

Training a neural network



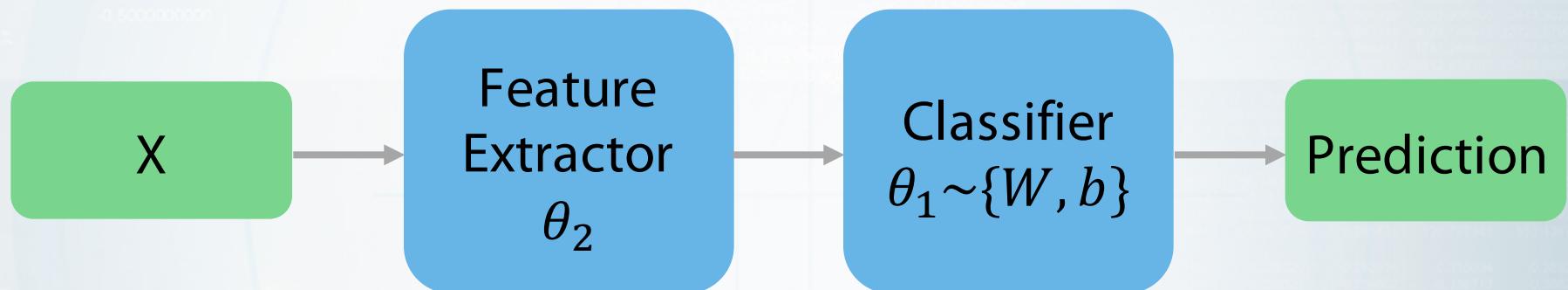
In this lecture

Learn how to train a deep learning model:

- Computing the gradient
- Backpropagation



Training?



Training:

$$\underset{\theta_1, \theta_2}{\operatorname{argmin}} L(y, P(y|x))$$



Optimization!

Deep model is just a (composite!) function

- $f_1(f_2(f_3(x, w_3), w_2), w_1) \rightarrow \text{loss}$
- x – features vector
- $w_{1,2,3}$ – model parameters

Remember those optimization methods from the last week?
But they need gradient...



Differentiation recap

$$f: \mathbb{R} \rightarrow \mathbb{R}$$

You know what to do



Differentiation recap

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

X is now a vector $[x_1, x_2, \dots, x_n]$

f' is now a vector of partial derivatives $\frac{df}{dx} = \left[\frac{df}{dx_1}, \frac{df}{x_2}, \dots, \frac{df}{x_n} \right]$



Differentiation recap

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

X is a vector $[x_1, x_2, \dots, x_n]$

f is also vector $[f_1, f_2, \dots, f_m]$

f' is now a Jacobian matrix

$$\frac{df}{dx} = \left[\frac{df}{dx_1} \cdots \frac{df}{dx_n} \right] = \begin{bmatrix} \frac{df_1}{dx_1} & \cdots & \frac{df_1}{dx_n} \\ \vdots & \ddots & \vdots \\ \frac{df_m}{dx_1} & \cdots & \frac{df_m}{dx_n} \end{bmatrix}$$



Matrix derivative

$$F: \mathbb{R}^{m \cdot n} \rightarrow \mathbb{R}^{p \cdot q}$$

X is a matrix

F is also a matrix

F' is now

- a matrix of matrices
- fourth-rank tensor
- four-dimensional array

$$\frac{\partial F}{\partial X} = \begin{bmatrix} \frac{\partial F}{\partial X_{11}} & \cdots & \frac{\partial F}{\partial X_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial F}{\partial X_{m1}} & \cdots & \frac{\partial F}{\partial X_{mn}} \end{bmatrix}$$

$$\frac{\partial F}{\partial X_{ij}} = \begin{bmatrix} \frac{\partial F_{11}}{\partial X_{ij}} & \cdots & \frac{\partial F_{1q}}{\partial X_{ij}} \\ \vdots & & \vdots \\ \frac{\partial F_{p1}}{\partial X_{ij}} & \cdots & \frac{\partial F_{pq}}{\partial X_{ij}} \end{bmatrix}$$

is a $p \times q$ matrix



Matrix derivative example

$$W^2 = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix}^2 = \begin{bmatrix} w_{11}^2 + w_{12}w_{21} & w_{11}w_{12} + w_{12}w_{22} \\ w_{21}w_{11} + w_{22}w_{21} & w_{21}w_{12} + w_{22}^2 \end{bmatrix}$$

$$\frac{dW^2}{dW} = \begin{bmatrix} \begin{bmatrix} 2w_{11} & w_{12} \\ w_{21} & 0 \\ w_{12} & 0 \\ w_{11} + w_{22} & w_{12} \end{bmatrix} & \begin{bmatrix} w_{21} & w_{11} + w_{22} \\ 0 & w_{21} \\ 0 & w_{12} \\ w_{21} & 2w_{22} \end{bmatrix} \end{bmatrix}$$

4 Dimensions



Chain rule

$$L(x) = g(f(x))$$

For $f, g : \mathbb{R} \rightarrow \mathbb{R}$ $\frac{dL}{dx} = \frac{dg}{df} \cdot \frac{df}{dx}$

$$\frac{dL}{dx} \Big|_{x=x_0} = \frac{dg}{df} \Big|_{u=f(x_0)} \cdot \frac{df}{dx} \Big|_{x=x_0}$$



Chain rule example

$$\frac{d}{d\vec{x}} \sin(|\vec{x}|) \Big|_{\vec{x}=\vec{a}} = \frac{d}{dy} \sin(y) \Big|_{y=\vec{a}} \cdot \frac{d}{d\vec{x}} |\vec{x}| \Big|_{\vec{x}=\vec{a}} =$$

$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$\vec{a} = \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}$$

$$= \cos(|\vec{a}|) \cdot \frac{d\sqrt{x_1^2 + x_2^2}}{d \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}} \Big|_{\vec{x}=\vec{a}} =$$
$$= \cos(|\vec{a}|) \cdot \frac{1}{\sqrt{a_1^2 + a_2^2}} \cdot \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}$$



Chain rule quiz (технический слайд)



Backpropagation

Chain rule can be evaluated numerically!

1. Compute network output and the loss value
2. Compute $\frac{d\text{Loss}}{d\text{Activation_of_the_last_layer}}$
3. For each layer, starting from the last:
 1. Compute $\frac{d\text{Activation}}{d\text{Layer_parameters}}, \frac{d\text{Activation}}{d\text{Layer_input}}$
 2. Multiply it by $\frac{d\text{Loss}}{d\text{Activation}}$, get $\frac{d\text{Loss}}{\dots}$
 3. Make optimization step for the parameters



Summary

- Neural network training is an optimization problem
- Commonly solved via gradient-based methods
- Gradients are computed via backpropagation

