



UCL

Бейсов подход в Дълбокото
Подсилено Обучение Бейсов подход в
Дълбокото Машинно Обучение

Firstname Surname

A thesis presented for the degree of
Doctor of Philosophy

Supervised by:
Professor Louis Fage
Captain J. Y. Cousteau

University College London, UK
January 2015

I, AUTHORNAME confirm that the work presented in this thesis is my own. Where information has been derived from other sources, I confirm that this has been indicated in the thesis.

Абстракт

Целта на настоящата работа е да създаде математически модел на самообучаващ се агент и да реализира агент, който напълно изследва дадена среда на Графичен Потребителски Интерфейс (ГПИ). Като входни данни за модела се използват изображения от ГПИ и информация за неговата сегментация. Агентът може да извършва действия, променяйки състоянието на ГПИ средата. Адекватността на извършените действия се оценява с награда, която средата предоставя. Наградата се определя от различни фактори, един от които е процентът покрит нов програмен код. Агентът успешно е изследвал средата, когато покритието на код на средата е пълно.

На базата на този модел ще създадем агент, който ще намира грешки в програмни продукти, ще синтезира пакети от автоматични тестове за качеството на софтуера и ще изпълнява посочени от потребителя задачи.

Благодарности

Съдържание

Абстракт	i
Благодарности	ii
List of tables	iv
Речник	v
1 Увод	1
1.1 Мотивация	1
1.2 Цели на дисертацията	4
1.3 Структура на дисертационния труд	5
2 Литературен обзор	7
2.1 Подсилено обучение	7
2.1.1 Дълбоко обучение	9
2.1.2 Дълбоко подсилено обучение	12
2.2 Бейсова статистика	14
2.2.1 Монте Карло алгоритми за Марковски Вериги (MCMC)	15
2.2.2 Извод със свободни вариационни параметри	16
2.2.3 Бейсови Невронни Мрежи	17
2.3 Автоматизирано тестване на ГПИ	17
2.3.1 Автоматизирано тестване на Android приложения	18
2.3.2 Проверка на качеството	20
3 Избор на действие	21
3.1 Литературен обзор	22
3.1.1 Обучение с учител (supervised learning)	23
3.1.2 Оптимизация	25

3.1.3	Изкуствени невронни мрежи (ИНМ)	26
3.1.4	Бейсово моделиране	27
3.1.5	Бейсови изкуствени невронни мрежи (БИНМ)	29
3.2	Пример	30
3.3	Модел	31
3.3.1	Бейсова изкуствена невронна мрежа (БИНМ)	32
3.4	Експерименти	33
3.4.1	Данни	33
3.4.2	Обучение	34
3.4.3	Резултати	34
3.5	Заключение	36
4	Среда за изучаване (RL exploration) на ГПИ приложения	48
4.1	Android специфична среда	49
4.2	Web специфична среда	50
5	Изучаване на ГПИ среди	51
5.1	Related Work	51
5.2	Дадено на агента	52
5.3	Задачи	53
5.4	Модел на агента	53
5.4.1	Пространство на състоянието (State space)	54
5.4.2	Пространство на действията (Action space)	54
5.4.3	Представяне на изображенията	55
5.4.4	Модел за избор на действия	55
5.4.5	Определяне на награди	55
5.4.6	Архитектура на модела	56
5.4.7	Цел/Оптимизация/Тренировка/Обучение	56
5.4.8	Памет	56
5.4.9	Оценка на модела за избор на действия	57
5.4.10	Намиране на апостериорно вероятно разпределение на действията	57
5.5	Експерименти	57
6	Заключение	58
6.1	Нерешени проблеми	58

6.2	Бъдеща работа	58
6.3	Дискусия	58
Приложение 1: Някои важни вероятностни разпределения		59
Приложение 2: Фигури		60
Литература		61

Списък на фигурите

3.1	<i>Примерен ГПИ на мобилно приложение предоставящо възможност за поръчка на цветя</i>	38
3.2	<i>Решетка приложена върху състояние на средата</i>	38
3.3	<i>Закодиране на състояние Start</i>	39
3.4	<i>Преходи на средата описани в T</i>	39
3.5	<i>Преходи при изучаване на AutoMath Photo Calculator</i>	40
3.6	<i>Закодирано представяне на AutoMath Photo Calculator</i>	40
3.7	<i>Преходи при изучаване на Memrise</i>	40
3.8	<i>Закодирано представяне на Memrise</i>	41
3.9	<i>Промяна на грешката по време на обучение върху примерната среда</i>	41
3.10	<i>Закодиране на последното състояние представено в примерната среда</i>	41
3.11	<i>Апостериорно разпределение за всяко действие</i>	42
3.12	<i>Апостериорно разпределение за всяко действие</i>	42
3.13	<i>Промяна на грешката по време на обучение върху AutoMath</i>	43
3.14	<i>Апостериорно разпределение за всяко действие</i>	43
3.15	<i>Апостериорно разпределение за всяко действие</i>	44
3.16	<i>Промяна на грешката по време на обучение върху Memrise</i>	44
3.17	<i>Апостериорно разпределение за всяко действие</i>	45
3.18	<i>Апостериорно разпределение за всяко действие</i>	45
3.19	<i>Промяна на грешката по време на обучение върху всички данни</i>	46
3.20	<i>Апостериорно разпределение за всяко действие</i>	46
3.21	<i>Апостериорно разпределение за всяко действие</i>	47

List of tables

Table 5.1 This is an example table . . .	pp
Table x.x Short title of the figure . . .	pp

Речник

API	A pplication P rogramming I nterface
JSON	J ava S cript O bject N otation

Глава 1

Увод

1.1 Мотивация

Последните години донесоха бързо развитие в сферата на Изкуственият Интелект (ИИ) и по-конкретно в дълбокото самообучение (Deep Learning). Моделите, използващи тези методи, имат разнообразни приложения и понякога се представят сходно или по-добре от човек, на същата задача:

- здравеопазване (Rajkumar et al. 2018), (Poplin et al. 2018), (Lee et al. 2017)
- самоуправляващи се автомобили (Bojarski et al. 2016), (Huval et al. 2015), (Bojarski et al. 2017)
- търгуване на стоковата борса
- създаване на изкуство

Една от основните цели на Изкуствения Интелект (ИИ) е да създаде агенти, които разбират и взаимодействат със света около нас. Значителен прогрес в тази насока беше постигнат през последните години благодарение на развитието на изчислителната техника (графични ускорители), наличието на голямо количество данни, нови начини за събиране и съхранението им и нови алгоритми. Бързият напредък в сферата на подсиленото обучение доведе до разработката на интелигентни агенти, взаимодействащи с все по-сложни среди (Mnih et al. 2015), (Silver & Hassabis 2016), (Levine et al. 2016) и (Silver et al. 2017). Критич-

ни за това са обучаващите алгоритми, техники за скалирането им и симулационни среди, които предоставят начини за оценка и сравняване на различни агенти (Bellemare et al. 2013), (Todorov et al. 2012) и (Johansson et al. 2016).

Хората се справят лесно с редица задачи, които изискват комплексно разбиране на визуалния свят, разпознаване на различни обекти в него и взаимодействие с тях. Например (екран от ГПИ среда и разпознаване на обекти в него). Би било лесно за човек да изучи подобна визуална среда.

Агенти, действащи в симулирани среди, са фундаментално ограничени - те никога не се сблъскват със сложността на реалния свят, поради което не могат да използват семантично знание и достигнат интелигентност. В роботиката агентите действат в реална среда, но процесът на обучение е бавен и скъп, дори и за тясно дефинирани задачи (Levine et al. 2016).

За справянето с тези проблеми могат да се използват среди, базирани на ГПИ приложения (Shi et al. 2017). Те предоставят разнообразни задачи, възможност за бързо итериране и обучение. Агентите получават същите сензорни данни, които получава човек, взаимодействащ с тези среди. Те предоставят възможност за изграждане на знание, невъзможно за придобиване в симулации.

Предизвикателства Тъй като подобни възможности изглеждат естествени за човек, може да забравим колко трудни са те за един агент. Изображенията са представени като голям масив от числа, които представят яркостта за всяка позиция. Едно такова изображение може да съдържа милиони такива *пиксели*, които агентът трябва да трансформира до семантични концепции на високо ниво, като например “текст” или “бутон”. При това, различни форми и цветове на даден бутон, също трябва да се класифицират като такъв, независимо от възможността за наличие на напълно различни шаблони (patterns) в яркостта на пикселите.

Концептуалното разбиране на дадено изображение е само първата стъпка за създаването на подобни агенти. Основна задача е създаване на модел на агент, който избира действия, които доведат до постигането на поставена задача. Трудността тук се изразява в липсата на пълна информация (fully observed) за средата в която агента действа. Например, натискането на един и същи бутон в различни състояния на средата, може да доведе до наблюдаването на две на-

пълно различни състояния на средата. Това означава, че е необходимо знание за конкретното състояние на средата.

Агентът няма предварителен модел на средата, която изучава. Той я “опознава”, чрез опити и грешки, като се опитва да приложи различни комбинации от действия в дадено състояние, за да постигне оптимална награда. На всяка стъпка даден агент трябва да избере дали да използва вече наученото или да избере действие, което не е изпълнил в конкретното състояние. Тази дилема се нарича: компромис на изучаване и използване на наученото (exploration exploitation tradeoff) и е основна задача за решаване от всеки агент. ГПИ средите могат да предоставят голям брой действия за дадено състояние (например меню системата на Microsoft Word), което прави пълното им изучаване неприложимо в кратки интервали от време.

Обнадеждаващ прогрес Въпреки сложността на задачата, през последните години се наблюдава значителен прогрес в областта на подсиленото обучение. По конкретно, развитието на Изкуствените Невронни Мрежи (ИНМ) (Artificial Neural Networks) и методи за създаване на Бейсови модели с милиони параметри доведоха до значително разширение на областите в които подсиленото обучение е приложимо. Алгоритми като Дълбоко Q-обучение (Deep Q-learning) допринасят за създаване на агенти, които надхвърлят възможностите на хората в тясно дефинирани задачи (Mnih et al. 2015), както и по-широко приложими такива (Silver et al. 2017)

Неотговорени въпроси Основният подход, използван в много от тези приложения е създаване на модел, който работи в добре дефинирани среди или такива в които се наблюдава пълна информация. Допълнително, агентите взимат решения на базата на изчислени точкови оценки. Възможно е това да намали ефективността на тези модели, както и обяснението на взетите решения. Открит остава въпросът дали добавянето на вероятностно разбиране за средата може да се справи с тези проблеми (Bellemare et al. 2017) (Anonymous 2018).

Принос В тази работа добавяме вероятностно разбиране за средата и разработваме модели и техники за ефективно Бейсово изучаване на ГПИ среди. Също така създаваме конкретна среда, която Агентът изучава. Например, агентът трябва да наблюдава дадено изображение и избере действие, което максимизира вероятността за постигане на поставена цел. Този модел ще опише процесът

за взимане на решения използвайки вероятностни разпределения. С други думи, целта на работата е създаване на агент, който ефективно изучава визуални среди.

Дългосрочна мотивация Основен стремеж на работата е да направи принос към изграждането и развитието на мислещи машини, както и приложи създадените модели в конкретни приложения. Техниките предложени тук са стъпка към достигането на бъдеще, в което агентите могат ефективно да взаимодействат с реалния свят (или по-сложни виртуални среди) и изпълняват комплексни задачи.

Краткосрочна мотивация Създаване на автоматизирани тестове за оценка на качеството на софтуерен продукт използвайки агент. Търсене за семантични и логически грешки в дадена програмна ГПИ среда. Възпроизвеждане на стъпки, необходими за възпроизвеждане на грешка.

Съществуващите методи не дават възможност за споделяне на наученото от друго приложение, намиране на аномалии във функционалността, бързо научаване на промени (премахване на старо знание) и оценка на несигурността при изпълнение на действие.

1.2 Цели на дисертацията

Целта на дисертацията е да дефинира политика π , определяща поведението на агента.

- Дефиниране на множество от възможни действия на агента (пространство на действията)
- Създаване на модел, който избира действията на агента
- Създаване на среда (за тестване и използване на агента) в която ще работи агентът
 - Подбиране на приложения (applications)
 - Измерване на покритието на код
 - Създаване на изображение и сегментация от текущото състояние на средата

- Избор на метрики за оценка на действията на агента
- Предварително обучение (с учител) на агента с данни от хора, взаимодействащи със средата (imitation learning)
- Провеждане на експерименти и анализ на постигнатите резултати

Целта на настоящата дисертационна работа е да създаде система за автоматизирано тестване на ГПИ, която използва за входни данни само визуалния изход на тестваното приложение. За постигане на целта трябва да се изпълнят следните задачи:

- Избор на подходящи оценъчни функции и награди, които да мотивират максималното покритие на програмен код по време на тестване
- Създаване на структури, в които да се запазват поредиците от стъпки, необходими за повтаряне на тестови случаи
- Създаване на модел, който генерира поредица от действия, използвани за играждане на тестовите случаи
- Създаване на модел, който намира аномалии по време на изпълнение на програмата
- Автоматично именуване на отделни екрани и действия с цел улесняване на разбирането
- Създаване и провеждане на експерименти, които да сравнят предложения модел с вече съществуващи такива
- По подадено изображение, йерархия на изгледите и действия, да се определи кои действия са валидни върху кои елементи

1.3 Структура на дисертационния труд

Глава 2 дава познания върху Дълбокото подсилено обучение и Бейсовото моделиране. **Глава 3** поставя целите и задачите на текущата работа. В **Глава 4** се създава среда за тестване на Android мобилни приложения. **Глава 5** представя модел за генериране на входни данни за тестови случаи. В **Глава 6** се представя модел за намиране на аномалии в тестови случаи по подадено изображение. Цялостната система е представена в **Глава 7** заедно с емпирични сравнения

спрямо други решения. Нерешени проблеми, бъдещи подобрения и дискусия се намират в последната глава.

Глава 2

Литературен обзор

2.1 Подсилено обучение

Една от основните задачи в сферата на изкуствения интелект е взимане на поредица от решения в стохастична среда. Един конкретен пример за взимане на решения в стохастична среда е агент, който изучава ГПИ. Тази задача се състои в избор на редица от решения, които да максимизират броят на разгледаните състояния на текущото приложение. Това е по-сложно от задачи, в които трябва да се направи само едно решение. Оценката за представянето на агента може да се даде само след много извършени от него стъпки. Това означава, че той може да избере неправилно действие сега и да разбере за това много по-късно, т.е. имаме *забавяне на последствията*. Допълнително, не може да наблюдаваме точното състояние на средата, поради липсата на точен модел на средата, която се изучава.

Основният начин за моделиране на такива среди са Марковски вериги.

