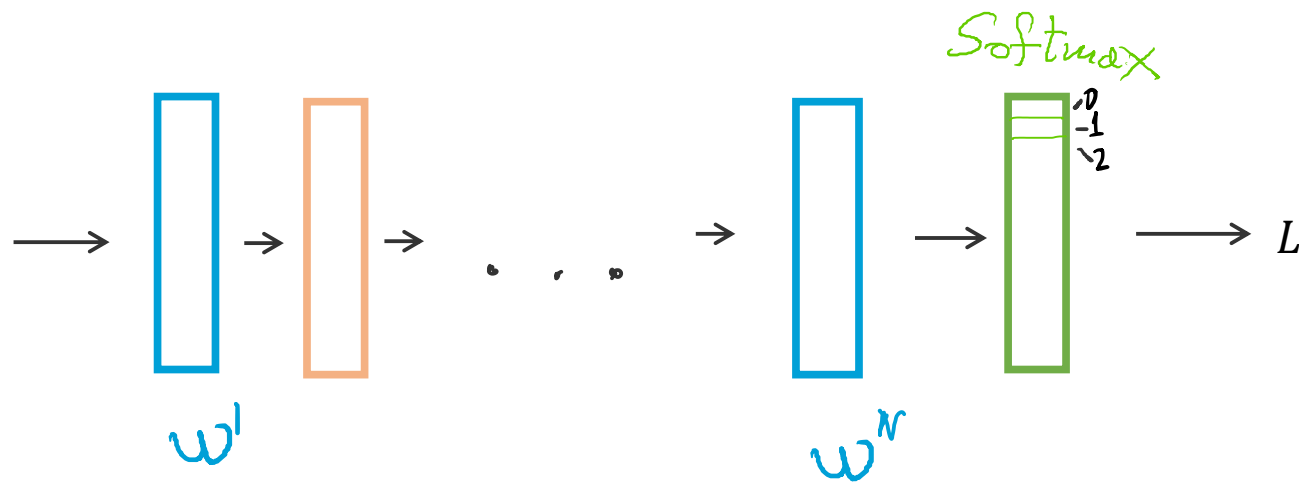


5

Нейронные сети —  
практика

# Общая схема тренировки

Прямой проход (forward pass) – посчитать loss



Обратный проход (backward pass) – посчитать градиент

Обновить параметры

$$\overrightarrow{w^1} = \overrightarrow{w^1} - \eta \overrightarrow{\nabla_{w^1} L}$$

...

$$\overrightarrow{w^n} = \overrightarrow{w^n} - \eta \overrightarrow{\nabla_{w^n} L}$$





“In theory, there is no difference between theory and practice. But, in practice, there is.”

---

Jan L. A. van de Snepscheut as quoted by Doug Rosenberg and Matt Stephens (2007) *Use Case Driven Object Modeling with UML Theory and Practice* p. xxvii; the quote is also cited without attribution in Doug Rosenberg and Kendall Scott (2001) *Applying Use Case Driven Object Modeling With UML*, p. 1

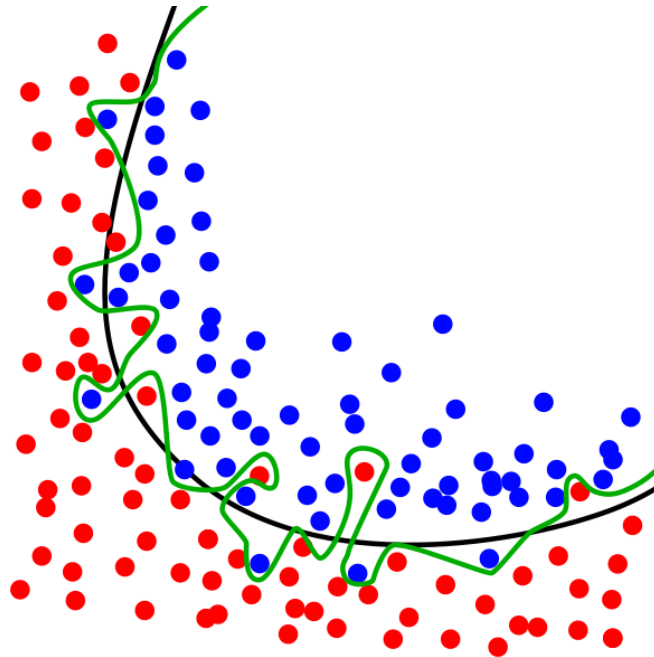
— This quote has also been attributed to Yogi Berra, to Chuck Reid, to William T. Harbaugh, to Manfred Eigen, to Alvin E. Roth *Wer hat das gesagt?* and to Karl Marx

[Source](#)

