



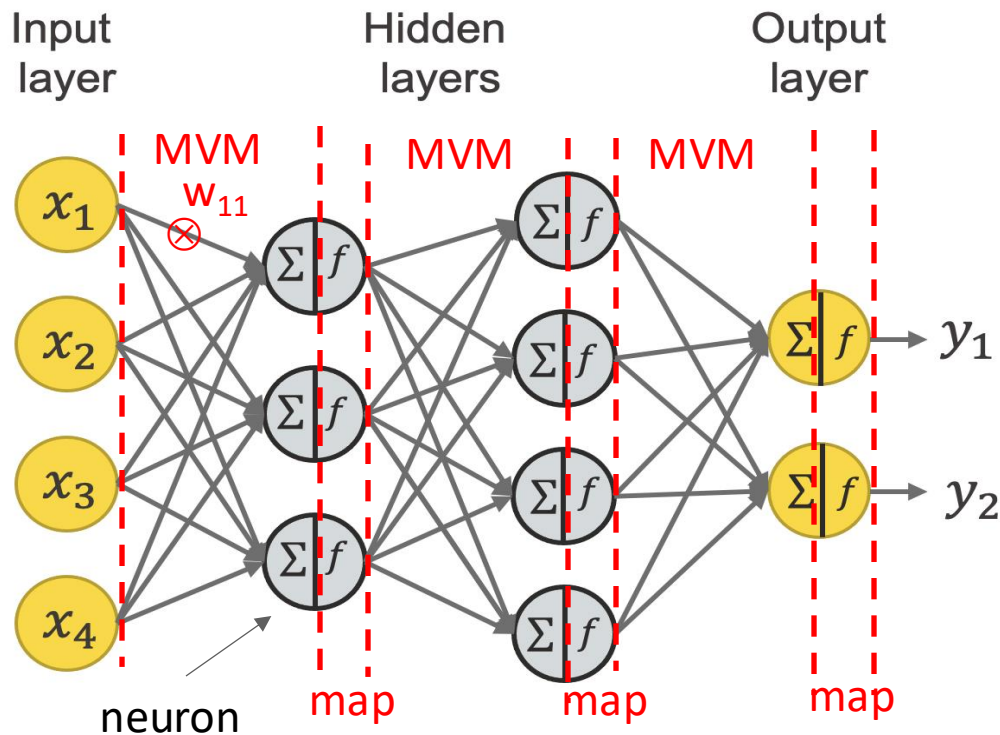
Back Propagation Made Easy

Keshav Pingali
The University of Texas at Austin

Organization

- **Problem**
 - Parameter estimation in neural networks
 - Non-linear optimization problem
 - Solved using gradient descent
 - Usual presentations of gradient descent are hard to understand and difficult to extend to irregular network connections
- **Parameterized programs**
 - Abstraction for neural networks
- **Gradient computation in parameterized programs**
 - Backward dataflow analysis
- **Advantages**
 - Compositional algorithm for gradient computation
 - Handles weight-sharing
 - Easy to extend to skip connections and other irregular networks
 - Natural extension to higher dimensional data such as tensors

Multilayer Perceptron (MLP) Example



Frank Rosenblatt (Cornell)
Inventor of Perceptron

- **Type:** Inputs: $x_1..x_4$, outputs: y_1, y_2 (all \mathbb{R})
- **Scalar view:**
 - Each edge performs a multiplication with a real-valued **parameter**
 - Each neuron adds its input values and applies a non-linear operation f such as tanh, ReLU etc. (known as **activation functions**)
- **Vector view:**
 - Each layer performs a dense matrix vector multiplication
 - Followed by pointwise (map) non-linear operation f
- Abstraction of MLP and more complex neural networks: **parameterized programs**

Parameterized program: running example

- Type of desired function: real x real \rightarrow real
- Training data: set of N 3-tuples $\{(p_i, q_i, t_i)\}$
- Model
 - **Composed from**
 - > **base functions** f_i (may be **nonlinear** such as tanh, sigmoid, sin, cos, ...)
 - > **parameters** W_i : *real* (assume no weight sharing so each weight occurs just once)
- Notation: capital letters for variable names, small letters for variable values
 - Function written as $R(w; p_i, q_i)$ where w is (w_0, w_1, w_2)
- Parameter optimization
 - Square error for training sample $(p_i, q_i, t_i) = (t_i - R(w; p_i, q_i))^2$
 - Goal: choose (w_0, w_1, w_2) to minimize mean square error
$$\text{Loss}(w_0, w_1, w_2) = \frac{1}{N} \sum_{i=1}^N (t_i - R(w; p_i, q_i))^2$$

Function $R(P, Q)$ {
 *$A = W_0 * P$*
 $B = f_0(A, Q)$
 *$C = W_1 * B$*
 *$D = W_2 * B$*
 $E = f_1(C)$
 $F = f_2(D)$
 $R = f_3(E, F)$
return R }

