

# Computer Vision: Training neural networks

Siheng Chen 陈思衡

# Training neural networks

Stochastic gradient descent

Dropout

Initialization

Hyperparameter tuning

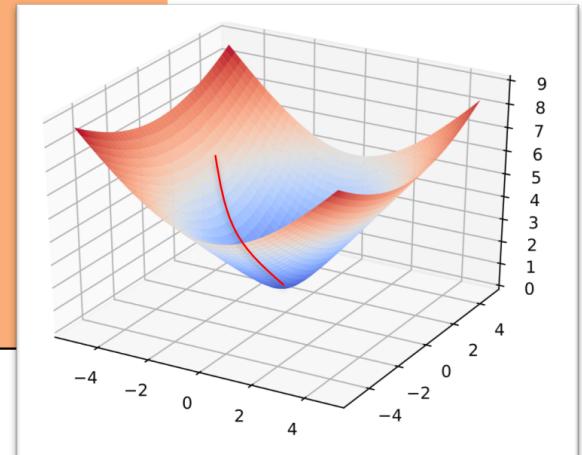
Batch normalization

# Stochastic gradient descent

## Gradient descent

### Algorithm 1 Gradient Descent

```
1: procedure GD( $\mathcal{D}$ ,  $\theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$ 
3:   while not converged do
4:      $\theta \leftarrow \theta - \gamma \nabla_{\theta} J(\theta)$ 
5:   return  $\theta$ 
```

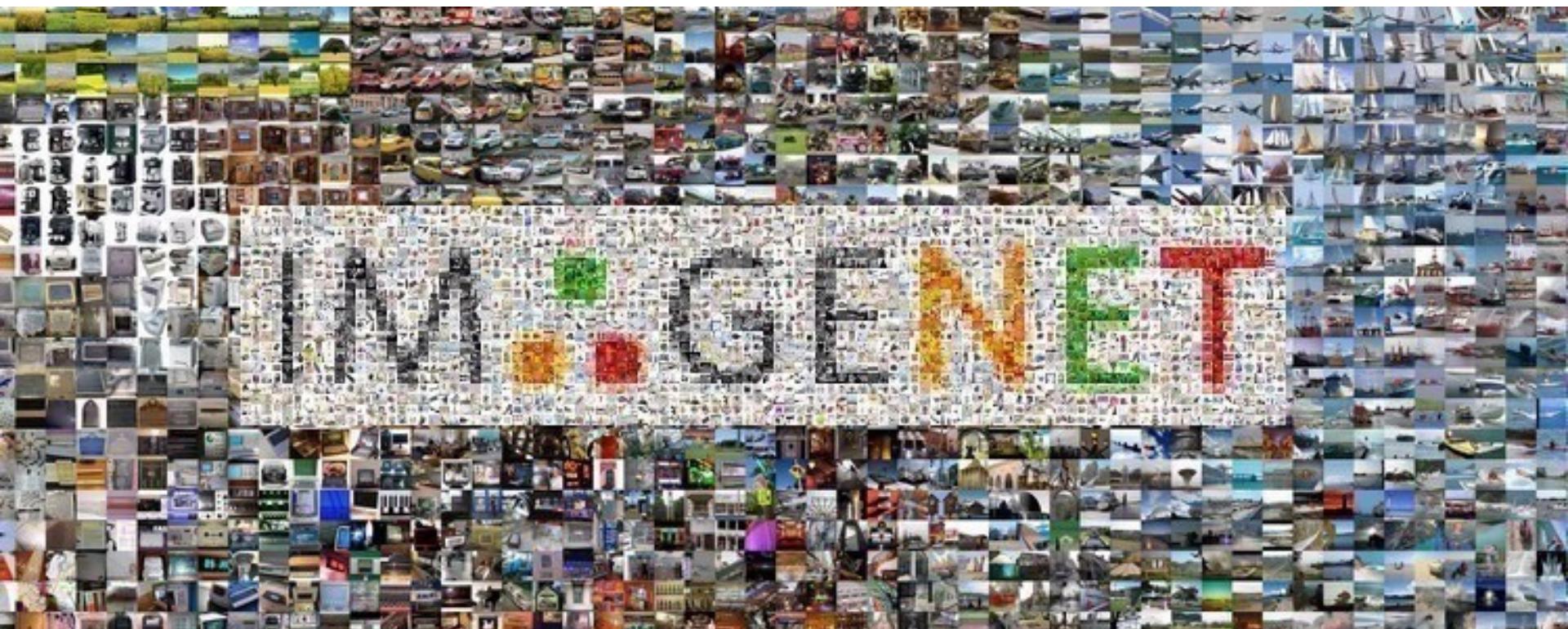


original objective:

$$J(\theta) = \sum_{i=1}^N J^{(i)}(\theta)$$

# Stochastic gradient descent

Gradient descent



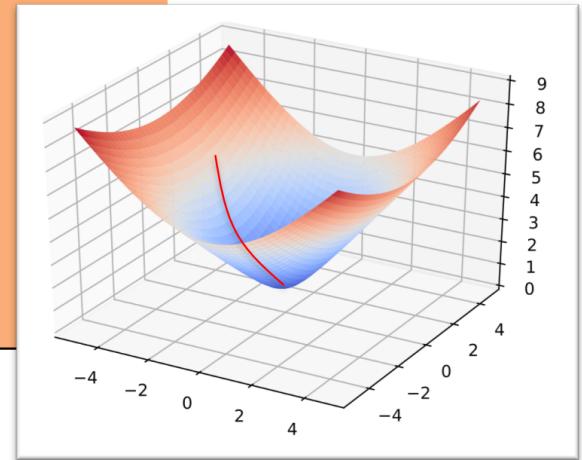
14,197,122 images

# Stochastic gradient descent

## Gradient descent

### Algorithm 1 Gradient Descent

```
1: procedure GD( $\mathcal{D}$ ,  $\theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$ 
3:   while not converged do
4:      $\theta \leftarrow \theta - \gamma \nabla_{\theta} J(\theta)$ 
5:   return  $\theta$ 
```



original objective:

$$J(\theta) = \sum_{i=1}^N J^{(i)}(\theta)$$

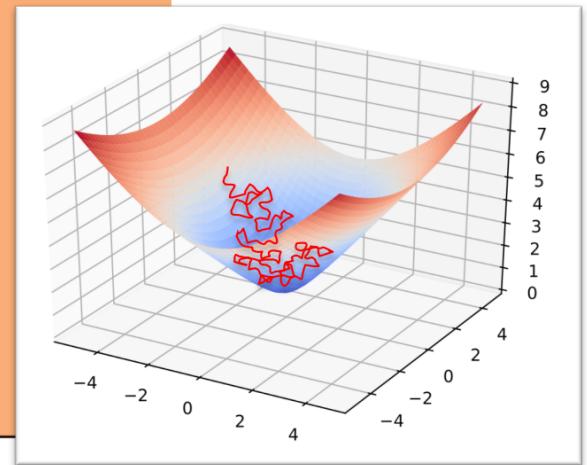
Need to enumerate all the training samples every iteration. Expensive!

# Stochastic gradient descent

## Stochastic gradient descent

### Algorithm 2 Stochastic Gradient Descent (SGD)

```
1: procedure SGD( $\mathcal{D}$ ,  $\theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$ 
3:   while not converged do
4:      $i \sim \text{Uniform}(\{1, 2, \dots, N\})$ 
5:      $\theta \leftarrow \theta - \gamma \nabla_{\theta} J^{(i)}(\theta)$ 
6:   return  $\theta$ 
```



per-example objective:

$$J^{(i)}(\theta)$$

original objective:

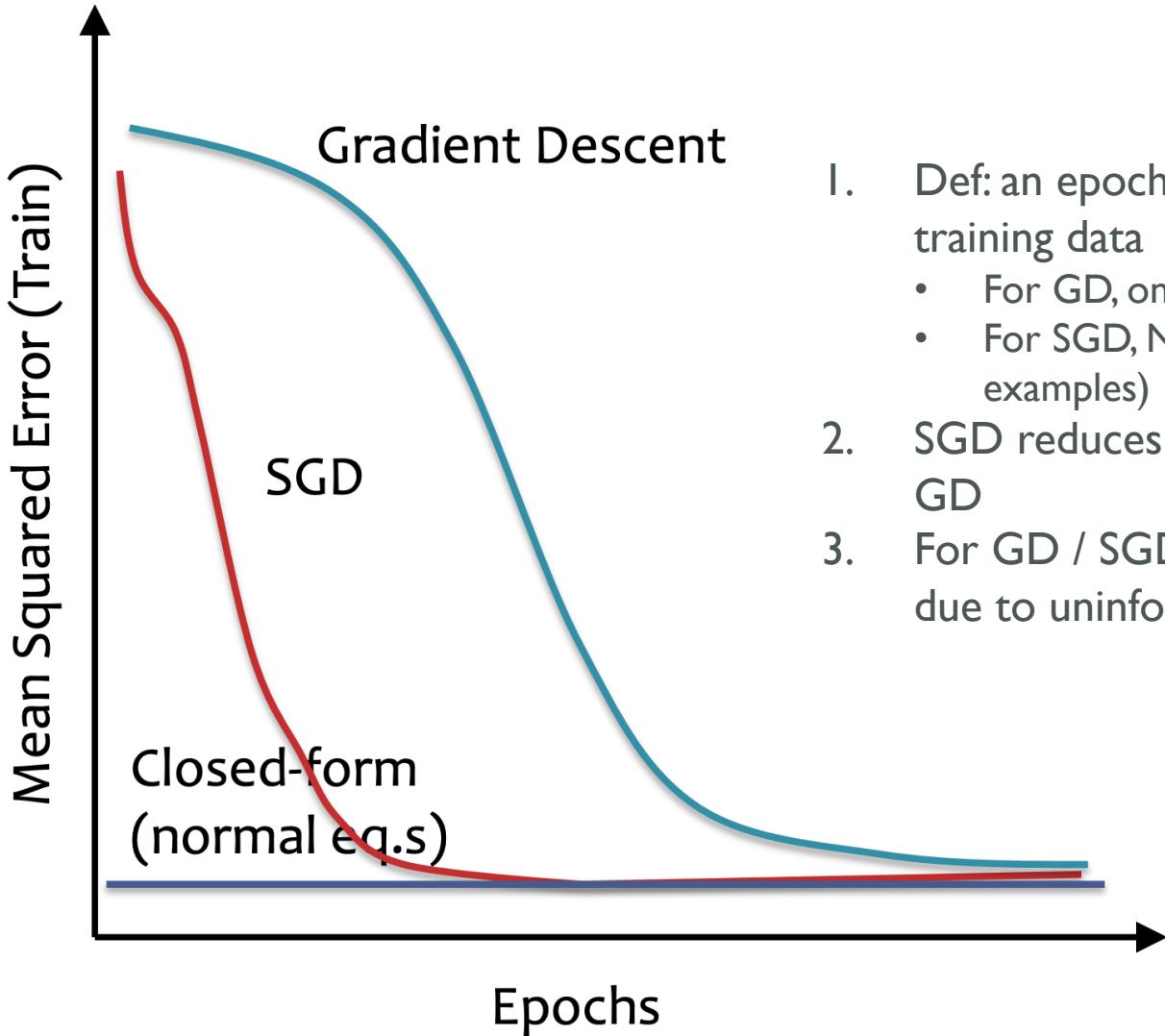
$$J(\theta) = \sum_{i=1}^N J^{(i)}(\theta)$$

# Stochastic gradient descent

Stochastic gradient descent

$$\frac{\partial J(\vec{\theta})}{\partial \theta_j} = \frac{1}{N} \sum_{i=1}^N \frac{\partial}{\partial \theta_j} (J_i(\vec{\theta}))$$
$$\nabla J(\vec{\theta}) = \begin{bmatrix} \vdots \\ j^{th} \\ \vdots \\ \end{bmatrix} = \frac{1}{N} \sum_{i=1}^N \nabla J_i(\vec{\theta})$$

# Stochastic gradient descent



1. Def: an epoch is a single pass through the training data
  - For GD, only one update per epoch
  - For SGD,  $N$  updates per epoch  $N = (\# \text{ train examples})$
2. SGD reduces MSE much more rapidly than GD
3. For GD / SGD, training MSE is initially large due to uninformed initialization

# Stochastic gradient descent

## Batch gradient descent

computes the gradient of the cost function to the parameters for the entire training dataset

$$\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta)$$

```
for i in range(nb_epochs):
    params_grad = evaluate_gradient(loss_function, data, params)
    params = params - learning_rate * params_grad
```

# Stochastic gradient descent

## Stochastic gradient descent

computes the gradient of the cost function to the parameters for each training example

$$\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta; x^{(i)}; y^{(i)})$$

```
for i in range(nb_epochs):
    np.random.shuffle(data)
    for example in data:
        params_grad = evaluate_gradient(loss_function, example, params)
        params = params - learning_rate * params_grad
```

# Stochastic gradient descent

## Mini-batch gradient descent

computes the gradient of the cost function to the parameters for every mini-batch of  $n$  training examples

$$\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta; x^{(i:i+n)}; y^{(i:i+n)})$$

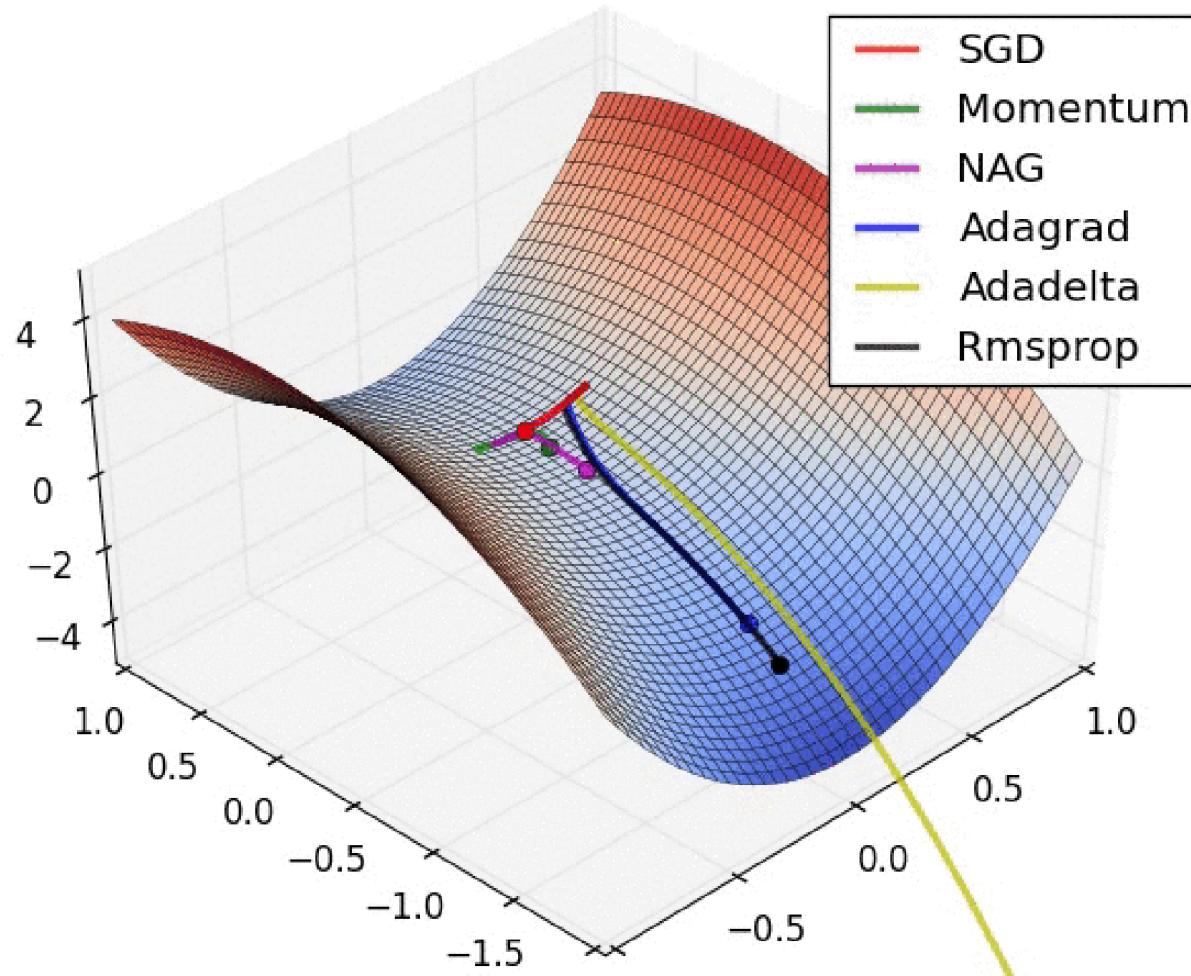
```
for i in range(nb_epochs):
    np.random.shuffle(data)
    for batch in get_batches(data, batch_size=50):
        params_grad = evaluate_gradient(loss_function, batch, params)
        params = params - learning_rate * params_grad
```

# Stochastic gradient descent



- Local optimum?
- How to choose a proper learning rate?
- Same learning rate for all parameter updates?

# Stochastic gradient descent



Local optimum is not a major concern, but there are many saddle points

# Stochastic gradient descent

---

## An overview of gradient descent optimization algorithms\*

---

**Sebastian Ruder**

Insight Centre for Data Analytics, NUI Galway  
Aylien Ltd., Dublin  
`ruder.sebastian@gmail.com`

### Abstract

Gradient descent optimization algorithms, while increasingly popular, are often used as black-box optimizers, as practical explanations of their strengths and weaknesses are hard to come by. This article aims to provide the reader with intuitions with regard to the behaviour of different algorithms that will allow her to put them to use. In the course of this overview, we look at different variants of gradient descent, summarize challenges, introduce the most common optimization algorithms, review architectures in a parallel and distributed setting, and investigate additional strategies for optimizing gradient descent.

