

[Live] 실전! TensorFlow로 배우는 딥러닝

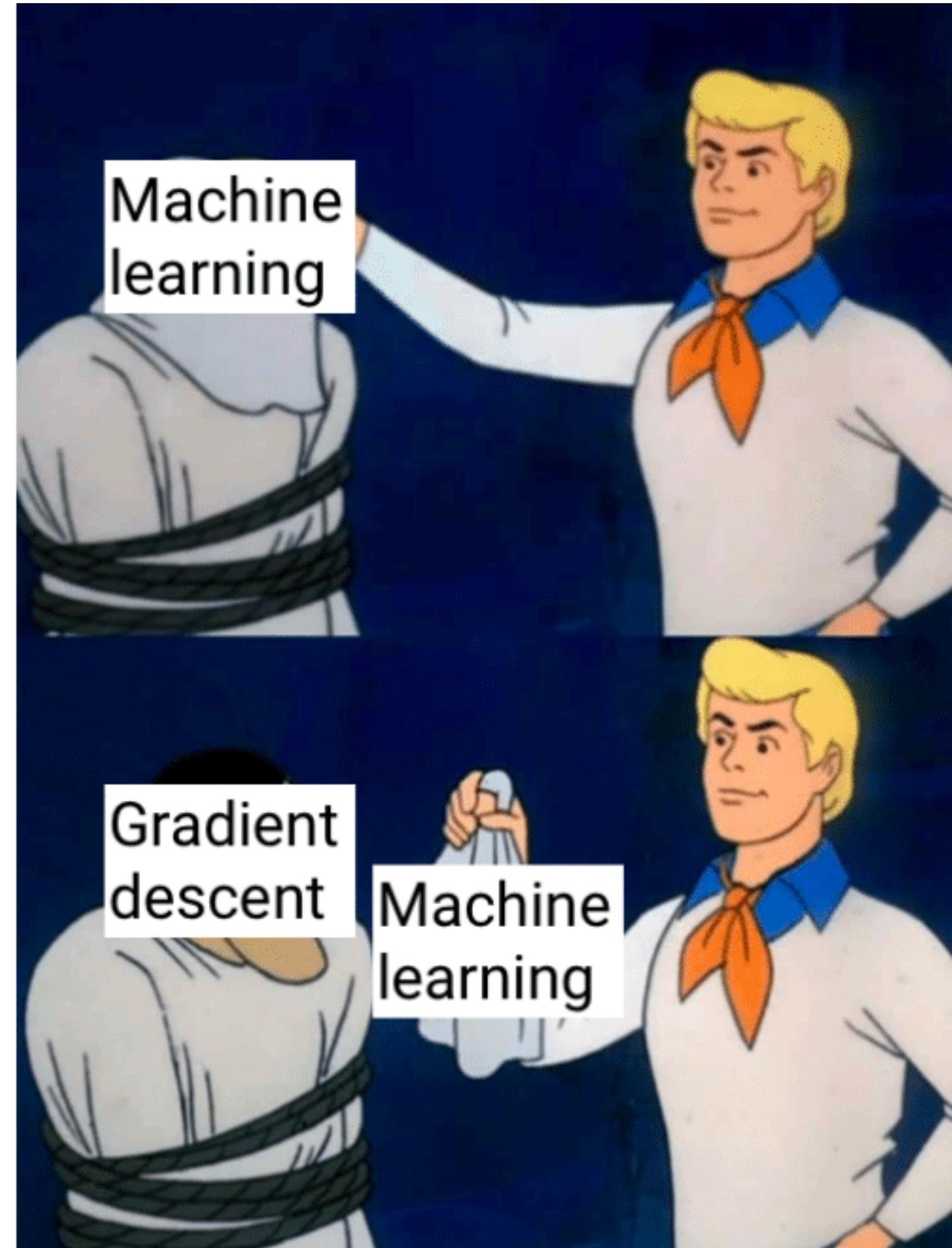
Day2

- Optimization & Regularization
 - CNN

목차

1. Optimization
 - ✓ Optimization
 - ✓ Gradient Descent
 - ✓ SGD + Momentum
 - ✓ 다양한 최적화 함수들
2. Generalization
 - ✓ Initialization
 - ✓ Over-fitting
 - ✓ Early stopping
 - ✓ Regularization
 - ✓ Dropout
 - ✓ Data Augmentation
 - ✓ Batch Normalization
3. Tensorflow를 이용한 DNN (Fashion MNIST Classification)
4. CNN
5. Tensorflow를 이용한 CNN (Fashion MNIST Classification)

1. Optimization optimization

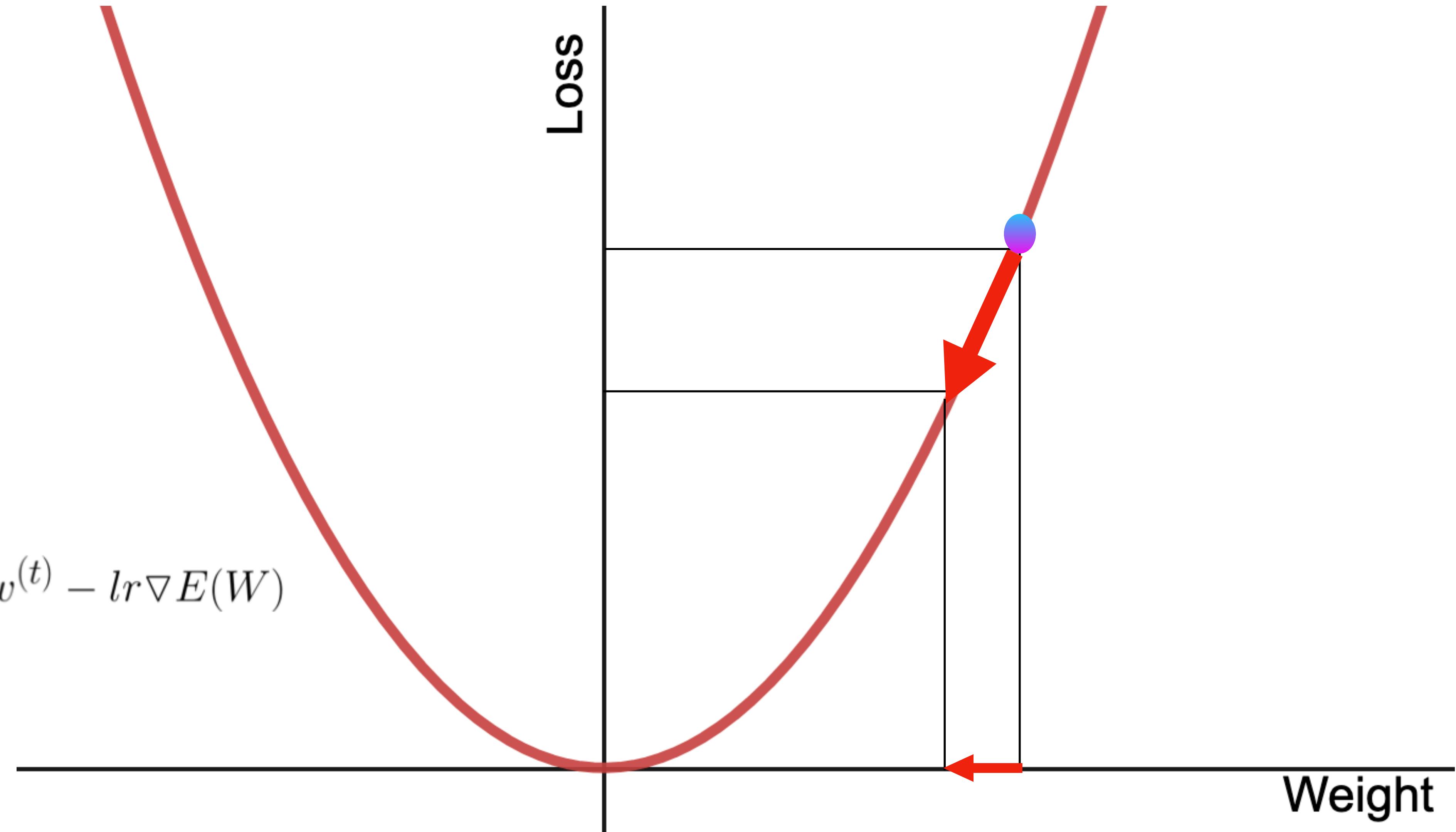


Machine learning behind the scenes

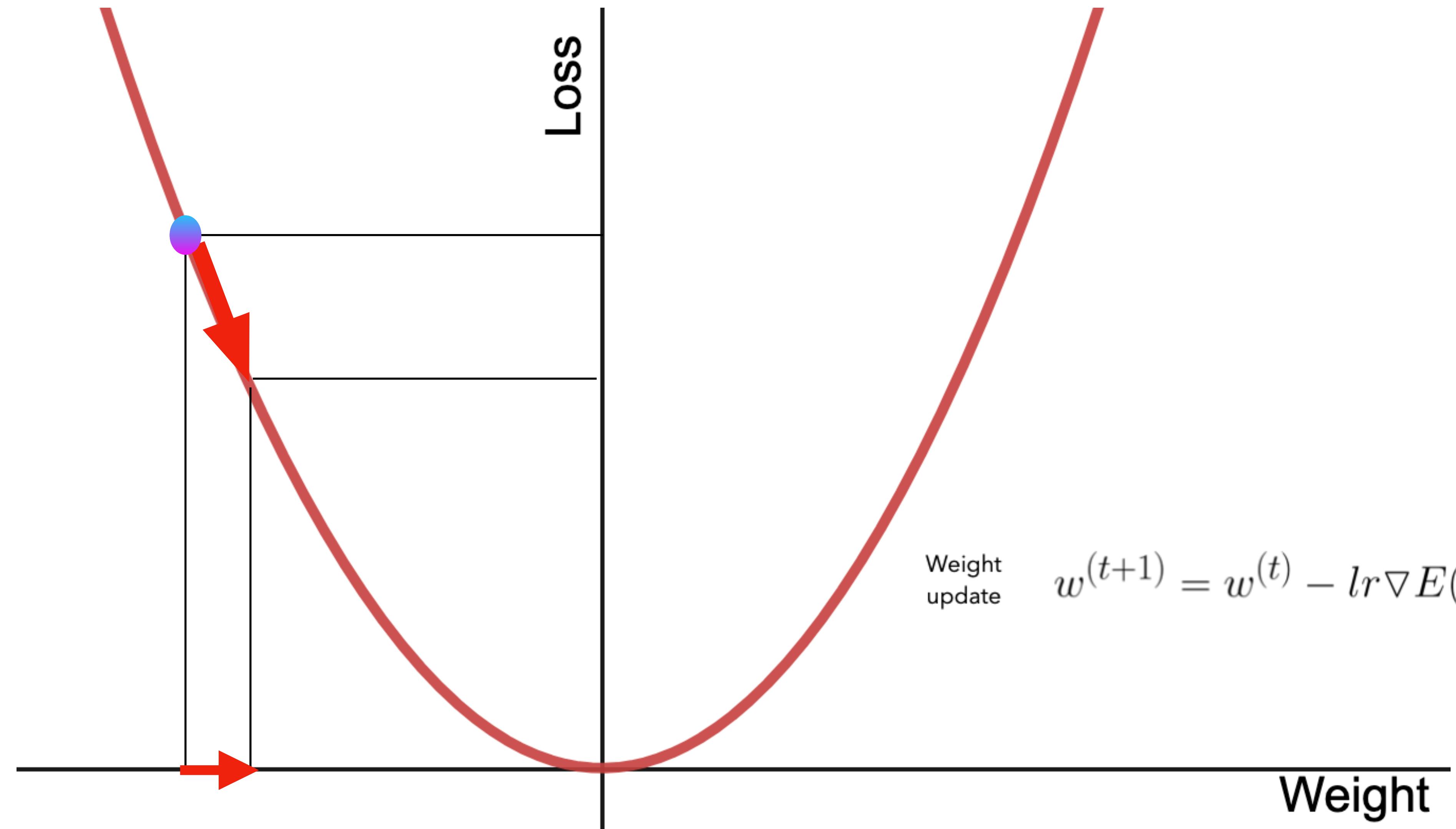
1. Optimization optimization

Weight update

$$w^{(t+1)} = w^{(t)} - lr \nabla E(W)$$



1. Optimization optimization



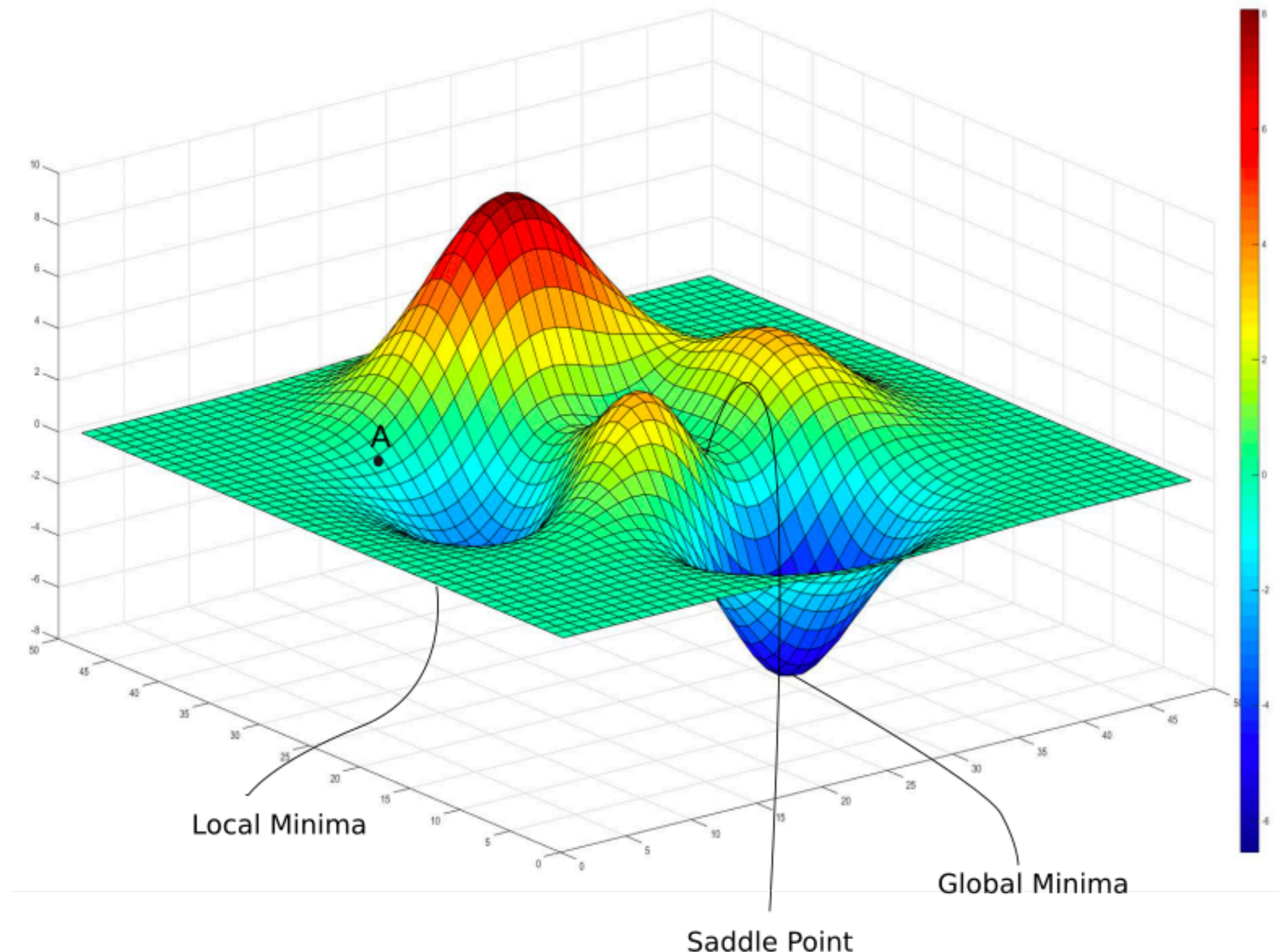
1. Optimization optimization

그런데,

정말 Loss function은 weight space에서
이쁜 빗살무늬토기 모양일까??



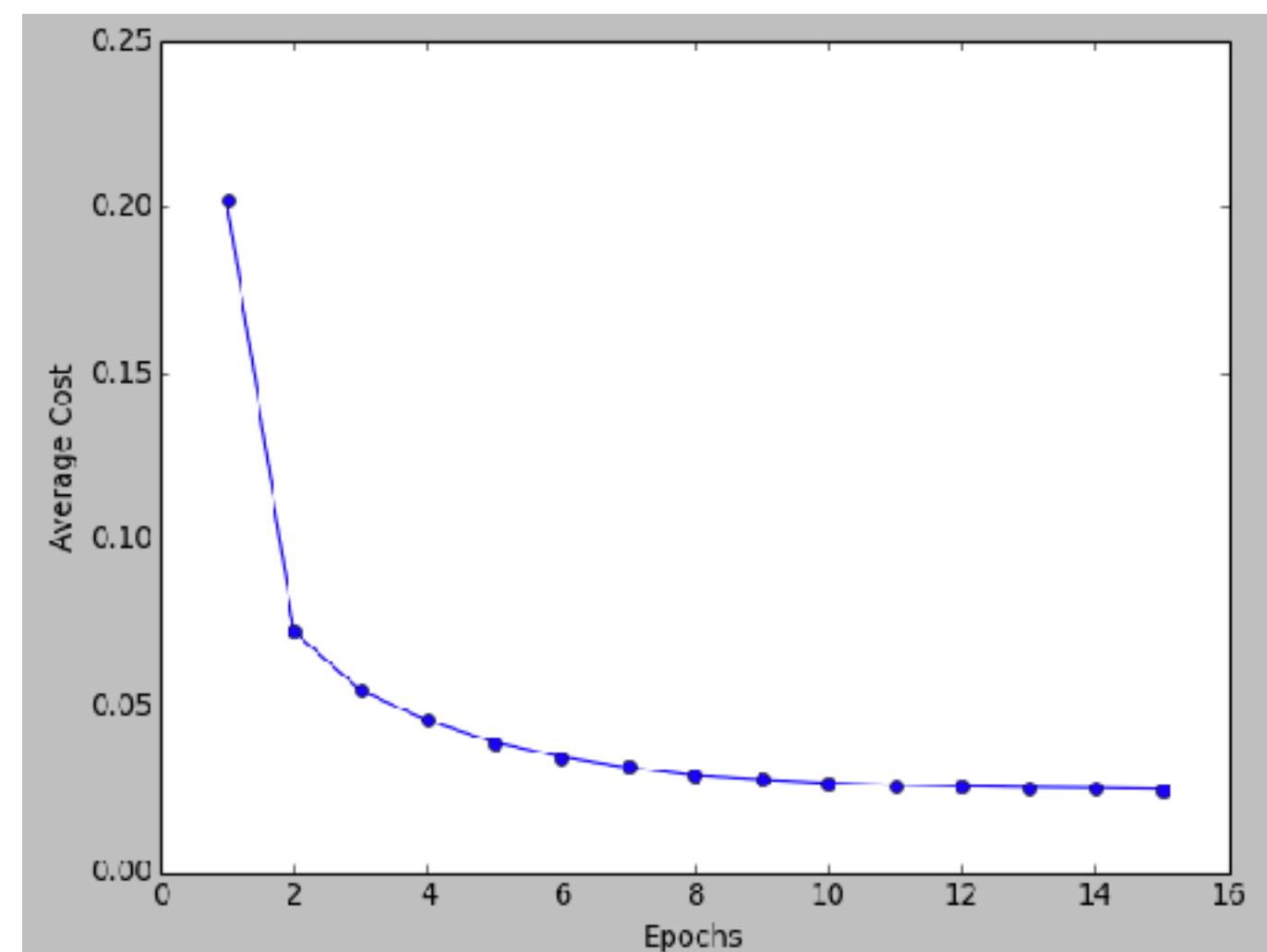
1. Optimization optimization



1. Optimization gradient descent

Batch Gradient descent

$$E(W) = \frac{1}{2} \sum_{n=1}^N (t - \hat{y}_n)^2$$

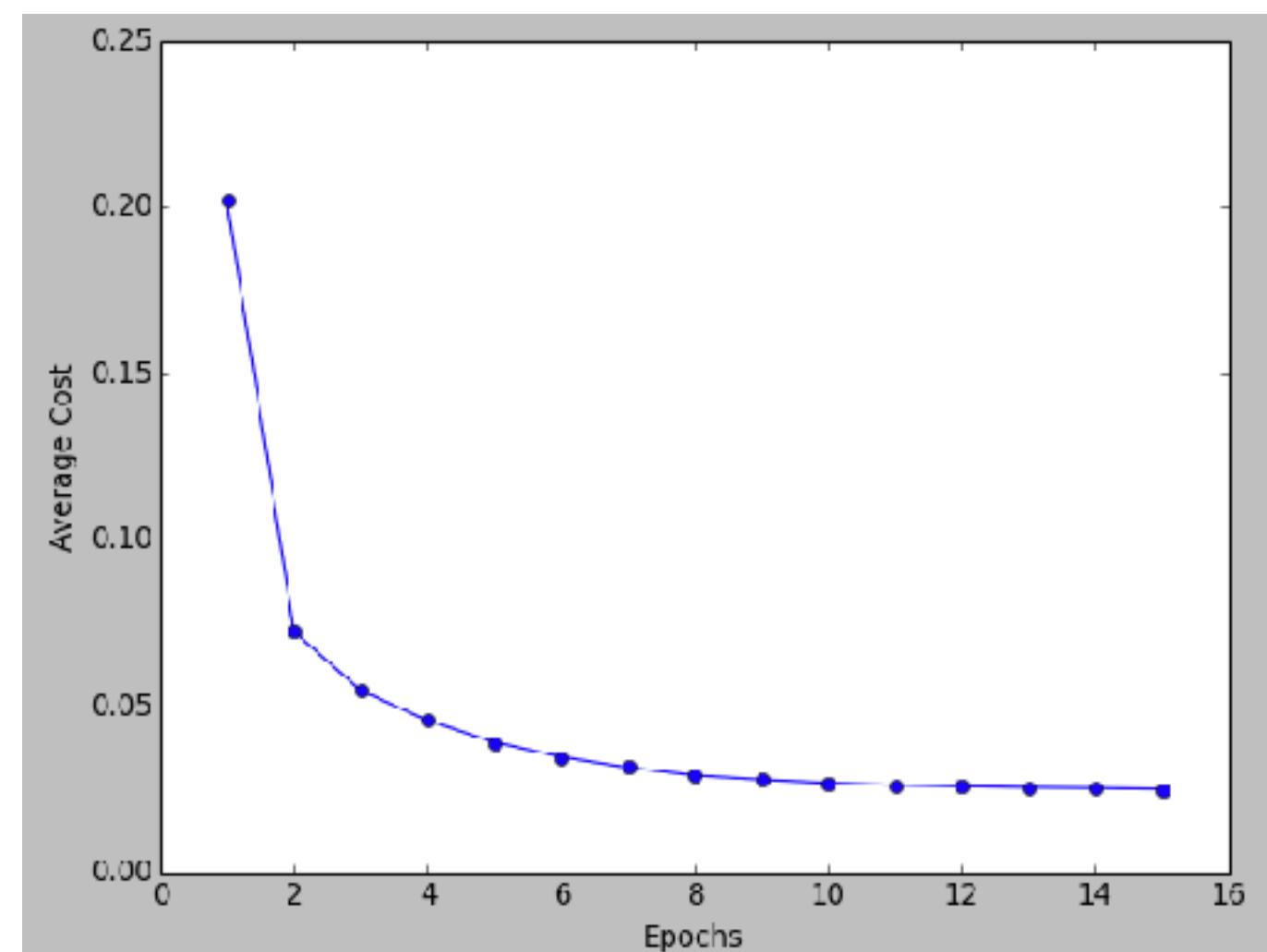


1. 모든 data에 대해서 loss를 다 구한다.
2. 개별 input에 대하여 구한 loss를 더하여 데이터 개수로 나눠준다.(평균을 구한다.)
3. Weight에 대한 gradient를 구한다.
4. backpropagation으로 한 번 weight가 update된다.

1. Optimization gradient descent

Batch Gradient descent

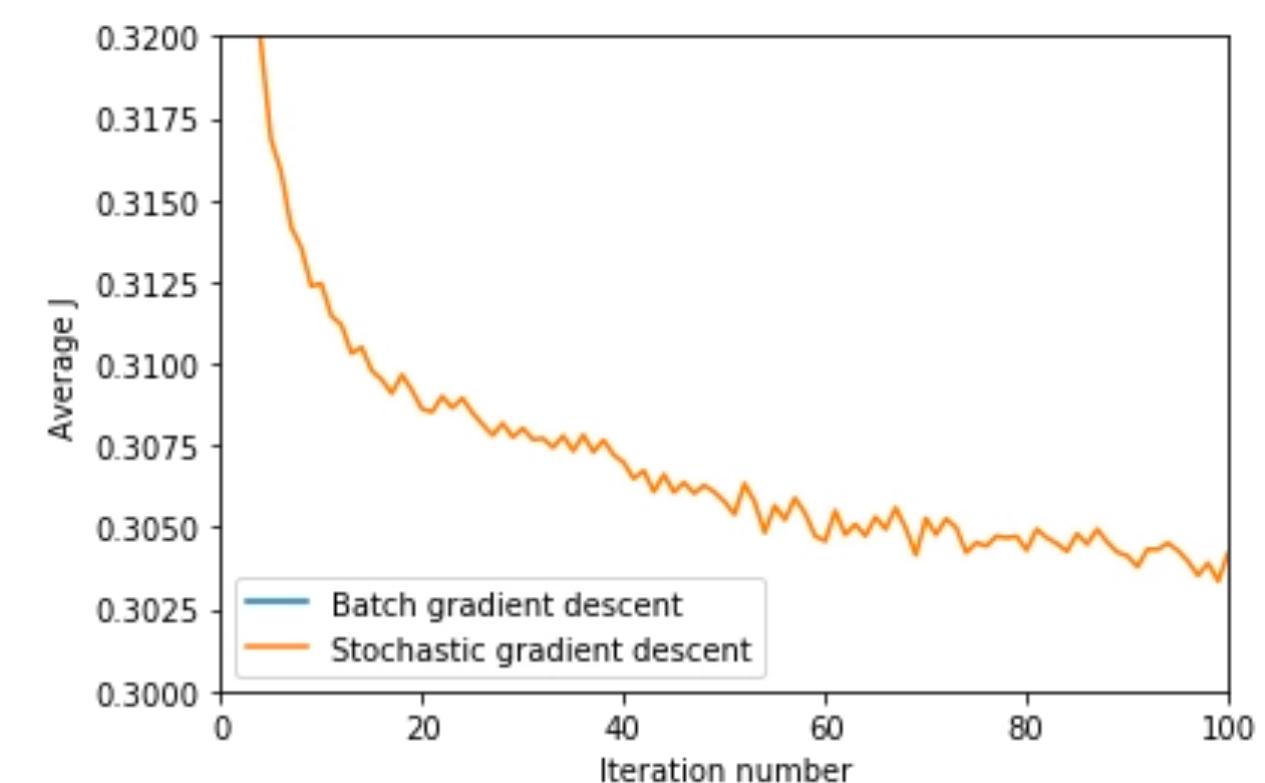
$$E(W) = \frac{1}{2} \sum_{n=1}^N (t - \hat{y}_n)^2$$



1. 모드에서 업데이트를 다룬다.
만약, 데이터가 50만개라면?
구한다.)
2. 개별 input
3. 50만번의 계산을 해야 가중치가 한 번 update된다.
4. backpropagation으로 한 번 weight가 update된다.

