

Lec: 10
02-18-24

"Resampling Method"

2 types of Resampling:

- i) Cross-validation
- ii) Bootstrap

overfitting : memorizing but not learning \rightarrow doing bad

underfitting : complex model

optimum fit: between overfitting and underfitting

① Cross-validation

* validation set approach
(Test-set, val-set, train-set)

Training
 \downarrow
enough data
to learn

Testing

80%

70%

90%

20%

30%

10%

valid set

validatⁿ - set → we will use
to test the model.
During training period

Train	val set	test
70 :	10 :	20

training → 200 times

first step ↴ checking / monitoring
iteration

2nd point at 20%, then model
change needed

classical → Logistic and Linear

↓
Training process time is less
multiple iteration

validatⁿ is applicable in deep-NH.

Advantages: faster, simple and easy to
understand

Training and testing error

Disadvantages: (i) split reduced the size
of the data (losing data)
(ii) random split → easy

(iii) sample in the test data
↓ overestimate

2. k -fold validation:

instead of splitting

we divide data according to k value.

model 1 is going to be tested with
2, 3, 4, 5 and validated by ①



Avg accuracy of all 5-model

if it is classified \rightarrow 3 cat 2 dog

↓
cat

variant of k -fold validation:

→ LOOCV



Leave one out cross-validation



splitting data according to
the no of sample.

If $n = 6$

$n = \boxed{|||||}$

Test with ①

Disadvantages \rightarrow slow

Advantages \rightarrow → accurate

→ all the data are getting
tested

best value of $k = 3, 5, 7, 10$
5 and 7 \rightarrow good accuracy

5 models \rightarrow unique test
rest - train
Loo \rightarrow (one is getting left-out)

② Bootstrap approaches:

General statistical technique
(Decision tree)

We want to quantify an association
between estimation

Sampling with replacement
randomly select
then new dataset

disadvantage: \rightarrow same sample can be
in the same all dataset
 \rightarrow performance low
might not reflect the real
dataset underestimate prediction

How to remove overlapping?

Sampling without replacement

Summary: Resampling

↓
cross-validation (2 methods)

k-fold
validation

Bootstrap

used in decision trees
** in resampling → not good

Vee 12
02-26-24

Model select and regularized

#

Dataset

Resampling method



Focuses on sample

	x_1	x_2	$\dots x_p$	y
1				
2				
⋮				
⋮				

But in model select



Focusing on the Predictors

Subset select
shrinkage

$$\hat{y} = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

Dataset → more predictors than the sample

$p = 50$
 $n = 20$

Difficult to control
the variance

↓
Feature select needed

i) Subset-select method:

Selecting some of the predictors:

Best subset-selection;
1st technique:

Best \rightarrow

① Model with no predictors

$k=0$

↓
 $M_0 = \beta_0$ (intercept/ mean of rand)
null model

② Fitting all combinations of P and k

$$k=1 \quad M_1 = \beta_0 + \beta_1 X_1$$

$$M_2 = \beta_0 + \beta_1 X_2$$

$$M_3 = \beta_0 + \beta_1 X_3 \quad Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2$$

n_1	n_2	n_3	y
-------	-------	-------	-----

$$k=2 \quad M_4 = \beta_0 + \beta_1 X_1 + \beta_2 X_2$$

$$M_5 = X_1 + X_3$$

$$M_6 = X_2 + X_3$$

} finding
the best one
if this, $k=3$

$$k=3 \quad M_7 = X_1 + X_2 + X_3$$

③ Step 3 → finding best model

M2, M6 and M7 comparison
by cross-validation

↓
selecting the best one

p = 50 predictors

k = 1 (all possible combination)

↓

dis: i) Time consuming

Advantage: i) Best possible model
is guaranteed
(as we are trying all possible
combination)

second technique of
subset-selection:

Getting ideas from sequential
moving

A lot of predictors

so overfitting may occur

so instead of moving at a time

- Best subset select
- forward stepwise select

$$\textcircled{1} \quad M_0 = \beta_0$$

first predictor $x_1 \rightarrow$

$$M_1 = \beta_0 + \beta_1 x_1$$

if this
is good

$$M_2 = \beta_0 + x_2 \rightarrow \begin{array}{l} \text{if accuracy} \\ \downarrow \text{than } M_1 \end{array}$$

$$M_3 = \beta_0 + x_2 + x_1$$

$k=1$

$$M_4 = \beta_0 + x_3 + x_2$$

moving to
 M_3

$k=3$

$$M_5 = \beta_0 + \cancel{x_3} + \cancel{x_1} + x_2 + x_1$$

if this is
good

$k=3$

$$\textcircled{1} \quad M_0 = \beta_0$$

$$\textcircled{2} \quad M_1 = \begin{cases} k=1 \\ M_1 = \beta_0 + x_1 \end{cases}$$

$$M_6 = \cancel{\beta_0} + x_2 + x_1 + x_3$$

not at this point

M_4 is good we found

Step 3 :

M_2 vs M_4

By cross-validation
best model select

Advantage: Best subset selection is good

Backward stepwise selection:
model with all predictors

$$\textcircled{1} \quad M_P = \beta_0 + x_1 + x_2 + x_3$$

$$\textcircled{2} \quad k = P - 1, \quad k = 2$$

better \checkmark $M_1 = \beta_0 + x_1 + x_2$ (without 1 Predictor from M_P)

$$\begin{array}{c} M_1 \\ \cancel{M_2} \\ \cancel{M_3} = \cancel{x_2 + x_3} \end{array}$$

$$k = 1, \quad M_3 = x_2 + x_3$$

$$M_4 = x_1 \quad \checkmark \text{ better}$$

$$M_5 = x_2$$

$$k_0 \quad M_6 = \beta_0$$

③ M_1 vs M_4 y cross-validation

↓
best selected

i) ~~Dis~~ [fitting all the predictors]
[at a time \rightarrow not viable]

[Differences are in forward and backward processes]

ii) ~~Dis~~ might loss the best one model

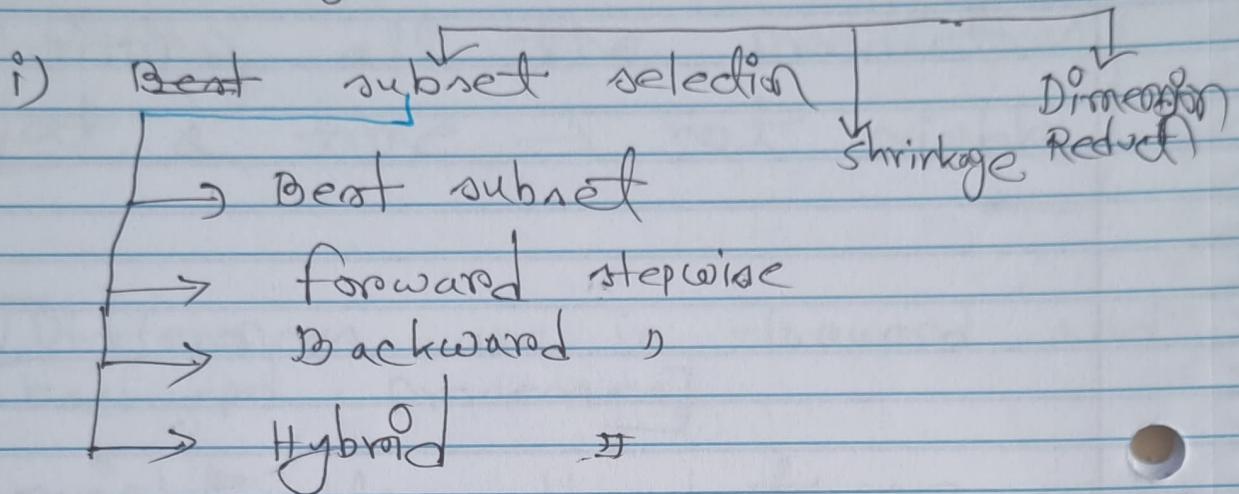
Shrinkage methods:

Vec: 10
02-28-24

we need essential predictors

↓
Systematic way
↓

3 - classes are:



Hybrid: $x_1 \ x_2 \ x_3 \ x_4$

Evaluation criteria

forward \rightarrow \downarrow \rightarrow then backwards \rightarrow
evaluation \rightarrow then forward

$F \rightarrow E \rightarrow B \rightarrow F \rightarrow E \rightarrow B \rightarrow E \rightarrow F$

