

Linear Algebra:

- ℓ_2 norm:

$$\|x\|_2 = \sqrt{\sum_{i=1}^N x_i^2}$$

- linearly dependent - vector is linear combination of other vectors -

$$x_n = \sum_{i=1}^{n-1} \alpha_i x_i$$

- column rank - largest number of columns that are linearly independent set
- rank - assumed col rank = row rank for a matrix $A \in \mathbb{R}^{m \times n}$
- inverse of a square matrix, A is invertible if this exists:

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

- orthogonal vectors $x^T y = 0$ | normalized vectors $\|x\|_2 = 1$

- orthogonal matrices - square matrix with columns orthogonal to each other and normalized (orthonormal columns)

$$U^T U = I = U U^T$$

↑
orthogonal square matrix with orthonormal columns also $U^T = U^{-1}$ inverse of orthogonal matrix is its transpose. *why? square

- span of a set of vectors is all vectors expressed as linear combination.
- projection of a vector onto the span of matrices

$$\text{Proj}(y; \{x_1, \dots, x_n\}) = \underset{v \in \text{span}(\{x_1, \dots, x_n\})}{\text{argmin}} \underbrace{\|y - v\|_2}_{\text{euclidean norm}}$$

↑
vector y
projected onto span of
 x_1, \dots, x_n

- range is span of cols of matrix A
- nullspace of a matrix is all vectors $= 0$ when multiplied by A
- determinant of a square matrix is a function. Absolute

value of determinant of square matrix A is volume of set S^2 .

↑
if 2×2 matrix \rightarrow
volume is area in 2D

• determinant of identity is 1 $\det I = 1$, volume of

a unit hypercube = 1

• gradient ($\mathbb{R}^n \rightarrow \mathbb{R}^1$) $\rightarrow \nabla f(x, y) = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix}$

• jacobian ($\mathbb{R}^n \rightarrow \mathbb{R}^m$) vector to vector of same or diff shape

• $f(x, y) = \begin{bmatrix} 2x + y^3 \\ e^y - 13x \end{bmatrix}$

$$J = \begin{matrix} \mathbb{R}^2 & \mathbb{R}^2 \end{matrix} \begin{bmatrix} \frac{\partial f_1}{\partial x} & \frac{\partial f_1}{\partial y} \\ \frac{\partial f_2}{\partial x} & \frac{\partial f_2}{\partial y} \end{bmatrix} = \begin{bmatrix} 2 & 3y^2 \\ -13 & e^y \end{bmatrix} \begin{matrix} \text{on} \nearrow \\ \nabla f_i^T \\ \nabla f_i^T \end{matrix}$$

• jacobian chain rule

$$f(x, y) = \begin{bmatrix} \sin(x^2 + y) \\ \ln(y^3) \end{bmatrix}$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial x} \text{ where } f(x) = \begin{bmatrix} \sin(g) \\ \ln(g) \end{bmatrix}, g(x) = \begin{bmatrix} x^2 + y \\ y^3 \end{bmatrix}$$

$$= \begin{bmatrix} \cos g_1 & 0 \\ 0 & \frac{1}{g_2} \end{bmatrix} \begin{bmatrix} 2x \\ 0 \end{bmatrix} \begin{bmatrix} 1 \\ 3y^2 \end{bmatrix}$$

$$= \begin{bmatrix} 2x \cos g_1 & \cos g_1 \\ 0 & \frac{3y^2}{g_2} \end{bmatrix}$$

$$= \begin{bmatrix} 2x \cos(x^2 + y) & \cos(x^2 + y) \\ 0 & \frac{3}{y} \end{bmatrix}$$

