

# Neural Networks

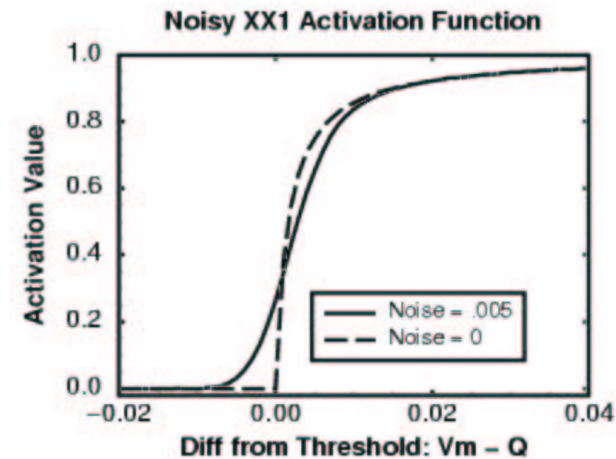
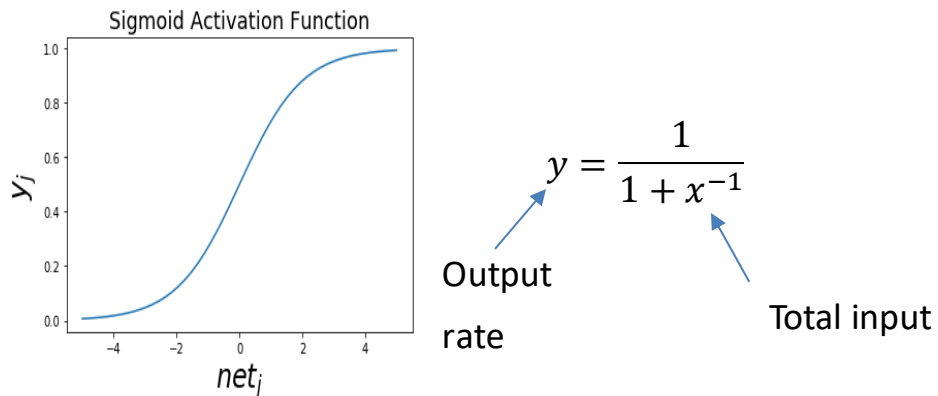
## *Theoretical and Mathematical Foundations*

---

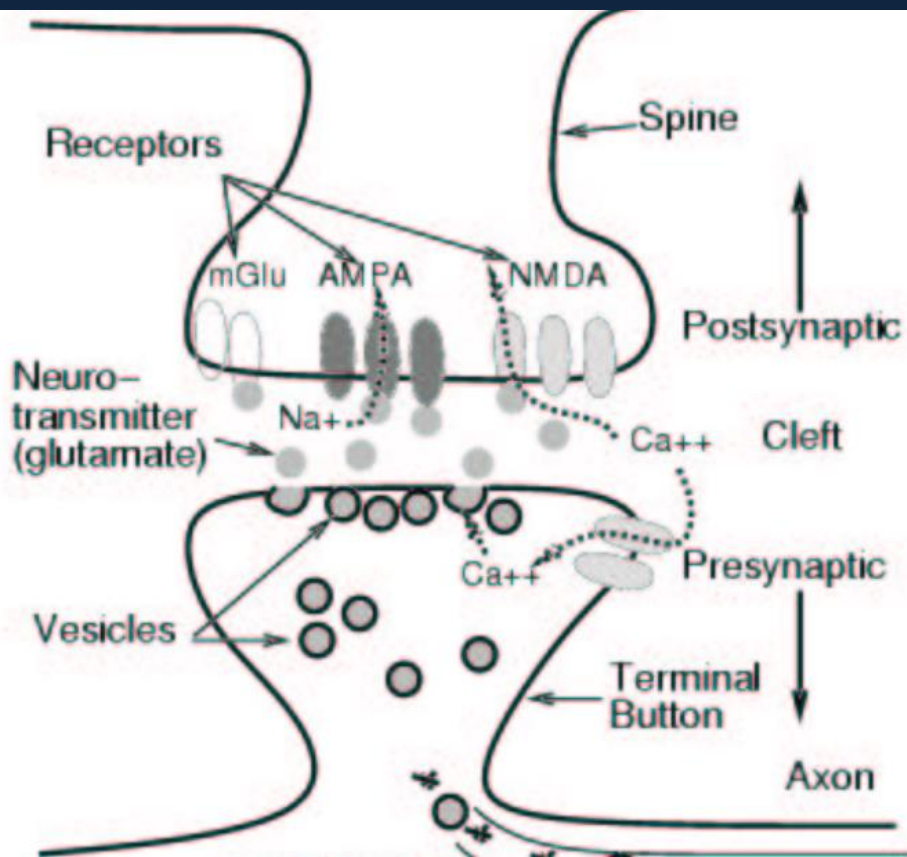
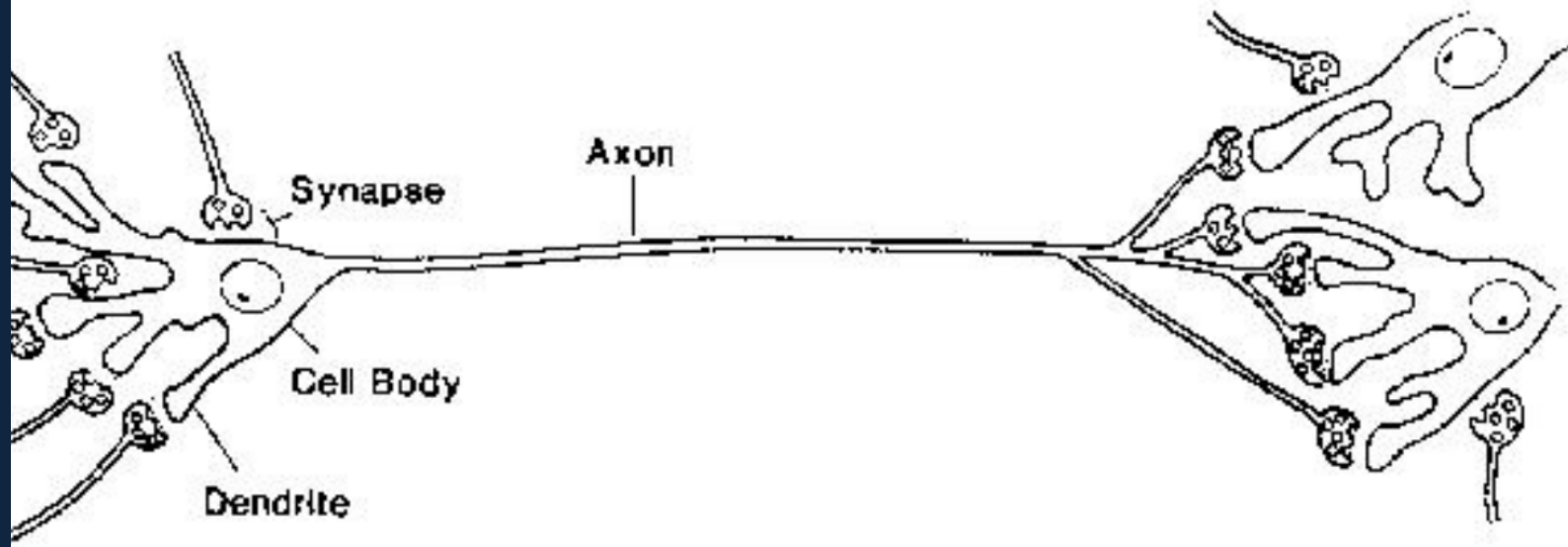
CSCI 4850/5850

# Rate Coding Activation Functions

- Even though action potentials are discrete events, they are often summarized together as a firing rate
- Zero means no APs and one means firing at the maximum possible rate
- A little noisy as well...



