

1. project proposal due on Mar 16 on grade scope

2. HW 2:

problem 1.2. normalized margin =  $y \cdot \bar{f}_t(x) \leftarrow f_t$

problem 4:  $F_m \rightarrow F_n$

$\sigma: \sigma(w_i \cdot x)$

$\sigma(v_1, \dots, v_d) = (\sigma(v_1) \dots \sigma(v_d))$

$\langle w_i, x \rangle \rightarrow w_i \cdot x$

$l_2$ -SVM.

min  $w$   $\|w\|_2 \rightarrow \frac{1}{2} \|w\|_2^2$

s.t.  $\forall i, y_i \cdot \langle w, x_i \rangle \geq 1$



$\alpha \hat{w} \quad w = \alpha \cdot \hat{w} \quad \alpha > 0, \|\hat{w}\| = 1$

min  $\alpha$   
s.t.  $\forall i, \alpha > 0, \|\hat{w}\| = 1$

s.t.  $\forall i, y_i \cdot \alpha \cdot \langle \hat{w}, x_i \rangle \geq 1$



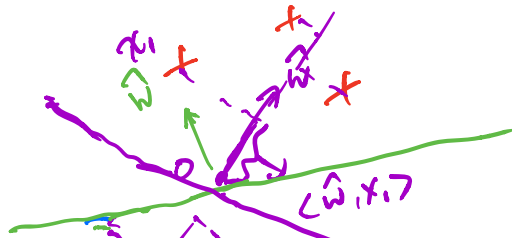
max  $\frac{1}{\alpha}$   
s.t.  $\forall i, \alpha > 0, \|\hat{w}\| = 1$

s.t.  $\forall i, y_i \cdot \langle \hat{w}, x_i \rangle \geq \frac{1}{\alpha}$

for any fixed  $\hat{w}$ , the optimal choice of  $d$  is such that

$$\frac{1}{d} = \min_i y_i \langle \hat{w}, x_i \rangle$$

$$\Leftrightarrow \max_{\hat{w}: \|\hat{w}\|=1} \left[ \min_i y_i \langle \hat{w}, x_i \rangle \right]$$



Generalization error bounds for  $l_2$ -SVMs ( $l_2$ -bounded linear predictors)


Thm: fix  $B_2, R_2 > 0$ ,  $S = (x_1, y_1) \dots (x_m, y_m) \stackrel{i.i.d.}{\sim} D$ .

$D$  supported on  $\{x \in \mathbb{R}^d: \|x\|_2 \leq R_2\} \times \{\pm 1\}$ . Fix  $\theta \in (0, B_2 R_2]$

Then, w.p.  $1-\delta$ , for all  $w: \|w\|_2 \leq B_2$ :

$$P_D (y \langle w, x \rangle \leq 0) \leq P_S (y \langle w, x \rangle \leq \theta) + O\left(\frac{B_2 R_2}{\theta} \sqrt{\frac{\ln 1/\delta}{m}}\right)$$

Pf: Same strategy as the  $l_1$  ( $l_\infty$  margin bound):

- introduce ramp loss 
- uniform concentration of ramp losses
- contraction inequality of Rademacher complexity

Pf comes down to show: given  $S = (x_1, y_1) \dots (x_m, y_m)$ .

s.t.  $\forall i, \|x_i\|_2 \leq R_2$ .

