# Utility Maximization in Cognitive Radio Networks using POMDP Approximate Value Iteration methods

Bharath Keshavamurthy, Nicolò Michelusi

**[NM: title is a bit too long..]**
*Abstract*—

*Index Terms*—**Hidden Markov Model, POMDP, and the PERSEUS Algorithm**

## I. INTRODUCTION

## II. SYSTEM MODEL

**[NM: add figure with system model]**

### A. ~~The Observation~~*Signal* Model

**[NM: This is still not an observation model, since the SU observes only a subset of frequencies..]** We consider a network consisting of $P$ licensed users termed the Primary Users (PUs) and one cognitive radio node termed the Secondary User (SU) equipped with a spectrum sensor. The objective of the SU is to opportunistically access portions of the spectrum left unused by the PUs in order to maximize its own throughput. To this end, the SU should learn how to intelligently access spectrum holes (white-spaces) intending to maximize its throughput while maintaining strict non-interference compliance with incumbent transmissions. The wideband signal ~~observed~~received at the SU receiver at time $n$ is denoted as $y(n)$ and is given by

$$y(n) \;=\; \sum_{p=1}^{P} \sum_{l=0}^{L_p-1} h_p(l) x_p(n-l) + v(n), \quad (1)$$

where $y(n)$ is expressed as a convolution of the signal $x_p(n)$ of the $p$th PU with the channel impulse response $h_p(n)$, and $v(n)$ denotes additive white Gaussian noise (AWGN) with variances $\sigma_v^2$. Eq. (1) can be written in the frequency domain by taking a $K$-point DFT which decomposes the observed wideband signal into $K$ discrete narrow-band components as

$$Y_k(i) \;=\; \sum_{p=1}^{P} H_{p,k}(i) X_{p,k}(i) + V_k(i), \quad (2)$$

where $i \in \{1,2,3,\ldots,T\}$ represents the time index ~~of the observation~~; $k \in \{1,2,3,\ldots,K\}$ represents the index of the components in the frequency domain; $V_k(i) \sim \mathcal{CN}(0,\sigma_V^2)$ represents a circularly symmetric additive complex Gaussian noise sample, i.i.d across channel indices and across time indices; $X_{p,k}(i)$ is the signal of the $p$th PU in the frequency domain, and $H_{p,k}(i)$ is its frequency domain channel. The

noise samples are assumed to be independent of the occupancy state of the channels. We further assume that the $P$ PUs employ an orthogonal access to the spectrum (e.g., OFDMA) so that $X_{p,k}(i)X_{q,k}(i) = 0$, $\forall p \neq q$. Thus, letting $p_k$ be the index of the PU that contributes to the signal in the $k$th spectrum band (possibly, $p_k = 0$ if no PU is transmitting in the $k$th spectrum band), and letting $H_k = H_{p_k,k}$ and $X_k(i) = X_{p_k,k}(i)$, we can rewrite (2) as

$$Y_k(i) \;=\; H_k(i)X_k(i) + V_k(i). \quad (3)$$

Thus, $H_k(i)$ represents the $k$th DFT coefficient of the impulse response $h_{p_k,k}(n)$ of the channel in between the PU operating on the $k$th spectrum band and the SU, in time index $i$; we model it as a zero-mean circularly symmetric complex Gaussian random variable with variance $\sigma_H^2$, $H_k \sim \mathcal{CN}(0,\sigma_H^2)$, ~~independent~~i.i.d. across frequency bands, over time, and independent of the occupancy state of the channels.

### B. PU Spectrum occupancy model

We now introduce the model of PU occupancy over time and across the frequency domain. We model each $X_k(i)$ as

$$X_k(i) = \sqrt{P_{tx}} B_k(i) S_k(i), \quad (4)$$

where $P_{tx}$ is the transmission power of the PUs, $S_k(i)$ is the transmitted symbol, modeled as a constant amplitude signal, $|S_k(i)| = 1$, i.i.d. over time and across frequency bands;[1] $B_k(i) \in \{0,1\}$ is the binary spectrum occupancy variable, with $B_k(i) = 1$ if the $k$th spectrum band is occupied by a PU at time $i$, and $B_k(i) = 0$ otherwise. Therefore, the PU occupancy behavior in the entire wideband spectrum of interest at time index $i$, discretized into narrow-band frequency components, ~~at time index $i$~~ can be modeled as the vector

$$\vec{B}(i) \;=\; [B_1(i),B_2(i),B_3(i),\cdots,B_K(i)]^T \in \{0,1\}^K. \quad (5)$$

PUs join and leave the spectrum at random times. To capture this temporal correlation in the spectrum occupancy dynamics of PUs, we model the spectrum occupancy dynamics as a Markov process: given $\vec{B}(i)$, the spectrum occupancy state at time index $i$, $\vec{B}(i+1)$ is independent of the past, $\vec{B}(j)$, $j < i$; $j,i \in \{1,2,3,\ldots,T\}$, i.e.

$$\mathbb{P}(\vec{B}(i+1)|\vec{B}(j), \; \forall j \leq i) \;=\; \mathbb{P}(\vec{B}(i+1)|\vec{B}(i)). \quad (6)$$

