



## Back Propagation: training neural network

- gradient descent + chain rule
- optimize weights to minimize objective function
- prediction no longer linear
- early stopping, regularization = prevent overfitting
- No guarantees of convergence; many reach local minimum initial weights
- to avoid local minima: several trials w/ diff random initial weights
- many epochs for adequate training (10k+)
- Trainable criterion: # of epochs, threshold on train error, nodes/epoch in error
- Dropout: non-linear error surface
- regularization, adding noise, data augmentation, dropout
- Many hidden units = overfitting: cross-validation techniques for # of hidden units
- Dropout: randomly deactivates neurons during training

NN can learn feature representation

- Representation Learning:
  - Enables models to automatically extract useful features from raw data.
  - Examples: Embeddings in natural language processing (e.g., Word2Vec, BERT) or convolutional layers in image classification.

## Mathematics:

- Loss Function for Binary Classification:

$$L = -\frac{1}{n} \sum_{i=1}^n [y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)]$$

where  $\hat{y}_i = \sigma(\mathbf{w}^T \mathbf{x}_i + b)$ , and  $\sigma$  is the sigmoid function.

### Chain Rule in Backpropagation:

- For a neural network layer:  $\frac{\partial L}{\partial w_j} = \frac{\partial L}{\partial z_k} \frac{\partial z_k}{\partial w_j}$

## General Guide for Forward Pass

### 1. Start with the Loss

- Compute the derivative of the loss with respect to the output ( $\hat{y}$ ):

$$\frac{\partial L}{\partial \hat{y}}$$

### 2. Backpropagate Through the Output Layer

For the output layer:

- 1. Compute the gradient of the activation function:

$$\frac{\partial f}{\partial z_j} = f'(z_j)$$

- For sigmoid:  $f'(z) = \sigma(z)(1 - \sigma(z))$

- For ReLU:  $f'(z) = 1 \text{ if } z > 0, \text{ else } 0$

- 2. Combine with the gradient from the loss:

$$\frac{\partial L}{\partial z_j} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_j}$$

- 3. Compute gradients for weights and biases:

where  $\eta$  is the learning rate and  $L$  is the loss function.

### 3. Deep Learning, Backpropagation, and Representation Learning

#### Concepts

- Deep Learning:
  - A subfield of machine learning focused on learning hierarchical representations of data using deep neural networks (DNNs).

#### Backpropagation:

- A gradient-based optimization algorithm for training neural networks.

#### Steps:

- Forward Pass: Compute predictions and loss.
- Backward Pass: Compute gradients of the loss with respect to parameters using the chain rule.

#### 3. Update Parameters: Use gradient descent to update weights:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\partial L}{\partial \mathbf{w}}$$

## General Guide for Backward Pass

### 1. Start with the Loss

- Compute the derivative of the loss with respect to the output ( $\hat{y}$ ):

$$\frac{\partial L}{\partial \hat{y}}$$

### 2. Backpropagate Through the Output Layer

For the output layer:

- 1. Compute the gradient of the activation function:

$$\frac{\partial f}{\partial z_j} = f'(z_j)$$

- For sigmoid:  $f'(z) = \sigma(z)(1 - \sigma(z))$

- For ReLU:  $f'(z) = 1 \text{ if } z > 0, \text{ else } 0$

- 2. Combine with the gradient from the loss:

$$\frac{\partial L}{\partial z_j} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_j}$$

- 3. Compute gradients for weights and biases:

where  $\eta$  is the learning rate and  $L$  is the loss function.

### 4. Update Weights and Biases

Using gradient descent:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\partial L}{\partial \mathbf{w}}$$

$$b \leftarrow b - \eta \frac{\partial L}{\partial b}$$

$\eta$ : Learning rate (controls step size).

## Representation Learning

### What is Representation Learning?

- Representation learning refers to automatically learning features or representations from raw data that are useful for downstream tasks.
- It reduces the need for manual feature engineering and finds patterns in data that are often not apparent.

## Key Concepts

### 1. Hierarchical Features:

- Learn simple features first (e.g., edges in images) and combines them to form complex representations (e.g., faces).
- Deep learning models (e.g., CNNs, RNNs) excel in hierarchical representation learning.

### 2. Latent Representations:

- Encodes raw data into a lower-dimensional space while retaining meaningful information (e.g., embeddings in NLP).

### 3. Transferability:

- Learned representations can be used for different tasks (e.g., pre-trained BERT for various NLP tasks).

