Neural Networks Class 6

CNN Architectures & Advanced Topics

Learning Objectives:

- Master complete CNN architectures and mathematical flow
- Understand famous CNN models and their innovations
- Learn advanced CNN concepts and modern developments

Duration: 50 minutes

Complete CNN Architecture

Conv → ReLU → Pool → FC Pipeline

Typical CNN Structure:

$$\operatorname{Input} \to [\operatorname{Conv} \to \operatorname{ReLU} \to \operatorname{Pool}]^n \to [\operatorname{FC} \to \operatorname{ReLU}]^m \to \operatorname{FC} \to \operatorname{Softmax}$$

Mathematical Flow Example (32×32×3 → 10 classes):

- 1. Conv1: $32 \times 32 \times 3 \rightarrow 32 \times 32 \times 32$ (5×5 filters, pad=2)
- 2. Pool1: $32 \times 32 \times 32 \rightarrow 16 \times 16 \times 32$ (2×2 max pool)
- 3. Conv2: $16 \times 16 \times 32 \rightarrow 16 \times 16 \times 64$ (5×5 filters, pad=2)
- 4. Pool2: $16 \times 16 \times 64 \rightarrow 8 \times 8 \times 64$ (2×2 max pool)
- 5. FC1: $8 \times 8 \times 64 = 4096 \rightarrow 128$

Mathematical Flow Through CNN

Complete Forward Pass

Layer-by-layer Mathematical Description:

Convolutional Layer *l*:

$$\mathbf{Z}_{i,j,k}^{(l)} = \sum_{c=1}^{C^{(l-1)}} \sum_{u=1}^F \sum_{v=1}^F \mathbf{W}_{u,v,c,k}^{(l)} \mathbf{A}_{i+u-1,j+v-1,c}^{(l-1)} + b_k^{(l)}$$

Activation:

$$\mathbf{A}_{i,j,k}^{(l)} = \mathrm{ReLU}(\mathbf{Z}_{i,j,k}^{(l)})$$

Mathematical Flow Through CNN

Complete Forward Pass

Pooling Layer:

$$\mathbf{P}_{i,j,k}^{(l)} = \max_{u,v \in \mathrm{pool}} \mathbf{A}_{2i+u,2j+v,k}^{(l)}$$

Flattening for FC:

$$\mathbf{h} = ext{flatten}(\mathbf{P}^{(L)}) \in \mathbb{R}^{H imes W imes C}$$

Receptive Field Through Network

Growing Perception

Receptive Field Growth:

Each layer sees increasingly larger input regions

Calculation Through Network:

- Conv1 (3×3): $RF = 3 \times 3$
- Pool1 (2×2, s=2): $RF = 4 \times 4$
- Conv2 (3×3): RF = 8×8
- Pool2 (2×2, s=2): RF = 10×10
- Conv3 (3×3): RF = 18×18

Receptive Field Through Network

Growing Perception

General Formula:

$$RF_l = RF_{l-1} + (K_l - 1) imes \prod_{i=1}^{l-1} S_i$$

Insight: Deep networks see large input regions while maintaining parameter efficiency

Parameter Counting Exercise

Complete Network Analysis

Example Network Architecture:

```
Input: 224\times224\times3 Conv1: 64 filters, 7\times7, stride=2, pad=3 \rightarrow 112\times112\times64 Pool1: 2\times2, stride=2 \rightarrow 56\times56\times64 Conv2: 128 filters, 3\times3, stride=1, pad=1 \rightarrow 56\times56\times128 Pool2: 2\times2, stride=2 \rightarrow 28\times28\times128 FC1: 1024 neurons FC2: 10 classes (output)
```

Parameter Counting Exercise

Complete Network Analysis

Parameter Calculations:

- Conv1: $7 \times 7 \times 3 \times 64 + 64 = 9,472$
- Conv2: $3 \times 3 \times 64 \times 128 + 128 = 73,856$
- FC1: $28 \times 28 \times 128 \times 1024 + 1024 = 103, 219, 200$
- FC2: $1024 \times 10 + 10 = 10,250$

Total: ~103.3M parameters (FC layers dominate!)

