

Compressão de Redes Neurais

Vítor Yeso Fidelis Freitas

Motivação

- Aplicações de modelos de Deep Learning em ambientes com recursos computacionais limitados (Embarcados, celulares, FPGA, computadores pessoais)
- Eficiência Energética (Grandes modelos consomem muita energia)

Redes Neurais geralmente são treinadas utilizando muitos recursos.

- Pesos em ponto flutuante 32 ou 64 bits
- Muitos valores de pesos por rede (por exemplo: **GPT-3 tem 175 bilhões de parâmetros**)

Pruning e Quantização de Redes Neurais

Pruning: Processo que tenta **aumentar a quantidade de zeros** de uma Rede Neural

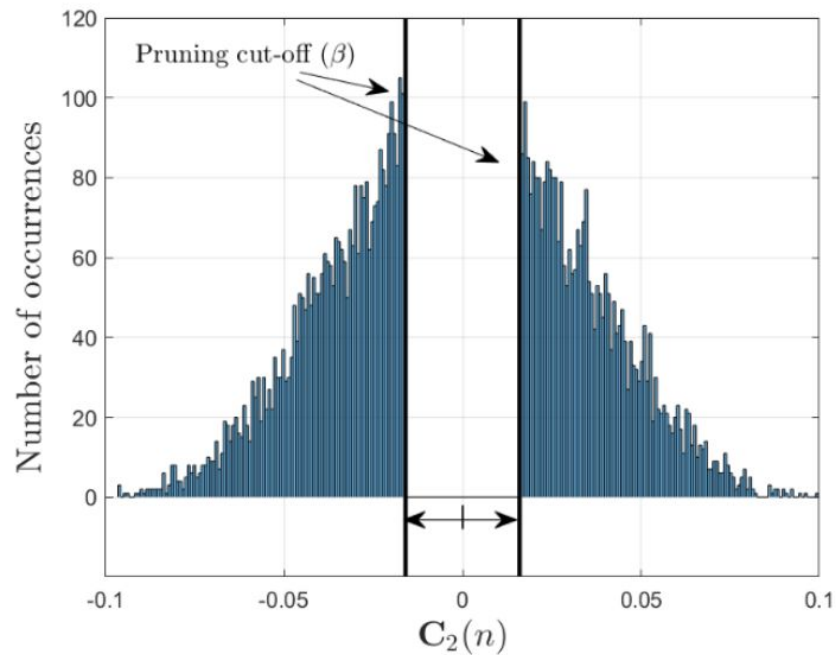
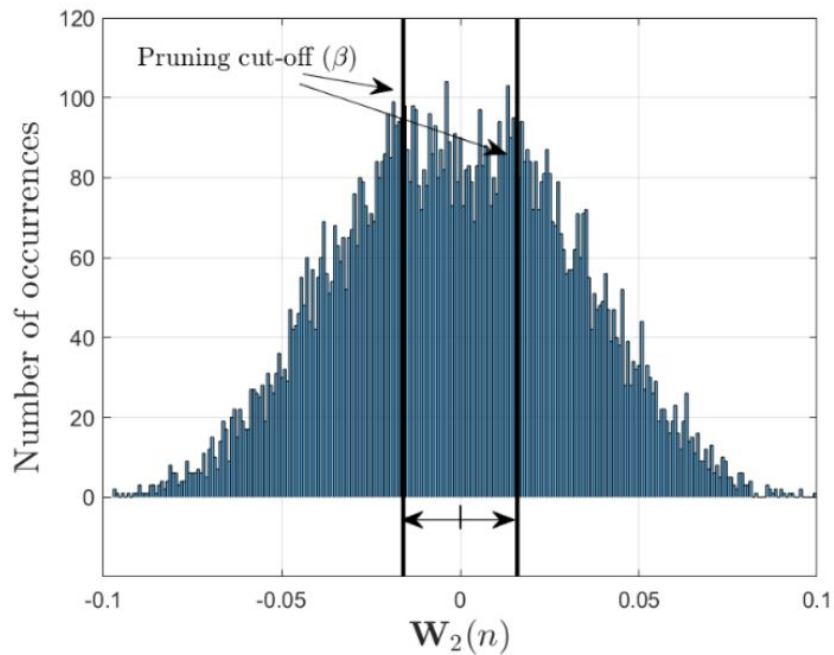
Quantização: Processo que tenta **diminuir a resolução de bits** dos parâmetros de uma Rede Neural

