

Midterm Exam 2 Review

Linear Regression (Closed Form)

$$\mathcal{X} = \mathbb{R}^{d+1}, \mathcal{Y} = \mathbb{R} \text{ regression}$$

$$\mathcal{F} \subset \{f \mid f: \mathbb{R}^{d+1} \rightarrow \mathbb{R}\}$$

$$D = ((\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n))$$

$$\hat{L}(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(\vec{x}_i), y_i) \quad \text{an estimate of } L(f) \text{ for all } f \in \mathcal{F}$$

we want

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

pick $\mathcal{F} = \{f \mid f: \mathcal{X} \Rightarrow \mathcal{Y} \text{ and } f(\vec{x}) = \vec{x}^T \vec{w} \text{ where } \vec{w} \in \mathbb{R}^{d+1}\}$

Learner

$$A(D) = \hat{f} \in \mathcal{F}$$

Depends on D

$$\text{where } \hat{f}(\vec{x}) = \vec{x}^T \hat{\vec{w}} \quad \text{and} \quad \hat{\vec{w}} = A^{-1} \vec{b}$$

Algorithm: Closed form linear regression Learner

input: $D = ((\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n))$

$$A \leftarrow \sum_{i=1}^n \vec{x}_i \vec{x}_i^T$$

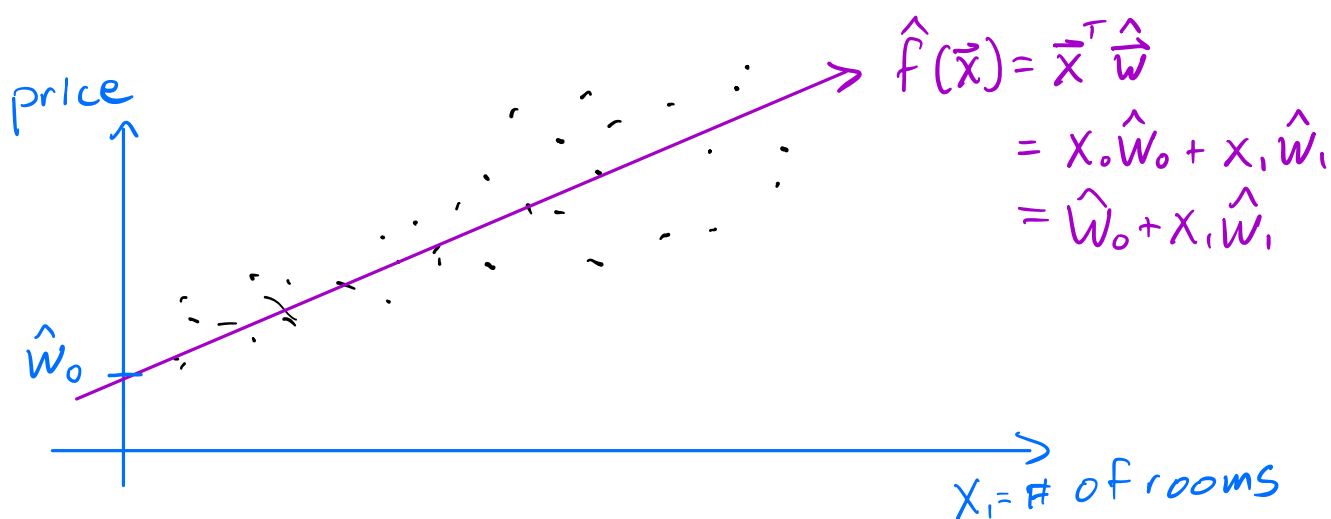
$$b \leftarrow \sum_{i=1}^n \vec{x}_i y_i$$

$$\hat{\vec{w}} \leftarrow A^{-1} b$$

return $\hat{f}(\vec{x}) = \vec{x}^T \hat{\vec{w}}$

Ex: $d=1$, $\mathcal{X} = \mathbb{R}^{1+1} = \mathbb{R}^2$, $\mathcal{Y} = \mathbb{R}$, $\hat{\vec{w}} = (\hat{w}_0, \hat{w}_1)^T$

$D = ((\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n))$



Gradient Descent

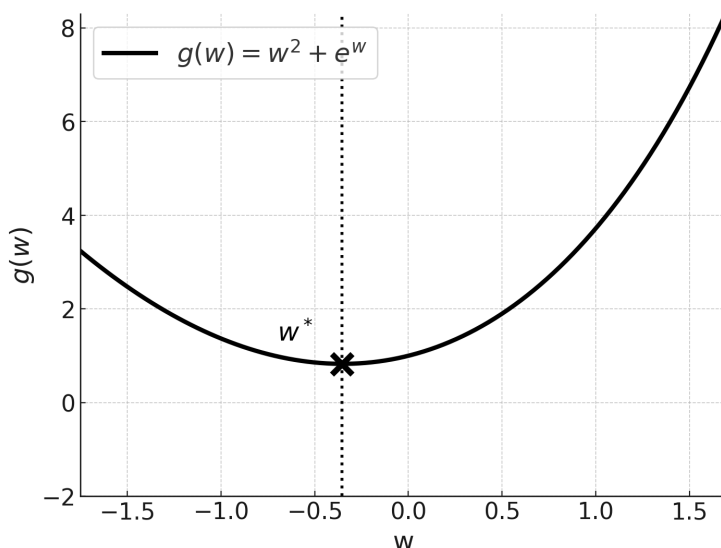
Ex: $g(w) = w^2 + e^w$, $g'(w) = 2w + e^w$, $g''(w) = 2 + e^w \geq 0$
Convex

