

Фамилия, имя, номер группы:

.....

Внесите сюда ответы на тест:

Вопрос	1	2	3	4	5	6	7	8	9	10
Ответ										

Табличка для проверяющих работу:

Задачи	1	2	3	4	5	6	7	Итого

Позитивная мотивация — явно не мой конёк, и мы все умрём.
 Всё уже было до нас, можно выдохнуть страх, и уставить глаза в небосклон.
 Если этой контрольной и сопротивляться, то не с печальным лицом.
 Всё повторится не раз, но мы живы сейчас, нами рано удобрять чернозём.
Охххутирон и Тося Чайкина про мидтёрм по ML (2022)

Работа состоит из трёх частей: тестовая, задачи и ответы на открытые вопросы. Списывание карается обнулением работы. Удачи!

Часть первая: тестовая

Дайте ответ на 10 тестовых вопросов. Каждый вопрос стоит 3 балла. Никакие дополнительные пояснений в этой части работы от вас не требуются.

Вопрос 1. Архимед хочет отбить римскую атаку на Сиракузы. С помощью беспилотника Архимед делает съёмки римских кораблей, а затем специальным алгоритмом компьютерного зрения пытается предсказать, сколько римских солдат находится на каком корабле. Какую задачу решает Архимед?

- | | | |
|---|---|---|
| <input type="checkbox"/> A Регрессия | <input type="checkbox"/> C Кластеризация | <input type="checkbox"/> E Рекомендации |
| <input type="checkbox"/> B Классификация | <input type="checkbox"/> D Ранжирование | <input type="checkbox"/> F Нет верного ответа. |

Вопрос 2. Что такое гиперпараметр?

- | | | |
|---|--|--|
| <input type="checkbox"/> A Параметр, который оказывает ключевое влияние на производительность модели | <input type="checkbox"/> C Параметр, чьё оптимальное значение нельзя подобрать по обучающей выборке | <input type="checkbox"/> E Параметр, от которого выход модели зависит нелинейно |
| <input type="checkbox"/> B Ровно то же самое, что и параметр | <input type="checkbox"/> D Параметр, который является многомерным | <input type="checkbox"/> F Нет верного ответа. |

Вопрос 3. В чём потенциальный недостаток кросс-валидации по двум блокам?

- | | | |
|---|--|--|
| <p><input type="checkbox"/> A Оценка качества модели на новых данных может оказаться очень заниженной из-за того, что тестирование проводится на слишком маленькой выборке</p> | <p><input type="checkbox"/> C Оценка качества модели на новых данных может оказаться очень заниженной из-за того, что обучение проводится на выборке сильно больше исходной</p> | <p><input type="checkbox"/> E Нужно обучать слишком много моделей, вычислительных мощностей для этого может не хватить, новые сервера в Россию из-за санкций не поставляют, поэтому кросс-валидация по двум блокам не имеет никакого смысла</p> |
| <p><input type="checkbox"/> B Оценка качества модели на новых данных может оказаться очень заниженной из-за того, что тестирование проводится на слишком большой выборке</p> | <p><input type="checkbox"/> D Оценка качества модели на новых данных может оказаться очень заниженной из-за того, что обучение проводится на выборке сильно меньше исходной</p> | <p><input type="checkbox"/> F Нет верного ответа.</p> |

Вопрос 4. Что из этого формула для шага в градиентном спуске?

- | | | |
|--|--|--|
| <p><input type="checkbox"/> A $w_t = w_{t-1} - \eta \cdot \nabla L(w_t)$</p> | <p><input type="checkbox"/> C $w_t = w_{t-1} - \eta \cdot \nabla L(w_{t-1})$</p> | <p><input type="checkbox"/> E $w_t = w_{t-1} + \eta \cdot \nabla L(w_{t-1})$</p> |
| <p><input type="checkbox"/> B $w_t = w_{t-1} + \eta \cdot \nabla L(w_t)$</p> | <p><input type="checkbox"/> D $w_t = w_{t-1} - \eta \cdot \nabla L(w_0)$</p> | <p><input type="checkbox"/> F Нет верного ответа.</p> |

Вопрос 5. Какие из способов приведённых ниже можно использовать для работы с пропусками в категориальных переменных при обучении линейных моделей?

