

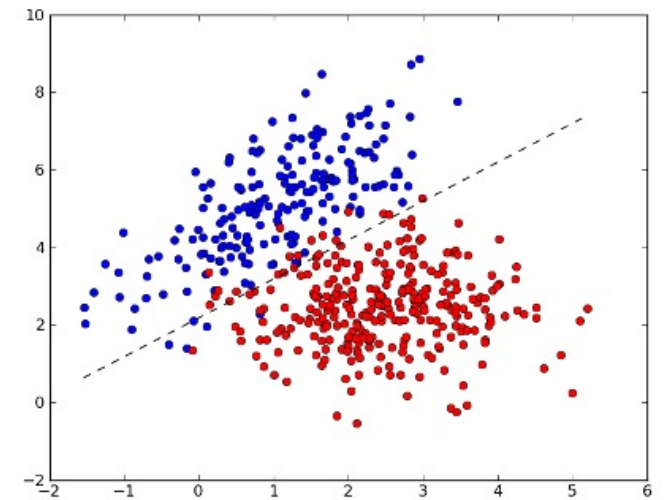
Classification model

X
Y

(Features,
Independent var,
Input var)
(class label,
Dependent var,
Output var)

Mean radius	Mean texture	Mean perimeter	Mean area	...	Tumor type
17	10	122	1001		0
20	17	132	1326		0
19	21	130	1203		0
11	20	77	386		1
...
21	22	142	1479		1
20	28	131	1261		0

N
(Rows,
Samples)



X

Y

Train {

Mean radius	Mean texture	Mean perimeter	Mean area	...	Tumor type
17	10	122	1001		0
20	17	132	1326		0
19	21	130	1203		0
11	20	77	386		1

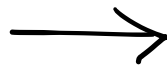
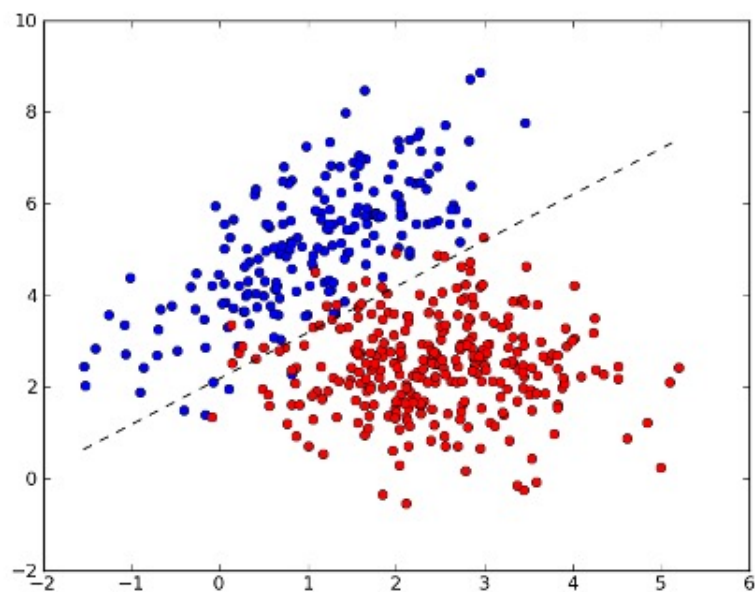
Test {