$$w \in \mathbb{R} = \mathcal{W}$$

$$g'(w) = 2w + e^w = 0$$

$$2w = -e^w$$

No closed form
solution



Second-order Gradient Descent:

$$g(w) \approx g_{w^{(0)}}(w) = g(w^{(0)}) + g'(w^{(0)})(w - w^{(0)}) + \frac{g''(w^{(0)})}{2} (w - w^{(0)})^2$$

$$g'_{w^{(0)}}(w) = g'(w^{(0)}) + g''(w^{(0)})(w - w^{(0)}) = 0$$

$$\Rightarrow w = w^{(0)} - \frac{g'(w^{(0)})}{g''(w^{(0)})}$$

General: $w^{(t+1)} = w^{(t)} - \frac{g'(w^{(t)})}{g''(w^{(t)})}$

$w^{(t+1)}$ approaches $w^* = \arg \min_{w \in \mathcal{W}} g(w)$
as $t \rightarrow \infty$

First-order Gradient Descent

$$w^{(t+1)} = w^{(t)} - \eta^{(t)} g'(w^{(t)})$$

Multivariate Gradient Descent

Objective:

$$\vec{w}^* = (w_1^*, \dots, w_d^*)^T = \arg \min_{\vec{w} \in \mathcal{W}} g(\vec{w})$$

gradient descent update rule:

$$\vec{w}^{(t+1)} = \vec{w}^{(t)} - \eta^{(t)} \nabla g(\vec{w}^{(t)})$$

where $\nabla g(\vec{w}^{(t)}) = \left(\frac{\partial g}{\partial w_1}(\vec{w}^{(t)}), \dots, \frac{\partial g}{\partial w_d}(\vec{w}^{(t)}) \right)^T \in \mathbb{R}^d$

Selecting the step size

1. constant: $\eta^{(t)} = \eta \in (0, \infty)$

2. Inverse decaying: $\eta^{(t)} = \frac{\eta}{1+\lambda t}$ where $\eta, \lambda \in (0, \infty)$

3. Exponential decaying: $\eta^{(t)} = \eta \frac{1}{e^{\lambda t}}$ where $\eta, \lambda \in (0, \infty)$

4. Normalized gradient: $\eta^{(t)} = \frac{\eta}{\epsilon + \|\nabla g(\vec{w}^{(t)})\|}$ where $\eta, \epsilon \in (0, \infty)$
 ϵ is small
(ex: $\epsilon = 10^{-8}$)

$$\|\nabla g(\vec{w}^{(t)})\| = \sqrt{\sum_{j=1}^d \left(\frac{\partial g}{\partial w_j}(\vec{w}^{(t)}) \right)^2}$$

(Batch) Gradient Descent (BGD)

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

$$\nabla \hat{L}(\vec{w}^{(t)}) = \frac{2}{n} \sum_{i=1}^n (\vec{x}_i^T \vec{w} - y_i) \vec{x}_i$$

Algorithm: BGD Linear Regression Learner
(with a constant size)

input: $D = ((\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n))$, η , T number of
"epochs"

