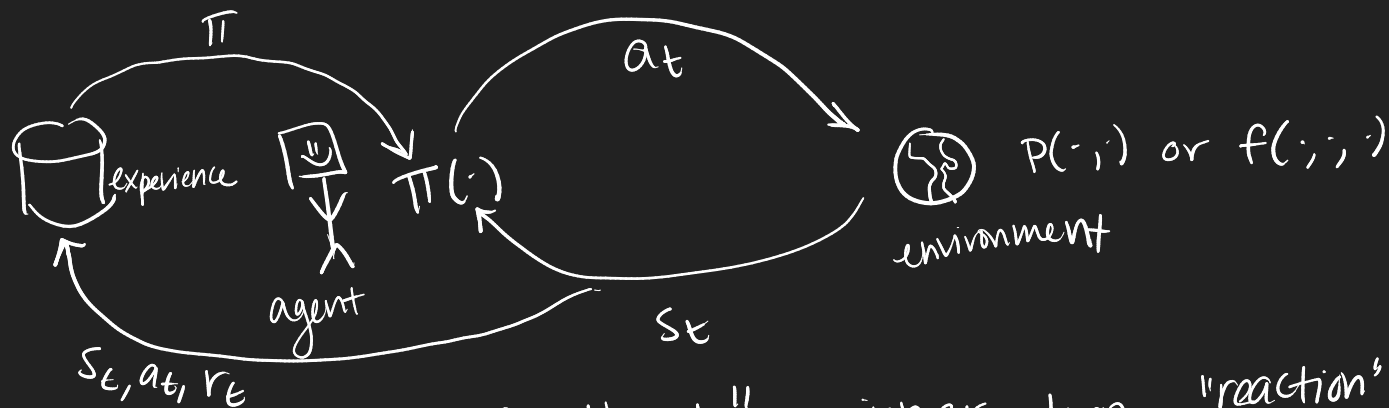


1) Types of Feedback in RL

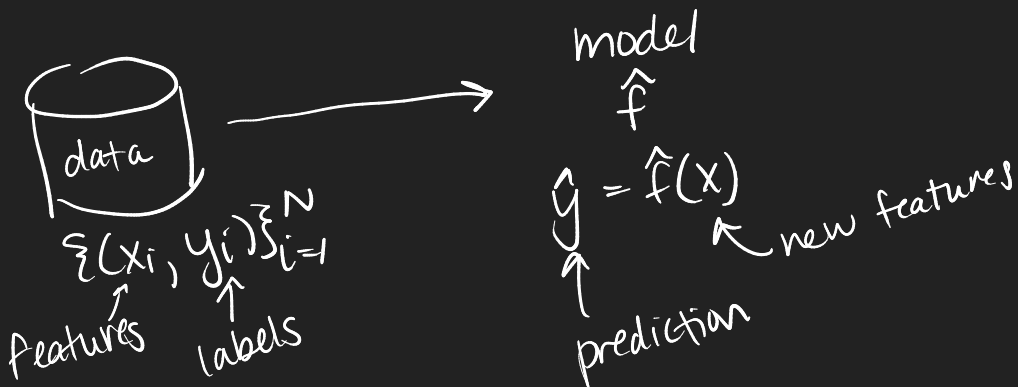


1) "control feedback" - inner loop "reaction"
 ex - thermostat regulating temperature

Unit 1
 2) "data feedback" - "outer" loop "adaptation"
 ex - smart thermostat learn preferences
 Unit 2+

From now on, transition / dynamics are unknown
 $P(-, \cdot)$ $f(-, \cdot, \cdot)$
 (often $r(-, \cdot)$ also unknown)
 reward

2) Supervised Learning (SL)



e.g. classification of images x - image (pixels) y - cat vs. dog
 regression x - financial history y - probability of loan repayment

Analogy between SL and special case RL

$$\min_f \mathbb{E} [\ell(y, \hat{y}) \mid \hat{y} = f(x)]$$

$(x, y) \sim \mathcal{D}$

features / states
no impact predictions / actions dependence
model / policy offline / online
loss / cost distribution / transition probabilities
iid / dependent

