

A predictive model for the Crowds protocol

Marco Tinacci

May 6, 2015

All the material of this work, included this document, can be found in the public repository <https://github.com/marcotinacci/predictive-crowds>.

1 Introduction

In this work we deal with the anonymity protocol CROWDS [3], we formulate the problem in three different ways extending the canonical set of rules with the possibility to give some inputs to tune the protocol's behavior. In Section 2 we will introduce the protocol and the preliminary concepts that we are going to use, in Section 3 we describe the three models and we show the respective achieved results and finally in Section 4 we conclude discussing possible evolution of the work.

2 Preliminaries

2.1 CROWDS protocol

The crowds protocol has been invented to preserve the anonymity of the sender of a message inside a network of communicating nodes. This approach is useful when corrupted nodes are present in the network and an external attacker may exploit information leaks to trace the origin of the message. For this purpose the protocol intuitively needs to show the sender as the origin of the message not more likely than the other nodes. To do so we use the following randomized algorithm:

1. the sender selects a crowd member at random (possibly itself), and forwards the message to it,
2. the selected router flips a biased coin. With probability $1 - p_f$, where p_f (forwarding probability) is given and globally known, it delivers the message directly to the destination. With probability p_f , it selects a crowd member at random (possibly itself) as the next router in the path, and forwards the message to it, re-encrypted with the appropriate pairwise key. The next router then repeats this step.

2.2 Discrete-time Markov chains

The following definition describe the basic model used in all the three models analyzed in this work.

Definition 2.1. A discrete-time Markov chain (DTMC) \mathcal{D} is defined as follows

$$\mathcal{D} = \langle S, \pi, T \rangle$$

where

- S is a finite set of states,
- $\pi \in \Delta(S)$ is the initial distribution over states,
- $T : S \rightarrow \Delta(S)$ is the transition function.

$\Delta(S)$ is used to denotes the set of all the discrete probability distributions with domain in S , and we denote the number of states $|S|$ as n . A DTMC \mathcal{D} describe an stochastic dynamic system whose dynamics can be summarized with a stochastic matrix $A \in [0, 1]^{n \times n}$. A matrix is stochastic when every column is a discrete probability distribution. We can estimate future states of the system by multiplying the current distribution, a Dirac distribution pointing to the current state, by A many times, one for every step forward. If we denote the current distribution $x_0 \in [0, 1]^n$ we can describe the canonical dynamic of a DTMC with the following equation:

$$x(k+1) = A \cdot x(k)$$

Using this rule we are able to predict how the probability mass is going to spread after a number of discrete steps $T \in \mathbb{N}$, indeed this measure, also known as transient probability, can be computed in the following way

$$x(T) = A^T \cdot x_0$$

The common target of the three models we considered is to tune the dynamics of the DTMC with some input that can affect state probabilities, edge probabilities or probability flow in order to obtain a transient probability that will preserve the anonymity of the sender.

3 Models

3.1 Input control array

Probabilistic behavior We model the protocol as a DTMC, states represent nodes of the network and edges represent transitions of the message from a node to another. Message forwarding is a probabilistic choice among nodes of the network. We can describe the dynamics of the protocol by using

a stochastic matrix $A \in \mathbb{R}^{n \times n}$, with n the number of nodes of the network. Since A is stochastic the following condition holds

$$a_{ij} \in [0, 1], \quad \forall i, j = 1, \dots, n \quad (1)$$

$$\sum_{i=1}^n a_{ij} = 1, \quad \forall j = 1, \dots, n \quad (2)$$

Let a_{ij} be an element of A , we can describe how the dynamics modify the state using the following system

$$x_i(k+1) = \sum_{j=1}^n a_{ij} x_j(k), \quad \forall i = 1, \dots, n \quad (3)$$

where the state $x(k) \in \mathbb{R}^n$ with $k \geq 0$ represents a probability distribution over states, then it holds $\sum_{i=1}^n x_i(k) = 1$ and $x_i(k) \in [0, 1]$ for any $i = 1, \dots, n$.

Predictive control We introduce a vector of parameters $u(k) \in \mathbb{R}^n$ at time k that affect the system dynamics in the following way

$$x(k+1) = Ax(k) + Bu(k) \quad (4)$$

for $k \in \mathbb{N}_0$ and $B \in \mathbb{R}^{n \times n}$. B is a constant real square matrix that describe how parameters affect the state evolution.

