

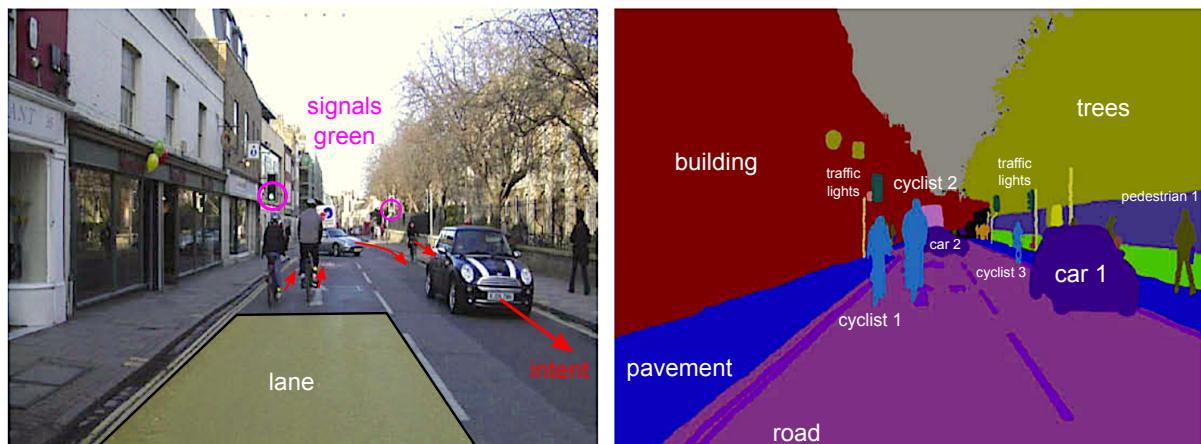
Deep Learning

Part 1: From single neurons to multi-layer perceptrons

Dr. Richard E. Turner (ret26@cam.ac.uk)

Engineering Tripos Part IB
Paper 8: Information Engineering

Example perception problems for a self-driving car



Very hard to hand-design systems to do this reliably
fairly easy to collect lots of labelled data
⇒ learn systems from data instead

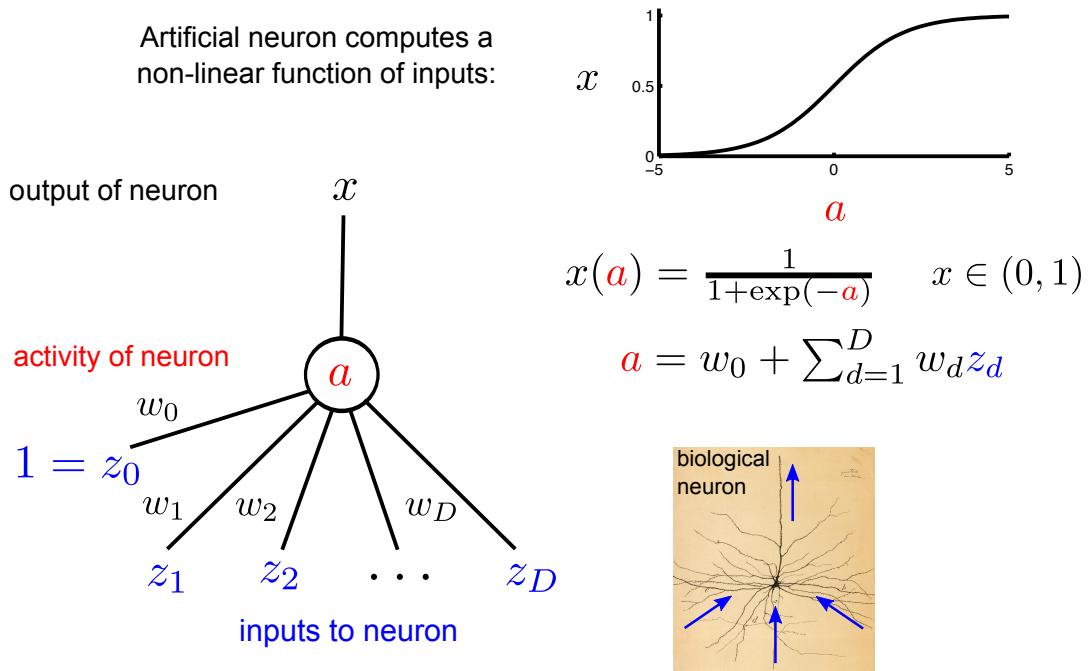
Object recognition

- Focus on **objection recognition**: predict objects present in an image
- A canonical classification problem: approaches to object recognition can be generalised to more complex tasks
- All state-of-the-art approaches use **neural networks**
- Outline of lectures:
 - single neurons and logistic classification (2 lectures)
 - neural networks and multi-layer perceptrons (1 lecture)
 - convolutional neural networks (2 lectures)
 - worked examples integrated through out (no specific examples class)

Current state-of-the-art

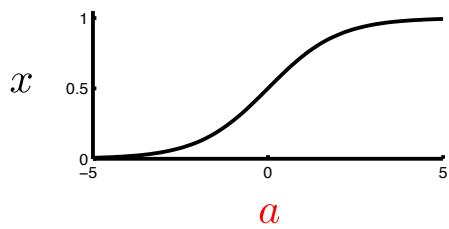
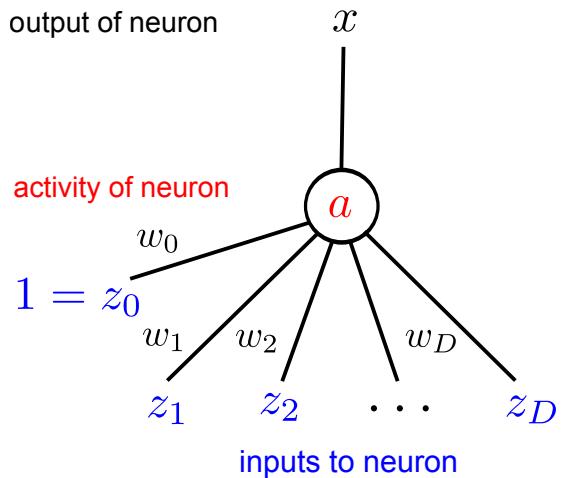
<https://www.clarifai.com/demo>

A single neuron

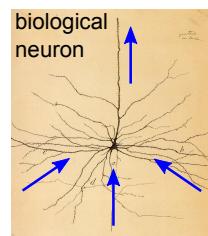


A single neuron

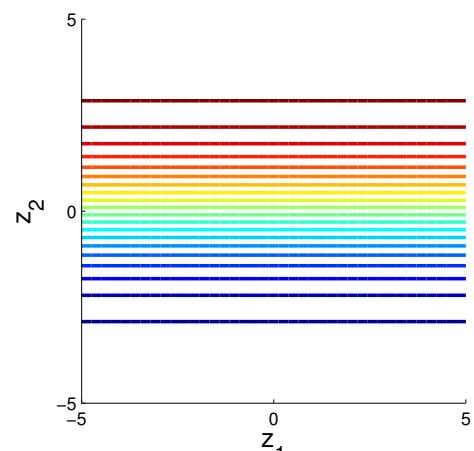
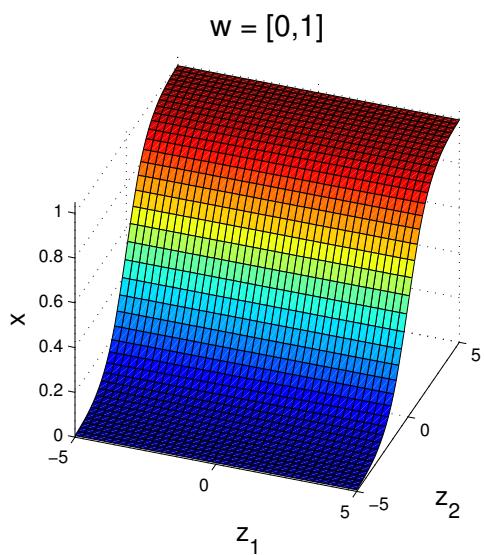
Artificial neuron computes a non-linear function of inputs:



$$a = \sum_{d=0}^D w_d z_d$$

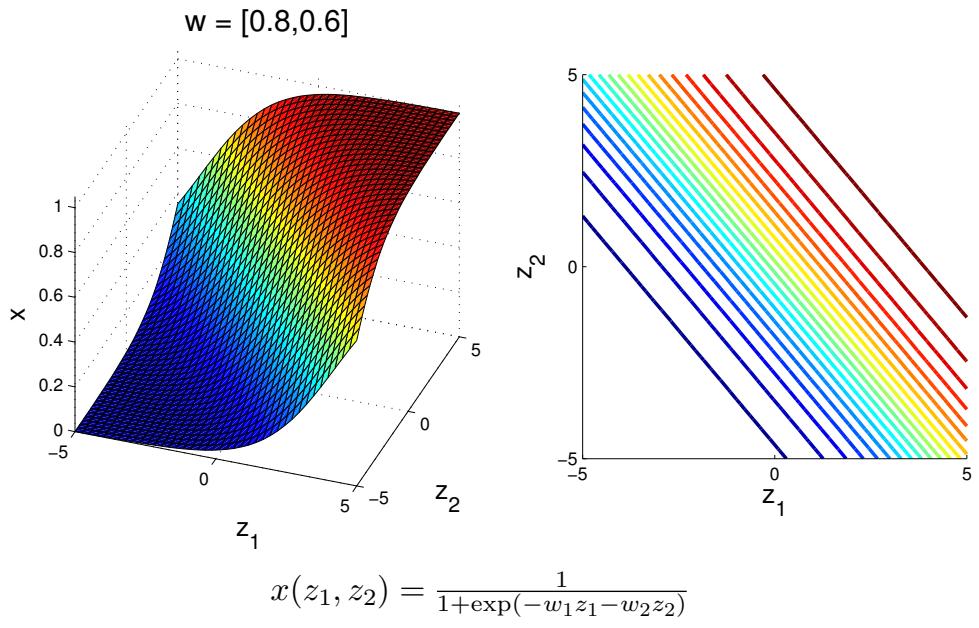


Input-output function of a single neuron

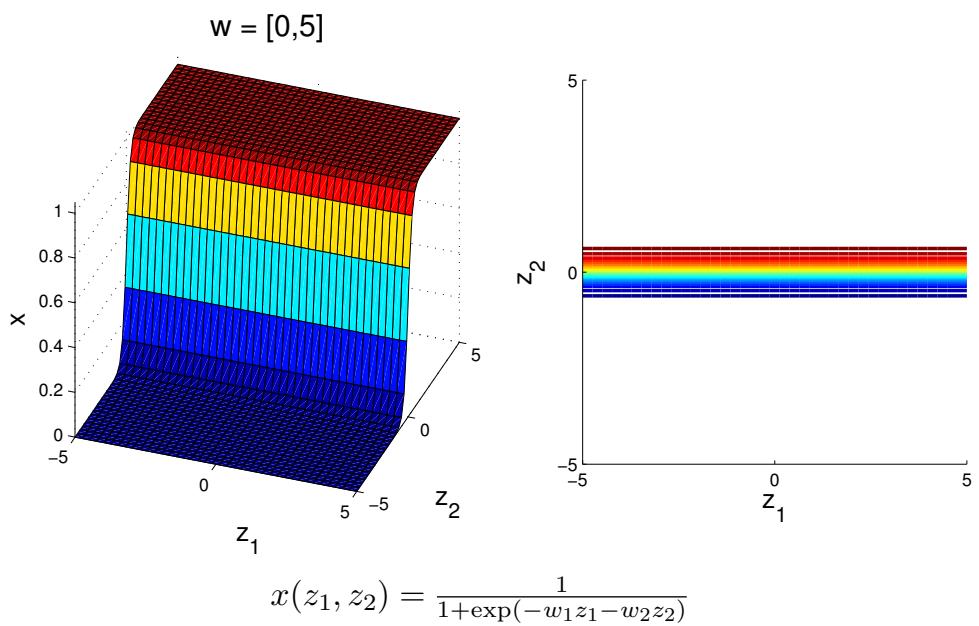


$$x(z_1, z_2) = \frac{1}{1+\exp(-w_1 z_1 - w_2 z_2)}$$

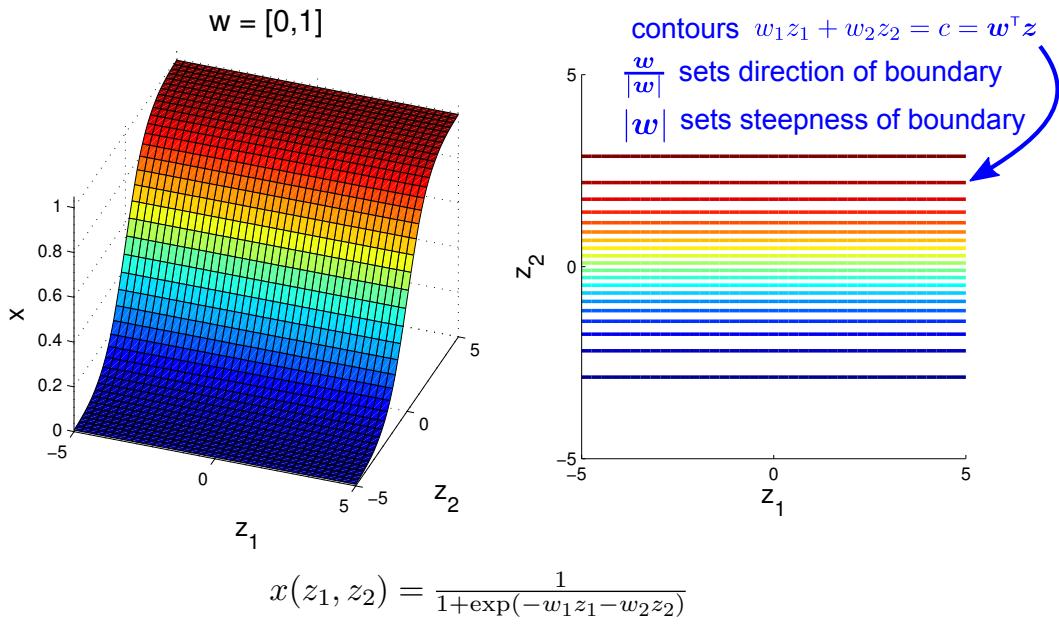
Input-output function of a single neuron



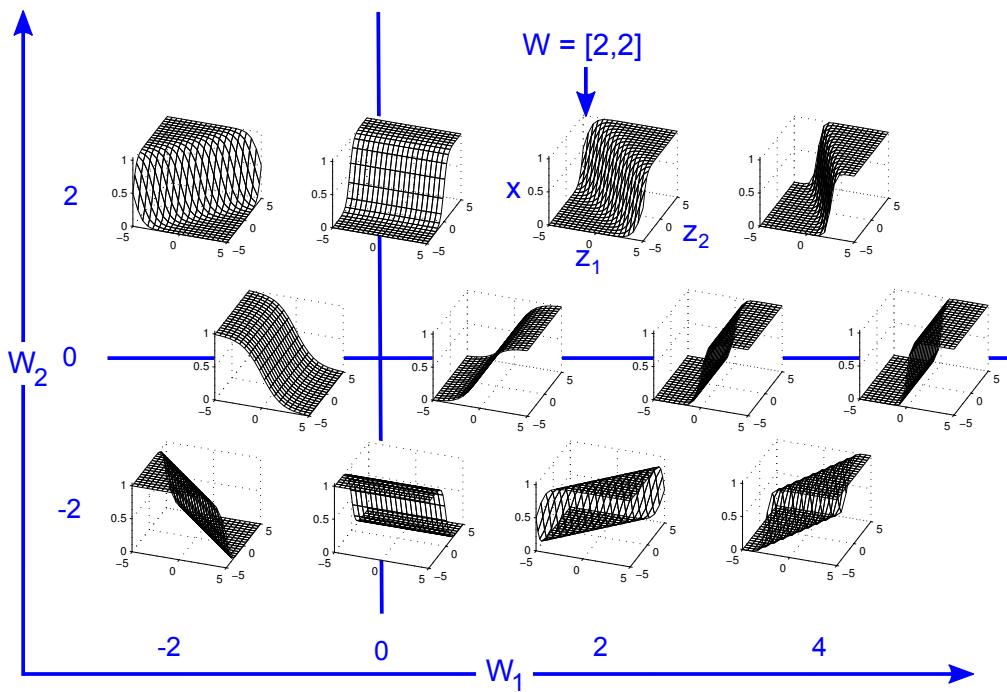
Input-output function of a single neuron



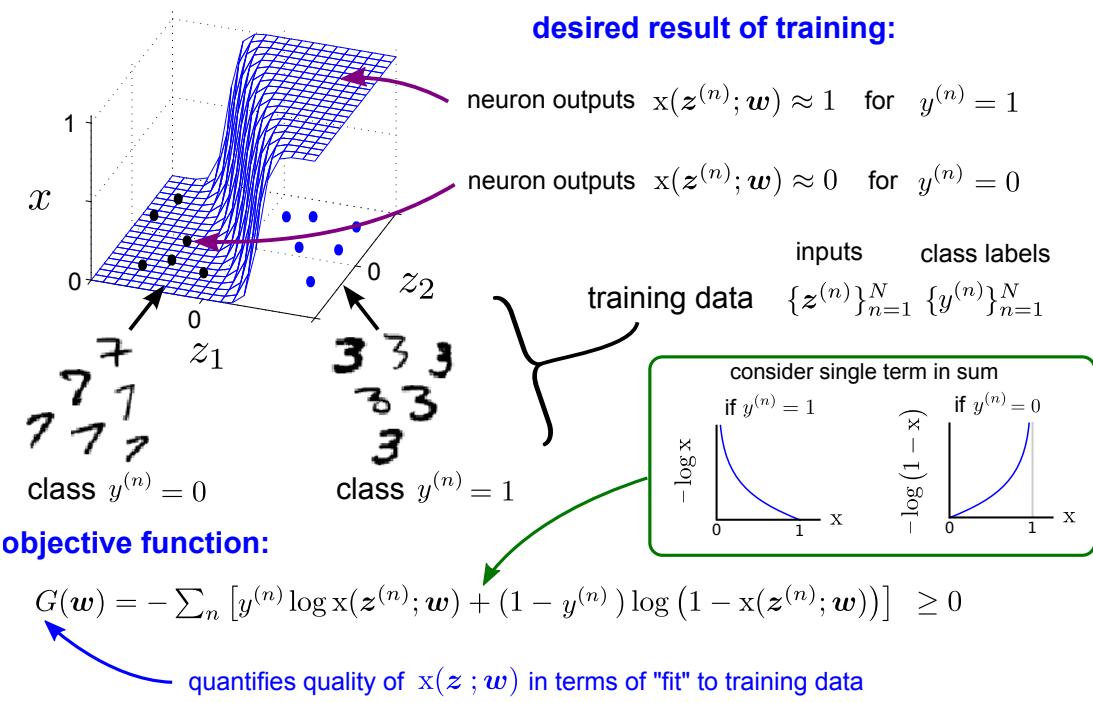
Input-output function of a single neuron



