# HW Assignment 8 (Due by 1:30 pm on Dec 3)

## 1 Theory (150 points)

1.  [**Kernel Nearest Neighbor**, 50 points]
    The nearest-neighbour classifier 1-NN assigns a new input vector $\mathbf{x}$ to the same class as that of the nearest input vector $\mathbf{x}_n$ from the training set, where in the simplest case, the distance is defined by the Euclidean metric $\|\mathbf{x} - \mathbf{x}_n\|^2$. By expressing this rule in terms of dot products and then making use of kernel substitution, formulate the nearest-neighbour classifier for a general nonlinear kernel.

2.  [**Distance-Weighted Nearest Neighbor**, 50 points]
    We have seen how to use kernels to formulate a distance-weighted nearest neighbor algorithm, when the labels are binary. Formulate a kernel-based, distance-weighted nearest neighbor that works for K classes, where $K \geq 2$.

3.  [**Naive Bayes**, 50 points]
    The Naive Bayes algorithm for text categorization presented in class treats all sections of a document equally, ignoring the fact that words in the title are often more important than words in the text in determining the document category. Describe how you would modify the Naive Bayes algorithm for text categorization to reflect the constraint that words in the title are $K$ times more important than the other words in the document for deciding the category, where $K$ is an input parameter (include pseudocode).

4.  [**Logistic Regression (\*)**, 50 points]
    Assume that a binary feature $x_i$ is equal to 1 for all training examples $\mathbf{x}$ belonging to a particular class $C_k$, and zero otherwise (i.e. $x_i$ perfectly separates examples from class $C_k$ from all other examples). Show that in this case the magnitude of the ML solution for $\mathbf{w}_k$ goes to infinity, thus motivating the use of a prior over the parameters (Hint: use the fact that the gradient on slide 24 must vanish at the solution).

## 2 Submission

Turn in a hard copy of your homework report at the beginning of class on the due date. On this theory assignment, **clear and complete explanations and proofs of your results are as important as getting the right answer**.