Vrije Universiteit Brussel

1
2
3

4

5

# Deep Q-learning for the selection of optimal isocratic scouting runs in liquid chromatography

8

Alexander Kensert[1], Gilles Collaerts[1], Kyriakos Efthymiadis[2], Gert Desmet[3] and Deirdre Cabooter[1,*]

11

12

[(1)]University of Leuven (KU Leuven), Department for Pharmaceutical and Pharmacological Sciences, Pharmaceutical Analysis, Herestraat 49, 3000 Leuven, Belgium

[(2)]Vrije Universiteit Brussel, Department of Computer Science, Artificial Intelligence Lab, Pleinlaan 9, 1050 Brussel, Belgium

[(3)]Vrije Universiteit Brussel, Department of Chemical Engineering, Pleinlaan 2, 1050 Brussel, Belgium

18
19

(*) corresponding author:

e-mail: deirdre.cabooter@kuleuven.be

tel.: (+) 32 (0)16.32.34.42

23
24
25
26
27
28

29
30  **ABSTRACT**
31

32  An important challenge in chromatography is the development of adequate separation methods.
33  Accurate retention models can significantly simplify and expedite the development of adequate
34  separation methods for complex mixtures. The purpose of this study was to introduce
35  reinforcement learning to chromatographic method development, by training a double deep $Q$-
36  learning algorithm to select optimal isocratic scouting runs to generate accurate retention
37  models. These scouting runs were fit to the Neue-Kuss retention model, which was then used
38  to predict retention factors both under isocratic and gradient conditions. The quality of these
39  predictions was compared to experimental data points, by computing a mean relative
40  percentage error (MRPE) between the predicted and actual retention factors. By providing the
41  reinforcement learning algorithm with a reward whenever the scouting runs led to accurate
42  retention models and a penalty when the analysis time of a selected scouting run was too high
43  ($> 1h$); it was hypothesized that the reinforcement learning algorithm should by time learn to
44  select good scouting runs for compounds displaying a variety of characteristics. The
45  reinforcement learning algorithm developed in this work was first trained on simulated data,
46  and then evaluated on experimental data for 57 small molecules – each run at 10 different
47  fractions of organic modifier (0.05 to 0.90) and four different linear gradients. The results
48  showed that the MRPE of these retention models (3.77% for isocratic runs and 1.93% for
49  gradient runs), mostly obtained via 3 isocratic scouting runs for each compound, were
50  comparable in performance to retention models obtained by fitting the Neue-Kuss model to all
51  (10) available isocratic datapoints (3.26% for isocratic runs and 4.97% for gradient runs) and
52  retention models obtained via a "chromatographer's selection" of three scouting runs (3.86%
53  for isocratic runs and 6.66% for gradient runs). It was therefore concluded that the
54  reinforcement learning algorithm learned to select optimal scouting runs for retention
55  modeling, by selecting 3 (out of 10) isocratic scouting runs per compound, that were
56  informative enough to successfully capture the retention behavior of each compound.
57
58

# 1 INTRODUCTION

One of the main challenges in (liquid) chromatography is to obtain adequate methods to separate complex samples. Method development is still often done based on a trial-and-error approach, leading to lengthy optimization procedures. Decision-algorithms can be used to help find adequate separations for complex samples (1–8). These algorithms need to be able to predict future runs and select the correct runs to produce the most desirable outcome, which means accurate retention models need to be incorporated into the algorithms. Retention models can generally be divided into two types: quantitative structure retention relationship (QSRR) models, that model retention times of compounds for specific separation conditions based on their molecular descriptors (9–16) and empirical models, that model retention times of compounds based on experimental instrument or system variables, e.g. the fraction of strong eluent in the mobile phase (17–24). The two types of retention models are complementary and usually answer different questions. For instance, QSRR models are commonly developed to improve the identification of known analytes, while empirical retention models are rather developed to improve the separation of complex mixtures of known or unknown compounds. As the present study focuses on improving separation methods for compounds for which molecular structures are not always available, empirical retention models will be used.

An often employed empirical retention model is the Neue-Kuss model (21,26–28):

$$k(\varphi) = k_w(1 + S_2\varphi)^2 \exp\left(-\frac{S_1\varphi}{1 + S_2\varphi}\right) \tag{1}$$

Wherein $\varphi$ is the fraction of strong eluent in the mobile phase, $k(\varphi)$ the retention factor for a specific $\varphi$ and $k_w$, $S_1$ and $S_2$ are retention parameters that need to be fit to experimental retention factors to yield a retention model for a given compound. More specifically, $k_w$ refers to the retention factor of a compound for $\varphi=0$ (in reversed-phase liquid chromatography (RPLC) this corresponds to a purely aqueous mobile phase), $S_1$ refers to the slope and $S_2$ to the curvature of the correlation between k and $\varphi$. Since the Neue-Kuss model requires three parameters for the prediction of the retention factor ($k_w, S_1 \ and \ S_2$), at least three datapoints (*i.e.*, experimental retention factors obtained at three different mobile phase compositions) are needed as input. In RPLC, the polarity of a compound is one of the main factors influencing its retention behaviour. Since compounds can have large differences in polarity, the selection of the scouting runs, *i.e.*, the selection of the mobile phase compositions that need to be selected to obtain adequate retention models, is an important problem as the most informative conditions can strongly differ for each compound. This importance is emphasized making the economic consideration that the number of scouting runs should preferably be kept as low as possible. Thus, there is a clear need to develop algorithms that can select optimal scouting runs leading to accurate retention models.

Decision algorithms for chromatographic applications, such as method development, have been studied since the 1980's. These algorithms, also sometimes referred to as *expert systems,* focused on different aspects of the chromatographic workflow: method selection (28,29), mobile phase selection (30,31), mobile phase/selectivity optimization (32,33), system optimization (34,35) and retention optimization (36–39). All these algorithms aimed at improving the decision-making process for method development and reducing the extent of trial-and-error that is commonly practiced. However, in the late 90's these expert systems lost their popularity due to the fact that the high expectations for these systems were not met, and the computational power available at that time was inadequate.

3

Traditional approaches, like those used in the 80's and 90's, usually take the form of predefined rules and if-else statements, which are directly programmed/handcrafted by the software engineer and experts in the field. The inherent problem of this approach is its great difficulty to generalize to new data or scale to large, diverse datasets. Algorithms like machine learning (40) can, to some extent, solve the problem of generalization by letting the machine program itself (learn its own rules) via highly parameterized and regularized function approximators that are trained on a so-called training dataset. Supervised machine learning algorithms have been around for many decades (41–44) but came to proliferate in the late 2000's when the availability of data in both private and public repositories grew considerably and sufficiently powerful computer hardware became readily available. This family of algorithms only needs to be supplied with input and output pairs, and everything in-between (mapping from input to output) is learned by the algorithm itself. Given an input, *supervised learning* algorithms learn how to minimize the error between their prediction and the target (the associated output for that input). A typical application of supervised learning is classification, in which, given an input, the algorithm outputs the correct class. However, the limitation of supervised learning is the need to label every input, which is often either very time-consuming or impossible to actualize. *Unsupervised learning* deals with situations where no labels are required, with a typical application being clustering. Instead of directly getting a prediction from some input, it clusters the input data and finds patterns in it. Either of these two paradigms is important, but neither fulfills the criteria for intelligent decision making as required in method development. In method development, especially when dealing with novel or unknown compounds, a prediction of the next step to perform is needed, but the effects of the steps are unknown and impossible to plan.

*Reinforcement learning* is a machine learning paradigm that deals with sequential decision making (45,46). In contrast to supervised learning, where each prediction is *independent* of the others, decisions are sequential and *dependent* on each other. This aligns well with method development, which is clearly a sequential procedure (cfr. the trial-and-error approach when method development is carried out in its most simple form). Furthermore, in reinforcement learning, there is no need for input-output pairs as in supervised learning, because no label is needed to teach the algorithm to perform well on a given task (which is of great importance as labels are commonly unavailable in method development). Instead of a label, a scalar reward signal reinforces a desirable behavior; *i.e.* the reinforcement learning algorithm learns a mapping from some input to some desirable output by receiving feedback based on the output. The goal of the reinforcement learning algorithm is not to optimize itself towards accurate label predictions, but instead maximize the cumulative future reward (*i.e.* learn to optimally perform a specific task). Due to its generality, reinforcement learning algorithms have been widely utilized and studied in many applications and areas of research and industry, including robotics (44–49). Besides its many applications, reinforcement learning is possibly most known for its achievements in attaining superhuman level performance in board games, such as Chess and Go (53–55), and in video and computer games such as Atari (56), Starcraft II (57,58) and Dota 2 (59).

