



Optimization and Training Techniques

Qinliang Su (苏勤亮)

Sun Yat-sen University

suqliang@mail.sysu.edu.cn

Outline

- Optimization Algorithms
- Training Techniques

Stochastic Gradient Descent (SGD)

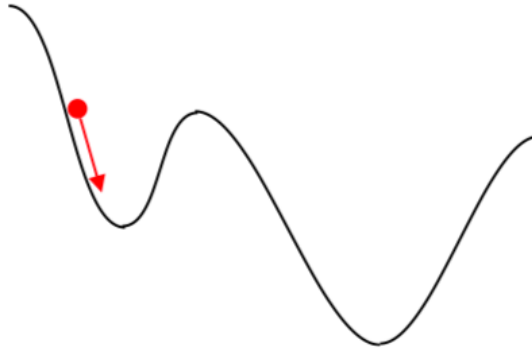
- Suppose \mathbf{w} is the parameter to optimize. The updating in SGD is

$$\mathbf{w}_{t+1} = \mathbf{w}_t - lr * \nabla f(\mathbf{w}_t)$$

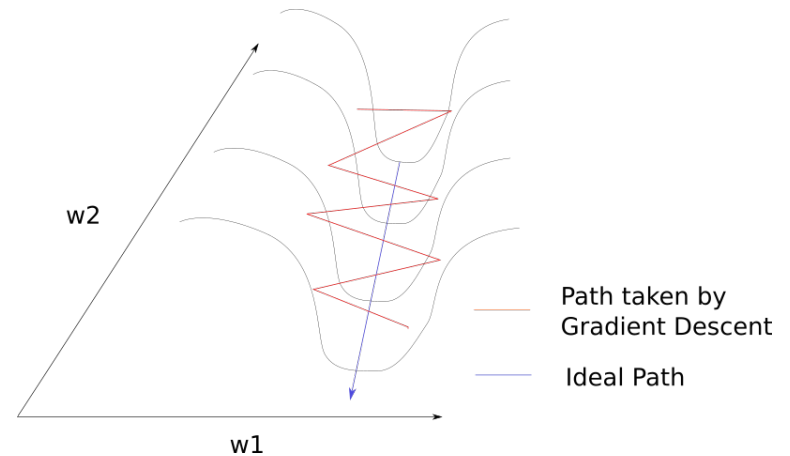
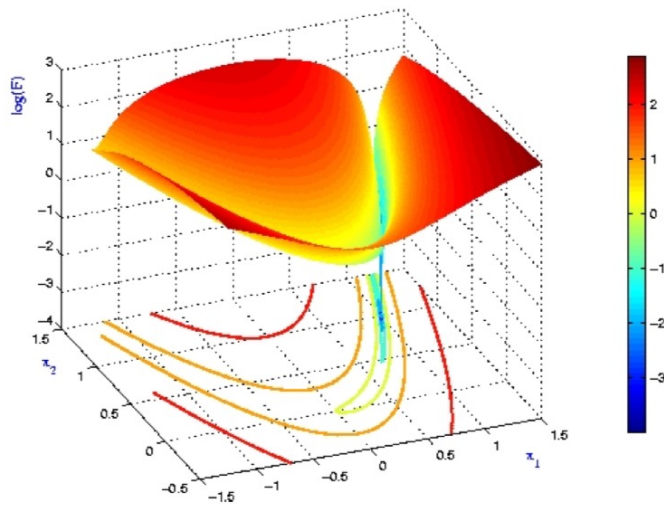
- $\nabla f(\mathbf{w})$ is the stochastic gradient of the loss \mathcal{L} , estimated with only partial training data
- lr is the learning rate

