- BGD: $\nabla \hat{L}(\vec{w}) = \left( \frac{\partial}{\partial w_0} \hat{L}(\vec{v}), \ldots \right.$

- $\mathcal{O}(n^3)$ for inverse calculation

- MBGD      $\nabla \hat{L}_m(\vec{w})$

# Evaluating Predictors/Models
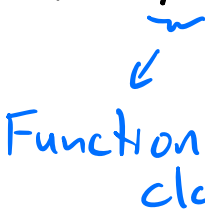
## Objective (formal):

Define a learner $\mathcal{A}: (\mathcal{X} \times \mathcal{Y})^n \Rightarrow \{f \mid f: \mathcal{X} \Rightarrow \mathcal{Y}\}$

such that $\mathbb{E}\left[L(\mathcal{A}(D))\right]$ 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 $\hat{L}(f)$

↳ 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 $L(\hat{f})$
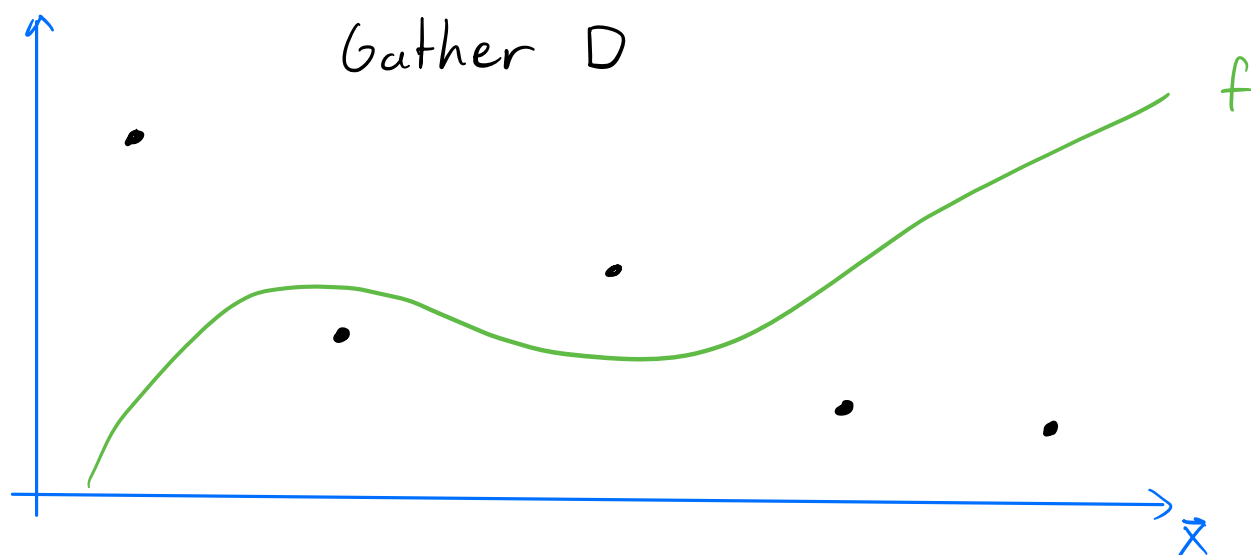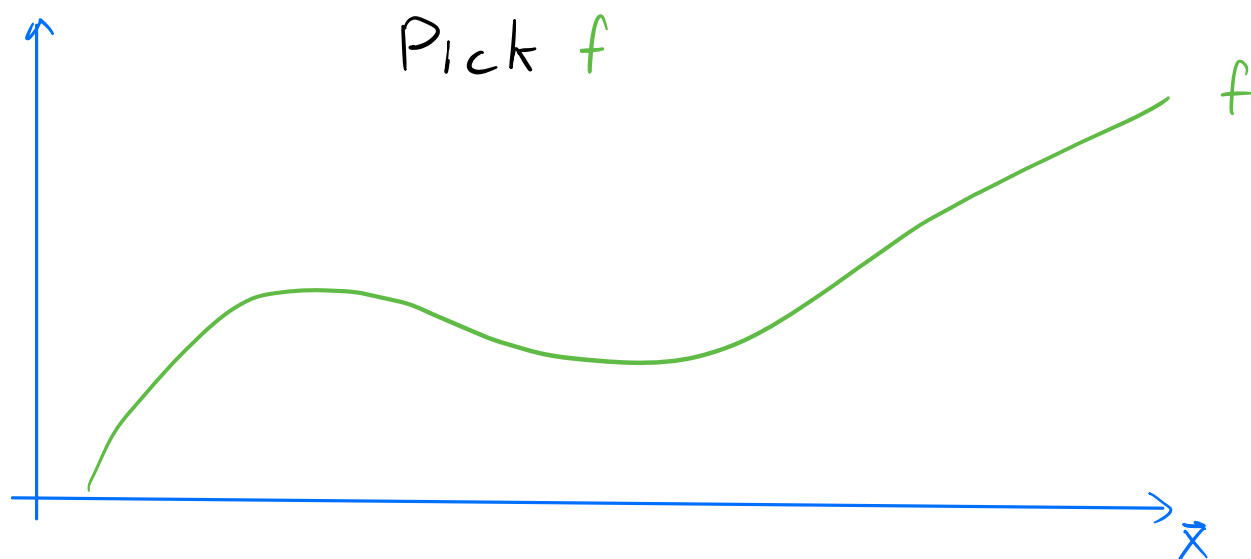
# Evaluating Predictors/Models

