

# CS6421: Deep Neural Networks

---

**Gregory Provan**

Spring 2020

Lecture 4: Mathematical Foundations

Based on notes from Marc Diesenroth

# Overview

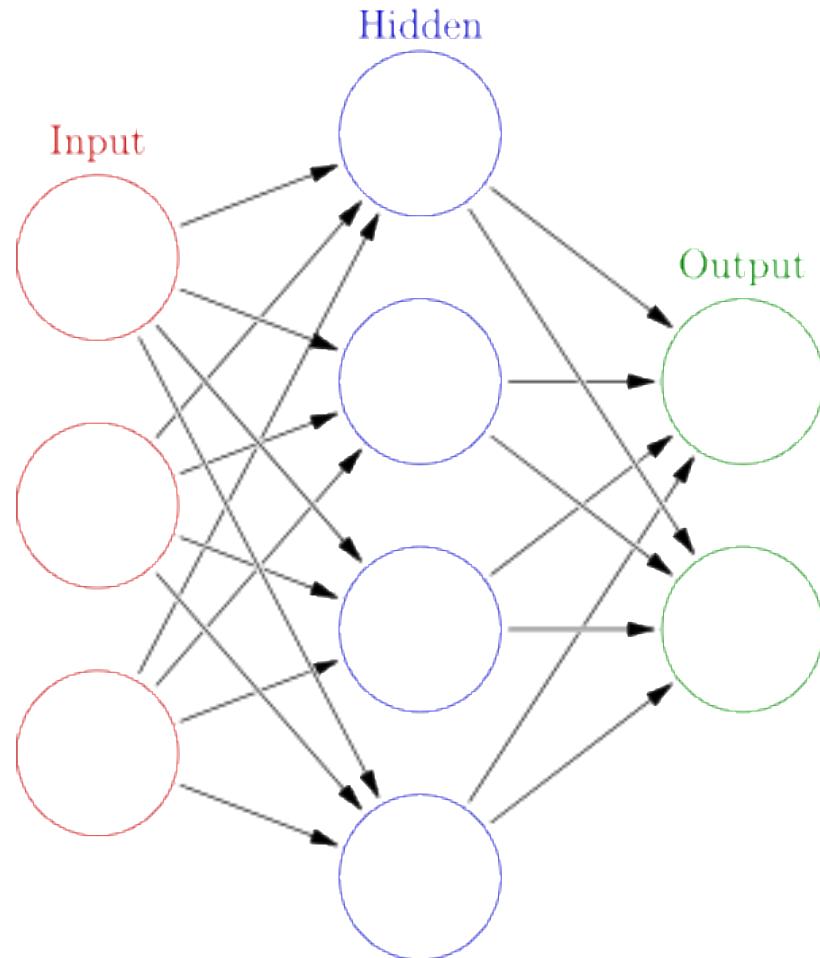
- Mathematical Basis of Deep Learning
- Linear Algebra
- Multivariate Calculus

# Key Mathematical Concepts

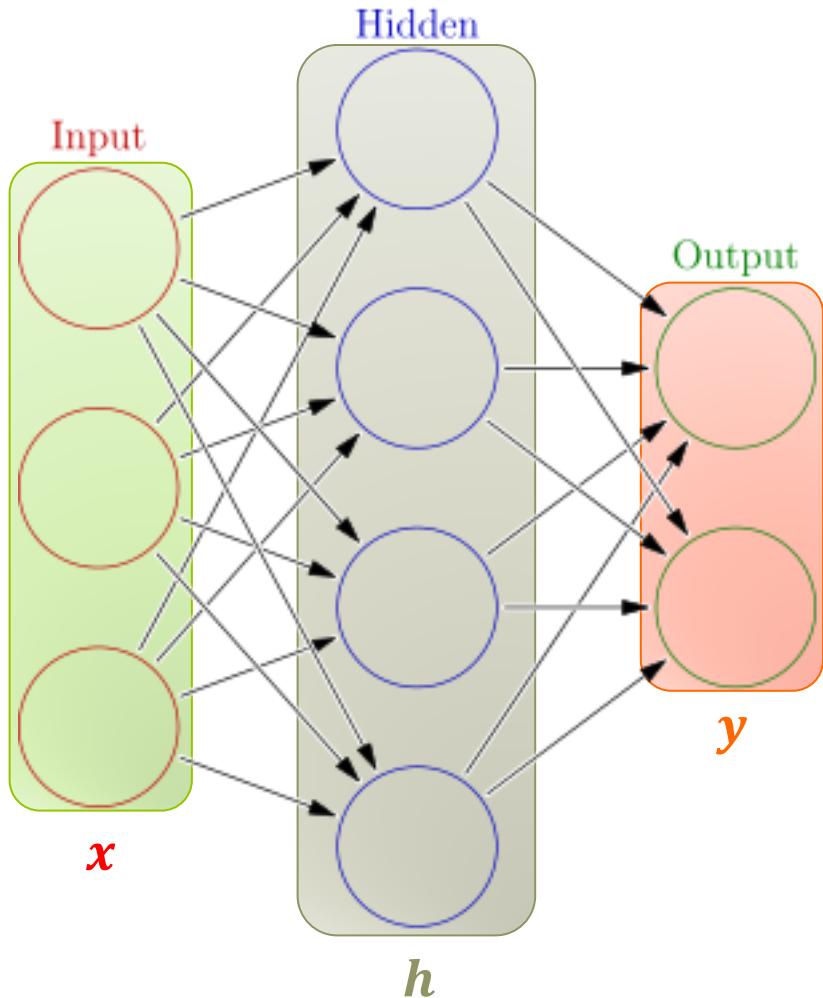
- Linear algebra and matrix decomposition
- **Differentiation**
- **Integration**
- Optimization
- Probability theory and Bayesian inference
- Functional analysis

