

Нейронные сети

Лекция 1. Введение в машинное обучение

Бабин Д.Н., **Иванов И.Е.**, Петюшко А.А.

MaTIC

13 сентября 2021г.



- 1 Организационные вопросы
- 2 Постановка основных задач машинного обучения
- 3 Тестирование моделей, выбор лучшей
- 4 Декомпозиция ошибки, недообучение и переобучение





Руководитель курса: д.ф.-м.н. Бабин Дмитрий Николаевич



Лектор: к.ф.-м.н. Иванов Илья Евгеньевич



Лектор: к.ф.-м.н. Петюшко Александр Александрович



- Авторы имеют более 15 лет опыта участия в проектах, связанных с машинным обучением и компьютерным зрением
- Являются постоянными участниками группы распознавания образов кафедры МаТИС
- В качестве научных консультантов работают или работали с такими крупнейшими российскими и международными компаниями как Нейроком, LSI Research, Fotonation, Huawei и др.



- В данный момент времени авторы ведут исследования в области компьютерного зрения в московском научно-исследовательском центре Хуавэй
- Данный курс является гибридом нескольких курсов программы **SHARE**
 - **SHARE** = School of Huawei Advanced Research Education
 - **SHARE** = Школа опережающего научного образования Хуавэй
- Приглашаем всех желающих принять участие в более продвинутых курсах SHARE, по окончании которых реально устроиться на стажировку в Huawei



Почему стоит уделить время этому курсу

- 1 Это возможность получить знания, которые пригодятся в работе



Почему стоит уделить время этому курсу

- 1 Это возможность получить знания, которые пригодятся в работе
- 2 Специалисты по компьютерному зрению/машинному обучению сейчас очень востребованы



Почему стоит уделить время этому курсу

- 1 Это возможность получить знания, которые пригодятся в работе
- 2 Специалисты по компьютерному зрению/машинному обучению сейчас очень востребованы
- 3 Это шанс максимально использовать своё образование



Почему стоит уделить время этому курсу

- 1 Это возможность получить знания, которые пригодятся в работе
- 2 Специалисты по компьютерному зрению/машинному обучению сейчас очень востребованы
- 3 Это шанс максимально использовать своё образование
- 4 Это просто очень интересно и затягивает



Что же такое искусственный интеллект?

Естественный интеллект (человек)

- Может мыслить, принимать решения, анализировать информацию



Что же такое искусственный интеллект?

Естественный интеллект (человек)

- Может мыслить, принимать решения, анализировать информацию

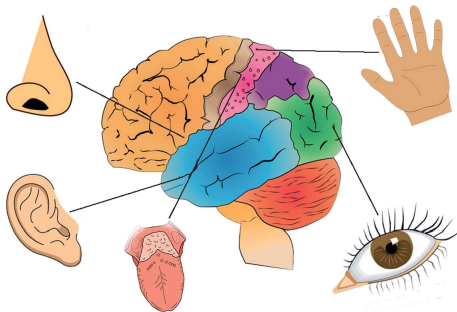
Искусственный интеллект

- (в широком смысле) то же самое, что и естественный, только с использованием компьютера вместо человека
- (в узком смысле) алгоритмы способные сами обучаться, чтобы выполнять задачи вместо человека



Взаимодействие со средой

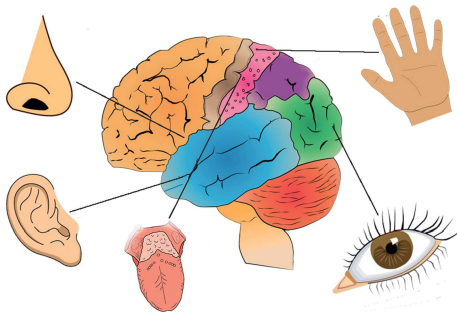
- Около 90 % информации поступает через **зрение**¹
- Около 9 % информации поступает через слух



¹https://www.rlsnet.ru/books_book_id_2_page_40.htm

Взаимодействие со средой

- Около 90 % информации поступает через **зрение**¹
- Около 9 % информации поступает через слух



Вывод

Чтобы построить интеллектуальную систему, её необходимо научить взаимодействовать со средой

¹https://www.rlsnet.ru/books_book_id_2_page_40.htm

1 Машинное обучение

- Необходимые основы

2 Введение в современные нейронные сети

- Методы построения и обучения современных нейронных сетей

3 Приложения нейронных сетей: компьютерное зрение

- Извлечение информации из визуальных образов (изображений и видео)



