

Homework 2 (Linear Algebra) Problem 1

Read: A guide to convolution arithmetic for deep learning.

Solution: Complete



Homework 2 (Linear Algebra) Problem 2

- (a) What is the arithmetic intensity for matrix matrix multiplication with sizes
- M, N, K
- .

Solution:

$$\text{Compute} = MNK \quad (1)$$

$$\text{Spatial} = KN + MK + MN \quad (2)$$

When multiplying $M \times K$ and $K \times N$ matrices we will produce a $M \times N$ output matrix where each element was the result of a linear combination of K elements. As such, we must compute $M \times N$ sums of K elements. Similarly, for memory operations we must bring in the original matrices and output the resulting matrix, leading to a sum over the dimensions of these matrices. ■

- (b) Prove that arithmetic intensity for matrix matrix multiplication is maximized when
- $M = N = K$
- .

Proof: Choose N such that it is the largest matrix dimension. Let $I(M, N, K)$ designate the total arithmetic intensity. We can express M, K in terms of a factor of N .

$$M = C_m N \quad K = C_k N \quad C_m, C_k \in \mathbb{R}^+ \quad (3)$$

Using the result of part a we can express the total arithmetic intensity as a sum of the compute and memory components.

$$I(M, N, K) = MNK + (KN + MK + MN) \quad (4)$$

$$= C_m C_k N^3 + C_m C_k N^2 + C_m N^2 + C_k N^2 \quad (5)$$

Since N was chosen to be the largest of the terms, the coefficients C_m, C_k must be on the interval

$$0 < C_m, C_k \leq 1 \quad (6)$$

From the form of $I(M, N, K)$ given above and the interval on which C_m, C_k fall, it is clear that arithmetic intensity will be maximized for

$$C_m = 1 \quad C_k = 1 \quad (7)$$

So

$$M = C_m N = N \quad K = C_k N = N \quad (8)$$

$$\therefore M = N = K \quad (9)$$

□

Homework 2 (Linear Algebra) Problem 3

What is the complexity (MACs and memory) of a $N_o = N_i$ dense layer applied to vectorized versions of the following inputs:

(a) MNIST $1 \times 28 \times 28$

Solution:

$$\text{Compute} = (1 * 28 * 28)^2 = 6.15 \text{E}^6 \quad (10)$$

$$\text{Spatial} = (1 * 28 * 28)^2 + 2(1 * 28 * 28) = 6.16 \text{E}^6 \quad (11)$$

■

(b) CIFAR $3 \times 32 \times 32$

Solution:

$$\text{Compute} = (3 * 32 * 32)^2 = 9.44 \text{E}^6 \quad (12)$$

$$\text{Spatial} = (3 * 32 * 32)^2 + 2(3 * 32 * 32) = 9.44 \text{E}^6 \quad (13)$$

■

(c) ImageNet $3 \times 224 \times 224$

Solution:

$$\text{Compute} = (3 * 224 * 224)^2 = 2.27 \text{E}^{10} \quad (14)$$

$$\text{Spatial} = (3 * 224 * 224)^2 + 2(3 * 224 * 224) = 2.27 \text{E}^{10} \quad (15)$$

■

(d) Quasi 1/4 HD $3 \times 512 \times 1024$

Solution:

$$\text{Compute} = (3 * 512 * 1024)^2 = 2.47 \text{E}^{12} \quad (16)$$

$$\text{Spatial} = (3 * 512 * 1024)^2 + 2(3 * 512 * 1024) = 2.47 \text{E}^{12} \quad (17)$$

■

(e) Quasi HD $3 \times 1024 \times 2048$

Solution:

$$\text{Compute} = (3 * 1024 * 2048)^2 = 3.96 \text{E}^{13} \quad (18)$$

$$\text{Spatial} = (3 * 1024 * 2048)^2 + 2(3 * 1024 * 2048) = 3.96 \text{E}^{13} \quad (19)$$

■

Vectorizing our input (using L_d to denote feature map depth) gives us

$$N_i = L_d \times L_r \times L_c \quad (20)$$

The dense layer transforms an input to an output vector of the same length. As such, our transformation matrix will be N_i^2 . We are then left with matrix vector multiplication with the following dimensions under BLAS notation.

$$M = N_i \quad K = N_i \quad N = 1 \quad (21)$$

$$M = L_d \times L_r \times L_c \quad K = L_d \times L_r \times L_c \quad N = 1 \quad (22)$$

