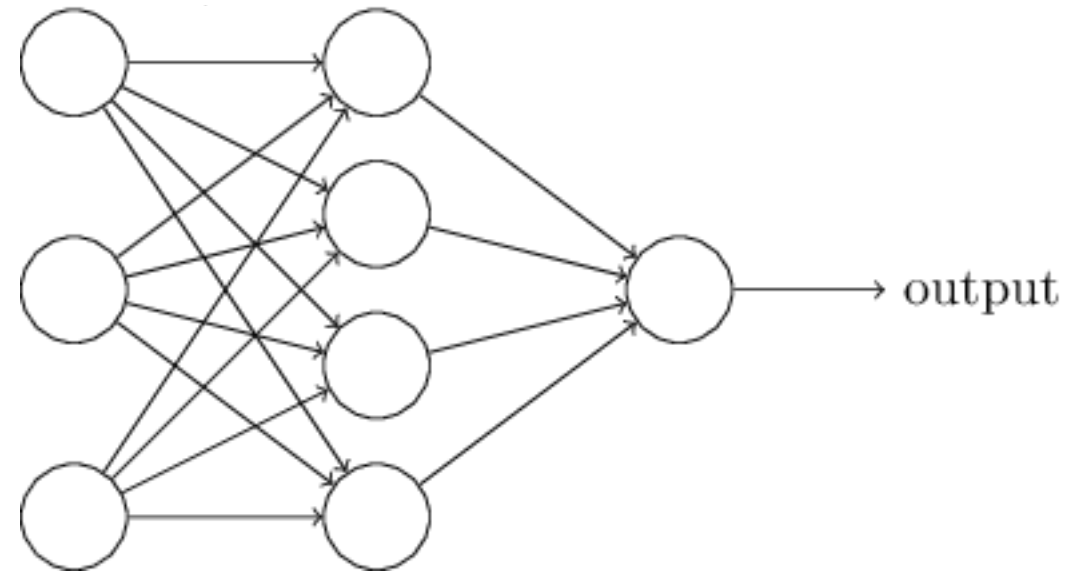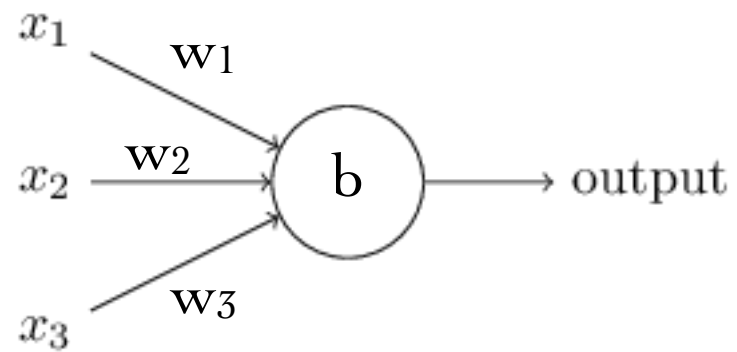## Perceptron

$$\text{output} = \begin{cases} 0 & \text{if } w \cdot x + b \leq 0 \\ 1 & \text{if } w \cdot x + b > 0 \end{cases}$$

step function
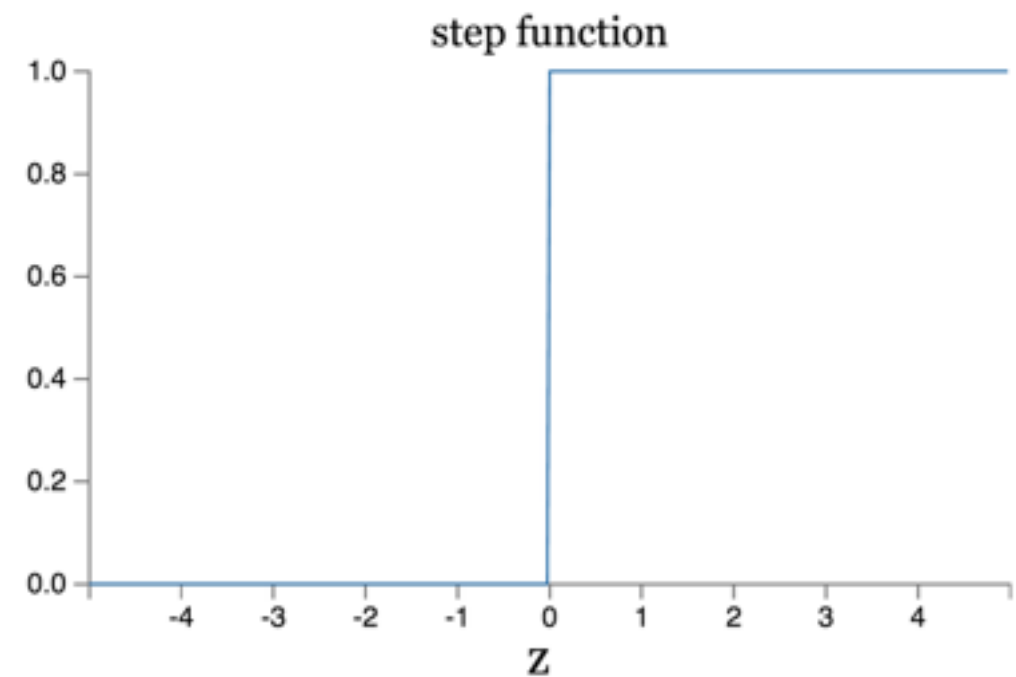
# Neural networks basics

Fully connected

hidden layer
(n = 15 neurons)

output layer

input layer
(784 neurons)

0
1
2
3
4
5
6
7
8
9

output
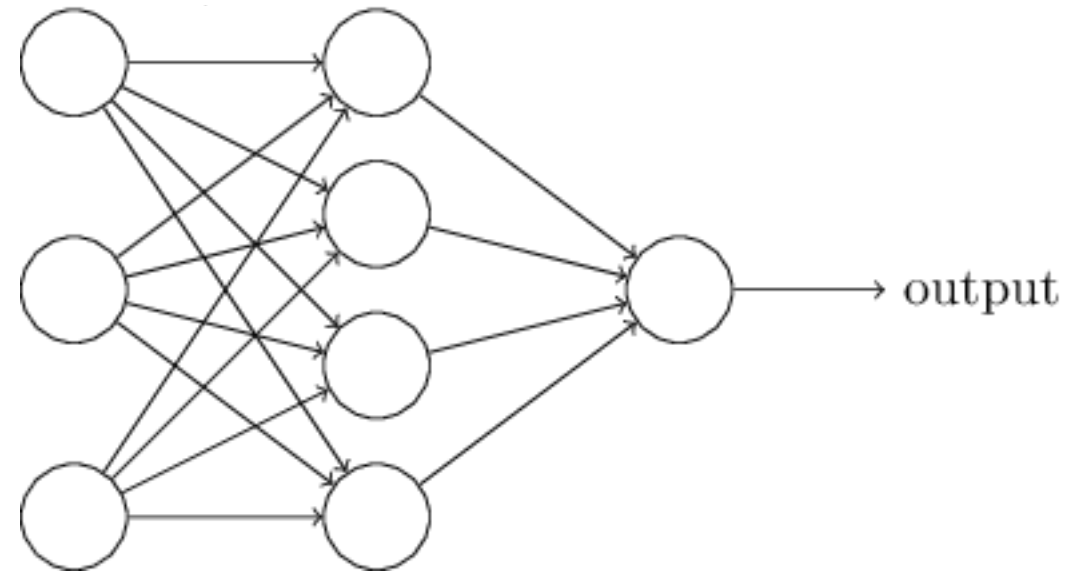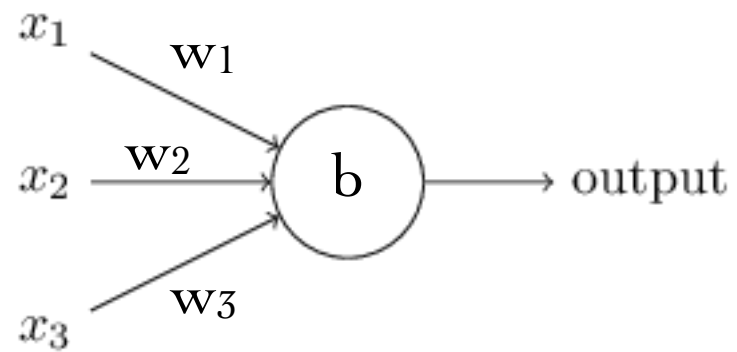
Perceptron

$$\text{output} = \begin{cases} 0 & \text{if } w \cdot x + b \leq 0 \\ 1 & \text{if } w \cdot x + b > 0 \end{cases}$$

