

lecture 6: differentiation

deep learning for vision

Yannis Avrithis

Inria Rennes-Bretagne Atlantique

Rennes, Nov. 2017 – Jan. 2018



outline

gradient descent

gradient computation

automatic differentiation: units

automatic differentiation: functions

fun

gradient descent

gradient descent

- a **first-order** Taylor approximation of $f : \mathbb{R}^d \rightarrow \mathbb{R}$ at \mathbf{x}_0 is

$$f_{\mathbf{x}_0}^{(1)}(\mathbf{x}) := f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^\top \nabla f(\mathbf{x}_0)$$

i.e. , the gradient points in the direction of the greatest increase rate

- a **second-order** approximation needs the Hessian matrix Hf

$$f_{\mathbf{x}_0}^{(2)}(\mathbf{x}) := f_{\mathbf{x}_0}^{(1)}(\mathbf{x}) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^\top (Hf(\mathbf{x}_0))(\mathbf{x} - \mathbf{x}_0)$$

- assuming f is locally convex with isotropic $Hf(\mathbf{x}_0) = \frac{1}{\epsilon}I$, the gradient of $f_{\mathbf{x}_0}^{(2)}$ is

$$\nabla f_{\mathbf{x}_0}^{(2)}(\mathbf{x}) = \nabla f(\mathbf{x}_0) + \frac{1}{\epsilon}(\mathbf{x} - \mathbf{x}_0)$$

- so if we were to minimize this approximation instead of f ,

$$\arg \min_{\mathbf{x}} f_{\mathbf{x}_0}^{(2)}(\mathbf{x}) = \mathbf{x}_0 - \epsilon \nabla f(\mathbf{x}_0)$$

gradient descent

- a **first-order** Taylor approximation of $f : \mathbb{R}^d \rightarrow \mathbb{R}$ at \mathbf{x}_0 is

$$f_{\mathbf{x}_0}^{(1)}(\mathbf{x}) := f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^\top \nabla f(\mathbf{x}_0)$$

i.e. , the gradient points in the direction of the greatest increase rate

- a **second-order** approximation needs the Hessian matrix Hf

$$f_{\mathbf{x}_0}^{(2)}(\mathbf{x}) := f_{\mathbf{x}_0}^{(1)}(\mathbf{x}) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^\top (Hf(\mathbf{x}_0))(\mathbf{x} - \mathbf{x}_0)$$

- assuming f is locally convex with isotropic $Hf(\mathbf{x}_0) = \frac{1}{\epsilon}I$, the gradient of $f_{\mathbf{x}_0}^{(2)}$ is

$$\nabla f_{\mathbf{x}_0}^{(2)}(\mathbf{x}) = \nabla f(\mathbf{x}_0) + \frac{1}{\epsilon}(\mathbf{x} - \mathbf{x}_0)$$

- so if we were to minimize this approximation instead of f ,

$$\arg \min_{\mathbf{x}} f_{\mathbf{x}_0}^{(2)}(\mathbf{x}) = \mathbf{x}_0 - \epsilon \nabla f(\mathbf{x}_0)$$

gradient descent

- a **first-order** Taylor approximation of $f : \mathbb{R}^d \rightarrow \mathbb{R}$ at \mathbf{x}_0 is

$$f_{\mathbf{x}_0}^{(1)}(\mathbf{x}) := f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^\top \nabla f(\mathbf{x}_0)$$

i.e. , the gradient points in the direction of the greatest increase rate

- a **second-order** approximation needs the Hessian matrix Hf

$$f_{\mathbf{x}_0}^{(2)}(\mathbf{x}) := f_{\mathbf{x}_0}^{(1)}(\mathbf{x}) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^\top (Hf(\mathbf{x}_0))(\mathbf{x} - \mathbf{x}_0)$$

- assuming f is locally convex with isotropic $Hf(\mathbf{x}_0) = \frac{1}{\epsilon}I$, the gradient of $f_{\mathbf{x}_0}^{(2)}$ is

$$\nabla f_{\mathbf{x}_0}^{(2)}(\mathbf{x}) = \nabla f(\mathbf{x}_0) + \frac{1}{\epsilon}(\mathbf{x} - \mathbf{x}_0)$$

- so if we were to minimize this approximation instead of f ,

$$\arg \min_{\mathbf{x}} f_{\mathbf{x}_0}^{(2)}(\mathbf{x}) = \mathbf{x}_0 - \epsilon \nabla f(\mathbf{x}_0)$$

gradient descent

- this yields the update rule

$$\mathbf{x}^{(\tau+1)} = \mathbf{x}^{(\tau)} - \epsilon \nabla f(\mathbf{x}^{(\tau)})$$

i.e. , we are moving in the direction of the greatest decrease rate such that locally (depending on ϵ)

$$\begin{aligned}f_{\mathbf{x}^{(\tau)}}^{(1)}(\mathbf{x}^{(\tau+1)}) &= f(\mathbf{x}^{(\tau)}) + (\mathbf{x}^{(\tau+1)} - \mathbf{x}^{(\tau)})^\top \nabla f(\mathbf{x}_0) \\&= f(\mathbf{x}^{(\tau)}) - \epsilon \nabla f(\mathbf{x}_0)^\top \nabla f(\mathbf{x}_0) \\&\leq f(\mathbf{x}^{(\tau)})\end{aligned}$$

- the step size ϵ is inversely proportional to the curvature we assume for f at the local minimum

gradient descent

- this yields the update rule

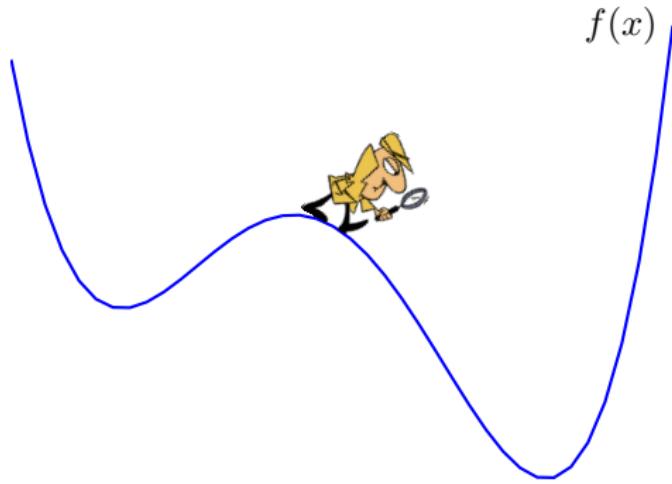
$$\mathbf{x}^{(\tau+1)} = \mathbf{x}^{(\tau)} - \epsilon \nabla f(\mathbf{x}^{(\tau)})$$

i.e. , we are moving in the direction of the greatest decrease rate such that locally (depending on ϵ)

$$\begin{aligned}f_{\mathbf{x}^{(\tau)}}^{(1)}(\mathbf{x}^{(\tau+1)}) &= f(\mathbf{x}^{(\tau)}) + (\mathbf{x}^{(\tau+1)} - \mathbf{x}^{(\tau)})^\top \nabla f(\mathbf{x}_0) \\&= f(\mathbf{x}^{(\tau)}) - \epsilon \nabla f(\mathbf{x}_0)^\top \nabla f(\mathbf{x}_0) \\&\leq f(\mathbf{x}^{(\tau)})\end{aligned}$$

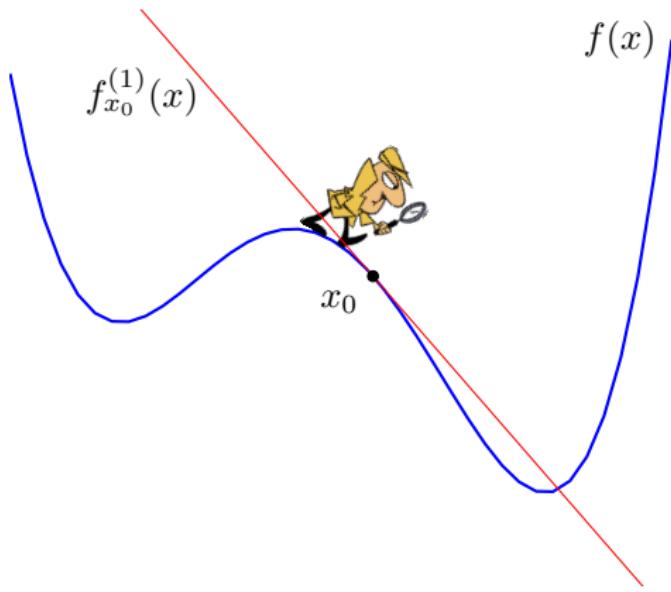
- the step size ϵ is inversely proportional to the curvature we assume for f at the local minimum

gradient descent in one dimension



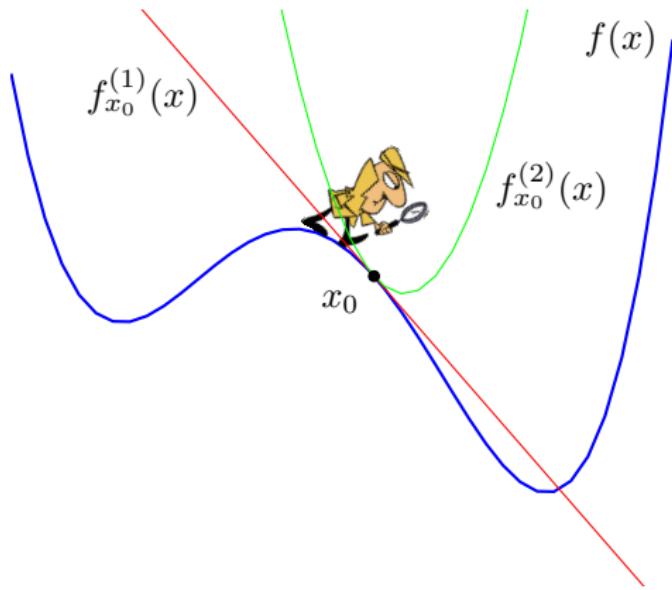
- $\epsilon = 0.05$: converges to local minimum

gradient descent in one dimension



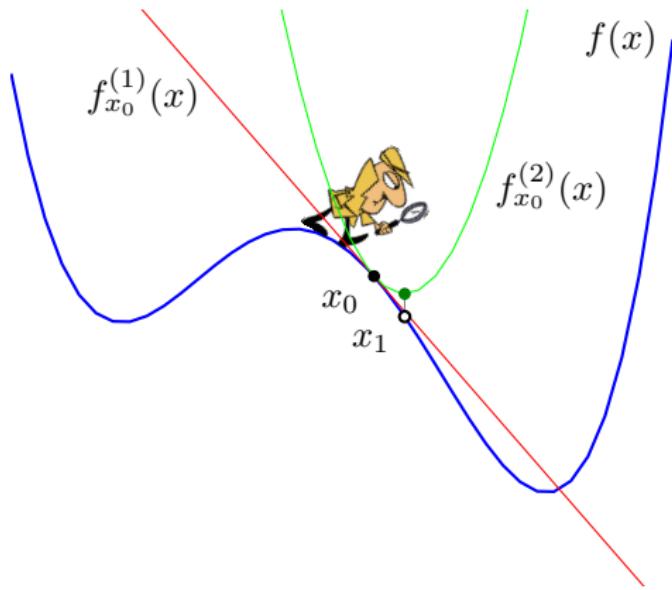
- $\epsilon = 0.05$: converges to local minimum

gradient descent in one dimension



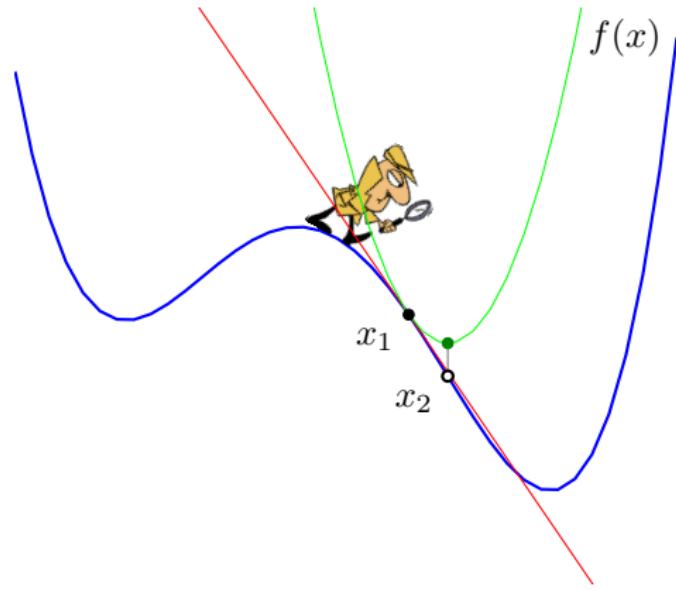
- $\epsilon = 0.05$: converges to local minimum

gradient descent in one dimension



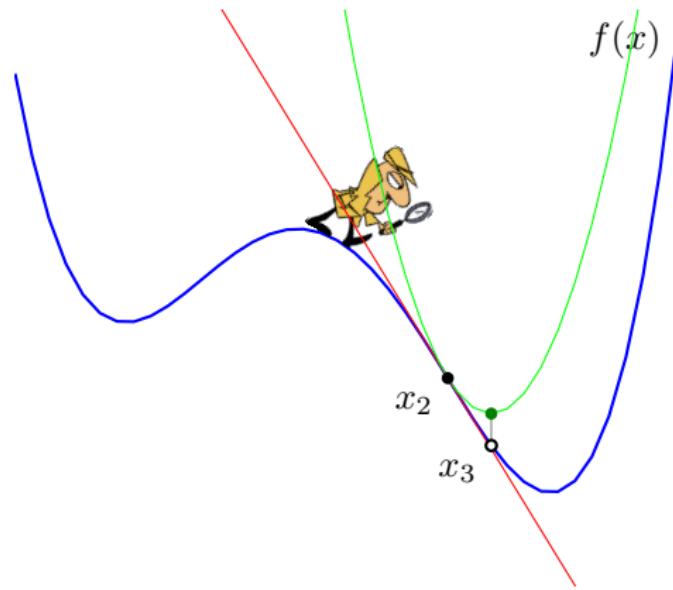
- $\epsilon = 0.05$: converges to local minimum

gradient descent in one dimension



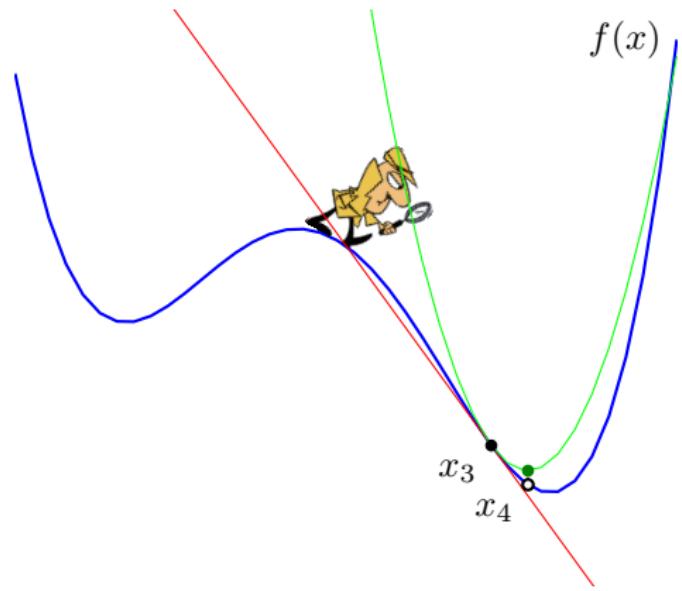
- $\epsilon = 0.05$: converges to local minimum

gradient descent in one dimension



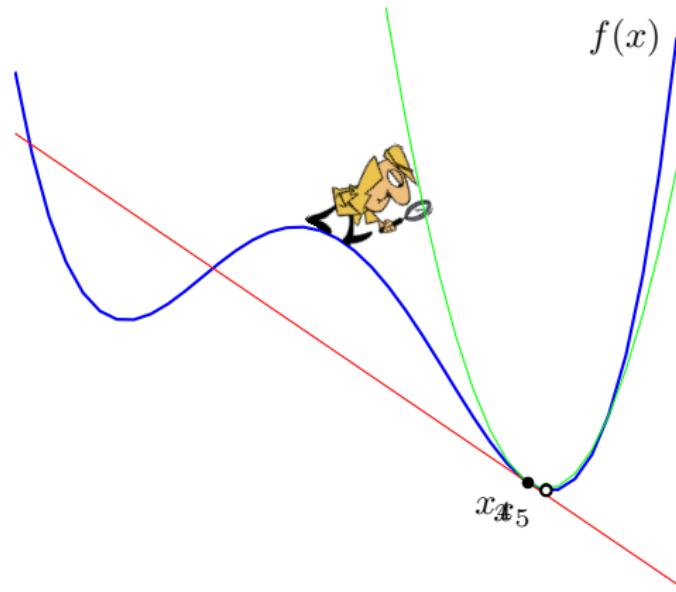
- $\epsilon = 0.05$: converges to local minimum

gradient descent in one dimension



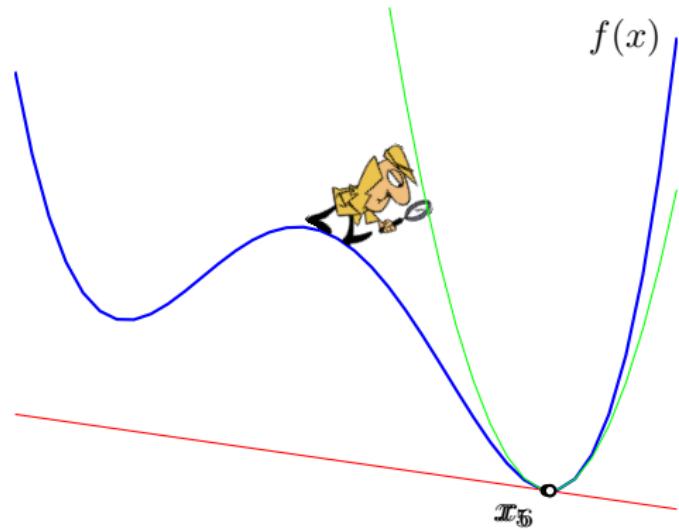
- $\epsilon = 0.05$: converges to local minimum

gradient descent in one dimension



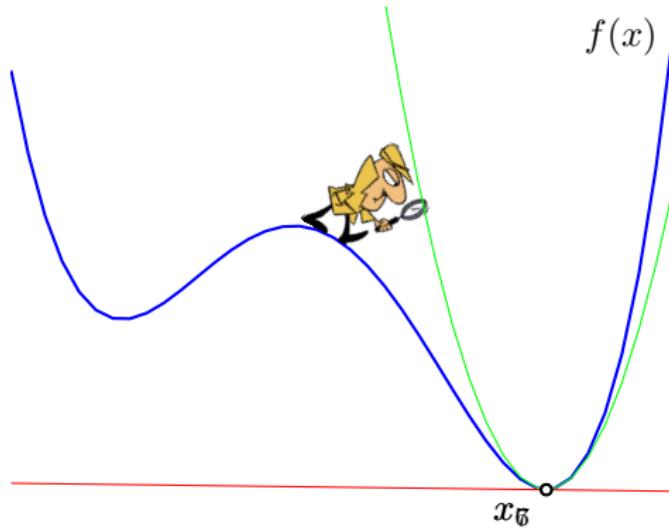
- $\epsilon = 0.05$: converges to local minimum

gradient descent in one dimension



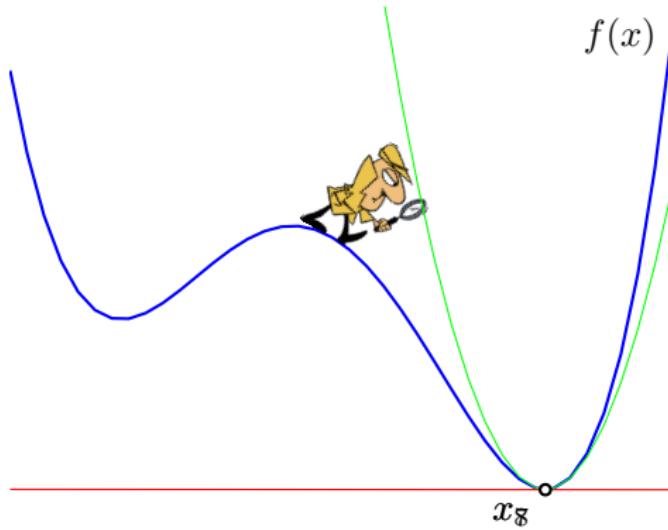
- $\epsilon = 0.05$: converges to local minimum

gradient descent in one dimension



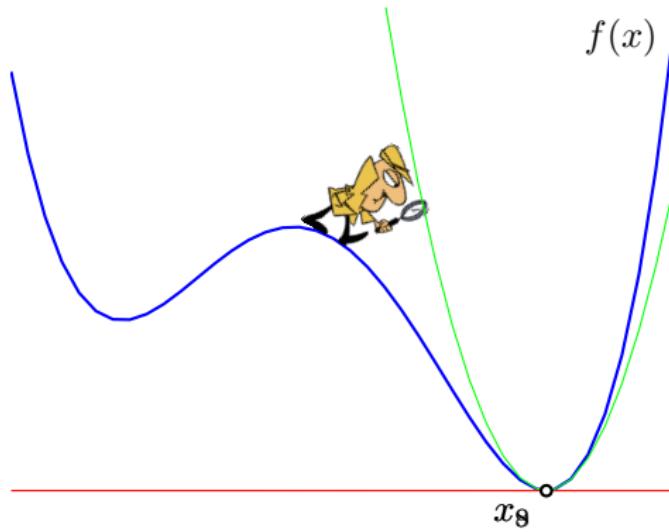
- $\epsilon = 0.05$: converges to local minimum

gradient descent in one dimension



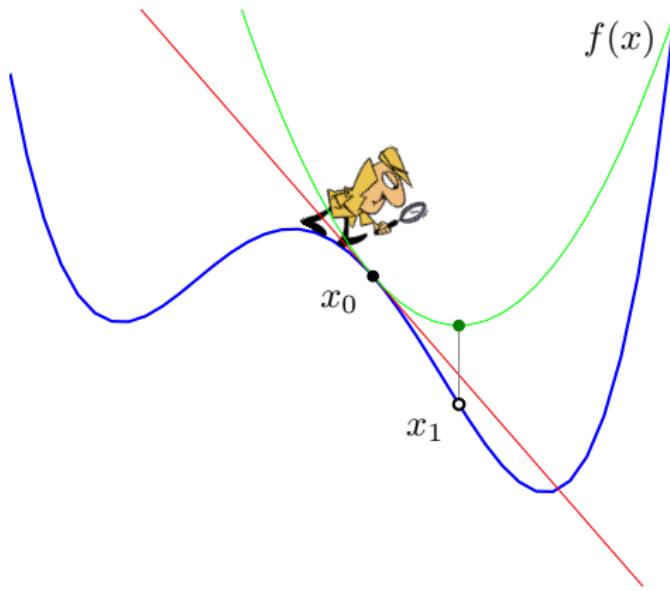
- $\epsilon = 0.05$: converges to local minimum

gradient descent in one dimension



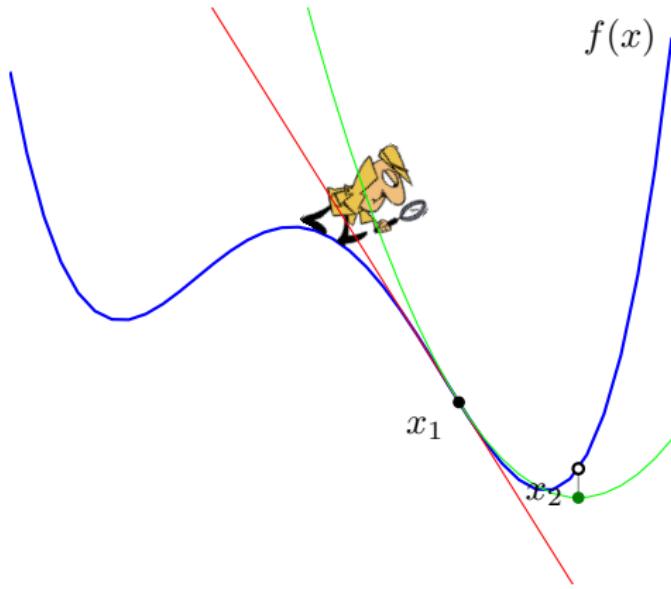
- $\epsilon = 0.05$: converges to local minimum

gradient descent in one dimension



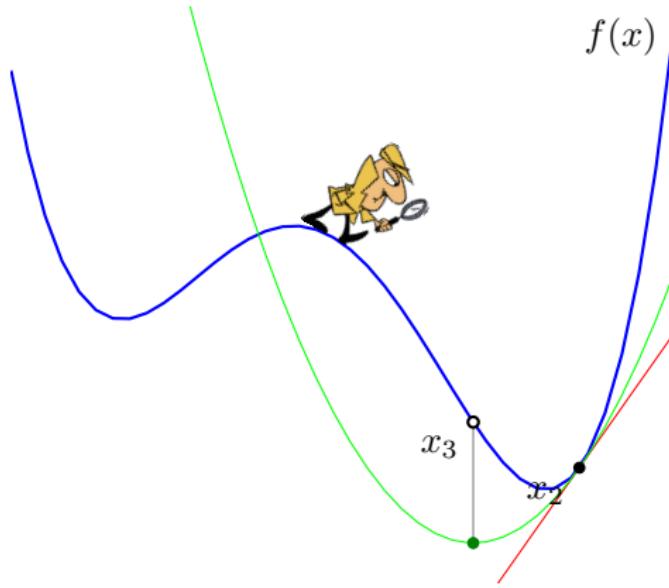
- $\epsilon = 0.14$: $1/\epsilon$ less than actual curvature, does not converge

gradient descent in one dimension



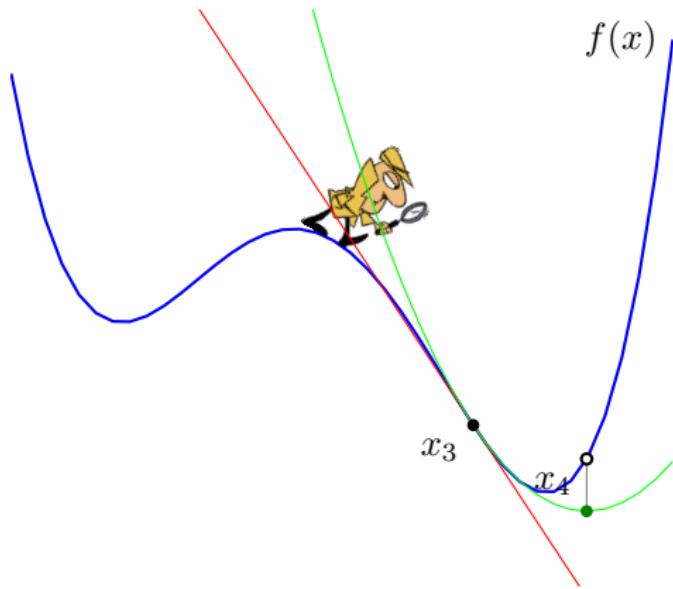
- $\epsilon = 0.14$: $1/\epsilon$ less than actual curvature, does not converge

gradient descent in one dimension



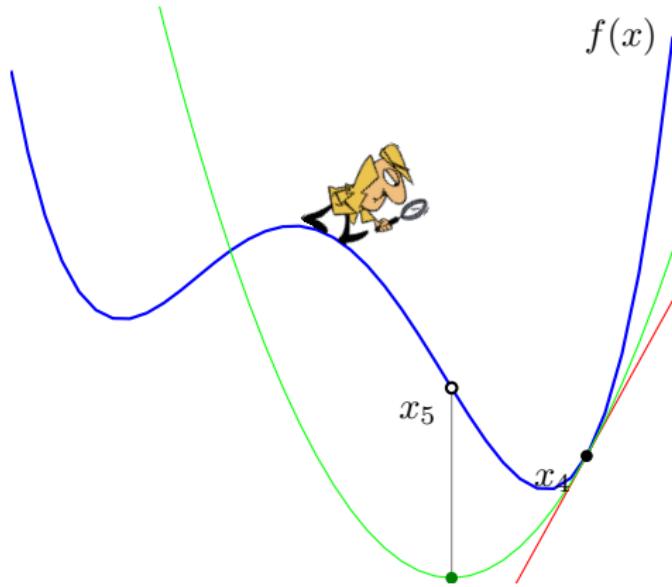
- $\epsilon = 0.14$: $1/\epsilon$ less than actual curvature, does not converge

gradient descent in one dimension



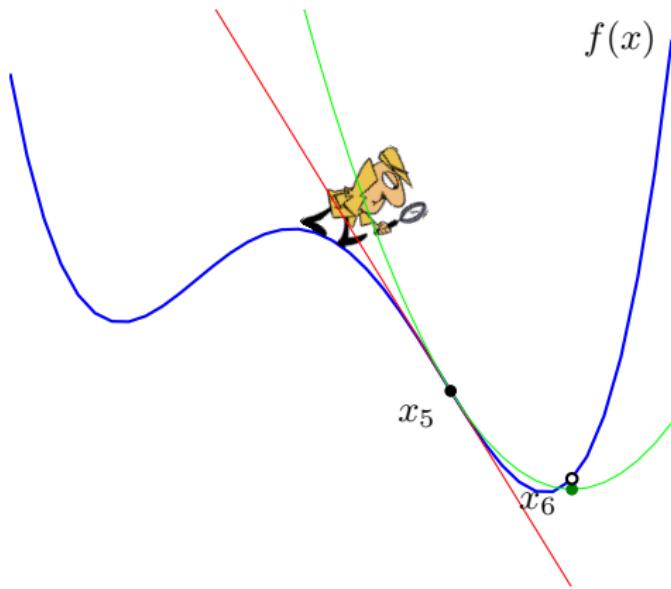
- $\epsilon = 0.14$: $1/\epsilon$ less than actual curvature, does not converge

gradient descent in one dimension



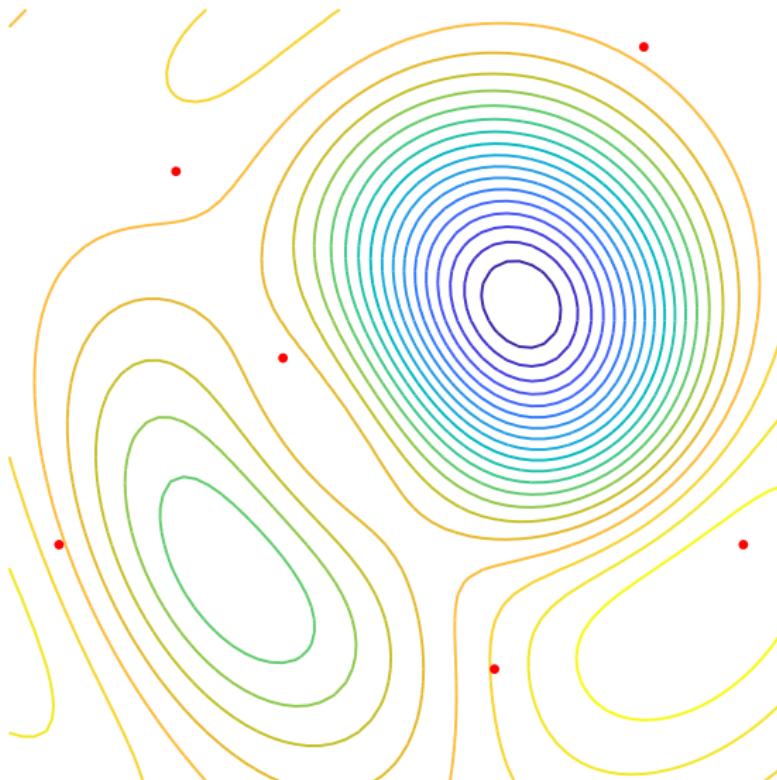
- $\epsilon = 0.14$: $1/\epsilon$ less than actual curvature, does not converge

gradient descent in one dimension



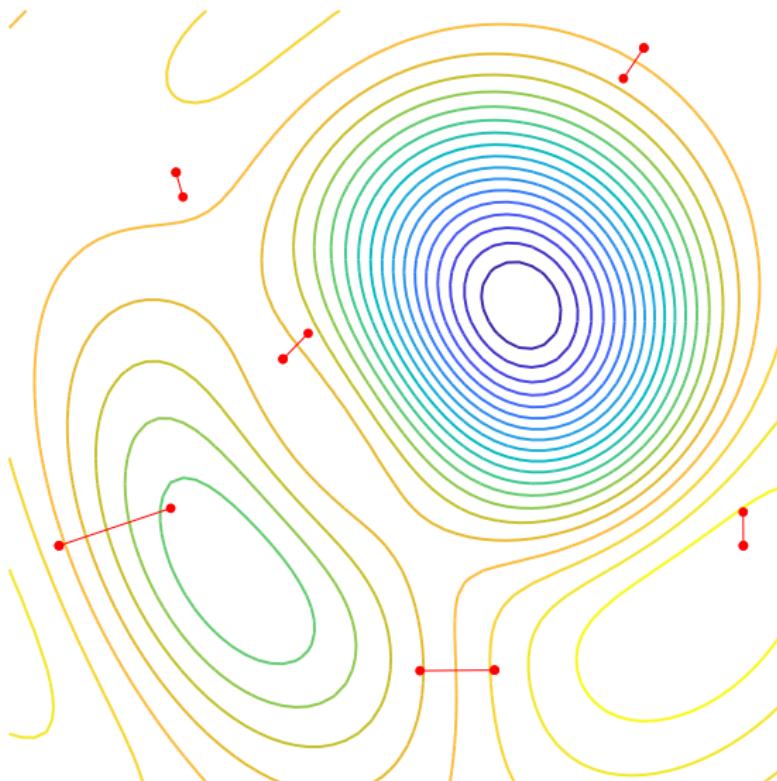
- $\epsilon = 0.14$: $1/\epsilon$ less than actual curvature, does not converge