Constraints To preserve the probability distribution structure of the state we must ensure that, given any correct $x(k)$, the respective evolution $x(k+1)$ is still correct, that is $x(k+1) \in [0, 1]^n$ and $\sum_{i=1}^n x_i(k+1) = 1$. From the former condition we derive the following constraints

$$-1 \leq -\sum_j^n a_{ij} x_j(k) \leq \sum_j^m b_{ij} u_j(k) \leq 1 - \sum_j^n a_{ij} x_j(k) \leq 1 \quad (5)$$

while, from the latter, we can derive

$$\sum_{i=1}^n \sum_{j=1}^m b_{ij} u_j(k) = 0 \quad (6)$$

Objective We use the canonical performance index (7) for quadratic programming and we customize values in $P, Q \in \mathbb{R}^{n \times n}$ and $R \in \mathbb{R}^{m \times m}$ matrices to suit our target.

$$J(x(0), U) = x(N)^T P x(N) + \sum_{k=0}^{N-1} x(k)^T R x(k) + u(k)^T Q u(k) \quad (7)$$

We are interested in maximizing the probability of being in a healthy node, with a priority for the receiver, and in minimizing the probability of using a corrupted node.

Results We developed a MATLAB script that implement the crowds scenario we defined in this section. We used YALMIP [2] to formulate the problem, generalizing it on the dimension of the network, on the number of corrupted nodes, on the number of parameters employed, on the length of the horizon and on the forward probability. In Figure 1 the predicted state behavior is represented (above) with the respective control actions (below) over time. We can see that healthy nodes (in blue) will be reached with a higher probability with respect to corrupted nodes (in red). The receiver (in green) has the highest probability due to the relatively high forward probability p_f . If we increase the number of control inputs, that in this test is 3, the blu lines tend to move away from the red lines, still remaining close to each other, while red lines tend to move close to zero.

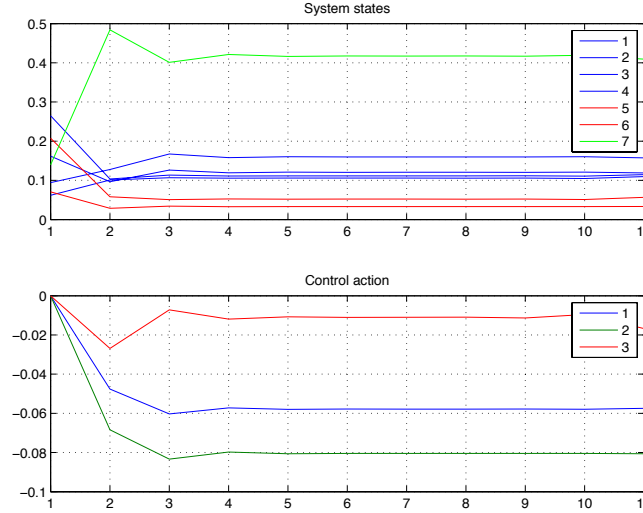


Figure 1: Control of a crowds network of 7 nodes with 3 parameters with model predictive control

The results represented in Figure 1 derive from an experiment with the following settings:

- the network contains 4 healthy nodes, 2 corrupted nodes and 1 receiver;
- the prediction horizon is 10;
- 3 input parameters are used.

To maximize the probability of using a healthy node instead of a corrupted one we set weights in the performance index (7) in the following way

$$P = Q = 10 \cdot I \quad (8)$$

$$r_{ij} = \begin{cases} 1 & \text{if } i = j \text{ and } i \text{ is a healthy node} \\ 10 & \text{if } i = j \text{ and } i \text{ is a corrupted node} \\ 0 & \text{if } i \neq j \end{cases} \quad (9)$$

where r_{ij} is an element of R . The MATLAB code used to run this experiment can be found in the github repository https://github.com/marcotinacci/predictive-crowds/blob/master/mpc_crowds_yalmip.m.

3.2 Input control matrix

We consider a variation of the control applied in the first model from Section 3.1, since every input parameter may affect more than one transition probability and, on the other side, every probability could be affected by more than one input parameter. This kind of formulation does not make clear what control input actually are in the concrete scenario, thus we reformulate the problem in order to have a one-to-one relation between control inputs and transition probabilities. We consider the same stochastic matrix A and constraints (1) and (2) to represent the system dynamics, but we change the control parameters in order to act on the specific transition probabilities.

Predictive control We introduce the possibility to tune probabilities over the edges of the DTMC preserving the distribution constraints (1) and (2). Let $U \in \mathbb{R}^{n \times n}$ be the input matrix that will be used to tune the probabilities.

