

M6. Support Vector Machines

1. What would happen when you use very large value of C(C->infinity)?

We can still classify data correctly for given setting of hyper parameter C

2. In SVM, if the number of input features is 3, then the hyperplane is a _plane_____.

3. Suppose you have trained an SVM with linear decision boundary after training SVM, you correctly infer that your SVM model is under fitting. Which of the following option would you more likely to consider iterating SVM next time?

Increase Data Points

4. Soft margin SVM is overly sensitive to noise as compared to Hard Margin SVM.

False

5. Given the primal constrained optimization problem as following.

$$L(x,y,\lambda)=f(x,y)-\lambda g(x,y)$$

Select the correct order of step for achieving dual formulation of the primal problem.

1. Formulate the function only in terms of lagrange-multipliers.
2. Eliminate x, y by taking the partial derivative of primal function w.r.t x and y
3. Maximize the dual function to obtain constrained-minimum of primal function.

6. Suppose you have a dataset with n = 10 features and m = 5000 examples.

After training your logistic regression classifier with gradient descent, you find that it has underfit the training set and does not achieve the desired performance on the training or cross validation sets. Which of the following might be promising steps to take? Check all that apply.

Try using a neural network with a large number of hidden units

Create / add new polynomial features

Use an SVM with a Gaussian Kernel.

M7. Kernel Methods

1. One of the most commonly used kernels in SVM is the Gaussian RBF kernel: $k(x_i, x_j) = \exp(-\|x_i - x_j\|^2 / 2\sigma)$. Suppose we have three points, z_1 , z_2 , and x . z_1 is geometrically very close to x , and z_2 is geometrically far away from x . What is the value of $k(z_1, x)$ and $k(z_2, x)$?. Choose one of the following:

$k(z_1, x)$ will be close to 1 and $k(z_2, x)$ will be close to 0.

2. You are training an RBF SVM with the following parameters: C (slack penalty) and $\gamma = 1/2(\sigma^2)$ (where σ^2 is the variance of the RBF kernel). How should you tweak the parameters to reduce overfitting?

Reduce C and/or reduce γ

3. The kernel trick

exploits the fact that in many learning algorithms, the weights can be written as a linear combination of input points

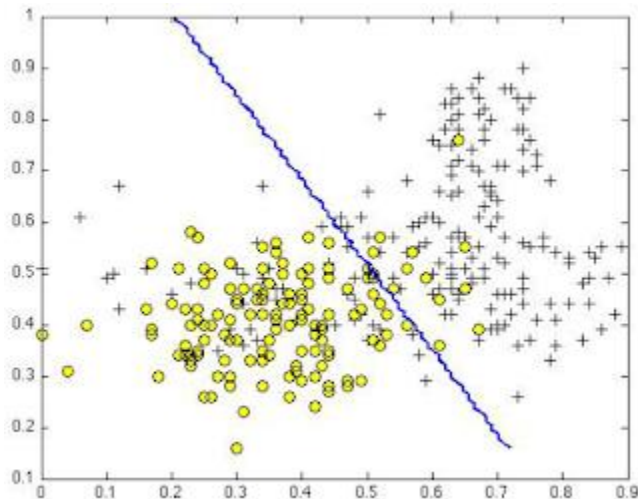
4. Which of the following can help to reduce overfitting in an SVM classifier?

Use of slack variables

5. Suppose you are using RBF kernel in SVM with high Gamma value. What does this signify?

The model would consider only the points close to the hyperplane for modeling

6. Suppose you have trained an SVM classifier with a Gaussian kernel, and it learned the following decision boundary on the training set:



You suspect that the SVM is underfitting your dataset. Should you try increasing or decreasing C ? Increasing or decreasing σ^2 ?

It would be reasonable to try increasing C . It would also be reasonable to try decreasing σ^2

M9. Ensemble Methods

1. Compared to a simpler model, a more complex model is more likely to be

More prone to overfitting, have a higher variance.

2. Which of the following is true about bootstrapping and bagging?

Bagging works best on "unstable learners", i.e. the models that are sensitive to small changes in the input data.