- | | | |
|--|---|---|
| <p><input type="checkbox"/> A Если пропусков очень много, выкинуть переменную</p> | <p><input type="checkbox"/> C Выделить пропуски в отдельную категорию и сделать ONE-преобразование</p> | <p><input type="checkbox"/> E Заполнить пропуски медианами, посчитанными по каждой колонке</p> |
| <p><input type="checkbox"/> B Заполнить пропуски аномальным значением</p> | <p><input type="checkbox"/> D Заполнить нулями</p> | <p><input type="checkbox"/> F Нет верного ответа.</p> |

Вопрос 6. Бог плодородия Дионис спустился с Олимпа вкушать вина свежего урожая. Зевс пытается понять, сколько дней будет идти кутёж Диониса. Для этого он использует линейную регрессию, обученную на предыдущих кутежах:

$$y_i = 7 + 0.5 \cdot x_1 + 0.2 \cdot x_2,$$

где x_1 — качество вина по десятибалльной шкале, x_2 — количество собутульников. Выберите все верные утверждения об этой модели.

- | | | |
|---|---|--|
| <p><input type="checkbox"/> A Каждые дополнительные 5 собутульников будут затягивать кутёж на день</p> | <p><input type="checkbox"/> C Кутёж затянется минимум на 14 дней</p> | <p><input type="checkbox"/> E Каждый дополнительный собутульник будет затягивать кутёж на один день</p> |
| <p><input type="checkbox"/> B In vino veritas, in aqua sanitas</p> | <p><input type="checkbox"/> D Кутёж затянется минимум на 7 дней</p> | <p><input type="checkbox"/> F Нет верного ответа.</p> |

Вопрос 7. Что из этого можно использовать для регуляризации? Под регуляризацией мы понимаем штрафование моделей за сильно отличающиеся от нуля веса.

☐ A $|w_1| + |w_2| + \dots + |w_d|$

☐ C $w_1^2 + w_2^2 + \dots + w_d^2$

☐ E $|w_1|^3 + |w_2|^3 + \dots + |w_d|^3$

☐ B $\frac{1}{|w_1|} + \frac{1}{|w_2|} + \dots + \frac{1}{|w_d|}$

☐ D $w_1^3 + w_2^3 + \dots + w_d^3$

☐ F Нет верного ответа.

Вопрос 8. У нас есть 2 класса. Классификатор предсказывает, что объект равновероятно относится к каждому из них. Какое значение принимает logloss на этом объекте?

☐ A $-\log 2$

☐ C $-2 \log 2$

☐ E $\log 2$

☐ B $-0.5 \log 2$

☐ D $-2 \log 0.5$

☐ F Нет верного ответа.

Вопрос 9. Леонид предсказывает цены на квартиры в Спарте с помощью метода ближайших соседей. При построении предсказания он хочет учитывать расстояние до соседей. Какая из формул ниже поможет Леониду корректно построить прогноз? Все суммы ищутся по ближайшим соседям, y_j — цена квартиры, ρ_j — расстояние от объекта для которого строится предсказание до соответствующего соседа.

☐ A $\frac{1}{k} \sum_{j=1}^k y_j$

☐ C $\frac{1}{k} \sum_{j=1}^k \frac{y_j}{\rho_j}$

☐ E $\frac{\sum_{j=1}^k \frac{1}{\rho_j} \cdot y_j}{\sum_{j=1}^k \rho_j}$

☐ B $\frac{\sum_{j=1}^k \rho_j \cdot y_j}{\sum_{j=1}^k \rho_j}$

☐ D $\frac{\sum_{j=1}^k \rho_j \cdot y_j}{\sum_{j=1}^k \frac{1}{\rho_j}}$

☐ F Нет верного ответа.

Вопрос 10. Какие из метрик перечисленных ниже используются для решения задачи классификации?

☐ A MSE

☐ C MAPE

☐ E ROC-AUC

☐ B f-мера

☐ D logloss

☐ F Нет верного ответа.

Часть вторая: открытые вопросы

Эта часть состоит из открытых вопросов. На них необходимо дать краткие, но ёмкие ответы. За каждый ответ вы можете получить 10 баллов.

Вопрос 11. Предложите для каждой из перечисленных ниже задач, как сформулировать их в терминах машинного обучения: укажите, что будет являться объектом и целевой переменной, а также напишите тип задачи.

1. В заповеднике голод — белки часто остаются некормленными, так как зайцы оказываются слишком прожорливыми и съедают чужую еду. Мы приняли решение поставить кормушки с фотоловушками, которые бы открывались для того животного, которое подошло к кормушке, чтобы ограничить ненасытных зайцев.
2. По статистике некоторая доля сотрудников одной компании ежедневно опаздывает из-за пробок, сокращая свой рабочий день на некоторое количество времени. Мы решили разработать приложение, которое будет подсказывать сотрудникам, на какое время им стоит запланировать выход из дома, чтобы не опоздать на работу.

