

Evaluating Predictors/Models

Objective (formal):

Define a Learner $\mathcal{A}: (\mathcal{X} \times \mathcal{Y})^n \rightarrow \{f \mid f: \mathcal{X} \rightarrow \mathcal{Y}\}$
such that $\mathbb{E}[L(\mathcal{A}(D))]$ is small

Defining $\mathcal{A}(D)$: Empirical Risk Minimization (ERM)

Estimation:

Use D to estimate $L(f)$ for all $f \in \mathcal{F} \subset \{f \mid f: \mathcal{X} \rightarrow \mathcal{Y}\}$
call the estimate $\hat{L}(f)$

Optimization:

pick \hat{f} to be the $f \in \mathcal{F}$ that minimizes $\hat{L}(f)$
 \downarrow
Function class

When should we expect ERM to work well?

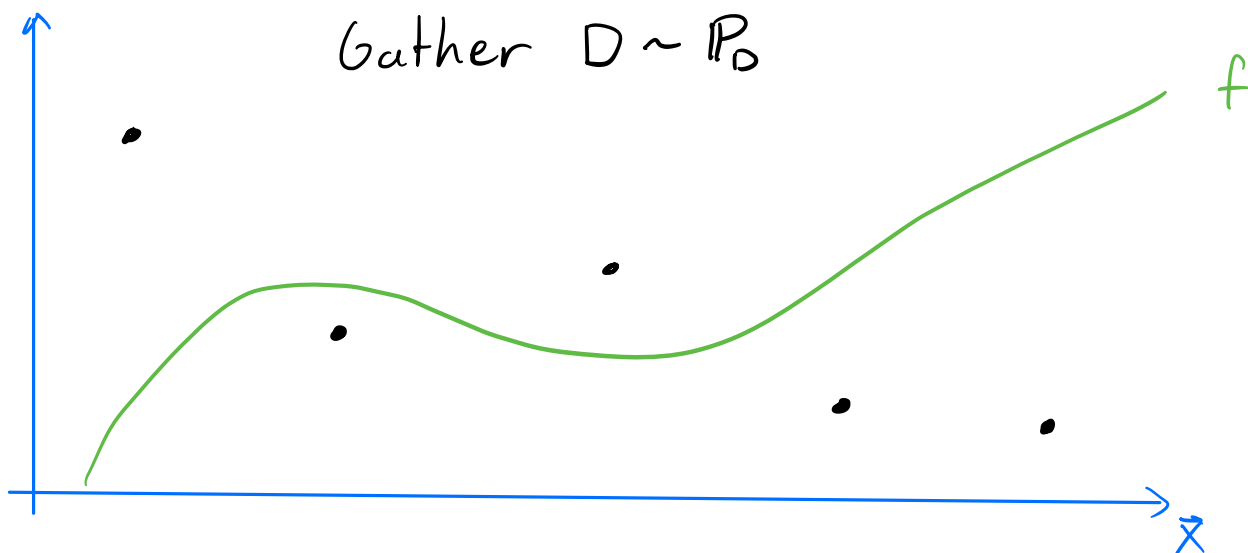
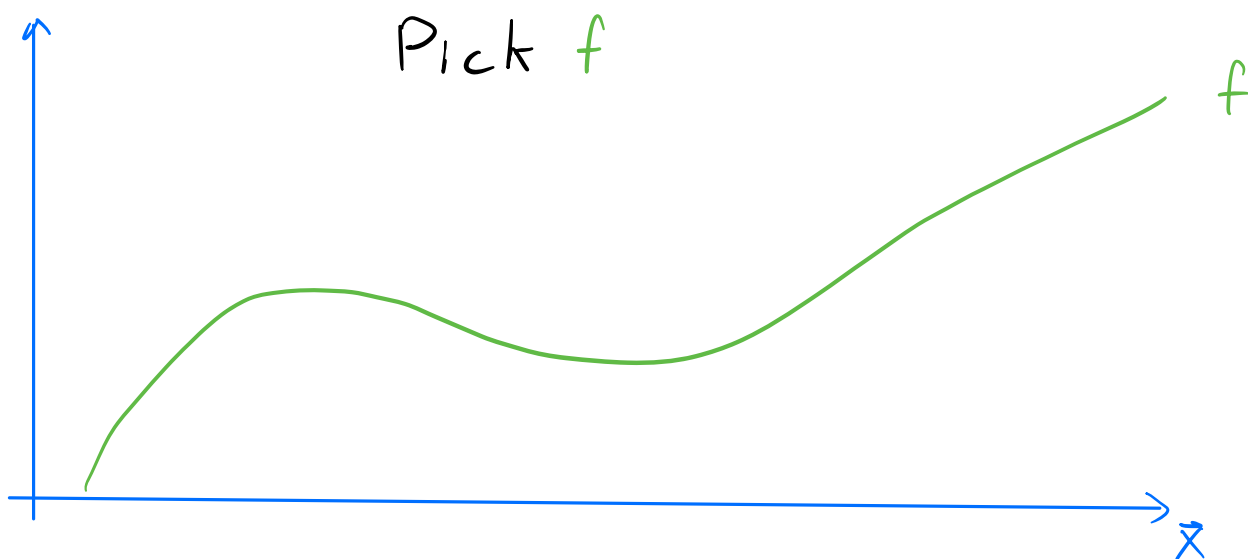
- When \mathcal{F} contains an f that can make $L(f)$ small
- When $\hat{L}(\hat{f})$ is a good estimate of $L(\hat{f})$

Is $\hat{L}(\hat{f}_D)$ really a good estimate of $L(\hat{f}_D)$?

where $A(D) = \hat{f}_D \in \mathcal{F}$ D is a r.v. $(\vec{x}, Y) \sim P_{\vec{x}, Y}$

$$\hat{L}(\hat{f}_D) = \frac{1}{n} \sum_{i=1}^n \ell(\hat{f}_D(\vec{x}_i), Y_i) \quad L(\hat{f}_D) = E[\ell(\hat{f}_D(\vec{x}), Y)]$$

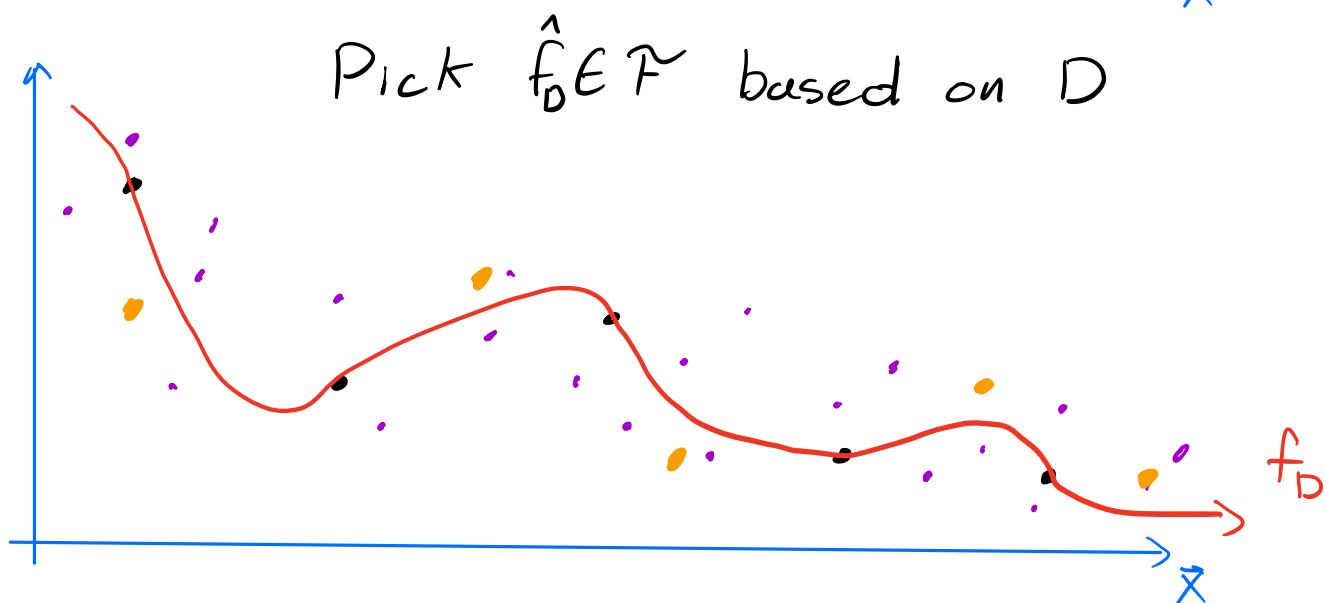
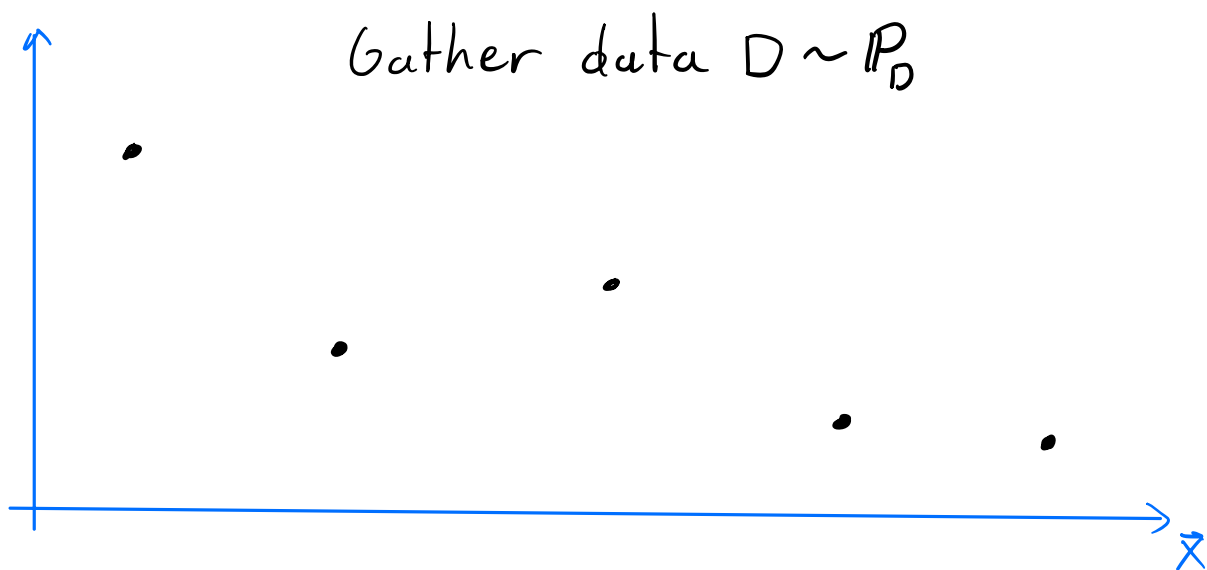
If we pick $f \in \mathcal{F}$ and then gather D
(i.e. f is chosen independently of D)