- Оценки за курс будут выставляться в соответствии с набранными баллами за выполнение домашних заданий.
- По курсу будут предложены домашние задания трёх видов:
 - теоретические
 - практические
 - соревнования
 - дополнительные
- В конце второй недели состоится экзамен, на котором при желании можно будет повысить свою оценку
- Предварительная шкала оценок:

Оценка	Процент выполненных заданий
Отлично	80 %
Хорошо	60 %
Зачет	40 %



- Оценки 'хорошо' и 'отлично' можно получить, выполняя задания в течение семестра
- Если по каким-то причинам вам необходимо перенести дедлайн, то об этом необходимо сообщить заранее
- Для тех, кто в течение семестра не набрал баллов на 'удовлетворительно', необходимо будет сдавать зачет по первым двум частям курса



- Оценки 'хорошо' и 'отлично' можно получить, выполняя задания в течение семестра
- Если по каким-то причинам вам необходимо перенести дедлайн, то об этом необходимо сообщить заранее
- Для тех, кто в течение семестра не набрал баллов на 'удовлетворительно', необходимо будет сдавать зачет по первым двум частям курса



- Списывать категорически запрещается!

- Списывать категорически запрещается!
- При подозрении на списанную работу ставится 0 баллов:
 - Списавшему
 - Давшему списать



- Списывать категорически запрещается!
- При подозрении на списанную работу ставится 0 баллов:
 - Списавшему
 - Давшему списать
- При использовании дополнительных источников (ресурсы в Интернете, учебники) обязательно ссылаться на них



- Страница курса: TBD
 - Телеграмм-канал: <https://t.me/joinchat/NiezPdBwWEF10DMy>
 - Группа обсуждения: TBD
 - Почта курса: TBD
- Пожалуйста в теме указывайте [Ташкент-2021]



- Предсказание стоимости недвижимости
- Предсказание платёжеспособности клиента
- Предсказание оттока клиентов
- Классификация заболеваний
- Предсказание клика пользователя по рекламному баннеру
- И многие другие задачи. . .



- Видеонаблюдение: нахождение и слежка за объектами
- Медицина: обнаружение злокачественных опухолей
- Улучшение качества фотографий
- Зрение роботов
- Поиск по фотографиям
- Дополненная реальность
- Автомобили без водителей
- Магазины без продавцов



Способы машинного обучения

Определения

- X — множество объектов
- Y — множество ответов
- $y : X \rightarrow Y$ — неизвестная зависимость

Основные способы машинного обучения

- **С учителем**
 - Достаточное количество обучающего материала, т.е. пар (x_i, y_i)
- **Частичное обучение**
 - Малое количество размеченных данных и много неразмеченных примеров x_i
- **Без учителя**
 - Нет размеченных пар, только примеры x_i
- **С подкреплением**
 - Формирование отклика на основе взаимодействия со средой

Постановка задачи обучения с учителем

- Дано:
 - $\{(x_1, y_1), \dots, (x_n, y_n)\} \subset X \times Y$ — обучающая выборка
- Найти
 - Решающую функцию $a : X \rightarrow Y$, которая приближает целевую зависимость y .
- Необходимо детализировать:
 - Как определяются объекты
 - Как задаются ответы
 - Что значит, что одна зависимость приближает другую



Как определяются объекты

Определение

Объект = совокупность признаков

Типы признаков

- Бинарный признак
- Категориальный признак
- Порядковый признак
- Количественный признак



Задачи классификации

- Бинарная классификация $Y = \{-1, 1\}$ или $Y = \{0, 1\}$
- Многоклассовая классификация $Y = \{0, 1, \dots, M - 1\}$
- Многозначная бинарная классификация $Y = \{0, 1\}^M$

Задачи восстановления регрессии

$Y = \mathbb{R}$ или $Y = \mathbb{R}^n$



Функция потерь

Определение

Функция потерь (loss function) $\mathcal{L}(a, x)$ — величина ошибка алгоритма a на объекте x

Функции потерь для задачи классификации

$\mathcal{L}(a, x) = [a(x) \neq y]$ — индикатор ошибки

Функции потерь для задач регрессии

$\mathcal{L}(a, x) = (a(x) - y)^2$ — квадратичная ошибка



Сравнение моделей машинного обучения

Как понять, что одна модель лучше другой?

Для этого используют независимое от **обучающего** множества множество, которое называется **тестовым**

Зачем вообще это понимать?