GPU



over/underfit

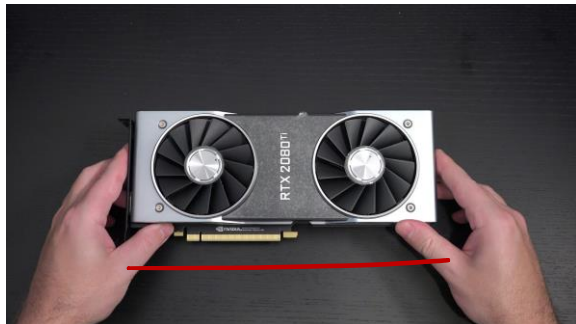


prod



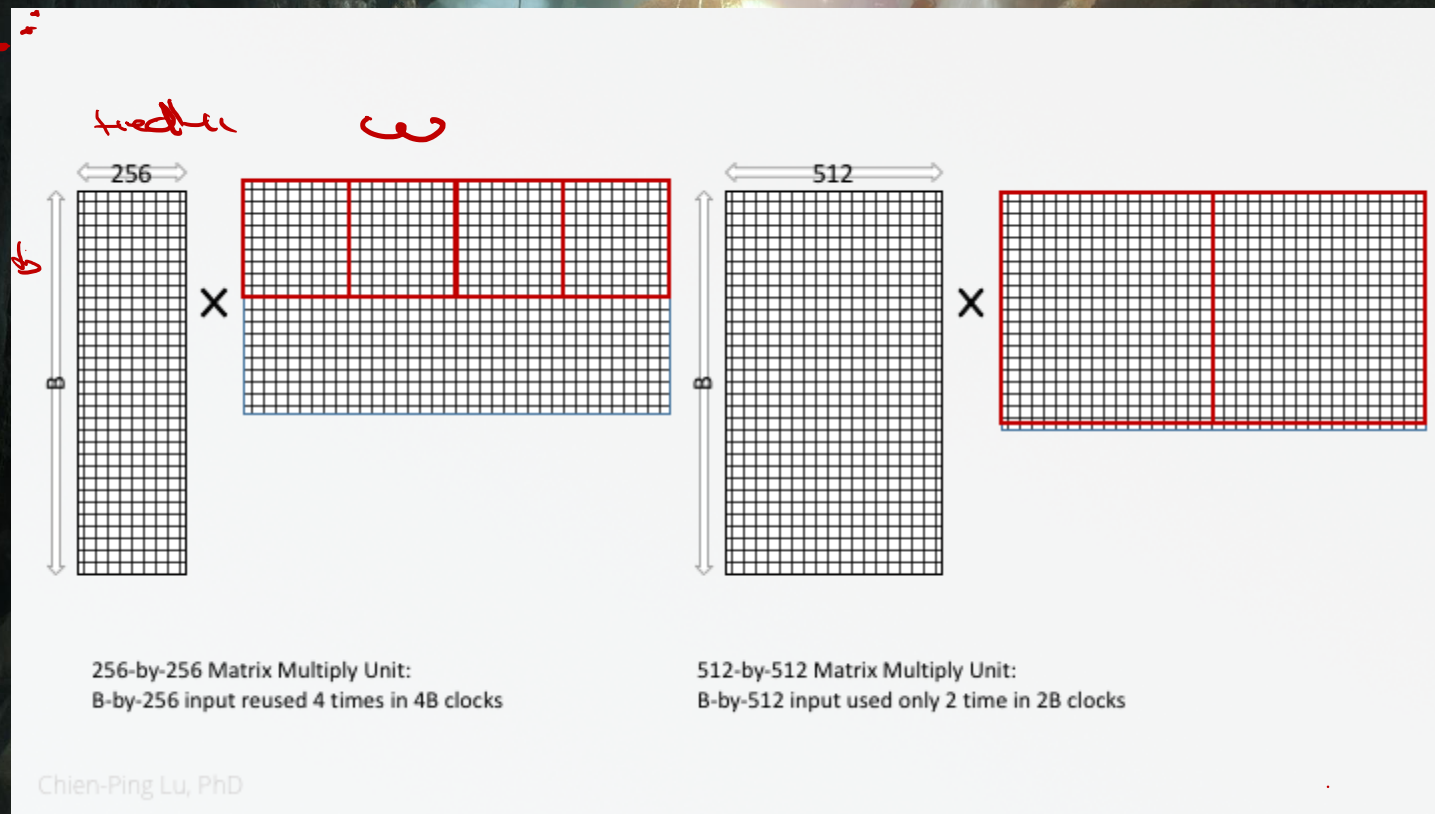
# GPU

## Graphics Processing Unit



	Цена	Flops	Количество ядер	Скорость памяти
<u>Intel Core i9-7980XE</u>	<u>\$2000</u>	1.8 TFlops	16	<u>70 Gb/sec</u>
<u>Nvidia Geforce 2080 Ti</u>	<u>\$1000</u>	<u>11.7 Tflops</u> <u>107 Tensor Tflops</u>	<u>4352</u>	<u>616 Gb/sec</u>

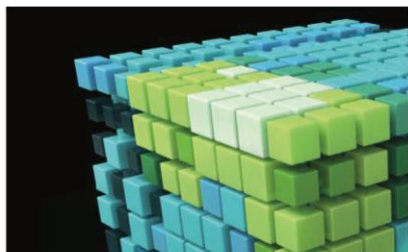




# Как программировать на GPU



cuBLAS



```
// Kernel definition
__global__ void VecAdd(float* A, float* B, float* C)
{
    int i = threadIdx.x;
    C[i] = A[i] + B[i];
}

int main()
{
    ...
    // Kernel invocation with N threads
    VecAdd<<<1, N>>>(A, B, C);
    ...
}
```





Caffe

theano

*dmlc*  
***mxnet***



PYTORCH

THE  
*secret ingredient*  
IS  
training  
process

# Machine Learning Flow



Ошибка на train

большая

underfitting

- Более мощную модель
- Больше ресурсов для тренировки
- Другой подход

маленькая ↓

Ошибка на val

большая

overfitting

- Больше данных
- Больше регуляризации
- Другой подход

маленькая ↓

Ошибка на test

большая

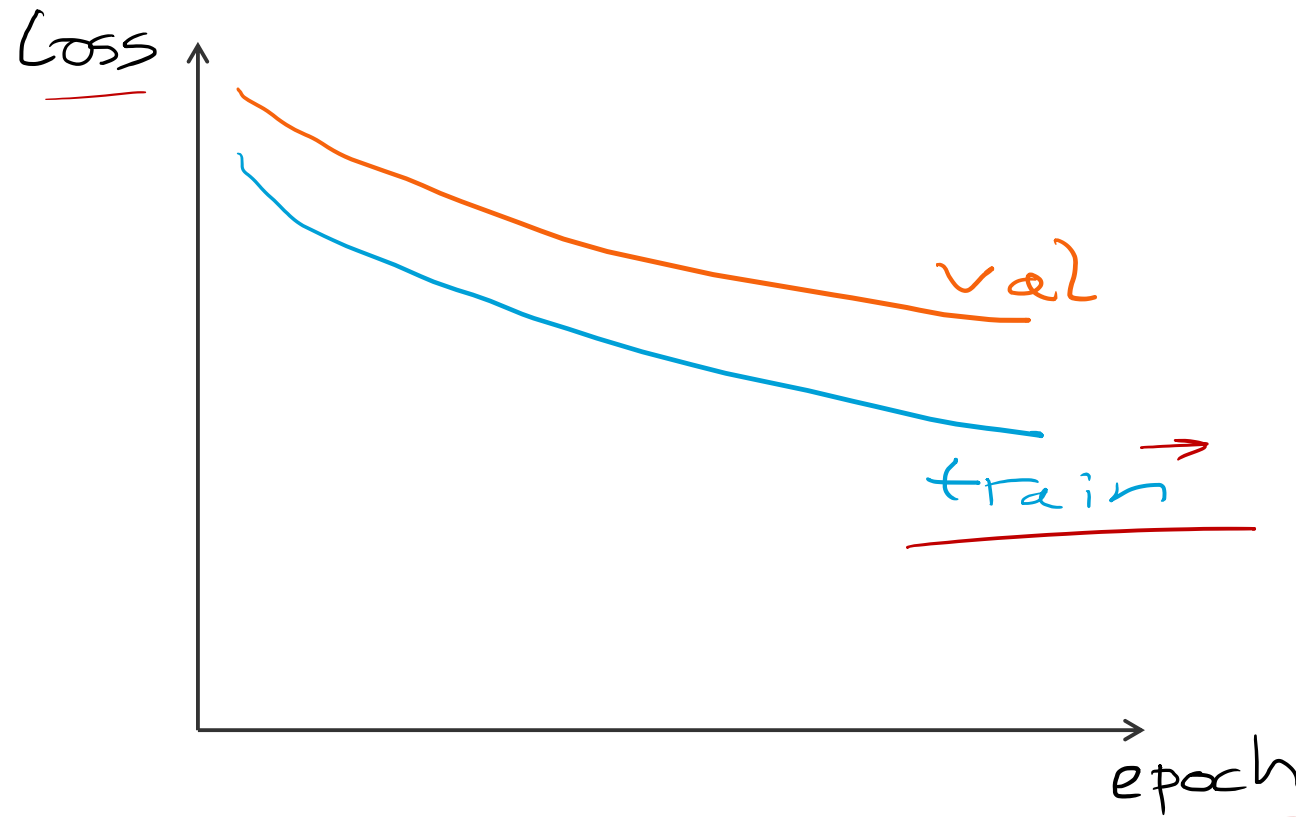
- Отличаются train и test
- Больше данных, таких как test