3. In the case of an imbalanced dataset with two classes, the decision boundary would be closer to ___Data points of the minority class___ (if we don't take any additional measures to correct for the class imbalance).

4. Which of the following is NOT True about Ensemble Techniques?

Bagging and Boosting are the only available ensemble techniques.

5. Suppose there are 25 base classifiers. Each classifier has error rates of $e = 0.35$. Suppose you are using averaging as ensemble technique. What will be the probabilities that ensemble of above 25 classifiers will make a wrong prediction?

Note: All classifiers are independent of each other

$$e = 1 - \text{accuracy} = 0.35$$

$$\sum_{i=13}^{25} \binom{25}{i} e^i (1-e)^{(25-i)} = 0.06$$

M10. Neural Networks

1.

GPA	STUDIED	PASSED
Low	No	No
Low	Yes	Yes
Medium	No	No
Medium	Yes	Yes
High	No	Yes
High	Yes	Yes

Find the entropy of H(Passed).

$$H(\text{passed}) = -\left(\frac{2}{6}\log_2\frac{2}{6} + \frac{4}{6}\log_2\frac{4}{6}\right) = \log_2 3 - \frac{2}{3} \approx 0.92$$

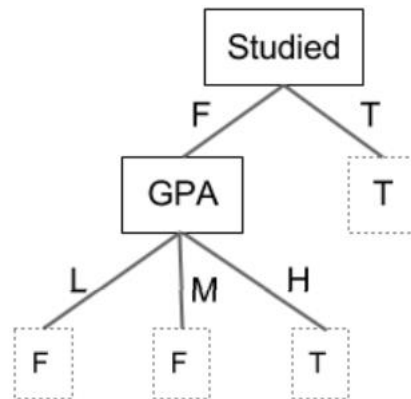
What is the entropy $H(\text{Passed} \mid \text{GPA})$

$$H(\text{passed} \mid \text{GPA}) = -\frac{1}{3}\left(\frac{1}{2}\log_2\frac{1}{2} + \frac{1}{2}\log_2\frac{1}{2}\right) - \frac{1}{3}\left(\frac{1}{2}\log_2\frac{1}{2} + \frac{1}{2}\log_2\frac{1}{2}\right) - \frac{1}{3}(1\log_2 1) = \frac{2}{3} \approx 0.66$$

What is the entropy $H(\text{Passed} \mid \text{Studied})$

$$H(\text{passed} \mid \text{GPA}) = -\frac{1}{2}\left(\frac{1}{3}\log_2\frac{1}{3} + \frac{2}{3}\log_2\frac{2}{3}\right) - \frac{1}{2}(1\log_2 1) = \frac{1}{2}\log_2 3 - \frac{1}{3} \approx 0.46$$

Draw the full decision tree that would be learned for this dataset. You do not need to show any calculations.



2. A 4-input neuron has weights 1, 2, 3, and 4. The transfer function is linear with the constant of proportionality being equal to 2. The inputs are 4, 3, 2, and 1 respectively. What will be the output?

$$\text{Output} = 2 * (1*4 + 2*3 + 3*2 + 4*1) = 40$$

3. If a neural network has q units in layer j , r units in layer $j+1$, then the layer governing the mapping from layer j to $j+1$ will be of the dimension

$$r \times (q + 1)$$

4. If $w=[w_1...w_m]^T, x=[x_1...x_m]^T$ are linear vectors which the decision function $g(z)$ of the perceptron takes then, the value of the decision function z is given by:

$$Z=w_1x_1+...+w_mx_m$$

5. $\begin{cases} ax & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$ This is the equation for which of the following activation function?

Leaky ReLU

6. Which of the following statements are true? Check all that apply.

If a neural network is overfitting the data, one solution would be to increase the regularization parameter.

In a neural network with many layers, we think of each successive layer as being able to use the earlier layers as features, so as to be able to compute the increasingly complex functions.