In the present study, the possibilities of reinforcement learning are investigated for a basic chromatographic problem: selecting optimal isocratic scouting runs in RPLC to fit retention models. Specifically, the aim is to train a reinforcement learning algorithm to select a minimum number of scouting runs that allow to obtain accurate retention models for a diverse set of compounds. The present study is considered a simple first step in introducing reinforcement learning to the area of chromatography.

## 2 THEORY

Reinforcement Learning is a goal directed paradigm in which a computer program, called the *agent,* is continuously interacting with the *environment*. The environment is the world in which the agent resides and encapsulates the task the agent tries to solve. By performing actions, the agent receives a numerical feedback from the environment based on the quality of the chosen action. Over time, the agent learns to optimize the actions it takes by trying to maximize the future rewards. Typically, reinforcement learning uses a Markov Decision Process (MDP) as its mathematical model. An MDP is a 4-tuple $<S, A, T, R>$, where $S$ is the state space, $A$ the action space, $T(s, a, s') = \Pr(s' | s, a)$ the probability that action $a$ in state $s$ will lead to state $s'$, and R($s, a, s'$) the immediate reward ($r$) received when action ($a$) taken in state ($s$) results in a transition to state ($s'$). The problem of solving an MDP is to find a policy ($\pi$) (i.e., mapping from states to actions) that maximizes the accumulated reward (i.e., the sum of all scalar reward signals). A typical interaction goes as follows: given a state $s$, the agent takes an action $a$, which causes the environment to change its state to a new state $s'$. In addition to the new state $s'$, the environment provides a specific reward signal, indicating how good it is in the new state. Moving across states by performing an action is known as a *transition*, and a sequence of transitions is called a *trajectory*. Specifically, in this study, a trajectory (also referred to as an episode) is defined as a sequence of transitions and ends with a *termination* signal indicating that no more scouting runs are to be run for a given compound.

The state space and action space of the environment define the space in which the reinforcement learning algorithm (or the agent) can perform intelligently. The state of the environment can be represented in many ways, everything from video/image input (a continuous state space), to a vector containing the positions of marks on a tic-tac-toe board (a discrete state space). Similarly, the action space may also be continuous and/or discrete. The environment also sets restrictions for what actions are legal – for instance, is the agent allowed to adjust the percentage of modifier to any percentage? Or is it only allowed to change it between certain discrete, predefined, ranges?

To formalize the learning procedure, the agent tries to learn a policy – what action to take given a certain state. The ultimate goal of the agent is to learn a policy $\pi$ that maximizes the expected cumulative future reward (known as the expected return) over a trajectory $\tau$. This optimal policy denoted as $\pi^*$ is defined as:

$$\pi^* = \arg\max_\pi E_{\tau \sim \pi}[R(\tau)] , \tag{2}$$

$$\text{where } R(\tau) = \sum_{t=0}^{\infty} \gamma^t r_t \tag{3}$$

$E_{\tau \sim \pi}[R(\tau)]$ denotes the expected return of trajectory $\tau$ following policy $\pi$ and $\arg\max_\pi$ specifies the policy maximizing the expected return. $\gamma \in (0,1)$ is known as the *discount factor* that discounts the reward $r$ depending on how far off in the future the reward was obtained. If the discount factor is smaller than one, it is more valuable to get rewards early on rather than later.

For the agent to judge the quality of a given state in an environment, it learns its value through a *value function*. The value function denotes the desirability of being in a state at a given point in time. In the present study, the *optimal action-value function* $Q^*$ will be used – which expresses the expected return if the agent starts in a state s, takes an action a, and then acts according to the optimal policy forever after:

208 $$Q^*(s,a) = \max_\pi E_{\tau \sim \pi}[R(\tau)|s,a] \qquad (4)$$

209

210  $\max_\pi$ indicates the policy yielding the maximum expected return and $\tau \sim \pi$ indicates that the
211  trajectory is sampled according to the policy. Importantly, $Q^*$ can be rewritten as a recursive
212  function, called the Bellman equation for $Q^*$ (60):

213

214 $$Q^*(s,a) = E_{s' \sim P}\left[r(s,a,s') + \gamma \max_{a'} Q^*(s',a')\right] \qquad (5)$$

215

216  $s' \sim P$ indicates that the next state is sampled according to the environment's transition rules,
217  $r(s,a,s')$ is the reward received in the next state $s'$ by taking action $a$ in state $s$, and
218  $\gamma \max_{a'} Q^*(s',a')$ is the discounted maximum expected return in the next state $s'$. The
219  expected return $Q^*$, also referred to as the $Q$-value, of each state-action pair can be computed
220  – allowing for the optimal action $a^*$ in a given state to be obtained via:

221

222 $$a^*(s) = \arg\max_a Q^*(s,a) \qquad (6)$$

223

224  $\arg\max_a$ specifies the action which maximizes $Q^*(s,a)$ and yields the highest $Q$-value.

225

226  The optimal action-value function ($Q^*$) is not known at the start and needs to be solved. A
227  common approach to solve $Q^*$ is via *Q-learning* (61). The *Q*-learning algorithm incorporates
228  the Bellman equation for $Q^*$ to iteratively update $Q$:
229

230 $$Q_{i+1}(s,a) = E_{s'}\left[r + \gamma \max_{a'} Q_i(s',a') |s,a\right] \qquad (7)$$

231  where $Q_i \rightarrow Q^*$ when $i \rightarrow \infty$

232

233  $Q$ is known as the $Q$-function, which approaches or represents $Q^*$ at any given time. $i$ indicates
234  iteration $i$, or transition $i$. By experiencing states and actions, $Q$ can iteratively be updated until
235  it reaches $Q^*$ or is close enough to it. The agent will select actions based on an action selection
236  mechanism. A typical method for action selection is $\epsilon$-*greedy* and is defined as:

237

238 $$a \leftarrow \begin{cases} \arg\max_a Q(s,a) & \text{with probability } 1 - \epsilon \\ \text{random action} & \text{with probability } \epsilon \end{cases}$$

239

240  In essence, the agent will choose to perform a random action with probability $\epsilon$, and will select
241  the optimal action (based on what the agent knows about the environment so far) with a
242  probability 1- $\epsilon$. Through this continual interaction, the agent will eventually learn which
243  actions to perform at each state.

244

245  However, a problem with this tabular approach is that it is impractical, because the state space
246  and action space in most interesting problems are often too large. This approach may not
247  converge to $Q^*$ and lacks generalization capabilities (only able to yield accurate expected
248  returns given actions and states that it has experienced previously). In reinforcement learning,
249  generalization is typically solved by some form of function approximation by representing the
250  value function as a linear function or as a neural network. A major breakthrough in
251  reinforcement learning was the use of the representational power of deep neural networks
252  instead of tabular value functions which led to the creation of Deep Q Networks (DQN) (56)
253  and spurred an entire research field called Deep Reinforcement Learning.

254

255     Specifically, a function approximator, in this case a deep neural network, approximates
256     $Q^*: Q_\theta \approx Q^*$. The parameterized action-value function $Q_\theta$ is thus a neural network, referred to
257     as the $Q$-network, which receives as input a state $s$ and outputs the expected return for each
258     possible action in that state. The parameters $\theta$ of the $Q$-network are estimated by *minimizing*
259     an *objective function* [with respect to $\theta$], which aims to yield an accurate approximation of $Q^*$:
260

261    
$$\min_{\theta_i} \left[ \left( r(s_i, a_i, s_i') + \gamma \max_{a_i'} Q(s_i', a_i', \theta_i^-) \right) - Q(s_i, a_i, \theta_i) \right]^2 \tag{8}$$

262

263     $\min_{\theta_t}$ indicates that $\theta$ is being adjusted at each iteration ($i$) to minimize the objective function.
264     The objective function is the mean squared error between the target (left term between brackets)
265     and the prediction (right term between brackets). For stability purposes, the *target* network has
266     its parameters $\theta^-$ held fixed when minimizing the objective function to not have a continually
267     moving target. The method is referred to as *DQN*, and the minimization procedure referred to
268     as *training* – which is an iterative procedure that iteratively tries to minimize the objective
269     function over the course of thousands of episodes.
270

