

Construindo um Framework Completo de Deep Learning em C

Do Zero ao Producao: Matematica, Algoritmos e
Arquitetura

Guia Master Avancado

12 de fevereiro de 2026

Conteúdo

1 Fundamentos Matematicos e Teoricos do Aprendizado	9
1.1 O Problema Fundamental do Aprendizado	9
1.1.1 Espaco de Hipoteses e Capacidade	9
1.2 Deducao Detalhada da Funcao de Custo	10
1.2.1 Mean Squared Error (MSE)	10
1.2.2 Cross-Entropy para Classificacao	10
1.3 Otimizacao e Gradiente Descendente	11
1.3.1 Deducao do Gradiente	11
1.3.2 Analise de Convergencia	11
1.4 Regra da Cadeia Vetorial e Jacobianos	12
1.4.1 Diferenca entre Gradiente e Jacobiano	12
2 Implementacao Base em C: Estruturas e Gerenciamento de Memoria	13
2.1 Gerenciamento de Memoria para Tensores	13
2.2 Funcoes de Criacao e Destruicao	14
2.3 Operacoes Basicas Otimizadas	15
2.4 Produto Matricial Otimizado	15
2.5 Estrategias de Otimizacao de Memoria	16
2.5.1 Pooling de Alocacao	16
3 Rede Neural: Matematica e Implementacao Detalhada	17

3.1	Arquitetura de uma Rede Neural	17
3.1.1	Neuronio Artificial: Modelo Matematico	17
3.2	Backpropagation: Deducao Passo a Passo	18
3.2.1	Caso Simples: 2 Camadas	18
3.2.2	Implementacao Completa do Backprop	19
3.3	Inicializacao de Pesos	20
3.3.1	Xavier/Glorot Initialization	20
3.3.2	He Initialization	20
4	Sistema de Autograd Completo	23
4.1	Arquitetura do Grafo Computacional	23
4.2	Implementacao de Operacoes com Grafo	24
4.3	Topological Sort e Propagacao Reversa	25
5	Camadas Avancadas e Funcoes de Ativacao	27
5.1	Camada Linear (Fully Connected)	27
5.2	Funcoes de Ativacao com Autograd	29
5.2.1	ReLU e Leaky ReLU	29
5.2.2	Softmax Estavel Numericamente	30
6	Redes Convolucionais: Implementacao Profunda	33
6.1	Fundamentos da Convolucao Discreta	33
6.1.1	Implementacao Ingênuas vs Otimizada	33
6.2	Backpropagation em Convolucao	35
6.3	Otimizacoes: Im2Col e GEMM	37
7	Otimizadores: Teoria e Implementacao	41
7.1	Stochastic Gradient Descent (SGD)	41
7.1.1	Com e sem Momentum	41
7.2	Adam: Adaptive Moment Estimation	43

<i>CONTEÚDO</i>	5
8 Treinamento e Validacao	45
8.1 DataLoader e Batching	45
8.2 Loop de Treinamento Completo	47
9 Serializacao e Persistencia de Modelos	49
9.1 Salvando e Carregando Modelos	49
10 Estudos de Caso e Projetos Praticos	53
10.1 Exemplo 1: MNIST do Zero	53
10.2 Exemplo 2: CNN para CIFAR-10	54
11 Otimizacoes Avancadas e Performance	57
11.1 BLAS-Level Optimizations	57
11.1.1 Tiling para Cache L1/L2	57
11.1.2 SIMD Intrinsics	58
11.2 Parallelism with OpenMP	59
12 Proximos Passos: Alem do Framework	61
12.1 Suporte a GPU com CUDA	61
12.2 Implementacao de Transformers	62
12.3 Mixed Precision Training	62

Prefacio

Este livro nao e um tutorial superficial. E uma jornada profunda atraves de todos os aspectos do deep learning, desde os fundamentos matematicos mais abstratos ate os detalhes de implementacao de baixo nivel em C puro.

A cada capitulo, voce nao apenas aprendera conceitos, mas implementara codigo funcional e testado. Ao final, voce tera construido um framework completo equivalente a um PyTorch simplificado, totalmente funcional e capaz de treinar redes neurais complexas.

Pre-requisitos:

- Conhecimento solido de C (ponteiros, alocacao dinamica, structs)
- Matematica basica de graduacao (calculo, algebra linear)
- Vontade de entender cada detalhe

Quanto tempo leva?

Implementar todo o framework do zero pode levar de 2 a 3 meses de trabalho dedicado. Cada capitulo contem exercicios que consolidam o aprendizado.

Capítulo 1

Fundamentos Matematicos e Teoricos do Aprendizado

1.1 O Problema Fundamental do Aprendizado

O aprendizado de maquina supervisionado pode ser formalizado como um problema de otimizacao. Dado um conjunto de dados $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ onde $\mathbf{x}_i \in \mathbb{R}^d$ e $y_i \in \mathbb{R}^k$, queremos encontrar uma funcao $f : \mathbb{R}^d \rightarrow \mathbb{R}^k$ parametrizada por θ que minimize:

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \sim p_{\text{data}}} [\mathcal{L}(f(\mathbf{x}; \theta), y)]$$

Na pratica, como p_{data} e desconhecida, aproximamos pelo risco empirico:

$$\theta^* \approx \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N \mathcal{L}(f(\mathbf{x}_i; \theta), y_i)$$

1.1.1 Espaco de Hipoteses e Capacidade

Uma rede neural define um espaco de hipoteses \mathcal{H} . O teorema da aproximacao universal (Cybenko, 1989) estabelece que uma rede feedforward com uma camada oculta pode aproximar qualquer funcao

continua em um compacto de \mathbb{R}^d .

Teorema 1.1 (Aproximação Universal). *Seja σ uma função de ativação não constante, limitada e monótona. Então o conjunto de funções da forma:*

$$f(\mathbf{x}) = \sum_{j=1}^m v_j \sigma(\mathbf{w}_j^T \mathbf{x} + b_j)$$

e denso em $C([0, 1]^d)$ no espaço das funções contínuas.

1.2 Dedução Detalhada da Função de Custo

1.2.1 Mean Squared Error (MSE)

Para problemas de regressão, assumimos que a variável alvo é gerada por:

$$y = f(\mathbf{x}; \theta) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

A verossimilhança dos dados é:

$$p(y|\mathbf{x}, \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - f(\mathbf{x}; \theta))^2}{2\sigma^2}\right)$$

Tomando o log negativo:

$$-\log p(y|\mathbf{x}, \theta) = \frac{1}{2} \log(2\pi\sigma^2) + \frac{(y - f(\mathbf{x}; \theta))^2}{2\sigma^2}$$

Ignorando constantes, minimizar MSE é equivalente a maximizar a verossimilhança sob ruido Gaussiano.

1.2.2 Cross-Entropy para Classificação

Para classificação binária, assumimos $y \in \{0, 1\}$ e modelamos:

$$p(y=1|\mathbf{x}) = \sigma(f(\mathbf{x}; \theta)) = \frac{1}{1 + e^{-f(\mathbf{x}; \theta)}}$$

A verossimilhança:

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

Log-verossimilhanca negativa:

$$\mathcal{L} = -y \log(\sigma(f)) - (1 - y) \log(1 - \sigma(f))$$

Para classificacao multiclass com K classes, generalizamos via softmax:

$$p(y = c | \mathbf{x}) = \frac{e^{f_c(\mathbf{x})}}{\sum_{j=1}^K e^{f_j(\mathbf{x})}}$$

$$\mathcal{L} = -\sum_{c=1}^K y_c \log p(y = c | \mathbf{x})$$

1.3 Otimizacao e Gradiente Descendente

1.3.1 Deducao do Gradiente

Queremos minimizar $J(\theta) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(f(\mathbf{x}_i; \theta), y_i)$. A expansao em serie de Taylor de primeira ordem:

$$J(\theta + \Delta\theta) \approx J(\theta) + \nabla_\theta J(\theta)^T \Delta\theta$$