- Problems with SGD

- Very easy to get stuck at local optima



- Slow convergence due to the **pathological curvature** in NNs

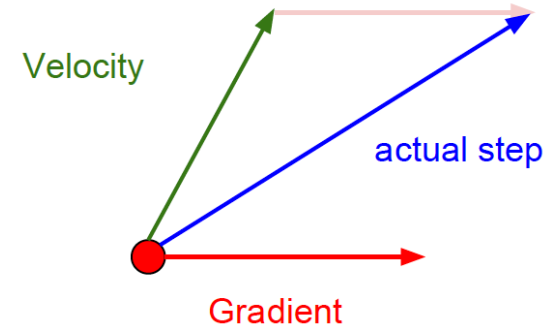


SGD + Momentum

- Updating equation

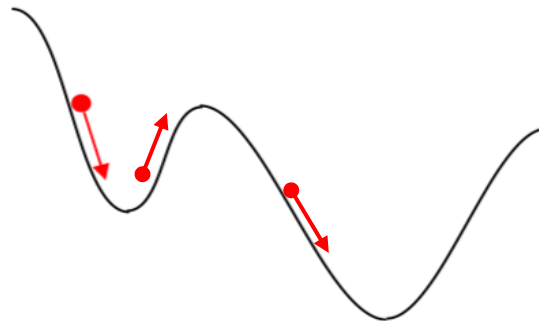
$$\mathbf{v}_t = \rho \mathbf{v}_{t-1} + \nabla f(\mathbf{w}_t)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - lr \cdot \mathbf{v}_t$$



where $\rho \in (0, 1)$ is the decay rate, often chosen as 0.9, 0.95, 0.99

The updating process can be *understood as a ball falling down from a rugged hillside*, with the friction strength proportional to ρ



RMSProp

- Updating equation

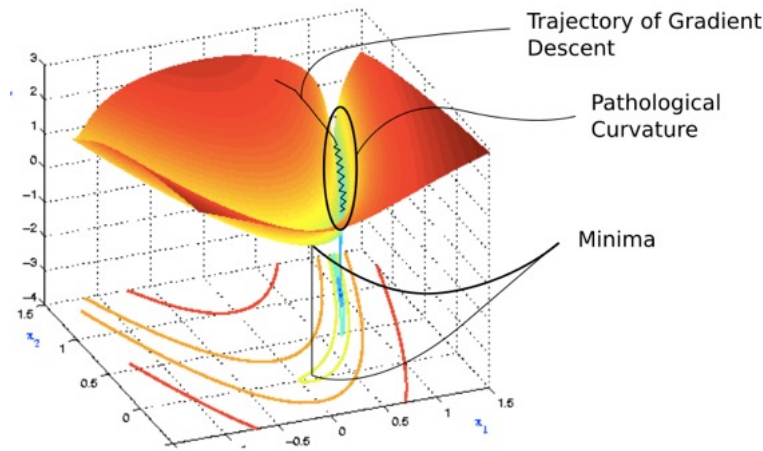
Moving average

$$\mathbf{s}_t = \rho \cdot \mathbf{s}_{t-1} + (1 - \rho)(\nabla f(\mathbf{w}_t))^2$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - lr \cdot \nabla f(\mathbf{w}_t) \oslash \sqrt{\mathbf{s}_t}$$

where \mathbf{s}_t is the weighted average of squared gradients until the current time step; \oslash denotes elementwise division

It *rescales the gradients of different dimensions to the same level*, alleviating the pathological curvature problem



Supplement: Moving average

$$\mathbf{s}_1 = \rho \cdot \mathbf{s}_0 + (1 - \rho)(\nabla f(\mathbf{w}_1))^2$$

$$\mathbf{s}_2 = \rho^2 \cdot \mathbf{s}_0 + \rho(1 - \rho)(\nabla f(\mathbf{w}_1))^2 + (1 - \rho)(\nabla f(\mathbf{w}_2))^2$$

$$\mathbf{s}_3 = \rho^3 \cdot \mathbf{s}_0 + \rho^2(1 - \rho)(\nabla f(\mathbf{w}_1))^2 + \rho(1 - \rho)(\nabla f(\mathbf{w}_2))^2 + (1 - \rho)(\nabla f(\mathbf{w}_3))^2$$

