

IF867 - Introdução à Aprendizagem Profunda

February 13, 2025

2º Lista Teórica

Discente: - Gabriel D'assumpção de Carvalho - gdc2@cin.ufpe.br

Curso: - Ciências Atuariais - 7º Período

LISTA TEÓRICA 3:

① As redes VGS são consideradas redes densas, onde foi feito a estrutura de blocos, contendo um parâmetro de utilização da n filtros de kernel (3×3) , seguido de MAX-POOLING (2×2) e stride de 2, mas ao longo dos blocos o número de filtros vão duplicando, aumentando o custo computacional.

O uso de kernel (3×3) é vantajoso pois com 2 filtros vamos ter que estimar apenas 18 parâmetros que é menos custoso que os 25 e 49 parâmetros dos kernel (5×5) e (7×7) respectivamente. Além disso, quando utilizamos filtros menores conseguimos ter uma hierarquia mais característica, devido a diminuição da conta das dimensões da imagem e utilização de mais filtros para capturar características diferentes.

② Um dos tipos mais famosos nos modelos de CNN era a maximizar o número de parâmetros que devem ser estimados pelas camadas totalmente conectadas. Por exemplo, se no último bloco de convolução temos um tensor $(16 \times 16 \times 512)$ e na primeira camada conectada temos 1000 neurônios, precisamos estimar 131.073.000 parâmetros. Por conta disso foi proposto o modelo NETWORK IN NETWORK (NIN) que implementa uma convolução de n filtros (1×1) , com isso é introduzindo não linearidade a rede, e diminuição de canais enquanto o output continua com a mesma resolução, fazendo com que não seja necessário camadas totalmente conectadas e diminuindo o custo computacional por se estimar n parâmetros.

O bloco inception proposto para GoogLeNet possui 4 ramificações, a primeira (P1) contendo apenas um kernel (1×1) , a segunda e a terceira iniciam com um filtro (1×1) seguido de uma convolução (3×3) , (5×5) e padding de 1 e 2 respectivamente e na última ramificação é feito um MAX-POOLING.

(3x3) e padding de 1, seguido por uma convolução de (1x1), depois todos os output das camadas são agrupados. Portanto, essa abordagem contribui para a detecção de diferentes características ao utilizar diferentes tamanhos de kernels e filtros, tudo isso sem ter um alto custo computacional devido as convoluções dos kernel (1x1) que diminui os canais de input.

③ O modelo ResNet não tem a proposta de blocos residuais, que é um bloco que tem duas convoluções e o seu output é somado com o input pela conexão residual, isso faz com que a rede se tenha que adicionar $g(x)$, dando que $y = x + g(x)$. Com essa estrutura a ResNet consegue construir uma rede mais profunda sem que aconteça o degradation problem que é o esquecimento da informação e causando o vanish gradient.

Com base nisso o modelo da ResNet consegue ter estrutura mais profunda, geralmente tendo uma precisão melhor que o Inception, mais por ser mais profunda o custo computacional é maior.

As funções mais relevantes para análise dos modelos são a ResU que ajuda no problema do vanish gradient e como mostrado no paper do AlexNet tem uma convergência mais rápida, também pode ser utilizado a Leaky Relu que mitiga o problema de neurônios mortos.

④ Quando temos poucos observações para treinar um modelo, podemos aplicar técnicas de data augmentation, pedindo as observações e rotando elas simultaneamente. Implementação do padding, espelhamento, aumento ou diminuição, a rotação de um canal de cor, e o zoom na imagem e dependendo do problema podemos quebrar o tensor em 2 partes e misturar como ordens.

A estrutura mais adequada para esse tipo de problema é a VGG, VGG inception ou ResNet, mas deve ser implementado técnicas para diminuir o risco de overfitting ao usar modelos profundos.

Para evitar o overfitting podemos utilizar técnicas como Transfer Learning, que é quando pegamos um modelo que já foi treinado e alteramos a sua saída para ser compatível com o problema a ser resolvido, também pode usar o Fine-tuning, que é o retreinamento dessa rede adquirida, geralmente as primeiras camadas são congeladas por serem mapas de características simples, que detectam borda e texturas. A aplicação de dropout é amplamente utilizada no treinamento para introduzir ruído e evitar o overfitting.

⑤

1) O modelo é o que possui a maior generalização, pois possui a loss de treino mais baixa, quando a sua acurácia de treinamento e teste são 92% e 90% respectivamente.

2) O modelo 1, possivelmente não tem overfitting, pois a sua loss de treinamento foi de 0,2 e a de teste foi 0,6, mostrando que o modelo não consegue generalizar para novos dados, pois se aplicarmos a exponencial no loss temos ($e^{0,2} \approx 0,82$) e ($e^{0,6} \approx 0,55$), portanto para 1000 vezes a probabilidade média dos acertos é 1,48 vezes maior, tendo um modelo quase aleatório para problemas binário, dado que a acurácia média 55% dos novos observações.