Stochastic Gradient descent

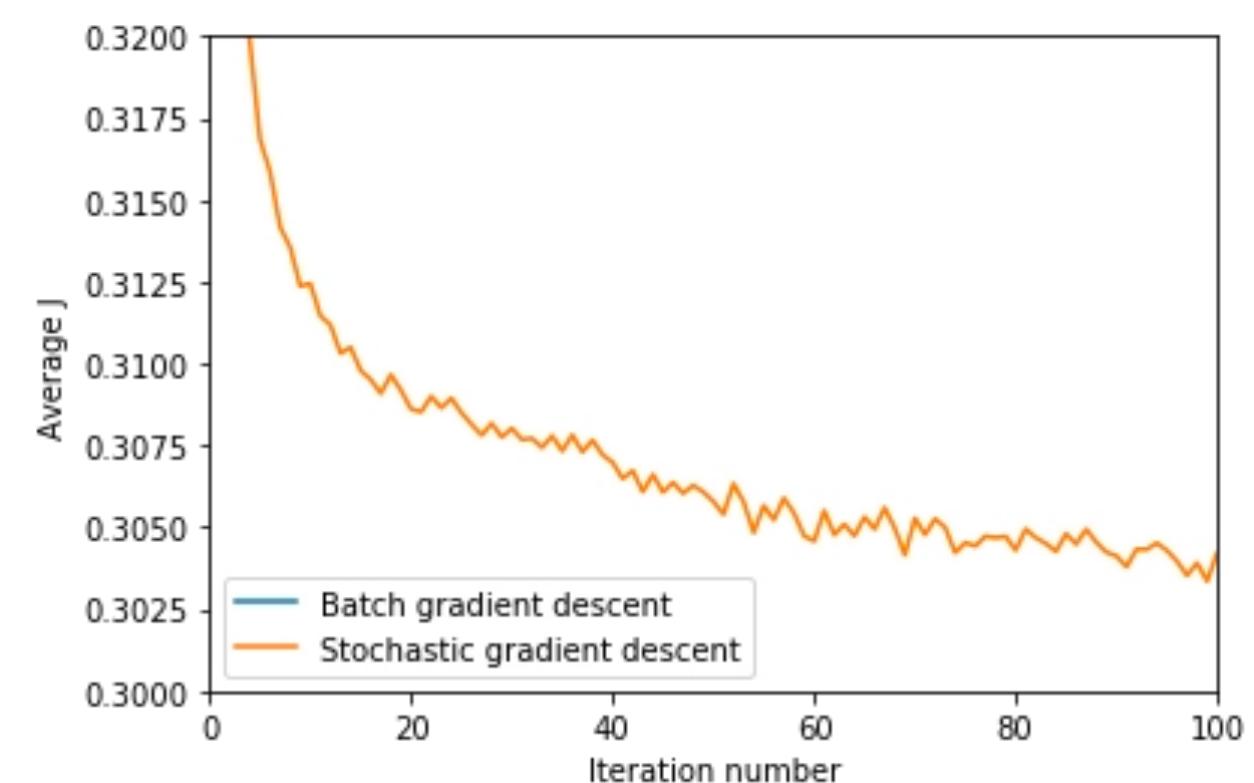
$$E_n(W) = \frac{1}{2}(t - \hat{y}_n)^2$$
$$n = 1 \text{ to } N$$



1. Single sample로 loss를 계산한다.
2. Weight에 대한 gradient를 구한다.
3. backpropagation으로 한 번 weight가 update된다.
4. Updated weight로 다음 sample data에 대한 loss를 구한다.
5. 2~4번을 전체 데이터에 대하여 반복한다.

Stochastic Gradient descent

$$E_n(W) = \frac{1}{2}(t - \hat{y}_n)^2$$
$$n = 1 \text{ to } N$$



- 1.
 - 2.
 - 3.
 - 4.
 - 5.
- 속도는 빠르지만,
data의 noise 때문에 학습이 불안정하고
BGD 보다 성능이 떨어진다.
Updated weight도 나음 sample data에 내안 loss를 구안나.
2~4번을 전체 데이터에 대하여 반복한다.

Mini-Batch Gradient descent

$$E_k(W) = \frac{1}{2} \sum_{m=1}^M (t - \hat{y}_m)^2$$

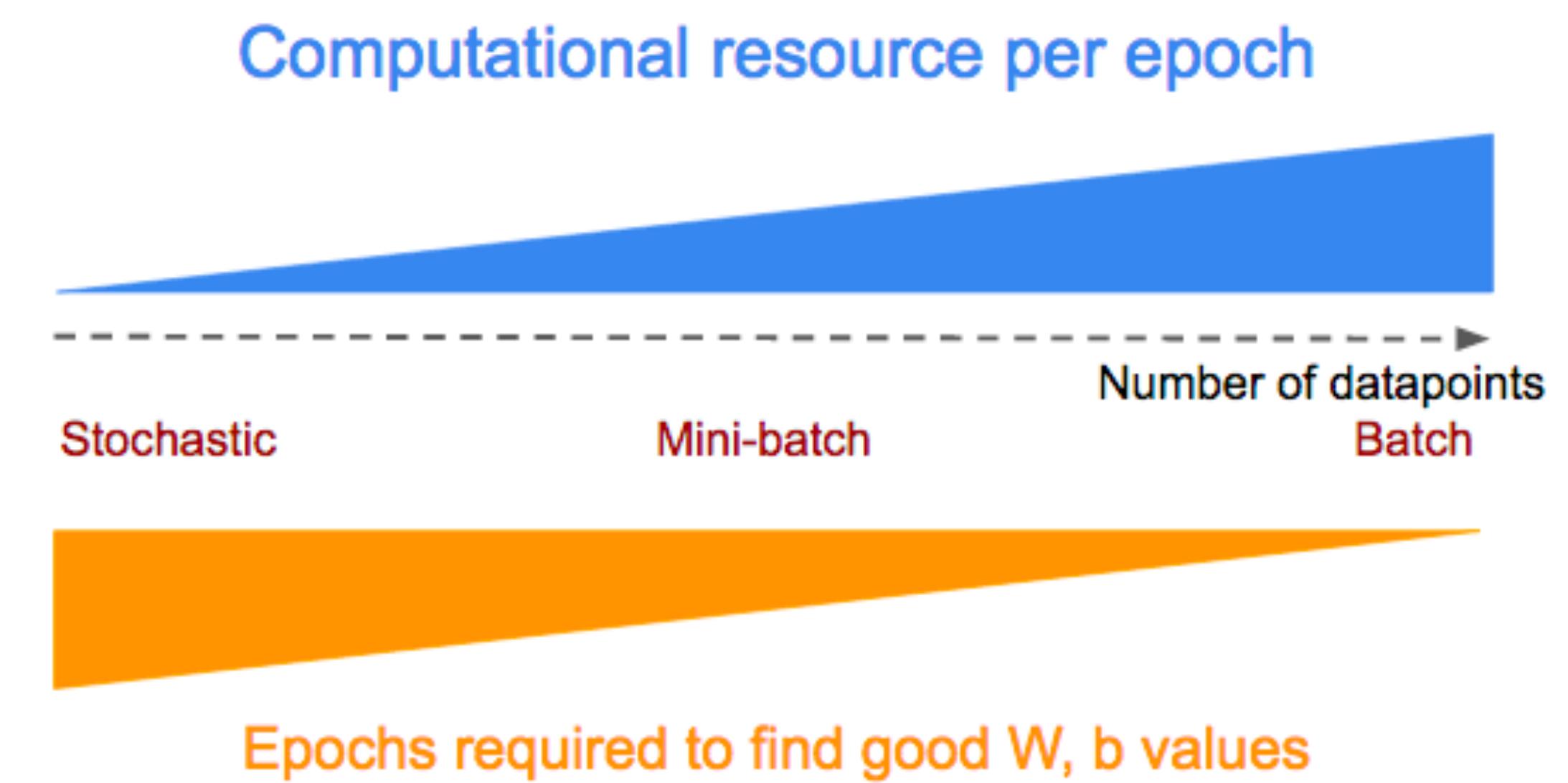
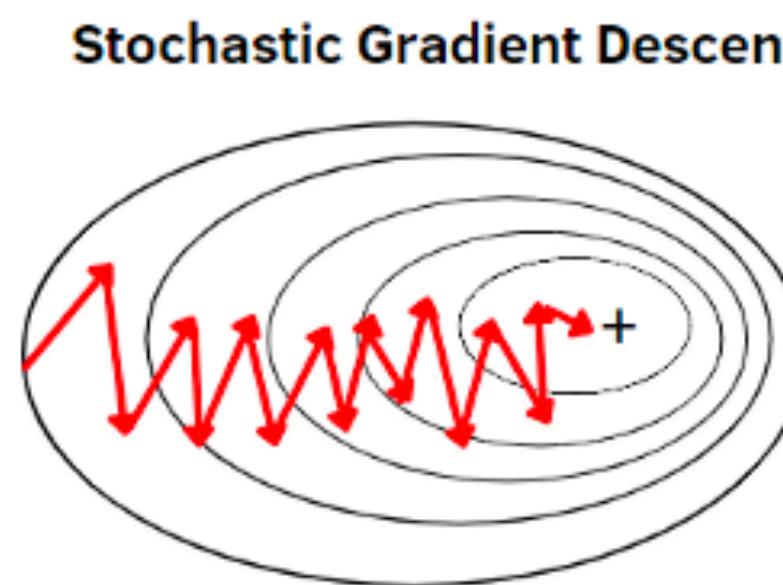
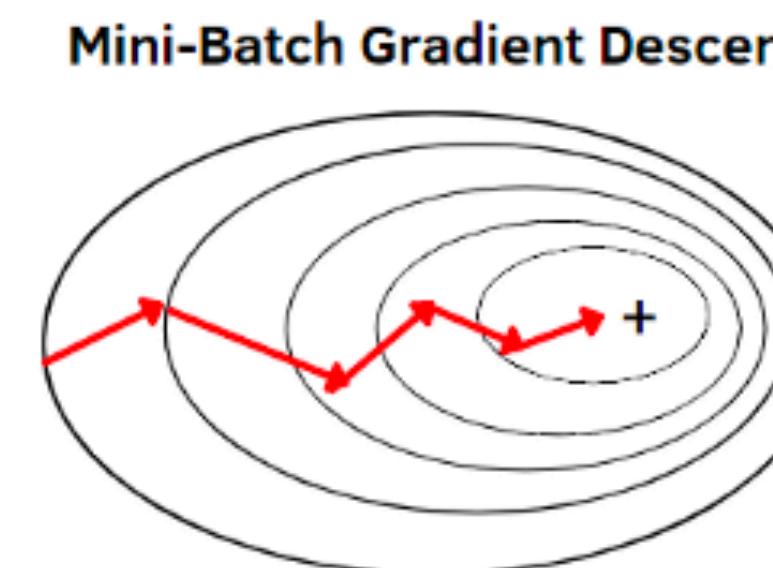
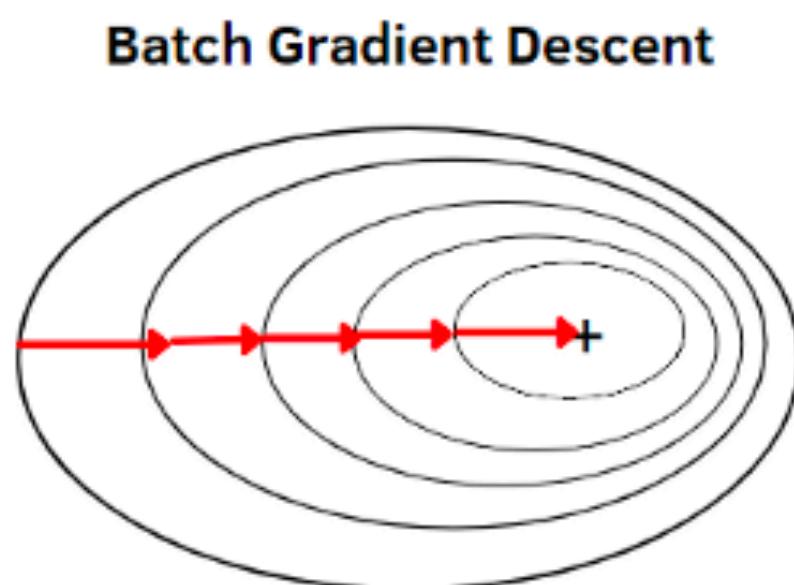
$k = 1$ to K (number of mini-batch)

$M = N/K$ (size of single mini-batch)

1. Mini batch size 만큼 loss를 계산한다.
2. Weight에 대한 gradient를 구한다.
3. backpropagation으로 한 번 weight가 update된다.
4. Updated weight로 다음 mini batch samples에 대한 loss를 구한다.
5. 2~4번을 전체 데이터에 대하여 반복한다.

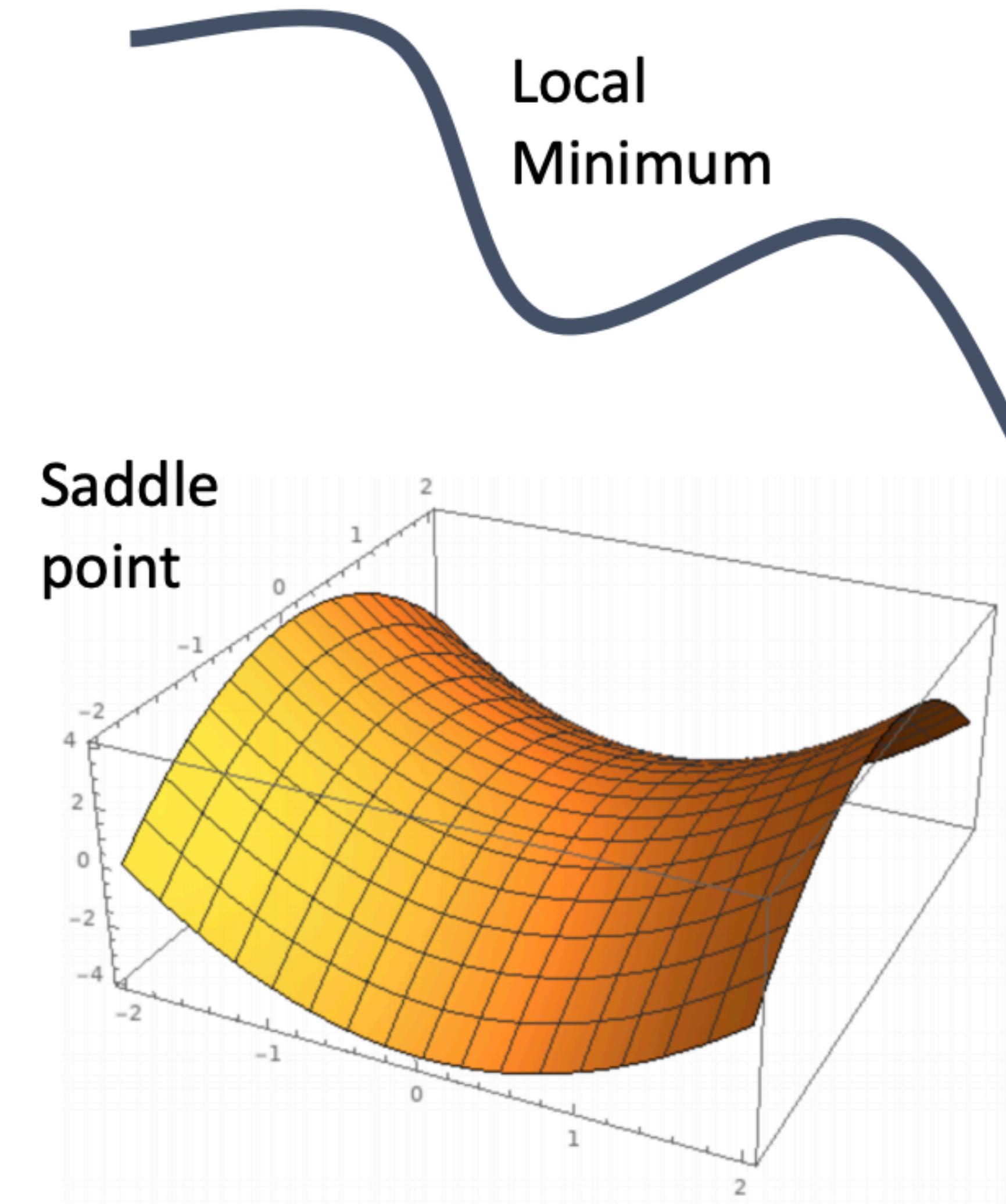
1. Optimization gradient descent

Put all GDs together



Problems with SGD

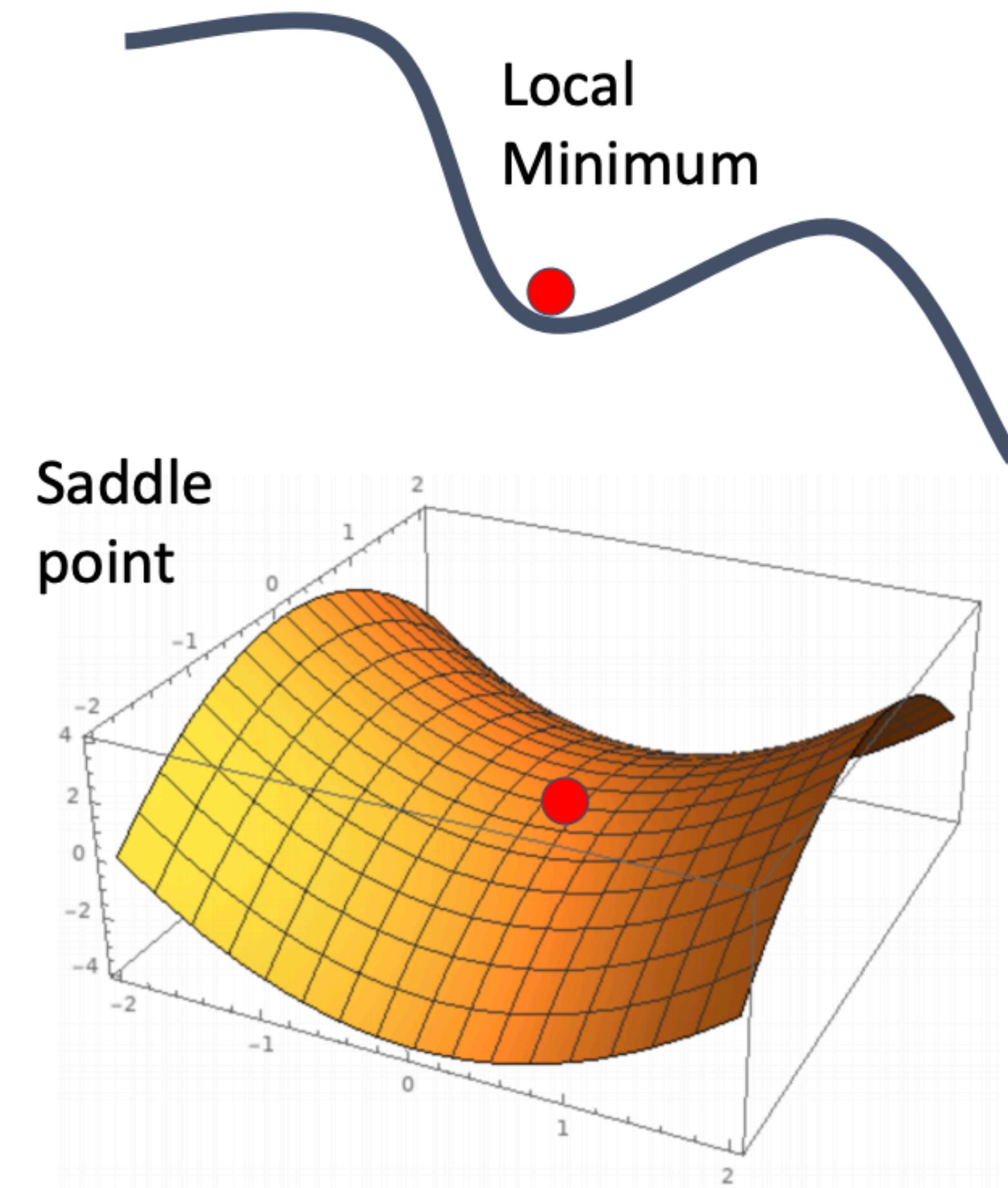
What if the loss function has a **local minimum** or **saddle point**?



Problems with SGD

What if the loss function has a **local minimum** or **saddle point**?

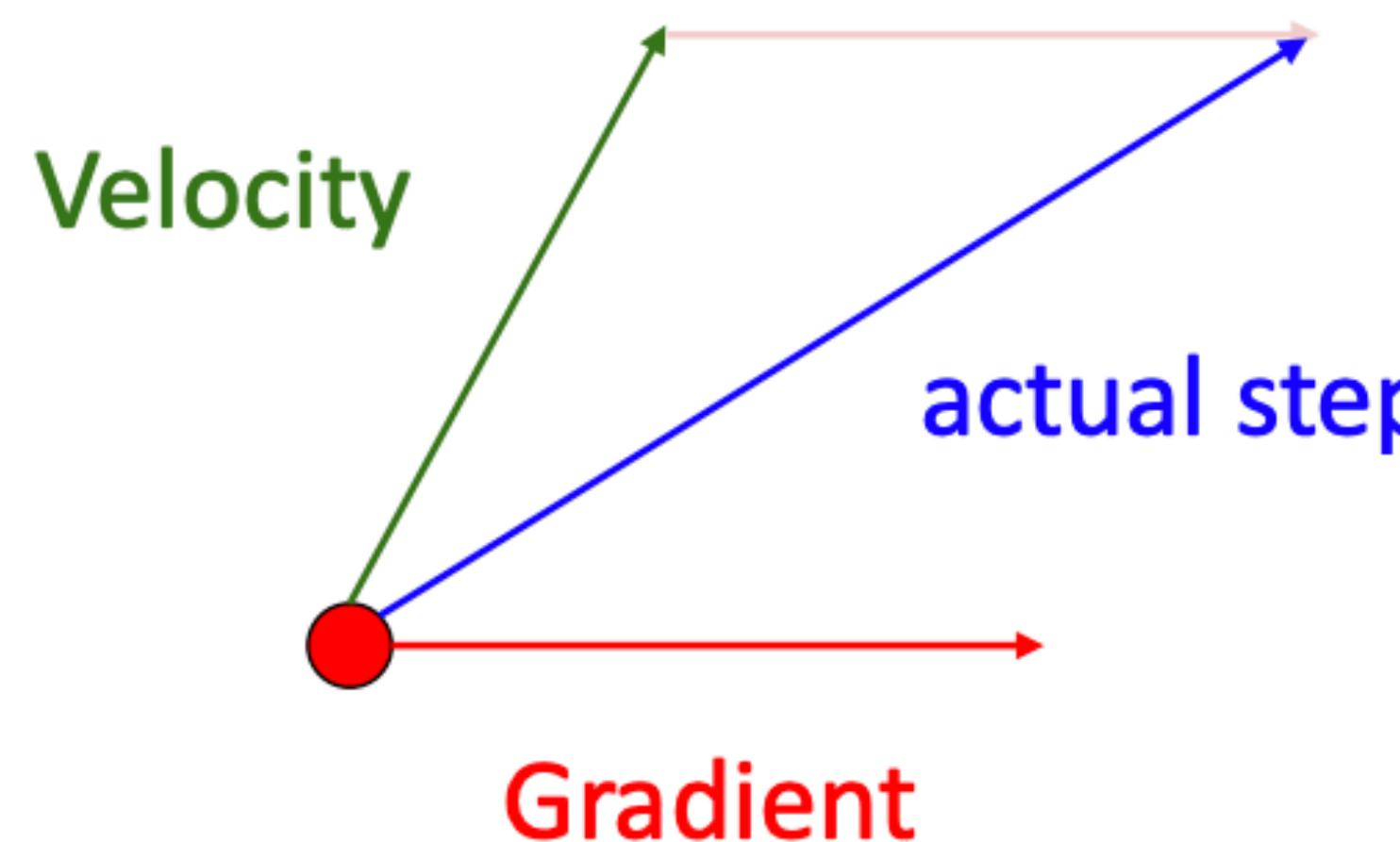
Zero gradient, gradient descent gets stuck



1. Optimization

Better optimizer

Momentum update:



Combine gradient at current point with velocity to get step used to update weights

SGD+Momentum

$$v_{t+1} = \rho v_t + \nabla f(x_t)$$
$$x_{t+1} = x_t - \alpha v_{t+1}$$

```
v = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    v = rho * v + dw
    w -= learning_rate * v
```

