

Lec: 16

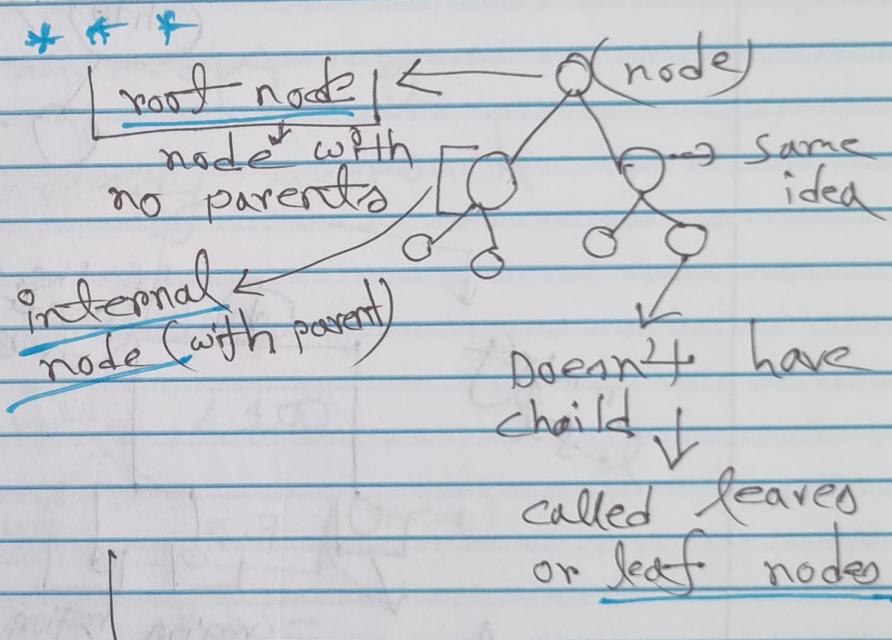
Tree-Based Methods

Decision tree → relationship between predictors and response

Tree-based methods

can be applied both in linear and classification

Stratify / segmenting the prediction

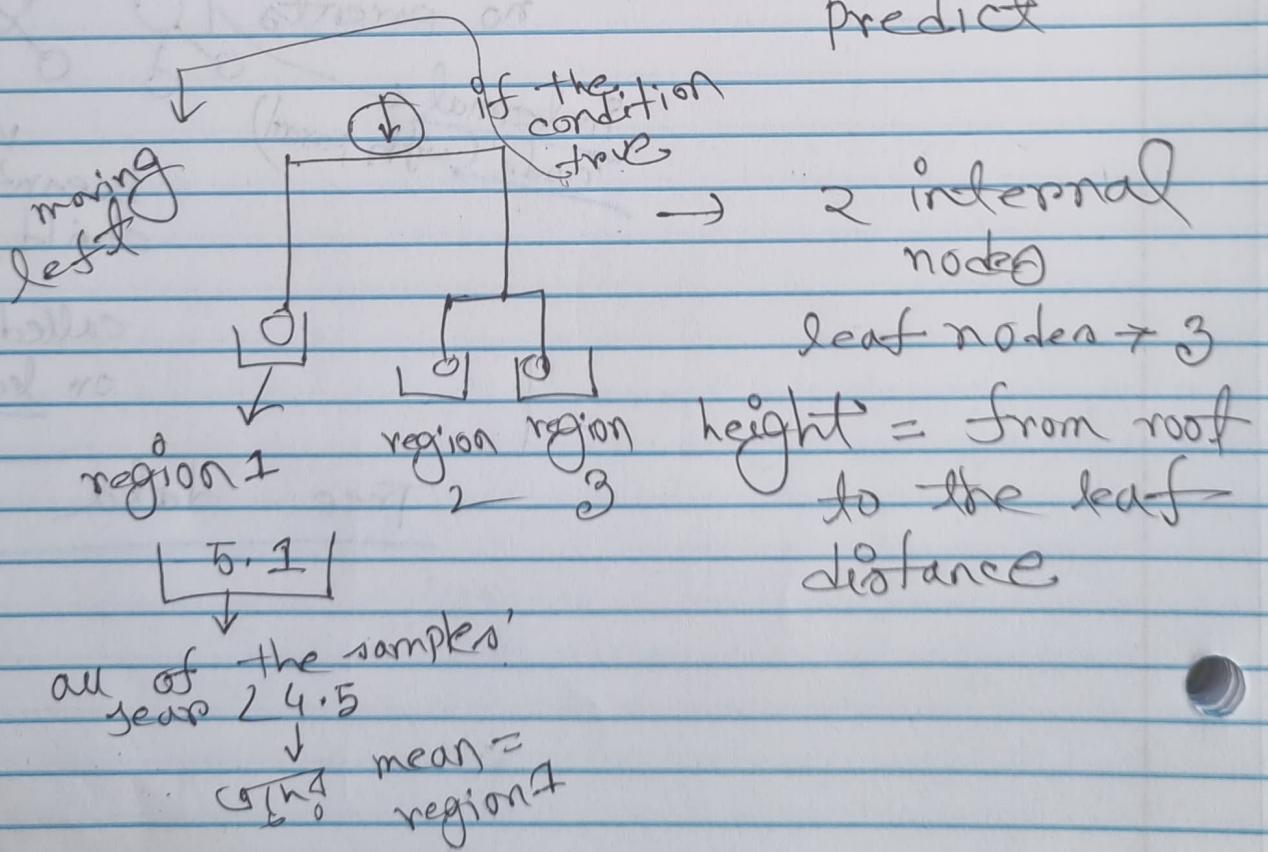
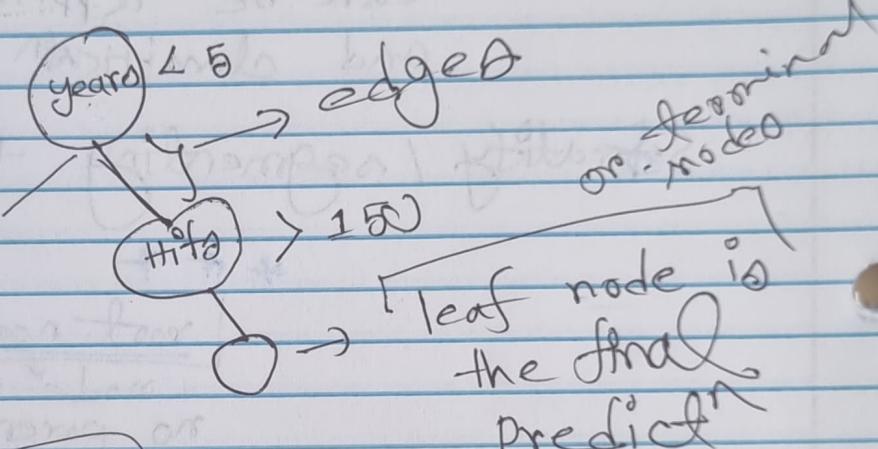


Lower the num of years

↓
lower hits

↓
Lower salary

Hits + Years → features



separating predictors space
called decision tree

How to build?

Step 1 : non-overlapping region or pr

" 2 : Sample mean in every slice region

" 3 : MSE

↓
minimizing by
↓
Greedy approach

#① selecting predictors

↓
H₀ < 100

4.0 5.7 (mean)