# Neural Networks Basics

- Operations
  - Forward prop: makes predictions
  - Backward prop: adjusts parameters



# Neural Network: Definition



Weights

$$h = \sigma(W_1 x + b_1)$$
$$y = \sigma(W_2 h + b_2)$$

Activation functions

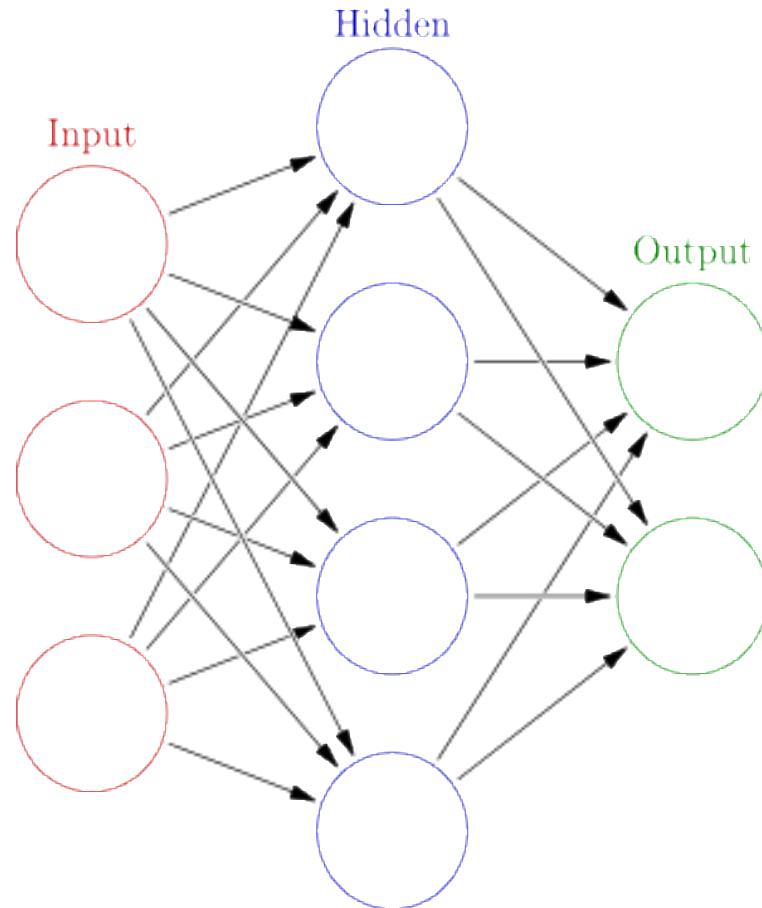
What are the underlying mathematical operations?

# Forward Propagation

$$\vec{X} \rightarrow \vec{H} \rightarrow \hat{\vec{Y}}$$

$$\vec{H} = \text{relu} \left( \left( [W] * \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} \right) + [b] \right)$$

$$\hat{\vec{Y}} = \text{relu} \left( \left( [U] * \begin{bmatrix} H_1 \\ H_2 \\ H_3 \\ H_4 \end{bmatrix} \right) + [c] \right)$$



# Forward Propagation

$$\vec{H} = \text{relu} \left( \left( [W] * \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} \right) + [b] \right)$$