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

where $\mathcal{A}(D) = f_D \in \mathcal{F}$   $D$ is a r.v.   $(\vec{X}, Y) \sim \mathbb{P}_{\vec{X}, Y}$

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

(estimate of $L(f)$)   (expected loss)

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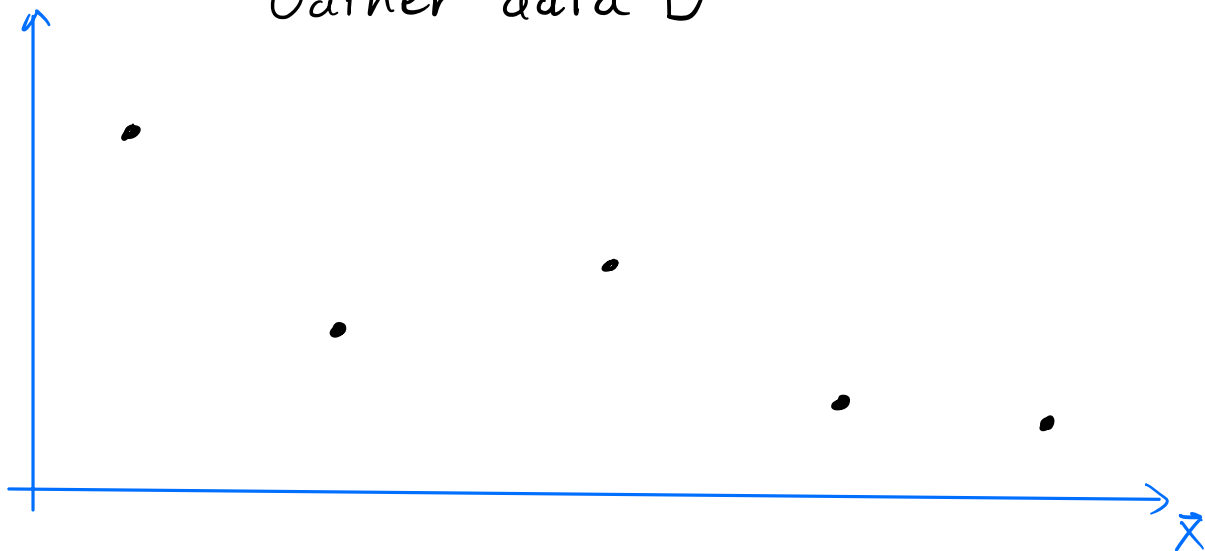
(i.e. $f$ is chosen independently of $D$)

Pick $f$



Gather $D$

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

$$\text{Var}[\hat{L}(f)] = \text{Var}\left[\frac{1}{n}\sum_{i=1}^{n}\ell(f(\vec{X}_i), Y_i)\right]$$

Since $\ell(f(\vec{X}_i), Y_i)$ are independent for all $i \in \{1, ..., n\}$

$$= \frac{1}{n^2}\sum_{i=1}^{n}\text{Var}\left[\ell(f(\vec{X}_i), Y_i)\right]$$

$$= \frac{1}{n}\text{Var}\left[\ell(f(\vec{X}_1), Y_1)\right]$$

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

Gather data $D$



Pick $\hat{f}_D \in \mathcal{F}$ based on $D$

problem:
$$L(\hat{f}_D) \geq \hat{L}(\hat{f}_D) \approx 0$$



$\hat{f}_D$

Then:
$$\mathbb{E}[\hat{L}(\hat{f}_D)] = \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\ell(f(\vec{X}_i), Y_i)\right] = \frac{1}{n}\sum_{i=1}^{n}\mathbb{E}[\ell(f(\vec{X}_i), Y_i)]$$
$$\neq \mathbb{E}[\ell(f(\vec{X}_1), Y_1)]$$

$$\mathrm{Var}[\hat{L}(\hat{f}_D)] = \mathrm{Var}\left[\frac{1}{n}\sum_{i=1}^{n}\ell(\hat{f}_D(\vec{X}_i, Y_i))\right]$$
$$\neq \frac{1}{n^2}\sum_{i=1}^{n}\mathrm{Var}\left[\ell(\hat{f}_D(\vec{X}_i), Y_i)\right]$$