Parameter optimization

- Find derivatives of Loss wrt W_0, W_1, W_2

$$\text{Loss}(w_0, w_1, w_2) = \frac{1}{N} \sum_{i=1}^N (t_i - R(w; p_i, q_i))^2$$

$$\frac{\partial \text{Loss}}{\partial W_0}(w_0, w_1, w_2) = -\frac{2}{N} \sum_{i=1}^N (t_i - R(w; p_i, q_i)) \frac{\partial R}{\partial W_0}(w; p_i, q_i)$$

$$\frac{\partial \text{Loss}}{\partial W_1}(w_0, w_1, w_2) = -\frac{2}{N} \sum_{i=1}^N (t_i - R(w; p_i, q_i)) \frac{\partial R}{\partial W_1}(w; p_i, q_i)$$

$$\frac{\partial \text{Loss}}{\partial W_2}(w_0, w_1, w_2) = -\frac{2}{N} \sum_{i=1}^N (t_i - R(w; p_i, q_i)) \frac{\partial R}{\partial W_2}(w; p_i, q_i)$$

$$\nabla_W R(w; p_i, q_i)$$

Function $R(P, Q)$ {

$A = W_0 * P$

$B = f_0(A, Q)$

$C = W_1 * B$

$D = W_2 * B$

$E = f_1(C)$

$F = f_2(D)$

$R = f_3(E, F)$

return R }

Derivatives are complicated, non-linear functions
so use gradient-descent for parameter optimization

