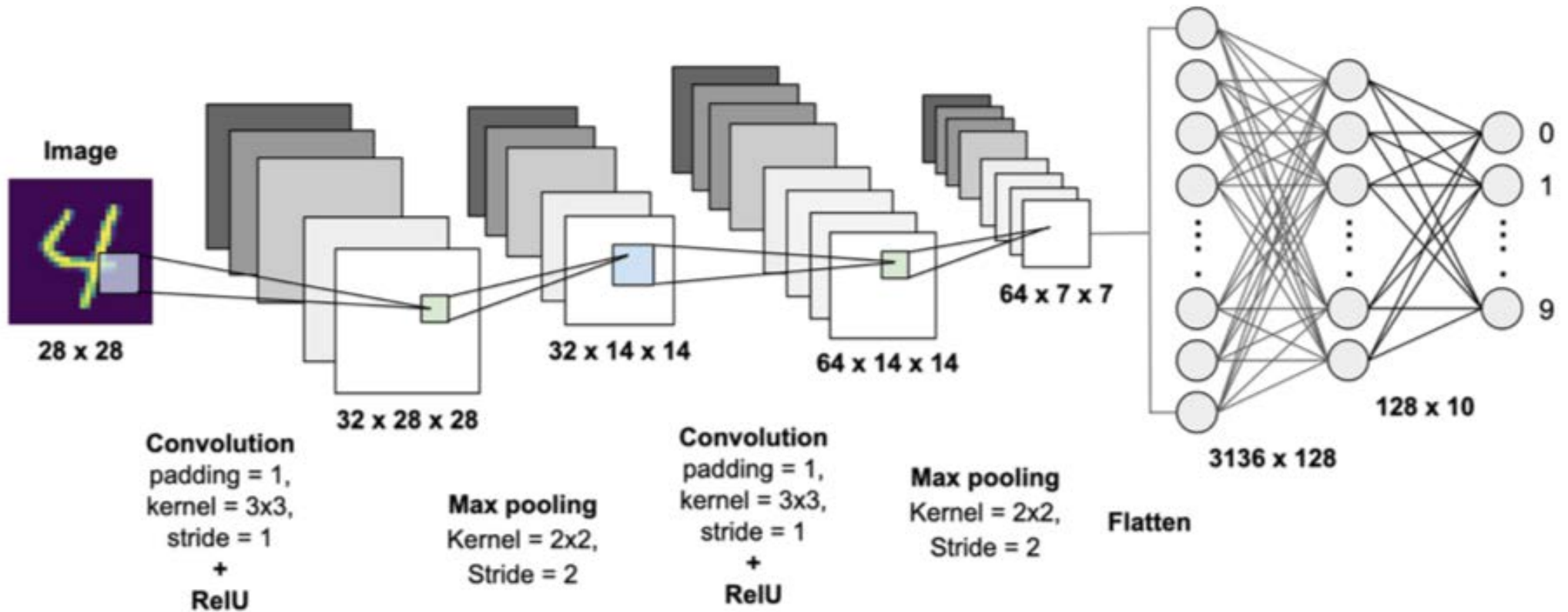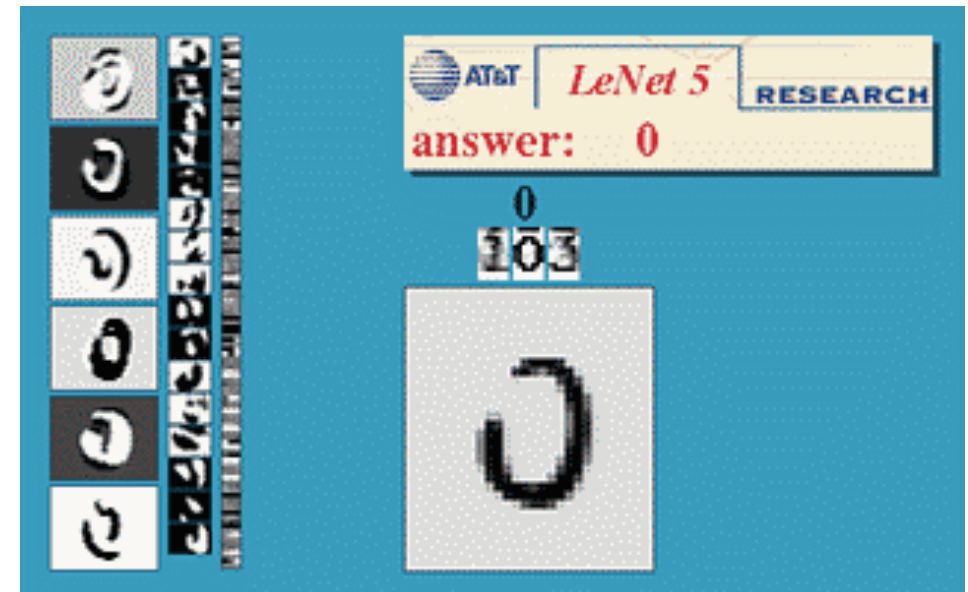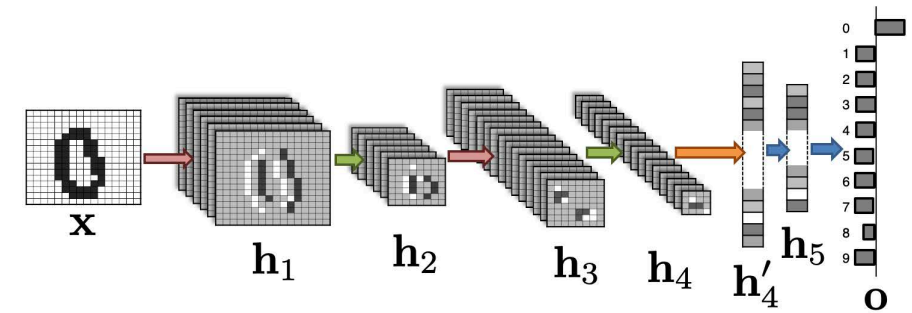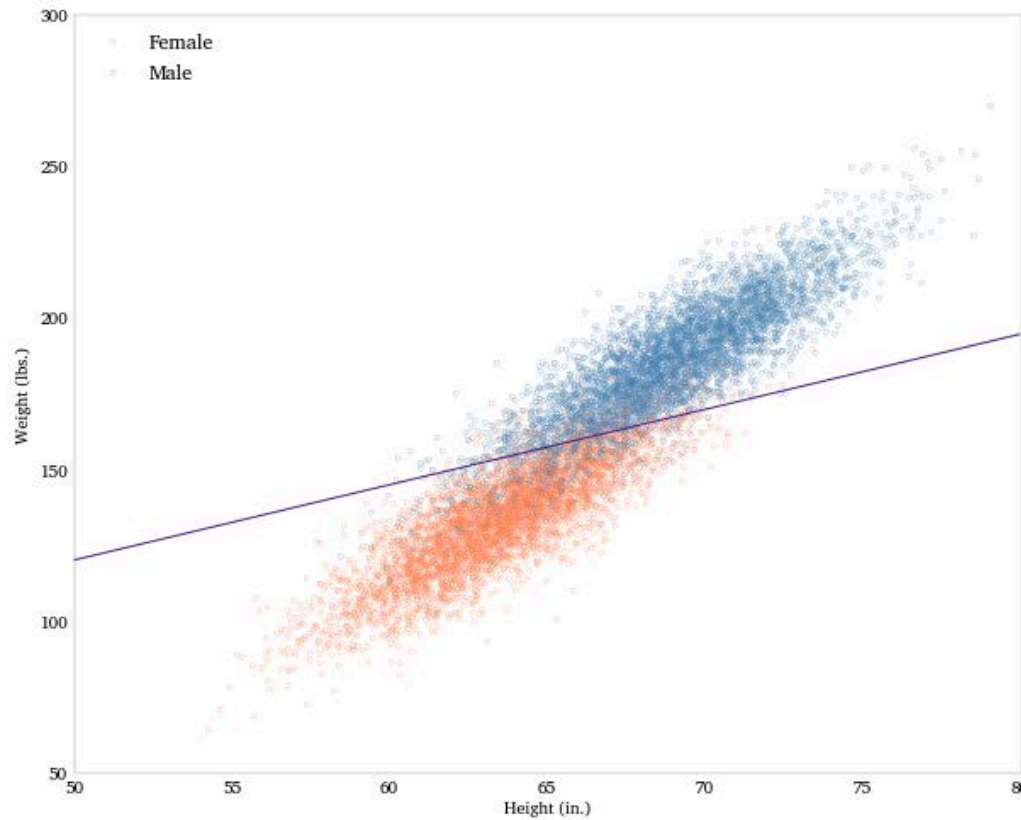# Deep Learning Crash Course



- Single Layer Perceptrons
- Multiple Layer Perceptrons
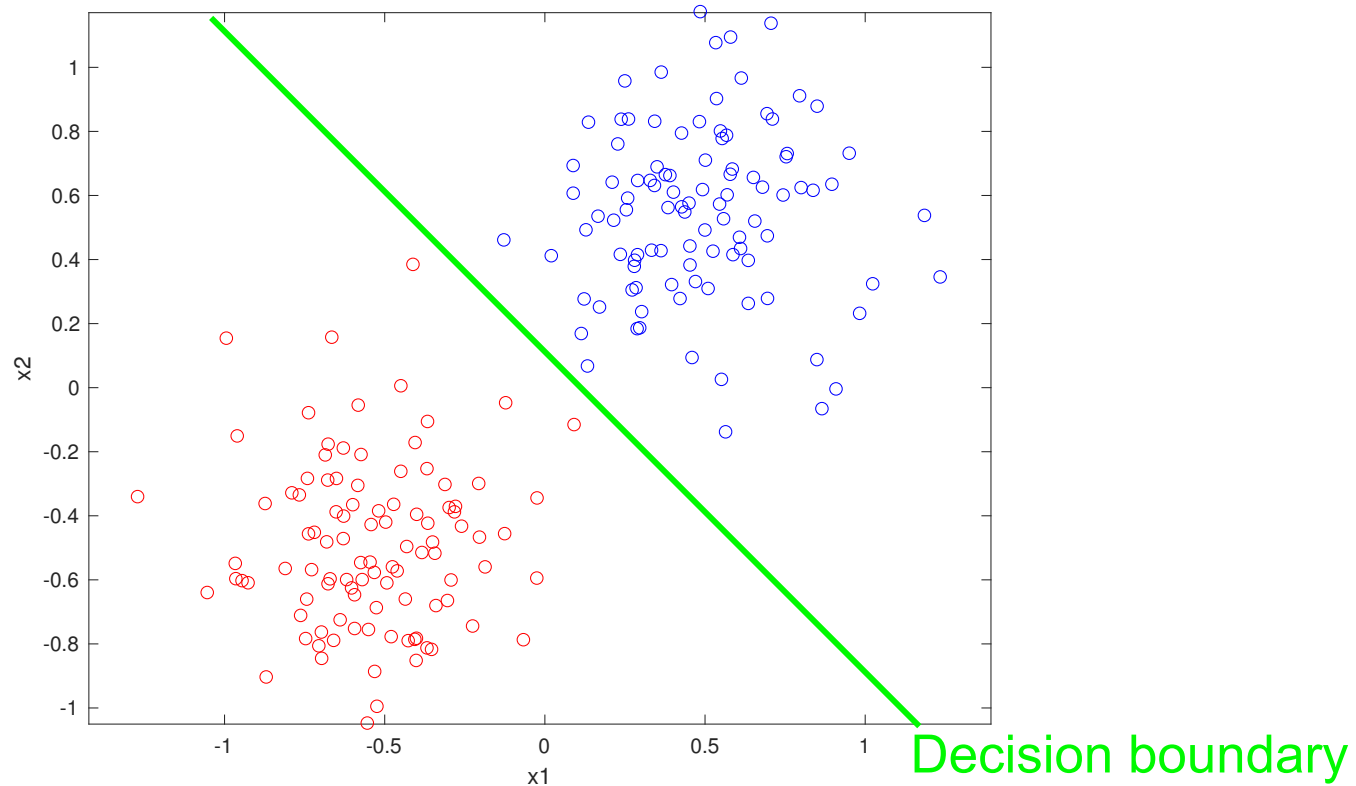- Convolutional Neural Nets

# Binary vs Multi-Class Classification



Logistic Regression

LeNet

# Linear Binary Classification



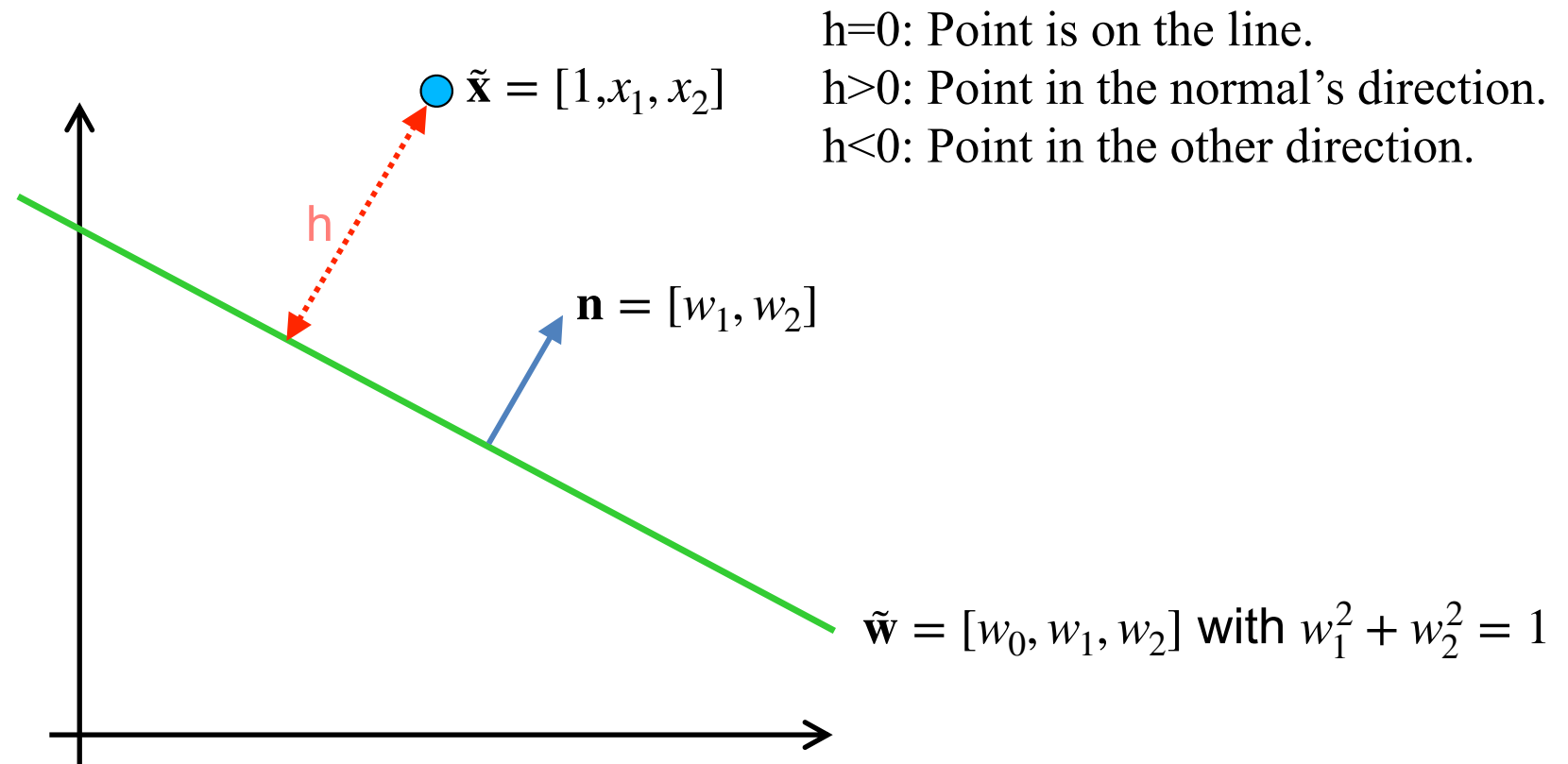Decision boundary

Two classes shown as different colors:

- The label $y \in \{-1, 1\}$ or $y \in \{0, 1\}$.
- The samples with label 1 are called positive samples.
- The samples with label -1 or 0 are called negative samples.

# Signed Distance

h=0: Point is on the line.
h>0: Point in the normal's direction.
h<0: Point in the other direction.



$\tilde{\mathbf{x}} = [1, x_1, x_2]$

h

$\mathbf{n} = [w_1, w_2]$

$\tilde{\mathbf{w}} = [w_0, w_1, w_2]$ with $w_1^2 + w_2^2 = 1$
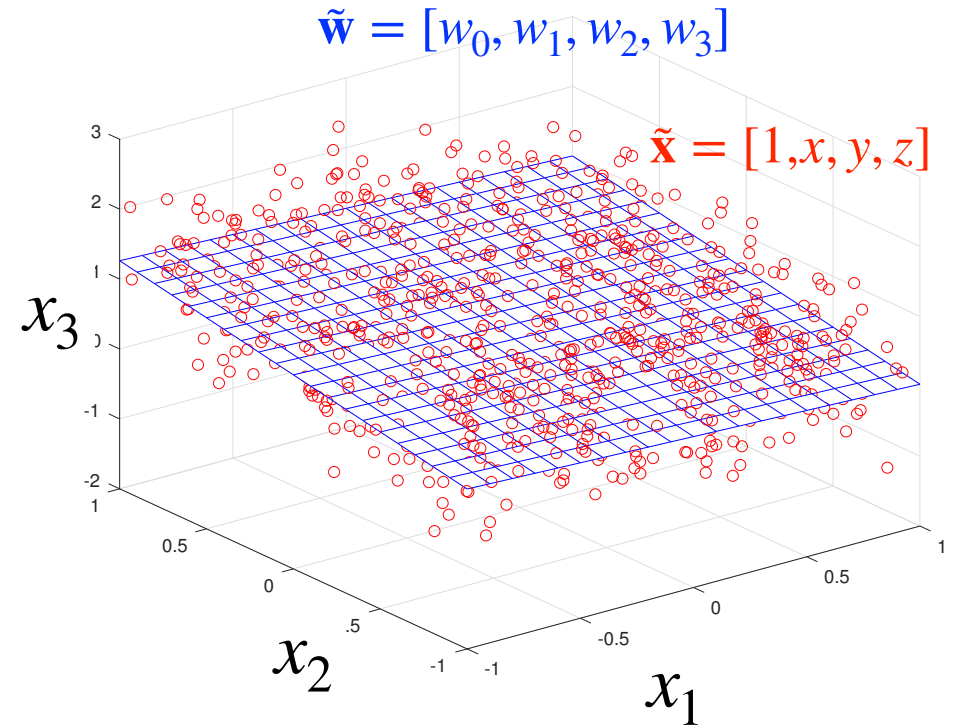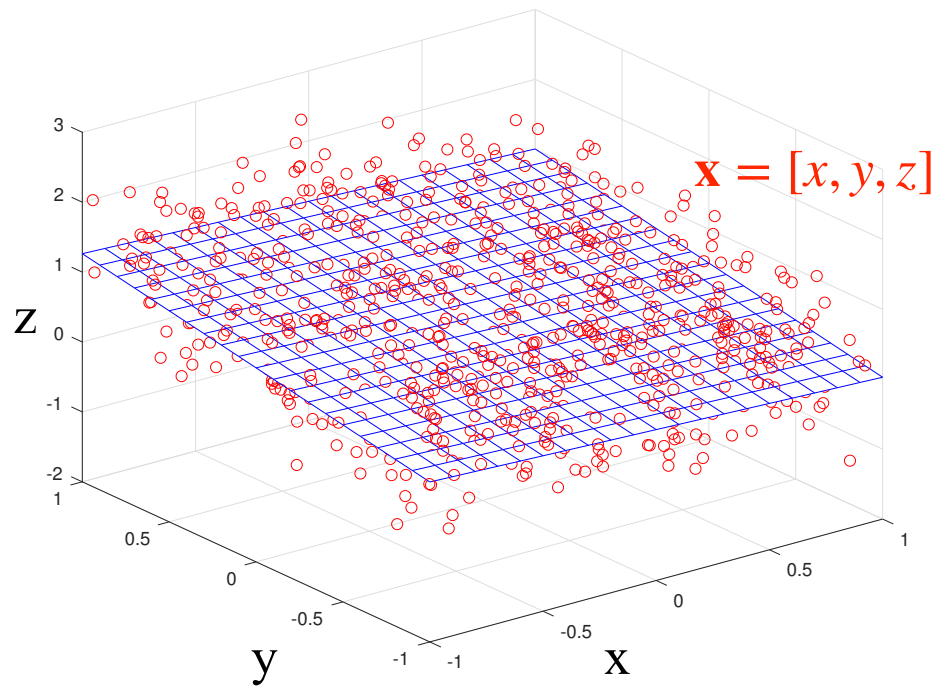
Notation:

$\mathbf{x} = [x_1, x_2]$

$\tilde{\mathbf{x}} = [1, x_1, x_2]$

Signed distance:

$h = w_0 + w_1 x_1 + w_2 x_2$

$\quad = \tilde{\mathbf{w}} \cdot \tilde{\mathbf{x}}$
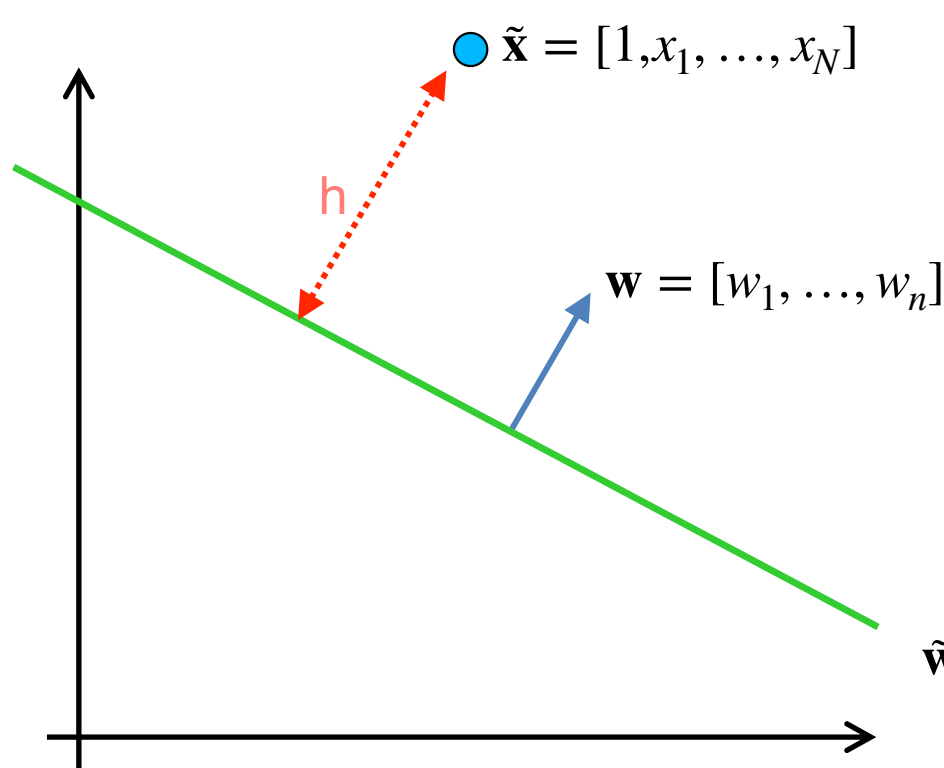
# Signed Distance in 3D



$$\mathbf{x} \in R^3, \; 0 = ax + by + cz + d$$

$$\tilde{\mathbf{x}} \in R^4, \; \tilde{\mathbf{w}} \cdot \tilde{\mathbf{x}} = 0$$

Signed distance $h = \tilde{\mathbf{w}} \cdot \tilde{\mathbf{x}}$ if $w_1^2 + w_2^2 + w_3^2 = 1$.

# Signed Distance in N Dimensions

$\tilde{\mathbf{x}} = [1, x_1, \ldots, x_N]$

$h$

$\mathbf{w} = [w_1, \ldots, w_n]$

$\tilde{\mathbf{w}} = [w_0, w_1, \ldots, w_n]$ with $\sum_{i=1}^{N} w_i^2 = 1$

h=0: Point is on the decision boundary.
h>0: Point on one side.
h<0: Point on the other side.

Notation:
$$\mathbf{x} = [x_1, \ldots, x_n]$$
$$\tilde{\mathbf{x}} = [1, x_1, \ldots, x_n]$$

Hyperplane:
$$\mathbf{x} \in R^n, \quad 0 = \tilde{\mathbf{w}} \cdot \tilde{\mathbf{x}}$$
$$= w_0 + w_1 x_1 + \ldots w_n x_n$$

Signed distance:
$$h = \tilde{\mathbf{w}} \cdot \tilde{\mathbf{x}}$$
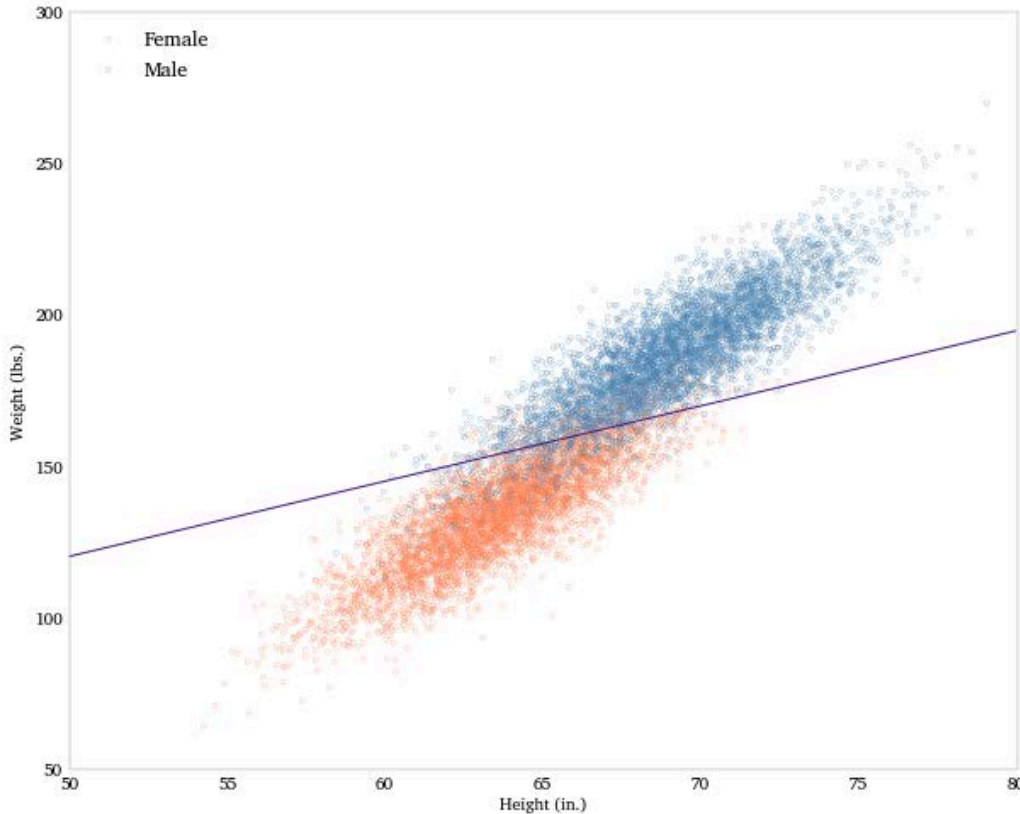
# Binary Classification in N Dimensions

**Hyperplane:** $\mathbf{x} \in R^N$, $\tilde{\mathbf{w}} \cdot \tilde{\mathbf{x}} = 0$, with $\tilde{\mathbf{x}} = [1 \,|\, \mathbf{x}]$.

**Signed distance:** $\tilde{\mathbf{w}} \cdot \tilde{\mathbf{x}}$, with $\tilde{\mathbf{w}} = [w_0 \,|\, \mathbf{w}]$ and $||\mathbf{w}|| = 1$.
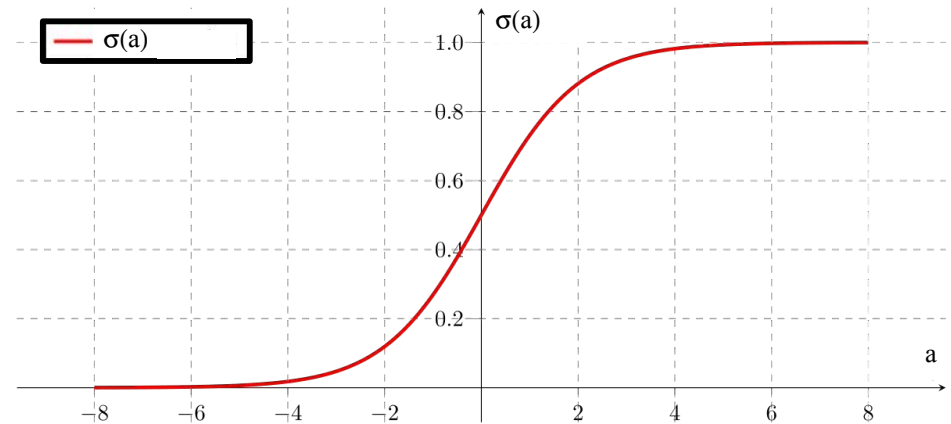
**Problem statement:** Find $\tilde{\mathbf{w}}$ such that

- for all or most positive samples $\tilde{\mathbf{w}} \cdot \tilde{\mathbf{x}} > 0,$
- for all or most negative samples $\tilde{\mathbf{w}} \cdot \tilde{\mathbf{x}} < 0.$

# Logistic Regression



$$y(\mathbf{x}; \tilde{\mathbf{w}}) = \sigma(\tilde{\mathbf{w}} \cdot \tilde{\mathbf{x}})$$

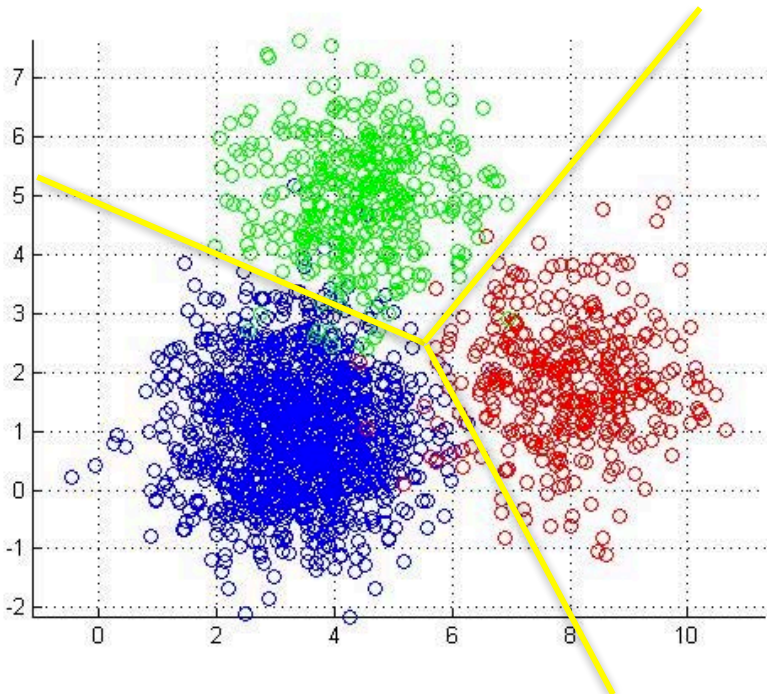$$= \frac{1}{1 + \exp(-\tilde{\mathbf{w}} \cdot \tilde{\mathbf{x}})}$$



Given a **training** set $\{(\mathbf{x}_n, t_n)_{1 \leq n \leq N}\}$ minimize

$$-\sum_n (t_n \ln y(\mathbf{x}_n) + (1 - t_n)\ln(1 - y(\mathbf{x}_n))$$

with respect to $\tilde{\mathbf{w}}$.

- When the noise is Gaussian, this is the maximum likelihood solution.
- $y(\mathbf{x}; \tilde{\mathbf{w}})$ can be interpreted at the probability that $\mathbf{x}$ belongs to positive class.

# Multi-Class Logistic Regression
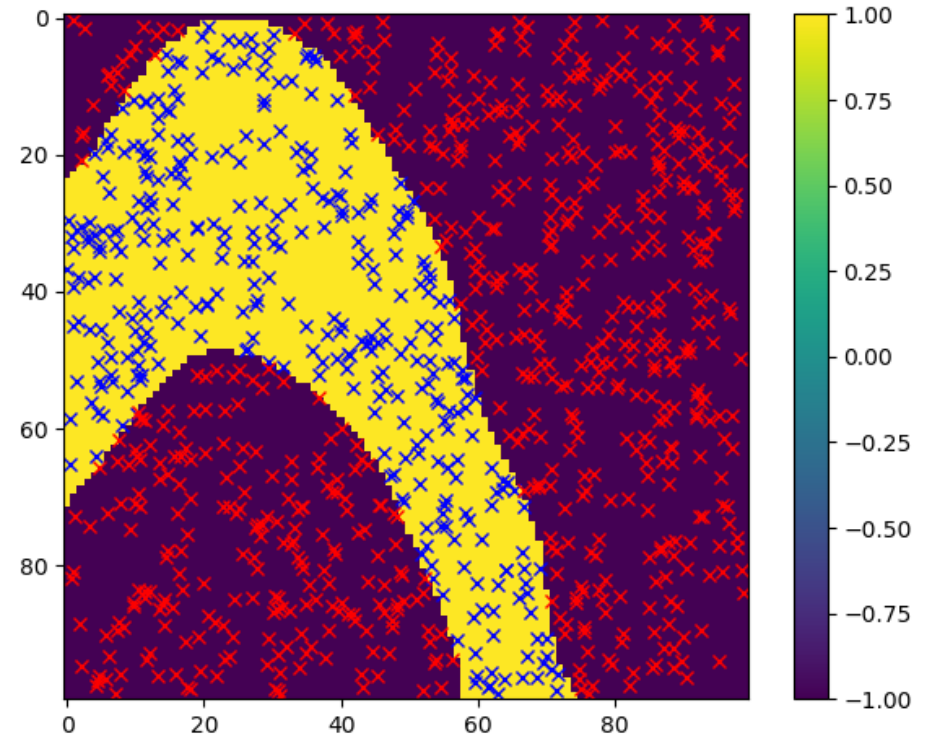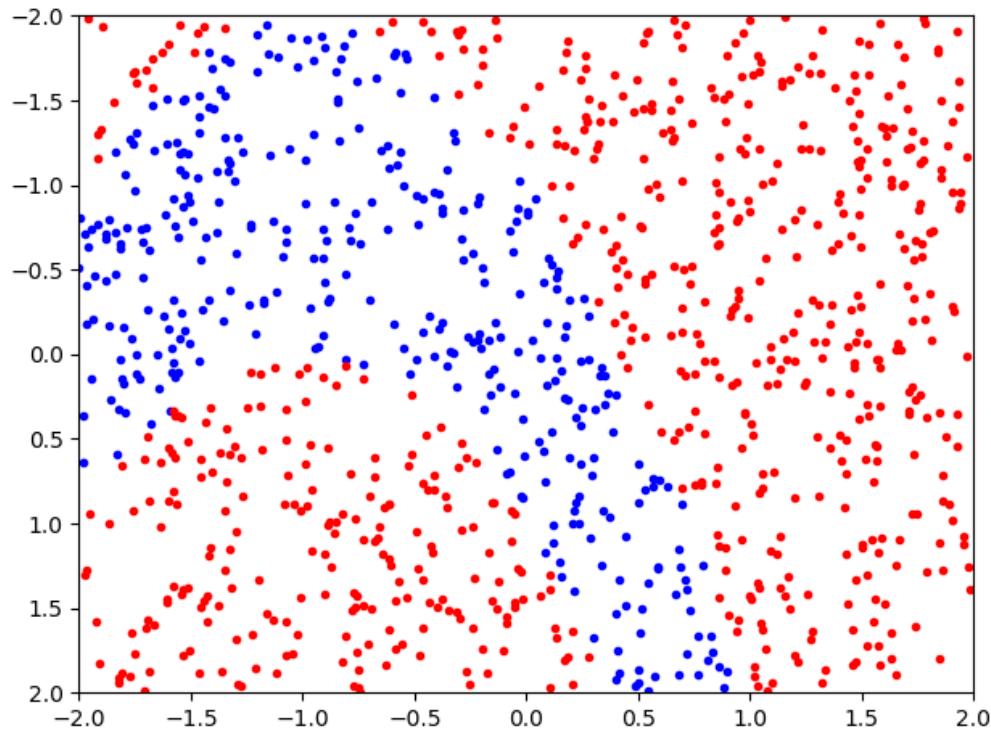


$$k = \arg\max_j y_k(\mathbf{x})$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_K \end{bmatrix} = \begin{bmatrix} \tilde{\mathbf{w}}_1^T \\ \vdots \\ \tilde{\mathbf{w}}_K^T \end{bmatrix} \tilde{\mathbf{x}}$$

$$k = \arg\max_j y_j$$

- K linear classifiers of the form $y^k(\mathbf{x}) = \sigma(\mathbf{w}_k^T \mathbf{x})$.

- Assign x to class k if $y^k(\mathbf{x}) > y^l(\mathbf{x}) \, \forall l \neq k$.

- Because the sigmoid function is monotonic, the formulation is almost unchanged.
- Only the objective function being minimized need to be reformulated.

# Non Separable Distribution



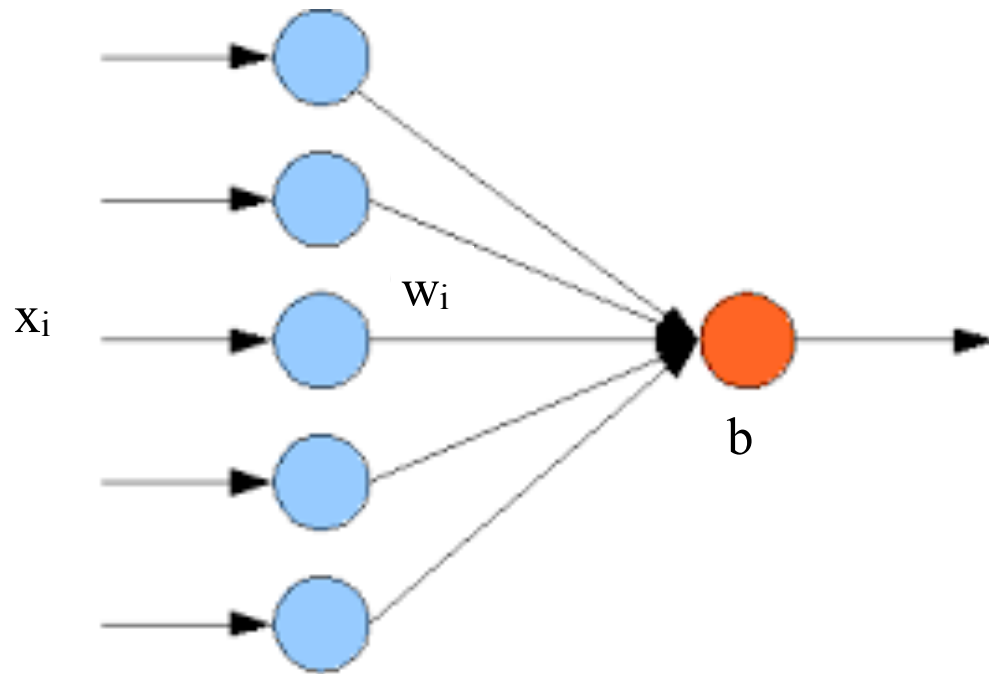Positive: $100(x_2 - x_1^2)^2 + (1 - x_1)^2 < 0.5$

Negative: Otherwise

$y(\mathbf{x}; \tilde{\mathbf{w}})$ must be a non-linear function.

- Logistic regression can handle a few outliers but not a complex non-linear boundary.