From the result of Problem 2 we have

$$\text{Compute} = MNK \quad \text{Spatial} = KN + MK + MN \quad (23)$$

$$= (L_d \times L_r \times L_c)^2 \quad = (L_d \times L_r \times L_c)^2 + 2 * (L_d \times L_r \times L_c) \quad (24)$$

Homework 2 (Linear Algebra) Problem 4

In practice, why can't you flatten a quasi HD input image to a vector for input to a dense layer

Solution: The results of Problem 3 illustrate how arithmetic intensity for vectorized image inputs to a dense layer grows rapidly, reaching tens of trillions of operations. This level of arithmetic intensity for large inputs is impractical on today's hardware. Furthermore, the use of matrix vector multiplication (as opposed to matrix matrix for CNN style 2D convolution layers) creates a constant factor ratio of compute to memory intensity. This leads to a memory wall, as memory movement operations are much slower than compute operations. It is also worth noting that dense layers acting on a vectorized 3D input will have no information about the topology of the true input, and will be forced to look for relationships between features that are separated by a great distance in the original input. ■

Homework 2 (Linear Algebra) Problem 5

Can a dense layer trained on an input of 1024×1 be applied to an input of size 2048×1 or 512×1

Solution: No to both. Dense layers learn weights such that a component in the output vector has a variable dependency on each of items in the input vector. The alterations needed to accommodate inputs of varying size would likely reduce the validity of the layer's output. ■

Homework 2 (Linear Algebra) Problem 6

Prove that CNN style 2D convolution with $N_o \times N_i \times F_r \times F_c$ filter, $N_i \times L_r \times L_c$ input and $N_o \times M_r \times M_c$ output can be lowered to the sum of $(F_r * F_c)$ matrix matrix multiplications.

Proof: Let the CNN style 2D convolution layer be of the usual form

$$Y^{3D} = H^{4D} X^{3D} + V^{3D} \quad (25)$$

If we consider a single output point $m_r \in M_r, m_c \in M_c$ of a single output feature map $n_o \in N_o$ we can express this as

$$Y(n_o, m_r, m_c) = \sum_{n_i}^{N_i} \sum_{f_r}^{F_r} \sum_{f_c}^{F_c} H(n_o, n_i, f_r, f_c) X(n_i, m_r + f_r, m_c + f_c) \quad (26)$$

If we expand the summation over n_i only we get a vector vector multiplication as follows

$$Y(n_o, m_r, m_c) = \sum_{f_r}^{F_r} \sum_{f_c}^{F_c} \begin{bmatrix} H(n_o, 0, f_r, f_c) \\ H(n_o, 1, f_r, f_c) \\ \dots \\ H(n_o, N_i, f_r, f_c) \end{bmatrix}^T \begin{bmatrix} X(0, f_r, f_c) \\ X(1, f_r, f_c) \\ \dots \\ X(N_i, f_r, f_c) \end{bmatrix} \quad (27)$$

We must next generalize this to all output points of all feature maps. If we extend row vector H into a matrix by adding N_o rows for $n_o \in N_o$, and if we extend column vector X into a matrix by adding $M_r \times M_c$ columns for the Cartesian product combinations of m_r, m_c , we get the following

$$Y = \sum_{f_r}^{F_r} \sum_{f_c}^{F_c} \underbrace{\begin{bmatrix} H(0, 0, f_r, f_c) & \dots & H(0, N_i, f_r, f_c) \\ \vdots & \ddots & \vdots \\ H(N_o, 0, f_r, f_c) & \dots & H(N_o, N_i, f_r, f_c) \end{bmatrix}}_{N_o \times N_i} \underbrace{\begin{bmatrix} X(0, f_r, f_c) & \dots & X(0, f_r + M_r, f_c + M_c) \\ \vdots & \ddots & \vdots \\ X(N_i, f_r, f_c) & \dots & X(N_i, f_r + M_r, f_c + M_c) \end{bmatrix}}_{N_i \times M_r * M_c} \quad (28)$$

From here we need only trivially reshape Y from a 2D matrix into a 3D tensor. The expression above indicates that CNN style 2D convolution can be lowered to a sum of $F_r * F_c$ matrix matrix multiplications. Furthermore, this result suggests that there are many alternative lowerings possible where each is determined by the which summations are expanded.