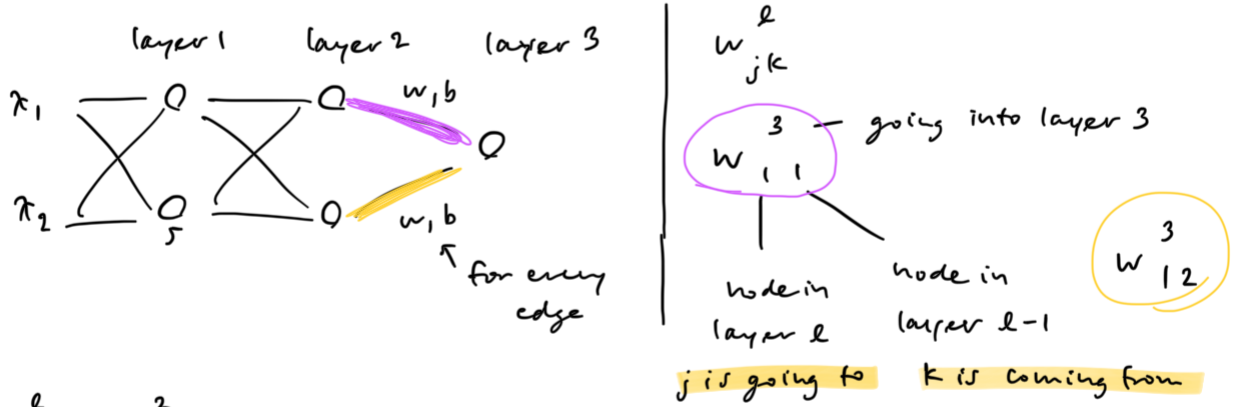
* if intermediate
can't be chain
rule-d then put 1

$$\begin{matrix} x_1 & w_1 \\ x_2 & w_2 \\ x_3 & w_3 \\ x_4 & w_4 \\ x_5 & w_5 \end{matrix} \quad \textcircled{N} \rightarrow \text{scalar}$$

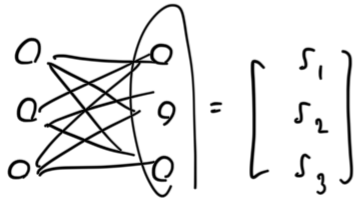
ReLU \swarrow vector to

$$\sigma \left(\underbrace{\sum_{i=1}^N x_i w_i}_{x^T w} \right) + b$$

Scalar function.

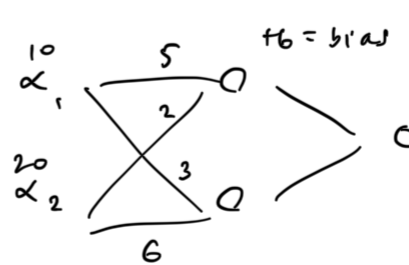
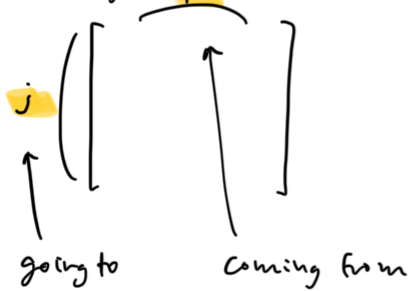


$$b_j^l \quad b_i^2 = 5$$



each node passes weights and biases to next node. output vector with scalars is passed to next layer.

weight matrix: w^l for entire layer l (feed forward)



backprop improves weights after fully trained network.

$$\sigma(w^l \cdot a^{l-1} + b^l)$$

$\sigma(z')$

OR

$$a_0 = \sigma(w^1 a^0 + b^1)$$

$$\sigma(w^2 a^1 + b^2) = \text{etc.}$$

$$w^1 = \begin{bmatrix} 1,1 & 1,2 \\ 2,1 & 2,2 \\ 3 & 6 \end{bmatrix} \begin{bmatrix} 10 \\ 20 \end{bmatrix} = \begin{bmatrix} 90 \\ 150 \end{bmatrix} + 6$$

$$(2,2) \times (2,1) = (2 \times 1)$$

$$z' = \begin{bmatrix} 96 \\ 156 \end{bmatrix}$$

is output vector, now put into activation function that's non-linear

Formulas for all ML topics:

- Unparametrized - ex k-NN, no assumptions on size of input data.
- Regularization - applying penalties to parameters of a model.

$$L(w) = \frac{1}{2} \sum_{n=1}^N (y_n - \underbrace{w^T \phi_n}_\text{spanning weight})^2 + \frac{\lambda}{2} w^T w$$

↑

spanning weight

larger values of λ produce