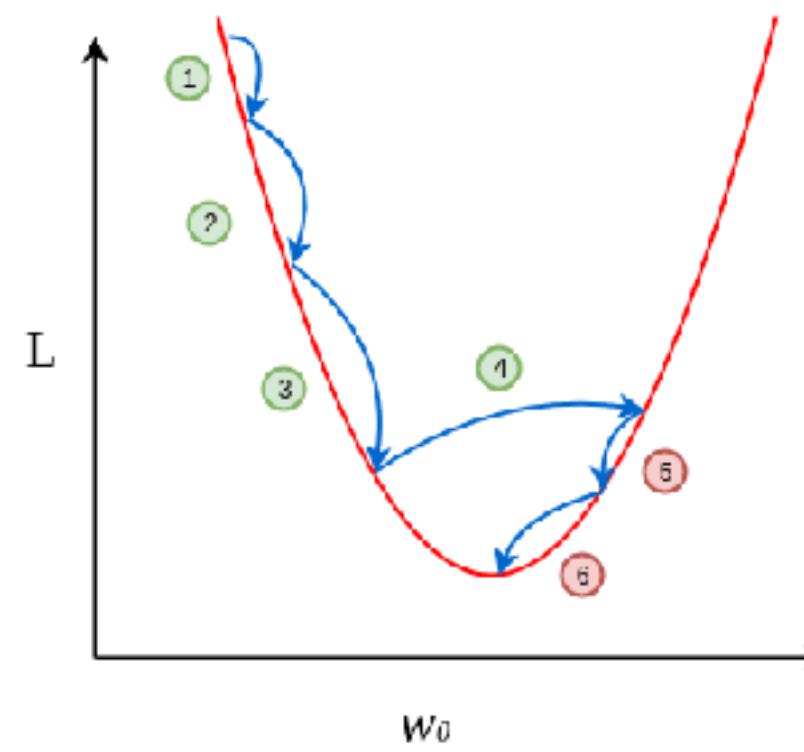
- Build up “velocity” as a running mean of gradients
- Rho gives “friction”; typically rho=0.9 or 0.99

Sutskever et al, “On the importance of initialization and momentum in deep learning”, ICML 2013

Nesterov Momentum

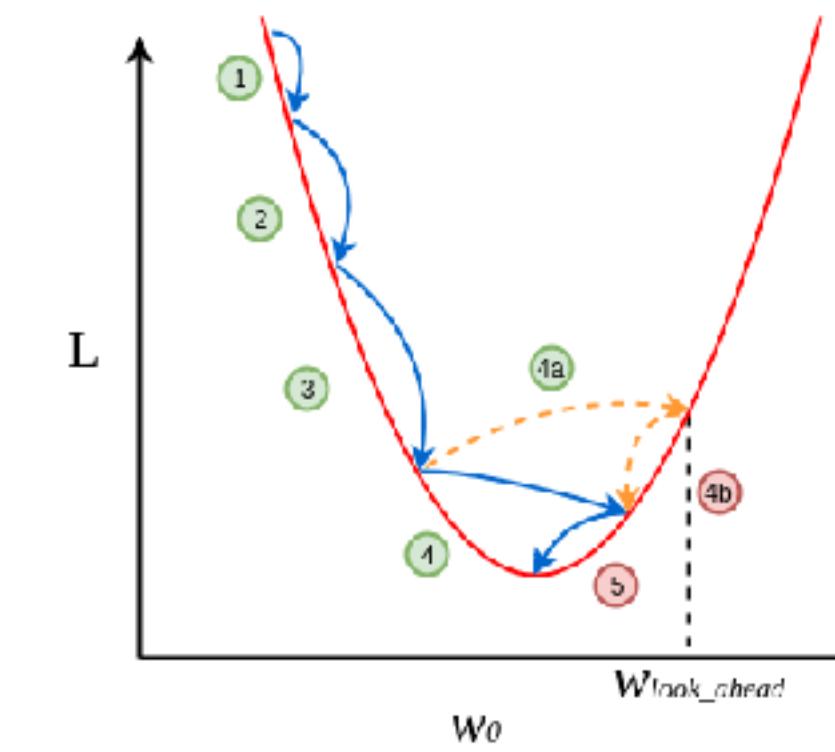
$$v_{t+1} = \rho v_t - \alpha \nabla f(x_t + \rho v_t)$$

$$x_{t+1} = x_t + v_{t+1}$$



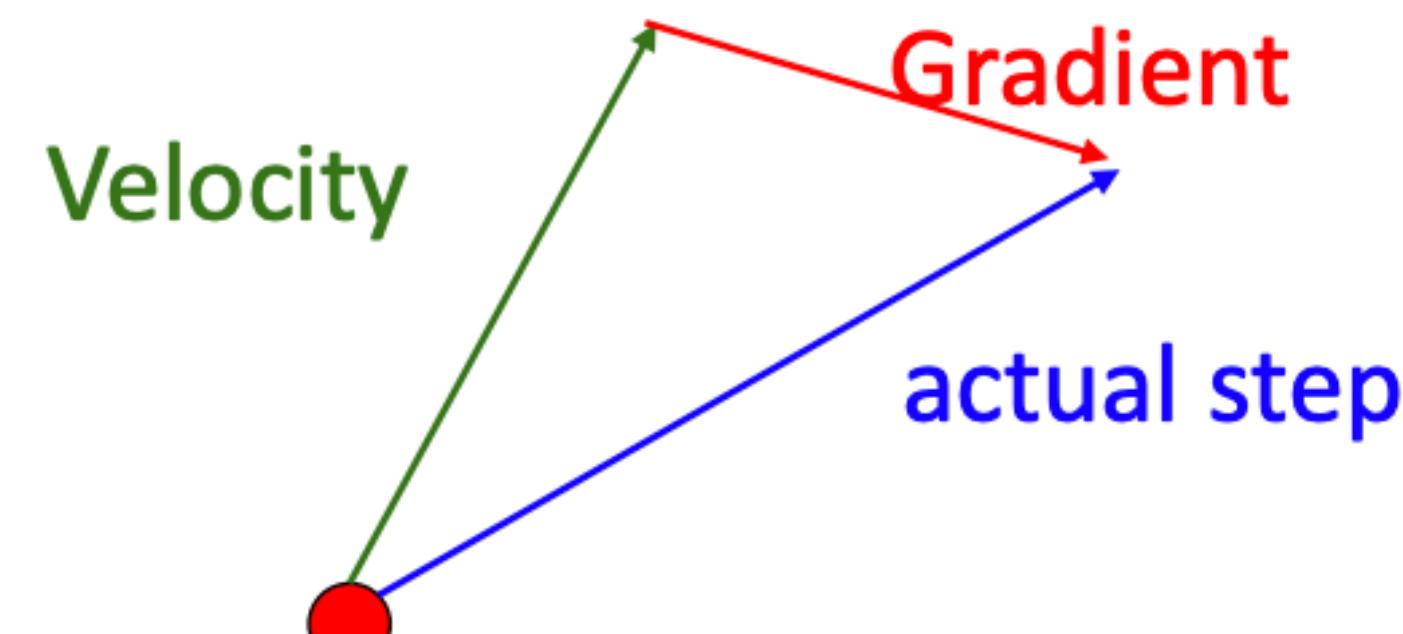
(a) Momentum-Based Gradient Descent

$$\text{○} \Rightarrow \frac{\partial L}{\partial w_0} - \frac{\text{Negative}(-)}{\text{Positive}(+)}$$



(b) Nesterov Accelerated Gradient Descent

$$\text{○} \Rightarrow \frac{\partial L}{\partial w_0} - \frac{\text{Negative}(-)}{\text{Negative}(-)}$$



“Look ahead” to the point where updating using velocity would take us; compute gradient there and mix it with velocity to get actual update direction

1. Optimization

Better optimizer

Nesterov Momentum

$$v_{t+1} = \rho v_t - \alpha \nabla f(x_t + \rho v_t)$$

$$x_{t+1} = x_t + v_{t+1}$$

Annoying, usually we want update in terms of $x_t, \nabla f(x_t)$

Change of variables $\tilde{x}_t = x_t + \rho v_t$
and rearrange:

$$\begin{aligned} v_{t+1} &= \rho v_t - \alpha \nabla f(\tilde{x}_t) \\ \tilde{x}_{t+1} &= \tilde{x}_t - \rho v_t + (1 + \rho)v_{t+1} \\ &= \tilde{x}_t + v_{t+1} + \rho(v_{t+1} - v_t) \end{aligned}$$

```
v = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    old_v = v
    v = rho * v - learning_rate * dw
    w -= rho * old_v - (1 + rho) * v
```

1. Optimization

Better optimizer

AdaGrad

$$\begin{aligned} h &\leftarrow h + \frac{\partial L}{\partial W} \odot \frac{\partial L}{\partial W} \\ W &\leftarrow W - lr \frac{1}{\sqrt{h}} \frac{\partial L}{\partial W} \end{aligned}$$

AdaGrad

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\varepsilon I + \text{diag}(G_t)}} \cdot g_t, \quad (1)$$

$$g_t = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \mathcal{L}(x^{(i)}, y^{(i)}, \theta_t), \quad (2)$$

$$G_t = \sum_{\tau=1}^t g_{\tau} g_{\tau}^{\top}. \quad (3)$$

$$\begin{bmatrix} \theta_{t+1}^{(1)} \\ \theta_{t+1}^{(2)} \\ \vdots \\ \theta_{t+1}^{(m)} \end{bmatrix} = \begin{bmatrix} \theta_t^{(1)} \\ \theta_t^{(2)} \\ \vdots \\ \theta_t^{(m)} \end{bmatrix} - \begin{bmatrix} \frac{\eta}{\sqrt{\varepsilon + G_t^{(1,1)}}} & 0 & \cdots & 0 \\ 0 & \frac{\eta}{\sqrt{\varepsilon + G_t^{(2,2)}}} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{\eta}{\sqrt{\varepsilon + G_t^{(m,m)}}} \end{bmatrix} \cdot \begin{bmatrix} g_t^{(1)} \\ g_t^{(2)} \\ \vdots \\ g_t^{(m)} \end{bmatrix}. \quad (5)$$

1. Optimization

Better optimizer

AdaGrad

$$\begin{bmatrix} \theta_{t+1}^{(1)} \\ \theta_{t+1}^{(2)} \\ \vdots \\ \theta_{t+1}^{(m)} \end{bmatrix} = \begin{bmatrix} \theta_t^{(1)} \\ \theta_t^{(2)} \\ \vdots \\ \theta_t^{(m)} \end{bmatrix} - \begin{bmatrix} \frac{\eta}{\sqrt{\varepsilon + G_t^{(1,1)}}} g_t^{(1)} \\ \frac{\eta}{\sqrt{\varepsilon + G_t^{(2,2)}}} g_t^{(2)} \\ \vdots \\ \frac{\eta}{\sqrt{\varepsilon + G_t^{(m,m)}}} g_t^{(m)}. \end{bmatrix} \quad (6)$$

$$\begin{bmatrix} \theta_{t+1}^{(1)} \\ \theta_{t+1}^{(2)} \\ \vdots \\ \theta_{t+1}^{(m)} \end{bmatrix} = \begin{bmatrix} \theta_t^{(1)} \\ \theta_t^{(2)} \\ \vdots \\ \theta_t^{(m)} \end{bmatrix} - \begin{bmatrix} \eta g_t^{(1)} \\ \eta g_t^{(2)} \\ \vdots \\ \eta g_t^{(m)}. \end{bmatrix}, \quad (7)$$

1. Optimization

Better optimizer

RMSProp: “Leaky Adagrad”

```
grad_squared = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    grad_squared += dw * dw
    w -= learning_rate * dw / (grad_squared.sqrt() + 1e-7)
```

AdaGrad



```
grad_squared = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    grad_squared = decay_rate * grad_squared + (1 - decay_rate) * dw * dw
    w -= learning_rate * dw / (grad_squared.sqrt() + 1e-7)
```

RMSProp

Tieleman and Hinton, 2012

Adam

Algorithm 1: Adam, our proposed algorithm for stochastic optimization. See section 2 for details, and for a slightly more efficient (but less clear) order of computation. g_t^2 indicates the elementwise square $g_t \odot g_t$. Good default settings for the tested machine learning problems are $\alpha = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$. All operations on vectors are element-wise. With β_1^t and β_2^t we denote β_1 and β_2 to the power t .

Require: α : Stepsize
Require: $\beta_1, \beta_2 \in [0, 1]$: Exponential decay rates for the moment estimates
Require: $f(\theta)$: Stochastic objective function with parameters θ
Require: θ_0 : Initial parameter vector

$m_0 \leftarrow 0$ (Initialize 1st moment vector)
 $v_0 \leftarrow 0$ (Initialize 2nd moment vector)
 $t \leftarrow 0$ (Initialize timestep)

while θ_t not converged **do**

- $t \leftarrow t + 1$
- $g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$ (Get gradients w.r.t. stochastic objective at timestep t)
- $m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$ (Update biased first moment estimate)
- $v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$ (Update biased second raw moment estimate)
- $\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$ (Compute bias-corrected first moment estimate)
- $\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$ (Compute bias-corrected second raw moment estimate)
- $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$ (Update parameters)

end while

return θ_t (Resulting parameters)

Adam (almost) : RMSProp + Momentum

```
moment1 = 0
moment2 = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    moment1 = beta1 * moment1 + (1 - beta1) * dw
    moment2 = beta2 * moment2 + (1 - beta2) * dw * dw
    w -= learning_rate * moment1 / (moment2.sqrt() + 1e-7)
```

Adam

Momentum

```
v = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    v = rho * v + dw
    w -= learning_rate * v
```

SGD+Momentum

Adam (almost) : RMSProp + Momentum

```
moment1 = 0
moment2 = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    moment1 = beta1 * moment1 + (1 - beta1) * dw
    moment2 = beta2 * moment2 + (1 - beta2) * dw * dw
    w -= learning_rate * moment1 / (moment2.sqrt() + 1e-7)
```

Adam

Momentum

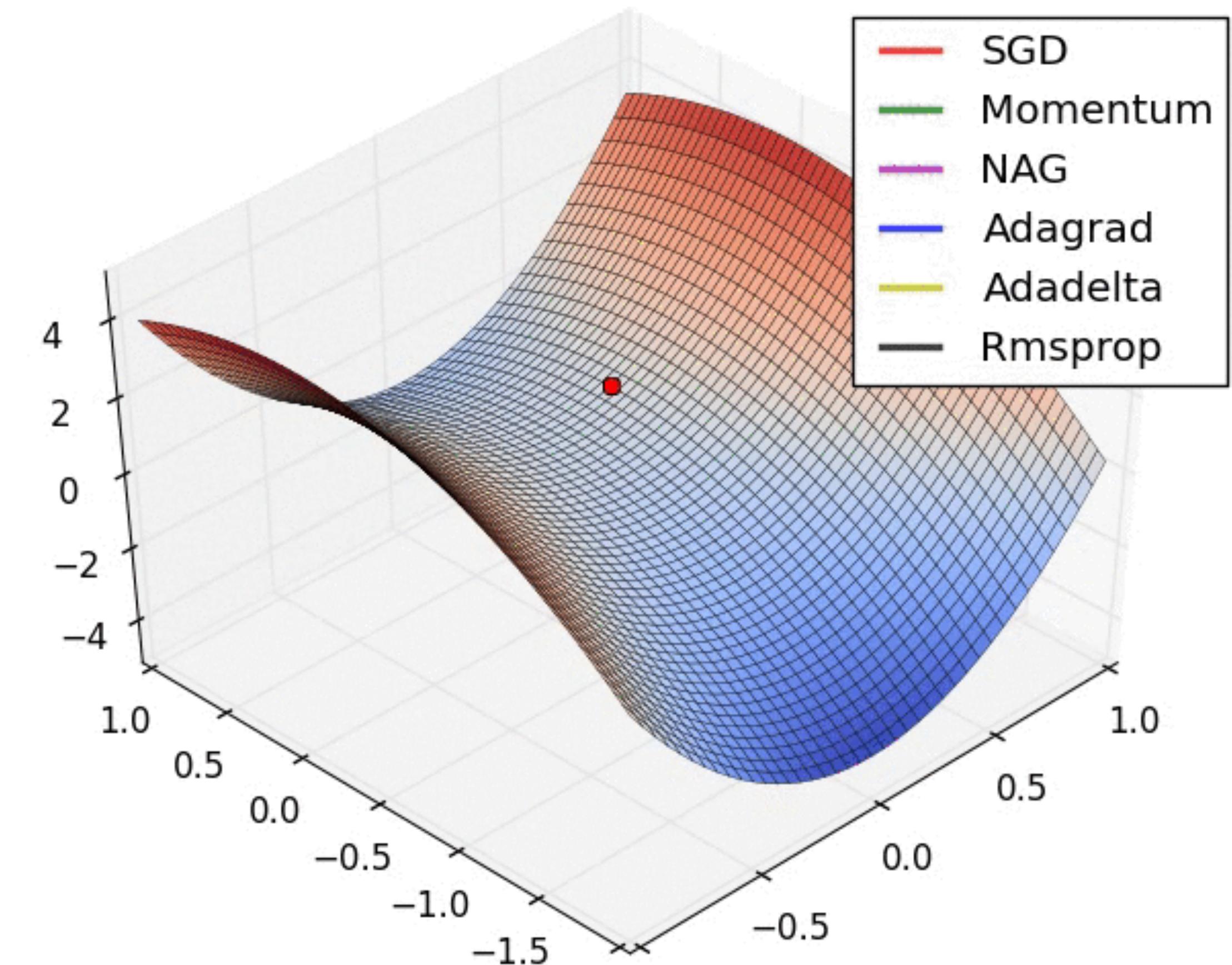
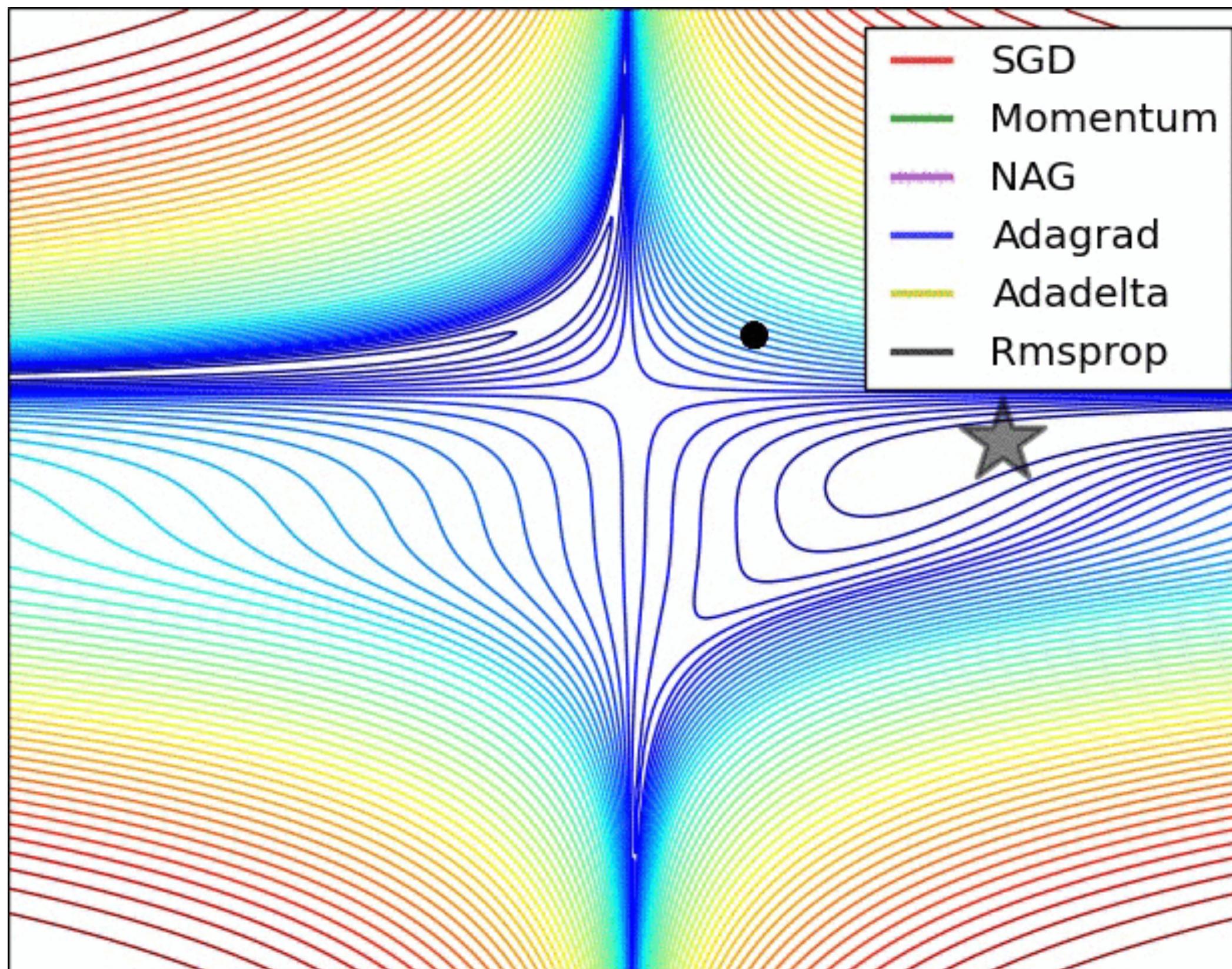
AdaGrad / RMSProp

```
grad_squared = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    grad_squared = decay_rate * grad_squared + (1 - decay_rate) * dw * dw
    w -= learning_rate * dw / (grad_squared.sqrt() + 1e-7)
```

RMSProp

1. Optimization

Better optimizer



1. Optimization

Better optimizer

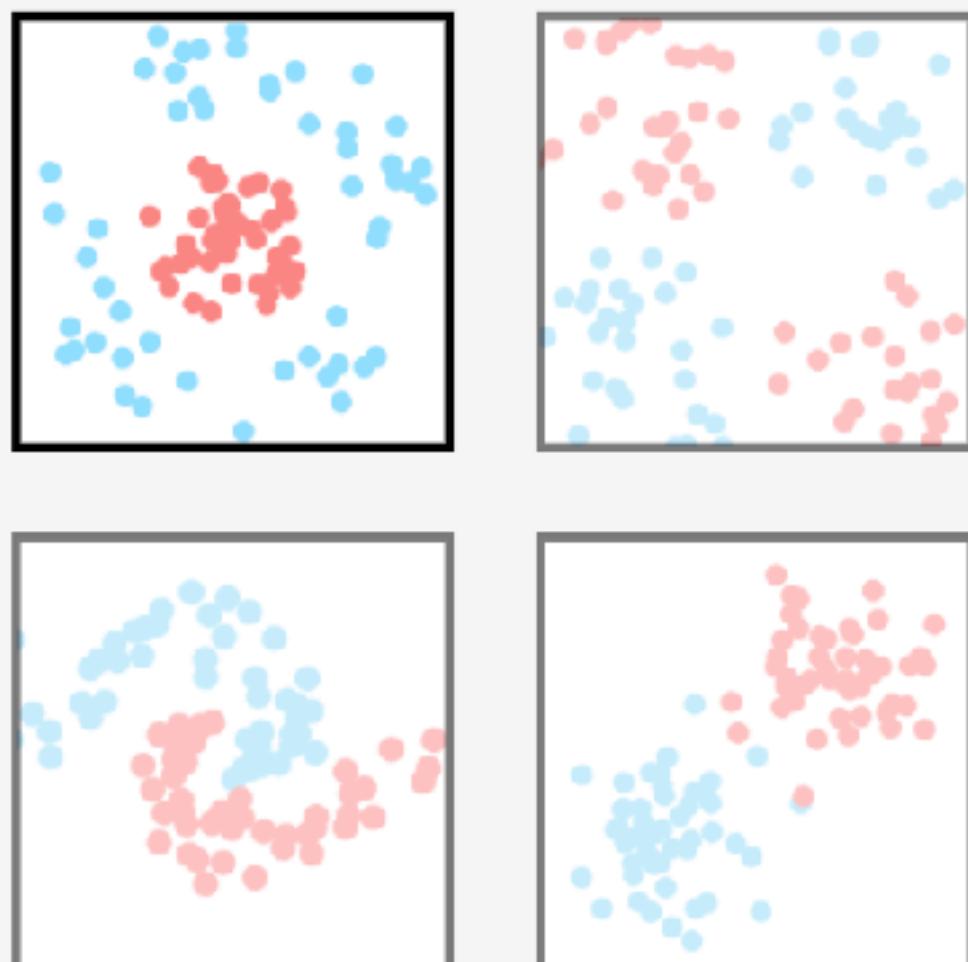
Optimization Algorithm Comparison

Algorithm	Tracks first moments (Momentum)	Tracks second moments (Adaptive learning rates)	Leaky second moments	Bias correction for moment estimates
SGD	✗	✗	✗	✗
SGD+Momentum	✓	✗	✗	✗
Nesterov	✓	✗	✗	✗
AdaGrad	✗	✓	✗	✗
RMSProp	✗	✓	✓	✗
Adam	✓	✓	✓	✓

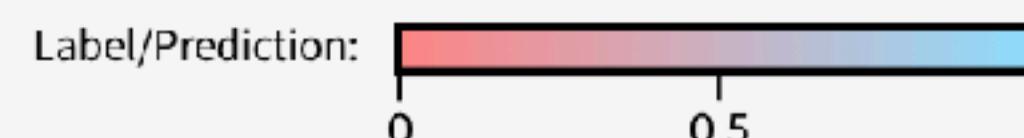
2. Generalization initialization

1. Choose input dataset

Select a training dataset.



This legend details the color scheme for labels, and the values of the weights/gradients.

