

**Министерство образования и науки Российской Федерации**  
ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ

**«САНКТ-ПЕТЕРБУРГСКИЙ НАЦИОНАЛЬНЫЙ  
ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ ИНФОРМАЦИОННЫХ  
ТЕХНОЛОГИЙ, МЕХАНИКИ И ОПТИКИ»**

**ПОЯСНИТЕЛЬНАЯ ЗАПИСКА  
ВЫПУСКНОЙ КВАЛИФИКАЦИОННОЙ РАБОТЫ**

**«Глубокие самообучающиеся агенты для мультиагентной системы  
маршрутизации»**

Автор: Мухутдинов Дмитрий Вадимович \_\_\_\_\_

Направление подготовки (специальность): 01.03.02 Прикладная математика и  
информатика

Квалификация: Бакалавр

Руководитель: Фильченков А.А., к.ф.-м.н \_\_\_\_\_

**К защите допустить**

Зав. кафедрой Васильев В.Н., докт. техн. наук, проф. \_\_\_\_\_

« \_\_\_\_ » \_\_\_\_\_ 20 \_\_\_\_ г.

Санкт-Петербург, 2017 г.

**Студент** Мухутдинов Д.В. **Группа** М3438 **Кафедра** компьютерных технологий **Факультет** информационных технологий и программирования

**Направленность (профиль), специализация** Математические модели и алгоритмы разработки программного обеспечения

**Консультанты:**

а) Вяткин В.В., PhD, Luleå University of Technology \_\_\_\_\_

Квалификационная работа выполнена с оценкой \_\_\_\_\_

Дата защиты « \_\_\_\_ » \_\_\_\_\_ 20 \_\_\_\_ г.

Секретарь ГЭК \_\_\_\_\_

Листов хранения \_\_\_\_\_

Демонстрационных материалов/Чертежей хранения \_\_\_\_\_

## ОГЛАВЛЕНИЕ

ВВЕДЕНИЕ.....	5
1. Обзор предметной области .....	8
1.1. Постановка задачи.....	8
1.2. Существующие алгоритмы маршрутизации .....	8
1.2.1. Дистанционно-векторные алгоритмы.....	8
1.2.2. Алгоритмы состояния канала связи .....	10
1.2.3. Q-routing.....	11
1.2.4. Другие подходы .....	12
1.2.5. Выводы по обзору существующих решений.....	12
1.3. Постановка задачи в терминах обучения с подкреплением ..	13
1.3.1. Термины и понятия .....	13
1.4. Обзор методов обучения нейросетей с подкреплением .....	14
Выводы по главе 1.....	14
ЗАКЛЮЧЕНИЕ.....	15
СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ.....	16

## ВВЕДЕНИЕ

Задача пакетной маршрутизации - это задача поиска кратчайшего пути в графе в условиях, когда за каждый узел графа отвечает отдельный вычислительный процесс. Это означает, что каждый отдельный узел графа должен принять решение о том, какому из соседей следует отправить очередной пакет, чтобы тот достиг пункта назначения как можно быстрее.

Задача пакетной маршрутизации впервые обрела актуальность с появлением компьютерных сетей. Первые алгоритмы сетевой маршрутизации появились в процессе разработки сети ARPANet. Именно тогда были изобретены такие подходы к пакетной маршрутизации, как distance-vector[1] и link-state[2], которые и по сей день лежат в основе таких стандартных и широко применяемых алгоритмов сетевой маршрутизации, как Routing Information Protocol (RIP)[3] и Open Shortest Path First (OSPF)[4].

Оба этих подхода основаны на идее вычисления кратчайшего пути между текущим узлом сети и пунктом назначения пакета. Однако существуют и другие подходы к решению задачи маршрутизации, основанные на идее обучения с подкреплением (reinforcement learning). Первым таким алгоритмом стал алгоритм Q-routing[5], основанный на методе Q-learning[6]. Этот алгоритм, как и его модификации[7, 8], благодаря обучению с подкреплением оказался способен лучше адаптироваться к изменениям в интенсивности сетевого трафика, чем алгоритмы, основанные на вычислении кратчайшего пути, но из-за использования большего количества служебных сообщений применение таких алгоритмов в реальных сетях ограничено.

