

Grado en Matemáticas

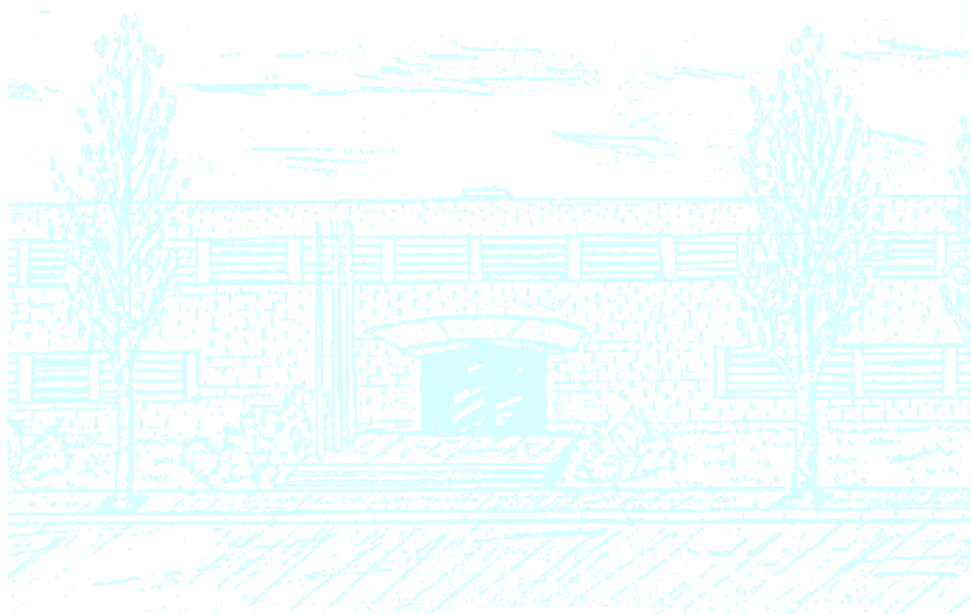
Título: Wordnet y Deep Learning: Una posible unión

Autor: Raquel Leandra Pérez Arnal

Directores: Dario Garcia Gasulla y Claudio Ulises Cortés García

Departamento: Ciencias de la Computación

Año Académico: 2017/2018



Universidad Politécnic de Catalunya
Facultad de Matemáticas i Estadística

Grado en Matemáticas
Trabajo de fin de grado

Wordnet y Deep Learning: Una posible unión

Raquel Leandra Pérez Arnal

Supervisado por Dario Garcia Gasulla y Claudio Ulises Cortés García

Enero, 2017

Gracias a Dario y Ulises por ser mis supervisores. Gracias a Ferran y Armand por ayudarme cuando tenía dudas. Gracias a mis amigos por apoyarme cuando estaba desbordada de trabajo y por leerse los distintos borradores que hice cuando estaba redactando el documento.

Abstract

Convolutional neural networks (CNN) are representation learning techniques that achieve state-of-the-art performance on almost every image-related, machine learning task. Applying the representation languages build by these models to tasks beyond the one they were originally trained for is a field of interest known as transfer learning for feature extraction. In this work we will study a Full-Network Embedding made using transfer learning from a deep convolutional neural network trained for image classification using the dataset of Imagenet and its relation with different synsets of Wordnet.

Keywords

Convolutional Neural Network, Wordnet, Imagenet, Transfer Learning

Índice

1. Introducción	3
2. Conocimientos previos	3
2.1. Redes neuronales	3
2.1.1. Una neurona	3
2.1.2. Organización por capas	5
2.1.3. Entrenar la red: Backpropagation	6
2.2. Redes convolucionales	10
2.3. Transfer Learning	12
3. Trabajo Relacionado	14
3.1. Full-Network embedding	14
3.2. Wordnet	15
3.3. Imagenet	16
4. Enfoque	17
4.1. Objetivos	18
4.2. Estadísticas	18
4.2.1. Synsets	18
4.2.2. Visión de conjunto del Embedding e Hipótesis iniciales	19
5. Análisis	22
5.1. De Wordnet al Full Network Embedding	22
5.1.1. Distribución por tipo de capa	22
5.1.2. Comportamiento respecto a la profundidad	22
5.1.3. Comportamiento de los synsets	23
5.2. Del full network embedding a wordnet	29
6. Conclusiones	35
Referencias	37
A. Gráficas de las distribuciones de las capas por <i>synset</i>	38

1. Introducción

Uno de los campos de estudio en alza del momento es el de *Machine Learning*, una rama de la inteligencia artificial que da a ciertos programas la habilidad de aprender autónomamente a partir de datos. Al no programar de forma directa los programas, sino entrenados, dentro de este campo una de las mayores incógnitas es descubrir por qué algunos algoritmos toman ciertas decisiones en vez de otras.

En este trabajo echaremos una ojeada dentro de una red convolucional profunda y estudiaremos las diferentes relaciones que se pueden encontrar entre una red convolucional entrenada con el conjunto de datos de imagenet y los *synsets* de *wordnet*. Todo esto basándonos en el trabajo presentado en [11].

Para ello empezaremos explicando a nivel general que son una red neuronal, una red convolucional y *transfer learning* en la sección 2. Seguidamente explicaremos el contenido del paper [11] y que son *Wordnet* e *Imagenet* (Sección 3). Una vez explicadas las bases en las que se fundamenta el trabajo, haremos un análisis de un *Full Network Embedding* generado a partir de una red convolucional entrenada con un conjunto de datos de *Imagenet* y estudiaremos las relaciones encontradas entre el *embedding* y unos *synsets* seleccionados (Secciones 4 y 5). Finalmente explicaremos las conclusiones obtenidas a partir de este estudio.

2. Conocimientos previos

2.1 Redes neuronales

Una **red neuronal** es un modelo computacional de aprendizaje inspirado en la forma de procesar información de las neuronas del cerebro. Las redes neuronales han generado una significativa cantidad de investigación e industria gracias a sus resultados en reconocimiento del habla, visión por computador y procesamiento de texto. Para dejar claro su funcionamiento explicaremos primero el funcionamiento de una sola neurona, para después unirlo con la organización por capas de varias neuronas y finalmente explicar el algoritmo de *backpropagation*, mediante el cual se entrenan estos modelos de aprendizaje automático.

2.1.1 Una neurona

Definición 2.1. La unidad de computación básica de una red neuronal es una **neurona**, también llamada nodo o unidad. Ésta, recibe un conjunto de entradas sobre la que se aplica una función para generar una salida.

Cada entrada (\mathbf{x}) tiene asociado un peso (\mathbf{w}), que es aprendido en base a su importancia relativa a otras entradas, a la salida y al objetivo de la red. La neurona aplica una función (f) a la suma ponderada de las entradas.

Además de los elementos ya comentados tendremos un término de sesgo (w_0), cuyo objetivo es proveer cada neurona con un valor entrenable que no dependa de la entrada y facilitar el ajuste del modelo a los datos.

La salida de la neurona (\mathbf{y}) se calculará de la siguiente manera:

$$\mathbf{y}(\mathbf{x}, \mathbf{w}) = f \left(\sum_{i=1}^N \omega_i x_i + w_0 \right)$$

Llamaremos a f la **función de activación**, es una función (no lineal), generalmente diferenciable, cuyo propósito es introducir no linealidad a la salida de la neurona. Ésto permite adaptar el modelo a problemas reales, puesto a que estos generalmente requieren relaciones no lineales.

En la figura 1 podemos ver un ejemplo visual del comportamiento de una sola neurona.

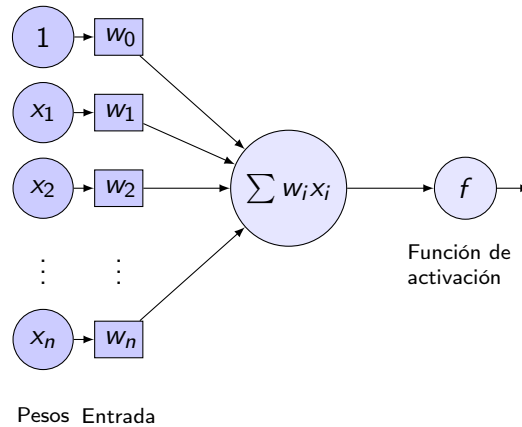
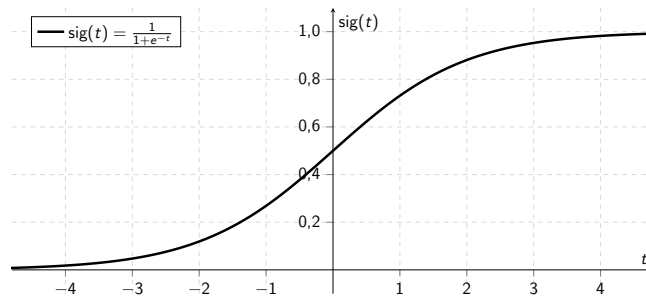


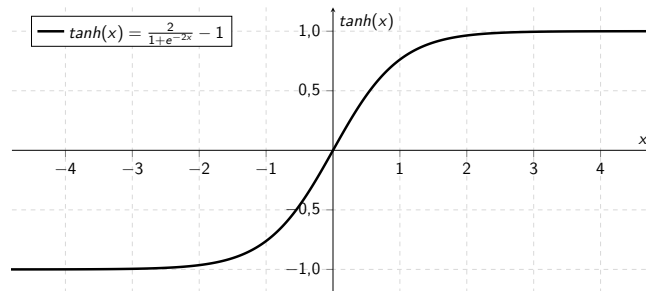
Figura 1: Ejemplo de una neurona

Las funciones de activación más utilizadas son:

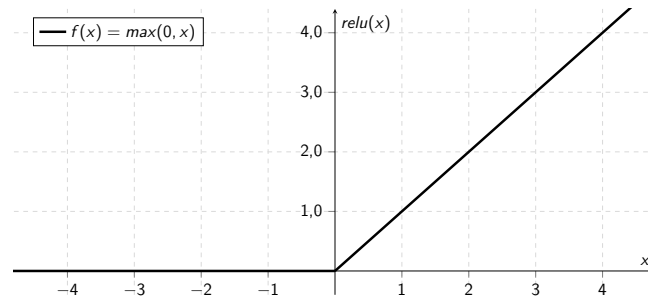
- Sigmoide: $f(x) = \frac{1}{1+e^{-x}}$



- Tanh: $f(x) = \frac{2}{1+e^{-2x}} - 1$



- ReLU: $f(x) = \max(0, x)$. Puesto que no es diferenciable en el cero, para poder utilizar *back-propagation* se suele tomar $f'(0) = x_0$ siendo x_0 un valor arbitrario.



Dotando una neurona de una función de pérdida en la salida podríamos convertirla en un clasificador lineal, pero la verdadera potencia del método viene dada al unir diversas neuronas en lo que llamaremos capas.

2.1.2 Organización por capas

Las redes neuronales se modelizan como una colección de neuronas conectadas en un grafo acíclico, comúnmente organizado por capas. El tipo más común de capa es la capa completa (*fully-connected layer*), en el que todas las neuronas de una capa están conectadas con todas las neuronas de la capa posterior a ésta, mientras que no comparten ninguna conexión con las de su propia capa, como muestra la figura 2.

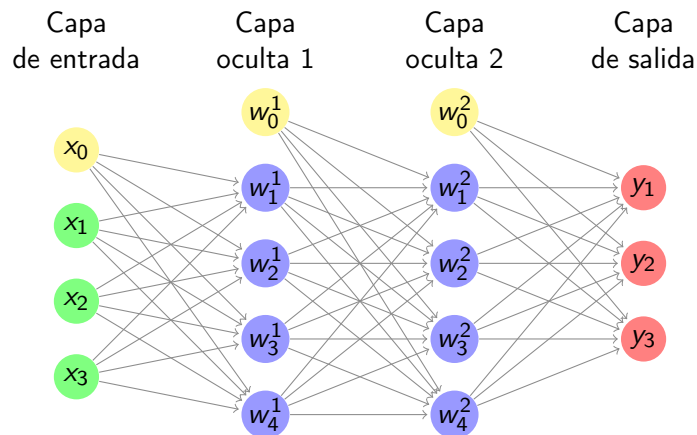


Figura 2: Ejemplo de red neuronal compuesta por capas completas

La red neuronal más básica la podemos dividir en:

1. Neuronas de entrada: Proveen la red de información del exterior, consisten en los datos de entrada, todo su conjunto se conoce como la capa de entrada.
2. Neuronas ocultas: Las neuronas ocultas no tienen conexión directa con el exterior. Solo calculan y transfieren información de la entrada a la salida. La colección de las neuronas ocultas se denomina las capas ocultas.
3. Neuronas de salida: Colectivamente denominados capa de salida y son responsables de generar la salida de la red neuronal calculada a partir de la entrada y procesada por varias capas.

Matemáticamente la notación (tomada de [1]) varía ligeramente al añadir capas. Si tenemos una muestra de tamaño N y una red con M capas, definimos las activaciones ($a_j^{(k)}$) para cada neurona j de la capa k como:

$$a_j^{(k)}(\mathbf{x}, \mathbf{w}) = \sum_{i=1}^N w_{ji}^{(k)} x_i + w_{j0}^{(k)}$$

Ecuación 1

Donde w_{ji} es el peso que va de la neurona i a la neurona j y w_{j0} es el sesgo.

Si tomamos $z_j^{(k)} = f(a_j^{(k)})$ podemos expresar las activaciones de la siguiente capa como:

$$a_j^{(k+1)}(\mathbf{x}, \mathbf{w}) = \sum_{i=1}^N w_{ji}^{(k+1)} z_i^{(k)} + w_{j0}^{(k+1)}$$

Ecuación 2

Después de procesar todas las capas la salida de la red sería:

$$\mathbf{y}(\mathbf{x}, \mathbf{w}) = \mathbf{a}^{(M)} = \sum_{i=1}^N w_{ji}^{(M-1)} z_j^{(M-1)} + w_{j0}^{(M-1)} = \sum_{i=1}^N w_{ji}^{(M-1)} \left(\sum_{i=1}^N w_{ji}^{(M-2)} z_j^{(M-2)} + w_{j0}^{(M-2)} \right) + w_{j0}^{(M-1)}$$

Ecuación 3

Donde, para calcular la salida, habría que calcular todos los sumatorios encadenados hasta llegar a la entrada.

