

Due December 2, 2016.

## Question 1

Let us consider the training data  $\{(\mathbf{x}^1, y^1), \dots, (\mathbf{x}^n, y^n)\}$  where  $\mathbf{x}^i \in \mathbb{R}^p$  and  $y^i \in \{-1, +1\}$ . A soft-margin SVM solves

$$\begin{aligned} \arg \min_{\mathbf{w} \in \mathbb{R}^p, b \in \mathbb{R}, \xi \in \mathbb{R}^n} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i & (1) \\ \text{s. t. } & y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) \geq 1 - \xi_i \\ & \xi_i \geq 0 & \forall i \in \{1, \dots, n\} \end{aligned}$$

- (a) Is a soft-margin SVM more likely to overfit if  $C$  is large or small?

**Solution:** If  $C$  is large, more importance is given to the error on the training set, and the SVM is more likely to overfit.

- (b) Give one way of choosing  $C$  in practice.

**Solution:** By cross-validation.

- (c) What does this mean for a feature  $j$  if the solution  $w_j$  is close to 0?

**Solution:** That this feature is uninformative. The class won't depend on this feature. (Note: we're talking about a feature weight  $w_j$ , not a (support) vector coefficient  $\alpha_i$ .)

- (d) Give an interpretation of the two terms of Equation (1):  $\|\mathbf{w}\|^2$  and  $\sum_{i=1}^n \xi_i$ .

**Solution:**  $\|\mathbf{w}\|$  is the inverse of the margin.  
 $\sum_{i=1}^n \xi_i$  is the sum of slacks  $\xi_i$ , which quantify the error for each misclassified training point.