



Universidade de Brasília

Instituto de Ciências Exatas
Departamento de Ciência da Computação

Compressão de imagens com perda usando Redes Neurais

Raphael Soares Ramos

Monografia apresentada como requisito parcial
para conclusão do Bacharelado em Ciência da Computação

Orientador
Prof. Dr. Teófilo Emidio de Campos

Brasília
2019



Universidade de Brasília

Instituto de Ciências Exatas
Departamento de Ciência da Computação

Compressão de imagens com perda usando Redes Neurais

Raphael Soares Ramos

Monografia apresentada como requisito parcial
para conclusão do Bacharelado em Ciência da Computação

Prof. Dr. Teófilo Emidio de Campos (Orientador)
CIC/UnB

Dr. Edson Mitsu Hung Prof. Dr. Bruno L. Macchiavello Espinoza
FT/UnB CIC/UnB

Prof. Dr. Edison Ishikawa
Coordenador do Bacharelado em Ciência da Computação

Brasília, 8 de julho de 2019

Dedicatória

Dedico esse trabalho a todas pessoas que sempre me apoiaram e acreditaram em mim.

Agradecimentos

Agradeço a todos que me ajudaram a chegar até aqui.

Resumo

Os recursos requeridos para armazenar e transmitir imagens são imensos, o que torna a sua compressão necessária. Todos os esforços feitos em algoritmos de compressão de imagens clássicos abordam o problema de compressão de um ponto de vista empírico: humanos desenvolvem várias heurísticas para reduzir a quantidade de informação necessária para representar a imagem explorando imperfeições no sistema visual humano, de modo que seja possível reconstruí-la sem muita perda de informação.

O sucesso das redes neurais convolucionais profundas (CNNs) em visão computacional tem inspirado pesquisadores da comunidade de compressão de imagens a tentar desenvolver algoritmos que aprendam com os dados, em vez de confiar no conhecimento de especialistas. Até agora, esses algoritmos não levaram a uma melhoria inquestionável em relação aos codecs clássicos.

O objetivo do presente trabalho é estudar e explorar soluções usando *autoencoders* convolucionais para o desafio de comprimir imagens, buscando propor um método que possa ser competitivo em relação aos codecs clássicos. Para isso, foram testados e avaliados os métodos de compressão clássico *JPEG* e *JPEG2000*. Esses métodos foram comparados com o framework de compressão de imagens usando autoencoders convolucionais desenvolvido em bases de dados comumente usadas para este tipo de problema.

Palavras-chave: codificação de imagens, redes neurais, aprendizado profundo, processamento de imagens, aprendizado de máquina, processamento de sinais

Abstract

The resources required to store and transmit images are huge, making their compression essential. All the efforts made in classical image compression algorithms address the problem from an empiric point of view: experts develop several heuristics to reduce the amount of information needed to represent the image by exploiting imperfections in the human visual system. This way, it is possible to reconstruct it without much perceptible information loss.

The success of deep convolutional neural networks (CNNs) in computer vision application has been inspiring researchers from the image compression community to try to develop algorithms that learn from data, rather than relying on expert knowledge. So far, these algorithms have not lead to an unquestionable improvement over classical codecs.

The objective of the present work is to study and explore neural networks solutions to the challenge of compressing images, aiming to propose a method that can be competitive to classical codecs in some scenarios. To achieve this goal, the classical compression method JPEG and JPEG2000 is evaluated in databases commonly used for this type of problem. Next, some experiments on encoder-decoder image compression framework have been performed.

Keywords: image coding, neural networks, deep learning, image processing, machine learning, signal processing

Sumário

1 Introdução	1
1.1 Motivação	1
1.2 Hipótese	5
1.3 Objetivos	5
1.4 Resultados Esperados	6
2 Trabalhos similares e conceitos fundamentais	7
2.1 JPEG	7
2.1.1 Conversão do espaço de cor	7
2.1.2 Aplicação da transformada direta de cossenos	8
2.1.3 Quantização	10
2.1.4 Codificação	13
2.2 Redes Neurais	13
2.2.1 Taxa de aprendizagem (<i>learning rate</i>)	14
2.2.2 Redes Neurais Convolucionais	16
2.2.3 Autoencoders	17
3 Metodologia	24
3.1 Bases de Dados	24
3.2 Modelos desenvolvidos	26
3.2.1 Modelo 1	29
3.2.2 Modelo 2	30
4 Experimentos e Resultados	33
4.1 JPEG	33
4.2 Modelo 1	34
4.3 Modelo 2	35
4.4 Modelo 3	35
4.5 Modelo 4	38
4.6 Modelo 5	48

4.7 Ganhos obtidos pelo GZIP	48
4.7.1 Modelo 4	48
5 Conclusão	51
5.1 Limitações do Trabalho	51
5.2 Análise Crítica	51
5.3 Trabalhos Futuros	52
Referências	53

Listas de Figuras

2.1	Diagram do método de compressão JPEG. Fonte: [1].	8
2.2	Imagen original sem compressão retirada de [2] e imagem (direita), gerada a partir da imagem original, com dimensionalidade nos canais de crominância reduzida por um fator de 8 nas duas direções. Pode-se jogar informação da imagem original fora e então usar a imagem da direita para armazenamento ou transmissão, que terá novos valores em seus canais de cores (aumenta-se a dimensionalidade apenas quando for necessário exibi-la). Fonte: [2].	9
2.3	Versão com zoom das imagens mostradas na Figura 2.2. A região mostrada é exatamente a mesma (espacialmente) para as duas imagens, com a mesma quantidade de pixels. Imagem original sem compressão (esquerda) e imagem com 64 vezes menos cor (direita).	9
2.4	64 (8 por 8) ondas base de cossenos com frequências variadas. Fonte: [3]. .	11
2.5	Comparação imagem de texto com e sem compressão. Fonte: [4].	12
2.6	Sequência zig-zague usada para melhorar codificação. Fonte: [1].	12
2.7	Um modelo linear aplicado diretamente à entrada original não pode implementar a função XOR. Para isso é necessário transformar o espaço original usando uma função de ativação. Fonte: [5].	14
2.8	Efeitos de várias taxas de aprendizagem no treinamento. Fonte: [6]).	15
2.9	Política <i>exp_range</i> de <i>learning rate</i> cíclica. Fonte: [7].	16
2.10	Ilustração da operação de filtragem no domínio do espaço (convolução). w é o kernel do filtro e f é a área da imagem coberta pelo filtro. Fonte: [8]. .	18
2.11	Ilustração de um autoencoder.	19
2.12	Um <i>autoencoder</i> residual <i>fully-connected</i> . Esta figura mostra uma arquitetura com dois níveis (dois <i>autoencoders</i> empilhados). O primeiro nível codificada a imagem original. O resíduo da reconstrução é passado para o segundo nível. Cada nível produz 4 bits [9].	21
2.13	O autoencoder residual convolucional. A loss é aplicada nos resíduos [9]. .	21
3.1	Histograma da base de dados completa formada por 6,231,440 de patches. .	26

3.2	Histograma da BD0	27
3.3	Histograma da BD1	28
3.4	Histograma da BD2	29
3.5	Histograma da BD3	30
3.6	Histograma da BD4	31
3.7	Ilustração do <i>autoencoder</i> mais básico desenvolvido.	31
3.8	Ilustração do segundo modelo desenvolvido.	32
4.1	Imagen original (esquerda) e <i>patch</i> reconstruído pelo Modelo 2 (direita).	36
4.2	Comparação do Modelo 3 com o JPEG na métrica PSNR em diferentes taxas.	37
4.3	Comparação do Modelo 3 com o JPEG na métrica SSIM em diferentes taxas.	37
4.4	Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica PSNR em diferentes taxas para a base Kodak [2].	39
4.5	Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica SSIM em diferentes taxas para a base Kodak [2].	40
4.6	Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica MS-SSIM em diferentes taxas para a base Kodak [2].	40
4.7	Imagen Kodim05 [2].	41
4.8	Imagen jason-leem [10].	41
4.9	Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica PSNR em diferentes taxas para a imagem kodim05 [2].	42
4.10	Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica SSIM em diferentes taxas para a imagem kodim05 [2].	42
4.11	Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica MS-SSIM em diferentes taxas para a imagem kodim05 [2].	43
4.12	Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica PSNR em diferentes taxas para a imagem de [10].	43
4.13	Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica SSIM em diferentes taxas para a imagem de [10].	44
4.14	Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica MS-SSIM em diferentes taxas para a imagem de [10].	44
4.15	Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica PSNR em diferentes taxas para 47 imagens com muito conteúdo de alta frequência retiradas das bases [10] e [2].	45
4.16	Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica SSIM em diferentes taxas para 47 imagens com muito conteúdo de alta frequência retiradas das bases [10] e [2].	45

4.17	Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica MS-SSIM em diferentes taxas para 47 imagens com muito conteúdo de alta frequência retiradas das bases [10] e [2].	46
4.18	Ganho percentual médio na taxa por nível ao usar o codificador de entropia <i>gzip</i> nos <i>bitstreams</i> de cada nível para a base Kodak [2].	49
4.19	Imagen Kodim20 [2].	49
4.20	Ganho percentual na taxa por nível ao usar o codificador de entropia <i>gzip</i> nos <i>bitstreams</i> de cada nível para a Figura 4.7.	50
4.21	Ganho percentual na taxa por nível ao usar o codificador de entropia <i>gzip</i> nos <i>bitstreams</i> de cada nível para a Figura 4.19.	50

Listas de Tabelas

4.1	Tabela contendo médias obtidas pelo JPEG em cada uma das bases de teste utilizadas.	33
4.2	Tabela contendo o valor da PSNR, em decíbeis, dos testes do Modelo 1 com o uso do otimizador <i>Adam</i> e <i>learning rate</i> fixa. As linhas denotam a base de treino utilizada. As colunas denotam as bases de teste usadas para avaliação do modelo. O índice “todas” se refere ao uso de todas as imagens de todas as bases BD para treino.	34
4.3	Tabela contendo o valor da PSNR, em decíbeis, dos testes do Modelo 1 com o uso do otimizador <i>Adam</i> e <i>learning rate</i> fixa. As linhas denotam a base de treino utilizada. As colunas denotam as bases de teste usadas para avaliação do modelo. O índice “todas” se refere ao uso de todas as imagens de todas as bases BD para treino. O uso de $x + y$ denota o uso de todas as imagens do conjunto x e do conjunto y para treinamento.	34
4.4	Tabela contendo os resultados do Modelo 2 para as métricas visuais PSNR, SSIM e MS-SSIM a uma taxa nominal de 8 bits por pixel.	35
4.5	Tabela contendo os resultados do Modelo 3.	36
4.6	Tabela contendo os valores da taxa (BPP) e PSNR (decíbeis) do Modelo 4, JPEG e JPEG2000 para a base [2].	38
4.7	Tabela contendo os valores da taxa (BPP) e SSIM do Modelo 4, JPEG e JPEG2000 para a base [2].	38
4.8	Tabela contendo os valores da taxa (BPP) e MSSSIM do Modelo 4, JPEG e JPEG2000 para a base [2].	39
4.9	Tabela contendo os valores da taxa (BPP) e PSNR do Modelo 4, JPEG e JPEG2000 para 47 imagens da base [2] e [10] com muito conteúdo de alta frequência.	46
4.10	Tabela contendo os valores da taxa (BPP) e SSIM do Modelo 4, JPEG e JPEG2000 para 47 imagens da base [2] e [10] com muito conteúdo de alta frequência.	47

4.11 Tabela contendo os valores da taxa (BPP) e MSSSIM do Modelo 4, JPEG e JPEG2000 para 47 imagens da base [2] e [10] com muito conteúdo de alta frequência.	47
---	----

Lista de Abreviaturas e Siglas

AE Autoencoder.

BPP Bits por Pixel.

CLIC Challenge on Learned Image Compression.

CNN Convolutional Neural Network.

dB Decibéis.

DC Direct Current Coefficient.

DCT Discrete cosine transform.

DFT Discrete Fourier Transform.

DIV2K Diverse 2k high resolutonal quality images.

EYE UDH and HD images Eye tracking dataset.

FDCT Forward Discrete cosine transform.

GRU Gated Recurrent Unit.

IDCT Inverse Discrete cosine transform.

JFIF JPEG File Interchange Format.

LSTM Long Short-Term Memory.

MS-SSIM Multi-Scale Structural Similarity Index.

MSE Mean squared error.

PNG Portable Network Graphics.

PSNR Peak Signal-to-Noise Ratio.

ReLU Rectified Linear Unit.

RNN Recurrent Neural Network.

SGD Stochastic Gradient Descent.

SSIM Structural Similarity Index.

XOR Exclusive OR.

Capítulo 1

Introdução

Codificação e compressão de imagens é um grande desafio no campo de processamento de imagens. Para entender melhor a complexidade e a necessidade de se resolver esse desafio, este capítulo apresenta uma motivação do tema na seção 1.1; na seção 1.2 é apresentada a hipótese; os objetivos gerais e específicos são propostos na seção 1.3 e resultados esperados na seção 1.4.

1.1 Motivação

A codificação de dados é a transformação feita nos dados para atingir um certo objetivo, como compressão ou criptografia. O principal objetivo dos algoritmos de compressão é a redução do comprimento da mensagem (codificação da fonte), enquanto a criptografia tem como foco transformar os dados para proteger sigilo ou integridade daquilo que eles significam, e/ou acesso a tais dados, durante a sua transmissão através de um canal vulnerável.

Compressão de dados é o processo de codificar uma determinada informação utilizando uma menor representação. Os dois principais benefícios trazidos pela compressão de dados são o aumento significativo na capacidade de armazenamento de um sistema e menor largura de banda necessária para transmití-los.

De forma sucinta, compressão de dados é arte ou ciência de representar informação de forma compacta [11]. Nós criamos essas representações compactas identificando e usando estruturas que existem nos dados para que seja possível extrair redundância dos dados e descrevê-la em forma de um modelo que será usado como base para a codificação [11].

O desenvolvimento de algoritmos de compressão de dados podem ser divididos em duas fases [11]. A primeira fase é geralmente chamada de modelagem. Nesta fase, tentamos extrair informações sobre qualquer redundância existente nos dados e descrevemos a redundância na forma de um modelo. A segunda fase é chamada de codificação. Uma

descrição do modelo e uma “descrição” de como os dados diferem do modelo são codificados, geralmente usando um alfabeto binário. A diferença entre os dados e o modelo é frequentemente referida como resíduo.

Existem dois tipos de compressão: com perdas e sem perdas. A compressão com perdas (*lossy*) potencializa uma melhor taxa de compressão em troca de perda de informação enquanto na compressão sem perdas (*lossless*) não há perda de informação. Esta última é requerida em algumas aplicações, como sinais biomédicos. No contexto de imagens digitais, a compressão sem perdas permite que, após a codificação da imagem, a imagem decodificada seja idêntica à original, enquanto na compressão com perdas a imagem decodificada não é idêntica à original e há perda de qualidade visual.

Durante um processo de compressão devemos balancear dois pontos: a capacidade de compressão, isto é, o tamanho final em bits da imagem comprimida, e a distorção, que é a diferença entre a imagem original e a imagem reconstruída. Essa otimização é amiúde representada pela equação:

$$J = D(B) + \lambda R(B),$$

onde $D(B)$ representa a distorção entre a imagem original e $R(B)$ a quantidade de bits usada para representar a imagem, e λ é um parâmetro adimensional.

A princípio, uma maior quantidade de bits implicaria numa distorção $D(B)$ menor, mas resultaria numa taxa $R(B)$ maior. Contudo, esse é na verdade um problema intratável, como mostrado por [12]. Para uma dada taxa, existem diversas representações com distorções melhores e piores.

O motivo pelo qual precisamos de usar compressão de dados é porque estamos gerando e usando cada vez mais dados digitais. O número de *bytes* necessários para representar dados multimídia pode ser enorme. Por exemplo, para representar digitalmente 1 segundo de vídeo sem compressão usando o formato CCIR 601 [11], é necessário mais do que 20 *megabytes* para armazenar ou 160 megabits para transmitir [11]. Considerando o número de segundos em um filme, é fácil ver porque compressão é necessária à determinadas aplicações. Para serviços de streaming de mídia como Netflix, não usar compressão não é uma opção.

Os recursos necessários para armazenar e transmitir imagens são imensos, o que torna a sua compressão necessária. O objetivo em codificar uma imagem é representá-la com o menor número possível de bits, preservando a qualidade e a inteligibilidade necessárias à sua aplicação de modo a facilitar sua transmissão e armazenamento. São utilizadas medidas de desempenho para a codificação sem perdas e com perdas que diz respeito a taxa de compressão e distorção. Uma das formas de medir distorção comumente utilizada

em processamento de imagens é o MSE¹:

$$\frac{1}{n} \sum_{n=1}^N (x(n) - \hat{x}(n))^2, \text{ onde } x \text{ representa a imagem original e } \hat{x} \text{ a imagem decodificada}$$
(1.1)

Algoritmos de compressão de imagens aproveitam da percepção visual e propriedades estatísticas de dados da imagem para fornecer resultados superiores quando comparados com métodos de compressão de dados genéricos, que são usados para outros dados digitais. A tarefa de compressão de imagens foi cuidadosamente examinada durante anos por pesquisadores e times como o *Joint Pictures Experts Group* que desenvolveram os métodos de compressão de imagens JPEG [1] e JPEG2000 [13]. Mais recentemente, o algoritmo WebP [14] foi proposto para melhorar as taxas de compressão em imagens de alta resolução, que vem sendo cada vez mais utilizadas. O codec (codificador e decodificador) do estado da arte atual é o BPG [15].

Assim como os outros codecs, o *JPEG* explora as características imperfeitas da nossa percepção. Ele foi o primeiro padrão internacional de compressão para imagens monocromáticas e coloridas. Até hoje é um padrão bastante utilizado e possui métodos para compressão com perdas (método baseado em transformada discreta de cosseno) e sem perdas (método preditivo). Para criar um arquivo JPEG, primeiro a imagem é convertida para outro espaço de cor: o *YCbCr*. Este espaço, que é usado em vários vídeos de alta definição, codifica a cor de uma forma diferente do RGB, apesar de cobrir as mesmas cores. Os componentes Cb e Cr (crominância azul e vermelha, respectivamente) são altamente compressíveis, enquanto o componente de luminância indica quão brilhante o pixel é.

Na superfície da retina existem dois tipos de células que contêm pigmentos: os cones e os bastonetes. Os bastonetes existem em maior quantidade na periferia da retina e são estimulados com luz de baixa intensidade. Eles servem para dar um quadro geral do campo de visão e não estão envolvidos com a visão colorida (em baixos níveis de luz, nós praticamente não vemos cor, pois a iluminação é muito baixa para estimular os cones da nossa retina). Os cones, por sua vez, são muito sensíveis às cores e ocorrem principalmente na região central da retina. Seu estímulo depende de altas intensidades luminosas. É nessa região que a imagem é formada com maior nitidez, pois são estimulados pela luz mais intensa. Os cones são especializados na acuidade da visão diurna e em reconhecer a cor. O cérebro interpreta os sinais recebidos por esses cones, o que permite processar a diferenciação das cores. O olho humano é capaz de discriminar o brilho de uma imagem muito mais que sua informação de cor, visto que existem cerca de 120 milhões de bastonetes distribuídos sobre a superfície da retina contra apenas 6 à 7 milhões de

¹MSE também é bastante utilizada como função de loss para modelos de aprendizado profundo baseados em redes neurais

cones². Isto significa que os valores dos componentes de luminância precisam de muito mais fidelidade que os componentes de crominância (o mesmo vale para o componente de cor verde no espaço RGB, por exemplo).

Os algoritmos de compressão existentes atualmente podem estar longe de serem os ideais para os novos formatos de mídia como vídeos em 360 graus ou conteúdos de realidade virtual. Enquanto um desenvolvimento de um novo codec pode levar anos, um *framework* de compressão de imagens mais geral baseado em redes neurais pode ser capaz de se adaptar mais rápido à estas diferentes tarefas e ambientes.

Algoritmos padrão de compressão de imagens tendem a fazer suposições sobre a escala da imagem. Por exemplo, usualmente assume-se que um *patch* (pedaço retangular da imagem) de uma imagem natural de alta resolução irá conter muita informação redundante. De fato, quanto maior a resolução da imagem, mais provável que a maior parte dos *patches* que a compõem conterão informação de baixa frequência onde não há muita variação nos valores dos pixels. Esse fato é explorado pela maior parte dos codecs de imagens, de modo que eles tendem a ser muito eficientes em comprimir imagens de alta resolução. Entretanto, tais suposições são invalidadas ao criar miniaturas de imagens naturais de alta resolução, visto que um *patch* obtido de uma miniatura pode conter informação de alta frequência que é mais difícil de ser comprimida por estes algoritmos.

