

CS 6350: Machine Learning Fall 2016

Gopal Menon
Homework 5

November 11, 2016

1 Warm up: Margins

1. [5 points] Suppose we want to use an SVM to learn the XOR function in two dimensions. We know that XOR is not linearly separable, so we apply a feature transformation. In order to do so, we map the input $[x_1, x_2]$ into a space consisting of two features: x_1 and x_1x_2 . All examples are Boolean (1 for positive and -1 for negative) What is the maximal margin? Draw the separating line back in original Euclidean input space.

x_1	x_2	x_1x_2	<i>Label</i>
-1	-1	+1	-1
-1	+1	-1	+1
+1	-1	-1	+1
+1	+1	+1	-1

XOR Function truth table

The XOR function truth table shows the original inputs x_1 and x_2 along with the new feature x_1x_2 . The label column shows $x_1 \oplus x_2$. Figure 1 shows the new space with features x_1 and x_1x_2 . The points in this input space are shown along with their labels. The axis for the feature x_1 is the linear separator with the maximum margin and is shown as the bold line. The four features in the new space become support vectors and the margin is half the distance between the lines joining the support vectors.

The input features x_1 and x_2 in the original feature space are shown in figure 2. The lines joining the support vectors now intersect each other. The line of maximum margin in figure 1 was the one where x_1 varied from $+\infty$ to $+\infty$ and x_1x_2 was 0. The corresponding line in the original feature space would be the one where x_1 varies from $+\infty$ to $+\infty$ and x_2 is 0, as $x_2 = 0$ corresponds to the case where $x_1x_2 = 0$ in the new feature space. This line of maximum separation is shown in bold in figure 2.

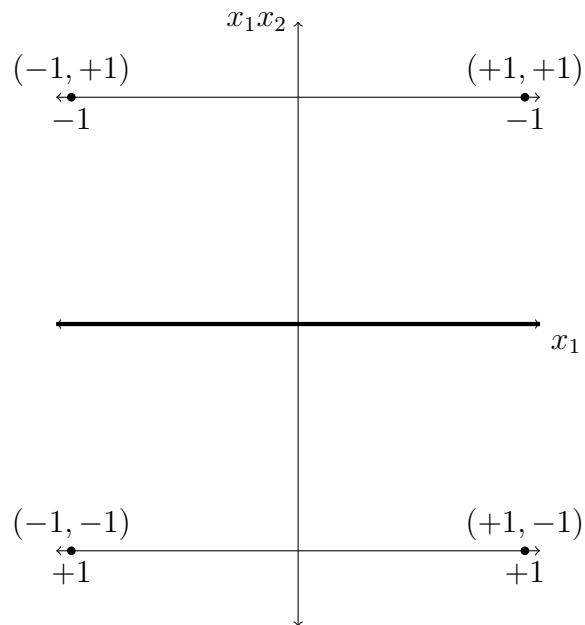


Figure 1: Separation of XOR function in feature space x_1 and x_1x_2 .

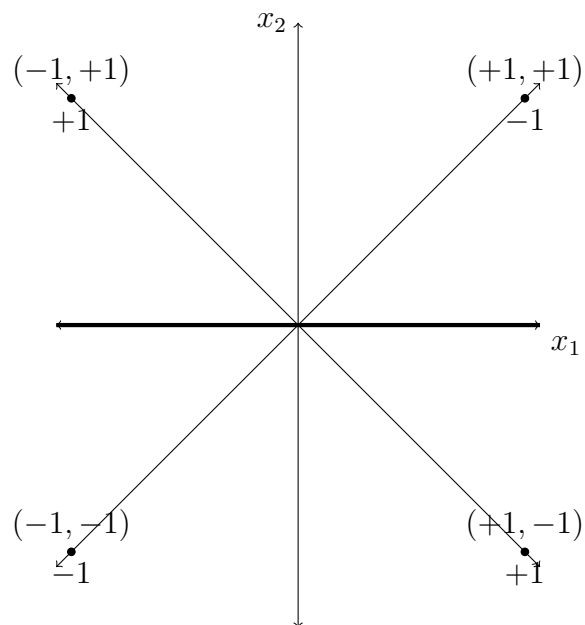


Figure 2: Separation of XOR function in original feature space x_1 and x_2 .