gradient descent in two dimensions



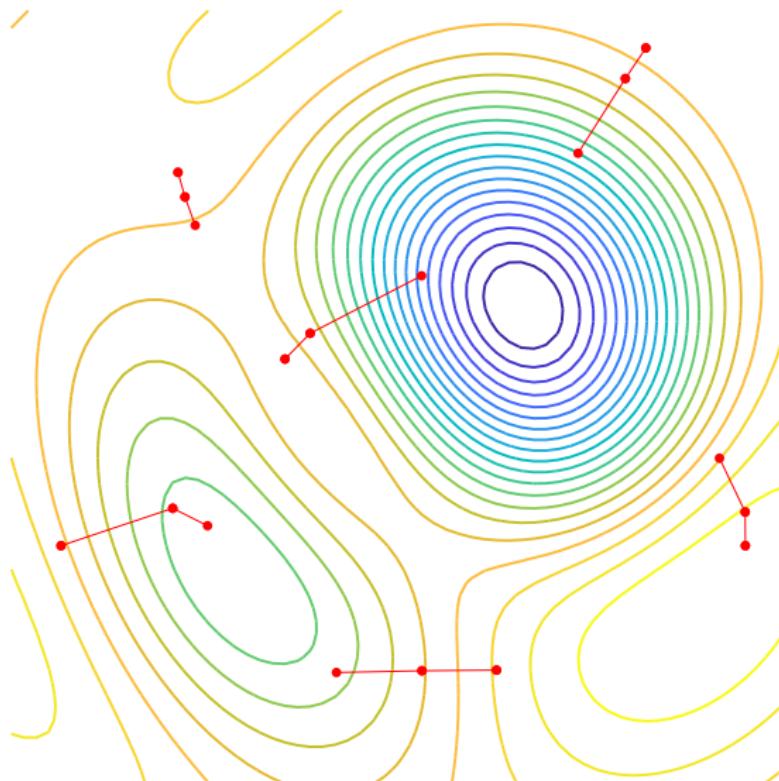
$\epsilon = 0.14$, iteration 0

gradient descent in two dimensions



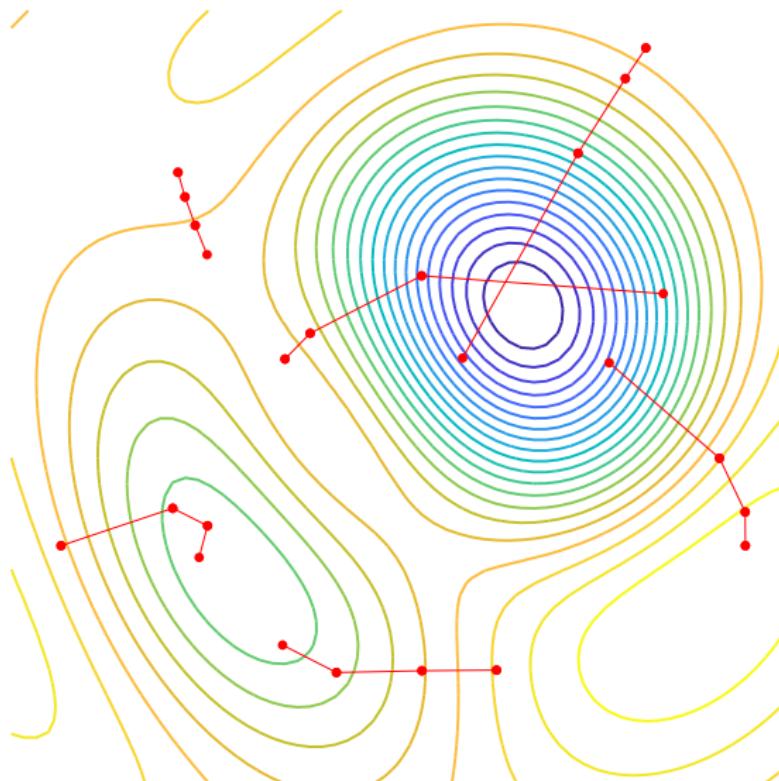
$\epsilon = 0.14$, iteration 1

gradient descent in two dimensions



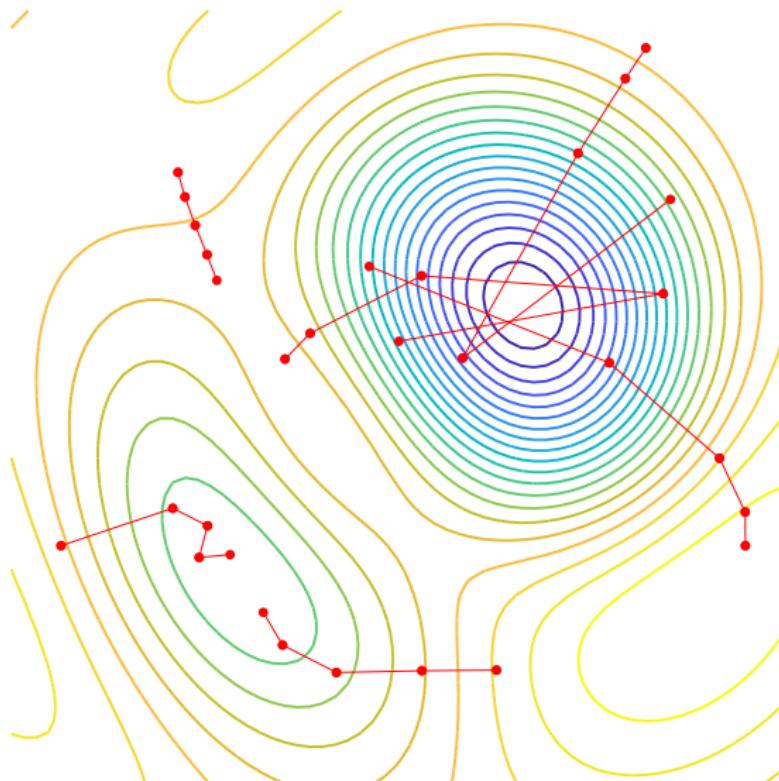
$\epsilon = 0.14$, iteration 2

gradient descent in two dimensions



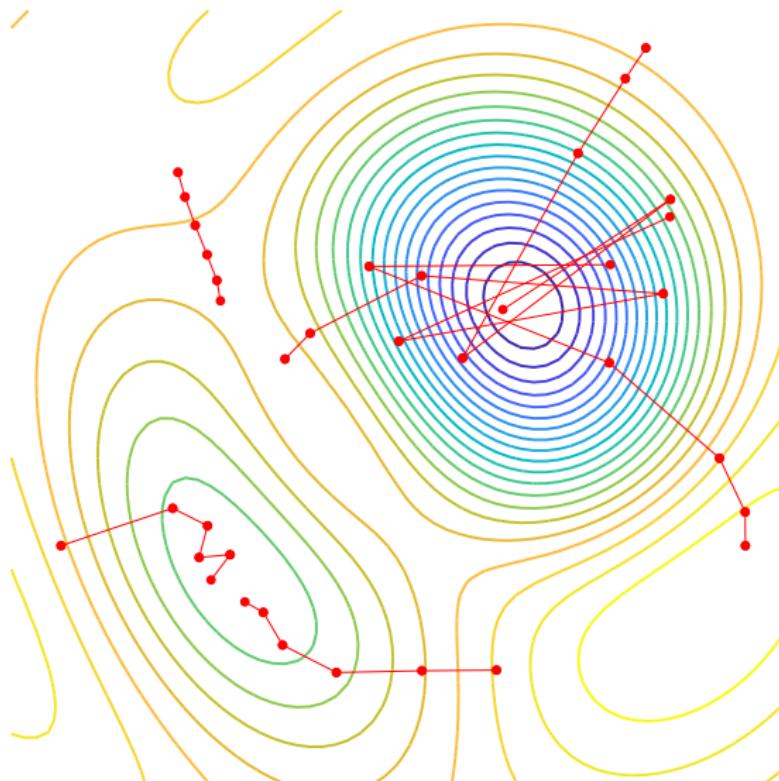
$\epsilon = 0.14$, iteration 3

gradient descent in two dimensions

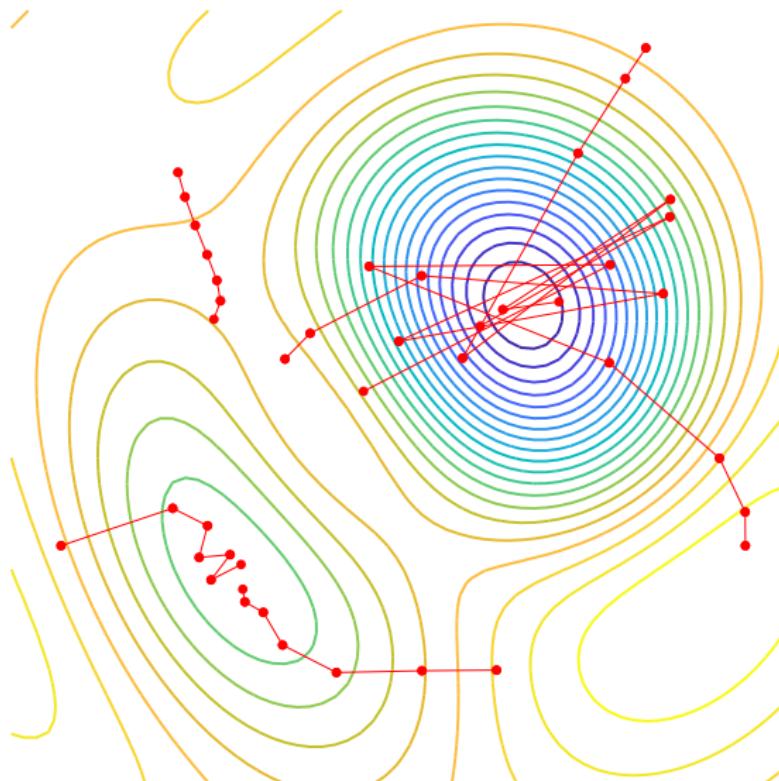


$\epsilon = 0.14$, iteration 4

gradient descent in two dimensions

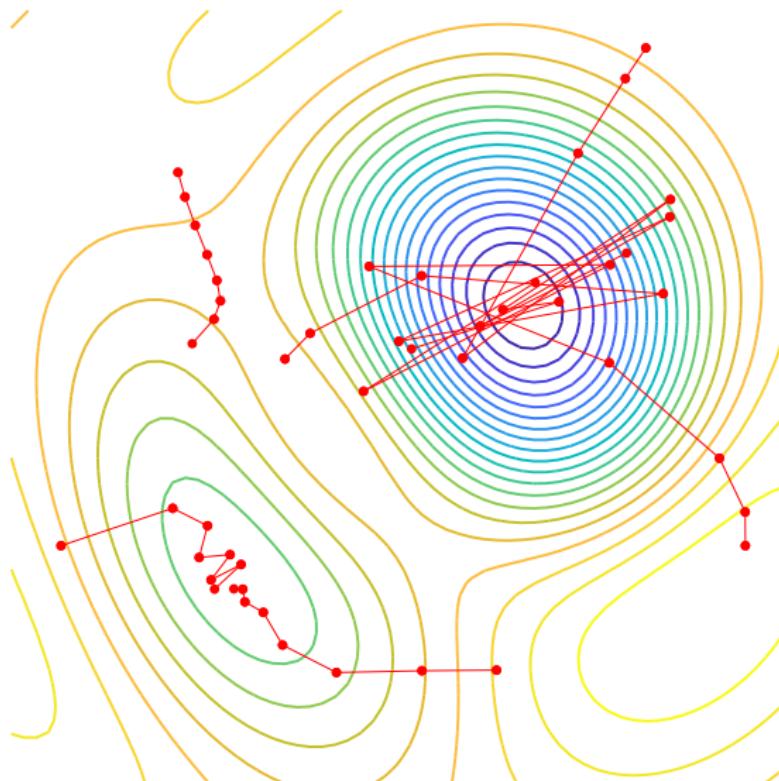


gradient descent in two dimensions



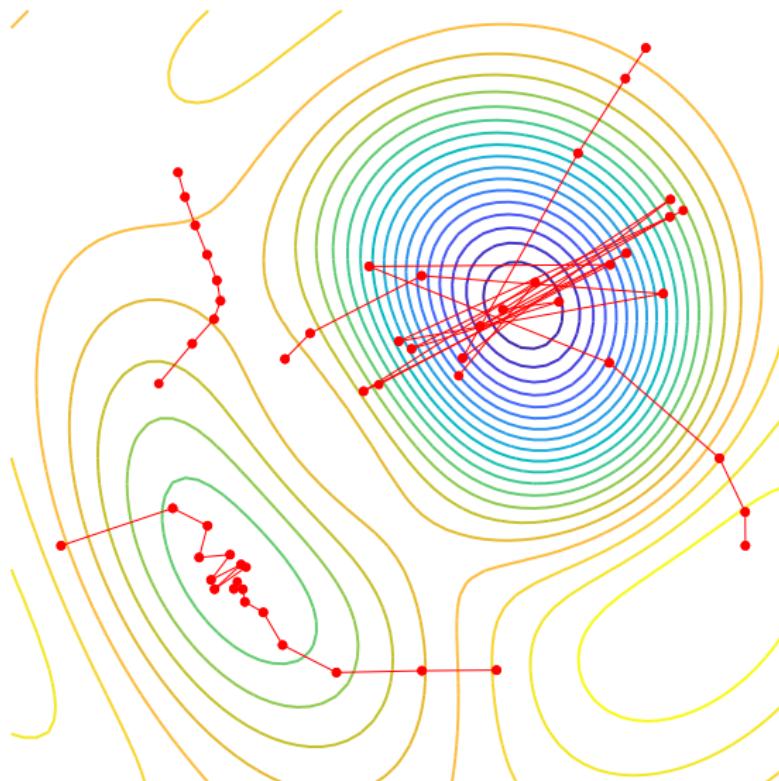
$\epsilon = 0.14$, iteration 6

gradient descent in two dimensions



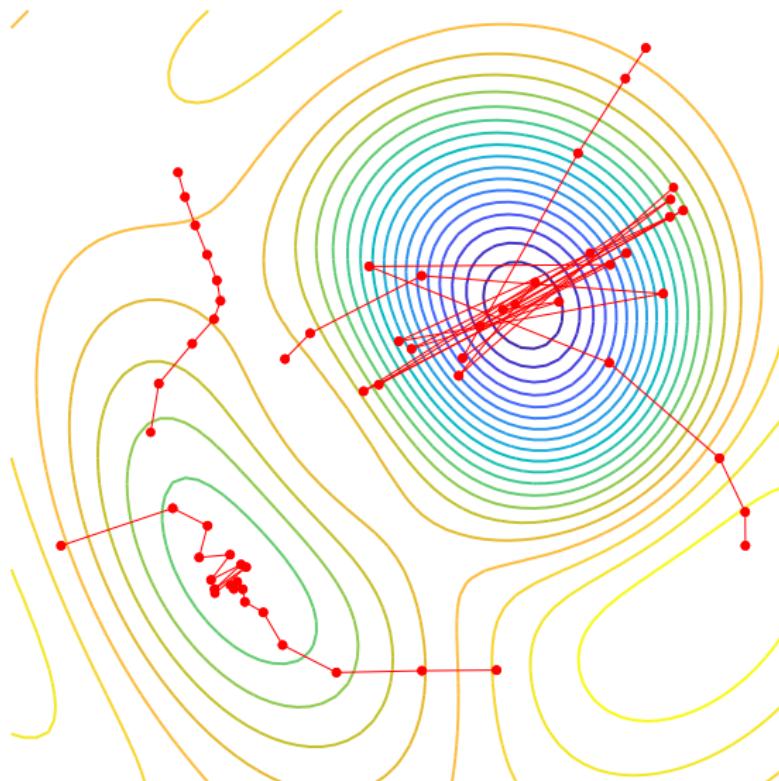
$\epsilon = 0.14$, iteration 7

gradient descent in two dimensions



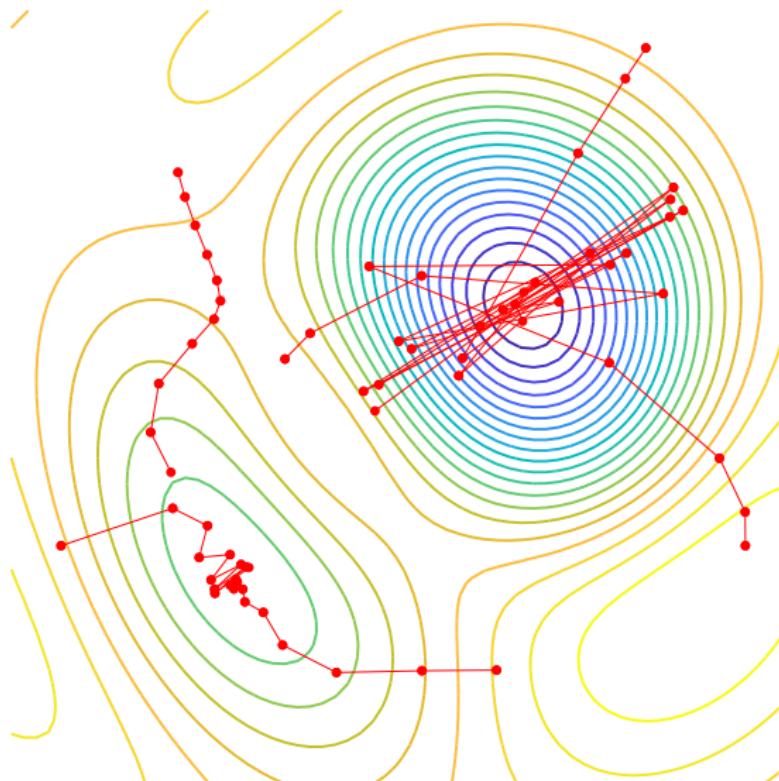
$\epsilon = 0.14$, iteration 8

gradient descent in two dimensions



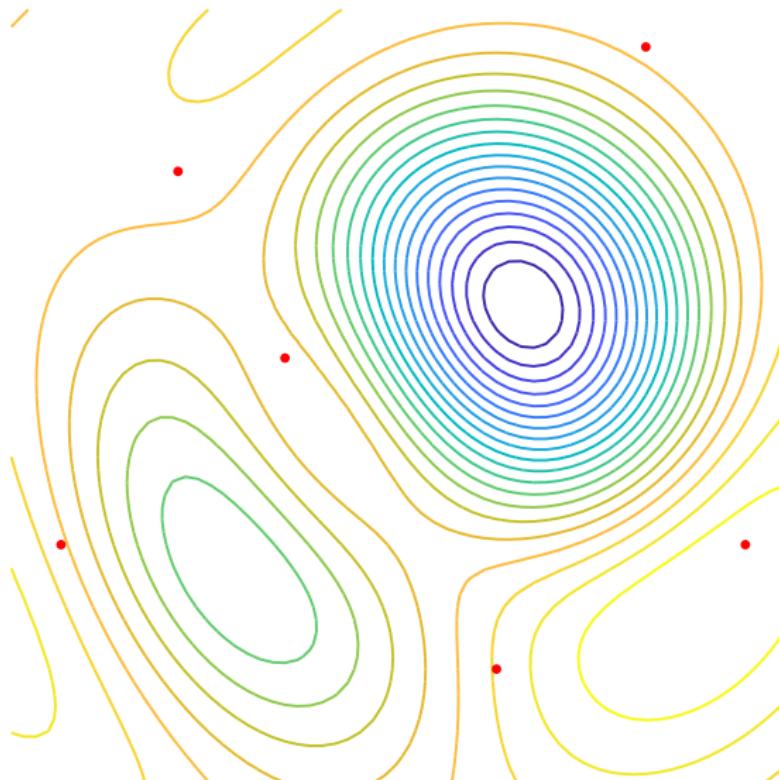
$\epsilon = 0.14$, iteration 9

gradient descent in two dimensions



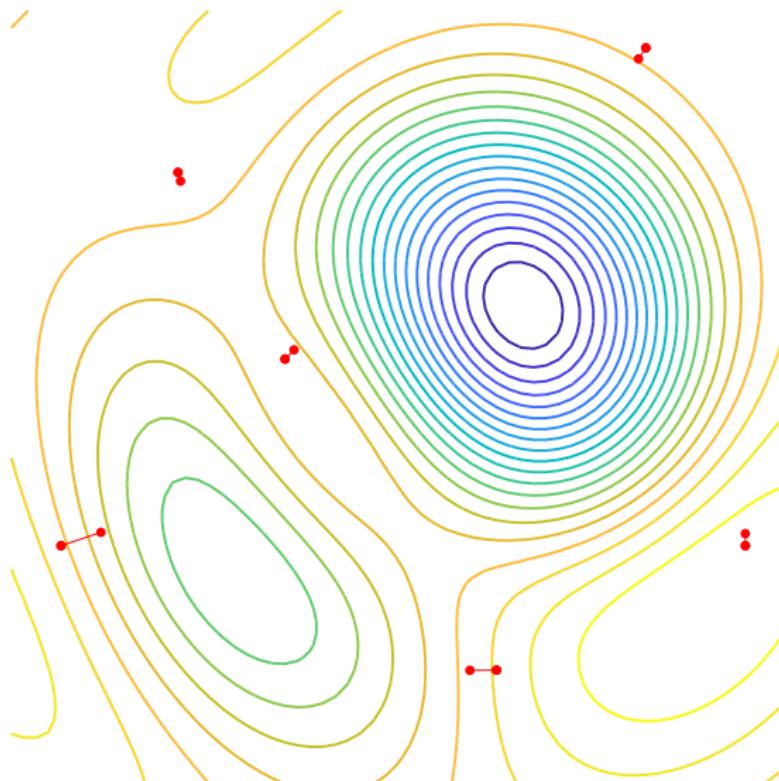
$\epsilon = 0.14$, iteration 10

gradient descent in two dimensions



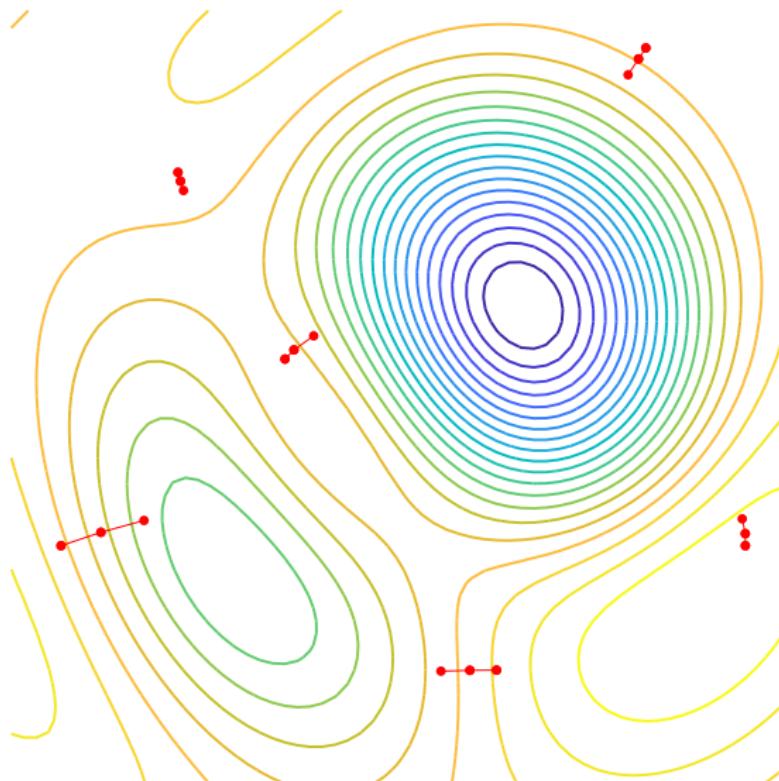
$\epsilon = 0.05$, iteration 0

gradient descent in two dimensions



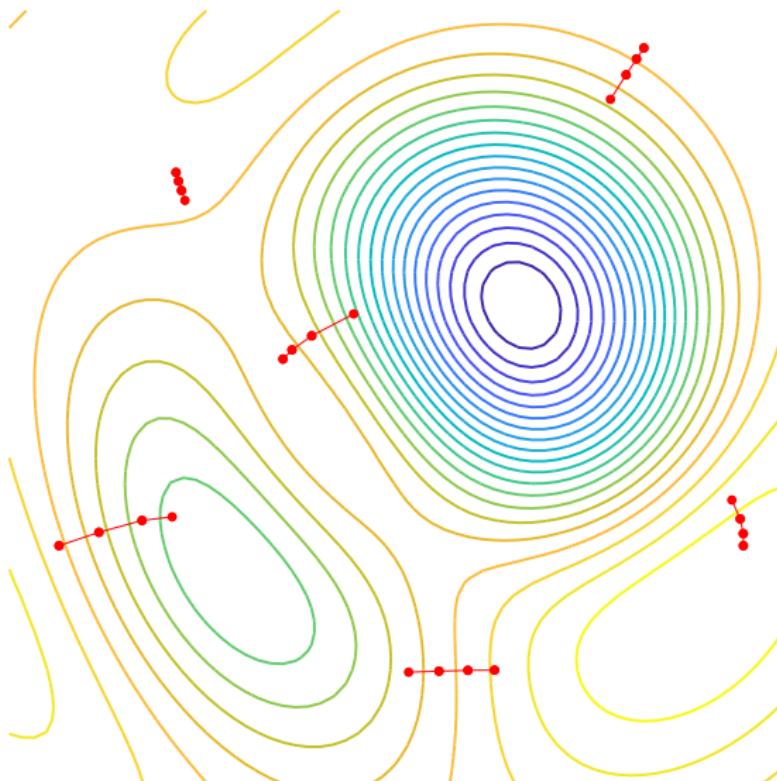
$\epsilon = 0.05$, iteration 1

gradient descent in two dimensions



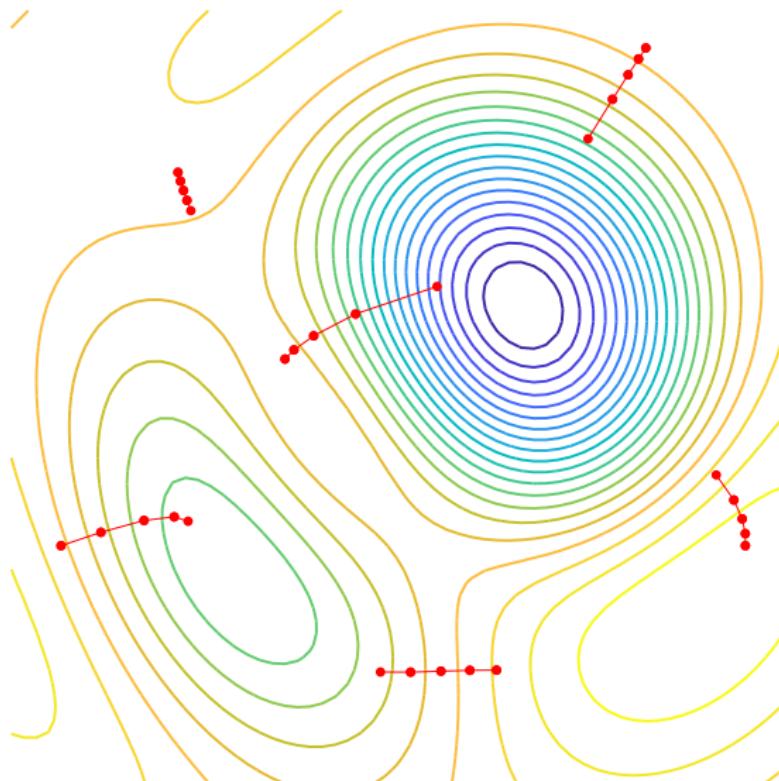
$\epsilon = 0.05$, iteration 2

gradient descent in two dimensions

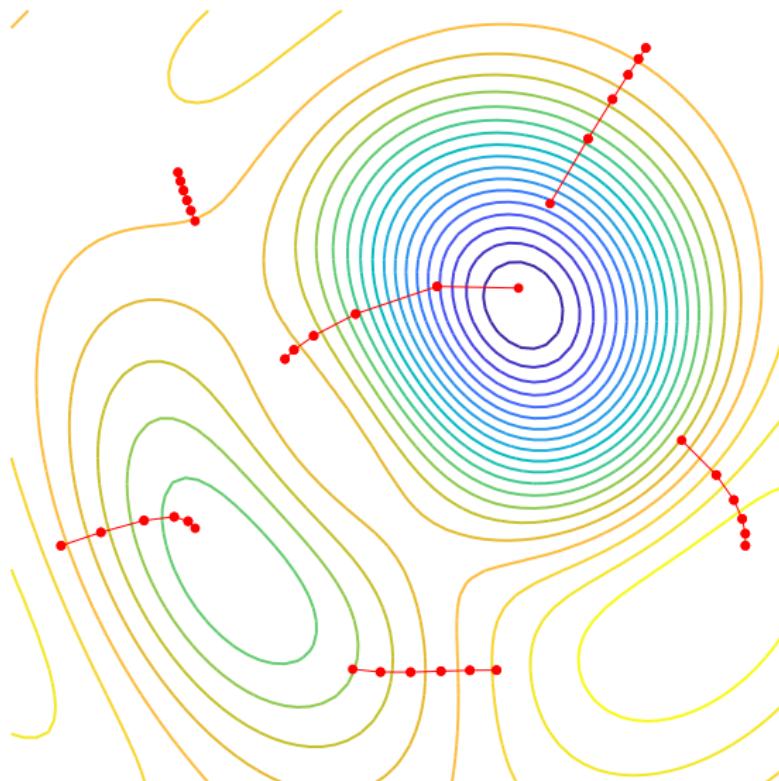


$\epsilon = 0.05$, iteration 3

gradient descent in two dimensions

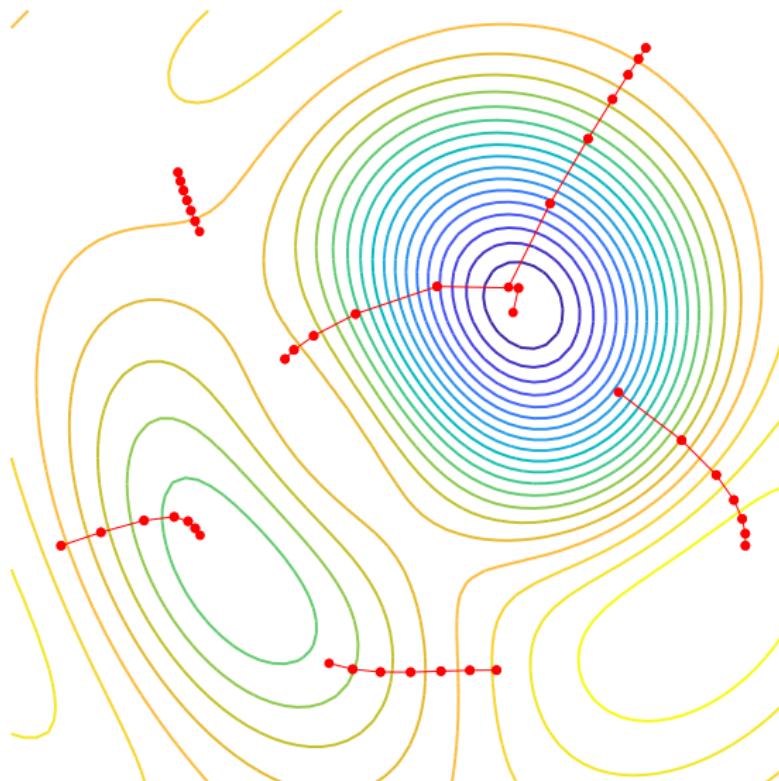


gradient descent in two dimensions



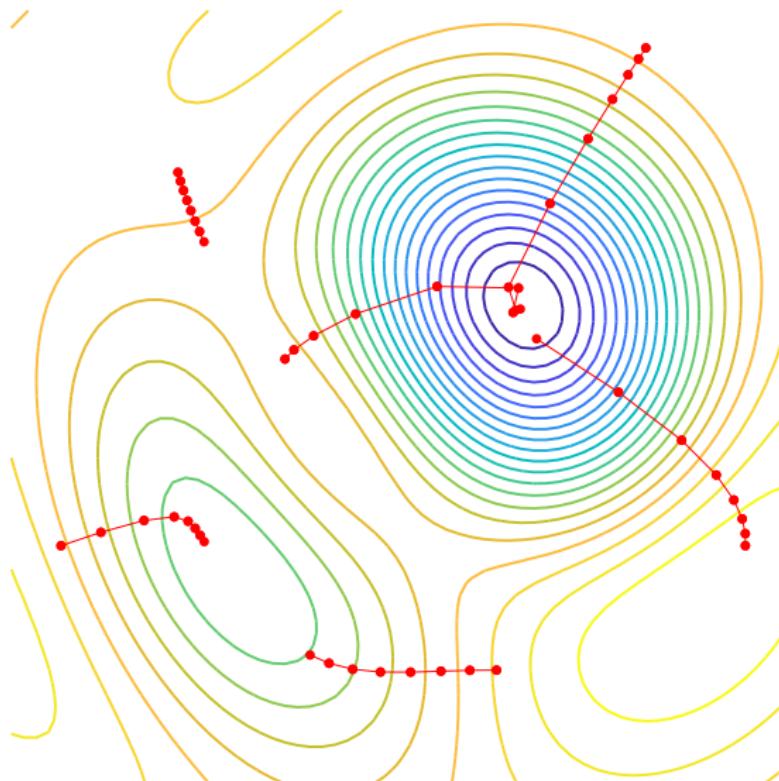
$\epsilon = 0.05$, iteration 5

gradient descent in two dimensions



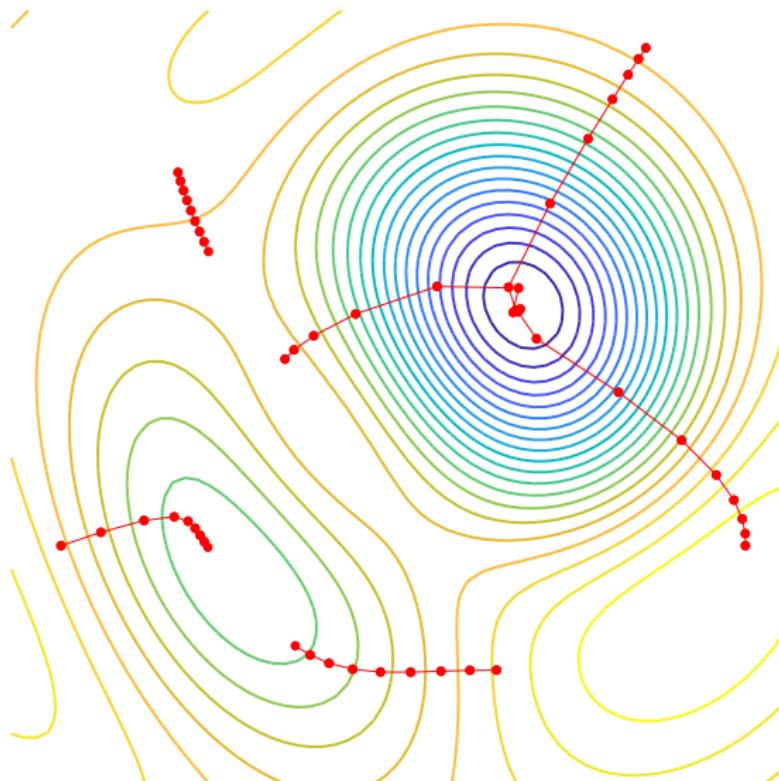
$\epsilon = 0.05$, iteration 6

gradient descent in two dimensions



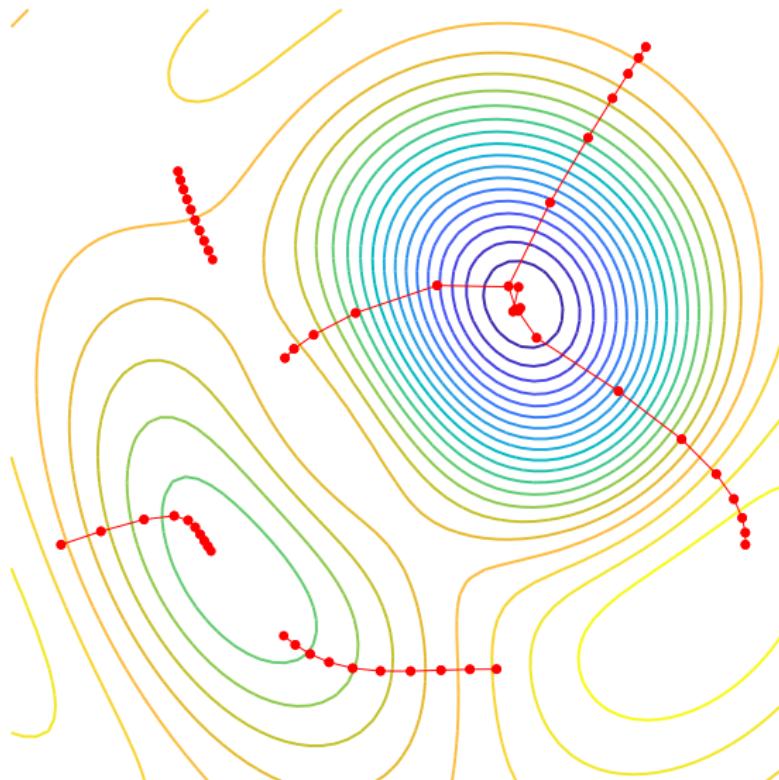
$\epsilon = 0.05$, iteration 7

gradient descent in two dimensions



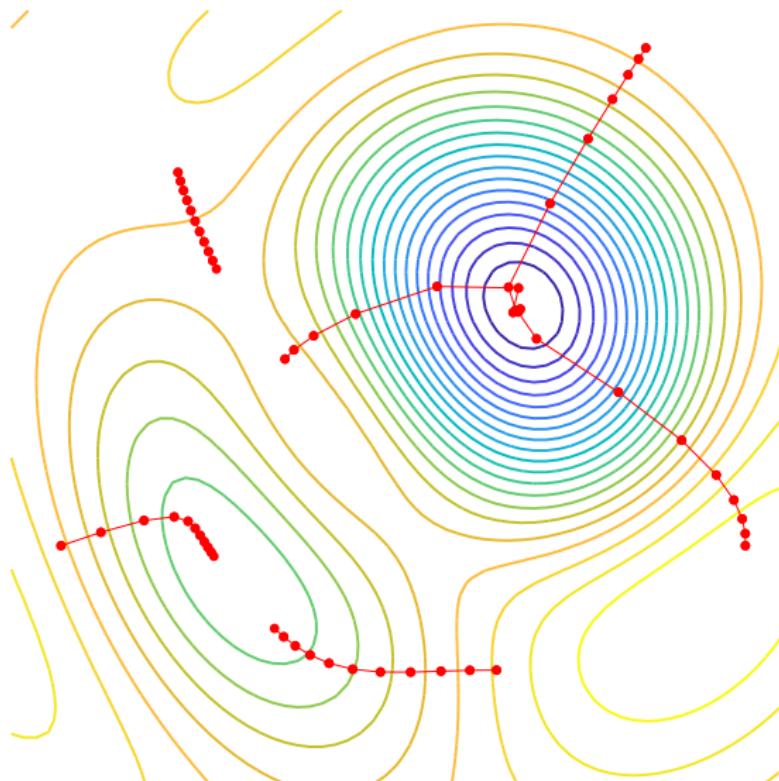
$\epsilon = 0.05$, iteration 8

gradient descent in two dimensions



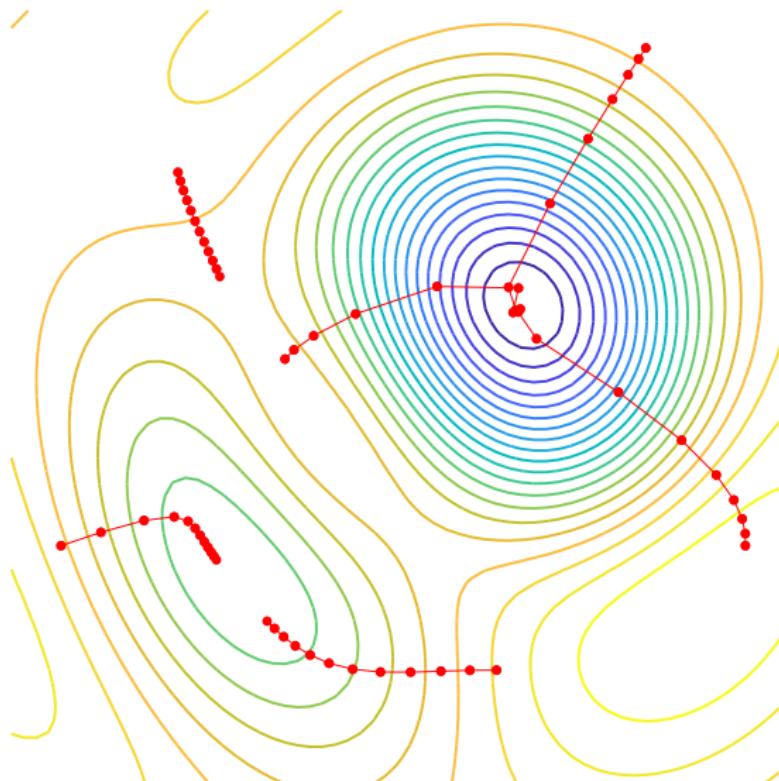
$\epsilon = 0.05$, iteration 9

gradient descent in two dimensions



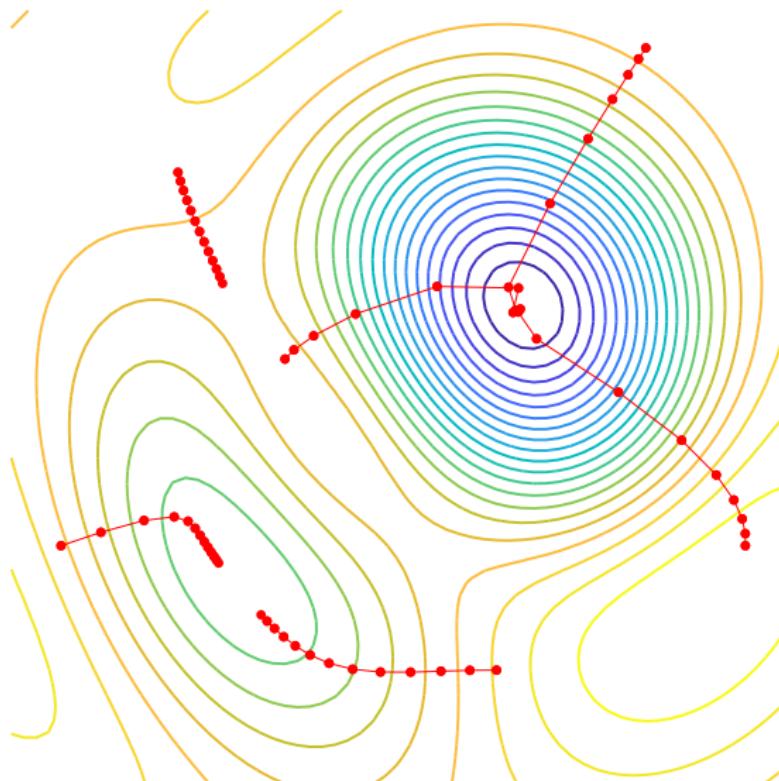
$\epsilon = 0.05$, iteration 10

gradient descent in two dimensions



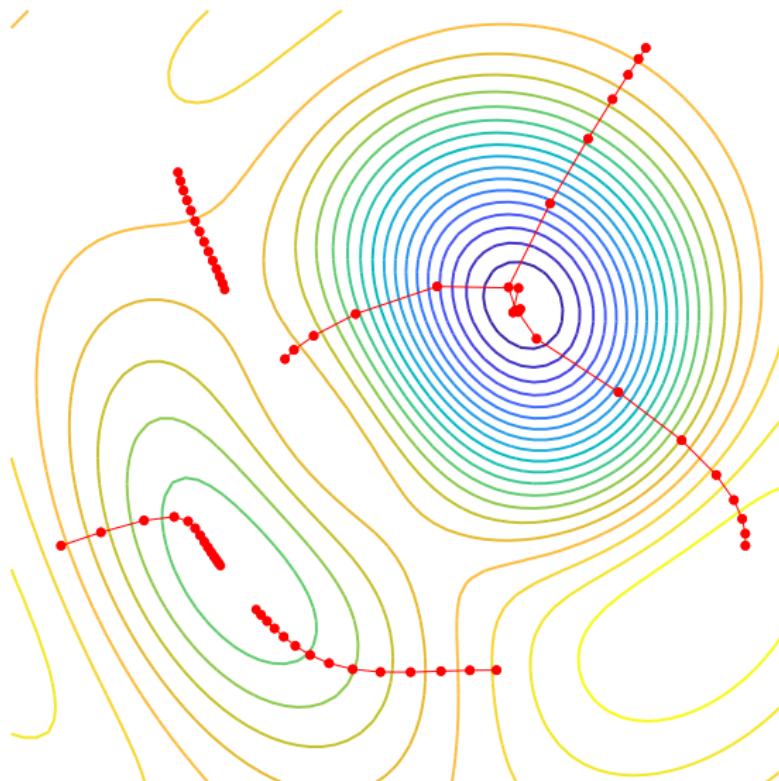
$\epsilon = 0.05$, iteration 11

gradient descent in two dimensions



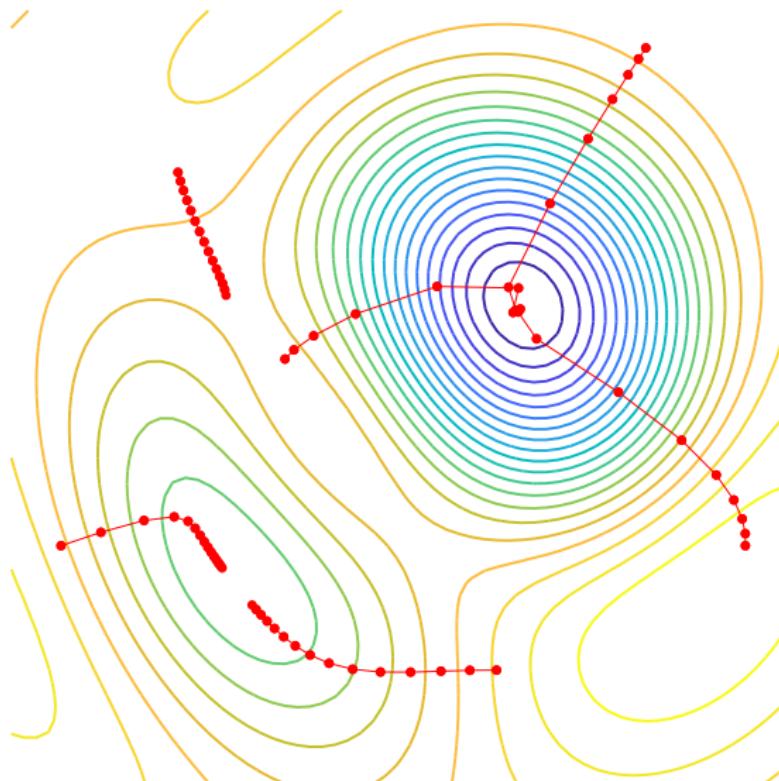
$\epsilon = 0.05$, iteration 12

gradient descent in two dimensions



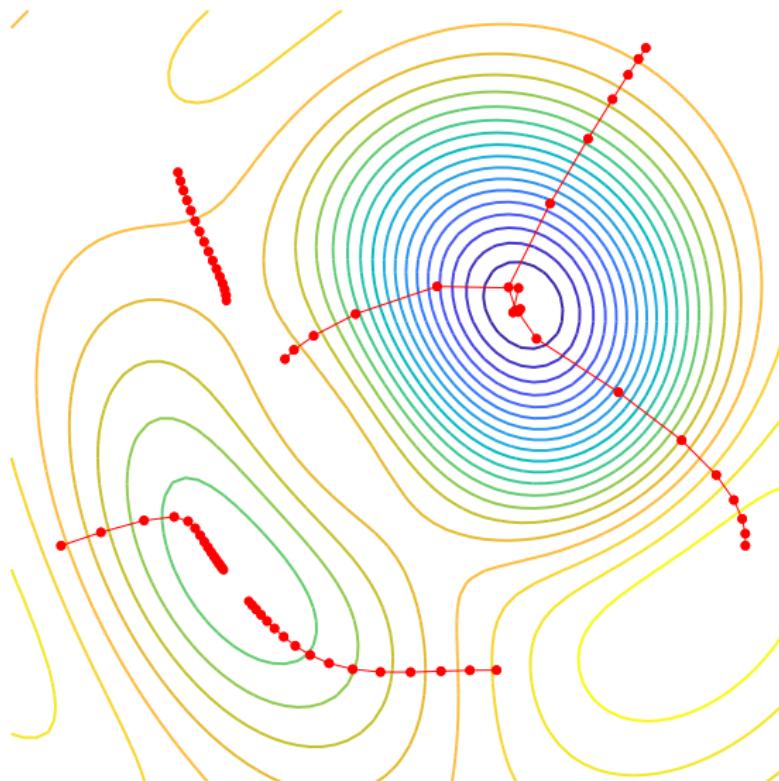
$\epsilon = 0.05$, iteration 13

gradient descent in two dimensions



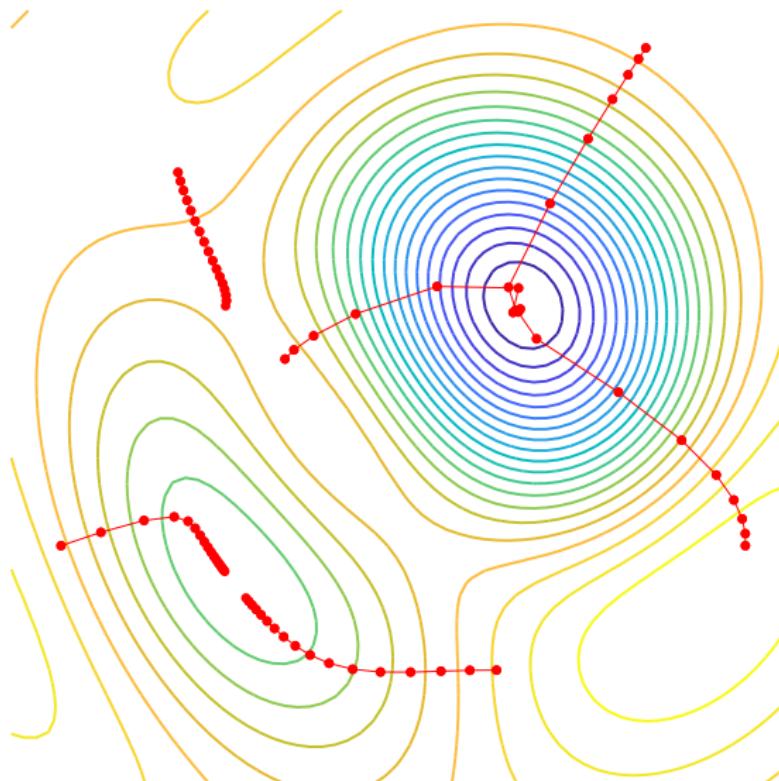
$\epsilon = 0.05$, iteration 14

gradient descent in two dimensions



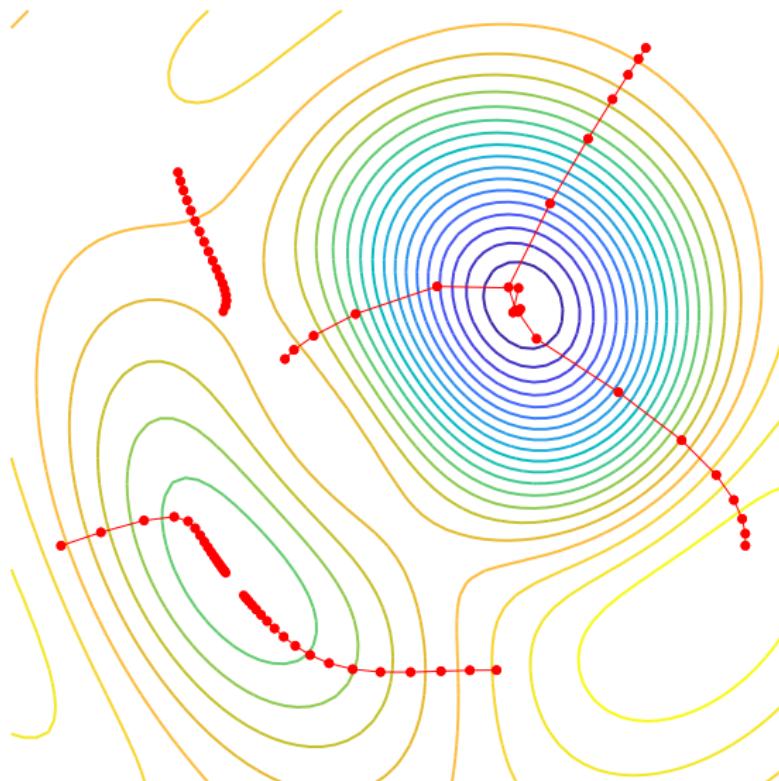
$\epsilon = 0.05$, iteration 15

gradient descent in two dimensions



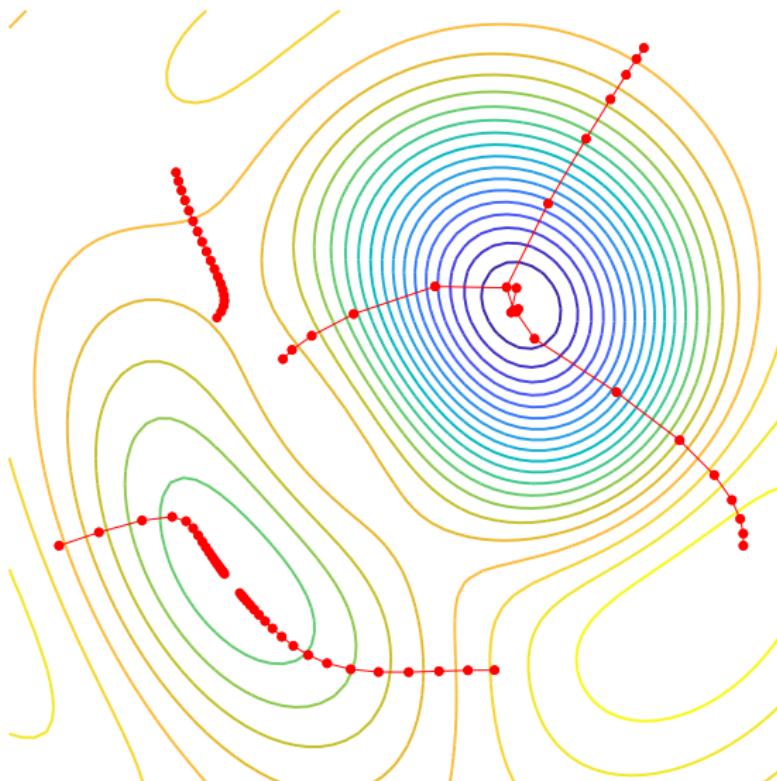
$\epsilon = 0.05$, iteration 16

gradient descent in two dimensions



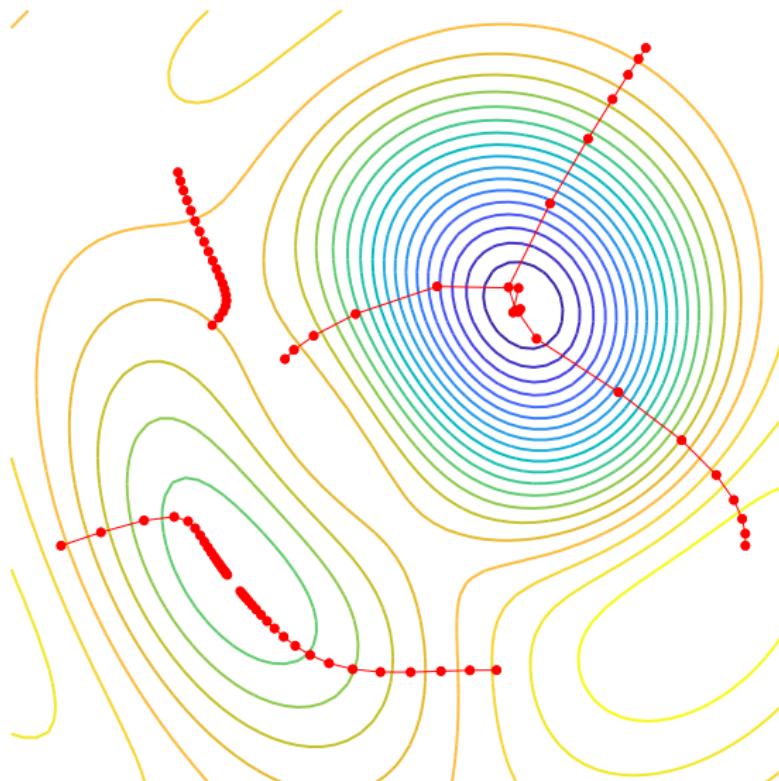
$\epsilon = 0.05$, iteration 17

gradient descent in two dimensions



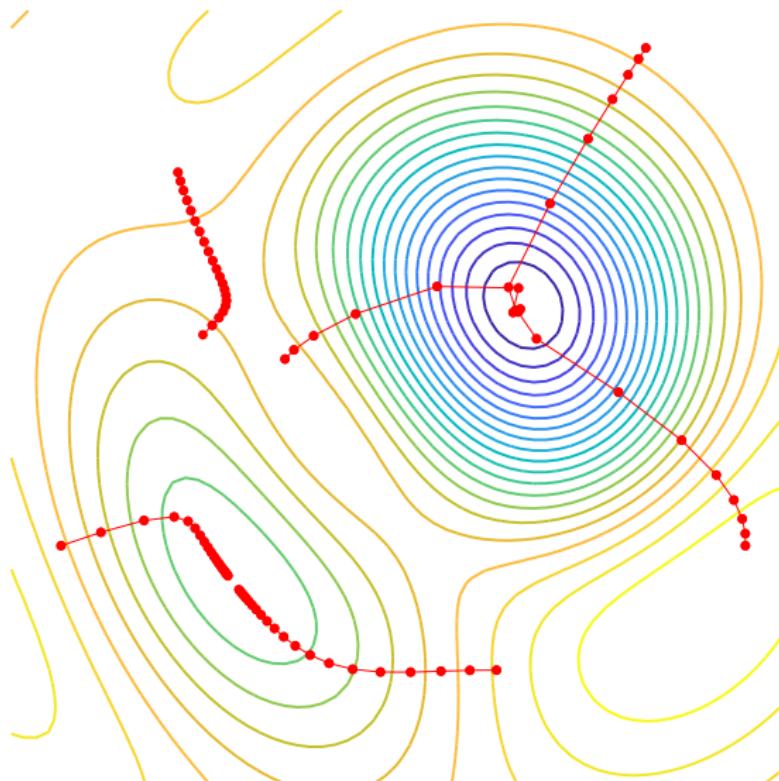
$\epsilon = 0.05$, iteration 18

gradient descent in two dimensions



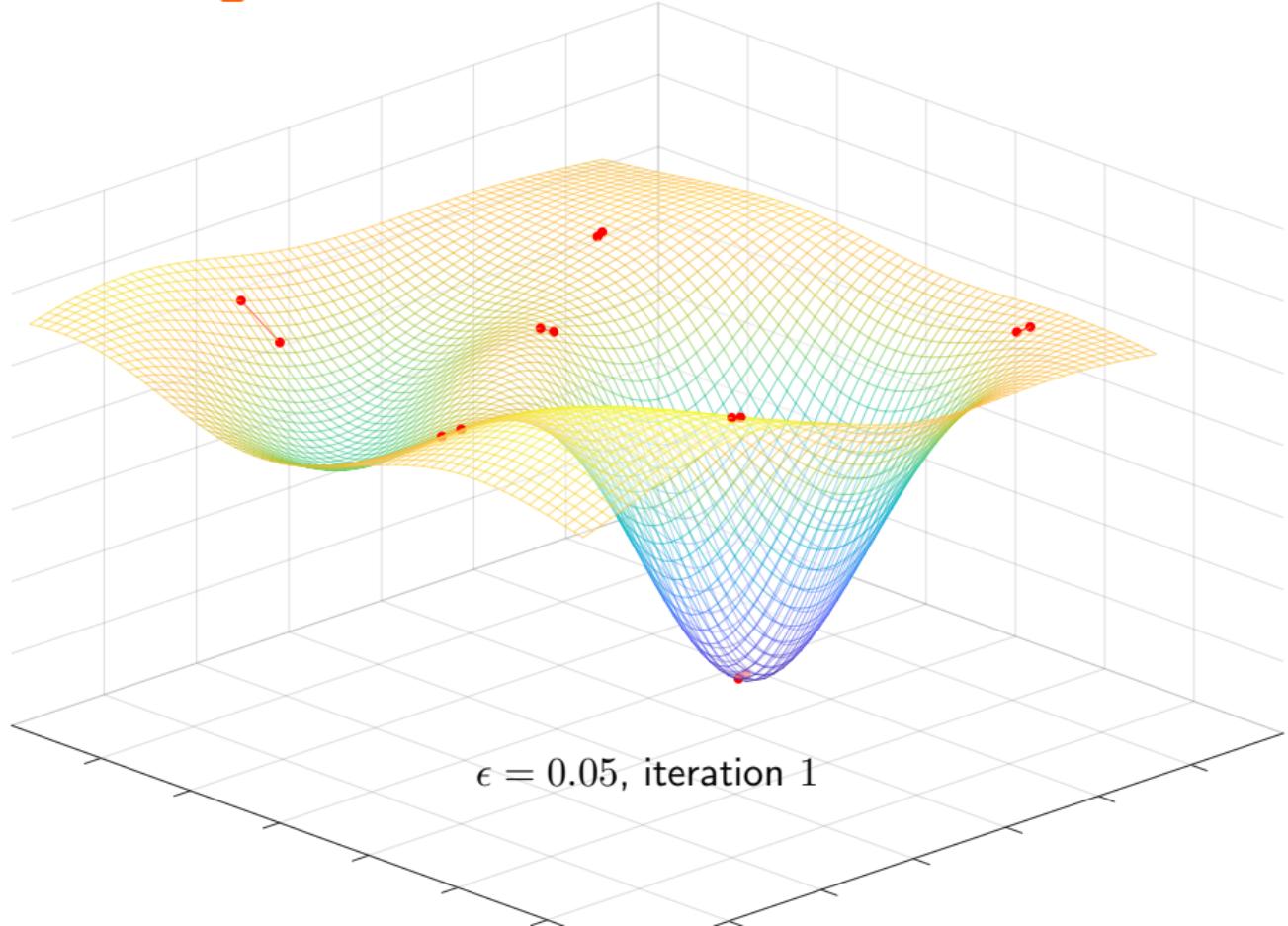
$\epsilon = 0.05$, iteration 19

gradient descent in two dimensions

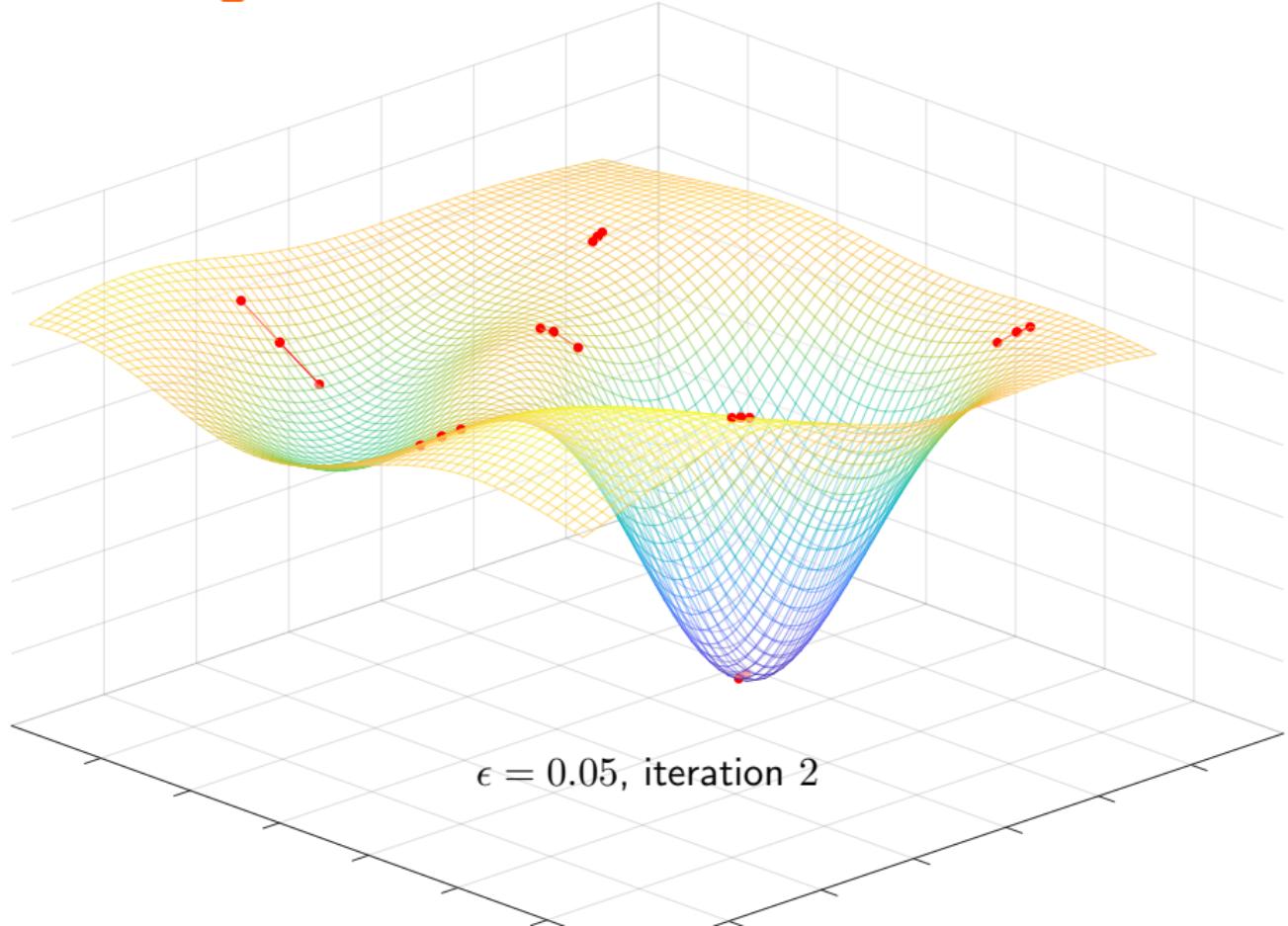


$\epsilon = 0.05$, iteration 20

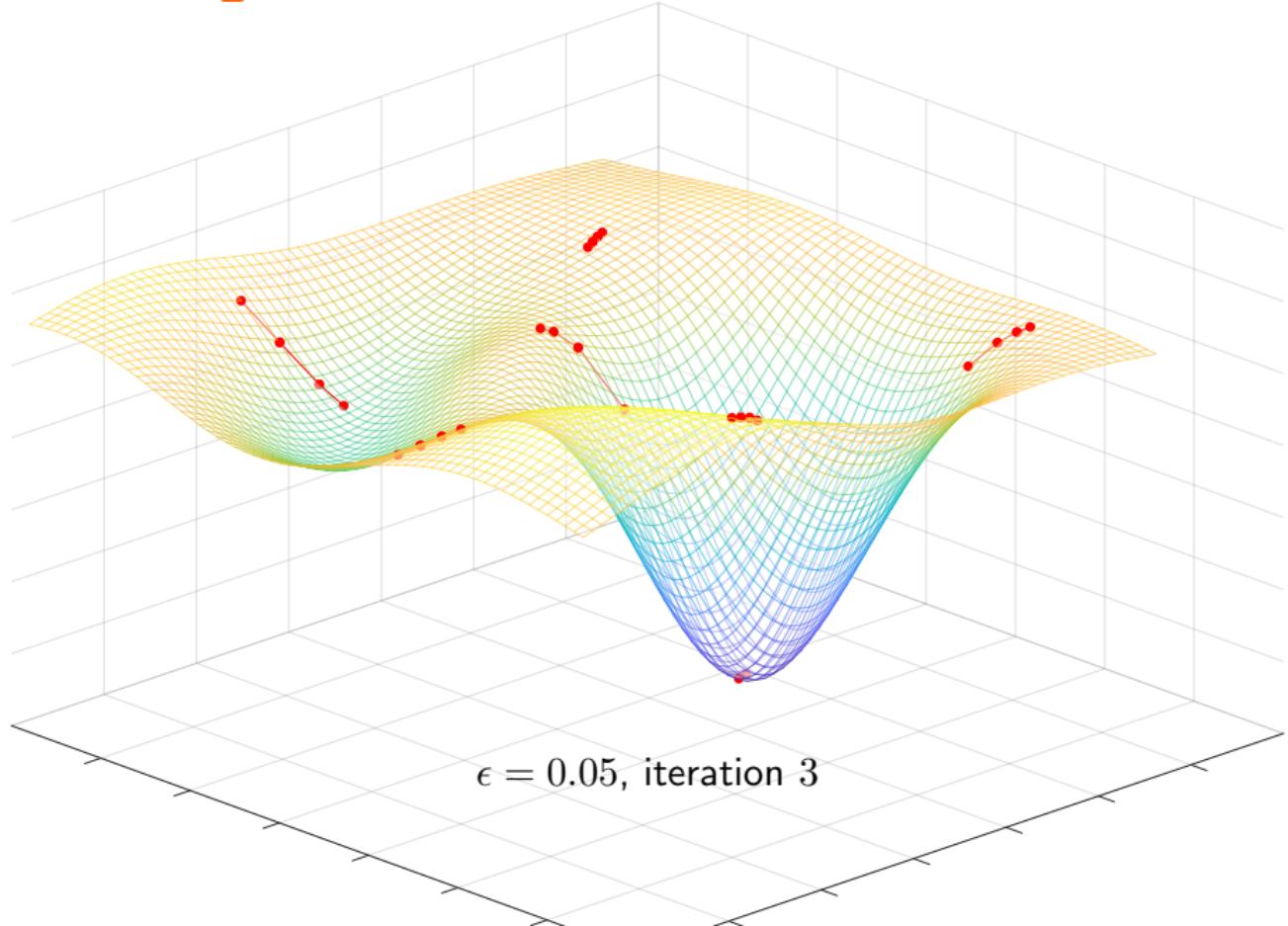
gradient descent in two dimensions



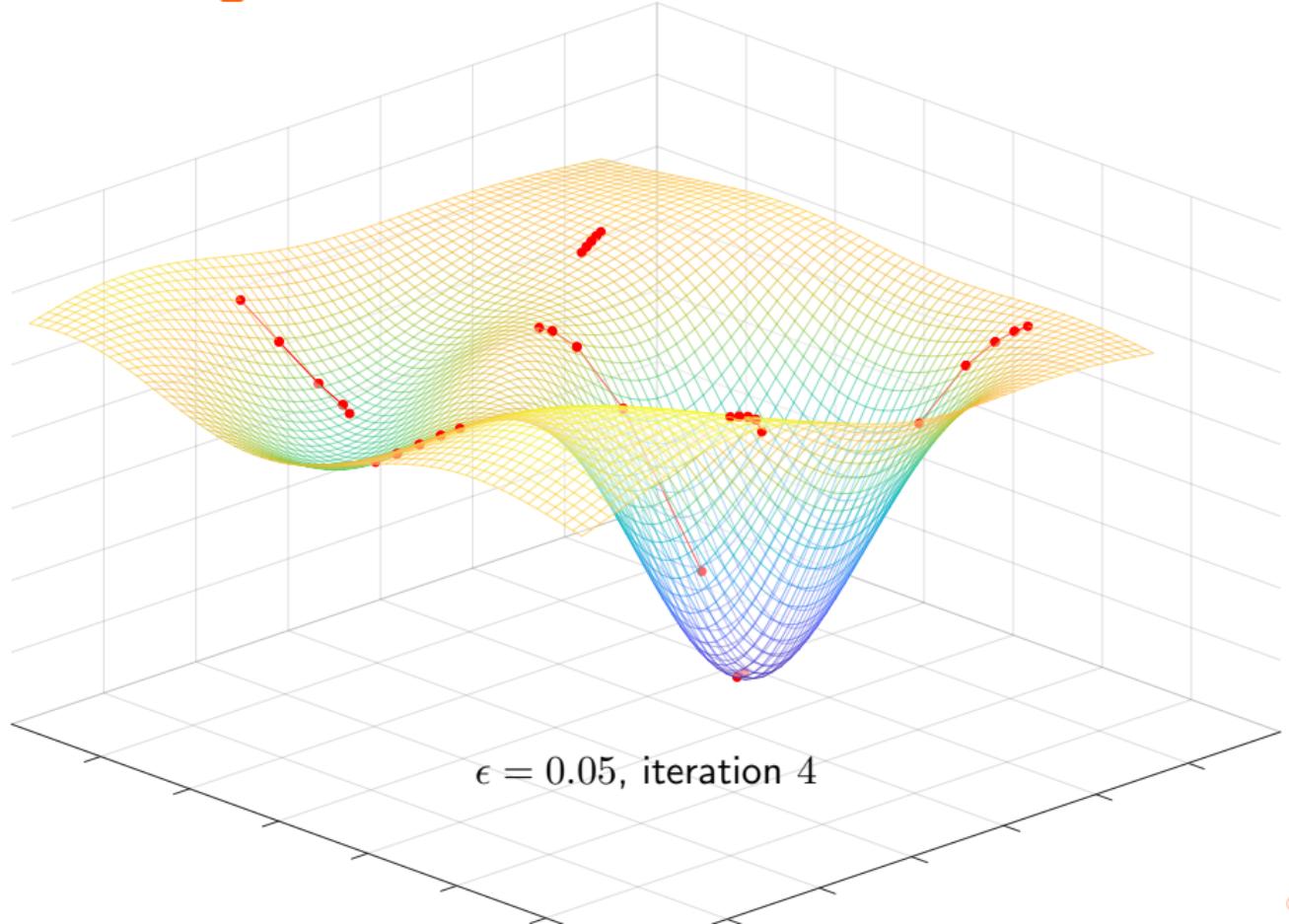
gradient descent in two dimensions



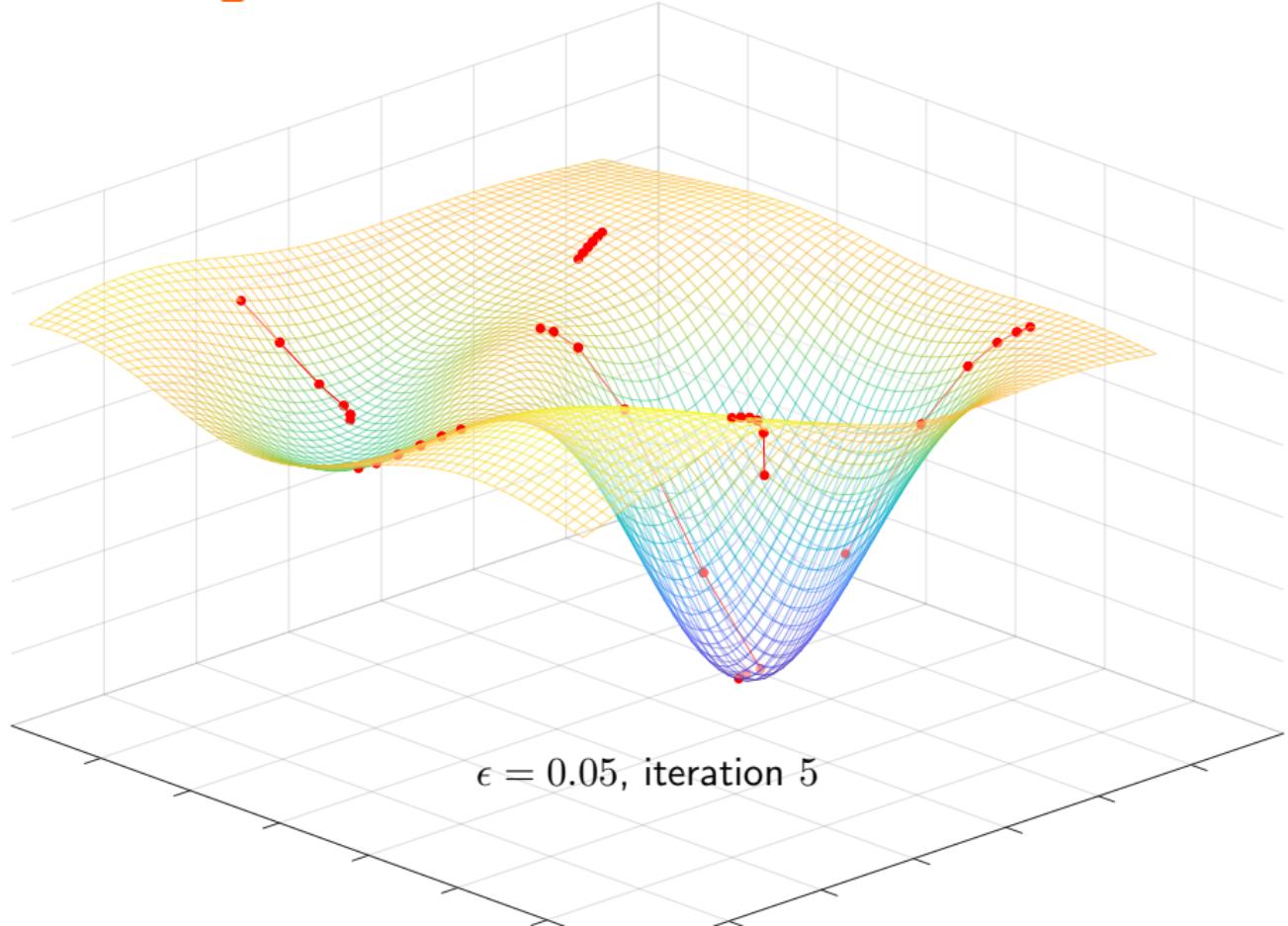
gradient descent in two dimensions



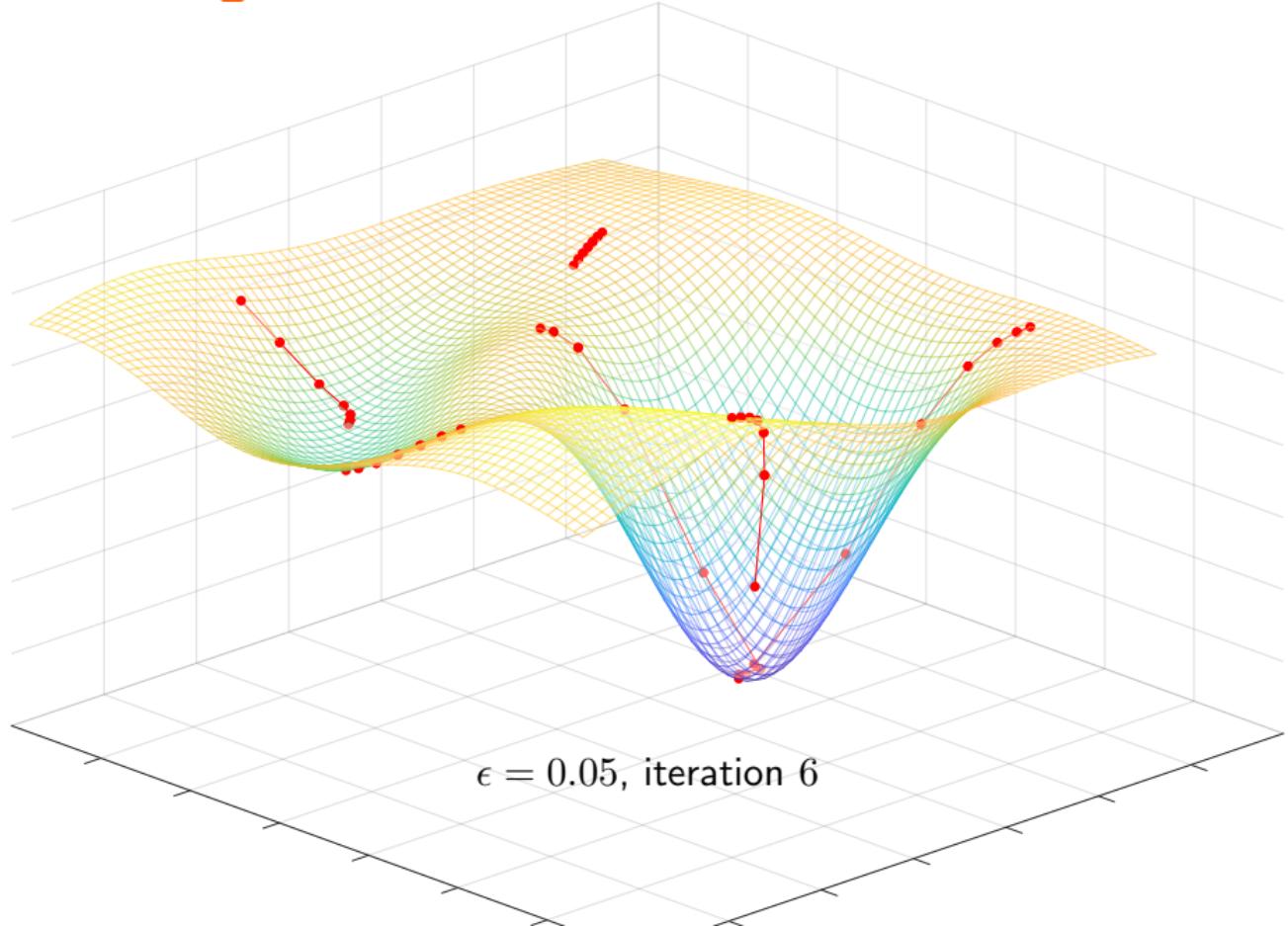
gradient descent in two dimensions



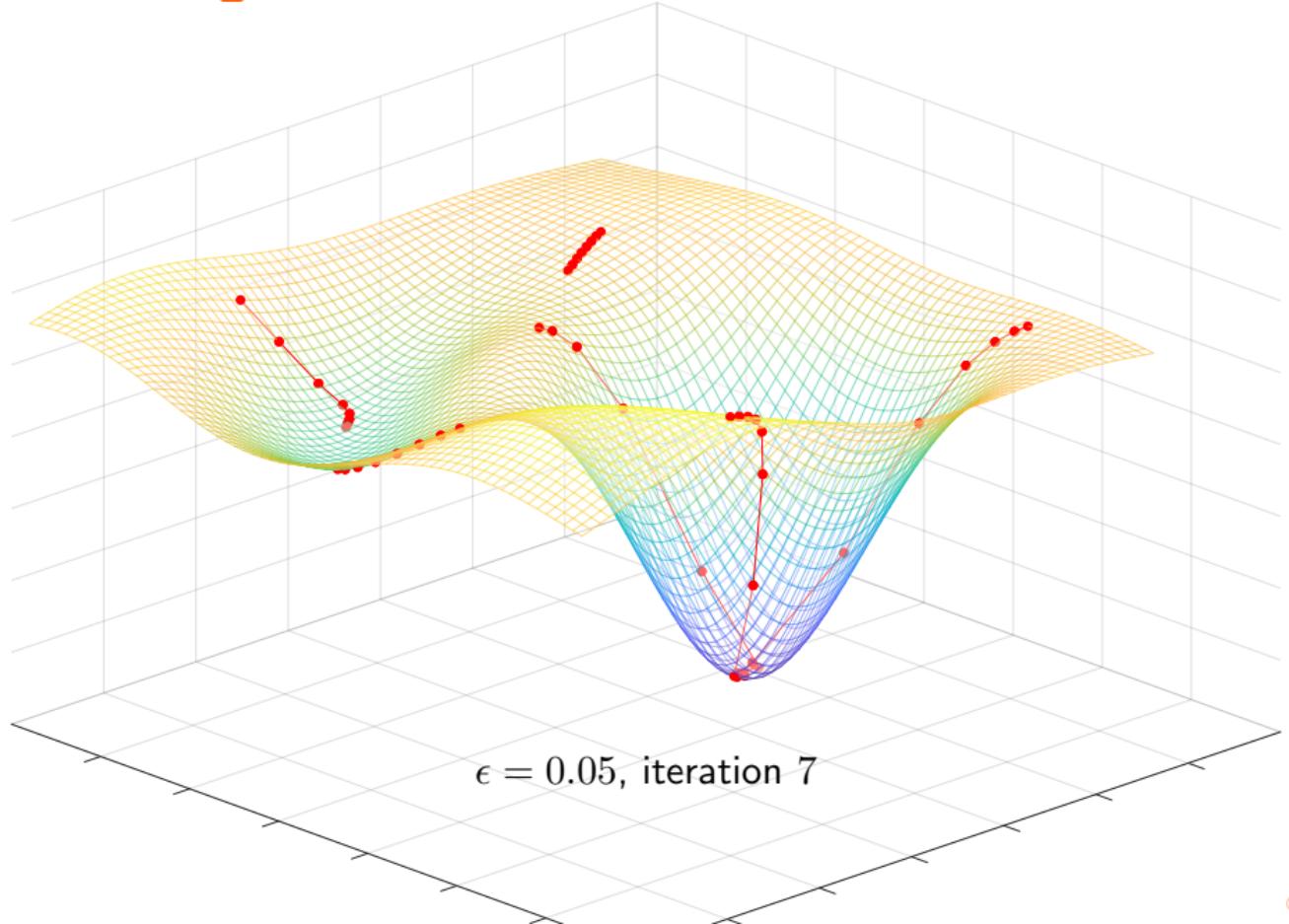
gradient descent in two dimensions



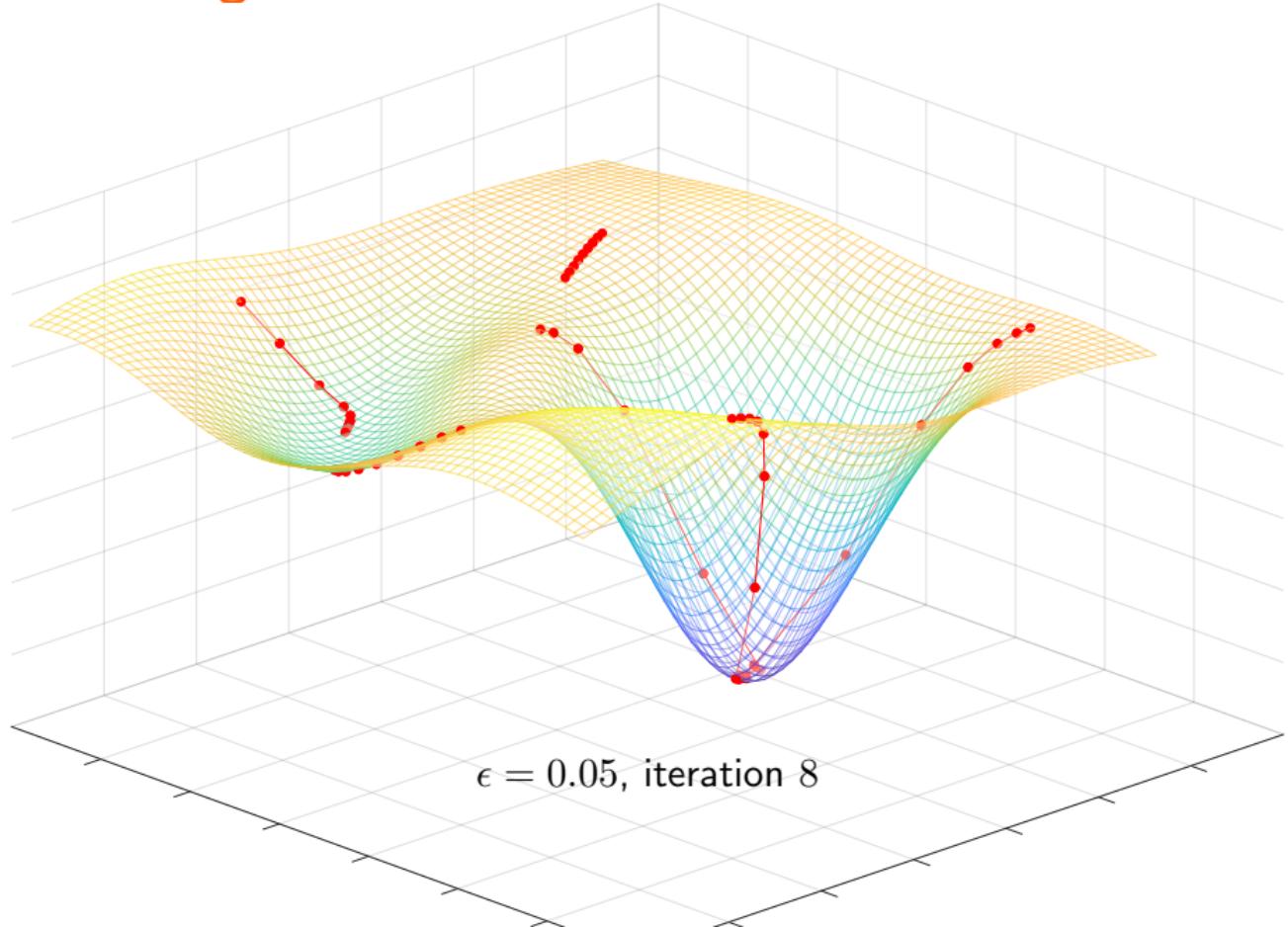
gradient descent in two dimensions



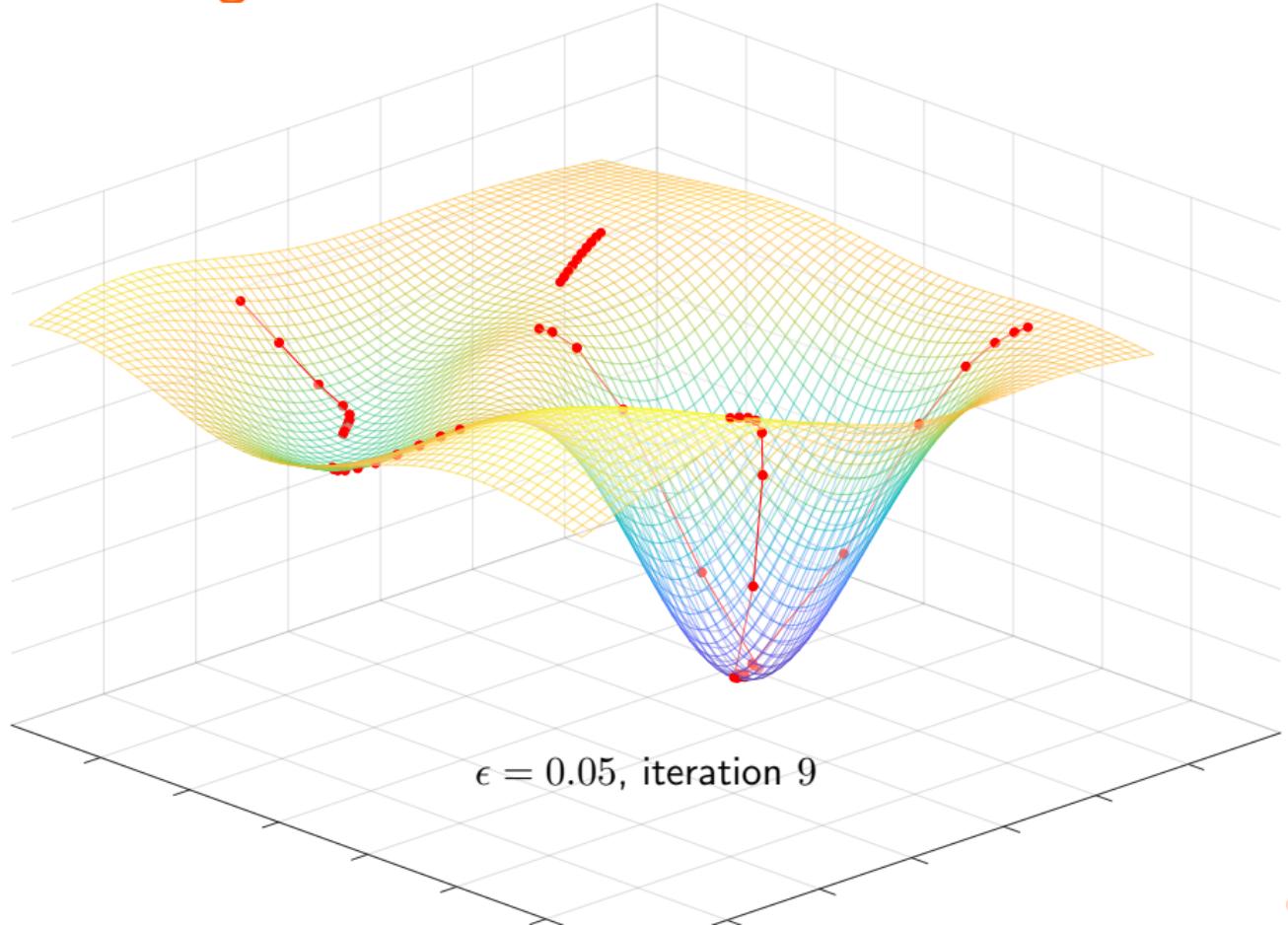
gradient descent in two dimensions



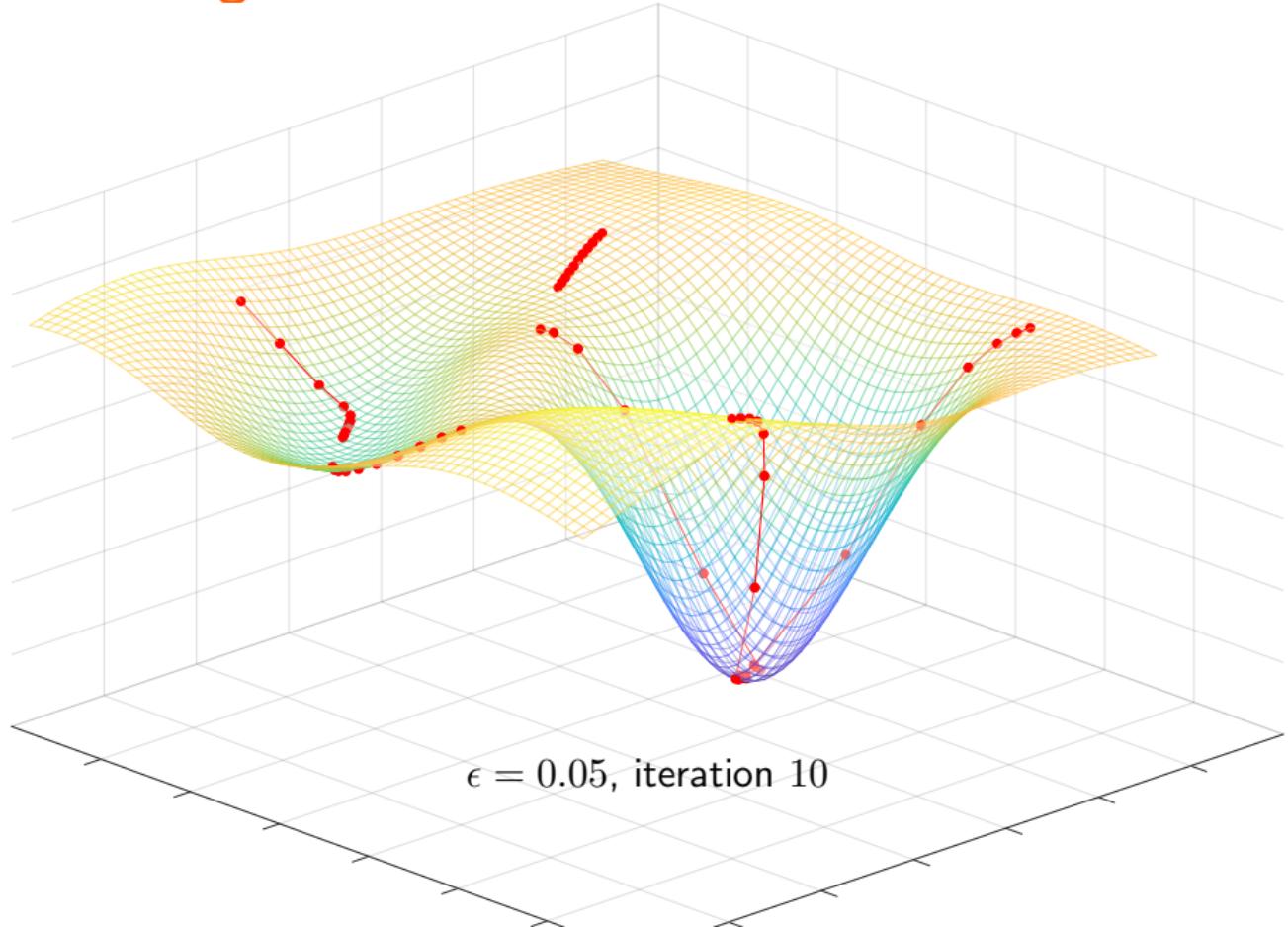
gradient descent in two dimensions



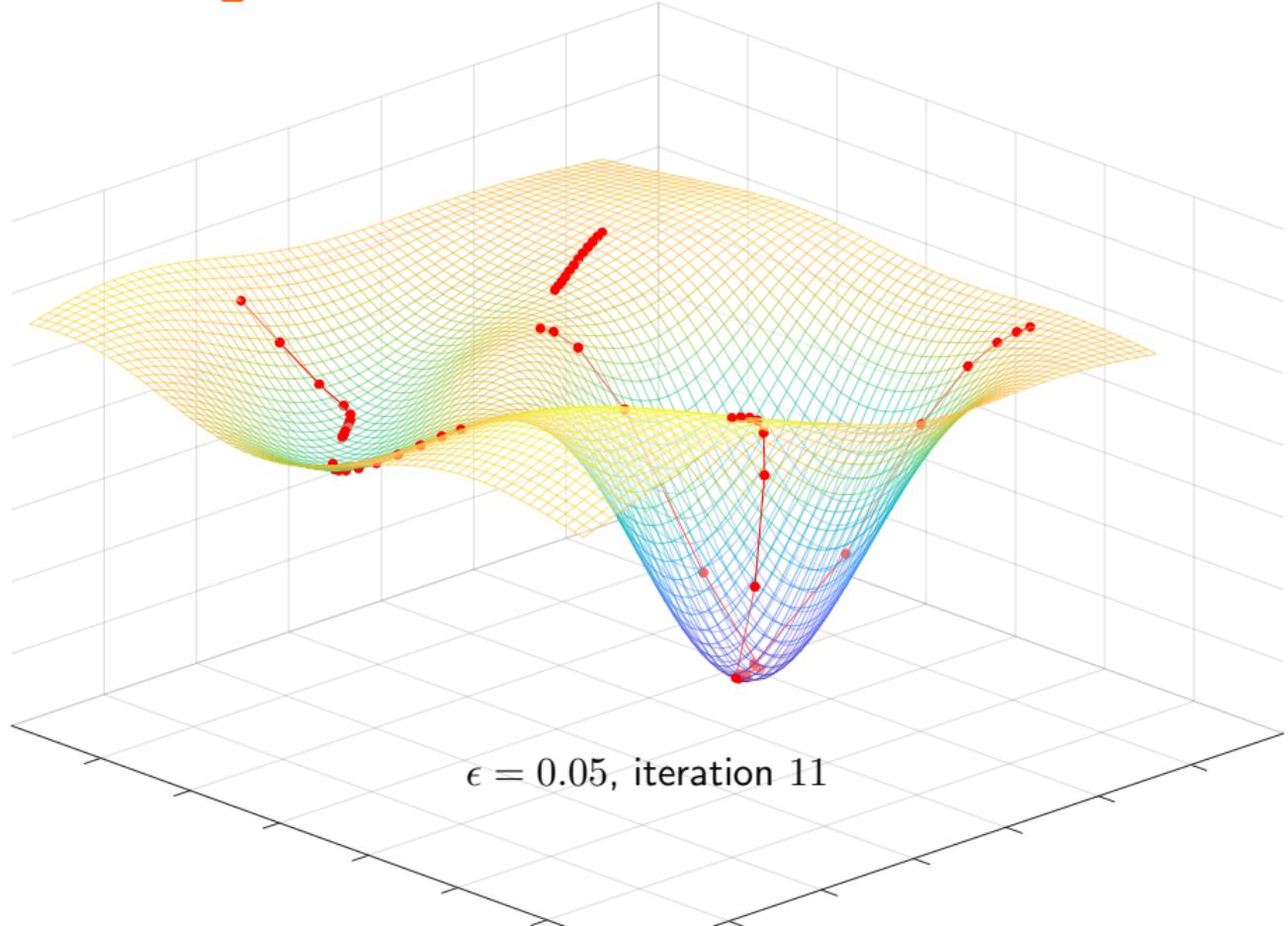
gradient descent in two dimensions



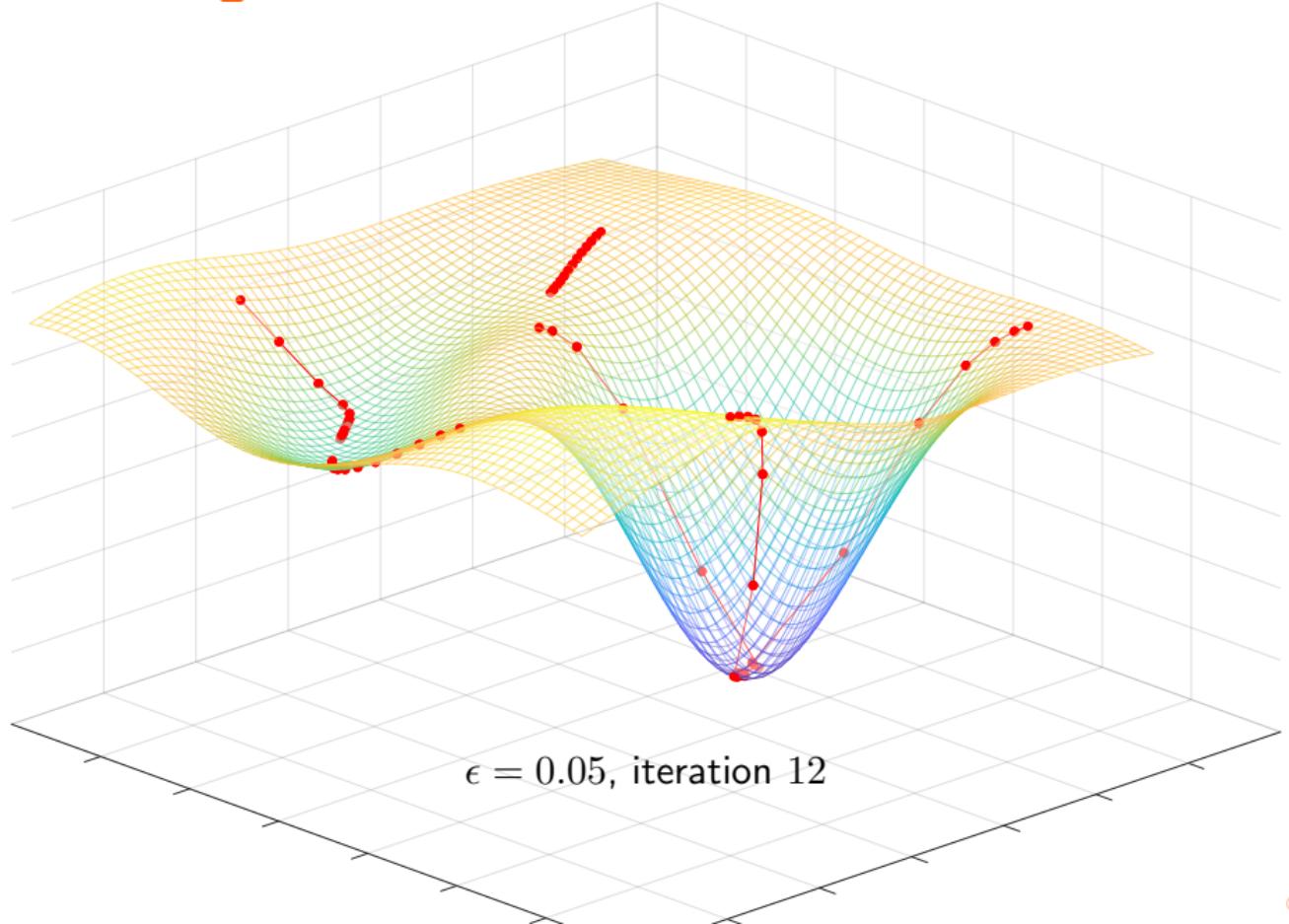
gradient descent in two dimensions



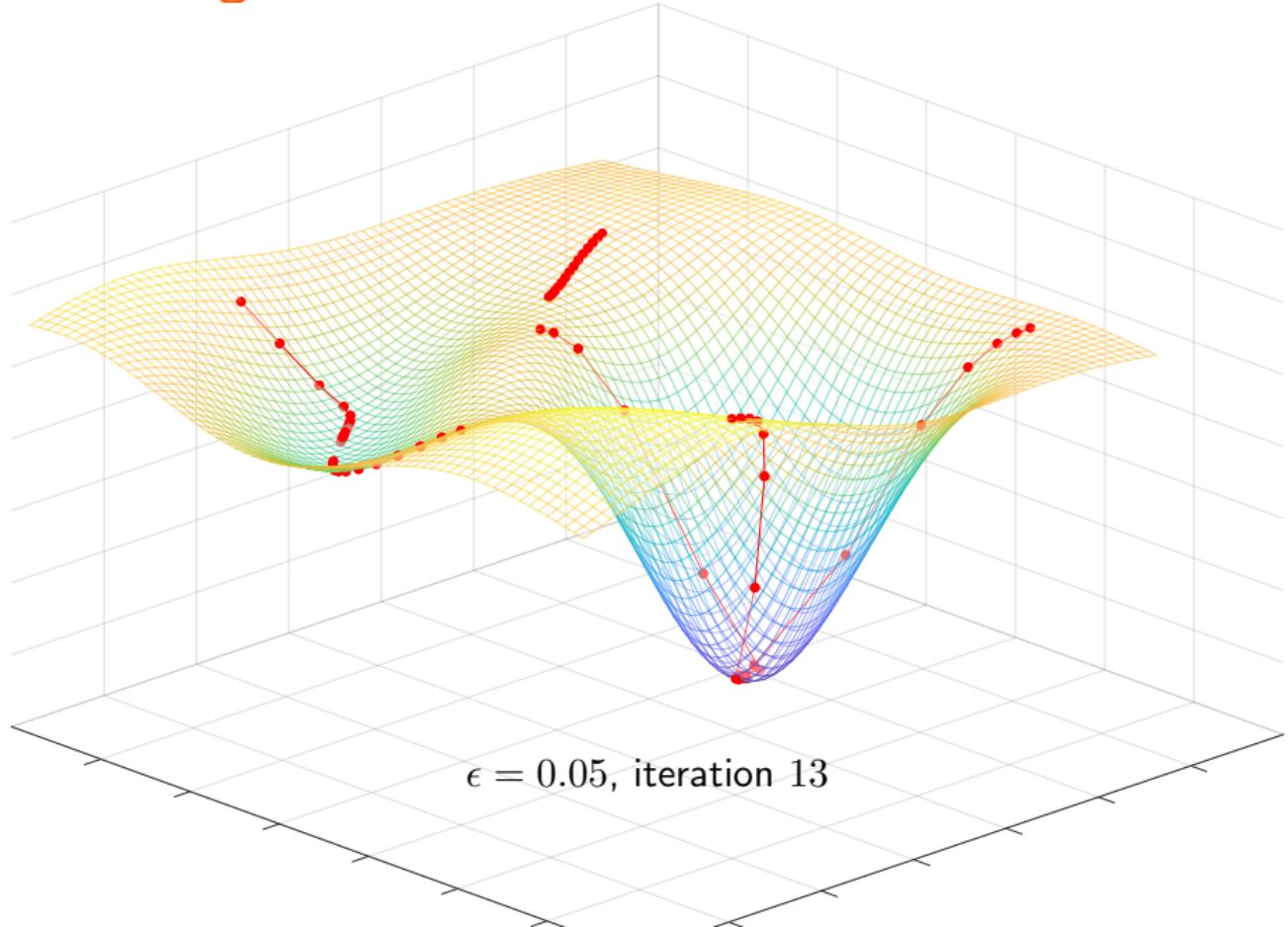
gradient descent in two dimensions



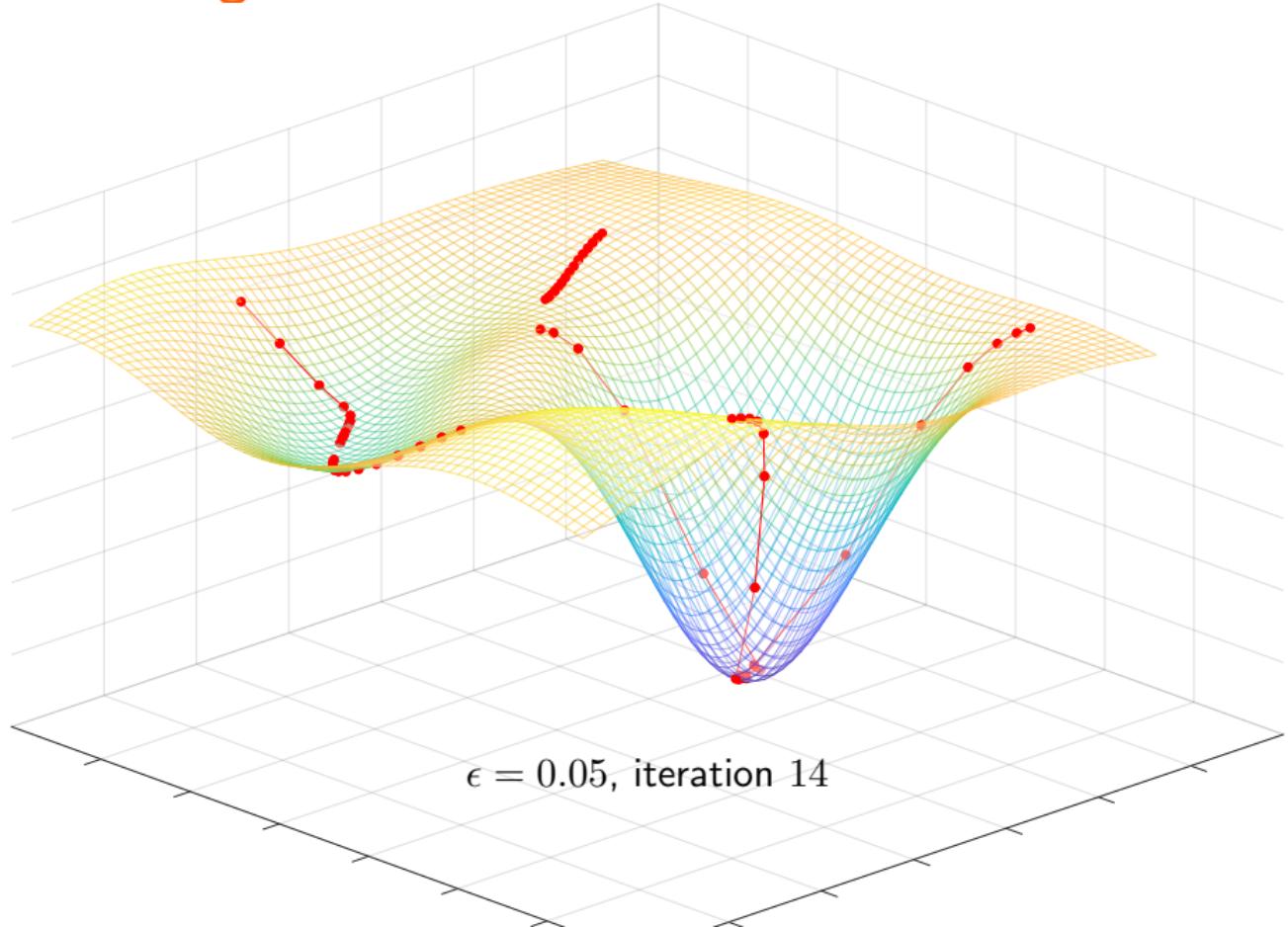
gradient descent in two dimensions



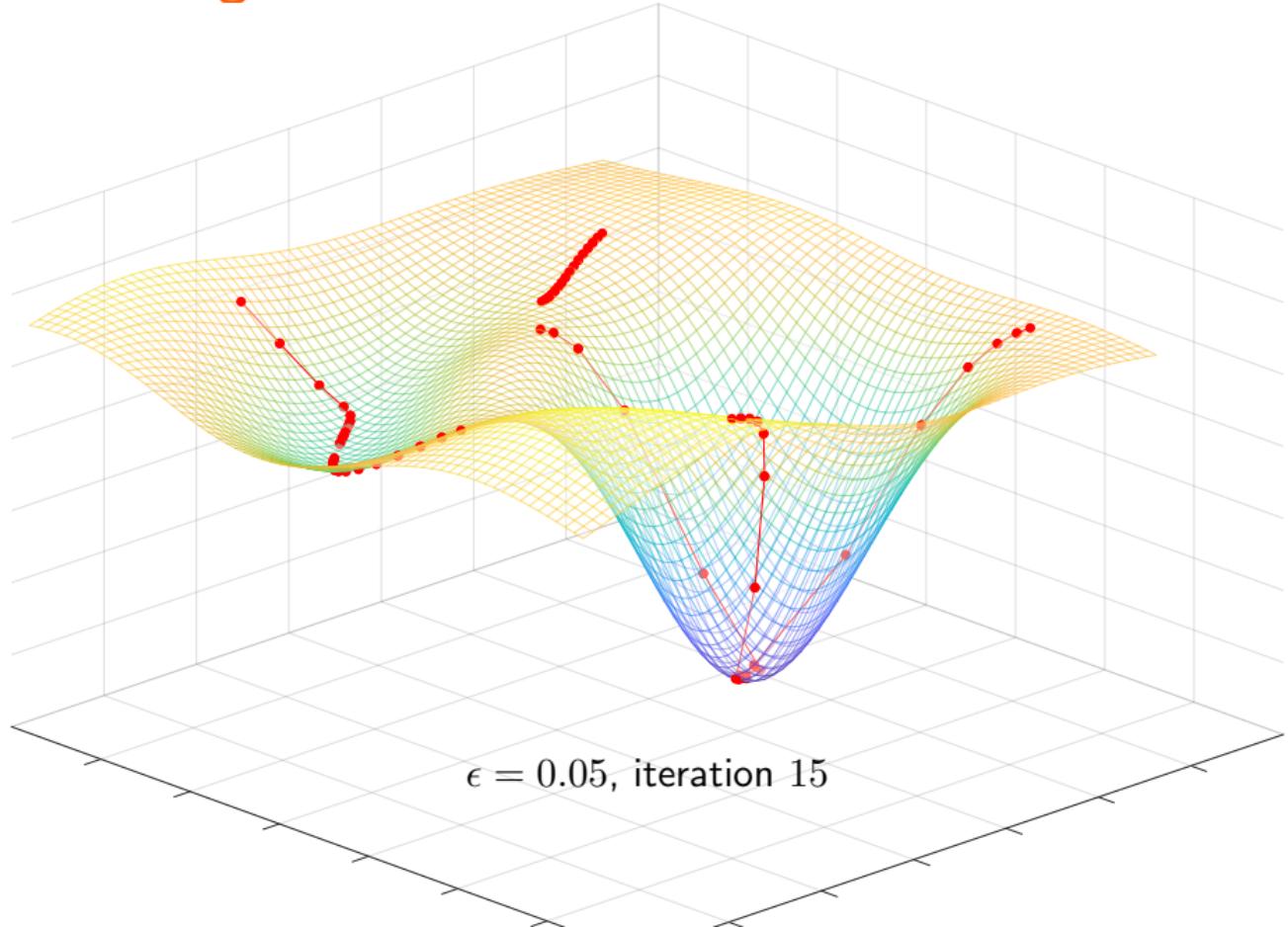
gradient descent in two dimensions



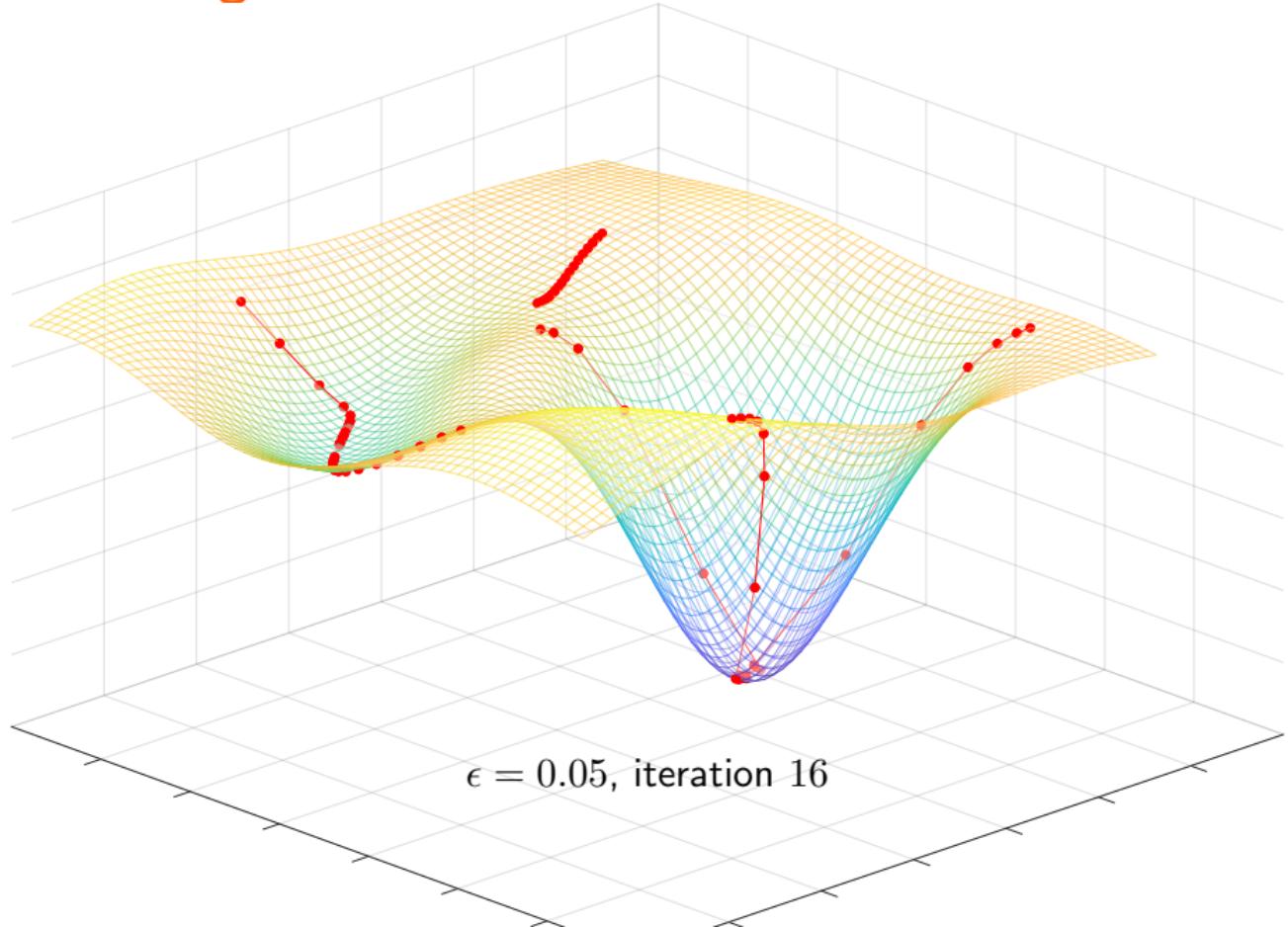
gradient descent in two dimensions



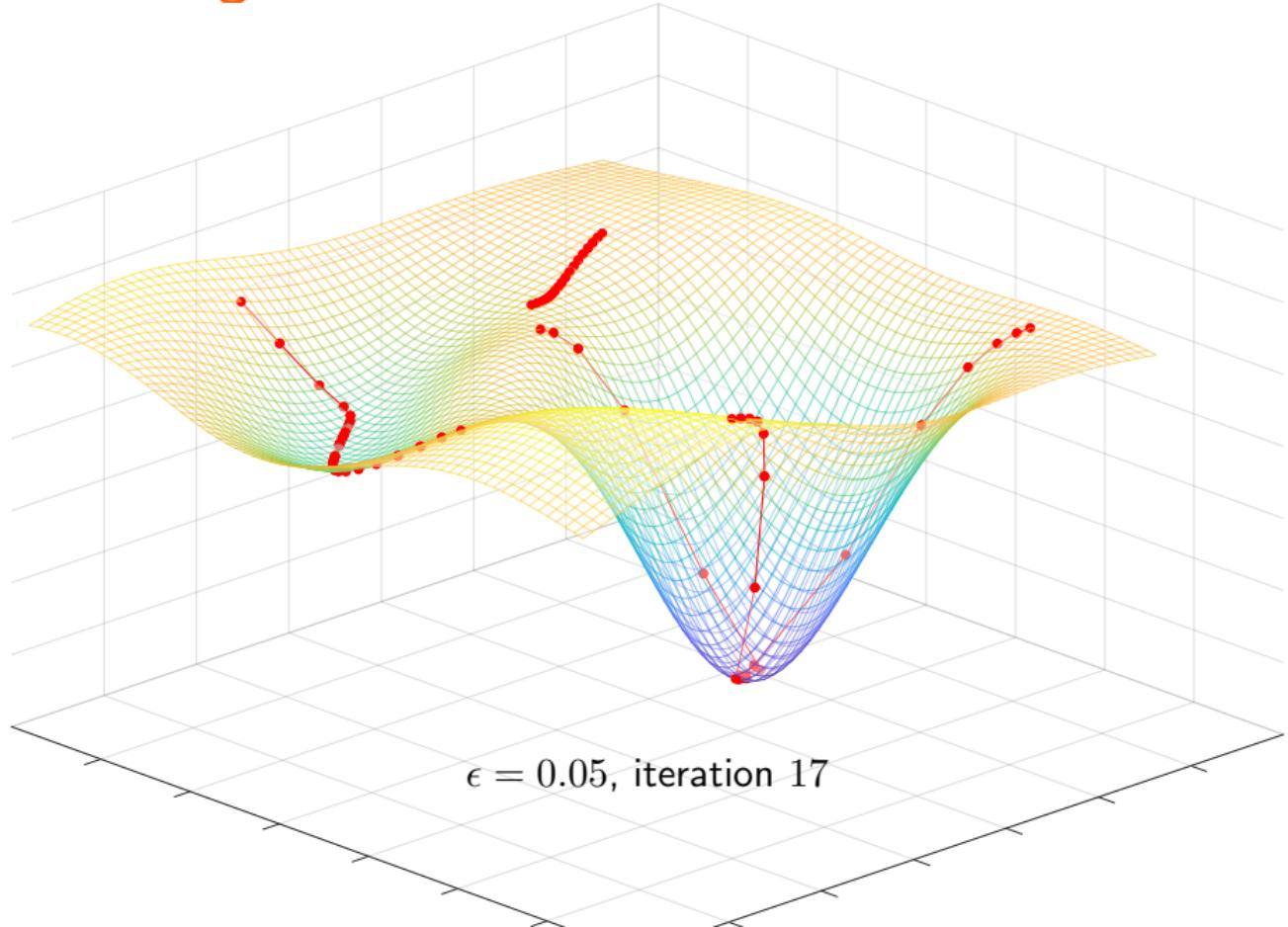
gradient descent in two dimensions



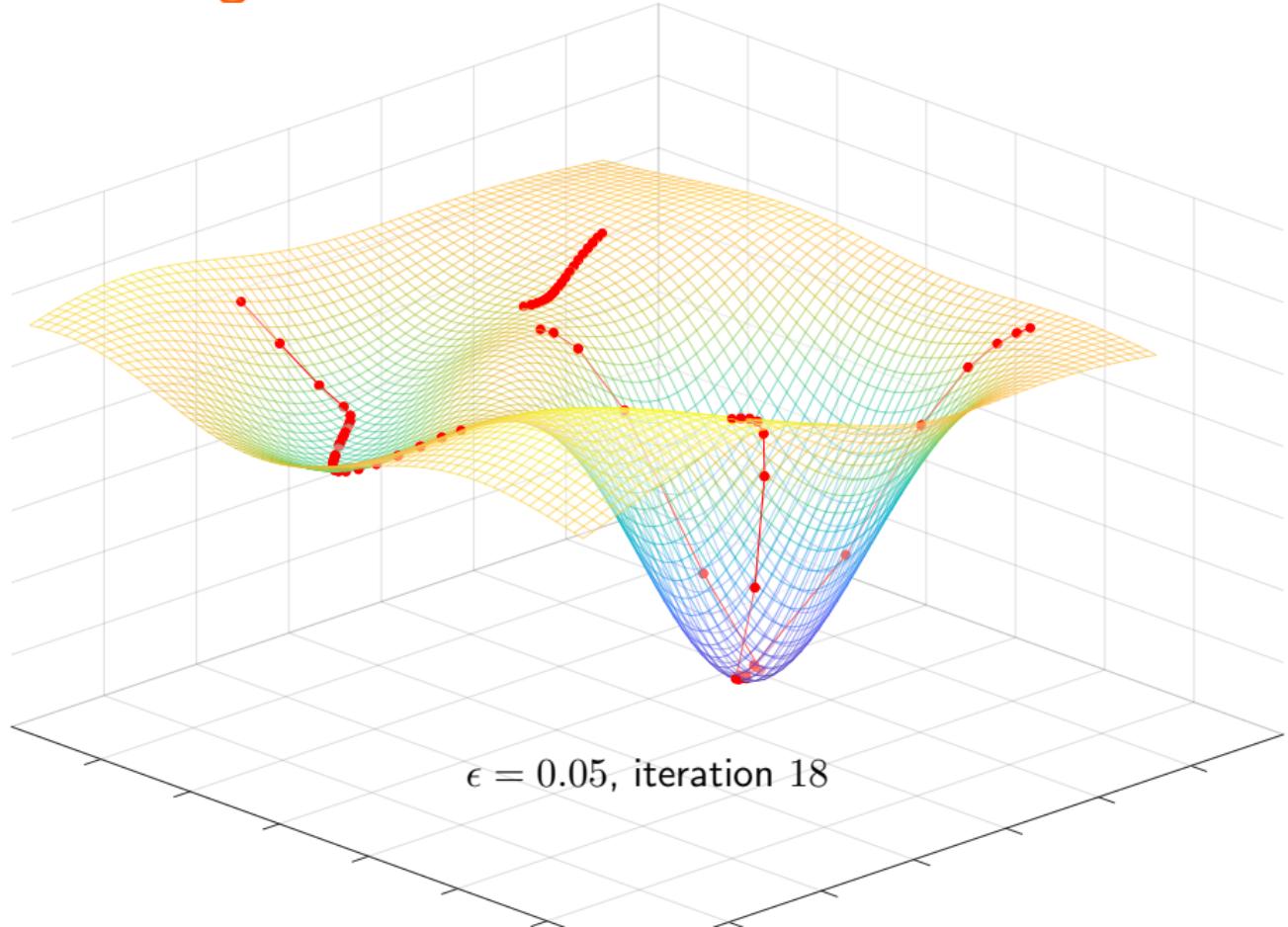
gradient descent in two dimensions



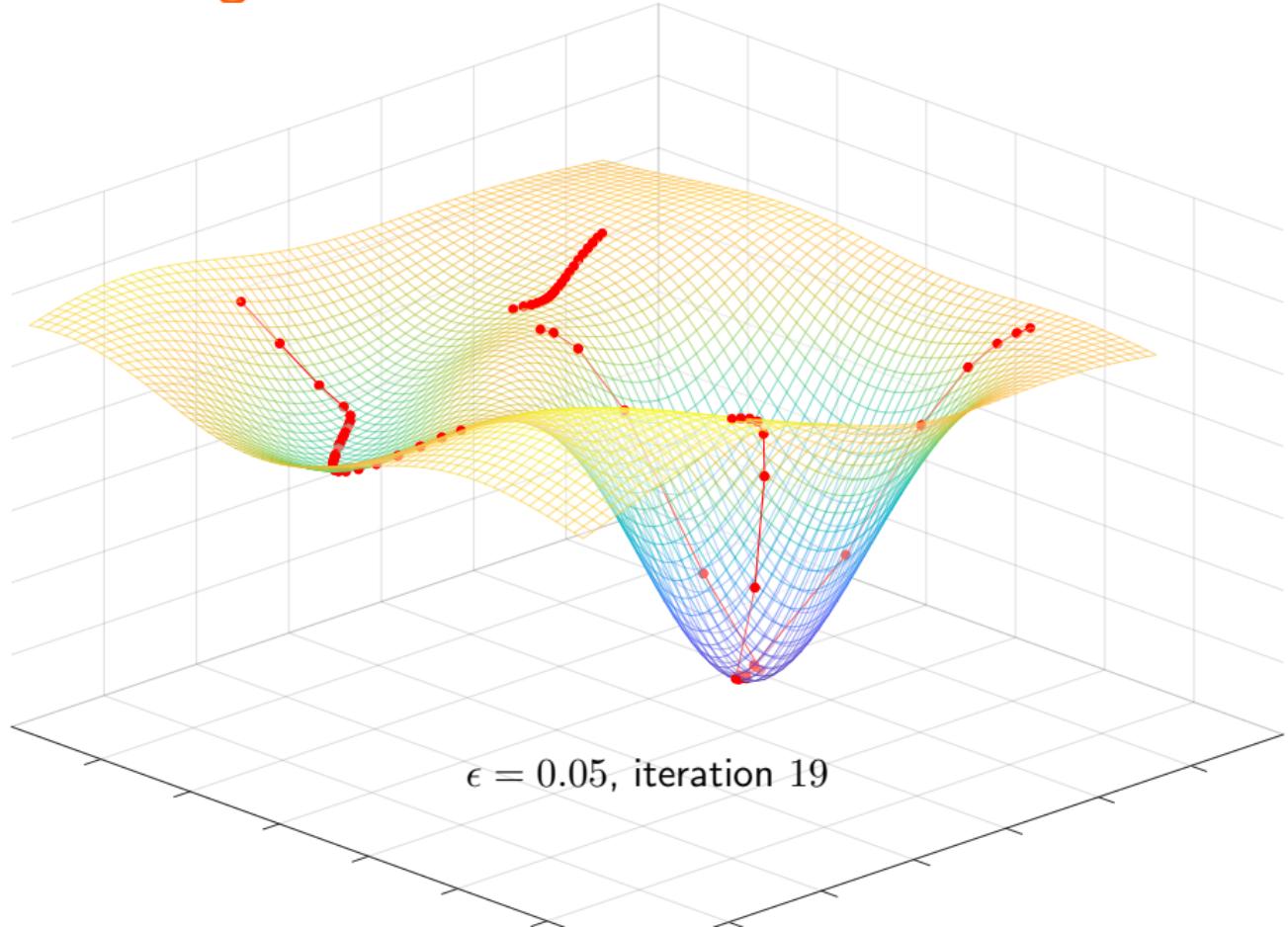
gradient descent in two dimensions



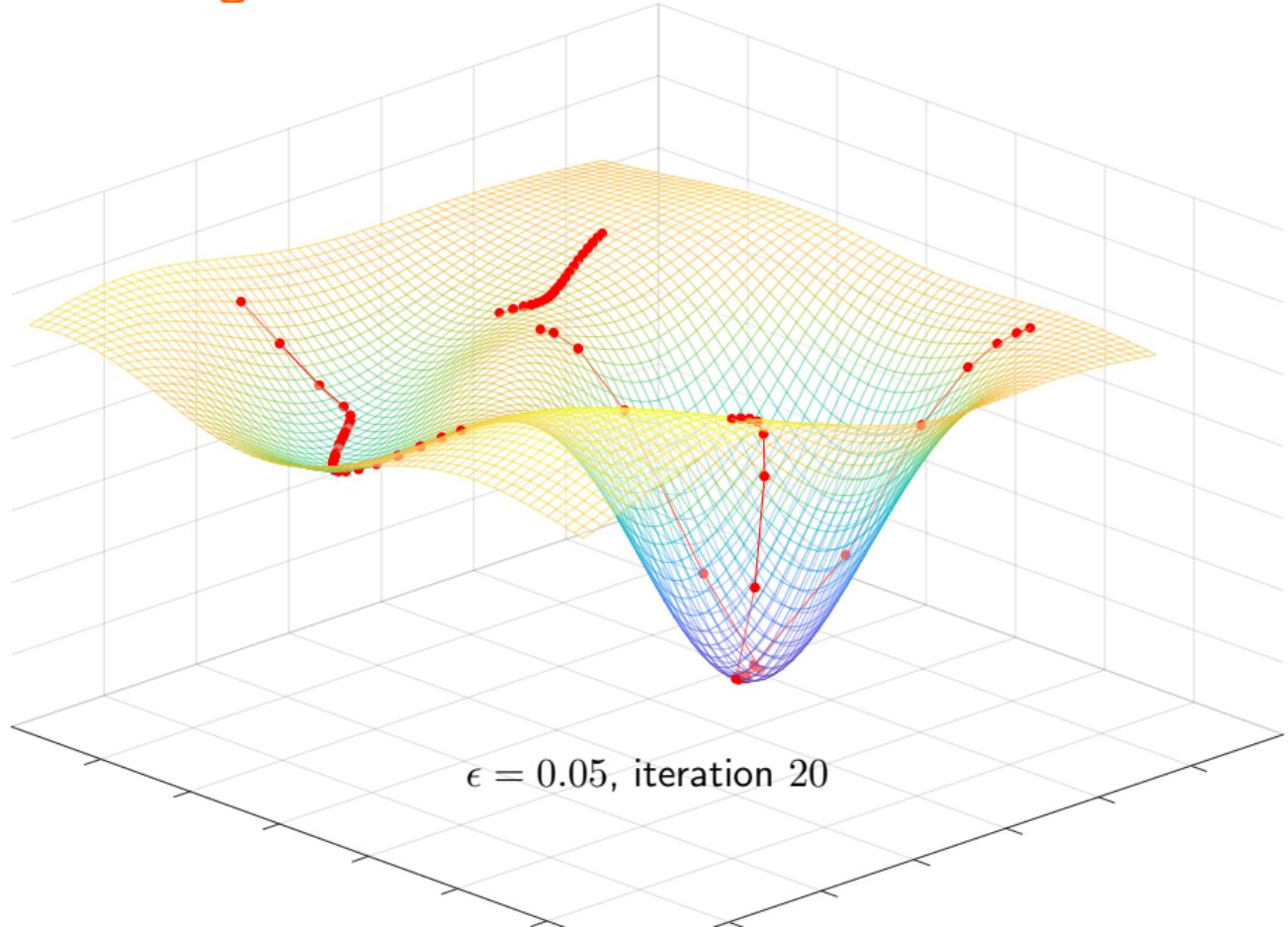
gradient descent in two dimensions



gradient descent in two dimensions



gradient descent in two dimensions



problems

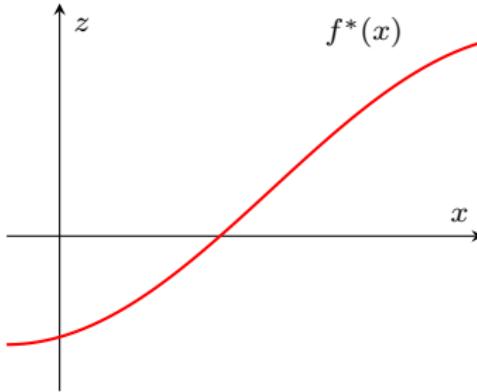
- f non-convex: local minima
- $d \times d$ Hessian matrix too expensive (d can be millions): unknown curvature
- high condition number: elongated regions
- plateaus, saddle points: no progress
- $\nabla f = \sum_{i=1}^n \nabla f_i$ itself too expensive (n can also be millions)

problems

- f non-convex: local minima
- $d \times d$ Hessian matrix too expensive (d can be millions): unknown curvature
- high condition number: elongated regions
- plateaus, saddle points: no progress
- $\nabla f = \sum_{i=1}^n \nabla f_i$ itself too expensive (n can also be millions)

sequential estimation

[Robbins and Monro 1951]



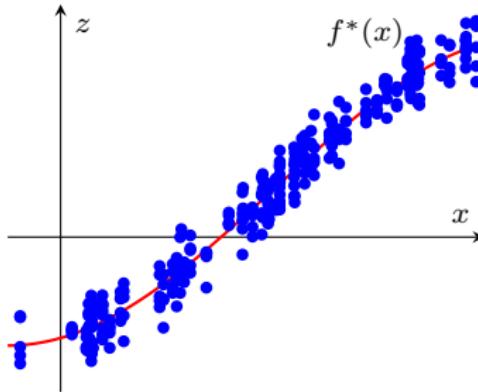
- suppose f^* is the expectation of random variable z conditional on x , and f is its empirical estimate on n samples

$$f^*(x) := \mathbb{E}[z|x] \quad f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x)$$

- we would like to estimate a root x^* of f where $f(x^*) = 0$

sequential estimation

[Robbins and Monro 1951]



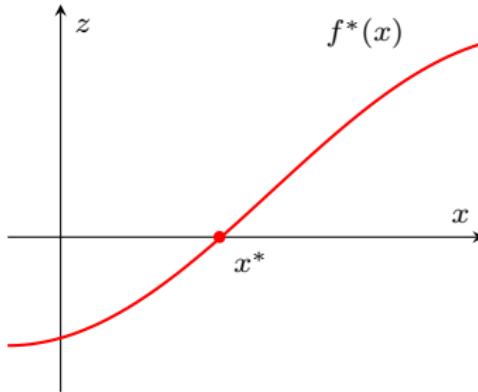
- suppose f^* is the expectation of random variable z conditional on x , and f is its empirical estimate on n samples

$$f^*(x) := \mathbb{E}[z|x] \quad f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x)$$

- we would like to estimate a root x^* of f where $f(x^*) = 0$

sequential estimation

[Robbins and Monro 1951]



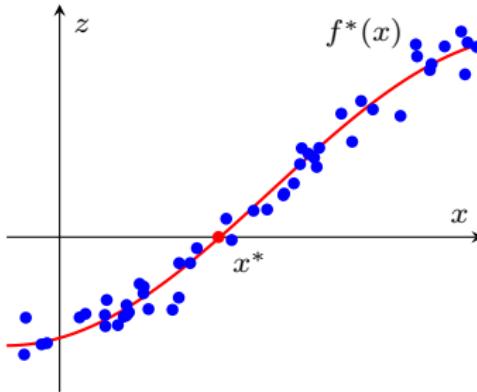
- suppose f^* is the expectation of random variable z conditional on x , and f is its empirical estimate on n samples

