

AUTOMATIC DIFFERENTIATION

DEEP LEARNING KU (DAT.C302UF)

Simon Hitzgiger

Oct 29, 2025

Institute of Machine Learning and Neural Computation
Graz University of Technology, Austria

THE GOAL OF SUPERVISED LEARNING

Training Data

- Given training dataset $\mathcal{D}_{\text{train}} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$
- Features $\mathbf{x} \in \mathbb{R}^d$

Training Data

- Given training dataset $\mathcal{D}_{\text{train}} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$
- Features $\mathbf{x} \in \mathbb{R}^d$

Classification

$$y \in \mathcal{C}$$

e.g., $\mathcal{C} = \{0, 1, \dots, 9\}$

Regression

$$y \in \mathbb{R} \quad \text{or} \quad \mathbf{y} \in \mathbb{R}^k$$

Model

- We wish to find a **model** f_{θ} that predicts y from \mathbf{x}
 - $\theta \in \mathbb{R}^n$ are the **model parameters**
- Can be a **point estimate** (e.g. if $y \in \mathbb{R}$):

$$f_{\theta}(\mathbf{x}) = \hat{y}$$

- Or a **distribution** over outputs (e.g. if $y \in \mathcal{C} = \{c_1, \dots, c_k\}$):

$$f_{\theta}(\mathbf{x}) = (p_{\theta}(y = c_1 | \mathbf{x}), \dots, p_{\theta}(y = c_k | \mathbf{x}))^{\top}$$

Model

- We wish to find a **model** f_{θ} that predicts y from x
 - $\theta \in \mathbb{R}^n$ are the **model parameters**
- Can be a **point estimate** (e.g. if $y \in \mathbb{R}$):

$$f_{\theta}(x) = \hat{y}$$

- Or a **distribution** over outputs (e.g. if $y \in \mathcal{C} = \{c_1, \dots, c_k\}$):

$$f_{\theta}(x) = (p_{\theta}(y = c_1 | x), \dots, p_{\theta}(y = c_k | x))^{\top}$$

Loss Function

- We specify a **loss function**

$$\ell(f_{\theta}(x), y) \in \mathbb{R}$$

that measures **deviation** between the prediction $f_{\theta}(x)$ and the “true” solution y

Training Loss 💪

- We can now compute the **average loss** over all training examples:

$$\mathcal{L}_{\text{train}}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \ell \left(f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}), y^{(i)} \right)$$

- Notice that $\mathcal{L}_{\text{train}} : \mathbb{R}^n \rightarrow \mathbb{R}$!

Question 🤔

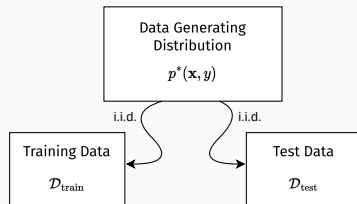
Is the final goal of supervised learning to minimize the average training loss? That is, finding

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \mathcal{L}_{\text{train}}(\boldsymbol{\theta})$$

- Statistical Learning Theory can answer this!
- If we had access to $p^*(\mathbf{x}, y)$, we want to minimize the **generalization error**:

$$\mathcal{L}^*(\boldsymbol{\theta}) = \mathbb{E}_{p^*(\mathbf{x}, y)} [\ell(f_{\boldsymbol{\theta}}(\mathbf{x}), y)]$$

- But since we don't know p^* , we minimize $\mathcal{L}_{\text{train}}$
- However, we **truly care about the average test loss** $\mathcal{L}_{\text{test}}$, since this is an unbiased estimate of the generalization error



MINIMIZING TRAINING LOSS

- Let's minimize the training loss

$$\operatorname{argmin}_{\theta} \frac{1}{N} \sum_{i=1}^N \ell \left(f_{\theta}(\mathbf{x}^{(i)}), \mathbf{y}^{(i)} \right)$$

- In *Deep Learning*, there is generally no closed form solution for this 😞

- Let's minimize the training loss

$$\operatorname{argmin}_{\boldsymbol{\theta}} \frac{1}{N} \sum_{i=1}^N \ell \left(f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}), y^{(i)} \right)$$

- In *Deep Learning*, there is generally no closed form solution for this 😞
- Hence, we rely on **iterative optimization algorithms** that use **gradient information**:

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\text{train}}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \nabla_{\boldsymbol{\theta}} \ell \left(f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}), y^{(i)} \right)$$

- Is there a problem with this ?

- Let's minimize the training loss

$$\operatorname{argmin}_{\boldsymbol{\theta}} \frac{1}{N} \sum_{i=1}^N \ell \left(f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}), y^{(i)} \right)$$

- In *Deep Learning*, there is generally no closed form solution for this 😞
- Hence, we rely on **iterative optimization algorithms** that use **gradient information**:

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\text{train}}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \nabla_{\boldsymbol{\theta}} \ell \left(f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}), y^{(i)} \right)$$

- Is there a problem with this ?
- For large N , this is computationally expensive. Instead, let's approximate this gradient with a **Mini-Batch Gradient** (with $B \ll N$):

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\text{train}}(\boldsymbol{\theta}) \approx \nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{B}}(\boldsymbol{\theta}) = \frac{1}{B} \sum_{i=1}^B \nabla_{\boldsymbol{\theta}} \ell \left(f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}), y^{(i)} \right)$$

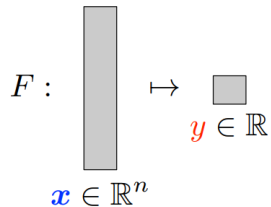
- Thankfully, we don't need to compute the gradients by hand 🙏
- Deep Learning Frameworks offer **Automatic Differentiation** (autodiff) 🥰