$\vec{w} \leftarrow$ random vector in \mathbb{R}^{d+1}

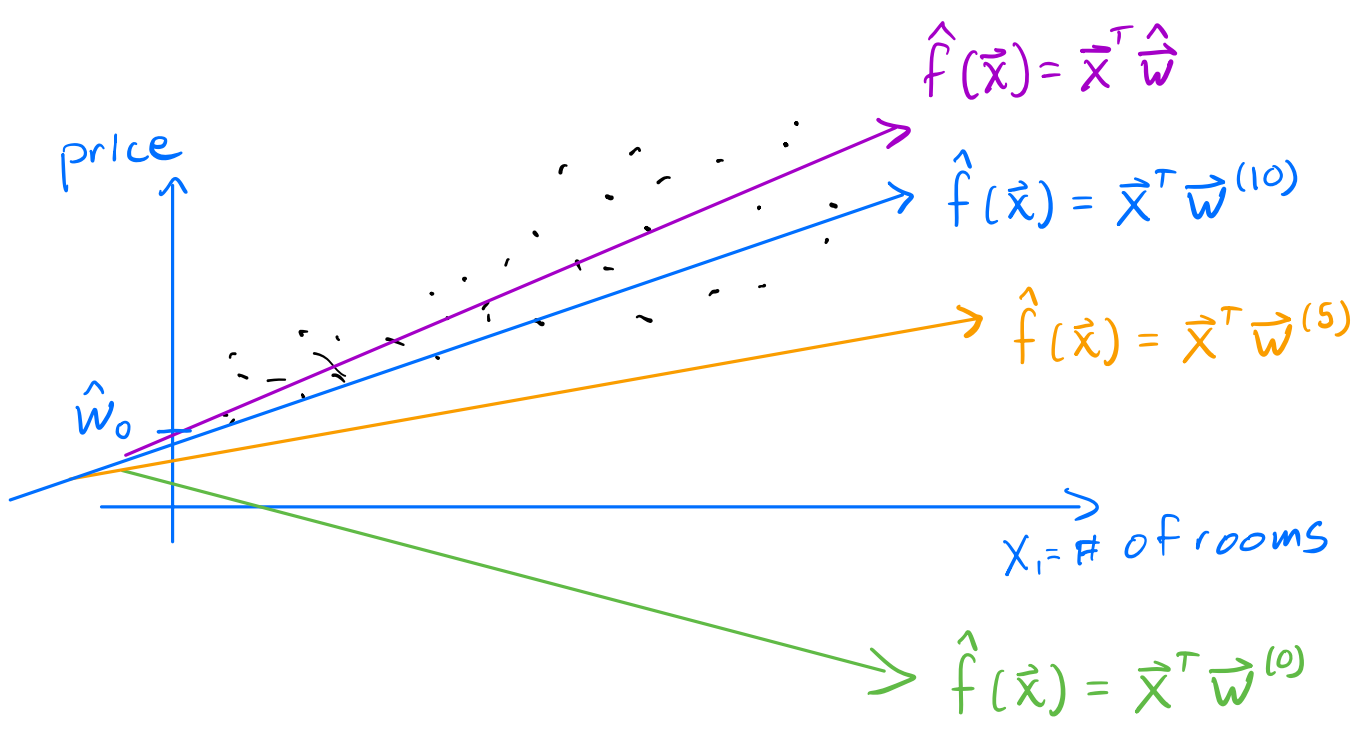
for $t = 1, \dots, T$

$$\nabla \hat{L}(\vec{w}) \leftarrow \frac{2}{n} \sum_{i=1}^n (\vec{x}_i^T \vec{w} - y_i) \vec{x}_i$$

$$\vec{w} = \vec{w} - \eta \nabla \hat{L}(\vec{w})$$

return $\hat{f}(\vec{x}) = \vec{x}^T \vec{w}$

\vec{x} :



Computation

Closed form: $O(nd^2 + d^3)$

BGD: $O(ndT)$

if $T < d$: $O(ndT) < O(nd^2) \leq O(nd^2 + d^3)$

BGD is more computationally efficient
if $T < d$

Mini-Batch Gradient Descent (MBGD)

$$D = \left(\begin{array}{l} (\vec{x}_1, y_1), \dots, (\vec{x}_b, y_b), \\ (\vec{x}_{b+1}, y_{b+1}), \dots, (\vec{x}_{2b}, \dots, y_{2b}), \\ \vdots \\ (\vec{x}_{(M-1)b+1}, y_{(M-1)b+1}), \dots, (\vec{x}_{Mb}, y_{Mb}) \end{array} \right)$$

b : mini-batch size

$M = \text{floor}\left(\frac{n}{b}\right)$: number of mini batches

We don't use datapoints at Mb

so we don't use $n - Mb \leq b$ datapoints

$$\hat{L}_m(\vec{w}) = \frac{1}{b} \sum_{i=(m-1)b+1}^{mb} (\vec{x}_i^T \vec{w} - y_i)^2 \quad m \in \{1, \dots, M\}$$

Algorithm: MBbD Linear Regression Learner
(with a constant size)

input: $D = ((\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n))$, η , T , b

$\vec{w} \leftarrow$ random vector in \mathbb{R}^{d+1}

$M \leftarrow \text{floor}(\frac{n}{b})$

for $t = 1, \dots, T$

Randomly shuffle D

for $m = 1, \dots, M$

$$\nabla \hat{L}(\vec{w}) \leftarrow \frac{2}{b} \sum_{i=(m-1)b+1}^{mb} (\vec{x}_i^T \vec{w} - y_i) \vec{x}_i$$

$$\vec{w} = \vec{w} - \eta \nabla \hat{L}(\vec{w})$$

return $\hat{f}(\vec{x}) = \vec{x}^T \vec{w}$

Advantage is MT is now the number
of gradient steps

Setting $b=n \Rightarrow M=1$ gives back BGD

$b=1 \Rightarrow M=n$ gives

"Stochastic GD (SGD)"

Polynomial Regression

$\phi_p: \mathcal{X} \rightarrow \mathcal{Z}$ is a degree p polynomial "feature map"

