

7.75

2

1.3

0.5

1

0.25

(N, D+1)

Assuming we have: X_{exp} , a matrix, where every row is the vector starting with a 1 and then the D features

predict: $\text{def predict}(X_{exp}, \theta):$
 $y_{pred} = \theta^T X$
 return $\theta^T X_{exp}$ using $\text{np.matmul}(X.T, \theta)$

2. update: $X = \begin{pmatrix} 1 & 1 & 2 & 2 \\ 1 & 2 & 2 & 3 \end{pmatrix}$ $w_0 = (1, 1)^T$ $\eta = 1$
 $N = 4$

We extend X : $X' = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 2 & 2 \end{bmatrix}$
 $J(w, X) = \frac{1}{N} \sum_{i=1}^N (\vec{w} \cdot \vec{x}_i - y_i)^2 = \vec{w}^T X X^T \vec{w} - 2 \vec{w}^T X \vec{y} + \vec{y}^T X X^T \vec{y}$
 $\nabla J(w, X) = 2 \cdot \frac{1}{N} \sum_{i=1}^N (\vec{w} \cdot \vec{x}_i - y_i) \cdot \vec{x}_i$

UPDATE RULE: $w_{t+1} = w_t - \eta \cdot \nabla J(w_t, X) = w_t - \frac{2}{N} \sum_{i=1}^N (\vec{w} \cdot \vec{x}_i - y_i) \cdot \vec{x}_i$

① $i=1$ $f(x_1) = [1 \ 1] \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 2 \rightarrow f(x_1) - y_1 = 1$
 $i=2$ $f(x_2) = [1 \ 1] \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 2 \rightarrow f(x_2) - y_2 = 0$
 $i=3$ $f(x_3) = [1 \ 1] \cdot \begin{bmatrix} 1 \\ 2 \end{bmatrix} = 2 \rightarrow f(x_3) - y_3 = 1$
 $i=4$ $f(x_4) = [1 \ 1] \cdot \begin{bmatrix} 1 \\ 2 \end{bmatrix} = 2 \rightarrow f(x_4) - y_4 = 0$
 $w_1 = w_0 - \frac{1}{2} \cdot \begin{bmatrix} 1 \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} + 0 \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} + 1 \cdot \begin{bmatrix} 1 \\ 2 \end{bmatrix} + 0 \cdot \begin{bmatrix} 1 \\ 2 \end{bmatrix} \end{bmatrix}$
 $= w_0 - \frac{1}{2} \cdot \begin{bmatrix} 2 \\ 3 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 1 \\ 1.5 \end{bmatrix} = \begin{bmatrix} 0 \\ -0.5 \end{bmatrix}$

3. relation not linear. we could apply a feature map to the input, ~~increase~~ the dimensionality of the data.

4. a) $J(w, \vec{x}_{train}) = \frac{1}{N} \sum_{i=1}^N (\vec{w} \cdot \vec{x}_i - y_i)^2 + \frac{\lambda}{2} \|\vec{w}\|^2$ Gradient descent. Larg^2
 $\nabla J = \frac{2}{N} \sum_{i=1}^N (\vec{w} \cdot \vec{x}_i - y_i) \cdot \vec{x}_i + \lambda \cdot \vec{w}$

$\vec{w}_{t+1} = \vec{w}_t - \eta \cdot \left(\frac{2}{N} \sum_{i=1}^N (\vec{w}_t \cdot \vec{x}_i - y_i) \cdot \vec{x}_i + \lambda \cdot \vec{w}_t \right)$
 $= \vec{w}_t - \eta \cdot \left(\frac{2}{N} \sum_{i=1}^N (\vec{w}_t \cdot \vec{x}_i - y_i) \cdot \vec{x}_i \right) - \lambda \cdot \vec{w}_t$

b) There are no changes in the predict function, the function to predict is the same, we just change the ~~predict~~ cost function. We will change the fit function, we will recalculate J as the expression above and the update of theta, according to above too.