$E \rightarrow R^2$ and b(B)

$$F: \begin{cases} y = x_1 \\ y = x_2 \\ y = x_3 \\ y = x_4 \end{cases}$$

M_1

$$\begin{cases} y = x_2 + x_3 \\ y = x_2 + x_4 \\ y = x_2 + x_1 \end{cases}$$

(M)
best

$$\begin{cases} y = x_2 + x_3 + x_4 \\ y = x_2 + x_3 + x_1 \end{cases}$$

(M)
best

P values $\rightarrow x_2, x_3, x_4 \rightarrow > 0.03$

$$R^2 \rightarrow \rightarrow$$

If the model meets the criteria
↓

then stop

Otherwise we will go
backwards



$$y = x_2 + x_3 \downarrow$$

$$y = x_1 + x_2 \downarrow$$

$$y = x_1 + x_3 \downarrow$$

M₁

$$y = x_2 \downarrow$$

$$y = x_3 \downarrow$$

M₂

Evaluation \rightarrow if good

then stop

Problem solving qua

Time $\downarrow \downarrow \rightarrow$ forward

Shrinkage method:

Forward

Accuracy \rightarrow hybrid

If we have time \rightarrow best subset

sales target \rightarrow Dec; 2nd

target \rightarrow next year

best subset

best model is guaranteed

if plenty time 2010 fit fast

time \rightarrow 2 (m)
forward fit best
Accuracy \rightarrow Hybrid
in case of

* Limited time and resources

backward won't be a good model

Shrinkage method:

↓
Get rid of unnecessary predictors

↓
those that aren't statistically significant

we want to shrink the less coefficient values.

- Ridge
- Lasso

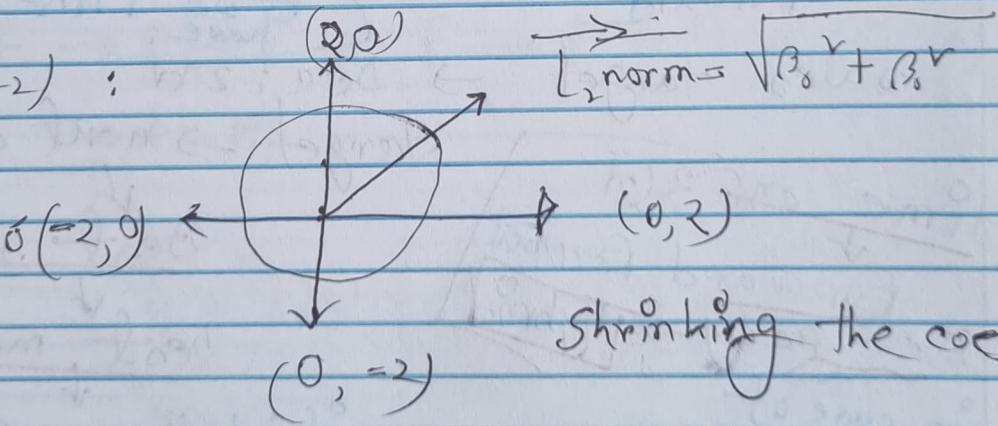
Ridge regression

$$y = \beta_0 + \beta_1 x_1 \quad | \quad y_2 = \beta_0 + 5.6(x_1)$$

$$x_1 = 0.5$$

$\beta_1 = 0.5$, $P > 0.05$
 $\beta_2 = 5.6 \rightarrow$ significant impact

Ridge (L_2):



Dis adva : \rightarrow still have all the predictors

Lasso : shrinking the co-efficient to 0

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

$$\beta_1 > 0$$

$$y = \beta_0 + \beta_2 x_2$$

pushing towards 0 \rightarrow very small numbers will get to the zero

but in ridge \rightarrow 0 এর বেশি কিম্বা কম কিম্বা যাবে but not exactly 0.

