

FOUNDATIONS OF MACHINE LEARNING [FOML]

Mohri et al.

[select portions]

$$\mathcal{D} = \left\{ (x^{(i)}, y^{(i)}) \mid x^{(i)} \in \mathcal{X}, y^{(i)} \in \{-1, 1\} \right. \\ \left. i \in [N] \right\}$$

$$\mathcal{D} \sim \mathcal{P}^{(N)}$$

Assume that there exist
 ω^* such that

$$\text{sign}((\omega^*)^T x^{(i)}) = y^{(i)}$$

$$i \in [N].$$

$$\|x^{(i)}\| \leq R$$

$$\omega^{(i)} \neq 0$$

$$i = 1, \dots, N$$

$$\Rightarrow y^{(i)} (\omega^{*T} x^{(i)}) > 0$$

$$\gamma = \min_i \frac{y^{(i)} (\omega^*)^T x^{(i)}}{\|\omega^*\|}$$

$\omega^{(n)}$ be the current estimate.

Let $(x^{(n)}, y^{(n)}) \in \mathcal{D}$ such that

$$\text{sign}((\omega^{(n)})^T x^{(n)}) \neq y^{(n)}$$

$$\Rightarrow y^{(n)} (\omega^{(n)T} x^{(n)}) < 0$$

$$\omega^{(n+1)} = \begin{cases} \omega^{(n)} + y^{(n)} x^{(n)} [\text{update}] \\ \omega^{(n)} & \text{otherwise} \end{cases}$$

On a linearly separable dataset
 the perceptron algorithm
 terminates after making
 at most $\frac{R^2}{\gamma^2}$ updates

$$\frac{R^2}{\gamma^2} = \frac{R^2 \|w^*\|^2}{\min_i \{y_i (w^{*T} x^{(i)})\}^2}$$

\mathcal{D} be a sample of size N .

Linearly separable [There exists $\gamma > 0$].

Let Perceptron Algorithm return
 a Classifier $h_{\mathcal{D}}^{(p)}$.

$$R(h_{\mathcal{D}}^{(p)}) = \mathbb{P}_{X, Y \sim \mathbf{P}} (h_{\mathcal{D}}^{(p)}(X) \neq Y)$$

An algorithm A acting on a sample \mathcal{D} of size m return a classifier $h_{\mathcal{D}}^{(A)}$.

The leave one out error of A on \mathcal{D} is

$$\bar{R}_{\mathcal{D}}^{\text{LOO}}(A) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{h_{\mathcal{D}^{(i)}}^{(A)}(x^{(i)}) \neq y^{(i)}\}}$$

$$\mathcal{D}^{(i)} : \mathcal{D} - \{(x^{(i)}, y^{(i)})\}$$

$$E_{\mathcal{D} \sim \mathcal{P}^{(m)}} \bar{R}_{\mathcal{D}}^{\text{LOO}}(A)$$

Expected LOO error. on a random sample of size m .

$$E_{\mathcal{Q} \sim P(m)} \bar{R}_{\mathcal{Q}}^{LOO}(A)$$

$$= E_{\mathcal{Q} \sim P(m-1)} R(h_{\mathcal{Q}}^{(A)})$$

Reading : Sec 5.2.4.
 Lemma 5.3.
 FOML

Perceptron

$$E_{\mathcal{D} \sim P(N)} R(h_{\mathcal{D}}^{(P)})$$

$$= E_{\mathcal{D} \sim P(N+1)} \bar{R}_{\mathcal{D}}^{\text{LOO}}(P)$$

$$E_{\mathcal{D} \sim P(N)} R(h_{\mathcal{D}}^{(P)})$$

$$\leq E_{\mathcal{D} \sim P(N+1)} \frac{\min(M(\mathcal{D}), \frac{R^2(\mathcal{D})}{\gamma^2(\mathcal{D})})}{N+1}$$

[Theorem 8.9]
FOML]