ANNÉE :

UNITÉ D'ENSEIGNEMENT :

ÉPREUVE DE :

NOTE	Coefficient	Note affectée du coefficient
18,25 / 20		

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

NOM: [REDACTED]

Prénoms: [REDACTED]

NO D'INSCRIPTION OU DE TABLE

SALE D'EXAMEN:

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

(4)

0.5

X

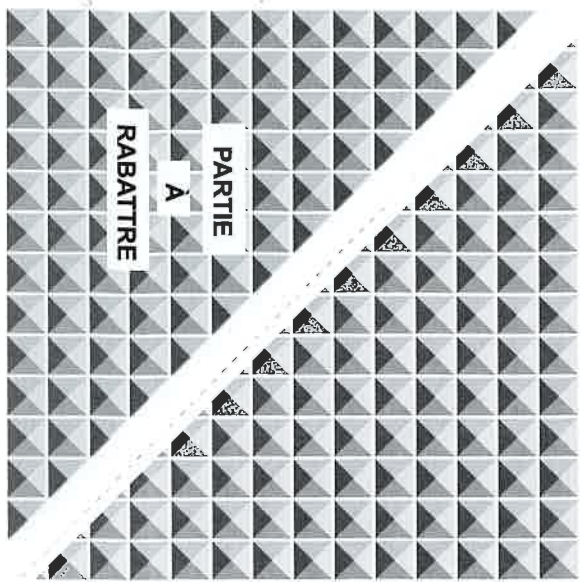
prick clearly!

ANALYSIS OF EXPERIMENTAL RESULTS

1. regularization

a) To assess expressivity, we should take a look at the accuracy, in this case, the validation score. For the correct regularization, both models have similar expressivity, but in general, the right one has higher expressivity, for smaller λ has better accuracy. The variations in a) are explained thanks to regularization. For small λ , there's not much regularization, our model is probably too complex for the data and is overfitting, that's why we have such a big gap between train and validation scores. However, for big λ , when regularization is applied, we reduce the complexity of our model, we reduce overfitting and improve generalization, validation score and the gap between train and validation.

b) If we take into account only train data we would say the model on the left fits it better (it overfits and has better train score). However, taking both training and validation sets into account, will make us conclude that, ~~even~~ with regularization, there's still a difference between train and validation score, ~~so~~ we might still be overfitting our data even for λ , we can see a small improvement with λ validation with big regularization! Consequently, we are probably fitting quite good the data and not generalizing that well.



2. Figure 2. In this case a complex model might probably be a high-degree polynomial while a simpler one might have lower degree. In case c the line has more noticeable variations, different slopes, that makes us think it was generated with a more complex function (even if regularization makes it smoother, to not overfit). Whereas in d the function is more smooth, simpler, removing of a low-degree polynomial. So, c: complex with regularization / d: less complex no regularization.

3. Figure 3. In these case, for small model complexity,

our validation error is bigger than the train, but as complexity increases the error decreases. This means that, at first, our model wasn't expressive enough, we were underfitting the data. When our model gets complex enough both errors are similar, there is no underfitting because both errors remain similar. b) The train error is not 0 because we are not perfectly fitting the data, otherwise we would have overfitting. Validation = 0 would mean a perfect model that can generalize perfectly, not very likely.

4. Figure 4:

a) The solid line is the border, the line that separates both sides (or all sides) of the plane. b) The dashed lines indicates where the margin starts (distance from the line to the border), indicating that all the points inside will generate support vectors. c) Numbering from 1 to 6, from left-right and top-bottom $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$, 5 and 6 are clear cases of overfitting the borders are very complex and we can have more than one. We will probably have good performance in train but it won't generalize well. Number 4 and 3 show really big margins, what might be a sign for overfitting as well. Consequently, 2 looks a good model simple margin and probably good generalization. Model 1 will also generalize well, even if the performance won't be as good as in 2. d) I will take the risk of model 2 probably, the size of margin and the amount of points inside seems reasonable for both classes.

c) Since train error and validation error are similar (for complex models), we can expect a test error similar to both other values, there's no overfitting. However, for less complex, test error will be similar to validation error.

5. Benji. What he calls test set might start becoming a second train set, since he is using it repeatedly to optimize hyper-parameters. We could do something better, like a cross-validation, even with this test set, but just to optimize the hyper parameters, making sure he's not overfitting, that there's no data leakage in his model.

6. Reliability people

We could add some prior of our model lets us. We could try to use the maximum a posteriori, add a prior probability estimation for the parameters (about the split in cases/suxs).