- How can we learn a function y such that $y(\mathbf{x}; \tilde{\mathbf{w}})$ is close to 1 for positive samples and close to 0 or -1 for negative ones?
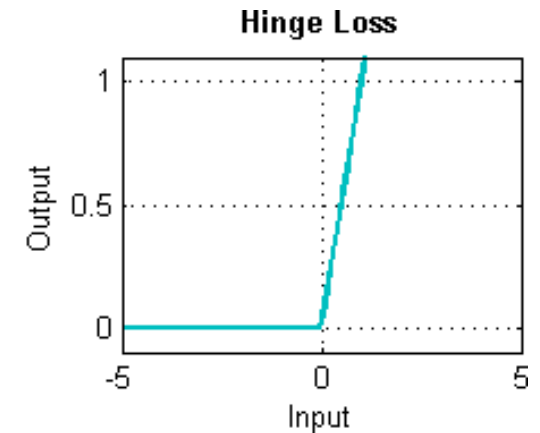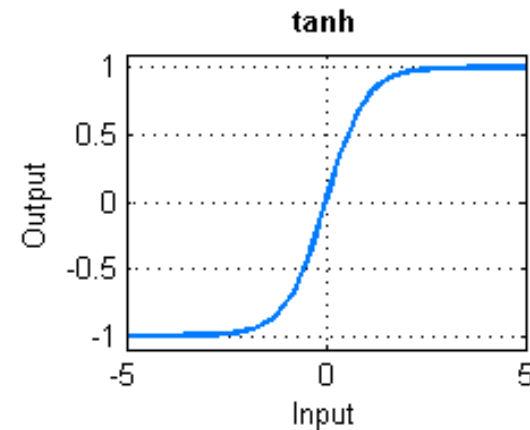
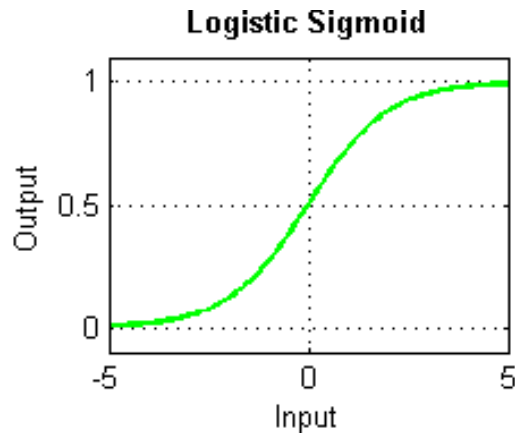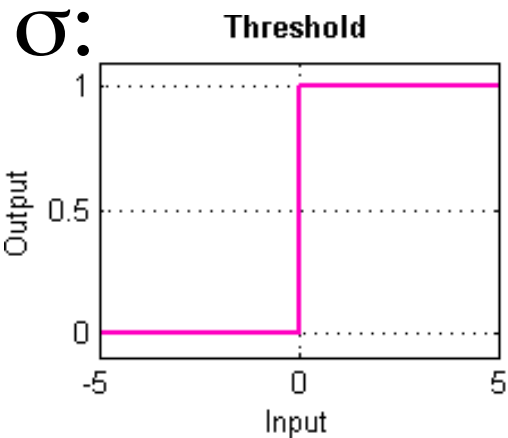$\longrightarrow$ Use LOTS of hyperplanes.

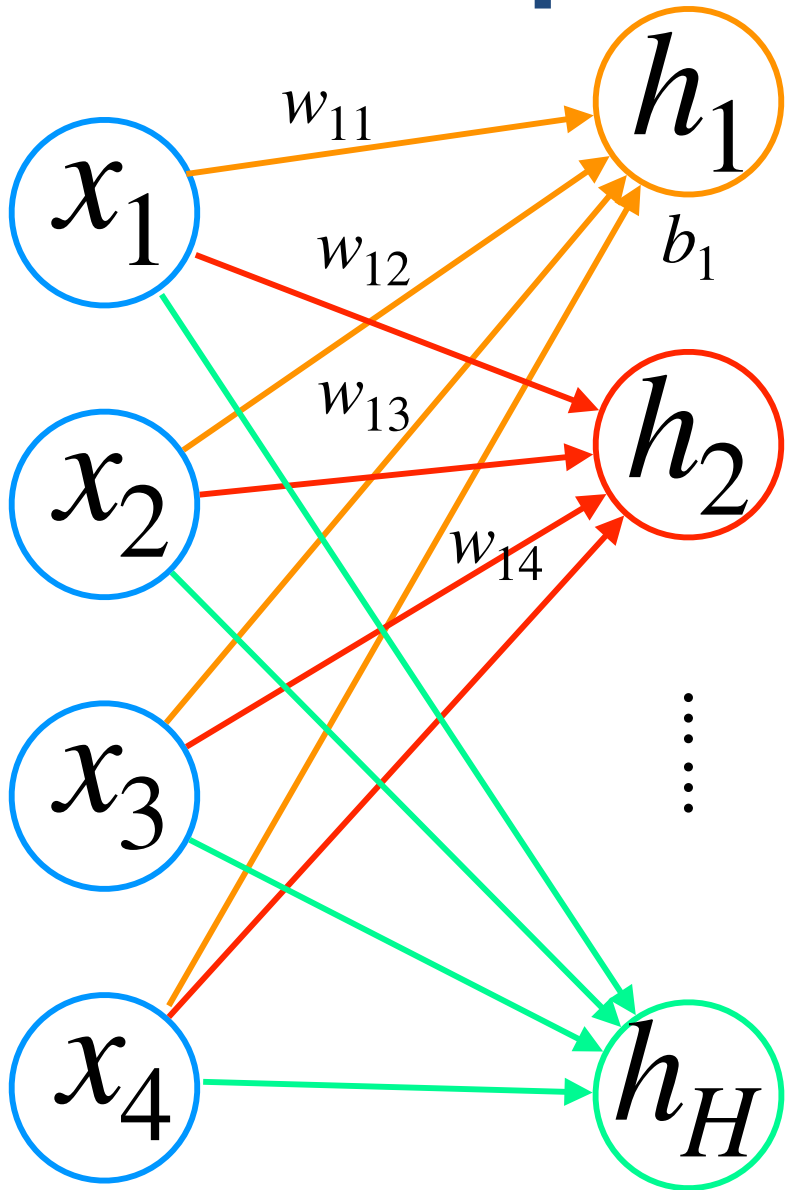# Reformulating Logistic Regression



$$y(\mathbf{x}) = \sigma(\mathbf{w} \cdot \mathbf{x} + b)$$

$$\mathbf{x} = \left[x_1, x_2, \ldots, x_n\right]^T$$

$$\mathbf{w} = \left[w_1, w_2, \ldots, w_n\right]^T$$

$\sigma$:

# Repeating the Process



$$h_1 = \sigma(\mathbf{w}_1 \cdot \mathbf{x} + b_1)$$
$$\mathbf{w}_1 = \begin{bmatrix} w_{11}, w_{12}, w_{13}, w_{14} \end{bmatrix}^T$$

$$h_2 = \sigma(\mathbf{w}_2 \cdot \mathbf{x} + b_2)$$
$$\mathbf{w}_2 = \begin{bmatrix} w_{21}, w_{22}, w_{23}, w_{24} \end{bmatrix}^T$$

$$\vdots$$

$$h_H = \sigma(\mathbf{w}_H \cdot \mathbf{x} + b_H)$$
$$\mathbf{w}_H = \begin{bmatrix} w_{H1}, w_{H2}, w_{H3}, w_{H4} \end{bmatrix}^T$$
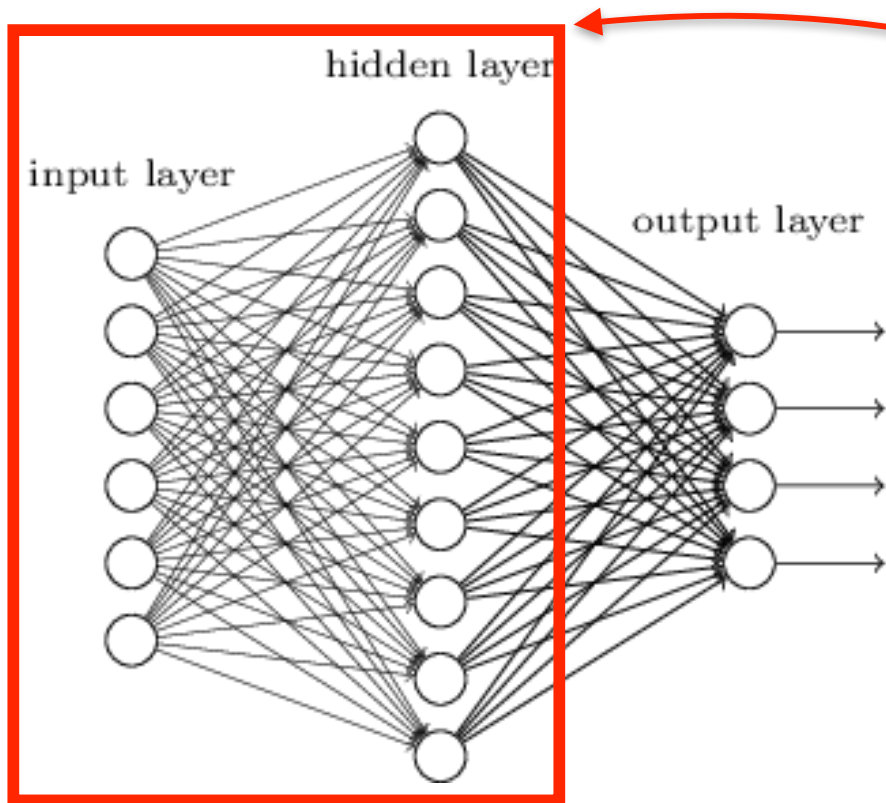
# Repeating the Process



$$\mathbf{h} = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b}),$$

$$\text{with } \mathbf{W} = \begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \\ \vdots \\ \mathbf{w}_H \end{bmatrix}$$

$$\text{and } \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_H \end{bmatrix}.$$

# Multi-Layer Perceptron



$$\mathbf{h} = \sigma_1(\mathbf{W}_1\mathbf{x} + \mathbf{b}_1)$$
$$\mathbf{y} = \sigma_2(\mathbf{W}_2\mathbf{h} + \mathbf{b}_2)$$

- The process can be repeated several times to create a vector **h**.
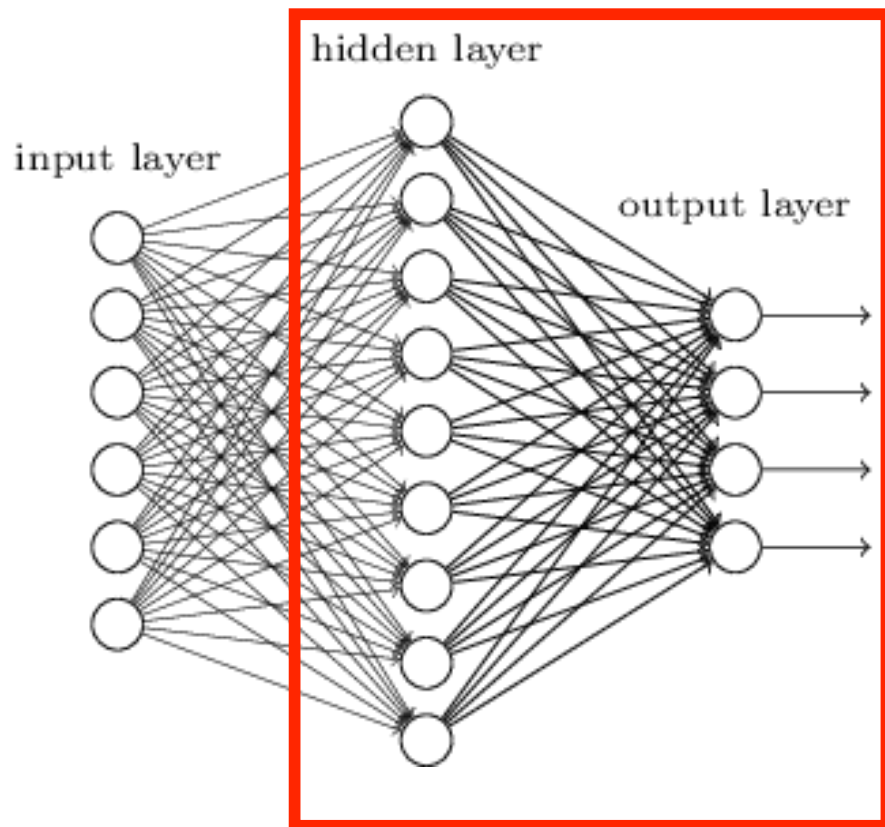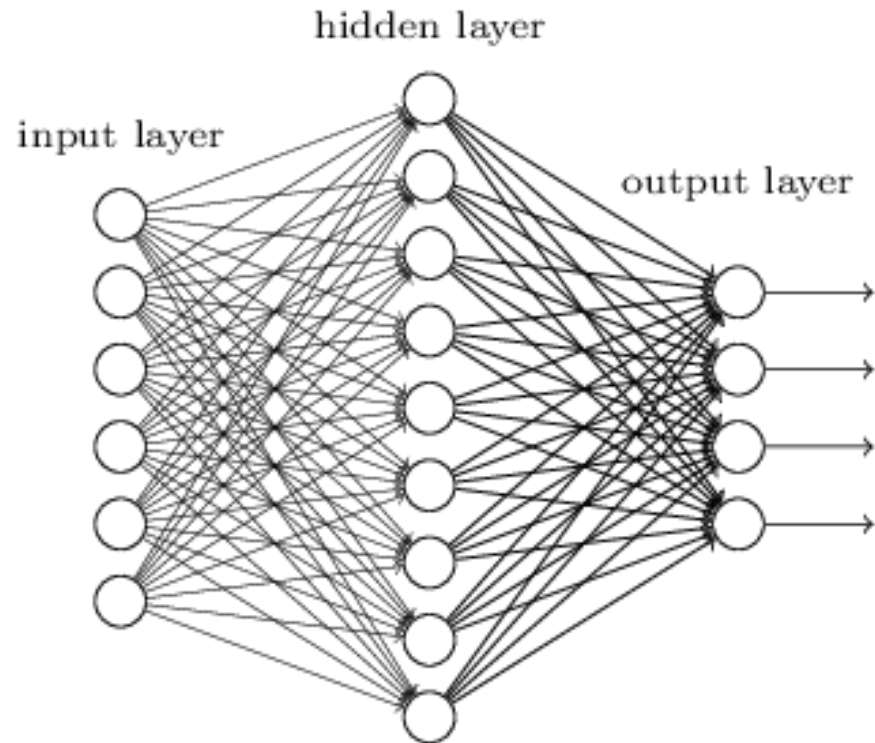
# Multi-Layer Perceptron



$$\mathbf{h} = \sigma_1(\mathbf{W}_1\mathbf{x} + \mathbf{b}_1)$$
$$\mathbf{y} = \sigma_2(\mathbf{W}_2\mathbf{h} + \mathbf{b}_2)$$

- The process can be repeated several times to create a vector **h**.
- It can then be done again to produce an output **y**.

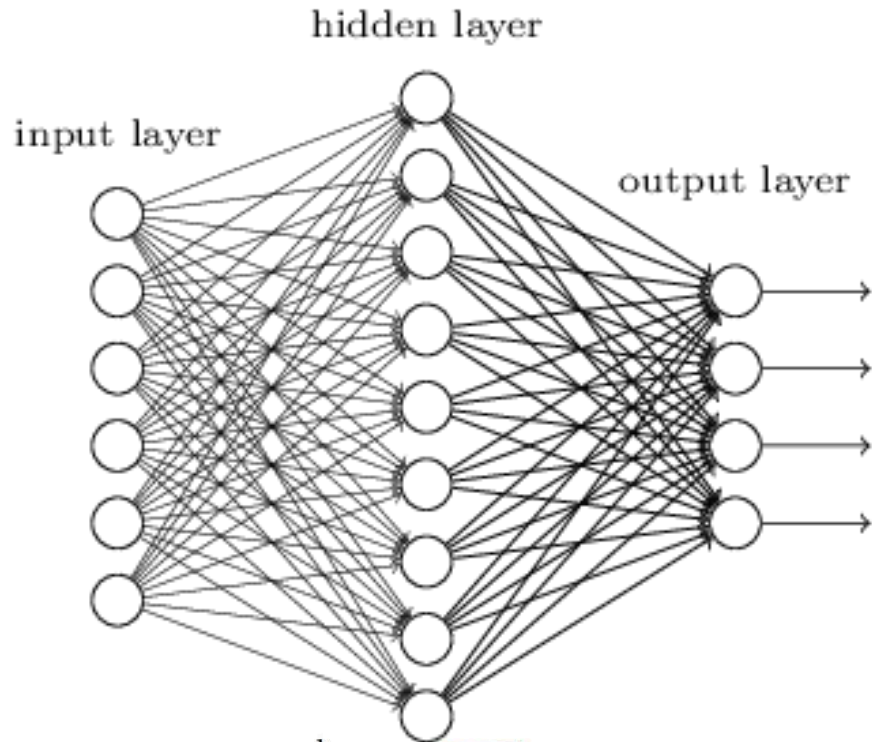—> This output is a **differentiable** function of the weights.

# ReLU



$$\mathbf{h} = \sigma_1(\mathbf{W}_1\mathbf{x} + \mathbf{b}_1)$$

$$\mathbf{y} = \sigma_2(\mathbf{W}_2\mathbf{h} + \mathbf{b}_2)$$
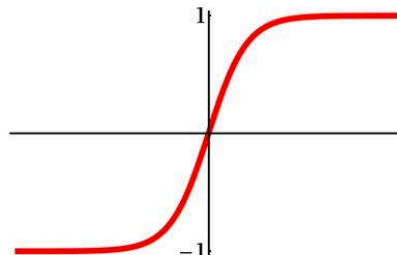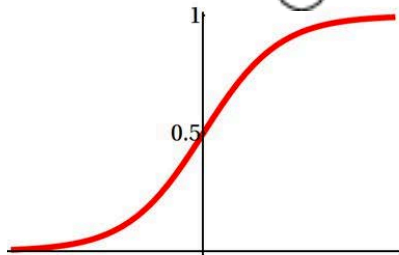
$$\sigma(\mathbf{x}) = \max(0, \mathbf{x})$$

- Each node defines a hyperplane.
- The resulting function is piecewise linear affine and continuous.

# Sigmoid and Tanh



$$\mathbf{h} = \sigma(\mathbf{W}_1\mathbf{x} + \mathbf{b}_1)$$
$$\mathbf{y} = \sigma(\mathbf{W}_2\mathbf{h} + \mathbf{b}_2)$$
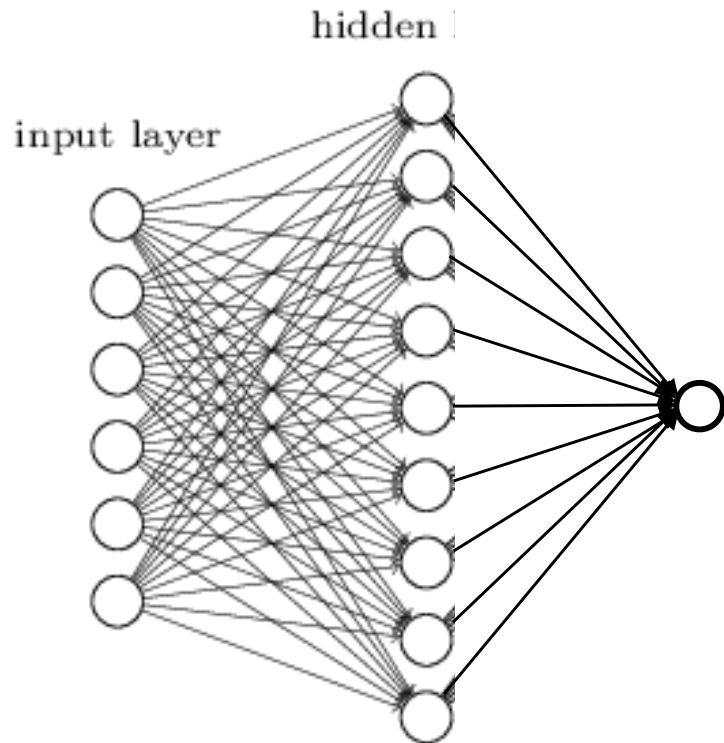
$$\text{sigm:} \quad \sigma(x) = \frac{1}{1 + \exp(-x)}$$

$$\text{tanh:} \quad \sigma(x) = \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)}$$

- Each node defines a hyperplane.
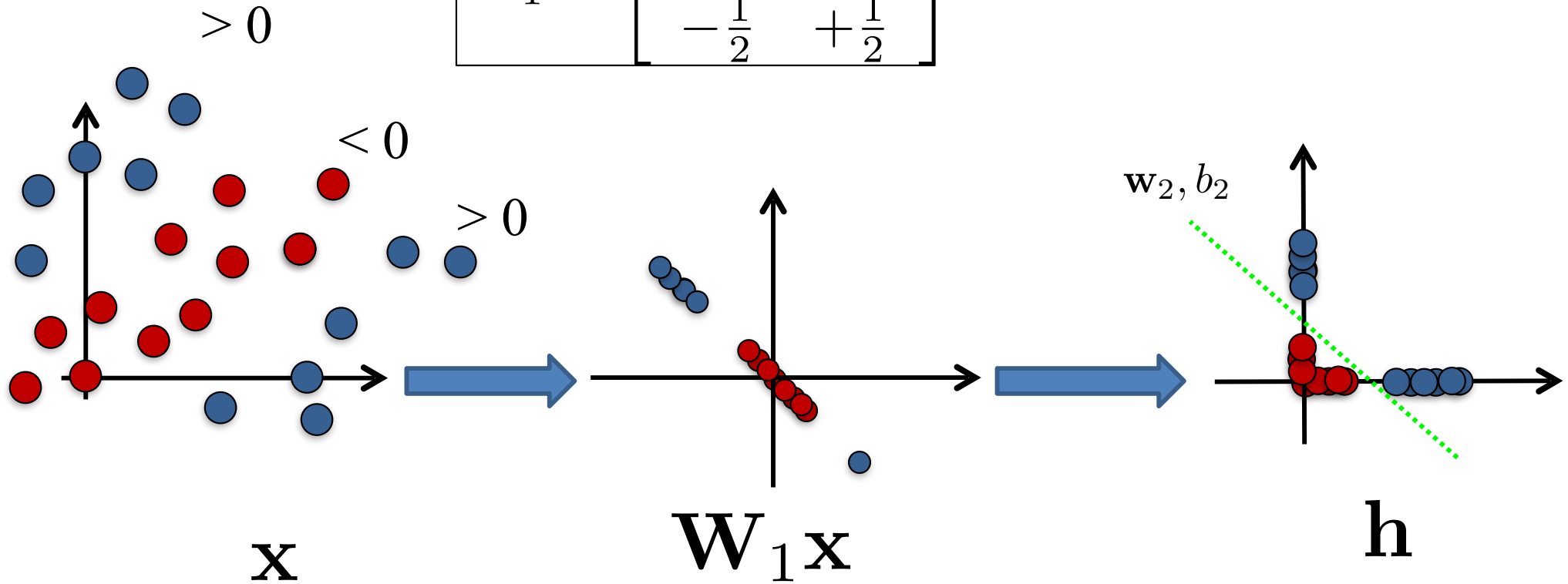- The resulting function is continuously differentiable.