Но стоит отметить, что задача маршрутизации встречается не только в компьютерных сетях, но также и в более сложных средах, таких как киберфизические системы. Примером такой задачи, частично решаемой с помощью алгоритмов маршрутизации, является управление конвейерной системой (в частности, системой распределения багажа в аэропорту). Задача управления такой системой в большинстве случаев решается с помощью централизованных алгоритмов, и децентрализованное решение (основанное на алгоритма Беллмана-Форда) было пред-

ложено совсем недавно[9]. Такая задача отличается от задачи сетевого роутинга, с одной стороны, тем, что служебные сообщения и целевые сообщения являются разными сущностями (“целевое сообщение” — это чемодан, а служебные сообщения являются цифровыми), и служебные сообщения передаются мгновенно по сравнению с целевыми. С другой стороны, состояние каждого конвейера и всей системы в целом задается большим количеством параметров — такими как скорости конвейеров, положение, количество и масса чемоданов на каждом конвейере, и так далее. С третьей стороны, желательно, чтобы система оптимизировала не только скорость доставки чемоданов до точки назначения, но и, к примеру, собственное энергопотребление.

Первое обстоятельство снимает технические ограничения на количество служебных сообщений, что позволяет в полной мере применять алгоритмы на основе обучения с подкреплением. Второе и третье обстоятельства усложняют написание оптимального детерминированного алгоритма и наталкивают на идею реализации приближенного решения, например, с использованием нейронных сетей.

В настоящее время в области обучения с подкреплением с применением нейронных сетей достигнуты впечатляющие успехи. Всплеск активности в этой области произошел после выхода статьи команды Google DeepMind об обучении глубокой сверточной нейронной сети игре на консоли Atari 2600[10]. Применению полученной модели к различным задачам в различных условиях посвящено множество исследований, часть из них посвящена проблеме мультиагентного обучения с подкреплением (multi-agent reinforcement learning). Так как задачу маршрутизации можно сформулировать как задачу мультиагентного обучения с подкреплением, имеет смысл применить данные наработки для ее решения.

В данной работе будет предложен алгоритм маршрутизации, основанный на методе Q-routing, но использующий нейронную сеть в качестве обучающегося агента. На данный момент не существует алгоритма маршрутизации, построенного по такому принципу.

В главе 1 будет сформулирована обобщенная постановка задачи маршрутизации в терминах мультиагентного обучения с подкреплением.

ем. Будут рассмотрены существующие алгоритмы маршрутизации, их сильные и слабые стороны. Также будут рассмотрены существующие методы обучения нейронных сетей с подкреплением, в том числе в мультиагентном случае.

В главе 2 будет рассмотрен предложенный алгоритм и обоснованы решения, принятые в ходе его разработки.

В главе 3 будут приведены экспериментальные результаты работы алгоритма для задач пакетной маршрутизации в компьютерной сети и управления системой багажных конвейеров. Будут приведены результаты работы в условиях неравномерной нагрузки на сеть (или конвейерную систему) и изменения топологии сети. Также будет проведено сравнение с существующими алгоритмами маршрутизации.

## ГЛАВА 1. ОБЗОР ПРЕДМЕТНОЙ ОБЛАСТИ

### 1.1. Постановка задачи

Пусть задана сеть в виде графа:  $G = (V, E)$ . Каждому узлу и каждому ребру в сети приписывается некоторое *состояние*: узел  $v$  имеет состояние  $s_v \in S_V$ , ребро  $e$  имеет состояние  $s_e \in S_E$ . Состояние всей сети, таким образом, описывается как  $s = (s_{v_1}, \dots, s_{v_n}, s_{e_1}, \dots, s_{e_m}) \in S$ .