Sigmoid

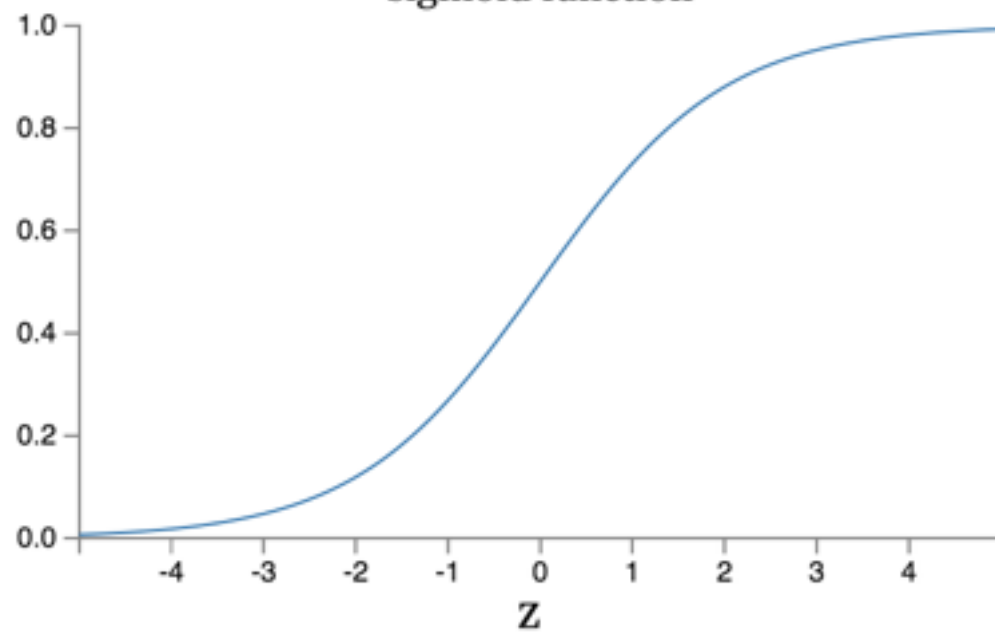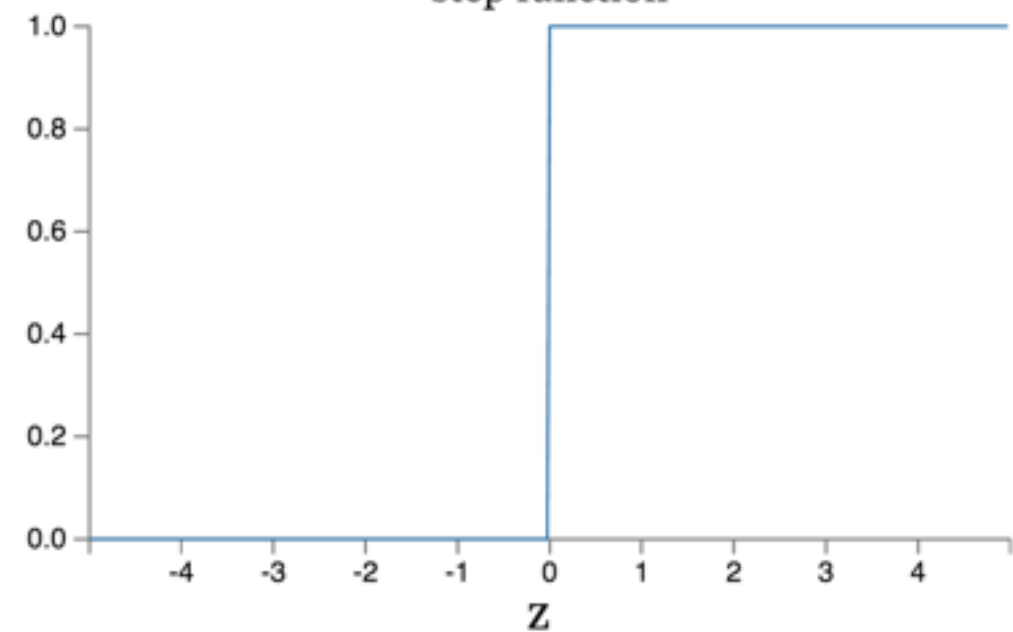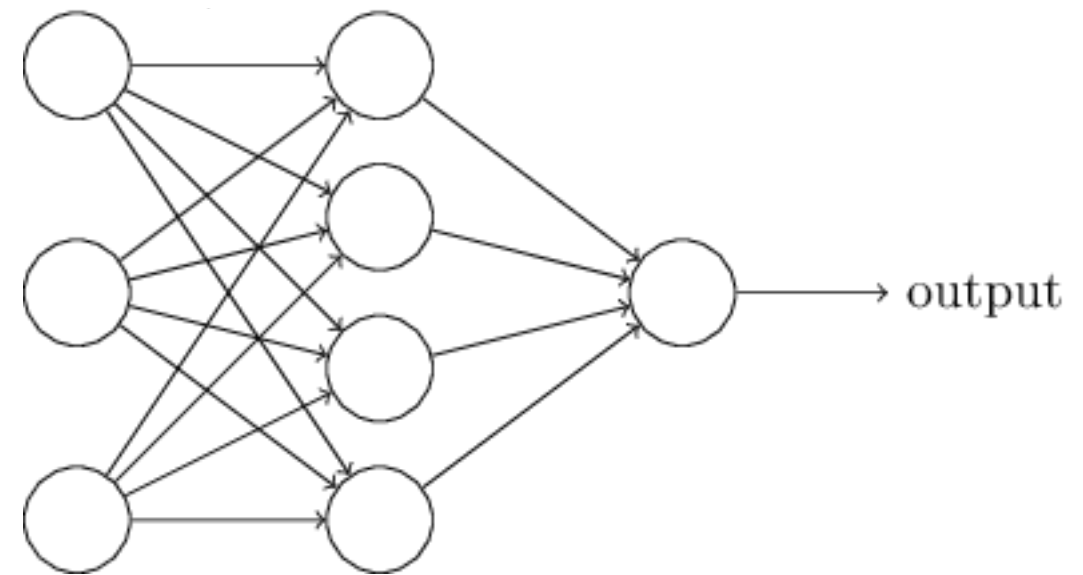$$\sigma(z) \equiv \frac{1}{1 + e^{-z}} \qquad z \equiv w \cdot x + b$$

sigmoid function

step function

# Neural networks basics



Fully connected

output

old output layer

hidden layer

new output layer

input layer
(784 neurons)

Fully connected (deep-ish)

# Training a Neural Net



small change in any weight (or bias)
causes a small change in the output

$w + \Delta w$

output+$\Delta$output

# Training a Neural Net

small change in any weight (or bias)

causes a small change in the output

$w + \Delta w$

output$+\Delta$output

$$\Delta\text{output} \approx \sum_j \frac{\partial\,\text{output}}{\partial w_j} \Delta w_j + \frac{\partial\,\text{output}}{\partial b} \Delta b$$

# Training a Neural Net

$$C(w, b) \equiv \frac{1}{2n} \sum_{x} \|y(x) - a\|^2$$

$x_1$

$w_1$

$w_2$

$x_2 \longrightarrow \quad b \longrightarrow a = \sigma(z)$

$w_3$

$x_3$

small change in any weight (or bias)

causes a small change in the output

$w + \Delta w$

$\longrightarrow$ output$+\Delta$output

$$\Delta\text{output} \approx \sum_{j} \frac{\partial\,\text{output}}{\partial w_j} \Delta w_j + \frac{\partial\,\text{output}}{\partial b} \Delta b$$

# Training a Neural Net

$$C(w, b) \equiv \frac{1}{2n} \sum_{x} \|y(x) - a\|^2$$

training examples
(mini-batches)

