Introduction to Neural Networks

November 17th 2020 KMI-2020

Nagoya University



Kazuhiro Terao SLAC National Lab



Introduction to Neural Networks Perceptron

 \overrightarrow{x}

The basic unit of a neural net is the *perceptron* (loosely based on a real neuron)

Takes in a vector of inputs (*x*). Commonly inputs are summed with weights (*w*) and offset (*b*) then run through activation.



Imagine using two features to separate cats and dogs



$$\sigma(\vec{x}) = \begin{cases} \mathbf{1} & \vec{w}_i \cdot \vec{x} + b_i \ge 0\\ \mathbf{0} & \vec{w}_i \cdot \vec{x} + b_i < 0. \end{cases}$$



By picking a value for **w** and **b**, we define a boundary between the two sets of data

Introduction to Neural Networks Perceptron

What if we have a new data point?





Introduction to Neural Networks Perceptron

What if we have a new data point?



from wikipedia



We can add **another neuron** to help (but does not yet solve the problem)

Introduction to Neural Networks Multi-Layer Perceptron

What if we have a new data point?



based on preceding layer's output (of **non-linear activation**)



Another layer can classify

SL AC

Introduction to Neural Networks Multi-Layer Perceptron



"Traditional neural net" Fully-Connected Multi-Layer Perceptrons



A traditional neural network consists of a stack of layers of such neurons where each neuron is *fully connected* to other neurons of the neighbor layers

Training Feed-forward Neural Networks Gradient-based Method





Introduction to Neural Networks Training (Optimization)

N-layer MLP

 $x_0 \dots$ input data $F_i \dots i$ -th hidden layer $x_i \dots i$ -th layer output $w_i \dots i$ -th layer weights $y \dots$ correct answer $\mathcal{L} \dots$ loss (cost) function

🖡 Xn $F_n(x_{n-1}, w_n)$ Forward Path $F_i(x_{i-1}, w_i)$ **↑** X₂ $F_{2}(x_{1}, w_{2})$ X_1 $F_{1}(x_{0}, w_{1})$

Introduction to Neural Networks Training (Optimization)

Loss: a cost to minimize (e.g. error).

N-layer MLP

 $x_0 \dots$ input data $F_i \dots i$ -th hidden layer $x_i \dots i$ -th layer output $w_i \dots i$ -th layer weights $y \dots$ correct answer $\mathcal{L} \dots$ loss (cost) function



Xo

Introduction to Neural Networks Training (Optimization)

Loss: a cost to minimize (e.g. error).

Training: an optimization of parameters by minimizing loss.



Loss: a cost to minimize (e.g. error).

Training: an optimization of parameters by minimizing loss.

Gradient Descend: an iterative approach to optimize weights in order to minimize the loss.

- Define & measure the error
- Compute $\partial \mathcal{L} / \partial w_n = \nabla W$
- Update weights: $W_{new} = W \lambda \nabla W$... where λ is called a "learning rate"



How can we compute the gradient for the *i*-th layer?

SLAC

 $\mathcal{L}(x_n, y) = \mathcal{L}(F_n(x_{n-1}, w_n), y) = ... = \mathcal{L}(F_n(...F_i(x_{i-1}, w_i)...), y)$

How can we compute the gradient for the *i*-th layer?

 $\mathcal{L}(x_n, y) = \mathcal{L}(F_n(x_{n-1}, w_n), y) = ... = \mathcal{L}(F_n(...F_i(x_{i-1}, w_i)...), y)$



How can we compute the gradient for the *i*-th layer?

 $\mathcal{L}(x_n, y) = \mathcal{L}(F_n(x_{n-1}, w_n), y) = ... = \mathcal{L}(F_n(...F_i(x_{i-1}, w_i)...), y)$

 $\frac{\partial \mathcal{L}}{\partial w_{i}} = \frac{\partial \mathcal{L}}{\partial x_{n}} \cdot \frac{\partial x_{n}}{\partial x_{n-1}} \cdot \frac{\partial x_{n-1}}{\partial x_{n-2}} \cdots \frac{\partial x_{i+1}}{\partial x_{i}} \cdot \frac{\partial x_{i}}{\partial w_{i}} \quad \text{(chain rule)}$

SLAC

Example: quadratic loss = $\sum_{m=1}^{m} (x_n - y)^2$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{X}_{n}} = \sum_{m=1}^{m} (2\mathbf{X}_{n} - 2\mathbf{y})$$

... or $\frac{\partial \mathcal{L}}{\partial x_n} = \sum_{m=1}^{m} (2x_n - 2y)^T$ if x_n and y are column vector

How can we compute the gradient for the *i*-th layer?

$$\mathcal{L}(x_n, y) = \mathcal{L}(F_n(x_{n-1}, w_n), y) = ... = \mathcal{L}(F_n(...F_i(x_{i-1}, w_i)...), y)$$



How can we compute the gradient for the *i*-th layer? *i*-th layer



How can we compute the gradient for the *i*-th layer? *i*-th layer **Input:** $X_{i-1} = \frac{\partial \mathcal{L}}{\partial x_i}$ (given) $\begin{array}{c|c} \text{Homogeneous}\\ \text{Homogeneous}\\ \text{Homogeneous}\\ \text{X}_{i} & \begin{array}{c} & \frac{\partial \mathcal{L}}{\partial X_{i}} \\ \text{Fi} (X_{i-1}, W_{i}) \\ \text{Ki-1} & \begin{array}{c} & \frac{\partial \mathcal{L}}{\partial X_{i-1}} \\ \frac{\partial \mathcal{L}}{\partial X_{i-1}} \end{array} \end{array}$ **Output** (we compute): $\mathbf{X}_{i} = \mathbf{F}_{i}(\mathbf{X}_{i-1}, \mathbf{W}_{i})$ $\frac{\partial \mathcal{L}}{\partial x_{i-1}} = \frac{\partial \mathcal{L}}{\partial x_i} \cdot \frac{\partial F_i(x_{i-1}, w_i)}{\partial x_{i-1}}$

 $\begin{array}{lll} \textbf{Weight Update:} & w_i - \lambda \nabla w_i \\ \textbf{where:} & \nabla w_i = \frac{\partial \mathcal{L}}{\partial w_i} = \frac{\partial \mathcal{L}}{\partial x_i} \cdot \frac{\partial F_i(x_{i-1}, w_i)}{\partial w_i} \end{array} \\ \end{array}$

How can we compute the gradient for the *i*-th layer?