# Binary Case



input layer

hidden l...

$$\mathbf{h} = \sigma(\mathbf{W}_1 \mathbf{x}_n + \mathbf{b}_1)$$

$$y = \sigma(\mathbf{w}_2 \mathbf{h} + b_2)$$

In this case $w_2$ is vector.

# ReLu Behavior

$$\mathbf{W}_1 = \begin{bmatrix} +\frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & +\frac{1}{2} \end{bmatrix}$$

$> 0$

$< 0$

$> 0$

$\mathbf{w}_2, b_2$

$\mathbf{x}$

$\mathbf{W}_1\mathbf{x}$

$\mathbf{h}$

$$\mathbf{h} = \mathrm{ReLu}(\mathbf{W}_1\mathbf{x})$$

$$\mathbf{y} = \mathbf{w}_2^T\mathbf{h} + b_2$$

# Binary Case

- Let the training set be $\{(\mathbf{x}_n, t_n)_{1 \leq n \leq N}\}$ where $t_n \in \{0,1\}$ is the class label and let us consider a neural net with a 1D output.

- We write

$$y_n = \sigma(\mathbf{w}_2(\sigma(\mathbf{W}_1\mathbf{x}_n + \mathbf{b}_1)) + \mathbf{b}_2) \in [0,1] \ .$$

- We want to minimize the binary cross entropy

$$E(\mathbf{W}_1, \mathbf{w}_2, \mathbf{b}_1, \mathbf{b}_2) = \frac{1}{N} \sum_{n=1}^{N} E_n(\mathbf{W}_1, \mathbf{w}_2, \mathbf{b}_1, \mathbf{b}_2) \ ,$$

$$E_n(\mathbf{W}_1, \mathbf{w}_2, \mathbf{b}_1, \mathbf{b}_2) = -(t_n \ln(y_n) + (1 - t_n)\ln(1 - y_n)) \ ,$$

  with respect to the coefficients of $\mathbf{W}_1$, $\mathbf{w}_2$, $\mathbf{b}_1$, and $\mathbf{b}_2$.

- E can be minimized using a gradient-based technique.

# Binary Case

Given a training set $\{\mathbf{x}_n, t_n\}_{1 \leq n \leq N}$ where $t_n \in \{0, 1\}$, minimize
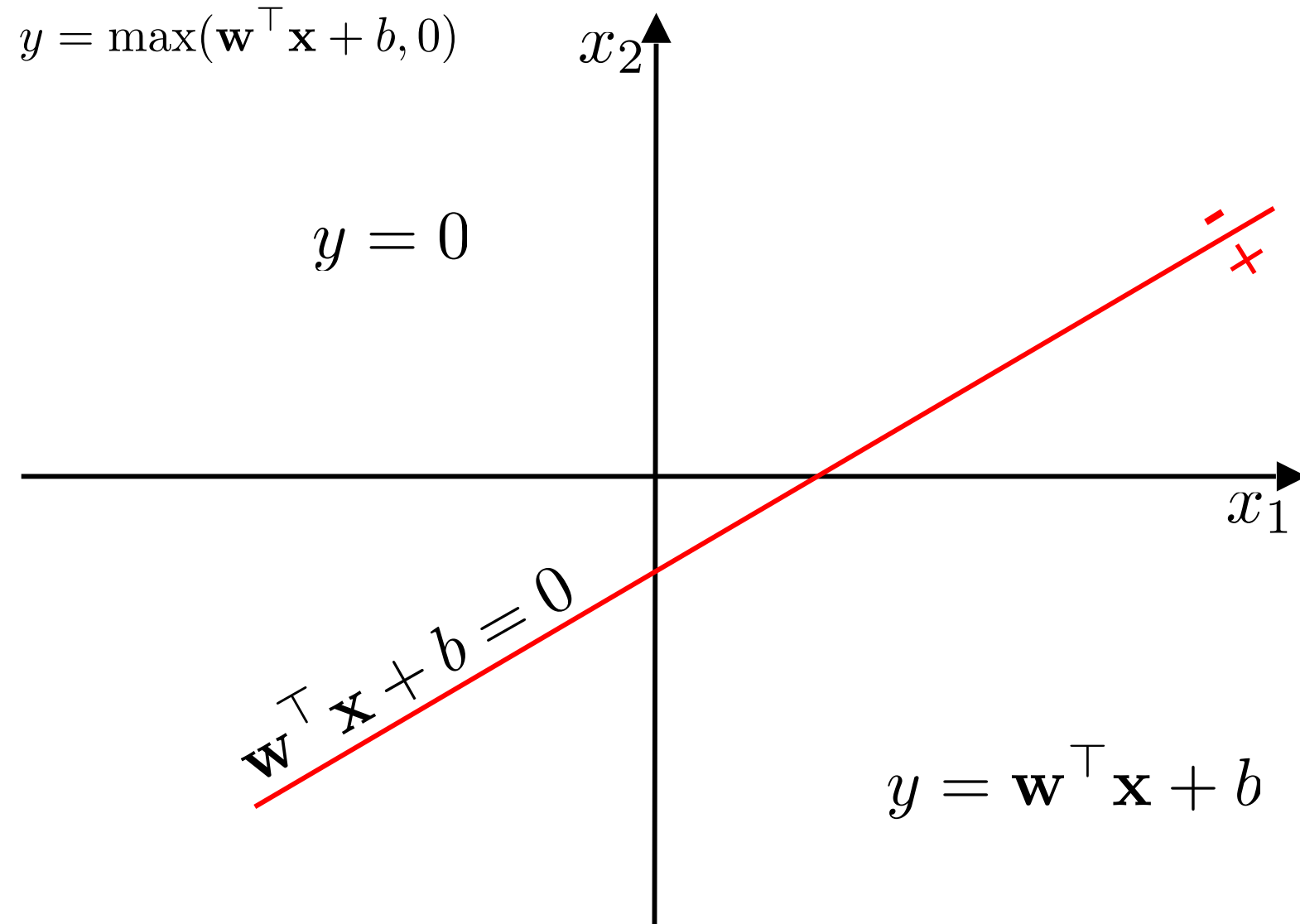
$$E(\mathbf{W}, \mathbf{b}) = -\frac{1}{N} \sum_{n=1}^{N} E_n(\mathbf{W}, \mathbf{b})$$

$$E_n(\mathbf{W}, \mathbf{b}) = t_n \log(y_n) + (1 - t_n) \log(1 - y_n) \ ,$$

$$y_n = f(\mathbf{x}_n)$$
$$= \sigma(\mathbf{W}_2(\sigma(\mathbf{W}_1 \mathbf{x}_n + \mathbf{b}_1)) + \mathbf{b}_2),$$

Since E is a differentiable function of **W** and **b,** this can be done using a gradient-based technique.

# One Single Hyperplane

$$y = \max(\mathbf{w}^\top \mathbf{x} + b, 0)$$

$x_2$

$y = 0$

$\mathbf{w}^\top \mathbf{x} + b = 0$

$x_1$

$y = \mathbf{w}^\top \mathbf{x} + b$
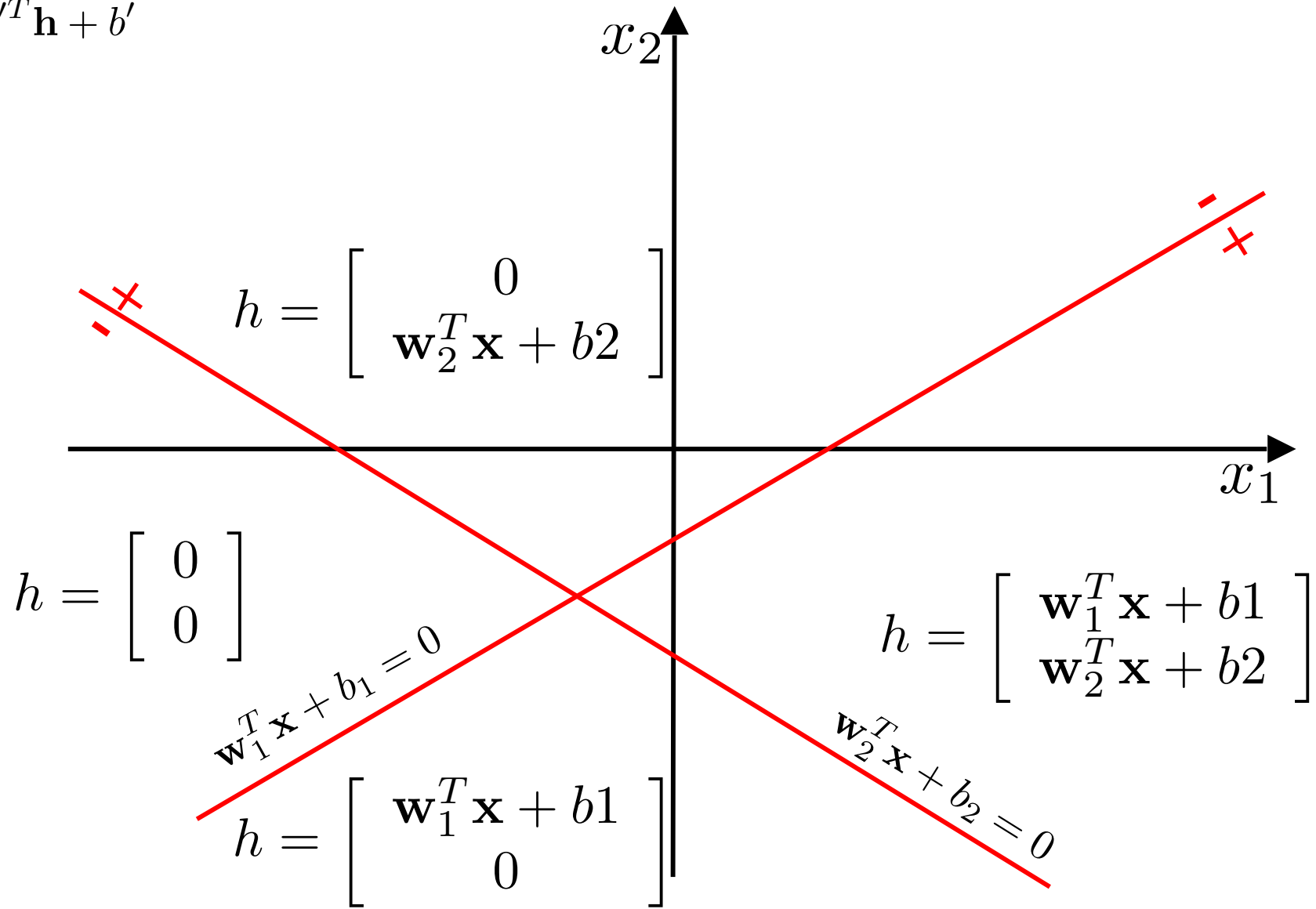
# Two Hyperplanes

$$\mathbf{h} = \max(\mathbf{W}\mathbf{x} + \mathbf{b}, 0) \text{ with } \mathbf{W} = \left[ \begin{array}{c} \mathbf{w}_1^T \\ \mathbf{w}_2^T \end{array} \right] \text{ and } \mathbf{b} = \left[ \begin{array}{c} b_1 \\ b_2 \end{array} \right]$$

$$y = \mathbf{w'}^T \mathbf{h} + b'$$

$$h = \left[ \begin{array}{c} 0 \\ \mathbf{w}_2^T \mathbf{x} + b2 \end{array} \right]$$

$$h = \left[ \begin{array}{c} 0 \\ 0 \end{array} \right]$$

$$h = \left[ \begin{array}{c} \mathbf{w}_1^T \mathbf{x} + b1 \\ \mathbf{w}_2^T \mathbf{x} + b2 \end{array} \right]$$

$$\mathbf{w}_1^T \mathbf{x} + b_1 = 0$$

$$\mathbf{w}_2^T \mathbf{x} + b_2 = 0$$

$$h = \left[ \begin{array}{c} \mathbf{w}_1^T \mathbf{x} + b1 \\ 0 \end{array} \right]$$
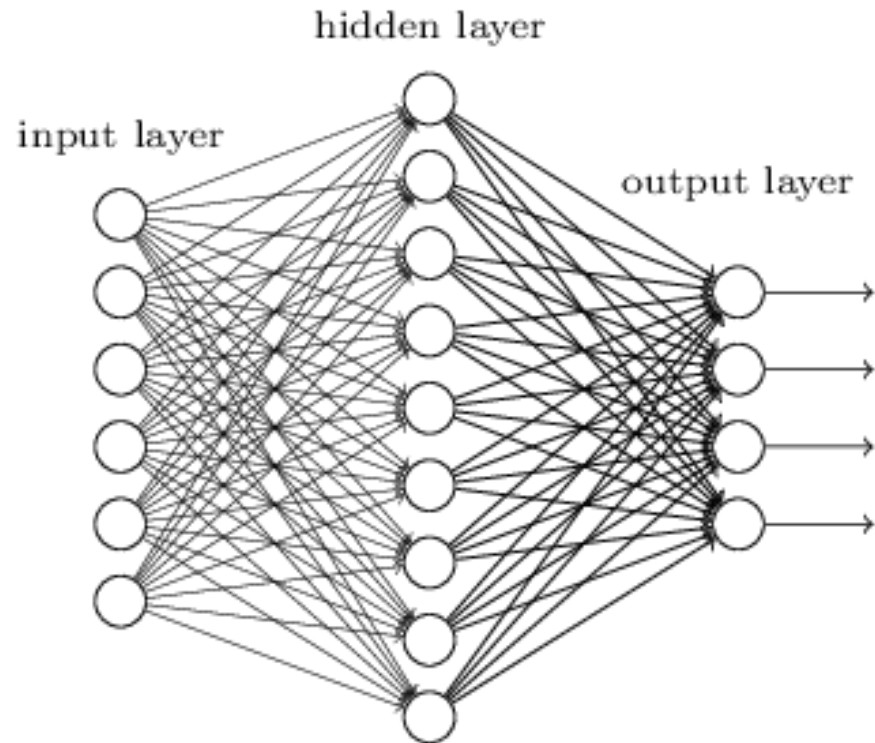
$x_2$

$x_1$

# Three Hyperplanes

$$\begin{cases} \mathbf{h} = \max(\mathbf{W}\mathbf{x} + \mathbf{b}, 0) \\ y = \mathbf{w'}^{\top}\mathbf{h} \end{cases}$$

with $\dim(\mathbf{h}) = 3$



$x_2$

$x_1$

$\mathbf{w}_3^T \mathbf{x} + b_3 = 0$

$\mathbf{w}_1^T \mathbf{x} + b_1 = 0$

$\mathbf{w}_2^T \mathbf{x} + b_2 = 0$

# Multi-Class Case

hidden layer

input layer

output layer

$$\mathbf{h} = \sigma(\mathbf{W}_1 \mathbf{x}_n + \mathbf{b}_1)$$

$$y = \sigma(\mathbf{W}_2 \mathbf{h} + \mathbf{b}_2)$$

In this case $\mathbf{W}_2$ is a matrix.

# Multi-Class Case

Let the training set be $\{(\mathbf{x}_n, [t_n^1, \ldots, t_n^K])_{1 \leq n \leq N}\}$ where $t_n^k \in \{0,1\}$ is the probability that sample $\mathbf{x}_n$ belongs to class k.

- We write

$$\mathbf{y}_n = \sigma(\mathbf{W}_2(\sigma(\mathbf{W}_1\mathbf{x}_n + \mathbf{b}_1)) + \mathbf{b}_2) \in R^K$$

$$p_n^k = \frac{\exp(\mathbf{y}_n[k])}{\sum_j \exp(\mathbf{y}_n[j])}$$

- We minimize the cross entropy

$$E(\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}_1, \mathbf{b}_2) = \frac{1}{N} \sum_{n=1}^{N} E_n(\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}_1, \mathbf{b}_2) \,,$$

$$E_n(\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}_1, \mathbf{b}_2) = - \sum t_n^k \ln(p_n^k) \,,$$

with respect to the coefficients of $\mathbf{W}_1$, $\mathbf{W}_2$, $\mathbf{b}_1$, and $\mathbf{b}_2$.

# Non-Linear Binary Classification



Positive: $100(x_2 - x_1^2)^2 + (1 - x_1)^2 < 0.5$
Negative: Otherwise

$y(\mathbf{x}; \tilde{\mathbf{w}})$ is now a non-linear function implemented by the network.

**Problem statement:** Find $\tilde{\mathbf{w}}$ such that

- for all or most positive samples $y(\tilde{\mathbf{x}}; \tilde{\mathbf{w}}) > 0.0$,
- for all or most negative samples $y(\tilde{\mathbf{x}}; \tilde{\mathbf{w}}) < 0.0$.
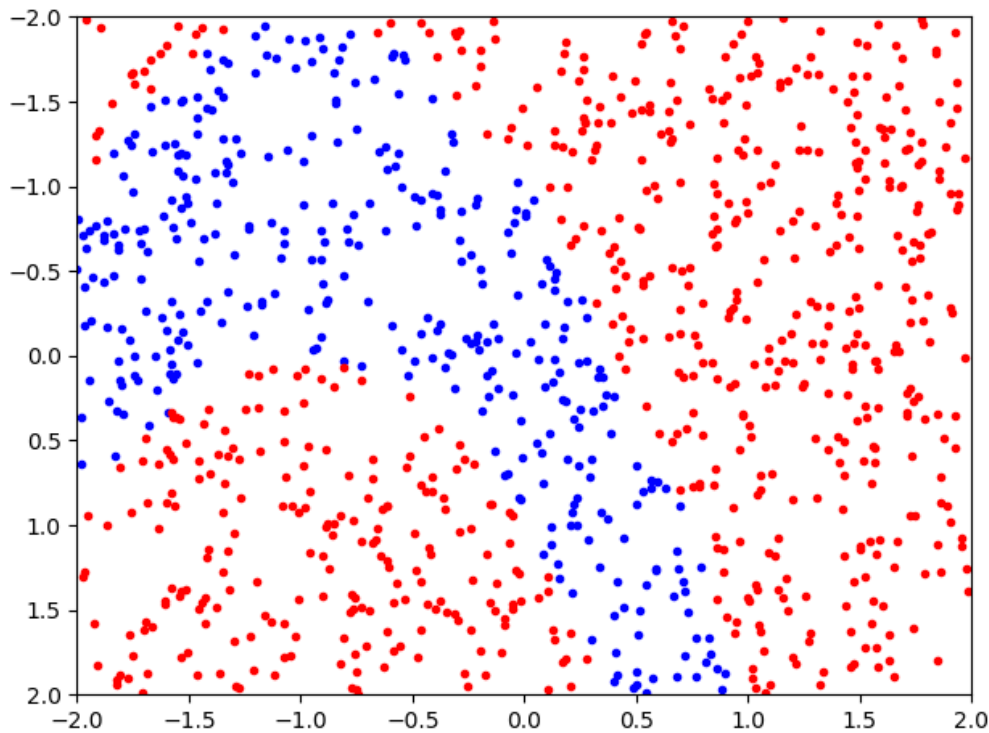
# Non-Linear Regression



$$z = f(\mathbf{x})$$
$$= 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

**Problem statement:** Given $(\{\mathbf{x}_1, z_1\}, \ldots, \{\mathbf{x}_n, z_n\})$, minimize
$$\sum_i (z_i - f(\mathbf{x}_i, \tilde{\mathbf{w}})^2)$$
w.r.t. $\tilde{\mathbf{w}}$.