Nos últimos anos, redes neurais profundas se tornaram a base dos resultados do estado da arte para reconhecimento de imagens [16], detecção de objetos [17], reconstrução tridimensional de objetos [18], reconhecimento de faces [19], reconhecimento de discurso [20], *machine translation* [21], geração de legendas de imagens [22], tecnologia de carros autônomos [23], entre outros.

É natural buscar usar essa poderosa classe de métodos para melhorar a tarefa de compressão de imagens, especialmente para imagens que não são cuidadosamente desenvolvidas para codecs otimizados por heurísticas, como miniaturas. Considerando um codificador de imagem como um problema de análise/síntese com uma camada de gargalo no meio, é possível encontrar uma grande área de pesquisa que usa redes neurais para encontrar representações comprimidas. Muito desse trabalho se concentra em uma classe de redes neurais conhecida como *autoencoders* [24]. Alguns resultados já existentes para compressão com perdas usando autoencoders se mostraram promissores: [25, 26, 9], e redes neurais já atingiram o estado da arte para compressão sem perdas [27, 28].

²Esta diferença no número de bastonetes se dá por motivos evolutivos pois era mais importante identificar possíveis predadores ou presas durante a noite do que identificar cor

1.2 Hipótese

Compressão de imagens usando redes neurais tem sido uma área ativa de pesquisa em tempos recentes com vários desafios a serem enfrentados para que essas técnicas sejam competitivas com os codificadores clássicos. Serão verificadas as seguintes hipóteses:

- Deseja-se verificar se modelos baseados em *autoencoders* convolucionais são competitivos com os clássicos codecs *JPEG* e *JPEG2000* para imagens com muito conteúdo de alta frequência;
- Deseja-se verificar se usar imagens que codificadores clássicos têm dificuldade para comprimir para treinamento do modelo é benéfico para os resultados dele;
- Deseja-se verificar se o codificador de entropia utilizado no latente será capaz de comprimir em proporções semelhantes para todos os níveis de resíduos;
- Deseja-se verificar se os resultados dos *autoencoders* serão melhores ao trabalhar com um diferente domínio para os dados de entrada, um que beneficie o processo de aprendizagem via atualização de gradientes. Ou seja, se é benéfico para o desempenho do modelo aplicar uma transformação nos dados de entrada que evite problemas como gradientes muito pequenos;

Estas verificações serão dadas usando-se métricas visuais objetivas e as imagens serão trabalhadas à baixas taxas.

1.3 Objetivos

Deseja-se construir e avaliar o desempenho de um *framework* de compressão de imagens ponta a ponta (modelagem e codificação) usando *autoencoders* empilhados convolucionais recorrentes.

Para isto serão estudadas propostas de compressão de imagens na literatura, de modo a obter estatísticas e informações para guiar a escolha e implementação de codificadores baseados em *autoencoders* que estendam a estrutura básica de um *autoencoder*, gerando uma representação binária para a imagem ao quantizar a camada de gargalo ou as variáveis latentes correspondentes. Esta representação binária será ainda comprimida usando um codificador de entropia. Também serão construídas bases de dados próprias, a partir de bases de dados comumente utilizadas para esse problema, para que seja avaliado o desempenho, à diferentes taxas, de *autoencoders* convolucionais e dos codecs *JPEG* e *JPEG2000* nas mais variadas imagens.

1.4 Resultados Esperados

Acredita-se que será obtido desempenho superior aos métodos de compressão *JPEG* e *JPEG2000* para imagens com muito conteúdo de alta frequência, visto que estes métodos assumem que baixas frequências são mais visualmente e perceptualmente relevantes. Além disso, visto que existem trabalhos semelhantes na literatura acredita-se que será obtido melhor desempenho com o uso de imagens mais difíceis de comprimir para treinamento e com o uso de modelos recorrentes.

Capítulo 2

Trabalhos similares e conceitos fundamentais

Este capítulo descreve trabalhos similares e conceitos fundamentais de codecs clássicos e arquiteturas baseadas em redes neurais utilizadas neste trabalho.

2.1 JPEG

Arquivos JPEG normalmente são descritos no formato JFIF [29], que é uma limitação do padrão JPEG completo. O método de compressão JPEG descrito na Figura ??, assim como os outros métodos, exploram a imperfeição do sistema visual humano. Os 5 passos usados para codificação no padrão JFIF são descritos nas subseções seguintes.

2.1.1 Conversão do espaço de cor

Para o padrão JFIF¹, primeiro, é feita a conversão do espaço RGB da imagem de entrada para o espaço de cor $YCbCr$. Os valores dos pixels estarão no intervalo de 0 à 255 (mesmo do RGB). Uma vez que o espaço de cor é transformado para $YCbCr$, é necessário decidir qual será o fator usado para reduzir a quantidade de pixels nos componentes de crominância, visto que o olho humano é muito mais sensível ao brilho do que a cor. Normalmente é usado um fator de 2 nas duas direções, o que dá 4 vezes menos cor (para cada 4 pixels Y só terá 1 pixel Cb e 1 Cr). Esse fator é determinado pelo argumento *quality* passado como parâmetro para codificação da imagem (com *quality* máxima, não haverá redução e a imagem possuirá a mesma resolução de cor).

¹JPEG permite você usar qualquer espaço de cor que queira, mas a maior parte das pessoas seguiram com o padrão JFIF por praticidade

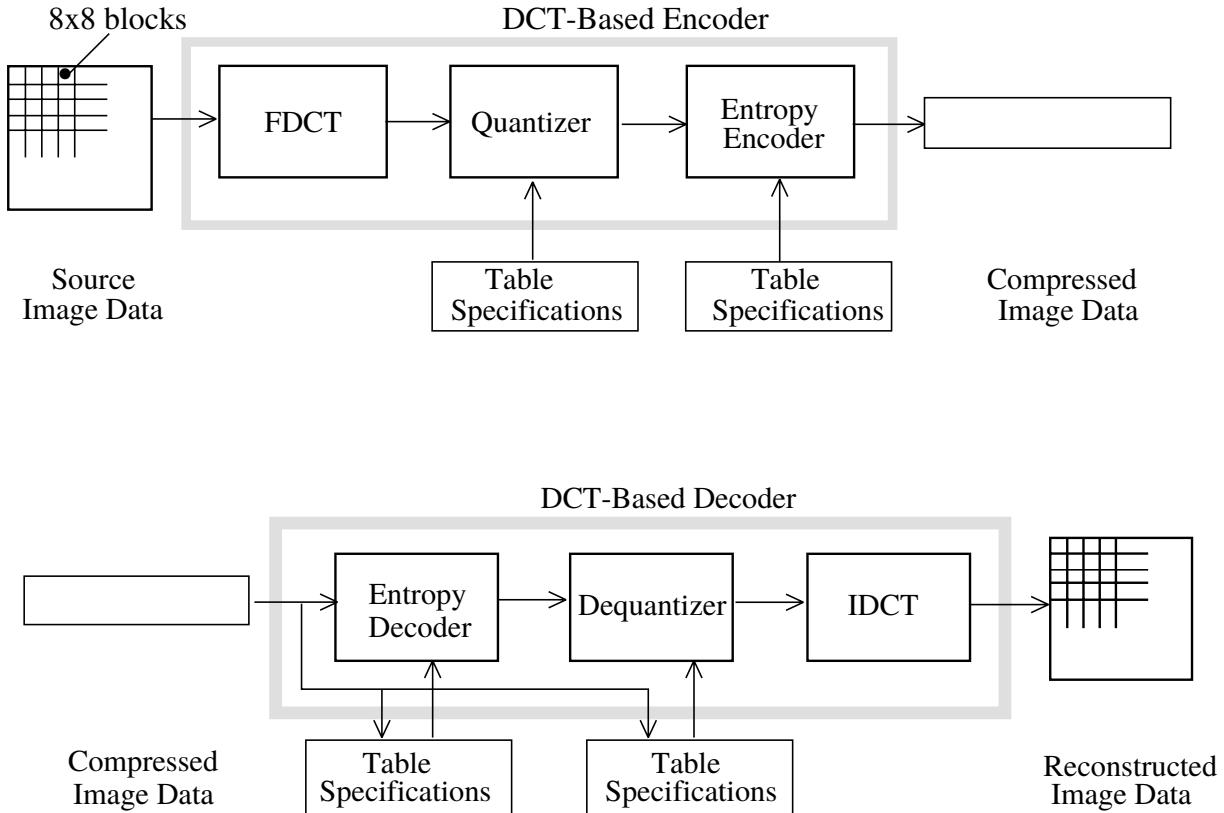


Figura 2.1: Diagram do método de compressão JPEG. Fonte: [1].

Nota-se que na Figura 2.2 praticamente não há diferença visual olhando a um nível normal de zoom, mesmo a codificação da imagem da direita possuindo 64 vezes menos cor do que a imagem da esquerda. Entretanto, olhando a Figura 2.3 é possível notar certa discrepância nas bordas da arara vermelha.

2.1.2 Aplicação da transformada direta de cossenos

Aqui é feita a divisão da imagem em blocos com 8 pixels de largura e 8 de altura que serão convertidos em uma nova matriz com o auxílio de uma transformada discreta de cossenos DCT [30] (*JPEG* sempre usa *DCT-II*). Essa transformação, que é similar a transformada de *Fourier*, analisa as frequências dos valores originais dos pixels da imagem ao longo de cada linha e coluna usando um conjunto de ondas cossenos oscilando em diferentes frequências e amplitudes. Cada um dos blocos serão codificados separadamente com sua própria transformada discreta de cossenos e podem ser exatamente replicados por ondas de cossenos 8 por 8, onde varia-se frequências e amplitudes de cada uma delas.

A representação de um sinal em DCT tende a ter maior parte da sua energia concentrada em um número menor de coeficientes quando comparado com outras transformações



Figura 2.2: Imagem original sem compressão retirada de [2] e imagem (direita), gerada a partir da imagem original, com dimensionalidade nos canais de crominância reduzida por um fator de 8 nas duas direções. Pode-se jogar informação da imagem original fora e então usar a imagem da direita para armazenamento ou transmissão, que terá novos valores em seus canais de cores (aumenta-se a dimensionalidade apenas quando for necessário exibi-la). Fonte: [2].



Figura 2.3: Versão com zoom das imagens mostradas na Figura 2.2. A região mostrada é exatamente a mesma (espacialmente) para as duas imagens, com a mesma quantidade de pixels. Imagem original sem compressão (esquerda) e imagem com 64 vezes menos cor (direita).

como a transformação discreta de *Fourier*, o que é mais desejável para algoritmos de compressão. Além disso, a DFT tem coeficientes complexos, o que dobra a representação. Definimos energia de um sinal como sendo:

$$\sum_{k=-\infty}^{\infty} |x_k|^2 \quad (2.1)$$

A Figura 2.4 mostra as 64 funções de cosseno que podem ser combinadas para formar qualquer imagem 8 por 8. Nota-se que a partir do bloco superior esquerdo, as frequências

das ondas de cosseno crescem tanto na direção horizontal quanto na vertical. Além disso, o bloco inferior direito constituído de um padrão de xadrez, é o que possui maior frequência. Para criar qualquer imagem 8 por 8, basta combinar todos esses blocos ao mesmo tempo. Cada um será ponderado baseado em um número denominado coeficiente que representará a contribuição de cada um desses blocos individuais para o todo. Assim, se a contribuição de um bloco for zero não haverá nenhuma parte desta função de cossenos na imagem 8 por 8 buscada.

Basicamente, na transformada discreta de cossenos é calculado os coeficientes das ondas de cossenos. Os coeficientes podem ser considerados como a quantidade relativa das frequências espaciais 2D contidas no sinal de entrada. O coeficiente com frequência zero nas duas dimensões é chamado de coeficiente corrente direto (DC) e os 63 restantes são chamados de coeficientes correntes alternados (AC). O decodificador reverte este passo usando a função inversa da DCT (IDCT) que pega os 64 coeficientes *Forward DCT* (FDCT) do codificador já quantizados e reconstrói o sinal da imagem de 64 pontos somando os sinais base. Se a FDCT e a IDCT pudessem ser computadas com acurácia perfeita e se os coeficientes da DCT não fossem quantizados no codificador, o sinal original de 64 pontos poderia ser exatamente recuperado.

Mudanças de altas frequências podem ser minimizadas ou zeradas, visto que nós não percebemos suas mudanças na imagem tão bem quanto componentes de baixas frequência. Ou seja, blocos de imagens cujos valores de pixels mudam de intensidade muito rápido (normalmente perto das bordas da imagem) podem ser borrados sem perda significativa de qualidade visual, o que economiza uma quantidade enorme de espaço. Por isso, os coeficientes das ondas de cossenos de alta frequência não contribuem muito para a imagem final. Entretanto, isso não é verdade para textos, o que faz com que o JPEG não seja uma boa escolha quando o objetivo é comprimir imagens de texto, conforme mostra a Figura 2.5.

2.1.3 Quantização

Quantização é definida dividindo cada coeficiente pelo passo do quantizador especificado na tabela de quantização, seguido por um arredondamento para o inteiro mais próximo. Na dequantização é feito o processo inverso multiplicando pelo passo do quantizador, com o auxílio da mesma tabela usada para quantização. Quantização é o passo em que há a maior perda de informação na imagem. É dada pela seguinte equação:

$$\left[F^Q(u, v) = \frac{F(u, v)}{Q(u, v)} \right] \quad (2.2)$$

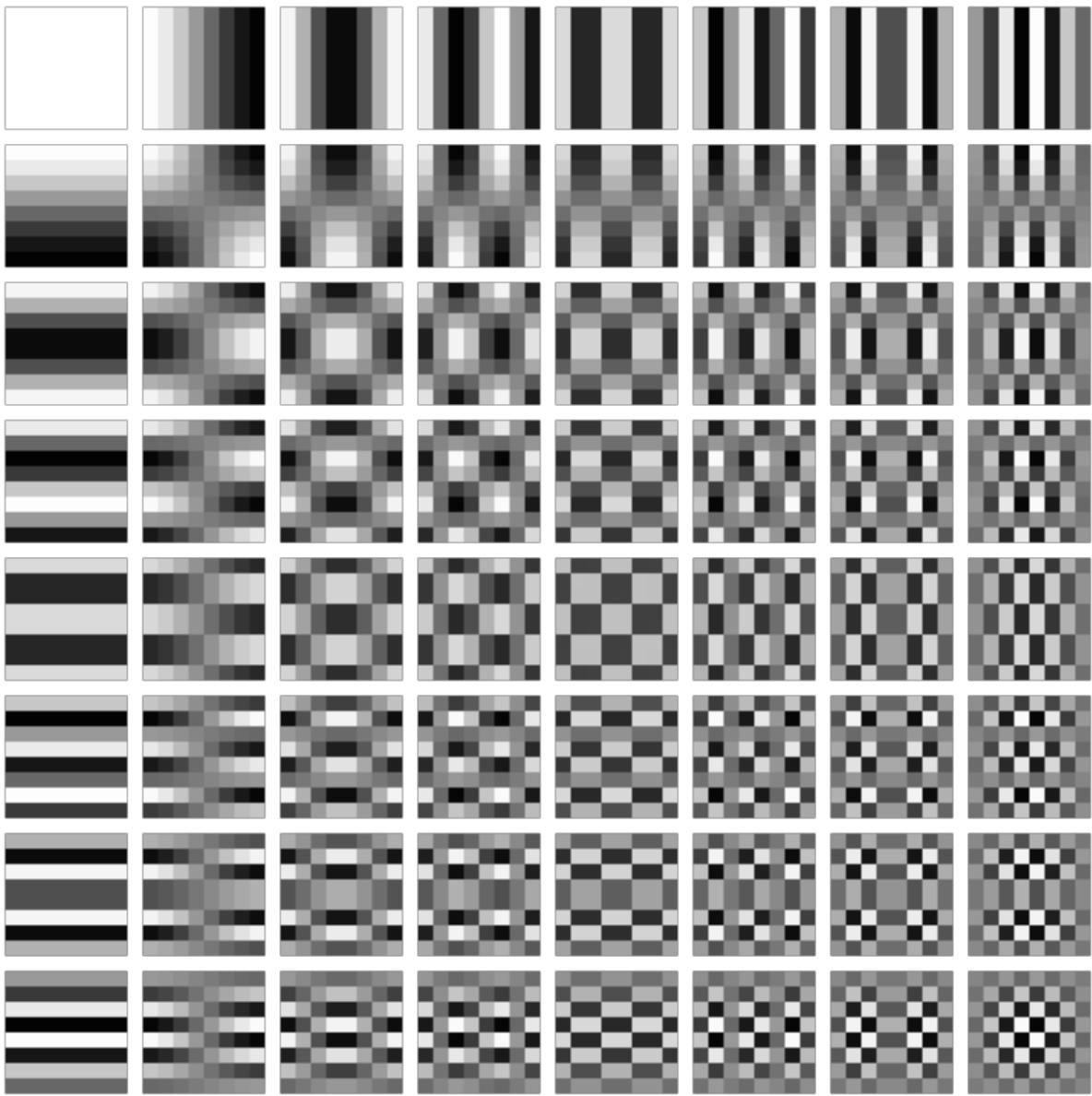


Figura 2.4: 64 (8 por 8) ondas base de cossenos com frequências variadas. Fonte: [3].

$F^Q(u, v)$ será o novo valor do coeficiente, $F(u, v)$ é o valor atual e $Q(u, v)$ é o passo de quantização que controla a qualidade da imagem definido para a aplicação (altas frequências são removidas usando um valor maior para $Q(u, v)$). Cada um dos 64 coeficientes DCT são uniformemente quantizados em conjunto com uma tabela de quantização de 64 elementos, que é definida pelo nível de qualidade escolhido para a aplicação. O propósito da quantização é alcançar mais compressão ao representar os coeficientes DCT com a menor precisão possível para alcançar a qualidade da imagem especificada. Isso é feito descartando informação que não é visualmente significante.



Figura 2.5: Comparação imagem de texto com e sem compressão. Fonte: [4].

Após a quantização, os coeficientes DC são tratados separadamente devido à alta correlação destes coeficientes em blocos 8 por 8 adjacentes da imagem, considerando que eles geralmente possuem maior valor e muito impacto na imagem. Assim, eles são codificados como a diferença do coeficiente DC do bloco anterior na ordem de codificação. Por fim, todos os coeficientes quantizados são ordenados em uma sequência zig-zag conforme mostra a Figura 2.6 com o objetivo de facilitar a codificação (que será usada no próximo passo) ao ordenar os coeficientes de frequência similares próximos uns dos outros.

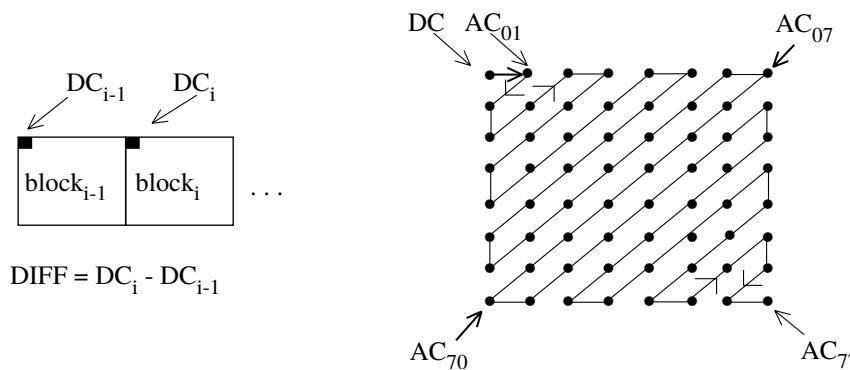


Figura 2.6: Sequência zig-zague usada para melhorar codificação. Fonte: [1].

Vale ressaltar que como existe apenas um controle de qualidade pelas matrizes de quantização e não uma otimização taxa-distorção, algumas vezes um ponto com quantização menor, com mais bits necessários para a representação e mais precisão em relação ao valor original, pode levar a um ponto que pode ser pior do ponto de vista de compressão.

2.1.4 Codificação

O passo final do codificador é a codificação de entropia responsável por gerar o *bitstream* comprimido que representará a imagem. Este passo permite compressão sem perdas adicional ao codificar os coeficientes DCT quantizados de forma mais compacta baseando-se em suas características estatísticas. O *JPEG* propõe o uso de dois métodos de codificação de entropia: Codificador de Huffman [31] e Codificador Aritmético [32].

No final, para cada bloco 8 por 8 da imagem original, teremos três matrizes 8 por 8 quantizadas, onde as matrizes correspondentes aos canais Cb e Cr serão as mais comprimidas.

Terminada a codificação, o papel do decodificador será de reverter os passos do codificador para reconstruir a imagem, conforme mostra a Figura ??: primeiro, a matriz quantizada será obtida decodificando os blocos comprimidos. Depois, é possível obter a matriz DCT multiplicando a matriz quantizada pela matriz de quantização utilizada pelo codificador. Depois, essa matriz é transformada usando a IDCT que resultará na matriz no espaço de cor $YCbCr$.