3) O modelo 2 apresenta fortes indícios de overfitting, pois a acurácia de treinamento e teste é de 70% e 68% e tendo uma probabilidade média de classe correta de aproximadamente 45% e 47%, indicando que o modelo não escolhe a classe correta com confiança.

4) O batch normalization é uma técnica usada para controlar a distribuição da ativação de uma rede neural, transformando a sua média para zero e variância 1. Essa manipulação ajuda a rede a não ter grandes explosões e vanishing gradient, e também ajuda na velocidade de convergência do modelo, por permitir taxa de aprendizagem mais rápida.

Em modelos de CNN o batch normalization vai atuar em cada canal de cada o mini-batch.

Como a normalização é feita usando o mini-batch, essa normalização introduz ruído na rede, tornando o modelo mais generalista, evitando o overfitting.

5) As CNN são modelos propostos para a área de visão computacional, pois muitos casos são empregados em problemas de classificação. Com isso o modelo é projetado para estimar a probabilidade da observação pertencer a cada categoria, sendo escolhida a classe de maior probabilidade.

Quando utilizado a função de MSE, os valores ter valores de 0 e 1. O modelo escolhe a classe independente do nível de confiança e, caso a previsão esteja errada, entretanto ao utilizar a cross-entropy loss, passa a modelar a probabilidade da classe certa. Por exemplo, se a classe correta é 1, o modelo estimou uma probabilidade de 0,9, mesmo adotando a loss vai ser de aproximadamente 0,105. Por conta disso o modelo tem uma convergência e melhora na função

bicidade das suas estimativas.

6) A função de ativação sigmoid é considerada uma função não linear saturada, isso implica que a sua derivada em muitos pontos é igual a zero e o valor máximo é 0,25. Por conta disso a aplicação dessa função em redes profundas acaba causando o vanishing gradient e como resultado no paper do AlexNet a sua convergência é mais lenta.

7) Devido a uma possível overfitting do modelo um padrão utilizar recursos de data augmentation e dropout para adicionar ruído ao modelo, também pode ser proposta uma estrutura menos profunda.

O modelo 2 que apresentou underfitting pode ser revisado se a função de ativação não está causando vanishing gradient, também é recomendado usar dropout e learning rates que tiver uma base de dados pequena, visando aumentar a profundidade da rede ou aumento na quantidade de épocas. É importante implementar métodos como adam que ajusta a taxa de aprendizagem ao decorrer do treinamento, trabalhar com taxas mais elevadas.

O modelo 3 apresentou um bom ajuste porém ambos os conjuntos de dados, mas ainda pode ser proposta implementação como dropout ou data augmentation para tornar o modelo mais estável ao longo do tempo.

① O TRANSFER LEARNING É UMA TÉCNICA QUE CONSISTE EM REUSAR O ALGORITMO E OS PARÂMETROS QUE JÁ FORAM CONSTRUÍDOS E ESTIMADOS EM UM TAREJA MAIS ROBUSTA E APLICAR AQUELE PARA OUTRA TAREJA EM ESPECÍFICO.

QUANDO ESTAMOS COM UM PROBLEMA DE CLASSIFICAÇÃO DE IMAGENS DEVEMOS AJUSTAR A ÚLTIMA CAMADA DO ALGORITMO PARA TER UM OBJECTIVO COM A MESMA SINTAXE DAS CLASSES DO NOSSO PROBLEMA, OUTRO PUNTO A SER USADO EM CONTRA É QUE MUITAS VEZES AS ÚLTIMAS CAMADAS APENAS CARACTERÍSTICAS ESPECÍFICAS, SENDO RECOMENDADO RETREINAR A REDE MAIS COMBILANDO AS CAMADAS INICIAIS, FAZENDO COM QUE O TREINAMENTO FIQUE MAIS EFICIENTE.

QUANDO SE TEM MUITAS OBSERVAÇÕES E TREINAMENTO É RECOMENDADO UTILIZAR O TRANSFER LEARNING APENAS PARA A INICIALIZAÇÃO DOS PESOS, SENDO FEITO UM FINE-TUNING, QUE É O AJUSTE EM TODOS OS PARÂMETROS DO MODELO.

② O VARIATIONAL AE (VAE) É UM MODELO QUE REDUZ A DIMENSIONALIDADE DO INPUT PARA DETECTAR AS VARIÁVEIS LATENTES, REMOVENDO RUÍDO E REDUZINDO A VARIÂNCIA DOS DADOS, ISSO AJUDA PARA MELHORAR A EFICIÊNCIA EM FOMAS DE LATENTES INDIVÍDUAS, PORTANTO AS SÃO MODELOS DETERMINÍSTICOS E INCLINAM-SE A GERAR NOVAS OBSERVAÇÕES.

