Trying to make the right decision when you can't see everything

Yoni Nazarathy

The University of Queensland

Sep 21, 2015



State: $X(t) \longrightarrow X(t+1) \longrightarrow X(t+2) \longrightarrow \cdots$









- Linear Systems: A Success Story
- Stabilising Control of Queues with Hidden Environments
- Reward Observing Restless Bandits

(Deterministic) Linear Systems

(Deterministic) Linear Systems Setup



The Luenberger Observer

$$\hat{X}(t) = A\hat{X}(t) + BU(t) - K_o(\hat{Y}(t) - Y(t))$$
 with $\hat{Y}(t) = C\hat{X}(t)$

$$\hat{X}(t) = A\hat{X}(t) + BU(t) - K_o(\hat{Y}(t) - Y(t))$$
 with $\hat{Y}(t) = C\hat{X}(t)$

$$e(t+1) = X(t+1) - \hat{X}(t+1)$$

= $AX(t) + BU(t) - (A\hat{X}(t) + BU(t) - K_o(C\hat{X}(t) - CX(t)))$
= $AX(t) - A\hat{X}(t) - K_oC(X(t) - \hat{X}(t))$
= $(A - K_oC)e(t)$

$$\hat{X}(t) = A\hat{X}(t) + BU(t) - K_o(\hat{Y}(t) - Y(t))$$
 with $\hat{Y}(t) = C\hat{X}(t)$

$$e(t+1) = X(t+1) - \hat{X}(t+1)$$

= $AX(t) + BU(t) - (A\hat{X}(t) + BU(t) - K_o(C\hat{X}(t) - CX(t)))$
= $AX(t) - A\hat{X}(t) - K_oC(X(t) - \hat{X}(t))$
= $(A - K_oC)e(t)$

Would like $e(t) \rightarrow 0$

$$\hat{X}(t) = A\hat{X}(t) + BU(t) - K_o(\hat{Y}(t) - Y(t))$$
 with $\hat{Y}(t) = C\hat{X}(t)$

$$e(t+1) = X(t+1) - \hat{X}(t+1)$$

= $AX(t) + BU(t) - (A\hat{X}(t) + BU(t) - K_o(C\hat{X}(t) - CX(t)))$
= $AX(t) - A\hat{X}(t) - K_oC(X(t) - \hat{X}(t))$
= $(A - K_oC)e(t)$

Would like $e(t) \rightarrow 0$

Observer Design: Select K_o so that sp $(A - K_o C) < 1$. This can always be done if the pair (A, C) satisfies a rank condition (observability). As controller would like to take $U(t) = -K_f X(t)$ so that,

$$egin{aligned} X(t+1) &= AX(t) + BU(t) \ &= (A-B\,\mathcal{K}_f)X(t) \end{aligned}$$

If (A, B) satisfy a rank condition (controllability) then can choose K_f so as to have arbitrary eigenvalues of $(A - B K_f)$ (e.g. stabilise).

As controller would like to take $U(t) = -K_f X(t)$ so that,

$$egin{aligned} X(t+1) &= AX(t) + BU(t) \ &= (A-B\,\mathcal{K}_f)X(t) \end{aligned}$$

If (A, B) satisfy a rank condition (controllability) then can choose K_f so as to have arbitrary eigenvalues of $(A - B K_f)$ (e.g. stabilise).

But we don't have X(t), so instead use $U(t) = -K_f \hat{X}(t)$.

The "Separation Principle": Can design the observer (K_o) and the controller (K_f) separately to achieve desired behaviour.

Making the "Right Decisions"

First choose K_o for a "good" observer. Now the separation principle allows to focus on finding K_f :

Stability

Choose K_f so that $X(t) \rightarrow 0$

Making the "Right Decisions"

First choose K_o for a "good" observer. Now the separation principle allows to focus on finding K_f :

Stability

Choose K_f so that $X(t) \rightarrow 0$

Quadratic Regulation

$$\min_{U} \quad \sum_{t=1}^{\infty} X(t)' Q X(t) + U(t)' R U(t)$$

Solution to this problem (LQR):

$$U(t)=-K_fX(t),$$

with K_f based on a solution of a Riccati equation.

