Deep Learning Gradient-based optimization

Caio Corro

Université Paris-Saclay

Table of contents

Recall: neural networks

The training loop

Backpropagation

Vanishing gradient, activation functions and initialization

Regularization

Better optimizers

Recall: neural networks

Neural network

- ► x: input features
- > $z^{(1)}, z^{(2)}$: hidden representation
- ▶ w: output logits or class weights
- **p**: probability distribution over classes
- $\theta = \{ A^{(1)}, b^{(1)}, ... \}$: parameters
- σ : non-linear activation function

Neural network

- ► x: input features
- > $z^{(1)}, z^{(2)}$: hidden representation
- ▶ w: output logits or class weights
- **p**: probability distribution over classes
- $\theta = \{ \boldsymbol{A}^{(1)}, \boldsymbol{b}^{(1)}, \ldots \}$: parameters
- σ : non-linear activation function

$$z^{(1)} = \sigma \left(\mathbf{A}^{(1)} \mathbf{x} + \mathbf{b}^{(1)} \right)$$

$$z^{(2)} = \sigma \left(\mathbf{A}^{(2)} \mathbf{z}^{(1)} + \mathbf{b}^{(2)} \right)$$

$$\mathbf{w} = \sigma \left(\mathbf{A}^{(3)} \mathbf{z}^{(2)} + \mathbf{b}^{(3)} \right)$$

$$\mathbf{p} = \text{Softmax}(\mathbf{w}) \quad \text{i.e.} \quad p_i = \frac{\exp(\mathbf{w}_i)}{\sum_j \exp(\mathbf{w}_j)}$$

Neural network

- ► x: input features
- > $z^{(1)}, z^{(2)}$: hidden representation
- ▶ w: output logits or class weights
- *p*: probability distribution over classes
 θ = {*A*⁽¹⁾, *b*⁽¹⁾, ...}: parameters
- σ : non-linear activation function

$$z^{(1)} = \sigma \left(\mathbf{A}^{(1)} \mathbf{x} + \mathbf{b}^{(1)} \right)$$

$$z^{(2)} = \sigma \left(\mathbf{A}^{(2)} z^{(1)} + \mathbf{b}^{(2)} \right)$$

$$\mathbf{w} = \sigma \left(\mathbf{A}^{(3)} z^{(2)} + \mathbf{b}^{(3)} \right)$$

$$\mathbf{p} = \text{Softmax}(\mathbf{w}) \quad \text{i.e.} \quad p_i = \frac{\exp(\mathbf{w}_i)}{\sum_j \exp(\mathbf{w}_j)}$$



Representation learning: Computer Vision [Lee et al., 2009]



Representation learning: Natural Language Processing [Voita et al., 2019]



The training loop

The big picture

Data split and usage

- ▶ Training set: to learn the parameters of the network
- Development (or dev or validation) set: to monitor the network during training
- Test set: to evaluate our model at the end

Generally you don't have to split the data yourself: there exists standard splits to allow benchmarking.

The big picture

Data split and usage

- Training set: to learn the parameters of the network
- Development (or dev or validation) set: to monitor the network during training
- Test set: to evaluate our model at the end

Generally you don't have to split the data yourself: there exists standard splits to allow benchmarking.

Training loop

- 1. Update the parameters to minimize the loss on the training set
- 2. Evaluate the prediction accuracy on the dev set
- 3. If not satisfied, go back to 1
- 4. Evaluate the prediction accuracy on the test set with the best parameters on dev

function TRAIN (f, θ, T, D)

```
function TRAIN(f, \theta, \mathcal{T}, \mathcal{D})
bestdev = -\infty
for epoch = 1 to E do
Shuffle \mathcal{T}
for x, y \in \mathcal{T} do
loss = \mathcal{L}(f(x; \theta), y)
\theta = \theta - \epsilon \nabla loss
```

```
function TRAIN(f, \theta, T, D)
    bestdev = -\infty
    for epoch = 1 to E do
         Shuffle \mathcal{T}
         for x, y \in \mathcal{T} do
              loss = \mathcal{L}(f(x; \theta), y)
              \theta = \theta - \epsilon \nabla \log \theta
         devacc = EVALUATE(f, \mathcal{D})
         if devacc > bestdev then
              \hat{\theta} = \theta
              bestdev = devacc
```

return $\hat{\theta}$

 $\begin{aligned} & \textbf{function } \operatorname{Train}(f,\theta,\mathcal{T},\mathcal{D}) \\ & \texttt{bestdev} = -\infty \\ & \textbf{for } \texttt{epoch} = 1 \texttt{ to } E \texttt{ do} \\ & \texttt{Shuffle } \mathcal{T} \\ & \textbf{for } x, y \in \mathcal{T} \texttt{ do} \\ & \texttt{loss} = \mathcal{L}(f(x;\theta),y) \\ & \theta = \theta - \epsilon \nabla \texttt{loss} \end{aligned}$

devacc =EVALUATE(f, D)if devacc > bestdev then $\hat{\theta} = \theta$ bestdev = devacc function EVALUATE(f, D) n = 0for $x, y \in D$ do $\hat{y} = \arg \max_y f(x; \theta)_y$ if $\hat{y} = y$ then n = n + 1return n/|D|

return $\hat{\theta}$

Further details

Sampling without replacement

- shuffle the training set
- loop over the new order

Experimentally it works better than "true" sampling and it seems to also have good theoretical properties [Nagaraj et al., 2019]

Verbosity

At each epoch, it is useful to display:

- mean loss
- accuracy on training data
- accuracy on dev data
- timing information
- (sometimes) evaluate on dev several times by epoch

$$\theta^{(t+1)} = \theta^{(t)} - \epsilon^{(t)} \nabla \text{loss} \qquad \Rightarrow \qquad \text{How to choose the step size } \epsilon^{(t+1)}?$$

 $\theta^{(t+1)} = \theta^{(t)} - \epsilon^{(t)} \nabla \text{loss} \implies \text{How to choose the step size } \epsilon^{(t+1)}?$ Convex optimization

▶ Nonsummable diminishing step size: $\sum_{t=1}^{\infty} \epsilon^{(t)} = \infty$ and $\lim_{t\to\infty} \epsilon^{(t)} = 0$

Backtracking/exact line search

 $\theta^{(t+1)} = \theta^{(t)} - \epsilon^{(t)} \nabla \text{loss} \implies \text{How to choose the step size } \epsilon^{(t+1)}?$ Convex optimization

- ▶ Nonsummable diminishing step size: $\sum_{t=1}^{\infty} \epsilon^{(t)} = \infty$ and $\lim_{t\to\infty} \epsilon^{(t)} = 0$
- Backtracking/exact line search

Simple neural network heuristic

- 1. Start with a small value, e.g. $\epsilon=0.01$
- 2. If dev accuracy did not improve during the last N epochs: decay the learning rate by a small value α , e.g. $\epsilon = \alpha * \epsilon$ with $\alpha = 0.1$

 $\theta^{(t+1)} = \theta^{(t)} - \epsilon^{(t)} \nabla \text{loss} \implies \text{How to choose the step size } \epsilon^{(t+1)}?$ Convex optimization

- ▶ Nonsummable diminishing step size: $\sum_{t=1}^{\infty} \epsilon^{(t)} = \infty$ and $\lim_{t\to\infty} \epsilon^{(t)} = 0$
- Backtracking/exact line search

Simple neural network heuristic

- 1. Start with a small value, e.g. $\epsilon=0.01$
- 2. If dev accuracy did not improve during the last N epochs: decay the learning rate by a small value α , e.g. $\epsilon = \alpha * \epsilon$ with $\alpha = 0.1$

Step-size annealing

▶ Step decay: multiple ϵ by $\alpha \in [0, 1]$ every N epochs

• Exponential decay:
$$\epsilon^{(t)} = \epsilon^{(0)} \exp(-\alpha \cdot t)$$

►
$$1/t$$
 decay: $\epsilon^{(t)} = \frac{\epsilon^{(0)}}{1+\alpha}$

Backpropagation

Derivative

Let $f : \mathbb{R} \to \mathbb{R}$ be a function and $x, y \in \mathbb{R}$ be variables such that:

y = f(x).

For a given x, how does an infinitesimal change of x impact y?

Derivative

Let $f : \mathbb{R} \to \mathbb{R}$ be a function and $x, y \in \mathbb{R}$ be variables such that:

y = f(x).

For a given x, how does an infinitesimal change of x impact y?

$$\frac{dy}{dx} = f'(x) = \lim_{\epsilon \to 0} \frac{f(x+\epsilon) - f(x)}{\epsilon}$$

Derivative

Let $f : \mathbb{R} \to \mathbb{R}$ be a function and $x, y \in \mathbb{R}$ be variables such that:

y = f(x).

For a given x, how does an infinitesimal change of x impact y?

$$\frac{dy}{dx} = f'(x) = \lim_{\epsilon \to 0} \frac{f(x+\epsilon) - f(x)}{\epsilon}$$

Linear approximation

Let $\tilde{f} : \mathcal{R} \to \mathcal{R}$ be function parameterized by $a \in \mathcal{R}$ defined as follows:

$$\widetilde{f}(x; a) = f(a) + f'(a) \cdot (x - a)$$

Then, $\tilde{f}(x; a)$ is an approximation of f at a.

Example $f(x) = x^{2} + 2$ f'(x) = 2x $\tilde{f}(x; a) = f(a) + f'(a) \cdot (x - a)$ $= a^{2} + 2 + 2a(x - a)$ $= 2ax + 2 - a^{2}$

Example

$$f(x) = x^{2} + 2$$

$$f'(x) = 2x$$

$$\tilde{f}(x; a) = f(a) + f'(a) \cdot (x - a)$$

$$= a^{2} + 2 + 2a(x - a)$$

$$= 2ax + 2 - a^{2}$$

Intuition: the sign of f'(a) gives the slope of the approximation, we can use this information to move closer to the minimum of f(x).



