

Thursday, 10/1/2020

EE 660

MACHINE LEARNING  
FROM SIGNALS:  
FOUNDATIONS AND METHODS

Prof. B. Keith Jenkins

**Lecture 12**

---

**Lecture 12****EE 660****Oct 1, 2020**

---

**Announcements**

- Homework 4 (Week 5) is due tomorrow.
- Homework 5 (Week 6) will be posted.

---

**Today's Lecture**

- Overfitting (part 2)
- Regularization (AML view)

## Continue examples / experiments of over fitting (from last lecture)

Now consider:

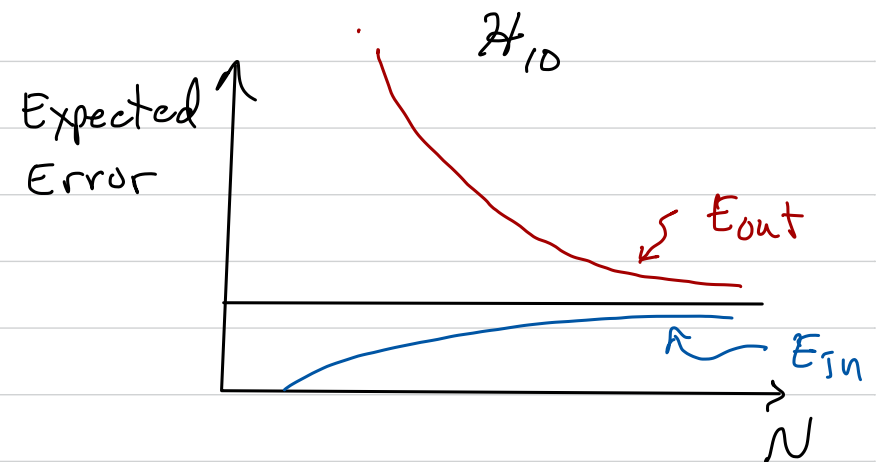
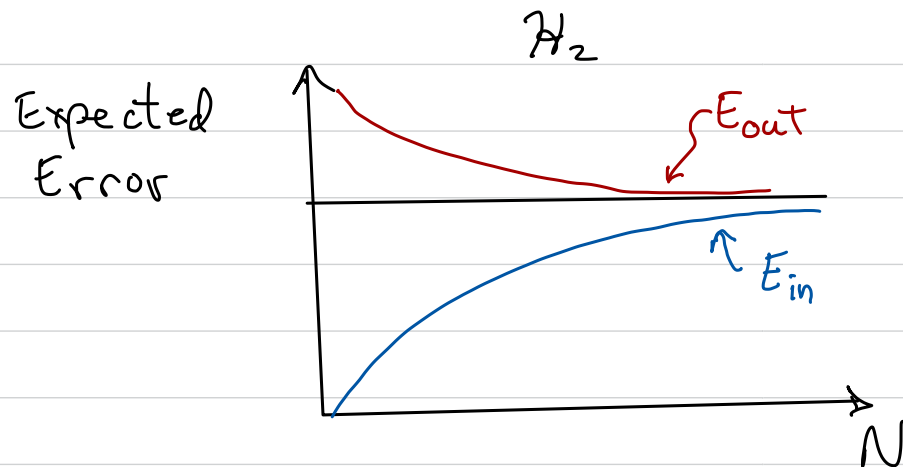
Target fcn.  $f(x)$  is  $d=50$  order polynomial, no noise.

$\mathcal{H}_2$  vs.  $\mathcal{H}_{10}$ :  $\mathcal{H}_{10}$  overfits! Even though  $\mathcal{H}_{10}$  is much simpler than  $f(x)$ .  
 $\rightarrow$  More data would help  $\mathcal{H}_{10}$ .

[AML fig. and table pp. 120-121]

