

# ERROR DECOMPOSITION

①

Classification Task

$(x, y)$ ,  $x \in \mathbb{R}^d$ ,  $y \in \{1, 2, \dots, K\}$   
K-way classification.

feature vector in d-dim      class label ground truth

$p(x, y)$  Unknown joint distribution  
d i/p  $x$  & output  $y$ .

$D$  sampled from  $p(x, y)$  ||  $D = \{D_{\text{train}}, D_{\text{test}}\} \sim p(x, y)$

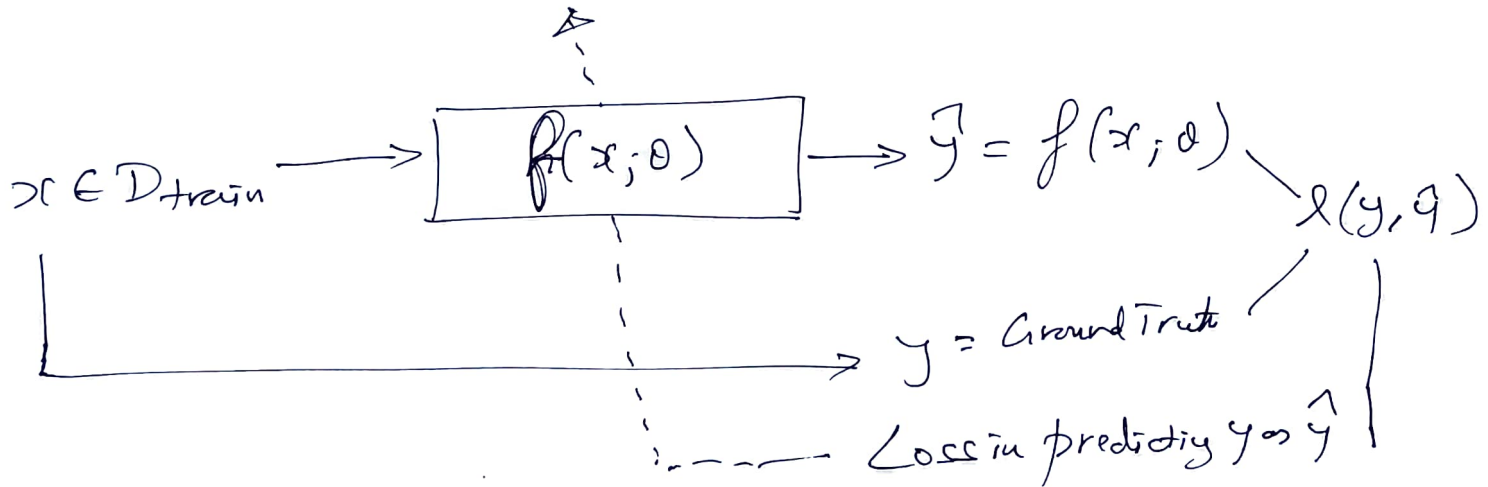
$N$  sample training data  $\{(x_i, y_i)\}_{i=1}^N$

$\{x^{\text{test}}\}$

②  $\exists$  optimal hypothesis/model  $f^+$ :  $f^+(x) \Rightarrow y$

But within a hypothesis/model space  $F$ , ML/DL learning  
seeks to discover  $f^+$  by fitting  $D_{\text{train}}$   $\approx$  testing on  $D_{\text{test}}$ .

$\Rightarrow$   
 $f(\cdot; \theta)$ :  $f(x; \theta) = \hat{y}$ : Predicted  
Label of  $x$ .



Expected Loss/Risk

Given  $p(x, y)$

$$R(f) = \int l(f(x), y) dp(x, y)$$

$$= \mathbb{E} [l(f(x), y)]$$

Expected Loss/Risk Minimizer

$$f^\dagger = \underset{f}{\operatorname{argmin}} R(f)$$

Given  $F$

$$f^* = \underset{f \in F}{\operatorname{argmin}} R(f)$$

Empirical Loss/Risk. (3)

Given  $D_{\text{train}} = \{(x_i, y_i)\}_{i=1}^N$

$$R_N(f) = \frac{1}{N} \sum_{i=1}^N l(f(x_i), y_i)$$

Empirical Loss/Risk Minimizer

$$f_N = \underset{f \in F}{\operatorname{argmin}} R_N(f)$$

To know how well  $f_N$  performs on test:

Measure  $R(f_N)$

(4)

$$f(x, y)$$

 $D_{\text{train}}$ 
 $D_{\text{test}}$ 

$$\{x_i, y_i\}_{i=1}^N$$

Expected Loss/Risk

$$R(f) = \mathbb{E}_{p(x, y)} [\ell(f(x), y)]$$

optimal

$$f^* = \underset{f}{\operatorname{argmin}} R(f)$$

- Minimize Expected Risk
- NO CONSTRAINTS on  $f$

within  $F$

$$f_F^* = \underset{f \in F}{\operatorname{argmin}} R(f)$$

- Minimize Expected Risk.
- $f \in F$
- $p(x, y)$ : UNKNOWN.

