

# Introduction to Deep Networks

It starts with regression.

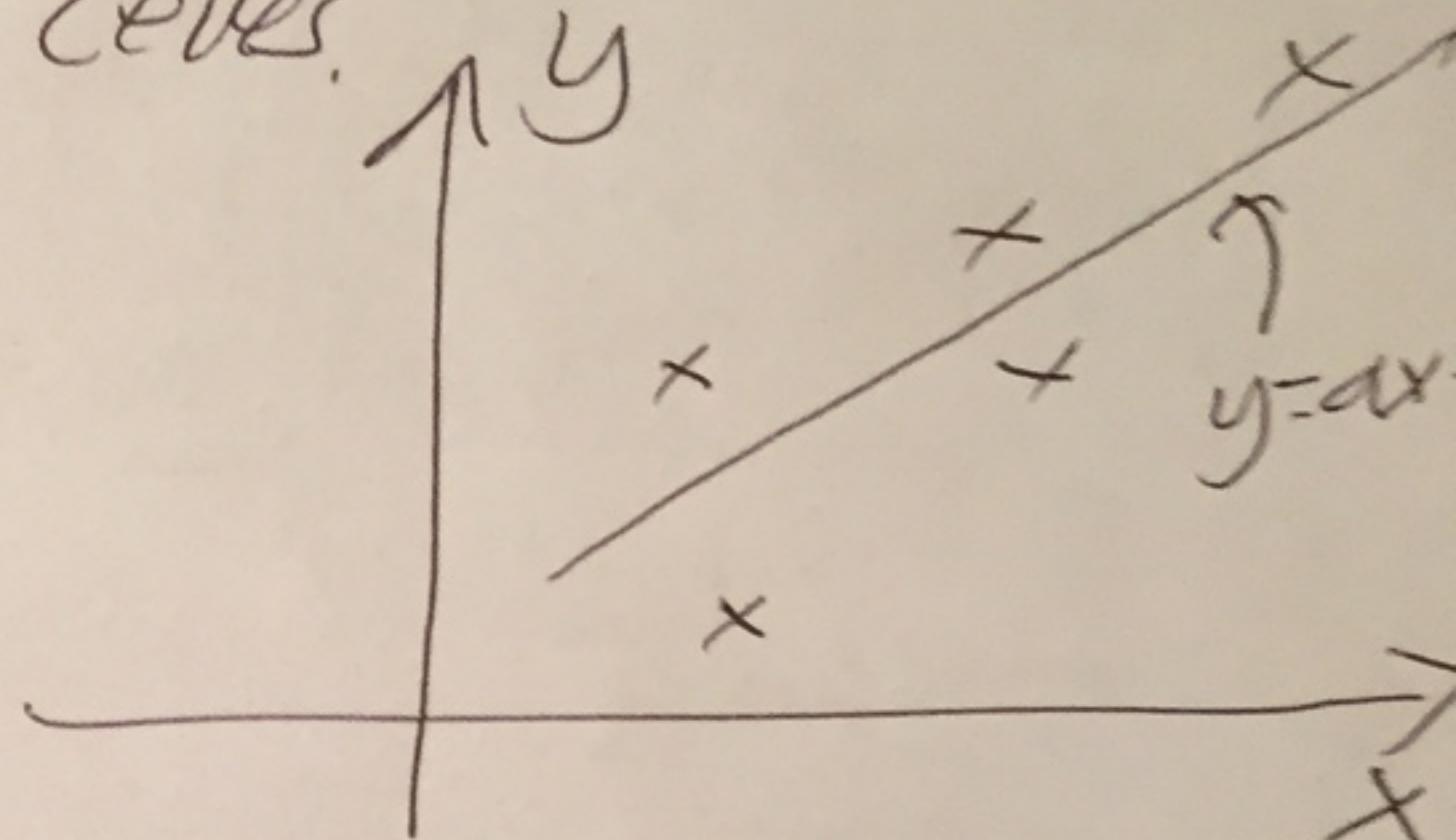
Gauss ~ 1800 - predicting the position of the planetoid Ceres.

set of measurements

$$\{(x_n, y_n) : n=1 \text{ to } N\}$$

model

$$y = ax + b + \underset{\substack{\leftarrow \\ \text{noise}}}{e}$$



Fit data by least squares.

