

DESENVOLVIMENTO DE UM APLICATIVO IOS UTILIZANDO OS *FRAMEWORKS CORE ML E VISION* PARA APLICAÇÃO DE TÉCNICAS DE DETECÇÃO DE OBJETOS PREDOMINANTES PRESENTES EM IMAGENS

Alan Henrique Pégoli¹ – Faculdade de Tecnologia de Carapicuíba

Profa. Orientadora Dra. Silvia Maria Farani Costa² – Faculdade de Tecnologia de
Carapicuíba

RESUMO

Este artigo trata sobre o desenvolvimento de um aplicativo móvel para o sistema operacional de dispositivos *iPhone* e *iPad* (*iOS*) que faz uso de dois recém-lançados *frameworks* da *Apple Inc.*, são eles, *Core ML* e *Vision*, para aplicar técnicas de análise de imagem de alto desempenho e de visão computacional para identificar objetos em imagens e vídeos. O objetivo deste artigo é introduzir e difundir tecnologias recém-chegadas ao mercado a fim de fomentar o aprendizado, a reflexão e o debate sobre o desenvolvimento *mobile* e as tendências em torno disso no meio acadêmico. O aplicativo, que é o produto deste artigo, contém uma interface gráfica pela qual o usuário é capaz de interagir com o sistema. Essa interface se compõe pela tela de câmera, onde o usuário pode enxergar o mundo real, e uma caixa de texto, que é onde o usuário receberá respostas responsivas do sistema sobre o objeto predominante que está à frente da câmera.

Palavras-chave: Aplicativo móvel. Aprendizado de máquina. Visão computacional.

ABSTRACT

This article discusses about the development of a mobile application for iPhone and iPad (iOS) devices operating systems that make use of two newly released frameworks from Apple Inc., namely Core ML and Vision, to apply high-performance image analysis techniques to identify objects in images and videos. The core purpose of this article is to introduce and disseminate new technologies to the market in order to foster learning, reflection and debate about the mobile development and the trends in the academic world. The application, which is the product of this article, contains a graphical interface through which the user is able to interact with the system. And this interface is composed by the camera screen, where the user can see the real world, and a text box, which is where the user will receive responsive responses from the system on the predominant object that is in front of the camera.

Keywords: Computer vision. Machine learning. Mobile application.

¹ - Aluno do CST em Análise e Desenvolvimento de Sistemas – e-mail: alanpegoli@icloud.com

² - Doutora em Engenharia Elétrica – e-mail: silvia.costa01@fatec.sp.gov.br

1 INTRODUÇÃO

Neste capítulo apresentam-se alguns conceitos básicos para posicionar o leitor das tecnologias aqui abordadas.

1.1 Aprendizado de Máquina

Aprendizado de máquina é uma área da ciência da computação que surgiu do estudo do reconhecimento de padrões e de teorias de inteligência artificial e aprendizados automáticos e que consiste na implementação de *softwares* que podem aprender autonomamente (HOSCH, 2009).

Samuel (1959) definiu o aprendizado de máquina como a habilidade de computadores aprenderem e/ou executarem uma tarefa para o qual não foram explicitamente programados.

Dentre as diversas abordagens e aplicações do aprendizado de máquina, encontra-se o uso de redes neurais artificiais e algoritmos genéticos. E um método de aprendizado de máquina muito divulgado é chamado de aprendizagem supervisionada, que consiste em uma rede neural que recebe interações já conhecidas de um agente externo (geralmente humano) de modo que com o passar das interações a rede se torne cada vez mais ajustada e assertiva (INTERNATIONAL CONFERENCE ON COMPUTATIONAL CREATIVITY, 2016). Por exemplo, um modelo que tenha sido treinado nos preços históricos das casas de uma região pode ser capaz de prever o preço de uma casa, dado o número de quartos e banheiros.

1.2 Visão Computacional

Visão computacional é uma área da ciência da computação que se dedica a aplicar conceitos e métodos de aprendizagem de máquina para conseguir que sistemas de softwares e hardwares, por meio de representações 2D do mundo físico, possam compreender uma cena 3D real (DAVIES, 2015). Em suma, trata-se de uma réplica computacional do processo de visão humana (embora o potencial dessa tecnologia vá muito além), podendo incorporar outras tecnologias, como sensores capazes de ver no escuro ou através de paredes.