Para garantir decrescimo, escolhemos $\Delta\theta = -\eta \nabla_\theta J(\theta)$:

$$J(\theta + \Delta\theta) \approx J(\theta) - \eta \|\nabla_\theta J(\theta)\|^2 \leq J(\theta)$$

1.3.2 Analise de Convergencia

Sob hipoteses de suavidade de Lipschitz do gradiente:

$$\|\nabla J(\theta_1) - \nabla J(\theta_2)\| \leq L \|\theta_1 - \theta_2\|$$

Podemos provar que com $\eta \leq 1/L$:

$$J(\theta_{t+1}) \leq J(\theta_t) - \frac{\eta}{2} \|\nabla J(\theta_t)\|^2$$

1.4 Regra da Cadeia Vetorial e Jacobianos

O backpropagation é uma aplicação eficiente da regra da cadeia para funções compostas vetoriais. Seja $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ e $g : \mathbb{R}^n \rightarrow \mathbb{R}^p$. Então:

$$\frac{\partial(g \circ f)}{\partial \mathbf{x}} = \frac{\partial g}{\partial \mathbf{f}} \cdot \frac{\partial \mathbf{f}}{\partial \mathbf{x}}$$

Onde $\frac{\partial g}{\partial \mathbf{f}} \in \mathbb{R}^{p \times n}$ e $\frac{\partial \mathbf{f}}{\partial \mathbf{x}} \in \mathbb{R}^{n \times m}$.

1.4.1 Diferença entre Gradiente e Jacobiano

Para uma função escalar $L : \mathbb{R}^n \rightarrow \mathbb{R}$, ∇L é um vetor. Para uma função vetorial $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, a derivada é uma matriz Jacobiana.

Capítulo 2

Implementacao Base em C: Estruturas e Gerenciamento de Memoria

2.1 Gerenciamento de Memoria para Ten- sores

A base de qualquer framework é o gerenciamento eficiente de arrays multidimensionais. Vamos implementar do zero:

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <string.h>
4 #include <math.h>
5 #include <assert.h>
6
7 typedef struct Tensor {
8     float* data;           // Dados contiguos em memoria
9     float* grad;          // Gradientes acumulados
10    int* shape;           // Dimensoes do tensor (ex:
11        [3,224,224])
12    int ndim;             // Numero de dimensoes
13    int size;              // Produto das dimensoes
14    int refcount;          // Contagem de referencias para
15        garbage collection
16
17 // Para autograd (sera expandido depois)
18 struct Tensor* grad_fn;
```

```

17     void (*backward)(struct Tensor*) ;
18
19     char requires_grad;    // Se calcula gradientes
20 } Tensor;

```

Listing 2.1: Estrutura base do Tensor

2.2 Funções de Criação e Destruição

```

1 Tensor* tensor_create(int* shape, int ndim) {
2     Tensor* t = (Tensor*)malloc(sizeof(Tensor));
3     t->ndim = ndim;
4     t->shape = (int*)malloc(ndim * sizeof(int));
5
6     t->size = 1;
7     for(int i = 0; i < ndim; i++) {
8         t->shape[i] = shape[i];
9         t->size *= shape[i];
10    }
11
12    t->data = (float*)calloc(t->size, sizeof(float));
13    t->grad = (float*)calloc(t->size, sizeof(float));
14    t->refcount = 1;
15    t->requires_grad = 1;
16    t->grad_fn = NULL;
17
18    return t;
19 }
20
21 void tensor_free(Tensor* t) {
22     if(!t) return;
23
24     t->refcount--;
25     if(t->refcount > 0) return;
26
27     free(t->data);
28     free(t->grad);
29     free(t->shape);
30     free(t);
31 }

```

Listing 2.2: Alocacão e liberação de memória

2.3 Operacoes Basicas Otimizadas

Operacoes vetorizadas usando loops otimizados:

```

1 void tensor_add_scalar(Tensor* t, float scalar) {
2     for(int i = 0; i < t->size; i++) {
3         t->data[i] += scalar;
4     }
5 }
6
7 void tensor_multiply_elementwise(Tensor* a, Tensor* b,
8     Tensor* out) {
9     assert(a->size == b->size && a->size == out->size);
10
11    for(int i = 0; i < a->size; i++) {
12        out->data[i] = a->data[i] * b->data[i];
13    }
14
15    if(a->requires_grad || b->requires_grad) {
16        // Guardar para backpropagation
17        out->requires_grad = 1;
18    }
19 }
```

Listing 2.3: Operacoes element-wise

2.4 Produto Matricial Otimizado

Implementacao de multiplicacao de matrizes com tres loops aninhados ($O(n^3)$) - posteriormente otimizaremos com blocagem e SIMD:

```

1 void tensor_matmul(Tensor* a, Tensor* b, Tensor* out) {
2     // Assume a: [m, k], b: [k, n], out: [m, n]
3     int m = a->shape[0];
4     int k = a->shape[1];
5     int n = b->shape[1];
6
7     assert(a->shape[1] == b->shape[0]);
8
9     for(int i = 0; i < m; i++) {
10         for(int j = 0; j < n; j++) {
11             float sum = 0.0f;
12             for(int l = 0; l < k; l++) {
13                 sum += a->data[i * k + l] * b->data[l * n
+ j];
14             }
15         }
16     }
17 }
```

```

14         }
15         out->data[i * n + j] = sum;
16     }
17 }
18 }
```

Listing 2.4: Multiplicação de matrizes

2.5 Estratégias de Otimização de Memória

2.5.1 Pooling de Alocacão

Para evitar fragmentação de memória, implementamos um pool de alocação:

```

1 typedef struct MemoryPool {
2     void* memory;
3     size_t size;
4     size_t used;
5 } MemoryPool;
6
7 MemoryPool* pool_create(size_t size) {
8     MemoryPool* pool = malloc(sizeof(MemoryPool));
9     pool->memory = malloc(size);
10    pool->size = size;
11    pool->used = 0;
12    return pool;
13 }
14
15 void* pool_alloc(MemoryPool* pool, size_t size) {
16     if(pool->used + size > pool->size) return NULL;
17     void* ptr = pool->memory + pool->used;
18     pool->used += size;
19     return ptr;
20 }
```

Listing 2.5: Alocador de memória customizado

Capítulo 3

Rede Neural: Matematica e Implementacao Detalhada

3.1 Arquitetura de uma Rede Neural

3.1.1 Neuronio Artificial: Modelo Matematico

Um neuronio artificial implementa:

$$z = \mathbf{w}^T \mathbf{x} + b = \sum_{i=1}^n w_i x_i + b$$

$$a = \phi(z)$$

Onde ϕ e a funcao de ativacao. As derivadas sao cruciais para back-propagation:

- Sigmoid: $\sigma(z) = \frac{1}{1+e^{-z}}$, $\sigma'(z) = \sigma(z)(1 - \sigma(z))$
- Tanh: $\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$, $\tanh'(z) = 1 - \tanh^2(z)$
- ReLU: $\text{ReLU}(z) = \max(0, z)$, $\text{ReLU}'(z) = \begin{cases} 1 & z > 0 \\ 0 & z \leq 0 \end{cases}$

3.2 Backpropagation: Deducao Passo a Passo

3.2.1 Caso Simples: 2 Camadas

Considere uma rede com uma camada oculta:

$$\begin{aligned} z_1 &= W_1 x + b_1 \\ a_1 &= \sigma(z_1) \\ z_2 &= W_2 a_1 + b_2 \\ \hat{y} &= z_2 \quad (\text{regressao}) \end{aligned}$$

Funcao de custo MSE:

$$L = \frac{1}{2}(\hat{y} - y)^2$$

Derivadas da Camada de Saída

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

$$\frac{\partial L}{\partial W_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_2} \cdot \frac{\partial z_2}{\partial W_2} = (\hat{y} - y) \cdot 1 \cdot a_1^T$$

$$\frac{\partial L}{\partial b_2} = (\hat{y} - y) \cdot 1 \cdot 1$$

Derivadas da Camada Oculta

$$\frac{\partial L}{\partial a_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_2} \cdot \frac{\partial z_2}{\partial a_1} = (\hat{y} - y) \cdot 1 \cdot W_2^T$$

$$\frac{\partial L}{\partial z_1} = \frac{\partial L}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} = \frac{\partial L}{\partial a_1} \odot \sigma'(z_1)$$

$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial z_1} \cdot \frac{\partial z_1}{\partial W_1} = \frac{\partial L}{\partial z_1} \cdot x^T$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial z_1}$$

3.2.2 Implementacao Completa do Backprop