$$E[a, b] = \frac{1}{N} \sum_{n=1}^N \{y_n - ax_n - b\}^2$$

$$\text{minimize } [\hat{a}, \hat{b}] = \arg \min_{[a, b]} E[a, b]$$

Probabilistic Regression.

$$(y - ax - b)^2 / 2\sigma^2$$

$$P(y | x; b, a) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y - ax - b)^2}{2\sigma^2}}$$

Estimate parameters by Maximum Likelihood (ML)

$$(\hat{a}, \hat{b}, \hat{\sigma}) = \arg \max_{(a, b, \sigma)} \prod_{n=1}^N P(y_n | x_n; b, a, \sigma)$$

$$= \arg \min_{(a, b, \sigma)} \left\{ - \sum_{n=1}^N \log P(y_n | x_n; b, a, \sigma) \right\}$$

$$\text{reduces to } [\hat{a}, \hat{b}] = \arg \min_{a, b} E[a, b], \text{ as before, } \hat{\sigma}^2 = \frac{1}{N} \sum_{n=1}^N (y_n - \hat{a}x_n - \hat{b})^2$$

# Introduction to Deep Networks

Perception → Rosenblatt → 1950's

$$y \in \{\pm 1\}$$

Task: classification.

Data  $\{(y^n, x^n) : n=1 \text{ to } N\}$ ,

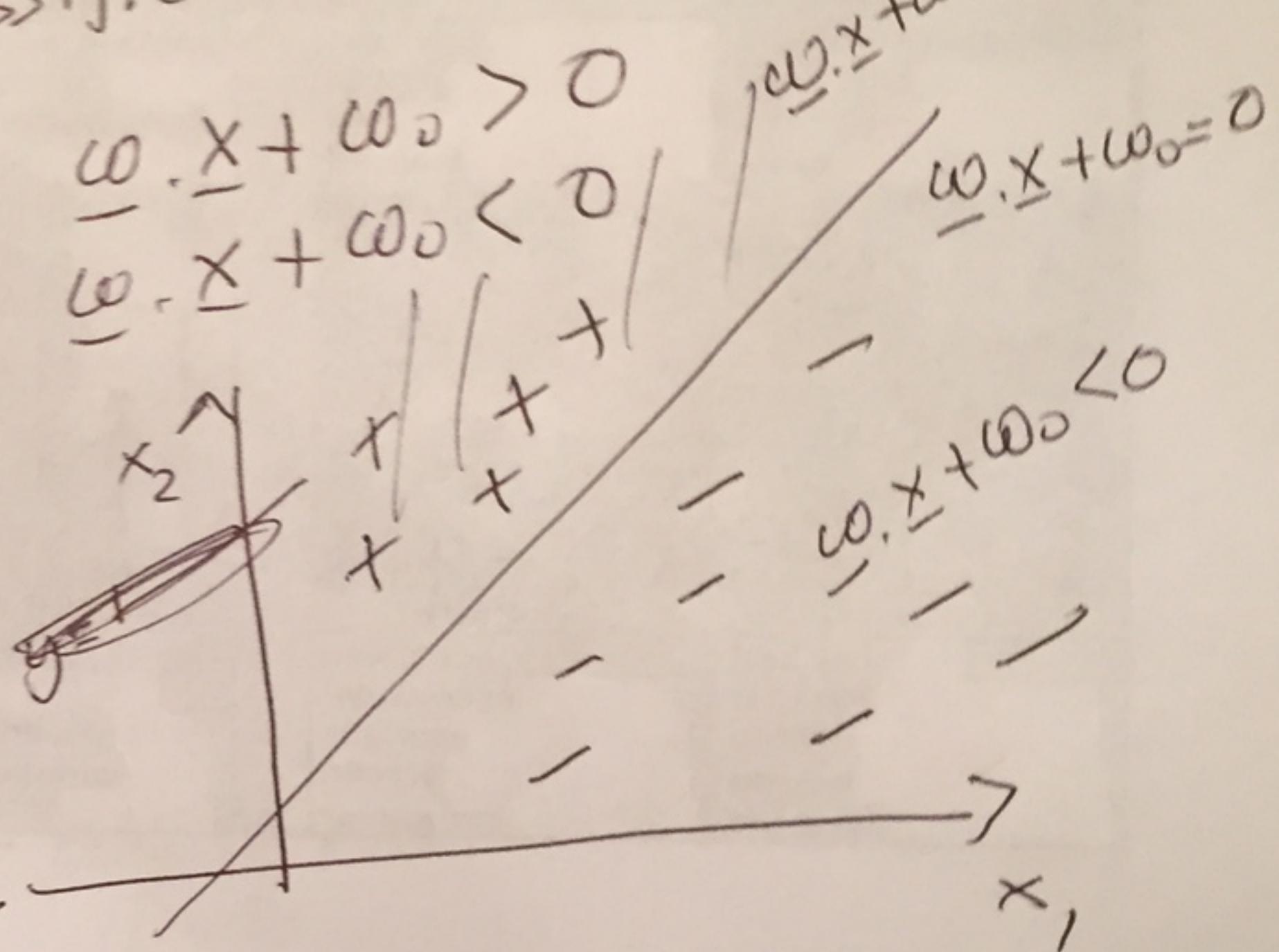
Try to learn a classifier / decision rule.