$x_1$
$w_1$
$w_2$
$x_2$
$w_3$
$b$
$a = \sigma(z)$
$x_3$

small change in any weight (or bias)
causes a small change in the output

$w + \Delta w$

output$+\Delta$output

$$\Delta\text{output} \approx \sum_{j} \frac{\partial\, \text{output}}{\partial w_j} \Delta w_j + \frac{\partial\, \text{output}}{\partial b} \Delta b$$

$$C(w, b) \equiv \frac{1}{2n} \sum_{x} \|y(x) - a\|^2$$

training examples
(mini-batches)

$x_1$
$w_1$

$w_2$
$x_2$ $b$ $\rightarrow a = \sigma(z)$
$w_3$

$x_3$

$$w_k \rightarrow w_k' = w_k - \eta \frac{\partial C}{\partial w_k}$$

$$b_l \rightarrow b_l' = b_l - \eta \frac{\partial C}{\partial b_l}.$$

small change in any weight (or bias)
causes a small change in the output
$w + \Delta w$

output$+\Delta$output

$$\Delta \text{output} \approx \sum_{j} \frac{\partial \, \text{output}}{\partial w_j} \Delta w_j + \frac{\partial \, \text{output}}{\partial b} \Delta b$$
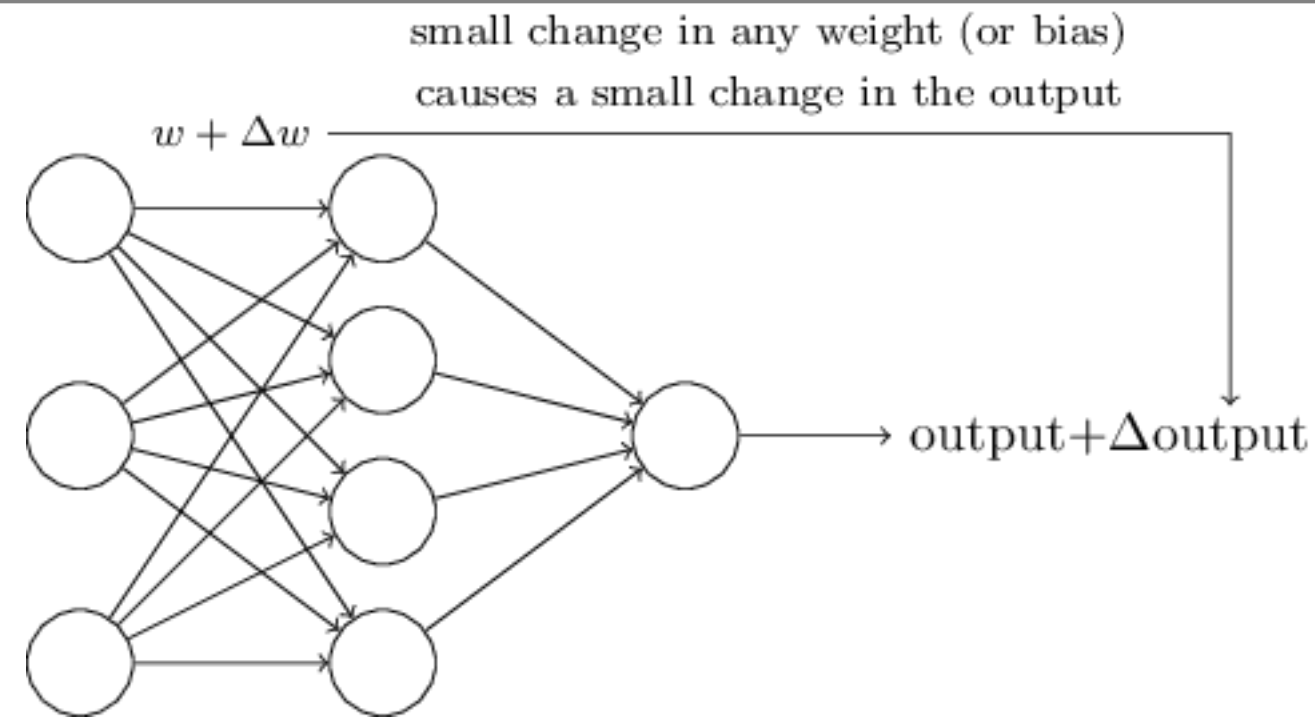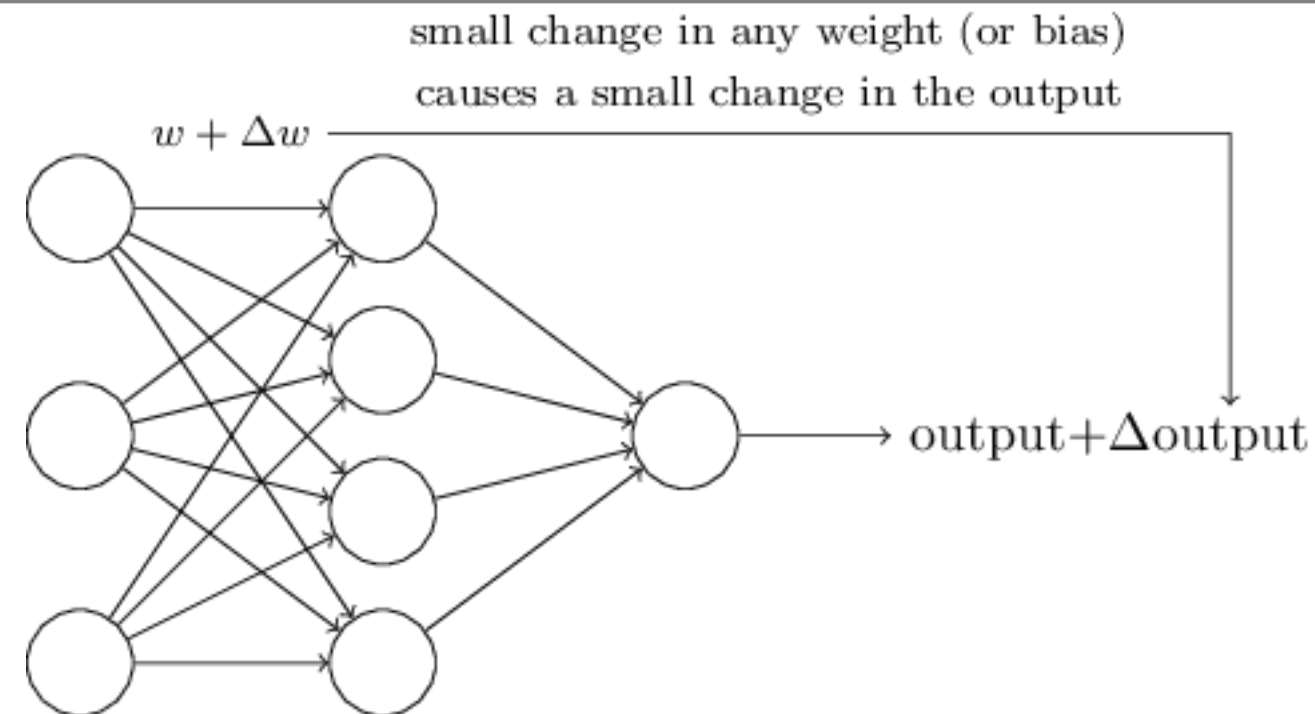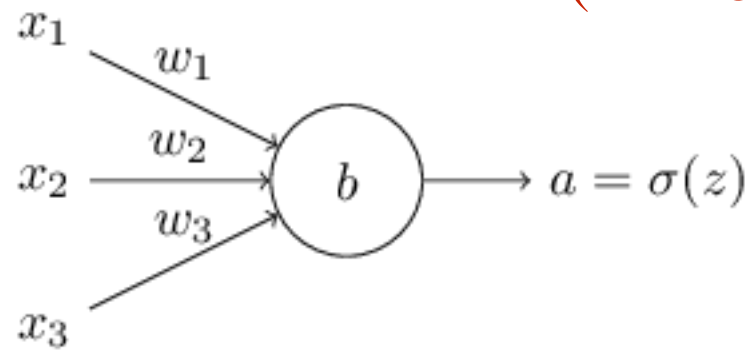
# Training a Neural Net

$$C(w, b) \equiv \frac{1}{2n} \sum_x \|y(x) - a\|^2$$
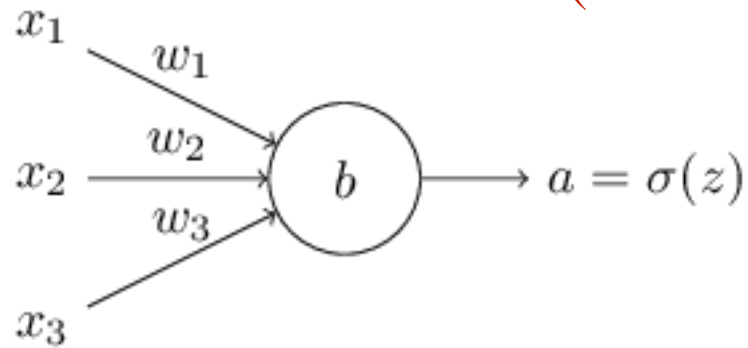
<span style="color:red">training examples (mini-batches)</span>