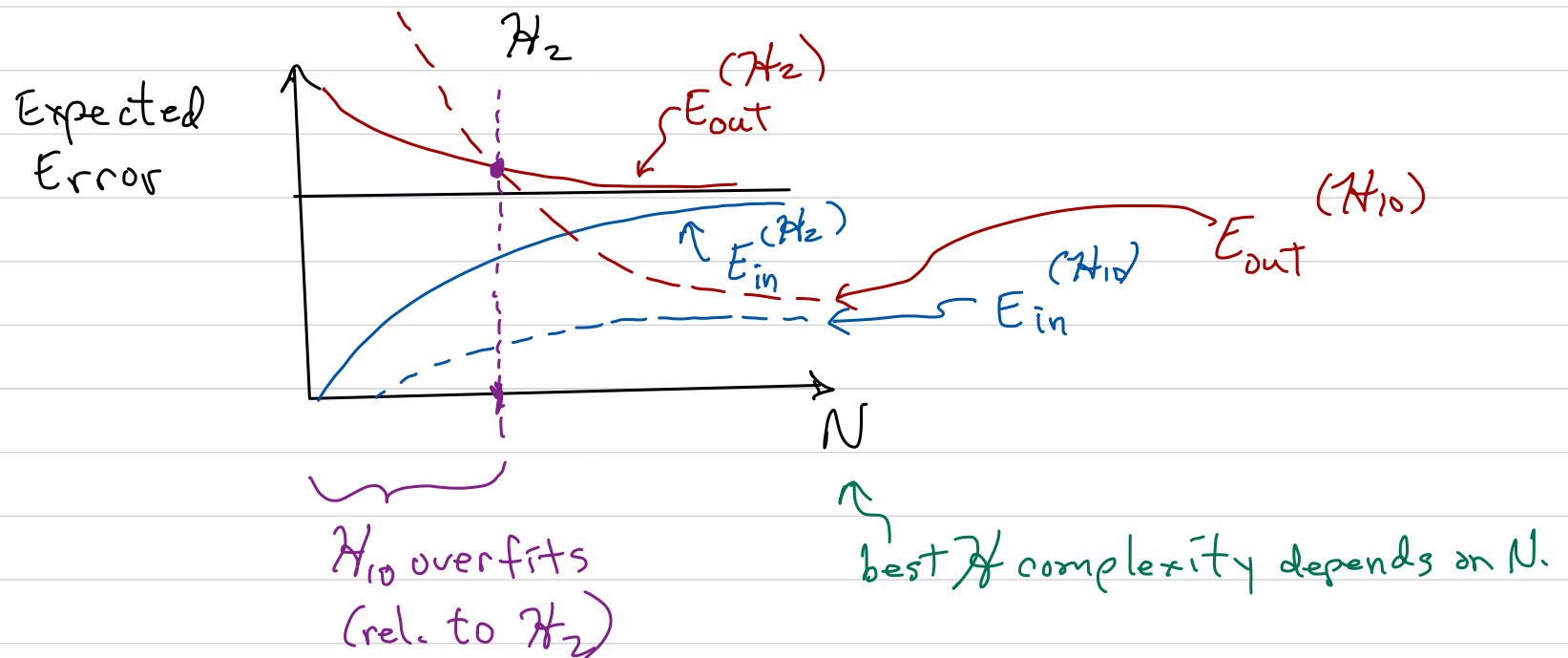
$\Rightarrow$  Best hypothesis set complexity depends on quantity of data.

### Learning Curves



Expected Error:  $\mathbb{E}_{\mathcal{D}} \{ E(h_g^{(\mathcal{D})}) \}$

If  $E_{in}(h_k) \leq E_{in}(h_j)$ , and  $E_{out}(h_k) > E_{out}(h_j)$ ,  
then  $h_k$  overfits the data (relative to  $h_j$ ).



What parameters affect amount of overfitting?

$N$ ,  $\sigma$ , complexity of  $H$ , complexity of  $f(\underline{x})$ .

$\uparrow$   
noise on data

[A good example / experiment will be covered in Discussion 6.]

# Regularization and Complexity [AML 4.2]

VC bound view:

$$E_{\text{out}}(h) \leq E_{\mathcal{D}}(h) + \Omega(\mathcal{H}, N, \delta) \quad \forall h \in \mathcal{H}$$

Learning alg. finds  $h_g = \operatorname{argmin}_{h \in \mathcal{H}} E_{\mathcal{D}_{\text{Tr}}}(h)$

Note:  $\Omega$  depends on  $\mathcal{H}$  but not on  $h_g$ .

With regularization: ( $\underline{w} = \underline{w}^{(0)}$ )

$$\text{Let } f_{\text{obj}}^{(r)}(h_{\underline{w}}) = E_{\mathcal{D}}(h_{\underline{w}}) \quad \underbrace{\text{subject to } \underline{w}^T \underline{w} \leq C}_{\text{"soft order constraint"}}$$

$$(i) \quad h_g = \operatorname{argmin}_{h \in \mathcal{H}} E_{\mathcal{D}}(h) \quad \text{subject to } \underline{w}^T \underline{w} \leq C \quad (C \text{ is a parameter})$$

equivalent to:

$$(ii) \quad h_g = \operatorname{argmin}_{h \in \mathcal{H}'} E_{\mathcal{D}}(h) \quad \text{in which } \mathcal{H}' = \{h \mid h \in \mathcal{H} \text{ and } \underline{w}^T \underline{w} \leq C\}$$

$\mathcal{H}'$  is different than  $\mathcal{H}$  (in general)

Note:  $\mathcal{Q}_{vc}(\mathcal{H}') \leq \mathcal{Q}_{vc}(\mathcal{H})$ .