Dentro de visão computacional, tem-se o que é chamado de reconhecimento de imagem, que trata a respeito de análise pixel a pixel em busca de padrões que formam objetos ou símbolos (DAVIES, 2015).

O objetivo final deste trabalho foi aplicar técnica de aprendizado de máquina a dispositivos iPhone por meio dos *frameworks Core ML e Vision*, e assim introduzir e difundir tecnologias recém-chegadas ao mercado para fomentar o aprendizado, a reflexão e o debate sobre o desenvolvimento mobile e as tendências em torno disso no meio acadêmico.

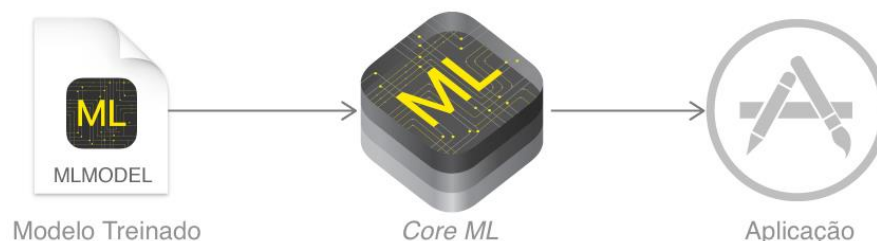
2 METODOLOGIA E DESENVOLVIMENTO

Para este trabalho foi usado um modelo treinado para detectar os objetos dominantes presentes em uma imagem de um conjunto de 1000 categorias, como árvores, animais, alimentos, veículos, pessoas, dentre outros. Para que fosse possível usar tal modelo, fez-se uso dos frameworks Core ML, que é a base para trabalhar com aprendizado de máquina intrinsecamente no iOS, e Vision, que torna a aplicação de aprendizado de máquina em aplicativos que implementam visão computacional menos complexa.

2.1 Core ML

O Core ML é um framework recém- lançado pela Apple Inc. que foi desenvolvido para transformar a aplicação de aprendizado de máquina em dispositivos móveis. Ele permite que modelos treinados sejam integrados e trabalhem intrinsecamente nos aplicativos, conforme Figura 1, o que garante a privacidade do usuário final, a responsividade e disponibilidade das funcionalidades do aplicativo, mesmo off-line, e reduz custos de transferência de dados (para o usuário final) e anula custos com servidores dedicados (para o desenvolvedor), uma vez que desloca o processamento direto para o dispositivo móvel (APPLE INC., 2017).

Figura 1 – Framework Core ML

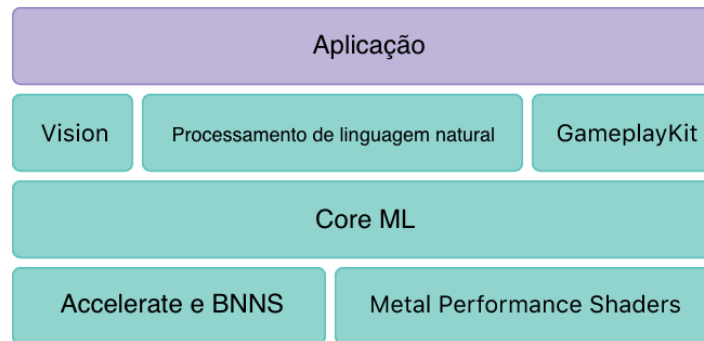


Fonte: Reprodução adaptada de Apple Inc.

O Core ML é a base para outros frameworks e funcionalidades direcionados ao aprendizado de máquina dentro do ecossistema do iOS. Ele suporta o framework Vision para análise de imagens, o Foundation para processamento de linguagem natural (por exemplo, a

classe `NSLinguisticTagger`) e o `GameplayKit` para avaliar árvores de decisão aprendidas. Por sua vez, o `Core ML` se constrói sobre outros frameworks ainda mais primitivos e de baixo nível do iOS, como `Accelerate` e `BNNS`, bem como o `Metal Performance Shaders`, conforme demonstrado na Figura 2.

