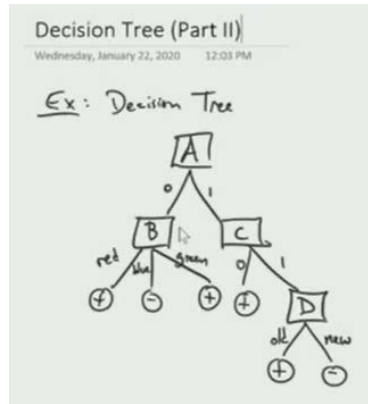


بسم الله الرحمن الرحيم

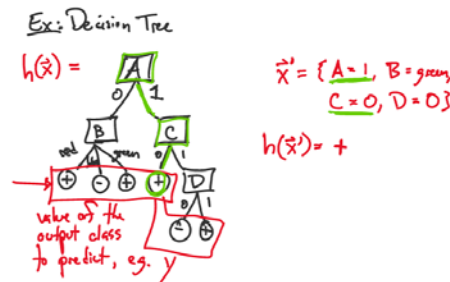
المحاضره ديه هنتكلم علي أول الجورنم هيبقا مفيد في ال decision trees .. المهم ان HW2 نزل .. فابدا فيه علي طول .. تعال نفكر في ال decision tree كإنها hypothesis .. تعال نرسم مثال .. ال Decision tree اللي هنتغل عليها بيتكون من شوية Attributes : بص ع الرسمه



المهم احنا هنتغل علي الشجره اللي فوق ديه .. تعال نشوف ايه اللي هيحصل لو دخل مثال جديد علي الشجره ديه .. هنسميها إكس برايم .. هيبقا ليه شوية قيم لل Attributes اللي عندنا ...

$$\vec{x}' = \{A=1, B=green, C=0, D=0\}$$

المهم بيقولك احنا في الشجره بتاعتنا ... عند ال A .. احنا عندنا  $h(x)$  وديه فانكشن بت map ال x values ل y values .... فاحنا هنبدأ في المثال الجديد اللي عندنا من أول A .. هنلاقي قيمتها 1 .. وبعدها نخش علي C .. هل القيمه صفر ولا واحد .. هتلاقي انها صفر .. فهترجع "+" ... فهتلاقي ال return



### Algorithm 3: Decision Tree

def  $h(\vec{x})$ : (prediction step)

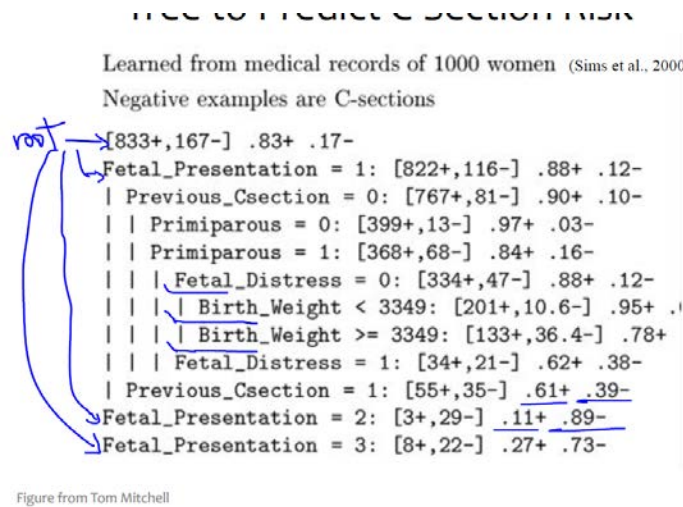
3 cases

① internal node: test an attribute  $X_m$

② branch for node: select branch correspondingly to the value of the attribute

③ leaf node: predict stored value of y  
(variant: return  $p(y|\vec{x})$ )

تعال نفهم هو ليه مستخدمين ال variant الأخير في الحالة رقم 3 .. بص ع المثال ده .. احنا بنبيص علي حالات من النسوان الحامل .. بنبيص علي ال medical record .. و لكل مريض عندنا different attributes .. فاحنا عندنا ال attributes ديه واحنا بنحاول ن predict ال C-Section Risk



9

المهم تعال نعمل تريننج ألجورزم .. وده هنا هنعرف recursive algorithm .. انت عندك داتا سيت D .. احنا هنا بس بنعرف generic version of the algorithm ... بس في حالات خاصه من ال algorithm ده زي مثلاً ID3, CAR, ... بس في الآخر الفكره واحده و الاختلاف مش بيبقا كبير يعني ..