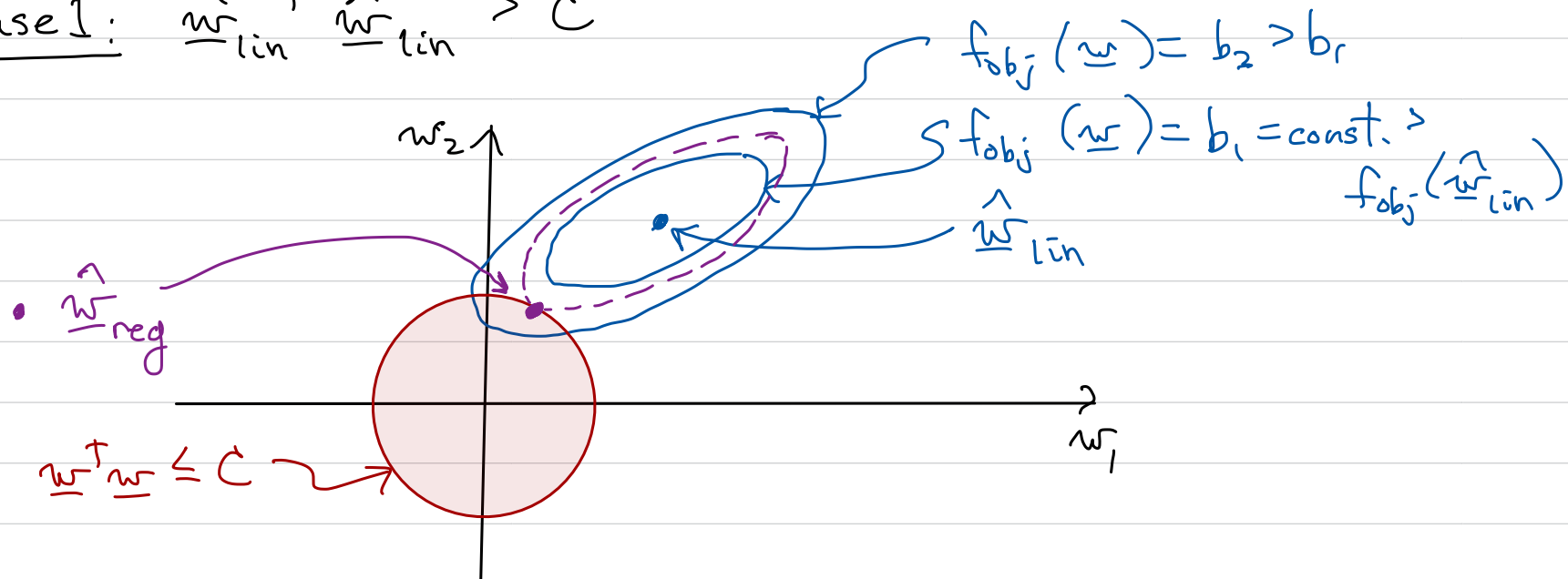
Graphical View:

$$f_{obj}(\underline{w}) \triangleq E_{\mathcal{D}}(h_{\underline{w}}) \stackrel{\text{For example}}{=} \frac{1}{N} \sum_{n=1}^N (\underline{w}^T \underline{x}_n + w_0 - y_n)^2$$

$$= \text{min. at } \hat{\underline{w}}_{lin} \quad \leftarrow \begin{array}{l} \text{(lin for linear)} \\ \text{(unconstrained optimum)} \end{array}$$

Using (i) above:

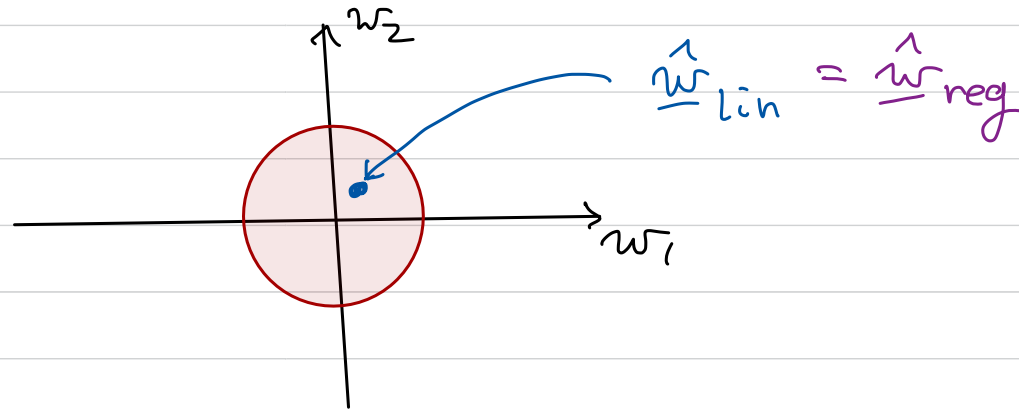
Case 1:  $\hat{\underline{w}}_{lin}^T \hat{\underline{w}}_{lin} > C$



$\hat{\underline{w}}_{\text{reg}}$  is the min. of  $f_{\text{obj}}$  subject to  $\underline{w}^T \underline{w} \leq C$ .

→ In Case 1,  $\hat{\underline{w}}_{\text{reg}}^T \hat{\underline{w}}_{\text{reg}} = C$ , and  $\hat{\underline{w}}_{\text{reg}} \neq \hat{\underline{w}}_{\text{lin}}$ .

Case 2:  $\hat{\underline{w}}_{\text{lin}}^T \hat{\underline{w}}_{\text{lin}} \leq C$



In Case 2,  $\hat{\underline{w}}_{\text{reg}}^T \hat{\underline{w}}_{\text{reg}} \leq C$ , and  $\hat{\underline{w}}_{\text{lin}} = \hat{\underline{w}}_{\text{reg}}$ .