$$f^*(x) := \mathbb{E}[z|x] \quad f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x)$$

- we would like to estimate a root x^* of f where $f(x^*) = 0$

sequential estimation

[Robbins and Monro 1951]



- then we can estimate x^* sequentially

$$x^{(\tau+1)} = x^{(\tau)} - \epsilon_\tau z(x^{(\tau)}) = x^{(\tau)} - \epsilon_\tau f_i(x^{(\tau)})$$

where $z(x^{(\tau)})$ is an observation of z when $x = x^{(\tau)}$ and i is a random index in $\{1, \dots, n\}$

sufficient conditions for convergence

- successive corrections decrease in magnitude

$$\lim_{\tau \rightarrow \infty} \epsilon_\tau = 0$$

- the algorithm does not converge short of the root

$$\sum_{\tau=1}^{\infty} \epsilon_\tau = \infty$$

- the accumulated “noise” has finite variance

$$\sum_{\tau=1}^{\infty} \epsilon_\tau^2 < \infty$$

online gradient descent

- now, replace x by the parameters θ of our model, and f by ∇E , the gradient of our empirical risk
- the update rule becomes

$$\theta^{(\tau+1)} \leftarrow \theta^{(\tau)} - \epsilon_\tau \nabla E_i(\theta^{(\tau)})$$

- and, under the same conditions, it converges to a root of

$$\nabla E(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla E_i(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla L(f(\mathbf{x}_i; \theta), t_i)$$

that is, to a local minimum of E

- mini-batch gradient descent is similar but with less “stochastic noise”

online gradient descent

- now, replace x by the parameters θ of our model, and f by ∇E , the gradient of our empirical risk
- the update rule becomes

$$\theta^{(\tau+1)} \leftarrow \theta^{(\tau)} - \epsilon_\tau \nabla E_i(\theta^{(\tau)})$$

- and, under the same conditions, it converges to a root of

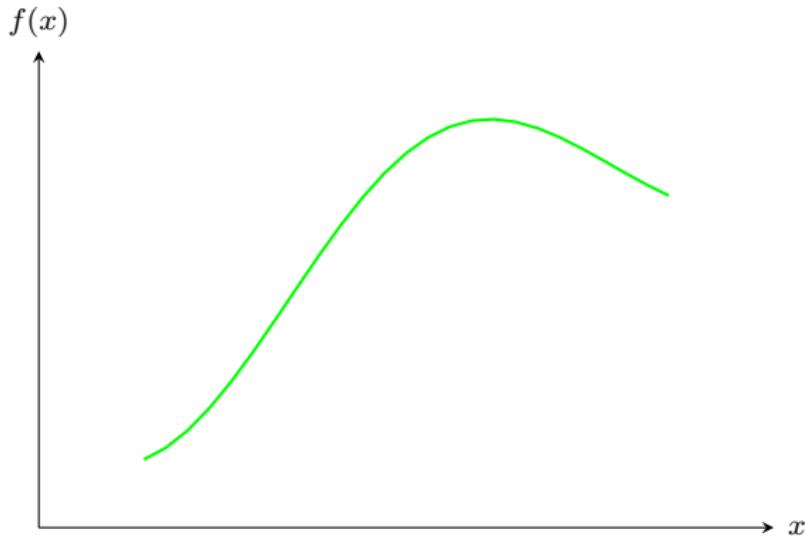
$$\nabla E(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla E_i(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla L(f(\mathbf{x}_i; \theta), t_i)$$

that is, to a local minimum of E

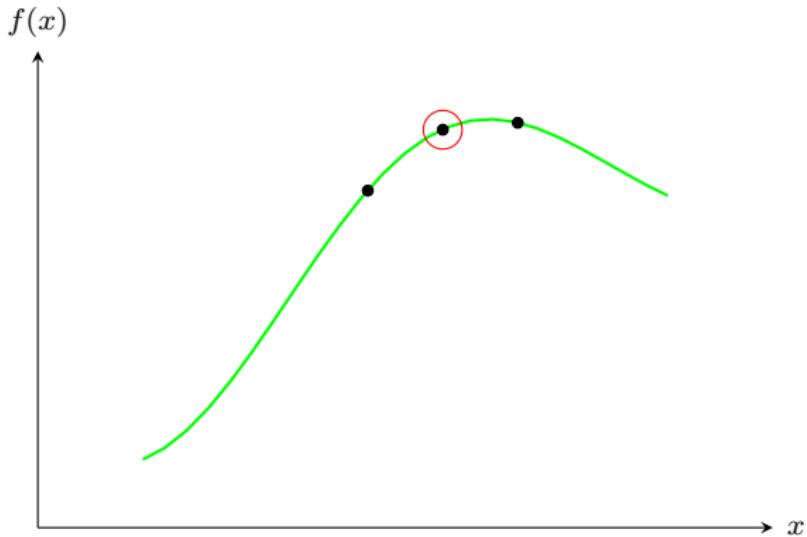
- mini-batch gradient descent is similar but with less “stochastic noise”

gradient computation

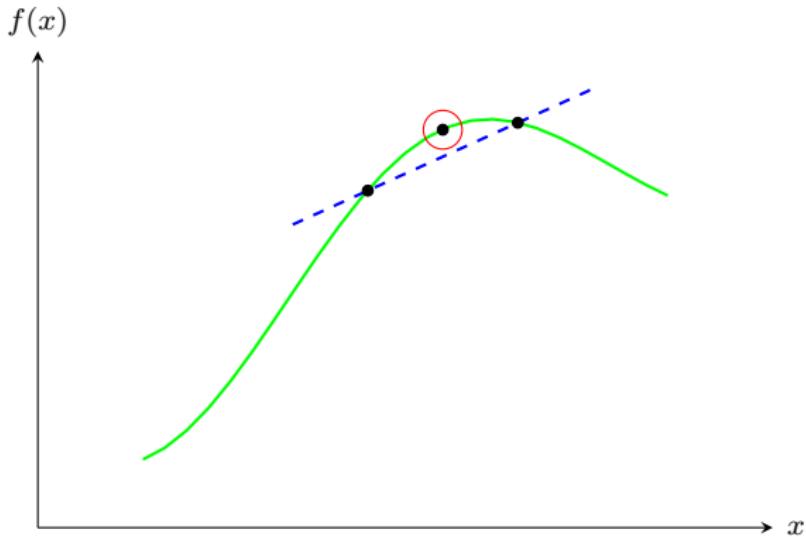
numerical approximation



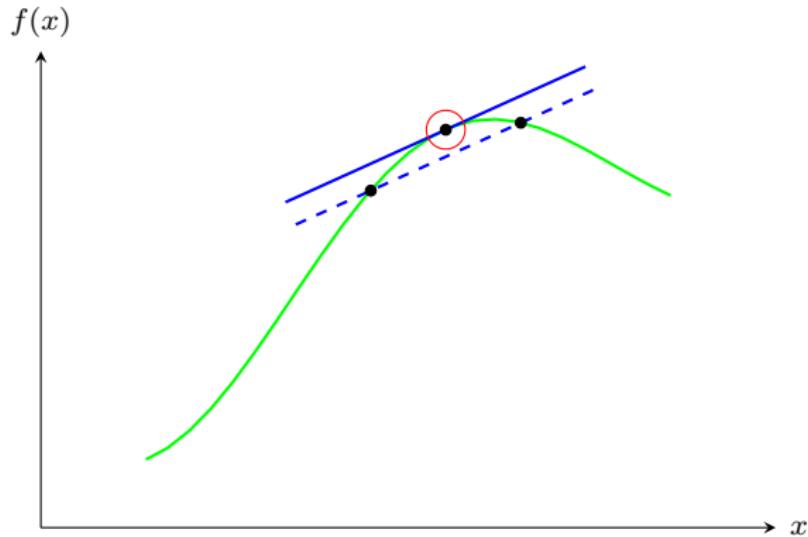
numerical approximation



numerical approximation



numerical approximation



$$\frac{df}{dx}(x) \approx \frac{f(x + \delta) - f(x - \delta)}{2\delta}$$

numerical approximation

- given $f : \mathbb{R}^p \rightarrow \mathbb{R}$, its gradient is the vector function

$$\nabla f := \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_p} \right)$$

- each partial derivative $\frac{\partial f}{\partial x_i}$ can be approximated at \mathbf{x} by the **symmetric difference** formula

$$\Delta_i f(\mathbf{x}; \delta) := \frac{f(\mathbf{x} + \delta \mathbf{e}_i) - f(\mathbf{x} - \delta \mathbf{e}_i)}{2\delta}$$

for small $\delta > 0$, where \mathbf{e}_i the i standard basis vector of \mathbb{R}^m

- in practice, the smallest δ should be used that does not cause numerical issues, e.g. $\delta \in [10^{-10}, 10^{-5}]$ for double-precision arithmetic

numerical approximation

- given $f : \mathbb{R}^p \rightarrow \mathbb{R}$, its gradient is the vector function

$$\nabla f := \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_p} \right)$$

- each partial derivative $\frac{\partial f}{\partial x_i}$ can be approximated at \mathbf{x} by the **symmetric difference** formula

$$\Delta_i f(\mathbf{x}; \delta) := \frac{f(\mathbf{x} + \delta \mathbf{e}_i) - f(\mathbf{x} - \delta \mathbf{e}_i)}{2\delta}$$

for small $\delta > 0$, where \mathbf{e}_i the i standard basis vector of \mathbb{R}^m

- in practice, the smallest δ should be used that does not cause numerical issues, e.g. $\delta \in [10^{-10}, 10^{-5}]$ for double-precision arithmetic

numerical approximation

- given $f : \mathbb{R}^p \rightarrow \mathbb{R}$, its gradient is the vector function

$$\nabla f := \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_p} \right)$$

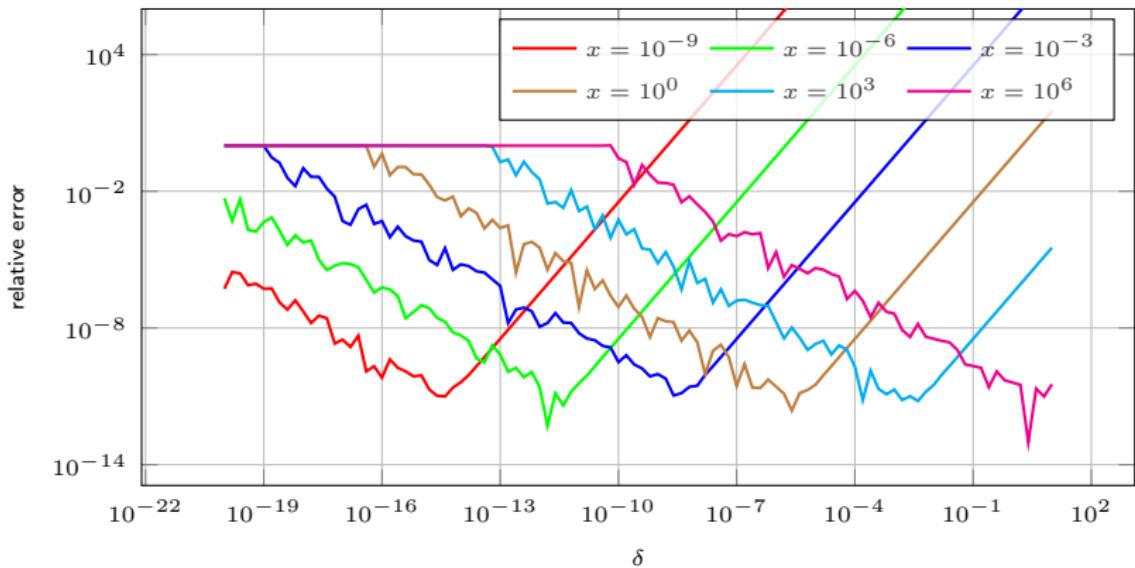
- each partial derivative $\frac{\partial f}{\partial x_i}$ can be approximated at \mathbf{x} by the **symmetric difference** formula

$$\Delta_i f(\mathbf{x}; \delta) := \frac{f(\mathbf{x} + \delta \mathbf{e}_i) - f(\mathbf{x} - \delta \mathbf{e}_i)}{2\delta}$$

for small $\delta > 0$, where \mathbf{e}_i the i standard basis vector of \mathbb{R}^m

- in practice, the smallest δ should be used that does not cause numerical issues, e.g. $\delta \in [10^{-10}, 10^{-5}]$ for double-precision arithmetic

example



- relative error for $f(x) = x^3$, $\nabla f(x) = 3x^2$

$$\frac{|\Delta f(x; \delta) - \nabla f(x)|}{\nabla f(x)}$$

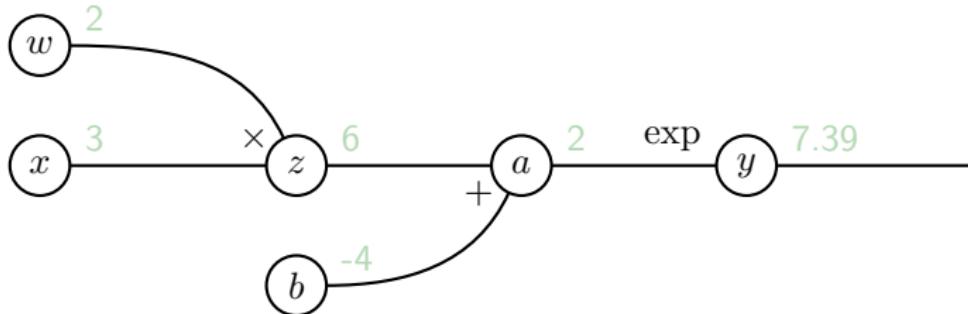
numerical vs. analytical

- apart from accuracy issues, the numerical approximation is impractical in high dimensions: one evaluation of Δf requires $2p$ evaluations of f , and dimension p is easily in the order of millions
- we turn to analytical computation of the gradient, which costs roughly as much as one evaluation of f
- but the numerical approximation always remains useful for double-checking

analytical computation

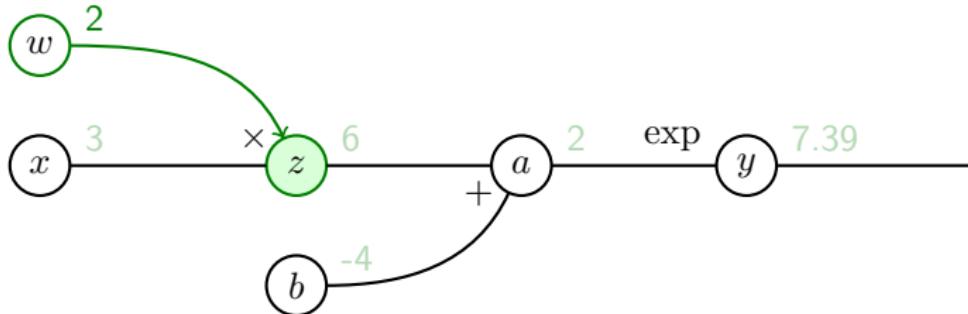
- all derivatives we care about are the derivatives of the error function with respect to the model parameters: the error function is **scalar** and we need its **gradient**
- we are going to write the error function as a composition of simpler functions, and use the **chain rule** to compute the gradient efficiently
- the error function can be as complex as a program with **control flow** statements
- each component function, called a **unit**, is assumed to be at least piecewise differentiable with a known formula for its derivative
- a unit may be a vector function, so we need **Jacobian matrices** in general, not just gradients

example



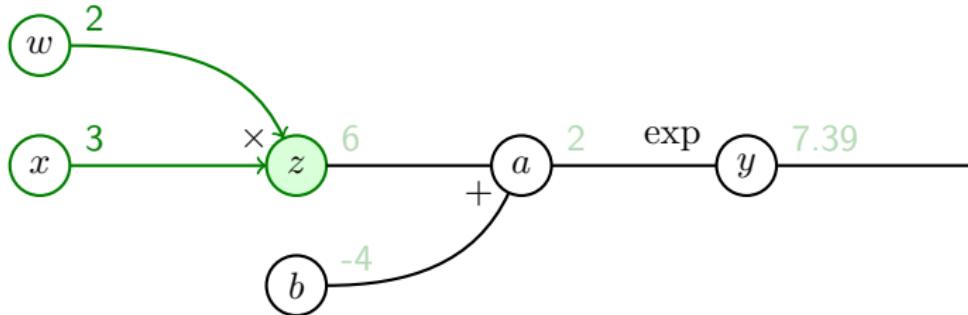
- $y = e^{wx+b}$ is broken down into $y = e^a$, $a = z + b$, $z = wx$
- we seek $\frac{da}{dw}$, $\frac{da}{dx}$, $\frac{da}{db}$ for $(w, x, b) = (2, 3, -1)$
- $\frac{dy}{dy} = 1$, $\frac{dy}{da} = e^a = 7.39$, $\frac{dy}{db} = \frac{dy}{da} = \frac{dy}{da} = 7.39$,
 $\frac{dy}{dz} = \frac{dy}{da} \frac{da}{dz} = \frac{dy}{da} = 7.39$, $\frac{dy}{dx} = \frac{dy}{dz} \frac{dz}{dx} = \frac{dy}{dz} w = 22.17$,
 $\frac{dy}{dw} = \frac{dy}{dz} \frac{dz}{dw} = \frac{dy}{dz} x = 14.78$

example



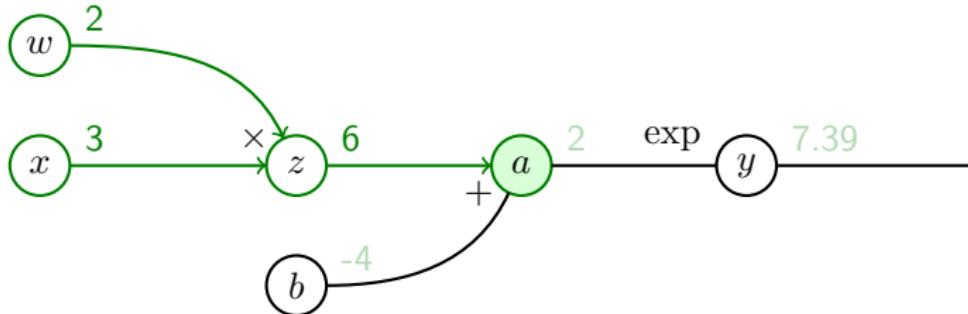
- $y = e^{wx+b}$ is broken down into $y = e^a$, $a = z + b$, $z = wx$
- we seek $\frac{da}{dw}$, $\frac{da}{dx}$, $\frac{da}{db}$ for $(w, x, b) = (2, 3, -1)$