## Matrix Multiplication

- Changes Dimensionality

$$(m * n) * (n * p) = (m * p)$$

$$(m * n) * (3 * 1) = (4 * 1)$$

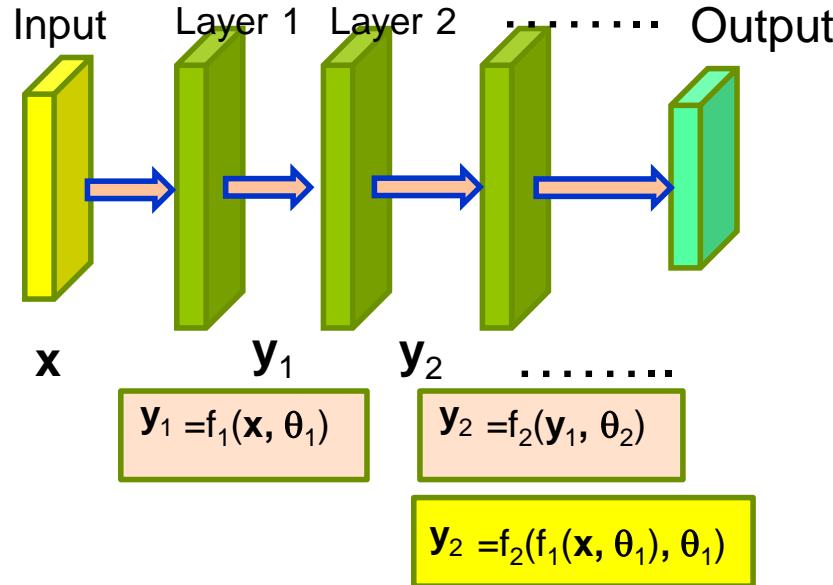
$$(4 * 3) * (3 * 1) = (4 * 1)$$

$$\vec{H} = \text{relu} \left( \left( \begin{bmatrix} W_{11} & W_{12} & W_{13} \\ W_{21} & W_{22} & W_{23} \\ W_{31} & W_{32} & W_{33} \\ W_{41} & W_{42} & W_{43} \end{bmatrix} * \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} \right) + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix} \right)$$

# Deep Network: Generic Definition

- A family of parametric, non-linear and hierarchical representation learning functions
  - massively optimized with stochastic gradient descent to encode domain knowledge, i.e. domain invariances, stationarity.
- $a^L(x; w^1, \dots, w^L) = h^L(h^{L-1} \dots h^1(x, w^1), w^{L-1}), w^L)$ 
  - $x$ : input,
  - $w^l$ : parameters for layer  $l$ ,
  - $a^l = h^l(x, w^l)$ : (non-) linear function
- Given training corpus  $\{X, Y\}$  find optimal parameters
  - $w^* \leftarrow \operatorname{argmin}_w \sum_{(x,y) \subseteq (X,Y)} L(y, a^L(x))$

# Forward Propagation: Deep Network



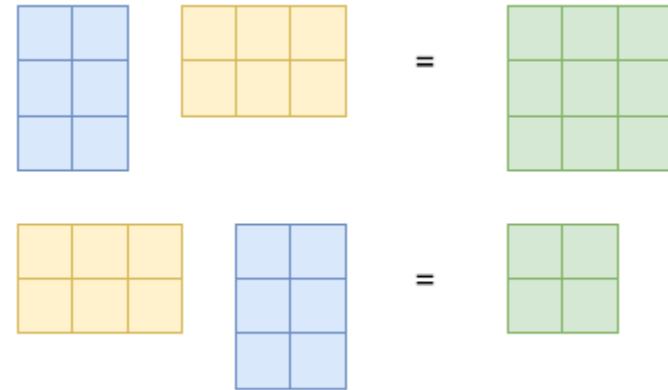
**Deep Network:** family of parametric, non-linear, hierarchical representations

$$y_N(x, \theta_1, \theta_2, \dots, \theta_N) = f_N(f_{N-1}(\dots(f_1(x, \theta_1), \theta_2), \dots, \theta_{N-1}), \theta_N)$$

Implemented as matrix/vector operations

# Matrix Multiplication

- Matrix multiplication is not commutative, i.e.,  $\mathbf{AB} \neq \mathbf{BA}$
- When multiplying matrices, the “neighbouring” dimensions have to fit:



$$\underbrace{\mathbf{A}}_{n \times k} \underbrace{\mathbf{B}}_{k \times m} = \underbrace{\mathbf{C}}_{n \times m}$$