# Memory

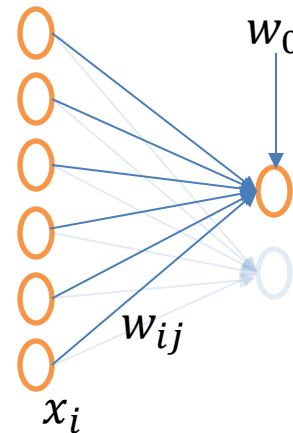


- Information or memory is stored in the connections between neurons (synapses)
- Summarized as the “weight” of the connection (large/small connections)

# Abstracting Details: Net Input

- The *net input* for unit  $j$  can be “folded” into a single value,  $net_j$ , which weights the activation of the sending neuron  $x_i$  by the strength of the synaptic connection between  $i$  and  $j$ ,  $w_{ij}$ .
- The sum of all weighted activations into a unit plus a *bias weight*,  $w_0$ , is the total net input.

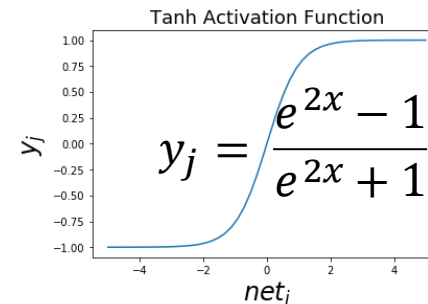
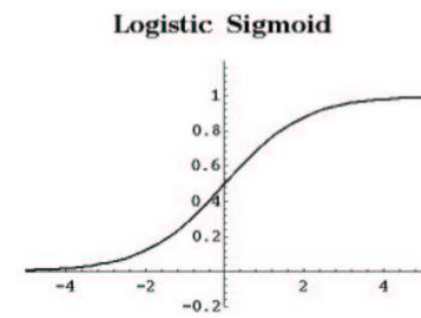
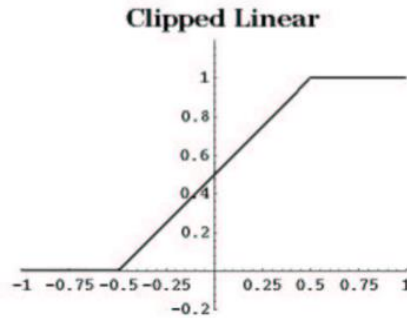
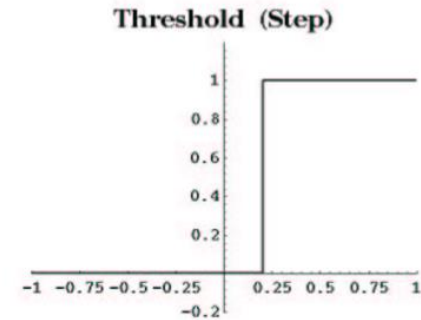
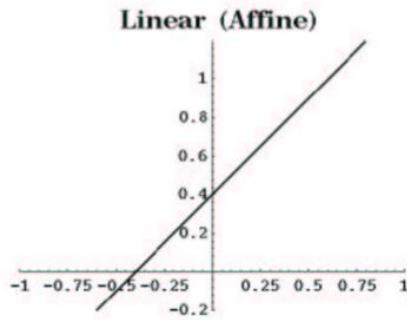
$$net_j = w_0 + \sum_{i=1}^n x_i w_{ij}$$



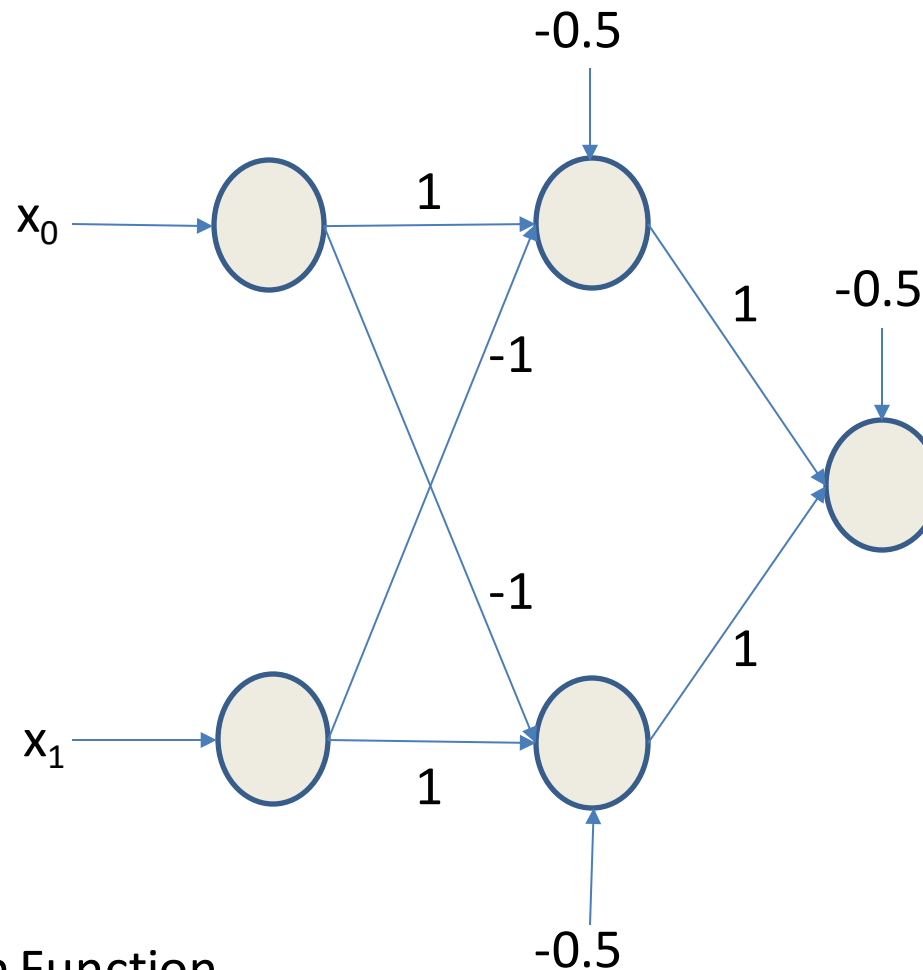
# Abstracting Details: Activation Functions

- The *net input* can then be transformed into a *rate code* using an *activation function*
- A commonly used function is the logistic sigmoid:

$$f(\text{net}_j) = \frac{1}{1 + e^{-\text{net}_j}}$$



# Simple Hardwired ANN



Vectors

$X \rightarrow F(X)$

0 0  $\rightarrow$  ?