أول خطوه هي إننا هنبدأ ال root node و ديه هتأخذ data = D .. هتجيبها ال Full training set D .. وبعدين تبدأ ال recursion implementation .. return(train\_tree(root)) ... الفانكشن اللي اسمها train\_tree ديه اللي هتعمل ال recursion علي ال root ديه .. تعال نشوف الفانكشن ديه بتشتغل ازاى ...

أول خطوه إنك تحسب attribute  $X_m$  .. و ده اللي هو أحسن attribute نقدر ن split عليه ال nodes data .. فال node ديه أوبجكت ... جواها شوية data .. ثاني خطوه هو إنك أول ما انت خلاص بقي عندك ال  $X_m$  .. خليه هو ال decision attribute لل node الحاليه .. وبعدين انت عاوز تبني شجره .... فهتقول .. لكل قيمه من قيم ال attribute  $X_m$  .. هحتاج نعمل branch ليه اسم باسم القيمه اللي داخلين من ناحيتها .. رابع خطوه هو انك تقسم ال node's data لل branches اللي داخلين عليها اللي هو ال children بتوع ال current node .. فهتلاقي نفسك بتعرف داتا سيت اسمها D .. بس المره ديه بتتكون من X and Y pairs بس دول اللي موجودين جوال node الحاليه .. بحيث إن قيمه ال  $X_m$  هي قيمة ال value ... كذا نخش علي خطوه ال recursion .. recurse on each branch ... لكل قيمة من قيم ال values اللي موجوده في ال branch الحالي .. هنعمل new node وديه هتبقا مخصصه ل value الحاليه يعني .. وديه هتأخذ القيمه بتاعت ال recursion بتاع ال train\_tree ثاني .. بس المره ديه هتبعثها ال new\_node .. والداتا بتاعت ال node اللي هتبعثها ثاني .. هي ال sub dataset اللي كوتأها فوق .... فاضل مشكله واحده .. احنا هنقف امتي؟؟ .. فاحنا محتاجين base case .. اللي هي هتقولك ده ال situation اللي عندو هتخرج من ال recursion .. في الحالة ديه .. هتقول لو النود داتا خلاص بقت perfectly classified .. اخرج .. بس وانت بتخرج ابعثلي ال label بتاع ال majority vote علي ال current node's data .. بس هنا في مشكله إيه معني ان ال node data be perfect classified? .. perfect classification معناها لما يكون ال error rate بصفر . اللي هو يا كلو + يا كلو - .. ده اللي هو ال majority vote perfectly classified

(return: return  $p(y|x)$ )

def train(D): (generic learning algorithm, lots of specific cases ID3, CART)

root = new Node(data=D)

train\_tree(root)

def train\_tree(node)

one that minimizes some splitting criterion

①  $X_m$  = best attribute on which to split nodes data

② Let  $X_m$  be decision attribute for node

③ For each value  $v$  of attribute  $X_m$ :

create a branch labeled with  $v$

④ Partition the nodes data into descendants

$D_{X_m=v} = \{(\vec{x}, y) \in D \mid x_m = v\}, \forall v$

⑤ Recurse on each branch:

For each value  $v$  and branch  $b$ , add a new node

node<sub>v</sub> = train\_tree(new Node( $D_{X_m=v}$ ))

⑥ If node's data is perfectly classified, stop and return node with label = majority vote (nodes data)

error rate on data is zero using the majority label.

Ex: Decision Tree Learning

المهم في الصورة اللي فوق في حاجه مهمه .. وهي إنك ازاي تختار ال best attributes .. اختار اللي ت maximize 1- error rate .. او ت maximize Gini gain .. او انك ت maximize mutual information

تعال ناخذ مثال: ..

## Decision Tree Learning Example

### Dataset:

Output Y, Attributes A, B, C

Y	A	B	C
-	1	0	0
-	1	0	1
-	1	0	0
+	0	0	1
+	1	1	0
+	1	1	1
+	1	1	0
+	1	1	1

### In-Class Exercise

Using **error rate** as the splitting criterion, what decision tree would be learned?

