

# Лекция 4. Оптимизация. Регуляризация

Денис Деркач, Дмитрий Тарасов

Слайды от А. Маевского, М. Гущина, А. Кленицкого, М Борисяка

02 февраля 2026 года



# Функции потерь



# Оценка максимального правдоподобия

Модель  $f(x, \theta)$  определяет параметры распределения вероятности  $P(y|x)$ :

$$P(y|x) = P(y|f(x, \theta))$$

Функция правдоподобия:

$$L(\theta) = P(D|\theta) = \prod_{i=1}^N P(y_i|f(x_i, \theta))$$

Метод максимального правдоподобия (maximum likelihood estimation):

$$\theta = \operatorname{argmax}_{\theta} L(\theta) = \operatorname{argmax}_{\theta} \left[ \prod_{i=1}^N P(y_i|f(x_i, \theta)) \right]$$

# Оценка максимального правдоподобия

$$\begin{aligned}\theta &= \operatorname{argmax}_{\theta} \left[ \log \left( \prod_{i=1}^N P(y_i | f(x_i, \theta)) \right) \right] \\ &= \operatorname{argmax}_{\theta} \left[ \sum_{i=1}^N \log (P(y_i | f(x_i, \theta))) \right] \\ &= \operatorname{argmin}_{\theta} \left[ - \sum_{i=1}^N \log (P(y_i | f(x_i, \theta))) \right]\end{aligned}$$

Negative log-likelihood loss:

$$\mathcal{L} = - \sum_{i=1}^N \log (P(y_i | f(x_i, \theta)))$$

# Ординарная регрессия

Предположим, что шум распределен нормально с центром в нуле:

$$y = f(x, \theta) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

то есть

$$\begin{aligned} P(y|x, \theta, \sigma) &= \mathcal{N}(y|f(x, \theta), \sigma^2) \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{(y - f(x, \theta))^2}{2\sigma^2} \right] \end{aligned}$$

$$\begin{aligned} \mathcal{L} &= -\sum_{i=1}^N \log \left( \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{(y - f(x, \theta))^2}{2\sigma^2} \right] \right) \\ &= \sum_{i=1}^N \left[ -\log \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) + \frac{(y - f(x, \theta))^2}{2\sigma^2} \right] \rightarrow (y - f(x, \theta))^2 \end{aligned}$$

Получили обычный метод наименьших квадратов

# Обобщённые линейные модели и другие

- Можем предсказывать и среднее, и дисперсию:

$$\theta = \operatorname{argmin}_{\theta} \sum_{i=1}^N \left[ -\log \left( \frac{1}{\sqrt{2\pi f_2(x, \theta)^2}} \right) + \frac{(y - f_1(x, \theta))^2}{2f_2(x, \theta)^2} \right]$$

- Если использовать распределение Лапласа, то получим MAE
- Если использовать распределение Пуассона, то можем предсказывать счетчики событий
- Если использовать бета-распределение, то можем предсказывать пропорции

# Бинарная классификация

Распределение Бернулли

$$P(y|\lambda) = \lambda^y(1 - \lambda)^{1-y} = \begin{cases} \lambda, & y = 1 \\ 1 - \lambda, & y = 0 \end{cases}$$

Чтобы оценивать  $\lambda$  - вероятность от 0 до 1, используем сигмоиду:

$$f(x, \theta) = \sigma(x, \theta)$$

$$P(y|f(x, \theta)) = f(x, \theta)^y(1 - f(x, \theta))^{1-y}$$

Функция потерь - binary cross-entropy

$$\mathcal{L}(\theta) = \sum_{i=1}^N [-y_i \log f(x_i, \theta) - (1 - y_i) \log(1 - f(x_i, \theta))]$$

# Многоклассовая классификация

Softmax:

$$\text{softmax}_k(z) = \frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}}$$

$$P(y = k | f(x, \theta)) = \text{softmax}_k [f(x, \theta)]$$

Функция потерь - cross-entropy (aka softmax loss)

$$\mathcal{L}(\theta) = - \sum_{i=1}^N \log (\text{softmax}_{y_i} [f(x, \theta)])$$



# Softmax поведение на границах

$$\text{softmax}_k(z) = \frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}}$$

Из-за численного переполнения

- При больших отрицательных значениях  $z$  можем получить 0 в знаменателе
- При больших положительных значениях  $z$  можем получить бесконечность

Трюк:

$$\frac{e^{z_k+c}}{\sum_{j=1}^K e^{z_j+c}} = \frac{e^{z_k} e^c}{\sum_{j=1}^K e^{z_j} e^c} = \frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}}$$

$$z_k \rightarrow z_k - \max_j z_j$$

# Другие функции потерь

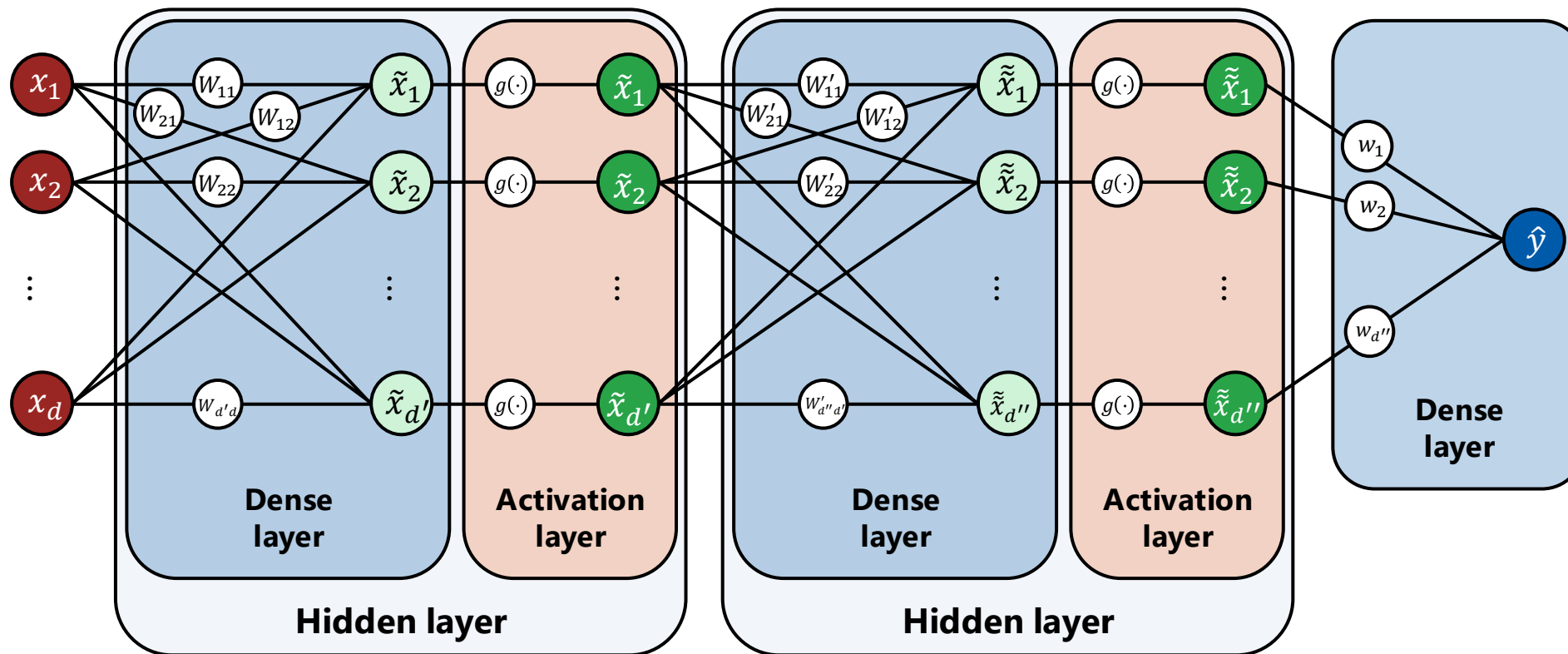
- Hinge loss
- Focal loss
- Dice loss
- Triplet loss
- Contrastive loss

И многие другие..

# Оптимизация



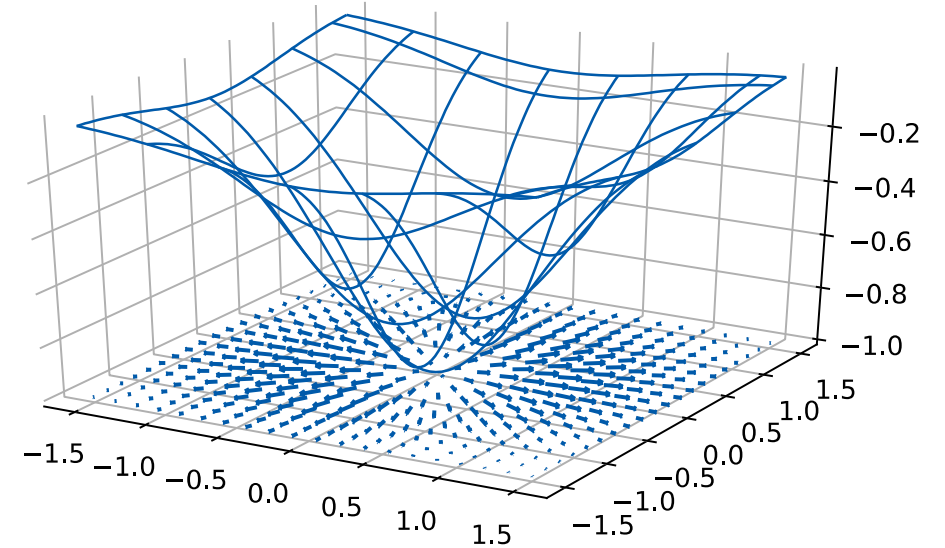
# Нейронная сеть как функция



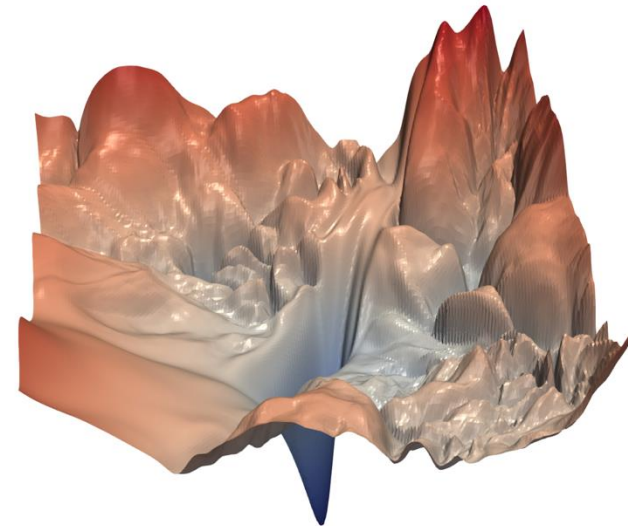
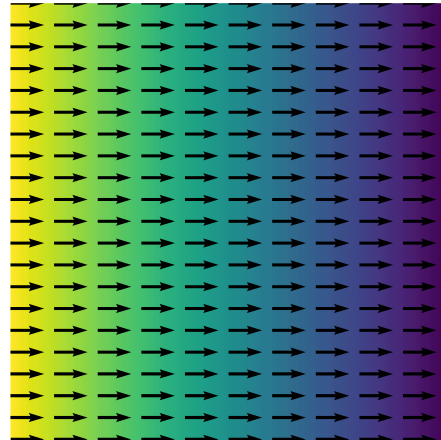
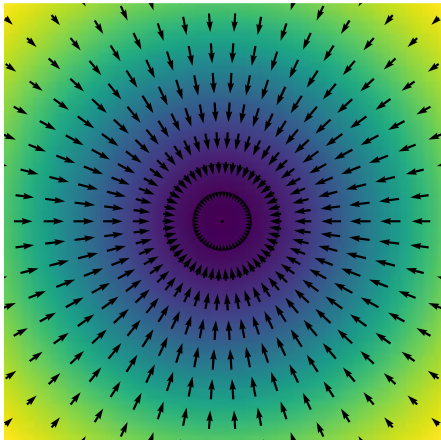
$$\hat{y} = W_{out} \cdot \dots \cdot W_{h2} \cdot W_{h1}x$$

# Градиент

- ▶ Градиент:  $\nabla_x f(x) \equiv \left( \frac{\partial f(x)}{\partial x_1}, \dots, \frac{\partial f(x)}{\partial x_d} \right)$
- ▶ Указывает на самый быстрый рост функции



NB: Neural NetLoss function landscape  
(in 2 random directions around an optimum)



