

# Chapter 1

## Background

### 1.1 Machine Learning

#### 1.1.1 Introduction

Machine Learning (ML) is a field which is raised out of Artificial Intelligence (AI). Applying AI, we wanted to build better and intelligent machines. But except for few mere tasks such as finding the shortest path between point A and B, we were unable to program more complex and constantly evolving challenges. There was a realisation that the only way to be able to achieve this task was to let machine learn from itself. This sounds similar to a child learning from its self. So machine learning was developed as a new capability for computers. And now machine learning is present in so many segments of technology, that we don't even realise it while using it.

Finding patterns in data on planet earth is possible only for human brains. The data being very massive, the time taken to compute is increased, and this is where Machine Learning comes into action, to help people with large data in minimum time.

There are three kinds of Machine Learning Algorithms :

1. Supervised Learning
2. Unsupervised Learning
3. Reinforcement Learning

#### Supervised Learning

A majority of practical machine learning uses supervised learning. In supervised learning, the system tries to learn from the previous examples that are given. Speaking mathematically, supervised learning is where you have both input variables ( $x$ ) and output variables ( $Y$ ) and can use an algorithm to derive

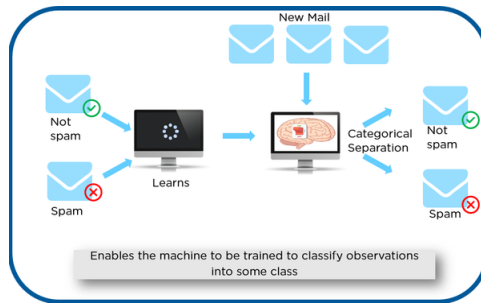


Figure 1.1: Example of supervised Learning

the mapping function from the input to the output. The mapping function is expressed as  $Y = f(x)$ .

As shown in Figure 1.1, we have initially taken some data and marked them as 'Spam' or 'Not Spam'. This labeled data is used by the training supervised model, in order to train the model. Once it is trained, we can test our model by testing it with some test new mails and checking of the model is able to predict the right output.

Supervised learning problems can be further divided into two parts, namely **classification**, and **regression**.

**Classification** : A classification problem is when the output variable is a category or a group, such as "black" or "white" or "spam" and "no spam".

**Regression** : A regression problem is when the output variable is a real value, such as "Rupees" or "height."

Some Supervised learning algorithms include:

- Decision trees
- Support-vector machine
- Naive Bayes classifier
- k-nearest neighbors
- linear regression

## Unsupervised Learning

In unsupervised learning, the algorithms are left to themselves to discover interesting structures in the data. Mathematically, unsupervised learning is when you only have input data ( $X$ ) and no corresponding output variables. This is called unsupervised learning because unlike supervised learning above, there are no given correct answers and the machine itself finds the answers. In Figure 1.2, we have given some characters to our model which are 'Ducks' and 'Not

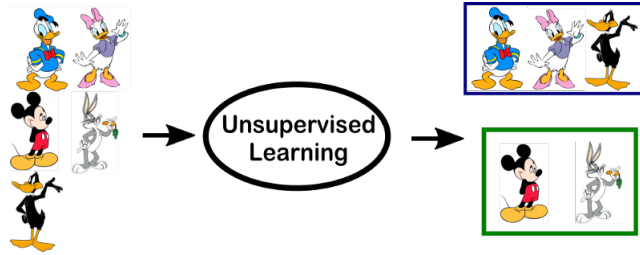


Figure 1.2: Example of unsupervised Learning

Ducks'. In our training data, we don't provide any label to the corresponding data. The unsupervised model is able to separate both the characters by looking at the type of data and models the underlying structure or distribution in the data in order to learn more about it. Unsupervised learning problems can be further divided into **association** and **clustering** problems.

**Association** : An association rule learning problem is where you want to discover rules that describe large portions of your data, such as "people that buy X also tend to buy Y".

**Clustering** : A clustering problem is where you want to discover the inherent groupings in the data, such as grouping customers by purchasing behaviour.