# Stochastic gradient descent

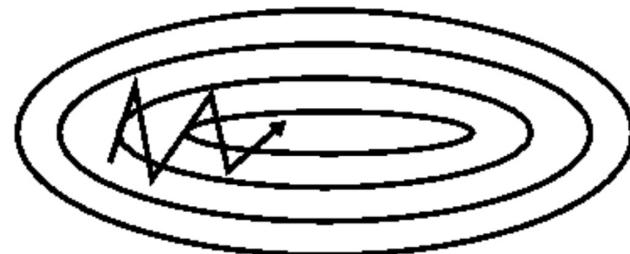
## Momentum

accelerate SGD in the relevant direction and dampen oscillations

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta)$$
$$\theta = \theta - v_t$$



(a) SGD without momentum



(b) SGD with momentum

# Stochastic gradient descent

## Nesterov accelerated gradient (NAG)

an approximation of the next position of the parameters

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta - \gamma v_{t-1})$$
$$\theta = \theta - v_t$$

look ahead by calculating the gradient not w.r.t. to our current parameters  
but w.r.t. the approximate future position of our parameters

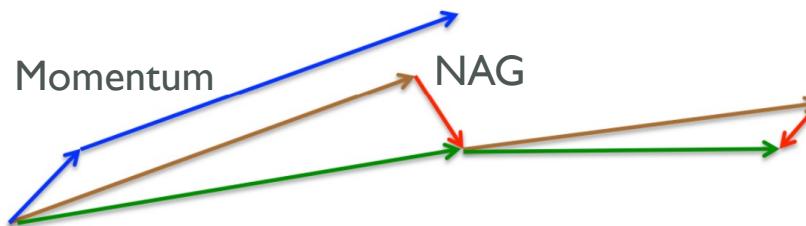


Figure 3: Nesterov update (Source: G. Hinton's lecture 6c)

# Stochastic gradient descent

## Adagrad

adapts the learning rate to the parameters, performing larger updates for infrequent and smaller updates for frequent parameters

$$g_{t,i} = \nabla_{\theta_t} J(\theta_{t,i})$$

element-wise     $\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,ii} + \epsilon}} \cdot g_{t,i}$



each diagonal element ( i, i ) is the sum of the squares of the gradients up to time step t

vectorize

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t$$

# Stochastic gradient descent

## Adaptive Moment Estimation (Adam)

computes adaptive learning rates for each parameter

first moment (the mean)  $m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$

second moment (the uncentered variance)  $v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$

bias-correction

$$\left\{ \begin{array}{l} \hat{m}_t = \frac{m_t}{1 - \beta_1^t} \\ \hat{v}_t = \frac{v_t}{1 - \beta_2^t} \\ \theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t \end{array} \right.$$

Adam might be the best overall choice

# Stochastic gradient descent

## Additional strategies for optimizing SGD

- Shuffling:  
shuffle the training data after every epoch to avoid bias
- Curriculum learning  
solve progressively harder problems
- Early stopping  
monitor error on a validation set during training and stop (with some patience) if your validation error does not improve enough
- Gradient noise:  
adding this noise makes networks more robust to poor initialization and helps training particularly deep and complex networks

# Training neural networks

Stochastic gradient descent

**Dropout**

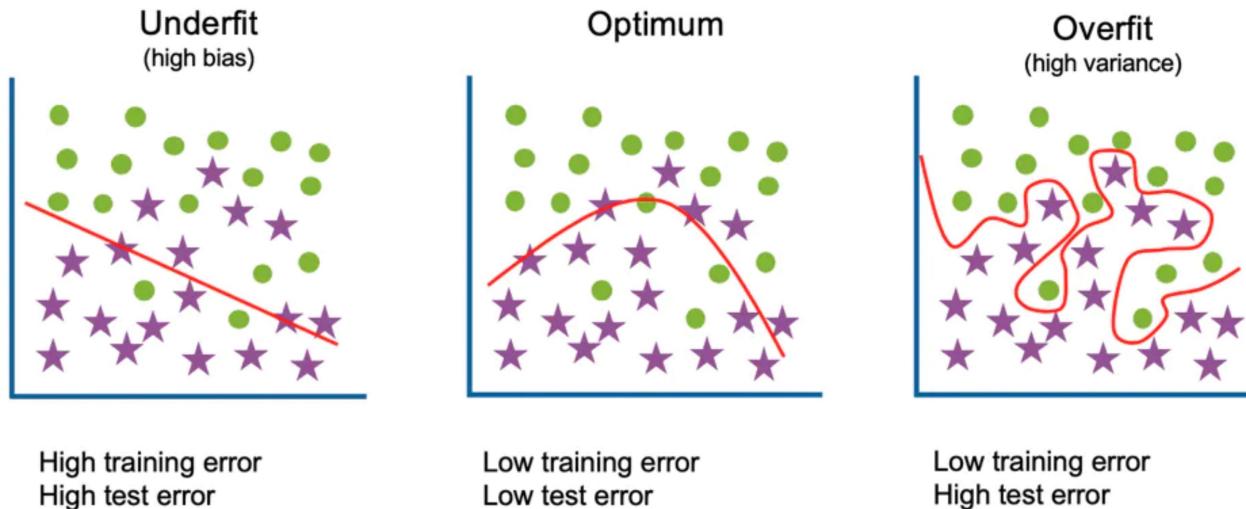
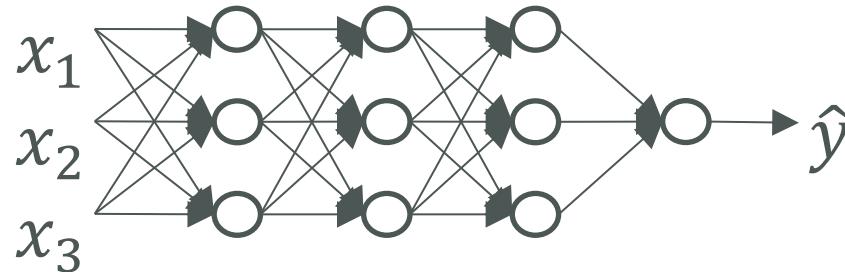
Initialization

Hyperparameter tuning

Batch normalization

# Dropout regularization

The term “dropout” refers to dropping out units in a neural network



# Dropout regularization

## Assumption:

A fully connected layer occupies most of the parameters, and hence, neurons develop co-dependency amongst each other during training which curbs the individual power of each neuron leading to over-fitting of training data.

**Limit the power of the model parameters**

# Dropout regularization

## Regularization

**Given** objective function:  $J(\theta)$

**Goal** is to find:  $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} J(\theta) + \lambda r(\theta)$

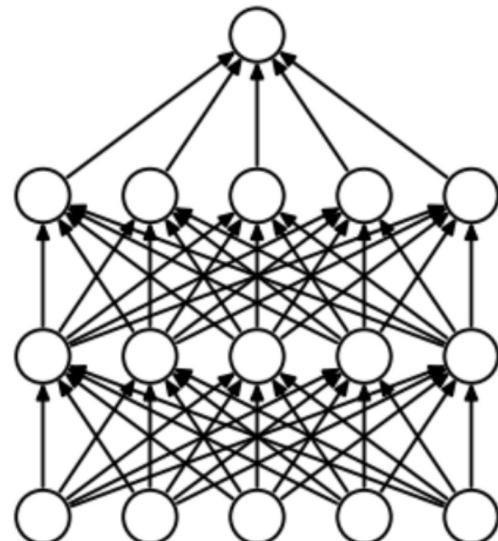
tradeoff between fitting the data and keeping the model simple  
*(push some parameters smaller)*

$q$	$r(\theta)$	yields parameters that are...	name	optimization notes
0	$  \theta  _0 = \sum \mathbb{1}(\theta_m \neq 0)$	zero values	Lo reg.	no good computational solutions
1	$  \theta  _1 = \sum  \theta_m $	zero values	L1 reg.	subdifferentiable
2	$(  \theta  _2)^2 = \sum \theta_m^2$	small values	L2 reg.	differentiable

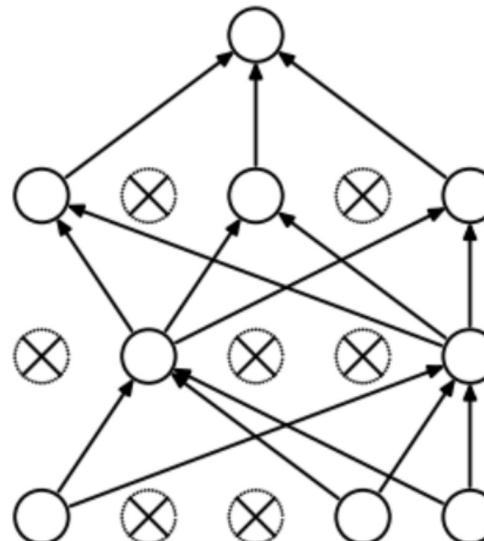
# Dropout regularization

The term “dropout” refers to dropping out units in a neural network

Intuition: cannot rely on any one feature, so have to spread out weights



(a) Standard Neural Net



(b) After applying dropout.

# Dropout regularization

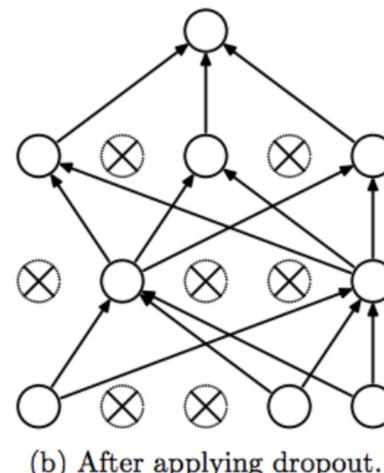
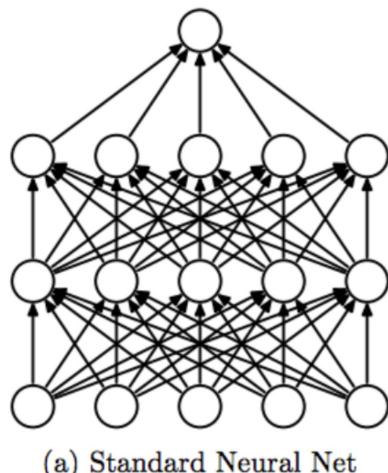
The term “dropout” refers to dropping out units in a neural network

## Training Phase:

For each hidden layer, for each training sample, for each iteration, ignore (zero out) a random fraction,  $p$ , of nodes (and corresponding activations).