* Dimension Reduction methods :

Before the creation of model,

feature selection

$$n \text{ rows} \times \text{columns} \rightarrow 4/3$$

$$P = 50 \\ n = 200$$

$$n \times P \\ 200 \times 50$$

We want to
↓ the number
of predictors
columns

= Data के variability को कम करने के लिए Predictor removal.

$$200 \times 50$$

$$200 \times 20$$

* PCR (Principal component's Regression):
पर PCA:

way of putting together \rightarrow predictors

Step 1: Normalize the data,

①

Same range in terms of numbers

Age salary (per year)

200 \rightarrow Highest 100,000

0 \rightarrow newborn 0

Model will give significant

STAGE
but age को underestimate करता है
numerical values को normalize

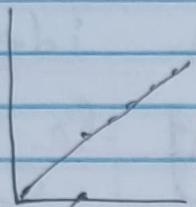
जहि तरी

200,

2nd step 2: finding the largest variance

highest coefficient of principle components finding

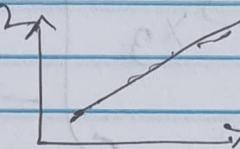
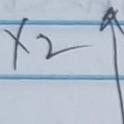
Wei	height	age	salary	/
↓				
pc1	age			



x_2 - features highly correlated

will be put together

so to form one component



will be the same

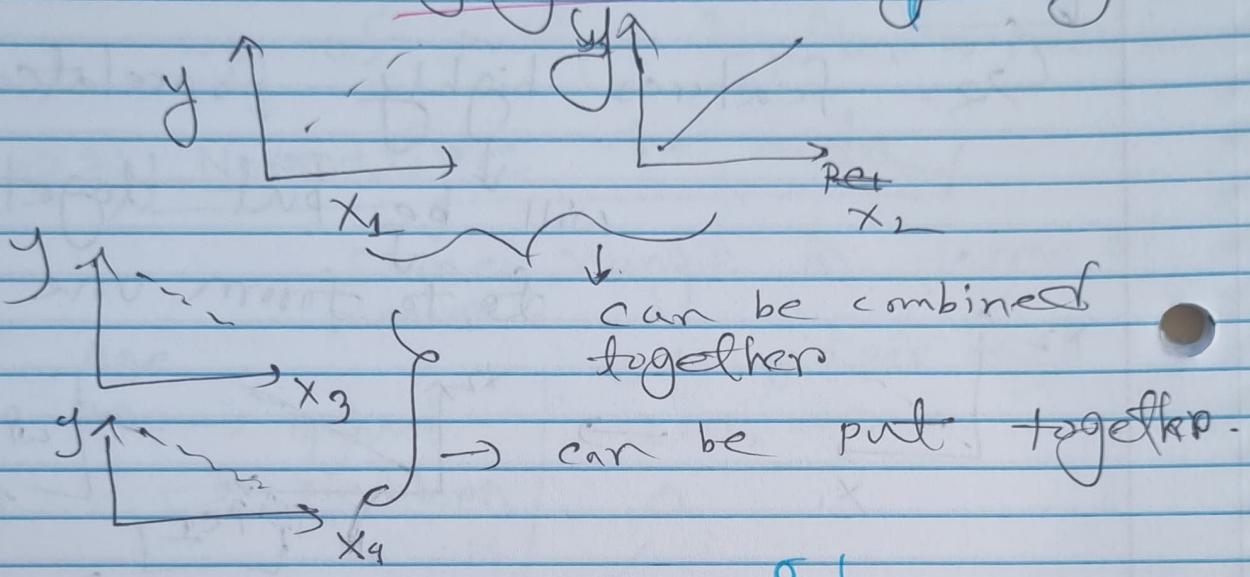
$\underbrace{PC1}_{(x_1 + x_2)}$

we are
not missing
any important
data)

↔ Partial Least squares:

↳ PCA → we do not use \vec{y}

* We go through the PCA idea while paying attention to \vec{y}



PCA → unsupervised

PLS → supervised

Lec: 14
03-04-24

"Coding"

Grid-search

~~best model~~
~~searching with the~~
~~best lambda value~~

↑↑ log value ($\log(100)$)
provided → ↓
log(10)

-10 -5 0

Here coefficients
are 0

↑↑ λambda
↓
avg coefficients
will be lower
as it's punishing
the model to shrink
the co-efficients

normalizing the data
and providing to model

++
Lasso isn't stable
ridge is stable

default → PCA → ③

CV MSE:

6 com and 17 com
↓

good

but 6 com
better

+ otherwise fine ↑

PCA PLS

6 12

Accuracy
+ efficiency

if
have the
computational
power

+ not
enough
computation