▶ Black: f(x)
 ▶ Red: *f*(x; a = -6)

Example

$$f(x) = x^{2} + 2$$

$$f'(x) = 2x$$

$$\widetilde{f}(x; a) = f(a) + f'(a) \cdot (x - a)$$

$$= a^{2} + 2 + 2a(x - a)$$

$$= 2ax + 2 - a^{2}$$

Intuition: the sign of f'(a) gives the slope of the approximation, we can use this information to move closer to the minimum of f(x).



• Red:
$$\tilde{f}(x; a = 3)$$

Chain rule

Let $f : \mathbb{R} \to \mathbb{R}$ and $g : \mathbb{R} \to \mathbb{R}$ be two functions and x, y, z be variables such that:

$$z = f(x),$$

 $y = g(z)$ i.e. $y = g(f(x)) = g \circ f(x).$

For a given x, how does an infinitesimal change of x impact y?

Chain rule

Let $f : \mathbb{R} \to \mathbb{R}$ and $g : \mathbb{R} \to \mathbb{R}$ be two functions and x, y, z be variables such that:

$$z = f(x),$$

 $y = g(z)$ i.e. $y = g(f(x)) = g \circ f(x).$

For a given x, how does an infinitesimal change of x impact y?

$$\frac{dy}{dx} = \frac{dy}{dz} \cdot \frac{dz}{dx}$$

Example: explicit differentiation

$$f(x) = (2x+1)^2$$

Example: explicit differentiation

$$f(x) = (2x+1)^2 = 4x^2 + 4x + 1$$

Example: explicit differentiation

$$f(x) = (2x + 1)^2 = 4x^2 + 4x + 1$$

$$f'(x) = 8x + 4$$

Example: explicit differentiation

$$f(x) = (2x + 1)^2 = 4x^2 + 4x + 1$$

$$f'(x) = 8x + 4$$

Example: differentiation using the chain rule

$$z = 2x + 1$$
$$y = z^2 = f(x)$$

Example: explicit differentiation

$$f(x) = (2x + 1)^2 = 4x^2 + 4x + 1$$

$$f'(x) = 8x + 4$$

Example: differentiation using the chain rule

$$z = 2x + 1$$

$$y = z^{2} = f(x)$$

$$\frac{dz}{dx} = 2$$

$$\frac{dy}{dz} = 2z$$

1_

Example: explicit differentiation

$$f(x) = (2x + 1)^2 = 4x^2 + 4x + 1$$

$$f'(x) = 8x + 4$$

1-

Example: differentiation using the chain rule

$$z = 2x + 1$$

$$y = z^{2} = f(x)$$

$$\frac{dy}{dx} = \frac{dy}{dz} \cdot \frac{dz}{dx}$$

Example: explicit differentiation

$$f(x) = (2x + 1)^2 = 4x^2 + 4x + 1$$

$$f'(x) = 8x + 4$$

1-

Example: differentiation using the chain rule

$$z = 2x + 1$$

$$y = z^{2} = f(x)$$

$$\frac{dy}{dx} = \frac{dy}{dz} \cdot \frac{dz}{dx} = 2z * 2$$
Scalar input

Example: explicit differentiation

$$f(x) = (2x + 1)^2 = 4x^2 + 4x + 1$$

$$f'(x) = 8x + 4$$

,

Example: differentiation using the chain rule

$$z = 2x + 1$$

$$y = z^{2} = f(x)$$

$$\frac{dy}{dx} = \frac{dy}{dz} \cdot \frac{dz}{dx} = 2z * 2 = 4(2x + 1)$$

Scalar input

Example: explicit differentiation

$$f(x) = (2x + 1)^2 = 4x^2 + 4x + 1$$

$$f'(x) = 8x + 4$$

Example: differentiation using the chain rule

$$z = 2x + 1$$

$$y = z^{2} = f(x)$$

$$\frac{dz}{dx} = 2$$

$$\frac{dy}{dz} = 2z$$

$$\frac{dy}{dx} = \frac{dy}{dz} \cdot \frac{dz}{dx} = 2z \cdot 2 = 4(2x+1) = 8x + 4 = f'(x)$$

Let $f : \mathbb{R}^m \to \mathbb{R}$ be a function and $\mathbf{x} \in \mathbb{R}^m, y \in \mathbb{R}$ be variables such that:

 $y = f(\mathbf{x}).$

Let $f : \mathbb{R}^m \to \mathbb{R}$ be a function and $\mathbf{x} \in \mathbb{R}^m, y \in \mathbb{R}$ be variables such that:

 $y = f(\mathbf{x}).$

Partial derivative

For a given \boldsymbol{x} , how does an infinitesimal change of x_i impact y?

$$\frac{\partial y}{\partial x_i}$$

i.e. each input $x_j, j \neq i$ is considered as a constant.

Let $f : \mathbb{R}^m \to \mathbb{R}$ be a function and $\mathbf{x} \in \mathbb{R}^m, y \in \mathbb{R}$ be variables such that:

 $y = f(\mathbf{x}).$

Partial derivative

For a given \boldsymbol{x} , how does an infinitesimal change of x_i impact y?

$$\frac{\partial y}{\partial x_i}$$

i.e. each input $x_j, j \neq i$ is considered as a constant.

Gradient

For a given x, how does an infinitesimal change of x impact y?

$$\nabla_{\mathbf{x}} y = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \dots \end{bmatrix}$$

Chain rule

Let $f : \mathbb{R}^m \to \mathbb{R}^n$ and $g : \mathbb{R}^n \to \mathbb{R}$ be two functions and x^m, z^n, y be variables such that:

 $m{z} = f(m{x}),$ $m{y} = m{g}(m{z})$

For a given x_i , how does an infinitesimal change of x_i impact y?

Chain rule

Let $f : \mathbb{R}^m \to \mathbb{R}^n$ and $g : \mathbb{R}^n \to \mathbb{R}$ be two functions and x^m, z^n, y be variables such that:

$$m{z} = f(m{x}),$$

 $m{y} = m{g}(m{z})$

For a given x_i , how does an infinitesimal change of x_i impact y?

$$\frac{\partial y}{\partial x_i} = \sum_j \frac{\partial y}{\partial z_j} \cdot \frac{\partial z_j}{\partial x_i}$$

$$z = Wx + b$$
 or $z_j = \sum_i W_{j,i}x_i + b_j$
 $y = \sum_j z_j$

$$m{z} = m{W} x + b$$
 or $m{z}_j = \sum_i W_{j,i} x_i + b_j$ $rac{\partial z_j}{x_i} = W_{j,i}$
 $y = \sum_j z_j$ $rac{\partial y}{z_j} = 1$

$$m{z} = m{W} x + b$$
 or $m{z}_j = \sum_i W_{j,i} x_i + b_j$ $rac{\partial z_j}{x_i} = W_{j,i}$
 $y = \sum_j z_j$ $rac{\partial y}{z_j} = 1$

$$\frac{\partial y}{\partial x_i} = \sum_j \frac{\partial y}{\partial z_j} \cdot \frac{\partial z_j}{\partial x_i} = \sum_j 1 * W_{j,i}$$

$$z^{(1)} = ... x...$$

 $z^{(2)} = ... z^{(1)}...$
 $y = ... z^{(2)}...$



$$z^{(1)} = ... x ...$$

 $z^{(2)} = ... z^{(1)} ...$
 $y = ... z^{(2)} ...$

$$\frac{\partial y}{\partial x_i} = \sum_k \frac{\partial y}{\partial z_k^{(2)}} \cdot \frac{\partial z_k^{(2)}}{\partial x_i}$$

$$z^{(1)} = ... x ...$$

 $z^{(2)} = ... z^{(1)} ...$
 $y = ... z^{(2)} ...$

$$\frac{\partial y}{\partial x_i} = \sum_k \frac{\partial y}{\partial z_k^{(2)}} \cdot \frac{\partial z_k^{(2)}}{\partial x_i} = \sum_k \frac{\partial y}{\partial z_k^{(2)}} \cdot \sum_j \frac{\partial z_k^{(2)}}{\partial z_j^{(1)}} \cdot \frac{\partial z_j^{(1)}}{\partial x_i}$$

 \Rightarrow It is starting to get annoying!

