

# Important Announcements

Merch 4

- BGD:  $\nabla \hat{L}(\vec{w}) = (\frac{\partial}{\partial w_0} \hat{L}(\vec{w}), \dots$
- $O(n^3)$  for inverse calculation
- MBGD  $\nabla \hat{L}_m(\vec{w})$

# Evaluating Predictors/Models

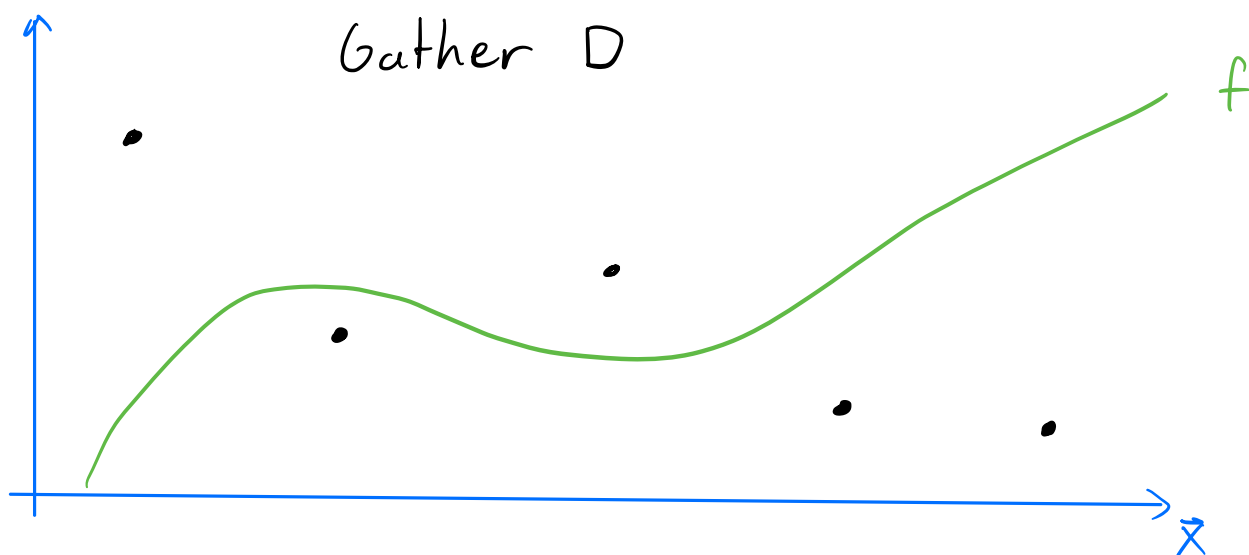
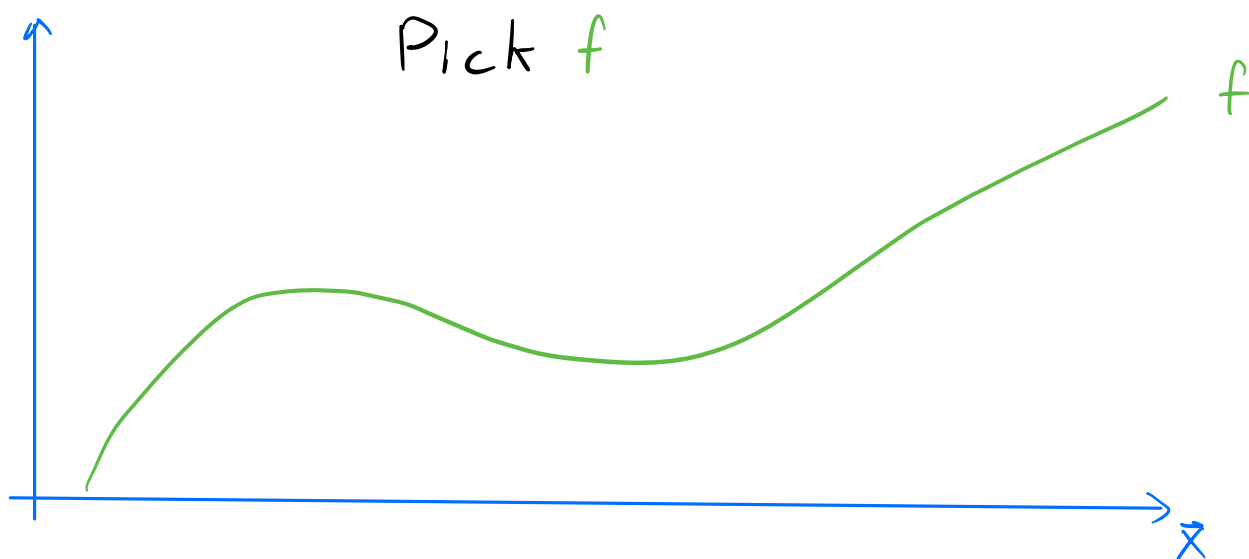
Is  $\hat{L}(\hat{f}_0)$  really a good estimate of  $L(\hat{f}_0)$ ?