```

1  typedef struct Layer {
2      Tensor* weights;
3      Tensor* bias;
4      Tensor* input;
5      Tensor* output;
6      Tensor* z;    // Pre-ativacao
7      char activation; // 's' sigmoid, 'r' relu, 't' tanh,
8      'n' none
9  } Layer;
10
11 typedef struct MLP {
12     Layer** layers;
13     int n_layers;
14     Tensor* (*forward)(struct MLP*, Tensor* );
15     void (*backward)(struct MLP*, Tensor* loss);
16 } MLP;
17
18 void mlp_backward(MLP* net, Tensor* loss_grad) {
19     Tensor* grad = loss_grad; // dL/dy
20
21     for(int l = net->n_layers - 1; l >= 0; l--) {
22         Layer* layer = net->layers[l];
23
24         // Gradiente dos parametros da camada atual
25         if(layer->weights->requires_grad) {
26             // dL/dW = input^T * grad
27             tensor_matmul_transpose(layer->input, grad,
28                                     layer->weights->grad);
29         }
30
31         if(layer->bias->requires_grad) {
32             // dL/db = sum(grad, axis=0)
33             tensor_sum_axis(grad, 0, layer->bias->grad);
34         }
35
36         // Propagar gradiente para camada anterior
37         if(l > 0) {
38             // dL/d(input) = grad * W^T
39             Tensor* new_grad = tensor_create(...);
40             tensor_matmul(grad, layer->weights, new_grad,
41                           1); // Transpose W
42
43             // Multiplicar pela derivada da ativacao
44             if(layer->activation == 's') {
45                 tensor_sigmoid_prime(layer->z, layer->
z_prime);

```

```

44         tensor_multiply_elementwise(new_grad ,
45             layer->z_prime , new_grad);
46         // ... outras ativações
47
48         grad = new_grad;
49     }
50 }

```

Listing 3.1: Backpropagation para MLP de 2 camadas

3.3 Inicializacao de Pesos

A inicializacao correta dos pesos é critica para evitar gradientes desvanecentes/explosivos.

3.3.1 Xavier/Glorot Initialization

Para ativações sigmoid/tanh:

$$W \sim \mathcal{U} \left(-\sqrt{\frac{6}{n_{in} + n_{out}}}, \sqrt{\frac{6}{n_{in} + n_{out}}} \right)$$

3.3.2 He Initialization

Para ReLU:

$$W \sim \mathcal{N} \left(0, \sqrt{\frac{2}{n_{in}}} \right)$$

```

1 void init_weights_xavier(Tensor* w, int fan_in, int
2     fan_out) {
3     float limit = sqrtf(6.0f / (fan_in + fan_out));
4     for(int i = 0; i < w->size; i++) {
5         w->data[i] = ((float)rand() / RAND_MAX) * 2 *
6             limit - limit;
7     }
8
9 void init_weights_he(Tensor* w, int fan_in) {
    float std = sqrtf(2.0f / fan_in);

```

```
10  for(int i = 0; i < w->size; i++) {  
11      // Box-Muller transform  
12      float u1 = (float)rand() / RAND_MAX;  
13      float u2 = (float)rand() / RAND_MAX;  
14      w->data[i] = sqrtf(-2.0f * logf(u1)) * cosf(2 *  
15          M_PI * u2) * std;  
16 }
```

Listing 3.2: Inicializadores de peso

Capítulo 4

Sistema de Autograd Completo

4.1 Arquitetura do Grafo Computacional

O autograd constroi um grafo aciclico dirigido (DAG) das operacoes.

```
1 typedef struct AutogradNode {
2     Tensor* output;
3     struct AutogradNode** parents;
4     int n_parents;
5
6     // Funcao que calcula gradientes dos pais
7     void (*backward)(struct AutogradNode*);
8
9     // Metadados da operacao
10    int op_type;   // 0: add, 1: mul, 2: matmul, 3: relu,
11    ...
12    void* op_ctx; // Contexto especifico da operacao
13 } AutogradNode;
14
15 typedef struct Tensor {
16     // ... campos anteriores ...
17
18     AutogradNode* grad_node;
19     char grad_computed; // Para evitar computacao
                           duplicada
} Tensor;
```

Listing 4.1: Estrutura do no do grafo

4.2 Implementacao de Operacoes com Grafo

```

1 Tensor* tensor_mul(Tensor* a, Tensor* b) {
2     Tensor* out = tensor_create_like(a); // Assume
3     broadcasting implementado
4
5     // Forward
6     for(int i = 0; i < out->size; i++) {
7         out->data[i] = a->data[i] * b->data[i];
8     }
9
10    // Setup autograd
11    if(a->requires_grad || b->requires_grad) {
12        out->requires_grad = 1;
13        out->grad_node = (AutogradNode*)malloc(sizeof(
14            AutogradNode));
15        out->grad_node->output = out;
16        out->grad_node->parents = (AutogradNode**)malloc
17            (2 * sizeof(AutogradNode*));
18        out->grad_node->parents[0] = a->grad_node;
19        out->grad_node->parents[1] = b->grad_node;
20        out->grad_node->n_parents = 2;
21        out->grad_node->op_type = 1; // multiplication
22        out->grad_node->backward = mul_backward;
23    }
24
25    return out;
26}
27
28 void mul_backward(AutogradNode* node) {
29     Tensor* out = node->output;
30     Tensor* a = node->parents[0]->output;
31     Tensor* b = node->parents[1]->output;
32
33     // dL/da = dL/dout * b
34     if(a->requires_grad) {
35         for(int i = 0; i < a->size; i++) {
36             a->grad[i] += out->grad[i] * b->data[i];
37         }
38     }
39
40     // dL/db = dL/dout * a
41     if(b->requires_grad) {
42         for(int i = 0; i < b->size; i++) {
43             b->grad[i] += out->grad[i] * a->data[i];
44         }
45     }
46 }
```

43 }

Listing 4.2: Operacao de multiplicacao com autograd

4.3 Topological Sort e Propagacao Reversa

```

1 void tensor_backward(Tensor* loss) {
2     if(!loss->requires_grad) return;
3
4     // Inicializa gradiente da loss como 1 (dL/dL = 1)
5     for(int i = 0; i < loss->size; i++) {
6         loss->grad[i] = 1.0f;
7     }
8
9     // Coleta todos os nos do grafo via DFS
10    AutogradNode** nodes = NULL;
11    int n_nodes = 0;
12    collect_nodes(loss->grad_node, &nodes, &n_nodes);
13
14    // Ordenacao topologica (DFS reversa)
15    topological_sort(nodes, n_nodes);
16
17    // Backward em ordem reversa
18    for(int i = n_nodes - 1; i >= 0; i--) {
19        if(nodes[i]->backward) {
20            nodes[i]->backward(nodes[i]);
21        }
22    }
23}
24
25 void collect_nodes(AutogradNode* node, AutogradNode*** nodes, int* n) {
26     if(!node) return;
27
28     // Evita duplicatas
29     for(int i = 0; i < *n; i++) {
30         if((*nodes)[i] == node) return;
31     }
32
33     // Adiciona a lista
34     *nodes = realloc(*nodes, (*n + 1) * sizeof(
35     AutogradNode*));
36     (*nodes)[*n] = node;
37     (*n)++;

```

```
38 // Recursao nos pais
39 for(int i = 0; i < node->n_parents; i++) {
40     collect_nodes(node->parents[i], nodes, n);
41 }
42 }
```

Listing 4.3: Backward recursivo com ordenacao topologica

Capítulo 5

Camadas Avancadas e Funcoes de Ativacao

5.1 Camada Linear (Fully Connected)