2. [10 points] Consider the following collection of points:

Point	coordinate	label	Point	coordinate	label
x_1	(0, 0)	+	x_5	(1, 0)	-
x_2	(0, 1)	+	x_6	$(\frac{1}{2}, \frac{\sqrt{3}}{2})$	-
x_3	(1, 1)	+	x_7	$(\frac{3}{2}, 0)$	-
x_4	$(\frac{1}{2}, 0)$	+	x_8	$(1, \frac{1}{2})$	-

Table 1: A collection of points

Suppose we have three training sets comprising of subsets of these points. We have

$$D_1 = \{x_1, x_2, x_3, x_5, x_7\}$$

$$D_2 = \{x_1, x_5, x_6, x_8\}$$

$$D_3 = \{x_3, x_4, x_5, x_7\}$$

- (a) [6 points] Give the maximum possible margin for D_1 , D_2 and D_3 .

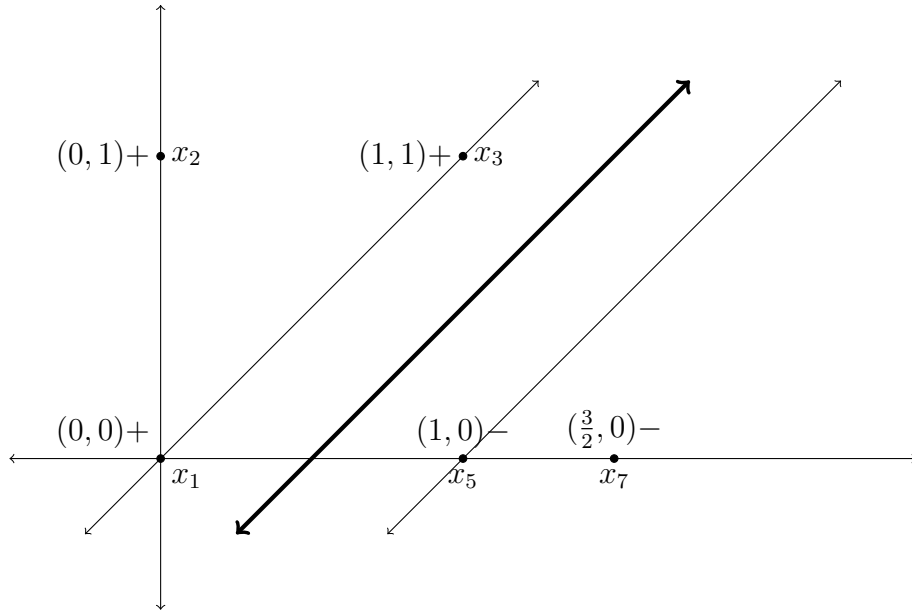


Figure 3: Points in set D_1 .

The line separating the positive and negative examples with the maximum margin for D_1 is shown by the bold line in figure 3. This bold line is in the middle of the widest strip separating the positive and negative labels. This widest strip is shown by the two lines through x_1, x_3 and through x_5 . The bold line intersects the horizontal axis at $(\frac{1}{2}, 0)$. All the separating lines are at 45° to the horizontal axis since one separating line passes through $(0, 0)$ and $(1, 1)$ and the three separating lines are parallel. So the distance between the bold line and the line through x_1 is $\frac{1}{2} \cos(45^\circ) = \frac{1}{2} \frac{1}{\sqrt{2}} = \frac{1}{2\sqrt{2}} = 0.3536$

