

UNIVERSITY OF PRINCE EDWARD ISLAND

HONOURS THESIS

**Fashion Forward-Propagation:  
Machine Learning Tackles Fashion  
Curation**

*Hailey LeClair*

supervised by  
Dr. Andrew GODBOUT

November 9, 2017

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Images . . . . .	3
1.1.1	Objects in Images . . . . .	4
1.2	Other Difficulties in Object Recognition in images . . . . .	6
<b>2</b>	<b>Filters and Feature Maps</b>	<b>9</b>
<b>3</b>	<b>Machine Learning</b>	<b>13</b>
3.1	Training Sets and Classifiers . . . . .	13
3.2	Neural Networks . . . . .	15
3.3	PUT A FIGURE OF MAX POOLING HERE . . . . .	18

# List of Figures

1.1	A simple black and white image represented by pixel values 0 and 1 . . . . .	4
1.2	Clothing shown as a Deformable Object . . . . .	5
1.3	Two photos of Nike Rosch One sneaker. One being worn by a person on cobblestone, and one with a white background . . . . .	7
2.1	Image from Viola and Jones well known paper, Robust Real-Time Face Detection[5]	10
2.2	Sobel Operator: Convolution Kernels $G_x$ and $G_y$ [2] . . . . .	10
2.3	A Convolution kernel being applied to an image [?] . . . . .	11
3.1	What is returned from Imagenet after searching for dress: <a href="http://image-net.org/synset?wnid=n03201638">http://image-net.org/synset?wnid=n03201638</a> . . . . .	14
3.2	A simple neural network with 3 input neurons, 1 hidden layer with 4 neurons, and 2 output neurons . . . . .	15

# Chapter 1

## Introduction

Imagine an up and coming clothing designer with a small boutique. They have a small inventory of about 100 different items, and only carry a few different types of items. Say dresses, shirts, shoes, and jeans. The designer wants to provide their customers with an easy way to find the clothing they make. To do this, they want to use an app in which the user can take a photo of any clothing item, and if the boutique has an item that looks similar, that item will be shown to the customer. Is this feasible for a small boutique with minimal technical resources? Can a computer that has five images of different dresses take another image of one of these dresses and find its match? This paper will attempt at answering these questions by exploring how an app like this might be possible using convolutional neural networks, and how a designer with minimal technical knowledge could use this sort of app.

To look at this problem on a very high level, let us assume we have already created an app for these purposes. The designer would need to provide the app with several images of each of their items, so that this app, can learn things about these items within the image. As humans, we can easily distinguish the defining features of an object by looking at the object itself, or by looking at an image of the object. Without turning this into a biology paper, the neurons in our brain that are connected to our eyes are very intelligent and complex, and can recognize certain features of objects and objects themselves[7]. For example, in two different images, we can easily recognize a similar red dress on two different women with different body shapes, in different positions, with different scenery around them in the images. For a computer, it is not so easy.

### 1.1 Images

To understand how it might be possible for a computer to recognize a dress in an image we must first understand how an image is represented on a computer. Computers see images as two dimensional data structures of pixels, as in the image on the left of Figure 1. Each pixel has a numeric value. In this simple image, each pixel is either 0 or 1, representing either white or black pixels, respectively. A greyscale image can similarly be represented, where the value is a number between 0 and 255, but larger, and images in colour are obviously a

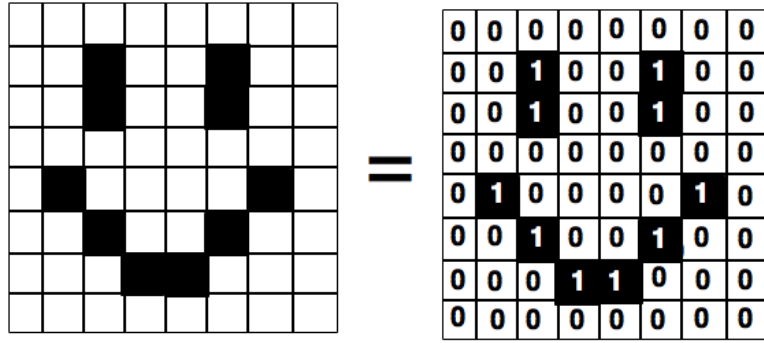


Figure 1.1: A simple black and white image represented by pixel values 0 and 1

lot more complicated than this trivial example. They often have hundreds or thousands of pixels and each pixel is represented by an RGB (red, green, and blue) value. RGB values are a set of 3 values each of which is a number from 0 to 255. Because an RGB pixel has 3 values, and an image is already a two dimensional structure, RGB images are seen as 3 dimensional by a computer. For this reason, images take more processing power than many other types of data. [11].