$x_1$
$w_1$
$w_2$
$x_2$ $\quad b \quad \rightarrow a = \sigma(z)$
$w_3$
$x_3$

<span style="color:red">learning rate</span>

$$w_k \rightarrow w_k' = w_k - \eta \frac{\partial C}{\partial w_k}$$

$$b_l \rightarrow b_l' = b_l - \eta \frac{\partial C}{\partial b_l}.$$

small change in any weight (or bias)
causes a small change in the output

$w + \Delta w$

output$+\Delta$output

$$\Delta \text{output} \approx \sum_j \frac{\partial \, \text{output}}{\partial w_j} \Delta w_j + \frac{\partial \, \text{output}}{\partial b} \Delta b$$
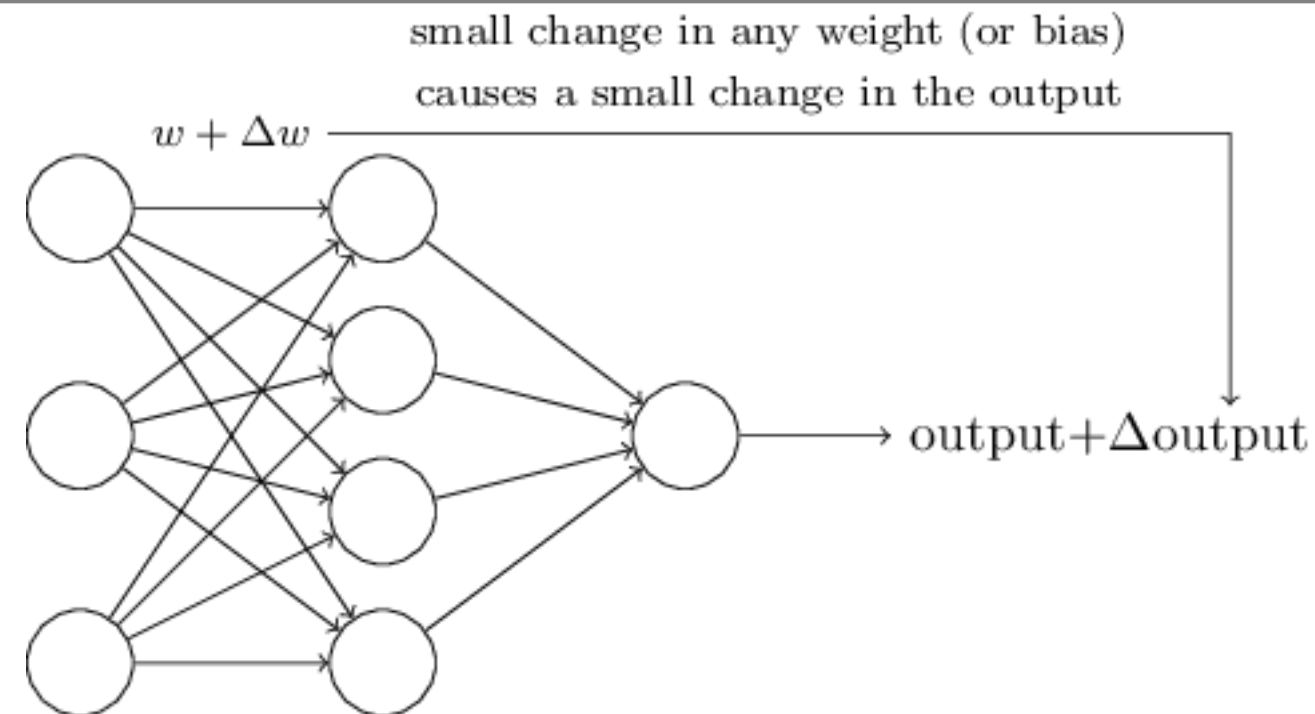
# Training a Neural Net

$$C(w, b) \equiv \frac{1}{2n} \sum_{x} \|y(x) - a\|^2$$

training examples
(mini-batches)



small change in any weight (or bias)
causes a small change in the output

$w + \Delta w$
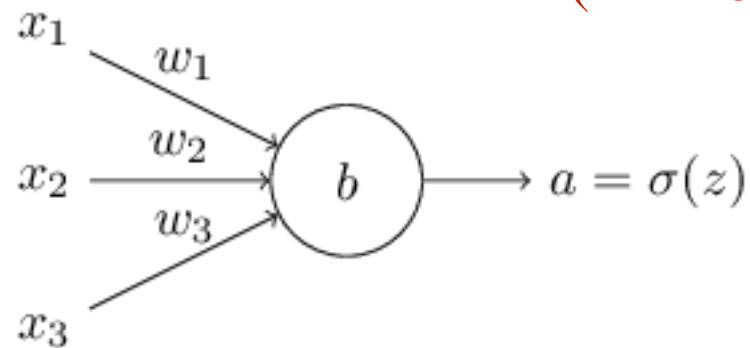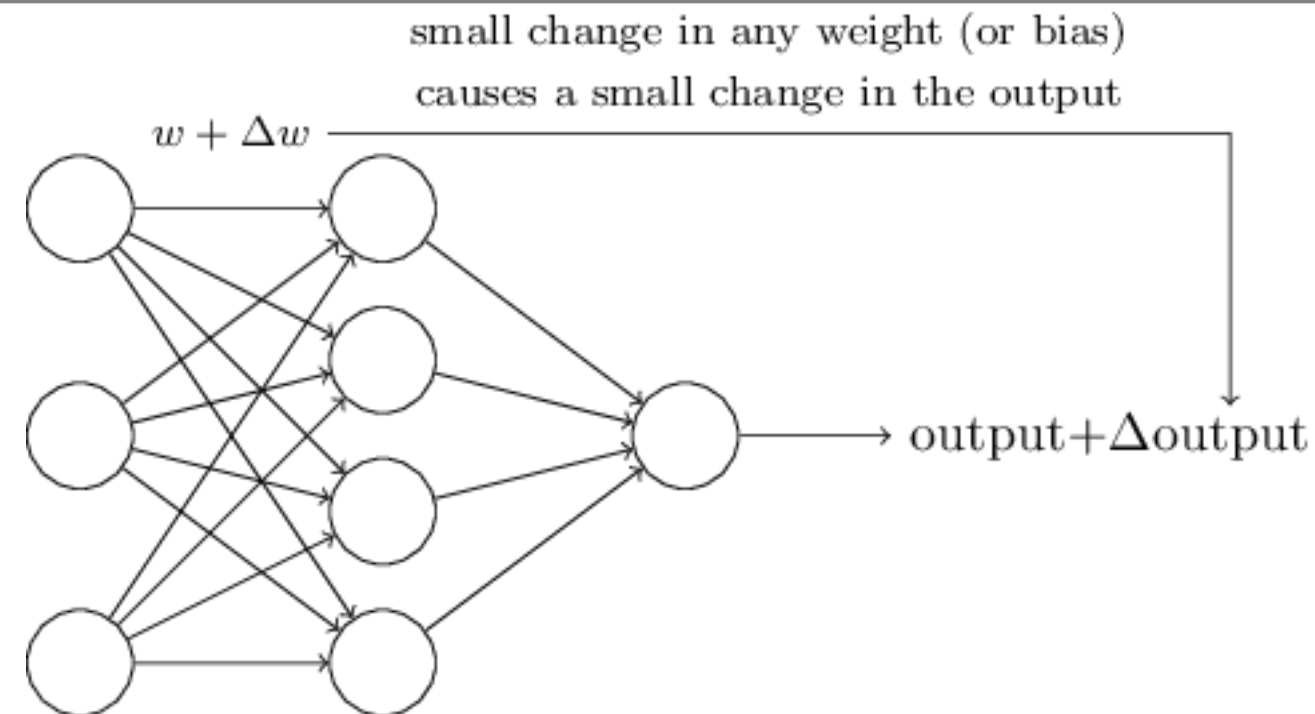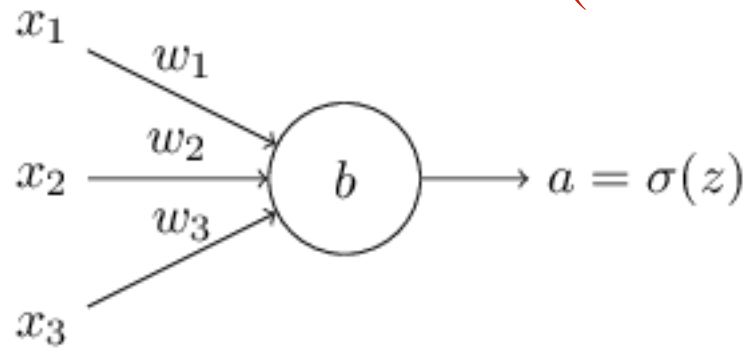
output$+\Delta$output

