

Decisions trees

1.2) IO3

Find the feature with the most Information gains

Feature	$\frac{ V_{a=\tau} }{ V }$	$\frac{ V_{a=\epsilon} }{ V }$	$H(V_{a=\tau})$	$H(V_{a=\epsilon})$	$IG(V,a) - H(V)$
A	$\frac{3}{4}$	$\frac{1}{4}$	$H(\frac{2}{3})$	$H(\frac{0}{1})$	$-\frac{3}{4}H(\frac{2}{3}) - \frac{1}{4}H(\frac{0}{1})$
B	$\frac{2}{4}$	$\frac{2}{4}$	$H(\frac{1}{2})$	$H(\frac{1}{2})$	$-\frac{2}{4}H(\frac{1}{2}) - \frac{2}{4}H(\frac{1}{2})$
C	$\frac{2}{4}$	$\frac{2}{4}$	$H(\frac{2}{2})$	$H(\frac{2}{2})$	$-\frac{2}{4}H(\frac{1}{2}) - \frac{2}{4}H(\frac{1}{2})$

$$IG(V,a) - H(V) = -\frac{|V_{a=\tau}|}{|V|} \cdot H(V_{a=\tau}) - \frac{|V_{a=\epsilon}|}{|V|} \cdot H(V_{a=\epsilon})$$

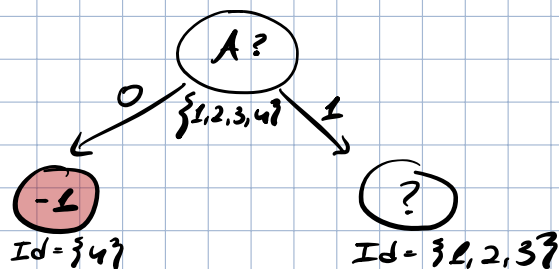
$$\bullet A \quad -\frac{3}{4}H(\frac{2}{3}) - \frac{1}{4}H(\frac{0}{1}) = -\frac{3}{4} \cdot 0.918 - 0 = -0.688$$

$$H(\frac{2}{3}) = -\frac{2}{3} \log_2(\frac{2}{3}) - \frac{1}{3} \log_2(\frac{1}{3}) = 0.918$$

$$H(0) = -\log_2(1) = 0$$

$$\bullet B, C \quad -H(\frac{1}{2}) = -1$$

(A) (B,C)
 $-1 > -0.688$ ולכן נבחר במידת A לחלוקה הראשונה

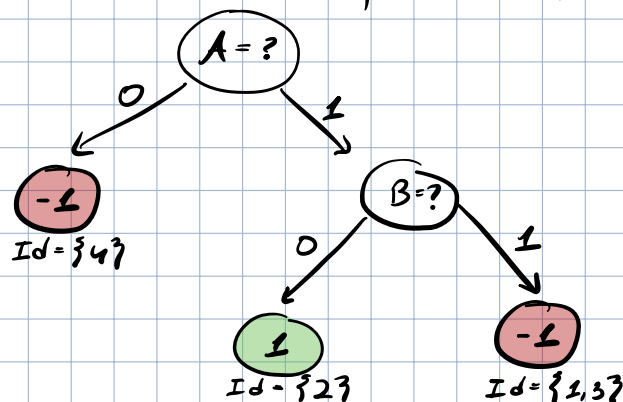


בהינתן A=0 ישנו Sample בודד ולכן לא נחלק אותו יותר.

נבחר במידת לחלוקה השנייה $8A = 1$

Feature	$\frac{ V_{a=\tau} }{ V }$	$\frac{ V_{a=\epsilon} }{ V }$	$H(V_{a=\tau})$	$H(V_{a=\epsilon})$	$IG(V,a) - H(V)$
B	$\frac{2}{3}$	$\frac{1}{3}$	$H(\frac{1}{2})$	$H(0)$	$-\frac{2}{3}H(\frac{1}{2}) - \frac{1}{3}H(0)$
C	$\frac{1}{3}$	$\frac{2}{3}$	$H(1)$	$H(\frac{1}{2})$	$-\frac{1}{3}H(1) - \frac{2}{3}H(\frac{1}{2})$

הפיצול לפי הפיצורים ב-1 ו-2 סקום עבור היחס הנ"ל. מאחר ואנחנו רוצים
 המקסימום (הוא 2, אין חצי חלוקות והמשקל ימנען אין חלוקות) לבחירת הפיצור
 [ע"פ סך המינימום שניתן לט, כמובן]. לבחור את הנושא ב-1 ו-2, ונראה שזה הנושא הנכון



סיכום היחס נקבע ע"פ
 התחית היחס בחקירה
 של "tie".

