

Linear Classification

Score function: $s = f(x; W) = Wx + b$

Softmax classifier: $P(y_i|x_i) = \frac{e^{sy_i}}{\sum_j e^{sj}}$

Cross-entropy loss: $L_i = -\log(P(y_i|x_i))$

- Min loss: 0 (when $P(y_i|x_i) = 1$)
- Max loss: ∞ (when $P(y_i|x_i) \approx 0$)
- Random initialization: $\log(C)$ for C classes

SVM loss: $L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + \Delta)$

- Δ is margin parameter (typically $\Delta = 1$)
- Wants correct class score higher than incorrect class scores by at least Δ
- Geometric interpretation: Linear hyperplanes separating classes

Key concepts:

- Any FC network can be expressed as a CNN (with 1x1 filters) and vice versa
- Loss gradients flow from softmax loss to weight matrix proportional to input
- Linear classifiers can't solve XOR problems (need non-linearities)

Regularization

Full loss: $L = \frac{1}{N} \sum_{i=1}^N L_i + \lambda R(W)$

Types:

- L2: $R(W) = \sum_k \sum_l W_{k,l}^2$ (prefers diffuse weights)
- L1: $R(W) = \sum_k \sum_l |W_{k,l}|$ (promotes sparsity)
- Elastic Net: $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$
- Dropout: Randomly zero outputs during training (scale by p at test)
- Normalization adds regularization effect due to noise in batch stats

Key concepts:

- Early stopping: end training when val error increases
- Dropout forces redundant representations, acts like ensemble averaging
- Regularizing bias terms is generally avoided (mainly regularize weights)
- Data augmentation: Adding transformed training examples

Optimization Algorithms

SGD: $w_{t+1} = w_t - \alpha \nabla L(w_t)$

SGD+Momentum:

$$v_{t+1} = \rho v_t + \nabla L(w_t) \quad (\text{typically } \rho = 0.9 \text{ or } 0.99)$$

$$w_{t+1} = w_t - \alpha v_{t+1}$$

RMSProp:

$$\text{grad_squared} = \beta \cdot \text{grad_squared} + (1 - \beta) \cdot (\nabla L(w_t))^2$$

$$w_{t+1} = w_t - \frac{\alpha \cdot \nabla L(w_t)}{\sqrt{\text{grad_squared} + \epsilon}}$$

Adam:

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) \nabla L(w_t) \quad (\text{momentum})$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) (\nabla L(w_t))^2 \quad (\text{RMSProp})$$

$$\hat{m}_t = m_t / (1 - \beta_1^t) \quad (\text{bias correction})$$

$$\hat{v}_t = v_t / (1 - \beta_2^t) \quad (\text{bias correction})$$

$$w_{t+1} = w_t - \alpha \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$$

Learning Rate Decay:

- Decrease LR over time (step, cosine, linear, etc.)
- Linear warmup: increase LR from 0 over first few steps, prevent exploding loss

Key concepts:

- SGD issues: poor conditioning, getting stuck in local minima/saddle points, noisy
- Momentum overcomes oscillations and escape poor local minima, continues moving in prev direction
- RMSProp adds per-parameter learning rate, addresses AdaGrad's decaying learning rate issue
- Adam combines momentum and RMSProp
- AdamW separates weight decay from gradient update for better regularization
- Second-order methods: $\theta^* = \theta_0 - \alpha H^{-1} \nabla L(\theta_0)$
Better updates, $O(N^2)$ mem and $O(N^3)$ time to invert

Neural Networks

MLP: $f = W_2 \max(0, W_1 x + b_1) + b_2$

$x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H \times D}, b_1 \in \mathbb{R}^H, W_2 \in \mathbb{R}^{C \times H}, b_2 \in \mathbb{R}^C$

Activation Functions:

- ReLU: $f(x) = \max(0, x)$
- Leaky ReLU: $f(x) = \max(\alpha x, x)$ with small α

$$\text{ELU: } f(x) = \begin{cases} x & \text{if } x \geq 0 \\ \alpha(e^x - 1) & \text{if } x < 0 \end{cases}$$

$$\text{GELU: } f(x) = x \cdot \Phi(x)$$

$$\text{Sigmoid: } \sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\text{Tanh: } \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

Key concepts:

- ReLU has zero grad with negative inputs (dying ReLU)
- Leaky ReLU (and variants) always has non-zero slope
- Tanh has zero-centered outputs (sigmoid has $\mu = 0.5$)
- Without nonlinear activations, deep networks reduce to linear models
- Deeper networks can represent more complex functions with fewer parameters (more non-linearities)

Backpropagation

Gradient flow: Upstream \times Local = Downstream

Vector derivatives:

- $d_x f(x)$ has same shape as x
- Apply chain rule using matrix calculus
- Matmul: $\frac{\partial}{\partial X}(XW) = W^T$ and $\frac{\partial}{\partial W}(XW) = X^T$
- Each element $X_{n,d}$ affects the whole row Y_n
- Backprop: X: [N, D], W: [D, M], Y: [N, M]

$$\frac{\partial L}{\partial X} = \left(\frac{\partial L}{\partial Y} \right)^T W^T \in [N, D] \quad \frac{\partial L}{\partial W} = X^T \left(\frac{\partial L}{\partial Y} \right) \in [D, M]$$

Special derivatives:

- Sigmoid: $d_x \sigma(x) = \sigma(x)(1 - \sigma(x))$
- Tanh: $d_x \tanh(x) = 1 - \tanh^2(x)$
- ReLU: $d_x \text{ReLU}(x) = 1(x > 0)$
- Max: $d_x \max(x, y) = 1(x > y)$
- Softmax: $\frac{dp_i}{ds_j} = p_i(1(i = j) - p_j)$ where $p = \frac{e^s}{\sum_k e^{sk}}$
- Cross-entropy: $d_{s_i} \left(-\sum_j y_j \log(p_j) \right) = p_i - y_i$
- Huber Loss: $d_x L_\delta(x, y) = \begin{cases} x - y & \text{if } |x - y| < \delta \\ \delta \cdot \text{sign}(x - y) & \text{otherwise} \end{cases}$
- L1 Loss: $d_x |x - y| = \text{sign}(x - y)$

Key Backpropagation Concepts:

- Vanishing gradients: gradients become too small in deep networks (esp. with sigmoid/tanh)
- Exploding gradients: gradients become too large (common in RNNs)
- Gradient clipping: Cap gradient magnitude to prevent explosion

Convolutional Neural Networks

Conv Layer Summary:

- Hyperparameters:
 - Kernel size: $K_H \times K_W$
 - Number filters: C_{out}
 - Padding: $P = (K - 1)/2$ (same padding)
 - Stride: S
- Weight matrix: $C_{out} \times C_{in} \times K_H \times K_W$
- Bias vector: C_{out}
- Input: $C_{in} \times H \times W$
- Output activation: $C_{out} \times H' \times W'$ where

$$[H', W'] = \frac{[H, W] - K_{[H, W]} + 2P}{S} + 1$$

Advantages:

- Parameter sharing: Same filter applied across image
- Sparse connectivity: Each output depends on small local region
- Translation equivariance: Shifting input shifts output

Pooling layers:

- Given input $C \times H \times W$, downsample each $1 \times H \times W$ plane
- Max pooling: Take maximum value in window
- Average pooling: Average values in window
- Reduces spatial dimensions, increases receptive field

Receptive field: Region of input that affects output

- For $K \times K$ filters, RF grows by $(K - 1)$ per layer
- With L layers and $S = 1$, RF is $1 + L * (K - 1)$
- In general, $R_0 = 1$ and $R_l = R_{l-1} + (K - 1) \times S$

Key concepts:

- Multiple 3×3 filters better than single large filter: fewer parameters, more nonlinearities
- 1×1 conv: same H/W , dim reduction across channels
- Dilated convolutions: expand receptive field without increasing parameters

Normalization Techniques

Batch Normalization:

μ_c = mean of feature values across batch for channel c

σ_c = standard deviation across batch for channel c

$$y = \gamma_c \cdot \frac{x - \mu_c}{\sigma_c} + \beta_c \quad (\text{scale and shift})$$

Normalizes across batch dimension for each channel

Layer Normalization:

μ_n = mean across all channels for sample n

σ_n = standard deviation across all channels for sample n

$$y = \gamma_c \cdot \frac{x - \mu_{n,c}}{\sigma_n} + \beta_c \quad (\text{scale and shift})$$