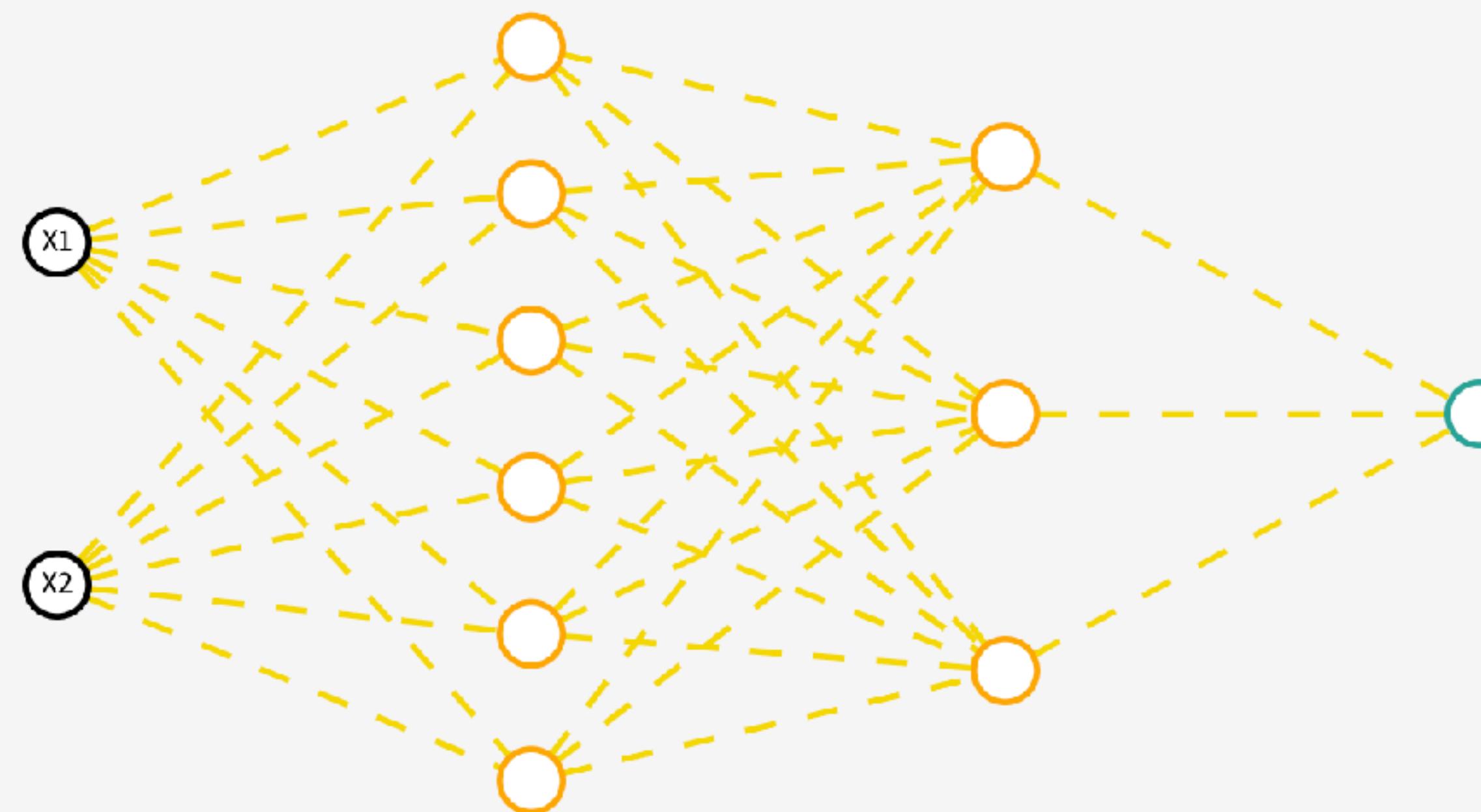


Node Type: Input Relu Sigmoid

2. Choose initialization method

Select an initialization method for the values of your neural network parameters¹.

- Zero
- Too small
- Appropriate
- Too large

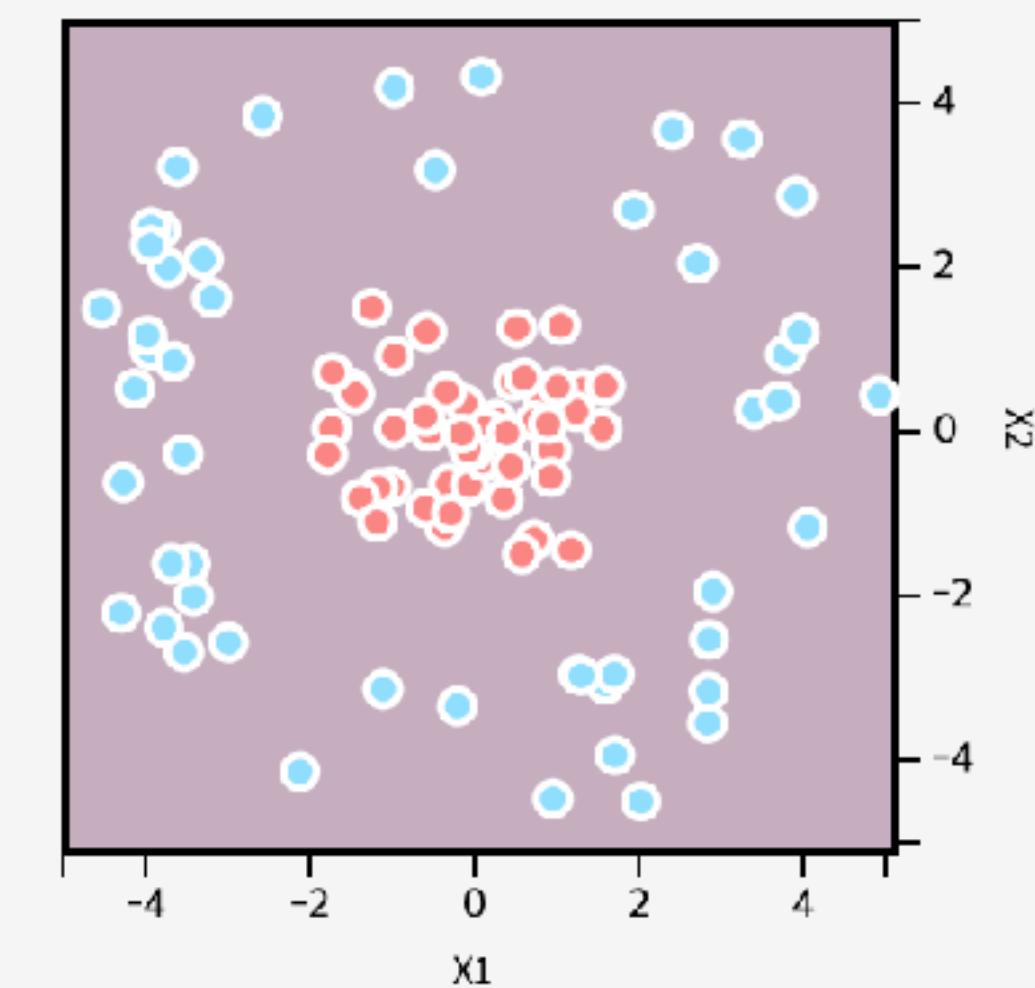
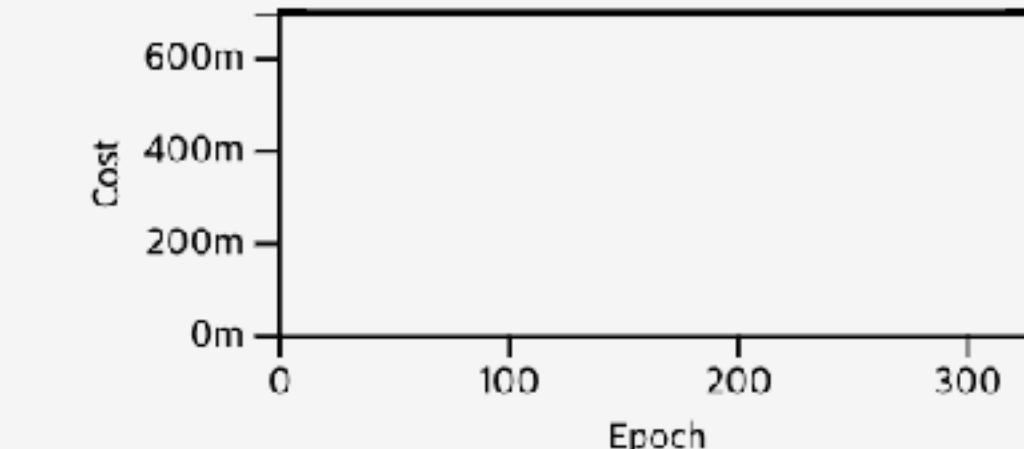


Select whether to visualize the weights or gradients of the network above.

- Weight
- Gradient

3. Train the network.

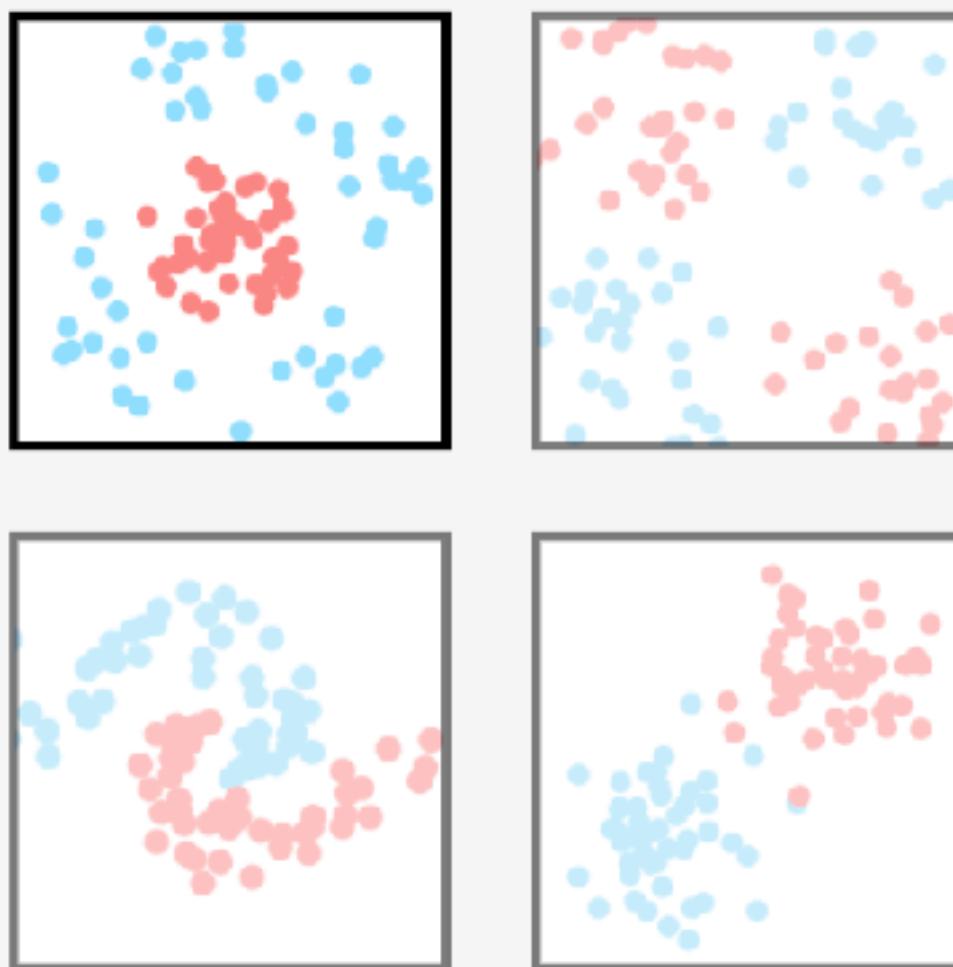
Observe the cost function and the decision boundary.



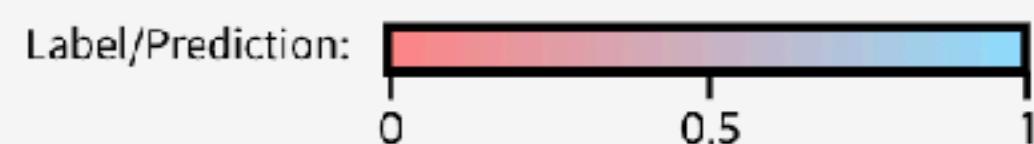
2. Generalization initialization

1. Choose input dataset

Select a training dataset.



This legend details the color scheme for labels, and the values of the weights/gradients.



Node Type: Input Relu Sigmoid

2. Choose initialization method

Select an initialization method for the values of your neural network parameters¹.

- Zero
- Too small
- Appropriate
- Too large

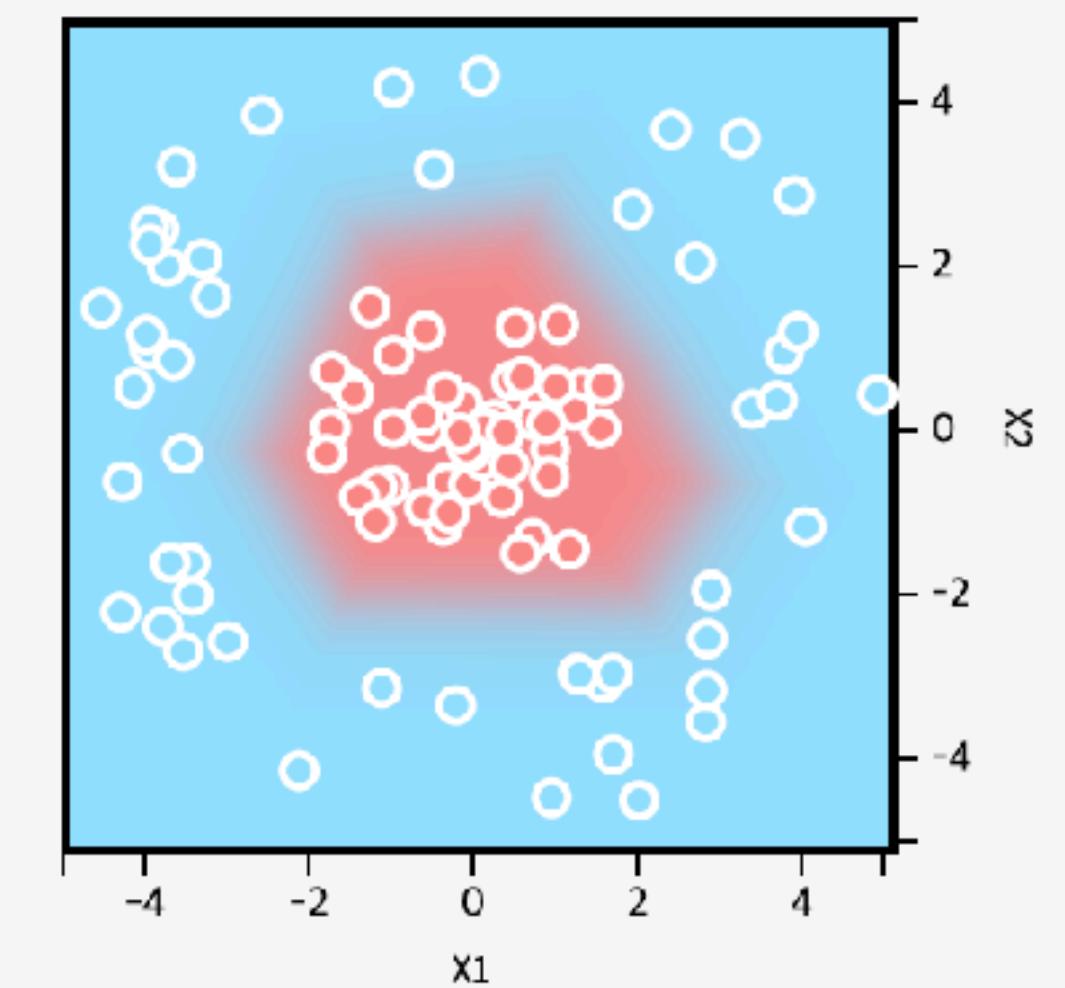
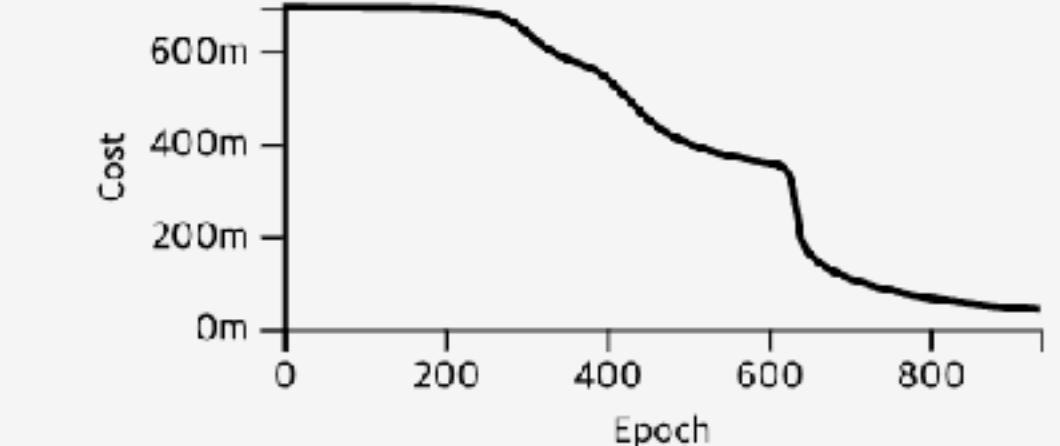


Select whether to visualize the weights or gradients of the network above.

- Weight
- Gradient

3. Train the network.

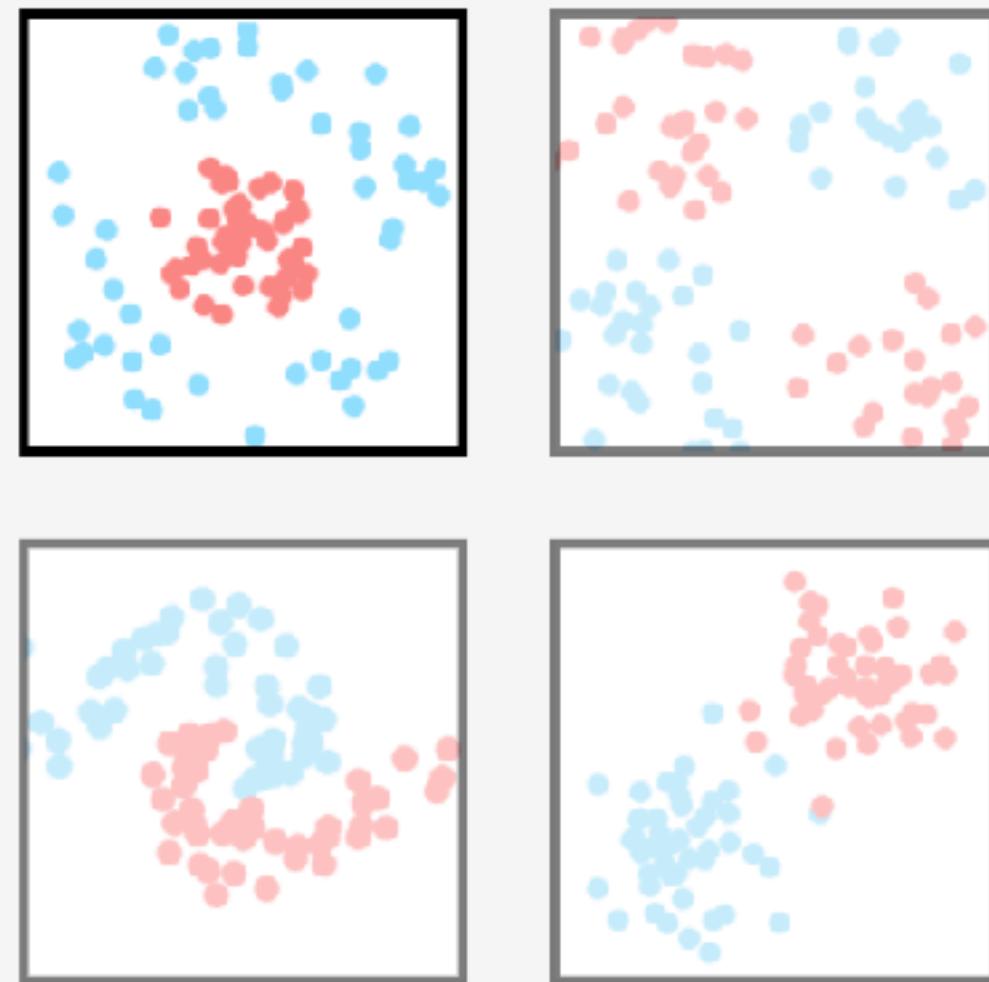
Observe the cost function and the decision boundary.



2. Generalization initialization

1. Choose input dataset

Select a training dataset.



This legend details the color scheme for labels, and the values of the weights/gradients.



Node Type: Input Relu Sigmoid

2. Choose initialization method

Select an initialization method for the values of your neural network parameters¹.

- Zero
- Too small
- Appropriate
- Too large

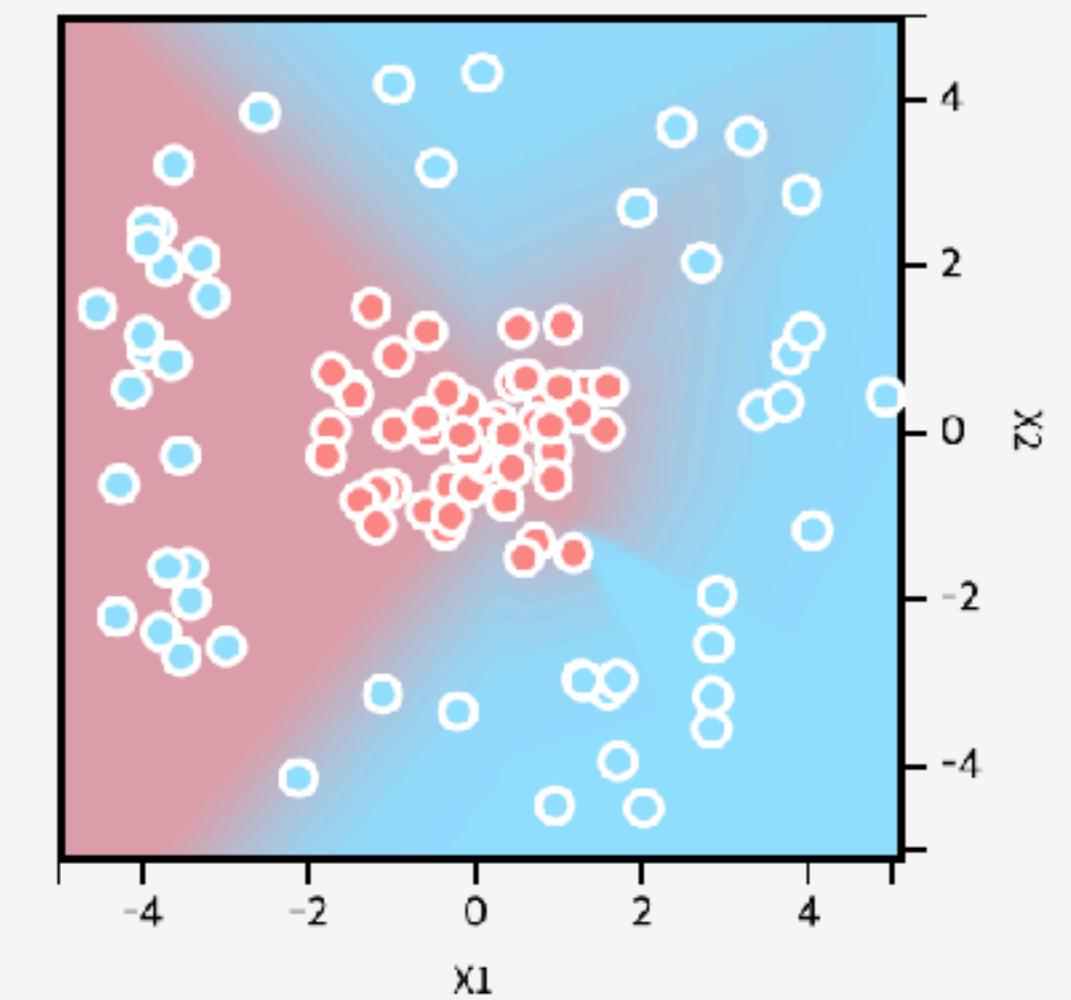
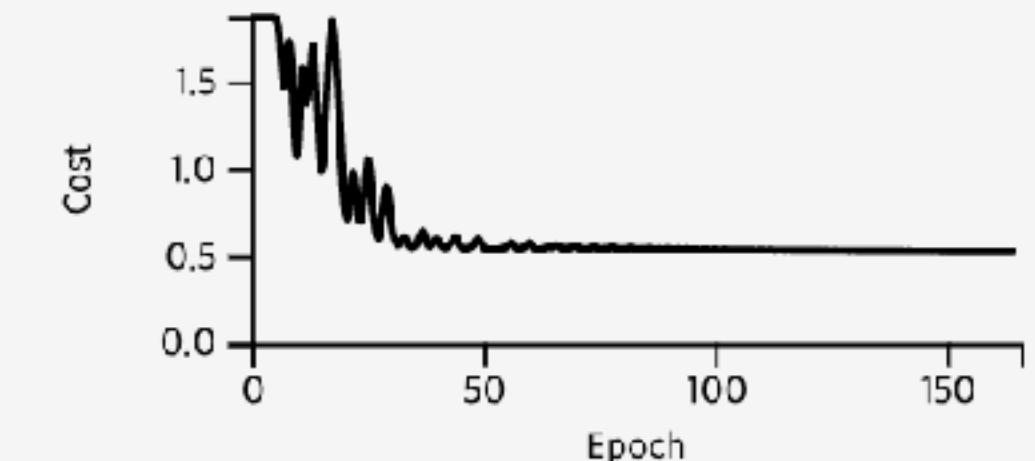


Select whether to visualize the weights or gradients of the network above.