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

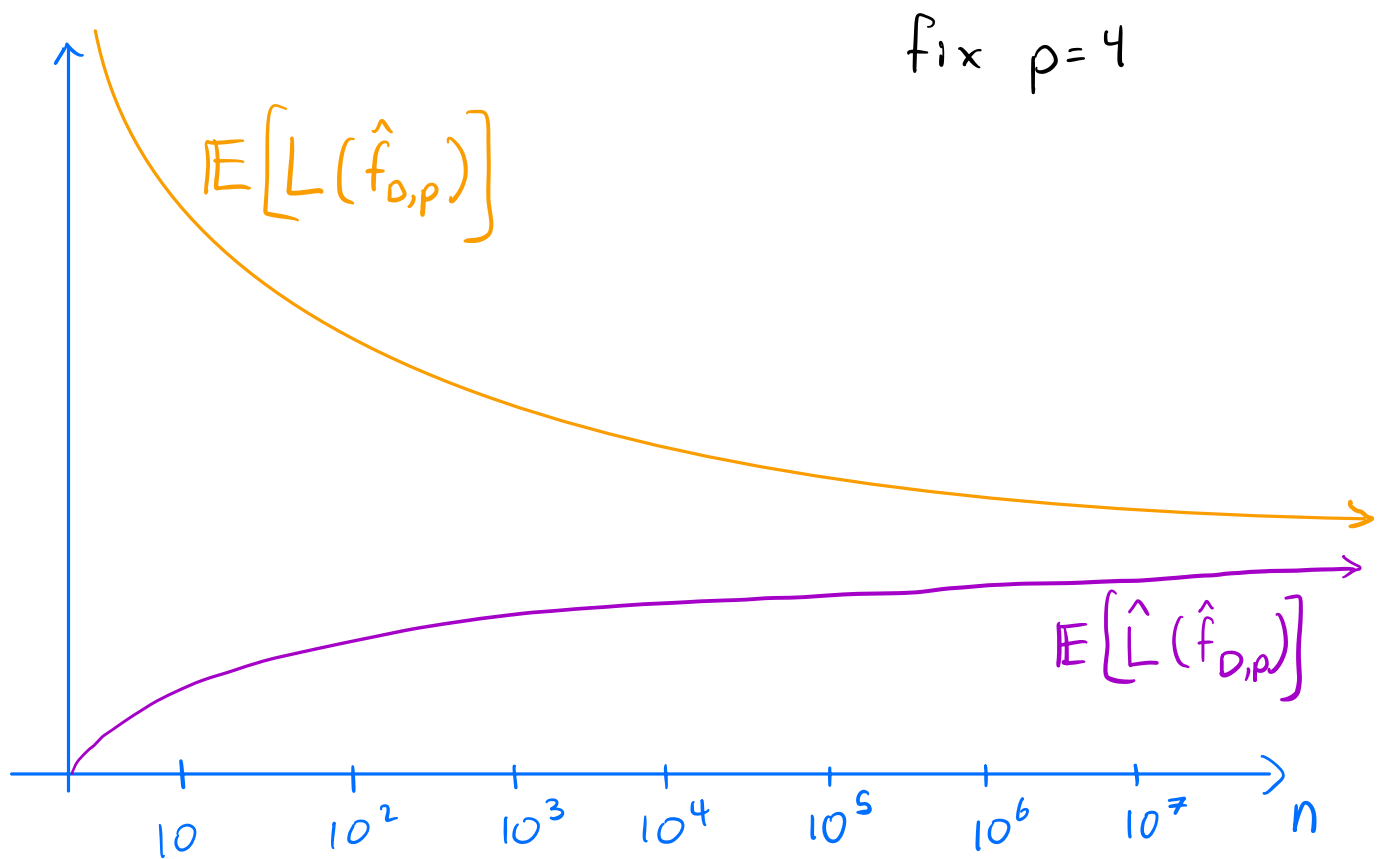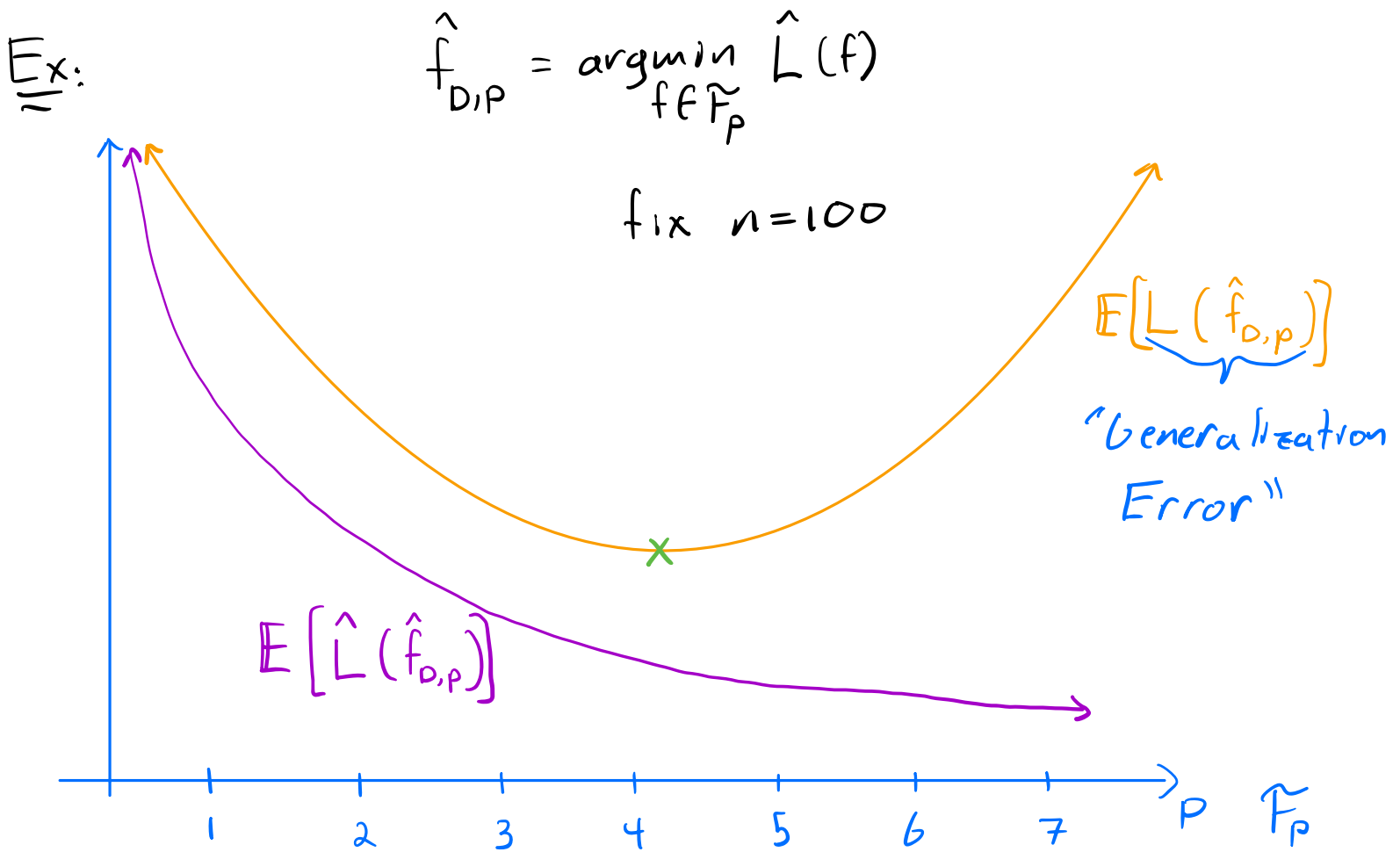$\hat{f}_D$ depends on $(\vec{X}_1, Y_1), \ldots, (\vec{X}_n, Y_n)$!