LeNet-5 Architecture

The Pioneer (1998)

Yann LeCun's LeNet-5 for MNIST:

Architecture:

$$32 \times 32 \times 1 \rightarrow C1 \rightarrow S2 \rightarrow C3 \rightarrow S4 \rightarrow C5 \rightarrow F6 \rightarrow \text{Output}$$

LeNet-5 Architecture

The Pioneer (1998)

Layer Details:

- C1: 6 filters, $5\times5 \rightarrow 28\times28\times6$ (156 params)
- **S2**: 2×2 avg pool → 14×14×6
- C3: 16 filters, $5 \times 5 \rightarrow 10 \times 10 \times 16$ (2,416 params)
- **S4**: 2×2 avg pool $\rightarrow 5\times 5\times 16$
- C5: 120 filters, $5 \times 5 \rightarrow 1 \times 1 \times 120$ (48,120 params)
- **F6:** FC layer → 84 neurons (10,164 params)
- Output: 10 classes (850 params)

AlexNet Architecture

The ImageNet Breakthrough (2012)

Key Innovations:

- 1. ReLU activation instead of tanh/sigmoid
- 2. **Dropout regularization** in FC layers
- 3. Data augmentation and multiple crops
- 4. GPU implementation for parallel training

AlexNet Architecture

The ImageNet Breakthrough (2012)

Architecture Highlights:

- Input: 224×224×3 (ImageNet scale)
- **Conv1:** 96 filters, 11×11, stride=4 → 55×55×96
- Conv2: 256 filters, $5\times5 \rightarrow 27\times27\times256$
- Conv3-5: Multiple 3×3 convolutions
- **FC layers:** 4096 → 4096 → 1000 neurons
- Total: ~60M parameters

Mathematical Impact: Showed CNNs could scale to large datasets

AlexNet Mathematical Improvements

Why It Dominated

ReLU vs Sigmoid/Tanh:

$$\operatorname{ReLU}(x) = \max(0,x) ext{ vs. } anh(x) = rac{e^x - e^{-x}}{e^x + e^{-x}}$$

Advantages:

- No saturation for positive inputs
- Sparse activation (many zeros)
- Fast computation (simple threshold)
- Better gradient flow (derivative is 0 or 1)

AlexNet Mathematical Improvements

Why It Dominated

Dropout Mathematics:

$$ilde{\mathbf{h}} = \mathbf{r} \odot \mathbf{h} ext{ where } \mathbf{r} \sim \operatorname{Bernoulli}(0.5)$$

Regularization Effect: Prevents co-adaptation of neurons

VGG Architecture

Deeper with Smaller Filters (2014)

Key Innovation: Small filters (3×3) instead of large ones

Mathematical Advantage:

Two 3×3 convolutions = One 5×5 convolution in receptive field

- ullet Parameters: $2 imes (3^2 imes C^2) = 18C^2 ext{ vs } 5^2 imes C^2 = 25C^2$
- Non-linearity: 2 ReLU activations vs 1
- 28% parameter reduction with more expressiveness

VGG Architecture

Deeper with Smaller Filters (2014)

VGG-16 Structure:

```
Conv3-64 × 2 → Pool → Conv3-128 × 2 → Pool → Conv3-256 × 3 → Pool Conv3-512 × 3 → Pool → Conv3-512 × 3 → Pool → FC-4096 × 2 → FC-1000
```

Total: ~138M parameters (mostly in FC layers)

ResNet Innovation

Solving the Degradation Problem

Problem: Deeper networks perform worse than shallow ones

- Not due to overfitting (training error also increases)
- Degradation problem: Optimization difficulty

Solution: Residual connections

$$\mathbf{H}(\mathbf{x}) = \mathbf{F}(\mathbf{x}) + \mathbf{x}$$

ResNet Innovation

Solving the Degradation Problem

Mathematical Insight:

Instead of learning $\mathbf{H}(\mathbf{x})$, learn residual $\mathbf{F}(\mathbf{x}) = \mathbf{H}(\mathbf{x}) - \mathbf{x}$

Identity Mapping:

If optimal function is identity, easier to learn ${f F}({f x})=0$ than ${f H}({f x})={f x}$

ResNet Mathematical Analysis

Gradient Flow Through Skip Connections

Forward Pass:

$$\mathbf{x}_{l+1} = \mathbf{x}_l + \mathcal{F}(\mathbf{x}_l, \mathcal{W}_l)$$

Backward Pass:

$$rac{\partial \mathcal{E}}{\partial \mathbf{x}_l} = rac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}} igg(1 + rac{\partial \mathcal{F}(\mathbf{x}_l, \mathcal{W}_l)}{\partial \mathbf{x}_l} igg)$$

Key Insight: The "+1" ensures gradient flow even if $\frac{\partial \mathcal{F}}{\partial \mathbf{x}_l} \to 0$