Parameterized program as flow graph

Function $R(P,Q)$ {

$A = W0 * P$

$B = f0(A,Q)$

$C = W1 * B$

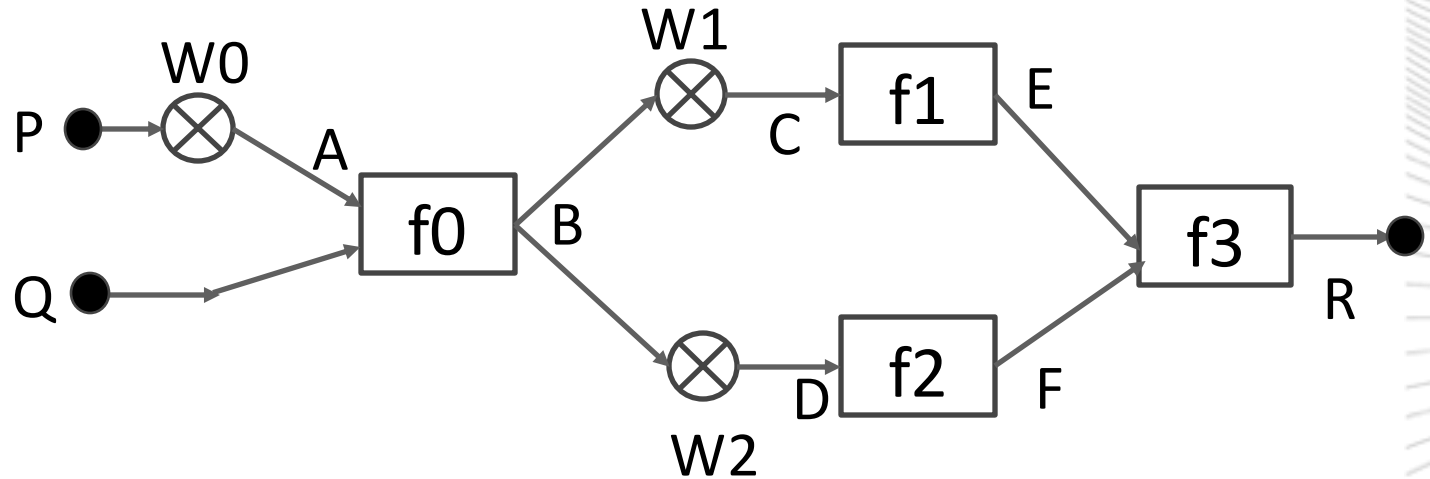
$D = W2 * B$

$E = f1(C)$

$F = f2(D)$

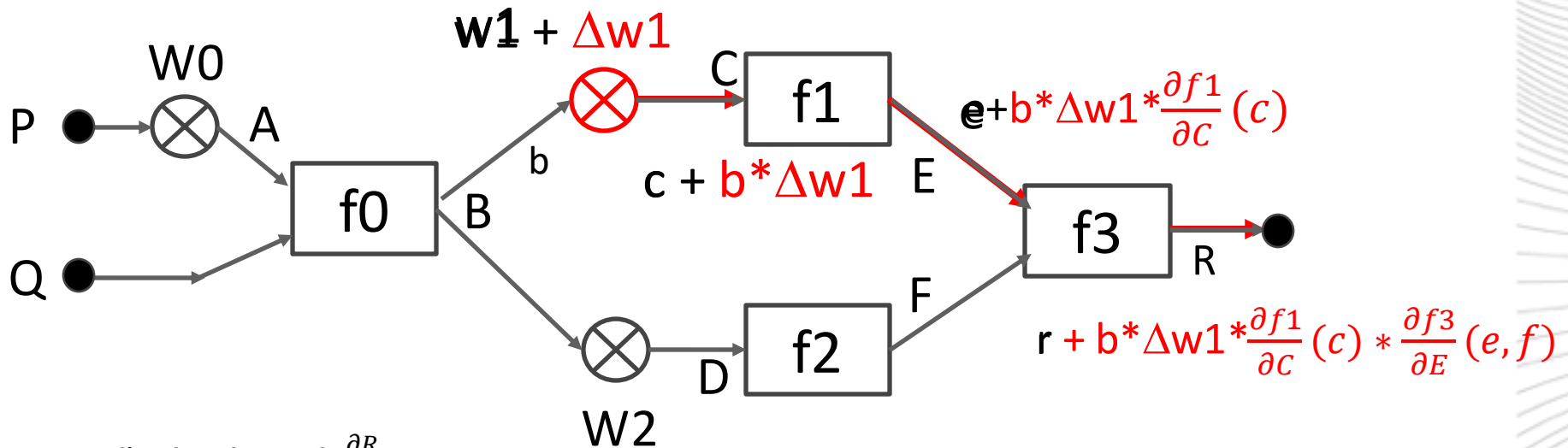
$R = f3(E,F)$

return R }



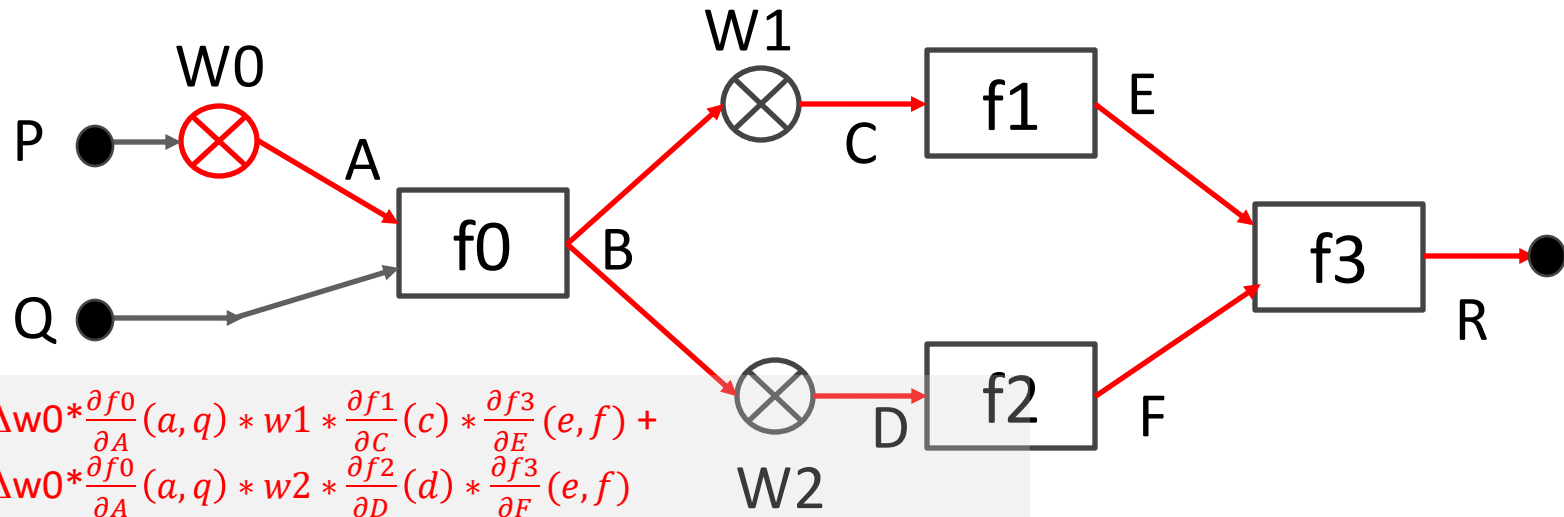
- Useful to represent multiplication by weights differently from functions f_i
 - Functions f_i are fixed but weights change during training
- Execution models
 - Sequential:
 - > Execute nodes sequentially in any topological order
 - Parallel:
 - > Asynchronous dataflow: node executes when inputs are available
 - > **Forward propagation**: level-by-level schedule of vertices
- All values at intermediate points (A, B, C, \dots) are stored
 - Needed for gradient computations