\vdots

Weighting

The weights are decaying for earlier records

It can be verified that all weights are summed to be 1, that is

$$(1 - \rho) + \rho(1 - \rho) + \rho^2(1 - \rho) + \rho^3(1 - \rho) + \dots$$

$$= \frac{1 - \rho^n}{1 - \rho} (1 - \rho) \rightarrow 1$$

Adam

- Updating equation

$$\mathbf{m}_t = \beta_1 \cdot \mathbf{m}_{t-1} + (1 - \beta_1) \nabla f(\mathbf{w}_t)$$

$$\mathbf{s}_t = \beta_2 \cdot \mathbf{s}_{t-1} + (1 - \beta_2) (\nabla f(\mathbf{w}_t))^2$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - lr \cdot \mathbf{m}_t \oslash \sqrt{\mathbf{s}_t}$$

Same as RMSProp,
except the \mathbf{m}_t

where \mathbf{m}_t is the moving average of the past gradients

A modified version of RMSProp, with *the the exact gradient $\nabla f(\mathbf{w}_t)$ replaced by its moving average \mathbf{m}_t*

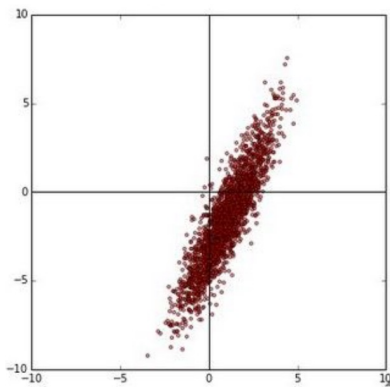
- Default settings: $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $lr = 10^{-3}$
- By setting $\beta_1 = 0$, Adam is degenerated to RMSProp

Outline

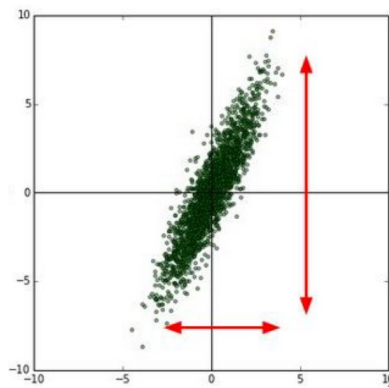
- Optimization Algorithms
- Training Techniques

Preprocessing

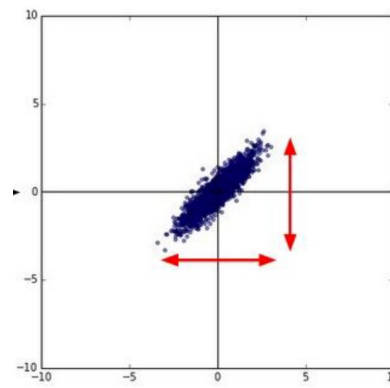
- Commonly used data preprocessing methods:
 - 1) **Centering** the data, *i.e.* subtracting the data mean
 - 2) **Normalizing the variance** for each dimension
 - 3) **Whitening** the data



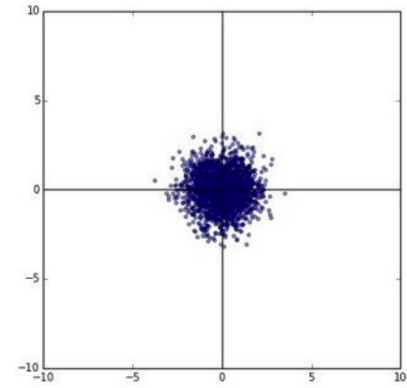
Original



Centered



Variance Normalized



Whitened

- Up to the characteristics of data, we may use one or several methods above to preprocess the data

- Why centering and normalizing operations are helpful?

Let's take the linear regression problem as the example

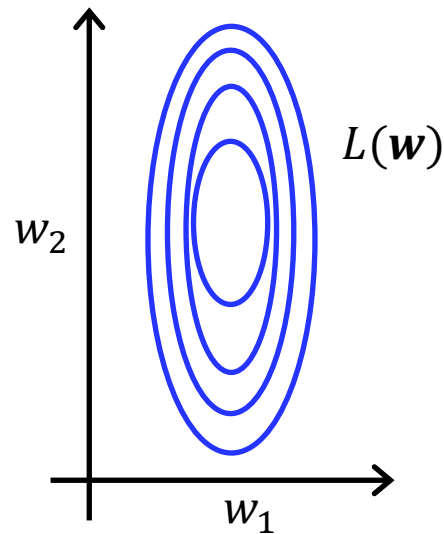
- The features in original data have very different scales