Moral: For such systems: Don't worry about the fact that "you can't see everything" when you choose an optimal decision

(Deterministic) Linear Systems Summary



Linear Systems with Noise (Stochastic)

(Stochastic) Linear Systems Setup



 $X(t+1) = AX(t) + BU(t) + \xi_{x}(t)$ $Y(t) = CX(t) + \xi_{y}(t)$

The noise components $\xi_x(\cdot)$ and $\xi_y(\cdot)$ are i.i.d. Gaussian

Kalman Filtering

Ignore the control and consider:

$$X(t+1) = AX(t) + \xi_x(t)$$
$$Y(t) = CX(t) + \xi_y(t)$$

Ignore the control and consider:

$$egin{aligned} X(t+1) &= AX(t) + \xi_{x}(t) \ Y(t) &= CX(t) + \xi_{y}(t) \end{aligned}$$

Given $\underline{Y} = (Y(1), \dots, Y(T))$ and X(1), it is straightforward to compute the MMSE $\underline{\hat{X}} = h(\underline{X})$ of $\underline{X} = (X(1), \dots, X(T))$:

 $\operatorname{argmin}_{h}\mathbb{E}\big[||\underline{X} - h(\underline{X})||^{2}\big] = \mathbb{E}\big[X(1), \dots, X(T) \mid Y(1), \dots, Y(T)\big]$

Ignore the control and consider:

$$egin{aligned} X(t+1) &= AX(t) + \xi_{x}(t) \ Y(t) &= CX(t) + \xi_{y}(t) \end{aligned}$$

Given $\underline{Y} = (Y(1), \dots, Y(T))$ and X(1), it is straightforward to compute the MMSE $\underline{\hat{X}} = h(\underline{X})$ of $\underline{X} = (X(1), \dots, X(T))$:

$$\operatorname{argmin}_{h}\mathbb{E}\big[||\underline{X} - h(\underline{X})||^2\big] = \mathbb{E}\big[X(1), \dots, X(T) \mid Y(1), \dots, Y(T)\big]$$

The Kalman Filter is a way to do this recursively (online)

Kalman Filtering (cont)

The (steady state) Kalman filter is a Luenberger observer with parameters optimised for solving the MMSE:

$$\hat{X}(t+1) = A\hat{X}(t) - \mathcal{K}_k \big(C A \hat{X}(t) - Y(t+1) \big).$$

The (steady state) Kalman filter is a Luenberger observer with parameters optimised for solving the MMSE:

$$\hat{X}(t+1) = A\hat{X}(t) - K_k (CA\hat{X}(t) - Y(t+1)).$$

The parameter K_k is calculated as follows:

$$S = \lim_{t \to \infty} Cov \left(X(t+1) - \hat{X}(t+1) \mid X(t), X(t-1), \dots, X(1) \right)$$
$$= A \left(S - SC' \left(CSC' + \Sigma_y \right)^{-1} CS \right) A' + \Sigma_x$$

Now,

$$K_k = SC' (CSC' + \Sigma_y)^{-1}$$

Making The Right Decisions (in a Stochastic Setting)

The Separation Principle generalises:

Certainty Equivalence Principle (holds for such systems, but not always)

In making optimal decisions use $\hat{X}(t) = \mathbb{E}[X(t) \mid \text{observations}]$ as though it was X(t).

Making The Right Decisions (in a Stochastic Setting)

The Separation Principle generalises:

Certainty Equivalence Principle (holds for such systems, but not always)

In making optimal decisions use $\hat{X}(t) = \mathbb{E}[X(t) \mid \text{observations}]$ as though it was X(t).

This allows to solve the LQG problem:

Quadratic Regulation with Gaussian Noise (LQG)

$$\min_{U} \quad \lim_{T \to \infty} \frac{1}{T} \mathbb{E} \big[\sum_{t=1}^{T} X(t)' Q X(t) + U(t)' R U(t) \big]$$

Solution:

$$U(t)=-K_fX(t),$$

with K_f based on a solution of a Riccati equation.

Moral: For such systems don't worry about the fact that "you can't see everything" when you choose an optimal decision.

(Stochastic) Linear Systems Summary



Stabilising Control of Queues with Hidden Server States

(preliminary results from a conference paper with T. Taimre, A. Asanjarani, J. Kuhn, B. Patch and A. Vuorinen)

Who Should Serve What?



Hidden server states: $X(t) = (x_1(t), \ldots, x_M(t))$ with $x_j(t)$ following a 2 state MC:

$$P^{j} = \begin{bmatrix} \bar{p} & p \\ q & \bar{q} \end{bmatrix} = \begin{bmatrix} 1 - \gamma \bar{\rho} & \gamma \bar{\rho} \\ \bar{\gamma} \bar{\rho} & 1 - \bar{\gamma} \bar{\rho} \end{bmatrix}$$

Service rates are $\mu_{i,j}(x_j(t))$. The observations are service successes/failures as well as queue lengths.

The Setup



A Simpler Problem: Who Should Serve the Queue?