Then: $E[\hat{L}(f)] = L(f)$

$$\text{Var}[\hat{L}(f)] = \frac{1}{n} \text{Var}[\ell(f(\vec{x}_i), Y_i)]$$

We are gathering data D and then pick $\hat{f}_D \in \mathcal{F}$! (i.e. \hat{f}_D depends on D)



Then, in general:

$$\mathbb{E}[\hat{L}(\hat{f}_0)] \neq \mathbb{E}[L(\hat{f}_0)]$$

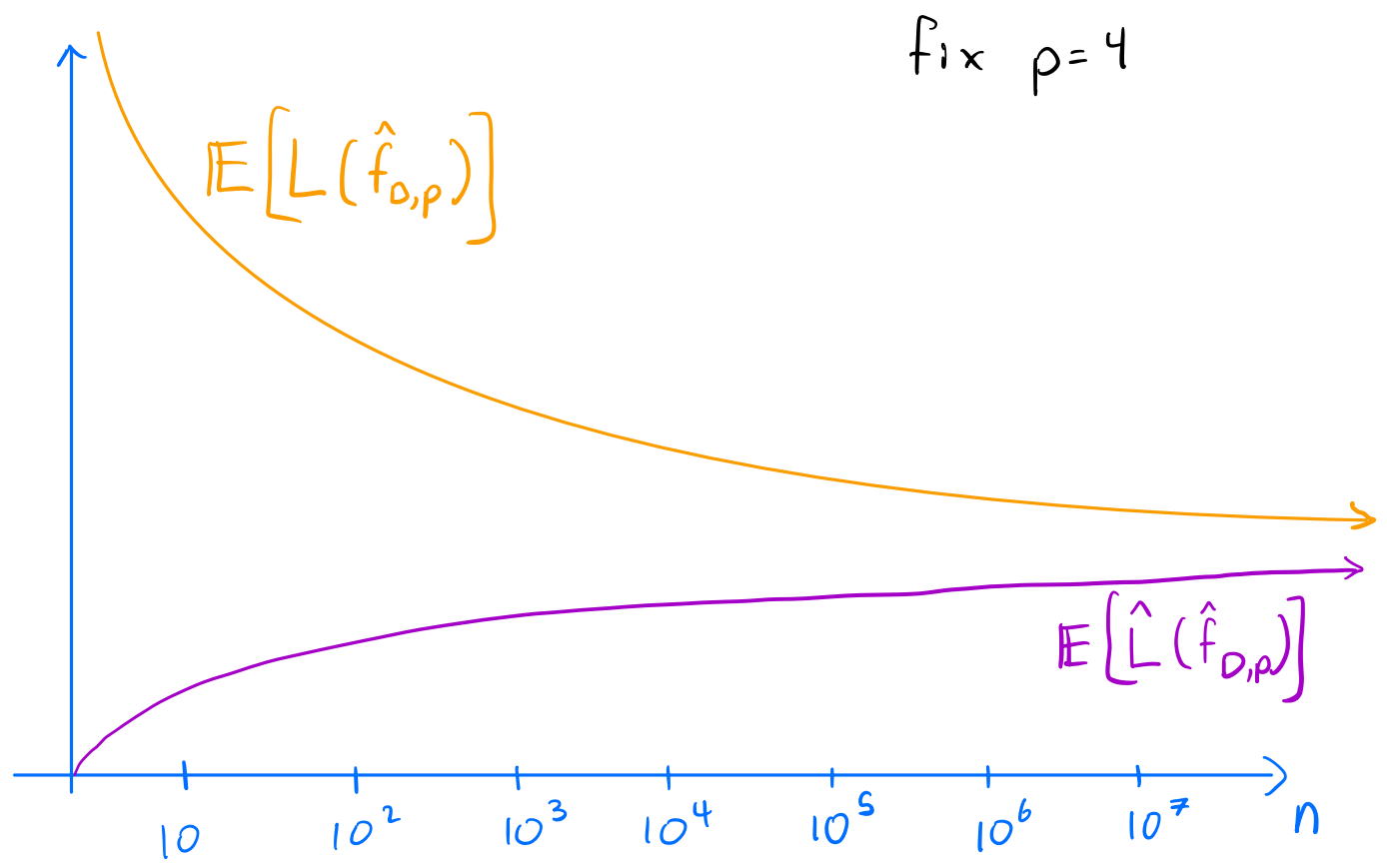
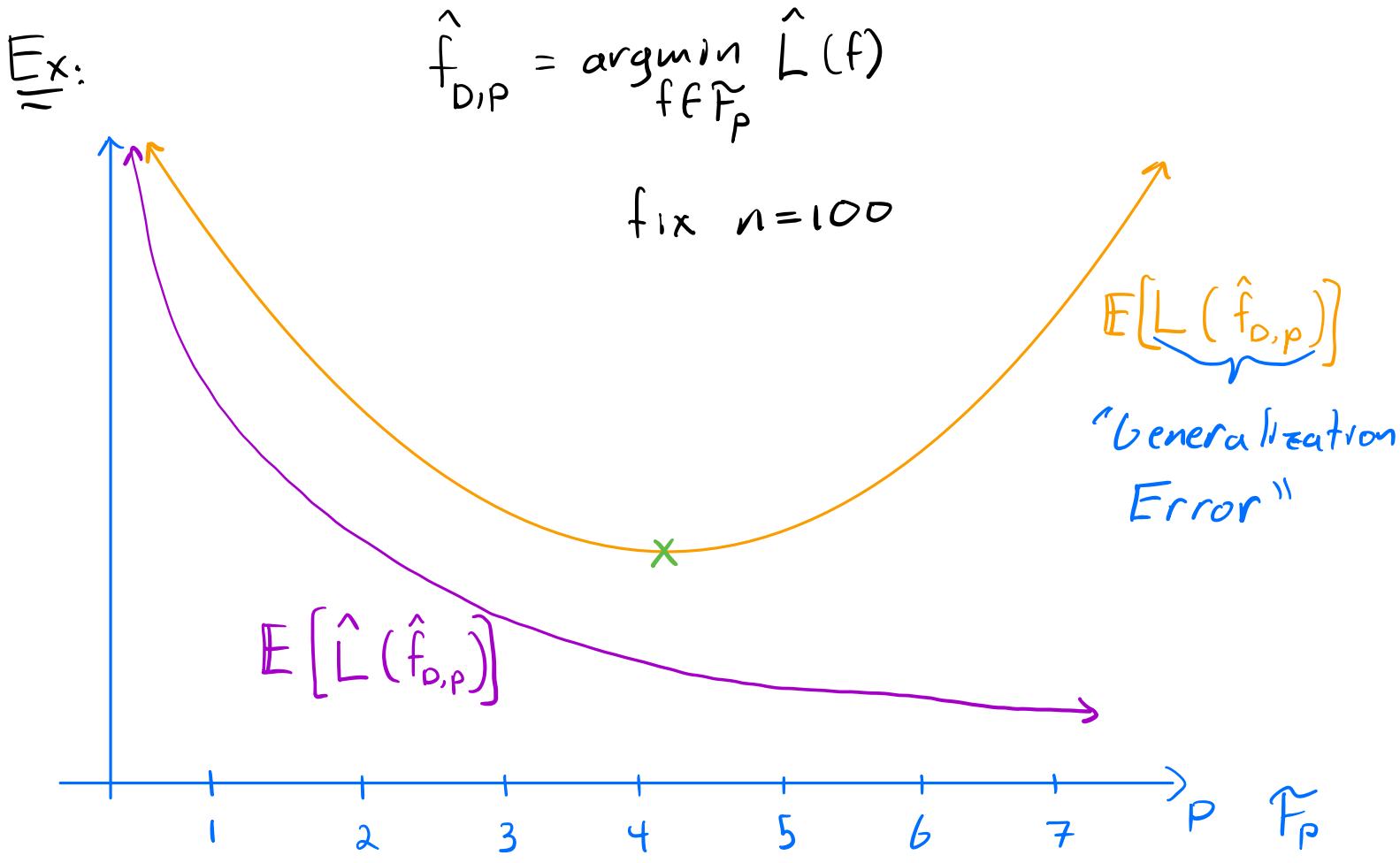
$l(\hat{f}_0(\vec{X}_i), Y_i)$ are not i.i.d.