Size (feet) x_1	# bedrooms x_2	# floors x_3	# years (Ages) x_4	Price (\$ 1000) y
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178

- Under the case, the loss function is

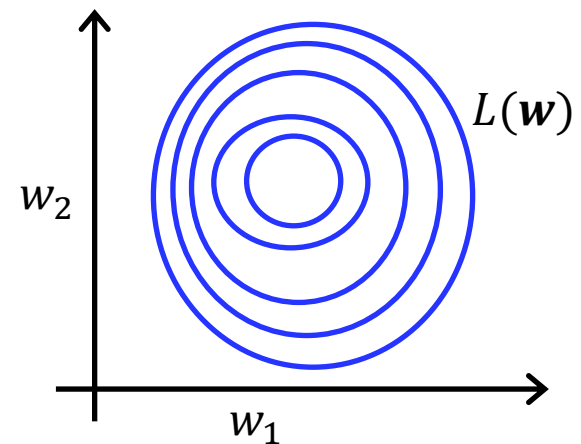
$$L(\mathbf{w}) = \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

- Obviously, the cost function $L(\mathbf{w})$ decreases at different rates for different dimensions



This results in the *slow convergence* of gradient descent. *Why??*

- If the features are preprocessed to similar scales, then the loss function looks like



Initialization

- 1) Initializing the weights of all layers by the random samples drawn from fixed Gaussian distribution, *e.g.*

$$\mathbf{W} \sim \mathcal{N}(0, 0.01^2)$$

- It works well for neural networks that are not too deep (*e.g.* layer # < 8)
- If ReLU is used as the activation function, some people also suggest to use the truncated Gaussian distribution, *i.e.*

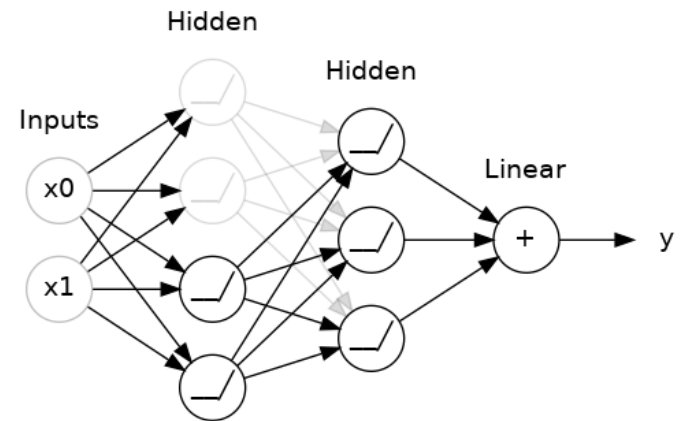
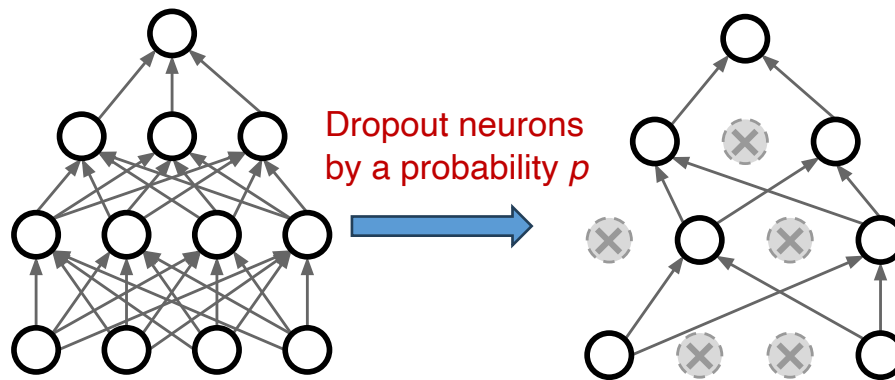
$$\mathbf{W} \sim \mathcal{N}_T(0, 0.01^2)$$