Empirical Loss/Risk

$$R_N(f) = \frac{1}{N} \sum_{i=1}^N \ell(f(x_i), y_i)$$

Within  $F$ , using  $D_{\text{train}}$

$$f_N = \underset{f \in F}{\operatorname{argmin}} R_N(f)$$

- Minimize Empirical risk on  $D_{\text{train}} \sim p(x, y)$  of  $N$  samples
- $f \in F$

$$f^{\dagger} = \underset{f}{\operatorname{argmin}} R(f)$$

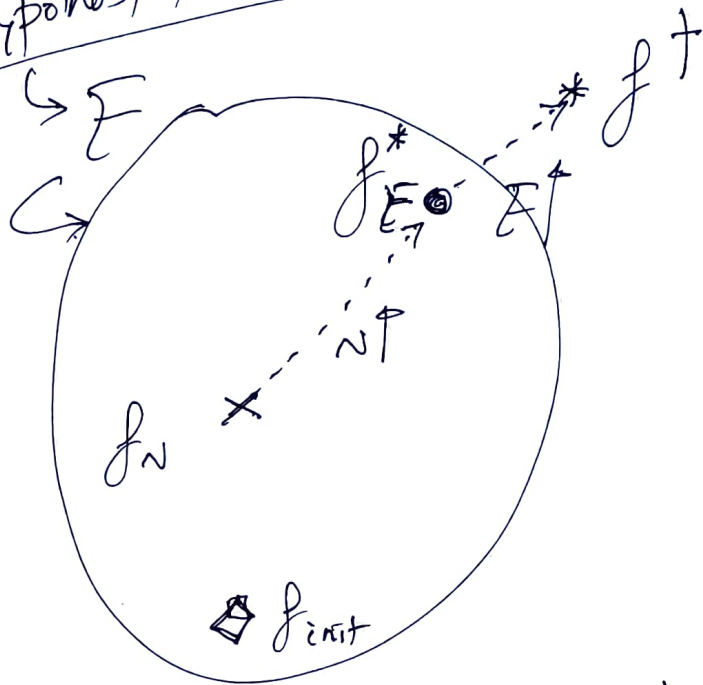
$$f_E^* = \underset{f \in E}{\operatorname{argmin}} R(f)$$

$$f_N = \underset{f \in E}{\operatorname{argmin}} R_N(f)$$

where  $R(f) = \mathbb{E}_{p(x,y)} [\ell(f(x), y)]$

$$\text{or } R_N(f) = \frac{1}{N} \sum_{i=1}^N \ell(f(x_i), y_i)$$

Hypothesis / Model space



As  $\mathbb{P} : f_E^* \rightarrow f^{\dagger}$

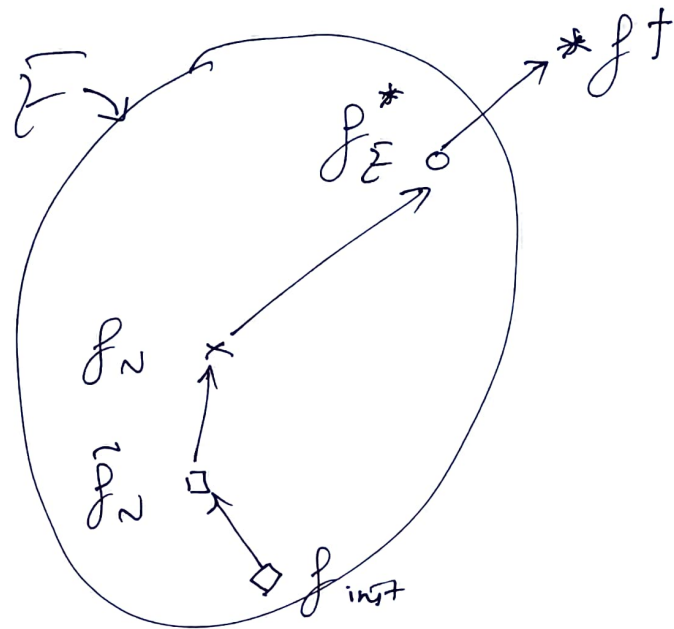
As  $N \uparrow : f_N \rightarrow f_E^*$

