



**Universidade Federal do Amazonas – UFAM**  
**Instituto de Computação – IComp**  
**Programa de Pós-Graduação em Informática – PPGI**

**Uma Abordagem para Reconhecimento de Emoção por Expressão Facial baseada  
em Redes Neurais de Convolução**

Anderson Araújo da Cruz

Manaus – AM

Julho, 2018



Anderson Araújo da Cruz

Uma Abordagem para Reconhecimento de Emoção por Expressão Facial baseada em Redes  
Neurais de Convolução

Texto de qualificação submetido ao Programa de Pós-Graduação em Informática do Instituto de Computação da Universidade Federal do Amazonas para obtenção de pré-requisito para o título de mestrado *stricto sensu*.

Orientador: Raimundo da Silva Barreto,  
D.Sc.

Manaus – AM

Julho, 2018



PODER EXECUTIVO  
MINISTÉRIO DA EDUCAÇÃO  
UNIVERSIDADE FEDERAL DO AMAZONAS  
INSTITUTO DE COMPUTAÇÃO  
PROGRAMA DE PÓS-GRADUAÇÃO EM INFORMÁTICA

---

## **FOLHA DE APROVAÇÃO**

“Uma Abordagem para Reconhecimento de Emoção por Expressão Facial baseada em Redes Neurais de Convolução”

Anderson Araújo da Cruz

Qualificação defendida e aprovada pela banca examinadora constituída pelos Professores:

---

**Raimundo da Silva Barreto, D.Sc.–  
Presidente**

---

**Elaine Harada Teixeira de Oliveira,  
D.Sc.– Membro**

---

**Daniel Lins da Silva, D.Sc.– Membro**

Manaus – AM  
Julho, 2018

*“Os loucos que acham que podem mudar o mundo,  
são os que efetivamente o fazem”*

Comercial Apple - “Pense Diferente”, 1997



# Resumo

Desenvolver a percepção dos computadores é uma tendência tecnológica. O reconhecimento de emoção contribui para a composição dos sistemas cognitivos e inteligentes possuindo aplicabilidade em diversas áreas. Esta proposta apresenta uma abordagem para reconhecer emoções por expressão facial utilizando redes neurais de convolução. Por meio das expressões faciais é possível classificar o grupo das emoções básicas. A expressão facial é uma maneira efetiva para reconhecer emoções, sobretudo por não ser intrusiva na coleta de dados e pela facilidade de obter imagens da face. As redes neurais de convolução tem sido o estado da arte para classificação de imagens. Um estudo experimental foi conduzido a fim de comparar três arquiteturas: AlexNet, Inception e ResidualNet. Este estudo consistiu no reconhecimento de emoções em dois cenários distintos: imagens fotografadas em laboratório e na natureza. A ResidualNet avaliada pelas métricas de precisão, revocação, f1-score e acurácia, obteve os melhores resultados. Em imagens laboratoriais, a ResidualNet aproximou-se do baseline. O que não aconteceu em imagens oriundas da natureza. De um modo geral, as emoções felicidade e surpresa foram as que tiveram maiores taxas de precisão. Em contrapartida, as emoções medo, raiva e tristeza, alcançaram as menores taxas. Contudo, o objetivo é ter uma solução robusta para lidar com situações complexas de classificação.

**Palavras-chaves:** Reconhecimento de Emoção, Expressão Facial, Redes Neurais de Convolução.





# Lista de ilustrações

Figura 1 – Solução Proposta . . . . .	6
Figura 2 – Expressão facial emocional . . . . .	8
Figura 3 – Detecção Facial . . . . .	10
Figura 4 – Rede Neural Artificial . . . . .	12
Figura 5 – Camada de Convolução com campos locais de recepção . . . . .	14
Figura 6 – Conectividade esparsa. É destacada a entrada $x_3$ e a saída em S que são afetadas por $x_3$ . (Cima) Quando S recebe a convolução com um <i>kernel</i> de tamanho 3, somente três saídas são afetadas por $x_3$ . (Baixo) Quando S é gerado por rede neural tradicional, todos são afetados por $x_3$ . . . . .	14
Figura 7 – Camada de <i>Max Pooling</i> . . . . .	15
Figura 8 – Módulo <i>inception</i> . . . . .	17
Figura 9 – Bloco residual . . . . .	19
Figura 10 – Extração dos pontos faciais para características geométrica . . . . .	25
Figura 11 – Concatenação dos pontos faciais com uma rede neural de convolução . . . . .	25
Figura 12 – Extração das sub-regiões faciais para características aparente . . . . .	26
Figura 13 – Algoritmo de Alinhamento da Face: (a) Imagem Original; (b) Face Alinhada. . . . .	33
Figura 14 – Pré-Processamento fluxo . . . . .	33
Figura 15 – Exemplos da técnica de aumento de dados: (a) Imagem Original; (b) Redução de Contraste; (c) Aumento de Contraste; (d) Perspectiva; (e) <i>Crop</i> ; (f) <i>Shear</i> . . . . .	34
Figura 16 – Solução Proposta . . . . .	36
Figura 17 – Arquitetura Inception-V3 . . . . .	43
Figura 18 – Arquitetura ResNet-34 . . . . .	44
Figura 19 – Gráfico de Acurácia na Base de Teste. As linhas verticais indicam o melhor modelo. . . . .	44
Figura 20 – Gráfico da Função de Perda na Base de Teste. As linhas verticais indicam o melhor modelo. . . . .	45
Figura 21 – Artigos por ano retornados pela <i>string</i> de busca . . . . .	67
Figura 22 – Gráfico representando a frequência dos critérios de inclusão para os artigos aceitos do segundo filtro . . . . .	68
Figura 23 – Gráfico representando a frequência dos critérios de exclusão para os artigos rejeitados do segundo filtro . . . . .	68



# Lista de tabelas

Tabela 1 – Arquitetura AlexNet . . . . .	16
Tabela 2 – Arquiteturas VGGNet . . . . .	18
Tabela 3 – Principais arquiteturas de Redes Neurais de Convolução para reconhecimento de expressões faciais. (*) Significa que a rede foi treinada ( <i>fine-tuning</i> ) por duas vezes. . . . .	28
Tabela 4 – Resultado da correlação de Pearson para cada emoção detectada e a entropia contra os atributos das questões . . . . .	40
Tabela 5 – Bases de dados localizadas na revisão sistemática da literatura . . . . .	42
Tabela 6 – As bases de dados foram concatenadas e divididas em três bases: treino, teste e validação. Na seguinte porcentagem: 50% para treino e 25% para teste e validação. . . . .	42
Tabela 7 – Distribuição das classes (emoções) nas bases de treino, teste e validação. As classes também foram divididas em: 50% para treino e 25% para teste e validação. . . . .	43
Tabela 8 – Resultados experimentais das redes neurais de convolução avaliando a base de validação geral. . . . .	45
Tabela 9 – Resultados experimentais das redes neurais de convolução avaliando a base de validação CK . . . . .	46
Tabela 10 – Resultados experimentais das redes neurais de convolução avaliando a base de validação FER . . . . .	47
Tabela 11 – Cronograma de Atividades . . . . .	54
Tabela 12 – Objetivos da Revisão Sistemática . . . . .	63
Tabela 13 – Principais arquiteturas de redes neurais de convolução e os trabalhos que utilizaram . . . . .	70
Tabela 14 – Bases de Dados . . . . .	71
Tabela 15 – Resultados da Arquitetura AlexNet avaliando a base de validação geral	76
Tabela 16 – Resultados da Arquitetura InceptionV3 avaliando a base de validação geral . . . . .	77
Tabela 17 – Resultados da Arquitetura ResNet avaliando a base de validação geral	78
Tabela 18 – Resultados experimentais das redes neurais de convolução avaliando a base de validação CIFE-Test . . . . .	80
Tabela 19 – Resultados experimentais das redes neurais de convolução avaliando a base de validação CIFE-Train . . . . .	81
Tabela 20 – Resultados experimentais das redes neurais de convolução avaliando a base de validação CK . . . . .	82

Tabela 21 – Resultados experimentais das redes neurais de convolução avaliando a base de validação FER . . . . .	83
Tabela 22 – Resultados experimentais das redes neurais de convolução avaliando a base de validação JAFFE . . . . .	84
Tabela 23 – Resultados experimentais das redes neurais de convolução avaliando a base de validação KDEF . . . . .	85
Tabela 24 – Resultados experimentais das redes neurais de convolução avaliando a base de validação NovaEmotions . . . . .	86
Tabela 25 – Resultados experimentais das redes neurais de convolução avaliando a base de validação RAFD . . . . .	87

# Lista de abreviaturas e siglas

RNC	Rede Neural de Convolução.
ILSVRC	ImageNet Large Scale Visual Recognition Challenge.



# Sumário

<b>1</b>	<b>INTRODUÇÃO</b>	<b>1</b>
<b>1.1</b>	<b>Contexto</b>	<b>1</b>
<b>1.2</b>	<b>Motivação</b>	<b>1</b>
<b>1.3</b>	<b>Definição do Problema</b>	<b>2</b>
<b>1.4</b>	<b>Objetivos</b>	<b>3</b>
1.4.1	Objetivo Geral	3
1.4.2	Objetivos Específicos	3
<b>1.5</b>	<b>Hipótese</b>	<b>3</b>
<b>1.6</b>	<b>Abordagem Proposta</b>	<b>3</b>
1.6.1	Visão Geral da Solução Proposta	3
1.6.2	Computação embarcada e em nuvem	4
<b>1.7</b>	<b>Organização do Trabalho</b>	<b>5</b>
<b>2</b>	<b>REFERENCIAL TEÓRICO</b>	<b>7</b>
<b>2.1</b>	<b>Reconhecimento de Emoção</b>	<b>7</b>
<b>2.2</b>	<b>Expressão Facial Emocional</b>	<b>8</b>
<b>2.3</b>	<b>Aprendizagem de Máquina</b>	<b>9</b>
<b>2.4</b>	<b>Processo de Classificação de Imagem</b>	<b>9</b>
<b>2.5</b>	<b>Técnicas para pré-processamento de imagens em reconhecimento de emoção por expressão facial</b>	<b>10</b>
2.5.1	Detecção Facial	10
2.5.2	Histograma de Equalização (Normalização do Brilho)	10
<b>2.6</b>	<b>Rede Neural Artificial</b>	<b>11</b>
<b>2.7</b>	<b>Rede Neural de Convolução</b>	<b>12</b>
2.7.1	Camada de Convolução	13
2.7.2	Camada de <i>Pooling</i>	13
2.7.3	Regressão <i>Softmax</i>	15
<b>2.8</b>	<b>Arquiteturas de Redes Neurais de Convolução</b>	<b>16</b>
2.8.1	AlexNet	16
2.8.2	GoogLeNet	17
2.8.3	VGGNet	18
2.8.4	Residual Network	19
2.8.5	<i>Ensemble</i>	19
<b>2.9</b>	<b>Métricas de Avaliação de Desempenho para Classificadores</b>	<b>20</b>
<b>2.10</b>	<b>Resumo</b>	<b>21</b>

<b>3</b>	<b>TRABALHOS CORRELATOS</b>	<b>23</b>
<b>3.1</b>	<b>Preparação dos dados</b>	<b>23</b>
<b>3.2</b>	<b>Extração de Característica</b>	<b>24</b>
3.2.1	Extração Geométrica	24
3.2.2	Extração Aparente	25
<b>3.3</b>	<b>Arquiteturas</b>	<b>26</b>
3.3.1	AlexNet	26
3.3.2	VGG	26
3.3.3	GoogLeNet	27
3.3.4	Ensemble	27
<b>3.4</b>	<b>Aplicações</b>	<b>27</b>
<b>3.5</b>	<b>Resumo</b>	<b>29</b>
<b>4</b>	<b>ABORDAGEM PROPOSTA</b>	<b>31</b>
<b>4.1</b>	<b>Monitoramento do Indivíduo</b>	<b>31</b>
<b>4.2</b>	<b>Detecção de Face e Recorte</b>	<b>31</b>
<b>4.3</b>	<b>Pré-Processamento</b>	<b>32</b>
<b>4.4</b>	<b>Rede Neural de Convolução</b>	<b>33</b>
4.4.1	Treinamento	34
4.4.2	Extração de Características e Classificação	34
<b>4.5</b>	<b>Resumo</b>	<b>35</b>
<b>5</b>	<b>RESULTADOS PARCIAIS</b>	<b>37</b>
<b>5.1</b>	<b>Prova de Conceito: Coleta e Inferência de Estados Emocionais de Estudantes</b>	<b>37</b>
5.1.1	Síntese	37
5.1.2	Objetivos	38
5.1.3	Metodologia Experimental	38
5.1.3.1	Planejamento	38
5.1.3.2	Execução	39
5.1.4	Resultados e Discussões	39
<b>5.2</b>	<b>Avaliação Experimental de Redes Neurais de Convolução</b>	<b>40</b>
5.2.1	Preparação dos Dados	40
5.2.2	Materiais	41
5.2.3	Arquiteturas	42
5.2.4	Treinamento	43
5.2.5	Discussões	46
<b>5.3</b>	<b>Resumo</b>	<b>50</b>
<b>6</b>	<b>CONSIDERAÇÕES FINAIS</b>	<b>53</b>



<b>6.1</b>	<b>Limitações do Trabalho</b>	<b>53</b>
<b>6.2</b>	<b>Trabalhos Futuros</b>	<b>53</b>
<b>6.3</b>	<b>Cronograma</b>	<b>54</b>

	<b>Referências</b>	<b>55</b>
--	--------------------	-----------

## **ANEXOS** **61**

	<b>ANEXO A – REVISÃO SISTEMÁTICA DA LITERATURA</b>	<b>63</b>
<b>A.1</b>	<b>Protocolo da Revisão Sistemática da Literatura</b>	<b>63</b>
A.1.1	Objetivo	63
A.1.2	Questões de Pesquisa	63
A.1.3	Biblioteca Digital	64
A.1.4	CrITÉrios de Inclusão e Exclusão dos Artigos	64
A.1.5	Formulário de Extração de Informação	65
A.1.6	<i>String</i> de Busca	66
<b>A.2</b>	<b>Condução da Revisão Sistemática da Literatura</b>	<b>67</b>
A.2.1	Primeiro Filtro	67
A.2.2	Segundo Filtro	67
<b>A.3</b>	<b>Resultados</b>	<b>69</b>
A.3.1	Q1: Quais emoções têm sido reconhecidas por meio da expressão facial utilizando redes neurais de convolução?	69
A.3.2	Q2: Quais tipos de pré-processamento tem sido realizado na imagem?	69
A.3.3	Q3: Quais arquiteturas de redes neurais de convolução têm sido mais utilizadas?	70
A.3.4	Q4: Quais técnicas, métodos e abordagens têm sido utilizados para tratar problemas na imagem como iluminação, rotação, obstrução e escala?	70
A.3.5	Q5: Quais bases de dados têm sido utilizadas?	70
A.3.6	Q6: Quais aplicações podem utilizar o reconhecimento de emoção por expressão facial?	71
A.3.7	Questão Principal: Como reconhecer emoções por meio da expressão facial utilizando redes neurais de convolução em uma imagem estática?	72
<b>A.4</b>	<b>Resumo</b>	<b>73</b>
	<b>ANEXO B – RESULTADOS POR BASE DE VALIDAÇÃO GERAL</b>	<b>75</b>
	<b>ANEXO C – RESULTADOS POR BASES DE VALIDAÇÃO</b>	<b>79</b>



# 1 Introdução

## 1.1 Contexto

Há décadas a comunidade científica está interessada no reconhecimento de emoções. As diversas maneiras de expressar as emoções humanas têm sido investigadas, tais como sinais fisiológicos, textos, envio de *emojicons*, dispositivo padrão de entrada de dados (e.g. teclado e mouse), voz e as expressões faciais. Esta última surgiu pelas anotações de [Darwin \(1965\)](#) e experiências de [Ekman and Davidson \(1994\)](#). Foi percebido que todas as culturas emitem emoção pela expressão facial, de tal forma que há um grupo de emoções básicas (raiva, felicidade, tristeza, desprezo, medo e surpresa) que possuem a mesma expressão facial independente da cultura dos indivíduos. Apesar de tantos anos de pesquisa, a comunidade continua interessada neste assunto, pois reconhecer emoção tem sido desafiador. Além de ser um campo promissor para a interação humano-computador e humano-robô. Vale ressaltar que ainda é considerado um problema em aberto, inclusive com a realização de concursos com premiação, como foi o caso do ICML'2013 e anualmente como o EmotiW dos anos de 2014 a 2017.

O progresso da área de aprendizado profundo ocasionou o surgimento de diversas técnicas poderosas de reconhecimento de padrões, gerando grande destaque para as redes neurais de convolução. Esta técnica foi projetada para processamento, extração de característica e classificação de imagem. Ultimamente, as redes neurais de convolução têm sido bastantes populares e utilizadas em diversos contextos dominando amplamente os trabalhos realizados pela comunidade em classificação de imagens. Possibilitando, inclusive, o reconhecimento automático de emoção por meio da expressão facial, sendo que tal reconhecimento está próximo do que um humano reconheceria ([Kim et al., 2016a](#)). Estes resultados expressivos têm motivado pesquisadores a continuar aprimorando estas técnicas. Sobretudo na expressão facial que tem se tornado uma abordagem eficaz para reconhecer emoções, pois não é uma abordagem intrusiva de coleta de dados quando comparada aos sensores fisiológicos. Estes sensores não são uma computação ubíqua resultando no incomodo do usuário quando seus sinais fisiológicos são monitorados. Além disso, a expressão facial possui a facilidade de ser obtida em uma captura de imagem devido a popularidade de dispositivos que possuem câmeras fotográficas ([Cruz et al., 2017](#)).

## 1.2 Motivação

O reconhecimento de emoção tem aplicação em muitas áreas. Destacamos alguns campos promissores. Na educação, por exemplo, segundo [Jaques and Nunes \(2013\)](#), estu-

dantes durante o seu processo de aprendizagem emitem constantemente diversas emoções. Portanto, sistemas educacionais como Ambientes Virtuais de Aprendizagem (AVA) e Sistemas de Tutores Inteligentes (STI) podem monitorar as emoções durante a interação com uma plataforma educacional em uma aula. Com intuito de fornecer *feedback* personalizado para o estudante através da recomendação de objetos de aprendizagem apropriados para aquele estado emocional e até mesmo, realizar ações que estimulem emoções positivas a fim de motivar os estudantes quando estes estiverem em um estado negativo. Outra área de aplicação para utilizar o reconhecimento de emoção é em realidade virtual. Segundo Riva et al. (2007), a realidade virtual pode estimular propositalmente emoções permitindo maior imersão do usuário à aplicação. Desta forma, o reconhecimento de emoção pode medir o quão efetivo tem sido o método de estímulo de emoções ao usuário e, caso não seja satisfatório, o método pode ser alterado. Para Li et al. (2015a), o reconhecimento de emoção pode auxiliar na construção de tecnologias assistivas para deficientes visuais que, quando possuem elevado grau de deficiência, apresentam dificuldades em reconhecer emoções na interação interpessoal. Em geral, é possível aplicar o reconhecimento de emoção na interação humano computador (Barsoum et al., 2016a; Chen et al., 2017a; Liu et al., 2016a; Wen et al., 2017a), e interação humano robô (Jung et al., 2015a; Shin et al., 2016a), criando a expectativa de que computadores do futuro possam reconhecer a emoção do usuário e realizar algum procedimento que ocasione maior aproximação entre homem e máquina.

### 1.3 Definição do Problema

Trata-se de um problema de classificação de imagem digital no qual há uma imagem  $\omega$  formada por um conjunto de *pixels* (RGB)  $\alpha$  pertencente a um conjunto de classes  $\sigma = \{\text{neutralidade, raiva, felicidade, tristeza, desprezo, medo e surpresa}\}$ , que são as emoções básicas definidas por (Ekman and Davidson, 1994), tal que haja uma função  $\phi$  que saiba mapear  $\omega$  por meio de  $\alpha$  para  $\sigma$ .

Embora existam trabalhos que classifiquem emoções em imagens (Barsoum et al., 2016a; Kim et al., 2016a; Yu et al., 2016b), pouca atenção tem sido dada aos problemas clássicos em imagens como: (i) ausência de iluminação no ambiente; (ii) rotação do objeto principal, neste caso a face, e (iii) escala do objeto principal (face). Abordagens que tratam estes problemas em imagens são mais apropriadas para o uso em cenários reais, no qual a exigência para classificação é maior devido as condições adversas do ambiente e pelas diferentes variações das características da face humana.

O problema considerado neste trabalho pode ser expresso na seguinte questão: *Como aprimorar os métodos de reconhecimento de emoções por meio da expressão facial a fim de permitir a classificação independente das características do ambiente e de indivíduos*

*para o alcance de maior generalização?*

## 1.4 Objetivos

### 1.4.1 Objetivo Geral

Propor um método para reconhecer emoção humana por expressão facial para classificar emoções básicas em múltiplas faces de uma imagem e comparar a eficácia em cenários de uso real.

### 1.4.2 Objetivos Específicos

- Propor técnicas de eliminação de ruídos e detecção com recorte das diversas faces de uma imagem;
- Classificar cada face detectada separadamente estimando a probabilidade para cada emoção básica;
- Avaliar experimentalmente a solução proposta visando a comparação da eficácia.

## 1.5 Hipótese

As emoções básicas emitidas por expressão facial podem ser reconhecidas por uma rede neural de convolução, desde que esteja treinada e validada por instâncias representativas do problema (veja Seção 1.3). Além disso, este trabalho apoia-se na combinação entre a rede neural de convolução com técnicas de pré-processamento, e também, o aumento de dados durante o treinamento. Essas técnicas aplicam nas imagens: normalizações da iluminação, eliminação de ruídos, geração de novas instâncias representativas e o alinhamento da face. Com intuito de gerar modelos que maximizam a acurácia e a generalização do aprendizado relacionada ao reconhecimento de emoção em diferentes ambientes e variações da face humana.

## 1.6 Abordagem Proposta

### 1.6.1 Visão Geral da Solução Proposta

Uma visão geral da solução proposta é apresentada na Figura 1. Há algum dispositivo de monitoramento que contém câmera fotográfica (e.g. smartphone, notebook, televisão e outros) que periodicamente captura imagens. Toda fotografia capturada pelo dispositivo é salva em um repositório de entrada de dados. Este repositório é consultado

para recuperar uma imagem e enviá-la para a classificação. No processo de classificação, é verificada a existência de uma face na imagem e, caso não exista, é encerrada a execução. Pois se trata de uma imagem que não contém uma face, portanto, não existe uma expressão facial para classificar. Caso exista uma face, uma função para recortá-la é chamada. Este procedimento é valioso por dois aspectos: o primeiro por excluir o *background* da imagem, sendo assim o classificador não necessita aprender a diferenciar o que é face e *background*, e o segundo é pela possibilidade de haver múltiplas faces na imagem realizando a classificação de cada face individualmente, reduzindo a complexidade do problema, pois é mais fácil classificar uma face por vez do que várias ao mesmo tempo. Posteriormente, a face recortada é enviada a um conjunto de filtros de pré-processamento que por sua vez operam sobre a imagem para eliminação de ruídos, normalizações do contraste e alinhamento. Finalmente, a imagem pré-processada é enviada a uma rede neural de convolução para a extração de características e classificação. Esta técnica pode ter uma característica em particular que é interessante: em vez de retornar à classe a que a expressão facial (ou uma instância) pertence, pode retornar às estimativas de probabilidade para cada classe da expressão facial (e.g. neutralidade: 0.95, felicidade: 0.025, medo: 0.025,...) possuindo assim uma propriedade em que a soma das probabilidades de todas as classes é igual a 1. Sendo assim, a classificação da instância em análise seria a classe com maior probabilidade estimada (neutralidade com 0.95 de certeza). Em nossa abordagem, após a geração das estimativas de probabilidades, as mesmas são salvas em um repositório de saída de dados disponibilizando o resultado para aplicação solicitante. O processo anteriormente descrito deve ser repetido enquanto houver faces para classificar, isto é, quando uma imagem há múltiplas faces, e cada face é classificada uma por vez.

### 1.6.2 Computação embarcada e em nuvem

Esta proposta visa fornecer soluções para reconhecimento de emoção que contempla dois tipos de computação: em nuvem e embarcada. Tais computações estão em alta na academia, indústria e mercado. A primeira pelo crescimento da internet havendo bilhões de dispositivos conectados, e também, pela a evolução da infraestrutura com aumento considerável de recursos computacionais e velocidade de conexão. Já a computação embarcada é principalmente pela explosão de dispositivos presentes em nosso cotidiano. Além disso, os dispositivos embarcados de hoje tem uma autonomia energética periódica, hardware semelhante a desktops, sistemas operacionais e sensores embutidos, formando um recurso independente e poderoso. Há no mercado vários smartphones e smartwatches com processadores octa-core e dual-core, respectivamente, com a memória RAM chegando a 8GB, e até mesmo com placas de vídeos embutidas para aceleração da computação.

