

# FA Prediction

---

Mon

---

---

---

---



# Supervised Learning

\* Samples arrive in the form of tuples of input and target  
output drawn i.i.d.:

$$\left( \underbrace{\underline{x}_k}_{\substack{\text{considered to be a} \\ \text{vector here.} \\ \text{input}}} \text{, } \underbrace{y_k}_{\substack{\text{As RV.} \\ \text{the target} \\ \text{output.} \\ \text{input 来源地)}} \right) \sim \text{true distribution. } P_{\underline{x}, y}(\underline{o}, o)$$

outcome  
 for first RV ( $\underline{x}$ )  
 output for second  
 RV ( $y$ ).  
 take two inputs.

Here  $y \in \mathbb{R}$  is a scalar and  $\underline{x}_k \in \mathbb{R}^n$   
 $\text{n-dimensional column}$   
 is any vector denoted by an underline  
 (boldface in the text book).

Also,  $P_{\underline{x}, y}(\underline{a}, c) = P(\underbrace{\underline{x} = \underline{a}, y = c}_{\substack{x \text{ takes value } a. \\ y \text{ takes value } c}})$   
 the joint distribution of  $\underline{x}, y$ .

Dynamic Programming: all off-line. Although can do it online.

Reinforcement learning: expert interacts with world, continue over time. (Storing things in table is not flexible).

## (Approximation)

- \* The learner wants to predict the target output accurately based on the observed input using a parameterized function  $f$ :

function  $f$  maps input  $\underline{x}$  to the output  $\hat{y}$ .  
the learner observes  $\underline{x}$ , and predicts  $\hat{y}$ .

$$f(\underline{x}, \underline{w}) = \hat{y} \approx y$$

from by adding weight to input  $\underline{x}$ .  $\therefore \underline{x}$  is a vector.  
can be a table, a linear function, polynomial ...

use input to approximate output by this function. (try to predict the reward).

Here, the function  $f$  takes  $\underline{x}$  as an input, is parameterized by a weight vector  $\underline{w}$ , and produces an output  $\hat{y}$ , which is a prediction or approximation of target output  $y$ .



Value functions ( $V_{\pi}, q_{\pi}$ ) depend on the policy  $\pi$ . policy change, tiny change.

# (objective)

Formalize the goal:

- \* The learner's prediction accuracy depends on the value of  $\underline{w}$ .  
 The learner wants to minimize the prediction error  $y - f(\underline{x}, \underline{w})$  over different samples. This goal can be well formulated as the minimization of the Mean Square Error (MSE) objective:

$\text{MSE}(\underline{w}) = \mathbb{E}_{p_{\underline{x}, y}} \left[ (y - f(\underline{x}, \underline{w}))^2 \right]$

What the agent have  
A function of the weight vector, gonna give different values to different  $\underline{w}$ .

Expectation is the true quantity.

It is the expected squared error between the target and predicted outputs.

→ expectation of the square error.

→ neural network the  $\underline{w}$  that minimize this square error.

- \* Note that

$$\text{MSE}(\underline{w}) = \sum p_{\underline{x}, y}(\underline{x}, c) \underbrace{(c - f(\underline{x}, \underline{w}))^2}_{\text{takes this outcome:}}$$

according to (which law?).

# (Stochastic gradient descent) or SGD

SGD forms an unbiased estimate of the objective. For example, if the learner observes  $(\underline{x}, Y)$ , then the squared prediction error is an unbiased estimate of MSE:  $(Y - f(\underline{x}, \underline{w}))^2$ . <sup>2</sup> the sample objective.

SGD then updates the weight vector  $\underline{w}$  incrementally using the gradient of the squared error:

$$\underline{w}_{k+1} \doteq \underline{w}_k - \frac{1}{2} \alpha \nabla_{\underline{w}_k} \left( Y_k - f(\underline{x}_k, \underline{w}_k) \right)^2$$

the gradient of  $(Y_k - f(\underline{x}_k, \underline{w}_k))^2$   
unbiased estimate of the MSE.

$$= \underline{w}_k + \alpha \left( Y_k - f(\underline{x}_k, \underline{w}_k) \right) \nabla_{\underline{w}_k} f(\underline{x}_k, \underline{w}_k).$$

the gradient of the function HSEF.