←  
маленькая



# Смотрим на графики тренировки

Train loss и val loss уменьшаются медленно  
Находятся близко друг к другу



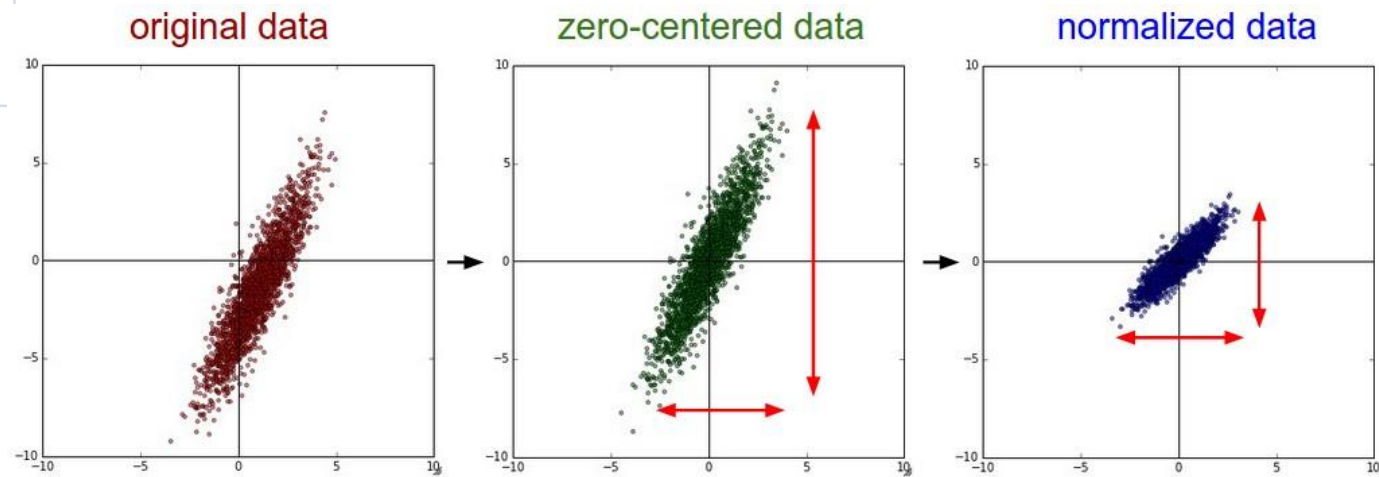
*underfitting*





# То, о чем мы говорили раньше

## Preprocessing



$$x' = \frac{x - \bar{x}}{\sigma}$$
$$\sigma = \frac{\sqrt{(x - \bar{x})^2}}{N - 1}$$

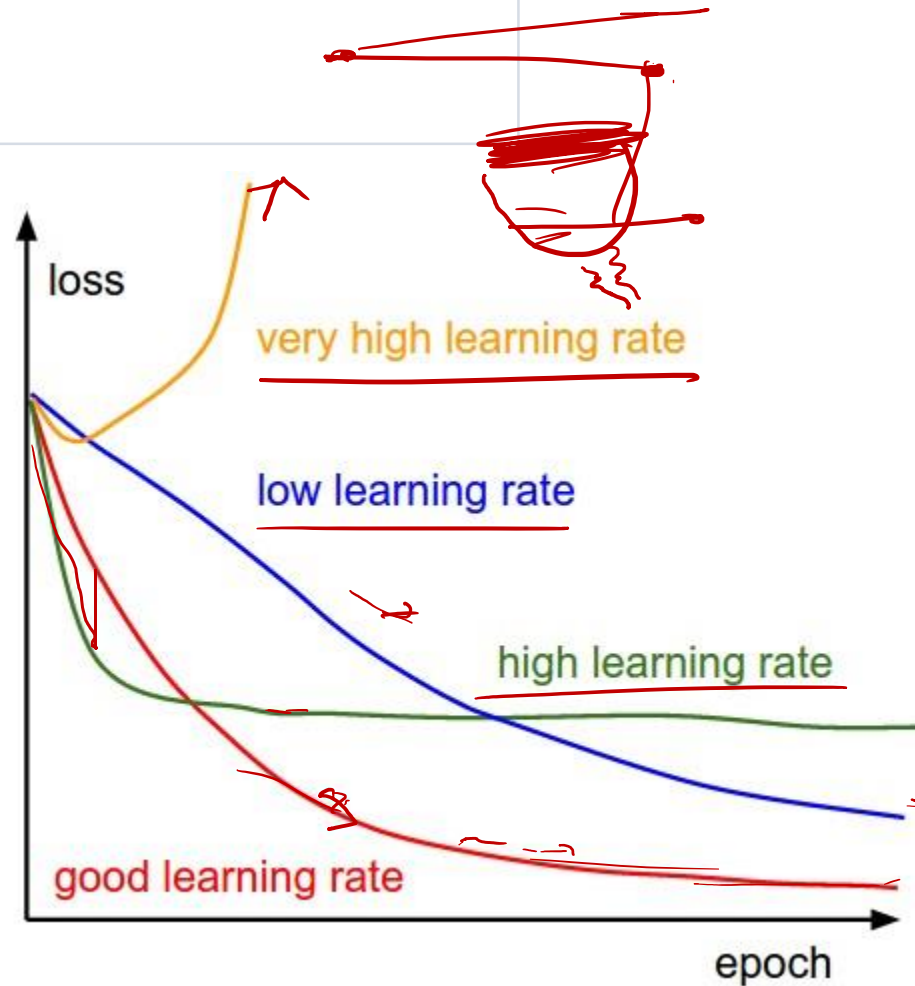
## Adam

```
velocity = beta1 * velocity + (1-beta1) * gradient
accumulated = beta2 * accumulated + (1-beta2) * gradient ** 2

adaptive_learning_rate = learning_rate / sqrt(accumulated)

w = w - adaptive_learning_rate * velocity
```