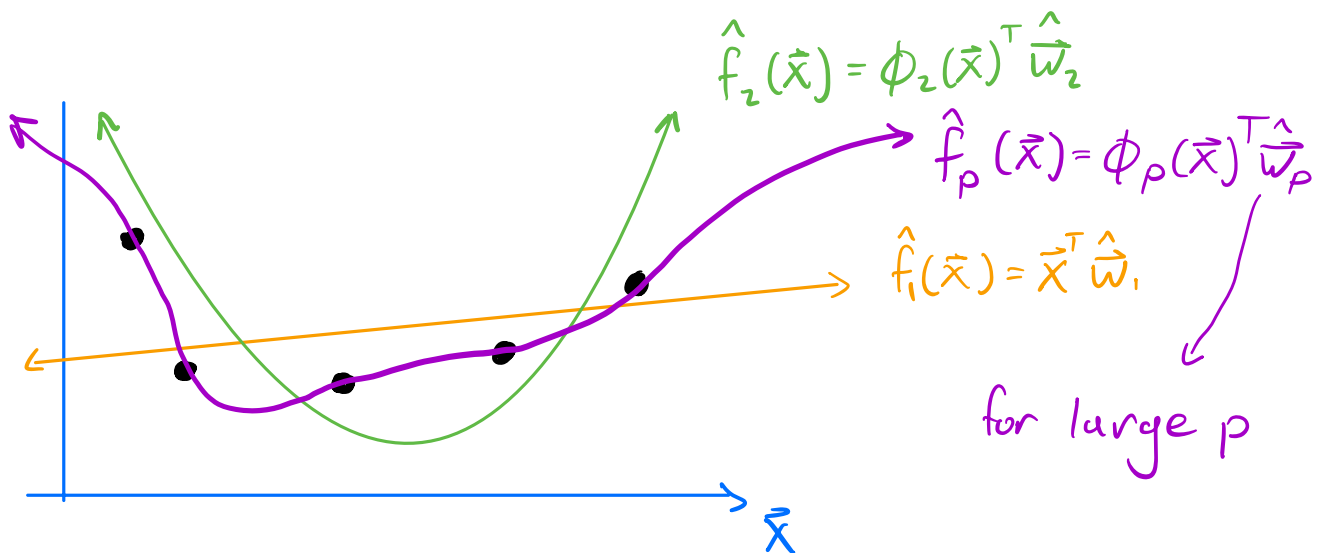
For $\mathcal{X} = \mathbb{R}^{d+1}$, $\mathcal{Z} = \mathbb{R}^{\bar{p}}$ where $\bar{p} = \binom{(d+1)+p-1}{p} = \binom{d+p}{p}$

$$\tilde{\mathcal{F}}_p = \left\{ f \mid f: \mathbb{R}^{d+1} \rightarrow \mathbb{R} \text{ and } f(\bar{x}) = \phi_p(\bar{x})^T \bar{w}, \bar{w} \in \mathbb{R}^{\bar{p}} \right\}$$

$$\tilde{\mathcal{F}}_1 \subset \tilde{\mathcal{F}}_2 \subset \dots \subset \tilde{\mathcal{F}}_p$$

$$\hat{L}(\hat{f}_1) \geq \hat{L}(\hat{f}_2) \geq \dots \geq \hat{L}(\hat{f}_p) \approx 0$$

Ex:



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 $L(\hat{f})$

If $A(D) = f_0$ depends on D then

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

$l(\hat{f}_0(\vec{X}_i), Y_i)$ 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 $E[L(A(D))]$