De esta forma nos queda una estructura formada por pesos ordenados en distintos tipos de capas. Cada peso de la red es un parámetro que necesitamos que la red aprenda, y de esta forma corresponda con el resultado esperado (con mayor o menor error). Para ello se utiliza el algoritmo de *backpropagation*, que explicaremos en el siguiente apartado.

2.1.3 Entrenar la red: Backpropagation

Finalmente, queda explicar como entrenar los distintos parámetros (i.e. pesos) para que se adapten a los datos, para ello se utiliza **backpropagation**.

La idea general del algoritmo consiste en obtener las derivadas parciales de la función de error respecto a los pesos. Como el error sólo se puede computar a la salida, para poder obtener las de los parámetros intermedios utiliza la regla de la cadena. El objetivo del algoritmo es, una vez tengamos las derivadas, utilizar un algoritmo de minimización sobre la función de error para obtener los valores entrenados de los pesos.

Sea $E(\mathbf{w})$ nuestra función de error, que depende de los pesos \mathbf{w} , y \mathbf{t} el vector objetivo, es decir, los valores que debería dar la red neuronal dada la entrada $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$.

Si tomamos el error de mínimos cuadrados tenemos:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \|\mathbf{y}(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}_n\|^2$$

Ecuación 4

Que, como vimos en la ecuación 3, depende de los pesos de todas las capas. Esta estructura nos dificulta minimizar la función de una forma eficiente.

Para facilitar la notación, definimos $z_0 = 1$, $y_{nk} = y_k(\mathbf{x}_n, \mathbf{w})$ y $E_n(\mathbf{w}) = \frac{1}{2} \|\mathbf{y}(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}_n\|^2$, de esta forma tenemos:

$$a_j^{(k)} = \sum_{i=0}^N w_{ji} z_i^{(k)}$$

Ecuación 5. Activación de la neurona j de la capa k

$$E_n = \frac{1}{2} \sum_h (y_{nh} - t_{nh})^2$$

Ecuación 6. Error en función de la muestra

$$E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w})$$

Ecuación 7. Error total

Donde $a_j^{(k)}$ es la activación de la neurona j de la capa k y el índice h corresponde a cada una de las coordenadas de los vectores \mathbf{y} y \mathbf{t} .

Para calcular las derivadas de $E(\mathbf{w})$, calcularemos cada una de las de E_n que definimos en la ecuación 6.

Como vemos en la siguiente ecuación:

$$\frac{\partial E_n}{\partial w_{ji}^{(k)}} = (y_{nj} - t_{nj}) \frac{\partial y_{nj}}{\partial w_{ji}^{(k)}}$$

Ecuación 8

Donde n es el índice de la muestra, j es el índice de la neurona y k es el índice de la capa.

Utilizando la ecuación 3, podemos calcular el valor que corresponde a la derivada del error respecto a los pesos de la última capa (que depende de la penúltima), sin embargo el cálculo de la derivada respecto a los pesos de capas anteriores se dificulta.

$$\frac{\partial E_n}{\partial w_{ji}^{(M)}} = (y_{nj} - t_{nj})z_j^{(M-1)}$$

Ecuación 9

Para continuar desarrollando introducimos la siguiente notación:

$$\delta_j^{(k)} = \frac{\partial E_n}{\partial a_j^{(k)}}$$

Ecuación 10

Donde las δ se suelen llamar errores. Además si calculamos la derivada de las activaciones respecto a los pesos utilizando la simplificación de la fórmula 5 tenemos:

$$\frac{\partial a_j^{(k)}}{\partial w_{ji}^{(k)}} = z_i^{(k)}$$

Ecuación 11

Utilizando la regla de la cadena sobre $\frac{\partial E_n}{\partial w_{ji}^{(k)}}$ y utilizando 10 y 11 tenemos:

$$\frac{\partial E_n}{\partial w_{ji}^{(k)}} = \frac{\partial E_n}{\partial a_j^{(k)}} \frac{\partial a_j^{(k)}}{\partial w_{ji}^{(k)}} = \delta_j^{(k)} z_i^{(k)}$$

Ecuación 12

Utilizando 10 con los nodos de salida y 4, tenemos:

$$\delta_h^{(M)} = y_h - t_h$$

Ecuación 13

El resto de δ s se pueden calcular utilizando la regla de la cadena sobre 10 y el valor de delta de la capa final 13. De esta forma tenemos:

$$\delta_j = \frac{\partial E_n}{\partial a_j} = \sum_l \frac{\partial E_n}{\partial a_l} \frac{\partial a_l}{\partial a_j}$$

Ecuación 14

Donde la l incluye todas las neuronas a las que está conectada j . Finalmente, uniendo todo obtenemos la **fórmula de backpropagation**:

$$\delta_j = f'(a_j) \sum_l w_{lj} \delta_l$$

Ecuación 15

En resumen, dada una inicialización aleatoria o con una distribución dada de la red, el algoritmo de *backpropagation* empieza con un *forward pass* (una evaluación en la que se procesa la entrada desde el principio hasta generar la salida de acuerdo con el estado actual de la red) por toda la red, con el que se calcula la y_n y todas las activaciones $a_{ij}^{(k)}$, la predicción dada por la red se compara con la salida esperada (t_n) y se calcula un error (δ) para los pesos de la salida, luego se propaga este error utilizando 15. Finalmente utilizando 12 obtenemos los valores de las derivadas para todos los pesos.

Una vez calculadas las derivadas, se utiliza el *Stochastic Gradient Descent*, o otros métodos de minimización, para obtener los valores de los pesos que minimizan el error. El *Stochastic Gradient Descent* [1] es un método iterativo que, dada una función a minimizar $E(w)$ obtiene un mínimo local.

Cada nuevo valor para $w^{(i+1)}$ (nótese que en este caso utilizamos el super-índice para indicar la iteración a la que corresponde w , no como índice de la capa) se calcula de la siguiente forma (partiendo de una inicialización aleatoria):

$$w^{(i+1)} := w^{(i)} - \lambda \nabla E(w^{(i)})$$

Donde λ es un hiper-parámetro a añadir a nuestro modelo y $\nabla E(w^{(i)})$ es el gradiente de E , cuyos términos hemos calculado utilizando *backpropagation*. Nótese que al ser una optimización en un espacio multi-dimensional (por cada capa un peso por neurona, siendo cada peso un parámetro a optimizar), se requiere muchas muestras y ver varias veces cada muestra para minimizar el error.

2.2 Redes convolucionales

Definición 2.2. Una **red convolucional** es una secuencia de capas de diferentes tipos, donde al menos una de las capas es una capa convolucional.

Las redes convolucionales son muy similares a las redes neuronales explicadas en la sección anterior, las mayores diferencias son que contienen un tipo especial de neuronas (las neuronas convolucionales) y que asume explícitamente que la entrada son datos con una estructura concreta (como por ejemplo imágenes), lo que permite añadir ciertas propiedades a la arquitectura. Esto permite hacer un *forward pass* más eficiente de implementar y reducir significativamente la cantidad de parámetros.

Una neurona convolucional es aquella que opera con una convolución[2].

Definición 2.3. Una **convolución** es una aplicación matemática cuya entrada son dos funciones de valores reales (x y w) y cuya salida es otra función (s).

Usualmente se denota como $s(t) = (x * w)(t)$.

En el contexto de redes convolucionales llamaremos a x *entrada*, a w *kernel* y a la salida s *feature map*.

Nótese que la definición no depende de la dimensión del espacio, es decir, que podemos definir una convolución para tantas dimensiones como necesitemos. En el caso concreto de una imagen I (de dos dimensiones) y un kernel K , la convolución quedaría como:

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(i + m, j + n) K(m, n)$$

Ecuación 16. Convolución

Definición 2.4. Una **función de kernel** es aquella que cumple:

$$K(x, x') = \Phi(x)^T \Phi(x')$$

Donde $x, x' \in X$, siendo X un conjunto, y Φ es una función de base tal que $\Phi(x)^T \Phi(x') \in \mathbb{R}$. En el caso de las redes convolucionales se suele utilizar un *kernel* lineal.

Típicamente las neuronas convolucionales tienen una entrada limitada, es decir están conectadas solo a un subconjunto de las neuronas de la capa anterior, de esta forma cada neurona se centra solo en parte de la entrada, lo que significa el aprendizaje de las neuronas individuales cuando hay entradas de muchas dimensiones. En contraste, las neuronas de una capa *fully-connected* están conectadas a todas las neuronas de la capa anterior. De esta forma las redes convolucionales reducen la cantidad de parámetros de la red.

Por su conectividad limitada, las capas convolucionales pueden centrarse en una parcela particular de la entrada. El conjunto de pesos (que sería equivalente al kernel comentado en la definición 2.3) aprendido por esta parcela puede ser relevante para las otras parcelas de la entrada, por tanto, podemos definir neuronas similares que utilicen el mismo kernel, pero que se enfoquen en otra parte de la entrada. Esta idea se conoce como "weight sharing", por que varias neuronas de la misma capa están definidas por un conjunto común de pesos, esto permite tener una considerable cantidad de neuronas utilizando los mismos parámetros. Esto puede ayudar a detectar los mismos patrones en diferentes partes de la entrada. Por ejemplo si tenemos una neurona que detecta ojos, al pasarla por una foto con una cara se activará dos veces.

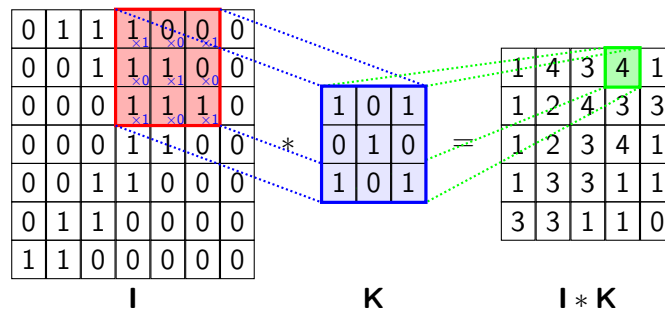


Figura 3: Ejemplo de la forma de operar una convolución sobre una imagen.

En la figura 3 podemos ver un ejemplo del funcionamiento de una convolución con un kernel de tamaño 9 sobre una imagen, podemos observar como se comparten los pesos a la hora de calcular el resultado. Dada una sección de la red, el *kernel* multiplica cada uno de sus valores por el valor correspondiente de ésta, para finalmente, sumar el resultado de todos los productos.

Finalmente queda comentar los tipos de capa más comunes en una red convolucional:

- Capa convolucional, capas formadas por neuronas convolucionales, la mayoría de capas de una red convolucional serán de este tipo.
- Capa *fully-connected*, es una capa donde todas las neuronas están conectadas entre ellas. Se suele utilizar como capa de clasificación.
- *Pooling Layer*, consiste en una capa que aplica una reducción matemática a su input (como una media o un max). El objetivo es aportar un cierto grado de invariancia espacial (imágenes equivalentes en diferentes posiciones), además reducen el tamaño de la salida de la capa, lo que provoca una reducción de la complejidad de la red. Un ejemplo de *pooling layer* sería la figura 4.

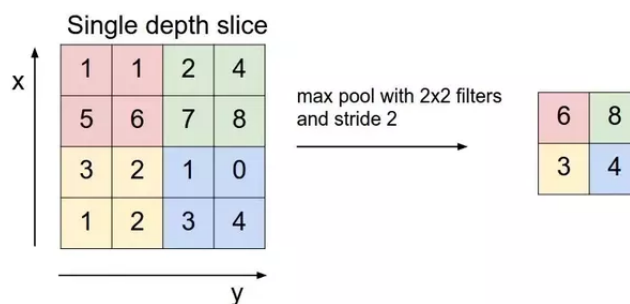


Figura 4

Otro concepto que utilizaremos a lo largo del documento es el de *feature*.

Definición 2.5. Llamaremos **feature** a los distintos subconjuntos de neuronas que ayudan a una red a cumplir con la tarea para la que fue diseñada, se generan al entrenar los pesos y pueden tener distintos grados de abstracción.

En redes convolucionales, la disposición por capas permite representar la información más compleja a partir de otra más simples. Por tanto puede acabar extrayendo *features* abstractas. En la figura 5 podemos

ver un ejemplo de *features* de distinto nivel de abstracción. Por ejemplo, si queremos clasificar personas, la presencia de *dos ojos* sería una *feature*.

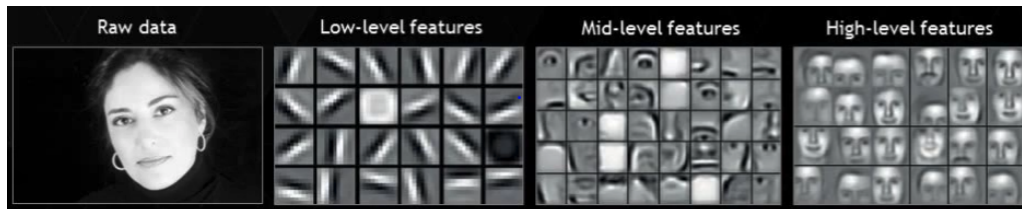


Figura 5: Ejemplo de features

2.3 Transfer Learning

En la práctica pocas personas entrenan una red profunda desde cero (con inicialización aleatoria), puesto a que éstas tienen unos requisitos bastante exigentes a la hora de ser entrenadas.