0 1  $\rightarrow$  ?

1 0  $\rightarrow$  ?

1 1  $\rightarrow$  ?

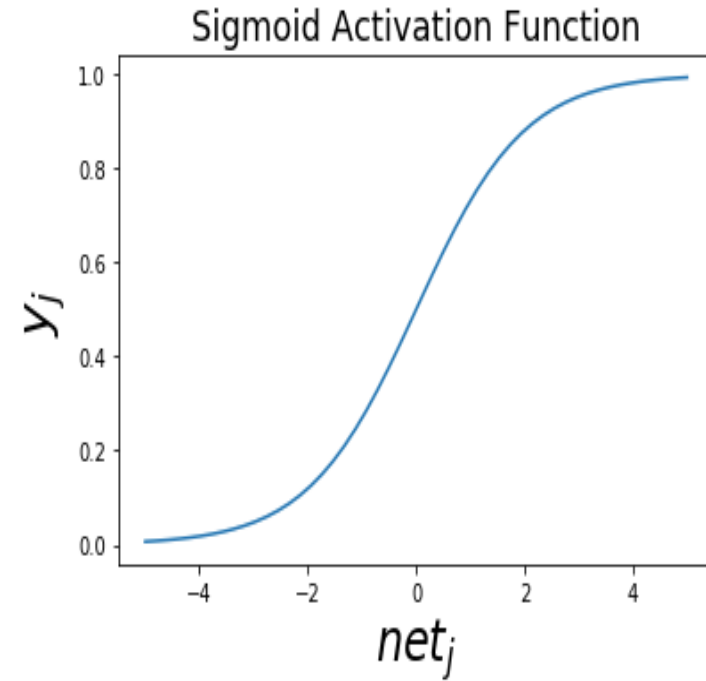
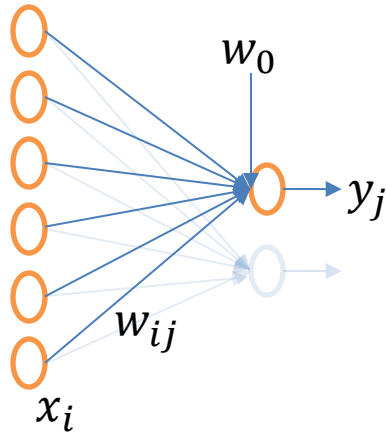
Threshold Activation Function

$f(\text{net}) = 1$  if  $\text{net} > 0$

$f(\text{net}) = 0$  if  $\text{net} \leq 0$

# Together: Simple Neural Units

$$net_j = w_0 + \sum_{i=1}^n x_i w_{ij}$$

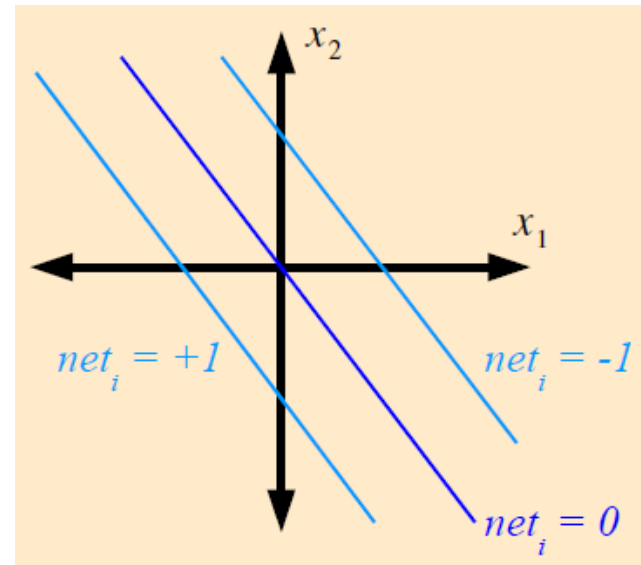


$$y_j = \frac{1}{1 + e^{-net_j}}$$

# What does Net Input mean?

- The change in postsynaptic membrane potential is reflected in the *net input*
- Net input is an **affine function** of the presynaptic cell activity (input)
  - In the field of machine learning, we often refer this as a *linear* relationship even though it's technically affine (due to the bias weight,  $w_0$ )
- Geometrically, there is a **region** where the net input is *positive* and a **region** where the net input is *negative* which is separated by the n-dimensional **hyperplane** where the net input is zero...

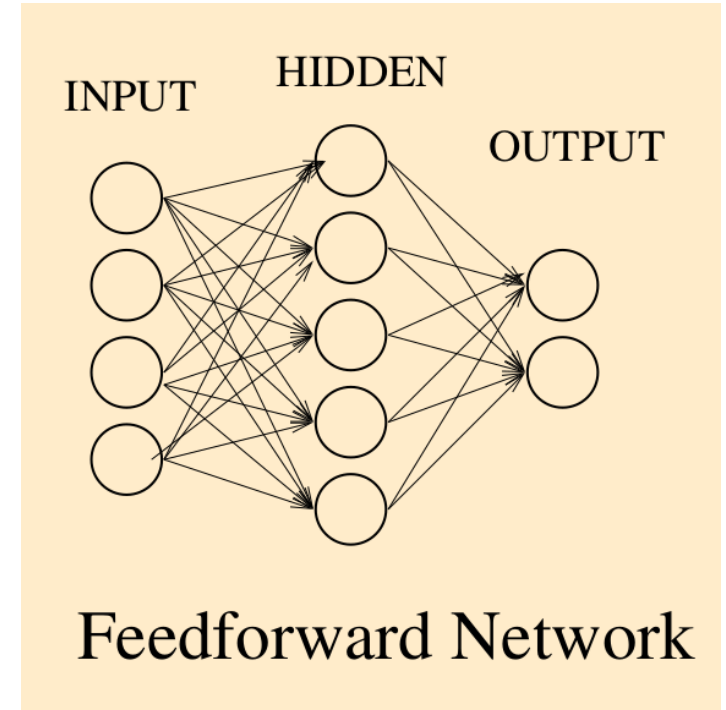
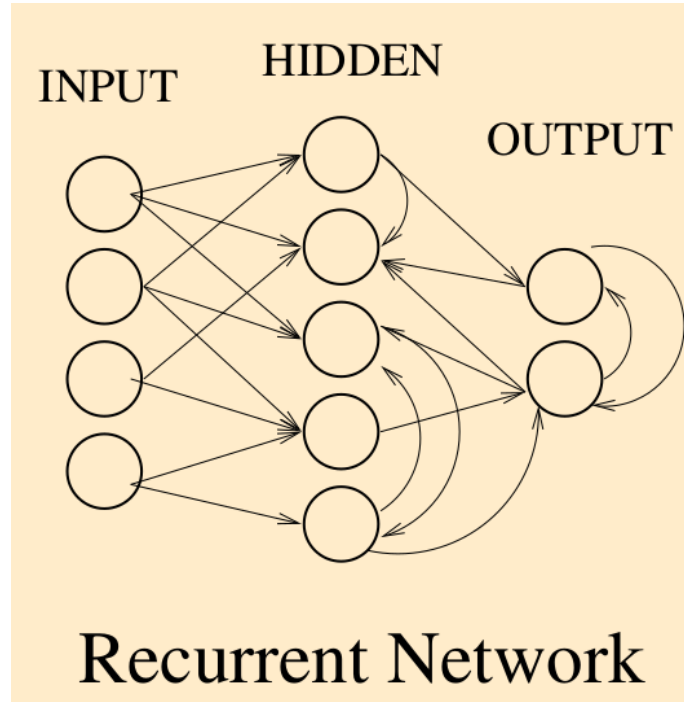
$$net_j = w_0 + \sum_{i=1}^n x_i w_{ij}$$