271     To better understand the concept of reinforcement learning from a higher level, and to tie it to
272     the process of method development, let us consider the agent-environment interaction. This
273     can be thought of as an analyst performing liquid chromatography (LC). If the environment is
274     the LC with its internal physics and chemistry (for example the interactions of analytes with
275     the stationary phase inside the column), then the agent would be the analyst, interacting with
276     the LC. The analyst merely sees a representation of the interaction between a mixture of
277     compounds and the mobile/stationary phase. This representation will take the form of a
278     chromatogram (or retention times), which is the output of the LC. Based on this output, the
279     analyst decides the next action, and as a consequence of this new action, a new representation
280     of a state (new output) is presented to the analyst, and so on. These abstractions are important
281     to understand the roles of the agent and the environment and how they interact, as well as the
282     important differences between them.
283
284

## 3    MATERIALS AND METHODS

### 3.1 Chemicals

The compounds used in this study, together with their structure, Log P and pKa values, stock solvent and supplier, are mentioned in Table S-1 in the Supporting Information. Structures were drawn in Marvin Sketch (v20.9.0, 2020, ChemAxon (www.chemaxon.com)) and values were gathered from PubChem (62). The solvents used to prepare stock solutions and mobile phases were: acetonitrile (ACN, HPLC grade) purchased from Fisher Chemical (Loughborough, UK); formic acid (99%) purchased from Acros Organics (Geel, Belgium); and ammonium formate ($\geq$ 99.995%) purchased from Sigma-Aldrich (Diegem, Belgium). Ultra-pure water ($H_2O$, conductivity = 0.1 $\mu$S/cm, pH 6.00) was produced in the laboratory using a Milli-Q gradient purification system from Millipore (Bedford, MA, USA).

### 3.2 Apparatus

All measurements were done on an Agilent Infinity 1290 system from Agilent Technologies (Waldbronn, Germany) that consisted of the following modules: a quaternary pump (G4204A), an autosampler (G4226A), a thermostatted column compartment (G1316C) and a diode array detector with a 1.0 $\mu$L flow cell (G4212A). OpenLab software (Agilent Technologies) was used to operate the system and acquire and analyse the data. Further data treatment was done in Microsoft Excel. The maximum pressure of the system was 1200 bar. The injection volume was kept constant at 0.50 $\mu$L and the flow rate was 0.500 mL/min. The detector was set at an acquisition rate of 40 Hz to measure the retention times of the compounds and at 80 Hz to measure the column void and extra-column times. The considered wavelengths were 240 nm, 254 nm and 275 nm. The column was placed in an oven with a constant temperature of 30.0 °C. Analyses were performed on an Xbridge $C_8$ column (2.1 x 50 mm; 2.5 $\mu$m) from Waters (Wexford, Ireland). Connections between the autosampler, column and detector were made with nanoViper tubing (inner diameter: 75 $\mu$m, total length: 1000 mm) from Thermo Scientific (Germering, Germany).

### 3.3 Stock and working solutions

Stock solutions of all compounds were prepared in pure organic (ACN) or aqueous solvent, depending on their solubility, as mentioned for each compound in Table S-1, in a concentration up to 10,000 $\mu$g/mL. Working solutions were made by diluting the stock solutions to 20 $\mu$g/mL in a final solvent mixture of 1.000 mL 5/95 ACN/$H_2O$.

### 3.4 Dataset

Retention factors were collected for 82 representative small molecules (see Table S-1), covering a wide range of physiochemical properties, at ten different isocratic strengths wherein the fraction of ACN ($\varphi_{ACN}$) was varied between 0.05 and 0.90 (in 0.10 intervals between 0.10 and 0.90). The aqueous component of the mobile phase was an ammonium formate buffer (brought to pH 3.00 with formic acid). Each compound was injected in triplicate and obtained retention times were averaged. The obtained retention times were converted into retention factors using the following equation:

$$k = \frac{t_r - t_0}{t_0 - t_e} \tag{9}$$

Wherein $t_r$ is the retention time of a particular compound, $t_0$ is the elution time of an unretained marker (thiourea in this study) and $t_e$ is the extra-column dead time, obtained by replacing the column in the system by a zero dead-volume union and injecting thiourea under the same experimental conditions as for $t_r$ and $t_0$.

Some compounds in the dataset were so strongly retained at low $\varphi_{ACN}$ (e.g., 0.05 and 0.10) that it became impractical to record their retention factors in a reasonable amount of time. Therefore, the retention factors at low $\varphi_{ACN}$ for these compounds were estimated by fitting the Neue-Kuss model (Eq. 1) to the available retention factors, and then using the fitted Neue-Kuss model to determine the missing retention factors. Some other compounds were highly polar and did not retain adequately (defined as having a retention factor below 2.5 at 0.05 $\varphi_{ACN}$ or having less than five retention factor values above 0.001) resulting in unreliable data. These compounds were therefore discarded in this study. The final remaining number of compounds was 57. Furthermore, all datapoints were clipped to have a lowest k-value of 0.001, which was necessary because inherent measurement errors caused low retention factor values to be unreliable, with fluctuations between small positive (<0.001) and small negative (> -0.001) values. All experimentally obtained retention factors are shown in Table S-2 in the supporting information. In addition to isocratic retention data, experimental gradient retention data were additionally collected for the 57 compounds. For this purpose, four different gradient profiles were experimentally run (see Table 1).

## 3.5 Compound simulator

To train the *Q*-network, a compound simulator was created to generate sufficient training data for the *Q*-learning algorithm. By fitting the Neue-Kuss model to all the experimentally obtained retention data for the 57 compounds described above, ranges of parameter values ($S_1$, $S_2$ and $k_w$) were deduced, and these ranges and their mutual relations laid the basis to generate retention parameters for simulated compounds (see Figure S-1 and S-2 in the Supporting Information). The simulator generated a simulated compound in five steps:

1. $S_1$ was randomly sampled between $10^{0.9}$ and $10^{1.8}$;
2. $S_2$ was sampled from its relationship with $S_1$: $S_2 = 2.5010 \cdot \log S_1 - 2.0822 + r_1$, where $r_1$ is a random number sampled from a uniform distribution $U \sim (-0.35, 0.35)$;
3. $k_w$ was sampled from its relationship with $S_1$: $k_w = 10^{0.0839 \cdot S_1 + 0.5054 + r_2}$, where $r_2$ is a random number sampled from a uniform distribution $U \sim (-1.2, 1.2)$;
4. the resulting parameter values were input to the Neue-Kuss model to output ten datapoints (retention factors) between $\varphi_{ACN} = 0.05$ and 0.90;
5. random noise sampled from a normal distribution $N \sim (1.0, 0.1)$ was added to each datapoint via multiplication.

The last step was included to mimic the noisy nature of the experimental data. A total of 10,000 compounds were simulated according to the five-step procedure above; each of which would occupy a single episode in the complete training of the deep *Q*-learning algorithm. Although not directly evaluated, 10,000 simulated compounds were considered sufficient to capture the compound space for the specific HPLC setup.

382

383

384

### 3.6 Double deep *Q*-learning algorithm

3.6.1   Environment

The dynamics of the reinforcement learning environment aimed to mimic that of a real chromatographic workflow: sequentially selecting [isocratic] scouting runs for *one* compound at a time (see Figure 1).

For consistency, and to allow for evaluation on experimental data later on, the state and action space of the environment was restricted to the fractions of organic modifier ($\varphi_{ACN}$) that had been experimentally run (namely 0.05 and 0.10 to 0.90 in steps of 0.10; Table S-2). Specifically, a state was defined as an array of 10 elements, each of which was a placeholder for a retention factor value at a certain $\varphi_{ACN}$ (e.g. the first element was a placeholder for the retention factor value at 0.05 $\varphi_{ACN}$, the second element was a placeholder for the retention factor value at 0.10 $\varphi_{ACN}$, and so on). The default value of the elements in the array when a given placeholder was not holding any retention factor value was set to -1, as is the case when no scouting runs had been run (see Fig. 1). The value was set to -1 and not 0 in an attempt to make the *Q*-network more easily distinguish between a non-retention factor value and a low retention factor value. Similar to the state space, the action space was defined to select from the different $\varphi_{ACN}$ that had been run experimentally, with an addition of a *STOP* action, allowing the agent to stop performing scouting runs (see Fig. 1). At *each* step, the agent was given the choice to select *one* specific scouting run, or to stop.