На узел сети  $s \in V$  приходят пакеты  $p = (d, prop)$ , где  $d \in V$  - узел назначения пакета, а  $prop$  - произвольная дополнительная информация о нем. Обозначим множество всевозможных свойств пакета как  $Prop$  ( $prop \in Prop$ ), а множество пакетов в целом, таким образом, как  $P = (V, Prop)$ . Также заданы некоторые функции *стоимости* прохождения данного пакета через ребра и узлы сети, зависящие от их состояний:  $R_v : V \times S_V \times P \rightarrow \mathbb{R}^+$  - стоимость прохождения пакета через узел,  $R_e : E \times S_E \times P \rightarrow \mathbb{R}^+$ . Мы явно указываем, что стоимости неотрицательны, чтобы избежать таких постановок задач, в которых существуют пути стоимости  $-\infty$  (циклы отрицательной стоимости).

Задача пакетной маршрутизации заключается в том, чтобы определить, какому из соседей  $n \in \{v | (s, v) \in E\}$  узел  $s$  должен перенаправить пакет  $p = (d, prop)$ , чтобы ожидаемая стоимость пути от  $s$  до  $d$  была минимальной.

### 1.2. Существующие алгоритмы маршрутизации

Почти все существующие алгоритмы маршрутизации были разработаны для маршрутизации пакетов в компьютерных сетях. Большинство алгоритмов маршрутизации в компьютерных сетях, применяемых на практике, относятся к одному из двух семейств алгоритмов — дистанционно-векторные (distance-vector)[1] или состояния каналов связи (link-state)[2]. Концептуально все алгоритмы внутри каждого из этих семейств одинаковы, и различаются только техническими деталями реализации, обусловленными спецификой конкретной узкой сферы применения. Поэтому мы не будем рассматривать алгоритмы каждого семейства по отдельности, а рассмотрим только концепции в целом.

#### 1.2.1. Дистанционно-векторные алгоритмы

Идея дистанционно-векторных алгоритмов (distance-vector algorithms) заключается в следующем:

- Каждый маршрутизатор  $s$  в сети хранит таблицу, в которой для каждого другого узла сети  $d$  хранится следующая информация:
  - Предполагаемое кратчайшее расстояние от  $s$  до  $d$
  - Сосед  $n$ , которому нужно отправить пакет, чтобы пакет прошел по кратчайшему пути до узла  $d$ .
- Периодически каждый маршрутизатор рассылает свою версию таблицы кратчайших расстояний всем своим соседям
- При получении вектора кратчайших расстояний от соседа маршрутизатор  $s$  сравнивает его поэлементно с текущей версией. Если оказывается, что наименьшая стоимость пути от соседа  $n$  до узла  $d$ , сложенная с оценкой стоимости ребра  $(s, d)$  меньше, чем наименьшая стоимость пути от  $s$  до  $d$  в текущей таблице, то значение в текущей таблице обновляется, и наилучшим соседом для отправки пакета в узел  $d$  становится сосед  $n$ .

Как можно видеть, дистанционно-векторный алгоритм является, в сущности, распределенной версией алгоритма Беллмана-Форда поиска кратчайшего пути в графе[11].