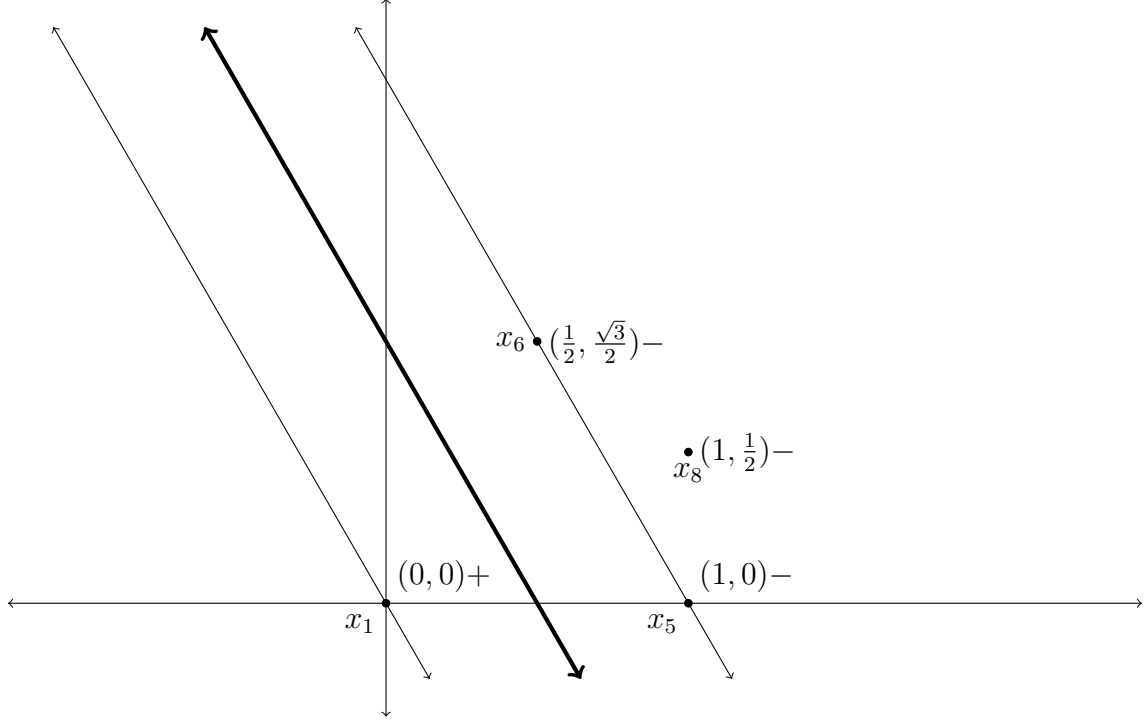


Figure 4: Points in set D_2 .

The widest strip separating the two labels in D_2 is shown in figure 4, with the solid line marking the separating margin between the two labels. One side of the widest strip passes through x_5 and x_6 and the other is parallel and passes through $(0,0)$. The slope of the separating lines is $\frac{\frac{\sqrt{3}}{2}-0}{\frac{1}{2}-1} = -\sqrt{3}$ (the slope of the line joining x_5 and x_6). Based on the value of the slope, we can say that the lines make an angle of 120° with the horizontal. Which means that $\sin(180 - 120) = \sin(60) = \frac{\sqrt{3}}{2} = \frac{\text{margin}}{0.5}$ since the distance between the horizontal intercepts of the bold line and the line through x_6 and x_5 is 0.5. This means that $\text{margin} = \frac{\sqrt{3}}{4} = 0.4330$.

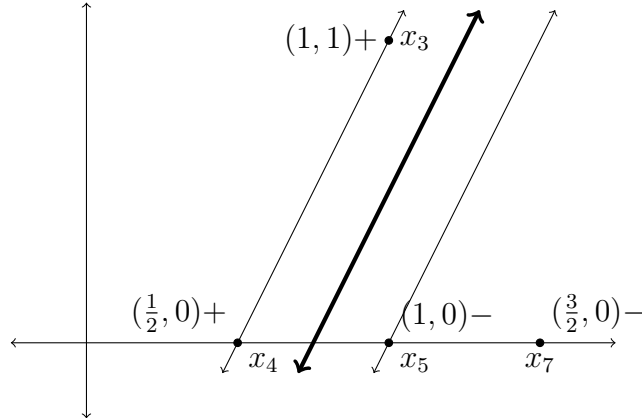


Figure 5: Points in set D_3 .