Jacobian

Let $f : \mathbb{R}^m \to \mathbb{R}^n$ be a function and $\mathbf{x} \in \mathbb{R}^m$, $\mathbf{y} \in \mathbb{R}^n$ be variables such that:

 $\mathbf{y} = f(\mathbf{x}).$

Jacobian

Let $f : \mathbb{R}^m \to \mathbb{R}^n$ be a function and $\mathbf{x} \in \mathbb{R}^m$, $\mathbf{y} \in \mathbb{R}^n$ be variables such that:

 $\boldsymbol{y} = f(\boldsymbol{x}).$

Gradient

For a given \boldsymbol{x} , how does an infinitesimal change of \boldsymbol{x} impact y_i ?

$$\nabla_{\mathbf{x}} y_j = \begin{bmatrix} \frac{\partial y_j}{\partial x_1} \\ \frac{\partial y_j}{\partial x_2} \\ \dots \end{bmatrix}$$

Jacobian

Let $f : \mathbb{R}^m \to \mathbb{R}^n$ be a function and $\mathbf{x} \in \mathbb{R}^m, \mathbf{y} \in \mathbb{R}^n$ be variables such that:

 $\boldsymbol{y} = f(\boldsymbol{x}).$

Gradient

For a given x, how does an infinitesimal change of x impact y_j ?

$$\nabla_{\mathbf{x}} y_j = \begin{bmatrix} \frac{\partial y_j}{\partial x_1} \\ \frac{\partial y_j}{\partial x_2} \\ \dots \end{bmatrix}$$

Jacobian

For a given **x**, how does an infinitesimal change of **x** impact **y**?

$$\mathbf{J}_{\mathbf{x}}\mathbf{y} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \dots \\ \dots & \dots & \dots \end{bmatrix}$$

Chain rule using the Jacobian notation

Let $f : \mathbb{R}^m \to \mathbb{R}^n$ and $g : \mathbb{R}^n \to \mathbb{R}$ be two functions and $\mathbf{x}^m, \mathbf{z}^n, y$ be variables such that:

$$\boldsymbol{z} = f(\boldsymbol{x}), \quad \boldsymbol{y} = g(\boldsymbol{z})$$

Chain rule using the Jacobian notation

Let $f : \mathbb{R}^m \to \mathbb{R}^n$ and $g : \mathbb{R}^n \to \mathbb{R}$ be two functions and $\mathbf{x}^m, \mathbf{z}^n, y$ be variables such that:

$$\boldsymbol{z} = f(\boldsymbol{x}), \quad \boldsymbol{y} = g(\boldsymbol{z})$$

Partial notation

$$\frac{\partial y}{\partial x_i} = \sum_j \frac{\partial y}{\partial z_j} \cdot \frac{\partial z_j}{\partial x_i}$$

Chain rule using the Jacobian notation

Let $f : \mathbb{R}^m \to \mathbb{R}^n$ and $g : \mathbb{R}^n \to \mathbb{R}$ be two functions and $\mathbf{x}^m, \mathbf{z}^n, y$ be variables such that:

$$\boldsymbol{z} = f(\boldsymbol{x}), \quad \boldsymbol{y} = g(\boldsymbol{z})$$

Partial notation

$$\frac{\partial y}{\partial x_i} = \sum_j \frac{\partial y}{\partial z_j} \cdot \frac{\partial z_j}{\partial x_i}$$

 $\begin{array}{l} \mbox{Gradient+Jacobian notation} \\ \mbox{Let } \langle \cdot, \cdot \rangle \mbox{ be the dot product operation:} \end{array}$

$$\nabla_{\boldsymbol{x}} \boldsymbol{y} = \langle \mathbf{J}_{\boldsymbol{x}} \boldsymbol{z}, \nabla_{\boldsymbol{z}} \boldsymbol{y} \rangle$$

$$\nabla_{\mathbf{x}} \mathbf{y} = \begin{bmatrix} \frac{\partial \mathbf{y}}{\partial \mathbf{x}_1} \\ \frac{\partial \mathbf{y}}{\partial \mathbf{x}_2} \\ \dots \end{bmatrix} \in \mathbb{R}^m \qquad \mathbf{J}_{\mathbf{x}} \mathbf{z} = \begin{bmatrix} \frac{\partial \mathbf{z}_1}{\partial \mathbf{x}_1} & \frac{\partial \mathbf{z}_1}{\partial \mathbf{x}_2} & \dots \\ \frac{\partial \mathbf{z}_2}{\partial \mathbf{x}_1} & \frac{\partial \mathbf{z}_2}{\partial \mathbf{x}_2} & \dots \\ \dots & \dots & \dots \end{bmatrix} \in \mathbb{R}^{n \times m} \qquad \nabla_{\mathbf{z}} \mathbf{y} = \begin{bmatrix} \frac{\partial \mathbf{y}}{\partial \mathbf{z}_1} \\ \frac{\partial \mathbf{y}}{\partial \mathbf{z}_2} \\ \dots \end{bmatrix} \in \mathbb{R}^n$$

Forward pass $z^{(1)} = f^{(1)}(x ; \theta^{(1)})$

Forward pass $z^{(1)} = f^{(1)}(x ; \theta^{(1)})$ \downarrow $z^{(2)} = f^{(2)}(z^{(1)}; \theta^{(2)})$

Forward pass $z^{(1)} = f^{(1)}(x ; \theta^{(1)})$ \downarrow $z^{(2)} = f^{(2)}(z^{(1)}; \theta^{(2)})$ \downarrow $z^{(3)} = f^{(3)}(z^{(2)}; \theta^{(3)})$

Forward pass $z^{(1)} = f^{(1)}(x ; \theta^{(1)})$ $\mathbf{z}^{(2)} = f^{(2)}(\mathbf{z}^{(1)}; \theta^{(2)})$ $z^{(3)} = f^{(3)}(z^{(2)}; \theta^{(3)})$ $z^{(4)} = f^{(4)}(z^{(3)}; \theta^{(4)})$

Forward pass $z^{(1)} = f^{(1)}(x ; \theta^{(1)})$ $z^{(2)} = f^{(2)}(z^{(1)}; \theta^{(2)})$ $z^{(3)} = f^{(3)}(z^{(2)}; \theta^{(3)})$ $z^{(4)} = f^{(4)}(z^{(3)}; \theta^{(4)})$ $y = f^{(5)}(z^{(4)}; \theta^{(5)})$

Forward pass $z^{(1)} = f^{(1)}(x ; \theta^{(1)})$ $z^{(2)} = f^{(2)}(z^{(1)}; \theta^{(2)})$ $z^{(3)} = f^{(3)}(z^{(2)}; \theta^{(3)})$ $z^{(4)} = f^{(4)}(z^{(3)}; \theta^{(4)})$ $y = f^{(5)}(z^{(4)}; \theta^{(5)})$

Backward pass

 $\nabla_{\theta^{(1)}} y$

 $\nabla_{\theta^{(2)}} y$

 $\nabla_{\theta^{(3)}} y$

 $\nabla_{\theta^{(4)}} y$

Forward pass $z^{(1)} = f^{(1)}(x ; \theta^{(1)})$ $z^{(2)} = f^{(2)}(z^{(1)}; \theta^{(2)})$ $z^{(3)} = f^{(3)}(z^{(2)}; \theta^{(3)})$ $z^{(4)} = f^{(4)}(z^{(3)}; \theta^{(4)})$ $v = f^{(5)}(\mathbf{z}^{(4)}; \theta^{(5)})$

Backward pass

 $\nabla_{\theta^{(1)}} y$

 $\nabla_{\theta^{(2)}} y$

 $\nabla_{\theta^{(3)}} y$

$$abla_{ heta^{(4)}} y = \langle \mathbf{J}_{ heta^{(4)}} \mathbf{z^{(4)}},
abla_{\mathbf{z^{(4)}}} y
angle$$

 $\nabla_{\theta^{(5)}} y$

Forward pass $z^{(1)} = f^{(1)}(x ; \theta^{(1)})$ $z^{(2)} = f^{(2)}(z^{(1)}; \theta^{(2)})$ $z^{(3)} = f^{(3)}(z^{(2)}; \theta^{(3)})$ $z^{(4)} = f^{(4)}(z^{(3)}; \theta^{(4)})$ $v = f^{(5)}(z^{(4)}; \theta^{(5)})$

 $\nabla_{z^{(4)}} y$

Backward pass

↑

 $\nabla_{\theta^{(1)}} y$

 $\nabla_{\theta^{(2)}} y$

 $\nabla_{\theta^{(3)}} y$

$$abla_{ heta^{(4)}} y = \langle \mathsf{J}_{ heta^{(4)}}, \nabla_{\mathbf{z}^{(4)}} y
angle$$
 $abla_{ heta^{(5)}} y$

Forward pass $z^{(1)} = f^{(1)}(x ; \theta^{(1)})$ $z^{(2)} = f^{(2)}(z^{(1)}; \theta^{(2)})$ $z^{(3)} = f^{(3)}(z^{(2)}; \theta^{(3)})$ $z^{(4)} = f^{(4)}(z^{(3)}; \theta^{(4)})$ $y = f^{(5)}(z^{(4)}; \theta^{(5)})$

 $\nabla_{z^{(4)}} y$

Backward pass

ᠰ

 $\nabla_{\theta^{(1)}} y$

 $\nabla_{\theta^{(2)}} y$

$$abla_{ heta^{(3)}} y = \langle \mathbf{J}_{ heta^{(3)}} \mathbf{z}^{(3)},
abla_{\mathbf{z}^{(3)}} y
angle$$

$$abla_{ heta^{(4)}} y = \langle \mathsf{J}_{ heta^{(4)}}, \nabla_{\mathsf{z}^{(4)}} y \rangle$$
 $abla_{ heta^{(5)}} y$

Forward pass **Backward pass** $z^{(1)} = f^{(1)}(x ; \theta^{(1)})$ $\nabla_{\theta^{(1)}} y$ $z^{(2)} = f^{(2)}(z^{(1)}; \theta^{(2)})$ $\nabla_{\theta^{(2)}} y$ $z^{(3)} = f^{(3)}(z^{(2)}; \theta^{(3)})$ $\nabla_{\theta^{(3)}} y = \langle \mathsf{J}_{\theta^{(3)}} \mathsf{z}^{(3)}, \nabla_{\mathsf{z}^{(3)}} \mathsf{y} \rangle$ $z^{(4)} = f^{(4)}(z^{(3)}; \theta^{(4)})$ $\nabla_{\mathbf{z}^{(3)}} y = \langle \mathbf{J}_{\mathbf{z}^{(3)}} \mathbf{z}^{(4)}, \nabla_{\mathbf{z}^{(4)}} y \rangle$ $\nabla_{\rho(4)} y = \langle \mathbf{J}_{\rho(4)} \mathbf{z}^{(4)}, \nabla_{\mathbf{z}^{(4)}} y \rangle$ 个 $v = f^{(5)}(\mathbf{z}^{(4)}; \theta^{(5)})$ $\nabla_{\mathbf{z}^{(4)}} y$ $\nabla_{\theta^{(5)}} y$

