| NOTE | **19,25** | /20 | Coefficient | Note affectée du coefficient |
|------|------|------|-------------|------------------------------|
|      |      |      |             |                              |

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

---

**(6)** / **(4)**

**2.2** $X = (1,1,2,2)^T$  $y = (1,2,2,3)^T$  $w = \binom{1}{1}$

(note: I take the transpose because I feel more comfortable with those dimensions)

$x = $ ~~$X$~~ addconstant $(x)$

$w = $ ~~$w$~~ ~~matmul (matmul(X,~~

$y$-pred ~~$= matmul (X, w) = \binom{1\ 1}{2\ 2}\binom{1}{1}$~~

$x = $ addconstant $(X)$

$y$-pred $= $ matmul $(X, w) = \begin{pmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 2 \\ 1 & 2 \end{pmatrix}\binom{1}{1} = \begin{pmatrix} 2 \\ 2 \\ 3 \\ 3 \end{pmatrix}$

$diff\text{-}y = y\text{-}pred - y = \begin{pmatrix} 2 \\ 2 \\ 3 \\ 3 \end{pmatrix} - \begin{pmatrix} 1 \\ 2 \\ 2 \\ 3 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 1 \\ 0 \end{pmatrix}$

$w = w - \left(\frac{2}{N} \cdot matmul(X^T, diff\text{-}y)\right) = \binom{1}{1} - \frac{2}{4} \cdot \binom{2}{3} = \binom{0}{-0.5}$

$\begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 2 & 2 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 0 \\ 1 \\ 0 \end{pmatrix} = \binom{2}{3}$

✓

**1,5**

**2.3** We could create feature maps in order to add more features to the data so we can take that dependence into account. ✗

**0,5**

**2.4.a**

$\nabla J(w,x) = $ ~~$\nabla \nabla J ...$~~

$\nabla_w J(w,x) = \overline{2} (\overline{w} \cdot \overline{x} - y)^3 \cdot \overline{x} + \lambda \|\overline{w}\|^3$  ✗

**0,8**

**2.4.b** ~~You~~ would just change the gradient in the function with the new formula.

$w = w - \mu \nabla_w J(w,x)$

**0,25**

**2.4.c** The error would be higher using MAE ~~too~~ because it's the MSE but squared. For this exact reason, the only change with this would be that because the error is greater, the steps of the gradient descent would be higher as well, which we could control anyway with ~~$\mu$~~ $\mu$ so there is no reason one should be better than the other one.
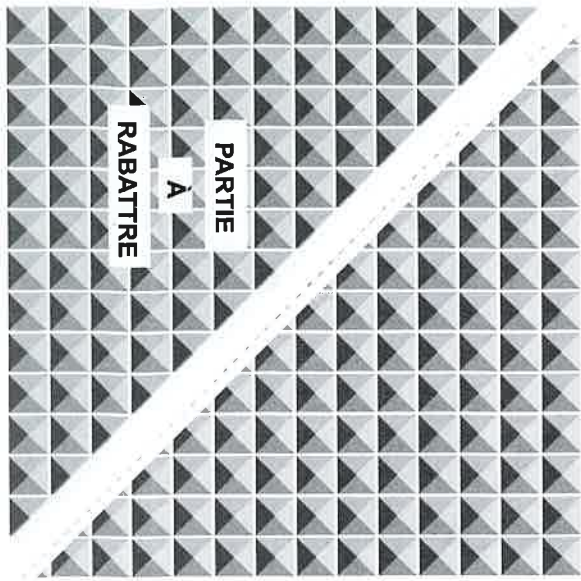
**✗**

---

**(5)**

**1.1.a** I would say that the model of the left is more expressive because by incrementing the value of $\lambda$ you get more improvement in the validation score, which means that the model was overfitting, ~~more~~ sign of ~~a~~ a complex model. The variation is explained by the fact that with little to no regularization, the model was overfitting the training set and by incrementing the regularization term, ~~yes~~ it reduced the overfitting and therefore incrementing the val. score.