Figura 2 – Camadas do *Core ML*



Fonte: Reprodução adaptada de Apple Inc.

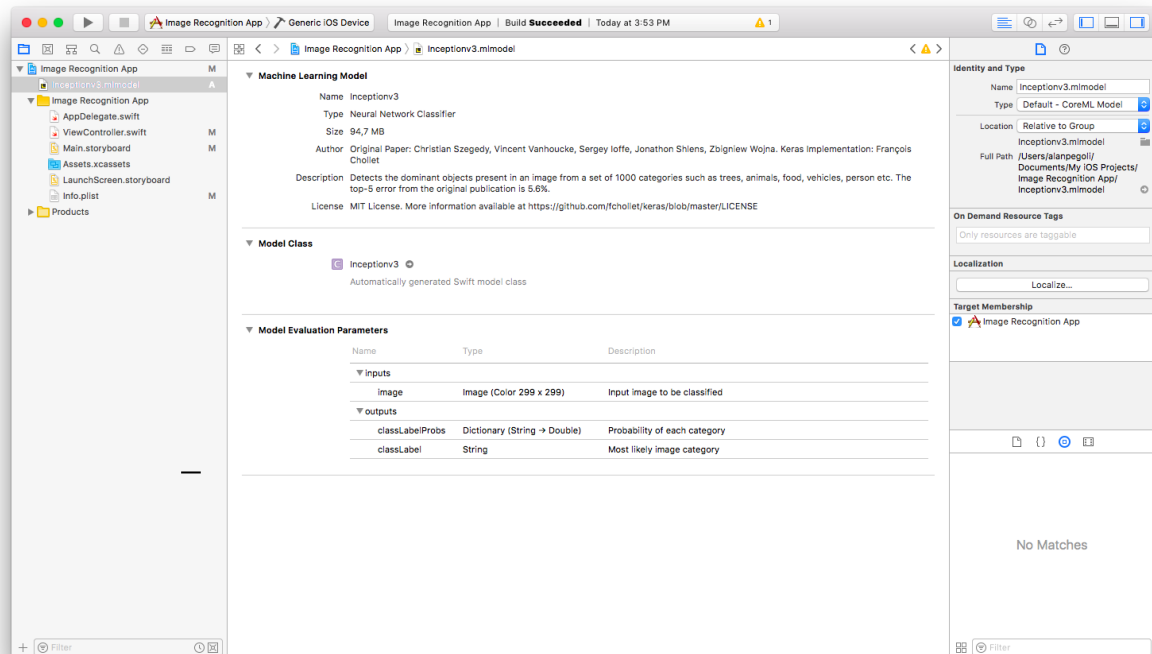
O `Core ML` suporta uma variedade de modelos de aprendizado de máquinas, incluindo redes neurais, conjuntos de árvores, modelos lineares generalizados etc. O `Core ML` requer um formato específico de modelo (arquivo com extensão `.mlmodel`). A Apple Inc. fornece alguns desses modelos populares de código aberto que já estão no formato do modelo `Core ML`. No entanto, modelos e dados de treinamento que não estão no formato do modelo `Core ML` podem ser convertidos através de ferramentas do próprio framework (APPLE INC., 2017).

2.2 Implementando o *Core ML*

Neste trabalho foi usado um modelo fornecido pela Apple Inc., o `Inception v3`, que foi originalmente convertido de um modelo de classificação de imagem treinado para `Keras`, uma biblioteca de rede neural de código aberto escrita em `Python`. Este modelo detecta os objetos dominantes presentes em uma imagem de um conjunto de 1000 categorias, como árvores, animais, alimentos, veículos, pessoa etc. A margem de erro da publicação original é de 5,6%.

Na Imagem 1 é demonstrado uma captura da tela de informações no `Xcode` (IDE para desenvolvimento iOS) sobre o modelo `Core ML` adicionado ao projeto da aplicação, incluindo o tipo de modelo, sua entrada e saídas esperadas. A entrada para o modelo é uma imagem colorida de tamanho 299p X 299p. A produção do modelo são a categoria do objeto dominante na imagem e a probabilidade da categoria.