The field of computer vision emerged in the early 1970's as an attempt to get a computer to process and view digital images in the same way our brain and eyes do. In 1966, a professor at MIT asked a student to spend the summer with a camera connected to a computer, and to get the computer to explain what it was seeing through the camera. This kind of experiment was bound to fail with the state of technology in 1966, but led to much more research into how images can be represented digitally, how scenes are perceived visually, and how objects are recognized in a scene or image. Much of this research was done using 3D mathematical modelling, algorithms for digital image processing, and image and scene analysis. [11]

### 1.1.1 Objects in Images

Since the idea here is to take an image, classify it, and find its match amongst a set of similar photos, this means that we need to be able to recognize the object in another context in another image, which is exactly what is done in the field of computer vision. Rigid, articulated, and deformable objects all need to be considered when doing object detection, recognition, or classification in images. In terms of physics, a rigid object is a body in which "the distance between any two given points on [the] rigid body remains constant in time regardless of external forces exerted on it "[9]. In terms of images, this means that a rigid object will always have the same dimensions, or some translation of the same dimensions, in any given photo taken of this object. There exist many methods discovered within the field of computer vision to recognize rigid objects in an image, one of which is Scale Invariant Feature Transform (SIFT)[4].

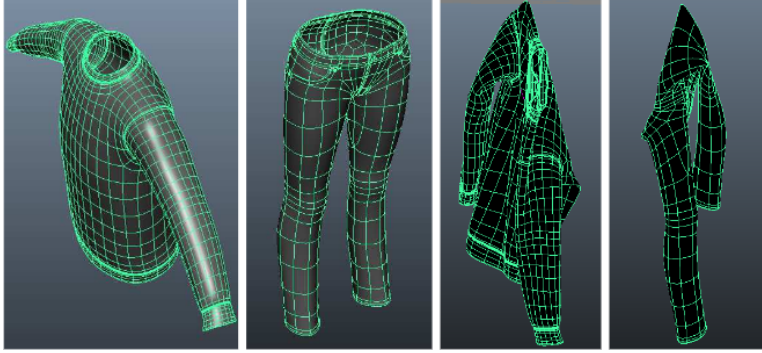


Figure 1.2: Clothing shown as a Deformable Object

SIFT can recognize any objects which do not vary, based on image rotation, scaling, or translation (rigid objects). It also does well with different kinds of lighting and illumination, as well as 3D projection[4]. The "key" to SIFT, is something called "key localization". Key localization chooses "key locations at maxima and minima of a difference of Gaussian function applied in scale space." [4]. The term "scale space" is used here, because any scaling of the image will not effect the locations selected during this process. Without delving too deeply into SIFT, this process is possible by building image pyramids at different levels of scaling. These pyramids make the locations very stable, which aids greatly in recognizing objects in an image.[4].

This is a step in the right direction in the attempt to recognizing clothing in images. Theoretically, if we already have something in place that can take an image and tell us what category of clothing an object belongs to, we can then use these methods like SIFT to match rigid objects like shoes. It is tremendously useful that we have access to methods, but is the same type of object detection possible for clothing that changes shape because of the articulated nature of the human body? Can we recognize a deformable dress that looks completely different on the hanger than it does while being worn? To attempt to answer these questions, we need to look at the possibilities surrounding these difficult types of objects.

Articulated objects may have rigid parts, but like human body parts, are connected at a basic joint which can move[11]. Multiple images of the same articulated object may show the object in different formations or positions. Although parts of these objects are rigid, possible movements from the joints mean that many different positions are possible for one articulated object. Think of all the possibilities of position for a human arm. An arm can be up in the air, straight out to the side, hanging down by the side of the body, bent at the elbow, behind the back, just to list a few. Even though articulated objects are rigid to a point, or, a joint, the degrees of freedom involved with these objects mean they are much more difficult to detect in images.

A deformable object, in contrast with the others, "changes its shape and/or volume while being acted upon by any kind of external force"[8]. Depending on the degree to which an object deforms, or the amount of force applied, many images of an object could

show it as having completely different dimensions and shape. As with articulated objects, the uncertainty of the general shape or dimensions makes deformable objects especially challenging. This problem, and the problem with articulated objects is particularly relevant with clothing.

With some exceptions (like shoes, jewelry, some accessories, for example), many everyday clothing items are deformable to some degree. For instance, a pair of skinny jeans may have a general shape when lifeless on a hanger, but take on very different dimensions and possibly even a totally different shape when being worn by different people. Like these jeans, of course, many other clothing items will take on the shape of the person wearing them. An example of this is shown in figure 2.[?] In terms of shape and dimension of these two clothing items being worn and not being worn, there are barely any similarities. We can also consider changes in an environment that could change the shape of an object. If it is a particularly windy day, a long flowing dress may take on a completely different shape in the wind.