Generalization error of SVM

$$\min_{\{\omega, b\}} \quad \frac{1}{2} \|\omega\|^2$$

$$y_i (\omega^T x^{(i)} + b) \geq 1 \quad i=1, \dots, N$$

$$\omega^* = \sum_{i=1}^N \lambda_i y_i x^{(i)}$$

$$\lambda_i \geq 0$$

$$SV = \{i \mid \lambda_i > 0\}$$

let $b=0$

$$\min_{\{\omega, b\}} \quad \frac{1}{2} \|\omega\|^2$$

$$y_i (\omega^T x^{(i)}) \geq 1 \quad i=1, \dots, N$$

Optimum attained
at ω^*

$$\Rightarrow \min_{\omega, \gamma} \quad \gamma$$

$$y_i (\omega^T x^{(i)}) \geq \gamma \|\omega\| \quad i \in [N]$$

optimum (γ^*, ω^*)

$$\gamma^* = \frac{1}{\|\omega^*\|}$$

$$h_{\mathcal{D}}^{\text{SVM}}(x) = \text{sign}(\omega^{*T} x)$$

$$R(h_{\mathcal{D}}^{\text{SVM}}) = P(\text{sign}(\omega^{*T} x) \neq y)$$

$$(x, y) \sim P$$

[Theorem 5.4 FOM2]

$$E_{\mathcal{D} \sim P(N)} R(h_{\mathcal{D}}^{\text{svm}})$$

$$\leq E_{\mathcal{D} \sim P(N+1)} \frac{|SV(\mathcal{D})|}{N+1}$$

$|SV(\mathcal{D})|$ = Number of Support vectors

Proof: Take $\mathcal{D} \sim P(N+1)$ to be separable. solve SVM.

Argue $\bar{R}_{\mathcal{D}}^{\text{Loo}}(\text{SVM}) \leq \frac{|SV(\mathcal{D})|}{N+1}$

Weak: Average generalization error.

Relate Training set error to generalization error.

If Training set error is 0 then is the generalization error 0?

$a \leq Z \leq b$ $Z_i, \text{iid } P$ $E(Z) = \mu$

$$P\left(\left|\frac{1}{n} \sum_{i=1}^n Z_i - \mu\right| > \epsilon\right) \leq 2e^{-\frac{2n\epsilon^2}{(b-a)^2}}$$

$$\therefore \frac{1}{n} \sum_{i=1}^n Z_i - \Delta \leq \mu \leq \frac{1}{n} \sum_{i=1}^n Z_i + \Delta$$

with prob $1-\delta$.

$$2e^{-\frac{2n\Delta^2}{(b-a)^2}} \leq \delta$$

A VC dimension approach

$$R(h) \leq R_{\text{emp}}(h) + \sqrt{\frac{V}{N} (\log N + \log \frac{1}{\delta})}$$

with prob $1-\delta$.

$$R_{\text{emp}}(h) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}\{h(x^{(i)}) \neq y^{(i)}\}$$