$$\Delta \text{output} \approx \sum_{j} \frac{\partial \, \text{output}}{\partial w_j} \Delta w_j + \frac{\partial \, \text{output}}{\partial b} \Delta b$$

$$w_k \rightarrow w_k' = w_k - \frac{\eta}{m} \sum_{j} \frac{\partial C_{X_j}}{\partial w_k}$$

$$b_l \rightarrow b_l' = b_l - \frac{\eta}{m} \sum_{j} \frac{\partial C_{X_j}}{\partial b_l},$$

Stochastic Gradient Descent

# Training a Neural Net

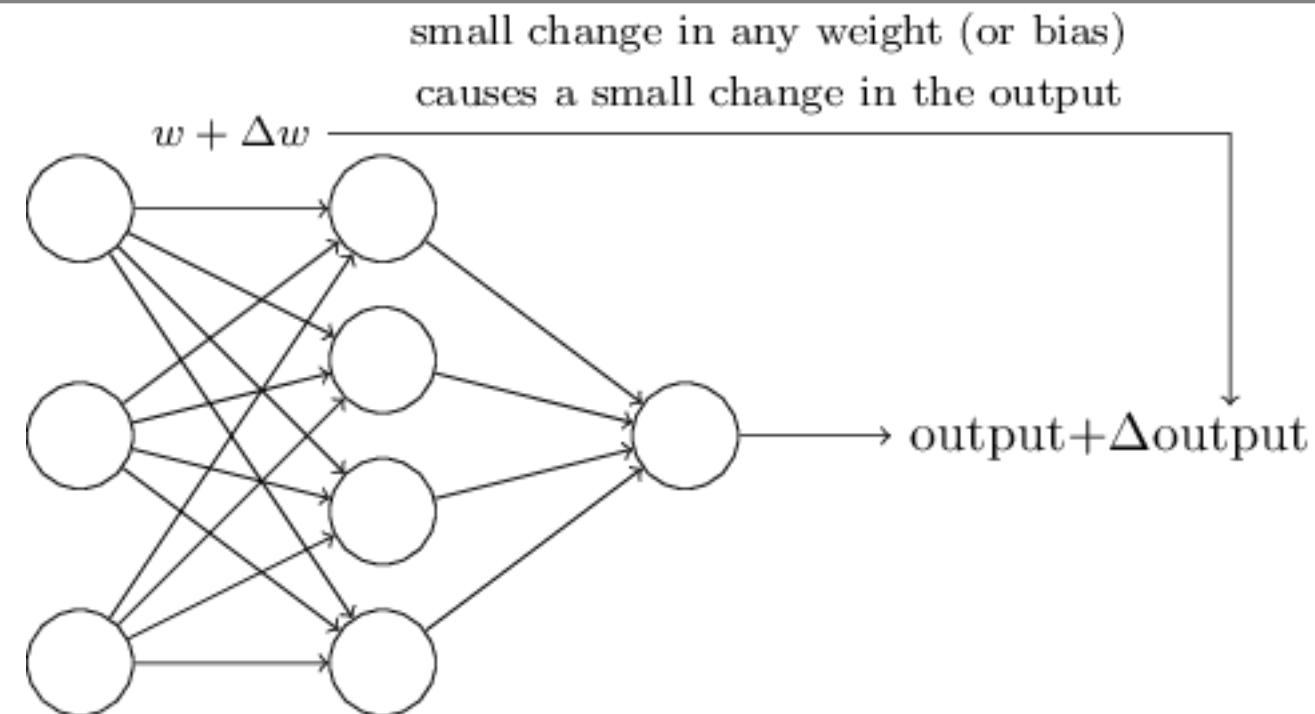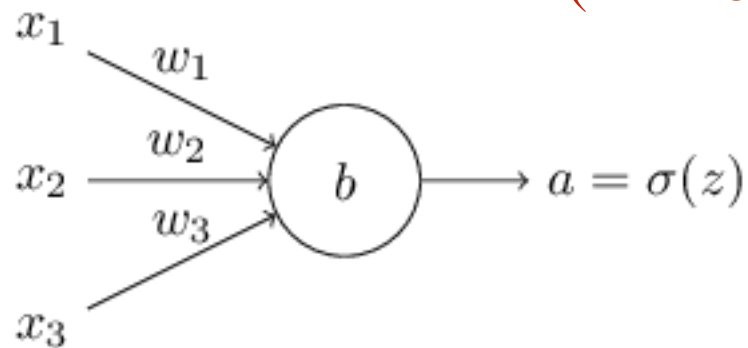$$C(w, b) \equiv \frac{1}{2n} \sum_{x} \|y(x) - a\|^2$$

training examples
(mini-batches)



$$w_k \to w_k' = w_k - \frac{\eta}{m} \sum_j \frac{\partial C_{X_j}}{\partial w_k}$$

$$b_l \to b_l' = b_l - \frac{\eta}{m} \sum_j \frac{\partial C_{X_j}}{\partial b_l},$$