□

Homework 2 (Linear Algebra) Problem 7

How many MACs are required to compute each output point in a CNN style 2D convolution layer $N_o \times N_i \times F_r \times F_c$

Solution: In an intermediate step of Problem 6 we expressed a single output point as a summation as follows

$$Y(n_o, m_r, m_c) = \sum_{n_i}^{N_i} \sum_{f_r}^{F_r} \sum_{f_c}^{F_c} H(n_o, n_i, f_r, f_c) X(n_i, m_r + f_r, m_c + f_c) \quad (29)$$

Recall that this summation can be expanded to an inner product with vectors of length $N_i * F_r * F_c$. We know trivially that the number of MACs for this operation will be given by

$$\text{Compute}_1 = F_r * F_c * N_i \quad (30)$$

■

Homework 2 (Linear Algebra) Problem 8

How does CNN style 2D convolution complexity (MACs and memory) scale as a function of

(a) Product of image rows and columns ($L_r * L_c$)

Solution:

$$\text{Compute} \rightarrow \text{Linearly} \quad (31)$$

$$\text{Memory} \rightarrow \text{Linearly} \quad (32)$$

$$(33)$$

■

(b) Product of filter rows and columns ($F_r * F_c$)

Solution:

$$\text{Compute} \rightarrow \text{Linearly} \quad (34)$$

$$\text{Memory} \rightarrow \text{Linearly} \quad (35)$$

$$(36)$$

■

(c) Product of input and output feature maps ($N_o * N_i$)

Solution:

$$\text{Compute} \rightarrow \text{Linearly} \quad (37)$$

$$\text{Memory} \rightarrow \text{Linearly} \quad (38)$$

$$(39)$$

■

From previous problems we know that MACs are determined by

$$\text{Compute} = N_i * F_r * F_c * N_o * M_r * M_c \quad (40)$$

and memory is determined by

$$\text{Memory} = N_i L_r L_c + N_o M_r M_c + N_i N_o F_r F_c \quad (41)$$

Both of these

Homework 2 (Linear Algebra) Problem 9

Consider a CNN style 2D convolution layer with filter size $N_o \times N_i \times F_r \times F_c$.

(a) How many padding 0s such that output feature map is the same size as input?

Solution: Padding to maintain size across the CNN style 2D convolution layer is given by

$$P_l + P_r = F_c - 1 \qquad P_t + P_b = F_r - 1 \qquad (42)$$

This can easily be confirmed by realizing that

$$M_r = (L_r + P_r) - F_r + 1 \qquad (43)$$

$$(44)$$

So total padding along rows is given by

$$M_r = L_r \implies L_r = L_r + P_r - F_r + 1 \qquad (45)$$

$$\implies P_r = F_r - 1 \qquad (46)$$

■

(b) What is the size of the border of 0s for $F_r = F_c = 1$?

Solution: 0

There is no need for padding since a 1×1 filter will map each point in the input feature maps to a point in an output feature map. ■

(c) What is the size of the border of 0s for $F_r = F_c = 3$?

Solution:

$$P_l + P_r = 3 - 1 \qquad P_t + P_b = 3 - 1 \qquad (47)$$

$$P_l = P_r = 1 \qquad P_t = P_b = 1 \qquad (48)$$

$$(49)$$

■

(d) What is the size of the border of 0s for $F_r = F_c = 5$?

Solution:

$$P_l + P_r = 5 - 1 \qquad P_t + P_b = 5 - 1 \qquad (50)$$

$$P_l = P_r = 2 \qquad P_t = P_b = 2 \qquad (51)$$

$$(52)$$

■

Homework 2 (Linear Algebra) Problem 10

Consider a CNN style 2D convolution layer with filter size $N_o \times N_i \times F_r \times F_c$, input $N_i \times L_r \times L_c$, $P_r = F_r - 1$ and $P_c = F_c - 1$.

(a) What is the size of the output feature map with striding $S_r = S_c = 1$

Solution:

$$M_r = L_r \qquad M_c = L_c \qquad (53)$$

Using a modified form of Relationship 6 in the reading we can write

$$M_r = \left\lfloor \frac{L_r + P_r - F_r}{S_r} \right\rfloor + 1 \qquad (54)$$

$$= \left\lfloor \frac{L_r + (F_r - 1) - F_r}{S_r} \right\rfloor + 1 \qquad (55)$$

$$= \left\lfloor \frac{L_r - 1}{S_r} \right\rfloor + 1 \qquad (56)$$

Symmetry in the given input values allows this result to be extended to M_c . For a stride of one and padding given by filter dimension minus one, we get the expected result of $M_r = L_r$

