

11 - Differentiating the Residual Neural Network

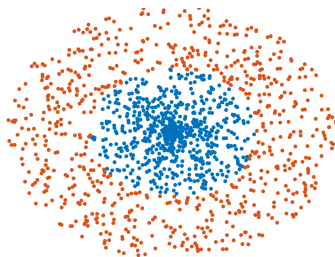
Numerical Methods for Deep Learning

March 19, 2018

Residual Network as a Path Planning Problem

$$\dot{\mathbf{Y}}(t) = \sigma(\mathbf{K}(t)\mathbf{Y}(t) + b(t)) \quad \mathbf{Y}(0) = \mathbf{Y}_0$$

The goal is to plan a path such that the transformed features, $\mathbf{Y}(T)$, can be linearly separated.

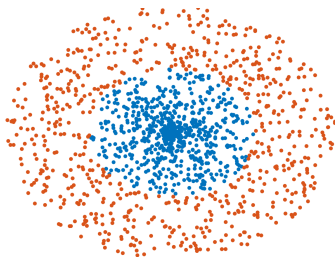


input features, $\mathbf{Y}(0)$

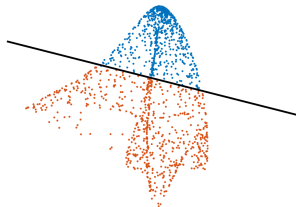
Residual Network as a Path Planning Problem

$$\dot{\mathbf{Y}}(t) = \sigma(\mathbf{K}(t)\mathbf{Y}(t) + b(t)) \quad \mathbf{Y}(0) = \mathbf{Y}_0$$

The goal is to plan a path such that the transformed features, $\mathbf{Y}(T)$, can be linearly separated.



input features, $\mathbf{Y}(0)$



transformed features $\mathbf{Y}(T)$

Residual Network - Forward Propagation

Idea: Obtain forward propagation by discretizing the ODE

$$\dot{\mathbf{Y}} = \sigma(\mathbf{KY} + b) \quad \mathbf{Y}(0) = \mathbf{Y}_0$$

Residual Network - Forward Propagation

Idea: Obtain forward propagation by discretizing the ODE

$$\dot{\mathbf{Y}} = \sigma(\mathbf{K}\mathbf{Y} + b) \quad \mathbf{Y}(0) = \mathbf{Y}_0$$

Example: Use forward Euler method

$$\mathbf{Y}_{j+1} = \mathbf{Y}_j + h\sigma(\mathbf{K}_j\mathbf{Y}_j + b_j)$$

Here: \mathbf{Y}_j is called the *state*, b_j are *controls*, and $h > 0$ is time step size

Residual Network - Forward Propagation

Idea: Obtain forward propagation by discretizing the ODE

$$\dot{\mathbf{Y}} = \sigma(\mathbf{K}\mathbf{Y} + b) \quad \mathbf{Y}(0) = \mathbf{Y}_0$$

Example: Use forward Euler method

$$\mathbf{Y}_{j+1} = \mathbf{Y}_j + h\sigma(\mathbf{K}_j\mathbf{Y}_j + b_j)$$

Here: \mathbf{Y}_j is called the *state*, j , b_j are *controls*, and $h > 0$ is time step size

More general forward propagation

$$\mathbf{Y}_{j+1} = \mathbf{P}_j\mathbf{Y}_j + h\sigma(\mathbf{K}_j\mathbf{Y}_j + b_j)$$

Allows for changing resolution and width.

Residual Network - Classification Problem

Classification with final state by solving

$$\min_{\mathbf{W}, \mathbf{K}_{0,\dots,N-1}, b_{0,\dots,N-1}} E(\mathbf{Y}_N(\mathbf{K}_{0,\dots,N-1}, b_{0,\dots,N-1})\mathbf{W}, \mathbf{C}^{\text{obs}})$$

Residual Network - Classification Problem

Classification with final state by solving

$$\min_{\mathbf{W}, \mathbf{K}_{0,\dots,N-1}, b_{0,\dots,N-1}} E(\mathbf{Y}_N(\mathbf{K}_{0,\dots,N-1}, b_{0,\dots,N-1})\mathbf{W}, \mathbf{C}^{\text{obs}})$$

Need to differentiate

- ▶ E w.r.t \mathbf{W} (linear classifier \leadsto Lecture 3)
- ▶ \mathcal{S} w.r.t \mathbf{Y}_N (single layer \leadsto Lecture 8)
- ▶ \mathbf{Y}_N w.r.t control variables $(\mathbf{K}_{0,\dots,N-1}, b_{0,\dots,N-1})$