The widest strip separating the two labels in D_3 is shown in figure 5, with the solid line marking the separating margin between the two labels. The separating lines have a slope of $\frac{1-0}{1-\frac{1}{2}} = \frac{1}{\frac{1}{2}} = 2$. The angle that the separating lines make with the horizontal is $\arctan(2)$. Since the horizontal component of the margin is $\frac{1}{4}$, the margin is $\frac{1}{4} \sin(\arctan(2)) = \frac{1}{4} \sin(1.1071) = \frac{1}{4} \times 0.8944 = 0.2236$.

- (b) [2 points] What is the Perceptron mistake bound for these dataset. Which has the greatest Perceptron mistake bound.

Data set	Farthest point	R	Margin γ	Mistake Bound $\frac{R^2}{\gamma^2}$
D_1	x_7	1.5000	0.3536	17.9952
D_2	x_8	1.1180	0.4330	6.6667
D_3	x_7	1.5000	0.2236	45.0027

Table 2: Perceptron mistake bound for data sets D_1 , D_2 and D_3

Data set D_3 has the greatest Perceptron mistake bound.

- (c) [2 points] Rank the datasets in terms of “ease of learning”. Justify your answer. The mistake bound for a perceptron is the most number of mistakes a learning algorithm will make in finding a separator. So data sets with the minimum mistake bound should be easiest to learn as they have a lower limit on the number of mistakes that can be made. Therefore the data sets ranked by the ease of learning, starting with the easiest first are $\{D_2, D_1, D_3\}$.

2 Kernels

1. [15 points] If $K_1(\mathbf{x}, \mathbf{z})$ and $K_2(\mathbf{x}, \mathbf{z})$ are both valid kernel functions. In this question, you will prove that certain functions of kernels are valid kernels.

[Hint: For both the proofs below, use the the definition of a kernel as a dot product in a high dimensional space.]

- (a) [5 points] Show that the product of two kernels is a kernel. That is, show that K in the expression below is a kernel:

$$K(\mathbf{x}, \mathbf{z}) = K_1(\mathbf{x}, \mathbf{z})K_2(\mathbf{x}, \mathbf{z})$$

Since K_1 and K_2 are kernels, according to Mercer’s condition, for every finite set $\{x_1, x_2, \dots\}$, for any real valued choice of c_1, c_2, \dots , we have

$$\sum_i \sum_j c_i c_j K_1(\mathbf{x}_i, \mathbf{x}_j) \geq 0 \quad (1)$$

$$\sum_i \sum_j c_i c_j K_2(\mathbf{x}_i, \mathbf{x}_j) \geq 0 \quad (2)$$

For $K_1(\mathbf{x}, \mathbf{z})K_2(\mathbf{x}, \mathbf{z})$ to be a Kernel, the following needs to be true

$$\sum_i \sum_j c_i c_j K_1(\mathbf{x}_i, \mathbf{x}_j) K_2(\mathbf{x}_i, \mathbf{x}_j) \geq 0 \quad (3)$$

Let $c_j K_1(\mathbf{x}_i, \mathbf{x}_j) = k_j$. Since c_j can take any real values, it follows that k_j can take any real value. So we can now reduce the requirement of the inequality 3 to

$$\sum_i \sum_j c_i k_j K_2(\mathbf{x}_i, \mathbf{x}_j) \geq 0 \quad (4)$$

However requirement 4 is equivalent to inequality 2, which is true. This means that $K_1(\mathbf{x}, \mathbf{z})K_2(\mathbf{x}, \mathbf{z})$ is a valid kernel.