A ideia é que a computação em nuvem hospedaria o melhor modelo gerado a partir das arquiteturas AlexNet, VGGNet, GoogLeNet e Residuais, considerando as mé-

tricas de avaliação de desempenho como acurácia, precisão, revocação e f1-score. Apesar de que o melhor modelo possa exigir elevada utilização de recursos computacionais por ser uma rede neural profunda, entende-se que um serviço em nuvem possuiria um hardware robusto capaz de suportar a demanda. Visto que estamos diante do crescimento da computação cognitiva, isto é, a capacidade de computadores tomar decisões como humanos. Este trabalho quer apoiar estas aplicações oferecendo um reconhecedor de emoções para a integração em um serviço em nuvem, onde a aplicação solicitante independente da linguagem de programação esteja implementada ou em qual plataforma está rodando, possam comunicar-se via padrão REST com intuito de enviar imagens e receber as emoções detectadas.

Em aplicações com características de computação embarcada, este trabalho propõe uma rede neural de convolução baseada na arquitetura MobileNet para funcionar nativamente em um dispositivo embarcado. Essa arquitetura tem características de ser mais enxuta sendo projetada para consumir menos recursos computacionais, e além disso, com o enxugamento possui o compromisso de perder o mínimo de precisão comparada as arquiteturas robustas. Portanto, ideal para dispositivos embarcados que dispõe de menos recursos computacionais. Inclusive podendo funcionar nativamente no sistema operacional Android que é amplamente usado por *smartphones*, *smartwatches* e *tablets*. Além disso, a arquitetura MobileNet pode ser também embutida em placas de desenvolvimento como Raspberry, Nvidia Jetson, Drones e outras.

Vale destacar que as maiores taxas de ocupação de recursos computacional de uma rede neural de convolução estão na fase de treinamento, isto é, durante a geração do modelo. Logo, a fase de classificação exige menores taxas de ocupação de hardware, pois o principal procedimento que demanda recursos computacionais consiste em carregar o modelo na memória. Para classificar uma imagem, considerando que o modelo está carregado na memória, a rede neural opera sobre a imagem aplicando os pesos recuperados do modelo. Os pesos configuram-se como o aprendizado retido. Em contrapartida, a fase de treinamento exige muito mais processamento, pois é executado o algoritmo de otimização gradiente descendente para minimizar a função de perda e adquirir aprendizado, realizando uma alta quantidade de cálculos vetoriais.

## 1.7 Organização do Trabalho

Este trabalho está dividido nos capítulos a seguir. O Capítulo 2 aborda os conceitos e definições necessários para o entendimento deste trabalho. O Capítulo 3 analisa os trabalhos relacionados. O Capítulo 4 apresenta a abordagem proposta. O Capítulo 5 discute os resultados parciais obtidos, enquanto o Capítulo 6 enfatiza as considerações finais, limitações do trabalhos e os trabalhos futuros.

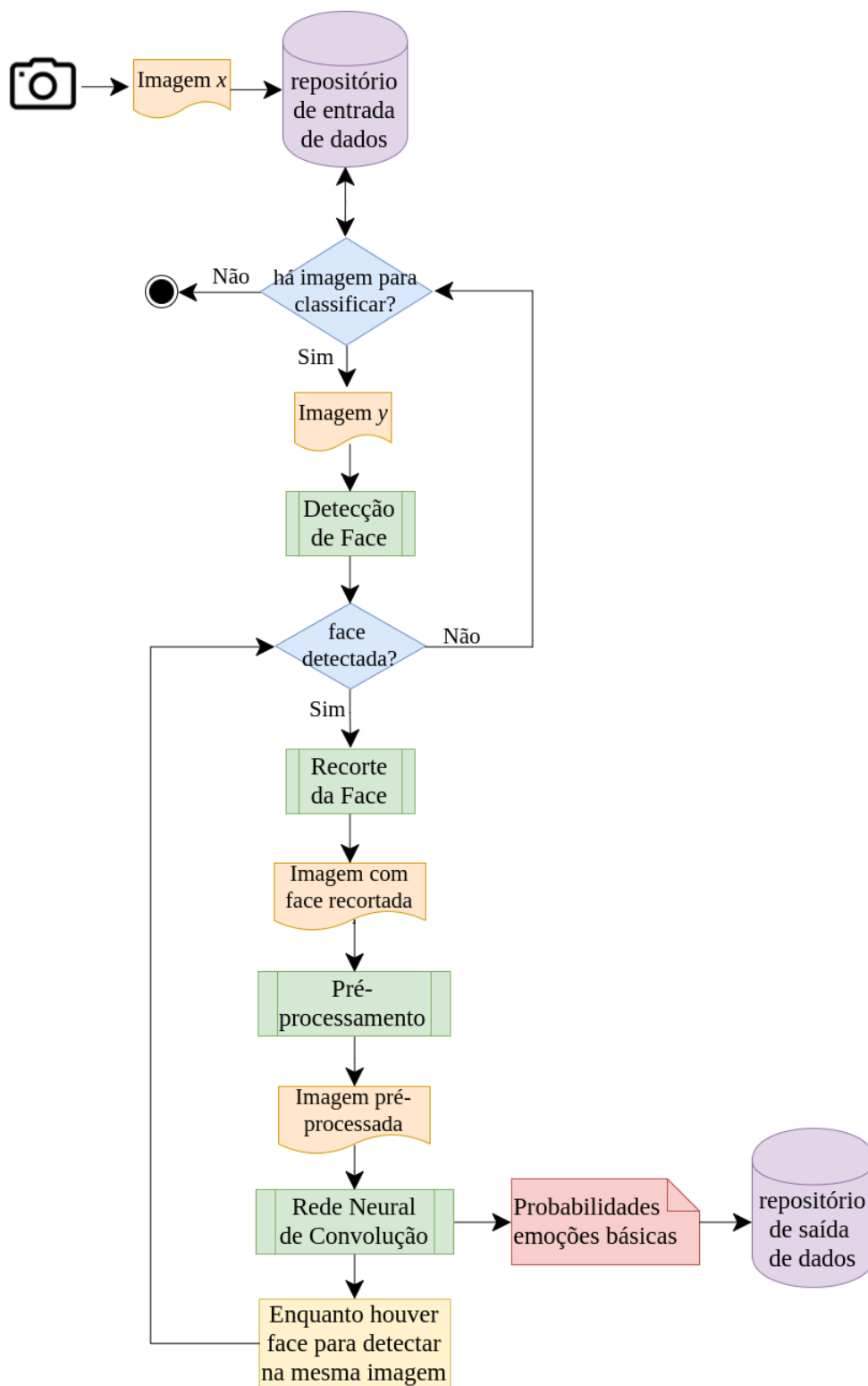


Figura 1 – Solução Proposta



## 2 Referencial Teórico

Neste capítulo são introduzidos os conceitos necessários para o entendimento deste trabalho e está organizado da seguinte forma. A Seção 2.1 define o que é o reconhecimento de emoção e os tipos de reconhecimento. A Seção 2.2 fundamenta a expressão facial emocional e mostra as expressões básicas existentes. A Seção 2.3 conceitua a área de aprendizagem de máquina. A Seção 2.4 apresenta o processo de classificação de imagem. A Seção 2.5 revisa as técnicas de pré-processamento que são utilizadas para eliminação de ruídos de uma imagem. As Seções 2.6 e 2.7 conceituam os aspectos básicos das redes neurais artificiais e das redes neurais de convolução, respectivamente. A Seção 2.8 descreve as arquiteturas de redes neurais de convolução utilizadas por este trabalho. A Seção 2.9 retrata as principais métricas para avaliar classificadores e para concluir a Seção 2.10 faz um resumo acerca deste capítulo.

### 2.1 Reconhecimento de Emoção

As emoções podem ser definidas como breves e intensas e são disparadas pela avaliação de um evento (Scherer, 2000). O reconhecimento de emoção tem sido explorado há algumas décadas e as emoções que tem sido frequentemente investigadas são as básicas como a raiva, alegria, tristeza, desgosto, medo e surpresa (Ekman and Davidson, 1994). Com objetivo de reconhecer emoção, os pesquisadores têm utilizado comumente as técnicas de reconhecimento de padrões, que por sua vez podem encontrar as características que são importantes para diferenciar as emoções, isto é, a busca pelo padrão relevante de uma entrada de dados, por exemplo uma imagem contendo uma expressão facial, de modo que a técnica consiga diferenciar a felicidade de neutralidade. Além disso, a comunidade tem gerado várias heurísticas para reconhecer emoções, entretanto com emprego somente em ambientes muito controlados.

É possível reconhecer as emoções por diversas formas (Nasoz et al., 2004): (i) sensores capturando os sinais fisiológicos; (ii) análise de expressões faciais; (iii) análise da variação da fala por microfone; (iv) movimento corporal por meio da captura de dados por dispositivos padrões de entrada (i.e. mouse e teclado) e (v) análise do texto ao escrever uma opinião.

Este trabalho está limitado ao uso de expressões faciais para o reconhecimento de emoção devido às justificativas a seguir: (i) a popularidade de dispositivos que possuem câmeras fotográficas (e.g. smartphone, tablet, smart TV e notebook) facilitam a captura da expressão facial do usuário; (ii) a evolução das técnicas de classificação de imagens que estão alcançando a taxa de reconhecimento a nível humano e (iii) por não causar qualquer



Figura 2 – Expressão facial emocional

intrusão ao usuário, pois fotografias da expressão facial podem ser capturadas dentro do cotidiano das pessoas, não causa incômodo, não necessita o uso de instrumentos especiais, além de ser imperceptível aos usuários.

## 2.2 Expressão Facial Emocional

[Darwin \(1965\)](#) verificou que fenômenos emocionais idênticos, principalmente relacionados as expressões faciais, podiam ser encontrados em diferentes culturas. Posteriormente, o trabalho de [Ekman and Davidson \(1994\)](#) apontou a existência de um conjunto de expressões faciais universais que representam as mesmas emoções em diferentes culturas e estão exemplificadas na Figura 2. Essas expressões faciais universais pertencem ao grupo das emoções básicas, portanto é possível reconhecer as seguintes emoções por expressão facial: raiva, alegria, tristeza, desgosto, medo e surpresa. Cada emoção que é emitida por um indivíduo possui a sua própria movimentação muscular facial, sendo assim, há a caracterização de vários padrões que são chamados de unidades de ação, como o movimento da sobrancelha, dos olhos fechando, ao levantar as bochechas e entre outros ([Ekman and Friesen, 1977](#)). As unidades de ação são os padrões relevantes para diferenciar cada tipo de emoção.

## 2.3 Aprendizagem de Máquina

Com intuito de criar sistemas que possuem a capacidade de reconhecer emoção de forma automática é muito comum o emprego de técnicas de aprendizagem de máquina, que é um ramo da inteligência artificial em que máquinas aprendem a partir de uma experiência e são habilitadas para reconhecer padrões. Uma definição de aprendizagem de máquina foi dada por [Alpaydin \(2014\)](#): “*É a programação de computadores para otimizar um critério de desempenho usando dados de exemplo ou experiência passada*”. Na prática, isto pode ser entendido como a existência de um modelo definido com alguns parâmetros, no qual a ação de aprender consiste na execução de uma função de otimização, cujos parâmetros do modelo são otimizados, a partir dos dados de treinamento ou experiência passada. O modelo pode ser preditivo, para fazer previsões do futuro; descritivo, para obter conhecimento dos dados realizando a classificação; ou ambos.

Em aprendizagem de máquina, se as instâncias são conhecidas, isto é, cada instância possui o seu rótulo, então o aprendizado é supervisionado. Caso contrário, se as instâncias são desconhecidas, a aprendizagem é não supervisionada. Neste trabalho, o aprendizado utilizado é o supervisionado no qual os métodos possuem uma fase de treinamento e outra de teste. A primeira consiste na utilização de um conjunto de características com instâncias previamente rotuladas, também conhecido como base de treino, com objetivo de encontrar padrões nos exemplos e, assim, produzir um modelo para armazenar o aprendizado. Por fim, na fase de teste, o algoritmo deve classificar dados desconhecidos, por meio da base de teste, a partir dos padrões encontrados na fase anterior e mensurar o desempenho obtido. Além disso, o modelo gerado pode ser testado por outra base chamada de validação para averiguar com mais confiança que não há aprendizagem viciosa. Caso o desempenho esteja satisfatório o método está apto para a produção, senão deve voltar para fase de treinamento ([Géron, 2017](#); [Kotsiantis et al., 2007](#)).

## 2.4 Processo de Classificação de Imagem

Nos problemas de classificação de imagem, quando a abordagem é por meio da aprendizagem de máquina, geralmente é seguido o processo: (i) a etapa inicial consiste em uma fase de pré-processamento em que são aplicadas várias técnicas com a intenção de eliminar o ruído da imagem, resultando em sua melhora considerável para as fases posteriores; (ii) a etapa de extração de característica foca em destacar ou retirar as principais formas da imagem que são importantes para a separação das classes e (iii) as características extraídas são enviadas para um classificador determinar qual a classe que a imagem pertence.



Figura 3 – Detecção Facial

## 2.5 Técnicas para pré-processamento de imagens em reconhecimento de emoção por expressão facial

### 2.5.1 Detecção Facial

Este procedimento consiste na utilização de técnicas que verificam a existência de uma face em uma imagem, seja em uma fotografia ou *frame* de vídeo. Geralmente é um problema difícil, pois dado uma imagem o método deve detectar em qual região há faces. No entanto, uma imagem pode ter diferentes objetos, *backgrounds* e ruídos. A união desses elementos pode induzir o método a identificar erroneamente uma face, causando confusão e a ocorrência de falsos positivos, isto é, objetos sendo incorretamente identificados como face.

O algoritmo Viola-Jones ([Viola and Jones, 2001](#)) é amplamente usado pela comunidade para detecção facial. Esta técnica possui vantagens como alta taxa de precisão, rapidez na execução e baixa taxa de falsos positivos. Este algoritmo é utilizado por [Chen et al. \(2017b\)](#), [Shan et al. \(2017\)](#), [Shin et al. \(2016b\)](#), [Vo and Le \(2016\)](#), [Mayya et al. \(2016\)](#), [Ng et al. \(2015\)](#) e [Li et al. \(2015b\)](#) para detectar a face em uma imagem e realizar o recorte com o intuito de excluir o *background* da imagem, reduzindo assim a complexidade do problema e diminuindo a carga de aprendizado da rede na qual, neste caso, o classificador não necessita mais aprender a diferenciar o que é *background* e face.

### 2.5.2 Histograma de Equalização (Normalização do Brilho)

A técnica de histograma de equalização é utilizada para normalizar o brilho da imagem ([Jain, 1989](#)). Obviamente há diferentes tipos de ambiente e consequentemente a iluminação pode variar bastante. Esta técnica permite que a intensidade das cores seja melhor distribuída. A atuação desta técnica equilibra o contraste da imagem nas regiões

em que há ausência.

Esta técnica foi utilizada por [Shan et al. \(2017\)](#), [Kim et al. \(2016b\)](#) e [Shin et al. \(2016b\)](#). Para todos os casos resultou no aumento da taxa de acurácia comparada a não aplicação desta técnica. Provavelmente o histograma de equalização funciona principalmente porque busca equalizar as cores da imagem retirando os ruídos de iluminação, isto é, realçando todos os pontos da imagem ocasionando maior facilidade para aprendizagem do problema devido a maximização da visibilidade nas regiões importantes que separam as classes, além disso, há a diminuição da carga de aprendizado do classificador, pois é menos necessário aprender a separar as classes nos casos de iluminação adversa.

## 2.6 Rede Neural Artificial

O ser humano inspirou-se nos pássaros para construir aeronaves e voar. A natureza também inspirou outras invenções da humanidade, por exemplo a dianteira de um trem-bala. Da mesma forma, as redes neurais artificiais tiveram a mesma inspiração, especificamente no cérebro, tendo como objetivo construir máquinas inteligentes ([Géron, 2017](#); [Goodfellow et al., 2016](#)).

Um *perceptron*, que é uma simples arquitetura inspirada em um neurônio biológico, possui conexões de entrada e saída para conectar a outros neurônios. Cada conexão de entrada é associada a um peso, que recebe um sinal para o neurônio realizar uma computação baseada em uma função de ativação, gerando um sinal de saída que serve de entrada para outro neurônio ([Géron, 2017](#)).

A Figura 4 ilustra uma rede neural *perceptron* multicamadas. A rede neural *perceptron* possui uma estrutura em que consiste de uma camada de entrada, várias camadas intermediárias denominadas ocultas e uma camada de saída. Todas as camadas são compostas por neurônios *perceptron*. A camada de entrada recebe os dados oriundos de uma instância, por exemplo, os *pixels* de uma imagem, e encaminha os dados recebidos para a próxima camada. A camada oculta caracteriza-se por ser completamente conectada, isto é, cada neurônio conecta-se com todos da camada anterior e posterior. A camada de saída é responsável em fornecer o resultado da rede neural, por isso, nos casos de classificação, a quantidade de neurônio da camada de saída é a mesma das classes do problema. Além disso, cada neurônio da camada de saída está associada a uma classe e quando uma instância desconhecida for processada para classificação, o neurônio que terminar ativado da camada de saída representa a classificação desta instância ([Géron, 2017](#); [Goodfellow et al., 2016](#)).



Figura 4 – Rede Neural Artificial

## 2.7 Rede Neural de Convolução

As redes neurais de convolução (RNC) surgiram dos estudos do córtex visual do cérebro e têm sido usadas em reconhecimento de imagem desde 1980 (Géron, 2017). Nos últimos anos, uma série de fatores contribuíram para a evolução das RNCs, principalmente relacionados ao aumento do poder de computação (hardware), ao surgimento da web, que proporcionou o aumento da quantidade de dados para treinamento e à evolução das técnicas de treinamento de uma rede neural. Este cenário favorável permitiu que as RNCs alcançassem nível super-humano em alguns problemas complexos de visão computacional. As RNCs têm sido utilizadas em larga escala tanto pela indústria como pelos pesquisadores, sobretudo em problemas como máquinas de busca, carros autônomos, sistemas de classificação automática de vídeo e imagens, entre outras tarefas.

Os trabalhos de Hubel (1959) e Hubel and Wiesel (1959) realizaram uma série de experimentos em gatos em 1958, e posteriormente, em macacos (Hubel and Wiesel, 1968), para encontrar intuições do funcionamento do córtex visual, que é a parte cerebral responsável em processar informação visual. Estes trabalhos levaram os autores a receberem o Prêmio Nobel em Fisiologia e Medicina em 1981. Seus trabalhos mostraram que muitos neurônios do córtex visual tem um pequeno campo de recepção local, ou seja, os neurônios reagem somente a um estímulo localizado na região limitada pelo campo visual. O campo de recepção local dos diferentes neurônios podem sobrepor um ao outro e a sua combinação gera o campo visual. Os autores mostraram que alguns neurônios somente reagem às imagens com padrões de linhas horizontais enquanto outros reagem às linhas com diferentes orientações. Notou-se que alguns neurônios têm um campo grande de recepção local, conseqüentemente, reagindo aos padrões mais complexos. Entretanto, os neurônios estão combinando um ao outro para gerar padrões menos complexos. Estas observações são evidências de que os neurônios são baseados na saída do vizinho.

O poderoso funcionamento do córtex visual está habilitado a detectar todos os padrões complexos em qualquer área do campo visual (Géron, 2017). Todos os estudos relacionados ao córtex visual foram gradualmente inseridos nas redes neurais artificiais para gerar a rede neural de convolução. A primeira RNC foi apresentada por LeCun et al. (1998), uma arquitetura denominada LeNet-5 que foi utilizada para reconhecer dígitos escritos no papel.

### 2.7.1 Camada de Convolução

A camada de convolução é o mais importante bloco de uma RNC e consiste em uma operação matemática que desliza uma função sobre a outra calculando a integral entre a multiplicação de duas funções. A Figura 5 ilustra o que cada camada de convolução recebe em seu campo visual dada uma imagem como entrada na RNC. Vale ressaltar que cada neurônio da primeira camada de convolução está conectado somente a alguns campos visuais de recepção da imagem, diferentemente da abordagem tradicional que é conectada a todos os *pixels*. Os neurônios da segunda camada de convolução estão conectados somente aos localizados no pequeno campo visual (retângulo) da primeira camada, novamente diferente da rede neural *perceptron* em que todos os neurônios são conectados com todos da camada anterior, e assim por diante (Géron, 2017). As vantagens da RNC sobre a abordagem tradicional são explicadas pela diferença entre ambas e são enfatizadas a seguir:

- (i) A RNC tem característica esparsa como ilustrado na Figura 6, por isso, a RNC possui menor quantidade de parâmetros para serem treinados do que as redes neurais tradicionais, isto é, requer menos tempo de treinamento e recursos computacionais (Goodfellow et al., 2016);
- (ii) A estrutura da RNC é comum no mundo real, por exemplo no córtex visual. Essa inspiração biológica é uma das razões para RNC funcionar tão bem no reconhecimento de imagens (Géron, 2017);
- (iii) E por fim, o compartilhamento de parâmetros resulta no aprendizado de rotações dos objetos da imagem e os neurônios aprendem em conjunto ao invés de separados (Goodfellow et al., 2016).

### 2.7.2 Camada de *Pooling*

Esta camada tem como foco principal realizar subamostra, isto é, diminuir o tamanho de entrada da imagem entre as camadas de convolução com objetivo de reduzir a carga computacional e o número de parâmetros, ocasionando a diminuição do risco de





Figura 5 – Camada de Convolução com campos locais de recepção



Figura 6 – Conectividade esparsa. É destacada a entrada  $x_3$  e a saída em S que são afetadas por  $x_3$ . (Cima) Quando S recebe a convolução com um *kernel* de tamanho 3, somente três saídas são afetadas por  $x_3$ . (Baixo) Quando S é gerado por rede neural tradicional, todos são afetados por  $x_3$ .



Figura 7 – Camada de *Max Pooling*

*overfitting* e o consumo de memória (Géron, 2017). Além disso, reduzir o tamanho de entrada da imagem faz a rede neural mais tolerável a variação do objeto principal, como a rotação da face.

Geralmente uma camada de *pooling* é implementada logo após uma camada de convolução, portanto recebe como entrada uma imagem processada anteriormente por uma camada de convolução com intuito de realizar subamostra da imagem e encaminhar o resultado para uma próxima camada de convolução (Goodfellow et al., 2016).

Existem duas principais funções de *pooling*. Por exemplo, o *max pooling* que considera o valor máximo de um campo de recepção e está exemplificado na Figura 7 por um *kernel* 2x2 que elimina 75% dos valores de entrada. Há também o *pooling* pela média de um campo de recepção e seu cálculo consiste na distância entre o *pixel* central e seus vizinhos (Goodfellow et al., 2016). A camada de *max pooling* é a mais comum operação de *pooling* utilizada em redes neurais de convolução.

### 2.7.3 Regressão *Softmax*

Geralmente uma arquitetura de RNC é composta em sua maior parte por camadas de convolução e *pooling*. Estas camadas processam a imagem com intuito de extrair os principais padrões para um classificador e determinar a classe dela. Por isso, ao final de uma arquitetura de RNC é necessário um classificador que tradicionalmente tem sido o *softmax*. Este classificador é um modelo generalizado de uma regressão logística, que por sua vez é normalmente usada para estimar probabilidades de uma instância pertencer a uma classe em particular, por exemplo, qual é a probabilidade de um email ser spam? (Géron, 2017). Se a estimativa de probabilidade for maior que 50%, então o modelo prevê que a instância pertence a classe positiva, caso contrário, pertence à classe negativa. Portanto, isto faz do regressor logístico um classificador binário. Um *softmax* é capaz de estimar probabilidades para múltiplas classes e não há a necessidade de treinar e combinar

Tabela 1 – Arquitetura AlexNet

Camada	Tipo	Mapas	Tamanho	Kernel	Ativação
Saída	Completamente Conectada	-	1000	-	Softmax
F9	Completamente Conectada	-	4096	-	ReLU
F8	Completamente Conectada	-	4096	-	ReLU
C7	Convolução	256	13 x 13	3 x 3	ReLU
C6	Convolução	384	13 x 13	3 x 3	ReLU
C5	Convolução	384	13 x 13	3 x 3	ReLU
S4	<i>Max Pooling</i>	256	13 x 13	3 x 3	-
C3	Convolução	256	27 x 27	5 x 5	ReLU
S2	<i>Max Pooling</i>	96	27 x 27	3 x 3	-
C1	Convolução	96	55 x 55	11 x 11	ReLU
Entrada	Entrada	3 (RGB)	224 x 224	-	-

múltiplos classificadores binários para tal tarefa. Um *softmax* está habilitado a estimar probabilidades para uma imagem processada em uma RNC e, para todos os efeitos, a imagem é uma expressão facial que pertence às emoções básicas: neutralidade, felicidade, surpresa, medo, raiva, tristeza ou desgosto.