Los mayores problemas que puedes encontrar son:

1. A causa de la gran cantidad de parámetros a entrenar que tiene una red profunda necesitas un **conjunto de datos de gran tamaño**.
2. El **coste computacional** de entrenar la red, determinado por el número de parámetros, que en el caso de redes convolucionales profundas pueden ser del orden de cientos de miles.
3. La búsqueda de **hiper-parámetros** (valores de los que depende el modelo) óptimos.

La forma de obtener los hiper-parámetros es, entrenar la red para cada uno de los hiper-parámetros que hayamos considerado y, utilizando una partición de validación (o el método de *cross-validation*) comparar los distintos resultados.

Esto incrementa los problemas 1, al necesitar una partición de validación necesitaremos más datos, y 2, el coste computacional se multiplica por cada vez que haya que entrenar la red, que puede ser varias veces por cada hiper-parámetro que tengamos.

Por tanto, es común pre-entrenar una red convolucional en un conjunto de datos significativamente grande y usar la red convolucional como inicialización o como extractor fijo de características de la tarea de interés (como vemos en [3]).

Definición 2.6. Llamaremos **transfer learning** al campo de estudio que reutiliza el lenguaje de representación de un problema (que llamaremos problema origen o *Source*) para resolver otro (que llamaremos objetivo o *Target*).

A la hora de utilizar **transfer learning** tenemos dos componentes base: Un dominio \mathcal{D} definido por un conjunto de instancias de datos con una distribución de probabilidades y una tarea \mathcal{T} , definida a partir de un conjunto de clases y una función objetivo.

Para un problema de *transfer learning* usaremos la notación (introducida en [4]):

- Problema origen: $(\mathcal{T}_S, \mathcal{D}_S)$
- Problema objetivo: $(\mathcal{T}_T, \mathcal{D}_T)$

Los dos casos más comunes de transfer learning son:

- **Fine-tuning**([5]): La primera estrategia consiste en inicializar los datos desde un estado no aleatorio (tomando los pesos ya entrenados), y a continuación, entrenar la red sobre estos datos. De esta manera puedes reducir significativamente el conjunto de datos necesario para entrenarla, sin embargo, sigue necesitando tiempo para optimizar los múltiples hiper-parámetros involucrados en el proceso y una cantidad significativa de recursos computacionales.
- **Feature Extraction**([4]) Consiste en procesar un conjunto de datos a través de una red neuronal ya entrenada y extraer valores de activación para que puedan ser utilizados por otro mecanismo de aprendizaje. Este método es aplicable a conjuntos de datos de cualquier tamaño([6], [7]), puesto que cada dato es procesado independientemente. Además tiene un menor coste computacional, ya que no tiene que entrenar la red y no requiere la optimización de hiper-parámetros. Por estos motivos las aplicaciones de *transfer learning for feature extraction* están limitadas solo a las capacidades de los métodos que utilices encima de la representación profunda obtenida.

En nuestro caso nos centraremos en un tipo concreto de *transfer learning for features extraction* [8], que explicaremos más adelante en la sección 3.1.

3. Trabajo Relacionado

3.1 Full-Network embedding

En esta sección explicaré el *full-network embedding* presentado en [11].

En general en *transfer learning for feature extraction* es común tomar los valores de activación de una sola capa cercana a la salida, como podemos ver en la figura 6, donde el proceso de *feature extraction* solo se hace en las últimas capas, *fully-connected*, mientras que las capas convolucionales (todas las anteriores) permanecen intactas. El resto de capas se descartan por "ser poco probable que contengan una representación mejor" [9], sin embargo es conocido que todas las capas de una red profunda pueden contribuir a caracterizar los datos de diferentes maneras. Esto implica que la representación más versátil y rica que puede ser generada por un proceso de *features extraction* debe incluir todas las capas de la red, es decir, debe definir un *full-network embedding*.

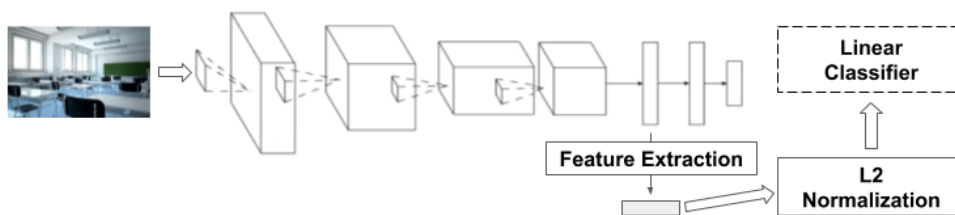


Figura 6: Estructura básica que se suele utilizar en *feature extraction*

Dado un conjunto de datos $t1$, queremos representarlo en el lenguaje aprendido para una tarea $t0$. Para ello el *full-network embedding* se divide en 4 pasos:

- El primer paso es hacer un **forward pass** de cada instancia de datos de $t1$ a través del modelo entrenado en $t0$, guardando todos los valores de activación de cada capa de la red, (tanto las convolucionales como las *fully-connected*).
- El segundo paso consiste en un **spatial average pooling** en los filtros convolucionales. Como se ha explicado en el capítulo sobre redes convolucionales, un filtro convolucional genera varias activaciones para una entrada, para darle la capacidad de obtener información espacial. El objetivo de este paso es obtener un solo valor de activación para cada filtro, evitando así un aumento de la dimensionalidad. Los valores resultantes se concatenan a los de las capas *fully-connected* en un solo vector, para generar el *embedding* completo.
- El tercer paso es una **estandarización de características**. Puesto a que los valores del *embedding* provienen de distintos tipos de neuronas en distintos puntos de la red, la distribución de las activaciones pueden variar mucho y si se concatenasen todas algunas dominarían sobre las otras.

El valor estandarizado de cada categoría se obtiene calculando la media y la desviación típica del conjunto de datos de entrenamiento, es decir, que se estandariza según el contexto de la feature en todo el conjunto de datos. De esta forma se pueden integrar los distintos valores en el *embedding*.

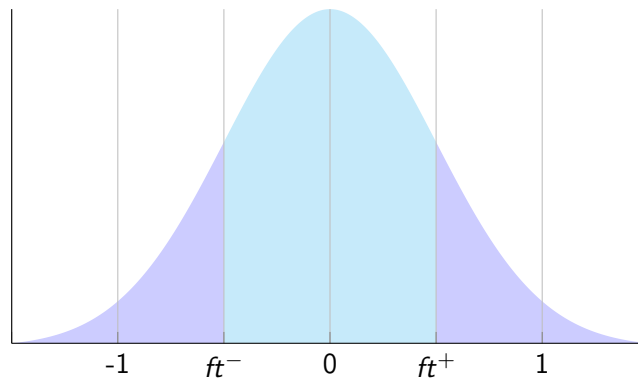


Figura 7

- Finalmente se aplica una **discretización de características** para evitar problemas al intentar explorar un espacio de dimensión tan grande (una sola capa convolucional puede llegar a tener del orden de 100,000,000 parámetros, por ejemplo [10]). Ésta discretización consiste en tomar unos límites que dependen de los datos, ft^- y ft^+ y reemplazar los valores que tienen un valor atípicamente bajo (todos los valores $x < ft^-$) por un -1, los que tienen un valor típico (todos los valores $ft^- < x < ft^+$) por un 0 y los que tienen un valor atípicamente alto (todos los valores $ft^+ < x$) por un 1, como podemos ver en la representación gráfica 7.

Para encontrar los límites, los autores se basan en [8], donde dan una visión estadística sobre como evaluar la importancia de las diferentes características de una red convolucional comparando las activaciones para una clase concreta respecto al resto de clases del conjunto de datos. De esta forma los autores de [8] separan las características en tres conjuntos: *característico por presencia*, *no característico* y *característico por ausencia*.

Un ejemplo de la separación comentada en la *discretización de características* sería: si tenemos como datos objetivo un conjunto de coches tendremos que las features que se activen con las ruedas no serán representativas, puesto a que todos los coches tienen, las categorías que se activen con el techo de un coche serán representativas por ausencia para clasificar descapotables mientras que las que se activen para la capota serán representativas por presencia.

En la figura 8 podemos ver un diagrama de como se desarrollarían los diferentes pasos. Nótese la diferencia con la figura 6. Mientras que en la figura 6 solo se modificaban las capas finales, en esta se modifican también las convolucionales, dando lugar al *full-network embedding* explicado.

3.2 Wordnet

Wordnet es una base de datos que contiene nombres, verbos, adjetivos y adverbios en conjuntos de sinónimos (que llamaremos *synsets*). Los *synsets* están conectados entre ellos por medio de relaciones conceptuales, semánticas y léxicas. Utilizando los *synsets* y sus relaciones, se puede generar un grafo que puede ser utilizado para distintos objetivos, como lingüística computacional y procesamiento del lenguaje natural.

En concreto utilizaremos las relaciones de:

- Sinonimia** : dos palabras son sinónimos si tienen el mismo significado (ej, pelo y cabello, figura 9a).

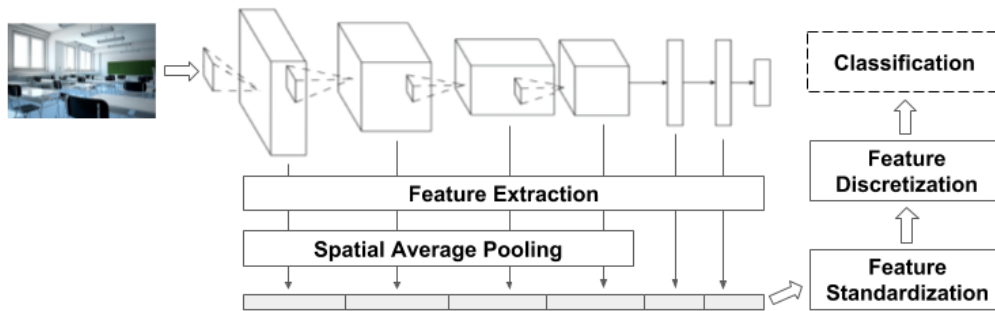


Figura 8: Estructura del *full-network embedding*

- **Hiponimia:** una palabra es hipónimo de otra si su significado es más específico que el de ésta (ej, mamífero es hipónimo de ser vivo, figura 9b).
- **Hipernimia:** una palabra es hiperónimo de otra si su significado es menos específico que el de ésta (ej, perro es hiperónimo de dalmata, figura 9c).

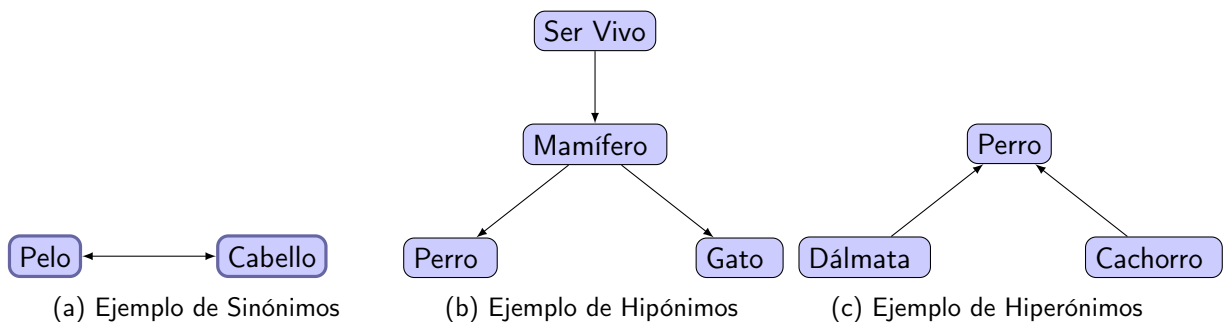


Figura 9

3.3 Imagenet

Imagenet es una base de datos de imágenes organizada utilizando la jerarquía de *wordnet*. Su principal objetivo es dotar a los investigadores en campos relacionados con visión artificial de una base de datos a gran escala con la que poder trabajar. Actualmente consta de 14,197,122 imágenes y 21,841 *synsets* indexados.

Para el *full-network embedding* utilizaron el subconjunto correspondiente al reto de *imagenet* de 2012 de reconocimiento de imágenes.

Éste consta de:

- Un conjunto de datos de entrenamiento de 1.2 millones de muestras y 1,000 categorías diferentes.
- Un conjunto de datos de validación de 50,000 muestras y 1,000 categorías.

Entre las posibles categorías podemos encontrar *synsets* de diferentes niveles de especificación, por ejemplo tenemos las categorías: perro, dalmata, pastor alemán...

4. Enfoque

En nuestro caso estudiaremos un *full-network embedding* obtenido a partir de la tarea origen (\mathcal{T}_S = Clasificación de imágenes, \mathcal{D}_S = 1.2 M imágenes de *Imagenet* con sus correspondientes clases) y una tarea de destino (\mathcal{T}_T = Clasificación de imágenes, \mathcal{D}_T = 50,000 imágenes de *Imagenet* con sus correspondientes clases). Es decir, partiendo de una red convolucional entrenada con todo el conjunto de datos de entrenamiento de *Imagenet*, hemos generado un lenguaje de representación para el conjunto de datos de validación de *Imagenet*. Lo que resulta en el *full-network embedding* con el que trabajaremos.

El *full-network embedding* consiste en una matriz de tamaño 50,000 muestras por 12,416 características, como observamos en la figura 10, cuyos valores están en $\{-1, 0, 1\}$. Para cada muestra también tenemos su correspondiente clasificación (un valor entre 0 y 999 que representa la clase de *Imagenet* con la que fue clasificado).

Es importante remarcar el significado de los diferentes valores que puede tomar una *feature* para una imagen.

- Si el valor para una imagen concreta es 1 significa que esta es representativa por presencia. Es decir, que dentro del contexto del conjunto de datos *Target*, es una *feature* que se activa con una potencia significativamente superior al resto.
- Si el valor para una imagen concreta es 0 significa que esa *feature* genera un valor de activación no significativo.
- Si el valor para una imagen concreta es -1 significa que se activa con una potencia significativamente inferior al resto.