21	22	142	1479		1
20	28	131	1261		0

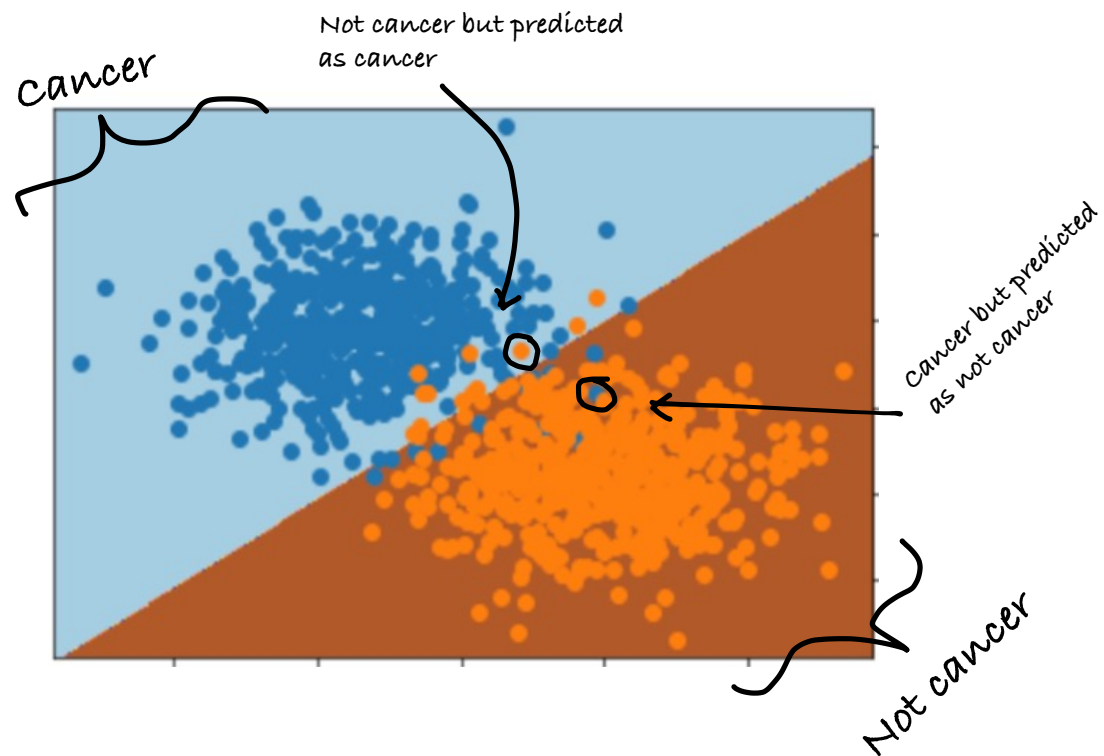


	X	Y
Train	x_train	y_train
Test	x_test	y_test

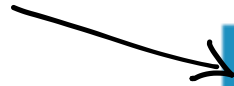
Real data



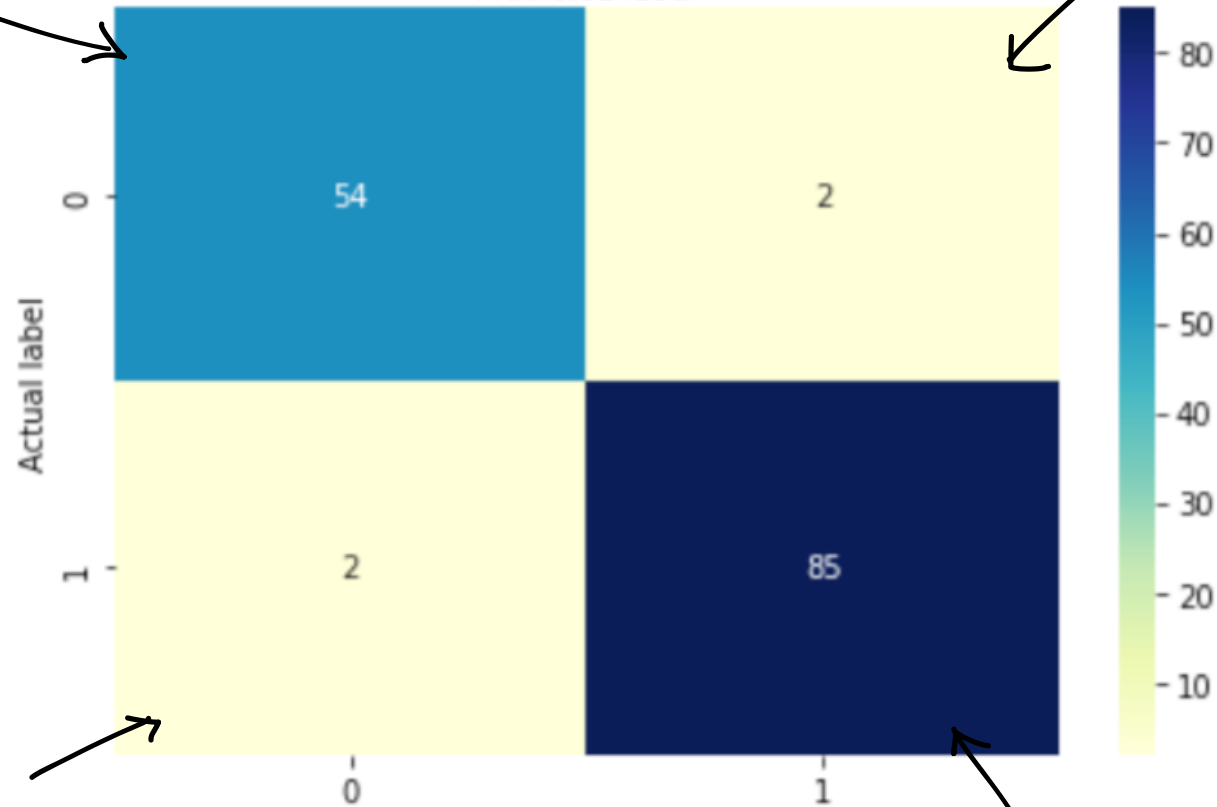
Predictions



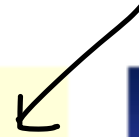
Not cancer and predicted
as not cancer



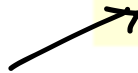
Confusion matrix
Predicted label