## Testing Phase:

Use all activations, but reduce them by a factor  $p$  (to account for the missing activations during training).



# Dropout regularization

The term “dropout” refers to dropping out units in a neural network

- Dropout forces a neural network to learn more robust features that are useful in conjunction with many different random subsets of the other neurons.
- Dropout roughly doubles the number of iterations required to converge. However, training time for each epoch is less.
- With  $H$  hidden units, each of which can be dropped, we have  $2^H$  possible models. In testing phase, the entire network is considered and each activation is reduced by a factor  $p$ .

**Apply only when you face overfitting issue (lack data)**

# Training neural networks

Stochastic gradient descent

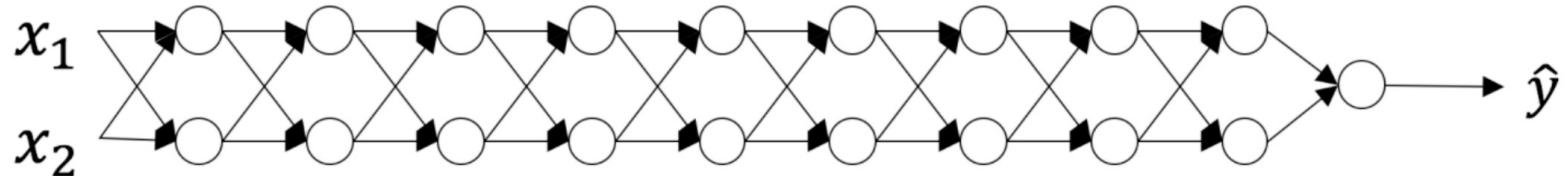
Dropout

**Initialization**

Hyperparameter tuning

Batch normalization

# Weight initialization



$$\hat{y} = a^{[L]} = W^{[L]}W^{[L-1]}W^{[L-2]} \dots W^{[3]}W^{[2]}W^{[1]}x$$

**Case 1:** A too-large initialization leads to exploding gradients

$$W^{[1]} = W^{[2]} = \dots = W^{[L-1]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$$

**Case 2:** A too-small initialization leads to vanishing gradients

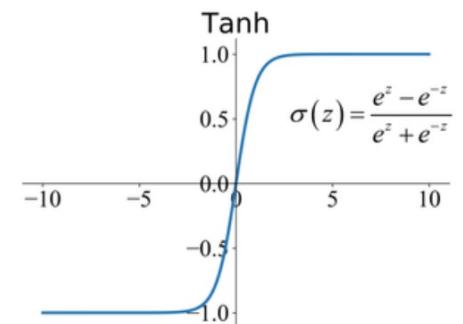
$$W^{[1]} = W^{[2]} = \dots = W^{[L-1]} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$$

# Weight initialization

Key assumptions:

1. Mean of the activations should be zero
2. Variance of the activations should stay the same across layers

$$z^{[l]} = W^{[l]} a^{[l-1]}$$
$$a^{[l]} = \tanh(z^{[l]})$$



Find right way to initialize  $W$ , such that :

$$\text{Var}(a^{[l-1]}) = \text{Var}(a^{[l]})$$

# Weight initialization

$$z^{[l]} = W^{[l]}a^{[l-1]} + b^{[l]}$$

$$a^{[l]} = \tanh(z^{[l]})$$

Additional assumptions:

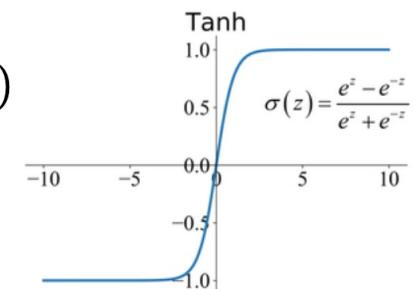
1. Early on in the training, we are in the *linear regime* of  $\tanh$   $\tanh(z^{[l]}) \approx z^{[l]}$
2. Weights are independent and identically distributed

$$\text{Var}(w_{kj}^{[l]}) = \text{Var}(w_{11}^{[l]}) = \text{Var}(w_{12}^{[l]}) = \dots = \text{Var}(W^{[l]})$$

3. Inputs are independent and identically distributed

$$\text{Var}(a_j^{[l-1]}) = \text{Var}(a_1^{[l-1]}) = \text{Var}(a_2^{[l-1]}) = \dots = \text{Var}(a^{[l-1]})$$

4. Weights and inputs are mutually independent



**Vectorize:**  $\text{Var}(a^{[l]}) = \text{Var}(z^{[l]})$

**Element-wise:**  $\text{Var}(a_k^{[l]}) = \text{Var}(z_k^{[l]}) = \text{Var}\left(\sum_{j=1}^{n^{[l-1]}} w_{kj}^{[l]} a_j^{[l-1]}\right) = \sum_{j=1}^{n^{[l-1]}} \text{Var}(w_{kj}^{[l]} a_j^{[l-1]})$

$$\text{Var}(XY) = E[X]^2 \text{Var}(Y) + \text{Var}(X)E[Y]^2 + \text{Var}(X)\text{Var}(Y)$$

$$\text{Var}(w_{kj}^{[l]} a_j^{[l-1]}) = E[\mathbf{0}_{kj}^{[l]}]^2 \text{Var}(a_j^{[l-1]}) + \text{Var}(w_{kj}^{[l]})E[\mathbf{0}_j^{[l-1]}]^2 + \text{Var}(w_{kj}^{[l]})\text{Var}(a_j^{[l-1]})$$

$$\text{Var}(z_k^{[l]}) = \sum_{j=1}^{n^{[l-1]}} \text{Var}(w_{kj}^{[l]}) \text{Var}(a_j^{[l-1]}) = \sum_{j=1}^{n^{[l-1]}} \text{Var}(W^{[l]}) \text{Var}(a^{[l-1]}) = n^{[l-1]} \text{Var}(W^{[l]}) \text{Var}(a^{[l-1]})$$

# Weight initialization

From previous slide, we end up with

$$Var(a^{[l]}) = n^{[l-1]} Var(W^{[l]}) Var(a^{[l-1]})$$

$$n^{[l-1]} Var(W^{[l]}) \begin{cases} < 1 & \Rightarrow \text{Vanishing Signal} \\ = 1 & \Rightarrow Var(a^{[L]}) = Var(x) \\ > 1 & \Rightarrow \text{Exploding Signal} \end{cases}$$

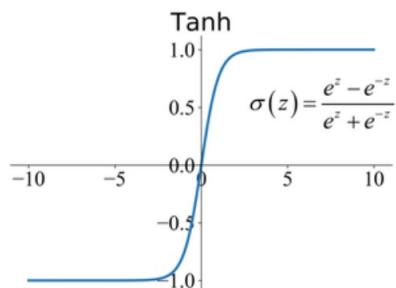
To make  $Var(a^{[l-1]}) = Var(a^{[l]})$ , we need to set