USANDO O MODELO VARIATIONAL AUTOENCODERS (VAE) ADICIONA UM ALCORÇO NA VARIÁVEL LATENTE, PARA QUE QUANDO ATÁ A SUA DECODIFICAÇÃO SEJA GERADA UMA NOVA OBSERVAÇÃO, TAMBÉM PODE FAZER PERTURBAÇÕES NA VARIÁVEL LATENTE E VERIFICAR QUAL ALTERAÇÃO É FEITA.

O MODELO GAN VAI ATUAR COM DUAS REDES PROFUNDAS, UMA VAI SER A GERADORA DE DADOS A PARTIR DE UM RUÍDO E ESSA OBSERVAÇÃO DEVE ENGANAR A OUTRA REDE QUE VAI DISCRIMINAR A OBSERVAÇÃO DE PARA

COM A VERDADEIRA. EM RESUMO, O MODELO TENTA IDENTIFICAR A DISTRIBUIÇÃO DE UMA VARIÁVEL PARA NÃO SER POSSÍVEL IDENTIFICAR A OBSERVAÇÃO VERDADEIRA DA MESMA. PORÉM, ESSE MODELO TEM UM CUSTO COMPUTACIONAL GRANDE E INSTABILIDADE NO TREINAMENTO.

⑧ MODELOS DISCRIMINATIVOS SÃO MÉTODOS UTILIZADOS PARA DISTINGUIR QUAL O VALOR DE UMA VARIÁVEL Y A PARTIR DE UM CONJUNTO DE VARIÁVEIS PREDITORAS X , SENDO UM TIPO DE MODELO UTILIZADO PARA REGRESSÃO E CLASSIFICAÇÃO. TENHO EM VISTA ISSO, O MODELO TENTA ESTIMAR A PROBABILIDADE DE Y A PARTIR DE X , SENDO $P(Y|X)$.

OS MODELOS GERATIVOS SÃO MÉTODOS QUE SÃO APLICADOS PARA GERAR DE NOVOS DADOS. O SEU ALGORITMO RECEBE COMO ENTRADA UMA VARIÁVEL X E O MODELO TENTA ESTIMAR A SUA PROBABILIDADE E GERAR UM \hat{X} . INTERESSANTE NESSES MODELOS GERATIVOS NÃO MODELAM $P(X|\hat{X})$, AO INVÉS, APROXIMAR ESSA PROBABILIDADE VAMOS ESTAR APROXIMANDO \hat{X} DOS CASOS REAIS. (10)

AS VANTAGENS DOS MODELOS DISCRIMINATIVOS É QUE MUITOS POSSUEM UMA PODER MAIS CAPACITIVOS E CONSEGUEM SER AJUSTADA COM UMA BASE DE DADOS MENOR, PORÉM MUITOS MODELOS NECESSITAM QUE SUAS SUPORTES SEJAM INTERMEDIÁRIAS. (11)

SÃO OS MODELOS GERATIVOS É VANTAJOSO QUANDO NECESSITAMOS GERAR SIMULAÇÕES E NÃO PRECISA QUE SUPORTES SEJAM ATENDIDAS, MAS A SUA GRANDE DESVANTAGEM É QUE É PORCOSA DE UMA FORÇA COMPUTACIONAL MAIOR E UM BANCO DE DADOS GRANDE.

9) A rede U-Net foi proposta para fazer segmentação semântica, onde cada pixel pertence a uma classe.

Por conta desse tipo de segmentação ao aplicar CNNs profundas as convoluções e pooling fazem com que a cada camada o tensor de entrada diminua, perdendo informações do input devido ao downsampling.

Porém a cada perda de informação, a U-Net introduz a conexão de salto que concatena o tensor da etapa de downsampling que está em paralelo ao tensor de upsampling, fazendo com que as camadas mais profundas tenham informações com uma resolução mais alta, ajudando na preservação de bordas e características locais.

10)

1) supondo que $W = 5 \times 5$, $K_1 = 3 \times 3$, $K_2 = 2 \times 2$. O tamanho da imagem após a primeira convolução é dado por $(5 - 3 + 1) = (3 \times 3)$. Com a segunda convolução passa a ser $(3 - 2 + 1) = (2 \times 2)$.

2) A convolução é uma operação matemática que leva em conta a multiplicação de duas funções, passando por 1 ponto, porém tendo a principal função invertida, sendo dado por $(f * g)(x) = \sum_k f(k) \cdot g(x - k)$. Mas em CNN as convoluções são aplicadas sem a inversão, sendo aplicadas diretamente, isso acontece porque se os kernels forem invertidos a rede não estima os mesmos valores independentemente da rotação feita no filtro.