- Существует множество алгоритмов машинного обучения и важно понимать, какой из них более применим в конкретной задаче
- Даже в рамках одной модели есть много параметров



Как выбрать лучшую модель

Наивный подход

Обучить модели с различными параметрами и посмотреть, что будет на тесте

Минусы наивного подхода

- Так как тест состоит из случайной выборки теста, то результат на теста тоже является некоторым приближение случайной величины
- Если все модели тестировать на тестовом датасете и выбирать лучшую таким образом, то будет происходить неявное обучение на тесте и на другом независимом тесте возможны сюрпризы

Что же делать?

Чтобы неявно не обучиться на тестовых данных – надо использовать кросс-валидацию

Общая идея

Основная идея кросс-валидации состоит в разбиении обучающего множества на два непересекающихся множества (возможно многократном):

$$X^{learn} = X^{train} \sqcup X^{val}$$

На одном из них происходит обучение, а на другом происходит валидация модели.

Частные случаи

- 1 Простейшая кросс-валидация (Holdout) — однократное разделение множества



Частные случаи

2 k-fold валидация^а:

Split 1	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 1
Split 2	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 2
Split 3	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 3
Split 4	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 4
Split 5	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 5

Training data Test data

- 3 Leave one out (LOO) валидация — частный случай k-fold валидации если k равно мощности обучающего множества
- 4 Многократная k-fold валидация — повторение k-fold валидации несколько раз с разными разбиениями.

^а<https://towardsdatascience.com>

Переобучение

Определение

Переобучение (overfitting) — нежелательное явление, возникающее при решении задач обучения по прецедентам, когда вероятность ошибки обученного алгоритма на объектах тестовой выборки оказывается существенно выше, чем средняя ошибка на обучающей выборке. Переобучение возникает при использовании избыточно сложной модели

Причины возникновения

Одной из причин переобучения избыточная сложность пространства параметров модели, "лишние" степени свободы используются для точной подгонки на обучающую выборку

Методы обнаружения

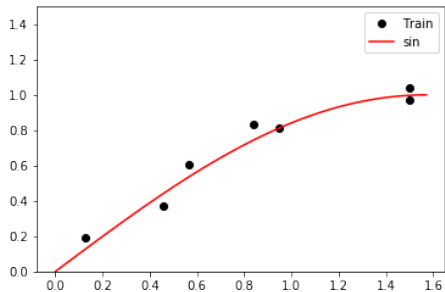
Основным методом обнаружения является использование кросс-валидации

Определение

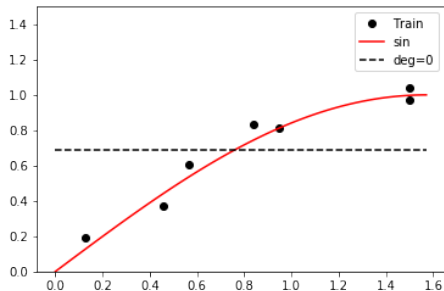
Недообучение (underfitting) – нежелательное явление, возникающее при решении задач обучения по прецедентам, когда алгоритм обучения не обеспечивает достаточно малой величины средней ошибки на обучающей выборке. Недообучение возникает при использовании недостаточно сложных моделей



Примеры недообучения и переобучения

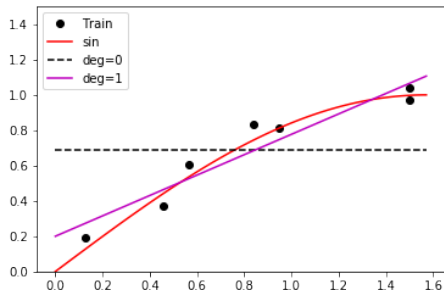


Примеры недообучения и переобучения



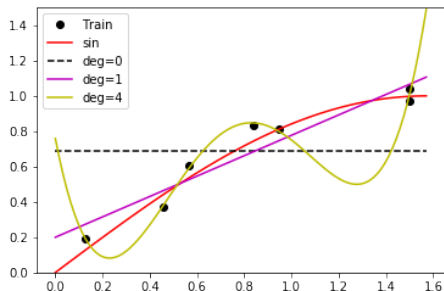
- Полином нулевой степени не может хорошо приближать зависимость в силу ограниченности параметров модели

Примеры недообучения и переобучения



- Полином нулевой степени не может хорошо приближать зависимость в силу ограниченности параметров модели
- Линейная и квадратичная модели адекватно описывают закономерность