$$Var(W^{[l]}) = \frac{1}{n^{[l-1]}} \quad \mathcal{N}(0, \frac{1}{n^{[l-1]}})$$

Xavier initialization

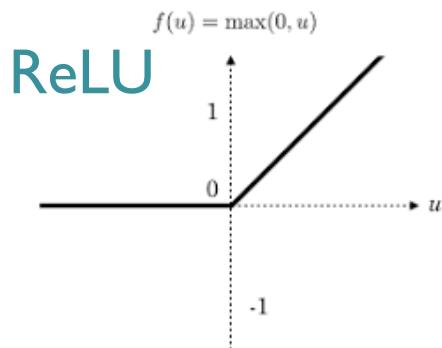
# Weight initialization

## Xavier initialization



$$\mathcal{N}(0, \frac{1}{n^{[l-1]}})$$

`W = tf.Variable(np.random.randn(node_in, node_out)) / np.sqrt(node_in)`

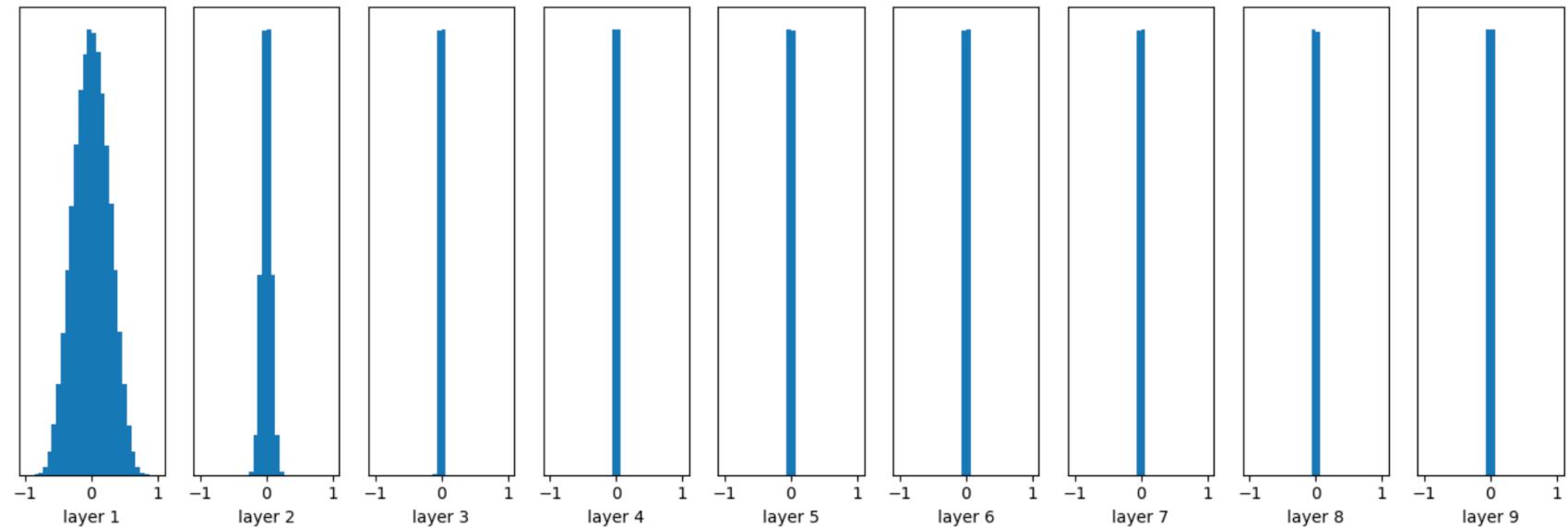


Only half of nodes will be activated, so use  $n^{[l-1]}/2$

`W = tf.Variable(np.random.randn(node_in, node_out)) / np.sqrt(node_in/2)`

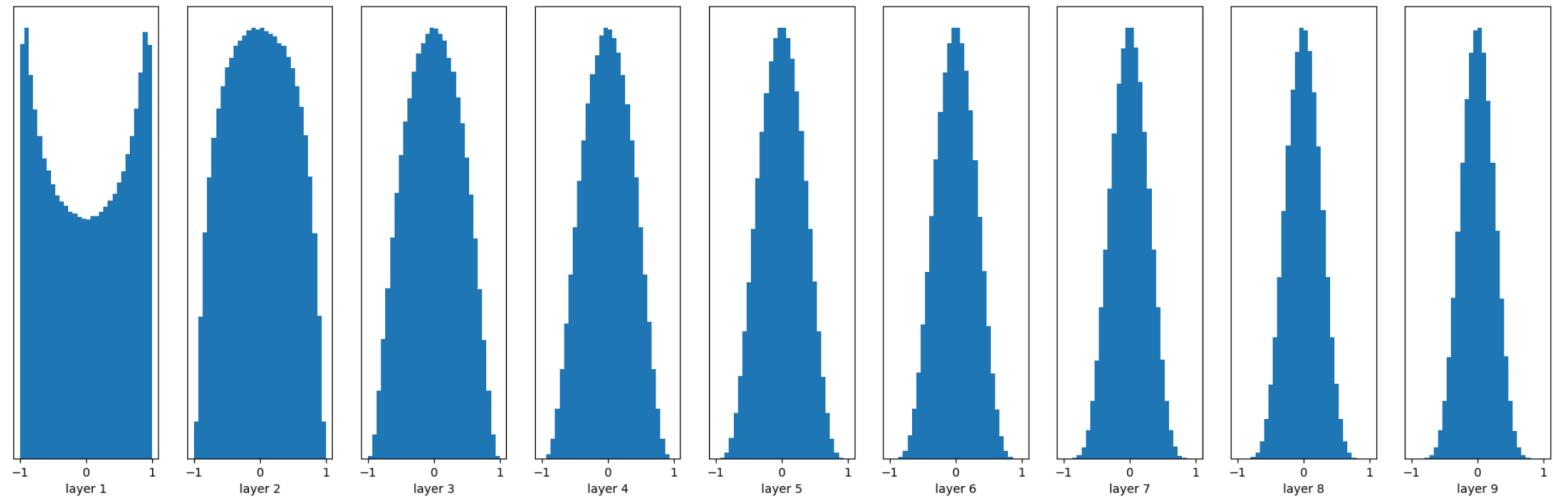
# Weight initialization

```
W = tf.Variable(np.random.randn(node_in, node_out)) * 0.01
```



# Weight initialization

```
W = tf.Variable(np.random.randn(node_in, node_out)) / np.sqrt(node_in)
```

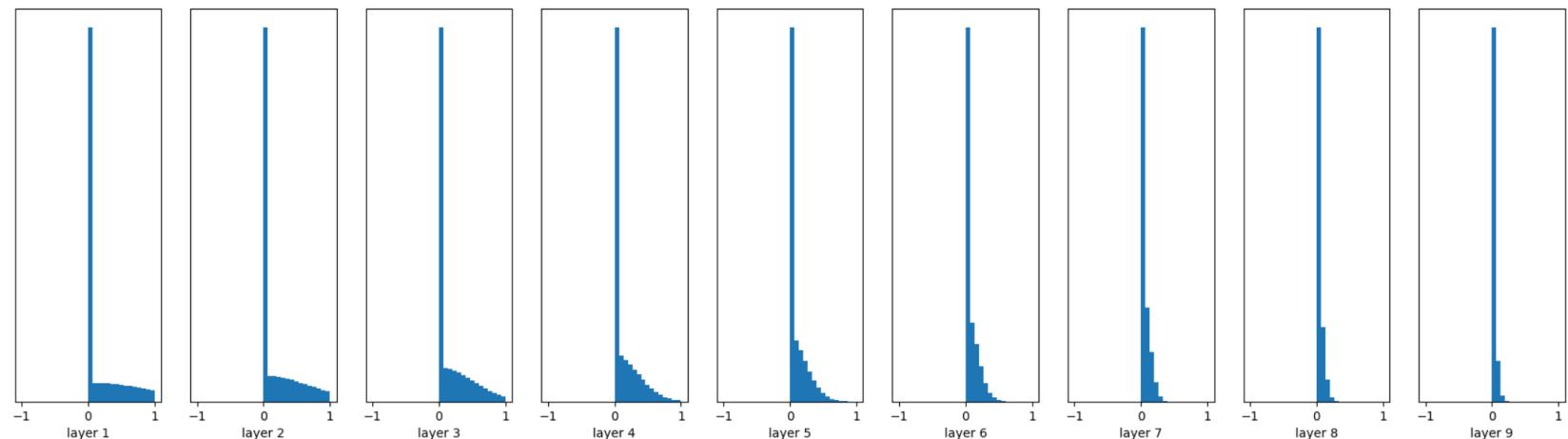


# Weight initialization

```
W = tf.Variable(np.random.randn(node_in, node_out)) / np.sqrt(node_in)
```

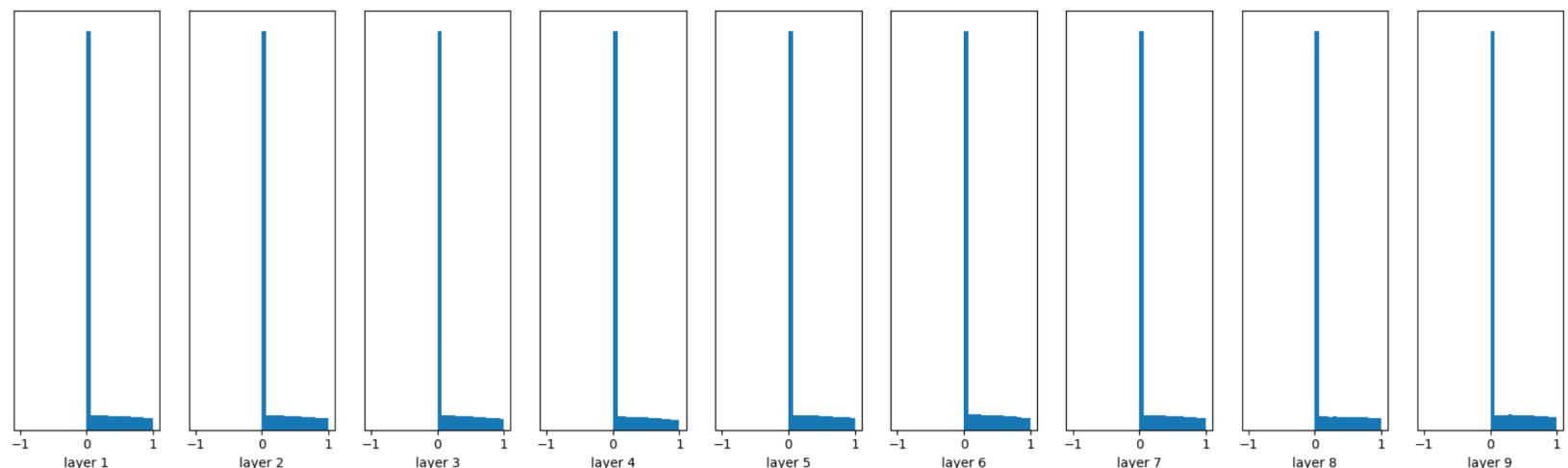
.....

```
fc = tf.nn.relu(fc)
```



# Weight initialization

```
W = tf.Variable(np.random.randn(node_in, node_out)) / np.sqrt(node_in / 2)  
.....  
fc = tf.nn.relu(fc)
```



