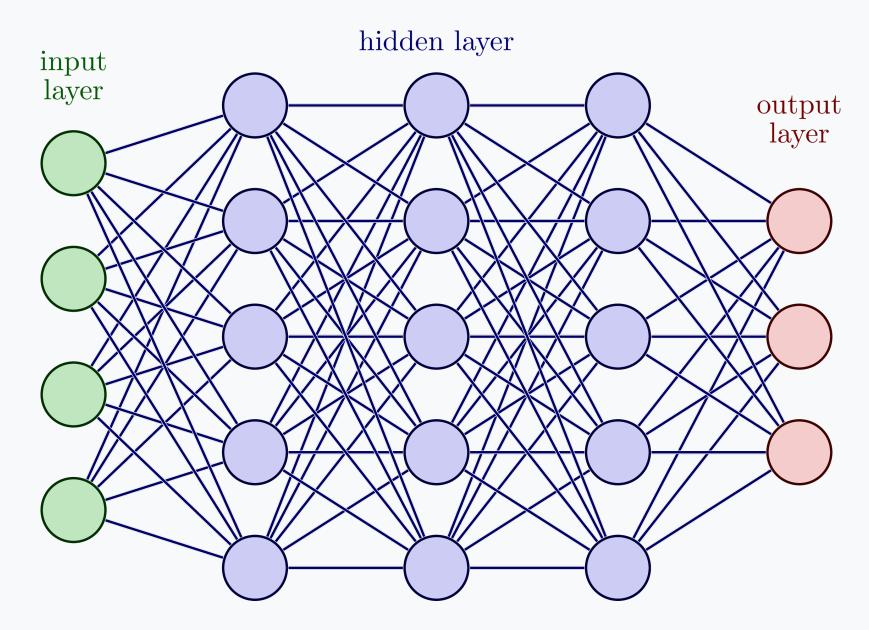
Neural Networks Class 4

Advanced Architectures & Regularization

Learning Objectives:

- Explore deep network architectures and residual connections
- Master advanced regularization techniques mathematically
- Understand various loss functions and their applications



2

Deep Networks - Depth vs Width

Classical Result (Width):

Single hidden layer with sufficient neurons can approximate any continuous function

Modern Understanding (Depth):

- Exponential advantage: Deep networks can represent certain functions exponentially more efficiently
- Hierarchical features: Each layer learns increasingly complex representations
- Computational efficiency: Depth reduces required parameters

Mathematical Insight:

Function with k layers and n neurons per layer: $O(kn^2)$ parameters

Same function with 1 layer: $O(n^k)$ parameters potentially needed

Residual Connections

The Vanishing Gradient Solution

Problem with Very Deep Networks:

$$rac{\partial L}{\partial \mathbf{W}^{(1)}} = rac{\partial L}{\partial \mathbf{A}^{(L)}} \prod_{l=2}^L \mathbf{W}^{(l)} \sigma'(\mathbf{Z}^{(l-1)})$$

As L increases, this product tends to vanish or explode.

Residual Block Mathematical Formulation:

$$\mathbf{A}^{(l+1)} = \sigma(\mathbf{A}^{(l)} + F(\mathbf{A}^{(l)}, \mathbf{W}^{(l)}))$$

Where $F(\mathbf{A}^{(l)}, \mathbf{W}^{(l)})$ is a residual function (typically 2-3 layers)

Advan**Key**r**Innovation:John tity** mapping $\mathbf{A}^{(l)}$ + learned residual $F(\cdot)$

Residual Connections - Gradient Flow Analysis

Forward Pass:

$$\mathbf{A}^{(l+1)} = \mathbf{A}^{(l)} + F(\mathbf{A}^{(l)}, \mathbf{W}^{(l)})$$

Backward Pass:

$$rac{\partial L}{\partial \mathbf{A}^{(l)}} = rac{\partial L}{\partial \mathbf{A}^{(l+1)}} \Bigg(1 + rac{\partial F(\mathbf{A}^{(l)}, \mathbf{W}^{(l)})}{\partial \mathbf{A}^{(l)}} \Bigg)$$

Key Insight: The "+1" term ensures gradient flow even if $\frac{\partial F}{\partial \mathbf{A}^{(l)}} o 0$