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

\mathcal{F} any function class

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

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

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

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

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

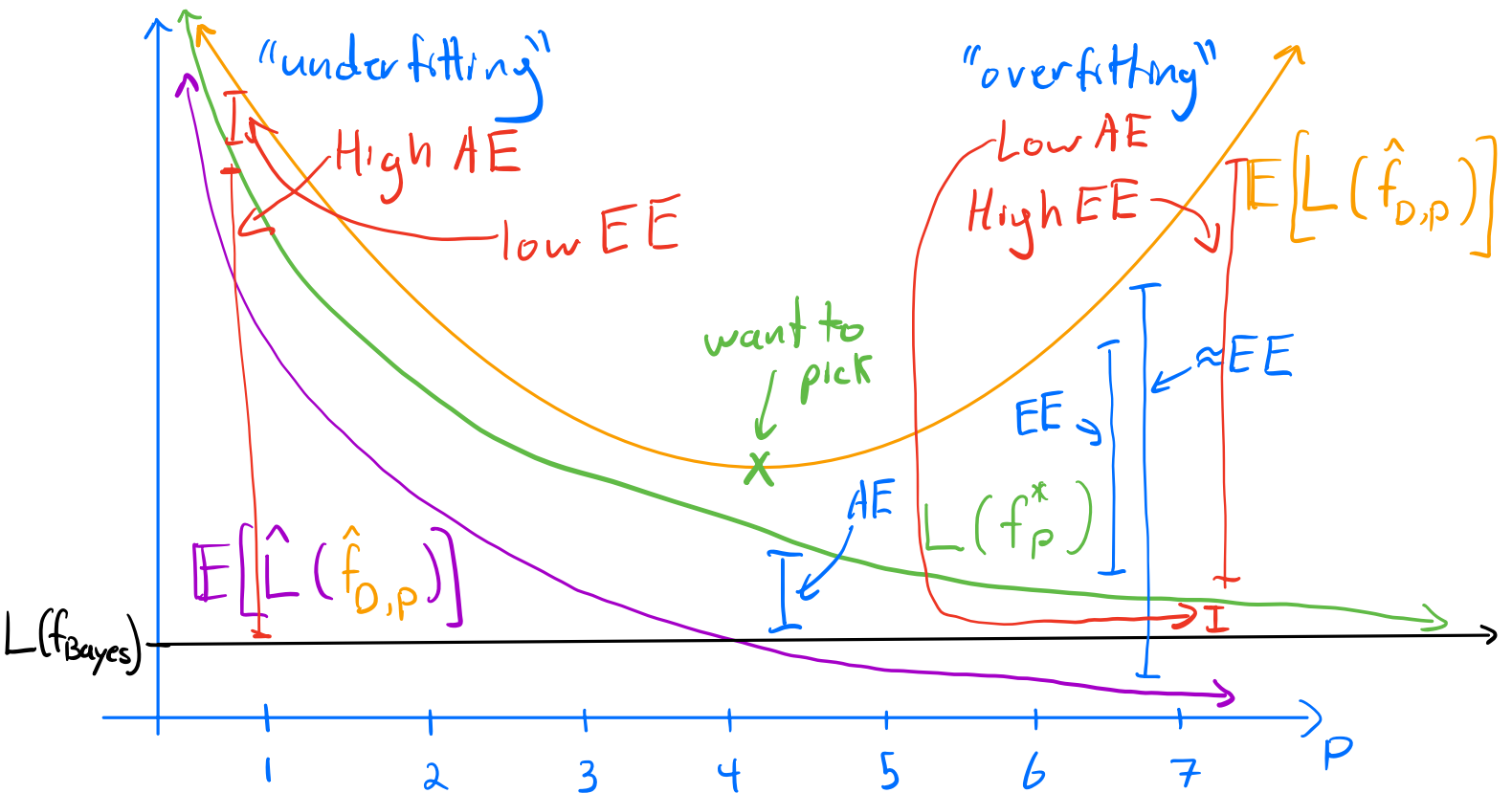
Estimation Error: Due to random dataset D

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

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

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

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



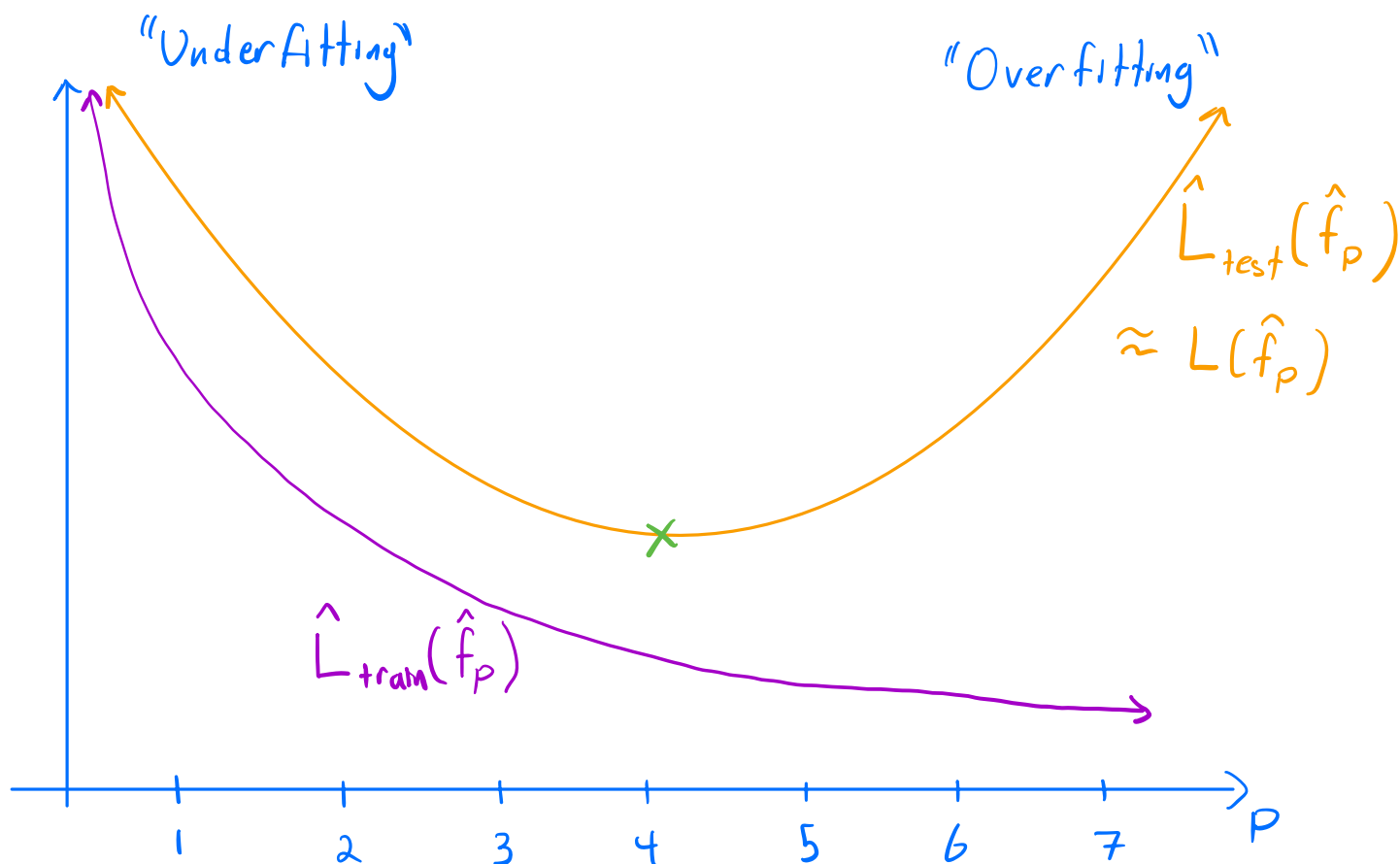
Can't calculate L , so use train, test split