```
1 typedef struct Linear {
2     Tensor* weight;
3     Tensor* bias;
4     Tensor* output;
5     Tensor* input_cache;
6
7     int in_features;
8     int out_features;
9 } Linear;
10
11 Linear* linear_create(int in_features, int out_features)
12 {
13     Linear* layer = (Linear*)malloc(sizeof(Linear));
14     layer->in_features = in_features;
15     layer->out_features = out_features;
16
17     int w_shape[] = {out_features, in_features};
18     layer->weight = tensor_create(w_shape, 2);
19     init_weights_he(layer->weight, in_features);
20
21     int b_shape[] = {out_features};
22     layer->bias = tensor_create(b_shape, 1);
23     for(int i = 0; i < out_features; i++) layer->bias->
data[i] = 0;
```

28 CAPÍTULO 5. CAMADAS AVANÇADAS E FUNÇÕES DE ATIVAÇÃO

```
24     return layer;
25 }
26
27 Tensor* linear_forward(Linear* layer, Tensor* input) {
28     // input shape: [batch_size, in_features]
29     // output shape: [batch_size, out_features]
30
31     int batch_size = input->shape[0];
32     int out_shape[] = {batch_size, layer->out_features};
33
34     if(!layer->output) {
35         layer->output = tensor_create(out_shape, 2);
36     }
37
38     // y = x @ W^T + b
39     tensor_matmul(input, layer->weight, layer->output, 1,
40     0); // Transpose W?
41
42     // Adiciona bias (broadcasting)
43     for(int i = 0; i < batch_size; i++) {
44         for(int j = 0; j < layer->out_features; j++) {
45             layer->output->data[i * layer->out_features +
46 j] +=
47                 layer->bias->data[j];
48         }
49     }
50
51     // Cache para backward
52     layer->input_cache = input;
53
54     return layer->output;
55 }
56
55 void linear_backward(Linear* layer, Tensor* grad_output)
56 {
57     // grad_output: [batch_size, out_features]
58     int batch_size = grad_output->shape[0];
59
60     // grad_weight = input^T @ grad_output
61     if(layer->weight->requires_grad) {
62         // grad_weight shape: [out_features, in_features]
63         tensor_matmul_transpose(layer->input_cache,
64         grad_output,
65                                     layer->weight->grad, 1, 0)
66     }
67
68     // grad_bias = sum(grad_output, dim=0)
```

```

67     if(layer->bias->requires_grad) {
68         for(int j = 0; j < layer->out_features; j++) {
69             float sum = 0;
70             for(int i = 0; i < batch_size; i++) {
71                 sum += grad_output->data[i * layer->
72                 out_features + j];
73             }
74             layer->bias->grad[j] += sum;
75         }
76     }
77
78     // grad_input = grad_output @ weight
79     if(layer->input_cache->requires_grad) {
80         tensor_matmul(grad_output, layer->weight,
81                         layer->input_cache->grad, 0, 1); // //
82         Transpose weight
83     }
84 }
```

Listing 5.1: Implementacao completa da camada Linear

5.2 Funcoes de Ativacao com Autograd

5.2.1 ReLU e Leaky ReLU

```

1 Tensor* relu_forward(Tensor* x) {
2     Tensor* out = tensor_create_like(x);
3
4     for(int i = 0; i < x->size; i++) {
5         out->data[i] = x->data[i] > 0 ? x->data[i] : 0;
6     }
7
8     if(x->requires_grad) {
9         out->requires_grad = 1;
10        out->grad_node = autograd_node_create(OP_RELU, 1,
11                                              &x);
12        out->grad_node->backward = relu_backward;
13        out->grad_node->op_ctx = (void*)x; // Salva
14        input para mascara
15    }
16
17    return out;
18 }
```

```

18 void relu_backward(AutogradNode* node) {
19     Tensor* out = node->output;
20     Tensor* x = (Tensor*)node->op_ctx;
21     Tensor* grad_out = out->grad; // dL/d(out)
22
23     if(x->requires_grad) {
24         for(int i = 0; i < x->size; i++) {
25             if(x->data[i] > 0) {
26                 x->grad[i] += grad_out->data[i];
27             }
28         }
29     }
30 }
```

Listing 5.2: ReLU com autograd

5.2.2 Softmax Estável Numericamente

Softmax pode sofrer overflow/underflow. A versão estável subtrai o máximo:

$$p_i = \frac{e^{z_i - \max(z)}}{\sum_j e^{z_j - \max(z)}}$$

```

1 Tensor* softmax_forward(Tensor* x, int dim) {
2     // Assume x shape: [batch_size, n_classes]
3     int batch_size = x->shape[0];
4     int n_classes = x->shape[1];
5
6     Tensor* out = tensor_create_like(x);
7
8     for(int b = 0; b < batch_size; b++) {
9         // Encontra máximo para estabilidade numérica
10        float max_val = x->data[b * n_classes];
11        for(int c = 1; c < n_classes; c++) {
12            if(x->data[b * n_classes + c] > max_val) {
13                max_val = x->data[b * n_classes + c];
14            }
15        }
16
17        // Calcula exponenciais e soma
18        float sum = 0;
19        for(int c = 0; c < n_classes; c++) {
20            int idx = b * n_classes + c;
```

```

21         out->data[idx] = expf(x->data[idx] - max_val)
22     ;
23     sum += out->data[idx];
24 }
25
26 // Normaliza
27 for(int c = 0; c < n_classes; c++) {
28     out->data[b * n_classes + c] /= sum;
29 }
30
31 return out;
32 }
33
34 Tensor* softmax_backward(Tensor* out, Tensor* grad_out) {
35     // Jacobiana da softmax: p_i * (delta_ij - p_j)
36     // dL/dx_i = sum_j (dL/dp_j * p_i * (delta_ij - p_j))
37     // = p_i * (dL/dp_i - sum_j(dL/dp_j * p_j))
38
39     int batch_size = out->shape[0];
40     int n_classes = out->shape[1];
41
42     Tensor* grad_in = tensor_create_like(out);
43
44     for(int b = 0; b < batch_size; b++) {
45         // Calcula dot product: sum(dL/dp_j * p_j)
46         float dot = 0;
47         for(int j = 0; j < n_classes; j++) {
48             int idx = b * n_classes + j;
49             dot += grad_out->data[idx] * out->data[idx];
50         }
51
52         for(int i = 0; i < n_classes; i++) {
53             int idx = b * n_classes + i;
54             grad_in->data[idx] = out->data[idx] * (
55                 grad_out->data[idx] - dot);
56         }
57     }
58
59     return grad_in;
}

```

Listing 5.3: Softmax estavel e seu gradiente

Capítulo 6

Redes Convolucionais: Implementacao Profunda

6.1 Fundamentos da Convolucao Discreta

Convolucao 2D discreta para imagens:

$$(I * K)(i, j) = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} I(i + m, j + n) \cdot K(m, n)$$

6.1.1 Implementacao Ingênuua vs Otimizada

```
1 typedef struct Conv2D {
2     Tensor* weight; // [out_channels, in_channels,
3                      kernel_h, kernel_w]
4     Tensor* bias;   // [out_channels]
5     int stride;
6     int padding;
7
8     // Cache para backward
9     Tensor* input;
10    Tensor* output;
11 } Conv2D;
12
13 Tensor* conv2d_forward(Conv2D* layer, Tensor* input) {
14     // input: [batch, in_channels, h, w]
15     int batch = input->shape[0];
16     int in_c = input->shape[1];
```



```

55         }
56     }
57
58     int out_idx = ((n * out_c + oc) *
59         out_h + oh) * out_w + ow;
60         out->data[out_idx] = sum;
61     }
62 }
63
64 if(layer->padding > 0) {
65     tensor_free(input_padded);
66 }
67
68 return out;
70}

```

Listing 6.1: Convolucao 2D com 6 loops

6.2 Backpropagation em Convolucao

O backward da convolucao e notoriamente complexo. Podemos implementar como convolucao com rotacao de kernels.

