

YE JIANGHENG

2016.8.22

4) Binary Classification Task

it is a supervised task which domain X , a space which all data set belong it, given $x \in X$, it is called sample and we do prediction on it. LABOL set $y = \{1, -1\}$, it's a discrete set (in this case, this set contains only 2 elements). Given training set $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$, we can learn MAP , or H , (Hypothesis class), which define all model possible.

a model
We choose $h \in H$ and we can use loss function to measure the performance.

loss function $l(h, (x_i, y_i)) = Hx \rightarrow R$ where $Z = X \times Y$. The

final goal is to minimize the generalization error

where the samples used for the test are not used in training error.

(Also in this case the sample are generated by a probability distribution function D , which is unknown).

In binary classification problem, we can use 0-1 Loss and if it is defined as:

$$\text{ERROR}_{(h, (x_i, y_i))} = \begin{cases} 1 & \text{if } (h(x_i) \neq y_i) \\ 0 & \text{else} \end{cases}$$

$$\text{Loss} = \frac{1}{m} \sum_{i=1}^m \text{error}_i$$

$m = \text{tot number of sample of test-set}$

H is to be chosen! \circlearrowleft

4.2)

the training error $L_s(h)$: is the error measured by loss function on the group of sample which is all elements of training set

S

$L_D(h)$ is the result of loss function on samples (generate by the same pdf).
 D (random used for testing error), which are not used in training of model. The main goal is measure the performance of the choose model $h \in H$ on DATA that it have seen before.

Proof? $\ominus 1$

③ ERM procedure take the model $h^* \in H$ such that: $\min_{h \in H} L_S(h) = L_S(h^*)$
 When the training data is large enough, we can define the hypothesis class H more precisely, so taking the model $h^* = \text{ERM model}$ the $L_D(h^*) \approx L_S(h^*)$

$\ominus 2$

④ Given a training set S , we define S is $\frac{\epsilon}{2}$ representative if

$$\forall h \in H, |L_S(h) - L_D(h)| < \frac{\epsilon}{2}.$$

taking its ERM model:

$$|L_S(h_S) - L_D(h_S)| < \frac{\epsilon}{2} \quad \leftarrow \text{by definition } \frac{\epsilon}{2} \text{ representative}$$

$$L_S(h_S) \leq L_D(h_S) + \frac{\epsilon}{2}$$

$\ominus 3$

$$L_S(h_S) \leq \min_{h \in H} L_D(h) + \frac{\epsilon}{2} \quad \leftarrow \text{because } \min_{h \in H} L_D(h) < L_D(h_S)$$

$$L_S(h_S) \leq \min_{h \in H} L_D(h) + \frac{\epsilon}{2} + \frac{\epsilon}{2} \quad \leftarrow \text{Adding } \frac{\epsilon}{2}$$

ϵ

Exercise 3

Linear Regression is a supervised task with domain X , label $Y = \mathbb{R}$. Consider

h :

$$L^D = \{ \text{all } \langle w, x \rangle + b \} \quad \text{Linear affine space, the hypothesis class } H$$

of Linear regression is $L^D \circ \phi(z)$ where $\phi(z) = z$.

h :

$$\text{loss function: } l(h_i(x_i, y_i)) = (h_i(x_i) - y_i)^2 \quad + \text{ is the squared root function}$$

✓ ✓

We use also Tikhonov regularization function for select a model with less complexity. (Tikhonov measure the complexity and balance the error of $l_s(h)$ and complexity of h) ✓

In this case, the model h is defined by the chose of w and b

This regression problem with Tikhonov regularization is called also l_2 -regularization:

$$\underset{w}{\operatorname{argmin}} (l_s(h) + \lambda \|w\|^2) \quad \text{where } \|w\|^2 = \sum_i^m w_i^2$$

2)

We can write it also in matrix:

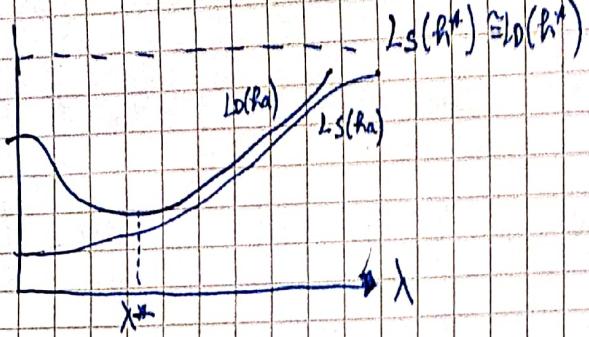
$$\underset{w}{\operatorname{argmin}} ((Y - Xw)^T (Y - Xw)) \quad Y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \end{bmatrix} \quad X = \begin{bmatrix} \cdot & x_1 & \cdot \\ \cdot & x_2 & \cdot \\ \cdot & x_3 & \cdot \end{bmatrix} \quad w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

↓ Compute Gradient and compare to 0: $2\lambda w - 2X^T(Y - Xw) = 0$

and solving we can get the best w

$$w = [I + X^T X]^{-1} X^T Y \quad \checkmark$$

- 3) λ regulates the trade-off between ERM and the complexity of the model. We can plot the error $l_s(h_\lambda)$, $l_D(h_\lambda)$ where h_λ is model ~~the~~ given by algorithm by increase or decrease λ



But how do you choose λ ?