- (b) [10 points] Show that a polynomial over a kernel that is constructed using positive coefficients is a kernel. That is, if P is any polynomial with positive coefficients, show that K below is a kernel:

$$K(\mathbf{x}, \mathbf{z}) = P(K_1(\mathbf{x}, \mathbf{z}))$$

Hint: You may need show $K(\mathbf{x}, \mathbf{z}) = \alpha K_1(\mathbf{x}, \mathbf{z}) + \beta K_2(\mathbf{x}, \mathbf{z})$ is a valid kernel and use the conclusion in the previous question.

We have already shown in the previous answer that the product of two kernels is also a kernel. That is, if $K_1(\mathbf{x}, \mathbf{z})$ is a valid kernel, then $K_1(\mathbf{x}, \mathbf{z})K_1(\mathbf{x}, \mathbf{z})$. From this it follows that $(K_1(\mathbf{x}, \mathbf{z}))^3$ is a kernel since $(K_1(\mathbf{x}, \mathbf{z}))^3 = (K_1(\mathbf{x}, \mathbf{z})K_1(\mathbf{x}, \mathbf{z})) K_1(\mathbf{x}, \mathbf{z}) = K'_1(\mathbf{x}, \mathbf{z})K_1(\mathbf{x}, \mathbf{z})$, where $K'_1(\mathbf{x}, \mathbf{z}) = K_1(\mathbf{x}, \mathbf{z})K_1(\mathbf{x}, \mathbf{z})$. We already know that $K'_1(\mathbf{x}, \mathbf{z})K_1(\mathbf{x}, \mathbf{z})$ is a valid kernel since its the product of two kernels. This means that $(K_1(\mathbf{x}, \mathbf{z}))^{n+1}$ is a kernel if $(K_1(\mathbf{x}, \mathbf{z}))^n$ is a kernel. Using induction we can therefore show that $(K_1(\mathbf{x}, \mathbf{z}))^n$ is a kernel for any whole number n .

A polynomial over a kernel that is constructed using positive coefficients can be shown as a combination of many terms of the form $\alpha K_1(\mathbf{x}, \mathbf{z}) + \beta K_2(\mathbf{x}, \mathbf{z}) + \dots$. Using Mercer's condition, we need to show that

$$\sum_i \sum_j c_i k_j (\alpha K_1(\mathbf{x}_i, \mathbf{z}_j) + \beta K_2(\mathbf{x}_i, \mathbf{z}_j) + \dots) \geq 0 \quad (5)$$

$$\alpha \sum_i \sum_j c_i k_j K_1(\mathbf{x}_i, \mathbf{z}_j) + \beta \sum_i \sum_j c_i k_j K_2(\mathbf{x}_i, \mathbf{z}_j) + \dots \geq 0 \quad (6)$$

Since $K_1(\mathbf{x}, \mathbf{z})$ and $K_2(\mathbf{x}, \mathbf{z})$ and other terms are kernels, and coefficients α , β and others are positive, relation 6 is true. Which means that a polynomial over a kernel that is constructed using positive coefficients is a kernel.

2. [10 points] Given two examples $\mathbf{x} \in \mathbb{R}^2$ and $\mathbf{z} \in \mathbb{R}^2$, let

$$K(\mathbf{x}, \mathbf{z}) = 15 (\mathbf{x}^T \mathbf{z})^2 \exp(-\|\mathbf{x} - \mathbf{z}\|^2) \quad (7)$$

Prove that this is a valid kernel function.

According to Mercer's condition, for every finite set $\{x_1, x_2, \dots\}$, for any real valued choice of c_1, c_2, \dots , the following needs to be true

$$\sum_i \sum_j c_i c_j 15 (\mathbf{x}_i^T \mathbf{z}_j)^2 \exp(-\|\mathbf{x}_i - \mathbf{z}_j\|^2) \geq 0$$