potentially difficult to compute

Stochastic Gradient Descent

small change in any weight (or bias)
causes a small change in the output

$$w + \Delta w$$

output+$\Delta$output

$$\Delta \text{output} \approx \sum_j \frac{\partial \, \text{output}}{\partial w_j} \Delta w_j + \frac{\partial \, \text{output}}{\partial b} \Delta b$$
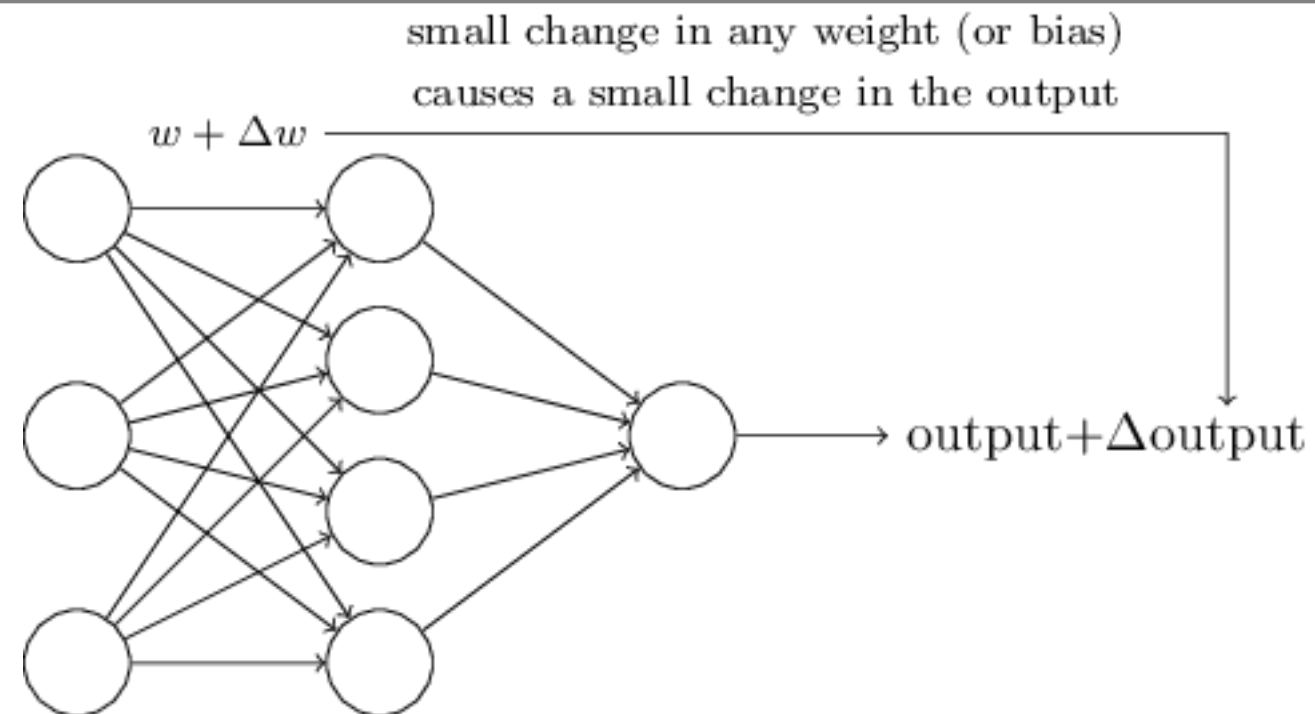
$x_1$
$w_1$
$x_2$
$w_2$
$w_3$
$x_3$
$b$
$a = \sigma(z)$

# Backpropagation

$$C(w, b) \equiv \frac{1}{2n} \sum_x \|y(x) - a\|^2$$

# Backpropagation

$$C(w, b) \equiv \frac{1}{2n} \sum_x \|y(x) - a\|^2 \qquad \text{Quadratic}$$

# Backpropagation

$$C(w, b) \equiv \frac{1}{2n} \sum_x \|y(x) - a\|^2 \qquad \text{Quadratic}$$

$$C = -\frac{1}{n} \sum_x [y \ln a + (1 - y) \ln(1 - a)] \qquad \text{Cross entropy}$$

# Backpropagation

$$C(w, b) \equiv \frac{1}{2n} \sum_x \|y(x) - a\|^2$$

Quadratic

$$C = -\frac{1}{n} \sum_x [y \ln a + (1-y) \ln(1-a)]$$

Cross entropy



small change in any weight (or bias) causes a small change in the output
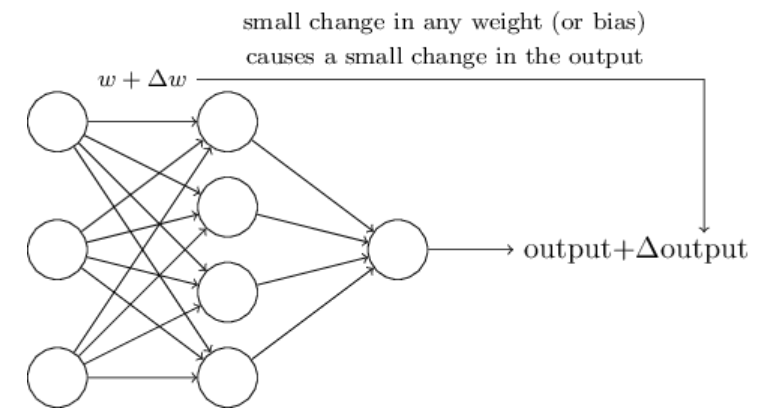
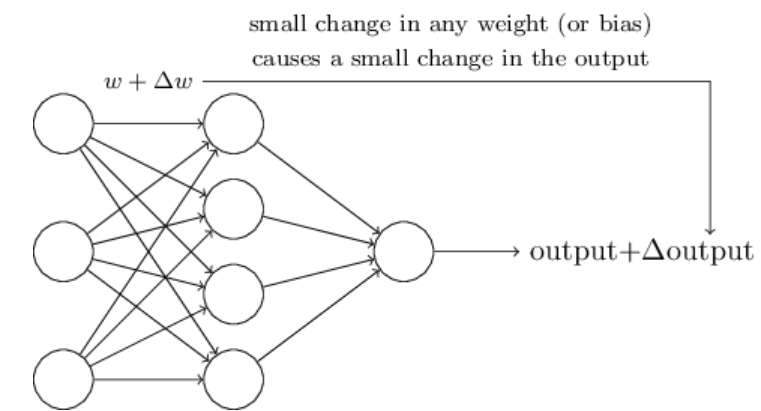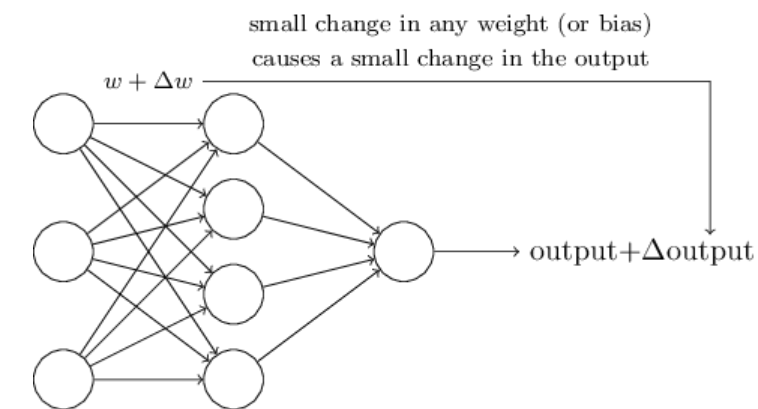$w + \Delta w$

output$+\Delta$output

# Backpropagation

$$C(w, b) \equiv \frac{1}{2n} \sum_x \|y(x) - a\|^2$$   Quadratic

$$C = -\frac{1}{n} \sum_x [y \ln a + (1 - y) \ln(1 - a)]$$   Cross entropy



small change in any weight (or bias)
causes a small change in the output

$w + \Delta w$

output$+\Delta$output

$$\delta^L = \nabla_a C \odot \sigma'(z^L)$$

$$\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$$

$$\frac{\partial C}{\partial b_j^l} = \delta_j^l$$

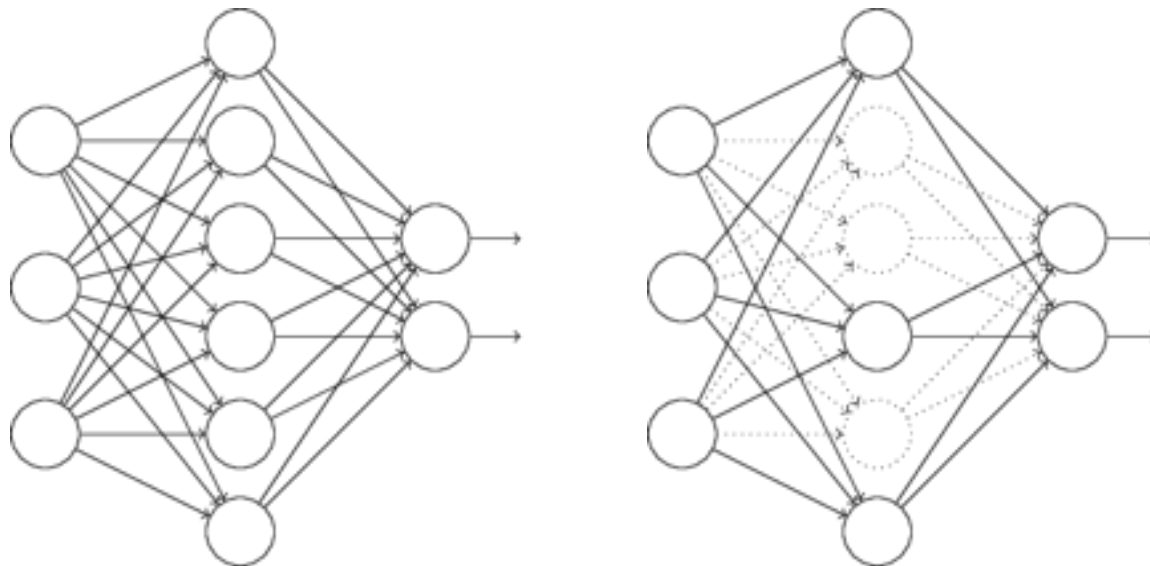$$\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$$