A Similar (Simplest) Problem: Go "Safe" or "Bandit"?



- Independent i.i.d. Bernoulli sequences, $\tilde{Y}_1(t)$, $\tilde{Y}_s(t)$, $\tilde{Y}_2(t)$ with means μ_1 , μ_2 and μ_s , respectively, and $\mu_1 < \mu_s < \mu_2$
- X(t) is a 2 state Markov chain
- Observations: $Y(t) = \mathbb{1}\{U(t) = \mathsf{'s'}\}\tilde{Y}_s(t) + \mathbb{1}\{U(t) = \mathsf{'b'}\}\tilde{Y}_{X(t)}$
- Causal policy U(t) should maximize: $\lim_{T\to\infty} \frac{1}{T} \mathbb{E}\left[\sum_{t=0}^{T} Y(t)\right]$

Belief States (POMDP)

Instead of state estimate, $\hat{X}(t)$, keep the conditional distribution (or parameters thereof) of X(t) given history:

 $\omega(t) = \mathbb{P}(X(t) = 2 \mid \text{previous observations and actions})$

Belief States (POMDP)

Instead of state estimate, $\hat{X}(t)$, keep the conditional distribution (or parameters thereof) of X(t) given history:

 $\omega(t) = \mathbb{P}(X(t) = 2 \mid \text{previous observations and actions})$

Now based on "no observation", "observation of failure" or "observation of success", update $\omega(t)$ using one of:

$$\tau(\omega) = \omega \rho + \gamma(1 - \rho)$$

$$\tau_0(\omega) = \frac{\bar{q}\bar{\mu}_2\omega + p\bar{\mu}_1\bar{\omega}}{\bar{\mu}_2\omega + \bar{\mu}_1\bar{\omega}}$$

$$\tau_1(\omega) = \frac{\bar{q}\mu_2\omega + p\mu_1\bar{\omega}}{\mu_2\omega + \mu_1\bar{\omega}}$$

(no observation on bandit)

(observation of failure on bandit)

(observation of success on bandit)

Value Function and Optimal Decision

For simplicity consider the discounted case with factor $\beta \in (0, 1)$:

$$\max \mathbb{E}[\sum_{t=0}^{\infty} \beta^t Y(t)]$$



The Bellman equation for average costs (no discounting) is similar.

Some Structural Results (Average Costs)

Numerical Observations

- The optimal policy is a threshold policy: For ω < ω* choose 's', otherwise choose 'b'.
- 2 In comparison to the "myopic" threshold,

$$\omega^m = (\mu_s - \mu_1)/(\mu_2 - \mu_1)$$
, we have $\omega^* \leq \omega^m$

Some Structural Results (Average Costs)

Numerical Observations

- The optimal policy is a threshold policy: For ω < ω* choose 's', otherwise choose 'b'.
- 2 In comparison to the "myopic" threshold,

$$\omega^m = (\mu_s - \mu_1)/(\mu_2 - \mu_1)$$
, we have $\omega^* \leq \omega^m$



Optimal threshold values with $\mu_s = 0.5$, $\mu_1 = 0.2$, $\mu_2 = 0.8$, $\rho = 0.4$. The attraction region of $\tau_i(\cdot)$ is marked by the vertical dotted lines.

Some Structural Results (Average Costs)

Numerical Observations

- The optimal policy is a threshold policy: For ω < ω* choose 's', otherwise choose 'b'.
- 2 In comparison to the "myopic" threshold,

$$\omega^m = (\mu_s - \mu_1)/(\mu_2 - \mu_1)$$
, we have $\omega^* \leq \omega^m$



Optimal threshold values with $\mu_s = 0.5$, $\mu_1 = 0.2$, $\mu_2 = 0.8$, $\rho = 0.4$.

The attraction region of $\tau_i(\cdot)$ is marked by the vertical dotted lines.

Moral: Myopic is typically not the best. Need to take exploration into account in optimal decision.

