

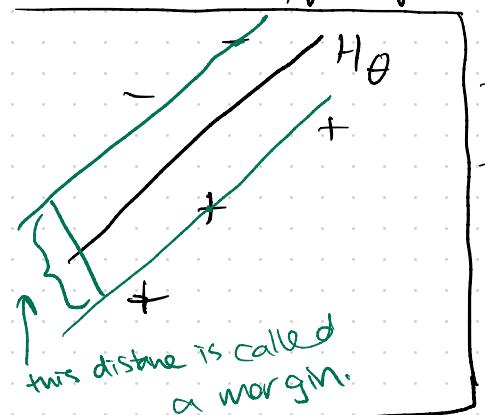
Recall from last lecture,

$$(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^D \times \underbrace{\{-1, 1\}}_N$$

binary classification

If data is linearly separable, that is

$\exists$  Hyperplane  $H = \{z \mid f_\theta(z) = 0\}$  with



$$f_\theta(z) = \langle z, \theta' \rangle + \theta_0 \text{ so}$$

that  $y_i = +1$  on one side,  
and  $y_i = -1$  on other side of  $H$ .

Linearly separable data admits a solution  
to the Hard margin SVM:

$$\max_{\theta \in \mathbb{R}^{D+1}} \left\{ r \geq 0 \mid \text{dist}(x_k, H_\theta) \geq r \right\}_{k=1, \dots, n}$$
$$\|\theta'\| = 1$$

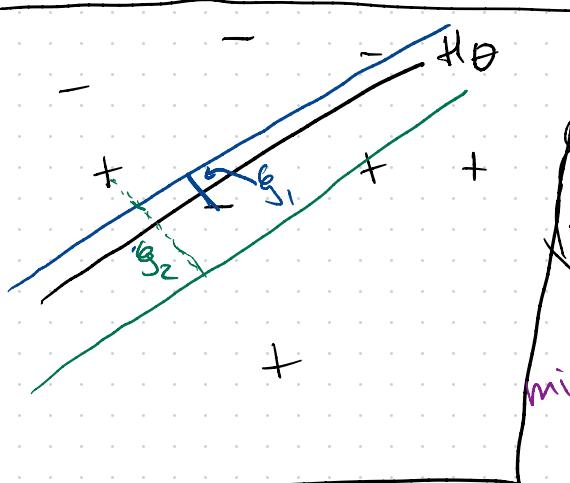
Dividing by  $r$ , equivalently, Hard margin SVM is

$$\min_{\theta \in \mathbb{R}^{D+1}} \left\{ \|\theta'\|^2 \mid y_k(\langle \theta', x_k \rangle + \theta_0) \geq 1 \right\}_{k=1, \dots, n}$$

## Soft margin SVM

still linear structure,  
but adds "optimal"  
band of error.

$C > 0$  is a fixed regularization  
variable.



$$\min_{\theta \in \mathbb{R}^{n+1}} \|\theta\|^2 + C \sum_{k=1}^n \xi_k$$

$\begin{cases} \xi_k \geq 0 \\ y_k(\langle x_k, \theta \rangle + \theta_0) \geq 1 - \xi_k \\ k = 1, \dots, n \end{cases}$

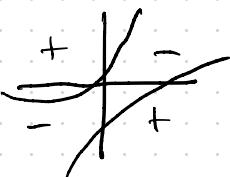
min  $\|\theta\| \rightarrow$  maximizes the margin

$\min \sum \xi_k \rightarrow$  minimizes loss

$C \rightarrow$  controls trade-off.

Benefits of SVM: Most data is not linearly separable.

Problem with SVM: Data might have non linear structure.



Solution: add a kernel function

$$K(x, y) = \langle \phi(x_k), \phi(x_e) \rangle$$

to reduce non-linear to linear problems

# Dual SVM

## Algorithm

1. Input: • Training data  $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^D \times \{-1, 1\}$

- Kernel map  $K(x_i, y)$

- Regularization parameter  $C$

2. Output:  $f: \mathbb{R}^D \rightarrow \{-1, 1\}$  of the form

$$f(x) = \text{sgn}(\langle \theta^*, \phi(x) \rangle + \theta_0)$$

3. Compute Kernel matrix  $G = (K(x_i, x_j))_{i,j=1}^n$

4. Solve Dual SVM problem by finding  $\alpha \in \mathbb{R}^n$

$$\max_{\alpha \in \mathbb{R}^n} \left\{ -u^T Gu + \sum_{k=1}^n \alpha_k \mid \begin{array}{l} \sum_{k=1}^n \alpha_k y_k = 0 \\ 0 < \alpha_i \leq C \end{array} \right\}$$

$$\text{where } u = \left( \frac{1}{2} \alpha_k y_k \right)_{i=1}^n$$

5. Define  $\gamma(x) = \frac{1}{2} \sum_{i=1}^n y_i \alpha_i K(x_i, x)$

6. Set  $b = \text{median}\{y_k - \gamma(x_k) \mid \alpha_k \neq 0\}$

7. Return  $f(x) = \text{sgn}(\gamma(x) + b)$ .

Notebook 5  $\rightarrow$  classification

1. we'll use kernel map  $k(x, y) = \langle x, y \rangle$   
the standard inner product

4.  $\star$  JuMP  $\rightarrow$  package for model

$\star$  NLopt  $\rightarrow$  non linear optimization

algorithm: AUGLAG = "Augmented Lagrangian"

$\star$  globally convergent augmented  
Lagrangian algorithm.

local optimizer:  $\underbrace{\text{LD-LBFGS}}$   
Low storage algorithm  
minimizes memory at each step

\* Maximization function looks different. we use the fact that if  $\alpha$  maximizes a function  $L$

then  $\alpha$  also maximizes  $\gamma L$ , so multiply Dual Vector problem by  $\gamma$ .

