

0,5

3° We could use a feature map or a kernel.

0,5

$$4^{\circ}(a) \nabla_{\vec{w}} J = \nabla_{\vec{w}} \frac{1}{4} \sum_n (\vec{w} \cdot \vec{x}_n - y_n)^2 + \frac{1}{4} \lambda \|\vec{w}\|_2^2$$

$$= \sum_n (\vec{w} \cdot \vec{x}_n - y_n) \begin{pmatrix} x_{n1} \\ \vdots \\ x_{nd} \end{pmatrix} + \lambda \vec{w}$$

0,25

(b) for the fit function we will change the  $w = \dots$  with this new formula.  
Predict world not change.

0,75

(c) I would think that LinReg<sup>2</sup> would perform better because it is updated with the same and evaluated with the same idea.

1

5° (a) Instead of doing an epoch on  $N$  points we would only do it on  $[1, \dots, \text{batch-Size}]$ .  
fit (...):

for epoch in  $[1, \dots, \text{Max-epoch}]$  do  
for  $n$  in  $[1, \dots, \text{batch-Size}]$  do.

$$\text{Thus } \vec{\nabla}_{\vec{w}} J = \vec{\nabla}_{\vec{w}} \frac{1}{2} \sum_n (\dots)$$

0,5

(b) Using mini-batch GD would help us avoid overfitting on outlier data because it would have fewer example of them.

0,5

(c) The fit time will be smaller.

④

ANNÉE : 2020

UNITÉ D'ENSEIGNEMENT : TCO

ÉPREUVE DE : Examen

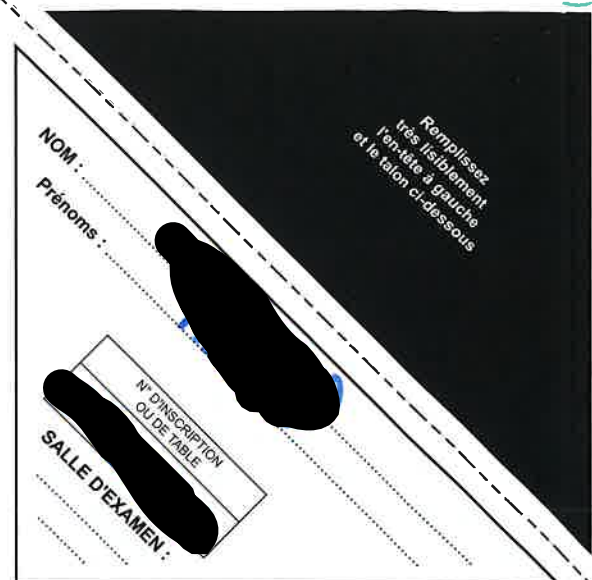
NOTE

/20

Coefficient

Note affectée  
du coefficient

19,5



Si votre composition comporte plusieurs feuilles, numérotez-les 1./n.

⑥

Analysis of experimental results:

1,5

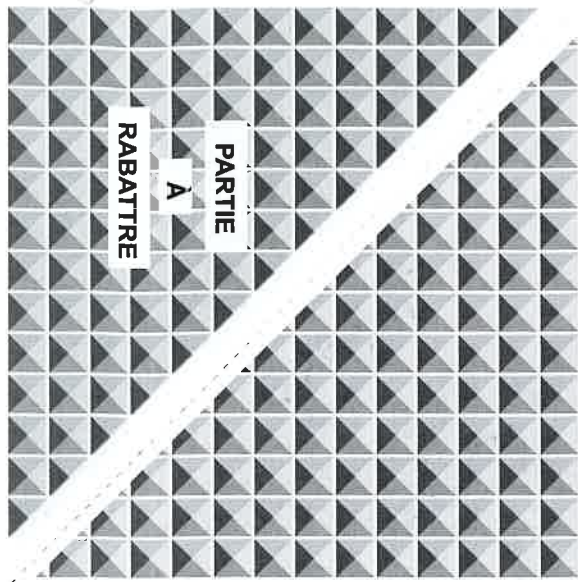
1° a- The left one is the richer cause when increase  $\lambda$  the train score decreases.  
The train score decreases because the regularization becomes stronger and it can not learn by heart the train set anymore.  
The valid score increases because the model generalizes better when not learning by heart.

b- The second one because we have the same result in validation when the model is less complexed.  
Thus we gain computation time.