Pruning

$$P(\mathbf{W}_k(n), \beta_k) = \begin{cases} \mathbf{W}_k(n) & \text{if } |\mathbf{W}_k(n)| \geq \beta_k \\ 0 & \text{if } |\mathbf{W}_k(n)| < \beta_k \end{cases} \quad (1)$$

$$\beta_k = \gamma \times \sigma_k$$

Pruning

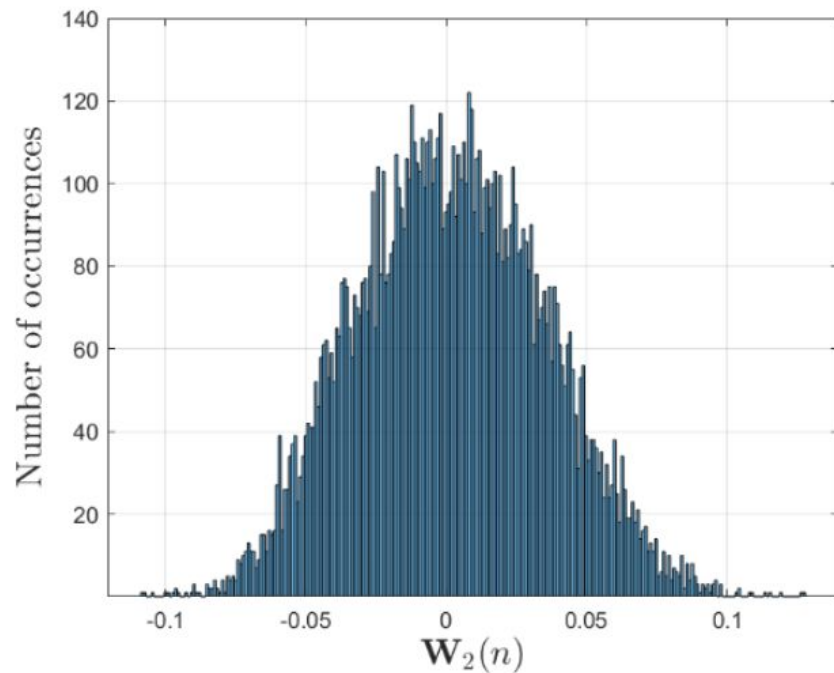


Quantização

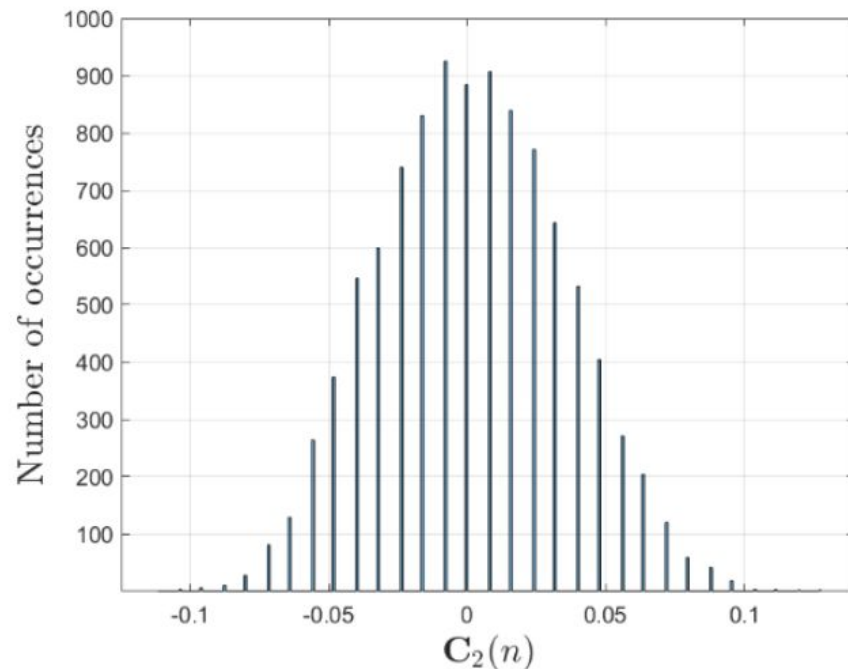
$$\mathbf{C}_k(n) = Q(\mathbf{W}_k(n), q_k) = \left\lceil \frac{\mathbf{W}_k(n)}{q_k} \right\rceil \times q_k$$

$$q_k = \frac{\max \{|\mathbf{W}_k(n)|\}}{2^{b-1} - 1}.$$

Quantização



(a) Before quantization.