```

1 void conv2d_backward(Conv2D* layer, Tensor* grad_output)
2 {
3     // grad_output: [batch, out_c, out_h, out_w]
4     int batch = grad_output->shape[0];
5     int out_c = grad_output->shape[1];
6     int out_h = grad_output->shape[2];
7     int out_w = grad_output->shape[3];
8
9     // grad_weight: rotacionar input e convoluir com
10    // grad_output
11    if(layer->weight->requires_grad) {
12        // Inicializa grad_weight com zeros
13        for(int i = 0; i < layer->weight->size; i++) {
14            layer->weight->grad[i] = 0;
15        }
16
17        for(int n = 0; n < batch; n++) {
18            for(int oc = 0; oc < out_c; oc++) {
19                for(int ic = 0; ic < layer->weight->shape
20                    [1]; ic++) {

```

```

18         for(int kh = 0; kh < layer->weight->
19             shape[2]; kh++) {
20             for(int kw = 0; kw < layer->
21                 weight->shape[3]; kw++) {
22                 float sum = 0;
23
24                 for(int oh = 0; oh < out_h;
25                     oh++) {
26                     for(int ow = 0; ow <
27                         out_w; ow++) {
28                         int ih = oh * layer->
29                             stride + kh;
30                         int iw = ow * layer->
31                             stride + kw;
32
33                         int in_idx = ((n *
34                             layer->weight->shape[1] + ic) *
35                             layer->
36                             input->shape[2] + ih) *
37                             layer->
38                             input->shape[3] + iw;
39                         int out_idx = ((n *
40                             out_c + oc) * out_h + oh) * out_w + ow;
41
42                         sum += layer->input->
43                             data[in_idx] *
44                             grad_output->
45                             data[out_idx];
46
47                     }
48                 }
49             }
50         }
51
52         int w_idx = ((oc * layer->
53             weight->shape[1] + ic) *
54             layer->weight->
55             shape[2] + kh) *
56             layer->weight->
57             shape[3] + kw;
58
59         layer->weight->grad[w_idx] +=
60             sum;
61
62     }
63 }
64 }
65 }
66 }
67 }
68 }
69 }
70 }
71 }
72 }
73 }
74 }
75 }
76 }
77 }
78 }
79 }
80 }
81 }
82 }
83 }
84 }
85 }
86 }
87 }
88 }
89 }
90 }
91 }
92 }
93 }
94 }
95 }
96 }
97 }
98 }
99 }
100 }
101 }
102 }
103 }
104 }
105 }
106 }
107 }
108 }
109 }
110 }
111 }
112 }
113 }
114 }
115 }
116 }
117 }
118 }
119 }
120 }
121 }
122 }
123 }
124 }
125 }
126 }
127 }
128 }
129 }
130 }
131 }
132 }
133 }
134 }
135 }
136 }
137 }
138 }
139 }
140 }
141 }
142 }
143 }
144 }
145 }
146 }
147 }
148 }
149 }

// grad_bias: soma sobre batch, altura, largura
if(layer->bias->requires_grad) {

```

```

50     for(int oc = 0; oc < out_c; oc++) {
51         float sum = 0;
52         for(int n = 0; n < batch; n++) {
53             for(int oh = 0; oh < out_h; oh++) {
54                 for(int ow = 0; ow < out_w; ow++) {
55                     int idx = ((n * out_c + oc) *
56 out_h + oh) * out_w + ow;
57                     sum += grad_output->data[idx];
58                 }
59             }
60             layer->bias->grad[oc] += sum;
61         }
62     }
63
64 // grad_input: convolucao transposta (full
65 convolution)
66 if(layer->input->requires_grad) {
67     // Implementar conv2d_transpose
68     Tensor* grad_input = conv2d_transpose(grad_output,
69     , layer->weight,
70                                         layer->
71 stride, layer->padding);
72     // Acumular em layer->input->grad
73     tensor_add_to(grad_input, layer->input->grad);
74     tensor_free(grad_input);
75 }
76 }
```

Listing 6.2: Backward da convolucao

6.3 Otimizacoes: Im2Col e GEMM

A abordagem de 6 loops é muito lenta. Podemos converter convolução em multiplicação de matrizes:

```

1 Tensor* im2col(Tensor* im, int k_h, int k_w, int stride,
2     int pad) {
3     // im: [batch, channels, h, w]
4     int batch = im->shape[0];
5     int channels = im->shape[1];
6     int h = im->shape[2];
7     int w = im->shape[3];
8
9     int out_h = (h + 2*pad - k_h) / stride + 1;
```

```

9   int out_w = (w + 2*pad - k_w) / stride + 1;
10
11 // Matriz de saída: [batch * out_h * out_w, channels
12 * k_h * k_w]
13 int col_shape[] = {batch * out_h * out_w, channels *
14 k_h * k_w};
15 Tensor* col = tensor_create(col_shape, 2);
16
17 for(int b = 0; b < batch; b++) {
18     for(int oh = 0; oh < out_h; oh++) {
19         for(int ow = 0; ow < out_w; ow++) {
20             int row_idx = (b * out_h + oh) * out_w +
21             ow;
22
23             for(int c = 0; c < channels; c++) {
24                 for(int kh = 0; kh < k_h; kh++) {
25                     for(int kw = 0; kw < k_w; kw++) {
26                         int ih = oh * stride + kh -
27                         pad;
28                         int iw = ow * stride + kw -
29                         pad;
30
31                         float val = 0;
32                         if(ih >= 0 && ih < h && iw >=
33                         0 && iw < w) {
34                             int im_idx = ((b *
35                             channels + c) * h + ih) * w + iw;
36                             val = im->data[im_idx];
37                         }
38
39                         int col_idx = row_idx * col->
40                         shape[1] +
41                         ((c * k_h + kh)
42                         * k_w + kw);
43                         col->data[col_idx] = val;
44                     }
45                 }
46             }
47         }
48     }
49
50     return col;
51 }
52
53 Tensor* conv2d_im2col(Conv2D* layer, Tensor* input) {
54     // Converte para multiplicação de matrizes
55     Tensor* col = im2col(input, layer->weight->shape[2],
56

```

```
48         layer->weight->shape[3] ,
49         layer->stride , layer->padding) ;
50
51 // weight: [out_c, in_c * k_h * k_w]
52 int w_shape[] = {layer->weight->shape[0] ,
53                  layer->weight->shape[1] *
54                  layer->weight->shape[2] *
55                  layer->weight->shape[3]} ;
56 Tensor* w_flat = tensor_reshape(layer->weight ,
57                                 w_shape , 2) ;
58
59 // output = col @ w_flat^T
60 int out_shape[] = {input->shape[0] ,
61                    layer->weight->shape[0] ,
62                    col->shape[0] / input->shape[0] ,
63                    1}; // Ajustar
64 Tensor* out = tensor_matmul(col , w_flat , 0 , 1) ;
65
66 // Adicionar bias e reshape
67
68 tensor_free(col);
69 return out;
}
```

Listing 6.3: Algoritmo im2col para convolucao otimizada

Capítulo 7

Otimizadores: Teoria e Implementacao

7.1 Stochastic Gradient Descent (SGD)

7.1.1 Com e sem Momentum

Momentum acumula gradientes passados:

$$v_t = \mu v_{t-1} + \eta \nabla L(\theta_t)$$

$$\theta_{t+1} = \theta_t - v_t$$

```
1 typedef struct SGD {
2     float lr;
3     float momentum;
4     float weight_decay;
5
6     // Velocidades para cada tensor
7     Tensor** velocities;
8     int n_params;
9 } SGD;
10
11 SGD* sgd_create(float lr, float momentum, float
12     weight_decay) {
13     SGD* opt = (SGD*)malloc(sizeof(SGD));
14     opt->lr = lr;
15     opt->momentum = momentum;
16     opt->weight_decay = weight_decay;
```

42 CAPÍTULO 7. OTIMIZADORES: TEORIA E IMPLEMENTACAO