- Weight
- Gradient

3. Train the network.

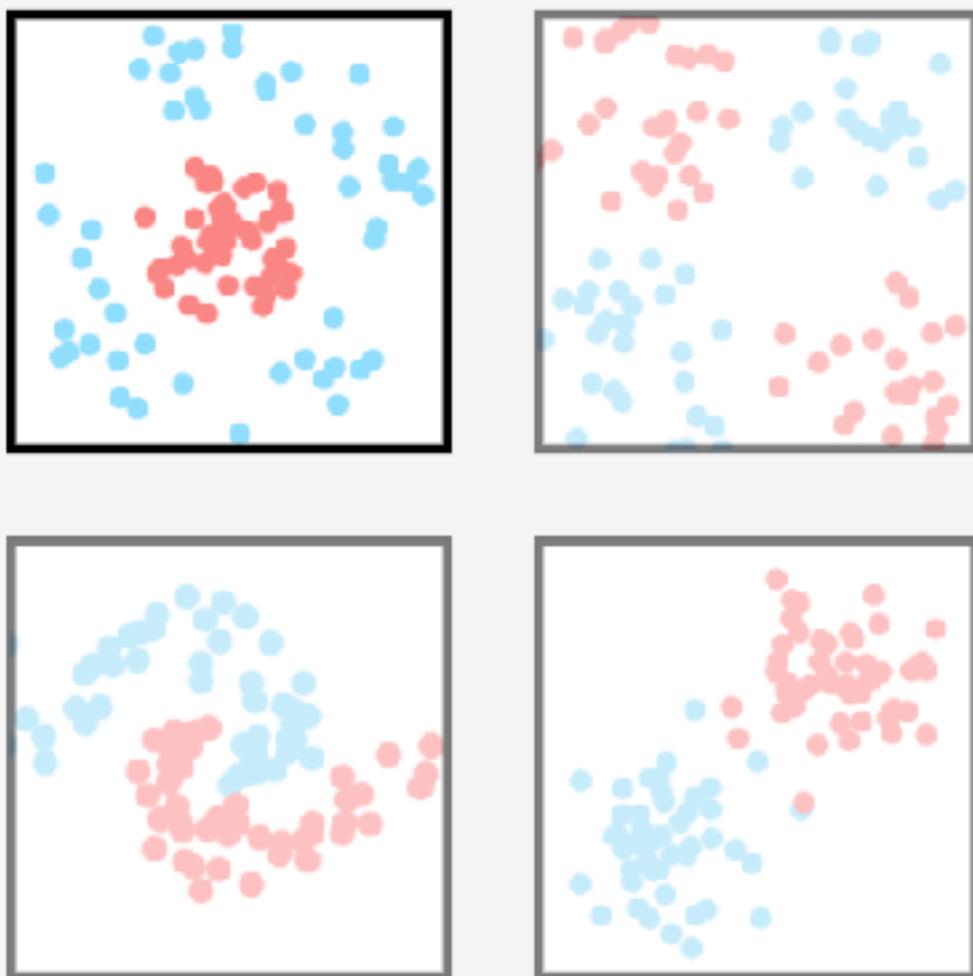
Observe the cost function and the decision boundary.



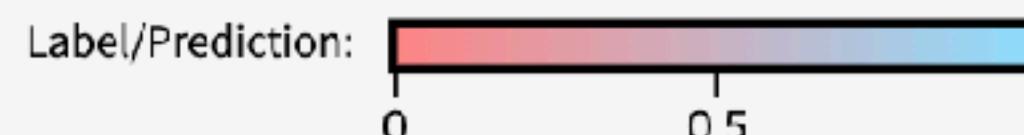
2. Generalization initialization

1. Choose input dataset

Select a training dataset.



This legend details the color scheme for labels, and the values of the weights/gradients.



Node Type: Input Relu Sigmoid

2. Choose initialization method

Select an initialization method for the values of your neural network parameters¹.

- Zero Too small Appropriate Too large

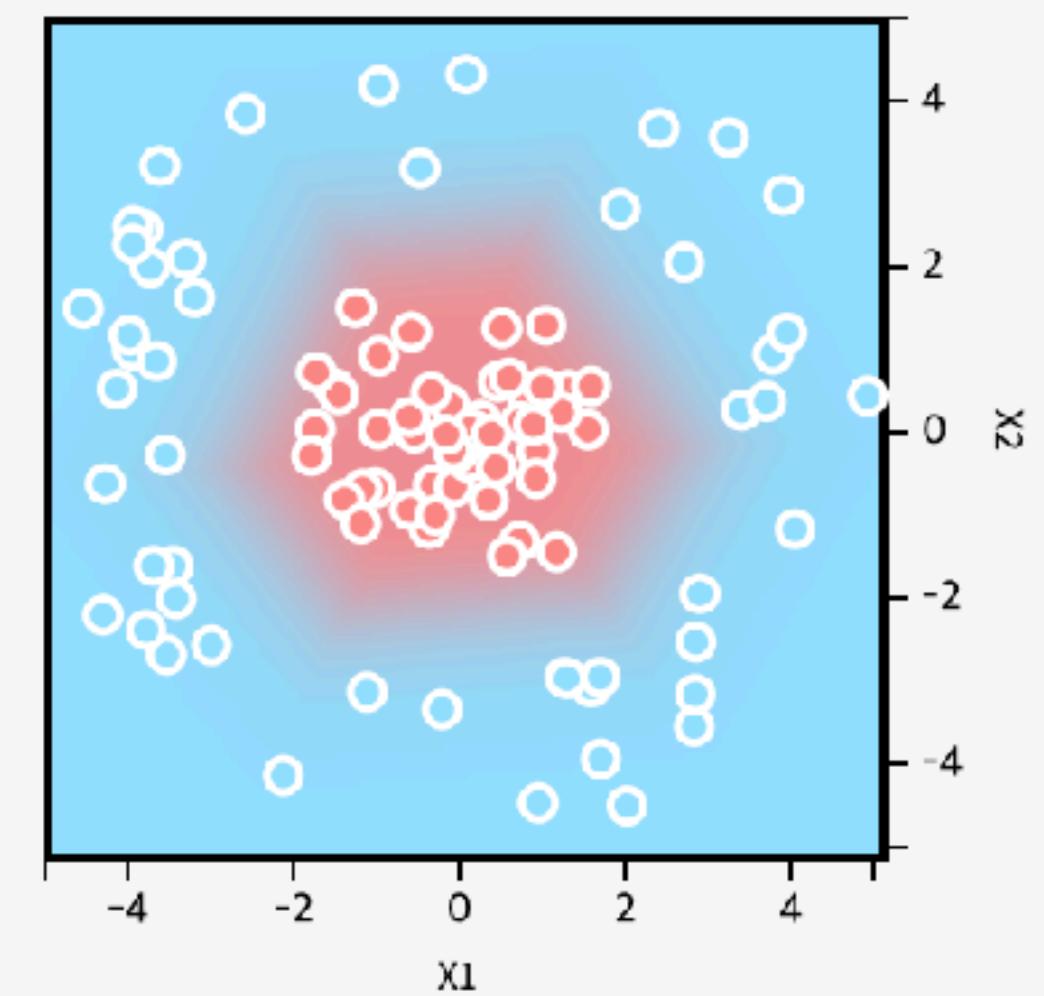
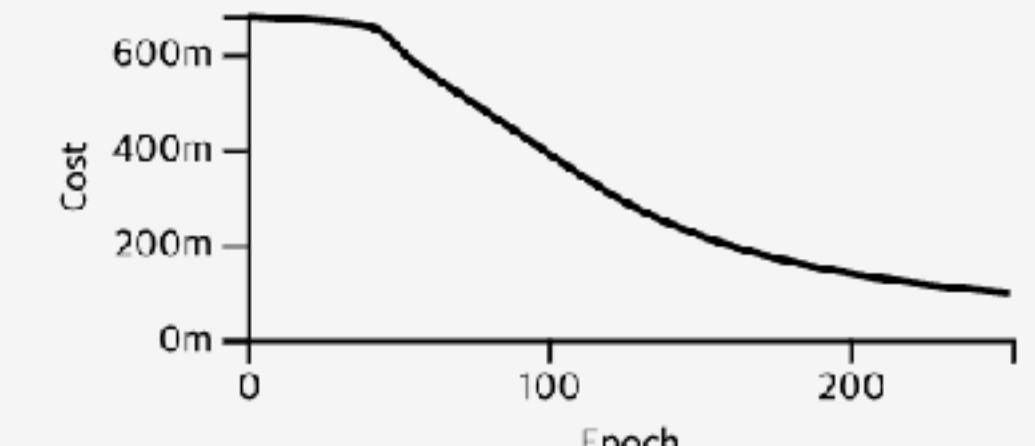


Select whether to visualize the weights or gradients of the network above.

- Weight Gradient

3. Train the network.

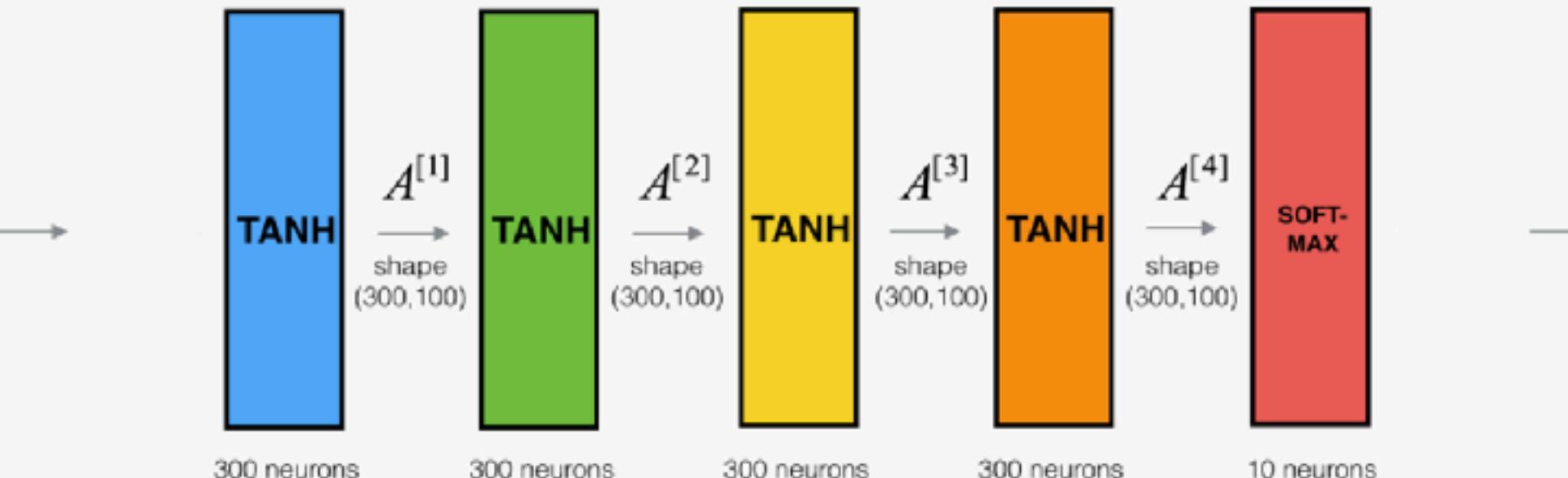
Observe the cost function and the decision boundary.



2. Generalization initialization

$$X = A^{[0]}$$

Batch of 100 grayscale images of shape 28x28
 $X.shape = (784, 100)$ because $784 = 28 \times 28$



$$\hat{y} = A^{[5]}$$

output probability over 10 classes for a batch of 100 images
 $\hat{y}.shape = (10, 100)$

1. Load your dataset

Load 10,000 handwritten digits images (MNIST).

Load MNIST (100%)

2. Select an initialization method

Among the below distributions, select the one to use to initialize your parameters³.

- Zero Uniform Xavier Standard Normal

3. Train the network and observe

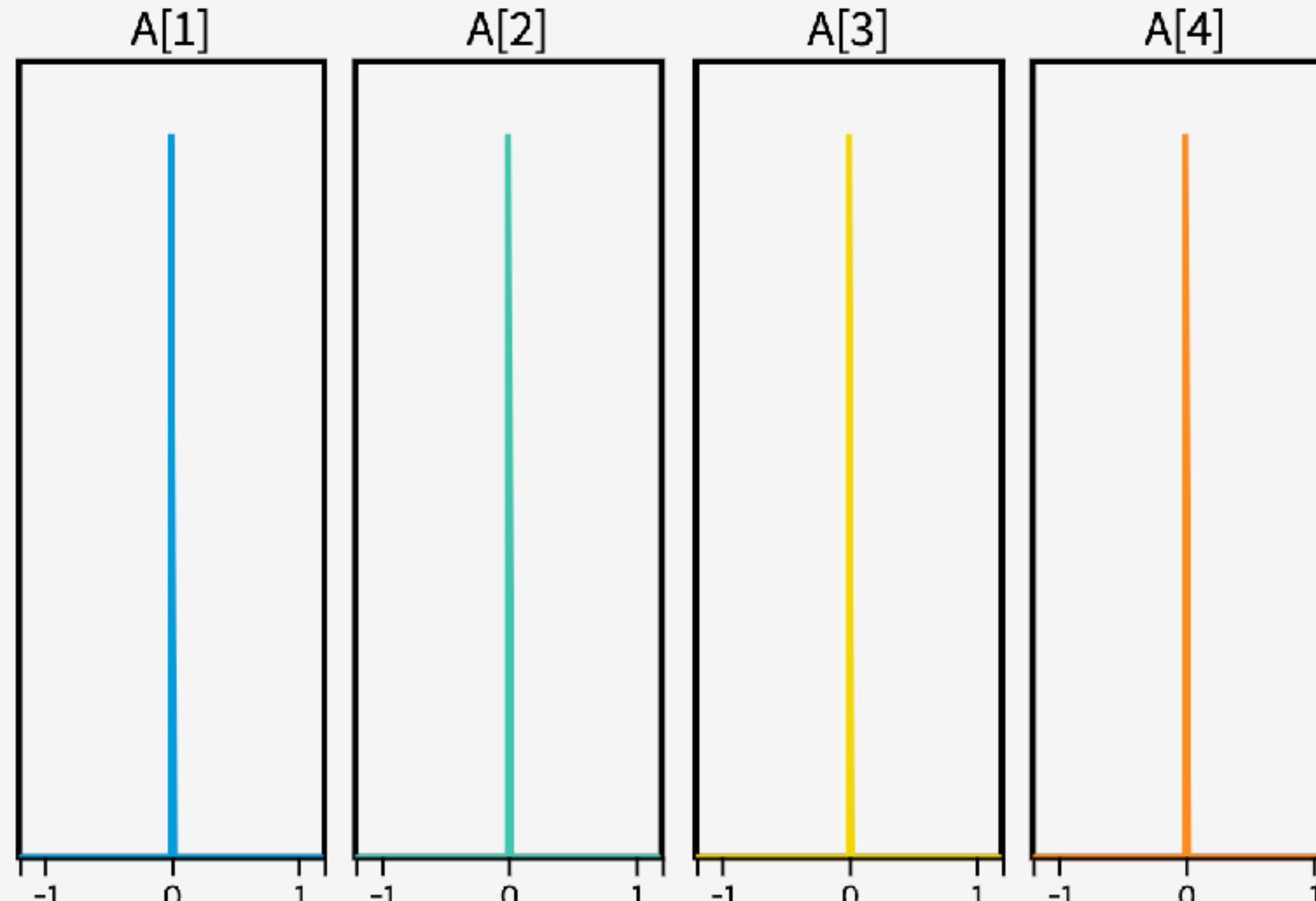
The grid below refers to the input images, Blue squares represent correctly classified images. Red squares represent misclassified images.



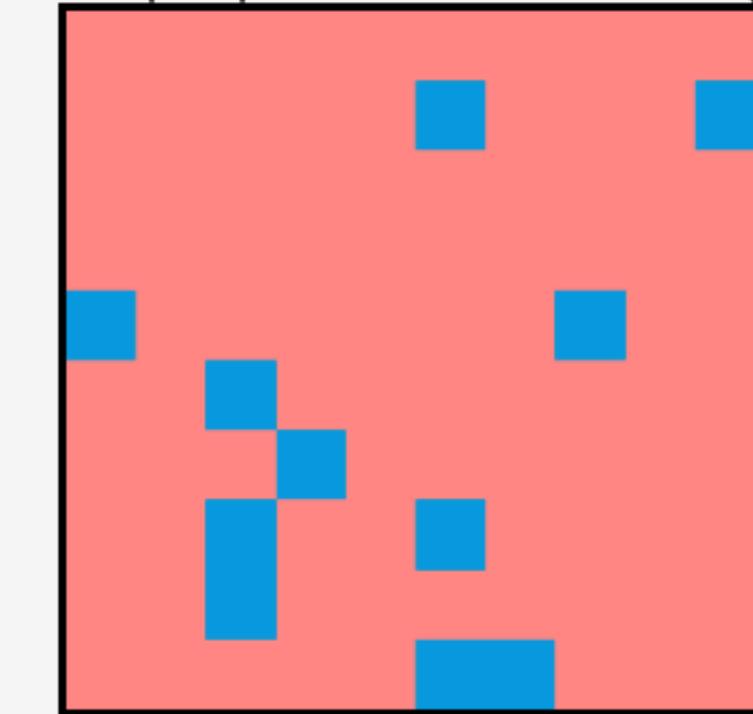
Input batch of 100 images

6 8 6 8 5 7 8 6 0 2
4 0 2 2 3 1 9 7 5 1
0 8 4 6 2 4 7 9 3 2
9 8 2 2 9 2 7 3 5 9
1 8 0 2 0 5 1 3 7
6 7 1 2 5 8 0 3 7 1
4 0 9 1 8 6 7 7 4 3
9 1 9 3 1 7 3 9 7
6 9 1 3 7 8 3 3 6 7
2 8 5 8 5 1 1 4 4 3

Batch: 6 Epoch: 0



Output predictions of 100 images

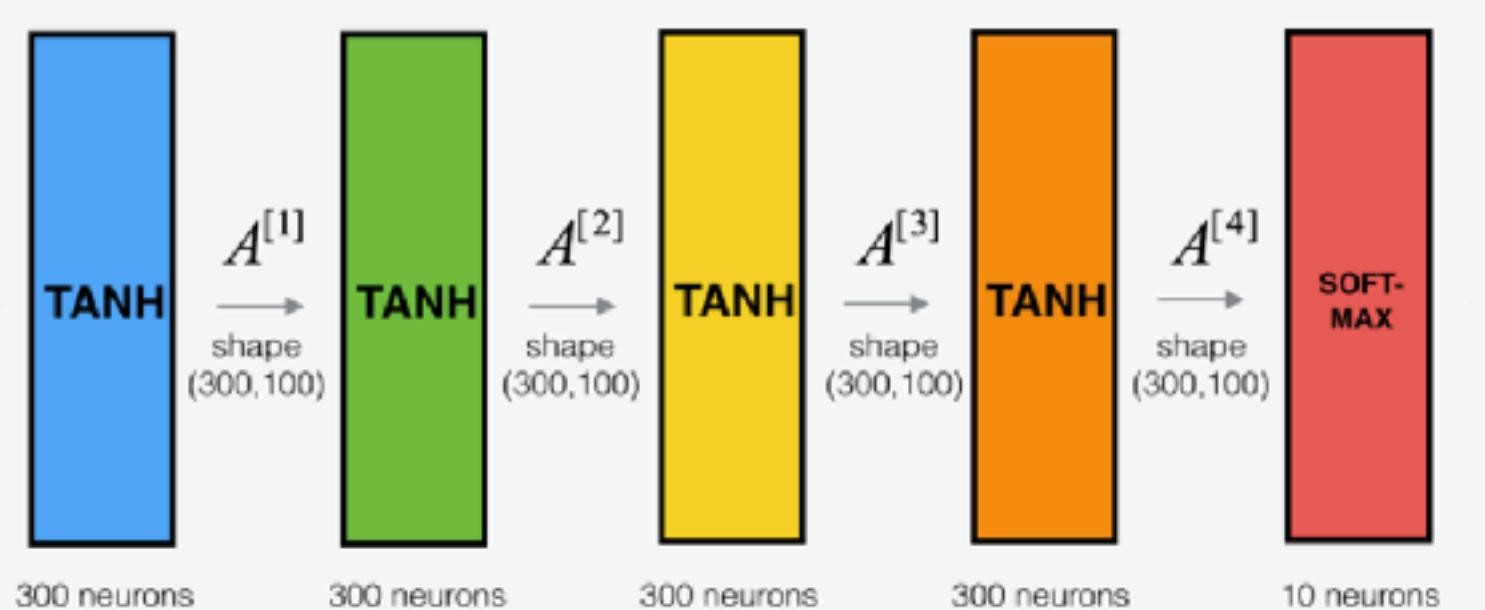


Misclassified: 89/100 Cost: 2.30

2. Generalization initialization

$$X = A^{[0]}$$

Batch of 100 grayscale images of shape 28x28
 $X.shape = (784, 100)$ because $784 = 28 \times 28$



$$\hat{y} = A^{[5]}$$

output probability over 10
classes for a batch of
100 images
 $\hat{y}.shape = (10, 100)$

1. Load your dataset

Load 10,000 handwritten digits images
([MNIST](#)).

Load MNIST (100%)

2. Select an initialization method

Among the below distributions, select the one to use to initialize your parameters³.

Zero

Uniform

Xavier

Standard Normal

Input batch of 100 images

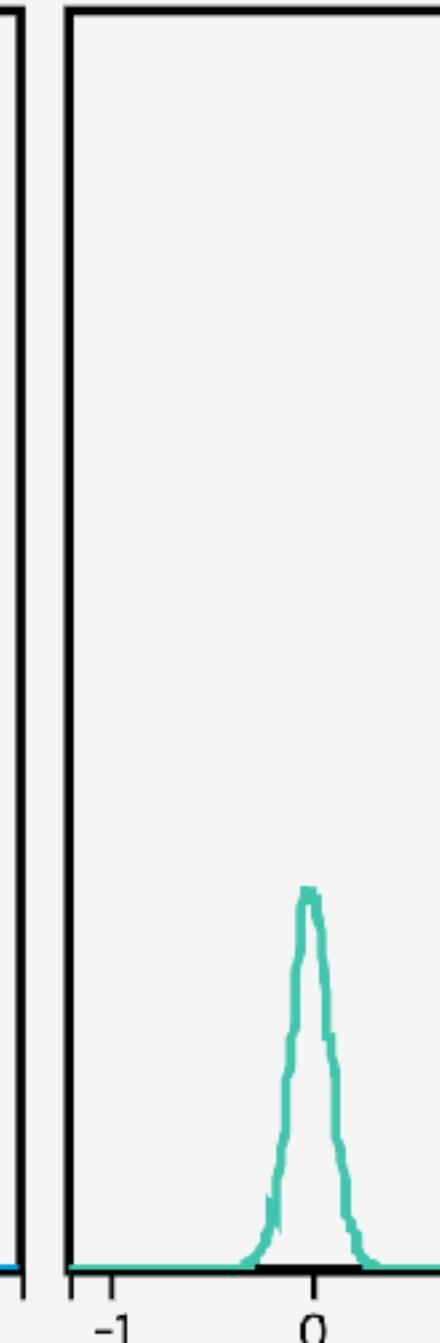
4 7 1 2 4 0 2 7 4 3
3 0 0 3 1 9 6 5 0 5
9 2 9 3 0 4 L 0 7 1
1 2 1 S 3 3 9 7 8 6
5 6 1 3 8 1 0 5 1 3
1 5 5 6 / 8 5 1 7 9
4 6 2 2 5 0 6 5 6 3
7 2 0 8 8 5 4 1 1 4
0 3 3 7 6 1 6 2 1 9
2 8 6 1 9 5 2 5 4 4

Batch: 3 Epoch: 0

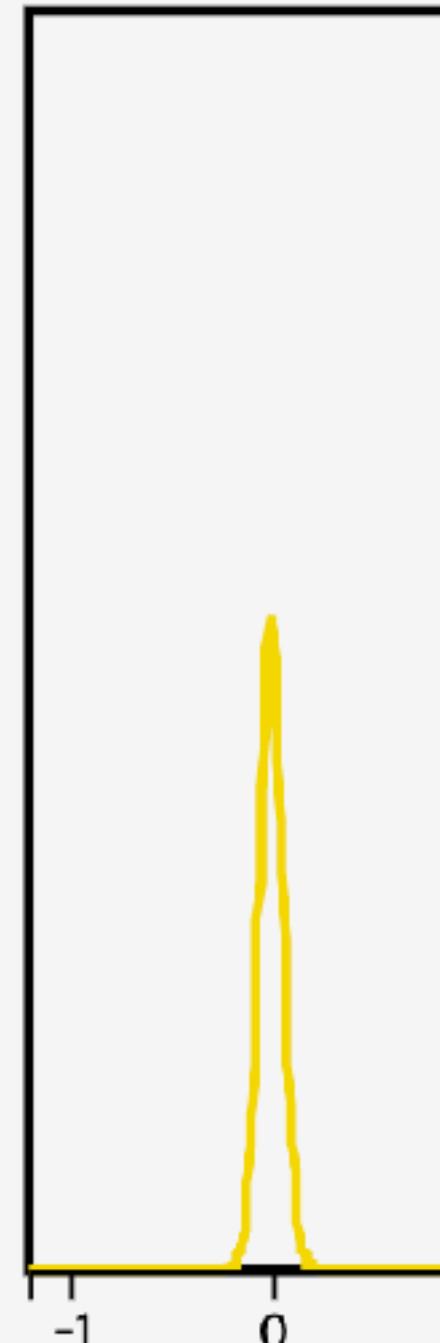
$A[1]$



$A[2]$



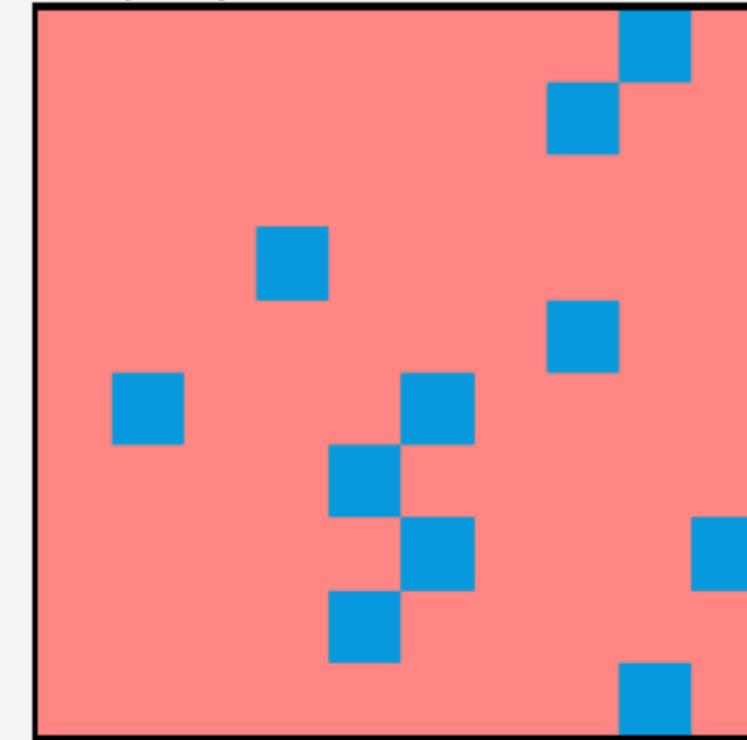
$A[3]$



$A[4]$



Output predictions of 100 images



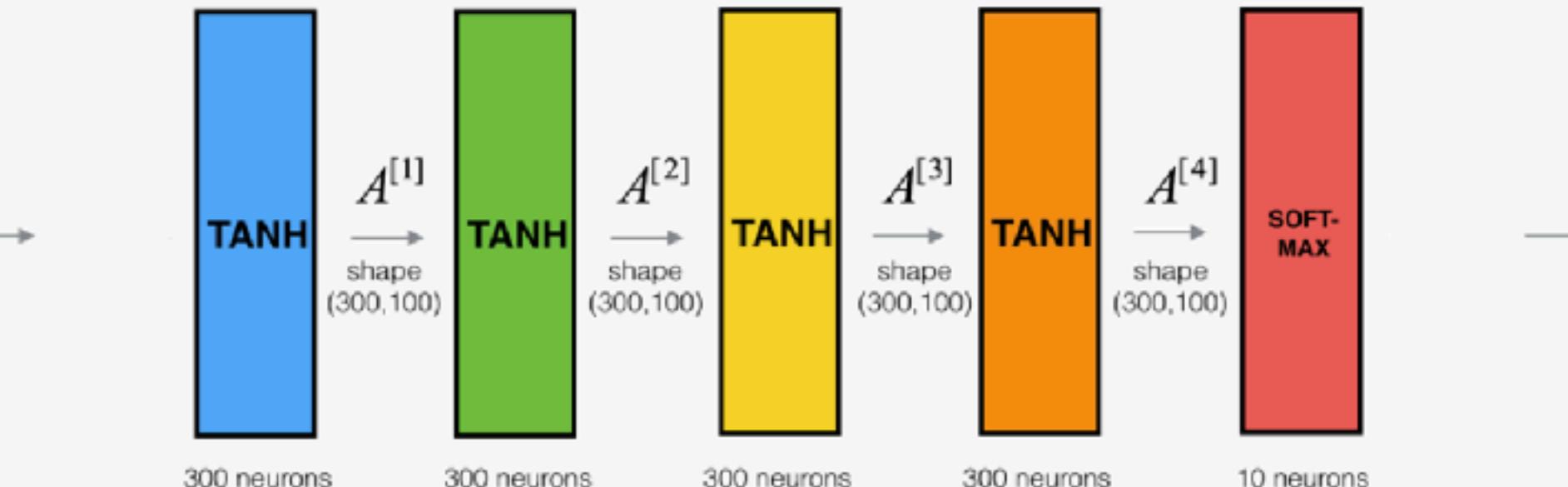
Misclassified: 89/100 Cost: 2.31

G ▶ ▶

2. Generalization initialization

$$X = A^{[0]}$$

Batch of 100 grayscale images of shape 28x28
 $X.shape = (784, 100)$ because $784 = 28 \times 28$



$$\hat{y} = A^{[5]}$$

output probability over 10
classes for a batch of
100 images
 $\hat{y}.shape = (10, 100)$

1. Load your dataset

Load 10,000 handwritten digits images
([MNIST](#)).