Imagem 1 – Captura de tela das informações do modelo *Core ML*



Fonte: Produção própria.

O Xcode usa essas informações sobre as entradas e saídas do modelo para gerar automaticamente uma interface programática personalizada para o modelo, que é usada para interagir com o próprio modelo no código. Para o modelo “Inceptionv3.mlmodel”, o Xcode cria classes para representar o próprio modelo (Inceptionv3), as entradas do modelo (Inceptionv3Input) e as saídas (Inceptionv3Output). O Trecho de Código 1 traz um exemplo de instanciação de um objeto da classe Inceptionv3:

Trecho de Código 1 – Instanciando um objeto da classe *Inceptionv3*

```
var model: Inceptionv3!  
model = Inceptionv3()
```

A classe Inceptionv3 tem um método de predição gerado (prediction(image: CVPixelBuffer)) que é usado para prever a categoria e sua probabilidade do objeto dominante numa certa imagem, que é usada como valor de entrada do modelo. O resultado desse método

é uma instância Inceptionv3Output, no código chamamos o resultado prediction, conforme Trecho de Código 2:

Trecho de Código 2 – Recuperando o resultado do método de predição

```
guard let prediction = try? model.prediction(image: pixelBuffer!) else {  
return }
```

Para acessar a informação contida no resultado que traz a categoria do objeto dominante na imagem, usa-se a propriedade classLabel da classe Inceptionv3Output, conforme Trecho de Código 3:

Trecho de Código 3 – Utilizando o resultado do método de predição

```
self.classifierLabel.text = "I think this is a \(prediction.classLabel)."
```

A predição gerada pode resultar um erro. O tipo mais comum de erro encontrado ao trabalhar com o Core ML ocorre quando o tipo de dados de entrada passado para o método não corresponde ao tipo de entrada que o modelo espera – no caso, uma imagem no formato errado poderia resultar num erro.

Quaisquer tipos de desajustes são capturados em tempo de compilação, e o aplicativo gera um erro fatal se algo der errado.

O Xcode compila o modelo Core ML em um recurso para otimizar a execução no dispositivo. Essa representação otimizada do modelo está incluída no pacote do aplicativo e é usada para fazer predições enquanto o aplicativo está sendo executado no dispositivo.

2.3 *Vision*

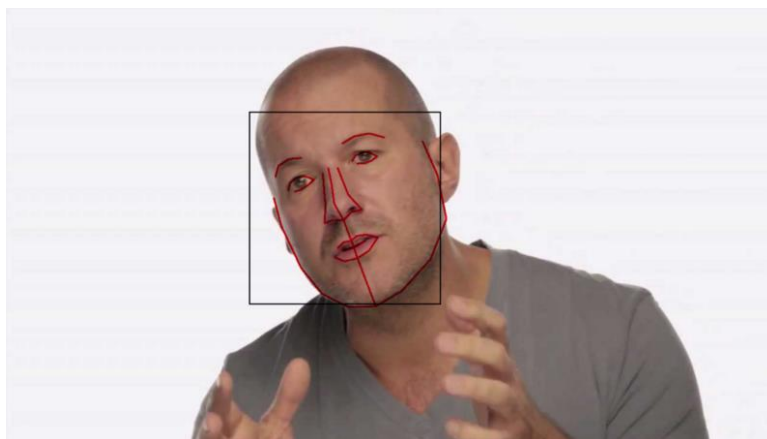
Vision é um framework que trabalha em conjunto com o Core ML para tornar a aplicação de aprendizado de máquina em aplicativos que implementam visão computacional menos complexa. Ele fornece o pipeline de imagens para suportar os modelos treinados automaticamente, evitando que seja necessário manipular as imagens manualmente, o que seria trabalhoso e demandaria tempo em termos de desenvolvimento.

O Vision é capaz de encontrar a superfície das regiões retangulares projetadas, detectar os retângulos e marcos da face, como contorno do rosto, linha mediana, boca, lábios, nariz, olhos, pupilas, sobrancelhas, dentre outras (vide Imagem 2), determinar o ângulo do horizonte em uma imagem, encontrar regiões de texto visíveis, encontrar e reconhecer códigos de barra, acompanhar o movimento de um objeto arbitrário previamente identificado em várias imagens ou quadros de vídeo, detectar transformações necessárias para alinhar o conteúdo de duas imagens, e, obviamente, processar imagens de um modelo Core ML (APPLE INC., 2017).