As mentioned earlier, clothing items are typically being worn by someone in images. This simple fact adds one extra dimension to the recognition of any given clothing item, because the human body is an articulated body. A loose denim jacket, which may have a similar shape on anyone who wears it, can be easily deformed by placing the arms in a different position. Different articulations of the joints in the human body (especially our knees, shoulders, and elbows) can dramatically affect the shape and dimensions of a piece of clothing, and this carries over to the digital image. These obstacles add a layer to object recognition in clothing items, which would not be present in recognizing various other types of objects in images. Unfortunately, we cannot rely on something like SIFT[4] to recognize these sorts of items in an image, as it struggles to detect these nuances.

Along with the problem of clothing being deformable and worn on an articulated body, we also have to consider the other contextual aspects of an image. As intelligent humans, we can normally determine if a dress in an image is the same as another dress in another image, regardless of what other objects are in the background. When it comes to digital recognition, though, an image containing an explicit clothing item could be convoluted with other objects or people in the same image, or could have an explicit, plain background which becomes a defining feature of an image. When introduced to another image of the same object with a similarly noisy background, a computer may then think this is a different object altogether, because of the dramatic difference in backgrounds between this photo and the previous one(s).

## 1.2 Other Difficulties in Object Recognition in images

Figure 3 shows the same nike shoe in two different images. We can even speculate here that the photo has been taken from almost the same angle in both images. With a quick glance at both of these images, we can confidently say that the shoes in these images are in fact the same, but as learned by many in a first year computer science course: a computer does exactly what you tell it to do. This is relevant because unless a computer knows to look for a background in an image, it will assume these two images are of completely different objects.



Figure 1.3: Two photos of Nike Rosch One sneaker. One being worn by a person on cobblestone, and one with a white background



The angle that a photo was taken and the location of a particular item in the image are also important considerations. All known images of an object could be front-facing, and if a side view of this object is introduced, a computer will have nothing to compare with this side view. In this event, the computer will not be able to correctly classify or recognize the object, since it has no previous knowledge of its side view.

When worn, clothing items like shoes and hats are almost always near the bottom and top of an image, respectively (for obvious reasons). If a computer has only seen images of shoes in which they are located near the bottom edge of an image, the computer may think that this location is a common, fundamental characteristic of shoes themselves! This must be considered when choosing images for a computer to classify or match. A shoe, of course, can be anywhere in an image, and recognizing this when choosing images lets the focus be on features of the shoes themselves, rather than the object location.

Even if articulation, deformability, angle, and rotation have all been accounted for when attempting to classify and recognize objects, colour and illumination must also be added to this long list of considerations. Depending on the lighting conditions and the camera itself, the same red shirt may look pink in one image and orange in another (in fact, there was a viral internet sensation over a blue and black, or white and gold dress, a few years ago, due to this very perceptual issue). An image must be recognizable under different lighting conditions and in light of minor variations in the object's colour.

Yet, what if an image of this shirt is in greyscale? This is an important question. If the goal is to recognize a red shirt in a given image, greyscale images cannot be considered. If they are considered, they must be accepted as either grey or colourless, even if the shirt in the image is actually red. In this case, the greyscale image is irrelevant. On the other hand, If the goal is to recognize a generic T-shirt, colour becomes unimportant, and a greyscale image can be as relevant as an RGB image. When it comes to recognizing clothing in images, it is naturally almost impossible to account for every possible variant from amongst these images. However, if we keep all of these degrees in freedom in mind when attempting to classify or match images, we can avoid many of the problems that come with such variations.

# Chapter 2

## Filters and Feature Maps

Where do we start in terms of finding similarities in different images? One approach, and a good approach, is to find and identify specific features of an object which are common to almost every image of the object. Viola and Jones did this in 2004, with Real-Time Robust Face Detection. They classified faces in images, based on the computed values of simple features specific to the faces.[5]. They found that simple rectangles (figure of Viola and Jones rectangles) could locate areas of brightness in a photo. If these specific areas of brightness are present, then the image contains a face.

For our purposes, Viola and Jones' specific facial recognition, by design, only works with images of faces. Naturally, though, one would imagine that a similar process to this facial recognition could be used in attempting to recognize other objects in images, like clothing. In theory, we should be able to find prominent features that belong to any given type of clothing item, and use this feature to consistently classify the image.

As stated above, in images, we count on an item to have one or more identifying features, like a human face. One of the ways that we can find the parts in an image that are likely to be important, and make them more prominent, is by applying a filter to the image.