# Classification / Regression



Positive: $f(\mathbf{x}) < 0.5$
Negative: Otherwise

$y(\mathbf{x}; \tilde{\mathbf{w}})$ is now a non-linear function implemented by the network.

Classification can be understood as finding $\tilde{\mathbf{w}}$ such that

$$y(\mathbf{x}; \tilde{\mathbf{w}}) \approx f(x)$$

# From Classification to Regression

x,y

Linear 2 -> n

**tanh**

Linear n -> 1

**sigm**

z

Minimize $\sum_i [t_i \log(\text{sigm}(f(x_i, y_i))) + (1 - t_i) \log(1 - \text{sigm}(f(x_i, y_i)))$ with respect to $\mathbf{W}_1, \mathbf{w}_2, b_x, b_y, b_z$.

Minimize $\sum_i (z_i - f(x_i, y_i))^2$, with respect to $\mathbf{W}_1, \mathbf{w}_2, b_x, b_y, b_z$.

# Approximating a Surface



$$x,y$$

| Linear 2 -> n | $W_1, b_1$ |

tanh

| Linear n -> 1 | $w_2, b_2$ |

$$z$$

$$z = f(x, y)$$

$$= w_2 \sigma\left(W_1 \begin{bmatrix} x \\ y \end{bmatrix} + b_1\right) + b_2$$

Given $(\{x_1, z_1\}, \ldots, \{x_n, z_n\})$, minimize

$$\sum_i (z_i - f(x_i))^2$$

with respect to $W_1, w_2, b_1, b_2$.

$$z = f(\mathbf{x})$$

$$= 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

# Interpolating a Surface



$$z = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

3-node hidden layer

# Interpolating a Surface



$$z = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

4-node hidden layer

loss: 1.089789e+00

# Adding more Nodes



2 nodes -> loss 3.02e-01

3 nodes -> loss 2.08e-02

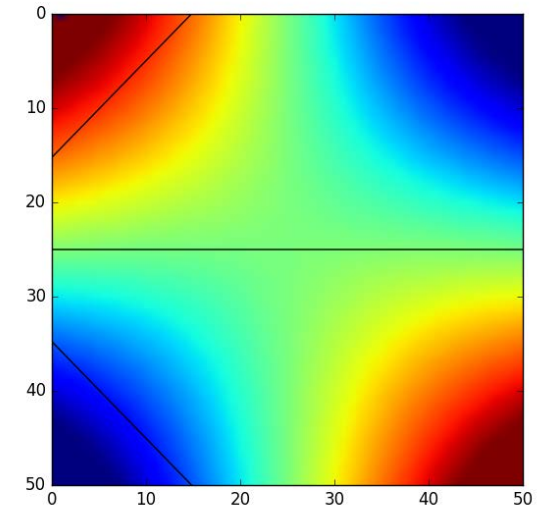4 nodes -> loss 8.27e-03

$$z = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

# Adding more Nodes
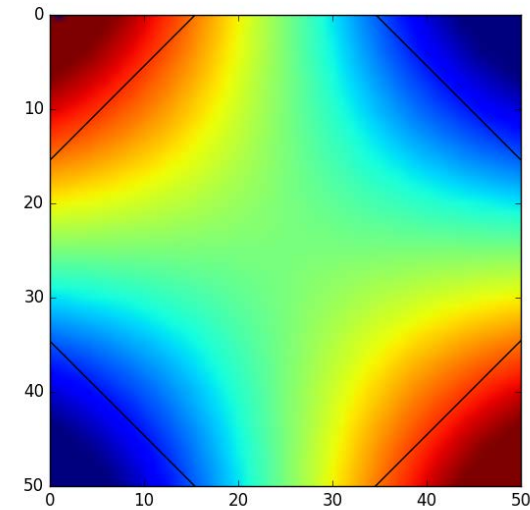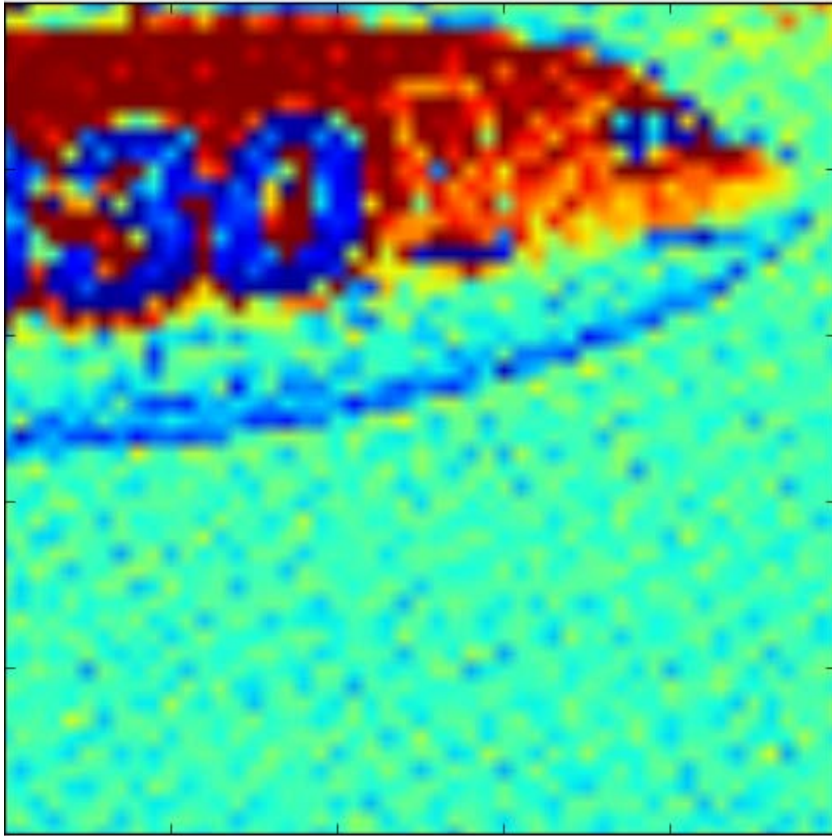


$$z = sin(x)sin(y)$$

2 nodes -> loss 2.61e-01

3 nodes -> loss 2.51e-04

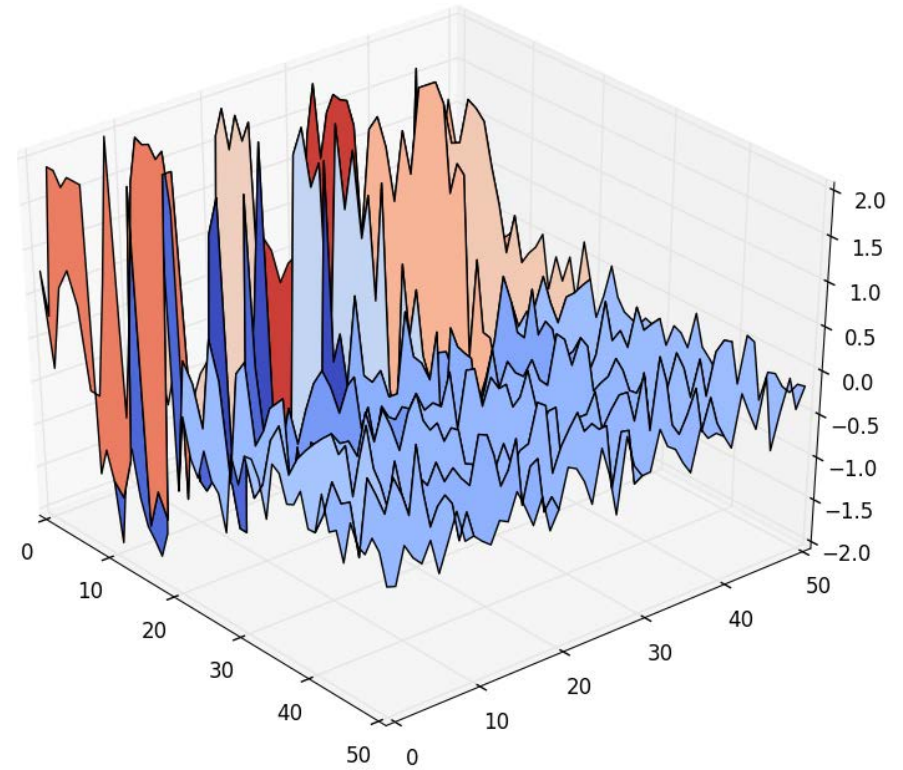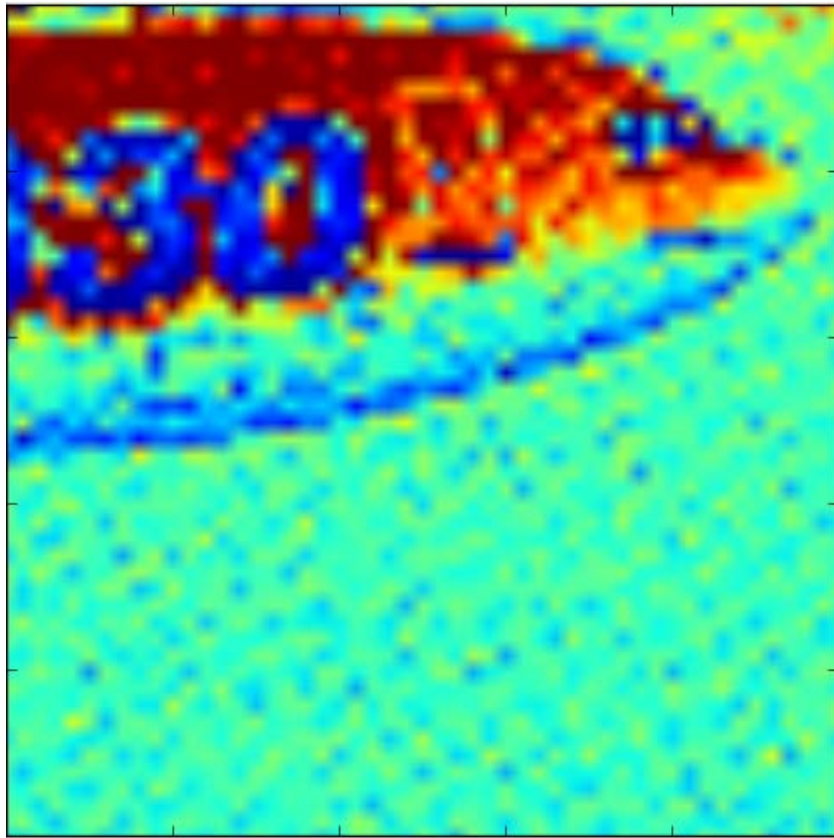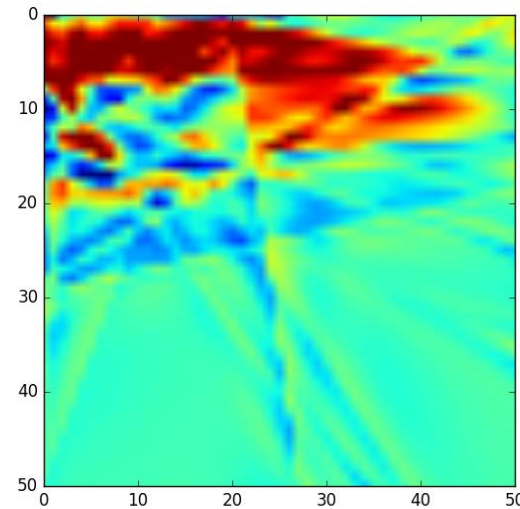4 nodes -> loss 3.07e-07

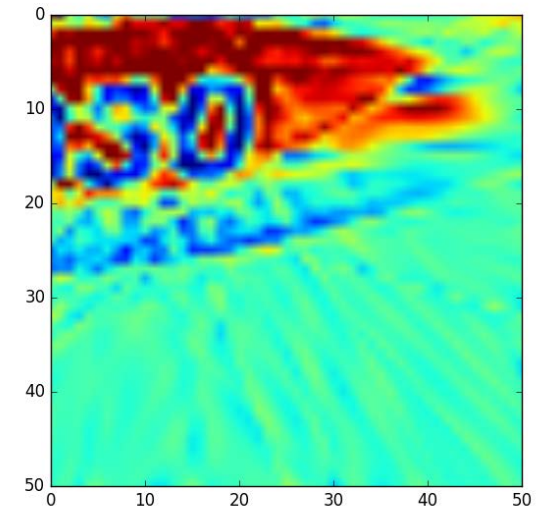# More Complex Surface



$$I = f(x, y)$$
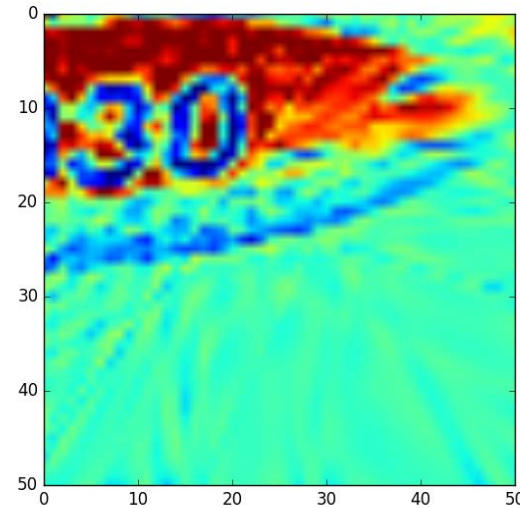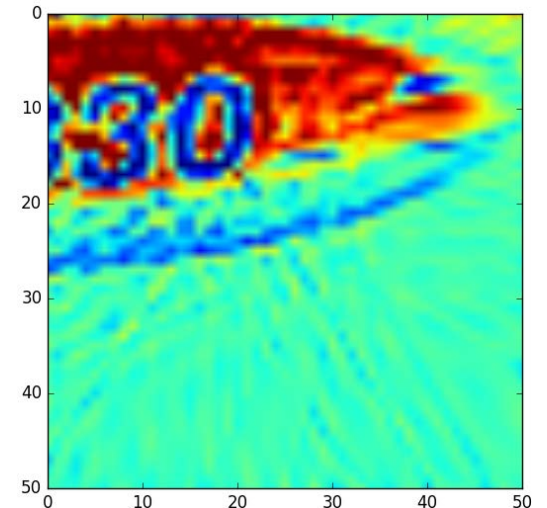
# More Complex Surface



$$I = f(x, y)$$

50 nodes -> loss 3.65e-01

100 nodes -> loss 2.50e-01

125 nodes -> loss 2.40e-01

300 nodes -> loss 1.92e-01

# Universal Approximation Theorem

A feedforward network with a linear output layer and at least one hidden layer with any 'squashing' activation function (e.g. logistic sigmoid) can approximate any Borel measurable function (from one finite-dimensional space to another) with any desired nonzero error.

Any continuous function on a closed and bounded set of $R^n$ is Borel-measurable.

—> In theory, any reasonable function can be approximated by a one-hidden layer network as long as it is continuous.

[Hornik et al, 1989; Cybenko, 1989]

# In Practice



- It may take an exponentially large number of parameters for a good approximation.
- The optimization problem becomes increasingly difficult.

—> The one hidden layer perceptron may not converge to the best solution!

# MNIST



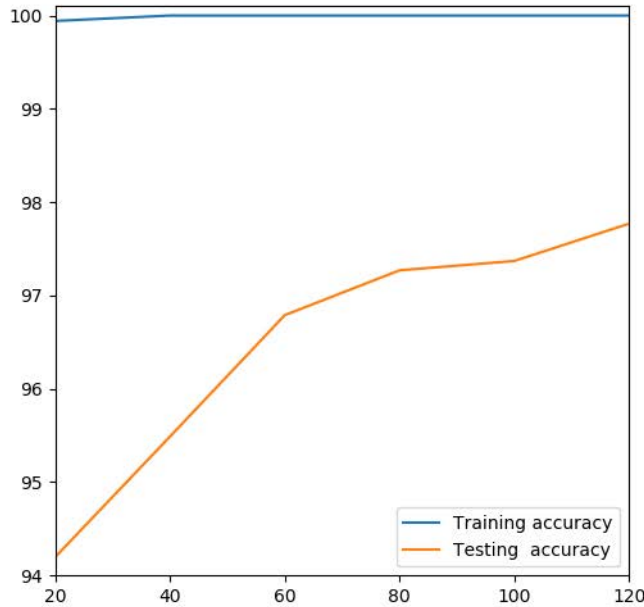- The network takes as input 28x28 images represented as 784D vectors.
- The output is a 10D vector giving the probability of the image representing any of the 10 digits.
- There are 50'000 training pairs of images and the corresponding label, 10'000 validation pairs, and 5'000 testing pairs.

# MNIST Results



nIn =784

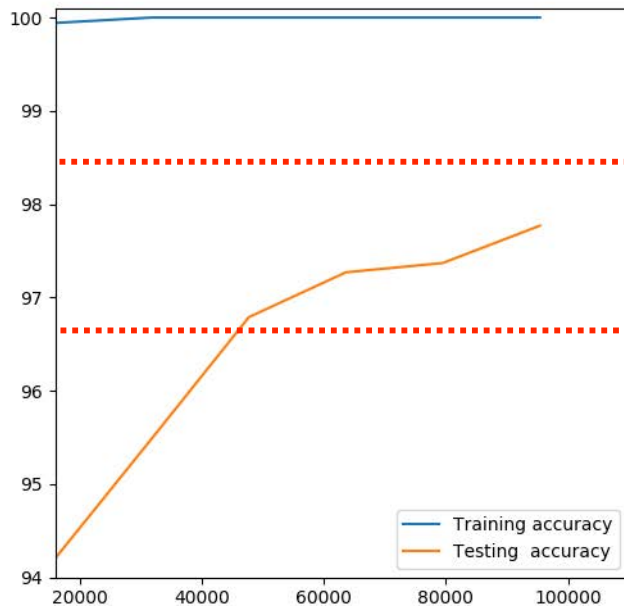nOut = 10

20 < hidden layer size < 120

- MLPs have **many** parameters.
- This has long been a major problem.
—> Was eventually solved by using GPUs.

- Around 2005, SVMs were often felt to be superior to neural nets.
- This is no longer the case ….

# Deep Learning



$$\mathbf{h}_1 = \sigma_1(\mathbf{W}_1\mathbf{x}_1 + \mathbf{b}_1)$$
$$\mathbf{h}_2 = \sigma_2(\mathbf{W}_2\mathbf{h}_2 + \mathbf{b}_2)$$
$$\cdots$$
$$\mathbf{y} = \sigma_n(\mathbf{W}_n\mathbf{h}_n + \mathbf{b}_n)$$

- The descriptive power of the net increases with the number of layers.

- In the case of a 1D signal, it is roughly proportional to $\prod_n W_n$ where $w_n$ is the width of layer n.