## Reinforcement Learning

A computer program will interact with a dynamic environment in which it must perform a particular goal (such as playing a game with an opponent or driving a car). The program is provided feedback in terms of rewards and punishments as it navigates its problem space. Using this algorithm, the machine is trained to make specific decisions. It works this way: the machine is exposed to an environment where it continuously trains itself using trial and error method. In Figure 1.3, we can see that the agent is given 2 options i.e. a path with water or a path with fire. A reinforcement algorithm works on reward a system i.e. if the agent uses the fire path then the rewards are subtracted and agent tries to learn that it should avoid the fire path. If it had chosen the water path or the safe path then some points would have been added to the reward points, the agent then would try to learn what path is safe and what path isn't

### 1.1.2 Neural Networks

Neural Networks are a class of models within the general machine learning literature. Neural networks are a specific set of algorithms that have revolutionized the field of machine learning. They are inspired by biological neural networks and the current so called deep neural networks have proven to work

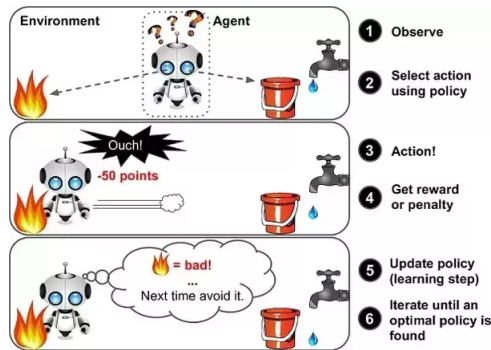


Figure 1.3: Example of Reinforcement Learning

quite very well. Neural Networks are themselves general function approximations, that is why they can be applied to literally almost any machine learning problem where the problem is about learning a complex mapping from the input to the output space.

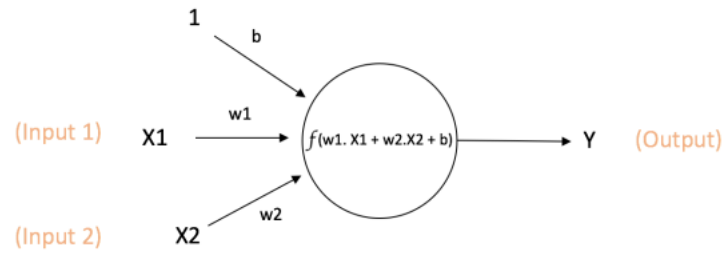
### 1.1.3 A single Neuron

The basic unit of computation in a neural network is the neuron, often called a **node** or **unit**. It receives input from some other nodes, or from an external source and computes an output. In purely mathematical terms, a neuron in the machine learning world is a placeholder for a mathematical function, and its only job is to provide an output by applying the function on the inputs provided. Each input has an associated weight ( $w$ ), which is assigned on the basis of its relative importance to other inputs. The node applies a function  $f$  (defined below) to the weighted sum of its inputs as shown in Figure 1.4. The network takes numerical inputs  $X1$  and  $X2$  and has weights  $w1$  and  $w2$  associated with those inputs. Additionally, there is another *input 1* with weight  $b$  (called *Bias*) associated with it. The main function of Bias is to provide every node with a trainable constant value (in addition to the normal inputs that the node receives). The output  $Y$  from the neuron is computed as shown in the Figure 1.4. The function  $f$  is non-linear and is called **Activation Function**. The purpose of the activation function is to introduce non-linearity into the output of a neuron. This is important because most real world data are non linear and we want neurons to learn these non-linear representations.

#### Activation Functions

Every activation function (or non-linearity) takes a single number and performs a certain fixed mathematical operation on it. There are several activation functions:

**Sigmoid** : takes a real-valued input and squashes it to range between 0 and



$$\text{Output of neuron} = Y = f(w1.X1 + w2.X2 + b)$$

Figure 1.4: An example of a single Neuron

1. Its formula is:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

It is easy to understand and apply but it has major reasons which have made it fall out of popularity:

- Vanishing gradient problem
- Its output isn't zero centered. It makes the gradient updates go too far in different directions.
- Sigmoids saturate and kill gradients.
- Sigmoids have slow convergence.