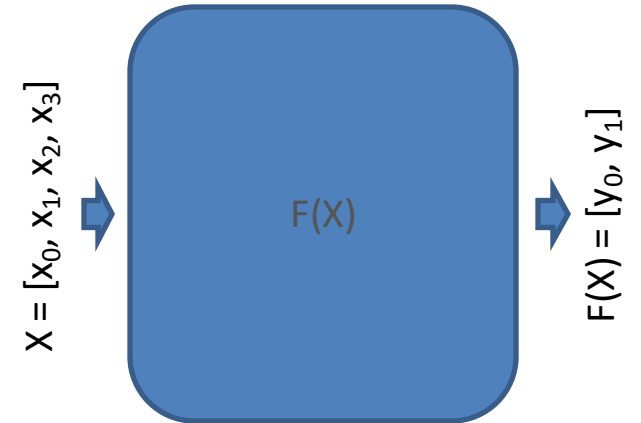
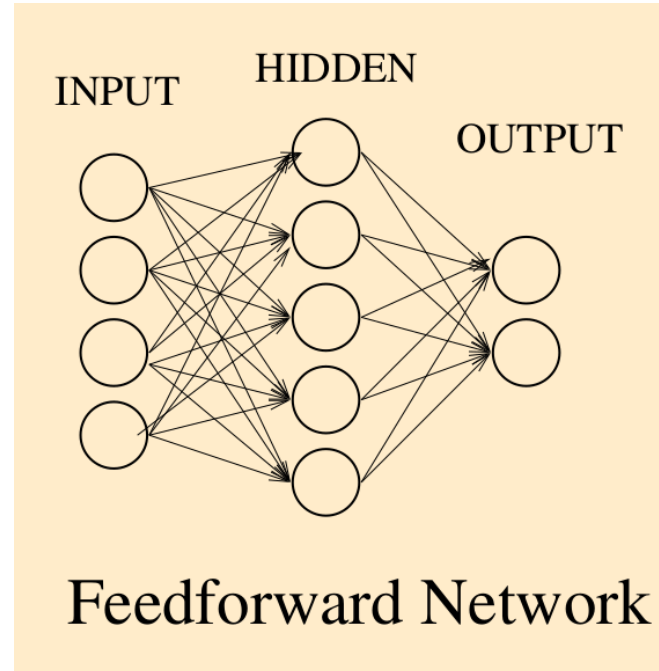


# Abstracting Details: Architecture

- ANNs typically consist of input, output, and hidden layers of artificial neural units
- May be “single layer” or “multi-layer”, but the pattern of connections is called the ANNs *architecture*



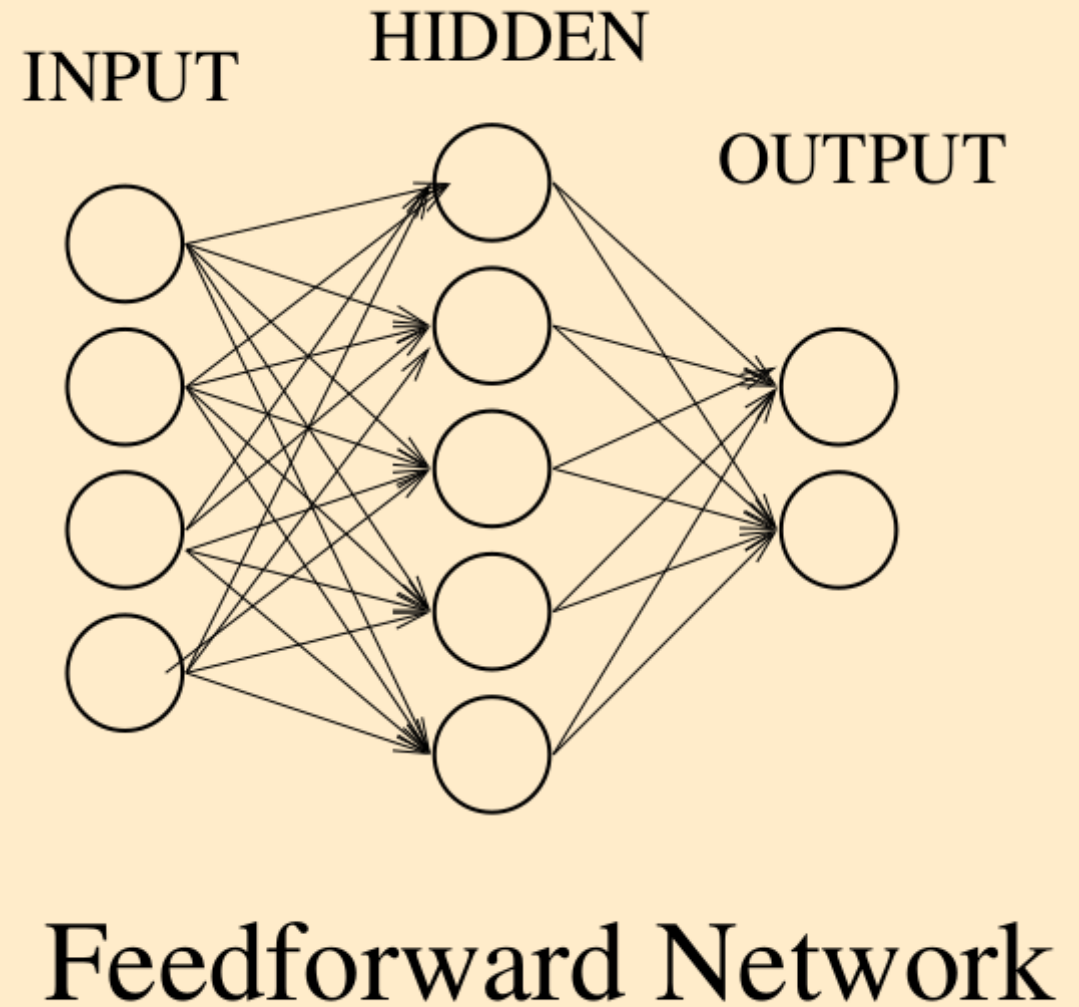
# Standard Artificial Neural Network: Feed-forward Networks