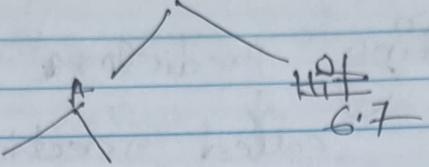
Years < 4.5 RSS =

5.2 6.5

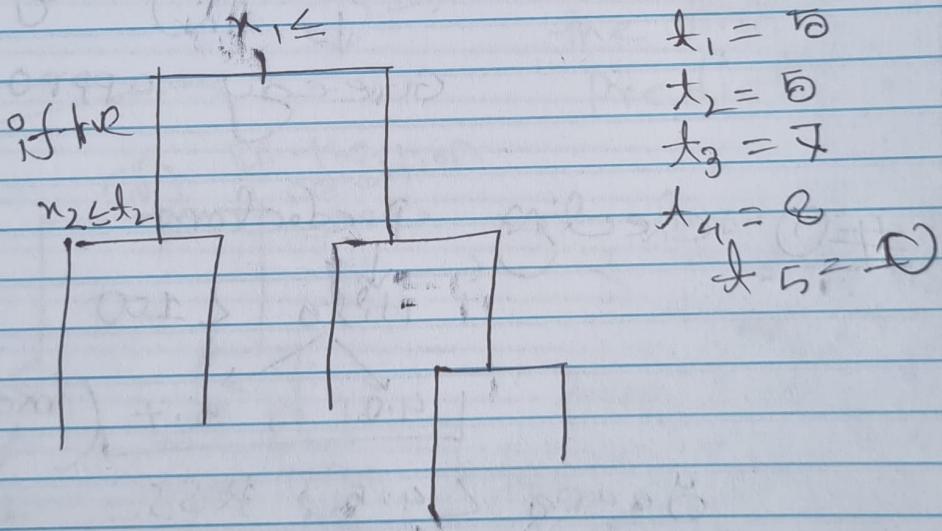
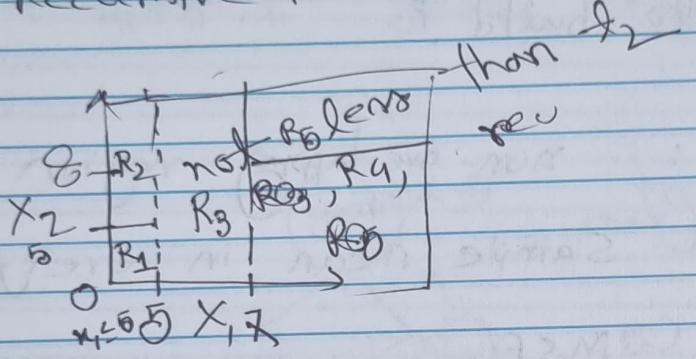
RSS =

based on lower num
root node selection

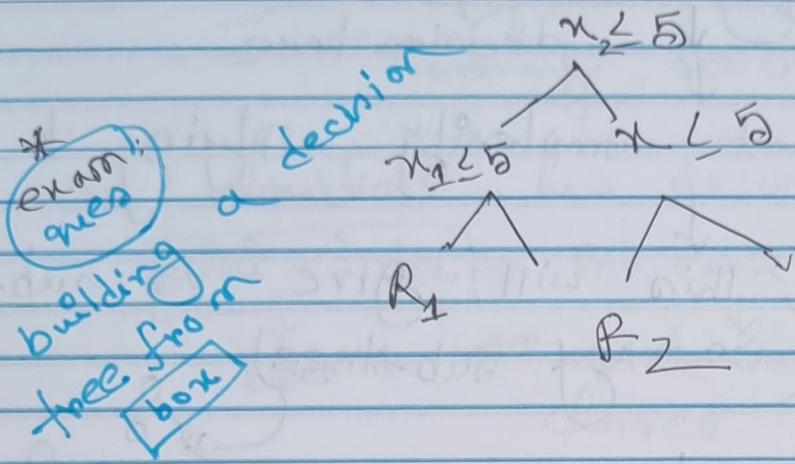
years < 4.5



Recursive process



root node = 0



Pruning a tree: cutting the unwanted parts

cont complexity pruning \rightarrow or weakest link pruning.

$$T = \text{Height}$$
$$= 50$$
$$\{\alpha = 0.01\}$$

need to be tested

During cutting \rightarrow too much pruning
useful branch can be lost

optimal α value using cross-validation

~~training~~ very big model (50 predictors)



decision tree

cost complexity applying to pruning



this will give us subset
(so many sub-tree)



Then cross-validation
on each individual tree
best one with the
minimum error

Baseball example \rightarrow 3rd is good
instead of 10

Lec: 16
03-13-24

classification tree

commonly occurring classes

In classification, how do we calculate errors?

Gini-index → for each internal node

Lower the better Total variance across all the classes

purity of node

root node

↑↑ Gini-index

but at the bottom of decision tree

Gini index ↓

Cross-entropy:

Gini index → underfitting → NO value

→ overfitting like root node

too low (overconfident)

yes no

Disadvantages → not much accurate

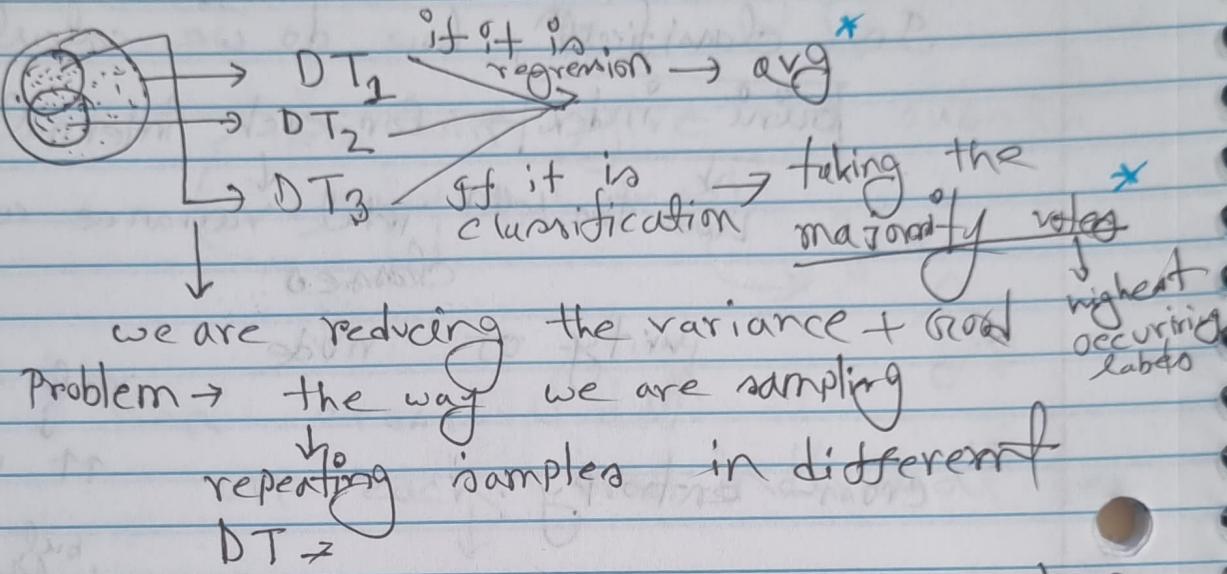
Advantages: → easy to explain

→ "understand graphically representation + easy interpretation"

→ can easily handle qualitative Predictors

Bagging decision tree :

Bagging: $B = 3$ sample randomly from the dataset 3 times



3 features:

x_1	x_2	\dots	x_{10}
0			
1			
n			

random sampling with replacement

$0, 1, 2, 3, 4 \rightarrow DT_1$

$0, 2, 4, 5, 6 \rightarrow DT_2$

repeating sampling

can't reduce the variance

0, 1, 2, 3, 4

Others than this sampling called
out of bag

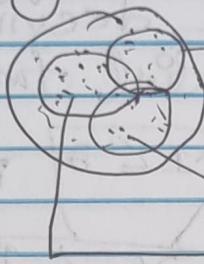
improving bagging:

Random forests,

wanna handle randomness which is common in
bagging

↓
→ introducing randomness

data →



DT, randomly selecting features

x_1, x_2, x_4

↓

though repeating occurs, but
they are unique

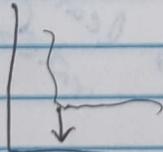
random sampling and } creating forest
↳ features

16 predictors : $\sqrt{16} = 4$

$m \approx 4$

small trees coming together to
build a forest

↓
error is minimized in
thin than the bagging and
other method

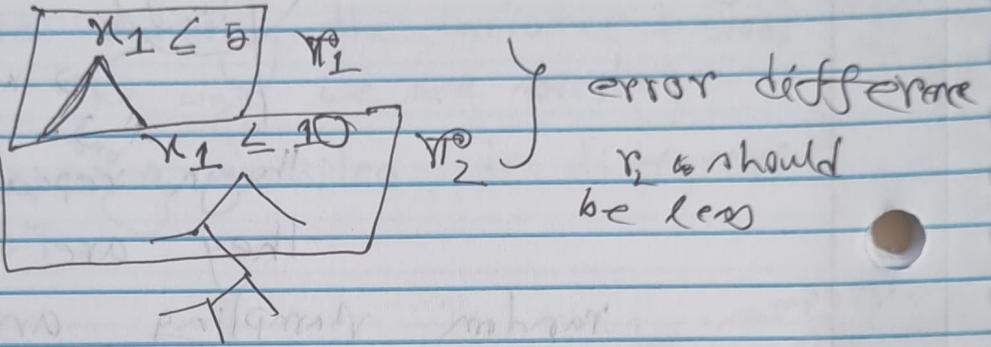


at go error is minimum

90 trees will be better than the
300 trees
disadvan.

