# Chain rule

✳ <u>Chain Rule</u>

$$\frac{d}{dx}[f(g(x))] = f'(g(x)) \cdot g'(x)$$

$$\frac{d}{dx}\left[\ln\left(\underbrace{\sin(x)}_{g(x)}\right)\right]$$

$\underbrace{\qquad}_{f(g(x))}$

$$f'(g(x)) = \frac{1}{g(x}$$

$$\frac{d}{dx}\left[\underbrace{\ln(x)}_{f(x)}\,\underbrace{\sin(x)}_{g(x)}\right]$$
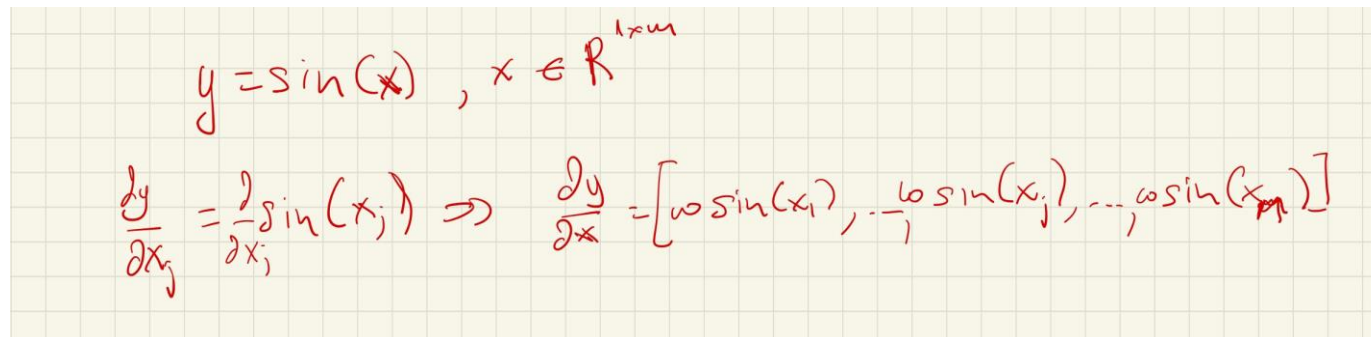
Product Rule

# Gradient

o Assuming our input is a row vector, that is $x \in \mathbb{R}^{1 \times M}$

o The gradient is a vector containing all partial derivatives

$$\frac{dh}{dx} = \nabla_x h = [\frac{\partial h}{\partial x_1}, \dots, \frac{\partial h}{\partial x_M}]$$

o Generalization of the derivative, defined on a univariate function ($M = 1$)

# Example

o Often, easier to write things out explicitly

o Let's say $y = \sin(\boldsymbol{x})$, where $\dim(\boldsymbol{x}) = 1 \times M$

$$y = \sin(x) \quad, \quad x \in R^{1 \times M}$$

$$\frac{\partial y}{\partial x_j} = \frac{\partial}{\partial x_j} \sin(x_j) \Rightarrow \frac{\partial y}{\partial x} = \left[ \omega \sin(x_1), \ldots \frac{\omega}{j} \sin(x_j), \ldots, \omega \sin(x_m) \right]$$

# Jacobian

o Generalization of the gradient for vector-valued functions $\boldsymbol{h}(\boldsymbol{x})$
  ◦ all input dimensions contribute to all output dimensions

$$J = \nabla_{\boldsymbol{x}} \boldsymbol{h} = \frac{d\boldsymbol{h}}{d\boldsymbol{x}} = \begin{bmatrix} \dfrac{\partial h_1}{\partial x_1} & \cdots & \dfrac{\partial h_1}{\partial x_M} \\ \vdots & \ddots & \vdots \\ \dfrac{\partial h_N}{\partial x_1} & \cdots & \dfrac{\partial h_N}{\partial x_M} \end{bmatrix}$$

o Single input, single output → ▪

o Multiple input, single output → ▪

o Single input, multiple output → ▪

o Multiple input, multiple output → ▪

# Taking gradients with index notation for matrices/vectors…

- Often, output is a vector/matrix/tensor depend on matrix/vector/tensor input

- We still want to see what is the effect of the output w.r.t. the input. How?

- Better use index notation
  - Assign input and output indices and take derivatives with scalar quantities
  - E.g., $y = x \cdot w^T$, where $\dim(y) = S \times N$, $\dim(x) = S \times M$, $\dim(w) = N \times M$

# Jacobians, gradients, intuitively

o The Jacobians, gradients and the likes ($\frac{dh}{dx}$) qualitatively capture the same thing

 ◦ $\underbrace{\text{Change in the output}}_{dh}$ with respect to $\underbrace{\text{change in the input}}_{dx}$

o That is, the final Jacobian/gradient/… is simply a <u>tensor</u> $\nabla$ with the shape

 ◦ $\dim(\nabla) = \text{shape}_{\text{out}} \times \text{shape}_{\text{in}}$

 ◦ If our 'in' is a vector, then we append that shape to the tensor gradient

 ◦ The [Einstein notation](#) can be useful ([np.einsum](#)) for the computations

# Jacobian, geometrically

o The Jacobian represents the best local approximation of how the space changes under a (non-linear) transformation
  ◦ Not unlike derivative being the best linear approximation of a curve (tangent)

o The Jacobian determinant (for square matrices) measures the ratio of areas
  ◦ Similar to what the 'absolute slope' measures in the 1d case (derivative)
  ◦ Used in change of variables (integration by substitution), normalizing flows