$$\mathbf{y} = \mathbf{Ax}$$

$$\mathbf{y} = \mathbf{A}.\text{dot}(\mathbf{x})$$

$$y_i = \sum_j A_{ij} x_j$$

$$\mathbf{y} = \text{np.einsum('ij, j', A, \mathbf{x})}$$

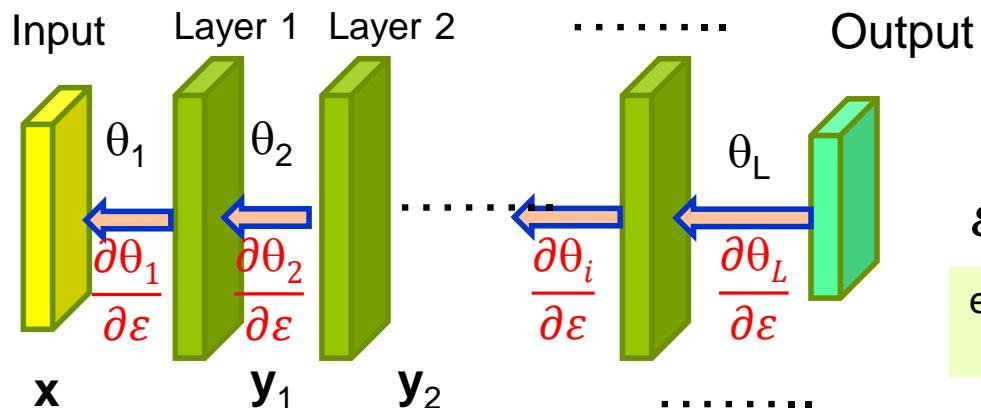
$$\mathbf{C} = \mathbf{AB}$$

$$\mathbf{C} = \mathbf{A}.\text{dot}(\mathbf{B})$$

$$C_{ij} = \sum_k A_{ik} B_{kj}$$

$$\mathbf{C} = \text{np.einsum('ik, kj', A, B)}$$

# Back Propagation: Deep Network



$$\varepsilon = t - y$$

error = true label – computed label

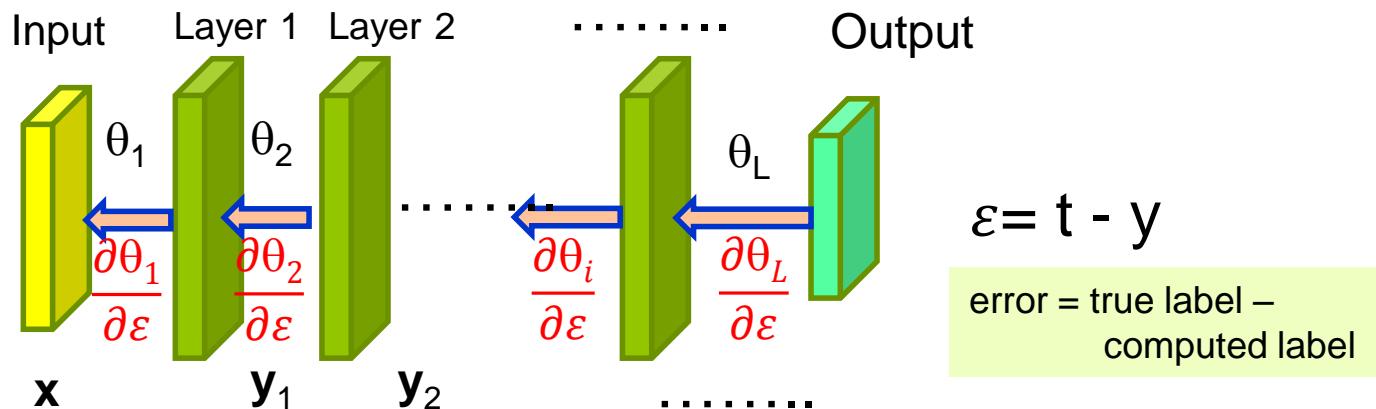
Compute how weights contribute to the error  $\varepsilon$

$$\frac{\partial \theta_i}{\partial \varepsilon}$$

**Training:** optimize network parameters to minimise loss over training set

$$\theta^* = \arg \min \sum_{(x,y) \in (X,Y)} J[y, f^L(x, \theta_1, \dots, \theta_L)]$$

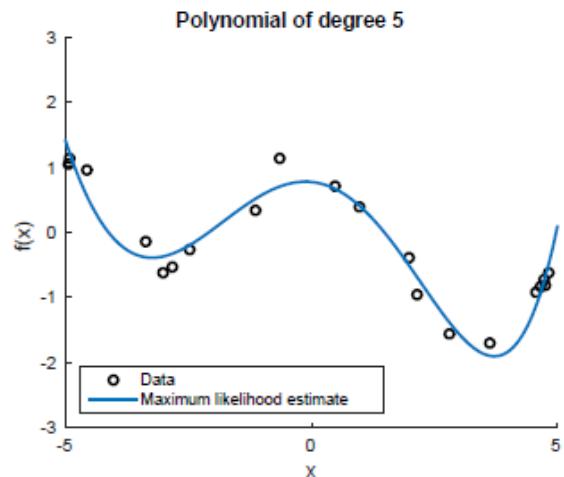
# Training: Computes Derivatives



- Must use partial derivative plus chain rule
- Shows that backpropagation can be done layer by layer

# Network Training

- ▶ Training data, e.g.,  $N$  pairs  $(x_i, y_i)$  of inputs  $x_i$  and observations  $y_i$
- ▶ Training the model means finding parameters  $\theta^*$ , such that  $f(x_i, \theta^*) \approx y_i$
- ▶ Define a loss function, e.g.,  $\sum_{i=1}^N (y_i - f(x_i, \theta))^2$ , which we want to optimize
- ▶ Typically: Optimization based on some form of gradient descent
  - ▶ Differentiation required



# Requirement for Differentiation

- Differentiate arbitrary formulae
  - Loss function:  $L = \Psi(t-y)^2$
  - Activation function:  $\tanh$
  - Algebraic operation:  $Wx$
- Rules of differentiation must be used
  - Sum, product, chain, etc.

