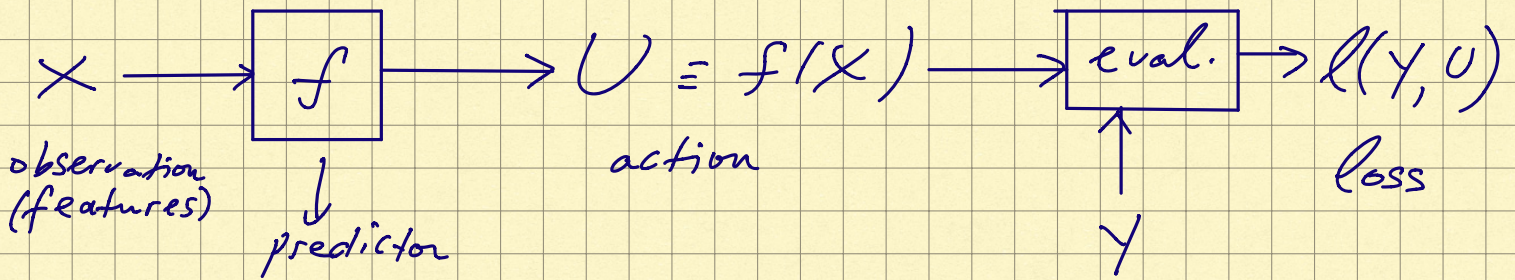


# Recap: prediction, known distribution

$$(X, Y) \sim P$$



$$\text{Loss (or risk) of } f: L_P(f) := \mathbb{E}_P[l(Y, f(X))]$$

$$\text{Optimality: } f_P^* = \underset{f}{\operatorname{argmin}} L_P(f)$$

## Examples:

1) binary classification

$Y, U$  take values in  $\{0, 1\}$

$$l(y, u) = \mathbb{1}_{\{y \neq u\}}$$

$$L_P^* = \min_f L_P(f) = \mathbb{E}_P[\min\{1 - \eta(x), \eta(x)\}]$$

$$f_P^*(x) = \mathbb{1}_{\{\eta(x) \geq 1/2\}}$$

$$\text{where } \eta(x) := \mathbb{P}[Y=1 | X=x] = \mathbb{E}[Y | X=x]$$

2) MMSE estimation

$$X \in \mathbb{R}^d; \quad Y, U \in \mathbb{R}$$

$$l(y, u) = (y - u)^2$$

$$L_P^* = \mathbb{E}_P \left( Y - \underbrace{\mathbb{E}(Y|X)}_{f_P^*(X)} \right)^2$$



3) "soft" binary classification

$$y \in \{0, 1\}$$

$$u \in [0, 1]$$

$$l(y, u) = |y - u|$$

Takeaway message: if  $P$  (dist. of  $(x, y)$ ) is known, no learning takes place

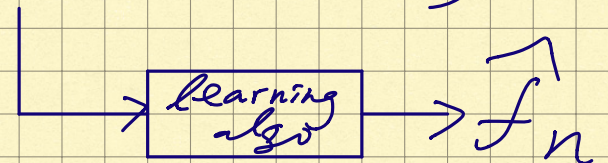
Learning:  $P$  unknown

$(x_1, y_1), \dots, (x_n, y_n) \stackrel{iid}{\sim} P$  - data

$$Z = (x, y)$$

$$z_1, \dots, z_n \stackrel{iid}{\sim} P$$

$$Z^n := (z_1, \dots, z_n)$$



$\hat{f}_n$  is a random function from features  $X$  into actions  $U$

$$L_P(\hat{f}_n) = \int_{X \times Y} l(y, \hat{f}_n(x)) P(dx, dy)$$

$$= \mathbb{E}[l(Y, \hat{f}_n(X)) | Z^n]$$

where  $(x, y) \sim P$ , indep. of  $Z^n$

$$\underbrace{z_1, \dots, z_n, z_{n+1}}_{\text{data}} \stackrel{iid}{\sim} P$$

$$\Rightarrow L_P(\hat{f}_n) = \mathbb{E}[l(Y_{n+1}, \hat{f}_n(X_{n+1})) | Z^n]$$



Important:  $L_P(\hat{f}_n)$  cannot be computed w/o knowledge of  $P$ !

$$\underbrace{z_1, \dots, z_n}_{\text{training data}} \xrightarrow{LA} \hat{f}_n$$

$z'_1, \dots, z'_n \stackrel{iid}{\sim} P, \perp\!\!\!\perp z^n$  (validation/test set)

$$\hat{L}_P(\hat{f}_n) = \frac{1}{n} \sum_{i=1}^n l(y'_i, \hat{f}_n(x'_i))$$

Cond. on  $z^n$ ,  $\hat{L}_P(\hat{f}_n) \approx L_P(\hat{f}_n)$

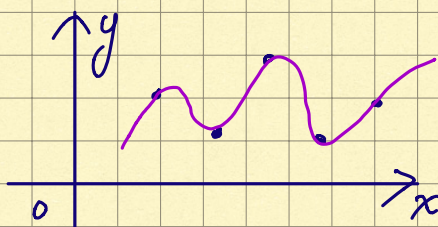
Goals:  $L_P(\hat{f}_n)$  vs.  $L_P^* \equiv \min_f L_P(f)$

$$L_P(\hat{f}_n) \geq L_P^*$$

- want to make excess loss  $L_P(\hat{f}_n) - L_P^*$  small

Interpolation/memorization?

$$\hat{f}_n(x_i) = y_i \quad i=1, \dots, n$$



overfitting:  $\hat{f}_n$  interpolates but  $L_P(\hat{f}_n) \gg L_P^*$

- inductive bias: learning algo is restricted to some structured class of predictors

$\hat{f}_n \in \mathcal{H}$  (hypothesis space)  
- typically restricted in comp.



$$L_p(\hat{f}_n) - L_p^*(\mathcal{H}) =: \text{excess risk relative to } \mathcal{H}$$

where  $L_p^*(\mathcal{H}) := \min_{f \in \mathcal{H}} L_p(f)$

$$L_p(\hat{f}_n) - L_p^*$$

$$= \underbrace{L_p(\hat{f}_n) - L_p^*(\mathcal{H})}_{\text{estimation error}} + \underbrace{L_p^*(\mathcal{H}) - L_p^*}_{\text{appr. error}}$$

Belkin-Rakhtin-Tsybakov (AISTATS 2019):  
 "Does interpolation contradict statistical optimality?"

Exist settings where  $\hat{f}_n$  interpolates the data, but is still asymptotically consistent:

$$\mathbb{E}(\hat{f}_n(x) - Y)^2 \rightarrow \min_{f \in \mathcal{H}} \mathbb{E}(f(x) - Y)^2$$

provided  $\eta(x) = \mathbb{E}[Y|X=x] \in \mathcal{H}$  as  $n \rightarrow \infty$

Double descent (Belkin et al.)