# Including "Optimization" Process & Corresponding 'Error'

$f_N$  is determined from some  $f_{\text{init}}$  by an optimization process to yield a sub-optimal  $\tilde{f}_N$

$$\text{i.e. } f_N = \text{opt} \left[ \underset{f \in \mathcal{F}}{\text{argmin}} R_N(f) \right]$$

$$f^\dagger \xrightarrow{f \in \mathcal{F}} f_\mathcal{E}^* \xrightarrow{D_{\text{train}}} f_N \xrightarrow{\text{opt}} \tilde{f}_N$$



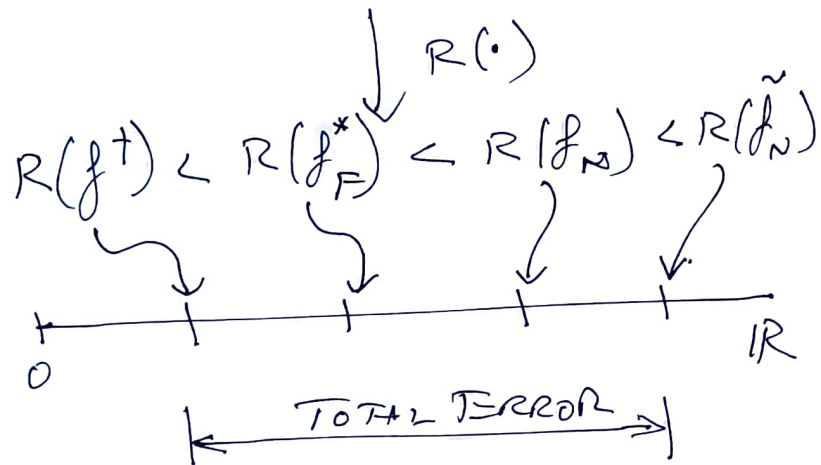
Generalization on Test:  $p(x, y)$

$$R(f^\dagger) \leq R(f_\mathcal{E}^*) \leq R(f_N) \leq R(\tilde{f}_N)$$

Note: Models evaluated on  $R(\cdot)$

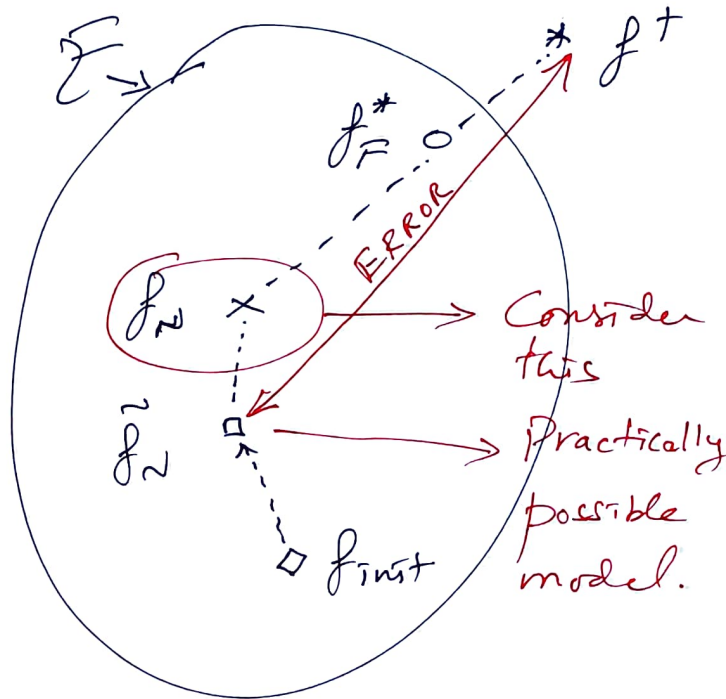


$$f^+ \longrightarrow f_F^* \longrightarrow f_N \longrightarrow \tilde{f}_N$$



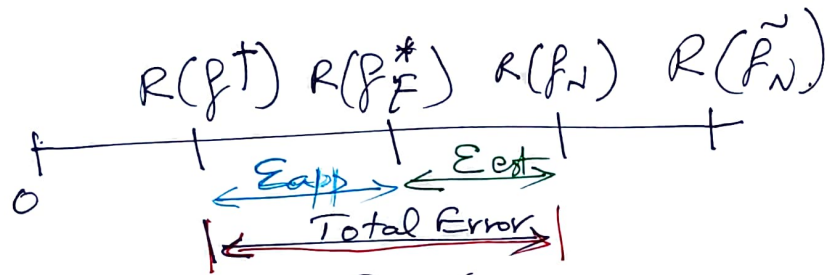
$$= \mathbb{E} [R(\tilde{f}_N) - R(f^+)]$$

on  
sampling  $d$   
D<sub>train</sub>.