Forward pass **Backward** pass $z^{(1)} = f^{(1)}(x ; \theta^{(1)})$ $z^{(2)} = f^{(2)}(z^{(1)}; \theta^{(2)})$ $\nabla_{\mathbf{z}^{(1)}} \mathbf{y} = \langle \mathbf{J}_{\mathbf{z}^{(1)}} \mathbf{z}^{(2)}, \nabla_{\mathbf{z}^{(2)}} \mathbf{y} \rangle$ $z^{(3)} = f^{(3)}(z^{(2)}; \theta^{(3)})$ $\nabla_{\mathbf{z}^{(2)}} \mathbf{y} = \langle \mathbf{J}_{\mathbf{z}^{(2)}} \mathbf{z}^{(3)}, \nabla_{\mathbf{z}^{(3)}} \mathbf{y} \rangle$ $z^{(4)} = f^{(4)}(z^{(3)}; \theta^{(4)})$ $\nabla_{\mathbf{z}^{(3)}} y = \langle \mathbf{J}_{\mathbf{z}^{(3)}} \mathbf{z}^{(4)}, \nabla_{\mathbf{z}^{(4)}} \mathbf{v} \rangle$ $v = f^{(5)}(\mathbf{z}^{(4)}; \theta^{(5)})$ $\nabla_{z^{(4)}} y$

 $\nabla_{\theta^{(1)}} y = \langle \mathbf{J}_{\theta^{(1)}} \mathbf{z}^{(1)}, \nabla_{\mathbf{z}^{(1)}} y \rangle$

 $\nabla_{\theta^{(2)}} y = \langle \mathbf{J}_{\theta^{(2)}} \mathbf{z}^{(2)}, \nabla_{\mathbf{z}^{(2)}} y \rangle$

 $\nabla_{\theta^{(3)}} y = \langle \mathbf{J}_{\theta^{(3)}} \mathbf{z}^{(3)}, \nabla_{\mathbf{z}^{(3)}} y \rangle$

 $\nabla_{\rho(4)} y = \langle \mathbf{J}_{\rho(4)} \mathbf{z}^{(4)}, \nabla_{\mathbf{z}^{(4)}} y \rangle$

 $\nabla_{\theta^{(5)}} y$

x






















Computation Graph (CG) implementation

CG construction / Eager forward pass

The computation graph is built in **topological order** (~order execution of operations):

- ► $x, z^{(1)}, z^{(2)}, ..., \mathcal{L}$: Expression nodes
- ▶ $W^{(1)}, b^{(1)}, ...$: Parameter nodes

Expression node

- Values
- Gradient
- Backward operation
- Backpointer(s) to antecedents

The backward operation and $\mathsf{backpointer}(s)$ are null for input operations



- Persistent values
- Gradient

Non-linear activation function:

$$\mathbf{z}' = \operatorname{relu}(\mathbf{z})$$

Non-linear activation function:

 $\mathbf{z}' = \mathsf{relu}(\mathbf{z})$

function RELU(*z*) ▷ Create node z' = ExpressionNode()▷ Compute forward value $z'.value = \begin{bmatrix} \max(0, z_1) \\ \max(0, z_2) \\ \dots \end{bmatrix}$ ▷ Set backward operation $\mathbf{z}' \cdot \mathbf{d} = \mathbf{d}$ relu ▷ Set backpointers z'.backptrs = [z]

return z'

Non-linear activation function:

 $\mathbf{z}' = \operatorname{relu}(\mathbf{z})$

function RELU(*z*) ▷ Create node z' = ExpressionNode()▷ Compute forward value $\mathbf{z}'.\mathbf{value} = \begin{bmatrix} \max(0, \mathbf{z}_1) \\ \max(0, \mathbf{z}_2) \\ \dots \end{bmatrix}$ ▷ Set backward operation $\mathbf{z}' \cdot \mathbf{d} = \mathbf{d}$ relu ▷ Set backpointers z'.backptrs = [z]

return z'

Projection operation $\boldsymbol{z} = \boldsymbol{W}\boldsymbol{x} + \boldsymbol{b}$:

$$z = Linear(x, W, b)$$

Non-linear activation function:

 $\mathbf{z}' = \operatorname{relu}(\mathbf{z})$

function RELU(z) ▷ Create node z' = ExpressionNode()▷ Compute forward value $\mathbf{z}'.value = \begin{bmatrix} \max(0, \mathbf{z}_1) \\ \max(0, \mathbf{z}_2) \end{bmatrix}$ ▷ Set backward operation $\mathbf{z}' \cdot \mathbf{d} = \mathbf{d}$ relu ▷ Set backpointers \mathbf{z}' .backptrs = $[\mathbf{z}]$

Projection operation $\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b}$:

 $\textbf{\textit{z}} = \mathsf{Linear}(\textbf{\textit{x}}, \textbf{\textit{W}}, \textbf{\textit{b}})$

function LINEAR(x, W, b) ▷ Create node z = ExpressionNode()▷ Compute forward value $z_{value} = Wx + b$ ▷ Set backward operation z.d = d linear ▷ Set backpointers z.backptrs = [W, b]

return z

return z'

Backward pass

Execution of the backward pass

Nodes are visited in **reverse topological order** (reverse order of creation):

- ▶ The gradient of the loss (last created node) is set to 1
- ► For each node, we call it's derivative function
- The derivative functions will backpropagate gradient to antecedents Gradient must be accumulated (expressions can be used several times)

Backward pass

Execution of the backward pass

Nodes are visited in reverse topological order (reverse order of creation):

- ▶ The gradient of the loss (last created node) is set to 1
- For each node, we call it's derivative function
- The derivative functions will backpropagate gradient to antecedents Gradient must be accumulated (expressions can be used several times)

```
function BACKWARD(nodes, \mathcal{L})
\mathcal{L}.grad = 1
```

```
for n \in reversed(nodes) do
```

> Call the derivative functions n.d(n.backptrs)

$$\operatorname{relu}(x) = \begin{cases} 0 & \text{if } \leq 0 \\ x & \text{otherwise} \end{cases} \qquad \frac{\partial}{\partial x} \operatorname{relu}(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x > 0 \\ \text{undefined otherwise, } = 0 \end{cases}$$

$$\operatorname{relu}(x) = \begin{cases} 0 \text{ if } \leq 0 \\ x \text{ otherwise} \end{cases} \qquad \frac{\partial}{\partial x} \operatorname{relu}(x) = \begin{cases} 0 \text{ if } x < 0 \\ 1 \text{ if } x > 0 \\ \text{undefined otherwise, } = 0 \end{cases}$$

$$\nabla_{\boldsymbol{z}} \mathcal{L} = \langle \mathbf{J}_{\boldsymbol{z}} \boldsymbol{z}', \nabla_{\boldsymbol{z}'} \mathcal{L} \rangle$$

$$\mathbf{J}_{\mathbf{z}}\mathbf{z}' = \begin{bmatrix} \frac{\partial z_1'}{\partial z_1} & \frac{\partial z_1'}{\partial z_2} & \dots \\ \frac{\partial z_2'}{\partial z_1} & \frac{\partial z_2'}{\partial z_2} & \dots \\ \dots & \dots & \dots \end{bmatrix}$$

$$\operatorname{relu}(x) = \begin{cases} 0 \text{ if } \leq 0 \\ x \text{ otherwise} \end{cases} \qquad \frac{\partial}{\partial x} \operatorname{relu}(x) = \begin{cases} 0 \text{ if } x < 0 \\ 1 \text{ if } x > 0 \\ \text{undefined otherwise, } = 0 \end{cases}$$

 $abla_{\mathbf{z}}\mathcal{L} = \langle \mathbf{J}_{\mathbf{z}}\mathbf{z}',
abla_{\mathbf{z}'}\mathcal{L}
angle$

$$\mathbf{J}_{\mathbf{z}}\mathbf{z}' = \begin{bmatrix} \frac{\partial z'_1}{\partial z_1} & \frac{\partial z'_1}{\partial z_2} & \dots \\ \frac{\partial z'_2}{\partial z_1} & \frac{\partial z'_2}{\partial z_2} & \dots \\ \dots & \dots & \dots \end{bmatrix}$$

$$\frac{\partial z'_i}{\partial z_j} = \begin{cases} 0, \text{ if } i \neq j & \text{(piecewise function!)} \\ \frac{\partial}{\partial z_i} f(z_i), \text{ if } i = j \end{cases}$$

$$\operatorname{relu}(x) = \begin{cases} 0 \text{ if } \leq 0 \\ x \text{ otherwise} \end{cases} \qquad \frac{\partial}{\partial x} \operatorname{relu}(x) = \begin{cases} 0 \text{ if } x < 0 \\ 1 \text{ if } x > 0 \\ \text{undefined otherwise, } = 0 \end{cases}$$