Residual Network - Classification Problem

Classification with final state by solving

$$\min_{\mathbf{W}, \mathbf{K}_{0,\dots,N-1}, b_{0,\dots,N-1}} E(\mathbf{Y}_N(\mathbf{K}_{0,\dots,N-1}, b_{0,\dots,N-1})\mathbf{W}, \mathbf{C}^{\text{obs}})$$

Need to differentiate

- ▶ E w.r.t \mathbf{W} (linear classifier \leadsto Lecture 3)
- ▶ \mathcal{S} w.r.t \mathbf{Y}_N (single layer \leadsto Lecture 8)
- ▶ \mathbf{Y}_N w.r.t control variables $(\mathbf{K}_{0,\dots,N-1}, b_{0,\dots,N-1})$

Having these, apply chain rule to get, e.g.,

$$\nabla_{\mathbf{K}_j} E = (\nabla_{\mathbf{K}_j} \mathbf{Y}_N)^\top \nabla_{\mathbf{Y}_N} E$$

Computing Derivatives - Sensitivity Equation

Idea: Differentiate the forward propagation (forward Euler).

Let $0 \leq i \leq N$ be fixed. Note that

$$\nabla_{\kappa_i} \mathbf{Y}_j = 0, \quad \text{for } j \leq i.$$

Computing Derivatives - Sensitivity Equation

Idea: Differentiate the forward propagation (forward Euler).
Let $0 \leq i \leq N$ be fixed. Note that

$$\nabla_{\mathbf{K}_i} \mathbf{Y}_j = 0, \quad \text{for } j \leq i.$$

Next, note that

$$\nabla_{\mathbf{K}_i} \mathbf{Y}_{i+1} = h \text{diag}(\sigma'(\mathbf{K}_i \mathbf{Y}_i + b_i))(\mathbf{Y}_i^\top \otimes \mathbf{I})$$

Computing Derivatives - Sensitivity Equation

Idea: Differentiate the forward propagation (forward Euler).
Let $0 \leq i \leq N$ be fixed. Note that

$$\nabla_{\mathbf{K}_i} \mathbf{Y}_j = 0, \quad \text{for } j \leq i.$$

Next, note that

$$\nabla_{\mathbf{K}_i} \mathbf{Y}_{i+1} = h \text{diag}(\sigma'(\mathbf{K}_i \mathbf{Y}_i + b_i)) (\mathbf{Y}_i^\top \otimes \mathbf{I})$$

Continuing like this, gives for the final state:

$$\begin{aligned} \nabla_{\mathbf{K}_i} \mathbf{Y}_N &= \mathbf{P}_{N-1} \nabla_{\mathbf{K}_j} \mathbf{Y}_{N-1} \\ &+ h \text{diag}(\sigma'(\cdots)) ((\mathbf{I} \otimes \mathbf{K}_{N-1}) \nabla_{\mathbf{K}_j} \mathbf{Y}_{N-1}) \end{aligned}$$

Computing Derivatives - Sensitivity Equation

Idea: Differentiate the forward propagation (forward Euler).
Let $0 \leq i \leq N$ be fixed. Note that

$$\nabla_{\mathbf{K}_i} \mathbf{Y}_j = 0, \quad \text{for } j \leq i.$$

Next, note that

$$\nabla_{\mathbf{K}_i} \mathbf{Y}_{i+1} = h \text{diag}(\sigma'(\mathbf{K}_i \mathbf{Y}_i + b_i)) (\mathbf{Y}_j^\top \otimes \mathbf{I})$$

Continuing like this, gives for the final state:

$$\begin{aligned} \nabla_{\mathbf{K}_i} \mathbf{Y}_N &= \mathbf{P}_{N-1} \nabla_{\mathbf{K}_j} \mathbf{Y}_{N-1} \\ &+ h \text{diag}(\sigma'(\cdots)) ((\mathbf{I} \otimes \mathbf{K}_{N-1}) \nabla_{\mathbf{K}_j} \mathbf{Y}_{N-1}) \end{aligned}$$

Next: Write this as a block triangular **linear** system.

