

Vision Transformer

Procesamiento Avanzado de Imágenes

Alumno: Sebastián Tinoco
Profesor: Javier Ruiz del Solar
Auxiliar: Patricio Loncomilla
Fecha de entrega: 19 de diciembre de 2022
Santiago de Chile

Índice de Contenidos

1. Introducción	1
2. Marco Teórico	2
2.1. Transformers	2
2.2. Vision Transformer	3
3. Datos	4
4. Implementación	5
5. Entrenamiento	6
6. Conclusiones	11
Referencias	12
Anexo A. Código implementado	13
A.1. Pregunta 1	13

Índice de Figuras

1. Arquitectura de Transformer	2
2. Mecanismo de Atención	3
3. Arquitectura de ViT	4
4. Distribución de clases de Yoga-82	5
5. Entrenamiento del modelo en el nivel superior en las primeras 25 épocas	7
6. Entrenamiento del modelo en el nivel superior en las últimas 39 épocas	7
7. Entrenamiento del modelo en el nivel intermedio	8
8. Entrenamiento del modelo en el nivel intermedio con <i>Fine Tunning</i>	9
9. Entrenamiento del modelo en el nivel inferior	9
10. Entrenamiento del modelo en el nivel inferior con <i>Fine Tunning</i>	10

Índice de Tablas

1. Accuracy en Entrenamiento y Validación por jerarquía de clases	10
---	----

Índice de Códigos

A.1. Generar patches	13
A.2. Multi-Head Self Attention	13
A.3. Multi Layer Perceptron	14
A.4. Bloque de Encoder	15
A.5. Modelo ViT	15

1. Introducción

La visión computacional ocupa un pilar fundamental tanto a nivel industrial como en la cotidianidad de las personas. Desde aplicaciones como el procesamiento de imágenes médicas (por ejemplo, detección de tumores) hasta el reconocimiento facial como medida de seguridad, el impacto de esta área ha ido creciendo cada vez más a través del tiempo, a tal punto que resulta difícil imaginar un escenario donde no se disponga de esta tecnología.

Con tal de dar continuidad al aprendizaje y entendimiento de esta disciplina, el presente trabajo tiene por objetivo implementar la red neuronal **Vision Transformer** (ViT) con Pytorch, un framework deep learning para Python. De esta forma, se buscará entrenar esta red para resolver la tarea de clasificación de imágenes en el dataset Yoga-82.

Para cumplir esto de forma efectiva, se divide el trabajo en diferentes apartados. La sección 1 introduce de forma general el objetivo del informe y da una idea de lo que se espera en este. La sección 2 contiene la descripción general del modelo, dando especial énfasis en las capas de atención y las diferencias de este modelo con otros modelos de tipo transformer. La sección 3 contiene una breve caracterización del dataset Yoga-82. La sección 4 detalla la implementación del modelo en Pytorch y cada uno de sus componentes. La sección 5 presenta los resultados de entrenar el modelo con el set de datos. Finalmente, la sección 6 finaliza el documento con los principales hallazgos del trabajo realizado, denotando los principales aprendizajes y dificultades.

2. Marco Teórico

2.1. Transformers

Los modelos *transformers* son un tipo de arquitectura de red neuronal la cual tiene la particularidad de implementar el mecanismo de *self-attention*, lo que ha permitido lograr importantes avances en diferentes áreas relevantes a la ciencia de datos. Esta arquitectura y mecanismo fueron introducidos formalmente en la investigación “*Attention is all you need*” [1], marcando un antes y un después en el los modelos basados en aprendizaje supervisado. A continuación, se presenta la arquitectura general de un Transformer:

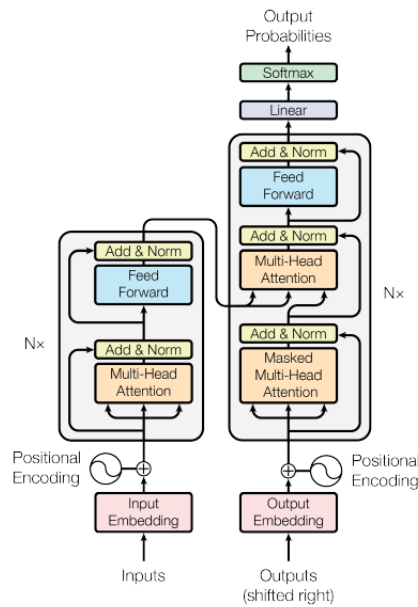


Figura 1: Arquitectura de Transformer

Para explicar esta arquitectura, es necesario distinguir sus dos componentes principales, *Encoder* y *Decoder*:

- *Encoder* (dibujo izquierdo): Su principal función es obtener una mejor representación de los datos de entrada, obteniendo como salida un arreglo de datos de la misma dimensión que la entrada. Para lograr lo anterior, se proyectan los datos de entrada a su representación en *embeddings*, se suma el encoding posicional, para luego aplicar N capas de *Multi-Head Attention* y *Feed Forward*, con una capa de *Layer Normalization* tras cada una de las operaciones anteriores y *skip connection* cada 2 capas.
- *Decoder* (dibujo derecho): De forma similar al *Encoder*, se utilizan embeddings posicionales, capas de *Multi-Head Attention*, *Layer Normalization* y *skip connection*, con la gran diferencia de que recibe los datos “ground truth” para aplicar el mecanismo de atención¹. Luego, se combina esta información con la salida del *Encoder*, volviendo a aplicar el mecanismo de atención sobre

¹ Donde se hace de vital importancia diferenciar entre los datos disponibles y los datos futuros a través de *Masked Multi-Head Attention*

este conjunto de datos. Finalmente, se proyectan los datos al espacio de salida y se aplica *Softmax* para devolver las probabilidades de cada clase.