A common example of a filter is the use of a neighbourhood operator, or local operator on pixels in an image. A neighbourhood operator uses an area of pixels around a specific pixel in an image to compute an output value for that pixel. When applied to all pixels across the image, the operator "filters" the image. A neighbourhood operator takes a weighted sum of all pixels in a neighbourhood. It is also known as a correlation, where  $h(k, l)$  are labelled the filter coefficients.[11]

$$g(i, j) = \sum_{k, l} f(i + k, j + l)h(k, l) \quad (2.1)$$

More relevant for the purposes of this paper is the convolutional operator. Convolution is a variation on correlation(filtering) and is both commutative and associative. It is a reversal of the sign of the offsets of the correlation function  $g = f * h$  where  $h$  is now an impulse response function.[11]. To explain this not just in formulas, but with respect to pixels in images, a convolution on a single pixel is done by taking a weighted sum of the pixels in the neighbourhood of this pixel. The weights, and size of the neighbourhood depend on the

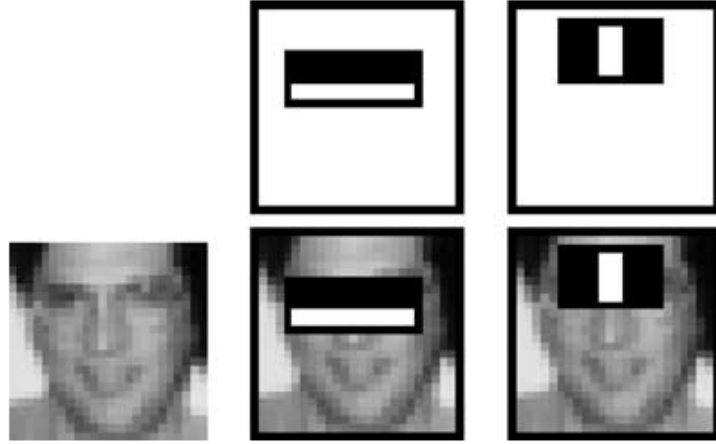


Figure 2.1: Image from Viola and Jones well known paper, Robust Real-Time Face Detection[5]

-1	0	+1
-2	0	+2
-1	0	+1

**G<sub>x</sub>**

+1	+2	+1
0	0	0
-1	-2	-1

**G<sub>y</sub>**

Figure 2.2: Sobel Operator: Convolution Kernels  $G_x$  and  $G_y$  [2]

convolution kernel[1].

$$g(i, j) = \sum_{k, l} f(k, l)h(i + k, j + l) \quad (2.2)$$

A convolution kernel can be any size, but is usually quite small. Since we only want to use pixels within a certain neighbourhood around a pixel, this makes sense. A convolution kernel too small will not consider any pixels around it, and too large will end up losing the importance of the pixel itself. A common convolutional operator is the sobel operator, whose 2 convolution kernels are shown in figure 5. This filter is used for edge detection and is very helpful in finding shapes, objects, and defining features in images. Two 3x3 convolution kernels are used to perform this edge detection[2]. On their own,  $G_x$  highlights vertical lines in an image and  $G_y$  highlights horizontal lines in an image. To apply one convolution kernel across an image, for instance  $G_x$ , we need to inversely apply the weights in the kernel around the neighbourhood of a given pixel (a convolution) and take the sum of these values. In the

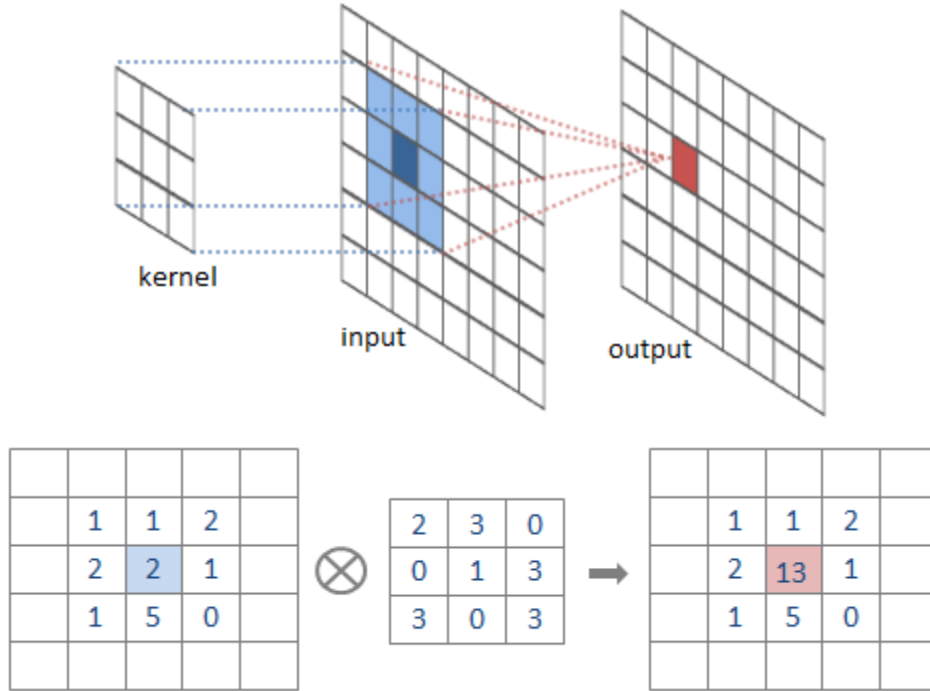


Figure 2.3: A Convolution kernel being applied to an image [?]

case of Sobel's  $G_x$  kernel, the pixel in question will be at (2, 2) or in the middle of the kernel. This weighted sum then becomes the value for that pixel. This is done across all pixels in an image[?]. An example of this can be seen in figure 6. The image on the top shows a high level view of a convolution being applied, and underneath shows a convolution kernel being applied to a set of pixels in the neighbourhood of a number 2 (in the middle of the 3x3 section of pixels).