We know that $\exp(-\|\mathbf{x}_i - \mathbf{z}_j\|^2)$ will always be positive. Let $\exp(-\|\mathbf{x}_i - \mathbf{z}_j\|^2) = k_{ij}$. So the inequality that must be true is

$$15k_{ij} \sum_i \sum_j c_i c_j (\mathbf{x}_i^T \mathbf{z}_j)^2 \geq 0$$

If we define $(\mathbf{x}_i^T \mathbf{z}_j)^2$ as $(K(\mathbf{x}_i, \mathbf{z}_j))^2$, then we can represent kernel $(K(\mathbf{x}_i, \mathbf{z}_j))^2$ as another kernel $K'(\mathbf{x}_i, \mathbf{z}_j)$ since we have already shown that any power of a kernel is also a kernel. We now have to show that

$$15k_{ij} \sum_i \sum_j c_i c_j K'(\mathbf{x}_i, \mathbf{z}_j) \geq 0$$

which is true by Mercer's condition since $K'(\mathbf{x}_i, \mathbf{z}_j)$ is a kernel.

3. (**For 6350 students**) [10 points] An valid kernel can always be expressed as inner product. Prove that the Gaussian kernel

$$K(\mathbf{x}, \mathbf{z}) = \exp\left(\frac{-\|\mathbf{x} - \mathbf{z}\|^2}{2\sigma^2}\right)$$

can be written down as the inner product of an feature space with infinite dimension. Hint: You may do some expansion and then show the middle factor can be expanded as a power series.

$$\begin{aligned}
K(\mathbf{x}, \mathbf{z}) &= \exp\left(\frac{-\|x - z\|^2}{2\sigma^2}\right) \\
&= \exp\left(\frac{-x^2 + 2xz - z^2}{\sigma^2}\right) \\
&= \exp\left(\frac{-x^2}{\sigma^2}\right) \exp\left(\frac{2xz}{\sigma^2}\right) \exp\left(\frac{-z^2}{\sigma^2}\right) \\
&= \exp\left(\frac{-x^2}{\sigma^2}\right) \left(1 + \frac{2xz}{\sigma^2} + \frac{\left(\frac{2xz}{\sigma^2}\right)^2}{2!} + \frac{\left(\frac{2xz}{\sigma^2}\right)^3}{3!} + \dots\right) \exp\left(\frac{-z^2}{\sigma^2}\right) \\
&= \exp\left(\frac{-x^2}{\sigma^2}\right) \left(1 + \frac{1}{\sqrt{1!}} \frac{\sqrt{2}x}{\sigma} \frac{1}{\sqrt{1!}} \frac{\sqrt{2}z}{\sigma} + \frac{1}{\sqrt{2!}} \left(\frac{\sqrt{2}x}{\sigma}\right)^2 \frac{1}{\sqrt{2!}} \left(\frac{\sqrt{2}z}{\sigma}\right)^2 + \dots\right) \exp\left(\frac{-z^2}{\sigma^2}\right)
\end{aligned}$$

This can be seen to be a product in a different dimension between

$$\phi(x) = \exp(-x^2) \left(1, \frac{\sqrt{2}x}{\sqrt{1!}\sigma}, \frac{\sqrt{2^2}x^2}{\sqrt{2!}\sigma^2}, \frac{\sqrt{2^3}x^3}{\sqrt{3!}\sigma^3}, \dots\right)$$

and

$$\phi(z) = \exp(-z^2) \left(1, \frac{\sqrt{2}z}{\sqrt{1!}\sigma}, \frac{\sqrt{2^2}z^2}{\sqrt{2!}\sigma^2}, \frac{\sqrt{2^3}z^3}{\sqrt{3!}\sigma^3}, \dots\right)$$

Which means $K(\mathbf{x}, \mathbf{z})$ can be written down as the inner product of a feature space with infinite dimension.

3 Experiments

In this question, you will implement the support vector machine (SVM) and a variant of random forest which combine SVMs and decision trees.

We will use two datasets for this question:

1. **semelion handwritten digits data:** This dataset contains 1593 handwritten digits from around 80 persons were scanned, stretched in a rectangular box 16x16 in a gray scale of 256 values. Our goal is implement svm in this data set to determine whether is a number 6.
2. **madelon:** This is one of five datasets used in the NIPS 2003 feature selection challenge. There are 2000 examples in training set and 600 examples in test set.