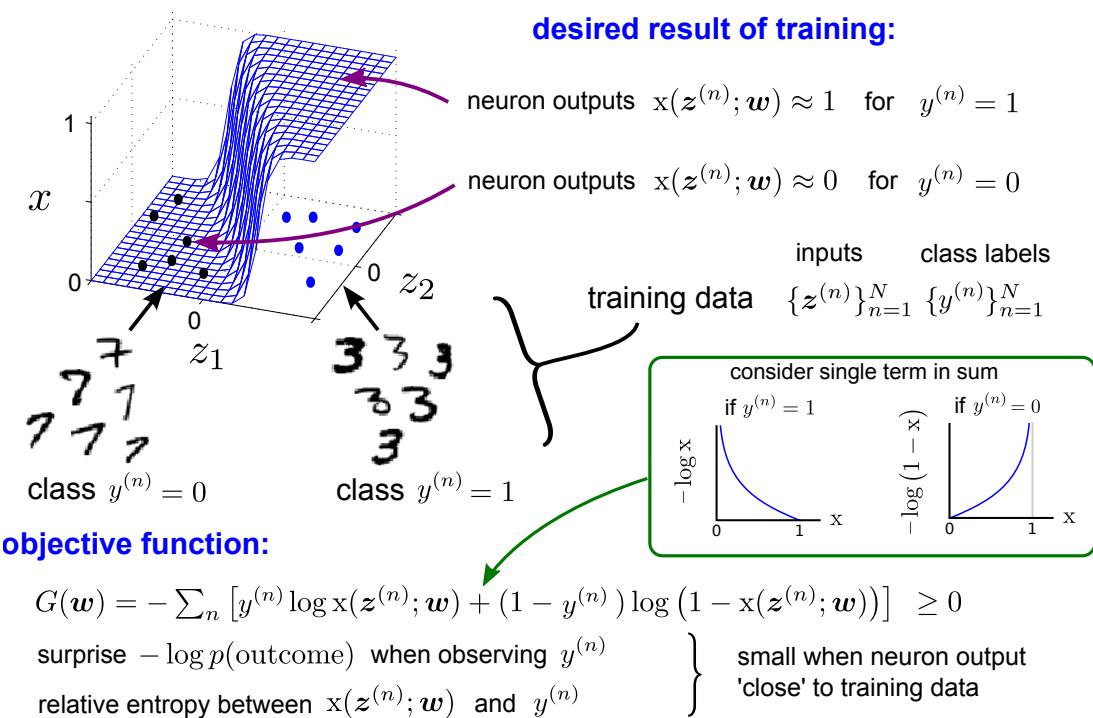
Weight space of a single neuron



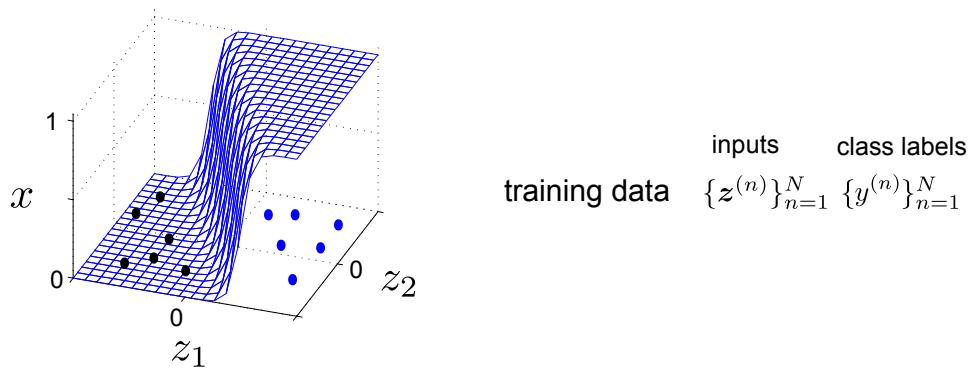
Training a single neuron



Training a single neuron



Training a single neuron: a problem



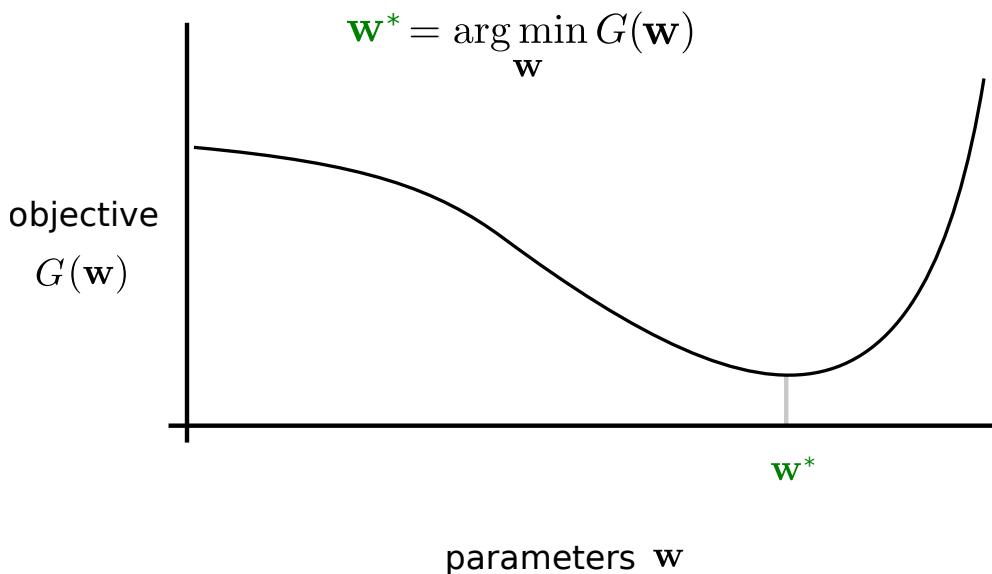
objective function:

$$G(\mathbf{w}) = - \sum_n [y^{(n)} \log x(\mathbf{z}^{(n)}; \mathbf{w}) + (1 - y^{(n)}) \log (1 - x(\mathbf{z}^{(n)}; \mathbf{w}))] \geq 0$$

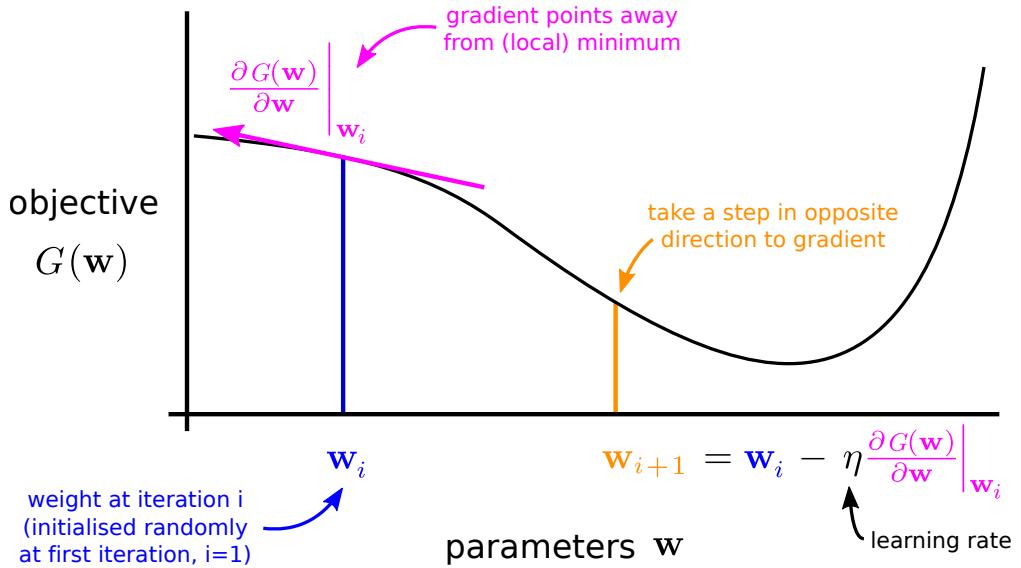
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} G(\mathbf{w}) \quad \text{choose weights that minimise network's surprise about training data}$$

There is no analytic solution to this optimisation problem: what should we do?