Вопрос 12. Аполоний обожает логистическую регрессию. Поэтому он решил закодить её самостоятельно, без всяких пакетов. Он записал модель следующим образом:

$$b(x_i) = \sigma(b + \langle x_i, w \rangle)$$
$$\sigma(t) = \frac{1}{1 - \exp(t)}$$

Для обучения Аполоний использует функцию потерь

$$L(b(x_i), y_i) = b(x_i) \cdot \ln y_i + (1 - b(x_i)) \cdot \ln(1 - y_i)$$

Какие ошибки вы тут видите? Для каждой объясните, к каким последствиям и почему она приведёт, а также как это исправить.

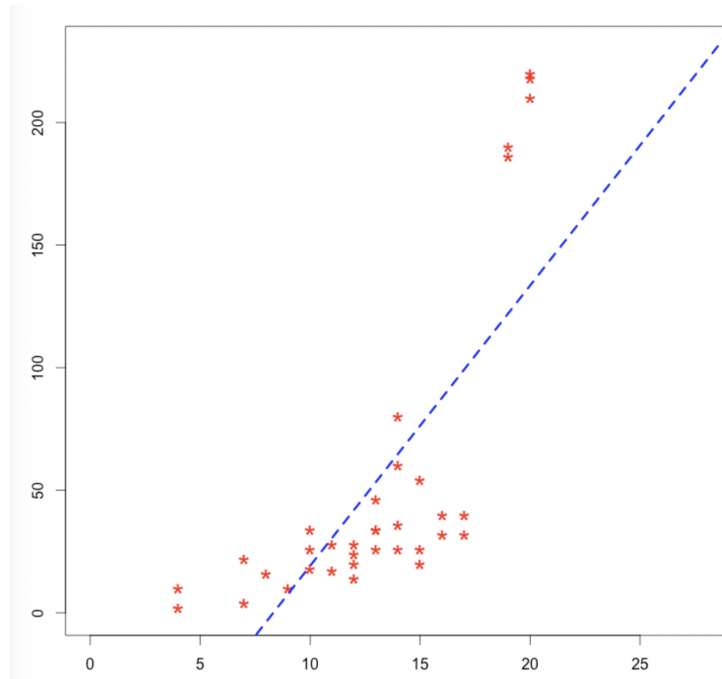
Вопрос 13. Объясните мем

Вопрос 14. Пигмалион оживил Галатею, вырезанную из слоновой кости! Теперь они хотят обучить метод ближайших соседей для классификации горных пород. В качестве признаков используются: вес пород, размер, цвет и тп.

Галатея хочет строить свои прогнозы с учётом того, каким получилось расстояние между объектами. Пигмалион хочет строить свои прогнозы с помощью метода большинства. Выпишите формулы, которые будут использованы для расчёта прогнозов. Объясните в них каждую компоненту и обозначения.

Вопрос 15. Аргонавты плавают вокруг нимф и пытаются понять, сколько их нужно слушать, чтобы весь экипаж зачаровало. По собранным данным оценивается линейная регрессия. В качестве объясняющей переменной, x используется длина песни, в качестве объясняемой число зачарованных аргонавтов, y . В качестве функции потерь используется MSE.

Ясон отложил по оси x длину песни, а по оси y число зачарованных аргонавтов. Поверх облака точек он нарисовал линейную регрессию. Получилась такая картинка:



С какими проблемами при обучении модели столкнулся Ясон? Предложите как минимум два способа исправить возникшую проблему.

Часть третья: задачи

Решите все задания. Все ответы должны быть обоснованы. Решения должны быть прописаны для каждого пункта. Рисунки должны быть чёткими и понятными. Все линии должны быть подписаны. За решение каждой задачи вы можете получить 10 баллов.

Вопрос 16. Мойры предсказывают судьбу. Клото использует для этого метод ближайших соседей. Лахеси использует линейную регрессию, а Атропо случайный лес. В тестовой выборке у них есть три наблюдения y_i . Для каждого из них мойры построили прогнозы.

настоящие y_i	1	2	3
KNN	2	3	1
линейная регрессия	2	3	4
случайный лес	1	1	1

1. Найдите для прогнозов MAE, MSE, RMSE и MAPE.
2. Объясните, зачем от MSE обычно переходят к RMSE.
3. Объясните, почему MAE считается более устойчивой к выбросам.

Вопрос 17. Константин основал Константинополь, а затем решил задачу классификации. У него получились следующие прогнозы.

y_i	\hat{p}_i
1	0.9
1	0.1
0	0.75
1	0.56
0	0.2
0	0.37
0	0.25

1. Бинаризируйте ответ по порогу t и посчитайте точность и полноту для $t = 0.3$ и для $t = 0.8$.
2. Постройте ROC-кривую и найдите площадь под ней.