 $\Gamma \partial C \neg$

$$\nabla_{\mathbf{z}} \mathcal{L} = \langle \mathbf{J}_{\mathbf{z}} \mathbf{z}', \nabla_{\mathbf{z}'} \mathcal{L} \rangle \qquad \qquad \nabla_{\mathbf{z}} \mathcal{L} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{z}_1} \\ \frac{\partial \mathcal{L}}{\partial \mathbf{z}_2} \\ \dots \end{bmatrix}$$

$$\mathbf{J}_{\mathbf{z}}\mathbf{z}' = \begin{bmatrix} \frac{\partial z'_1}{\partial z_1} & \frac{\partial z'_1}{\partial z_2} & \dots \\ \frac{\partial z'_2}{\partial z_1} & \frac{\partial z'_2}{\partial z_2} & \dots \\ \dots & \dots & \dots \end{bmatrix}$$

$$\operatorname{relu}(x) = \begin{cases} 0 \text{ if } \leq 0 \\ x \text{ otherwise} \end{cases} \qquad \frac{\partial}{\partial x} \operatorname{relu}(x) = \begin{cases} 0 \text{ if } x < 0 \\ 1 \text{ if } x > 0 \\ \text{undefined otherwise, } = 0 \end{cases}$$

$$\operatorname{relu}(x) = \begin{cases} 0 & \text{if } \leq 0 \\ x & \text{otherwise} \end{cases} \qquad \frac{\partial}{\partial x} \operatorname{relu}(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x > 0 \\ \text{undefined otherwise, } = 0 \end{cases}$$

$$abla_{\mathbf{z}}\mathcal{L} = \langle \mathbf{J}_{\mathbf{z}}\mathbf{z}',
abla_{\mathbf{z}'}\mathcal{L}
angle$$

$$\nabla_{\mathbf{z}} \mathcal{L} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial z_1} \\ \frac{\partial \mathcal{L}}{\partial z_2} \\ \cdots \end{bmatrix}$$

$$\mathbf{J}_{\mathbf{z}}\mathbf{z}' = \begin{bmatrix} \frac{\partial z_1'}{\partial z_1} & \frac{\partial z_1'}{\partial z_2} & \dots \\ \frac{\partial z_2'}{\partial z_1} & \frac{\partial z_2'}{\partial z_2} & \dots \\ \dots & \dots & \dots \end{bmatrix}$$

$$\frac{\partial \mathcal{L}}{\partial z_i} = \sum_j \frac{\partial \mathcal{L}}{\partial z'_j} \cdot \frac{\partial z'_j}{\partial z_i}$$

- - - -

$$\operatorname{relu}(x) = \begin{cases} 0 & \text{if } \leq 0 \\ x & \text{otherwise} \end{cases} \qquad \frac{\partial}{\partial x} \operatorname{relu}(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x > 0 \\ \text{undefined otherwise, } = 0 \end{cases}$$

$$abla_{\mathbf{z}}\mathcal{L} = \langle \mathbf{J}_{\mathbf{z}}\mathbf{z}',
abla_{\mathbf{z}'}\mathcal{L}
angle$$

$$\mathbf{J}_{\mathbf{z}}\mathbf{z}' = \begin{bmatrix} \frac{\partial z_1'}{\partial z_1} & \frac{\partial z_1'}{\partial z_2} & \dots \\ \frac{\partial z_2'}{\partial z_1} & \frac{\partial z_2'}{\partial z_2} & \dots \\ \dots & \dots & \dots \end{bmatrix}$$

$$\nabla_{\boldsymbol{z}} \mathcal{L} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \boldsymbol{z}_1} \\ \frac{\partial \mathcal{L}}{\partial \boldsymbol{z}_2} \\ \cdots \end{bmatrix}$$

- - - -

$$\begin{split} \frac{\partial \mathcal{L}}{\partial z_i} &= \sum_j \frac{\partial \mathcal{L}}{\partial z'_j} \cdot \frac{\partial z'_j}{\partial z_i} \\ &= \frac{\partial \mathcal{L}}{\partial z'_i} \cdot \frac{\partial z'_i}{\partial z_i} \quad \text{(piecewise function!)} \end{split}$$

$$\operatorname{relu}(x) = \begin{cases} 0 & \text{if } \leq 0 \\ x & \text{otherwise} \end{cases} \qquad \frac{\partial}{\partial x} \operatorname{relu}(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x > 0 \\ \text{undefined otherwise, } = 0 \end{cases}$$

$$abla_{\mathbf{z}}\mathcal{L} = \langle \mathbf{J}_{\mathbf{z}}\mathbf{z}',
abla_{\mathbf{z}'}\mathcal{L}
angle$$

$$\mathbf{J}_{\mathbf{z}}\mathbf{z}' = \begin{bmatrix} \frac{\partial z_1'}{\partial z_1} & \frac{\partial z_1'}{\partial z_2} & \dots \\ \frac{\partial z_2'}{\partial z_1} & \frac{\partial z_2'}{\partial z_2} & \dots \\ \dots & \dots & \dots \end{bmatrix}$$

$$\frac{\partial z'_i}{\partial z_j} = \begin{cases} 0, \text{ if } i \neq j & \text{(piecewise function!)} \\ \frac{\partial}{\partial z_i} f(z_i), \text{ if } i = j \end{cases}$$

$$\nabla_{\boldsymbol{z}} \mathcal{L} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \boldsymbol{z}_1} \\ \frac{\partial \mathcal{L}}{\partial \boldsymbol{z}_2} \\ \dots \end{bmatrix}$$

$$\begin{split} \frac{\partial \mathcal{L}}{\partial z_{i}} &= \sum_{j} \frac{\partial \mathcal{L}}{\partial z'_{j}} \cdot \frac{\partial z'_{j}}{\partial z_{i}} \\ &= \frac{\partial \mathcal{L}}{\partial z'_{i}} \cdot \frac{\partial z'_{i}}{\partial z_{i}} \quad \text{(piecewise function!)} \\ &= \frac{\partial \mathcal{L}}{\partial z'_{i}} \cdot \mathbb{1}[z_{i} > 0] \end{split}$$

$$\operatorname{relu}(x) = egin{cases} 0 & \operatorname{if} & \leq 0 \\ x & \operatorname{otherwise} \end{cases}$$
 $\operatorname{relu}'(x) =$

$$f(x) = \begin{cases} 0 \text{ if } x < 0\\ 1 \text{ if } x > 0\\ \text{undefined otherwise} \end{cases}$$

1

function RELU(z)

$$z' = \text{ExpressionNode}()$$

 $z'.value = \begin{bmatrix} \max(0, z_1) \\ \max(0, z_2) \\ \dots \end{bmatrix}$
 $z'.d = d_relu$
 $z'.backptrs = [z]$

return z'

$$rac{\partial \mathcal{L}}{\partial z_i} = rac{\partial \mathcal{L}}{\partial z'_i} \cdot \mathbb{1}[z_i > 0]$$

function D_RELU(z', [z]) for $i \in \{1...n\}$ do \triangleright If the value is positive, \triangleright we copy the gradient if $z_i > 0$ then \mathbf{z} .grad_i = \mathbf{z} .grad_i + \mathbf{z}' .grad_i

$$oldsymbol{z} = oldsymbol{W} oldsymbol{x} + oldsymbol{b}$$
 \Leftrightarrow $oldsymbol{z}_j = \sum_k W_{j,k} x_k + b_j$

$$oldsymbol{z} = oldsymbol{W} oldsymbol{x} + oldsymbol{b} \qquad \Leftrightarrow \qquad oldsymbol{z}_j = \sum_k W_{j,k} x_k + b_j$$

 $\frac{\partial \mathcal{L}}{\partial b_i} = \sum_j \frac{\partial \mathcal{L}}{\partial z_j} \cdot \frac{\partial z_j}{\partial b_i}$

$$\boldsymbol{z} = \boldsymbol{W}\boldsymbol{x} + \boldsymbol{b} \qquad \Leftrightarrow \qquad z_j = \sum_k W_{j,k} x_k + b_j$$
$$\frac{\partial \mathcal{L}}{\partial b_i} = \sum_j \frac{\partial \mathcal{L}}{\partial z_j} \cdot \frac{\partial z_j}{\partial b_i}$$
$$\frac{\partial z_j}{\partial b_i} = \frac{\partial}{\partial b_j} \sum_k W_{j,k} x_k + b_j$$

$$\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b} \qquad \Leftrightarrow \qquad z_j = \sum_k W_{j,k} x_k + b_j$$
$$\frac{\partial \mathcal{L}}{\partial b_i} = \sum_j \frac{\partial \mathcal{L}}{\partial z_j} \cdot \frac{\partial z_j}{\partial b_i}$$
$$\frac{\partial z_j}{\partial b_i} = \frac{\partial}{\partial b_i} \sum_k W_{j,k} x_k + b_j = \begin{cases} 1, \text{ if } i = j\\ 0, \text{ otherwise} \end{cases}$$

$$\begin{aligned} \boldsymbol{z} &= \boldsymbol{W}\boldsymbol{x} + \boldsymbol{b} \quad \Leftrightarrow \quad z_j = \sum_k W_{j,k} x_k + b_j \\ \frac{\partial \mathcal{L}}{\partial b_i} &= \sum_j \frac{\partial \mathcal{L}}{\partial z_j} \cdot \frac{\partial z_j}{\partial b_i} \\ &= \sum_j \frac{\partial \mathcal{L}}{\partial z_j} \cdot \mathbb{1}[j = i] \\ &= \frac{\partial \mathcal{L}}{\partial z_i} \quad (\text{copy incoming gradient!}) \end{aligned}$$

$$\frac{\partial z_j}{\partial b_i} = \frac{\partial}{\partial b_i} \sum_k W_{j,k} x_k + b_j = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{otherwise} \end{cases}$$

$$z = Wx + b \quad \Leftrightarrow \quad z_j = \sum_k W_{j,k} x_k + b_j$$

$$\frac{\partial \mathcal{L}}{\partial b_i} = \sum_j \frac{\partial \mathcal{L}}{\partial z_j} \cdot \frac{\partial z_j}{\partial b_i}$$

$$= \sum_j \frac{\partial \mathcal{L}}{\partial z_j} \cdot \mathbb{1}[j = i]$$

$$= \frac{\partial \mathcal{L}}{\partial z_i} \quad (\text{copy incoming gradient!})$$

$$\frac{\partial z_j}{\partial b_i} = \frac{\partial}{\partial b_i} \sum_k W_{j,k} x_k + b_j = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{otherwise} \end{cases}$$

$$z = Wx + b \quad \Leftrightarrow \quad z_j = \sum_k W_{j,k} x_k + b_j$$

$$\frac{\partial \mathcal{L}}{\partial b_i} = \sum_j \frac{\partial \mathcal{L}}{\partial z_j} \cdot \frac{\partial z_j}{\partial b_i}$$

$$= \sum_j \frac{\partial \mathcal{L}}{\partial z_j} \cdot \mathbb{1}[j = i]$$

$$= \frac{\partial \mathcal{L}}{\partial z_i} \quad (\text{copy incoming gradient!})$$

$$\frac{\partial z_j}{\partial b_i} = \frac{\partial}{\partial b_i} \sum_k W_{j,k} x_k + b_j = \begin{cases} 1, \text{ if } i = j \\ 0, \text{ otherwise} \end{cases}$$

$$\frac{\partial z_j}{\partial W_{i,l}} = \frac{\partial}{\partial W_{i,l}} \sum_k W_{j,k} x_k + b_j$$

