

بسم الله الرحمن الرحيم .. نبدأ محاضره جديده

هنتكلم شويه علي high level topics .. هنتكلم علي ال overfitting and K-Nearest neighbors .. الدكتور اتكلم شويه علي ال staff .. ناس جميله والله ربنا يوفقهم و يوفقنا جميعاً .. الدكتور قال ان في حد سأل سؤال في الكلاس ..

When and how do we decide to stop growing trees? What if the set of values an attribute could take was really large or even infinite?

انهارده هنتكلم انهارده يعني علي حته ال discrete attributes .. بس لو كان في attribute و كان real value .. كمية الارقام اللي ممكن تاخذها هتبقا كثيره جداً .. بس في حركة كذا بت consider ال order بتاع ال split اللي بيحصل في ال tree بيقا  $O(L)$  ... و ال L هنا هو عدد ال values الي ال attribute بياخذها في ال training set .. فبيقالك فكر ايه هي ال trick اللي ممكن تبقا موجوده؟

....

آخر مره احنا اتكلمنا عن ال splitting criteria for decision trees .. وقلنا ان ال splitting criteria هي ال how we measure the effectiveness of a particular split .. في الصورة اللي جايه ديه :: احنا قلنا في المحاضره ان كان في تعادل ما بين ال split on A or B ...

## 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

الحل اهو::

## 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



Misclass. Rate  
 $r(A) = 2/8$   
 $r(B) = 2/8$   
*r(.) treats attributes as equally good*

16

وبعدين اتكلمنا علي ال Gini impurity و إنها different splitting criterion .. اللي هي Gini gain .. وبعدين الدكتور كان قال اننا نعملها علي الداتا سبت بتاعتنا .. بس انا ماعملتش حاجه .. ففعال نشوف الدكتور عمل ايه ..

أول حاجه عند ال Root node ..  $G(Y) = 1 - P(\text{Getting it right as a plus}) - P(\text{Getting it right as a minus})$  فديه هي قيمة ال impurity

فالسؤال هنا .. how much can we reduce that impurity that exists in the root by splitting on either attribute A or B??

تعال نشوف: لو احنا فصلنا علي A .. هنقول إن  $P(A=0) = 0$  وهو إن صفر من ال Training examples عندها الحالة ديه ... و  $P(A=1) = 1$  اللي هو كلو هيخش في الناحيه ديه .. طيب لو جينا نقول بقي:

$$G(Y|A=1) = G(Y \text{ and } A) / A = G(Y)$$

كدا بقت نفس الرووت نوود .. انما لو قلت الكوندشن عند  $A=0$  .. انت معندكش داتا اصلا لل  $A=0$  .. فالدكتور سأل سؤال .. ازاي تقدر ن

Estimate a quantity about the random variable Y if we have no data about Y?????

فبتسمي الحالة ديه undefined ...

انت لما بتيجي تحسب ال Gini gain بتعمل ايه .. بتقول آدي عندك القيمه اللي هي 0.375 عند الرووت .. و تطرح منها ال weighted average of the gini impurity at the 2 leaf nodes

ده اللي هو undefined \* 0 و ال  $0.375 * 1$  و ده طلعك صفر Gini gain .. انت كدا مقلتش ال impurity نهائي ..

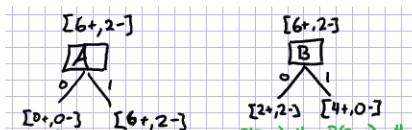
تعال نشوف حالة ال B ...  $P(B=1) = P(B=0)$  .. ولما جيت تتحرك زي A .. لاقيت ان عندك  $\text{reduction in impurity} = 0.125$  لما جيت تفصل عند B .. فعند ال root ال impurity كانت 0.375 .. بس ال weighted average ما بين ال 2 leafs .. بقي 0.25 .. فال Gini gain هنا عمل ايه .. هو خد some probabilistic notion لل remaining decisions اللي انت عاوز تعملها .. فانت بتقول ايه .. بتقول انك عاوز attribute بتقلل ال impurity on both sides ... وفي فرق ما بين انك تفصل جزئياً ال impurity و انك تفصلها كلياً ... و الدكتور بيقول اننا بن capture الحته ديه باستخدام ال expected miss classification rate ... هو قال كمان ان ال error rate احنا بس بنبص علي pure error rate .. انما هنا بنتكلم علي probabilistic notion of error rate ... فهنا في المثال ده .. في attribute فاز اللي هو B ...