What is important to note is that this operation cannot be done on the pixels around the edges or corners of an image, because their neighbourhood is limited. To deal with this problem, often zero padding is added, which adds a border of pixels all with the value 0 around the image. In this way, the original edge pixels in an image can then be part of the convolution. [?, ?, ?] Another approach is to avoid applying the convolution kernel to the outer pixels, and in this case these pixels will essentially be thrown away after the convolution. This should only be done if one is assuming that the pixels around the edge of the image are unimportant, which is often the case.

After applying a convolution operator, correlation operator, or any type of filter to an image, we are left with a feature map. A feature map is essentially a filtered image. Since each filter highlights a certain feature, or set of features in the image, a feature map can provide us with a lot of information about what is in an image, and what is important in an image[7]. Having multiple feature maps, which can show us if images have a number of different features, these can then be used to help us learn things about an image.

As outlined above, filters can alter an image in dramatic ways. They can be used for

edge detection, to blur photos, to adjust colour, or to sharpen images.[11]. Filters can play a crucial part in detecting any object in an image, as a feature map helps a computer to recognize certain features in an image. Even so, to manually find which features are important and which filters or convolution kernels work for each object is not always a feasible task. Filters are helpful in discovering features in an image and making them more prominent, but when dealing with 1000 or more images, how do we know which filters to use and what features are important?

# Chapter 3

## Machine Learning

Machine Learning can help us "learn" things about these images, namely which filters or convolution kernels work best, without having to manually go through and picking filters for each clothing type. Simply put, Machine learning is programming a computer to learn something based on the data itself. People use machine learning to analyze and monitor data much more efficiently than executing it manually. It recognizes patterns, and can be used to learn what people are buying, which emails in your inbox are spam, and in this case, if a pair of jeans is in an image, and indeed which pair of jeans it is. [7].

### 3.1 Training Sets and Classifiers

Using machine learning algorithms, we can program a computer to learn about almost anything as long as we provide it with enough data, and the right data. For instance, if we want a computer to use machine learning to "learn" if a pair of shoes are present in an image, we have to make sure that we provide it with many images of many different kinds and brands of shoes, in different positions and orientations. These images are called our training set.[7]. The images in the training set are input into a machine learning algorithm called a "classifier"[10]. The algorithm then goes through the entire set of images and looks for similar features among all or most of the images.

Choosing a training set can be difficult, to properly train a machine learning classifier we need (as stated earlier) a lot of the right data. That might be impossible for a small boutique with only 100 items in total. Because most people do not have a database with thousands of photos, there must be some resource providing photos to take advantage. And there is. Imagenet [?]Imagenet) is a website that gives access to millions of pre-classified images, 14,197,122 at this moment, to be exact. An entire database of images of a specific classification can be downloaded, or all images can be accessed by urls. Something like this will be useful for the clothing designer in question. Instead of having to personally take thousands of photographs to have a decent sized training set for the clothing items that this designer has, Imagenet can be used to easily get whatever images are needed, whenever they are needed. Figure 7 shows a screen shot of imagenet in which 1227 images are available

## Dinner dress, dinner gown, formal, evening gown

A gown for evening wear

1227  
pictures

82.76%  
Popularity  
Percentile

Wordnet  
IDs

Numbers in brackets: (the number of synsets in the subtree).

- ImageNet 2011 Fall Release (32326)
  - plant, flora, plant life (4486)
  - geological formation, formation (1112)
  - natural object (1112)
  - sport, athletics (176)
  - artifact, artefact (10504)
    - instrumentality, instrumentation (1405)
    - structure, construction (1405)
    - paving, pavement, paving material (650)
    - sheet, flat solid (115)
    - layer, bed (13)
    - facility (4)
    - lemon, stinker (0)
    - fabric, cloth, material, textile (21013)
    - covering (1013)
      - thumb (0)
      - imbrication, overlapping, lap (0)
      - finger (0)
      - folder (2)
      - upholstery (0)
      - artificial skin (0)
      - mask (2)
      - cover plate (0)
      - paddle box, paddle-box (0)