# Немного интуиции

- ▶ Для простоты давайте пока опустим функции активации.
- ▶ Тогда выход нейронной сети, состоящей только из плотных слоев, будет:

$$\hat{y} = W_{out} \cdot \dots \cdot W_{h2} \cdot W_{h1}x$$

- ▶ Обратите внимание, что градиент относительно любой из весовых матриц  $W_{hk}$  пропорционален **произведению** всех других матриц
- ▶ Например, для матриц  $1 \times 1$ , если все они имеют масштаб  $S \in \mathbb{R}$ , градиент  $g$ :

$$g \sim S^{m-1},$$

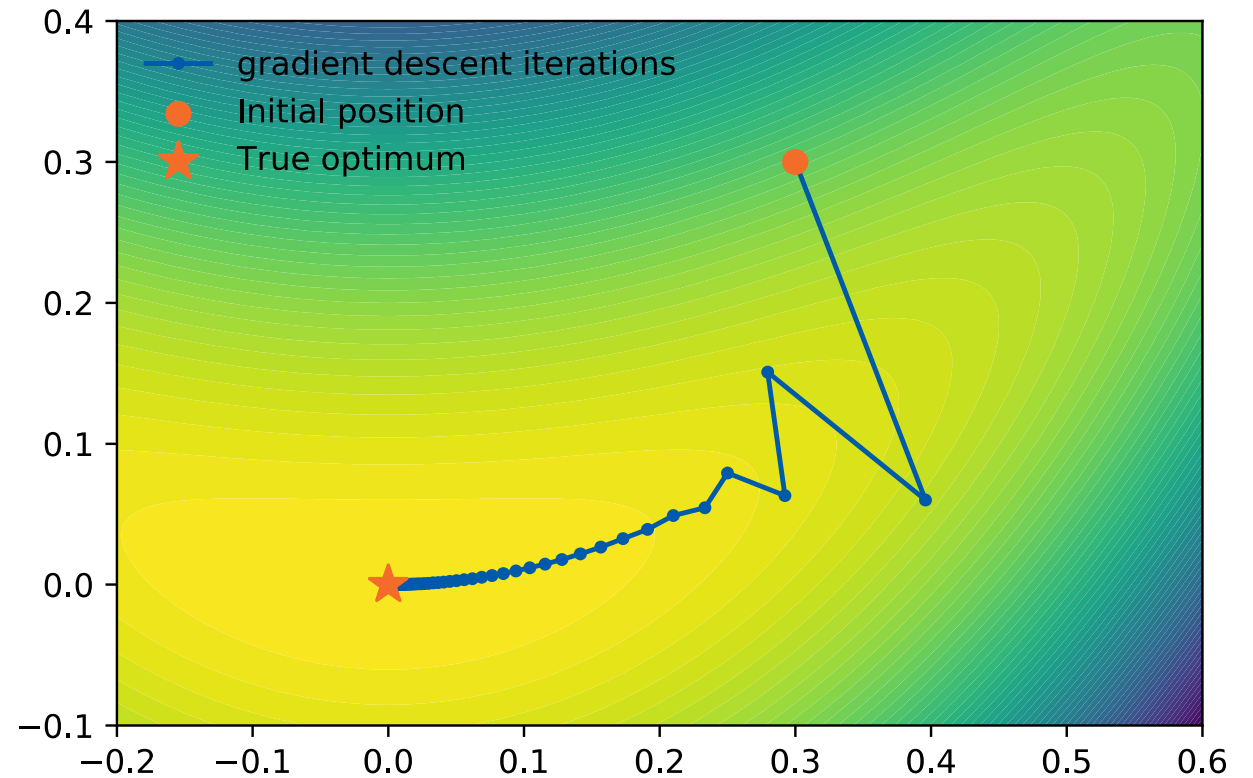
где  $m$  — глубина сети

- ▶ Если  $S$  слишком большое, градиенты **взорвутся**;
- ▶ Если  $S$  слишком маленькое, они **исчезнут**.

# Оптимизация градиентного спуска

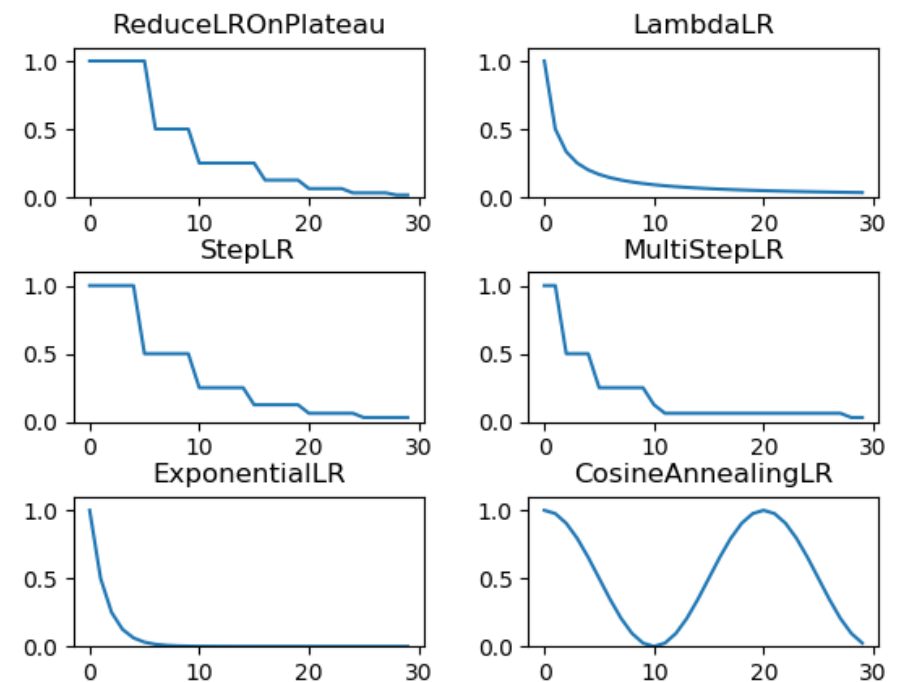
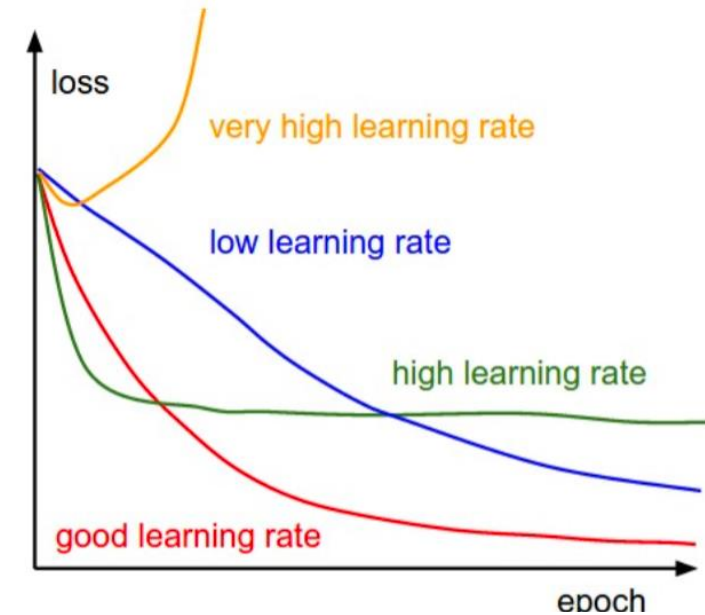
- ▶ Можно оптимизировать функции, начиная с некоторой начальной точки  $x^{(0)}$  и двигаясь против градиента
$$x^{(k)} \leftarrow x^{(k-1)} - \alpha \nabla_x f(x^{(k-1)})$$
где  $\alpha \in \mathbb{R}$ ,  $\alpha > 0$  – **темп обучения** (learning rate).
- ▶ Для гладких выпуклых функций с единственным минимумом  $x^*$ :

$$f(x^{(k)}) - f(x^*) = \mathcal{O}\left(\frac{1}{k}\right)$$



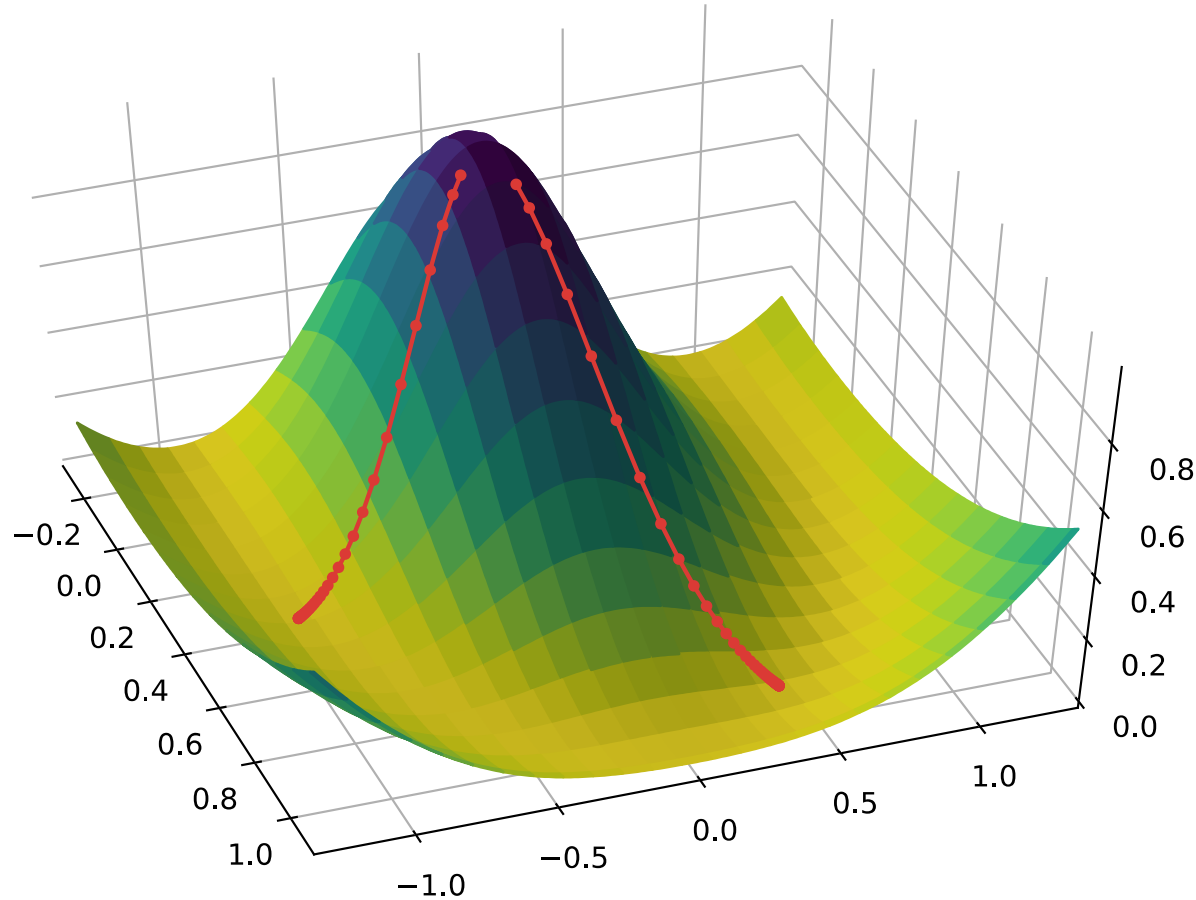
# Сценарии

- ▶ Важно правильно настроить learning rate и придумать расписание
- ▶ Можно уменьшать learning rate со временем по расписанию
  - В начале хотим активно исследовать пространство параметров, чтобы найти хорошую область
  - В конце хотим более детально исследовать найденную область
  - Используем ReduceLROnPlateau для зависимости от ошибок на валидации.