## 2.8 Arquiteturas de Redes Neurais de Convolução

### 2.8.1 AlexNet

A arquitetura AlexNet foi desenvolvida por Alex Krizhevsky (por isso, o nome da mesma), Ilya Sutskever e Geoffrey Hinton. Destacando-se por ser grande e muito profunda, a AlexNet foi a primeira RNC a empilhar camadas de convolução diretamente em cima da outra, ao invés da tradicional conexão entre uma camada de convolução e a camada de *pooling*. Esta arquitetura pode ser consultada na Tabela 1.

A AlexNet usa uma normalização bastante eficiente entre as camadas C1 e C3 denominada *local response normalization*. Essa forma de normalização causa um efeito que contribui para inibição dos neurônios em ativar mais fortemente no mesmo local, ocasionando que outros mapas de características se tornem também especialistas em determinada região da imagem. Este comportamento também é observado nos neurônios biológicos (Géron, 2017).

O seu grande sucesso foi devido ao desafio de 2012 do ImageNet ILSVRC, que tem sido o principal concurso de classificação de imagens organizado pela comunidade científica, em que venceu por uma margem bastante grande: alcançou 17% no top-5 da taxa de erro, enquanto o segundo melhor alcançou somente 26%!

Figura 8 – Módulo *inception*

### 2.8.2 GoogLeNet

A arquitetura GoogLeNet (Szegedy et al., 2015) foi desenvolvida por Christian Szegedy do *Google Research*, e venceu o desafio do ILSVRC 2014 por alcançar a taxa de erro no top-5 abaixo de 7%. Este grande desempenho foi devido em grande parte pelo fato de que a rede foi muito mais profunda do que as anteriores. O aumento de profundidade está diretamente relacionado à criação de sub-redes chamadas de *inception* e está ilustrada na Figura 8. Essas sub-redes permitiram que a GoogLeNet utilizasse os parâmetros de forma mais eficiente comparada às arquiteturas anteriores e a dimensão desta eficiência pode ser elucidada pela diferença de parâmetros entre a GoogLeNet e a AlexNet, em que a primeira possui 10 vezes menos parâmetros do que a segunda (Géron, 2017).

Um módulo *inception*, que está ilustrado na Figura 8, possui uma notação " $3 \times 3 + 2$  (S)". Isto significa que a camada usa um *kernel*  $3 \times 3$ , *stride* 2 e *SAME padding*. O sinal de entrada é primeiramente copiado e alimenta as camadas do módulo *inception* que estão divididas em dois conjuntos. O primeiro conjunto recebe o sinal para processar e encaminhar para o segundo. Vale ressaltar que todas as camadas de convolução utilizam a ReLU como função de ativação. É interessante observar que o segundo conjunto usa diferentes tamanhos de *kernel* ( $1 \times 1$ ,  $3 \times 3$  e  $5 \times 5$ ), permitindo que a rede possa capturar diferentes padrões de escala (Géron, 2017). No fim do módulo, há uma camada de concatenação, isto é, combina todas as saídas do segundo conjunto de camadas de convolução e encaminha um único sinal que é o resultado do módulo para uma próxima camada da rede. A rede GoogLeNet é composta por módulos *inception*, camadas de convolução, *max pooling*, camadas de *local normalization response*, camadas completamente conectadas e *softmax*.

Tabela 2 – Arquiteturas VGGNet

VGGNet configuração					
A	A-LRN	B	C	D	E
11 camadas	11 camadas	13 camadas	16 camadas	16 camadas	19 camadas
camada de entrada (224 x 224 imagem RGB)					
conv3-64	conv3-64 <b>LRN</b>	conv3-64 <b>conv3-64</b>	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 <b>conv3-128</b>	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 <b>conv1-256</b>	conv3-256 conv3-256 <b>conv3-256</b>	conv3-256 conv3-256 conv3-256 <b>conv3-256</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
camada completamente conectada - 4096					
camada completamente conectada - 4096					
camada completamente conectada - 1000					
softmax					

### 2.8.3 VGGNet

Proposta por [Simonyan and Zisserman \(2014\)](#), a VGGNet foi a vice-campeã do desafio ILSVRC 2014, tendo alcançado 6.8% na taxa de erro top-5. A sua principal contribuição foi uma avaliação exaustiva de seis RNCs, que consistiu no aumento de profundidade enfatizando a utilização de filtros de convolução com tamanho muito pequeno (3 x 3), promovendo o aumento da profundidade da rede que passou de 16 para 19 camadas. Esta abordagem mostrou um aumento de eficiência significativo comparado às técnicas anteriores.

A VGGNet é composta principalmente por camadas de convolução, *max pooling*, completamente conectadas e *softmax*, e está ilustrada na Tabela 2. O fato de usar filtros de convolução com tamanho pequeno ocasionou no não aumento do número de parâmetros a serem ajustados a medida que a rede cresce, promovendo assim a eficiência.

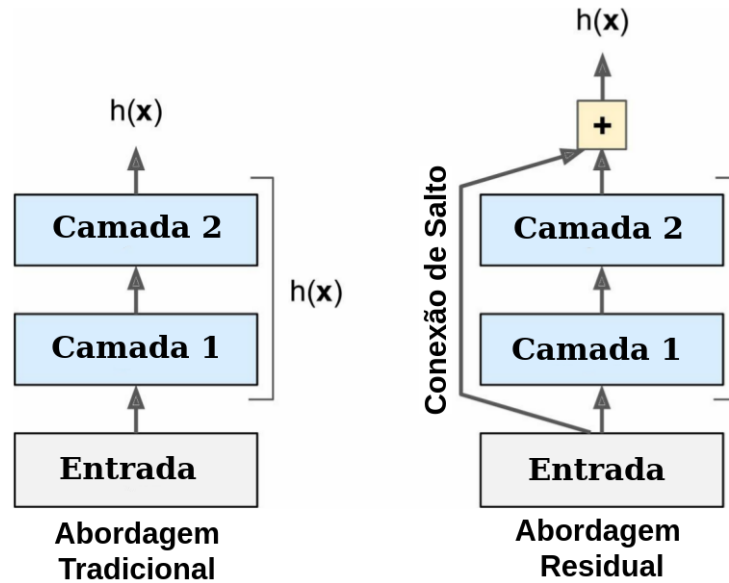


Figura 9 – Bloco residual

#### 2.8.4 Residual Network

A Residual Network (ou ResNet) foi desenvolvida por [He et al. \(2016\)](#) e foi a técnica vencedora do desafio ILSVRC 2015. A ResNet foi avaliada na base de teste do ImageNet que tinha 1000 classes e alcançou 3.6% em taxa de erro no top-5. Este tipo de avaliação consiste na verificação das 5 classes com maiores probabilidades após uma classificação, e caso esta imagem de fato pertencer ao grupo das 5 classes é caracterizado uma classificação correta. A ResNet proposta para o concurso tinha 152 camadas caracterizando uma rede muito profunda.

O segredo desta técnica consiste na composição dos blocos residuais. A novidade desta arquitetura é possibilitar que um bloco residual faça o mapeamento de outros blocos distantes através de conexões de salto, isto é, os sinais emitidos pelos neurônios tanto na direção para frente como para trás podem ser propagados diretamente de um bloco para qualquer outro bloco como ilustrado na Figura 9. Este tipo de conexão trouxe impactos positivos na aprendizagem, principalmente atuando na função objetiva, pois quando a mesma está próxima de convergir é aumentado consideravelmente a velocidade do treino comparada as redes tradicionais sem as conexões de salto.

#### 2.8.5 Ensemble

Suponha uma questão complexa perguntada aleatoriamente para milhões de pessoas, então, todas as respostas são agregadas para obter o resultado. Em muitos casos, a resposta agregada é melhor do que a de um especialista. Isto é conhecido como sabedoria popular. Similarmente, se agregar as classificações de uma determinada instância oriunda de um grupo de RNCs e, frequentemente, tem mais acertos do que a classificação de uma

única RNC. Um grupo de RNCs ou outros tipos de classificadores e regressores são chamados de *ensemble* (Géron, 2017). Geralmente, quem trabalha com *ensemble* no problema de reconhecimento de emoção tem calculada a média das probabilidades estimadas pelas diversas RNCs para decidir qual a classificação final da expressão facial.

## 2.9 Métricas de Avaliação de Desempenho para Classificadores

As principais métricas utilizadas para avaliação de desempenho de classificadores são: acurácia, precisão, revocação e f1-score. Por meio dessas métricas é possível comparar os classificadores determinando os pontos fortes e fracos, investigar quais são as classes que o modelo tem maiores desempenho, verificar a existência de *overfitting* em uma classe específica, e também, se houve aprendizado em classes com menores amostras. É relevante tais verificações, pois quando a base de treino é desbalanceada pode gerar um modelo que consiste de chutar as classes com maiores amostras a fim de maximizar as taxas de acertos, no entanto, claramente caracterizando um modelo com *overfitting*.

A acurácia é a proporção dos casos corretamente classificados e é calculada a partir da Equação 2.1, em que TP é a taxa de verdadeiros positivos, TN é a taxa de verdadeiros negativos, FP é a taxa de falsos positivos e FN é a taxa de falsos negativos. A precisão nos diz a fração de instâncias que são verdadeiras positivas de um grupo que o classificador preveu ser positivo e a fórmula pode ser consultada em 2.2. A revocação mede a fração de instâncias positivas que o classificador esqueceu de classificar corretamente (Equação 2.3).

Classificadores com uma taxa de revocação alta não tem muitas instâncias positivas classificadas incorretamente. Vale destacar que facilmente é possível construir um classificador que alcança altas taxas de precisão ou revocação mas não ambos. Caso o classificador prevê todas as instâncias da classe positiva, o modelo teria uma taxa perfeita de revocação, entretanto uma baixa precisão. Por isso, criar um classificador que maximiza tanto precisão como revocação é um desafio.

É frequentemente conveniente combinar precisão e revocação em uma única métrica chamada de f1-score, principalmente quando é necessário comparar de uma simples maneira dois classificadores. A f1-score é a média harmônica de precisão e revocação (Equação 2.4). Enquanto a média tradicional trata todos os valores igualmente, isto é, sem pesos, a média harmônica favorece os valores mais baixos. Portanto, caso tenha alta revocação e baixa precisão, o resultado de f1-score será mais próximo de precisão, pois esta métrica adiciona pesos nos valores mais baixos, que neste caso é o valor de precisão.

$$acurácia = \frac{TP + TN}{TP + TN + FP + FN} \quad (2.1)$$

$$precisão = \frac{TP}{TP + FP} \quad (2.2)$$

$$revocação = \frac{TP}{TP + FN} \quad (2.3)$$

$$f1 - score = \frac{2}{\frac{1}{precisão} + \frac{1}{revocação}} = 2 * \frac{precisão * revocação}{precisão + revocação} \quad (2.4)$$

## 2.10 Resumo

Neste capítulo foram apresentados os principais conceitos utilizados por esta proposta. Foi visto que uma RNC, que é uma técnica de aprendizagem de máquina, foi projetada para a classificação de imagem, justificando a sua escolha para a classificação de emoções em expressão facial. Esta técnica é muito poderosa e tem sido inspirada no cérebro dos mamíferos. A expressão facial tem se destacado como uma forma eficaz de coleta de dados para o reconhecimento de emoção, principalmente, por ser ubíqua, não intrusiva e pela popularidade de dispositivos que contém câmeras fotográficas. Entretanto, somente as emoções básicas são possíveis de reconhecer por meio da expressão facial. Isso significa que reconhecer tédio, frustração e confusão se torna bastante difícil por esse meio. Além disso, técnicas de pré-processamento, tais como histograma de equalização e detecção com recorte da face são úteis durante o processo de classificação de imagens, justamente por diminuir a carga de aprendizado da rede. Por fim, foram conceituadas algumas das principais arquiteturas de RNCs em que cada uma possui sua característica em particular que precisam ser experimentadas em diferentes cenários.





## 3 Trabalhos Correlatos

Os trabalhos relacionados são discutidos neste capítulo e está organizado da seguinte forma. Na Seção 3.2 analisa as principais extrações de características encontradas na literatura para o processamento da expressão facial. Na Seção 3.3 avalia os trabalhos correlatos categorizados pelas arquiteturas de RNC. Na Seção 3.4 é comentado as aplicações para o reconhecimento de emoção por expressão facial e enquanto a Seção 3.5 faz um resumo a respeito deste capítulo.

### 3.1 Preparação dos dados

Os problemas de classificação em geral, seja de imagem, vídeo, áudio ou qualquer tipo, tradicionalmente sofrem pela ausência de dados. Algoritmos de aprendizagem de máquina requerem quantidade de dados expressivos para apresentar soluções com desempenho satisfatório, especificamente as redes neurais profundas. Raramente há dados disponíveis e que sejam suficientes para treinar e validar uma rede neural de convolução, vale ressaltar que cada problema tem sua particularidade, isto é, quanto maior a complexidade mais dados são necessários.

Contudo, a comunidade de reconhecimento de emoção para amenizar esse problema utiliza a técnica de aumento de dados e multiplicação de imagens. Essa técnica consiste na geração de cópias de uma imagem original, que contém uma expressão facial, para gerar imagens duplicadas. Entretanto, tais imagens duplicadas são diferentes da imagem original, justamente por possuir alterações na posição da face com leves rotações da mesma, variação da intensidade das cores e redimensionamento com aplicação de *zoom*. As imagens aumentadas são usadas durante o treinamento contribuindo para a rede neural aprender a reconhecer emoção em diferentes rotações, intensidade de iluminação e escala.

Os trabalhos de Barsoum et al. (2016b); Huang and Lu (2016); Kim et al. (2016b); Shin et al. (2016b); Yu et al. (2016c) e Li et al. (2015b) utilizaram a técnica de aumento de dados. A técnica foi configurada para aumentar entre 5 a 10 vezes cada imagem original. Sendo assim, a base de dados original foi ampliada em até 10 vezes, gerando um ganho considerável dos dados. Tais trabalhos alcançaram boas taxas de reconhecimento em que o *fine tuning* dos modelos possuem generalização adequada e não apresentando *overfitting* e *underfitting*. O resultado expressivo foi viabilizado pela técnica de aumento de dados justamente pela rede ser treinada e validada com maiores quantidades de dados.

## 3.2 Extração de Característica

Uma etapa essencial durante o processo de classificação de imagem é a extração de característica. A extração de característica é sucintamente enfatizada na Seção 2.4 e tem como finalidade destacar ou retirar as formas mais relevantes da imagem que são cruciais para a separação das classes. A seguir, os principais tipos de extração de características empregados para o reconhecimento de emoção por expressão facial são analisados.

### 3.2.1 Extração Geométrica

A extração de características geométrica consiste na obtenção de pontos faciais ilustradas pela Figura 10. As características geométricas tem como finalidade capturar as deformações na face causadas pela ativação dos músculos a partir dos pontos faciais (Yu et al., 2016c). Esses pontos faciais podem ser mapeados pelos seguintes métodos: Yu et al. (2016a) e Yu et al. (2014). A extração geométrica é uma abordagem que realiza medições entre diversas partes da face tais como:

- (i) Altura da sobrancelha esquerda/direita (distância vertical entre o ponto mais superior da sobrancelha e centro do olho);
- (ii) Altura da pálpebra esquerda/direita (distância vertical entre o ponto mais superior do olho e parte inferior do olho);
- (iii) Altura do nariz (distância vertical entre o ponto mais inferior do olho para o nariz e centro de ambos os olhos);
- (iv) Largura do nariz (distância horizontal entre os pontos do nariz mais à esquerda e à direita);
- (v) Altura do lábio superior (distância vertical entre o ponto mais superior e o centro da boca);
- (vi) Altura do lábio inferior (distância vertical entre o ponto mais inferior e o centro da boca);
- (vii) A distância do ponto da boca mais a esquerda para o centro da boca;
- (viii) E por fim, a distância do ponto da boca mais a direita para o centro da boca.

A extração geométrica é amplamente empregada nas abordagens tradicionais de aprendizado de máquina, isto é, abordagens que não utilizam as redes neurais de convolução. Todavia, o trabalho de Yu et al. (2016c) realiza a extração geométrica concatenando com a RNC e obteve um pequeno ganho na taxa de precisão de 1%. Um diagrama da

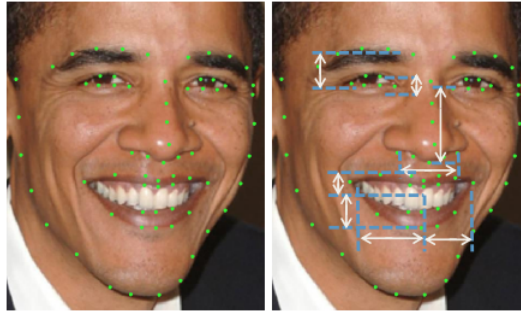


Figura 10 – Extração dos pontos faciais para características geométrica

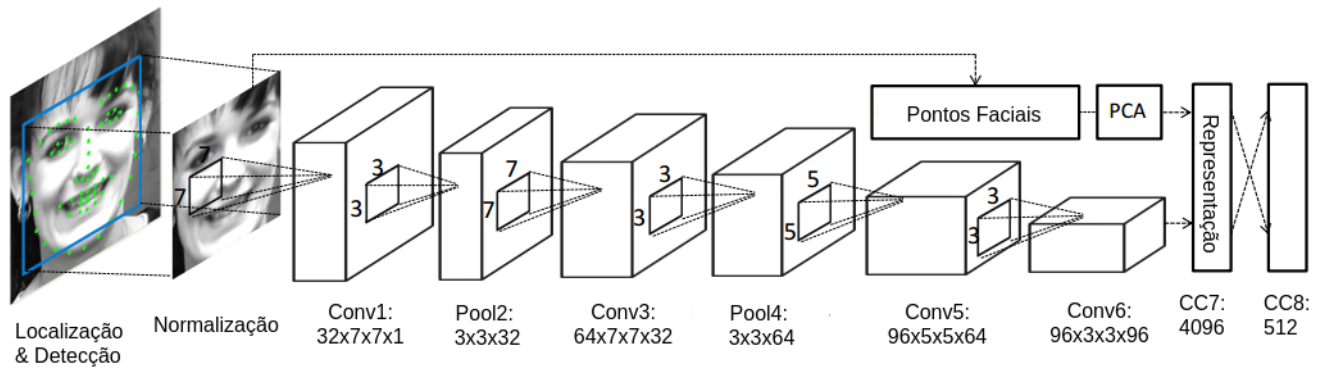


Figura 11 – Concatenação dos pontos faciais com uma rede neural de convolução

sua abordagem é ilustrada na Figura 11. No entanto, a combinação entre a extração geométrica e RNC, obviamente, aumenta o custo computacional devido a outros algoritmos serem executados como o mapeamento dos pontos faciais e as suas distâncias. Caso o foco do reconhecedor de emoções for aplicações em cenários reais, provavelmente, não é viável a concatenação devido o aumento na taxa de reconhecimento ser baixo, portanto não compensada pelo aumento do custo computacional, pois tais cenários requerem classificação instantânea e em tempo real.

### 3.2.2 Extração Aparente

A extração de característica aparente considera as sub-regiões faciais, principalmente próximas da boca e dos olhos, como características essenciais para a classificação, diferentemente, das características geométricas que foca na captura das deformações dos pontos faciais e possui como desvantagem não considerar as mudanças aparentes causadas por essas deformações capturadas (Yu et al., 2016c). A extração aparente foi amplamente estudada por Ekman and Davidson (1994), encontrando 96 unidades de ações (ou sub-regiões) na face correspondentes a movimentação de diversos músculos relacionada a uma determinada emoção. A RNC está habilitada naturalmente a realizar a extração de característica aparente, sendo assim, o processo de aprendizado está associado a descoberta de

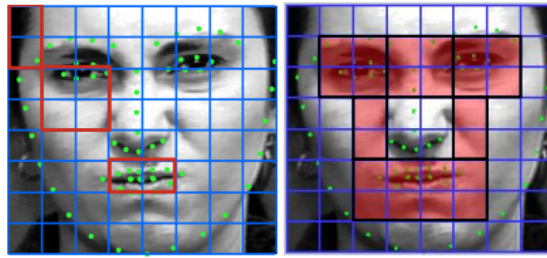


Figura 12 – Extração das sub-regiões faciais para características aparente

quais as sub-regiões da face mais relevantes para determinar qual a emoção está emitida em uma face. A Figura 12 ilustra as sub-regiões de uma face que são interessantes para a classificação depois de um processo de aprendizado.

### 3.3 Arquiteturas

#### 3.3.1 AlexNet

A arquitetura AlexNet foi a RNC mais encontrada na literatura para reconhecimento de emoção, principalmente por ter sido a pioneira da família de métodos conhecido como aprendizado profundo que fez grande sucesso, inclusive esta rede venceu o desafio ILSVRC-2012. O trabalho de Kim et al. (2016b) utiliza AlexNet para reconhecer emoções, a sua abordagem é destacada por combinar a AlexNet com uma rede profunda *autoencoder* para alinhamento de face, esta foi uma grande contribuição do autor que apresentou uma solução funcional para o problema de rotação ou desalinhamento da face. O trabalho de Shan et al. (2017) propõe uma abordagem que consiste no emprego da técnica de equalização de histograma durante a fase de pré-processamento com intuito de resolver o problema da iluminação, desta forma, foi mostrado que a aplicação da técnica melhorava o aprendizado da rede. Nos trabalhos que utilizaram a AlexNet é notável que para alcançar maiores taxas de reconhecimento foi necessário o emprego de técnicas de pré-processamento.

#### 3.3.2 VGG

Barsoum et al. (2016b) utilizou uma VGG de 13 camadas, isto é, uma arquitetura bastante profunda inclusive com camadas de *dropout* para alcançar maior generalização. Utilizou a técnica de aumento de dados, isto é, gerou novas imagens multiplicando por 10 a imagem original, as imagens geradas tinham variações da pose da face, pequenas rotações que contribuem para o aprendizado da rede. No trabalho de Ng et al. (2015), abordou a transferência de aprendizado para melhorar a performance de classificação, a VGG utilizada anteriormente tinha sido pré-treinada com a base do *ImageNet* (do tradicional

desafio ILSVRC), que é uma base que contém dezenas de objetos, e então, foi realizado um segundo auto-ajuste (treinamento) para a classificação das expressões faciais.

### 3.3.3 GoogLeNet

Guo et al. (2016) realizou um trabalho comparando a arquitetura *GoogLeNet* com a *AlexNet*. A primeira alcançou melhores resultados, principalmente por possuir camadas *inceptions* que são arquiteturas mais otimizadas para a classificação de imagens. Neste trabalho, também foi testado o classificador *kNN* na última camada, que alcançou melhor resultado que o tradicional *softmax*. Entretanto, vale ressaltar que, o *kNN* é uma técnica baseada em instâncias, isto é, não aprende e consequentemente não gera um modelo. Em outras palavras, caso haja  $100k$  imagens na base de treino, o *kNN* para classificar qualquer instância aleatória da base de teste ou validação necessita realizar um conjunto de cálculos de distâncias entre as  $100k$  imagens da base de treino para encontrar a classificação da instância, implicando em uma grande desvantagem pois este processo é repetido para cada instância que se deseja classificar.

### 3.3.4 Ensemble

Os trabalhos que utilizaram *Ensemble* combinaram CNN com outras técnicas ou com outras CNN. No trabalho realizado por Wen et al. (2017b) foram combinadas até 100 CNNs, no qual obteve melhor resultado que uma CNN sozinha. Entretanto, foi verificado que poucas CNNs, isto é, menos de 10, já alcançam o mesmo resultado que 100 CNNs juntas, ou seja, o aprendizado fica estável se continuar aumentando a quantidade de CNN a partir de um limiar. Isto implica que menos de 10 CNNs combinadas, é o suficiente para aprender o problema de reconhecimento de emoção por expressão facial. Liu et al. (2016b) implementou um *ensemble* de 3 CNNs com um único classificador *softmax* que recebia a extração de características das 3 CNNs. No trabalho de Shin et al. (2016b) foi treinada 20 redes diferentes com 5 entradas diferentes, também, foi testado a utilização do classificador Support Vector Machine (SVM) ao invés do tradicional *softmax* na última camada, e o Support Vector Machine alcançou resultados melhores.

## 3.4 Aplicações

Há diversas aplicações para o reconhecimento de emoção no mundo real, foi percebido que os pesquisadores de reconhecimento de emoção por expressão facial utilizando RNC, ultimamente concentraram seus esforços mais no desenvolvimento de reconhecedores de emoção do que a aplicação em cenários reais, mesmo assim, está aberto para trabalhos futuros inúmeras aplicações desses reconhecedores em diversas áreas, tendo destaque principalmente para:

Tabela 3 – Principais arquiteturas de Redes Neurais de Convolução para reconhecimento de expressões faciais. (\*) Significa que a rede foi treinada (*fine-tuning*) por duas vezes.