**1**

**1.1.b** I would say that the first model fit better the data because it's a more complex model which ~~can~~ probably can predict more accurately future data. By adding ~~a~~ regularization, you ~~make~~ try to reduce the overfitting and allowing yourself to use such ~~a~~ complexity in the model.

**0,5**

**1.2** The model of case d is the one which is complex ~~and~~ with regularization because even if it doesn't ~~fit~~ ~~ever~~ overfit the points (because of the regularization), we can see more curves on the line which is a sign of complexity.

?

**No ✗**

**(3.a)** By incrementing the complexity of the model, you allowed the model to learn the patterns of the data, ~~turning i~~ lowering both errors. Also, in the least complex models, ~~the overfitting was k~~ there was overfitting due to ~~ilst~~ lack of expressivity. We can see in the more complex models that the overfitting disapear.

**(3.b)** Because it's normally very hard to achieve a 100% accuracy because of outlayers or even the complexity of the problem. The same happens with lowering the error to 0.

1

**(3.c)** For that model, we can expect the same error that in the validation set, assuming the validation set hasn't been overfitted by an extensive parameter-tuning process.

**(4.a)** The solid line is the ~~decision line~~ hyperplane that serves as the decision boundary to classify the examples. Examples in different sides are going to be classified ~~as~~ in ~~the~~ the opposite class.

**(4.b)** The two dash lines are the lines in which the points are called support vector because it's the line from which there is no error classifying a point in ~~it~~ its truth ground class.

1,5

**(4.c)** The top 3 seem to be doing a reasonably good job meanwhile the bottom 3 are overfitting the data due to the complexity of the lines. The 4th one could be arguable but the last two are clearly overfitting. ~~Eve~~ There could be controversy on the 3rd one as well due to the line in the left-top corner

**(4.d)** ~~I would choose the~~ that shows that the model it's creating a circle around the data. However, it's not enough information to conclude there is overfitting.

---

0,25 ✓

**(4.d)** I would pick the econd one due to it's simplicity because it makes understandable the process of predicting, while having a good score at the same time.

0,5

**(5)** The test set should only be used after all the work because otherwise it could be the case that Benji is overfitting the data to the test set by making many predictions on it.

~ 0,25

**(1.6)** ~~We could~~ Seems at first sight that the problem could be the underrepresentation of women and people from different ethnics in the dataset so we could test this hypothesis and rebalance the dataset according to it.

2

**(2.1)** Linear regression is used for predicting continuous data, that's why it's called regression.

Cost function:

$$J(x) = \frac{1}{N} \sum_{n=1}^{N} (x_n w - y_n)^2 \rightarrow \text{Assuming the following dimension: } (X_{N,D}) (W_{D,1}) (Y_{N,1})$$

The optimization algorithm is gradient descent which subtract the gradient of the cost function respect to $w$ to $w$.

```
def fit (n_iters, learning_rate, X, y
    w = initialize_w();  X = add_constant(x) # used to assume the bias
    for i=0... n_iters :
        y_pred = matmul (X, w)
        w = w - (learning_rate · 2/N · matmul(X^T, y_pred - y)
    return w

def predict (X, w)
    X = add_constant(X)
    return matmul (X, w)
```

Notez avec exactitude votre numéro de table ou d'inscription.
Il est interdit au candidat de signer sa copie ou d'y mettre un signe quelconque pouvant identifier l'auteur de la copie.

Notez avec exactitude votre numéro de table ou d'inscription.
Il est interdit au candidat de signer sa copie ou d'y mettre un signe quelconque pouvant identifier l'auteur de la copie.

④ Ⓐ The goal is to predict the variable 'edible/poisonous'.

Ⓑ It's a supervised task ~~vare~~ and the sub-category is classification.

Ⓒ We have a matrix which rows are each mushroom and in the columns are the features (possibly including the predicting one or in other vector y)

Ⓓ ~~A naive bayes~~ A gaussian naive bayes ~~seems~~ good for this task because it seems simple at first sight, also it would help us understand the model and not use a black box so we can learn more about mushrooms :)

Ⓔ We should be worried about the imbalance of the classes because of the fact that there are ~~~~ more edible mushrooms. Also, we care more about predicting correct the poisonous ones since otherwise we could die ;) We should use a performance metric which takes into account more the false negatives (wrongly predict 'edible') than the false positives.