O *JPEG2000* tem um funcionamento geral similar ao *JPEG*, entretanto é utilizado uma transformada discreta *wavelet* no lugar da DCT. Esta transformada é aplicada na imagem inteira, o que elimina o efeito de *blocking* causado pelo *JPEG* mas causa o efeito de *ring*.

2.2 Redes Neurais

Um algoritmo simples de aprendizado de máquina (*machine learning*) chamado regressão logística pode determinar quando recomendar cesariana para uma paciente [33]. O desempenho desses algoritmos dependem muito da representação dos dados fornecidos. Cada pedaço de informação incluída na representação dos dados é conhecida como *feature* (característica) [5]. Algumas tarefas em inteligência artificial podem ser resolvidas desenvolvendo características a serem extraídas dos dados. Entretanto, para muitas tarefas é difícil saber quais *features* devem ser extraídas. Uma solução para esse problema é descobrir não apenas a função que mapeará a entrada para a saída mas também a própria representação. Essa abordagem é conhecida como *representation learning* (aprendizado de representações) [5].

Deep learning (aprendizado profundo) resolve o problema de *representation learning* introduzindo representações que são expressadas em termos de outras representações mais simples [5]. Por exemplo, usando um modelo de *deep learning* é possível representar o conceito de uma imagem de um carro combinando conceitos mais simples, como bordas e contornos.

Redes neurais são modelos de *deep learning* capazes de fazer previsões aprendendo uma função que relaciona as características dos dados à respostas observadas/desejadas. Redes neurais são consideradas aproximadores universais de funções, o que significa que elas podem computar e são capazes de aproximar qualquer função (não só lineares) [34]. Para isso, é necessário que elas sejam profundas o suficiente e possuam funções de ativações (funções não-lineares), visto que a saída de uma rede sem funções de ativação seria apenas uma função linear (polinômio de grau um) que não é capaz de representar algumas funções como a função XOR [Figura 2.7]. As ativações permitem que o modelo aprenda funções mais complexas, ao introduzir transformações não-lineares nas saídas das camadas. A *Rectified Linear Unit* (ReLU), definida como $f(x) = \max(0, x)$, é uma das funções de ativações não-lineares mais comuns e recomendadas pois ela é quase linear, o que faz com que o modelo seja fácil de otimizar com métodos comumente usados como descida de gradiente [35] (método que atualiza os pesos com base no gradiente, de modo que a função de erro será minimizada dando passos proporcional ao negativo do gradiente em direção ao ponto mínimo) e preserva várias propriedades existentes em modelos lineares que permitem que eles generalizem bem [5].

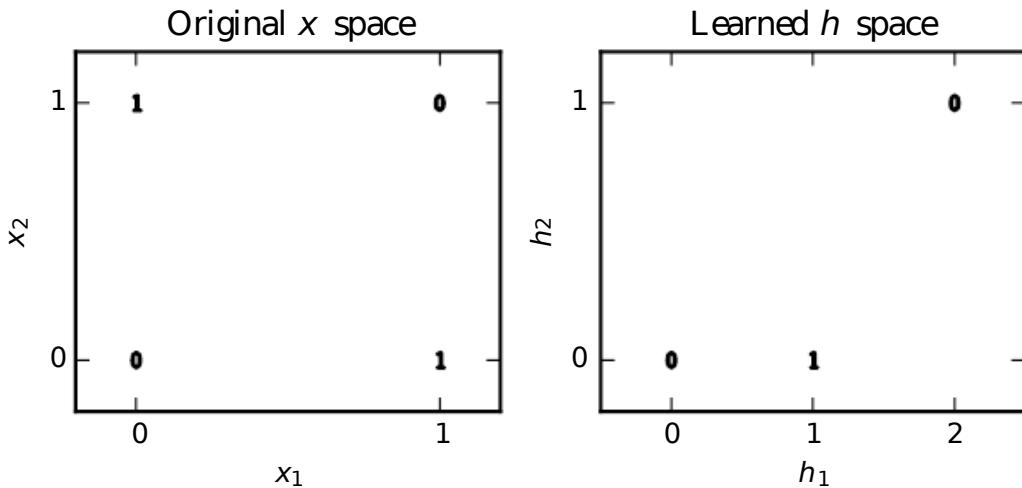


Figura 2.7: Um modelo linear aplicado diretamente à entrada original não pode implementar a função XOR. Para isso é necessário transformar o espaço original usando uma função de ativação. Fonte: [5].

2.2.1 Taxa de aprendizagem (*learning rate*)

A *learning rate* é um hiperparâmetro que controla o quanto nós ajustamos os parâmetros aprendíveis da nossa rede com respeito ao gradiente. Quanto menor o valor, menor o

passo dado ao longo do declive em direção ao mínimo da função de custo. A *learning rate* é um dos hiperparâmetros que devem ser escolhidos com cuidado, pois ela pode ter uma grande influência na convergência do seu modelo. O gráfico Figura 2.3 mostra os diferentes cenários de *learning rate* e como ela afeta o treinamento. Leslie N. Smith

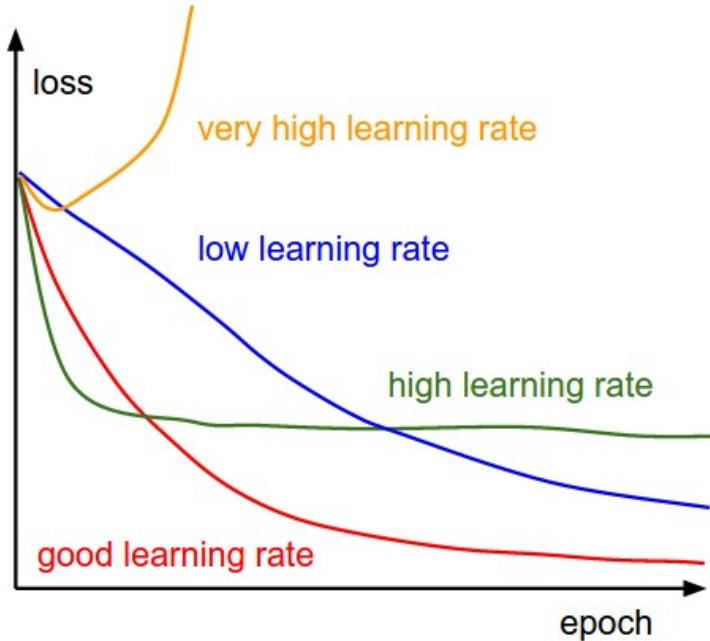


Figura 2.8: Efeitos de várias taxas de aprendizagem no treinamento. Fonte: [6]).

argumenta em [7] que é benéfico para a aprendizagem variar a *learning rate* de forma cíclica durante o treinamento para evitar cair em mínimos locais não ótimos (pontos de sela). Também é mostrado que é possível atingir resultados iguais ou superiores com menos iterações de treinamento usando este método quando comparado com uma rede que foi treinada com *learning rate* fixa ou usando outro método padrão de variação (este fenômeno ficou conhecido como “superconvergência”).

Leslie propõe ciclos para variar a *learning rate*. Um ciclo é definido como o número de iterações necessárias para *learning rate* ir do valor mínimo até o máximo definido no ciclo e voltar ao mínimo. Dadas as constantes *baselr*, *maxlr*, *step* e γ que representam, respectivamente, *learning rate* inicial, *learning rate* máxima, número de iterações correspondente à metade de um ciclo e constante responsável por diminuir limite superior do ciclo; e as variáveis *itr* que representa a iteração atual no treinamento e $cycle = \left\lfloor 1 + \frac{itr}{2 \cdot step} \right\rfloor$ o ciclo atual, a *learning rate* *lr* para uma *itr* qualquer na política *exp_range* [Figura 2.9], é calculada pela seguinte equação:

$$lr = baselr + (maxlr - baselr)max(0, 1 - |itr/step - 2cycle + 1|)\gamma^{itr}. \quad (2.3)$$

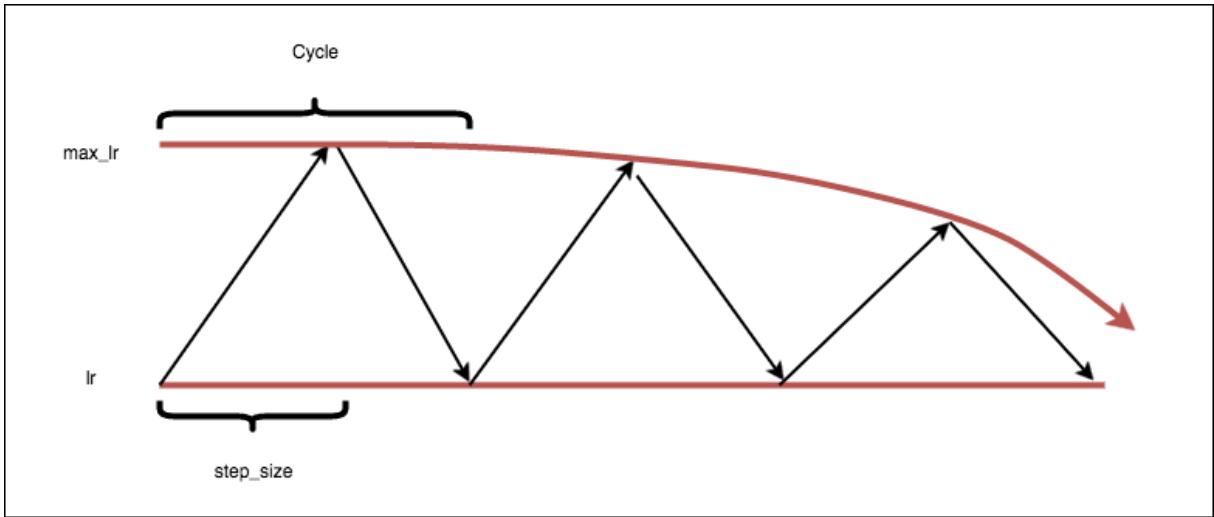


Figura 2.9: Política *exp_range* de *learning rate* cíclica. Fonte: [7].

O objetivo do treinamento das redes neurais é minimizar a função de custo/erro (*loss*) definida para a aplicação. Pode-se usar um conjunto de validação para salvar e usar os pesos do modelo que obter a melhor métrica no conjunto de validação, pois o modelo que obtiver a menor função de erro no treinamento não necessariamente será aquele que obterá a melhor métrica no conjunto de teste usado para avaliação do modelo. Para que este objetivo seja atingido, normalmente são usados otimizadores. O método clássico de descida do gradiente estocástico (SGD) consiste basicamente em usar amostras aleatórias (*mini-batch* para aproximar o verdadeiro gradiente que levará à minimização da função de custo/erro escolhida. O SGD mantém uma única *learning rate* (não muda durante o treino) para todas as atualizações de peso. O *Adam* [36] é um otimizador que adapta a *learning rate* baseando-se na média do primeiro momento e do segundo momento dos gradientes e na média móvel exponencial do gradiente e da raiz quadrada dele. Os parâmetros utilizados neste otimizador controlam as taxas de decaimento destas médias móveis.

2.2.2 Redes Neurais Convolucionais

Em *deep learning* (aprendizado profundo), uma rede neural convolucional (CNN) [37] é uma classe de redes neurais profundas que usa convoluções para detectarem características e padrões presentes nas imagens. As primeiras camadas detectam características que podem ser reconhecidas e interpretadas de maneira relativamente fácil. Camadas posteriores detectam características mais abstratas e usualmente presentes em muitas das características detectadas por camadas anteriores. A arquitetura de uma CNN é análoga

ao padrão de conectividade dos neurônios no cérebro humano e foi inspirada pela organização do córtex visual. Neurônios individuais respondem à estímulos apenas em uma região restrita do campo visual que é conhecida como campo receptivo. Uma coleção de sobreposição desses campos cobrem toda a área visual [38].

Uma CNN é capaz de capturar dependências espaciais e temporais na imagem de forma bem-sucedida através da aplicação de filtros relevantes e dispensa a necessidade de engenharia de características. A arquitetura realiza um melhor ajuste ao conjunto de imagens devido à redução do número de parâmetros envolvidos e a reusabilidade dos pesos. Em outras palavras, a rede pode ser treinada para entender melhor a complexidade da imagem e suas características relevantes. Esta extração de características é feita por meio da aplicação de filtros² no domínio do espaço denominados convoluções. Para uma máscara (filtro) de tamanho $m \times n$, $m = 2a+1$ e $n = 2b+1$, onde a e b são inteiros positivos. Para qualquer ponto (x, y) na imagem f , a resposta do filtro é a soma dos produtos dos coeficientes do filtro e dos pixels da imagem englobados pelo filtro (Equação 2.4).

$$g(x, y) = \sum_{s=-a}^a \sum_{t=-b}^b w(s, t) f(x + s, y + t) \quad (2.4)$$

Conforme mostra a Figura 2.10, x e y variam de modo que cada pixel em w visite todos os pixels em f .

2.2.3 Autoencoders

Um Autoencoder (AE) é uma classe de redes neurais que são formados por duas redes conectadas: um **encoder** e um **decoder**.

- O **encoder** tem como função converter a informação da entrada em uma representação menor e mais densa chamada de espaço latente. Pode ser representado como uma função de x , $f(x) = h$.
- O **decoder**, por sua vez, tenta reconstruir a informação original, passando do espaço latente criado pelo encoder para o espaço original da informação. Pode ser representado como uma função de h , $g(h) = \hat{x}$

Uma rede do tipo AE é treinada de forma não supervisionada e pode ser descrita como $g(f(x)) = \hat{x}$. Normalmente, o objetivo é apenas diminuir a diferença entre x e \hat{x} (nesse caso, a função de custo a ser minimizada normalmente é a $MSE(x, \hat{x})$ [Equação 1.1]). A camada entre o *encoder* e o *decoder* que contém menos neurônios [Figura 2.11] e força o

²Filtro é um termo que vem de processamento no domínio da frequência e se refere a aceitar ou rejeitar certos componentes de frequência

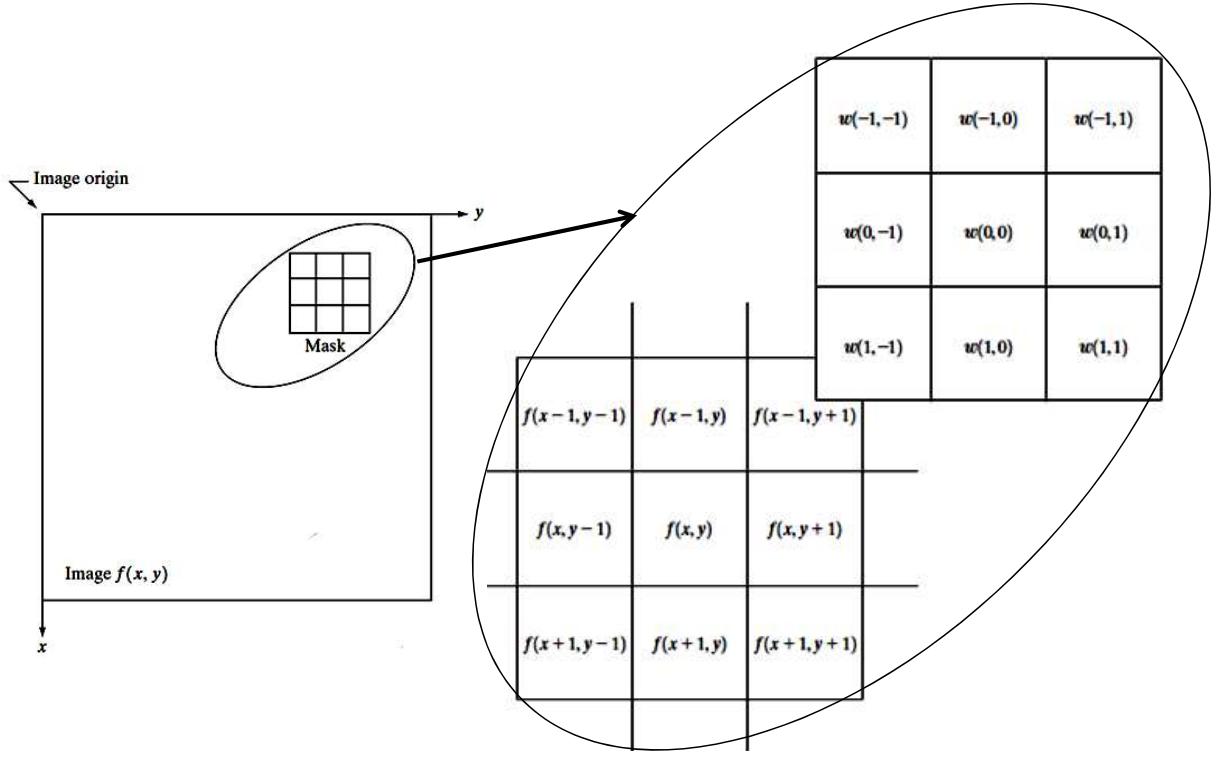


Figura 2.10: Ilustração da operação de filtragem no domínio do espaço (convolução). w é o kernel do filtro e f é a área da imagem coberta pelo filtro. Fonte: [8].

encoder a comprimir informação da representação original da entrada original gerando o espaço latente é denominada *bottleneck* (camada de gargalo).

Para o problema de compressão de imagens, normalmente usa-se um binarizador na camada de gargalo com o objetivo de binarizar o latente gerado pelo *encoder*. Assim, o *encoder* é forçado a comprimir informação e o *decoder* a diminuir a distorção usando menos informação. O binarizador transforma os valores em ponto flutuante (representação limitada dos números reais no computador) em inteiros que serão binarizados. Ele é necessário pois os códigos práticos precisam ter entropia finita, logo valores contínuos precisam ser quantizados para um conjunto finito de valores discretos, o que introduz erro. Com a binarização também é possível reduzir o espaço consumido pela imagem codificada, visto que números em pontos flutuante com precisão simples ocupam 32 bits o que levaria a uma alta taxa de bits por pixel. A avaliação dos modelos usados para este tipo de problema é dada considerando não só a taxa, mas também métricas visuais que calculam o nível de distorção da imagem reconstruída. Existem três tipos de métricas visuais comumente utilizadas:

1. *PSNR* definida por

$$20 \cdot \left(\frac{\text{MAX}_I}{\sqrt{\text{MSE}}} \right), \quad (2.5)$$

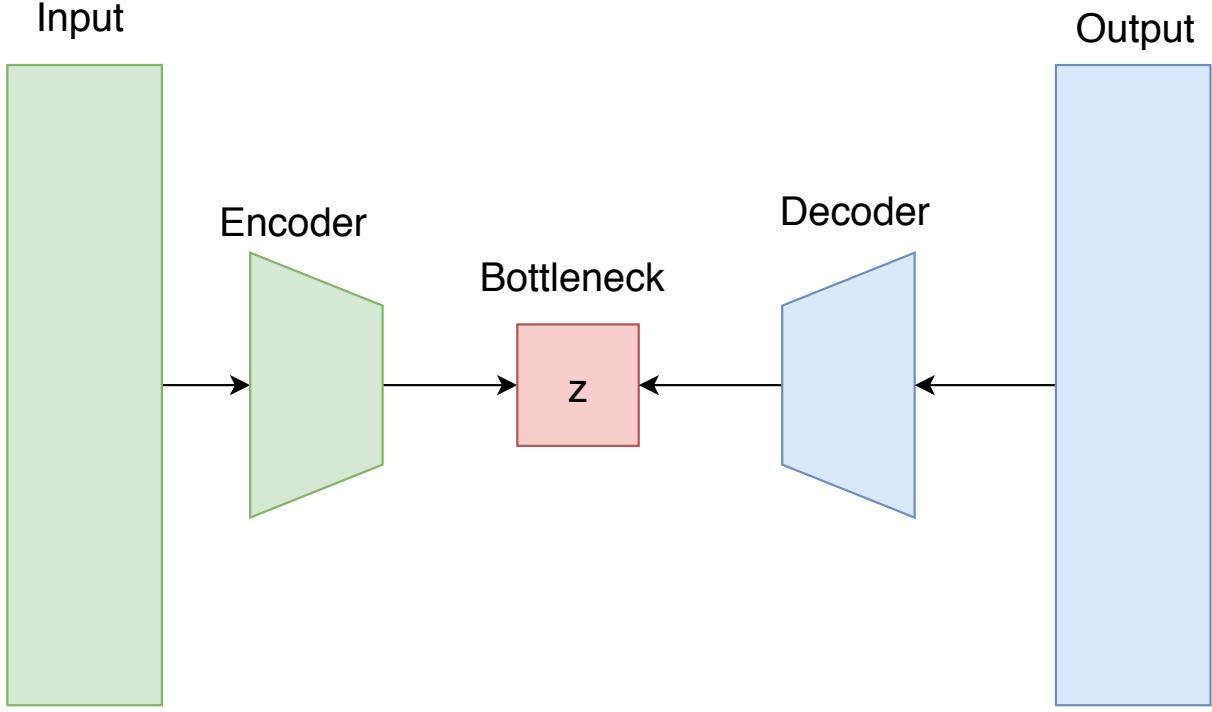


Figura 2.11: Ilustração de um autoencoder.

onde MAX indica o maior valor possível para o pixel em uma imagem. Quando estes são representados em bits, usa-se $MAX_I = 2^B - 1$;

2. *SSIM* [39]. Seja $x = \{x_i | i = 1, 2, \dots, N\}$ e $y = \{y_i | i = 1, 2, \dots, N\}$ dois sinais discretos não negativos e μ_x, σ_x^2 e σ_{xy} serem a média de x , a variância de x e a covariância de x e y , respectivamente. A SSIM é dada por:

$$SSIM(x, y) = \frac{(2\mu_x\mu_y + C_1)(2\mu_x\mu_y + C_2)}{(\mu_x^2 + \mu_y^2 + C_1)(\mu_x^2 + \mu_y^2 + C_2)}, \quad (2.6)$$

onde $C_1 = (K_1 L)^2$, $C_2 = (K_2 L)^2$ e $C_3 = C_2/2$. L é o intervalo dinâmico dos valores dos pixels ($L = 255$ para 8 bits por pixel) e K_1, K_2 são constantes.