## 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



- $G(Y) = 1 - (6/8)^2 - (2/8)^2 = 0.375$
- $P(A=1) = 8/8 = 1$
- $P(A=0) = 0/8 = 0$
- $G(Y|A=1) = G(Y)$
- $G(Y|A=0) = \text{undef}$
- $\text{GiniGain}(Y|A) = 0.375 - 0(\text{undef}) - 1(0.375) = 0$
- $P(B=1) = 4/8 = 0.5$
- $P(B=0) = 4/8 = 0.5$
- $G(Y|B=1) = 1 - (4/4)^2 - (0/4)^2 = 0$
- $G(Y|B=0) = 1 - (2/4)^2 - (2/4)^2 = 0.5$
- $\text{GiniGain}(Y|B) = 0.375 - 0.5(0) - 0.5(0.5) = 0.125$

20

الدكتور دلوقت هيتكلم علي mutual information ... فكر في 2 راندم فاريلز إكس و واي .. المهم لما جيه يحسب ال entropy .. حسبته كانه expected log probability of a random variable .. الدكتور بيقول إن ال entropy هو شبه اللي كنا بنعملو في ال Gini impurity .. هو بيقيس the sort of stochasticity in a collection of data .. فاللي هتحاول نعملو اننا نقلل ال Entropy عن طريق اننا نختار attribute معين .. أول حاجه ممكن نعملها هو اننا نحدد a specific conditional entropy .. وديه شبه ال Gini impurity on one of the leaf nodes .. الدكتور بيقول اننا هنستخدم ال Specific conditional entropy كانه entropy بعد ما انت خلاص عملت ال condition علي ال attribute انو ياخذ قيمه معينه .....

ال Conditional entropy ده بقا هو بياخد weighted average of specific conditional entropy .. بمعني تاني ديه احتمالية إن ال attribute ياخذ قيمه x مضروبه في ال x for that specific conditional entropy .. اللي هو فكر فيها كإن لكل leaf node احنا عملنا weight علي حسب عندنا كام مثال مشي في ال branch ده .. وبعدين خد ال specific conditional entropy .... بعد كذا ال mutual information .. بيتحسب بال entropy at root node تطرح منو ال weighted average of entropies after conditioning on X ... والدكتور قال ان ده اللي كنا بنعملو في ال Gini Gain .. احنا هناخد ال Impurity at root – weighted average of the impurities .. فلو بنتكلم علي ال decision tree احنا نقدر نستخدم ال definition ده بتاع ال mutual information ك splitting criterion

## Mutual Information

Let  $X$  be a random variable with  $X \in \mathcal{X}$ .  
 Let  $Y$  be a random variable with  $Y \in \mathcal{Y}$ .

Entropy:  $H(Y) = - \sum_{y \in \mathcal{Y}} P(Y=y) \log_2 P(Y=y)$

Specific Conditional Entropy:  $H(Y | X=x) = - \sum_{y \in \mathcal{Y}} P(Y=y | X=x) \log_2 P(Y=y | X=x)$

Conditional Entropy:  $H(Y | X) = \sum_{x \in \mathcal{X}} P(X=x) H(Y | X=x)$

Mutual Information:  $I(Y; X) = H(Y) - H(Y|X) = H(X) - H(X|Y)$

- For a decision tree, we can use **mutual information** of the output class  $Y$  and some attribute  $X$  on which to split as a **splitting criterion**
- Given a dataset  $D$  of training examples, we can estimate the required probabilities as...

$$P(Y=y) = N_{Y=y}/N$$

$$P(X=x) = N_{X=x}/N$$

$$P(Y=y|X=x) = N_{Y=y, X=x}/N_{X=x}$$

where  $N_{Y=y}$  is the number of examples for which  $Y=y$  and so on.