# Differentiation: Rules

- ▶ Sum Rule

$$(f(x) + g(x))' = f'(x) + g'(x) = \frac{df}{dx} + \frac{dg}{dx}$$

- ▶ Product Rule

$$(f(x)g(x))' = f'(x)g(x) + f(x)g'(x) = \frac{df}{dx}g(x) + f(x)\frac{dg}{dx}$$

- ▶ Chain Rule

$$(g \circ f)'(x) = (g(f(x)))' = g'(f(x))f'(x) = \frac{dg}{df} \frac{df}{dx}$$

# Example

$$(g \circ f)'(x) = (g(f(x)))' = g'(f(x))f'(x) = \frac{dg}{df} \frac{df}{dx}$$

$$g(z) = \tanh(z)$$

$$z = f(x) = x^n$$

$$(g \circ f)'(x) = \underbrace{(1 - \tanh^2(x^n))}_{dg/df} \underbrace{nx^{n-1}}_{df/dx}$$

# Partial Derivatives

$$f : \mathbb{R}^N \rightarrow \mathbb{R}$$

$$y = f(\mathbf{x}), \quad \mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} \in \mathbb{R}^N$$

- ▶ Partial derivative (change one coordinate at a time):

$$\frac{\partial f}{\partial x_i} = \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_{i-1}, x_i + h, x_{i+1}, \dots, x_N) - f(\mathbf{x})}{h}$$

- ▶ Jacobian vector (gradient) collects all partial derivatives:

$$\frac{df}{d\mathbf{x}} = \left[ \frac{\partial f}{\partial x_1} \quad \dots \quad \frac{\partial f}{\partial x_N} \right] \in \mathbb{R}^{1 \times N}$$

Note: This is a row vector.

# Example

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}$$

$$f(x_1, x_2) = x_1^2 x_2 + x_1 x_2^3 \in \mathbb{R}$$

- ▶ Partial derivatives:

$$\frac{\partial f(x_1, x_2)}{\partial x_1} = 2x_1 x_2 + x_2^3$$

$$\frac{\partial f(x_1, x_2)}{\partial x_2} = x_1^2 + 3x_1 x_2^2$$

- ▶ Gradient:

$$\frac{df}{dx} = \begin{bmatrix} \frac{\partial f(x_1, x_2)}{\partial x_1} & \frac{\partial f(x_1, x_2)}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2x_1 x_2 + x_2^3 & x_1^2 + 3x_1 x_2^2 \end{bmatrix} \in \mathbb{R}^{1 \times 2}.$$

# Example: Chain Rule

- Consider the function

$$L(\mathbf{e}) = \frac{1}{2} \|\mathbf{e}\|^2 = \frac{1}{2} \mathbf{e}^\top \mathbf{e}$$

$$\mathbf{e} = \mathbf{y} - \mathbf{A}\mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^N, \mathbf{A} \in \mathbb{R}^{M \times N}, \mathbf{e}, \mathbf{y} \in \mathbb{R}^M$$

- Compute  $dL/d\mathbf{x}$ . What is the dimension/size of  $dL/d\mathbf{x}$ ?
- $dL/d\mathbf{x} \in \mathbb{R}^{1 \times N}$

$$\frac{dL}{d\mathbf{x}} = \frac{dL}{d\mathbf{e}} \frac{d\mathbf{e}}{d\mathbf{x}}$$

$$\frac{dL}{d\mathbf{e}} = \mathbf{e}^\top \in \mathbb{R}^{1 \times M}$$

$$\frac{d\mathbf{e}}{d\mathbf{x}} = -\mathbf{A} \in \mathbb{R}^{M \times N}$$

$$\Rightarrow \frac{dL}{d\mathbf{x}} = \mathbf{e}^\top (-\mathbf{A}) = -(y - \mathbf{A}\mathbf{x})^\top \mathbf{A} \in \mathbb{R}^{1 \times N}$$

# Jacobian Matrix

$$f : \mathbb{R}^N \rightarrow \mathbb{R}^M$$

$$\mathbf{y} = f(\mathbf{x}) \in \mathbb{R}^M, \quad \mathbf{x} \in \mathbb{R}^N$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} f_1(x_1, \dots, x_N) \\ \vdots \\ f_M(x_1, \dots, x_N) \end{bmatrix}$$

- Jacobian matrix (collection of all partial derivatives)

$$\begin{bmatrix} \frac{dy_1}{d\mathbf{x}} \\ \vdots \\ \frac{dy_M}{d\mathbf{x}} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_N} \\ \vdots & & \vdots \\ \frac{\partial f_M}{\partial x_1} & \dots & \frac{\partial f_M}{\partial x_N} \end{bmatrix} \in \mathbb{R}^{M \times N}$$