- $\frac{dy}{dy} = 1$, $\frac{dy}{da} = e^a = 7.39$, $\frac{dy}{db} = \frac{dy}{da} = \frac{dy}{da} = 7.39$,
 $\frac{dy}{dz} = \frac{dy}{da} \frac{da}{dz} = \frac{dy}{da} = 7.39$, $\frac{dy}{dx} = \frac{dy}{dz} \frac{dz}{dx} = \frac{dy}{dz} w = 22.17$,
 $\frac{dy}{dw} = \frac{dy}{dz} \frac{dz}{dw} = \frac{dy}{dz} x = 14.78$

example



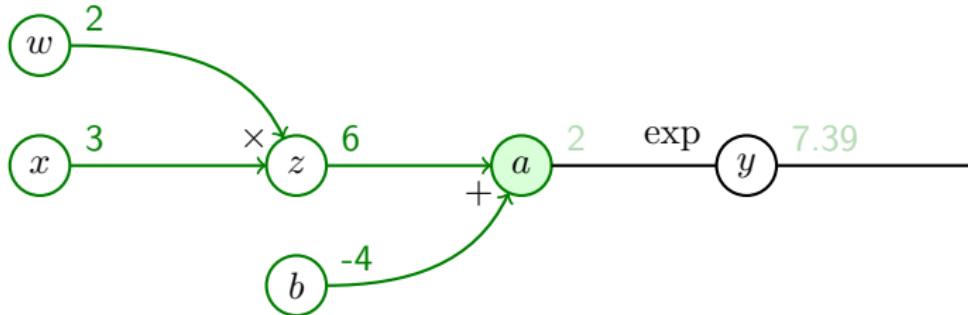
- $y = e^{wx+b}$ is broken down into $y = e^a$, $a = z + b$, $z = wx$
- we seek $\frac{da}{dw}$, $\frac{da}{dx}$, $\frac{da}{db}$ for $(w, x, b) = (2, 3, -1)$
- $\frac{dy}{dy} = 1$, $\frac{dy}{da} = e^a = 7.39$, $\frac{dy}{db} = \frac{dy}{da} = \frac{dy}{da} = 7.39$,
 $\frac{dy}{dz} = \frac{dy}{da} \frac{da}{dz} = \frac{dy}{da} = 7.39$, $\frac{dy}{dx} = \frac{dy}{dz} \frac{dz}{dx} = \frac{dy}{dz} w = 22.17$,
 $\frac{dy}{dw} = \frac{dy}{dz} \frac{dz}{dw} = \frac{dy}{dz} x = 14.78$

example



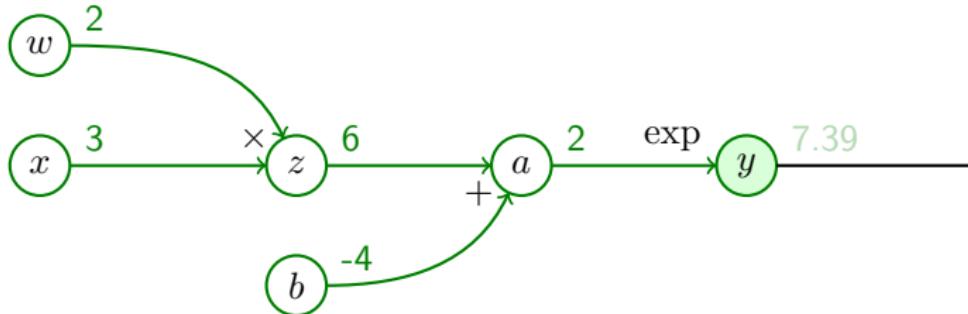
- $y = e^{wx+b}$ is broken down into $y = e^a$, $a = z + b$, $z = wx$
- we seek $\frac{da}{dw}$, $\frac{da}{dx}$, $\frac{da}{db}$ for $(w, x, b) = (2, 3, -1)$
- $\frac{dy}{dy} = 1$, $\frac{dy}{da} = e^a = 7.39$, $\frac{dy}{db} = \frac{dy}{da} = \frac{dy}{da} = 7.39$,
 $\frac{dy}{dz} = \frac{dy}{da} \frac{da}{dz} = \frac{dy}{da} = 7.39$, $\frac{dy}{dx} = \frac{dy}{dz} \frac{dz}{dx} = \frac{dy}{dz} w = 22.17$,
 $\frac{dy}{dw} = \frac{dy}{dz} \frac{dz}{dw} = \frac{dy}{dz} x = 14.78$

example



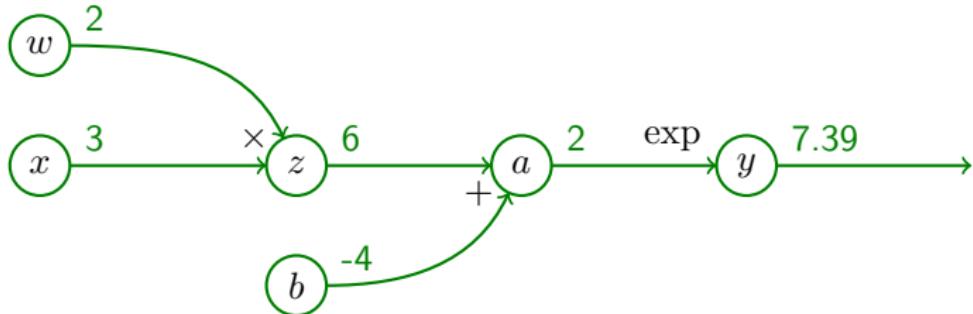
- $y = e^{wx+b}$ is broken down into $y = e^a$, $a = z + b$, $z = wx$
- we seek $\frac{da}{dw}$, $\frac{da}{dx}$, $\frac{da}{db}$ for $(w, x, b) = (2, 3, -1)$
- $\frac{dy}{dy} = 1$, $\frac{dy}{da} = e^a = 7.39$, $\frac{dy}{db} = \frac{dy}{da} = \frac{dy}{da} = 7.39$,
 $\frac{dy}{dz} = \frac{dy}{da} \frac{da}{dz} = \frac{dy}{da} = 7.39$, $\frac{dy}{dx} = \frac{dy}{dz} \frac{dz}{dx} = \frac{dy}{dz} w = 22.17$,
 $\frac{dy}{dw} = \frac{dy}{dz} \frac{dz}{dw} = \frac{dy}{dz} x = 14.78$

example



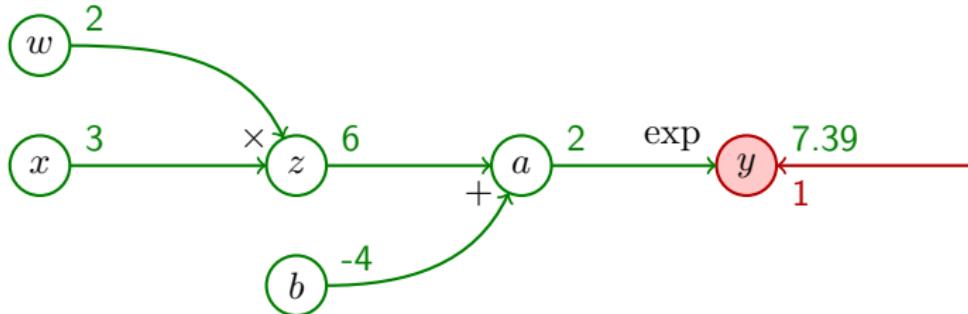
- $y = e^{wx+b}$ is broken down into $y = e^a$, $a = z + b$, $z = wx$
- we seek $\frac{da}{dw}$, $\frac{da}{dx}$, $\frac{da}{db}$ for $(w, x, b) = (2, 3, -1)$
- $\frac{dy}{dy} = 1$, $\frac{dy}{da} = e^a = 7.39$, $\frac{dy}{db} = \frac{dy}{da} = \frac{dy}{da} = 7.39$,
 $\frac{dy}{dz} = \frac{dy}{da} \frac{da}{dz} = \frac{dy}{da} = 7.39$, $\frac{dy}{dx} = \frac{dy}{dz} \frac{dz}{dx} = \frac{dy}{dz} w = 22.17$,
 $\frac{dy}{dw} = \frac{dy}{dz} \frac{dz}{dw} = \frac{dy}{dz} x = 14.78$

example



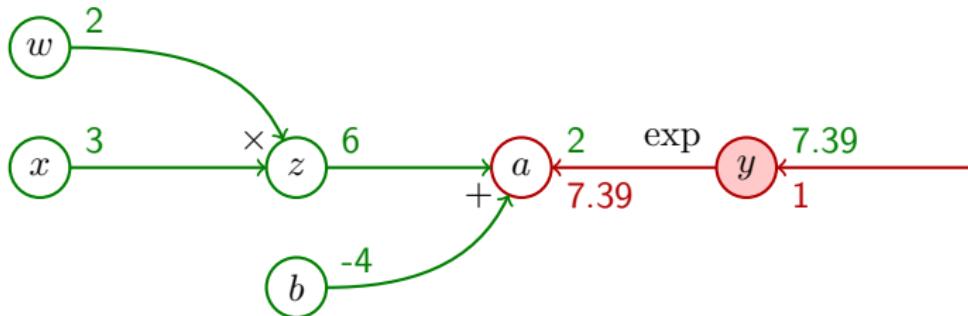
- $y = e^{wx+b}$ is broken down into $y = e^a$, $a = z + b$, $z = wx$
- we seek $\frac{da}{dw}$, $\frac{da}{dx}$, $\frac{da}{db}$ for $(w, x, b) = (2, 3, -1)$
- $\frac{dy}{dy} = 1$, $\frac{dy}{da} = e^a = 7.39$, $\frac{dy}{db} = \frac{dy}{da} = \frac{dy}{da} = 7.39$,
 $\frac{dy}{dz} = \frac{dy}{da} \frac{da}{dz} = \frac{dy}{da} = 7.39$, $\frac{dy}{dx} = \frac{dy}{dz} \frac{dz}{dx} = \frac{dy}{dz} w = 22.17$,
 $\frac{dy}{dw} = \frac{dy}{dz} \frac{dz}{dw} = \frac{dy}{dz} x = 14.78$

example



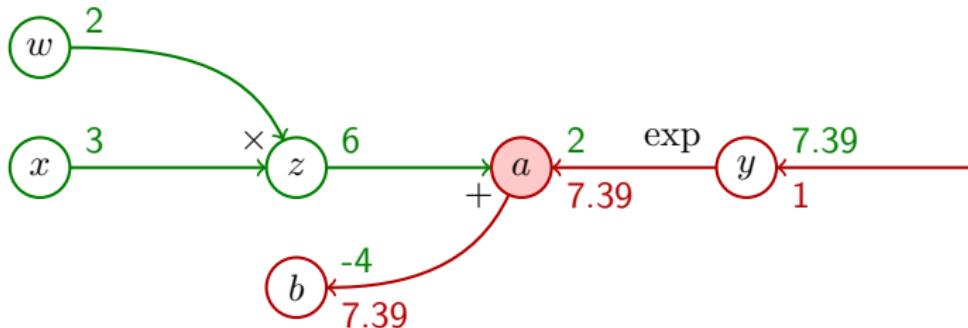
- $y = e^{wx+b}$ is broken down into $y = e^a$, $a = z + b$, $z = wx$
- we seek $\frac{da}{dw}$, $\frac{da}{dx}$, $\frac{da}{db}$ for $(w, x, b) = (2, 3, -1)$
- $\frac{dy}{dy} = 1$, $\frac{dy}{da} = e^a = 7.39$, $\frac{dy}{db} = \frac{dy}{da} = \frac{dy}{da} = 7.39$,
 $\frac{dy}{dz} = \frac{dy}{da} \frac{da}{dz} = \frac{dy}{da} = 7.39$, $\frac{dy}{dx} = \frac{dy}{dz} \frac{dz}{dx} = \frac{dy}{dz} w = 22.17$,
 $\frac{dy}{dw} = \frac{dy}{dz} \frac{dz}{dw} = \frac{dy}{dz} x = 14.78$

example



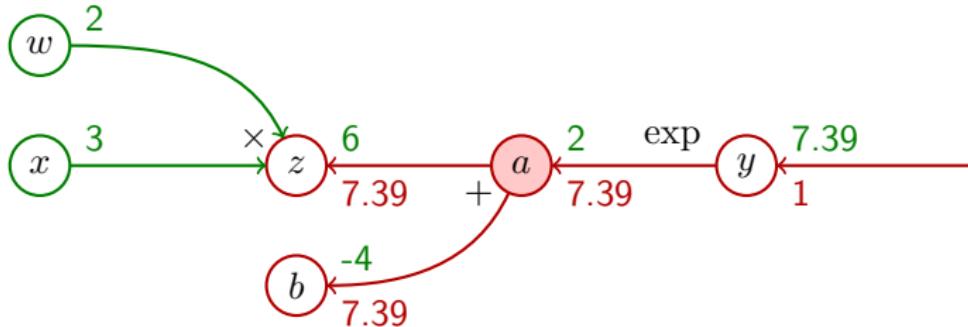
- $y = e^{wx+b}$ is broken down into $y = e^a$, $a = z + b$, $z = wx$
- we seek $\frac{da}{dw}$, $\frac{da}{dx}$, $\frac{da}{db}$ for $(w, x, b) = (2, 3, -1)$
- $\frac{dy}{dy} = 1$, $\frac{dy}{da} = e^a = 7.39$, $\frac{dy}{db} = \frac{dy}{da} = \frac{dy}{da} = 7.39$,
 $\frac{dy}{dz} = \frac{dy}{da} \frac{da}{dz} = \frac{dy}{da} = 7.39$, $\frac{dy}{dx} = \frac{dy}{dz} \frac{dz}{dx} = \frac{dy}{dz} w = 22.17$,
 $\frac{dy}{dw} = \frac{dy}{dz} \frac{dz}{dw} = \frac{dy}{dz} x = 14.78$

example



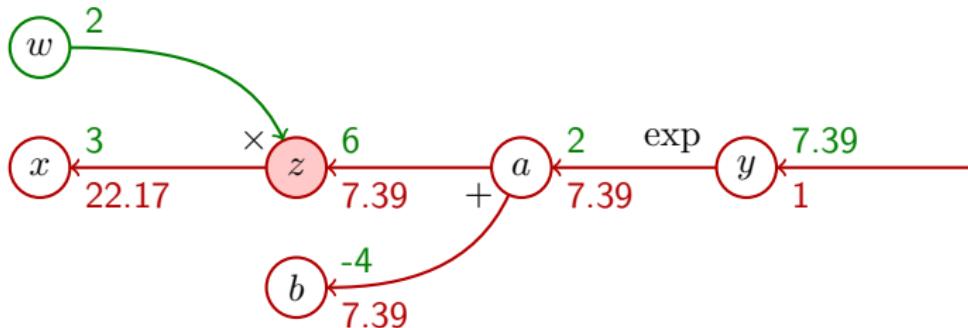
- $y = e^{wx+b}$ is broken down into $y = e^a$, $a = z + b$, $z = wx$
- we seek $\frac{da}{dw}$, $\frac{da}{dx}$, $\frac{da}{db}$ for $(w, x, b) = (2, 3, -1)$
- $\frac{dy}{dy} = 1$, $\frac{dy}{da} = e^a = 7.39$, $\frac{dy}{db} = \frac{dy}{da} = \frac{dy}{da} = 7.39$,
 $\frac{dy}{dz} = \frac{dy}{da} \frac{da}{dz} = \frac{dy}{da} = 7.39$, $\frac{dy}{dx} = \frac{dy}{dz} \frac{dz}{dx} = \frac{dy}{dz} w = 22.17$,
 $\frac{dy}{dw} = \frac{dy}{dz} \frac{dz}{dw} = \frac{dy}{dz} x = 14.78$

example



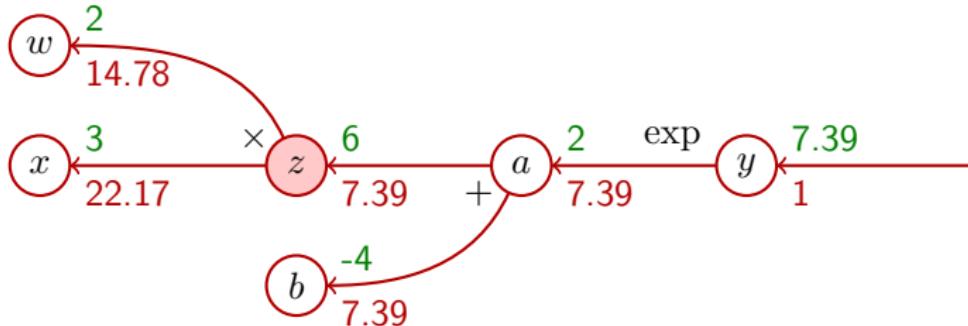
- $y = e^{wx+b}$ is broken down into $y = e^a$, $a = z + b$, $z = wx$
- we seek $\frac{da}{dw}$, $\frac{da}{dx}$, $\frac{da}{db}$ for $(w, x, b) = (2, 3, -1)$
- $\frac{dy}{dy} = 1$, $\frac{dy}{da} = e^a = 7.39$, $\frac{dy}{db} = \frac{dy}{da} = \frac{dy}{da} = 7.39$,
 $\frac{dy}{dz} = \frac{dy}{da} \frac{da}{dz} = \frac{dy}{da} = 7.39$, $\frac{dy}{dx} = \frac{dy}{dz} \frac{dz}{dx} = \frac{dy}{dz} w = 22.17$,
 $\frac{dy}{dw} = \frac{dy}{dz} \frac{dz}{dw} = \frac{dy}{dz} x = 14.78$

example



- $y = e^{wx+b}$ is broken down into $y = e^a$, $a = z + b$, $z = wx$
- we seek $\frac{da}{dw}$, $\frac{da}{dx}$, $\frac{da}{db}$ for $(w, x, b) = (2, 3, -1)$
- $\frac{dy}{dy} = 1$, $\frac{dy}{da} = e^a = 7.39$, $\frac{dy}{db} = \frac{dy}{da} = \frac{dy}{da} = 7.39$,
 $\frac{dy}{dz} = \frac{dy}{da} \frac{da}{dz} = \frac{dy}{da} = 7.39$, $\frac{dy}{dx} = \frac{dy}{dz} \frac{dz}{dx} = \frac{dy}{dz} w = 22.17$,
 $\frac{dy}{dw} = \frac{dy}{dz} \frac{dz}{dw} = \frac{dy}{dz} x = 14.78$

example



- $y = e^{wx+b}$ is broken down into $y = e^a$, $a = z + b$, $z = wx$
- we seek $\frac{da}{dw}$, $\frac{da}{dx}$, $\frac{da}{db}$ for $(w, x, b) = (2, 3, -1)$
- $\frac{dy}{dy} = 1$, $\frac{dy}{da} = e^a = 7.39$, $\frac{dy}{db} = \frac{dy}{da} = \frac{dy}{da} = 7.39$,
 $\frac{dy}{dz} = \frac{dy}{da} \frac{da}{dz} = \frac{dy}{da} = 7.39$, $\frac{dy}{dx} = \frac{dy}{dz} \frac{dz}{dx} = \frac{dy}{dz} w = 22.17$,
 $\frac{dy}{dw} = \frac{dy}{dz} \frac{dz}{dw} = \frac{dy}{dz} x = 14.78$

vector functions: derivative

- a function $f : \mathbb{R}^p \rightarrow \mathbb{R}^q$ is **differentiable** at \mathbf{x} if there is a $q \times p$ matrix A such that

$$\frac{f(\mathbf{x} + \mathbf{h}) - f(\mathbf{x}) - A \cdot \mathbf{h}}{|\mathbf{h}|} \rightarrow \mathbf{0}$$

as $\mathbf{h} \rightarrow \mathbf{0}$; matrix A is the **derivative** of f at \mathbf{x} , denoted as $Df(\mathbf{x})$