31 / 64

$$\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b} \quad \Leftrightarrow \quad z_j = \sum_k W_{j,k} x_k + b_j$$

$$\frac{\partial \mathcal{L}}{\partial b_i} = \sum_j \frac{\partial \mathcal{L}}{\partial z_j} \cdot \frac{\partial z_j}{\partial b_i}$$

$$= \sum_j \frac{\partial \mathcal{L}}{\partial z_j} \cdot \mathbb{1}[j = i]$$

$$= \frac{\partial \mathcal{L}}{\partial z_i} \quad (\text{copy incoming gradient!})$$

$$\frac{\partial z_j}{\partial b_i} = \frac{\partial}{\partial b_i} \sum_k W_{j,k} x_k + b_j = \begin{cases} 1, \text{ if } i = j \\ 0, \text{ otherwise} \end{cases}$$

$$\mathbf{z}_j = \sum_k W_{j,k} x_k + b_j$$

$$\frac{\partial \mathcal{L}}{\partial W_{i,l}} = \sum_j \frac{\partial \mathcal{L}}{\partial z_j} \cdot \frac{\partial z_j}{W_{i,l}}$$

$$= \sum_j \frac{\partial \mathcal{L}}{\partial z_j} \cdot x_l \cdot \mathbb{1}[i = j]$$

$$\frac{\partial z_j}{\partial W_{i,l}} = \frac{\partial}{\partial W_{i,l}} \sum_k W_{j,k} x_k + b_j$$

31 / 64

j]

$$z = Wx + b \quad \Leftrightarrow \quad z_j = \sum_k W_{j,k} x_k + b_j$$

$$i]$$

$$i]$$

$$i] \quad \frac{\partial \mathcal{L}}{\partial W_{i,l}} = \sum_j \frac{\partial \mathcal{L}}{\partial z_j} \cdot \frac{\partial z_j}{W_{i,l}}$$

$$= \sum_j \frac{\partial \mathcal{L}}{\partial z_j} \cdot x_l \cdot \mathbb{1}[i = j]$$

$$\nabla_W \mathcal{L} = (\nabla_z \mathcal{L})(\mathbf{x}^{\top}) \quad (\text{outer product})$$

$$\frac{\partial z_j}{\partial W_{i,l}} = \frac{\partial}{\partial W_{i,l}} \sum_k W_{j,k} x_k + b_j$$

$$= \begin{cases} x_l, \text{ if } i = j \\ 0, \text{ otherwise} \end{cases}$$

$$\begin{split} \frac{\partial \mathcal{L}}{\partial b_i} &= \sum_{j} \frac{\partial \mathcal{L}}{\partial z_j} \cdot \frac{\partial z_j}{\partial b_i} \\ &= \sum_{j} \frac{\partial \mathcal{L}}{\partial z_j} \cdot \mathbb{1}[j=i] \\ &= \frac{\partial \mathcal{L}}{\partial z_i} \quad (\text{copy incoming gradie}) \end{split}$$

$$\frac{\partial z_j}{\partial b_i} = \frac{\partial}{\partial b_i} \sum_k W_{j,k} x_k + b_j = \begin{cases} 1, \text{ if } i = j \\ 0, \text{ otherwise} \end{cases}$$

function LINEAR(x, W, b)
z = ExpressionNode()
z.value = Wx + b
z.d = d_linear
z.backptrs = [W, b]

$$\frac{\partial \mathcal{L}}{\partial b_i} = \frac{\partial \mathcal{L}}{\partial z_i} \qquad \nabla_{\boldsymbol{W}} \mathcal{L} = (\nabla_{\boldsymbol{z}} \mathcal{L})(\boldsymbol{x}^{\top})$$

function D_LINEAR(z, [x, W, b]) b.grad = b.grad + z.grad $W.grad = W.grad + z.grad @ x^{\top}$

return z

function LINEAR(x, W, b)
z = ExpressionNode()
z.value = Wx + b
z.d = d_linear
z.backptrs = [W, b]

return z

$$\frac{\partial \mathcal{L}}{\partial b_i} = \frac{\partial \mathcal{L}}{\partial z_i} \qquad \nabla_{\mathbf{W}} \mathcal{L} = (\nabla_{\mathbf{z}} \mathcal{L})(\mathbf{x}^\top)$$

function D_LINEAR(z, [x, W, b]) b.grad = b.grad + z.grad $W.grad = W.grad + z.grad @ x^T$ x.grad = x.grad + ...

Missing gradient?

Why don't we backpropagate to x?!

We do not need it for today's lab exercises, you will see how to do that next week.

Summary

Computation graph

- Forward pass: compute values
- Backward pass: compute gradient for each parameter
- ► Gradient initialization: you should be careful with that because gradient is accumulated

First lab exercises

- Simple linear model: don't build a computation graph, explicitly apply forward and backward operations
- d_Linear: return a tuple with gradient of W and b instead of writing into a node
- ► Do not need to worry about gradient initialization / accumulation :-)

Pytorch

In Pytorch, expression nodes used to be of type Variable. Nowadays, autodiff is directly implemented in the Tensor class. Vanishing gradient, activation functions and initialization

Experimental observations

The MNIST database 82944649709295159103 23591762822507497832 11836103100112730465 26471899307102035465

Comparison of different depth for feed-forward architecture



Hidden layers have a sigmoid activation function.

The output layer is softmax.
Experimental observations: http://neuralnetworksanddeeplearning.com/chap5.html

- Without hidden layer: $\approx 88\%$ accuracy
- ▶ 1 hidden layer (30): \approx 96.5% accuracy
- ▶ 2 hidden layer (30): \approx 96.9% accuracy
- ▶ 3 hidden layer (30): \approx 96.5% accuracy
- ▶ 4 hidden layer (30): \approx 96.5% accuracy

Experimental observations: http://neuralnetworksanddeeplearning.com/chap5.html

- Without hidden layer: $\approx 88\%$ accuracy
- ▶ 1 hidden layer (30): \approx 96.5% accuracy
- ▶ 2 hidden layer (30): \approx 96.9% accuracy
- ▶ 3 hidden layer (30): \approx 96.5% accuracy
- ▶ 4 hidden layer (30): \approx 96.5% accuracy



Intuitive explanation 1/2

Let consider the simplest deep neural network, with just a single neuron in each layer.



 w_i, b_i are resp. the weight and bias of neuron *i* and *C* some loss function.

Intuitive explanation 1/2

Let consider the simplest deep neural network, with just a single neuron in each layer.



 w_i, b_i are resp. the weight and bias of neuron i and C some loss function. Compute the gradient of C w.r.t the bias b_1

$$\frac{\partial C}{\partial b_1} = \frac{\partial C}{\partial y_4} \times \frac{\partial y_4}{\partial a_4} \times \frac{\partial a_4}{\partial y_3} \times \frac{\partial y_3}{\partial a_3} \times \frac{\partial a_3}{\partial y_2} \times \frac{\partial y_2}{\partial a_2} \times \frac{\partial a_2}{\partial y_1} \times \frac{\partial y_1}{\partial a_1} \times \frac{\partial a_1}{\partial b_1}$$
(1)
$$= \frac{\partial C}{\partial y_4} \times \sigma'(a_4) \times w_4 \times \sigma'(a_3) \times w_3 \times \sigma'(a_2) \times w_2 \times \sigma'(a_1)$$
(2)

Intuitive explanation 2/2



$$\sigma(x) = rac{1}{1 + \exp(-x)}$$

 $\sigma'(x) = \sigma(x)(1 - \sigma(x))$

Vanishing gradient

▶ if the last layer are well trained (and outputs "strong values" close to 0 or 1),

• early layers receive a really small incoming gradient.

In the "best case", we successive multiplications by 0.25!

Other activation functions



Hyperbolic tangent

$$tanh(x) = \frac{1 - \exp(-2x)}{1 + \exp(-2x)} \quad tanh'(x) = 1 - tanh(x)^2$$

Better gradient than sigmoid around 0

Popular in Natural Language Processing

Other activation functions



Hyperbolic tangent

$$anh(x) = rac{1-\exp(-2x)}{1+\exp(-2x)}$$
 $anh'(x) = 1- anh(x)^2$

- Better gradient than sigmoid around 0
- Popular in Natural Language Processing

Rectified Linear Unit

$$\operatorname{relu}(x) = \begin{cases} 0 \text{ if } \leq 0 \\ x \text{ otherwise} \end{cases} \quad \operatorname{relu}'(x) = \begin{cases} 0 \text{ if } x < 0 \\ 1 \text{ if } x > 0 \\ \text{undefined otherwise} \end{cases}$$

- No vanishing gradient issue
- "Dead units" problem (i.e. $b_i \ll 0$)
- Popular in Computer Vision (very deep networks)

Parameters initialization

What do we want?