Existem 4 categorias de classes abstratas do Vision. São elas:

- VNRequest: cuida das requisições de análise de imagem. Possui um completion handler da requisição e um conjunto de resultados;
- VNObservation: cuida dos resultados das análises de imagem;
- VNImageRequestHandler , VNSequenceRequestHandler: processam um ou mais VNRequest em determinada imagem;
- VNErrorDomain: cuida dos erros das análises de imagem.

Imagem 2 – Exemplo de detecção de retângulos e marcos da face pelo *Vision*



Fonte: Reprodução de Apple Inc.

2.4 Implementando o *Vision*

Neste trabalho o Vision foi utilizado para realizar todo o tratamento de imagens, do modo natural e humanamente familiar para a forma correta de entrada do modelo Core ML.

O fluxo de trabalho padrão do Vision é criar um modelo, fazer uma ou mais requisições e, em seguida, criar e executar um completion handler da requisição, ou, em tradução livre, um manipulador de conclusão da requisição, que serve justamente para trabalhar com serviços assíncronos no Swift.

No Trecho de Código 4, tem-se o exemplo de um desempacotamento de modelo Core ML para um modelo Vision. Uma vez que o desempacotamento pode retornar um erro, usa-se o comando try:

Trecho de Código 4 – Desempacotando o modelo
Core ML para um modelo Vision

```
let model = try! VNCoreMLModel(for: Inceptionv3().model)
```

VNCoreMLRequest é uma solicitação de análise de imagem que usa um modelo Core ML para processar. Seu completion handler recebe os objetos de solicitação e erro. O resultado dessa solicitação no caso deste trabalho é um VNClassificationObservation, que é o que o Vision retorna quando o modelo Core ML é um classificador, ao invés de um preditor ou um processador de imagem. E o Inception v3 é um classificador porque ele prediz apenas um recurso: a categoria do objeto predominante na imagem.

A VNClassificationObservation tem duas propriedades: identifier (uma String que classifica a categoria do objeto) e confidence (um número entre 0 e 1, é a probabilidade da classificação estar correta). Neste trabalho, considerou-se válido apenas aqueles resultados com confidence acima de 50%.

2.5 Inception v3

Neste trabalho fez-se uso da rede neural Inception V3, pois ela está classificada como uma rede neural convolucional, que é um tipo muito usado para reconhecimento de imagens, uma vez que seus neurônios individuais são organizados de modo a responder regiões de sobreposição no campo visual (SZEGEDY, 2015).

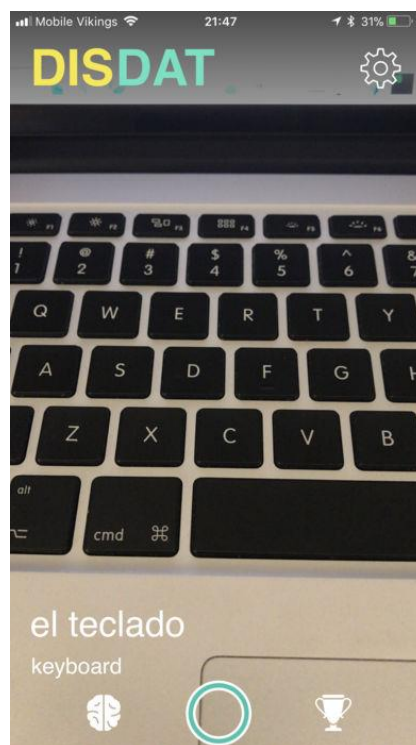
3 REFERENCIAL TEÓRICO

Foram encontrados alguns aplicativos já fazendo uso dos frameworks recém-lançados pela Apple Inc. logo após a abertura oficial da nova App Store, loja de e-commerce de aplicativos para dispositivos Apple em geral, que ocorreu no dia 19 de setembro de 2017. Ressaltam-se aqui dois deles, que inclusive usam o mesmo modelo Core ML que o usado neste trabalho.