3. *MS-SSIM* [40] é baseada na SSIM:

$$MSSSIM(x, y) = [l_M(x, y)]^{\alpha M} \cdot \prod_{j=1}^M [c_j(x, y)]^{\beta_j} [s_j(x, y)]^{\gamma_j}, \quad (2.7)$$

onde os expoentes são usados para ajustar a importância relativa de cada componente e l , c e s são componentes de luminância, contraste e estrutura, respectivamente.

Compressão de Imagens com Taxa Variável usando Redes Neurais Recorrentes

Nesse trabalho [9] é feita uma abordagem de ponta a ponta (modelagem e codificação) para compressão de imagens. Esta abordagem é baseada em várias redes neurais empilhadas, especificamente, uma pilha de *autoencoders*, o que possibilita a transmissão de informação incremental. Cada *autoencoder* comprime sua entrada e tenta reconstruí-la. O erro residual é propagado para o próximo *autoencoder* que, novamente, tenta reconstruir seu resíduo de entrada. Considerando E como o *encoder*, D como o *decoder*, B como a função de binarização, e x como a entrada, um *autoencoder* pode ser representado como

$$\hat{x} = D(B(E(x))). \quad (2.8)$$

Esta equação pode ser usada para compor uma pilha de *autoencoders* residuais pelo seguinte conjunto de equações:

$$\begin{aligned} b_t &= B(E_t(r_{t-1})), \hat{x}_t = D_t(b_t) + \gamma \hat{x}_{t-1}, \\ r_t &= x - \hat{x}_t, r_0 = x, \hat{x}_0 = 0 \end{aligned} \quad (2.9)$$

onde $t \in \{1 \dots n\}$, para um modelo com n níveis. O par (E_t, D_t) é o *autoencoder* para o t -ésimo nível com r_{t-1} como entrada. A entrada inicial é a imagem original. O binarizador B , nesta formulação, é o mesmo para todos os níveis de iterações. Usa-se $\gamma = 1$ para a reconstrução aditiva, de modo que a reconstrução final para um modelo com t níveis será igual a \hat{x}_t , que será a soma das saídas de todos os níveis.

Redes Neurais não-recorrentes

Toderici et al. [9] primeiramente propôs um esquema de compressão com esta estratégia usando diferentes arquiteturas de redes neurais. Este esquema é baseado em um framework de codificação aditivo que restringe o número de bits de codificação. A Figura 2.12 mostra um exemplo deste empilhamento de *autoencoders* completamente convolucionais. Para esta arquitetura, cada rede é composta de um *encoder*, que produz uma representação latente de 8 bits e um *decoder* que tenta reconstruir a entrada a partir do latente. Cada camada possui 512 *kernels* com tangente hiperbólica como função de ativação. Os resíduos serão as entradas para as próximas redes. Considerando 16 níveis, um total de 128 bits é usado, o que dá uma taxa de bits por pixel de $\frac{2 \cdot 2 \cdot 32}{32 \cdot 32} = 0.125$ para representar a imagem de entrada de tamanho 32x32.

Uma parte interessante deste trabalho reside na binarização B (no intervalo $[-1, 1]$) dos valores em pontos flutuantes da *bottleneck*. A derivada da função chão é zero em todos os lugares, exceto nos inteiros, onde ela é indefinida. Por isso, a derivada na passagem de volta (*backward pass*) usada no algoritmo de backpropagation [35] é substituída com a

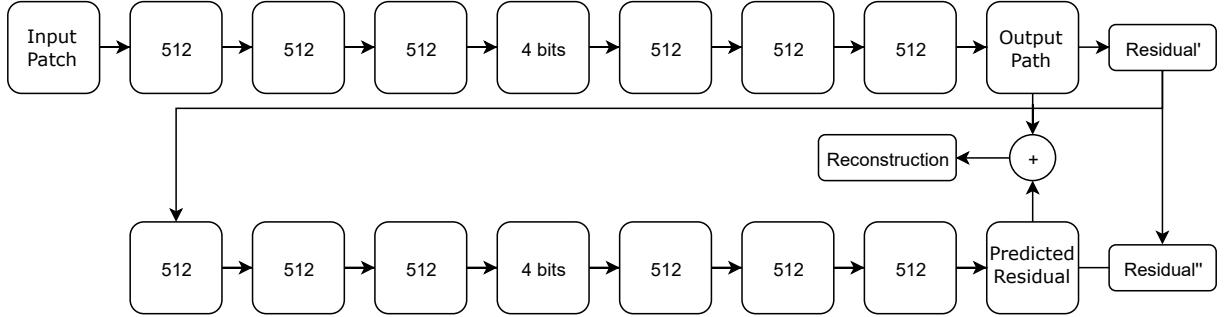


Figura 2.12: Um *autoencoder* residual *fully-connected*. Esta figura mostra uma arquitetura com dois níveis (dois *autoencoders* empilhados). O primeiro nível codificada a imagem original. O resíduo da reconstrução é passado para o segundo nível. Cada nível produz 4 bits [9].

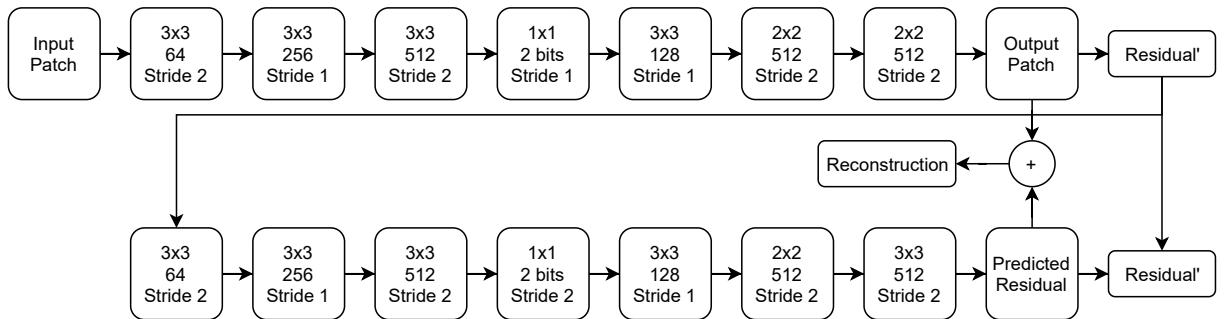


Figura 2.13: O *autoencoder* residual convolucional. A loss é aplicada nos resíduos [9].

derivada da esperança:

$$\frac{d}{dy}\{y\} := \frac{d}{dy}E[\{y\}] = \frac{d}{dy} = 1. \quad (2.10)$$

A função de arredondamento também é não-diferenciável, por isso na *backward pass* os gradientes são passados sem modificação do *decoder* para o encoder *encoder*. A quantização é realizada normalmente na *forward pass*, visto que substituir o arredondamento por uma aproximação completamente pode levar o *decoder* a aprender à inverter essa aproximação, removendo a informação da *bottleneck* (camada de gargalo) que força a rede a comprimir informação.

O trabalho [9] propõe uma função de binarização em dois passos. O primeiro passo é geração de valores no intervalo $[-1, 1]$. A segunda parte envolve converter esses valores para o conjunto $-1, 1$. Para este propósito, uma camada completamente convolucional com ativações de tangente hiperbólica são usadas para produzir as saídas no intervalo desejado. Em seguida, uma abordagem estocástica é aplicada. Estes passos podem ser

representados pelo seguinte conjunto de equações:

$$b(x) = x + \epsilon \in \{-1, 1\}, \quad (2.11)$$

$$\epsilon \sim \begin{cases} 1 - x & \text{com probabilidade } \frac{1+x}{2}, \\ -1 - x & \text{com probabilidade } \frac{1-x}{2}, \end{cases} \quad (2.12)$$

onde ϵ , corresponde ao ruído de quantização.

O encoder pode ser summarizado pela seguinte equação:

$$B(x) = b(\tanh(W^{bin}x + b^{bin})). \quad (2.13)$$

W^{bin} e b^{bin} são os pesos e bias padrões da rede. Esta formulação para a binarização é usada em todos os modelos para a *forward pass* da rede. Para a *backward pass* da *back-propagation*, é usada a derivada da esperança. Visto que a esperança de $b(x)$ será igual a x para todo x , os gradientes serão passados por b sem mudanças. Essa binarização é aplicada somente em tempo de treino. Em tempo de teste, b é substituída por

$$b^{inf}(x) = \begin{cases} -1, & \text{se } x < 0, \\ 1, & \text{caso contrário.} \end{cases} \quad (2.14)$$

Os autores propõe uma versão iterativa desta rede, mas usando convoluções em vez de redes *fully-connected*. A Figura 2.13 ilustra essa arquitetura. Neste caso, é utilizado 2 bits por pixel na *bottleneck* binarizada.

O conjunto de dados é composto de 216 milhões de images aleatórias coletadas da internet. As imagens escolhidas são redimensionadas para terem tamanho 32 por 32 e salvas sem perda de informação como arquivos PNG. Então, 90% destas imagens são usadas para treinamento, enquanto 10% são usadas para teste. Esse conjunto de dados simula o cenário de compressão de miniaturas de imagens.

O treinamento foi realizado com o otimizador Adam usando *learning rates* 0.1, 0.3, 0.5, 0.8, 1 e uma *loss* normalizada. O número de níveis varia de 8 à 16. A SSIM é usada para avaliar o desempenho em tempo de teste.

Redes Neurais recorrentes

Além das operações usadas pelas redes não recorrentes, a abordagem recorrente usa camadas recorrentes. Uma rede neural recorrente (RNN) é um tipo de rede neural com memória para salvar informação sobre entradas passadas. Para isto, estas unidades de memória possuem conexões para elas mesmas. Esta informação temporal muda o comportamento da camada para a entrada atual [41].

Para manter informação temporal, muitas implementações das camadas foram propostas. Uma das primeiras abordagens efetivas se espelham na *Long Short-Term Memory (LSTM)*. Melhorias posteriores incluem a *Gated Recurrent Unit (GRU)*, que simplifica a componente recorrente e alcança *performance* similar em alguns cenários [41].

A camada LSTM possui a seguinte formulação, onde x_t , c_t e h_t denotam a *input*, *cell* e *hidden states* na iteração t [26]:

$$\begin{aligned} [f, i, o, j]^\top &= [\sigma, \sigma, \sigma, \tanh]^\top((Wx_t + Uh_{t-1}) + b) \\ c_t &= f \odot c_{t-1} + i \odot j \\ h_t &= o \odot \tanh(c_t), \end{aligned} \tag{2.15}$$

onde \odot é a multiplicação elemento a elemento, b é o bias, e σ é a função sigmoide. A saída da camada é h_t . A GRU, mencionada anteriormente, é descrita pelas seguintes equações:

$$\begin{aligned} z_t &= \sigma(W_z x_t + U_z h_{t-1}) \\ r_t &= \sigma(W_r x_t + U_r h_{t-1}) \\ h_t &= (1 - z_t) \odot h_{t-1} + z_t \odot \tanh(Wx_t + U(r_t \odot h_{t-1})). \end{aligned} \tag{2.16}$$

Se a versão convolucional é desejada, as multiplicações são substituídas por convoluções. Considerando as camadas de memória descritas por essas equações, a abordagem recursiva pode ser sumarizada usando a seguinte formulação [26]:

$$\begin{aligned} b_t &= B(E_t(r_{t-1})), \hat{x}_t = D_t(b_t) \\ r_t &= x - \hat{x}_t, r_0 = x, \hat{x}_0 = 0, \end{aligned} \tag{2.17}$$

onde D_t e E_t são o *encoder* e o *decoder* com os estados na iteração t , e b_t é a representação binária progressiva. \hat{x}_t é reconstrução “one-shot” e r_t é o resíduo entre x e a construção \hat{x}_t . Neste caso, B irá produzir uma representação binarizada de tamanho fixo. Uma das arquiteturas propostas em [9] usa a abordagem recorrente, além da arquitetura não recorrente descrita anteriormente.

No trabalho posterior, Toderici [26] mostrou em seus resultados que o efeito de usar um conjunto de treinamento de alta entropia³ é benéfico para a rede, dada a importância de treinar modelos de *machine learning* com exemplos difíceis. No seu trabalho, ele definiu o conjunto de alta entropia como sendo o conjunto formado por imagens que o PNG tem mais dificuldade de comprimir, ou seja, os arquivos no formato PNG com o maior número de bytes.

³Em processamento de imagens entropia diz respeito à quantidade de informação na imagem. Alta entropia, pode ser visto como maior variância nos valores dos pixels

Capítulo 3

Metodologia

Este capítulo apresenta a metodologia seguida para a verificação das hipóteses, incluindo o modelo construído e os *datasets* utilizados. Os modelos foram construídos usando o *framework* PyTorch [42] e avaliados usando métricas visuais: PSNR, SSIM e MS-SSIM. Também foi avaliada a quantidade de bits por pixel da imagem decodificada.

3.1 Bases de Dados

Foram utilizadas, principalmente, cinco bases de dados para o treinamento do modelo proposto, construídas usando as imagens das seguintes bases de dados:

1. *CLIC* [10]. Deste *dataset* foram pegos 4 conjuntos com os seguintes nomes e tamanhos:
 - (a) Professional valid: 41 imagens;
 - (b) Professional train: 585 imagens;
 - (c) Mobile valid: 61 imagens;
 - (d) Mobile train: 1048 imagens;
2. *DIV2K* [43]. Deste *dataset* foram pegos 2 conjuntos com os seguintes nomes e tamanhos:
 - (a) Train: 800 imagens;
 - (b) Valid: 100 imagens;
3. *EYE* [44]. Deste *dataset* foram pegos 2 conjuntos com os seguintes nomes e tamanhos:
 - (a) HD: 38 imagens;

(b) UHD: 40 imagens;

Também foi utilizada a base [2] e o conjunto *Mobile test* da base CLIC para teste. Todas as imagens usadas são de alta qualidade (sem ruído, boa iluminação, alta nitidez), sendo que as imagens *professional* possuem qualidade maior que as mobile. A base EYE consiste de imagens naturais de alta qualidade e resolução adquiridas usando várias câmeras. As imagens cobrem uma quantidade variada de cenas, incluindo cenas ao ar livre e interiores, imagens da natureza, pessoas, animais e cenas históricas retratadas em pinturas. As imagens das bases DIV2K e EYE têm resolução maior que as do CLIC, entretanto isso não faz muita diferença visto que são usados patches com 32 pixels de largura e altura na rede.

Primeiramente, todas as imagens foram separadas em *patches* com 32 pixels de largura e altura, resultando em 6,231,440 *patches*. Cada imagem foi codificada sem perdas no formato PNG, e o tamanho de cada arquivo é usado como critério para a entropia do *patch* (*patches* com tamanhos menores são considerados como sendo de “baixa entropia”). O histograma da base de dados completa é mostrado na Figura 3.1. Foram geradas cinco base de dados com cerca de 1.25 milhões de *patches* em cada uma. Para cada base é pego um subconjunto do total de *patches*.

Cada base de dados tem características específicas e entendê-las é um fator essencial para avaliar o modelo proposto e o impacto de métodos e hiperparâmetros diferentes nos resultados. As bases, nomeadas BD_i , $i \in \{0, \dots, 4\}$, possuem as seguintes características:

- **BD0:** formada por 1248978 de *patches* que pertencem ao grupo dos 20% com menor entropia;
- **BD1:** formada por 1251421 de *patches* que pertencem ao grupo dos que estão na faixa 40% à 60% (porcentagem dada de acordo com o *patch* com maior entropia);
- **BD2:** formada por 1248725 de *patches* que pertencem ao grupo dos 20% com maior entropia;
- **BD3:** formada por 1247033 de *patches* pegos de forma aleatória. Correspondem à 20% do total.
- **BD4:** formada por 1246698 de *patches*. 20% do total retirados aleatoriamente dos 50% dos *patches* com maior entropia.

Por construção, não há sobreposição entre as bases de dado 0, 1 e 2, mas existe sobreposição destas bases com as bases 3 e 4. Um histograma de cada base [Figuras 3.2 a 3.6] é dado.

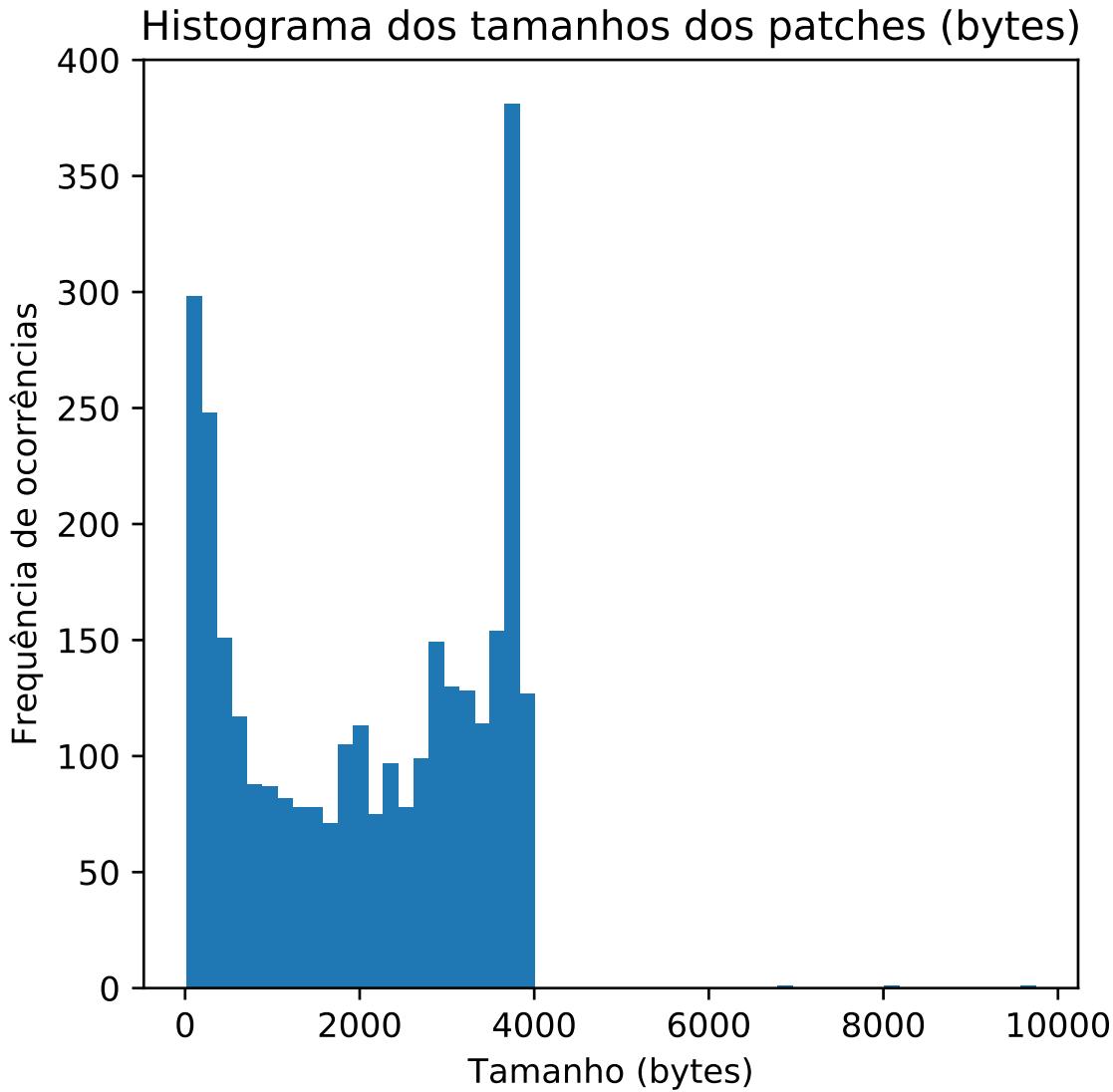


Figura 3.1: Histograma da base de dados completa formada por 6,231,440 de patches.