2. LINEAR REGRESSION

1. Linear regression are model to perform a regression task, we want to predict continuous labels. In a regression task we have: • Input features $X = \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix}$, N samples, N in total. • Their corresponding output, or ground truth (y) . • Model, with some parameters $\theta = (\theta_0, \theta_1, b)$ that tries to predict the label \hat{y} for each data point \vec{x} inputted. • To do the prediction we use a function $f(\vec{x})$ that can be either linear $(f(\vec{x}) = \vec{\theta} \cdot \vec{x} + b)$ or polynomial $(f(\vec{x}) = \sum \theta_i \cdot x^i)$. In order to have a good prediction, our model needs to define its parameters, so the predicted labels are as similar as possible to ground truths. To control the similarity we have a cost function, least squares for example $(LSE: J(\theta, X) = \frac{1}{2} \sum_{i=1}^N (f(\vec{x}_i) - y_i)^2)$. Ideally, we need to minimize this cost function and to do so we want our parameters θ to be:

$\theta^* = \text{argmin}_{\theta} \{J(\theta, X)\}$. To minimize the function we use the Gradient Descent algorithm. We start with an initialization of $\vec{\theta} = \vec{\theta}_0$ and we set a learning rate η , small enough to assure convergence. We will iterate through the dataset several times (called epochs) and after each step we will update the parameters with the following rule: $\vec{\theta}_{t+1} = \vec{\theta}_t - \eta \cdot \vec{\nabla}_{\theta} J$ where $\vec{\nabla}_{\theta} J$ is the gradient of J with respect to θ . (Instead of the whole dataset, we can update every often every batch, a part of the dataset).

fit function: (taking \vec{x} and y)
 $\theta = \text{random.normal}(x.\text{shape}[-1])$ # we initialize theta with small random values with $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$
 $\eta = 0.1$
 $N_{\text{epochs}} = 10$ # 10 examples
 $X_{\text{test}} = \text{np.stack}(\text{range}(1, X))$ # we extend every sample x_i to have 1 on top, to simplify.
 for epoch in N_{epochs} :
 $y_{\text{pred}} = \text{predict}(X, \theta)$
 $J = \frac{1}{2} \sum (y_{\text{pred}} - y)^2$ # all this using matrix multiplications if needed, and transpose to get correct dimensions
 $\theta = \theta - \eta \cdot \nabla_{\theta} J$
 return θ

6. $\hat{y} = \text{sign} [\sigma(\vec{w}_{k1} \cdot \vec{x}) - t_{nkk}]$
 we will have $K(K-1)$ predictions / hypotheses, we need to be updated.
 OK
 1
 0.5
 3
 5. Nothing will happen to that plane, no parameters will be updated.

clear $K \geq 2$ we have the same idea, classes with $K \geq 2$ will be linear class

number 2:

$$\sigma(\vec{w}_{k1} \cdot \vec{x}_n) - (t_{n1} - \sum_{k=2}^K t_{nkk})$$

so $t_{nkk} = \begin{bmatrix} +\delta_{k1} \\ -\delta_{k2} \\ \vdots \\ -\delta_{kK} \end{bmatrix} \Rightarrow \text{PP} \text{ is } \dots$

2. $K \geq 2$ case, parameters in the model?

$$|\Theta| = \frac{K(K-1)}{2}$$

we have K classes, each one has a plane for all the others, but they're reciprocal.

K classes, each one a plane for the rest $(K-1)$

0.75 $|\Theta|$ is $K(K-1)$, but we know that $w_{ij} = -w_{ji}$, we only have half of the parameters:

$$|\Theta| = \frac{K(K-1)}{2} \quad \times D!$$

3. $J(\Theta, X) = \sum_{i=1}^N \sum_{k=1}^K \sum_{k'=1}^K [\sigma(\vec{w}_{kk'} \cdot \vec{x}_i) - t_{nikk'}]^2$

4. $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$ $x: [-\infty, +\infty]$ $y: [-1, 1]$ $\frac{\partial}{\partial x} \tanh(x) = \frac{4u'(x)}{\cosh^2(x)}$ $\cosh(x) = \frac{e^x + e^{-x}}{2}$ $\geq 1, \forall x$

