

## Bayesian Learning / Posterior Evolution / Bayesian Bandits

Source: Bandit Algorithms, L&S., Chap. 34

Setting: Prior information on Arms ; prior characterized in terms of a given arm marginal distribution. The approach is to evolve the posterior distribution of the arms based on history of samples, and use this posterior to guide future actions.

Posterior dist: Conditional pdf / pmf of an arm reward distribution based on past samples.

(Several technical issues regarding existence, uniqueness.  
Please read Sec. 34.2 in text). In these notes,  
we will implicitly assume everything is well defined, and  
not worry about technical issues.

In a Bayesian setting, information about each arm fully characterized by the posterior dist. In general, computing the posterior is tricky. However, in several important cases, the form of the posterior distribution is the same as that of the prior,

only with new parameters. Such distributions are called **Conjugate pairs** and the associated prior is a **conjugate prior**.

Conjugate Prior | Pairs

(1) Gaussian model:  $\{x_i, i=1, 2, 3, \dots\}$  a conditionally iid sequence of Gaussian r.v.s with unknown mean but known variance. Specifically, let

$\Theta$  be a proper random variable. Then,

$$P(x_1, \dots, x_n | \theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_s^2}} e^{-\frac{(x_i - \theta)^2}{2\sigma_s^2}}$$

$\underbrace{\phantom{\dots}}_{\substack{\text{known} \\ \text{variance}}}$

Further,  $\Theta$  itself is a Gaussian r.v., with known mean and variance.

$$\Theta \sim N(\mu_p, \sigma_p^2) \quad \text{the prior distribution}$$

i.e.,  $P_\theta(\theta) = \frac{1}{\sqrt{2\pi\sigma_p^2}} e^{-\frac{(\theta - \mu_p)^2}{2\sigma_p^2}}$

Then, the posterior distribution on the (unknown) mean  $\theta$  is given by:

$$P_{\theta|x}(\theta|x) \sim N\left(\frac{\frac{\mu_p}{\sigma_p^2} + \frac{x}{\sigma_s^2}}{\frac{1}{\sigma_p^2} + \frac{1}{\sigma_s^2}}, \frac{1}{\frac{1}{\sigma_p^2} + \frac{1}{\sigma_s^2}}\right)$$

$\underbrace{\frac{\mu_p}{\sigma_p^2} + \frac{x}{\sigma_s^2}}$        $\underbrace{\frac{1}{\sigma_p^2} + \frac{1}{\sigma_s^2}}$   
updated mean      updated variance

Since the prior and posterior are both Gaussian, and the posterior simply updates the mean and variance of  $\theta$  as above, the posterior evolution given samples  $x_1, x_2, \dots, x_n$  is easy to describe in closed form. Specifically,

$$P_{\theta|x_1, \dots, x_n}(\theta|x_1, \dots, x_n) \sim N\left(\frac{\frac{\mu_p}{\sigma_p^2} + \frac{\sum_{i=1}^n x_i}{\sigma_s^2}}{\frac{1}{\sigma_p^2} + \frac{n}{\sigma_s^2}}, \frac{1}{\frac{1}{\sigma_p^2} + \frac{n}{\sigma_s^2}}\right)$$

$\underbrace{\frac{\mu_p}{\sigma_p^2} + \frac{\sum_{i=1}^n x_i}{\sigma_s^2}}$        $\underbrace{\frac{1}{\sigma_p^2} + \frac{n}{\sigma_s^2}}$

(2) Beta-Bernoulli Model:  $\{x_i, i=1, 2, \dots\}$  are conditionally iid Bernoulli( $\theta$ ) r.v.s., given  $\tilde{\theta} = \theta$  (i.e., Bernoulli with unknown mean).

Initial Prior (before seeing any samples):

$$\theta \sim \text{Uniform}[0, 1]$$

$$P(x_1 = x_1, x_2 = x_2, \dots, x_n = x_n \mid \theta = \theta) \quad x_i \in \{0, 1\} \\ \theta \in [0, 1]$$

$$= \prod_{i=1}^n P(x_i = x_i \mid \theta = \theta)$$

Let  $S_n = \sum_{i=1}^n x_i$  ( $S_n = \# \text{ of } 1s$ )  
 $F_n = (n - S_n)$ . ( $F_n = \# \text{ of } 0s$ )

$$= \theta^{S_n} (1-\theta)^{F_n} \quad \xrightarrow{\text{Simplification}} \quad S_n = \sum_{i=1}^n x_i$$

Posterior Dist. of  $\theta$  given  $(x_1, \dots, x_n)$ :

$$P_{\theta \mid (x_1, \dots, x_n)}(\theta \mid x_1, \dots, x_n)$$

$$= P(x_1 = x_1, \dots, x_n = x_n \mid \theta = \theta) P_{\theta}(\theta)$$

$$\frac{1}{\int_0^1 P(x_1 = x_1, \dots, x_n = x_n \mid \theta = z) P_{\theta}(z) dz}$$

$$= C \theta^{s_n} (1-\theta)^{n-s_n} \rightarrow \textcircled{1}$$

normalizing constant.