To teach the agent to select desirable scouting runs that can be used to accurately model the retention behavior of a given compound, a reward function was defined. The reward function was defined in such a way as to teach, or *reinforce*, the agent to select optimal scouting runs. After the agent had finished performing scouting runs (and at least three different scouting runs had been selected and run), the resulting retention factors (from those three scouting runs) were used to fit the Neue-Kuss model (see appendix A for more information). Subsequently, the reward function computed a reward based on the inverse of the mean relative error ($MRE^{-1}$) between the target retention factors ($y$) and the retention factors predicted by the Neue-Kuss model ($\hat{y}$):

$$reward = \; MRE^{-1}, \tag{10}$$

$$where \; MRE \; = \frac{1}{10} \sum_{i=1}^{10} \frac{|y_i - \hat{y}_i|}{|y_i|} \tag{11}$$

The targets $y$ are assumed to be the *true* retention factors which we want the predictions $\hat{y}$ to be as close as possible to. In the best-case scenario, true retention factors would be obtained from real experiments. However, as mentioned before, because *Q*-networks need a lot of training data, simulated compounds were used instead; and the relatively few (57) experimental compounds were held-out as a *test set*, and later used to evaluate the agent after it had been trained. Although the training data were randomly sampled using the simulator, which used ranges of Neue-Kuss parameter values ($S_1$, $S_2$ and $k_w$) that were based on those of the compounds in the held-out test set, these parameters were merely roughly deduced from the overall trends of the parameter values of the test compounds and should therefore not undermine the purpose of the held-out test set, which is to evaluate the agent on real 'unseen' data. Importantly, the MRE was always calculated based on all (10) retention factors ($0.05 \rightarrow 0.90$ $\varphi_{ACN}$), as the purpose was to obtain models that predicted well for all (10) $\varphi_{ACN}$.

434 Furthermore, if the agent stopped before three scouting runs had been selected, the MRE
435 evaluation (Eq. 10) was not carried out and the agent was rewarded zero (0). Finally, the agent
436 was penalized (negatively rewarded) in two ways:

438     1. by selecting scouting runs that had already been selected before, penalized -5;
439     2. by selecting scouting runs that resulted in long analysis times (> 1 hour), penalized
440        based on a sigmoidal (s-shaped) function $f_s$:

442 $$f_s(k) = \frac{1}{1 - e^{-(k \cdot 0.00425 - 4.0)}} \times 20 \tag{12}$$

444 where $k$ is the retention factor. The function $f_s$ is adapted in such a way as to have a maximum
445 penalty close to the maximum reward that can be obtained (Eq. 10), and a minimum penalty at
446 zero (0). These penalties were implemented to enable the agent to quickly learn not to select
447 the same scouting run multiple times, and to avoid selecting scouting runs that would result in
448 high analysis times, respectively. The analysis time penalty (Eq. 12) was important because the
449 $Q$-network had to learn to tailor the choice of scouting runs depending on the retention behavior
450 of the compound. For example, if the agent had to select scouting runs for a highly retained
451 compound, it was expected to avoid running scouting runs at low $\varphi_{ACN}$.

453 3.6.2   Agent

455 In this study, as an extension to the deep $Q$-learning algorithm introduced in section 2, a *double*
456 *deep Q-learning algorithm* (63) with *experience replay* (64,65) was implemented.
457 Additionally, an *epsilon-greedy policy* (60) was incorporated for selecting actions (see
458 appendix A, B and C for details on these implementations respectively).

460 Instead of having a single $Q$-network for both the target and the prediction together, the double
461 deep $Q$-learning algorithm utilizes two separate neural networks with two separate $Q$-functions
462 to approximate $Q^*$: a *target* network and a *prediction* network, for the target and the prediction
463 respectively. The minimization procedure is based on the procedure of a single $Q$-network as
464 described in (Eq 8). The target network has fixed parameters $\theta^-$ which are updated every $N$
465 episodes by having the prediction network copying over its parameters $\theta$, while the prediction
466 network is updated every episode via the minimization procedure (Eq. 8).

468 It is of great importance to have a separate network to produce the target, because it counters
469 the issue of overestimating and biasing the target. The overestimation (and bias) of the target
470 naturally occurs for standard $Q$-learning and DQN due to the max operator (cfr. $\max_a$ in eq. 5
471 and 6) both selecting and evaluating the actions. Hasselt et al. have shown how Double deep
472 $Q$-learning produces more accurate estimations of the expected return as well as better policies
473 (63,66).

475 A crucial quality of the minimization procedure (Eq. 8) is that it can be formulated to accept
476 diverse batches of non-correlated examples. Specifically, the minimization step can be done
477 batch-wise, where each batch contains transitions from different trajectories. In other words,
478 each batch may contain tuples (s,a,s') from entirely different scouting runs. To be able to
479 sample such batches, a memory unit called *replay memory* was implemented to store transitions
480 that were experienced by the agent while it was performing scouting runs. These stored
481 transitions were scheduled to be sampled (in batches) after each episode, and to be directly
482 used to train the agent (minimize the objective function). Although we did not test any

483 algorithm without experience replay (replay memory), it is believed that this technique greatly
484 stabilizes and smoothens the training process (64,65). The discount factor $\gamma$ was set to 0.95 as
485 an attempt to give slightly higher weight to early rewards.
486
487
488 **3.7 Evaluation**
489
490 After the double deep $Q$-learning agent was trained, it was evaluated on real, experimental data.
491 Similar to how the agent selected scouting runs for simulated data during training, the agent
492 selected scouting runs for isocratic experimental compounds. The resulting Neue-Kuss model
493 (fitted on the retention factors resulting from the scouting runs selected by the agent) were then
494 used to predict the retention factors for *all* 10 isocratic scouting runs. Additionally, the same
495 Neue-Kuss parameters were also used in the Neue-Kuss equation for gradient elution, which
496 was used to predict retention data under four different gradient conditions (Table 1). The model
497 for gradient elution is defined as follows:
498

499
$$t_{r,gradient} = \frac{1}{S\beta} \frac{(1+S_2\varphi_0)^2 \ln(1+S_1\beta k_w \exp\left(\frac{-S_1\varphi_0}{1+S_2\varphi_0}\right)\left(t_0 - \frac{t_D}{k_0}\right))}{1 - S_2\left(\frac{1+S_2\varphi_0}{S_1}\right)\ln(1+S_1\beta k_w \exp\left(\frac{-S_1\varphi_0}{1+S_2\varphi_0}\right)\left(t_0 - \frac{t_D}{k_0}\right))} + t_0 + t_D \qquad (13)$$

500
501 where $S_1$, $S_2$ and $k_w$ are the same parameters as for the isocratic model (cfr Eq. 1), $t_{r,gradient}$
502 is the gradient retention time, $\varphi_0$ the fraction of strong eluting solvent in the mobile phase at
503 the start of the gradient, $\beta$ the slope of the gradient ($= \Delta\varphi/t_G$), $k_0$ the retention factor at $\varphi_0$, $t_0$
504 the elution time of an unretained compound and $t_D$ the dwell time of the system (the time the
505 mobile phase needs to flow from the pump to the head of the column). The predicted retention
506 factors (for both isocratic and gradient runs) were then compared to *all* experimental data
507 available for that compound by calculating the mean relative percentage error (MRPE) with
508 modification, as follows:
509
510
$$\text{MRPE} = \frac{1}{10} \sum_{i=1}^{10} \frac{|y_i - \hat{y}_i|}{|1 + y_i|} \times 100 \qquad (14)$$

511
512 The (modified) MRPE was used as it was less sensitive to small k-values (<0.1), creating a
513 more stable and robust evaluation metric. To further illustrate how well the resulting model
514 (based on the selected scouting runs) performed, a comparison was made to a model fit on all
515 available isocratic datapoints ($0.05 \rightarrow 0.90$ $\varphi_{ACN}$), a model fit on a random selection of three
516 isocratic scouting runs, as well as a model fit on a chromatographer's selection of three scouting
517 runs, namely $\varphi_{ACN} = 0.1$, 0.5 and 0.8.
518
519
520

## 4   RESULTS AND DISCUSSION

### 4.1 Learning curves: training progression

To illustrate how the agent learned to perform optimal scouting run selection, the *rewards* and *number of selected scouting runs* were collected during the training procedure. Figure 2 illustrates both how the reward and the number of selected actions changed over the first 3500 episodes. In addition, for comparison, the theoretical optimum (highest possible reward) was also calculated and plotted. Notice that the theoretical optimum does not reflect the best possible performance by an agent, but rather illustrates what in theory could be obtained. Specifically, the theoretical optimum was obtained by selecting, for each and every compound separately, the combination of (three) scouting runs (out of the 120 combinations that exist for 3 actions and 10 possible states) that resulted in the highest reward. Due to random action selection, leading to poor retention models and high penalties, the reward given to the agent started off at high negative values (around -5 to -15). Specifically, the low reward as well as the high number of selected scouting runs in the first 100 episodes confirmed that the selected actions were highly random (due to the high $\epsilon$ for the epsilon-greedy policy). Between episode 100 and 200 however, the agent started to select scouting runs non-randomly, which resulted in the agent selecting much fewer scouting runs (around 2 on average). Notice that, although not yielding a retention model due to too few scouting runs ($< 3$), the agent had at this point learned to select better actions than before, because no or fewer penalties were given. After 200 episodes, the agent went from an insufficient number of scouting runs (which resulted in a reward around 0) to selecting enough scouting runs for a retention model, resulting in a positive reward outweighing potential penalties. In other words, after 200 episodes, the agent started to learn how to optimize the selection of scouting runs for retention modelling – selecting around four scouting runs on average, which is on average one scouting run more than the minimum required to fit a three-parameter model. This number was overestimated due to the epsilon-greedy policy ($\epsilon > 0$), which is likely to cause the agent to select scouting runs sub-optimally (i.e. not selecting actions with the greatest *Q*-value). Notice, although not visualized, if only the greedy policy was followed (epsilon = 0), the agent would select closer to three scouting runs on average.