Training set error

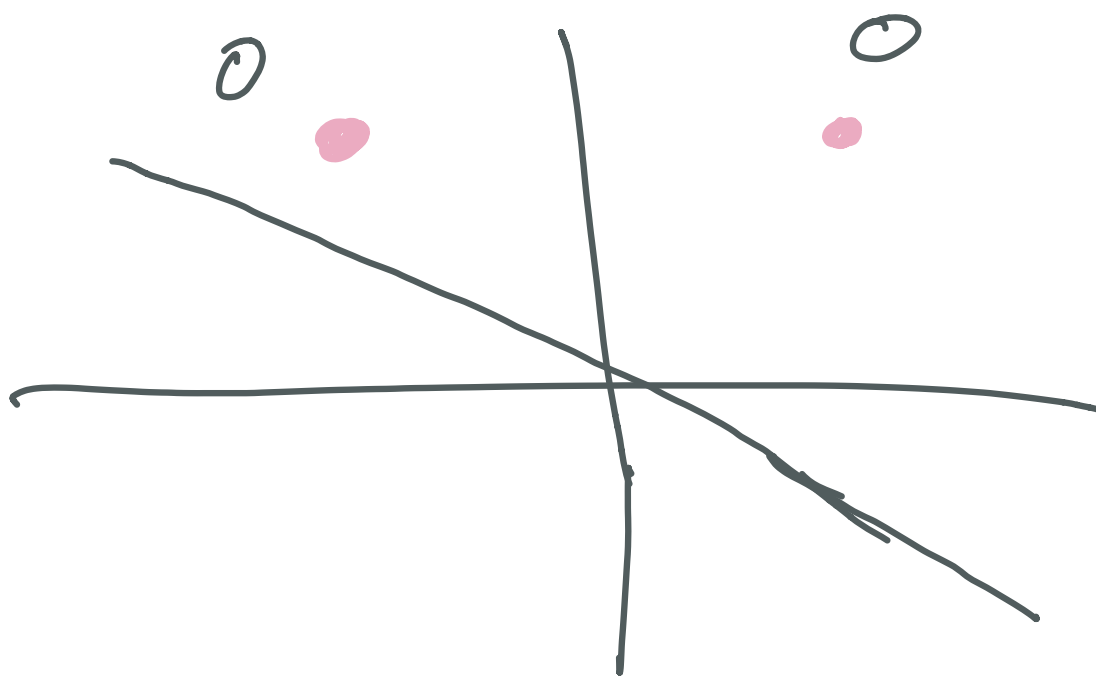
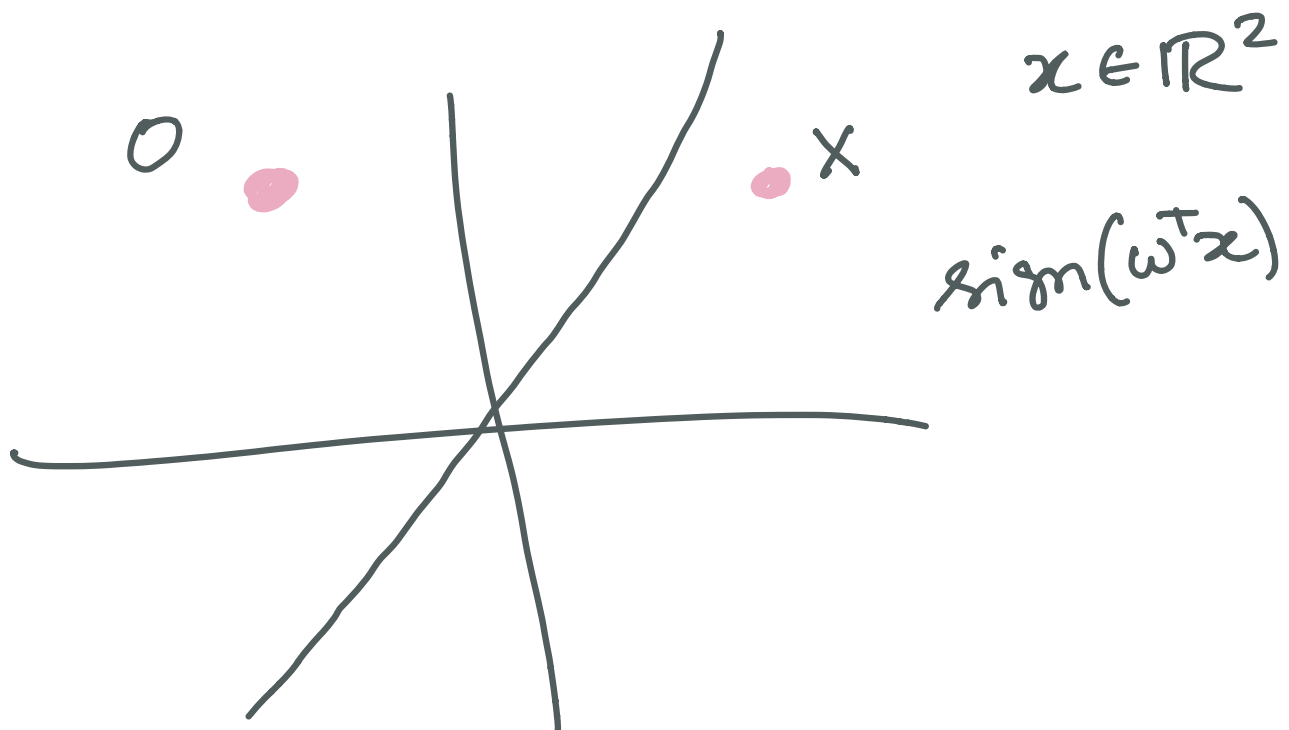
$V \rightarrow$ VC Dimension

