

Statistical Learning

$$P_{XY}(X = x, Y = y)$$

Goal: Predict Y given X

Find a function η that minimizes

$$\mathbb{E}_{XY}[(Y - \eta(X))^2]$$

a metric
to evaluate
a prediction function

\uparrow
 P_{XY}

quadratic \rightarrow regression
other losses: l1
logistic
d1 loss

Thus far, we've been using η which is a:

- Linear functions of X
- Degree p polynomials of X
- Linear "generalization" of X

Statistical Learning

best / optimal predictor

$$P_{XY}(X = x, Y = y)$$

Goal: Predict Y given X

Find a function η that minimizes

$$\mathbb{E}_{XY}[(Y - \eta(X))^2] = \mathbb{E}_X \left[\mathbb{E}_{Y|X}[(Y - \eta(x))^2 | X = x] \right]$$

$$\eta(x) = \arg \min_c \mathbb{E}_{Y|X}[(Y - c)^2 | X = x] = \mathbb{E}_{Y|X}[Y | X = x]$$

Then:

Under LS loss, optimal predictor: $\eta(x) = \mathbb{E}_{Y|X}[Y | X = x]$

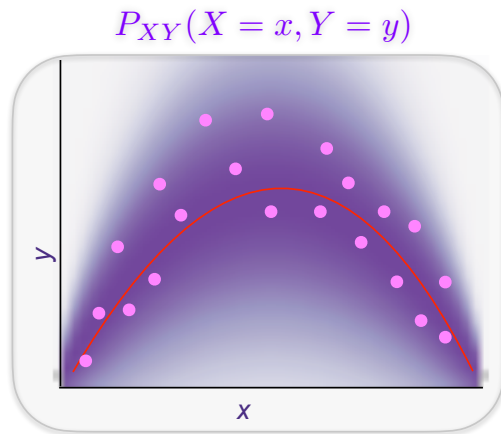
Optimal Prediction

$$\mathbb{E}_{XY}[(Y - \eta(X))^2] = \mathbb{E}_X \left[\mathbb{E}_{Y|X}[(Y - \eta(x))^2 | X = x] \right]$$

Under LS loss, optimal predictor: $\eta(x) = \mathbb{E}_{Y|X}[Y | X = x]$

$$\begin{aligned} \text{pf: } 0 &= \frac{d}{d\eta(x)} \mathbb{E}_{Y|X} [(Y - \eta(x))^2 | X=x] \\ &= \mathbb{E}_{Y|X} \left[\frac{d}{d\eta(x)} (Y - \eta(x))^2 \mid X=x \right] \\ &= \mathbb{E}_{Y|X} [-2 (Y - \eta(x)) \mid X=x] \\ &= -2 \mathbb{E}_{Y|X} [Y | X=x] + 2 \eta(x) \\ \Rightarrow \quad \eta(x) &= \mathbb{E}_{Y|X} [Y | X=x] \end{aligned}$$

Statistical Learning



Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

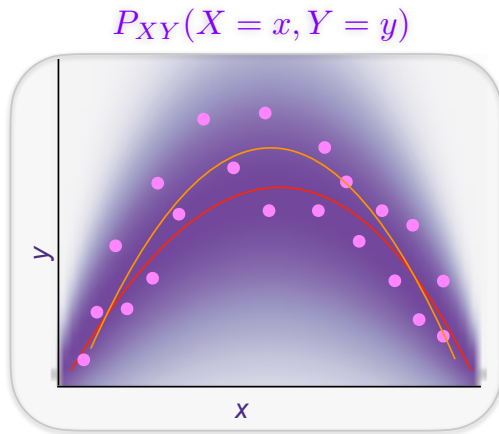
But we only have samples:

$$(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY} \quad \text{for } i = 1, \dots, n$$

$$\{(x_i, y_i)\}_{i=1}^n$$

\Rightarrow only estimate $\eta(x)$

Statistical Learning



Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

But we only have samples:

$$(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY} \quad \text{for } i = 1, \dots, n$$

and are restricted to a function class (e.g., linear) so we compute:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