# Скорость обучения Learning rate

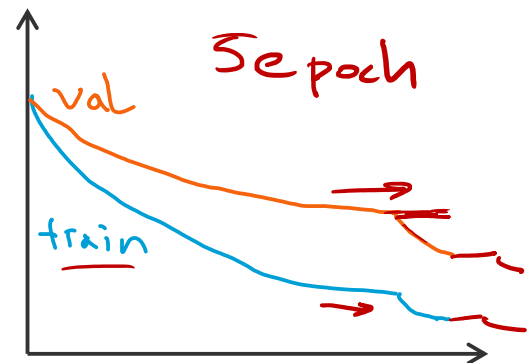


Annealing (decay)

$$\gamma = 1, \lambda = 0.99 \quad \gamma = 30 \quad \lambda = 0.5^{0.1}$$

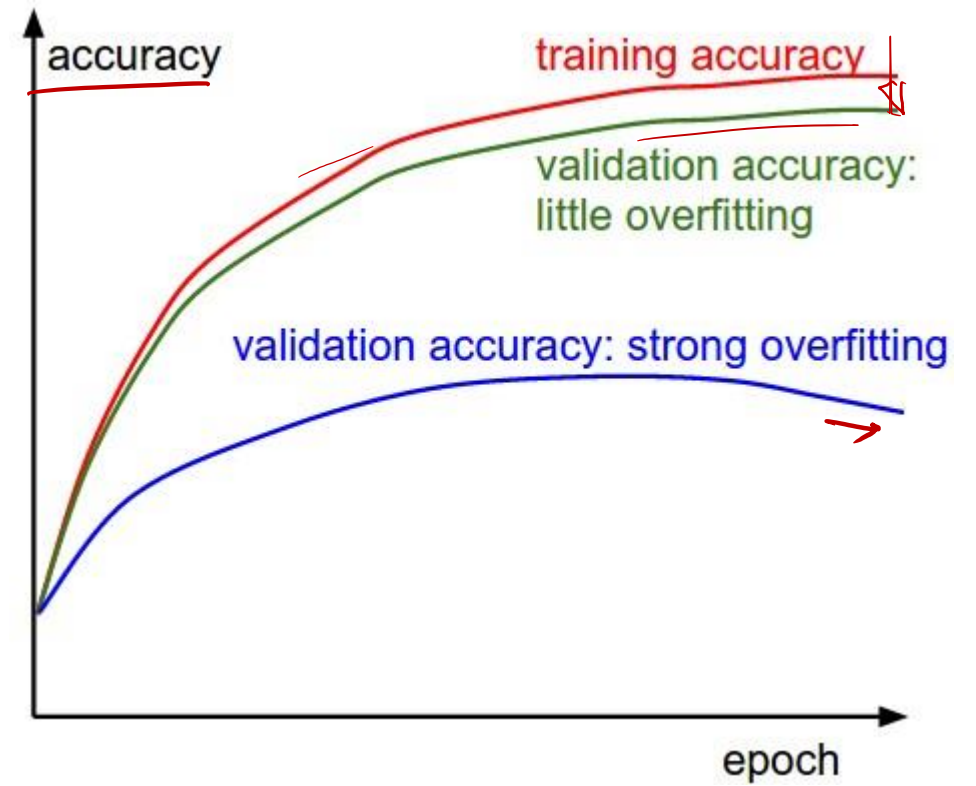
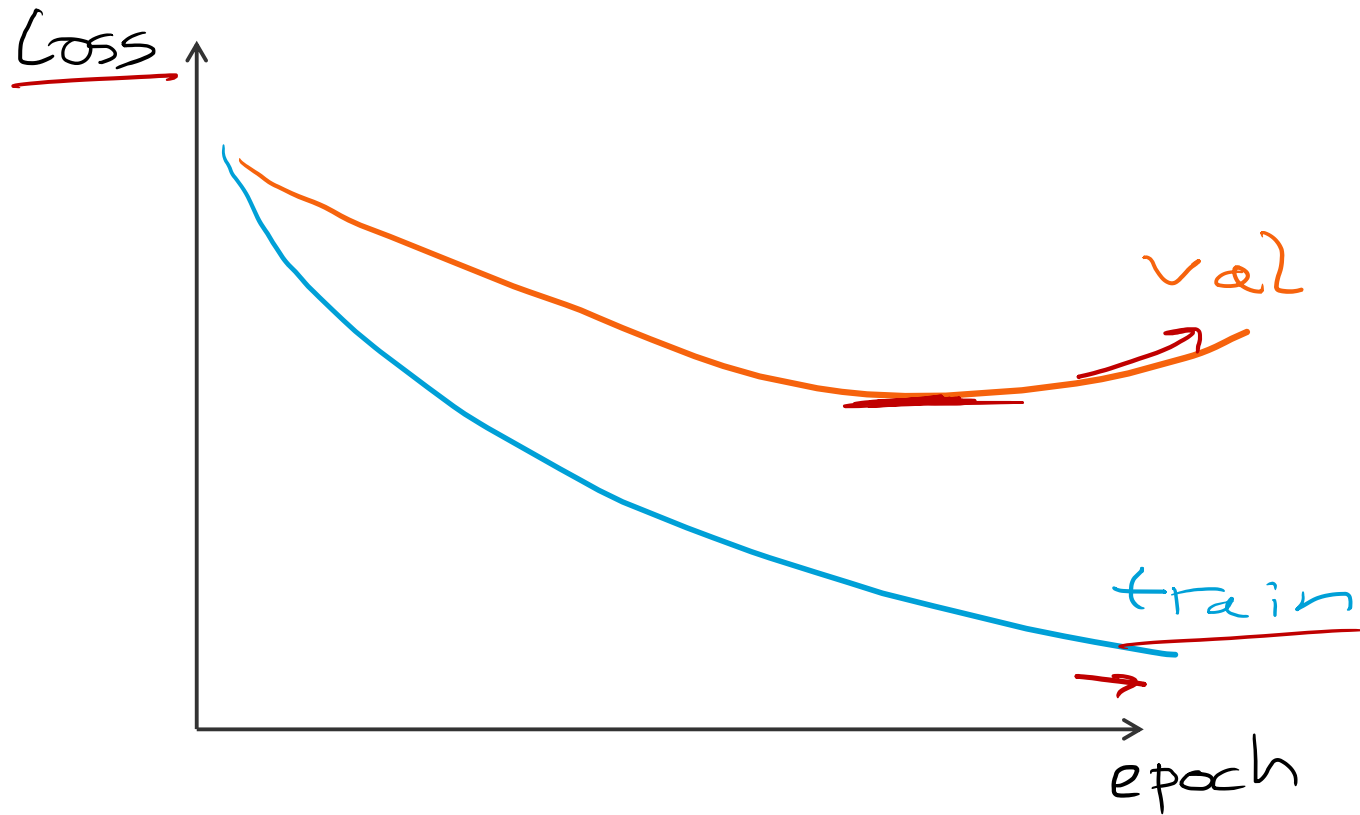
- Каждые  $X$  эпох умножать на  $\lambda$
- Умножать на  $\lambda$  на плато