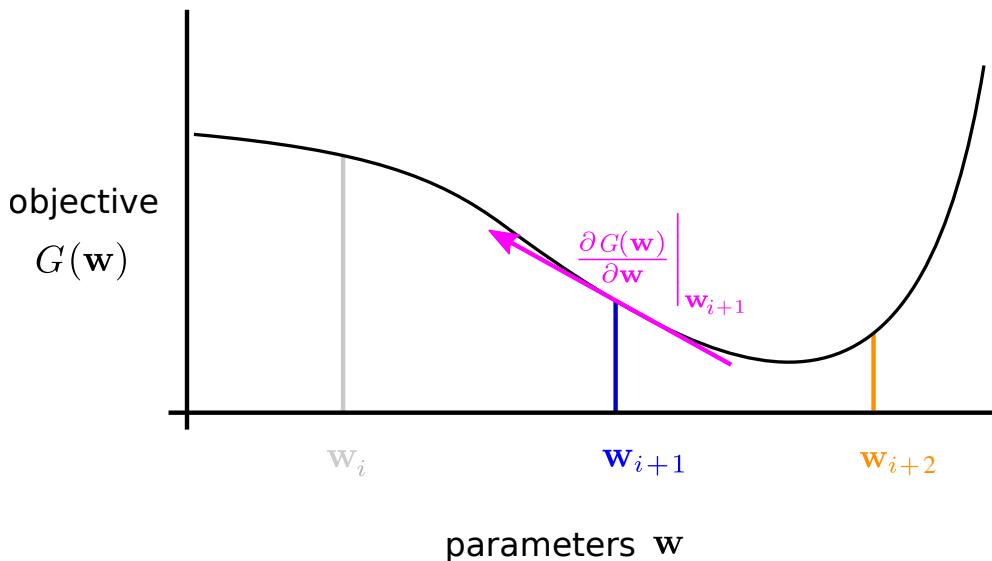
Goal of training



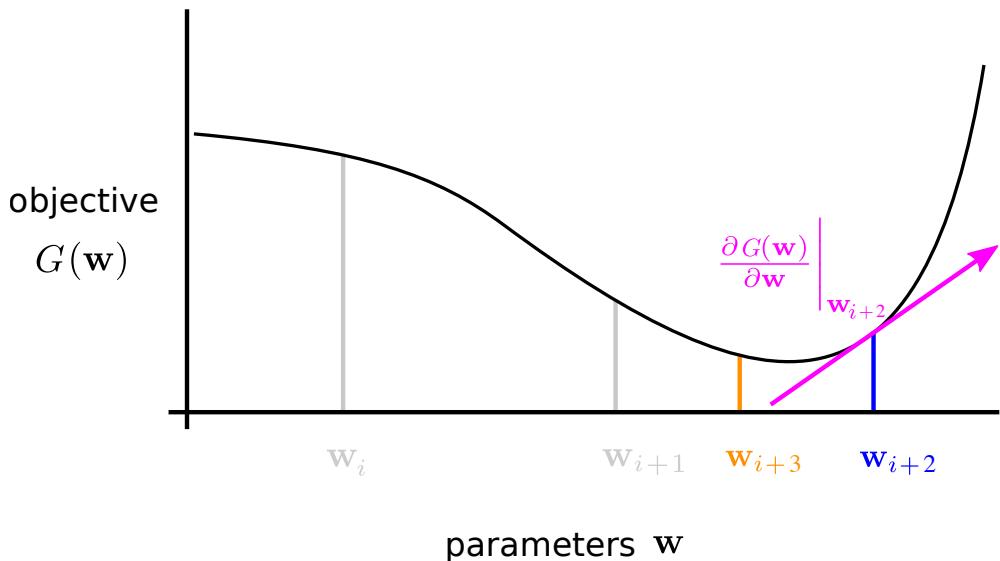
Training using gradient descent



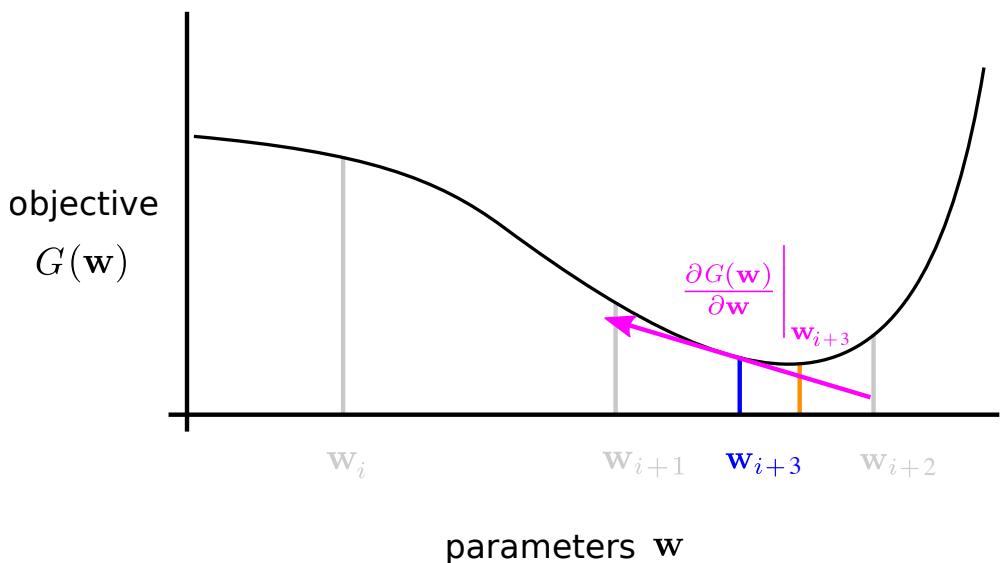
Training using gradient descent



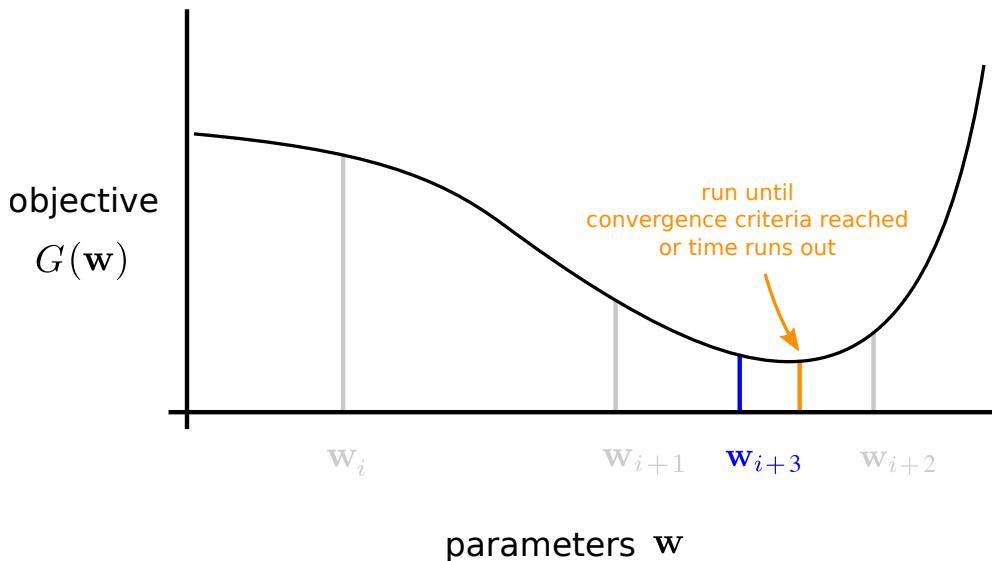
Training using gradient descent



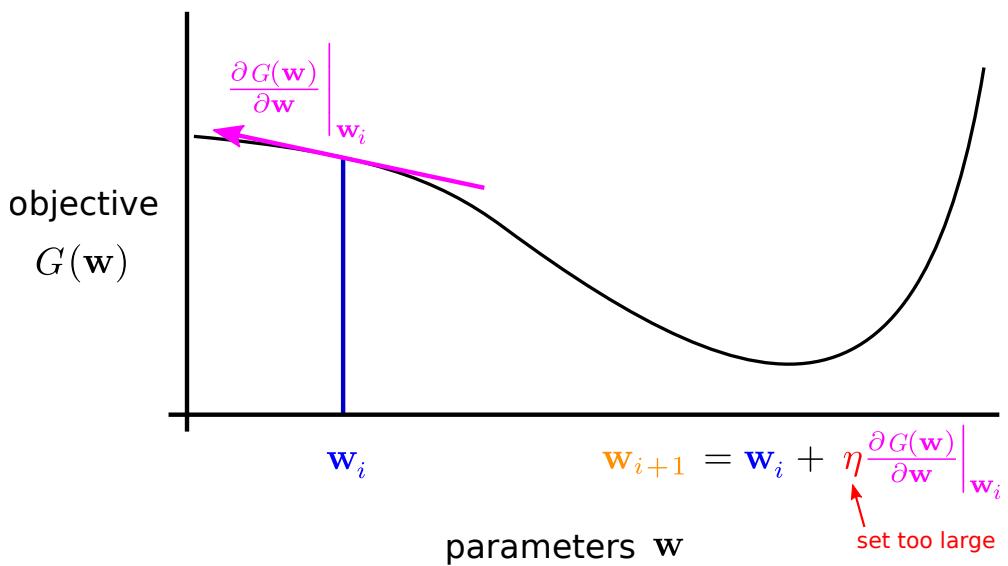
Training using gradient descent



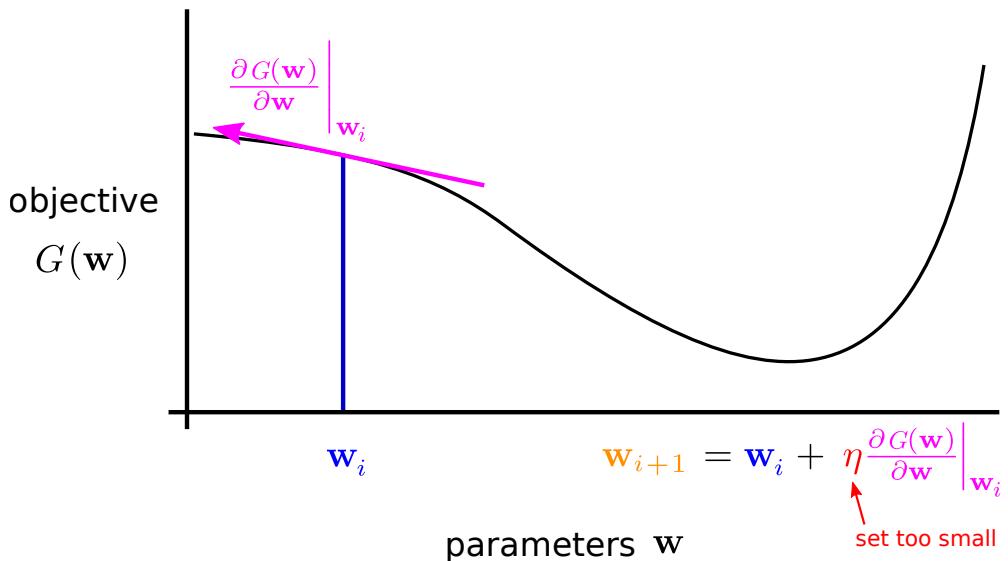
Training using gradient descent



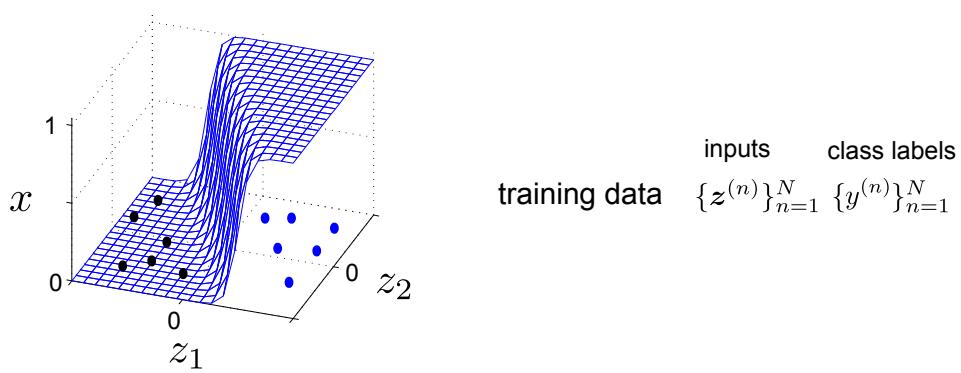
What happens if the learning rate η is set too large?



What happens if the learning rate η is set too small?



Training a single neuron: a solution



objective function:

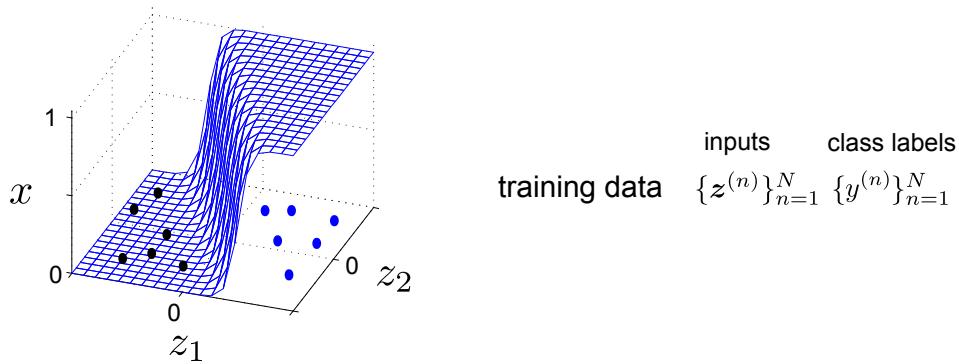
$$G(\mathbf{w}) = -\sum_n [y^{(n)} \log x(\mathbf{z}^{(n)}; \mathbf{w}) + (1 - y^{(n)}) \log (1 - x(\mathbf{z}^{(n)}; \mathbf{w}))] \geq 0$$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} G(\mathbf{w}) \quad \text{choose weights that minimise network's surprise about training data}$$

There is no analytic solution to this optimisation problem: what should we do?

Answer: use gradient descent

Training a single neuron



objective function:

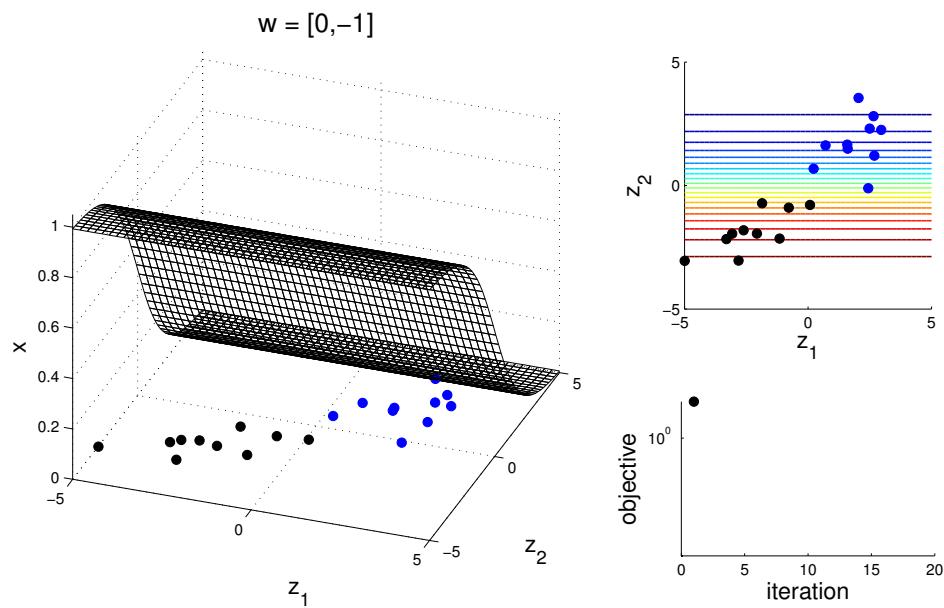
$$G(\mathbf{w}) = - \sum_n [y^{(n)} \log x(\mathbf{z}^{(n)}; \mathbf{w}) + (1 - y^{(n)}) \log (1 - x(\mathbf{z}^{(n)}; \mathbf{w}))] \geq 0$$

$\mathbf{w}^* = \arg \min_{\mathbf{w}} G(\mathbf{w})$ choose weights that minimise network's surprise about training data