# One Layer: Two Hyperplanes

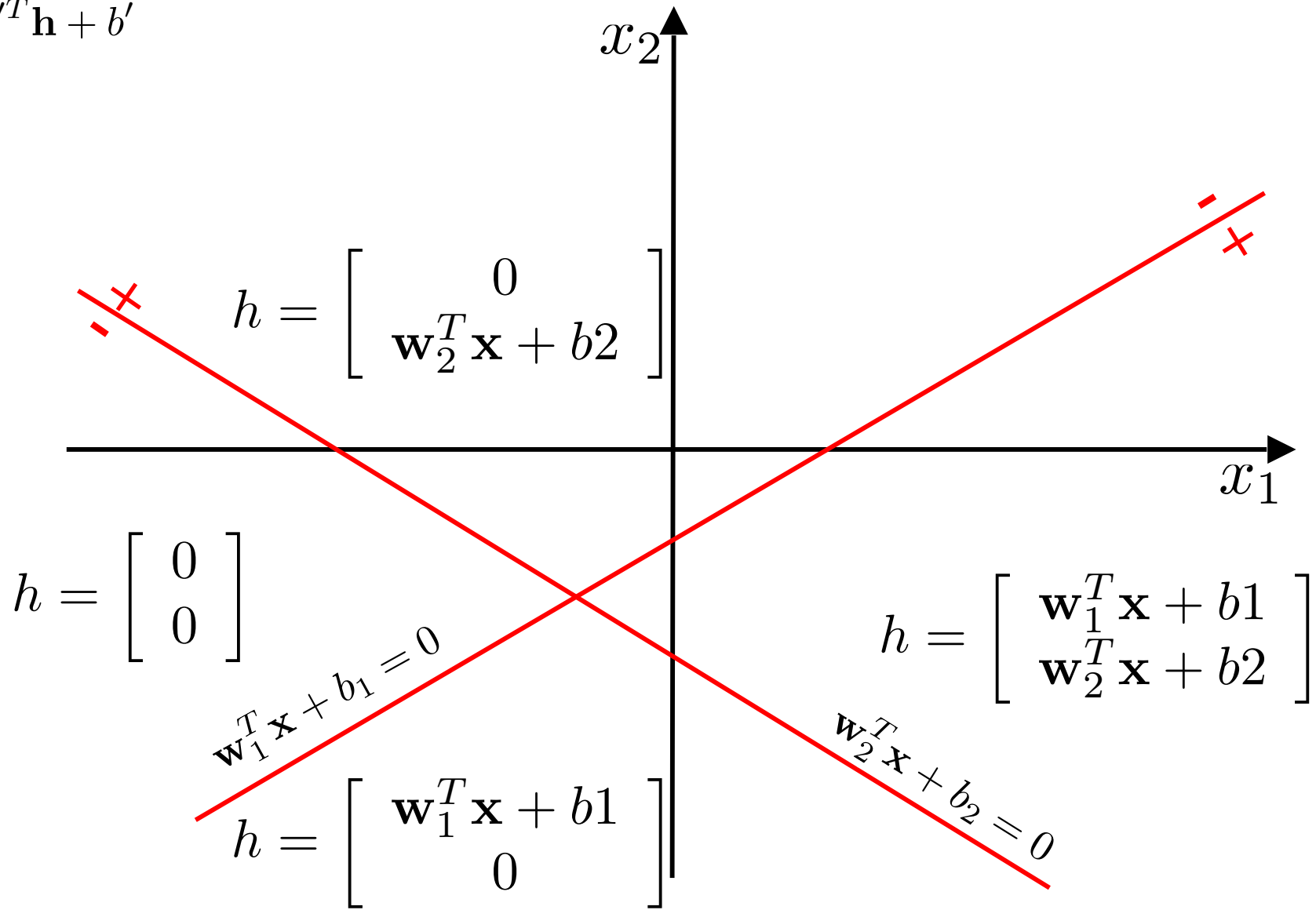$$\mathbf{h} = \max(\mathbf{W}\mathbf{x} + \mathbf{b}, 0) \text{ with } \mathbf{W} = \begin{bmatrix} \mathbf{w}_1^T \\ \mathbf{w}_2^T \end{bmatrix} \text{ and } \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$

$$y = \mathbf{w}'^T\mathbf{h} + b'$$

$$h = \begin{bmatrix} 0 \\ \mathbf{w}_2^T\mathbf{x} + b2 \end{bmatrix}$$

$$h = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$h = \begin{bmatrix} \mathbf{w}_1^T\mathbf{x} + b1 \\ \mathbf{w}_2^T\mathbf{x} + b2 \end{bmatrix}$$

$$\mathbf{w}_1^T\mathbf{x} + b_1 = 0$$

$$\mathbf{w}_2^T\mathbf{x} + b_2 = 0$$

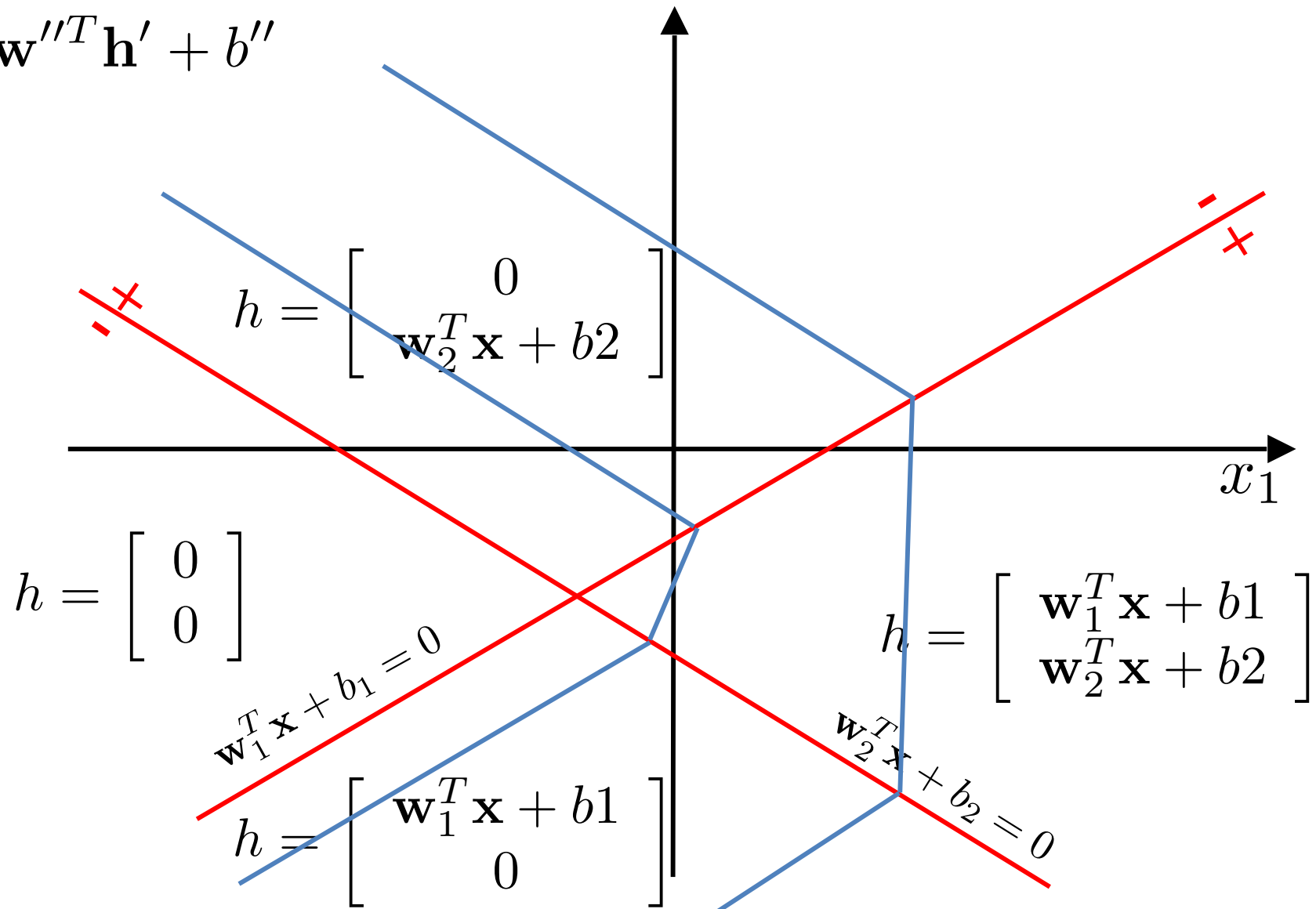$$h = \begin{bmatrix} \mathbf{w}_1^T\mathbf{x} + b1 \\ 0 \end{bmatrix}$$
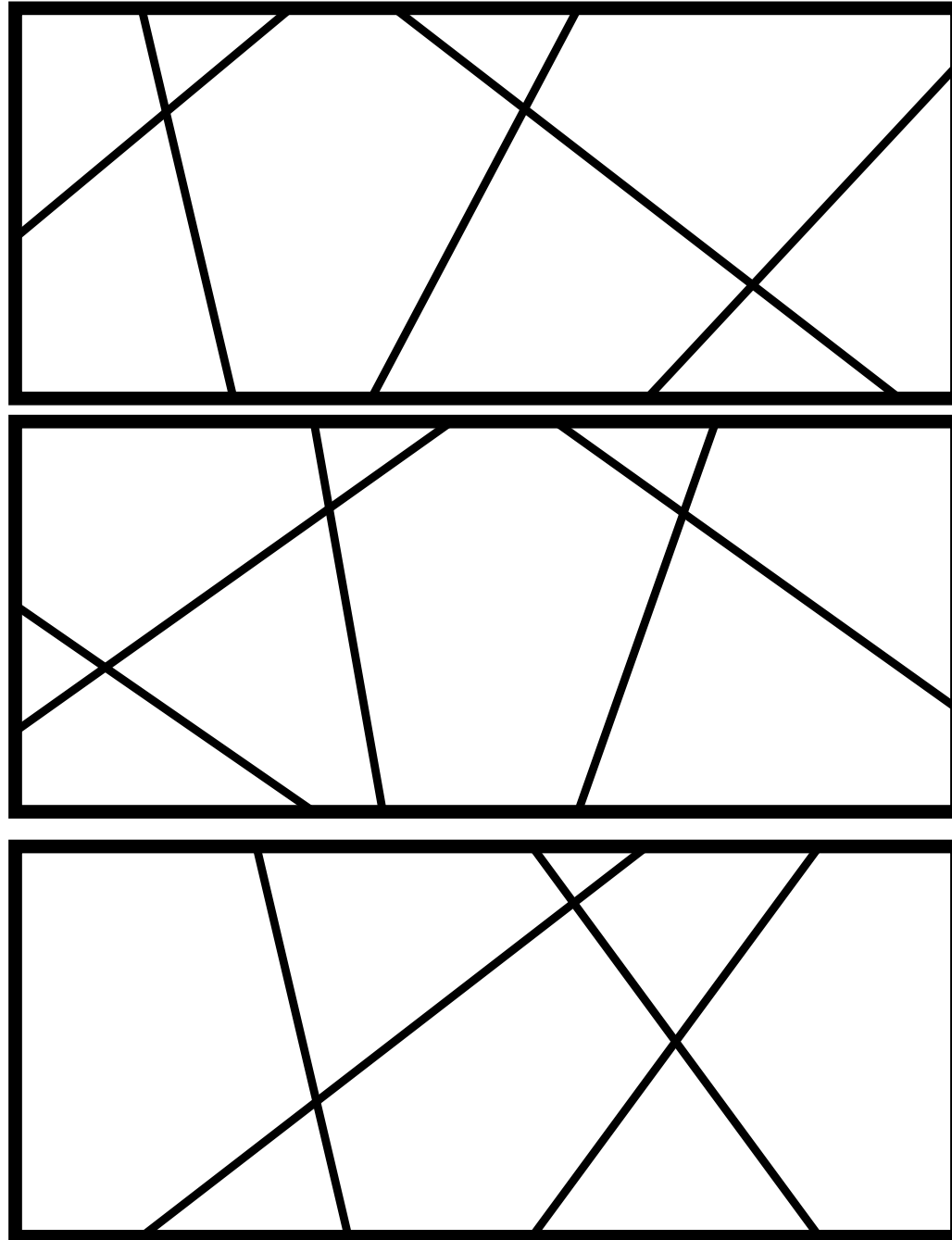
# Two Layers: Two Hyperplanes

$$\mathbf{h} = \max(\mathbf{Wx} + \mathbf{b}, 0)$$

$$\mathbf{h}' = \max(\mathbf{W}'\mathbf{h} + \mathbf{b}', 0)$$

$$y = \mathbf{w}''^T \mathbf{h}' + b''$$



$$h = \begin{bmatrix} 0 \\ \mathbf{w}_2^T \mathbf{x} + b2 \end{bmatrix}$$

$$h = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\mathbf{w}_1^T \mathbf{x} + b_1 = 0$$

$$h = \begin{bmatrix} \mathbf{w}_1^T \mathbf{x} + b1 \\ \mathbf{w}_2^T \mathbf{x} + b2 \end{bmatrix}$$

$$\mathbf{w}_2^T \mathbf{x} + b_2 = 0$$

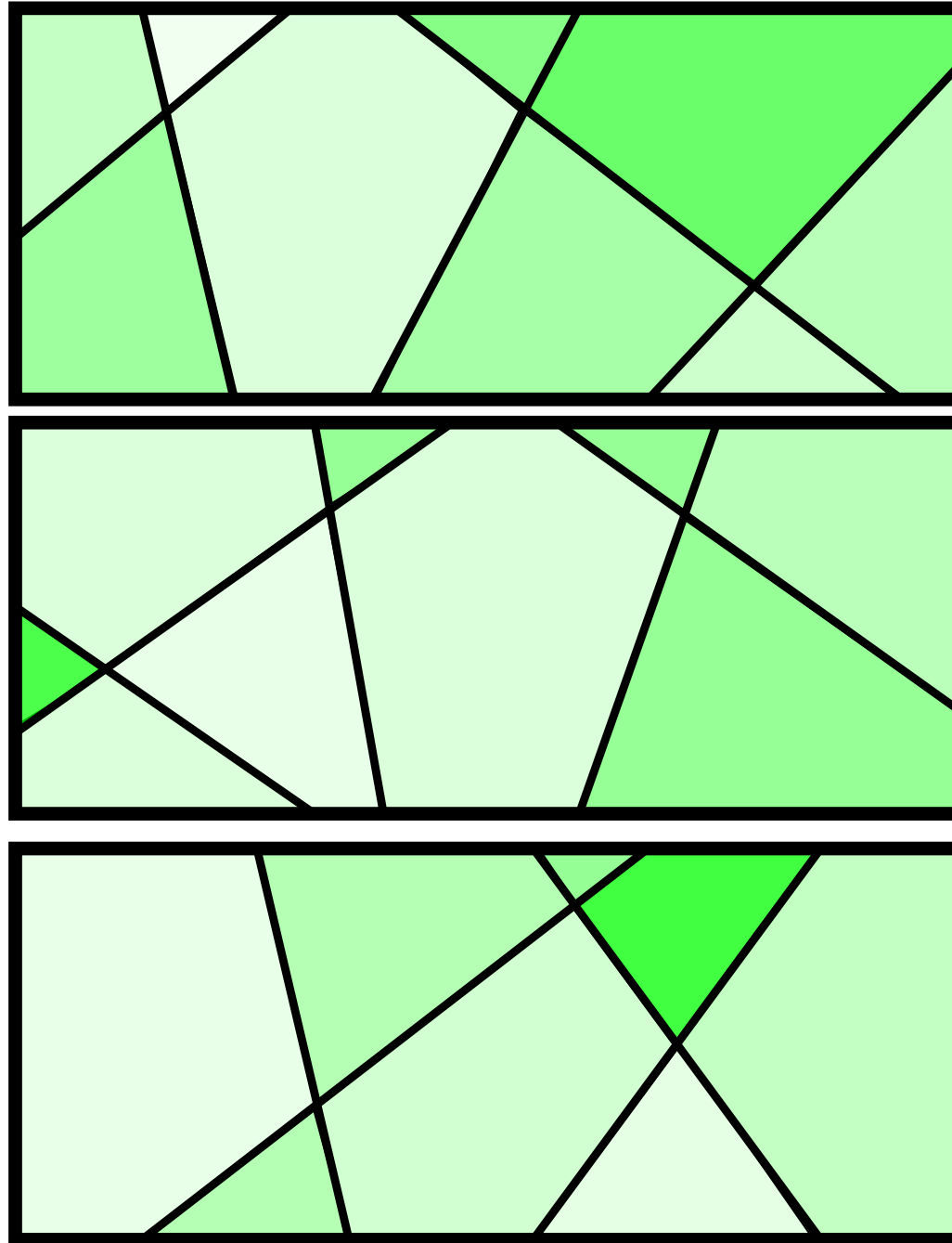$$h = \begin{bmatrix} \mathbf{w}_1^T \mathbf{x} + b1 \\ 0 \end{bmatrix}$$
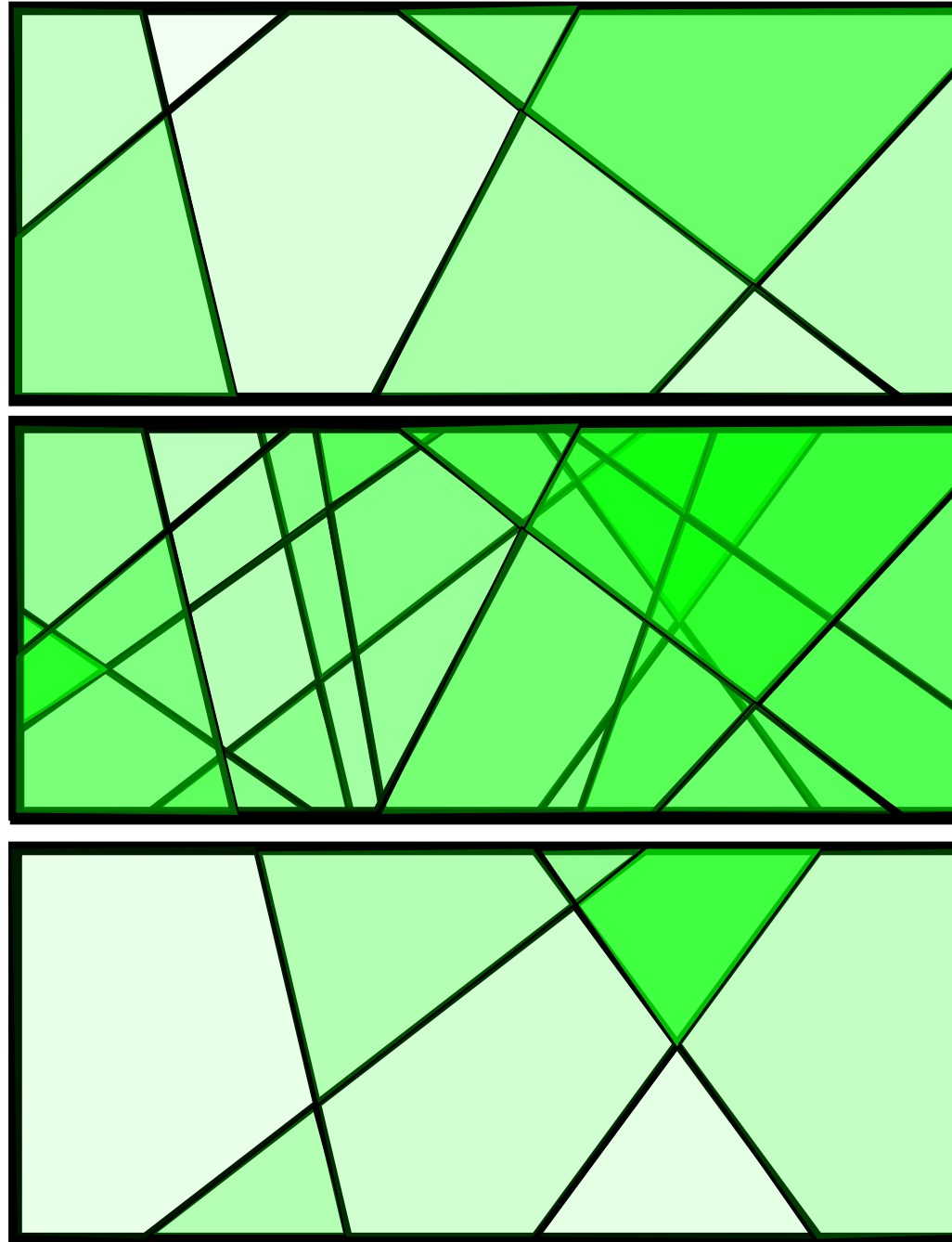
$x_1$

# Graphical Interpretation



Hyperplanes at every level of the network split the space.
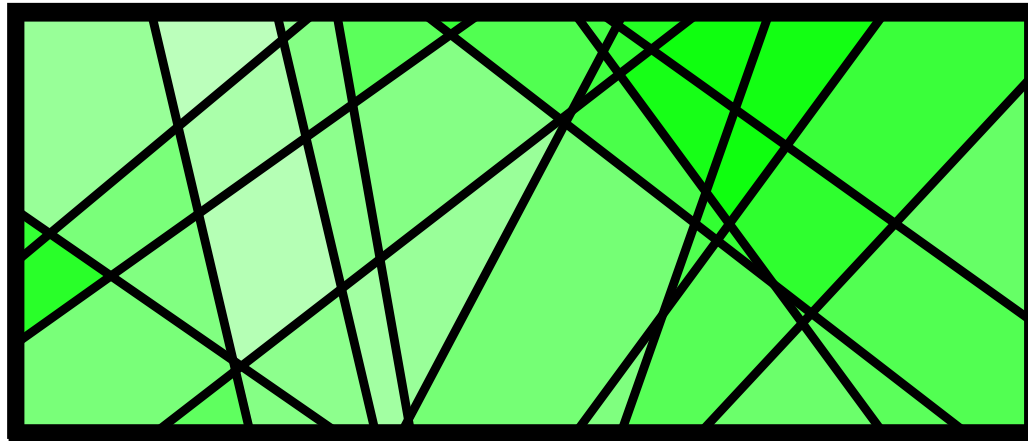
# Graphical Interpretation



Hyperplanes at every level of the network split the space.