0.1



`torch.optim.lr_scheduler.`

# overfitting



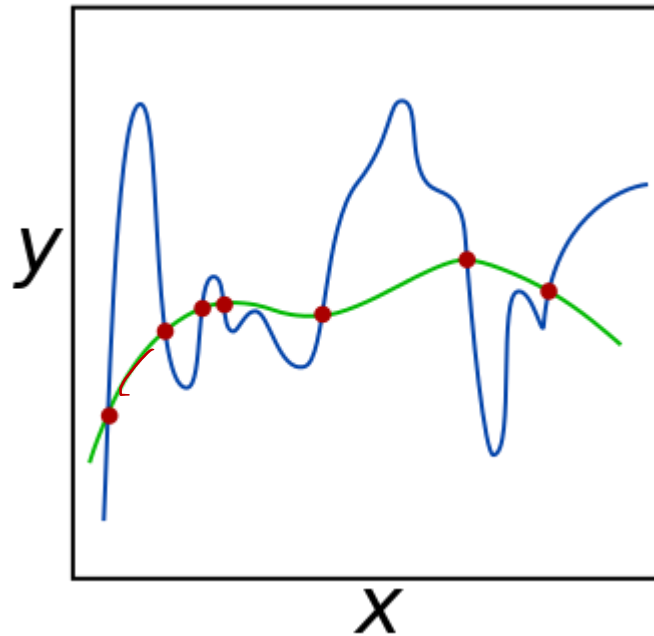
# Регуляризация Regularization

Tensor board

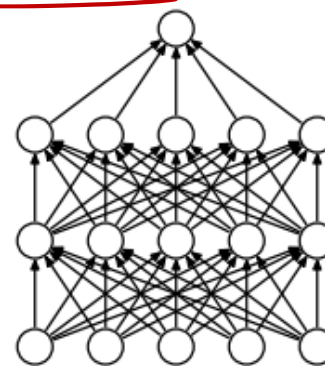
L2 regularization (или weight decay)

$$+ \lambda \|w\|_2^2$$

Dropout



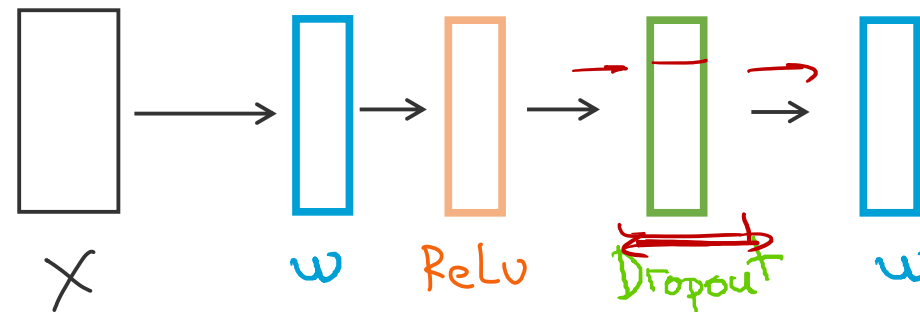
[Wikipedia](#)



(a) Standard Neural Net



(b) After applying dropout.



tail  $\rightarrow$  cat  
+ ears  $\rightarrow$

0.3

0.5

0.7



# Batch Normalization

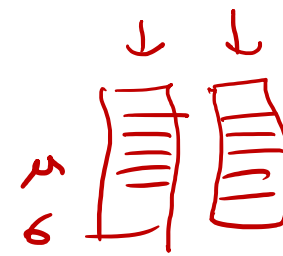
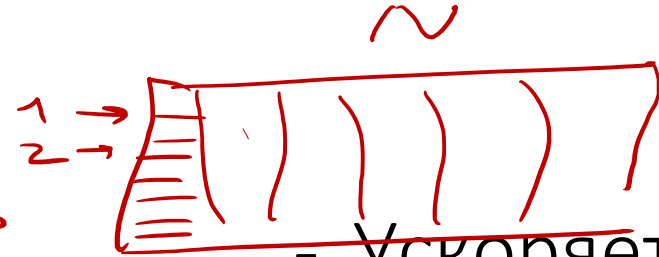
$$\hat{x}^{(k)} = \frac{x^{(k)} - \overline{E[x^{(k)}]}}{\sqrt{\overline{\text{Var}[x^{(k)}]}}}$$

bs

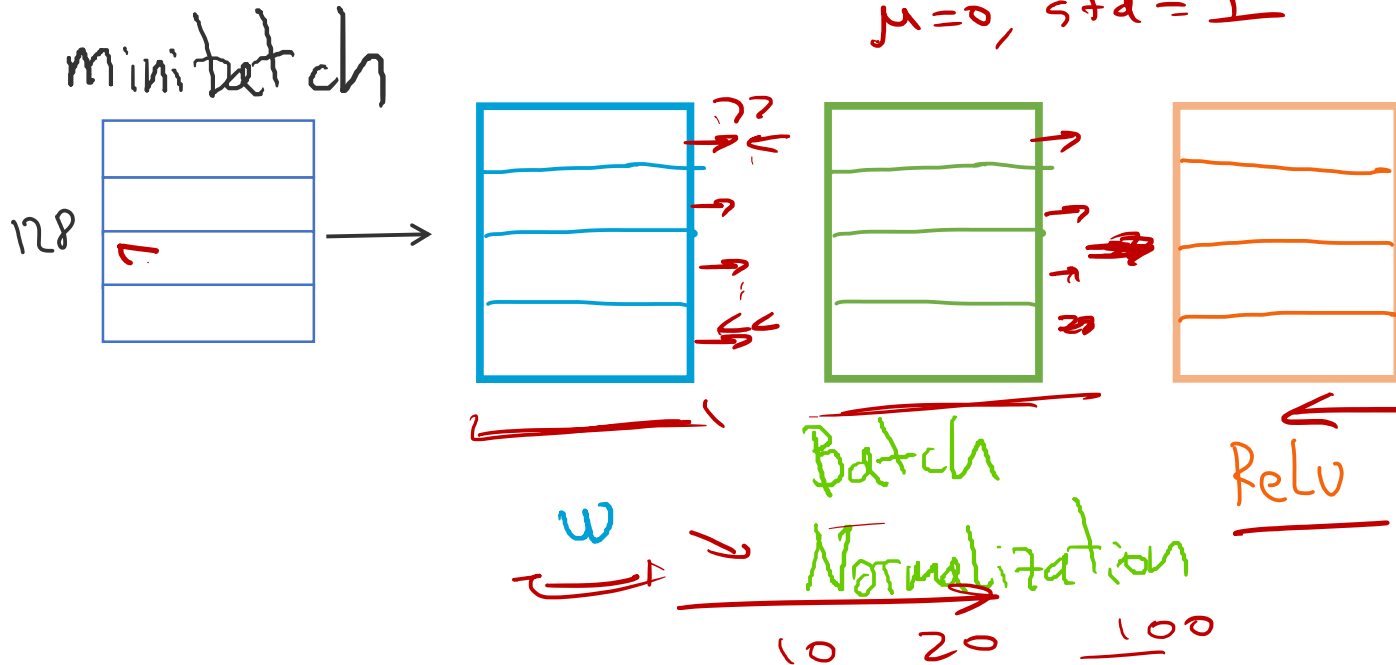
[batch-size, ~]

$$\hat{y}^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)}$$

$\mu = 0, \text{std} = 1$



- Ускоряет и стабилизирует тренировку
- Регуляризует
- Не так важна инициализация
- Счастье какое-то