- Values close to 0 prevent gradient vanishing (or gradient exploding/disappearing in the case of relu)
- Gradient magnitude approximately similar for all layers (to prevent that a subset of layers do all the works while others are useless)

Parameters initialization

What do we want?

- Values close to 0 prevent gradient vanishing (or gradient exploding/disappearing in the case of relu)
- Gradient magnitude approximately similar for all layers (to prevent that a subset of layers do all the works while others are useless)

Hyperbolic tangent

Let
$$\boldsymbol{W} \in \mathbb{R}^{m \times n}$$
 and $\boldsymbol{b} \in \mathbb{R}^{m}$:
 $\boldsymbol{W} \sim \mathcal{U} \left[-\frac{\sqrt{6}}{\sqrt{m+n}}, +\frac{\sqrt{6}}{\sqrt{m+n}} \right]$
 $\boldsymbol{b} = 0$

Usually called Xavier or Glorot initialization [Glorot and Bengio, 2010]

Rectified Linear Unit

Let
$$\boldsymbol{W} \in \mathbb{R}^{m imes n}$$
 and $\boldsymbol{b} \in \mathbb{R}^m$:

•
$$\boldsymbol{W} \sim \mathcal{U}\left[-\frac{\sqrt{6}}{\sqrt{n}}, +\frac{\sqrt{6}}{\sqrt{n}}\right]$$

 \blacktriangleright **b** = 0

(or $\boldsymbol{b}=0.01$ to prevent dying units)

Usually called Kaiming or He initialization [He et al., 2015]

Regularization

Generalization

Overparameterized neural networks

Networks where the number of parameters exceed the training dataset size.

Can learn by heart the dataset,

i.e. $\textbf{overfit}\ \textbf{the}\ \textbf{data} \rightarrow \textbf{does}\ \textbf{not}\ \textbf{generalize}\ \textbf{well}\ \textbf{to}\ \textbf{unseen}\ \textbf{data}$

Are easier to optimize in practice

Generalization

Overparameterized neural networks

Networks where the number of parameters exceed the training dataset size.

- Can learn by heart the dataset,
 - i.e. $overfit\ the\ data$ \rightarrow does not generalize well to unseen data
- Are easier to optimize in practice

Monitoring the training process

- ▶ Loss should go down \Rightarrow otherwise your step-size is probably too big!
- Training accuracy should go up
- Dev accuracy should go up \Rightarrow otherwise the network is overfitting!

Generalization

Overparameterized neural networks

Networks where the number of parameters exceed the training dataset size.

- Can learn by heart the dataset,
 - i.e. overfit the data \rightarrow does not generalize well to unseen data
- Are easier to optimize in practice

Monitoring the training process

- Loss should go down \Rightarrow otherwise your step-size is probably too big!
- Training accuracy should go up
 - Dev accuracy should go up \Rightarrow otherwise the network is overfitting!

Regularization

Techniques to control parameters during learning and prevent overfitting

model	# params	random crop	weight decay	train accuracy	test accuracy
Inception	1,649,402	yes	yes	100.0	89.05
		yes	no	100.0	89.31
		no	yes	100.0	86.03
		no	no	100.0	85.75
(fitting random labels)		no	no	100.0	9.78
Inception w/o BatchNorm	1,649,402	no	yes	100.0	83.00
		no	no	100.0	82.00
(fitting random labels)		no	no	100.0	10.12
Alexnet	1,387,786	yes	yes	99.90	81.22
		yes	no	99.82	79.66
		no	yes	100.0	77.36
		no	no	100.0	76.07
(fitting random labels)		no	no	99.82	9.86
MLP 3x512	1,735,178	no	yes	100.0	53.35
		no	no	100.0	52.39
(fitting random labels)		no	no	100.0	10.48
MLP 1x512	1,209,866	no	yes	99.80	50.39
		no	no	100.0	50.51
(fitting random labels)		no	no	99.34	10.61

Learning with random inputs and labels 1/2 [Zhang et al., 2017]

Learning with random inputs and labels 2/2 [Zhang et al., 2017]



Regularization L2 or Gaussian prior or weight decay 1/3

$$\hat{\theta} = \arg\min_{\theta} \mathcal{L}(f(x;\theta), y) + \frac{\lambda}{2} ||\theta||^2$$
$$= \arg\min_{\theta} \mathcal{L}(f(x;\theta), y) + \mathcal{R}(\theta; \lambda)$$

Regularization term

The second term $\mathcal{R}(\theta; \lambda)$ is a L2 regularization term which can be equivalently interpreted as:

- ▶ a soft constraint on the magnitude of parameters
- ▶ a Gaussian prior on parameters: $\mathcal{N}(0, 1/\lambda)$
- re-scaling the parameters after each update (weight decay)

Regularization L2 or Gaussian prior or weight decay 2/3

$$\hat{\theta} = \arg\min_{\theta} \mathcal{L}(f(x;\theta), y) + \frac{\lambda}{2} ||\theta||^2$$

=
$$\arg\min_{\theta} \mathcal{L}(f(x;\theta), y) + \mathcal{R}(\theta; \lambda)$$

Gradient update

$$oldsymbol{ heta} = oldsymbol{ heta} - \epsilon
abla_{ heta} \mathcal{L} - \epsilon
abla_{ heta} \mathcal{R} \ = oldsymbol{ heta} - \epsilon (
abla_{ heta} \mathcal{L} -
abla_{ heta} \mathcal{R})$$

Regularization L2 or Gaussian prior or weight decay 2/3

$$\hat{\theta} = \arg\min_{\theta} \mathcal{L}(f(x;\theta), y) + \frac{\lambda}{2} ||\theta||^2$$
$$= \arg\min_{\theta} \mathcal{L}(f(x;\theta), y) + \mathcal{R}(\theta; \lambda)$$

Gradient update

$$egin{aligned} m{ heta} &= m{ heta} - \epsilon
abla_{ heta} \mathcal{L} - \epsilon
abla_{ heta} \mathcal{R} \ &= m{ heta} - \epsilon (
abla_{ heta} \mathcal{L} -
abla_{ heta} \mathcal{R}) \end{aligned}$$

What does the gradient of the regularizer look like? Let *b* be a a parameter of the network:

$$\frac{\partial}{\partial b}\mathcal{R} = \frac{\partial}{\partial b}\frac{\lambda}{2}||\theta||^2 = \frac{\lambda}{2}2b = \lambda b$$

Regularization L2 or Gaussian prior or weight decay 3/3

```
Implementation from Pytorch (slightly modified):
class SGD(Optimizer):
  def step(self, closure=None):
    """Performs a single optimization step."""
    for group in self.param groups:
      for p in group['params']:
        if p.grad is None:
          continue
      d p = p.grad.data # get gradient
      weight_decay = group['weight decav']
      if weight decay != 0:
        d p.add (weight decay, p.data) # add weight decay to the gradient
```

```
p.data.add_(-group['lr'], d_p) # update parameters
```

Dropout 1/4 [Hinton et al., 2012, Srivastava et al., 2014] How does dropout work?

- During training, we randomly "turn off" neurons, i.e. we randomly set elements of hidden layers z to 0
- During test, we do use the full network



Intuition

- prevents co-adaptation between units
- equivalent to averaging different models that have different structure but share parameters

Dropout 2/4 [Hinton et al., 2012]



Dropout 2/4 [Hinton et al., 2012]





Dropout layer

A dropout layer is parameterized by the probability of "turning off" a neuron $p \in [0, 1]$:

z' = Dropout(z; p = 0.5)

Dropout layer

A dropout layer is parameterized by the probability of "turning off" a neuron $p \in [0, 1]$:

z' = Dropout(z; p = 0.5)

Implementation

- $z \in \mathbb{R}^n$: output of a hidden layer
- ▶ $p \in [0, 1]$: dropout probability
- ▶ $m \in \{0,1\}^n$: mask vector
- z': hidden values after dropout application

The mask \boldsymbol{m} is a vector of booleans stating if neurons z_i is kept $(m_i = 1)$ or "turned off" $(m_i = 0)$.

Dropout layer

A dropout layer is parameterized by the probability of "turning off" a neuron $p \in [0, 1]$:

z' = Dropout(z; p = 0.5)

Implementation

- $z \in \mathbb{R}^n$: output of a hidden layer
- ▶ $p \in [0, 1]$: dropout probability
- ▶ $m \in \{0,1\}^n$: mask vector
- z': hidden values after dropout application

Forward pass:

$$m{m} \sim {\sf Bernoulli}(1-p) \ z'_i = rac{z_i * m_i}{1-p}$$

The mask \boldsymbol{m} is a vector of booleans stating if neurons z_i is kept $(m_i = 1)$ or "turned off" $(m_i = 0)$.

Dropout layer

A dropout layer is parameterized by the probability of "turning off" a neuron $p \in [0, 1]$:

z' = Dropout(z; p = 0.5)

Implementation

- $z \in \mathbb{R}^n$: output of a hidden layer
- ▶ $p \in [0, 1]$: dropout probability
- ▶ $m \in \{0,1\}^n$: mask vector
- z': hidden values after dropout application

Forward pass: $m \sim \text{Bernoulli}(1-p)$ $z'_i = \frac{z_i * m_i}{1-p}$ $\Rightarrow \text{ no gradient for}$ "turned off" neurons

The mask \boldsymbol{m} is a vector of booleans stating if neurons z_i is kept $(m_i = 1)$ or "turned off" $(m_i = 0)$.

Dropout 4/4

Where do you apply dropout?

- On the input of the neural network \boldsymbol{x}
- After activation functions ($\sigma(0) \neq 0$)
- Do not apply dropout on the output logits

Dropout 4/4

Where do you apply dropout?