Stochastic gradient  $\Rightarrow$   
update: (we get  
the main point like  
 $\frac{1}{n} \sum (Y_i - f(\underline{x}_i, \underline{w}))^2$ )

For linear function approximation,  
 $\underline{w} \in \mathbb{R}^n$  is a column vector of the  
 same size as  $\underline{x}$ :

$$\underline{f}(\underline{x}, \underline{w}) = \underline{x}^T \underline{w} = \sum_{i=1}^n x_i w_i,$$

where  $x_i$  and  $w_i$  are  $i$ th elements  
 of vectors  $\underline{x}$  and  $\underline{w}$ , respectively.

Then the SGD update reduces to

$$\underline{w}_{k+1} = \underline{w}_k + \alpha (Y_k - \underline{f}(\underline{x}_k, \underline{w}_k)) \underline{x}_k,$$

because

$$\nabla_{\underline{w}} \underline{f}(\underline{x}, \underline{w}) = \nabla_{\underline{w}} (\underline{x}^T \underline{w}) = \underline{x}.$$

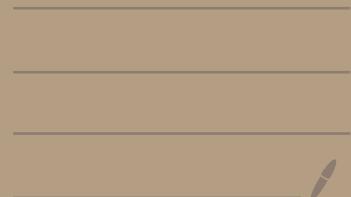
the input vector  $\underline{x}$ .

# Prediction with

---

## Approximation

### Wed



Deepak Ranganatha Sastry Mamillapalli 下午1:09

对所有人说

DR can rewards be modelled as a function  
of states and actions? 



## [Mean Squared Value Error]

What can be an objective for policy evaluation? Let's consider function approximation as described in Section 9.1:

$\hat{v}(s, w) \approx \underline{v_{\pi}(s)}$ . Then the policy evaluation / value prediction problem can be seen as a supervised learning problem with input, target output pairs as:  $(s_t, \underline{v_{\pi}(s_t)})$ .

Then we can use the MSE objective we introduced before, which we call Mean Squared Value Error (VE):

$$\overline{VE}(w) = E_{s_t \sim \mu} \left[ (v_{\pi}(s_t) - \hat{v}(s_t, w))^2 \right],$$

where  $\mu$  is some state distribution  $\text{iid}$

## [ Mean Squared Return Error ]

$\overline{VE}$  can be re-written as:

$$\overline{VE}(\underline{w}) = \sum_s \mu(s) \left( v_{\pi}(s) - \hat{v}(s, \underline{w}) \right)^2.$$

Interestingly, minimizing  $\overline{VE}$  is equivalent to minimizing the following objective we call Mean Squared Return Error (MSRE):

$$MSRE(\underline{w}) = E_{\underbrace{S_t \sim \mu, A \sim \pi}_{\text{Assume $S_t$ is not iid}}} \left[ (G_t - \hat{v}(s_t, \underline{w}))^2 \right].$$

agent takes state \$S\_t\$ as input and predict expected return \$G\_t\$ as output.

It is a different objective with the same solution / minimum.

## [ Deriving Monte Carlo Method as SGD ]

Then an unbiased estimate can be easily obtained as  $(Q_t - \hat{v}(s, w))^2$  if we can sample states from  $\mu$ . If  $\mu$  is the on-policy distribution, then we can just use the sequential experience of  $S_t$  and make SGD updates:

in tabular case:

$$w \leftarrow w + \alpha [Q_t - \hat{v}(s_t, w)] \nabla_w \hat{v}(s_t, w)$$

the Monte Carlo error.

gradient of the approximation.

which is the gradient MC pseudocode from Section 9.3.

it's for tabular case,  
we don't have gradient!!