Por otro lado, el mecanismo de atención queda descrito por la siguiente ecuación:

$$Attention(Q, K, V) = \frac{QK^T}{\sqrt{d_k}} \quad (1)$$

Donde las matrices Q , K y V corresponden a proyecciones de los datos en un espacio de dimensión d_k . De esta forma, *Multi-Head Attention* corresponde a aplicar la ecuación anterior M veces, concatenar estos resultados y luego aplicar volver a la dimensión original de los datos de entrada a través de una proyección lineal. La Figura 2 presenta la ilustración de este mecanismo:

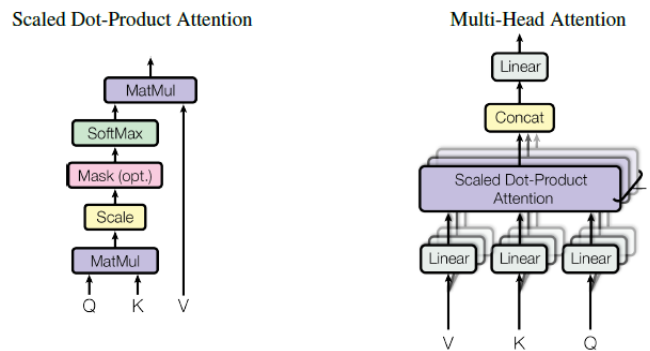


Figura 2: Mecanismo de Atención

De esta forma, se observa que el mecanismo de atención logra capturar el *contexto* de los datos a través del producto punto, lo que induce a resaltar los *tokens* con fuerte relación entre sí.

Para complementar la discusión y el entendimiento de *self-attention*, se presentan algunas ventajas y desventajas de este mecanismo con respecto al operador convolución. Entre ellas, se encuentran:

- Ventajas de *self-attention*:
 - Logra capturar el “contexto” de cada input asignando pesos distintos a los datos dependiendo de la información de entrada
 - Es capaz de identificar la relación entre los datos (en este caso, entre los patches) de forma simultánea
- Desventajas de *self-attention*:
 - Es mas costosa de computar debido a su gran cantidad de parámetros, aumentando el tiempo de entrenamiento de manera significativa.
 - Requiere de una cantidad mucho mayor de datos para converger a buenos resultados

2.2. Vision Transformer

El modelo Vision Transformer es un modelo de red neuronal basado en la arquitectura Transformer y sus mecanismos de *self-attention*. Este modelo presentado formalmente en la investigación *An image*

is worth 16×16 words: Transformers for Image Recognition at scale [2] como una alternativa a las redes neuronales convolucionales para la clasificación de imágenes. A diferencia de los Transformers tradicionales, este modelo usa sólo el *Encoder* de la arquitectura, utilizando una capa *Feed-forward* para hacer la clasificación. La Figura 3 presenta la arquitectura de este modelo:

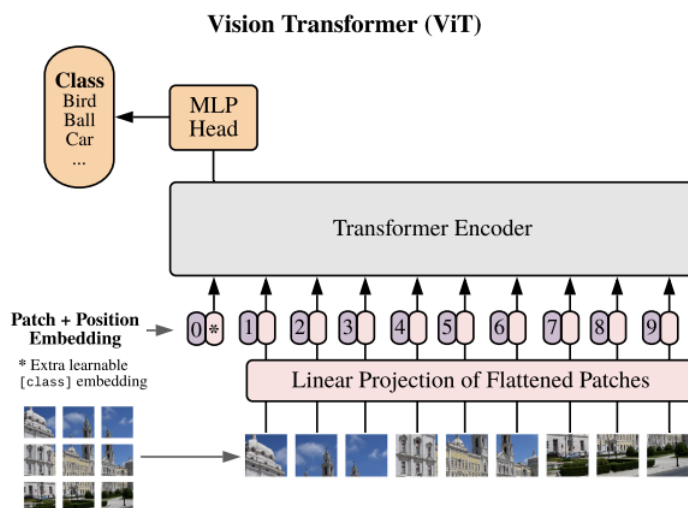


Figura 3: Arquitectura de ViT

Recordando la sección anterior, para hacer uso del mecanismo de atención es necesario contar con secuencias de entrada para así rescatar las relaciones importantes entre sus elementos. Para esto, el modelo ViT requiere transformar la imagen en *patches*, es decir, secciones de una imagen de un tamaño definido (usualmente 16×16). Estas secciones son después proyectadas en una capa lineal, donde además se le suma el encoding posicional. Es importante añadir que en esta etapa también se añade el *token* de clasificación, el cual se utiliza después para realizar la clasificación. Habiendo realizado todo lo anterior, se transforma la información mediante múltiples capas de *Multi-Head Attention* (similar a lo explicado en el apartado de Transformers), para finalmente utilizar el *token* de clasificación como *input* para la capa lineal de clasificación.

Finalmente, es importante mencionar que el modelo ViT tiene dificultades para manejar imágenes con alta resolución debido a que el costo computacional y de memoria es cuadrático debido al mecanismo de atención, específicamente en la multiplicación QK^T de la ecuación 1. Para resolver esto, los autores redimensionan las imágenes a 224×224 , permitiendo que los datos de entrada sean computacionalmente tratables en la gran mayoría de las máquinas virtuales a usar.

3. Datos

El dataset Yoga-82 fue presentado en la investigación “Yoga-82: A New Dataset for Fine-grained Classification of Human Poses” [3] con el propósito de contribuir al avance en la inferencia de posiciones humanas mediante procesamiento de imágenes. De esta forma, se publica este set de datos para ser utilizado como benchmark del ajuste de diferentes modelos a la clasificación de las diferentes

posiciones de yoga (en otras palabras, posiciones humanas). El dataset cuenta con 28.450 imágenes etiquetadas en 82 posiciones distintas de yoga, disponiendo desde 64 hasta 1.133 imágenes por cada clase. Una particularidad del dataset es que cuenta con imágenes etiquetadas en 3 niveles (jerarquías): Primer nivel (6 clases), Segundo nivel (20 clases) y Tercer nivel (82 clases). De esta manera, se buscará usar el modelo *Vision Transformer* para clasificar las imágenes en las 3 jerarquías nombradas.