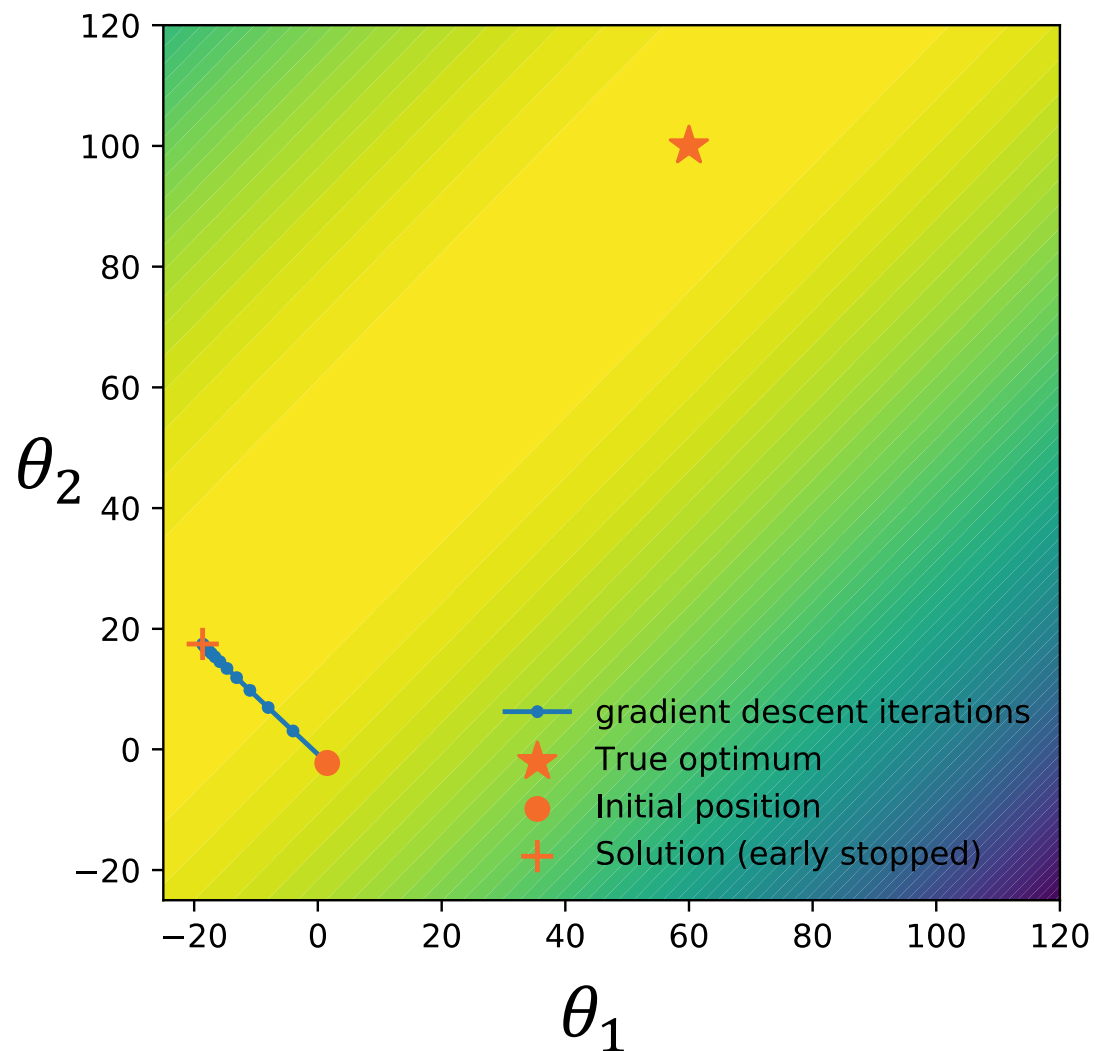


# Градиентный спуск для невыпуклых функций



- ▶ Может достичь минимума, который не является глобальным
- ▶ Результат зависит от начальной точки

# Градиентный спуск как средство регуляризации



- ▶ Большие значения параметров обычно означают переобучение
- ▶ Можно избежать этой проблемы, инициализируя параметры малыми значениями и останавливая градиентный спуск на ранней стадии (early stopping)

# Стохастический градиентный спуск (SGD)

- ▶ В машинном обучении мы оптимизируем функции потерь, которые обычно являются средними по объектам:

$$L = \frac{1}{N} \sum_{i=1 \dots N} \mathcal{L}(y_i, \hat{f}_{\theta}(x_i))$$

- ▶ Для больших  $N$ , градиентный спуск неэффективен с точки зрения вычислений и может быть неосуществим с точки зрения потребления памяти
- ▶ Альтернатива:
  - На каждом шаге  $k$  взять случайный  $l_k \in \{1, \dots, N\}$
  - Оптимизировать:  $\theta^{(k)} \leftarrow \theta^{(k-1)} - \eta \nabla_{\theta} \mathcal{L}(y_{l_k}, \hat{f}_{\theta}(x_{l_k})) \Big|_{\theta = \theta^{(k-1)}}$

# SGD и мини-батчи

## ► SGD:

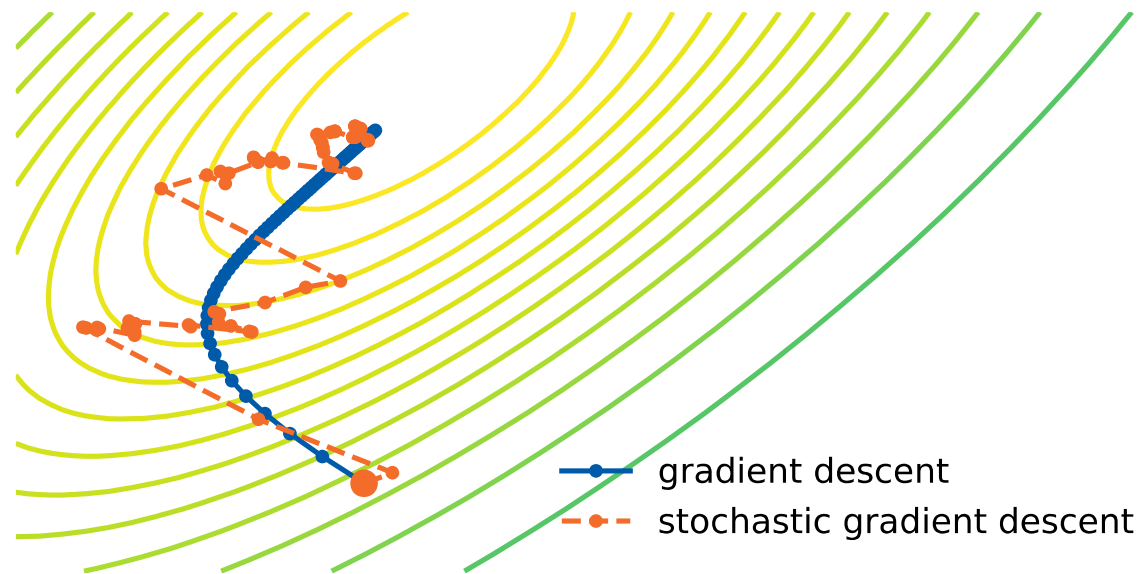
- На каждом шаге  $k$  взять случайный  $l_k \in \{1, \dots, N\}$ , тогда изменение:
- $\theta^{(k)} \leftarrow \theta^{(k-1)} - \eta \nabla_{\theta} \mathcal{L}(y_{l_k}, \hat{f}_{\theta}(x_{l_k})) \Big|_{\theta = \theta^{(k-1)}}$

## ► Мини-пакетный (батч) SGD:

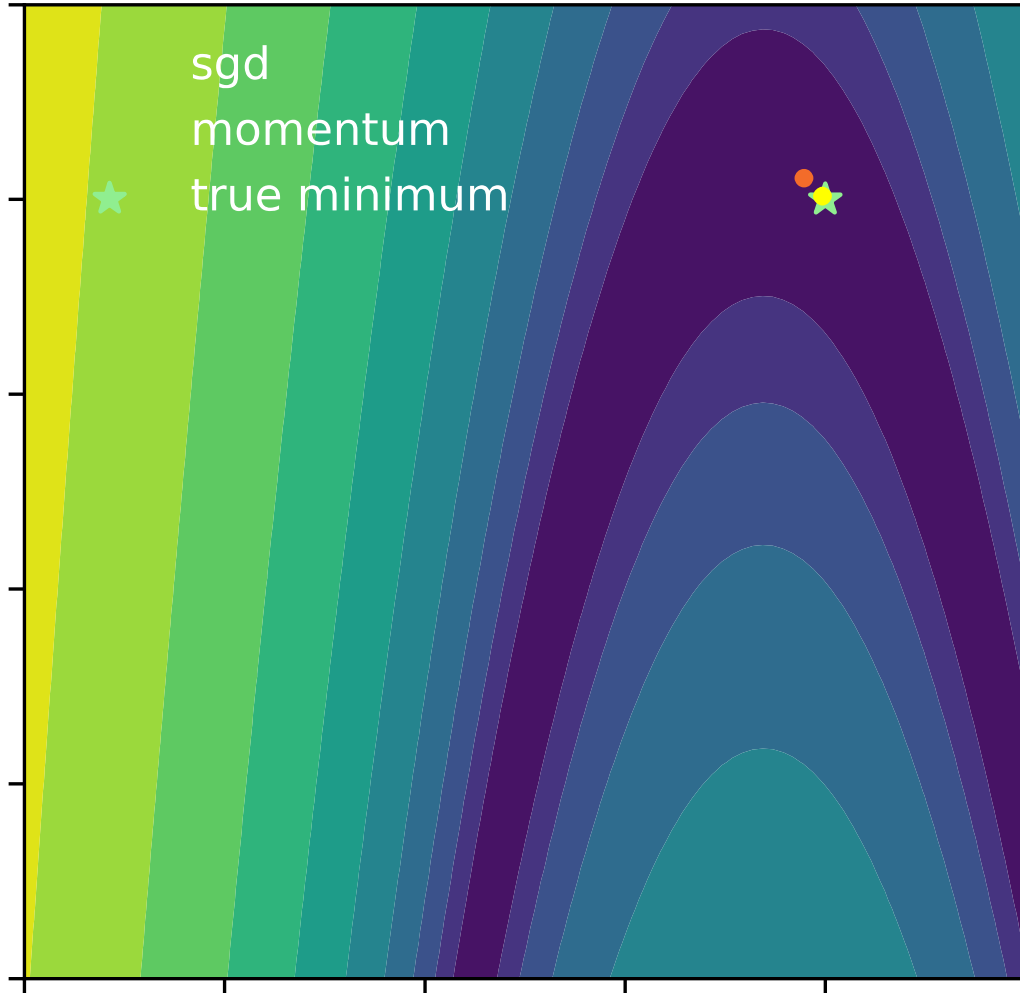
- Перемешиваем обучающий НД, затем итерируем по нему пакетами (батчами) фиксированного размера.
- На каждой итерации оцениваем градиенты потерь на данном участке.

$$B: g = \sum_{i \in B} \nabla_{\theta} \mathcal{L}(y_i, \hat{f}_{\theta}(x_i))$$

- Модифицируем параметры:  $\theta^{(k)} \leftarrow \theta^{(k-1)} - \eta \cdot g$



# Импульсный (momentum) SGD



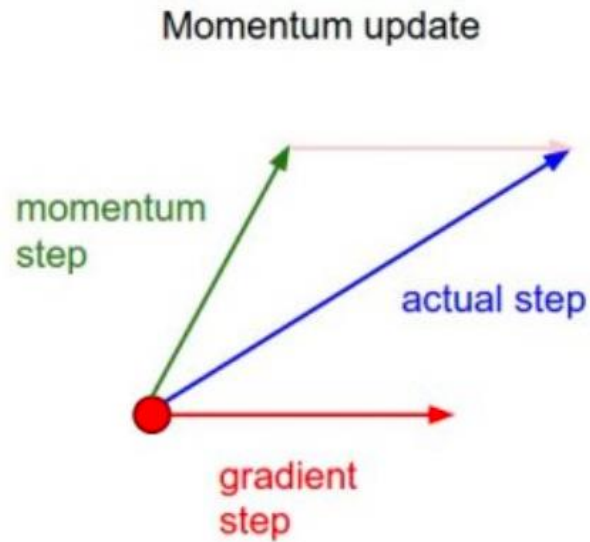
- ▶ Идея: ввести инерцию (как у мяча, катящегося с горы)
  - Помогает выйти из небольших локальных минимумов
  - Позволяет использовать более широкий диапазон скоростей обучения\*

$$m^{(k)} \leftarrow \beta \cdot m^{(k-1)} + (1 - \beta) \cdot \left. \frac{\partial L}{\partial \theta} \right|_{\theta = \theta^{(k-1)}}$$

$$\theta^{(k)} \leftarrow \theta^{(k-1)} - \eta \cdot m^{(k)}$$

- Фактически, экспоненциальное сглаживание.
- Дороже и капризнее обычного SGD

# Метод Нестерова



- ▶ На самом деле мы уже знаем, что сдвинемся на  $\beta m_{t-1}$  на промежуточном шаге
- ▶ Можем вычислить градиент в этой точке, на полпути

$$m_t = \beta m_{t-1} + (1 - \beta) \nabla_{\theta} L(\theta_t - \beta m_{t-1})$$