# Basic rules of partial differentiation

○ Product rule

○ $\frac{\partial}{\partial x}\big(f(x) \cdot g(x)\big) = f(x) \cdot \frac{\partial}{\partial x} g(x) + g(x) \cdot \frac{\partial}{\partial x} f(x)$

○ Sum rule

○ $\frac{\partial}{\partial x}\big(f(x) + g(x)\big) = \frac{\partial}{\partial x} f(x) + \frac{\partial}{\partial x} g(x)$

# Computing gradients in complex functions: Chain rule

o Assume a composite function, $h = h_L\left(h_{L-1}\left(\ldots\left(h_1\left(\boldsymbol{x}\right)\right)\right)\right)$, or

$$h = h_L \circ h_{L-1} \circ \cdots \circ h_1\left(\boldsymbol{x}\right)$$

o To compute the derivative/gradient, we can use the chain rule
  ◦ Intuitively, similar to matrix multiplications

$$\frac{dh}{dx} = \frac{dh}{dh_L} \cdot \frac{dh_L}{dh_{L-1}} \cdot \ldots \cdot \frac{dh_1}{dx}$$

o Each $\frac{dh_i}{dh_{i-1}}$ is a Jacobian/gradient/… vector/matrix/tensor

o Make sure each component matches dimensions

# Chain rule and tensors, intuitively



o What does the chain rule stand for with high-dimensional tensors

o Let's keep it simple: $\frac{d\boldsymbol{h}}{d\boldsymbol{x}} = \frac{d\boldsymbol{h}}{d\boldsymbol{g}} \cdot \frac{d\boldsymbol{g}}{d\boldsymbol{x}}$

  ◦ $\boldsymbol{h}(\boldsymbol{g})$ has $M$ inputs, $N$ outputs
  ◦ $\boldsymbol{g}(\boldsymbol{x})$ has $K$ inputs (because of $\boldsymbol{x}$), $M$ outputs

o We can think of the chain rule as
  ◦ summing over all possible changes
  ◦ caused to $\boldsymbol{h}$ by each element in $\boldsymbol{x}$ via all possible $\boldsymbol{g}$'s

o For high-dim tensors, $\boldsymbol{h}, \boldsymbol{g}, \boldsymbol{x}$, we apply the same logic
  ◦ Replace shape of the vector with shape of tensor
  ◦ Do the summations keeping those shapes fixed
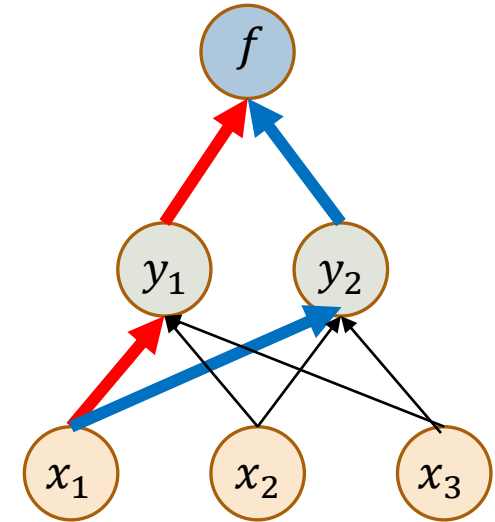  ◦ Think it in terms of indices, again [Einstein notation](Einstein notation)

# Example

o For h $= f \circ y(x)$

$$\frac{dh}{dx} = \frac{df}{dy}\frac{dy}{dx} = \begin{bmatrix} \frac{\partial f}{\partial y_1} & \frac{\partial f}{\partial y_2} \end{bmatrix} \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \frac{\partial y_1}{\partial x_3} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \frac{\partial y_2}{\partial x_3} \end{bmatrix}$$

o Focusing on one of the partial derivatives: $\frac{\partial h}{\partial x_1}$

$$\frac{\partial h}{\partial x_1} = \frac{\partial f}{\partial y_1}\frac{\partial y_1}{\partial x_1} + \frac{\partial f}{\partial y_2}\frac{\partial y_2}{\partial x_1}$$

o The partial derivative depends on all paths from $f$ to $x_i$

# Example

o For h $= f \circ y(x)$

$$\frac{dh}{dx} = \frac{df}{dy}\frac{dy}{dx} = \begin{bmatrix} \frac{\partial f}{\partial y_1} & \frac{\partial f}{\partial y_2} \end{bmatrix} \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \frac{\partial y_1}{\partial x_3} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \frac{\partial y_2}{\partial x_3} \end{bmatrix}$$

o Focusing on one of the partial derivatives: $\frac{\partial h}{\partial x_2}$

$$\frac{\partial h}{\partial x_2} = \frac{\partial f}{\partial y_1}\frac{dy_1}{dx_2} + \frac{\partial f}{\partial y_2}\frac{\partial y_2}{\partial x_2}$$

o The partial derivative depends on all paths from $f$ to $x_i$