# Batch Normalization – test time

$$\hat{x}^{(k)} = \frac{x^{(k)} - \underline{E[x^{(k)}]}}{\underline{\sqrt{\text{Var}[x^{(k)}]}}}$$

Используем значения среднего и дисперсии, накопленные во время тренировки

$$y^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)}$$

mini batch



w



Batch Normalization



Relu

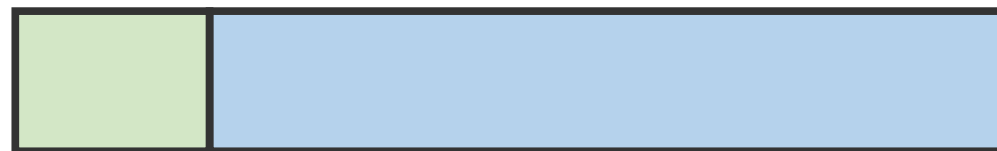
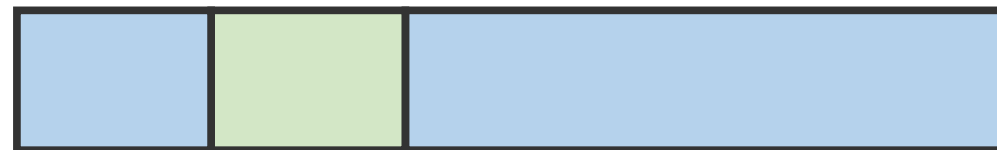
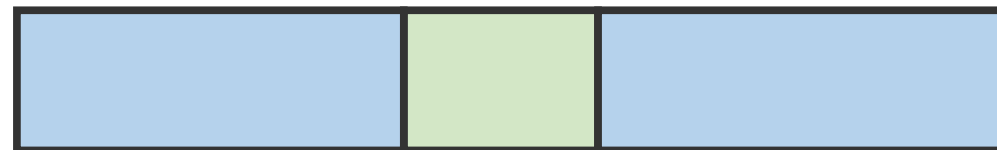
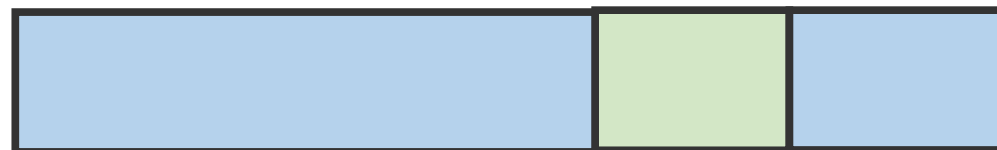
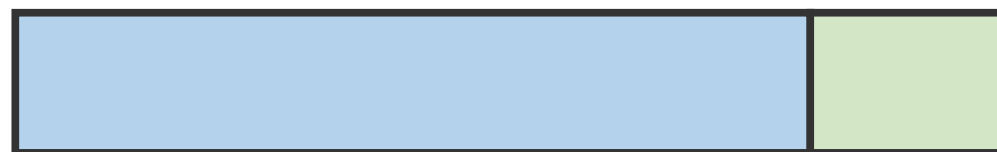
До или после ReLU?

Почитать  
Почитать еще

# Поиск гиперпараметров

## Hyperparameter search

train Val



(cross  
validation  
!)

Начальный learning rate

Коэфф. annealing

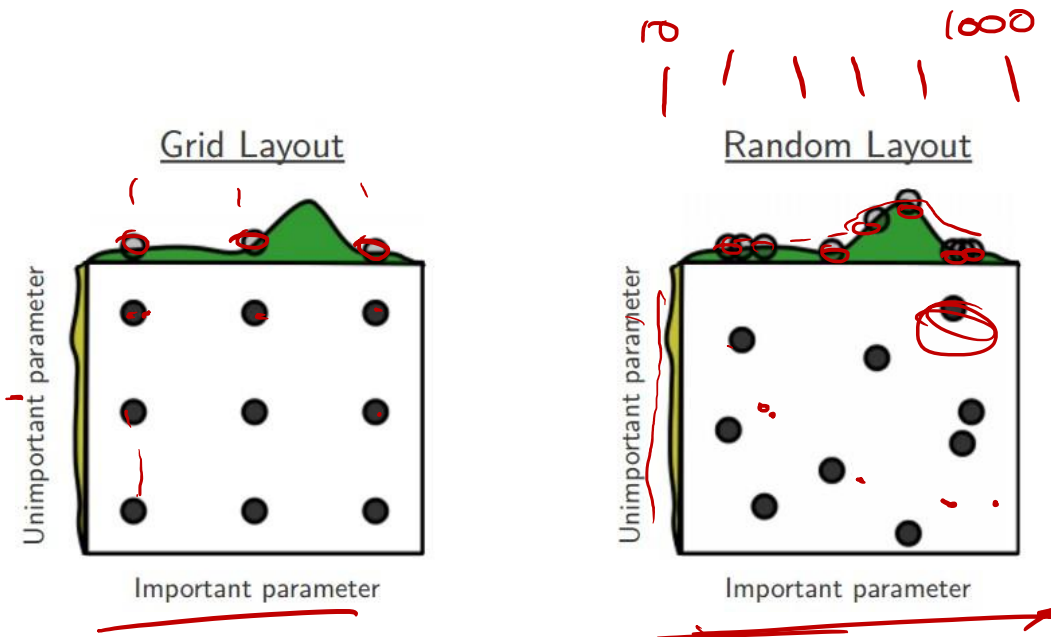
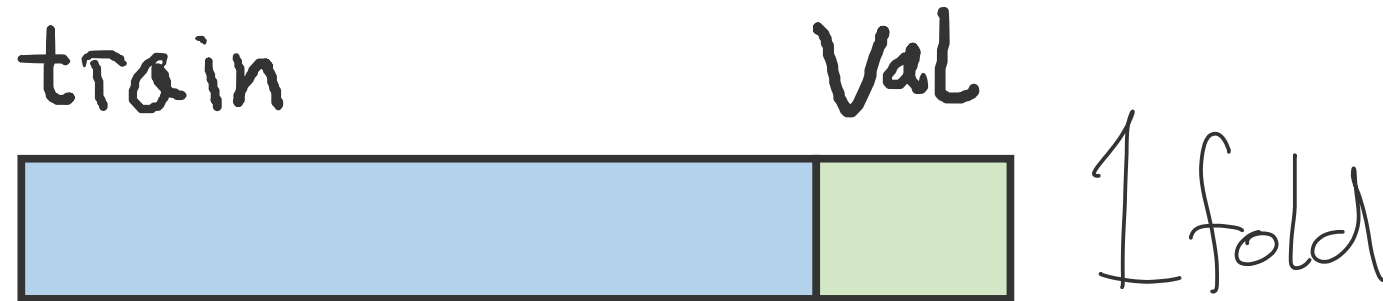
Коэфф. L2 регуляризации

Размер батча

....