Consider only  $f_N$  (not  $\tilde{f}_N$ )

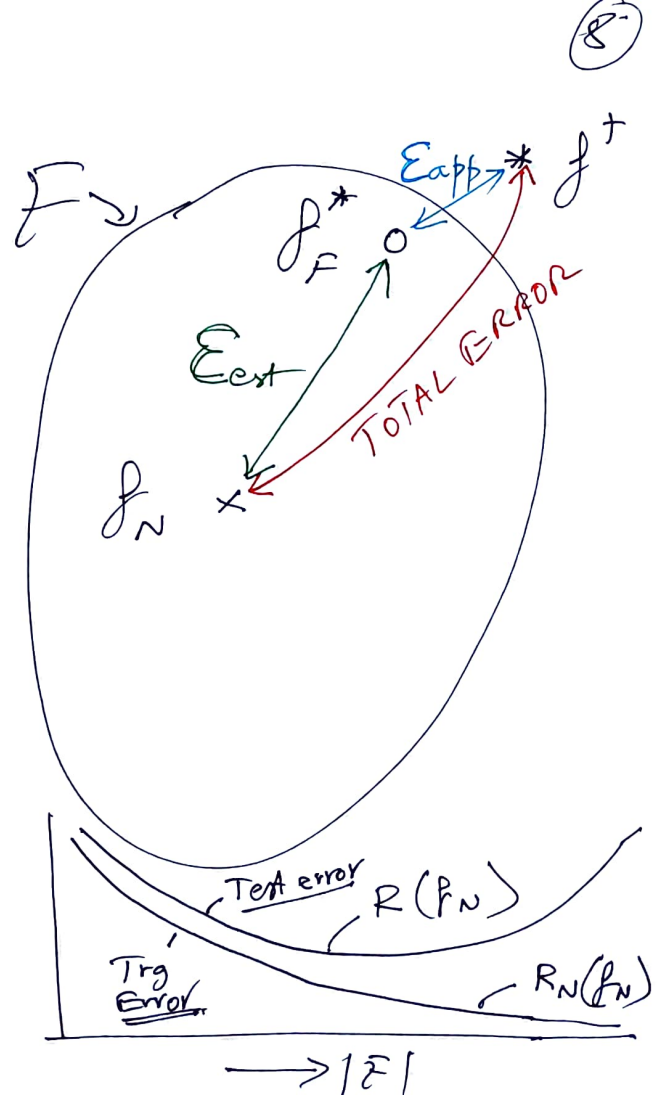


Generalization Error

$$= \mathbb{E}[R(f_N) - R(f^+)]$$

$$= \mathbb{E}[R(f_N) - R(f_F^*)]$$

$$+ \mathbb{E}[R(f_F^*) - R(f^+)]$$



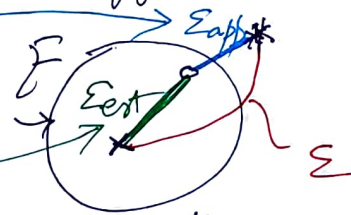


(8.7)

$$\text{Error } \underline{\Sigma} = \mathbb{E} [R(f_N) - R(f_{\mathcal{F}}^*)] \leadsto \Sigma_{\text{estimation}}$$

$$+ \mathbb{E} [R(f_{\mathcal{F}}^*) - R(f^*)] \leadsto \Sigma_{\text{approximation}}$$

$$= \Sigma_{\text{estimation}} + \Sigma_{\text{approximation}}$$



As  $|\mathcal{F}| \rightarrow \infty$  or  $|\mathcal{F}| \uparrow$ ,  $\Sigma_{\text{app}} \rightarrow 0$  i.e.  $f_{\mathcal{F}}^* \rightarrow f^*$

For a given  $\mathcal{F}$ ,

As  $N \rightarrow \infty$  or  $N \uparrow$ ,  $\Sigma_{\text{est}} \rightarrow 0$  i.e.  $f_N \rightarrow f_{\mathcal{F}}^*$

(9)

Then  
Error  
(GE)  $= \mathbb{E} \left[ R(f_N) - R(f_F^*) \right] + \epsilon_F$

where  $\epsilon_F = \Sigma_{\text{app}}$  [a function of  $F$  only]