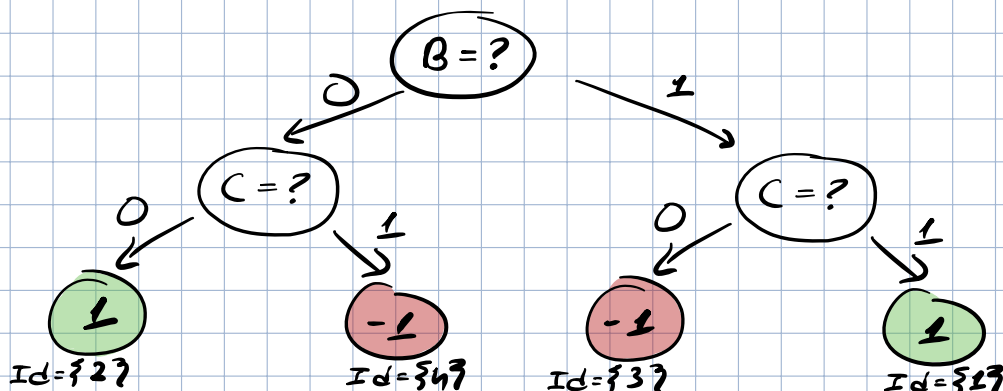
Training Error:

$$\frac{\# \text{incorrect_predictions}}{\# \text{all_predictions}} = \frac{1}{4}$$

Misclassified tuples: (1, 1)

sample's id
sample's label

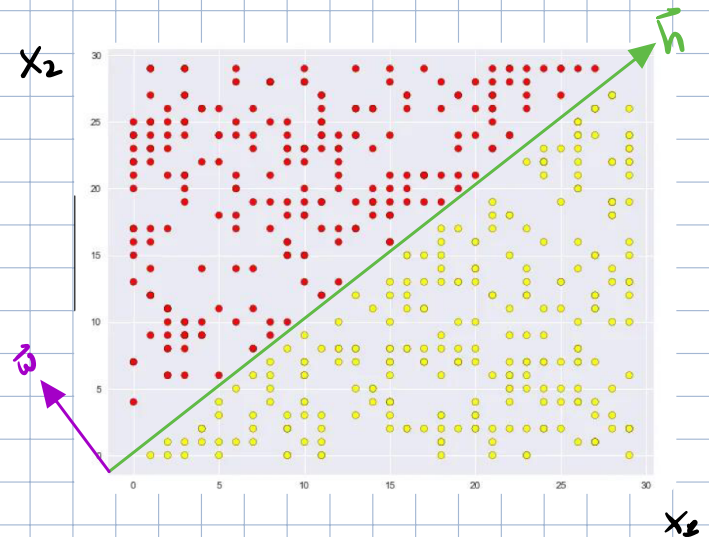
1.2



Decision trees strike again

2.1) $h(x_1, x_2) = \text{Sigh}(w^T x + b)$

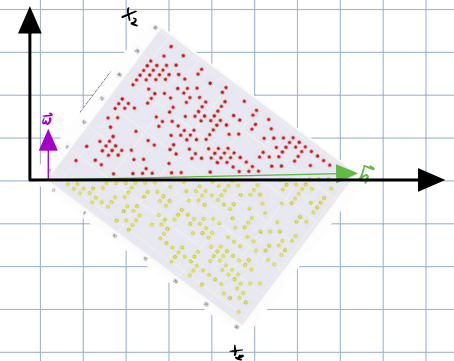
Yes. There are linear classifiers that achieve Perfect Accuracy on the dataset. For example we can look at the following \vec{h} defined by \vec{w} .