$$y = \{ \underline{m}_w : \|w\|_2 \leq \beta_2 \}$$

$$m_w(x, y) = y \langle w, x \rangle.$$

$$\text{Bound} \quad \text{Rad}_S(y) = \frac{1}{m} \mathbb{E}_{\sigma \sim U(\pm 1)^m} \sup_{w: \|w\|_2 \leq \beta_2} \sum_{i=1}^m \sigma_i \langle w, x_i \rangle$$

$$= \frac{1}{m} \mathbb{E}_{\sigma \sim U(\pm 1)^m} \sup_{w: \|w\|_2 \leq \beta_2} \left[ \sum_{i=1}^m \sigma_i \langle w, x_i \rangle \right]$$

$$\downarrow m$$

$$\langle w, \sum_{i=1}^m \sigma_i x_i \rangle$$

Cauchy Schwarz:

$$\langle \alpha, \beta \rangle \leq \|\alpha\|_2 \cdot \|\beta\|_2$$

$$\leq \|w\|_2 \left\| \sum_{i=1}^m \sigma_i x_i \right\|_2$$

$$\leq \beta_2 \left\| \sum_{i=1}^m \sigma_i x_i \right\|_2$$

$$(\mathbb{E} z)^2 \leq \mathbb{E}[z^2]$$

$$\leq \frac{\beta_2}{m} \mathbb{E}_{\sigma} \left[ \left\| \sum_{i=1}^m \sigma_i x_i \right\|_2 \right]$$

$$\leq \frac{\beta_2}{m} \sqrt{\mathbb{E}_{\sigma} \left\| \sum_{i=1}^m \sigma_i x_i \right\|_2^2}$$

$$= \frac{\beta_2}{m} \sqrt{\mathbb{E}_{\sigma} \left[ \left\langle \sum_{i=1}^m \sigma_i x_i, \sum_{j=1}^m \sigma_j x_j \right\rangle \right]}$$

linearity of inner product

$$= \frac{\beta_2}{m} \sqrt{\mathbb{E}_{\sigma} \left[ \sum_{i=1}^m \sum_{j=1}^m \sigma_i \sigma_j \langle x_i, x_j \rangle \right]}$$

$$= \frac{\beta_2}{m} \sqrt{\sum_i \sum_j \mathbb{E}_{\sigma} [\sigma_i \sigma_j] \langle x_i, x_j \rangle}$$

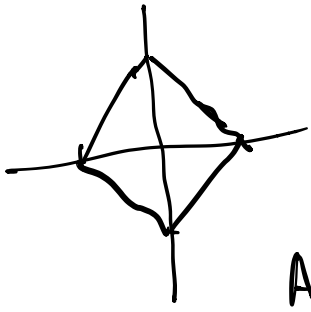
$$\begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases}$$

$$= \frac{B_2}{m} \cdot \sqrt{\sum_{i=1}^m \langle x_i, x_i \rangle} = \frac{B_2}{m} \cdot \sqrt{\|x\|_2^2} \leq R_2$$

$$\leq \frac{B_2}{m} \cdot \sqrt{m \cdot R_2^2} = \frac{B_2 R_2}{\sqrt{m}} \quad *$$

Comparison b/w  $l_1/l_\infty$  &  $l_2/l_2$  margin bounds

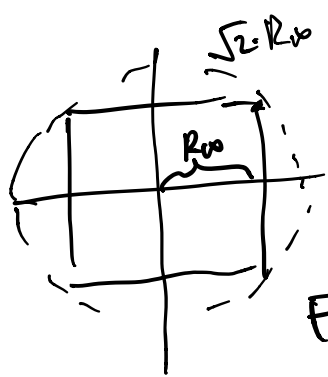
	constraints on $x$	constraint on $w$	error bound
$l_1/l_\infty$	$\frac{\ x\ _\infty \leq R_\infty}{\max_i  x_i }$	$\ w\ _1 \leq B_1$	$\tilde{O}\left(\frac{B_1 R_\infty}{\theta} \sqrt{\frac{1}{m}}\right)$
$l_2/l_2$	$\ x\ _2 \leq R_2$	$\ w\ _2 \leq B_2$	$\tilde{O}\left(\frac{B_2 R_2}{\theta} \sqrt{\frac{1}{m}}\right)$



incomparable in general

Applying  $l_2/l_2$  generalization bound to  $l_1/l_\infty$  setting:

$$\|x\|_\infty \leq R_\infty \Rightarrow \|x\|_2 \leq \underline{R_2} ?$$



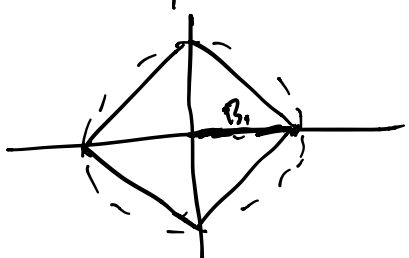
$$R_2 = \sqrt{d} \cdot R_\infty$$

$$B_2 = B_1$$

Fact  $\|w\|_2 \leq \|w\|_1$

will get a bound in terms of

$$R_2 B_2 = \sqrt{d} \cdot R_\infty \cdot B_1 \quad \text{a factor of } \sqrt{d}$$



worse than the original  $l_1/l_\infty$  bound.

Applying  $l_1/l_\infty$  bound to the  $l_2/l_2$  setting.

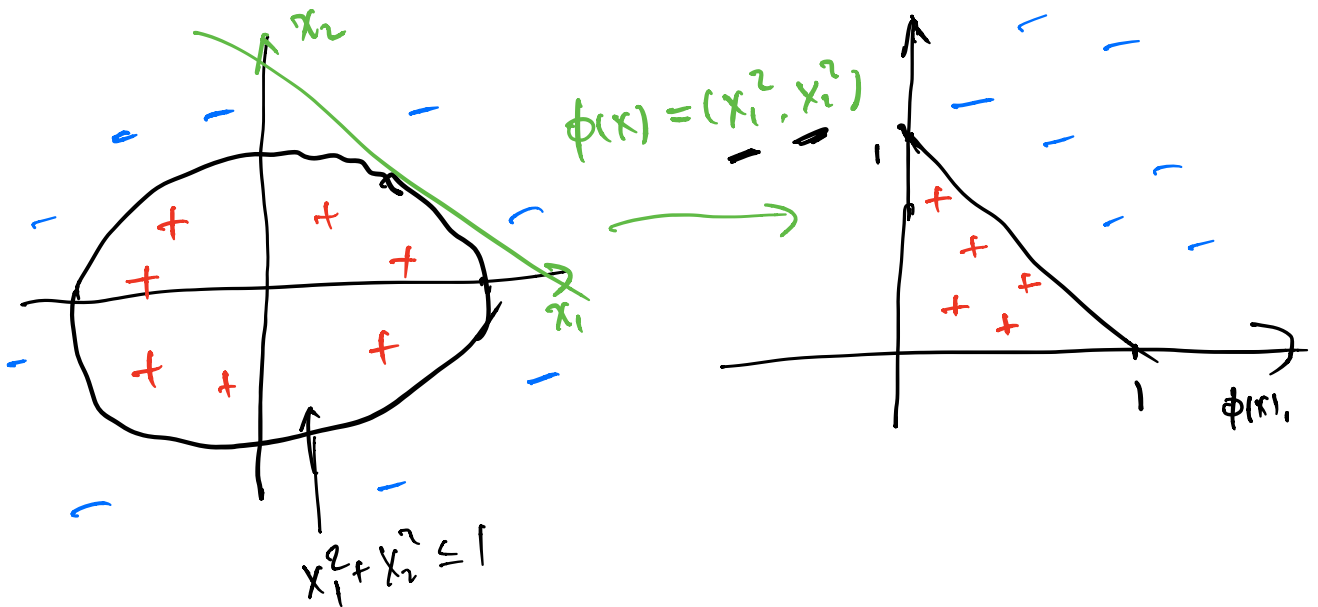
(will be  $\sqrt{d}$  factor worse, left as exercise).

Copying of linear non separability in SVMs.

Ideas:

1. introduce nonlinear feature maps (basis fns)
2. relax the SVM formulation: allowing some samples to be incorrectly classified.

1.