## Unsupervised Learning

- simply data, identify patterns, learn meaningful representations, identify outliers?
- Descriptive models surround; cluster analysis, density estimation

### Cluster Analysis:

- partition data into groups
- inter-group similarity low, intra-group similarity high
- well-separated clusters, center based clusters, contiguous based clusters
- density based, conceptual

### Clustering Algorithms:

- partition based models
- hierarchical clustering: agglomerative (bottom up), divisive (top-down)

### Partition-based clustering: k-means

- data:  $n$   $k$  clusters; Evaluation: score ( $L$ ) maximized/minimized

#### Most approaches used iterative improvement algorithms

#### Cluster Score Functions: score ( $L$ ) = measures quality of clustering $C$ for dataset $D$

- compact clusters: minimize within cluster distance ( $wc$ )

- separated clusters: maximize between cluster distance ( $bc$ )

#### Cluster centroid: $\mathbf{c}_k = \frac{1}{N_k} \sum_{i=1}^{N_k} \mathbf{x}_i$

- between-cluster dist:  $bc(L) = \sum_{k \neq l} d(c_k, c_l)^2$

- within-cluster dist:  $wc(L) = \sum_{k=1}^K \sum_{i=1}^{N_k} d(x_i, c_k)^2$

#### X-means: select $K$ centroids randomly

- Assign each data point to nearest centroid

- Update centroids by avg. mean of all points in each cluster

- Repeat steps 2-3 until centroids are stable

- Knowledge:  $K$  clusters are defined by canonical members (centroids)

- Model space: all possible partitions of the examples into  $K$  groups

- Score: minimize within cluster distance, search: iterative greedy

- Strengths: efficient, finds spherical clusters

- Weakness: terminates at local minima, sensitive to outliers, specifying  $K$ , applicable only when mean defined

- Variations: selection of initial centroids, elbow modification? allows for merge, split & combine

- can w/ multiple random selections, pick one w/ best score

### Probability mixture model:

- $f(x) = \sum_k p_k f_k(x; \theta)$  = prob. of observing  $x$

### Gaussian Mixture Models: assume data is generated from mixture of gaussian distributions

- uses expectation-maximization (EM) to estimate parameters

- $p(x|z) = \sum_k p_k(z) p(x|z \sim N(\mu_k, \Sigma_k))$

- Pros: soft clustering; works computationally intensive

### 2. Clustering Evaluation

#### Internal Metrics:

- Silhouette Score: Measures how similar a point is to its cluster compared to others. Range: -1 (bad clustering) to 1 (perfect clustering).

- Inertia (SSE): Sum of squared distances of points to their cluster center. Lower is better.

#### External Metrics (if ground truth exists):

- Adjusted Rand Index (ARI): Measures the similarity between predicted clusters and true labels.

- Normalized Mutual Information (NMI): Measures shared information between clusters and ground truth.

## Recent neural network:

- RNN version of a language model

- utilizes n-gram, RNN conditions current word on all previous words

- efficient time + space

- we define prob. dist. over V

- input vector:  $\mathbf{x}_1, \dots, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_n$

- Loss:  $L_{\text{RNN}}(\mathbf{x}) = -\sum_{i=1}^n \log p_{\text{out}}(\mathbf{x}_i)$

- $p_{\text{out}}(\mathbf{x}_i) = \sigma(\mathbf{W}_{\text{out}} \mathbf{x}_i + \mathbf{b}_{\text{out}})$

- $p_{\text{out}}(\mathbf{x}_i) = \text{softmax}(\mathbf{W}_{\text{out}} \mathbf{x}_i)$

- $p_{\text{out}}(\mathbf{x}_i) = \frac{e^{\mathbf{W}_{\text{out}} \mathbf{x}_i + \mathbf{b}_{\text{out}}}}{\sum_j e^{\mathbf{W}_{\text{out}} \mathbf{x}_j + \mathbf{b}_{\text{out}}}}$

- loss function

- chain Rule:  $\frac{\partial L}{\partial \mathbf{W}_{\text{out}}} = \frac{\partial L}{\partial \mathbf{x}_i} \frac{\partial \mathbf{x}_i}{\partial \mathbf{W}_{\text{out}}}$

- to avoid local minima: several trials w/ diff random initial weights

- many epochs for adequate training (10k+)

- Trainable criterion: # of epochs, threshold on train error, nodes/epoch in error

- Dropout: non-linear error surface

- regularization, adding noise, data augmentation, dropout

- Many hidden units = overfitting: cross-validation techniques for # of hidden units

- Dropout: randomly deactivates neurons during training

## General Guide for Forward Pass

### 1. Inputs to the Network

- Identify the input features ( $x_1, x_2, \dots, x_n$ ).