Un primer problema encontrado fue que no todas las imágenes estaban disponibles para ser utilizadas, explicándose tanto en enlaces de descarga caídos o que las imágenes descargadas contenían información corrupta que no podía ser leída. De esta manera, del total de 28.450 imágenes originales postuladas en [3], sólo se pudieron utilizar 20.264. A modo de comparación, la Figura 4 presenta la distribución de las clases en el dataset haciendo la distinción entre el dataset original y las imágenes habilitadas para ser usadas:

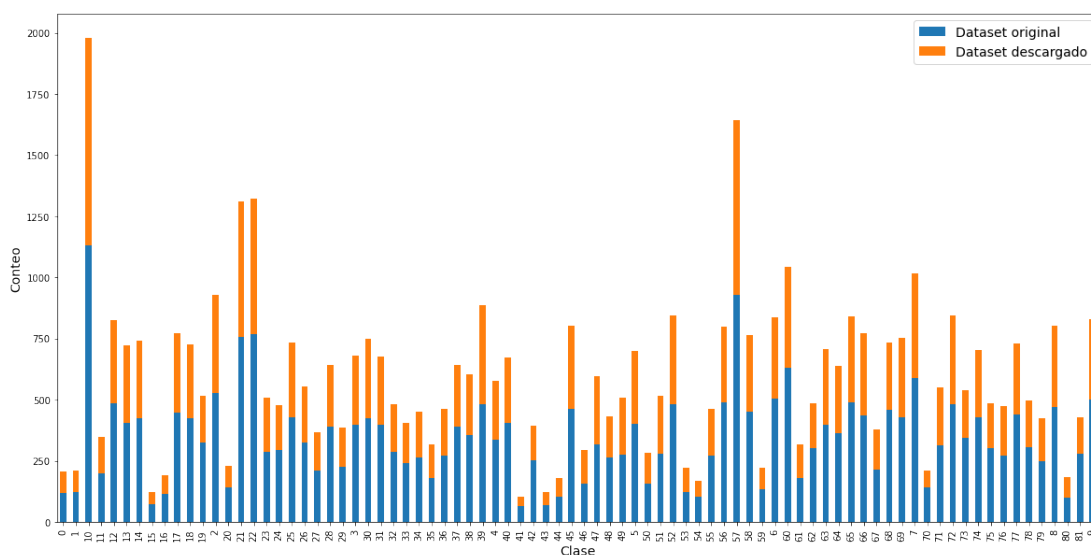


Figura 4: Distribución de clases de Yoga-82

De esta manera, se observa que la distribución de las clases se mantiene relativamente estable, por lo que no sería necesario incurrir en técnicas de rebalanceo de clases.

En cuanto a la distribución de conjuntos de entrenamiento, validación y test, se debió crear el conjunto de validación usando el 20 % de los datos de entrenamiento, manteniendo la distribución de clases vista en 4. De esta forma, el conjunto de entrenamiento quedó conformado por 12.094 observaciones (59.68 %), el conjunto de validación por 3.024 observaciones (14.92 %) y el conjunto de test por 5.146 observaciones (25.39 %).

4. Implementación

Para implementar el modelo ViT se trabajó de forma modular (es decir, dividiendo la implementación en diferentes módulos o clases) para facilitar la comprensión y el diseño de la arquitectura del transformer. En concreto, se diseñaron 5 módulos distintos:

- **gen_patches**: Función que recibe imágenes, devuelve las mismas imágenes segmentadas en

patches cuadrados (ver Anexo A.1).

- **MultiHeadSelfAttention:** Módulo que implementa el “Multi-Head Self Attention” de la arquitectura transformers (ver Anexo A.2).
- **MultiLayerPerceptron:** Módulo que implementa la capa de MLP referente al bloque de Encoder de ViT transformer (ver Anexo A.3).
- **EncoderBlock:** Módulo que implementa un bloque o “capa” del Encoder, referente al modelo ViT transformer. Se apoya de las clases *MultiHeadSelfAttention* y *MultiLayerPerceptron* (ver Anexo A.4)
- **ViT:** Módulo que implementa el modelo ViT transformer. Se apoya de la clase *EncoderBlock* y la función *gen_patches* (ver Anexo A.5).

La totalidad de este código puede ser encontrado en el siguiente [repositorio de github](#).

5. Entrenamiento

Entrenar el modelo *Vision Transformer* con el dataset Yoga-82 no fue una tarea fácil, especialmente considerando que los recursos computacionales² disponibles eran limitados. Dado lo anterior, las principales dificultades encontradas en el entrenamiento del modelo fueron el tamaño del dataset, la gran cantidad de parámetros del modelo y la resolución de algunas imágenes, cayendo reiteradas veces en problemas de memoria RAM y/o de agotar el tiempo máximo³ para utilizar GPU. Para solucionar lo anterior, se abordó el entrenamiento con tres enfoques:

- Se utilizó la arquitectura de *Vision Transformer* en su versión minimalista (**ViT-Ti/16**), la cual cuenta con 12 capas de Encoder, embeddings de dimensión 192, capas lineales MLP de tamaño 768 y 3 cabezas, reduciendo así la cantidad de parámetros. Además, trabaja sobre *patches* de 16x16 píxeles.
- Para resolver el problema del tiempo máximo de ejecución, se incurrió en el guardado automático de los pesos en el disco duro (en este caso, *Google Drive*) cada vez que se obtenía un *accuracy* mayor en el conjunto de validación. De esta manera cada vez que se agotaba el tiempo máximo de ejecución en una cuenta, se cargaba el *checkpoint* en otra cuenta con tiempo disponible y se seguía entrenando.
- Con el fin de obtener mayor eficiencia en el uso del tiempo disponible de GPU, se redimensionaron todas las imágenes del dataset a 256x256 antes del entrenamiento, logrando disminuir el tiempo de procesamiento por iteración a más de la mitad. De esta forma, la única transformación implementada por la librería *torchvision* fue el *crop* a 224x224, volviendo a las dimensiones propuestas por [2].

