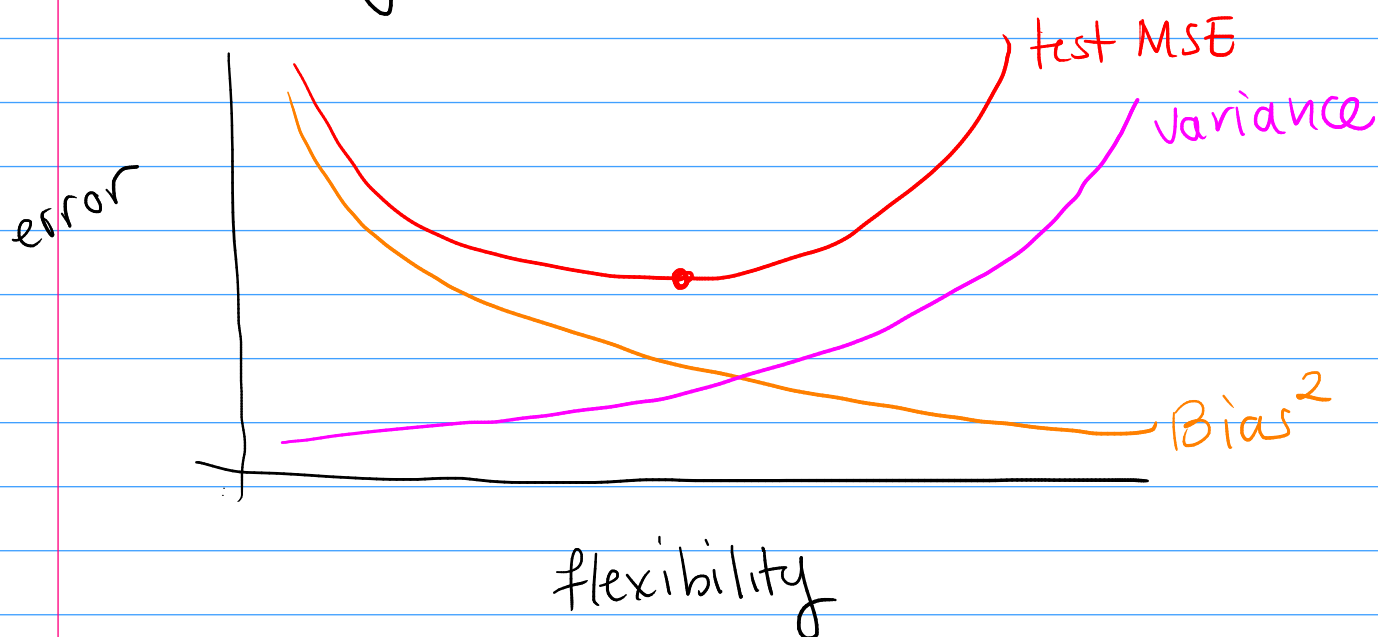


Lecture 8: Risk Minimization

Mathematically: $MSE = Bias^2 + Var + Irreducible$



Mathematically:

$$Bias(\hat{f}(x_0)) = E[\hat{f}(x_0)] - f(x_0)$$

$$Var(\hat{f}(x_0)) = E[(\hat{f}(x_0) - E[\hat{f}(x_0)])^2]$$

$$Err(x_0) = E[(Y_0 - \hat{f}(x_0))^2] \quad \leftarrow MSE$$

$$= E[(Y_0 - E[\hat{f}(x_0)])^2] \quad (1)$$

$$+ E[(\hat{f}(x_0) - E[\hat{f}(x_0)])^2] \quad (2)$$

$$+ 2E[(Y_0 - E[\hat{f}(x_0)])(\hat{f}(x_0) - E[\hat{f}(x_0)])] \quad (3)$$

$$Y_0 = f(x_0) + \varepsilon_0$$

$$\begin{aligned}
 (1) &= \mathbb{E}[(\underbrace{f(x_0)}_a + \underbrace{\varepsilon_0}_b - \mathbb{E}[\hat{f}(x_0)])^2] = \mathbb{E}[(a+b)^2] \\
 &= \mathbb{E}[(\underbrace{f(x_0) - \mathbb{E}[\hat{f}(x_0)]}_{\text{deterministic}})^2] = \text{Bias}^2 \\
 &\quad + \mathbb{E}[\varepsilon_0^2] = \sigma^2 \\
 &\quad + 2\mathbb{E}[\varepsilon_0(\underbrace{f(x_0) - \mathbb{E}[\hat{f}(x_0)]}_{\text{deterministic}})] = 0 \\
 &\quad \mathbb{E}[\varepsilon_0] = 0 \\
 &\rightarrow \text{Bias}^2 + \sigma^2
 \end{aligned}$$