3.2 Modelos desenvolvidos

Foram testados *autoencoders* convolucionais com o objetivo de avaliar o potencial de cada um nas bases de dados propostas. O primeiro não possui binarização/quantização no latente gerado pelo *encoder*. O segundo incorpora o binarizador (representado com um trapézio nas figuras) baseado no do *Toderici* usando $\hat{b}(x)$. O últimos são os modelos convolucionais recursivos propostos por *Toderici* em [9] com a utilização de $\hat{b}(x)$ e de um codificador de entropia, de modo que se alguma redundância ainda for encontrada no latente a lógica é que o codificador irá explorá-la de alguma forma, reduzindo o tamanho

Histograma dos tamanhos dos patches (bytes)

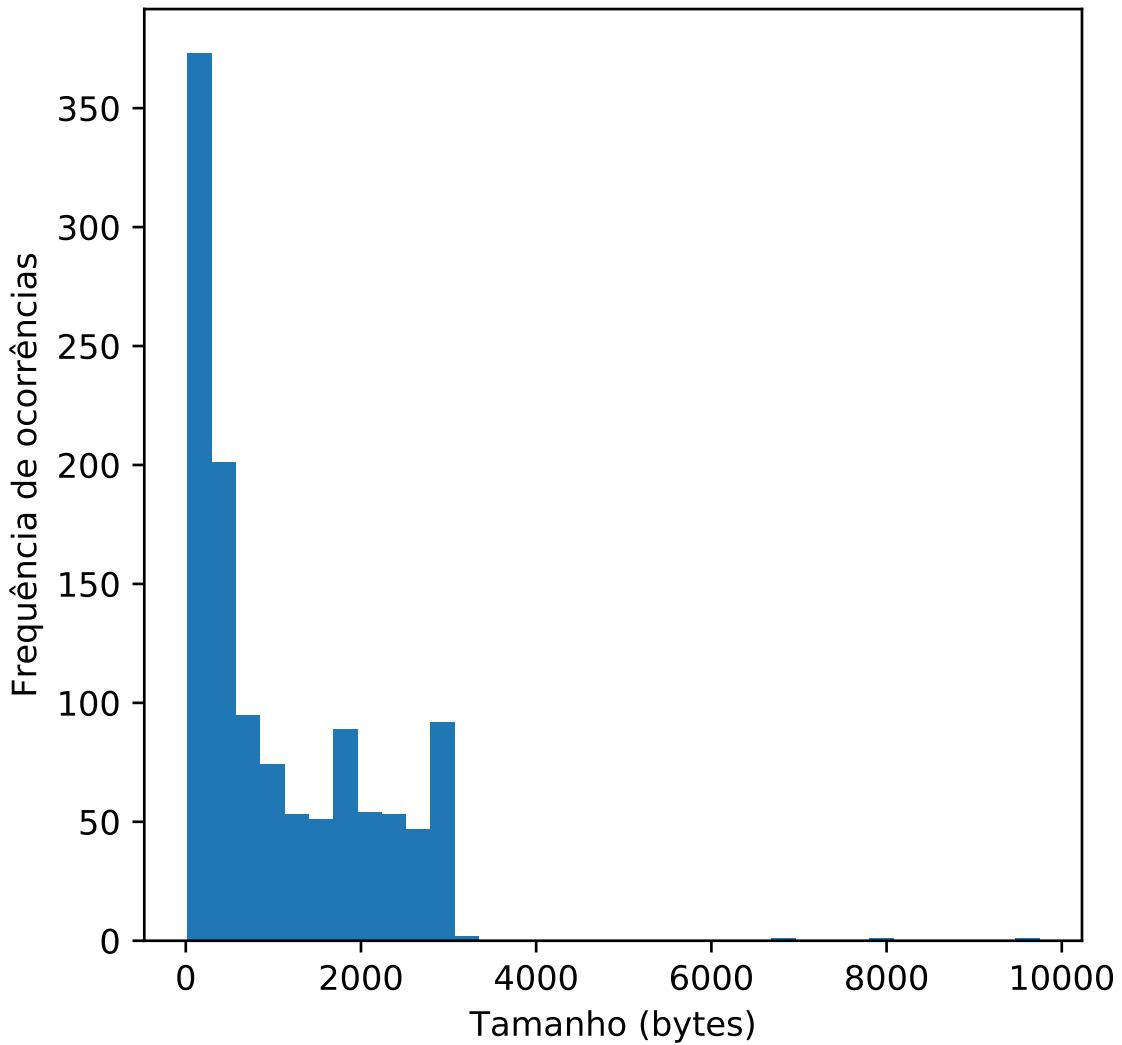


Figura 3.2: Histograma da BD0.

do *bitstream* comprimido. O Modelo 3 corresponde ao modelo não recursivo convolucional e o 4 é o modelo recursivo convolucional usando LSTM.

Pode-se pensar que estes últimos modelos realizam o trabalho de transformação e quantização, empacotando e descartando informação, e que o codificador de entropia irá finalizar o trabalho, explorando alguma redundância ainda não explorada.

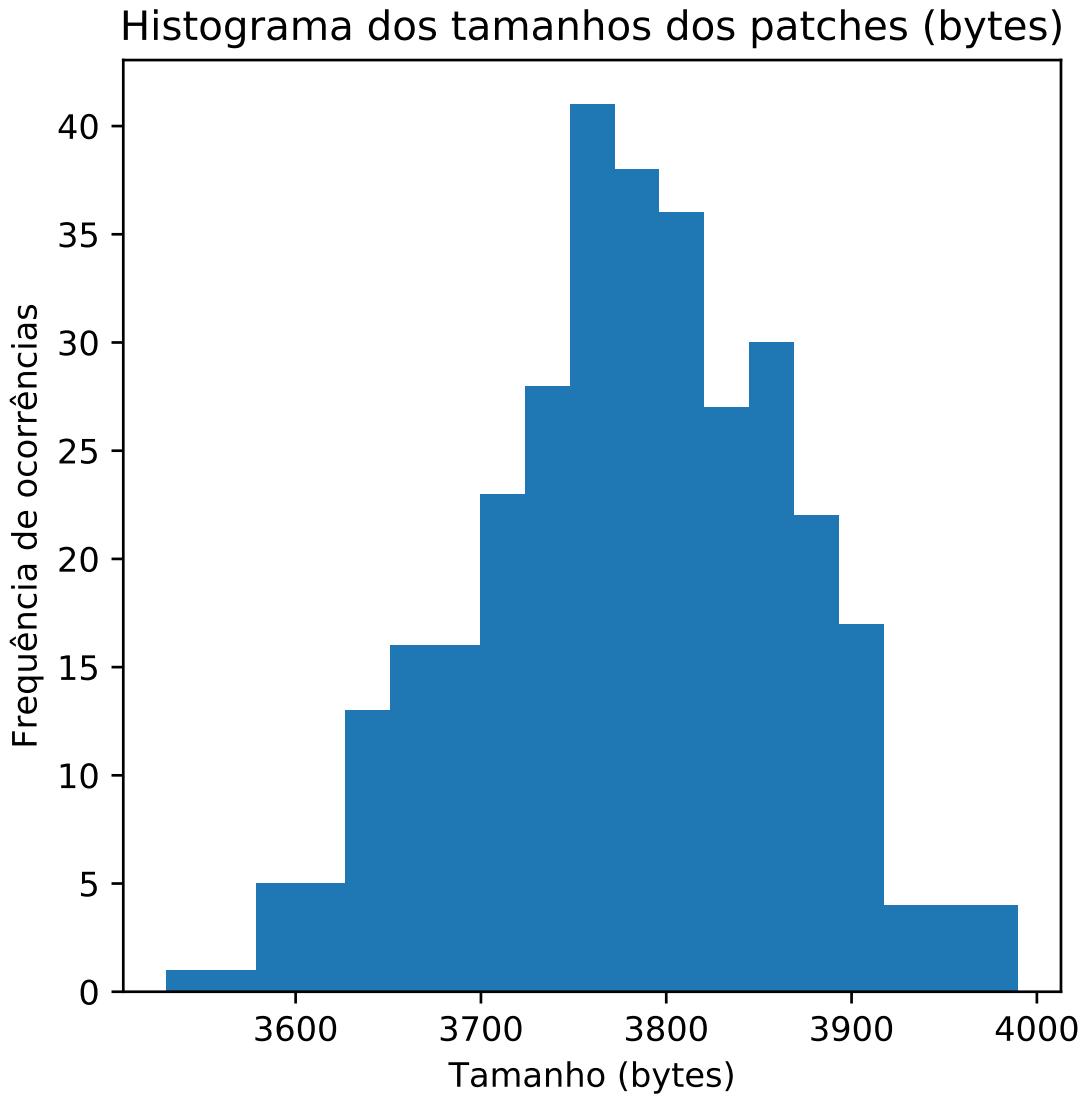


Figura 3.3: Histograma da BD1.

3.2.1 Formando o Bitstream

É importante descrever como o *bitstream* é formado quando se aplica o binarizador B . Em todos os modelos, a imagem é dividida em *patches* de tamanho 32x32 *pixels*. Assim, se uma imagem possuir 9 *patches*, ela terá um *bitstream* para cada latente de cada *patch* e o *bitstream* da imagem será dado pela concatenação dos *bitstreams* de cada *patch*.

Para cada um desses *patches*, o modelo é treinado de maneira residual conforme explicado no capítulo anterior. Para os modelos que usam mais de 1 nível de resíduo, a taxa nominal é acrescida de 0.125 bits por pixel para cada nível.

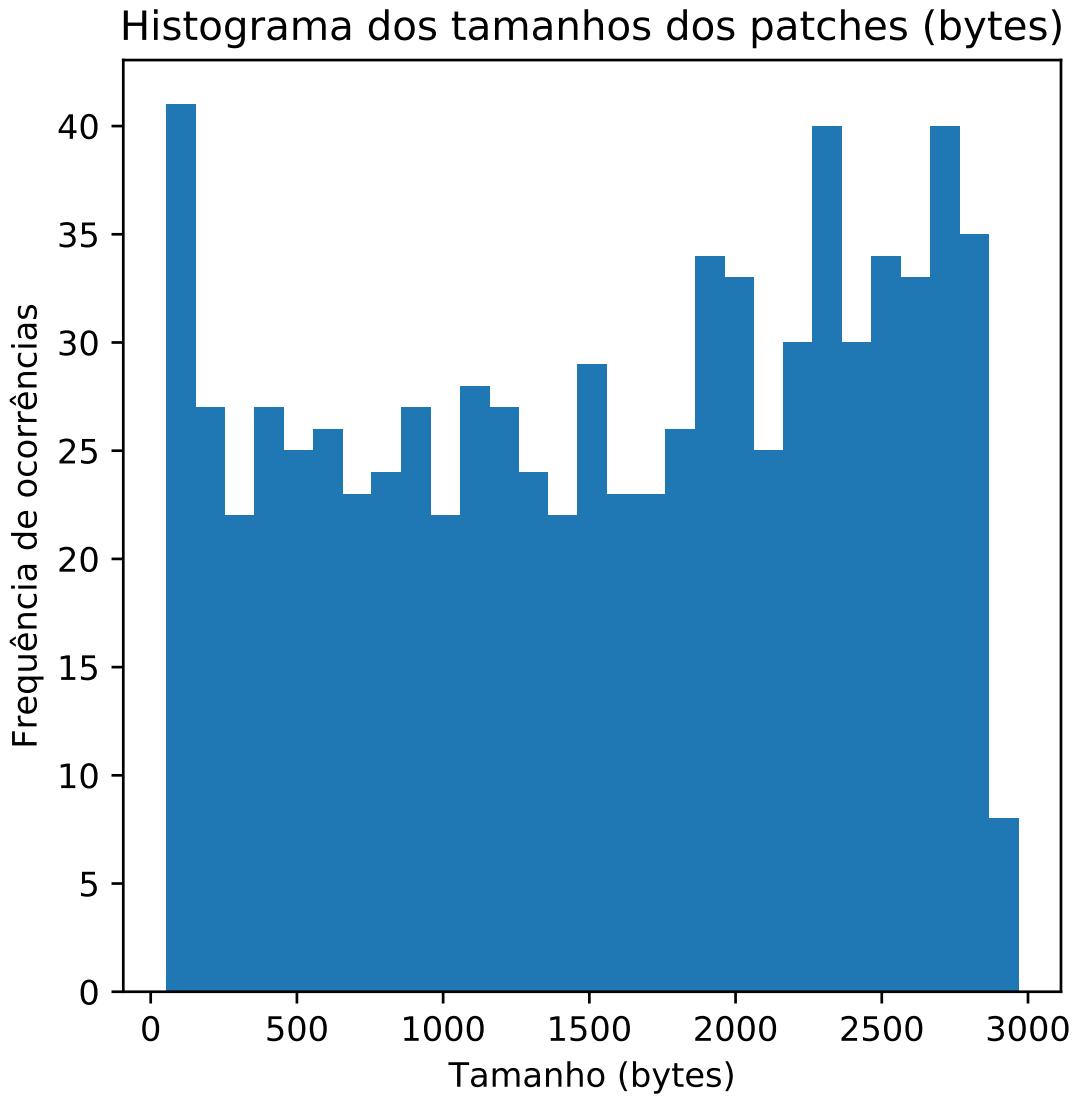


Figura 3.4: Histograma da **BD2**.

Os latentes são concatenados de modo que para o nível n tem-se n latentes binarizados concatenados. Assim, para o nível 1 tem-se o *bitstream* b_1 correspondente ao latente binarizado do nível 1, para o nível 2 tem-se o *bitstream* $b_1 + b_2$, e assim em diante. De maneira geral, para o nível n tem-se $b_1 + b_2 + \dots + b_n$, onde $+$ denota a concatenação de *bitstreams*.

Para os últimos três modelos foi usado o *gzip* (codificador de entropia) no latente com o objetivo de reduzir a taxa de bits por pixel.

Nas seguintes subseções são apresentadas as ilustrações de algumas arquiteturas utilizadas. Os retângulos indicam as convoluções e os retângulos arredondados indicam as

Histograma dos tamanhos dos patches (bytes)

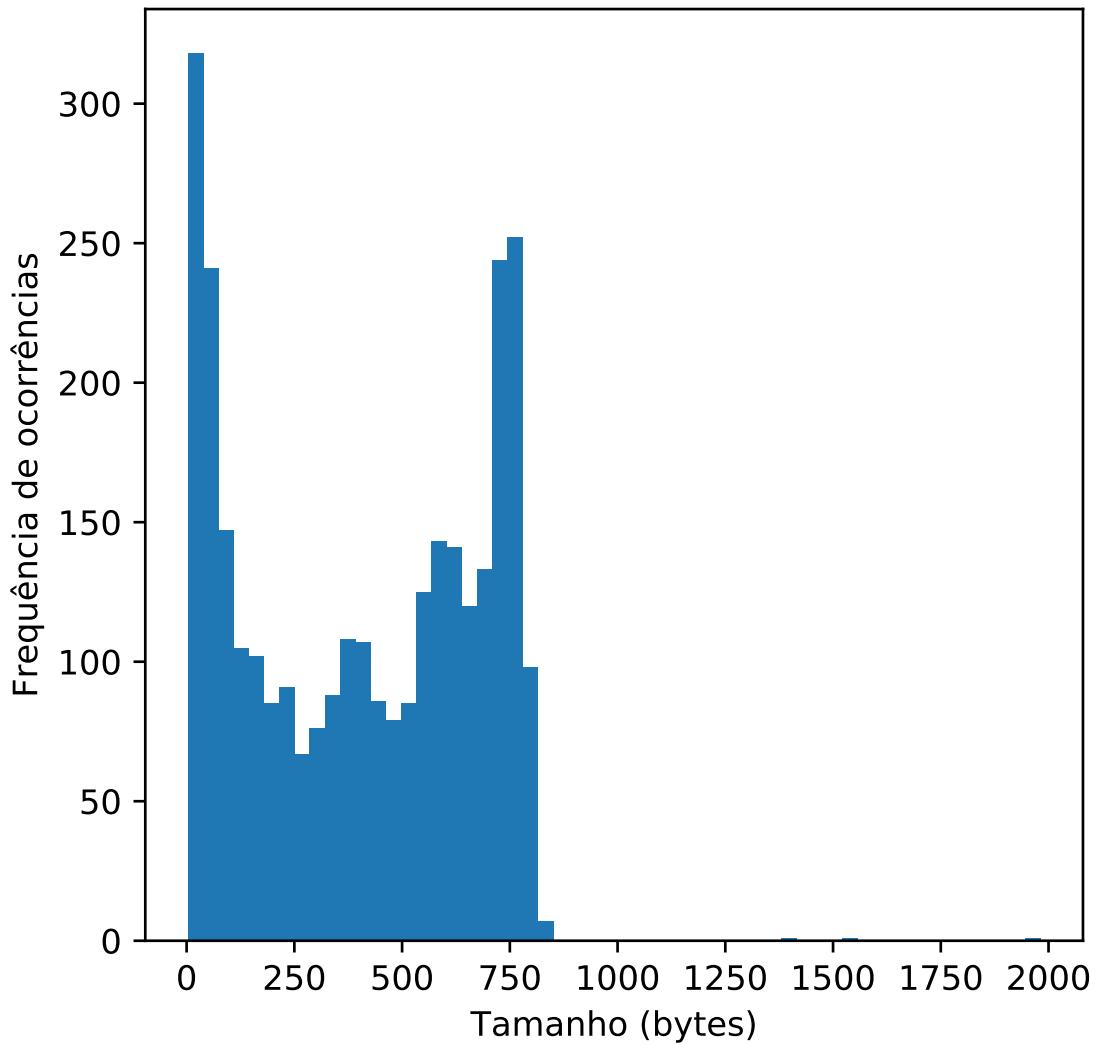


Figura 3.5: Histograma da **BD3**.

convoluções transpostas. O tamanho do *kernel w* é indicado na primeira linha do retângulo. A segunda linha informa o número de filtros (canais de saídas). A última linha indica o tamanho (igual nas duas direções) do *stride* utilizado e a função de ativação. A última camada do *decoder* (convolução transposta) de cada um dos modelos é usada para recuperar a informação de cor.

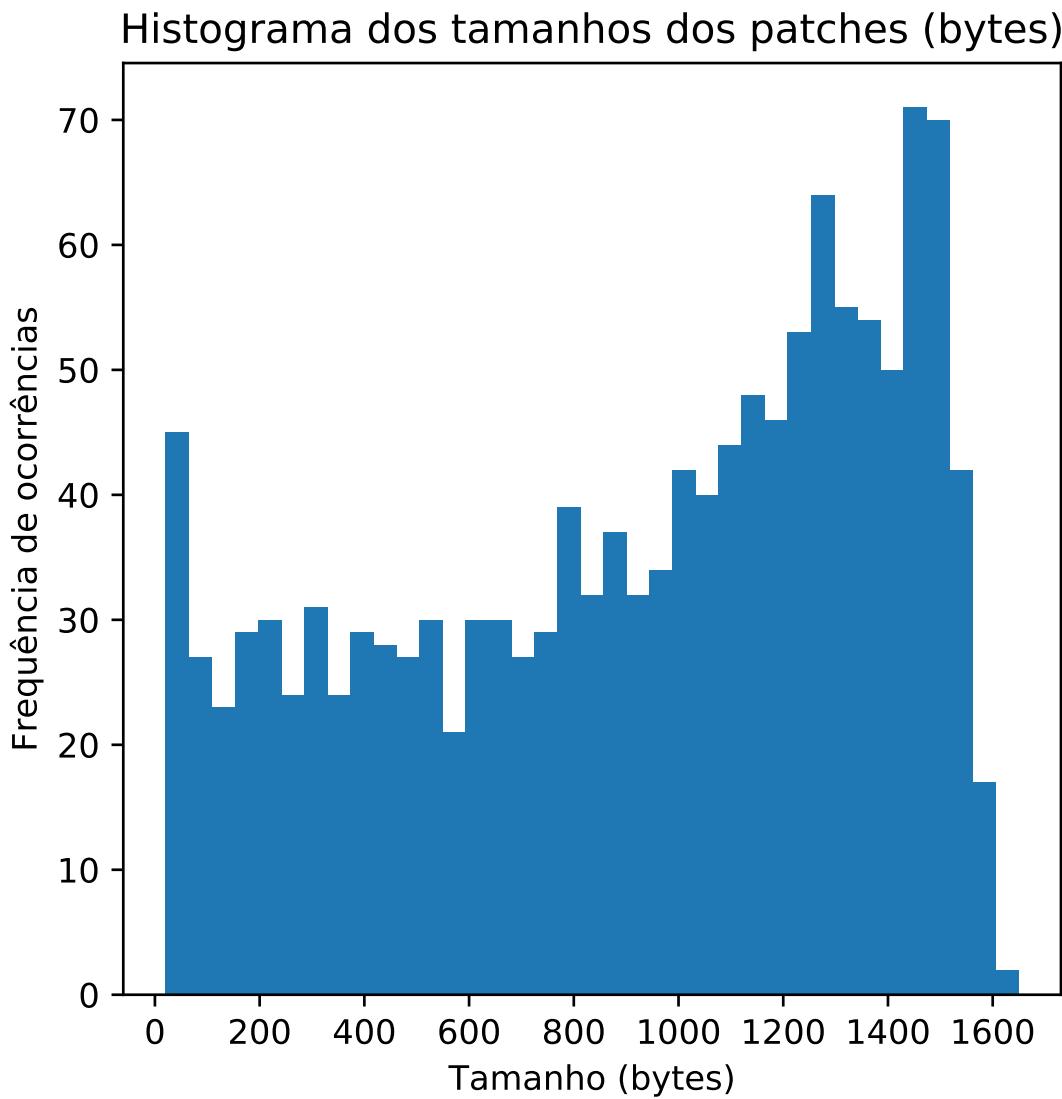


Figura 3.6: Histograma da BD4.