Марковските вериги за вземане на решения (MDP) моделират системи, които искаме да контролираме. Във всяка времева стъпка t , системата се намира в дадено състояние s . Например, описаният агент може да се намира на даден екран от приложението, след като е натиснал определен бутон. Системата преминава през различни състояния като резултат от действията, които сме избрали. Задачата ни е да избираме действия, които са добри и да минимизираме броя на

тези, които не са. Разнообразни проблеми са моделирани чрез Марковски вериги (MDP формализма). Някои примери за използване на марковски вериги са системи за препоръки (Joachims et al. 1997), рутиране на мрежи (Boyan et al. 1994), управление на асансьори (Crites & Barto 1996), навигация на роботи (Sutton & Barto 1998).

Подсиленото обучение (RL) (Sutton & Barto 1998) дава начини за решаване на задачи, дефинирани чрез MDP формализма. Самообучаващ се агент с подсилено обучение (RL) взаимодейства със средата за определено време. На всяка времева стъпка t , агентът получава състояние s_t от пространството на състоянията S и избира действие a_t от пространство с действия A , следвайки политика $\pi(a_t|s_t)$. Политиката π определя поведението на агента, т.е. в определено състояние s_t , какво действие агентът трябва да избере. Тя дава функция за преобразуване на състояние s_t до състояние s_{t+1} чрез действие a_t . Използвайки дадена политика, агентът получава скаларна награда r_t и преминава в следващо състояние s_{t+1} , което се определя от функцията за награди $R(s, a)$ и функцията, даваща вероятности за преминаване в друго състояние $P(s_{t+1}|s_t, a_t)$. Когато моделът, който моделира поведението на агента е дискретен, т.е. може да се разглежда като отделни епизоди, описаният процес продължава докато агентът не достигне до крайно състояние. Тогава агентът се рестартира за започване на ново обучение. Общата награда е дефинирана като:

$$R_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$$

представлява обезценена стойност с фактор $\gamma \in (0, 1]$. Агентът се опитва да максимизира очакваната награда във всяко състояние.

Функция на стойностите $Q^\pi(s, a)$ дава предсказана обща бъдеща награда, която измерва до колко са добри дадено състояние или двойка състояние-действие. Стойността на дадено действие $Q^\pi(s, a) = E[R_t|s_t = s, a_t = a]$ ни дава очакваната награда за избиране на действие a в състояние s и следвайки фиксирана политика π . Оптимална стойностна функция $Q^*(s, a)$ предоставя действие a , което максимизира стойността на наградата за дадено състояние s . Може да дефинираме функция даваща стойност на състоянията $V^\pi(s)$, както и оптималната ѝ версия $V^*(s)$ по сходен начин.

Извод: ние подсилваме обучението с въвеждане на награда.

2.1.1 Дълбоко обучение

Нека разгледаме един от най-простите статистически модели - линейната регресия (Gauss 1809; Legendre 1805). Нека е дадено множество от N входно-изходни двойки $\{(x_1, y_1), \dots, (x_n, y_n)\}$. Например, нека x да е тегло в кг, а y - височина в см на N човека. Линейната регресия прави предположението, че съществува линейна функция, която преобразува всяко $x_i \in \mathbb{R}^Q$ към $y_i \in \mathbb{R}^D$. Тогава нашият модел е линейна трансформация на входните данни:

$$f(x) = xW + b$$

където W е $Q \times D$ матрица и b е вектор от D елемента. Тогава, задачата се свежда до намиране на такива параметри W и b , които минимизират средната квадратична грешка:

$$e = \frac{1}{N} \sum_i ||y_i - (x_i W + b)||^2$$

В общия случай, връзката между x и y може да не е линейна. Тогава искаме да дефинираме нелинейна функция $f(x)$, която преобразува входните данни до изходни. За тази цел може да приложим linear basis function regression (превод?) (Bishop 2007; Gergonne 1815), където входните данни x се подават на K фиксирани скаларни нелинейни трансформации $\phi_k(x)$ за създаване на свойствен вектор $\Phi(x) = [\phi_1(x), \dots, \phi_k(x)]$. Трансформациите ϕ_k наричаме базисни функции. Върху така създадения вектор се прилага линейна регресия. LBFR може да се сведе до линейна регресия, когато $\phi_k(x) := x_k$ и $K = Q$. Този тип функции се смятат за фиксирани и взаимно ортогонални. Когато тези ограничения се пропуснат говорим за *параметризирани* базисни функции.

2.1.1.1 Изкуствени невронни мрежи

Когато подредим параметризиран базисни функции в йерархия, може да говорим за изкуствени невронни мрежи. Всеки свойствен вектор в тази йерархия ще наричаме слой. Композицията от подобни слоеве води до голямата гъвкавост на тези модели. Често те постигат високи резултати на различни задачи и могат да се приложат върху реални проблеми, работещи върху терабайти от данни.

Feed-forward neural networks. Нека разгледаме модел с един *скрит слой* (Rumelhart et al. 1985). Нека x е вектор с Q елемента, представящ входните данни. Трансформираме го с афинна трансформация до вектор с K елемента. Отбелязваме с W_1 линейната преобразуваща матрица (матрица на теглата) и с b трансляцията използвана за трансформиране на x за да получим $xW_1 + b$. Върху всеки елемент на получената матрица се прилага нелинейна функция $\sigma(\cdot)$. Резултатът е т. нар. *скрит слой*, а всеки елемент се нарича *мрежова единица*. Върху резултата се прилага втора линейна трансформация с матрица на теглата W_2 , която преобразува скрития слой до изходен вектор с D елемента. Имаме $Q \times K$ матрица W_1 , $K \times D$ матрица W_2 и b - вектор от K елемента. Резултат от дадена невронна мрежа би бил:

$$\hat{y} = \sigma(xW_1 + b)W_2$$

при дадени входни данни x .

Когато използваме невронната мрежа за решаване на регресионна задача, може да минимизираме Евклидовата грешка:

$$e^{W_1, W_2, b}(X, Y) = \frac{1}{2N} \sum_{i=1}^N \|y_i - \hat{y}_i\|^2$$

където $\{y_1, \dots, y_n\}$ са N наблюдавани изходни стойности, $\{\hat{y}_1, \dots, \hat{y}_n\}$ са изходни данни от модела, а $\{x_1, \dots, x_n\}$ са входните данни. Предполагаме, че минимизирайки тази грешка спрямо W_1, W_2, b ще получим модел, който генерира добре при нови данни $X_{\text{test}}, Y_{\text{test}}$.

Когато задачата е да се предскаже класът, към който x принадлежи, от множес-

твото $\{1, \dots, D\}$, използваме същия модел. Промяната се състои в това, че прилагаме softmax функция върху получения резултат. Тази функция ни дава нормализирани оценки за всеки клас:

$$\hat{p}_i = \frac{\exp(\hat{y}_i)}{\sum d' \exp(\hat{y}_i')}$$

Когато вземем логаритъма от горната функция, получаваме softmax грешка:

$$e^{W_1, W_2, b}(X, Y) = -\frac{1}{N} \sum_{i=1}^N \log(\hat{p}_{i, c_i})$$

където $c_i \in \{1, \dots, D\}$ е наблюдаваният клас за вход i .

Описаният по-горе модел има проста структура, но може да бъде разширен за по-специализирани задачи. Този тип по-сложни модели се използват, когато задачите изискват обработка на поредици или изображения.

Convolutional Neural Networks CNN CNN е архитектура (LeCun et al. 1989), която се използва при изображения. Задачи, които до скоро се смятаха за нерешими, имат решения посредством този тип модели (Hinton et al. 2012). Моделът е създаден чрез рекурсивно приложение на конволуции и обединяващи слоеве. Конволуционният слой е линейна трансформация, която запазва пространствена информация от входното изображение.

Recurrent neural networks (RNN) RNN е модел (Rumelhart et al. 1985; Werbos 1988), базиран на поредици от данни, който се използва за обработка на текст, обработка на видео и други (Kalchbrenner & Blunsom 2013; Sundermeyer et al. 2012). Входните данни за RNN са поредица от символи. За всяка времева стъпка t , проста невронна мрежа е приложена върху единствен символ, както и изходните данни от мрежата от предишната стъпка.

Конкретно, при дадена редица от входни данни $x = [x_1, \dots, x_t]$ с дължина T , прост RNN модел е създаден чрез повтарящо се приложение на функция f_h . Така се генерира скрито състояние h_t за времева стъпка t :

$$h_t = f_h(x_t, h_{t-1}) = \sigma(x_t W_h + h_{t-1} U_h + b_h)$$

за някаква нелинейна функция σ . Изходните данни от модела може да бъдат дефинирани като:

$$\hat{y} = f_y(h_T) = h_T W_y + b_y$$

Съществуват и по-сложни RNN модели, като LSTM (Hochreiter & Schmidhuber 1997) и GRU (Cho et al. 2014).

2.1.2 Дълбоко подсилено обучение

Този тип методи се класифицират, когато използваме дълбоки невронни мрежи за апроксимиране на някой от компонентите на подсиленото обучение: функция на стойностите $V(s; \theta)$, политика $\pi(a|s; \theta)$ или модела за промяна на състояние и награди. Параметрите θ представляват тегла в дълбоки невронни мрежи. Когато използваме “плитки” модели, като например линейна регресия, дървета за вземане на решения и др. като апроксиматори на функция, имаме “плитко” подсилено обучение с параметри θ за съответния модел. Основната разлика между дълбокото и плиткото подсилено обучение се състои в апроксиматора на функцията, която използват. Когато се използва извън политикова апроксимация - например на нелинейни функции, може да се наблюдават нестабилност и разходимост (Tsitsiklis et al. 1997). Въпреки това, скорошната работа върху дълбоки Q -мрежи (Mnih et al. 2015) и *AlphaGo* (Silver & Hassabis 2016) стабилизират процеса на обучение и постигат много добри резултати.

Дълбокото подсилено обучение започна рязкото си развитие с работата на (Mnih et al. 2015). Преди това, RL даваше нестабилни резултати, когато се използваха нелинейни апроксиматори като невронни мрежи. Дълбоките Q мрежи (DQN) направиха няколко важни приноса: 1) стабилизиране на обучението, използвайки дълбоки невронни мрежи (Lin 1992) 2) подход за цялостно обучение без почти никакво познание за областта 3) обучаване на гъвкава невронна мрежа с еднакъв алгоритъм за изпълняване на различни задачи, например 49 Atari игри (Bellemare et al. 2013), на които се представят по-добре от всеки известен алгоритъм до момента.

2.1.2.1 Double DQN

(Van Hasselt et al. 2016) предложиха Double DQN (D-DQN) за справяне с проблема на прекалена увереност (overestimate?) на Q-learning алгоритъма. В базовият алгоритъм (както и в DQN), параметрите се обновяват според:

$$\theta_{t+1} = \theta_t + \alpha(y_t^\theta - Q(s_t, a_t; \theta_t))\Delta_{\theta_t} Q(s_t, a_t; \theta_t)$$

където

$$y_t^Q = r_{t+1} + \gamma \max_a Q(s_{t+1}, a; \theta_t)$$

така че оператора \max използва еднакви стойности, за да избере и оцени дадено действие. Като следствие от това, е по-вероятно да избере недостатъчно добри стойности. Double DQN предлага да оцени алчната политика спрямо невронна мрежа, но използва друга, за да оцени стойността ѝ. Това може да се постигне с малка промяна на DQN алгоритъма, заменяме y_t^Q с:

$$y_t^{D-DQN} = r_{t+1} + \gamma Q(s_{t+1}, \max_a Q(s_{t+1}, a; \theta_t); \theta_{\bar{t}})$$

където θ_t е параметър за първата невронна мрежа, а $\theta_{\bar{t}}$ е параметър за целевата мрежа.

2.1.2.2 Асинхронни методи

(Mnih et al. 2016) предложи асинхронни методи за четири RL алгоритъма: Q-learning, SARSA, n -step Q-learning and advantage actor-critic и asynchronous advantage actor-critic (A3C). Този подход използва паралелни агенти, които използват различни политики за изучаване на средата. Асинхронните методи могат да се изпълняват върху многоядрени процесори. Те се изпълняват много по-бързо и предоставят по-бързо обучение от други известни методи.

2.2 Бейсова статистика

Избиране на следващо действие по време на създаване на тестов случай пряко зависи от увереността във взимането на правилното решение. Несигурността от избиране на действие може да бъде моделирана посредством Бейсов подход.

Нека θ е неизвестна стойност, която може да е скаларна, векторна или матрица. Методите за статистически извод (inference) могат да ни помогнат да я намерим. Класическият статистически подход третира θ като фиксирана стойност. Единствената информация, която използваме за намиране на неизвестната стойност, идва от данните, с които разполагаме. Изводът се базира на резултат, получен от функцията на правдоподобие на θ , която свързва стойности от $p(y|\theta)$ с всяка възможност на θ , където $y = (y_1, \dots, y_n)$ е вектор с наблюдавани стойности.

Бейсовият подход третира θ като случайна стойност. За достигане на извод се използва разпределението на параметри при дадени данни $p(\theta|y)$. Това разпределение се нарича апостериорно. Освен функцията на правдоподобие, Бейсовият подход включва априорно разпределение $p(\theta)$, което представя вярванията ни за θ преди да се разгледат данните.

Теоремата на Бейс дава връзка между функцията на правдоподобие и априорното разпределение:

$$p(\theta|y) = \frac{p(\theta|y)p(\theta)}{p(y)}$$

където:

$$p(y) = \int p(y|\theta)p(\theta)d\theta$$

Формулата на Бейс може да бъде пренаписана по следния начин:

$$(1) \quad p(\theta|y) \propto p(y|\theta)p(\theta)$$

тъй като $p(y)$ не зависи от θ

Когато θ е многомерна величина може да напишем уравнение (1) използвайки маргиналните апостериорни разпределения като например:

$$p(\theta_1|y) = \int p(\theta|y)d\theta_2$$

където $\theta = (\theta_1, \theta_2)$. В много случаи резултатите са многомерни и точни изводи може да бъдат направени само аналитично. Поради тази причина често се използват приближения.

2.2.1 Монте Карло алгоритми за Марковски Вериги (MCMC)

MCMC алгоритмите правят неявно интегриране като взимат извадки от апостериорното разпределение. По този начин се намират приближения на стойностите, от които се интересуваме.

В съществото си тези методи създават Марковска верига с апостериорното разпределение на параметрите като стационарно разпределение. Когато веригата е крайна и повтаряща се, стойността на θ може да бъде оценена от извадки на средни пътища. Генерираните извадки $\theta^{(t)}, t = 1, \dots, N$ от това разпределение дават представа за целевото разпределение.

2.2.1.1 Метрополис-Хастингс алгоритъм

Този алгоритъм е предложен от Metropolis (Metropolis et al. 1953) и по-късно генерализиран от Hastings (Hastings 1970). Методът създава Марковска верига с желаното стационарно разпределение. Алгоритъмът избира кандидат стойност θ' от предварително избрано разпределение $q(\theta, \theta')$, където $\theta' \neq \theta$. Избраната стойност θ' се проверява чрез приеми-откажи метод (accept-reject step), за да се подсили, че принадлежи на целевото разпределение.