3.1 *DISDAT*

DISDAT é um aplicativo educativo (Imagem 3), desenvolvido pela Balloon Inc., que faz uso do Core ML, do Vision e do NLP (Natural Language Processing, em tradução livre, Processamento de Linguagem Natural), para ensinar seus usuários novas línguas por meio da observação do mundo (BALLOON INC., 2017).

Imagem 3 – Captura da tela principal do aplicativo *DISDAT*

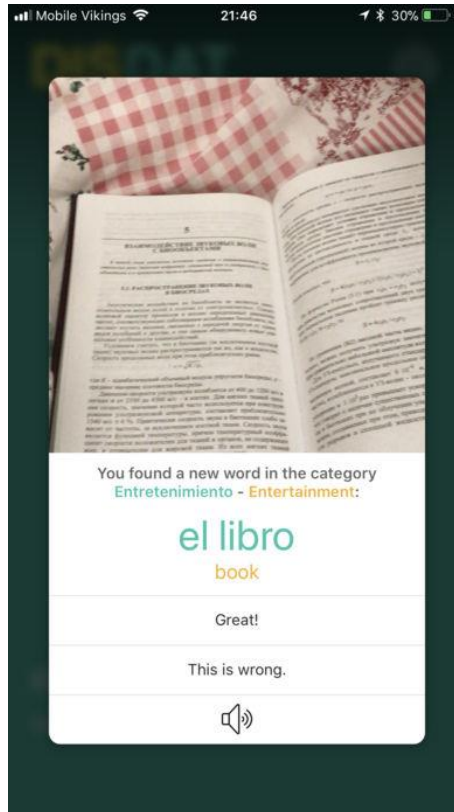


Fonte: Reprodução de Balloon Inc.

Segundo seus criadores, Balloon Inc. (2017), a ideia é reproduzir o processo pelo qual uma criança aprende os objetos do mundo e “gamificar” isso num aplicativo. O usuário deve apontar a câmera do dispositivo para o objeto e esperar pelo retorno, que é, devido ao Core ML, imediato.

São 120 palavras divididas em 22 categorias, em 6 línguas (inglês, espanhol, alemão, francês, italiano e russo). A medida que o usuário encontra os objetos recém-classificados, o aplicativo exibe um cartão, como um alerta, com a escrita e a pronúncia do objeto na língua desejada (Imagem 4). E também conta o progresso do usuário naquele idioma conforme o usuário for encontrando outras novas palavras.

Imagem 4 – Captura do cartão do aplicativo *DISDAT*



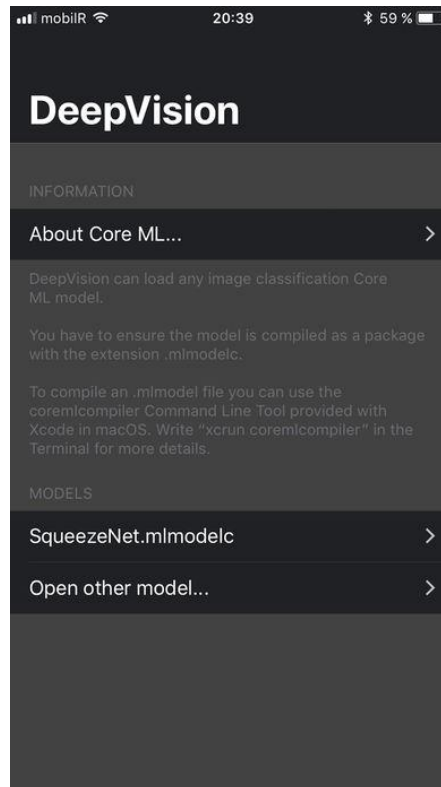
Fonte: Reprodução de Balloon Inc.

3.2 DEEPCONVISION

DeepVision é um aplicativo comercial, desenvolvido por Pedro Jose Pereira Vieito, tanto para iPhone como para Mac. Ele fornece uma plataforma para testes de modelos Core ML de classificação de imagem (VIEITO, 2017).