Multi-layer Gradient:

$$rac{\partial L}{\partial \mathbf{A}^{(l)}} = rac{\partial L}{\partial \mathbf{A}^{(L)}} \prod_{k=l+1}^{L} \left(1 + rac{\partial F^{(k)}}{\partial \mathbf{A}^{(k-1)}}
ight)$$

Skip Connections - Dense and Highway Networks

Dense Connections (DenseNet):

$$\mathbf{A}^{(l)} = \sigma([\mathbf{A}^{(0)}, \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(l-1)}]\mathbf{W}^{(l)} + \mathbf{b}^{(l)})$$

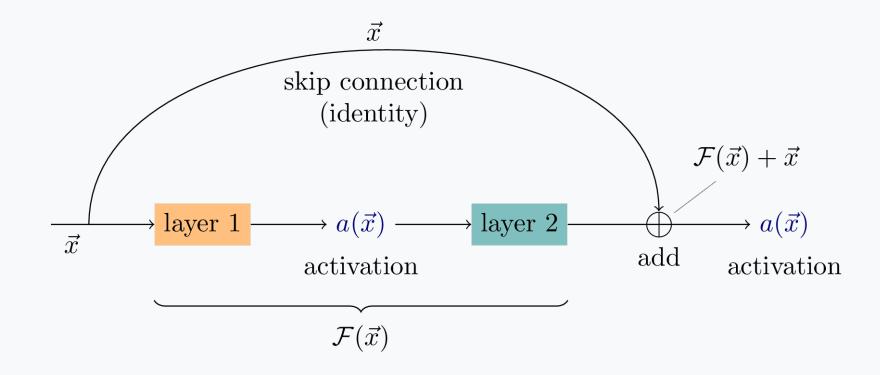
Where $[\cdot]$ denotes concatenation.

Highway Networks:

$$\mathbf{A}^{(l+1)} = \mathbf{T}(\mathbf{A}^{(l)}) \odot \sigma(\mathbf{W}_H^{(l)} \mathbf{A}^{(l)} + \mathbf{b}_H^{(l)}) + (1 - \mathbf{T}(\mathbf{A}^{(l)})) \odot \mathbf{A}^{(l)}$$

Where:

- $\mathbf{T}(\mathbf{A}^{(l)}) = \sigma(\mathbf{W}_T^{(l)}\mathbf{A}^{(l)} + \mathbf{b}_T^{(l)})$: Transform gate
- $1 \mathbf{T}(\mathbf{A}^{(l)})$: Carry gate



Dropout Deep Dive - Mathematical Analysis

Training Phase: For layer l with dropout probability p:

$$\mathbf{r}^{(l)} \sim \mathrm{Bernoulli}(1-p)$$

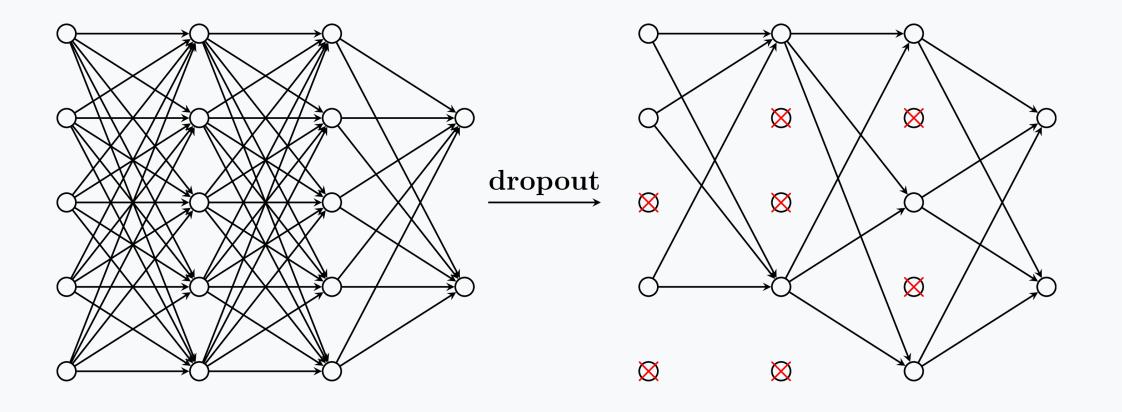
$$ilde{\mathbf{A}}^{(l)} = rac{1}{1-p} \mathbf{r}^{(l)} \odot \mathbf{A}^{(l)}$$

