
CSE 151A - Homework 06

Due: Wednesday, May 13, 2020

Write your solutions to the following problems by either typing them up or handwriting them on another piece of paper. Unless otherwise noted by the problem's instructions, show your work or provide some justification for your answer. Homeworks are due via Gradescope on Wednesday at 11:59 p.m.

Essential Problem 1.

A vector \vec{a} is a linear combination of vectors $\vec{v}_1, \dots, \vec{v}_k$ if there exist scalars c_1, \dots, c_k such that $\vec{a} = \sum_{i=1}^k c_i \vec{v}_i$. Suppose we run the perceptron train algorithm on training data $(\vec{x}^{(1)}, y_1), \dots, (\vec{x}^{(n)}, y_n)$ using an initial weight vector of $\vec{0}$, and we get an output \vec{w} after a single pass on the data. Is \vec{w} a linear combination of $\text{Aug}(\vec{x}^{(1)}), \dots, \text{Aug}(\vec{x}^{(n)})$? Justify your answer.

Solution:

Recall our update rule

$$\vec{w}^{(t)} = \vec{w}^{(t-1)} - \alpha \begin{cases} \text{Aug}(\vec{x}^{(i)}) & \vec{w}^{(t-1)} \cdot \text{Aug}(\vec{x}^{(i)}) \geq 0 \\ -\text{Aug}(\vec{x}^{(i)}) & \vec{w}^{(t-1)} \cdot \text{Aug}(\vec{x}^{(i)}) < 0 \end{cases}$$

The first case results from a misclassified positive label, meaning the true label is $y_i = -1$. The opposite is true for our second update case which occurs when the true label $y_i = +1$. From this, we can generalize an update as $\vec{w}^{(t-1)} - \alpha(-y_i \text{Aug}(\vec{x}^{(i)}))$.

Define an indicator function which will return 1 if a point is misclassified, and 0 otherwise.

$$\mathbf{1}_{\text{Miss}}(i) = \begin{cases} 1 & \text{sign}(\vec{w}^{(t-1)} \cdot \text{Aug}(\vec{x}^{(i)})) \neq \text{sign}(y_i) \\ 0 & \text{sign}(\vec{w}^{(t-1)} \cdot \text{Aug}(\vec{x}^{(i)})) = \text{sign}(y_i) \end{cases}$$

Since multiplying our update by zero is essentially the same as ignoring it, we can create an update rule which applies to the entire data set as

$$\vec{w}^{(t)} = \vec{w}^{(t-1)} + \mathbf{1}_{\text{Miss}}(i) \cdot (\alpha y_i \text{Aug}(\vec{x}^{(i)}))$$

We can reexpress the final output of all updates as a summation

$$\vec{w} = \sum_{i=1}^n \mathbf{1}_{\text{Miss}}(i) \cdot (\alpha y_i \text{Aug}(\vec{x}^{(i)}))$$

This is indeed of the form $\vec{w} = \sum_{i=1}^n c_i \text{Aug}(\vec{x}^{(i)})$, and we can note that c_i will be 0 when a point is properly classified, $-\alpha$ when a point is misclassified as positive (y_i negative), and $+\alpha$ when a point is misclassified as negative (y_i positive).

Essential Problem 2.

The perceptron learning algorithm is run on a data set. It converges after performing $n = a + b$ updates. Of these n updates, a are on data points with label $+1$ and b are on points whose label is -1 . What is the final value of the “bias weight”, w_0 ? Assume that the initial value of \vec{w} is $\vec{0}$, and the learning rate α is set to 1.

Solution: Notice from our update rule that a point misclassified as positive (with a true label of $y_i = -1$) will result in an update of

$$\vec{w}^{(t)} = \vec{w}^{(t-1)} - \alpha \text{Aug}(\vec{x}^{(i)})$$

The opposite update is true of points misclassified as negative (with true labels of $y_i = +1$), resulting in an update of

$$\vec{w}^{(t)} = \vec{w}^{(t-1)} - \alpha(-\text{Aug}(\vec{x}^{(i)}))$$

Since there are a points misclassified as negative with true $y_i = +1$, and b points misclassified as positive with true $y_i = -1$, then there are a updates resulting in an addition of $\alpha \text{Aug}(\vec{x}^{(i)})$ and b updates resulting in its subtraction.

Both the value of α and all $x_0^{(i)}$ are equal to 1.

Therefore the final value of $w_0 = a(1) + b(-1) = a - b$.

Essential Problem 3.

An SVM classifier is learned for a data set in \mathbb{R}^2 . The weight vector of this classifier is $\vec{w} = (-12, 3, 4)^T$.