Simple Server Selection Summary



Hope: Proving optimality of ω^* threshold (structural result) Less Hope: An explicit ω^* Hope: Rougher structural results for the more general problem

Reward Observing Restless Bandits

(based on some joint work with J. Kuhn)

Several (or Many) Servers to Choose From (No Queues)



- U(t): choose d < M servers, to maximise long term reward
- Server states, X_i(t), evolve independently and are fully observed when the server is selected, otherwise not observed
- Easy to handle server state models are:
 - Two state Markov chains
 - $\bullet\,$ Autoregressive processes of order 1
- Explicit (numerical) solution as a POMDP is hopeless
 ⇒ use index policies on the belief state (approximation)

Reward Observing Restless Bandits Summary



Index Policies

- On The belief state for each channel *i* is ω_i ∈ [0, 1] for 2-state MC channels and (μ_i, ν_i) ∈ ℝ × ℝ₊ for AR channels
- **②** For each channel set an index function $\mathcal{I}_i(\text{belief state}) \to \mathbb{R}$
- **O** Policy: U(t) indicates the *d* channels with the highest index

Index Policies

- On The belief state for each channel *i* is ω_i ∈ [0, 1] for 2-state MC channels and (μ_i, ν_i) ∈ ℝ × ℝ₊ for AR channels
- **②** For each channel set an index function $\mathcal{I}_i(\text{belief state}) \to \mathbb{R}$
- **O** Policy: U(t) indicates the *d* channels with the highest index

The Celebrated Gittins Index Result(s)

If d = 1 and channels freeze when not being selected, then an *optimal* policy is a specific index (Gittins et. al.).

Index Policies

- On The belief state for each channel *i* is ω_i ∈ [0, 1] for 2-state MC channels and (μ_i, ν_i) ∈ ℝ × ℝ₊ for AR channels
- **②** For each channel set an index function $\mathcal{I}_i(\text{belief state}) \to \mathbb{R}$
- **O** Policy: U(t) indicates the *d* channels with the highest index

The Celebrated Gittins Index Result(s)

If d = 1 and channels freeze when not being selected, then an *optimal* policy is a specific index (Gittins et. al.).

Restless Bandits and the Whittle Index

The Reward Observing Restless Bandits problem is a specific case of the *Restless Bandits* problem of P. Whittle (1988). A solution of a relaxed problem is an index policy. In certain cases as $M \rightarrow \infty$ the policy becomes optimal.

EXPLORATION VS. EXPLOITATION WITH PARTIALLY OBSERVABLE AR(1) ARMS

Julia Kuhn

I. Model and Framework

A dynamic decision problem under uncertainty We select k out of d restless reward observing one-armed bandits to play on, such as to maxreward. Rewards are collected and states are observed ONLY if an arm is played. Should we collect new information or ont for the immediate navoff? State processes are Gaussian AR(1).



 $X_i(t) = i \uparrow X_i(t-1) + \varepsilon_i(t)$ where $\alpha \in (0, 1)$ and $c \sim i i d N(0, \sigma^2)$. An application is channel selection in wireless networks.

- The belief states (µ.(t), µ.(t)), i.e. the means and variances condition on the available information, contain all relevant information available at
- At the same time, m(t) and $\nu_i(t)$ quantify the expected gain from exploiting an arm vs. the need for exploring it.

Updating the Belief States

A policy π maps the information available to actions $a_i(t) = 1$ ("play") or a = 0 ("do not play"), such that in total k out of d are played at every time t. With $Y_{\mu\nu} \sim N(\mu, \nu)$.

$$(\mu_i(t+1), \nu_i(t+1)) = \begin{cases} (\varphi \mu_i(t), \varphi^2 \nu_i(t) + \sigma^2), & a_i(t) = 0, \\ (\varphi Y_{\mu(t), \nu(t)}, \sigma^2), & a_i(t) = 1. \end{cases}$$



References

- 2 J. KURN, M. MANDERS and Y. NAZABATHY (2014). Exploration vs. Exploitation with Portially Ob-3 J. GITTINS, K. GLARREROOK and R. WEIRE (2011). Multi-armed Bandit Allocation Indices, 2nd
- Joss Hary & Ann.
 P. WHITTE (1988). Botion handle: Article Allocation in a Changing World. Journal of Arelied

II. Index Policies

$$x \text{ policy is of the form}$$

 $\pi_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \underset{a \sum_{i=1}^{d} \rightarrow a}{\operatorname{arg max}} \left\{ \sum_{i=1}^{d} \gamma(\mu_i, \nu_i) a_i \right\}$

The index function γ maps the belief state of each arm to some priority index.

$$\begin{array}{ll} \mathrm{Mycepic} & \gamma^{M}(\mu, \nu) = \mu \\ \mathrm{Parametric} & \gamma^{\theta}(\mu, \nu) = \mu + \delta \nu, \quad \theta > 0 \\ \mathrm{Whittle} & \gamma^{W}(\mu, \nu) = \mathrm{inf} \left\{ \lambda \mid \pi_{\mathrm{opt}}(\mu, \nu) = 0 \right\} \end{array}$$

Here π_{aut} is the optimal policy for a one-armed bandit problem with subsidu where the decision maker observes and collects the reward when playing, and obtains a subsidy λ



III Whittle Index: Structural Results

The Whittle index policy has been found to be asymptotically optimal in many cases (although no such result is known for our model) but no closed-form expression is known. The associated optimal value function can in principle be found using dynamic programming techniques. We can further prove the following.