Expected Value Preservation:

$$\mathbb{E}[ilde{A}_i^{(l)}] = rac{1}{1-p} \cdot (1-p) \cdot A_i^{(l)} = A_i^{(l)}$$

Variance:

$$ext{Var}[ilde{A}_i^{(l)}] = rac{p}{1-p}(A_i^{(l)})^2$$



Advanced Architectures & Regularization

Dropout Variants

Spatial and Structured Dropout

Standard Dropout: Independent for each neuron

 $\mathbf{r}_i \sim \mathrm{Bernoulli}(1-p)$ independently

Spatial Dropout (for CNNs): Entire feature maps

 $\mathbf{r}_c \sim \mathrm{Bernoulli}(1-p) ext{ for channel } c$

$$ilde{\mathbf{A}}_{i,j,c} = rac{1}{1-p} \mathbf{r}_c \cdot \mathbf{A}_{i,j,c}$$

DropConnect: Drop connections instead of neurons

$$\mathbf{M}_{i,j} \sim \mathrm{Bernoulli}(1-p)$$

Early Stopping - Mathematical Formulation

Validation Loss Monitoring:

Let $L_{\mathrm{val}}(t)$ be validation loss at epoch t

Stopping Criterion:

Stop training when:

$$L_{\rm val}(t) > L_{\rm val}(t-k)$$
 for k consecutive epochs

Patience Parameter: Number of epochs to wait before stopping

Mathematical Justification:

Minimizes expected test error by finding optimal bias-variance trade-off point

Optimal Stopping Time:

$$t^* = rg\min_t \mathbb{E}[L_{ ext{test}}(t)]$$

Data Augmentation - Mathematical Transformations

Geometric Transformations:

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

- Translation: $\mathbf{x}' = \mathbf{x} + \mathbf{t}$
- Scaling: $\mathbf{x}' = s \cdot \mathbf{x}$

Statistical Transformations:

- Gaussian Noise: $\mathbf{x}' = \mathbf{x} + \mathcal{N}(0, \sigma^2)$
- Normalization: $\mathbf{x}' = \frac{\mathbf{x} \mu}{\sigma}$

Theoretical Foundation:

Augmentation increases effective dataset size and reduces overfitting by making model

Loss Functions for Classification - Cross-Entropy

Binary Cross-Entropy: For probability $p = \sigma(\mathbf{w}^T\mathbf{x} + b)$ and true label $y \in \{0, 1\}$:

Likelihood:

$$P(y|\mathbf{x}) = p^y (1-p)^{1-y}$$

Log-Likelihood:

$$\log P(y|\mathbf{x}) = y \log p + (1-y) \log(1-p)$$

Cross-Entropy Loss:

$$L=-[y\log p+(1-y)\log(1-p)]$$

Derivative:

$$\frac{\partial L}{\partial p} = \frac{p - y}{p(1 - p)}$$

Loss Functions for Classification - Multi-class CE

Softmax Function:

$$p_k = rac{e^{z_k}}{\sum_{j=1}^K e^{z_j}}$$

Properties:

ullet $\sum_{k=1}^K p_k = 1$, $p_k > 0$ for all k, Differentiable everywhere

Multi-class Cross-Entropy:
$$L = -\sum_{k=1}^K y_k \log p_k$$

Derivative:

$$rac{\partial L}{\partial z_k} = p_k - y_k$$

Loss Functions for Regression - Beyond MSE

Mean Squared Error (MSE):

$$L_{ ext{MSE}} = rac{1}{2}(\hat{y}-y)^2$$

• Sensitive to outliers, Derivative: $\hat{y}-y$

Mean Absolute Error (MAE):

$$L_{ ext{MAE}} = |\hat{y} - y|$$

Robust to outliers, Non-differentiable at 0

Huber Loss (Smooth MAE):

$$L_{
m Huber} = egin{cases} rac{1}{2}(\hat{y}-y)^2 & ext{if } |\hat{y}-y| \leq \delta \ \delta(|\hat{y}-y|-rac{1}{2}\delta) & ext{otherwise} \end{cases}$$