Arquitetura	Trabalho	Base de Treino	Base de Validação	Acurácia
AlexNet	Chen et al. (2017b)	CK+	CK+	99.1%
		CK+	JAFfE	83.11%
		JAFfE	JAFfE	87.7%
	Shan et al. (2017)	JAFfE	JAFfE	76.7%
		CK+	CK+	80.3%
	Kim et al. (2016b)	FER	FER	73.73%
	Huang and Lu (2016)	FER	FER	76.9%
		CK+	CK+	97.3%
	Vo and Le (2016)	CK+	CK+	96.04%
	Yu et al. (2016c)	CK+	CK+	98.7%
		MMI	MMI	98.6%
	Ng et al. (2015)*	FER/EmotiW	FER/EmotiW	55.6%
	Jung et al. (2015b)	CK+/FER	CK+/FER	86.54%
VGG	Barsoum et al. (2016b)	CIFE	CIFE	81.5%
		CK+	CK+	83%
VGG	Barsoum et al. (2016b)	FER+	FER+	84.9%
	Ng et al. (2015)*	FER/EmotiW	FER/EmotiW	52.1%
GoogLeNet	Guo et al. (2016)	FER/SFEW2.0	FER/SFEW2.0	71.3%
Ensemble	Wen et al. (2017b)	FER	FER-Private	69.96%
		FER	CK+	76.05%
		FER	JAFfE	50.70%
		FER	EmotiW	34.09%
	Liu et al. (2016b)	FER	FER	65.03%
	Shin et al. (2016b)	FER/SFEW	FER-Test	66.67%
			SFEW	64.84%
			CK+	65.54%
			KDEF	50.66%
			JAFfE	49.17%

- Interação humano computador (Barsoum et al., 2016b; Chen et al., 2017b; Liu et al., 2016b; Wen et al., 2017b), onde pode ser possível projetar interfaces que se adaptam ao estado emocional do usuário;
- Psiquiatria e cuidados médicos (Chen et al., 2017b; Mayya et al., 2016; Wen et al., 2017b), no qual o reconhecedor de emoção deve monitorar constantemente o paciente ou usuário fornecendo dados emocionais que podem contribuir para diagnósticos;
- Deficiente visual (Li et al., 2015b), pois pessoas com alto grau de deficiência visual, tem dificuldades na interação entre pessoas para identificar qual a emoção que as pessoas em volta estão emitindo;
- Interação humano robô (Jung et al., 2015b; Shin et al., 2016b), fazendo com que robôs estejam habilitados a interagir com humanos podendo adaptar-se a emoção

dos humanos em volta, ou até mesmo emitir emoção se aproximando de um humanoide;

- Personagens virtuais e animação (Vo and Le, 2016; Yu et al., 2016c), habilitando avatares a copiar expressão humana que podem ser útil para gravações de filmes de animação, também pode ser usado em aplicações de animação como o popular aplicativo para *smartphone* o *Snapchat*, que identifica a expressão facial do usuário e retorna alguma animação sobrepondo a expressão anteriormente detectada do usuário.

## 3.5 Resumo

Neste capítulo foram apresentado os trabalhos relacionados, verificamos que o tema está em crescente investigação pela comunidade, visto que arquiteturas poderosas de RNC surgiram recentemente. Apesar de uma RNC ter embutido o pré-processamento em sua própria arquitetura, para o problema tratado por este trabalho, verificamos que pré-processamento adicionais (não originais da RNC) melhoraram consideravelmente a taxa de acurácia. É importante frisar que as principais arquiteturas empregadas tem sido as vencedoras ou com pouca variação do tradicional problema de aprendizado profundo: o desafio ILSVRC. Estas arquiteturas tem se saído bem no problema de reconhecimento de emoção alcançando resultados comparado a nível humano.

Um conjunto de aplicações foram identificadas para o reconhecimento de emoção por expressão facial, entretanto, investigando a literatura há um índice baixo de adesão em cenários de uso reais. Atualmente, os pesquisadores apenas tem falado que é possível usar em uma determinada área, mas de fato não o experimentaram na prática. Contudo, modelos de reconhecimento de emoção tem alcançado taxas de acurácia confiáveis para serem empregados na indústria e pesquisa.





## 4 Abordagem Proposta

Neste capítulo é descrita a abordagem proposta para reconhecer emoções por meio da expressão facial e está dividido da seguinte forma. A Seção 4.1 detalha o monitoramento do indivíduo. A Seção 4.2 descreve um módulo de detecção de face e recorte. A Seção 4.3 aborda as operações de pré-processamento aplicadas na imagem. Na Seção 4.4 é apresentada a atribuição da rede neural de convolução e por fim um resumo do capítulo na Seção 4.5.

### 4.1 Monitoramento do Indivíduo

Um dispositivo que possui uma câmera fotográfica periodicamente está fotografando o indivíduo. Por exemplo, em um ambiente educacional um *tablet* ou *notebook*, que executa uma plataforma educacional, pode estar monitorando o estudante por meio da câmera frontal ou *webcam*. Este cenário configura um ambiente favorecedor para o monitoramento, pois a câmera frontal ou *webcam* sempre estão capturando as reações do estudante com ângulo frontal. Em um outro exemplo, destacamos uma aplicação para deficientes visuais onde uma câmera apropriada para *wearables* pode estar anexada a roupa do usuário na região peitoral, com intuito de monitorar as reações das pessoas ao redor. Entretanto, este cenário, ao contrário do anterior, possui maiores desafios do ponto de vista de coleta de dados, pois o usuário que está com a câmera pode movimentar-se capturando imagens tremidas, desfocadas e borradas, além de faces com ângulos variáveis. Tais problemas, entretanto, podem ser tratados por equipamentos de qualidade e, portanto, estão fora do escopo deste trabalho. A Figura 16 ilustra o fluxo da solução proposta. Portanto, enquanto o monitoramento está ocorrendo, as imagens são enviadas para um repositório de entrada de dados.

### 4.2 Detecção de Face e Recorte

Este procedimento consiste na detecção de todas as faces de uma imagem por meio do algoritmo Viola Jones (veja Seção 2.5.1). Primeiramente, uma imagem é consumida do repositório de entrada de dados. A detecção consiste na geração de um conjunto de coordenadas que possibilita o desenho de um retângulo indicando a localização da face. Vale ressaltar que esta atividade possui complexidade moderada, pois uma imagem oriunda de cenários reais contém vários objetos com diferentes geometrias, inclusive podendo assemelhar-se a uma face ocasionando a geração de falsos positivos. Logo após a detecção de face, é realizada o recorte da mesma indicada pelo conjunto de coordenadas

definidas pela etapa anterior. Tal atividade é valiosa para exclusão do *background*. Desta forma, somente a face recortada é enviada para a fase seguinte, reduzindo a complexidade do problema, pois não há necessidade do classificador aprender a separar o *background* da face. Posteriormente ao recorte da face, a imagem original, que deve estar com uma face recortada, é mantida para nova averiguação de recorte de face. Caso exista outras faces na imagem, este processo é repetido até não existir mais faces para recortar. Obviamente caso seja enviada uma imagem para a etapa de detecção e recorte que não contém uma face (e.g. imagem de um avião) o processo é automaticamente encerrado, pois se não há uma face para detectar, logo não há uma expressão facial emocional para reconhecer.

### 4.3 Pré-Processamento

Uma face recortada é recebida pelo módulo de Pré-processamento. Nesta etapa, operações de pré-processamento são aplicadas com a finalidade de enaltecer as características próprias de cada expressão facial, com intuito de preparar a imagem para a classificação, facilitando a diferenciação das emoções. A Figura 14 mostra o fluxo desta etapa. Adicionalmente, uma função de redimensionamento é chamada para transformar a imagem em uma escala de 60x60 *pixels*. Neste trabalho, consideramos que a imagem é colorida, portanto a mesma possui 3 canais denominados RGB (do inglês: Red, Green e Blue), sendo que a imagem resultante possui 10.800 características que pode ser calculada por  $Qtd\_Caracteristica = N\_Pixels\_X * N\_Pixels\_Y * N\_Canais$ .

A meta principal desta proposta consiste em classificar emoções em qualquer ambiente. Obviamente que a variação do ambiente acarreta em diferentes níveis de intensidade da luz, ocorrendo a perda de características importantes da face que diferenciam as emoções, seja por excesso ou ausência de luz. Vale destacar que esta proposta é baseada principalmente em redes neurais de convolução que originalmente possui vários filtros de pré-processamento. Entretanto, a literatura tem mostrado que filtros clássicos aplicados antes da inserção de uma imagem em redes neurais tem sido eficazes na eliminação de ruídos, principalmente aqueles relacionados a iluminação e brilho (Kim et al., 2016b; Shan et al., 2017; Shin et al., 2016b). Portanto, as técnicas de normalização de brilho e iluminação são parte da etapa de pré-processamento, agregando valor para a solução minimizando a sensibilidade relacionada com a variação da luz do ambiente. Assim, a imagem normalizada por filtros de correção de iluminação ressaltará melhor os traços faciais, além da imagem transformada estar com maior nitidez para a continuação do *workflow*.

Este trabalho também foca em minimizar os efeitos negativos de rotações da face e variação da pose para treinamento e classificação. Graus elevados de rotações e poses dificultam a classificação e minimizam o aprendizado, caso a base de treino tenha muita imagem com face rotacionada. Por isso, no pré-processamento é feito o alinhamento da

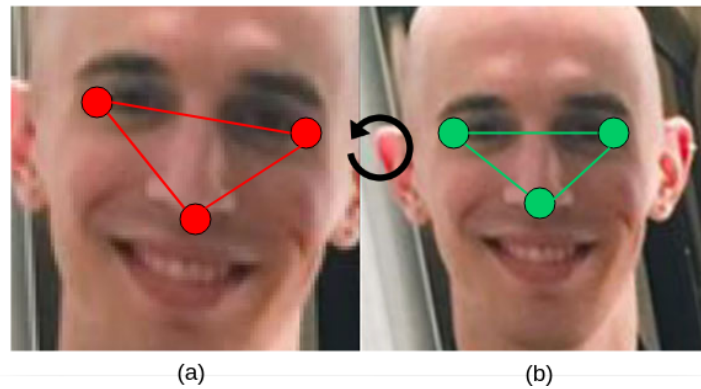


Figura 13 – Algoritmo de Alinhamento da Face: (a) Imagem Original; (b) Face Alinhada.



Figura 14 – Pré-Processamento fluxo

face. Este alinhamento é por meio da localização de três pontos: o canto do olho direito, o canto do olho esquerdo e um ponto central do nariz. Desta forma, esses três pontos formam um triângulo, sendo possível rotacionar a imagem até que não haja uma inclinação na linha dos olhos orientada pelo triângulo. A Figura 13 ilustra o exemplo de uma face não alinhada e após o alinhamento. Por fim, nesta etapa, é realizado a normalização da imagem dividindo cada *pixel* por 255, isto é, o valor máximo que um *pixel* pode possuir. Resultando na normalização dos valores dos *pixels* entre o intervalo de 0 a 1.

## 4.4 Rede Neural de Convolução

A rede neural de convolução é a parte central e com importante contribuição nesta abordagem. Por meio dela, durante o treinamento, as imagens são processadas a fim de aprender os contornos, padrões, formas e características relevantes para a classificação. Além disso, funções são aplicadas para redução de dimensionalidade, extração de características e normalização, ocasionando que a rede não seja sensível a rotações, posições e escala da imagem. Tais aptidões são requisitos para um classificador de imagens ser usado em cenários reais. Apesar de a rede neural ter embutida normalizações e rotações internas, vale a pena realizar as operações da Seção 4.3, pois são técnicas selecionadas especialmente para o problema de faces, e além disso, a comunidade científica tem descoberto que combinar o processamento interno da rede neural de convolução com outros processamentos externos alcançam os melhores resultados (Kim et al., 2016b; Shan et al., 2017; Shin et al., 2016b). Vale ressaltar que um desempenho satisfatório da rede neural de convolução, assim como de qualquer algoritmo supervisionado de aprendizagem de



Figura 15 – Exemplos da técnica de aumento de dados: (a) Imagem Original; (b) Redução de Contraste; (c) Aumento de Contraste; (d) Perspectiva; (e) *Crop*; (f) *Shear*.

máquina, está estritamente relacionado ao processo de treinamento e validação do modelo.

#### 4.4.1 Treinamento

O treinamento da rede neural de convolução é parte fundamental para o classificador de emoção funcionar bem. Nesta etapa, é buscado tanto a generalização satisfatória quanto a aquisição do aprendizado suficiente para funcionar em variados ambientes. Para isso, neste trabalho, o treinamento é apoiada pela técnica de aumento de dados com intuito de maximizar a generalização do aprendizado durante o treinamento. O aumento de dados consiste na multiplicação das imagens em tempo dinâmico modificando levemente a imagem e seu contexto, realizando alterações nas imagens aplicando *crop*, rotações, perspectiva, *shear*, diferentes níveis de contraste e dentre outros. Com o propósito de estimular maiores taxas de aprendizado e generalização, pois, assim, a rede neural observa a região de interesse que são as expressões faciais emocionais em diferentes cenários, sendo assim, gerando modelos capazes de classificar emoções em contextos variados. A Figura 15 ilustra alguns exemplos da técnica de aumento de dados.

#### 4.4.2 Extração de Características e Classificação

Em problemas de classificação de imagem, uma etapa fundamental é a extração de características, sendo responsável por retirar de uma entrada de dados as principais informações para enviar a um classificador e, assim, determinar a classe. Na rede neural de convolução, a extração de característica é um procedimento que procura identificar as zonas da imagem que são mais relevantes para a separação do problema, isto é, classificar uma expressão facial (*e.g* o sorriso humano é uma característica indicadora para a emoção felicidade). Este processo funciona com as camadas de convolução da rede operando sobre a imagem retirando as informações relevantes e com a camada de *pooling* aplicando a

redução de dimensionalidade. Algumas camadas de convolução tornam-se especialistas na extração dos padrões verticais, outras nos horizontais, em um padrão geométrico específico, até que seja gerado um conjunto de características para enviar a um classificador.

O classificador recebe as características extraídas, e este fica localizado na última camada da rede neural de convolução. O conjunto extraído é um vetor de uma dimensão que possui uma quantidade de elementos bastante inferior a imagem original. Seu conteúdo consiste no conjunto de informações representativas para diferenciação das emoções que a rede neural aprendeu a extrair no treinamento. Então, o classificador é treinado para separar as emoções. Neste trabalho, o classificador adotado (*softmax*) possui a característica de estimar a probabilidade para cada emoção: neutralidade, raiva, felicidade, tristeza, desprezo, medo e surpresa. E essa estimativa é distribuída entre as classes (*e.g* neutralidade: 0.95, felicidade: 0.25 e surpresa: 0.25), possuindo a propriedade em que o somatório das probabilidades é igual a 1, e a emoção eleita é a que tem maior probabilidade, nesse caso, a neutralidade com 0.95. A escolha do *softmax* é devido ao mesmo fazer parte da configuração padrão das arquiteturas testadas. Em nossa abordagem, após a classificação de uma imagem, as estimativas de probabilidades geradas são salvas em um repositório de saída de dados. É importante destacar que em uma imagem com múltiplas faces são classificadas uma por vez, pois é mais simples reconhecer a emoção de uma única face do que em todas ao mesmo tempo.

## 4.5 Resumo

Neste capítulo foi descrito a abordagem proposta. As principais etapas são: detecção e recorte da face, pré-processamento, extração de características e classificação. A detecção de face consiste na utilização do algoritmo Viola-Jones. Esta etapa gera dois benefícios valiosos que reduzem a complexidade do problema: a exclusão do *background* e o recorte individual de cada face. A etapa de pré-processamento é constituída da aplicação de vários filtros e funções de normalizações na imagem para limpeza, eliminação de ruídos e alinhamento da face. A extração de características é feita pela rede neural de convolução em que procura na imagem pré-processada as informações interessantes para a diferenciação da emoção. Por fim, a etapa de classificação, onde o classificador está localizado na última camada da rede neural de convolução e recebe as informações extraídas para estimar a probabilidade para cada emoção, determinando as emoções da imagem.

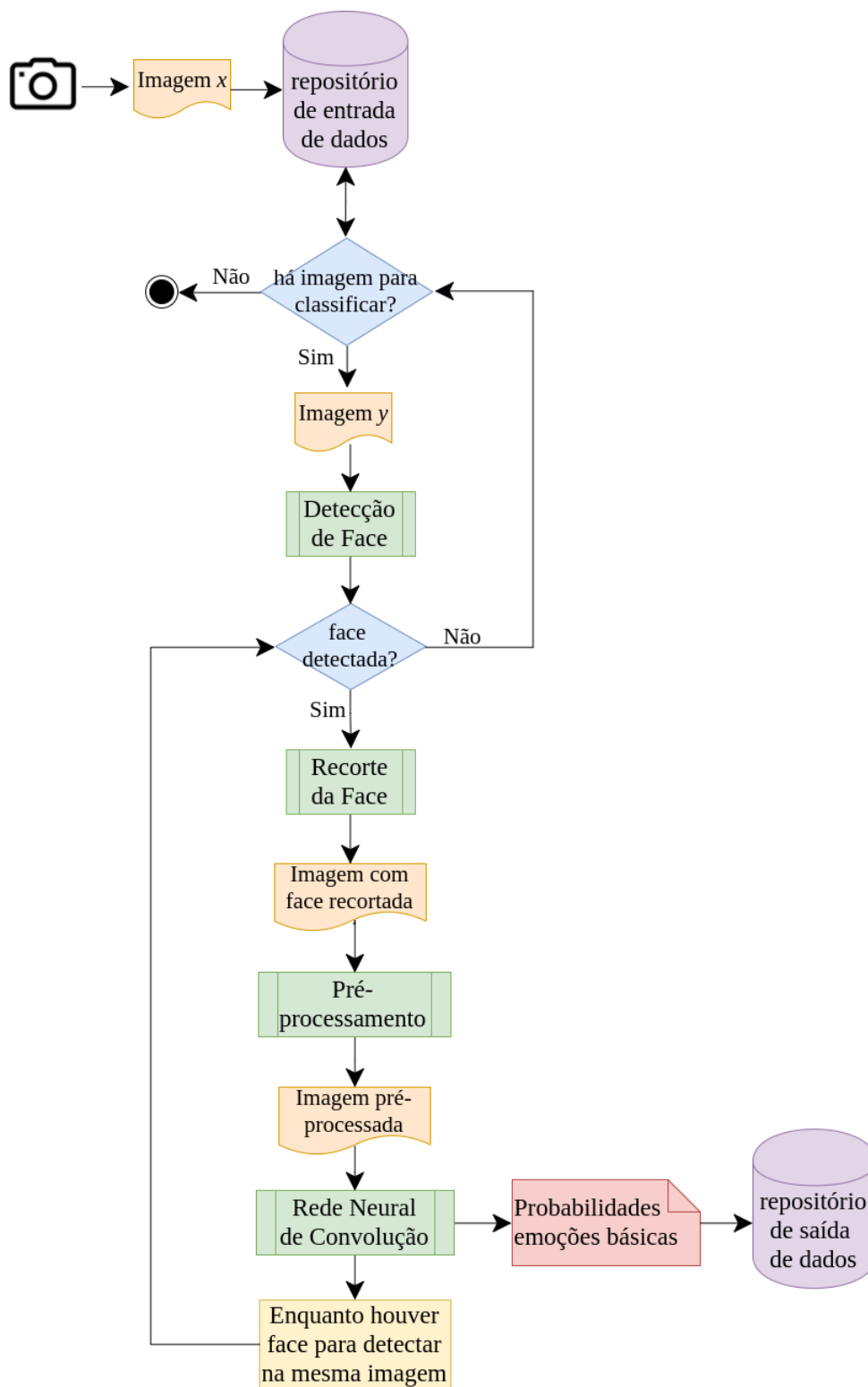


Figura 16 – Solução Proposta

## 5 Resultados Parciais

Neste capítulo os resultados parciais são apresentados divididos nas seguintes seções. A Seção 5.1 destaca um *framework* para coleta e inferência de estados emocionais de estudantes, enquanto a Seção 5.2 apresenta um estudo experimental sobre as arquiteturas de redes neurais de convolução envolvendo o domínio do problema.

### 5.1 Prova de Conceito: Coleta e Inferência de Estados Emocionais de Estudantes

Em (Cruz et al., 2017) foi proposto uma prova de conceito com a intenção de nortear o andamento deste projeto. Este trabalho consistiu na utilização da API da *Microsoft Cognitives Services*, um serviço na nuvem que oferece um conjunto de aplicações cognitivas, inclusive um reconhecedor de emoções via expressão facial. Esta experiência foi bastante positiva possibilitando ponderar as vantagens e desvantagens desta API, e assim, nortear a construção do reconhecedor de emoções proposto. A contribuição principal deste trabalho está em reconhecer emoções em tempo real, em um cenário real e de forma automática, correlacionando as emoções detectadas com o desempenho escolar em um teste.

A revisão sistemática (ver Capítulo A) mostrou que os pesquisadores da área de reconhecimento de emoção ainda estão tímidos para colocar esses reconhecedores no mundo real. Desta forma, a utilização em cenários reais passa a ser um ponto de contribuição, atestando que essas tecnologias estão amadurecidas possibilitando a integração com outros sistemas, principalmente para tomadas de decisão, recomendações, análises de comportamento dentre outros. Contudo, alcançar precisão no cenário real não é uma tarefa trivial, visto que o reconhecedor deve alcançar uma excelente taxa de generalização. O mundo real é desafiador, por haver muitas variações da face e do ambiente, dificultando a precisão do algoritmo.

#### 5.1.1 Síntese

Em (Cruz et al., 2017) foi proposto um *framework* para detectar estados emocionais de alunos baseado em reconhecimento de expressões faciais no contexto das plataformas digitais educacionais. Neste trabalho, foi analisado e discutido o uso de correlação e entropia entre os estados emocionais dos estudantes e o desempenho durante uma avaliação de múltipla escolha. Foi realizado um experimento com 27 estudantes a partir de um questionário avaliativo composto de 40 perguntas. Na análise dos dados, os estados emocionais

de neutralidade, tristeza, felicidade, raiva, desgosto, medo, desprezo e surpresa, foram correlacionados com o desempenho no teste. O experimento concluiu que as questões que ocorreram maior variabilidade das emoções tinham também as maiores proporções de acertos.

### 5.1.2 Objetivos

O artigo (Cruz et al., 2017) teve como objetivo: (i) propor uma arquitetura de detecção automática de emoções para ambientes educacionais digitais por meio de reconhecimento automático de expressões faciais, utilizando processamento de imagens, de modo a tornar possível a obtenção de dados emocionais dos alunos durante o processo de aprendizagem, e (ii) analisar por meio de estatística descritiva um estudo de caso com dados obtidos a partir desta arquitetura. Esta análise pretendeu investigar e medir correlações entre as emoções e o desempenho obtido nas questões, levantando hipóteses relevantes sobre a relação entre os estados emocionais e o desempenho, o qual é relevante para sistemas de recomendações, tutores inteligentes e heurísticas em geral, que se interagem de forma dinâmica com as necessidades de aprendizado de cada aluno.

### 5.1.3 Metodologia Experimental

Um experimento foi realizado com 27 alunos do Ensino Médio de uma escola de tempo integral que, na época, estavam se preparando para o Exame Nacional do Ensino Médio (ENEM) de 2017. O experimento consistiu em um simulado do exame contendo 40 questões de múltipla escolha.

#### 5.1.3.1 Planejamento

O experimento foi feito a partir de uma plataforma educacional para tornar possível a execução de questionários de múltipla escolha, coleta de cliques efetuados pelo estudante, captura automática de foto via câmera frontal do dispositivo, seja por *tablet*, *smartphone* ou *notebook*. Os assuntos escolhidos foram: matemática, língua portuguesa, química, raciocínio lógico, geografia e história. O simulado teve duração de duas horas e cada questão possuía dois níveis de dificuldade (fácil ou difícil), além de ter cinco respostas alternativas.

Para o componente classificador, foi selecionada a API da *Microsoft Cognitives Services*, justamente por classificar bem emoções como neutralidade, felicidade e tristeza. Acreditamos que, no contexto da educação, uma emoção bastante comum é a neutralidade, pelo fato da expressão facial do estado de concentração se assemelhar bastante com a expressão facial de neutralidade reconhecida por esta API. Constatamos que estudantes



quando estão pensando, estão concentrados, emitindo poucas movimentações intensas e variações de suas expressões faciais, assemelhando-se com a expressão de neutralidade.

#### 5.1.3.2 Execução

A seleção dos estudantes para participar do experimento ocorreu voluntariamente. O grupo final formado foi heterogêneo, onde 53% consideravam até o momento seu desempenho na escola como bom ou ótimo e, 47% como regular ou ruim; além disso, 30% deles consideravam a sua preparação para o vestibular como boa ou ótima e, 70% como regular ou fraca. Os alunos selecionados foram de turmas diferentes.

#### 5.1.4 Resultados e Discussões

Os seguintes atributos foram calculados como os valores médios por questão, utilizando os dados coletados dos vinte e sete alunos: (i) a proporção de acertos; (ii) o nível de dificuldade; (iii) a média das probabilidades para cada emoção detectada; e (iv) a entropia por questão. Assim, um total de onze atributos (resumidos na Tabela 4), onde cada um deles é representado por uma variável aleatória com quarenta valores.