3.2.2 Modelo 1

Este primeiro modelo foi desenvolvido com o objetivo de avaliar o potencial de um *autoencoder* convolucional simples sem camada de gargalo com binarizador/quantizador. Aqui, o armazenamento da imagem codificada pelo *encoder* seria custoso, visto que não há binarização. Portanto, o objetivo é apenas avaliar a distorção das imagens reconstruídas. A arquitetura é ilustrada na Figura 3.7.

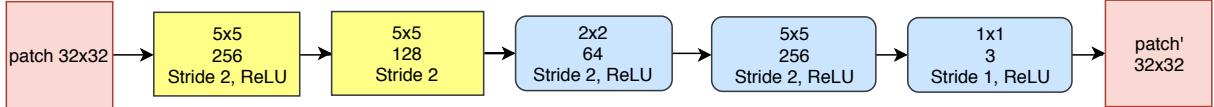


Figura 3.7: Ilustração do *autoencoder* mais básico desenvolvido.

3.2.3 Modelo 2

O segundo modelo adiciona o binarizador na camada de gargalo, o que força o *encoder* à comprimir informação visto que será gerada uma saída inteira discreta no conjunto $\{-1, 1\}$ a partir da entrada real contínua. Há uma grande perda de informação em troca de ganho em espaço, pois cada valor pode ser salvo usando apenas um bit agora. Nesta arquitetura [Figura 3.8] a taxa *nominal* de bits por pixel é $\frac{8 \cdot 8 \cdot 128}{32 \cdot 32} = 8$.

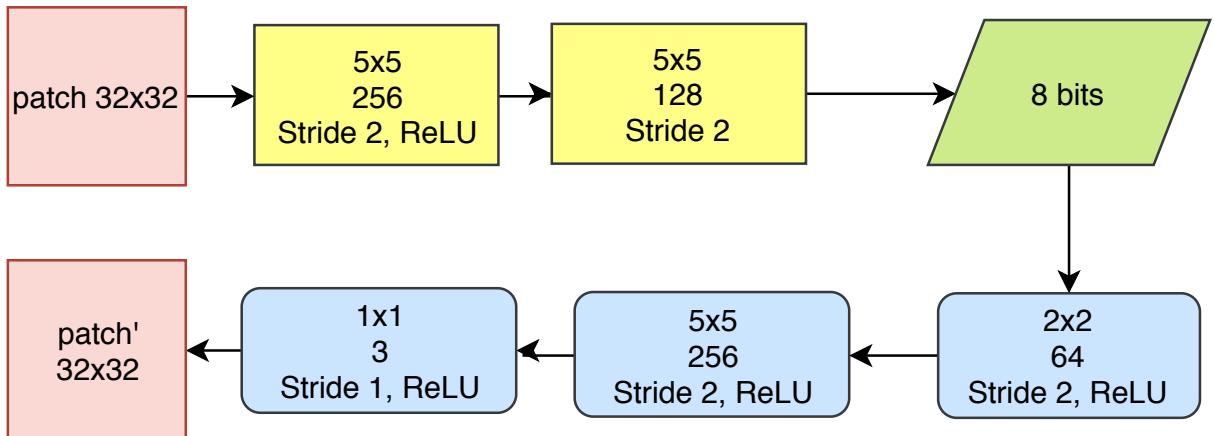


Figura 3.8: Ilustração do segundo modelo desenvolvido.

Capítulo 4

Experimentos e Resultados

Foram realizados vários experimentos com os modelos apresentados no capítulo anterior e com o **JPEG** e **JPEG2000** com o objetivo de avaliar o potencial de cada um deles nas bases de dados propostas. Os *patches* usados em todos os experimentos realizados neste capítulo possuem 32 pixels de largura e altura. Para os testes em que o conjunto de treino e de teste são a mesma base, foram gerados dois subconjuntos disjuntos: um para treino e um para teste. Este último possui cerca de 10% do tamanho total e, evidentemente, não faz parte do conjunto de treino. Os resultados do **JPEG** nas bases de teste utilizadas são apresentados na seção 4.1 e os resultados dos modelos apresentados no capítulo anterior são apresentados nas seções 4.2, 4.3, 4.4, 4.5. Os três primeiros modelos apresentados possuem um único nível de resíduo. Todas as taxas estão em bits por pixel.

4.1 JPEG

Nesta seção é apresentado o desempenho obtido pelo **JPEG** nas bases de teste.

Tabela 4.1: Tabela contendo médias obtidas pelo **JPEG** em cada uma das bases de teste utilizadas.

Bases	Taxa	PSNR	SSIM	MSSIM	Quality
CLIC Mobile test (patches 32)	8	44.23	0.98	0.99	94.06
CLIC Mobile test	2	39.58	0.96	0.99	88.69
Kodak	2	36.77	0.95	0.99	85.91
BD0	8	57.85	0.99	0.99	99.18
BD1	8	42.94	0.97	0.99	95.94
BD2	8	32.31	0.95	0.99	83.69
BD3	8	43.84	0.96	0.99	93.94
BD4	8	36.72	0.96	0.99	89.68

Tabela 4.2: Tabela contendo o valor da PSNR, em decíbeis, dos testes do Modelo 1 com o uso do otimizador *Adam* e *learning rate* fixa. As linhas denotam a base de treino utilizada. As colunas denotam as bases de teste usadas para avaliação do modelo. O índice “todas” se refere ao uso de todas as imagens de todas as bases **BD** para treino.

Treino (linhas) x Teste (colunas)	BD0	BD1	BD2	BD3	BD4
BD0	49.66	38.07	26.34	38.05	31.13
BD1	47.57	44.08	34.54	42.66	38.69
BD2	53.69	51.16	44.07	50.06	47.14
BD3	50.62	47.88	39.76	46.59	43.25
BD4	47.30	46.98	41.10	45.63	43.75
Todas	46.77	46.35	43.94	45.88	45.08

Tabela 4.3: Tabela contendo o valor da PSNR, em decíbeis, dos testes do Modelo 1 com o uso do otimizador *Adam* e *learning rate* fixa. As linhas denotam a base de treino utilizada. As colunas denotam as bases de teste usadas para avaliação do modelo. O índice “todas” se refere ao uso de todas as imagens de todas as bases **BD** para treino. O uso de $x + y$ denota o uso de todas as imagens do conjunto x e do conjunto y para treinamento.

Treino (linhas) x Teste (coluna)	CLIC Mobile test
BD0	34.77
BD1	44.11
BD2	48.97
BD3	45.34
BD4	42.42
Todas	51.27
Todas + CLIC Mobile train	55.78
Todas + Clic Mobile train + Clic Professional train	47.26
CLIC Mobile Train	46.63

4.2 Modelo 1

Primeiramente, o Modelo 1 [Figura 3.7] foi testado em todas as bases de dados usando *learning rate* fixa durante todo o treinamento e o otimizador *Adam*. Foi utilizada apenas uma época para treino em todos as avaliações realizadas. Os resultados são apresentados na Tabela 4.2. Com estes mesmos hiperparâmetros para treino, foram realizados alguns testes na base de teste (nomeada como *test*) do CLIC [10]. Alguns destes testes também usaram, em adição aos conjuntos de treino montados, os conjuntos de treino (*train*) e validação (*valid*) do CLIC. Os resultados são apresentados na Tabela 4.3.

É interessante notar que o **BD2** é a melhor base de dados para treino, o que reforça os resultados encontrados por *Toderici* em [26]. Outro resultado interessante obtido ocorre ao treinar na base de dados com menor entropia e testar na base com maior entropia.

Tabela 4.4: Tabela contendo os resultados do Modelo 2 para as métricas visuais PSNR, SSIM e MS-SSIM a uma taxa nominal de 8 bits por pixel.

Bases de Treino e Teste	BPP	PSNR	SSIM	MS-SSIM	Épocas
CLIC Mobile test	8	35.03	0.94	0.98	30
BD1	8	35.26	0.93	0.98	30
BD2	8	29.10	0.95	0.98	30

Pode-se notar que foi muito maléfico para a aprendizagem da rede treinar somente com exemplos fáceis para testar em exemplos difíceis. Posteriormente foram realizados alguns testes usando o método de atualização de *learning rate* explicado no capítulo anterior. A política utilizada foi a `exp_range`, pois foi a que obteve melhores resultados. Após a realização de vários testes, foi possível aumentar os dB obtidos anteriormente de 44.08 e 44.07 ao treinar e testar no BD0 e BD1 para 50.81 e 47.83, respectivamente. Considerando que todos esses treinamentos referentes aos resultados apresentados nas Tabelas 4.2 a 4.3 foram executados utilizando uma única época, esta melhora pode ser considerada como a “superconvergência” apontada em [7]. É interessante notar que, ao contrário do que foi observado por *Leslie* neste artigo, o uso do otimizador *Adam* com a `exp_range` apresentou melhorias significativas para o Modelo 1.

4.3 Modelo 2

Para o Modelo 2 [Figura 3.8], foram feitos testes nas bases **CLIC Mobile**, **BD1** e **BD2**. Os resultados são apresentados na Tabela 4.4. Conforme esperado, dentre as bases **BD1** e **BD2**, os piores resultados do *JPEG* e do *autoencoder* foram encontrados no **BD2** que é o com maior entropia (PNG teve mais dificuldade para comprimir). Nota-se que a menor diferença proporcional entre o resultado do modelo e do *JPEG* na métrica PSNR se dá no **BD2**, o que reforça a observação feita em [9] de que em baixas taxas e resoluções espaciais, os artefatos blocantes do *JPEG* (ruído causado pela perda de informação) se tornam mais comuns. Na Figura 4.1 é mostrado um *patch* reconstruído que obteve 36.69 dB de *PSNR*.

4.4 Modelo 3

Para o Modelo 3, foram feitos testes nas bases **CLIC Mobile**, **BD0**, **BD1**, **BD2**, **BD3**, **BD4** e **Kodak**. Os resultados são apresentados na Tabela 4.5. Uma comparação com o *JPEG* para diferentes taxas e métricas de distorção é apresentada na Figuras 4.3 a 4.2. Este modelo possui um único nível de resíduo.



Figura 4.1: Imagem original (esquerda) e *patch* reconstruído pelo Modelo 2 (direita).

Tabela 4.5: Tabela contendo os resultados do Modelo 3.

Bases	Taxa	PSNR	SSIM	MS-SSIM
CLIC Mobile test (patches 32)	2	33.75	0.92	0.97
Kodak (patches 32)	2	31.46	0.88	0.96
BD0	2	40.24	0.97	0.99
BD1	2	35.00	0.91	0.98
BD2	2	27.53	0.91	0.97
BD3	2	33.27	0.91	0.97
BD4	2	30.16	0.90	0.97

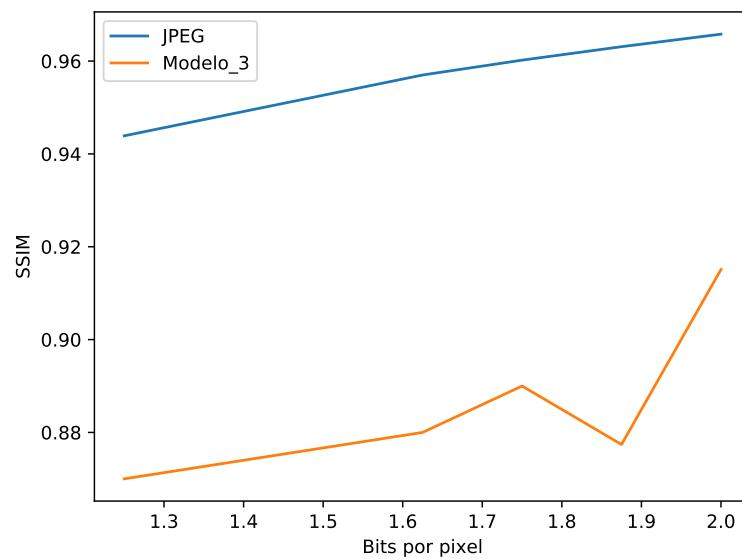


Figura 4.2: Comparação do Modelo 3 com o JPEG na métrica PSNR em diferentes taxas.

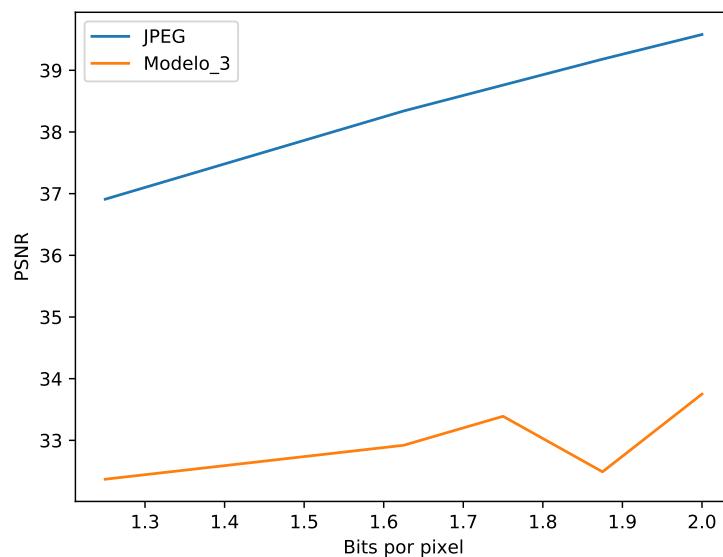


Figura 4.3: Comparação do Modelo 3 com o JPEG na métrica SSIM em diferentes taxas.

Tabela 4.6: Tabela contendo os valores da taxa (BPP) e PSNR (decíbeis) do Modelo 4, JPEG e JPEG2000 para a base [2].

Nível	Network		Network + GZIP		JPEG		JPEG 2000	
	Taxa	PSNR	Taxa	PSNR	Taxa	PSNR	Taxa	PSNR
1	0.125	24.44	0.10	24.44	0.17	21.52	0.10	26.23
2	0.250	26.81	0.23	26.81	0.23	24.51	0.23	28.29
3	0.375	28.32	0.35	28.32	0.36	27.53	0.35	29.54
4	0.500	29.46	0.48	29.46	0.48	29.08	0.47	30.66
5	0.625	30.44	0.60	30.44	0.60	30.23	0.60	31.68
6	0.750	31.26	0.73	31.26	0.73	31.14	0.72	32.55
7	0.875	31.98	0.85	31.98	0.85	31.92	0.85	33.43
8	1.000	32.63	0.98	32.63	0.98	32.60	0.97	34.10
9	1.125	33.19	1.10	33.19	1.10	33.21	1.10	34.81
10	1.250	33.65	1.23	33.65	1.22	33.72	1.21	35.37

Tabela 4.7: Tabela contendo os valores da taxa (BPP) e SSIM do Modelo 4, JPEG e JPEG2000 para a base [2].

Nível	Network		Network + GZIP		JPEG		JPEG 2000	
	Taxa	SSIM	Taxa	SSIM	Taxa	SSIM	Taxa	SSIM
1	0.125	0.65586	0.10	0.65586	0.17	0.56163	0.10	0.69056
2	0.250	0.74829	0.23	0.74829	0.23	0.66125	0.23	0.77412
3	0.375	0.80243	0.35	0.80243	0.36	0.76348	0.35	0.81706
4	0.500	0.83860	0.48	0.83860	0.48	0.81317	0.47	0.84603
5	0.625	0.86547	0.60	0.86547	0.60	0.84571	0.60	0.86805
6	0.750	0.88592	0.73	0.88592	0.73	0.86819	0.72	0.88374
7	0.875	0.90060	0.85	0.90060	0.85	0.88507	0.85	0.89826
8	1.000	0.91252	0.98	0.91252	0.98	0.89808	0.97	0.90900
9	1.125	0.92177	1.10	0.92177	1.10	0.90873	1.10	0.91981
10	1.250	0.92873	1.23	0.92873	1.22	0.91712	1.21	0.92769

4.5 Modelo 4

Para este modelo foram usadas as bases **BD2**, **BD3** e **BD4** como treino. O modelo foi testado em duas bases: Kodak [2] e [10]. O modelo possui 10 níveis de resíduos e foi treinado por 500 mil iterações. Os resultados referentes à base Kodak, são mostrados nas Figuras 4.4 a 4.6 e summarizados nas Tabelas 4.6 a 4.7. Os resultados referentes às 47 imagens com muito conteúdo de alta frequência são mostrados nas Figuras 4.15 a 4.17 e summarizados nas Tabelas 4.9 a 4.11.

Apesar da função de perda utilizada ser a MSE, o modelo conseguiu obter melhores resultados comparados aos codecs na SSIM e MS-SSIM.

Nota-se que para imagens (Figuras 4.7 a 4.8) com muito conteúdo de alta frequência (muitos detalhes), o modelo se sai melhor para todas as métricas visuais avaliadas. O

Tabela 4.8: Tabela contendo os valores da taxa (BPP) e MSSSIM do Modelo 4, JPEG e JPEG2000 para a base [2].

Nível	Network		Network + GZIP		JPEG		JPEG 2000	
	Taxa	MS	Taxa	MS	Taxa	MS	Taxa	MS
1	0.125	0.84270	0.10	0.84270	0.17	0.71702	0.10	0.88245
2	0.250	0.92192	0.23	0.92192	0.23	0.82584	0.23	0.93421
3	0.375	0.94893	0.35	0.94893	0.36	0.91058	0.35	0.95471
4	0.500	0.96244	0.48	0.96244	0.48	0.94160	0.47	0.96546
5	0.625	0.97054	0.60	0.97054	0.60	0.95843	0.60	0.97269
6	0.750	0.97581	0.73	0.97581	0.73	0.96788	0.72	0.97750
7	0.875	0.97993	0.85	0.97993	0.85	0.97420	0.85	0.98118
8	1.000	0.98274	0.98	0.98274	0.98	0.97840	0.97	0.98387
9	1.125	0.98502	1.10	0.98502	1.10	0.98158	1.10	0.98630
10	1.250	0.98665	1.23	0.98665	1.22	0.98381	1.21	0.98807

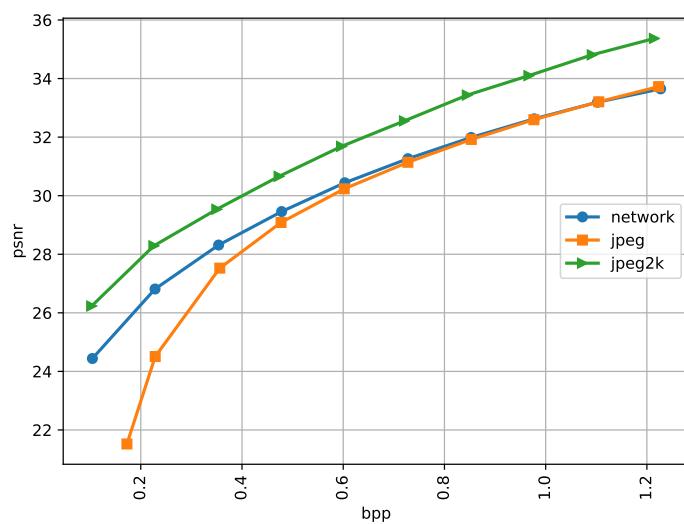


Figura 4.4: Comparaçāo do Modelo 4 com o JPEG e JPEG2000 na métrica PSNR em diferentes taxas para a base Kodak [2].