### 4.2 Scouting run selection

After the agent, or the *Q*-network, was optimized or trained to select scouting runs for simulated data, it was evaluated on experimental data. Figure 3 illustrates how the agent selected scouting runs for four representative compounds it had never seen before (see Table S-3 for complete results). For the four compounds considered in Figure 3, $\varphi_{ACN}$ of 0.20 and 0.90 were selected by the trained agent; and for all four compounds, three scouting runs were run in total (which is the minimum requirement). What differed between compounds of different retention behavior, specifically between less retained compounds (e.g. nitrobenzene and o-methylacetophenone) and more retained compounds (e.g. iodobenzene and anthracene), was the selection of the third scouting run ($\varphi_{ACN}$ of 0.05 for the less retained compounds and 0.50 for the more retained compounds), which indicates that the selection of scouting runs was tailored based on the retention behavior of the compounds. Because nitrobenzene and o-methylacetophenone are less retained compounds, they were both selected to be run at as low as $\varphi_{ACN}$=0.05. Intuitively, a scouting run at $\varphi_{ACN}$=0.05 holds important information for a retention model but requires significantly longer analysis time for highly apolar compounds, like iodobenzene and anthracene. For these highly retained compounds, $\varphi_{ACN}$=0.05was avoided and instead $\varphi_{ACN}$=0.50 was selected by the agent. Furthermore, in all cases, the scouting run

571 selection spanned a wide range of selectable $\varphi_{ACN}$, in a relatively evenly spaced-out manner,
572 which intuitively should result in more accurate retention models. Finally, these results
573 illustrate that the penalty given to the agent for selecting scouting runs resulting in high (>> 1
574 hour) analysis times (i.e. high retention factors), made the agent avoid the lower percentages
575 for compounds like iodobenzene and anthracene, but not for compounds like nitrobenzene and
576 o-methylacetophenone, which in regards to this study was highly desirable.
577
578
579 **4.3 Q-values: taking action**
580
581 While Figure 3 illustrates which actions the agent took depending on the retention behavior of
582 the compounds, Figure 4 illustrates how the agent, or the $Q$-network and its $Q$-values, decided
583 which actions to take, given a certain state. It also illustrates how the $Q$-values varied between
584 steps and between compounds with a different retention behavior (specifically, in Figure 4,
585 acetophenone and biphenyl). Interestingly, the penalty given for high analysis times, forced the
586 $Q$-values for $\varphi_{ACN}= 0.05$ and $\varphi_{ACN}= 0.10$ to differ significantly between the two compounds.
587 Specifically, biphenyl, which is more retained than acetophenone, was predicted by the $Q$-
588 network to result in a high penalty if run at $\varphi_{ACN}= 0.05$ or $\varphi_{ACN}= 0.10$ (the analysis time would
589 be too long), and therefore predicted higher $Q$-values at higher $\varphi_{ACN}$. Furthermore, the penalty
590 given when the same $\varphi_{ACN}$ was selected more than once, resulted in a significant lowering of
591 the $Q$-value in the succeeding steps for that $\varphi_{ACN}$.
592
593 Acetophenone had higher $Q$-values on average, suggesting the reward/penalty on average was
594 greater for less retained compounds (like acetophenone) than higher retained compounds (like
595 biphenyl). This suggests that (1) $\varphi_{ACN}= 0.05$ was an important datapoint to accurately model
596 retention behaviors (resulting in a low MRPE) – a datapoint which did not get selected for
597 highly retained compounds due to high penalties; (2) the continual higher penalties for highly
598 retained compounds forced down the $Q$-values.
599
600 Although not directly presented in this study, the *number* of selected scouting runs, as well as
601 the specific selection of scouting runs may differ between agents trained on slightly different
602 data. Specifically, Figure 4 illustrates how several actions have similar $Q$-values given a certain
603 state, which could easily nudge the agent in a different direction. This does not necessarily
604 mean that the agent (or the *double deep Q-learning algorithm*) is not robust, but rather that
605 there are several solutions to the problem – *i.e.* several directions (combinations of scouting
606 runs) that result in good retention models.
607
608 **4.4 Retention model performance: isocratic and gradient prediction errors**
609
610 To get a better understanding of how well the retention models resulting from the scouting runs
611 selected by the agent perform, the prediction accuracy of these retention models was compared
612 to the prediction accuracy of retention models obtained by fitting to (1) all 10 available isocratic
613 data points ($\varphi_{ACN}$ of 0.05 to 0.90); (2) a random selection of three isocratic scouting runs and
614 (3) a chromatographer's selection of three isocratic scouting runs, namely $\varphi_{ACN}= 0.1, 0.5$ and
615 0.8. These four differently obtained retention models were also used to predict gradient
616 retention factors for four different gradients (Table 2).
617
618 Table 2 illustrates the MRPE between the experimentally obtained data and the predictions, for
619 both isocratic and gradient data (see Table S-4 for complete results). The retention model

obtained by the agent (based on (mostly) 3 scouting runs), performed comparably and better to the retention models based on all experimental datapoints and the retention models based on the chromatographer's selection of three scouting runs, for the isocratic and gradient data respectively; and performed significantly better than the retention models based on the random scouting run selection. Importantly, the models resulting from the agent's selection of scouting runs were, compared to retention models based on all datapoints, obtained via seven fewer datapoints; *i.e.* seven fewer scouting runs, saving significant time and costs. These results indicate that the agent successfully learned to optimize the scouting run selection.

The reason why the retention model obtained by the agent had a significantly lower gradient prediction error compared to the retention model fitted on all datapoints and the chromatographer's selection, lies in the high prediction error for the lowly retained compounds. It is speculated that (1) having many datapoints at medium-high $\varphi_{ACN}$ for the lowly retained compounds, where most $\varphi_{ACN}$ have retention factors close to 0, is forcing the estimated parameters ($k_w$, $S_1$ and $S_2$) to model the behavior of the given lowly retained compound favorably in the regions of low $\varphi_{ACN}$, but highly unfavorably in the regions of high $\varphi_{ACN}$, which is more relevant for gradient retention modelling; and (2) not having $\varphi_{ACN}$ of 0.05 (in the case of the chromatographer's selection) misses out on valuable information for gradient predictions.

## 5   Conclusions

In this study, a reinforcement learning algorithm, specifically the double deep $Q$-learning algorithm, was shown to be able to learn to optimally select informative scouting runs in a fully self-taught way. Although only isocratic scouting runs were considered for a specific RPLC setup, these experiments illustrate how the agent learned to tailor the selection of scouting runs for different compounds depending on the retention behavior (mainly defined by its polarity, or $k_w$). The experiments also illustrate how the agent limited the number of selected scouting runs yet still yielding a retention model with low prediction error (MRPE of 3.77% and 1.93% for isocratic and gradient data, respectively). The strategy of the agent was to select relatively evenly spaced-out scouting runs (in terms of $\varphi_{ACN}$), including at least one scouting run at low $\varphi_{ACN}$ (as long as the analysis time was not too long), and at least one scouting run at high $\varphi_{ACN}$. Intuitively, selecting scouting runs as such will cover a greater space, better capturing the complete behavior of the compound. The retention models based on the agent's selection of scouting runs (MRPE of 3.77% for the isocratic data and 1.93% for the gradient data) compared well to both the retention models fitted on all datapoints (MRPE of 3.26% for isocratic data and 4.97% for gradient data) and the chromatographer's selection (3.86% for isocratic data and 6.66% for gradient data); and performed significantly better than retention models based on the random selection of three scouting runs (MRPE of 46.30% for isocratic data 7.60% for gradient data).