Load MNIST (100%)

2. Select an initialization method

Among the below distributions, select the one to use to initialize your parameters³.

- Zero Uniform Xavier Standard Normal

3. Train the network and observe

The grid below refers to the input images,
Blue squares represent correctly
classified images. Red squares represent
misclassified images.

G ▶ ▶

Input batch of 100 images

7 8 5 9 7 9 6 9 6 3
7 4 4 5 3 5 4 7 8 7
8 0 7 6 8 8 7 3 3 1
9 5 2 7 3 5 1 1 2 1
4 7 4 7 5 4 5 4 0 8
3 6 9 6 0 2 7 4 4 4
4 6 6 4 7 9 3 4 5 5
8 7 3 7 2 7 0 2 4 1
1 6 6 9 2 8 7 2 0 1
5 0 9 1 7 0 6 0 8 6

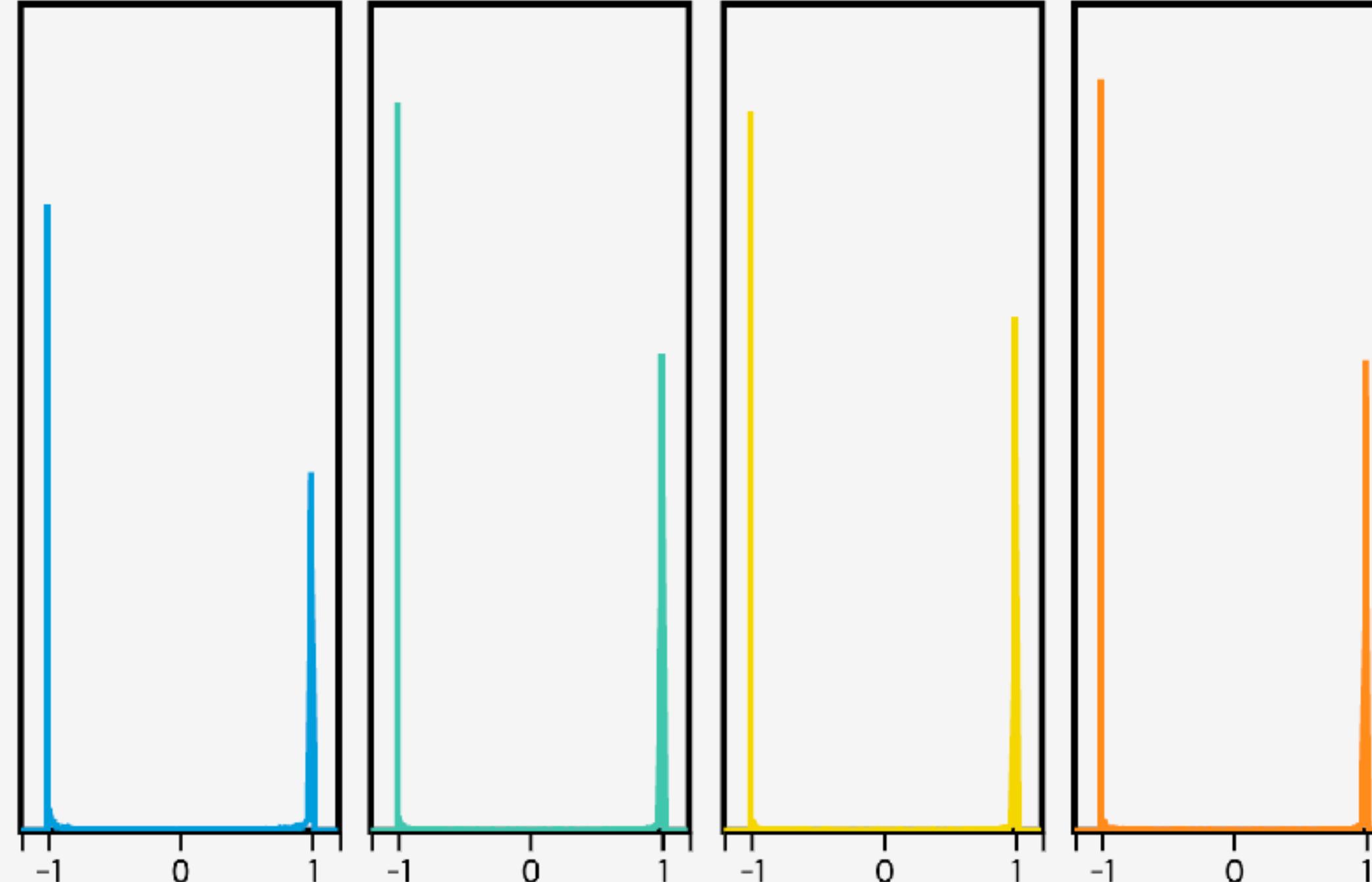
Batch: 11 Epoch: 0

$A[1]$

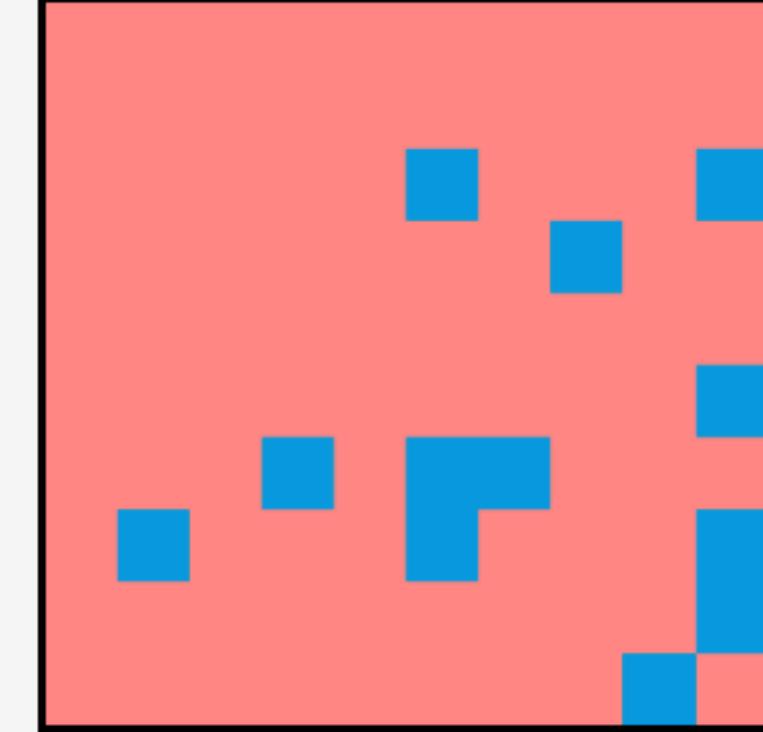
$A[2]$

$A[3]$

$A[4]$

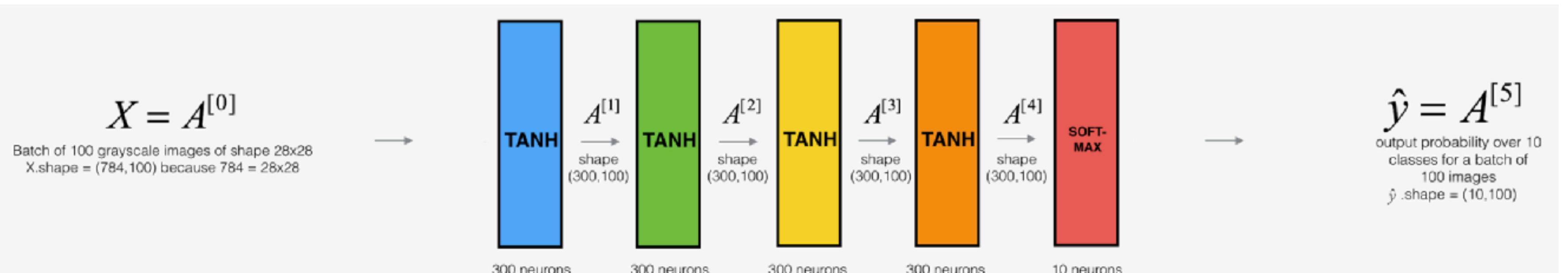


Output predictions of 100 images



Misclassified: 88/100 Cost: 26.92

2. Generalization initialization



1. Load your dataset

Load 10,000 handwritten digits images (MNIST).

Load MNIST (100%)

2. Select an initialization method

Among the below distributions, select the one to use to initialize your parameters³.

- Zero Uniform Xavier Standard Normal

3. Train the network and observe

The grid below refers to the input images, Blue squares represent correctly classified images. Red squares represent misclassified images.

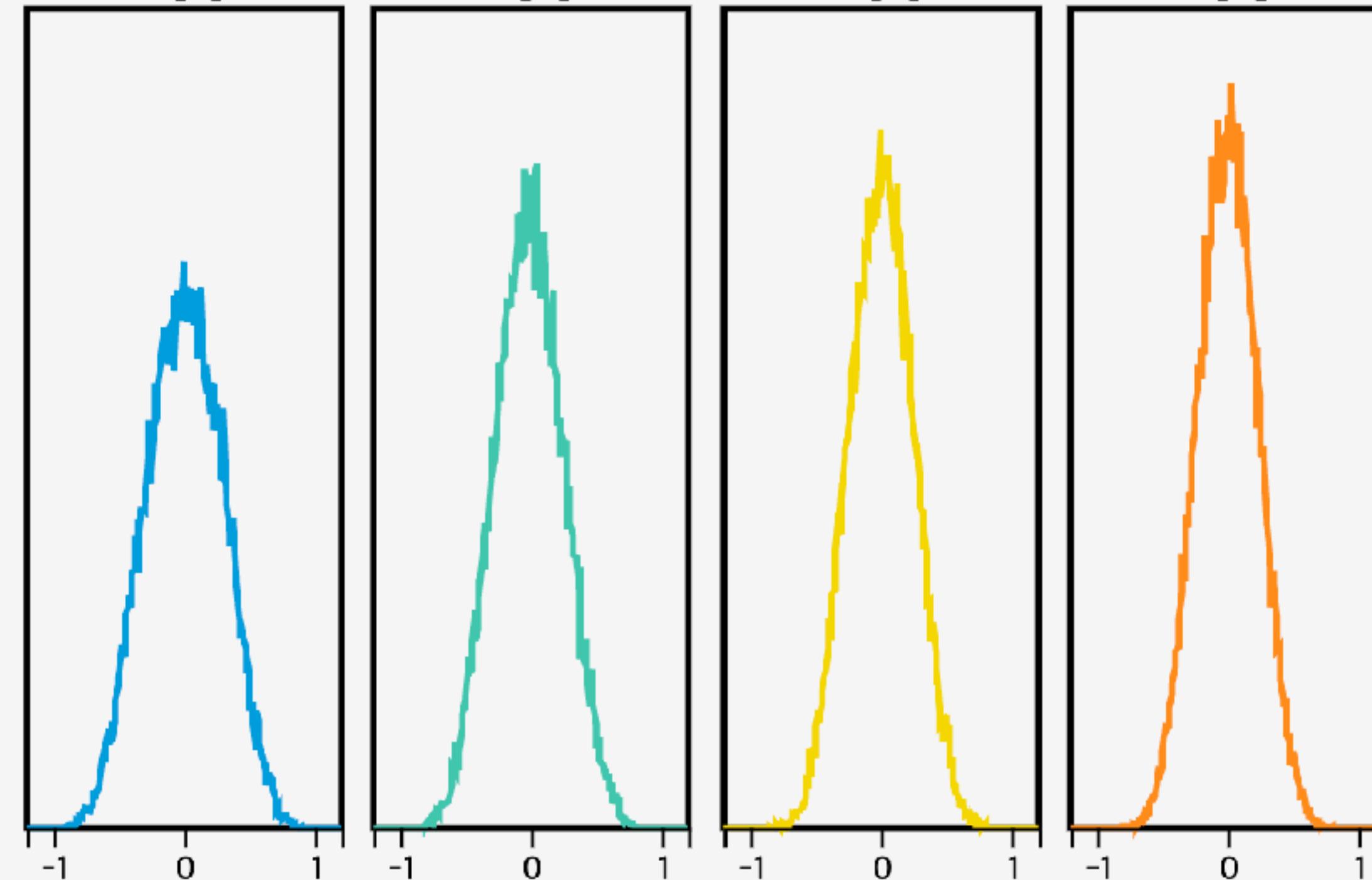


Input batch of 100 images

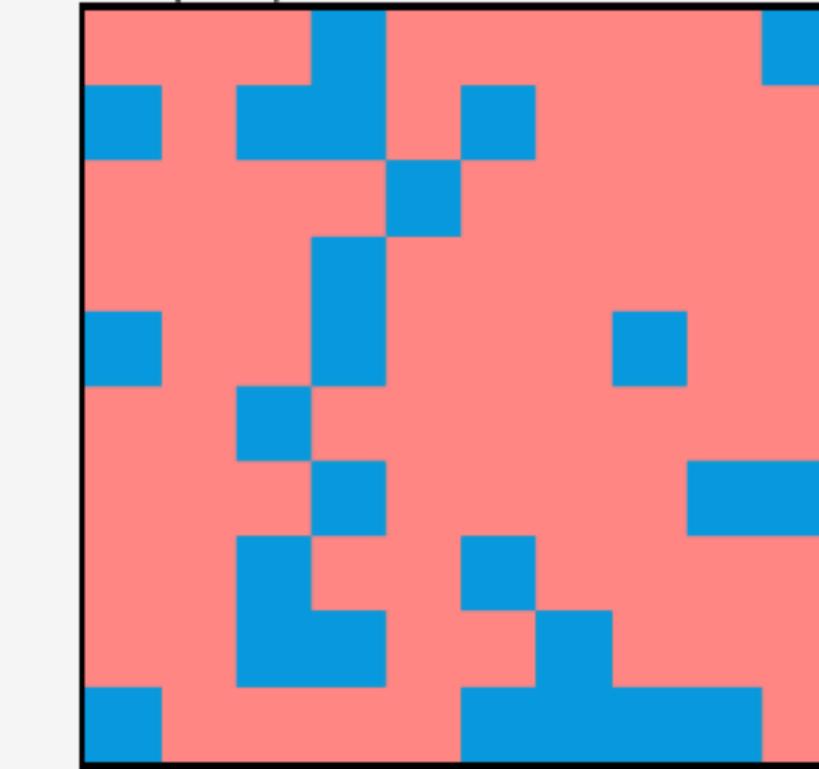
6 8 6 8 5 7 8 6 0 2
4 0 2 2 3 1 9 7 5 1
0 8 4 6 2 4 7 9 3 0
9 8 2 2 9 2 7 3 5 9
1 8 0 2 0 5 1 3 7
6 7 1 2 5 8 0 3 1 1
4 0 9 1 8 6 7 7 4 3
9 9 1 9 5 1 7 3 9 7
6 9 1 3 7 8 3 3 6 1
2 8 5 8 5 1 1 4 9 3

Batch: 6 Epoch: 0

A[1] A[2] A[3] A[4]



Output predictions of 100 images



Misclassified: 75/100 Cost: 2.22

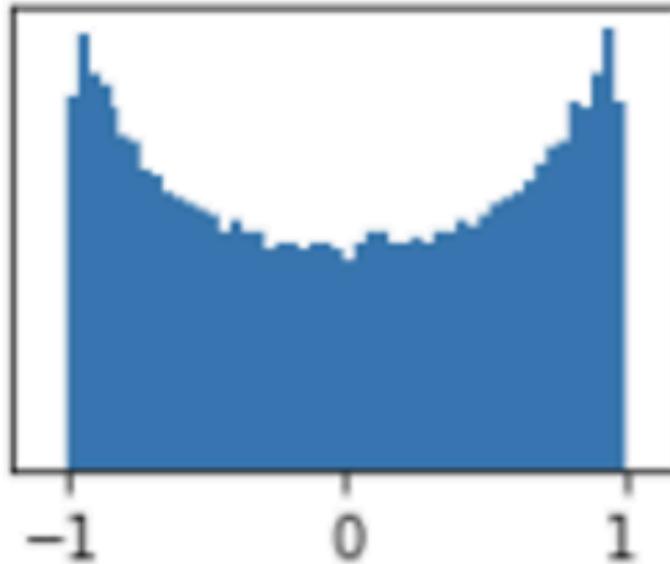
Weight Initialization: Xavier Initialization

```
dims = [4096] * 7          "Xavier" initialization:  
hs = []                      std = 1/sqrt(Din)  
x = np.random.randn(16, dims[0])  
for Din, Dout in zip(dims[:-1], dims[1:]):  
    W = np.random.randn(Din, Dout) / np.sqrt(Din)  
    x = np.tanh(x.dot(W))  
    hs.append(x)
```

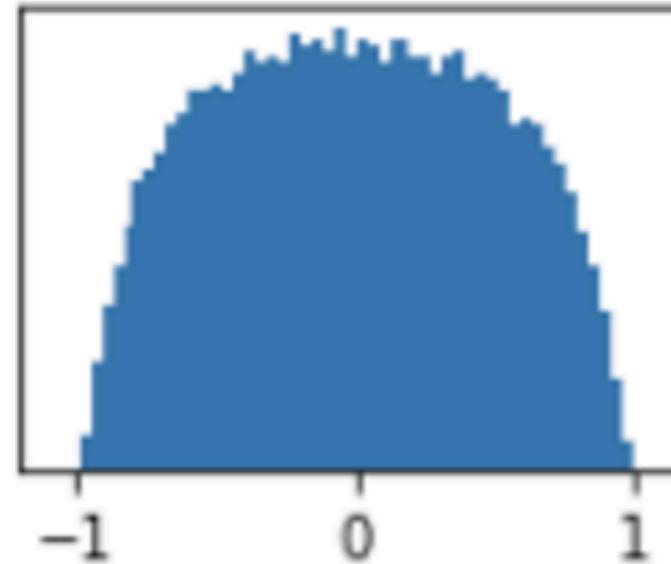
"Just right": Activations are nicely scaled for all layers!

For conv layers, Din is kernel_size² * input_channels

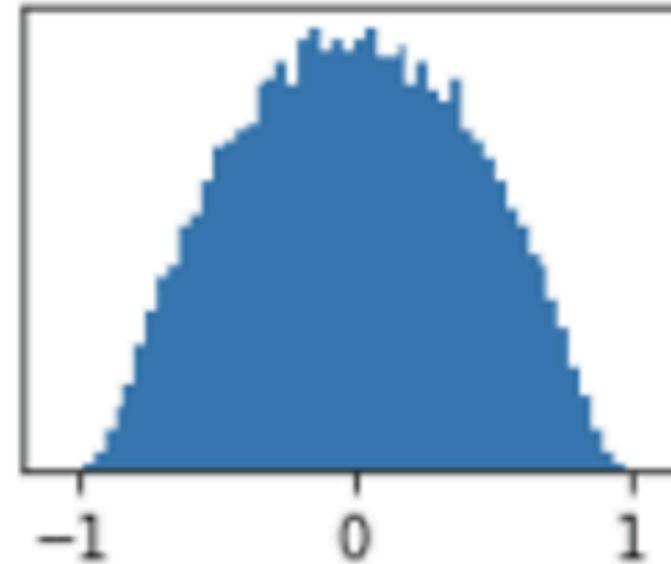
Layer 1
mean=-0.00
std=0.63



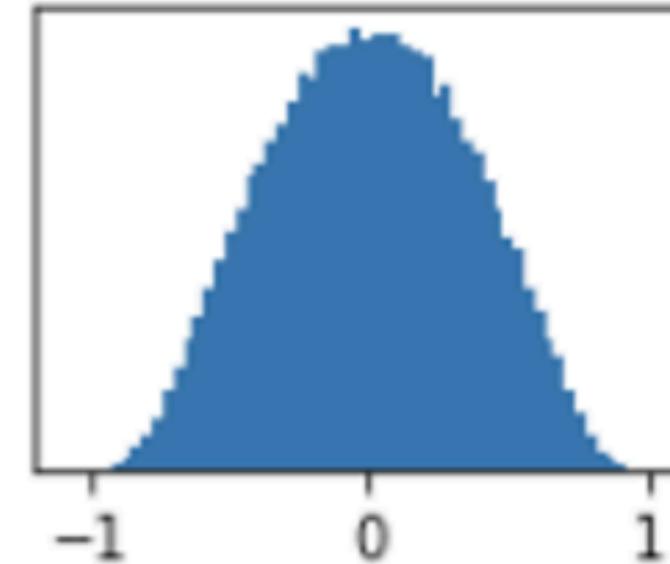
Layer 2
mean=-0.00
std=0.49



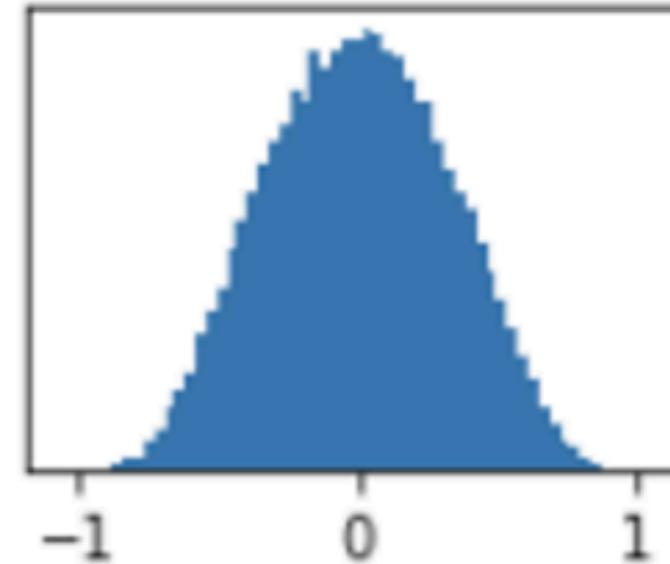
Layer 3
mean=0.00
std=0.41



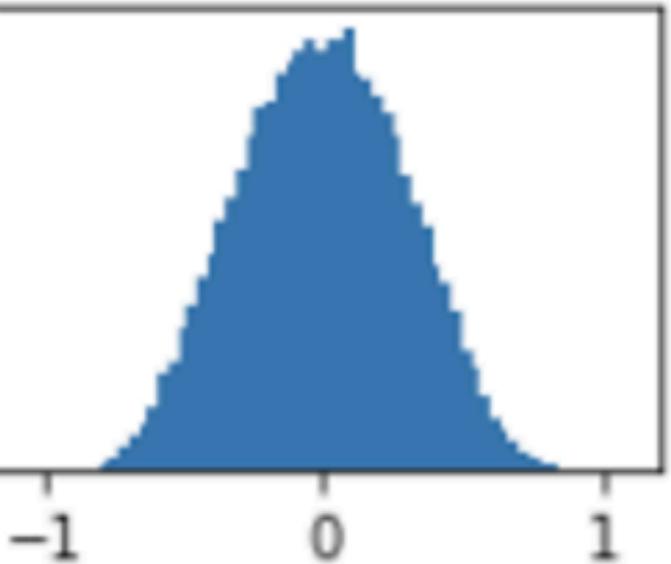
Layer 4
mean=0.00
std=0.36



Layer 5
mean=0.00
std=0.32



Layer 6
mean=-0.00
std=0.30



2. Generalization initialization

Weight Initialization: Xavier Initialization

“Xavier” initialization:
 $\text{std} = 1/\sqrt{D_{in}}$

Derivation: Variance of output = Variance of input

$$\mathbf{y} = \mathbf{W}\mathbf{x}$$

$$y_i = \sum_{j=1}^{D_{in}} x_j w_j$$

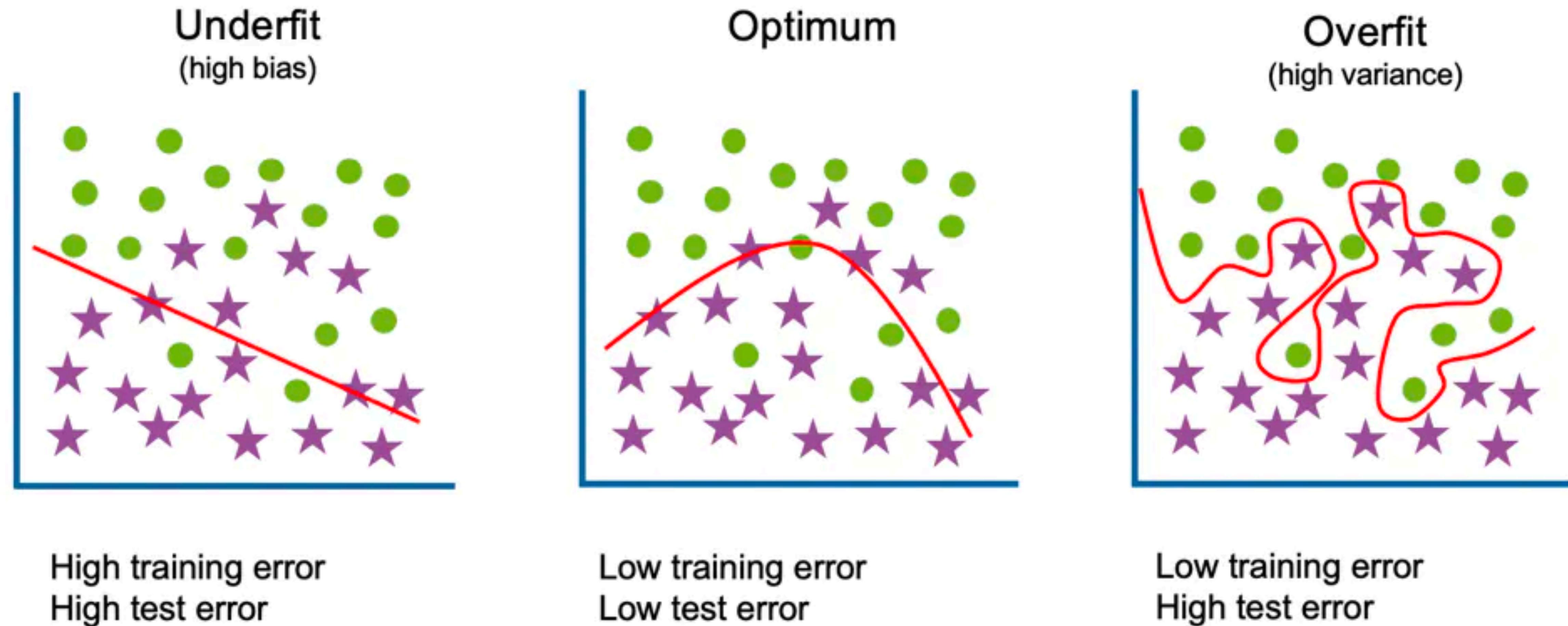
$$\text{Var}(y_i) = D_{in} * \text{Var}(x_i w_i) \quad [\text{Assume } x, w \text{ are iid}]$$

$$= D_{in} * (E[x_i^2]E[w_i^2] - E[x_i]^2 E[w_i]^2) \quad [\text{Assume } x, w \text{ independent}]$$