# Поиск гиперпараметров

## Hyperparameter search



В любом случае  $(1 = 100\%)$

- Coarse to fine
- Log space

$lr = \text{random}(10, 10000)$  ←  $100\%$   
 $lr = 10^{**}\text{random}(1, 4)$

топово





# Короче!

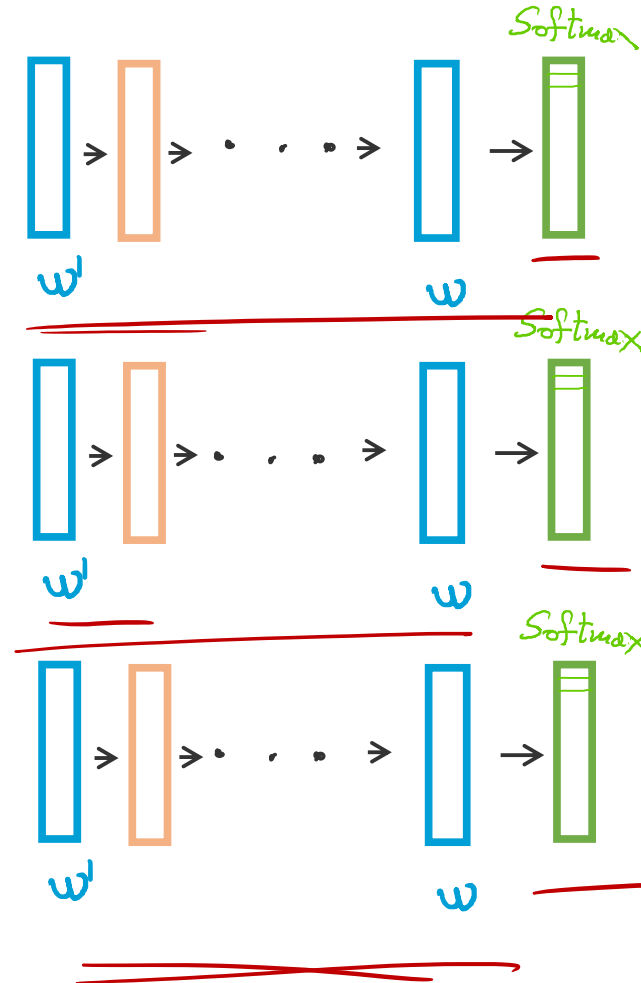
## Начинать с:

- Вычитаем среднее из данных
- Оптимизатор - Adam
- Batch Normalization – добро
- Learning rate annealing на плато
- Перебираем гиперпараметры
  - Важнее всего – learning rate
- Смотрим на графики тренировки!

# Ансамбль моделей Model ensemble

Откуда брать разнообразие?

- Разные модели и подходы
- Cross-validation folds
- Креатив!



Усреднить

↑ к точности!

# LR Range test + Cyclical Learning Rate

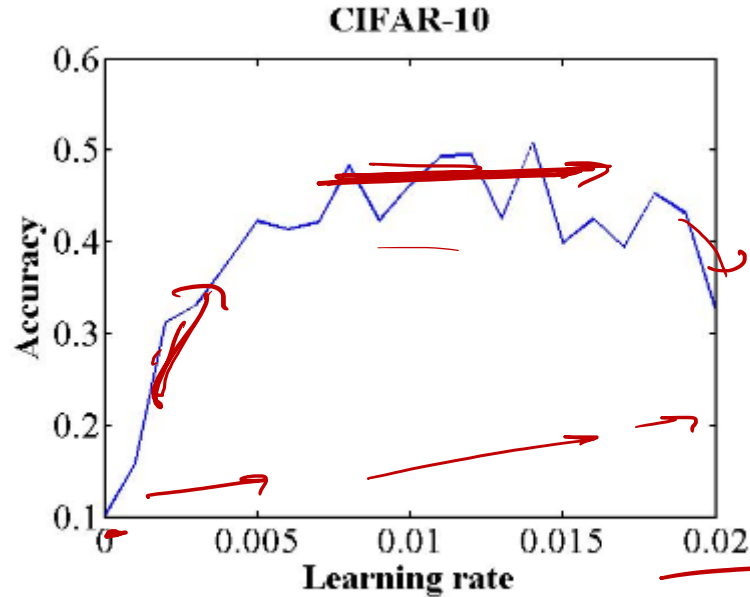


Figure 3. Classification accuracy as a function of increasing learning rate for 8 epochs (LR range test).

еур мохно  
угаубат

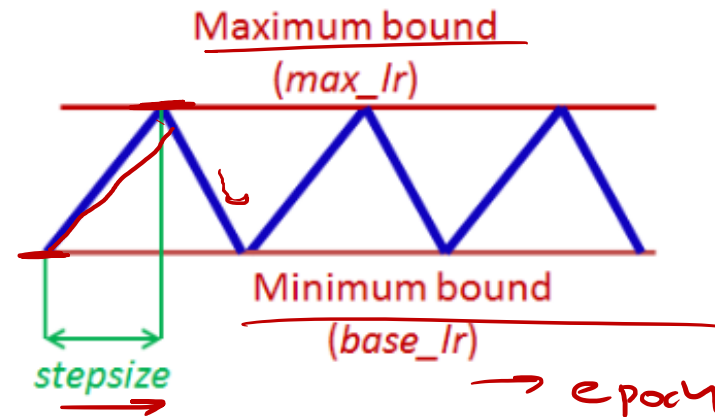


Figure 2. Triangular learning rate policy. The blue lines represent learning rate values changing between bounds. The input parameter *stepsize* is the number of iterations in half a cycle.