الانتروبي هو high level بنقيس بيه ال expected #of bits عشان نقدر ن code 1 random draw from  $X$  فكر فيها كإنها قياس ل level of uncertainty encoded in a random variable .. فمثلاً لو كان random variable كان بيتكلم علي coin flip و كان دائماً بيجي ال head .. فهنا مفيش random uncertainty يعني لو جيت ت flip تاني .. هيطلع head .. انما لو flipped the coin multiple times .. ونص الوقت بيطلع هيد و النص التاني بيطلع tails .. فهنا انت عندك high uncertainty .. فبالتالي high entropy لأنك مش عارف ايه اللي جاي المره الجايه .. فبالتالي احنا في ال decision tree .. علوزين نقل ال entropy of the random variable that we try to predict

## Mutual Information

Let  $X$  be a random variable with  $X \in \mathcal{X}$ .  
 Let  $Y$  be a random variable with  $Y \in \mathcal{Y}$ .



Entropy:  $H(Y) = - \sum_{y \in \mathcal{Y}} P(Y=y) \log_2 P(Y=y)$



Specific Conditional Entropy:  $H(Y | X=x) = - \sum_{y \in \mathcal{Y}} P(Y=y | X=x) \log_2 P(Y=y | X=x)$



Conditional Entropy:  $H(Y | X) = \sum_{x \in \mathcal{X}} P(X=x) H(Y | X=x)$



Mutual Information:  $I(Y; X) = H(Y) - H(Y|X) = H(X) - H(X|Y)$

- Entropy measures the **expected # of bits** to code one random draw from  $X$ .
- For a decision tree, we want to **reduce the entropy of the random variable we are trying to predict!**

**Conditional entropy** is the expected value of specific conditional entropy  
 $E_{P(X=X)}[H(Y | X=x)]$

**Informally**, we say that **mutual information** is a measure of the following:  
 If we know  $X$ , how much does this reduce our uncertainty about  $Y$ ?

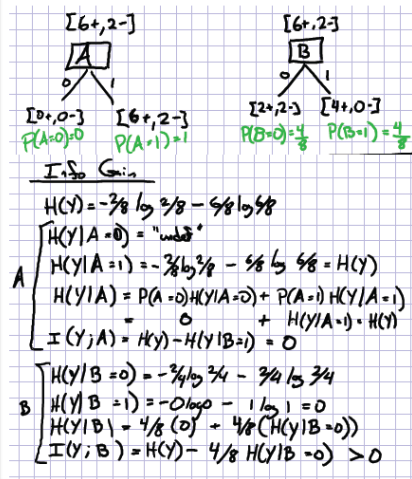
لو رجعنا لنفس الداتا سيت تاني .. نقدر تحسب ال mutual information as a splitting criterion

# 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



26

الدكتور دلوقة عاوز يتكلم علي Comparison of splitting criteria .. في ال practice .. انه في واحد من الطرق ديه احسن حاجه .. ازاي نقدر ن

work through which of these works best on datasets

الدكتور بيستعرض داتا سيت من ال 1992 .. كانت بتتكلم علي 4 splitting criterion on 12 datasets ....

## Experiments: Splitting Criteria

Bluntine & Niblett (1992) compared 4 criteria (random, Gini, mutual information, Marshall) on 12 datasets

### Medical Diagnosis Datasets: (4 of 12)

- hypo:** data set of 3772 examples records expert opinion on possible hypo- thyroid conditions from 29 real and discrete attributes of the patient such as sex, age, taking of relevant drugs, and hormone readings taken from drug samples.
- breast:** The classes are reoccurrence or non-reoccurrence of breast cancer sometime after an operation. There are nine attributes giving details about the original cancer nodes, position on the breast, and age, with multi-valued discrete and real values.
- tumor:** examples of the location of a primary tumor
- lymph:** from the lymphography domain in oncology. The classes are normal, metastases, malignant, and fibrosis, and there are nineteen attributes giving details about the lymphatics and lymph nodes

Table 1. Properties of the data sets

Data Set	Classes	Attr.s	Training Set	Test Set
hypo	4	29	1000	2772
breast	2	9	200	86
tumor	22	18	237	102
lymph	4	18	103	45
LED	10	7	200	1800
mush	2	22	200	7924
votes	2	17	200	235
votes1	2	16	200	235
iris	3	4	100	50
glass	7	9	100	114
xd6	2	10	200	400
pole	2	4	200	1647

ديه بتقارن ما بين 4 splitting criterion علي 12 داتا سيت .. ال 12 داتا سيتس مكتوبين عندك اهم .. احنا هنتكلم علي أول 4 .. عشان هم medical diagnosis datasets .. احنا هنركز علي Gini gain and Mutual information ....