- 2) When the network goes deeper, it can then be initialized by a Gaussian $\mathbf{W} \sim \mathcal{N}(0, \sigma_i^2)$ with the variance of each layer *inversely proportional to the number of neurons* in the previous layer, *i.e.*

$$\sigma_i = \frac{1}{\sqrt{\# \text{ neurons in } (i - 1)\text{th layer}}}$$

Dropout

- In each forward pass, randomly set some neurons to zero
 - Probability of dropping neurons is a hyperparameter; 0.5 is a commonly used value

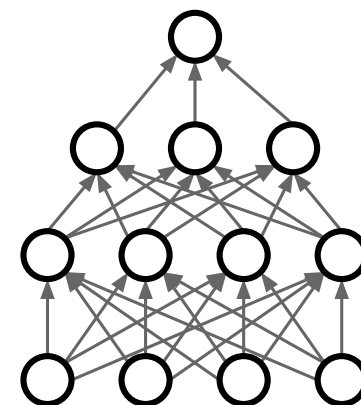


Dropout can effectively alleviate the overfitting issue commonly encountered in deep neural networks

- Dropout at testing time

- During testing, it is preferred to have model outputting a fixed value for an input

- Thus, the randomness on the neurons should be removed by turning on all of them at the testing stage



- To keep the input to each neuron consistent with that at training, we need to rescale the value with the dropout probability p

$$\begin{array}{ccc} \mathbf{h}_{\ell+1} = \sigma(\mathbf{W}_{\ell} \mathbf{h}_{\ell-1}) & \longrightarrow & \mathbf{h}_{\ell+1} = \sigma(\mathbf{p} \cdot \mathbf{W}_{\ell} \mathbf{h}_{\ell-1}) \\ \text{w/o dropout} & & \text{w/ dropout} \end{array}$$

Batch Normalization

- Although a good initialization could have the network locating at a normal state at the beginning, its influence dims as the training process proceeds
- A more effective way to keep the network locating at a normal state is to **intentionally normalize** the hidden states $\{\tilde{\mathbf{h}}_\ell^{(i)}\}_{i=1}^N$ to a standard normal distribution **during the training** as

$$\hat{\mathbf{h}}_\ell^{(i)} = \left(\tilde{\mathbf{h}}_\ell^{(i)} - \boldsymbol{\mu}_\ell \right) \oslash \boldsymbol{\sigma}_\ell$$

where

$$\boldsymbol{\mu}_\ell = \frac{1}{N} \sum_{i=1}^N \tilde{\mathbf{h}}_\ell^{(i)} \quad \boldsymbol{\sigma}_\ell^2 = \frac{1}{N} \sum_{i=1}^N \left(\tilde{\mathbf{h}}_\ell^{(i)} - \boldsymbol{\mu}_\ell \right)^2$$

- Then, the normalized $\hat{\mathbf{h}}_\ell^{(i)}$ is further transformed by a scaling coefficient γ_ℓ and an offset coefficient β_ℓ

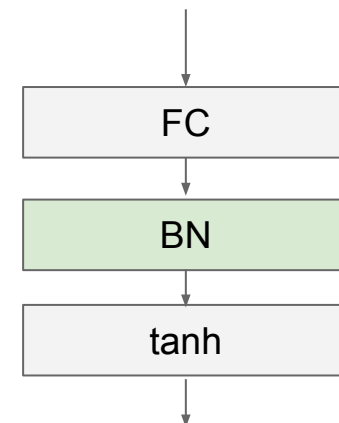
$$\mathbf{z}_\ell^{(i)} = \gamma_\ell \odot \hat{\mathbf{h}}_\ell^{(i)} + \beta_\ell$$

- It can be easily verified that $\mathbf{z}_\ell^{(i)}$ will recover the input state $\tilde{\mathbf{h}}_\ell^{(i)}$ if γ_ℓ and β_ℓ are set as σ_ℓ and μ_ℓ , respectively
- The coefficients γ_ℓ and β_ℓ are learned with the model weights simultaneously
- Computing the exact mean μ_ℓ and variance σ_ℓ^2 is expensive. In practice, they are estimated by the current mini-batch \mathcal{B}

$$\mu_\ell = \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \tilde{\mathbf{h}}_\ell^{(i)} \quad \sigma_\ell^2 = \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \left(\tilde{\mathbf{h}}_\ell^{(i)} - \mu_\ell \right)^2$$