Hemos de tener en cuenta que estas definiciones son con respecto a \mathcal{D}_T , es decir, respecto al nuevo espacio de representación.

Las características están ordenadas en capas de la siguiente manera:

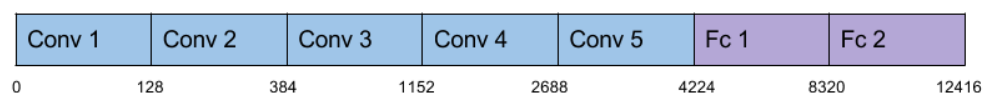


Figura 10: La disposición de las características por capas

Finalmente, de todos los *synsets* posibles tomamos dos subconjuntos diferentes:



Figura 11: Conjuntos de synsets que estudiaremos

Hemos tomado estos dos conjuntos por estar distribuidos de una forma bastante uniforme entre las clases, como veremos en el capítulo 4.2, y por ser de dos "ramas" lo más diferentes posibles, además están en sitios parecidos respecto a la jerarquía global, y por tanto representan niveles de abstracción parecidos.

4.1 Objetivos

Partiendo de los datos comentados y del trabajo publicado en [11] tomamos varios **objetivos** a partir de los que desarrollar el trabajo.

- Analizar el *embedding* dado y el comportamiento de las *features* en las distintas capas. Para ello se realizará un estudio estadístico inicial de la matriz de *embeddings*. De este estudio se analizarán los diferentes puntos de interés, y finalmente se hará un estudio más detallado de estos.
- Analizar si hay alguna relación entre los *embeddings* de las diferentes clases y *synsets* relacionados con ellas (sus hipónimos e hipérmimos). Para ello se realizarán las estadísticas del comportamiento de los *synsets* de la figura 11 con respecto a la matriz de *embeddings*, con intención de profundizar en los puntos de interés que se encuentren.

4.2 Estadísticas

La finalidad de este capítulo es hacer un estudio inicial de los datos de los que se dispone, es decir, de los correspondientes al *embedding* y a los *synsets* que hemos tomado para el estudio.

4.2.1 Synsets

Como comentamos en la sección anterior los *synsets* que hemos tomado son los especificados en la figura 11. El criterio que hemos utilizado para decidir que una clase de *Imagenet* pertenece a un *synset* es que sea hipónimo de esta, por tanto es coherente que los *synsets* más generales tengan más imágenes.

Además, ambas familias (la de seres vivos y objetos), se distribuyen de forma parecida (cada *synset* hijo tiene aproximadamente la mitad de imágenes que su *synset* padre), como podemos ver en la figura 12.

Cuadro 1: Cantidad de muestras de cada *synset*

Total	Ser Vivo	Mamífero	Perro	Perro de caza	Artefacto	Instrumento	Transporte	Vehículo
50000	20500	10900	5900	3150	26100	17900	3600	2200

Elegimos dos familias de *synsets* que cumplieran esto para que al comparar los resultados se distorsionaran lo menos posible por la cantidad de imágenes pertenecientes a ellas.

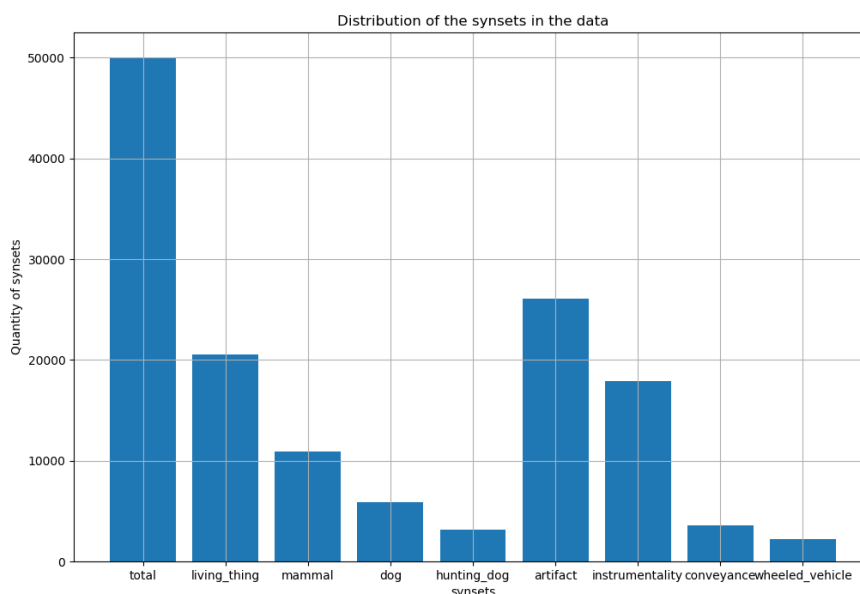


Figura 12: Cantidad de imágenes de cada *synset* respecto al embedding total

4.2.2 Visión de conjunto del Embedding e Hipótesis iniciales

Primero de todo estudiaremos las estadísticas generales del *embeddings*.

Para ello, dado nuestro *embedding* de 50,000 muestras por 12,416 *features*, calcularemos cuantos elementos del *embedding* total tienen valor 1, 0 o -1. A estos valores los llamaremos **categorías**.

En la figura 13 vemos la cantidad de *features* de cada categoría. Podemos observar que la cantidad total de -1 es significativamente mayor, esta distribución se conserva (en general) para los distintos *synsets*, independientemente de lo específicos que sean. Teniendo en cuenta que al hacer *transfer learning* se suele pasar de un conjunto de datos más general a uno más específico que haya más *features* características por ausencia es coherente, puesto a que, al ser más específicas las propiedades visuales del *Target* que del *Source*, la mayoría de las *features* se alejan más de la media de todas las imágenes. Más adelante veremos también como se distribuyen las categorías entre las diferentes capas del *embedding*.

Observamos también como se distribuyen las *features* respecto al conjunto de imágenes que tenemos.

En la figura 15 tenemos la distribución de las imágenes por categoría, es decir cuantas imágenes tienen una cierta cantidad de *features*. La idea de generar esta gráfica era ver si podíamos encontrar alguna

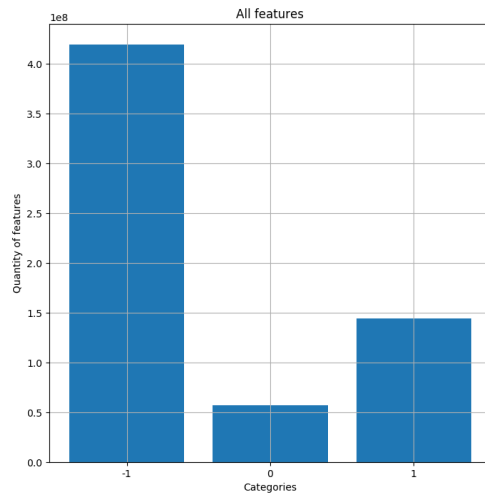
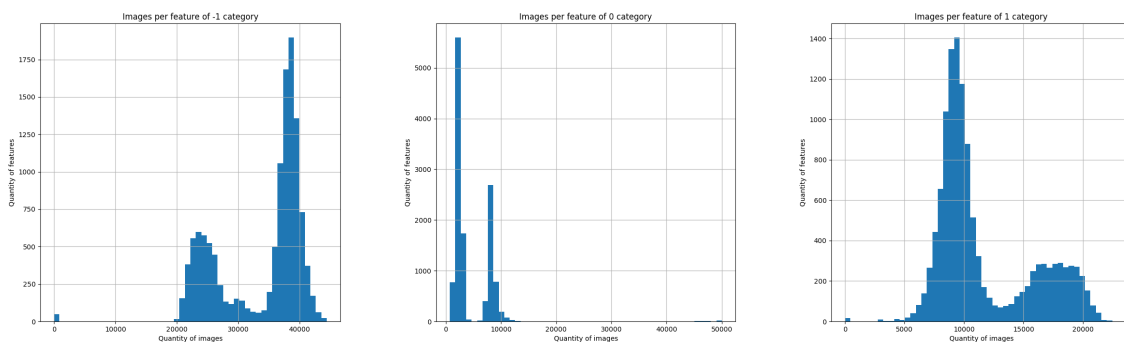


Figura 13: Cantidad de *features* de cada categoría



(a) Categoría -1

(b) Categoría 0

(c) Categoría 1

Figura 14: Distribución del número de *features* con los distintos valores categóricos, para las 50,000 imágenes

distribución concreta o patrón.

Al estudiar la matriz de *embeddings* nos planteamos las siguientes hipótesis:

1. Las características se distribuyen de diferente manera en las capas convolucionales y los completos. Esta hipótesis proviene de la gráfica 15, donde en las tres categorías se ven dos distribuciones bien diferenciadas. Esto lo analizaremos en la sección 5.1.1 del análisis.
2. Cuanto más profundo es el layer, debería haber más *features* representativas, tanto por ausencia como por presencia. Esta hipótesis viene de que en una red convolucional cuanto más profunda es una capa más concretas suelen ser las *features* con las que se activa, como vimos en la figura 5. Esta hipótesis la estudiaremos en 5.1.2
3. Cuanto más concreto es un *synset*, debería haber más *features* representativas, tanto por ausencia como por presencia, es decir, mayor proporción de -1 y 1 respecto a los 0. Debido a que al ser conjuntos más concretos de imágenes la cantidad de *features* no representativas debería ser inferior.

Esta hipótesis la estudiaremos en 5.1.3.

4. Se puede ver una relación entre los embeddings de synsets hipónimos. La idea sería que dada una imagen perteneciente a un *synset*, compartiría *features* características con sus hipónimos. Esta hipótesis la estudiaremos en 5.2.

5. Análisis

5.1 De Wordnet al Full Network Embedding

En este apartado explicaré el estudio hecho para intentar aceptar o revocar las hipótesis comentadas. Empezaré con la distribución por tipo de capa en el *embedding* general, para luego concretar estudiando el comportamiento por capa. Finalmente explicaré los descubrimientos relacionados con los *synsets*.

5.1.1 Distribución por tipo de capa

Para comprobar si la hipótesis 1 es cierta calculamos las gráficas de la figura 15 distinguiendo entre las capas convolucionales y *fully-connected*, obteniendo las siguientes gráficas:

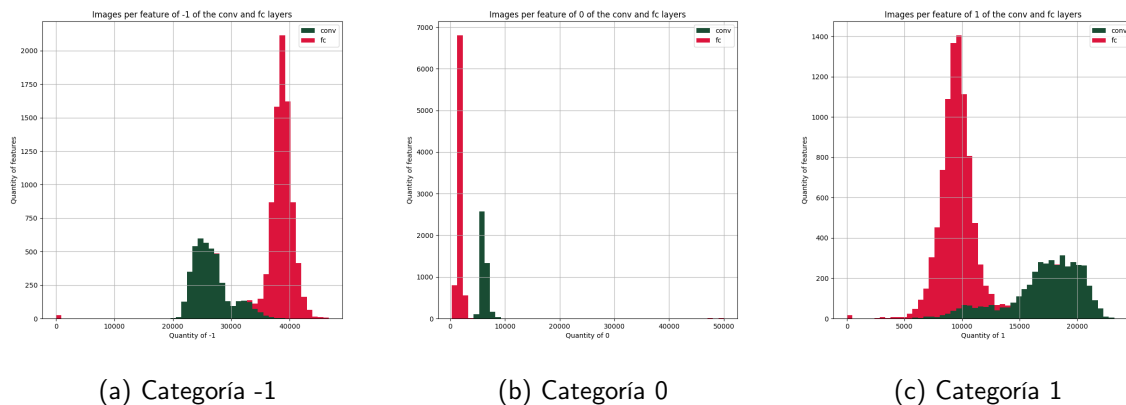


Figura 15: Distribución del número de *features* con los distintos valores categóricos distinguiendo las capas convolucionales de las *fully-connected*

En ellas podemos ver claramente la distinción que suponíamos. Sobretudo para la categoría 0, hay una clara separación entre el comportamiento de los diferentes tipos de capa. En parte esto se debe a que hay más *features* en las capas *fully-connected*, por eso sus gráficas correspondientes tienen mayor área, además en éstas podemos ver una menor variabilidad, que concuerda con el hecho de que solamente son 2 capas *fully-connected*, contra las 5 convolucionales. También vemos que el grueso de las *features* características por presencia se encuentran en las capas *fully-connected*, cosa que concuerda con el hecho de que sean estas capas las que dan la clasificación final.

En la sección (5.1.2) profundizaremos en las diferencias encontradas entre los dos tipos de capa, pero ya con lo estudiado vemos una cierta diferencia entre los dos tipos de capa, como suponíamos en nuestra hipótesis 1.

5.1.2 Comportamiento respecto a la profundidad

Una vez comprobado que los dos tipos de capa se comportan diferente, vamos a estudiar en que nivel se diferencian y a que profundidad. De esta forma veremos si podemos aceptar o desmentir la hipótesis 2.

Para ello hemos calculado la cantidad de *features* de cada tipo por capa respecto al *embedding* total.

Cuadro 2: Proporción de cada categoría por capa

	conv1	conv2	conv3	conv4	conv5	fc6	fc7
Proporción de -1	0.47	0.44	0.46	0.49	0.55	0.77	0.76
Proporción de 0	0.18	0.17	0.17	0.17	0.17	0.05	0.06
Proporción de 1	0.36	0.39	0.37	0.34	0.28	0.18	0.18

Como podemos observar en la gráfica 16. De esta forma se esperaba ver como se distribuyen las *features* por cada capa.