# Example

$$f(\mathbf{x}) = \mathbf{A}\mathbf{x}, \quad f(\mathbf{x}) \in \mathbb{R}^M, \quad \mathbf{A} \in \mathbb{R}^{M \times N}, \quad \mathbf{x} \in \mathbb{R}^N$$

- Compute the gradient  $df/d\mathbf{x}$

- Dimension of  $df/d\mathbf{x}$ :

Since  $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$ , it follows that  $df/d\mathbf{x} \in \mathbb{R}^{M \times N}$

- Gradient:

$$f_i = \sum_{j=1}^N A_{ij}x_j \quad \Rightarrow \quad \frac{\partial f_i}{\partial x_j} = A_{ij}$$

$$\Rightarrow \frac{df}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_N} \\ \vdots & & \vdots \\ \frac{\partial f_M}{\partial x_1} & \dots & \frac{\partial f_M}{\partial x_N} \end{bmatrix} = \begin{bmatrix} A_{11} & \dots & A_{1N} \\ \vdots & & \vdots \\ A_{M1} & \dots & A_{MN} \end{bmatrix} = \mathbf{A}$$

# Example

- Consider  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $x : \mathbb{R} \rightarrow \mathbb{R}^2$

$$f(x) = f(x_1, x_2) = x_1^2 + 2x_2,$$

$$x(t) = \begin{bmatrix} x_1(t) \\ x_2(t) \end{bmatrix} = \begin{bmatrix} \sin(t) \\ \cos(t) \end{bmatrix}$$

- The dimensions  $df/dx$  and  $dx/dt$  are  $1 \times 2$  and  $2 \times 1$ , respectively
- Compute the gradient  $df/dt$  using the chain rule.

$$\begin{aligned}\frac{df}{dt} &= \left[ \frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \right] \begin{bmatrix} \frac{\partial x_1}{\partial t} \\ \frac{\partial x_2}{\partial t} \end{bmatrix} \\ &= \begin{bmatrix} 2 \sin t & 2 \end{bmatrix} \begin{bmatrix} \cos t \\ -\sin t \end{bmatrix} \\ &= 2 \sin t \cos t - 2 \sin t = 2 \sin t(\cos t - 1)\end{aligned}$$

# Derivatives with Matrices: Summary

- Re-cap: Gradient of a function  $f : \mathbb{R}^D \rightarrow \mathbb{R}^E$  is an  $E \times D$ -matrix:

# target dimensions  $\times$  # parameters

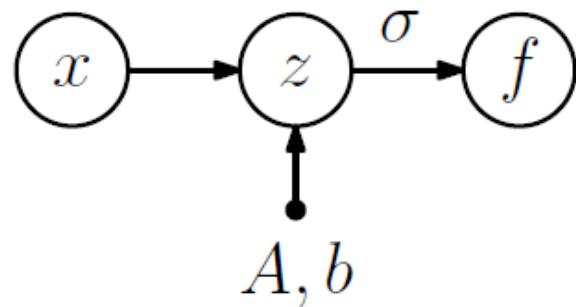
with

$$\frac{df}{dx} \in \mathbb{R}^{E \times D}, \quad df[e, d] = \frac{\partial f_e}{\partial x_d}$$

- Generalization to cases, where the parameters ( $D$ ) or targets ( $E$ ) are matrices, apply immediately
- Assume  $f : \mathbb{R}^{M \times N} \rightarrow \mathbb{R}^{P \times Q}$ , then the gradient is a  $(P \times Q) \times (M \times N)$  object (tensor) where

$$df[p, q, m, n] = \frac{\partial f_{pq}}{\partial X_{mn}}$$

# Gradients of Single-Layer Network



$$f = \tanh(\underbrace{Ax + b}_{=:z \in \mathbb{R}^M}) \in \mathbb{R}^M, \quad x \in \mathbb{R}^N, A \in \mathbb{R}^{M \times N}, b \in \mathbb{R}^M$$

$$\frac{\partial f}{\partial b} = \underbrace{\frac{\partial f}{\partial z}}_{M \times M} \underbrace{\frac{\partial z}{\partial b}}_{M \times M} \in \mathbb{R}^{M \times M}$$

$$\frac{\partial f}{\partial A} = \underbrace{\frac{\partial f}{\partial z}}_{M \times M} \underbrace{\frac{\partial z}{\partial A}}_{M \times (M \times N)} \in \mathbb{R}^{M \times (M \times N)}$$