# AdaGrad

- ▶ **Идея:** давайте быстрее двигаться по тем параметрам, которые не сильно меняются, и медленнее по быстро меняющимся параметрам:

- ▶ Обозначим  $g_t = \nabla_{\theta} L(\theta_t)$

- ▶ Накапливаем историю градиентов

$$v_t = \sum_{\tau=1}^t g_{\tau}^2$$

и учитываем ее

$$\theta_{t+1} = \theta_t - \frac{\alpha}{\sqrt{v_t + \epsilon}} g_t$$

- ▶ Learning rate все время уменьшается, но с разной скоростью для разных параметров

# RMSprop

- ▶ Идея: настроить скорость обучения отдельно для разных компонентов вектора параметров, избегая **проблемы** Adagrad с исчезающим градиентом
- ▶ Считать скользящее среднее

$$v_t = \beta v_{t-1} + (1 - \beta) g_t^2$$

- ▶ обновлять параметры с его использованием

$$\theta_{t+1} = \theta_t - \frac{\alpha}{\sqrt{v_t + \epsilon}} g_t$$



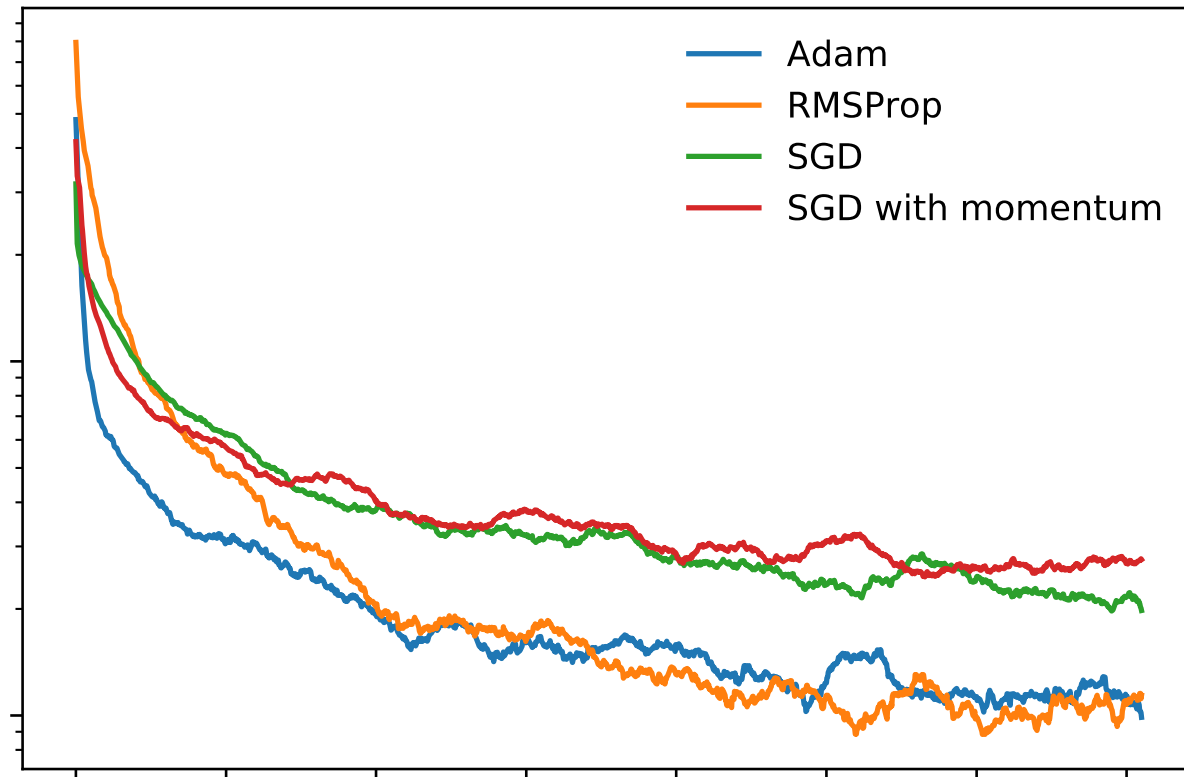
# Adam

- ▶ Скомбинируем идеи (momentum + RMSprop)
- ▶ Обычно хороший алгоритм для начала оптимизации

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

$$\theta_{t+1} = \theta_t - \frac{\alpha}{\sqrt{v_t} + \epsilon} m_t$$



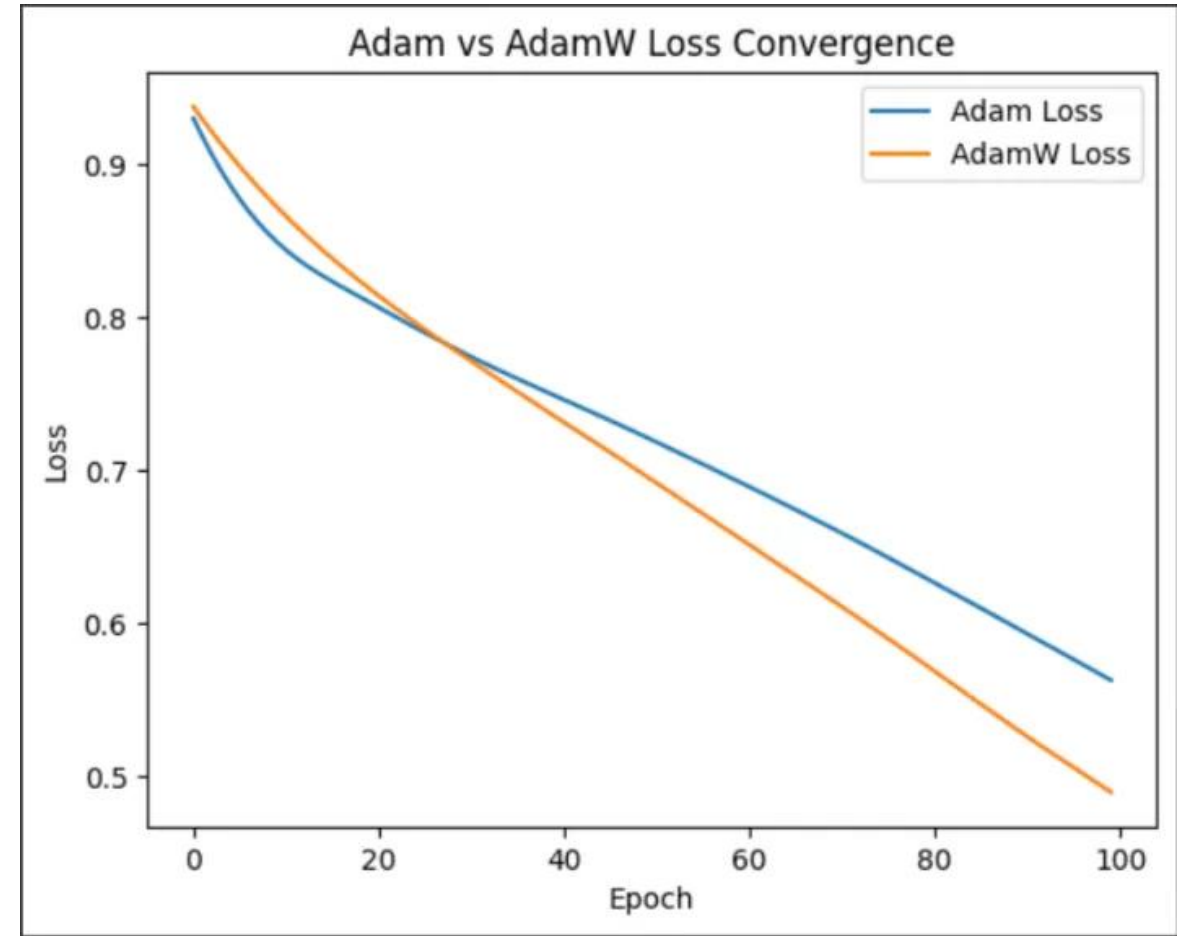
# AdamW

- Вместо добавления механизма затухания весов к функции потерь, этот механизм применяет затухание весов непосредственно во время обновления параметров, что приводит к более стабильной регуляризации и лучшей обобщающей способности.

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) \nabla f(\theta_t)$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) (\nabla f(\theta_t))^2$$

$$\theta_{t+1} = \theta_t - \eta \frac{m_t}{\sqrt{v_t} + \epsilon}$$



# Schedule-free optimization

$$\begin{aligned}y_t &= (1 - \beta)z_t + \beta x_t, \\z_{t+1} &= z_t - \gamma \nabla f(y_t, \zeta_t), \\x_{t+1} &= (1 - c_{t+1})x_t + c_{t+1}z_{t+1}, \\c_{t+1} &= \frac{1}{t+1}\end{aligned}$$

$x$  - последовательность, в которой должны происходить оценки потерь test/val, которая отличается от первичных итераций  $z$  и мест оценки градиента  $y$ . Обновления в  $z$  соответствуют базовому оптимизатору, в данном случае простому градиентному спуску / SGD.

Полученный алгоритм не должен иметь расписания уменьшения LR, при этом ведёт себя оптимальнее других.

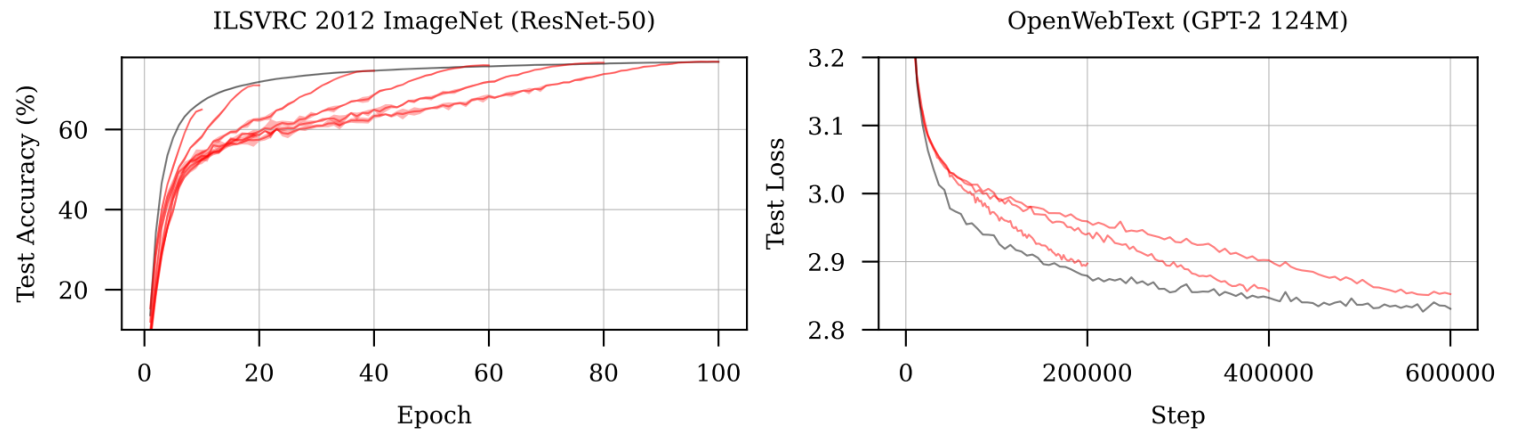


Figure 1: Schedule-Free methods (black) closely track the Pareto frontier of loss v.s. training time in a single run. Both Schedule-Free SGD (left) and AdamW (right) match or exceed the performance of cosine learning rate schedules of varying lengths (red).

# Muon (MomentUm Orthogonalized by Newton-Schulz)

Шаг Ньютона-Шульца:

приблизительная ортогонализация матрицы обновления, то есть в применении следующей операции:

$$\text{Ortho}(G) = \arg \min_O \{\|O - G\|_F : \text{either } O^\top O = I \text{ or } OO^\top = I\}$$

Ортогонализация эффективно увеличивает масштаб «редких направлений» в оптимизации, которые имеют небольшую величину в обновлении, но тем не менее важны для обучения.