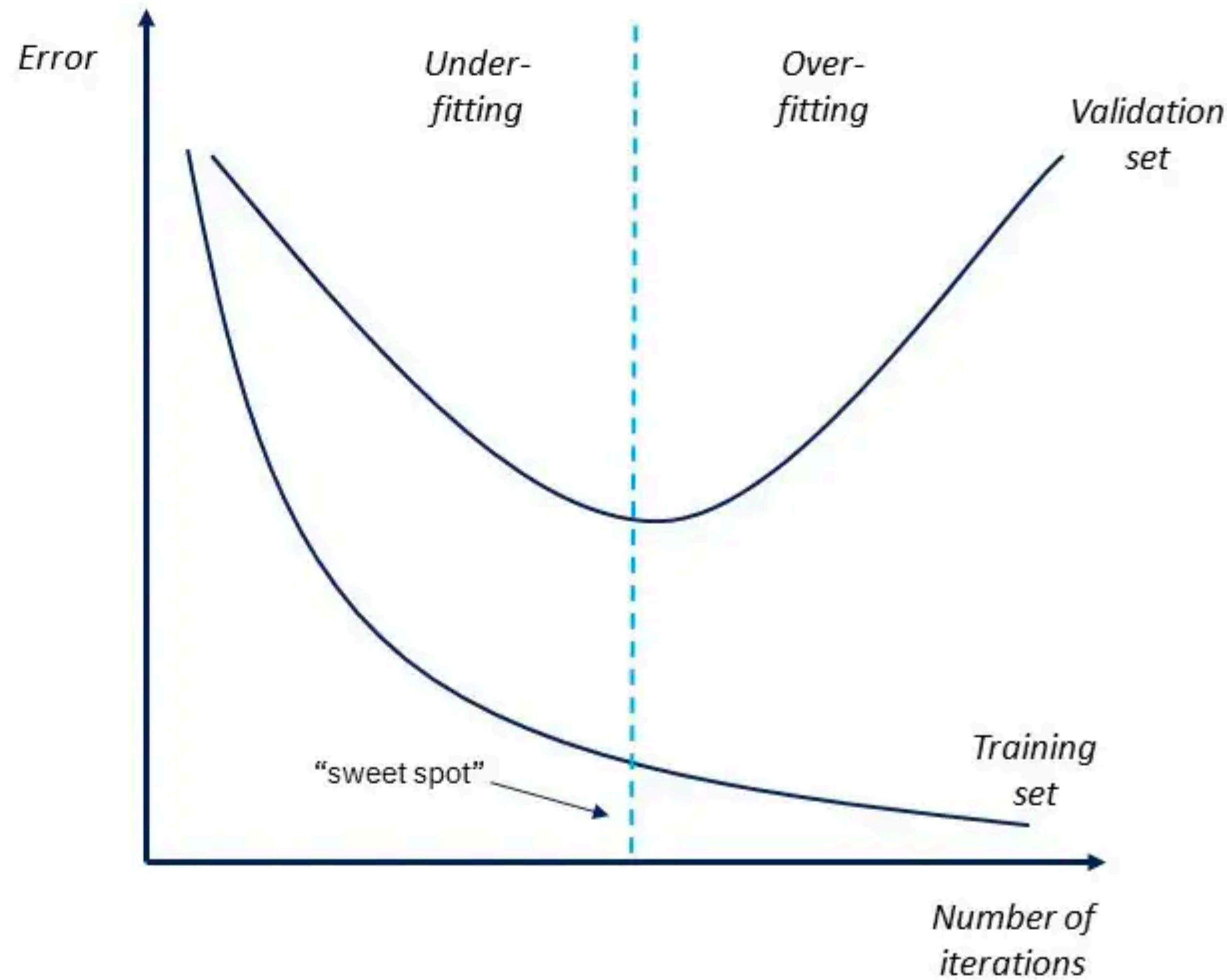
$$= D_{in} * \text{Var}(x_i) * \text{Var}(w_i) \quad [\text{Assume } x, w \text{ are zero-mean}]$$

If $\text{Var}(w_i) = 1/D_{in}$ then $\text{Var}(y_i) = \text{Var}(x_i)$

2. Generalization Overfitting



2. Generalization Overfitting



2. Generalization

Early Stopping

```
>>> callback = tf.keras.callbacks.EarlyStopping(monitor='loss', patience=3)
>>> # This callback will stop the training when there is no improvement in
>>> # the loss for three consecutive epochs.
>>> model = tf.keras.models.Sequential([tf.keras.layers.Dense(10)])
>>> model.compile(tf.keras.optimizers.SGD(), loss='mse')
>>> history = model.fit(np.arange(100).reshape(5, 20), np.zeros(5),
...                      epochs=10, batch_size=1, callbacks=[callback],
...                      verbose=0)
>>> len(history.history['loss']) # Only 4 epochs are run.
4
```

2. Generalization Regularization

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i) + \lambda R(W)$$

Data loss: Model predictions should match training data

Regularization: Prevent the model from doing *too well* on training data

λ is a hyperparameter giving regularization strength

L1 Regularization(Lasso)

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left(h_{\theta}(x^{(i)}) - y^{(i)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^n |\theta_j|$$

L2 Regularization(Ridge)

$$x = [1, 1, 1, 1] \quad J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left(h_\theta(x^{(i)}) - y^{(i)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^n \theta_j^2$$
$$w_1 = [1, 0, 0, 0]$$

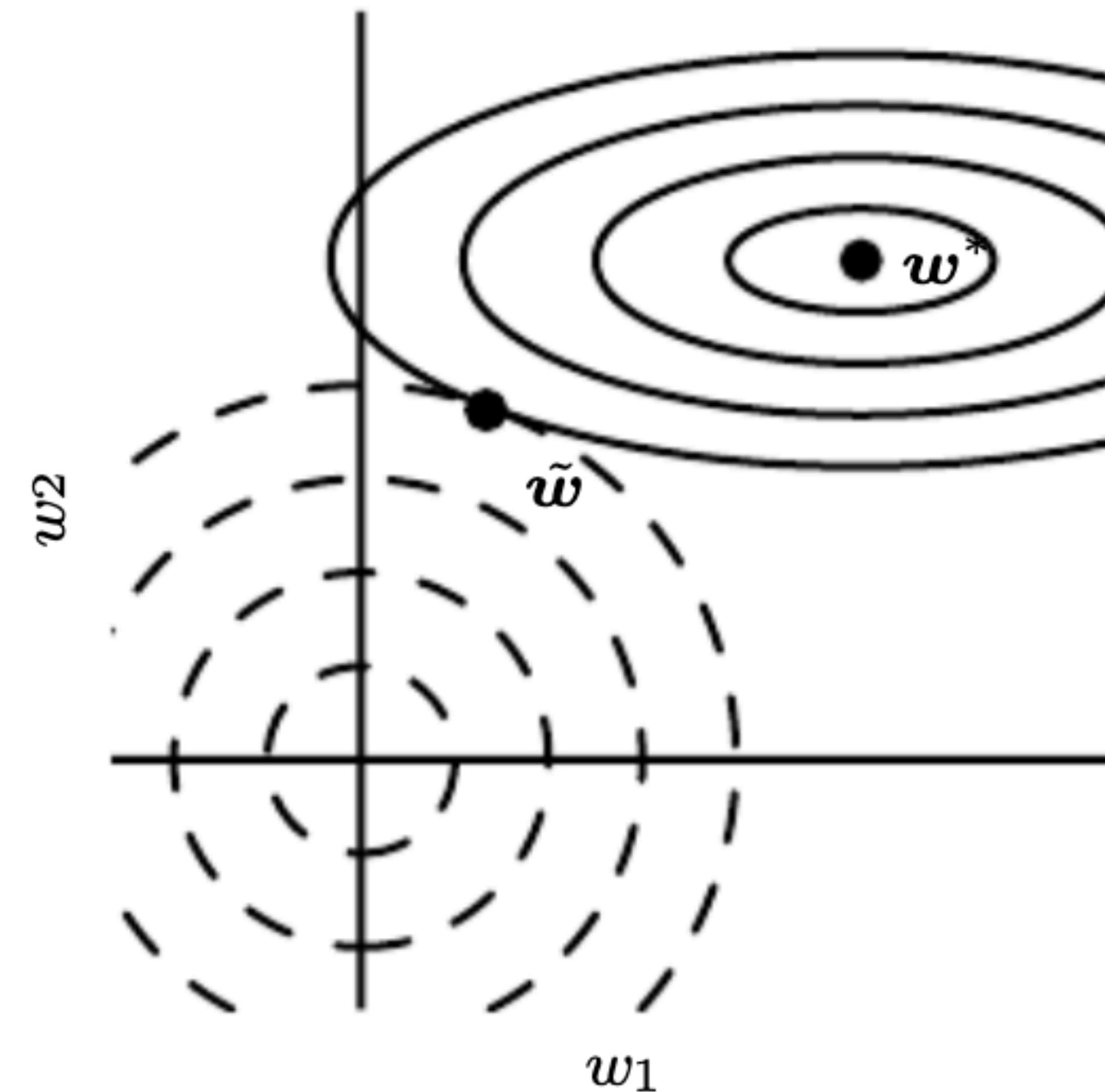
$$w_2 = [0.25, 0.25, 0.25, 0.25]$$

L2 regularization prefers weights to be “spread out”

$$w_1^T x = w_2^T x = 1$$

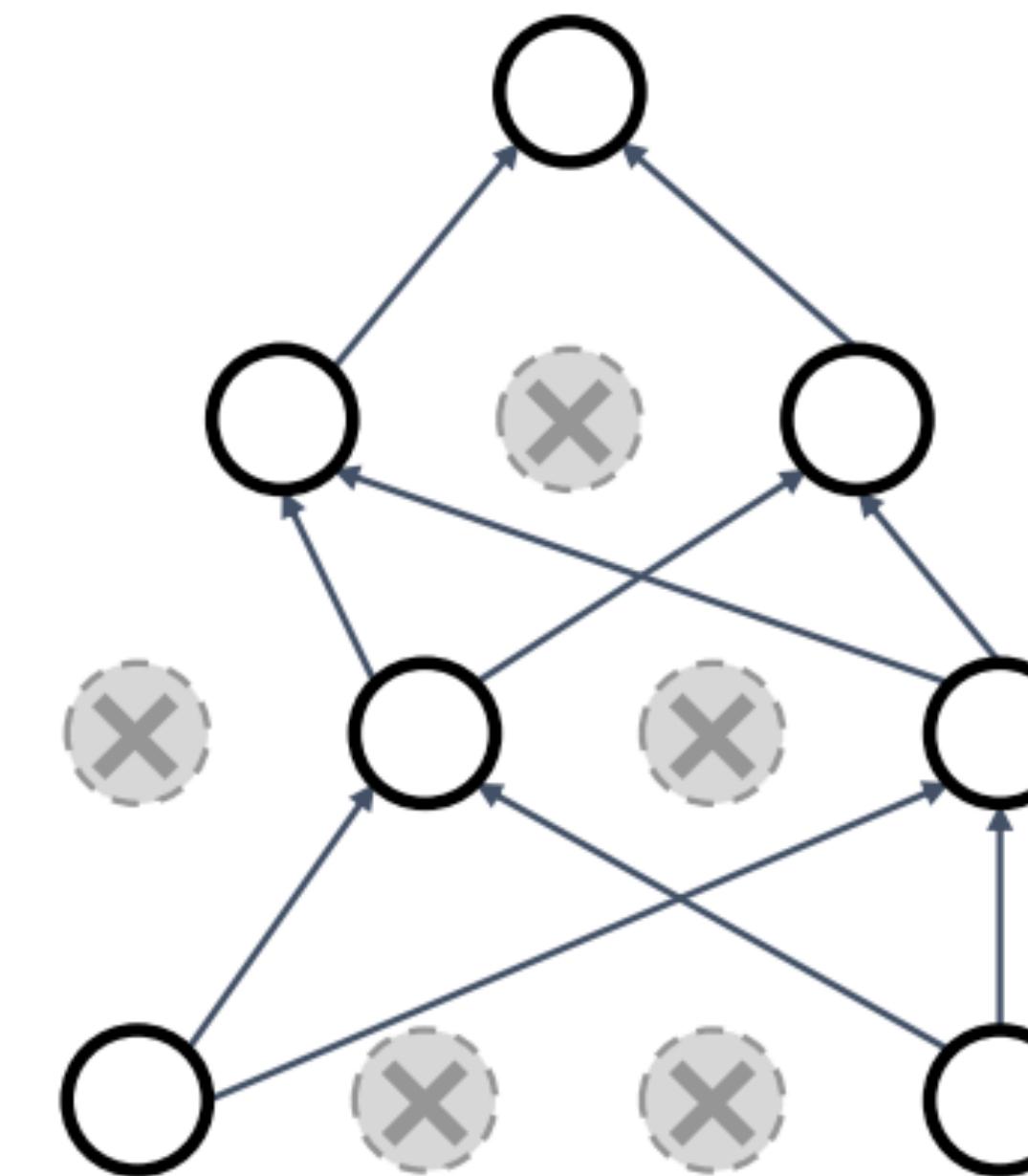
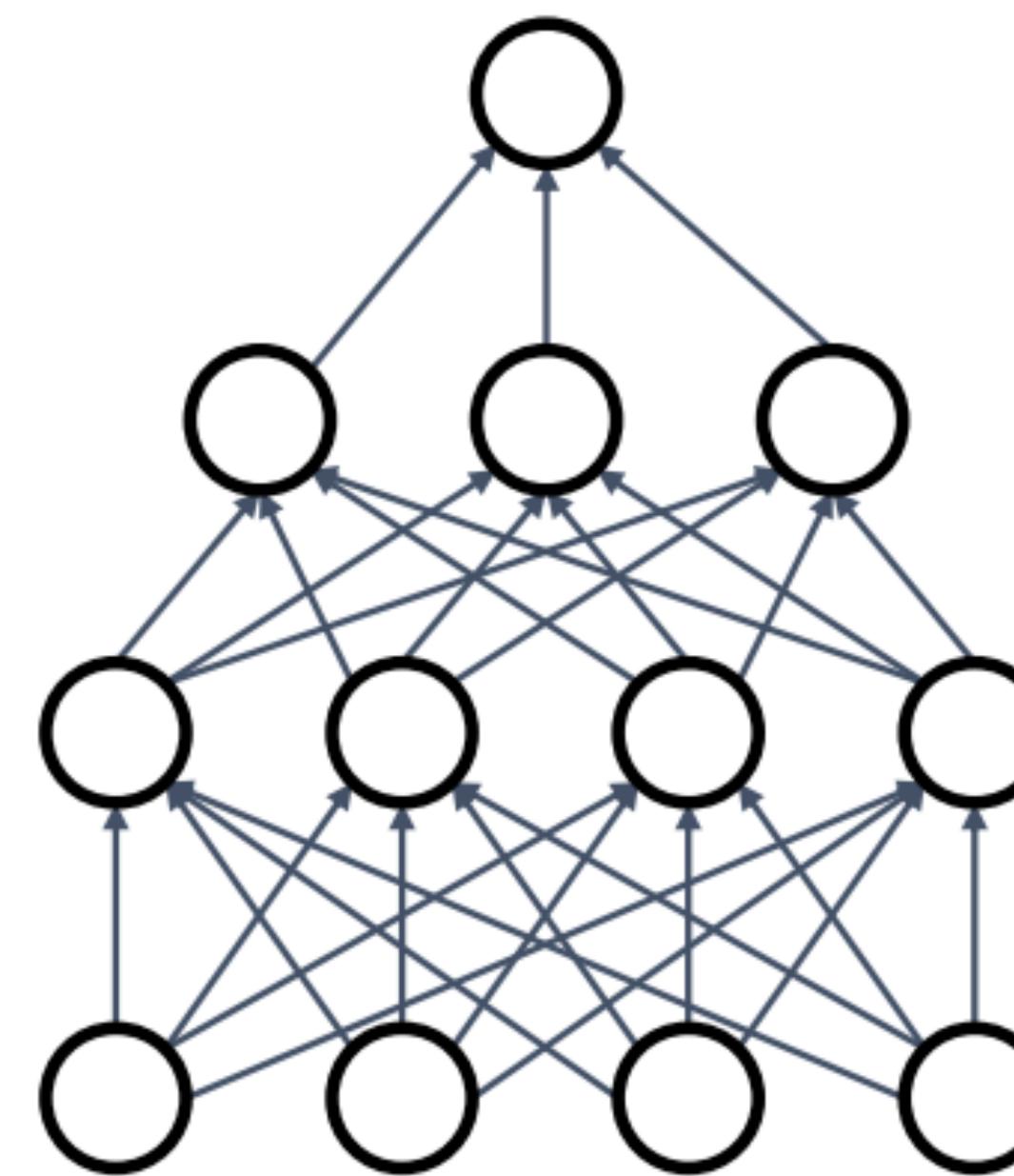
Same predictions, so data loss will always be the same

L2 Regularization(Ridge)



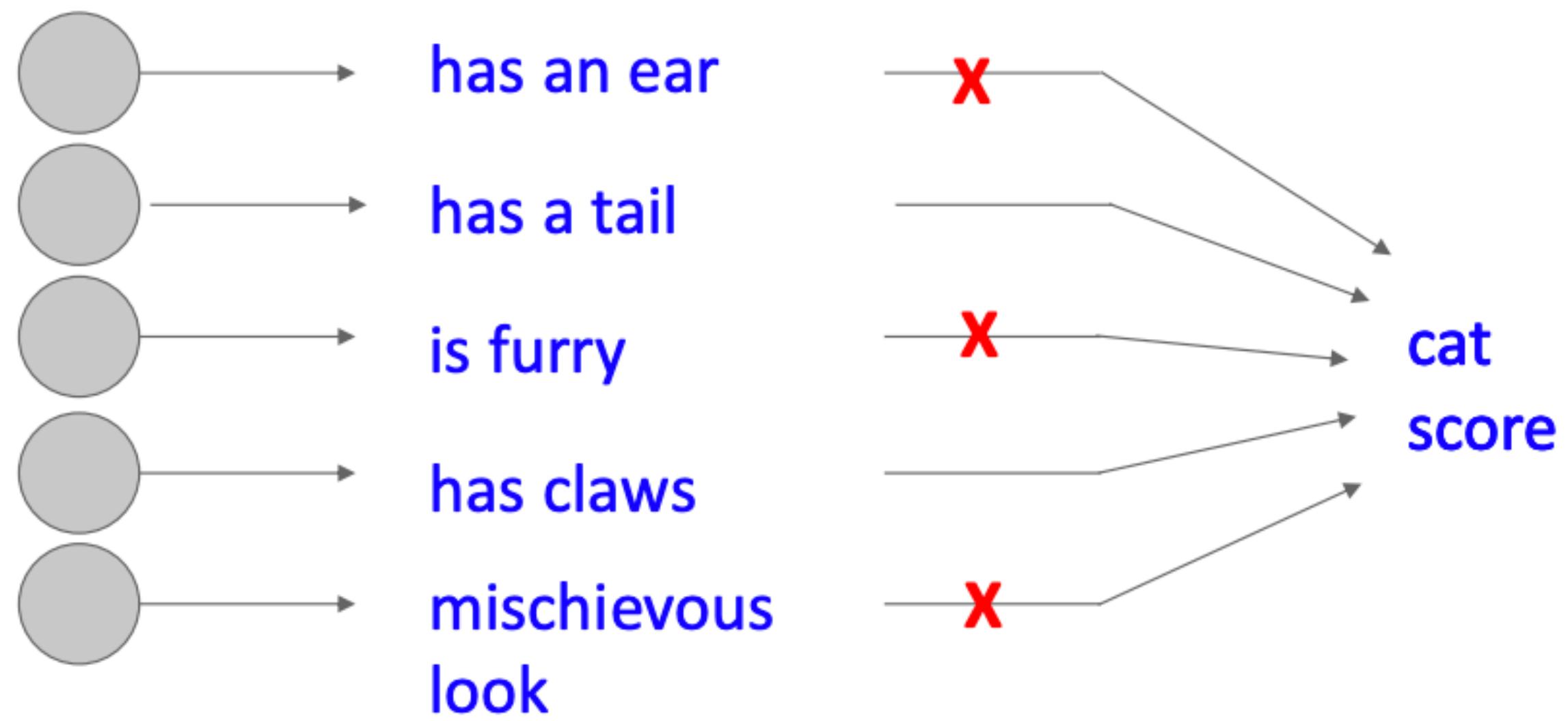
Regularization: Dropout

In each forward pass, randomly set some neurons to zero
Probability of dropping is a hyperparameter; 0.5 is common



2. Generalization

Dropout



Another interpretation:

Dropout is training a large **ensemble** of models (that share parameters).

Each binary mask is one model

Dropout: Test Time

Dropout makes our output random!

