



Agenda

Linear Algebra

Calculus

A Brief History of Neural Networks

Neural Networks in Spark

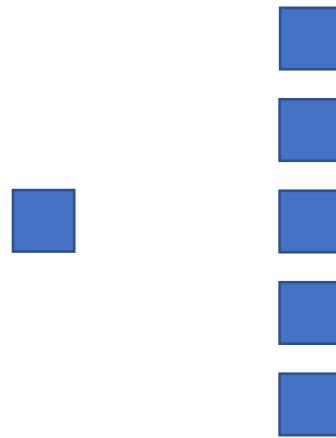


Math Is Fun

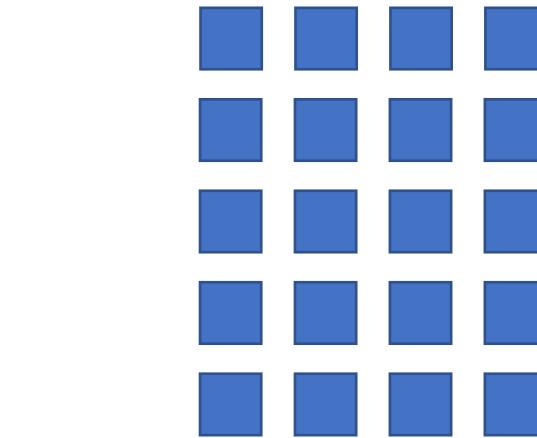
Linear Algebra and Calculus

- Scalars, Vectors, Matrices, and Tensors
- Operations
- Norms
- Special Matrices and Inversion
- Derivatives

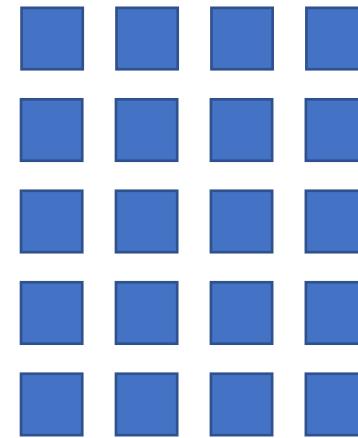
Scalars, Vectors, Matrices, and Tensors



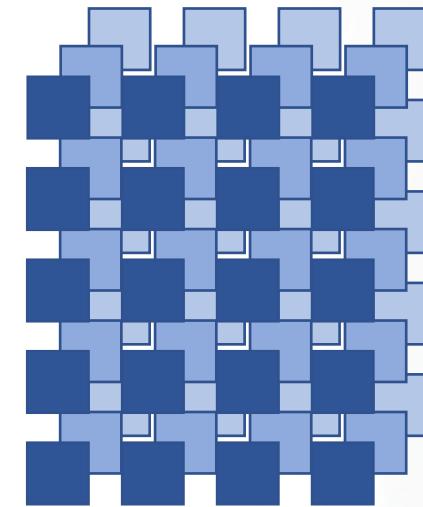
Scalar



Vector



Matrix



Tensor



Operations

- Addition
- Transpose
- Multiplication

Addition

Scalar: $2 + 3 = 5$

Vector: $\begin{bmatrix} 1 \\ 2 \end{bmatrix} + \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 4 \\ 6 \end{bmatrix}$

Matrix: $\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} + \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} = \begin{bmatrix} 6 & 8 \\ 10 & 12 \end{bmatrix}$

Transpose

Scalar: $2^T = 2$

Vector: $\begin{bmatrix} 1 \\ 2 \end{bmatrix}^T = [1 \quad 2]$

Matrix: $\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}^T = \begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix}$

Multiplication

Two types of multiplication

- Inner Product
- Outer Product

Inner Product

- Dimension reducing

Scalar: $(2)(3) = 6$

Vector: $\begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} 3 \\ 4 \end{bmatrix} = [(1)(3) + (2)(4)] = [11]$

Matrix/Vector: $\begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} 5 \\ 6 \end{bmatrix} = [(1)(5) + (2)(6)]$
 $\qquad\qquad\qquad [(3)(5) + (4)(6)]$

Matrix: $\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} = \begin{bmatrix} [(1)(5) + (2)(7)] & [(1)(6) + (2)(8)] \\ [(3)(5) + (4)(7)] & [(3)(6) + (4)(8)] \end{bmatrix}$

Outer Product

- Dimension increasing

$$\text{Vector}/\text{Vector}: [1 \quad 2] \otimes \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 3 & 6 \\ 4 & 8 \end{bmatrix}$$

Vector \otimes Matrix \rightarrow Tensor(3d)

Norms

Length of a vector

$$\|X\| = \sqrt{x_1^2 + x_2^2 + \cdots + x_n^2} = \sqrt{X^T X}$$

$$\|X\|^2 = X^T X$$

Distance

$$\|X - Y\|$$

Identity and Inverses

Scalar: $(1)(5) = 5 \quad ab = 1 \rightarrow a = 1/b$

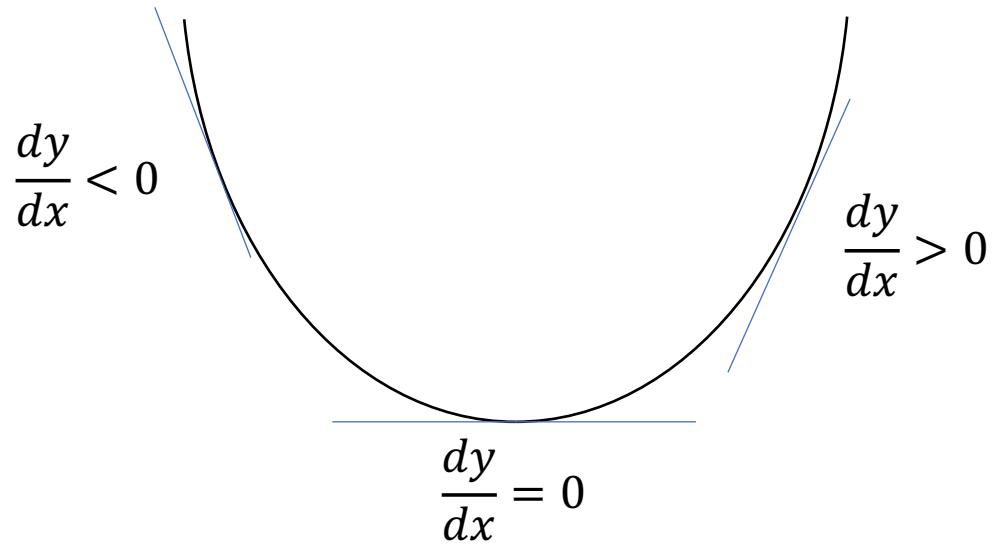
Matrices

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 2 & 3 \\ 4 & 5 \end{bmatrix} = \begin{bmatrix} 2 & 3 \\ 4 & 5 \end{bmatrix}$$

$$I = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad IA = AI = A$$

$$AA^{-1} = I$$

Calculus



Slope is zero at minimum

Calculus

Gradient Decent

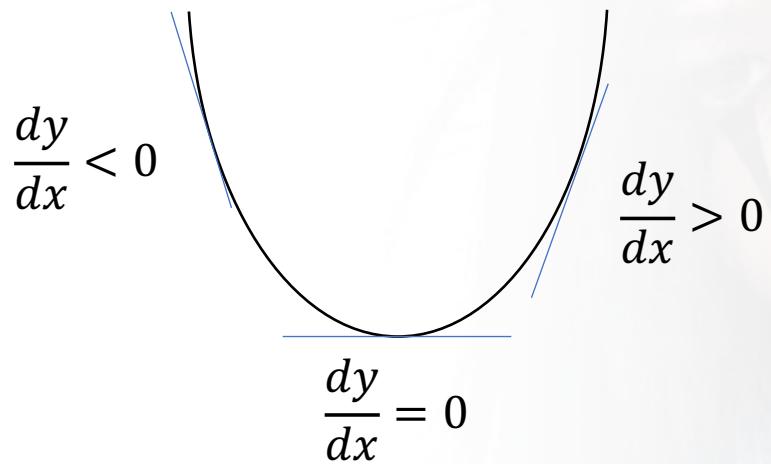
Find minimum of $f(x)$

pick x

while $0 < f(x)$

let $\Delta x = -\frac{df(x)}{dx}$

update $x = x + \Delta x$



Calculus

Chain Rule

$$z = f(y) \quad y = g(x)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

Linear Algebra and Calculus

- Scalars, Vectors, Matrices, and Tensors
- Operations
- Norms
- Special Matrices and Inversion
- Derivatives