- if

$$f(\mathbf{x}) = A\mathbf{x}$$

then

$$Df(\mathbf{x}) = A$$

vector functions: derivative

- a function $f : \mathbb{R}^p \rightarrow \mathbb{R}^q$ is **differentiable** at \mathbf{x} if there is a $q \times p$ matrix A such that

$$\frac{f(\mathbf{x} + \mathbf{h}) - f(\mathbf{x}) - A \cdot \mathbf{h}}{|\mathbf{h}|} \rightarrow \mathbf{0}$$

as $\mathbf{h} \rightarrow \mathbf{0}$; matrix A is the **derivative** of f at \mathbf{x} , denoted as $Df(\mathbf{x})$

- if

$$f(\mathbf{x}) = A\mathbf{x}$$

then

$$Df(\mathbf{x}) = A$$

vector functions: derivative vs. Jacobian

- given $f = (f_1, \dots, f_q) : \mathbb{R}^p \rightarrow \mathbb{R}^q$ whose its partial derivatives exist at \mathbf{x} , and $\mathbf{y} = f(\mathbf{x})$, its **Jacobian matrix** at \mathbf{x} can be written as

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \frac{\partial f}{\partial \mathbf{x}} := \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_p} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_q}{\partial x_1} & \cdots & \frac{\partial f_q}{\partial x_p} \end{pmatrix}$$

- if f is differentiable at \mathbf{x} , its **derivative** at \mathbf{x} is

$$Df(\mathbf{x}) := \begin{pmatrix} D_1 f_1 & \cdots & D_p f_1 \\ \vdots & \ddots & \vdots \\ D_1 f_q & \cdots & D_p f_q \end{pmatrix} (\mathbf{x})$$

- if f is differentiable at \mathbf{x} , the derivative $Df(\mathbf{x})$ equals the Jacobian at \mathbf{x} ; but the Jacobian may exist without any derivative defined

vector functions: derivative vs. Jacobian

- given $f = (f_1, \dots, f_q) : \mathbb{R}^p \rightarrow \mathbb{R}^q$ whose its partial derivatives exist at \mathbf{x} , and $\mathbf{y} = f(\mathbf{x})$, its **Jacobian matrix** at \mathbf{x} can be written as

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \frac{\partial f}{\partial \mathbf{x}} := \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_p} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_q}{\partial x_1} & \cdots & \frac{\partial f_q}{\partial x_p} \end{pmatrix}$$

- if f is differentiable at \mathbf{x} , its **derivative** at \mathbf{x} is

$$Df(\mathbf{x}) := \begin{pmatrix} D_1 f_1 & \cdots & D_p f_1 \\ \vdots & \ddots & \vdots \\ D_1 f_q & \cdots & D_p f_q \end{pmatrix} (\mathbf{x})$$

- if f is differentiable at \mathbf{x} , the derivative $Df(\mathbf{x})$ equals the Jacobian at \mathbf{x} ; but the Jacobian may exist without any derivative defined

scalar functions: derivative vs. gradient

- the gradient of a scalar $f : \mathbb{R}^p \rightarrow \mathbb{R}$ with respect to an input vector \mathbf{x} is a **column** vector in \mathbb{R}^p , the same size as \mathbf{x}

$$\nabla f := \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_p} \right)$$

- in contrast, the derivative is an $1 \times p$ **row** vector

$$Df(\mathbf{x}) := \begin{pmatrix} \frac{\partial f}{\partial x_1} & \dots & \frac{\partial f}{\partial x_p} \end{pmatrix} = (\nabla f)^\top$$

- the following analysis uses derivatives/Jacobians, so we will **transpose** them to make them compatible with \mathbf{x}

chain rule

- if $f : \mathbb{R}^p \rightarrow \mathbb{R}^q$ is differentiable at \mathbf{x} and $g : \mathbb{R}^q \rightarrow \mathbb{R}^r$ is differentiable at $\mathbf{y} = f(\mathbf{x})$, then $g \circ f : \mathbb{R}^p \rightarrow \mathbb{R}^r$ is differentiable at \mathbf{x} and

$$D(g \circ f)(\mathbf{x}) = Dg(\mathbf{y}) \cdot Df(\mathbf{x})$$

where \cdot denotes matrix multiplication

- how to use it:

$$\frac{\partial z}{\partial x_1} = \frac{\partial z}{\partial x_2} \cdot \frac{\partial x_2}{\partial x_1}$$



chain rule

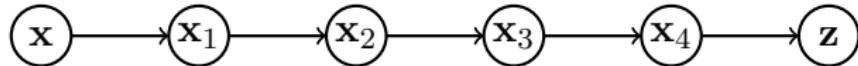
- if $f : \mathbb{R}^p \rightarrow \mathbb{R}^q$ is differentiable at \mathbf{x} and $g : \mathbb{R}^q \rightarrow \mathbb{R}^r$ is differentiable at $\mathbf{y} = f(\mathbf{x})$, then $g \circ f : \mathbb{R}^p \rightarrow \mathbb{R}^r$ is differentiable at \mathbf{x} and

$$D(g \circ f)(\mathbf{x}) = Dg(\mathbf{y}) \cdot Df(\mathbf{x})$$

where \cdot denotes matrix multiplication

- how to use it:

$$\frac{\partial \mathbf{z}}{\partial \mathbf{x}_1} = \frac{\partial \mathbf{z}}{\partial \mathbf{x}_2} \cdot \frac{\partial \mathbf{x}_2}{\partial \mathbf{x}_1}$$



chain rule

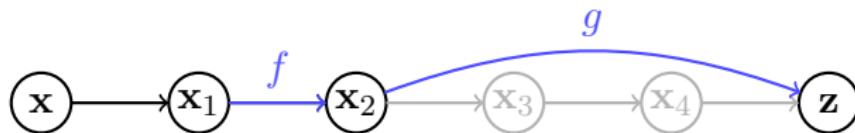
- if $f : \mathbb{R}^p \rightarrow \mathbb{R}^q$ is differentiable at \mathbf{x} and $g : \mathbb{R}^q \rightarrow \mathbb{R}^r$ is differentiable at $\mathbf{y} = f(\mathbf{x})$, then $g \circ f : \mathbb{R}^p \rightarrow \mathbb{R}^r$ is differentiable at \mathbf{x} and

$$D(g \circ f)(\mathbf{x}) = Dg(\mathbf{y}) \cdot Df(\mathbf{x})$$

where \cdot denotes matrix multiplication

- how to use it:

$$\frac{\partial \mathbf{z}}{\partial \mathbf{x}_1} = \frac{\partial \mathbf{z}}{\partial \mathbf{x}_2} \cdot \frac{\partial \mathbf{x}_2}{\partial \mathbf{x}_1}$$



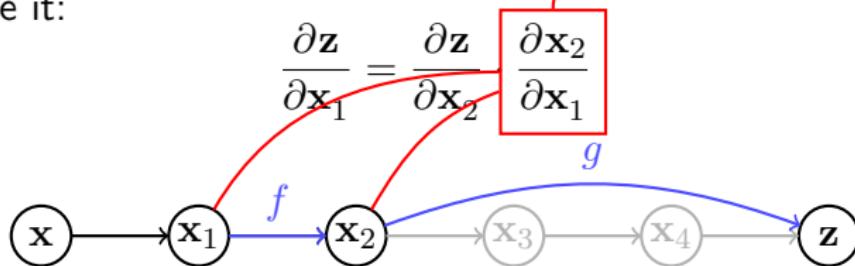
chain rule

- if $f : \mathbb{R}^p \rightarrow \mathbb{R}^q$ is differentiable at \mathbf{x} and $g : \mathbb{R}^q \rightarrow \mathbb{R}^r$ is differentiable at $\mathbf{y} = f(\mathbf{x})$, then $g \circ f : \mathbb{R}^p \rightarrow \mathbb{R}^r$ is differentiable at \mathbf{x} and

$$D(g \circ f)(\mathbf{x}) = Dg(\mathbf{y}) \cdot Df(\mathbf{x})$$

where \cdot denotes matrix multiplication

- how to use it:



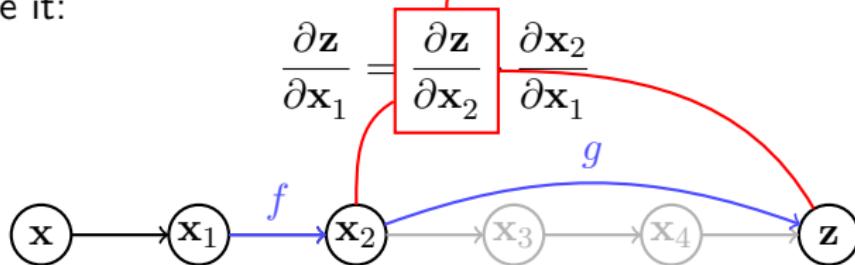
chain rule

- if $f : \mathbb{R}^p \rightarrow \mathbb{R}^q$ is differentiable at \mathbf{x} and $g : \mathbb{R}^q \rightarrow \mathbb{R}^r$ is differentiable at $\mathbf{y} = f(\mathbf{x})$, then $g \circ f : \mathbb{R}^p \rightarrow \mathbb{R}^r$ is differentiable at \mathbf{x} and

$$D(g \circ f)(\mathbf{x}) = \boxed{Dg(\mathbf{y})} \cdot Df(\mathbf{x})$$

where \cdot denotes matrix multiplication

- how to use it:



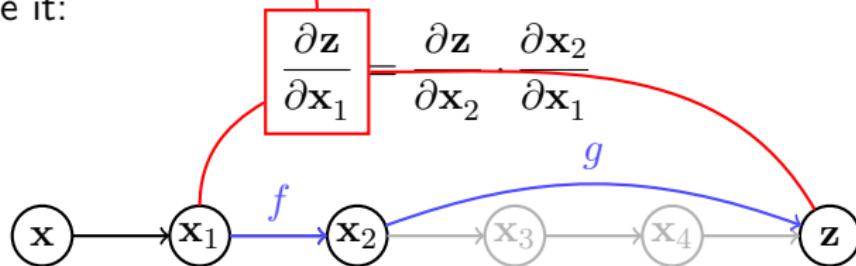
chain rule

- if $f : \mathbb{R}^p \rightarrow \mathbb{R}^q$ is differentiable at \mathbf{x} and $g : \mathbb{R}^q \rightarrow \mathbb{R}^r$ is differentiable at $\mathbf{y} = f(\mathbf{x})$, then $g \circ f : \mathbb{R}^p \rightarrow \mathbb{R}^r$ is differentiable at \mathbf{x} and

$$D(g \circ f)(\mathbf{x}) = Dg(\mathbf{y}) \cdot Df(\mathbf{x})$$

where \cdot denotes matrix multiplication

- how to use it:



chain rule

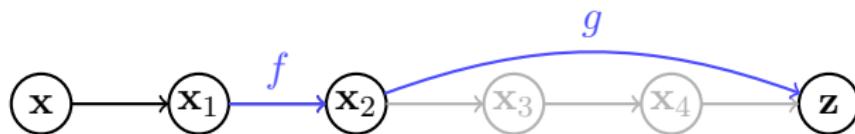
- if $f : \mathbb{R}^p \rightarrow \mathbb{R}^q$ is differentiable at \mathbf{x} and $g : \mathbb{R}^q \rightarrow \mathbb{R}^r$ is differentiable at $\mathbf{y} = f(\mathbf{x})$, then $g \circ f : \mathbb{R}^p \rightarrow \mathbb{R}^r$ is differentiable at \mathbf{x} and

$$D(g \circ f)(\mathbf{x}) = Dg(\mathbf{y}) \cdot Df(\mathbf{x})$$

where \cdot denotes matrix multiplication

- how to use it:

$$\frac{\partial \mathbf{z}}{\partial \mathbf{x}_1} = \frac{\partial \mathbf{z}}{\partial \mathbf{x}_2} \cdot \frac{\partial \mathbf{x}_2}{\partial \mathbf{x}_1}$$



- now, for all i , let us call the partial derivatives

$$d\mathbf{x}_i^\top := \frac{\partial \mathbf{z}}{\partial \mathbf{x}_i}$$

chain rule

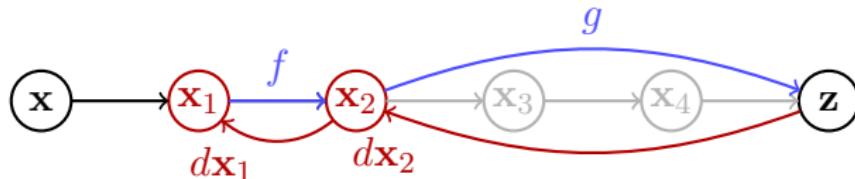
- if $f : \mathbb{R}^p \rightarrow \mathbb{R}^q$ is differentiable at \mathbf{x} and $g : \mathbb{R}^q \rightarrow \mathbb{R}^r$ is differentiable at $\mathbf{y} = f(\mathbf{x})$, then $g \circ f : \mathbb{R}^p \rightarrow \mathbb{R}^r$ is differentiable at \mathbf{x} and

$$D(g \circ f)(\mathbf{x}) = Dg(\mathbf{y}) \cdot Df(\mathbf{x})$$

where \cdot denotes matrix multiplication

- how to use it:

$$\frac{\partial \mathbf{z}}{\partial \mathbf{x}_1} = \frac{\partial \mathbf{z}}{\partial \mathbf{x}_2} \cdot \frac{\partial \mathbf{x}_2}{\partial \mathbf{x}_1}$$

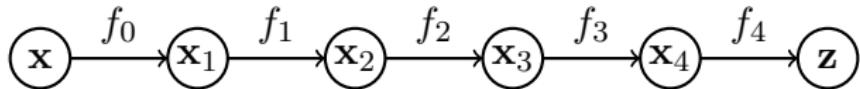


- then, we are **back-propagating** from $d\mathbf{x}_2$ to $d\mathbf{x}_1$

$$d\mathbf{x}_1^\top = d\mathbf{x}_2^\top \cdot Df(\mathbf{x}_1)$$

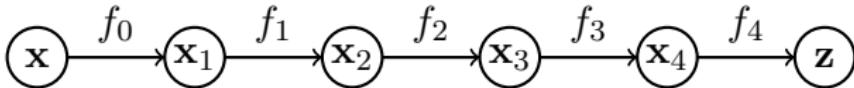
chaining

- let $f = f_4 \circ f_3 \circ f_2 \circ f_1 \circ f_0$ and $\mathbf{z} = f(\mathbf{x})$



chaining

- let $f = f_4 \circ f_3 \circ f_2 \circ f_1 \circ f_0$ and $\mathbf{z} = f(\mathbf{x})$

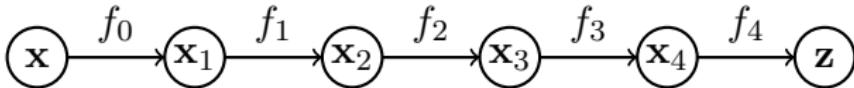


- we apply the chain rule

$$\begin{aligned}\frac{\partial \mathbf{z}}{\partial \mathbf{x}} &= Df(\mathbf{x}) = D(f_4 \circ f_3 \circ f_2 \circ f_1)(\mathbf{x}_1) \cdot Df_0(\mathbf{x}) \\ &= D(f_4 \circ f_3 \circ f_2)(\mathbf{x}_2) \cdot Df_1(\mathbf{x}_1) \cdot Df_0(\mathbf{x}) \\ &= D(f_4 \circ f_3)(\mathbf{x}_3) \cdot Df_2(\mathbf{x}_2) \cdot Df_1(\mathbf{x}_1) \cdot Df_0(\mathbf{x}) \\ &= Df_4(\mathbf{x}_4) \cdot Df_3(\mathbf{x}_3) \cdot Df_2(\mathbf{x}_2) \cdot Df_1(\mathbf{x}_1) \cdot Df_0(\mathbf{x}) \\ &= d\mathbf{x}_4^\top \cdot Df_3(\mathbf{x}_3) \cdot Df_2(\mathbf{x}_2) \cdot Df_1(\mathbf{x}_1) \cdot Df_0(\mathbf{x}) \\ &= d\mathbf{x}_3^\top \cdot Df_2(\mathbf{x}_2) \cdot Df_1(\mathbf{x}_1) \cdot Df_0(\mathbf{x}) \\ &= d\mathbf{x}_2^\top \cdot Df_1(\mathbf{x}_1) \cdot Df_0(\mathbf{x}) \\ &= d\mathbf{x}_1^\top \cdot Df_0(\mathbf{x}) \\ &= d\mathbf{x}^\top\end{aligned}$$

chaining

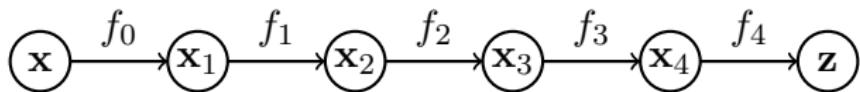
- let $f = f_4 \circ f_3 \circ f_2 \circ f_1 \circ f_0$ and $\mathbf{z} = f(\mathbf{x})$



- we apply the chain rule, then collect back into factors $d\mathbf{x}_i$

$$\begin{aligned}\frac{\partial \mathbf{z}}{\partial \mathbf{x}} &= Df(\mathbf{x}) = D(f_4 \circ f_3 \circ f_2 \circ f_1)(\mathbf{x}_1) \cdot Df_0(\mathbf{x}) \\&= D(f_4 \circ f_3 \circ f_2)(\mathbf{x}_2) \cdot Df_1(\mathbf{x}_1) \cdot Df_0(\mathbf{x}) \\&= D(f_4 \circ f_3)(\mathbf{x}_3) \cdot Df_2(\mathbf{x}_2) \cdot Df_1(\mathbf{x}_1) \cdot Df_0(\mathbf{x}) \\&= Df_4(\mathbf{x}_4) \cdot Df_3(\mathbf{x}_3) \cdot Df_2(\mathbf{x}_2) \cdot Df_1(\mathbf{x}_1) \cdot Df_0(\mathbf{x}) \\&= d\mathbf{x}_4^\top \cdot Df_3(\mathbf{x}_3) \cdot Df_2(\mathbf{x}_2) \cdot Df_1(\mathbf{x}_1) \cdot Df_0(\mathbf{x}) \\&= d\mathbf{x}_3^\top \cdot Df_2(\mathbf{x}_2) \cdot Df_1(\mathbf{x}_1) \cdot Df_0(\mathbf{x}) \\&= d\mathbf{x}_2^\top \cdot Df_1(\mathbf{x}_1) \cdot Df_0(\mathbf{x}) \\&= d\mathbf{x}_1^\top \cdot Df_0(\mathbf{x}) \\&= d\mathbf{x}^\top\end{aligned}$$

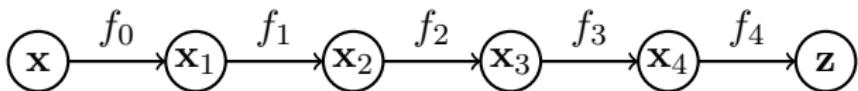
back-propagation



forward pass

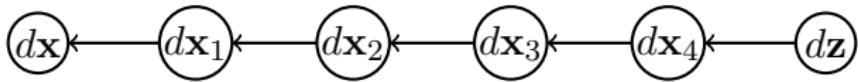
$$\mathbf{x}_1 = f_0(\mathbf{x}) \quad \mathbf{x}_2 = f_1(\mathbf{x}_1) \quad \mathbf{x}_3 = f_2(\mathbf{x}_2) \quad \mathbf{x}_4 = f_3(\mathbf{x}_3) \quad \mathbf{z} = f_4(\mathbf{x}_4)$$

back-propagation



forward pass

$$\mathbf{x}_1 = f_0(\mathbf{x}) \quad \mathbf{x}_2 = f_1(\mathbf{x}_1) \quad \mathbf{x}_3 = f_2(\mathbf{x}_2) \quad \mathbf{x}_4 = f_3(\mathbf{x}_3) \quad \mathbf{z} = f_4(\mathbf{x}_4)$$



backward pass

$$d\mathbf{z}^\top = I \qquad d\mathbf{x}_4^\top = d\mathbf{z}^\top \cdot Df_4(\mathbf{x}_4) \quad d\mathbf{x}_3^\top = d\mathbf{x}_4^\top \cdot Df_3(\mathbf{x}_3)$$
$$d\mathbf{x}_2^\top = d\mathbf{x}_3^\top \cdot Df_3(\mathbf{x}_3) \quad d\mathbf{x}_1^\top = d\mathbf{x}_2^\top \cdot Df_1(\mathbf{x}_1) \quad d\mathbf{x}^\top = d\mathbf{x}_1^\top \cdot Df_0(\mathbf{x})$$

back-propagation is dynamic programming

- we need to store all the x_i that we compute in the forward pass before the backward pass begins
- the dx_i can be computed on the fly in reverse order on a chain, but may need to be all stored on a general network structure
- that's exactly what we do in **dynamic programming**: break the problem down into a collection of smaller, overlapping subproblems, store their solutions and save computation time at the expense of a (hopefully) modest expenditure in storage space
- as in all dynamic programming problems, there is a **bottom-up** approach that we have just described, and a **top-down** approach coming out of the **recursive** formulation through **memoization**; this can be useful if we are looking for the derivative with respect to only few parameters

back-propagation is dynamic programming

- we need to store all the x_i that we compute in the forward pass before the backward pass begins
- the dx_i can be computed on the fly in reverse order on a chain, but may need to be all stored on a general network structure
- that's exactly what we do in **dynamic programming**: break the problem down into a collection of smaller, overlapping subproblems, store their solutions and save computation time at the expense of a (hopefully) modest expenditure in storage space
- as in all dynamic programming problems, there is a **bottom-up** approach that we have just described, and a **top-down** approach coming out of the **recursive** formulation through **memoization**; this can be useful if we are looking for the derivative with respect to only few parameters

partial derivatives

- in the following, for any vector \mathbf{x} appearing in our function, we will use the symbol

$$d\mathbf{x}^\top := \frac{\partial}{\partial \mathbf{x}}$$

for the partial derivative operator of any quantity with respect to \mathbf{x}

- in practice, we will apply this to the quantity we want to optimize, *i.e.* the error
- the error gradient will consist of the partial derivatives with respect to the model parameters, but we still need to compute partial derivatives with respect to all variables appearing in back-propagation

partial derivatives

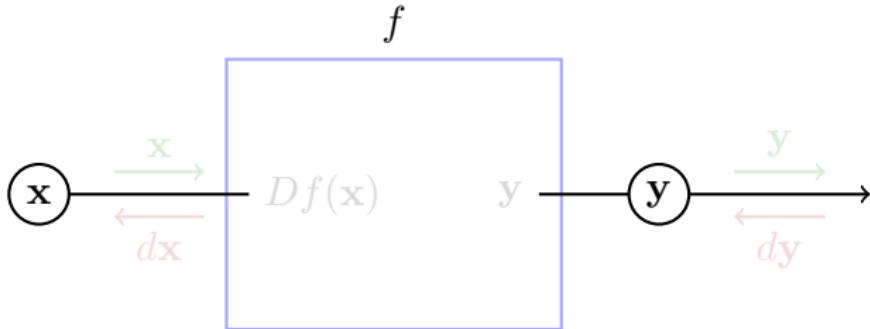
- in the following, for any vector \mathbf{x} appearing in our function, we will use the symbol

$$d\mathbf{x}^\top := \frac{\partial}{\partial \mathbf{x}}$$

for the partial derivative operator of any quantity with respect to \mathbf{x}

- in practice, we will apply this to the quantity we want to optimize, *i.e.* the error
- the error gradient will consist of the partial derivatives with respect to the model parameters, but we still need to compute partial derivatives with respect to all variables appearing in back-propagation

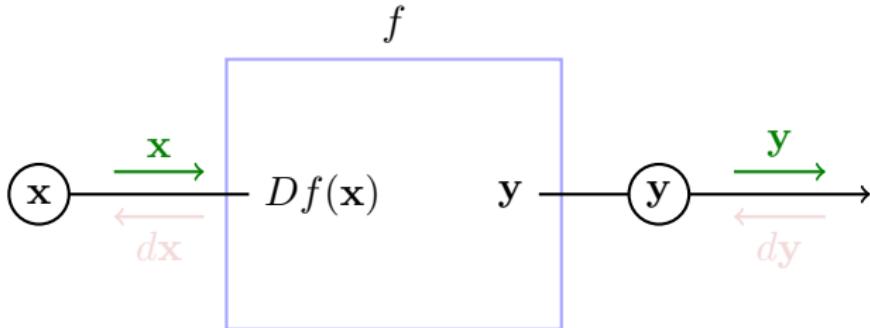
nodes



- to every variable y is associated a node with the function f that produces it, from input variable x
- given x , derivative $Df(x)$ is “stored”, and output y is computed and flows forward
- given dy , partial derivative dx is computed and flows backward

$$d\mathbf{x}^\top = d\mathbf{y}^\top \cdot Df(\mathbf{x}) \quad \text{or} \quad \frac{\partial}{\partial \mathbf{x}} = \frac{\partial}{\partial \mathbf{y}} \cdot \frac{\partial \mathbf{y}}{\partial \mathbf{x}}$$

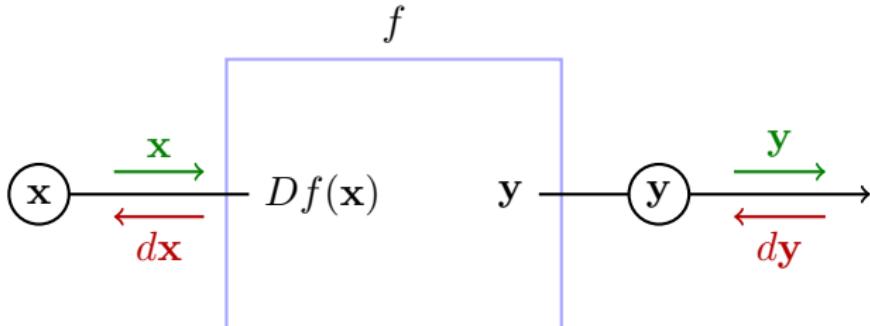
nodes



- to every variable y is associated a node with the function f that produces it, from input variable x
- given x , derivative $Df(x)$ is “stored”, and output y is computed and flows forward
- given dy , partial derivative dx is computed and flows backward

$$d\mathbf{x}^\top = d\mathbf{y}^\top \cdot Df(\mathbf{x}) \quad \text{or} \quad \frac{\partial}{\partial \mathbf{x}} = \frac{\partial}{\partial \mathbf{y}} \cdot \frac{\partial \mathbf{y}}{\partial \mathbf{x}}$$

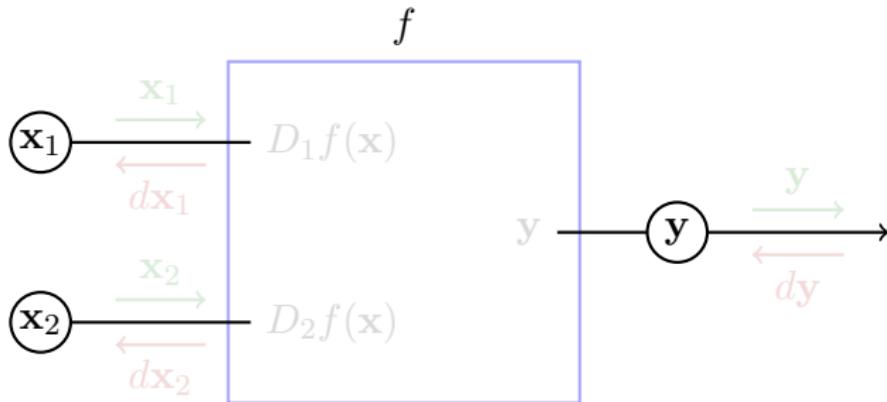
nodes



- to every variable y is associated a node with the function f that produces it, from input variable x
- given x , derivative $Df(x)$ is “stored”, and output y is computed and flows forward
- given dy , partial derivative dx is computed and flows backward

$$d\mathbf{x}^\top = d\mathbf{y}^\top \cdot Df(\mathbf{x}) \quad \text{or} \quad \frac{\partial}{\partial \mathbf{x}} = \frac{\partial}{\partial \mathbf{y}} \cdot \frac{\partial \mathbf{y}}{\partial \mathbf{x}}$$

splitting the input



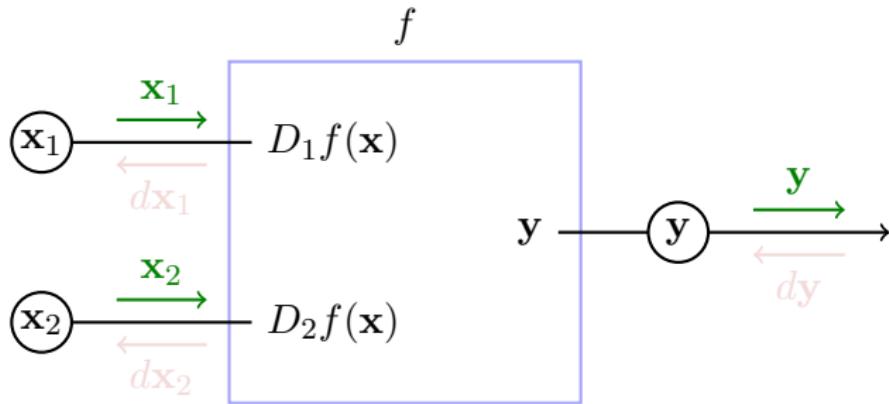
- we split input vector \mathbf{x} into subvectors as $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$
- then, the derivative consists of blocks stacked horizontally

$$Df(\mathbf{x}) = (D_1 f \ D_2 f)(\mathbf{x}) \quad \text{or} \quad \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial \mathbf{y}}{\partial \mathbf{x}_1} & \frac{\partial \mathbf{y}}{\partial \mathbf{x}_2} \end{pmatrix}$$

- $d\mathbf{x}$ is also split as $d\mathbf{x} = (d\mathbf{x}_1, d\mathbf{x}_2)$ and $d\mathbf{x} = d\mathbf{y} \cdot Df(\mathbf{x})$ becomes

$$d\mathbf{x}_i^\top = d\mathbf{y}^\top \cdot D_i f(\mathbf{x}) \quad \text{or} \quad \frac{\partial}{\partial \mathbf{x}_i} = \frac{\partial}{\partial \mathbf{y}} \cdot \frac{\partial \mathbf{y}}{\partial \mathbf{x}_i}$$

splitting the input



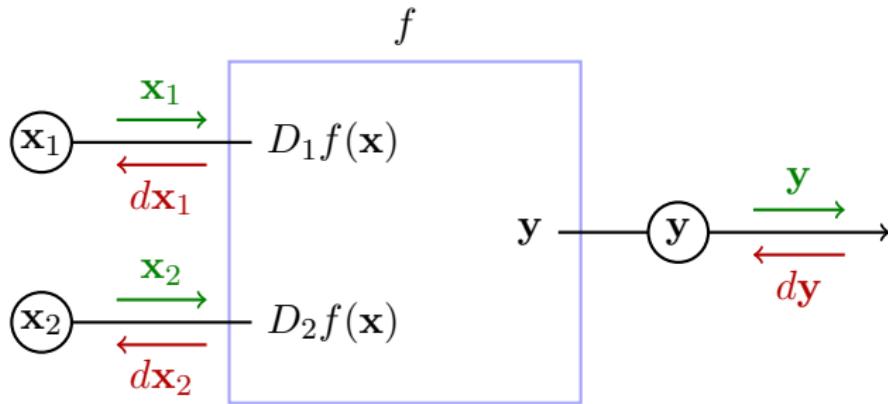
- we split input vector \mathbf{x} into subvectors as $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$
- then, the derivative consists of blocks stacked horizontally

$$Df(\mathbf{x}) = (D_1 f \ D_2 f)(\mathbf{x}) \quad \text{or} \quad \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial \mathbf{y}}{\partial \mathbf{x}_1} & \frac{\partial \mathbf{y}}{\partial \mathbf{x}_2} \end{pmatrix}$$

- $d\mathbf{x}$ is also split as $d\mathbf{x} = (dx_1, dx_2)$ and $d\mathbf{x} = dy \cdot Df(\mathbf{x})$ becomes

$$d\mathbf{x}_i^\top = dy^\top \cdot D_i f(\mathbf{x}) \quad \text{or} \quad \frac{\partial}{\partial \mathbf{x}_i} = \frac{\partial}{\partial \mathbf{y}} \cdot \frac{\partial \mathbf{y}}{\partial \mathbf{x}_i}$$

splitting the input



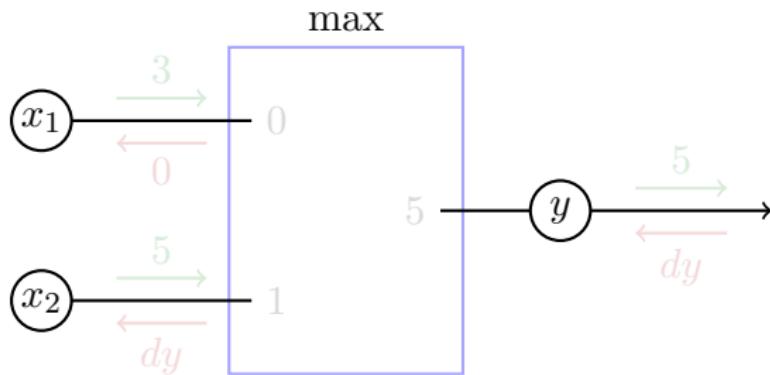
- we split input vector \mathbf{x} into subvectors as $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$
- then, the derivative consists of blocks stacked horizontally

$$Df(\mathbf{x}) = (D_1 f \ D_2 f)(\mathbf{x}) \quad \text{or} \quad \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial \mathbf{y}}{\partial \mathbf{x}_1} & \frac{\partial \mathbf{y}}{\partial \mathbf{x}_2} \end{pmatrix}$$

- $d\mathbf{x}$ is also split as $d\mathbf{x} = (d\mathbf{x}_1, d\mathbf{x}_2)$ and $d\mathbf{x} = d\mathbf{y} \cdot Df(\mathbf{x})$ becomes

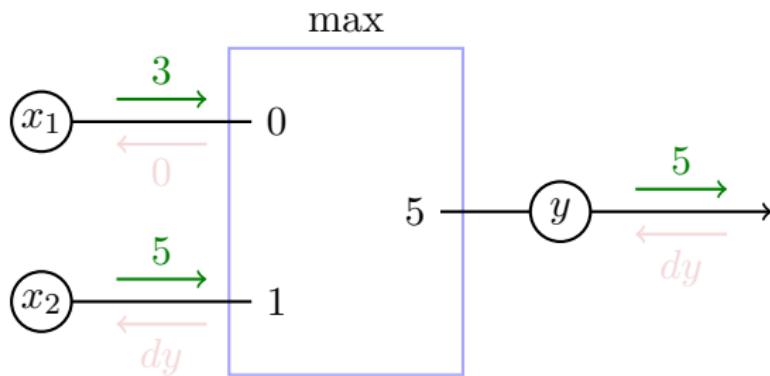
$$d\mathbf{x}_i^\top = d\mathbf{y}^\top \cdot D_i f(\mathbf{x}) \quad \text{or} \quad \frac{\partial}{\partial \mathbf{x}_i} = \frac{\partial}{\partial \mathbf{y}} \cdot \frac{\partial \mathbf{y}}{\partial \mathbf{x}_i}$$

example: maximum



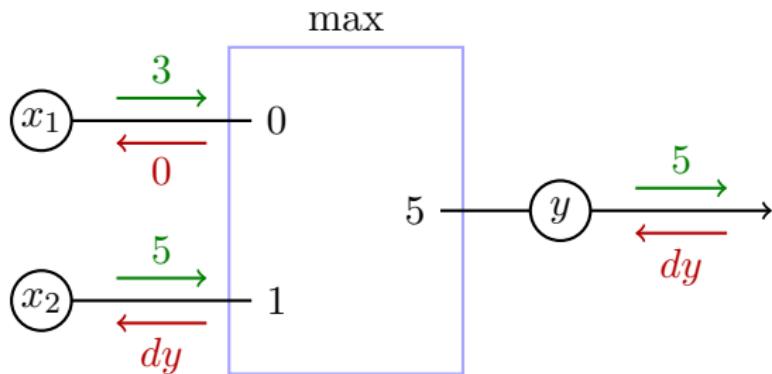
- if $f(x_1, x_2) = \max(x_1, x_2)$, then $D_i f(x_1, x_2) = \mathbb{1}[x_i = \max(x_1, x_2)]$
- and dy is **routed** into the branch of the maximum input

example: maximum



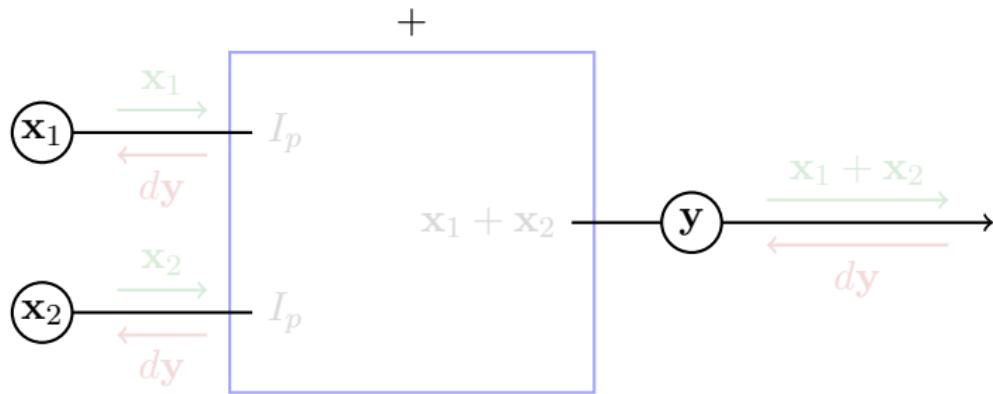
- if $f(x_1, x_2) = \max(x_1, x_2)$, then $D_i f(x_1, x_2) = \mathbb{1}[x_i = \max(x_1, x_2)]$
- and dy is **routed** into the branch of the maximum input

example: maximum



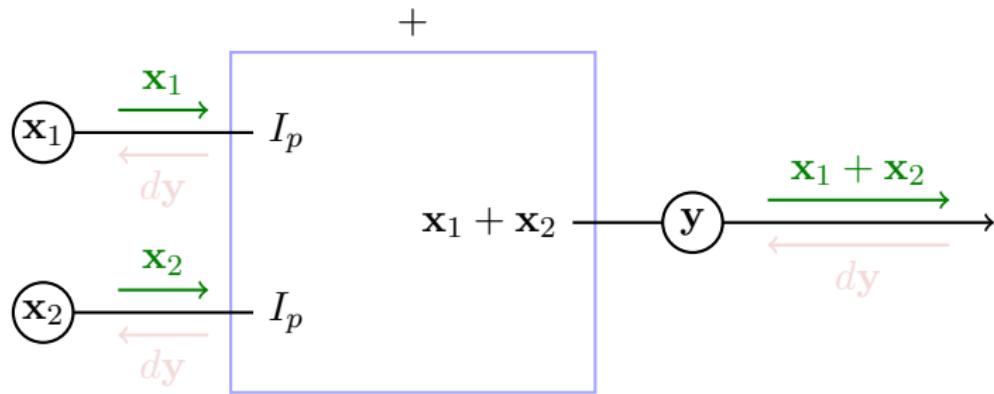
- if $f(x_1, x_2) = \max(x_1, x_2)$, then $D_i f(x_1, x_2) = \mathbb{1}[x_i = \max(x_1, x_2)]$
- and dy is **routed** into the branch of the maximum input

example: sum



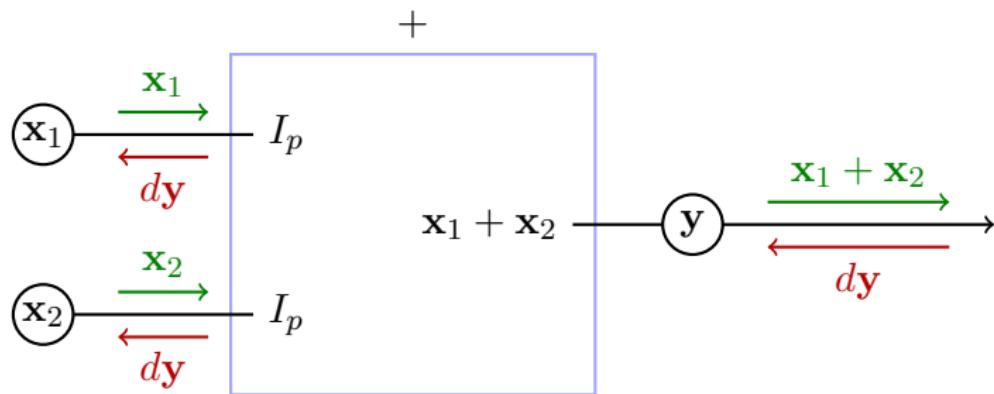
- if $f(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1 + \mathbf{x}_2$ and $\mathbf{x}_i \in \mathbb{R}^p$, then $D_i f(\mathbf{x}_1, \mathbf{x}_2) = I_p$
- and $d\mathbf{y}$ is **distributed** to both branches

example: sum



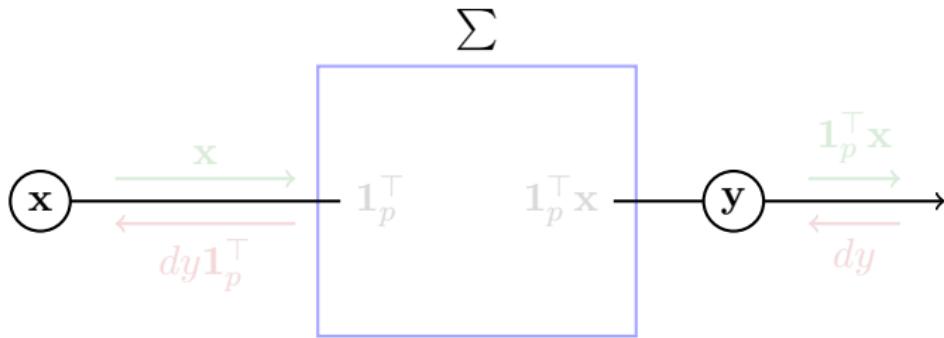
- if $f(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1 + \mathbf{x}_2$ and $\mathbf{x}_i \in \mathbb{R}^p$, then $D_i f(\mathbf{x}_1, \mathbf{x}_2) = I_p$
- and dy is distributed to both branches

example: sum



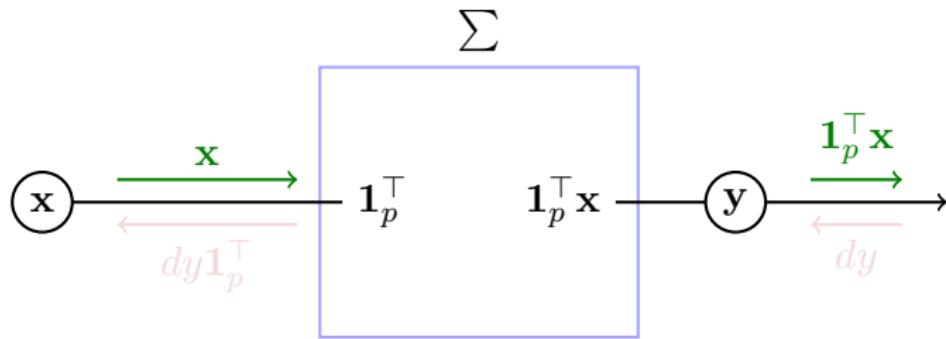
- if $f(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1 + \mathbf{x}_2$ and $\mathbf{x}_i \in \mathbb{R}^p$, then $D_i f(\mathbf{x}_1, \mathbf{x}_2) = I_p$
- and dy is distributed to both branches

example: vector sum



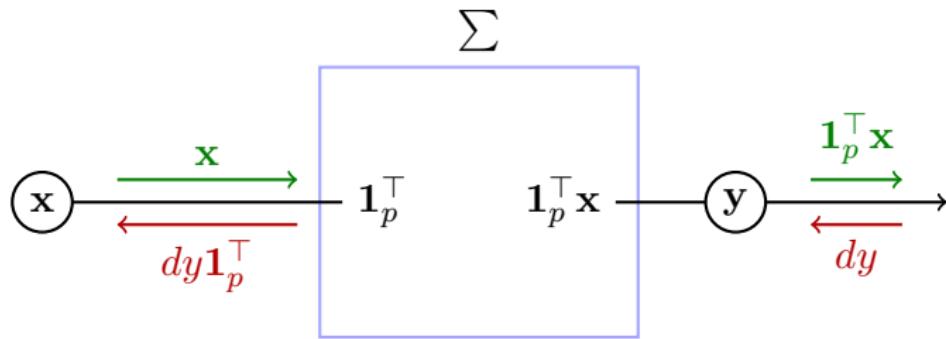
- if $f(\mathbf{x}) = \mathbf{1}_p^\top \mathbf{x} = \sum_{i=1}^p x_i$ and $\mathbf{x} \in \mathbb{R}^p$, then $Df(\mathbf{x}) = \mathbf{1}_p^\top$
- and dy is **distributed** to every element

example: vector sum



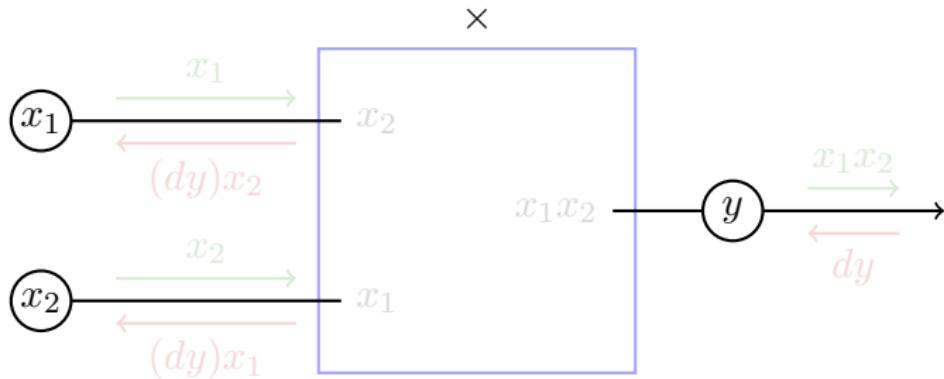
- if $f(\mathbf{x}) = \mathbf{1}_p^\top \mathbf{x} = \sum_{i=1}^p x_i$ and $\mathbf{x} \in \mathbb{R}^p$, then $Df(\mathbf{x}) = \mathbf{1}_p^\top$
- and $d\mathbf{y}$ is distributed to every element

example: vector sum



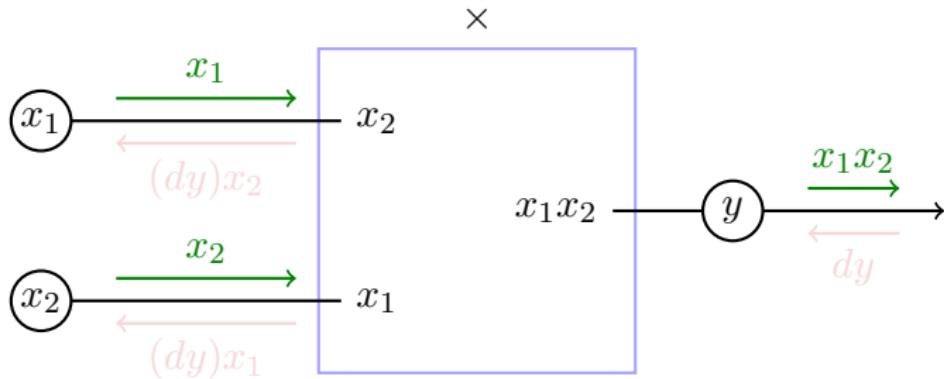
- if $f(\mathbf{x}) = \mathbf{1}_p^\top \mathbf{x} = \sum_{i=1}^p x_i$ and $\mathbf{x} \in \mathbb{R}^p$, then $Df(\mathbf{x}) = \mathbf{1}_p^\top$
- and $d\mathbf{y}$ is **distributed** to every element