(-2)

if $\lambda=0$: the error of model $L_s(h_a)$ is given by Empirical Risk, but may over fitting

if $\lambda=\infty$: When $\lambda \rightarrow \infty$, each error converge to $L_s(h^*) \approx L_D(h^*)$ where h^* is minimum complexity model of regularization

then we can take λ^* , when $L_D(h_a)$ is Minimum.

2) Exercise 2

1)?

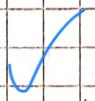
2.2) the SGD algorithm is learning algorithm for updating the vector \vec{w} . Since \vec{w} represents model, modifies \vec{w} iteratively means to fix model for making less error.



Assume the loss function is $j(w)$, the GD :

$\vec{w} = \text{Random initialization}$

while (not convergence):



Compute $\frac{\partial j(w)}{\partial w}$

← we compute the gradient of loss, which represents the direction of growth of error

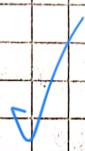
$w = w - \eta \cdot \frac{\partial j(w)}{\partial w}$

← we take the opposite direction

return w

and we move according η (learning rate)

Step.



Since computing gradient for all point of whole set is requires lot of time,
We use Gradient for saves time. ✓
Statistics

SGD take a batch of sample and we compute gradient of them

B = batch of sample

$$\frac{\partial J(w)}{\partial w} = \frac{1}{B} \sum_{k=0}^B \frac{\partial j_k(w)}{\partial w}$$

→ average gradient

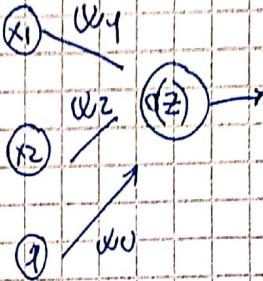
$$w = w - \eta \cdot \frac{\partial J(w)}{\partial w}$$



$$T(z) = e^z$$

2.3) We need compute $\frac{\partial \text{loss}}{\partial w_1}$

$$\frac{\partial \text{loss}}{\partial w_1} = \frac{\partial \text{loss}}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$



$$\frac{\partial \text{loss}}{\partial z} = \frac{\partial (e^z - y)^4}{\partial z} = 4e^z \cdot 4(e^z - y)^3 e^z$$

$$z = w_1 \cdot x_1 + w_2 \cdot x_2 + w_0$$

$$\begin{aligned} \frac{\partial z}{\partial w_1} &= w_1 \cdot x_1 \\ &= 4(e^z - y) \cdot e^z \cdot x_1 \quad \boxed{\text{A}} \quad \checkmark \end{aligned}$$

$$\text{loss} (T(z) - y)^4$$

$$(e^z - y)^4$$

$$(e^z - y)^4$$

$$\frac{\partial \text{loss}}{\partial w_2} = \frac{\partial \text{loss}}{\partial z} \cdot \frac{\partial z}{\partial w_2} = 4(e^z - y)^3 \cdot e^z \cdot x_2 \quad \boxed{\text{A}} \quad \checkmark$$

$$\frac{\partial \text{loss}}{\partial w_3} = \frac{\partial \text{loss}}{\partial z} \cdot \frac{\partial z}{\partial w_3} = 4(e^z - y)^3 \cdot e^z \cdot 1 \quad \boxed{\text{A}} \quad \checkmark$$

$$-0.5$$

So in update we write:

$$w^{t+1} = w^t - \eta \left[\begin{array}{c} \boxed{\text{A}} \\ \boxed{\text{A}} \\ \boxed{\text{A}} \end{array} \right]$$

1) exercise 1

INPUT: cluster

output: DECODEGRAPH



LINKAGE-based clustering is a clustering method where input are clusters and merge them according the rule (single clustering, average clustering) based on the closest clusters.



General algorithm:

divide each point in to a cluster (each point is a cluster)

NEAR

Merge to the closest clusters



Common termination:

- ALL cluster is merged in A single BIG cluster



- AFTER T-Step (T is number of step MAX Step)



- the clusters are in convergence situation \rightarrow

99



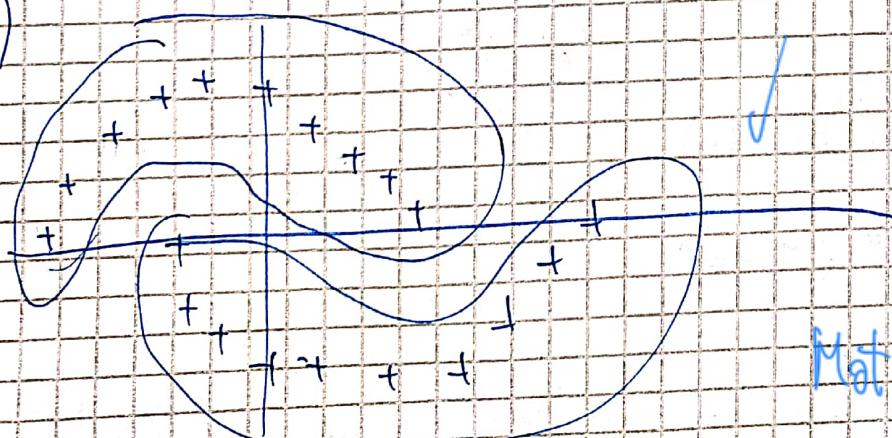
2) Given A | B 2 cluster

merge

it takes $\min \{ \text{dist}(A, B) \}$



3)



Motivation?