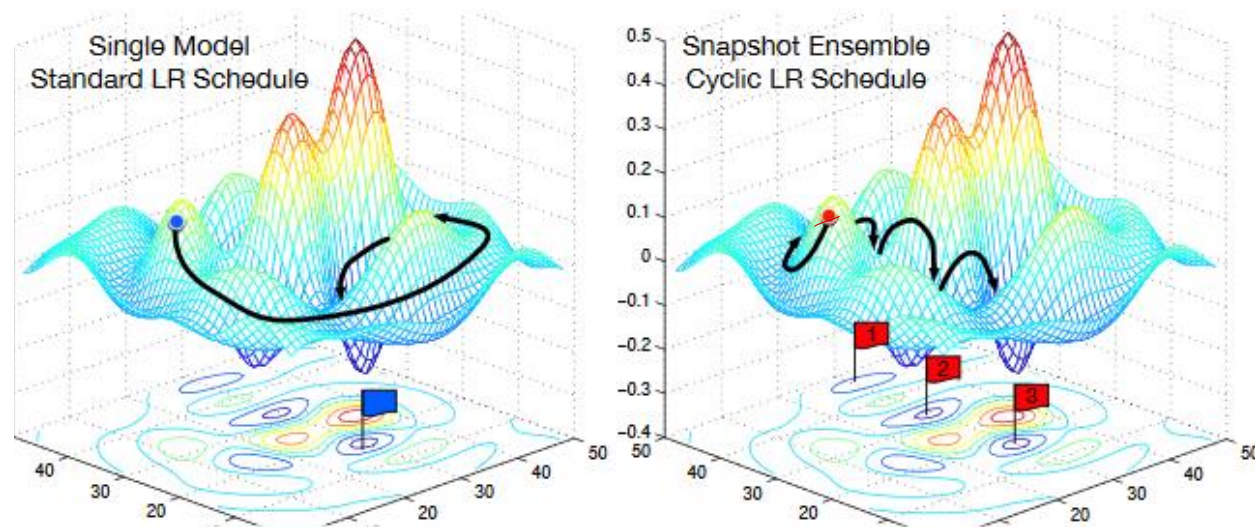
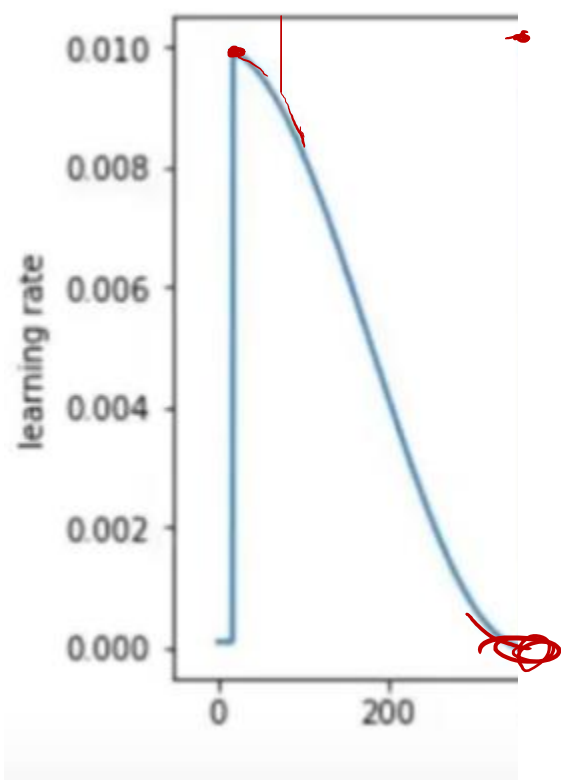
5





# Cosine learning rate

- + [Stochastic Gradient Descent with Restarts \(SGDR\)'16](#)
- + [Snapshot Ensembles'17](#)



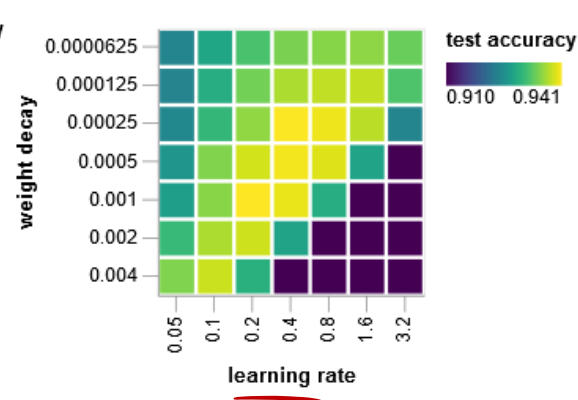
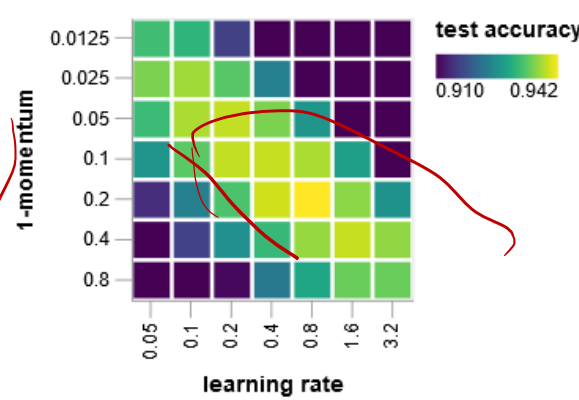
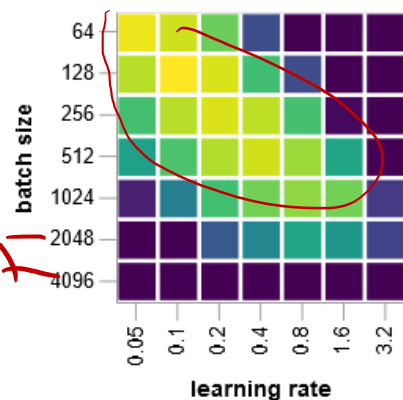
# Инженерный подход к перебору гиперпараметров

$\alpha = 125$

$\rho = \text{moment}$

$\lambda = 15$

$(\frac{\lambda\alpha}{1-\rho}, \rho, \alpha)$



train run	$\lambda$	$\rho$	$\alpha$	$\frac{\lambda\alpha}{1-\rho}$	test acc
0	0.001	0.5	.01	.00002	30.6%
10	<b>0.256</b>	0.5	.01	.00512	93.0%
13	<b>0.128</b>	<b>0.75</b>	.01	.00512	93.2%
17	<b>1.024</b>	0.75	<b>.00125</b>	.00512	93.9%
20	<b>0.512</b>	0.75	.00125	.00256	94.1%
22	<b>0.256</b>	<b>0.875</b>	.00125	.00256	94.2%

**BACK**  
**TO THE PRESENT**

**STAND BACK**

**WE'RE TRYING THIS IN PRODUCTION**

# Research vs Deployment

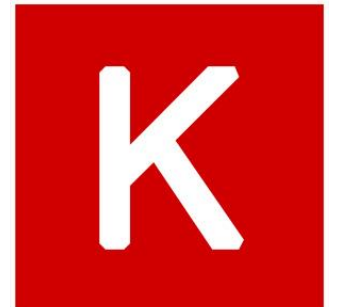
Caffe

theano



*dmlc*  
***mxnet***

Microsoft  
CNTK



PYTORCH

# Deployment



- Интегрируется не только с Python
- Высокие требования по стабильности и надежности
- Инструменты для запуска модели на сервере
  - Например, [Tensorflow Serving](#)
- Оптимизировано под железо, где запускается модель
  - Например, телефоны



Что делать, если тренировал в одном фреймворке, а в прод надо другой?

Конвертуем



ONNX