$$\hat{y}(x) = +1, \text{ if}$$

$$\hat{y}(x) = -1, \text{ if}$$

Here  $+$  is a datapoint  $(y^n, x^n)$  with  $y^n = 1$

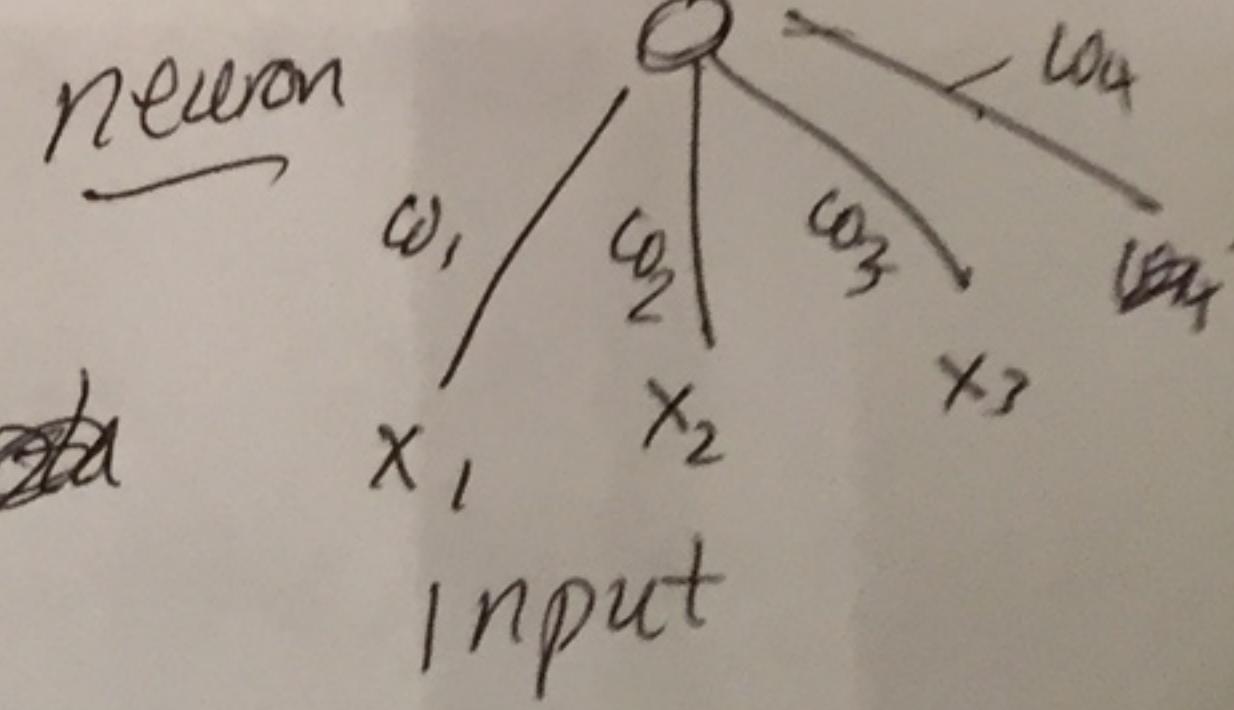
$-$  is a datapoint  $(y^n, x^n)$  with  $y^n = -1$ .



The perceptron algorithm learns the parameters  $w_0, w_0$  from the data  $\{(y^n, x^n) : n=1 \text{ to } N\}$

The perceptron algorithm was motivated by simplified models of neurons

Support Vector Machines (SVM)  
are descended from Perceptrons.  
But now we follow another path.