En cuanto a los hiperparámetros y técnicas de regularización utilizadas, se siguió la metodología expuesta en “*How to train your ViT? Data, Augmentation, and Regularization in Vision Transformers*” [4], donde se utiliza un optimizador *Adam* con *learning rate* de 0.0002, $\beta_1 = 0.9$ y $\beta_2 = 0.999$.

² Se usó *Google Colab* para el entrenamiento de los modelos

³ Google Colab no publica estos límites dado que puede variar en función de la intensidad de las operaciones. Para el caso del trabajo, Google Colab hacía disposición de casi 2 horas de GPU diarias por cuenta.

El tamaño de cada batch se fijó en 256 pues fue la potencia de 2 más grande que podía sostener la memoria RAM otorgada. En cuanto a técnicas de regularización, se utiliza *RandAugment* con $\alpha = 2$ para realizar *data augmentation*, *gradient clipping* de 5 como norma máxima del gradiente, *scheduler* lineal y dropout de 0.1 entre las capas lineales del *Encoder*. Finalmente, se utilizó la técnica *Early Stopping* con *patience* de 12 épocas para evitar el sobreajuste del modelo a los datos de entrenamiento, lo que tuvo como consecuencia que el modelo se entrenara entre 20 y 40 épocas por cada ejecución.

Con la finalidad de realizar experimentos con *transfer learning* entre modelos, se decide empezar entrenando con las clases del nivel superior debido a su menor número de clases (sólo 6) y así facilitar el ajuste del modelo. A continuación, se presentan las curvas de *loss* y *accuracy* para la primera y última ejecución⁴ del modelo:

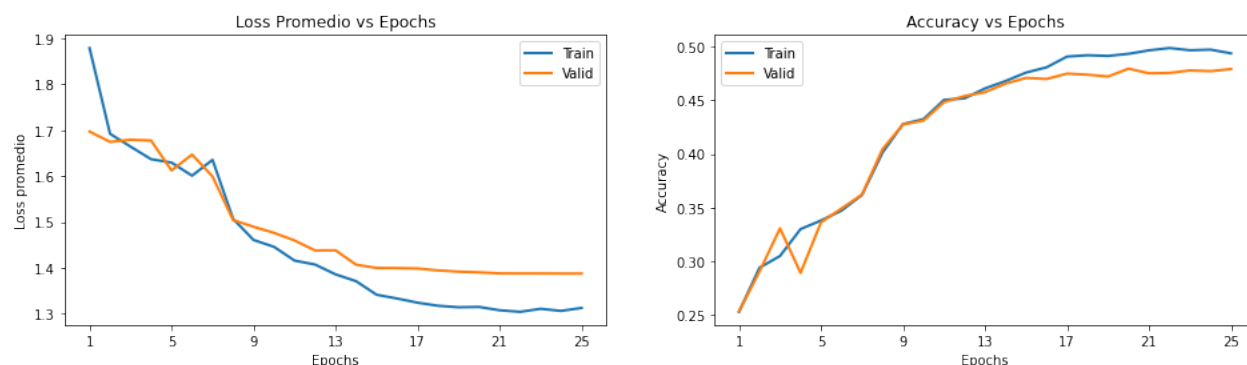


Figura 5: Entrenamiento del modelo en el nivel superior en las primeras 25 épocas

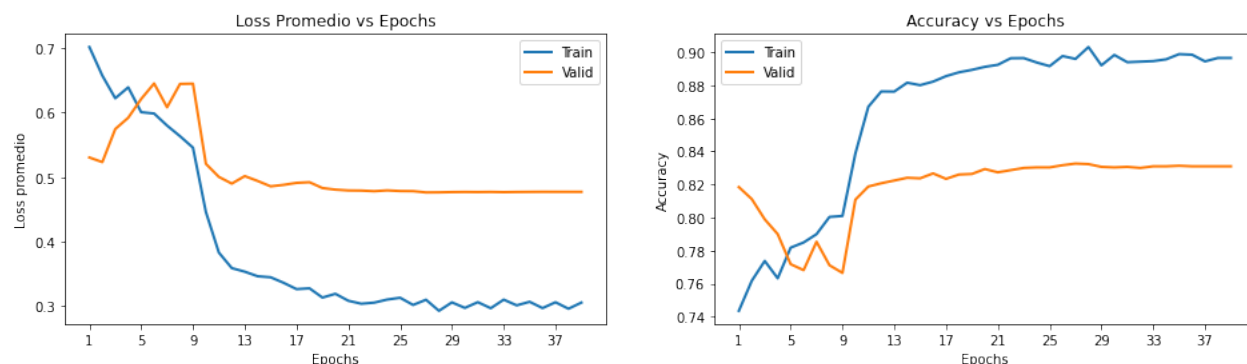


Figura 6: Entrenamiento del modelo en el nivel superior en las últimas 39 épocas

De la Figura 5, se observa que el modelo aprende rápidamente de los datos, posicionándose en el tramo de 40-50 % de *accuracy* en las primeras 10 épocas. Para el resto de las épocas, se observa como el modelo desacelera su aprendizaje, aprendiendo apenas pasar el 45 % de *accuracy*. Este fenómeno es algo que se observó en todos los experimentos realizados y se puede explicar en el *scheduler* lineal usado, el cual reduce el *learning rate* del optimizador cada cierto número de iteraciones, explicando

⁴ Por disponer de recursos computacionales limitados, se debió entrenar el modelo por medio de *checkpoints* y en diferentes ejecuciones

como el modelo casi no obtiene mejoras en su ajuste para las últimas épocas. Otro punto interesante es que debido a las técnicas de regularización aplicadas, el sobreajuste de los datos se mantiene en niveles mínimos, lo que se evidencia en la cercanía entre las curvas de entrenamiento y validación.

Por otro lado, se observa como en la Figura 6 el modelo encuentra dificultades para obtener mejoras en su rendimiento, logrando mejoras apenas notables en el conjunto de validación. Además, se observa como las curvas de *loss* y *accuracy* comienzan desde un punto con mayor ajuste a los datos, lo que es fruto de haber entrenado este modelo bajo el enfoque de *checkpoints*. Por último, se estima⁵ que el modelo se entrenó con 120 épocas para alcanzar el ajuste observado en esta Figura.