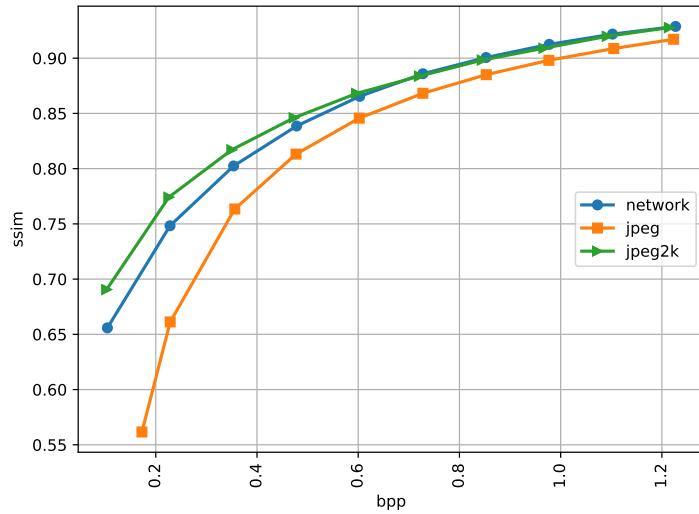


Figura 4.5: Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica SSIM em diferentes taxas para a base Kodak [2].

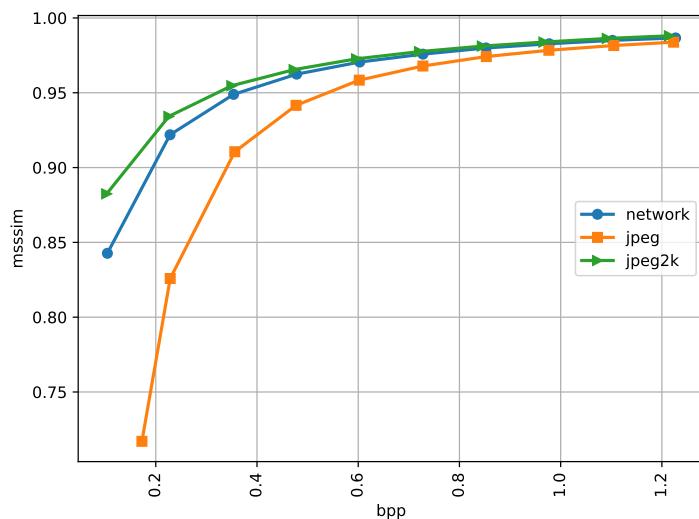


Figura 4.6: Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica MS-SSIM em diferentes taxas para a base Kodak [2].

que era esperado, visto que o JPEG e o JPEG2000 assumem que sinais de alta frequência não importam muito (assumem que maior parte energia da imagem estará contida em coeficientes de baixa frequência) e faz com que os coeficientes de baixa frequência tenham maior precisão ao quantizar.



Figura 4.7: Imagem Kodim05 [2].

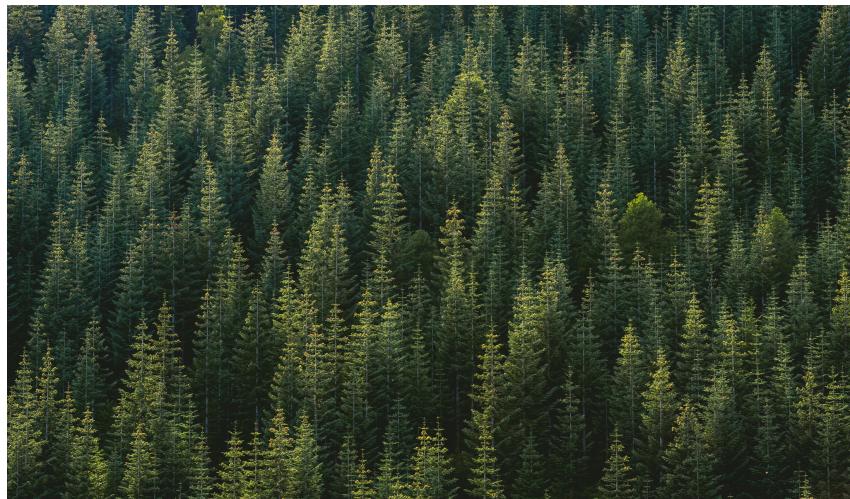


Figura 4.8: Imagem jason-leem [10].

Considerando os resultados anteriores obtidos para as imagens kodim05 e kodim09, o Modelo 4 também foi testado em um conjunto de 47 imagens com muito conteúdo de alta frequência retiradas da [2] e da [10]. Os resultados são mostrados nas Figuras 4.15 a 4.17 e sumarizados nas Tabelas 4.9 a 4.11.

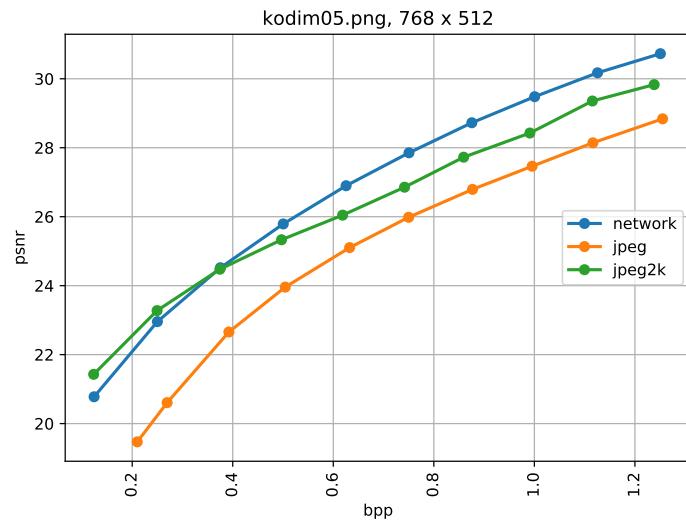


Figura 4.9: Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica PSNR em diferentes taxas para a imagem kodim05 [2].

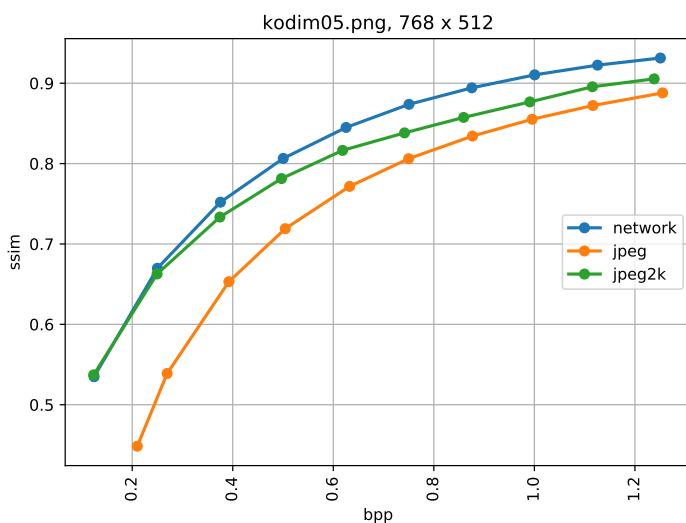


Figura 4.10: Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica SSIM em diferentes taxas para a imagem kodim05 [2].

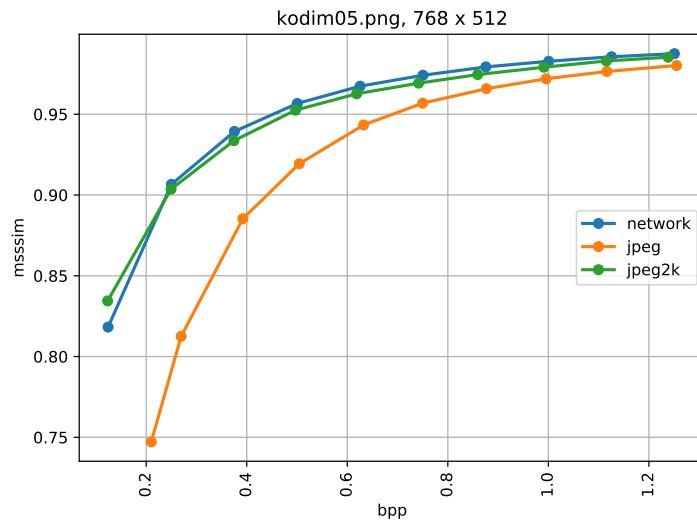


Figura 4.11: Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica MS-SSIM em diferentes taxas para a imagem kodim05 [2].

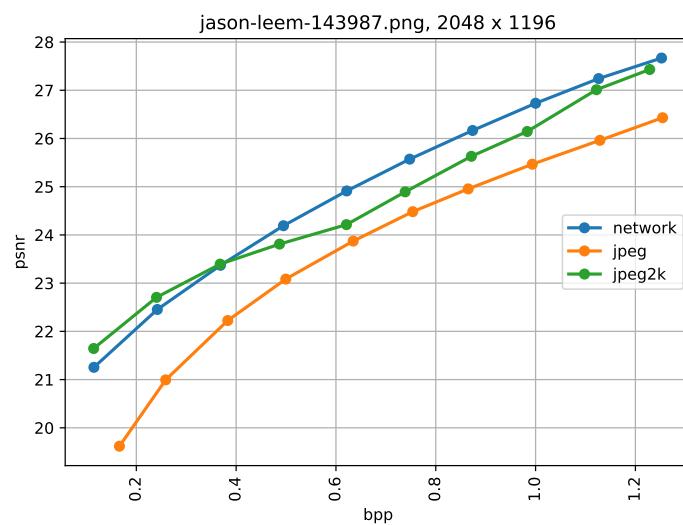


Figura 4.12: Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica PSNR em diferentes taxas para a imagem de [10].

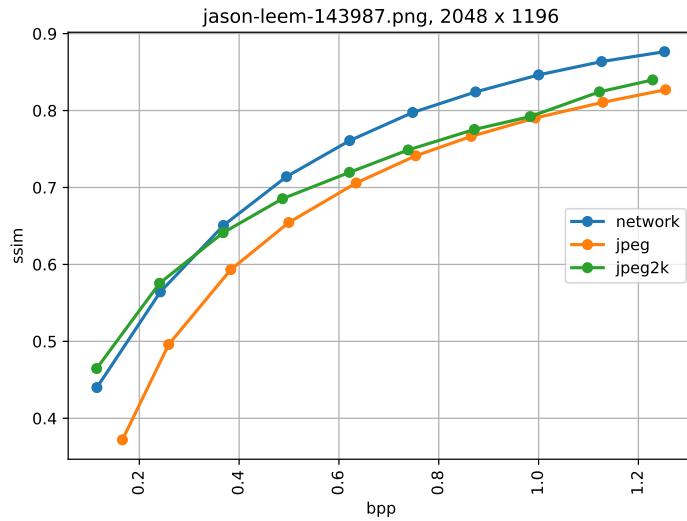


Figura 4.13: Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica SSIM em diferentes taxas para a imagem de [10].

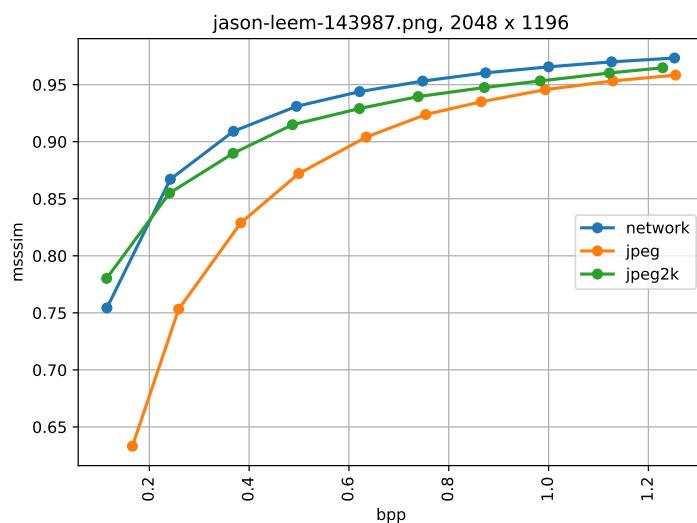


Figura 4.14: Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica MS-SSIM em diferentes taxas para a imagem de [10].

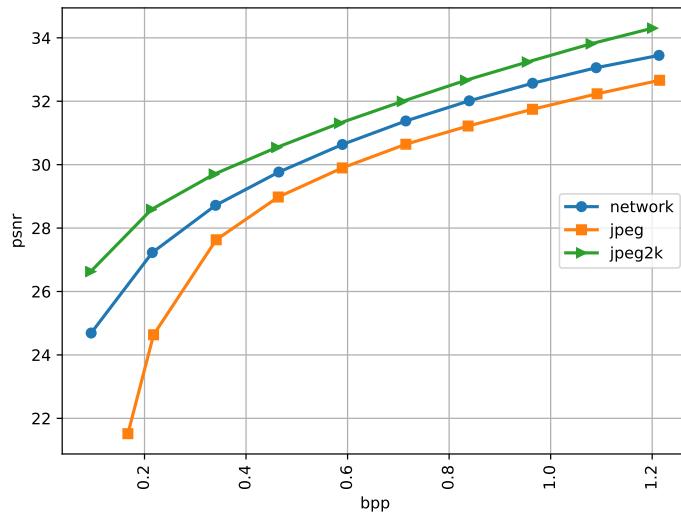


Figura 4.15: Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica PSNR em diferentes taxas para 47 imagens com muito conteúdo de alta frequência retiradas das bases [10] e [2].

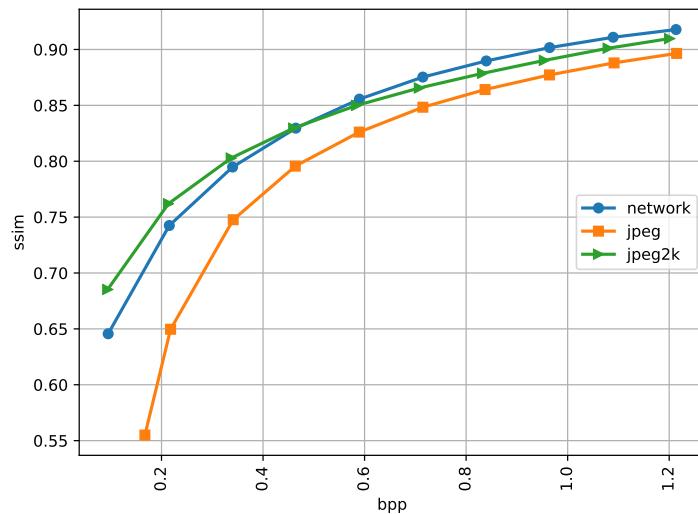


Figura 4.16: Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica SSIM em diferentes taxas para 47 imagens com muito conteúdo de alta frequência retiradas das bases [10] e [2].

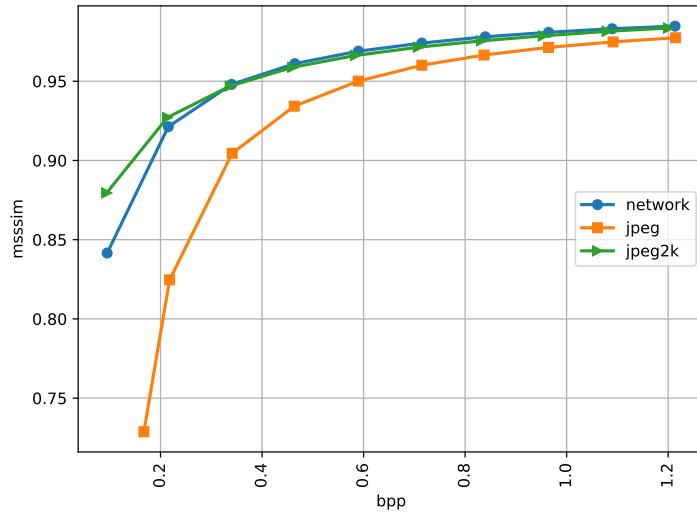


Figura 4.17: Comparação do Modelo 4 com o JPEG e JPEG2000 na métrica MS-SSIM em diferentes taxas para 47 imagens com muito conteúdo de alta frequência retiradas das bases [10] e [2].

Tabela 4.9: Tabela contendo os valores da taxa (BPP) e PSNR do Modelo 4, JPEG e JPEG2000 para 47 imagens da base [2] e [10] com muito conteúdo de alta frequência.

Nível	Network		Network + GZIP		JPEG		JPEG 2000	
	Taxa	PSNR	Taxa	PSNR	Taxa	PSNR	Taxa	PSNR
1	0.125	24.69	0.10	24.69	0.17	21.51	0.10	26.63
2	0.250	27.23	0.21	27.23	0.22	24.64	0.21	28.59
3	0.375	28.72	0.34	28.72	0.34	27.63	0.33	29.70
4	0.500	29.77	0.46	29.77	0.46	28.98	0.46	30.54
5	0.625	30.64	0.59	30.64	0.59	29.90	0.58	31.31
6	0.750	31.38	0.71	31.38	0.71	30.64	0.71	31.99
7	0.875	32.01	0.84	32.01	0.84	31.22	0.83	32.66
8	1.000	32.57	0.96	32.57	0.96	31.75	0.95	33.23
9	1.125	33.06	1.09	33.06	1.10	32.23	1.08	33.81
10	1.250	33.45	1.20	33.45	1.21	32.67	1.20	34.30

Tabela 4.10: Tabela contendo os valores da taxa (BPP) e SSIM do Modelo 4, JPEG e JPEG2000 para 47 imagens da base [2] e [10] com muito conteúdo de alta frequência.

Nível	Network		Network + GZIP		JPEG		JPEG 2000	
	Taxa	SSIM	Taxa	SSIM	Taxa	SSIM	Taxa	SSIM
1	0.125	0.64562	0.10	0.64562	0.17	0.55499	0.10	0.68526
2	0.250	0.74250	0.21	0.74250	0.22	0.64966	0.21	0.76202
3	0.375	0.79480	0.34	0.79480	0.34	0.74767	0.33	0.80283
4	0.500	0.82961	0.46	0.82961	0.46	0.79561	0.46	0.82961
5	0.625	0.85558	0.59	0.85558	0.59	0.82607	0.58	0.84989
6	0.750	0.87530	0.71	0.87530	0.71	0.84832	0.71	0.86553
7	0.875	0.88973	0.84	0.88973	0.84	0.86414	0.83	0.87862
8	1.000	0.90168	0.96	0.90168	0.96	0.87722	0.95	0.89020
9	1.125	0.91092	1.09	0.91092	1.10	0.88800	1.08	0.90109
10	1.250	0.91784	1.20	0.91784	1.21	0.89648	1.20	0.90973

Tabela 4.11: Tabela contendo os valores da taxa (BPP) e MSSSIM do Modelo 4, JPEG e JPEG2000 para 47 imagens da base [2] e [10] com muito conteúdo de alta frequência.

Nível	Network		Network + GZIP		JPEG		JPEG 2000	
	Taxa	MS	Taxa	MS	Taxa	MS	Taxa	MS
1	0.10	0.84156	0.10	0.84156	0.17	0.72882	0.10	0.87953
2	0.21	0.92131	0.21	0.92131	0.22	0.82462	0.21	0.92713
3	0.34	0.94802	0.34	0.94802	0.34	0.90454	0.33	0.94699
4	0.46	0.96112	0.46	0.96112	0.46	0.93426	0.46	0.95877
5	0.59	0.96897	0.59	0.96897	0.59	0.95006	0.58	0.96628
6	0.71	0.97408	0.71	0.97408	0.71	0.96010	0.71	0.97148
7	0.84	0.97806	0.84	0.97806	0.84	0.96655	0.83	0.97552
8	0.96	0.98088	0.96	0.98088	0.96	0.97134	0.95	0.97863
9	1.09	0.98315	1.09	0.98315	1.10	0.97483	1.08	0.98151
10	1.20	0.98477	1.20	0.98477	1.21	0.97746	1.20	0.98356

4.6 Ganhos obtidos pelo GZIP

Nas seguintes subseções será analisado o ganho na taxa obtido ao usar o codificador de entropia *gzip*.

4.6.1 Modelo 4

A Figura 4.18 mostra o ganho percentual médio por nível do Modelo 4 [4.5] para a base da Kodak [2]. Nota-se um decaimento exponencial no ganho. Considerando o funcionamento do modelo com resíduos, isto sugere uma especialização e robustez cada vez maior nos *encoders* dos níveis mais superiores ao explorar estruturas dos dados e comprimir informação nos latentes. Assim, os *bitstreams* gerados a partir dos latentes se tornam cada vez mais complexos (os resíduos possuirão muita variância nos valores) e com menos redundâncias fazendo com que o *gzip* não consiga explorar redundâncias bem ou que não haja muitas redundâncias.

Isto acontece na maior parte das imagens, principalmente nas que possui muito conteúdo de alta frequência e onde há muita informação que, consequentemente, geram latentes muito complexos, o que faz com que os *bitstreams* contenham muita informação já que os resíduos dos *patches* contém muita informação e geram muitas ativações no modelo.