# Training neural networks

Stochastic gradient descent

Dropout

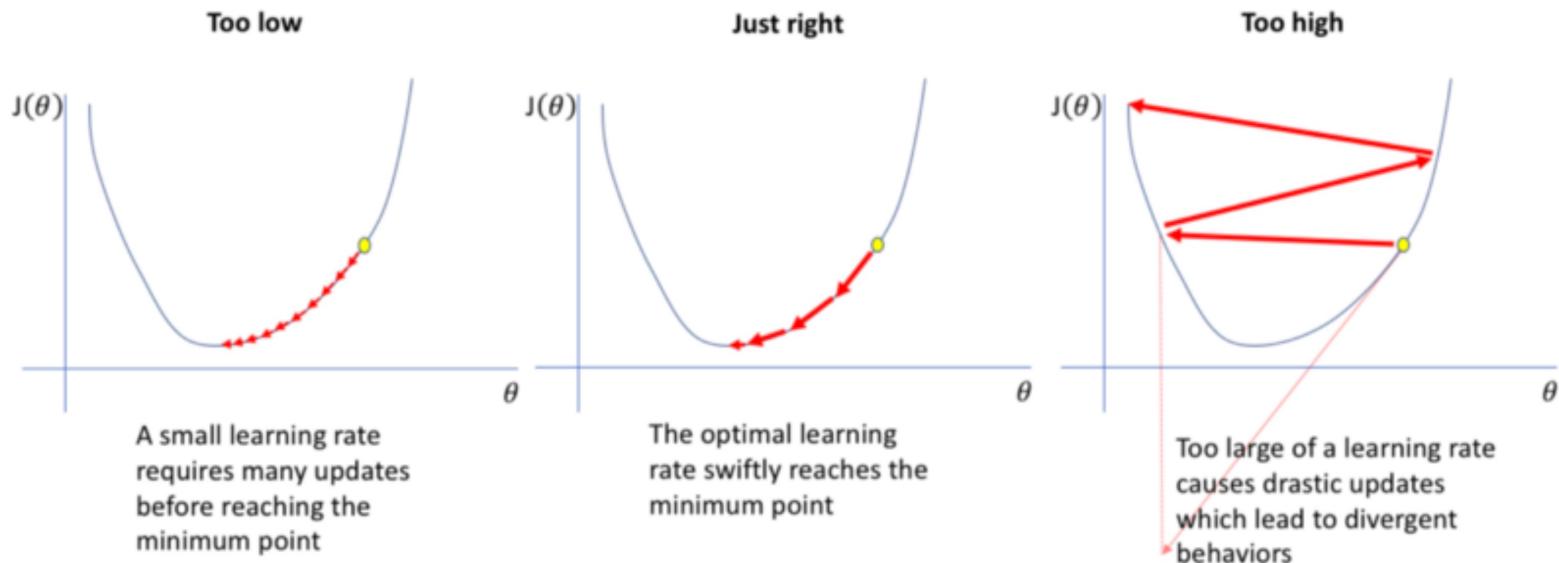
Initialization

**Hyperparameter tuning**

Batch normalization

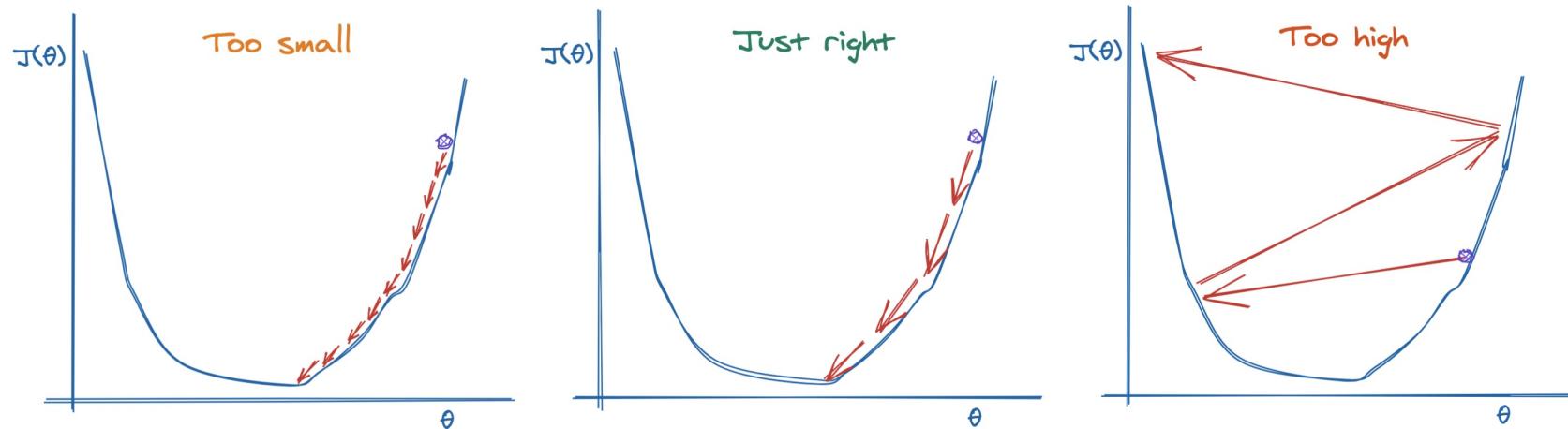
# Hyperparameter tuning

Learning rate, dropout probability, number of layers, neurons per layer...



# Hyperparameter tuning

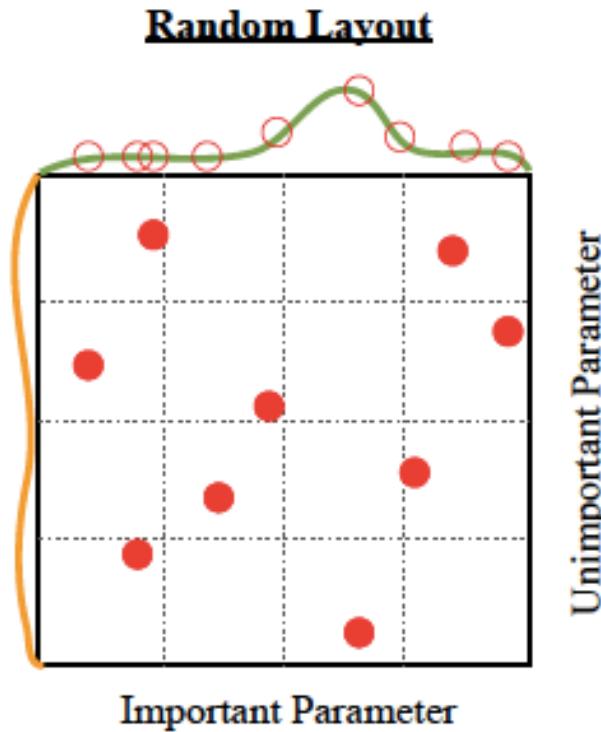
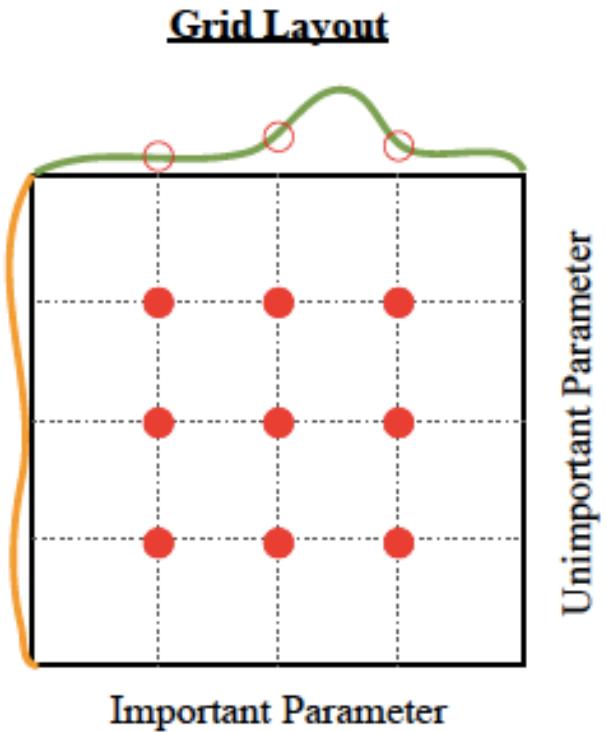
Learning rate, dropout probability, number of layers, neurons per layer...