11

احنا عاوزين نحدد ال root node .. وده هيقا new node object .. وهنا هيقا عندنا ال DataSet كلها .. السؤال هنا .. فيها كام واحد + وكام واحد - .... هتلاقي عندك 5 بوزيتيفو 3 نيجاتيف .. انت عندك 3 attributes ABC .. و دول كلهم binary .. لكل split عندنا .. هنقول مثلاً لو A كانت بصفر .. انت عندك واحد + و معنديش - ... لو ال A كانت بواحد .. يبقا انت عندك 3 -ve و 5 +ve .. نفس الكلام بالنسبة ل B

Ex: Decision Tree Learning w/ Error Rate

root = new Node(D with [5+, 3-])

Y A B C

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

h<sub>A</sub> [A] splits into [4+, 0-] and [1-, 3-]

error(h<sub>A</sub>, D) = 3/8

h<sub>B</sub> [B] splits into [1+, 3-] and [4+, 0-]

error(h<sub>B</sub>, D) = 1/8

h<sub>C</sub> [C] splits into [2+, 2-] and [3+, 1-]

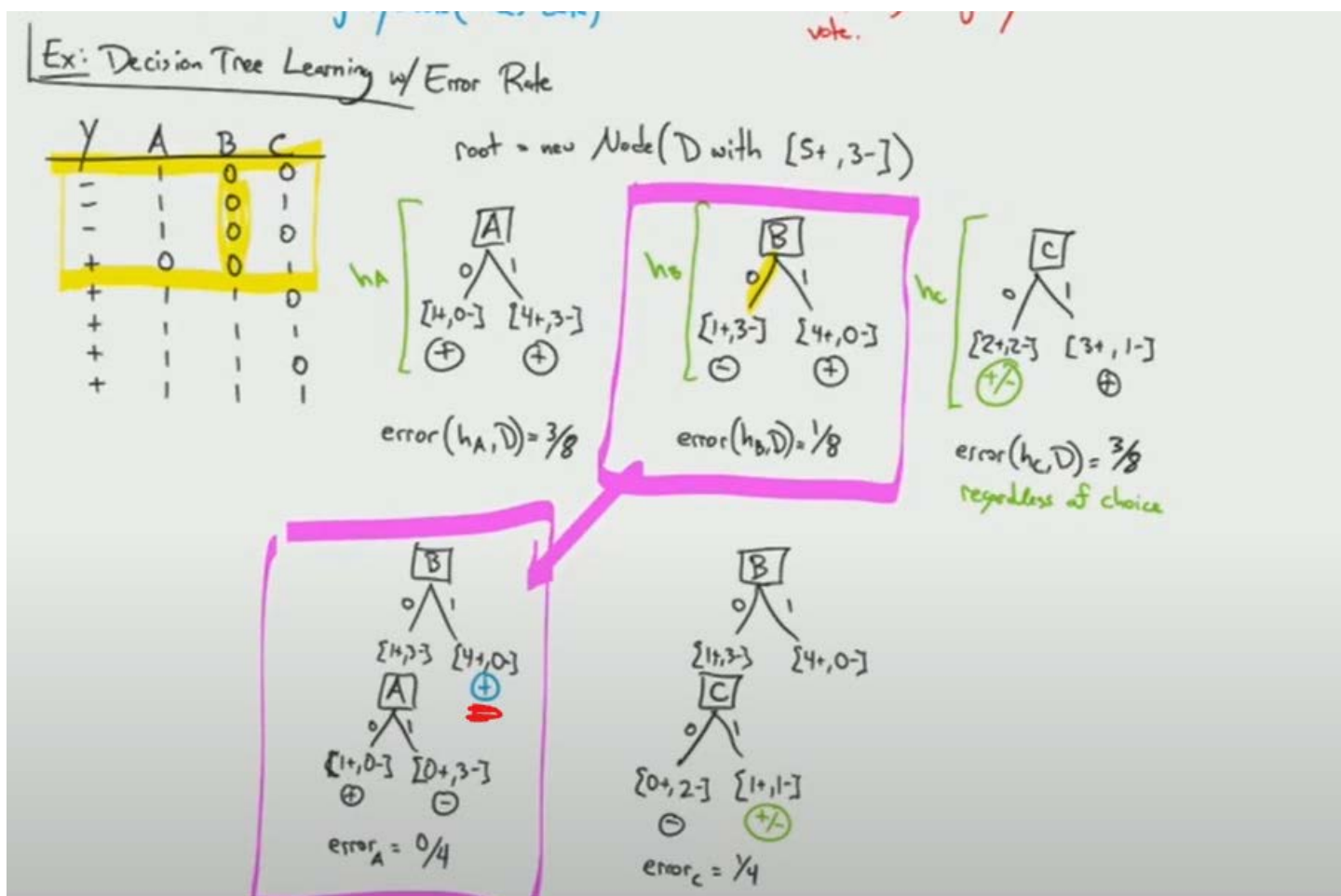
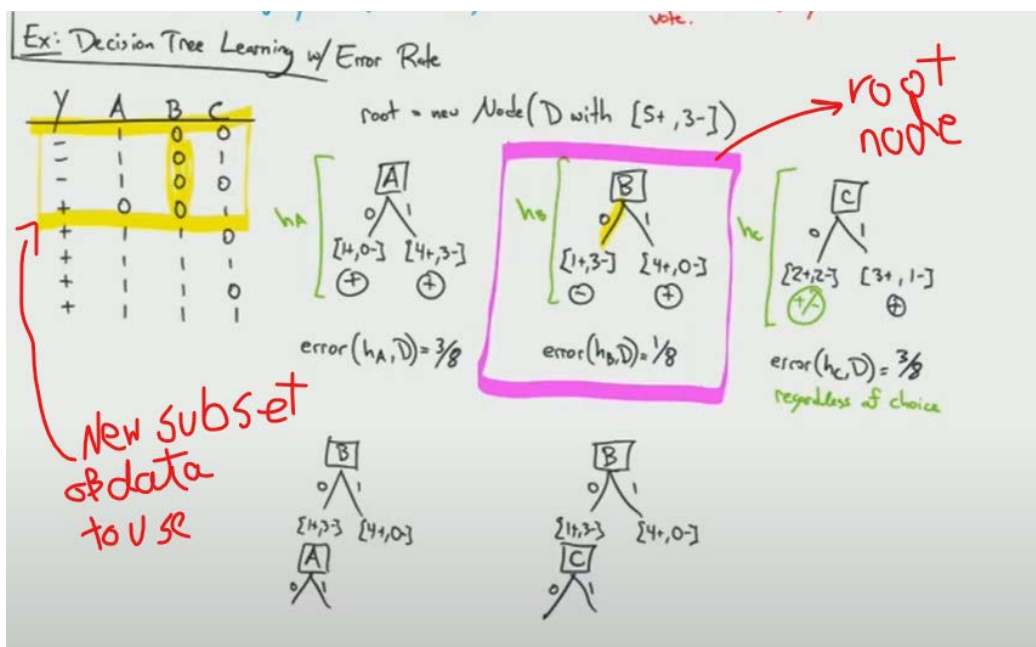
error(h<sub>C</sub>, D) = 3/8

