

SUPPORT VECTOR MACHINES

Idea: explicitly search for large-margin classifiers.

EASY CASE: DATA IS SEPARABLE

$$\min_{\vec{x}, b} \frac{1}{f(\vec{x}, b)}$$

$$\text{s.t.} \quad y_i (\langle x_i, \vec{x} \rangle + b) \geq 1$$

(How do we optimize this?)

What if there's noise? Give each point some slack, try to optimize combination of slack and margins.

$$\min_{\vec{x}, b, \epsilon} \frac{1}{f(\vec{x}, b)} + C \sum_i \epsilon_i$$

$$\text{subject to } y_i (\langle x_i, \vec{x} \rangle + b) \geq 1 - \epsilon_i$$

$$\epsilon_i \geq 0$$

(How do you compute margin there?)

FROM LARGE MARGIN TO SMALL WEIGHTS

$$d^+ = \frac{1}{\|w\|} w \cdot x^+ + b - 1$$

$$d^- = -\frac{1}{\|w\|} w \cdot x^- - b + 1$$

We can then compute the margin by algebra:

$$\begin{aligned}\gamma &= \frac{1}{2} [d^+ - d^-] \\&= \frac{1}{2} \left[\frac{1}{\|w\|} w \cdot x^+ + b - 1 - \left(-\frac{1}{\|w\|} w \cdot x^- - b + 1 \right) \right] \\&= \frac{1}{2} \left[\frac{1}{\|w\|} w \cdot x^+ - \frac{1}{\|w\|} w \cdot x^- \right] \\&= \frac{1}{2} \left[\frac{1}{\|w\|} (+1) - \frac{1}{\|w\|} (-1) \right] \\&= \frac{1}{\|w\|}\end{aligned}$$

$$\min_{\vec{x}, b, \xi} \quad \|\vec{x}\|^2 + C \sum_i \xi_i$$

$$\text{subject to } y_i (\langle x_i, \vec{x} \rangle + b) \geq 1 - \xi_i$$

$$\xi_i \geq 0$$

WHAT ARE THE VALUES OF THE OPTIMAL ξ_i ?

$$\xi_n = \begin{cases} 0 & \text{if } y_n(\mathbf{w} \cdot \mathbf{x}_n + b) \geq 1 \\ 1 - y_n(\mathbf{w} \cdot \mathbf{x}_n + b) & \text{otherwise} \end{cases}$$

In other words, the optimal value for a slack variable is *exactly* hinge loss on the corresponding example! Thus, we can write SVM objective as an *unconstrained* optimization problem:

$$\min_{\mathbf{w}, b} \underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{large margin}} + C \underbrace{\sum_n \ell^{(\text{hin})}(y_n, \mathbf{w} \cdot \mathbf{x}_n + b)}_{\text{small slack}}$$

As a consequence, Algorithm 24 finds a large-margin classifier!

Algorithm 24 HINGEREGULARIZEDGD($\mathbf{D}, \lambda, \text{MaxIter}$)

```
1:  $\mathbf{w} \leftarrow \langle 0, 0, \dots, 0 \rangle$  ,  $b \leftarrow 0$  // initialize weights and bias
2: for  $\text{iter} = 1 \dots \text{MaxIter}$  do
3:    $\mathbf{g} \leftarrow \langle 0, 0, \dots, 0 \rangle$  ,  $g \leftarrow 0$  // initialize gradient of weights and bias
4:   for all  $(\mathbf{x}, y) \in \mathbf{D}$  do
5:     if  $y(\mathbf{w} \cdot \mathbf{x} + b) \leq 1$  then
6:        $\mathbf{g} \leftarrow \mathbf{g} + y \mathbf{x}$  // update weight gradient
7:        $g \leftarrow g + y$  // update bias derivative
8:     end if
9:   end for
10:   $\mathbf{g} \leftarrow \mathbf{g} - \lambda \mathbf{w}$  // add in regularization term
11:   $\mathbf{w} \leftarrow \mathbf{w} + \eta \mathbf{g}$  // update weights
12:   $b \leftarrow b + \eta g$  // update bias
13: end for
14: return  $\mathbf{w}, b$ 
```