Not cancer but predicted
as cancer



Cancer but predicted
as not cancer

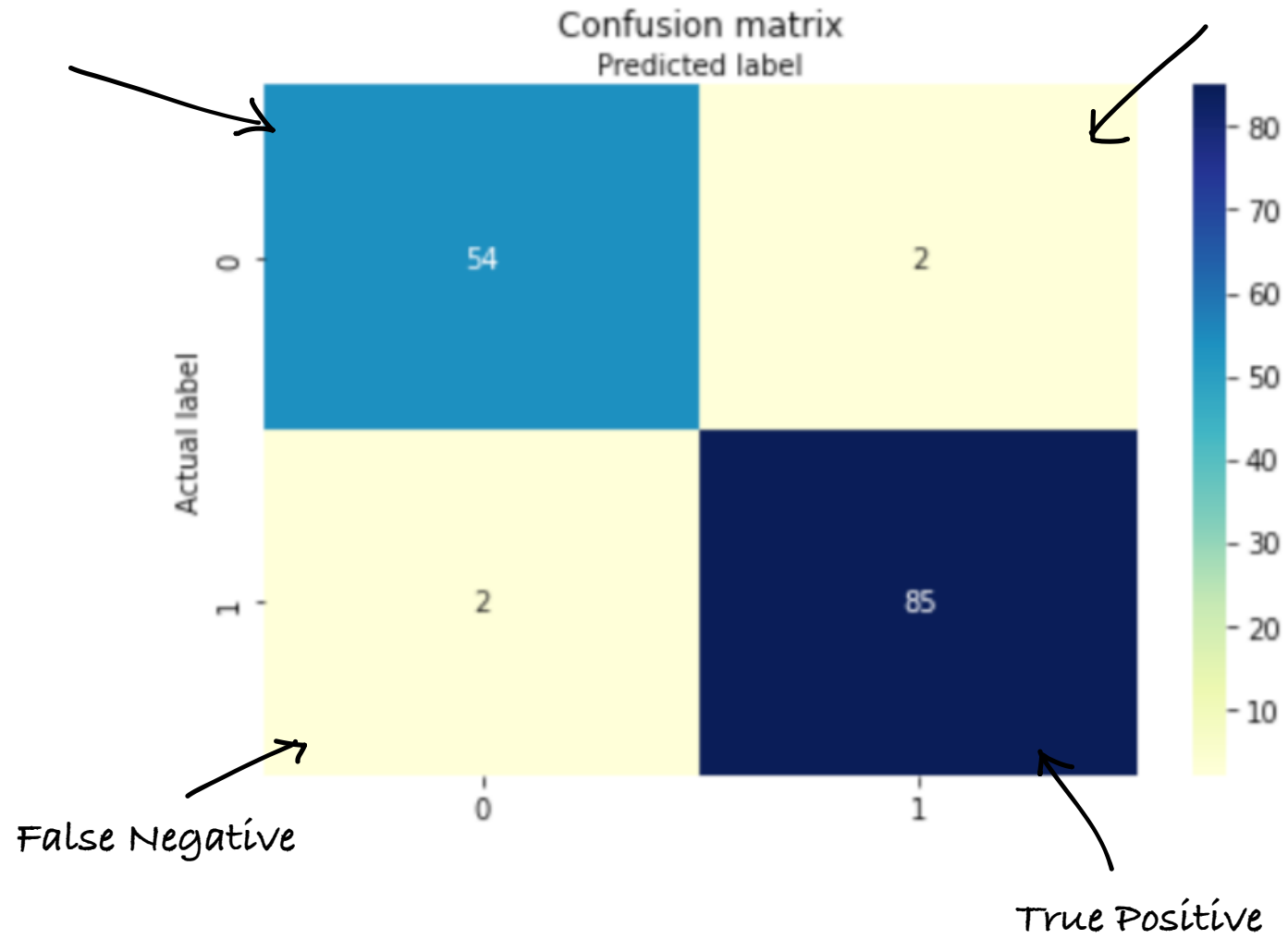


Cancer and predicted
as cancer



True negative

False positive



What is CV

	X				Y	
	Mean radius	Mean texture	Mean perimeter	Mean area	...	Tumor type
Train	17	10	122	1001		0
	20	17	132	1326		0
	19	21	130	1203		0
	11	20	77	386		1
Test	21	22	142	1479		1
	20	28	131	1261		0

Test data

Fold 1

Mean radius	Mean texture	Mean perimeter	Mean area	...	Tumor type
17	10	122	1001		0
20	17	132	1326		0
19	21	130	1203		0
11	20	77	386		1
14	23	143	556		0
21	22	142	1479		1
20	28	131	1261		0

Score 1

Fold 2

Mean radius	Mean texture	Mean perimeter	Mean area	...	Tumor type
17	10	122	1001		0
20	17	132	1326		0
19	21	130	1203		0
11	20	77	386		1
14	23	143	556		0
21	22	142	1479		1
20	28	131	1261		0

