

NOM BOUKHELOUA

Prénom Rémi

Promo L2S0

Date



20150257: BOUKHELOUA Rémi

M1:

ST2ML2-DE (19/04/2019)

Amphi jaune

13

MATIÈRE ML and TM (1/2)

6. Naïve Bayes algorithm is based on probabilities and the formula following formula: $P(C|c|x) = \frac{P(x|C) * P(C)}{P(x)}$.

It is a classification method that performs well with multi classes dataset and better than others in general.

For a given input, it will compute and output a probability that the input is in a class or another.

Its main problem is that it assumes independence of features what is rarely the case in real life problem. To counter that we can remove correlated features. We can't boost or bagg it and tuning is very restricted.

7. A minimal edit distance is the minimum number of operations we have to make to transform one word to another one.

It exists three operations: add a letter, delete a letter and change a letter. The two first ones cost 1 while the second, as a combination of the two others, costs 2.

For example: AND D.
 ↓ ↓ L
 * N O D E

$$\text{Edit}(\text{"AND"}, \text{"NODE"}) = 1 \text{ delete} + 2 \text{ additions} = 3$$

11. From our corpora, we will first tokenise each text or speech.
Tokenisation first begins with ~~transforming~~ the text into words, splitting

then you remove stop words such as "or", "and" which ~~are~~ ~~meanless~~ have no impact but also non-English words.

On the set of words, we can decide of lowercasing or keeping caps. The advantage of it is that it will reduce number of columns (~~As~~ ~~the~~ "YEAR" and "year" become "year") but we will lose some sense of caps or ~~meaning~~ showing writer's sentiments.

After that, we decide to lemmatise which is change a word into a closer one ("has", "have", "had" become "have") or to stem which is keeping only root of word ~~the new base form~~

We can now construct the matrix counting how much time a word appears or only if it appears (for Binary Naive Bayes). ~~But having few words~~ We may prefer to work with frequencies to avoid having too big numbers. But words appearing a lot of time, may have ~~less~~ impact than ~~other~~ less frequent ones. That's why we introduce log function. ~~we consider~~ We will obtain a TF-IDF matrix. As it's now only composed of numbers, we can use this matrix on a machine learning method.

13. A general bagging framework generates ~~a lot~~ a lot of classifier on ~~sets~~^{some algorithm} of subsets of the train dataset. It will then do the mean of them.

For random forest, we will create a lot of decision trees based on different subsets of dataset and then, instead of doing a mean, we will combine them by addition to have a bigger tree. As it is based on nodes we just can't ~~make~~ make a mean.

Ex

14. Boosting is based on training several time increasing the importance of ~~failed~~ failed set and reducing ~~the~~ on the succeed one. To do so, we introduce a weight at each row and identical at start. Higher it will be, higher will be the importance. The classifier will then ~~focus~~ focus on failed ones.

AdaBoost is based on creating stump (1 node decision tree), with boosting method α . It will generate ~~a lot of stumps~~^{a stump at each boosting} iteration and combine them, giving each one a weight, in order to have a good model.

15. As bagging is based on generating lot of model on different part of train dataset so each little model has a different train and then merging them, it will reduce the impact of training set on the model and so ~~the variance~~ will reduce the variance.

Boosting is based on focusing on failed rows, so it will ~~adapt~~ adapt the model each ~~time~~ iterations so it will reduce variance -

Bias?

16. Latent Dirichlet Allocation (LDA) is a method to analyse text with a certain repetition of words from different topics. For example, if you have 3 topics A, B and C, you can use LDA to select words ~~that~~ with a certain repetition (20% of A, 80% of B).

It is based on Dirichlet coefficient which is a statistical variable representing the dispersion.

It could also be used to group text with same topic. Measuring Dirichlet you would be able to group ~~topic~~ text from topic A together, topic B together and so on (based on their word repetition)

8. A bag of word of a text is words contained in the text ~~and~~ ~~that~~ ~~words~~ ~~are~~ ~~repeated~~ ~~in~~ ~~a~~ ~~text~~ with binary

We could easily represent a text based on it ~~as~~ as a matrix. If the word is ~~there~~ in the bag, we put a one, if not a zero.
or nb of times it appears

It is easy to do and to use in a machine learning method.

The problem is that we do not keep the order of the words

which means that "I like math and hate physic" and

"I like physic and hate math" would have the same representation. Moreover we do it based on a dictionary so it ~~has to~~ will not adapt on ~~unknown~~ ^{new and} words. Finally,

if we do it on long text, as we put a ~~feature for each~~ column for each word, it will ~~be~~ create a high number of dimension which ~~causes~~ longer computational time.