# Introduction to Deep Networks

Yuille

Yuille

Yuille

~~Task~~ Reformulate perceptrons as logistic regression.

$$y \in \{0, 1\}$$

$$P(y | \underline{x}) = \frac{e^{\underline{w} \cdot \underline{x} + w_0}}{e^{\underline{w} \cdot \underline{x} + w_0} + e^{-(\underline{w} \cdot \underline{x} + w_0)}}$$

input data

$$\{(y^n, \underline{x}^n) : n=1 \dots N\}$$

Estimate  $\underline{w}, w_0$  by Maximum likelihood

$$\hat{\underline{w}}, \hat{w}_0 = \arg \min_{\underline{w}, w_0} \left\{ - \sum_{n=1}^N \log P(y^n | \underline{x}^n, \underline{w}) \right\}$$

This The parameters  $\underline{w}, w_0$  can be estimated by steepest descent / gradient descent

$$\underline{w}^{t+1} = \underline{w}^t - \eta_t \frac{\partial}{\partial \underline{w}} \left( - \sum_{n=1}^N \log P(y^n | \underline{x}^n, \underline{w}) \right)$$

BATCH MODE.

or. ~~stochastic~~ stochastic gradient descent.

At time  $t$  pick ~~n~~  $(\underline{x}^n, y^{nH})$

$$\underline{w}^{t+1} = \underline{w}^t - \eta_t \frac{\partial}{\partial \underline{w}} \left( - \log P(y^n | \underline{x}^n, \underline{w}) \right)$$

# Introduction to Deep Networks

Yann LeCun

Alternative perspective:

This is regression with

$$P(y | \underline{x}; \underline{\omega}, \underline{b}) = \frac{e^{y \{ \underline{\omega} \cdot \underline{z}(\underline{x}, \underline{b}) \}}}{e^{\underline{\omega} \cdot \underline{z}(\underline{x}, \underline{b})} + e}$$

where  $\underline{z}(\underline{x}, \underline{b}) = \tanh(\underline{R} \cdot \underline{x})$ ,

Learn the parameters  $\underline{\omega}, \underline{b}$  by maximum likelihood (ML).

$$\hat{\underline{\omega}}, \hat{\underline{b}} = \arg \min_{\underline{\omega}, \underline{b}} \left\{ - \sum_{n=1}^N \log P(y_n | \underline{x}_n; \underline{\omega}, \underline{b}) \right\}$$

stochastic gradient descent

$$\underline{\omega}^{t+1} = \underline{\omega}^t - \eta_t \frac{\partial}{\partial \underline{\omega}} \left\{ - \log P(y_{n(t)} | \underline{x}_{n(t)}; \underline{\omega}, \underline{b}) \right\}$$

where  $(y_{n(t)}, \underline{x}_{n(t)})$  are the training pair selected at time  $t$ .

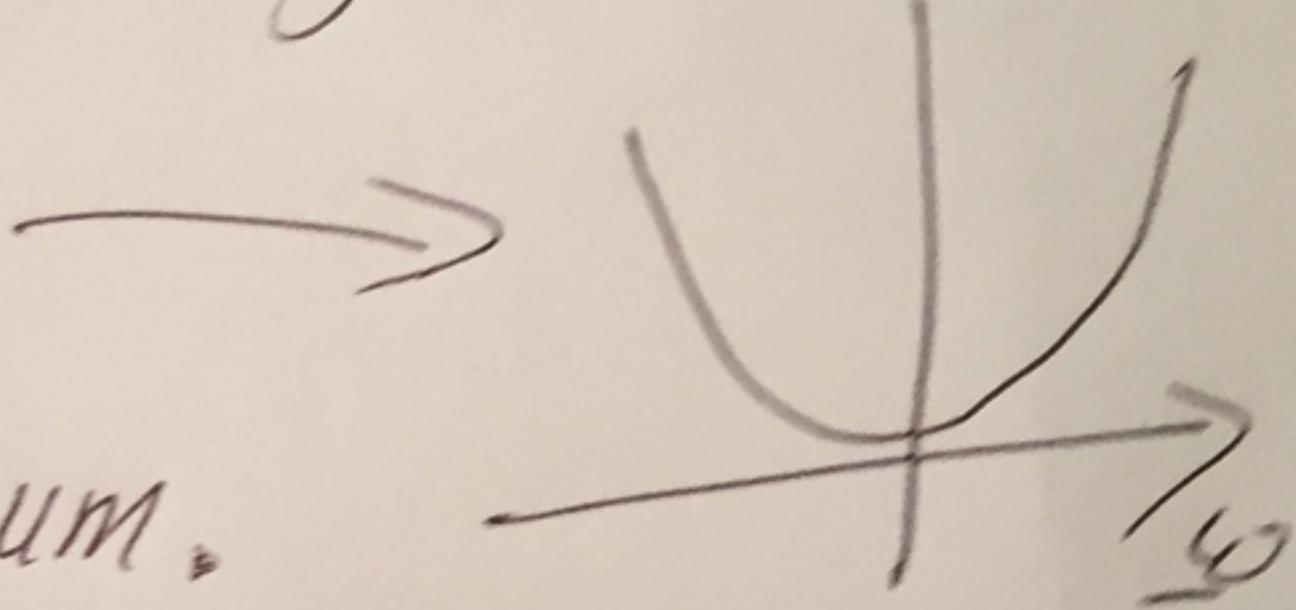
or gradient descent.