Con la red entrenada, se procede a evaluar su desempeño en los conjuntos de validación y test, usando el modelo con mejor ajuste en validación conseguido durante la fase de entrenamiento. De esta manera, se consigue un 83.27 % de *accuracy* en el conjunto de validación, en contraste con el 65.66 % de *accuracy* en el conjunto de test.

Se repite el experimento con las clases del nivel intermedio (20 clases en total). A diferencia del experimento anterior (y considerando el tiempo y recursos limitados para implementar la totalidad del proyecto) el modelo sólo se entrena una ejecución, alcanzando las 37 épocas efectuadas sin hacer uso de *checkpoints*. La Figura 7 presenta las curvas de *loss* y *accuracy* de este experimento:

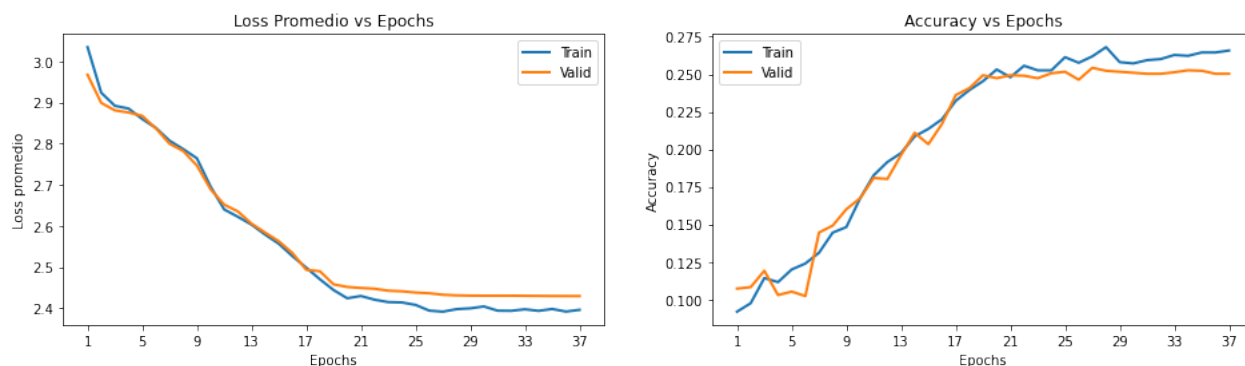


Figura 7: Entrenamiento del modelo en el nivel intermedio

Al igual que en el experimento anterior, se observa como el modelo posee poco sobreajuste, lo que se ilustra en la cercanía de las curvas de entrenamiento y validación. Además, se observa como el mayor número de clases impacta de manera negativa en el *accuracy* del modelo, el cual se estabiliza en el 25-27 % en las últimas épocas del entrenamiento. Esto último se puede explicar en que un mayor número de clases aumenta la complejidad del problema, haciendo que sea mas difícil para el modelo distinguir entre una clase y otra.

Con el motivo de entender mejor el funcionamiento de la red *Vision Transformer* y su ajuste a los datos disponibles, se utiliza *transfer learning* del modelo entrenado con las clases de nivel superior con la intención de lograr un ajuste mayor a los datos de las clases del intermedio. En concreto, se hace uso de la técnica *fine tuning*, es decir, entrenar empezando desde los pesos pre entrenados, con la salvedad de cambiar la dimensión de salida de la red (en este caso, pasar de predecir 6 a 20 clases).

⁵ Se perdió el número exacto de épocas efectuadas resultado de aplicar *Early Stopping* y el enfoque de *checkpoints*

La Figura 8 presenta los resultados de haber realizado este experimento:

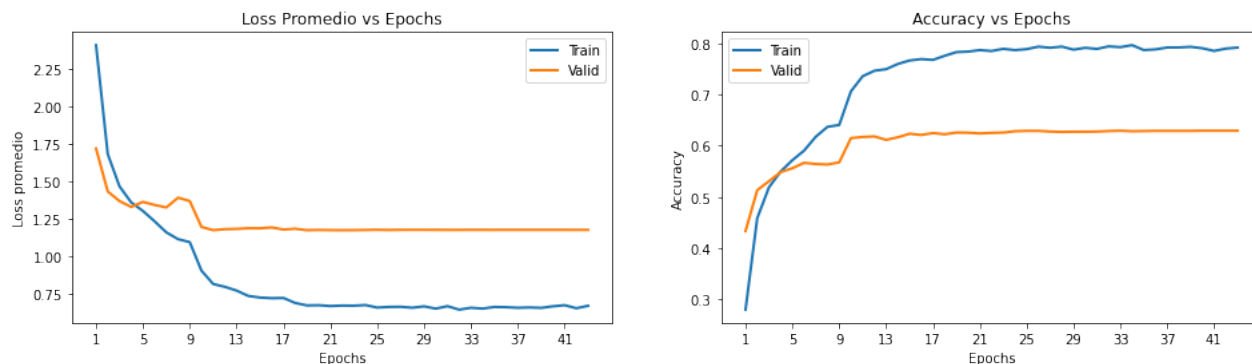


Figura 8: Entrenamiento del modelo en el nivel intermedio con *Fine Tuning*

A partir de los resultados, se observa como la red obtiene un mejor desempeño de forma inmediata, alcanzando el 60 % de *accuracy* en el conjunto de validación en las primeras 10 épocas. Esta diferencia con el entrenamiento “desde cero” es notable, pues es prueba de que el *transfer learning* es una técnica efectiva para alcanzar un buen ajuste en un menor número de épocas. Finalmente, se evalúa el ajuste del algoritmo en los datos de validación y test, consiguiendo un *accuracy* de 62.96 % en los datos de validación y 53.46 % en los datos de test.