We care about future predictions: $\mathbb{E}_{XY}[(Y - \hat{f}(X))^2]$

\mathcal{F} : predictor class

① linear

② quadratic

③ degree p poly

④ generalized linear

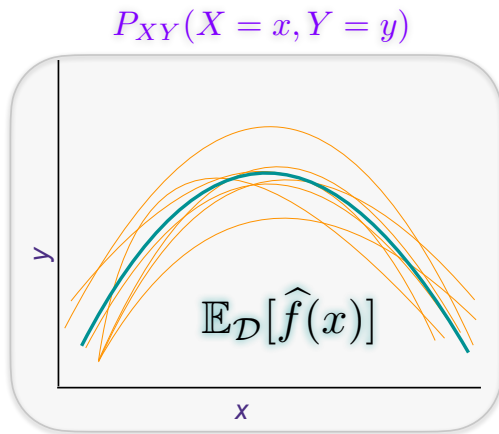
⑤ kernel

⑥ neural network

.....

Q: is \hat{f} a random or deterministic?

Statistical Learning



Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

But we only have samples:

$$(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY} \quad \text{for } i = 1, \dots, n$$

and are restricted to a function class (e.g., linear) so we compute:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

f is deterministic given D

Each draw $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ results in different \hat{f}

Bias-Variance Tradeoff

$$D = \left\{ (X_i, Y_i) \right\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} P_{XY}$$

new

$$\eta(x) = \mathbb{E}_{Y|X} [Y|X=x]$$

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

What we care

$$\begin{aligned} & \mathbb{E}_{XY} [(Y - \hat{f}_D(X))^2] \\ \Rightarrow & \mathbb{E}_{XY} [\mathbb{E}_D [(Y - \hat{f}_D(X))^2]] \\ = & \mathbb{E}_X [\mathbb{E}_{Y|X} \mathbb{E}_D [(Y - \hat{f}_D(X))^2] | X=x] \\ = & \mathbb{E}_X [\mathbb{E}_{Y|X} \mathbb{E}_D [(Y - \eta(x) + \eta(x) - \hat{f}_D(X))^2] | X=x] \\ = & \mathbb{E}_X [\mathbb{E}_{Y|X} \mathbb{E}_D [(Y - \eta(x))^2 + 2(Y - \eta(x))(\eta(x) - \hat{f}_D(X)) + (\eta(x) - \hat{f}_D(X))^2] | X=x] \\ = & \mathbb{E}_X [\mathbb{E}_{Y|X} [(Y - \eta(x))^2] + \mathbb{E}_{Y|X} \mathbb{E}_D [2(Y - \eta(x))(\eta(x) - \hat{f}_D(X)) | X=x] + \mathbb{E}_{Y|X} \mathbb{E}_D [(\eta(x) - \hat{f}_D(X))^2] | X=x] \end{aligned}$$

$$\begin{aligned} & \mathbb{E}_{Y|X} \mathbb{E}_D [(Y - \eta(x))(\eta(x) - \hat{f}_D(X)) | X=x] \\ & = \mathbb{E}_{Y|X} [(Y - \eta(x)) (\mathbb{E}_D [\eta(x) - \hat{f}_D(X)] | X=x)] \\ & \quad (\text{we know } \mathbb{E}_{Y|X} [Y] = \eta(x)) \end{aligned}$$

$$= 0$$

Bias-Variance Tradeoff

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X=x]$$

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

(1) $\mathbb{E}_{XY}[(Y - \eta(x))^2]$: irreducible error
independent of data
caused by stochastic

(2) $\mathbb{E}_{XY}[\mathbb{E}_D (Y(x) - \hat{f}_D(x))^2]$
 $= \mathbb{E}_X [\mathbb{E}_D (Y(x) - \hat{f}_D(x))^2]$: learning error
depend on D
· \mathcal{F} is too simple
· limited data

