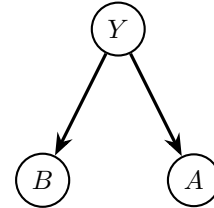


CSCD84: Artificial Intelligence

Worksheet: Naive Bayes, Perceptron, and Logistic Regression

Q1. (Naive Bayes) In this question, we will train a Naive Bayes classifier to predict class labels Y as a function of input features A and B . Y , A , and B are all binary variables, with domains 0 and 1. We are given 10 training points from which we will estimate our distribution.

A	1	1	1	1	0	1	0	1	1	1
B	1	0	0	1	1	1	1	0	1	1
Y	1	1	0	0	0	1	1	0	0	0



- 1.a.** What are the maximum likelihood estimates for the tables $\mathbb{P}(Y)$, $\mathbb{P}(A | Y)$, and $\mathbb{P}(B | Y)$?
- 1.b.** Consider a new data point $(A = 1, B = 1)$. What label would this classifier assign to this sample?
- 1.c.** Compute the new distribution for $\mathbb{P}(A | Y)$ given Laplace Smoothing with $k = 2$.

Q2. (Perceptrons) For all parts of this question perceptrons should output 1 if $w_n + \sum_{i=0}^{n-1} w_i x_i \geq 0$ and 0 otherwise. The weight w_n is called the bias weight.

- 2.a.** Assuming there will be two inputs x_0 and x_1 , each with possible values in $\{0, 1\}$, give values for a triple of weights (w_0, w_1, w_2) such that the corresponding perceptron would act as a NAND gate for the two inputs. (Weight w_2 is the bias weight.)
Note that a NAND gate outputs 1 when at least one of the inputs is 0; it outputs 0 otherwise.
- 2.b.** Draw a perceptron, with weights, that accepts a single integer x and outputs 1 if and only if the input is greater than or equal to -8 . Be sure to include the bias input of value 1 and its weight in your diagram. Draw another perceptron that outputs 1 if and only if the input is less than or equal to 8.
- 2.c.** Using the previous perceptrons, create a two-layer perceptron that outputs 1 if $|x| \leq 8$, and 0 otherwise.
- 2.d.** Suppose we want to train a perceptron to compare two numbers x_0 and x_1 and produce output $y = 1$ provided that x_1 exceeds x_0 by at least 5. Assume that the initial weight vector is: $(w_0, w_1, w_2) = (-1, -1, 1)$. Consider a first training example: $((x_0, x_1), y) = ((3, 8), 1)$. This says that with inputs 3, and 8, the output y should be 1, since 8 exceeds 3 by 5. What will be the new values of the weights after this training example has been processed one time? Assume the learning rate is 2.
- 2.e.** Continuing with the last example, now suppose that the next step of training involves a different training example: $((2, 5), 0)$. The output for this example should be 0, since 5 does not exceed 2 by at least 5. Starting with the weights already learned in the first step, determine what the adjusted weights should be after this new example has also been processed once.

Q3. (Multiclass Perceptrons) Consider the problem of classifying text items into three topic areas: AI, Medicine, and Hiking. We have four training examples:

e_1 : Hill climbing informs gradient descent. (AI)

e_2 : Alzheimers is a gradual descent. (Medicine)

e_3 : Hill climbing helps. (Medicine)

e_4 : Gorp helps hill climbing. (Hiking)

3.a. Convert the four training examples into the (vector, category) form indicated by the table. (You'll use the reference vocabulary given in the table.)

e_i	Alzheimers	climbing	descent	gorp	gradient	gradual	helps	hill	(category)
e_1									
e_2									
e_3									
e_4									





3.b. We will start with a weight vector for each of the 3 categories (\mathbf{W}_A for AI, \mathbf{W}_M for Medicine, and \mathbf{W}_H for Hiking), with all weights 0 except 1 for the bias on the AI vector.

\mathbf{W}	bias	Alzheimers	climbing	descent	gorp	gradient	gradual	helps	hill
\mathbf{W}_A	1	0	0	0	0	0	0	0	0
\mathbf{W}_M	0	0	0	0	0	0	0	0	0
\mathbf{W}_H	0	0	0	0	0	0	0	0	0

Perform one epoch of training showing the resulting changed vectors whenever there is a change made to a weight vector. When each training example is processed, if any weight vector does not change, do not rewrite that vector. Whenever a vector is updated, write the new vector on its own line below, clearly indicating which category it belongs to. For example, $\mathbf{W}_H = \dots$

3.c. Training could take many epochs to converge. See if you can skip the training and manually provide a set of three weight vectors that will correctly handle the training examples.