Aside: Beta Distribution:  $\gamma \sim \text{Beta}(\alpha, \beta)$  if  
 support ( $\gamma$ ) =  $[0, 1]$  and

$$P_\gamma(y) = \left( \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \right) y^{\alpha-1} (1-y)^{\beta-1}$$

$$\beta(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} \quad \Gamma(m) = (m-1)!$$

Comparing  $\textcircled{1}$  with above,

Note:  $\text{Beta}(1, 1) = \text{Unif}[0, 1]$

$$P_\theta |(x_1, \dots, x_n)(\theta | x_1, \dots, x_n) \sim \text{Beta}(s_n+1, f_n+1).$$

$\sum_{i=1}^n x_i$        $n - \sum_{i=1}^n x_i$   
 = # of 1's      = # of 0's

(3) Exponential family of Distributions (both of the above are a special case of this). Read sec. 34.5.1 in book for a discussion of this family.

## Bayesian Bandit Environment

K-armed bandit environment:  $(\mathcal{E}, \mathcal{G}, Q, P)$

the unknown parameter space over k arms.

e.g. in Bernoulli bandit with unknown means  $(\theta_1, \theta_2, \dots, \theta_k)$ ,

$$\mathcal{E} = [0, 1]^k$$

$$\mathcal{G} = \text{Borel-}\sigma\text{-algebra } [0, 1]^k$$

prior dist on

Parameters:  
Example is:

$$Q \sim \text{Beta}(\alpha, \beta)$$

$$\alpha, \beta \geq 1$$

reward dist. for a fixed environment,  
i.e., fixed  $\theta = (\theta_1, \dots, \theta_k)$ .

- $(\mathcal{E}, \mathcal{G})$  the parameter space.

- $Q$  a prob. dist. over  $(\mathcal{E}, \mathcal{G})$ , and is the prior distribution. For any fixed  $v \in \mathcal{E}$ ,

$$P(v \in B) = Q(B).$$

↪ subset of environments (e.g. subset of  $[0, 1]^k$  for Bernoulli bandit).

- Policy for playing arms  $\pi = (\pi_t, t=1, 2, \dots)$  generating actions  $\{A_1, A_2, \dots\}$  and corresponding arm rewards  $(x_1, x_2, \dots)$ , where the conditional dist. of  $A_t$  given  $v, A_1, x_1, A_2, x_2, \dots, A_{t-1}, x_{t-1}$  is:

$$\pi_t(\cdot | a_1, x_1, a_2, x_2, \dots, a_{t-1}, x_{t-1})$$

depends on actions and observed rewards

- $P = (P_{v_i}, v \in \mathcal{E}, i \in \{1, 2, \dots, k\})$  a family of arm reward distributions, for each fixed environment  $v \in \mathcal{E}$  and arm  $i \in \{1, 2, \dots, k\}$ .

The dist. of the reward given  $v, a_1, x_1, \dots, a_{t-1}, x_{t-1}$  is

$$x_t \sim P_{v a_t}(\cdot)$$

depends on environment  $v$  and action  $a_t$

Posterior dist. on (Action, Reward) Sequence:

$(\mathcal{E}, G, Q, P)$  bandit with policy  $\pi$ .

$$P_{v \pi}(a_1, x_1, a_2, x_2, \dots, a_n, x_n)$$

$$= \prod_{t=1}^n \pi_t(a_t | a_1, x_1, \dots, a_{t-1}, x_{t-1}) \cdot P_{v a_t}(x_t)$$

Bayesian Regret

Recall (frequentist / environment-dependent) regret is:

$$R_n(\pi, \nu) = n\mu^* - E \left[ \sum_{i=1}^n x_i \right]$$

↓  
 $\max_{1 \leq i \leq k} M_i$

expectation over randomization  
 due to policy and arm  
 rewards for a FIXED environment  $\nu$ .

$$BR_n(\pi, Q) = \int_{\nu \in \Sigma} R_n(\pi, \nu) dQ(\nu)$$

↓  
 Bayesian Regret  
 $\equiv$  Average over environments w.r.t prior  $Q$ .

$$= \mathbb{E}_{\nu \sim Q} [R_n(\pi, \nu)]$$