Posteriormente, a correlação de Pearson foi aplicada para analisar se há qualquer grau de correlação entre os pares dos atributos mencionados. A Tabela 4 apresenta a correlação entre a média das probabilidades de cada emoção com o nível de dificuldade e a proporção de acertos. Os principais resultados estão destacados em negrito.

É possível verificar na Tabela 4, a expressão facial neutra possui uma correlação negativa com a proporção de acertos dos alunos (segunda coluna da Tabela 4). Isto indica que estimular emoções diferentes da neutralidade durante a avaliação favorece o desempenho dos alunos. Um segundo indicativo de que isto ocorre, é dado pela correlação positiva entre o desprezo e a felicidade com a proporção de acertos e, de forma discreta, também ocorre com tristeza, surpresa e medo.

A entropia é calculada a partir do valor das probabilidades das emoções detectadas. Logo, percebemos que quando a neutralidade é baixa, a entropia aumenta, o que significa que outras emoções estão sendo detectadas com maior probabilidade, ocorrendo a dispersão dos estados emocionais. Portanto, o fato de existir uma correlação positiva entre o aumento da entropia e a proporção de acertos reforça a observação constatada no parágrafo anterior. Adicionalmente, a emoção mais frequente foi a neutralidade, devido aos alunos passarem a maior parte do tempo concentrados analisando as questões para a busca de soluções. Assim, quando o nível de dificuldade da questão aumenta, a neutralidade também aumenta, isto pode ser um indício de que questões mais difíceis tem tendências de exigir maiores níveis de concentração do estudante.

É possível considerar a hipótese de que, quando o estudante está respondendo

Tabela 4 – Resultado da correlação de Pearson para cada emoção detectada e a entropia contra os atributos das questões

	Nível de Dificuldade	Proporção de Acertos
<b>Tristeza</b>	<b>-0.33</b>	0.27
<b>Neutralidade</b>	<b>0.36</b>	<b>-0.48</b>
<b>Desprezo</b>	-0.15	<b>0.30</b>
Desgosto	-0.13	0.07
Raiva	-0.14	-0.08
Surpresa	0.07	0.24
Medo	-0.06	0.14
<b>Felicidade</b>	-0.14	<b>0.31</b>
<b>Entropia</b>	-0.12	<b>0.36</b>

uma questão, ao selecionar uma resposta, o mesmo tem uma percepção se acertou ou errou e, nesse momento, há possibilidade de emitir emoções positivas como felicidade e surpresa, ou emoções negativas como tristeza ou desprezo. Portanto, há uma variação dos estados emocionais durante o tempo de resposta de cada questão que deve ser considerado como um problema de mudança de estados. Este resultado é reforçado por questões que ocasionaram maior entropia, ou seja, quanto maior a dispersão das emoções, maior é o índice de proporções de acertos.

A emoção desprezo aumenta a medida que as questões têm maiores proporções de acertos, isto pode ser explicado, pela mudança de estados durante o tempo de resposta de cada questão ou pelo fato da expressão facial de desprezo se assemelhar com a expressão facial de felicidade. Neste caso, é bem provável estar ocorrendo confusão por parte do classificador em diferenciar felicidade e desprezo.

Finalmente, percebeu-se que o *framework* para classificação automática de emoções a partir de imagens necessita de um maior volume de dados, para reduzir eventuais erros de classificação, e permitir a classificação de novos tipos de emoções.

## 5.2 Avaliação Experimental de Redes Neurais de Convolução

Um estudo experimental foi realizado a fim de comparar as seguintes arquiteturas de redes neurais de convolução: Alexnet, Residual Net e Inception-V3. Este estudo tem como objetivo eleger a arquitetura que gerou o melhor modelo avaliando a precisão, revocação, f1-score e a acurácia. Lembrando que a métrica de f1-score significa a média harmônica entre precisão e revocação (veja Seção 2.9).

### 5.2.1 Preparação dos Dados

A preparação dos dados consistiu na formação de três bases de dados: treino, teste e validação. A base de treino foi utilizada para treinar os algoritmos e, geralmente, é a base que tem mais instâncias. Uma base de treino formada erroneamente reflete nos modelos

gerados ocasionando *underfitting* (não aprendeu a resolver o problema) ou *overfitting* (não aprendeu a generalizar o problema ou apresenta super vício na base de treino), dessa forma, provendo modelos não confiáveis que apresentam problemas de acurácia. Diante disso, a preparação dos dados é essencial para obtenção de modelos com bons resultados. A base de teste é usada para validar a rede neural durante o treinamento, em que após cada interação é calculada a função de perda. O valor de perda é importante para diagnosticar como está o treinamento, isto é, se o gradiente descendente está convergindo para o ponto de parada, ou seja, o momento certo de interromper o treinamento ou verificar se a rede neural apresenta *underfitting* (ainda falta treinar) ou *overfitting* (treinou muito, vício na base de treino). A base de validação é composta por instâncias totalmente desconhecidas pelo modelo, pois são imagens que não fizeram parte do treinamento. A base de validação seria uma espécie de teste para verificar como determinado modelo se comportaria no mundo real.

Um dos objetivos da revisão sistemática (ver Seção A) foi encontrar as bases de dados mais populares utilizadas pela comunidade científica. A Tabela 5 apresenta as bases de dados localizadas na revisão sistemática. Os dados advindos dessas bases foram usadas neste estudo experimental. Partindo do pressuposto que as bases de dados aplicadas em aprendizagem de máquina devem ser bastantes diversificadas para gerar modelos com bons resultados. As bases de treino, teste e validação, foram formadas a partir da concatenação de todas as bases a fim de gerar a diversificação nos dados na seguinte divisão: 50% para treino e 25% para teste e 25% validação. Uma etapa de limpeza foi executada em todas as bases com intuito de retirar amostras que possuem baixa representatividade para o problema. Essa etapa de limpeza consistiu na verificação da face seguindo a abordagem proposta na Seção 4.2, pois há várias imagens com forte ruídos, por exemplo, rotações com grau elevado, faces incompletas e imagens com faces de personagens de animação. Após o processo de limpeza, as imagens resultantes foram distribuídas conforme a Tabela 6, sendo divididas por base de treino, teste e validação. A Tabela 7 indica a distribuição das imagens referentes as classes (emoções). Vale ressaltar que a emoção felicidade é a classe que tem mais instâncias representativas, ao contrário da emoção medo que possui a menor amostra.

### 5.2.2 Materiais

Para a realização desse estudo experimental foi necessário o uso de alguns materiais. Para implementação da rede neural de convolução foi utilizado o *framework TensorFlow* (<https://www.tensorflow.org/>) e a API de alto nível para *tensorflow* a *TFLearn* (<http://tflearn.org/>). A biblioteca *OpenCV 3.0* (<https://opencv.org/>) para realizar processamento na imagem, e também, para detectar face, pois nesta biblioteca tem uma implementação do algoritmo Viola Jones que é o algoritmo selecionado para a abor-

Tabela 5 – Bases de dados localizadas na revisão sistemática da literatura

Bases de Dados	Trabalhos que utilizaram a base para treinamento ou validação
CK+	Chen et al. (2017b), Shan et al. (2017), Wen et al. (2017b), Shin et al. (2016b), Huang and Lu (2016), Vo and Le (2016), Yu et al. (2016c), Mayya et al. (2016), Jung et al. (2015b), Li et al. (2015b)
JAFPE	Chen et al. (2017b), Shan et al. (2017), Wen et al. (2017b), Shin et al. (2016b), Mayya et al. (2016)
FER	Wen et al. (2017b), Kim et al. (2016b), Liu et al. (2016b), Shin et al. (2016b), Huang and Lu (2016), Guo et al. (2016), Ng et al. (2015), Jung et al. (2015b)
FER+	Barsoum et al. (2016b)
SFEW2.0	Shin et al. (2016b), Guo et al. (2016)
KDEF	Shin et al. (2016b)
MMI	Yu et al. (2016c)
CIFE	Li et al. (2015b)
EmotiW2015	Wen et al. (2017b), Ng et al. (2015)

Tabela 6 – As bases de dados foram concatenadas e divididas em três bases: treino, teste e validação. Na seguinte porcentagem: 50% para treino e 25% para teste e validação.

Base de Dados	B. de Treino	B. de Teste	B. de Validação	Total de Imagens
RAFD	2408	1206	1205	4819
CIFE-TRAIN	4086	2042	2042	8170
CIFE-TEST	1759	879	878	3516
CK	1509	754	755	3018
KDEF	1466	735	733	2934
JAFPE	105	53	55	213
NOVAEMOTIONS	16840	8418	8417	33675
FER	11782	5892	5891	23565
Total de Imagens	39955	19979	19976	79910

dagem proposta no Capítulo 4. Os experimentos foram realizados em um computador com a seguinte configuração: *GPU NVIDIA GEFORCE 930, Intel Core-i7 e 16 GB de RAM DDR4*.

### 5.2.3 Arquiteturas

Este estudo experimental utilizou três arquiteturas de redes neurais de convolução: Alexnet, Inception e Residual Net. A arquitetura Alexnet foi implementada de forma original descrita na Tabela 1. A Inception possui diversas variantes inclusive fazendo parte da família GoogLeNet, porém neste estudo consideramos a Inception-V3 que está ilustrada na Figura 17. A Residual Net também tem várias variantes. A selecionada para o estudo foi a que possui 34 camadas denominada ResNet-34 mostrada na Figura 18. O motivo da escolha tanto para Inception-V3 como para ResNet-34 foi devido a limitação de hardware que este projeto dispõe para execução de experimentos, portanto foram selecionadas as variantes menos complexas destas arquiteturas que por natureza são complexas

Tabela 7 – Distribuição das classes (emoções) nas bases de treino, teste e validação. As classes também foram divididas em: 50% para treino e 25% para teste e validação.

Classe	B. de Treino	B. de Teste	B. de Validação	Total de Imagens
Raiva	3299	1650	1650	6599
Desgosto	2453	1226	1226	4905
Medo	2821	1411	1410	5642
Felicidade	13943	6971	6971	27885
Tristeza	4349	2175	2174	8698
Surpresa	6311	3156	3155	12622
Neutralidade	6779	3390	3390	13559
Total de Imagens	39955	19979	19976	79910

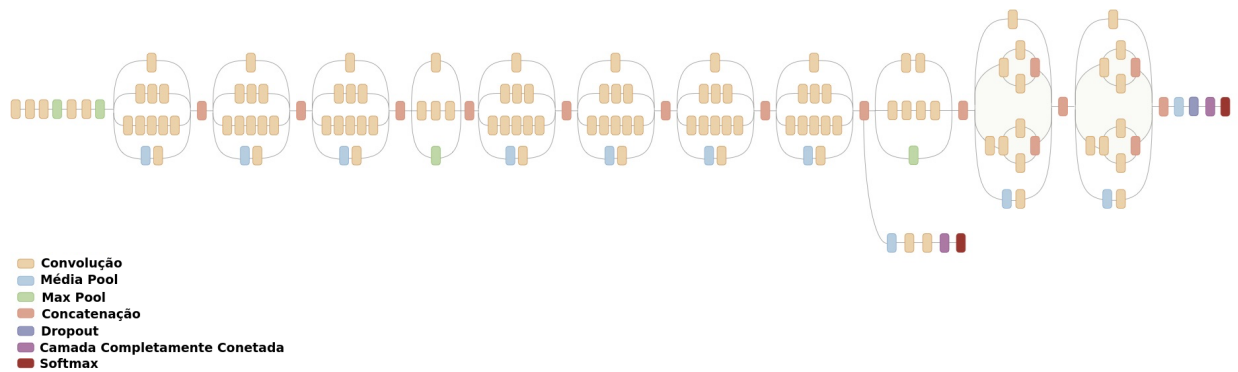


Figura 17 – Arquitetura Inception-V3

e profundas.

#### 5.2.4 Treinamento

O treinamento consistiu no uso dos materiais descritos na Seção 5.2.2 e das bases de treino e teste descritas nas Tabelas 6 e 7. A estratégia adotada durante a fase de treinamento foi salvar vários modelos enquanto a função de perda convergia para serem utilizadas em comparações. Para todos os experimentos o tamanho de *batch* foi de 64 imagens e a taxa de aprendizado foi de 0.001 para Alexnet e Inception-V3, no entanto para ResNet-34 foi de 0.1. Obviamente que isso resultou que a ResNet-34 fizesse o treinamento mais rápido que as demais, fato que pode ser constatado pela Figura 20. Os parâmetros de *batch* e taxa de aprendizado definidos são os valores padrões de cada arquitetura, por isso não houve qualquer alteração.

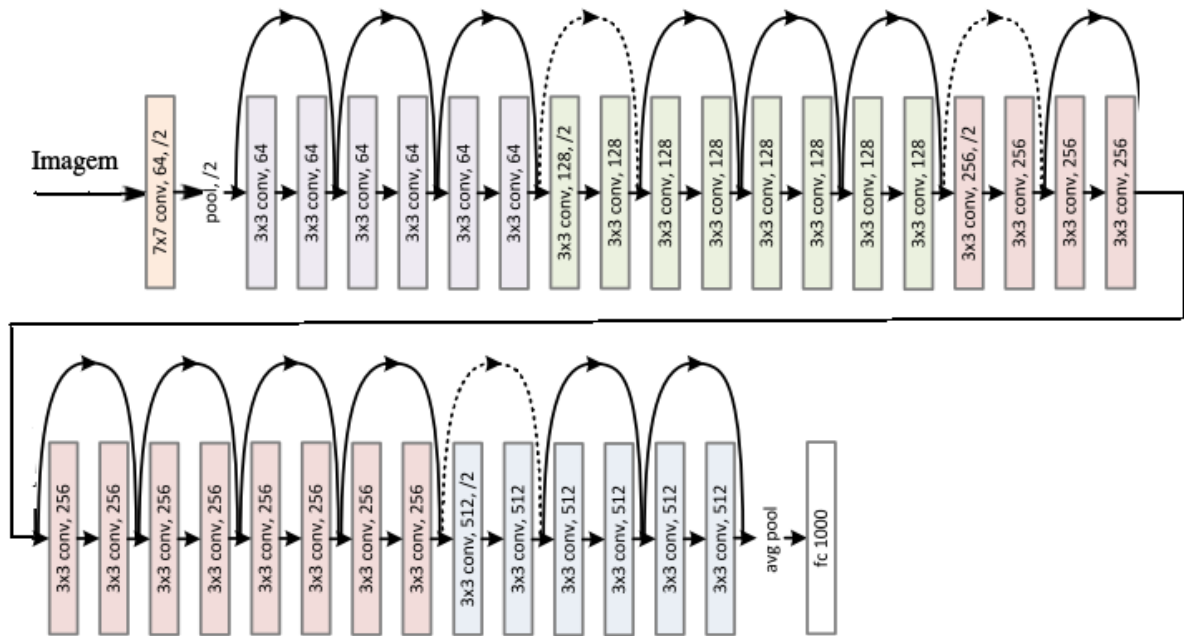


Figura 18 – Arquitetura ResNet-34

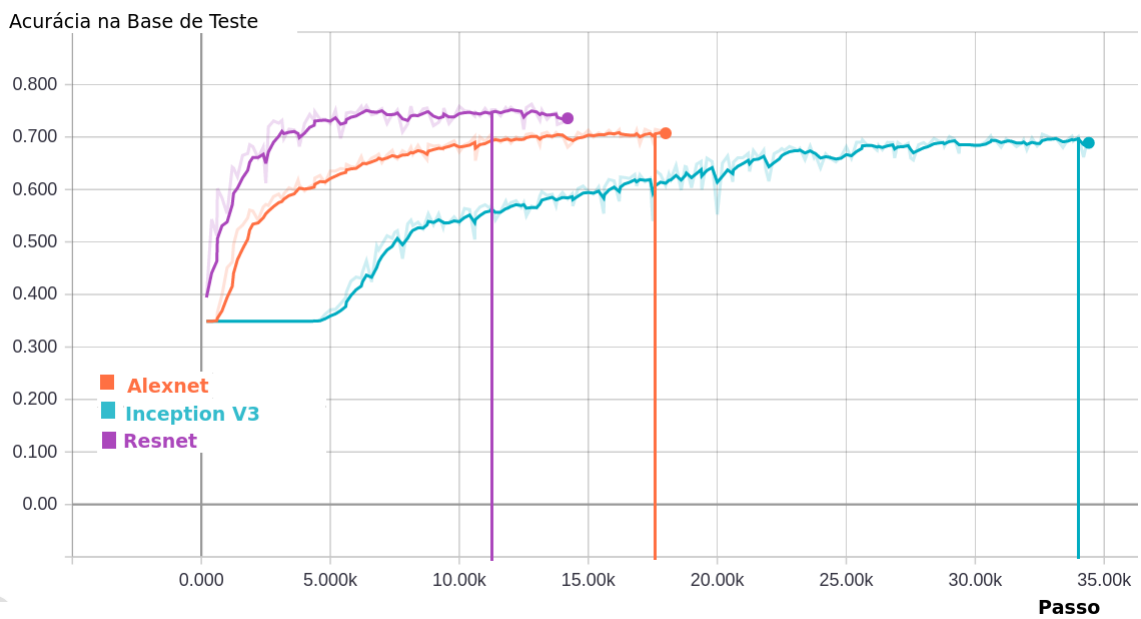


Figura 19 – Gráfico de Acurácia na Base de Teste. As linhas verticais indicam o melhor modelo.

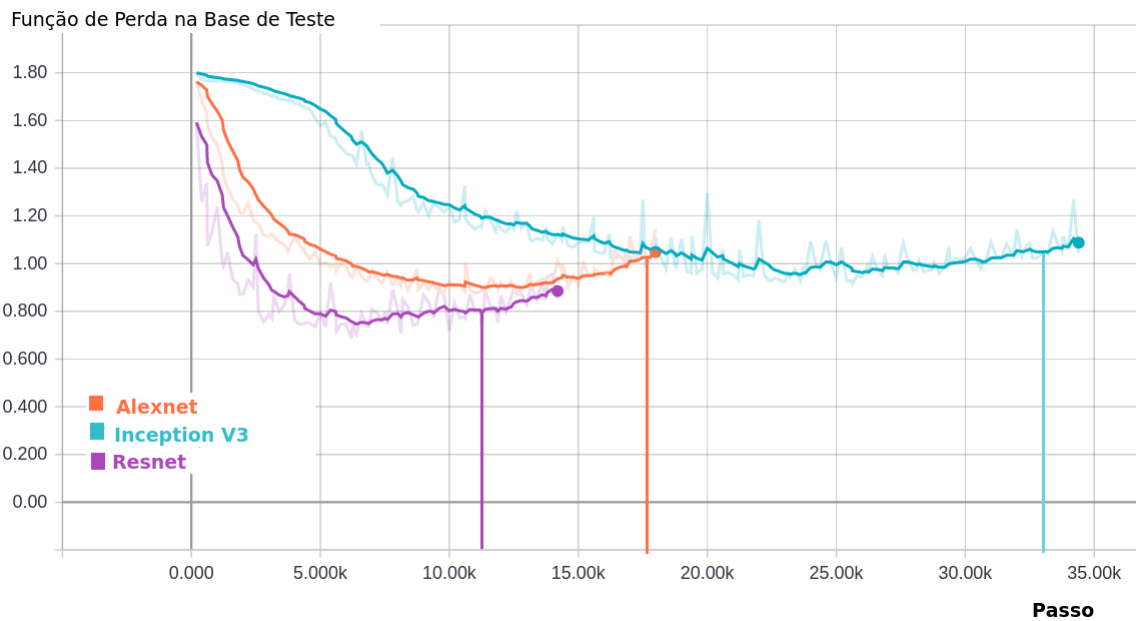


Figura 20 – Gráfico da Função de Perda na Base de Teste. As linhas verticais indicam o melhor modelo.

Tabela 8 – Resultados experimentais das redes neurais de convolução avaliando a base de validação geral.

Arquitetura	Emoção	Precisão	Revocação	F1-score	Acurácia
Alexnet	Raiva	0.51	0.60	0.55	0.712
	Desgosto	0.62	0.64	0.63	
	Medo	0.47	0.41	0.44	
	Felicidade	0.84	0.89	0.86	
	Tristeza	0.64	0.50	0.56	
	Surpresa	0.84	0.77	0.80	
	Neutralidade	0.62	0.64	0.63	
	Média/Total	0.71	0.71	0.71	
Inception-V3	Raiva	0.54	0.51	0.52	0.701
	Desgosto	0.56	0.57	0.56	
	Medo	0.47	0.42	0.44	
	Felicidade	0.88	0.88	0.88	
	Tristeza	0.47	0.53	0.50	
	Surpresa	0.85	0.79	0.82	
	Neutralidade	0.59	0.62	0.61	
	Média/Total	0.70	0.70	0.70	
ResNet-34	<b>Raiva</b>	<b>0.69</b>	<b>0.57</b>	<b>0.62</b>	<b>0.757</b>
	<b>Desgosto</b>	<b>0.79</b>	<b>0.66</b>	<b>0.72</b>	
	<b>Medo</b>	<b>0.45</b>	<b>0.50</b>	<b>0.47</b>	
	<b>Felicidade</b>	<b>0.90</b>	<b>0.89</b>	<b>0.90</b>	
	<b>Tristeza</b>	<b>0.60</b>	<b>0.65</b>	<b>0.63</b>	
	<b>Surpresa</b>	<b>0.82</b>	<b>0.86</b>	<b>0.84</b>	
	<b>Neutralidade</b>	<b>0.67</b>	<b>0.68</b>	<b>0.68</b>	
	<b>Média/Total</b>	<b>0.76</b>	<b>0.76</b>	<b>0.76</b>	

Tabela 9 – Resultados experimentais das redes neurais de convolução avaliando a base de validação CK

Arquitetura	Emoção	Precisão	Revocação	F1-score	Acurácia
Alexnet	Raiva	0.91	1	0.95	0.96
	Desgosto	0.98	0.97	0.98	
	Medo	0.89	0.96	0.92	
	Felicidade	0.99	0.99	0.99	
	Tristeza	0.98	0.84	0.91	
	Surpresa	1	0.94	0.97	
	Neutralidade	0	0	0	
	Média/Total	0.97	0.96	0.96	
Inception-V3	Raiva	0.93	0.94	0.93	0.954
	Desgosto	0.96	0.94	0.95	
	Medo	0.89	0.96	0.92	
	Felicidade	0.99	0.98	0.99	
	Tristeza	0.91	0.97	0.94	
	Surpresa	1	0.94	0.97	
	Neutralidade	0	0	0	
	Média/Total	0.96	0.95	0.96	
ResNet-34	<b>Raiva</b>	<b>0.97</b>	<b>0.96</b>	<b>0.97</b>	<b>0.969</b>
	<b>Desgosto</b>	<b>1</b>	<b>0.92</b>	<b>0.96</b>	
	<b>Medo</b>	<b>0.91</b>	<b>0.99</b>	<b>0.95</b>	
	<b>Felicidade</b>	<b>0.98</b>	<b>0.99</b>	<b>0.99</b>	
	<b>Tristeza</b>	<b>0.94</b>	<b>0.96</b>	<b>0.95</b>	
	<b>Surpresa</b>	<b>0.98</b>	<b>0.99</b>	<b>0.99</b>	
	<b>Neutralidade</b>	<b>0</b>	<b>0</b>	<b>0</b>	
	<b>Média/Total</b>	<b>0.97</b>	<b>0.97</b>	<b>0.97</b>	

## 5.2.5 Discussões

Durante o treinamento vários modelos foram salvos. Adotamos o melhor modelo gerado por cada arquitetura para esta seção. Logo, os três modelos eleitos a partir de uma avaliação na base de validação. Os resultados referentes a todos os modelos podem ser consultados no Anexo B. Analisando o anexo somente por arquitetura podemos concluir que o melhor modelo da Alexnet foi gerado no passo 17800, enquanto que a ResidualNet-34 foi no passo 11250 e por fim a Inception-V3 em 34000.