Value of $\frac{\partial R}{\partial w_1}(w; p_i, q_i)$



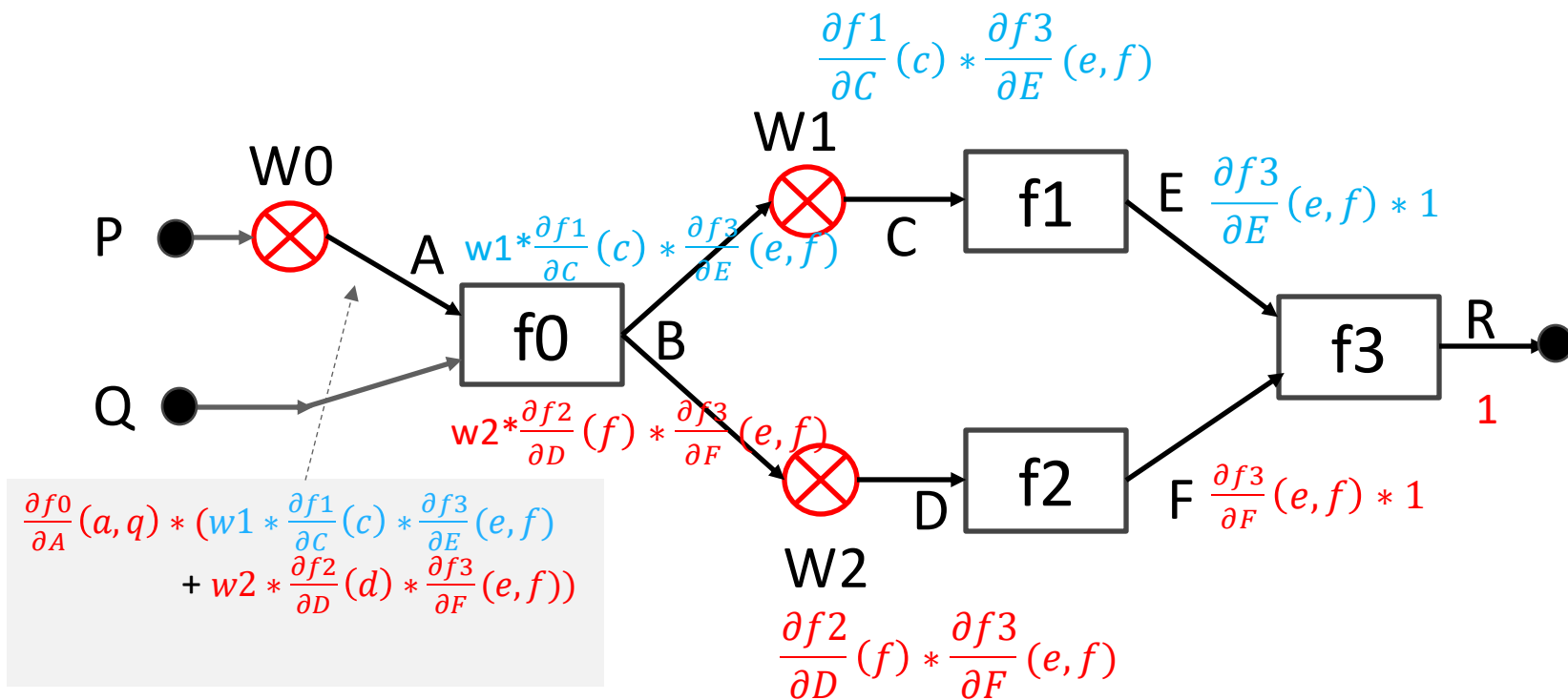
- To find value of $\frac{\partial R}{\partial w_1}(w; p_i, q_i)$
 - Multiply values of partial derivatives of all vertices on path from W1 to R
 - Multiply result by value of input to W1 (i.e., b)
 - Result: $b * \frac{\partial f_1}{\partial c}(c) * \frac{\partial f_3}{\partial E}(e, f)$
 - Compute the product either forwards or backwards along path
- In general, given path $\rho : X \xrightarrow{*} Y$
 - $\pi(\rho)$ = product of derivatives of nodes on path excluding X and Y (*path derivative*)
= 1 for empty path or if there are no intermediate nodes
- $\frac{\partial R}{\partial w_1}(w; p_i, q_i) = b * \pi(W1 \xrightarrow{*} R)$

Value of $\frac{\partial R}{\partial w_0}(w; p_i, q_i)$



- In general, there is a DAG from weight to the output
 - Value of partial derivative:
 - Enumerate all paths from weight to output and add up the contributions of all paths
- $$\frac{\partial R}{\partial w_0}(w_0; p_i, q_i) = p_i * \sum_{\rho \in \mathcal{P}(W_0)} \pi(\rho)$$
- (where $\mathcal{P}(W_0)$ is the set of paths from W_0 to exit)
- Intuition: derivatives make this a linear problem, so *superposition of paths* works
- Problems:
 - Treats DAG like tree so could do exponential computation in size of DAG. More efficient solution?
 - What order should we compute $\frac{\partial R}{\partial w_0}(w; p_i, q_i)$, $\frac{\partial R}{\partial w_1}(w; p_i, q_i)$ and $\frac{\partial R}{\partial w_2}(w; p_i, q_i)$?