■

(b) What is the size of the output feature map with striding $S_r = S_c = 2$

Solution:

$$M_r = \left\lfloor \frac{L_r - 1}{2} \right\rfloor + 1 \qquad M_c = \left\lfloor \frac{L_c - 1}{2} \right\rfloor + 1 \qquad (57)$$

We use the result of the previous part with $S_r = S_c = 2$

■

(c) How does this change the shape of the equivalent lowered matrix equation

Solution: In the equivalent lowered matrix operation, striding alters the input matrix of shape $N_i \times (M_r * M_c)$. Values M_r, M_c are divided by S_r, S_c respectively. This in turn alters the shape of the output matrix in a similar way. Intuitively, striding is a downsampling operation where the input feature map is roughly downsampled by a factor of S_r, S_c along L_r, L_c respectively.

■

Homework 2 (Linear Algebra) Problem 11

Can a CNN style 2D convolution layer trained on an input of size $3 \times 1024 \times 2048$ be applied to an input of size:

- (a) $3 \times 512 \times 1024$

Solution: Yes, there are upsampling techniques like deconvolution or interpolation that could be used.

This is the strength of CNN style 2D convolution layers. Since each output point is determined by input points exposed to the filter kernel at a given position, there will be no dependencies between distant points in the input feature map. Upsampling and downsampling can be applied to an input without compromising the validity of outputs from the CNN layer. ■

- (b) $3 \times 512 \times 512$

Solution: Again, upsampling techniques could be used, however care must be taken since $U_r \neq U_c$. ■

Homework 2 (Linear Algebra) Problem 12

In a standard RNN, if the state update matrix is constrained to a diagonal, what does this do for the mixing of the previous state with new inputs.

Solution: It forces an output element y_i to depend only on the value of y_i at time $t - 1$, rather than all outputs at time $t - 1$.

An RNN layer is given by

$$\mathbf{y}_t = \mathbf{H}\mathbf{X} + \mathbf{G}\mathbf{y}_{t-1} + \mathbf{V} \quad (58)$$

The state update matrix \mathbf{G} acts on vector \mathbf{y}_{t-1} which was produced by the previous RNN state. For non-diagonal \mathbf{G} this will produce a vector of length N where each element is a linear combination of all outputs from the previous state, as follows

$$\mathbf{G}\mathbf{y}_{t-1} = \begin{pmatrix} G_{00}y_0^{t-1} + G_{01}y_1^{t-1} + \cdots + G_{0N_o}y_{N_o}^{t-1} \\ G_{10}y_0^{t-1} + G_{11}y_1^{t-1} + \cdots + G_{1N_o}y_{N_o}^{t-1} \\ \vdots \\ G_{N_o0}y_0^{t-1} + G_{N_o1}y_1^{t-1} + \cdots + G_{N_oN_o}y_{N_o}^{t-1} \end{pmatrix} \quad (59)$$

Constraining \mathbf{G} to the diagonal means that a given output $y_i \in \mathbf{y}$ will be determined only be a scaled factor of $y_i \in \mathbf{y}_{t-1}$ as follows

$$\mathbf{G}\mathbf{y}_{t-1} = \begin{pmatrix} G_{00}y_0^{t-1} \\ G_{11}y_1^{t-1} \\ \vdots \\ G_{N_oN_o}y_{N_o}^{t-1} \end{pmatrix} \quad (60)$$

It is clear that in the addition of column vectors there will be no interplay between different output features. ■

Homework 2 (Linear Algebra) Problem 13

The size of the input to the global average pooling layer is $1024 \times 16 \times 32$.

(a) What is the size of the output

Solution: 1024×1

The global average pooling layer computes a global average for each of the input feature maps. As such, we will average feature maps of size 16×32 to a single value, repeated for each of the 1024 given feature maps. ■

(b) What is the complexity (MACs) of the layer

Solution:

$$\text{Compute} = N_i * L_r * L_c \quad (61)$$

$$= 1024 * 16 * 32 \quad (62)$$

$$= 524288 \quad (63)$$

$$(64)$$

Global average pooling can be thought of as a matrix vector multiplication with the following BLAS dimensions

$$M = N_i \quad K = L_r * L_c \quad N = 1 \quad (65)$$

where the $K \times N$ matrix has each value set to $1/(L_r * L_c)$ for averaging. Based on this interpretation we can determine the complexity of the layer. ■