\hat{f}_0 depends on $(\vec{X}_1, Y_1), \dots, (\vec{X}_n, Y_n)$!

It can be shown that:

$$\mathbb{E}[L(\hat{f}_0)] - \mathbb{E}[\hat{L}(\hat{f}_0)]$$

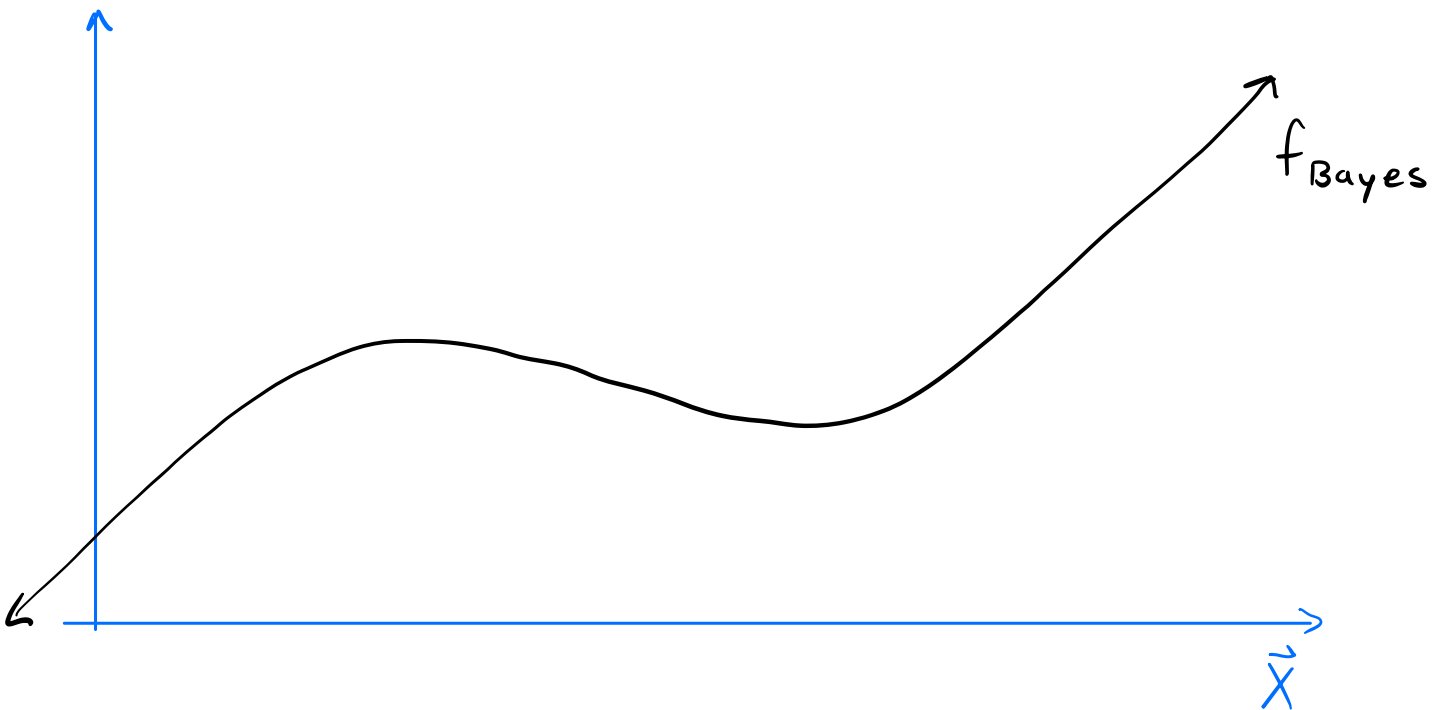
- increases as \mathcal{F} gets more complex
- decreases as n increases



Decomposing $E[L(A(D))]$ into Error Terms

Suppose we knew $P_{\mathcal{X}, \mathcal{Y}}$ what would we choose for A ?

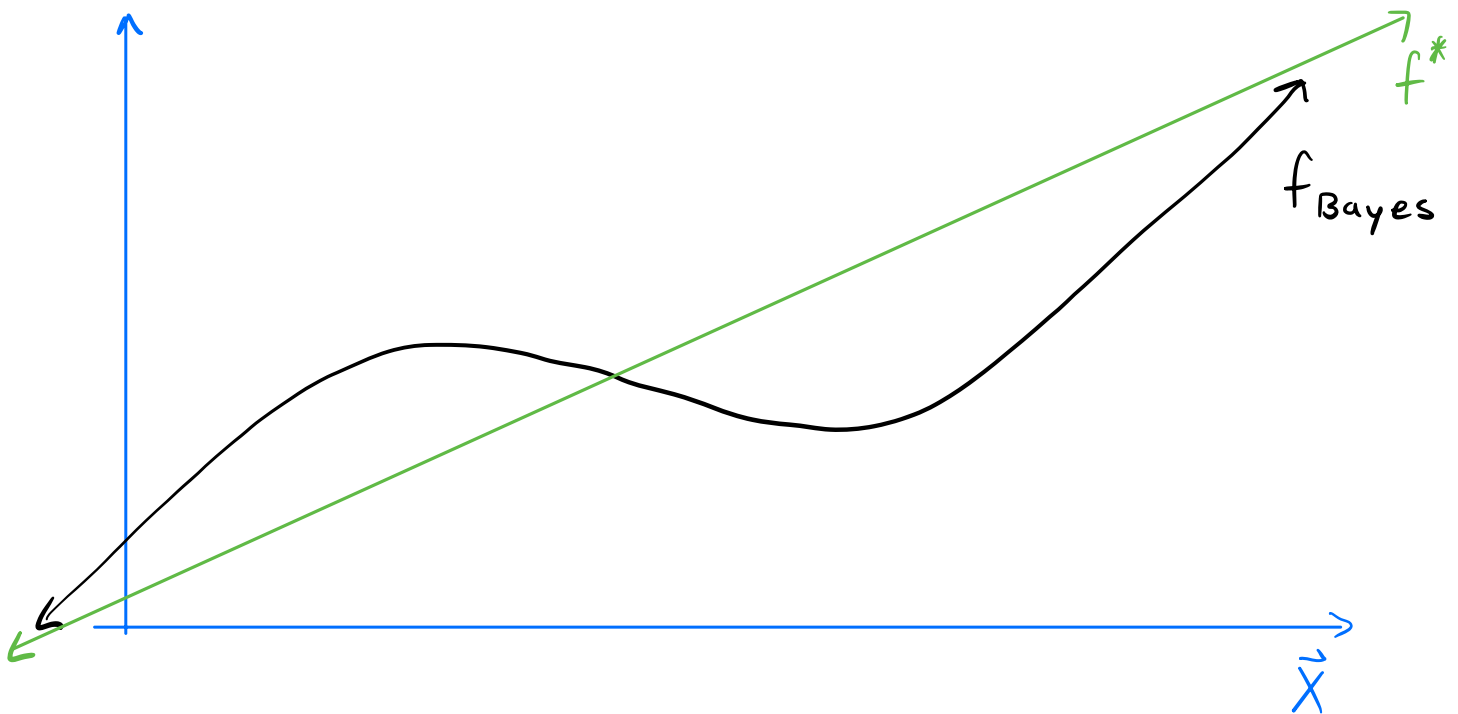
$$f_{\text{Bayes}} = \underset{f \in \{f | f: \mathcal{X} \rightarrow \mathcal{Y}\}}{\text{argmin}} L(f) \quad \text{"Bayes optimal predictor"}$$



Suppose we knew $P_{\vec{x}, y}$ but $\mathcal{A}: (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathcal{F}_l$