O aplicativo funciona basicamente da mesma forma, usando Core ML e Vision, com a ressalva de que seu modelo Core ML não é fixo. Ele tem um modelo predefinido, que é o SqueezeNet, mas seu grande diferencial é que também permite que usuários incluam seus próprios modelos Core ML de classificação de imagem para testarem (Imagem 5). Assim, ele é capaz de executar o modelo em tempo real na câmera de vídeo do dispositivo, seja um iPhone ou um Mac, e trazer os resultados do modelo, bem como suas informações descritivas, como criadores, licença e descrição.

Imagem 5 – Captura da tela do DeepVision para iPhone



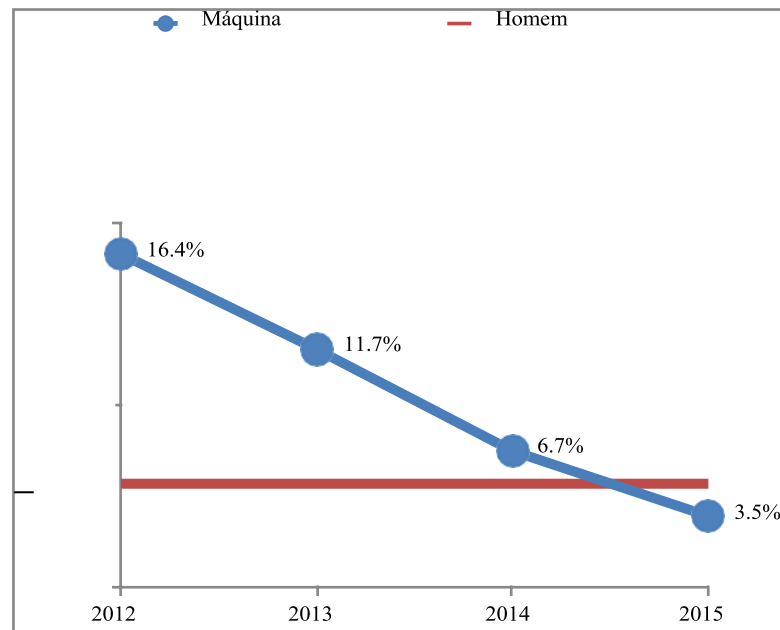
Fonte: Reprodução de Apple Inc.

4 RESULTADOS E DISCUSSÃO

No desafio de reconhecimento de imagem da IMAGENet de 2014, o Google Inc. se apresentou com uma abordagem de rede neural convolucional para reconhecimento de objetos em imagens, que teve taxa de erro de 6,6%, quase metade da taxa de 11,7% do ano anterior. No entanto, um humano foi capaz de cumprir o mesmo desafio com uma taxa de erro de apenas 5,1%.

Já em 2015, a Microsoft Inc. anunciou que havia conseguido bater o record humano com uma taxa de erro de apenas 4,94%. E em alguns meses depois, novamente no desafio da IMAGENet, a Microsoft Inc. quebrou seu próprio record com uma taxa de erro de 3,5%. (Vide Gráfico 1) Foi então que a visão computacional se tornou novamente popular e o tópico central em diversas discussões sobre inteligência artificial, aprendizado de máquina e deep learning (RUSSAKOVSKY et al., 2014).

Gráfico 1 – Teste *IMAGENet*



Fonte: IMAGENet Inc.

Este trabalho colocou à disposição informações e conhecimentos acerca do que tem sido encarado como a próxima tendência em termos de tecnologia para desenvolvimento de aplicações para dispositivos iPhone. Com os frameworks Core ML e Vision, bem como NLP, é possível afirmar que a Apple Inc. está colocando uma pedra fundamental para as inúmeras aplicações que regerão os próximos anos de desenvolvimento em sua plataforma (DAVIES et al., 2016).

Um grande ponto a favor dessa nova tecnologia é que ela foi projetada para ser executada em dispositivos existentes, o que significa que os usuários finais não terão que atualizar para um hardware mais caro e especializado para aproveitar os benefícios do aprendizado de máquinas em seus dispositivos. Isso dá uma abertura ainda maior para a sua adoção.

O software desenvolvido em conjunto com o trabalho trouxe uma breve nuance da capacidade dos frameworks desenvolvidos e recém-lançados pela Apple Inc. , e ajuda a tanger as possibilidades de aplicações que podem vir a implementar aprendizados de máquina de uma forma mais segura, barata, off-line e muito menos complexa em termos de desenvolvimento.