2.2.1.2 Извадки на Гибс

Този метод, предложен от Geman и Geman (Geman & Geman 1984), често се представя като специален случай на Метрополис-Хастингс алгоритъма.

2.2.2 Извод със свободни вариационни параметри

Variational Inference (VI) методите обикновено предлагат по-добри резултати спрямо МСМС, когато времето за изпълнение е ограничено. Допълнително предимство на тези подходи е, че те са детерминирани. Систематичната грешка и дисперсията се приближават до 0 при МСМС методите, за колкото повече време бъдат оставени да се изпълняват те. Тези свойства правят МСМС алгоритмите много ефективни на теория. В практиката обаче, времето за изпълнение и изчислителната мощ са ограничени. Това налага търсенето на по-бързо методи дори когато това намаля точността на получените резултати.

Този тип методи дефинират приближено вариационно разпределение $q_{\omega}(\theta)$, параметризирано от ω , с лесна за оценяване структура. Искаме приближеното разпределение да е максимално близко до това на апостериорното. За целта свеждаме задачата до оптимизационна и минимизираме Kullback-Leibler (KL) (Kullback & Leibler 1951) отклонението спрямо ω . Интуитивно, KL измерва приликата между две разпределения:

$$KL(q_{\omega}(\theta) || p(\theta|x, y)) = \int q_{\omega}(\theta) \log \frac{q_{\omega}(\theta)}{p(\theta|x, y)} d\theta$$

(Define x, y - dataset)

Този интеграл е дефиниран, когато $q_{\omega}(\theta)$ е непрекъснатата спрямо $p(\theta|x, y)$. Нека $q_{\omega}^*(\theta)$ е минимизираща точка (може да е локален минимум). Тогава KL може да ни даде приближение на апостериорното разпределение:

$$p(y^*|x^*, x, y) \approx \int p(y^*|x^*, \theta) q_{\omega}^*(\theta) d\theta =: q_{\omega}^*(y^*|x^*)$$

VI методите заменят изчисляването на интеграли с такова на производни. То-

ва е много подобно на оптимизационните методи използвани в DL. Основната разлика се състои в това, че оптимизацията е върху разпределения, а не точкови оценки. Този подход запазва много от предимствата на Бейсовото моделиране и представя вероятностни модели, които дават оценка на несигурността в изводите си.

2.2.3 Бейсови Невронни Мрежи

Един от големите недостатъци на съществуващите архитектури на невронни мрежи е, че изводите, които получаваме, са оценки на точки. Моделите не казват до колко са сигурни в предложените резултати. Когато например един лекар получи резултат от даден модел, той трябва да знае защо и как моделът е стигнал до него. Бейсовата статистика може да даде отговор на тези въпроси (Gal & Ghahramani 2015). Дори при модели използващи RNN, Бейсова интерпретация на задачата дава по-добри резултати от съществуващи такива (Gal & Ghahramani 2016).

Бейсови невронни мрежи, предложени в края на 80-те години (Kononenko 1989) и задълбочено изучавани по-късно (MacKay 1992a; Neal 2012), предлагат вероятностна интерпретация на моделите за дълбоко обучение, като представят теглата им като вероятностни разпределения. Този тип модели са устойчиви на пренастройване (overfitting), предлагат оценки на несигурността и могат да се тренират върху малко на брой данни.

2.2.3.1 Оценка на несигурността

2.3 Автоматизирано тестване на ГПИ

Проверката за правилно поведение на софтуер продукт е неизменна част от създаването му. Откриване и поправяне на всички потенциални проблеми преди той да бъде доставен до крайния потребител може да се сметне за най-добър случай.

2.3.1 Автоматизирано тестване на Android приложения

Мобилните приложения също имат нужда от проверка на качеството. Поради тази причина, в последните години засилено се разглеждат начини за автоматизацията на подобен вид тестове. Много голяма част от извършената работа до момента се състои в създаване на входни данни за приложения за мобилната операционна система Android. Подходите използвани до момента, се различават по начина, по който създават входни данни и изучават и използват евристики за приложението.

2.3.1.1 Съществуващи системи

Dynodroid (Machiry et al. 2013) е инструмент, който се базира на случайно изучаване. Предлага се и ръчен начин за въвеждане на входни данни, когато системата е заседнала.

MobiGUITAR (Amalfitano et al. 2015) строи модел на приложението по време на тестване. За всяко ново състояние се поддържа списък с възможни действия, които се изпълняват използвайки DFS (depth first search) стратегия.

SwiftHand (Choi et al. 2013) се опитва да максимизира покритието на код за тестваното приложение. Допълнително, инструментът се старее да минимизира броя рестартирания на приложението. SwiftHand генерира единствено докосвания и скролвания.

PUMA (Nao et al. 2014) предлага генерална среда за автоматизиране на ГПИ. Инструментът предлага рамка за програмиране, в която могат да бъдат имплементирани различни стратегии за изучаване на тестваното приложение.

2.3.1.2 Покритие на програмен код

Една от основните цели на системите за автоматизирано тестване на софтуер е да постигнат максимално покритие на програмния код. Няколко решения се опитват да постигнат това и за операционната система на Android.

BBoxTester (Zhauniarovich et al. 2015) е рамка за изготвяне на доклади относно

покритието на програмния код, без той да бъде наличен. За разлика от други подобни системи, BBoxTester предлага детайлни метрики за покритието на отделни класове, методи и т.н. В основата на системата се използва друг софтуерен продукт - Emma (Roubtsov & others 2005). BBoxTester е система с отворен код, намираща се на <https://github.com/zyrikby/BBoxTester>. За съжаление системата е неподдържана (от 2015г.) и несъвместима с нови версии на Android.

CovDroid (Yeh & Huang 2015) е друга система за тестване посредством подход на черната кутия (black-box testing). Програмният код на продукта не е наличен. CovDroid изчислява покритието на код като инструментираща кода на приложението и използва Android Debug Bridge (adb) за да наблюдава изхода от изпълнение на програмата.

ABCA (Huang et al. 2015) използва подход, много близък до този на CovDroid. Софтуерният пакет може да бъде намерен на <http://cc.ee.ntu.edu.tw/~farn/tools/abca/>. По време на този обзор, страницата на инструмента не беше активна. Авторите на статията не отговориха на запитването за активен адрес за изтегляне на ABCA.

GUITracer (Molnar 2015) представя иновативен подход за визуализация на покритието на код, когато приложението е базирано на ГПИ. Основен недостатък на системата е ограничението за работа върху Java AWT, SWING или SWT рамки за изграждане на ГПИ.

GroddDroid (Abraham et al. 2015) предлага автоматично намиране и изпълнение на зловреден софтуер (malware). Системата предлага и измерване на покрит код. Софтуерът може да бъде намерен на <http://kharon.gforge.inria.fr/>. Програмният код е ясно документиран и лесен за употреба. Един недостатък е използването на Logcat монитора за извличане на метрики за покритие на кода. Повечето от горепосочените системи използват този подход.

2.3.1.3 Текущо състояние (State of the art?)

(Choudhary et al. 2015)

2.3.2 Проверка на качеството

- Достатъчно бързо ли е? (Model should monitor for speed exec anomalies or report just slow parts)
- Как да повторя грешката? (Provide/execute steps for reproduction)
- Има ли разлики в изходните данни? (Change in hierarchy/image screenshot)

Глава 3

Избор на действие

Проверката за качество на софтуерни продукти често се извършва посредством автоматизирани, полуавтоматизирани или ръчно изпълняващи се тестове. Основна цел при създаване и изпълнение на тези тестове е постигане на високо или пълно покритие на създадения програмен код (Zhu et al. 1997). От своя страна, това покритие повишава възможността програмата да не достига непредвидени състояния и да притежава желаната функционалност (Ohba 1982).

Създаването на тестове, които проверяват цялостната функционалност на системата често се извършва от специалисти по проверка на качеството (QA). Те създават автоматизирани тестове или изпълняват проверката ръчно, спрямо предварително създадени спецификации. Част от тестовете, обхващащи целия софтуерен продукт се извършват спрямо графичния потребителски интерфейс (ГПИ), който този софтуер предоставя. Тези тестове (наречени ГПИ тестове (GUI tests)) симулират взаимодействието на потребител с програмата.

Създаването на автоматизирани ГПИ тестове е обвързано с трудности, като често променящи се визуални елементи, забавено изпълнение, достигане на непредвидени състояния на средата и др. (Memon 2002). Често, поради тази причина подобен вид тестове се изпълняват изцяло ръчно или полуавтоматизирано, което изисква взаимодействие с експерт.

QA експертът взаимодейства с ГПИ чрез поредица от действия (извършени чрез мишка, клавиатура, докосвания върху екран и/или др.), които променят

ГПИ и водят до друго нейно състояние (в частност, нов екран). Когато това състояние не е наблюдавано до момента, покритието на програмен код се увеличава, поради нуждата от изпълнение му за създаване на самото състояние.

Тогава, целта при създаване на ГПИ тестове може да се определи като посещаване на всяко състояние на визуалната среда поне веднъж. Повторно наблюдение на дадено състояние може да е необходимо поради допълнителни възможни действия. Действията, които се избират, определят последователността на наблюдаваните състояния, както и бързодействието на текущия тест (минимален брой на взети действия за постигане на целта).

В тази част от работата ще създадем модел, който избира следващо действие, когато средата се намира в определено състояние. Това действие трябва да бъде избрано, така че да максимизира увеличението на покритие на програмен код и минимизира нуждата възможността за попадане във вече наблюдавано състояние. Състоянието на средата ще бъде закодирано чрез матрица, отговаряща на елементите в нея. Това е опростен подход към решаване на поставената задача, като той ще бъде разширен в следващата глава.

За решаване на задачата ще използваме подход управляван от наличните данни (data-driven approach). Конкретно, създаваме БИНМ (Бейсова Изкуствена Невронна Мрежа), която приема състоянието на средата като входен параметър и изчислява апостериорните разпределения на вероятностите за възможните действия за да оценим до колко добро е всяко от тях. Обучението на БИНМ изисква предварително събрани данни.

3.1 Литературен обзор

Съществуват различни подходи за автоматизирано създаване на ГПИ тестове за мобилни и уеб приложения: (Amalfitano et al. 2015), (Choi et al. 2013), (Moreira & Paiva 2014), (Salvesen et al. 2015), (Moreira et al. 2017) (Memon 2002), но тяхната практическа употреба и ефективност са незадоволителни (Choudhary et al. 2015).

Предложеният подход е вдъхновен от работата представена в:

- (Mnih et al. 2015) - използват се ИНМ за обучение на агент, който играе игри, надвишавайки възможностите на човек в някои от тези игри
- (Shi et al. 2017) - среда, предоставяща възможност за създаване и обучение на агенти, изпълняващи задачи в уеб среди (напр. закупуване на самолетен билет)
- (Chang et al. 2010) - визуален скриптов език за създаване на тестове, който използва изображения за определяне на следващо действие

Текущият подход се различава по това, че:

- автоматизира напълно (или в голяма степен) създаването на ГПИ тестове
- предоставя среда, даващата информация за новото покритие на код при взимане на действие
- оценя несигурността (uncertainty) за избиране на действие, което може да е полезно за:
 - Състояния в които е необходима допълнителна информация за да бъде продължено изучаването на средата (напр. екран изискваш потребителско име и парола)
 - достигнато е неочаквано състояние (аномалия), което може да е свързано с грешка в програмният код

3.1.1 Обучение с учител (supervised learning)

Много практически проблеми могат да бъдат формулирани като намиране на функция $f : X \mapsto Y$, където X е пространство на входните данни, а Y е пространство на изходните данни. Често, дефинирането на f е трудно или невъзможно. Например, каква е функцията, която намира позицията на бутон в изображение от ГПИ?

Обучението с учител предлага подход, който използва примерни данни имащи вида: $(x, y) \in X \times Y$, за да намери функция, която предоставя добри приближения на резултатите на f .

Цел. Нека имаме обучителна извадка E от вероятностно разпределение D съдържаща n примера $\{(x_1, y_1), \dots, (x_n, y_n)\}$, които са независими и еднакво разпределени. *Обучение* наричаме търсенето на такава функция $f : X \mapsto Y$, която дава най-близки резултати до тези от обучителната извадка. Обучението се състои в избиране на функция на загубата/грешката (loss) $L(\hat{y}, y)$ измерваща несъгласието между предсказаната стойност $\hat{y} = f(x)$ и истинската стойност y . Целта на обучението е да намери $f^* \in \mathcal{F}$, която удовлетворява уравнението:

$$f^* = \arg \min_{f \in \mathcal{F}} E_{(x,y) \sim D} L(f(x), y)$$

където \mathcal{F} е някакво множество от възможни функции. Обучението се свежда до търсене на такава f^* , която минимизира очаквана грешка над D .

Поставеният оптимизационен проблем е нерешим, защото нямаме достъп до всички $d \in D$. Следователно, не може да намерим очаквана грешка или да я опростим, без да наложим силни предположения относно D , L или f . Ако използваме предположението за независимост и еднаква разпределеност може да получим приближение на очакваната грешка като използваме извадки от осреднената грешка върху обучителната извадка:

$$f^* \approx \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n L(f(x_i), y_i)$$

Предполагаме, че минимизирането на грешката върху обучителната извадка E ще ни даде таква функция f , която минимизира грешката върху всички стойности от D .

Нека разгледаме един от най-простите статистически модели - линейната регресия (Gauss 1809; Legendre 1805). Нека е дадено множество от n входно-изходни двойки $\{(x_1, y_1), \dots, (x_n, y_n)\}$. Например, нека x да е тегло в кг, а y - височина в см на n човека. Линейната регресия прави предположението, че съществува линейна функция, която преобразува всяко $x_i \in \mathbb{R}$ към $y_i \in \mathbb{R}$. Тогава нашият модел е линейна трансформация на входните данни:

$$f(x) = W^T x + b$$

където W и b са параметри на модела, които трябва да оценим. Често използвана функция оценява грешката е квадрата от разликата между предсказаната и истинската стойност:

$$L(\hat{y}, y) = (\hat{y} - y)^2$$

Тогава, задачата свеждаме до намиране на такива стойности за W и b , които минимизират грешката:

$$f^* = \arg \min_{w,b} \left[\frac{1}{n} \sum_{i=1}^n (W^T x_i + b - y_i)^2 \right]$$

Обучението на модел се свежда до оптимизационен проблем, който има следната обща форма $\theta^* = \arg \min_{\theta} g(\theta)$, където θ са параметри на модела, а $g(\theta) = \frac{1}{n} \sum_{i=1}^n L(f_{\theta}(x_i), y_i)$.