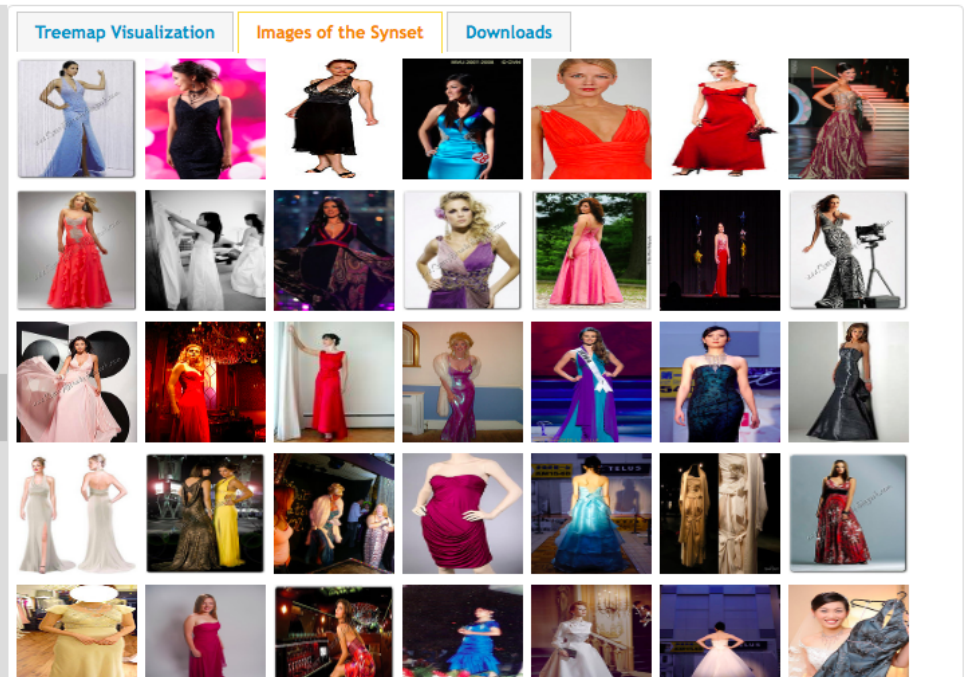


Figure 3.1: What is returned from Imagenet after searching for dress: <http://image-net.org/synset?wnid=n03201638>

and were classified as "Dinner dress, dinner gown, formal, evening gown". If this clothing designer specializes in something like fancy dresses, they can take these images as part of their training set that represents fancy dresses, or any of the labels previously mentioned. Even with the tremendous help of Imagenet[?]Imagenet) providing a good training set, we need a classifier to recognize (or correctly classify) a clothing item in an image. Once a classifier is trained, it gives us a probability that the image belongs to each of the classes it has been trained on. Whichever probability is the highest, is the predicted class of an item[10]. If the designer with a boutique has photos of their clothing ready on hand, the theoretical app and classifier should be able to take images of all of the designer's clothes, and give each image a label, or class. Before being able to do this though, the classifier needs to learn the way in which it will get the best results with the data (in this case, images) available. This data is how a machine learning algorithm will be trained, hence the name "training set". After being trained, an algorithm should be able to predict the class of an object in an image, but only for objects that it has already been trained on. In other words, as long as the classifier was trained on dresses, the predicted class, or label of a photo of a dress, should be one that refers to a dress.

In this case, supervised learning will be used. With supervised learning, the algorithm

already knows the class, or label, that each specific data element in the training set belongs to, and makes use of this for training purposes. If a classifier predicts that a data element belongs to the wrong class during the training phase, supervised learning lets the algorithm run again with the same data to attempt to correct its' previous mistakes[7]. Thus, a classifier learns by making predictions on the same training set continuously until it can accurately predict the classes of all or most data elements in the training set.

How can supervised learning help in predicting which clothing item is in an image? Since we do not have to get everything right the first time with supervised learning, we can use a method that makes mistakes each time (or each epoch) that it iterates through our dataset. When a classifier is in its' training phase, an epoch refers to one iteration through the training data. Depending on the way in which it is trained and with which and how much data, a classifier may need many many training epochs to be properly trained. The number of training epochs needed will depend on the type of classifier (neural network, nearest neighbour, SVM, etc), the number of possible classes, and the data. So how can we tell when a classifier is properly trained? To measure how well a classifier has been trained we use a testing set of data (again, which in this case will be images) in which all elements are different from those in the training set. The classifier will iterate through the testing set and will keep track of the actual prediction it makes for the element, and the correct prediction. After the classifier has made a prediction for each element in the testing set, we can see how well it has done by calculating precision, accuracy, and recall. We can look at each of these as a percentage.

## 3.2 Neural Networks

Neural networks are an effective way to do this. Like most machine learning algorithms or in this case deep learning algorithms, neural networks take one data element as input, and outputs a probability that this element belongs to each class in a list of classes. What is unique about neural networks is that they are inspired by the human brain. Neural networks are made up of a series of layers, and each layer is made up of one or more neurons, like the neurons in a human brain. Every neural network has an input layer made up of input neurons, an output layer made up of output neurons, and 1 or more hidden layers composed of hidden neurons.