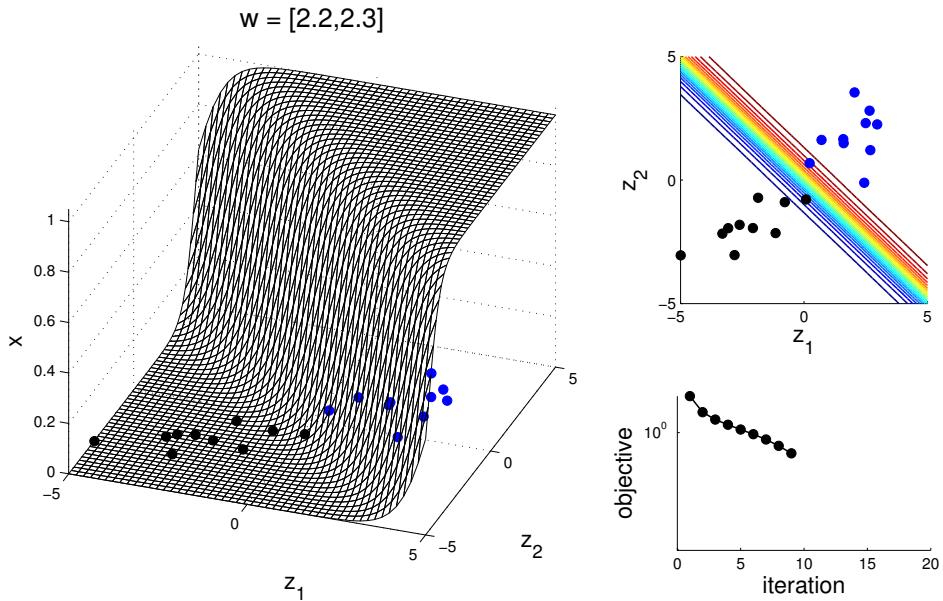
$$\frac{d}{d\mathbf{w}} G(\mathbf{w}) = \sum_n \frac{dG(\mathbf{w})}{dx^{(n)}} \frac{dx^{(n)}}{d\mathbf{w}} = - \sum_n (\underbrace{y^{(n)} - x^{(n)}}_{\text{prediction error}}) \mathbf{z}^{(n)} \leftarrow \begin{array}{l} \text{feature} \\ \text{prediction error} \end{array}$$

$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{d}{d\mathbf{w}} G(\mathbf{w})$ iteratively step down the objective (gradient points up hill)

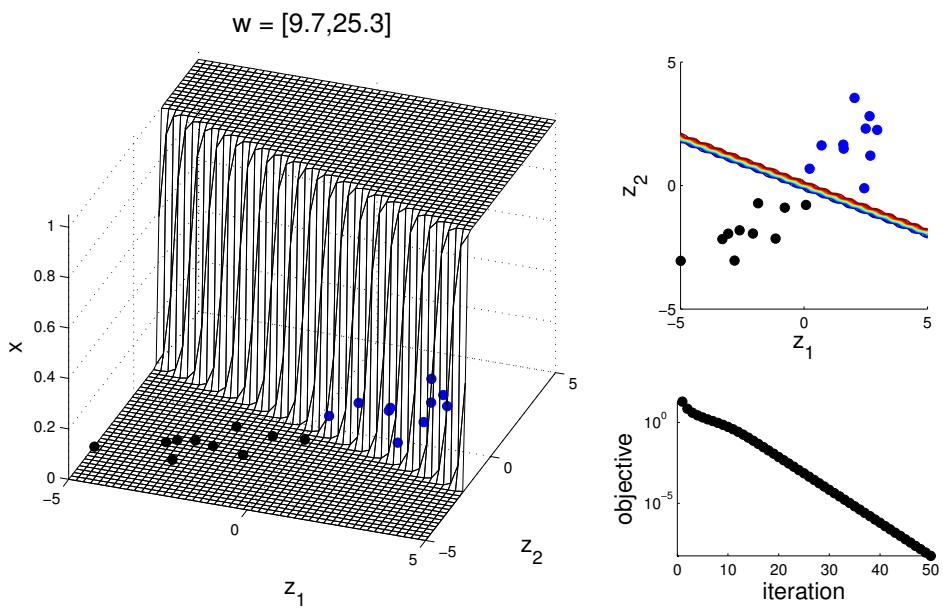
Training a single neuron



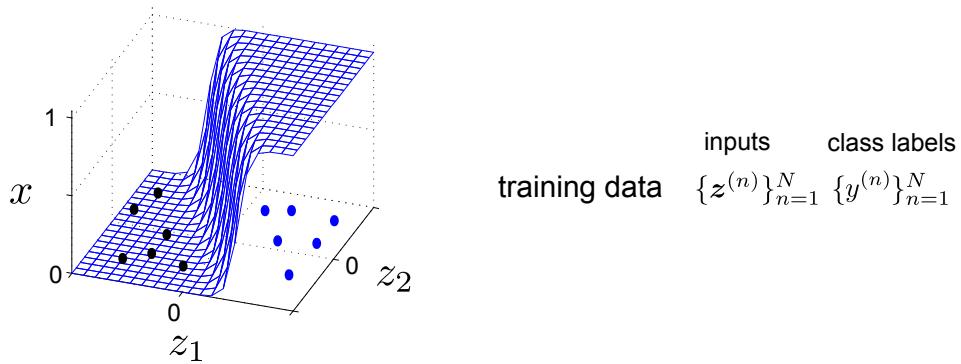
Training a single neuron



Training a single neuron



Over-fitting and weight decay



objective function:

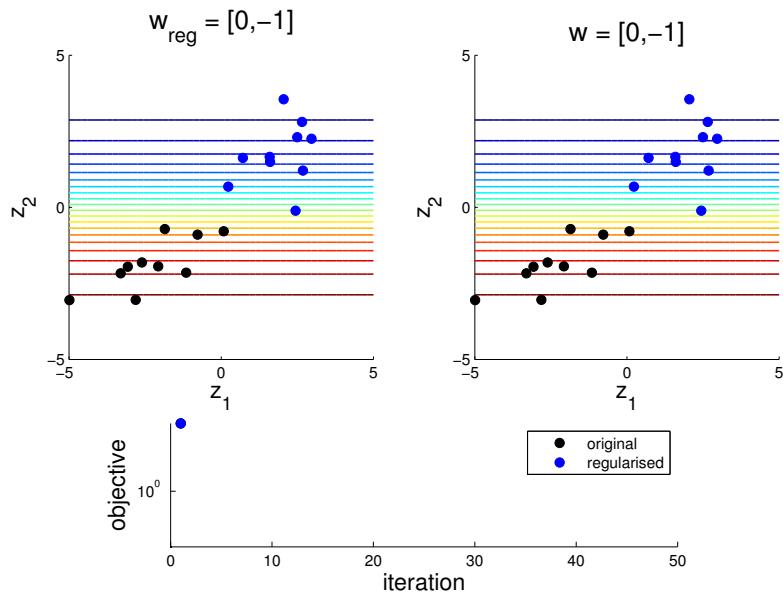
$$G(\mathbf{w}) = - \sum_n [y^{(n)} \log x(\mathbf{z}^{(n)}; \mathbf{w}) + (1 - y^{(n)}) \log (1 - x(\mathbf{z}^{(n)}; \mathbf{w}))] \geq 0$$

$$E(\mathbf{w}) = \frac{1}{2} \sum_i w_i^2 \quad \text{"regulariser": penalise magnitude of the weights and discourages the network from using extreme weights}$$

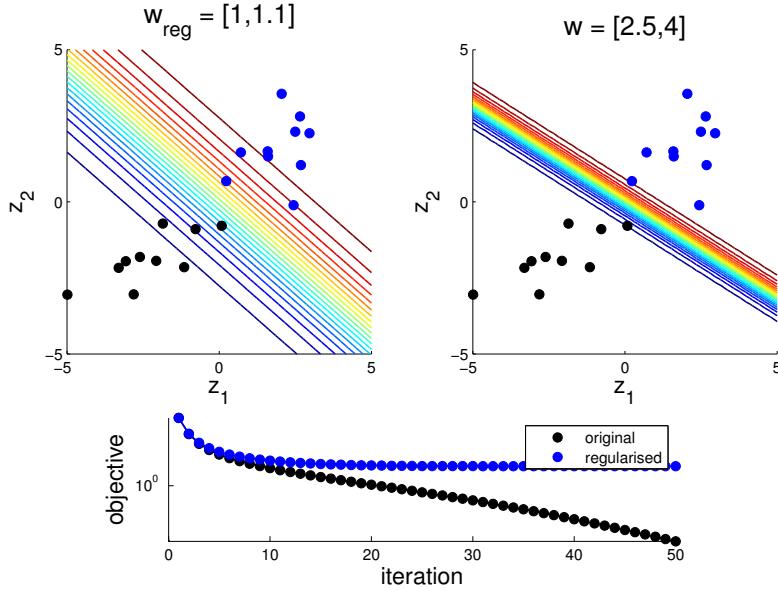
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} M(\mathbf{w}) = \arg \min_{\mathbf{w}} [G(\mathbf{w}) + \alpha E(\mathbf{w})]$$

$$\frac{d}{d\mathbf{w}} M(\mathbf{w}) = - \sum_n (y^{(n)} - x^{(n)}) \mathbf{z}^{(n)} + \alpha \mathbf{w} \quad \text{weight decay - shrinks weights towards zero}$$

Training a single neuron



Training a single neuron



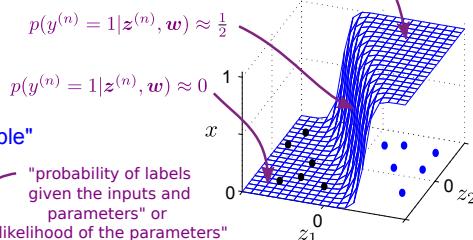
Probabilistic interpretation of the single neuron

Interpret the output of the network as a probability:

$$x(\mathbf{z}^{(n)}; \mathbf{w}) = p(y^{(n)} = 1 | \mathbf{z}^{(n)}, \mathbf{w})$$

$$p(y^{(n)} = 1 | \mathbf{z}^{(n)}, \mathbf{w}) \approx 1$$

Fitting method 1: maximum-likelihood learning



"choose setting of w that makes observed data most probable"

$$\mathbf{w}_{\text{ML}} = \arg \max_{\mathbf{w}} p(\{y^{(n)}\}_{n=1}^N | \{\mathbf{z}^{(n)}\}_{n=1}^N, \mathbf{w})$$

"probability of labels given the inputs and parameters or likelihood of the parameters"