3.1.2 Оптимизация

Нека предположим, че g е диференцируема. Тогава може да намерим градиента на g :

$$\nabla_{\theta} g = \theta \frac{\delta g}{\delta \theta}$$

Градиентът е вектор от частни производни даващ ни наклона на g от всяко измерение за θ . Той може да бъде използван за посока на търсене - може да подобрим оценката на θ като добавим малка стойност от негативната посока на градиента (тъй като търсим минимум на функция). Това мотивира създаването на алгоритъмът за постепенно спускане (gradient descent) (Cauchy 1847), който работи на две стъпки:

1. намери градиента

2. обнови параметрите, като използваш малка стъпка в посока на негативния градиент

Процесът продължава докато не се достигне желаната грешка. Например, намирането на стойностите на W се свежда до продължителното прилагане на:

$$W := W - \alpha \frac{\delta g}{\delta W}$$

Където α е размерът на стъпката. Ако стойността е прекалено висока, може да не се достигне сходимост, когато стойността на α е прекалено ниска обучението отнема прекалено дълго време.

3.1.3 Изкуствени невронни мрежи (ИНМ)

В предишните раздели видяхме, че може да дефинираме произволна диференцируема функция, която съпоставя входни данни x на предсказани стойности \hat{y} . За обучение на моделът използвахме постепенно спускане. Нека разгледаме по-подробно функцията f и намиране на градиента, необходим за изпълнения на алгоритъма.

Построяването на изкуствена невронна мрежа се състои в повторение на умножение на матрици и прилагане на активационни функции. Активационните функции (напр. sigmoid, tanh, ReLU) позволяват на ИНМ да апроксимират нелинейни функции.

Например, построяването на двуслойна невронна мрежа може да се представи като $f(x) = W_2\sigma(W_1x)$, където W_1 и W_2 са матрици, а σ е активационна функция. Ако σ е identity функция, ИНМ е линейна функция.

Теорема Универсална теорема за апроксимация изкуствена невронна мрежа с поне 1 скрит слой може да апроксимира произволна функция с произволна точност (Cybenko 1989), (Hornik et al. 1989).

Биологично вдъхновение. Изкуствените невронни мрежи са вдъхновени от груб модел на биологичния неврон. Всеки ред от матрицата с теглата W моделира един неврон и силата на връзките от входните данни. Претеглената сума

от входните данни пристига в тялото на клетката и се прилага активационна функция, което се интерпретира като скорост на предаване (firing rate) на неврона.

3.1.3.1 Backpropagation

3.1.4 Бейсово моделиране

Избиране на следващо действие по време на създаване на тестов случай пряко зависи от увереността във взимането на правилното решение. Несигурността от избиране на действие може да бъде моделирана посредством Бейсов подход.

Нека θ е неизвестна стойност, която може да е скаларна, векторна или матрица. Методите за статистически извод (inference) могат да ни помогнат да я намерим. Класическият статистически подход третира θ като фиксирана стойност. Единствената информация, която използваме за намиране на неизвестната стойност, идва от данните, с които разполагаме. Изводът се базира на резултат, получен от функцията на правдоподобие на θ , която свързва стойности от $p(y|\theta)$ с всяка възможност на θ , където $y = (y_1, \dots, y_n)$ е вектор с наблюдавани стойности.

Бейсовият подход третира θ като случайна стойност. За достигане на извод се използва разпределението на параметри при дадени данни $p(\theta|y)$. Това разпределение се нарича апостериорно. Освен функцията на правдоподобие, Бейсовият подход включва априорно разпределение $p(\theta)$, което представя вярванията ни за θ преди да се разгледат данните.

Теоремата на Бейс дава връзка между функцията на правдоподобие и априорното разпределение:

$$p(\theta|y) = \frac{p(\theta|y)p(\theta)}{p(y)}$$

където:

$$p(y) = \int p(y|\theta)p(\theta)d\theta$$

Формулата на Бейс може да бъде пренаписана по следния начин:

$$p(\theta|y) \propto p(\theta)y p(\theta)$$

тъй като $p(y)$ не зависи от θ .

Когато θ е многомерна величина може да напишем уравнение 3.1.4 използвайки маргиналните апостериорни разпределения като например:

$$p(\theta_1|y) = \int p(\theta|y) d\theta_2$$

където $\theta = (\theta_1, \theta_2)$. В много случаи резултатите са многомерни и точни изводи не могат да бъдат направени, дори аналитично. Поради тази причина често използваме извадки от апостериорното разпределение.

Пример. Нека имаме обучителна извадка с размер n : $X = \{x_1, \dots, x_n\}$ и $Y = \{y_1, \dots, y_n\}$. Искаме да намерим параметрите θ на функцията $y = f^\theta(x)$, които са вероятно са използване за генериране на обучителната извадка.

Следвайки Бейсовият подход прилагаме априорно разпределение над параметрите, $p(\theta)$ и дефинираме функция на правдоподобие $p(y|x, \theta)$.

За класификационни задачи може да използваме softmax функция на правдоподобие:

$$p(y = d|x, \theta) = \frac{\exp(f_d^\theta(x))}{\sum_{i=1}^n \exp(f_i^\theta(x))}$$

или Гаусова функция на правдоподобие за регресионни задачи:

$$p(y|x, \theta) = \mathcal{N}(y; f^\theta(x), \tau^{-1}I)$$

където τ определя прецизността на модела.

При дадена обучителна извадка (X, Y) може да намерим апостериорното разпределение за параметър θ използвайки правилото на Бейс:

$$p(\theta|X, Y) = \frac{p(Y|X, \theta)p(\theta)}{p(Y|X)}$$

Използвайки апостериорното разпределение може да направим изводи относно ненаблюдавани данни x^* :

$$p(y^*|x^*, X, Y) = \int p(y^*|x^*, \theta)p(\theta|X, Y)d\theta$$

Намиране на апостериорно разпределение. Основен компонент при намирането на апостериорното разпределение е знаменателят (нормализатор) в правилото на Бейс:

$$p(Y|X) = \int p(Y|X, \theta)p(\theta)d\theta$$

Това интегриране се нарича още маргинализация на функцията на правдоподобие над θ . Маргинализацията може да бъде извършена аналитично за прости модели като Бейсовата линейна регресия. В подобни модели вероятностното разпределение на функцията на правдоподобие е спрегнато (conjugate) с това на вероятностното разпределение на функцията на апостериорното разпределение, което позволява аналитично решение. Такова решение не може да бъде намерено, когато моделите са по-сложни, защото искаме да приложим маргинализация върху всички възможни стойности на параметъра θ . За оценяване на параметри в по-сложни модели може апроксимиращи методи.

3.1.5 Бейсови изкуствени невронни мрежи (БИНМ)

Бейсовите изкуствени невронни мрежи, предложени в (Tishby et al. 1989) и подробно разгледани в (MacKay 1992a), (MacKay 1992b) и (Neal 2012), предоставят вероятностна интерпретация чрез представяне теглата на изкуствените неврони като вероятностни разпределения. Тези модели са издръжливи на прекомерно нагаждане (overfitting), предоставят оценка на несигурността при взимане

на решения и могат да се обучават с малки извадки.

БИНМ поставят априорно разпределение върху теглата на невронната мрежа. Най-често се използва Гаусово вероятностно разпределение, приложено върху матрицата на параметрите $p(W_i) = \mathcal{N}(0, 1)$. Функцията на правдоподобие се дефинира използвайки уравнение 3.1.4 или 3.1.4.

3.1.5.1 MC Dropout

3.2 Пример

Ще разгледаме мобилно приложение предоставящо възможност за поръчка на цветя с ГПИ представен на 3.1.

ГПИ се състои от 4 различни състояния, като началното е маркирано със *Start*.

Ще опростим задачата, като приложим “решетка”, която разделя изображенията на средата на 4 правоъгълника с равни лица, визуализирани на 3.2.

Това представяне ни позволява да изпълняваме следните 5 различни действия:

- a_1 - клик горе в ляво
- a_2 - клик горе в дясно
- a_3 - клик долу в ляво
- a_4 - клик долу в дясно
- a_5 - връщане назад

Ще закодираме съдържанието на всяка клетка в решетката като:

- w - бял цвят, върху който не могат да бъдат предприемани действия (изображение или празно пространство)
- b - син цвят, който представя текстова информация
- g - зелен цвят, който представя бутон

Състояние *Start* може да закодираме като 3.3

Нека след първоначално обучение от специалист качество на софтуер (QA expert) имаме матрицата на преходите T , дефинирана като:

s_{x_1}	s_{x_2}	s_{x_3}	s_{x_4}	action
b	w	g	g	a_3
b	w	w	b	a_5
b	w	g	g	a_4

където $s = (s_{x_1}, s_{x_2}, s_{x_3}, s_{x_4})$ е вектор от характеристиките на състоянието

Преходите от матрица T може да се представят като:

Получаваме ново състояние, което не е описано в T . Свеждаме задачата до пресмятане на апостериорните разпределения на действията и намиране на оптималното действие.

3.3 Модел

Нека имаме среда E , намираща се в състояние $s \in \mathbb{S}$, върху което могат да бъдат изпълнени действия от множеството от действия \mathbb{A} . При избор на действие $a \in \mathbb{A}$, средата E преминава в ново състояние s' (в частност, $s' = s$, т.е. средата може да не премине в ново състояние). Множеството \mathbb{A} е ненаредено и всяко $a \in \mathbb{A}$ може да се обозначи с единствено цяло число, като по този начин въвеждаме наредба в \mathbb{A} . Всяко състояние на средата S позволява изпълнението на действия \mathbb{A} , които са предварително дефинирани. Множеството от всички възможни състояния на средата \mathbb{S} е неизвестно.

Нека след първоначално обучение от специалист имаме матрица на преходите T с размерност $n \times 2$, където n е броя на преходите. Всеки ред от T дефинира наредена двойка $(,)$, която описва оптималните действия за състоянията.

Нека имаме състояние s' , за което T не съдържа информация. В този случай, целта е да намерим наредено подмножество от подходящи действия.

3.3.1 Бейсова изкуствена невронна мрежа (БИНМ)

Вероятностното разпределение над всички възможни действия би позволило оценяване на несигурността при избор на действие. С тази информация може да решим кога да използваме знанията за средата и кога да я изучаваме (Azizzadenesheli et al. 2018). Когато вероятността е по-голяма ще е по-вероятно да изберем конкретното действие.

Ще използваме бейсова невронна мрежа със следната архитектура:

- входен слой: слой с 12 неврона
- първи скрит слой: пълно свързан слой (fully-connected) с 20 неврона
- втори скрит слой: пълно свързан слой (fully-connected) с 15 неврона
- изходен слой: слой с 5 неврона (броя на възможните действия)

Прилагаме ReLU активизационна функция, предложена в (Hahnloser et al. 2000) и отпадане (dropout), предложен в (Srivastava et al. 2014), с вероятност за отпадане на неврон $p = 0.2$ върху първи и втори скрит слой. Допълнително, прилагаме нормализация на група от данни (batch normalization), предложена в (Ioffe & Szegedy 2015), след втори скрит слой. Описаната БИНМ може да се представи в PyTorch (Paszke et al. 2017) като:

```
class Model(nn.Module):

    def __init__(self):
        super(Model, self).__init__()
        self.fc1 = nn.Linear(12, 20)
        self.drop1 = nn.Dropout(p=0.2)
        self.fc2 = nn.Linear(20, 15)
        self.bn = nn.BatchNorm1d(15)
        self.drop2 = nn.Dropout(p=0.2)
        self.fc3 = nn.Linear(15, 5)

    def forward(self, x):
        x = F.relu(self.drop1(self.fc1(x)))
```

```
x = F.relu(self.bn(self.drop2(self.fc2(x))))
return self.fc3(x)
```

3.4 Експерименти

В следващите експерименти ще приложим описания БИНМ модел върху 3 различни ГПИ среди. Към всяка среда имаме данни предоставени от QA експерт, закодирани по схемата представена в таблица на преходите T .

Моделът ще оценим като използваме кръстосана ентропия (cross entropy) дефинирана като:

$$H(p, q) = E_p[-\log q] = H(p) + D_{KL}(p||q)$$

където p и q са вероятностни разпределения, съответно на предсказани и наблюдавани стойности, $H(p)$ е ентропията на p и $D_{KL}(p||q)$ е Кулбак-Лейблер разстоянието до q от p . За дискретни p и q имаме:

$$H(p, q) = - \sum_x p(x) \log q(x)$$

В проведените експерименти $H(p, q)$ има вида

$$H(p, q) = - \sum_{x=1}^5 p(x) \log q(x)$$

3.4.1 Данни

Първата среда създадена за целите на текущата работа и е представена в описания пример по-горе.

Втората среда представя мобилното приложение *AutoMath Photo Calculator* и използва данни от Risco представени в (Deka et al. 2017). Приложението дава решения на математически задачи след заснемането им с камерата на устройс-

твото. Част от визуалната среда е представена на 3.5

Прилагайки предложения по-горе подход за закодиране върху тази среда получаваме 3.6

Третата среда е базирана на част от мобилното приложение *Memrise*. То предоставя флашкарти за изучаване на чужди езици и материали създадени от потребителя.

Закодираната версия е представена на 3.8

3.4.2 Обучение

Целта е да минимизираме грешката $H(p, q)$. Ще използваме *SGD* оптимизатор със скорост на обучение (learning rate) $lr = 0.1$ и движеща сила (momentum) $m = 0.5$. Описаният модел обучаваме за $1,000$ епохи върху всяка среда по отделно. Записваме грешката върху тренировъчните данни след обучението на модела във всяка епоха.

3.4.3 Резултати

За извличане на извадки от апостериорното разпределение на действията използваме алгоритъм Монте Карло отпадане (MC dropout) и прилагаме отпадане по време на тестване $n = 10,000$ пъти. Стойностите във всяка извадка се нормализират - отнемат всички стойности с минималната стойност от извадката и сумата от всички стойности приравняваме на единица.

3.4.3.1 Примерна среда

От графиката се вижда как след първите 50 епохи от обучението, средната грешка намалява до стойност близка до 0.5 и след това остава постоянна. В идеалния случай грешката трябва да намалява постепенно и да достигне стойности близки до 0. Ако графиката много бързо достига до 0, моделът се влияе прекалено силно от обучителната извадка.

Матрицата T предоставя малко данни за примерната среда (само 3 реда), но въпреки това предложеният модел успява да намали тренировъчната грешка. Разбира се, възможно е модела да се е нагодил (overfit) спрямо данните, които използваме за обучение. Средната стойност и стандартното отклонение на апостериорното разпределение на вероятността действието да е оптимално са дадени на 3.4.3.1.

	a_1	a_2	a_3	a_4	a_5
mean	0.006	0.025	0.427	0.434	0.107
std	0.010	0.018	0.036	0.031	0.073

От таблица 3.4.3.1 виждаме, че действие a_4 е оптималното действие за последното наблюдавано състояние.