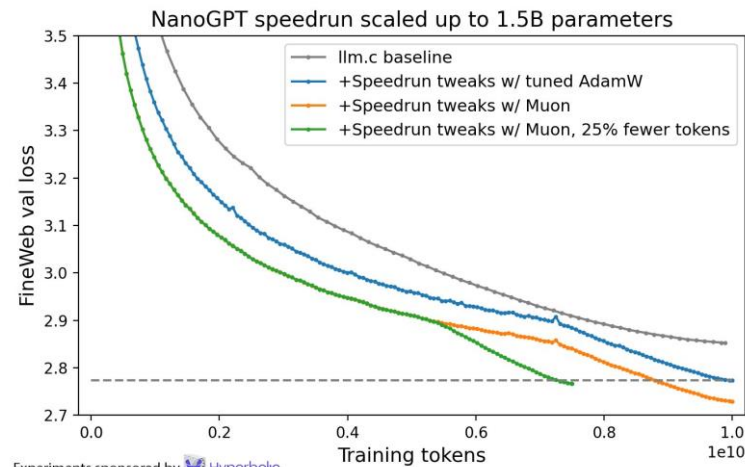
---

## Algorithm 2 Muon

---

**Require:** Learning rate  $\eta$ , momentum  $\mu$

- 1: Initialize  $B_0 \leftarrow 0$
  - 2: **for**  $t = 1, \dots$  **do**
  - 3:     Compute gradient  $G_t \leftarrow \nabla_{\theta} \mathcal{L}_t(\theta_{t-1})$
  - 4:      $B_t \leftarrow \mu B_{t-1} + G_t$
  - 5:      $O_t \leftarrow \text{NewtonSchulz5}(B_t)$
  - 6:     Update parameters  $\theta_t \leftarrow \theta_{t-1} - \eta O_t$
  - 7: **end for**
  - 8: **return**  $\theta_t$
- 



<https://kellerjordan.github.io/posts/muon/>

# Саммари

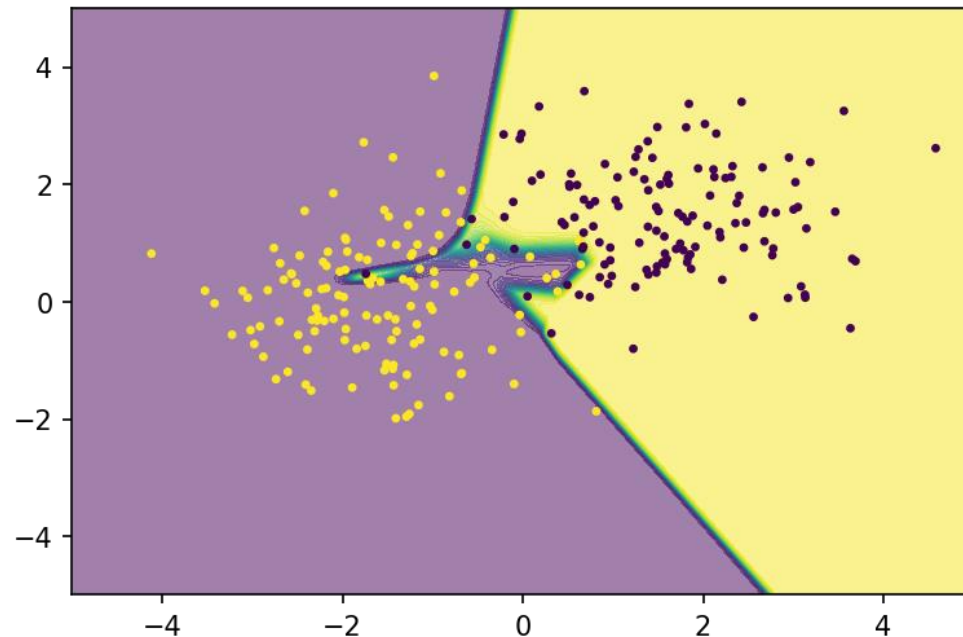
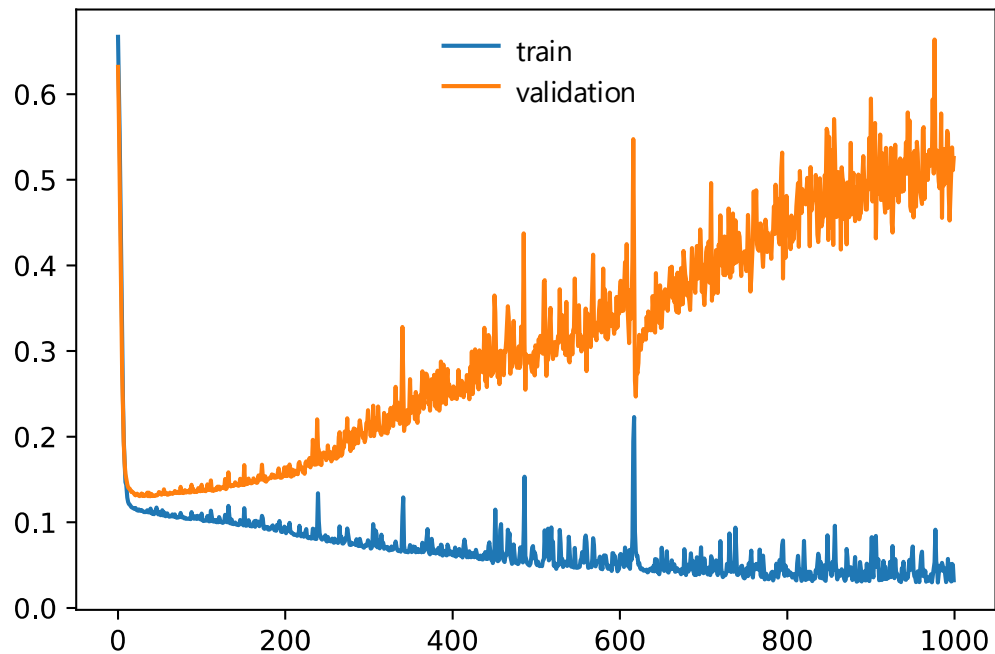
- ▶ Множество разнообразных алгоритмов оптимизации
- ▶ Хороший старт с Adam
- ▶ Muon – лидер в современных приложениях

# Переобучение



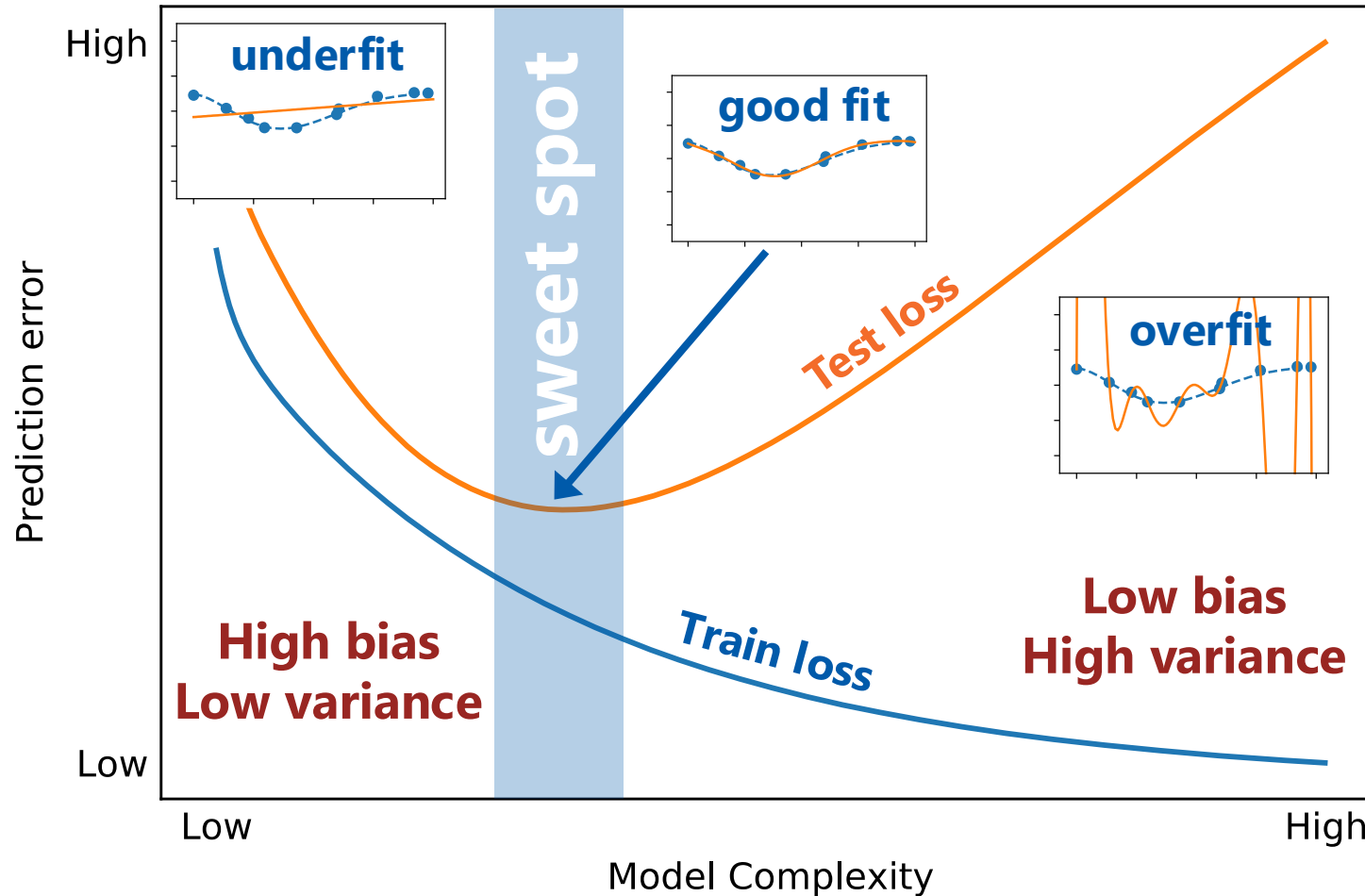
# Проблема переобучения

- Будучи очень сложными моделями, нейронные сети склонны к переобучению



# Bias-variance

- Будучи очень сложными моделями, нейронные сети склонны к переобучению

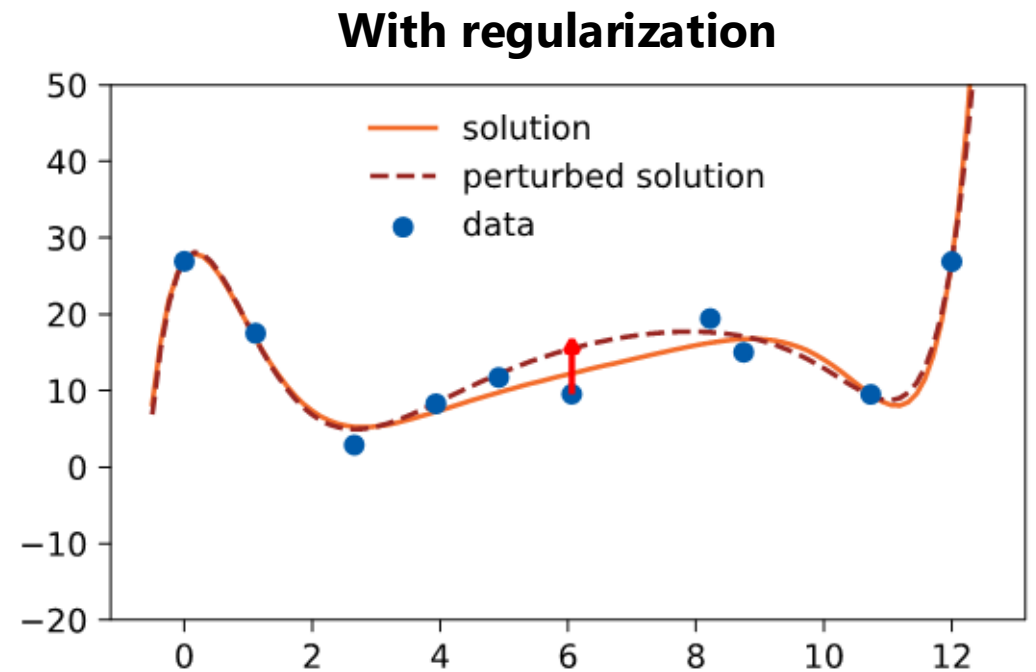
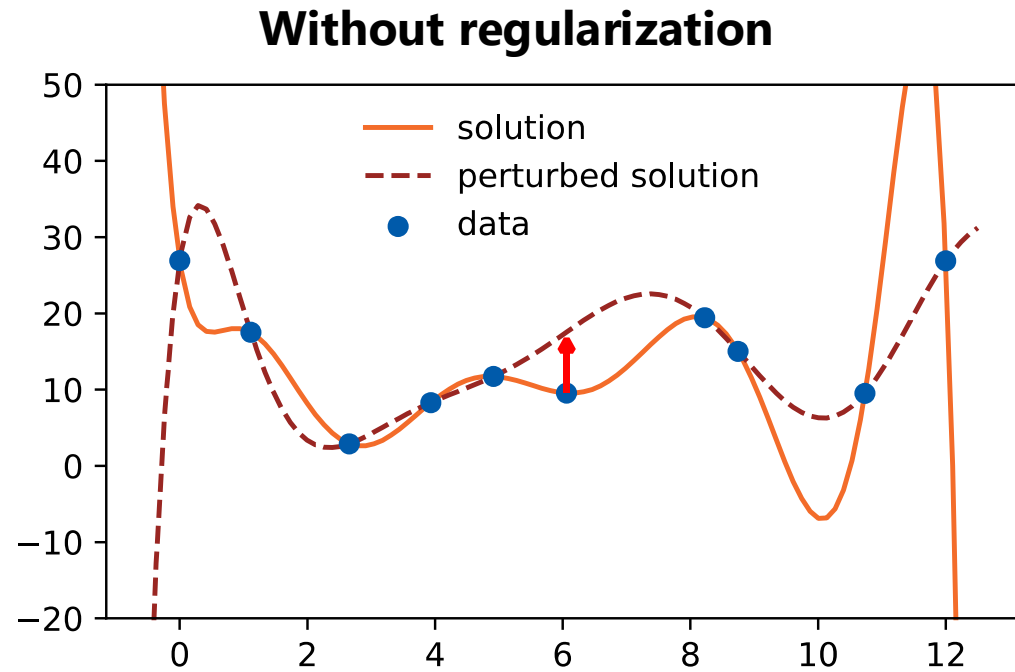


Обычно есть выбор между смещением и дисперсией.