Instead:

This difference
→ increases as $\mathcal{F}$ gets more complex

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

→ decreases as $n$ increases

$$\mathrm{Var}[\hat{L}(\hat{f}_D)] \leq \frac{1}{n}\max_{f\in\mathcal{F}}\mathrm{Var}[\ell(f(\vec{X}_1), Y_1)] + \mathrm{Var}[L(\hat{f}_D)]$$

As $\mathcal{F}$ gets larger, $\mathrm{Var}[\hat{L}(\hat{f}_D)]$ increases

$\Rightarrow$ i.e. $\hat{L}(\hat{f}_D)$ becomes a worse estimate of $L(\hat{f}_D)$

As $n$ gets larger, $\mathrm{Var}[\hat{L}(\hat{f}_D)]$ decreases

$\Rightarrow$ $\hat{L}(\hat{f}_D)$ becomes a better estimate of $L(\hat{f}_D)$

Ex:

$$\hat{f}_{D,P} = \underset{f \in \tilde{F}_P}{\arg\min}\ \hat{L}(f)$$

fix $n = 100$

$\mathbb{E}\left[L(\hat{f}_{D,P})\right]$

"Generalization Error"

$\mathbb{E}\left[\hat{L}(\hat{f}_{D,P})\right]$

$P$   $\tilde{F}_P$

1   2   3   4   5   6   7

fix $P = 4$

$\mathbb{E}\left[L(\hat{f}_{D,P})\right]$

$\mathbb{E}\left[\hat{L}(\hat{f}_{D,P})\right]$

$10$   $10^2$   $10^3$   $10^4$   $10^5$   $10^6$   $10^7$   $n$

## Objective (formal):

Define a learner $\mathcal{A} : (\mathcal{X} \times \mathcal{Y})^n \Rightarrow \{f \mid f : \mathcal{X} \Rightarrow \mathcal{Y}\}$