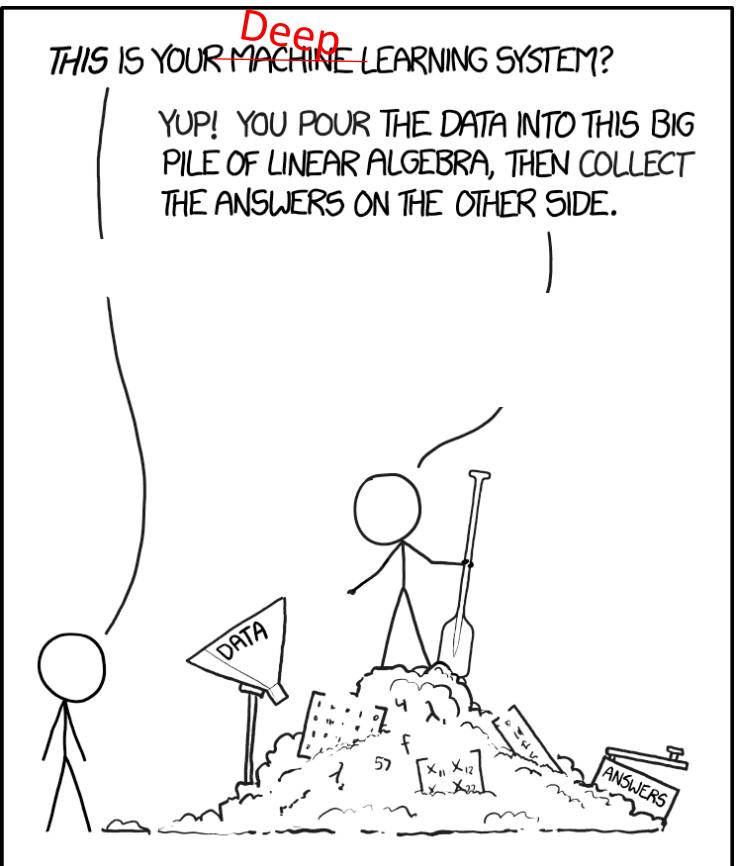
A close-up, low-angle shot of a character from the video game Cyberpunk 2077. The character has short, light-colored hair and is wearing a detailed, purple-toned leather jacket over a dark top. They are looking upwards and to the right. The background is dark and filled with glowing blue and red lights, suggesting a futuristic or cyberpunk setting.

A Brief History of Neural Networks

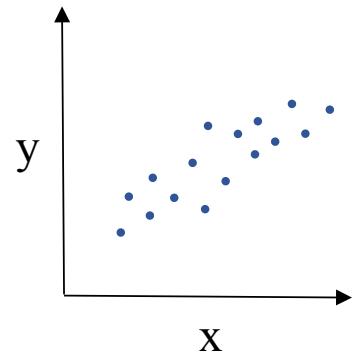
Neural Networks

- Linear Regression
- Logistic Regression
- Perceptrons
- Multilayered Perceptrons
- Support Vector Machines
- Deep Learning

After 120 Years of Research

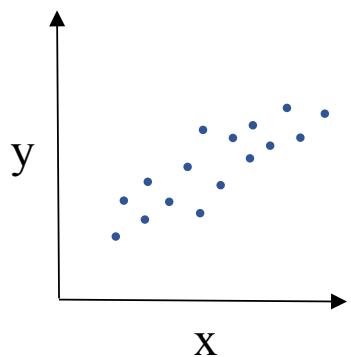


In The Beginning There Was Data



Data

Linear Regression (1894)



$$y = \sum_{i=1}^n w_i x_i$$

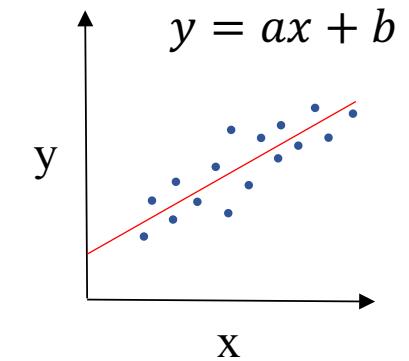
\mathbf{x}

\mathbf{w} Σ y

$$y = \mathbf{w}^T \mathbf{x}$$

Minimize
 $\|\mathbf{y} - \mathbf{X}\mathbf{w}\|$

Closed form solution
for \mathbf{w}

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X} \mathbf{y}$$


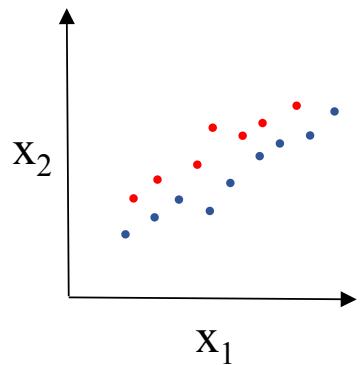
Data

Mechanism

Learning

Model

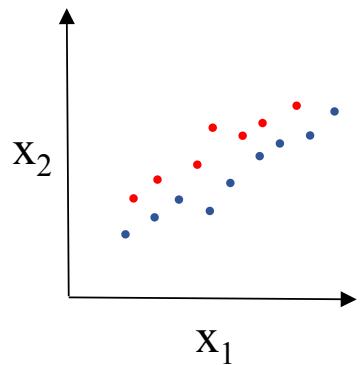
Logistic Regression (1958)



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

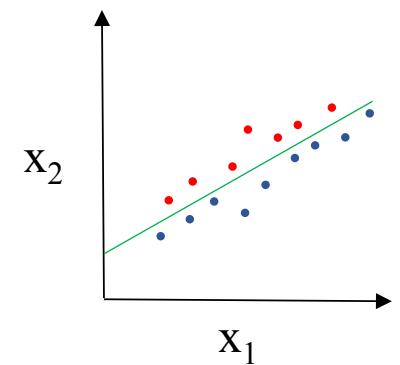
Data

Logistic Regression (1958)



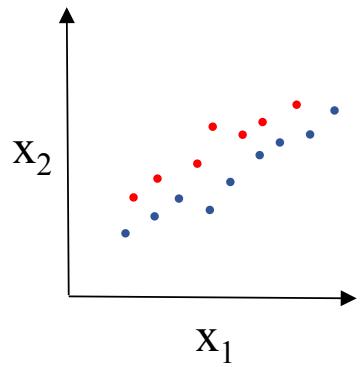
$$y \in \{0,1\}$$

Data



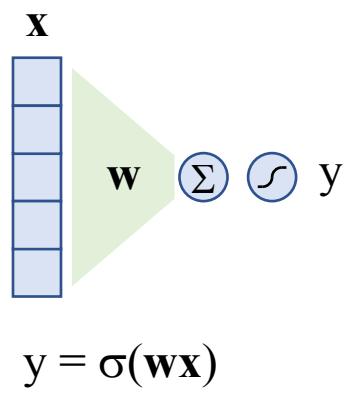
Model

Logistic Regression (1958)