For the purpose of explanation, an example will be given, with only one hidden layer. In the figure below, there is three input neurons, or signals:  $x_1$ ,  $x_2$ , and  $x_3$ . Each neuron in the input layer is connected to every neuron in first hidden layer. As can be seen in figure 8, every neuron in one layer is connected to every neuron in the next layer. Each one of these connections has its' own unique weight (weights are important) and a neuron in the first hidden layer takes a weighted sum of all input neurons. This weighted sum is then given to an activation function in each hidden layer. The output of this activation function is the output value of this hidden neuron. This is done for each neuron in the first hidden layer. If there are more hidden layers, the output of the previous hidden layer's neurons are fed into the neurons in these layers as input. This continues until the output layer is reached, where



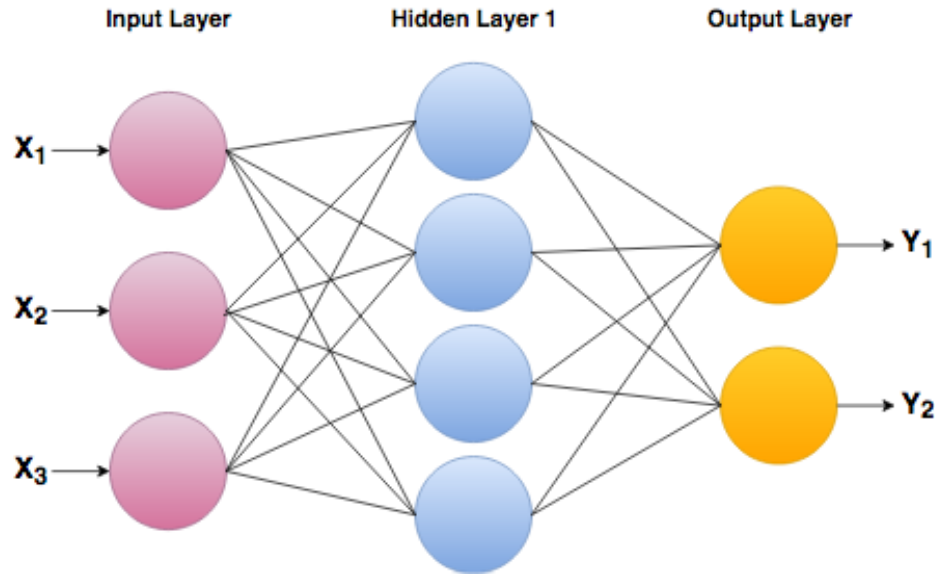


Figure 3.2: A simple neural network with 3 input neurons, 1 hidden layer with 4 neurons, and 2 output neurons

each output neuron represents a specific class and calculates the probability that the input belongs to this class.[10]

A Neural network has input neurons and output neurons, but how does it learn, and how is it supervised? Back-propagation takes the same input through a network to correct any mistakes made in previous predictions (probabilities of classification). This is where the training set comes in.

During each Epoch of a training session, each element in the training set is forward-propagated through the neural network. Since we know the correct classification of all elements in our training set, after each epoch we can compute the error in the classification. This can be done using the weights. Since all neurons contribute independently to the classification of the input, each input or hidden neuron in a particular layer has a different weight to carry it to each neuron in the next layer. All weights should be updated by adjusting for the error in the previous epoch. A neural network should improve at predicting the actual class of its input after every training epoch. The error of an entire output layer can be calculated using accuracy, precision, or recall. Each one of these calculations can tell us how close the network was to making a correct prediction. Once a neural network is properly trained, the same weights can be used over and over to make predictions. This process of using error to back-propagate through a network using weights can help a network to learn almost anything.[10]

Since neural networks can learn almost anything, one would expect that they could learn which objects are in an image, or if an image contains a particular clothing item. They can, but with a bit of an adjustment. A specific type of deep neural network is used for image classification called a convolutional neural network.

Convolutional neural networks have a similar flow to other neural networks. The input is all pixel values in an image and the output is again the probability that the input belongs to a certain class. In this case the output will be the probability that the input image contains a certain object. The hidden layers are where CNNs get more complicated. There are a few possible types of hidden layers. The first and probably most important to the network is a convolutional layer. A convolutional layer applies a number of convolutions to an image each return a feature map[7]. Convolutional layers try to mimic the way that our eyes perceive what is around us. When we look at something, the neurons in our visual cortex are only really focused on a small area of the visual field called the local receptive field. Each neuron has a different local receptive field, although they may overlap. These local receptive fields of specific neurons, which are how our eyes process what they perceive, can be similar to filters in many ways. Some neurons have larger local receptive fields than others, whereas some can focus on and react to lines or other shapes. This is exactly what a convolutional layer attempts to do.[7]