```

16     opt->velocities = NULL;
17     opt->n_params = 0;
18     return opt;
19 }
20
21 void sgd_add_param(SGD* opt, Tensor* param) {
22     opt->n_params++;
23     opt->velocities = realloc(opt->velocities,
24                               opt->n_params * sizeof(
25                                     Tensor*));
26
27     int idx = opt->n_params - 1;
28     opt->velocities[idx] = tensor_create_like(param);
29     // Inicializa velocidades com zero
30     memset(opt->velocities[idx]->data, 0,
31            opt->velocities[idx]->size * sizeof(float));
32 }
33
34 void sgd_step(SGD* opt) {
35     for(int i = 0; i < opt->n_params; i++) {
36         Tensor* param = opt->velocities[i]; // Na verdade
37         Tensor* vel = opt->velocities[i];
38
39         for(int j = 0; j < param->size; j++) {
40             // weight decay: adiciona penalidade L2
41             float grad = param->grad[j];
42             if(opt->weight_decay > 0) {
43                 grad += opt->weight_decay * param->data[j]
44             ];
45         }
46
47         // momentum
48         vel->data[j] = opt->momentum * vel->data[j] +
49         opt->lr * grad;
50
51         // update
52         param->data[j] -= vel->data[j];
53     }
54 }
55 }
```

Listing 7.1: SGD com momentum

7.2 Adam: Adaptive Moment Estimation

Adam combina momentum com adaptacao de taxa de aprendizado por parametro:

$$\begin{aligned} m_t &= \beta_1 m_{t-1} + (1 - \beta_1) g_t \\ v_t &= \beta_2 v_{t-1} + (1 - \beta_2) g_t^2 \\ \hat{m}_t &= \frac{m_t}{1 - \beta_1^t} \\ \hat{v}_t &= \frac{v_t}{1 - \beta_2^t} \\ \theta_{t+1} &= \theta_t - \eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} \end{aligned}$$

```

1  typedef struct Adam {
2      float lr;
3      float beta1;
4      float beta2;
5      float eps;
6
7      Tensor** m;    // Primeiro momento
8      Tensor** v;    // Segundo momento
9      int* t;        // Timestep para cada parametro
10     int n_params;
11 } Adam;
12
13 Adam* adam_create(float lr, float beta1, float beta2,
14                     float eps) {
15     Adam* opt = (Adam*)malloc(sizeof(Adam));
16     opt->lr = lr;
17     opt->beta1 = beta1;
18     opt->beta2 = beta2;
19     opt->eps = eps;
20     opt->m = NULL;
21     opt->v = NULL;
22     opt->t = NULL;
23     opt->n_params = 0;
24     return opt;
25 }
26
27 void adam_add_param(Adam* opt, Tensor* param) {
28     opt->n_params++;
29     opt->m = realloc(opt->m, opt->n_params * sizeof(Tensor));
30 }
```

```

29     opt->v = realloc(opt->v, opt->n_params * sizeof(
30         Tensor*));
31     opt->t = realloc(opt->t, opt->n_params * sizeof(int))
32     ;
33
34     int idx = opt->n_params - 1;
35     opt->m[idx] = tensor_create_like(param);
36     opt->v[idx] = tensor_create_like(param);
37     memset(opt->m[idx]->data, 0, opt->m[idx]->size *
38         sizeof(float));
39     memset(opt->v[idx]->data, 0, opt->v[idx]->size *
40         sizeof(float));
41     opt->t[idx] = 0;
42 }
43
44 void adam_step(Adam* opt) {
45     for(int i = 0; i < opt->n_params; i++) {
46         Tensor* param = /* ... */;
47         Tensor* m = opt->m[i];
48         Tensor* v = opt->v[i];
49         opt->t[i]++;
50
51         float bias_correction1 = 1.0f - powf(opt->beta1,
52             opt->t[i]);
53         float bias_correction2 = 1.0f - powf(opt->beta2,
54             opt->t[i]);
55
56         for(int j = 0; j < param->size; j++) {
57             float g = param->grad[j];
58
59             // Atualiza momentos
60             m->data[j] = opt->beta1 * m->data[j] + (1 -
61             opt->beta1) * g;
62             v->data[j] = opt->beta2 * v->data[j] + (1 -
63             opt->beta2) * g * g;
64
65             // Correcao de vies
66             float m_hat = m->data[j] / bias_correction1;
67             float v_hat = v->data[j] / bias_correction2;
68
69             // Atualiza parametro
70             param->data[j] -= opt->lr * m_hat / (sqrtf(
71                 v_hat) + opt->eps);
72         }
73     }
74 }

```

Listing 7.2: Implementacao completa do Adam

Capítulo 8

Treinamento e Validação

8.1 DataLoader e Batching

```
1  typedef struct Dataset {
2      Tensor* data;
3      Tensor* targets;
4      int size;
5  } Dataset;
6
7  typedef struct DataLoader {
8      Dataset* dataset;
9      int batch_size;
10     int* indices;
11     int current_idx;
12     int shuffle;
13 } DataLoader;
14
15 DataLoader* dataloader_create(Dataset* dataset, int
16                                batch_size, int shuffle) {
17     DataLoader* dl = (DataLoader*)malloc(sizeof(DataLoader));
18     dl->dataset = dataset;
19     dl->batch_size = batch_size;
20     dl->shuffle = shuffle;
21     dl->current_idx = 0;
22
23     dl->indices = (int*)malloc(dataset->size * sizeof(int));
24     for(int i = 0; i < dataset->size; i++) {
25         dl->indices[i] = i;
26     }
```

```
26
27     if(shuffle) {
28         for(int i = dataset->size - 1; i > 0; i--) {
29             int j = rand() % (i + 1);
30             int temp = dl->indices[i];
31             dl->indices[i] = dl->indices[j];
32             dl->indices[j] = temp;
33         }
34     }
35
36     return dl;
37 }
38
39 int dataloader_next_batch(DataLoader* dl, Tensor** batch_data, Tensor** batch_targets) {
40     if(dl->current_idx >= dl->dataset->size) {
41         dl->current_idx = 0;
42         if(dl->shuffle) {
43             // Reshuffle
44         }
45         return 0;
46     }
47
48     int batch_end = dl->current_idx + dl->batch_size;
49     if(batch_end > dl->dataset->size) {
50         batch_end = dl->dataset->size;
51     }
52     int current_batch_size = batch_end - dl->current_idx;
53
54     // Alocar tensores do batch
55     int data_shape[] = {current_batch_size,
56                         dl->dataset->data->shape[1]}; // Simplificado
57     *batch_data = tensor_create(data_shape, 2);
58
59     int target_shape[] = {current_batch_size,
60                          dl->dataset->targets->shape[1]};
61     *batch_targets = tensor_create(target_shape, 2);
62
63     // Copiar dados
64     for(int i = 0; i < current_batch_size; i++) {
65         int idx = dl->indices[dl->current_idx + i];
66
67         // Copiar features
68         for(int j = 0; j < dl->dataset->data->shape[1]; j++)
69         {
```

```

69         (*batch_data)->data[i * dl->dataset->data->
70         shape[1] + j] =
71             dl->dataset->data->data[idx * dl->dataset-
72             ->data->shape[1] + j];
73     }
74
75     // Copiar targets
76     for(int j = 0; j < dl->dataset->targets->shape
77         [1]; j++) {
78         (*batch_targets)->data[i * dl->dataset->
79         targets->shape[1] + j] =
80             dl->dataset->targets->data[idx * dl->
81             dataset->targets->shape[1] + j];
82     }
83
84     dl->current_idx = batch_end;
85     return current_batch_size;
86 }
```

Listing 8.1: DataLoader com shuffle

8.2 Loop de Treinamento Completo

```

1 typedef struct Trainer {
2     MLP* model;
3     Optimizer* optimizer;
4     float (*loss_fn)(Tensor*, Tensor*);
5     void (*loss_backward)(Tensor*, Tensor*);
6
7     float train_loss;
8     float val_loss;
9     float train_acc;
10    float val_acc;
11 } Trainer;
12
13 void trainer_train_epoch(Trainer* trainer, DataLoader*
14     train_loader) {
15     float total_loss = 0;
16     int total_correct = 0;
17     int total_samples = 0;
18
19     while(1) {
20         Tensor* batch_x;
21         Tensor* batch_y;
```