or  $R(f_N) - R(f_F^*) = GE - \epsilon_F$

For finite  $F$ ,  $\Sigma_{\text{app}}$  (or  $\epsilon_F$ )  $> 0$

Then  $R(f_N) - R(f_F^*) < GE$

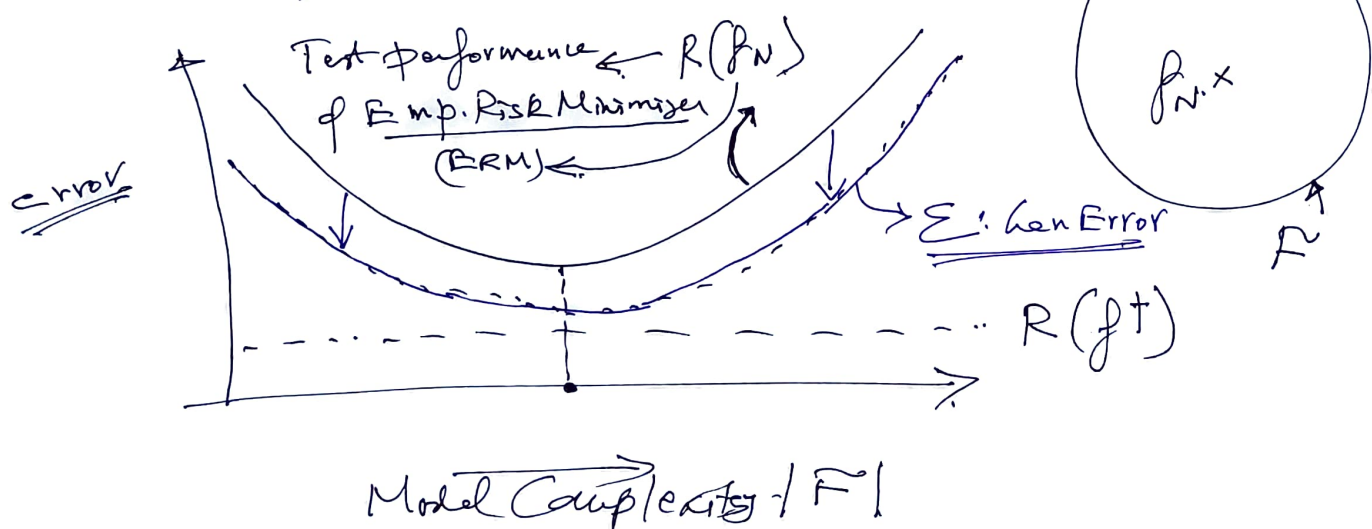
or  $\boxed{R(f_N) < R(f_F^*) + GE}$

Error

$$\varepsilon = \varepsilon_{\text{est}} + \varepsilon_{\text{app}}$$

$$\mathbb{E}(R(f_N) - R(f_F^*)) \quad \mathbb{E}[R(f_F^*) - R(f^\dagger)]$$

$$\varepsilon = \mathbb{E}[R(f_N) - R(f^\dagger)]$$



Why does  $R(f_N)$  behave so?

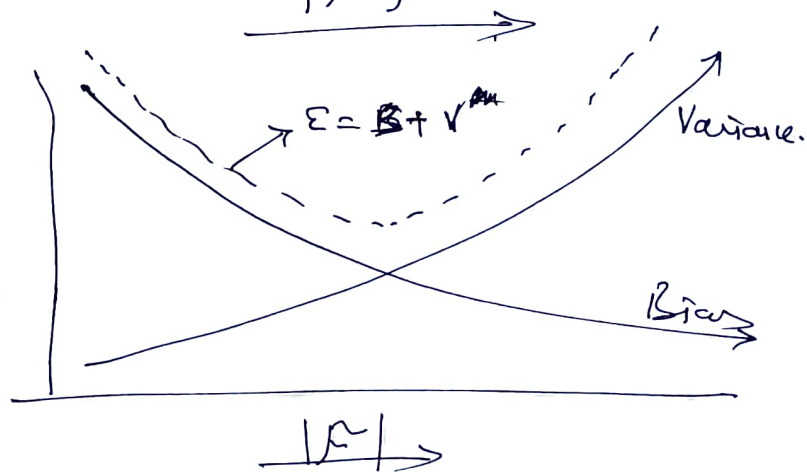
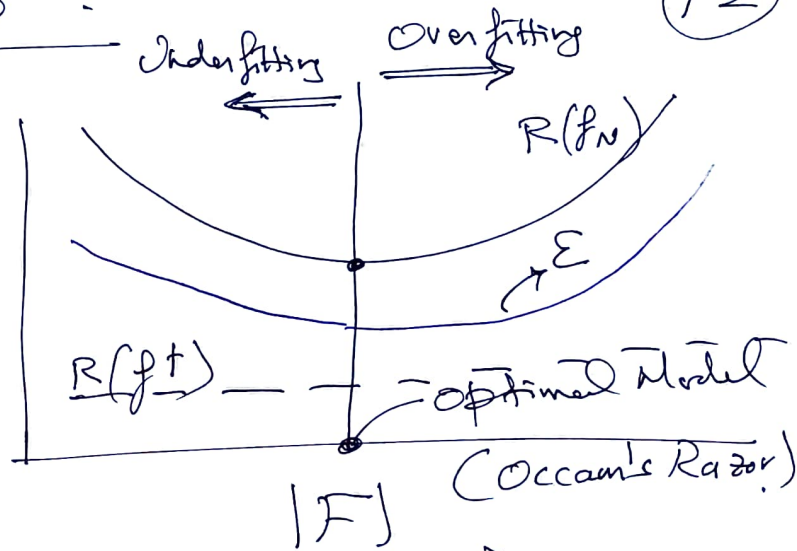
(9.2)

