

In this section, we consider a more general problem setting.

An agent applies a sequence of actions  $x_1, x_2, x_3, \dots$  to a system, selecting from a set  $X$ .  $X$  could be either finite or infinite. After applying  $x_t$ , the agent observes  $y_t$ , which is generated from a conditional probability measure  $q_{\theta}(y|x_t)$ . The agent then collects a reward  $r_t = r(y_t)$ , where  $r$  is a known function. The agent is initially uncertain about the value of  $\theta$ , and he represents his uncertainty using a prior distribution  $\mathcal{P}$ .

We first list the Greedy Algorithm and TS algorithm and then explain.

Greedy ( $X, p, q, r$ )

for  $t=1, 2, \dots$ , do

# estimate  $\theta$  from prior

$$\hat{\theta} = \mathbb{E}_{\mathcal{P}}[\theta]$$

# select and apply action

$$x_t = \operatorname{argmax}_{x \in X} \mathbb{E}_{q_{\hat{\theta}}} [r(y_t) | x_t = x]$$

Apply action  $x_t$  and get  $y_t$ .

# update prior

$$\mathcal{P} = \operatorname{Pr}_{p, q}(\theta \in \cdot | x_t, y_t)$$

Thompson Sampling ( $X, p, q, r$ )

for  $t=1, 2, \dots$  do

# sample prior

$$\hat{\theta} \sim \mathcal{P}$$

# select and apply action

$$x_t = \operatorname{argmax}_{x \in X} \mathbb{E}_{q_{\hat{\theta}}} [r(y_t) | x_t = x]$$

# update prior

$$\mathcal{P} = \operatorname{Pr}_{p, q}(\theta \in \cdot | x_t, y_t)$$

For Greedy.

\*  $\hat{\theta}$  is the optimal estimator from  $\theta$  with MSE.  $\Rightarrow \mathbb{E}_p(\theta)$

\* Suppose  $y \in \mathcal{Y}$ , then by going through all possible outcomes, and figure out which action, on average, rewards the most,

$$\mathbb{E}_{q_{\hat{\theta}}} [r(y_t) | x_t = x] = \sum_{y \in \mathcal{Y}} q_{\hat{\theta}}(y | x) \cdot r(y)$$

\* The prior distribution  $p$  is updated based on observation  $y_t$ .

If  $\theta$  is sampled from a finite set, apply Bayes Rule, the estimation of  $p$  is updated follow: (for each value of  $u$ )

$$P(\theta = u | x_t, y_t) = \frac{P(u) q_u(y_t | x_t)}{\sum_v P(v) q_v(y_t | x_t)}$$

old prior

Example 1. (Independent Travel Times)

Background:

An agent commutes from home to work every morning. She would like to minimize the travel time. How can she learn efficiently and minimize the expected travel time?

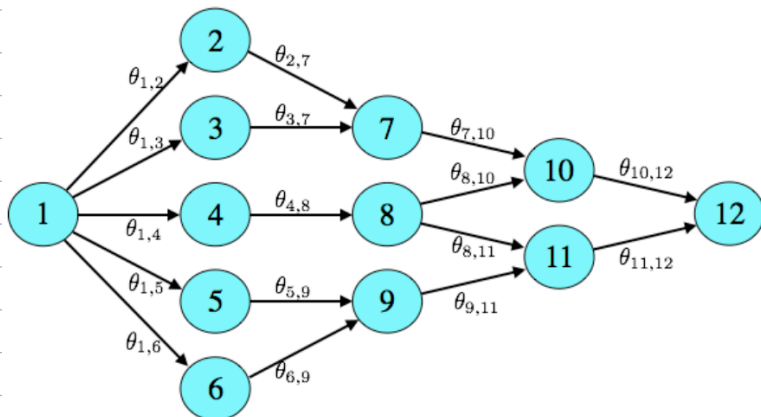


Figure 1.1: Shortest path problem.

Let  $G = (V, E)$  be the graph representing the route.

$$V = [N], \quad E = \{(i, j) \mid i, j \in V, \text{edge}(i, j) \text{ in graph}\}.$$

Vertex 1 is the source

Vertex  $N$  is the destination.

An action is a sequence of distinct edges from 1 to  $N$ , or a path.

After applying action  $\kappa_t$ , the travel time at edge  $e \in \kappa_t$  is

$Y_{t,e}$  and  $Y_{t,e}$  is independently sampled from a distribution with mean  $\theta_e$ .

The cost (negative of the reward) is  $\sum_{e \in \kappa_t} Y_{t,e}$ .

---

Consider a prior for  $\theta_e$  to be Log-Gaussian with  $\mu_e$  and  $\sigma_e^2$ , i.e.  $\ln(\theta_e) \sim \mathcal{N}(\mu_e, \sigma_e^2)$ , then.

$$f_{\theta}(\theta) = \frac{1}{\kappa} \cdot \frac{1}{\sqrt{2\pi} \sigma_e} \cdot \exp\left(-\frac{(\ln \theta - \mu_e)^2}{2\sigma_e^2}\right)$$

$$\theta \Rightarrow \mathbb{E}[\theta] = \exp\left(\mu_e + \frac{\sigma_e^2}{2}\right)$$

$$\text{Var}[\theta] = \left[\exp(\sigma_e^2) - 1\right] \left[\exp(2\mu_e + \sigma_e^2)\right].$$

In the problem setting, we assume  $Y_{t,e} | \theta$  is independent across  $e \in E$ , so,  $\mathbb{E}[Y_{t,e} | \theta_e] = \theta_e \leftarrow$  problem setting.

and  $Y_{t,e} | \theta_e \sim \text{log-Gaussian}$  with parameter  $(\ln \theta_e - \frac{\sigma_e^2}{2}, \sigma_e^2)$  where  $\sigma_e$  is known.

