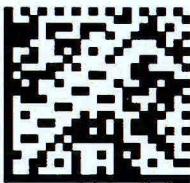


NOM BOUKHELOUA

Prénom Rémi

Promo 2020

Date 10.01.19



20150257: BOUKHELOUA Rémi

M1:

ST2ML1-DE (10/01/2019)

Amphi orange

15,5

MATIÈRE Machine Learning (1/2)

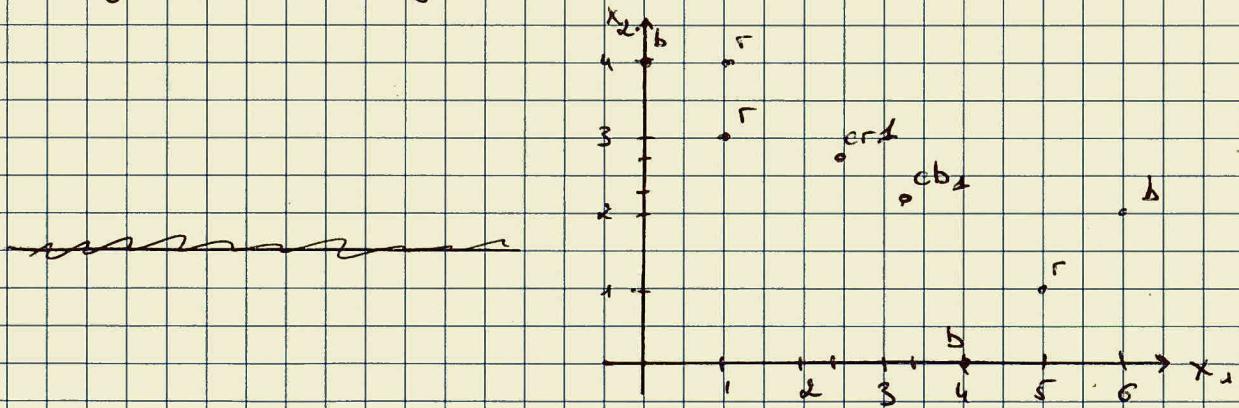
1. kNN regression is a method used to predict values based on nearest neighbors. ($\frac{1}{k} \sum_{i=1}^k y_i$) meanwhile kNN classification is a method to ~~show~~ separate the dataset into groups & based on distance to the neighbors. Its output will be a label (red, blue or green color).

2.3. With logistic regression, we have a higher ~~one~~ error rate than in kNN classification in both training and testing set.

~~Moreover, we have different terms in logistic method which means that~~ So, kNN classification has a better accuracy. But, logistic has a higher interpretability, it is easier to understand.

3. Any way, I would use kNN classification because of its ~~secondary~~ error rates.

S.



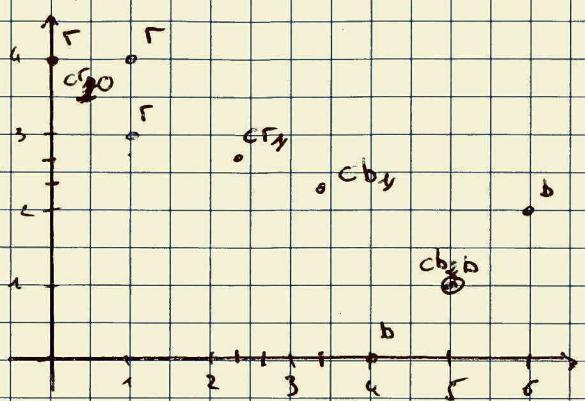
we randomly assign ~~class~~ label to the points (b or r)

we calculate the center of the two groups:

$$cr_1 = \left(\frac{1+1+5}{3}, \frac{4+3+1}{3} \right) = \left(\frac{7}{3}, \frac{8}{3} \right) \text{ and } cb_1 = \left(\frac{10}{3}, \frac{2}{3} \right).$$

Now we will assign label of points based on closest center point.

If ~~you are~~ you are closer to cr_1 we label you r (same with cb_1).



now we calculate new centers based

on new labels:

$$c_{r_2} = \left(\frac{2}{3}, \frac{11}{3} \right) \text{ and } c_{b_2} = (5, 1)$$

now we assign labels based on
 c_{r_2} and c_{b_2} and we find same
groups so we can stop.

we will have as result

Obs x_1 x_2 Label

1	1	4	r
2	1	3	r
3	0	4	r
4	5	1	b
5	6	2	b
6	4	0	b

6. That means that ~~of~~ the value of the first component variable
counts as 10% in the calculation of the estimated value.

If we add 1 to the value of the variable, the estimated
value will gain 10%.

7. k-fold cross validation is based on train-test sets. It
will cut the set into k-folds and will perform train
on $k-1$ fold and test on the one not used in train.

