| House | Furniture | Nr rooms | New kitchen | Acceptable |
|-------|-----------|----------|-------------|------------|
| 1 | No | 3 | Yes | Yes |
| 2 | Yes | 3 | No | No |
| 3 | No | 4 | No | Yes |
| 4 | No | 3 | No | No |
| 5 | Yes | 4 | No | Yes |

1. Formalize the learning problem.

2. Describe how variables are chosen when building a Decision Tree.

3. Execute the ID3 algorithm on the data set above and generate the corresponding output tree.

   Note: point 3 can be answered even if point 2 is not properly addressed, by using any invented method (or invented numbers) for the selection of the variables.

Furniture → rooms → now Kitchen → acceptable

The target function is $f = \{0,1\} \times \mathbb{N} \times \{0,1\} \rightarrow \{0,1\}$

The entropy of a dataset is $-P_\oplus \ln P_\oplus - (1-P_\oplus) \ln(1-P_\oplus)$ where $P_\oplus =$ ratio of samples labeled as 1.

IF D is the set and A an attribute of the input, we have:

$$GAIN(D,A) = ent(D) - \sum_{v \in A} \frac{|D_v|}{|D|} ent(D_v) \quad \text{where} \quad D_v = \{\text{samples with } A = v\}$$
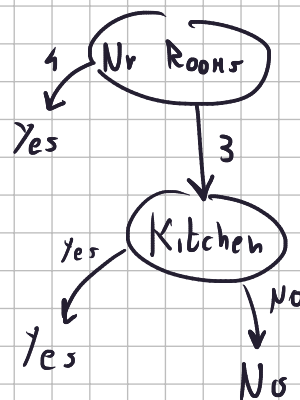
We choose the attribute that maximize the gain. For this problem:

$P_\oplus$ of D is $^3/_5$ so $ent(D) = 0.673$

Alt: Furniture. $\begin{cases} |D_o| = 3 & en(D_o) = 0.636 \\ |D_1| = 2 & en(D_1) = 0.693 \end{cases}$ ⇒ GAIN = 0.0142

Alt: Nr Rooms $\begin{cases} |D_3| = 3 & en(D_3) = 0.636 \\ |D_4| = 2 & en(D_4) = 0 \end{cases}$ ⇒ GAIN = 0.2914 ⇒

Alt: Kitchen $\begin{cases} |D_o| = 4 & en(D_o) = 0.693 \\ |D_1| = 1 & en(D_1) = 0 \end{cases}$ ⇒ GAIN = 0.1186

(Tree: Nr Rooms — 4 → Yes; 3 → Kitchen — Yes → Yes; No → No)

---

Question 2. (5 points)

1. What is a maximum likelihood (ML) hypothesis?

2. Comment the following statement: in a classification problem, the class returned by the ML hypothesis on a new instance $x$ is always the most probable class.

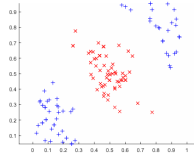Let $h^* = \arg\max_h \mathbb{P}(h|D)$, we can rewrite $h^* = \arg\max_h \mathbb{P}(D|h)\mathbb{P}(h)$. IF we assume that the hypothesis have the same probability, we get the ML hypothesis: $h_{ML} = \arg\max_h \mathbb{P}(D|h)$

The ML doesn't always return the most probable class, since the hypothesis usually have different probabilities $\mathbb{P}(h)$. The most probable is given by the Bayes Optimal Classifier: $c^* = \arg\max_c \sum_h \mathbb{P}(c|x,h)\mathbb{P}(h|D)$

1. Explain which kernel function you would choose to obtain perfect separability.
2. How SVMs classification can be applied if the data are note perfectly separable?



We notice how the red samples are the one closer to $(\frac{1}{2} \ \frac{1}{2})^T$, so i consider a kernel that compares the distances from that point.

$$K(x, x') = |\ \|x - c\| - \|x' - c\|\ | \quad \text{where} \quad c = \begin{pmatrix} \frac{1}{2} \\ \frac{1}{2} \end{pmatrix}.$$

If the dataset is not perfectly separable, i add some slack variables $\xi_i$

$$\text{argmin } \|w\|^2 \qquad \text{becomes} \qquad \text{argmin } \|w\|^2 + \gamma \sum_i \xi_i \quad \leftarrow \gamma \in \mathbb{R}$$
$$\text{s.t. } t_i \zeta(x_i) \geq 1 \ \forall i \qquad \qquad \text{s.t. } t_i \zeta(x_i) \geq 1 - \xi_i \ \forall i$$
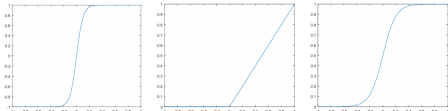$$\xi_i > 0 \ \forall i$$

**Question 4.** (5 points) Let $x \in \mathcal{X} = \{(1.8, 2.1)^T, (3.2, 1.4)^T, (2.5, 4.1)^T, (1.6, 7.7)^T, (3.1, 9.1)^T\}$ and $t \in T = \{0, 1, 1, 0, 1\}$. Consider designing an artificial neural network for learning the function $t = f(x, \theta)$.