La última ola de experimentos fue intentar clasificar las posiciones de yoga de las clases del nivel inferior (82 clases en total). Similar al experimento anterior de 20 clases, el modelo se entrena sólo una ejecución, alcanzando las 34 épocas efectuadas sin hacer uso de *checkpoints*. La Figura 9 presenta las curvas de *loss* y *accuracy* de este experimento:

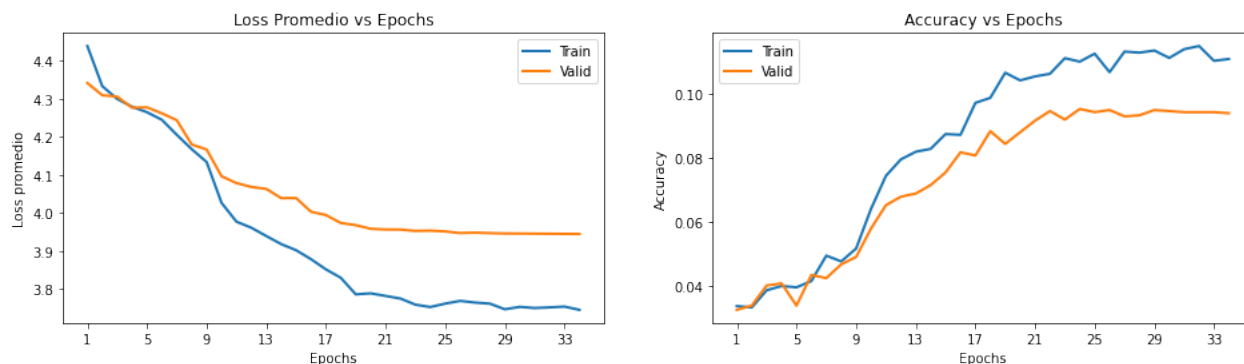


Figura 9: Entrenamiento del modelo en el nivel inferior

Similar a los dos experimentos anteriores donde el modelo se entrena “desde cero”, se observa como los niveles de sobreajuste se mantienen en el mínimo. Además, se observa nuevamente como un mayor número de clases impacta de forma negativa sobre el *accuracy* del modelo, el cual se estabiliza entre 9-10 % en el conjunto de validación en las últimas épocas de entrenamiento.

Por otro lado, se repite el experimento de usar los pesos del modelo entrenado con las clases del nivel superior para hacer *transfer learning* (en específico, *fine tuning*) y así obtener un mayor rendimiento en la predicción de las 82 clases. La Figura 10 presenta los resultados de este experimento:

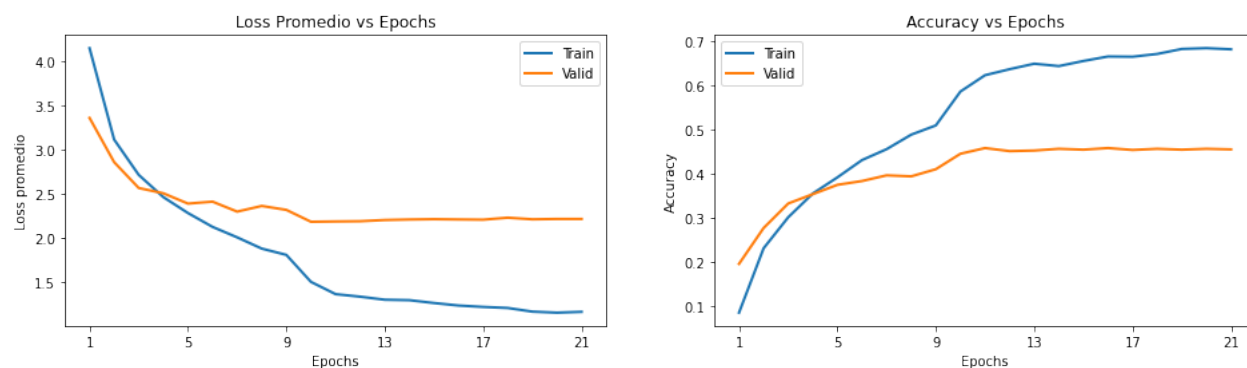


Figura 10: Entrenamiento del modelo en el nivel inferior con *Fine Tunning*

Al igual que en el caso anterior, se observa como la red obtiene un mejor desempeño de forma inmediata, posicionando su *accuracy* sobre el 40 % en las primeras 13 épocas. Es interesante notar también que el aprendizaje de la red se estanca rápidamente, fenómeno que puede ser explicado en la red habiendo extraído todo el conocimiento “útil” de las clases superiores en las primeras épocas, para luego verse en la obligación de aprender de los datos para seguir mejorando su ajuste.

La Tabla 1 resume los experimentos realizados a través de esta sección:

Tabla 1: Accuracy en Entrenamiento y Validación por jerarquía de clases

Clases	Transfer Learning	Validación	Test
6	No	83.27 %	65.66 %
20	Sí	62.96 %	53.46 %
82	Sí	45.77 %	42.54 %
20	No	25.43 %	25.77 %
82	No	9.52 %	9.09 %

6. Conclusiones

La visión computacional tiene un rol fundamental en el funcionamiento de la sociedad, tanto en la cotidianidad de las personas como a nivel industrial. En el trabajo expuesto fue posible dilucidar el diseño y entrenamiento de la red neuronal *Vision Transformer* para resolver el problema de clasificación de imágenes del dataset Yoga-82 por medio de Pytorch. Además, se estudió el impacto de realizar *transfer learning* sobre el rendimiento en datos dentro y fuera de la muestra.

En la sección 2 se da una explicación del funcionamiento del modelo *Vision Transformer*, haciendo un contraste del mecanismo de atención con la operación de convolución y explicando también el problema con la dimensión de los *patches* de entrada. En la sección 3 se da una explicación breve de la composición del dataset Yoga-82. En la sección 4 se explican los principales módulos programados para implementar el modelo *Vision Transformer* en *Pytorch*. En la sección 5 se muestran los resultados de haber entrenado el modelo con el dataset Yoga-82.