# Example

$$f = \tanh(\underbrace{Ax + b}_{=: z \in \mathbb{R}^M}) \in \mathbb{R}^M, \quad x \in \mathbb{R}^N, A \in \mathbb{R}^{M \times N}, b \in \mathbb{R}^M$$

$$\frac{\partial f}{\partial b} = \underbrace{\frac{\partial f}{\partial z}}_{M \times M} \underbrace{\frac{\partial z}{\partial b}}_{M \times M} \in \mathbb{R}^{M \times M}$$

$$\frac{\partial f}{\partial z} = \text{diag}(1 - \tanh^2(z)) \in \mathbb{R}^{M \times M}$$

$$\frac{\partial z}{\partial b} = I \in \mathbb{R}^{M \times M} \tag{5}$$

$$\frac{\partial f}{\partial b}[i, j] = \sum_{l=1}^M \frac{\partial f}{\partial z}[i, l] \frac{\partial z}{\partial b}[l, j]$$

```
dfdb = np.einsum('il, lj', dfdz, dzdb)
```

# Example (continued)

$$f = \tanh(\underbrace{Ax + b}_{=:z \in \mathbb{R}^M}) \in \mathbb{R}^M, \quad x \in \mathbb{R}^N, A \in \mathbb{R}^{M \times N}, b \in \mathbb{R}^M$$

$$\frac{\partial f}{\partial A} = \underbrace{\frac{\partial f}{\partial z}}_{M \times M} \underbrace{\frac{\partial z}{\partial A}}_{M \times (M \times N)} \in \mathbb{R}^{M \times (M \times N)}$$

$$\frac{\partial f}{\partial z} = \text{diag}(1 - \tanh^2(z)) \in \mathbb{R}^{M \times M} \tag{6}$$

$$\frac{\partial z}{\partial A} \quad \blacktriangleright \text{ See (4)}$$

$$\frac{\partial f}{\partial A}[i, j, k] = \sum_{l=1}^M \frac{\partial f}{\partial z}[i, l] \frac{\partial z}{\partial A}[l, j, k]$$

```
dfdA = np.einsum('il, ljk', dfdz, dzdA)
```

# Summary: Single Layer NN

- Inputs  $x$ , observed outputs  $y = f(z, \theta) + \epsilon$ ,  $\epsilon \sim \mathcal{N}(\mathbf{0}, \Sigma)$
- Train single-layer neural network with

$$f(z, \theta) = \tanh(z), \quad z = Ax + b, \quad \theta = \{A, b\}$$

- Find  $A, b$ , such that the squared loss

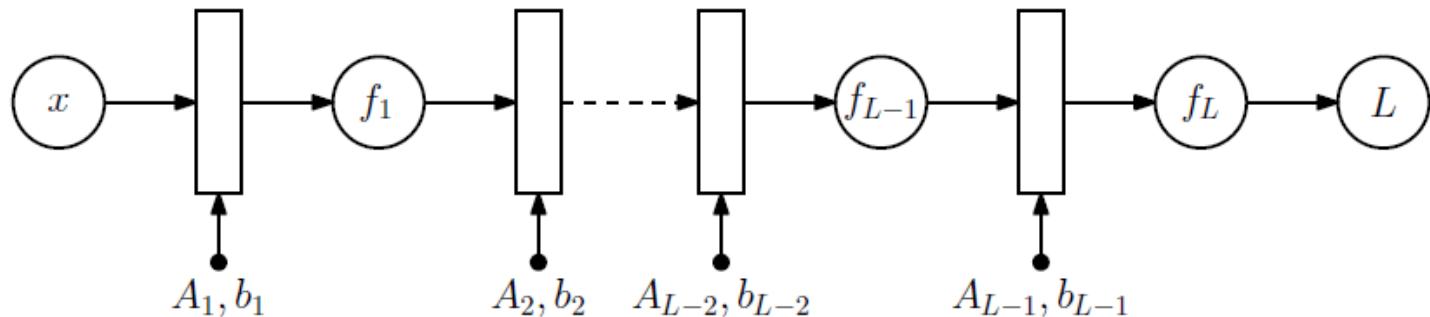
$$L(\theta) = \frac{1}{2} \|e\|^2, \quad e = y - f(z, \theta)$$

is minimized

- Partial derivatives:

$$\left. \begin{array}{l} \frac{\partial L}{\partial A} = \frac{\partial L}{\partial e} \frac{\partial e}{\partial f} \frac{\partial f}{\partial z} \frac{\partial z}{\partial A} \\ \frac{\partial L}{\partial b} = \frac{\partial L}{\partial e} \frac{\partial e}{\partial f} \frac{\partial f}{\partial z} \frac{\partial z}{\partial b} \end{array} \right\} \quad \begin{array}{lll} \frac{\partial L}{\partial e} \blacktriangleright (1) & \frac{\partial e}{\partial f} \blacktriangleright (2), (3) & \frac{\partial f}{\partial z} \blacktriangleright (6) \\ \frac{\partial z}{\partial A} \blacktriangleright (4) & \frac{\partial z}{\partial b} \blacktriangleright (5) & \end{array}$$