The optimal policy for the one-armed bandit problem with subsidy is a threshold policy

The Whittle index $\gamma^{W}(\mu, \nu)$ is monotone nondecreasing in u and u and senerally not constant



 $\gamma^W(\mu, \nu) - \mu.$ $\beta = 0.8, \ \varphi = 0.9, \ \sigma = 2.$

action is "play", below "do not play" IV. Parametric Index: Many-Arms Asymptotic Behaviour

Switching curves: above the curve the optimal

 Intuitively, as d → ∞, k_d/d → ρ, in the long-run the system approaches an equilibrium at which the proportion of arms associated with a certain belief state remains fixed. Then the action chosen for a certain arm is independent of the current helief state of any other arm as there is always the same proportion of arms associated with a certain belief state in the system.

2. We explicitly identify a measure-valued neursion that de-

scribes the many, arms behaviour of the system at conilibrium. Namely, the limiting proportion of arms that have been ob-

served h time steps ago and whose conditional mean falls in

 $m_k(x, t + 1) = \begin{cases} \sum_{k=0}^{\infty} \int_{t_0(t)}^{\infty} \Phi_{z,z^{(0)}}\left(\frac{x}{\tau}\right) m_k(dz, t), & h = 0, \\ m_{k-1}\left(\min\left\{\frac{x}{\tau}, t_{k-1}^*(t)\right\}, t\right), & h \ge 1, \end{cases}$

 $\ell^{*}(t) = \sup \left\{ \ell \mid \sum m_{h} \left(\left\{ \mu \mid \mu + \theta \nu^{(h)} \in [\ell, \infty) \right\}, t \right) = \rho \right\}$

policy activates all arms that are of age h and have conditional

where $\ell_h^*(t) := \ell^*(t) - \theta \nu^{(h)}(t)$ with $\ell^*(t)$ defined by

 $(-\infty, x]$ can be modeled as

3. Based on 1. and 2. we conjecical system at conilibrium is directly related to a one-armed process where dex exceeds a particular threshold $\overline{\ell}$, namely $\overline{\ell} = \ell^*$.

- 1. For large T determine \overline{t} such that $T^{-1} \sum_{i=0}^{T} a_i(t) = \rho$ is achieved for a parametric index policy applied to the one-armed process.
- 2. Use the sample path of Step 1 to obtain an estimate G for the expected average reward of the onearmed system.
- 3. Output $\overline{G}_d := d \overline{G}$ as an approx reward of the multiarmed system with d arms.



Expected average reward $\overline{G}(\theta)$ computed by the algorithm as a function of θ . $\sigma = 2$, $\varphi \in \{0.9, 0.925, 0.95, 0.975\}$, $\rho = 0.4$, $T = 2 \times 10^6$



Comparison of average rewards achieved per arm. θ is found by optimizing (i) the problem with d arms (θ_{d}^{*}), and (ii) the one-armed problem (θ^*), $\varphi = 0.9$, $\sigma = 2$, $\rho = 0.4$, $T = 10^{\circ}$

- P. Whittle, "Optimization Over Time: Dynamic Programming and Stochastic Control", John Wiley & Sons, 1982.
- P. Antsaklis and A. Michel, "A Linear Systems Primer", Birkhauser Boston, 2007.
- P. Whittle, "*Restless Bandits: Activity allocation in a changing world*.", Journal of Applied Probability, 1988.
- Y. Nazarathy, T. Taimre, A. Asanjarani, J. Kuhn, B. Patch and A. Vuorinen, *"The Challenge of Stabilizing Control for Queueing Systems with Unobservable Server States"*, Proceedings of AUCC Confrence, accepted for publication, 2015.
- J. Kuhn, Y. Nazarathy "Wireless Channel Selection with Reward-Observing Restless Multi-armed Bandits", Chapter to appear in "Markov Decision Processes in Practice", Editors: R. Boucherie and N. van Dijk.
- J. Kuhn, M. Mandjes and Y. Nazarathy "Exploration vs. Exploitation with Partially Observable Gaussian Autoregressive Arms", Proceedings of the Valuetools Conference, 2014.
- Y. Nazarathy, L. Rojas-Nandayapa and T. Salisbury "Non-existence of Stabilizing Policies for the Critical Push-Pull Network and Generalizations", Operations Research Letters, 2013.