```

21     int batch_size = dataloader_next_batch(
22         train_loader, &batch_x, &batch_y);
23         if(batch_size == 0) break;
24
25         // Forward pass
26         Tensor* output = mlp_forward(trainer->model,
27             batch_x);
28
29         // Compute loss
30         float loss = trainer->loss_fn(output, batch_y);
31         total_loss += loss * batch_size;
32
33         // Compute accuracy
34         int correct = compute_accuracy(output, batch_y);
35         total_correct += correct;
36         total_samples += batch_size;
37
38         // Backward pass
39         trainer->loss_backward(output, batch_y);
40         tensor_backward(output);
41
42         // Update weights
43         optimizer_step(trainer->optimizer);
44
45         // Zero gradients
46         optimizer_zero_grad(trainer->optimizer);
47
48         // Cleanup
49         tensor_free(batch_x);
50         tensor_free(batch_y);
51         tensor_free(output);
52     }
53
54     trainer->train_loss = total_loss / total_samples;
55     trainer->train_acc = (float)total_correct /
56     total_samples;
57 }
58
59 void trainer_validate(Trainer* trainer, DataLoader*
60 val_loader) {
61     // Similar ao train, mas sem gradientes e updates
62     // ...
63 }
```

Listing 8.2: Training loop profissional

Capítulo 9

Serializacao e Persistencia de Modelos

9.1 Salvando e Carregando Modelos

Formato binario portavel com metadados:

```
1 typedef struct ModelHeader {
2     char magic[4];    // "DLFW"
3     int version;
4     int n_layers;
5     int model_type;
6     int data_size;
7 } ModelHeader;
8
9 void model_save(MLP* model, const char* filename) {
10    FILE* fp = fopen(filename, "wb");
11
12    // Header
13    ModelHeader header = {
14        .magic = {'D', 'L', 'F', 'W'},
15        .version = 1,
16        .n_layers = model->n_layers,
17        .model_type = 0,
18        .data_size = 0
19    };
20
21    // Calcular tamanho total
22    for(int i = 0; i < model->n_layers; i++) {
23        if(model->layers[i]->weights) {
24            header.data_size += model->layers[i]->weights
```

50 CAPÍTULO 9. SERIALIZACAO E PERSISTENCIA DE MODELOS

```
    ->size * sizeof(float);
25        header.data_size += model->layers[i]->bias->
26            size * sizeof(float);
27    }
28}
29
30fwrite(&header, sizeof(ModelHeader), 1, fp);
31
32// Salvar arquitetura
33for(int i = 0; i < model->n_layers; i++) {
34    Layer* layer = model->layers[i];
35
36    // Tipo da camada
37    int layer_type = 0; // 0: linear, 1: conv2d, etc.
38    fwrite(&layer_type, sizeof(int), 1, fp);
39
40    // Parametros da camada
41    int in_features = layer->weights->shape[1];
42    int out_features = layer->weights->shape[0];
43    fwrite(&in_features, sizeof(int), 1, fp);
44    fwrite(&out_features, sizeof(int), 1, fp);
45
46    // Ativacao
47    fwrite(&layer->activation, sizeof(char), 1, fp);
48
49    // Pesos e bias
50    fwrite(layer->weights->data, sizeof(float),
51           layer->weights->size, fp);
52    fwrite(layer->bias->data, sizeof(float),
53           layer->bias->size, fp);
54}
55
56fclose(fp);
57}
58
59MLP* model_load(const char* filename) {
60    FILE* fp = fopen(filename, "rb");
61
62    // Ler header
63    ModelHeader header;
64    fread(&header, sizeof(ModelHeader), 1, fp);
65
66    if(header.magic[0] != 'D' || header.magic[1] != 'L'
67    ||
68        header.magic[2] != 'F' || header.magic[3] != 'W')
69    {
70        printf("Formato de arquivo invalido\n");
71        return NULL;
72    }
73}
```

```
69 }
70
71 // Reconstruir modelo
72 MLP* model = mlp_create();
73
74 for(int i = 0; i < header.n_layers; i++) {
75     int layer_type;
76     fread(&layer_type, sizeof(int), 1, fp);
77
78     int in_features, out_features;
79     fread(&in_features, sizeof(int), 1, fp);
80     fread(&out_features, sizeof(int), 1, fp);
81
82     char activation;
83     fread(&activation, sizeof(char), 1, fp);
84
85     Layer* layer = layer_create(in_features,
86         out_features);
87     layer->activation = activation;
88
89     // Carregar pesos
90     fread(layer->weights->data, sizeof(float),
91             layer->weights->size, fp);
92     fread(layer->bias->data, sizeof(float),
93             layer->bias->size, fp);
94
95     mlp_add_layer(model, layer);
96 }
97 fclose(fp);
98 return model;
99 }
```

Listing 9.1: Save/Load completo

Capítulo 10

Estudos de Caso e Projetos Práticos

10.1 Exemplo 1: MNIST do Zero

Implementação completa de reconhecimento de dígitos:

```
1 int main() {
2     // 1. Carregar dados
3     Dataset* train_dataset = mnist_load("train-images.
4                                         idx3-ubyte",
5                                         "train-labels.
6                                         idx1-ubyte");
7     Dataset* test_dataset = mnist_load("t10k-images.idx3-
8                                         ubyte",
9                                         "t10k-labels.idx1-
10                                        ubyte");
11
12     // 2. Criar modelo
13     MLP* model = mlp_create();
14     mlp_add_layer(model, linear_create(784, 256));
15     mlp_add_layer(model, relu_create());
16     mlp_add_layer(model, linear_create(256, 128));
17     mlp_add_layer(model, relu_create());
18     mlp_add_layer(model, linear_create(128, 10));
19     mlp_add_layer(model, softmax_create());
20
21     // 3. Configurar treinamento
22     Adam* optimizer = adam_create(0.001, 0.9, 0.999, 1e
23 -8);
24     mlp_add_parameters(model, optimizer);
```

```

20
21     Trainer* trainer = trainer_create(model, optimizer,
22                                         cross_entropy_loss,
23
24                                         cross_entropy_backward);
25
26     // 4. Loop de treinamento
27     DataLoader* train_loader = dataloader_create(
28         train_dataset, 64, 1);
29     DataLoader* test_loader = dataloader_create(
30         test_dataset, 64, 0);
31
32     for(int epoch = 0; epoch < 10; epoch++) {
33         trainer_train_epoch(trainer, train_loader);
34         trainer_validate(trainer, test_loader);
35
36         printf("Epoch %d: Train Acc: %.4f, Test Acc: %.4f
37             \n",
38             epoch, trainer->train_acc, trainer->
39             val_acc);
40     }
41
42     // 5. Salvar modelo treinado
43     model_save(model, "mnist_model.dlfw");
44
45     return 0;
46 }
```

Listing 10.1: MNIST pipeline completo

10.2 Exemplo 2: CNN para CIFAR-10

```

1 Model* create_cifar10_cnn() {
2     Model* model = model_create();
3
4     // Conv1: 32x32x3 -> 32x32x32
5     model_add(model, conv2d_create(3, 32, 3, 1, 1));
6     model_add(model, batch_norm_create(32));
7     model_add(model, relu_create());
8     model_add(model, maxpool2d_create(2, 2));
9
10    // Conv2: 16x16x32 -> 16x16x64
11    model_add(model, conv2d_create(32, 64, 3, 1, 1));
12    model_add(model, batch_norm_create(64));
13    model_add(model, relu_create());
```