Efficient computation of all derivatives

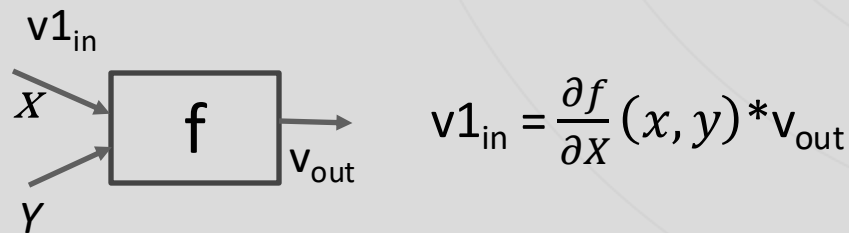
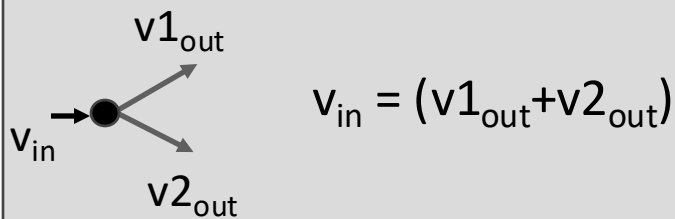
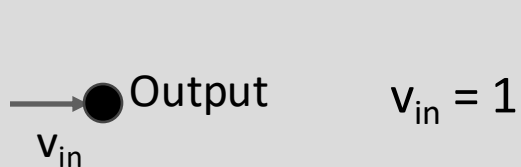
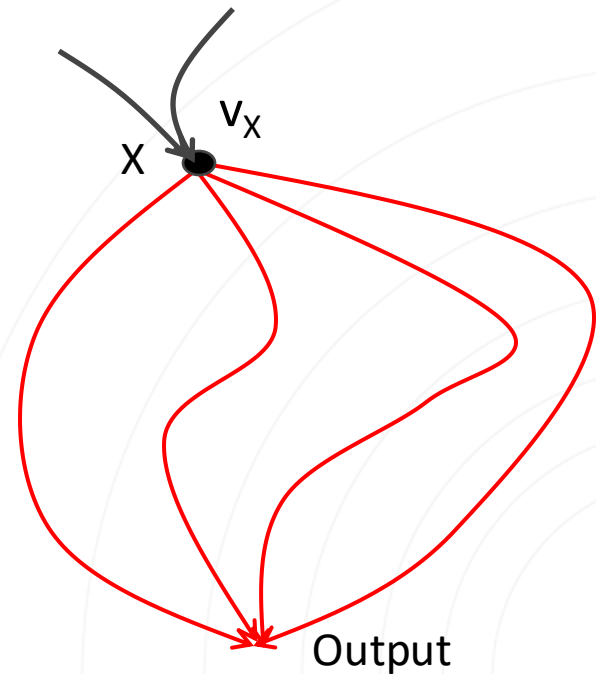


- Compute derivative of output wrt every variable/edge ($\frac{\partial R}{\partial A}, \frac{\partial R}{\partial C}, \frac{\partial R}{\partial F} \dots$)
 - Real number on each edge
 - Derivatives of output wrt weights can be computed from this
- Traverse DAG in reverse topological order of variables for computation
 - $\frac{\partial R}{\partial R} = 1$
 - **Transfer functions** to propagate derivative from function output to its inputs

Summary of derivatives computation

- At each point X , compute $v_X: \Re$ where

$$v_X = \frac{\partial \text{Output}}{\partial X}(w; p_i, q_i) = \sum_{\rho \in \mathcal{P}(X)} \pi(\rho)$$
- Small tweak to handle weight-sharing
- Called back-propagation in ML literature



Back to running example

- Optimization problem: choose W_i values to minimize loss

$$\text{Loss}(w_0, w_1, w_2) = \frac{1}{N} \sum_{i=1}^N (t_i - R(w; p_i, q_i))^2$$

Function $R(P, Q)$ {

$A = W_0 * P$

$B = f_0(A, Q)$

$C = W_1 * B$

$D = W_2 * B$

$E = f_1(C)$

$F = f_2(D)$

$R = f_3(E, F)$

return R }

$$\frac{\partial \text{Loss}}{\partial W_0}(w_0, w_1, w_2) = -\frac{2}{N} \sum_{i=1}^N (t_i - R(w; p_i, q_i)) \left(\frac{\partial R}{\partial W_0}(w; p_i, q_i) \right)$$