- define  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^m$   $(x_i, y_i) \rightarrow (\phi(x_i), y_i)$
- solve SVM on  $(\phi(x_i), y_i)$ 's.  $\Rightarrow \hat{w} \in \mathbb{R}^m$
- Final predictor:  $\text{sign}(\langle \phi(x), \hat{w} \rangle)$
- there are SVM training algs that has time complexity independent of  $m$ .

(If  $\langle \phi(x), \phi(y) \rangle$  can be evaluated in time independent of  $m$ ) . — kernel trick

## 2. soft margin SVMs.

$$\min_{w, \xi_1, \dots, \xi_m} \frac{\lambda}{2} \|w\|_2^2 + \sum_{i=1}^m \xi_i$$

$$\text{s.t. } \forall i. \quad y_i \langle w, x_i \rangle \geq 1 - \xi_i$$

$\lambda \downarrow$  : more likely to classify more correctly.  
 $\|w\|$  will likely to be large.

fix any  $w$ : the optimal choices of  $\xi_i$ 's are:

$$\left. \begin{array}{l} \xi_i \geq 1 - y_i \langle w, x_i \rangle \\ \xi_i \geq 0 \end{array} \right\} \Leftrightarrow \xi_i \geq \max(0, 1 - y_i \langle w, x_i \rangle)$$

the optimal choices of  $\xi_i$ 's =  $\max(0, 1 - y_i \langle w, x_i \rangle)$ .

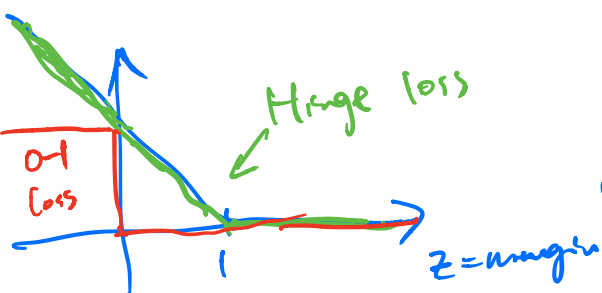
$$\min_w \frac{\lambda}{2} \|w\|_2^2 + \sum_{i=1}^m \max(0, 1 - y_i \langle w, x_i \rangle)$$

$\downarrow$   
complexity

$\downarrow$   
empirical risk

$$= \phi(y_i \langle w, x_i \rangle)$$

where  $\phi(z) = \max(0, 1 - z)$ .



Regularized loss minimization:

$$\min_w \lambda \cdot R(w) + \sum_{i=1}^m \phi(f_w(x_i), y_i)$$

$\downarrow$  free to change.

Notable examples:

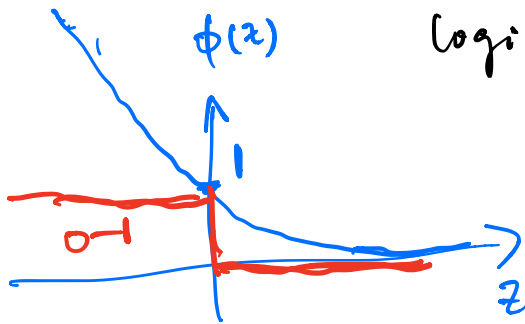
$R(w) = \frac{1}{2} \|w\|_2^2$ ,  $\phi(f_w(x_i), y_i) = (f_w(x_i) - y_i)^2$   
 $f_w(x) = \langle w, x \rangle$ .

$\lambda = 0$ : Ordinary least squares.

$\lambda > 0$ : Ridge regression.  
 $\ln(1+x) \approx x$

$R(w) = \|w\|_1$ : Lasso

$R(w) = \frac{1}{2} \|w\|_2^2$ .  $\phi(f_w(x_i), y_i) = \log_2(1 + e^{-y_i f_w(x_i)})$   
 $f_w(x) = \langle w, x \rangle$ . logistic loss.



logistic regression

Next class:

regularized loss minimization's soln's generalization performance from stability perspective