

How do Experiences Shape more and less Experienced Cops during Containment of Civil Violence?

Maikel Withagen (S1867733)^a Steven Bosch (S1861948)^a
Robin Kramer (S1970755)^a

^a *University of Groningen, Faculty of Mathematics and Natural Sciences,
Nijenborgh 9, 9747 AG, Groningen*

1 Introduction

1.1 Problem

The civil war in Bosnia and Herzegovina from 1992 to 1995, the Ferguson unrest in 2014, and more recently the civil war in the Ukraine; Civil violence has been an issue in many countries for many years. There is no single cause for these riots and wars. Some were the result of cultural differences, some because of the feeling of being treated unjustly, and others came about due to political reasons. "Each war is as different as the society producing it", and understanding the reasons for these riots to come about is of utmost importance (Ref5Goh) and, in turn, how the authorities attempt to handle these situations.

1.2 State of the Art

Much research has been done to understand how these riots emerge, including simulations of riots using a game theoretic approach (Ref8Goh) and social networks (Ref10Goh) among others. These investigations have shown that the behavior at a macroscopic level, that is, the behavior of an entire group, is the result of behavior on the microscopic level, that is, the individual agents of that group. In a more recent study, Goh and colleagues (Goh, Quek, Tan & Abbass, 2006) also studied how macroscopic behavior emerged using a game theoretic approach in a simulation. The most important issue the authors addressed was how experiences changed the behavior of the individual and how this learning affected macroscopic behavior.

Goh et al. included many different interactions in his simulation, such as a probability that civilians turned to active protesters, how jail time affected rehabilitation of arrested protesters, and how the amount and types of people affected individual decisions. The focus was thus mainly on the civilian and the protesters. What the authors failed to focus on, however, was how experience affected the learning rate, and how cops behaved during riots. It has been shown that many previous experiences ensure a more stable performance (ACT-R, SOAR), meaning that people switch less often between different strategies for the same goal. In other words, more experienced humans are less affected by new experiences compared to less experienced people.

1.3 New Idea

Knowing how experience influences group behavior may be of vital importance. It is very well possible that the 'correct' choice may result in a negative outcome in a particular situation. The

experienced person will be resilient to such an ‘accidental’ outcome, whereas a less experienced person may be affected more severely by that negative outcome. The latter person may therefore change to a less favourable course of action in a future occurrence of said situation. This may result in frailty of the group’s dynamics, ultimately leading to a drop in performance (REF?). Individual experience may thus have a great impact on individual performance and, therefore, group performance, but as yet this has not received much attention.

The current paper focusses on how experience will influence cop behavior. By using a simulation, we were able to compare how experienced cop agents, not experienced cop agents and a mix thereof perform against a group of trained hostile agents. In the followings sections, the simulation and, in turn, data acquisition will be descibed in more detail. Following the results of the different simulation will be presented. The paper concludes with a discussion of the implications of the results, and some shortcomings to this paper.

2 Method

2.1 Simulation model

Three types of agents will be included in this simulation: cops, hostiles and civilians. The goal of the cops is to keep (civil and cop) casualties as low as possible. The hostiles on the other hand have only one goal: to kill as many civilians as possible, to cause mayhem and despair. Civilians have no particular function, but are subject to the actions taken by cops and hostiles. The simulation takes place in a 20x20 2D matrix. In every box of the matrix at most 50 agents can reside. The agents that reside in this box can only see the other agents in that box. This allows for easier computation of interactions with visible agents. In every simulation, around 15.000 agents (± 5000 agents of each group) will be randomly spawned in the matrix.

2.1.1 Goals and Actions

The cops’ goal is to keep casualties as low as possible. To achieve this, every cop must choose one of two actions. The cop can shoot a hostile, which eliminates a treath of future killings, but at the same time there is the risk of killing a civilian. The second action is to to save a civilian, who is then safe from being killed, with a higher risk of being killed by a hostile. Hostiles only have one action: to kill. Depending on the amount of civilians and cops, one group has a larger possibility to get killed by the hostiles. When, for example, only cops and hostiles reside in a box, the possibility that cops are targeted is 100 percent. It depends on the amount of hostiles how many civilians and cops may die.