# Backpropagation

$$C(w, b) \equiv \frac{1}{2n} \sum_x \|y(x) - a\|^2 \qquad \text{Quadratic}$$

$$C = -\frac{1}{n} \sum_x [y \ln a + (1 - y) \ln(1 - a)] \qquad \text{Cross entropy}$$

small change in any weight (or bias) causes a small change in the output

$w + \Delta w$

output$+\Delta$output

| | |
|---|---|
| $\delta^L = \nabla_a C \odot \sigma'(z^L)$ | 1. **Input** $x$: Set the corresponding activation $a^1$ for the input layer. |
| $\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$ | 2. **Feedforward:** For each $l = 2, 3, \ldots, L$ compute $z^l = w^l a^{l-1} + b^l$ and $a^l = \sigma(z^l)$. |
| | 3. **Output error** $\delta^L$: Compute the vector $\delta^L = \nabla_a C \odot \sigma'(z^L)$. |
| $\frac{\partial C}{\partial b_j^l} = \delta_j^l$ | 4. **Backpropagate the error:** For each $l = L - 1, L - 2, \ldots, 2$ compute $\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$. |
| $\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$ | 5. **Output:** The gradient of the cost function is given by $\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$ and $\frac{\partial C}{\partial b_j^l} = \delta_j^l$. |

# Guarding against overfitting and saturation
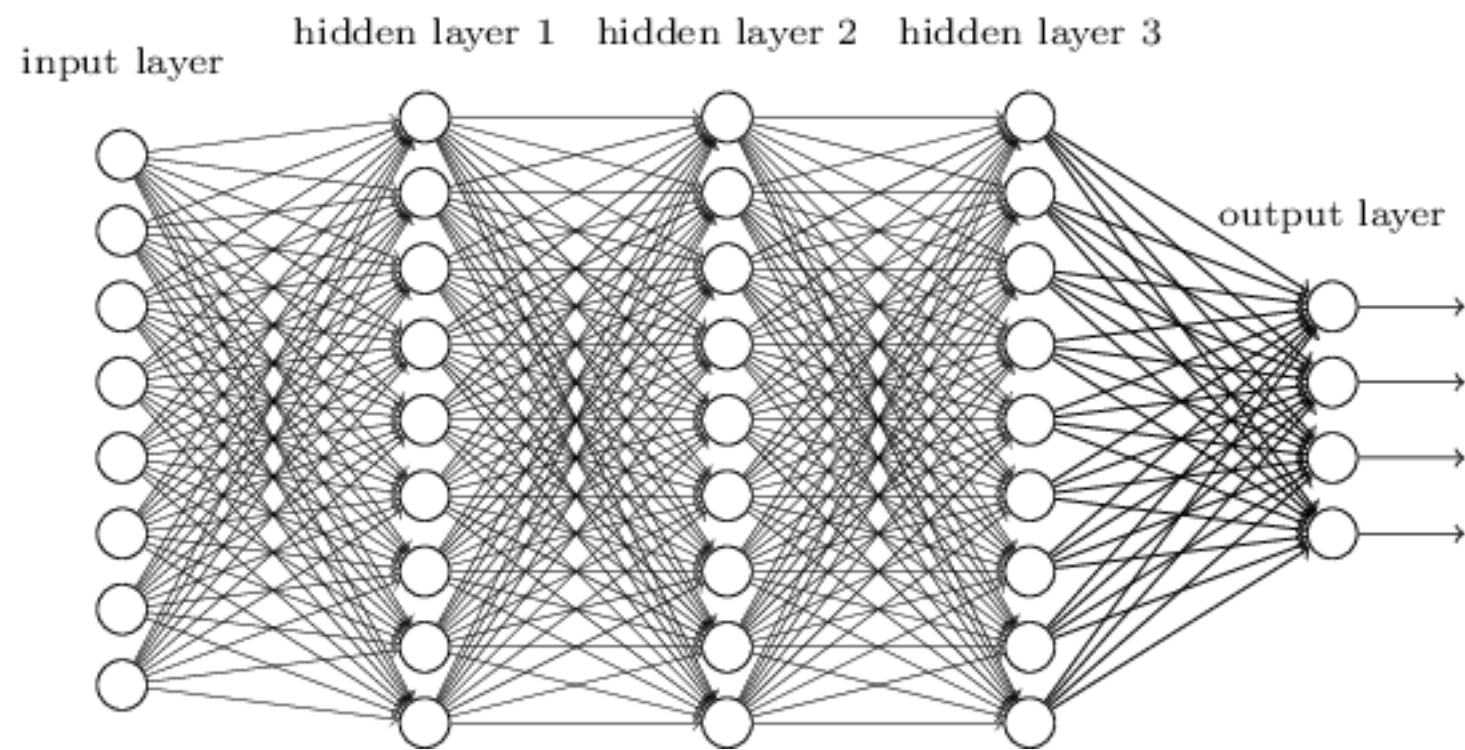
· Regularization (weight decay)

$$C = -\frac{1}{n} \sum_{xj} \left[ y_j \ln a_j^L + (1 - y_j) \ln(1 - a_j^L) \right] + \frac{\lambda}{2n} \sum_w w^2$$
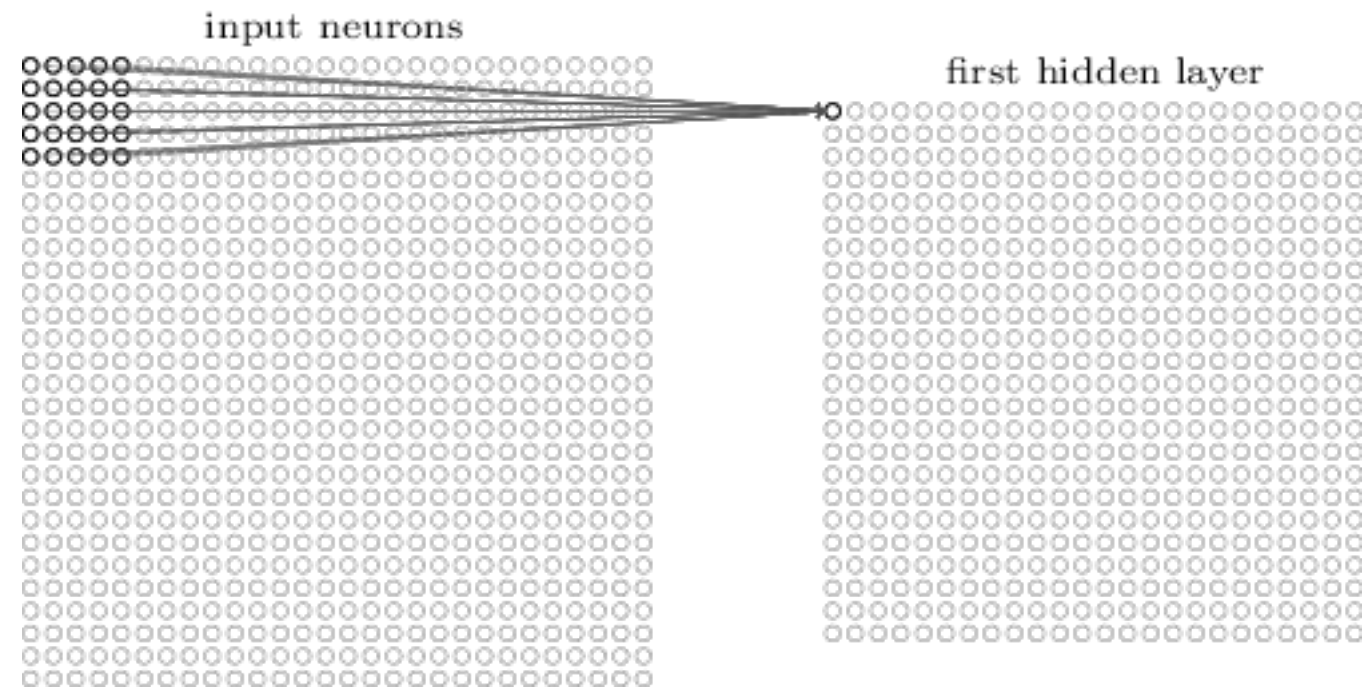
· Dropout