1. Explain what is a suitable activation function for the output layer of the network.
2. For the selected activation function explain if the output units saturate and how learning the parameters of the network is affected if this is the case.

Since it's binary class. For the output we consider the sigmoid function (in this way the model is differentiable) and we predict be considering $\zeta(x) = \text{Step}(f(x, \theta))$. In this way the output unit saturate only when the correct answer is given.
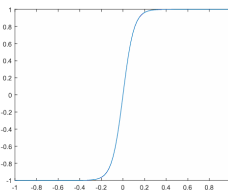
**Question 5.** (5 points)

1. Consider an image of $w_{im} \times h_{im}$ elements (pixels) and a convolution kernel of dimension $3 \times 3 \times 16$. What are the possible values of the stride and padding in order to convolve without skipping any pixels, while obtaining a feature map with the same dimensions with the input image.
2. Associate the correct name of the activation function to the plots above and provide the corresponding mathematical definition. The list of names is {ReLU, Sigmoid, Hyperbolic Tangent}.
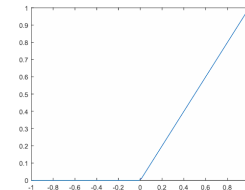


the following relations holds:

$$W_{out} = \frac{W_{im} - W_\kappa + 2P}{S} + 1 \qquad \text{in our case} \begin{cases} W_{im} = \frac{W_{im} - 3 + 2P}{S} + 1 \\ \\ h_{im} = \frac{h_{im} - 3 + 2P}{S} + 1 \end{cases} \begin{cases} P = \frac{1}{2}\left(W_{im}(S-1) - S + 3\right) \\ \\ S = \frac{h_{im} - 3 + 2P}{h_{im} - 1} \end{cases}$$
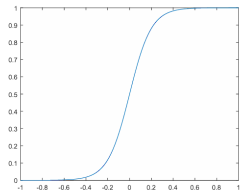
$$h_{out} = \frac{h_{im} - h_\kappa + 2P}{S} + 1$$

$$\begin{cases} P = \frac{1}{2}\left(W_{im}(S-1) - S + 3\right) \\ \\ S = \frac{1}{h_{im} - 1}\left(h_{im} - 3 + W_{im}(S-1) - S + 3\right) \end{cases} \Rightarrow \begin{cases} P = 1 \\ \\ S = 1 \end{cases}$$
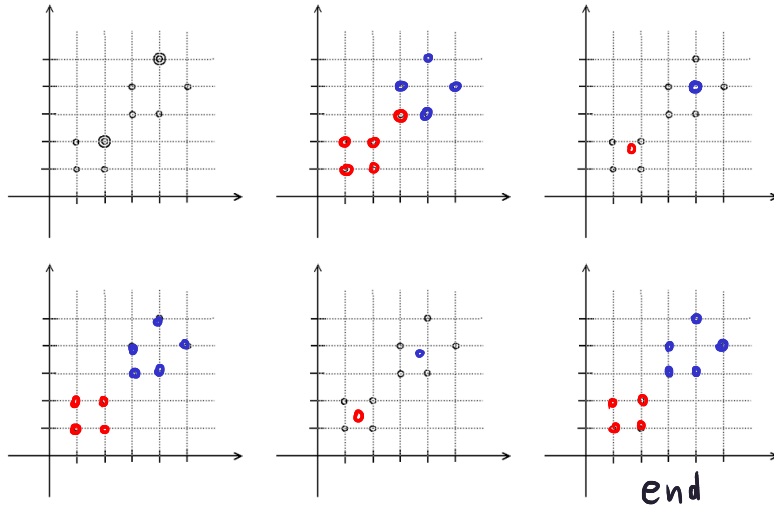


tanh
$2\sigma(2x) - 1$

ReLU
$\max(0, x)$

Sigmoid
$\sigma(x) = \frac{1}{1 + \exp(-x)}$

**Question 6.** (5 points)

Simulate the execution of K-means in this 2-D data set with k=2 and initial centroids indicated by double circles: use one diagram for each step of the algorithm. Describe explicitly how each step is obtained and what is the termination condition of the algorithm. Drawing only the steps is not sufficient.



at each steps we divide the samples in two groups (group 1: colser to centroid 1, group 2: the others) and then recompute the centroids by considering the center of mass of the groups.

we stop when the new centroids computed are equals to the one of the previous step.