# Universal Approximation

---

- It has been proven that given *enough* units and using only a single layer of hidden units, essentially *any* function can be approximated to an arbitrary degree of precision
- Hornik, Stinchcombe, and White, 1989



# Knowledge Retention/Representation

---

- Neural networks typically utilize 2 kinds of *memory* for storing or processing information:
  - **Weight-based:** long-term, non-volatile stored in the synaptic connections between units
  - **Activation-based:** short-term, volatile stored in the activation values of the neural units



# Learning



Notice that most of the information encoded by a feed-forward network is in the connection (and bias) weights

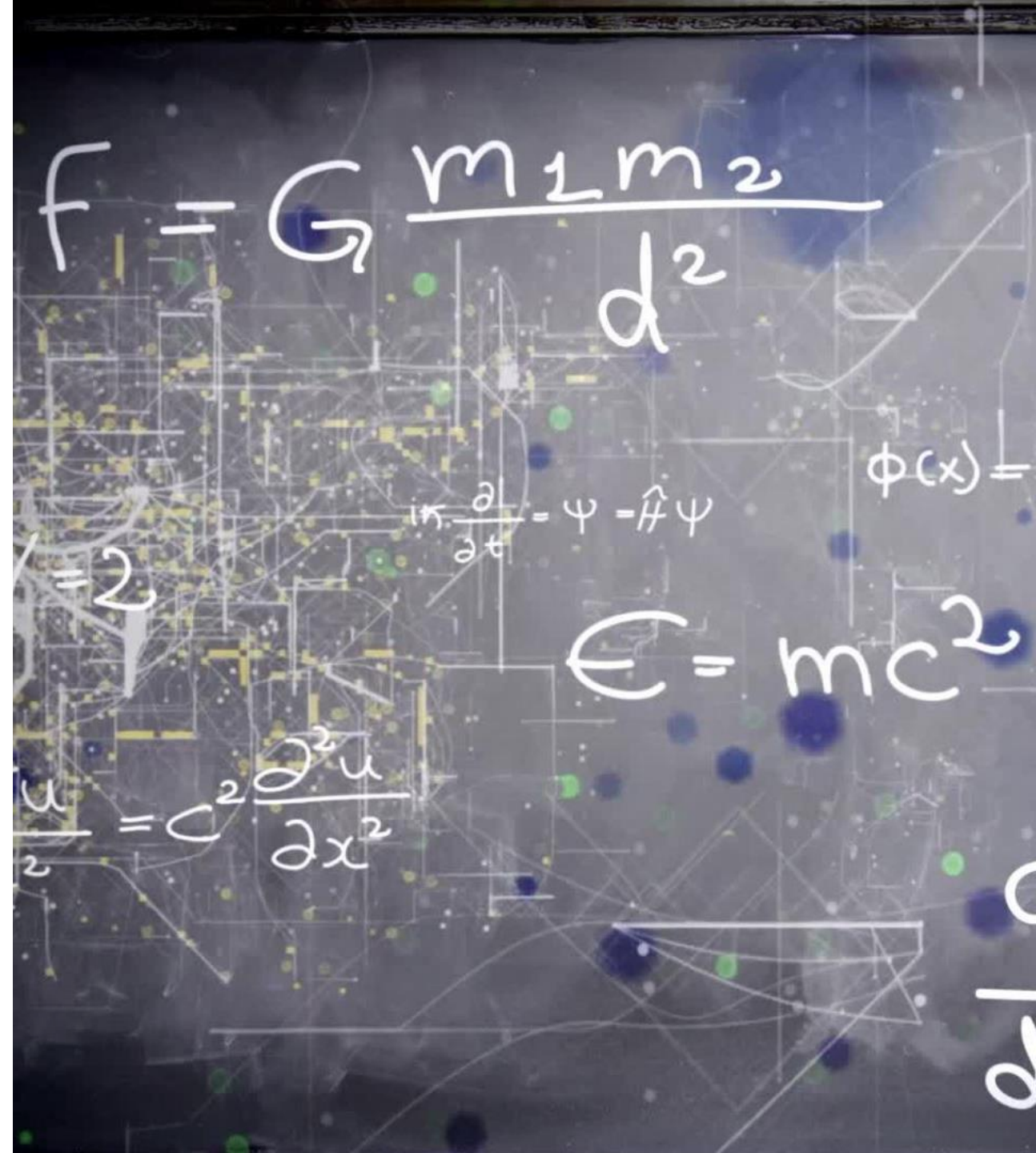


So, the question becomes, how does one calculate the weights?

Hand-wired?  
Trial-and-error?  
Random generation?  
Hmmm...

# Foundation: Calculus

- Learning often involves making **incremental changes** to a system, and **calculus is the mathematics of change**.
- Activation levels in **recurrent** neural architectures are naturally **dynamic**, and calculus is the foundational mathematics of **dynamical systems**.



# The Derivative

- A function is ***differentiable*** if its first partial derivative exists everywhere in its domain for each of its arguments. A function is ***smooth*** if all partial derivatives exist.