NOM BOUKHELOVA

Prénom Rimi

Promo L2de

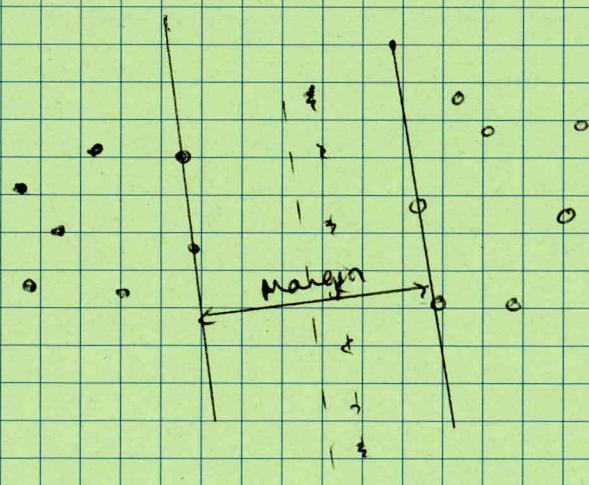
Date

MATIÈRE ML and TM(L12)

10. Embedding can be creating n-grams which are group of following words. Unlike bag of words, it will keep ~~a aspect~~ of order of words. ~~Creating~~ Creating bi-grams (2-gram) out of "I like math and hate physic" would be ("I", "like"), ("like", "math"), ("math", "and"), ("and", "hate"), ("hate", "physic") without ~~sharpening~~ inflectionalization / stemming. Based on this, we could create a Markov chain such as from one n-gram we forward probability of following one. And with a lot of dataset have a complex network from what we could generate text. or recognise author from a text (ex Shakespeare).

11. Word2vec is a way of creating ~~numerical~~ vectors from a word that will allow ~~do~~ common operation on them. As a result we could have : king - man + woman = queen.

3.



Support Vector Machine is used to create a hyperplane between classes. Without kernel trick, the hyperplane is linear. It will compute the plane ~~with the logistic~~ maximizing the margin. As we guess the following graph is 2-dimensional, I drew ~~on the~~ the hyperplane maximizing it.

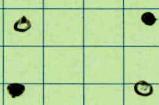
But with more dimensions and ~~at~~ hidden ones not scaled,
hyperplane could be very different!

1. F measure is frequency measures that we use when we
transform a corpora into matrix.

We could use $\frac{\text{nb of apparition}}{\text{nb of words}}$ but also a ~~more~~ interesting

one $\frac{\text{nb of apparition}}{\text{nb of apparition of most common word}}$. In this case in a lang, the
result of most common word would be 0.

2. A straight line in a two dimensional feature space
~~could only separate~~ can't even separate two classes as following.



Nb OF POINTS

So I would say that it's 1.

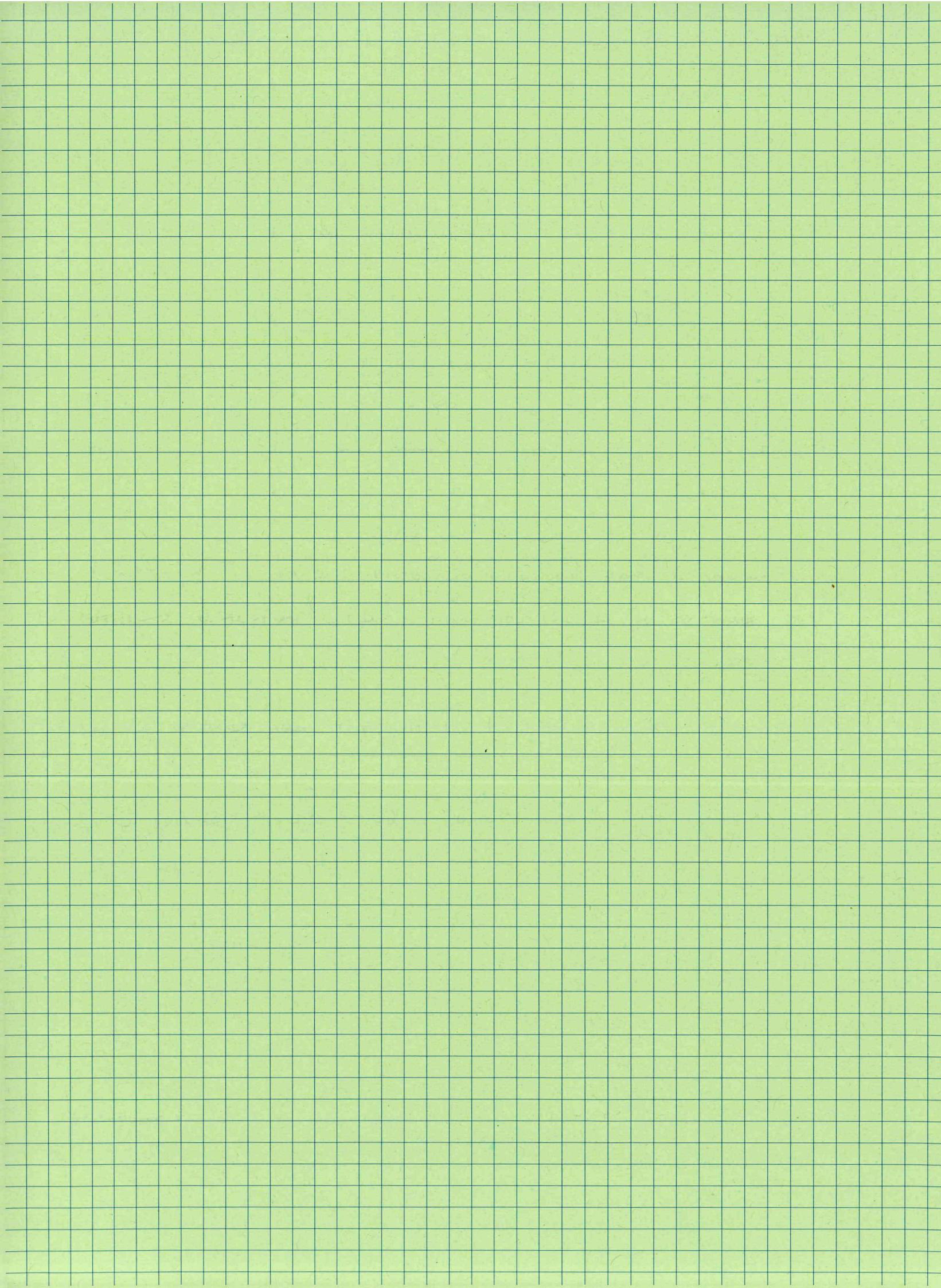
4. In a neural network, we calculate weight on training set.