O modelo de rede neural convolucional utilizado neste trabalho apresenta uma taxa de erro top-5 de 5.06% no conjunto de teste da IMAGENet (SZEGEDY, 2015). Há, porém, relato

de que com um conjunto de três resíduos e um Inception-v4, é possível alcançar 3.08% de erro top-5 neste mesmo conjunto de tese (INTERNATIONAL CONFERENCE ON COMPUTATIONAL CREATIVITY, 2016).

5 CONSIDERAÇÕES FINAIS

A implementação do Core ML se mostrou relativamente simples e eficaz. Ela permitiu que o foco no desenvolvimento estivesse todo voltado para o manejo dos dados e das informações no "pré e pós- processamento". A implementação do Vision, por sua vez, serviu para atuar justamente nesse primeiro quesito, onde haveria bastante complexidade no tratamento de imagens.

Com o uso de ambas as plataformas foi possível desenvolver um aplicativo totalmente funcional capaz de detectar objetos predominantes presentes à frente da câmera e imprimir na tela informações sobre esses tais objetos.

REFERÊNCIAS

APPLE INC. (Cupertino) (Org.). **Core ML Apple Developer Documentation**: Integrate machine learning models into your app. 2017. Disponível em: <<https://developer.apple.com/documentation/coreml>>. Acesso em: 30 set. 2017.

APPLE INC. (Cupertino) (Org.). **Vision Apple Developer Documentation**: Apply high-performance image analysis and computer vision techniques to identify faces, detect features, and classify scenes in images and video. 2017. Disponível em: <<https://developer.apple.com/documentation/vision>>. Acesso em: 30 set. 2017.

BALLOON INC. (Org.). **DISDAT**: iOS app to learn a language using machine learning. 2017. Disponível em: <<https://disdat.ai>>. Acesso em: 30 set. 2017.

DAVIES, Roy. **Machine Vision**: Theory, Algorithms, Practicalities. 3. ed. Amsterdã: Elsevier, 2015. 934 p.

DAVIES, Sam; RAMES, Jeff; TURTON, Rich. **IOS 10 by Tutorials**: Learning the new iOS APIs with Swift 3. Virgínia: Razeware Llc, 2016. 324 p.

HOSCH, William L.. **Machine Learning**. Chicago: Encyclopædia Britannica, Inc., 2009. Disponível em: <<https://global.britannica.com/technology/machine-learning>>. Acesso em: 30 set. 2017.

INTERNATIONAL CONFERENCE ON COMPUTATIONAL CREATIVITY, 7., 2016, Paris. **Proceedings Of The Seventh International Conference On Computational**

Creativity. Paris: Sony CSL, 2016. 403 p. Disponível em: <http://www.computationalcreativity.net/iccc2016/wp-content/uploads/2016/08/Proceedings_ICCC16.pdf>. Acesso em: 30 set. 2017.

RUSSAKOVSKY, Olga et al. **ImageNet Large Scale Visual Recognition Challenge.** 3. ed. Carolina do Norte: Ijcv, 2014. 43 p. Disponível em: <<https://arxiv.org/pdf/1409.0575.pdf>>. Acesso em: 30 set. 2017.

SAMUEL, Arthur Lee. **Some Studies in Machine Learning Using the Game of Checkers.** Ibm Journal Of Research And Development. EUA, p. 535-554. mar. 1959. Disponível em: <<https://www.cs.virginia.edu/~evans/greatworks/samuel1959.pdf>>. Acesso em: 30 set. 2017.

SZEGEDY, Christian. **Rethinking the Inception Architecture for Computer.** Vision. 2015. 10 f. Cornell University, Ithaca, 2015. Disponível em: <<https://arxiv.org/abs/1512.00567?context=cs>>. Acesso em: 30 set. 2017.

VIEITO, Pedro José Pereira. **PVIEITO.** Disponível em: <<https://pvieito.com>>. Acesso em: 30 set. 2017.

“O conteúdo expresso no trabalho é de inteira responsabilidade do(s) autor(es).”