## Experiments: Splitting Criteria

Table 3. Error for different splitting rules (pruned trees).

Data Set	Splitting Rule			
	GINI	Info. Gain	Marsh.	Random
hypo	1.01 ± 0.29	0.95 ± 0.22	1.27 ± 0.47	7.44 ± 0.53
breast	28.66 ± 3.87	28.49 ± 4.28	27.15 ± 4.22	29.65 ± 4.97
tumor	60.88 ± 5.44	62.70 ± 3.89	61.62 ± 3.98	67.94 ± 5.68
lymph	24.44 ± 6.92	24.00 ± 6.87	24.33 ± 5.51	32.33 ± 11.25
LED	33.77 ± 3.06	32.89 ± 2.59	33.15 ± 4.02	38.18 ± 4.57
mush	1.44 ± 0.47	1.44 ± 0.47	7.31 ± 2.25	8.77 ± 4.65
votes	4.47 ± 0.95	4.57 ± 0.87	11.77 ± 3.95	12.40 ± 4.56
votes1	12.79 ± 1.48	13.04 ± 1.65	15.13 ± 2.89	15.62 ± 2.73
iris	5.00 ± 3.08	4.90 ± 3.08	5.50 ± 2.59	14.20 ± 6.77
glass	39.56 ± 6.20	50.57 ± 6.73	40.53 ± 6.41	53.20 ± 5.01
xd6	22.14 ± 3.23	22.17 ± 3.36	22.06 ± 3.37	31.86 ± 3.62
pole	15.43 ± 1.51	15.47 ± 0.88	15.91 ± 1.15	26.38 ± 6.92

**Key Takeaway: GINI gain and Mutual Information are statistically indistinguishable!**

Info. Gain is another name for mutual information

ديه الإيرورز .. والصوره الجايه ديه بتعبر عن تفاصيل أكثر من ناحية ال statistical significance of the error difference between Gini gain and mutual information

## Experiments: Splitting Criteria

Table 4. Difference and significance of error for GINI splitting rule versus others.

Data Set	Splitting Rule		
	Info. Gain	Marsh.	Random
hypo	-0.06 (0.82)	0.26 (0.99)	6.43 (1.00)
breast	-0.17 (0.23)	-1.51 (0.94)	0.99 (0.72)
tumor	1.81 (0.84)	0.74 (0.39)	7.06 (0.99)
lymph	-0.44 (0.83)	0.11 (0.05)	7.89 (0.99)
LED	0.12 (0.17)	5.86 (0.00)	5.86 (0.00)
mush	0.00 (0.00)	5.86 (0.00)	5.86 (0.00)
votes	0.11 (0.55)	7.30 (0.00)	7.30 (0.00)
votes1	0.11 (0.55)	7.34 (0.00)	7.34 (0.00)
iris	0.11 (0.55)	7.50 (0.00)	7.50 (0.00)
glass	0.11 (0.55)	9.96 (0.00)	9.96 (0.00)
xd6	0.11 (0.55)	7.07 (0.00)	7.07 (0.00)
pole	0.11 (0.55)	4.43 (0.00)	4.43 (0.00)

**Key Takeaway: GINI gain and Mutual Information are statistically indistinguishable!**

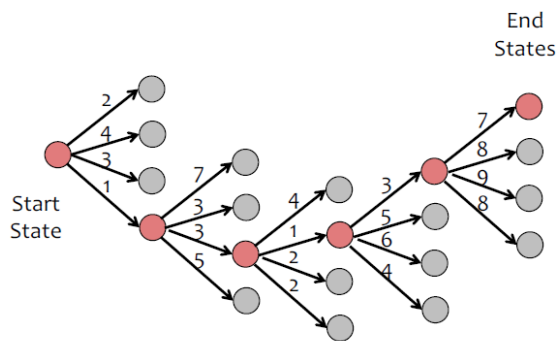
Results are of the form A.AA (B.BB) where:  
 1. A.AA is the average difference in errors between the two methods  
 2. B.BB is the significance of the difference according to a two-tailed paired t-test

أول رقم هو ال average difference و ثاني رقم هو ال significance of that difference اللي هو يعني  $1 - P$  .. أي حاجه أعلي من 0.95 هتبقا statistically significant difference .. لما هتيجي تبص علي أول عمود اللي هو Info. Gain .. هتلاقي ان مفيش ولا واحده من الداتا سيتس significant difference .. فاللي ناخذو من الحوار ده ان مفيش فرق ما بين الاتنين لما جم يقارنو ما بينهم ..