Initialize weights to random values

for #epochs do {

GradientVector = 0

 for each training sample (p_i, q_i, t_i) do {

 perform forward propagation and compute

 perform backpropagation and compute weight derivatives

 update GradientVector with products}

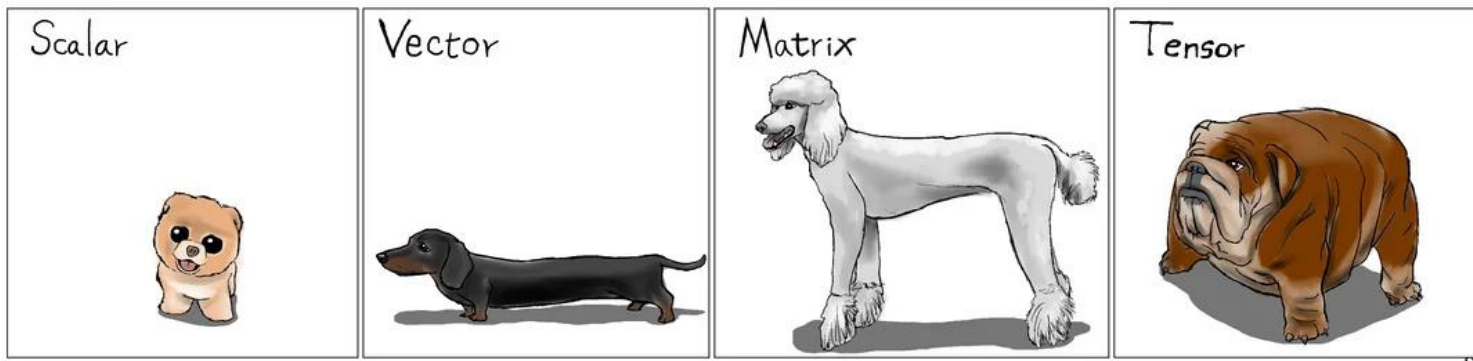
 scale GradientVector by $-2/N$

 use GradientVector to update weights using gradient descent

step

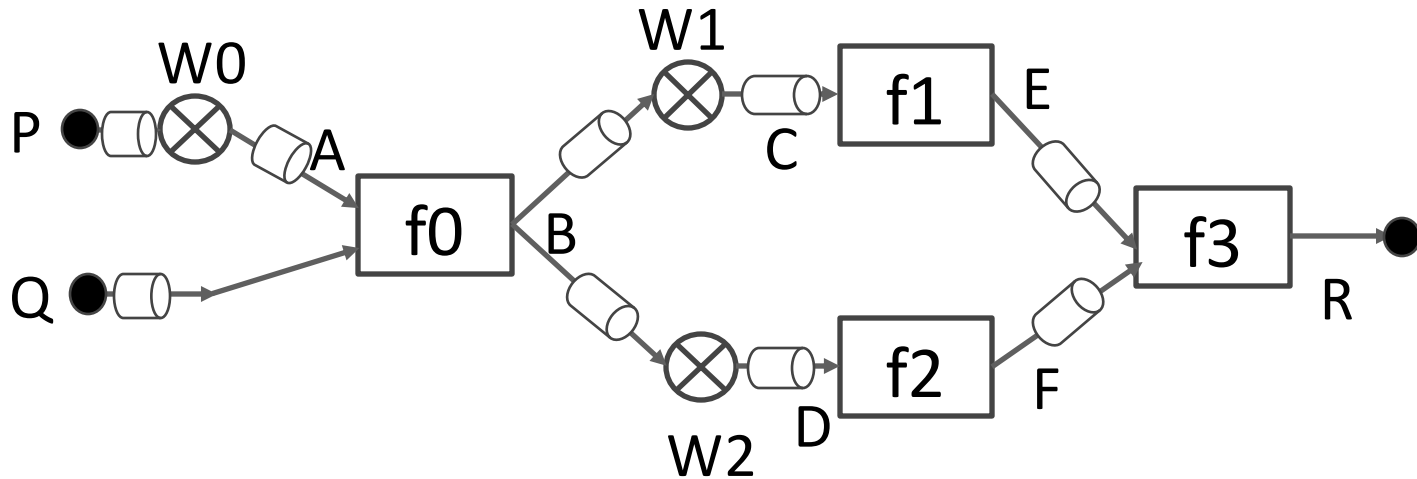
}

Generalization to vectors, matrices, tensors



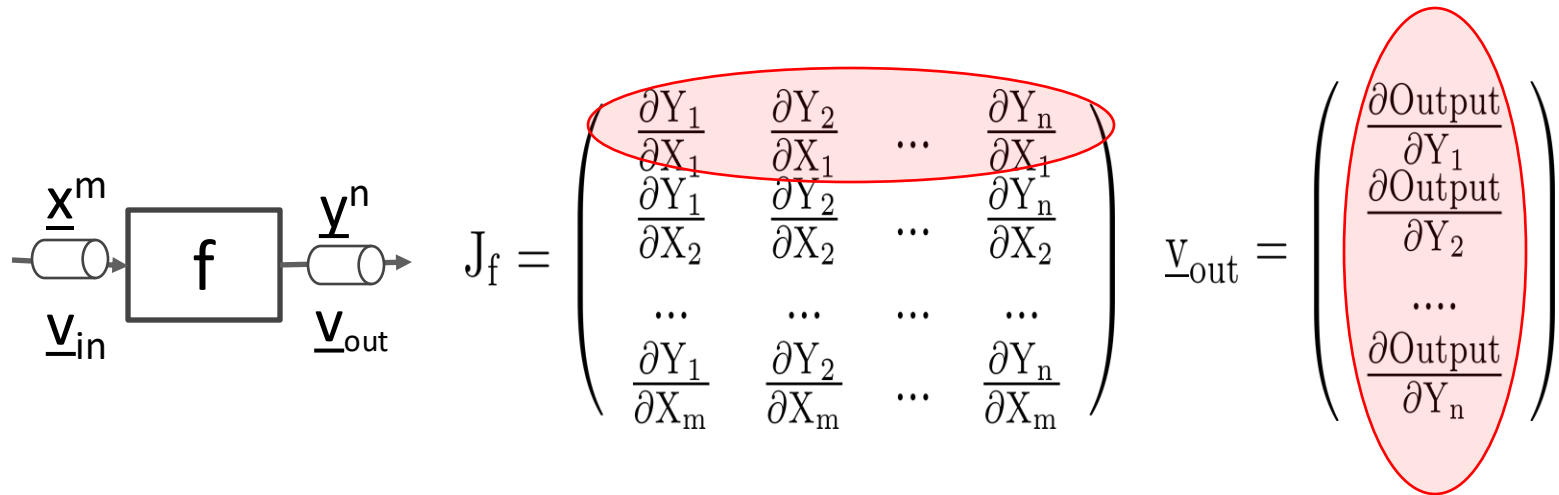
How to understand data with dogs
Karl Stratos ([Reddit post](#))

Generalization to vectors (matrices, tensors are similar)



- **Programs**
 - Inputs can be scalars or **vectors**, but **output is still a scalar** (loss)
 - Weights → Weight matrices
 - > Weight matrix need not be square