**Tanh** : takes a real-valued input and squashes it to the range  $[-1, 1]$ . Its formula is:

$$\tanh(x) = 2\sigma(2x) - 1$$

Now its output is zero centered because its range is between -1 to 1. Hence optimization is easier in this method and in practice it is always preferred over Sigmoid function. But still it suffers from Vanishing gradient problem.

**ReLU** : ReLU stands for *Rectified Linear Unit*. It takes a real-valued input and thresholds it at zero (replaces negative values with zero). So its formula is:

$$f(x) = \max(0, x)$$

It has become very popular in the past couple of years. It was recently proved that it had 6 times improvement in convergence from Tanh function. Seeing the mathematical form of this function we can see that it is very simple and efficient. A lot of times in Machine learning and computer science we notice that most simple and consistent techniques

and methods are only preferred and are best. Hence it avoids and rectifies vanishing gradient problem . Almost all deep learning Models use ReLu nowadays.

Figure 1.5 show each of the above activation functions.

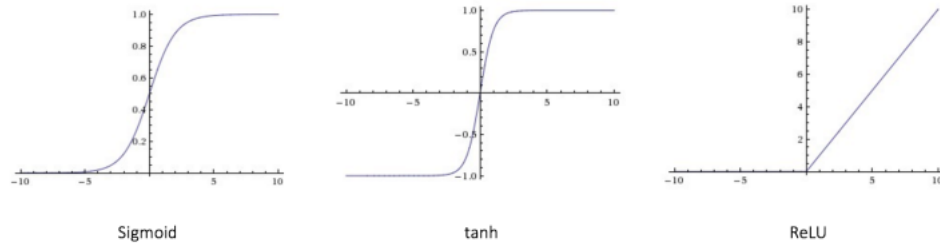


Figure 1.5: Plots of Activation functions

## Feedforward Neural Network

Till now we have covered neuron and activation functions which together for the basic building blocks of any neural network. The feedforward neural network was the first and simplest type of artificial neural network devised. It contains multiple neurons (nodes) arranged in layers. A layer is nothing but a collection of neurons which take in an input and provide an output. Inputs to each of these neurons are processed through the activation functions assigned to the neurons. Nodes from adjacent layers have connections or edges between them. All these connections have weights associated with them. An example of a feedforward neural network is shown in Figure 1.6. A feedforward neural network can consist of three types of nodes:

**Input Nodes** The Input nodes provide information from the outside world to the network and are together referred to as the “Input Layer”. No computation is performed in any of the Input nodes – they just pass on the information to the hidden nodes.

**Hidden Nodes** The Hidden nodes have no direct connection with the outside world (hence the name “hidden”). They perform computations and transfer information from the input nodes to the output nodes. A collection of hidden nodes forms a “Hidden Layer”. While a feedforward network will only have a single input layer and a single output layer, it can have zero or multiple Hidden Layers.

**Output Nodes** The Output nodes are collectively referred to as the “Output Layer” and are responsible for computations and transferring information from the network to the outside world.

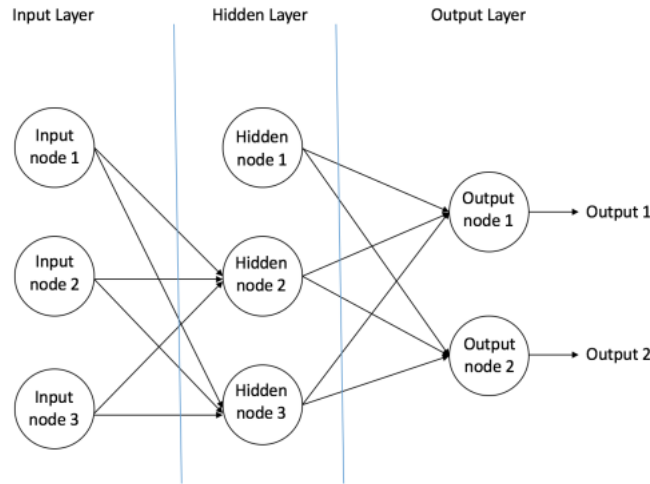


Figure 1.6: An example of a Feedforward Neural Network

In a feedforward network, the information moves in only one direction – forward – from the input nodes, through the hidden nodes (if any) and to the output nodes. There are no cycles or loops in the network (this property of feed forward networks is different from Recurrent Neural Networks in which the connections between the nodes form a cycle). Another important point to note here is that each of the hidden layers can have a different activation function, for instance, hidden layer1 may use a sigmoid function and hidden layer2 may use a ReLU, followed by a Tanh in hidden layer3 all in the same neural network. Choice of the activation function to be used again depends on the problem in question and the type of data being used.