# Multi-Layer NN: Gradients



- Inputs  $x$ , observed outputs  $y$
- Train multi-layer neural network with

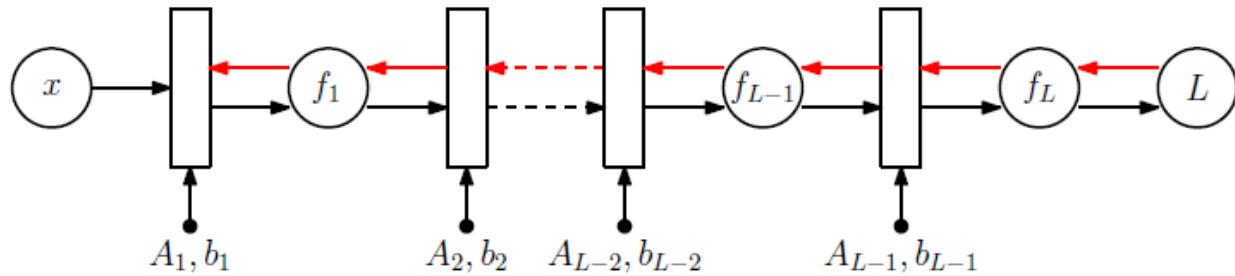
$$f_0 = x$$

$$f_i = \sigma_i(A_{i-1}f_{i-1} + b_{i-1}), \quad i = 1, \dots, L$$

- Find  $A_j, b_j$  for  $j = 0, \dots, L - 1$ , such that the squared loss

$$L(\theta) = \|y - f_L(\theta, x)\|^2$$

is minimized, where  $\theta = \{A_j, b_j\}, \quad j = 0, \dots, L - 1$



$$\frac{\partial L}{\partial \theta_{L-1}} = \frac{\partial L}{\partial f_L} \frac{\partial f_L}{\partial \theta_{L-1}}$$

$$\frac{\partial L}{\partial \theta_{L-2}} = \frac{\partial L}{\partial f_L} \boxed{\frac{\partial f_L}{\partial f_{L-1}} \frac{\partial f_{L-1}}{\partial \theta_{L-2}}}$$

$$\frac{\partial L}{\partial \theta_{L-3}} = \frac{\partial L}{\partial f_L} \frac{\partial f_L}{\partial f_{L-1}} \boxed{\frac{\partial f_{L-1}}{\partial f_{L-2}} \frac{\partial f_{L-2}}{\partial \theta_{L-3}}}$$

$$\frac{\partial L}{\partial \theta_i} = \frac{\partial L}{\partial f_L} \frac{\partial f_L}{\partial f_{L-1}} \dots \boxed{\frac{\partial f_{i+2}}{\partial f_{i+1}} \frac{\partial f_{i+1}}{\partial \theta_i}}$$

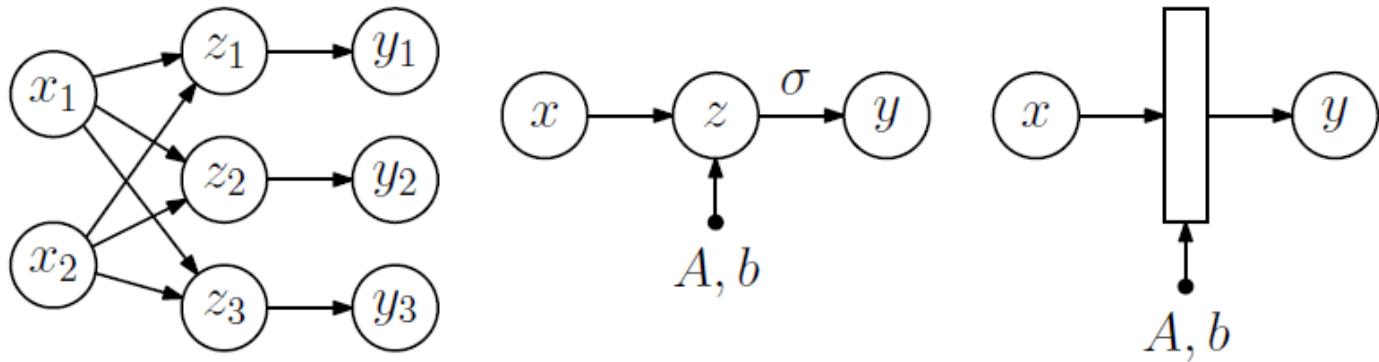
# Summary

- Deep Learning is based on core mathematical principles
- Review of
  - Linear algebra
  - Multivariable calculus (chain rule, etc.)

# Feed-Forward Neural Network

$$y = \sigma(z)$$

$$z = Ax + b$$



- Training a neural network means parameter optimization:
  - Typically via some form of gradient descent
- **Challenge 1: Differentiation.** Compute gradients of a loss function with respect to neural network parameters  $A$ ,  $b$ 
  - Computing statistics (e.g., means, variances) of predictions
- **Challenge 2: Integration.** Propagate uncertainty through a neural network