Multi-task Learning - Mathematical Framework

Shared Representation: $\mathbf{h} = f_{\mathrm{shared}}(\mathbf{x}; oldsymbol{ heta}_{\mathrm{shared}})$

Task-Specific Heads: $\hat{y}_i = f_i(\mathbf{h}; \boldsymbol{\theta}_i) ext{ for task } i$

Combined Loss: $L_{ ext{total}} = \sum_{i=1}^T \lambda_i L_i(\hat{y}_i, y_i)$

Gradient Update:

$$oldsymbol{ heta_{ ext{shared}} \leftarrow oldsymbol{ heta_{ ext{shared}} - \eta \sum_{i=1}^{T} \lambda_i rac{\partial L_i}{\partial oldsymbol{ heta_{ ext{shared}}}}$$

Challenge: Balancing task weights λ_i

Multi-label Classification - Mathematical Differences

Multi-class: One label per example, $\sum_k p_k = 1$

• Output layer: Softmax, Loss: Categorical cross-entropy

Multi-label: Multiple labels per example, Each p_k independent

Output layer: Multiple sigmoids, Loss: Binary cross-entropy per label

Multi-label Loss:
$$L = -\sum_{k=1}^K [y_k \log p_k + (1-y_k) \log (1-p_k)]$$

Threshold Decision:

$$\hat{y}_k = egin{cases} 1 & ext{if } p_k > au \ 0 & ext{otherwise} \end{cases}$$

Advanced Regularization

Spectral Normalization

Weight Matrix Singular Value Decomposition: $\mathbf{W} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$

Spectral Norm: Largest singular value

$$\sigma(\mathbf{W}) = \max_{\mathbf{h}: \|\mathbf{h}\| = 1} \|\mathbf{W}\mathbf{h}\|$$

Spectral Normalization:

$$\mathbf{W}_{\mathrm{SN}} = rac{\mathbf{W}}{\sigma(\mathbf{W})}$$

Effect: Constrains Lipschitz constant of network

$$\|f(\mathbf{x}_1)-f(\mathbf{x}_2)\|\leq L\|\mathbf{x}_1-\mathbf{x}_2\|$$

Layer Normalization

Batch Normalization: Normalize across batch dimension

$$\hat{x}_i = rac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$$

Layer Normalization: Normalize across feature dimension

$$\hat{x}_j = rac{x_j - \mu_L}{\sqrt{\sigma_L^2 + \epsilon}}$$

Where:

$$\mu_L = rac{1}{H} \sum_{j=1}^{H} x_j, \quad \sigma_L^2 = rac{1}{H} \sum_{j=1}^{H} (x_j - \mu_L)^2$$

Advantages:

Attention Mechanisms Preview

Attention Weights:

$$lpha_{i,j} = rac{\exp(e_{i,j})}{\sum_{k=1}^n \exp(e_{i,k})}$$

Where $e_{i,j} = a(\mathbf{s}_i, \mathbf{h}_j)$ is attention energy

Context Vector: $\mathbf{c}_i = \sum_{j=1}^n lpha_{i,j} \mathbf{h}_j$

Self-Attention:

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

Foundation for Transformers: Modern architecture revolutionizing NLP and beyond

Practical Architecture Design

Choosing Network Depth and Width

Depth Considerations: Shallow networks: Fast training, limited expressiveness, Deep

networks: More expressive, harder to train, need residual connections

Width Considerations: Narrow networks: Few parameters, may underfit, Wide networks:

Many parameters, may overfit, expensive

Rule of Thumb: Start with proven architectures (ResNet, DenseNet) and adapt

Architecture Search:

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

Where
$$\mathbf{w}^*(\alpha) = \arg\min_{\mathbf{w}} \mathcal{L}_{\text{train}}(\mathbf{w}, \alpha)$$

Summary: Class 4 Key Concepts

Deep Network Architectures:

- Residual connections: $\mathbf{A}^{(l+1)} = \mathbf{A}^{(l)} + F(\mathbf{A}^{(l)})$ solve vanishing gradients
- Skip connections: Enable very deep networks, improve gradient flow

Advanced Regularization:

- **Dropout variants:** Spatial, DropConnect with mathematical formulations
- Data augmentation: Mathematical transformations for robustness
- Spectral normalization: Control Lipschitz constant

Loss Functions:

• Classification: Binary/multi-class cross-entropy derivations