Computing Derivatives - Sensitivity Equations

Block triangular **linear** system for the gradients

$$\begin{pmatrix} \mathbf{I} & & & \\ -\mathbf{T}_{i+1} & \mathbf{I} & & \\ & \ddots & \ddots & \\ & & -\mathbf{T}_{N-1} & \mathbf{I} \end{pmatrix} \begin{pmatrix} \nabla_{\mathbf{K}_i} \mathbf{Y}_{i+1} \\ \\ \\ \nabla_{\mathbf{K}_i} \mathbf{Y}_N \end{pmatrix} = \begin{pmatrix} \mathbf{R}_i \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$\mathbf{T}_j = \mathbf{P}_j + h \text{diag}(f'(\mathbf{K}_j \mathbf{Y}_j + b_j))(\mathbf{I} \otimes \mathbf{K}_j)$$

and

$$\mathbf{R}_i = h \text{diag}(f'(\mathbf{K}_i \mathbf{Y}_i + b_i))(\mathbf{Y}_i^\top \otimes \mathbf{I}).$$

Computing Derivatives - Sensitivity Equation

Block triangular **linear** system for the gradients

$$\underbrace{\begin{pmatrix} \mathbf{I} & & & \\ -\mathbf{T}_{i+1} & \mathbf{I} & & \\ & \ddots & \ddots & \\ & & -\mathbf{T}_{N-1} & \mathbf{I} \end{pmatrix}}_{=\mathbf{T}} \underbrace{\begin{pmatrix} \nabla_{\mathbf{K}_i} \mathbf{Y}_{i+1} \\ \vdots \\ \nabla_{\mathbf{K}_i} \mathbf{Y}_N \end{pmatrix}}_{=\nabla_{\mathbf{K}_i} \mathbf{Y}} = \underbrace{\begin{pmatrix} \mathbf{R}_i \\ 0 \\ \vdots \\ 0 \end{pmatrix}}_{=\mathbf{R}}$$

To compute matrix-vector product $(\nabla_{\mathbf{K}_i} \mathbf{Y}_N) \mathbf{v}$

- ▶ Multiply $\mathbf{R} \mathbf{v}$
- ▶ Solve (forward propagate) $\mathbf{T} \nabla_{\mathbf{K}_j} \mathbf{Y}_N = \mathbf{R} \mathbf{v}$
- ▶ Extract the last time step

The sensitivity equation

Symbolically

$$\nabla_{\mathbf{K}_j} \mathbf{Y}_N = \mathbf{Q} \mathbf{T}^{-1} \mathbf{R}$$

where

$$\mathbf{Q} = [0, \dots, \mathbf{I}].$$

The transpose

$$(\nabla_{\mathbf{K}_i} \mathbf{Y}_N)^\top = \mathbf{R}^\top \mathbf{T}^{-\top} \mathbf{Q}^\top$$

The Sensitivity Equation

$$(\nabla_{\mathbf{\kappa}_i} \mathbf{Y}_N)^\top = \mathbf{R}^\top \mathbf{T}^{-\top} \mathbf{Q}^\top$$

$$(\mathbf{R}_i^\top \quad 0 \quad \dots \quad 0) \begin{pmatrix} \mathbf{I} & -\mathbf{T}_{i+1}^\top & & & \\ & \mathbf{I} & -\mathbf{T}_{i+2}^\top & & \\ & & \ddots & \ddots & \\ & & & \mathbf{I} & -\mathbf{T}_N^\top \\ & & & & \mathbf{I} \end{pmatrix}^{-1} \begin{pmatrix} 0 \\ \vdots \\ \mathbf{I} \end{pmatrix}$$

To multiply by the transpose

- ▶ Initialize with last step
- ▶ **solve backward** in time
- ▶ Extract the first step and multiply by $(\nabla_{\mathbf{\kappa}_j} f_j)^\top$

More about the sensitivity equation

To compute $(\nabla_{\mathbf{K}_i} \mathbf{Y}_N)^\top$ for all i 's note that the same quantities are recomputed. Can be evaluated in $\mathcal{O}(N)$

For gradient based method the transpose is sufficient

Newton based methods require both forward sensitivities and adjoint.

Testing Derivatives

Task 1: Programming the derivative test

as usual

Task 2: Programming the adjoint - the adjoint test

Code a - Computes $(\nabla_{\mathbf{K}_i} \mathbf{Y}_N) \mathbf{v}$

Code b - Computes $(\nabla_{\mathbf{K}_i} \mathbf{Y}_N)^\top \mathbf{u}$

Testing - for random \mathbf{u}, \mathbf{v}

$$\mathbf{u}^\top ((\nabla_{\mathbf{K}_i} \mathbf{Y}_N) \mathbf{v}) = \mathbf{v}^\top ((\nabla_{\mathbf{K}_i} \mathbf{Y}_N)^\top \mathbf{u})$$