Neural architecture search: too expensive

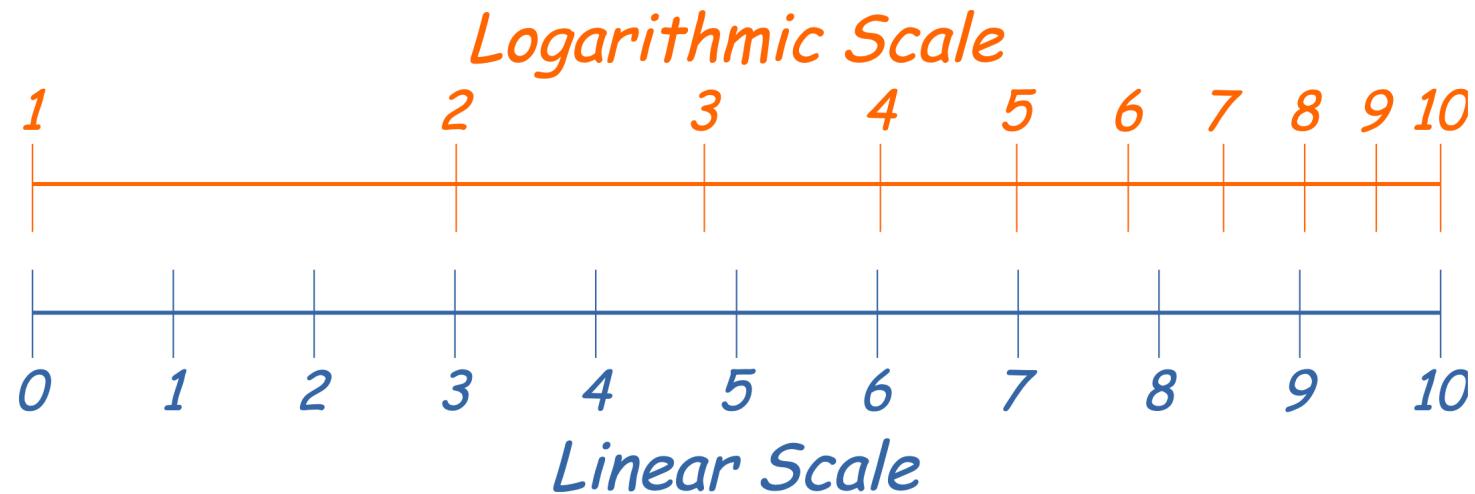
# Hyperparameter tuning

Trick I: Random search instead of grid search



# Hyperparameter tuning

Trick 2: log-scale search



# Training neural networks

Stochastic gradient descent

Dropout

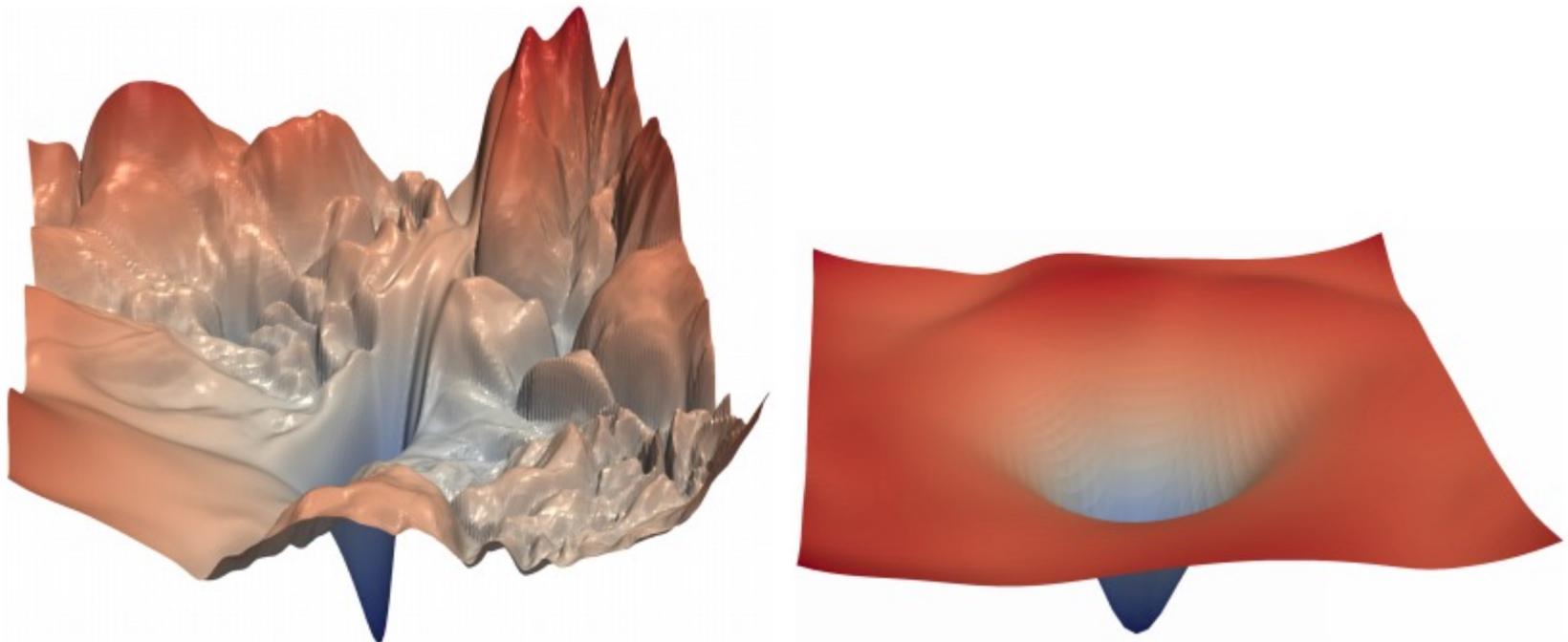
Initialization

Hyperparameter tuning

**Batch normalization**

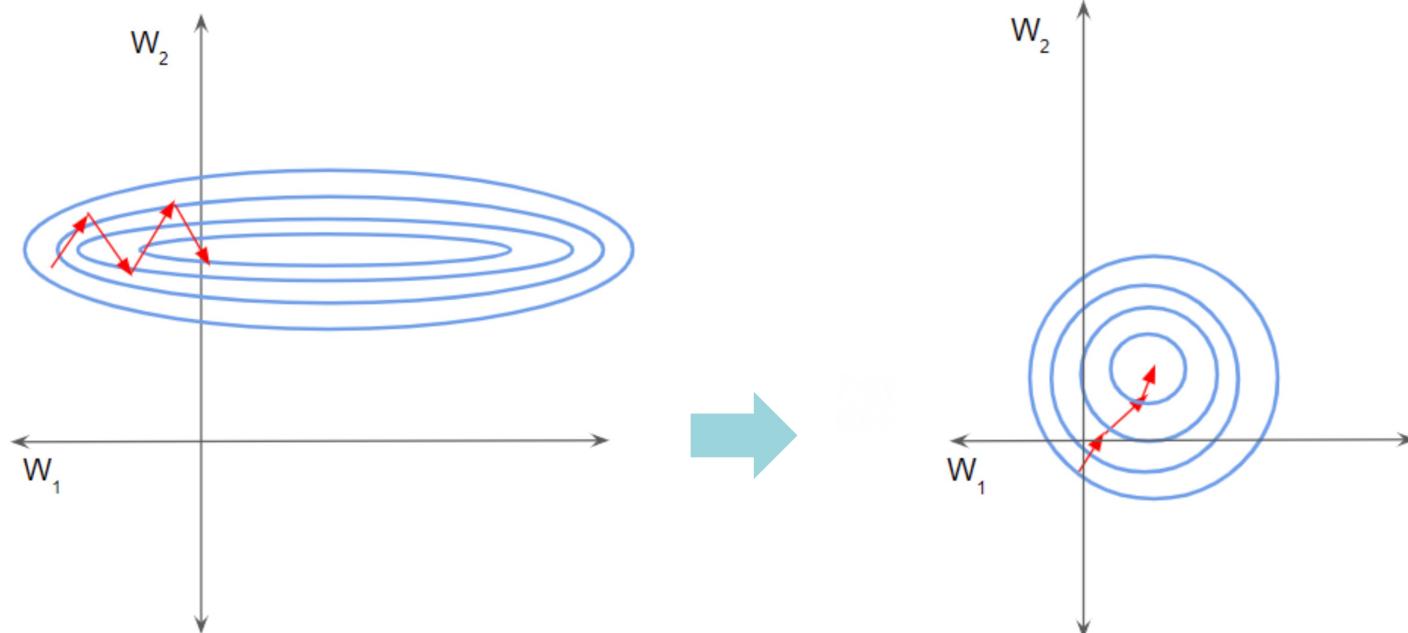
# Batch normalization

Make the loss landscape nicer



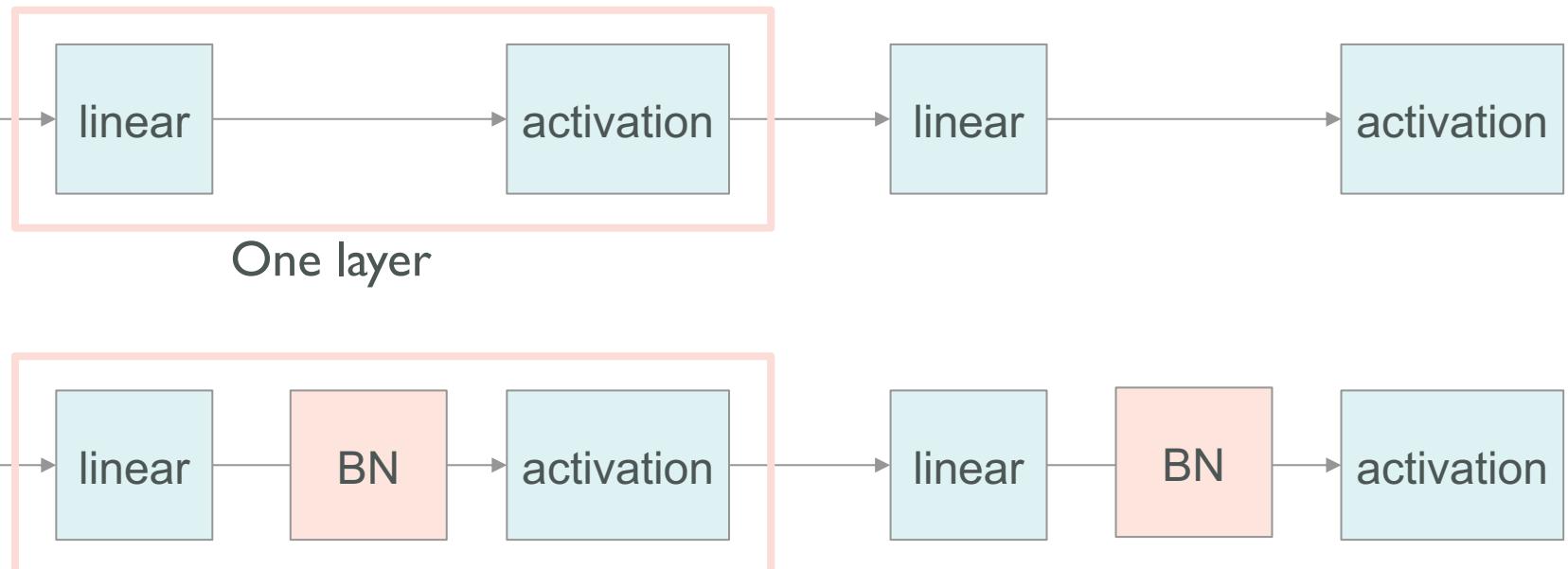
# Batch normalization

Make the loss landscape nicer



# Batch normalization

Batch normalization is a method used to make neural networks faster and more stable through normalization of the layers' inputs by re-centering and re-scaling.