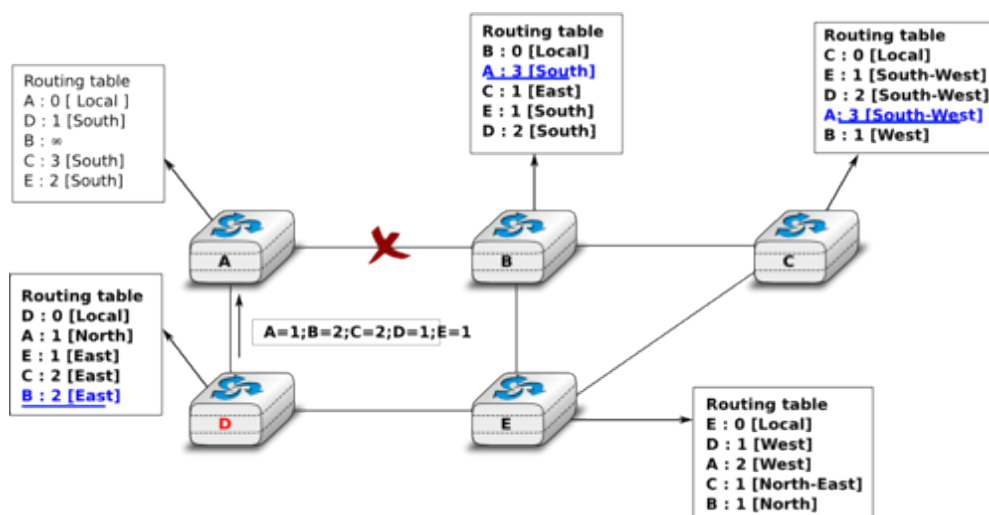


Рисунок 1 – Иллюстрация работы distance-vector алгоритма

Различные реализации дистанционно-векторного метода различаются, в частности, оценками стоимости соединений в сети. Так, например, протокол RIP[3] просто оценивает стоимость каждого соединения в 1, а IGRP[12] оценивает стоимость соединений исходя из оценок задержки и пропускной способности.



Преимуществами дистанционно-векторных алгоритмов являются простота реализации и низкие требования к памяти и вычислительной мощности. Недостатками же являются низкая скорость распространения информации по сети и сложности с приспособлением под изменяющуюся топологию (проблема count-to-infinity). Этих проблем удастся избежать при применении другого распространенного подхода - алгоритмов на основе состояния канала связи.

### 1.2.2. Алгоритмы состояния канала связи

В отличие от дистанционно-векторных алгоритмов, в алгоритмах состояния канала связи (link-state) каждый узел сети хранит у себя модель всей сети в виде графа. Рассмотрим шаги алгоритма подробнее:

- Каждый маршрутизатор периодически проверяет состояние соединений до соседей
- При обнаружении обрыва какого-либо соединения алгоритм удаляет это соединение из собственного графа и рассылает соседям новую версию состояния соединений до них
- Соседи обновляют собственные версии графов в соответствии с полученной информацией и пересылают сообщение дальше
- Чтобы избежать заикливания сообщений об обновлении состояния, каждое сообщение снабжается *номером версии*. Маршрутизатор  $n$  игнорирует сообщение от маршрутизатора  $n$ , если номер версии этого сообщения меньше или равен предыдущему.

Имея информацию обо всей сети в целом, маршрутизатор может рассчитать кратчайшие пути до всех остальных узлов. Обычно для этого используется алгоритм Дейкстры[13].

Link-state алгоритмы обладают способностью адаптироваться под изменения топологии сети гораздо быстрее, чем distance-vector алгоритмы за счет несколько более сложной реализации и чуть больших затрат по памяти и вычислительной мощности. Это обуславливает то, что на данный момент именно link-state протоколы, такие как OSPF[4], доминируют в сетевой маршрутизации. Однако даже в решении задачи сетевой маршрутизации link-state алгоритмы в чистом виде не лучшим образом адаптируются к повышению нагрузки в сети. Рассматриваемые в дальнейшем другие алгоритмы, основанные на принципе обучения