Na Figura 4.20 é possível notar um ganho praticamente insignificante (e até negativo) ao usar o *gzip* para a Figura 4.7, enquanto na Figura 4.21 nota-se um ganho mais significativo com o *gzip* mesmo nos níveis mais superiores. Considerando que a Figura 4.7 possui muito conteúdo de alta frequência, os *bitstreams* gerados pelo modelo se tornam muito complexos. Enquanto, a Figura 4.19 possui um céu que faz parte da maior parte da imagem. É possível que os resíduos desses *patches* se tornem pequenos muito rápido de modo que o *gzip* é capaz de explorar bem esta grande quantidade de zeros.

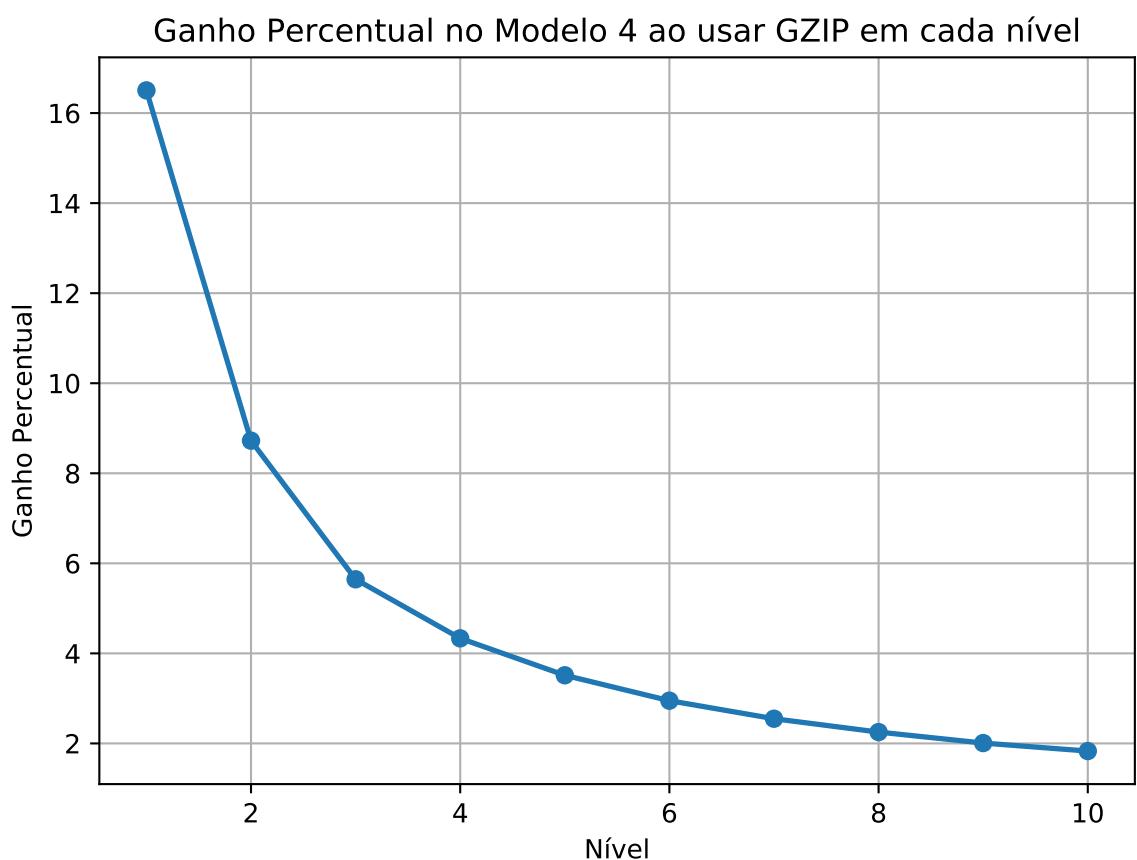


Figura 4.18: Ganho percentual médio na taxa por nível ao usar o codificador de entropia *gzip* nos *bitstreams* de cada nível para a base Kodak [2].



Figura 4.19: Imagem Kodim20 [2].

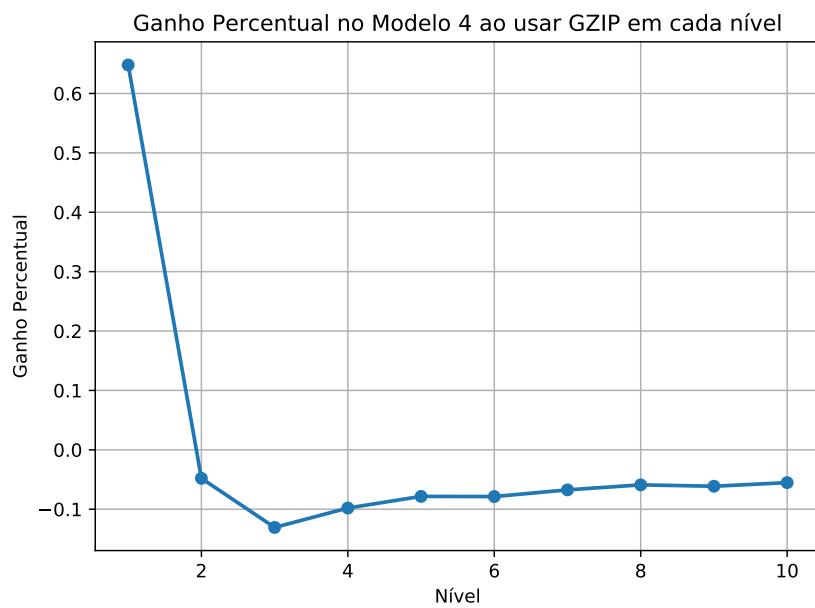


Figura 4.20: Ganho percentual na taxa por nível ao usar o codificador de entropia *gzip* nos *bitstreams* de cada nível para a Figura 4.7.

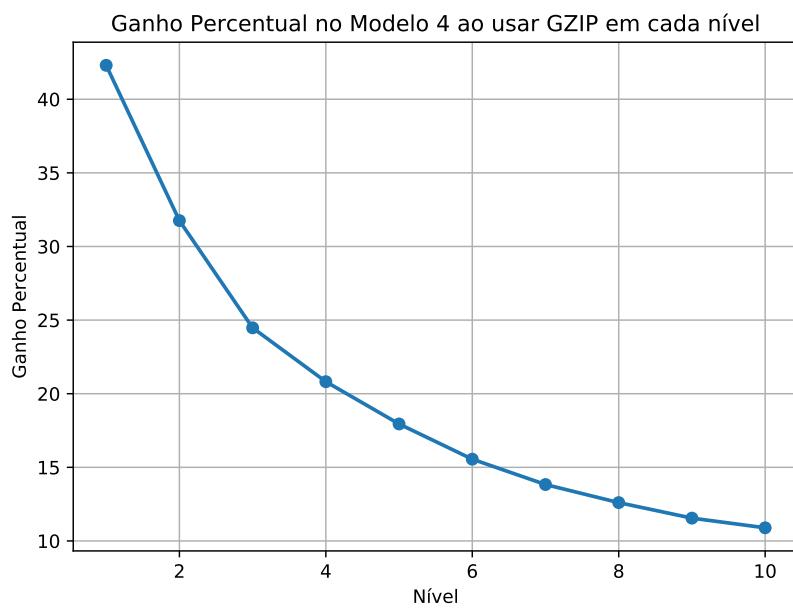


Figura 4.21: Ganho percentual na taxa por nível ao usar o codificador de entropia *gzip* nos *bitstreams* de cada nível para a Figura 4.19.

Capítulo 5

Conclusão

Este capítulo mostra o que se alcançou com os objetivos. É feita uma análise crítica na seção 5.2 para verificar a completude dos objetivos propostos pelo trabalho. A seção 5.3 indica as perspectivas futuras e os próximos passos a serem dados. A seção 5.1 aborda as limitações enfrentadas durante o desenvolvimento do trabalho.

5.1 Limitações do Trabalho

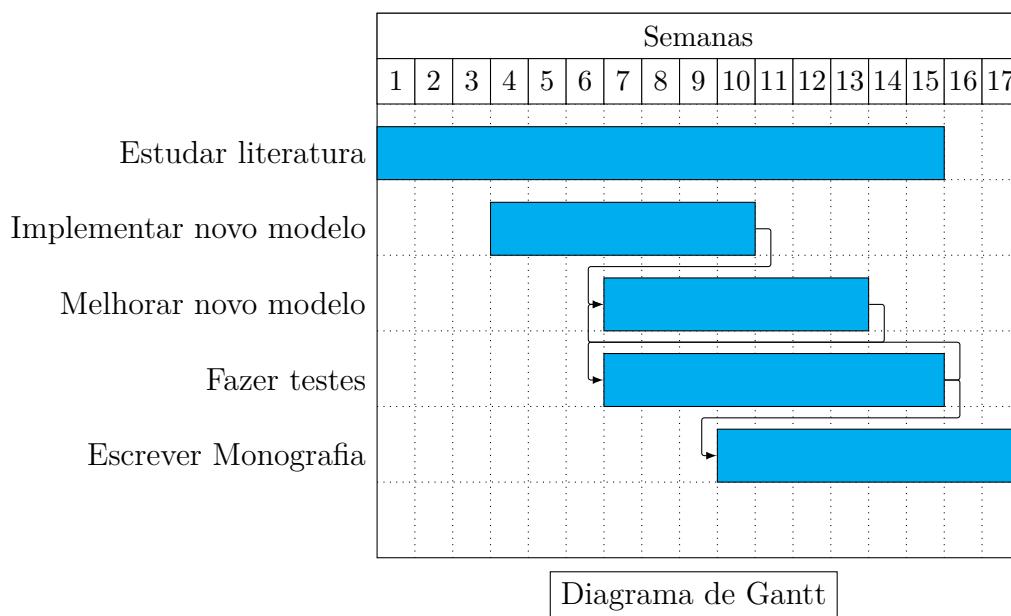
Modelos baseados em redes neurais para compressão de imagens precisam de muitas imagens para serem treinadas e possuem custo computacional superior aos codecs clássicos para codificar e decodificar imagens. Métodos baseados em redes neurais possuem o potencial de suprir uma necessidade crescente por algoritmos de compressão com perdas flexíveis. Entretanto, compressão com perdas é um problema não diferenciável. Em particular, quantização é uma parte integral do *pipeline* de compressão mas não é diferenciável, o que dificulta o trabalho de treinar redes neurais para esta tarefa.

5.2 Análise Crítica

Neste trabalho, não foi proposto um novo método para comprimir imagens mas foram replicadas técnicas já existentes na literatura alcançando resultados razoáveis para um trabalho que tinha como principal objetivo se familiarizar com a literatura e propostas existentes. No entanto, foi detectado uma necessidade de usar outros tipos de modelos para que seja possível superar os *codecs* clássicos no âmbito do tema deste trabalho. Esta necessidade é abordada na seção 5.3.

5.3 Trabalhos Futuros

Ainda há muito espaço para novas soluções no âmbito do problema abordado, portanto foram definidas algumas atividades chaves para dar sequência ao trabalho. Será estudado formas de melhorar o desempenho dos modelos usando técnicas propostas nos trabalhos apresentados em 2.2.3 com técnicas comumente usadas em redes neurais convolucionais. Segue um cronograma em formato de Diagrama de Gantt que organiza as atividades do próximo semestre letivo.



Referências

- [1] Wallace, G. K.: *The jpeg still picture compression standard*. IEEE Transactions on Consumer Electronics, 38(1):xviii–xxxiv, Feb 1992, ISSN 0098-3063. ix, 3, 8, 12
- [2] Kodak., E.: *Kodak lossless true color image suite*. <http://r0k.us/graphics/kodak/>, 2014. [Online; accessed 25-June-2019]. ix, x, xi, xii, xiii, 9, 25, 38, 39, 40, 41, 42, 43, 45, 46, 47, 48, 49
- [3] Commons, Wikimedia: *File:dct-8x8.png — wikimedia commons, the free media repository*. <https://commons.wikimedia.org/w/index.php?curid=10414002>, 2015. [Online; accessed 30-June-2019]. ix, 11
- [4] Computerphile: *The problem with jpeg - computerphile*. <https://youtu.be/yBX8GFqt6GA?t=48>, 2015. [Online; accessed 30-June-2019]. ix, 12
- [5] Goodfellow, Ian, Yoshua Bengio e Aaron Courville: *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>. ix, 13, 14
- [6] Materials, CS231n Course: *Cs231n convolutional neural networks for visual recognition*. <http://cs231n.github.io/neural-networks-3/>, 2019. [Online; accessed 30-June-2019]. ix, 15
- [7] Smith, Leslie N: *Cyclical learning rates for training neural networks*. Em *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*, páginas 464–472. IEEE, 2017. ix, 15, 16, 35
- [8] Zaghetto, Alexandre: *Intensity transformation and spatial filtering*. <https://github.com/zaghetto/ImageProcessing>, 2018. [Online; accessed 30-June-2019]. ix, 18
- [9] Toderici, George, Sean M. O’Malley, Sung Jin Hwang, Damien Vincent, David Minnen, Shumeet Baluja, Michele Covell e Rahul Sukthankar: *Variable Rate Image Compression with Recurrent Neural Networks*. Em *International Conference on Learning Representations (ICLR)*, páginas 1–12, nov 2015. <http://arxiv.org/abs/1511.06085>. ix, 4, 20, 21, 23, 26, 35
- [10] George Toderici, Michele Covell, Wenzhe Shi Radu Timofte Lucas Theis Johannes Ballé Eirikur Agustsson Nick Johnston Fabian Mentzer: *Challenge on learned image compression*. <http://www.compression.cc/challenge/>, 2018. [Online; accessed 25-June-2019]. x, xi, xii, xiii, 24, 34, 38, 41, 43, 44, 45, 46, 47

- [11] Sayood, Khalid: *Introduction to data compression*. Morgan Kaufmann, 2017. 1, 2
- [12] Shoham, Yair e Allen Gersho: *Efficient bit allocation for an arbitrary set of quantizers (speech coding)*. IEEE Transactions on Acoustics, Speech, and Signal Processing, 36(9):1445–1453, 1988. 2
- [13] Christopoulos, C., A. Skodras e T. Ebrahimi: *The jpeg2000 still image coding system: an overview*. IEEE Transactions on Consumer Electronics, 46(4):1103–1127, Nov 2000, ISSN 0098-3063. 3
- [14] Google, WebP: *Compression techniques*. <https://developers.google.com/speed/webp/docs/compression>. Acessed: 2019-06-16. 3
- [15] Bellard, Fabrice: *BPG image format*. <https://bellard.org/bpg/>. Accessed: 2019-06-13. 3
- [16] Simonyan, Karen e Andrew Zisserman: *Very deep convolutional networks for large-scale image recognition*. Relatório Técnico, Cornell University Library, 2014. arXiv:1409.1556. 4
- [17] Girshick, Ross, Jeff Donahue, Trevor Darrell e Jitendra Malik: *Rich feature hierarchies for accurate object detection and semantic segmentation*. Em *Proceedings of the IEEE conference on computer vision and pattern recognition*, páginas 580–587, 2014. 4
- [18] Choy, Christopher B, Danfei Xu, JunYoung Gwak, Kevin Chen e Silvio Savarese: *3d-r2n2: A unified approach for single and multi-view 3d object reconstruction*. Em *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016. 4
- [19] Taigman, Yaniv, Ming Yang, Marc'Aurelio Ranzato e Lior Wolf: *Deepface: Closing the gap to human-level performance in face verification*. Em *Proceedings of the IEEE conference on computer vision and pattern recognition*, páginas 1701–1708, 2014. 4
- [20] Graves, Alex e Navdeep Jaitly: *Towards end-to-end speech recognition with recurrent neural networks*. Em *International Conference on Machine Learning*, páginas 1764–1772, 2014. 4
- [21] Sutskever, Ilya, Oriol Vinyals e Quoc V Le: *Sequence to sequence learning with neural networks*. Em *Advances in neural information processing systems*, páginas 3104–3112, 2014. 4
- [22] Vinyals, Oriol, Alexander Toshev, Samy Bengio e Dumitru Erhan: *Show and tell: A neural image caption generator*. Em *Proceedings of the IEEE conference on computer vision and pattern recognition*, páginas 3156–3164, 2015. 4
- [23] Huval, Brody, Tao Wang, Sameep Tandon, Jeff Kiske, Will Song, Joel Pazhayampallil, Mykhaylo Andriluka, Pranav Rajpurkar, Toki Migimatsu, Royce Cheng-Yue et al.: *An empirical evaluation of deep learning on highway driving*. Relatório Técnico, Cornell University Library, 2015. arXiv:1504.01716. 4

- [24] Krizhevsky, Alex e Geoffrey E Hinton: *Using very deep autoencoders for content-based image retrieval*. Em *ESANN*, 2011. 4
- [25] Gregor, Karol, Frederic Besse, Danilo Jimenez Rezende, Ivo Danihelka e Daan Wierstra: *Towards conceptual compression*. Em *Advances In Neural Information Processing Systems*, páginas 3549–3557, 2016. 4
- [26] Toderici, George, Damien Vincent, Nick Johnston, Sung Jin Hwang, David Minnen, Joel Shor e Michele Covell: *Full resolution image compression with recurrent neural networks*. Em *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, páginas 5306–5314, 2017. Preprint available at <http://arxiv.org/abs/1608.05148>. 4, 23, 34
- [27] Mentzer, Fabian, Eirikur Agustsson, Michael Tschannen, Radu Timofte e Luc Van Gool: *Practical full resolution learned lossless image compression*. Em *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 4
- [28] Theis, Lucas e Matthias Bethge: *Generative image modeling using spatial lstms*. Em *Advances in Neural Information Processing Systems*, páginas 1927–1935, 2015. 4
- [29] Hamilton, Eric: *Jpeg file interchange format*. 2004. 7
- [30] Ahmed, Nasir, T_ Natarajan e Kamisetty R Rao: *Discrete cosine transform*. IEEE transactions on Computers, 100(1):90–93, 1974. 8
- [31] Huffman, David A: *A method for the construction of minimum-redundancy codes*. Proceedings of the IRE, 40(9):1098–1101, 1952. 13
- [32] Pennebaker, WB, JL Mitchell *et al.*: *Arithmetic coding articles*. IBM J. Res. Dev, 32(6):717–774, 1988. 13
- [33] Mor-Yosef, Shlomo, Arnon Samueloff, Baruch Modan, Daniel Navot e Joseph G Schenker: *Ranking the risk factors for cesarean: logistic regression analysis of a nationwide study*. Obstetrics and gynecology, 75(6):944–947, 1990. 13
- [34] Nielsen, Michael A.: *Neural Networks and Deep Learning*. Determination Press, 2015. 14
- [35] Rumelhart, David E, Geoffrey E Hinton, Ronald J Williams *et al.*: *Learning representations by back-propagating errors*. Cognitive modeling, 5(3):1, 1988. 14, 20
- [36] Kingma, Diederik P e Jimmy Ba: *Adam: A method for stochastic optimization*. Relatório Técnico, Cornell University Library, 2014. arXiv:1412.6980. 16
- [37] LeCun, Yann, Koray Kavukcuoglu e Clément Farabet: *Convolutional networks and applications in vision*. Em *Proceedings of 2010 IEEE International Symposium on Circuits and Systems*, páginas 253–256. IEEE, 2010. 16
- [38] Academy, Data Science: *Deep learning book*. <http://www.deeplearningbook.com.br>, 2019. [Online; accessed 30-June-2019]. 17

- [39] Wang, Zhou, Alan C Bovik, Hamid R Sheikh, Eero P Simoncelli *et al.*: *Image quality assessment: from error visibility to structural similarity*. IEEE transactions on image processing, 13(4):600–612, 2004. 19
- [40] Wang, Zhou, Eero P Simoncelli e Alan C Bovik: *Multiscale structural similarity for image quality assessment*. Em *The Thirly-Seventh Asilomar Conference on Signals, Systems & Computers, 2003*, volume 2, páginas 1398–1402. Ieee, 2003. 19
- [41] Ma, Siwei, Xinfeng Zhang, Chuanmin Jia, Zhenghui Zhao, Shiqi Wang e Shanshe Wang: *Image and video compression with neural networks: A review*. IEEE Transactions on Circuits and Systems for Video Technology, PP:1–1, abril 2019. 22, 23
- [42] Paszke, Adam, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga e Adam Lerer: *Automatic differentiation in pytorch*. 2017. 24
- [43] Agustsson, Eirikur e Radu Timofte: *Ntire 2017 challenge on single image super-resolution: Dataset and study*. Em *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, July 2017. 24
- [44] Group, Multimedia Signal Processing: *Ultra-eye: Uhd and hd images eye tracking dataset*. <https://mmspgr.epfl.ch/downloads/ultra-eye/>, 2014. [Online; accessed 25-June-2019]. 24