Последното състояние получено от примерната среда при изпълнение на действията описани в T е представено на 3.10

При неколккратно изпълнение за обучение и оценка на модела наблюдаваме, че получените извадки се различават и понякога a_3 има по-висока средна стойност от a_4 . Причината за това е малкото количество данни в T .

Апостериорното разпределение е представено на 3.11 и обобщено в 3.12

3.4.3.2 AutoMath

Резултатите от прилагането на предложеният модел за втората среда, използвайки реални данни от Riso, са представени в следващите графики:

От графика ?? се вижда как след първите 500 епохи от обучението, средната грешка клони към 0. От графиките 3.14 и 3.15 се вижда, че моделът е много сигурен в избора на действие a_4 като оптимално.

3.4.3.3 Memrise

От графика ?? се вижда как след първите 400 епохи от обучението, средната грешка клони към 0. От графиките 3.17 и 3.18 се вижда, че моделът е много сигурен в избора на действие a_4 като оптимално и в отхвърлянето на действие a_5 .

3.4.3.4 Използване на всички данни

Промяната на $H(p, q)$ по време на обучение използвайки всички налични данни е представена на 3.19.

Получихме желаното поведение на грешката.

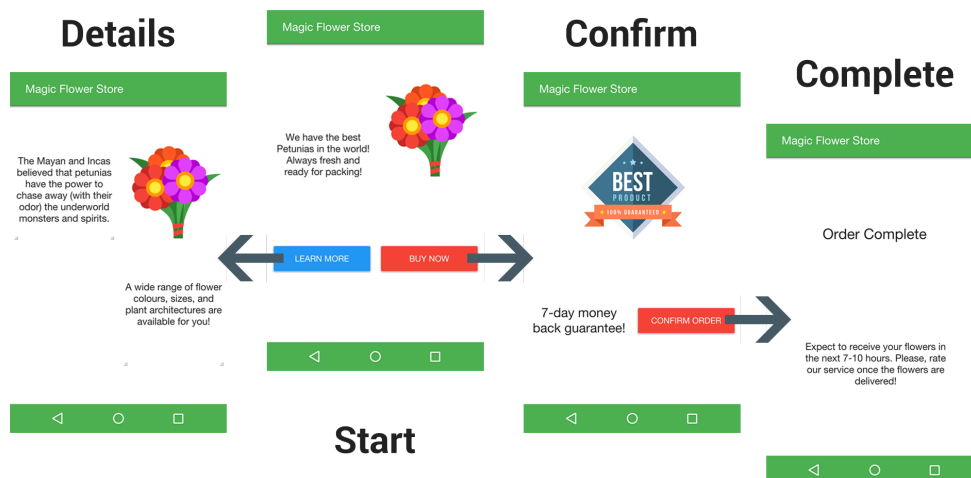
Графики 3.20 и 3.21 представят предсказани стойности от модела, когато използваме състояние ?? от примерната среда. Наблюдаваме, че вече моделът е много сигурен в избора на действие a_4

3.5 Заключение

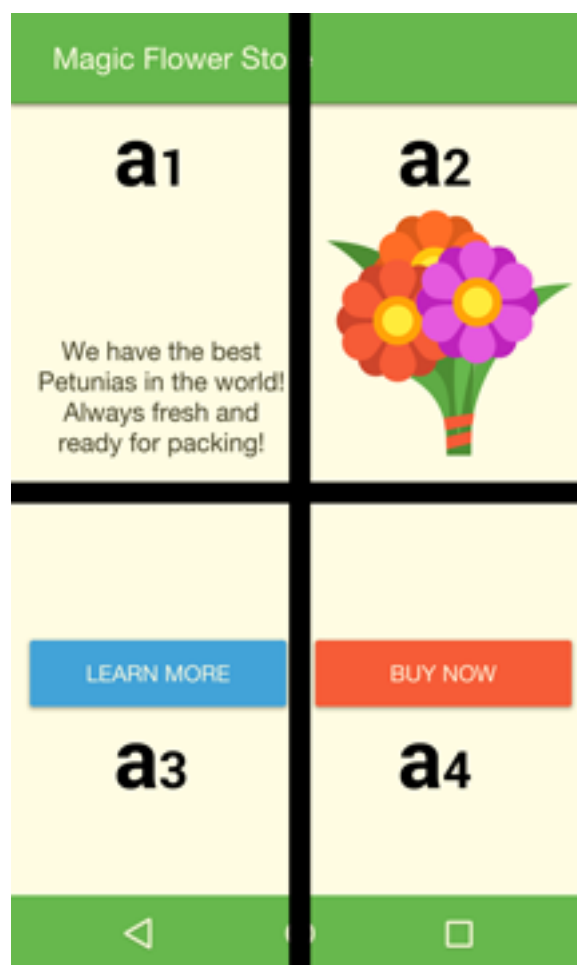
Разработихме модел, който може да избира оптимално следващо действие при представено закодирано състояние на ГПИ среда. Допълнително, моделът ни дава възможност да оценим до колко сигурен е в предсказаните стойности. Моделът се състои от дълбока Бейсова невронна мрежа, която получава закодирано състояние на средата като входни данни и е обучена използвайки кръстосана ентропия оценяща грешката. Извадки от апостериорното разпределение на действията извлякохме използвайки MC dropout. Оценката на обучените модели показва добра ефективност при избор на действие и генерализация между различните среди, когато използваме данни от всички среди.

Едно от ограниченията на предложения модел е предположението, че имаме външен агент, който закодира изображенията на средата до представяне, което е подходящо за обучение. В глава 5 ще разширим предложения модел, така че входните данни да са “сурови” - изображенията, които средата предоставя. Друго ограничение на нашия модел е, че не взима под предвид последователността

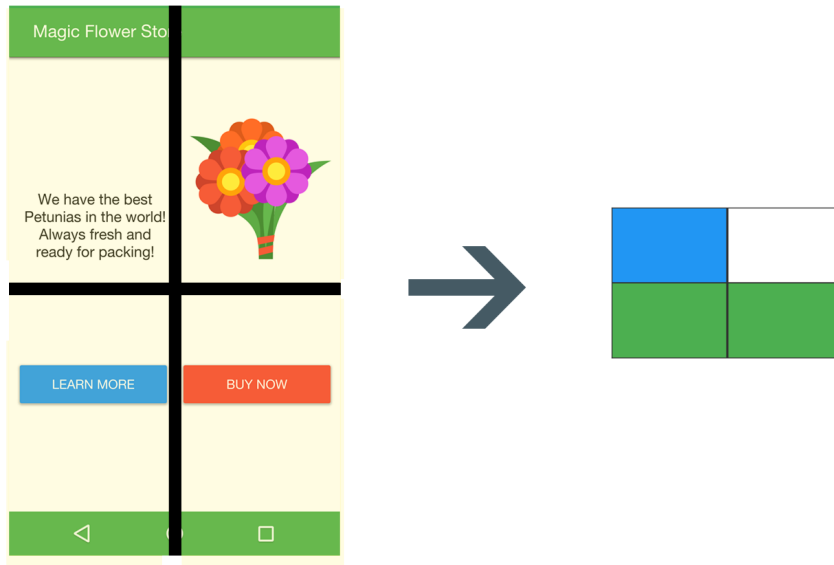
от взетите действия по време на изучаването на ГПИ средата. Следователно, агентът има възможност да опита да изучи състояния от средата, които вече е изучил.



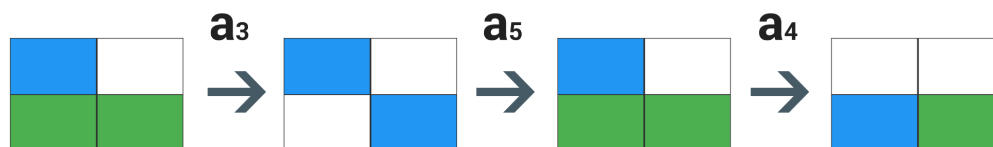
3.1:



3.2:

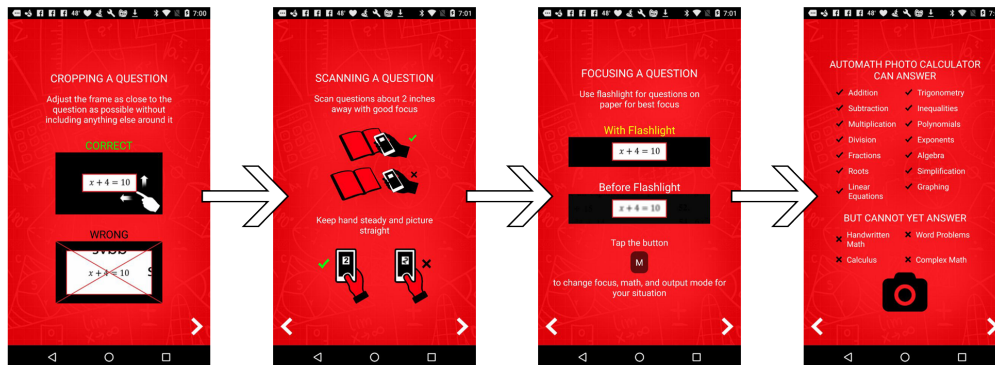


3.3: *Start*



3.4: *T*

Start

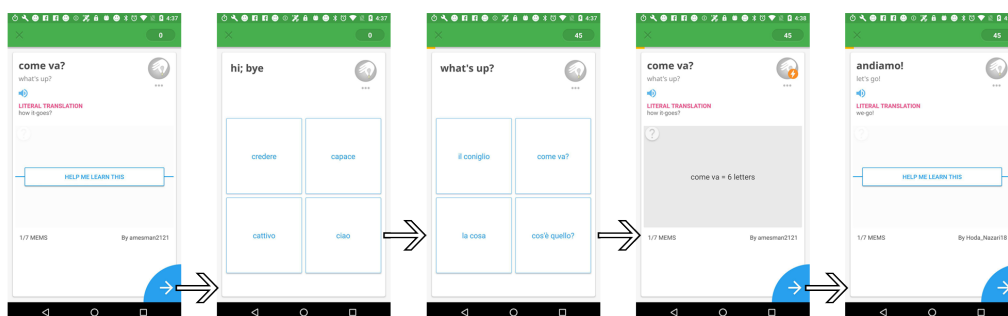


3.5: *AutoMath Photo Calculator*

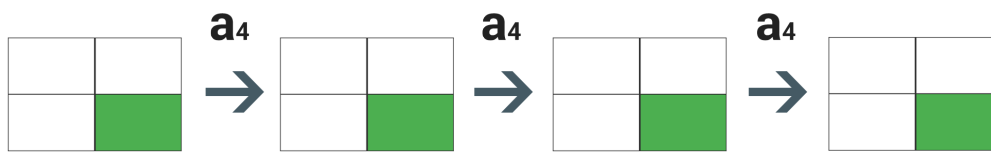


3.6: *AutoMath Photo Calculator*

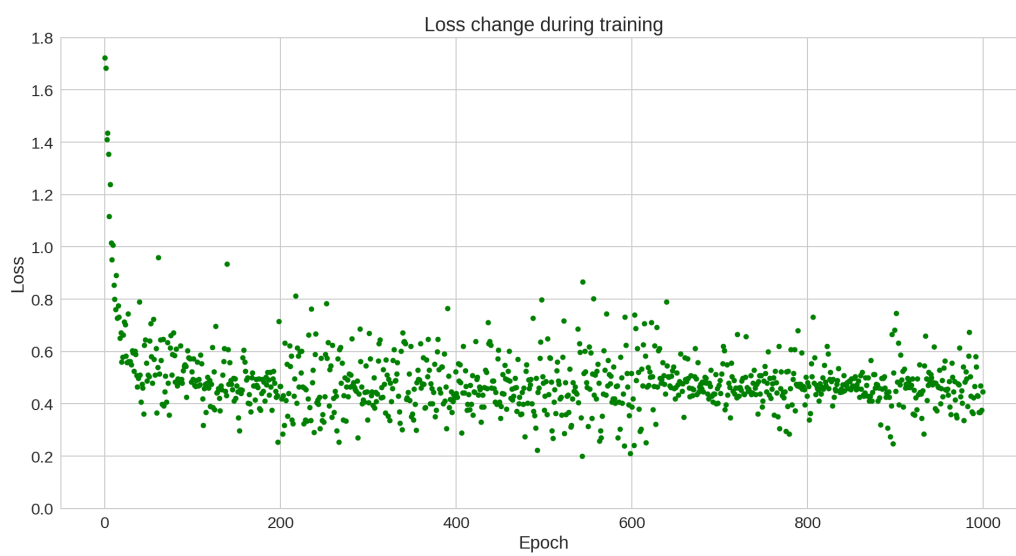
Start



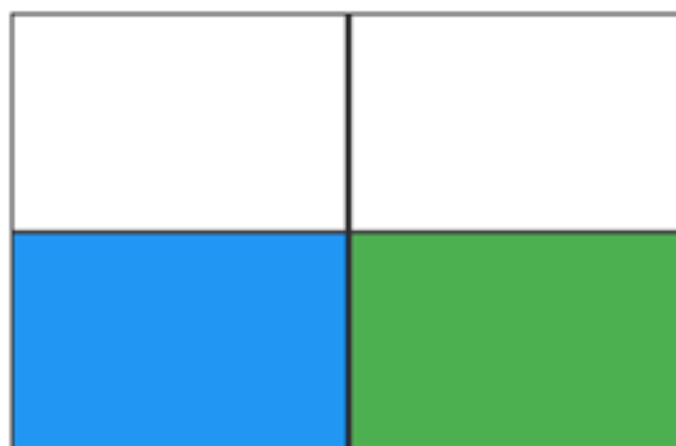
3.7: *Memrise*



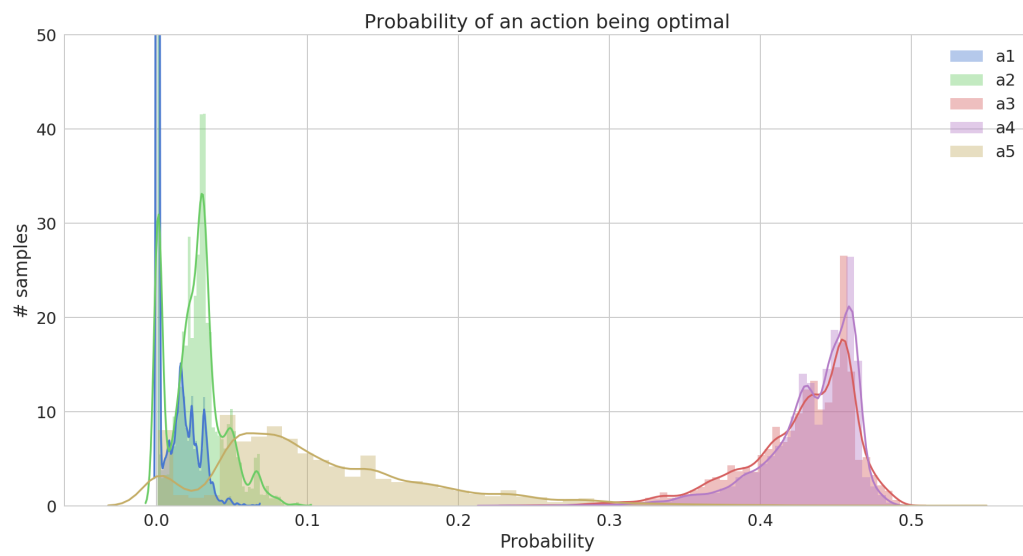
3.8: *Memrise*



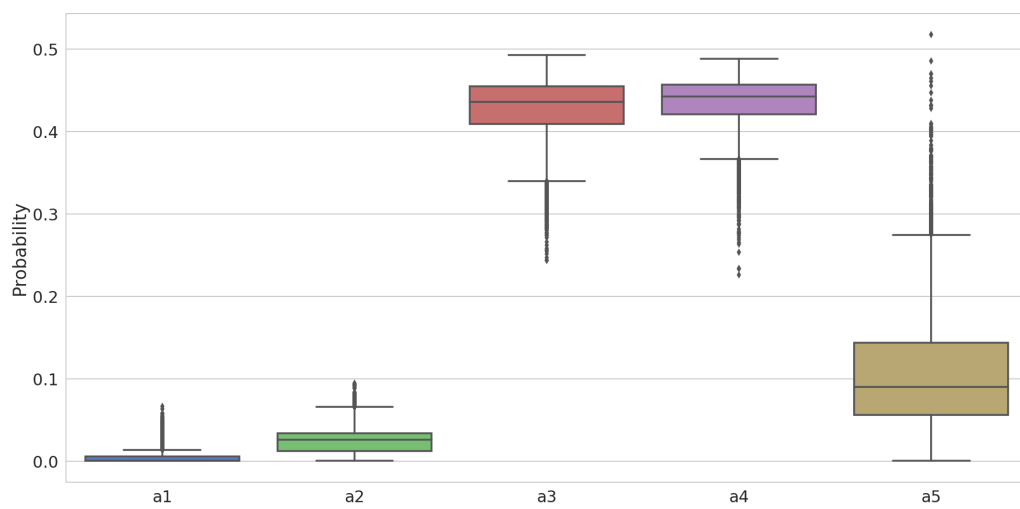
3.9:



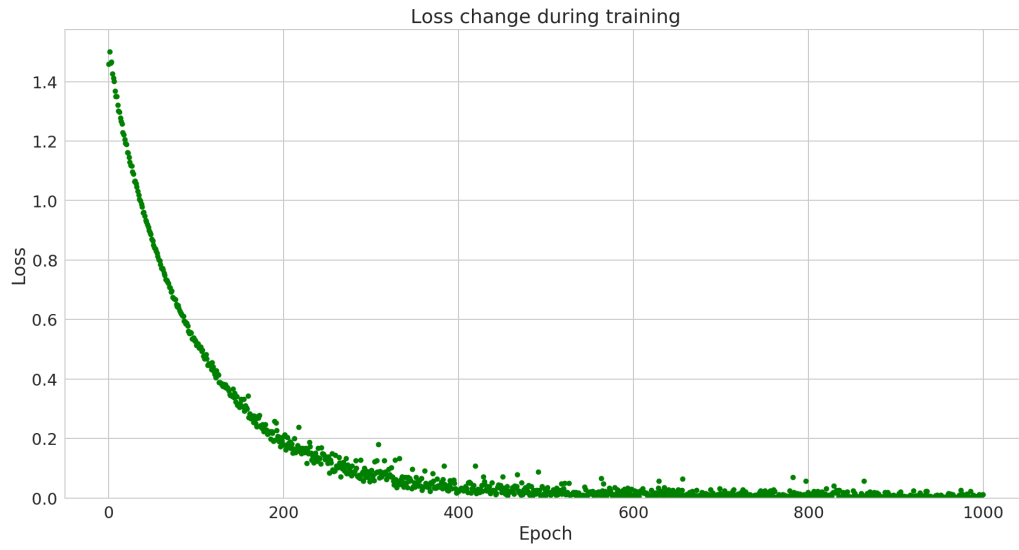
3.10:



3.11:

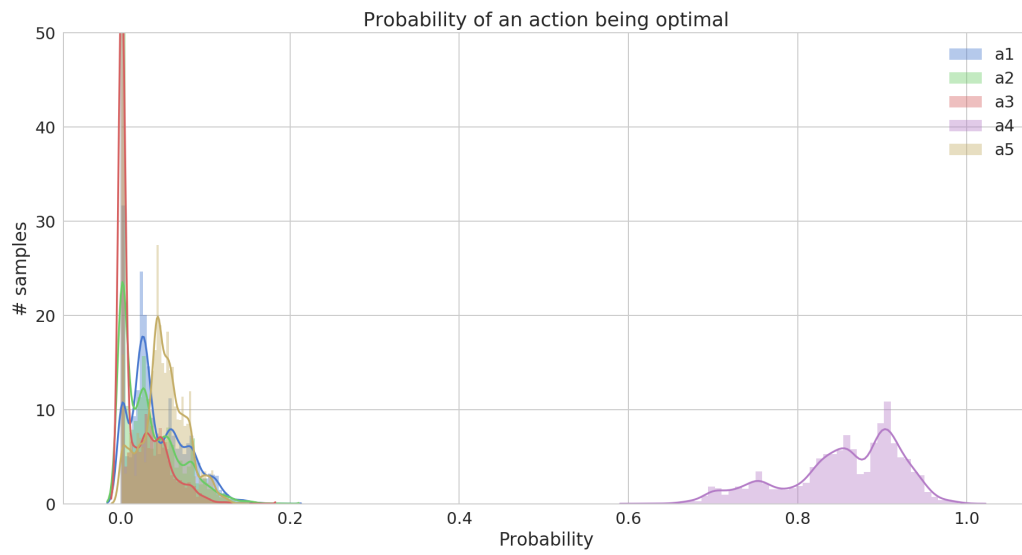


3.12:

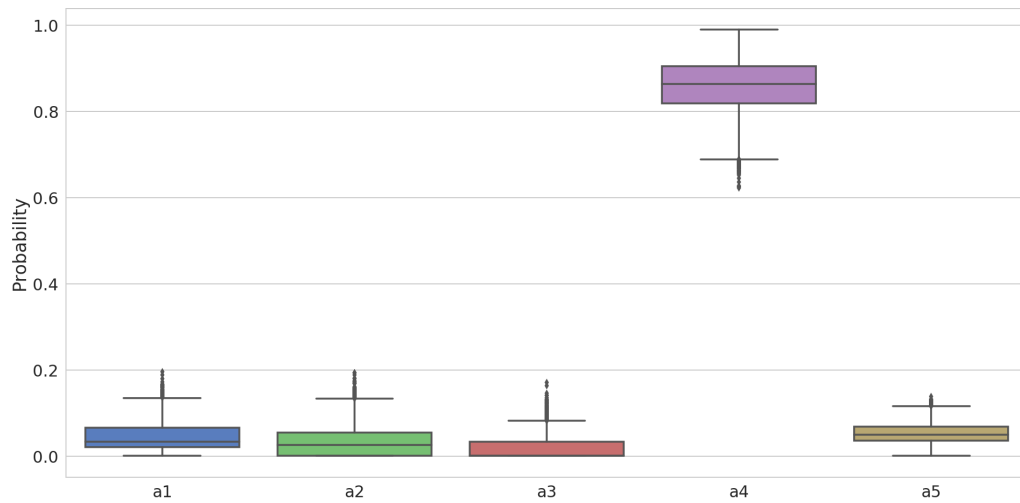


3.13:

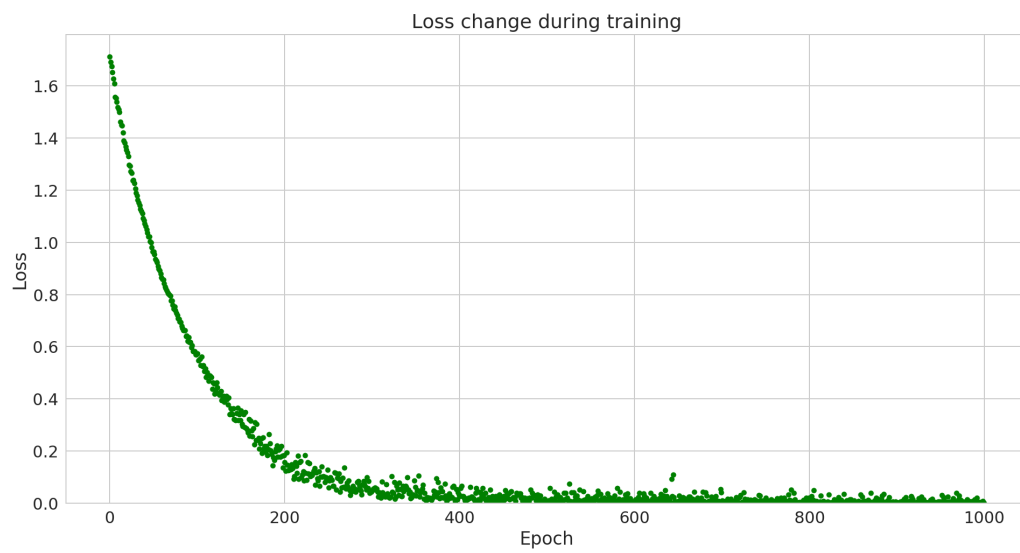
AutoMath



3.14:

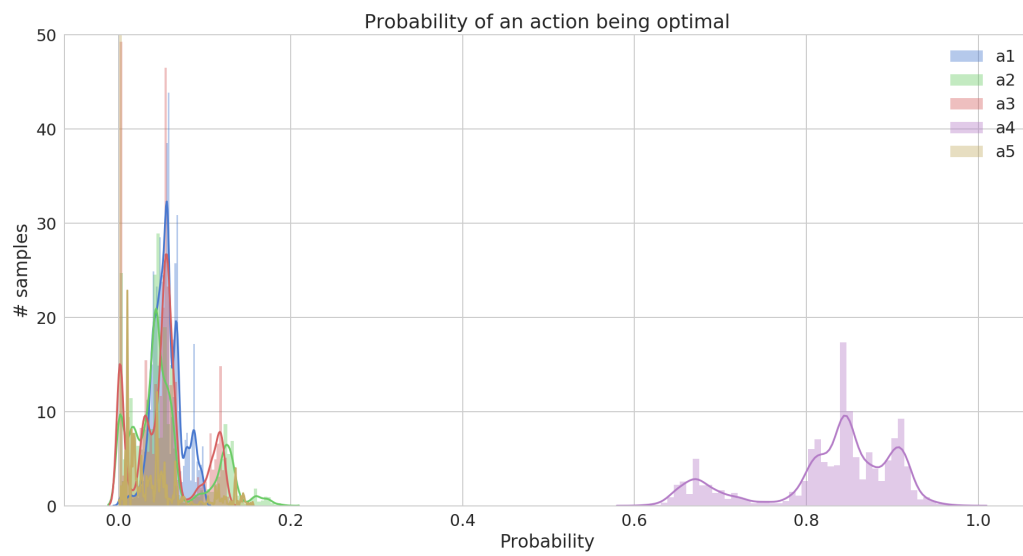


3.15:

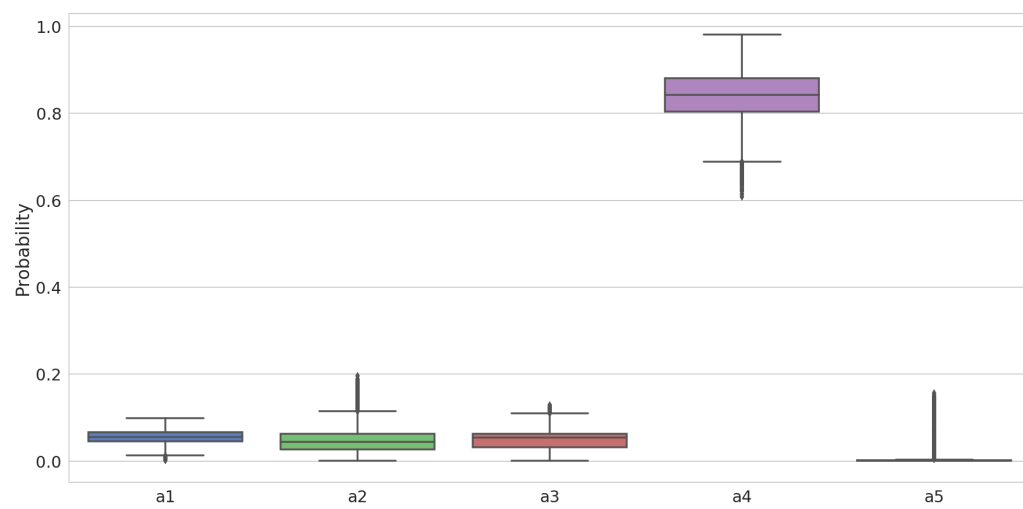


3.16:

Memrise



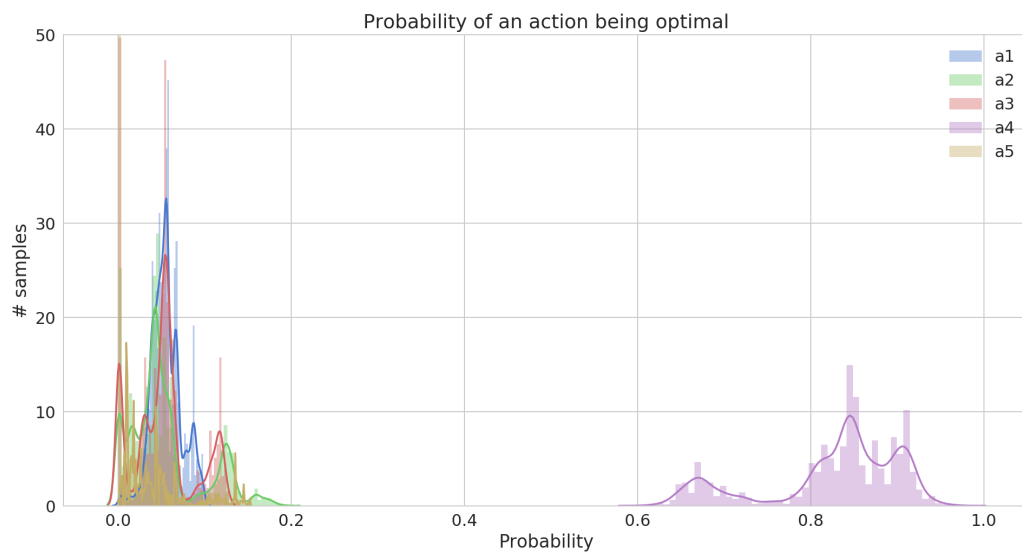
3.17:



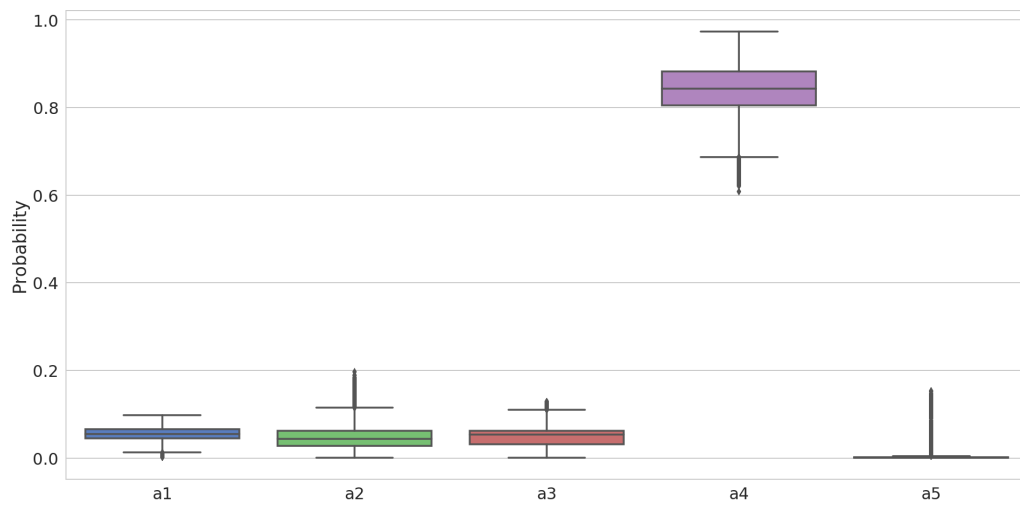
3.18:



3.19:



3.20:



3.21:

Глава 4

Среда за изучаване (RL exploration) на ГПИ приложения

Много от съществуващите системи за автоматизирано тестване на Android приложения се опитват да изградят решения, които взимат предвид недостатъците при тестване на приложения. Някои от трудностите повече не съществуват благодарение на напредъка на модерния компютърен хардуер, а други могат да бъдат решени много по-ефективно благодарение на нововведени инструменти за разработка за Android. Например, **SwiftHand** (Choi et al. 2013) се опитва да намали нуждата от преинсталиране на приложението върху устройството. В по-новите си версии, adb, предлага способ за изчистване на състоянието на дадено приложение, без нужда от преинсталирането му.

От особена важност за изграждане на алгоритъм в среда за подсиленото обучение е наличието на награда. Повечето от изградените системи се опитват да максимизират покритието на код. В практиката тази метрика е важна, но и недостатъчна. Фактът, че дадена част от програмния код се е изпълнила и не е предизвикала грешка в програмата не означава, че поведението на програмата е правилно (или не се е променило без това бъде желан ефект от разработчиците).

За нуждите на текущата работа и всеки желаещ да използва се предлага обща среда за изучаване/тестване на мобилни приложения. Поради липсата на други свободни инструменти (или такива, които са използвани). Системата е свобод-

на за използване, с отворен код и може да бъде намерена на <https://github.com/appgym/appgym>.

4.1 Android специфична среда

Средата се състои от два основни компонента - клиент и сървър. Сървърът работи върху Android устройството и предоставя данни за постигнатото покритие на код, изпълнение на действията, генерирани от модела и изображение за текущото състояние. Клиентът предоставя възможните действия на модела, както и комуникира със сървъра за да представи неговата функционалност.

Клиентът предоставя интерфейс към средата подобен на този на OpenAI gym (Brockman et al. 2016). Двата основни метода, които реализира са `reset()` и `step(action)`. `reset()` предоставя възможност на средата да се върне до първоначално състояние. Това се постига чрез спиране на приложението (ако то е стартирано), изтриване на данните поддържащи състоянието му, стартирането му и предоставяне на образ от екрана, както и възможните действия за състоянието. Изброената функционалност се реализира посредством `adb` команди и библиотеката `uiautomator` <https://github.com/xiacong/uiautomator>. Методът `step(action)` изпълнява избраното действие и предоставя новото състояние на средата, заедно с получената награда и новите възможни действия. Тук също се взима решение дали текущия епизод от обучението е приключил.

Множеството от възможните действия за текущото състояние се базират на броя и видовете графични елементи в него. Всеки елемент върху който може да се извърши докосване, задържане, скролиране, влачене и т.н. се превръща в действие. Множеството от графични елементи се извлича посредством библиотеката `uiautomator`.

Наградата за всяка избрана стъпка пряко се базира на покритието на код за текущия епизод. Стойността се изменя в интервала $[0; 1.0]$ и е нарастваща. Получаването на наградата след всяка стъпка е необходимо за обучение на модела. Скоростта на изпълнение пряко влияе на общото бързодействие на системата. За намиране на текущото покритие на код и изграждане на доклад се използва JaCoCo (Hoffmann et al. 2009). JaCoCo е интегриран в инструментите за

разработка на Android и се използва основно, когато е нужно покритие на код базирано на преминали тестове.

За целта на текущата работата бяха направени някои промени, които предоставят възможност за извличане на необходимите данни, докато програмата се изпълнява и не е в тестова среда. Няколко подхода бяха изпробвани за изграждане на крайните доклади. Първоначално бяха използвани adb команди и генериране на доклад посредством gradle задача. Бързодействието не беше задоволително - необходими бяха около 2 секунди на съвременен мобилен компютър. Около половината от времето се губеше в генериране на доклад, който предоставя повече от необходимата информация.

Крайното решение използва комбинация от HTTP сървър на устройството, клиент, специализиран начин за записване на данните за покритие на код и специализиран генератор за доклади. Допълнително бързодействие се постига чрез Nailgun <https://github.com/martylamb/nailgun> сървър, който изпълнява генератора за доклади. Така описаните оптимизации извършват необходимата работа за около 20 милисекунди (или 0.02 секунди) на същия компютър.

Взимането на текущото състояние се състои в направата на изображение на текущия екран на устройството. Тази задача отново се извършва от библиотеката `uiautomator`. Изображението може да бъде намалено до желани размери в зависимост от изискванията на задачата. Така полученото изображение се представя за текущо състояние под формата на тензор.

4.2 Web специфична среда

Web средата предоставя възможност за изучаване на мобилни web приложения. Тя използва библиотеката `puppeteer`, която предоставя автоматизиран достъп до Web Browser.

Системата AppGym е структурирана така че да може да се използва с други мобилни и настолни операционни системи. Интерфейсът е много сходен до този на OpenAI gym. Възможно е рамката да бъде променена, така че да бъде използвана за други задачи и други цели.

Глава 5

Изучаване на ГПИ среди

В тази глава изграждаме подход, който ни позволява да изучаваме среда, чието състояние се базира на изображения. По зададено изображение и множество от действия трябва да изберем действие, което максимизира изучената част от средата. По-конкретно, текущото състояние се определя от изображението на приложението и възможните действия с графичните елементи, а изучената част е моментното покритие на код.

Подходът, който използваме се базира изцяло на данните, които средата предоставя. В частност, разработваме модели, които представляват дълбока невронна мрежа, която приема изображения като входни данни и избира действие. Така създадения модел се тренира върху мобилни приложения за Android. Експерименталните резултати показват...

5.1 Related Work

Извличане на характеристики (features) от сензорна информация намалява нуждата от ръчното им закодиране и увеличава скоростта на изграждане на модели. Скоростните открития в областта на дълбокото самообучение доведоха до големи открития в компютърното зрение (Krizhevsky et al. 2012; Mnih 2013; Sermanet et al. 2013) и гласовото разпознаване (Dahl et al. 2012; Graves et al. 2013). Те използват различни архитектури за дълбоки невронни мрежи, вклю-

чително конволюции, многослойни персептрони (multilayer perceptrons) и рекурентни невронни мрежи. Използвани са в обучение с учител, без учител. Дълбока конволюционна невронна мрежа беше използвана за създаване на агент, който играе Atari 2600 компютърни игри (Mnih et al. 2013). За входни данни е използван видео вход с размери 84x84 пиксела и възможните действия. Работата използва повторение на преживяното (experience replay) (Lin 1992) и надгражда върху Neural Fitted Q-Learning (NFQ) (Riedmiller 2005). Още по-добри резултати бяха постигнати като се използва Монте Карло дървета за планиране, които бавно достигат извод, за обучение на дълбоки невронни мрежи, които са многократно по-бързи (Guo et al. 2014).

5.2 Дадено на агента

Предполагаме, че на стъпка t получаваме изображение и множество от възможни действия, определящи текущото състояние на средата. Искаме да създадем модел, който при подадени така описаните входни данни ни дава апостериорно разпределение, описващо вероятностите всяко от възможните действия да ни доведе до състояние с оптимално увеличение на покритието на програмен код.

Предизвикателството в така поставената задача се състои във факта, че изображенията са сложни многомерни обекти, а оптималното действие във всяко състояние може да е различно от това в което и да е друго. Нашият модел трябва да изгради вътрешно представяне на средата (напр. да научи какво е бутон, граници на отделните елементи и кои действия да използва върху тях) и да научи оптималните действия за всяко състояние.

1. Изображение на ГПИ средата - пример ГПИ. примери + изображения
2. Изучи ГПИ средата - да достигне всички възможни състояние на средата
3. Възможни действия -

Агентът щракване върху екрана (пример визуален) Агентът натиска и задържа (пример визуален) Агентът натиска и влачи курсора (пример визуален) Агентът scroll-ва нагоре и надолу

s_{x_1}	s_{x_2}	s_{x_3}	s_{x_4}	action	reward
b	w	g	g	a_3	0.25
b	w	w	b	a_5	0.00
b	w	g	g	a_4	0.25

5.3 Задачи

- Да постигне 100% покритие на програмния код на даден софтуерен продукт (СП)
- Да генерира поредица от действия с които преминава през всички клонове на дървото, описващо възможните състояния на СП

5.4 Модел на агента

Нека имаме среда E , намираща се в състояние $s \in \mathbb{S}$, върху което могат да бъдат изпълнени действия от множеството от действия \mathbb{A} . При избор на действие $a \in \mathbb{A}$, средата E предоставя награда r , която приема стойности в интервала $[0; 1]$ и преминава в ново състояние s' (в частност, $s' = s$, т.е. средата може да не премине в ново състояние). Множеството \mathbb{A} е ненаредено и всяко $a \in \mathbb{A}$ може да се обозначи с единствено цяло число, като по този начин въвеждаме наредба в \mathbb{A} . Всяко състояние на средата S позволява изпълнението на действия \mathbb{A} , които са предварително дефинирани. Множеството от всички възможни състояния на средата \mathbb{S} е неизвестно.

Нека след първоначално обучение от специалист имаме матрица на преходите D с размерност $n \times 3$, където n е броя на преходите. Всеки ред от D дефинира наредена тройка $(, ,)$, която описва получените награди при изпълнение на съответното действие за даденото състояние.

Нека имаме състояние s' , за което D не съдържа информация. В този случай, целта е да намерим подмножеството от действията, така че изпълнението им да води до получаване на оптимална награда от средата E .

Нека всяко състояние s е представено като изображение и имаме множество от действия A , които отговарят на различни графични елементи на екрана. Искаме да изберем действие a , което максимизира покритието на код. Предизвикателството се състои в намирането на възможно най-кратките поредици от действия, когато A може да е различно за всяко състояние s .

Подобно на (Mnih et al. 2013) създаваме дълбока конволюционна невронна мрежа, която използва изображения за входни данни. Директна работа с реалните размери на изображение взето от устройството може да изисква прекалено много изчисления (computationally demanding). Съвременните мобилни устройства достигат до 3840x2160 разделителна способност (Syny Xperia Z5 Premium). Въпреки това, физическият размер на екраните на смартфоните достигат до 5.5“. Това ги прави далеч по-малки от екраните на настолните и преносимите компютри. Основното взаимодействие с мобилните устройства се извършва чрез различни жестове (натискане, задържане, приплъзване и т.н.) с екрана. Предвид физическите размери на пръстите на човек, екраните не могат да съдържат голямо множество от елементи с които може да се взаимодейства.

5.4.1 Пространство на състоянието (State space)

Дадено състояние s на средата съдържа цветно изображение I и информация за сегментацията DOM модел D , т.е. $s = (I, D)$. Изображението има размер $W \times H \times 3$, където W е широчината на изображението в пиксели, H височината на изображението в пиксели и 3 - броя на цветовете в палитрата (червено, зелено и синьо (rgb)). DOM моделът е представен като списък от текстови елементи, а информацията за сегментацията на изображението е дадена от наредената четворка (x, y, w, h) x абцисната координата, y ординатната ос, w - широчината на сегмента, h - височина на сегмента.

5.4.2 Пространство на действията (Action space)

Позицията на курсора $m = (m_x, m_y) \in [0, W) \times [0, H)$ се моделира чрез мултиномиално разпределение върху възможните позиции. ГПИ средата не изисква наличие на клавиатура, защото . Възможните действия са: click, drag, scroll-

up, scroll-down.

5.4.3 Представяне на изображенията

Векторното представяне на изображения се базира на напредъка постигнат в компютърното зрение, където Конволюционните Невронни Мрежи (КНМ) са показали, че могат да превръщат сурови изображения в мощни представяния [cite], които позволяват създаването на модели, представящи се на човешко ниво в състезания като ImageNet classification challenge [cite]. КНМ може да се разглежда като функция $CNN_{\theta_c}(I)$, която извлича характеристики за изображение I и има параметри θ_c .

5.4.4 Модел за избор на действия

Агентът избира действие a във време t , когато се намира в състояние s . Решението на агента се взема благодарение на Дълбока Бейсова Невронна Мрежа (Deep Bayesian Neural Network)

5.4.5 Определяне на награди

Наградата r_t за всяка стъпка t се дефинира като:

$$r = -1 + C_a$$

където C_a е новият процент покритие на код след избора на действие a . “Наказанието”, количествено оценено с -1 за всяко следващо взето действие “мотивира” агента да се стреми изучи средата максимално бързо. Това спомага за намаляване на възможността за разглеждане на две съседни състояния в цикъл.

5.4.6 Архитектура на модела

- Входен слой - броят на невроните е равен на броят на пикселите в изображението, което представя средата в текущото време (обикновено 6400)
- 1-ви скрит слой - 2-мерен конволюционен слой с 3 входни канала и 16 изходни. Размерността на ядрото (kernel) е 5.
- 2-ри скрит слой - нормализиращ слой за 16-те изходни канала от предходния слой. ... Повтаряме предходните 2 слоя още 2 пъти
- 7-ми скрит слой - напълно свързан (fully-connected) слой, който сплесква (flatten?) броя на измеренията до 1.
- 8-ми скрит слой - отпадащ (dropout) слой, който изключва 20% от невроните в мрежата на всяка стъпка.
- Изходен слой - редуцира броя на невроните от предходния слой до броя на възможните действия на Агента

5.4.7 Цел/Оптимизация/Тренировка/Обучение

Целта на създадения модел е да научи апостериорното разпределение на действията в дадено състояние на средата s - формула.