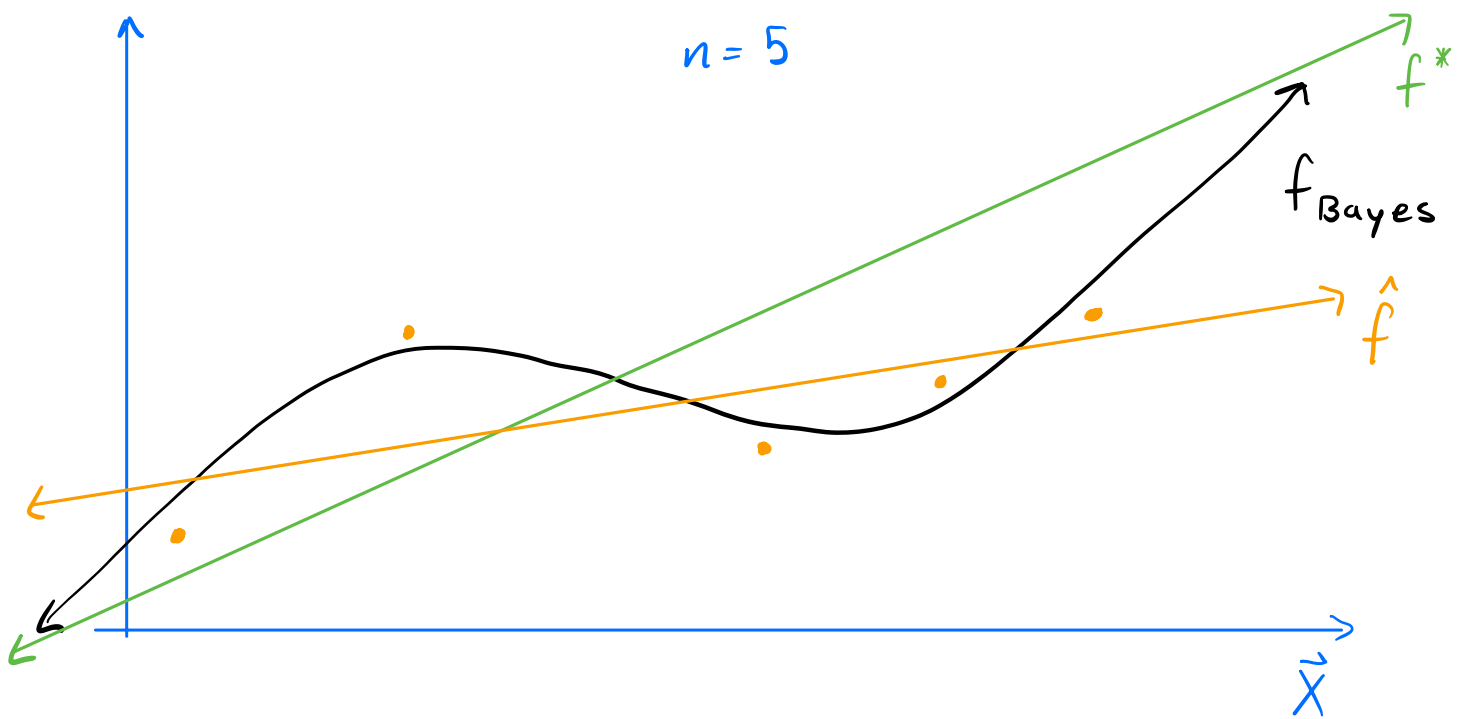
linear functions

$$f^* = \operatorname{argmin}_{f \in \mathcal{F}_l} L(f)$$



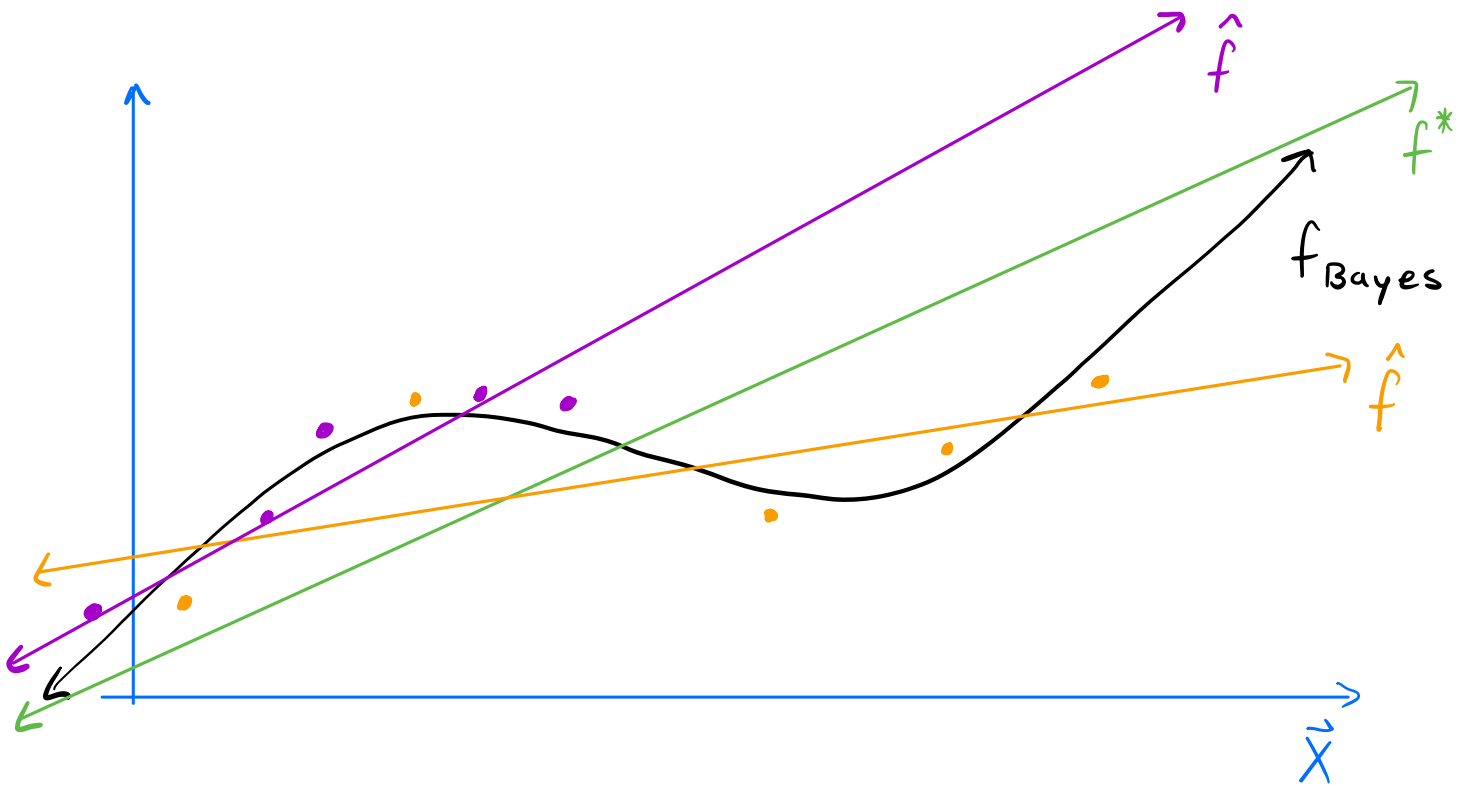
Suppose we didn't know $P_{\vec{x}, y}$ but we had a dataset $D_l = ((\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n))$ and $\mathcal{A}: (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathcal{F}_l$