Score 2

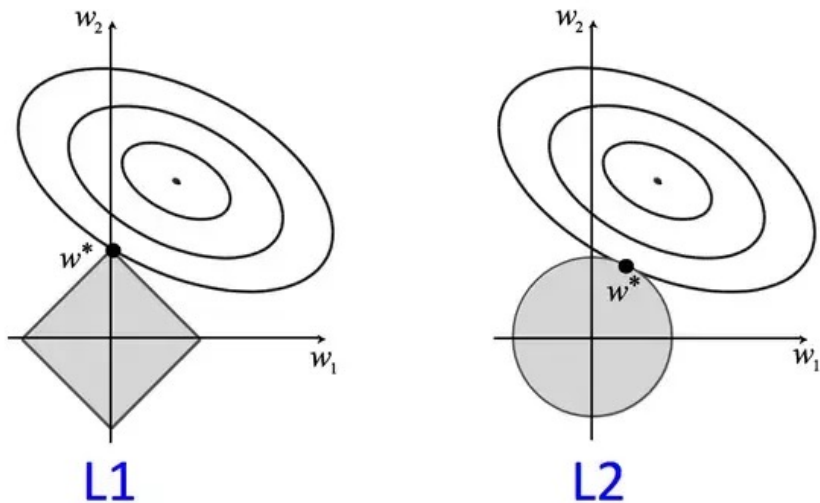
Fold 3

Mean radius	Mean texture	Mean perimeter	Mean area	...	Tumor type
17	10	122	1001		0
20	17	132	1326		0
19	21	130	1203		0
11	20	77	386		1
14	23	143	556		0
21	22	142	1479		1
20	28	131	1261		0

Score 3

Mean Score of Test data

What is grid search?

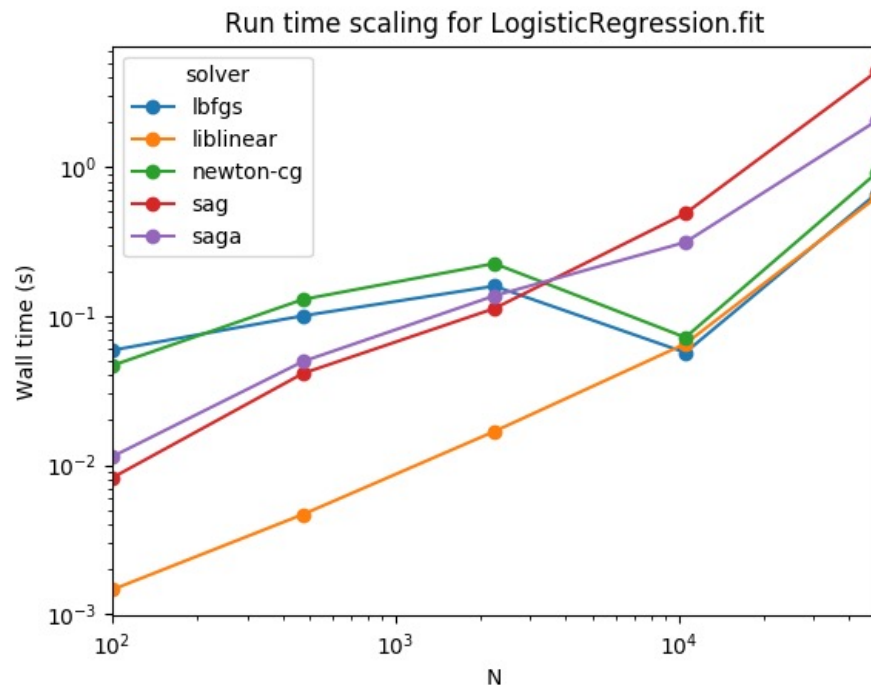


penalty:

- 'none': no penalty is added;
- 'l2': add a L2 penalty term and it is the default choice;
- 'l1': add a L1 penalty term;
- 'elasticnet': both L1 and L2 penalty terms are added.

4

X



Solver

{'newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga'},
default='lbfgs'

- For small datasets, 'liblinear' is a good choice, whereas 'sag' and 'saga' are faster for large ones;
- For multiclass problems, only 'newton-cg', 'sag', 'saga' and 'lbfgs' handle multinomial loss;
- 'liblinear' is limited to one-versus-rest schemes.

5

= 20

• 'none' {
newton-cg'
'lbfgs',
liblinear',
sag'
'saga'

• 'l2' {
newton-cg'
'lbfgs',
liblinear',
sag'
'saga'

• 'l1' {
newton-cg'
'lbfgs',
liblinear',
sag'
'saga'

• 'elasticnet' {
newton-cg'
'lbfgs',
liblinear',
sag'
'saga'

X 3

} Mean score of
the three splits