As Figuras 19 e 20 são a respeito da acurácia na base de teste e da função de perda durante o treinamento. Ambas possuem linhas verticais demarcando os modelos selecionados para a discussão ao longo da seção. Analisando a Figura 19, a taxa de acurácia na base de teste, a ResNet-34 no início do treino tem um crescimento semelhante a uma função exponencial e, a partir da interação 5k estabiliza alcançando 76%. A Alexnet também tem uma curva bastante crescente no início, entretanto foi inferior a ResNet-34 alcançando 70.5%, uma diferença considerável. A Inception-V3 teve um comportamento diferente, sendo que no início do treino a função ficou estática, e posteriormente, foi



Tabela 10 – Resultados experimentais das redes neurais de convolução avaliando a base de validação FER

Arquitetura	Emoção	Precisão	Revocação	F1-score	Acurácia
Alexnet	Raiva	0.39	0.5	0.44	0.543
	Desgosto	0.45	0.17	0.25	
	Medo	0.37	0.33	0.35	
	Felicidade	0.74	0.79	0.76	
	Tristeza	0.4	0.29	0.34	
	Surpresa	0.72	0.64	0.67	
	Neutralidade	0.49	0.52	0.51	
	Média/Total	0.54	0.54	0.54	
Inception-V3	Raiva	0.43	0.4	0.41	0.529
	Desgosto	0.13	0.25	0.17	
	Medo	0.39	0.36	0.37	
	Felicidade	0.81	0.77	0.79	
	Tristeza	0.29	0.38	0.32	
	Surpresa	0.72	0.64	0.68	
	Neutralidade	0.49	0.46	0.47	
	Média/Total	0.55	0.53	0.54	
ResNet-34	<b>Raiva</b>	<b>0.61</b>	<b>0.42</b>	<b>0.5</b>	<b>0.604</b>
	<b>Desgosto</b>	<b>0.69</b>	<b>0.28</b>	<b>0.39</b>	
	<b>Medo</b>	<b>0.35</b>	<b>0.47</b>	<b>0.4</b>	
	<b>Felicidade</b>	<b>0.87</b>	<b>0.81</b>	<b>0.84</b>	
	<b>Tristeza</b>	<b>0.41</b>	<b>0.42</b>	<b>0.41</b>	
	<b>Surpresa</b>	<b>0.71</b>	<b>0.76</b>	<b>0.73</b>	
	<b>Neutralidade</b>	<b>0.56</b>	<b>0.61</b>	<b>0.58</b>	
	<b>Média/Total</b>	<b>0.62</b>	<b>0.6</b>	<b>0.61</b>	

crescendo, no entanto assemelhando-se ao crescimento de uma função linear. A Inception-V3 ficou abaixo da Alexnet, alcançando 70% de acurácia na base de teste.

A função de perda definida foi a entropia cruzada. A entropia cruzada é apropriada para problemas de classificação. Analisando a Figura 20, quanto menor for o valor melhor está sendo o desempenho do classificador, pois esta função é um indicador do erro propagado na rede e o objetivo do gradiente descendente é minimizar estes valores. Logo, quem obteve o menor valor de função de perda foi a Resnet, e consequentemente, a que alcançou melhor desempenho em acurácia.

Verificando os pontos da função de perda da Inception-V3, é visto que nas primeiras épocas não há uma queda acentuada como nas outras duas. O que coincide com a demora na inclinação da medida de acurácia, obviamente também relacionada a Inception-V3. De acordo com Géron (2017), os melhores modelos são os gerados a partir do momento em que a função de perda na base de teste atinge seu valor mínimo até o momento em que começa a crescer. Considerando esta teoria o momento certo de interromper o treinamento é justamente este, quando a função de perda começa a crescer após encontrar o mínimo,

este conceito é conhecido como parada antecipada, pois os modelos gerados adiantes possuem *overfitting* (super vícios) na base de treino.

A Tabela 8 apresenta os resultados das redes neurais de convolução avaliando a base de validação geral. Vale ressaltar que a base de validação geral seria uma medição do desempenho dos modelos em cenário real, pois a base de validação não fez parte do treinamento, isto é, são instâncias desconhecidas. Analisamos que a ResNet-34 gerou o modelo que obteve os melhores resultados em todas as métricas: acurácia, precisão, revocação e f1-score, sendo destacada em negrito. A ResNet-34 alcançou 75.7% em acurácia, enquanto a Alexnet 71.2% e Inception-V3 70.1%. Percebemos uma diferença de 4.5% entre a ResNet-34 e Alexnet caracterizando uma diferença considerável.

A emoção que a ResNet-34 teve melhor desempenho para reconhecer foi a felicidade alcançando um f1-score em 90%. Provavelmente pela característica da felicidade ser atribuída a um sorriso. A segunda emoção com melhor resultado foi a surpresa com f1-score em 84%. Esta emoção tem como marca a boca aberta, outra característica que a rede aprendeu. Este resultado não surpreende, pois essas emoções são duas das três emoções que mais possuem amostragem (consultar Tabela 7). O número de amostra está associado ao aprendizado, teoricamente o cenário ideal é ter a base balanceada ou a quantidade de amostras proporcionais ao mundo real.

Os piores desempenhos foram nas emoções medo e raiva. Coincidentemente duas das três emoções que menos tem amostragem. Este resultado não necessariamente foi devido a menor amostragem dessas classes, pois realmente são expressões faciais complicadas. A expressão facial medo não há um padrão bem definido e envolve vários músculos do rosto. Vale ressaltar que a construção de um sistema que possui requisitos de alta precisão nessas duas emoções seja necessário outras fontes de dados além da expressão facial.

A neutralidade foi uma das três emoções que mais teve amostragem. Entretanto, alcançou somente 68% de f1-score, portanto um baixo valor. Este resultado pode ser explicado pelo fato da rede neural ter aprendido que a neutralidade consiste em nenhum movimento muscular da face. Caso haja um leve movimento facial a rede tende a classificar em uma das emoções que não seja neutralidade.

Analisando as emoções felicidade e surpresa percebemos que as três arquiteturas obtiveram resultados relativamente aproximados. Lembrando que são duas das três emoções que mais tem amostragem e possuem padrões menos variável como características próprias para diferenciá-las. Entretanto, as outras emoções, com exceção da emoção medo, a ResNet-34 obteve bastante superioridade no desempenho, no qual teoricamente são emoções mais difíceis de aprender, não apenas por ter menos amostragem, mas também por ter padrões menos definidos e variáveis. A emoção medo foi praticamente empatada com resultado ruim em todas as arquiteturas.

As Tabelas 20 e 21 apresentam os resultados das redes neurais de convolução avaliando as bases de validação: *CK* e *FER*. Essas duas bases foram escolhidas por serem as mais utilizadas na literatura e por uma possuir característica de imagens capturadas no laboratório (*CK*) e a outra na natureza (*FER*). No Anexo C pode ser consultado os resultados para as demais bases: *CIFE-Train*, *CIFE-Test*, *JAFFE*, *KDEF*, *NovaEmotions* e *RAFD*.

A base *CK* teve ótimos resultados nas três arquiteturas (ver Tabela 20). No entanto, a ResNet-34 foi a melhor chegando a 96.9% de acurácia. O que acarretou este resultado expressivo nas três arquiteturas foi o fato desta base ser laboratorial. Neste caso, configura um cenário perfeito onde as imagens não tem ruídos, pouca variabilidade, distância próxima e fixa da câmera, iluminação ideal, pouca angulação da face, os atores emitindo as emoções com alta intensidade e rotulação correta da base.

Os trabalhos de Chen et al. (2017b), Yu et al. (2016c) e Huang and Lu (2016) obtiveram 99.1%, 98.7% e 97.3%, respectivamente, contudo vale destacar que estes trabalhos treinaram e testaram somente com a base *CK*, que tende a ter melhor desempenho de fato, pois se especializaram em reconhecer imagens dessa base e laboratoriais. Outro ponto relevante é que não fica claro se esses resultados são referentes a base de teste ou validação. Nosso trabalho alcançou 96.9% na base de validação, ou seja, próximo dos resultados acima, com a combinação de várias bases inclusive de características de imagens na natureza. Além disso, nesta base não há amostragem da emoção neutralidade. Entretanto em nosso trabalho foi incluída esta emoção. Isto acarretou o aumento da chance de existir uma confusão entre a neutralidade e qualquer outra emoção da base, enquanto os outros trabalhos não apoiaram o reconhecimento para tal. Contudo, não houve qualquer confusão em atribuir como neutralidade qualquer imagem oriunda da base *CK*, o que está correto, pois nessa base não há neutralidade. Podemos concluir que imagens fotografadas na natureza para treinamento não influenciam negativamente o desempenho do classificador atuando em imagens laboratoriais estilo *CK*.

O resultado da base *FER* não foi tão bom (consultar Tabela 21). Analisando a acurácia a ResNet-34 alcançou 60.4%, enquanto a Alexnet 54.3% e a Inception-V3 52.9%. Percebemos que em ambos os casos há uma diferença de desempenho acima de 6% entre a ResNet-34 com as outras arquiteturas. Isto supõe que a base *FER* é uma base difícil para ser aprendida e classificada. O que não é surpresa, pois essa base é composta por imagens oriundas da natureza, no qual não há um padrão de iluminação, variação do fundo do ambiente, distância fixa com câmera, faces com fortes movimentações e emoções com grau de intensidades diferentes. É válido ressaltar que o fato desta base possuir emoções com variação na intensidade pode acarretar imagens com baixa emoção, caracterizando que essas instâncias se encontrem em uma região de confusão, onde não é tão evidente distinguir qual é a emoção havendo sobreposição. Por isso, esta base possui algumas

imagens com rótulos questionáveis. Tanto é que pesquisadores da Microsoft resolveram usar as plataformas da empresa para re-rotular a base FER gerando uma nova base denominada FER+ (Barsoum et al., 2016a).

Comparando os resultados com a comunidade científica, o trabalho de Huang and Lu (2016) alcançou 76.9%, Kim et al. (2016a) conseguiu 73.73% e Liu et al. (2016b) atingiu 65%, todos estes trabalhos utilizaram a base FER para treino e teste. Obviamente por treinarem e testarem somente nesta base os modelos gerados especializaram-se nas características dessa base. Um dos motivos para nosso trabalho ter ficado tão atrás foi que as imagens laboratoriais pode influenciar negativamente no ajuste de pesos para classificar imagens na natureza. Outro fator prejudicial foi a ausência de técnicas como normalização do brilho e contraste, alinhamento da face e aumento de dados. Estas técnicas fazem parte da solução proposta, entretanto ainda precisam ser implementadas. A ideia é tornar a abordagem mais robusta para lidar com cenários de uso reais, e consequentemente, com imagens advindas da natureza.

Portanto, o estudo comparativo mostrou que a ResNet-34 foi a melhor arquitetura. Verificamos que houve praticamente um empate entre a Alexnet e a Inception-V3, enquanto a ResNet-34 sempre esteve a frente em desempenho com uma diferença considerável. Vale destacar que na base CK, uma base de origem laboratorial, nosso trabalho obteve taxas próximas da comunidade científica. Entretanto, na base FER, uma base oriunda da natureza, ficamos distantes dos resultados da literatura. Porém, é de se considerar que a nossa abordagem ainda não está totalmente implementada. Esperamos que o desenvolvimento dos componentes faltantes tornem a nossa solução mais robusta para casos de usos complexos na natureza.

### 5.3 Resumo

Os resultados parciais foram iniciados com uma prova de conceito na área de educação. Esta prova de conceito foi fundamental para a evolução desta proposta. Há uma grande dificuldade para aplicar o método em cenários reais. Realizamos uma análise das emoções dos estudantes com seu desempenho em um teste escolar. As emoções foram reconhecidas usando um serviço em nuvem. O método proposto inclui o reconhecimento de emoções por meio de câmeras, uma coleta não intrusiva, e que o método pode ser útil para gerar sistemas inteligentes mais robustos. Entretanto, um reconhecedor de emoção está em desenvolvimento para não precisar mais usar reconhecedores da nuvem. O primeiro passo foi um estudo comparativo entre as arquiteturas Alexnet, Inception-V3 e ResNet. Concluímos que a ResNet obteve os melhores resultados. Em imagens laboratoriais, a ResNet foi muito bem, aproximado dos baselines. Há uma série de desafios na classificação de imagens da natureza, no qual a ResNet ainda não obteve os resultados desejados.

---

Porém, a abordagem proposta ainda não está totalmente desenvolvida. O objetivo é ter um sistema robusto e preciso para lidar com situações complexas de classificação.



## 6 Considerações Finais

Neste trabalho, foi apresentado uma abordagem para reconhecer emoção por meio da expressão facial utilizando redes neurais de convolução. O diferencial desta abordagem é reunir os principais elementos identificados na literatura para reconhecer emoções. Além disso, este trabalho pretende fornecer uma solução para reconhecer emoções em computação embarcada e em nuvem.

Uma prova de conceito foi aplicada para verificar a viabilidade do reconhecimento de emoção por expressão facial. Nesta prova de conceito, um serviço em nuvem foi utilizado para reconhecer as emoções. O foco foi monitorar estudantes por meio da captura de imagens em tempo real, enquanto faziam um teste escolar com questões de múltipla escolha. Foram correlacionadas as emoções com o desempenho no teste. Portanto, como resultado principal identificamos que as questões onde ocorreram diferentes emoções, foram as questões que tiveram maiores proporções de acertos.

Adicionalmente, para gerar nosso próprio reconhecedor de emoção, um estudo experimental foi conduzido com intuito de avaliar três arquiteturas de redes neurais de convolução: AlexNet, Inception-V3 e ResNet-34. Os resultados apontaram que a ResNet-34 obteve as melhores taxas de acurácia, precisão, revocação e f1-score. As imagens utilizadas no estudo foram provenientes da natureza e laboratório. Concluímos que um método treinado pelas imagens da natureza não influenciam negativamente para reconhecer imagens laboratoriais. Em contrapartida, um método treinado com as imagens laboratoriais não consegue reconhecer emoções no contexto da natureza.

### 6.1 Limitações do Trabalho

Este trabalho tem como objetivo reconhecer emoções humanas por meio da expressão facial. Todavia, os trabalhos de [Darwin \(1965\)](#) e [Ekman and Davidson \(1994\)](#) apontaram que somente o grupo das emoções básicas (raiva, alegria, tristeza, desgosto, medo e surpresa) são emitidas por meio da expressão facial. Portanto, este trabalho tem possibilidades de reconhecer exclusivamente as emoções básicas.

### 6.2 Trabalhos Futuros

- **Analisar sequência de imagens:** Atualmente, a abordagem analisa somente uma imagem para reconhecer as emoções. Desta forma, está sendo ignorada a característica temporal e apenas uma imagem é considerada para definir as emoções. Integrar

uma sequência de imagens, na forma de uma série temporal, pode ser relevante para melhorar o reconhecimento, justamente por analisar um conjunto de imagens amostrada em uma fração de tempo para definir as emoções.

- **Avaliar experimentalmente outros classificadores:** O *softmax* foi o único classificador adotado no estudo realizado. É interessante avaliar o desempenho de outros classificadores analisando diferentes famílias: SVM (funcional), Gaussiana (probabilística) e RandomForest (ensemble de árvores de decisões).
- **Implementar e avaliar a MobileNet:** A MobileNet é a arquitetura de rede neural de convolução apropriada para sistemas embarcados. É necessário um estudo para parametrizar a MobileNet, mensurar o consumo de recursos computacionais e, além disso, definir as especificações mínimas de hardware para a instalação e funcionamento.
- **Desenvolver e avaliar o componente de pré-processamento:** As técnicas de alinhamento de face, normalização de iluminação e aumento de dados, devem ser implementadas. É interessante um estudo para avaliar a contribuição e o impacto dessas técnicas. Principalmente, no cenário das imagens oriundas da natureza, em que nesse contexto, o método proposto não obteve bons resultados.
- **Avaliar em cenários de uso reais:** Nos seguintes cenários pretendemos aplicar a solução: na educação ou em tecnologia assistiva para deficientes visuais. Na educação é útil para auxiliar sistemas inteligentes que interagem com aluno de acordo com a emoção do mesmo. Na tecnologia assistiva pode auxiliar deficientes visuais que tem dificuldades em reconhecer emoção, inclusive este tipo de aplicação abre um campo desafiador para a interação humano computador.

## 6.3 Cronograma

O cronograma está descrito na Tabela 11.

Tabela 11 – Cronograma de Atividades

Atividades	2018					2019		
	ago	set	out	nov	dez	jan	fev	mar
Desenvolver e avaliar o componente pré-processamento	x							
Analisar sequência de imagens		x						
Avaliar experimentalmente outros classificadores			x					
Implementar e avaliar a MobileNet				x				
Avaliar em cenários de uso reais					x	x		
Escrita da dissertação				x	x	x	x	x



# Referências

- Alpaydin, E. (2014). *Introduction to machine learning*. MIT press.
- Barsoum, E., Zhang, C., Ferrer, C. C., and Zhang, Z. (2016a). Training deep networks for facial expression recognition with crowd-sourced label distribution. In *Proceedings of the 18th ACM International Conference on Multimodal Interaction*, pages 279–283. ACM.
- Barsoum, E., Zhang, C., Ferrer, C. C., and Zhang, Z. (2016b). Training deep networks for facial expression recognition with crowd-sourced label distribution. In *Proceedings of the 18th ACM International Conference on Multimodal Interaction*, pages 279–283. ACM.
- Basili, V. R., Caldiera, G., and Rombach, H. D. (1994). Experience factory. *Encyclopedia of software engineering*.
- Biolchini, J., Mian, P. G., Natali, A. C. C., and Travassos, G. H. (2005). Systematic review in software engineering. *System Engineering and Computer Science Department COPPE/UFRJ, Technical Report ES*, 679(05):45.
- Chen, X., Yang, X., Wang, M., and Zou, J. (2017a). Convolution neural network for automatic facial expression recognition. In *Applied System Innovation (ICASI), 2017 International Conference on*, pages 814–817. IEEE.
- Chen, X., Yang, X., Wang, M., and Zou, J. (2017b). Convolution neural network for automatic facial expression recognition. In *Applied System Innovation (ICASI), 2017 International Conference on*, pages 814–817. IEEE.
- Cruz, A., Leitão, G., Colonna, J., Silva, E., Barreto, R., and Primo, T. (2017). Framework para coleta e inferência de estados emocionais de alunos baseado em reconhecimento de expressões faciais. In *Brazilian Symposium on Computers in Education (Simpósio Brasileiro de Informática na Educação-SBIE)*, volume 28, page 997.
- Darwin, C. (1965). *The expression of the emotions in man and animals*, volume 526. University of Chicago press.
- Ekman, P. and Friesen, W. V. (1977). Facial action coding system.
- Ekman, P. E. and Davidson, R. J. (1994). *The nature of emotion: Fundamental questions*. Oxford University Press.

- Géron, A. (2017). Hands-on machine learning with scikit-learn and tensorflow: concepts, tools, and techniques to build intelligent systems.
- Goodfellow, I., Bengio, Y., and Courville, A. (2016). *Deep learning*. MIT press.
- Guo, Y., Tao, D., Yu, J., Xiong, H., Li, Y., and Tao, D. (2016). Deep neural networks with relativity learning for facial expression recognition. In *Multimedia & Expo Workshops (ICMEW), 2016 IEEE International Conference on*, pages 1–6. IEEE.
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.
- Huang, Y. and Lu, H. (2016). Deep learning driven hypergraph representation for image-based emotion recognition. In *Proceedings of the 18th ACM International Conference on Multimodal Interaction*, pages 243–247. ACM.
- Hubel, D. H. (1959). Single unit activity in striate cortex of unrestrained cats. *The Journal of physiology*, 147(2):226–238.
- Hubel, D. H. and Wiesel, T. N. (1959). Receptive fields of single neurones in the cat’s striate cortex. *The Journal of physiology*, 148(3):574–591.
- Hubel, D. H. and Wiesel, T. N. (1968). Receptive fields and functional architecture of monkey striate cortex. *The Journal of physiology*, 195(1):215–243.
- Jain, A. K. (1989). *Fundamentals of digital image processing*. Prentice-Hall, Inc.
- Jaques, P. A. and Nunes, M. A. S. (2013). Ambientes inteligentes de aprendizagem que inferem, expressam e possuem emoções e personalidade. *Jornada de Atualização em Informática na Educação*, 1(1):30–81.
- Jung, H., Lee, S., Park, S., Kim, B., Kim, J., Lee, I., and Ahn, C. (2015a). Development of deep learning-based facial expression recognition system. In *Frontiers of Computer Vision (FCV), 2015 21st Korea-Japan Joint Workshop on*, pages 1–4. IEEE.
- Jung, H., Lee, S., Park, S., Kim, B., Kim, J., Lee, I., and Ahn, C. (2015b). Development of deep learning-based facial expression recognition system. In *Frontiers of Computer Vision (FCV), 2015 21st Korea-Japan Joint Workshop on*, pages 1–4. IEEE.
- Kim, B.-K., Dong, S.-Y., Roh, J., Kim, G., and Lee, S.-Y. (2016a). Fusing aligned and non-aligned face information for automatic affect recognition in the wild: A deep learning approach. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 48–57.

- Kim, B.-K., Dong, S.-Y., Roh, J., Kim, G., and Lee, S.-Y. (2016b). Fusing aligned and non-aligned face information for automatic affect recognition in the wild: A deep learning approach. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 48–57.
- Kitchenham, B. (2004). Procedures for performing systematic reviews. *Keele, UK, Keele University*, 33(2004):1–26.
- Kotsiantis, S. B., Zaharakis, I., and Pintelas, P. (2007). Supervised machine learning: A review of classification techniques.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324.
- Li, W., Li, M., Su, Z., and Zhu, Z. (2015a). A deep-learning approach to facial expression recognition with candid images. In *Machine Vision Applications (MVA), 2015 14th IAPR International Conference on*, pages 279–282. IEEE.
- Li, W., Li, M., Su, Z., and Zhu, Z. (2015b). A deep-learning approach to facial expression recognition with candid images. In *Machine Vision Applications (MVA), 2015 14th IAPR International Conference on*, pages 279–282. IEEE.
- Liu, K., Zhang, M., and Pan, Z. (2016a). Facial expression recognition with cnn ensemble. In *Cyberworlds (CW), 2016 International Conference on*, pages 163–166. IEEE.
- Liu, K., Zhang, M., and Pan, Z. (2016b). Facial expression recognition with cnn ensemble. In *Cyberworlds (CW), 2016 International Conference on*, pages 163–166. IEEE.
- Mafrá, S. N. and Travassos, G. H. (2006). Estudos primários e secundários apoiando a busca por evidência em engenharia de software. *Relatório Técnico, RT-ES*, 687(06).
- Mayya, V., Pai, R. M., and Pai, M. M. (2016). Automatic facial expression recognition using dcnn. *Procedia Computer Science*, 93:453–461.
- Nasoz, F., Alvarez, K., Lisetti, C. L., and Finkelstein, N. (2004). Emotion recognition from physiological signals using wireless sensors for presence technologies. *Cognition, Technology & Work*, 6(1):4–14.
- Ng, H.-W., Nguyen, V. D., Vonikakis, V., and Winkler, S. (2015). Deep learning for emotion recognition on small datasets using transfer learning. In *Proceedings of the 2015 ACM on international conference on multimodal interaction*, pages 443–449. ACM.
- Riva, G., Mantovani, F., Capideville, C. S., Preziosa, A., Morganti, F., Villani, D., Gaggioli, A., Botella, C., and Alcañiz, M. (2007). Affective interactions using virtual reality: the link between presence and emotions. *CyberPsychology & Behavior*, 10(1):45–56.

- Scherer, K. R. (2000). Psychological models of emotion. *The neuropsychology of emotion*, 137(3):137–162.
- Shan, K., Guo, J., You, W., Lu, D., and Bie, R. (2017). Automatic facial expression recognition based on a deep convolutional-neural-network structure. In *Software Engineering Research, Management and Applications (SERA), 2017 IEEE 15th International Conference on*, pages 123–128. IEEE.
- Shin, M., Kim, M., and Kwon, D.-S. (2016a). Baseline cnn structure analysis for facial expression recognition. In *Robot and Human Interactive Communication (RO-MAN), 2016 25th IEEE International Symposium on*, pages 724–729. IEEE.
- Shin, M., Kim, M., and Kwon, D.-S. (2016b). Baseline cnn structure analysis for facial expression recognition. In *Robot and Human Interactive Communication (RO-MAN), 2016 25th IEEE International Symposium on*, pages 724–729. IEEE.
- Simonyan, K. and Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.
- Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., and Rabinovich, A. (2015). Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9.
- Viola, P. and Jones, M. (2001). Rapid object detection using a boosted cascade of simple features. In *Computer Vision and Pattern Recognition, 2001. CVPR 2001. Proceedings of the 2001 IEEE Computer Society Conference on*, volume 1, pages I–I. IEEE.
- Vo, D. M. and Le, T. H. (2016). Deep generic features and svm for facial expression recognition. In *Information and Computer Science (NICS), 2016 3rd National Foundation for Science and Technology Development Conference on*, pages 80–84. IEEE.
- Wen, G., Hou, Z., Li, H., Li, D., Jiang, L., and Xun, E. (2017a). Ensemble of deep neural networks with probability-based fusion for facial expression recognition. *Cognitive Computation*, pages 1–14.
- Wen, G., Hou, Z., Li, H., Li, D., Jiang, L., and Xun, E. (2017b). Ensemble of deep neural networks with probability-based fusion for facial expression recognition. *Cognitive Computation*, pages 1–14.
- Yu, X., Huang, J., Zhang, S., and Metaxas, D. N. (2016a). Face landmark fitting via optimized part mixtures and cascaded deformable model. *IEEE transactions on pattern analysis and machine intelligence*, 38(11):2212–2226.