End

Compares Dual SVM as above to a Neural Network

Recall activation functions:

$$\text{ReLU } \sigma(z) = (\max(0, z_i))_{i=1}^k$$

$$\text{Sigmoid } \sigma(z) = \left( \frac{1}{1 + \exp(-z_i)} \right)_{i=1}^k$$

$$\text{Softmax } \sigma(z) = \left( \frac{\exp(z_i)}{\sum_{j=1}^k \exp(z_k)} \right)_{i=1}^k$$

\* On Neural Networks vs. SVM

Both can do linear +

non linear classification

NN  $\rightarrow$  eventually achieves higher accuracy  
SVM  $\rightarrow$  faster to train, more reliable since guarantees convergence to global minimum

After Notebook 5

Soft Margin  $\rightarrow$  dual

plus why steps 5 + 6 of algorithm work-

$$\min_{\theta \in \mathbb{R}^{D+1}, \beta \in \mathbb{R}^n} \max_{\alpha \in \mathbb{R}^n, \alpha_k \beta_k \geq 0}$$

$$\|\theta'\|^2 + C \sum_{k=1}^n \xi_k - \sum_{k=1}^n \beta_k \xi_k$$

$$- \sum_{k=1}^n \alpha_k \left[ y_k (\langle \theta', x_k \rangle + \beta_0) - (1 - \xi_k) \right]$$

= soft margin over

$$y_k (\langle \theta', x_k \rangle + \beta_0) - (1 - \xi_k) \geq 0$$

↑  
this is called  
the Lagrange function

$\mathcal{L}(\theta', \beta_0, \beta, \alpha, \xi)$   
for soft margin SVMs.

KKT = Karush-Kuhn-Tucker conditions  
1939

1951

KKT  $\Rightarrow$  min max problem above is solved  
when

$$\frac{\partial \mathcal{L}}{\partial \theta^1} = \frac{\partial \mathcal{L}}{\partial \theta_0} = \frac{\partial \mathcal{L}}{\partial \gamma} = 0.$$

\* The "kernel trick" comes from the fact that the conditions for these 3 conditions to be 0 are independent of  $\theta$  and  $\gamma$ , so we can first solve maximum function!

$$\frac{\partial \mathcal{L}}{\partial \theta^1} = 2\theta^1 - \sum_{k=1}^n \alpha_k y_k x_k = 0$$

$$\boxed{\theta_*^1 = \frac{1}{2} \sum_{k=1}^n \alpha_k y_k x_k}$$

For algorithm step 5  $\boxed{5}$   $\gamma(x) = \langle \theta_*^1, x \rangle$  = portion of  $x$  in the optimal direction.

Rest of KKT gives [b] of algorithm

$$\frac{\partial \mathcal{L}}{\partial \theta_0} = -\sum_{k=1}^n \alpha_k y_k = 0$$

$$\Rightarrow \sum_{k=1}^n \alpha_k y_k = 0$$

$$\frac{\partial \mathcal{L}}{\partial \theta_j} = \left( C - (\alpha + \beta) \right)^n = 0$$

$$\Rightarrow C = \alpha + \beta = 0 \leq \alpha_k \leq C.$$

Plugging back into  $\mathcal{L}$  the fact that  $\alpha + \beta = C$   
 $\sum \alpha_k y_k = 0$

$$\mathcal{L} = \|\theta'\|^2 - \sum_{k=1}^n \alpha_k (y_k \langle \theta', x_k \rangle + \theta_0) - 1$$

Maximizing over  $\alpha$  means  $(y_k \langle \theta', x_k \rangle + \theta_0) - 1 \leq 0$   
 $\Rightarrow \alpha_k = 0$

$$\mathcal{L} = \|\theta'\|^2 - \sum_{\alpha_k > 0} \alpha_k (y_k \langle \theta' + x_k \rangle + \theta_0) - 1$$

Minimizing  $\mathcal{L}$  over  $\theta \Rightarrow$  maximize  $\alpha_k = C$

$$\mathcal{L} = \|a\|^2 - C \sum_{\alpha_k > 0} \underbrace{g_k(\langle \theta_0^\top x_k \rangle + \theta_0) - 1}_{\geq 0}$$

So  $\theta_0^*$  minimizes  $\sum_{\alpha_k > 0} |1 - g_k(\langle \theta_0^\top x_k \rangle + \theta_0)|$

Multiplying by  $(y_k) = 1$ , and  $y_k^2 = 1$

$\theta_0^*$  minimizes  $\sum_{\alpha_k > 0} |y_k - \langle \theta_0^\top x_k \rangle - \theta_0|$

By exercise '3.7' minimizer  $\hat{\theta}_0$   
 $\text{median}_{\substack{1 \leq k \leq n \\ \alpha_k > 0}} \underbrace{\{y_k - \langle \theta_0^\top x_k \rangle\}}_{\parallel}$   
 $y_k - \hat{\gamma}(x_k)$