• Model Complexity  $\propto \downarrow$  (low) error

- No learning of complex patterns in data
- High Bias, low Variance.

• Model Complexity  $\propto \uparrow$  (High)

- Learns complex patterns
- Memorizes Training Data.
- Performs poorly on Unseen data
- Low Bias, high Variance.



•  $\mathcal{E}_{est} = \mathbb{E} [R(p_N) - R(p_F^*)]$

•  $\mathcal{E}_{app} = \mathbb{E} [R(p_F^*) - R(p^\dagger)]$

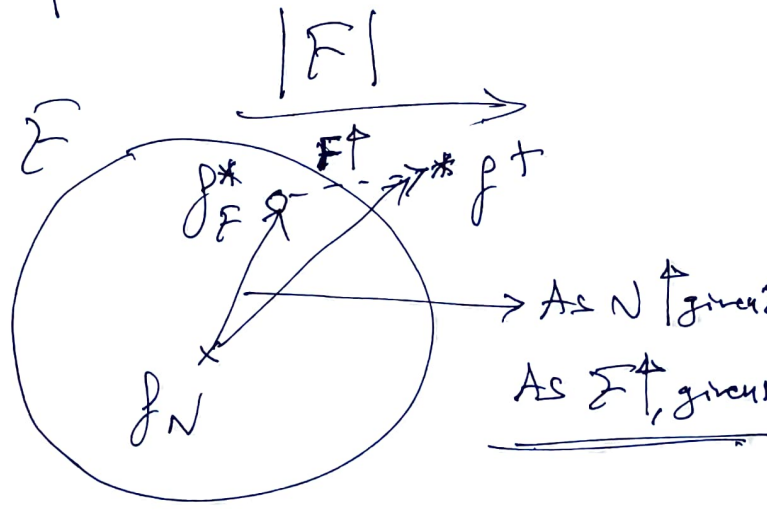
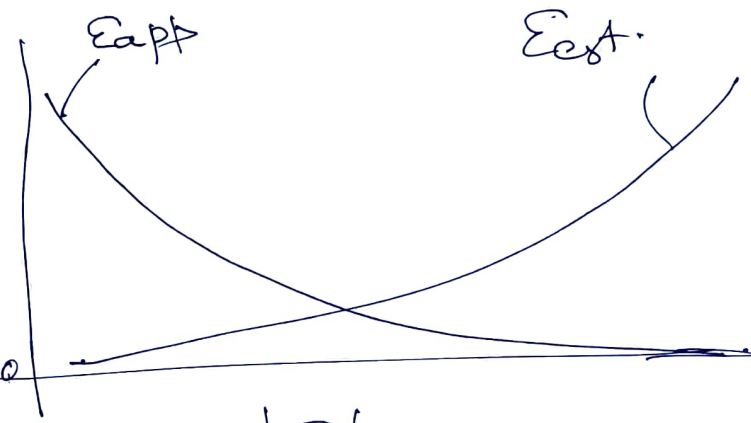
$\xrightarrow{\mathcal{E}_{app}}$  as  $|\mathcal{F}| \uparrow \rightarrow [p_F^* \rightarrow p^\dagger]$