0,5

2° The left one seems to be the one complexed and with regularization because it looks like we have small pikes but small due to the regularization.  
In the right one, we have a nice curve without any pikes.





3° - The train set and valid set have the same score, but it does not look like we have overfitting because the two are close and the validation does not decrease.

5 - Maybe because weird data points that can not fit even with a complex model. For the validation it's the same idea.

c - We could expect a test score of  $< 0.05$ .

4° (a) As being the separator between the two groups.

(b) As being the margin which describes which points are being the support vectors.

(c) Case one to 4 does not seem to overfit but of course 1 and 2 are underfitting as they are not ok. But 1 seems to be a bit too simple. So 2 to 4 are ok and 1 is too simple and for 5 and 6 we overfit.

(d) I would pick 2 because even if we have some errors, I would expect that the red points could exist to the bottom right and it's the one that expects it the most.

5° He seems to learn on the test set about which model and hyperparameters are the best, whereas we should only do this on the validation set.

6° We should try to gather data that is not biased or we could increase the penalty when the model becomes racist/sexist.

Yes!

6.5

## Linear Regression

1° Linear regression is used when we want to predict a scalar result using our features.

$$J = \frac{1}{2} \sum_{n=1}^n (f(x^{(n)}) - y)^2 \quad \text{with } f(x^{(n)}) = \vec{w} \cdot \vec{x}^{(n)}$$

We are doing a gradient descent as optimization algorithm.

def fit(X, y, η, w₀, maxEpoch):

$$\vec{w} = \vec{w}_0$$

for epoch in range(1, maxEpoch + 1):

$$\vec{w} = \vec{w} + \eta \sum_{n=1}^n (y - \vec{w} \cdot \vec{x}^{(n)}) \vec{x}^{(n)}$$

for n in range(1, size(X) + 1):

$$\vec{w} = \vec{w} + \eta (X[n] \cdot \vec{w} - y[n]) \vec{x}^{(n)}$$

return  $\vec{w}$

def predict(X, w):

$$\text{return } \vec{x} \cdot \vec{w}$$

$$2^\circ w_0 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$X = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 2 & 2 \end{pmatrix} \quad \text{with } X[0] = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$w_1 = w_0 + \eta (X[0] \cdot \vec{w}_0 - y[0]) X[0] = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$w_2 = w_1 + \eta (X[1] \cdot \vec{w}_1 - y[1]) X[1] = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$$

$$w_3 = w_2 + \eta (X[2] \cdot \vec{w}_2 - y[2]) X[2] = \begin{pmatrix} -3 \\ -6 \end{pmatrix}$$

$$w_4 = w_3 + \eta (X[3] \cdot \vec{w}_3 - y[3]) X[3] = \begin{pmatrix} -56 \\ -104 \end{pmatrix}$$

Online algo

1



0,25 5° If we have a plane  $w_{hh}^T$ , when we visit  $h$ , we are going to have an inner product being equal to 0.

6° For  $x^T$  test we are going to have (1) or (0) with (1) saying that  $x$  is more likely to be  $k$  and (0) saying that  $x$  is more likely to be  $k'$ .

Thus we have  $D_x \sum_{k=1}^K$  predictions in total.

Predict will look like

predict =  $\frac{1}{R \cdot R'' \cdot R'''} \sum y_{pred}$

7°

8° As hyperparameter we have:

- ✓ - the definition of  $\sigma$
- ✗ - the shape of our data set.
- ✓ - how many classes  $k$  do we have

ANNÉE : 2020

UNITÉ D'ENSEIGNEMENT : TCO

ÉPREUVE DE : Examen

NOTE	/20	Coefficient	Note affectée du coefficient

Remplissez  
l'en-tête à gauche  
et le talon ci-dessous

NOM : \_\_\_\_\_

Prénoms : \_\_\_\_\_

N° D'INSCRIPTION  
OU DE TABLE

SALLE D'EXAMEN : \_\_\_\_\_

Si votre composition comporte plusieurs feuilles,  
numérotez-les 2/2

④

#### 4 Identifying the task

1° (a) The goal is to classify mushrooms as being poisonous or not.

(b) We should do supervised classification

(c) Our label is the feature "edible/poisonous" and our features would be encoded as "one-hot".

So for each features such as "cap-shape" we would create  $n$  features where  $n$  is the number of possibilities.

(d) A decision Tree Classifier could work.

(e) We should make a metric where if a poisonous mushrooms is being predict as edible, it should be punished 10 times more than the inverse.