$$\underline{\omega}^{t+1} = \underline{\omega}^t - \eta_t \frac{\partial}{\partial \underline{\omega}} \left\{ - \sum_{n=1}^N \log P(y_n | \underline{x}_n; \underline{\omega}, \underline{b}) \right\}$$

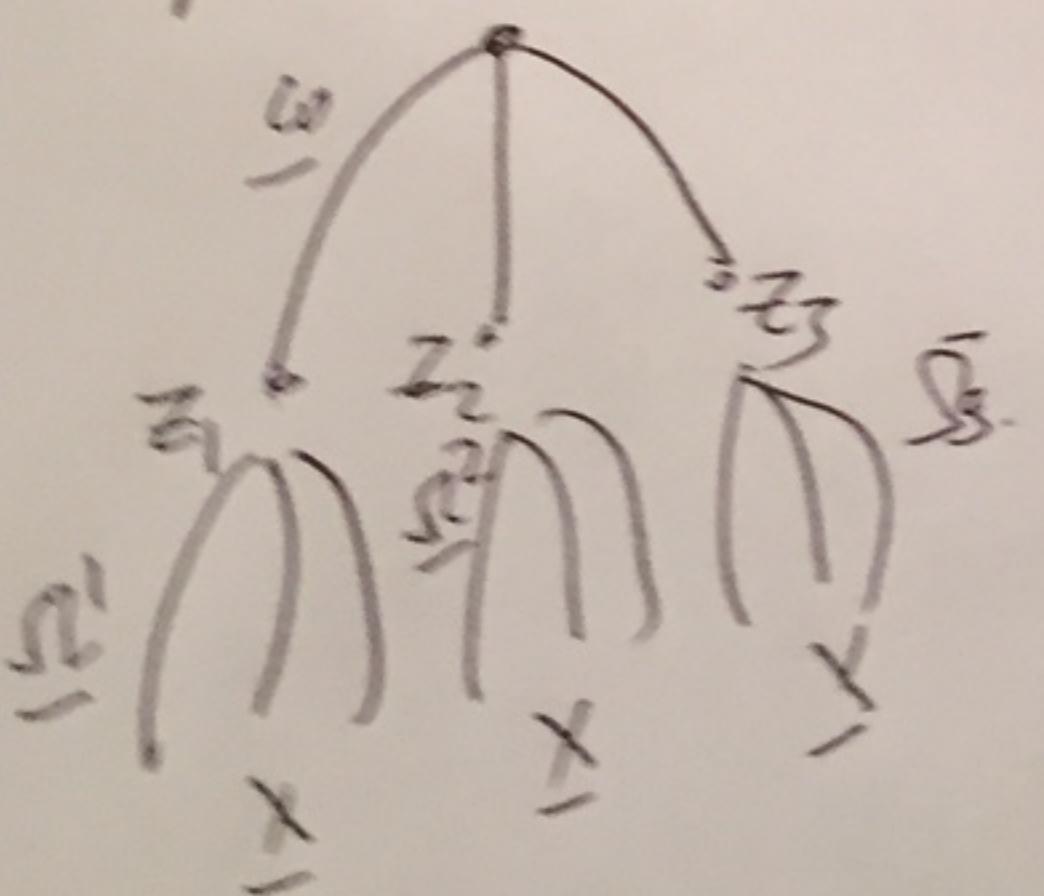
# Introduction to Deep Networks

Part 4.