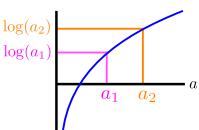
$$p(\{y^{(n)}\}_{n=1}^N | \{\mathbf{z}^{(n)}\}_{n=1}^N, \mathbf{w}) = \prod_{n=1}^N p(y^{(n)} | \mathbf{z}^{(n)}, \mathbf{w})$$

$$p(y^{(n)} | \mathbf{z}^{(n)}, \mathbf{w}) = \begin{cases} x(\mathbf{z}^{(n)}; \mathbf{w}) & \text{if } y^{(n)} = 1 \\ 1 - x(\mathbf{z}^{(n)}; \mathbf{w}) & \text{if } y^{(n)} = 0 \end{cases}$$

computationally and
mathematically simpler
to work with log-likelihoods

$$\mathbf{w}_{\text{ML}} = \arg \max_{\mathbf{w}} \log p(\{y^{(n)}\}_{n=1}^N | \{\mathbf{z}^{(n)}\}_{n=1}^N, \mathbf{w})$$

monotonicity of
logarithm preserves
location of optimum



$$\log p(\{y^{(n)}\}_{n=1}^N | \{\mathbf{z}^{(n)}\}_{n=1}^N, \mathbf{w}) = \sum_{n=1}^N y^{(n)} \log x(\mathbf{z}^{(n)}; \mathbf{w}) + (1 - y^{(n)}) \log(1 - x(\mathbf{z}^{(n)}; \mathbf{w})) = -G(\mathbf{w})$$

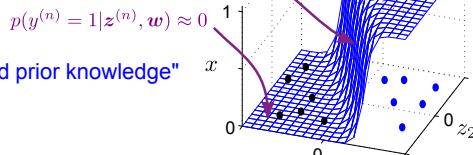
Probabilistic interpretation of the single neuron

Interpret the output of the network as a probability:

$$x(\mathbf{z}^{(n)}; \mathbf{w}) = p(y^{(n)} = 1 | \mathbf{z}^{(n)}, \mathbf{w})$$

$$p(y^{(n)} = 1 | \mathbf{z}^{(n)}, \mathbf{w}) \approx 1$$

Fitting method 2: maximum a posteriori (MAP)



"choose setting of w that is most probable given the data and prior knowledge"

$$\mathbf{w}_{\text{MAP}} = \arg \max_{\mathbf{w}} p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N, \{y^{(n)}\}_{n=1}^N)$$

Bayes' rule $p(A|B,C) = \frac{1}{p(B|C)} p(A|C)p(B|A,C)$ with $A = \mathbf{w}$, $B = \{y^{(n)}\}_{n=1}^N$ and $C = \{\mathbf{z}^{(n)}\}_{n=1}^N$

$$p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N, \{y^{(n)}\}_{n=1}^N) = \frac{1}{p(\{y^{(n)}\}_{n=1}^N | \{\mathbf{z}^{(n)}\}_{n=1}^N)} p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N) p(\{y^{(n)}\}_{n=1}^N | \{\mathbf{z}^{(n)}\}_{n=1}^N, \mathbf{w})$$

posterior distribution
(what we know after seeing data) normalising constant prior distribution
(what we knew before seeing data) likelihood of the parameters
(what the data tell us)

use a zero mean, independent Gaussian prior $p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N) = \frac{1}{(2\pi\sigma_w^2)^{D/2}} \exp(-\frac{1}{2\sigma_w^2} \sum_{d=1}^D w_d^2)$

$$\begin{aligned} \mathbf{w}_{\text{MAP}} &= \arg \max_{\mathbf{w}} \log p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N, \{y^{(n)}\}_{n=1}^N) \\ &= \arg \max_{\mathbf{w}} \left[\log p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N) + \log p(\{y^{(n)}\}_{n=1}^N | \{\mathbf{z}^{(n)}\}_{n=1}^N, \mathbf{w}) - \log p(\{y^{(n)}\}_{n=1}^N | \{\mathbf{z}^{(n)}\}_{n=1}^N) \right] \end{aligned}$$

does not depend on w

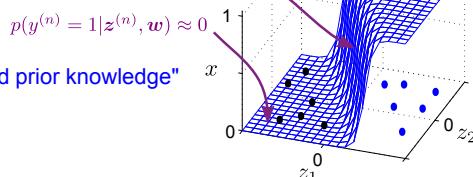
Probabilistic interpretation of the single neuron

Interpret the output of the network as a probability:

$$x(\mathbf{z}^{(n)}; \mathbf{w}) = p(y^{(n)} = 1 | \mathbf{z}^{(n)}, \mathbf{w})$$

$$p(y^{(n)} = 1 | \mathbf{z}^{(n)}, \mathbf{w}) \approx 1$$

Fitting method 2: maximum a posteriori (MAP)



"choose setting of w that is most probable given the data and prior knowledge"

$$\mathbf{w}_{\text{MAP}} = \arg \max_{\mathbf{w}} p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N, \{y^{(n)}\}_{n=1}^N)$$

Bayes' rule $p(A|B,C) = \frac{1}{p(B|C)} p(A|C)p(B|A,C)$ with $A = \mathbf{w}$, $B = \{y^{(n)}\}_{n=1}^N$ and $C = \{\mathbf{z}^{(n)}\}_{n=1}^N$

$$p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N, \{y^{(n)}\}_{n=1}^N) = \frac{1}{p(\{y^{(n)}\}_{n=1}^N | \{\mathbf{z}^{(n)}\}_{n=1}^N)} p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N) p(\{y^{(n)}\}_{n=1}^N | \{\mathbf{z}^{(n)}\}_{n=1}^N, \mathbf{w})$$

posterior distribution
(what we know after seeing data) normalising constant prior distribution
(what we knew before seeing data) likelihood of the parameters
(what the data tell us)

use a zero mean, independent Gaussian prior $p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N) = \frac{1}{(2\pi\sigma_w^2)^{D/2}} \exp(-\frac{1}{2\sigma_w^2} \sum_{d=1}^D w_d^2)$

$$\begin{aligned} \mathbf{w}_{\text{MAP}} &= \arg \max_{\mathbf{w}} \log p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N, \{y^{(n)}\}_{n=1}^N) \\ &= \arg \max_{\mathbf{w}} \left[\log p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N) + \log p(\{y^{(n)}\}_{n=1}^N | \{\mathbf{z}^{(n)}\}_{n=1}^N, \mathbf{w}) - \log p(\{y^{(n)}\}_{n=1}^N | \{\mathbf{z}^{(n)}\}_{n=1}^N) \right] \end{aligned}$$

$\curvearrowleft = -G(\mathbf{w})$

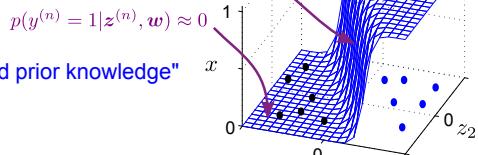
Probabilistic interpretation of the single neuron

Interpret the output of the network as a probability:

$$x(\mathbf{z}^{(n)}; \mathbf{w}) = p(y^{(n)} = 1 | \mathbf{z}^{(n)}, \mathbf{w})$$

$$p(y^{(n)} = 1 | \mathbf{z}^{(n)}, \mathbf{w}) \approx 1$$

Fitting method 2: maximum a posteriori (MAP)



"choose setting of w that is most probable given the data and prior knowledge"

$$\mathbf{w}_{\text{MAP}} = \arg \max_{\mathbf{w}} p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N, \{y^{(n)}\}_{n=1}^N)$$

Bayes' rule $p(A|B,C) = \frac{1}{p(B|C)} p(A|C)p(B|A,C)$ with $A = \mathbf{w}$, $B = \{y^{(n)}\}_{n=1}^N$ and $C = \{\mathbf{z}^{(n)}\}_{n=1}^N$

$$p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N, \{y^{(n)}\}_{n=1}^N) = \frac{1}{p(\{y^{(n)}\}_{n=1}^N | \{\mathbf{z}^{(n)}\}_{n=1}^N)} p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N) p(\{y^{(n)}\}_{n=1}^N | \{\mathbf{z}^{(n)}\}_{n=1}^N, \mathbf{w})$$

posterior distribution
(what we know after seeing data) normalising constant prior distribution
(what we knew before seeing data) likelihood of the parameters
(what the data tell us)

use a zero mean, independent Gaussian prior $p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N) = \frac{1}{(2\pi\sigma_w^2)^{D/2}} \exp(-\frac{1}{2\sigma_w^2} \sum_{d=1}^D w_d^2)$

$$\begin{aligned} \mathbf{w}_{\text{MAP}} &= \arg \max_{\mathbf{w}} \log p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N, \{y^{(n)}\}_{n=1}^N) \\ &= \arg \max_{\mathbf{w}} \left[\log p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N) - G(\mathbf{w}) \right] \end{aligned}$$

substitute in log prior

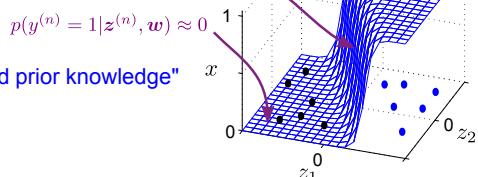
Probabilistic interpretation of the single neuron

Interpret the output of the network as a probability:

$$x(\mathbf{z}^{(n)}; \mathbf{w}) = p(y^{(n)} = 1 | \mathbf{z}^{(n)}, \mathbf{w})$$

$$p(y^{(n)} = 1 | \mathbf{z}^{(n)}, \mathbf{w}) \approx 1$$

Fitting method 2: maximum a posteriori (MAP)



"choose setting of w that is most probable given the data and prior knowledge"

$$\mathbf{w}_{\text{MAP}} = \arg \max_{\mathbf{w}} p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N, \{y^{(n)}\}_{n=1}^N)$$

Bayes' rule $p(A|B,C) = \frac{1}{p(B|C)} p(A|C)p(B|A,C)$ with $A = \mathbf{w}$, $B = \{y^{(n)}\}_{n=1}^N$ and $C = \{\mathbf{z}^{(n)}\}_{n=1}^N$

$$p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N, \{y^{(n)}\}_{n=1}^N) = \frac{1}{p(\{y^{(n)}\}_{n=1}^N | \{\mathbf{z}^{(n)}\}_{n=1}^N)} p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N) p(\{y^{(n)}\}_{n=1}^N | \{\mathbf{z}^{(n)}\}_{n=1}^N, \mathbf{w})$$

posterior distribution
(what we know after seeing data) normalising constant prior distribution
(what we knew before seeing data) likelihood of the parameters
(what the data tell us)

use a zero mean, independent Gaussian prior $p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N) = \frac{1}{(2\pi\sigma_w^2)^{D/2}} \exp(-\frac{1}{2\sigma_w^2} \sum_{d=1}^D w_d^2)$

$$\begin{aligned} \mathbf{w}_{\text{MAP}} &= \arg \max_{\mathbf{w}} \log p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N, \{y^{(n)}\}_{n=1}^N) \\ &= \arg \max_{\mathbf{w}} \left[-\frac{D}{2} \log(2\pi\sigma_w^2) - \frac{1}{2\sigma_w^2} \sum_{d=1}^D w_d^2 - G(\mathbf{w}) \right] \end{aligned}$$

\mathbf{w} does not depend on w

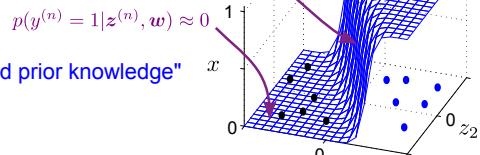
Probabilistic interpretation of the single neuron