ResNet Mathematical Analysis

Gradient Flow Through Skip Connections

Multi-hop Gradient:

$$rac{\partial \mathcal{E}}{\partial \mathbf{x}_l} = rac{\partial \mathcal{E}}{\partial \mathbf{x}_L} \Biggl(1 + \sum_{i=l}^{L-1} rac{\partial \mathcal{F}_i}{\partial \mathbf{x}_l} \Biggr)$$

Result: Enables training of 152+ layer networks

Transfer Learning Mathematics

Feature Reuse and Fine-tuning

Pre-trained Feature Representation:

$$\mathbf{f}_{\mathrm{pre}} = CNN_{\mathrm{pre}}(\mathbf{x};oldsymbol{ heta}_{\mathrm{pre}})$$

Transfer Learning Strategies:

1. Feature Extraction:

$$\mathbf{y} = ext{Classifier}(\mathbf{f}_{ ext{pre}}; oldsymbol{ heta}_{ ext{new}})$$

Only train $oldsymbol{ heta}_{
m new}$, freeze $oldsymbol{ heta}_{
m pre}$

Transfer Learning Mathematics

Feature Reuse and Fine-tuning

2. Fine-tuning:

$$\mathbf{y} = CNN_{\mathrm{full}}(\mathbf{x}; [oldsymbol{ heta}_{\mathrm{pre}}, oldsymbol{ heta}_{\mathrm{new}}])$$

Train all parameters with smaller learning rate for $oldsymbol{ heta}_{
m pre}$

Mathematical Justification: Lower layers learn general features (edges, textures), higher layers learn task-specific features

Data Augmentation for Computer Vision

Mathematical Transformations

Geometric Transformations:

Rotation:

$$\mathbf{R}_{ heta} = egin{bmatrix} \cos heta & -\sin heta \ \sin heta & \cos heta \end{bmatrix}$$

Translation:

$$\mathbf{T}_{t_x,t_y}(\mathbf{x}) = \mathbf{x} + egin{bmatrix} t_x \ t_y \end{bmatrix}$$

Data Augmentation for Computer Vision

Mathematical Transformations

Scaling:

$$\mathbf{S}_s(\mathbf{x}) = s \cdot \mathbf{x}$$

Horizontal Flip:

$$\operatorname{Flip}(I[i,j]) = I[i,W-1-j]$$

Mathematical Effect: Increases effective dataset size by factor of augmentations used

Batch Normalization in CNNs

Channel-wise Normalization

For Convolutional Layers:

Normalize across batch and spatial dimensions, per channel

Statistics Computation:

$$\mu_c = rac{1}{N \cdot H \cdot W} \sum_{n=1}^N \sum_{h=1}^H \sum_{w=1}^W x_{n,h,w,c}$$

$$\sigma_c^2 = rac{1}{N \cdot H \cdot W} \sum_{n=1}^{N} \sum_{h=1}^{H} \sum_{w=1}^{W} (x_{n,h,w,c} - \mu_c)^2$$

Batch Normalization in CNNs

Channel-wise Normalization

Normalization:

$$\hat{x}_{n,h,w,c} = rac{x_{n,h,w,c} - \mu_c}{\sqrt{\sigma_c^2 + \epsilon}}$$

Scale and Shift (per channel):

$$y_{n,h,w,c} = \gamma_c \hat{x}_{n,h,w,c} + eta_c$$

Modern CNN Innovations

Depthwise Separable Convolutions

Standard Convolution:

$$ext{Params} = K imes K imes C_{ ext{in}} imes C_{ ext{out}}$$
 $ext{FLOPs} = K imes K imes C_{ ext{in}} imes C_{ ext{out}} imes H imes W$

Depthwise Separable:

1. Depthwise: Each input channel convolved separately

Params =
$$K \times K \times C_{\rm in}$$

2. Pointwise: 1×1 convolution across channels

$${
m Params} = 1 imes 1 imes C_{
m in} imes C_{
m out}$$

Modern CNN Innovations

Depthwise Separable Convolutions

Total Reduction:

$$rac{K imes K imes C_{
m in} + C_{
m in} imes C_{
m out}}{K imes K imes C_{
m in} imes C_{
m out}} = rac{1}{C_{
m out}} + rac{1}{K^2}$$

Attention in Computer Vision

Spatial Attention Mechanisms

Spatial Attention Map:

 $\mathbf{A} = \operatorname{sigmoid}(\operatorname{Conv}(\operatorname{AdaptiveAvgPool}(\mathbf{F}) + \operatorname{AdaptiveMaxPool}(\mathbf{F})))$

Channel Attention:

$$\mathbf{M}_c = \sigma(\mathbf{W}_2(\mathrm{ReLU}(\mathbf{W}_1(\mathrm{GAP}(\mathbf{F}_c)))))$$

Attention in Computer Vision

Spatial Attention Mechanisms

Self-Attention for Images:

$$\operatorname{Attention}(\mathbf{Q},\mathbf{K},\mathbf{V}) = \operatorname{softmax}\left(rac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}
ight)\mathbf{V}$$

Where $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ are linear projections of spatial features

Vision Transformer (ViT): Pure attention-based architecture for images

Performance Metrics

Beyond Accuracy

Classification Metrics:

Top-1 Accuracy:

$$Acc_1 = \frac{correct\ predictions}{total\ predictions}$$

Top-5 Accuracy:

$$Acc_5 = \frac{\text{samples where true label in top 5 predictions}}{\text{total predictions}}$$

Performance Metrics

Beyond Accuracy

Precision and Recall (per class):

$$ext{Precision}_c = rac{TP_c}{TP_c + FP_c}, \quad ext{Recall}_c = rac{TP_c}{TP_c + FN_c}$$

F1-Score:

$$F1_c = rac{2 imes ext{Precision}_c imes ext{Recall}_c}{ ext{Precision}_c + ext{Recall}_c}$$

Model Efficiency Metrics

Speed and Memory Considerations

FLOPs (Floating Point Operations):

For convolution: $\mathrm{FLOPs} = K^2 imes C_{\mathrm{in}} imes H_{\mathrm{out}} imes W_{\mathrm{out}} imes C_{\mathrm{out}}$

Parameters vs FLOPs:

- Parameters: Storage/memory requirement
- FLOPs: Computational requirement

Model Efficiency Metrics

Speed and Memory Considerations

Inference Time:

- Latency: Time for single prediction
- Throughput: Predictions per second

Memory Usage:

- Model size: Parameter storage
- Activation memory: Intermediate feature maps
- Peak memory: Maximum during forward/backward pass

CNN Design Principles

Architecture Guidelines

Depth vs Width Trade-offs:

- Deeper networks: More representational power, harder to train
- Wider networks: More parameters per layer, easier to train

Filter Size Evolution:

- Early layers: Larger filters (7×7, 5×5) for low-level features
- Later layers: Smaller filters (3×3, 1×1) for complex combinations

CNN Design Principles

Architecture Guidelines

Channel Progression:

- Input → Output: Gradually increase channels while decreasing spatial size
- Typical pattern: $3 \rightarrow 64 \rightarrow 128 \rightarrow 256 \rightarrow 512 \rightarrow 1024$

Mathematical Intuition: Preserve information flow while reducing spatial dimensions

Future Directions

Beyond Traditional CNNs

Neural Architecture Search (NAS):

$$lpha^* = rg \min_{lpha} \mathcal{L}_{ ext{val}}(\mathbf{w}^*(lpha), lpha)$$

EfficientNet Scaling:

$$\mathrm{depth} = lpha^{\phi}, \quad \mathrm{width} = eta^{\phi}, \quad \mathrm{resolution} = \gamma^{\phi}$$

Subject to: $\alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$ and $\alpha, \beta, \gamma \geq 1$

Future Directions

Beyond Traditional CNNs

Vision Transformers:

Moving from convolution-based to attention-based architectures

3D CNNs:

Extension to video: (T, H, W, C) with temporal convolutions

Summary: Class 6 Key Concepts

Complete CNN Architectures:

- Full pipeline: Conv → ReLU → Pool → FC with mathematical flow
- Parameter counting: Understanding computational complexity
- Receptive field growth: How networks see larger input regions

Famous Architectures:

- **LeNet-5**: Pioneer for digit recognition
- AlexNet: ImageNet breakthrough with ReLU and dropout
- VGG: Deeper networks with small filters
- **ResNet:** Skip connections solve degradation problem

Summary: Class 6 Key Concepts

Advanced Topics:

- Transfer learning: Feature reuse and fine-tuning mathematics
- Batch normalization: Stabilizing training in CNNs
- Modern innovations: Attention, efficient architectures, NAS

The journey from simple perceptrons to sophisticated CNN architectures represents the power of mathematical foundations combined with computational innovation.