- Yu, X., Lin, Z., Brandt, J., and Metaxas, D. N. (2014). Consensus of regression for occlusion-robust facial feature localization. In *European Conference on Computer Vision*, pages 105–118. Springer.
- Yu, X., Yang, J., Luo, L., Li, W., Brandt, J., and Metaxas, D. (2016b). Customized expression recognition for performance-driven cutout character animation. In *Applications of Computer Vision (WACV), 2016 IEEE Winter Conference on*, pages 1–9. IEEE.
- Yu, X., Yang, J., Luo, L., Li, W., Brandt, J., and Metaxas, D. (2016c). Customized expression recognition for performance-driven cutout character animation. In *Applications of Computer Vision (WACV), 2016 IEEE Winter Conference on*, pages 1–9. IEEE.



## Anexos





# ANEXO A – Revisão Sistemática da Literatura

Neste anexo é descrito e discutido uma Revisão Sistemática da Literatura acerca do tema deste trabalho. Na Seção A.1 é descrito o protocolo seguido para a realização da revisão sistemática. Na Seção A.2 está o processo de condução da revisão sistemática e quantos artigos veio em cada filtro. Na Seção A.3 contém os resultados obtidos, assim como, as respostas para as questões de pesquisa e por fim na Seção A.4 o resumo e outras discussões sobre esta revisão sistemática da literatura.

## A.1 Protocolo da Revisão Sistemática da Literatura

Este protocolo foi elaborado conforme especificado em: [Biolchini et al. \(2005\)](#), [Mafra and Travassos \(2006\)](#), e [Kitchenham \(2004\)](#):

### A.1.1 Objetivo

O objetivo deste estudo será esquematizado a partir do paradigma GQM (goal, question, and metric) ([Basili et al., 1994](#)):

### A.1.2 Questões de Pesquisa

**Questão Principal:** Como reconhecer emoções por meio da expressão facial utilizando redes neurais de convolução em uma imagem estática?

- **Q1:** Quais emoções têm sido reconhecidas por meio da expressão facial utilizando redes neurais de convolução?
- **Q2:** Quais tipos de pré-processamento tem sido realizado na imagem?
- **Q3:** Quais arquiteturas de redes de convolução têm sido mais utilizadas?

Tabela 12 – Objetivos da Revisão Sistemática

<b>Analisar</b>	Reconhecimento de emoções por meio da expressão facial em uma imagem estática.
<b>Com o propósito de</b>	Identificar técnicas, métodos, abordagens, arquiteturas, base de dados e aplicações.
<b>No que diz respeito a</b>	Utilização de redes neurais de convolução.
<b>Do ponto de vista do</b>	Pesquisador.
<b>No contexto</b>	Acadêmico.

- **Q4:** Quais técnicas, métodos e abordagens têm sido utilizados para tratar problemas na imagem como iluminação, rotação, obstrução e escala?
- **Q5:** Quais bases de dados têm sido utilizadas?
- **Q6:** Quais aplicações podem utilizar o reconhecimento de emoção por expressão facial?

### A.1.3 Biblioteca Digital

Scopus: <http://www.scopus.com/> - Contempla as principais conferências da área (foi verificado por meio de busca manual)

### A.1.4 Critérios de Inclusão e Exclusão dos Artigos

#### Critérios de Inclusão:

- **CI1:** Reconhecimento de emoção por expressão facial usando somente CNN com abordagem que funciona para classificação em imagem;
- **CI2:** Reconhecimento de emoção por expressão facial combinando CNN com várias arquiteturas de redes neurais com abordagem que funciona para classificação em imagem;
- **CI3:** Reconhecimento de emoção por expressão facial combinando CNN com outros métodos de aprendizado de máquina que funciona para classificação em imagem;
- **CI4:** Reconhecimento de emoção por expressão facial combinando CNN com técnicas de pré-processamento que não são originais da arquitetura CNN com abordagem que funciona para classificação em imagem.

#### Critérios de Exclusão:

- **CE1:** Trabalho somente apresenta teoria ou discussão relacionada ao reconhecimento de emoções;
- **CE2:** Não apresenta reconhecimento de emoção por expressão facial para classificação em imagens;
- **CE3:** Não utiliza redes neurais de convolução;
- **CE4:** Trabalho anterior ao ano de 2013;
- **CE5:** Reconhecimento por vídeo ou *streaming* de imagens;

- **CE6:** Trabalho utiliza durante a metodologia experimental uma base de dados não disponível para a comunidade científica;
- **CE7:** Publicação não disponível;
- **CE8:** Reconhecimento de emoção multimodal.

**Observação:** O *CE4* foi definido devido o surgimento das (atuais) redes neurais de convolução ter sido a partir de 2013.

### A.1.5 Formulário de Extração de Informação

Inicialmente, no primeiro filtro serão analisados e considerados os seguintes itens:

- Título;
- Resumo;
- Palavras-chaves.

Posteriormente, no segundo filtro serão extraídas as seguintes informações:

- Autores do trabalho;
- Fonte: local que o trabalho foi publicado;
- Ano de Publicação;
- Emoções que foram reconhecidas;
- Aplicações para o reconhecimento de emoções por expressão facial;
- Arquiteturas de redes neurais de convolução utilizadas;
- Metodologia utilizada para o treinamento da rede neural de convolução;
- Base de dados utilizadas para treino e validação;
- Perspectivas futuras;
- Comentários.

### A.1.6 *String* de Busca

As *strings* de busca foram definidas a partir das questões de pesquisa e do padrão PICO (*population, intervention, comparison, outcomes*) (KITCHENHAM e CHARTERS, 2007), conforme a estrutura abaixo:

- **População:** Reconhecimento de emoção por expressão facial;
- **Intervenção:** Por meio de redes neurais de convolução;
- **Comparação:** Não há;
- **Resultados:** Técnicas, métodos, arquiteturas, base de dados, aplicações e abordagens.

("emotion recognition"OR "emotion detection"OR "emotion identification"OR "emotion analysis"OR "emotion classification"OR "affect recognition"OR "affect detection"OR "affect analysis"OR "affect classification"OR "facial expression")

AND

("convolutional neural network"OR "CNN"OR "long short term memory"OR "LSTM"OR "recurrent neural network"OR "RNN")

AND

("technique"OR "method"OR "architecture"OR "database"OR "application"OR "approach")

Para a montagem da string de busca foi testado cada termo da população com todos os termos da intervenção mais resultado, desta forma, verificando e validando que todos os termos da população realmente contribuem para os artigos retornados. Este mesmo procedimento foi utilizado para verificar e validar os termos do resultado, no entanto foi verificado cada termo do resultado com todos os termos da intervenção e população.

A presença dos seguintes termos na intervenção: "long short term memory", "LSTM", "recurrent neural network" e "RNN", podem ser explicados devido à intuição do autor deste protocolo acreditar que a comunidade estava utilizando essas técnicas, que também são redes neurais profundas, combinadas com as redes neurais de convolução para classificação de expressões faciais em imagens.

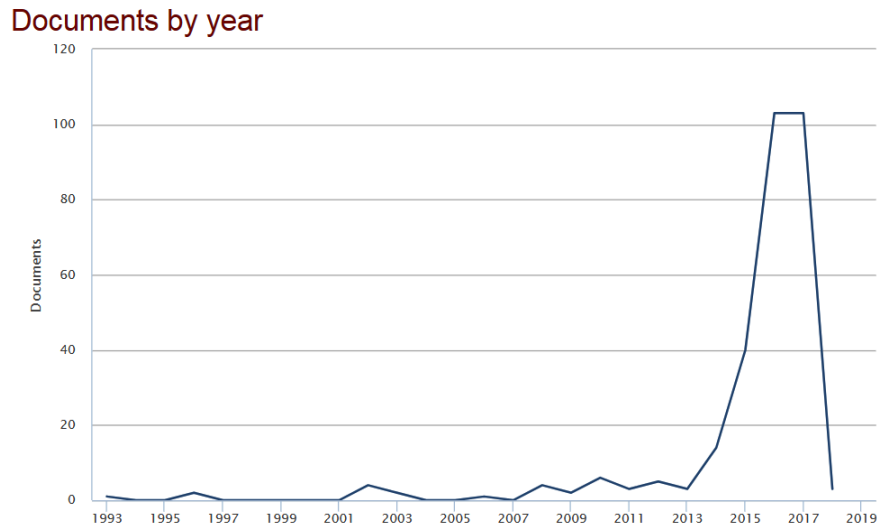


Figura 21 – Artigos por ano retornados pela *string* de busca

## A.2 Condução da Revisão Sistemática da Literatura

### A.2.1 Primeiro Filtro

No primeiro filtro foi lido somente o título, resumo e as palavras-chaves do artigo. A *string* de busca retornou 281 artigos para classificar no primeiro filtro. Foram aceitos 99 (35%) para o segundo filtro, 3 (1%) duplicados e 179 (64%) rejeitados.

Na Seção A.1.6, o autor do protocolo esclarece porque utilizou os seguintes termos na *string* de busca: "long short term memory", "LSTM", "recurrent neural network" e "RNN", depois do primeiro filtro, realmente comprovou-se que estas técnicas são combinadas com as redes neurais de convolução para classificação de emoção em expressão facial, porém, somente em vídeos ou streaming de imagens. Portanto, os artigos retornados por essas palavras receberam a classificação de rejeitado devido a esta revisão focar em trabalhos com classificação em imagens estática sem streaming.

### A.2.2 Segundo Filtro

Para a realização do segundo filtro foi lido o artigo completo para a extração dos dados e, conseqüentemente a obtenção dos resultados. No segundo filtro tinham 99 artigos para classificar, onde 34 foram aceitos, 1 duplicados e 64 rejeitados.

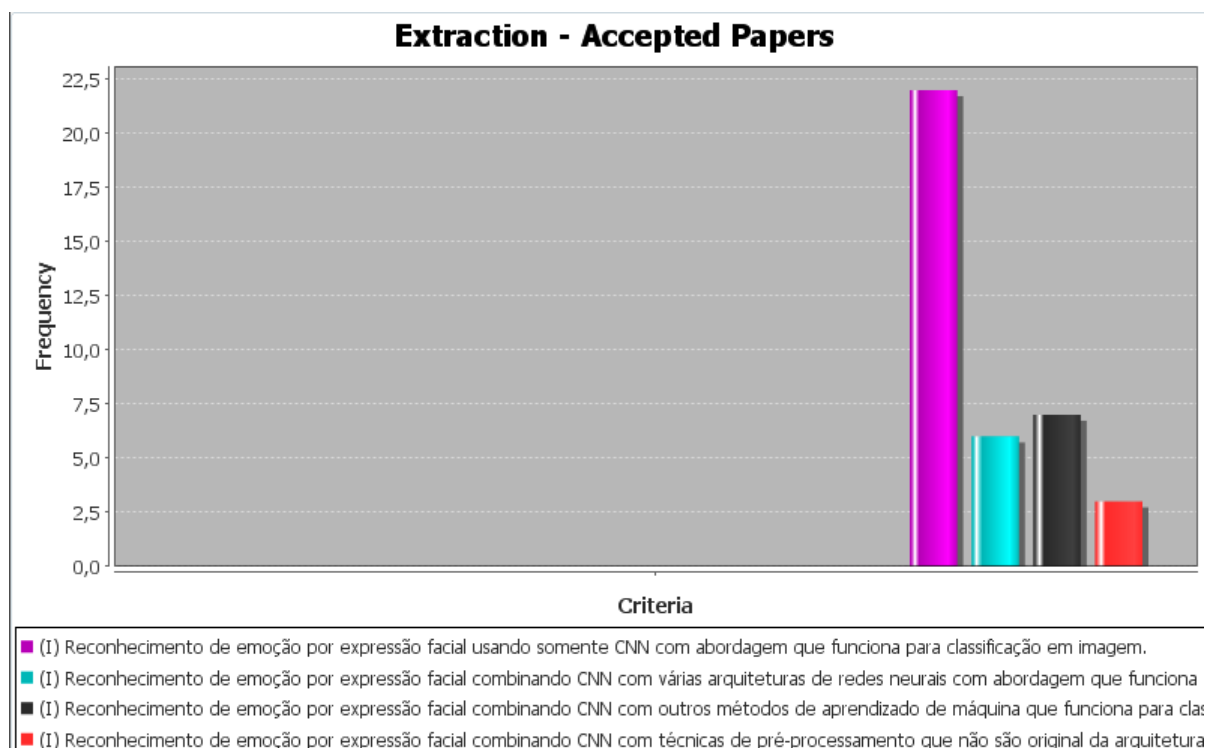


Figura 22 – Gráfico representando a frequência dos critérios de inclusão para os artigos aceitos do segundo filtro

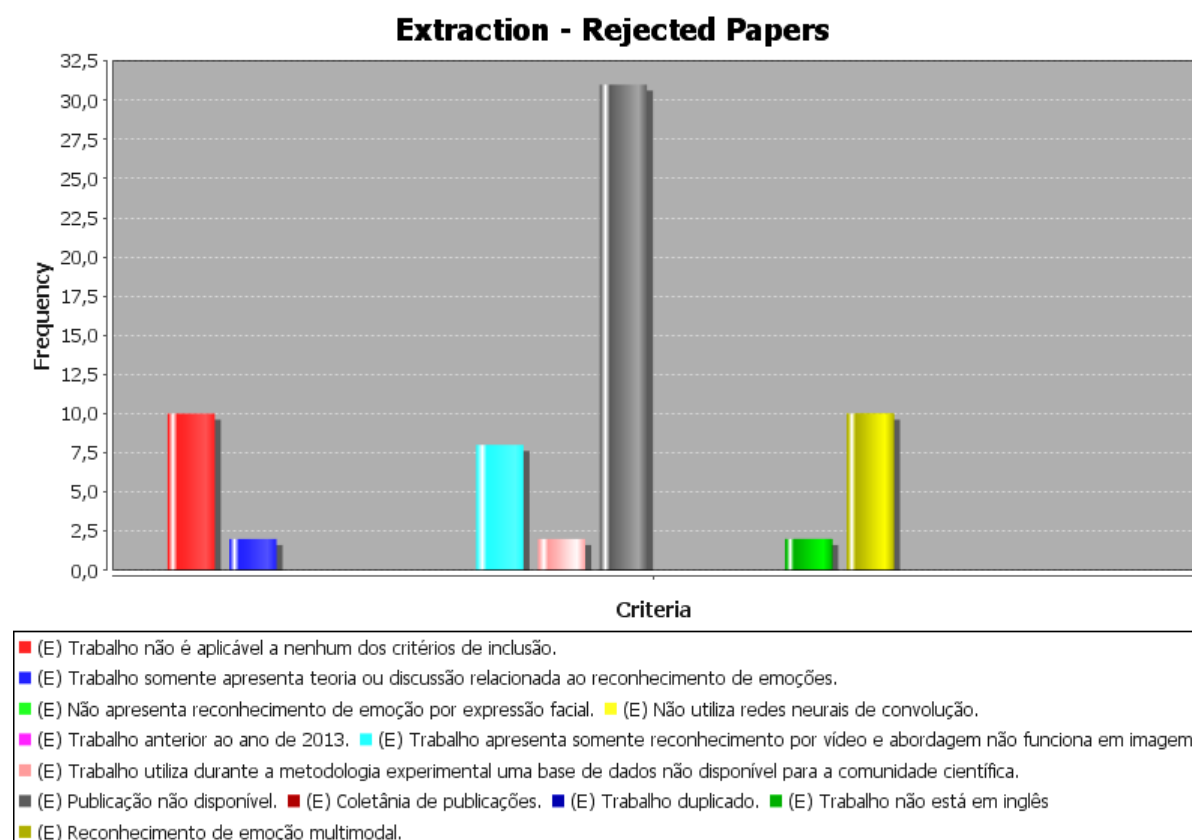


Figura 23 – Gráfico representando a frequência dos critérios de exclusão para os artigos rejeitados do segundo filtro

## A.3 Resultados

### A.3.1 Q1: Quais emoções têm sido reconhecidas por meio da expressão facial utilizando redes neurais de convolução?

As emoções reconhecidas obviamente são as emoções que são representadas e contidas em base de dados, logo são as emoções básicas: neutralidade, felicidade, surpresa, tristeza, raiva, desgosto e medo. Alguns trabalhos como Barsoum et al. (2016b) reconhece também nojo.

### A.3.2 Q2: Quais tipos de pré-processamento tem sido realizado na imagem?

Abaixo segue uma lista com o pré-processamento e os trabalhos que fizeram utilização da técnica.

- **Detector de Face:** Consiste na detecção e recorte da face (Chen et al., 2017b; Li et al., 2015b; Mayya et al., 2016; Ng et al., 2015; Shan et al., 2017; Shin et al., 2016b; Vo and Le, 2016);
- **Normalização de Brilho (Equalização de histograma):** Transformada para realce do contraste (Kim et al., 2016b; Shan et al., 2017; Shin et al., 2016b);
- **Normalização Min e Max:** Transformação linear baseada no valor mínimo e máximo da imagem (Kim et al., 2016b);
- **Pontos da Face (pontos geométricos):** Extração de pontos da face e a distância entre os pontos (Yu et al., 2016c);
- **Escala de Cinza:** Transformação da imagem para escala de cinza (Mayya et al., 2016);
- **Diferença Gaussiana:** Detecta as bordas do objeto, neste caso, evidencia as bordas da face (Shin et al., 2016b);
- **Filtro de Difusão Isotópica:** (Shin et al., 2016b);
- **Normalização DCT:** Transformada discreta do cosseno utilizada em compressão de dados e eventualmente evidenciando informações relevantes da imagem (Shin et al., 2016b);
- **Alinhamento da Face:** Utilização de uma rede neural *autoencoder* para alinhar a face no centro (Kim et al., 2016b).

Arquitetura	Trabalhos que utilizaram a arquitetura
AlexNet	Chen et al. (2017b), Shan et al. (2017), Kim et al. (2016b), Huang and Lu (2016), Vo and Le (2016), Yu et al. (2016c), Ng et al. (2015), Jung et al. (2015b), Li et al. (2015b)
GoogLeNet	Guo et al. (2016)
VGG	Barsoum et al. (2016b), Ng et al. (2015)
Ensemble	Wen et al. (2017b), Liu et al. (2016b), Shin et al. (2016b)

Tabela 13 – Principais arquiteturas de redes neurais de convolução e os trabalhos que utilizaram

### A.3.3 Q3: Quais arquiteturas de redes neurais de convolução têm sido mais utilizadas?

A Tabela 13 contém as arquiteturas encontradas e os trabalhos que a utilizaram, como podemos verificar a arquitetura AlexNet foi a mais utilizada.

### A.3.4 Q4: Quais técnicas, métodos e abordagens têm sido utilizados para tratar problemas na imagem como iluminação, rotação, obstrução e escala?

A principal solução encontrada na literatura foi a ênfase na generalização adequada do aprendizado da rede neural de convolução, isto é, durante a fase de treinamento. Obviamente, as técnicas aplicadas no pré-processamento da imagem (ver Seção A.3.2), contribuem para resolver problemas de iluminação por meio da equalização de histograma e de rotação com o alinhamento de face. Entretanto, um achado bastante interessante foi a utilização da técnica de aumento de dados, que consiste durante a fase de treinamento da rede neural de convolução em multiplicar por 10 vezes uma instância (imagem), isto é, gerando 10 novas imagens com pequenos giros da faces, e variações da rotação da pose, escala e iluminação. Acrescendo em 10 vezes o tamanho da base de treinamento, com inserção de variações da imagem resultando em melhor aprendizado da rede.

### A.3.5 Q5: Quais bases de dados têm sido utilizadas?

Esta seção tem enfoque nas bases de dados mapeadas para reconhecimento de emoção por expressão facial em uma imagem estática. Obviamente, é possível encontrar outras base de dados para reconhecimento de emoção que não seja por imagem estática, por exemplo, reconhecimento em vídeo, por sensores, em textos e outras.

As bases de dados para reconhecimento de emoção por expressão facial em uma imagem estática tem algo em comum, geralmente as amostras de expressões faciais são as mesmas emoções, as chamadas emoções básicas investigadas por Ekman and Davidson (1994) que são: neutralidade, felicidade, medo, desgosto, raiva, surpresa e tristeza, isto



Bases de Dados	Trabalhos que utilizaram a base para treinamento ou validação
CK+	Chen et al. (2017b), Shan et al. (2017), Wen et al. (2017b), Shin et al. (2016b), Huang and Lu (2016), Vo and Le (2016), Yu et al. (2016c), Mayya et al. (2016), Jung et al. (2015b), Li et al. (2015b)
JAFFE	Chen et al. (2017b), Shan et al. (2017), Wen et al. (2017b), Shin et al. (2016b), Mayya et al. (2016)
FER	Wen et al. (2017b), Kim et al. (2016b), Liu et al. (2016b), Shin et al. (2016b), Huang and Lu (2016), Guo et al. (2016), Ng et al. (2015), Jung et al. (2015b)
FER+	Barsoum et al. (2016b)
SFEW2.0	Shin et al. (2016b), Guo et al. (2016)
KDEF	Shin et al. (2016b)
MMI	Yu et al. (2016c)
CIFE	Li et al. (2015b)
EmotiW2015	Wen et al. (2017b), Ng et al. (2015)

Tabela 14 – Bases de Dados

significa que são essas as emoções que a comunidade tem reconhecido por expressão facial. As bases de dados mapeadas podem ser consultadas na Tabela ??.

### A.3.6 Q6: Quais aplicações podem utilizar o reconhecimento de emoção por expressão facial?

Há diversas aplicações para o reconhecimento de emoção no mundo real, foi percebido que os pesquisadores de reconhecimento de emoção por expressão facial utilizando rede neural de convolução, ultimamente concentraram seus esforços mais no desenvolvimento de reconhecedores de emoção do que a aplicação em cenários reais, mesmo assim, está aberto para trabalhos futuros inúmeras aplicações desses reconhecedores em diversas áreas, tendo destaque principalmente para:

- **Interação humano computador:** no qual pode ser possível projetar interfaces que se adaptam ao estado emocional do usuário (Barsoum et al., 2016b; Chen et al., 2017b; Liu et al., 2016b; Wen et al., 2017b);
- **Psiquiatria e cuidados médicos:** no qual o reconhecedor de emoção deve monitorar constantemente o paciente ou usuário fornecendo dados emocionais que podem contribuir para diagnósticos (Chen et al., 2017b; Mayya et al., 2016; Wen et al., 2017b);
- **Deficiente visual:** pois pessoas com alto grau de deficiência visual, tem dificuldades na interação entre pessoas para identificar qual a emoção que as pessoas em volta estão emitindo (Li et al., 2015b);
- **Interação humano robô:** fazendo com que robôs estejam habilitados a interagir com humanos podendo adaptar-se a emoção dos humanos em volta, ou até mesmo

emitir emoção se aproximando de um humanoide (Jung et al., 2015b; Shin et al., 2016b);

- **Personagens virtuais e animação:** habilitando avatares a copiar expressão humana que podem ser útil para gravações de filmes de animação, também pode ser usado em aplicações de animação como o popular aplicativo para *smartphone* o *Snapchat*, que identifica a expressão facial do usuário e retorna alguma animação sobrepondo a expressão anteriormente detectada do usuário (Vo and Le, 2016; Yu et al., 2016c).

### A.3.7 Questão Principal: Como reconhecer emoções por meio da expressão facial utilizando redes neurais de convolução em uma imagem estática?

Diante das fichas de extração, foi percebido que a comunidade explorou diversas estratégias para processar imagens de expressão facial e reconhecer emoção. Algumas abordagens se destacam como: a técnica para aumentar os dados de treinamento e teste, utilizando a técnica flip fazendo até 10 pequenas rotações na imagem, para a CNN aprender a generalizar melhor sendo treinada e testada com uma base de dados maior. A técnica de normalização de brilho (equalização) no pré-processamento, no qual todos os trabalhos que utilizaram esta técnica aumentaram a acurácia do reconhecimento. Também merece destaque o trabalho de Kim et al. (2016a) em que na sua abordagem, a rede de convolução recebe duas expressões faciais de entrada: a saída de um autoencoder que alinha a face e a imagem original sem alinhamento da face. Esta abordagem melhorou bastante o reconhecimento.

Com relação à arquitetura da rede neural de convolução, quem utilizou um SVM como classificador ao invés de um tradicional softmax obteve maior acurácia. Teve trabalhos que utilizou uma rede com camadas inceptions, hipergrafo, ensembles e concatenação de redes, e todas essas abordagens superaram uma CNN simples. Neste caso, falta um trabalho que possa dizer experimentalmente qual dessas arquiteturas é a melhor.

Percebemos que existem várias bases de dados disponíveis para a comunidade. As bases de dados que foram mais exploradas foram a CK+, FER2013 e a JAFFE. A base CK+ é composta por expressões faciais capturadas em laboratório, por isso, tem altas taxas de reconhecimento, pois, sua dificuldade para o reconhecimento diminui. Já a base FER2013 foi capturada na “natureza”, por isso sua taxa de reconhecimento é mais baixa sendo uma base bastante complexa para classificação.