Bias-Variance trade-off



# Регуляризация



NB: регуляризованная модель больше не является несмещённой!

То есть мы увеличили смещение, чтобы уменьшить дисперсию.

# Разные методы регуляризации

L2 regularization (Ridge):

$$\mathcal{L} = \|X\theta_\tau - y_\tau\|^2 + \alpha\|\theta_\tau\|^2$$

L1 regularization (Lasso):

$$\mathcal{L} = \|X\theta_\tau - y_\tau\|^2 + \alpha\|\theta_\tau\|_1$$

Elastic net:

$$\mathcal{L} = \|X\theta_\tau - y_\tau\|^2 + \alpha\|\theta_\tau\|^2 + \beta\|\theta_\tau\|_1$$

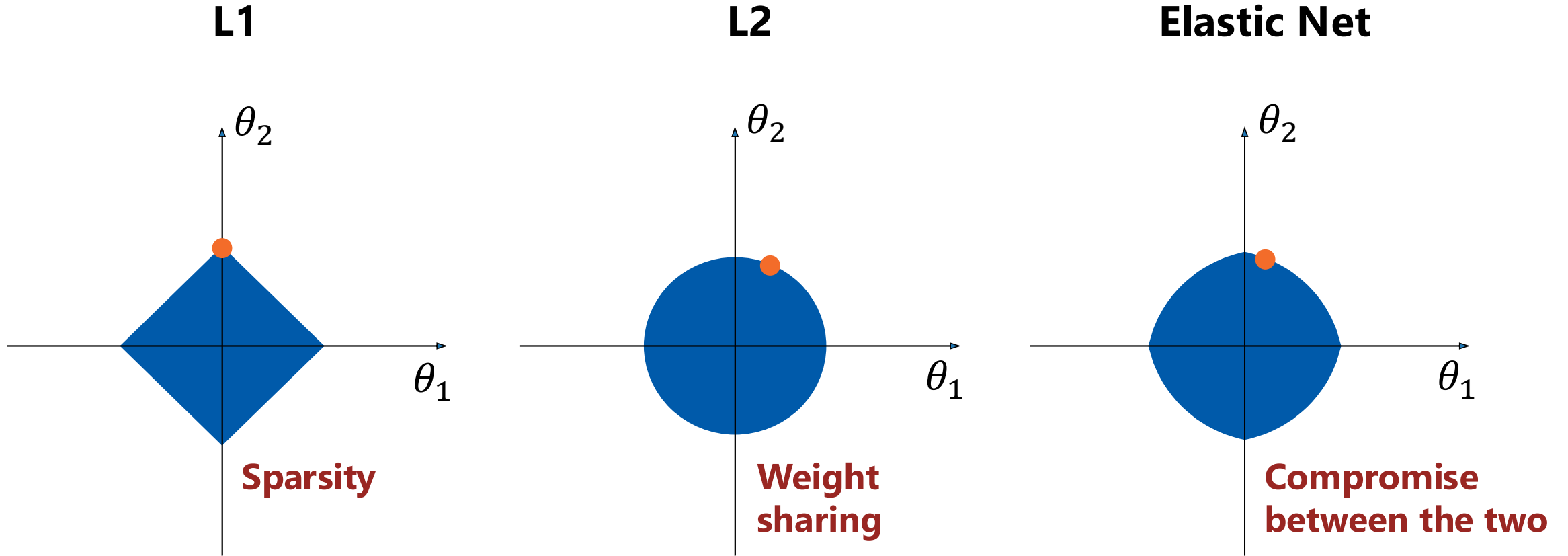
**L2 norm:**

$$\|x\|^2 \equiv \sum_{i=1\dots d} x_i^2$$

**L1 norm:**

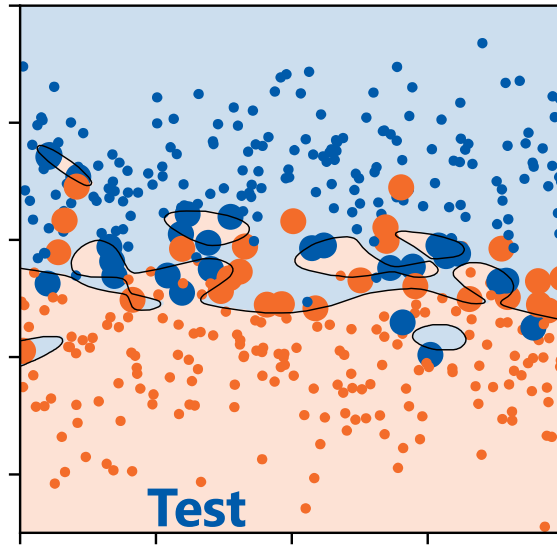
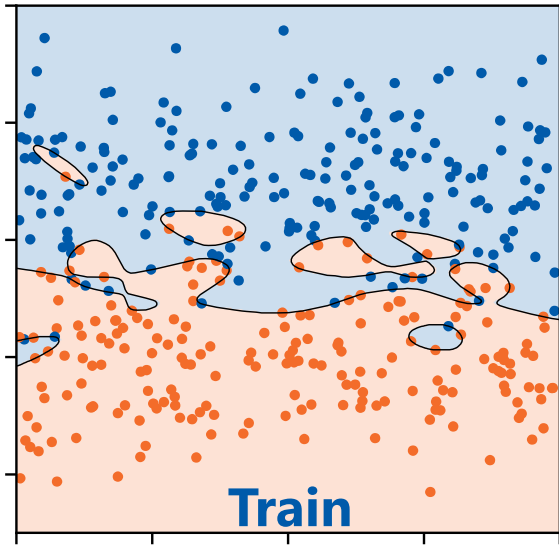
$$\|x\|_1 \equiv \sum_{i=1\dots d} |x_i|$$

# Свойства регуляризации



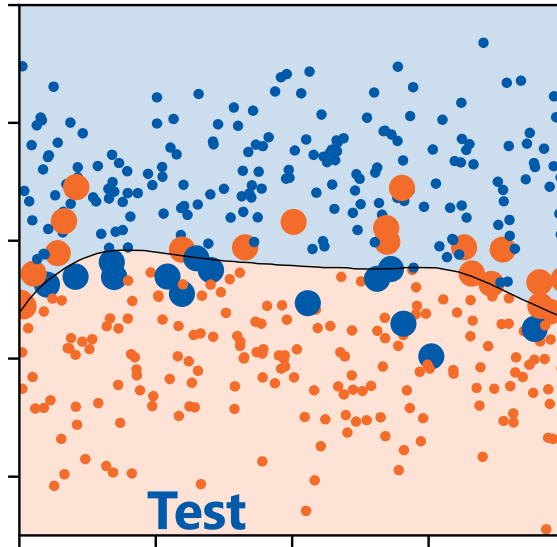
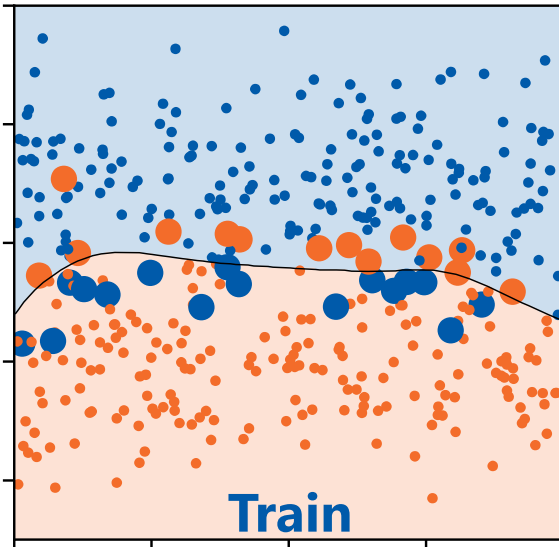
Все они смещают вес в сторону меньших значений.  
Однако они вызывают различные свойства решения.

# Пример: L2-регуляризованная классификация



**Without regularization**

Регуляризируя модель, мы увеличиваем функцию потерь в обучении и уменьшаем потери при тестировании.



**With regularization**

Это улучшает обобщаемость модели.

# Регуляризация и подходы

**MSE loss  $\Leftrightarrow$  Prob.  
model with normal  
label noise!**

Частотный подход:

**Регуляризация эквивалентна  
предположению в модели о  
распределении данных**

L2 регуляризация это разрешение  
модели учитывать нормальный  
шум

**Normal prior  $\Leftrightarrow$  L2  
regularization**

Байесовский подход:

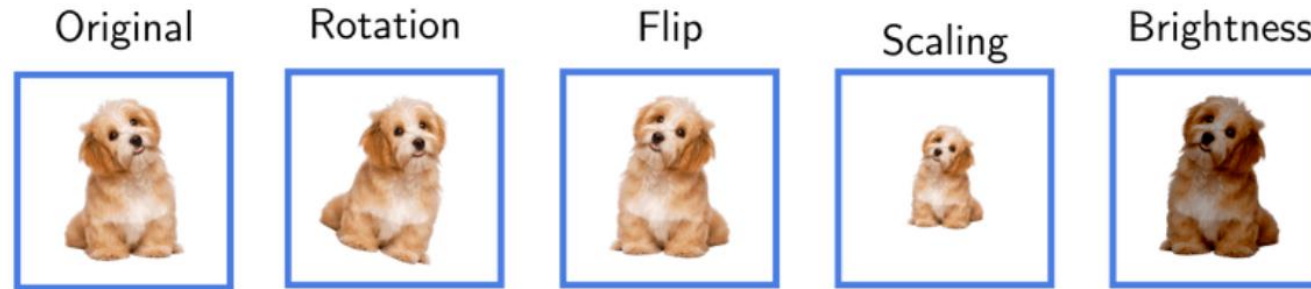
**Регуляризация эквивалентна  
априорному предположению о  
распределении данных**

Байесовские нейронные сети =  
регуляризация

L2 регуляризация это априорное  
нормальное распределение

# Аугментация данных

Что если брать регуляризацию из данных, зная некоторые свойства



<https://www.baeldung.com/cs/ml-data-augmentation>

Примеры аугментаций:

- Для изображений - flip, rotation, crop, rescale,...
- Для звука - добавляем фоновый шум, меняем высоту звука
- Для текстов - замена на синонимы, backtranslation
- Для последовательностей - пропускаем элементы, меняем местами