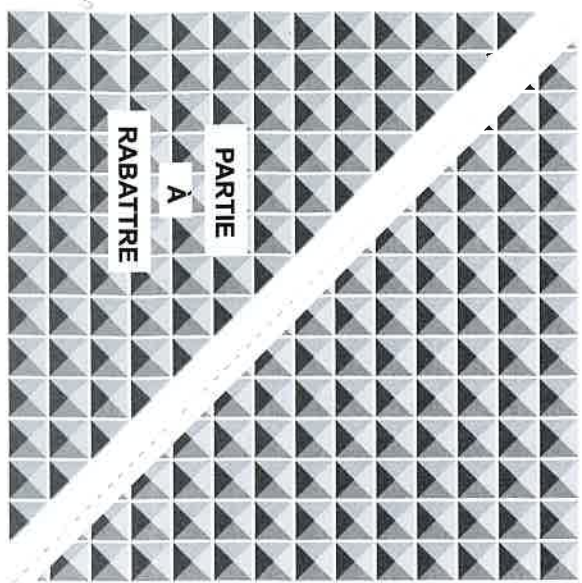
2° (a) The goal is which team will win.

(b) We should do supervised classification

(c) Our label will be the winner.

Our features are for the most numerical but for the encoding of the categorical one (= "name" and "type") we will do hot encoding (except if the name creates too many features).





(d) A perceptron could do the work.  
(e) Maybe if have the too many we could have one issue.  
And depending on which team attack  $f_i$

3° (a) Our goal is to predict temperatures  
(b) We would have supervised regression.  
(c) We could not encode it

As For the temperature, longitude and latitude they are already numerical so we have nothing to do, but for the time we could encode them as being the hour since the beginning of 1900.

(d) We could use linear regression  
(e) We could have a feature where we take the date (hour/month) and (the one not being transformed yet) and use the day and month to add a feature of the season.

4° (a) Our goal is to classify players into a role.

(b) We would have supervised classification.

(c) We could remove the name because it has nothing to do with the position and encode with week foot and preferred foot as ... foot right as ~~two~~ (hot encoding).  
For the rest as it's numerical we have nothing to do.

(d) A Multi-Class perceptron could work.

(e) We could create labels being the combination of two positions which are similar.

### 3 OVO Classification

1° (a) If  $x$  is of class 1 then we have

$$\sigma(\vec{w}_{12} \cdot \vec{x}_n) = (1, 0)^T$$

If  $x$  is of class 2 then we have

$$\sigma(\vec{w}_{12} \cdot \vec{x}_n) = (0, 1)^T$$

Thus we could have  $t_{n,12}$  as being either  $(1, 0)^T$  or  $(0, 1)^T$  but not  $(0, 0)^T$  or  $(1, 1)^T$ .

(b)  $\vec{w}_{21} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \vec{w}_{12}$  and  $t_{n,21} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} t_{n,12}$

(c)  $\sigma(\vec{w} \cdot \vec{x}_n) = \sum_k t_{n,k}$   
 $t_{n,kk'} = \begin{cases} (1, 0)^T & \text{if } k \text{ is more likely} \\ (0, 1)^T & \text{if } k' \text{ is more likely} \end{cases}$

2° We have  $D \times K$  parameters.

3°  $J = \frac{1}{N} \sum_n \left( \sum_d \left( x_n \cdot \vec{w}_d - y_n \right)^2 \right)$

2° We have  $D \times \left( \sum_{k=1}^K R_k \right)$  parameters, because we can retrieve  $w_{2,1}$  from  $w_{1,2}$ .

3°  $J = \frac{1}{N} \sum_n \left( \sum_{k=1}^K \left( \sum_{k'=k+1}^K ((x_n \cdot w_{kk'}) - t_{n,kk'})^2 \right) \right)$

4°  $\nabla_{\theta} J(\theta, x) = \nabla_{\theta} \frac{1}{N} \sum_n \left( \sum_{k=1}^K \left( \sum_{k'=k+1}^K (\tanh(x_n \cdot w_{kk'}) - t_{n,kk'})^2 \right) \right)$

$= \frac{1}{N} \sum_n \sum_{k=1}^K \sum_{k'=k+1}^K (2 (\tanh(x_n \cdot w_{kk'}) - t_{n,kk'}))$

$\frac{(4(x_n))}{(\cosh(x_n \cdot w_{kk'}))^2}$

(b)  $\frac{8}{N} \tanh(W^T X) = (T_{n,kk'})^T$