(b) After quantization.

Formas de aplicar Pruning e Quantização

- No final do treinamento
- No final de cada época de treinamento
- No final de cada iteração de treinamento

Implementação

```
def pruning(model, gamma):  
    if gamma > 0.0:  
        with torch.no_grad():  
            betas = []  
            for p in model.parameters():  
                # transformando o layer em um vetor unidimensional  
                flattened_weights = torch.flatten(p)  
  
                # calculando o desvio padrao  
                std = torch.std(flattened_weights, -1)  
  
                # calculando o beta para a camada  
                beta = gamma * std  
                betas.append(beta)  
  
                # mascara booleana, que indica qual peso deveria ser mantido  
                mask = torch.gt(p.abs(), torch.ones_like(p) * beta)  
  
                # multiplicando os pesos da camada pelo tensor booleano  
                p.multiply_(mask)
```

Implementação

```
def quantization(model, b, betas=[]):
    if b >= 1:
        if len(betas) > 0:
            with torch.no_grad():
                for p, beta in zip(model.parameters(), betas):
                    flatted_weights = torch.flatten(p)

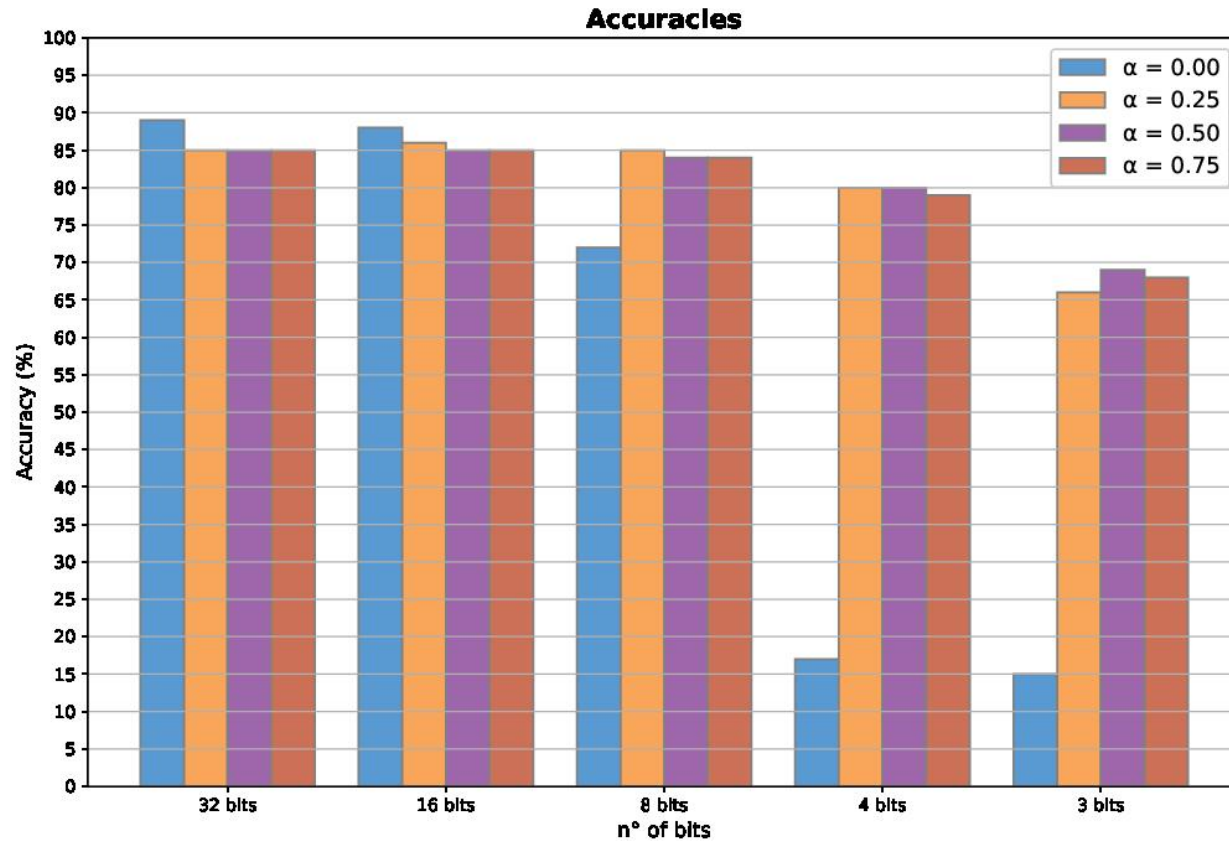
                    qk_prime = (torch.max(torch.abs(flatted_weights)) - beta) / ( (2**(b - 1)) - 1 )
                    torch.round(p/qk_prime, out=p)

                    p.multiply_(qk_prime)
        else:
            with torch.no_grad():
                for p in model.parameters():
                    flatted_weights = torch.flatten(p)
                    beta = 0.0;
                    qk_prime = (torch.max(torch.abs(flatted_weights)) - beta) / ( (2**(b - 1)) - 1 )
                    torch.round(p/qk_prime, out=p)
                    p.multiply_(qk_prime)
```