Interpret the output of the network as a probability:

$$x(\mathbf{z}^{(n)}; \mathbf{w}) = p(y^{(n)} = 1 | \mathbf{z}^{(n)}, \mathbf{w})$$

$$p(y^{(n)} = 1 | \mathbf{z}^{(n)}, \mathbf{w}) \approx 1$$

Fitting method 2: maximum a posteriori (MAP)



"choose setting of w that is most probable given the data and prior knowledge"

$$\mathbf{w}_{\text{MAP}} = \arg \max_{\mathbf{w}} p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N, \{y^{(n)}\}_{n=1}^N)$$

Bayes' rule $p(A|B,C) = \frac{1}{p(B|C)} p(A|C)p(B|A,C)$ with $A = \mathbf{w}$, $B = \{y^{(n)}\}_{n=1}^N$ and $C = \{\mathbf{z}^{(n)}\}_{n=1}^N$

$$p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N, \{y^{(n)}\}_{n=1}^N) = \frac{1}{p(\{y^{(n)}\}_{n=1}^N | \{\mathbf{z}^{(n)}\}_{n=1}^N)} p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N) p(\{y^{(n)}\}_{n=1}^N | \{\mathbf{z}^{(n)}\}_{n=1}^N, \mathbf{w})$$

posterior distribution
(what we know after seeing data) normalising constant prior distribution
(what we knew before seeing data) likelihood of the parameters
(what the data tell us)

use a zero mean, independent Gaussian prior $p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N) = \frac{1}{(2\pi\sigma_w^2)^{D/2}} \exp(-\frac{1}{2\sigma_w^2} \sum_{d=1}^D w_d^2)$

$$\begin{aligned} \mathbf{w}_{\text{MAP}} &= \arg \max_{\mathbf{w}} \log p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N, \{y^{(n)}\}_{n=1}^N) \\ &= \arg \max_{\mathbf{w}} \left[-\frac{1}{2\sigma_w^2} \sum_{d=1}^D w_d^2 - G(\mathbf{w}) \right] = \alpha E(\mathbf{w}) \end{aligned}$$

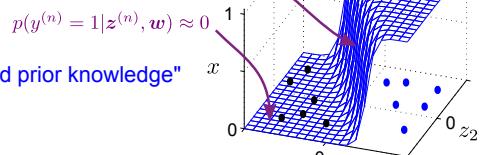
Probabilistic interpretation of the single neuron

Interpret the output of the network as a probability:

$$x(\mathbf{z}^{(n)}; \mathbf{w}) = p(y^{(n)} = 1 | \mathbf{z}^{(n)}, \mathbf{w})$$

$$p(y^{(n)} = 1 | \mathbf{z}^{(n)}, \mathbf{w}) \approx 1$$

Fitting method 2: maximum a posteriori (MAP)



"choose setting of w that is most probable given the data and prior knowledge"

$$\mathbf{w}_{\text{MAP}} = \arg \max_{\mathbf{w}} p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N, \{y^{(n)}\}_{n=1}^N)$$

Bayes' rule $p(A|B,C) = \frac{1}{p(B|C)} p(A|C)p(B|A,C)$ with $A = \mathbf{w}$, $B = \{y^{(n)}\}_{n=1}^N$ and $C = \{\mathbf{z}^{(n)}\}_{n=1}^N$

$$p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N, \{y^{(n)}\}_{n=1}^N) = \frac{1}{p(\{y^{(n)}\}_{n=1}^N | \{\mathbf{z}^{(n)}\}_{n=1}^N)} p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N) p(\{y^{(n)}\}_{n=1}^N | \{\mathbf{z}^{(n)}\}_{n=1}^N, \mathbf{w})$$

posterior distribution
(what we know after seeing data) normalising constant prior distribution
(what we knew before seeing data) likelihood of the parameters
(what the data tell us)

use a zero mean, independent Gaussian prior $p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N) = \frac{1}{(2\pi\sigma_w^2)^{D/2}} \exp(-\frac{1}{2\sigma_w^2} \sum_{d=1}^D w_d^2)$

$$\begin{aligned} \mathbf{w}_{\text{MAP}} &= \arg \max_{\mathbf{w}} \log p(\mathbf{w} | \{\mathbf{z}^{(n)}\}_{n=1}^N, \{y^{(n)}\}_{n=1}^N) \\ &= \arg \max_{\mathbf{w}} \left[-\alpha E(\mathbf{w}) - G(\mathbf{w}) \right] \quad \text{where } \alpha = \frac{1}{\sigma_w^2} \end{aligned}$$

Machine learning: rebranding statistics?

Machine Learning terminology

single neuron for classification

$$x(z^{(n)}; w)$$

relative entropy (data fit term)

$$G(w)$$

regularisation (weight decay term)

$$\alpha E(w)$$

minimise of relative entropy

$$\arg \min_w G(w)$$

minimise relative entropy + regulariser

$$\arg \min_w [G(w) + \alpha E(w)]$$

Statistical terminology

logistic regression or logistic classification

$$p(y^{(n)} = 1 | z^{(n)}, w)$$

log-likelihood function

$$\log p(\{y^{(n)}\}_{n=1}^N | \{z^{(n)}\}_{n=1}^N, w)$$

regularisation (log-prior)

$$\log p(w | \{z^{(n)}\}_{n=1}^N)$$

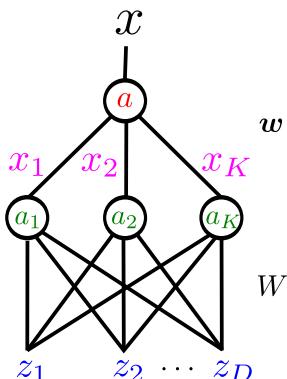
maximum likelihood learning

$$\arg \max_w \log p(\{y^{(n)}\}_{n=1}^N | \{z^{(n)}\}_{n=1}^N, w)$$

maximum a posteriori fitting

$$\arg \max_w p(w | \{z^{(n)}\}_{n=1}^N, \{y^{(n)}\}_{n=1}^N)$$

Single hidden layer neural networks



output of network

$$x(a) = \frac{1}{1 + \exp(-a)}$$

activity of output neuron

$$a = \sum_{k=1}^K w_k x_k$$

hidden neuron outputs

$$x(a_k) = \frac{1}{1 + \exp(-a_k)} = x_k$$

hidden neuron activities

$$a_k = \sum_{d=1}^D W_{k,d} z_d$$

inputs

output weights

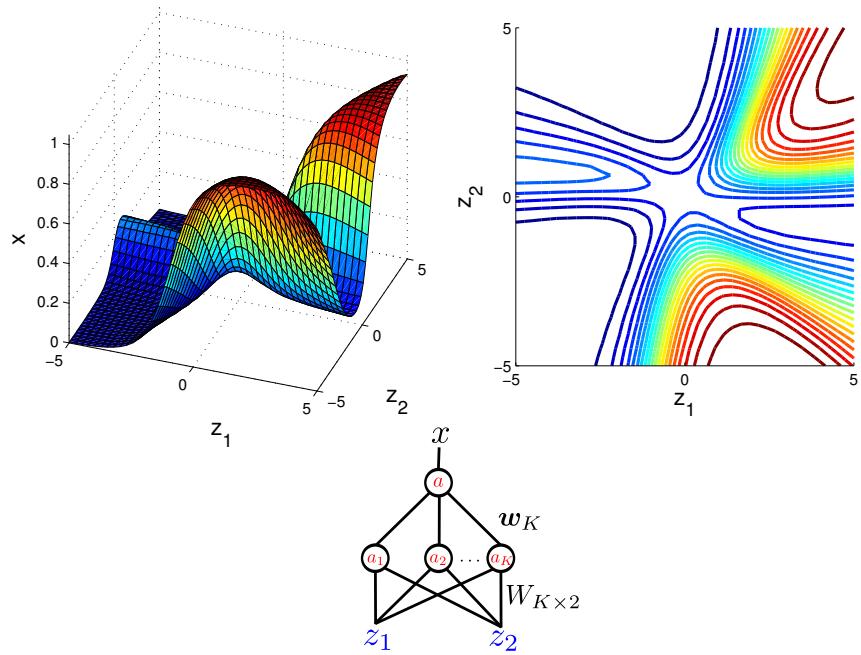
dimensionality(w) = K

input weights

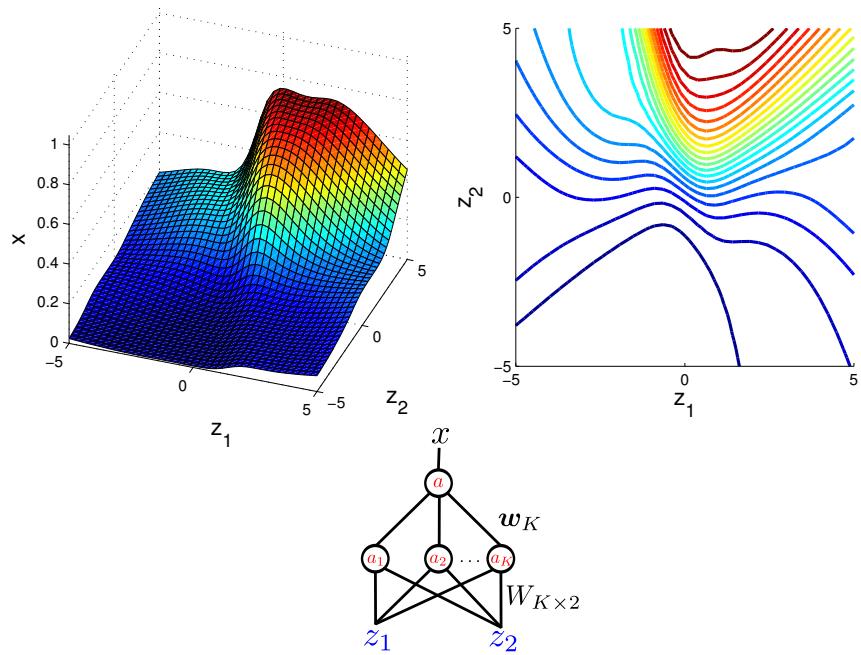
dimensionality(W) = $K \times D$



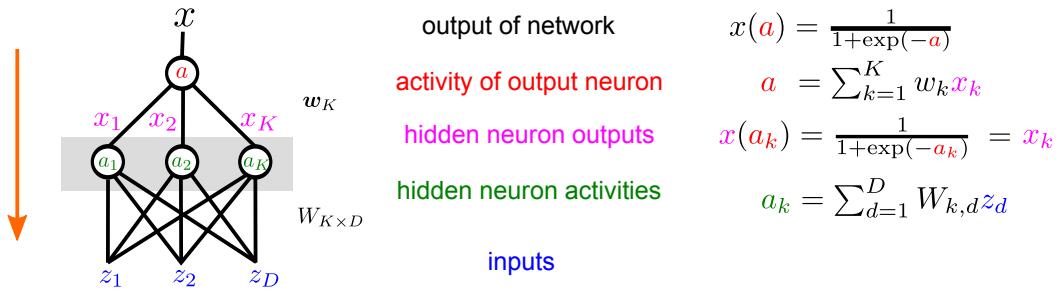
Sampling random neural network classifiers