#### 1.1.4 2D Convolutional Neural Network

A Convolutional Neural Network (ConvNet/CNN) is one of the variants of neural networks used heavily in the field of Computer Vision. It derives its name from the type of hidden layers it consists of. The hidden layers of a CNN typically consist of convolutional layers, pooling layers, fully connected layers, and normalization layers. Here it simply means that instead of using the normal activation functions defined above, convolution and pooling functions are used as activation functions. It can take in an input image, assign importance (learning weights and biases) to various aspects/objects in the image and be able to differentiate one from the other. The pre-processing required in a ConvNet is much lower as compared to the other classification algorithms. While in primitive method filters are hand-engineered, with enough training, ConvNets have the ability to learn these filters/characteristics.

The architecture of a ConvNet is analogous to that of the connectivity pattern of Neurons in the Human Brain and was inspired by the structure of

the Visual Cortex. However, most ConvNets consist mainly in 2 parts:

- **Feature extractor :**

This part of the network takes as input the image and extracts features that are meaningful for its classification. It amplifies aspects of the input that are important for discrimination and suppresses irrelevant variations. Usually, the feature extractor consists of several layers. For instance, an image which could be seen as an array of pixel values. The first layer often learns representations that represent the presence or absence of edges at particular orientations and locations in the image. The second layer typically detects motifs by spotting particular arrangements of edges, regardless of small variations in the edge positions. Finally, the third may assemble motifs into larger combinations that correspond to parts of familiar objects, and subsequent layers would detect objects as combinations of these parts.

- **Classifier :**

This part of the network takes as input the previously computed features and uses them to predict the correct label.

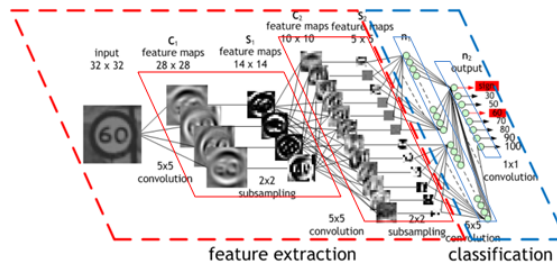


Figure 1.7: Typical structure of a ConvNet

**Convolutional Layers** In order to extract such features, ConvNets use 2D convolution operations. These operations take place in convolutional layers. Convolutional layers consist of a set of learnable filters. Every filter is small spatially (along width and height), but extends through the full depth of input. During forward pass, we slide (more precisely, convolve) each filter across the width and height of the input volume and compute dot products between the entries of the filter and the input at any position (as Figure 1.8 shows). The objective of the Convolution Operation is to extract the high-level features such as edges, from the input image. ConvNets need not be limited to only one Convolutional Layer. Conventionally, the first ConvLayer is responsible for capturing the Low-Level features such as edges, color, gradient orientation, etc. With added layers, the architecture adapts to the High-Level features as well, giving us a network which has the wholesome understanding of images in the dataset, similar to how we would.





Figure 1.8: Convolution with kernel of 3, stride of 2 and padding of 1

**Pooling Layers** Pooling Layers are also referred to as downsampling layers and are used to reduce the spatial dimensions, but not depth, on a convolution neural network. The intuitive reasoning behind this layer is that once we know that a specific feature is in the original input volume (there will be a high activation value), its exact location is not as important as its relative location to the other features. The main advantages of pooling layer are:

- We gain computation performance since the amount of parameters is reduced.
- Less parameters also means we deal with overfitting situations.