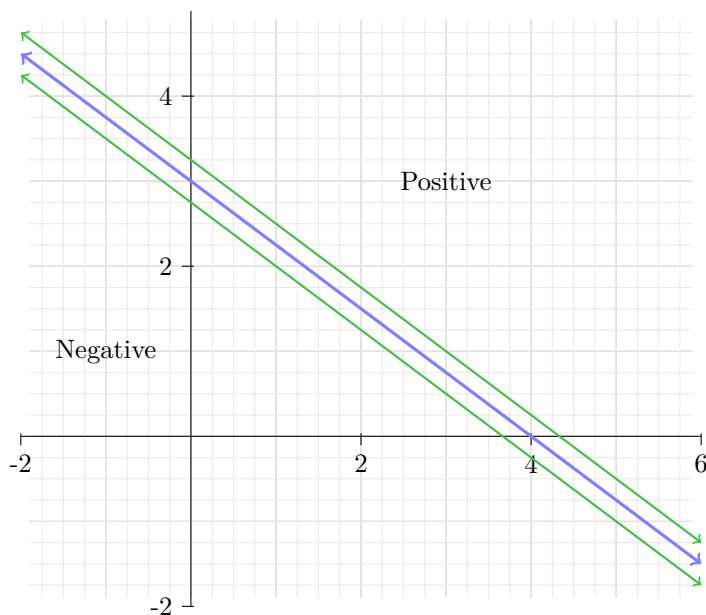
- a) Draw the decision boundary and mark where it intersects the axes.

Solution: See b)

- b) Draw the margin boundaries (the lines where the dot product is 1 and -1). Also mark where these intersect the axes.

Solution: Let's solve $\vec{w} \cdot \text{Aug}(\vec{x}) = 0$ for x_2 in order to plot the decision boundary. This yields $x_2 = -\frac{3}{4}x_1 + 3$ for the equation of our decision boundary.

We can likewise solve for the margin boundaries in terms of x_2 . Solving $\vec{w} \cdot \text{Aug}(\vec{x}) \pm 1$ yields $x_2 = -\frac{3}{4}x_1 + \frac{13}{4}$ and $x_2 = -\frac{3}{4}x_1 + \frac{11}{4}$ as the equations for our margin boundaries.



- c) What is the (shortest) distance between the margin lines? Show your work.

Solution:

$$distance = 2 \cdot \frac{1}{5} = 0.4$$

- d) How would the point $(2, 2)$ be classified by the SVM? Show your work.

Solution: Recall that a prediction can be made by determining which side of the decision boundary a new point lies on.

$$\vec{w} \cdot \text{Aug}((2, 2)^T) = -12 + 6 + 8 \geq 0$$

Therefore we predict that $(2, 2)^T$ belongs to the positive class with label $+1$.

- e) It turns out that there were only two support vectors, and they both have the form $(1, ?)$. What are they? Describe how you know.

Solution: Since we have scaled our margin boundaries such that the dot product is $+1$ and -1 , we know that the dot product between \vec{w} and our two support vectors must be $+1$ and -1 .

$$\begin{aligned} \vec{w} \cdot \text{Aug}((1, ?)^T) &= -12 + 3 + 4(?) = \pm 1 \\ \Rightarrow \quad ? &= \frac{9 \pm 1}{4} \end{aligned}$$

Thus our two support vectors must be $(1, 2.5)^T$ and $(1, 2)^T$ on the positive and negative side, respectively.

Essential Problem 4.

Consider the following sentences:

- That is a big bug.
- No, seriously, that is a really big bug.
- Everybody run!

Encode each of the above sentences as bag-of-words feature vectors, using the set of all words in the three sentences as your dictionary. Sort the words alphabetically when creating your dictionary, and disregard capitalization and punctuation. You do not need to show your work (but doing so can help you get credit if you make a simple mistake).

Solution: Alphabetically arranged, the unique words that occur in our sentences (ignoring capitalization and punctuation) constitute our dictionary. For each sentence (labeled A, B, C respectively) we encode the count of each word in our dictionary.

Sentence	a	big	bug	everyb... is	no	really	run	seriously that	
A	1	1	1	0	1	0	0	0	1
B	1	1	1	0	1	1	0	1	1
C	0	0	0	1	0	0	1	0	0

Plus Problem 1. (16 plus points)

This is a legitimate machine learning mini-project that ties together several ideas, so it is worth 16 plus points.

The file <http://cse151a.com/data/yelp/train.csv> contains 10,000 Yelp reviews along with the score the user left (from 1 to 5, with 5 being the best). In this plus problem, you'll train an SVM to do sentiment analysis on these reviews and predict the sentiment of an unlabeled piece of text.

- a) Split the data 75%/25% into training and validation sets, encode the training data using a bag of words feature representation, and train a (linear, soft-margin) support vector machine. When training, consider any review with a score of 4 or higher to be a positive review, and anything with a smaller score to be a negative review. Find the value of C that minimizes the error of your classifier on the validation set and make a plot of the validation error as a function of C .

For this part, turn in four things:

1. the value of C that was best,
2. the training and validation error that corresponded to this choice of C ,
3. your plot, and
4. your code.

You can use whatever machine learning libraries you like in whatever language you'd like. Note that most languages have libraries which will do the bag-of-words encoding for you. For instance, `sklearn` has this feature (but I'll let you Google for it!).

- b) Is the data in `train.csv` linearly separable? How do you know?

- c) Give an example of:

- A sentence that you think is positive that your predictor got right.
- A sentence that you think is negative that your predictor got right.
- A sentence that you think is positive that your predictor got wrong.

Solution: See https://go.ucsd.edu/2Z7RBRj for a solution notebook.