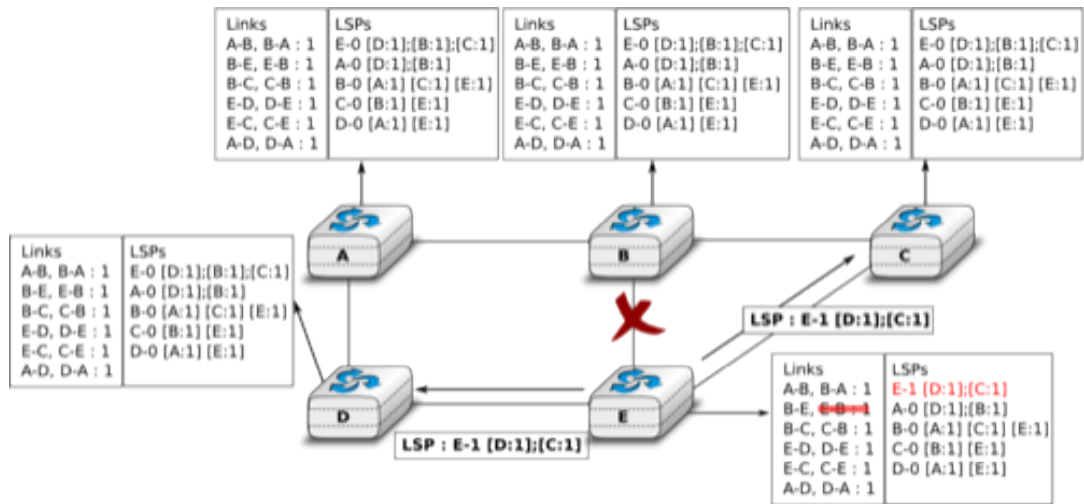


Рисунок 2 – Иллюстрация работы link-state алгоритма

с подкреплением, справляются с задачей адаптации к изменчивой нагрузке лучше.

### 1.2.3. Q-routing

Среди других подходов особый интерес представляют подходы на основе обучения с подкреплением. Первым алгоритмом маршрутизации, основанным на этой идее, стал алгоритм Q-routing[5]. Принцип его работы таков:

- Каждый маршрутизатор  $x$  хранит  $Q_x(d, y)$  — оценку минимального времени в пути до узла  $d$ , если следующим узлом на пути является сосед  $y$ . Очевидно, что  $\forall y : Q_x(x, y) = 0$
- Пакет, который необходимо доставить в узел  $d$ , отправляется соседу  $y = \operatorname{argmin}_{(x,y) \in E} Q_x(d, y)$
- При получении пакета узел  $y$  отправляет узлу  $x$  время получения  $t_r$  и собственную оценку оставшегося времени в пути  $t = \min_{(y,z) \in E} Q_y(d, z)$
- Зная время отправления пакета  $t_s$  и получив  $t_r$  и  $t$ , узел  $x$  обновляет собственную оценку по формуле:  $Q_x(d, y) = \alpha((t_r - t_s) + t - Q_x(d, y)) + Q_x(d, y)$ , где  $\alpha$  — это learning rate, параметр алгоритма.

Было показано, что этот алгоритм способен хорошо адаптироваться к изменениям в топологии сети и интенсивности трафика. Такие

его модификации, как dual Q-routing[8] и predictive Q-routing[7] еще более высокое качество маршрутизации. Однако по сравнению с distance-vector или link-state методами данные алгоритмы используют гораздо больше служебных сообщений (служебный пакет на каждую пересылку целевого пакета), что ограничивает их применение в реальных высоконагруженных компьютерных сетях.

Однако в задачах маршрутизации вне контекста компьютерных сетей это перестает быть проблемой, так как целевые “пакеты” (чемоданы на конвейере, автомобили на автостраде, etc.) и служебные сообщения в таких задачах являются объектами разной природы и передаются по разным каналам, причем служебные сообщения по сравнению с целевыми “пакетами” доставляются мгновенно. Эти обстоятельства делают применение алгоритмов обучения с подкреплением в таких задачах более привлекательным.

#### **1.2.4. Другие подходы**

Для полноты обзора приведем еще несколько примеров.

Идея использования нейросетей для решения задачи маршрутизации не нова. В работах [14, 15] для решения задачи поиска кратчайшего пути в графе используются нейронные сети Хопфилда. Однако, эти исследования преследовали цель ускорения вычисления кратчайшего пути за счет аппаратной реализации нейросети, что кардинально отличается от цели текущей работы.

Еще одним интересным подходом является AntNet[16]. Это алгоритм, построенный на идее исследования состояния сети с помощью специальных пакетов-“агентов”. Алгоритм показал хорошие результаты в ходе исследований, но не получил широкого применения, вероятно, в силу уже массового к тому времени распространения link-state и distance-vector протоколов.

#### **1.2.5. Выводы по обзору существующих решений**

Все алгоритмы маршрутизации, рассмотренные в статье, были разработаны для решения задачи маршрутизации именно в компьютерных сетях, но не для решения задачи маршрутизации в общей формулировке. Таким образом, есть потребность в разработке алгоритма, способного решать более общую задачу, что и является целью данной работы.