Although the double deep $Q$-learning algorithm presented in this study shows potential, it has only been tested on its ability to learn from isocratic data for a specific RPLC setup. This double deep $Q$-learning algorithm is also limited to a discrete, and preferably small, action space, which could potentially be a problem for tasks like selecting scouting runs for gradient elution – where the action space is either larger or continuous. The prospect is to develop a reinforcement learning algorithm that can deal with a more complex/larger state and action space (e.g. Branching Dueling Q-learning (67)) and/or continuous action space (e.g. Twin-Delayed Deep Deterministic Policy Gradient (TD3) or Soft Actor-Critic (SAC) algorithm (68,69)), to be able to perform scouting runs for a mixture of compounds (*i.e.* selecting scouting

670 runs for multiple compounds at the same time) in a more complex setting such as gradient
671 conditions. If such an algorithm can be developed successfully, it would be worthwhile
672 integrating it into normal practice.
673
674
675 **Availability**
676 All code used in this study (except the plots), including the complete implementation of the
677 agent and the environment, can be found at https://github.com/akensert/ddqn-isocratic-
678 scouting-runs. Due to stochasticity in training an artificial neural network, results may differ
679 somewhat from run to run. InChI for each experimental compound evaluated in this study can
680 be found in supplementary table 5 (Table S-5).
681
682 **Conflicts of interest**
683 There are no conflicts of interest.
684
692
693 **Credit Author Statement**
694 **Alexander Kensert:** Conceptualization; Data curation; Formal analysis; Investigation;
695 Methodology; Validation; Visualization; Writing - original draft.
696 **Gilles Collaerts:** Data curation; Formal analysis; Investigation; Validation; Writing – review
697 & editing.
698 **Kyriakos Efthymiadis:** Formal analysis; Investigation; Methodology; Writing – review &
699 editing.
700 **Gert Desmet:** Conceptualization; Funding acquisition; Investigation; Methodology;
701 Supervision; Writing - review & editing.
702 **Deirdre Cabooter:** Conceptualization; Funding acquisition; Formal analysis; Methodology;
703 Investigation; Supervision; Writing - review & editing.

**References**

1. Xu H, Yang L, Freitas MA. A robust linear regression based algorithm for automated evaluation of peptide identifications from shotgun proteomics by use of reversed-phase liquid chromatography retention time. BMC Bioinformatics. 2008 Aug 19;9:347.

2. Fong SS, Rearden P, Kanchagar C, Sassetti C, Trevejo J, Brereton RG. Automated peak detection and matching algorithm for gas chromatography-differential mobility spectrometry. Anal Chem. 2011 Mar 1;83(5):1537–46.

3. Samanipour S, Reid MJ, Bæk K, Thomas KV. Combining a Deconvolution and a Universal Library Search Algorithm for the Nontarget Analysis of Data-Independent Acquisition Mode Liquid Chromatography-High-Resolution Mass Spectrometry Results. Environ Sci Technol. 2018 17;52(8):4694–701.

4. Peters S, Vivó-Truyols G, Marriott PJ, Schoenmakers PJ. Development of an algorithm for peak detection in comprehensive two-dimensional chromatography. J Chromatogr A. 2007 Jul 13;1156(1–2):14–24.

5. Cramer JA, Hammond MH, Loegel TN, Morris RE. Evolving window factor analysis-multivariate curve resolution with automated library matching for enhanced peak deconvolution in gas chromatography-mass spectrometry fuel data. J Chromatogr A. 2018 Dec 21;1581–1582:125–34.

6. López-Ureña S, Torres-Lapasió JR, Donat R, García-Alvarez-Coque MC. Gradient design for liquid chromatography using multi-scale optimization. J Chromatogr A. 2018 Jan 26;1534:32–42.

7. Freier L, von Lieres E. Multi-objective global optimization (MOGO): Algorithm and case study in gradient elution chromatography. Biotechnol J. 2017 Jul;12(7).

8. Woldegebriel M, Gonsalves J, van Asten A, Vivó-Truyols G. Robust Bayesian Algorithm for Targeted Compound Screening in Forensic Toxicology. Anal Chem. 2016 Feb 16;88(4):2421–30.

9. Zapadka M, Kaczmarek M, Kupcewicz B, Dekowski P, Walkowiak A, Kokotkiewicz A, et al. An application of QSRR approach and multiple linear regression method for lipophilicity assessment of flavonoids. J Pharm Biomed Anal. 2019 Feb 5;164:681–9.

10. Taraji M, Haddad PR, Amos RIJ, Talebi M, Szucs R, Dolan JW, et al. Chemometric-assisted method development in hydrophilic interaction liquid chromatography: A review. Anal Chim Acta. 2018 Feb 13;1000:20–40.

11. Bouwmeester R, Martens L, Degroeve S. Comprehensive and Empirical Evaluation of Machine Learning Algorithms for Small Molecule LC Retention Time Prediction. Anal Chem. 2019 05;91(5):3694–703.

12. Maljurić N, Golubović J, Otašević B, Zečević M, Protić A. Quantitative structure - retention relationship modeling of selected antipsychotics and their impurities in green liquid chromatography using cyclodextrin mobile phases. Anal Bioanal Chem. 2018 Apr;410(10):2533–50.

744  13. Ramezani AM, Yousefinejad S, Shahsavar A, Mohajeri A, Absalan G. Quantitative
745      structure-retention relationship for chromatographic behaviour of anthraquinone
746      derivatives through considering organic modifier features in micellar liquid
747      chromatography. J Chromatogr A. 2019 Aug 16;1599:46–54.

748  14. Wen Y, Amos RIJ, Talebi M, Szucs R, Dolan JW, Pohl CA, et al. Retention Index
749      Prediction Using Quantitative Structure-Retention Relationships for Improving
750      Structure Identification in Nontargeted Metabolomics. Anal Chem. 2018
751      07;90(15):9434–40.

752  15. Wen Y, Talebi M, Amos RIJ, Szucs R, Dolan JW, Pohl CA, et al. Retention prediction
753      in reversed phase high performance liquid chromatography using quantitative structure-
754      retention relationships applied to the Hydrophobic Subtraction Model. J Chromatogr A.
755      2018 Mar 16;1541:1–11.

756  16. Zarei K, Atabati M, Ahmadi M. Shuffling cross-validation-bee algorithm as a new
757      descriptor selection method for retention studies of pesticides in biopartitioning micellar
758      chromatography. J Environ Sci Health B. 2017 May 4;52(5):346–52.

759  17. Cabooter D, Clicq D, De Boever F, Lestremau F, Szucs R, Desmet G. A variable
760      column length strategy to expedite method development. Anal Chem. 2011 Feb
761      1;83(3):966–75.

762  18. Tyteca E, De Vos J, Vankova N, Cesla P, Desmet G, Eeltink S. Applicability of linear
763      and nonlinear retention-time models for reversed-phase liquid chromatography
764      separations of small molecules, peptides, and intact proteins. J Sep Sci. 2016
765      Apr;39(7):1249–57.

766  19. Snyder LR, Dolan JW, Gant JR. Gradient elution in high-performance liquid
767      chromatography: I. Theoretical basis for reversed-phase systems. Journal of
768      Chromatography A. 1979 Mar 21;165(1):3–30.

769  20. Dolan JW, Gant JR, Snyder LR. Gradient elution in high-performance liquid
770      chromatography: II. Practical application to reversed-phase systems. Journal of
771      Chromatography A. 1979 Mar 21;165(1):31–58.

772  21. Neue UD. Nonlinear Retention Relationships in Reversed-Phase Chromatography.
773      Chroma. 2006 Jun 1;63(13):S45–53.

774  22. Tyteca E, Desmet G. On the inherent data fitting problems encountered in modeling
775      retention behavior of analytes with dual retention mechanism. J Chromatogr A. 2015 Jul
776      17;1403:81–95.

777  23. Tyteca E, Périat A, Rudaz S, Desmet G, Guillarme D. Retention modeling and method
778      development in hydrophilic interaction chromatography. J Chromatogr A. 2014 Apr
779      11;1337:116–27.

780  24. Tyteca E, Guillarme D, Desmet G. Use of individual retention modeling for gradient
781      optimization in hydrophilic interaction chromatography: separation of nucleobases and
782      nucleosides. J Chromatogr A. 2014 Nov 14;1368:125–31.