AUTOMATIC DIFFERENTIATION

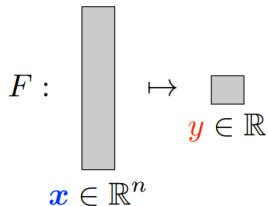
CREDIT: MATT JOHNSON'S GREAT TALK ON AUTODIFF

[JOHNSON, 2017]

$$F : \mathbb{R}^n \rightarrow \mathbb{R}$$



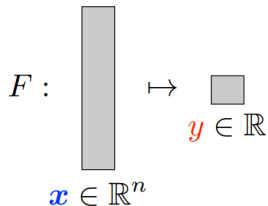
$$F : \mathbb{R}^n \rightarrow \mathbb{R}$$



$$F = D \circ C \circ B \circ A$$

$$y = F(\mathbf{x}) = D(C(B(A(\mathbf{x}))))$$

$$F : \mathbb{R}^n \rightarrow \mathbb{R}$$



$$F = D \circ C \circ B \circ A \qquad y = F(\mathbf{x}) = D(C(B(A(\mathbf{x}))))$$

$$y = D(\mathbf{c}), \quad \mathbf{c} = C(\mathbf{b}), \quad \mathbf{b} = B(\mathbf{a}), \quad \mathbf{a} = A(\mathbf{x})$$

$$y = D(c), \quad c = C(b), \quad b = B(a), \quad a = A(x)$$

$$F'(x) = \frac{\partial y}{\partial x} = \left[\frac{\partial y}{\partial x_1} \quad \cdots \quad \frac{\partial y}{\partial x_n} \right]$$

$$y = D(c), \quad c = C(b), \quad b = B(a), \quad a = A(x)$$

$$F'(x) = \frac{\partial y}{\partial x} = \left[\frac{\partial y}{\partial x_1} \quad \cdots \quad \frac{\partial y}{\partial x_n} \right]$$

$$F'(x) = \frac{\partial y}{\partial c} \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial x}$$

$$y = D(c), \quad c = C(b), \quad b = B(a), \quad a = A(x)$$

$$F'(x) = \frac{\partial y}{\partial x} = \left[\frac{\partial y}{\partial x_1} \quad \cdots \quad \frac{\partial y}{\partial x_n} \right]$$

$$F'(x) = \frac{\partial y}{\partial c} \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial x}$$

$$\frac{\partial y}{\partial c} = D'(c)$$

$$\frac{\partial c}{\partial b} = C'(b)$$

$$\frac{\partial b}{\partial a} = B'(a)$$

$$\frac{\partial a}{\partial x} = A'(x)$$



$$F'(\mathbf{x}) = \frac{\partial y}{\partial \mathbf{c}} \left(\underbrace{\frac{\partial \mathbf{c}}{\partial \mathbf{b}} \left(\frac{\partial \mathbf{b}}{\partial \mathbf{a}} \quad \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \right)}_{\text{Forward accumulation}} \right)$$

$$\frac{\partial \mathbf{b}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial b_1}{\partial x_1} & \dots & \frac{\partial b_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial b_m}{\partial x_1} & \dots & \frac{\partial b_m}{\partial x_n} \end{bmatrix}$$

Forward
accumulation

$$F'(\mathbf{x}) = \frac{\partial y}{\partial \mathbf{c}} \underbrace{\left(\frac{\partial \mathbf{c}}{\partial \mathbf{b}} \left(\frac{\partial \mathbf{b}}{\partial \mathbf{a}} \quad \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \right) \right)}_{\frac{\partial \mathbf{b}}{\partial \mathbf{x}}}$$

$$\frac{\partial \mathbf{b}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial b_1}{\partial x_1} & \cdots & \frac{\partial b_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial b_m}{\partial x_1} & \cdots & \frac{\partial b_m}{\partial x_n} \end{bmatrix}$$

Forward
accumulation

$$F'(\mathbf{x}) = \underbrace{\left(\left(\frac{\partial y}{\partial \mathbf{c}} \quad \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \right) \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \right)}_{\frac{\partial y}{\partial \mathbf{b}}} \frac{\partial \mathbf{a}}{\partial \mathbf{x}}$$

$$\frac{\partial y}{\partial \mathbf{b}} = \begin{bmatrix} \frac{\partial y}{\partial b_1} & \cdots & \frac{\partial y}{\partial b_m} \end{bmatrix}$$

Reverse
accumulation

Reverse-Mode Autodiff (“Backpropagation”)

- `loss.backward()` does reverse-mode autodiff !
 - Hence this will fail if `loss` is not a scalar
- Implemented using only Vector-Jacobian products $v^T J$
- Backward pass needs \approx same amount of compute than forward pass 🥰
- But we need to store all activations in the forward pass 😓

Reverse-Mode Autodiff (“Backpropagation”)

- `loss.backward()` does reverse-mode autodiff !
 - Hence this will fail if `loss` is not a scalar
- Implemented using only Vector-Jacobian products $v^T J$
- Backward pass needs \approx same amount of compute than forward pass 🥰
- But we need to store all activations in the forward pass 😓

Forward-Mode Autodiff

- Constant memory overhead 🥰
- Compute inefficient if we compute gradients of $\mathcal{L} : \mathbb{R}^n \rightarrow \mathbb{R}$ 😓
- In typical Deep Learning scenarios, you rarely care about forward-mode AD

AUTODIFF DEMO IN PyTorch



Johnson, M. J. (2017).

Automatic differentiation.

Presented at the Deep Learning Summer School 2017 in Montreal.