$\geq \mathcal{E}_{app} \rightarrow 0$

$\xrightarrow{\mathcal{E}_{est}}$   $\mathbb{E} [R(p_N) - R(p_F^*)]$

as  $|\mathcal{F}| \uparrow \rightarrow R(p_F^*) \rightarrow R(p^\dagger)$

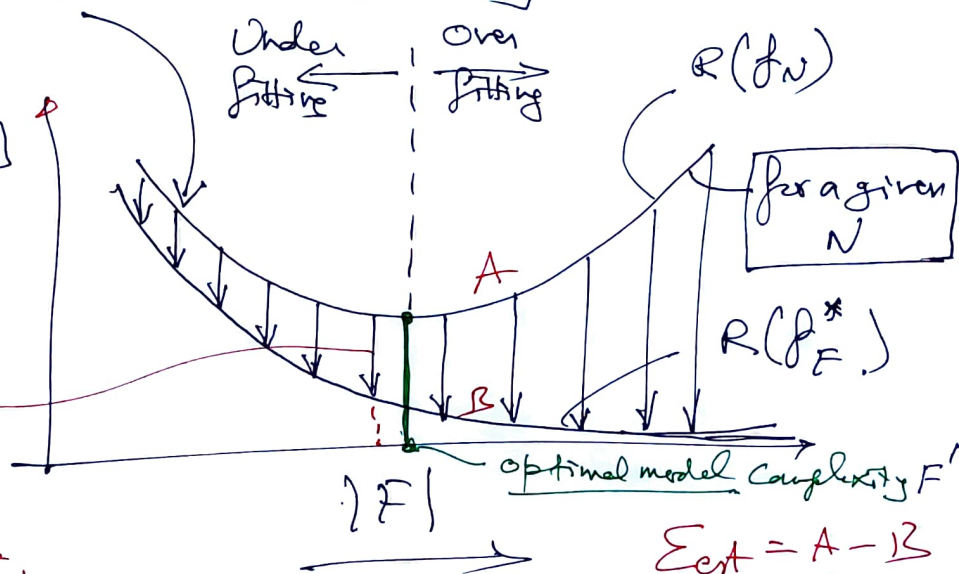
~~faster than  $R(p_N) \rightarrow R(p^\dagger)$~~



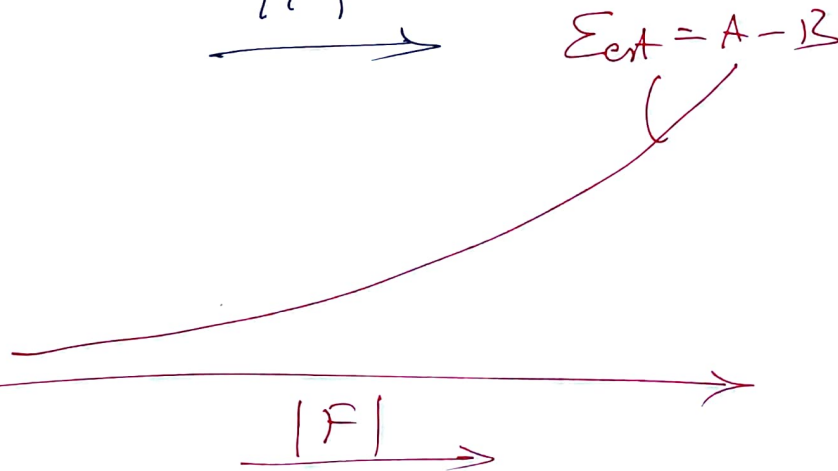
Analyze  $E_{\text{est}} = \mathbb{E} [R(f_N) - R(f_F^*)]$  (9.4)

- $A$  decreases faster than  $B$  till  $F = F'$  [optimal model]
- $A \uparrow$  for  $F > F'$

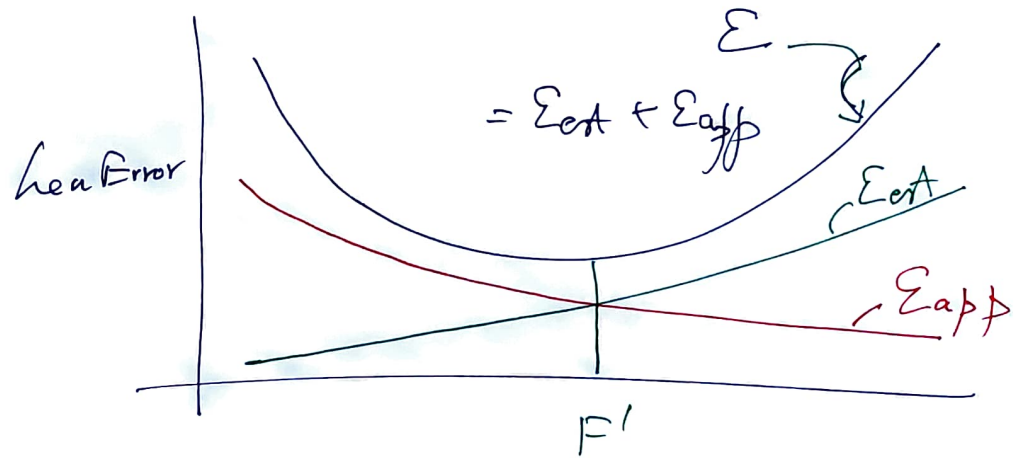
$E_{\text{est}}(F) \leftarrow$   
 $= A - B$