En vista y función de los resultados obtenidos a lo largo de este informe, se observa como el modelo *Vision Transformer* obtiene mejores predicciones al utilizar *transfer learning* entre modelos y jerarquías. Este hallazgo va acorde con lo esperado y hace sentido con la teoría, pues utilizar redes pre entrenadas es una forma ampliamente usada para aumentar el rendimiento de los modelos, en especial si los pesos provienen de una red entrenada con la misma data. Además, se observa como la red disminuye su ajuste a medida que aumenta el número de clases a predecir. Este resultado también va acorde a lo esperado, pues un mayor número de clases aumenta la complejidad del problema, siendo mas difícil para el modelo diferenciar entre las clases y hacer predicciones precisas.

En consideración del ajuste de la red *Vision Transformer* a los datos en donde se logra predecir de forma parcial las distintas jerarquías de clases, se concluye la obtención de resultados aceptables en el entrenamiento del modelo.

A pesar de haber implementado con éxito el modelo, el trabajo expuesto no estuvo exento de dificultades. La principal dificultad estuvo en lidiar con los recursos computacionales limitados para entrenar el modelo, donde se tuvo que recurrir al entrenamiento fragmentado por medio de *checkpoints*. Otra complejidad encontrada fue entender la compleja arquitectura y mecanismos subyacentes del modelo para su implementación.

En retrospectiva del trabajo realizado, se consolida el aprendizaje de el diseño y entrenamiento de la red *Vision Transformer* mediante código Python y la librería Pytorch.

Finalmente, queda como trabajo propuesto entrenar el modelo con recursos suficientes para aumentar la cantidad de parámetros del modelo y prescindir del enfoque de *checkpoints*. Se espera que con la inclusión de estos cambios se disponga de un mayor tiempo de entrenamiento y así conseguir un mejor ajuste a los datos.

Referencias

- [1] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., y Polosukhin, I., “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [2] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” *arXiv preprint arXiv:2010.11929*, 2020.
- [3] Verma, M., Kumawat, S., Nakashima, Y., y Raman, S., “Yoga-82: a new dataset for fine-grained classification of human poses,” en *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pp. 1038–1039, 2020.
- [4] Steiner, A., Kolesnikov, A., Zhai, X., Wightman, R., Uszkoreit, J., y Beyer, L., “How to train your vit? data, augmentation, and regularization in vision transformers,” *arXiv preprint arXiv:2106.10270*, 2021.

Anexo A. Código implementado

A.1. Pregunta 1

Código A.1: Generar patches

```

1 import torch
2
3 def gen_patches(images, patch_size: int = 16):
4
5     """
6     Función que recibe imágenes, devuelve las mismas imágenes segmentadas en patches cuadrados.
7     images: iterable con imágenes
8     patch_size: tamaño de cada patch (int)
9     """
10
11     batch_size, C, H, W = images.shape # dimensiones de la imagen
12
13     assert H == W # assert imagenes cuadradas
14     assert H % patch_size == 0 # assert imagenes divisibles por patch_size
15
16     n_patches = H // patch_size # cantidad de patches en una arista
17
18     output = [] # lista de imagenes en formato batches
19     for img in images:
20         patches = [] # batches de una imagen
21         for row in range(n_patches):
22             for col in range(n_patches):
23                 patches.append(img[:, patch_size * row : patch_size * (row + 1), patch_size * (col) : patch_size * (col + 1)].flatten())
24                 ↪ # tensor (1, patch_size)
25
26         output.append(torch.stack(patches))
27
28     output = torch.stack(output)
29     return output

```

Código A.2: Multi-Head Self Attention

```

1 import numpy as np
2 import torch
3 import torch.nn as nn
4
5 class MultiHeadSelfAttention(nn.Module):
6     def __init__(self, linear_dim: int, heads: int):
7
8         """
9         Clase que implementa el "Multi-Head Self Attention" de la arquitectura transformers.
10        linear_dim: dimensión lineal de los datos de entrada (int)
11        heads: número de "heads" por los que hacer "self-attention" en paralelo (int)
12        """
13
14        super(MultiHeadSelfAttention, self).__init__()
15
16        self.linear_dim = linear_dim
17        self.heads = heads
18
19        self.W_q = nn.ModuleList([nn.Linear(linear_dim, linear_dim) for _ in range(heads)]) # matriz de pesos Q
20        self.W_k = nn.ModuleList([nn.Linear(linear_dim, linear_dim) for _ in range(heads)]) # matriz de pesos K

```

```

21 self.W_v = nn.ModuleList([nn.Linear(linear_dim, linear_dim) for _ in range(heads)]) # matriz de pesos V
22
23 self.W_o = nn.Linear(linear_dim * heads, linear_dim) # h x d_v
24
25 def forward(self, x):
26
27     """
28     Función que recibe un tensor de entrada, computa "Multi-Head Self Attention" sobre este.
29     Devuelve un tensor con las mismas dimensiones.
30     x: tensor a transformar (torch.tensor)
31     """
32
33     output = [] # lista de resultados por head
34     for head in range(self.heads):
35
36         Q = self.W_q[head](x) # Proyección de x: Q
37         K = self.W_k[head](x) # Proyección de x: K
38         V = self.W_v[head](x) # Proyección de x: V
39
40         Z = (Q @ K.transpose(1, 2)) / np.sqrt(self.linear_dim)
41         Z = torch.softmax(Z, dim = 1) # softmax
42         Z = Z @ V # seq x d_v
43         output.append(Z) # append de resultados
44
45     output = torch.cat(output, dim = 2) # concatenamos resultados
46     output = self.W_o(output) # capa lineal para volver a dimension original
47
48     return output