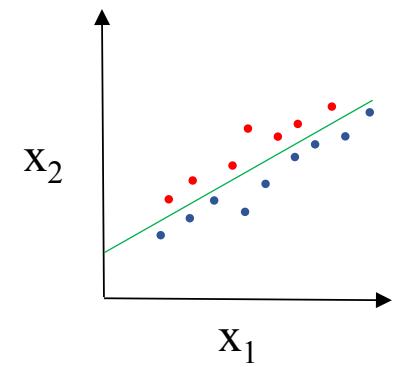


$$y \in \{0,1\}$$

Data

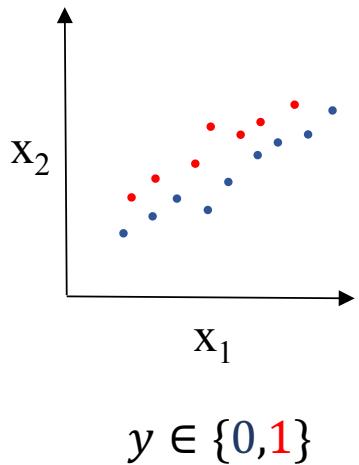


Mechanism

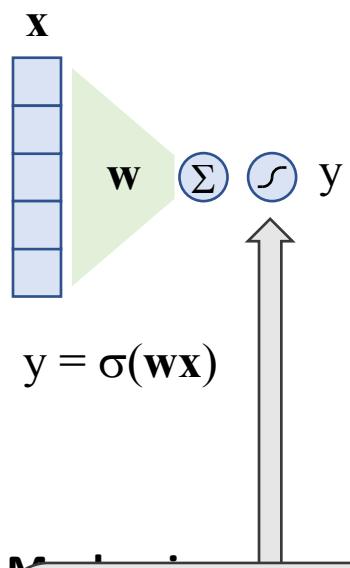


Model

Logistic Regression (1958)



Data

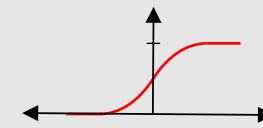
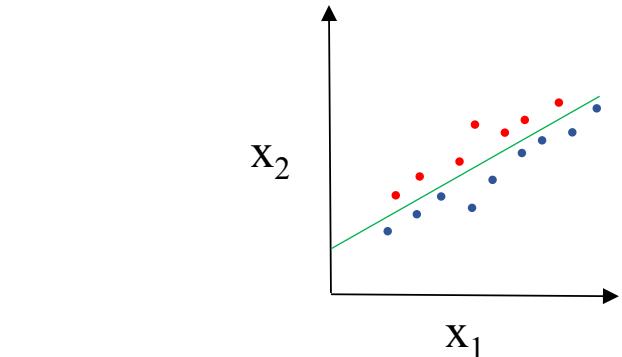


Activation Function

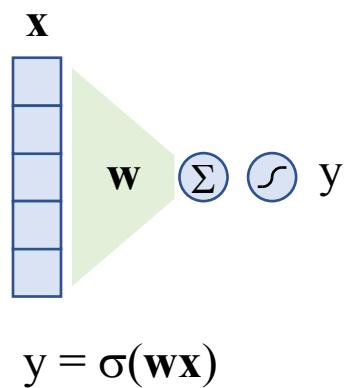
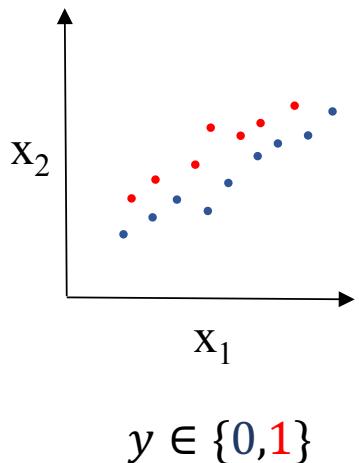
A nonlinearity to restrict range of output

Probability of classification:

$$\sigma(x) = 1/(1+e^{-x})$$



Logistic Regression (1958)

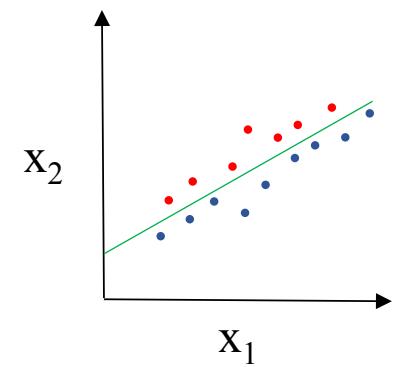


Data

Minimize
 $\|\mathbf{y} - \sigma(\mathbf{X}\mathbf{w})\|$

No closed form solution for \mathbf{w}

Use gradient descent
(Newton's Method)

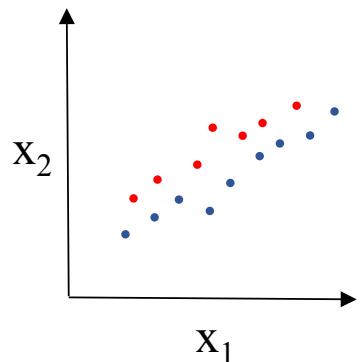


Mechanism

Learning

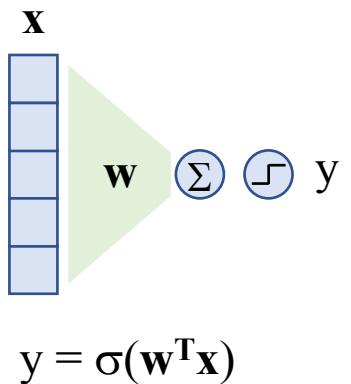
Model

Perceptron (1957)

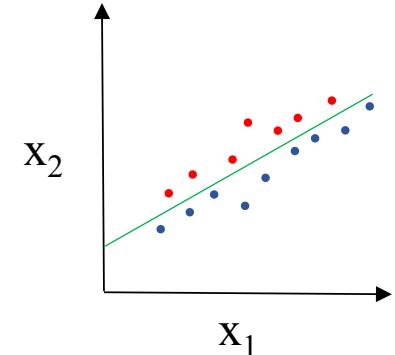


$$y \in \{0,1\}$$

Data

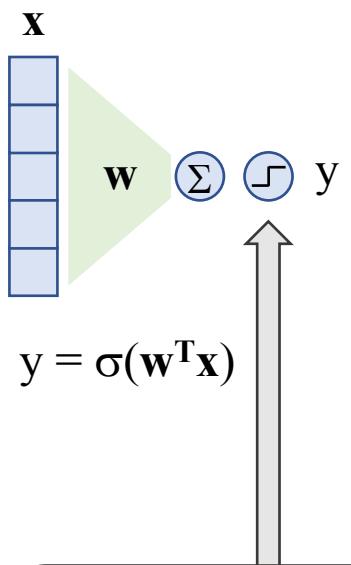
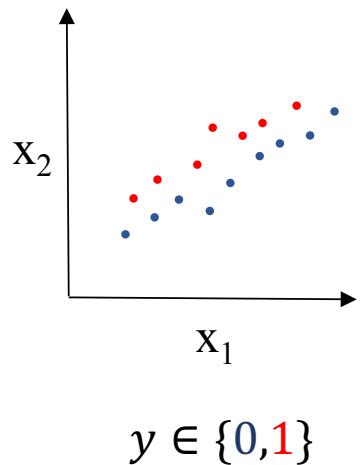