For logistic regression the steepest descent algorithm is guaranteed to converge to the optimal solution. Technically, because  $-\log P(y|X, \underline{\omega})$  is a convex function of  $\underline{\omega}$  which is bounded below. So it has a single minimum.



In the 1970's researchers extended perceptrons to multi-layer perceptron



$$y = f(\underline{\omega}, \underline{R}, \underline{x})$$

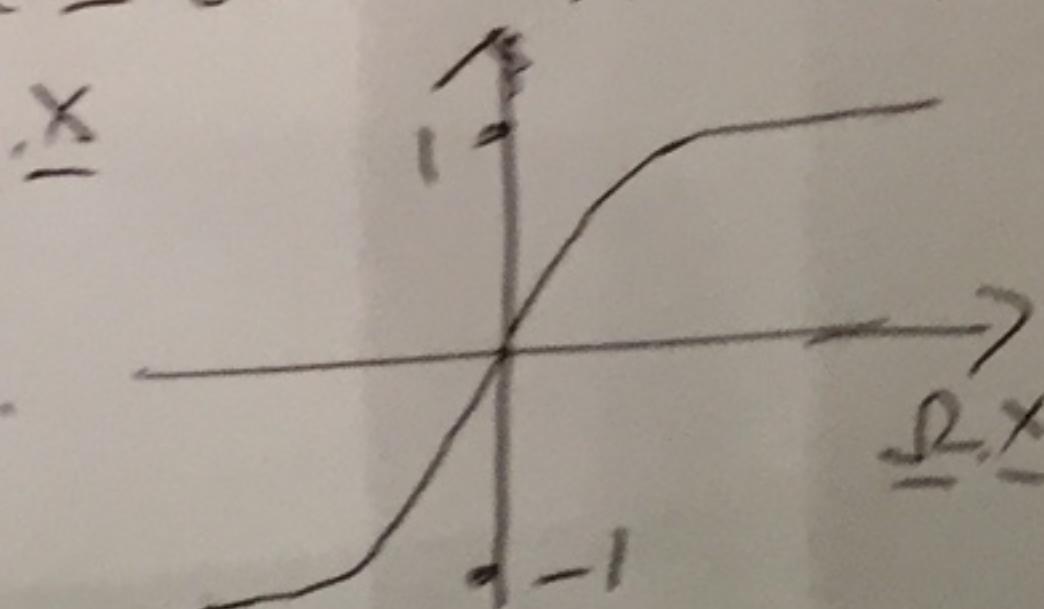
$$y = \tanh(\underline{\omega} \cdot \underline{z})$$

$$\underline{z} = (z_1, z_2, z_3)$$

$$z_i = \tanh(\underline{R}_i \cdot \underline{x})$$

$$\tanh(\underline{R}_i \cdot \underline{x}) = \frac{e^{\underline{R}_i \cdot \underline{x}} - e^{-\underline{R}_i \cdot \underline{x}}}{e^{\underline{R}_i \cdot \underline{x}} + e^{-\underline{R}_i \cdot \underline{x}}} = \tanh(\underline{R} \cdot \underline{x})$$

Key result: the output  $y$  is a differentiable function of the parameters  $\underline{\omega}$  &  $\underline{R}$ .