7. What is perceptron?

a single layer feed-forward neural network with pre-processing

M11. Unsupervised Learning

1. Targeted Marketing, Recommender System, Customer Segmentation are the applications in

Unsupervised Learning: Clustering

2. Assume, you want to cluster 7 observations into 3 clusters using K-Means clustering algorithm. After first iteration clusters, C1, C2, C3 has following observations:

C1: {(2,2), (4,4), (6,6)}

C2: {(0,4), (4,0)}

C3: {(5,5), (9,9)}

What will be the cluster centroids if you want to proceed for second iteration?

C1: (4,4), C2: (2,2), C3: (7,7)

3. If you are using Multinomial mixture models with the expectation-maximization algorithm for clustering a set of data points into two clusters, which of the assumptions are important:

All the data points follow two multinomial distribution

4. Which of the following is/are not true about Centroid based K-Means clustering algorithm and Distribution based expectation-maximization clustering algorithm:

- 1. Both starts with random initializations**
- 2. Both are iterative algorithms**

3. Both have strong assumptions that the data points must fulfill
4. Both are sensitive to outliers
5. Expectation maximization algorithm is a special case of K-Means
6. Both requires prior knowledge of the no. of desired clusters
7. The results produced by both are non-reproducible.

5 only

4. Which of the following statements about the K-means algorithm are correct?

The K-means algorithm is sensitive to outliers

M12. Anomaly Detection

1. Which of the following statement is TRUE?

The nature of our business problem determines how outliers are used.

2. When doing K-Nearest Neighbors classification, which of the following would make the model more resilient to outliers?

Using a smaller K.

3. Consider a list of data points.

The stats:

n : 15

min : -5

max : 15

mean : 5.333333333333333

median : 5

mode : (3, 3)

stdev : 4.376706016578628

q1 : 3

q3 : 7.5

iqr : 4.5

What would be the IQR based limits to detect the outliers.

$$\text{IQR} = Q3 - Q1 = 7.5 - 3 = 4.5$$

$$\text{Maximum_limit} = Q3 + 1.5 * \text{IQR}$$

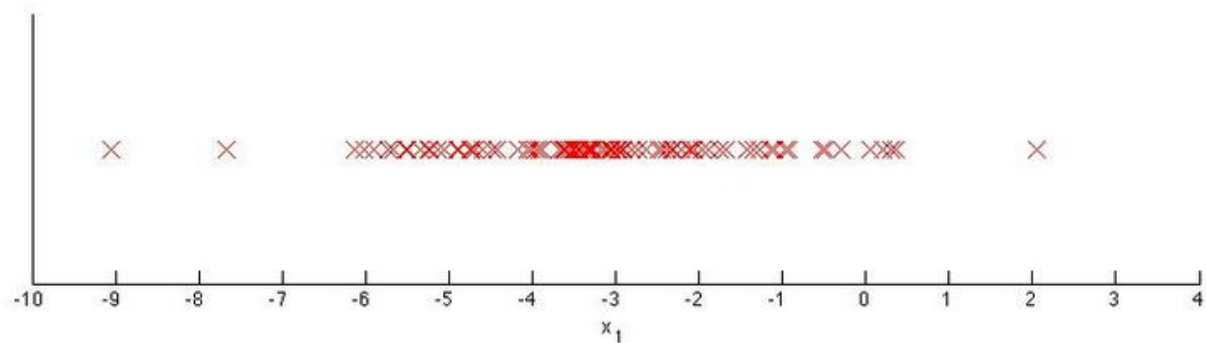
$$\text{Minimum_limit} = Q1 - 1.5 * \text{IQR}$$

$$[-1.5 * 4.5, 1.5 * 4.5] = [-6.75, 6.75]$$

4. Suppose you have trained an anomaly detection system for fraud detection, and your system that flags anomalies when $p(x)$ is less than ϵ , and you find on the cross-validation set that it is missing many fraudulent transactions (i.e., failing to flag them as anomalies). What should you do?

Increase ϵ

5. Suppose you fit the Gaussian Distribution parameters μ_1, σ_1 to the 1D dataset given below. Your task is to predict the outliers. Which of the following values for μ_1, σ_1 you might get?



-3, 4