- If there are multiple layers, treat the outputs of the previous layer as inputs for the next

### 2. Weighted Sum for Each Neuron

for a given neuron  $j$  in a layer:

## General Guide for Forward Pass

### 1. Inputs to the Network

- Identify the input features ( $x_1, x_2, \dots, x_n$ ).

- If there are multiple layers, treat the outputs of the previous layer as inputs for the next

### 2. Weighted Sum for Each Neuron

for a given neuron  $j$  in a layer:

$$\sum_{i=1}^n w_{ij} x_i + b_j$$

- $w_{ij}$ : Weight connecting input  $x_i$  to neuron  $j$ .

- $b_j$ : Bias of the neuron.

- $x_i$ : Input to the neuron (can be raw inputs or outputs from the previous layer).

### 3. Activation Function

Apply the chosen activation function  $f(z_j)$  to the weighted sum:

$$a_j = f(z_j)$$

- Common activation functions:

- Sigmoid:  $f(z) = \frac{1}{1 + e^{-z}}$

- ReLU:  $f(z) = \max(0, z)$

- Tanh:  $f(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$

## 4. Output Layer

- Use the same procedure as above to compute the final outputs. If it's a regression problem, the output might be a raw value. For classification:

- Binary Classification: Use sigmoid activation.

- Multiclass Classification: Use softmax activation:

$$\text{Softmax}(x_i) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$$

## 5. Loss Function

Compute the loss for a given target  $y$ :

- Common loss functions:

- Mean Squared Error (MSE):

$$L = \frac{1}{n} \sum_{i=1}^n [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

- Binary Cross-Entropy:

$$L = -\frac{1}{n} \sum_{i=1}^n [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

### Step 1: Compute Hidden Layer Activations ( $z_1, z_2$ )

For each hidden unit  $z_1, z_2$ :

- $z_1 = \sigma(\alpha_{11}x_1 + \alpha_{12}x_2 + \alpha_{13}x_3)$

$$z_1 = \sigma((0.5)(1) + (-0.4)(2) + (0.3)(3))$$

- $z_2 = \sigma(0.5 - 0.8 + 0.9) = \sigma(0.6) = \frac{1}{1 + \exp(-0.6)} \approx 0.645$

$$z_2 = \sigma(-0.2 + 0.2 + 1.8) = \sigma(1.8) = \frac{1}{1 + \exp(-1.8)} \approx 0.858$$

Thus:

$$z_1 \approx 0.645, \quad z_2 \approx 0.858$$

### Step 2: Compute Output ( $y$ )

Using the output weights:

- $y = \sigma(\beta_1 z_1 + \beta_2 z_2)$ :

$$y = \sigma((0.7)(0.645) + (-0.5)(0.858))$$

- $y = \sigma(0.4515 - 0.429) = \sigma(0.0225)$

$$y = \frac{1}{1 + \exp(-0.0225)} \approx 0.506$$

## Example General Procedure

### 1. Forward Pass:

- Compute  $z_j$  for each neuron using inputs, weights, and biases.

- Apply activation functions to get  $a_j$ .

- Compute final output ( $\hat{y}$ ).

- Compute loss using the loss function.

### 2. Backward Pass:

- Start with the gradient of the loss w.r.t output ( $\hat{y}$ ).

- Backpropagate through the output layer to compute gradients for  $z_j, w_{ij}$ , and  $b_j$ .

- Backpropagate through all hidden layers in reverse order.

- Update weights and biases using the computed gradients.

## Key Notes

- Activation Function and Derivative: Always know the activation function for each layer and its derivative.

- Chain Rule: Backpropagation heavily relies on applying the chain rule to compute gradients layer by layer.

- Learning Rate ( $\eta$ ): Ensure a proper learning rate to avoid overshooting or slow convergence.

## Probabilistic clustering:

- Model provides full distributional description

- Soft clustering (k-means = hard)

- Key cost: optimization of parametric model

## Mixture models:

- Knowledge: parametric

- Store: likelihood

- Search: EM

## Probabilistic clustering:

- Model provides full distributional description

- Soft clustering (k-means = hard)

- Key cost: optimization of parametric model

## Hierarchical Methods

- constructs a hierarchy of nested clusters rather than picking  $k$  representat

- Approaches: agglomerative, divisive

## Agglomerative:

- for  $i \in n$ :

- while