# Зашумление данных

Добавляем шум в данные

- на входе
- зашумление весов
- на выходе - зашумляем метки

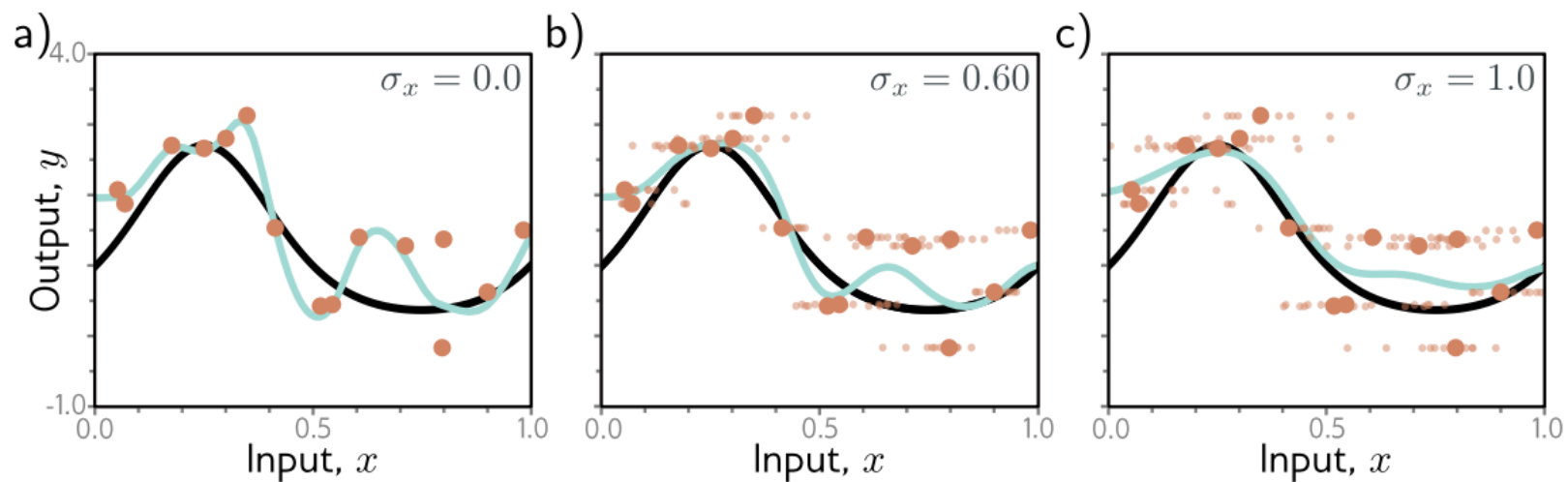
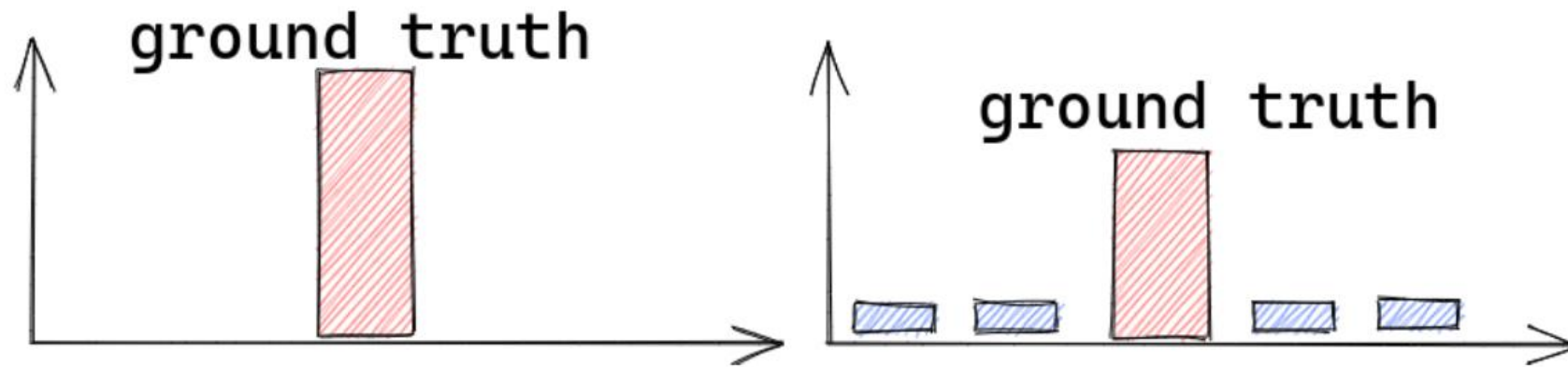


Image credit

# Зашумление меток



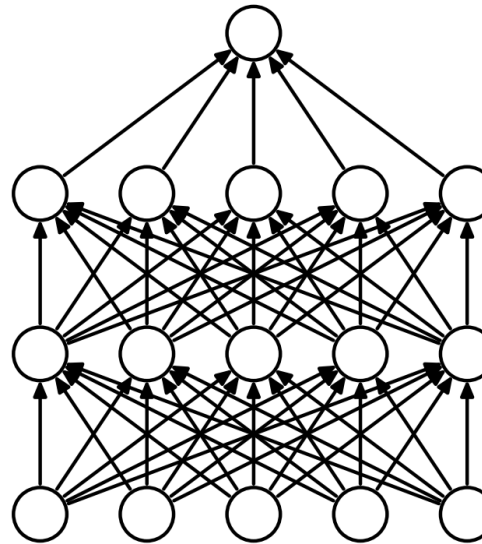
Вероятность правильной метки  $1 - \epsilon$ , всех остальных -  $\epsilon/(K - 1)$



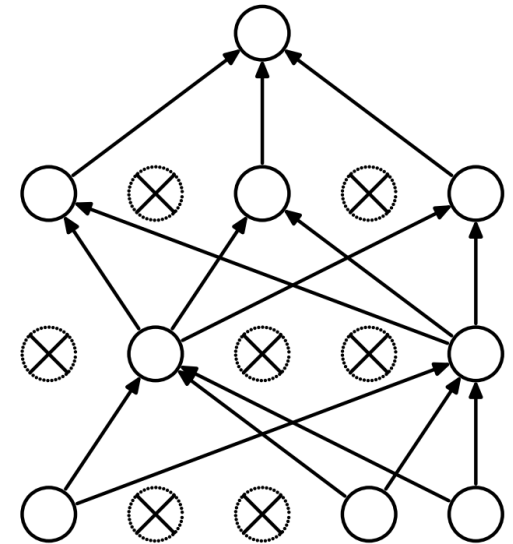
# Dropout

- ▶ При обучении – ставим активацию нейронов 0 с вероятностью  $p$
- ▶ При тестировании – умножаем на  $1 - p$ 
  - то есть приравниваем матожиданию
- ▶ Заставляет нейрон учиться работать со случайно выбранной выборкой других нейронов
- ▶ Заставляет его создавать полезные признаки, а не полагаться на другие нейроны для исправления своих ошибок.

Image from:  
<http://jmlr.org/papers/v15/srivastava14a.html>

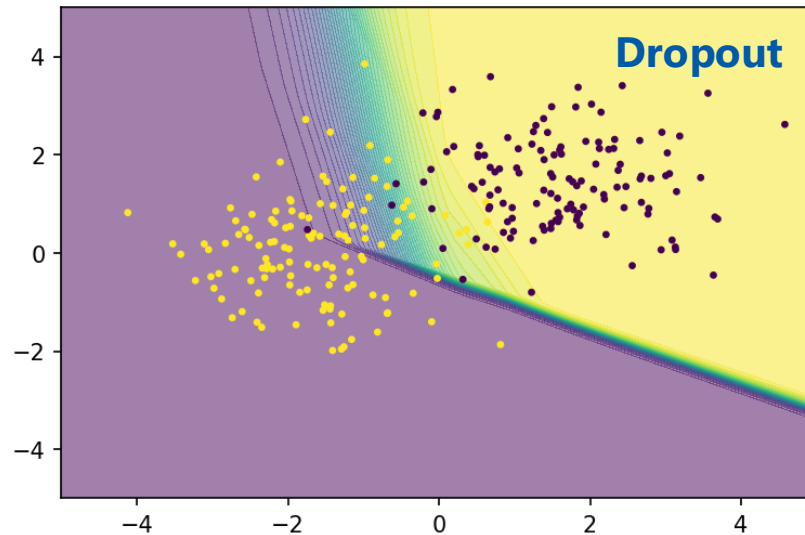
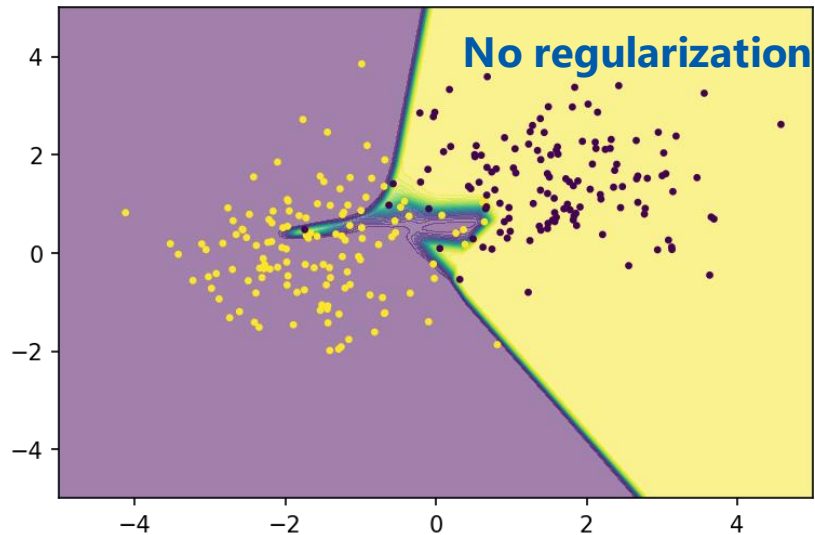
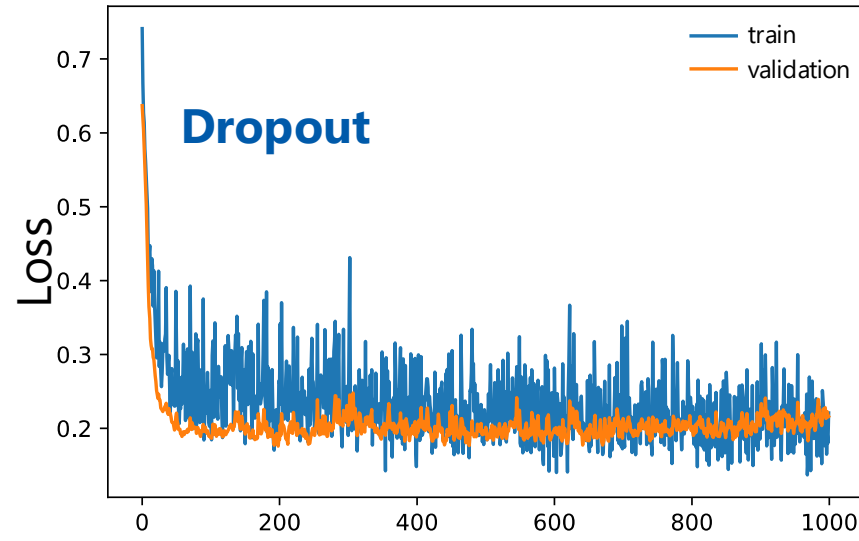
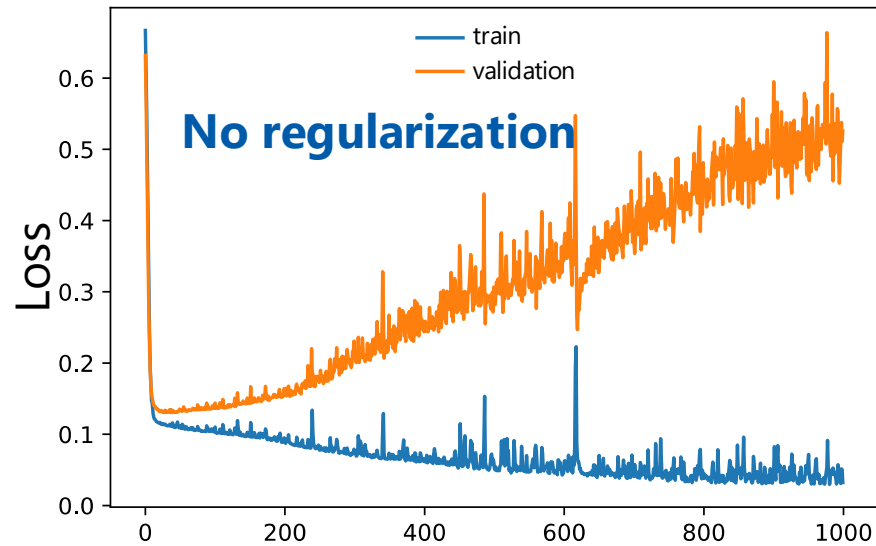


(a) Standard Neural Net



(b) After applying dropout.