such that $\mathbb{E}\left[L(\mathcal{A}(D))\right]$ is small over the datasets
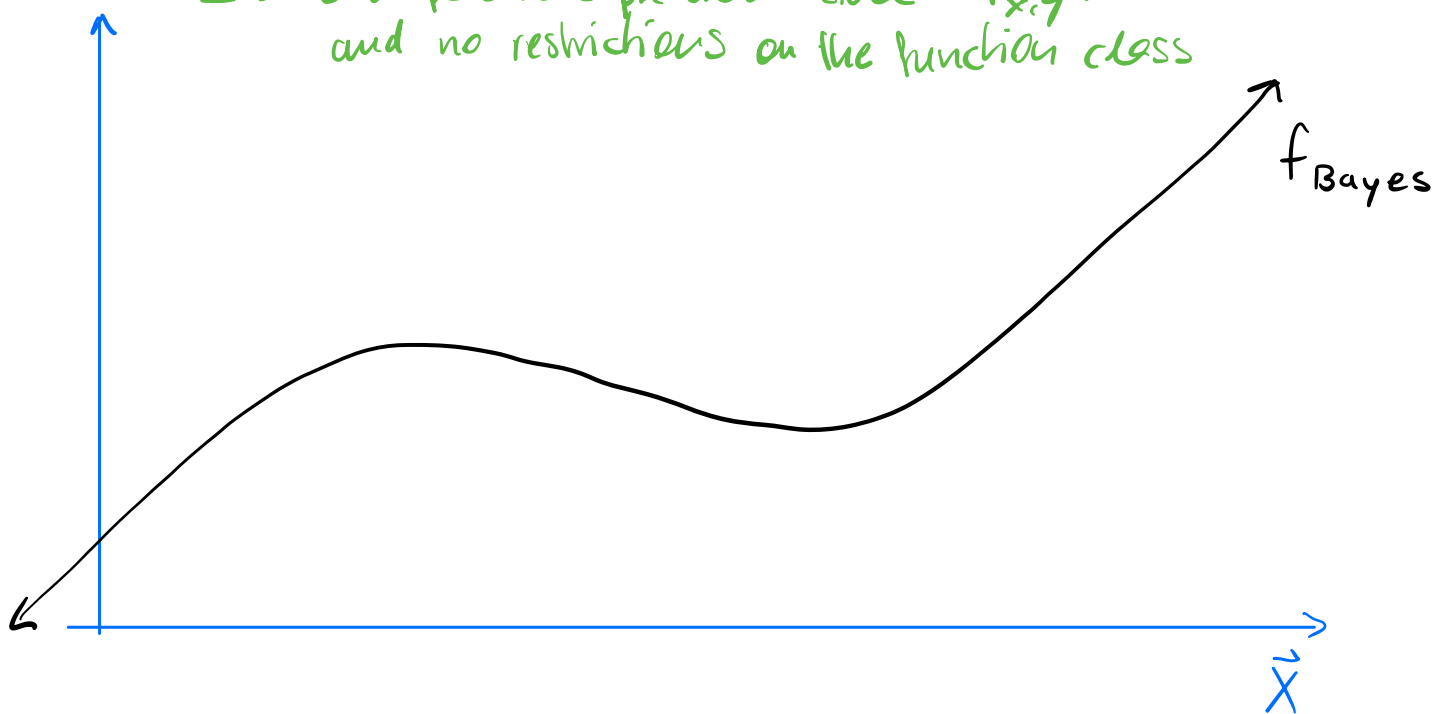
(expected loss)

Suppose we knew $\mathbb{P}_{\vec{X}, Y}$ what would we choose for $\mathcal{A}$?

(true distribution of $\vec{X}, Y$)

$$f_{Bayes} = \underset{f \in \{f \mid f : \mathcal{X} \Rightarrow \mathcal{Y}\}}{\arg\min} L(f) \qquad \text{"Bayes optimal predictor"}$$

(true expected loss)

no restriction on the function class

where $L(f) = \mathbb{E}\left[\ell(f(\vec{x}), Y)\right]$ and $(\vec{X}, Y) \sim \mathbb{P}_{\vec{X}, Y}$

$\Rightarrow$ best possible predictor since $\mathbb{P}_{\vec{X}, Y}$ known and no restrictions on the function class



$f_{Bayes}$

$\vec{X}$

Suppose we knew $\mathbb{P}_{\vec{X},Y}$ but $\mathcal{A} : (X \times Y)^n \to F_1$

linear functions

$f^* = \underset{f \in F_1}{\text{argmin}} \; L(f)$

best possible predictor in the function class



$f^*$

$f_{Bayes}$

$\vec{X}$

Suppose we didn't know $\mathbb{P}_{\vec{X},Y}$ but we had a dataset $\mathcal{D}_1 = ((\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n))$ and $\mathcal{A} : (X \times Y)^n \to F_1$

our original setting

$\hat{f} = \underset{f \in F_1}{\text{argmin}} \; \hat{L}(f)$

$n = 5$

(estimation of the true expected loss)

How about for a different dataset
$$D_2 = ((\vec{x}_1, y_1), \ldots, (\vec{x}_n, y_n))$$

# How about for $F_2$



$\hat{f}$

$f_{Bayes}$

$f^*$

$\vec{X}$

# How about if $F_{10}$



$\hat{f}$

$f_{Bayes}$

$f^*$

$\hat{f}$

$\vec{X}$

"Overfitting"

$\Rightarrow \hat{f}$ and $\hat{f}$ very different for different samples of the dataset

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

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

$\mathcal{F}$ any function class

$$f_{Bayes} = \underset{f \in \{f | f : x \to y\}}{\arg\min} L(f)$$

$$f^* = \underset{f \in \mathcal{F}}{\arg\min} L(f)$$

$$\hat{f}_D = \underset{f \in \mathcal{F}}{\arg\min} \hat{L}(f)$$

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

## 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

why EE↑ if $\nearrow$ $n$ decreases
$\searrow$ $\mathcal{F}$ more complex

## Understanding EE:

EE: $\mathbb{E}\left[L(\hat{f}_D)\right] - L(f^*)$

$$\mathcal{A}(D) = \hat{f}_{D,P} = \underset{f \in \tilde{F}_P}{\arg\min} \hat{L}(f) \qquad \tilde{F}_1 \subset \cdots \subset \tilde{F}_P$$

$$\mathbb{E}\left[L(\hat{f}_D)\right] - \mathbb{E}\left[\hat{L}(\hat{f}_D)\right]$$

