

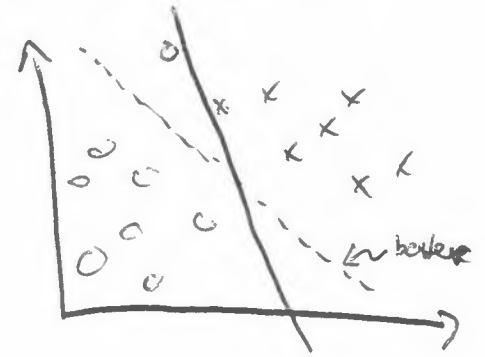
# SVM (Support Vector Machines)

- Finds separating, "soft", hyperplane w/ largest margin

- can kernelize to make non-linear

- typically used for classification,  $y = \{-1, 1\}$ , but can also formulate as regression

soft margin fixes the non-linearly separable cases, and the case where 1 or 2 errant training points drastically change the decision boundary



The problem is posed as a convex optimization problem:

$$\text{Max}_{w, b, \xi, M} \quad M \leftarrow \text{margin}$$

st:  $\|W\|_2^2 = 1$  ← ensures that  $b + w^T x_i$  gives signed distance to hyperplane

$$y_i(b + w^T x_i) \geq M(1 - \xi_i)$$

$$\forall i = 1, \dots, m$$

$$\forall i = 1, \dots, m$$

$$\xi_i \geq 0$$

allows for slack

$$\sum_{i=1}^m \xi_i \leq C$$

“soft” budget

ensures correct classification allowing for violations of margin

- actual optimization<sup>objective</sup> is a bit different
- can also rephrase w/ ERM using the hinge loss (no penalty for correct prediction, linear penalty for incorrect score)  
w/  $l_2$  regularization

- After working out the dual form of the problem, we notice that the entire objective function, as well as the prediction function can be written in terms of inner products of the data:  $\langle x^{(i)}, x^{(j)} \rangle (= x^{(i)T} x^{(j)})$   
thus, for more model expressiveness we can use a feature map  $\phi: \mathbb{R}^n \rightarrow \mathbb{R}^d$  ( $d > n$ )  
so that all occurrences of  $\langle x^{(i)}, x^{(j)} \rangle$  become  $\langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$ .

### Kernel Trick

- given a function  $k: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ , if  $\exists$  a feature map,  $\phi: \mathbb{R}^n \rightarrow \mathbb{R}^d$  st.  $k(x, y) = \langle \phi(x), \phi(y) \rangle$   
 $\forall x, y \in \mathbb{R}^n$ , then  $k$  is said to be a valid kernel.
- thus we can replace all occurrences of  $\langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$  w/  $k_{ij} = k(x^{(i)}, x^{(j)})$
- this has space + time savings
  - don't need to explicitly store  $\phi(x^{(i)})$  (which could be infinitely dimensional)
  - evaluating  $k(x^{(i)}, x^{(j)})$  is also typically  $O(n)$  time, whereas computing  $\langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$

is much worse

- Can formalize all ERM objectives of the form:

$$R(\|w\|) + \sum_{i=1}^n L(y_i, w^T x_i)$$

(Representer thm) ( $R$  is a nondecreasing function)

- can obtain useful kernel functions with Mercer theorem.

### Pros

- mathematically appealing, and is clear why it works well (i.e., has low generalization error) as compared to say the GBM algo

- memory efficient and predictions are fast, since prediction only depends on  $\alpha$ 's of "support vectors"

### Cons

- doesn't provide probability estimates (only scores)
- may take long to train w/ lots of data since entire  $n \times n$  kernel matrix must be computed