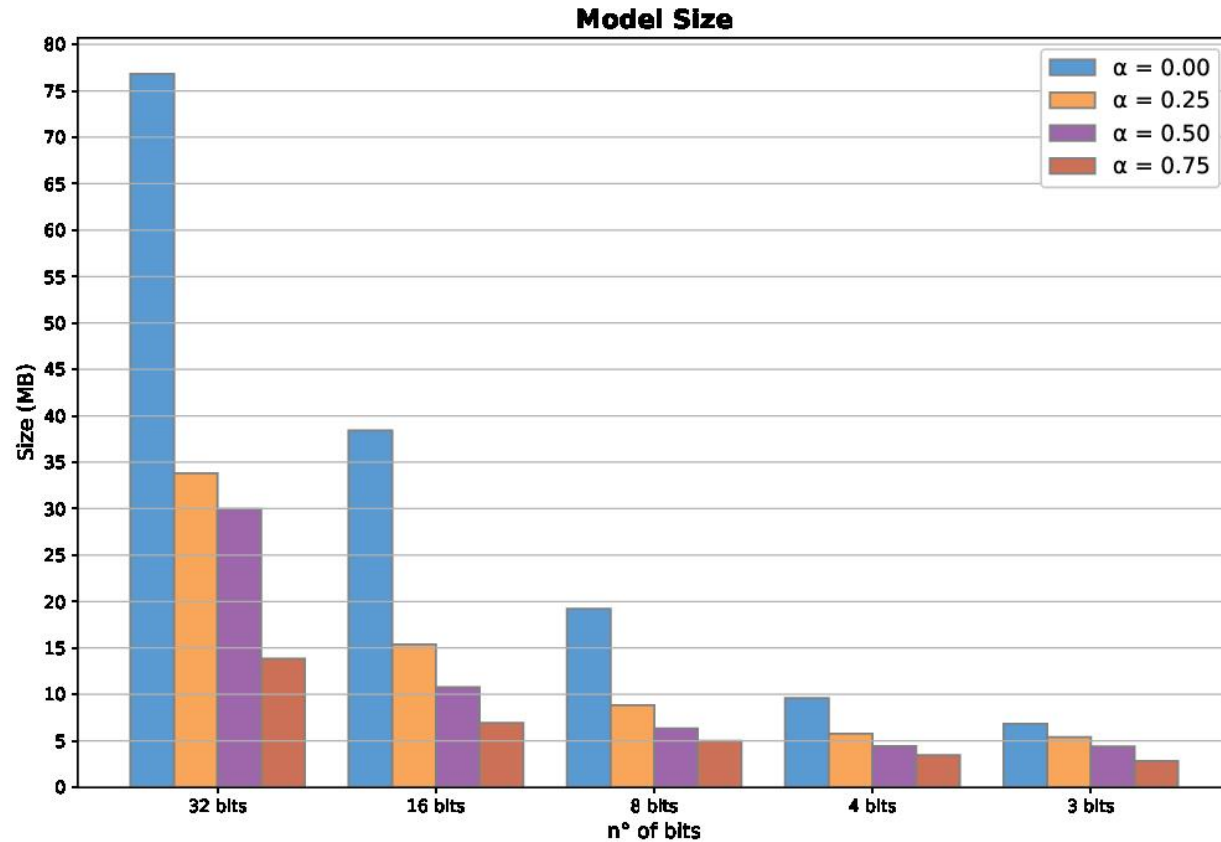
Dataset CIFAR10



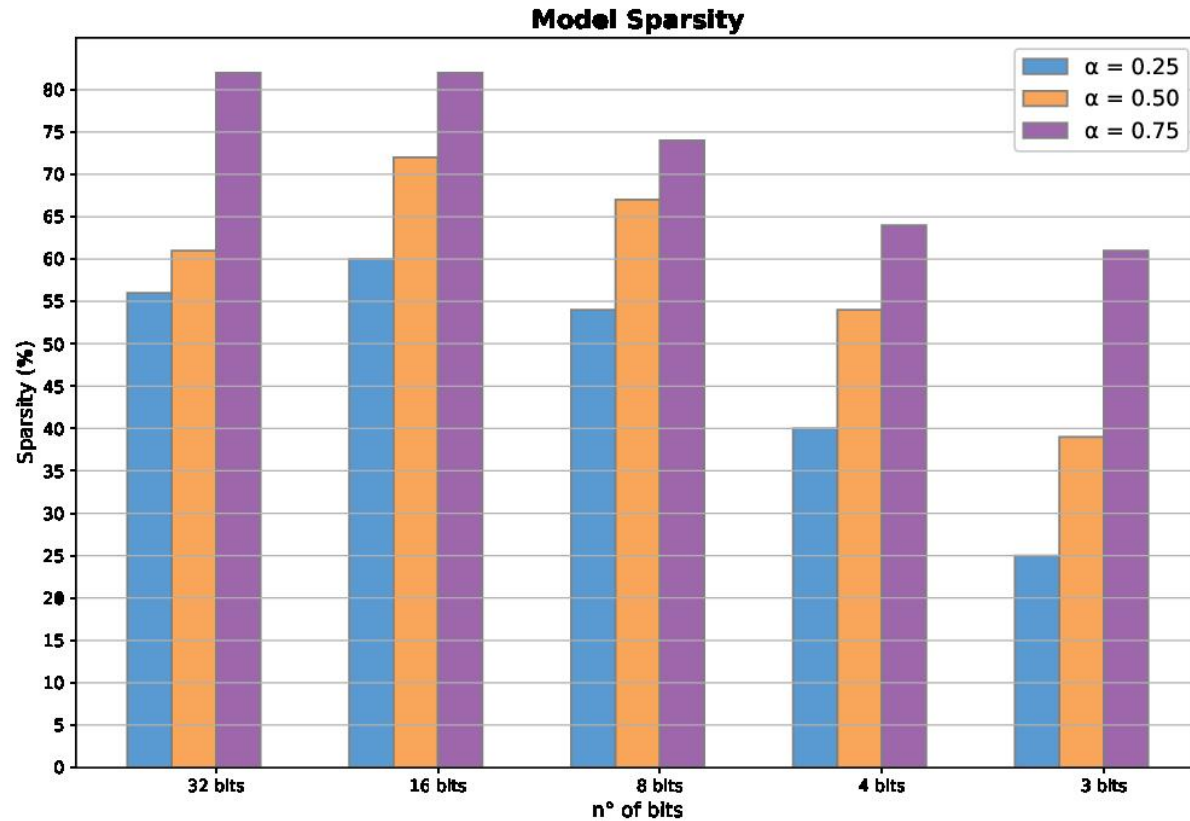
Resultados



Resultados



Resultados



Obrigado! 😁

Códigos e resultados disponíveis em:

https://github.com/vitoryeso/dl_model_compression