② Ⓐ The goal is predicting the winning team giving two pokemon teams.

Ⓑ It's a supervised task of classification.

Ⓒ We could have a 3dimensional structure with dimensions $N^{ex}_x \times N^{pok}_x \times D$ in which $N^{ex}$ is the number of fights, $N^{pok}$ is the number of pokemon per fight and $D$ is the attributes per pokemon. We could unroll the last two so we have a 2d matrix which is easier to train a model with, having a $N^{ex}_x \times (N^{pok} \cdot D)$ dimensional matrix.

Ⓓ This is a very difficult problem since the features are correlated and the same pokemon can appear in different indices in each fight. I would use a deep neural network which will take a lot of time to train but since it should take all possibilities into account, is a reasonable price.

Ⓔ We should particularly worry about the complexity of the problem to try to reduce it as much as we could so we are able to predict accurately the results. Also, ~~we~~ keep in mind the overfitting ~~so we can't~~ applying regularization or even dropout to the model.
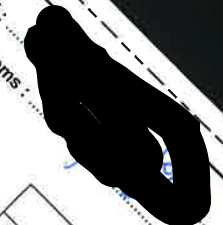
---

② .5 Ⓐ This means that ~~~~ in each epoch, we would ~~then~~ update the weight as many times as mini-batches exist.

```
def fit (...):
    for mini-batches = ✗ do batches (X)
    for i=0...n-iter:
        for mini-batch in mini batches:
            w = w - N ∇ J(mini-batch)
    return w
```
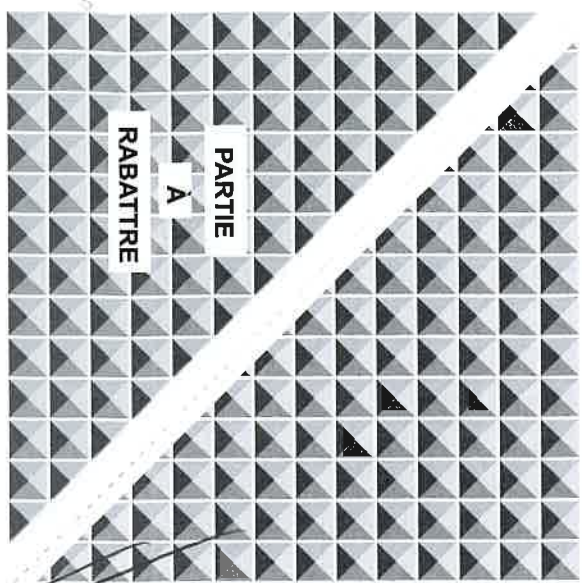
2.5.b First, because we update the weights more frequently, it's faster. ~~~~ And second, we mitigate the noise issue by creating mini-batches which are not representative of all the data since it is the case as well in the full-batch (because of the noise). This will give least-represented ~~points~~ more importance in their mini batches.

2.5.c As I said in the previous question it is faster because you update the weights #mini-batches times per epoch istead of only 1. Even if it sometimes it doesn't take a good direction, it converges faster usually.

PARTIE A
RABATTRE

**3.1.a** $\frac{\partial}{\partial}\left(w_{j2}^T \bar{x}\right) = \left(t_{n1} - t_{n2}\right) =$

When $k=1$: $(t_{n1} - t_{n2}) = 1 - 0 = 1$

When $k=2$: $(t_{n1} - t_{n2}) = 0 - 1 = -1$

We see that we have $1$ and $-1$ as labels so we can create the variable $t_{n,12}$ to be this new encoding. $t_{n,12} = 1$ if $k_n = 1$ else $t_{n,12} = -1$.

0,25

**3.1.b** $\bar{w}_{21}$ is the weight of the ~~the bias~~ the first feature for predicting the likelihood of being class 1.
$\bar{w}_{21}$ is the weight of the bias term to predict the second class. They don't seem to be related.