		Shoot	Save
# Cops > # Hostiles	Many civilians	Ux	Uy
	Few civilians	Ur	Ut
# Cops < # Hostiles	Many civilians	Ua	Ub
	Few civilians	Uc	Ud

Table 1: The different scenarios an agent can encounter and fictional success values for each action.

All the cops, during every epoch, can either save or shoot once. Based on the overall result in that block, a reward is given to the individual agent. This reward can be either positive or negative. If many civilians die, for instance, a negative reward will be given, whilst when many hostiles were killed, a positive reward will be given. The decision that rewards are given based

on group success, as opposed to individual success, is based on the idea that in a crowded situation it is not always clear what the results of individual actions are, but that a general idea of group performance can be perceived.

The reward is used to update the utility of the particular actions in a specific circumstance. These circumstances are shown in Table 1. The cops must thus decide, based on the amount of team members, opponents and civilians what action would be most successful, that is have the highest utility. In the beginning of the simulation, all the utilities are set to equal values, such that no biases may exist. Reward of an action can be calculated as following:

$$R = \frac{Kills + Saves - Deaths}{Kills + Saves + Deaths}$$

in which *Kills* are the amount of killed hostiles, *Saves* the amount of saved civilians, and *Deaths* the amount of killed cops. Based on these calculations the reward will lie between -1 (only deaths) and 1 (only kills and saves). Because hostiles only have one action, no learning is necessary.

As the Table 1 also shows, there is no situation in which the teams are of equal size. It is assumed that the agents have no perfect knowledge of the area and, therefore, decide which team has the overhand. This decision is made according to the following function:

$$\Omega = \frac{n_{cops} + (n_{cops} - n_{hostiles}) * \sigma}{n_{cops} + n_{hostiles}}$$

in which Ω is a value between 0 and 1, n is the amount of agents of a group in that box and σ is a random value between -1 and 1 to make the decision stochastic. If Ω is higher than 0.5, the cops considers to have the overhand.

After each epoch, the agents can decide to move to a neighboring box in the matrix. If the agent decides to move it will go to the place, in which the action with the highest utility can be applied. In other words, the cops will move to where he/she will thrive best. If several

Will – This – Be – Done – With – Some – Function – With – Randomness – Question

2.1.2 Learning

Many reinforcement learning algorithms, such as Q-learning or Temporal-Difference learning, have implemented the fact that newer experiences have a reduced influence (Watkins & Dayan, 1992: Q-Learning, technical notes; REF TD). Q-learning implements a decrease in learning rate by adjusting the learning rate variable λ over time, which allows the function to converge to the true Q-value, that is the utility value. In temporal difference learning this is achieved by storing memories (of actions and results) in time (REF). Over time these memories will decay (less often experienced ones quicker than more experienced ones) allowing a flexible learning of the best actions/strategies. For simplicity reasons, the learning rate variable is used, similar to the Q-learning algorithm.

The utility will be updated according to the following function:

$$U_{as}(t) = U_{as}(t - 1) + \lambda * R_{as}(t)$$

in which $U_x(t)$ is the new utility of an action a in situation s , t is the moment in time, λ is the learning rate factor ranging from 0 to 1, and $R(t)$ is the reward, as calculated before. After each epoch, in which one action could be taken by each agent, the utility functions will be updated.

To initialize the cop agents and set the learning rate variable, random values between 0 and 1 will be taken from a Gaussian distribution. This allows us to create a team with mixed

levels of experience, which increases the realism of the simulation. By setting the skewness of the Gaussian distribution we can manipulate the mean learning rate value to be low and high, which allows the comparison of highly experienced teams (low mean learning rate) with less experienced teams (high mean learning rate).

2.2 Experiment Design

Three simulations were run. In the first simulation, learning rate values were drawn from a normal Gaussian distribution, having a mean value of 0.5. For the other two simulations, the Gaussian distribution were positively and negative skewed, therefore moving the mean learning rate to the left and right respectively.

During the simulation, the following information was kept track of: (1) The number of killed hostiles per box, and the overall number of kills, (2) the number of overall dead cops and per box, and the amount of saved and killed civilians, (3) the number of saved civilians per box and the overall amount. One simulation is considered to be more successful if the mean value for success in the last 10 epochs was significantly higher compared to the other simulations.

3 Results