

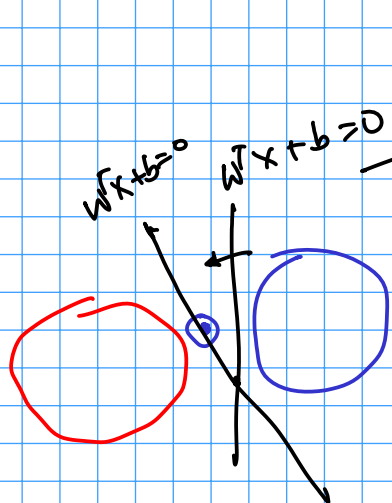
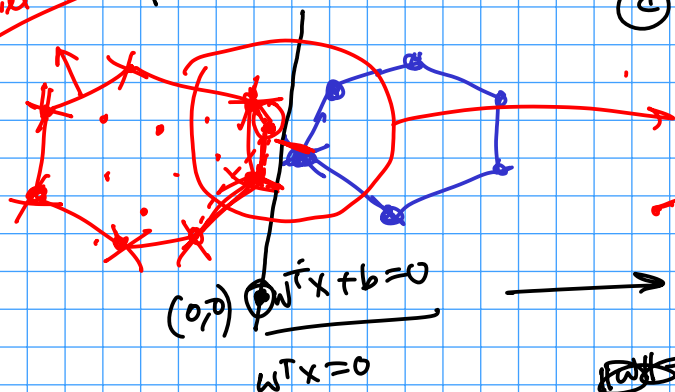
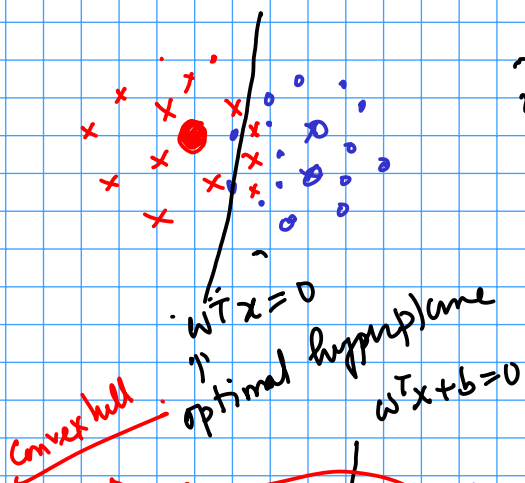
Stability or Robustness of Machine Learning

Focus on regression/classification.

$$\min_w \|w\|^V + C \sum_{i=1}^{|D|} \max(0, 1 - w^T x_i)$$

Q: In addition to the new point, do you expect that

- (a) the line $w^T x = 0$ will significantly shift
- (b) _____ Not shift
- (c) _____ shift to a small extent



$$\min_w \left(\|w\|^V + b \right) + C \sum_{i=1}^{|D|} \max(0, 1 - y_i (w^T x_i + b))$$

$$\min_w \|w\|^V + C \sum_{i=1}^{|D|} \max(0, 1 - y_i (w^T x_i))$$

Robustness \Rightarrow stability

\rightarrow additional point does not change learned params to a large extent

$$\frac{1}{n} \|w_1 - w_2\| + |b_1 - b_2|$$

\rightarrow overfitting \rightarrow model the noise as well.

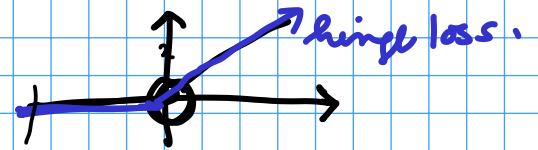
\rightarrow Differential privacy will be breached

$$A(I) \approx A(I \cup e')$$

in prob. sense

* $a_+ = \max(0, a) = \text{ReLU}(a)$

* $a_+ \rightarrow$ hinge loss on a



$$\min_w \|w\|^2 + C \sum_{i=1}^{|D|} (1 - y_i \cdot w^T x_i)_+ = \max(0, 1 - y_i \cdot w^T x_i)$$

$$\lambda = 1/C \cdot \min_w \lambda \|w\|^2 + \sum_{i=1}^{|D|} (1 - y_i \cdot w^T x_i)_+ \Rightarrow w^*(D, \lambda)$$

$D = \{(x_i, y_i)\}$

$L_D(w) = \sum_{i=1}^{|D|} (1 - y_i \cdot w^T x_i)_+$
 $L_{DUE}(w) = \sum_{i=1}^{|D|} (1 - y_i \cdot w^T x_i)^2$

Soham De

$$w^*(D, \lambda) = \underset{w}{\operatorname{argmin}} L_D(w)$$

$$w^*(D_{UE}, \lambda) = \underset{w}{\operatorname{argmin}} L_{DUE}(w)$$

(x, y)

If the learner/SVM is stable or robust
 what we expect

$$\|w^*(D, \lambda) - w^*(D_{UE}, \lambda)\| \rightarrow 0$$

$$\lambda \|w^*(D, \lambda)\|^2 \neq \lambda \|w^*(D_{UE}, \lambda)\|^2$$

$|D| \rightarrow 1000000$ 30000 classification

$$\lambda \|w\|^2 + \sum_{i=1}^{|D|} \ell(i, w)$$

\downarrow convex

$\ell(i, w) = \max(0, 1 - y_i \cdot w^T x_i)$
 Regression
 $\ell(i, w) = (y_i - w^T x_i)^2$

If ℓ is convex

$$\Rightarrow \ell(i, w) = \ell(i, w') + \left(\frac{\partial \ell}{\partial w'} \right)^T (w - w') + (w - w')^T \left[\frac{\partial^2 \ell}{\partial w'^2} \right] (w - w')$$

$$\ell(i, w) \geq \ell(i, w') + \left(\frac{\partial \ell}{\partial w'} \right)^T (w - w')$$

In theory $\min_w \lambda \|w\|^2 + \frac{1}{|D|} \sum_{i=1}^{|D|} \ell(i, w) \quad \lambda \approx \lambda_c / |D|$

In practice $\min_w \lambda \|w\|^2 + \frac{1}{|D|} \sum_{i=1}^{|D|} \ell(i, w)$

Stability + Generalization by Olivier Bosquet
2002
← Understanding HL
141-144.