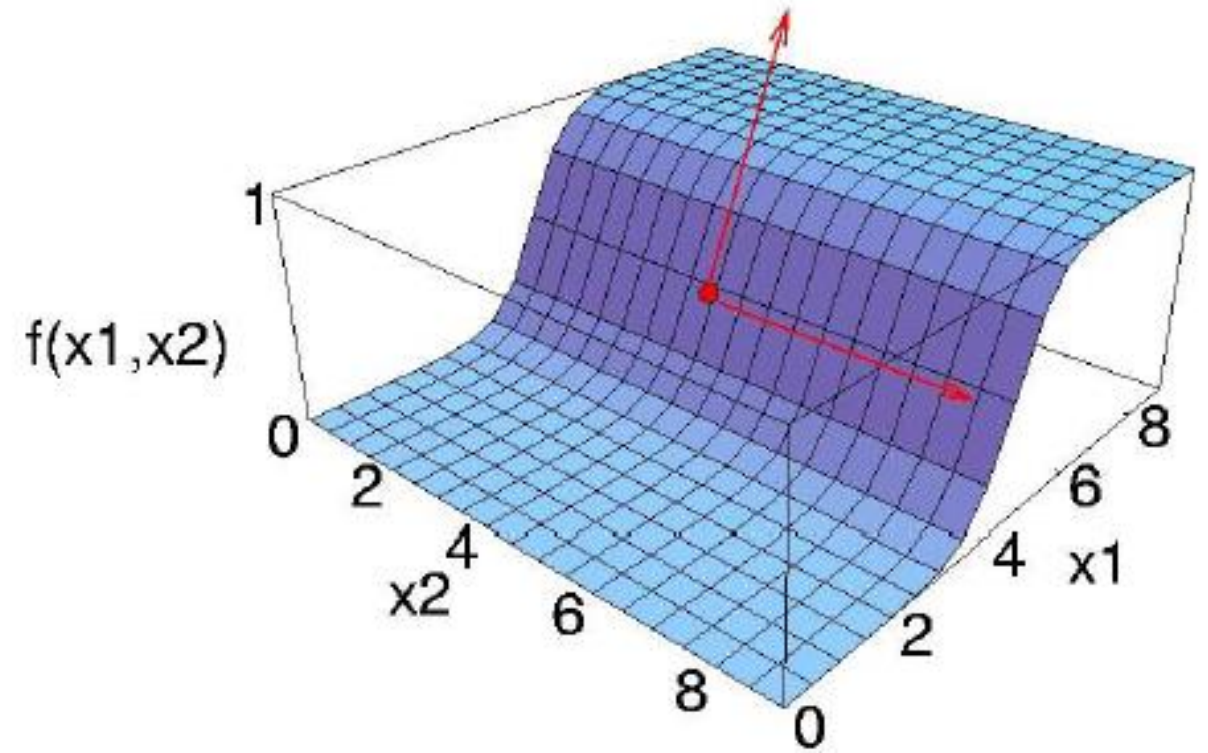
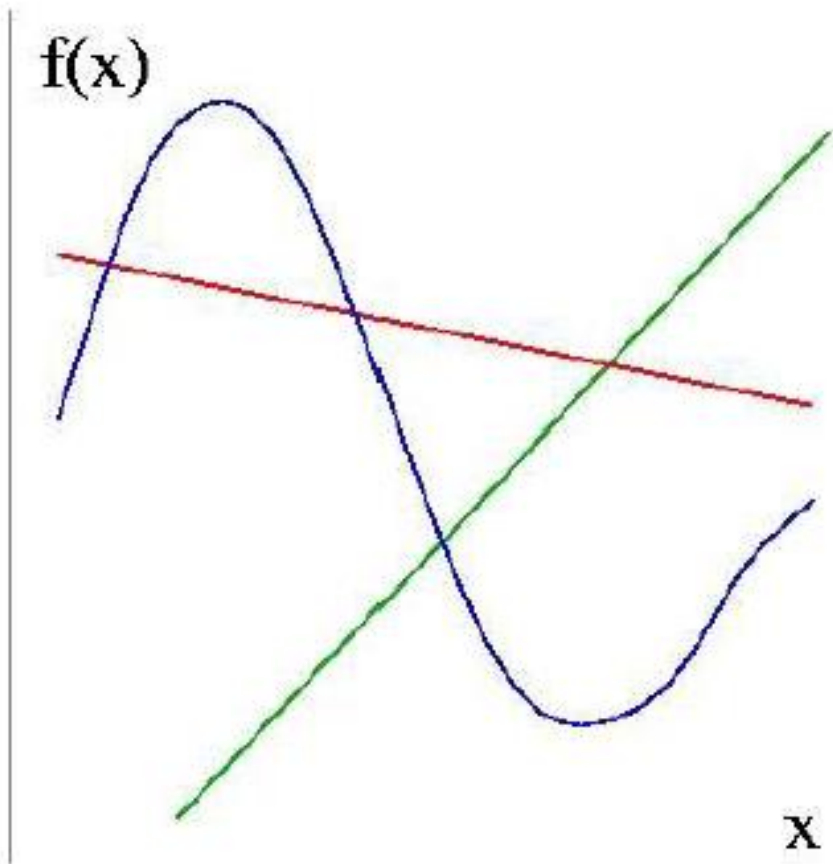
$$f'(x) = \frac{df(x)}{dx} = \lim_{\epsilon \rightarrow 0} \frac{f(x+\epsilon) - f(x)}{\epsilon}$$

$$f''(x) = \frac{d}{dx} \left( \frac{df(x)}{dx} \right)$$

$$\frac{\partial f(x_1, \dots, x_i, \dots, x_n)}{\partial x_i} = \lim_{\epsilon \rightarrow 0} \frac{f(x_1, \dots, x_i + \epsilon, \dots, x_n) - f(x_1, \dots, x_i, \dots, x_n)}{\epsilon}$$



# Geometric Interpretation of Derivatives





# The Chain Rule

$$\frac{d f(g(x))}{d x} = \frac{d f(u)}{d u} \frac{d g(x)}{d x}$$

Also works for partial derivatives...

$$f(x) = f(s(x), t(x))$$

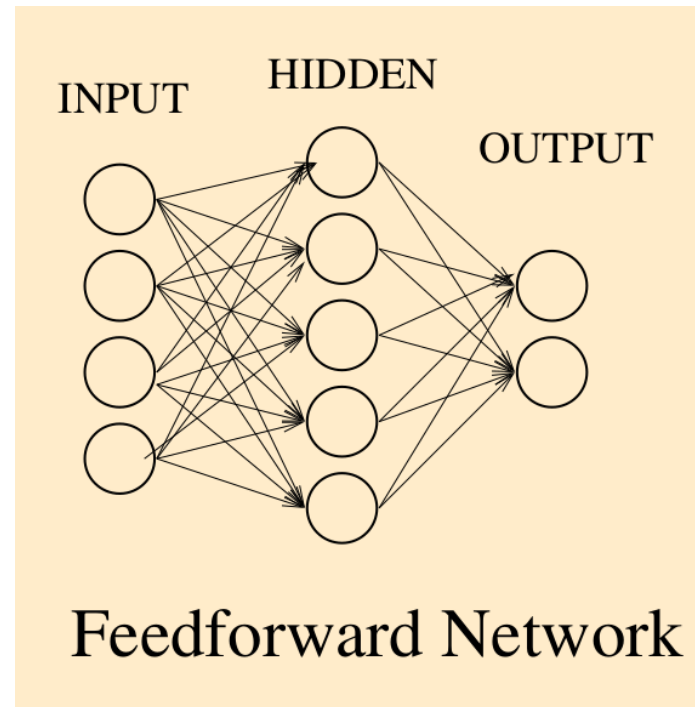
$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial s} \frac{\partial s}{\partial x} + \frac{\partial f}{\partial t} \frac{\partial t}{\partial x}$$

$$f(x) = f(g_1(x), g_2(x), \dots, g_n(x))$$

$$\frac{\partial f}{\partial x} = \sum_{i=1}^n \frac{\partial f}{\partial g_i} \frac{\partial g_i}{\partial x}$$