Bias-Variance Tradeoff

$$D = P_{XY}^n$$

D : diff over $\{(x_i, y_i)\}_{i=1}^n$
 X is independent of D

fix x

$$\mathbb{E}_D [f_0(x)]$$

$$= \int_{\{(x_i, y_i)\}_{i=1}^n} \hat{f}_{\{(x_i, y_i)\}_{i=1}^n}(x) dD(\{(x_i, y_i)\}_{i=1}^n)$$

$$\eta(x) = \mathbb{E}_{Y|X} [Y|X=x]$$

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

learning error
 fix x :

$$\mathbb{E}_D [(y(x) - \hat{f}_0(x))^2]$$

$$= \mathbb{E}_D [(y(x) - \mathbb{E}_D [\hat{f}_0(x)] + \mathbb{E}_D [\hat{f}_0(x)] - \hat{f}_0(x))^2]$$

$$= \mathbb{E}_D [(y(x) - \mathbb{E}_D [\hat{f}_0(x)]^2 + \underbrace{2(y(x) - \mathbb{E}_D [\hat{f}_0(x)])(\mathbb{E}_D [\hat{f}_0(x)] - \hat{f}_0(x))}_{=0} + (\mathbb{E}_D [\hat{f}_0(x)] - \hat{f}_0(x))^2]$$

$$= (y(x) - \mathbb{E}_D [\hat{f}_0(x)])^2 + \mathbb{E}_D [(\mathbb{E}_D [\hat{f}_0(x)] - \hat{f}_0(x))^2]$$

bias squared
variance
a random variable

independent of y

$$P_{XY} \Rightarrow \eta(x) = E[Y | X=x]$$

$$D = P_{XY}^n$$

$$\{(x_i, y_i)\}_{i=1}^n \sim D$$

$$\Rightarrow \hat{f}_D$$

$$E, [f_D]$$

Bias-Variance Tradeoff

$$\mathbb{E}_{Y|X}[\mathbb{E}_{\mathcal{D}}[(Y - \hat{f}_{\mathcal{D}}(x))^2] | X = x] = \underbrace{\mathbb{E}_{Y|X}[(Y - \eta(x))^2 | X = x]}_{\text{irreducible error}}$$

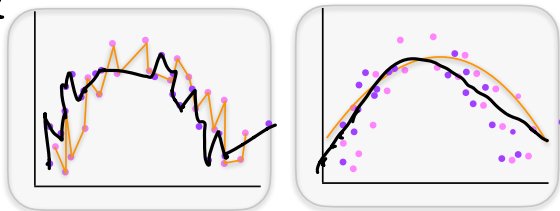
irreducible error

$$+ \underbrace{(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2}_{\text{bias squared}} + \underbrace{\mathbb{E}_{\mathcal{D}}[(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]}_{\text{variance}}$$

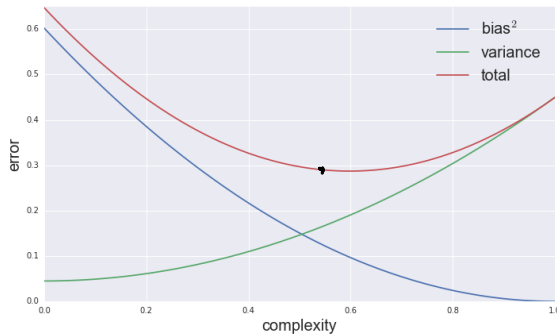
bias squared

variance

less complex
→ smaller
variance



If we re-drew our data, what the LS training error estimator look like for generalized linear functions in small p/large p dimensions?



more complex model
small bias
 $\hat{f} \leftarrow \arg \min_{\tilde{f}} \frac{1}{n} \sum_{i=1}^n (\tilde{f}(x_i) - y_i)^2$
 $\exists \tilde{f} \approx y$

How many points do I use for training/testing?

- > **Very hard question to answer!**

- Too few training points, learned model is bad
- Too few test points, you never know if you reached a good solution

- > **Bounds, such as Hoeffding's inequality can help:**

$$P(|\hat{\theta} - \theta^*| \geq \epsilon) \leq 2e^{-2N\epsilon^2}$$

- > **More on this later the quarter, but still hard to answer**

- > **Typically:**

- If you have a reasonable amount of data 90/10 splits are common
- If you have little data, then you need to get fancy (e.g., bootstrapping)