

# Классификация и CNN

---

Маша Шеянова, [masha.shejanova@gmail.com](mailto:masha.shejanova@gmail.com)

# Задачи NLP в МО

---

# Как можно?

- на уровне текстов
- на уровне предложений
- на уровне слов
- speech, OCR, image captioning

# На уровне текстов

Классификация текста:

- spam detection
- жанры
- sentiment analysis (тональность)
- предсказание темы

Кластеризация текстов:

- НОВОСТИ
- topic modelling

# Предложения

- Paraphrase
- Textual entailment
- QA systems
- machine translation

# Слова

- POS-tagging
- named entity recognition

(и то и другое решается как задача sequence tagging)

# Dropout

---

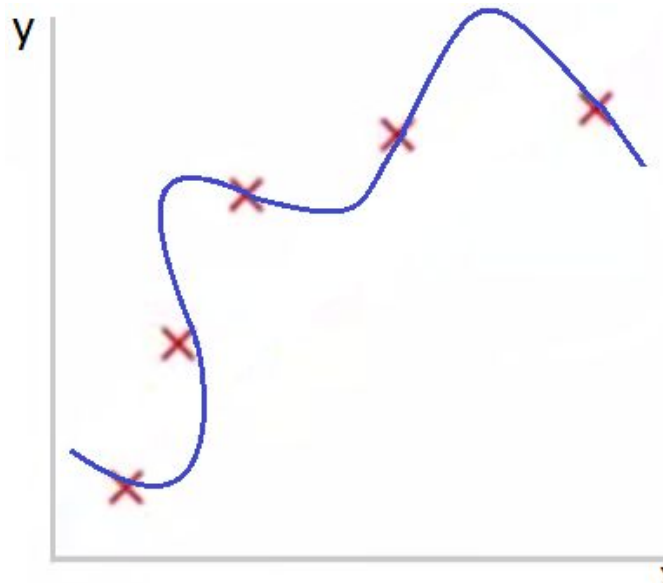
# Что такое переобучение?

Переобучение означает, что модель обращает слишком много внимания на незначительные признаки, чтобы “слишком хорошо” подстроиться под данные на обучающей выборке.

При этом модель становится излишне сложной (посмотрите на эти ненужные загогулины).

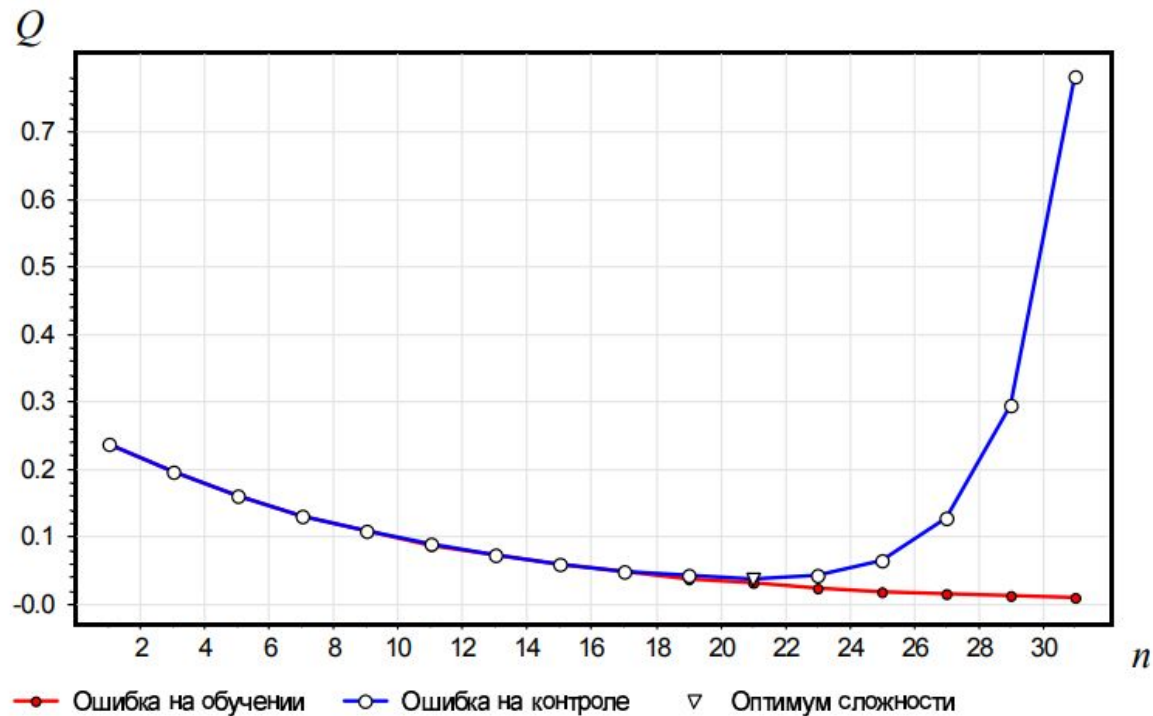
Скорее всего, она покажет плохой результат на тестовой выборке.