# Graphical Interpretation



The splits are combined by the hierarchical nature of the network.

# Graphical Interpretation



The splits are combined due to the sequential nature of the network.
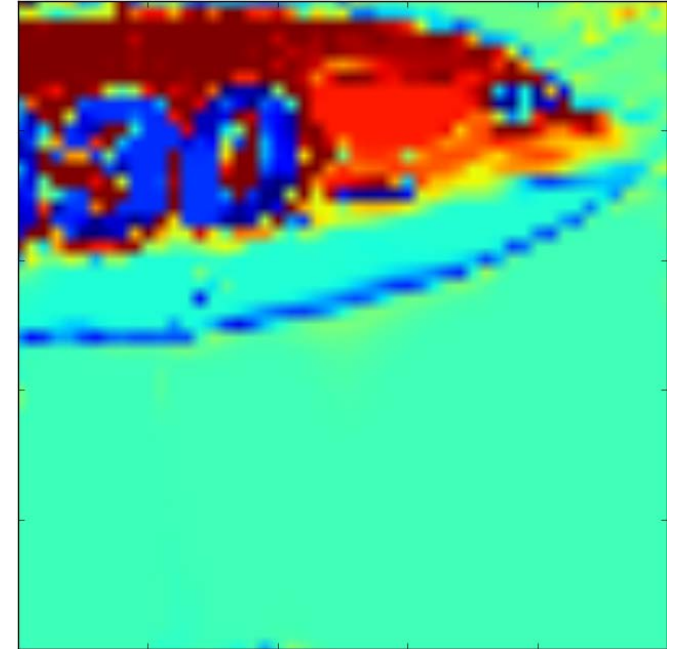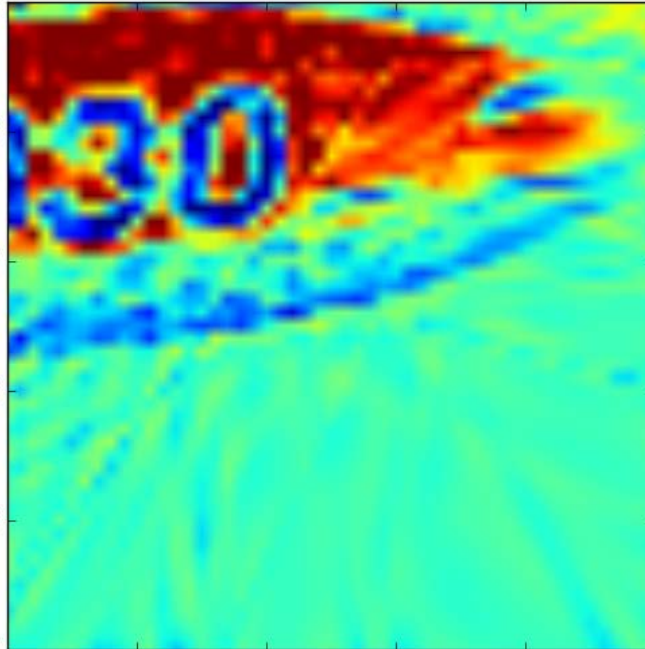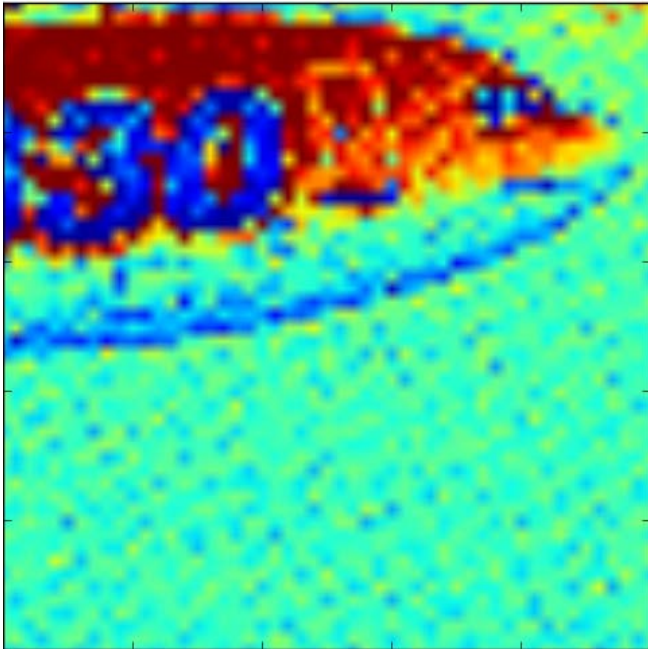
# Multi Layer Perceptrons

The function learned by an MLP using the ReLU, Sigmoid, or Tanh operators is:

- piecewise affine or smooth;
- continuous because it is a composition of continuous functions.

Each region created by a layer is split into smaller regions:

- Their boundaries are correlated in a complex way.
- Their descriptive power is larger than shallow networks for the same number of parameters.
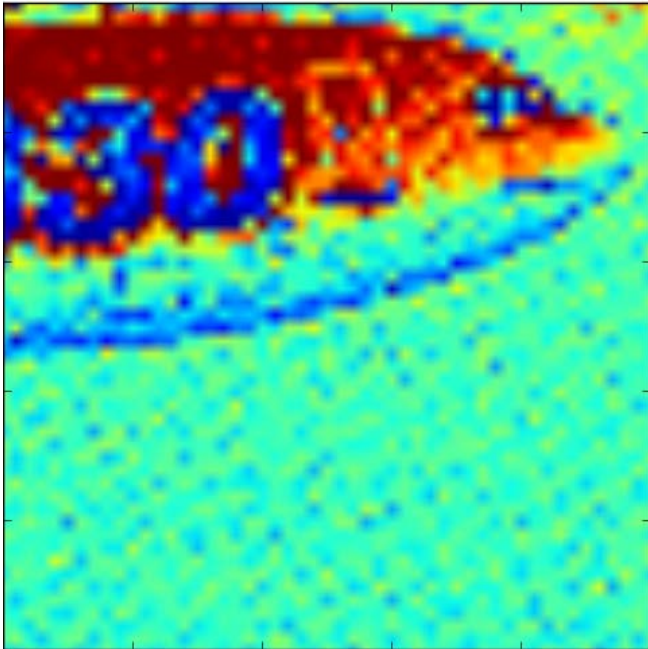
# Second Layer for Approximation
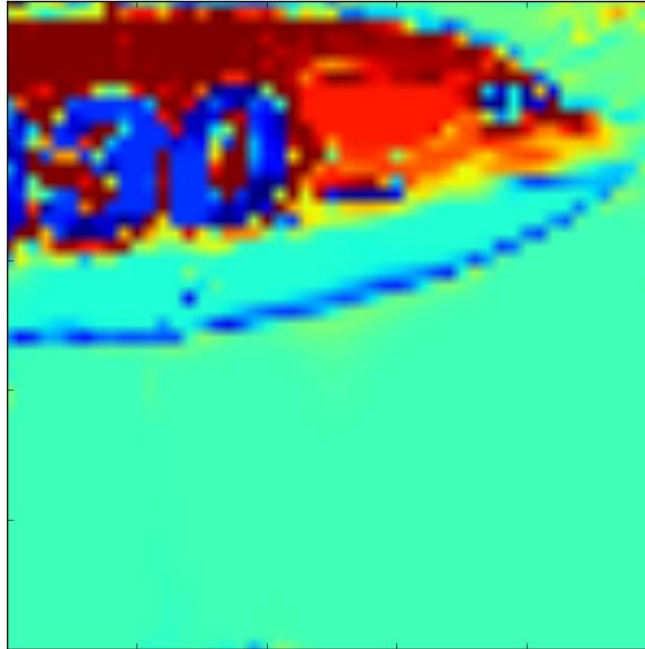


$$I = f(x, y)$$

1 Layer: 125 nodes -> loss 2.40e-01     2 Layers: 20 nodes -> loss 8.31e-02
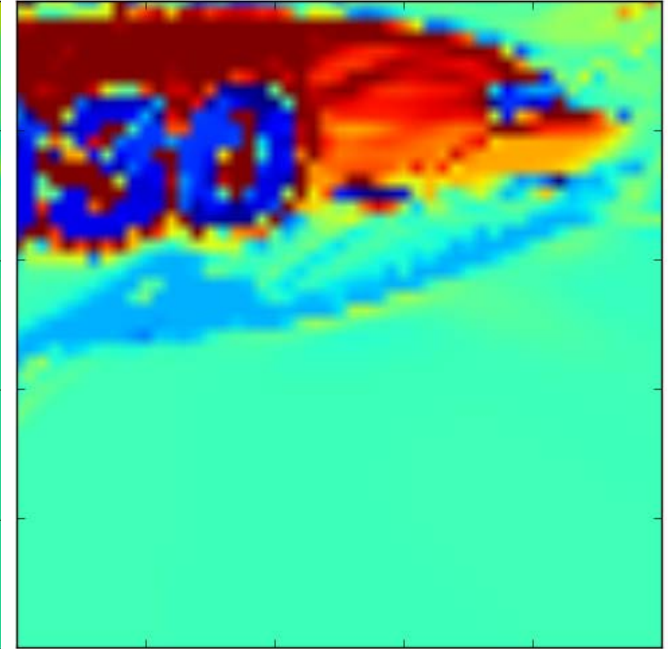
501 weights in both cases
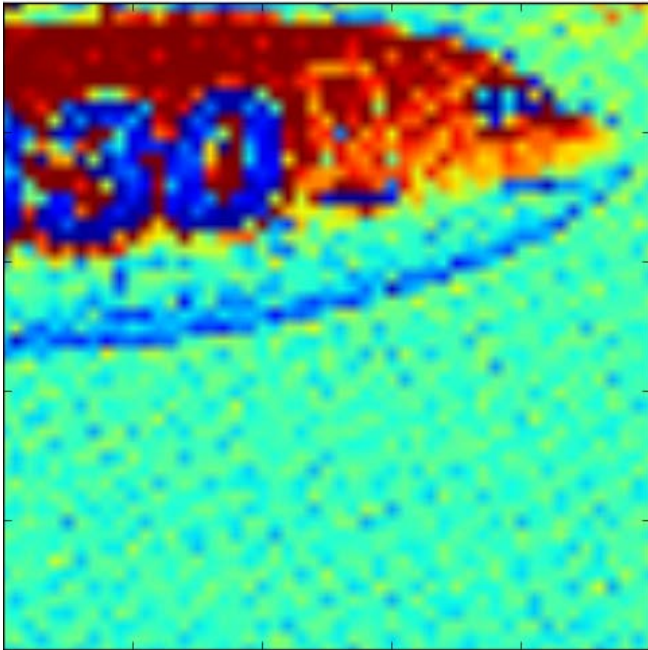
# Adding a Third Layer



$$I = f(x, y)$$

2 Layers: 20 nodes -> loss 8.31e-02
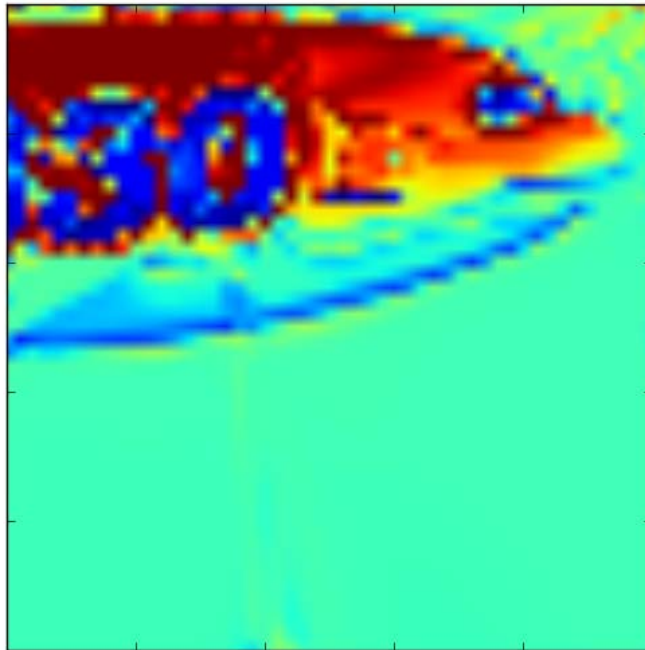
501 weights
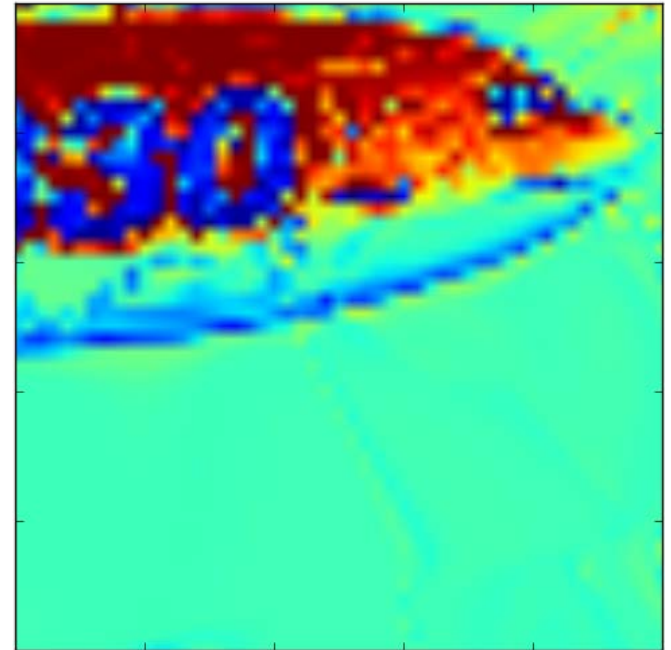
3 Layers: 14 nodes -> loss 7.55e-02

477 weights

# Adding a Third Layer



$$I = f(x, y)$$
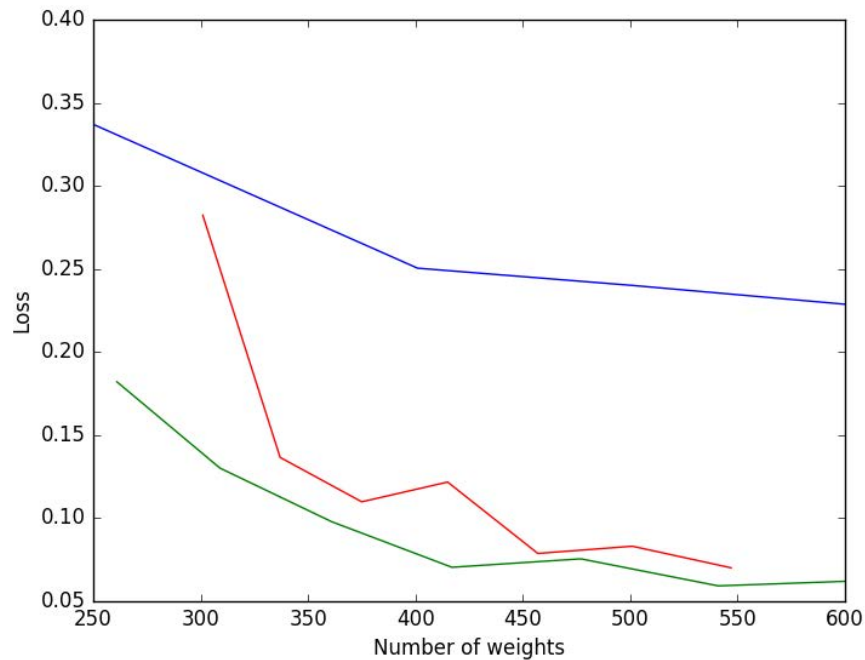
3 Layers: 15 nodes -> loss 5.93e-02     3 Layers: 19 nodes -> loss 4.38e-02

541 weights                                             837 weights

# Multi Layer Perceptrons



- Adding layers tends to deliver better convergence properties.
- In current practice, deeper is usually better.

# MLP to ResNet

Linear 2 -> n

tanh

Linear n -> n

tanh

Linear n -> n

tanh

Linear n -> 1

Linear 2 -> n

Block n -> n

Block n -> n

Linear n -> 1

x

Linear n -> n

$\sigma$

Linear n -> n

$x + l_2(\sigma(l_1(x))$

Further improvements in the convergence properties have been obtained by adding a bypass, which allows the final layers to only compute residuals.

# Improving the Network



Original 51x51 image:
2601 gray level values.

MLP 10/20/10 Interpolation:
471 weights, loss 6.43e-02.

MLP 10/20/10/10 Interpolation:
581 weights, loss 5.30e-2.

# Digital Images



```
136 134 161 159 163 168 171 173 173 171 166 159 157 155
152 145 136 130 151 149 151 154 158 161 163 163 159 151
145 149 149 145 140 133 145 143 145 145 145 146 148 148
148 143 141 145 145 145 141 136 136 135 135 136 135 133
131 131 129 129 133 136 140 142 142 138 130 128 126 120
115 111 108 106 106 110 120 130 137 142 144 141 129 123
117 109 098 094 094 094 100 110 125 136 141 147 147 145
136 124 116 105 096 096 100 107 116 131 141 147 150 152
152 152 137 124 113 108 105 108 117 129 139 150 157 159
159 157 157 159 135 121 120 120 121 127 136 147 158 163
165 165 163 163 163 166 136 131 135 138 140 145 154 163
166 168 170 168 166 168 170 173 145 143 147 148 152 159
168 173 173 175 173 171 170 173 177 178 151 151 153 156
161 170 176 177 177 179 176 174 174 176 177 179 155 157
161 162 168 176 180 180 180 182 180 175 175 178 180 180
```
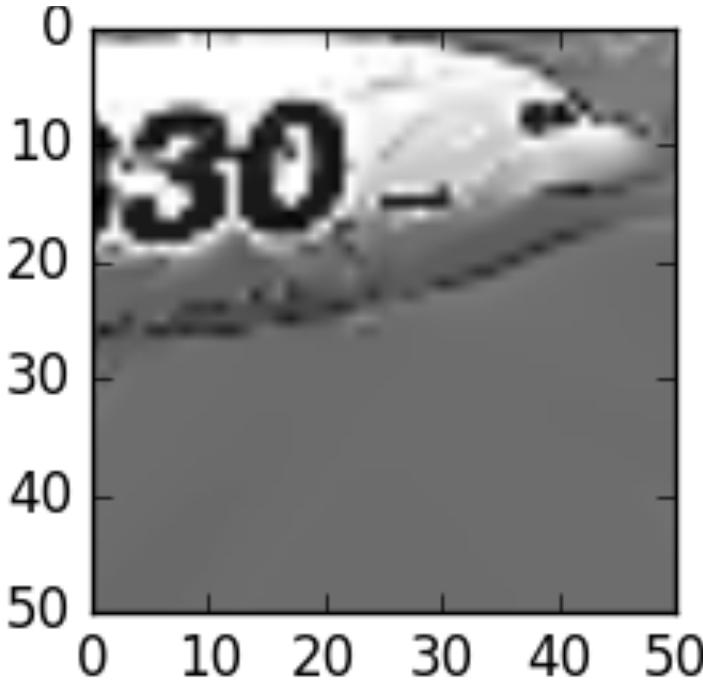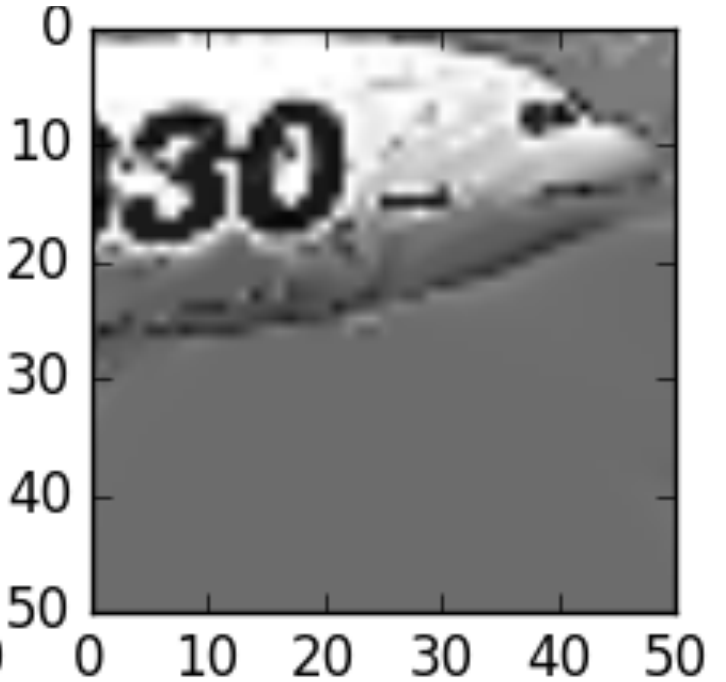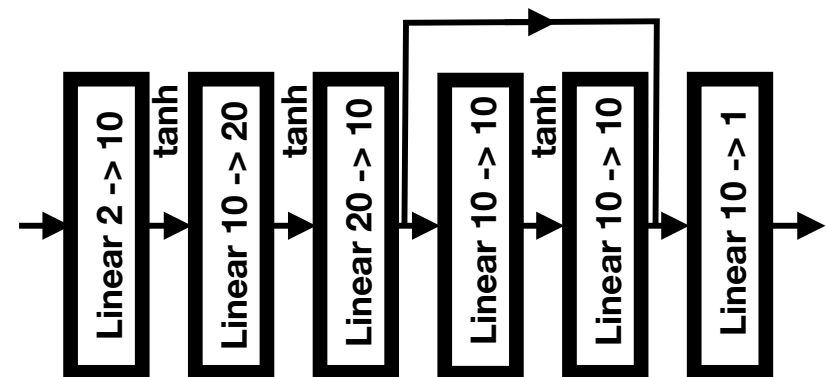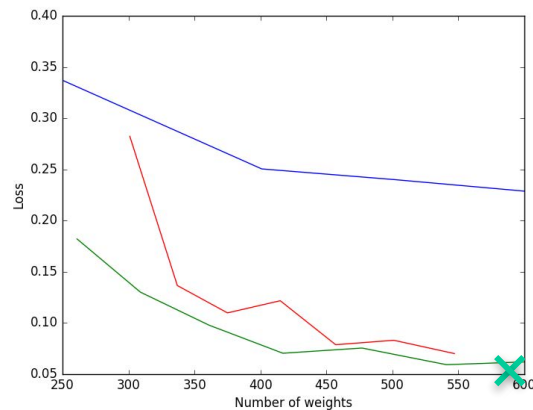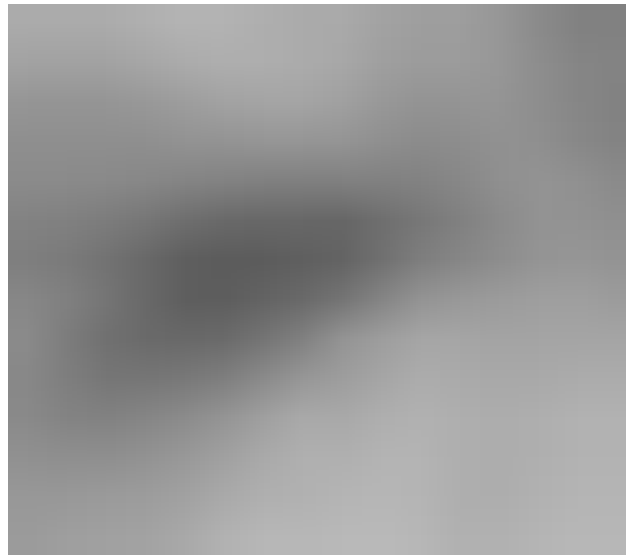
- A MxN image can be represented as an MN vector, in which case neighborhood relationships are lost.
- By contrast, treating it as a 2D array preserves neighborhood relationships.