$\sigma(\cdot) = \tanh(\cdot)$

a) Derivate GD for this $J(\Theta, X)$

we call $(\vec{w}_{kk'} \cdot \vec{x}_n) = u(x)$

$$\nabla_{\vec{w}} J = \nabla_{\vec{w}} \sum_{i=1}^N \sum_{k=1}^K \sum_{k'=1}^K [\sigma(u) - t_{nikk'}]^2 = \checkmark \text{ ok}$$

$$= \frac{1}{2} \sum_{i=1}^N \sum_{k=1}^K \sum_{k'=1}^K [\sigma(u) - t_{nikk'}] \cdot \frac{4 \cdot u'(\cdot)}{\cosh^2(u)} = 4 \sum_{i=1}^N \sum_{k=1}^K \sum_{k'=1}^K [\sigma(u) - t_{nikk'}] \cdot \frac{u'}{\cosh^2(u)}$$

$$\nabla_{\vec{w}} J = 4 \sum_{i=1}^N \sum_{k=1}^K \sum_{k'=1}^K [\sigma(\vec{w}_{kk'} \cdot \vec{x}_i) - t_{nikk'}] \cdot \frac{\vec{x}_i}{\cosh^2(\vec{w}_{kk'} \cdot \vec{x}_i)}$$

$$\vec{w}_{kk'} = \vec{w}_{kk'} - \eta \cdot 4 \sum_{i=1}^N \sum_{k=1}^K \sum_{k'=1}^K [\sigma(\vec{w}_{kk'} \cdot \vec{x}_i) - t_{nikk'}] \cdot \frac{\vec{x}_i}{\cosh^2(\vec{w}_{kk'} \cdot \vec{x}_i)}$$

$$\vec{w}_{kk'} = \vec{w}_{kk'} - \eta \cdot 4 \cdot \dots$$

$$\vec{w}_{kk'} = \vec{w}_{kk'} - \eta \cdot 4 \cdot \sum_{i=1}^N \sum_{k=1}^K \sum_{k'=1}^K \sigma(\vec{w}_{kk'} \cdot \vec{x}_i) \cdot \vec{x}_i$$

$$= \vec{w}_{kk'} - \eta \cdot 4 \cdot \sigma(\vec{w}_{kk'} \cdot \vec{x}_i) \cdot \vec{x}_i$$

Noter 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.

ANNÉE :

UNITÉ D'ENSEIGNEMENT :

ÉPREUVE DE :

NOTE	/20	Coefficient	Note affectée du coefficient

NOM :
Prénoms :

N° D'INSCRIPTION OU DE TABLE
SALLE D'EXAMEN :

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

0.75
✓
thanks!

1

✓

② LINEAR REGRESSION (II)

4. (c) In our second model our model is minimizing a cost function that minimizes the MSE, as a result we can expect that this model LinReg² will perform better than LinReg that was trained to minimize squared errors. (As long as the regularization is well defined). The higher the degree of the polynomial, for a measure of over, the higher the penalty for big errors, so a LinReg² will be better trained for big errors than LinReg, probably the error will be smaller.

5. Mini-batches

a) Now we will iterate over each epoch but also over batches.

(as before)
for each m epochs:
batches = split-in-batches (X, y)

for batch m batches:

$y_{\text{pred}} = \text{predict}(X, \theta)$

$J = \text{cost}$ according to the cost function decided (for each mini-batch)

$\theta = \theta - \eta \cdot \text{grad} J$

Grad $J = \nabla_{\theta} J$

return θ .

we update θ after each mini-batch

Mathematically it means that now the cost function will be:

$$J(\Theta, X) = \frac{1}{M} \sum_{n=1}^M (f(\vec{x}_n) - y_n)^2$$

where M is the number of samples in the minibatch.

after we would update

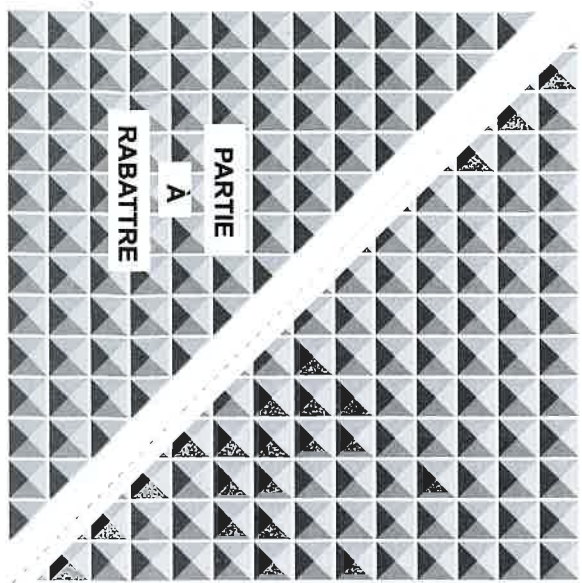
$$\theta = \theta - \eta \cdot \nabla_{\theta} J$$

we would repeat this

So after each iteration for a mini-batch, we would continue with the following one