Notoriamente os trabalhos utilizam o algoritmo Viola Jones para detecção de face pelo programa OpenCV e fazem o recorte da face excluindo o background. Desta forma, elimina o trabalho da rede em aprender a separar o que é background e o que é face, diminuindo a complexidade da classificação.

Portanto, para reconhecer emoção em uma imagem estática, é necessário o treinamento de uma rede de convolução com um classificador na última camada, com a maior quantidade de dados possível, realizando o recorte da face e utilizar técnicas de normalização na imagem, ocasionando um aumento da taxa de reconhecimento.

## A.4 Resumo

Neste anexo apresentou uma revisão sistemática da literatura que investigou o estado-da-arte sobre o reconhecimento de emoção por expressão facial por meio de redes neurais de convolução. Verificamos que o tema está bem quente na comunidade, pois, antes de 2013 a String de busca retornou 33 artigos, e em 2013 (3 artigos), 2014 (14 artigos), 2015 (40 artigos), 2016 (103 artigos), 2017 (103 artigos) e 2018 (3 artigos), isso demonstra o crescimento exponencial da área.

Foram mapeadas as principais técnicas de pré-processamento, arquitetura de rede neural de convolução, base de dados, metodologias de treinamento e aplicações do reconhecimento de emoção. A impressão que fica é que a comunidade ainda não está utilizando esses classificadores no mundo real, e o amadurecimento rápido da área depois do surgimento do aprendizado profundo, nos levar acreditar que esses sistemas já estão prontos para ser posto em prática apoiando outras aplicações de interação humano computador, interação humano robô, educação, segurança, computação afetiva e etc.



## ANEXO B – Resultados por Base de validação Geral

Neste anexo são apresentados os resultados da Arquitetura AlexNet, InceptionV3 e ResNet avaliando a base de validação geral. Vale ressaltar que na primeira coluna está arquitetura e o passo (ou interação) em que foi gerado o modelo.

Tabela 15 – Resultados da Arquitetura AlexNet avaliando a base de validação geral

Arquitetura	Emoção	Precisão	Revocação	F1-score	Acurácia
Alexnet-13750	Raiva	0.56	0.52	0.54	0.708
	Desgosto	0.59	0.63	0.61	
	Medo	0.59	0.27	0.37	
	Felicidade	0.82	0.89	0.85	
	Tristeza	0.75	0.4	0.52	
	Surpresa	0.81	0.8	0.81	
	Neutralidade	0.55	0.74	0.63	
	Média/Total	0.71	0.71	0.7	
Alexnet-13800	Raiva	0.56	0.52	0.54	0.704
	Desgosto	0.69	0.53	0.6	
	Medo	0.43	0.42	0.42	
	Felicidade	0.87	0.86	0.86	
	Tristeza	0.54	0.6	0.57	
	Surpresa	0.88	0.73	0.8	
	Neutralidade	0.57	0.7	0.63	
	Média/Total	0.72	0.7	0.71	
Alexnet-16000	Raiva	0.6	0.49	0.54	0.702
	Desgosto	0.72	0.53	0.61	
	Medo	0.4	0.43	0.42	
	Felicidade	0.84	0.89	0.86	
	Tristeza	0.47	0.68	0.56	
	Surpresa	0.79	0.81	0.8	
	Neutralidade	0.69	0.52	0.59	
	Média/Total	0.71	0.7	0.7	
Alexnet-17800	<b>Raiva</b>	<b>0.51</b>	<b>0.6</b>	<b>0.55</b>	<b>0.712</b>
	<b>Desgosto</b>	<b>0.62</b>	<b>0.64</b>	<b>0.63</b>	
	<b>Medo</b>	<b>0.47</b>	<b>0.41</b>	<b>0.44</b>	
	<b>Felicidade</b>	<b>0.84</b>	<b>0.89</b>	<b>0.86</b>	
	<b>Tristeza</b>	<b>0.64</b>	<b>0.5</b>	<b>0.56</b>	
	<b>Surpresa</b>	<b>0.84</b>	<b>0.77</b>	<b>0.8</b>	
	<b>Neutralidade</b>	<b>0.62</b>	<b>0.64</b>	<b>0.63</b>	
	<b>Média/Total</b>	<b>0.71</b>	<b>0.71</b>	<b>0.71</b>	
Alexnet-18000	Raiva	0.53	0.56	0.54	0.704
	Desgosto	0.76	0.47	0.59	
	Medo	0.43	0.37	0.4	
	Felicidade	0.82	0.89	0.85	
	Tristeza	0.77	0.38	0.51	
	Surpresa	0.74	0.84	0.79	
	Neutralidade	0.58	0.69	0.63	
	Média/Total	0.71	0.7	0.69	

Tabela 16 – Resultados da Arquitetura InceptionV3 avaliando a base de validação geral

Arquitetura	Emoção	Precisão	Revocação	F1-score	Acurácia
InceptionV3-25600	Raiva	0.5	0.51	0.5	0.686
	Desgosto	0.63	0.46	0.53	
	Medo	0.39	0.41	0.4	
	Felicidade	0.86	0.9	0.88	
	Tristeza	0.45	0.44	0.44	
	Surpresa	0.82	0.79	0.8	
	Neutralidade	0.58	0.6	0.59	
	Média/Total	0.69	0.69	0.69	
InceptionV3-26800	Raiva	0.55	0.45	0.49	0.687
	Desgosto	0.57	0.56	0.57	
	Medo	0.41	0.41	0.41	
	Felicidade	0.91	0.85	0.88	
	Tristeza	0.43	0.54	0.48	
	Surpresa	0.85	0.78	0.81	
	Neutralidade	0.56	0.66	0.6	
	Média/Total	0.7	0.69	0.69	
InceptionV3-28400	Raiva	0.55	0.47	0.5	0.694
	Desgosto	0.57	0.54	0.55	
	Medo	0.42	0.4	0.41	
	Felicidade	0.87	0.9	0.88	
	Tristeza	0.46	0.48	0.47	
	Surpresa	0.85	0.78	0.81	
	Neutralidade	0.58	0.63	0.6	
	Média/Total	0.69	0.69	0.69	
InceptionV3-30200	Raiva	0.59	0.43	0.5	0.683
	Desgosto	0.69	0.43	0.53	
	Medo	0.33	0.38	0.35	
	Felicidade	0.89	0.87	0.88	
	Tristeza	0.49	0.44	0.47	
	Surpresa	0.72	0.85	0.78	
	Neutralidade	0.56	0.64	0.6	
	Média/Total	0.69	0.68	0.68	
InceptionV3-34000	<b>Raiva</b>	<b>0.54</b>	<b>0.51</b>	<b>0.52</b>	<b>0.701</b>
	<b>Desgosto</b>	<b>0.56</b>	<b>0.57</b>	<b>0.56</b>	
	<b>Medo</b>	<b>0.47</b>	<b>0.42</b>	<b>0.44</b>	
	<b>Felicidade</b>	<b>0.88</b>	<b>0.88</b>	<b>0.88</b>	
	<b>Tristeza</b>	<b>0.47</b>	<b>0.53</b>	<b>0.5</b>	
	<b>Surpresa</b>	<b>0.85</b>	<b>0.79</b>	<b>0.82</b>	
	<b>Neutralidade</b>	<b>0.59</b>	<b>0.62</b>	<b>0.61</b>	
	<b>Média/Total</b>	<b>0.7</b>	<b>0.7</b>	<b>0.7</b>	

Tabela 17 – Resultados da Arquitetura ResNet avaliando a base de validação geral

Arquitetura	Emoção	Precisão	Revocação	F1-score	Acurácia
ResNet-8600	Raiva	0.64	0.65	0.64	0.753
	Desgosto	0.73	0.63	0.68	
	Medo	0.46	0.52	0.49	
	Felicidade	0.93	0.85	0.89	
	Tristeza	0.73	0.58	0.65	
	Surpresa	0.84	0.83	0.84	
	Neutralidade	0.6	0.79	0.68	
	Média/Total	0.77	0.75	0.76	
ResNet-9600	Raiva	0.37	0.87	0.52	0.727
	Desgosto	0.7	0.65	0.67	
	Medo	0.58	0.47	0.52	
	Felicidade	0.91	0.87	0.89	
	Tristeza	0.71	0.53	0.61	
	Surpresa	0.82	0.82	0.82	
	Neutralidade	0.77	0.55	0.64	
	Média/Total	0.77	0.73	0.74	
ResNet-10625	Raiva	0.8	0.51	0.62	0.75
	Desgosto	0.84	0.54	0.66	
	Medo	0.63	0.34	0.44	
	Felicidade	0.89	0.9	0.89	
	Tristeza	0.51	0.78	0.61	
	Surpresa	0.75	0.88	0.81	
	Neutralidade	0.69	0.66	0.68	
	Média/Total	0.76	0.75	0.74	
ResNet-11250	<b>Raiva</b>	<b>0.69</b>	<b>0.57</b>	<b>0.62</b>	<b>0.757</b>
	<b>Desgosto</b>	<b>0.79</b>	<b>0.66</b>	<b>0.72</b>	
	<b>Medo</b>	<b>0.45</b>	<b>0.5</b>	<b>0.47</b>	
	<b>Felicidade</b>	<b>0.9</b>	<b>0.89</b>	<b>0.9</b>	
	<b>Tristeza</b>	<b>0.6</b>	<b>0.65</b>	<b>0.63</b>	
	<b>Surpresa</b>	<b>0.82</b>	<b>0.86</b>	<b>0.84</b>	
	<b>Neutralidade</b>	<b>0.67</b>	<b>0.68</b>	<b>0.68</b>	
	<b>Média/Total</b>	<b>0.76</b>	<b>0.76</b>	<b>0.76</b>	
ResNet-14200	Raiva	0.74	0.55	0.63	0.74
	Desgosto	0.86	0.52	0.65	
	Medo	0.56	0.48	0.52	
	Felicidade	0.83	0.92	0.87	
	Tristeza	0.61	0.69	0.65	
	Surpresa	0.68	0.9	0.77	
	Neutralidade	0.74	0.55	0.63	
	Média/Total	0.74	0.74	0.73	



## ANEXO C – Resultados por Bases de Validação

Neste anexo são apresentados os resultados por base de validação: CIFE-test, CIFE-train, CK, FER, JAFFE, KDEF, NovaEmotions e RAFD.

Tabela 18 – Resultados experimentais das redes neurais de convolução avaliando a base de validação CIFE-Test

Arquitetura	Emoção	Precisão	Revocação	F1-score	Acurácia
Alexnet	Raiva	0.58	0.7	0.63	0.626
	Desgosto	0.35	0.38	0.36	
	Medo	0.46	0.35	0.4	
	Felicidade	0.72	0.8	0.76	
	Tristeza	0.72	0.61	0.66	
	Surpresa	0.73	0.62	0.67	
	Neutralidade	0.5	0.51	0.5	
	Média/Total	0.63	0.63	0.62	
Inception-V3	Raiva	0.62	0.62	0.62	0.613
	Desgosto	0.26	0.3	0.28	
	Medo	0.47	0.46	0.46	
	Felicidade	0.86	0.79	0.83	
	Tristeza	0.56	0.54	0.55	
	Surpresa	0.7	0.61	0.65	
	Neutralidade	0.45	0.57	0.5	
	Média/Total	0.63	0.61	0.62	
ResNet-34	<b>Raiva</b>	<b>0.78</b>	<b>0.6</b>	<b>0.68</b>	<b>0.687</b>
	<b>Desgosto</b>	<b>0.65</b>	<b>0.34</b>	<b>0.45</b>	
	<b>Medo</b>	<b>0.4</b>	<b>0.42</b>	<b>0.41</b>	
	<b>Felicidade</b>	<b>0.84</b>	<b>0.87</b>	<b>0.85</b>	
	<b>Tristeza</b>	<b>0.66</b>	<b>0.71</b>	<b>0.69</b>	
	<b>Surpresa</b>	<b>0.62</b>	<b>0.81</b>	<b>0.7</b>	
	<b>Neutralidade</b>	<b>0.61</b>	<b>0.6</b>	<b>0.6</b>	
	<b>Média/Total</b>	<b>0.69</b>	<b>0.69</b>	<b>0.68</b>	

Tabela 19 – Resultados experimentais das redes neurais de convolução avaliando a base de validação CIFE-Train

Arquitetura	Emoção	Precisão	Revocação	F1-score	Acurácia
Alexnet	Raiva	0.55	0.6	0.58	0.624
	Desgosto	0.3	0.31	0.3	
	Medo	0.5	0.49	0.5	
	Felicidade	0.73	0.85	0.79	
	Tristeza	0.68	0.6	0.64	
	Surpresa	0.71	0.53	0.61	
	Neutralidade	0.54	0.53	0.53	
	Média/Total	0.62	0.62	0.62	
Inception-V3	Raiva	0.58	0.5	0.53	0.572
	Desgosto	0.21	0.26	0.23	
	Medo	0.4	0.47	0.43	
	Felicidade	0.83	0.79	0.81	
	Tristeza	0.49	0.49	0.49	
	Surpresa	0.69	0.54	0.61	
	Neutralidade	0.44	0.53	0.48	
	Média/Total	0.59	0.57	0.58	
ResNet-34	<b>Raiva</b>	<b>0.71</b>	<b>0.61</b>	<b>0.66</b>	<b>0.673</b>
	<b>Desgosto</b>	<b>0.44</b>	<b>0.22</b>	<b>0.3</b>	
	<b>Medo</b>	<b>0.46</b>	<b>0.52</b>	<b>0.49</b>	
	<b>Felicidade</b>	<b>0.82</b>	<b>0.83</b>	<b>0.83</b>	
	<b>Tristeza</b>	<b>0.67</b>	<b>0.74</b>	<b>0.7</b>	
	<b>Surpresa</b>	<b>0.6</b>	<b>0.75</b>	<b>0.67</b>	
	<b>Neutralidade</b>	<b>0.61</b>	<b>0.56</b>	<b>0.58</b>	
	<b>Média/Total</b>	<b>0.67</b>	<b>0.67</b>	<b>0.67</b>	

Tabela 20 – Resultados experimentais das redes neurais de convolução avaliando a base de validação CK

Arquitetura	Emoção	Precisão	Revocação	F1-score	Acurácia
Alexnet	Raiva	0.91	1	0.95	0.96
	Desgosto	0.98	0.97	0.98	
	Medo	0.89	0.96	0.92	
	Felicidade	0.99	0.99	0.99	
	Tristeza	0.98	0.84	0.91	
	Surpresa	1	0.94	0.97	
	Neutralidade	0	0	0	
	Média/Total	0.97	0.96	0.96	
Inception-V3	Raiva	0.93	0.94	0.93	0.954
	Desgosto	0.96	0.94	0.95	
	Medo	0.89	0.96	0.92	
	Felicidade	0.99	0.98	0.99	
	Tristeza	0.91	0.97	0.94	
	Surpresa	1	0.94	0.97	
	Neutralidade	0	0	0	
	Média/Total	0.96	0.95	0.96	
ResNet-34	<b>Raiva</b>	<b>0.97</b>	<b>0.96</b>	<b>0.97</b>	<b>0.969</b>
	<b>Desgosto</b>	<b>1</b>	<b>0.92</b>	<b>0.96</b>	
	<b>Medo</b>	<b>0.91</b>	<b>0.99</b>	<b>0.95</b>	
	<b>Felicidade</b>	<b>0.98</b>	<b>0.99</b>	<b>0.99</b>	
	<b>Tristeza</b>	<b>0.94</b>	<b>0.96</b>	<b>0.95</b>	
	<b>Surpresa</b>	<b>0.98</b>	<b>0.99</b>	<b>0.99</b>	
	<b>Neutralidade</b>	<b>0</b>	<b>0</b>	<b>0</b>	
	<b>Média/Total</b>	<b>0.97</b>	<b>0.97</b>	<b>0.97</b>	

Tabela 21 – Resultados experimentais das redes neurais de convolução avaliando a base de validação FER

Arquitetura	Emoção	Precisão	Revocação	F1-score	Acurácia
Alexnet	Raiva	0.39	0.5	0.44	0.543
	Desgosto	0.45	0.17	0.25	
	Medo	0.37	0.33	0.35	
	Felicidade	0.74	0.79	0.76	
	Tristeza	0.4	0.29	0.34	
	Surpresa	0.72	0.64	0.67	
	Neutralidade	0.49	0.52	0.51	
	Média/Total	0.54	0.54	0.54	
Inception-V3	Raiva	0.43	0.4	0.41	0.529
	Desgosto	0.13	0.25	0.17	
	Medo	0.39	0.36	0.37	
	Felicidade	0.81	0.77	0.79	
	Tristeza	0.29	0.38	0.32	
	Surpresa	0.72	0.64	0.68	
	Neutralidade	0.49	0.46	0.47	
	Média/Total	0.55	0.53	0.54	
ResNet-34	<b>Raiva</b>	<b>0.61</b>	<b>0.42</b>	<b>0.5</b>	<b>0.604</b>
	<b>Desgosto</b>	<b>0.69</b>	<b>0.28</b>	<b>0.39</b>	
	<b>Medo</b>	<b>0.35</b>	<b>0.47</b>	<b>0.4</b>	
	<b>Felicidade</b>	<b>0.87</b>	<b>0.81</b>	<b>0.84</b>	
	<b>Tristeza</b>	<b>0.41</b>	<b>0.42</b>	<b>0.41</b>	
	<b>Surpresa</b>	<b>0.71</b>	<b>0.76</b>	<b>0.73</b>	
	<b>Neutralidade</b>	<b>0.56</b>	<b>0.61</b>	<b>0.58</b>	
	<b>Média/Total</b>	<b>0.62</b>	<b>0.6</b>	<b>0.61</b>	

Tabela 22 – Resultados experimentais das redes neurais de convolução avaliando a base de validação JAFFE

Arquitetura	Emoção	Precisão	Revocação	F1-score	Acurácia
Alexnet	Raiva	0.62	0.62	0.62	0.618
	Desgosto	0.57	0.36	0.44	
	Medo	0	0	0	
	Felicidade	0.93	0.93	0.93	
	Tristeza	0.5	0.4	0.44	
	Surpresa	0.67	0.8	0.73	
	Neutralidade	0	0	0	
	Média/Total	0.65	0.62	0.63	
Inception-V3	Raiva	0.45	0.62	0.53	0.636
	Desgosto	0.6	0.27	0.37	
	Medo	0.5	0.5	0.5	
	Felicidade	0.93	0.93	0.93	
	Tristeza	0.44	0.4	0.42	
	Surpresa	0.82	0.9	0.86	
	Neutralidade	0	0	0	
	Média/Total	0.67	0.64	0.64	
ResNet-34	<b>Raiva</b>	<b>0.8</b>	<b>0.5</b>	<b>0.62</b>	<b>0.654</b>
	<b>Desgosto</b>	<b>0.71</b>	<b>0.45</b>	<b>0.56</b>	
	<b>Medo</b>	<b>1</b>	<b>0.5</b>	<b>0.67</b>	
	<b>Felicidade</b>	<b>1</b>	<b>0.64</b>	<b>0.78</b>	
	<b>Tristeza</b>	<b>0.58</b>	<b>0.7</b>	<b>0.64</b>	
	<b>Surpresa</b>	<b>0.53</b>	<b>1</b>	<b>0.69</b>	
	<b>Neutralidade</b>	<b>0</b>	<b>0</b>	<b>0</b>	
	<b>Média/Total</b>	<b>0.75</b>	<b>0.65</b>	<b>0.67</b>	

Tabela 23 – Resultados experimentais das redes neurais de convolução avaliando a base de validação KDEF

Arquitetura	Emoção	Precisão	Revocação	F1-score	Acurácia
Alexnet	Raiva	0.65	0.59	0.62	0.624
	Desgosto	0.6	0.79	0.68	
	Medo	0.48	0.43	0.45	
	Felicidade	0.8	0.93	0.86	
	Tristeza	0.52	0.4	0.45	
	Surpresa	0.81	0.61	0.7	
	Neutralidade	0.52	0.62	0.57	
	Média/Total	0.63	0.62	0.62	
Inception-V3	Raiva	0.51	0.53	0.52	0.587
	Desgosto	0.62	0.45	0.52	
	Medo	0.54	0.36	0.43	
	Felicidade	0.93	0.87	0.9	
	Tristeza	0.34	0.5	0.41	
	Surpresa	0.76	0.8	0.78	
	Neutralidade	0.54	0.6	0.57	
	Média/Total	0.61	0.59	0.59	
ResNet-34	<b>Raiva</b>	<b>0.7</b>	<b>0.7</b>	<b>0.7</b>	<b>0.751</b>
	<b>Desgosto</b>	<b>0.65</b>	<b>0.88</b>	<b>0.75</b>	
	<b>Medo</b>	<b>0.73</b>	<b>0.44</b>	<b>0.55</b>	
	<b>Felicidade</b>	<b>0.94</b>	<b>0.99</b>	<b>0.96</b>	
	<b>Tristeza</b>	<b>0.69</b>	<b>0.7</b>	<b>0.7</b>	
	<b>Surpresa</b>	<b>0.73</b>	<b>0.9</b>	<b>0.8</b>	
	<b>Neutralidade</b>	<b>0.86</b>	<b>0.66</b>	<b>0.75</b>	
	<b>Média/Total</b>	<b>0.76</b>	<b>0.75</b>	<b>0.74</b>	

Tabela 24 – Resultados experimentais das redes neurais de convolução avaliando a base de validação NovaEmotions

Arquitetura	Emoção	Precisão	Revocação	F1-score	Acurácia
Alexnet	Raiva	0.68	0.33	0.44	0.857
	Desgosto	0.69	0.63	0.66	
	Medo	0.57	0.25	0.34	
	Felicidade	0.89	0.94	0.92	
	Tristeza	0.86	0.73	0.79	
	Surpresa	0.9	0.87	0.89	
	Neutralidade	0.76	0.77	0.76	
	Média/Total	0.85	0.86	0.85	
Inception-V3	Raiva	0.3	0.22	0.25	0.853
	Desgosto	0.68	0.6	0.64	
	Medo	0.44	0.15	0.23	
	Felicidade	0.9	0.94	0.92	
	Tristeza	0.89	0.67	0.77	
	Surpresa	0.92	0.88	0.9	
	Neutralidade	0.71	0.77	0.74	
	Média/Total	0.85	0.85	0.85	
ResNet-34	<b>Raiva</b>	<b>0.54</b>	<b>0.33</b>	<b>0.41</b>	<b>0.87</b>
	<b>Desgosto</b>	<b>0.85</b>	<b>0.65</b>	<b>0.74</b>	
	<b>Medo</b>	<b>0.73</b>	<b>0.15</b>	<b>0.25</b>	
	<b>Felicidade</b>	<b>0.93</b>	<b>0.93</b>	<b>0.93</b>	
	<b>Tristeza</b>	<b>0.71</b>	<b>0.88</b>	<b>0.78</b>	
	<b>Surpresa</b>	<b>0.93</b>	<b>0.89</b>	<b>0.91</b>	
	<b>Neutralidade</b>	<b>0.75</b>	<b>0.79</b>	<b>0.77</b>	
	<b>Média/Total</b>	<b>0.87</b>	<b>0.87</b>	<b>0.87</b>	



Tabela 25 – Resultados experimentais das redes neurais de convolução avaliando a base de validação RAFD

Arquitetura	Emoção	Precisão	Revocação	F1-score	Acurácia
Alexnet	Raiva	0.66	0.77	0.71	0.641
	Desgosto	0.58	0.75	0.66	
	Medo	0.74	0.58	0.65	
	Felicidade	0.72	0.81	0.76	
	Tristeza	0.67	0.36	0.47	
	Surpresa	0.76	0.73	0.74	
	Neutralidade	0.45	0.38	0.41	
	Média/Total	0.65	0.64	0.63	
Inception-V3	Raiva	0.65	0.72	0.68	0.674
	Desgosto	0.71	0.65	0.68	
	Medo	0.82	0.55	0.66	
	Felicidade	0.89	0.87	0.88	
	Tristeza	0.46	0.68	0.55	
	Surpresa	0.82	0.81	0.81	
	Neutralidade	0.49	0.48	0.49	
	Média/Total	0.69	0.67	0.68	
ResNet-34	<b>Raiva</b>	<b>0.67</b>	<b>0.92</b>	<b>0.77</b>	<b>0.775</b>
	<b>Desgosto</b>	<b>0.8</b>	<b>0.81</b>	<b>0.81</b>	
	<b>Medo</b>	<b>0.81</b>	<b>0.65</b>	<b>0.72</b>	
	<b>Felicidade</b>	<b>0.87</b>	<b>0.94</b>	<b>0.9</b>	
	<b>Tristeza</b>	<b>0.74</b>	<b>0.7</b>	<b>0.72</b>	
	<b>Surpresa</b>	<b>0.76</b>	<b>0.94</b>	<b>0.84</b>	
	<b>Neutralidade</b>	<b>0.8</b>	<b>0.44</b>	<b>0.57</b>	
	<b>Média/Total</b>	<b>0.78</b>	<b>0.78</b>	<b>0.77</b>	