Conjugacy property inspires the following update when  $Y_{t,e}$  is observed.

$$(\mu_e, \sigma_e^2) \leftarrow \left( \frac{\frac{1}{\sigma_e^2} \mu_e + \frac{1}{\sigma_e^2} (\ln(Y_{t,e}) + \frac{\sigma_e^2}{2})}{\frac{1}{\sigma_e^2} + \frac{1}{\sigma_e^2}}, \frac{1}{\frac{1}{\sigma_e^2} + \frac{1}{\sigma_e^2}} \right)$$

Suppose the agent knows the distance of each edge  $d_e$ , but is not sure about the travel time  $Y_{t,e}$ . Initially, it may be a good idea to let  $Y_{t,e} \propto d_e$ , or even more brutally let  $\mathbb{E}[Y_{t,e}] = d_e$ . So initially, we have the averaged travel time

$$\mathbb{E}[Y_{t,e}] = \exp(\mu_e + \frac{\sigma_e^2}{2}) = d_e$$

$$\Rightarrow \mu_e = \ln(d_e) - \frac{\sigma_e^2}{2}$$

However, from a single  $d_e$ , it's impossible to determine  $\mu_e$  and  $\sigma_e^2$  simultaneously. Any value can be possible.

Note that,  $\text{Var}[Y_{t,e}] = (\exp(\sigma_e^2) - 1) \cdot d_e^2$ .

After that, Greedy algorithm as TS algorithm can be applied as follows:

At the beginning of each round, the agent has  $(\mu_e, \sigma_e^2)$  for each edge from previous trials.

① \* For Greedy Algorithm,  $\hat{\theta}_e = \mathbb{E}_p[\theta_e] = \exp(\mu_e + \frac{\sigma_e^2}{2})$

\* For Thompson Sampling,  $\hat{\theta}_e$  is drawn from a log-Gaussian with mean  $(\mu_e, \sigma_e^2)$

② Then, each algorithm select its path to maximize

$$\mathbb{E}_{g \in \mathcal{G}} [r(y_t) | x_t = x] = - \sum_{e \in x_t} \hat{\theta}_e$$

This can be solved efficiently using Dijkstra for example.

③ Apply  $x_t$ .

④ Update  $\mu_e$  and  $\sigma_e$  for each involved edge.

## Example 2. (Correlated Travel Time)

We modify the shortest path problem.

We change the observation time model to

$$y_{t,e} = \underbrace{\zeta_{t,e}}_{\substack{\text{the idiosyncratic factor} \\ \text{with edge } e.}} \cdot \underbrace{\eta_t}_{\substack{\text{a factor common} \\ \text{to all edges.}}} \cdot \underbrace{U_{t,l(e)}}_{\substack{l(e) \text{ indicates whether edge } e \\ \text{locates in the lower half} \\ \text{of the binomial bridge.} \\ l(e) \in \{0, 1\}.}} \cdot \theta_e$$

$\zeta_{t,e}$ .

the idiosyncratic factor  
with edge  $e$ .

a factor common  
to all edges.

$l(e)$  indicates whether edge  $e$   
locates in the lower half  
of the binomial bridge.  
 $l(e) \in \{0, 1\}$ .

$U_{t,1}$  : edge  $e \in$  lower half

$U_{t,0}$  :  $e \in$  upper half.

We let  $\zeta_{t,e}$ ,  $\eta_t$ ,  $U_{t,0}$ ,  $U_{t,1}$  to be independent log-Gaussian with parameter  $(-\frac{\tilde{\sigma}^2}{6}, \frac{\tilde{\sigma}^2}{3})$ . Their distributions are known, but  $\theta_e$  also subjects to log-Gaussian  $(\mu_e, \sigma_e^2)$  with unknown  $\mu_e$  and  $\sigma_e^2$ .

Given  $\zeta_{t,e}$ ,  $\eta_t$ ,  $U_{t,0}$ ,  $U_{t,1}$ , the marginal distribution

$y_{t,e} | \theta$  is identical to Example 1.

- Common factor  $\eta_t$  is used to model a global factor, say weather.

-  $U_{t,0}$ ,  $U_{t,1}$  reflect events affecting only half of the path. Say, a bridge locates in the middle of the path.

Conjugate properties benefits the updating process.

$$\text{Let } \phi_e = \ln(\theta_e), \text{ and } z_{t,e} = \begin{cases} \ln(y_{t,e}) & \text{if } e \in X_t \\ 0 & \text{o.w.} \end{cases}$$

$$z_t = \begin{bmatrix} z_{t,1} \\ z_{t,2} \\ \vdots \\ z_{t,N} \end{bmatrix}$$

We then formulate a covariance matrix  $\hat{\Sigma}_{e,e'} \in \mathbb{R}^{|X_t| \times |X_t|}$

$$\hat{\Sigma}_{e,e'} = \begin{cases} \hat{\sigma}^2 & \text{if } e=e', \\ \frac{2}{3} \hat{\sigma}^2 & \text{if } e \neq e' \text{ but } l(e) = l(e') \\ \hat{\sigma}^2/3 & \text{o.w.} \end{cases} \quad \text{for } e, e' \in X_t$$

$$\hat{C}_{e,e'} = \begin{cases} \hat{\Sigma}_{e,e'}^{-1} & \text{if } e, e' \in X_t \text{ for all } e, e' \in E \\ 0 & \text{o.w.} \end{cases}$$

$$\text{Then } (\mu, \Sigma) = \left( (\Sigma^{-1} + \hat{C})^{-1} (\Sigma^{-1} \mu + \hat{C} \cdot z_t), (\Sigma^{-1} + \hat{C})^{-1} \right)$$

Again, we can use Greedy as TS.