# Image Specificities



- In a typical image, the values of <span style="color:red">neighboring pixels</span> tend to be more highly correlated than those of distant ones.

- An image filter should be translation invariant.

—> These two properties can be exploited to drastically reduce the number of weights required by CNNs using so-called convolutional layers.

# Fully Connected Layers



- The descriptive power of the net increases with the number of layers.
- In the case of a 1D signal, it is roughly proportional to $\prod_n W_n$ where $W_n$ represents the width of a layer.

# Convolutional Layer



$$\sigma \left( b + \sum_{x=0}^{n_x} \sum_{y=0}^{n_y} w_{x,y} a_{i+x,j+y} \right)$$

# Feature Maps

28 × 28 input neurons    first hidden layer: 3 × 24 × 24 neurons



Filters:

# Filters



g(x,y)

d/dx g(x,y)    d/dy g(x,y)

d/dx2 g(x,y)    d/dy2 g(x,y)    d/dxdy g(x,y)

d/dx3 g(x,y)    d/dy3 g(x,y)    d/dx2dy g(x,y)    d/dxdy2 g(x,y)

Derivatives

(c)

Learned filters

# Pooling Layer



hidden neurons (output from feature map)

max-pooling units

- Reduces the number of inputs by replacing all activations in a neighborhood by a single one.
- Can be thought as asking if a particular feature is present in that neighborhood while ignoring the exact location.

# Adding the Pooling Layers



28 × 28 input neurons

3 × 24 × 24 neurons

3 × 12 × 12 neurons

The output size is reduced by the pooling layers.

# Adding a Fully Connected Layer



- Each neutron in the final fully connected layer is connected to all neurons in the preceding one.
- Deep architecture with many parameters to learn but still far fewer than an equivalent multilayer perceptron.

# LeNet (1989-1999)

# Lenet Results

# AlexNet (2012)



Task: Image classification
Training images: Large Scale Visual Recognition Challenge 2010
Training time: 2 weeks on 2 GPUs

Major Breakthrough: Training large networks has now been shown to be practical!!

# AlexNet Results



- At the 2012 ImageNet Large Scale Visual Recognition Challenge, AlexNet achieved a top-5 error of 15.3%, more than 10.8% lower than the runner up.
- Since 2015, networks outperform humans on this task.

# Feature Maps



First convolutional layer



Second convolutional layer

- Some of the convolutional masks are very similar to oriented Gaussian or Gabor filters.
- The trained neural nets compute oriented derivatives, which the brain is also **believed** to do.

# Reminder: Discrete 2D Convolution

Input image: f

Convolved image: m**f



Convolution mask m, also known as a *kernel*.

$$\begin{bmatrix} m_{11} & \dots & m_{1w} \\ \dots & \dots & \dots \\ m_{w1} & \dots & m_{ww} \end{bmatrix}$$

$$m**f(x,y) = \sum_{i=0}^{w} \sum_{j=0}^{w} m(i,j) f(x-i, y-j)$$

# Reminder: 3X3 Masks

x derivative         y derivative



$$\begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix} \text{ and } \begin{bmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix} \qquad \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix} \text{ and } \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{bmatrix}$$

Prewitt operator                    Sobel operator

# Filter Banks



Derivatives of order
0, 1, and 2.

Learned

(c)

# Bigger and Deeper



VGG19, 3 weeks of training.

GoogleLeNet.

"hibiscus"     "dahlia"

"It was demonstrated that the representation depth is beneficial for the classification accuracy, and that state-of-the-art performance on the ImageNet challenge dataset can be achieved using a conventional ConvNet architecture."

# Deeper and Deeper



$$x + l_2(\sigma(l_1(x)))$$

Resnet

# Without Max Pooling



| Accuracy | Train | Test |
|---|---|---|
| Conv 5x5, stride 1 Max pool 2x3 | 99.58 | 98.77 |
| Conv 5x5, stride 2 | 99.42 | 98.31 |
| Conv 5x5, stride 1 Conv 3x3, stride 2 | 99.38 | 98.57 |

Max pooling can be replaced by Gaussian convolutions with stride > 1 .

# Image Classification Taxonomy 1989 —2016

LSTM
(Hochreiter and Schmidhuber, 1997)

LeNet5
(LeCun et al., 1989)

Bigger + GPU

Deep hierarchical CNN
(Ciresan et al., 2012)

Bigger + ReLU + dropout

No recurrence

Fully convolutional

AlexNet
(Krizhevsky et al., 2012)

MLPConv

Bigger + small filters

Overfeat
(Sermanet et al., 2013)

Net in Net
(Lin et al., 2013)

Highway Net
(Srivastava et al., 2015)

Inception modules

VGG
(Simonyan and Zisserman, 2014)

GoogLeNet
(Szegedy et al., 2015)

No gating

Batch Normalization

ResNet
(He et al., 2015)

BN-Inception
(Ioffe and Szegedy, 2015)

Wider

Dense pass-through    Aggregated channels

Wide ResNet
(Zagoruyko and Komodakis, 2016)

DenseNet
(Huang et al., 2016)

ResNeXt
(Xie et al., 2016)

Inception-ResNet
(Szegedy et al., 2016)

# ResNet to U-Net



$$x+l_2(\sigma(l_1(x)))$$

**ResNet block**

**U-Net**

—> Add skip connection to produce an output of the same size as the input.

Ronneberger, MICCAI'15

# Training a U-Net

Train Encoder-decoder U-Net architecture using binary cross-entropy



Minimize

$$L_{bce}(\mathbf{x}, \mathbf{y}; \mathbf{w}) = -\frac{1}{i} \sum_{1}^{P} [y_{\mathbf{i}} \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

where

- $\hat{\mathbf{y}} = f_{\mathbf{w}}(\mathbf{x})$,

- $\mathbf{x}$ in an input image,

- $\mathbf{y}$ the corresponding ground truth.

# Network Output



Image

BCE Loss

Ground truth

- Good start but not the end of the story.
- We will discuss this again during the delineation lecture.

# Streets Of Toronto



False negatives
False positives

# Transformers



- Break up the images into square patches.
- Transform each path into a feature vector.
- Feed to a transformer architecture.

# Self Attention



Given $\mathbf{X} = [\mathbf{x_1}, \ldots, \mathbf{x_I}]$:

- $a[\mathbf{x_i}, \mathbf{x_j}]$ is the attention that $\mathbf{x_i}$ gives to $\mathbf{x_j}$ . It measures the influence of one on the other.
- It can be computed for all I and j using far fewer weights that in a fully connected layer.

—> Provides context.

# Transformer Layer



$$\mathbf{X} \leftarrow \mathbf{X} + Sa(\mathbf{X})$$

$$\mathbf{X} \leftarrow LayerNorm(\mathbf{X})$$

$$\mathbf{x}_i \leftarrow \mathbf{x}_i + mlp[\mathbf{x}_i] \quad \forall i$$

$$\mathbf{X} \leftarrow LayerNorm(\mathbf{X})$$

# Depth from Single Images



Input     MiDaS (MIX 6)     DPT-Hybrid     DPT-Large

- Pros: Good at modeling long range relationships.
- Cons: Flattening the patches looses some amount of information.

# U-NET + Transformers



(a) GroundTruth  (b) TransUNet  (c) R50-ViT-CUP  (d) AttnUNet  (e) UNet

- A CNN produces a low-resolution feature vector.
- A transformer operates on that feature vector.
- The upsampling is similar to that of U-Net

—> Best of both worlds?

# Regression



$$\min_{\mathbf{W}_l, \mathbf{B}_l} \sum_i ||\mathbf{F}(\mathbf{x}_i, \mathbf{W}_1, \ldots, \mathbf{W}_L, \mathbf{b}_1, \ldots, \mathbf{b}_L) - \mathbf{y}_i||^2$$

using
- stochastic gradient descent on mini-batches,
- dropout,
- hard example mining,
- ...........

# Body Pose Estimation

# People and their Clothes



- Regress the body pose.
- Drape the garments on the body.
- Enforce consistency.

Input          Ours *(raw)*      Ours *(post ref.)*

# Crowd Counting



EPFL at lunchtime: The colors denote crowd density.

Liu et al. , CVPR'22

# Brains vs Neural Networks

- Neural networks are said to "bio-inspired".
- An excellent marketing argument but how true is it?

# Monkey Cortex



Pink:    Feed forward.
Cyan:   Feed back.
Yellow: Horizontal

# Human Visual Cortex



Dorsal pathway

Ventral pathway

# Recognize And Classify: Animal /No Animal

Subjects must raise their hand if they see an animal:

- 60 images
- 1 image per second

→    Measure their reaction time.

# Reminder: Recurrent Pathways



"Shape stimuli are optimally reinforcing each other when separated in time by ~60 ms, suggesting an underlying recurrent circuit with a time constant (feedforward + feedback) of 60 ms."

# Adversarial Images

# Brains vs Neural Networks

- Neural networks are said to "bio-inspired".
- An excellent marketing argument but how true is it?

Not that good:

- Much feedback is involved in biological systems.
- We don't need large databases to learn.
- We are not as susceptible to adversarial examples.

Neural nets are powerful but not the final answer!

# XKCD's View On The Matter

# Deep Nets in Short

- Deep Neural Networks can handle huge training databases.
- When the objective function can be minimized, the results are outstanding.
- There are failure cases and performance is hard to predict.

—> Many questions are still open and there is much theoretical work left to do.

# Alpha Go



- Uses Deep Nets to find the most promising locations to focus on.

- Performs Tree based search when possible.

- Relies on reinforcement learning and other ML techniques to train.

—> Beat the world champion in 2017.

# Optional: ChatGPT



- At this point, the transformer layer is fed both the prompt and the already generated text.
- It uses this information to guess the next word.
- The process is then iterated.

Keys to ChatGPT success:
- The network looks back as far as needed.
- It uses a huge corpus.
- Human guided training.

Vaswani'17

# Optional: ChatGPT

**PV** How do I recover from a stall

If you are asking about a stall in an aircraft, the first thing to do is to identify the cause of the stall. This may involve checking the flight instruments, looking outside the aircraft for visual clues, and consulting the aircraft's flight manual. Once the cause of the stall has been determined, the appropriate recovery procedure can be initiated. This may involve reducing the angle of attack, increasing power, or taking other corrective action to restore the flow of air over the wings and regain control of the aircraft. It is important to follow the correct procedures for recovering from a stall, as attempting to recover improperly can make the situation worse. If you are not trained in how to recover from a stall, it is best to consult with a qualified instructor or pilot.

**A little late for that.**

**PV** How do I recover from a spin in a glider

Recovering from a spin in a glider involves a few key steps. First, you should try to stay calm and avoid panicking. This will help you think clearly and make the right decisions to recover from the spin.

Next, you should locate the spin recovery instructions in your glider's flight manual and follow them carefully. These instructions will be specific to your glider and will outline the correct procedure for recovering from a spin. In general, the steps for recovering from a spin in a glider will involve using the controls to stop the rotation of the aircraft and regain control.

**That's going to get you killed!!**

Here are some general guidelines for recovering from a spin in a glider:

1. Pull the control stick firmly and fully back to the neutral position. This will stop the rotation of the aircraft and help you regain control.

2. Apply full opposite rudder to the direction of the spin. This will help stop the rotation of the aircraft and straighten it out.

3. Once the rotation has stopped and the aircraft is flying straight and level again, gently ease the control stick forward to return to normal flight.

It's important to remember that recovering from a spin in a glider can be a challenging and potentially dangerous maneuver. If you're not sure how to recover from a spin, or

To recover from a stall in aviation, the pilot should follow these steps:

1. Reduce the angle of attack by pushing forward on the control column or yoke.
2. Apply full power to the engine(s).
3. Keep the wings level by using the ailerons or by turning the yoke in the direction opposite the wing that is dropping.
4. Once the airplane begins to pick up speed and the wing starts to produce lift again, gently pull back on the control column or yoke to level the wings and regain altitude.
5. Be aware of the altitude and terrain and take appropriate action to regain altitude if necessary.
6. Once the aircraft is out of stall, maintain a safe airspeed, and climb to a safe altitude.

It's important to note that these steps are general guidelines and may vary depending on the type of aircraft and the specific circumstances of the stall. It's also important to practice stall recovery in a simulated environment before attempting it in an actual flight.

**It is.**

Recovering from a spin in a glider typically involves the following steps:

1. Identify the spin and its direction.
2. Neutralize the controls by releasing any pressure on the ailerons, rudder and elevator.
3. Apply full opposite rudder to the spinning direction and hold it until rotation stops.
4. Apply forward stick to reduce the angle of attack and increase airspeed.
5. Level the wings and check for proper control response.
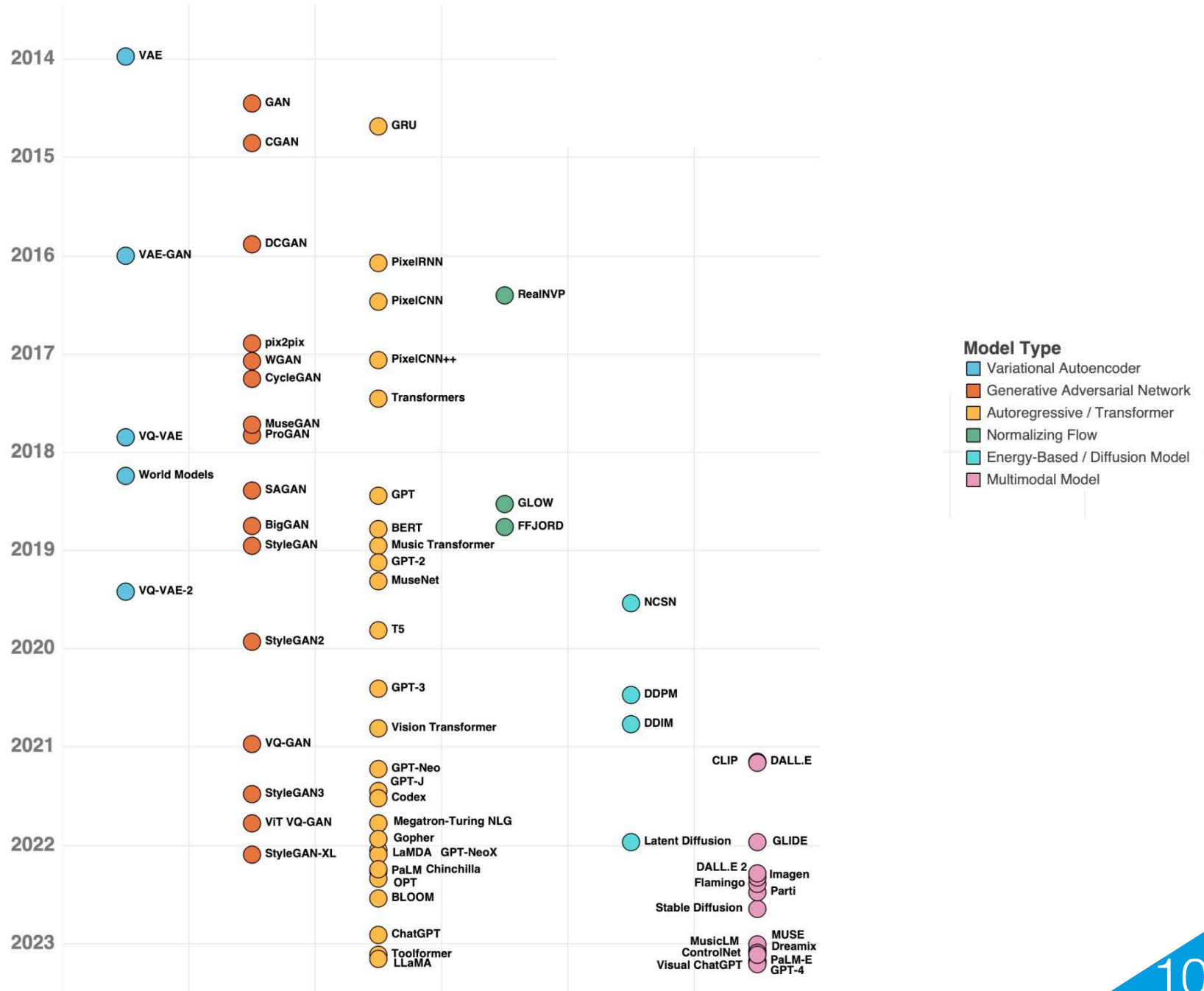
**Correct. You will live.**

It's important to note that a spin recovery procedure may vary depending on the type of glider and the specific training you have received. Always follow the instructions provided by the manufacturer or your instructor. Additionally, always make sure you are properly trained and familiar with the spin recovery procedure before attempting it.

27.01.2023

What changed? Presumably, enough people complained and the system was re-retrained with correct responses.

# Optional: Generative Models Taxonomy 2014 —2023



**Model Type**
- Variational Autoencoder
- Generative Adversarial Network
- Autoregressive / Transformer
- Normalizing Flow
- Energy-Based / Diffusion Model
- Multimodal Model

# What does it mean for Vision?

Two distinct approaches to increasing performance:

- Use ever larger training databases.
    - How do build them?
    - How do we tame the computational explosion?


- Use existing knowledge to reduce the need for training data.
    - Physics-based knowledge.
    - Geometrical knowledge.

—> Self-supervised methods.