# Bias - Variance Tradeoff

Perf. of method has two comps.:

$\underline{Bias}$ : err from approx. complicated real world phenom w/ a simple fn.

$\underline{Variance}$ : err b/c our $\hat{f}$ is sens. to $D_{train}$

Often, are opposition.

Can formulate mathy:

$$Y = f(\underset{\sim}{X}) + \varepsilon, \quad \varepsilon \text{ random}$$
$$E[\varepsilon] = 0$$
$$Var(\varepsilon) = \sigma^2$$
$$\underset{\sim}{X} \perp \varepsilon$$

gen. according to this model

Let $D_{train} = \{(X_n, y_n)\}_{n=1}^{N}$ be used to

Let $D_{train}$ $\{(x_n, y_n)\}_{n=1}$ be ????

fit $\hat{f} = \hat{f}_{D_{train}}$.

let $\underset{\sim}{X}_{new}, Y_{new}$ be indep draws from this model.

Defn
$$MSE = E\left[\left(Y_{new} - \hat{f}(X_{new})\right)^2\right]$$

Assume $X_{new}$ is fixed (not random)

Then if we define

$$Bias(\hat{f}) = E\left[\hat{f}(X_{new})\right] - f(X_{new})$$

$$Var(\hat{f}) = Var\left(\hat{f}_{D_{train}}(X_{new})\right)$$

Can show:

reducible (?)

$$MSE = Bias(\hat{f})^2 + Var(\hat{f}) + \sigma^2$$

irreducible

Called BV decomp.

Typ. bias and var are opposed.

low flex $\iff$ high bias/low variance
high flex $\iff$ low bias/high var.



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

$L$ = loss fn

e.g. $L(y, f(x)) = (y - f(x))^2$    (Sq. loss)

e.g. $L(y, f(x)) = |y - f(x)|$    (abs. loss)

$$f^* = \arg\min_{f} E\left[L(Y, f(x))\right]$$

$(X, Y) \sim P_{\wedge}$ joint dist of $X, Y$

**Can get answer;**

Consider $E\left[L(Y, f(x))\right]$

$$= E_X \underbrace{E\left[L(Y, f(x))|X\right]}_{A(x)}$$

$$= \int A(x) p(x) dx \qquad \text{dens of } \pi$$

Total Exp.

$E[A] = E_B E[A|B]$

$\llcorner$ depends on $f(x)$

Can choose $f(x)$ for each $x$ separately.

for each $x$ choose $f(x)$
to push down as much
as possible

$$\int A(x) p(x) dx$$

$A(x) p(x)$

$x$

Can choose $f$ to do this sep. for each $x$
Since $p(x)$ doesn't change based on $f(x)$
all I need to look at is $A(x)$.

To find $f^*$ just need to optimize

$$E[L(Y, f(x)) | X]$$

sep for each $X$.

$$f^*(x) = \arg\min_{f(x)} E[L(Y, f(x)) | X]$$

---

Consider Sq. loss $L(Y, f(x)) = (Y - f(x))^2$

$$f^*(x) = \arg\min_{c \in \mathbb{R}} E[(Y - c)^2 | X]$$

---

$\underline{Ex.}$  $\arg\min_{c} E[(Z - c)^2] = E[Z]$

$\underline{pf.}$  $E[(Z-c)^2] = E[Z^2 - 2Zc + c^2]$

$$= E[Z^2] - 2c E[Z] + c^2$$

take deriv wrt $c$ :
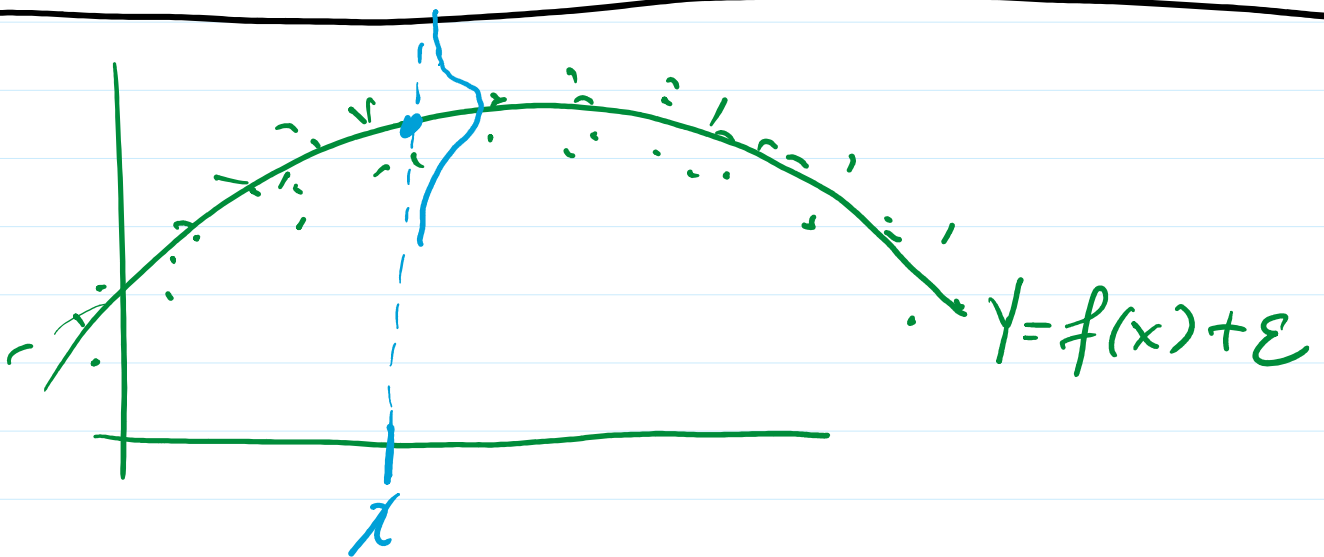
$$\frac{\partial}{\partial c}[\cdots] = -2E[Z] + 2c = 0$$

$$\frac{\partial}{\partial c}[\cdots] = -2E[Z] + 2C = 0$$

Solve for $c$ to get

$$C = E[Z].$$

So $\boxed{f^*(x) = E[Y \mid X = x]}$



$Y = f(x) + \varepsilon$

$x$

If $L(Y, f(x)) = |Y - f(x)|$

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

To realize one way to build $\hat{f}$ is as

In reality one way to build $\hat{f}$ is as

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

using training data.

e.g. $\hat{f}(\underline{x}) = $ avg. $y$s for $x$s near $\underline{x}$

Called KNN regression.

e.g. make some assumption about form of

$$E[Y | \underline{X} = \underline{x}]$$

maybe linear (?)

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

Called OLS lin. rgr.