 - Function type: vector → vector, vector x vector → vector, etc.
- **Compute gradients instead of derivatives**
 - Variable is vector of size n → gradient of output wrt this vector is also vector of size n
 - > Intuition: each dimension of gradient vector is derivative of output wrt value in that dimension
 - Transfer functions: Jacobians instead of derivatives
- Useful to know [matrix derivatives](#) notation

Transfer functions



- General function f

- $\underline{v}_{in} = J_f * \underline{v}_{out}$

- Linear function W

- $J_f = W^T$

- $\underline{v}_{in} = W^T * \underline{v}_{out}$

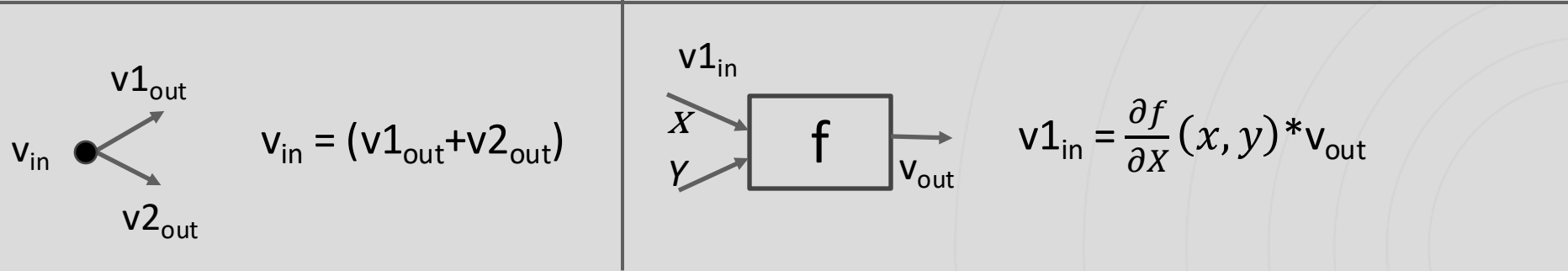
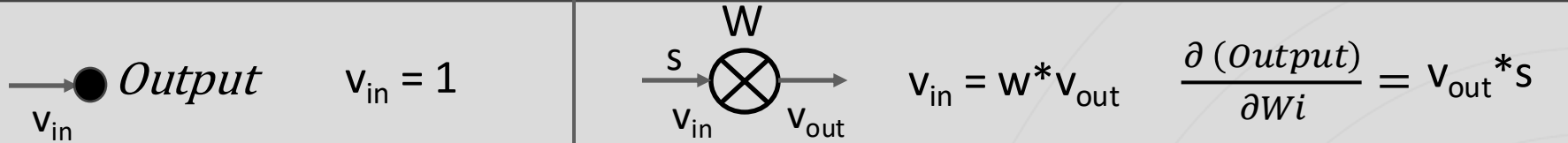
- Derivatives of output wrt weights in W

$$\frac{\partial(\text{Output})}{\partial W} = \underline{v}_{out} \otimes \underline{x} \quad (\otimes \text{ is outer-product})$$