It will be ~~a~~ highly dependent of it. We could use some
bagging method to generate ~~as~~ several neural networks and
at the end averaging them.

~~Average~~ We also have to care about the number of
neurons in the hidden layer that should not be over twice
number of input ~~size~~ layer neurons to reduce overfitting.

~~I'm not sure but~~ Generally, cross validation process also induces
~~over~~ train it dependence, ~~but I'm not sure~~ that could also
be a possibility even though regarding computational time, I'm
not sure it would be acceptable.

5. Instead of using the output of a threshold, we could use a sigmoid function with threshold to output only 0 or 1 depending on the threshold.





Prénom Rewmi.....

Ne rien inscrire dans ce cadre

Nom BOUKHEDDA.....

Promotion 2020.....

Groupe

Promotion M1

Advanced Machine Learning and Text Mining

ST2ML2

DE - 1h45 min

Date Horaire

Sujet proposé par :

Calculatrice autorisée : **OUI** **NON**

Documents autorisés : **OUI** **NON** **Type de documents :**

Ordinateur portable autorisé : **OUI** **NON**

Internet : **OUI** **NON**

Traducteur électronique, dictionnaire : **OUI** **NON**

Consigne :

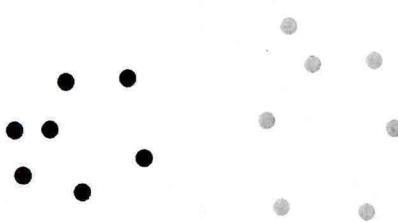
Merci de restituer uniquement : « **choisir un élément de la liste** »

Rappel :

- Tous les appareils électroniques (téléphones portables, ordinateurs, tablettes, montres connectées ...) doivent être éteints et rangés.
- Il est interdit de communiquer.
- Toute fraude ou tentative de fraude fera l'objet d'un rapport de la part du surveillant et sera sanctionnée par la note zéro, assortie d'une convocation devant le conseil de discipline. Aucune contestation ne sera possible. Tous les documents et supports utilisés frauduleusement devront être remis au surveillant.
- Aucune sortie de la salle d'examen ne sera autorisée avant la moitié de la durée de l'épreuve.

Ne rien inscrire dans ce cadre

1. What is "minimizing the empirical risk"? Why and when are we using it?
- + 2. What is VC dimension for a straight line in two dimensional feature space? Explain your answer.
- + 3. Show SVM separating hyperplane for the dataset plotted bellow and explain why (without the kernel trick)



- + 4. How to fight overfitting in a Neural Network learning?
- + 5. What difficulties with convergence of Neural Network fitting process?
- + 6. Explain Naïve Bayes algorithm
- + 7. What is minimal edit distance?
- + 8. Explain what is a bag of words model, its advantages and disadvantages
- + 9. What is F-measure?
- + 10. What is embedding? Example(s)?
- + 11. What are the ways to create a matrix representation from corpora to use as an input for machine learning models? (feature creation)
- + 12. Describe a schematic algorithm of skip-gram word2vec
- + 13. What is particular for the random forest algorithm comparing to a general bagging framework?
- + 14. What is a general concept of Adaboost algorithm
- + 15. Explain boosting/bagging in bias-variance set up
- + 16. What is the generative model behind Latent Dirichlet Allocation (LDA) ?