$$(2) = \text{Var}(\hat{f}(x_0)) \quad , \quad \text{Var}(z) = \mathbb{E}[(z - \mathbb{E}z)^2] \\
 z = \hat{f}(x_0)$$

$$\begin{aligned}
 (3) &\mathbb{E}[(\underbrace{Y_0 - \mathbb{E}[\hat{f}(x_0)]}_{\text{independent b/c } Y_0 \perp \text{ training data}})(\underbrace{\hat{f}(x_0) - \mathbb{E}[\hat{f}(x_0)]}_{\text{independent b/c } Y_0 \perp \text{ training data}})] \\
 &\downarrow \\
 &\mathbb{E}[Y_0 - \mathbb{E}[\hat{f}(x_0)]] \mathbb{E}[\underbrace{\hat{f}(x_0) - \mathbb{E}[\hat{f}(x_0)]}_0] = 0 \\
 &\quad \mathbb{E}[z - \mathbb{E}[z]] \quad z = \hat{f}(x_0) \\
 &= \mathbb{E}[z] - \mathbb{E}[z] = 0
 \end{aligned}$$

$A \perp B \Rightarrow \mathbb{E}[AB] = \mathbb{E}[A]\mathbb{E}[B]$

Q: what's the best \hat{f} theoretically?

$L = \text{loss}$ (prev. called M)

$$L(Y, f(x)) = (Y - f(x))^2 = \text{sq. loss}$$

$$L(Y, f(x)) = |Y - f(x)| = \text{abs. loss}$$

$$f^* = \underset{f}{\operatorname{argmin}} \underbrace{\mathbb{E}[L(Y, f(x))]}_{\text{Risk}}$$

Turns out can actually get an answer.

Let $(\underline{x}, Y) \sim p$ joint dist

then

$$\mathbb{E}[L(Y, f(\underline{x}))]$$

$$= \mathbb{E}_{\underline{x}} \underbrace{\mathbb{E}[L(Y, f(\underline{x})) | \underline{x} = \underline{x}]}_{A(\underline{x})}$$

want to minimize

Iterated Expectation

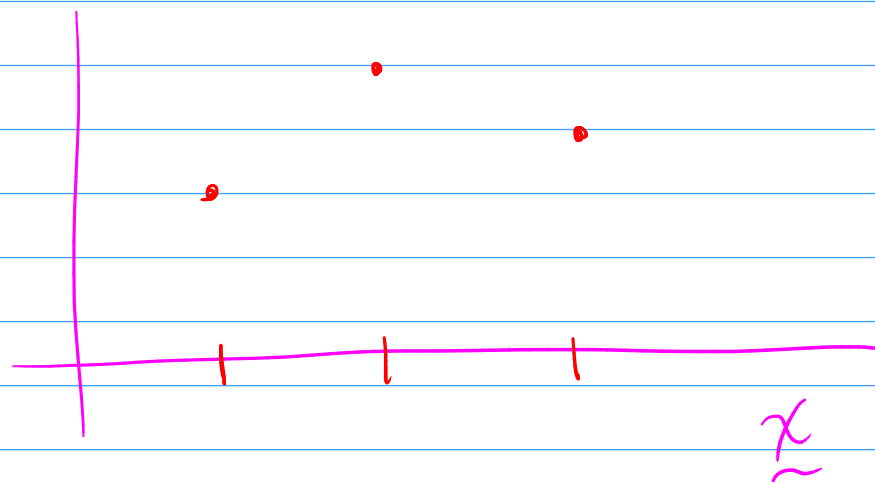
$$\mathbb{E}[A] = \mathbb{E}_B \mathbb{E}[A|B]$$