783    25.  Česla P, Vaňková N, Křenková J, Fischer J. Comparison of isocratic retention models
784          for hydrophilic interaction liquid chromatographic separation of native and fluorescently
785          labeled oligosaccharides. J Chromatogr A. 2016 Mar 18;1438:179–88.

786    26.  Tyteca E, De Vos J, Vankova N, Cesla P, Desmet G, Eeltink S. Applicability of linear
787          and nonlinear retention-time models for reversed-phase liquid chromatography
788          separations of small molecules, peptides, and intact proteins. J Sep Sci. 2016
789          Apr;39(7):1249–57.

790    27.  Roca LS, Schoemaker SE, Pirok BWJ, Gargano AFG, Schoenmakers PJ. Accurate
791          modelling of the retention behaviour of peptides in gradient-elution hydrophilic
792          interaction liquid chromatography. J Chromatogr A. 2020 Mar 15;1614:460650.

793    28.  Hindriks R, Maris F, Vink J, Peeters A, De Smet M, Massart DL, et al. Expert system
794          for the selection of initial high-performance liquid chromatographic conditions for the
795          analysis of pharmaceuticals. Journal of Chromatography A. 1989 Dec 27;485:255–65.

796    29.  Maris F, Hindriks R, Vink J, Peeters A, Vanden Driessche N, Massart L. Validation of
797          an expert system for the selection of initial high-performance liquid chromatographic
798          conditions for the analysis of basic drugs. Journal of Chromatography A. 1990 May
799          11;506:211–21.

800    30.  De Smet M, Peeters A, Buydens L, Massart DL. Expert system for the selection of high-
801          performance liquid chromatographic methods in pharmaceutical analysis: Validation of
802          the rules for the selection of the mobile phase. Journal of Chromatography A. 1988 Jan
803          1;457:25–42.

804    31.  Szepesi G, Valkó K. Prediction of initial high-performance liquid chromatographic
805          conditions for selectivity optimization in pharmaceutical analysis by an expert system
806          approach. Journal of Chromatography A. 1991 Jan 1;550:87–100.

807    32.  Gros N, Gorenc B. Expert system for the ion chromatographic determination of alkali
808          and alkaline earth metals in mineral waters. Journal of Chromatography A. 1995 Apr
809          21;697(1):31–43.

810    33.  Fell AF, Bridge TP, Williams MH. Design and application of an expert system for
811          mobile phase optimisation in reversed-phase liquid chromatography. Journal of
812          Pharmaceutical and Biomedical Analysis. 1988 Jan 1;6(6):555–64.

813    34.  Schoenmakers PJ, Peeters A, Lynch RJ. Optimization of chromatographic methods by a
814          combination of optimization software and expert systems. Journal of Chromatography
815          A. 1990 May 11;506:169–84.

816    35.  Schoenmakers PJ, Dunand N. Explanations and advice provided by an expert system for
817          system optimization in high-performance liquid chromatography. Journal of
818          Chromatography A. 1989 Dec 27;485:219–36.

819    36.  Smith RM, Burr CM. Retention prediction of analytes in reversed-phase high-
820          performance liquid chromatography based on molecular structure: I. Monosubstituted
821          aromatic compounds. Journal of Chromatography A. 1989 Jan 1;475(2):57–74.

822    37.  Valkó K, Szabó G, Röhricht J, Jemnitz K, Darvas F. Prediction of retention of
823          metabolites in high-performance liquid chromatography by an expert system approach.
824          Journal of Chromatography A. 1989 Dec 27;485:349–63.

825    38.  Hamoir T, Bourguignon B, Massart DL, Hindriks H. Model building for the prediction
826          of initial chromatographic conditions in pharmaceutical analysis using reversed-phase
827          liquid chromatography. Journal of Chromatography A. 1993 Feb 24;633(1):43–56.

828    39.  Fekete J, Morovján G, Csizmadia F, Darvas F. Method development by an expert
829          system advantages and limitations. Journal of Chromatography A. 1994 Feb
830          4;660(1):33–46.

831    40.  Domingos P. A Few Useful Things to Know About Machine Learning. Commun ACM.
832          2012 Oct 1;55:78–87.

833    41.  Cortes C, Vapnik V. Support-vector networks. Mach Learn. 1995 Sep 1;20(3):273–97.

834    42.  Tin Kam Ho. Random decision forests. In: Proceedings of 3rd International Conference
835          on Document Analysis and Recognition. 1995. p. 278–82 vol.1.

836    43.  Rumelhart DE, Hinton GE, Williams RJ. Learning representations by back-propagating
837          errors. Nature. 1986 Oct;323(6088):533–6.

838    44.  Hopfield JJ. Neural networks and physical systems with emergent collective
839          computational abilities. PNAS. 1982 Apr 1;79(8):2554–8.

840    45.  Kaelbling LP, Littman ML, Moore AW. An Introduction to Reinforcement Learning. In:
841          Steels L, editor. The Biology and Technology of Intelligent Autonomous Agents.
842          Berlin, Heidelberg: Springer; 1995. p. 90–127. (NATO ASI Series).

843    46.  Arulkumaran K, Deisenroth MP, Brundage M, Bharath AA. Deep Reinforcement
844          Learning: A Brief Survey. IEEE Signal Processing Magazine. 2017 Nov;34(6):26–38.

845    47.  Ertefaie A, Shortreed S, Chakraborty B. Q-learning residual analysis: application to the
846          effectiveness of sequences of antipsychotic medications for patients with schizophrenia.
847          Stat Med. 2016 15;35(13):2221–34.

848    48.  Shi W, Song S, Wu C, Chen CLP. Multi Pseudo Q-Learning-Based Deterministic Policy
849          Gradient for Tracking Control of Autonomous Underwater Vehicles. IEEE Trans Neural
850          Netw Learn Syst. 2019 Dec;30(12):3534–46.

851    49.  Mnih V, Kavukcuoglu K, Silver D, Graves A, Antonoglou I, Wierstra D, et al. Playing
852          Atari with Deep Reinforcement Learning. arXiv:13125602 [cs] [Internet]. 2013 Dec 19
853          [cited 2020 Jul 29]; Available from: http://arxiv.org/abs/1312.5602

854    50.  Hwangbo J, Lee J, Dosovitskiy A, Bellicoso D, Tsounis V, Koltun V, et al. Learning
855          agile and dynamic motor skills for legged robots. Science Robotics [Internet]. 2019 Jan
856          16 [cited 2020 Jul 29];4(26). Available from:
857          https://robotics.sciencemag.org/content/4/26/eaau5872

858    51.  Liu Z, Abbaszadeh S. Double Q-Learning for Radiation Source Detection. Sensors
859          (Basel). 2019 Feb 24;19(4).

860 52. Komorowski M, Celi LA, Badawi O, Gordon AC, Faisal AA. The Artificial Intelligence
861 Clinician learns optimal treatment strategies for sepsis in intensive care. Nature
862 Medicine. 2018 Nov;24(11):1716–20.

863 53. Silver D, Schrittwieser J, Simonyan K, Antonoglou I, Huang A, Guez A, et al.
864 Mastering the game of Go without human knowledge. Nature. 2017
865 Oct;550(7676):354–9.

866 54. Silver D, Hubert T, Schrittwieser J, Antonoglou I, Lai M, Guez A, et al. A general
867 reinforcement learning algorithm that masters chess, shogi, and Go through self-play.
868 Science. 2018 Dec 7;362(6419):1140–4.

869 55. Silver D, Huang A, Maddison CJ, Guez A, Sifre L, van den Driessche G, et al.
870 Mastering the game of Go with deep neural networks and tree search. Nature. 2016
871 Jan;529(7587):484–9.

872 56. Mnih V, Kavukcuoglu K, Silver D, Rusu AA, Veness J, Bellemare MG, et al. Human-
873 level control through deep reinforcement learning. Nature. 2015 Feb;518(7540):529–33.