2.2) No. Decision trees with depth=1 only split the dataset once with a single feature. Looking at the given dataset, we can see that there is no $\Theta_1 \in x_1$ or $\Theta_2 \in x_2$ s.t. the dataset is separable solely by this value.

2.3) No, for the same logic mentioned above. To simplify the explanation, we could look at $\Theta_1 \in x_1$ as a vertical line in the graph, and $\Theta_2 \in x_2$ as a horizontal line. We can see that there are no 2 values $\in \Theta_1 \cup \Theta_2$ for which the data is separable. (=No 2 vertical/horizontal lines)

2.4) Yes. Rotating the data 45° clockwise will make it separable by a horizontal line. Illustration:



Information Gain

3.1) The Problem with the chatGPT Proof is an Incorrect mathematical transition from the assumptions $\{ \text{Entropy}(D) \geq 0, \text{Entropy}(D_i) \geq 0 \}$ to the conclusion that the subtraction $\text{Entropy}(D) - \sum_i \frac{|D_i|}{|D|} \cdot \text{Entropy}(D_i)$ is also non negative.

3.2) Prove that $IG(V, \alpha) = H(V) - \frac{|V_{\alpha=T}|}{|V|} H(V_{\alpha=T}) - \frac{|V_{\alpha=F}|}{|V|} H(V_{\alpha=F}) \geq 0$

Given that $H(V) = H(P_V) = -P_V \log_2(P_V) - (1-P_V) \log_2(1-P_V)$ holds
 $\forall \beta_1, \beta_2, \alpha \in [0, 1] : -\alpha H(\beta_1) - (1-\alpha) H(\beta_2) \geq -H(\alpha\beta_1 + (1-\alpha)\beta_2)$

$IG(V, \alpha) =$ / definition of IG

$$H(V) - \frac{\overset{\alpha}{|V_{\alpha=T}|}}{|V|} H(\overset{\beta_1}{V_{\alpha=T}}) - \frac{\overset{(1-\alpha)}{|V_{\alpha=F}|}}{|V|} H(\overset{\beta_2}{V_{\alpha=F}}) \geq \text{ / Given Property}$$

$$H(V) - H\left(\frac{|V_{\alpha=T}|}{|V|} \cdot P_{V_{\alpha=T}} + \frac{|V_{\alpha=F}|}{|V|} \cdot P_{V_{\alpha=F}}\right) = \text{ / definition of } P_{V_{\alpha=T}} \text{ and } P_{V_{\alpha=F}}$$

$$H(V) - H\left(\frac{\cancel{|V_{\alpha=T}|}}{|V|} \cdot \frac{|\{(x,y) \in V_{\alpha=T} | y=1\}|}{\cancel{|V_{\alpha=T}|}} + \frac{\cancel{|V_{\alpha=F}|}}{|V|} \cdot \frac{|\{(x,y) \in V_{\alpha=F} | y=1\}|}{\cancel{|V_{\alpha=F}|}}\right) =$$

$$H(V) - H\left(\frac{|\{(x,y) \in V_{\alpha=T} | y=1\}| + |\{(x,y) \in V_{\alpha=F} | y=1\}|}{|V|}\right) = \text{ / } V_{\alpha=T} \cup V_{\alpha=F} = V, V_{\alpha=T} \cap V_{\alpha=F} = \emptyset$$

$$H(V) - H\left(\frac{|\{(x,y) \in V | y=1\}|}{|V|}\right) = \text{ / definition of } P_V$$

$$H(V) - H(P_V) = 0$$

Linear Classification

- For homogeneous: $h(x) = +1 \Leftrightarrow w^T x = \underbrace{\|w\| \|x\|}_{>0} \cos \angle(w, x) > 0 \Leftrightarrow \cos \angle(w, x) > 0$
invariant to the scale of w .
- The geometric (signed) margin of $x \in \mathbb{R}^d$ is $\frac{w^T x}{\|w\|}$
- For non-homogeneous: $h(x) = +1 \Leftrightarrow w^T x \geq -b$
 b being almost equivalent to the margin.