- At the testing stage, it is undesirable to see the output relying on which minibatch the input is associated with
- There are two solutions:
 - Computing the exact mean $\mu_\ell = \frac{1}{N} \sum_{i=1}^N \tilde{\mathbf{h}}_\ell^{(i)}$ and variance $\sigma_\ell^2 = \frac{1}{N} \sum_{i=1}^N \left(\tilde{\mathbf{h}}_\ell^{(i)} - \mu_\ell \right)^2$ over the whole training dataset. Only need to compute once
 - Using moving average to estimate the exact mean and variance from the mini-batch based mean and variance
- The batch normalization can be viewed as a layer denoted by $BN(\cdot)$, that is,

$$\gamma_\ell \odot \left(\tilde{\mathbf{h}}_\ell^{(i)} - \mu_\ell \right) \oslash \sigma_\ell + \beta_\ell = BN \left(\tilde{\mathbf{h}}_\ell^{(i)} \right)$$



- Batch normalization in MLP vs CNNs

Batch Normalization for
fully-connected networks

$$\begin{array}{l}
 \mathbf{x}: \mathbf{N} \times \mathbf{D} \\
 \text{Normalize} \quad \downarrow \\
 \boldsymbol{\mu}, \boldsymbol{\sigma}: 1 \times \mathbf{D} \\
 \boldsymbol{\gamma}, \boldsymbol{\beta}: 1 \times \mathbf{D} \\
 \mathbf{y} = \boldsymbol{\gamma}(\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta}
 \end{array}$$

Batch Normalization for
convolutional networks
(Spatial Batchnorm, BatchNorm2D)

$$\begin{array}{l}
 \mathbf{x}: \mathbf{N} \times \mathbf{C} \times \mathbf{H} \times \mathbf{W} \\
 \text{Normalize} \quad \downarrow \quad \downarrow \quad \downarrow \\
 \boldsymbol{\mu}, \boldsymbol{\sigma}: 1 \times \mathbf{C} \times 1 \times 1 \\
 \boldsymbol{\gamma}, \boldsymbol{\beta}: 1 \times \mathbf{C} \times 1 \times 1 \\
 \mathbf{y} = \boldsymbol{\gamma}(\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta}
 \end{array}$$

- Layer normalization

Batch Normalization for
fully-connected networks

$$\begin{array}{l}
 \mathbf{x}: \mathbf{N} \times \mathbf{D} \\
 \text{Normalize} \quad \downarrow \\
 \boxed{\boldsymbol{\mu}, \boldsymbol{\sigma}: 1 \times \mathbf{D}} \\
 \gamma, \beta: 1 \times \mathbf{D} \\
 \mathbf{y} = \gamma(\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \beta
 \end{array}$$



Layer Normalization for
fully-connected networks
Same behavior at train and test!
Can be used in recurrent networks

$$\begin{array}{l}
 \mathbf{x}: \mathbf{N} \times \mathbf{D} \\
 \text{Normalize} \quad \downarrow \\
 \boxed{\boldsymbol{\mu}, \boldsymbol{\sigma}: \mathbf{N} \times 1} \\
 \gamma, \beta: 1 \times \mathbf{D} \\
 \mathbf{y} = \gamma(\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \beta
 \end{array}$$

- Instance normalization

Batch Normalization for
convolutional networks

$$\begin{array}{l}
 \mathbf{x}: \mathbf{N} \times \mathbf{C} \times \mathbf{H} \times \mathbf{W} \\
 \text{Normalize} \quad \downarrow \quad \downarrow \quad \downarrow \\
 \boldsymbol{\mu}, \boldsymbol{\sigma}: \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1} \\
 \boldsymbol{\gamma}, \boldsymbol{\beta}: \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1} \\
 \mathbf{y} = \boldsymbol{\gamma}(\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta}
 \end{array}$$