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}_l} \hat{L}(f)$$

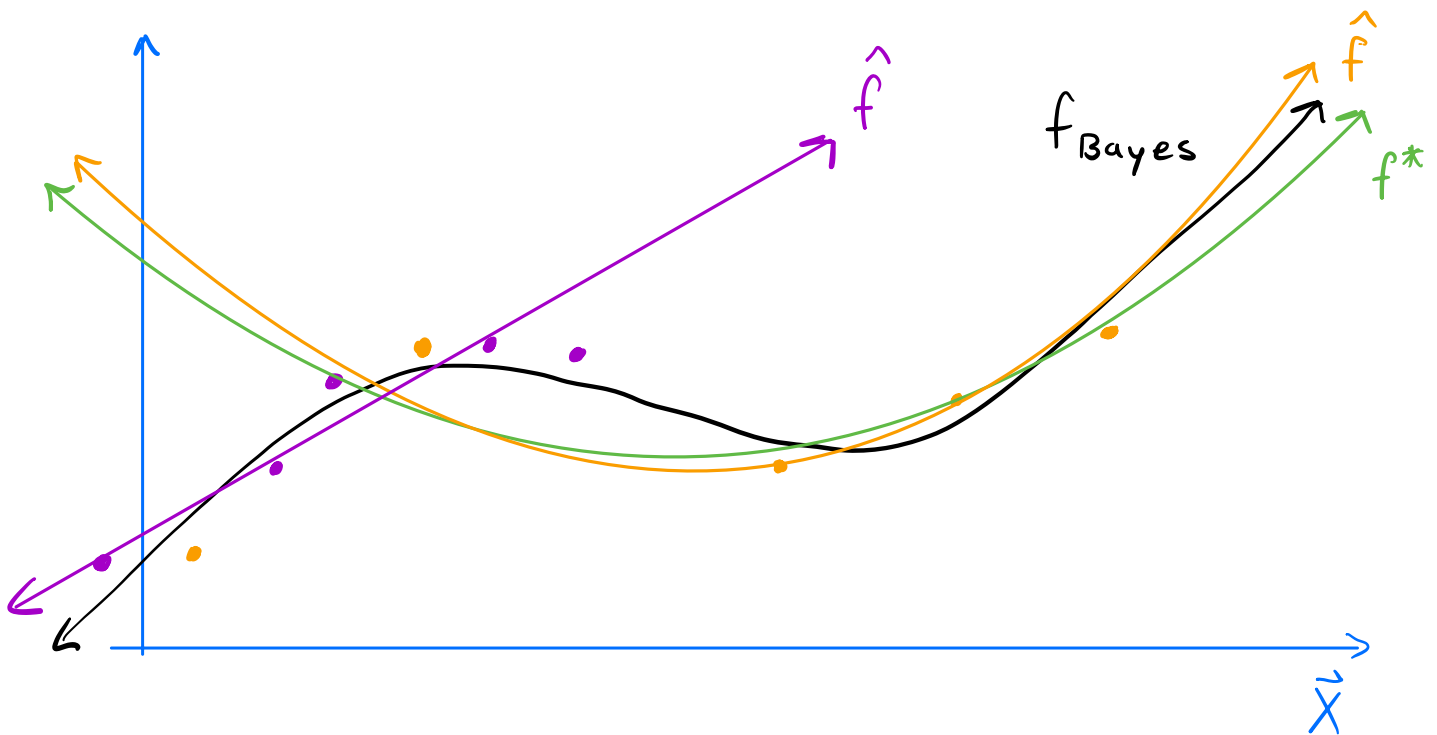


How about for a different dataset

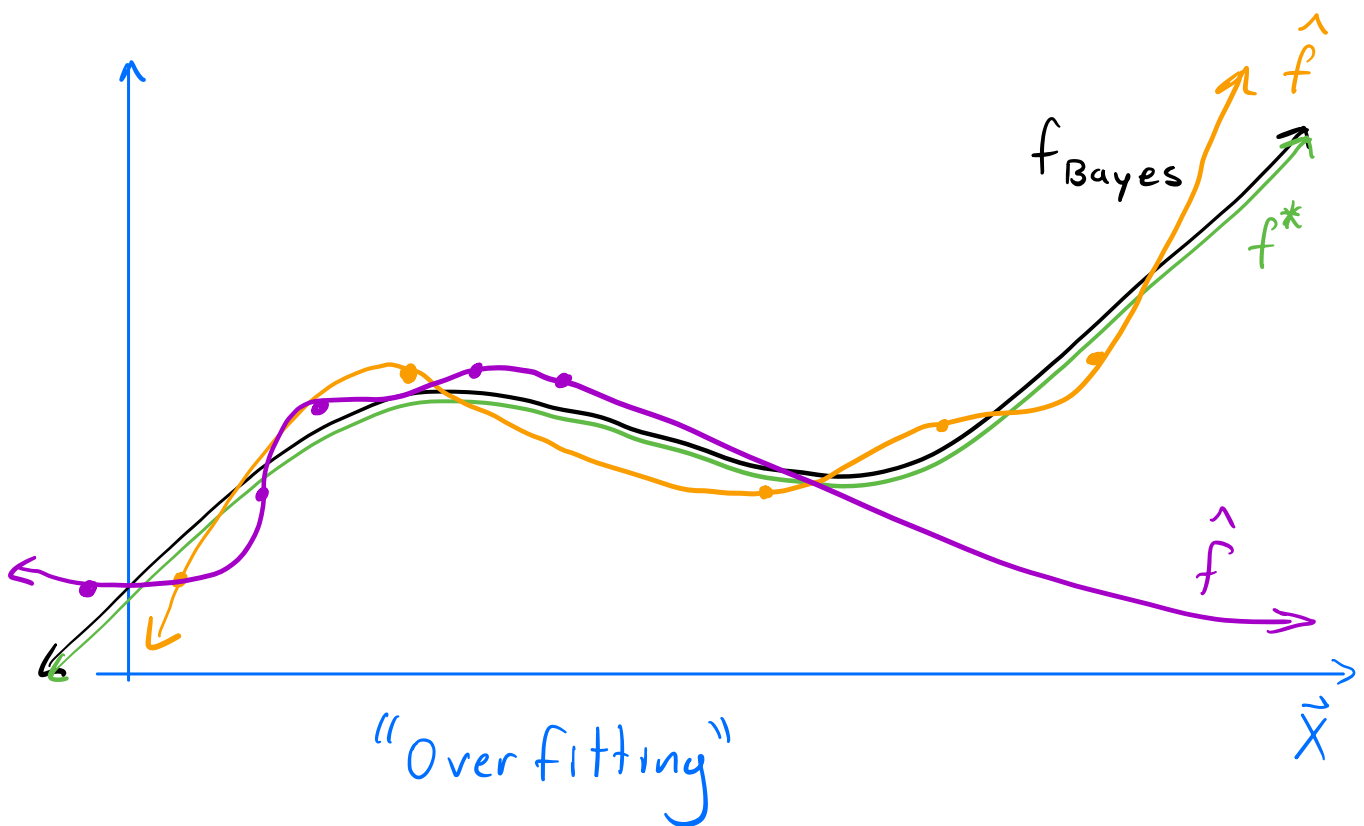
$$D_2 = ((\bar{x}_1, y_1), \dots, (\bar{x}_n, y_n))$$



How about for \tilde{F}_2



How about if \tilde{F}_{10}



Decomposing $E[L(\mathcal{A}(D))]$

$$\text{Let } \mathcal{A}(D) = \hat{f}_D$$

\mathcal{F} any function class

$$f_{\text{Bayes}} = \operatorname{argmin}_{f \in \{f | f: x \rightarrow y\}} L(f)$$

$$f^* = \operatorname{argmin}_{f \in \mathcal{F}} L(f)$$

$$\hat{f}_D = \operatorname{argmin}_{f \in \mathcal{F}} \hat{L}(f)$$

$$E[L(\hat{f}_D)] = \underbrace{E[L(\hat{f}_D)] - L(f^*)}_{\text{Estimation Error (EE)}} + \underbrace{L(f^*) - L(f_{\text{Bayes}})}_{\text{Approximation Error (AE)}} + \underbrace{L(f_{\text{Bayes}})}_{\text{Irreducible Error (IE)}}$$

What affects the different types of errors?

Irreducible Error: Due to inherent noise in labels

- Decreases if you gather more/better feature info
- Usually not possible to do "irreducible"

Approximation Error: Due to a small \mathcal{F}

- Decreases if you make \mathcal{F} larger

Estimation Error: Due to random dataset D

- Decreases if you increase n
- Increases if you increase \mathcal{F}

High EE: small n , large \mathcal{F}

High AE: f_{Bayes} complex, \mathcal{F} simple

Understanding EE:

$$\hat{L}(\hat{f}_0) \leq \hat{L}(f) \text{ for all } f \in \mathcal{F} \\ \leq 0$$

$$\text{EE: } \mathbb{E}[L(\hat{f}_0)] - L(f^*)$$

$$= \mathbb{E}[L(\hat{f}_0)] - \mathbb{E}[\hat{L}(\hat{f}_0)] + \mathbb{E}[\hat{L}(\hat{f}_0)] - \mathbb{E}[\hat{L}(f^*)] \\ + \mathbb{E}[\hat{L}(f^*)] - L(f^*)$$

$$\leq \underbrace{\mathbb{E}[L(\hat{f}_0)] - \mathbb{E}[\hat{L}(\hat{f}_0)]}_{\text{decreases as } n \text{ increases}} + \mathbb{E}[\hat{L}(f^*)] - L(f^*)$$