Как понять, что модель переобучилась?





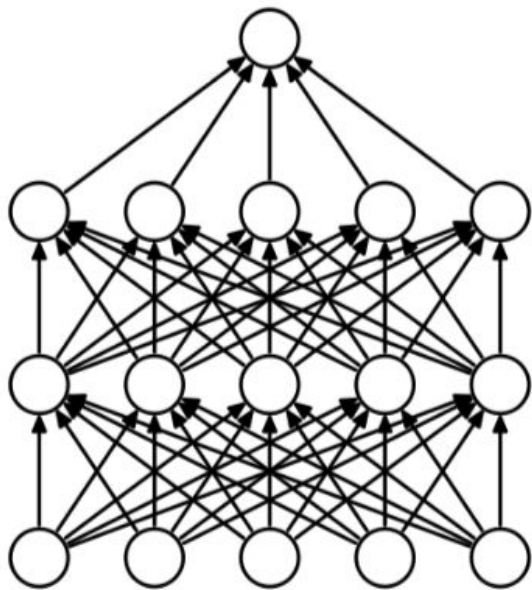
# Переобучение — это когда



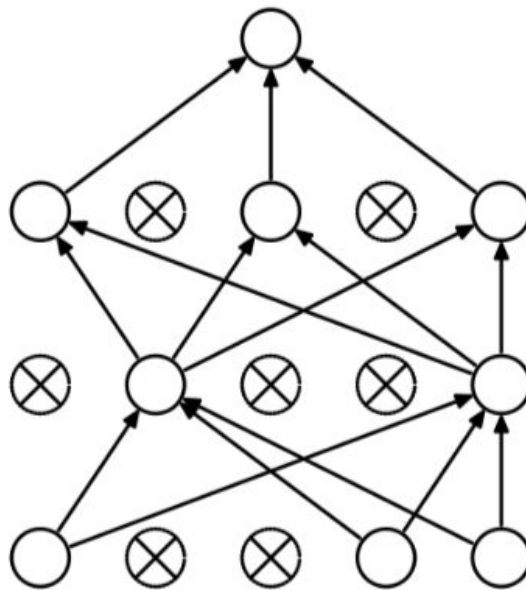
Источник картинки.

Иными словами, если качество модели на обучении продолжает расти, а на тесте — падает, модель переобучилась.

# Dropout



(a) Standard Neural Net



(b) After applying dropout.

CNN

---

# Свёртка (convolution)

1	1 <sub>x1</sub>	1 <sub>x0</sub>	0 <sub>x1</sub>	0
0	1 <sub>x0</sub>	1 <sub>x1</sub>	1 <sub>x0</sub>	0
0	0 <sub>x1</sub>	1 <sub>x0</sub>	1 <sub>x1</sub>	1
0	0	1	1	0
0	1	1	0	0

Image

4	3	

Convolved  
Feature

- проходимся по матрице скользящим окном (его называют ядро, kernel)
- перемножаем каждое попавшее в окно число матрицы на числа ядра
- складываем результат и записываем в ячейку новой матрицы

[анимация](#)

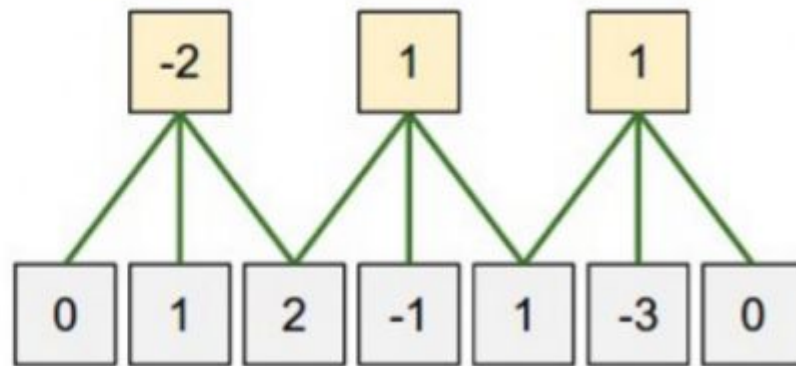
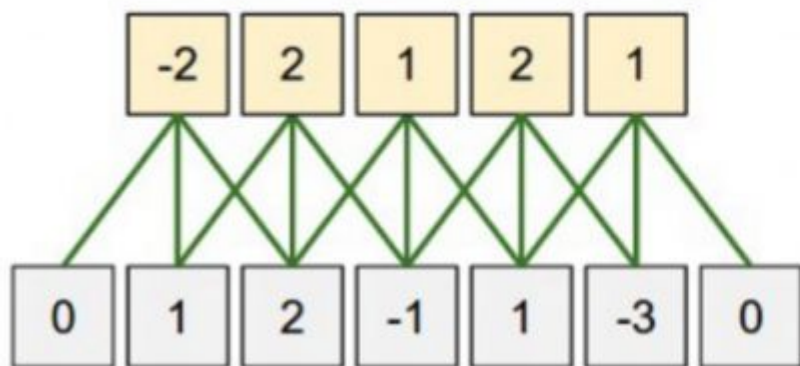
# CNN