Instance Normalization for
convolutional networks
Same behavior at train / test!

$$\begin{array}{l}
 \mathbf{x}: \mathbf{N} \times \mathbf{C} \times \mathbf{H} \times \mathbf{W} \\
 \text{Normalize} \quad \quad \quad \downarrow \quad \downarrow \\
 \boldsymbol{\mu}, \boldsymbol{\sigma}: \mathbf{N} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1} \\
 \boldsymbol{\gamma}, \boldsymbol{\beta}: \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1} \\
 \mathbf{y} = \boldsymbol{\gamma}(\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta}
 \end{array}$$

Hyper-parameter Tuning

- There are a lot of parameters that need to be set before starting the training, such as the learning rate, decay parameters in RMSprop and Aadm, dropout rate *etc*
- To set appropriate values for these hyper-parameters, we need to split the whole dataset into three parts, *i.e.* training set, validation set, and test set



Training set



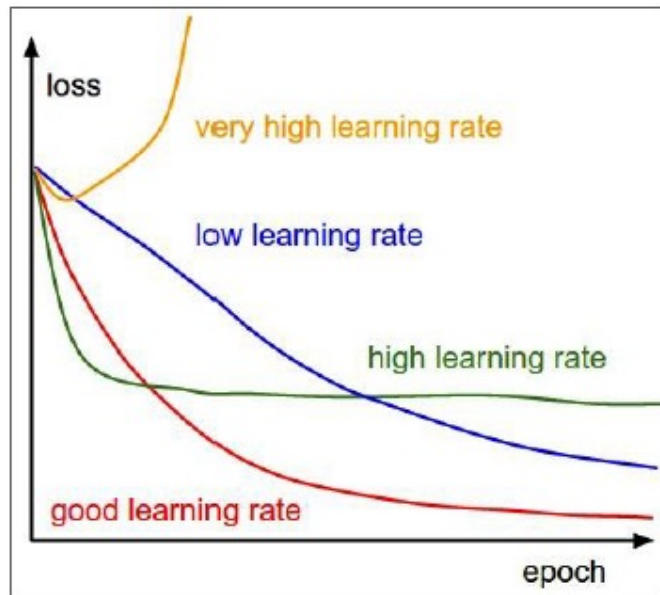
Validation set



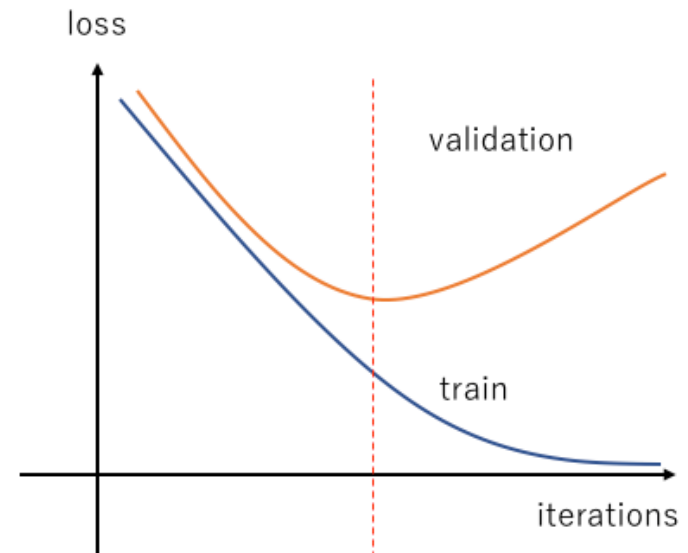
Test set

- **Train** the model on the **training set**
- **Tune** the hyper-parameters according to the performance on the **validation set**
- **Test** the performance on the **test set**, and report it as the final performance

- The hyper-parameters are not tuned blindly, but instead can be chosen according to some guidelines
 - 1) First try some default values that perform well on many tasks
 - 2) Choose values based on phenomenon observed during training



The training loss curve implies the possible range of good learning rate



The relation between performance on training and validation sets indicates the model is too flexible