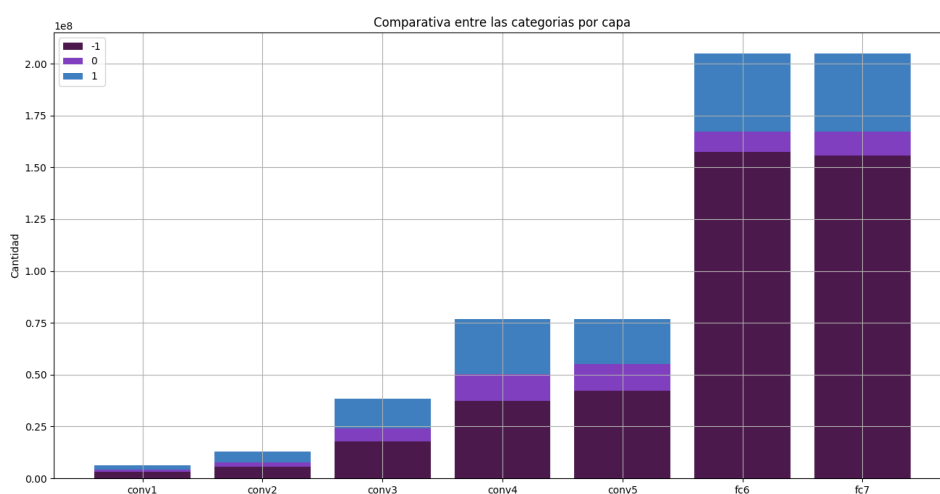


Figura 16: Cantidad de *features* de cada categoría por capa

Para complementar la información de la gráfica, podemos observar la proporción de las diferentes categorías por cada capa en la tabla 2.

En la gráfica esperábamos encontrar una mayor proporción de *features* características por presencia en las capas finales, sin embargo podemos ver que la mayor proporción se encuentra en la segunda capa convolucional, con un 39%. En esta gráfica podemos ver claramente que todas las capas, no solo las *fully-connected* contienen información importante de cara a caracterizar el espacio de representación, al contrario de lo que sostienen en [9].

Lo estudiado en esta sección nos muestra que la proporción de *features* no representativas es similar entre las capas convolucionales y entre las *fully-connected*, es decir, respecto al tipo de la capa, pero no se ve una diferencia significativa respecto a la profundidad. Por tanto revocamos la hipótesis 2. Sin embargo, vemos que la profundidad si que afecta a la proporción de *features* características por presencia en las capas convolucionales, como vimos en la tabla 2.

5.1.3 Comportamiento de los synsets

En esta sección estudiaremos las hipótesis 3.

A que a la hora de estudiar los *synsets* nos hemos topado con la dificultad de decidir el valor de una

feature para el conjunto de todas las muestras que pertenecen al *synset*, puesto a que para diferentes muestras puede tener distintos valores, como podemos ver en el ejemplo de la figura 17.

	<i>features</i>							
ser vivo	1	1	-1	1	0	0	-1	} muestras
	0	0	-1	-1	-1	0	0	
	0	0	0	1	-1	-1	0	
mamífero	1	1	-1	1	-1	0	0	
	0	0	-1	1	0	0	0	
	0	-1	1	0	0	0	1	
ser vivo	1	1	-1	-1	-1	0	0	
	-1	-1	0	0	0	0	0	

Figura 17: Muestra de una sección del embedding.

Para solucionar este problema hemos decidido tomar dos enfoques diferentes, por un lado hemos tomado la sub-matriz correspondiente a cada *synset*, por otro lado hemos tomado un vector representante para cada *synset* a estudiar.

Sub-matriz:

En este caso hemos tomado una matriz cuyas filas consisten en las muestras que pertenecen al *synset* que queremos estudiar y cuyas columnas corresponden a las *features*, como muestra el ejemplo de la figura 18.

1	-1	-1	1	0	0	-1
0	0	-1	-1	-1	0	0
0	0	0	1	-1	-1	0
0	1	-1	1	1	0	0
1	1	-1	0	-1	0	0
0	0	-1	1	0	0	0
0	-1	-1	0	0	0	1
1	1	-1	-1	-1	0	0
-1	-1	0	0	0	0	0

embedding

1	-1	-1	1	0	0	-1
1	1	-1	0	-1	0	0
1	1	1	-1	-1	0	0

sub-matriz

Figura 18: Ejemplo de una sub-matriz de un *synset*.

Este enfoque nos da información general sobre el comportamiento de los distintos *synsets* en el *embedding*.

Utilizando la sub-matriz asociada a cada uno de los *synsets* hemos generado gráficas equivalentes a la de la figura 15. En este caso, las distribuciones se mantienen bastante en los *synsets* más generales, como podemos ver en la figura 19, mientras que en los menos generales se empiezan a distorsionar. Como podemos observar en la figura 20. Las gráficas correspondientes al resto de *synsets* se encuentran en el apéndice A.

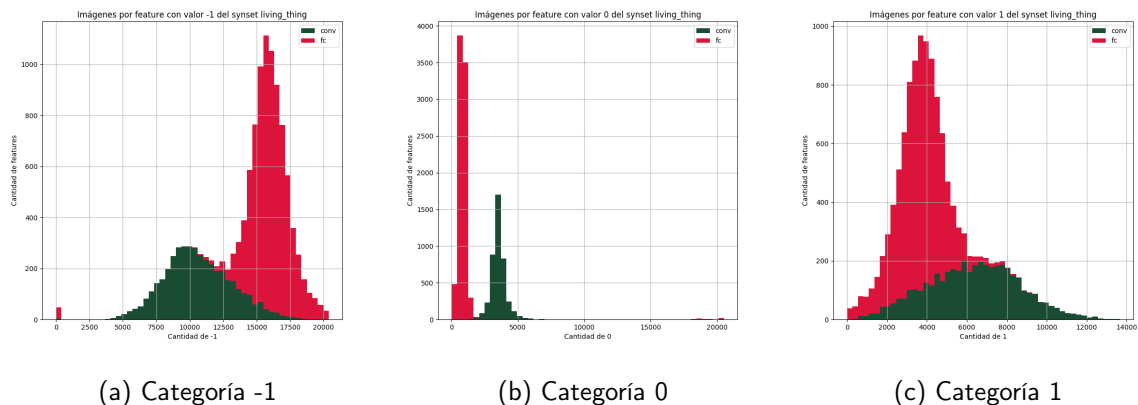


Figura 19: Distribución del número de *features* con los distintos valores categóricos distinguiendo las capas convolucionales de las *fully-connected* del *synset* seres vivos

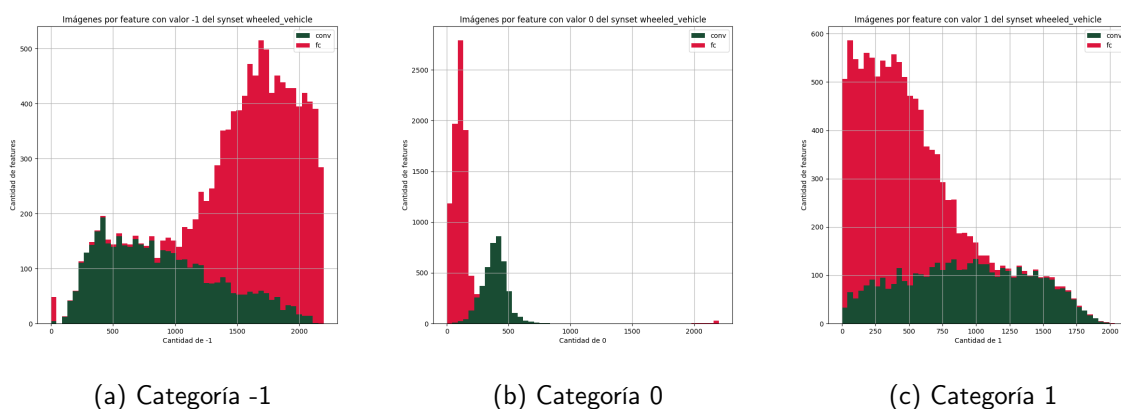


Figura 20: Distribución del número de *features* con los distintos valores categóricos distinguiendo las capas convolucionales de las *fully-connected* del *synset* vehículo

La diferencia entre los distintos grupos de gráficas (*synsets* más generales contra más específicos) se debe a que al ser una muestra más pequeña hay una mayor varianza. También se observa una disminución de la proporción de *features* no característica cuanto más concreto es el *synset*, esto nos da una idea de que la red convolucional se adapta al nivel de generalización que tiene un *synset*, por ejemplo, *perro* tendrá menos *features* no características (en proporción) que *mamífero*. Este hecho concuerda con nuestra hipótesis 3, en la que ahondaremos más adelante.

También se observan unos comportamientos simétricos entre las *features* características por presencia y por ausencia en todos los *synsets* y en la gráfica general. Este tipo de simetría es una buena señal, puesto a que la caracterización que hemos dado a las *features* representativas por ausencia y por presencia es contraria.

Otro dato que podemos obtener es la cantidad de *features* de cada una de las categorías por *synset*, como vemos en la figura 21, simplemente sumando los diferentes valores de las sub-matrices. Tanto en la figura como en la tabla 3, donde vemos la proporción de cada una de las categorías por *synset*, podemos ver que las proporciones se conservan entre *synsets*.

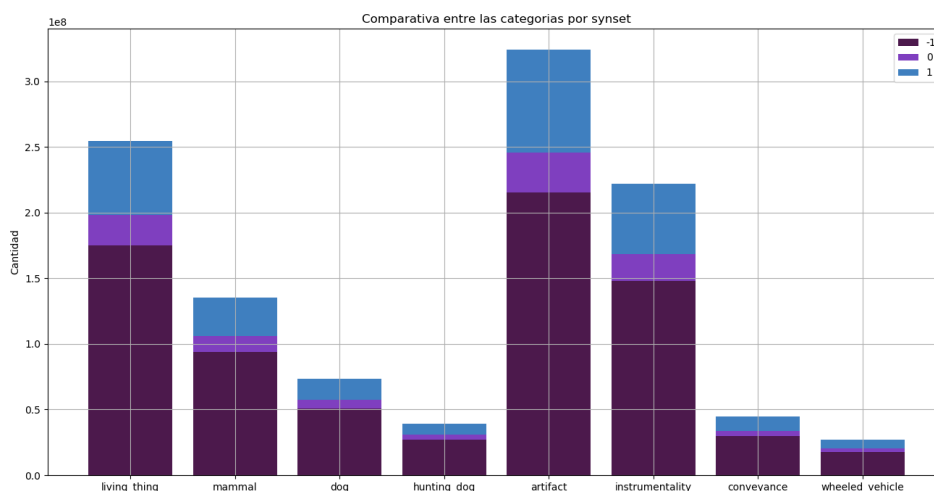


Figura 21: Cantidad de *features* de cada categoría por synset

Cuadro 3: Proporción de cada categoría por synset

	Ser Vivo	Mamífero	Perro	Perro de caza	Artefacto	Instrumento	Transporte	Vehículo
Prop. de -1	0.69	0.69	0.70	0.70	0.66	0.67	0.66	0.65
Prop. de 0	0.09	0.09	0.09	0.09	0.09	0.09	0.09	0.09
Prop. de 1	0.22	0.22	0.21	0.21	0.24	0.24	0.25	0.26

Además utilizando esta sub-matriz hemos calculado la cantidad de *features* de cada tipo por capa, de la misma forma que calculamos los valores de la gráfica 16. En general no se ve mucha diferencia entre las gráficas de cada *synset* y la general (como podemos observar en la figura 22), esto implica que la cantidad de *features* con los distintos valores por capa es parecida independientemente del *synset* que tomemos y de lo general que éste sea.

Mientras que la distribución de las imágenes por *feature* se difuminaba con *synsets* más concretos (figura 20), vemos que la proporción de *features* características se conserva, tanto por *synset* como por capa. Esto nos confirma que la difuminación que veíamos sea debida a una menor cantidad de muestras a la hora de calcular la gráfica y no a una distinción por *synset*.

Representante:

Con la sub-matriz hemos logrado una visión general del comportamiento de los *synsets* dentro del *embedding*, pero hay información que estamos perdiendo. Sería interesante saber si a nivel de *feature* diferentes *synsets* se comportan distinto, o de la misma forma, como parecían indicar los descubrimientos para sub-matrices. Por ejemplo que se concentrasen todos los 1 en un mismo conjunto de *features* para el mismo *synset*.

Para solventar este problema decidimos tomar un vector muestra para cada uno de los *synsets* de tamaño $1 \times 12,416$ donde cada componente se toma como la categoría más frecuente para esa *feature* concreta, es decir, hacemos la tomamos la moda de la matriz del *embedding* por columnas, como muestra la figura 23.

De esta obtenemos una información más concreta respecto a cada *feature*, además podemos comparar

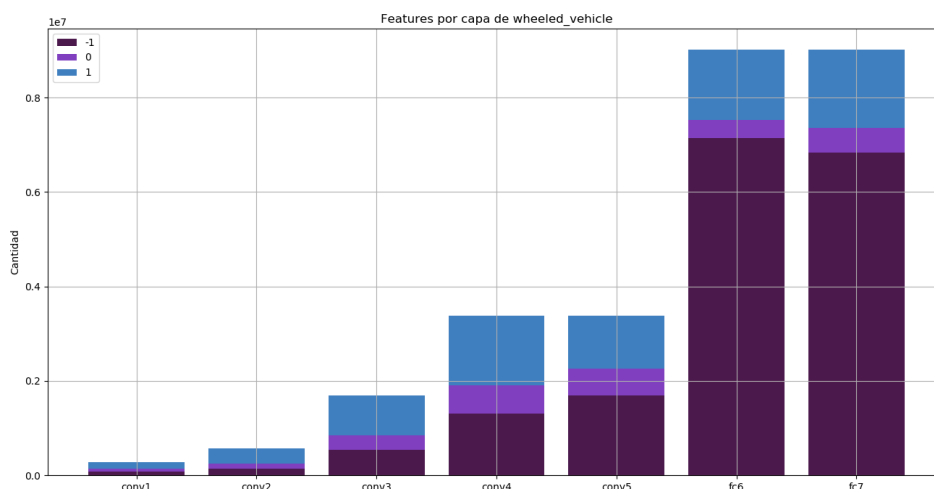


Figura 22: Cantidad de features por capa de vehículo

los distintos *synsets* de una forma más efectiva. Por ejemplo, si una *feature* suele ser característica para mamíferos, quedará marcada en el representante, mientras que si tomáramos los valores de la proporción de las características respecto a la sub-matriz correspondiente al *synset* esta información se perdería, ya que la proporción de *features* no características es muy superior a la del resto de opciones.