example: product



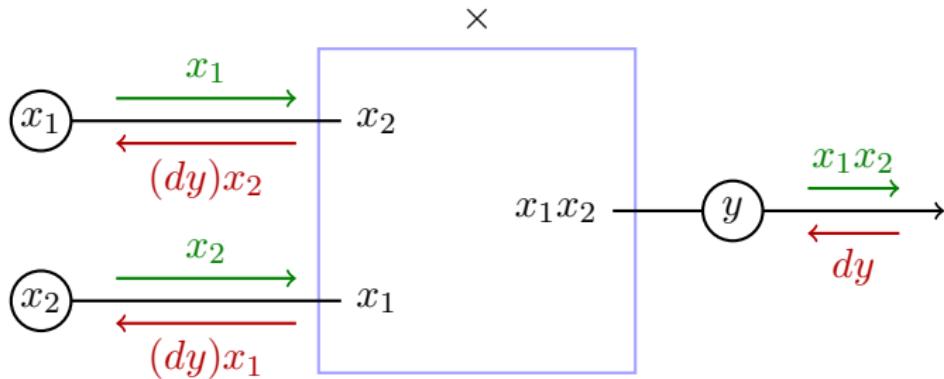
- if $f(x_1, x_2) = x_1 x_2$, then $D_1 f(x_1, x_2) = x_2$ and $D_2 f(x_1, x_2) = x_1$
- the derivative on each branch is multiplied by the input of the other

example: product



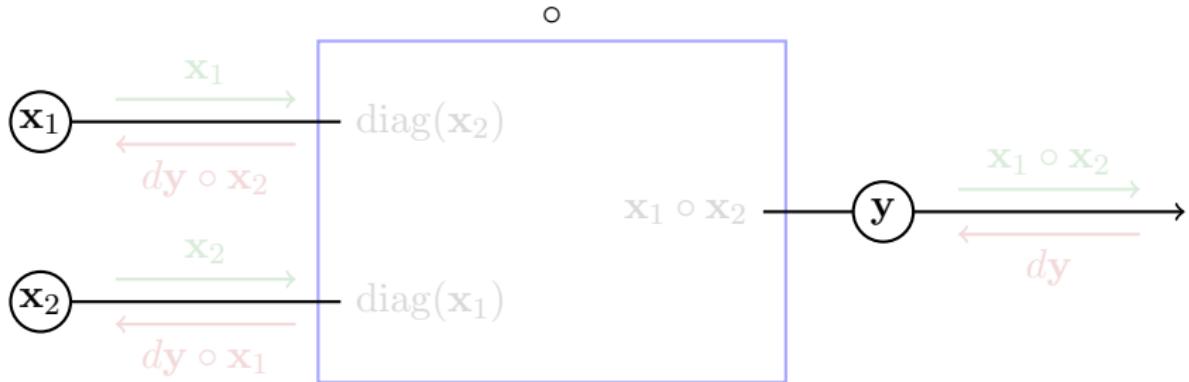
- if $f(x_1, x_2) = x_1x_2$, then $D_1f(x_1, x_2) = x_2$ and $D_2f(x_1, x_2) = x_1$
- the derivative on each branch is multiplied by the input of the other

example: product



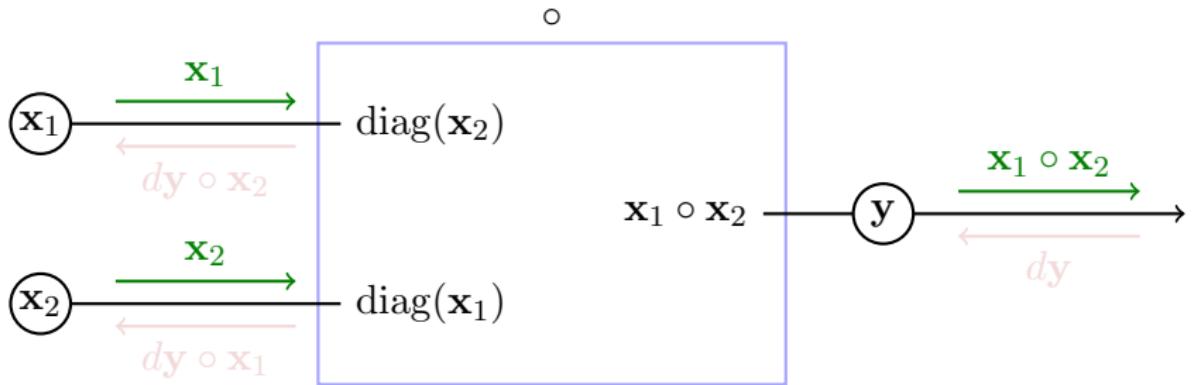
- if $f(x_1, x_2) = x_1x_2$, then $D_1f(x_1, x_2) = x_2$ and $D_2f(x_1, x_2) = x_1$
- the derivative on each branch is multiplied by the input of the other

example: Hadamard (element-wise) product



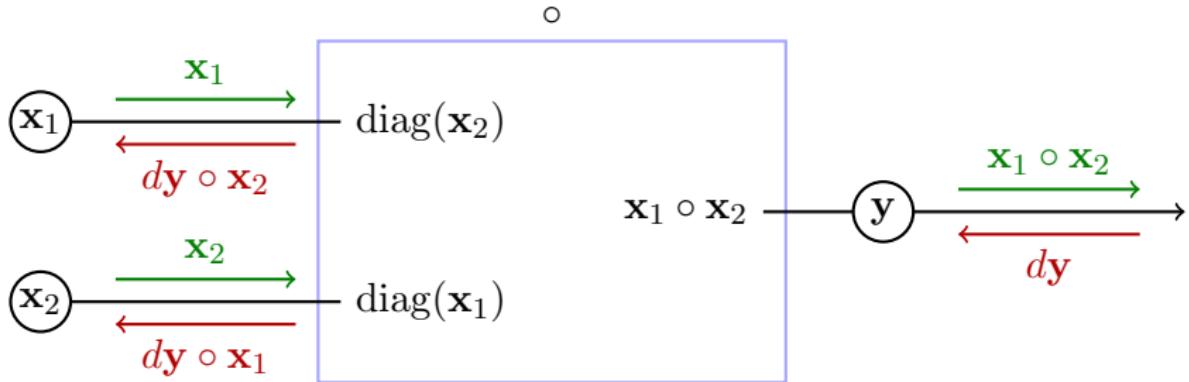
- if $f(x_1, x_2) = x_1 \circ x_2$, then $D_1 f(x_1, x_2) = \text{diag}(x_2)$ and $D_2 f(x_1, x_2) = \text{diag}(x_1)$
- the derivative on each branch is element-wise multiplied by the input of the other

example: Hadamard (element-wise) product



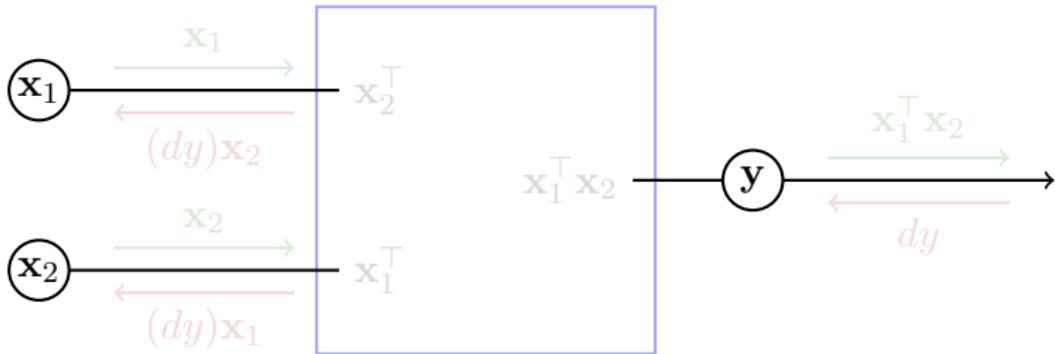
- if $f(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1 \circ \mathbf{x}_2$, then $D_1 f(\mathbf{x}_1, \mathbf{x}_2) = \text{diag}(\mathbf{x}_2)$ and $D_2 f(\mathbf{x}_1, \mathbf{x}_2) = \text{diag}(\mathbf{x}_1)$
- the derivative on each branch is element-wise multiplied by the input of the other

example: Hadamard (element-wise) product



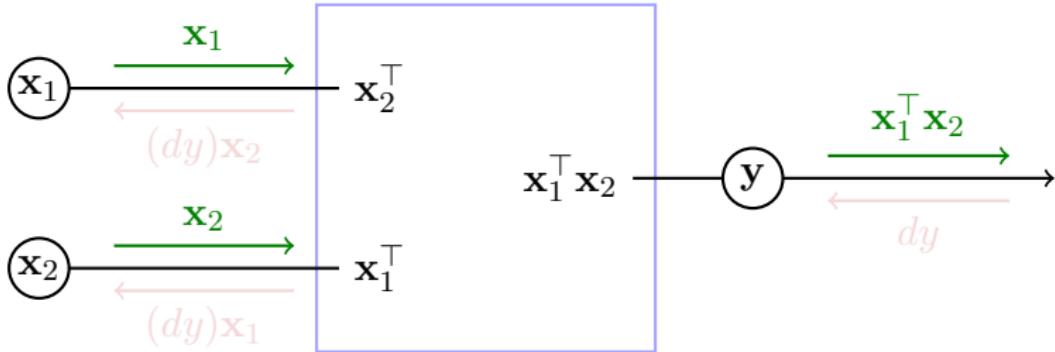
- if $f(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1 \circ \mathbf{x}_2$, then $D_1 f(\mathbf{x}_1, \mathbf{x}_2) = \text{diag}(\mathbf{x}_2)$ and $D_2 f(\mathbf{x}_1, \mathbf{x}_2) = \text{diag}(\mathbf{x}_1)$
- the derivative on each branch is element-wise multiplied by the input of the other

example: dot product



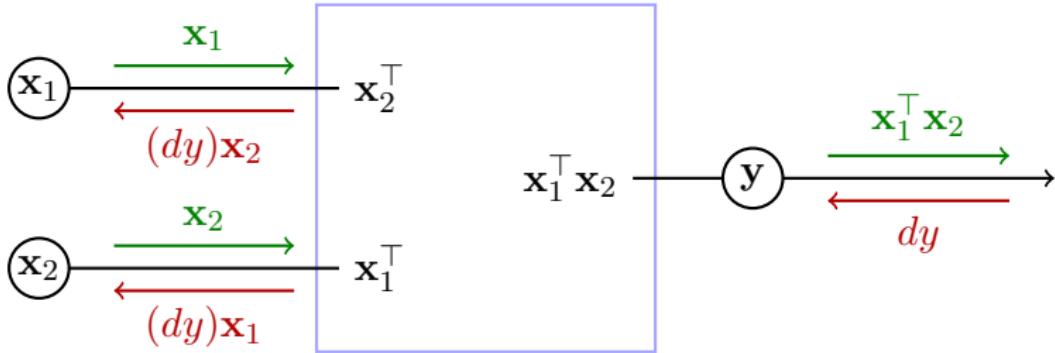
- if $f(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1 \cdot \mathbf{x}_2 = \mathbf{x}_1^\top \mathbf{x}_2$, then $D_1 f(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_2$ and $D_2 f(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1$
- the derivative on each branch is multiplied by the input of the other; this can be seen by composing an element-wise product with a vector sum

example: dot product



- if $f(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1 \cdot \mathbf{x}_2 = \mathbf{x}_1^\top \mathbf{x}_2$, then $D_1 f(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_2$ and $D_2 f(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1$
- the derivative on each branch is multiplied by the input of the other; this can be seen by composing an element-wise product with a vector sum

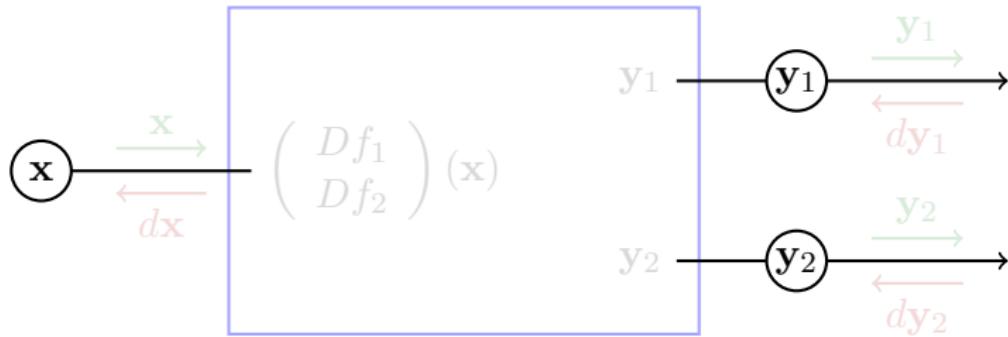
example: dot product



- if $f(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1 \cdot \mathbf{x}_2 = \mathbf{x}_1^\top \mathbf{x}_2$, then $D_1 f(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_2$ and $D_2 f(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1$
- the derivative on each branch is multiplied by the input of the other; this can be seen by composing an element-wise product with a vector sum

splitting the output

$$f = (f_1, f_2)$$



- we split output \mathbf{y} into subvectors as $\mathbf{y} = (\mathbf{y}_1, \mathbf{y}_2) = (f_1(\mathbf{x}), f_2(\mathbf{x}))$
- then, the derivative consists of blocks stacked vertically

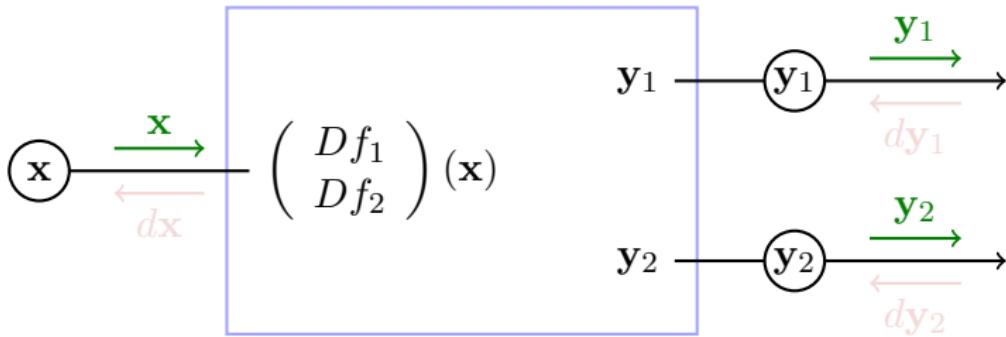
$$Df(\mathbf{x}) = (Df_1; Df_2)(\mathbf{x}) \quad \text{or} \quad \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \left(\frac{\partial \mathbf{y}_1}{\partial \mathbf{x}}; \frac{\partial \mathbf{y}_2}{\partial \mathbf{x}} \right)$$

- dy is also split as $dy = (dy_1, dy_2)$ and $d\mathbf{x} = dy \cdot Df(\mathbf{x})$ becomes

$$d\mathbf{x}^\top = \sum_i dy_i^\top \cdot Df_i(\mathbf{x}) \quad \text{or} \quad \frac{\partial}{\partial \mathbf{x}} = \sum_i \frac{\partial}{\partial \mathbf{y}_i} \cdot \frac{\partial \mathbf{y}_i}{\partial \mathbf{x}}$$

splitting the output

$$f = (f_1, f_2)$$



- we split output \mathbf{y} into subvectors as $\mathbf{y} = (\mathbf{y}_1, \mathbf{y}_2) = (f_1(\mathbf{x}), f_2(\mathbf{x}))$
- then, the derivative consists of blocks stacked vertically

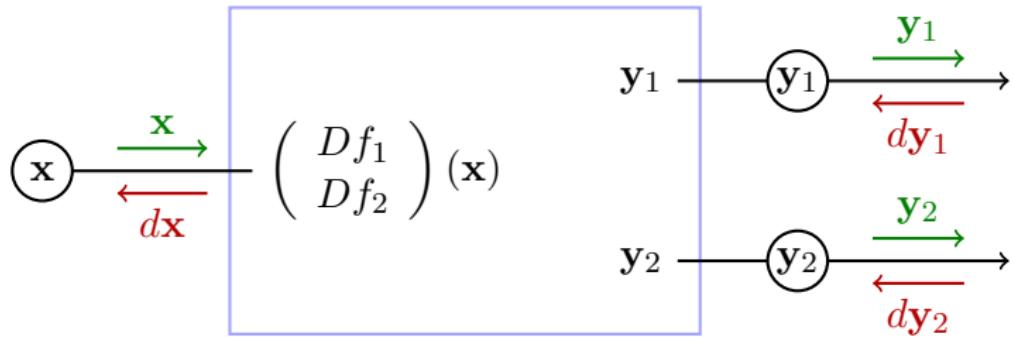
$$Df(\mathbf{x}) = (Df_1; Df_2)(\mathbf{x}) \quad \text{or} \quad \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \left(\frac{\partial \mathbf{y}_1}{\partial \mathbf{x}}; \frac{\partial \mathbf{y}_2}{\partial \mathbf{x}} \right)$$

- $d\mathbf{y}$ is also split as $d\mathbf{y} = (d\mathbf{y}_1, d\mathbf{y}_2)$ and $d\mathbf{x} = d\mathbf{y} \cdot Df(\mathbf{x})$ becomes

$$d\mathbf{x}^\top = \sum_i d\mathbf{y}_i^\top \cdot Df_i(\mathbf{x}) \quad \text{or} \quad \frac{\partial}{\partial \mathbf{x}} = \sum_i \frac{\partial}{\partial \mathbf{y}_i} \cdot \frac{\partial \mathbf{y}_i}{\partial \mathbf{x}}$$

splitting the output

$$f = (f_1, f_2)$$



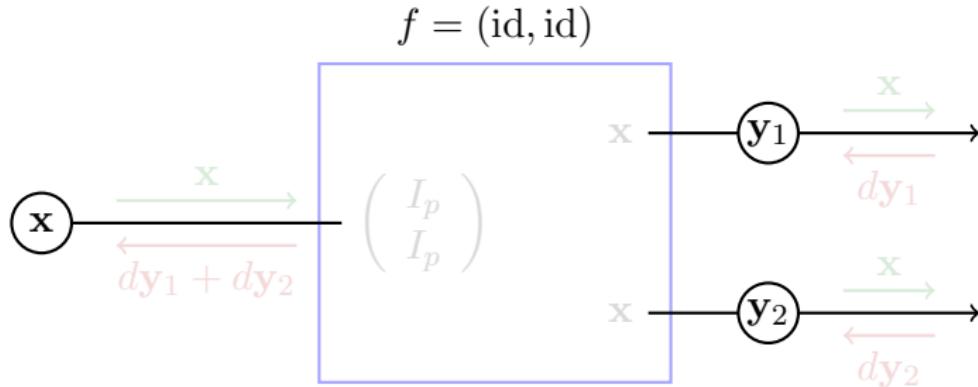
- we split output \mathbf{y} into subvectors as $\mathbf{y} = (\mathbf{y}_1, \mathbf{y}_2) = (f_1(\mathbf{x}), f_2(\mathbf{x}))$
- then, the derivative consists of blocks stacked vertically

$$Df(\mathbf{x}) = (Df_1; Df_2)(\mathbf{x}) \quad \text{or} \quad \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \left(\frac{\partial \mathbf{y}_1}{\partial \mathbf{x}}; \frac{\partial \mathbf{y}_2}{\partial \mathbf{x}} \right)$$

- $d\mathbf{y}$ is also split as $d\mathbf{y} = (d\mathbf{y}_1, d\mathbf{y}_2)$ and $d\mathbf{x} = d\mathbf{y} \cdot Df(\mathbf{x})$ becomes

$$d\mathbf{x}^\top = \sum_i d\mathbf{y}_i^\top \cdot Df_i(\mathbf{x}) \quad \text{or} \quad \frac{\partial}{\partial \mathbf{x}} = \sum_i \frac{\partial}{\partial \mathbf{y}_i} \cdot \frac{\partial \mathbf{y}_i}{\partial \mathbf{x}}$$

example: splitter (sharing)

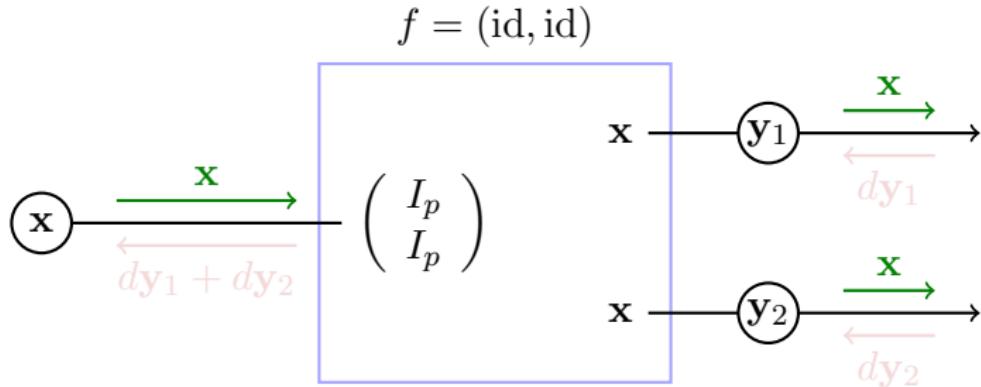


- if $f(\mathbf{x}) = (\mathbf{x}, \mathbf{x})$ and $\mathbf{x} \in \mathbb{R}^p$, then $Df(\mathbf{x}) = (I_p; I_p)$
- and the node behaves like **sum** backwards

$$d\mathbf{x} = d\mathbf{y}_1 + d\mathbf{y}_2 \quad \text{or} \quad \frac{\partial}{\partial \mathbf{x}} = \frac{\partial}{\partial \mathbf{y}_1} + \frac{\partial}{\partial \mathbf{y}_2}$$

- whenever a variable is shared (used more than once), we need to sum the gradients flowing from all paths where it appears

example: splitter (sharing)

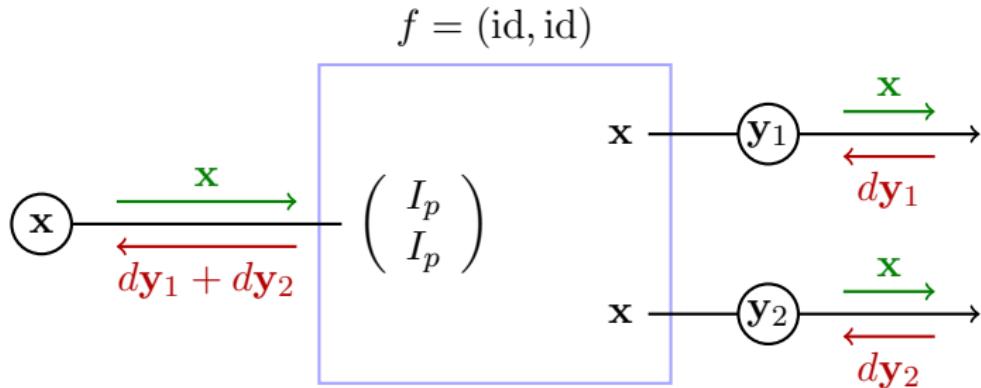


- if $f(\mathbf{x}) = (\mathbf{x}, \mathbf{x})$ and $\mathbf{x} \in \mathbb{R}^p$, then $Df(\mathbf{x}) = (I_p; I_p)$
- and the node behaves like **sum** backwards

$$d\mathbf{x} = dy_1 + dy_2 \quad \text{or} \quad \frac{\partial}{\partial \mathbf{x}} = \frac{\partial}{\partial \mathbf{y}_1} + \frac{\partial}{\partial \mathbf{y}_2}$$

- whenever a variable is shared (used more than once), we need to sum the gradients flowing from all paths where it appears

example: splitter (sharing)

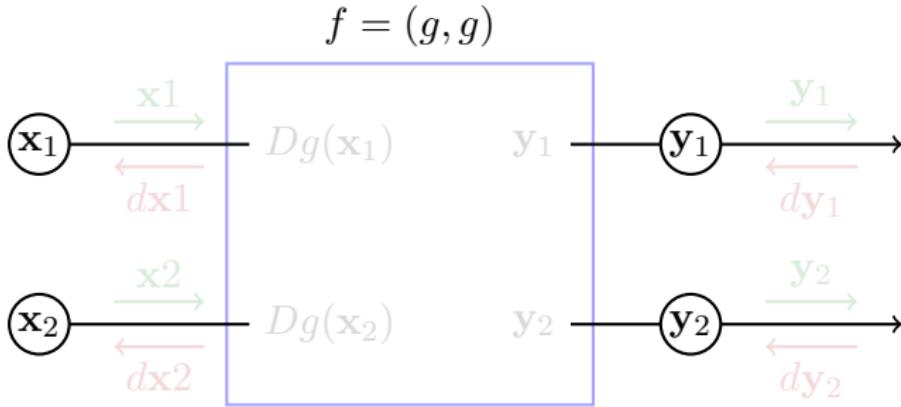


- if $f(\mathbf{x}) = (\mathbf{x}, \mathbf{x})$ and $\mathbf{x} \in \mathbb{R}^p$, then $Df(\mathbf{x}) = (I_p; I_p)$
- and the node behaves like **sum** backwards

$$d\mathbf{x} = d\mathbf{y}_1 + d\mathbf{y}_2 \quad \text{or} \quad \frac{\partial}{\partial \mathbf{x}} = \frac{\partial}{\partial \mathbf{y}_1} + \frac{\partial}{\partial \mathbf{y}_2}$$

- whenever a variable is shared (used more than once), we need to sum the gradients flowing from all paths where it appears

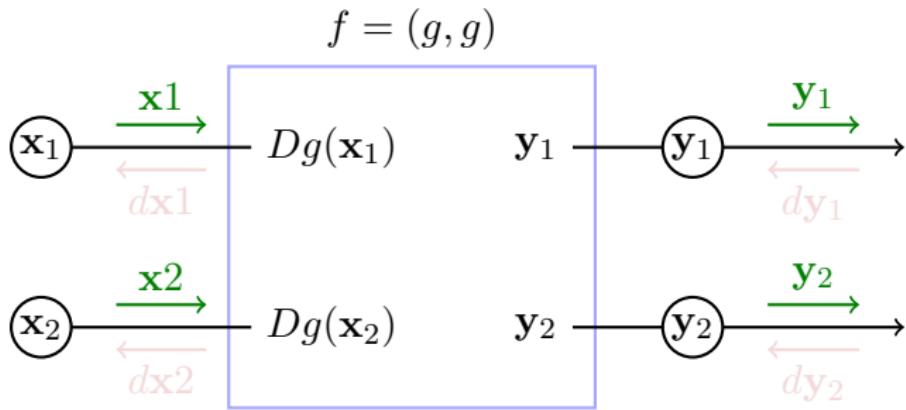
example: tuples



- if $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$, $\mathbf{y} = (\mathbf{y}_1, \mathbf{y}_2)$ and $f = (g, g)$, then $Df(\mathbf{x})$ is block-wise diagonal: $\text{diag}(Dg(\mathbf{x}_1), Dg(\mathbf{x}_2))$
- and the backward paths flow independently like the forward

$$d\mathbf{x}_i^\top = d\mathbf{y}_i^\top \cdot Dg(\mathbf{x}_i) \quad \text{or} \quad \frac{\partial}{\partial \mathbf{x}_i} = \frac{\partial}{\partial \mathbf{y}_i} \cdot \frac{\partial \mathbf{y}_i}{\partial \mathbf{x}_i}$$

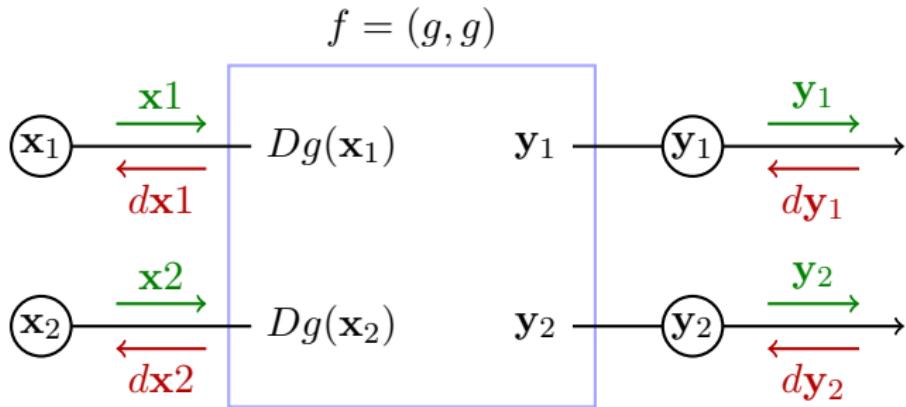
example: tuples



- if $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$, $\mathbf{y} = (\mathbf{y}_1, \mathbf{y}_2)$ and $f = (g, g)$, then $Df(\mathbf{x})$ is block-wise diagonal: $\text{diag}(Dg(\mathbf{x}_1), Dg(\mathbf{x}_2))$
- and the backward paths flow independently like the forward

$$d\mathbf{x}_i^\top = d\mathbf{y}_i^\top \cdot Dg(\mathbf{x}_i) \quad \text{or} \quad \frac{\partial}{\partial \mathbf{x}_i} = \frac{\partial}{\partial \mathbf{y}_i} \cdot \frac{\partial \mathbf{y}_i}{\partial \mathbf{x}_i}$$

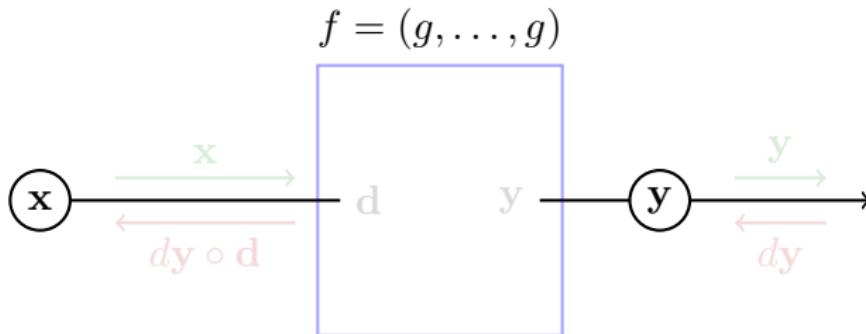
example: tuples



- if $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$, $\mathbf{y} = (\mathbf{y}_1, \mathbf{y}_2)$ and $f = (g, g)$, then $Df(\mathbf{x})$ is block-wise diagonal: $\text{diag}(Dg(\mathbf{x}_1), Dg(\mathbf{x}_2))$
- and the backward paths flow independently like the forward

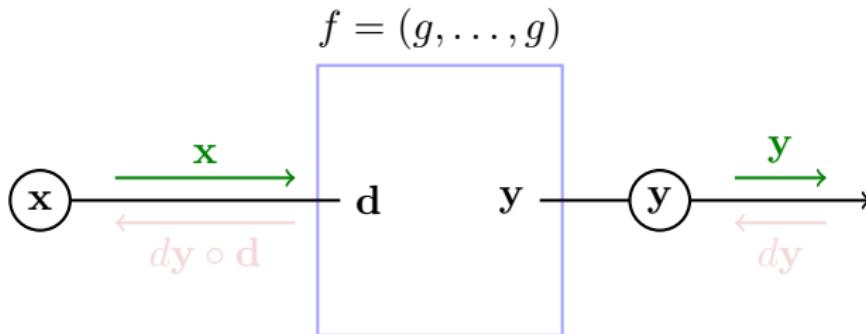
$$d\mathbf{x}_i^\top = d\mathbf{y}_i^\top \cdot Dg(\mathbf{x}_i) \quad \text{or} \quad \frac{\partial}{\partial \mathbf{x}_i} = \frac{\partial}{\partial \mathbf{y}_i} \cdot \frac{\partial \mathbf{y}_i}{\partial \mathbf{x}_i}$$

example: element-wise functions



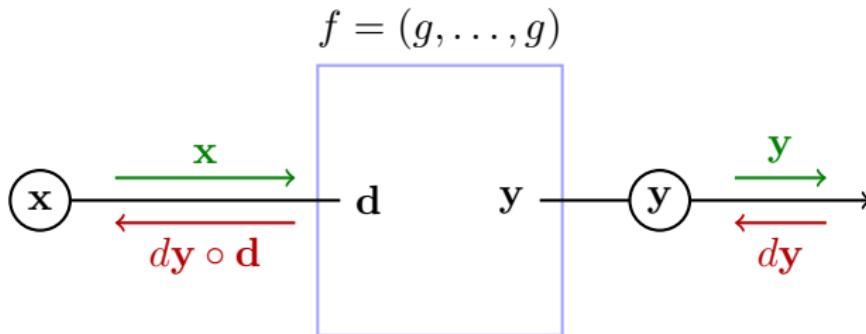
- if $\mathbf{x} \in \mathbb{R}^p$ and f is element-wise with $f(\mathbf{x}) = (g(x_1), \dots, g(x_p))$ where $g : \mathbb{R} \rightarrow \mathbb{R}$, then $Df(\mathbf{x}) = \text{diag } \mathbf{d}$ is diagonal, where $\mathbf{d} = (Dg(x_1), \dots, Dg(x_p))$
- and the partial derivatives are element-wise multiplied

example: element-wise functions



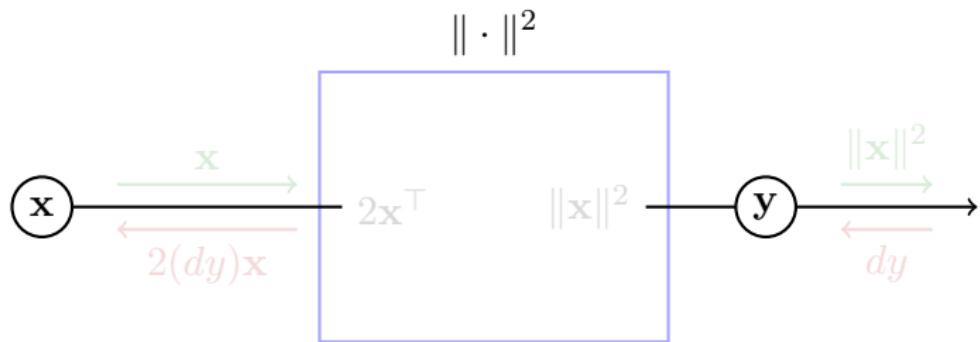
- if $\mathbf{x} \in \mathbb{R}^p$ and f is element-wise with $f(\mathbf{x}) = (g(x_1), \dots, g(x_p))$ where $g : \mathbb{R} \rightarrow \mathbb{R}$, then $Df(\mathbf{x}) = \text{diag } \mathbf{d}$ is diagonal, where $\mathbf{d} = (Dg(x_1), \dots, Dg(x_p))$
- and the partial derivatives are element-wise multiplied

example: element-wise functions



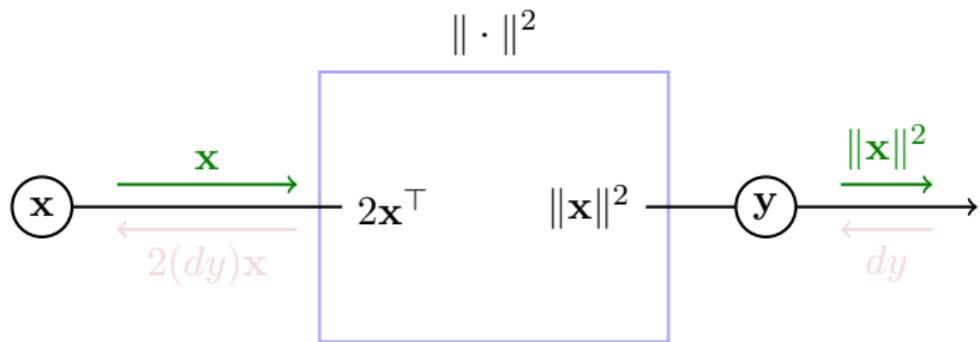
- if $\mathbf{x} \in \mathbb{R}^p$ and f is element-wise with $f(\mathbf{x}) = (g(x_1), \dots, g(x_p))$ where $g : \mathbb{R} \rightarrow \mathbb{R}$, then $Df(\mathbf{x}) = \text{diag } \mathbf{d}$ is diagonal, where $\mathbf{d} = (Dg(x_1), \dots, Dg(x_p))$
- and the partial derivatives are element-wise multiplied

example: squared norm



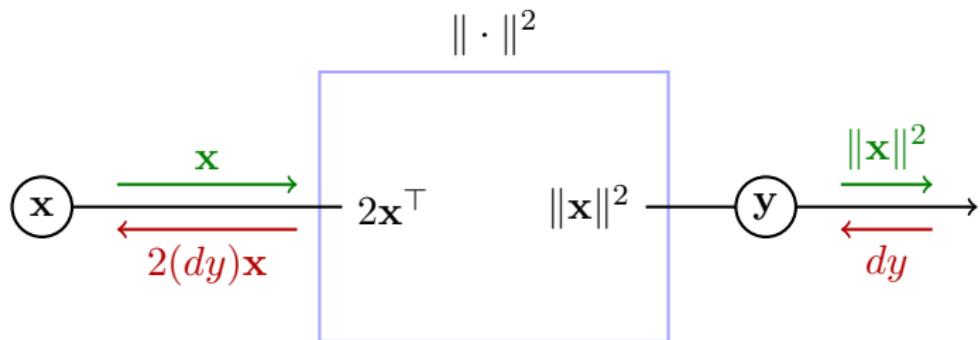
- if $f(x) = \|x\|^2$ then $Df(x) = 2x^\top$
- and dy is multiplied by $2x^\top$; this can be seen by composing a splitter (factor 2) with a dot product (factor x^\top)

example: squared norm



- if $f(\mathbf{x}) = \|\mathbf{x}\|^2$ then $Df(\mathbf{x}) = 2\mathbf{x}^\top$
- and dy is multiplied by $2\mathbf{x}^\top$; this can be seen by composing a splitter (factor 2) with a dot product (factor \mathbf{x}^\top)

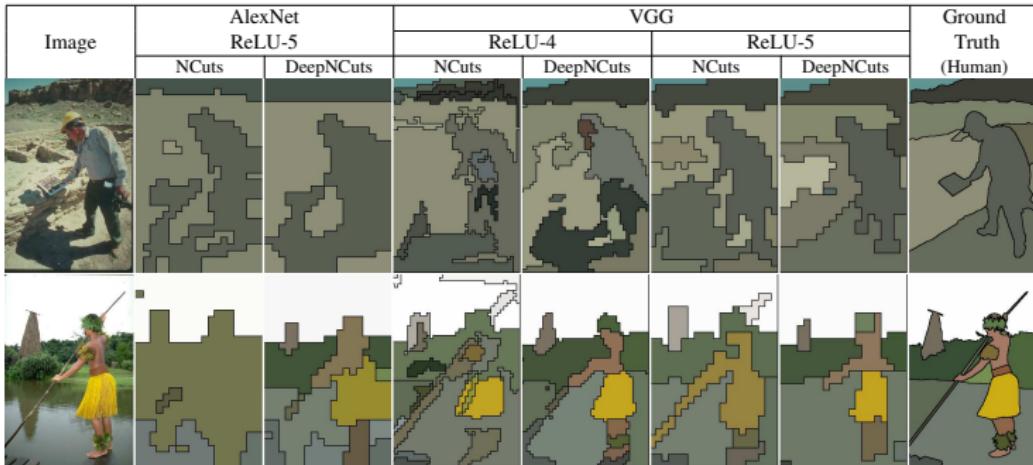
example: squared norm



- if $f(\mathbf{x}) = \|\mathbf{x}\|^2$ then $Df(\mathbf{x}) = 2\mathbf{x}^\top$
- and dy is multiplied by $2\mathbf{x}^\top$; this can be seen by composing a splitter (factor 2) with a dot product (factor \mathbf{x}^\top)

matrix derivatives

[Ionescu et al. 2015]



- derivatives for
 - SVD decomposition $A = U\Sigma V^\top$
 - eigenvalue decomposition $A = U\Sigma U^\top$
 - nonlinear matrix functions $f(A) = Uf(\Sigma)U^\top$
- application to spectral methods for image segmentation

matrix calculus

- results like these, and many more

$$\frac{\partial A\mathbf{x}}{\partial \mathbf{x}} = A$$

$$\frac{\partial \mathbf{x}^\top A\mathbf{x}}{\partial \mathbf{x}} = \mathbf{x}^\top (A + A^\top)$$

$$\frac{\partial \text{vec}(\mathbf{x}^\top A\mathbf{x})}{\partial \text{vec } A} = \mathbf{x}^\top \otimes \mathbf{x}^\top$$

$$\frac{\partial AXB}{\partial X} = B^\top \otimes A$$

$$\frac{dA^{-1}}{dA} = -(A^{-\top} \otimes A^{-1})$$

$$\frac{d \ln |A|}{dA} = \text{vec}(A^{-\top})^\top$$

$$\frac{\partial \text{tr}(AX)}{\partial X} = \text{vec}(A^\top)^\top$$

in general

- apparently, we do not need to store the Jacobian matrix $Df(\mathbf{x})$, which may be huge, but only what is needed to compute the partial derivatives in the backward pass
- our function can be decomposed into a **directed acyclic graph** (DAG) of nodes, called a **computational graph**
- each time we call the function in the forward pass, a new graph may be constructed if our program contains control flow statements like conditionals and loops; methods supporting this operation are called **dynamic**

automatic differentiation

[Wengert 1964]

- is the more general set of methods used to automatically evaluate the derivative of a given function at a given input; it is **not numerical** and **not symbolic**
- what we call back-propagation here is known as the **reverse accumulation** mode in this context and makes sense because we compute the gradient of a single scalar quantity with respect to maybe millions of parameters
- **forward accumulation** makes sense when we need the derivative of many variables with respect to few parameters
- we will use the term **automatic differentiation** to refer to the process of generating a computer program for the derivatives given the program for the original function

aside: higher-order derivatives

- the Hessian was assumed fixed and isotropic in gradient descent; if we knew it, we could use the **Newton method** instead and solve all curvature-related problems
- given $f : \mathbb{R}^p \rightarrow \mathbb{R}$, its **Hessian matrix** at \mathbf{x} is

$$Hf(\mathbf{x}) := \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_p} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_p \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_p^2} \end{pmatrix}(\mathbf{x}) = \nabla(Df)(\mathbf{x})$$

- unfortunately, this is a $p \times p$ matrix and with p in the order of millions, it is impractical even to store it, let alone compute it

aside: multiplication by Hessian

[Pearlmutter 1994]

- fortunately, in many cases what we need is only the **product** of the Hessian with a given vector \mathbf{v} , which is just a vector in \mathbb{R}^p

$$\mathbf{v}^\top \cdot Hf(\mathbf{x}) = \mathbf{v}^\top \cdot \nabla(Df)(\mathbf{x}) = \nabla_{\mathbf{v}}(Df)(\mathbf{x})$$

- here $\nabla_{\mathbf{v}}$ is the **directional derivative** operator

$$\nabla_{\mathbf{v}}(f) := \mathbf{v}^\top \cdot \nabla f$$

- remember that in back-propagation, for each variable \mathbf{x} , we defined a vector $d\mathbf{x}$, which was computed in the backward pass
- so all we need to do is allocate another vector $\nabla_{\mathbf{v}}(\mathbf{x})$ for the forward pass and another $\nabla_{\mathbf{v}}(d\mathbf{x})$ for the backward, and compute them by applying the chain rule in both passes!

aside: multiplication by Hessian

[Pearlmutter 1994]

- fortunately, in many cases what we need is only the **product** of the Hessian with a given vector \mathbf{v} , which is just a vector in \mathbb{R}^p

$$\mathbf{v}^\top \cdot Hf(\mathbf{x}) = \mathbf{v}^\top \cdot \nabla(Df)(\mathbf{x}) = \nabla_{\mathbf{v}}(Df)(\mathbf{x})$$

- here $\nabla_{\mathbf{v}}$ is the **directional derivative** operator

$$\nabla_{\mathbf{v}}(f) := \mathbf{v}^\top \cdot \nabla f$$

- remember that in back-propagation, for each variable \mathbf{x} , we defined a vector $d\mathbf{x}$, which was computed in the backward pass
- so all we need to do is allocate another vector $\nabla_{\mathbf{v}}(\mathbf{x})$ for the forward pass and another $\nabla_{\mathbf{v}}(d\mathbf{x})$ for the backward, and compute them by applying the chain rule in both passes!

automatic differentiation: units

automatic differentiation

forward

- evaluation is carried out by **units**, one calling another
- when invoked, each unit generates a **node** object
- each node holds the **gradient** with respect to its unit's inputs, including parameters
- it also holds any information needed for the backward pass

backward

- all gradients are set to **zero**, except for the gradient with respect to the scalar quantity that is to be optimized (the error), which is set to **one**
- the **back()** method is invoked on the node of this quantity
- this, in turn, triggers the same method on all units that have participated in the forward pass

automatic differentiation

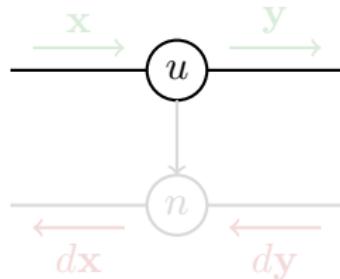
forward

- evaluation is carried out by **units**, one calling another
- when invoked, each unit generates a **node** object
- each node holds the **gradient** with respect to its unit's inputs, including parameters
- it also holds any information needed for the backward pass

backward

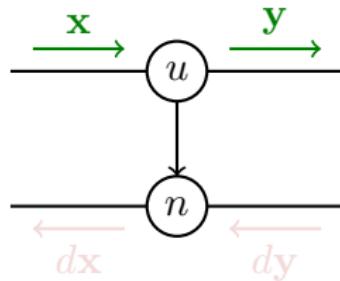
- all gradients are set to **zero**, except for the gradient with respect to the scalar quantity that is to be optimized (the error), which is set to **one**
- the **back()** method is invoked on the node of this quantity
- this, in turn, triggers the same method on all units that have participated in the forward pass

units and nodes



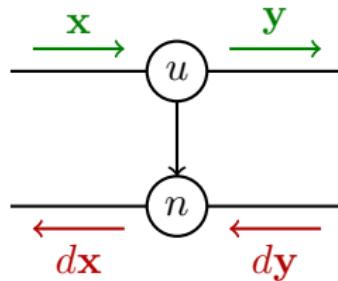
- unit u manually generates node n

units and nodes



- unit u manually generates node n

units and nodes



- unit u manually generates node n

units and nodes

- given a function f with derivative Df , a *unit* is a function of the form

```
def forward(x1, ..., xn):  
    y = f(x1, ..., xn)  
    def back(dy, dx1, ..., dxn):  
        dx1T += dyT · D1f(x1, ..., xn)  
        ...  
        dxnT += dyT · Dnf(x1, ..., xn)  
    return node(y, back)
```

- a *node* object:
 - holds y and an associated derivative dy of the same shape
 - exposes a method $back(x_1, \dots, x_n)$ where x_i can be *nodes*
 - automatically adds its own dy as first argument
 - if an input x_i is a *node*, extracts the derivative part dx_i
 - otherwise, dx_i is an object for which operation $+=$ is ignored

units and nodes

- given a function f with derivative Df , a *unit* is a function of the form

```
def forward(x1, ..., xn):  
    y = f(x1, ..., xn)  
    def back(dy, dx1, ..., dxn):  
        dx1T += dyT · D1f(x1, ..., xn)  
        ...  
        dxnT += dyT · Dnf(x1, ..., xn)  
    return node(y, back)
```

- a *node* object:

- holds y and an associated derivative dy of the same shape
- exposes a method $back(x_1, \dots, x_n)$ where x_i can be *nodes*
- automatically adds its own dy as first argument
- if an input x_i is a *node*, extracts the derivative part dx_i
- otherwise, dx_i is an object for which operation $+=$ is ignored

the affine unit

- input vectors are represented as rows of $m \times p$ **input matrix** X where m is the **mini-batch size** and p the input dimension
- parameters are represented by $p \times q$ **weight matrix** W and $1 \times q$ **bias vector** \mathbf{b} where q is the output dimension
- the unit is defined as

```
def affine(X, (W, b)):  
    A = dot(X, W) + b  
def back(dA, dX, (dW, db)):  
    dW += dot(XT, dA)  
    db += sum0(dA)  
    dX += dot(dA, WT)  
return node(A, back)
```

the affine unit

- input vectors are represented as rows of $m \times p$ **input matrix** X where m is the **mini-batch size** and p the input dimension
- parameters are represented by $p \times q$ **weight matrix** W and $1 \times q$ **bias vector** \mathbf{b} where q is the output dimension
- the unit is defined as

```
def affine(X, (W, b)):  
    A = dot(X, W) + b  
    def back(dA, dX, (dW, db)):  
        dW += dot(XT, dA)  
        db += sum0(dA)  
        dX += dot(dA, WT)  
    return node(A, back)
```

the affine unit in math

forward

- input $X \in \mathbb{R}^{m \times p}$, $W \in \mathbb{R}^{p \times q}$, $\mathbf{b} \in \mathbb{R}^q$, output $A \in \mathbb{R}^{m \times q}$

$$A = f(X; W, \mathbf{b}) := XW + \mathbf{1}_m \mathbf{b}^\top$$

observe that in the code, addition of \mathbf{b} is handled by **broadcasting**

backward

- if \mathbf{a}_i , \mathbf{w}_i is the i -th column of A , W ,

$$\frac{\partial \mathbf{a}_i}{\partial \mathbf{w}_i} = \frac{\partial (X\mathbf{w}_i)}{\partial \mathbf{w}_i} = X$$

and there are no other dependencies, so by the chain rule

$$d\mathbf{w}_i^\top := \frac{\partial}{\partial \mathbf{w}_i} = \frac{\partial}{\partial \mathbf{a}_i} \cdot \frac{\partial \mathbf{a}_i}{\partial \mathbf{w}_i} = d\mathbf{a}_i^\top \cdot X$$

- finally, the partial derivative with respect to W

$$dW = (dA^\top X)^\top = X^\top dA$$

the affine unit in math

forward

- input $X \in \mathbb{R}^{m \times p}$, $W \in \mathbb{R}^{p \times q}$, $\mathbf{b} \in \mathbb{R}^q$, output $A \in \mathbb{R}^{m \times q}$

$$A = f(X; W, \mathbf{b}) := XW + \mathbf{1}_m \mathbf{b}^\top$$

observe that in the code, addition of \mathbf{b} is handled by **broadcasting**

backward

- if \mathbf{a}_i , \mathbf{w}_i is the i -th column of A , W ,

$$\frac{\partial \mathbf{a}_i}{\partial \mathbf{w}_i} = \frac{\partial (X \mathbf{w}_i)}{\partial \mathbf{w}_i} = X$$

and there are no other dependencies, so by the chain rule

$$d\mathbf{w}_i^\top := \frac{\partial}{\partial \mathbf{w}_i} = \frac{\partial}{\partial \mathbf{a}_i} \cdot \frac{\partial \mathbf{a}_i}{\partial \mathbf{w}_i} = d\mathbf{a}_i^\top \cdot X$$

- finally, the partial derivative with respect to W

$$dW = (dA^\top X)^\top = X^\top dA$$

the affine unit in math

- by symmetry, writing $A^\top = W^\top X^\top + \mathbf{b}\mathbf{1}_m^\top$ and using the previous result for dW , we find $dX^\top = (W^\top)^\top dA^\top$ or

$$dX = (dA)W^\top$$

- again, by replacing X and W by $\mathbf{1}_m$ and \mathbf{b}^\top respectively in the previous result for dW ,

$$d\mathbf{b}^\top = (dA^\top \mathbf{1}_m)^\top = \mathbf{1}_m^\top dA$$

- observe that distributing \mathbf{b} in the forward yields a sum in the backward

the affine unit in math

- by symmetry, writing $A^\top = W^\top X^\top + \mathbf{b}\mathbf{1}_m^\top$ and using the previous result for dW , we find $dX^\top = (W^\top)^\top dA^\top$ or

$$dX = (dA)W^\top$$

- again, by replacing X and W by $\mathbf{1}_m$ and \mathbf{b}^\top respectively in the previous result for dW ,

$$d\mathbf{b}^\top = (dA^\top \mathbf{1}_m)^\top = \mathbf{1}_m^\top dA$$

- observe that distributing \mathbf{b} in the forward yields a sum in the backward

the logistic unit

- the input is an $m \times q$ activation matrix A and the $m \times k$ one-of- k encoded target matrix, where k is the number of classes
- there are no parameters
- the unit integrates softmax with average cross-entropy loss

```
def logistic(A, T):
    E = exp(A)
    Y = E/sum_1(E)
    C = -sum_1(T * log(Y))
    D = sum_0(C)/m
    def back(dD, dA, _):
        dA += dD * (Y - T)/m
    return node(D, back)
```

the logistic unit

- the input is an $m \times q$ activation matrix A and the $m \times k$ one-of- k encoded target matrix, where k is the number of classes
- there are no parameters
- the unit integrates softmax with average cross-entropy loss

```
def logistic(A, T):
    E = exp(A)
    Y = E/sum_1(E)
    C = -sum_1(T * log(Y))
    D = sum_0(C)/m
    def back(dD, dA, _):
        dA += dD * (Y - T)/m
    return node(D, back)
```

the logistic unit in math

forward

- E is given element-wise as $e_{ij} = \exp(a_{ij})$, and $m \times q$ matrix Y is row-normalized as

$$Y = (\text{diag}(E\mathbf{1}_k))^{-1}E$$

- the i -th row of Y is the softmax output of the i -th input sample representing the k posterior class probabilities
- C is actually a $m \times 1$ column vector and its i -th element represents the cross-entropy loss of the i -th input sample

$$c_i = - \sum_{j=1}^k t_{ij} \log(y_{ij})$$

- finally, $D = \frac{1}{m} \sum_{i=1}^m c_i$ is a scalar and represents the average cross-entropy (data) error over the mini-batch

the logistic unit in math

forward

- E is given element-wise as $e_{ij} = \exp(a_{ij})$, and $m \times q$ matrix Y is row-normalized as

$$Y = (\text{diag}(E\mathbf{1}_k))^{-1}E$$

- the i -th row of Y is the softmax output of the i -th input sample representing the k posterior class probabilities
- C is actually a $m \times 1$ column vector and its i -th element represents the cross-entropy loss of the i -th input sample

$$c_i = - \sum_{j=1}^k t_{ij} \log(y_{ij})$$

- finally, $D = \frac{1}{m} \sum_{i=1}^m c_i$ is a scalar and represents the average cross-entropy (data) error over the mini-batch

the logistic unit in math

backward

- if $\mathbf{a}_i^\top, \mathbf{y}_i^\top, \mathbf{t}_i^\top$ is the i -th row of A, Y, T , the derivative of the cross-entropy loss is, according to what we have seen,

$$\frac{\partial c_i}{\partial \mathbf{a}_i}(\mathbf{a}_i, \mathbf{t}_i) = (\sigma(\mathbf{a}_i) - \mathbf{t}_i)^\top = (\mathbf{y}_i - \mathbf{t}_i)^\top$$

- since D is the **average** of the individual sample losses c_i , the derivative of the total error, which is 1 by default, is **distributed** over the samples with a factor of $\frac{1}{m}$

$$dA^\top = \frac{1}{m}(Y - T) \cdot dD$$

why integrate softmax with cross-entropy?