4.1) In the homogeneous case we had a degree of freedom in scaling w . For the non-homogeneous case we still have a degree of freedom as long as we scale b along with w . We can consider the classifier to be a $w' \in \mathbb{R}^{d+1}$ vector s.t. $w' = \begin{pmatrix} w \\ b \end{pmatrix}$ and reframe the classification as such: $h = +1 \Leftrightarrow w'^T \begin{pmatrix} x \\ 1 \end{pmatrix} > 0$.

In that case we get the same classifier, and we have an homogeneous equation which is invariant to the scale of w' .

4.2) A is b_1 , B is b_4 , C is b_2 , D is b_3 .

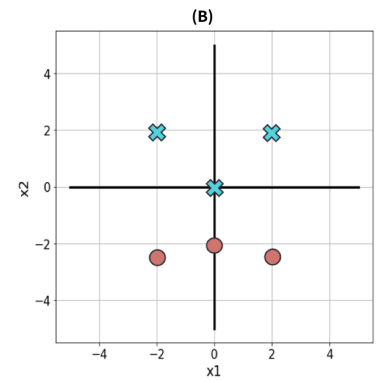
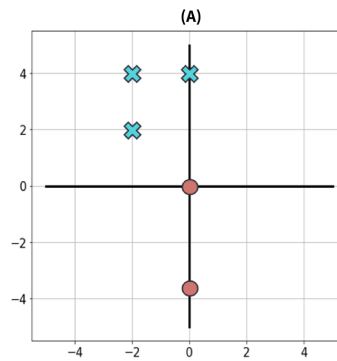
K-Nearest Neighbors

• A training Point is not considered a neighbor of itself.

i) $k=1$, $d(u,v) = \|u-v\|_2$

ii) $k=3$, $d(u,v) = \|u-v\|_2$

iii) $k=1$, $d(u,v) = \|u-v\|_1$



dataset \ model	A	B
i ($k=1$, $d(u,v)$)	NO. The red Point on the Axes Intersection is misclassified as blue.	NO. The blue Point on the Axes Intersection is misclassified as the red below it and vice versa.
ii ($k=3$, $d(u,v)$)	NO. red Points will be blue as there could only be 1 red neighbor on 2 blue ones.	Yes
iii ($k=1$, $\ u-v\ _1$)	Yes	NO. blue Point on intersection will be classified by the red Point below it and vice versa.

5.2) • Answers unchanged for models (i) & (ii).

In these models the distance is calculated by the L_2 norm which is an euclidian distance. This norm is invariant to rotations. [I'm not sure if calculations are considered 'brief' enough, so added below, ignore if not needed). If the distance calculation is unchanged, so is the classifier.

• For model (iii) answers might change.

This model uses the L_1 norm as a distance metrics, and it is affected by rotation as the distance is affected by the points relation to the axes.

- For dataset A8 rotation could decrease the L2 distance between the (0,0), (2,2) points s.t they'll be nearest neighbors.
- For dataset B8 a (45°) rotation could lead to a tie in the distances from the (0,0) point. ($2\sqrt{2}$)

Some calculations for 5.2)

$$X = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \quad X' = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} = \begin{pmatrix} X_1 \cos \theta - X_2 \sin \theta \\ X_1 \sin \theta + X_2 \cos \theta \end{pmatrix}$$

$$\|X\|_2 = \sqrt{X_1^2 + X_2^2}$$

$$\|X'\|_2 = \sqrt{(X_1 \cos \theta - X_2 \sin \theta)^2 + (X_1 \sin \theta + X_2 \cos \theta)^2}$$

$$= \sqrt{X_1^2 \cos^2 \theta + X_2^2 \sin^2 \theta - 2X_1 X_2 \cos \theta \sin \theta + X_1^2 \sin^2 \theta + X_2^2 \cos^2 \theta + 2X_1 X_2 \sin \theta \cos \theta}$$

$$= \sqrt{X_1^2 (\cos^2 \theta + \sin^2 \theta) + X_2^2 (\sin^2 \theta + \cos^2 \theta)} = \sqrt{X_1^2 + X_2^2} = \|X\|_2$$