· Expanding the training data
  - reflections
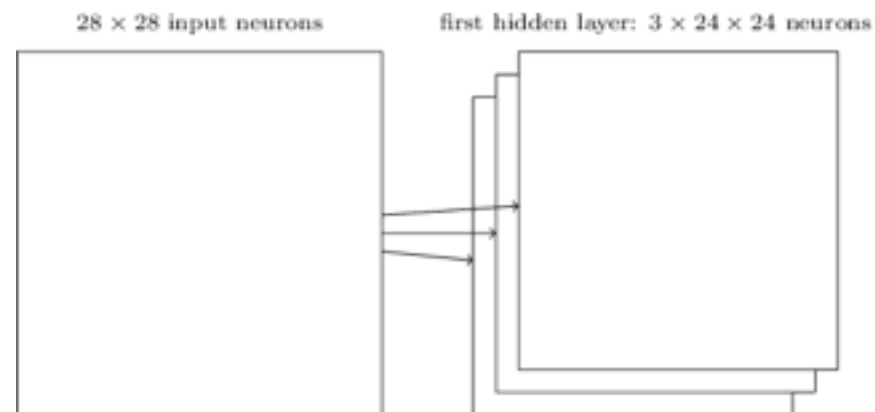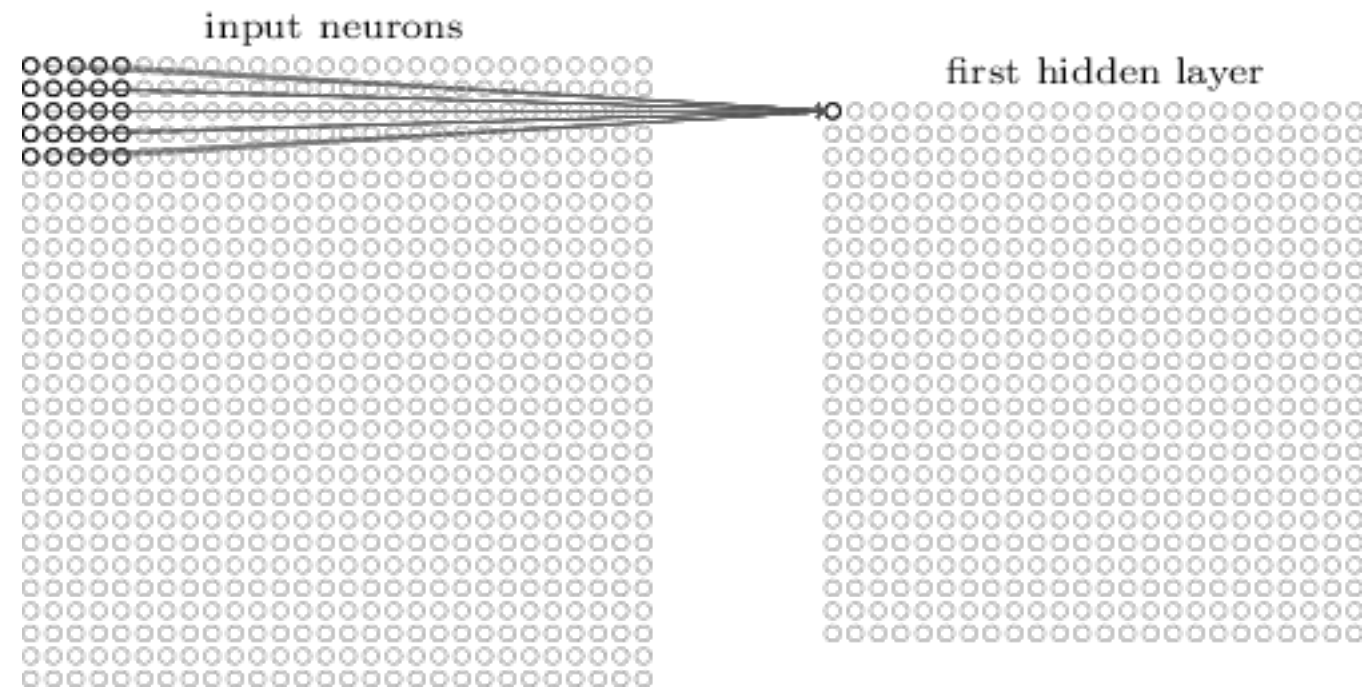  - translations
  - rotations
  - contrast
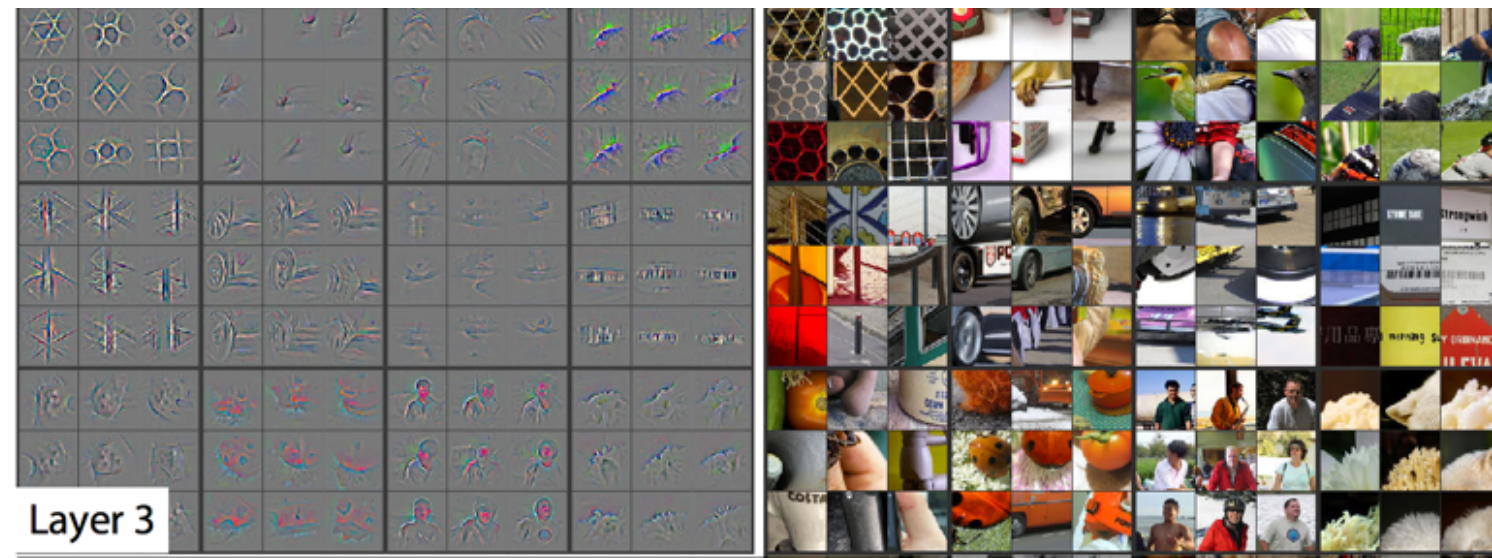  - skewness

# Convolutional nets

# Convolutional nets



input neurons

first hidden layer

# Convolutional nets

input neurons

first hidden layer

28 × 28 input neurons

first hidden layer: 3 × 24 × 24 neurons

# Convolutional nets



Layer 1

Layer 2

Layer 3

Zeiler & Fergus (2013)

# Convolutional nets

input neurons

first hidden layer

$28 \times 28$

convolutional layer
$20 \times 24 \times 24$

pooling layer
$20 \times 12 \times 12$

100 sigmoid
neurons

10 neuron
output layer
(softmax)

Layer 3

# Convolutional nets



input neurons

first hidden layer

"**AlexNet**":Krizhevsky, Sutskever, & Hinton (2012)