

## 1 K-nearest Neighbor (40pts)

**Solution.** 1. What is the role of the number of training instances to accuracy (hint: try different  $-limit$  and plot accuracy vs. number of training instances)? 2. What numbers get confused with each other most easily? 3. What is the role of  $k$  to training accuracy? 4. In general, does a small value for  $k$  cause overfitting or underfitting?

## 2 Cross Validation (30pts)

**Solution.** 1. What is the best  $k$  chosen from 5-fold cross validation with  $-limit$  500? 2. What is the best  $k$  chosen from 5-fold cross validation  $-limit$  5000? 3. Is the best  $k$  consistent with the best performance  $k$  in problem 1?

## 3 Bias-variance tradeoff (20pts)

**Solution.** Derive the bias-variance decomposition for  $k$ -NN regression in class. Specifically, assuming the training set is fixed  $S = (x_1, y_1), \dots, (x_n, y_n)$ , where the data are generated from the following process  $y = f(x) + \epsilon$ ,  $E(\epsilon) = 0$ ,  $\text{Var}(\epsilon) = 2$ .  $k$ -NN regression algorithm predict the value for  $x_0$  as  $\hat{h}_S(x_0) = \frac{1}{k} \sum_{l=1}^k y_{(l)}$ . Prove that  $\text{Err}(x_0)$ , where  $x_{(l)}$  is the  $l$ th nearest neighbor to  $x_0$ .  $\text{Err}(x_0)$  is defined as  $\text{Err}(x_0) = E((y_0 - \hat{h}_S(x_0))^2)$ .