# Batch normalization

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

**Algorithm 1:** Batch Normalizing Transform, applied to activation  $x$  over a mini-batch.

# Training neural networks

Stochastic gradient descent: Use Adam

Dropout: Use when overfitting

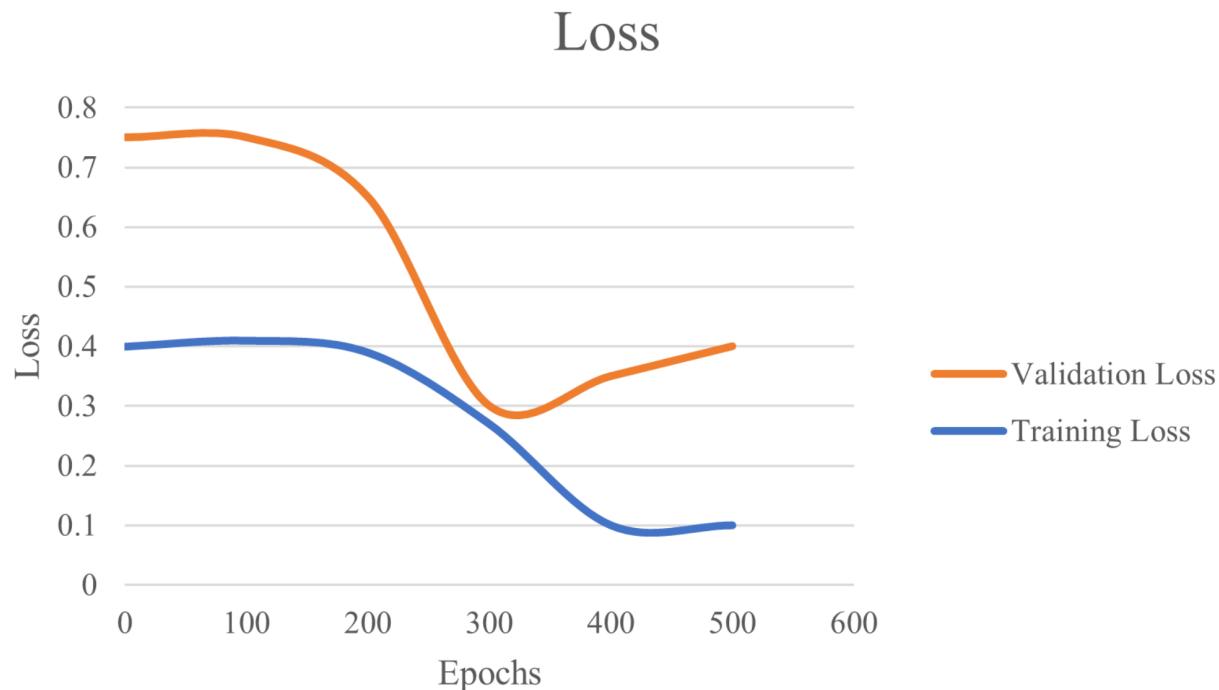
Initialization: Use Xavier initialization

Hyperparameter tuning: Random search

Batch normalization: Use all the time

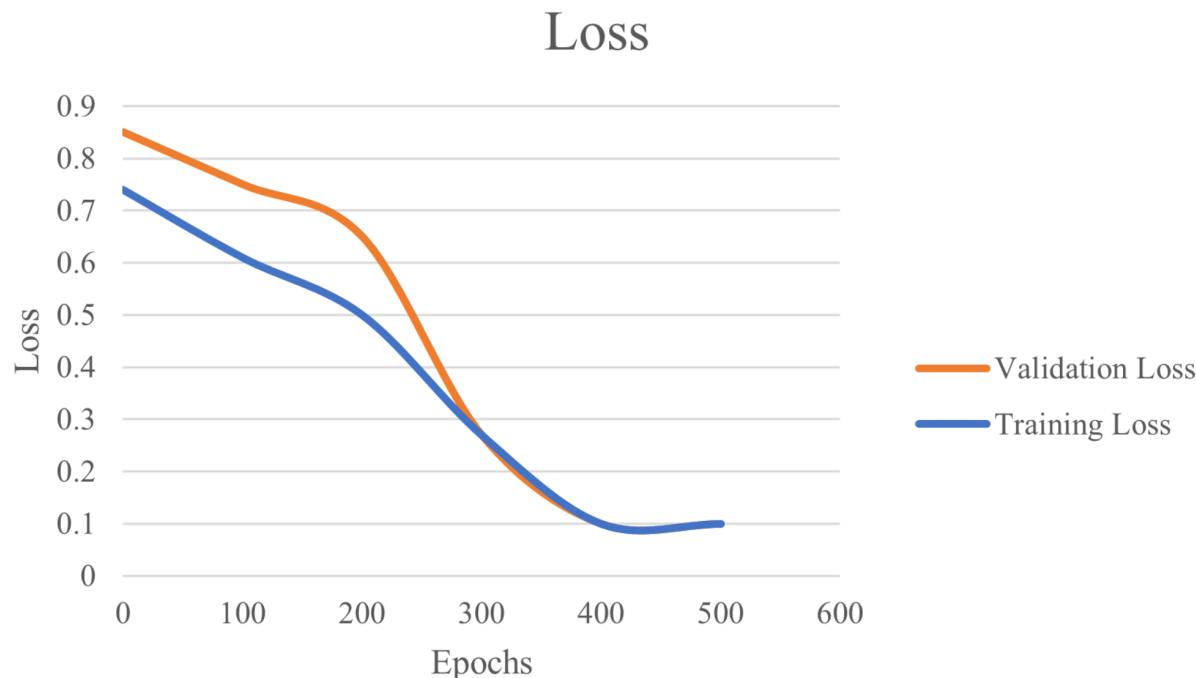
# Training neural networks

## I. Train on a small dataset (overfitting: loss and visualization)



# Training neural networks

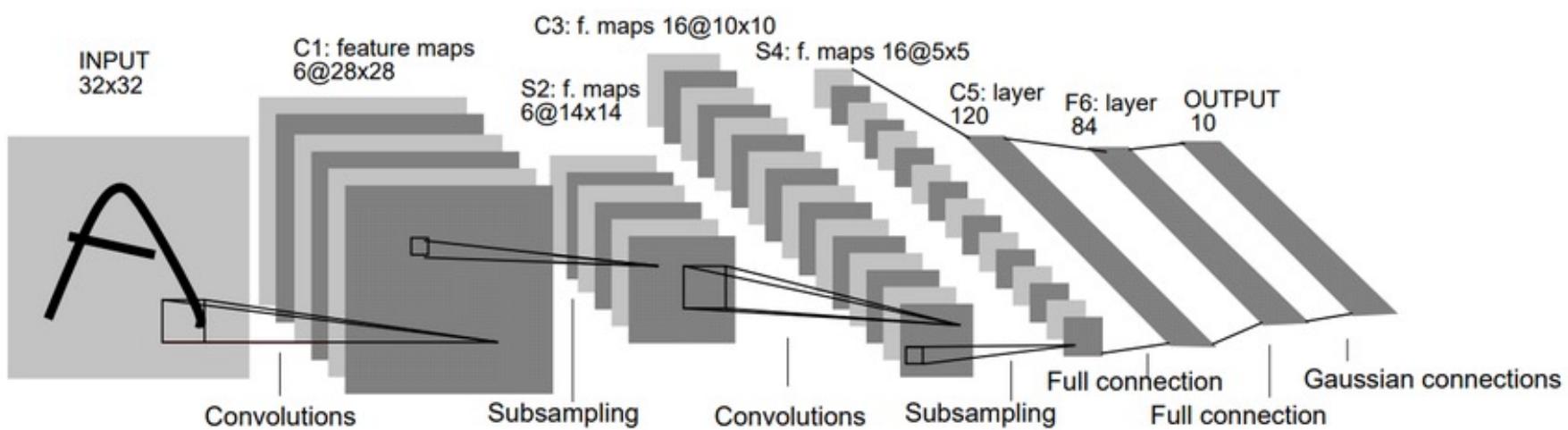
- I. Train on a small dataset (overfitting: loss and visualization)
2. Train on a large dataset (monitor your training/validation loss)



# Training neural networks

1. Train on a small dataset (overfitting: loss and visualization)
2. Train on a large dataset (monitor your training/validation loss)
3. Evaluation metric and test

# Next lecture: CNNs



Thank you very much!

sihengc@sjtu.edu.cn