```

Código A.3: Multi Layer Perceptron

```

1
2 import torch.nn as nn
3 import torch.nn.functional as F
4
5 class MultiLayerPerceptron(nn.Module):
6     def __init__(self, embedding_dim: int, mlp_dim: int, dropout: float, dropout_2: float = 0.0):
7
8         """
9         Clase que implementa la capa de MLP referente al bloque de Encoder de ViT transformer.
10        embedding_dim: Dimensión del embedding de los datos de entrada (int)
11        mlp_dim: Dimensión de la capa lineal MLP para transformar los datos (int)
12        dropout: Dropout a aplicar entre ambas capas lineales (float)
13        dropout_2: Dropout a aplicar después de la última capa lineal (float)
14        """
15
16        super(MultiLayerPerceptron, self).__init__()
17
18        self.fc1 = nn.Linear(embedding_dim, mlp_dim)
19        self.fc2 = nn.Linear(mlp_dim, embedding_dim)
20        self.dropout = nn.Dropout(dropout)
21        self.dropout_2 = nn.Dropout(dropout_2)
22
23    def forward(self, x):
24
25        """
26        Función que recibe un tensor de entrada, la transforma a través de capas lineales "Feed-forward".
27        x: tensor a transformar (torch.tensor)
28        """
29

```



```

30 x = self.fc1(x) # Capa 1
31 x = F.gelu(x) # GELU
32 x = self.dropout(x) # Dropout
33 x = self.fc2(x) # Capa 2
34 x = self.dropout_2(x) # Dropout
35
36 return x

```

Código A.4: Bloque de Encoder

```

1
2 from model.MultiHeadSelfAttention import MultiHeadSelfAttention
3 from model.MultiLayerPerceptron import MultiLayerPerceptron
4 import torch.nn as nn
5
6 class EncoderBlock(nn.Module):
7     def __init__(self, embedding_dim: int, mlp_dim: int, heads: int, dropout: float, dropout_2: float):
8
9         """
10         Clase que implementa un bloque o "capa" del Encoder, referente al modelo ViT transformer.
11         Se apoya de las clases MultiHeadSelfAttention y MultiLayerPerceptron.
12         embedding_dim: dimensión del embedding de los datos de entrada (int)
13         mlp_dim: Dimensión de la capa lineal MLP para transformar los datos (int)
14         heads: número de "heads" por los que hacer "self-attention" en paralelo (int)
15         dropout: Dropout a aplicar entre ambas capas lineales (float)
16         dropout_2: Dropout a aplicar después de la última capa lineal (float)
17         """
18
19         super(EncoderBlock, self).__init__()
20
21         self.attention = MultiHeadSelfAttention(linear_dim = embedding_dim, heads = heads)
22         self.mlp = MultiLayerPerceptron(embedding_dim = embedding_dim, mlp_dim = mlp_dim, dropout = dropout,
23                                         ↳ dropout_2 = dropout_2)
24         self.ln1 = nn.LayerNorm(embedding_dim)
25         self.ln2 = nn.LayerNorm(embedding_dim)
26
27     def forward(self, x):
28
29         """
30         Función que recibe un tensor de entrada, computa "Multi-Head Self Attention" y lo transforma con capas lineales.
31         x: tensor a transformar (torch.tensor)
32         """
33
34         res = x # skip connection
35         output = self.ln1(x) # LN
36         output = self.attention(output) + res # MHSA + skip connection
37
38         res = output # skip connection
39         output = self.ln2(output) # LN
40         output = self.mlp(output) + res # MLP + skip connection
41
42         return output

```

Código A.5: Modelo ViT

```

1
2 from model.EncoderBlock import EncoderBlock
3 from model.gen_patches import gen_patches
4 import torch
5 import torch.nn as nn

```

```

6
7 class ViT(nn.Module):
8     def __init__(self, img_size: int, patch_size: int, embedding_dim: int, mlp_dim: int, n_blocks: int, heads: int,
9         n_classes: int, n_channels: int = 3, dropout: float = 0.1, dropout_2: float = 0.1):
10
11         """
12         Clase que implementa el modelo ViT transformer. Se apoya de la clase EncoderBlock y la función gen_patches.
13         img_size: tamaño de las imágenes de entrada (int)
14         patch_size: tamaño de los patches (int)
15         embedding_dim: tamaño de los embeddings a ocupar (int)
16         mlp_dim: tamaño de la capa lineal del MLP a ocupar en el Encoder (int)
17         n_blocks: cantidad de bloques o "capas" en el Encoder
18         heads: número de "heads" por los que hacer "self-attention" en paralelo (int)
19         n_classes: cantidad de clases a predecir (int)
20         n_channels: número de canales de las imágenes de entrada (int)
21         dropout: Dropout a aplicar entre ambas capas lineales (float)
22         dropout_2: Dropout a aplicar después de la última capa lineal (float)
23         """
24
25         super(ViT, self).__init__()
26
27         self.img_size = img_size # dimension de H y W
28         self.patch_size = patch_size # dimension de patches
29         self.linear = nn.Linear(n_channels * self.patch_size ** 2, embedding_dim) # capa lineal de proyección de patches
30
31         self.cls = nn.Parameter(torch.rand(1, embedding_dim)) # token cls
32         self.pos = nn.Parameter(torch.rand((img_size // patch_size) ** 2 + 1, embedding_dim)) # embedding posicional
33
34         self.blocks = nn.ModuleList(
35             [EncoderBlock(embedding_dim = embedding_dim,
36                 mlp_dim = mlp_dim,
37                 heads = heads,
38                 dropout = dropout,
39                 dropout_2 = dropout_2) for _ in range(n_blocks)]) # bloques encoder
40
41         self.header = nn.Linear(embedding_dim, n_classes) # header de clasificación
42
43     def forward(self, x):
44
45         """
46         Función que recibe imágenes de entrada en formato de tensor, devuelve una predicción para cada una.
47         x: tensor a transformar (torch.tensor)
48         """
49
50         batch_size, C, H, W = x.shape # dimensiones de entrada
51
52         x = gen_patches(x, patch_size = self.patch_size) # generamos patches
53
54         x = self.linear(x) # proyección lineal de patches
55
56         token = self.cls.repeat(batch_size, 1, 1) # token cls
57         pos = self.pos.repeat(batch_size, 1, 1) # embeddings posicionales
58
59         x = pos + torch.cat((x, token), dim = 1) # sumamos embeddings posicionales + token cls
60
61         # capas de atención
62         for block in self.blocks:
63             x = block(x)
64
65         x = x[:, 0] # recuperamos cls token
66

```

```
67     x = self.header(x) # header de clasificación
68
69     return x
```