Mechanism



Model

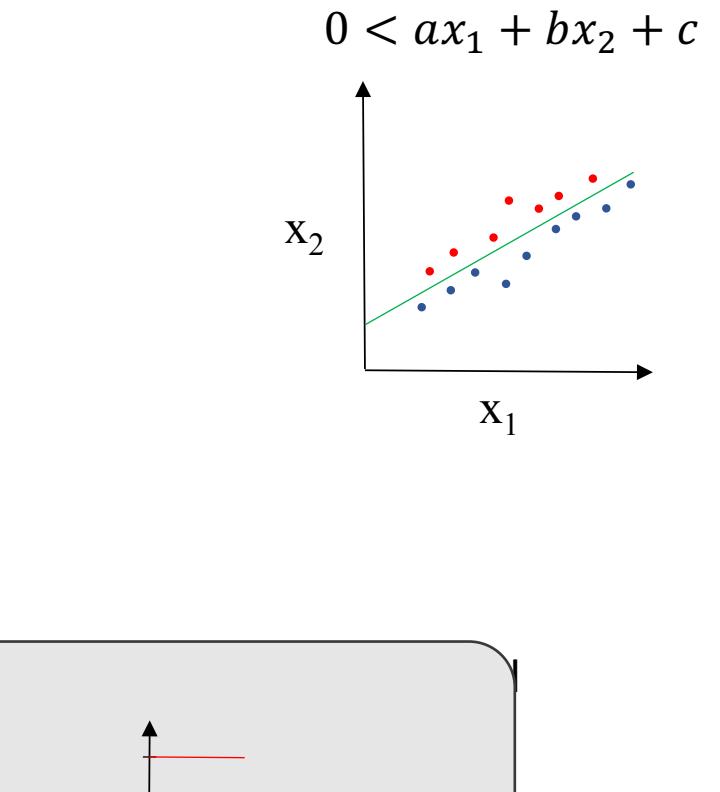
Perceptron (1957)



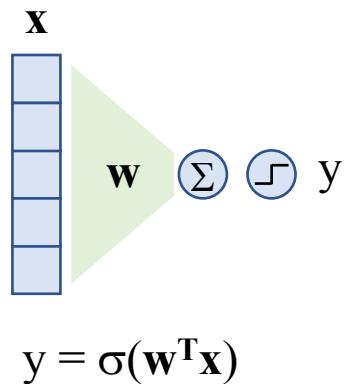
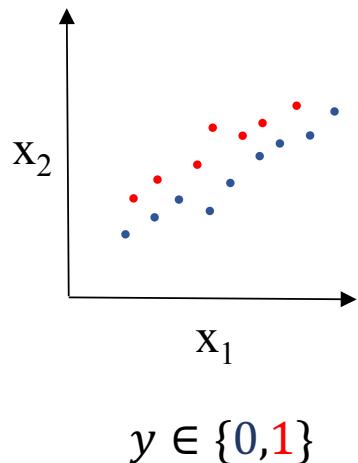
Data

Activation Function

$$\sigma(x) = 1 \text{ if } x > 0, 0 \text{ otherwise}$$



Perceptron (1957)

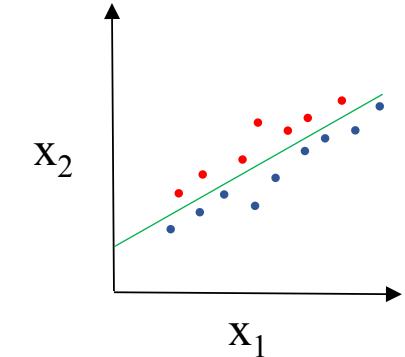


Iterate over input output pairs

$$\mathbf{w}_{t+1} = \mathbf{w}_t + (y - \sigma(\mathbf{w}_t \mathbf{x})) \mathbf{x}$$

Will converge if separable

$$0 < ax_1 + bx_2 + c$$



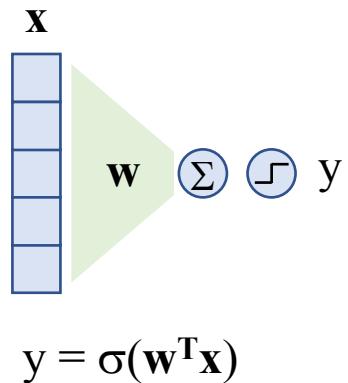
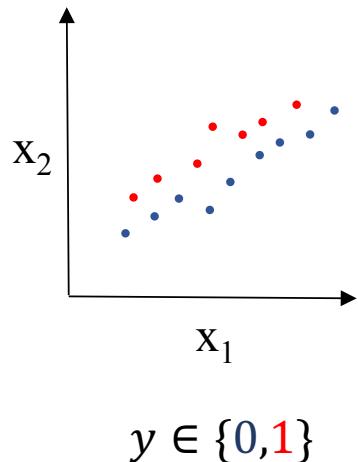
Data

Mechanism

Learning

Model

Perceptron (1957)



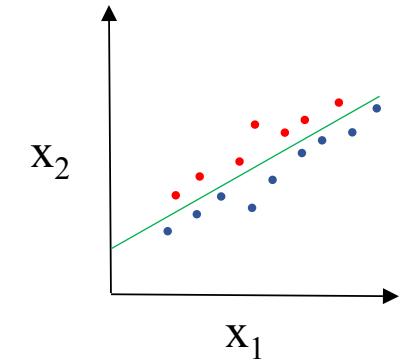
Data

Iterate over input output pairs

$$\begin{aligned}\mathbf{e} &= \mathbf{y} - \sigma(\mathbf{W}_t \mathbf{x}) \\ \Delta \mathbf{W}_t &= \mathbf{e} \otimes \mathbf{x}\end{aligned}$$

Will converge if separable

$$0 < ax_1 + bx_2 + c$$



Mechanism

Learning

Model

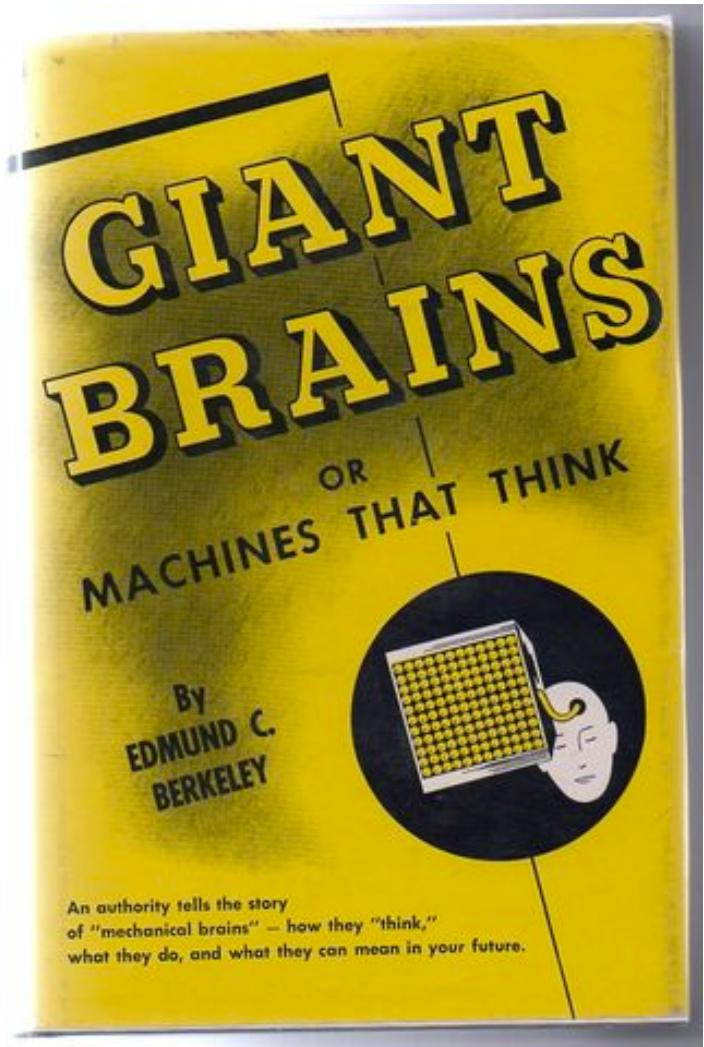
Perceptrons (1957)