$$= \int A(\underline{x}) p(\underline{x}) d\underline{x}$$

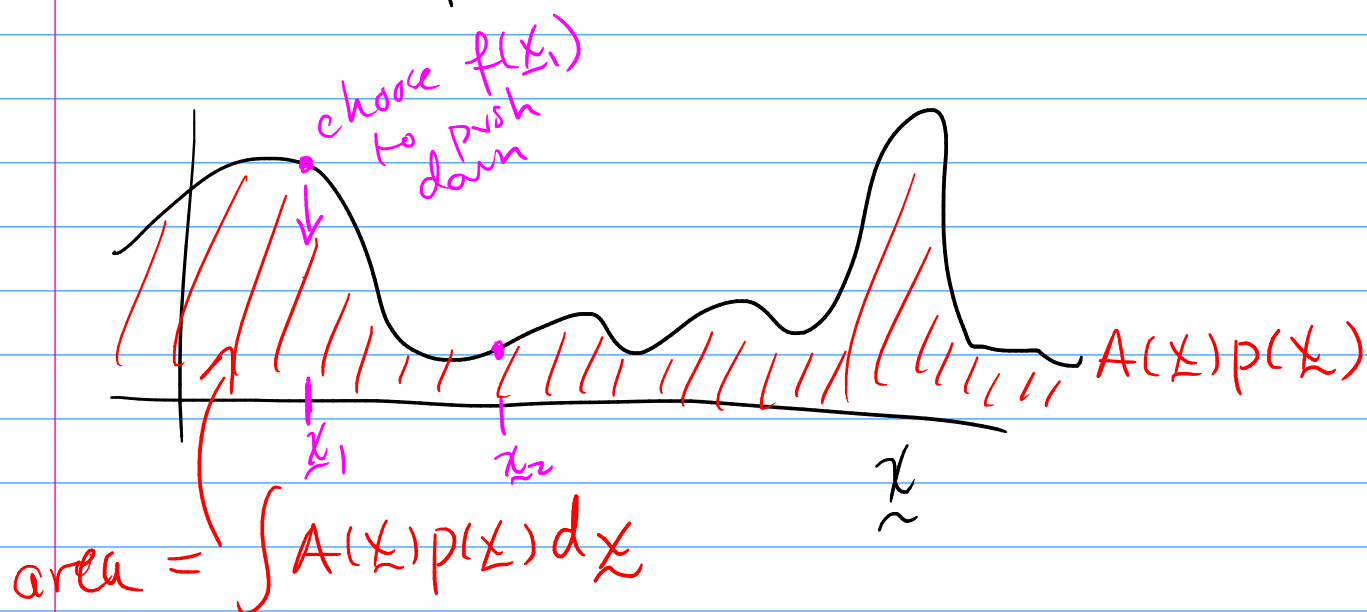
depends on f

need to choose f to minimize

I need to choose f , can do it separately for each \underline{x} .



Look at $A(\underline{x})p(\underline{x})$



because $p(x)$ doesn't depend on $f(x)$

only need to try to minimize $A(x)$

So, really just need to choose $f(x)$ at each x to minimize $A(x)$

i.e.

$$f^*(\underline{x}) = \underset{f(\underline{x})}{\operatorname{argmin}} \mathbb{E}[L(Y, f(\underline{x})) | \underline{X} = \underline{x}]$$

$$= \underset{c}{\operatorname{argmin}} \mathbb{E}[L(Y, c) | \underline{X} = \underline{x}]$$

more concrete, $L = \text{sq loss}$

$$f^*(\underline{x}) = \underset{c}{\operatorname{argmin}} \mathbb{E}[(Y - c)^2 | \underline{X} = \underline{x}]$$

Aside:

$$\underset{c}{\operatorname{argmin}} \mathbb{E}[(z - c)^2] = \mathbb{E}[z]$$

pf

$$\begin{aligned} \mathbb{E}[(z - c)^2] &= \mathbb{E}[z^2 + c^2 - 2zc] \\ &= \mathbb{E}[z^2] + c^2 - 2c \mathbb{E}[z] \end{aligned}$$