You may reuse your code in decision tree. If **you have problems in decision tree**, please contact with TA to get help. You may use Java, Python, Matlab, C/C++ for this assignment. If you want to use a different language, you must contact the instructor first. Any other language you may want to use **MUST** run on the CADE machines.

3.1 Support Vector Machines

1. [6 points] Implement the stochastic gradient descent algorithm for SVMs. Run it on the `handwriting` dataset with hyperparameter $C = 1$ and $\rho_0 = 0.01$. Report the accuracy in test set and training set.
2. [8 points] Run your SVM code on the `madelon` dataset and use 5-fold cross-validation to choose suitable parameters. At least attempt 6 different values for C and 3 different values for ρ . Report the average accuracy for each group of parameters. Report the accuracy in your test set as well as training set.

Hint: You should try out C in exponential steps, for example, $2^1, 2^{-1}, 2^{-2}, \dots$.

3. [6 points] Precision, recall and F_1 score are another metrics besides accuracy, which are useful if the dataset is unbalanced with respect to the positive and negative examples. To compute these quantities, you should count the number of true positives (that is, examples that your classifier predicts as positive and are truly positive), the false positives (i.e, examples that your classifier predicts as positive, but are actually labeled negative) and the false negatives (i.e., examples that are predicted as negative by your classifier, but are actually positive).

Denote true positives, false positive and false negative as TP, FP and FN respectively. The precision (p), recall (r) and f-value F_1 are defined as:

$$\begin{aligned} p &= \frac{TP}{TP + FP} \\ r &= \frac{TP}{TP + FN} \\ F_1 &= 2 \frac{p \cdot r}{p + r} \end{aligned}$$

Give precision, recall and F_1 score for your classifiers constructed in the previous two questions.

3.2 Ensemble of decision trees

Recall that a random forest is an ensemble based on bagging and decision tree. For bagging, we draw m samples *with replacement* from the training set. According to

$$\lim_{m \rightarrow \infty} (1 - \frac{1}{m})^m \rightarrow \frac{1}{e} \simeq 0.368$$

there are about 63.2 percent items may not appear that set. In random forest, we use this sampling method N times, to construct N training sets and grow N decision trees. Note that we build unpruned decision trees.

In each node, instead of using the best feature by ID3, we choose k features randomly and then use the ID3 heuristic to find the best feature to split on. Generally, $k = \log_2 d$ is a good choice where d is the number of features for our data.

Since there are N trees, there will be N predictions for each example. Generally, the final prediction is voted on by these trees. However, we would like to use SVM to combine these predictions for this question. Specifically, after growing the N decision trees, you should construct a new D consisting of transformed features. The feature transformation $\phi(x)$ is defined using the N trees as follows:

$$\phi(x) = [tree_1(x), tree_2(x), \dots, tree_N(x)]$$

In other words, you will build an N dimensional vector consisting of the prediction (1 or -1) of each tree that you created. Thus, you have a *learned* feature transformation.

You will finally train an SVM on this new dataset D .

1. [15 points] Using the method mentioned above, construct $N = 5$ decision trees for the **handwriting** dataset. For each node, select $k = \log_2 d = 8$ features randomly and then use the ID3 heuristic to find the best feature for splitting.

Train the SVM meta-classifier and report the accuracy for both training set and test set. (No cross-validation is needed but please choose good parameters for SVM, we will take out points for very low accuracy.)

2. [25 points] Implement same method on the **madelon** dataset. ($k = \log_2 d = 11$)
 - (a) [20 points] Try $N = 10, 30, 100$. For each N , report accuracy on training set and test set. (No cross-validation is needed, but please choose good parameters for the SVM.)
 - (b) [5 points] Choose the best N among those you have tried (you may try some new numbers). Report the accuracy, precision, recall and f_1 score on test set.