# Linear Algebra

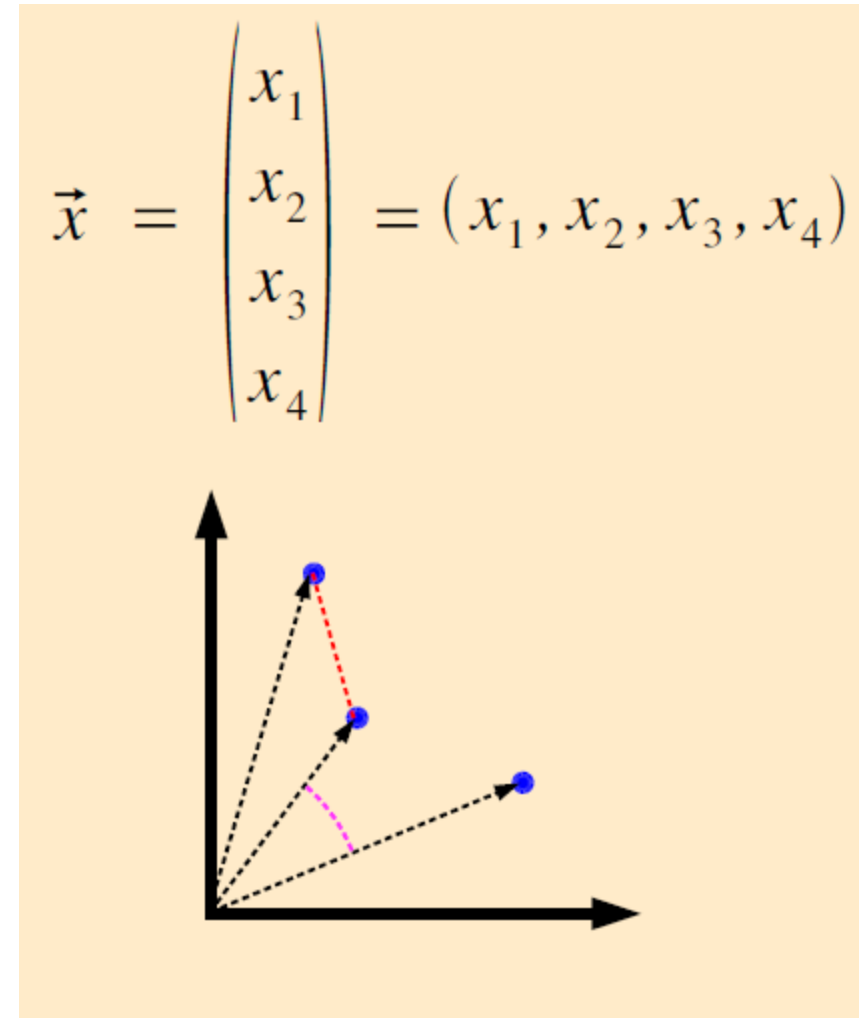
- Neural units are usually grouped into layers, so it's more natural to group variable values for all units in a layer into **vectors**.
- Connection weights to a single unit can be grouped into a **weight vector**, but when layers are completely interconnected, it makes more sense to group them into a **weight matrix**.



$$\mathbf{net} = [\mathit{net}_1, \mathit{net}_2, \mathit{net}_3, \mathit{net}_4, \mathit{net}_5]$$
$$\mathbf{act} = [f(\mathit{net}_1), f(\mathit{net}_2), f(\mathit{net}_3), f(\mathit{net}_4), f(\mathit{net}_5)]$$

# Linear Review

- A **vector** is similar to a 1-D array, but geometrically we think of the vector as a **point** in an  **$n$ -dimensional vector space**, where  $n$  is the length of the vector.
- Concepts:
  - Euclidean Distance
  - Cosine Similarity
  - Orthogonality
  - Orthonormality
  - Linear Independence



# Vector Arithmetic

- Multiply a vector by a scalar

$$ax = (ax_1, ax_2, \dots, ax_n)^T$$

- Adding two vectors

$$x + y = (x_1 + y_1, x_2 + y_2, \dots, x_n + y_n)^T$$

- Length of a vector (magnitude, norm)

$$|x| = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$$

# Vector Distances

- Euclidean distance between two vectors

$$|x - y| = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \cdots + (x_n - y_n)^2}$$

- Dot product (or inner product) of two vectors

$$x \cdot y = x_1 y_1 + x_2 y_2 + \cdots + x_n y_n$$

- Angular distance between two vectors

$$x \cdot y = |x||y| \cos \theta$$

- Cosine distance between two vectors

$$1 - \frac{x \cdot y}{|x||y|}$$

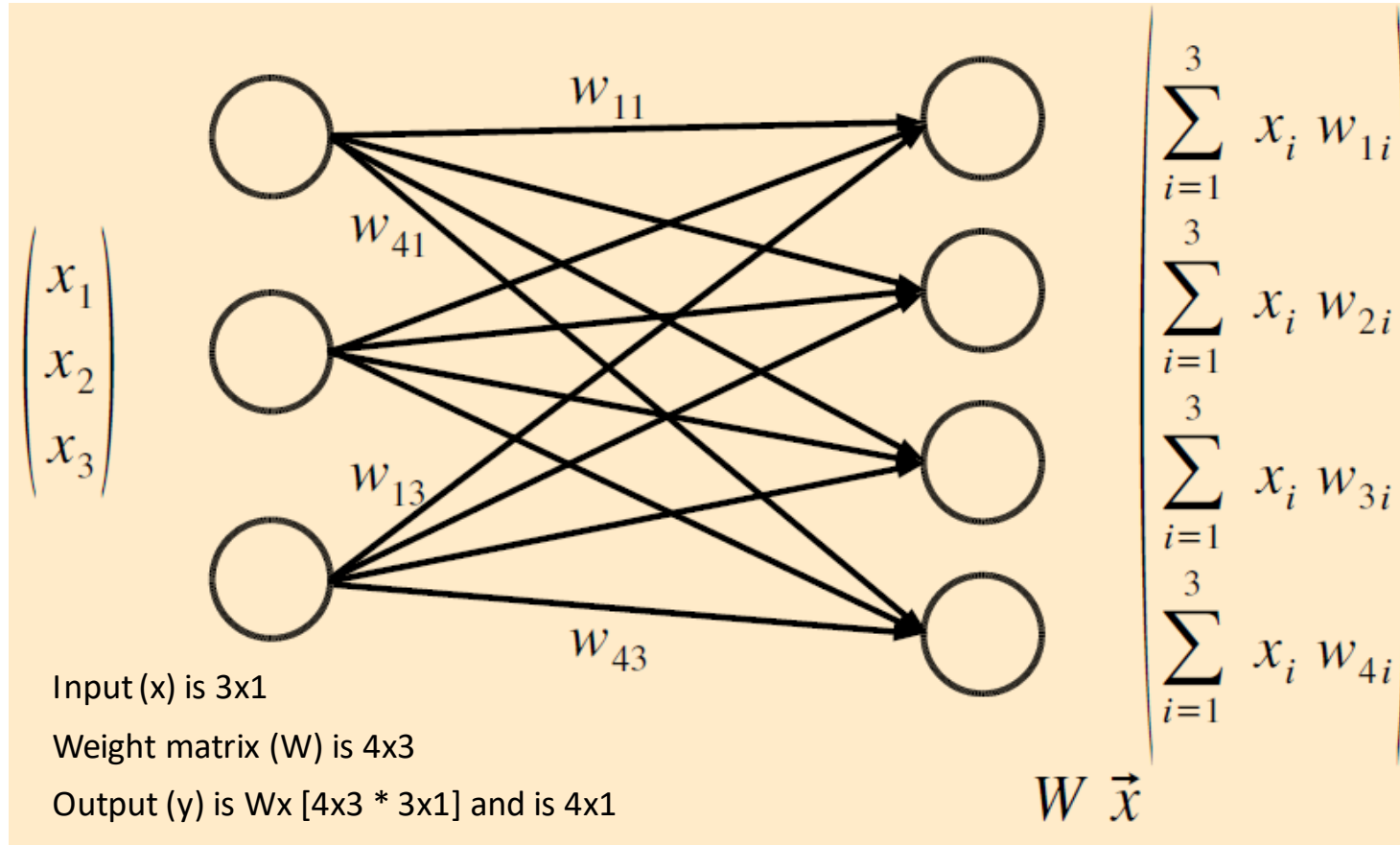
# Matrices

- A **matrix** is similar to a 2-D array
- Multiplying a *vector* by a *matrix* produces another **vector**...