Now we can extend the system dynamics (3) with parameters given by U in this way

$$x_i(k+1) = \sum_{j=1}^n a_{ij}x_j(k) + u_{ij}(k), \quad \forall i = 1, \dots, n \quad (10)$$

Constraints After applying (10) we need to preserve the distribution constraints, thus, for every k :

$$a_{ij}x_j(k) + u_{ij}(k) \in [0, 1], \quad \forall i, j = 1, \dots, n \quad (11)$$

$$\sum_{j=1}^n a_{ij}x_j(k) + u_{ij}(k) = 1, \quad \forall i = 1, \dots, n \quad (12)$$

Since parameters $U(k)$ can increase or decrease transition probabilities A out of the probability range $[0, 1]$ we need to fix individual constraints to

avoid that. Parameter $u_{ij}(k)$ is a local absolute modification of the probability mass flowing in a specific edge at time k , thus, starting from 11 we can derive the following constraints

$$-1 \leq -a_{ij}x_j(k) \leq u_{ij}(k) \leq 1 - a_{ij}x_j(k) \leq 1 \quad (13)$$

Objective We define a performance index with two features: firstly, minimizing the euclidean distance between the current state $x(k)$ and the ideal state $x_{id} \in [0, 1]^n$, that is the distribution we want to obtain from the system; secondly, penalizing high modifications of input controls with specific costs. We use a weight vector w^x to realize the former feature, while we use a weight matrix W^u for the latter, to measure the effort of every input control on every probabilistic transition. The performance index is defined as follow

$$J(x(0), U) = \sum_{k=0}^{N-1} \sum_{i=1}^n \sum_{j=1}^n w_i^x (x_i(k) - x_{id})^2 + w_{ij}^u u_{ij}^2(k) \quad (14)$$

Thus the problem can be summarize as the following optimization

$$\begin{aligned} \min_U \quad & J(x(0), U) \\ \text{subject to} \quad & (10)(12)(13) \end{aligned} \quad (15)$$

Results Again we used the same environment (MATLAB + YALMIP) to formulate the problem and to solve it with quadratic programming. The result is shown in Figure 2 using the same notation from Figure 1. The outcome is clearly more accurate than the previous experiment as expected, since we are using more parameters. To be specific we parametrize weights on edges instead of states. Moreover, we do not give a limit to the number of parameters that can be used. In Figure 3 are shown the parameters affecting the transition probabilities starting from the first state.

Since the formulation of the problem is different, so it is the parameters setting. In this case we have to chose values for w^x , w^u and x_{id} , while the rest of the values are the same as the previous experiment. The MATLAB source code to generate this experiment can be found at https://github.com/marcotinacci/predictive-crowds/blob/master/crowds_qp_yalmip.m.

$$\begin{aligned} x_{id} &= (0.05, 0.05, 0.05, 0.05, 0, 0, 0.8) \\ w^x &= (1, 1, 1, 1, 1, 1, 1) \\ w^u &= \text{rand}(N) \end{aligned} \quad (16)$$

We randomized w^u with values in $[0, 1]$ in order to distinguish values of the parameters like in Figure 3 that, in case of uniform costs, would assumed the same value as consequence of (15).

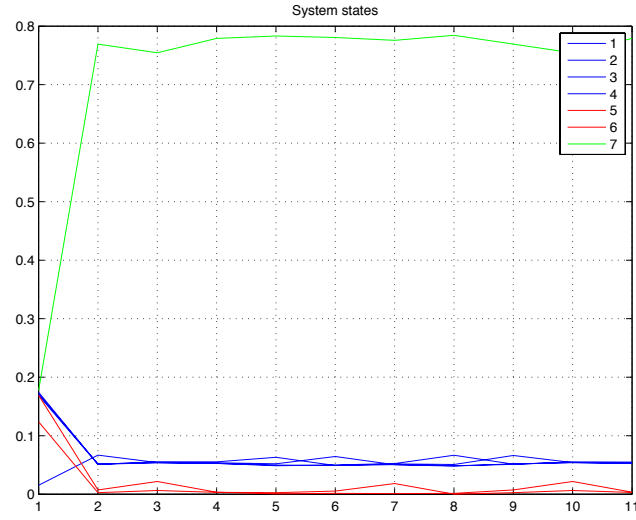


Figure 2: Control of a crowds network of 7 nodes with quadratic programming

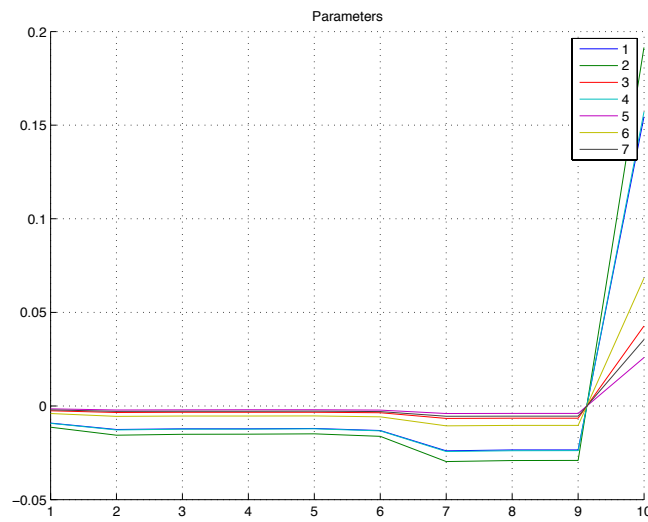


Figure 3: Parameter values on edges starting from the first state with quadratic programming