- the simplified formula is faster compared to blind application of back-propagation at the level of elementary functions
- if this is not convincing, try evaluating the binary cross-entropy loss

$$\ell(x) := \ln(1 + e^{-x})$$

- $\ell(-1) = 1.3133$
- $\ell(-2) = 2.1269$
- $\ell(-5) = 5.0067$
- $\ell(-10) = 10.0000$
- $\ell(-20) = 20.0000$
- $\ell(-50) = 50.0000$
- $\ell(-100) = 100.0000$
- $\ell(-200) = 200.0000$
- $\ell(-500) = 500.0000$
- $\ell(-1000) = \infty$
- $\ell(-2000) = \infty$

why integrate softmax with cross-entropy?

- the simplified formula is faster compared to blind application of back-propagation at the level of elementary functions
- if this is not convincing, try evaluating the binary cross-entropy loss

$$\ell(x) := \ln(1 + e^{-x})$$

- $\ell(-1) = 1.3133$
- $\ell(-2) = 2.1269$
- $\ell(-5) = 5.0067$
- $\ell(-10) = 10.0000$
- $\ell(-20) = 20.0000$
- $\ell(-50) = 50.0000$
- $\ell(-100) = 100.0000$
- $\ell(-200) = 200.0000$
- $\ell(-500) = 500.0000$
- $\ell(-1000) = \infty$
- $\ell(-2000) = \infty$

why integrate softmax with cross-entropy?

- the simplified formula is faster compared to blind application of back-propagation at the level of elementary functions
- if this is not convincing, try evaluating the binary cross-entropy loss

$$\ell(x) := \ln(1 + e^{-x})$$

- $\ell(-1) = 1.3133$
- $\ell(-2) = 2.1269$
- $\ell(-5) = 5.0067$
- $\ell(-10) = 10.0000$
- $\ell(-20) = 20.0000$
- $\ell(-50) = 50.0000$
- $\ell(-100) = 100.0000$
- $\ell(-200) = 200.0000$
- $\ell(-500) = 500.0000$
- $\ell(-1000) = \infty$
- $\ell(-2000) = \infty$

why integrate softmax with cross-entropy?

- the simplified formula is faster compared to blind application of back-propagation at the level of elementary functions
- if this is not convincing, try evaluating the binary cross-entropy loss

$$\ell(x) := \ln(1 + e^{-x})$$

- $\ell(-1) = 1.3133$
- $\ell(-2) = 2.1269$
- $\ell(-5) = 5.0067$
- $\ell(-10) = 10.0000$
- $\ell(-20) = 20.0000$
- $\ell(-50) = 50.0000$
- $\ell(-100) = 100.0000$
- $\ell(-200) = 200.0000$
- $\ell(-500) = 500.0000$
- $\ell(-1000) = \infty$
- $\ell(-2000) = \infty$

why integrate softmax with cross-entropy?

- the simplified formula is faster compared to blind application of back-propagation at the level of elementary functions
- if this is not convincing, try evaluating the binary cross-entropy loss

$$\ell(x) := \ln(1 + e^{-x})$$

- $\ell(-1) = 1.3133$
- $\ell(-2) = 2.1269$
- $\ell(-5) = 5.0067$
- $\ell(-10) = 10.0000$
- $\ell(-20) = 20.0000$
- $\ell(-50) = 50.0000$
- $\ell(-100) = 100.0000$
- $\ell(-200) = 200.0000$
- $\ell(-500) = 500.0000$
- $\ell(-1000) = \infty$
- $\ell(-2000) = \infty$

why integrate softmax with cross-entropy?

- the simplified formula is faster compared to blind application of back-propagation at the level of elementary functions
- if this is not convincing, try evaluating the binary cross-entropy loss

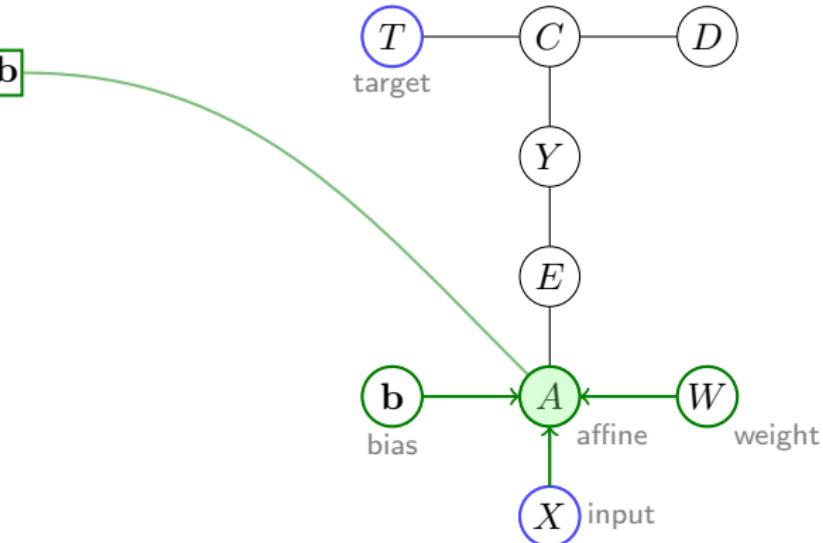
$$\ell(x) := \ln(1 + e^{-x})$$

- $\ell(-1) = 1.3133$
- $\ell(-2) = 2.1269$
- $\ell(-5) = 5.0067$
- $\ell(-10) = 10.0000$
- $\ell(-20) = 20.0000$
- $\ell(-50) = 50.0000$
- $\ell(-100) = 100.0000$
- $\ell(-200) = 200.0000$
- $\ell(-500) = 500.0000$
- $\ell(-1000) = \infty$
- $\ell(-2000) = \infty$

back-propagation

forward

$$A = \text{dot}(Z, W) + b$$



back-propagation

forward

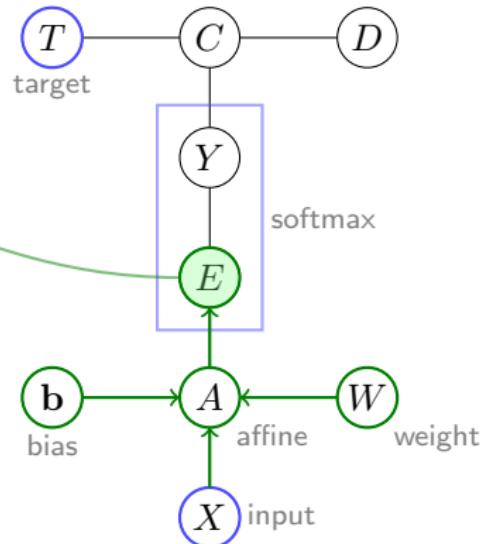
$$A = \text{dot}(Z, W) + b$$

$$E = \exp(A)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$



back-propagation

forward

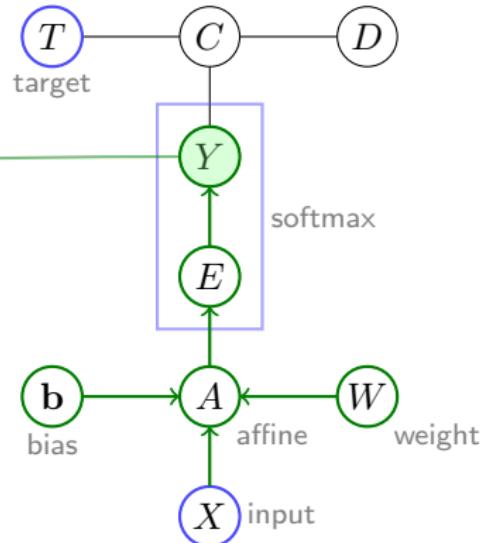
$$A = \text{dot}(Z, W) + b$$

$$E = \exp(A)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$



back-propagation

forward

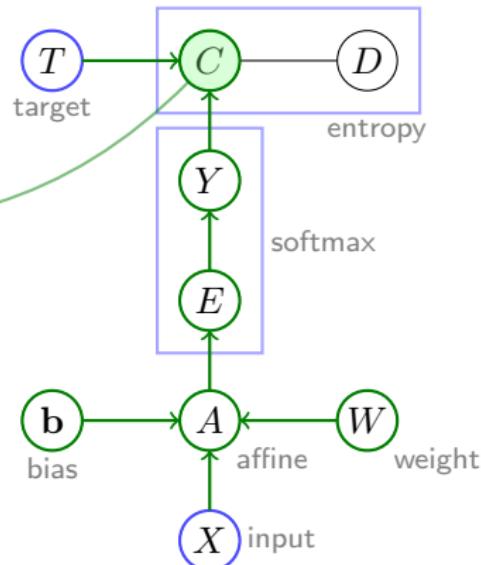
$$A = \text{dot}(Z, W) + b$$

$$E = \exp(A)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$



back-propagation

forward

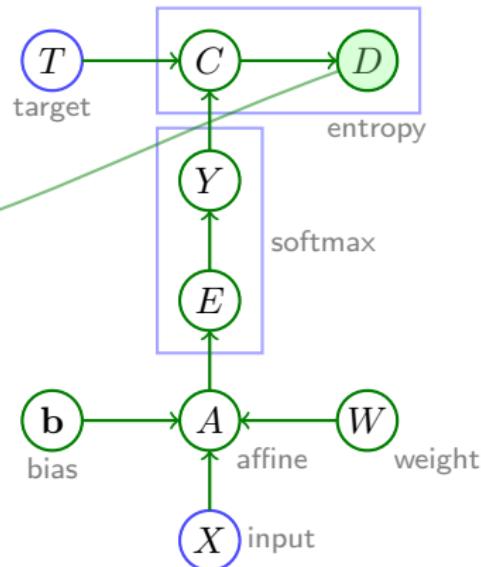
$$A = \text{dot}(Z, W) + b$$

$$E = \exp(A)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$



back-propagation

forward

$$A = \text{dot}(Z, W) + \mathbf{b}$$

$$E = \exp(A)$$

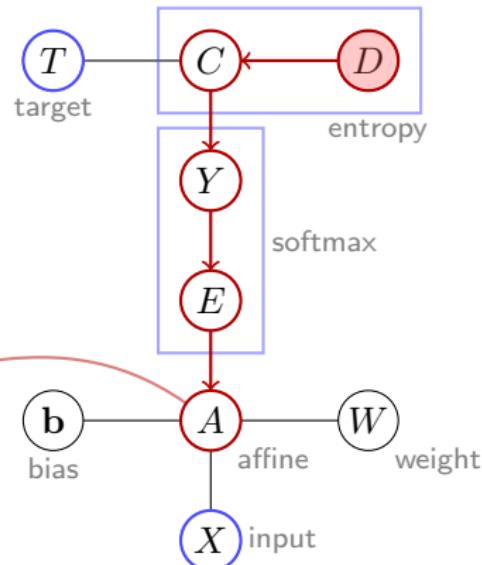
$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$

backward

$$dA = dD * (Y - T)/m$$



back-propagation

forward

$$A = \text{dot}(Z, W) + \mathbf{b}$$

$$E = \exp(A)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

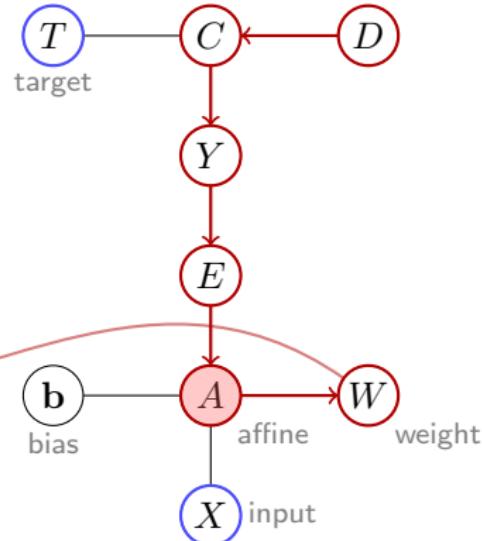
$$D = \text{sum}_0(C)/m$$

backward

$$dA = dD * (Y - T)/m$$

$$dW += \text{dot}(Z^\top, dA)$$

$$d\mathbf{b} = \text{sum}_0(dA)$$



back-propagation

forward

$$A = \text{dot}(Z, W) + b$$

$$E = \exp(A)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

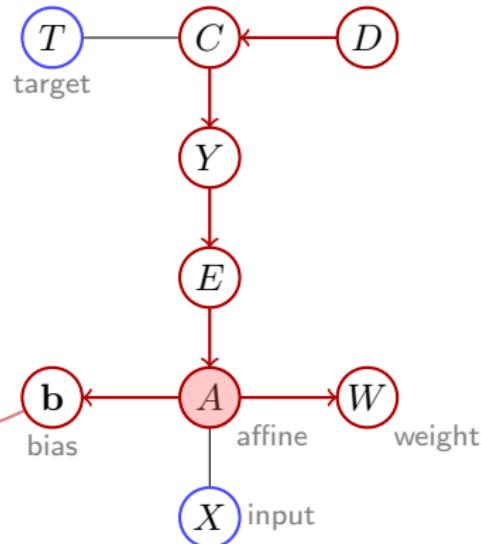
$$D = \text{sum}_0(C)/m$$

backward

$$dA = dD * (Y - T)/m$$

$$dW += \text{dot}(Z^\top, dA)$$

$$db = \text{sum}_0(dA)$$



automatic differentiation

forward

$$A = \text{dot}(Z, W) + \mathbf{b}$$

$$E = \exp(A)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$

backward

$$dA = dD * (Y - T)/m$$

$$dW += \text{dot}(Z^\top, dA)$$

$$d\mathbf{b} = \text{sum}_0(dA)$$

now we organize **forward** and **backward** code into units

automatic differentiation

forward

```
A = dot(Z, W) + b  
E = exp(A)  
Y = E/sum1(E)  
C = -sum1(T * log(Y))  
D = sum0(C)/m
```

backward

```
dA = dD * (Y - T)/m  
dW += dot(Z⊤, dA)  
db = sum0(dA)
```

```
def affine(X, (W, b)):  
    A = dot(X, W) + b  
    def back(dA, dX, (dW, db)):  
        dW += dot(X⊤, dA)  
        db += sum0(dA)  
        dX += dot(dA, W⊤)  
    return node(A, back)
```

automatic differentiation

forward

```
A = affine(Z, (W, b))  
E = exp(A)  
Y = E/sum1(E)  
C = -sum1(T * log(Y))  
D = sum0(C)/m
```

backward

```
dA = dD * (Y - T)/m  
A.back(Z, (W, b))
```

```
def affine(X, (W, b)):  
    A = dot(X, W) + b  
  
def back(dA, dX, (dW, db)):  
    dW += dot(X⊤, dA)  
    db += sum0(dA)  
    dX += dot(dA, W⊤)  
  
return node(A, back)
```

automatic differentiation

forward

```
A = affine(Z, (W, b))
```

```
E = exp(A)
```

```
Y = E/sum1(E)
```

```
C = -sum1(T * log(Y))
```

```
D = sum0(C)/m
```

backward

```
dA = dD * (Y - T)/m
```

```
A.back(Z, (W, b))
```

```
def logistic(A, T):
```

```
E = exp(A)
```

```
Y = E/sum1(E)
```

```
C = -sum1(T * log(Y))
```

```
D = sum0(C)/m
```

```
def back(dD, dA, _):
```

```
dA += dD * (Y - T)/m
```

```
return node(D, back)
```

automatic differentiation

forward

```
A = affine(Z, (W, b))
```

```
D = entropy(A, T)
```

```
def logistic(A, T):  
    E = exp(A)  
    Y = E/sum_1(E)  
    C = -sum_1(T * log(Y))  
    D = sum_0(C)/m
```

backward

```
D.back(A, T)  
A.back(Z, (W, b))
```

```
def back(dD, dA, _):  
    dA += dD * (Y - T)/m  
return node(D, back)
```

automatic differentiation: functions

the relu unit

- relu is an element-wise activation function; its input is **activation matrix** A and returns matrix Z of the same size
- its backward pass behaves like a **switch**

```
def relu(A):  
    Z = max(0, A)  
    def back(dZ, dA):  
        dA += dZ * (Z > 0)  
    return node(Z, back)
```

the relu unit

- relu is an element-wise activation function; its input is **activation matrix** A and returns matrix Z of the same size
- its backward pass behaves like a **switch**

```
def relu(A):
    Z = max(0, A)
    def back(dZ, dA):
        dA += dZ * (Z > 0)
    return node(Z, back)
```

the decay unit

- it takes as input a tuple or list W of weight matrices of any size and returns the weight decay error term $\frac{\lambda}{2}\|w\|^2$ for each $w \in W$, where $\|\cdot\|_F$ is the Frobenius norm
- the backward derivative is proportional to w , as for the ℓ_2 norm

```
def decay(W):
    R =  $\frac{\lambda}{2} * \text{sum}(\|w\|_F^2 \text{ for } w \text{ in } W)$ 
    def back(dR, dW):
        for (w, dw) in zip(W, dW):
            dw += dR *  $\lambda * w$ 
    return node(R, back)
```

the decay unit

- it takes as input a tuple or list W of weight matrices of any size and returns the weight decay error term $\frac{\lambda}{2} \|w\|^2$ for each $w \in W$, where $\|\cdot\|_F$ is the Frobenius norm
- the backward derivative is proportional to w , as for the ℓ_2 norm

```
def decay(W):
    R =  $\frac{\lambda}{2} * \text{sum}(\|w\|_F^2 \text{ for } w \text{ in } W)$ 
    def back(dR, dW):
        for (w, dw) in zip(W, dW):
            dw += dR *  $\lambda * w$ 
    return node(R, back)
```

the add unit

- it takes as input a tuple or list X of matrices (or vectors, or scalars) of the same size and returns their sum
- its backward pass **distributes** the derivative to all input branches

```
def add(X):
    S = sum(X)
    def back(dS, dX):
        for dx in dX:
            dx += dS
        return node(S, back)
```

- operator $+$ is overloaded for *nodes* such that $A + B$ means $\text{add}((A, B))$

the add unit

- it takes as input a tuple or list X of matrices (or vectors, or scalars) of the same size and returns their sum
- its backward pass **distributes** the derivative to all input branches

```
def add(X):
    S = sum(X)
    def back(dS, dX):
        for dx in dX:
            dx += dS
        return node(S, back)
```

- operator `+` is overloaded for *nodes* such that $A + B$ means `add((A, B))`

the add unit

- it takes as input a tuple or list X of matrices (or vectors, or scalars) of the same size and returns their sum
- its backward pass **distributes** the derivative to all input branches

```
def add(X):
    S = sum(X)
    def back(dS, dX):
        for dx in dX:
            dx += dS
        return node(S, back)
```

- operator `+` is overloaded for *nodes* such that $A + B$ means `add((A, B))`

the loss function

- it takes as input the activation matrix A , the target matrix T and the weight matrix list W
- it calls the logistic unit on (A, T) and the decay unit on W , and returns the sum of the two scalar terms

```
def loss(A, T, W):
    L = logistic(A, T) + decay(W)
    return block(L)
```

- addition is handled by add and the error derivative flows backward to both branches

the loss function

- it takes as input the activation matrix A , the target matrix T and the weight matrix list W
- it calls the logistic unit on (A, T) and the decay unit on W , and returns the sum of the two scalar terms

```
def loss(A, T, W):
    L = logistic(A, T) + decay(W)
    return block(L)
```

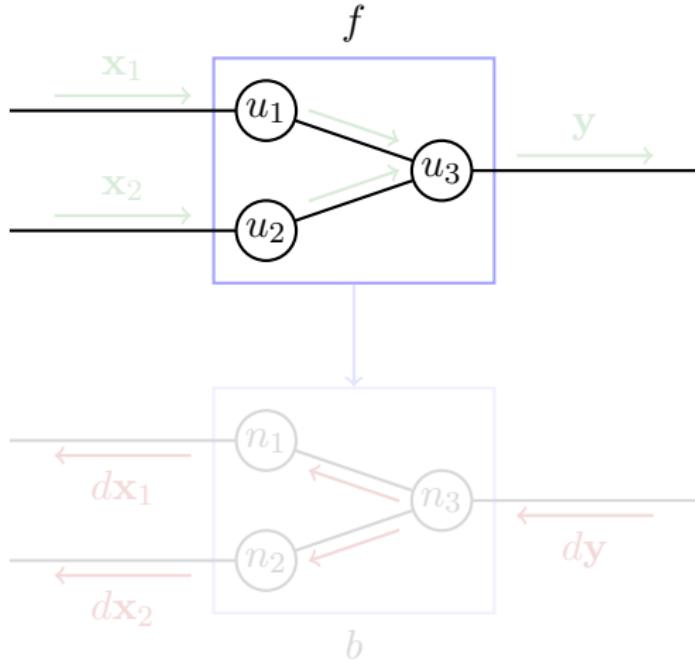
- addition is handled by add and the error derivative flows backward to both branches

the model function

- this is a two-layer network model where an affine layer is followed by a relu activation function and another affine layer
- the parameter tuple $U_i = (W_i, \mathbf{b}_i)$ for layer i contains a weight matrix W_i and a bias vector \mathbf{b}_i
- unit calls are nested like every other function

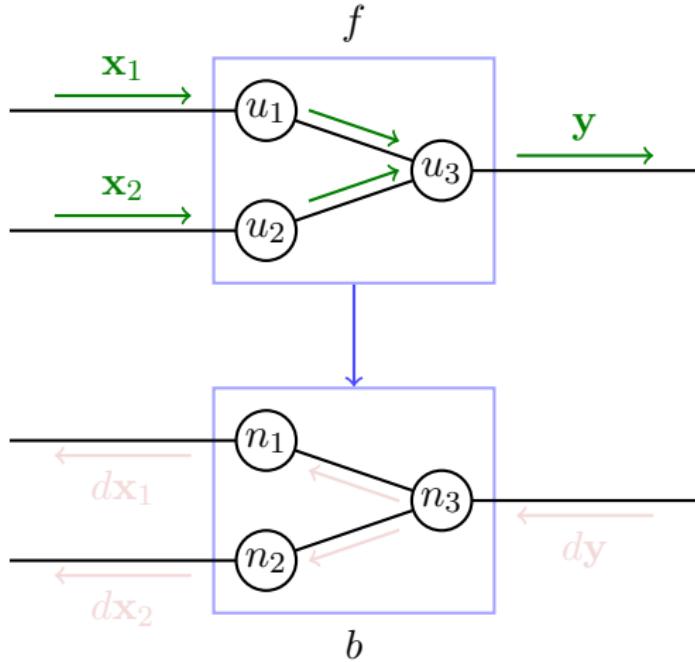
```
def model(X, (U1, U2)):  
    A = affine(relu(affine(X, U1)), U2)  
    return block(A)
```

functions and blocks



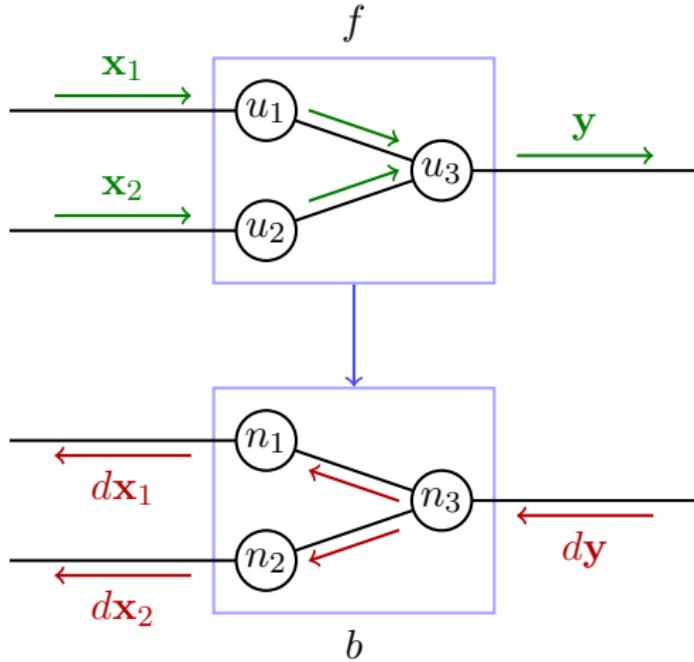
- function f containing units u_1, u_2, u_3
- f dynamically generates block b containing nodes n_1, n_2, n_3 , manually generated by u_1, u_2, u_3 respectively

functions and blocks



- function f containing units u_1, u_2, u_3
- f dynamically generates block b containing nodes n_1, n_2, n_3 , manually generated by u_1, u_2, u_3 respectively

functions and blocks



- function f containing units u_1, u_2, u_3
- f dynamically generates block b containing nodes n_1, n_2, n_3 , manually generated by u_1, u_2, u_3 respectively

functions and blocks

- a *function* is a function of the following form, where code is arbitrary but computation takes place through calls to *units* or *functions*

```
def name(x1, ..., xn):  
    <code generating the following>  
    r1 = call1(a1, ..., an1)  
    :  
    rs = calls(a1, ..., ans)  
    return block(rs)
```

- all calls are recorded as a list of *units* or *functions* by *call order*, each associated with a list of arguments
- a *block* object is a *node*, but
 - its method `back()` does not add its own derivative in the arguments
 - its method `back()` is *automatically generated* and its body calls the recorded functions with the same arguments in *reverse order*

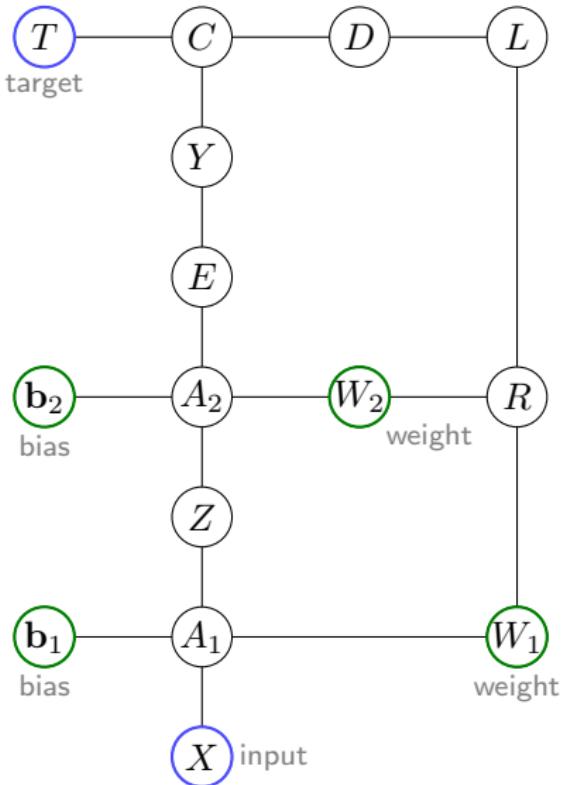
functions and blocks

- a *function* is a function of the following form, where code is arbitrary but computation takes place through calls to *units* or *functions*

```
def name(x1, ..., xn):  
    <code generating the following>  
    r1 = call1(a1, ..., an1)  
    :  
    rs = calls(a1, ..., ans)  
    return block(rs)
```

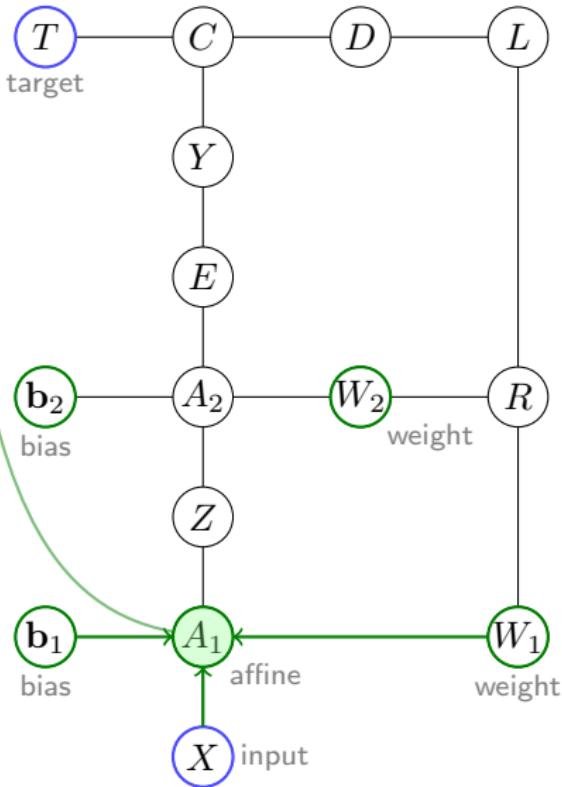
- all calls are recorded as a list of *units* or *functions* by *call order*, each associated with a list of arguments
- a *block* object is a *node*, but
 - its method `back()` does not add its own derivative in the arguments
 - its method `back()` is *automatically generated* and its body calls the recorded functions with the same arguments in *reverse order*

back-propagation



back-propagation

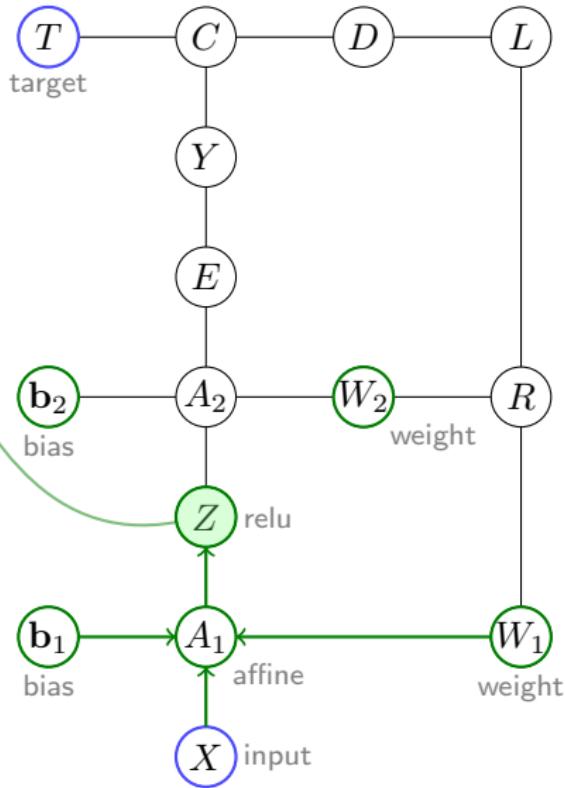
$$A_1 = \text{dot}(X, W_1) + \mathbf{b}_1$$



back-propagation

$$A_1 = \text{dot}(X, W_1) + b_1$$

$$Z = \max(0, A_1)$$

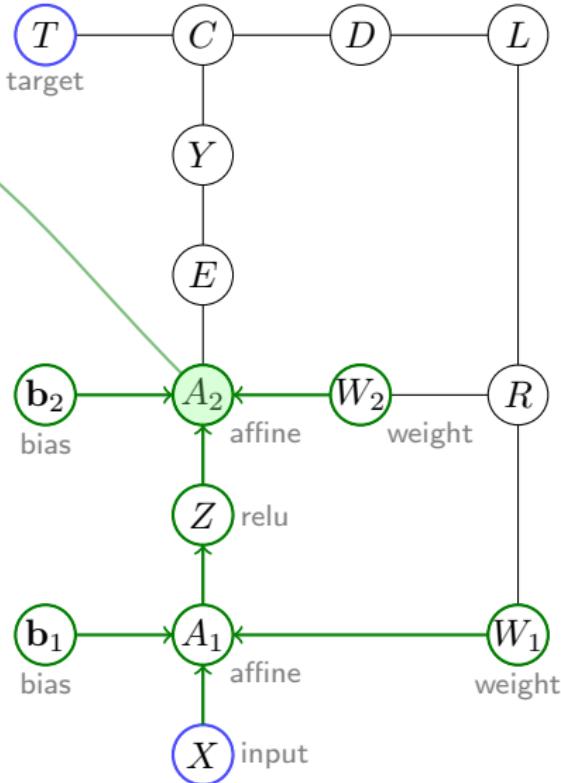


back-propagation

$$A_1 = \text{dot}(X, W_1) + b_1$$

$$Z = \max(0, A_1)$$

$$A_2 = \text{dot}(Z, W_2) + b_2$$



back-propagation

$$A_1 = \text{dot}(X, W_1) + b_1$$

$$Z = \max(0, A_1)$$

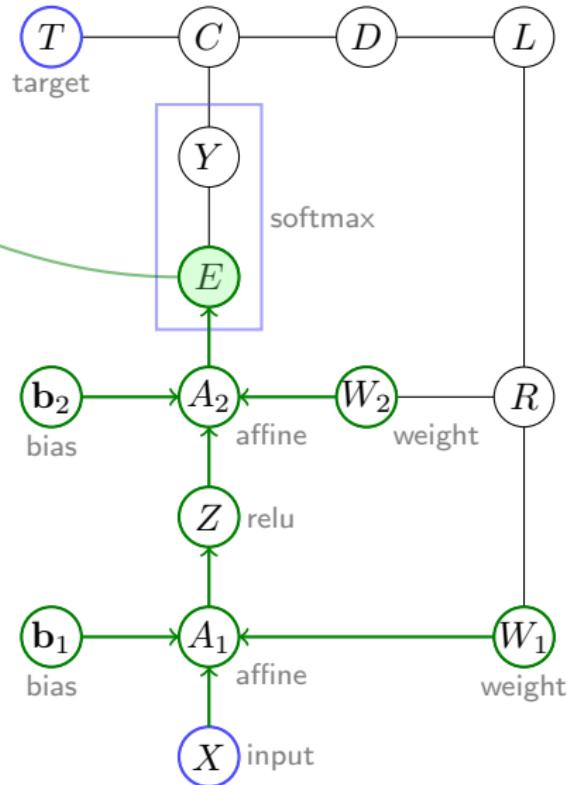
$$A_2 = \text{dot}(Z, W_2) + b_2$$

$$E = \exp(A_2)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$



back-propagation

$$A_1 = \text{dot}(X, W_1) + b_1$$

$$Z = \max(0, A_1)$$

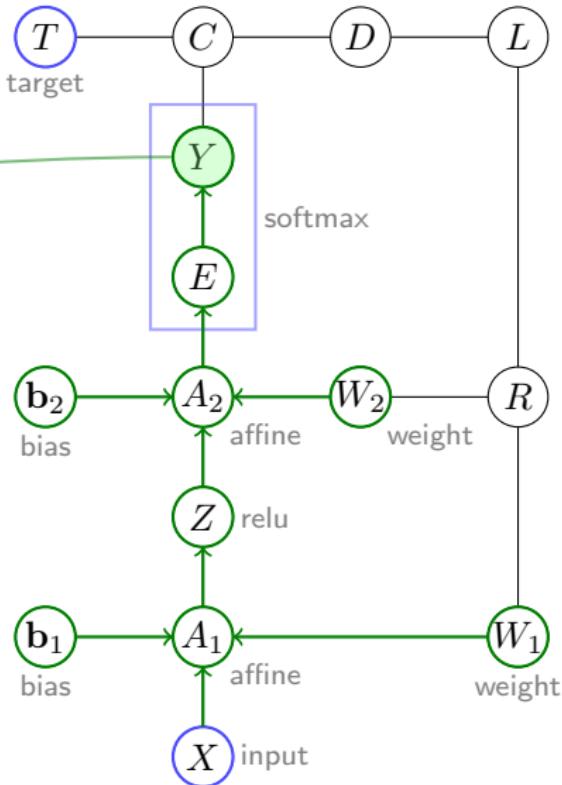
$$A_2 = \text{dot}(Z, W_2) + b_2$$

$$E = \exp(A_2)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$



back-propagation

$$A_1 = \text{dot}(X, W_1) + \mathbf{b}_1$$

$$Z = \max(0, A_1)$$

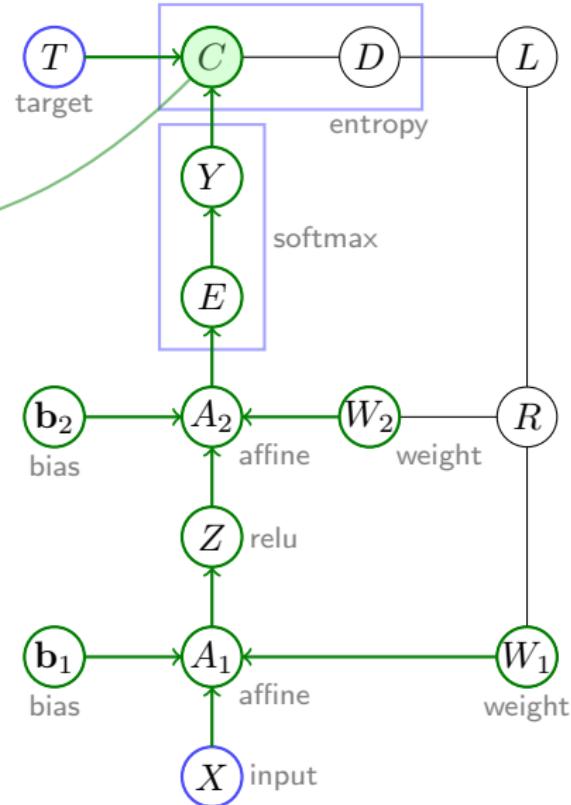
$$A_2 = \text{dot}(Z, W_2) + \mathbf{b}_2$$

$$E = \exp(A_2)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$



back-propagation

$$A_1 = \text{dot}(X, W_1) + \mathbf{b}_1$$

$$Z = \max(0, A_1)$$

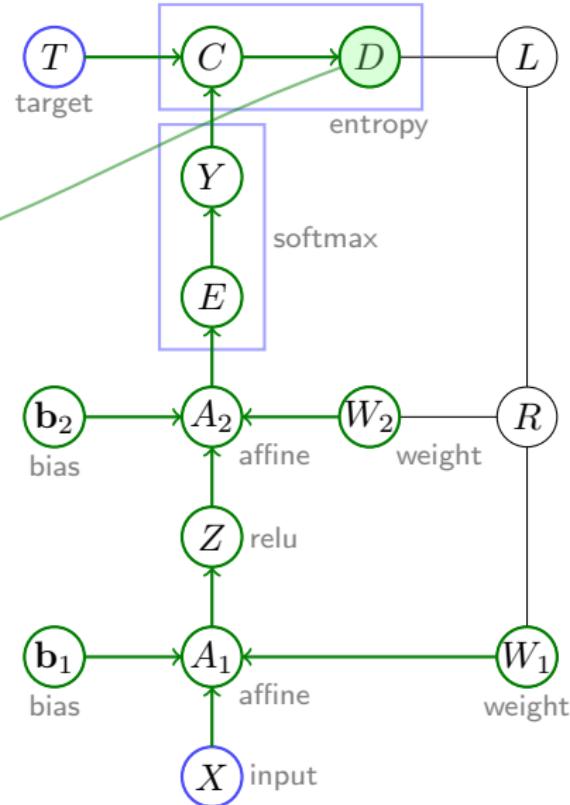
$$A_2 = \text{dot}(Z, W_2) + \mathbf{b}_2$$

$$E = \exp(A_2)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$



back-propagation

$$A_1 = \text{dot}(X, W_1) + \mathbf{b}_1$$

$$Z = \max(0, A_1)$$

$$A_2 = \text{dot}(Z, W_2) + \mathbf{b}_2$$

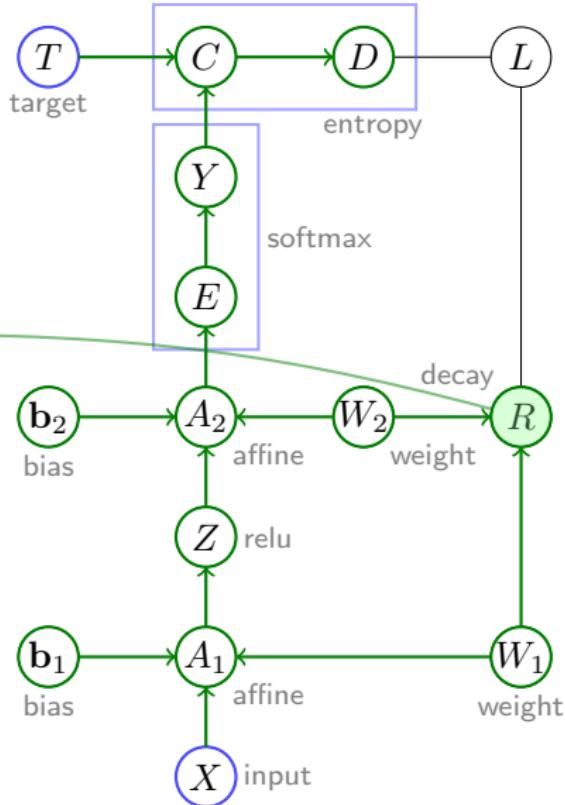
$$E = \exp(A_2)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$

$$R = \frac{\lambda}{2} * (\|W_1\|_F^2 + \|W_2\|_F^2)$$



back-propagation

$$A_1 = \text{dot}(X, W_1) + \mathbf{b}_1$$

$$Z = \max(0, A_1)$$

$$A_2 = \text{dot}(Z, W_2) + \mathbf{b}_2$$

$$E = \exp(A_2)$$

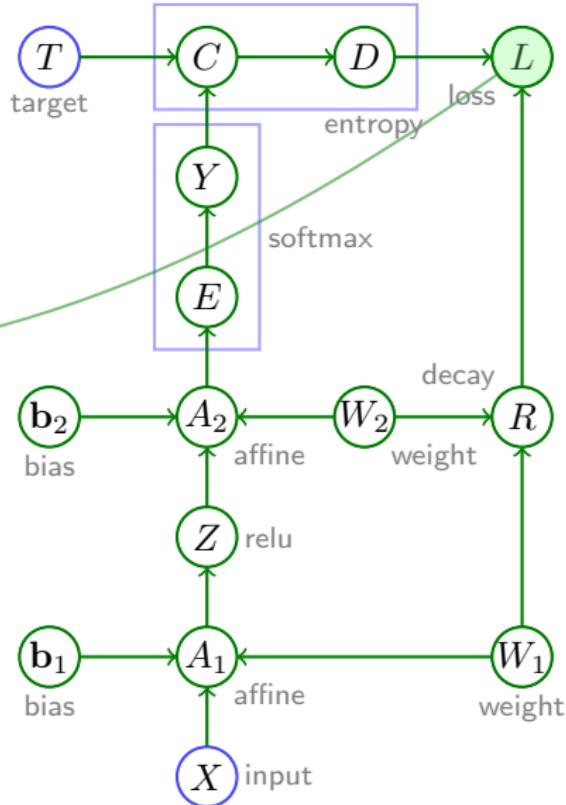
$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$

$$R = \frac{\lambda}{2} * (\|W_1\|_F^2 + \|W_2\|_F^2)$$

$$L = D + R$$



back-propagation

$$A_1 = \text{dot}(X, W_1) + \mathbf{b}_1$$

$$Z = \max(0, A_1)$$

$$A_2 = \text{dot}(Z, W_2) + \mathbf{b}_2$$

$$E = \exp(A_2)$$

$$Y = E / \text{sum}_1(E)$$

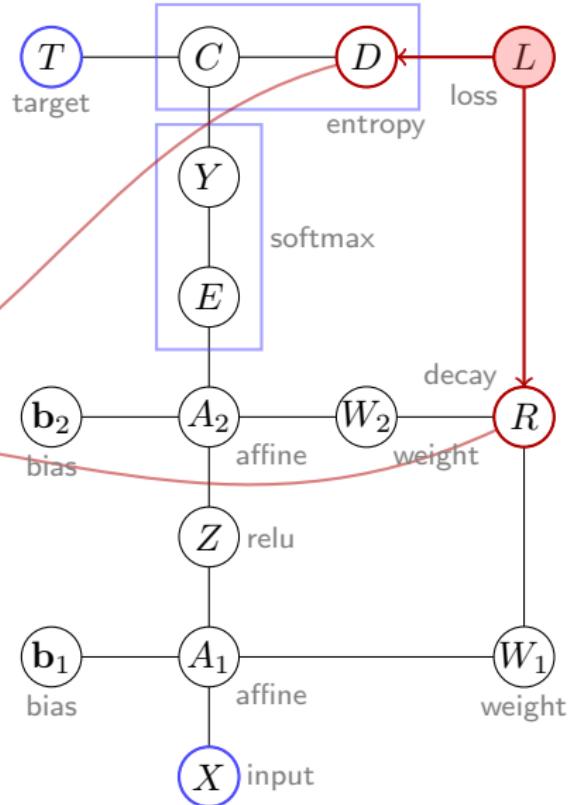
$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$

$$R = \frac{\lambda}{2} * (\|W_1\|_F^2 + \|W_2\|_F^2)$$

$$L = D + R$$

$$(dD, dR) = (dL, dL)$$



back-propagation

$$A_1 = \text{dot}(X, W_1) + \mathbf{b}_1$$

$$Z = \max(0, A_1)$$

$$A_2 = \text{dot}(Z, W_2) + \mathbf{b}_2$$

$$E = \exp(A_2)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$

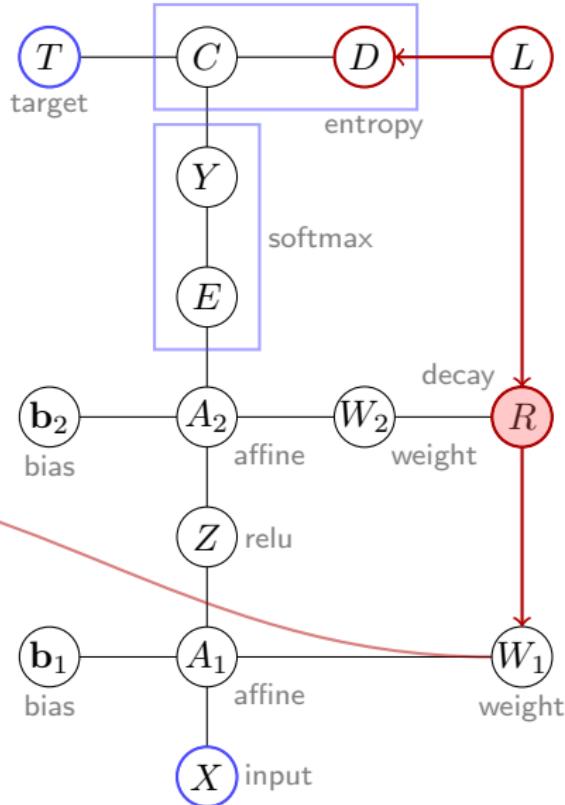
$$R = \frac{\lambda}{2} * (\|W_1\|_F^2 + \|W_2\|_F^2)$$

$$L = D + R$$

$$(dD, dR) = (dL, dL)$$

$$dW_1 = dR * \lambda * W_1$$

$$dW_2 = dR * \lambda * W_2$$



back-propagation

$$A_1 = \text{dot}(X, W_1) + \mathbf{b}_1$$

$$Z = \max(0, A_1)$$

$$A_2 = \text{dot}(Z, W_2) + \mathbf{b}_2$$

$$E = \exp(A_2)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$

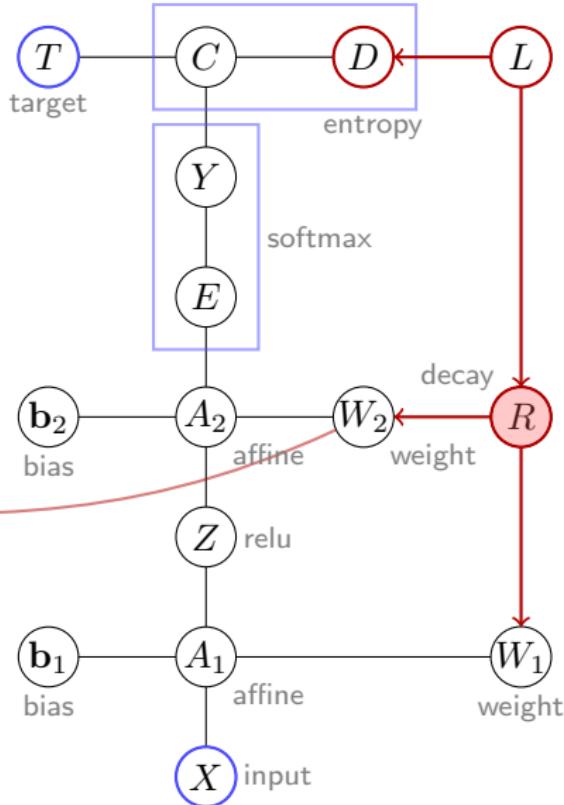
$$R = \frac{\lambda}{2} * (\|W_1\|_F^2 + \|W_2\|_F^2)$$