It will do it ~~a~~ k times in order to use ~~at least~~ one
every fold as a ~~train~~ test set. We commonly use $k=5$ or 10 .
~~or validate~~ something only in two groups

It will keep on and MSEs and choose best parameters.

a. ~~Dividing in two sets means that your model will fit b-fold validation will not have variance problem as it didn't use~~
b-fold performs a lot of train-test and picks the best so it will quite always have a better accuracy than just train-test once. But as it performs a lot of test it takes more times to execute than just dividing in two sets.

b. Leave one across means doing k-fold with only one value in folder. So on large dataset, it will be much longer to perform than a k-fold. Moreover it uses almost all data as train so model can have high variance we can't see because of small test set.

8. ~~Assume~~ It seems that we used a linear regression for a non-linear model. ~~because~~ because we have residual that are forming sometimes low, sometime high but "equally" spread on train (like sinus) I would change it to polynomial model ~~so use RNN instead~~

9. ~~Assume have too much points~~ For a low Y, we have low residuals so the model is fitting well for low Y. But for the bigger ones, we have ~~higher~~ a bigger residual ~~(absolute)~~ (absolute). I think that this is also due to the choice of linear regression. I would try to use a RNN method and compare it to this model.
~~or a polynomial~~

Or I would test a polynomial as residual gets bigger when Y grows like x compared to x^2 .

10 Logistic regression is a classification method. It assigns labels and separate dataset in two groups according to a threshold.

For a given threshold, we will obtain a matrix ~~as shown~~ -

		in reality	
		Positive	Negative
predicted	True	Positive	Negative
	Positive	TP	FP
Negative	FN	TN	

Based on this we can calculate. Sensitivity = $\frac{TP}{TP+FN}$.

that shows part of positive well predicted

and Specificity : $\frac{TN}{TN+FP}$ that shows part of negative well predicted.

For each threshold you can calculate both sensitivity and specificity. So, for each threshold you have one point of the ROC. The best case is ~~sensitivity~~ sensitivity = specificity = 1.

In the example, for specificity = 0,9 we have sensitivity = 0,35 so we have only 35% of positive predicted as positive and 90% negative as negative. In the

For specificity = 0,1, we have sensitivity = 0,97 so 97% of positive predicted as positive and 10% of negative as negative.

None of them is ~~the best~~ better, it depends on the case we always

study. If detecting positive, like detecting ebola is more important we would go for a high sensitivity. ☺

2. To calculate B_0 , we use $\hat{P}_0 = \bar{y} - \beta_1 \bar{x}$ (in 1 dimension)

so we have $\bar{y} = B_0 + \beta_1 \bar{x}$ with B_0 and β_1 coefficient of the linear regression. For any linear regression,

the line will ~~not~~ always pass by ($\bar{x} = \frac{1}{n} \sum x_i$; $\bar{y} = \frac{1}{n} \sum y_i$)

NOM BOEKHELOOA

Prénom Rami

Promo 2020

Date

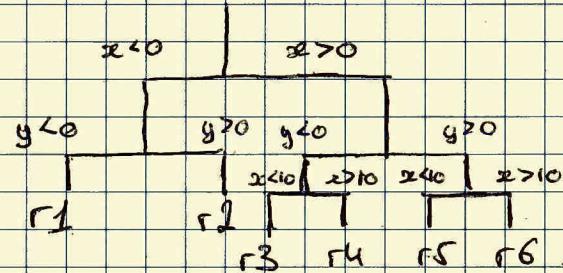
MATIÈRE Machine Learning (2/2)

11. Elbow method will help us to choose number of cluster. we will plot the graph of ~~number of cluster~~ MSE by number of cluster.

→ Thanks to the slope we will pick a number of cluster with low MSE but not that much which would result in too much clusters.

If we have

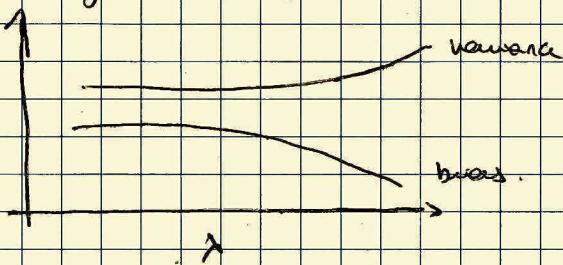
4.	x	y	label
-5	-3	r1	
-3	2	r2	
5	-2	r3	
12	-4	r4	
3	3	r5	
12	12	r6	



~~If we had a column having 3 different values "blue", "black" and "green", we could~~

12. If your bias is low, that means your model fits well the reality
~~the test set~~. If your variance is low that means
you are not flexible to test set, you are too dependent
on your train test. They are anti-correlated, if one is low,
the other is high.

In ridge method, we want to minimize $\lambda \sum (\beta_i)^2$ too



we will pick the point
where variance and bias are
closest.