874 57. Vinyals O, Ewalds T, Bartunov S, Georgiev P, Vezhnevets AS, Yeo M, et al. StarCraft
875 II: A New Challenge for Reinforcement Learning. arXiv:170804782 [cs] [Internet].
876 2017 Aug 16 [cited 2020 Jul 29]; Available from: http://arxiv.org/abs/1708.04782

877 58. Pang Z-J, Liu R-Z, Meng Z-Y, Zhang Y, Yu Y, Lu T. On Reinforcement Learning for
878 Full-length Game of StarCraft. arXiv:180909095 [cs, stat] [Internet]. 2019 Feb 3 [cited
879 2020 Jul 29]; Available from: http://arxiv.org/abs/1809.09095

880 59. OpenAI, Berner C, Brockman G, Chan B, Cheung V, Dębiak P, et al. Dota 2 with Large
881 Scale Deep Reinforcement Learning. arXiv:191206680 [cs, stat] [Internet]. 2019 Dec 13
882 [cited 2020 Jul 29]; Available from: http://arxiv.org/abs/1912.06680

883 60. Press TM. Reinforcement Learning, Second Edition | The MIT Press [Internet]. The
884 MIT Press; [cited 2020 Aug 17]. Available from:
885 https://mitpress.mit.edu/books/reinforcement-learning-second-edition

886 61. Watkins CJCH, Dayan P. Q-learning. Mach Learn. 1992 May 1;8(3):279–92.

887 62. Kim S, Chen J, Cheng T, Gindulyte A, He J, He S, et al. PubChem 2019 update:
888 improved access to chemical data. Nucleic Acids Res. 2019 Jan 8;47(D1):D1102–9.

889 63. Hasselt H van, Guez A, Silver D. Deep reinforcement learning with double Q-Learning.
890 In: Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence. Phoenix,
891 Arizona: AAAI Press; 2016. p. 2094–2100. (AAAI'16).

892 64. Mnih V, Kavukcuoglu K, Silver D, Graves A, Antonoglou I, Wierstra D, et al. Playing
893 Atari with Deep Reinforcement Learning. arXiv:13125602 [cs] [Internet]. 2013 Dec 19
894 [cited 2020 Aug 14]; Available from: http://arxiv.org/abs/1312.5602

895 65. Zhang S, Sutton RS. A Deeper Look at Experience Replay. arXiv:171201275 [cs]
896 [Internet]. 2018 Apr 30 [cited 2020 Aug 14]; Available from:
897 http://arxiv.org/abs/1712.01275

898   66. Hasselt H van. Double Q-learning. In: Proceedings of the 23rd International Conference
899       on Neural Information Processing Systems - Volume 2. Red Hook, NY, USA: Curran
900       Associates Inc.; 2010. p. 2613–2621. (NIPS'10).

901   67. Arash Tavakoli, Fabio Pardo, Petar Kormushev. Action Branching Architectures for
902       Deep Reinforcement Learning. In: 31st Conference on Neural Information Processing
903       Systems (NIPS 2017). California, USA; 2017.

904   68. Fujimoto S, Hoof H, Meger D. Addressing Function Approximation Error in Actor-
905       Critic Methods. In: International Conference on Machine Learning [Internet]. PMLR;
906       2018 [cited 2020 Oct 27]. p. 1587–96. Available from:
907       http://proceedings.mlr.press/v80/fujimoto18a.html

908   69. Haarnoja T, Zhou A, Abbeel P, Levine S. Soft Actor-Critic: Off-Policy Maximum
909       Entropy Deep Reinforcement Learning with a Stochastic Actor. In: International
910       Conference on Machine Learning [Internet]. PMLR; 2018 [cited 2020 Oct 27]. p. 1861–
911       70. Available from: http://proceedings.mlr.press/v80/haarnoja18b.html

912

913
914

**Figure Captions**

**Figure 1:** A simple schematic illustration of the interaction between agent and environment. A) the environment supplies the agent with a state *s,* which is an array of retention factors (k) where each element corresponds to a $\varphi_{ACN}$; based on this state *s* the agent produces an array of *Q*-values for each action. B) the agent takes the greedy action, namely the action with the highest *Q*-value, resulting in a new state *s'*. For readability, the figure excludes the reward and only covers one single transition (s,a,s').

**Figure 2:** Illustration of how the reward (top; blue) and the number of scouting runs selected (bottom; orange) by the agent increased over the number of episodes. To reduce the noise of these learning curves, a moving average filter was applied to the data (window size of 35). The lighter shaded areas represent a 95% confidence interval of the average lines. The black line in the top plot illustrates the theoretical optimum (the highest possible reward) that can be obtained by selecting the best combination of three scouting runs out of 120 different combinations.

**Figure 3:** Illustration of how the scouting runs were selected sequentially for four representative compounds with a significantly different retention behavior. The estimated Neue-Kuss parameter $k_w$ can be used to assess the retention behavior of the compounds, with a high $k_w$ indicating strong retention, and a low $k_w$ indicating weaker retention. The blue dots indicate the experimental data, the orange triangles pointing to the x-axis indicate the scouting runs selected (sequentially, $1 \rightarrow 3$), and the orange lines indicate the predicted retention model resulting from the scouting runs. MRPE stands for mean relative percentage error. The y-axis has been set to log scale to better visualize the prediction of lower $\varphi_{ACN}$ values.

**Figure 4:** Illustration of how the scouting runs were sequentially (from top to bottom) selected based on the highest Q-value at each step, for two representative compounds. Similar to Figure 3, the estimated Neue-Kuss parameter $k_w$ indicates the retainability of the compounds, with a high $k_w$ indicating strong retention, and a low $k_w$ indicating weaker retention. Orange bars indicate the selected actions.

# Appendix:

# Deep Q-learning for the selection of optimal isocratic scouting runs in liquid chromatography

**A: Retention Model**
The retention model used in this study was the Neue-Kuss model (Eq. 1). The model is obtained by *fitting* the model's parameters ($S_1$, $S_2$ and $k_w$) to the available datapoints for a given compound. The fitting is done via Nelder-Mead optimization (doi: https://doi.org/10.1093/comjnl/7.4.308), which is a type of simplex method.

**B: Double Deep Q-network**
In this study, the two $Q$-networks had a 10-unit input layer (corresponding to the state of the environment), two 1024-unit hidden layers with rectified linear unit activation and a dropout rate of 0.2, and a final 11-unit linear output layer (corresponding to the actions of the agent). While, the prediction network was updated each episode via the minimization procedure, the target network was updated (by copying over the prediction network's parameters) every 64 steps. The batch-size, which is the number of transitions used per minimization step (Eq. 8), was set to 128. To minimize the objective function (Eq. 8), a stochastic gradient descent optimizer was used, with a momentum and an initial learning rate of 0.9 and 0.001 respectively. The learning rate was scheduled to decay for 4096 iterations (episodes) until it reached a minimum learning rate of 0.0001.

Before selecting the hyperparameter values mentioned above, some preliminary testing was done (i.e. the $Q$-network was tested with different hyperparameter values). The number of hidden units was varied between 512 and 2048 (for which 1024 was equal or better to the other values and was therefore selected), the number of layers were varied between 1 and 3 (for which 2 was selected), and the dropout rate was varied between 0.1 and 0.5 (for which 0.2 was selected). Furthermore, the initial learning rate was varied between 0.01 and 0.001 (for which 0.001 was selected), and both rectified linear unit (ReLU) and sigmoid activation was tested (for which ReLU was selected). Finally, discount factor $\gamma$ was varied between 0.0 and 1.0 (for which 0.95 was selected).

**C: Replay Memory/Experience Replay**
Because the batch-size was set to 128, the number of samples sampled from the replay memory unit was also 128. The capacity (maximum number of transitions stored) for the replay memory unit was restricted to 2048. This restriction was added to avoid having the $Q$-network learn from old experiences (transitions). The sampling of transitions from the replay memory unit did not start until it had at least 512 transitions stored; because the $Q$-network (or specifically the prediction network) was updated every episode, having at least 4 times the batch-size of transitions helped to avoid having the $Q$-network train on similar batches multiple times in a row; and although no evidence is presented here, it was thought that this could potentially bias the $Q$-network at the beginning of training and hence slow down learning. As for the hyperparameters of the $Q$-network the batch size was varied between 32 and 512 (for which 128 was selected).

**D: Epsilon-greedy policy**

The epsilon-greedy policy is important for balancing the exploration-exploitation tradeoff. The epsilon-greedy policy in this study was divided into $\varepsilon_{initial}$, $\varepsilon_{decay}$, and $\varepsilon_{minimum}$, which were set to 1.0, 0.99, and 0.1 respectively. This made the agent take mostly random actions (exploring) in the beginning, and increasingly by time started to take actions based on the output of the $Q$-function (the $Q$-values), namely $\arg \max_a Q(s, a)$ (exploiting), The minimum epsilon was kept at 0.1 to force the agent to explore throughout the entire run. As for the hyperparameters of the $Q$-network, the initial epsilon, the decay rate, as well as the minimum decay was varied between 0.1-1.0, 0.99-0.999 and 0.0-1.0 respectively. For which 1.0, 0.99 and 0.1 was selected respectively.