### 1.3. Постановка задачи в терминах обучения с подкреплением

#### 1.3.1. Термины и понятия

**Обучение с подкреплением** (reinforcement learning) — вид машинного обучения, в котором *агент* (agent) каждый момент времени  $t$  взаимодействует со *средой* (environment), находящейся в *состоянии* (state)  $s_t \in \mathcal{S}$  путем выбора *действия* (action)  $a \in \mathcal{A}_{s_t}$  и получения *вознаграждения* (reward)  $r_{t+1} \in \mathbb{R}$  с переходом в новое *состояние*  $s_{t+1}$ .

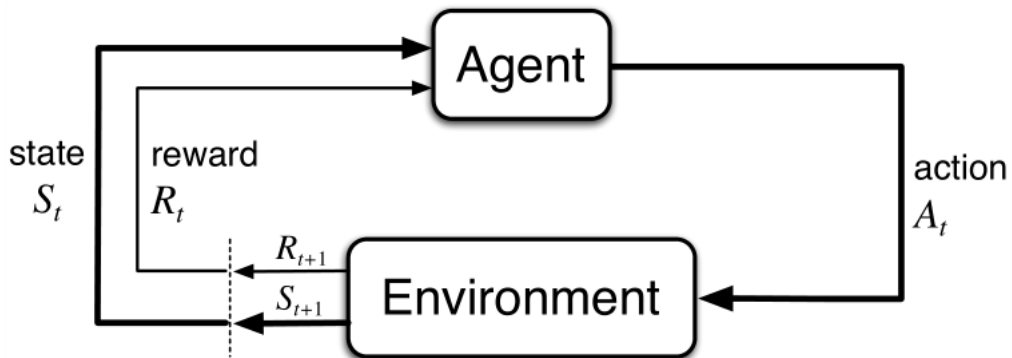


Рисунок 3 – Схема взаимодействия агента и среды в обучении с подкреплением

**Марковский процесс принятия решений** (Markov decision process, MDP) — это кортеж  $(\mathcal{S}, \mathcal{A}_s, P, R, \gamma)$ , где

- $\mathcal{S}$  — конечное множество состояний
- $\mathcal{A}_s$  — конечное множество действий, доступных из состояния  $s$
- $P(s'|s, a)$  — вероятность того, что действие  $a$  в состоянии  $s$  приведет к переходу в состояние  $s'$  в следующий момент времени.
- $R : \mathcal{S} \times \mathcal{A}_s \rightarrow \mathbb{R}$  — вознаграждение за действия  $a$  в состоянии  $s$
- $\gamma \in [0, 1]$  — *скидочный коэффициент* (discount factor), управляющий соотношением между важностью текущих вознаграждений и будущих вознаграждений.

**Оптимальной стратегией** для данного Марковского процесса принятия решений называется такая функция выбора действий  $\pi : \mathcal{S} \rightarrow \mathcal{A}_s$ , что взвешенная сумма вознаграждений  $\sum_{t=0}^{\infty} \gamma^t R_{\pi(s_t)}(s_t, s_{t+1})$  максимальна.

**Частично наблюдаемый Марковский процесс принятия решений** — это кортеж  $(\mathcal{S}, \mathcal{A}_s, P, R, \Omega, O, \gamma)$ , где

- $\mathcal{S}$  — конечное множество состояний

- $\mathcal{A}_s$  — конечное множество действий, доступных из состояния  $s$
- $P(s'|s, a)$  — вероятность перехода из  $s$  в  $s'$  при выполнении действия  $a$
- $R : \mathcal{S} \times \mathcal{A}_s \rightarrow \mathbb{R}$  — вознаграждение за действие  $a$  в состоянии  $s$ .
- $\Omega$  — множество *наблюдений*
- $O(o|s', a)$  — вероятность получения наблюдения  $o$  при переходе в истинное состояние  $s'$  в результате действия  $a$ .