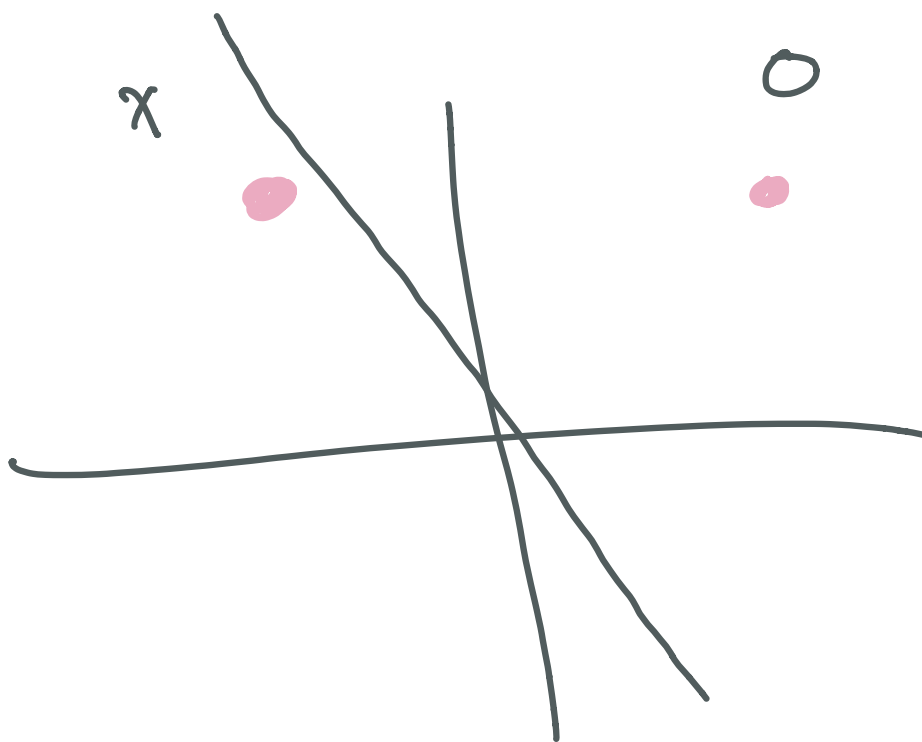
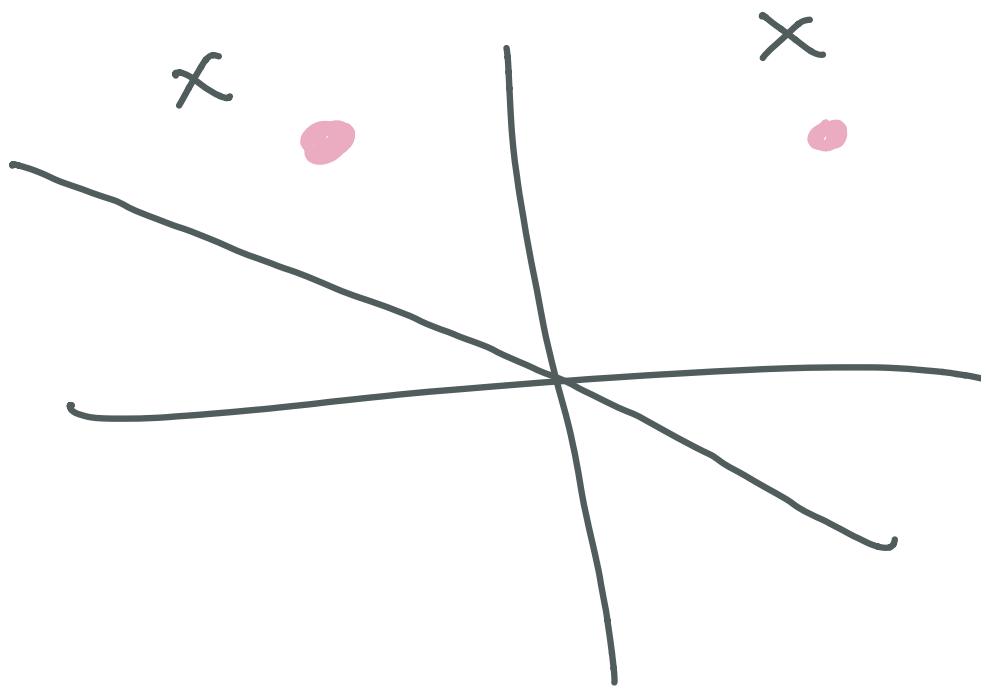
[Reading: Burges Tutorial]

$$H = \{h \mid h \text{ is a classifier}\}$$

$V(H)$ = maximum number of points on which classifiers from H can learn all possible Labellings.

$$H = \{h \mid h(x) = \text{sign}(\vec{w}^T x), w \in \mathbb{R}^2\}$$





Will also work for 3 points
But not for 4 points

$$\mathcal{H} \quad \|\omega\| \leq B. \quad \|x\| \leq R \quad x \in \mathbb{R}^d$$

$$H = \{h \mid h(x) = \text{sign}(\omega^T x + b)\}$$

$$V(H) \leq B^2 R^2$$

$$R(h) \leq R_{\text{emp}}(h) + O\left(\frac{RB}{\sqrt{n}}\right)$$

with prob $1 - \delta$.

$$\min_{\|\omega\| \leq B} R_{\text{emp}}(h)$$

$$h(x) = \text{sign}(\omega^T x + b)$$

Motivates the following
formulation

$$\min_{w, b} \quad C \cdot \sum_{i=1}^N \max(0, 1 - y_i (w^T x_i + b)) + \frac{1}{2} \|w\|^2$$