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

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

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

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



Bias-Variance Tradeoff

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

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

$$\begin{aligned} &= \overbrace{\mathbb{E}[\mathbb{E}[(\hat{f}_D(\vec{x}) - \bar{f}(\vec{x}))^2 | \vec{x}]]]}^{= EE} + \overbrace{\mathbb{E}[(\bar{f}(\vec{x}) - f_{\text{Bayes}}(\vec{x}))^2]}^{= AE} + \underbrace{L(f_{\text{Bayes}})}_{IE} \\ &\text{Variance} = \text{Var}[\hat{f}_D(\vec{x}) | \vec{x}] \qquad \text{Bias} \qquad \text{IE} \end{aligned}$$

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

Effects of changing \mathcal{F}, n on Bias, Variance

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

Variance \downarrow if $n \uparrow$

\uparrow if $\mathcal{F} \uparrow$

Regularization

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

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

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

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

Regularization parameter

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

Regularized estimated loss

$j=0$ not included "Regularizer"

Bias vs. Variance

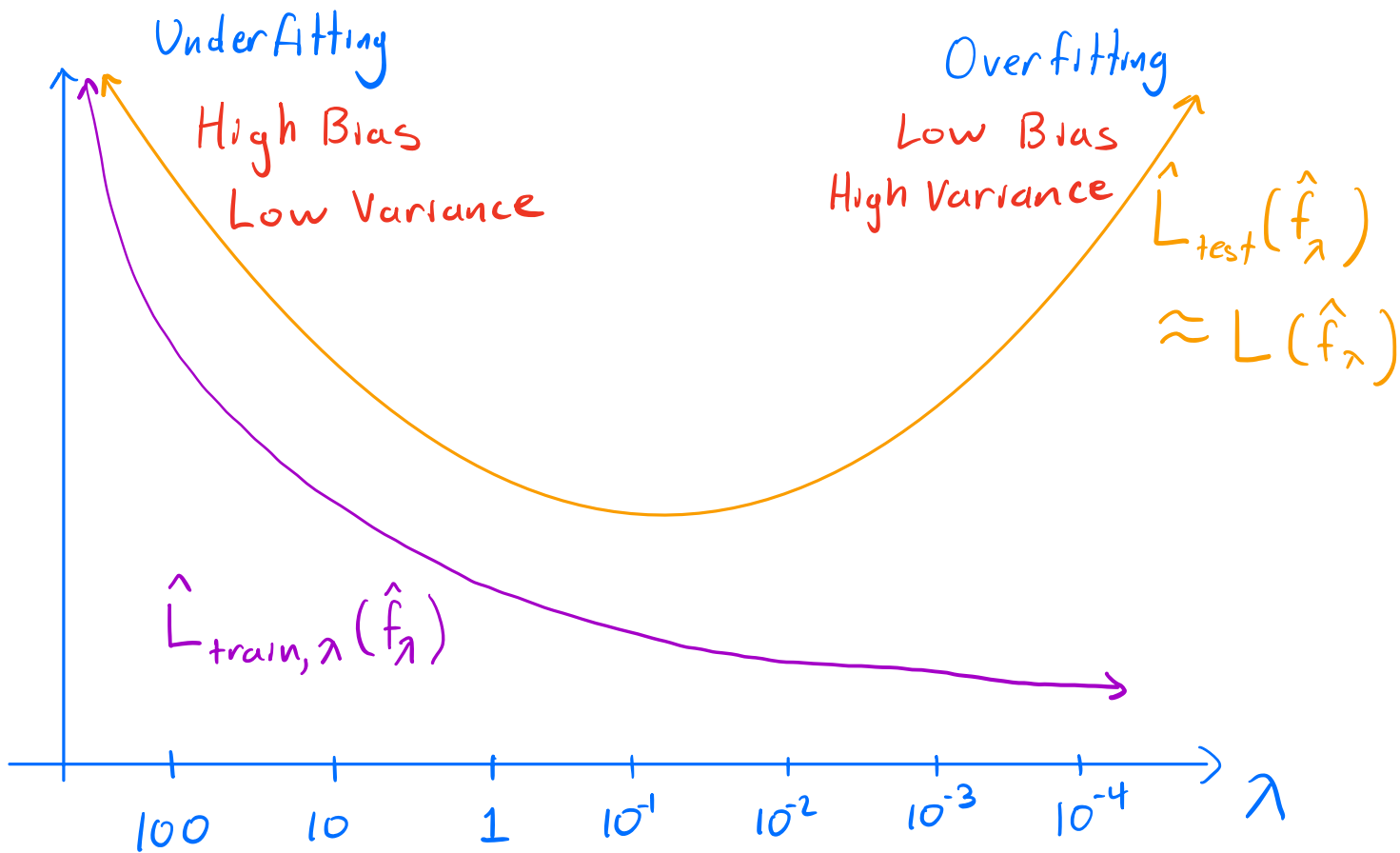
$$\text{Bias: } (\bar{f}_\lambda(\vec{x}) - f_{\text{Bayes}}(\vec{x}))^2$$