دلوقت هنعوز نتكلم علي inductive bias لل decision trees .. تعال دلوقت نتكلم علي ال algorithmic prespective الي ممكن تبقا مختلفه عن ايه اللي ممكن نعمله .. ايه هو ال greedy search .. مراجعه سريعه عليه .. عندك سيرش سببيس .. بتتكون من نودز و edges وال edges ديه معاها weights .. والهدف ان نلاقي lowest total weight path from a root node to a leaf node ... فال greedy search بت follow البروسيس ديه .. عند كل نوود اختار ال edge اللي ليها lowest immediate weight ديه و بعددين عبيد ثاني .. لو فضلت ماشي كدا هتوصل فعلاً للآخر بس بص علي اللي هيحصل في الصور الجايه .. وده حصل عشان انت طماع .. دائماً عاوز أقل حاجه



## Background: Greedy Search



**Goal:**

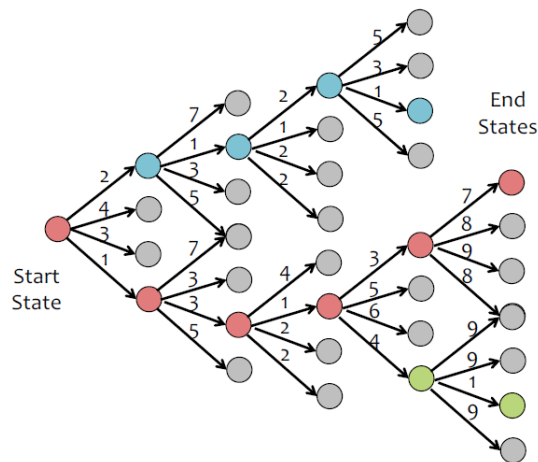
- Search space consists of nodes and weighted edges
- Goal is to find the lowest (total) weight path from root to a leaf

### Greedy Search:

- At each node, selects the edge with lowest (immediate) weight
- Heuristic method of search (i.e. does not necessarily find the best path)

38

## Background: Greedy Search



**Goal:**

- Search space consists of nodes and weighted edges
- Goal is to find the lowest (total) weight path from root to a leaf

### Greedy Search:

- At each node, selects the edge with lowest (immediate) weight
- Heuristic method of search (i.e. does not necessarily find the best path)

40

فال جريدي سيرش .. هي مش شرط تلاقي أحسن طريق .. عشان هي هترجعك الطريق الاحمر .. بس في طرق تانيه انت مدخلتهاش عشان انت علي طول طماع عاوز دايماً تكسب في اللحظة ديه .. انت عادي ممكن تخسر النورود ديه بس هتكسب مكانها قدام ... وده اللي بيحصل ف الدنيا عموماً ساعات بتقول هي ديه .. هي ديه .. بس مبتيقاش هي ... فأسلم حل .. سلمها الله و توكل علي الله .. شفت بقي لما هتفهم وتلاقي الكونسيبت شغال هتقتنع ازاى .. توكل علي الله وهو اللي هيفتحك الطريق اللي انت مش شايفو ..

دلوقت احنا عاوزين نفكر في ال decision tree learning as search .. عاوز نفهم ازاي نقدر ناخذ المشكله اللي شغالين عليها في ال decision tree  
و ن recast as a search problem .. اول حاجه هنفكر فيها هو ايه هو ال entire search space بتاعنا ..