```
14 model_add(model, maxpool2d_create(2, 2));  
15  
16 // Conv3: 8x8x64 -> 8x8x128  
17 model_add(model, conv2d_create(64, 128, 3, 1, 1));  
18 model_add(model, batch_norm_create(128));  
19 model_add(model, relu_create());  
20  
21 // Global average pooling  
22 model_add(model, global_avgpool2d_create());  
23  
24 // Classifier  
25 model_add(model, linear_create(128, 256));  
26 model_add(model, relu_create());  
27 model_add(model, dropout_create(0.5));  
28 model_add(model, linear_create(256, 10));  
29 model_add(model, softmax_create());  
30  
31 return model;  
32 }
```

Listing 10.2: CNN completa para CIFAR-10

Capítulo 11

Otimizações Avançadas e Performance

11.1 BLAS-Level Optimizations

11.1.1 Tiling para Cache L1/L2

```
1 #define BLOCK_SIZE 32
2
3 void matmul_tiled(Tensor* a, Tensor* b, Tensor* out,
4                     int M, int N, int K) {
5
6     for(int i = 0; i < M; i += BLOCK_SIZE) {
7         for(int j = 0; j < N; j += BLOCK_SIZE) {
8             for(int k = 0; k < K; k += BLOCK_SIZE) {
9
10                 // Process a block (i, i+BS) x (k, k+BS) * (
11                 k, k+BS) x (j, j+BS)
12                 for(int ii = i; ii < i + BLOCK_SIZE && ii
13 < M; ii++) {
14                     for(int jj = j; jj < j + BLOCK_SIZE
15 && jj < N; jj++) {
16                         float sum = out->data[ii * N + jj];
17
18                         for(int kk = k; kk < k +
19                             BLOCK_SIZE && kk < K; kk++) {
20                             sum += a->data[ii * K + kk] *
21                             b->data[kk * N + jj];
22                         }
23                     }
24                 }
25             }
26         }
27     }
28 }
```

```

18         out->data[ii * N + jj] = sum;
19     }
20   }
21 }
22 }
23 }
24 }
25 }
```

Listing 11.1: Multiplicacao de matrizes com tiling

11.1.2 SIMD Intrinsics

```

1 #include <immintrin.h>
2
3 void tensor_add_simd(float* a, float* b, float* out, int
4 n) {
5     int i = 0;
6
7 #ifdef __AVX__
8 // Processa 8 floats por vez com AVX
9 for(; i <= n - 8; i += 8) {
10    __m256 va = _mm256_loadu_ps(&a[i]);
11    __m256 vb = _mm256_loadu_ps(&b[i]);
12    __m256 vout = _mm256_add_ps(va, vb);
13    _mm256_storeu_ps(&out[i], vout);
14 }
15 #elif __SSE__
16 // Processa 4 floats por vez com SSE
17 for(; i <= n - 4; i += 4) {
18    __m128 va = _mm_loadu_ps(&a[i]);
19    __m128 vb = _mm_loadu_ps(&b[i]);
20    __m128 vout = _mm_add_ps(va, vb);
21    _mm_storeu_ps(&out[i], vout);
22 }
23 #endif
24
25 // Processa o restante
26 for(; i < n; i++) {
27     out[i] = a[i] + b[i];
28 }
```

Listing 11.2: Usando SSE/AVX para operacoes vetoriais

11.2 Parallelism with OpenMP

```
1 #include <omp.h>
2
3 void conv2d_parallel(Conv2D* layer, Tensor* input, Tensor*
4     * output) {
5     int batch = input->shape[0];
6     int out_c = output->shape[1];
7     int out_h = output->shape[2];
8     int out_w = output->shape[3];
9
10    #pragma omp parallel for collapse(4) schedule(dynamic
11    )
12    for(int n = 0; n < batch; n++) {
13        for(int oc = 0; oc < out_c; oc++) {
14            for(int oh = 0; oh < out_h; oh++) {
15                for(int ow = 0; ow < out_w; ow++) {
16                    // Calculo da convolucao
17                    // ...
18                }
19            }
20        }
}
```

Listing 11.3: Paralelizacao de operacoes

Capítulo 12

Proximos Passos: Alem do Framework

12.1 Suporte a GPU com CUDA

```
1 #ifdef USE_CUDA
2 #include <cuda_runtime.h>
3
4 __global__ void relu_kernel(float* x, float* out, int n)
5 {
6     int idx = blockIdx.x * blockDim.x + threadIdx.x;
7     if(idx < n) {
8         out[idx] = x[idx] > 0 ? x[idx] : 0;
9     }
10
11 void relu_cuda(Tensor* x, Tensor* out) {
12     float* d_x, d_out;
13
14     // Alocar memoria na GPU
15     cudaMalloc(&d_x, x->size * sizeof(float));
16     cudaMalloc(&d_out, out->size * sizeof(float));
17
18     // Copiar dados para GPU
19     cudaMemcpy(d_x, x->data, x->size * sizeof(float),
20               cudaMemcpyHostToDevice);
21
22     // Lancar kernel
23     int block_size = 256;
```

```

24     int num_blocks = (x->size + block_size - 1) /
25         block_size;
26     relu_kernel<<<num_blocks, block_size>>>(d_x, d_out, x
27         ->size);
28
29     // Copiar resultado de volta
30     cudaMemcpy(out->data, d_out, out->size * sizeof(float)
31     ),
32             cudaMemcpyDeviceToHost);
33
34     // Liberar memoria
35     cudaFree(d_x);
36     cudaFree(d_out);
37 }
38 #endif

```

Listing 12.1: Interface CUDA para operacoes

12.2 Implementacao de Transformers

Arquitetura completa de Transformer para processamento de sequencias:

- Multi-Head Self-Attention
- Positional Encoding
- Feed-Forward Networks
- Layer Normalization
- Residual Connections

12.3 Mixed Precision Training

Implementacao de treinamento com FP16 para maior performance:

```

1 void convert_to_fp16(float* fp32, __fp16* fp16, int n) {
2     for(int i = 0; i < n; i++) {
3         fp16[i] = (__fp16)fp32[i];
4     }
5 }

```

```
6
7 void convert_to_fp32(__fp16* fp16, float* fp32, int n) {
8     for(int i = 0; i < n; i++) {
9         fp32[i] = (float)fp16[i];
10    }
11 }
12
13 void matmul_mixed(Tensor* a_fp32, Tensor* b_fp32, Tensor*
14                     out_fp32) {
15     // Converter para FP16
16     __fp16* a_fp16 = malloc(a_fp32->size * sizeof(__fp16)
17 );
18     __fp16* b_fp16 = malloc(b_fp32->size * sizeof(__fp16)
19 );
20
21     convert_to_fp16(a_fp32->data, a_fp16, a_fp32->size);
22     convert_to_fp16(b_fp32->data, b_fp16, b_fp32->size);
23
24     // Multiplicacao em FP16 (usando hardware nativo)
25     __fp16* out_fp16 = malloc(out_fp32->size * sizeof(
26     __fp16));
27     matmul_fp16(a_fp16, b_fp16, out_fp16, /* dims */);
28
29     // Converter de volta
30     convert_to_fp32(out_fp16, out_fp32->data, out_fp32->
31     size);
32
33     free(a_fp16);
34     free(b_fp16);
35     free(out_fp16);
36 }
```

Listing 12.2: Cast de precisao mista

Conclusao

Voce agora tem em maos um framework completo de deep learning implementado em C puro, com:

- Sistema de tensores com gerenciamento de memoria
- Autograd com grafo computacional dinamico
- Camadas: Linear, Conv2D, BatchNorm, Dropout
- Ativacoes: ReLU, Sigmoid, Tanh, Softmax
- Otimizadores: SGD, Adam
- Funcoes de custo: MSE, Cross-Entropy
- DataLoaders e batching
- Serializacao de modelos
- Otimizacoes de performance: SIMD, OpenMP, CUDA

Este codigo e completamente funcional e pode ser usado para treinar redes neurais reais em problemas como MNIST, CIFAR-10, e ate mesmo tarefas mais complexas.

O proximo passo e expandir para arquiteturas modernas como Transformers, GANs, e modelos de difusao, sempre mantendo a filosofia de implementar cada detalhe do zero para comprehensao profunda.

A jornada continua...