The pooling operation is specified, rather than learned. Two common functions used in the pooling operation are:

**Average Pooling** Calculate the average value for each patch on the feature map.

**Maximum Pooling (or Max Pooling)** Calculate the maximum value for each patch of the feature map.

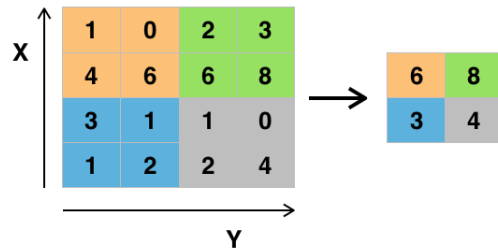


Figure 1.9: Example of Max pooling operation with a 2x2 filter and stride of 2

### 1.1.5 3D Convolutional Neural Network

Traditionally, ConvNets are targeting RGB images (3 channels). The goal of 3D CNN is to take as input a video and extract features from it. When ConvNets extract the graphical characteristics of a single image and put them in a vector (a low-level representation), 3D ConvNets extract the graphical characteristics of a set of images. 3D CNNs takes in to account a temporal dimension (the order of the images in the video). From a set of images, 3D CNNs find a low-level representation of a set of images, and this representation is useful to find the right label of the video (a given action is performed). In order to extract such features, 3D ConvNets use 3D convolution operations.

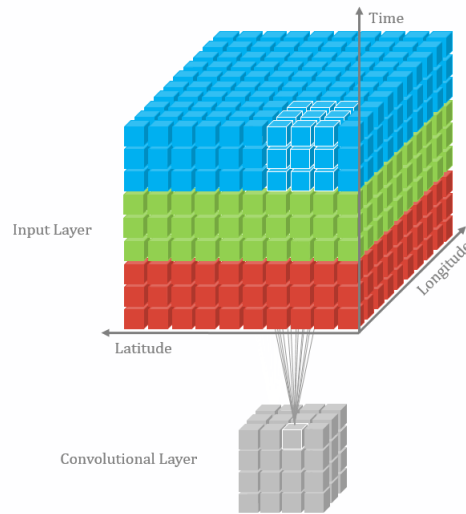


Figure 1.10: 3D Convolution operation

There are several existing approaches to tackle the video classification. This is a nonexhaustive list of existing approaches:

- **ConvNets + LSTM cell** : Extract features from each frame with a ConvNet, passing the sequence to an RNN
- **Temporal Relation Networks** : Extract features from each frame with a ConvNet and pass the sequence to an MLP
- **Two-Stream Convolutional Networks** : Use 2 CNN, 1 spatial stream ConvNet which process one single frame at a time, and 1 Temporal stream ConvNet which process multi-frame optical flow

## 1.2 Object Detection

Within the field of Deep Learning, the sub-discipline called “Object Detection” involves processes such as identifying the objects through a picture, video or a webcam feed. Object Detection is used almost everywhere these days. The use cases are endless such as Tracking objects, Video surveillance, Pedestrian detection etc. An object detection model is trained to detect the presence and location of multiple classes of objects. For example, a model might be trained with images that contain various pieces of fruit, along with a label that specifies the class of fruit they represent (e.g. an apple, a banana, or a strawberry), and data specifying where each object appears in the image.

The main process followed by most of CNN for Object Detection is:

1. Firstly, we do feature extraction using as backbone network, the first Convolutional Layers of a known pre-trained CNN such as AlexNet, VGG, ResNet etc.
2. Then, we propose regions of interest (ROI) in the image. These regions contain possibly an object, which we are looking for.
3. Finally, we classify each proposed ROI.

### 1.2.1 Region Proposal Network

From the 3 above steps, the 2nd step is considered to be very important. That is because, in this step, we should choose regions of the image, which will be classified. Poor choice of ROIs means that the CNN will pass by some object that are located in the image, because, they were not be proposed to be classified.