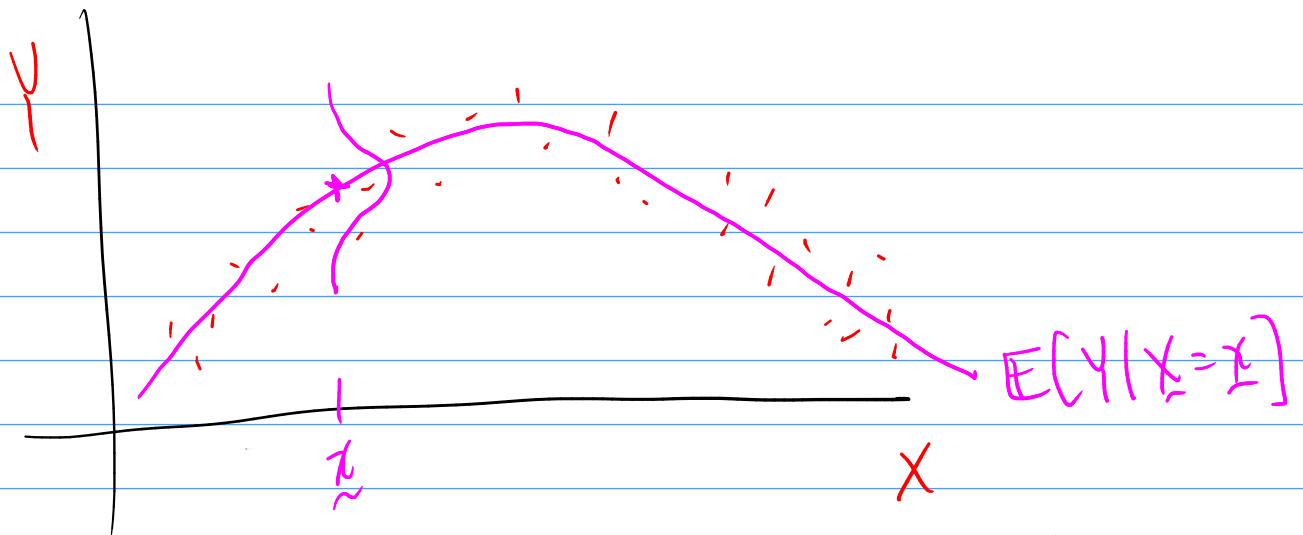
$$\frac{\partial}{\partial c}(\dots) = 2c - 2\mathbb{E}[z] = 0$$

Solve for c

$$\Rightarrow c = \mathbb{E}[z].$$

So

$$f^*(\underline{x}) = \mathbb{E}[Y | \underline{X} = \underline{x}]$$



What about $L(Y, f(x)) = |Y - f(x)|$?

$$f^*(x) = \text{Median}(Y | X = x)$$

Problem: I don't know underlying dists.

Can try to approx

$$\hat{f}(\underline{x}) \approx E[Y | \underline{X} = \underline{x}]$$

using my training data.

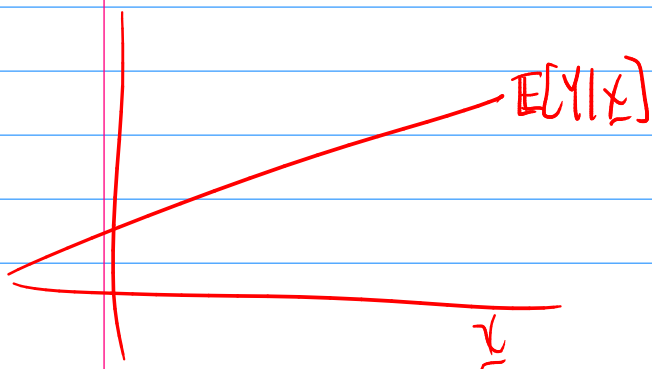
Ex. $\hat{f}(\underline{x}) = \text{avg. of } y_n\text{'s for } \underline{x}_n\text{'s near } \underline{x}$

KNN regression.