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



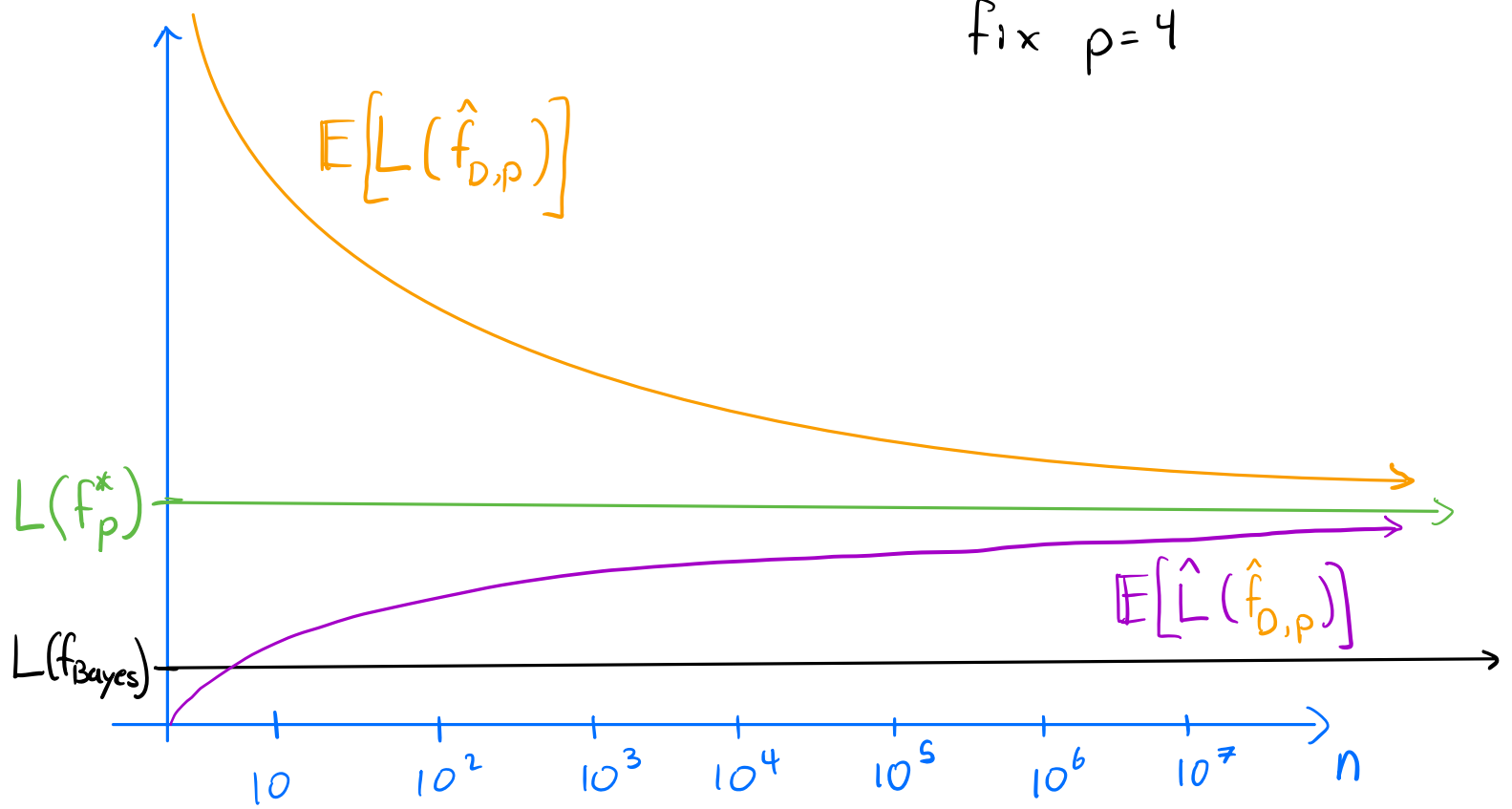<u>Underfitting:</u> $\mathcal{F}$ is too simple (small) compared to $n$

    — High AE, Low EE

<u>Overfitting:</u> $\mathcal{F}$ is too complex (large) compared to $n$

    — Low AE, High EE

fix $\rho = 4$

$E\left[L(\hat{f}_{D,\rho})\right]$

$L(f_\rho^*)$

$L(f_{Bayes})$

$E\left[\hat{L}(\hat{f}_{D,\rho})\right]$

$10 \quad 10^2 \quad 10^3 \quad 10^4 \quad 10^5 \quad 10^6 \quad 10^7 \quad n$

In practice we only have a fixed dataset $D$
How can we tell if we are overfitting or
underfitting if we can't calculate $L(\hat{f}_D)$?