---

[1]In the case where $S_k(i)$ does not have constant amplitude, we may approximate $H_k(i)S_k(i)$ as complex Gaussian with zero mean and variance $\sigma_H^2 \mathbb{E}[|S_k(i)|^2]$, without any modification to the subsequent analysis. ~~$B_k(i) \in \{0,1\}$ is the binary spectrum occupancy variable, with $B_k(i) = 1$ if the $k$th spectrum band is occupied by a PU at time $i$, and $B_k(i) = 0$ otherwise.~~

Additionally, when joining the spectrum pool, PUs occupy a number of adjacent spectrum bands, and may vary their spectrum needs depending on traffic demands, channel conditions, etc. To capture this behavior, we model $\vec{B}(i)$ as having Markovian correlation across sub-bands as, **[NM: use "align" instead of "equation"]**

$$\mathbb{P}(\vec{B}(i+1)|\vec{B}(i)) \tag{7}$$
$$= \mathbb{P}(B_1(i+1)|B_1(i)) \prod_{k=2}^{K} \mathbb{P}(B_k(i+1)|B_k(i), B_{k-1}(i+1)).$$

That is, the spectrum occupancy at time $i+1$ in frequency band $k$, $B_k(i+1)$, depends on the occupancy state of the adjacent spectrum band at the same time, $B_{k-1}(i+1)$, and that of the same spectrum band $k$ in the previous time index $i$, $B_k(i)$.

## C. Spectrum sensing model

In order to detect the available spectrum holes, the SU performs spectrum sensing. However, owing to physical design limitations at the SU's spectrum sensor **[NM: citation?]**, not all sub-bands in the discretized spectrum can be sensed. Therefore, due to limited sensing capabilities, the SU can sense only $\kappa$ out of $K$ spectrum bands at any given time, with $1 \leq \kappa \leq K$. ~~Let $\kappa$ with $1 \leq \kappa \leq K$ be the number of spectrum bands that can be sensed by the SU at any time, capturing the spectrum sensing constraint.~~ Let $\mathcal{K}_i \subseteq \{1, 2, \ldots, K\}$ with $|\mathcal{K}_i| \leq \kappa$ be the set of indices corresponding to the spectrum bands sensed by the SU at time $i$, which is part of our design. Then, we model the emission process of the HMM as

$$\vec{Y}(i) = [Y_k(i)]_{k \in \mathcal{K}_i}, \tag{8}$$

where $Y_k(i)$ is given by (3).

The true states $\vec{B}(i)$ encapsulate the actual occupancy behavior of the PU and the measurements at the SU are noisy observations of these true states which are modeled to be the observed states of a Hidden Markov Model (HMM). Given the spectrum occupancy vector $\vec{B}(i)$ and the set of sensed spectrum bands $\mathcal{K}_i$, the probability density function of $\vec{Y}(i)$ is expressed as

$$f(\vec{Y}(i)|\vec{B}(i), \mathcal{K}_i) = \prod_{k \in \mathcal{K}_i} f(Y_k(i)|B_k(i)), \tag{9}$$

owing to the independence of channels, noise and transmitted symbols across frequency bands. **[NM: it should not be too difficult to incorporate channel correlation across frequency bands]** Moreover,

$$Y_k(i)|B_k(i) \sim \mathcal{CN}(0, \sigma_H^2 P_{tx} B_k(i) + \sigma_V^2). \tag{10}$$

Now, we model the spectrum access scheme of the SU as a Partially Observable Markov Decision Process (POMDP) wherein the goal of the POMDP agent is to devise an optimal sensing and access policy in order to maximize its throughput while maintaining strict non-interference compliance with incumbent transmissions.

## D. The POMDP Agent Model

The agent's limited observational capabilities coupled with its noisy observations and limited sensing capabilities result in an increased level of uncertainty at the agent's end about the occupancy state of the spectrum under consideration and the exact effect of executing an action on the radio environment. The transition model of the underlying MDP as described in the previous subsection (transition model of $\vec{B}(i)$), is denoted by $A$ and is learnt by the agent by interacting with the radio environment. The emission model is denoted as and $B$ **[NM: change notation, you use $B$ already for occupancy vector]** is given by (9), with $f(Y_k(i)|B_k(i))$ given by (10). We model the POMDP as a tuple **[NM: dont change notation, keep using $B$ for your states]** $(\mathcal{B}, \mathcal{A}, \mathcal{Y}, A, B)$ where $\mathcal{B} \equiv \{0,1\}^K$ represents the state space of the underlying MDP with states ~~$\vec{x}$~~$\vec{B}$ given by all possible realizations of the spectrum occupancy vector as given by (3), $\mathcal{A}$ represents the action space of the agent, given by all $\binom{K}{\kappa}$ possible combinations in which the $\kappa$ spectrum bands are chosen to be sensed out of $K$ at any given time index; and $\mathcal{Y}$ represents the observation space of the agent based on the ~~Observation Model~~observation model**[NM: no need for capitals]** outlined in the previous subsection. ~~The run-time or interaction time of the agent is quantized into discrete time indices.~~ **[NM: confusing,, yu have already introduced $i$]** At the beginning of each time index $i$, the agent selects $\kappa$ spectrum bands out of $K$, thus defining the sensing set $\mathcal{K}_i$, performs spectrum sensing on these spectrum bands, observes ~~$\vec{y} \in \mathcal{Y}$~~$\vec{Y}(i) \in \mathcal{Y}$**[NM: keep consistent with the notation]**, and updates its belief of the current spectrum occupancy $\vec{B}(i)$ as ~~$\forall \vec{x}'$~~ as