$$L = D + R$$

$$(dD, dR) = (dL, dL)$$

$$dW_1 = dR * \lambda * W_1$$

$$dW_2 = dR * \lambda * W_2$$



back-propagation

$$A_1 = \text{dot}(X, W_1) + \mathbf{b}_1$$

$$Z = \max(0, A_1)$$

$$A_2 = \text{dot}(Z, W_2) + \mathbf{b}_2$$

$$E = \exp(A_2)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$

$$R = \frac{\lambda}{2} * (\|W_1\|_F^2 + \|W_2\|_F^2)$$

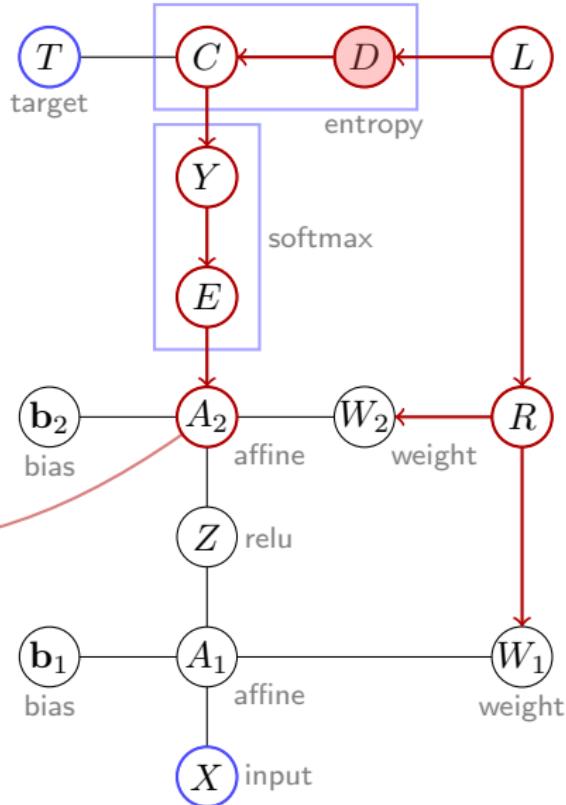
$$L = D + R$$

$$(dD, dR) = (dL, dL)$$

$$dW_1 = dR * \lambda * W_1$$

$$dW_2 = dR * \lambda * W_2$$

$$\boxed{dA_2 = dD * (Y - T)/m}$$



back-propagation

$$A_1 = \text{dot}(X, W_1) + \mathbf{b}_1$$

$$Z = \max(0, A_1)$$

$$A_2 = \text{dot}(Z, W_2) + \mathbf{b}_2$$

$$E = \exp(A_2)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$

$$R = \frac{\lambda}{2} * (\|W_1\|_F^2 + \|W_2\|_F^2)$$

$$L = D + R$$

$$(dD, dR) = (dL, dL)$$

$$dW_1 = dR * \lambda * W_1$$

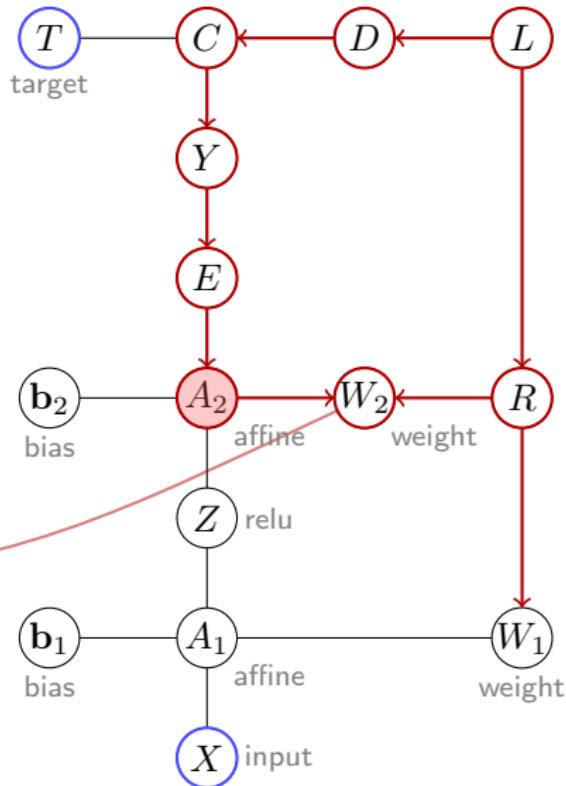
$$dW_2 = dR * \lambda * W_2$$

$$dA_2 = dD * (Y - T)/m$$

$$dW_2 += \text{dot}(Z^\top, dA_2)$$

$$d\mathbf{b}_2 = \text{sum}_0(dA_2)$$

$$dZ = \text{dot}(dA_2, W_2^\top)$$



back-propagation

$$A_1 = \text{dot}(X, W_1) + \mathbf{b}_1$$

$$Z = \max(0, A_1)$$

$$A_2 = \text{dot}(Z, W_2) + \mathbf{b}_2$$

$$E = \exp(A_2)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$

$$R = \frac{\lambda}{2} * (\|W_1\|_F^2 + \|W_2\|_F^2)$$

$$L = D + R$$

$$(dD, dR) = (dL, dL)$$

$$dW_1 = dR * \lambda * W_1$$

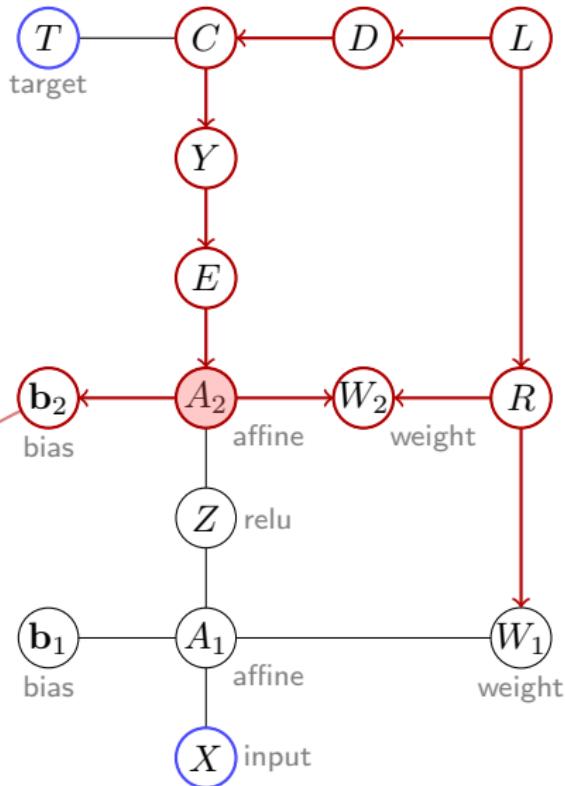
$$dW_2 = dR * \lambda * W_2$$

$$dA_2 = dD * (Y - T)/m$$

$$dW_2 += \text{dot}(Z^\top, dA_2)$$

$$\boxed{db_2 = \text{sum}_0(dA_2)}$$

$$\boxed{dZ = \text{dot}(dA_2, W_2^\top)}$$



back-propagation

$$A_1 = \text{dot}(X, W_1) + \mathbf{b}_1$$

$$Z = \max(0, A_1)$$

$$A_2 = \text{dot}(Z, W_2) + \mathbf{b}_2$$

$$E = \exp(A_2)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$

$$R = \frac{\lambda}{2} * (\|W_1\|_F^2 + \|W_2\|_F^2)$$

$$L = D + R$$

$$(dD, dR) = (dL, dL)$$

$$dW_1 = dR * \lambda * W_1$$

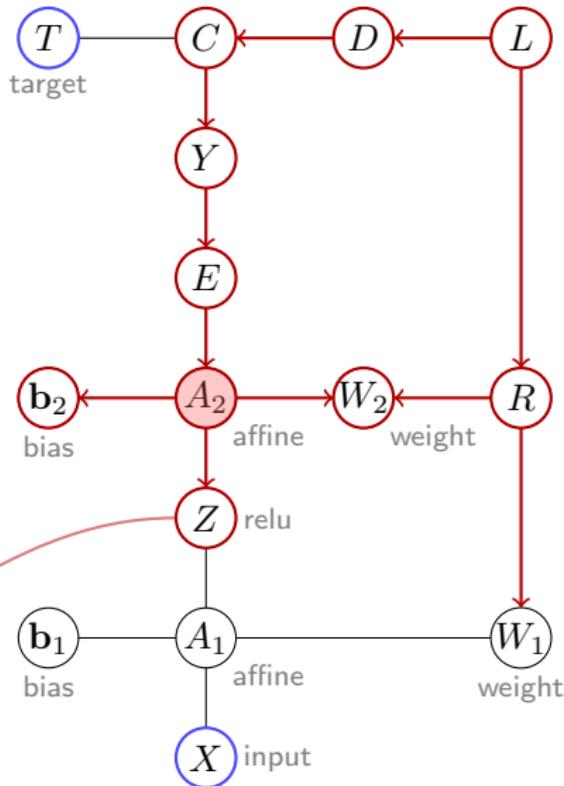
$$dW_2 = dR * \lambda * W_2$$

$$dA_2 = dD * (Y - T)/m$$

$$dW_2 += \text{dot}(Z^\top, dA_2)$$

$$\mathbf{d}\mathbf{b}_2 = \text{sum}_0(dA_2)$$

$$dZ = \text{dot}(dA_2, W_2^\top)$$



back-propagation

$$A_1 = \text{dot}(X, W_1) + \mathbf{b}_1$$

$$Z = \max(0, A_1)$$

$$A_2 = \text{dot}(Z, W_2) + \mathbf{b}_2$$

$$E = \exp(A_2)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$

$$R = \frac{\lambda}{2} * (\|W_1\|_F^2 + \|W_2\|_F^2)$$

$$L = D + R$$

$$(dD, dR) = (dL, dL)$$

$$dW_1 = dR * \lambda * W_1$$

$$dW_2 = dR * \lambda * W_2$$

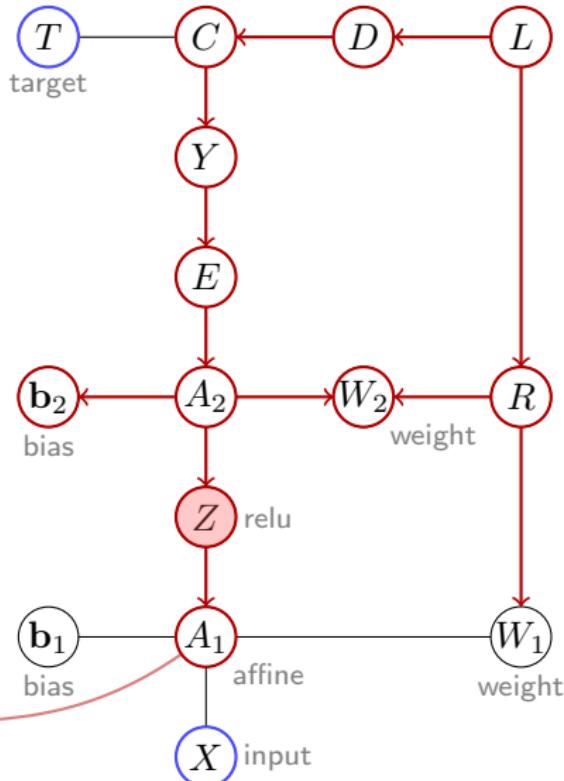
$$dA_2 = dD * (Y - T)/m$$

$$dW_2 += \text{dot}(Z^\top, dA_2)$$

$$d\mathbf{b}_2 = \text{sum}_0(dA_2)$$

$$dZ = \text{dot}(dA_2, W_2^\top)$$

$$dA_1 = dZ * (Z > 0)$$



back-propagation

$$A_1 = \text{dot}(X, W_1) + \mathbf{b}_1$$

$$Z = \max(0, A_1)$$

$$A_2 = \text{dot}(Z, W_2) + \mathbf{b}_2$$

$$E = \exp(A_2)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$

$$R = \frac{\lambda}{2} * (\|W_1\|_F^2 + \|W_2\|_F^2)$$

$$L = D + R$$

$$(dD, dR) = (dL, dL)$$

$$dW_1 = dR * \lambda * W_1$$

$$dW_2 = dR * \lambda * W_2$$

$$dA_2 = dD * (Y - T)/m$$

$$dW_2 += \text{dot}(Z^\top, dA_2)$$

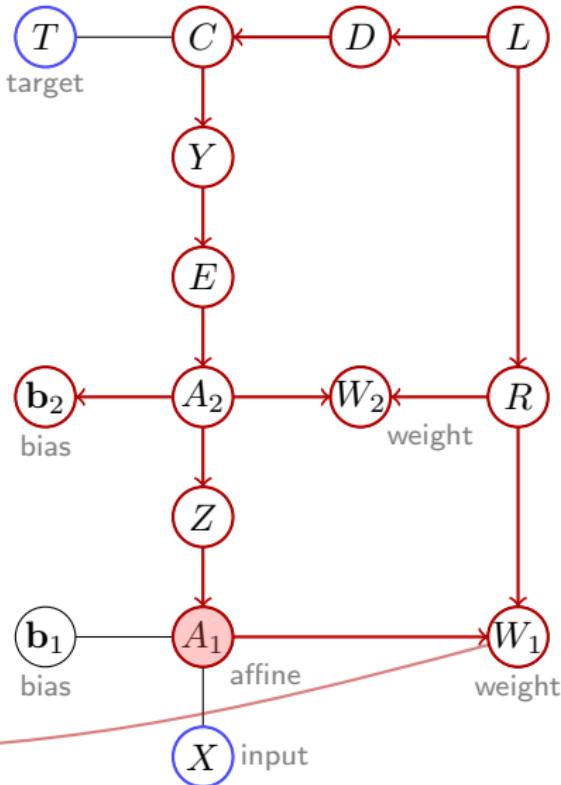
$$d\mathbf{b}_2 = \text{sum}_0(dA_2)$$

$$dZ = \text{dot}(dA_2, W_2^\top)$$

$$dA_1 = dZ * (Z > 0)$$

$$dW_1 += \text{dot}(X^\top, dA_1)$$

$$d\mathbf{b}_1 = \text{sum}_0(dA_1)$$



back-propagation

$$A_1 = \text{dot}(X, W_1) + b_1$$

$$Z = \max(0, A_1)$$

$$A_2 = \text{dot}(Z, W_2) + b_2$$

$$E = \exp(A_2)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C)/m$$

$$R = \frac{\lambda}{2} * (\|W_1\|_F^2 + \|W_2\|_F^2)$$

$$L = D + R$$

$$(dD, dR) = (dL, dL)$$

$$dW_1 = dR * \lambda * W_1$$

$$dW_2 = dR * \lambda * W_2$$

$$dA_2 = dD * (Y - T)/m$$

$$dW_2 += \text{dot}(Z^\top, dA_2)$$

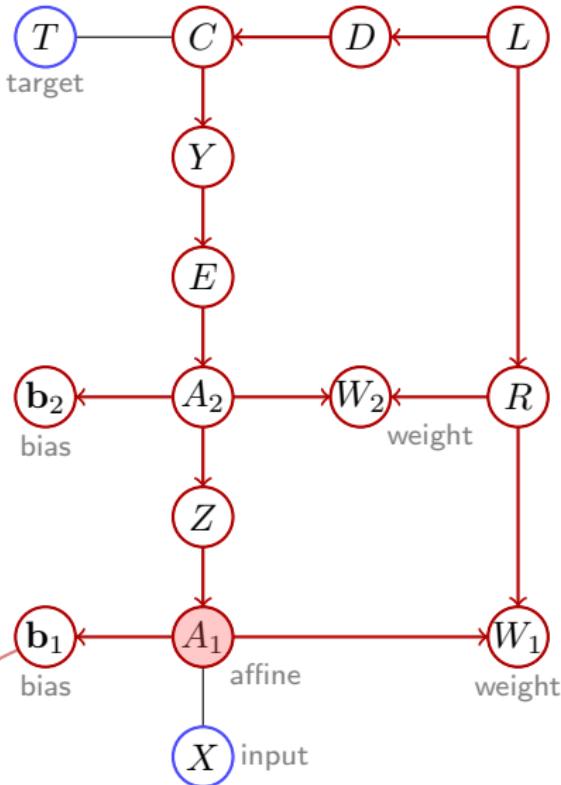
$$db_2 = \text{sum}_0(dA_2)$$

$$dZ = \text{dot}(dA_2, W_2^\top)$$

$$dA_1 = dZ * (Z > 0)$$

$$dW_1 += \text{dot}(X^\top, dA_1)$$

$$db_1 = \text{sum}_0(dA_1)$$



automatic differentiation

```
A1 = dot(X, W1) + b1
Z = max(0, A1)
A2 = dot(Z, W2) + b2
E = exp(A2)
Y = E / sum1(E)
C = -sum1(T * log(Y))
D = sum0(C) / m
R = λ / 2 * (||W1||F2 + ||W2||F2)
L = D + R
```

```
(dD, dR) = (dL, dL)
dW1 = dR * λ * W1
dW2 = dR * λ * W2
dA2 = dD * (Y - T) / m
dW2 += dot(ZT, dA2)
db2 = sum0(dA2)
dZ = dot(dA2, W2T)
dA1 = dZ * (Z > 0)
dW1 += dot(XT, dA1)
db1 = sum0(dA1)
```

now we organize **forward** and **backward** code into units and functions

automatic differentiation

```
A1 = dot(X, W1) + b1
Z = max(0, A1)
A2 = dot(Z, W2) + b2
E = exp(A2)
Y = E/sum1(E)
C = -sum1(T * log(Y))
D = sum0(C)/m
R = λ/2 * (||W1||F2 + ||W2||F2)
L = D + R
(dD, dR) = (dL, dL)
dW1 = dR * λ * W1
dW2 = dR * λ * W2
dA2 = dD * (Y - T)/m
dW2 += dot(ZT, dA2)
db2 = sum0(dA2)
dZ = dot(dA2, W2T)
dA1 = dZ * (Z > 0)
dW1 += dot(XT, dA1)
db1 = sum0(dA1)
```

```
def relu(A):
    Z = max(0, A)
    def back(dZ, dA):
        dA += dZ * (Z > 0)
    return node(Z, back)
```

automatic differentiation

$$A_1 = \text{dot}(X, W_1) + \mathbf{b}_1$$

$$Z = \text{relu}(A_1)$$

$$A_2 = \text{dot}(Z, W_2) + \mathbf{b}_2$$

$$E = \exp(A_2)$$

$$Y = E / \text{sum}_1(E)$$

$$C = -\text{sum}_1(T * \log(Y))$$

$$D = \text{sum}_0(C) / m$$

$$R = \frac{\lambda}{2} * (\|W_1\|_F^2 + \|W_2\|_F^2)$$

$$L = D + R$$

$$(dD, dR) = (dL, dL)$$

$$dW_1 = dR * \lambda * W_1$$

$$dW_2 = dR * \lambda * W_2$$

$$dA_2 = dD * (Y - T) / m$$

$$dW_2 += \text{dot}(Z^\top, dA_2)$$

$$d\mathbf{b}_2 = \text{sum}_0(dA_2)$$

$$dZ = \text{dot}(dA_2, W_2^\top)$$

$$Z.\text{back}(A_1)$$

$$dW_1 += \text{dot}(X^\top, dA_1)$$

$$d\mathbf{b}_1 = \text{sum}_0(dA_1)$$

```
def relu(A):
```

```
    Z = max(0, A)
```

```
def back(dZ, dA):
```

```
    dA += dZ * (Z > 0)
```

```
return node(Z, back)
```

automatic differentiation

```
A1 = dot(X, W1) + b1
Z = relu(A1)
A2 = dot(Z, W2) + b2
E = exp(A2)
Y = E / sum1(E)
C = -sum1(T * log(Y))
D = sum0(C) / m
R =  $\frac{\lambda}{2} * (\|W_1\|_F^2 + \|W_2\|_F^2)$ 
L = D + R
(dD, dR) = (dL, dL)
dW1 = dR * λ * W1
dW2 = dR * λ * W2
dA2 = dD * (Y - T) / m
dW2 += dot(Z⊤, dA2)
db2 = sum0(dA2)
dZ = dot(dA2, W2⊤)
Z.back(A1)
dW1 += dot(X⊤, dA1)
db1 = sum0(dA1)
```

```
def affine(X, (W, b)):
    A = dot(X, W) + b
    def back(dA, dX, (dW, db)):
        dW += dot(X⊤, dA)
        db += sum0(dA)
        dX += dot(dA, W⊤)
    return node(A, back)
```

automatic differentiation

```
A1 = affine(X, (W1, b1))  
Z = relu(A1)  
A2 = affine(Z, (W2, b2))  
E = exp(A2)  
Y = E / sum1(E)  
C = -sum1(T * log(Y))  
D = sum0(C) / m  
R =  $\frac{\lambda}{2} * (\|W_1\|_F^2 + \|W_2\|_F^2)$   
L = D + R  
(dD, dR) = (dL, dL)  
dW1 = dR * λ * W1  
dW2 = dR * λ * W2  
dA2 = dD * (Y - T) / m  
A2.back(Z, (W2, b2))
```

```
def affine(X, (W, b)):  
    A = dot(X, W) + b  
  
def back(dA, dX, (dW, db)):  
    dW += dot(X⊤, dA)  
    db += sum0(dA)  
    dX += dot(dA, W⊤)  
  
return node(A, back)
```

Z.back(A₁)

A₁.back(X, (W₁, b₁))

automatic differentiation

$A_1 = \text{affine}(X, (W_1, \mathbf{b}_1))$

$Z = \text{relu}(A_1)$

$A_2 = \text{affine}(Z, (W_2, \mathbf{b}_2))$

$E = \exp(A_2)$

$Y = E / \text{sum}_1(E)$

$C = -\text{sum}_1(T * \log(Y))$

$D = \text{sum}_0(C) / m$

$R = \frac{\lambda}{2} * (\|W_1\|_F^2 + \|W_2\|_F^2)$

$L = D + R$

$(dD, dR) = (dL, dL)$

$dW_1 = dR * \lambda * W_1$

$dW_2 = dR * \lambda * W_2$

$dA_2 = dD * (Y - T) / m$

$A_2.\text{back}(Z, (W_2, \mathbf{b}_2))$

def logistic(A, T):

$E = \exp(A)$

$Y = E / \text{sum}_1(E)$

$C = -\text{sum}_1(T * \log(Y))$

$D = \text{sum}_0(C) / m$

def back($dD, dA, _$):

$dA += dD * (Y - T) / m$

return node(D, back)

$Z.\text{back}(A_1)$

$A_1.\text{back}(X, (W_1, \mathbf{b}_1))$

automatic differentiation

$A_1 = \text{affine}(X, (W_1, \mathbf{b}_1))$

$Z = \text{relu}(A_1)$

$A_2 = \text{affine}(Z, (W_2, \mathbf{b}_2))$

$D = \text{logistic}(A_2, T)$

$$R = \frac{\lambda}{2} * (\|W_1\|_F^2 + \|W_2\|_F^2)$$

$$L = D + R$$

$$(dD, dR) = (dL, dL)$$

$$dW_1 = dR * \lambda * W_1$$

$$dW_2 = dR * \lambda * W_2$$

$D.\text{back}(A_2, T)$
 $A_2.\text{back}(Z, (W_2, \mathbf{b}_2))$

```
def logistic(A, T):
    E = exp(A)
    Y = E/sum_1(E)
    C = -sum_1(T * log(Y))
    D = sum_0(C)/m

def back(dD, dA, _):
    dA += dD * (Y - T)/m
    return node(D, back)
```

$Z.\text{back}(A_1)$
 $A_1.\text{back}(X, (W_1, \mathbf{b}_1))$

automatic differentiation

$A_1 = \text{affine}(X, (W_1, \mathbf{b}_1))$

$Z = \text{relu}(A_1)$

$A_2 = \text{affine}(Z, (W_2, \mathbf{b}_2))$

$D = \text{logistic}(A_2, T)$

$$R = \frac{\lambda}{2} * (\|W_1\|_F^2 + \|W_2\|_F^2)$$

$$L = D + R$$

$$(dD, dR) = (dL, dL)$$

$$dW_1 = dR * \lambda * W_1$$

$$dW_2 = dR * \lambda * W_2$$

$D.\text{back}(A_2, T)$

$A_2.\text{back}(Z, (W_2, \mathbf{b}_2))$

def decay(W):

$$R = \frac{\lambda}{2} * \text{sum}(\|w\|_F^2 \text{ for } w \text{ in } W)$$

def back(dR, dW):

for (w, dw) **in** $\text{zip}(W, dW)$:
 $dw += dR * \lambda * w$

return $node(R, \text{back})$

$Z.\text{back}(A_1)$

$A_1.\text{back}(X, (W_1, \mathbf{b}_1))$

automatic differentiation

$A_1 = \text{affine}(X, (W_1, \mathbf{b}_1))$

$Z = \text{relu}(A_1)$

$A_2 = \text{affine}(Z, (W_2, \mathbf{b}_2))$

$D = \text{logistic}(A_2, T)$

$R = \text{decay}((W_1, W_2))$

$L = D + R$

$(dD, dR) = (dL, dL)$

$R.\text{back}((W_1, W_2))$

$D.\text{back}(A_2, T)$

$A_2.\text{back}(Z, (W_2, \mathbf{b}_2))$

$Z.\text{back}(A_1)$

$A_1.\text{back}(X, (W_1, \mathbf{b}_1))$

```
def decay(W):
    R =  $\frac{\lambda}{2} * \text{sum}(\|w\|_F^2)$  for w in W
    def back(dR, dW):
        for (w, dw) in zip(W, dW):
            dw += dR *  $\lambda * w$ 
    return node(R, back)
```

automatic differentiation

$A_1 = \text{affine}(X, (W_1, \mathbf{b}_1))$

$Z = \text{relu}(A_1)$

$A_2 = \text{affine}(Z, (W_2, \mathbf{b}_2))$

$D = \text{logistic}(A_2, T)$

$R = \text{decay}((W_1, W_2))$

$L = D + R$

$(dD, dR) = (dL, dL)$

$R.\text{back}((W_1, W_2))$

$D.\text{back}(A_2, T)$

$A_2.\text{back}(Z, (W_2, \mathbf{b}_2))$

$Z.\text{back}(A_1)$

$A_1.\text{back}(X, (W_1, \mathbf{b}_1))$

```
def add(X):
    S = sum(X)
    def back(dS, dX):
        for dx in dX:
            dx += dS
    return node(S, back)
```

automatic differentiation

$A_1 = \text{affine}(X, (W_1, \mathbf{b}_1))$

$Z = \text{relu}(A_1)$

$A_2 = \text{affine}(Z, (W_2, \mathbf{b}_2))$

$D = \text{logistic}(A_2, T)$

$R = \text{decay}((W_1, W_2))$

$L = \text{add}((D, R))$

$L.\text{back}((D, R))$

$R.\text{back}((W_1, W_2))$

$D.\text{back}(A_2, T)$

$A_2.\text{back}(Z, (W_2, \mathbf{b}_2))$

$Z.\text{back}(A_1)$

$A_1.\text{back}(X, (W_1, \mathbf{b}_1))$

def add(X):

$S = \text{sum}(X)$

def back(dS, dX):

for dx **in** dX :

$dx += dS$

return node(S , back)

automatic differentiation

$A_1 = \text{affine}(X, (W_1, \mathbf{b}_1))$

$Z = \text{relu}(A_1)$

$A_2 = \text{affine}(Z, (W_2, \mathbf{b}_2))$

$D = \text{logistic}(A_2, T)$

$R = \text{decay}((W_1, W_2))$

$L = \text{add}((D, R))$

$L.\text{back}((D, R))$

$R.\text{back}((W_1, W_2))$

$D.\text{back}(A_2, T)$

$A_2.\text{back}(Z, (W_2, \mathbf{b}_2))$

$Z.\text{back}(A_1)$

$A_1.\text{back}(X, (W_1, \mathbf{b}_1))$

```
def loss(A, T, W):
    D = logistic(A, T)
    R = decay(W)
    L = add((D, R))
    def back(A, T, W):
        L.back((D, R))
        R.back(W)
        D.back(A, T)
    return block(L, back)
```

automatic differentiation

$A_1 = \text{affine}(X, (W_1, \mathbf{b}_1))$

$Z = \text{relu}(A_1)$

$A_2 = \text{affine}(Z, (W_2, \mathbf{b}_2))$

$L = \text{loss}(A_2, T, (W_1, W_2))$

$L.\text{back}(A_2, T, (W_1, W_2))$

$A_2.\text{back}(Z, (W_2, \mathbf{b}_2))$

$Z.\text{back}(A_1)$

$A_1.\text{back}(X, (W_1, \mathbf{b}_1))$

```
def loss(A, T, W):
    D = logistic(A, T)
    R = decay(W)
    L = add((D, R))

def back(A, T, W):
    L.back((D, R))
    R.back(W)
    D.back(A, T)

return block(L, back)
```

automatic differentiation

$A_1 = \text{affine}(X, (W_1, \mathbf{b}_1))$

$Z = \text{relu}(A_1)$

$A_2 = \text{affine}(Z, (W_2, \mathbf{b}_2))$

$L = \text{loss}(A_2, T, (W_1, W_2))$

$L.\text{back}(A_2, T, (W_1, W_2))$

$A_2.\text{back}(Z, (W_2, \mathbf{b}_2))$

$Z.\text{back}(A_1)$

$A_1.\text{back}(X, (W_1, \mathbf{b}_1))$

```
def loss(A, T, W):
    L = logistic(A, T) + decay(W)
    return block(L)
```

```
def loss(A, T, W):
    D = logistic(A, T)
    R = decay(W)
    L = add((D, R))
```

```
def back(A, T, W):
    L.back((D, R))
    R.back(W)
    D.back(A, T)
    return block(L, back)
```

automatic differentiation

```
A1 = affine(X, (W1, b1))  
Z = relu(A1)  
A2 = affine(Z, (W2, b2))  
L = loss(A2, T, (W1, W2))
```

$L.$.back($A_2, T, (W_1, W_2)$)

$A_2.$.back($Z, (W_2, b_2)$)

$Z.$.back(A_1)
 $A_1.$.back($X, (W_1, b_1)$)

def model($X, (U_1, U_2)$):

```
A1 = affine(X, U1)  
Z = relu(A)  
A2 = affine(Z, U2)
```

def back($X, (U_1, U_2)$):

```
A2.back(Z, U2)  
Z.back(A)  
A1.back(X, U1)
```

return block(A_2 , back)

automatic differentiation

```
A2 = model(X, ((W1, b1), (W2, b2)))
```

```
L = loss(A2, T, (W1, W2))
```

```
L.back(A2, T, (W1, W2))
```

```
A2.back(X, ((W1, b1), (W2, b2)))
```

```
def model(X, (U1, U2)):  
    A1 = affine(X, U1)  
    Z = relu(A)  
    A2 = affine(Z, U2)  
  
def back(X, (U1, U2)):  
    A2.back(Z, U2)  
    Z.back(A)  
    A1.back(X, U1)  
    return block(A2, back)
```

automatic differentiation

$A_2 = \text{model}(X, ((W_1, \mathbf{b}_1), (W_2, \mathbf{b}_2)))$

$L = \text{loss}(A_2, T, (W_1, W_2))$

$L.\text{back}(A_2, T, (W_1, W_2))$

$A_2.\text{back}(X, ((W_1, \mathbf{b}_1), (W_2, \mathbf{b}_2)))$

```
def model(X, (U1, U2)):  
    A = affine(relu(affine(X, U1)), U2)  
    return block(A)
```

```
def model(X, (U1, U2)):  
    A1 = affine(X, U1)  
    Z = relu(A)  
    A2 = affine(Z, U2)
```

```
def back(X, (U1, U2)):  
    A2.back(Z, U2)  
    Z.back(A)  
    A1.back(X, U1)  
    return block(A2, back)
```

what is not shown

- *units* and *functions* are actually **objects**
 - both expose a “call” operator () for the forward pass, which handles recording of calls
 - *function*’s operator also generates method **back** and constructs a *block* object for the backward pass
 - both expose method **init**, which handles **initialization** of parameters, if any

deep learning software

Caffe



Caffe²



PYTORCH



TensorFlow

theano

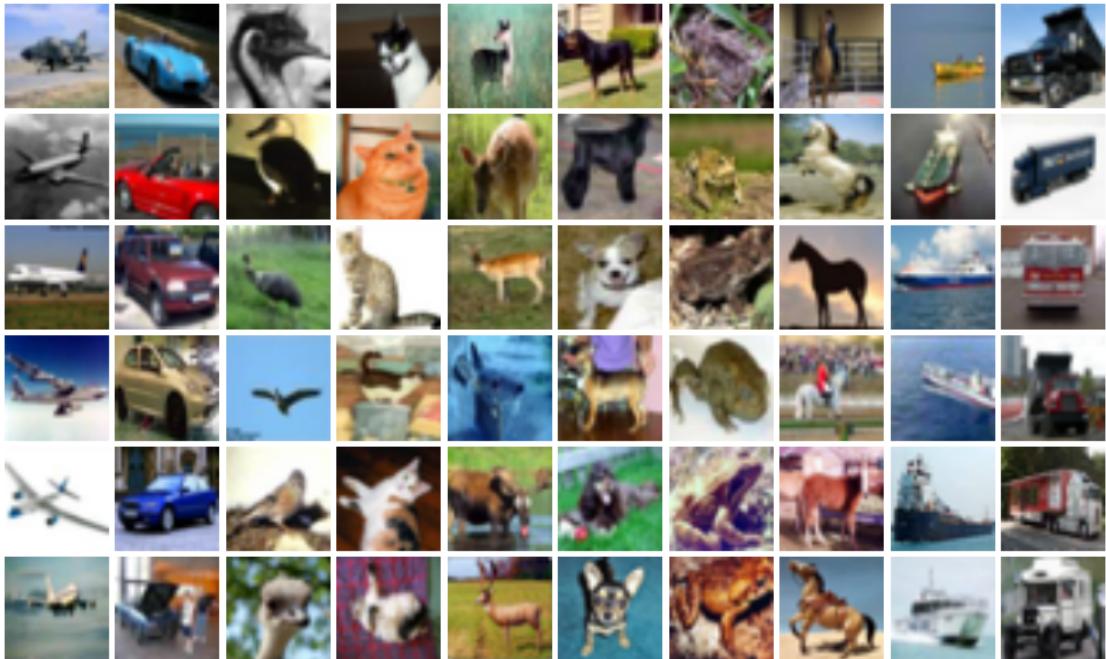


dmlc
mxnet

- automatically build computational graphs and compute derivatives
- run on GPU, multiple GPU, distributed
- component (unit, layer) libraries
- pre-trained models
- community

fun

CIFAR10 dataset



plane car bird cat deer dog frog horse ship truck

- 10 classes, 50k training images, 10k test images, 32×32 images

pipeline

prepare

- **vectorize** $32 \times 32 \times 3$ images into 3072×1
- **split** training set e.g. into $n_{\text{train}} = 45000$ training samples and $n_{\text{val}} = 5000$ samples to be used for validation
- **center** vectors by subtracting mean over the training samples
- **initialize** network weights as Gaussian with standard deviation 10^{-4}

learn

- train for a few iterations and evaluate accuracy on the **validation** set for a number of learning rates ϵ and regularization strengths λ
- **train** for 10 epochs on the full training set for the chosen hyperparameters
- evaluate accuracy on the **test** set

pipeline

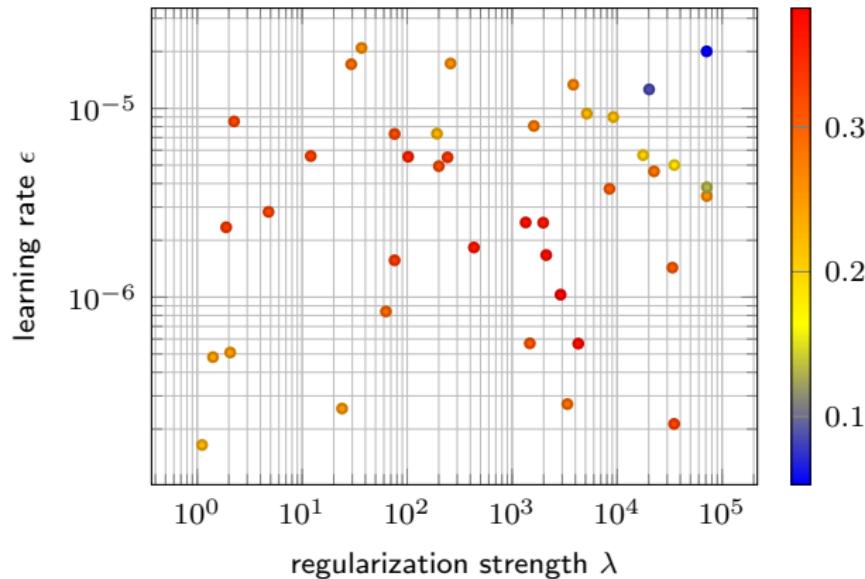
prepare

- **vectorize** $32 \times 32 \times 3$ images into 3072×1
- **split** training set e.g. into $n_{\text{train}} = 45000$ training samples and $n_{\text{val}} = 5000$ samples to be used for validation
- **center** vectors by subtracting mean over the training samples
- **initialize** network weights as Gaussian with standard deviation 10^{-4}

learn

- train for a few iterations and evaluate accuracy on the **validation** set for a number of learning rates ϵ and regularization strengths λ
- **train** for 10 epochs on the full training set for the chosen hyperparameters
- evaluate accuracy on the **test** set

linear classifier validation accuracy



- classes $k = 10$, samples $n_{\text{train}} = 45000, n_{\text{val}} = 5000$, mini-batch $m = 200$, learning rate $\epsilon = 10^{-6}$, regularization strength $\lambda = 5 \times 10^2$
- test accuracy: 38%

linear classifier weights



plane



car



bird



cat



deer



dog



frog



horse

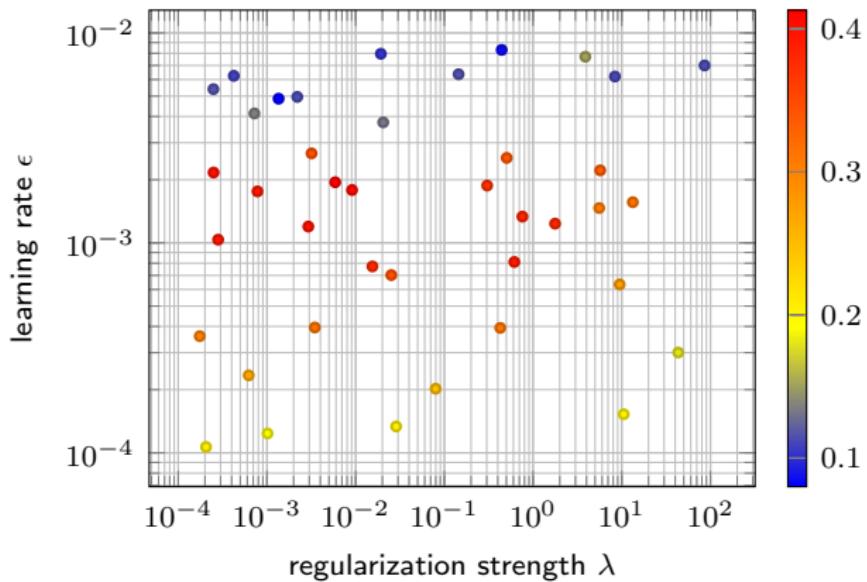


ship



truck

2-layer classifier validation accuracy



- classes $k = 10$, samples $n_{\text{train}} = 45000, n_{\text{val}} = 5000$, mini-batch $m = 200$, learning rate $2 \times \epsilon = 10^{-3}$, regularization strength $\lambda = 2 \times 10^{-1}$
- hidden layer width: 100; test accuracy: 51%

two-layer classifier weights

layer 1 weights 0-49



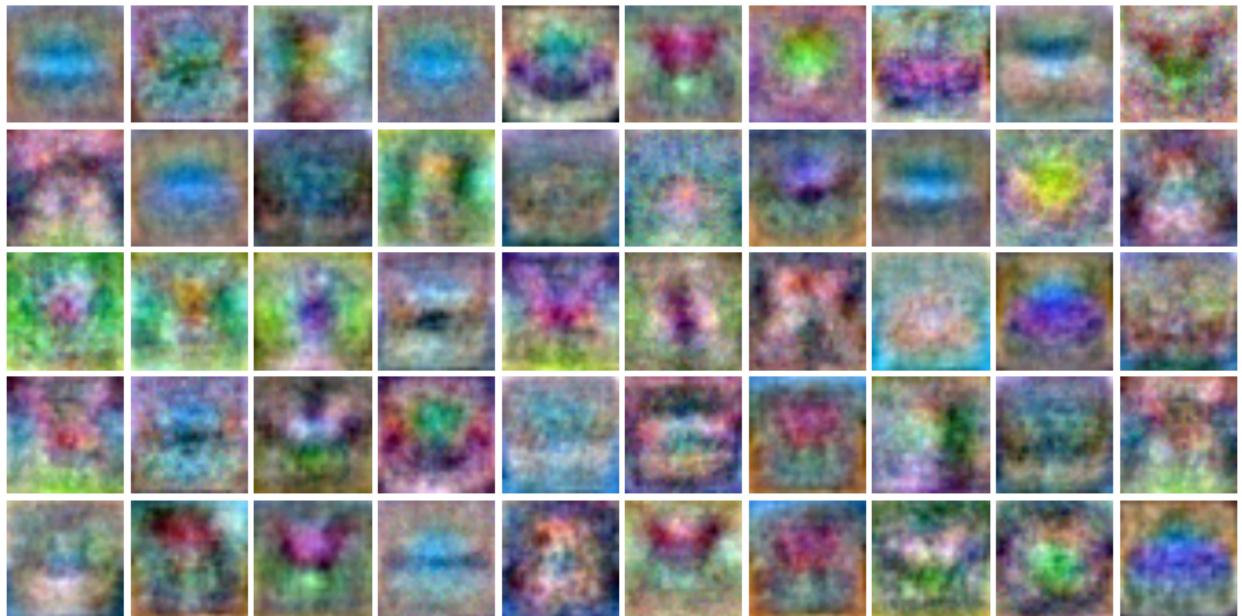
two-layer classifier weights

layer 1 weights 50-99



two-layer classifier weights

layer 1 weights 100-149

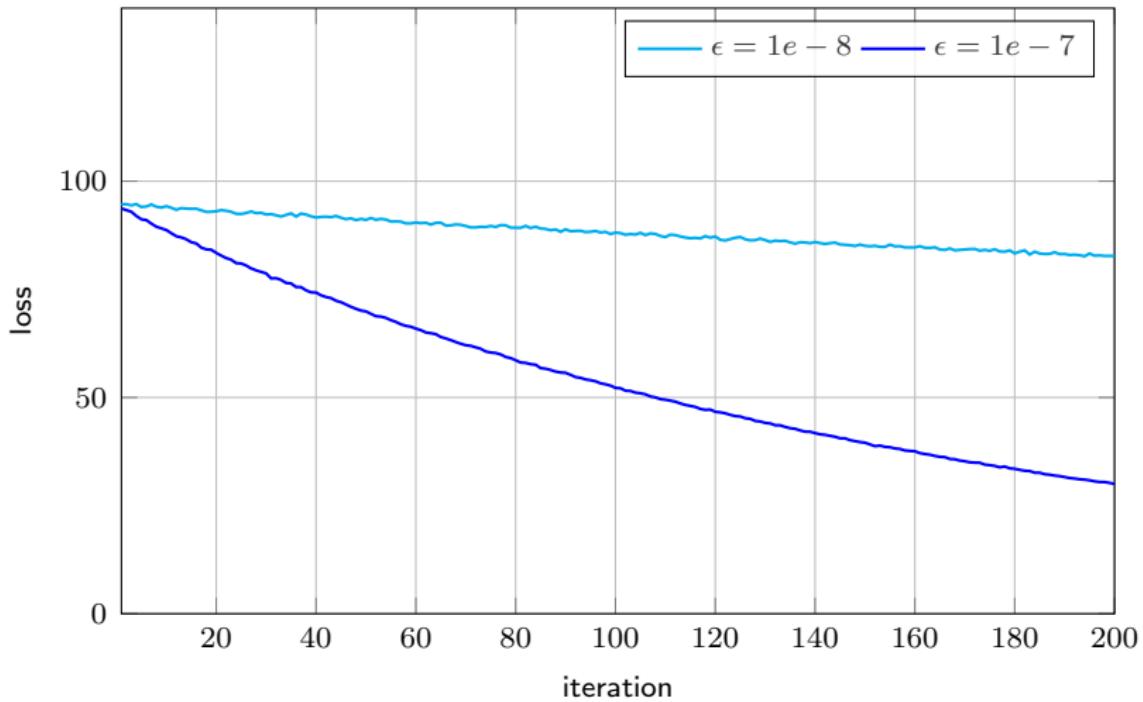


two-layer classifier weights

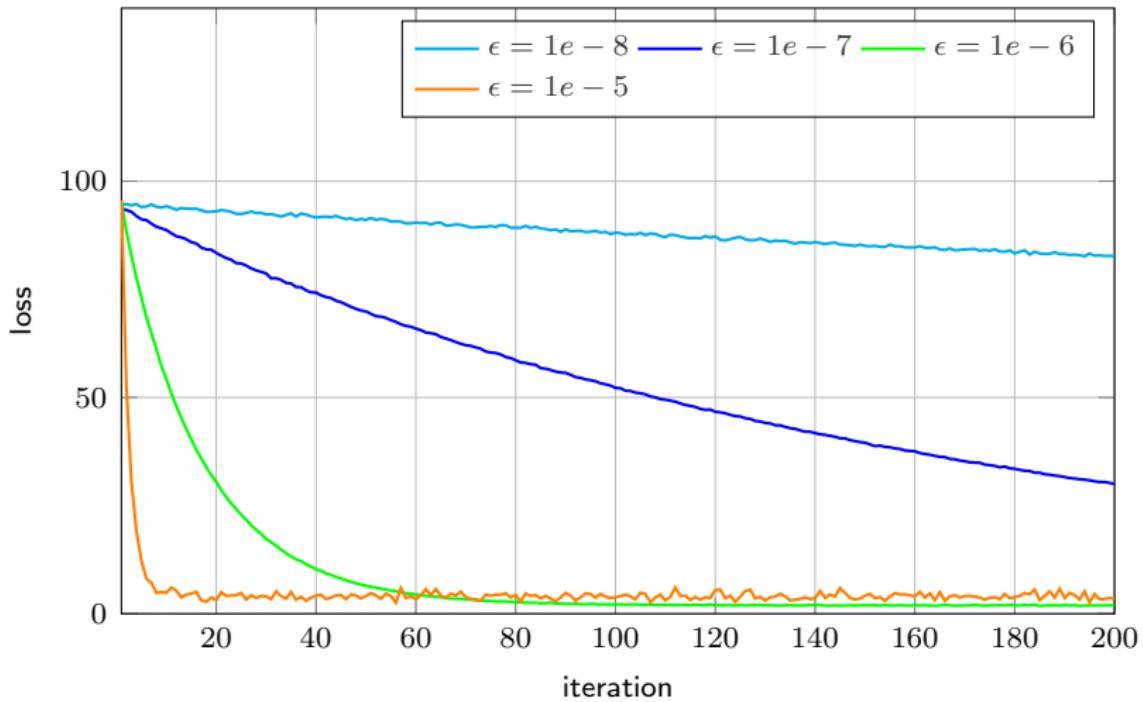
layer 1 weights 150-199



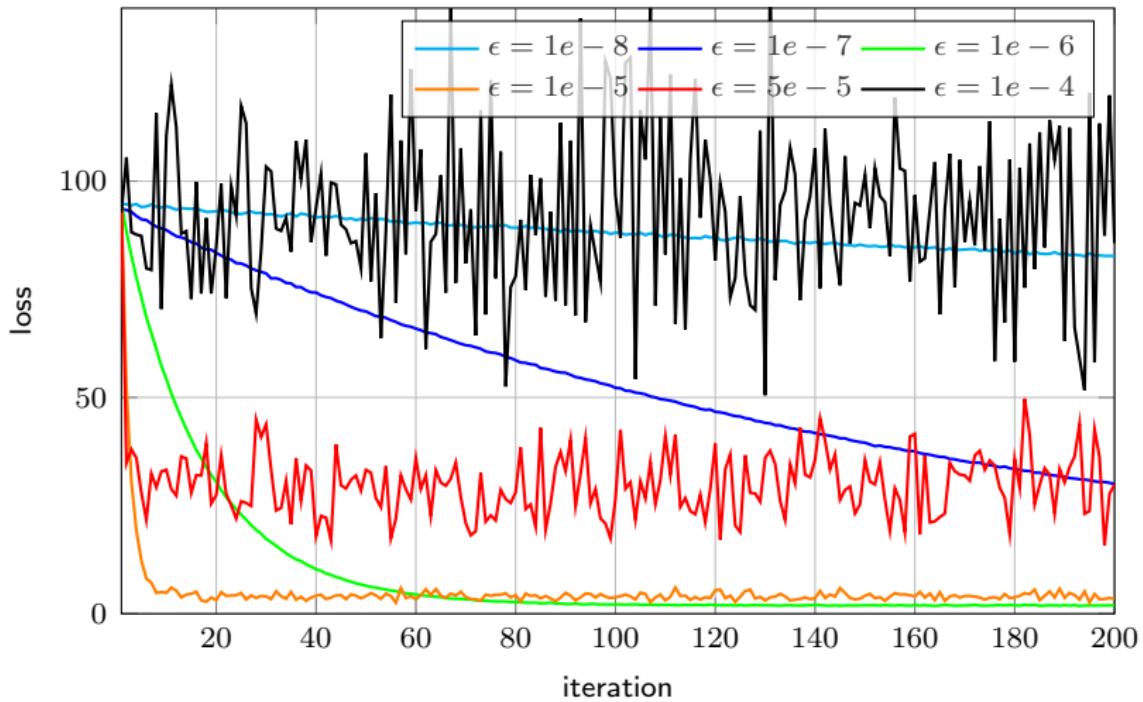
learning rate



learning rate

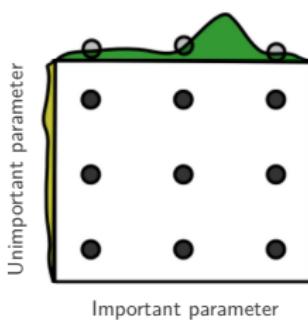


learning rate

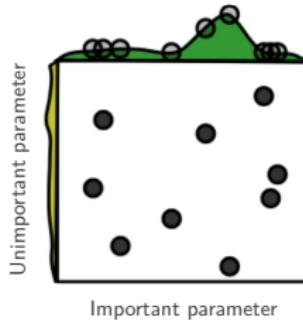


setting hyperparameters

Grid Layout



Random Layout



- compared to grid search, random search allows to explore more values of an important parameter regardless of unimportant parameters
- when the search spans orders of magnitude, draw samples uniformly at random in log space
- start with coarse range and few iterations, gradually move to finer range and more iterations

summary

- stochastic gradient descent and its limitations
- numerical gradient approximation
- analytical computation by decomposing and applying the chain rule
- back-propagation as dynamic programming
- chaining, splitting and sharing
- common patterns between forward and backward flow
- decomposition into unit (forward) and nodes (backward)
- grouping into functions (forward) and blocks (backward)
- learning CIFAR10, validation and hyperparameters