$$y = \sigma(w \cdot x)$$

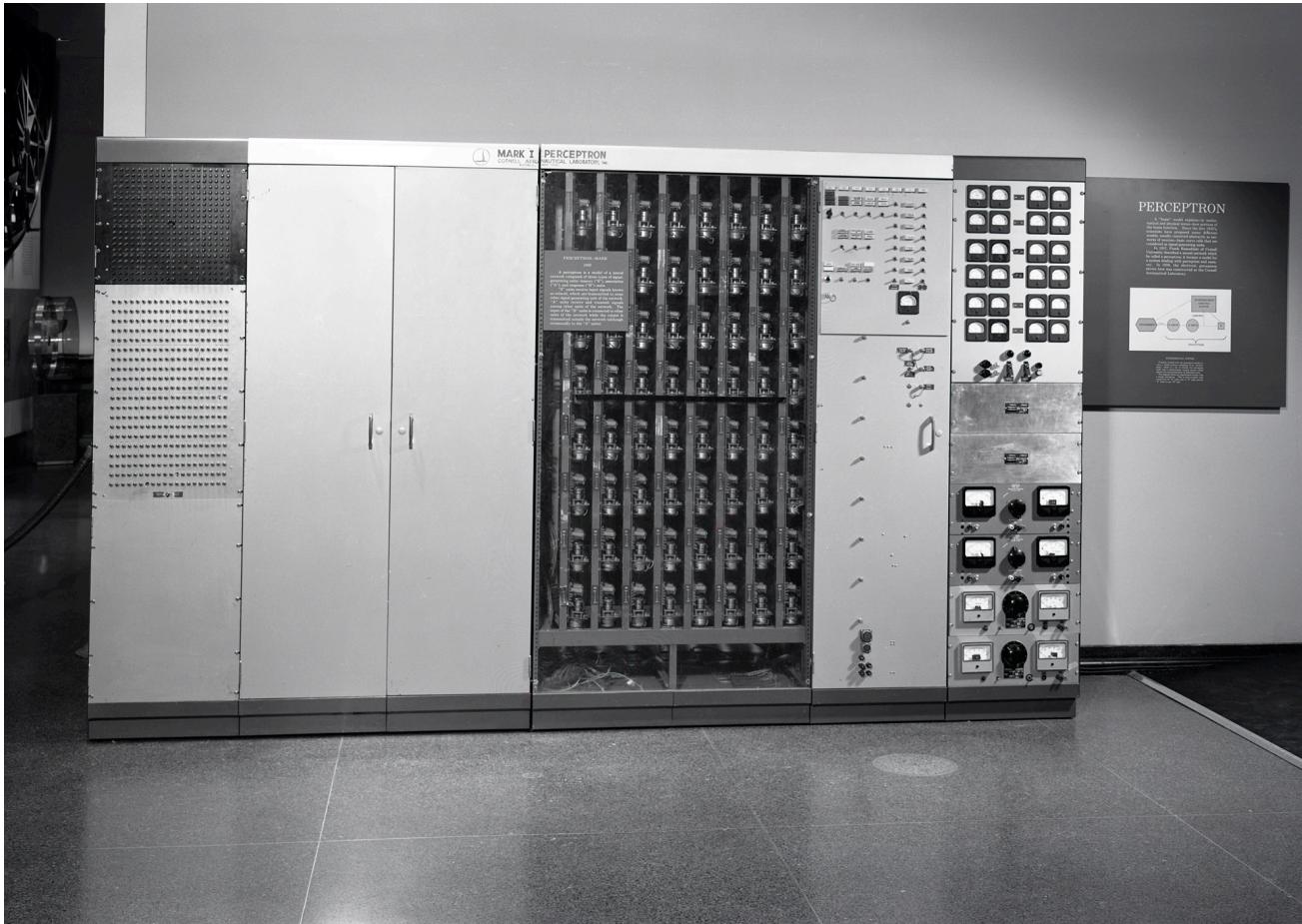
$$\Delta w$$

$$\Delta W = e \otimes x$$

Neural Networks Are Great!



Mark 1 Perceptron Hardware



Perceptrons (1960s)



Perceptrons (1960s)

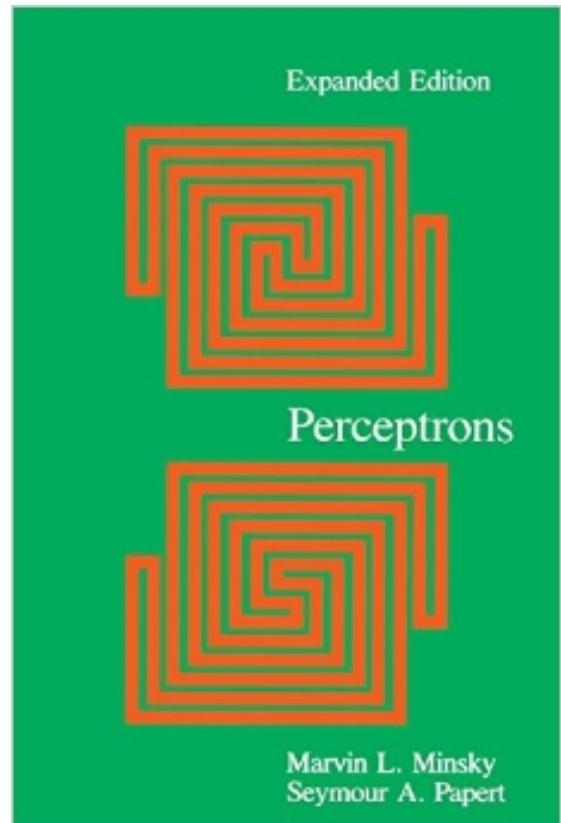


Cloudy



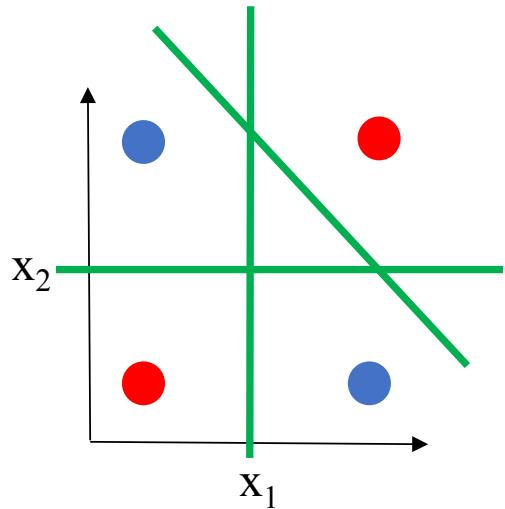
Sunny

Twenty Year Intermission



1969

Twenty Year Intermission



A failure of
functionality

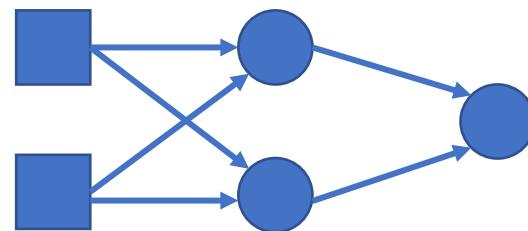
Rise of Symbolic AI (1970s & 1980s)

- Physical Symbol System Hypothesis
- Expert Systems
- Game playing (Chess & Checkers)
- Common sense knowledge bases (Cyc)
- ...