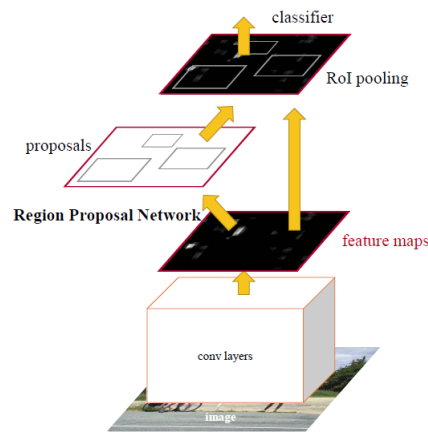


Figure 1.11: Region Proposal Network’s structure

The first Object-Detection CNNs use several algorithms for proposing ROIs. For example, R-CNN[?], and Fast R-CNN[?] used Selective Search Algorithm for extracting ROIs. One of novelties introduced by the Faster R-CNN[?] is **Region Proposal Network** (RPN). Its Function is to propose ROIs and its structure can be shown in 1.11. As we can see, RPN is consisted of:

- 1 2D Convolutional Layer
- 1 score layer
- 1 regression layer

Another basic element of RPN is the **anchors**. Anchors are predefined boxes used for extracting ROIs. In figure 1.12 is depicted an exaple of some anchors

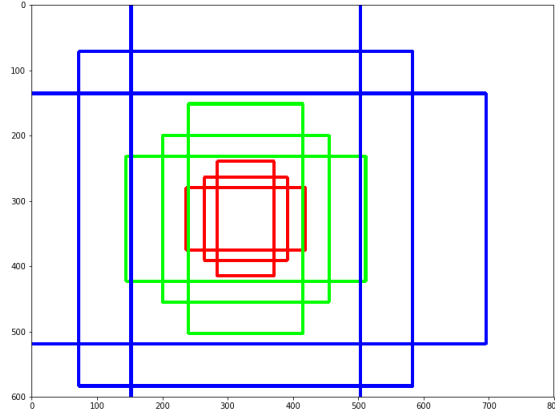


Figure 1.12: Anchors for pixel (320,320) of an image (600,800)

For each feature map's pixel corresponds **k** (**k=9**) anchors (3 different scales and 3 different ratios 1:1, 1:2, 2:1).

As a result, RPN gets as input feature maps extracted from the backbone CNN. Then performs 2D convolution over this input and passes the output to its scoring layer and regression layer. Those scores represent the confidence of existing an object in this specific position. On top of that, regression layer outputs 4k displacements, 4 for each anchor. Finally, we keep as output only the *n-best scoring* anchors.

### 1.3 Losses and Metrics

In order to train our model and check its performance, we use some known Loss functions and Metrics used in Object Detection systems.

### 1.3.1 Losses

For training our network, we use **Cross Entropy Loss** for classification layers and **smooth L1-loss** for regression.

#### Cross Entropy Loss

Cross-entropy loss, or log loss, measures the performance of a classification model whose output is a probability value between 0 and 1.

Entropy is the measure of uncertainty associated with a given distribution  $q(y)$  and its formula is:

$$H = - \sum_{i=1}^n p_i \cdot \log p_i$$

Intuitively, entropy tells us how surprised we are when some event  $E$  happened. When we are sure about an event  $E$  to happen ( $p_E = 1$ ) we have 0 entropy (we are not surprised) and vice versa.

On top of that, let's assume that we have 2 distributions, one known (our network's distribution)  $\mathbf{p}(\mathbf{y})$  and one unknown (the actual data's distribution)  $\mathbf{q}(\mathbf{y})$ . Cross-entropy tells us how accurate is our known distribution in predicting the unknown distribution's results. Respectively, Cross-entropy measures how accurate is our model in predicting the test data. Its formula is:

$$H_p(q) = - \sum_{c=1}^C q(y_c) \cdot \log(p(y_c))$$

#### Smooth L1-loss

Smooth L1-loss can be interpreted as a combination of L1-loss and L2-loss. It behaves as L1-loss when the absolute value of the argument is high, and it behaves like L2-loss when the absolute value of the argument is close to zero. It is usually used for doing box regression on some object detection systems like Fast-RCNN([?]), Faster-RCNN([?]) and it is less sensitive to outliers according to [?]. As shown in [?], its formula is:

$$smooth_{L1}(x) = \begin{cases} 0.5x^2 & \text{if } |x| < 1 \\ |x| - 0.5 & \text{otherwise} \end{cases}$$

It is similar to Huber loss whose formula is:

$$L_\delta(x) = \begin{cases} \frac{1}{2}a^2 & \text{for } |a| \leq \delta \\ \delta(|a| - \frac{1}{2}\delta), & \text{otherwise} \end{cases}$$

if we set  $\delta$  parameter equal 1.