Noter 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.

until we reach the desired number of epochs



b) for big datasets, as the ones used in Deep Learning, the samples might have great variance, the cost function can be very complex, have different minimums, so the gradient won't probably follow a clear descent direction, but on average above all the different ones, making convergence slower and probably leading to not so good local minimum.

c) mini-batches can be also useful when we have limited resources; if we can't handle big amounts of data, mini-batches allow us to process the data in parts, making the task less resource-demanding.

③ 0-10 classification
n, k=2 a) n is class 1?

④ 1) Mushroom poisons: a) Given features of a mushroom, say if it is edible or poisonous. b) Supervised task, classification (edible/poisonous). c) data are features, usually categorical that we will treat (poisonous or not) and some other numerical ones (height, width), each mushroom will be an array with a component for each feature. d) SVM for classification seems fine. e) we should take into account that predicting true = poisonous a false positive doesn't have the same consequences as a false negative. We will adapt our performance metric so our classifier has more false positive (we won't eat mushrooms that are actually edible) than false negative (we don't want to eat something that is poisonous).

2) Pokémon: a) Our task is ranking Pokémon, to say which one will win. We need to create a strength metric given its features. b) we should do a supervised regression. c) data will be an array of features for each Pokémon. d) we could try linear regression, but it seems a quite complicated task, so neural networks might work better. e) there's nothing special to worry about, maybe the fact that some Pokémon won't appear as much as others in the dataset, we might have unbalanced data, but it will probably not have any huge impact.

3. world temperatures: a) The goal is to predict future temperatures, give a prediction for temperature based on past records. b) It ^{will be} an unsupervised task, since we don't know which temperature will have in the future, but we can train our model to predict temperatures from the past or present, and then predict in the future (that part would be supervised). c) Data is an array for each city containing the numerical features. e) we could ~~include~~ a periodic function to manage this period of one year. we could generate new dimensions for each season, according to the time of recording and maybe latitude for the hemisphere, so each season is in a different component.

4. Footballers positions

a) we want to predict the label position given a football player's characteristics. b) It is a supervised task of classification. c) The data is an array for each player containing numerical and categorical data (that will be encoded), each one in a component. d) As a complex classification task SVM can work well, or even Neural Networks. e) To do this dimensional reduction, we need to cluster these labels, create groups, for this we can use an unsupervised method like K-means. This way, instead of 27 different labels, we can group our players in the desired number of different labels (9-12). afterwards, we could use and predict these new labels.

③ 0-10 classification.

$$1. k=2 \quad \sigma(\vec{w}_{1,2} \cdot \vec{x}_n) - (t_{n1} - t_{n2})$$

$$a) n=1? \quad \frac{\sigma(\vec{w}^*)}{K} - \left(\begin{bmatrix} 1 \\ 0 \end{bmatrix} - \begin{bmatrix} 0 \\ 0 \end{bmatrix} \right) = K - \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$n=2? \quad \sigma(\vec{w}^*) - \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix} - \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right) = K + \begin{bmatrix} 0 \\ 1 \end{bmatrix} = K - \begin{bmatrix} 0 \\ -1 \end{bmatrix}$$

So we need \vec{v} to be: $\vec{v} = \begin{bmatrix} 1 & \text{if } n=1 \\ -1 & \text{if } n=2 \end{bmatrix}$

b) $\vec{w}_{1,2}$ is $(-1) \cdot \vec{w}_{2,1}$, is the same plane, but pointing in opposite directions to have $\vec{w}_{1,2} \cdot \vec{x} > 0$ if x is class one and $\vec{w}_{2,1} \cdot \vec{x} > 0$ if x is class 2 (< 0 if class two)

$$t_{n,12} \text{ is also } (-1) \cdot t_{n,21}$$