Thinking is easy, perception is hard

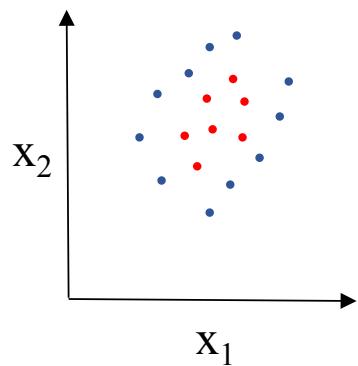
Multilayered Perceptrons

- Two layered networks can classify XOR data set
 - $\text{XOR}(x,y) = \text{OR}(\text{AND}(x,y), \text{NOR}(x,y))$



- No way to learn it

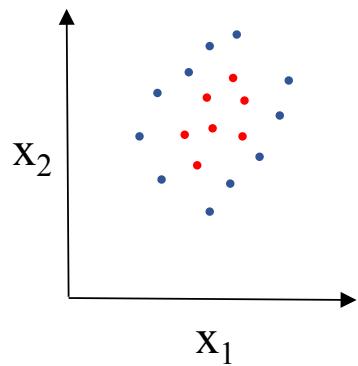
Multilayered Perceptrons (1986)



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

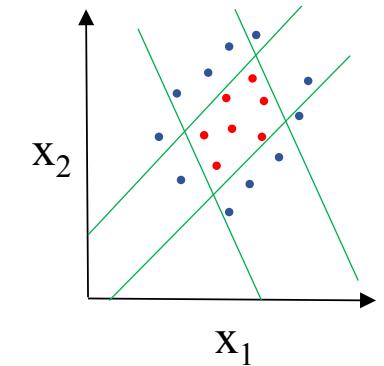
Data

Multilayered Perceptrons (1986)



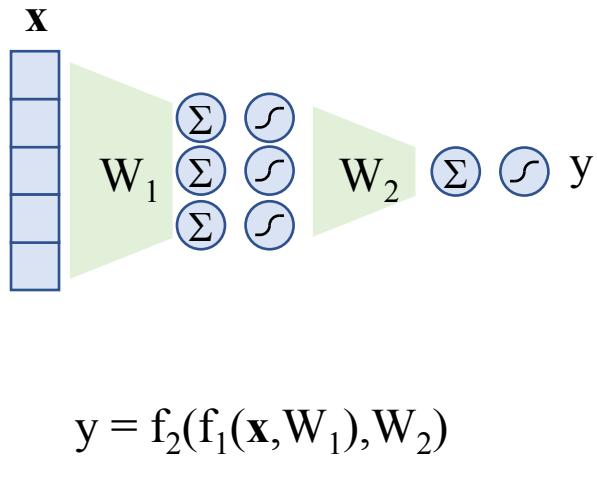
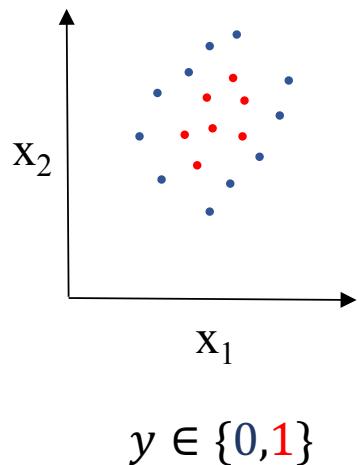
$$y \in \{0, 1\}$$

Data

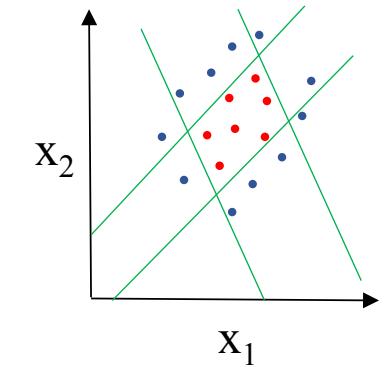


Model

Multilayered Perceptrons (1986)

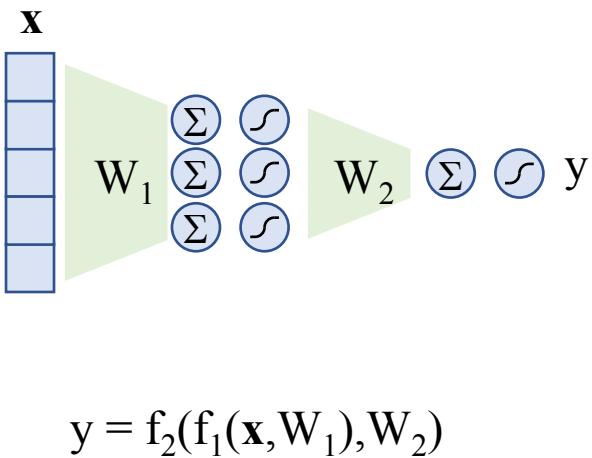
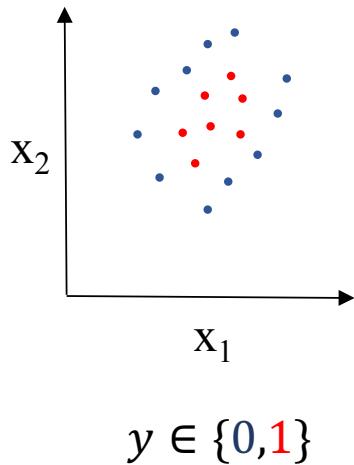


Data



Model

Multilayered Perceptrons (1986)

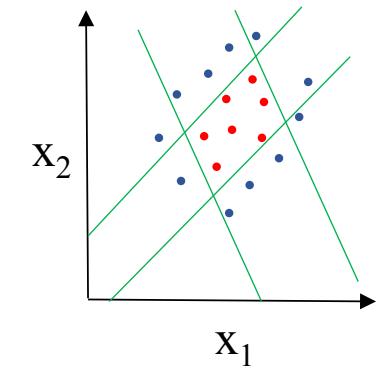


Data

Minimize
 $\|y - f_2(f_1(\mathbf{x}, W_1), W_2)\|$

No closed form solution for \mathbf{w}

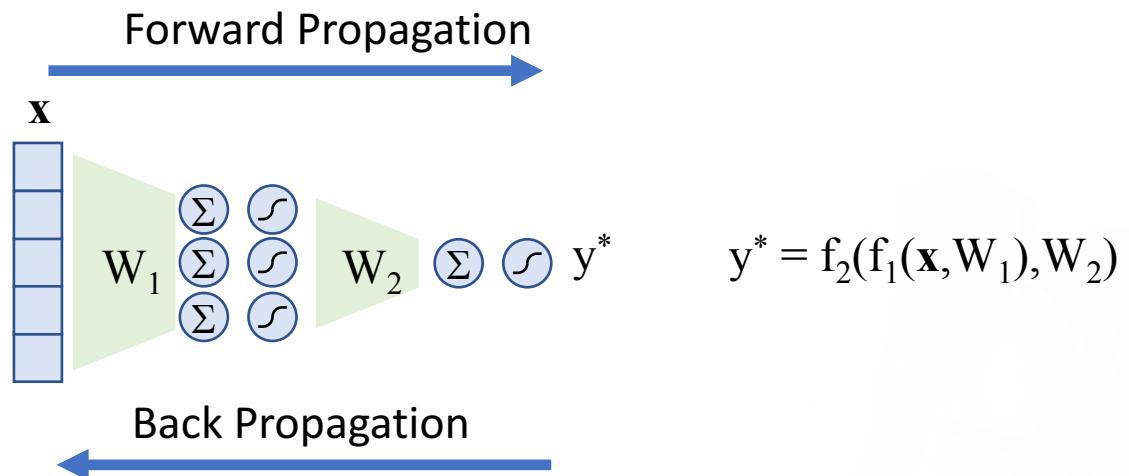
Minimize with gradient descent (back propagation)