Smooth L1-loss combines the advantages of L1-loss (steady gradients for large values of  $x$ ) and L2-loss (less oscillations during updates when  $x$  is small).

### 1.3.2 Metrics

Evaluating our machine learning algorithm is an essential part of any project. The way we choose our metrics influences how the performance of machine learning algorithms is measured and compared. They influence how to weight the importance of different characteristics in the results and finally, the ultimate choice of which algorithm to choose. Most of the times we use classification accuracy to measure the performance of our model, however it is not enough to truly judge our model.

At first, we introduce some basic evaluation metrics in order to present then those we use for our assesment.

### 1.3.3 Precision, Recall & F1 score

**Precision** measures how accurate is your predictions. i.e. the percentage of your predictions are correct.

**Recall** measures how good you find all the positives. For example, we can find 80% of the possible positive cases in our top K predictions.

Their definitions are:

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

$$F1 = 2 \cdot \frac{precision \cdot recall}{precision + recall}$$

where

- TP = True positive
- TN = True negative
- FP = False positive
- FN = False negative

**TODO describe how we use recall**

### 1.3.4 Intersection over Union

Intersection over Union (IoU) measures the overlap between 2 boundaries. We use that to measure how much our predicted boundary overlaps with the ground truth (the real object boundary). In some datasets, we predefine an IoU threshold (say 0.5) in classifying whether the prediction is a true positive or a false positive.

### 1.3.5 mAP

Precision and recall are single-value metrics based on the whole list of predictions. By looking their formulas, we can see that there is a trade-off between precision and recall performance. This trade-off can be adjusted by the softmax threshold, used in model's final layer. In order to have high precision performance, we need to decrease the number of FP. But this will lead to decrease recall performance and vice-versa.

As a result, these metrics fail to determine if a model is performing well in object detection tasks as well as action detection tasks. For that reason, we use mAP metric

**AP (Average precision)** is a popular metric in measuring the accuracy of object detectors like Faster R-CNN, SSD, etc. Average precision computes the average precision value for recall value over 0 to 1.

In our case, we use the metrics presented in [?] in order to quantify our results:

**frame-AP** measures the area under the precision-recall curve of the detections for each frame (similar to the PASCAL VOC detection challenge [?]). A detection is correct if the intersection-over-union with the ground truth at that frame is greater than and the action label is correctly predicted.

**video-AP** measures the area under the precision-recall curve of the action tubes predictions. A tube is correct if the mean per frame intersection-over-union with the ground truth across the frames of the video is greater than and the action label is correctly predicted.

**AUC** measures the area under the ROC curve, a metric previously used on this task. An action tube is correct under the same conditions as in video-AP. Following [?], the ROC curve is plotted until a false positive rate of 0.6, while keeping the top-3 detections per class and per video. Consequently, the best possible AUC score is 60%.

### 1.3.6 MABO

In order to evaluate more the quality of our proposals, both during TPN and connecting tube stages, recall metric isn't enough. That's because recall metric tells us only in how many known objects there was at least 1 proposal that satisfied the detection criterion. However, it doesn't tell us how much good. This can be better clarified if we consider figure .

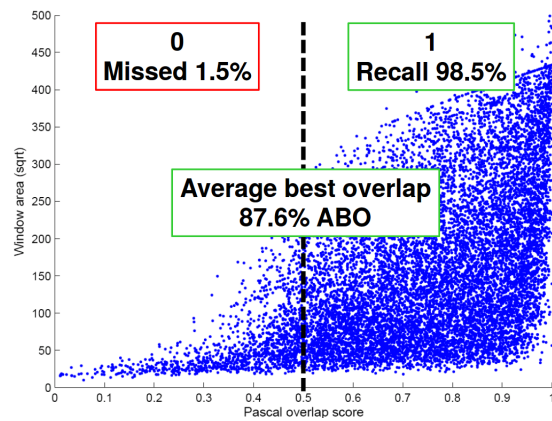


Figure 1.13: Example of Max pooling operation with a 2x2 filter and stride of 2