Al calcular las gráficas anteriores, pero del representante, esta vez si que obtenemos diferencias.

Primero de todo observamos que en la gráfica 24, cuanto más concreto es el *synset* mayor proporción de *features* características por presencia tenemos, además que las *features* no características desaparecen prácticamente por completo. Esto implica que las *features* no características no sólo están en minoría para la mayoría de imágenes y *features*, sino que además no se concentran en conjuntos concretos de *features*.

Si además comparamos la gráfica 24 con la de 21 vemos que las *features* características aparecen en mayor proporción. Si no estuvieran concentradas, las *features* tomarían todas el valor -1, puesto a que es el más común en el embedding.

Respecto a las proporciones por capa, para todos los *synsets* han salido resultados parecidos, siendo el más visible el correspondiente a Vehículo, como vemos en la 25. En este caso también vemos una desaparición de las *features* no representativas en todas las capas. Además la proporción de *features* características aumenta en las capas convolucionales, mientras que disminuye en las *fully-connected*. Lo que implica una dispersión de los valores característicos con respecto a los que no lo son.

Para comparar los *synsets* a un nivel todavía más concreto hemos generado diferentes matrices de cambio a partir de los representantes, éstas están definidas como vemos en la figura 26.

Donde $C_{(i,j)}$ es la cantidad de valores que en el primer *synset* valían i y en el segundo valen j , comparando elemento a elemento. La suma total de todas las $C_{(i,j)}$ es 12,416, es decir, la cantidad total de *features* de un representante. Además, si sumamos todas las $C_{(i,j)}$ de una columna concreta obtenemos el total de *features* con valor i del primer *synset*, y de forma simétrica, si sumamos las filas obtenemos el total de *features* con valor j del segundo *synset*.

En la diagonal contienen la cantidad de valores que se conservan, es decir, que para el primer *synset* y el

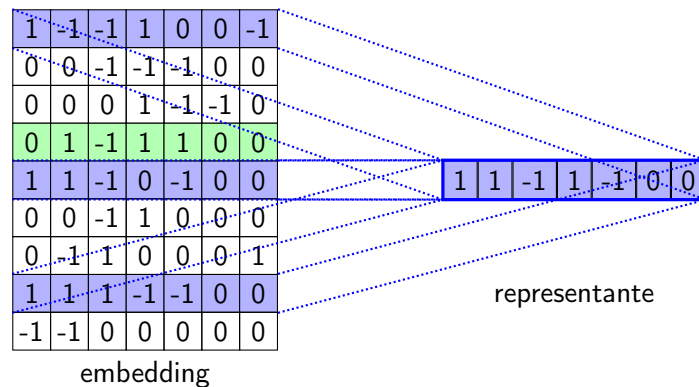


Figura 23: Ejemplo de un representante de *synset*.

segundo tienen el mismo valor. En el resto de posiciones contiene la cantidad de *features* que en el primer *synset* tenían un valor y en el segundo pasa a tener otro.

En nuestro primer ejemplo de matrices de cambio (figura 27) podemos observar que, tanto la columna como la fila correspondiente al cero (las *features* no características) tiene unos valores muy bajos con respecto al resto. Esto se repite en todas las matrices para todos los *synsets*, puesto a que el total de *features* con este valor es muy bajo. Respecto al resto de valores, los explicaré por columnas para facilitar la comprensión.

En la columna de las *features* características por ausencia se acumulan la mayoría de los valores, como era de esperar. En la figura 27a, vemos una respetable cantidad de *features* que eran representativas por ausencia para seres vivos pasan a ser representativas por presencia en *instrumentos*, además, este fenómeno se acentúa en *synsets* más concretos, como vemos en la figura 27b. Un ejemplo de este comportamiento podrían ser *features* que detectaran ruedas, a la hora de clasificar, sería representativo para seres vivos que no tuvieran ruedas, mientras que el *synset* de los instrumentos, contiene los vehículos, para los que esta *feature* sería característica por presencia. Para Perro de caza, sin embargo, la cantidad de *features* de valor -1 que pasan a valer 1, es significativamente menor. Las que hay podrían ser *features* bastante concretas, que se asociaran a perros de caza, pero no al conjunto general de seres vivos, como por ejemplo ser cierta raza de perro.

Finalmente en la columna de las *features* características por presencia, vemos en la figura en la que se conservan más es 27c, mientras que en el resto casi todas pasan a ser características por ausencia. De forma simétrica al comportamiento de la columna de las *features* características por ausencia, pero a menor escala.

El siguiente ejemplo (28) es para ver el comportamiento en *synsets* más concretos. Podemos ver claramente que la red está captando la similitud entre Perro y Perro de caza, puesto a que casi todas las *features* se mantienen en la diagonal. Se comporta de forma parecida con Transporte y Vehículo. Mientras que entre Perro de Caza y Vehículo la mayoría de las *features* pasan de características a no características y viceversa. Nótese la similitud de la figura 28c y la figura 27a, aun siendo en el primer caso los *synsets* más concretos y en el segundo los más generales, el comportamiento es parecido.

Lo estudiado en esta sección nos hace revocar la hipótesis 3, puesto a que la proporción de *features* representativas, tanto por ausencia como por presencia, se conserva entre *synsets*. Sin embargo, hemos visto que en *synsets* más concretos las *features* características por presencia están más concentradas, lo que nos

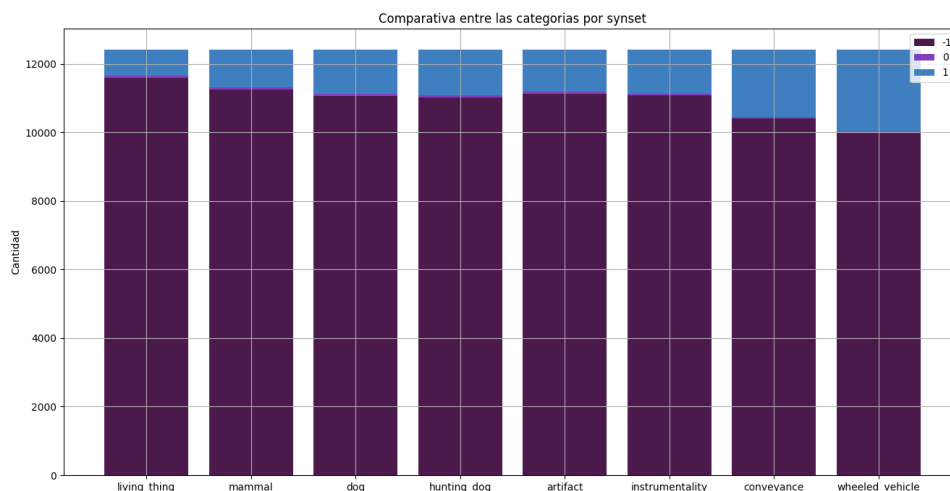


Figura 24: Cantidad de *features* de cada tipo de los representantes de los distintos *synsets*

indica que la red está distinguiendo el nivel de concreción que tienen los distintos *synsets*, sólo que no es de la manera que esperábamos en un primer momento.

5.2 Del full network embedding a wordnet

Partiendo de los resultados del apartado 5.1.3 hemos generado una distancia que representa lo diferentes que son cada uno de los *synsets* que aparecen en *imagenet* para la red convolucional entrenada. En este apartado explicaremos como se ha obtenido y como se relaciona con lo estudiado en el apartado anterior, para de esta forma decidir si podemos aceptar o revocar la hipótesis 4.

Para empezar, definimos los siguientes conjuntos.

Definición 5.1. Llamaremos \mathcal{C} al conjunto de *synsets* correspondientes a cada una de las 1,000 clases posibles en las que clasifica nuestra red convolucional.

Llamaremos \mathcal{I} al conjunto de todos los *synsets* que pertenecen a \mathcal{C} junto con todos los *synsets* de *wordnet* que tienen algún hipónimo en \mathcal{C} .

Al conjunto de *synsets* de *wordnet* lo denotaremos \mathcal{W} .

Por construcción tenemos:

$$\mathcal{C} \subset \mathcal{I} \subset \mathcal{W}$$

Nuestro objetivo es construir una métrica para \mathcal{I} que dependa de la red convolucional entrenada. Para de esta forma poder medir lo diferentes que son dos *synsets* en términos de la red.

Para ello recordamos la definición de métrica:

Definición 5.2. Sea X un conjunto, una **métrica** o distancia es una función:

$$d : X \times X \rightarrow [0, \infty) \subset \mathbb{R}$$

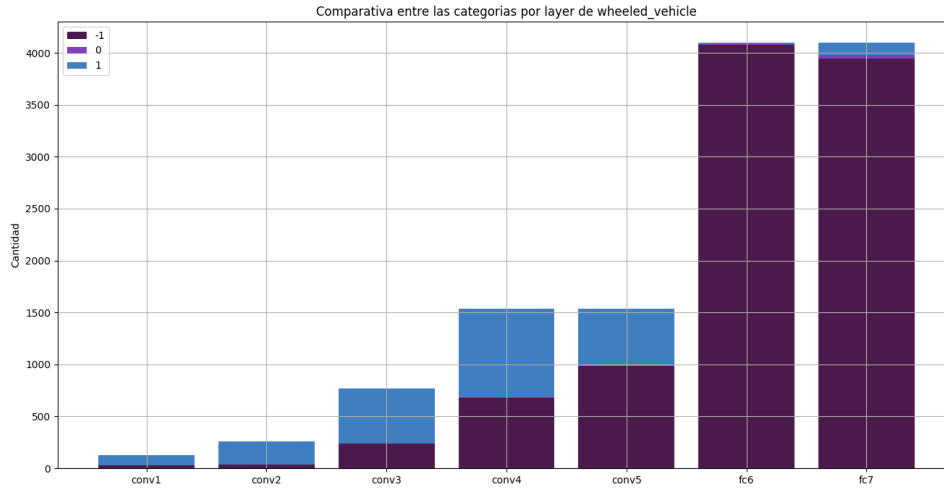


Figura 25: Cantidad de *features* de cada tipo del representante del *synset* Vehículo por capa.

Donde dados $x, y, z \in X$ se cumple:

1. $d(x, y) \geq 0$
2. $d(x, y) = 0 \iff x = y$
3. $d(x, y) = d(y, x)$
4. $d(x, z) \geq d(x, y) + d(y, z)$

La forma natural de conectar diferentes *synsets* sería a partir de un árbol utilizando las conexiones de hiponimia e hiperonimia, similar la de la figura 9b y utilizar la distancia discreta. Sin embargo, una distancia que utilizase solamente esta información no dependería de la red.

Para solventar este problema utilizamos la información obtenida en la sección 5.1.3, donde vimos que se podía definir matrices de cambio a partir de los representantes de cada *synset*. Estas matrices nos pueden dar una idea de lo parecidos que son dos *synsets* en términos de la red.

Otro problema que nos encontramos es que, de la forma en la que hemos definido \mathcal{J} se puede dar el caso de que dos *synsets* diferentes tengan las mismas *features* (y en consecuencia la misma matriz). Por ejemplo, si tomamos un *synset* que sólo contenga un hipónimo en \mathcal{C} y lo comparamos con éste, los dos tendrían exactamente el mismo vector representante. Esto hace que no podamos utilizar la definición de métrica estándar, puesto a que fallaría la propiedad 2, así que utilizaremos una **pseudométrica**.

Definición 5.3. Sea X un conjunto, una **pseudo-métrica** es una función:

$$d : X \times X \rightarrow [0, \infty) \subset \mathbb{R}$$

Donde dados $x, y, z \in X$ se cumple:

1. $d(x, y) \geq 0$
2. $d(x, x) = 0$

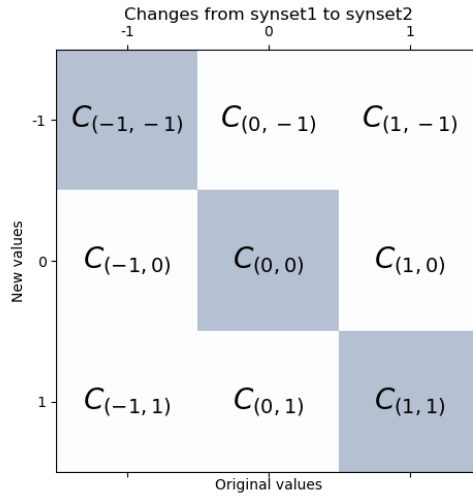


Figura 26: Matriz de cambios general

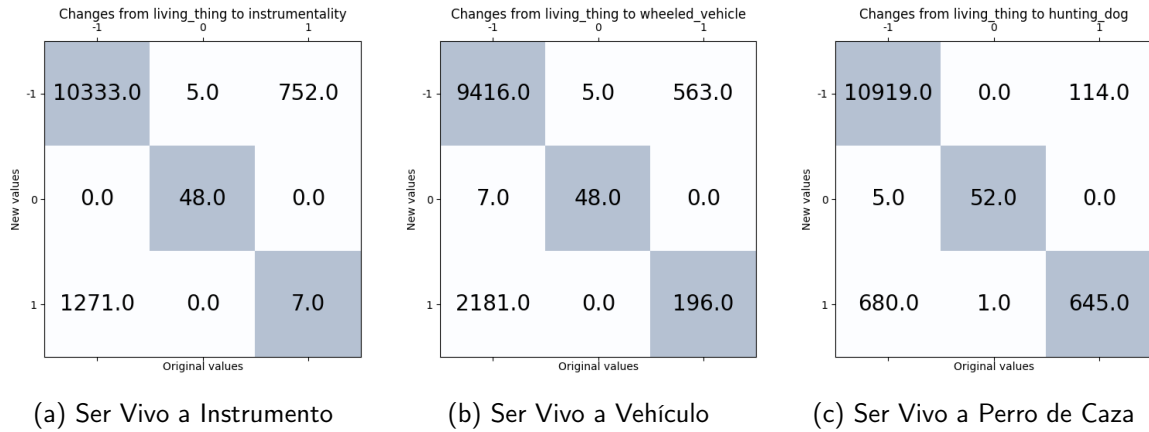


Figura 27: Matrices de cambio de Ser Vivo

3. $d(x, y) = d(y, x)$

4. $d(x, z) \geq d(x, y) + d(y, z)$

Como nuestro objetivo es medir lo diferentes que son los *synsets* para la red, una pseudo-métrica es suficiente. La definiremos utilizando la notación introducida en la figura 26:

$$d(s_1, s_2) = 1 - \frac{C_{(1,1)}(s_1, s_2)}{C_{(1,-1)}(s_1, s_2) + C_{(1,0)}(s_1, s_2) + C_{(1,1)}(s_1, s_2) + C_{(1,1)}(s_1, s_2) + C_{(0,1)}(s_1, s_2) + C_{(-1,1)}(s_1, s_2)}$$