If $t_{n,21}$ is $1$ then $t_{n,12} = -1$ so: $t_{n,ij} = -t_{n,ji}$

0,25

**3.1.c** $\frac{\partial}{\partial}(w^T x)$ $\frac{}{}$ $k_n =$

$t_{n,kk'}$

$$t_{n,kk'} = \begin{cases} 1 & \text{if } k_{true}^{(k)} = 1 \\ -1 & \text{if } k_{true}^{(k')} = 1 \\ 0 & \text{otherwise} \end{cases}$$

0,5 / 0,75

**3.2** There are the ~~features~~ of $X$ which is $D$ and the one-hot vectors of $Y$ which are $\frac{K \cdot (K-1)}{2}$ so:

$D + \frac{K(K-1)}{2}$

**3.3** $J(\Theta, X) = \sum_{n=1}^{N} \left(w^T \cdot x_n\right) - t_n\right)^2 \cdot t_n$

**3.4.a**

**3.5** Because $t_{n,kk'}$ would be $0$, there won't be any updates in ~~the to~~ $w$.

**3.6** ~~We have~~ The prediction is $k^A$ $k_{pred}^{tot} = w^T X$. We have $\frac{K(K-1)}{2}$ predictions and we take the max as the predicted class.

```
def predict (...)
    max (matmul (w^T, X))

    return preds
```

0,25

**3.7** I would choose the accuracy as performance metric.

$$Acc = \frac{TP + TN}{TP + TN + FP + FN} \cdot 100$$

In case of imbalanced classes we may choose the balanced accuracy instead.

**3.8** The hyper-parameters are: train ~~and~~ validation and test size of our data, the learning rate, the way we initilize the parameters, the stopping criterion, the error/loss function chosen, ~~the~~ and the activation function and the way to choose the prediction (max, softmax...).

0,5

0,25

| NOTE | /20 | Coefficient | Note affectée du coefficient |
|------|-----|-------------|------------------------------|
|      |     |             |                              |

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

(4) (3a) The goal is to predict future temperature given historical data.

(3b) It's a ~~regression~~ supervised regression task of time-series.

(3c) We have a vector of all the records of temperature in chronological order.

(3d) We could use a recurrent neural network since it would take into account previous temperature as well as the weights of the network to predict the future temperature.

(3e) We should worry about the fact that we lack a lot of variables to predict future temperature such as $CO_2$ emissions so ~~the~~ it would be very complicated to create a reliable model. Also, the variance between seasons could undercover the ~~temperature~~ increment of the temperature over the years due to climate change. I would propose a feature map taking the average of the whole year ~~and then creating a simpler model that given the avg temperature of and~~ the averaged mean of the last 4 months, giving a greater impact on the closest months and ~~a boolean with the curren~~ the substraction of the last 2 months temperature so we know if its increasing (positive) or decreasing (negative).

[+0.25]
1,75
TB

**4a)** The goal is to predict the position played of a player given attributes of the player.

**4b)** It's a supervised multilabel multilabel classification task.

**4c)** We have a matrix in which the rows are each player and the columns are their attributes. ~~offsite~~

**4c)** The non-standarized labels could be solved by hiring experts on the field to unify ver similar positions due to the limited set of classes and the (seemingly) ease of the task. Another approach would be applying a clustering algorithm to the players and then unify those positions that are very close to each other (this could turn into a weird result though). ~~An~~ We could also mix both approaches and suggest a lot of unification to the experts. Another problem to tackle would be the feature selection since probably not all the features are useful and the rating could also be biased due to the human-made task.

**4d)** A seemingly good algorithm for this task would be a decision tree or random forest classifier since they create clear set of rules to predict the ~~player~~ position, which could be very useful in this problem because certain positions needs a certain level of skills (a center player doesn't need to be as tall as a defense player).

Notez avec exactitude votre numéro de table ou d'inscription.
Il est interdit au candidat de signer sa copie ou d'y mettre un signe quelconque pouvant identifier l'auteur de la copie.

Notez avec exactitude votre numéro de table ou d'inscription.
Il est interdit au candidat de signer sa copie ou d'y mettre un signe quelconque pouvant identifier l'auteur de la copie.