Model

Back Propagation of Error

- Push error into nested functions
- Error metric
 - $E = \|y - f_2(f_1(x, W_1), W_2)\|$
- Continuous & differentiable
- Gradient descent on error with respect to weights



Back Propagation of Error

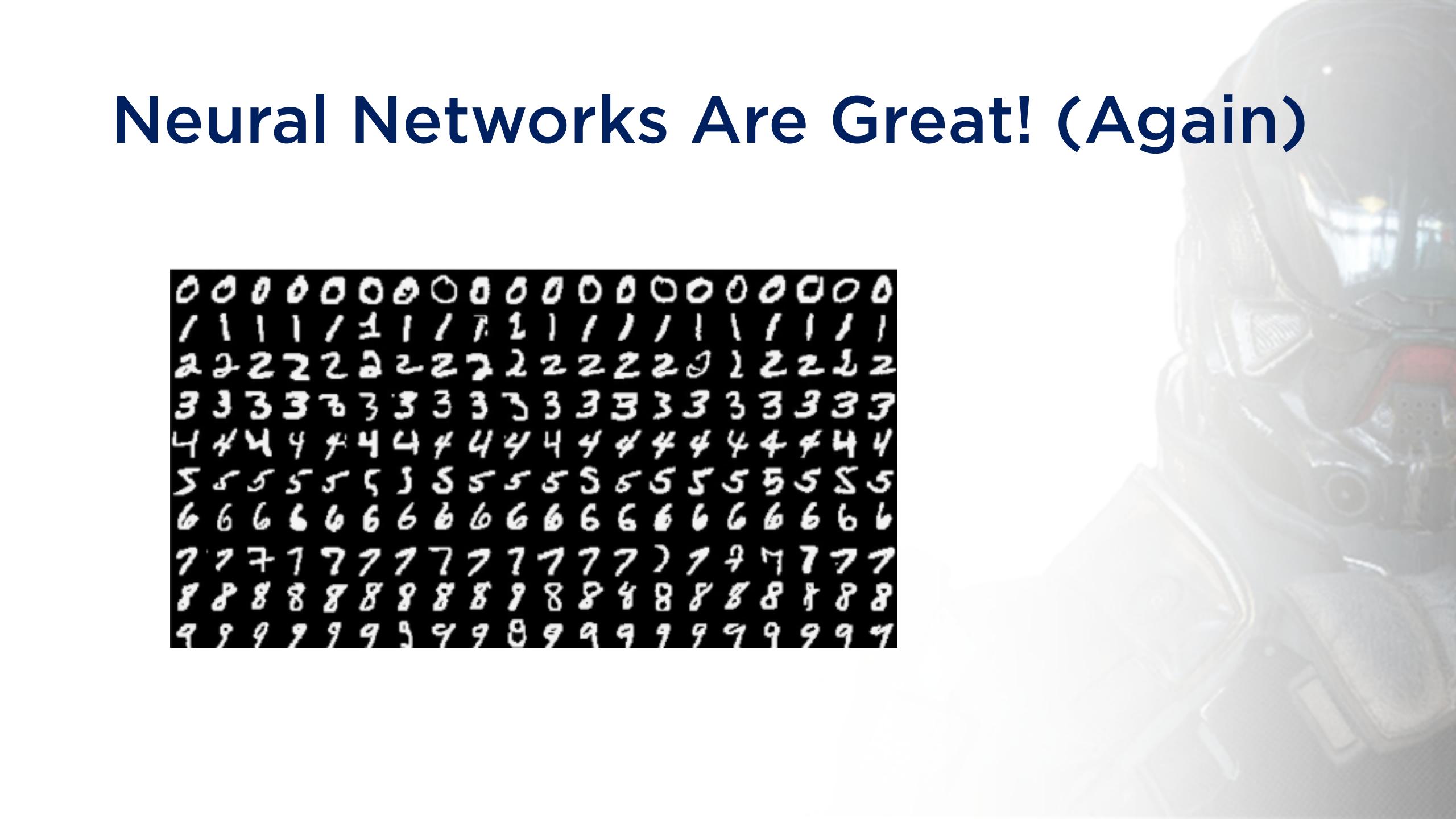
$$\sigma(W_1 | x) = f_1$$

$$\sigma(W_2 | f_1) = y^*$$

$$x \quad \Delta W_1 \quad f'_1 \quad e'_1$$

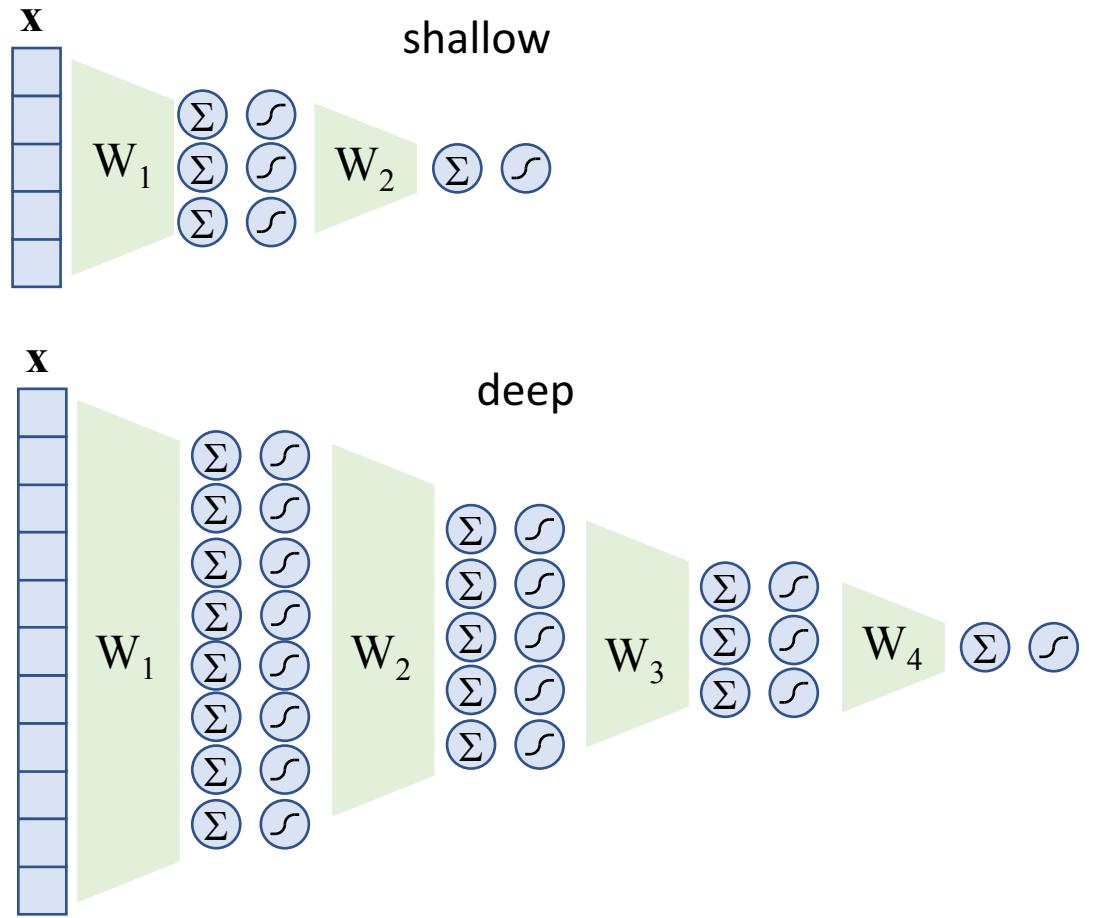
$$e^1 = W_2^T | f'_2 | y - y^*$$

Neural Networks Are Great! (Again)