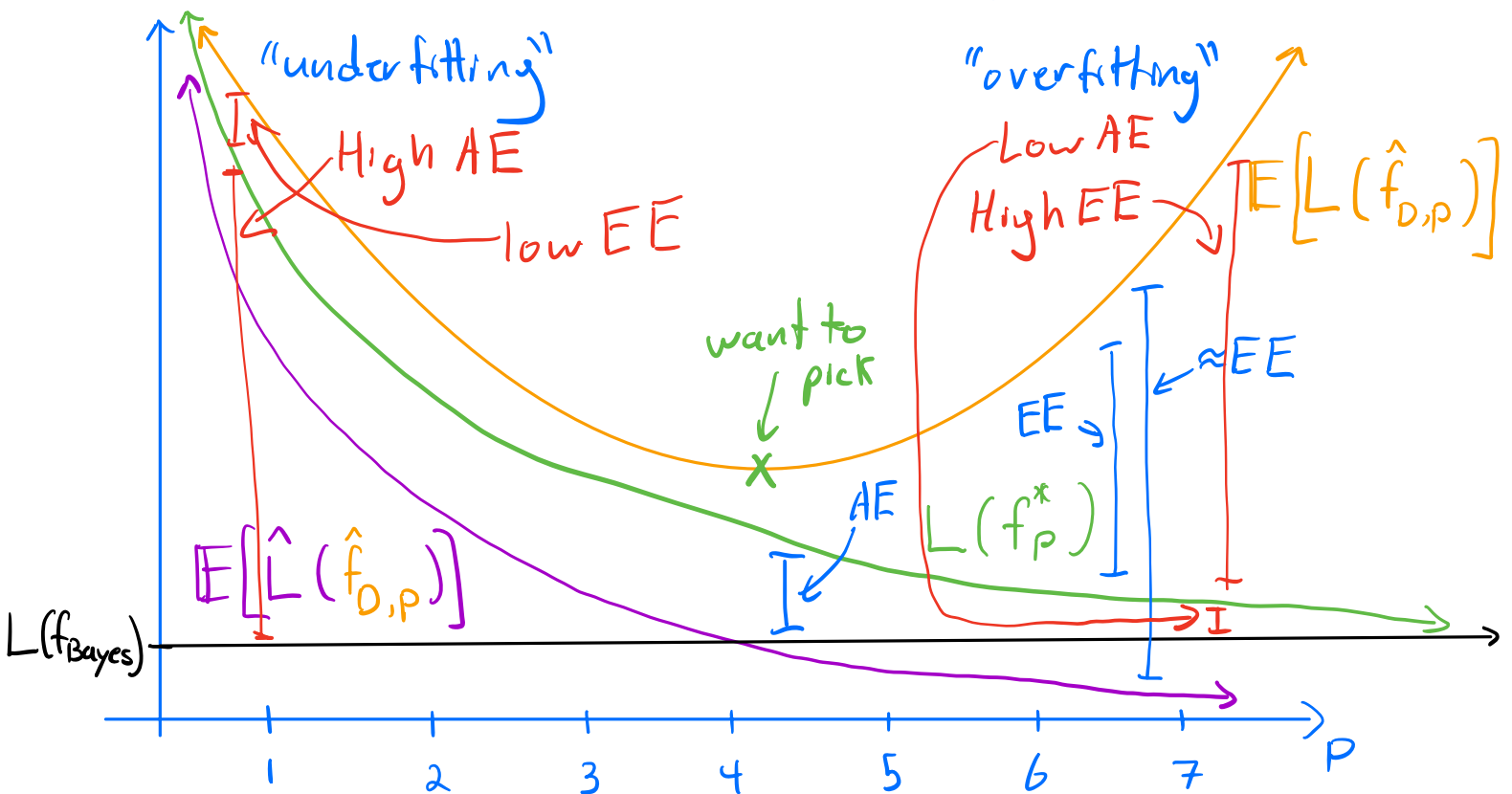
decreases as n increases

increases as \mathcal{F} gets larger

$$A(D) = \hat{f}_{D,p} = \underset{f \in \tilde{F}_p}{\operatorname{argmin}} \hat{L}(f) \quad \tilde{F}_1 \subset \dots \subset \tilde{F}_p$$

$$E[L(\hat{f}_0)] - E[\hat{L}(\hat{f}_0)]$$

$$E[L(\hat{f}_0)] = \underbrace{E[L(\hat{f}_0)] - L(f^*)}_{\text{Estimation Error (EE)}} + \underbrace{L(f^*) - L(f_{\text{Bayes}})}_{\text{Approximation Error (AE)}} + \underbrace{L(f_{\text{Bayes}})}_{\text{Irreducible Error (IE)}}$$



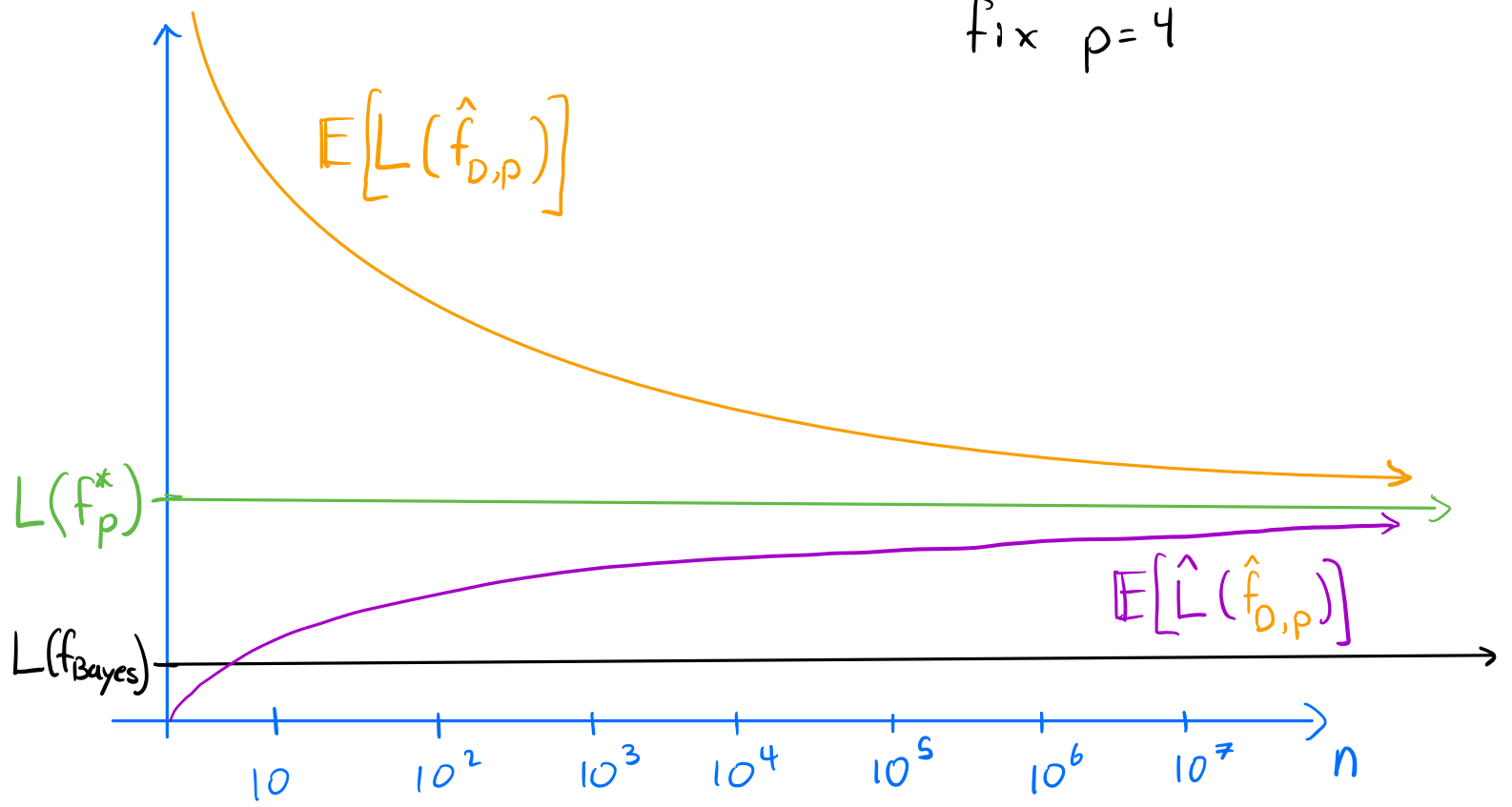
Underfitting: \tilde{F} is too simple (small) compared to n

- High AE, Low EE

Overfitting: \tilde{F} is too complex (large) compared to n

- Low AE, High EE

fix $\rho=4$



In practice we only have a fixed dataset D

How can we tell if we are overfitting or underfitting if we can't calculate $L(\hat{f}_0)$?

Estimate $L(\hat{f}_0)$ with a different dataset D_{test}

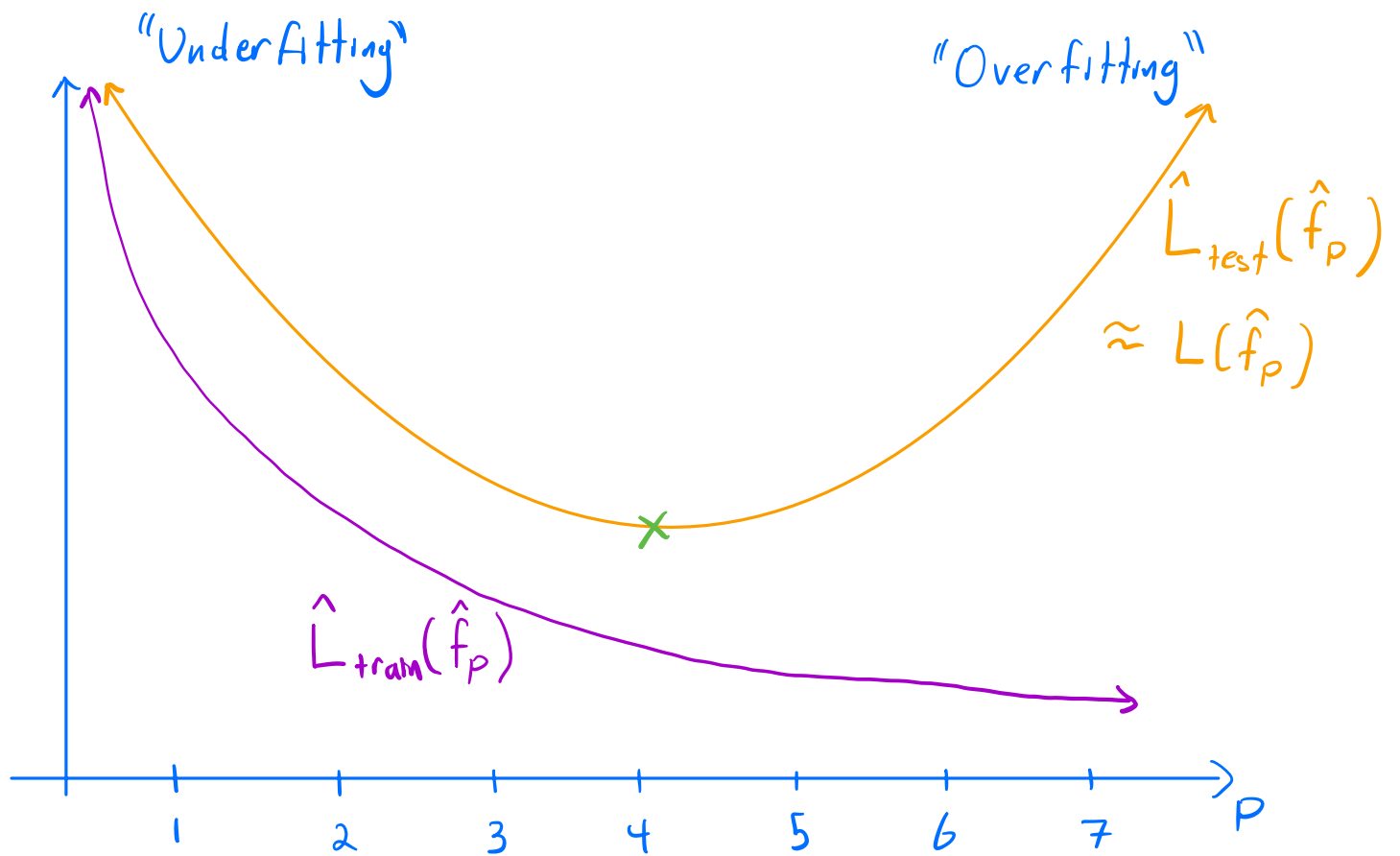
Since we can't gather new data we split D into $D_{\text{train}}, D_{\text{test}}$