- On the input of the neural network x
- After activation functions ($\sigma(0) \neq 0$)
- Do not apply dropout on the output logits

Which dropout probability should you use?

- Empirical question: you have to test!
- Dropout probability at different layers can be different (especially input vs. hidden layers)
- Usually $0.1 \le p \le 0.5$

Dropout variants

Dropout can be applied differently for special neural network architectures (e.g. convolutions, recurrent neural networks)

Better optimizers

Stochastic Gradient Descent (SGD)

$$\theta^{(t+1)} = \theta^{(t)} - \epsilon^{(t)} \nabla_{\theta} \mathcal{L}$$

Advantages

- Simple
- Single hyper-parameter: the step-size ϵ

Stochastic Gradient Descent (SGD)

$$\theta^{(t+1)} = \theta^{(t)} - \epsilon^{(t)} \nabla_{\theta} \mathcal{L}$$

Advantages

- Simple
- Single hyper-parameter: the step-size ϵ

Downsides

- Forget information about previous updates
- Require to search for the best step-size strategy
- Require step-size annealing in practice: how? what scaling factor?
- Based on first-order information only
 - (i.e. the curvature of the optimized function is ignored)

Momentum 1/3



Momentum 1/3



Momentum 1/3



Momentum 2/3

[Polyak, 1964]

γ: velocity of parameters, i.e. cumulative information about past gradients
 μ ∈ [0,1]: momentum, i.e. how much information must be preserved?

$$\gamma^{(t+1)} = \mu \gamma^{(t)} + \nabla_{\theta} \mathcal{L}$$
$$\theta^{(t+1)} = \theta^{(t)} - \epsilon \gamma^{(t+1)}$$

Variants

- Gradient dampening, i.e. diminish the contribution of the current gradient
- Nesterov's Accelerated Gradient [Sutskever et al., 2013]

Momentum 3/3

```
Implementation from Pytorch (slightly modified):
for group in self.param_groups:
  for p in group['params']:
    if p.grad is None:
      continue
  d_p = p.grad.data # get the gradient
  if momentum != 0:
    param state = self.state[p]
    if 'momentum buffer' not in param_state: # initialize velocity vector
      buf = param state['momentum buffer'] = torch.clone(d p).detach()
    else:
      buf = param_state['momentum_buffer'] # retrieve velocity vector
      buf.mul_(momentum).add_(d_p) # update velocity vector
  d_p = buf
```

```
p.data.add_(-group['lr'], d_p) # update parameters
```
Adagrad [Duchi et al., 2011]

- ▶ Replace global step-size with dynamic per parameter step-size + global learning rate
- ► The dynamic per parameter step-size is computed w.r.t. previous gradient l2-norm ⇒ parameters with small (resp. large) gradient will have a large (resp. small) step-size

Adagrad [Duchi et al., 2011]

- ▶ Replace global step-size with dynamic per parameter step-size + global learning rate
- ► The dynamic per parameter step-size is computed w.r.t. previous gradient l2-norm ⇒ parameters with small (resp. large) gradient will have a large (resp. small) step-size

Adadelta [Zeiler, 2012]

- Dynamic per parameter rate is computed with a fixed window of past gradients
- Approximate second-order information to incorporate curvature information ⇒ less sensitive to the learning rate hyper-parameter!

Adagrad [Duchi et al., 2011]

- ▶ Replace global step-size with dynamic per parameter step-size + global learning rate
- ► The dynamic per parameter step-size is computed w.r.t. previous gradient l2-norm ⇒ parameters with small (resp. large) gradient will have a large (resp. small) step-size

Adadelta [Zeiler, 2012]

- Dynamic per parameter rate is computed with a fixed window of past gradients
- ► Approximate second-order information to incorporate curvature information ⇒ less sensitive to the learning rate hyper-parameter!

	SGD	MOMENTUM	ADAGRAD
$\epsilon = 1e^0$	2.26%	89.68%	43.76%
$\epsilon = 1e^{-1}$	2.51%	2.03%	2.82%
$\epsilon = 1e^{-2}$	7.02%	2.68%	1.79%
$\epsilon = 1e^{-3}$	17.01%	6.98%	5.21%
$\epsilon = 1e^{-4}$	58.10%	16.98%	12.59%

Table 1. MNIST test error rates after 6 epochs of training forvarious hyperparameter settings using SGD, MOMENTUM,and ADAGRAD.

Adagrad [Duchi et al., 2011]

- ▶ Replace global step-size with dynamic per parameter step-size + global learning rate
- ► The dynamic per parameter step-size is computed w.r.t. previous gradient l2-norm ⇒ parameters with small (resp. large) gradient will have a large (resp. small) step-size

Adadelta [Zeiler, 2012]

- Dynamic per parameter rate is computed with a fixed window of past gradients
- ► Approximate second-order information to incorporate curvature information ⇒ less sensitive to the learning rate hyper-parameter!

	SGD	MOMENTUM	ADAGRAD
$\epsilon = 1e^0$	2.26%	89.68%	43.76%
$\epsilon = 1e^{-1}$	2.51%	2.03%	2.82%
$\epsilon = 1e^{-2}$	7.02%	2.68%	1.79%
$\epsilon = 1e^{-3}$	17.01%	6.98%	5.21%
$\epsilon = 1e^{-4}$	58.10%	16.98%	12.59%

Table 1. MNIST test error rates after 6 epochs of training forvarious hyperparameter settings using SGD, MOMENTUM,and ADAGRAD.

	$\rho = 0.9$	$\rho = 0.95$	$\rho = 0.99$
$\epsilon = 1e^{-2}$	2.59%	2.58%	2.32%
$\epsilon = 1e^{-4}$	2.05%	1.99%	2.28%
$\epsilon = 1e^{-6}$	1.90%	1.83%	2.05%
$\epsilon = 1e^{-8}$	2.29%	2.13%	2.00%

Table 2. MNIST test error rate after 6 epochs for varioushyperparameter settings using ADADELTA.

Adam [Kingma and Ba, 2015]

Combine dynamic per parameter learning rate and momentum

Initialization bias correction

Convergence issue but works very well in practice [Reddi et al., 2018] Variants: AdaMax, Nadam [Dozat, 2016], Radam [Liu et al., 2019], AMSGrad

Adam [Kingma and Ba, 2015]

Combine dynamic per parameter learning rate and momentum

Initialization bias correction

Convergence issue but works very well in practice [Reddi et al., 2018] Variants: AdaMax, Nadam [Dozat, 2016], Radam [Liu et al., 2019], AMSGrad

Rule of thumb

- Optimizers based on adaptive learning rates usually work out of the box e.g. Adam is really popular in Natural Language Processing
- Fine-tuned SGD with step-size annealing can provide better results at the cost of expensive hyper-parameter tuning

Regularization issue

Weight decay is not equivalent to I2-norm when using adaptive learning rates!

References I

Dozat, T. (2016).
 Incorporating nesterov momentum into adam.
 ICLR Workshop.

Duchi, J., Hazan, E., and Singer, Y. (2011).

Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159.

Glorot, X. and Bengio, Y. (2010).

Understanding the difficulty of training deep feedforward neural networks. In Teh, Y. W. and Titterington, M., editors, *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of Machine Learning Research*, pages 249–256, Chia Laguna Resort, Sardinia, Italy. PMLR.

References II

```
He, K., Zhang, X., Ren, S., and Sun, J. (2015).
Delving deep into rectifiers: Surpassing human-level performance on imagenet
classification.
```

In Proceedings of the 2015 IEEE International Conference on Computer Vision (ICCV), ICCV '15, pages 1026–1034, Washington, DC, USA. IEEE Computer Society.

```
Hinton, G. E., Srivastava, N., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2012).
```

Improving neural networks by preventing co-adaptation of feature detectors. CoRR, abs/1207.0580.

Kingma, D. P. and Ba, J. (2015).

Adam: A method for stochastic optimization. *ICLR*.

References III

- Lee, H., Grosse, R., Ranganath, R., and Ng, A. Y. (2009). Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. pages 609–616.
- Liu, L., Jiang, H., He, P., Chen, W., Liu, X., Gao, J., and Han, J. (2019). On the variance of the adaptive learning rate and beyond. *arXiv preprint arXiv:1908.03265*.
- Nagaraj, D., Jain, P., and Netrapalli, P. (2019).
 SGD without replacement: Sharper rates for general smooth convex functions. In Chaudhuri, K. and Salakhutdinov, R., editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 4703–4711, Long Beach, California, USA. PMLR.

References IV

Polyak, B. T. (1964).

Some methods of speeding up the convergence of iteration methods. USSR Computational Mathematics and Mathematical Physics, 4(5):1–17.

Reddi, S. J., Kale, S., and Kumar, S. (2018). On the convergence of adam and beyond. *ICLR*.

 Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2014).
 Dropout: A simple way to prevent neural networks from overfitting. Journal of Machine Learning Research, 15:1929–1958.

References V

- Sutskever, I., Martens, J., Dahl, G., and Hinton, G. (2013).
 On the importance of initialization and momentum in deep learning.
 In Dasgupta, S. and McAllester, D., editors, *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 1139–1147, Atlanta, Georgia, USA. PMLR.
- Voita, E., Talbot, D., Moiseev, F., Sennrich, R., and Titov, I. (2019). Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned.

In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 5797–5808, Florence, Italy. Association for Computational Linguistics.

Zeiler, M. D. (2012).

Adadelta: an adaptive learning rate method.

arXiv preprint arXiv:1212.5701.

References VI

Zhang, C., Bengio, S., Hardt, M., Recht, B., and Vinyals, O. (2017). Understanding deep learning requires rethinking generalization. *ICLR 2017.*