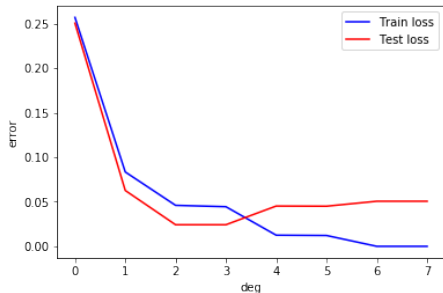
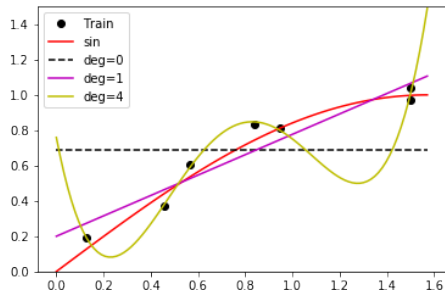
Примеры недообучения и переобучения



- Полином нулевой степени не может хорошо приближать зависимость в силу ограниченности параметров модели
- Линейная и квадратичная модели адекватно описывают закономерность
- Полиномы высоких степеней могут в точности пройти через точки обучающей выборки



Примеры недообучения и переобучения



Определения

Пусть $y = y(x) = f(x) + \varepsilon$, $\varepsilon \sim N(0, \sigma^2)$ — целевая зависимость, и $a(x)$ — алгоритм машинного обучения.

Разложение квадрата ошибки

$$\begin{aligned} E(y - a)^2 &= E(y^2 + a^2 - 2ya) = Ey^2 + Ea^2 - 2Eya = \\ &= Ey^2 + Ea^2 - 2E(f + \varepsilon))a = Ey^2 + Ea^2 - 2Efa = \\ &= Ey^2 - (Ey)^2 + (Ey)^2 + Ea^2 - (Ea)^2 + (Ea)^2 - 2fEa = \\ &= Dy + Da + (Ey)^2 + (Ea)^2 - 2fEa = \\ &= Dy + Da + (E(f - a))^2 = \sigma^2 + \text{variance}(a) + \text{bias}^2(f, a) \end{aligned}$$



Определение

Разброс (variance) – дисперсия ответов алгоритмов $a(x)$.
Характеризует разнообразие алгоритмов (из-за случайности обучающей выборки, шума, стохастичности обучения и т.д.)

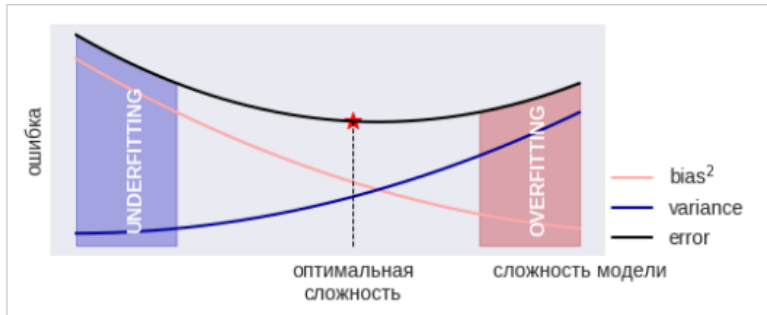
Определение

Смещение (bias) – матожидание разности между истинным ответом и выбранным алгоритмом. Характеризует способность модели настраиваться на целевую зависимость



Модель оптимальной сложности

- Для простых моделей характерно недообучение
- Для сложных моделей характерно переобучение
- Оптимальная сложность модели где-то между²



²<https://dyakonov.org>

- Две основные задачи машинного обучения — классификация и регрессия



- Две основные задачи машинного обучения — классификация и регрессия
- Для любой задачи машинного обучения очень важно контролировать и избегать переобучение. Именно для этого и нужна кросс-валидация



- Две основные задачи машинного обучения — классификация и регрессия
- Для любой задачи машинного обучения очень важно контролировать и избегать переобучение. Именно для этого и нужна кросс-валидация
- Следствие теоремы о декомпозиции ошибки: не всегда более сложная модель даёт результаты лучше



- Две основные задачи машинного обучения — классификация и регрессия
- Для любой задачи машинного обучения очень важно контролировать и избегать переобучение. Именно для этого и нужна кросс-валидация
- Следствие теоремы о декомпозиции ошибки: не всегда более сложная модель даёт результаты лучше



Спасибо за внимание!