# Пример прореживания (dropout)



В этом примере dropout приводит к гораздо лучшему (хотя все еще не идеальному) соответствию с меньшей ошибкой теста.

# Нормализационные слои



# Пакетная нормализация (Batch normalization)

- ▶ Первоначально этот метод был предложен для смягчения «внутреннего ковариационного сдвига»
- ▶ Работает следующим образом (входы слоя  $x_i$ , выходы  $y_i$ ):

## internal covariate shift

обновления в одном слое  
изменяют входные  
распределения последующих  
слоев

- рассчитать выборочное среднее значение и дисперсию входных данных для одного батча  $B$

$$\mu_B = \frac{1}{|B|} \sum_{i \in B} x_i \quad \sigma_B^2 = \frac{1}{|B|} \sum_{i \in B} (x_i - \mu_B)^2$$

- нормализуем входные данные, затем масштабируем и сдвигаем (с обучаемыми параметрами  $(\gamma, \beta)$ ):

$$y_i = \gamma \cdot \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} + \beta$$

# Пакетная нормализация

- ▶ Во многих случаях оказалось чрезвычайно эффективным.
  - Более быстрая и стабильная сходимость
- ▶ Позже было доказано, что это **не уменьшает** ковариационный сдвиг.
- ▶ По сути, удаляет степени свободы «сдвиг» и «масштабирование» из предыдущего слоя.

## internal covariate shift

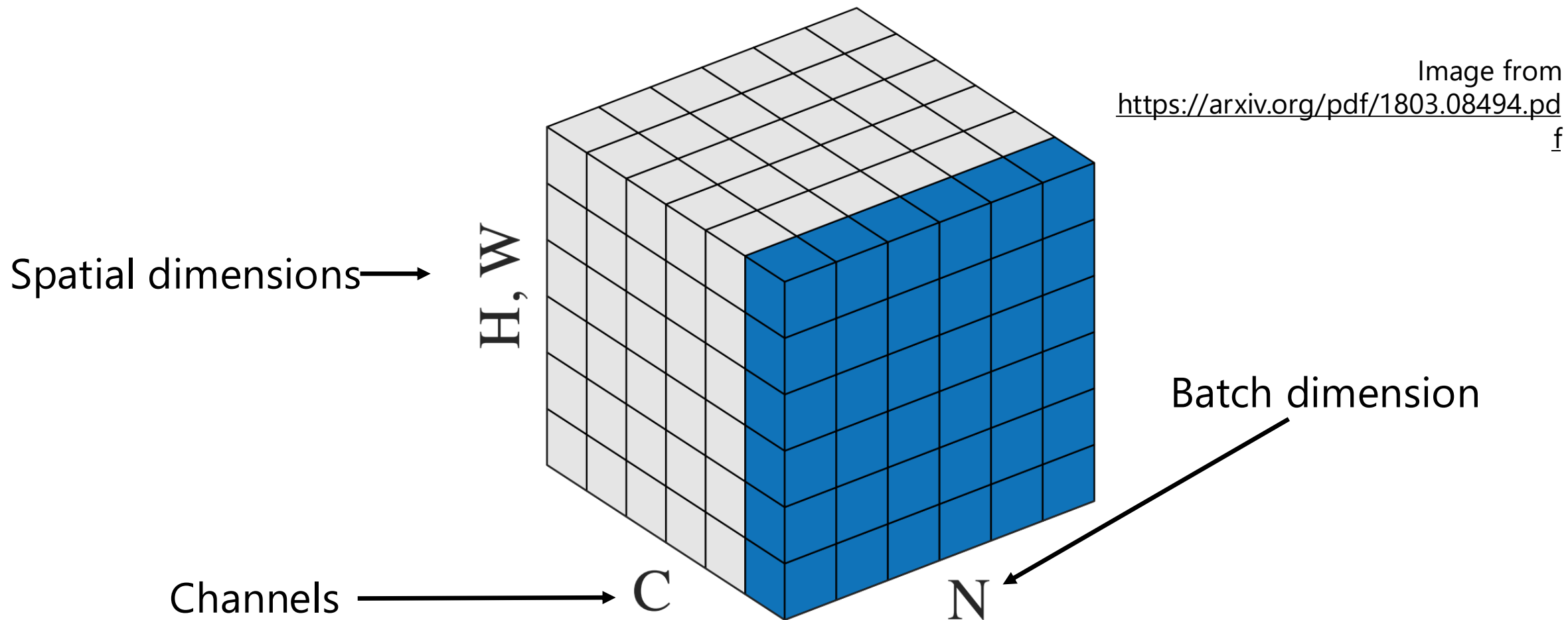
the updates in one layer  
change the input distributions  
of the subsequent layers

$$y_i = \gamma \cdot \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} + \beta$$

# Пакетная нормализация

- ▶ По какому измерению следует проводить нормализацию? Обычно это делается так:
  - Batch of 1D vectors [ $\text{Batch\_dim}$  x  $\text{Features\_dim}$ ]
    - separately for each component in  $\text{Features\_dim}$ , i.e. over  $\text{Batch\_dim}$
  - Batch of ND objects [ $\text{Batch\_dim}$  x  $\text{Spacial\_dim1}$  x ... x  $\text{Channel\_dim}$ ]
    - separately for each component in  $\text{Channel\_dim}$ , i.e. over  $\text{Batch\_dim}$  x  $\text{Spacial\_dim1}$  x ...

# Пакетная нормализация

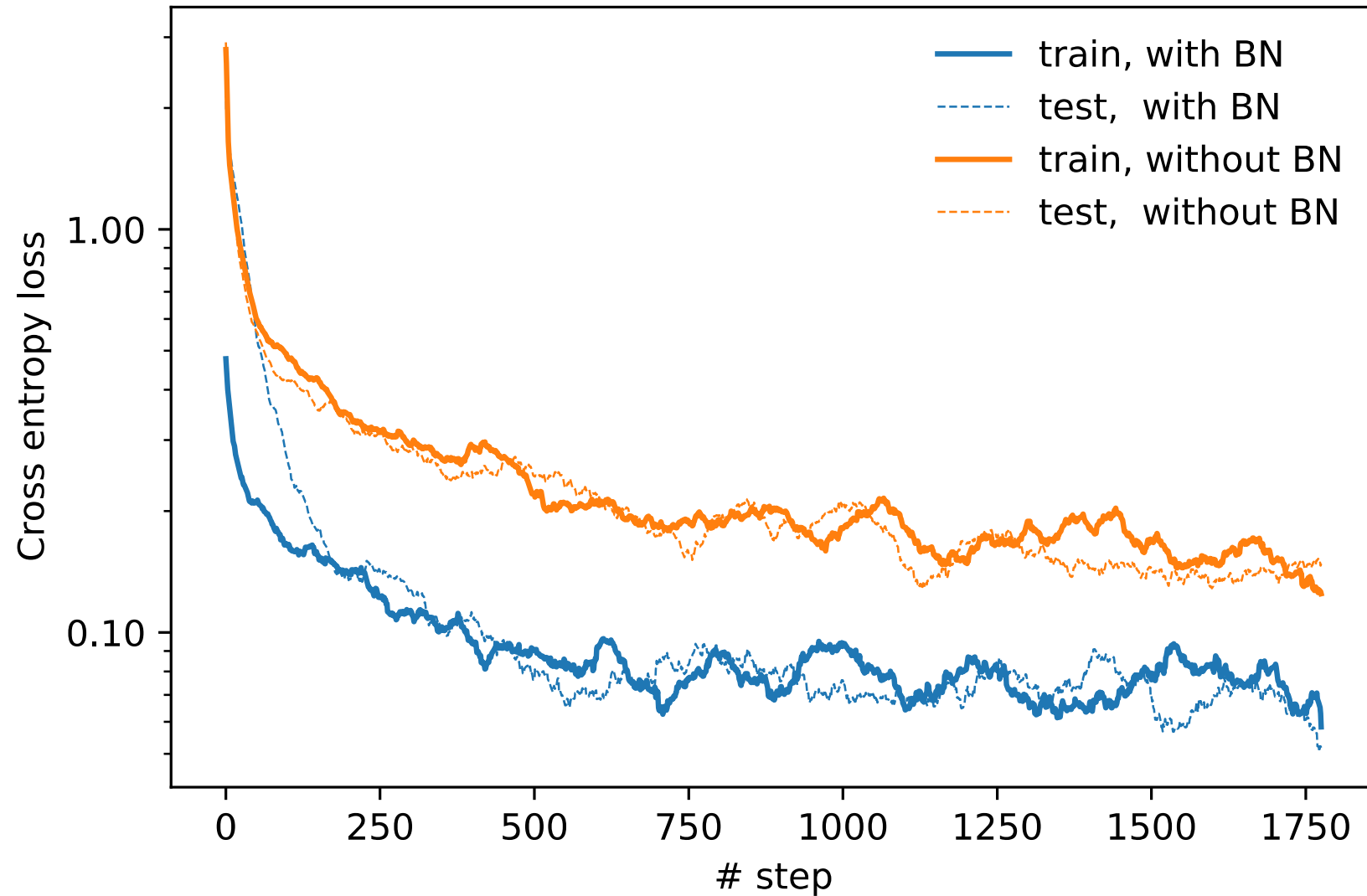


# Пакетная нормализация во время вывода результатов

- ▶ Расчет статистики по пакетам во время тестирования может быть проблематичным.
  - Например, когда есть только один объект
- ▶ Вместо этого: вычисляем скользящее среднее и дисперсию во время обучения, применяйте их во время тестирования.



# Example: CNN on MNIST



(shown: moving average loss)

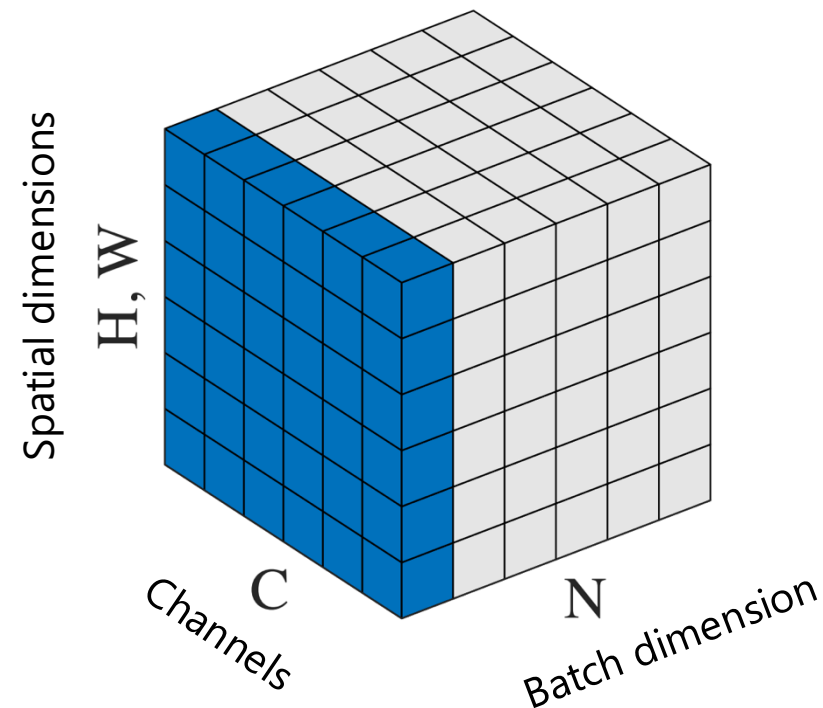
# Нормализация слоёв

- ▶ Нормализация пакетов устанавливает ограничения на размер пакета.
  - Если значение слишком мало, то дисперсия выборочных статистик будет слишком высокой.

# Layer Normalization

- ▶ Нормализация пакетов устанавливает ограничения на размер пакета
  - Если значение слишком мало, то дисперсия выборочных статистик будет слишком высокой
- ▶ Проблема для рекуррентных сетей
- ▶ Альтернатива: **Layer Normalization**
  - Математические расчеты те же, за исключением того, что статистика вычисляется по каналам, а не по элементам пакета.
  - Однако эффект совершенно иной
    - е.г. Нормализация слоев «запутывает» различные нейроны внутри слоя.

Image from  
<https://arxiv.org/pdf/1803.08494.pdf>



# Выводы

- ▶ Нейронные сети можно регуляризовать с помощью штрафов L1/L2 или ранней остановки.
- ▶ Дропаут заставляет нейроны создавать полезные признаки, а не полагаться на другие нейроны для исправления своих ошибок.
- ▶ Пакетная нормализация — чрезвычайно мощный метод регуляризации, хотя причина этого не совсем ясна.