$$D_{\text{train}} = ((\vec{x}_1, y_1), \dots, (\vec{x}_{n-m}, y_{n-m}))$$

$$D_{\text{test}} = ((\vec{x}_{n-m+1}, y_{n-m+1}), \dots, (\vec{x}_n, y_n))$$

$$|D_{\text{train}}| = n - m, \quad |D_{\text{test}}| = m$$

$$A(\mathcal{D}_{\text{train}}) = \hat{f}_P = \underset{f \in \tilde{\mathcal{F}}_P}{\text{argmin}} \hat{L}_{\text{train}}(f) \quad \tilde{\mathcal{F}}_1 \subset \dots \subset \tilde{\mathcal{F}}_P$$



Bias-Variance Tradeoff

$\mathbb{E}[L(\hat{f}_D)]$ Let l be squared loss

$$= \mathbb{E}[L(\hat{f}_D)] - L(f_{\text{Bayes}}) + L(f_{\text{Bayes}})$$

if $\bar{f} = f^*$

$$= \underbrace{\mathbb{E}[\mathbb{E}[(\hat{f}_D(\vec{x}) - \bar{f}(\vec{x}))^2 | \vec{x}]]]}_{= EE} + \underbrace{\mathbb{E}[(\bar{f}(\vec{x}) - f_{\text{Bayes}}(\vec{x}))^2]}_{= AE} + \underbrace{L(f_{\text{Bayes}})}_{IE}$$

Variance = $\text{Var}[\hat{f}_D(\vec{x}) | \vec{x}]$ Bias IE

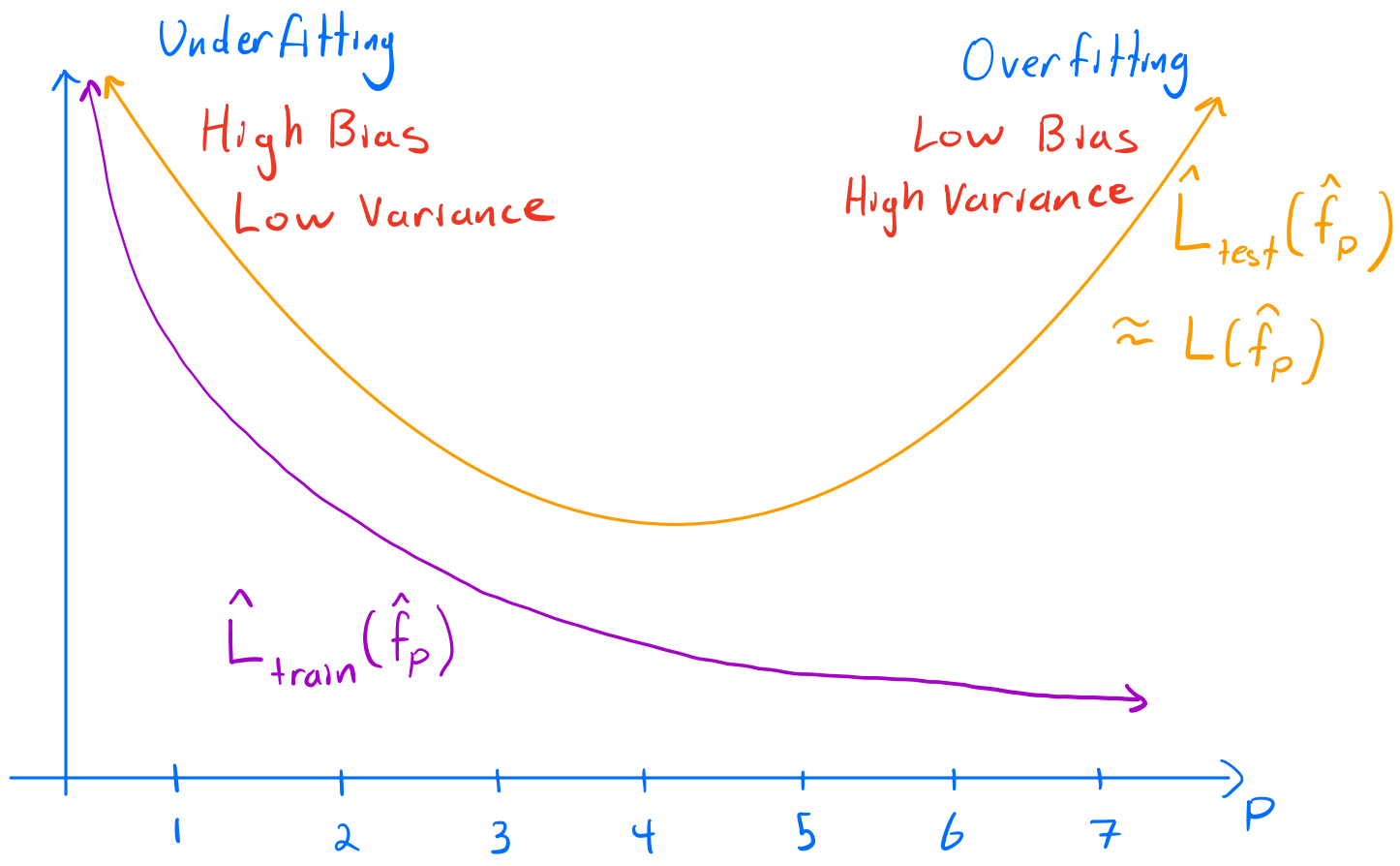
where $\bar{f}(\vec{x}) = \mathbb{E}[\hat{f}_D(\vec{x}) | \vec{x}]$ "expected predictor"

Effects of changing F, n on Bias, Variance follow the same trend as for AE, EE:

Bias \downarrow if $F \uparrow$

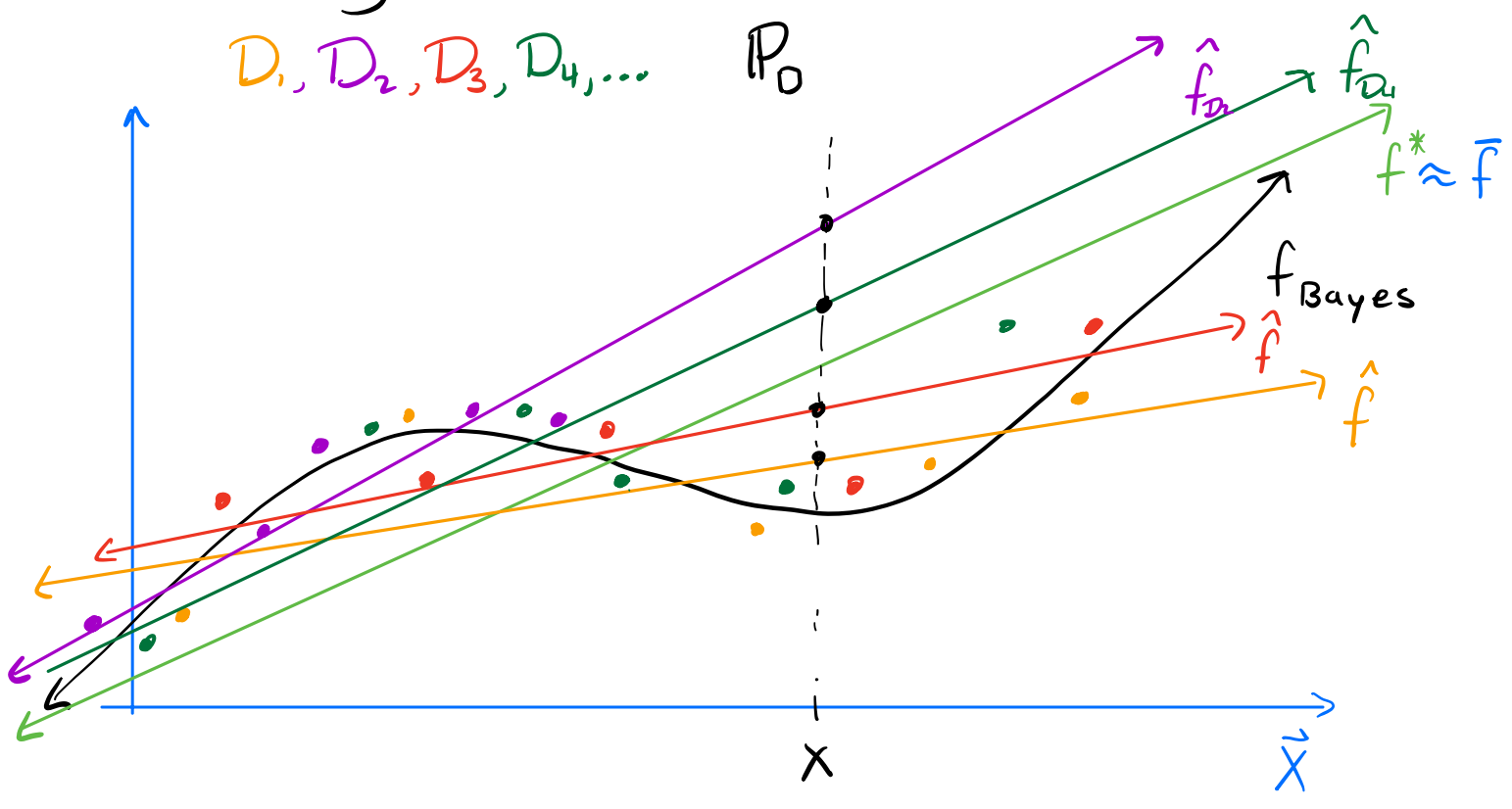
Variance \downarrow if $n \uparrow$

\uparrow if $F \uparrow$



Visualizing $\bar{f}(x) = \mathbb{E}[\hat{f}_D(x)|X]$

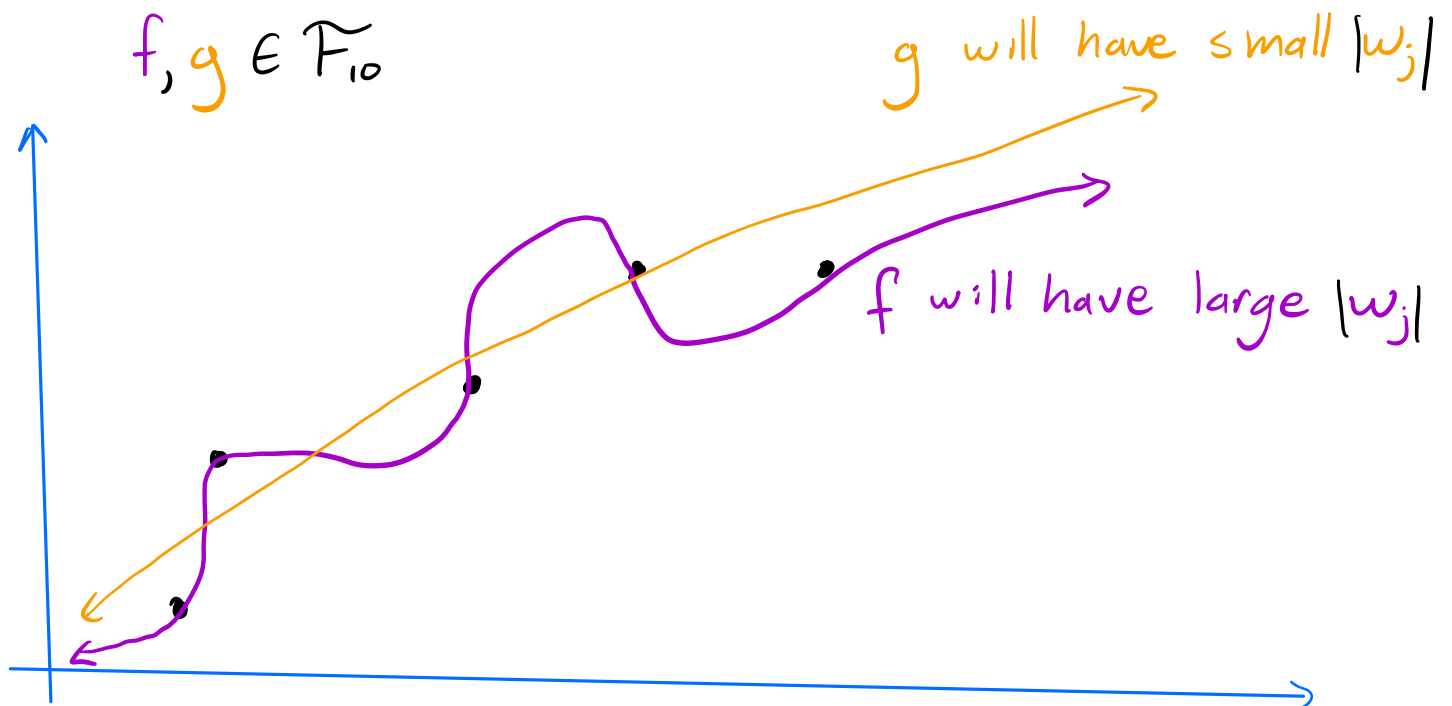
$D_1, D_2, D_3, D_4, \dots$ P_0



Regularization

let $\vec{w} = (w_0, w_1, \dots, w_{p-1})^T \in \mathbb{R}^p$

Observation: large values of $|w_0|, |w_1|, \dots, |w_{p-1}|$ leads to more complex $f_p(\vec{x}) = \phi_p(\vec{x})^T \vec{w}$



Regularization: penalize large weights

If $f_p \in \mathcal{F}_p$:

$$= \hat{L}(f_p)$$

$$\hat{L}_\lambda(f_p) = \frac{1}{n} \sum_{i=1}^n \ell(f_p(\vec{x}_i), y_i) + \frac{\lambda}{n} \sum_{j=1}^{p-1} w_j^2$$

Regularization parameter $\lambda \in (0, \infty)$

Regularized estimated loss

$j=0$ not included "Regularizer"

Minimizing $\hat{L}_\lambda(f)$ instead of $\hat{L}(f)$ is called
"Ridge Regression"

Let $\hat{f}_\lambda = \underset{f \in \mathcal{F}}{\operatorname{arg\,min}} \hat{L}_\lambda(f)$, $f^* = \underset{f \in \mathcal{F}}{\operatorname{arg\,min}} L(f)$

If λ increases, then \hat{f}_λ gets simpler fix \mathcal{F}

, \bar{f}_λ gets simpler

, but f^* does not change

$\bar{f}_\lambda \neq f^*$ unless $\lambda = 0$

Bias vs. Variance

Bias: $(\bar{f}_\lambda(\vec{X}) - f_{\text{Bayes}}(\vec{X}))^2$

- Decreases if λ decreases

Variance: $E[(\hat{f}_{0,\lambda}(\vec{X}) - \bar{f}_\lambda(\vec{X}))^2 | \vec{X}]$

- Increases if λ decreases
- Decreases if n increases

Minimizing $\hat{L}_\lambda(f)$

$$\hat{\vec{w}}_\lambda = \arg \min_{\vec{w} \in \mathbb{R}^{d+1}} \hat{L}_\lambda(\vec{w}) \quad \text{using squared loss, } \mathcal{F}_1$$

$$\text{where } \hat{L}_\lambda(\vec{w}) = \underbrace{\frac{1}{n} \sum_{i=1}^n (\vec{x}_i^T \vec{w} - y_i)^2}_{\hat{L}(\vec{w})} + \underbrace{\frac{\lambda}{n} \sum_{j=1}^d w_j^2}_{g(\vec{w})}$$

There is a closed form solution

but it is more complicated so we use gradient descent to find the minimum instead

$$\vec{w}^{(t+1)} = \vec{w}^{(t)} - \eta^{(t)} \nabla \hat{L}_\lambda(\vec{w}^{(t)})$$

$$\nabla \hat{L}_\lambda(\vec{w}^{(t)}) = \left(\frac{\partial \hat{L}_\lambda}{\partial w_0}(\vec{w}^{(t)}), \dots, \frac{\partial \hat{L}_\lambda}{\partial w_d}(\vec{w}^{(t)}) \right)^T$$

$$\frac{\partial \hat{L}_\lambda}{\partial w_k}(\vec{w}) = \frac{\partial \hat{L}}{\partial w_k}(\vec{w}) + \frac{\partial g}{\partial w_k}(\vec{w}) \quad k \in \{0, \dots, d\}$$

$$\frac{\partial g}{\partial w_0}(\vec{w}) = 0$$

$$\frac{\partial g}{\partial w_k}(\vec{w}) = \frac{\lambda}{n} \cdot 2w_k \quad k = \{1, \dots, d\}$$

$$\frac{\partial L_\lambda(\vec{w})}{\partial w_0} = \frac{2}{n} \sum_{i=1}^n (\vec{x}_i^T \vec{w} - y_i) x_{i0}$$

$$\frac{\partial L_\lambda(\vec{w})}{\partial w_k} = \frac{2}{n} \sum_{i=1}^n (\vec{x}_i^T \vec{w} - y_i) x_{ik} + \frac{2\lambda}{n} w_k \quad k \in \{1, \dots, d\}$$

$$\hat{f}_\lambda = \operatorname{argmin}_{f \in \mathcal{F}_\lambda} \hat{L}_\lambda(f)$$