Определение 1 (функция на загубата на Хубер). Казваме, че

Обучението ще се състои в минимизиране на т.нар функция на загубата на Хубер. Тя има свойствата на средна квадратична грешка, когато грешката е малка и тези на средна абсолютна грешка, когато грешката е голяма - това прави модела ни издръжлив (robust) на екстремни стойности (outliers) [cite].

5.4.8 Памет

Ще запазваме преходите от състояние s до s' при избрано действие a и получената награда r , които ще използваме за допълнително обучение на модела. Избирането на случайно подмножество от така запазените преходи

5.4.9 Оценка на модела за избор на действия

Адекватността на агента, т.е. адекватността на всички избрани от агента действия се измерва чрез т.нар мярка за “съжаление” (regret) - разликата между оптималната обща награда и получената обща награда. Оптималната награда може да се постигне, когато на всяка стъпка t , агентът избира оптимално действие a^* .

5.4.10 Намиране на апостериорно вероятностно разпределение на действията

Политиката, която използват агентите в [cite] е епсилон-лакома с намаляща стойност. Тя се отличава с това, че избира случайни действия, за да събере данни в началната фаза на обучението по-късно намаля вероятността за избиране на случайно действие. Обученият агент избира действие с максимална награда. Вместо това, може да опитаме да минимизираме несигурността (uncertainty) на нашата НМ. Това се оказва сравнително лесно ползвайки Thompson Sampling.

5.4.10.1 Thompson Sampling

Thompson sampling е политика, която насърчава агента да изследва средата в която действа, като избира действие с максимална награда, използвайки текущото си познание за средата. В нашият случай, това може да направим като симулираме стохастичен преход през НМ и изберем действието с най-висока очаквана награда.

5.5 Експерименти

Глава 6

Заключение

6.1 Нерешени проблеми

6.2 Бъдеща работа

6.3 Дискусия

Приложение 1: Някои важни вероятностни разпределения

Приложение 2: Фигури

Литература

Abraham, A. et al., 2015. GroddDroid: A gorilla for triggering malicious behaviors. In *Malicious and unwanted software (malware), 2015 10th international conference on*. IEEE, pp. 119–127.

Amalfitano, D. et al., 2015. MobiGUITAR: Automated model-based testing of mobile apps. *IEEE Software*, 32(5), pp.53–59.

Anonymous, 2018. Efficient exploration through bayesian deep q-networks. *International Conference on Learning Representations*. Available at: <https://openreview.net/forum?id=Bk6qQGWRb>.

Azizzadenesheli, K., Brunskill, E. & Anandkumar, A., 2018. Efficient Exploration through Bayesian Deep Q-Networks. *ArXiv e-prints*.

Bellemare, M.G., Dabney, W. & Munos, R., 2017. A distributional perspective on reinforcement learning. *arXiv preprint arXiv:1707.06887*.

Bellemare, M.G. et al., 2013. The arcade learning environment: An evaluation platform for general agents. *J. Artif. Intell. Res.(JAIR)*, 47, pp.253–279.

Bishop, C., 2007. Pattern recognition and machine learning (information science and statistics), 1st edn. 2006. Corr. 2nd printing edn. *Springer, New York*.

Bojarski, M. et al., 2016. End to end learning for self-driving cars. *arXiv preprint arXiv:1604.07316*.

Bojarski, M. et al., 2017. Explaining how a deep neural network trained with end-to-end learning steers a car. *arXiv preprint arXiv:1704.07911*.

Boyan, J.A., Littman, M.L. & others, 1994. Packet routing in dynamically changing networks: A reinforcement learning approach. *Advances in neural information processing systems*, pp.671–671.

Brockman, G. et al., 2016. OpenAI gym. *arXiv preprint arXiv:1606.01540*.

Cauchy, A., 1847. Méthode générale pour la résolution des systemes d'équations simultanées. *Comp. Rend. Sci. Paris*, 25(1847), pp.536–538.

Chang, T.-H., Yeh, T. & Miller, R.C., 2010. GUI testing using computer vision. In *Proceedings of the*

- sigchi conference on human factors in computing systems*. ACM, pp. 1535–1544.
- Cho, K. et al., 2014. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*.
- Choi, W., Necula, G. & Sen, K., 2013. Guided gui testing of android apps with minimal restart and approximate learning. In *ACM sigplan notices*. ACM, pp. 623–640.
- Choudhary, S.R., Gorla, A. & Orso, A., 2015. Automated test input generation for android: Are we there yet?(E). In *Automated software engineering (ase), 2015 30th ieee/acm international conference on*. IEEE, pp. 429–440.
- Crites, R.H. & Barto, A.G., 1996. Improving elevator performance using reinforcement learning. *Advances in neural information processing systems*, 8.
- Cybenko, G., 1989. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4), pp.303–314.
- Dahl, G.E. et al., 2012. Context-dependent pre-trained deep neural networks for large-vocabulary speech recognition. *IEEE Transactions on Audio, Speech, and Language Processing*, 20(1), pp.30–42.
- Deka, B. et al., 2017. Rico: A mobile app dataset for building data-driven design applications. In *Proceedings of the 30th annual acm symposium on user interface software and technology*. ACM, pp. 845–854.
- Gal, Y. & Ghahramani, Z., 2016. A theoretically grounded application of dropout in recurrent neural networks. In *Advances in neural information processing systems*. pp. 1019–1027.
- Gal, Y. & Ghahramani, Z., 2015. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. *arXiv preprint arXiv:1506.02142*, 2.
- Gauss, C.F., 1809. *Theoria motus corporum coelestium in sectionibus conicis solem ambientium auctore carolo friderico gauss*, sumtibus Frid. Perthes et IH Besser.
- Geman, S. & Geman, D., 1984. Stochastic relaxation, gibbs distributions, and the bayesian restoration of images. *IEEE Transactions on pattern analysis and machine intelligence*, (6), pp.721–741.
- Gergonne, J., 1815. Application de la méthode des moindres quarrés à l’interpolation des suites. *Annales de Math. Pures et Appl*, 6, pp.242–252.
- Graves, A., Mohamed, A.-r. & Hinton, G., 2013. Speech recognition with deep recurrent neural networks. In *Acoustics, speech and signal processing (icassp), 2013 ieee international conference on*. IEEE, pp. 6645–6649.
- Guo, X. et al., 2014. Deep learning for real-time atari game play using offline monte-carlo tree search planning. In *Advances in neural information processing systems*. pp. 3338–3346.

- Hahnloser, R.H. et al., 2000. Digital selection and analogue amplification coexist in a cortex-inspired silicon circuit. *Nature*, 405(6789), p.947.
- Hao, S. et al., 2014. Puma: Programmable ui-automation for large-scale dynamic analysis of mobile apps. In *Proceedings of the 12th annual international conference on mobile systems, applications, and services*. ACM, pp. 204–217.
- Hastings, W.K., 1970. Monte carlo sampling methods using markov chains and their applications. *Biometrika*, 57(1), pp.97–109.
- Hinton, G.E. et al., 2012. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*.
- Hochreiter, S. & Schmidhuber, J., 1997. Long short-term memory. *Neural computation*, 9(8), pp.1735–1780.
- Hoffmann, M. et al., 2009. Jacoco code coverage tool. Online, 2009.
- Hornik, K., Stinchcombe, M. & White, H., 1989. Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5), pp.359–366.
- Huang, S.-Y. et al., 2015. ABCA: Android black-box coverage analyzer of mobile app without source code. In *Progress in informatics and computing (pic), 2015 ieee international conference on*. IEEE, pp. 399–403.
- Huval, B. et al., 2015. An empirical evaluation of deep learning on highway driving. *arXiv preprint arXiv:1504.01716*.
- Ioffe, S. & Szegedy, C., 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*. pp. 448–456.
- Joachims, T. et al., 1997. Webwatcher: A tour guide for the world wide web. In *IJCAI (1)*. Citeseer, pp. 770–777.
- Johansson, F., Shalit, U. & Sontag, D., 2016. Learning representations for counterfactual inference. In *International conference on machine learning*. pp. 3020–3029.
- Kalchbrenner, N. & Blunsom, P., 2013. Recurrent continuous translation models. In *EMNLP*. p. 413.
- Kononenko, I., 1989. Bayesian neural networks. *Biological Cybernetics*, 61(5), pp.361–370.
- Krizhevsky, A., Sutskever, I. & Hinton, G.E., 2012. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*. pp. 1097–1105.
- Kullback, S. & Leibler, R.A., 1951. On information and sufficiency. *The annals of mathematical statistics*, 22(1), pp.79–86.
- LeCun, Y. et al., 1989. Backpropagation applied to handwritten zip code recognition. *Neural*

computation, 1(4), pp.541–551.

Lee, J.-G. et al., 2017. Deep learning in medical imaging: General overview. *Korean journal of radiology*, 18(4), pp.570–584.

Legendre, A.M., 1805. *Nouvelles méthodes pour la détermination des orbites des comètes*, F. Didot.

Levine, S. et al., 2016. End-to-end training of deep visuomotor policies. *Journal of Machine Learning Research*, 17(39), pp.1–40.

Lin, L.-J., 1992. Self-improving reactive agents based on reinforcement learning, planning and teaching. *Machine learning*, 8(3-4), pp.293–321.

Machiry, A., Tahiliani, R. & Naik, M., 2013. Dynodroid: An input generation system for android apps. In *Proceedings of the 2013 9th joint meeting on foundations of software engineering*. ACM, pp. 224–234.

MacKay, D.J., 1992a. A practical bayesian framework for backpropagation networks. *Neural computation*, 4(3), pp.448–472.

MacKay, D.J., 1992b. *Bayesian methods for adaptive models*. PhD thesis. California Institute of Technology.

Memon, A.M., 2002. GUI testing: Pitfalls and process. *Computer*, 35(8), pp.87–88.

Metropolis, N. et al., 1953. Equation of state calculations by fast computing machines. *The journal of chemical physics*, 21(6), pp.1087–1092.

Mnih, V., 2013. *Machine learning for aerial image labeling*. PhD thesis. University of Toronto.

Mnih, V. et al., 2016. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*.

Mnih, V. et al., 2013. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*.

Mnih, V. et al., 2015. Human-level control through deep reinforcement learning. *Nature*, 518(7540), pp.529–533.

Molnar, A.-J., 2015. Live visualization of gui application code coverage with guitracr. In *Software visualization (vissoft), 2015 ieee 3rd working conference on*. IEEE, pp. 185–189.

Moreira, R.M. & Paiva, A.C., 2014. A gui modeling dsl for pattern-based gui testing paradigm. In *Evaluation of novel approaches to software engineering (enase), 2014 international conference on*. IEEE, pp. 1–10.

Moreira, R. et al., 2017. Pattern-based gui testing: Bridging the gap between design and quality assurance.

Neal, R.M., 2012. *Bayesian learning for neural networks*, Springer Science & Business Media.

- Ohba, M., 1982. Software quality \approx test accuracy \times test coverage. In *Proceedings of the 6th international conference on software engineering*. IEEE Computer Society Press, pp. 287–293.
- Paszke, A. et al., 2017. PyTorch: Tensors and dynamic neural networks in python with strong gpu acceleration.
- Poplin, R. et al., 2018. Prediction of cardiovascular risk factors from retinal fundus photographs via deep learning. *Nature Biomedical Engineering*, p.1.
- Rajkomar, A. et al., 2018. Scalable and accurate deep learning for electronic health records. *arXiv preprint arXiv:1801.07860*.
- Riedmiller, M., 2005. Neural fitted q iteration—first experiences with a data efficient neural reinforcement learning method. In *European conference on machine learning*. Springer, pp. 317–328.
- Roubtsov, V. & others, 2005. Emma: A free java code coverage tool.
- Rumelhart, D.E., Hinton, G.E. & Williams, R.J., 1985. *Learning internal representations by error propagation*, DTIC Document.
- Salvesen, K. et al., 2015. Using dynamic symbolic execution to generate inputs in search-based gui testing. In *Proceedings of the eighth international workshop on search-based software testing*. IEEE Press, pp. 32–35.
- Sermanet, P. et al., 2013. Pedestrian detection with unsupervised multi-stage feature learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 3626–3633.
- Shi, T. et al., 2017. World of bits: An open-domain platform for web-based agents. In D. Precup & Y. W. Teh, eds. *Proceedings of the 34th international conference on machine learning*. Proceedings of machine learning research. International Convention Centre, Sydney, Australia: PMLR, pp. 3135–3144. Available at: <http://proceedings.mlr.press/v70/shi17a.html>.
- Silver, D. & Hassabis, D., 2016. AlphaGo: Mastering the ancient game of go with machine learning. *Research Blog*.
- Silver, D. et al., 2017. Mastering the game of go without human knowledge. *Nature*, 550(7676), pp.354–359.
- Srivastava, N. et al., 2014. Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1), pp.1929–1958.
- Sundermeyer, M., Schlüter, R. & Ney, H., 2012. LSTM neural networks for language modeling. In *Interspeech*. pp. 194–197.
- Sutton, R.S. & Barto, A.G., 1998. *Reinforcement learning: An introduction*, MIT press Cambridge.
- Tishby, N., Levin, E.D. & Solla, S.A., 1989. Consistent inference of probabilities in layered networks:

Predictions and generalizations. *International 1989 Joint Conference on Neural Networks*, pp.403–409 vol.2.

Todorov, E., Erez, T. & Tassa, Y., 2012. MuJoCo: A physics engine for model-based control. In *Intelligent robots and systems (iros), 2012 ieee/rsj international conference on*. IEEE, pp. 5026–5033.

Tsitsiklis, J.N., Van Roy, B. & others, 1997. An analysis of temporal-difference learning with function approximation. *IEEE transactions on automatic control*, 42(5), pp.674–690.

Van Hasselt, H., Guez, A. & Silver, D., 2016. Deep reinforcement learning with double q-learning. In *AAAI*. pp. 2094–2100.

Werbos, P.J., 1988. Generalization of backpropagation with application to a recurrent gas market model. *Neural networks*, 1(4), pp.339–356.

Yeh, C.-C. & Huang, S.-K., 2015. CovDroid: A black-box testing coverage system for android. In *Computer software and applications conference (compsac), 2015 ieee 39th annual*. IEEE, pp. 447–452.

Zhauniarovich, Y. et al., 2015. Towards black box testing of android apps. In *Availability, reliability and security (ares), 2015 10th international conference on*. IEEE, pp. 501–510.

Zhu, H., Hall, P.A. & May, J.H., 1997. Software unit test coverage and adequacy. *Acm computing surveys (csur)*, 29(4), pp.366–427.