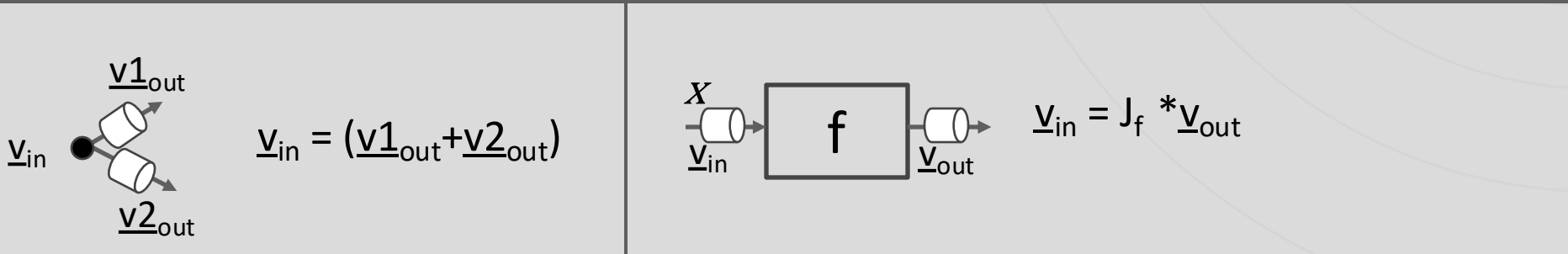
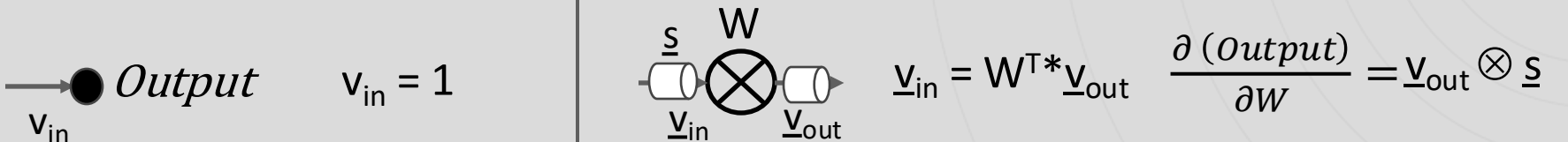
$$\frac{\partial(\text{Output})}{\partial W}(i, j) = \underline{v}_{out}(i) * \underline{x}(j) \quad (\text{If } w(i, j) \text{ changes a small amount, how much does the output change?})$$

Transfer Functions (scalar and vector)

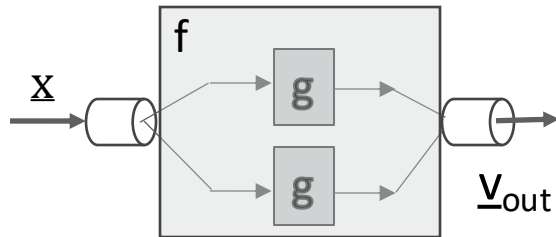
Scalar case



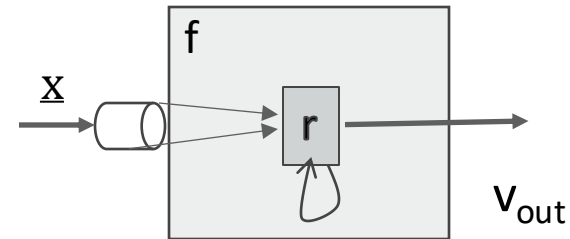
Vector case



Important special cases



map



reduce

- **f: map** ($g: \mathbb{R} \rightarrow \mathbb{R}$)

- J_f = diagonal matrix with $J_f(i,i) = g'(x(i))$

$$\underline{v}_{in} = \begin{pmatrix} g'(\underline{x}(1)) & 0 \\ 0 & g'(\underline{x}(2)) \end{pmatrix} \underline{v}_{out}$$

- **f: reduce** ($r: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$)

- Called **pooling** in ML literature

- $r = +, *, \dots$

> $J_f = (1, 1, \dots, 1)^T$ for +

- $r = \max, \min, \dots$

> If $\max(x) = x(j)$, $J_f = l(j)$ (where $l(j)$ is the indicator vector with 1 in j^{th} position)

- Standard presentations of gradients and back propagation
 - Biological metaphors like neurons and synapses come in the way
 - Properties of activation functions are distraction: enough to know we can compute value and gradient at any point
- Presentation in this lecture
 - **Abstraction** for neural networks: parameterized programs
 - Gradient computation
 - > Abstractly (what?): **sum over paths** (sum of products)
 - > Efficient computation (how?): compositional algorithm on dataflow graph representation
 - Handles complex neural networks with weight-sharing and irregular interconnections (such as “skip connections”) smoothly