The first convolutional layer will not have every input neuron connected to every other neuron in the layer like we see in the hidden layers of some other neural networks. Each neuron in the first convolutional layer is connected to a certain number of neurons in the input image which are referred to as the neuron's local receptive field. The weights applied here are similar to a small image in equal size to the local receptive field of a neuron. This small image of weights that are applied to an image is a convolution kernel. A convolution happens here when this kernel is applied to all neurons local receptive field. This process acts as a sort of inverse filter when applied throughout an entire image. This process of convolution on local receptive fields from a previous layer, helps the network focus on low level features in the beginning stages of forward propagation through the network, and higher level features near the end of the network.[7]

When all of these convolutions are applied in a convolutional layer, we are left with feature maps (or filtered images) of the previous layer. Because these layers return a number of feature maps, the network learns which features are best at predicting the class of a particular image using back-propagation. A feature which helps in predicting the class will be more heavily weighted in the network. whereas a feature which does not contribute to a prediction will be minimally weighted.

Since neural networks apply an activation function after all data has entered a neuron, we would expect CNNs to do the same. There are several different activation functions that a CNN can use on the image data of a layer, but the most common one is the ReLU function. In biological neurons, the sigmoid function gives the best learning results, but with deep neural networks, there are other activation functions that give better results and speed up the training time.[6]

For CNNs and some other deep neural networks, The ReLU (Rectified Linear Unit) function is preferred over other functions like sigmoid, because it does not saturate for positive values and gives 0 when any input is negative. Since the output of some neurons will undoubtedly be 0 after computing ReLU, we can assume that many neurons will end up dropping out, or having no effect on the classification. This is okay for our purposes, but

there are some good variations on ReLU that make sure neurons do not get killed in the training process. [7]

$$(x) = \max(0, x) \quad (3.1)$$

One variation on ReLU activation is LeakyReLU. It uses a hyperparameter (or tuneable parameter)  $\alpha$  to determine the "leakage" of the function. In this case, the neuron may become temporarily inactive, but may eventually contribute to the network again, which is an improvement over ReLU. Another variation on ReLU is PReLU (Parametric Leaky ReLU) which takes the hyperparameter in LeakyReLU and makes it a parameter learned with backpropagation during the network's training. These variations perform well when the training set is large, but may overfit on smaller datasets. Since the dataset being considered for this project is quite small, ReLU will be fine for our training purposes. The activation function is important to note here because each convolutional layer in a CNN is usually followed by a layer in which the activation function (ReLU) is applied to all neurons in the previous layer.[7]

The next hidden layer commonly used in convolutional neural networks, and other deep networks is a pooling layer. A pooling layer is used to summarize a large amount. In the case of images, this means reducing the size of the image, while still keeping certain parts or features of the image.[7]. What these features are will determine what type of pooling layer is used. In most CNNs, max pooling layers are used to keep the most prominent features in an image, which is usually desirable when attempting to recognize objects. Max pooling layers summarize the image data by taking the largest pixel value in a group of pixels and using this to represent the group. Doing this across an entire image gives a smaller image with less noise. Max pooling layers are usually in between or after convolutional layers (after ReLU has been computed for the convolutional layer). Throughout the network the classification of an image becomes more clear, so it makes sense that after each convolution, which gives a set of feature maps, we pool these maps which gets us closer and closer to finding the most important and defining features in an image. The last pooling layer is usually followed by one or more fully connected layers. Fully connected layers are actually quite simple and are often placed at the end of the network. Each neuron in a previous layer is connected to every neuron in the fully connected layer, hence the name "fully connected" [3].

$$\text{LeakyReLU}_{\alpha}(z) = \max(\alpha z, z) \quad (3.2)$$

### 3.3 PUT A FIGURE OF MAX POOLING HERE

# Bibliography

- [1] Bringing parallelism to the web with river trail.
- [2] Sobel edge detector.
- [3] Leonardo Aurajosantos. Artificial intelligence.
- [4] David G Lowe. Object recognition from local scale-invariant features. 2:1150–1157, 1999.
- [5] Paul Viola and Michael J Jones. Robust real-time face detection. *International journal of computer vision*, 57(2):137–154, 2004.
- [6] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [7] Geron Aurelien. *Hands-On Machine Learning with Scikit-Learn & TensorFlow*. O’Rielly Media, Inc., 2017.
- [8] <http://scienceworld.wolfram.com/physics/DeformableBody.html>. *DeformableBody*.
- [9] [https://en.wikipedia.org/wiki/Rigid\\_body](https://en.wikipedia.org/wiki/Rigid_body). *Rigid Body*.
- [10] Miroslav Kubat. *An Introduction to Machine Learning*. Springer International Publishing Switzerland, 2015.
- [11] Richard Szeliski. *Computer vision: algorithms and applications*. Springer Science & Business Media, 2010.