1. السيرش سبيس بتاعنا هو كل الاحتمالات المتاحة لكل ال decision trees
2. Node: : كل نوود عبارة عن decision tree
3. Edge: : كل إيدج هتوصل ل a child ال parent from which it could have been created with 1 additional attribute

تعال نبص ع الصورة ديه . الارقام بالسالب خد اقل رقم اللي هو lowest weight path عشان ت maximize the spletting at each step

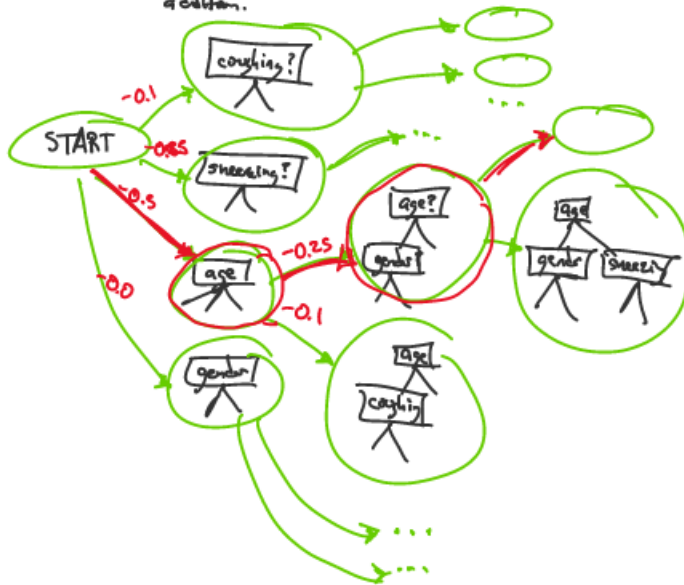
## Overfitting + KNN

9:27 AM Monday, January 27, 2020

### Decision Tree Learning as Search

- ① search space : all possible decision trees
- ② node : single decision tree
- ③ edge : connects a child to a parent from which it could have been created with 1 attrb. addition.

- ④ DT Learning = greedy search, maximizing or splitting criterion at each step
- ⑤ edge weight : **negative splitting criterion**



الصورة ديه هي هي نفس ال decision tree algorithm اللي اتشرح من اول محاضره ... هو هو بس الفرق بتبصلو ازاى .. احنا هنا بنبصلو ك greedy search maximizing our splitting criterion ... واحده من المميزات هنا في طريقة التفكير ديه ان في سؤال مهم جداً ... سؤال جنري في الماشين ..

How is it that your ML algorithm can generalize to unseen examples?

ده السؤال المهم جداً .. خلينا نأجل الرد عليه ثواني ...

تعال نتكلم دقايق بس :

هنحط ديفينيشن بس ... السؤال بتاعنا مرتبط بيه يعني ... ال inductive bias هو ال principal اللي من خلاله بيحصل generalization to unseen examples .. فلما الدكتور سأل انهى tree احنا هنلاقيها في الألجورزم بتاعنا .. هو كان يقصد ايه هو ال inductive bias بتاع ال decision tree learning .. ايه هو ال principle اللي هو بيشتغل بيه عشان يجيب شجرة كويسه .. من ضمن الشجر الكثير اللي فوق .. ازاى اختار ما بينهم ... حد رد وقالو اختار أقل عمق للشجرة ... نقدر نقول ان ال inductive bias هيحاول يختار اصغر شجرة ت match the data وليها أعلى mutual information attributes ...

## DT: Remarks

ID3 = Decision Tree  
Learning with Mutual  
Information as the  
splitting criterion

**Question:** Which tree does ID3 find?

### Definition:

We say that the **inductive bias** of a machine learning algorithm is the principal by which it generalizes to unseen examples

### Inductive Bias of ID3:

Smallest tree that matches the data with high mutual information attributes near the top

### Occam's Razor: (restated for ML)

Prefer the simplest hypothesis that explains the data

44

دلوقت الدكتور بيقولك عاوزنا نفكر في الجورزم مختلف شويه عن اللي احنا فيه .. احنا هنا هنعمل exhaustive global search .. بمعنى اننا هنعوز ن  
enumerate كل شجرة موجوده في السيرش سبيس .. كل احتمال موجود فكر فيه .. وبعدين انسي خالص ال edge weights ... وبعدين إسأل نفسك لكل  
شجرة انهي شجرة ليها أقل training error و ديه هنتبعا احسن حاجه .. بس ممكن يكون في اكثر من شجرة بيدوك 0 ايروور .. ازاي اختار احسن واحده  
منهم؟