where  $A(D) = f_0 \in \mathcal{F}$   $D$  is a r.v.  $(\vec{x}, y) \sim P_{\vec{x}, y}$

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

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

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



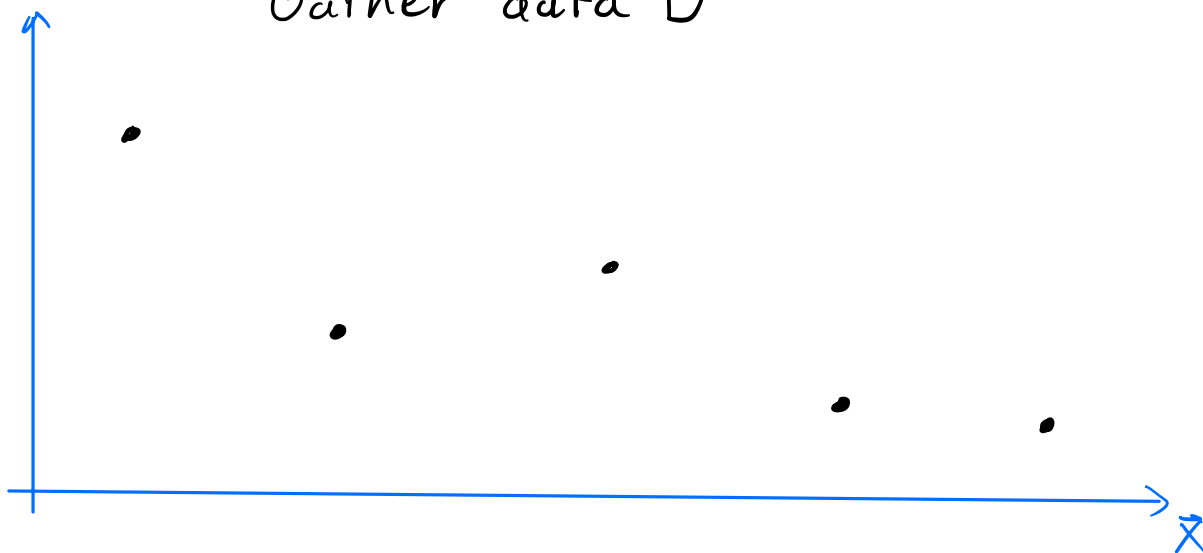
Then:  $\mathbb{E}[\hat{L}(f)] = L(f)$

$$\begin{aligned}\text{Var}[\hat{L}(f)] &= \text{Var}\left[\frac{1}{n} \sum_{i=1}^n \ell(f(\vec{X}_i), Y_i)\right] \\ &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}[\ell(f(\vec{X}_i), Y_i)] \\ &= \frac{1}{n} \text{Var}[\ell(f(\vec{X}_1), Y_1)]\end{aligned}$$