$$b_a^{\vec{Y}(i)}(\vec{b}') = \mathbb{P}(\vec{B}(i) = \vec{b}'|\vec{Y}(i), \mathcal{K}_i, a\text{[NM : whatisthis?]}, b\text{[NM : none}$$
$$\tag{11}$$

**[NM: fix notation as suggested]** where, $\mathbb{P}(\vec{y}|a, \vec{b})$ is ~~the normalization constant~~given by (9), ~~$\vec{b} \in \mathcal{B}$~~ $b \in \mathcal{B}$ represents the belief ~~vector~~ of the agent prior to the observation $\vec{Y}(i)$, i.e. a probability distribution over all states, in the previous time-step**[NM: previous time step or current one?]**, $b(\vec{x}) \in \vec{b}$ is termed the belief and it represents the degree of certainty assigned to world state $\vec{x} \in \mathcal{X}$ by the belief vector $\vec{b}$. Considering sub-band $k$, the set of available actions to the agent at time index $i$ is given by

$$a_k(i) = \begin{cases} 1, & \text{sense and access sub-band } k, \\ 0, & \text{do nothing with respect to sub-band } k. \end{cases} \tag{12}$$

Based on the number of truly idle sub-bands found accounting for the throughput maximization aspect of our end-goal and a penalty for missed detections accounting for the incumbent non-interference constraint, the reward to the agent is modelled as **[NM: you still did not explain how you decide which one are occupied and which ones are not based on the belief. That would be your spectrum access. Given the belief, how do you make that decision?]** Given the belief $b_i$, the SU uses a threshold based mechanism to access the spectrum: if $\mathbb{P}(B_k(i) = 0|b_i) > \tau$, then the $k$th spectrum band is detected

as idle and it is accessed by the SU; otherwise, it is detected as occupied and is left unused. Therefore, the probability of false alarm and misdetection are given by ..... **[NM: fill the rest]**

$$R(\vec{x}(i), a(i)) = (1 - P_{FA}(i)) + \lambda P_{MD}(i), \qquad (13)$$

where, $P_{FA}(i)$ represents the False Alarm Probability**[NM: no need to capitalize]** across all channels**[NM: how do you deinfe PFA and PMD? How do you compute them? Why do you assume they are the same across all channels?]** at time index $i$, $P_{MD}(i)$ represents the Missed Detection Probability across all channels at time index $i$, and $\lambda < 0$ represent the cost term penalizing the agent for missed detections, i.e. interference with the incumbent. The action policy of the agent $\pi : \mathcal{B} \to \mathcal{A}$ maps the belief vectors $\vec{b} \in \mathcal{B}$ to actions $a \in \mathcal{A}$ and is characterized by a Value Function

$$V^{\pi}(\vec{b}) = \mathbb{E}_{\pi}\Big[\sum_{i=0}^{\infty} \gamma^i R(\vec{b}_i, \ \pi(\vec{b}_i))|\vec{b}_0 = \vec{b}\Big], \qquad (14)$$

where, $0 < \gamma < 1$ is the discount factor, $\pi(\vec{b}_i)$ is the action taken by the agent at time index $i$ under policy $\pi$, and $\vec{b}_0$ is the initial belief vector. The optimal policy $\pi^*$ specifies the optimal action to take at the current time index assuming that the agent behaves optimally at future time indices as well. It is evident from equation (14) that we have an infinite-horizon discounted reward problem formulation and in order to solve for the optimal policy we need to solve the modified**[NM: why modified?]** Bellman equation given as

$$V^*(\vec{b}) = \max_{a \in \mathcal{A}}\Big[\sum_{\vec{x} \in \mathcal{X}} R(\vec{x}, a)b(\vec{x}) + \gamma \sum_{\vec{y} \in \mathcal{Y}} \mathbb{P}(\vec{y}|a, \vec{b})V^*(\vec{b}_a^{\vec{y}})\Big],$$
$$(15)$$

$\forall \vec{b} \in \mathcal{B}$. Given the high dimensionality of the spectrum sensing and access problem, i.e. the number of states of the underlying MDP scales exponentially with the number of sub-bands, solving equation (15) using Exact Value Iteration and Policy Iteration algorithms is computationally infeasible. Additionally, solving for the optimal policy from equation (15) requires prior knowledge about the underlying MDP's transition model. Therefore, in this paper we present a framework to estimate the transition model of the underlying MDP and then utilize this learned model to solve for the optimal policy by employing Randomized Point-Based Value Iteration techniques, namely, the PERSEUS algorithm. **[NM: All the citations should go into the file ref.bib with the specific bibtex format (you can typically get the entry from IEEE Xplore)]**