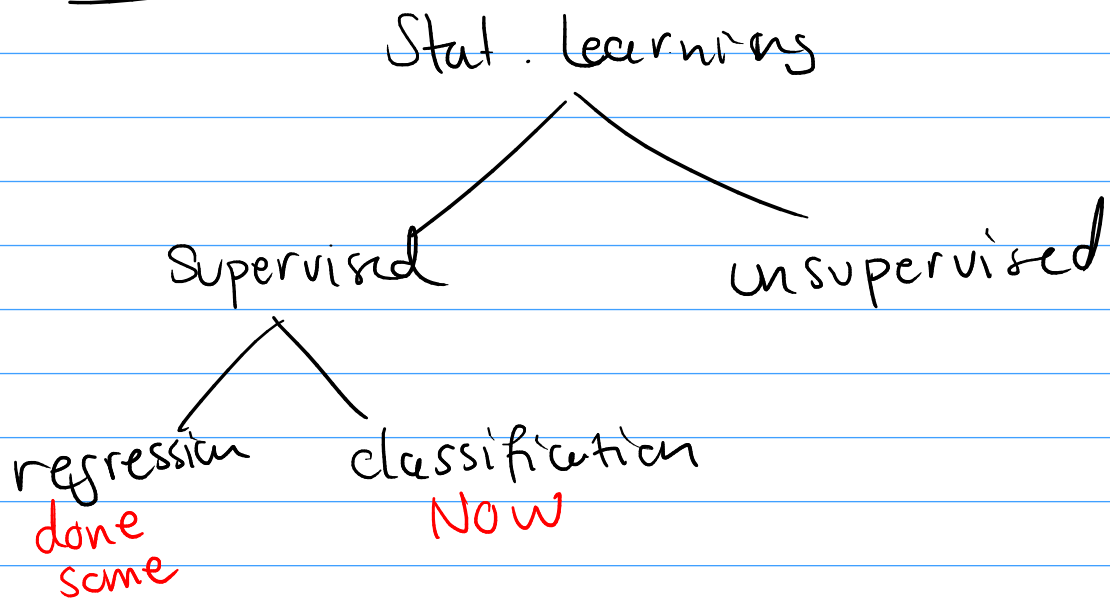
Ex Make some assumptions about structure of $E[Y | \underline{X} = \underline{x}]$

in particular assume $E[Y | \underline{X} = \underline{x}] = \underline{x}^T \beta$



OLS linear regression.

Outline:



Setup: Classification

$$\underline{x} \in \mathbb{R}^p, y \in \mathcal{C}$$

$$\mathcal{C} = \{c_1, c_2, \dots, c_K\}$$

↑ set of possible classes

Goal: find some \hat{f} so that $y \approx \hat{f}(\underline{x})$.

Loss function for classification: 0-1 loss

$$\begin{aligned} L(y, \hat{f}(\underline{x})) &= \mathbb{1}(y \neq \hat{f}(\underline{x})) \\ &= \begin{cases} 0 & \hat{f}(\underline{x}) = y \\ 1 & \hat{f}(\underline{x}) \neq y \end{cases} \end{aligned}$$

What is f^* ?

$$f^*(\underline{x}) = \operatorname{argmin}_c \mathbb{E}[L(Y, c) | \underline{X} = \underline{x}]$$
$$= \operatorname{argmin}_c \mathbb{E}[\mathbb{1}(Y \neq c) | \underline{X} = \underline{x}]$$

Claim! $\left(\mathbb{E}[\mathbb{1}(\overset{\text{statement}}{\text{---}})] = P(\overset{\text{statement}}{\text{---}}) \right)$

$$= \operatorname{argmin}_c P(Y \neq c | \underline{X} = \underline{x})$$
$$= \operatorname{argmin}_c 1 - P(Y = c | \underline{X} = \underline{x})$$

$$f^*(\underline{x}) = \operatorname{argmax}_c P(Y = c | \underline{X} = \underline{x})$$

Bayes classifier

Ex. 3-class problem: Cats, Dogs, People

$$P(Y = \text{cat} | \underline{X}), P(Y = \text{dog} | \underline{X}), P(Y = \text{people} | \underline{X})$$

predict class w/ largest prob.