- Given  $N$ 
  - $-F \uparrow \rightarrow [f_F^* \rightarrow f^*]$   $R(.)$  error  
unaffected by  $N$ .
  - $-F \uparrow \rightarrow f_N$  affected / fixed  
 $\rightarrow R(f_N) \begin{cases} \rightarrow \text{UF} [F < F'] \\ \rightarrow \text{OV} [F > F'] \end{cases}$



$$\frac{\text{Overall Gen Error}}{\text{Gen Error}} = \Sigma_{\text{est}} + \Sigma_{\text{app}}$$





# In Summary

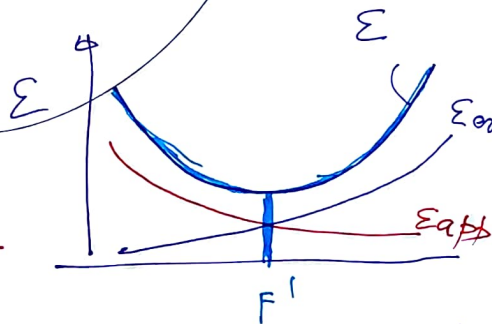
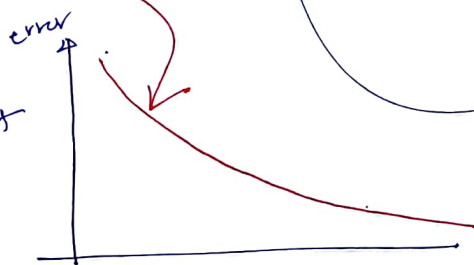
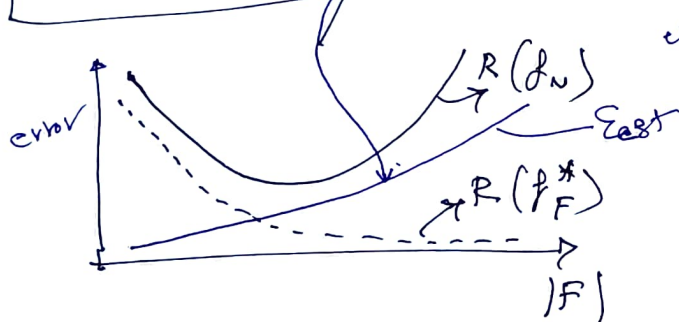
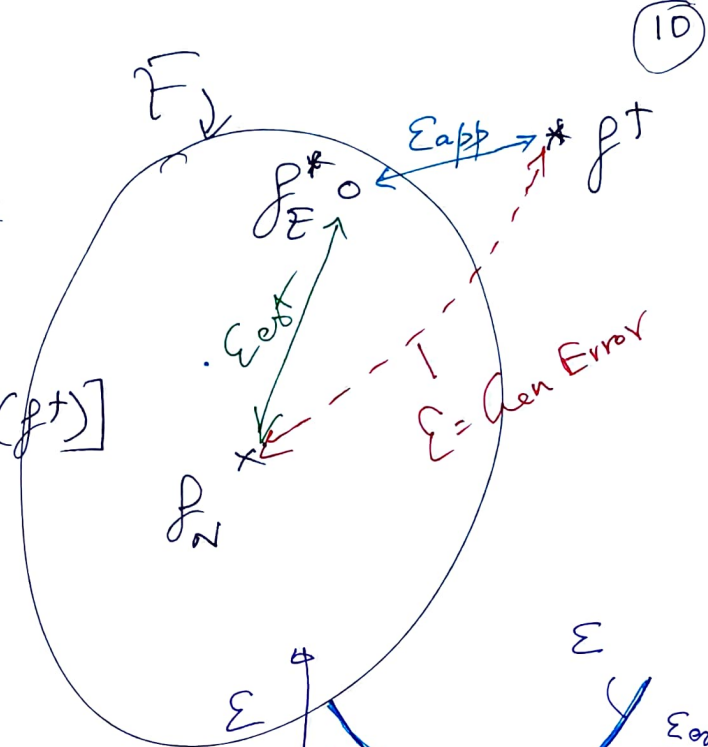
$\mathcal{E}$  = Generalization Error

$$= \mathbb{E}[R(p_N) - R(p^*)]$$

$$= \mathbb{E}[R(p_N) - R(p_F^*)] + \mathbb{E}[R(p_F^*) - R(p^*)]$$

$$\mathcal{E} = \mathcal{E}_{\text{est}}(F, N) + \mathcal{E}_{\text{app}}(F)$$

For a given  
 $|D_{\text{train}}| = N$



For finite  $F$ ,  $\mathcal{E}_{\text{app}}(F) > 0$

$$\therefore \text{G.E.} > R(p_N) - R(p_F^*)$$

$$R(p_N) < R(p_F^*) + \mathcal{E}$$

Generalization Error