What should we learn to solve (i.e. find near-optimal policy)

$$\mathcal{M} = \{ \mathcal{S}, \mathcal{A}, \underbrace{P}_{\text{unknown}}, \underbrace{r}_{\text{unknown}}, \gamma \}$$

- transition P , reward r
- value function V^π / Q function Q^π } plug in π & observe
- optimal value/Q function V^*, Q^*
- optimal policy π^*

supervision

- ✓ $s_{t+1} \sim P(s_t, a_t)$
 $r_t = r(s_t, a_t)$
- ✓ \approx after delay $\sum_0^T \gamma^t r_t$
- X
- X exception: imitation learning

3) Estimation & Prediction

A) Tabular Setting:

suppose $x_i \stackrel{iid}{\sim} \mathcal{D}$, $x_i \in \mathcal{X}$, $p(x) = \mathbb{P}_{\mathcal{D}}(x_i = x)$

$$\hat{p}(x) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}\{x_i = x\}$$

How good is \hat{p} ?

Lemma: (consistency)

$$\mathbb{E} [\hat{p}(x)] = p(x)$$

x_1, \dots, x_N

Proof:
$$\begin{aligned} \mathbb{E}[\hat{p}(x)] &= \frac{1}{N} \sum_{i=1}^N \mathbb{E}[\mathbb{1}\{x_i = x\}] \\ &= \frac{1}{N} \sum_{i=1}^N \mathbb{P}(x_i = x) \\ &= p(x) \end{aligned}$$

Theorem (concentration)

For all $x \in \mathcal{X}$, with probability $1 - \delta$

"pointwise"
$$|\hat{p}(x) - p(x)| \leq \sqrt{\frac{2 \log(2/\delta) |\mathcal{X}|}{N}} \approx O\left(\frac{1}{\sqrt{N}}\right)$$

Proof out of scope (Hoeffding's Inequality)

$x, y \sim \mathcal{D}$ $\hat{f}_N \approx y$ $\{(x_i, y_i)\}_{i=1}^N$
By similar approach,

$$\hat{f}(x) = \frac{\sum_{i=1}^N y_i \mathbb{1}\{x = x_i\}}{\sum_{i=1}^N \mathbb{1}\{x = x_i\}}$$

Details out of scope, but if $y = f^*(x) + w$ \downarrow iid noise

$\forall x \in \mathcal{X}$, w.p. $1 - \delta$

pointwise
$$|\hat{f}(x) - f^*(x)| \lesssim \sqrt{\frac{|\mathcal{X}| \log(1/\delta)}{N}}$$

applies when $|\mathcal{X}| < N$

B) Non-tabular:

Empirical Risk Minimization

$$\hat{f} = \underset{\substack{\text{function} \\ \text{class}}}{f \in \mathcal{F}} \operatorname{argmin} \frac{1}{N} \sum_{i=1}^N \underset{\text{loss}}{\ell} \left(\underset{\text{label}}{y_i}, \underset{\text{prediction}}{f(x_i)} \right)$$

1) Estimation (Parametric)

$$\mathcal{F} = \{ f_{\theta}(x) \mid \theta \in \mathbb{R}^d \}$$

eg. neural network θ -weights
 e.g. $f_{\theta}(x) = \theta^T \phi(x)$
 ↑ fixed

Suppose $y = f_{\theta^*}(x) + w$ ← iid noise
 ↑ true parameter

Estimation Error: $\|\theta^* - \hat{\theta}\|$ $\hat{f} = f_{\hat{\theta}}$

Details out of scope, estimation error bounded

parametric $\|\theta^* - \hat{\theta}\| \lesssim \sqrt{\frac{d \log(1/\delta)}{N}}$
 ↑ hiding constants

useful when $N > d$

Prediction

expected prediction error on $(x, y) \sim \mathcal{D}$

$$\mathbb{E}_{x, y \sim \mathcal{D}} [\ell(\hat{f}(x), y)]$$

Assume $\mathcal{D}: x \sim \mathcal{D}_x$ and $f^* \in \mathcal{F}$
 $y = f^*(x) + w$ ← noise

Details out of scope, but

$$\mathbb{E}_{x, y \sim \mathcal{D}} [\ell(\hat{f}(x), y)] \lesssim \sqrt{\frac{\log(1/\delta)}{N}}$$

Prediction error is about average case and fixed distribution