regardless of choice

1st split  
+ve or -ve  
هتتفرق  
هتتفرق

المهم انك هتختار أقل error .. احنا كدا هن split علي B ... ويلا بينا علي ال recursion .. احنا لما بن recurse .. بن recurse علي multiple options .. فا يابا هنفصل علي ال A أو هنفصل علي ال C .. المهم انك تاخذ بالك من حاجه مهمه ... انت لما بتبدأ ال recursion انت مش بتشتغل علي نفس الداتا سبت .. انت بتشتغل علي subset من الداتا سبت اللي كانت في الأول ..

مثلاً .. لما قلنا إن لو ال B كانت بصفر .. فانت دخلت في البرانش بتاع الصفر ... فدلوقت الداتا سبت بتاعتك اللي هتشتغل عليها هي .. اللي متلوونه باللون الاصفر .



## Decision Tree Learning

- Definition: a **splitting criterion** is a function that measures the effectiveness of splitting on a particular attribute
- Our decision tree learner **selects the “best” attribute** as the one that maximizes the splitting criterion
- Lots of options for a splitting criterion:
  - error rate (or *accuracy* if we want to pick the tree that *maximizes* the criterion)
  - Gini gain
  - Mutual information
  - random
  - ...

13

تعال نشوف مثال:

## Decision Tree Learning Example

### Dataset:

Output Y, Attributes A and B

Y	A	B
-	1	0
-	1	0
+	1	0
+	1	0
+	1	1
+	1	1
+	1	1
+	1	1

### In-Class Exercise

Which attribute would **error rate** select for the next split?

1. A
2. B
3. A or B (tie)
4. Neither

14

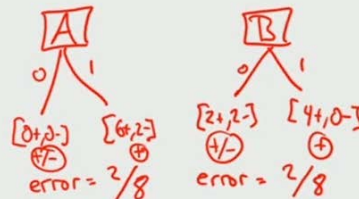
## Decision Tree Learning Example

### Dataset:

Output Y, Attributes A and B

Y	A	B
-	1	0
-	1	0
+	1	0
+	1	0
+	1	1
+	1	1
+	1	1
+	1	1

root : [6+, 2-]



TIE

از اي نعالج التعادل اللي عندنا ده؟ .. هنعرف حاجه اسمها Gini Impurity ... Given a random variable Y over K classes {1,2,...,k}

هنقول ان التعريف بتاع Gini Impuity هو عبارة عن  $G(y)$  .. مجموع حاصل ضرب احتمالية إن ال  $Y=k$  و ال  $y!=k$  ... من اول ال  $K \rightarrow k=1$

تعال نفكر في الآتي .. consider the case where  $y$  is the outcome of a weighted dice roll

$P(Y=k)$  is the prob of landing on side  $k$

$P(Y!=k)$  is the prob of landing on any other side besides  $k$

اللي عاوزين نعملو اننا conceptualize a gain

في اللعبة هنا هدفنا اننا predict the next dice roll .. و مديك ال weights of the die ... فانت عارف من الاول احتمالية كل side بكام .. ايه الحل بتاعك في الحالة ديه ... ازاى هتلعب اللعبة ديه ..

أول حل ... توقع ال most probable side every time .. في الحالة ديه ... expected error rate .. تقدر تقول إنها  $1-P(y=\text{most probable side})$  (اللي هي  $P(Y!=Y^*)$  .. ده عبارة عن  $y^* = \arg \max_k P(Y=k)$  .. Summation over all the values  $k \neq y^*$ )

.. ثاني حل .. هي إنك ... roll another die w/some weight on sides and predict whatever it lands on expected error rate الحقيقه مفهمنش الحته ديه .. ابقا اسمع آخر 10 دقائق ثاني ..

Def: Gini Impurity. Given a random variable  $Y$  over  $K$  classes  $\{1, 2, \dots, K\}$

$$G(Y) = \sum_{k=1}^K P(Y=k)P(Y \neq k)$$

probability of landing on  $k$       probability of prediction die landing on not  $k$

$$= \sum_{k=1}^K P(Y=k)(1 - P(Y=k))$$

$$= 1 - \sum_{k=1}^K [P(Y=k)]^2$$

prob. of both die landing on same side  
= expected misclassification error for Case 2.

Consider  $Y$  as a random dice roll.  
 $P(Y=k)$  is the probability of landing on side  $k$ .  $P(Y \neq k)$  is the prob. of any other side.

Suppose you are given side weights.  
Case 1: expected error rate if you classify which side lands next with majority vote. = 3

Case 2: Roll an equally weighted die to decide how to predict.  
 $G(Y)$

Given a dataset  $D$ , then let  $P(Y=k) = \frac{N_{y=k}}{N}$  (where  $N_{y=k}$  is # of examples w/  $Y=k$  and  $N$  is total # of examples)

Def: Gini Gain:

$$G(Y, X_M; D) = G(Y; D) - [P(X_M=0)G(Y; D_{X_M=0}) + P(X_M=1)G(Y; D_{X_M=1})]$$

impurity before split      impurity on left      impurity on right      weighted avg. impurity if we split on  $X_M$

$$G(Y; D) = 1 - (7/10)^2 - (3/10)^2$$

$$P(X_M=0) = 4/10 \quad P(X_M=1) = 6/10$$

$$G(Y | D_{X_M=0}) = 1 - (4/4)^2 - (0/4)^2$$

$$G(Y | D_{X_M=1}) = 1 - (3/6)^2 - (3/6)^2$$