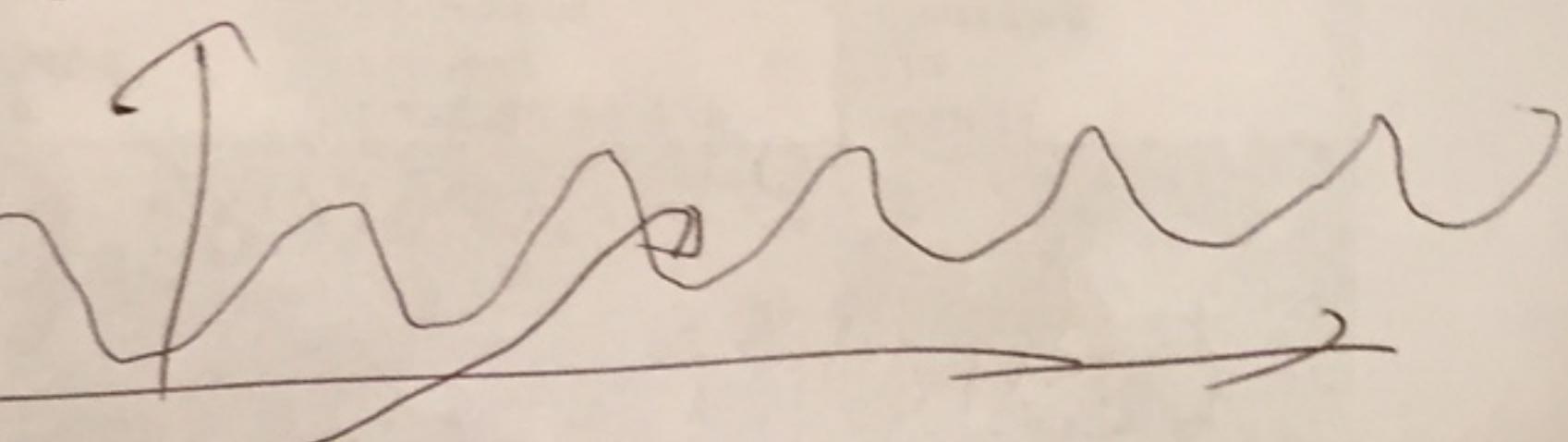


# Introduction to Deep Networks

Yville 66.

For multilayer perceptron, there is no guarantee of convergence to the global minimum of  $-\sum_{n=1}^N \log P(y_n | x_n; w, b)$ .

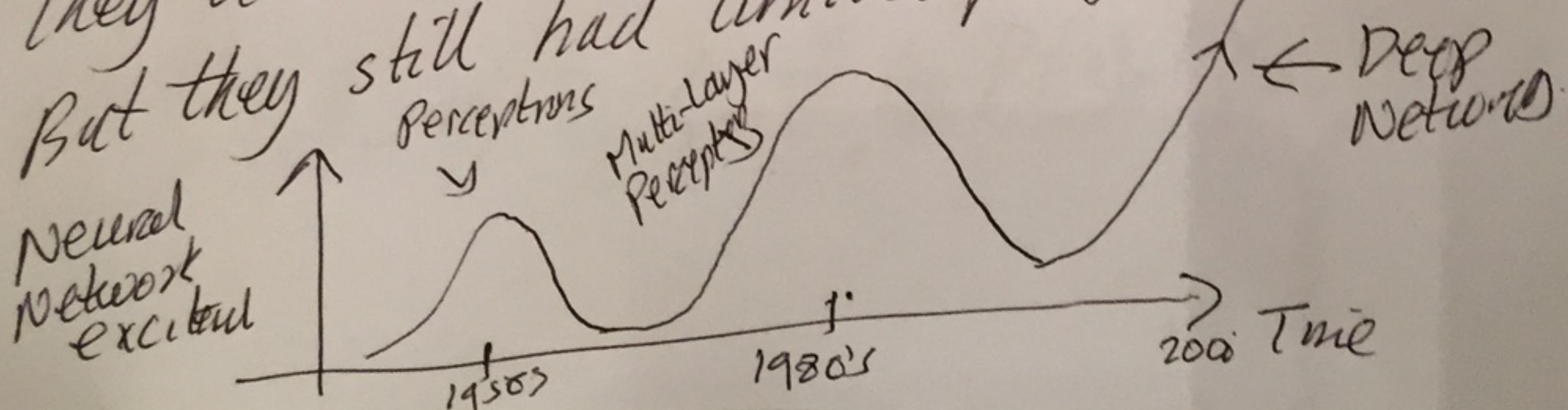
This function is not convex.

Easy for a learning algorithm to get  trapped in local minimum.

But, in practice, algorithms often converge to good solutions. Perhaps because there are many local minima which are all equally good solutions. (c.f. Prof. René Vidal)

Multilayer Perceptrons can represent many more decision rules than Perceptrons so they are more effective.

But they still had limited performance



## Introduction to Deep Networks

Yuike 7.

Deep Network extend multi layer perceptrons by having convolutional filters & filterbanks (see handout slides)

They are more effective than 1980's multi layer perceptrons because:

(i) they are much bigger

AlexNet  $\rightarrow$  650,000 neurons

60,000,000 parameters

630,000,000 connections.

(ii) they can be implemented in Graphical Processing Units (GPUs), which were developed for video games.

(iii) there is now enough data to train them  $\{(x_n, y_n) : n=1 \text{ to } N\}$

$N \approx 1,000,000$  or more.

Researchers in the 1980's did not have the technology - computers, data.