Example: linear transformation

SLAC

$$F(x_{in}, w) = w \cdot x_{in}$$

Let's generalize for a tensor input

How can we compute the gradient for the *i*-th layer?



Example: linear transformation $F(X_{in}, W) = W \cdot X_{in}$ Let's generalize for a tensor input $\frac{\partial \mathcal{L}}{\partial X_{in}} = \frac{\partial \mathcal{L}}{\partial X_{out}} \cdot W$ e.g.)

How can we compute the gradient for the *i*-th layer?



$$w_{\text{new}} = w - \lambda \left(\frac{\partial \mathcal{L}}{\partial w}\right)^{\mathrm{T}}$$

Example: linear transformation

SLAC

$$F(x_{in}, w) = w \cdot x_{in}$$

Let's generalize for a tensor input



Introduction to Neural Networks Activation Functions



Weights Initialization

• Random values to set different neuron states

si ac

• Not too large or too small: gradients become negligible for *tanh* and *sigmoid* activation layers! (positive large value not an issue for ReLU)

• Gaussian @ (μ, σ) = (0,1)?

- Slight modification: normalize such that the variance of a signal strength stays similar (0~1) across layers (depends on # of filters)
- See more details <u>here (CS231)</u>
- <u>Xavier initialization</u>, <u>He (MSRA) initialization</u>

Start: some random initial set of parameters
Goal: minimize the overall loss = sum over all samples
How: SGD from yesterday :)



Introduction to Neural Networks Gradient-based Optimization

Adaptive learning rate (LR) Ideally start with a larger LR to quickly converge, then decrease the LR to fine-tune near the minimum. Warning! too small LR rate is not only slow but can trap in false minimum.



Adaptive learning rate (LR) Ideally start with a larger LR to quickly converge, then decrease the LR to fine-tune near the minimum. Warning! too small LR rate is not only

slow but can trap in false minimum.

Momentum

The direction of the steepest descent can be almost perpendicular if the ellipse is elongated... A solution is to average the history of past updates. This accumulates small updates in the right direction

Introduction to Neural Networks Gradient-based Optimization

Popular Visualization of Optimizers (not mine)

SLAC

Image Credit: Alec Radford



- Excellent lecture <u>here (CS232)</u>
- Want to try by yourself ? You can get started <u>here</u>.
- Popular choices: vanila SGD and Adam

Neural Network

• Each neuron produces a "feature", activation function can add non-linearity, and additional layer can represent non-linear functional approximation.

Gradient-based Optimization

- Essentially a composite function = chain rule!
- Initialization should be non-zero random values
- Flavors of gradient-based optimizers to take into account for the error surface



Backup Slides

Cost Functions & Activation Functions





1. Mean Absolute Error (MAE, L1 loss)

$$MAE = \frac{\sum_{n=1}^{n} |x_{pred} - y_{true}|}{n}$$

- Constant magnitude gradient everywhere
- Adaptive LR critical near the minimum



1. Mean Absolute Error (MAE, L1 loss)

$$MAE = \frac{\sum_{n=1}^{n} |x_{pred} - y_{true}|}{n}$$

• Constant magnitude gradient everywhere

SLAC

• Adaptive LR critical near the minimum

2. Mean Square Error (MSE, L2 loss)

$$MSE = \frac{\sum_{n=1}^{n} (x_{pred} - y_{true})^2}{n}$$

- Adaptive LR not needed
- Huge gradient far from the minimum

1. Mean Absolute Error (MAE, L1 loss)

$$MAE = \frac{\sum_{n=1}^{n} |x_{pred} - y_{true}|}{n}$$

• Constant magnitude gradient everywhere

SLAC

• Adaptive LR critical near the minimum

2. Mean Square Error (MSE, L2 loss)

$$MSE = \frac{\sum_{n=1}^{n} (x_{pred} - y_{true})^2}{n}$$

- Adaptive LR not needed
- Huge gradient far from the minimum

3. Hubar Loss

Hubar =
$$\begin{cases} 0.5 (x_{\text{pred}} - y_{\text{true}})^2 & \dots \text{ if } |x_{\text{pred}} - y_{\text{true}}| < \delta \\ \delta |x_{\text{pred}} - y_{\text{true}}| - 0.5 \delta^2 & \dots \text{ otherwise} \end{cases}$$

Cross-Entropy Loss (entropy, softmax)

- Also called multinomial logistic loss
- Probability of being a class *i* in N categories

$$p_{i}(y_{i}) = \frac{\exp(x_{i})}{\sum_{j=1}^{N} \exp(x_{j})}$$

- Binomial is a special case where probability becomes **sigmoid**
- Entropy loss minimization is **maximum likelihood** method.
- Others: hinge (SVM) loss

• Cross-entropy loss

$$\mathcal{L}(x_{i}, y_{i}) = -\log \left[\frac{\exp(x_{i})}{\sum_{j=1}^{N} \exp(x_{j})} \right]$$