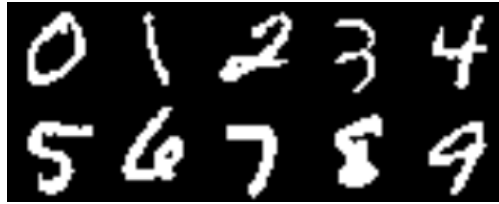


HW1: Decision trees and KNN

Please note that only PDF submissions are accepted. We encourage using L^AT_EX to produce your writeups. You'll need *mydefs.sty* and *notes.sty* which can be downloaded from the course page. You can use overleaf.com to write your latex file and save the result as PDF.

1. (not graded): The following are true/false questions. You don't need to answer the questions. Just tell us which ones you can't answer confidently in less than one minute. (You won't be graded on this.) If you can't answer at least 8, you should probably spend some extra time outside of class beefing up on elementary math. I would strongly suggest going through this math tutorial by Hal Daume: <https://web.cs.ucdavis.edu/%7Ehpirsiav/courses/MLf22/math4ml.pdf>
 - (a) $\log x + \log y = \log(xy)$
 - (b) $\log[ab^c] = \log a + (\log b)(\log c)$
 - (c) $\frac{\partial}{\partial x} \sigma(x) = \sigma(x) \times (1 - \sigma(x))$ where $\sigma(x) = 1/(1 + e^{-x})$
 - (d) The distance between the point (x_1, y_1) and line $ax + by + c$ is $(ax_1 + by_1 + c)/\sqrt{a^2 + b^2}$
 - (e) $\frac{\partial}{\partial x} \log x = -\frac{1}{x}$
 - (f) $p(a | b) = p(a, b)/p(b)$
 - (g) $p(x | y, z) = p(x | y)p(x | z)$
 - (h) $C(n, k) = C(n-1, k-1) + C(n-1, k)$, where $C(n, k)$ is the number of ways of choosing k objects from n
 - (i) $\|\alpha \mathbf{u} + \mathbf{v}\|^2 = \alpha^2 \|\mathbf{u}\|^2 + \|\mathbf{v}\|^2$, where $\|\cdot\|$ denotes Euclidean norm, α is a scalar and \mathbf{u} and \mathbf{v} are vectors
 - (j) $|\mathbf{u}^\top \mathbf{v}| \geq \|\mathbf{u}\| \times \|\mathbf{v}\|$, where $|\cdot|$ denotes absolute value and $\mathbf{u}^\top \mathbf{v}$ is the dot product of \mathbf{u} and \mathbf{v}
 - (k) $\int_{-\infty}^{\infty} dx \exp[-(\pi/2)x^2] = \sqrt{2}$
2. (not graded): Go through this Matlab tutorial by Stefan Roth:
<http://cs.brown.edu/courses/csci1950-g/docs/matlab/matlabtutorialcode.html>
3. In class, we looked at an example where all the attributes were binary (i.e., yes/no valued). Consider an example where instead of the attribute "Morning?", we had an attribute "Time" which specifies when the class begins.
 - (a) We can pick a threshold τ and use $(\text{Time} < \tau)?$ as a criteria to split the data in two. Explain how you might pick the optimal value of τ .
 - (b) In the decision tree learning algorithm discussed in class, once a binary attribute is used, the subtrees do not need to consider it. Explain why when there are continuous attributes this may not be the case.
4. Why memorizing the training data and doing table lookups is a bad strategy for learning? How do we prevent that in decision trees?
5. What does the decision boundary of 1-nearest neighbor classifier for 2 points (one positive, one negative) look like?
6. Does the accuracy of a kNN classifier using the Euclidean distance change if you (a) translate the data (b) scale the data (i.e., multiply the all the points by a constant), or (c) rotate the data? Explain. Answer the same for a kNN classifier using Manhattan distance¹.

¹http://en.wikipedia.org/wiki/Taxicab_geometry



7. Implement kNN in Matlab or Python for handwritten digit classification and submit all codes and plots:
- (a) Download MNIST digit dataset (60,000 training and 10,000 testing data points) and the starter code from the course page. Each row in the matrix represents a handwritten digit image. The starter code shows how to visualize an example data point in Matlab. The task is to predict the class (0 to 9) for a given test image, so it is a 10-way classification problem.
 - (b) Write a Matlab or Python function that implements kNN for this task and reports the accuracy for each class (10 numbers) as well as the average accuracy (one number).
 $[acc \text{ } acc_av] = kNN(images_train, labels_train, images_test, labels_test, k)$
where acc is a vector of length 10 and acc_av is a scalar. Look at a few correct and wrong predictions to see if it makes sense. To speed it up, in all experiments, you may use only the first 1000 testing images.
 - (c) For $k = 1$, change the number of training data points (30 to 10,000) to see the change in performance. Plot the average accuracy for 10 different dataset sizes. You may use command *logspace* in Matlab. In the plot, x-axis is for the number of training data and y-axis is for the accuracy.
 - (d) Show the effect of k on the accuracy. Make a plot similar to the above one with multiple colored curves on the top of each other (each for a particular k in [1 2 3 5 10].) You may use command *legend* in Matlab to name different colors.
 - (e) Choose the best k . First choose 2,000 training data randomly (to speed up the experiment). Then, split the training data randomly to two halves (the first for training and the second for cross-validation to choose the best k). Please plot the average accuracy wrt k on the validation set. You may search for k in this list: [1 2 3 5 10]. Finally, report the accuracy for the best k on the testing data.