$$W = \begin{pmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \\ w_{41} & w_{42} & w_{43} \end{pmatrix}$$

$$W \vec{x} = \begin{pmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \\ w_{41} & w_{42} & w_{43} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} x_1 w_{11} + x_2 w_{12} + x_3 w_{13} \\ x_1 w_{21} + x_2 w_{22} + x_3 w_{23} \\ x_1 w_{31} + x_2 w_{32} + x_3 w_{33} \\ x_1 w_{41} + x_2 w_{42} + x_3 w_{43} \end{pmatrix}$$

# Weight Matrices



Alternatively:  $x^T W^T = y^T$

# Matrix Multiplication

$$A \ B = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \\ a_{41} & a_{42} & a_{43} \end{pmatrix} \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \end{pmatrix}$$

Shapes:  
 $4 \times 3 * 3 \times 2 = 4 \times 2$

$$= \begin{pmatrix} (a_{11}b_{11} + a_{12}b_{21} + a_{13}b_{31}) & (a_{11}b_{12} + a_{12}b_{22} + a_{13}b_{32}) \\ (a_{21}b_{11} + a_{22}b_{21} + a_{23}b_{31}) & (a_{21}b_{12} + a_{22}b_{22} + a_{23}b_{32}) \\ (a_{31}b_{11} + a_{32}b_{21} + a_{33}b_{31}) & (a_{31}b_{12} + a_{32}b_{22} + a_{33}b_{32}) \\ (a_{41}b_{11} + a_{42}b_{21} + a_{43}b_{31}) & (a_{41}b_{12} + a_{42}b_{22} + a_{43}b_{32}) \end{pmatrix}$$



# Special Matrices

- The **identity matrix** (**I**) is zero everywhere except along its main diagonal, where it is exactly 1. This means that for a square matrix, A:

$$AI = A$$

- Some square matrices have an **inverse**, such that:

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

- Every matrix has a **transpose**, which swaps the rows and columns of the original matrix:

$$A^T$$

- A **symmetric** matrix is any matrix which is equal to its own transpose:

$$A = A^T$$

# The Gradient and The Hessian

- For a function of  $n$  variables, the **gradient** is a *vector* of the form:

$$\frac{d f(\vec{x})}{d \vec{x}} = \nabla_{\vec{x}} f(\vec{x}) = \left( \frac{\delta f}{\delta x_1}, \frac{\delta f}{\delta x_2}, \dots, \frac{\delta f}{\delta x_n} \right)^T$$

- For a function of  $n$  variables, the **Hessian** is *matrix* with each element (i,j) is:

$$\frac{\partial}{\partial x_i} \left( \frac{\partial f(\vec{x})}{\partial x_j} \right)$$

# Matrix Factorization (aka Decomposition)

- Eigen decomposition (square)
  - The nonzero vector,  $\mathbf{x}$ , is an eigenvector of the matrix  $\mathbf{A}$ , with corresponding eigenvalue,  $\lambda$ , if and only if:

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$$

- Alternatively, Eigen decomposition factors the square matrix  $\mathbf{A}$  as follows:

$$\mathbf{A} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$$

- $\mathbf{Q}$  is the square matrix of eigen vectors (one in each column), and  $\mathbf{\Lambda}$  is the diagonal matrix of eigen values

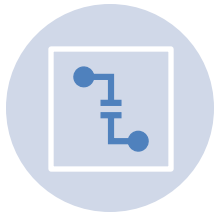
- Singular value decomposition (rectangular)

- The singular value decomposition of a matrix is:

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$$

- Where  $\mathbf{\Sigma}$  is a rectangular diagonal matrix of singular values, and  $\mathbf{U}$  and  $\mathbf{V}$  are the left-singular and right-singular vectors of  $\mathbf{A}$ , respectively.

# Probability Theory



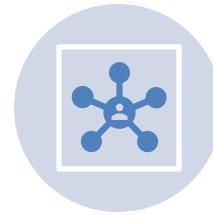
There is often *uncertainty* in the reliability of network inputs



There is often *uncertainty* in what to do with a given set of inputs



There is often *uncertainty* about our neural network performance



We often construct networks using noisy components, introducing *uncertainty* in behavior



Probability, information, and statistical theory provide mathematical formalisms for handling *uncertainty*

# Probability

- Probabilities are defined over a space of **events**.
  - Each **atomic event** has a **probability** assigned to it in the range  $[0,1]$
  - The sum of **all** atomic event probabilities is 1
  - Events are discrete (happen / don't happen)
  - Mutually exclusive and exhaustive
- A **discrete random variable** takes on one of a finite set of values depending on the relative probability of those values
- A **continuous random variable** takes on a real value in some (potentially infinite) range, events are defined by infinitesimal subranges.

# Concepts in Probability Theory

- Probability Distribution
- Joint Probability Distribution
- Unconditional Probability
- Conditional Probability
- Bayes' Rule
- Prior Probability
- Posterior Probability
- Likelihood

$$P(X)$$

$$P(X_1, X_2, \dots, X_n)$$

$$P(X = x_i) \text{ or } P(x_i)$$

$$P(a \mid b) = P(a \wedge b) / P(b)$$

$$P(a \mid b) = \frac{P(b \mid a)P(a)}{P(b)}$$

$$P(\textit{hypothesis})$$

$$P(\textit{hypothesis} \mid \textit{evidence})$$

$$P(\textit{evidence} \mid \textit{hypothesis})$$

# Common Statistics

- **Expected value** of a discrete random variable (commonly known as the mean or average):

$$E[x] = \mu_x = \sum_{i=1}^n x_i P(x = x_i)$$

- **Variance** of such a variable:

$$\text{var}[x] = \sigma_x^2 = \sum_{i=1}^n (x_i - E[x])^2 P(x = x_i)$$

- **Standard deviation** is the square-root of the variance

$$\sigma_x = \sqrt{\text{var}[x]} = \sqrt{\sigma_x^2}$$

# Covariance

- Elements within a vector may influence one another. A measure of the relationship between elements is the covariance:

$$\text{cov}[x_i, x_j] = \sigma_{ij} = E[(x_i - \mu_i)(x_j - \mu_j)]$$

- The covariance matrix is the square matrix whose (i,j) element is  $\text{cov}[x_i, x_j]$ .