

Lecture (6)

⇒ Gradient descent (optimization technique to minimize the loss function and thus fine tune the parameters)

Ex) Linear Regression with (SGD)

$$y = B_0 + B_1 X, \text{ assume } B_0 = 0, B_1 = 0$$

$$y = 0 + 0 \Rightarrow y = 0$$

① Calculate error $x=1, y=1$

$$\Rightarrow p(i) - y(i)$$

(predicted - expected)

* substitute (X) to find y then $(p(x) - 1)$

$$\Rightarrow e = 0 - 1 \Rightarrow e = -1$$

② update parameters

$$\Rightarrow B_0(t+1) = B_0(t) - \text{alpha} * \text{error}$$

$$B_0(t+1) = 0 - 0.01 * -1$$

$$B_0(t+1) = 0.01$$

* Repeat this until convergence

$$\Rightarrow B_1(t+1) = B_1(t) - \text{alpha} * \text{error}$$

$$B_1(t+1) = 0 - 0.01 * -1$$

$$B_1(t+1) = 0.01$$

Lecture (7) Logistic Regression

$$f = B_0 + B_1 x_1 + B_2 x_2, \quad B_0 = 0, B_1 = 0, B_2 = 0, y = 0$$

~> sigmoid $\frac{1}{1 + e^{-\text{output}}} = \frac{1}{1 + e^0} = \frac{1}{2} = 0.5$ $x_1 = 2.7, x_2 = 2.5$

use the update equation $(b = b + \text{alpha} * (y - \text{prediction}) * p * (1-p) * X$

$$B_0 = B_0 + \alpha * (0 - 0.5) * 0.5 * (1 - 0.5) * 1$$

$$B_0 = -0.6375$$

⋮

$$\text{Accuracy} = \frac{\text{Correct predictions} * 100\%}{\text{Total predictions}}$$

Lecture (8) classification

⇒ Binary classifiers

→ SGD classifier

→ SVM

⇒ Multiclass classifiers

→ Random Forest

→ Naive Bayes

⇒ multi-label classification

→ multiple classes for each instance

⇒ We can perform multiclass classification using Binary classifiers

→ One versus All (OVA)

→ We get a score from each classifier to decide on the correct prediction

→ One versus One (OVO) (SVM)

→ train Binary classifier for each pair of (digits)^{Ex}

of needed classifiers = $\frac{N(N-1)}{2}$ (45 for MNIST)

Each classifier needs to train on a subset of training set

* Performance measures

① Accuracy

→ $\frac{\text{correct predictions}}{\text{total \# of predictions}}$

→ measuring accuracy on certain # of tests (cross validation)

② confusion matrix

→ recall = $\frac{TP}{TP+FN}$, precision = $\frac{TP}{TP+FP}$

③ F1 score $F_1 = \frac{TP}{TP + \frac{FN+FP}{2}}$

* If we increase decision threshold ⇒ (↑ precision) eliminate FP

* // // Decrease // // ⇒ (↑ recall) eliminate FN

Lecture (9) Knn

⇒ Supervised learning, used for regression & classification
 (a lazy learner algorithm)

- How it works?

Select # of neighbors and calculate Euclidean distance between the new point & each point of (K) then take the nearest K's to the point and assigns the new point to the category in which number of neighbors is Maximum

⇒ pros & cons

pros: (easy to implement, robust to noisy data, efficient when training data is large)

cons: (Slow and lazy, high computational cost)

$$\text{Euclidean Distance} = \sqrt{\sum_{i=1}^n (a_i - b_i)^2} \quad x_1 = 8.0936, x_2 = 3.3677$$

x_1	x_2	y	$(x_1 - x_1)^2$	$(x_2 - y_1)^2$	Sum	Distance
3.3952	2.3312	0	22.0749	1.0701	23.1450	4.810

take x_1 from table - x_1 new
 same as x_1
 Sum the previous two columns
 Take the Sqrt of the previous column

then take the least distance according to # of neighbors and make the decision

Lecture (10) SVM (Binary classification)

⇒ high accuracy & low computation power, used for regression & classification, the objective is to find hyperplane in n -dimension which classify the data

$$\text{output} = w(B_1 x_1) + (B_2 x_2)$$

⇒ Optimum hyperplane (Maximum margin to ensure classifying with more confidence)

⇒ Kernel trick to transform the space into high dimensional space and use dot product without computing new higher space dimensions

→ Linear Kernel

→ Polynomial Kernel (more general)

$$\left\{ \begin{array}{l} \text{update if output} > 1, b = b(1 - \frac{1}{e}) \text{ ts iterates} \\ \text{update if output} < -1 \Rightarrow b = (1 - \frac{1}{e}) * b + \frac{1}{\text{lambda} * t} (y * x) \end{array} \right.$$

⇒ Output margin (greater than 0 ⇒ label 1, otherwise label -1)

Lecture (11) Clustering

Identify similar instances & assign them into one cluster, used for classification, anomaly detection and semi-supervised learning

K-means: # Initialize centroids & find the Euclidean distance between the centroid and point then assign the point to certain cluster according to (K)

Optimum # of clusters (K)

→ Inertia: Mean squared distance between instance & closest centroid

→ Silhouette coefficient: between (1 & -1)

→ 1 or close to it means that the instance is inside its own cluster and far from other clusters

→ 0 or close means that instance is close to boundary

→ -1 or close means that instance is assigned to wrong cluster

We plot the silhouette score with respect to K to find the optimum # of clusters

⇒ We can use clustering for preprocessing (dimensionality reduction)

Lecture (12) ANN

Perceptron is used for classification and regression, a single layer of perceptrons consists of (TLUs)

The slides are step function for activation (either 0 or 1)

DNN: When it has two or more hidden layers

⇒ Multilayer perceptron steps:-

① Forward pass ② back propagation ③ Fine tune (Gradient descent)

MLP's for Regression: (1 to 5 hidden layers, 10 to 100 neurons per hidden layer MSE for loss)

MLP's for classification: (uses softmax activation to ensure probabilities between 1 and $e^{-\infty}$ for each class, 1 to 5 hidden layers, 10 to 100 neurons per layer uses logistic function for binary & multilabel classification and softmax for multi-class, cross entropy for loss)

Q1) Advantages of loss function?

① helps in fine tune parameters

② measure the distribution between the true & predicted values (cross entropy)

Q2) What is soft max used for?

Ensure that the probabilities are between 1 and 0 for each class (useful in multiclass classification problems)

Q2.3) Multi independent task need?

→ In context of NN there are many answers for the question

① input layer, shared layers, task specific layer

② Loss Function ③ training process

The correct answer: (Chat Gpt) multi task learning needs validation data

Q4) The best algorithm for multiclass classification?
Neural networks

Q5) Knn euclidean distance?

Q6) How SVM classify with more confidence?
Choosing the optimum hyper plane which provide maximum margin

Q7) How PCA works?

Find correlated features and capture the highest variance (impact on target value) by calculating Eigen values of components and choose the highest value

أي سؤال تجاه المادة أو مساعدة،
تواصل مع الحساب المعطى 