- Decreases if λ decreases

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

- Increases if λ decreases
- Decreases if n increases

$$\hat{f}_\lambda = \operatorname{argmin}_{f \in \mathcal{F}_0} \hat{L}_\lambda(f)$$



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

Objective:

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

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

BGD:

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

MLE

$$f_{\text{Bayes}} = \underset{f \in \{f \mid f: \mathcal{X} \rightarrow \mathcal{Y}\}}{\text{argmin}} L(f)$$

assume squared loss and Regression

$$\begin{aligned} f_{\text{Bayes}}(\bar{x}) &= \mathbb{E}[Y \mid \bar{X} = \bar{x}] \\ &= \int_{\mathcal{Y}} y P_{Y|\bar{X}}(y|\bar{x}) dy \end{aligned}$$

Use D to estimate $P_{Y|\bar{X}}$

MLE Basics

$D = (z_1, \dots, z_n)$ and z_i are i.i.d. with p_z

Assume p_z is based on some parameter w^*

You are give a fixed data $D = (z_1, \dots, z_n)$

$$\begin{aligned} w_{\text{MLE}} &= \underset{w \in \mathcal{W}}{\text{argmax}} \underbrace{p(D|w)}_{\text{likelihood}} \\ &= \underset{w \in \mathcal{W}}{\text{argmin}} \underbrace{-\log(p(D|w))}_{\text{negative log-likelihood}} \end{aligned}$$

$$= \arg \min_{w \in \mathcal{W}} -\log \left(\prod_{i=1}^n p(z_i | w) \right)$$

$$= \arg \min_{w \in \mathcal{W}} -\sum_{i=1}^n \log(p(z_i | w))$$

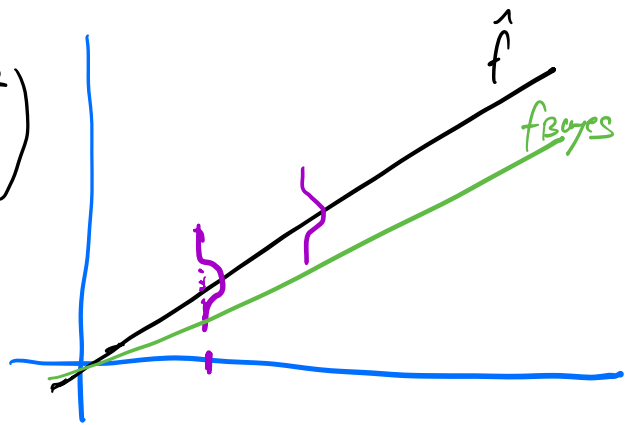
Estimating $P_{Y|\bar{X}}$

$$D = ((\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)) \in (\mathcal{X} \times \mathcal{Y})^n, P_D, p_D$$