Boosting!

Build the decision tree step by
step / sequential fashion



Baggin + random for

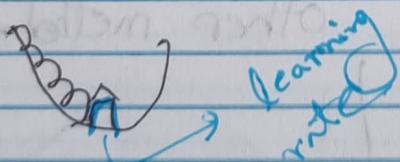
① \rightarrow number of trees (by guess)

different values of
B by cross-validation

② Shrinkage parameter in boosting

controls how boosting
learns

small positive numbers



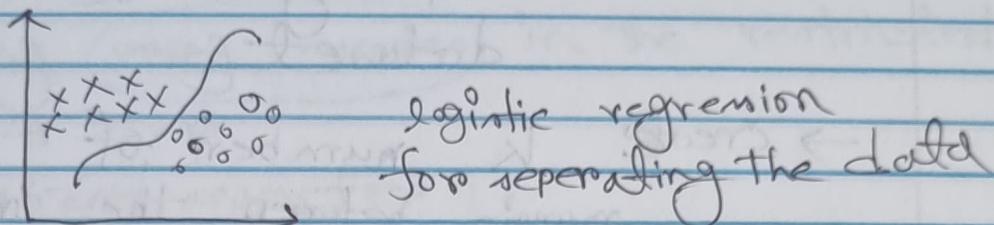
(ii) number of options

Final $\times m$ [20 - multiple question
2 → case study → per 5]

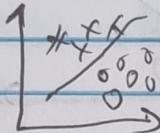
Lec: 17
04-01-24

"Support vector machine"

2 class / 2 labels within the dataset



But in SVM \rightarrow separating by using a line

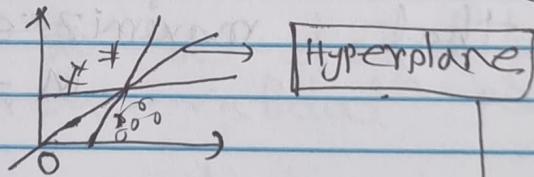


If gets tough to separate

\downarrow
we will use SVM

feature expansion \rightarrow polynomial lines to separate the data

hyperplanes:



coefficient doesn't matter

From many, which one can we select?

criteria for selecting best hyperplane:

- i) Maximal Margin classifier

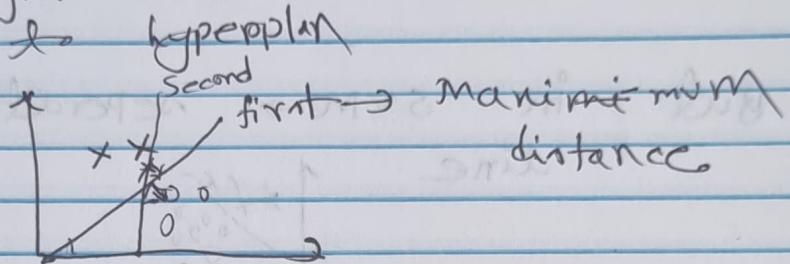
41:30
13-10-10

i) Maximal Margin classifier:

Finds all possible hyperplanes

↓
depending on the size of the dataset.

→ create k number of hyperplanes margin between the immediate sample



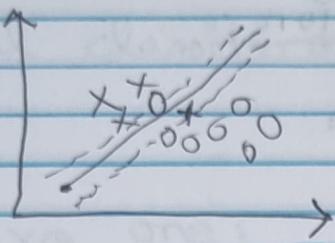
$$|\vec{w}| = 1$$



→ Select the best hyperplane that maximizes M
 $M = \text{Largest Margin}$

com: It can't be used if the data isn't perfectly separable.

ii) Support vector classifier:



allowing wrong samples in the restricted area \rightarrow SVM
 c -regularizer/parameter :