Ecuación 17. Distancia entre el *synset* s_1 y el s_2

Donde $s_1, s_2 \in \mathcal{S}$. Es trivial comprobar que esta función cumple las propiedades necesarias para ser pseudo-métrica. Puesto a que se basa en una proporción, los valores de la pseudo-métrica estarán entre cero

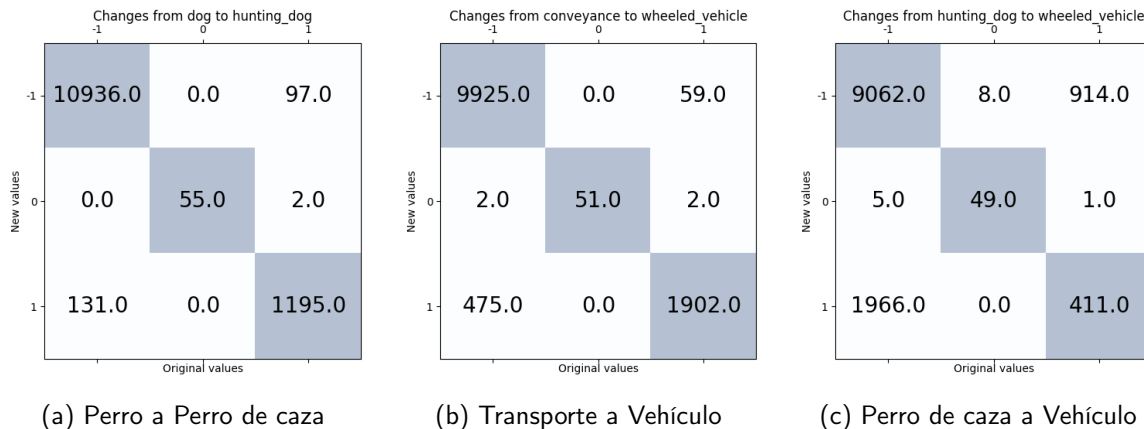


Figura 28: Matrices de cambio

y uno.

Esta distancia corresponde a la proporción de *features* características por presencia que se conservan entre s_1 y s_2 respecto a la suma de *features* características por presencia de los dos *synsets*. De esta manera podemos medir la similitud de dos *synsets* respecto al *embedding* y por tanto, respecto a la red convolucional entrenada.

Si la utilizamos para comparar los ejemplos de las figuras 27 y 28 tenemos las siguientes distancias:

$$\begin{aligned}
 d(\text{Ser vivo, Instrumento}) &= 0,9965 & d(\text{Ser vivo, Vehículo}) &= 0,9333 \\
 d(\text{Ser vivo, Perro de caza}) &= 0,5520 & d(\text{Perro, Perro de caza}) &= 0,1614 \\
 d(\text{Transporte, Vehículo}) &= 0,2201 & d(\text{Perro de caza, Vehículo}) &= 0,8751
 \end{aligned}$$

Donde vemos claramente que capta bien la similitud entre los distintos *synsets*, siendo los más diferentes entre si Ser vivo e Instrumento y los más similares perro y perro de caza. Esta información nos lleva a aceptar la hipótesis 4, puesto a que estamos viendo que *synsets* cercanos a nivel de hiponimia o hiperonimia comparten una significativa cantidad de *features*.

Para terminar de visualizar las distancias hemos generado el árbol de hipónimos de diferentes *synsets* utilizando la psedudo-métrica definida como distancia entre los nodos. Para ello hemos hecho un BFS en el grafo de *wordnet* tomando como aristas la relación de hiponimia y como nodo padre el *synset* del que queríamos generar el árbol. Se podrían obtener también las distancias entre el resto de nodos y generar un grafo completo, pero esto dificultaría la visualización, y en consecuencia, la obtención de conclusiones.

Un ejemplo de ello sería la figura 30, donde vemos que la distancia entre Perro con Perro de caza, y Perro con Perro de trabajo son inferiores que las de razas concretas de perro con sus respectivos hiperónimos. Lo que significa que, dentro del conjunto de perros, para la red se parecen más *synsets* más generales que razas concretas.

Para tener una idea de las relaciones entre razas concretas hemos calculado todas las distancias entre los diferentes *synsets* del conjunto de hipónimos de Perro de caza. Un ejemplo de ello serían las razas de la figura 29, donde vemos tres razas de perro claramente diferenciadas.



(a) Spaniel



(b) Greyhound



(c) Water Spaniel

Figura 29: Ejemplos de razas

Sus respectivas distancias son:

Cuadro 4: Distancias entre las razas del ejemplo

	Spaniel	Grayhound	Water Spaniel
Spaniel	0	0.7371	0.6442
Grayhound	0.7371	0	0.8330
Water Spaniel	0.6442	0.8330	0

En la tabla 4 podemos ver que el Spaniel se parece más al Water-spaniel que al Grayhound, cosa que es coherente con la imagen 29.

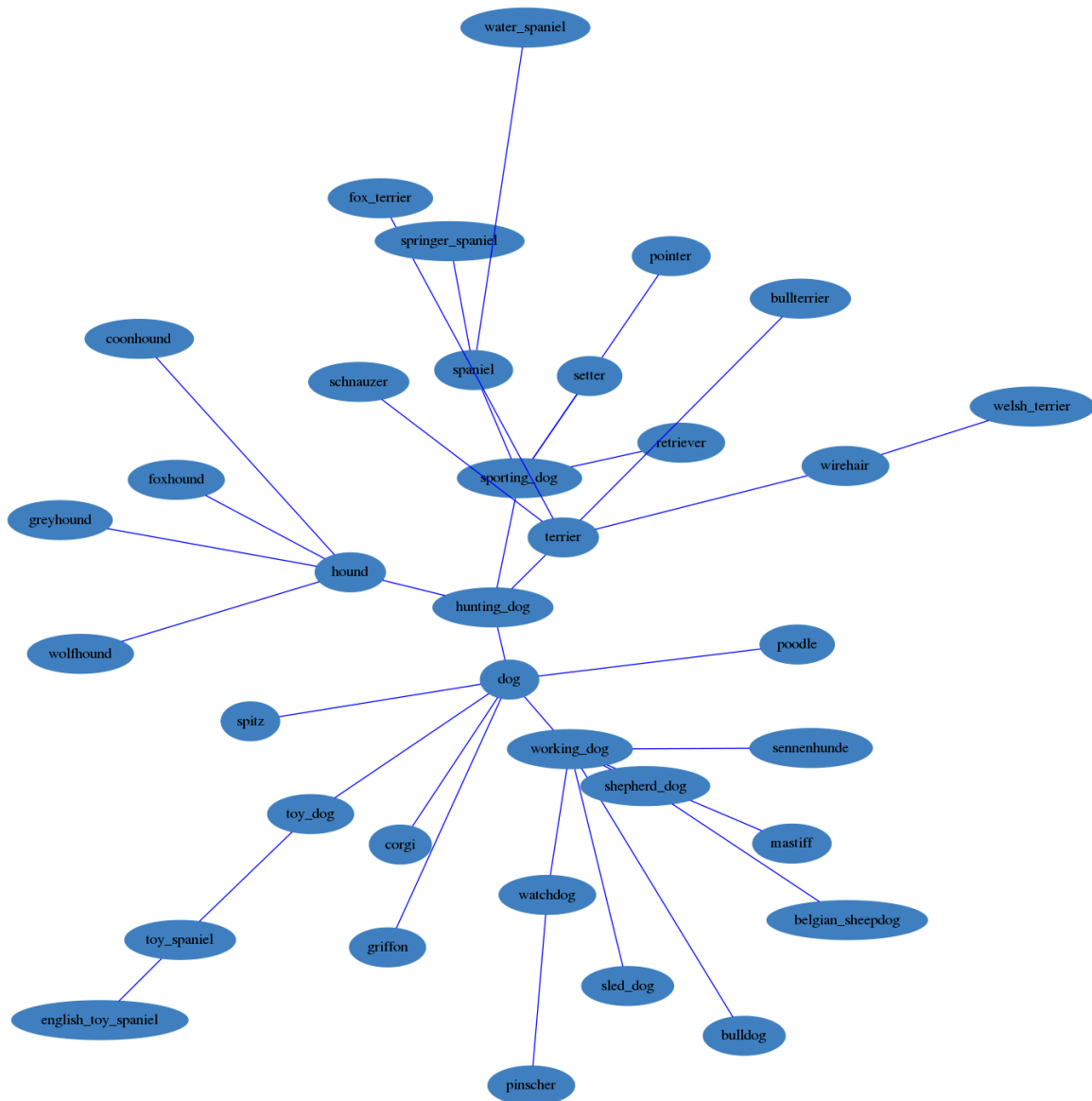


Figura 30: Árbol del *synset* perro donde las aristas están normalizadas utilizando la pseudo-métrica.

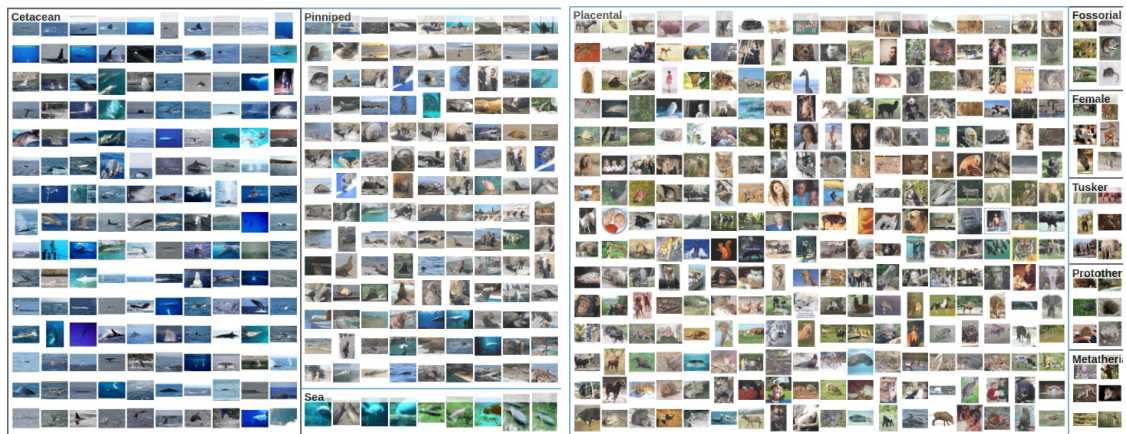
6. Conclusiones

En esta sección explicaremos las conclusiones obtenidas a partir de las secciones anteriores. Además revisaremos las hipótesis planteadas en 4. Los puntos base de la conclusión son los siguientes:

- Las capas *fully-connected* tienen una distribución diferente de las convolucionales, tanto a nivel general como por *synset*. Lo que implica que aun después de la estandarización hecha a las capas convolucionales, éstas mantienen sus diferencias respecto a las *fully-connected*. Además en todas las distribuciones se ve una cierta simetría entre las gráficas de la categoría de las *features* representativas por ausencia y las representativas por presencia. Por estos motivos aceptamos la hipótesis 1.
- Las *features* no características no solo son minoría, sino que además no están concentradas para ninguna *feature* concreta. De esta forma vemos que todas las *features* contienen información que caracteriza el espacio de representación.
- Las proporciones de las *features* características por ausencia se mantienen respecto a las profundidades de las capas. Por tanto vemos que la proporción se ve afectada por el tipo de capa, no por la profundidad de ésta. Lo que nos lleva a revocar la hipótesis 2.
- La proporción de las *features* características por ausencia, presencia o no características se mantiene respecto a los diferentes *synsets*. Sin embargo, en el estudio utilizando representantes hemos visto que estas se distribuyen en grupos, es decir, aunque en la proporción global las *features* características por ausencia siempre son mayoría, para ciertas *features* concretas hay una mayoría de *features* características por presencia, y por tanto las proporciones son significativamente diferentes a nivel de representante, donde se observan los siguientes fenómenos:
 - Las *features* representativas por presencia se concentran en las capas convolucionales, sobretudo en la cuarta y la quinta. Mientras que en las *fully-connected* se concentran las representativas por ausencia.
 - Cuanto más concreto es el *synset* mayor concentración de *features* representativas contiene, aunque no tenga la mayor proporción que habíamos supuesto en la hipótesis 3.

Puesto a que la proporción a nivel global era similar en los *synsets* (como vimos en la tabla 3) esto significa que, a diferencia de lo que pensábamos en un primer lugar, en las capas convolucionales cuanto más concreto es el *synset* o más profunda es la capa más se agrupan las *features* características por presencia en *features* concretas. Esto nos hace revocar la hipótesis 3.