Normalizes across channel dimension for each sample

Instance Normalization:

$\mu_{n,c}$ = spatial mean for channel c of sample n

$\sigma_{n,c}$ = spatial standard deviation for channel c of sample n

$$y = \gamma_c \cdot \frac{x - \mu_{n,c}}{\sigma_{n,c}} + \beta_c \quad (\text{scale and shift})$$

Normalizes each channel independently for each sample

Key concepts:

- BatchNorm: Used in CNNs, must track running stats for inference
- LayerNorm: Used in transformers, no dependence on batch statistics
- Instance Norm: Used in style transfer, normalizes each feature map
- Group Norm: Compromise between Layer and Instance Norm

CNN Architectures

VGG:

- Multiple 3×3 convs followed by max-pooling

- Multiple 3×3 filters have same receptive field as larger filter with fewer parameters
 - Uniform design: doubles channels after each pooling
- ResNet:
- Skip connections: output = $F(x) + x$
 - Allow deeper networks by learning residual mapping
 - Solves vanishing gradient problem in deep nets
 - Typical block: Conv \rightarrow BN \rightarrow ReLU \rightarrow Conv \rightarrow BN \rightarrow Add \rightarrow ReLU

Bottleneck layer:

- 1×1 conv to reduce channels, 3×3 conv, 1×1 conv to expand
- Reduces computation by decreasing channels in 3×3 conv

Key concepts:

- Network design trade-offs: depth vs. width vs. resolution
- Skip connections help gradient flow and enable deeper networks
- Deeper networks generally need more regularization
- ResNet skip connections: $O_l = I_l + F(I_l)$ not $O_l = I_l * F(I_l)$

Weight Initialization

Xavier/Glorot: $W \sim \mathcal{N}(0, \sqrt{\frac{2}{n_{in} + n_{out}}})$

- For tanh/sigmoid activations
- Maintains variance across linear layers

Kaiming initialization: $W \sim \mathcal{N}(0, \sqrt{\frac{2}{D_{in}}})$ for ReLU

- For ReLU activations (accounts for half being zeroed)
- For CNN: D_{in} = kernel.size² \times input_channels

Key concepts:

- Poor initialization can cause vanishing/exploding gradients
- Initialization in deep nets is crucial for trainability
- Even with good normalization, bad initialization slows training
- Initialization should match the activation function

Training Techniques

Data Augmentation:

- Increases dataset size/diversity without new data
- Robustness to variations (position, lighting)
- Common techniques: flips, crops, color jitter, rotations

Transfer Learning:

- Small dataset: Freeze pretrained model, retrain final layers
- Medium dataset: Freeze early layers, fine-tune later layers
- Large dataset: Initialize with pretrained weights, fine-tune all layers

Diagnostics:

- Underfitting: Low train/val accuracy, small gap
 - Overfitting: High train accuracy, low val accuracy, large gap
 - Not training enough: Low train/val accuracy with gap
- Hyperparameter selection:
- Random search usually better than grid search
 - Check initial loss, overfit small sample first
 - Find LR that makes loss decrease within 100 iterations

Loss Functions

Cross-entropy: $L = -\sum_i y_i \log(\hat{y}_i)$

- For classification problems
- Measures dissimilarity between two probability distributions

$$\text{KL Divergence: } D_{KL}(p||q) = \sum_i p_i \log \frac{p_i}{q_i} = \sum_i p_i (\log p_i - \log q_i)$$

- Rewrite as subtraction for numerical stability
- Not symmetric: $D_{KL}(p||q) \neq D_{KL}(q||p)$

Smooth L1/Huber Loss:

$$L_\delta(x, y) = \begin{cases} \frac{1}{2}(x - y)^2 & \text{if } |x - y| < \delta \\ \delta(|x - y| - \frac{1}{2}\delta) & \text{otherwise} \end{cases}$$

- Combines MSE (near zero) and L1 (for outliers)
 - Differentiable everywhere, robust to outliers
- Triplet margin loss: $L(a, p, n) = \max\{d(a, p) - d(a, n) + \text{margin}, 0\}$
- Used in contrastive learning
 - Pushes anchor (a) closer to positive (p) than negative (n)
 - Margin controls separation between positive and negative pairs

Recurrent Neural Networks

Vanilla RNN:

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$$

$$y_t = W_{hy}h_t + b_y$$

LSTM:

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

$$C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t$$

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t \odot \tanh(C_t)$$

RNN Applications:

- Language modeling: Predict next token in sequence
 - Image captioning: CNN feature extractor + RNN decoder
 - Sequence-to-sequence: Translation, summarization
- Training RNNs:**
- Backpropagation through time (BPTT)
 - Truncated BPTT for long sequences
 - Gradient clipping to prevent explosion

- Key concepts:**
- RNNs can process variable-length sequences
 - Vanishing/exploding gradients limit learning long-term dependencies
 - LSTM addresses vanishing gradients via cell state pathway
 - RNNs sequential processing limits parallelization

Attention Mechanism

Inputs:

- Input sequence: $X \in \mathbb{R}^{N \times D_{in}}$
- Query matrix: $W_Q \in \mathbb{R}^{D_{in} \times D_{out}}$
- Key matrix: $W_K \in \mathbb{R}^{D_{in} \times D_{out}}$
- Value matrix: $W_V \in \mathbb{R}^{D_{in} \times D_{out}}$

Computation:

$$\text{Queries: } Q = XW_Q \in \mathbb{R}^{N \times D_{out}}$$

$$\text{Keys: } K = XW_K \in \mathbb{R}^{N \times D_{out}}$$

$$\text{Values: } V = XW_V \in \mathbb{R}^{N \times D_{out}}$$

$$\text{Similarities: } E = \frac{QK^T}{\sqrt{D_{out}}} \in \mathbb{R}^{N \times N}$$

$$\text{Weights: } A = \text{softmax}(E, \text{dim}=1) \in \mathbb{R}^{N \times N}$$

$$\text{Output: } Y = AV \in \mathbb{R}^{N \times D_{out}}$$

Multi-Head Attention:

Inputs:

- Input vectors: $X \in \mathbb{R}^{N \times D}$
- For each head $h \in \{1, \dots, H\}$:
 - Key matrix: $W_K^h \in \mathbb{R}^{D \times D_H}$
 - Value matrix: $W_V^h \in \mathbb{R}^{D \times D_H}$
 - Query matrix: $W_Q^h \in \mathbb{R}^{D \times D_H}$
- Output matrix: $W_O \in \mathbb{R}^{H D_H \times D}$

Computation for each head h :

$$\text{Queries: } Q^h = XW_Q^h \in \mathbb{R}^{N \times D_H}$$

$$\text{Keys: } K^h = XW_K^h \in \mathbb{R}^{N \times D_H}$$

$$\text{Values: } V^h = XW_V^h \in \mathbb{R}^{N \times D_H}$$

$$\text{Similarities: } E^h = \frac{Q^h(K^h)^T}{\sqrt{D_H}} \in \mathbb{R}^{N \times N}$$

$$\text{Attention weights: } A^h = \text{softmax}(E^h, \text{dim}=1) \in \mathbb{R}^{N \times N}$$

$$\text{Head output: } Y^h = A^h V^h \in \mathbb{R}^{N \times D_H}$$

Combining heads:

$$\text{Concatenated output: } Y = [Y^1; Y^2; \dots; Y^H] \in \mathbb{R}^{N \times H D_H}$$

$$\text{Final output: } O = YW_O \in \mathbb{R}^{N \times D}$$

Types of Attention:

- Self-attention: Q, K, V from same sequence
- Cross-attention: Q from one sequence, K, V from another
- Masked attention: Future positions masked (decoder)

Key concepts:

- Time complexity: $O(n^2d)$ for sequence length n and dimension d
- Memory complexity: $O(n^2)$ for attention weights
- Attention weights computed from Q and K (not V)
- Scaling factor $\sqrt{d_k}$ prevents vanishing gradients with large dimensions
- Self-attention is permutation equivariant without posi-

tional encoding

Transformers

Transformer block:

1. Layer normalization
2. Multi-head self-attention
3. Residual connection
4. Layer normalization
5. Feed-forward network (MLP)
6. Residual connection

Parameters in transformer block:

- Self-attention: $4d^2$ (Q, K, V projections + output)
- Feed-forward: $2df$ (where f is FF dimension, typically $4d$)

Vision Transformer (ViT):