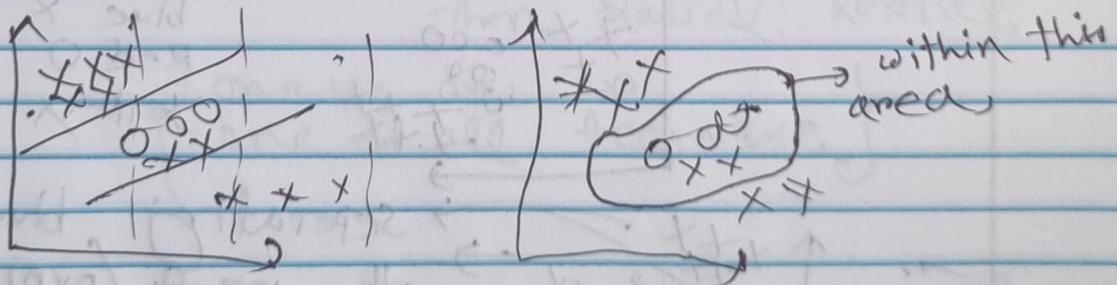
if $c \uparrow \uparrow \rightarrow$ we are giving more space for the wrong samples / errors
 if $c \downarrow \downarrow \rightarrow \downarrow \downarrow$ errors

$$c_1 \leq c_2$$

Tradeoff: if $c \uparrow \uparrow \rightarrow$ error $\uparrow \uparrow$
 but model will be robust

but if $c \downarrow \downarrow \rightarrow$ error $\downarrow \downarrow$
 but model will be less robust

so need to be checked by checking with the cross-validation parameter



feature-expansion:

amplifying the prediction:

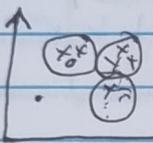
③ SVM:

kernels = hyperplane function

nonlinearities and kernels:

[Radial kernel]:

overfitting may occur:

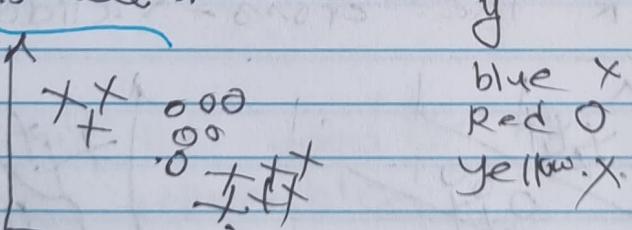


within radial kernel $\rightarrow \alpha$ -value is
for controlling the parameter

SVM \rightarrow works for 2 classes:

how can we extend it?

a) one vs all:



separating blue from
the rest (ova)

separating red from
the rest

L → yellow

final classification based on majority:

b) one vs one:

one vs others

Blue vs Red

Blue vs Yellow

Red vs Yellow

* combination will take time → in case of bigger sample.

separation:

Summary: maximal margin → if perfect is necessary

↓
if not, then allowing

for errors

↓
by

SVM classifier

controlling margin with C-

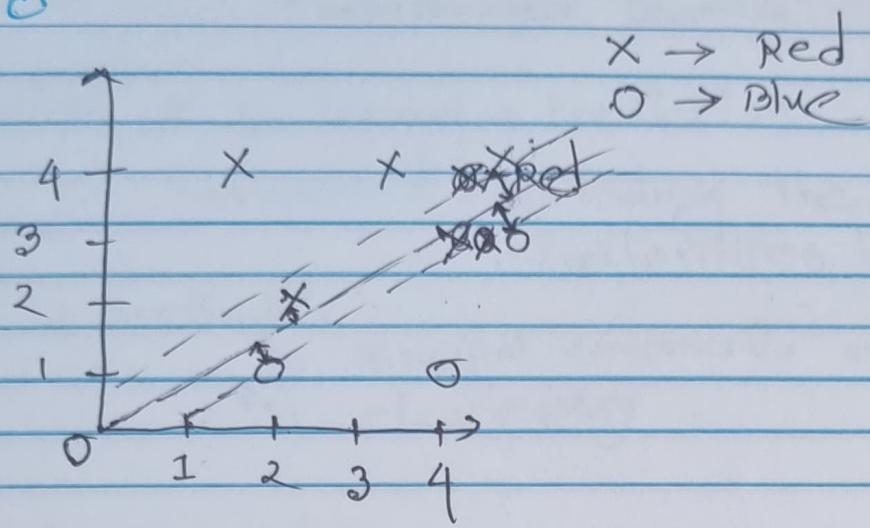
samples that are on the edges

{non-linear} → SVM → kernel

↓
then Radial kernel

{ one vs one or
one vs all technique }

class problem:



- ① sketch the optimal separating hyperplane
- ② Indicate the SV