Sampling random neural network classifiers



Backpropagation



objective function:

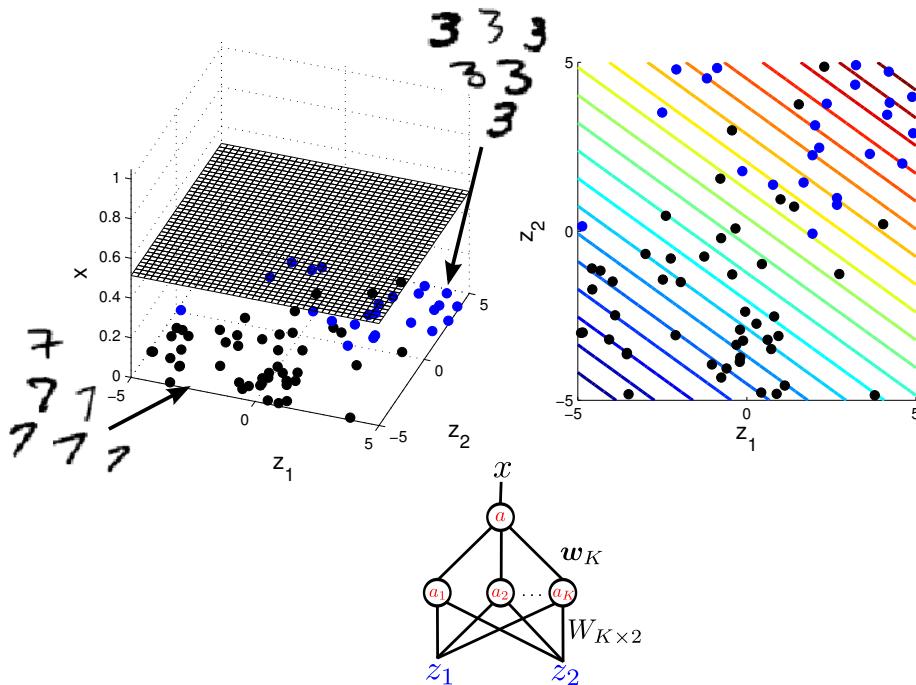
$$G(\mathbf{w}) = -\sum_n [y^{(n)} \log x^{(n)} + (1 - y^{(n)}) \log (1 - x^{(n)})] \quad \text{likelihood same as before}$$

$$E(W, \mathbf{w}) = \frac{1}{2} \sum_i w_i^2 + \frac{1}{2} \sum_{ij} W_{ij}^2 \quad \text{regulariser discourages extreme weights}$$

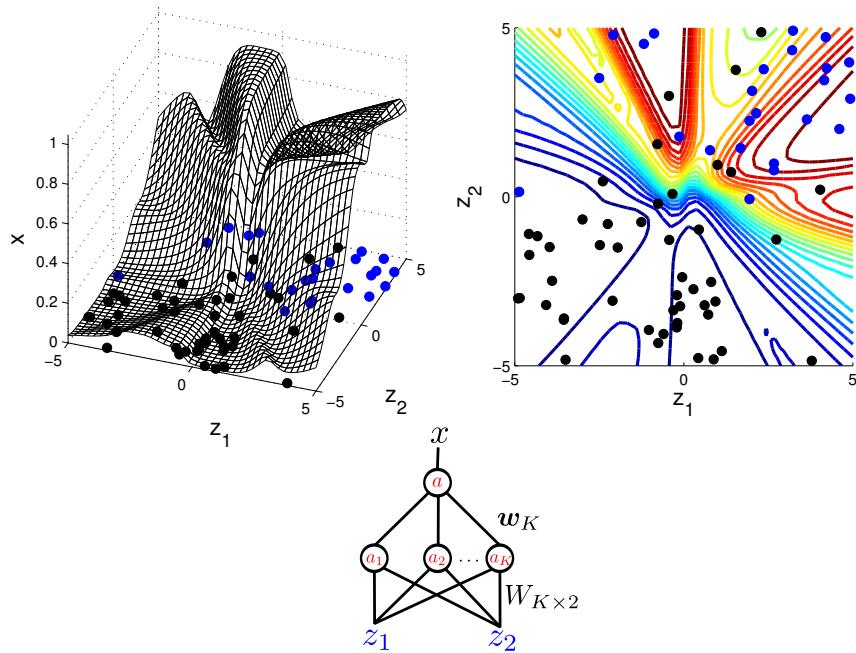
$$\begin{aligned} \{W, \mathbf{w}^*\} &= \arg \min_{W, \mathbf{w}} M(W, \mathbf{w}) = \arg \min_{W, \mathbf{w}} [G(W, \mathbf{w}) + \alpha E(W, \mathbf{w})] \\ \frac{dG(W, \mathbf{w})}{dW_{ij}} &= \sum_n \frac{dG(W, \mathbf{w})}{dx^{(n)}} \frac{dx^{(n)}}{dW_{ij}} = \sum_n \frac{dG(W, \mathbf{w})}{dx^{(n)}} \frac{dx^{(n)}}{da^{(n)}} \frac{da^{(n)}}{dW_{ij}} \\ &= \sum_n \frac{dG(W, \mathbf{w})}{dx^{(n)}} \frac{dx^{(n)}}{da^{(n)}} \frac{da^{(n)}}{dx_i^{(n)}} \frac{dx_i^{(n)}}{dW_{ij}} = \sum_n \frac{dG(W, \mathbf{w})}{dx^{(n)}} \frac{dx^{(n)}}{da^{(n)}} \frac{da^{(n)}}{dx_i^{(n)}} \frac{dx_i^{(n)}}{da_i^{(n)}} \frac{da_i^{(n)}}{dW_{ij}} \end{aligned}$$

chain rule = backward pass
through network = "backpropagation"

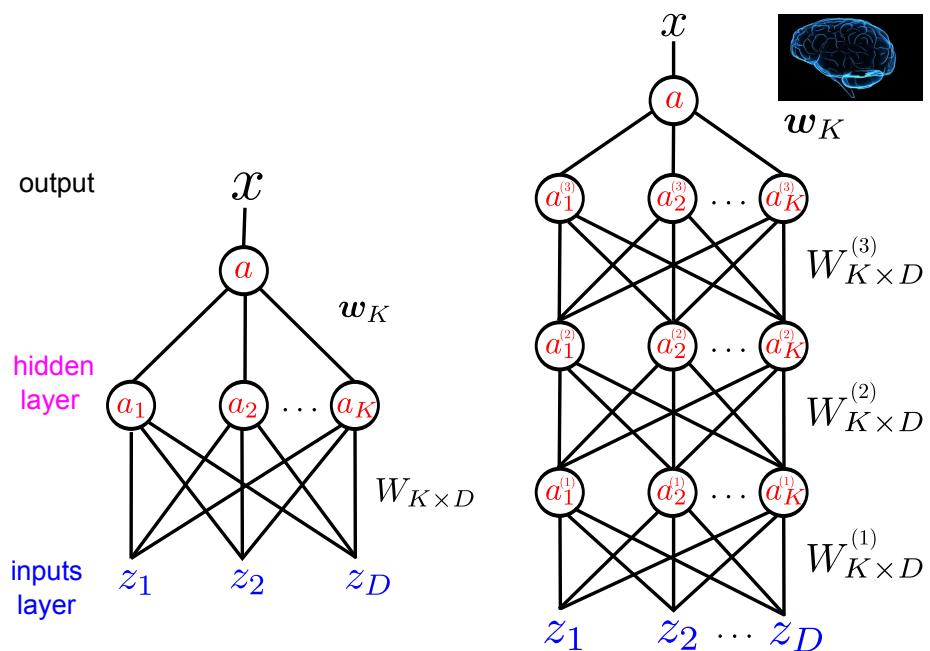
Training a neural network with a single hidden layer



Training a neural network with a single hidden layer



Multi-layer perceptrons: hierarchical models with many hidden layers



A Potted History of Neural Networks

1980s First wave of neural network development

- multi-layer perceptrons, convolutional neural networks, gradient based optimisation etc.
- limited training success caused progress to stall.

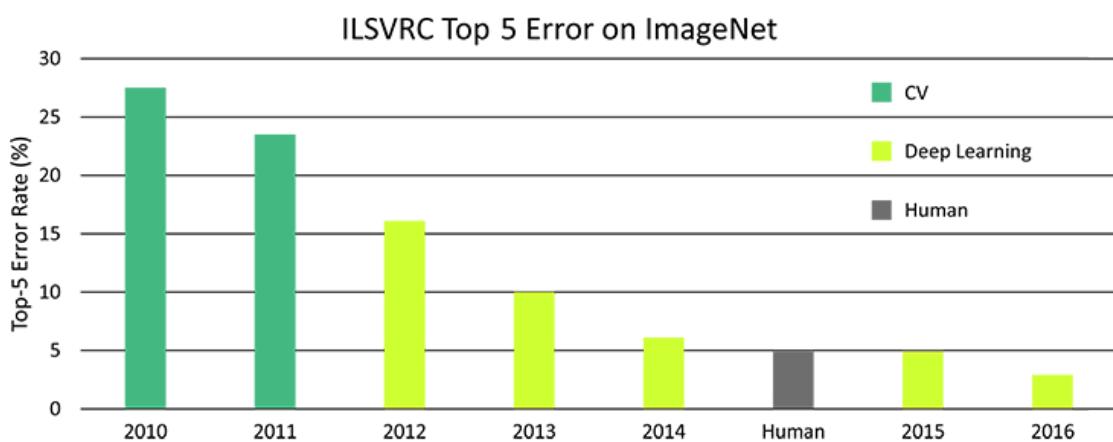
2010s Second wave of neural network development due to

1. more data (e.g. imageNet)
2. more compute (e.g. GPUs)
3. technical advances (e.g. initialisation, neuron non-linearities, model architecture)

2019 Bespoke software for designing & training neural networks

- automatic differentiation, simple deployment to GPUs etc.
- examples include PyTorch and TensorFlow
- software 2.0: neural network design = curating a programme with each layer being a step, now learn precise form from data

Progression of Performance on Object Recognition



Summary of Part I

- **Supervised artificial neural networks** fit a non-linear function x that maps from input features (z) to output targets (y), like a fancy form of curve fitting.
- Networks with several hidden layers are known as **multi-layer perceptrons** and can be fit by optimising an **objective function** by **gradient descent**.
- The objective function comprises a **data fit term** $G(w)$ that ensures the output of the network matches the training data and a **regularisation term** $E(w)$ to stop **over-fitting**. These terms can be related to the **likelihood function** and **prior** of the weights from the field of statistics.
- The gradient computation requires the chain rule and is known as **backpropagation**.
- **Smart initialisation** is required to stop learning falling into poor **local optima**.

Demo: <http://yann.lecun.com/exdb/lenet/>