- Split image into patches (16×16)
- Linear projection + position embeddings
- Standard transformer encoder architecture
- CLS token or pooling for classification

Key concepts:

- Transformers use LayerNorm, NOT BatchNorm
- Pre-norm vs. post-norm: affects training stability
- Transformers parallelize better than RNNs for sequences
- Positional encodings enable model to learn position information

Semantic Segmentation

Task: Classify each pixel in an image

Architectures:

- Fully Convolutional Networks (FCN)
- U-Net: Encoder-decoder with skip connections
- DeepLab: Atrous convolutions for dense predictions

Upsampling techniques:

- Unpooling: Reverse pooling operation
- Transposed convolution: Learnable upsampling
- Bilinear interpolation + 1×1 convs: Smoother results

Key concepts:

- Semantic segmentation: One label per pixel, no instance separation
- Downsampling followed by upsampling preserves context while maintaining resolution
- Skip connections help preserve spatial detail
- Dilated/atrous convolutions expand receptive field without losing resolution

Object Detection

Key architectures:

- R-CNN family: Region proposals + classification
- YOLO: Single-pass detection with grid cells
- DETR: Transformers with object queries

Region Proposal Network:

- Generate candidate boxes
- Binary classification (object vs. background)
- Bounding box regression

Evaluation metrics:

- IoU (Intersection over Union): $\frac{\text{area of intersection}}{\text{area of union}}$
- Precision: $\frac{TP}{TP + FP}$, Recall: $\frac{TP}{TP + FN}$
- AP: Area under PR curve for each class
- mAP: Mean AP across all classes

Key concepts:

- Two-stage detectors (R-CNN family): region proposal + classification
- One-stage detectors (YOLO, SSD): directly predict boxes from grid cells
- NMS (Non-Maximum Suppression): Remove duplicate detections
- Anchor boxes: Pre-defined box shapes to match during training

Instance Segmentation

Mask R-CNN:

- Extends Faster R-CNN with mask branch
- RoIAlign for accurate feature extraction
- Parallel heads for classification, box regression, mask prediction

Key concepts:

- RoIAlign: Keeps spatial information intact (avoids quantization)
- Instance segmentation separates individual instances of same class
- Panoptic segmentation: Combines semantic and instance segmentation

Video Understanding

Architectures:

- Single-frame CNN + temporal pooling
- Early fusion: Treat time as channels
- 3D CNN: 3D convolutions (C3D, I3D)
- CNN + RNN: CNN features fed to RNN
- Transformer: Space-time attention

3D convolution:

$$\text{Output: } F \times T' \times H' \times W'$$

$$\text{Filter size: } C \times k_t \times k \times k$$

Two-stream networks:

- Spatial stream: RGB frames
- Temporal stream: Optical flow
- Late fusion of predictions

Key concepts:

- 3D CNN receptive fields span space and time dimensions
- Early fusion builds temporal receptive field all at once
- Slow fusion gradually builds temporal receptive field
- 3D CNNs have temporal-shift invariance (early fusion doesn't)

Neural Network Visualization

Saliency maps:

- Compute gradient of class score w.r.t input pixels
- Highlights regions important for classification

Class Activation Mapping (CAM):

$$M_c(x, y) = \sum_k w_k^c \cdot f_k(x, y)$$

Grad-CAM:

- Generalizes CAM to any CNN architecture
- Global-average-pools gradients for importance weights
- Weighted combination of feature maps

Key concepts:

- Visualizations help debug network decisions
- CNN filters often detect edges, textures, patterns, and semantic concepts
- Attention maps in transformers provide built-in visualization

Evaluation Metrics

Classification:

- Accuracy: $\frac{\text{correct predictions}}{\text{total predictions}}$
- Precision: $\frac{TP}{TP + FP}$
- Recall: $\frac{TP}{TP + FN}$
- F1 Score: $\frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$

Segmentation:

- Pixel accuracy: $\frac{\text{correctly classified pixels}}{\text{total pixels}}$
- Mean IoU: Average IoU across all classes
- Dice coefficient: $\frac{2 \times \text{intersection}}{\text{sum of areas}}$

Point Cloud Processing:

- Translation equivariance: Output shifts when input shifts
- Rotation equivariance: Output rotates when input rotates
- Convs on grid structure not rotation-equivariant
- Continuous point convs use weight functions of relative positions