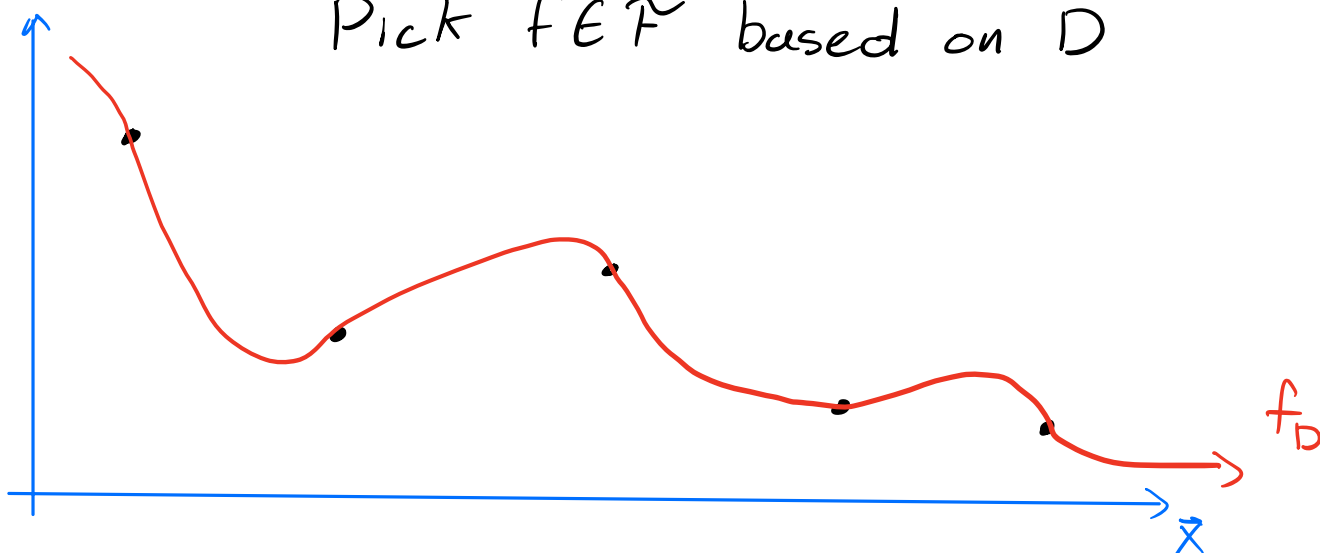
Since  $\ell(f(\vec{X}_i), Y_i)$  are independent for all  $i \in \{1, \dots, n\}$

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  $f \in \mathcal{F}$  based on  $D$



Then:

$$\mathbb{E}[\hat{L}(\hat{f}_0)] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n \ell(f(\vec{X}_i), Y_i)\right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\ell(f(\vec{X}_i), Y_i)] \\ \neq \mathbb{E}[\ell(f(\vec{X}_\cdot), Y_\cdot)]$$

$$\text{Var}[\hat{L}(\hat{f}_0)] = \text{Var}\left[\frac{1}{n} \sum_{i=1}^n \ell(\hat{f}_0(\vec{X}_i), Y_i)\right] \\ \neq \frac{1}{n^2} \sum_{i=1}^n \text{Var}[\ell(\hat{f}_0(\vec{X}_i), Y_i)]$$

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

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

Instead:

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

$$\text{Var}[\hat{L}(\hat{f}_0)] \leq \frac{1}{n} \max_{f \in \mathcal{F}} \text{Var}[\ell(f(\vec{X}_i), Y_i)] + \text{Var}[L(\hat{f}_0)]$$

As  $\mathcal{F}$  gets larger,  $\text{Var}[\hat{L}(\hat{f}_0)]$  increases  
 $\Rightarrow$  i.e.  $\hat{L}(\hat{f}_0)$  becomes a worse estimate of  $L(\hat{f}_0)$

As  $n$  gets larger,  $\text{Var}[\hat{L}(\hat{f}_0)]$  decreases  
 $\Rightarrow \hat{L}(\hat{f}_0)$  becomes a better estimate of  $L(\hat{f}_0)$