- идея за свёрткой — из сырых фичей извлечь полезные признаки
- вообще говоря, числа в ядре могут быть фиксированными
  - например, все единички
- однако, тогда придётся заранее их выбирать
- универсальный оптимальный вариант выбрать нельзя
- но можно обучить нейросеть подбирать цифры в ядре самостоятельно!

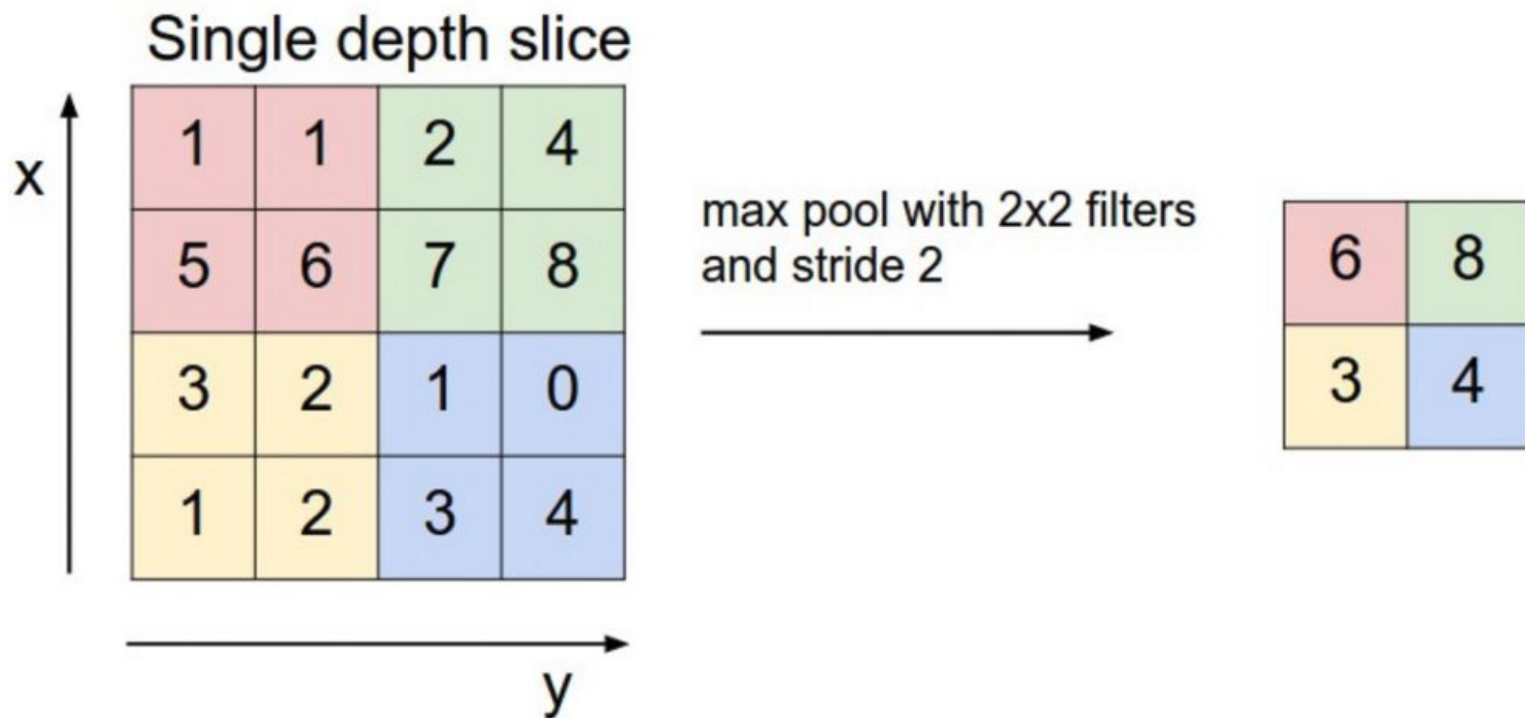
# Stride

Можно проходиться, как на картинке выше, каждый раз сдвигаясь на 1.

Но можно сдвигаться сильнее и делать меньшее перекрытие.



# Max pooling



# CNN для классификации предложений