Q4. (Hinge Loss) Here are two reviews of Perfect Blue, from Rotten Tomatoes:

	Panos Kotzathanasis <i>Asian Movie Pulse</i>	 "Perfect Blue" is an artistic and technical masterpiece; however, what is of utmost importance is the fact that Satoshi Kon never deteriorated from the high standards he set here, in the first project that was entirely his own.	January 26, 2020
		Full Review	
	Derek Smith <i>Cinematic Reflections</i>	 [An] nime thriller [that] often plays as an examination of identity and celebrity, but ultimately gets so lost in its own complex structure that it doesn't end up saying much at all.	August 19, 2006
		Full Review Original Score: 2/4	

Rotten Tomatoes has classified these reviews as “positive” and “negative,” respectively, as indicated by the intact tomato on the top and the splatter on the bottom.

In this assignment, you will create a simple text classification system that can perform this task automatically. We'll warm up with the following set of four mini-reviews, each labeled positive (+1) or negative (−1):

1. x_1 : not good; label: (−1)
2. x_2 : pretty bad; label: (−1)
3. x_3 : good plot; label: (+1)
4. x_4 : pretty scenery; label: (+1)

Each review x is mapped onto a feature vector $\phi(x)$, which maps each word to the number of occurrences of that word in the review. For example, the second review maps to the (sparse) feature vector $\phi(x) = \{\text{pretty} : 1, \text{bad} : 1\}$. The hinge loss is defined as

$$L_{\text{hinge}}(x, y, \mathbf{w}) = \max\{0, 1 - \mathbf{w} \cdot \phi(x)y\},$$

where x is the review text, y is the correct label, \mathbf{w} is the weight vector.

4.a. (Gradient Descent) Suppose we run stochastic gradient descent once for each of the 4 samples in the order given above, updating the weights according to

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla_{\mathbf{w}} L_{\text{hinge}}(x, y, \mathbf{w}).$$

After the updates, what are the weights of the six words (“pretty”, “good”, “bad”, “plot”, “not”, “scenery”) that appear in the above reviews?

- Use $\eta = 0.1$ as the step size.
- Initialize $\mathbf{w} = [0, 0, 0, 0, 0, 0]$.
- The gradient $\nabla_{\mathbf{w}} L_{\text{hinge}}(x, y, \mathbf{w}) = 0$ when margin is exactly 1 (i.e., when $\mathbf{w} \cdot \phi(x)y = 1$).

Present your answer as a weight vector that contains a numerical value for each of the tokens in the reviews (“pretty”, “good”, “bad”, “plot”, “not”, “scenery”).

4.b. (Linearly Inseparable) Given the following dataset of reviews:

1. x_1 : bad; label: (−1)
2. x_2 : good; label: (+1)
3. x_3 : not bad; label: (+1)
4. x_4 : not good; label: (−1)

Prove that no linear classifier using word features (i.e., word count) can get zero error on this dataset. Remember that this is a question about classifiers, not optimization algorithms; your proof should be true for any linear classifier of the form $f_{\mathbf{w}}(x) = \text{sign}(\mathbf{w} \cdot \phi(x))$, regardless of how the weights are learned. Propose a single additional feature for your dataset that we could augment the feature vector with that would fix this problem.

Q5. (Squared Loss) Suppose that we are now interested in predicting a numeric rating for movie reviews. We will use a non-linear predictor that takes a movie review x and returns $\sigma(\mathbf{w} \cdot \phi(x))$, where $\sigma(z) = (1 + e^{-z})^{-1}$ is the logistic function that squashes a real number to the range $(0, 1)$. For this problem, assume that the movie rating y is a real-valued variable in the range $[0, 1]$.

- 5.a.** Suppose that we wish to use squared loss. Write out the expression for $L(x, y, \mathbf{w})$ for a single datapoint (x, y) .
- 5.b.** Given $L(x, y, \mathbf{w})$ from the previous part, compute the gradient of the loss with respect to \mathbf{w} , $\nabla_{\mathbf{w}} L(x, y, \mathbf{w})$. Write the answer in terms of the predicted value $p = \sigma(\mathbf{w} \cdot \phi(x))$.
- 5.c.** Suppose there is one datapoint (x, y) with some arbitrary non-zero $\phi(x)$ and $y = 1$. Specify conditions for \mathbf{w} to make the magnitude of the gradient of the loss with respect to \mathbf{w} arbitrarily small (*i.e.*, minimize the magnitude of the gradient). Can the magnitude of the gradient with respect to \mathbf{w} ever be exactly zero? You are allowed to make the magnitude of \mathbf{w} arbitrarily large but not infinity.

Why does it matter? the reason why we're interested in the magnitude of the gradients is because it governs how far gradient descent will step. For example, if the gradient is close to zero when \mathbf{w} is very far from the optimum, then it could take a long time for gradient descent to reach the optimum (if at all). This is known as the vanishing gradient problem when training neural networks.