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

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

$$D_{train} = \left( (\vec{X}_1, y_1), \ldots, (\vec{X}_{n-m}, y_{n-m}) \right)$$

$$D_{test} = \left( (\vec{X}_{n-m+1}, y_{n-m+1}), \ldots, (\vec{X}_n, y_n) \right)$$

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

$$\mathcal{A}(D_{train}) = \hat{f}_P = \underset{f \in \widetilde{F_P}}{arg\,min} \; \hat{L}_{train}(f) \qquad \widetilde{F_1} \subset \cdots \subset \widetilde{F_P}$$



"Underfitting"

"Overfitting"

$\hat{L}_{test}(\hat{f}_P)$

$\approx L(\hat{f}_P)$

$\hat{L}_{train}(\hat{f}_P)$

1   2   3   4   5   6   7   $P$

# Bias-Variance Tradeoff

$$\mathbb{E}\left[L(\hat{f}_D)\right]$$
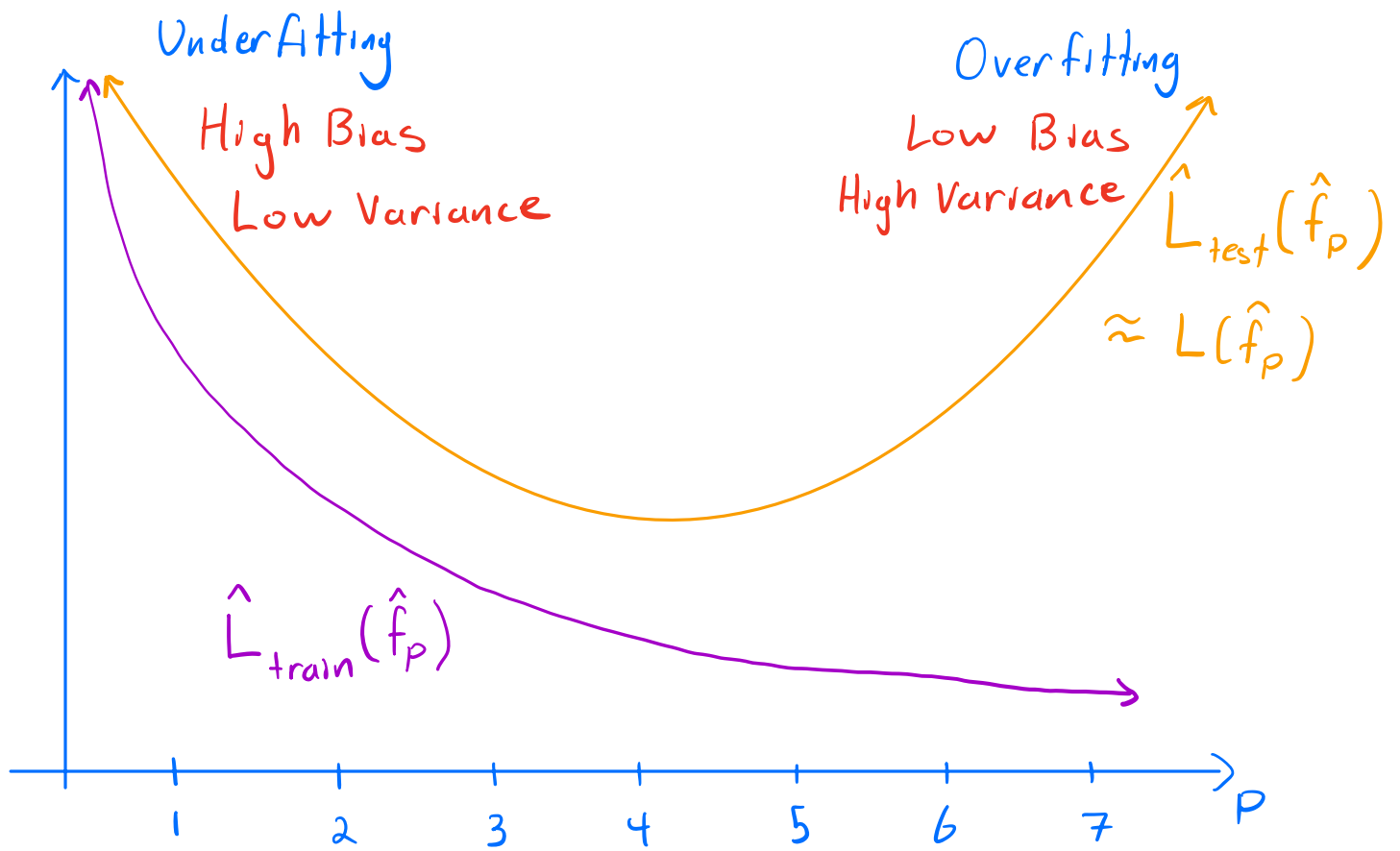
Effects of changing $\tilde{F}, n$ on Bias, Variance follow the same trend as for AE, EE:

Bias $\downarrow$      if $\tilde{F} \uparrow$

Variance $\downarrow$      if $n \uparrow$

$\uparrow$      if $\tilde{F} \uparrow$

Underfitting

High Bias
Low Variance

Overfitting

Low Bias
High Variance

$\hat{L}_{test}(\hat{f}_P)$

$\approx L(\hat{f}_P)$

$\hat{L}_{train}(\hat{f}_P)$

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

$D_1, D_2, D_3, D_4, \ldots \sim \mathbb{P}_D$

$\hat{f}_{D_1}$  $\hat{f}_{D_4}$

$f^* \approx \bar{f}$

$f_{Bayes}$

$\hat{f}_1$

$\hat{f}_3$

X

$\vec{X}$
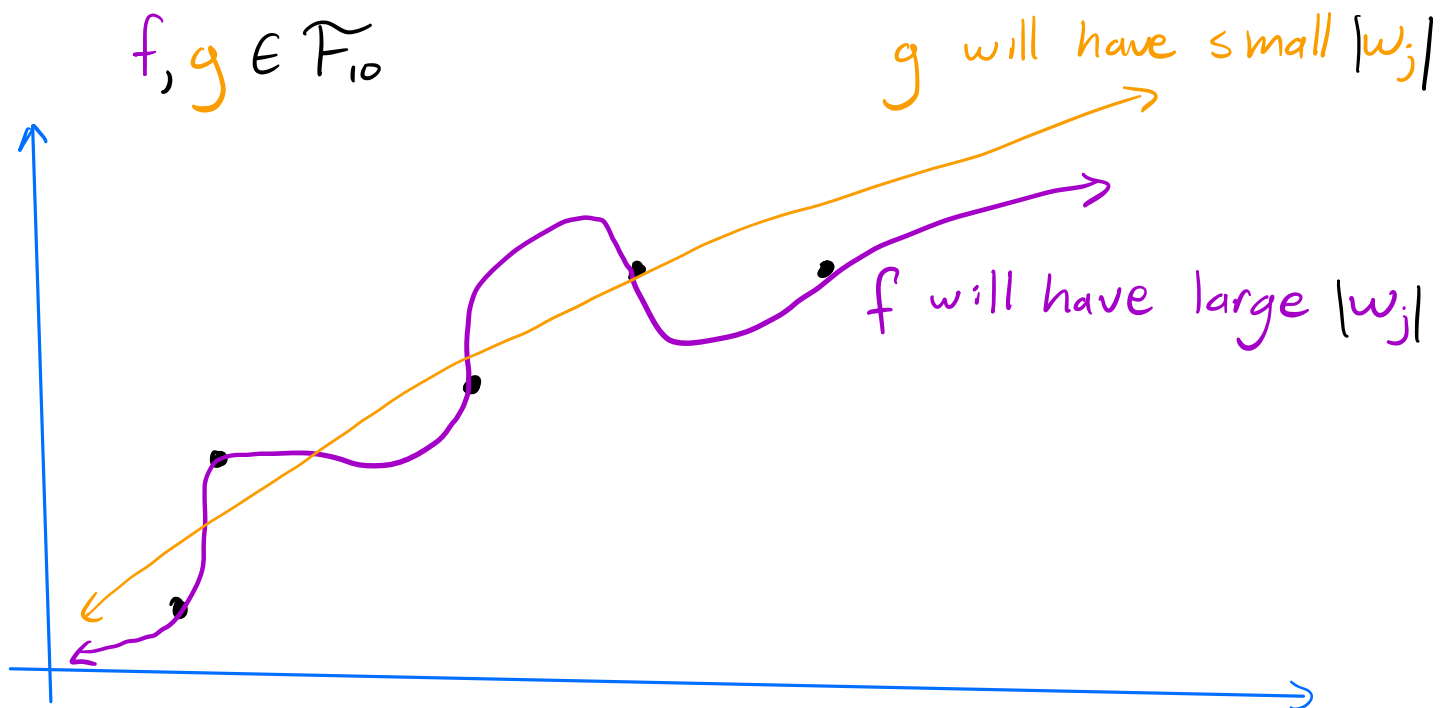
# Regularization

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

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

$f, g \in \widetilde{F}_{10}$

$g$ will have small $|w_j|$

$f$ will have large $|w_j|$

## Regularization: penalize large weights

If $f_p \in \widetilde{F}_{\bar{p}}$:

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

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

"Ridge Regression"

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

If $\lambda$ increases, then $\hat{f}_\lambda$ gets simpler

, $\bar{f}_\lambda$ gets simpler

, but $f^*$ does not change

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

## Bias vs. Variance

Bias: $\left(\bar{f}_\lambda(\vec{x}) - f_{Bayes}(\vec{x})\right)^2$

- Decreases if $\lambda$ decreases

Variance: $E\left[\left(\hat{f}_{D,\lambda}(\vec{x}) - \bar{f}_\lambda(\vec{x})\right)^2 \mid \vec{x}\right]$

- Increases if $\lambda$ decreases
- Decreases if $n$ increases

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

$\hat{\vec{w}}_\lambda = \underset{\vec{w} \in \mathbb{R}^{d+1}}{\arg\min} \hat{L}_\lambda (\vec{w})$    using squared loss, $\tilde{F}_1$

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

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

$$\hat{f}_\lambda = \underset{f \in \mathcal{F}_{10}}{\text{argmin}} \; \hat{L}_\lambda(f)$$

Underfitting

High Bias
Low Variance

Overfitting
Low Bias
High Variance

$\hat{L}_{test}(\hat{f}_\lambda)$
$\approx L(\hat{f}_\lambda)$

$\hat{L}_{train,\lambda}(\hat{f}_\lambda)$

$100 \quad 10 \quad 1 \quad 10^{1} \quad 10^{2} \quad 10^{3} \quad 10^{4} \qquad \lambda$