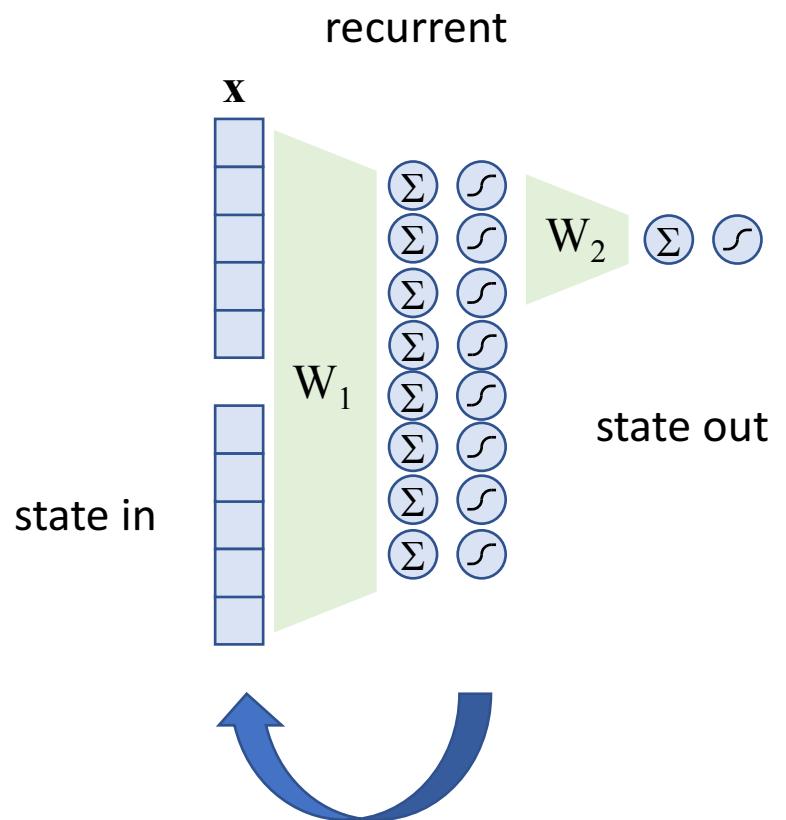


00000000000000000000
11111111111111111111
22222222222222222222
33333333333333333333
44444444444444444444
55555555555555555555
66666666666666666666
77777777777777777777
88888888888888888888
99999999999999999999

Deep Networks



Recurrent Networks



Three Problems

- Vanishing gradients in deep and recurrent networks
- High dimensional parameter spaces
- Long training times

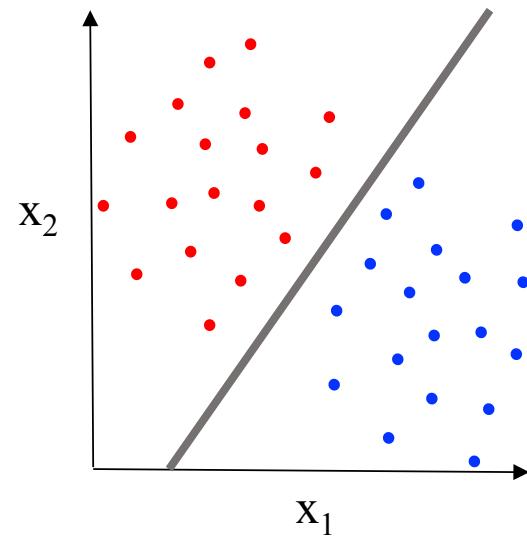
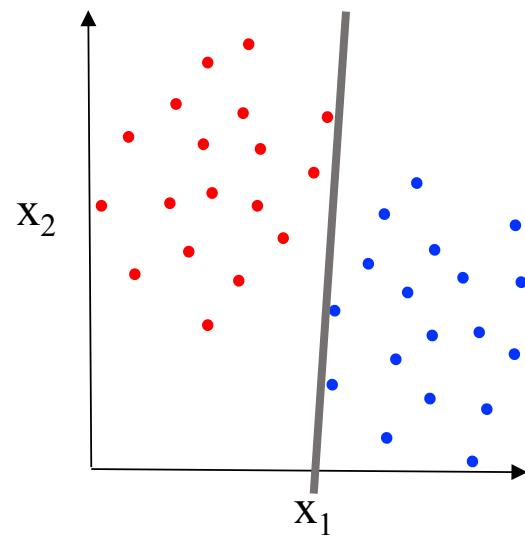


Support Vector Machines (1992)

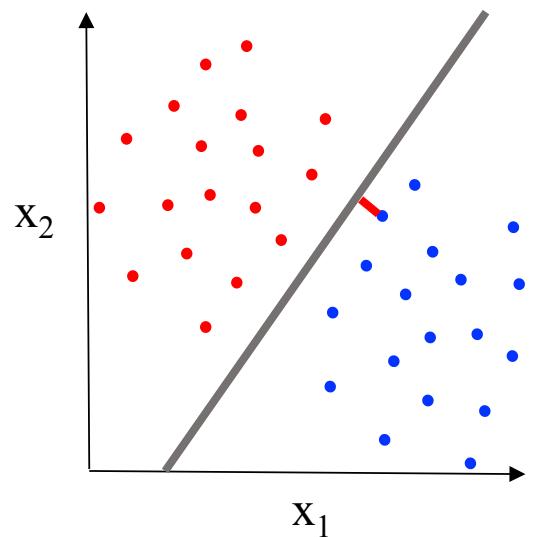
- Margin optimization (1963)
- Kernel Trick



Support Vector Machines (1992)



Support Vector Machines (1992)



Maximize Margin

Learning SVM
Write down constraints
Solve Quadratic Program

Support Vector Machines (1992)

- Learns linear mapping
- Non-linear mappings handled by “Kernel Trick”
 - Parabolic boundary
 - $[x,y] \rightarrow [x,y,x^2,y^2]$



Neural Networks on
Spark

Linear Regression

```
val lr = new LinearRegression().  
  setMaxIter(10).  
  setRegParam(0.3)
```

// Fit the model

```
val lrModel = lr.fit(training)
```

// Print the coefficients and intercept for logistic regression

```
val c = lrModel.coefficients  
val i = lrModel.intercept  
println(s"Coefficients: ${c} Intercept: ${i}")
```

Logistic Regression

```
val lr = new LogisticRegression().  
  setMaxIter(10).  
  setRegParam(0.3)
```

// Fit the model
val lrModel = lr.fit(training)

*// Print the coefficients and intercept for logistic
regression*

Support Vector Machine

```
val lsvc = new LinearSVC().
  setMaxIter(10).
  setRegParam(0.1)

// Fit the model
val lsvcModel = lsvc.fit(training)

// Print the coefficients and intercept for linear svc
val c = lsvcModel.coefficients
val i = lsvcModel.intercept
println(s"Coefficients: ${c} Intercept: ${i}")
```

Multilayered Perceptron

```
// specify layers for the neural network:  
// input layer of size 4 (features), two intermediate of size 5 and 4  
// and output of size 3 (classes)  
val layers = Array[Int](4, 5, 4, 3)  
  
// create the trainer and set its parameters  
val trainer = new MultilayerPerceptronClassifier().  
  setLayers(layers).  
  setBlockSize(128).  
  setSeed(1234L).  
  setMaxIter(100)  
  
// train the model  
val model = trainer.fit(train)  
  
// compute accuracy on the test set  
val result = model.transform(test)  
val predictionAndLabels = result.select("prediction", "label")  
val evaluator = new MulticlassClassificationEvaluator().  
  setMetricName("accuracy")  
  
println("Test set accuracy = " +  
  evaluator.evaluate(predictionAndLabels))
```