- El *embedding* detecta similitud a nivel de *synset*. Si comparamos representantes de diferentes *synsets* observamos que entre *synsets* cercanos a nivel de hiponimia e hipernimia coinciden muchas más *features* que en *synsets* lejanos. Por tanto aceptamos la hipótesis 4
- Utilizando la pseudo-métrica definida podemos medir esta distancia y representarla gráficamente. De esta forma obtenemos información sobre como está decidiendo que clase corresponde a que imagen. Por ejemplo, entre mamífero y mamífero acuático (figura 31) se veía una gran distancia. Al comparar ambos conjuntos de imágenes se puede ver que mientras los mamíferos acuáticos tienen fondo azul por el mar, los mamíferos suelen contener colores más marrones o verdes. De esta forma podemos ver que algo que está teniendo en cuenta para clasificar es si están en el mar o no.



(a) Mamífero Acuático

(b) Mamífero

Figura 31: Ejemplos de synsets

Referencias

- [1] Christopher M. Bishop. *Pattern recognition and machine learning*. Springer, New York :, 2006.
- [2] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [3] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- [4] Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2010.
- [5] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? In *Advances in Neural Information Processing Systems 27*, pages 3320–3328, 2014.
- [6] Hossein Azizpour, Ali Sharif Razavian, Josephine Sullivan, Atsuto Maki, and Stefan Carlsson. Factors of transferability for a generic convnet representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(9):1790–1802, 2016.
- [7] Ali Sharif Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. Cnn features off-the-shelf: an astounding baseline for recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 806–813, 2014.
- [8] Dario Garcia-Gasulla, Ferran Parés, Armand Vilalta, Jonathan Moreno, Eduard Ayguadé, Jesús Labarta, Ulises Cortés, and Toyotaro Suzumura. On the behavior of convolutional nets for feature extraction. *CoRR*, abs/1703.01127, 2017.
- [9] Jeff Donahue, Yangqing Jia, Oriol Vinyals, Judy Hoffman, Ning Zhang, Eric Tzeng, and Trevor Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. *CoRR*, abs/1310.1531, 2013.
- [10] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012.
- [11] Dario Garcia-Gasulla, Armand Vilalta, Ferran Parés, Jonathan Moreno, Eduard Ayguadé, Jesús Labarta, Ulises Cortés, and Toyotaro Suzumura. An out-of-the-box full-network embedding for convolutional neural networks. *CoRR*, abs/1705.07706, 2017.

A. Gráficas de las distribuciones de las capas por *synset*

En este apéndice añadiré las gráficas de la distribución del número de *features* con los varios valores categóricos para los distintos *synsets* estudiados.

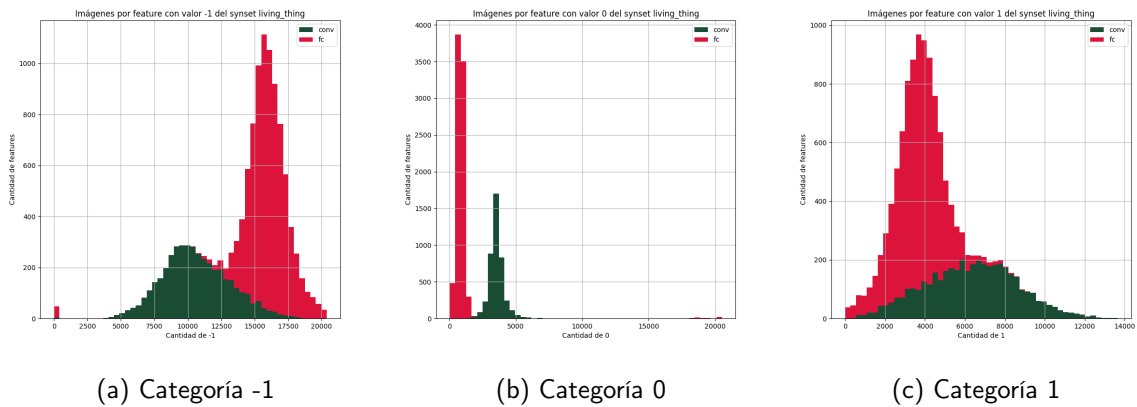


Figura 32: Ser Vivo

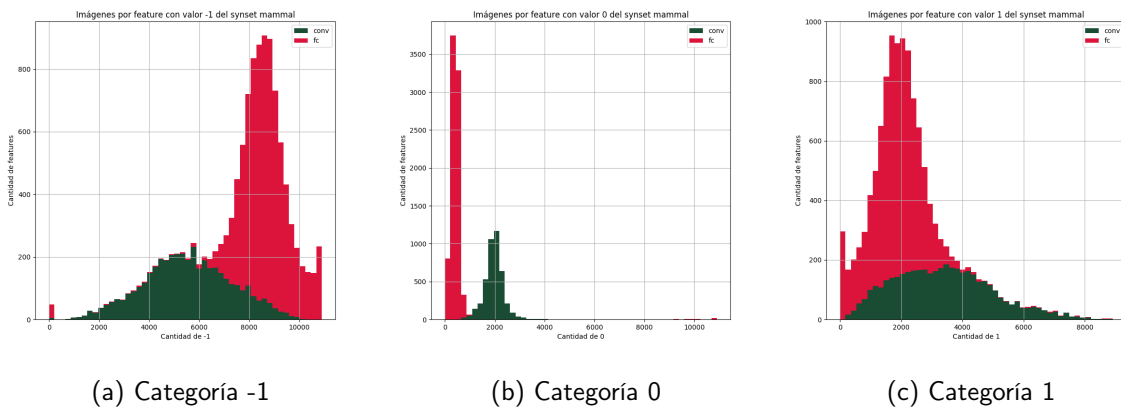
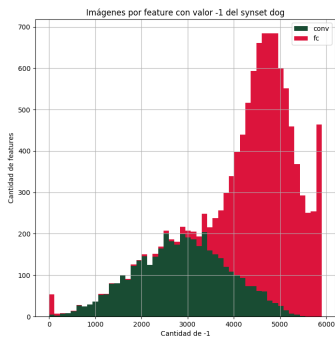
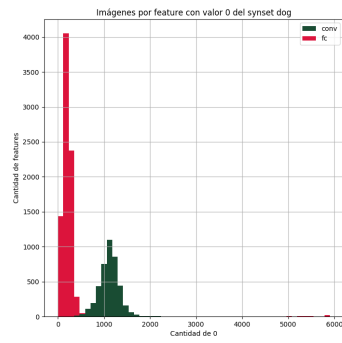


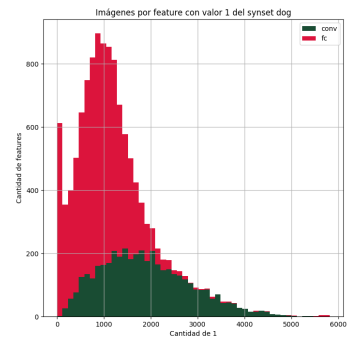
Figura 33: Mamífero



(a) Categoría -1

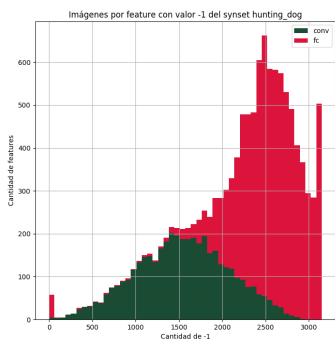


(b) Categoría 0

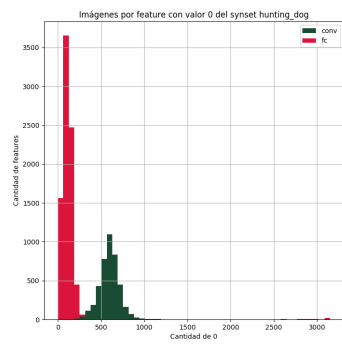


(c) Categoría 1

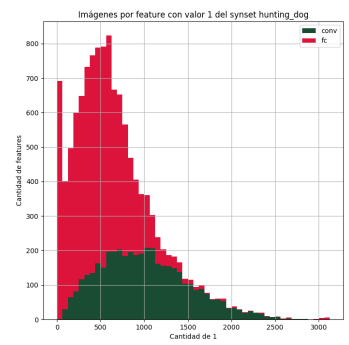
Figura 34: Perro



(a) Categoría -1

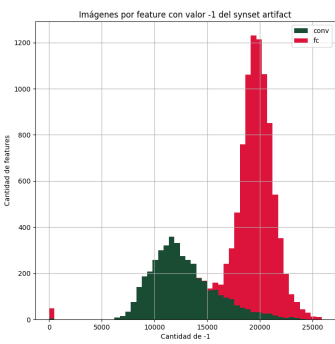


(b) Categoría 0

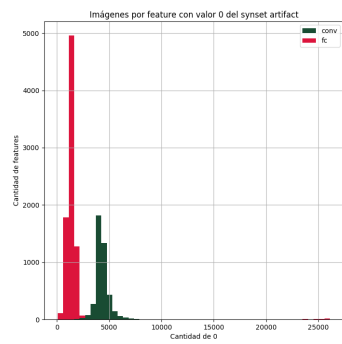


(c) Categoría 1

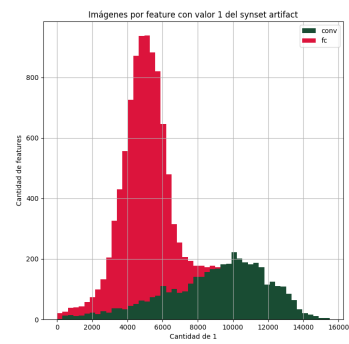
Figura 35: Perro de caza



(a) Categoría -1

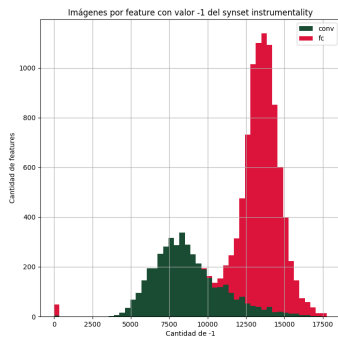


(b) Categoría 0

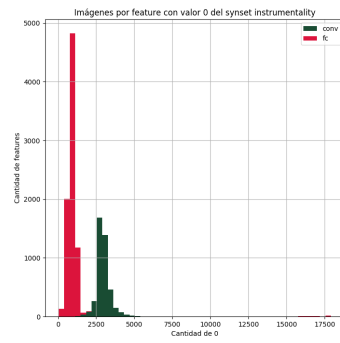


(c) Categoría 1

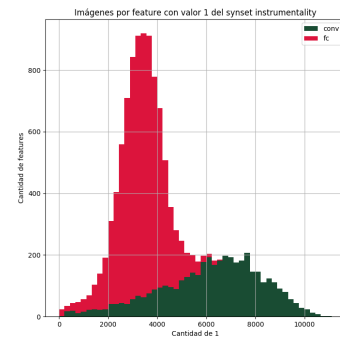
Figura 36: Artefacto



(a) Categoría -1

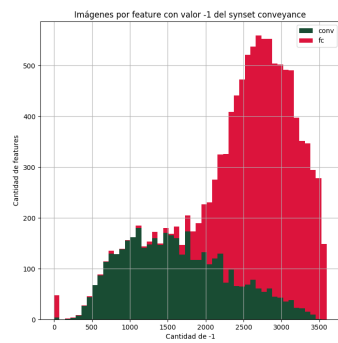


(b) Categoría 0

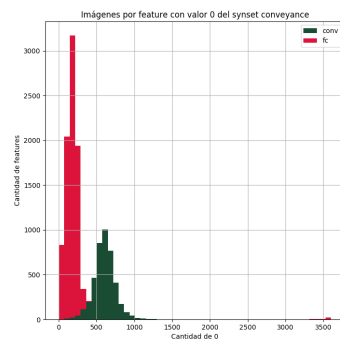


(c) Categoría 1

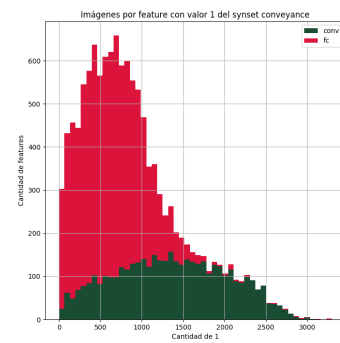
Figura 37: Instrumento



(a) Categoría -1

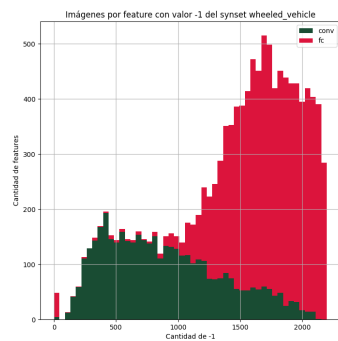


(b) Categoría 0

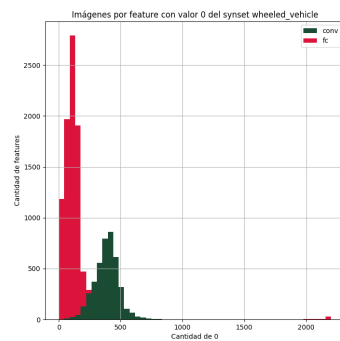


(c) Categoría 1

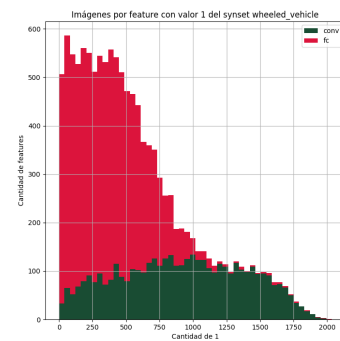
Figura 38: Transporte



(a) Categoría -1



(b) Categoría 0



(c) Categoría 1

Figura 39: Vehículo