3.3 Topological control matrix

We formulate the same problem a third time. In this case we want to tune the weights over transitions before the probability mass is affected by them, like in the work of Bartocci et al. [1]. To do so we formulate the problem as follows

$$x(k+1) = (A + U(k))x(k) \quad (17)$$

Again, we have the following straightforward constraints:

- A is a stochastic matrix,
- $x(k)$ is a discrete probability distribution $\forall k = 1 \dots H$, where H is the chosen horizon,
- U is the input matrix that must be zero-sum per columns.

The performance index is again (14). Using MATLAB and YALMIP we implemented the just described problem, the script can be found in the github repository https://github.com/marcotinacci/predictive-crowds/blob/master/crowds_nlp_yalmip.m. A nonlinear programming solver was used until the last iteration over 1000 cycles, but the results were quite accurate, as shown in Figure 4 and Figure 5.

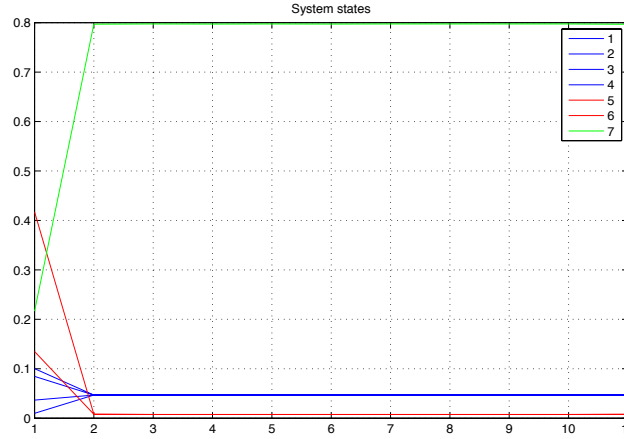


Figure 4: Control of a crowds network of 7 nodes with nonlinear programming

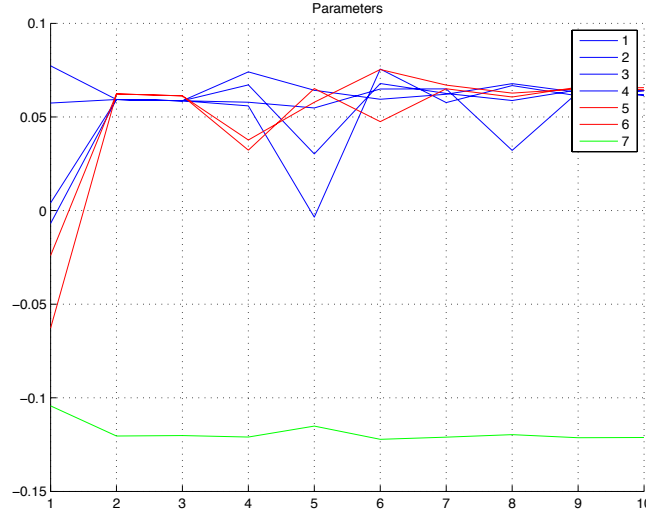


Figure 5: Parameter values on edges starting from the first state with non-linear programming

4 Future work

At the moment the scenario is not very realistic because we need to know which are the corrupted nodes before, but in that case any path avoiding those corrupted nodes would be a good solution. We could use the concept of partially observable model, like hidden Markov models, to model the perception of a leak of information. In this way penalties given to suspected corrupted nodes may be updated according to the perceived observations. Once the system has been modeled as a hidden Markov model, it can be converted into a DTMC and analyzed as we currently do. This kind of transformation is linear in the number of nodes times the number of observations.

In this work only the classical definition of the crowds protocol has been taken into account and the network is assumed to be strongly connected. Sparsity properties of different kind of topologies, block diagonal or band matrix structures for instance, may be exploited to improve analysis performance.

References

- [1] Parosh Aziz Abdulla and K. Rustan M. Leino, editors. *Tools and Algorithms for the Construction and Analysis of Systems - 17th International*

Conference, TACAS 2011, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2011, Saarbrücken, Germany, March 26-April 3, 2011. Proceedings, volume 6605 of *Lecture Notes in Computer Science*. Springer, 2011.

- [2] J. Lfberg. Yalmip : A toolbox for modeling and optimization in MATLAB. In *Proceedings of the CACSD Conference*, Taipei, Taiwan, 2004.
- [3] Michael K. Reiter and Aviel D. Rubin. Crowds: Anonymity for web transactions. *ACM Trans. Inf. Syst. Secur.*, 1(1):66–92, 1998.