# Knowledge Gradient (KG) Model.

$M$: a collection of distinct alternatives.

$\theta_i$: unknown mean. $\rightarrow \theta: (\theta_1, \dots \theta_M)'$.

$\lambda_i$: known variance.

$$\theta \sim N(\mu^0, \Sigma^0) \qquad (1)$$

$x^0, x^1, \dots, x^{N-1}$: a sequence of $N$ sampling decisions.

$x^n$ selects an alternative from $\{1, \dots M\}$.

$\varepsilon^{n+1} \sim N(0, \lambda_{x^n}) \rightarrow$ measurement error.

$$\hat{y}^{n+1} = \theta_{x^n} + \varepsilon^{n+1}, \quad \hat{y}^{n+1} \sim N(\theta_{x^n}, \lambda_{x^n}).$$

$\mathcal{F}^n$: $\sigma$-algebra generated by samples observed by time $n$ and the identities of their originating alternatives. $(x^0, \hat{y}^1, x^1, \hat{y}^2, \dots, x^{n-1}, \hat{y}^n)$.

$E_n \sim E[\cdot | F^n]$.

$\mu^n := E_n[\theta]$. $\Sigma^n := Cov[\theta | F^n]$.

$\Pi$: set of experimental designs, satisfying the sequential requirement.

$\Pi := \{(x^0, \dots, x^{N-1}): x^n \in F^n\}$.

$\pi = (x^0, \dots, x^{N-1}) \rightarrow$ a generic element of $\pi$.

Target: choose a measurement policy maximizing expected reward.

$$\sup_{\pi \in \Pi} E^{\pi}[\max_i \mu_i^N]. \qquad (2)$$

update equations.

$$\mu^{n+1} = \Sigma^{n+1}((\Sigma^n)^{-1}\mu^n + (\lambda_{x^n})^{-1}\hat{y}^{n+1}e_{x^n}). \qquad (3)$$

$$\Sigma^{n+1} = ((\Sigma^n)^{-1} + (\lambda_{x^n})^{-1}e_{x^n}(e_{x^n})')^{-1}. \qquad (4)$$

$x = x^n$

$$\mu^{n+1} = \mu^n + \frac{\hat{y}^{n+1} - \mu_x^n}{\lambda_x + \Sigma_{xx}^n} \Sigma^n e_x. \qquad (5)$$

$$\Sigma^{n+1} = \Sigma^n - \frac{\Sigma^n e_x e_x \Sigma^n}{\lambda_x + \Sigma_{xx}^n} \qquad (6)$$

$$\tilde{\sigma}(\Sigma, x) := \frac{\Sigma e_x}{\sqrt{\lambda_x + \Sigma_{xx}}}. \quad \tilde{\sigma}: \text{ a vector-valued function} \qquad (7)$$

$$\text{Var}[\hat{y}^{n+1} - \mu^n | F^n] = \text{Var}[\theta_x^n + \varepsilon^{n+1} | F^n] = \lambda_x^n + \Sigma_{x^n x^n}^n$$

Define random variables $(Z^n)_{n=1}^N$, $Z^{n+1} := (\hat{y}^{n+1} - \mu^n) / \sqrt{\text{Var}[\hat{y}^{n+1} - \mu^n | F^n]}$.

So $\mu^{n+1} = \mu^n + \tilde{\sigma}(\Sigma^n, x^n) Z^{n+1}$  (8)

$\Sigma^{n+1} = \Sigma^n - \tilde{\sigma}(\Sigma^n, x^n)(\tilde{\sigma}(\Sigma^n, x^n))' = \Sigma^n - \text{cov}[\mu^{n+1} | F^n]$.

$(V^n)_n \sim$ a sequence of value functions. one for each time $n$.

$V^n: S \to \mathbb{R}$.

$V^n(s) := \sup_{\pi \in \Pi} E^\pi [\max_i \mu_i^N | s^n = s]$ for every $s \in S$.

Resulting expression:

$V^N(s) = \max_{x \in [1,\ldots,M]} \mu_x$. for every $s = (\mu, \Sigma) \in S$.

$0 \leq n < N$

$V^n(s) = \max_{x \in [1,\ldots M]} \dot{E}[V^{n+1}(s^{n+1}) | s^n = s, x^n = x]$ for every $s \in S$.  (9)

Define Q factors. $Q^n: S \times \{1, \ldots, M\} \to \mathbb{R}$, as

$Q^n(s, x) := \dot{E}[V^{n+1}(s^{n+1}) | s^n = s, x^n = x]$. for every $s \in S$.  ~~(9)~~

$V^{\pi,n}(s) := E^\pi[V^N(s^N) | s^n = s]$ for every $s \in S$.

$\pi$ is said to be optimal if $V^n(s) = V^{\pi,n}(s)$ for every $s \in S$ and $n \leq N$.  (10)

$\pi^*$:

$X^{\pi^*,n}(s) \in \arg\max_{x \in [1,\ldots,M]} Q^n(s, x)$

for every $s \in S$, $n < N$, and $x \in [1,\ldots,M]$. is optimal.

Define KG Policy $\pi^{KG}$:

$$X^{KG}(S) \in \arg\max_{X} E_n\left[\max_i \mu_i^{n+1} \mid S^n = S, x^n = x\right] - \max_i \mu_i^n,$$

① $N = 1$, KG policy meets the requirements of (1).

② KG policy is the only stationary myopically optimal policy.

$$X^{KG}(S^h) \in \arg\max_{X} \tilde{\sigma}_X(\Sigma^n, x) f\left(\frac{-|\mu_x^n - \max_{i \neq x} \mu_i^n|}{\tilde{\sigma}_X(\Sigma^n, x)}\right).$$

$-\ f(z) := \varphi(z) + z\Phi(z)$

$\varphi$: normal probability density function.

$\Phi$: normal cumulative distribution function.

If $\Sigma^n$ is diagonal, then $\tilde{\sigma}_X(\Sigma^n, x) = \bar{\Sigma}_{xx}^n / \sqrt{\lambda_x + \Sigma_{xx}^n}$

In general case, $\Sigma^n$ is not diagonal.

$$X^{KG}(S^n)$$

$$= \arg\max_{X} E\left[\max_i \mu_i^n + \tilde{\sigma}_x(\Sigma^n, x^n) Z^{n+1} \mid S^n, x^n = x\right] - \max_i \mu_i^n$$

$$= \arg\max_{X} h\left(\mu^n, \tilde{\sigma}(\Sigma^n, x)\right).$$

$h: R^m \times R^m \to R$ defined by $h(a,b) = E\left[\max_i a_i + b_i Z\right] - \max_i a_i.$