(\vec{x}_i, y_i) are i.i.d with $P_{\vec{x}, y}$ and $p_{\vec{x}, y}$

Assume $Y_i | \vec{X}_i = \vec{x}_i \sim \mathcal{N}(\vec{x}_i^T \vec{w}^*, 1)$

$$P_{Y|\vec{X}=\vec{x}}(y|\vec{x}) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y - \vec{x}^T \vec{w}^*)^2}{2}\right)$$



$$\vec{w}_{MLE} = \arg \min_{\vec{w} \in \mathbb{R}^{d+1}} \sum_{i=1}^n \frac{(y_i - \vec{x}_i^T \vec{w})^2}{2}$$

$$= \arg \min_{\vec{w} \in \mathbb{R}^{d+1}} \hat{L}(\vec{w})$$

$$f_{Bayes}(\vec{x}) \approx \vec{x}^T \vec{w}_{MLE}$$

MAP

If P_{riX} is a function of some parameter $w \in W$ then we just need to estimate that parameter

$$\text{MLE: } \underset{w \in W}{\text{argmax}} \underbrace{p(D|w)} = \prod_{i=1}^n p(z_i|w)$$

"find w that maximizes the likelihood of the data"

$$\text{MAP: } \underset{w \in W}{\text{argmax}} \underbrace{p(w|D)} = \text{"posterior"}$$

"find w that is the most likely given the data"

$$w_{\text{MAP}} = \underset{w \in W}{\text{argmax}} p(w|D)$$

$$= \underset{w \in W}{\text{argmax}} \frac{p(w, D)}{P(D)}$$

prod rule \Downarrow

$$= \underset{w \in W}{\text{argmax}} \frac{p(D|w) p(w)}{P(D)}$$

$$= \underset{w \in W}{\text{argmax}} \underbrace{p(D|w)}_{\text{likelihood}} \underbrace{p(w)}_{\text{prior}}$$

$$= \arg \min_{w \in \mathcal{W}} -\log(p(D|w) p(w))$$

$$= \arg \min_{w \in \mathcal{W}} \left[-\log(p(D|w)) - \log(p(w)) \right]$$

if n large then:

$$w_{\text{MAP}} \approx w_{\text{MLE}}$$

if $p(w) = c$ a constant then:

$$w_{\text{MAP}} \approx w_{\text{MLE}}$$

if $\text{Var}[W]$ is large then:

$$w_{\text{MAP}} \approx w_{\text{MLE}}$$

small then:

$$w_{\text{MAP}} \approx \mathbb{E}[W]$$

Estimating \vec{w} for $P_{Y|\vec{X}}$

$$D = ((\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)) \in (\mathcal{X} \times \mathcal{Y})^n, P_D, p_D$$

(\vec{x}_i, y_i) are i.i.d with $P_{\vec{x}, y}$ and $p_{\vec{x}, y}$

$$\text{Assume } Y_i | \vec{X}_i = \vec{x}_i \sim \mathcal{N}(\vec{x}_i^T \vec{w}^*, 1)$$

$$P_{Y|\vec{X}=\vec{x}}(y|\vec{x}) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y - \vec{x}^T \vec{w}^*)^2}{2}\right)$$

Assume $w_j \sim \mathcal{N}(0, \frac{1}{\lambda})$ are i.i.d. for $j \in \{1, \dots, d\}$

and $w_0 \sim \mathcal{N}(0, \sigma^2)$ for very large σ

$\approx \text{Uniform}(-a, a)$ for large a

w_0 is independent of w_j for all $j \in \{1, \dots, d\}$

$$\vec{w}_{\text{MAP}} = \underset{\vec{w} \in \mathbb{R}^{d+1}}{\text{argmax}} p(\vec{w} | D)$$

$$= \underset{\vec{w} \in \mathbb{R}^{d+1}}{\text{argmin}} \left[\underbrace{\sum_{i=1}^n \frac{(y_i - \vec{x}_i^T \vec{w})^2}{2}}_{= \frac{n}{2} \hat{L}} + \underbrace{\frac{\lambda}{2} \sum_{j=1}^d w_j^2}_{\text{almost regularizer}} \right]$$

$$= \underset{\vec{w} \in \mathbb{R}^{d+1}}{\text{argmin}} \left[\hat{L}_{\lambda}(\vec{w}) \right]$$

$$f_{\text{Bayes}}(\vec{x}) \approx \vec{x}^T \vec{w}_{\text{MAP}}$$

Lasso regression

Assume $w_j \sim \text{Laplace}(0, 1/\lambda)$ are i.i.d for $j \in \{1, \dots, d\}$

and $w_0 \sim \text{Laplace}(0, b)$ for very large b

$\approx \text{Uniform}(-a, a)$ for large a

w_0 is independent of w_j for all $j \in \{1, \dots, d\}$

$$\vec{w}_{\text{MAP}} = \arg \max_{\vec{w} \in \mathbb{R}^{d+1}} p(\vec{w} | D)$$

$$= \arg \min_{\vec{w} \in \mathbb{R}^{d+1}} \left[\sum_{i=1}^n \frac{(y_i - \vec{x}_i^T \vec{w})^2}{2} + \lambda \sum_{j=1}^d |w_j| \right]$$