Определение оптимальной стратегии для частично наблюдаемого Марковского процесса аналогично таковому для обычного.

**Q-обучение** (Q-learning)[6] — это метод нахождения оптимальной стратегии для Марковского процесса принятия решений, который заключается в оценке *функции полезности* (action-value function)  $Q(s, a)$ . Функция полезности изменяется при каждом предпринятом действии  $a$  с переходом из состояния  $s$  в  $s'$  по следующей формуле:

$$Q(s, a) = Q(s, a) + \alpha \left( r + \gamma \cdot \max_{a \in \mathcal{A}_{s'}} - Q(s, a) \right)$$

Известно, что для любого конечного Марковского процесса принятия решений Q-обучение находит оптимальную стратегию, т. е.  $Q(s, a) \xrightarrow{t \rightarrow \infty} Q^*(s, a)$ , и  $\pi(s) = \operatorname{argmax}_{a \in \mathcal{A}_s} Q^*(s, a)$  — оптимальная стратегия.

## 1.4. Обзор методов обучения нейросетей с подкреплением

### Выводы по главе 1

TBD

## **ЗАКЛЮЧЕНИЕ**

TBD

## СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

- 1 *McQuillan J. M., Walden D. C.* The ARPA network design decisions // Computer Networks. — 1977. — Т. 1, № 5. — С. 243–289.
- 2 *McQuillan J. M., Richer I., Rosen E. C.* The New Routing Algorithm for the ARPANet // IEEE Trans. on Comm. — 1980. — Т. 28, № 5. — С. 711–719.
- 3 *Hendrick C.* Routing Information Protocol: RFC. — Июнь 1988. — № 1058.
- 4 *Moy J.* OSPF Version 2: RFC. — Апр. 1998. — № 1058.
- 5 *Boyan J. A., Littman M. L.* Packet routing in dynamically changing networks: a reinforcement learning approach // Advances in Neural Information Processing Systems. — 1994. — No. 6. — P. 671–678.
- 6 *Watking C.* Learning from Delayed Rewards: дис. ... канд. / Watking C. — Cambridge : King's College, 1989.
- 7 *Choi S. P. . M., Yeung D.-Y.* Predictive Q-Routing: A Memory-based Reinforcement Learning Approach to Adaptive Traffic Control // Advances in Neural Information Processing Systems. — 1996. — No. 8. — P. 945–951.
- 8 *Kumar S., Miikkulainen R.* Dual reinforcement Q-routing: An on-line adaptive routing algorithm // Artificial Neural Networks in Engineering. — 1997. — No. 7. — P. 231–238.
- 9 *Yan J., Vyatkin V.* Distributed Software Architecture Enabling Peer to Peer Communicating Controllers // IEEE Transactions on Industrial Informatics. — 2013. — Vol. 9, no. 4. — P. 2200–2209.
- 10 Human-level control through deep reinforcement learning / V. Mnih [et al.] // Nature. — 2015. — No. 518. — P. 529–533.
- 11 *Bellman R.* On a routing problem // Quarterly of Applied Mathematics. — 1958. — № 16. — С. 87–90.
- 12 *Bosack L.* Method and apparatus for routing communications among computer networks. — Февр. 1992. — URL: <https://www.google.com/patents/US5088032> ; US Patent 5,088,032.
- 13 *Dijkstra E. W.* A note on two problems in connexion with graphs // Numerische Mathematik. — 1959. — № 1. — С. 269–271.

- 14 *Ali M. K. M., Kamoun F.* Neural networks for shortest path computation and routing in computer networks // IEEE Transactions on Neural Networks. — 1993. — T. 4, № 6. — C. 941–953.
- 15 *Araujo F., Ribeiro B., Rodrigues L.* A neural network for shortest path computation // IEEE Transactions on Neural Networks. — 2001. — T. 12, № 5. — C. 1067–1073.
- 16 *Di Caro G., Dorigo M.* AntNet: Distributed stigmergetic control for communications networks // Journal of Artificial Intelligence Research. — 1998. — T. 9. — C. 317–365.