Output (label)	Input (image)	Random mask
$\mathbf{y} = f_W(\mathbf{x}, \mathbf{z})$		

Want to “average out” the randomness at test-time

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

2. Generalization

Dropout

```
def predict(X):
    # ensembled forward pass
    H1 = np.maximum(0, np.dot(W1, X) + b1) * p # NOTE: scale the activations
    H2 = np.maximum(0, np.dot(W2, H1) + b2) * p # NOTE: scale the activations
    out = np.dot(W3, H2) + b3
```

2. Generalization

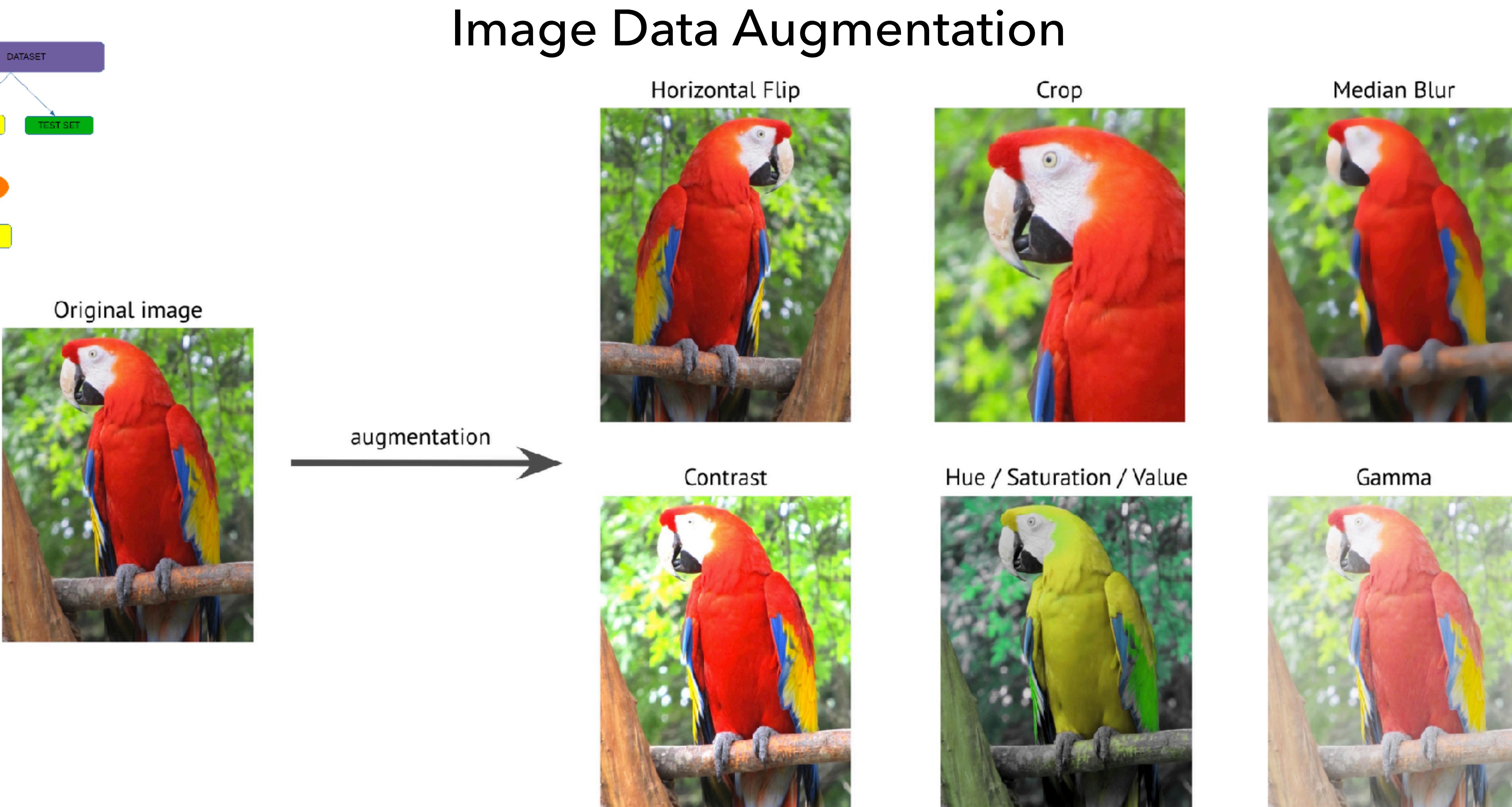
Dropout

Inverted dropout

```
def train_step(X):
    # forward pass for example 3-layer neural network
    H1 = np.maximum(0, np.dot(W1, X) + b1)
    U1 = (np.random.rand(*H1.shape) < p) / p # first dropout mask. Notice /p!
    H1 *= U1 # drop!
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    U2 = (np.random.rand(*H2.shape) < p) / p # second dropout mask. Notice /p!
    H2 *= U2 # drop!
    out = np.dot(W3, H2) + b3
```

2. Generalization

Data Augmentation



2. Generalization

Data Augmentation

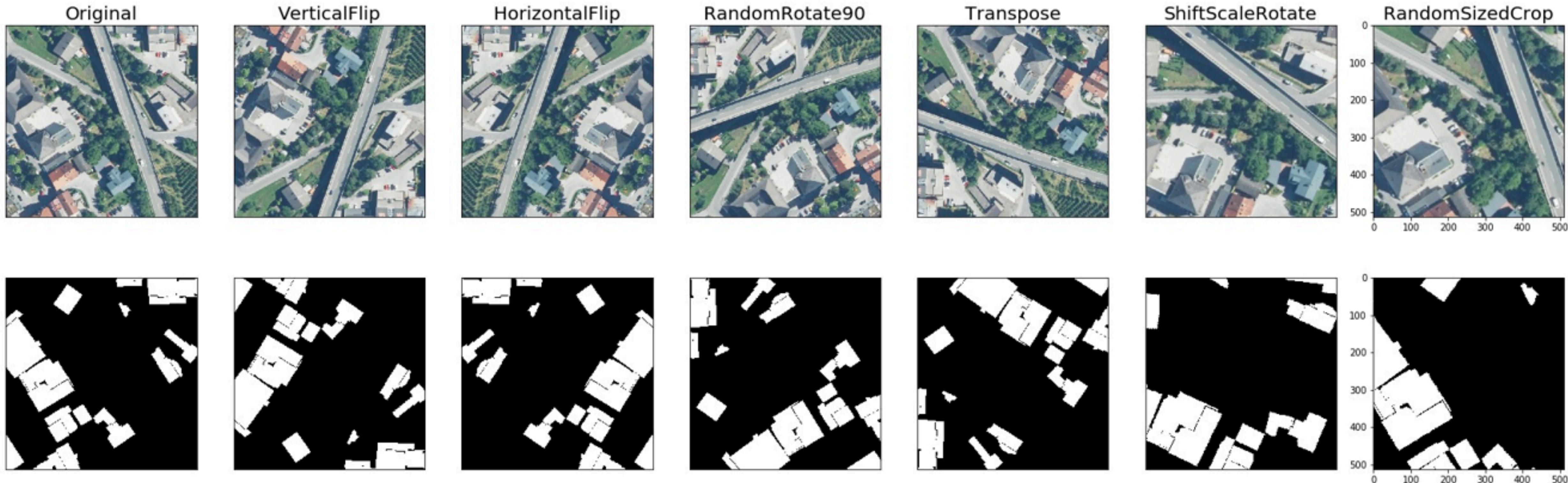
Pixel-Level Transforms



2. Generalization

Data Augmentation

Spatial-Level Transforms



PatchShuffle regularization

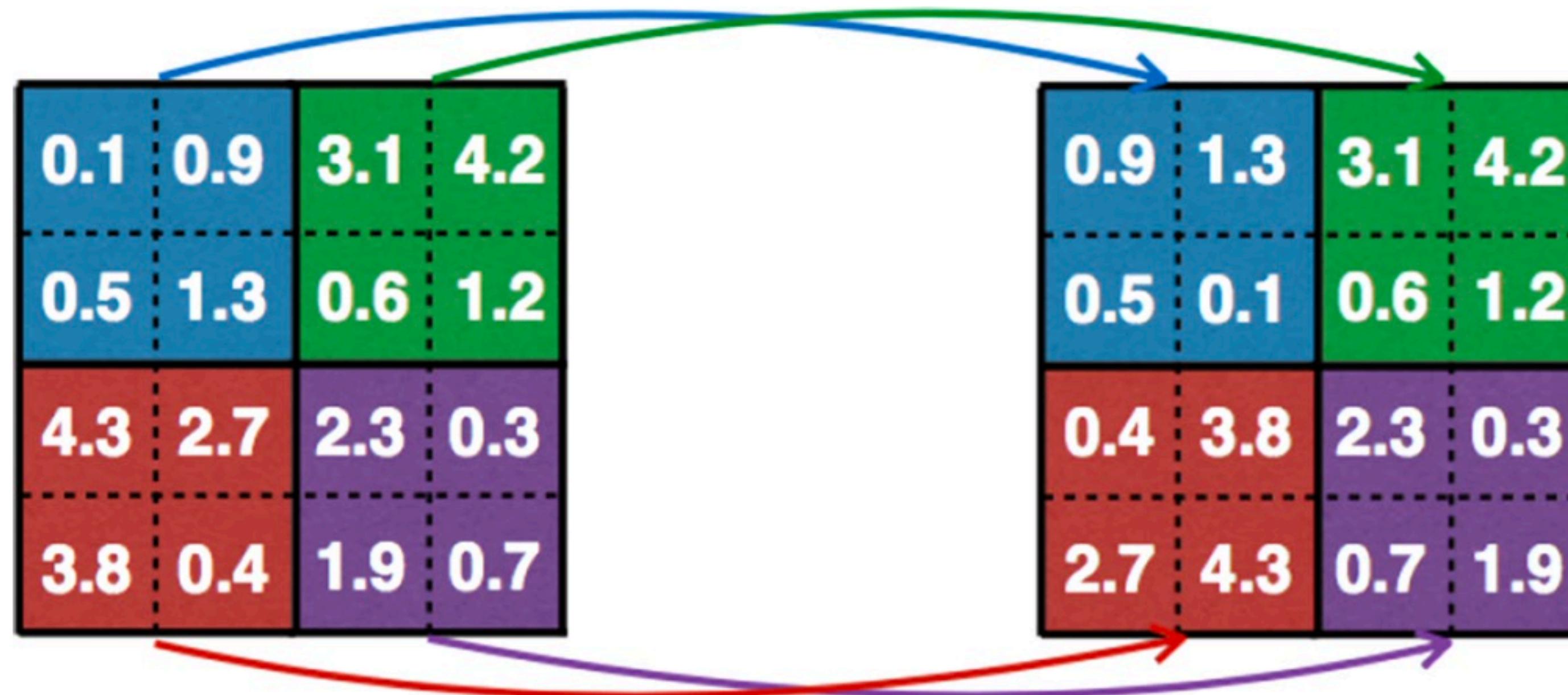


Fig. 5 Examples of applying the PatchShuffle regularization technique [64]

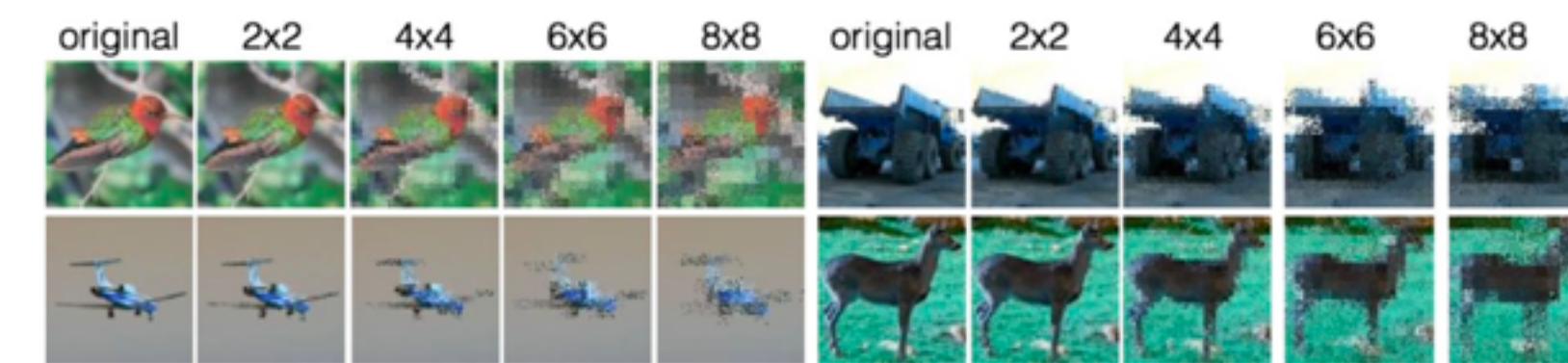


Fig. 6 Pixels in a $n \times n$ window are randomly shifted with a probability parameter p

2. Generalization

Data Augmentation

PatchShuffle regularization

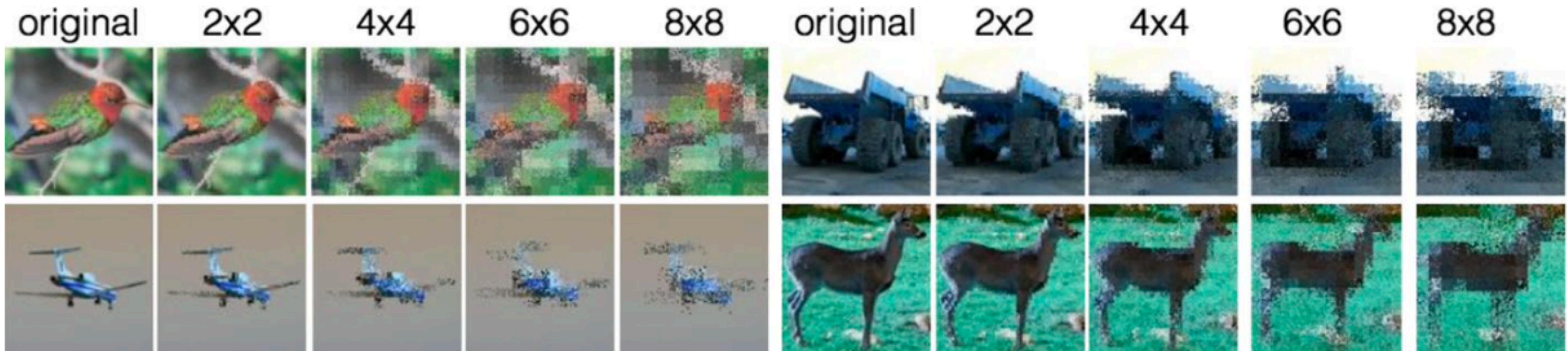
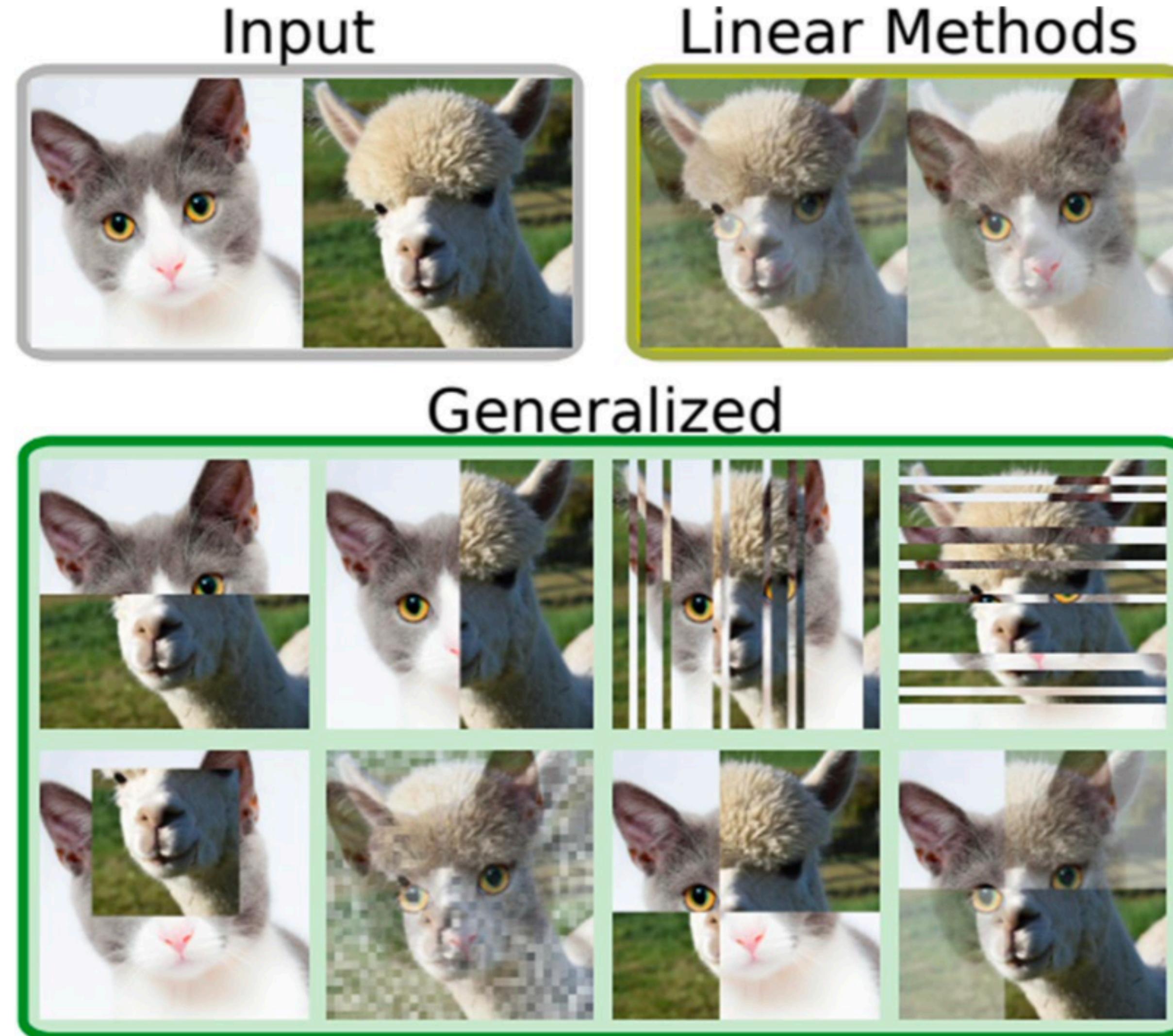


Fig. 6 Pixels in a $n \times n$ window are randomly shifted with a probability parameter p

2. Generalization

Data Augmentation

Non-linearly mixing images



Cutout



Figure 1: Cutout applied to images from the CIFAR-10 dataset.

2. Generalization

Data Augmentation

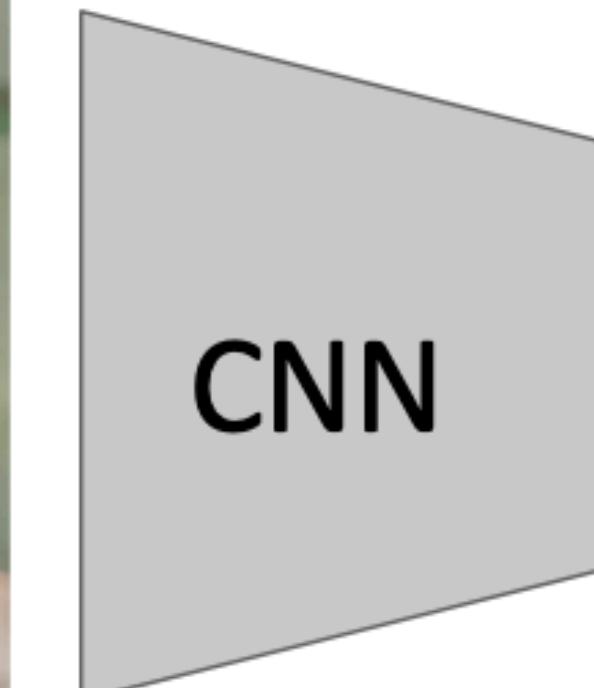
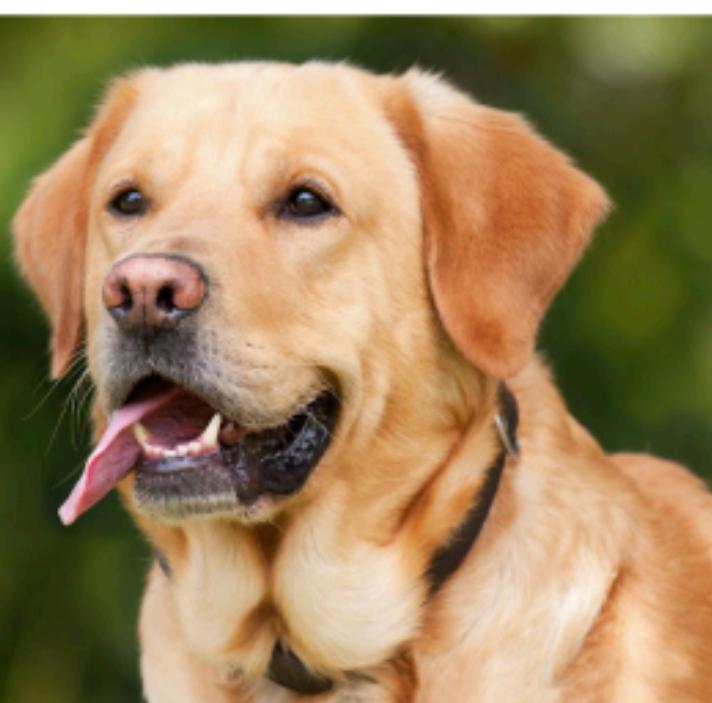
mixup

$$\tilde{x} = \lambda x_i + (1 - \lambda)x_j,$$

$$\tilde{y} = \lambda y_i + (1 - \lambda)y_j,$$

where x_i, x_j are raw input vectors

where y_i, y_j are one-hot label encodings



Target label:
cat: 0.4
dog: 0.6

Randomly blend the pixels of pairs of training images, e.g.
40% cat, 60% dog

[mixup: BEYOND EMPIRICAL RISK MINIMIZATION](#)

2. Generalization

Data Augmentation

NLP에는 Augmentation이 없을까?

2. Generalization

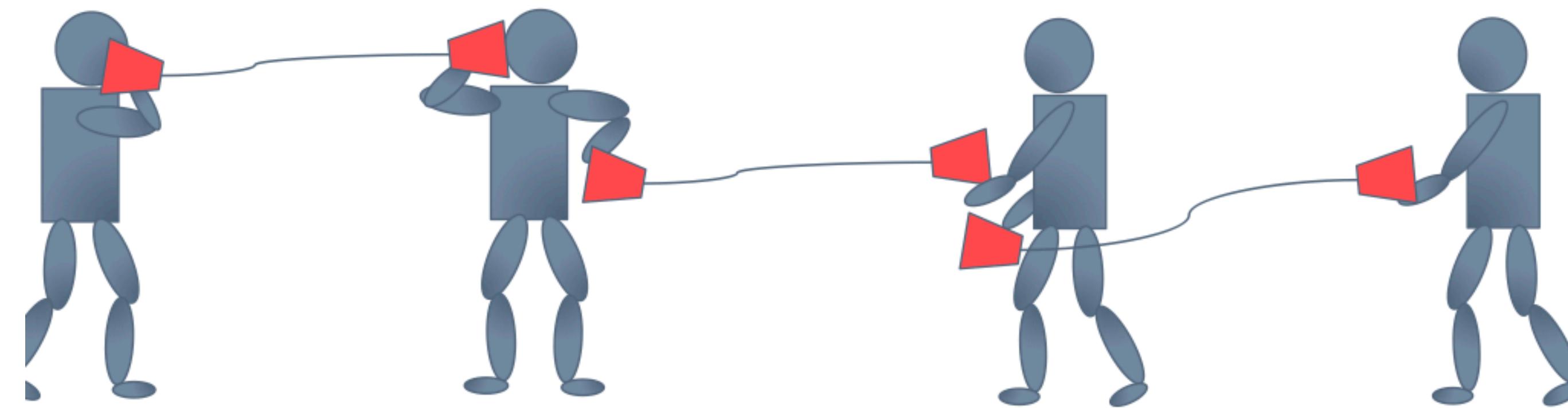
Data Augmentation

1. Back Translation
2. Easy Data Augmentation
 - Synonym Replacement
 - Random Insertion
 - Random Swap
 - Random Deletion
3. Applying Data Augmentation skill for CV to NLP

2. Generalization

Batch Normalization

Internal Covariate Shift

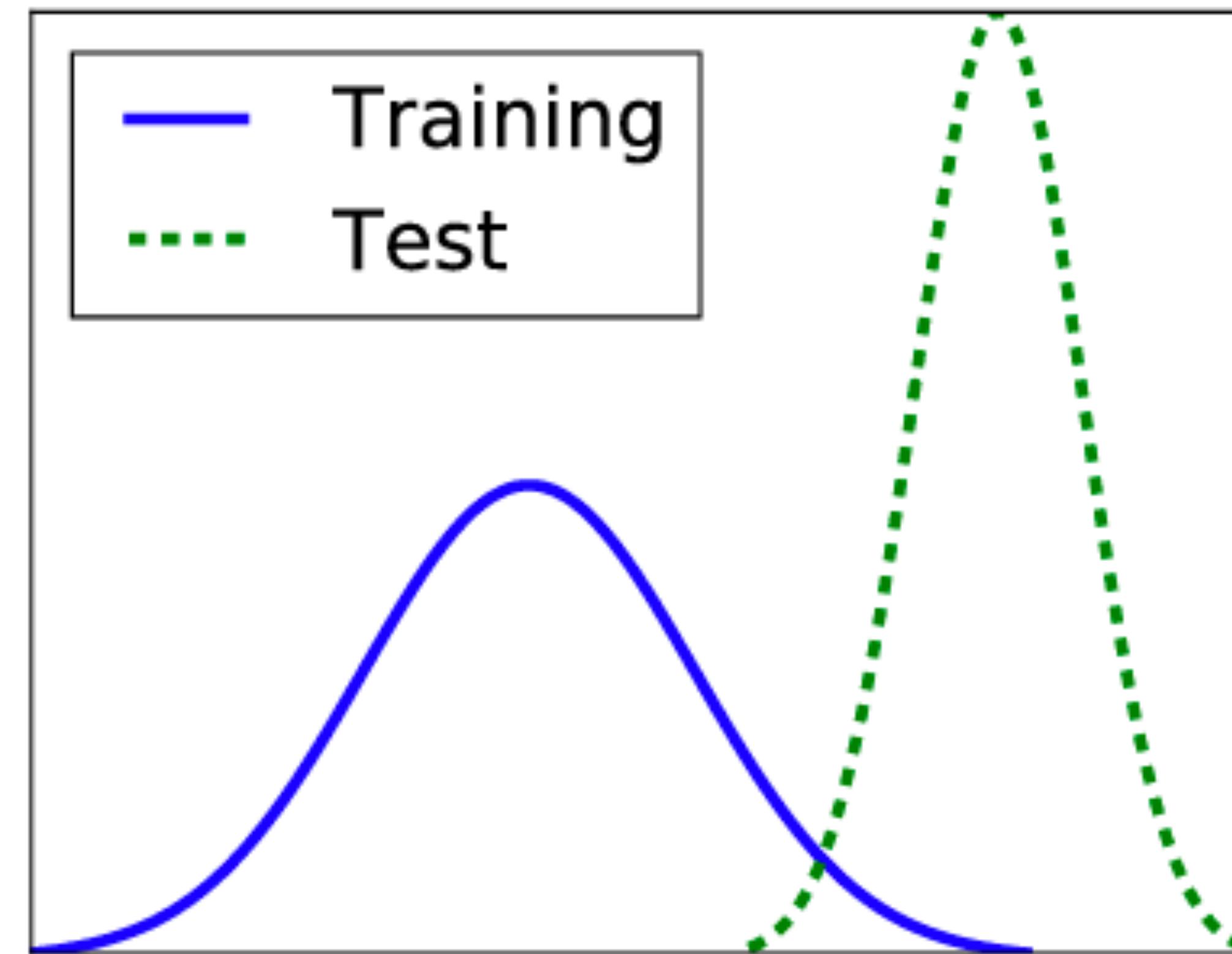


[Batch Normalization – What the hey?](#)

2. Generalization

Batch Normalization

Covariate Shift

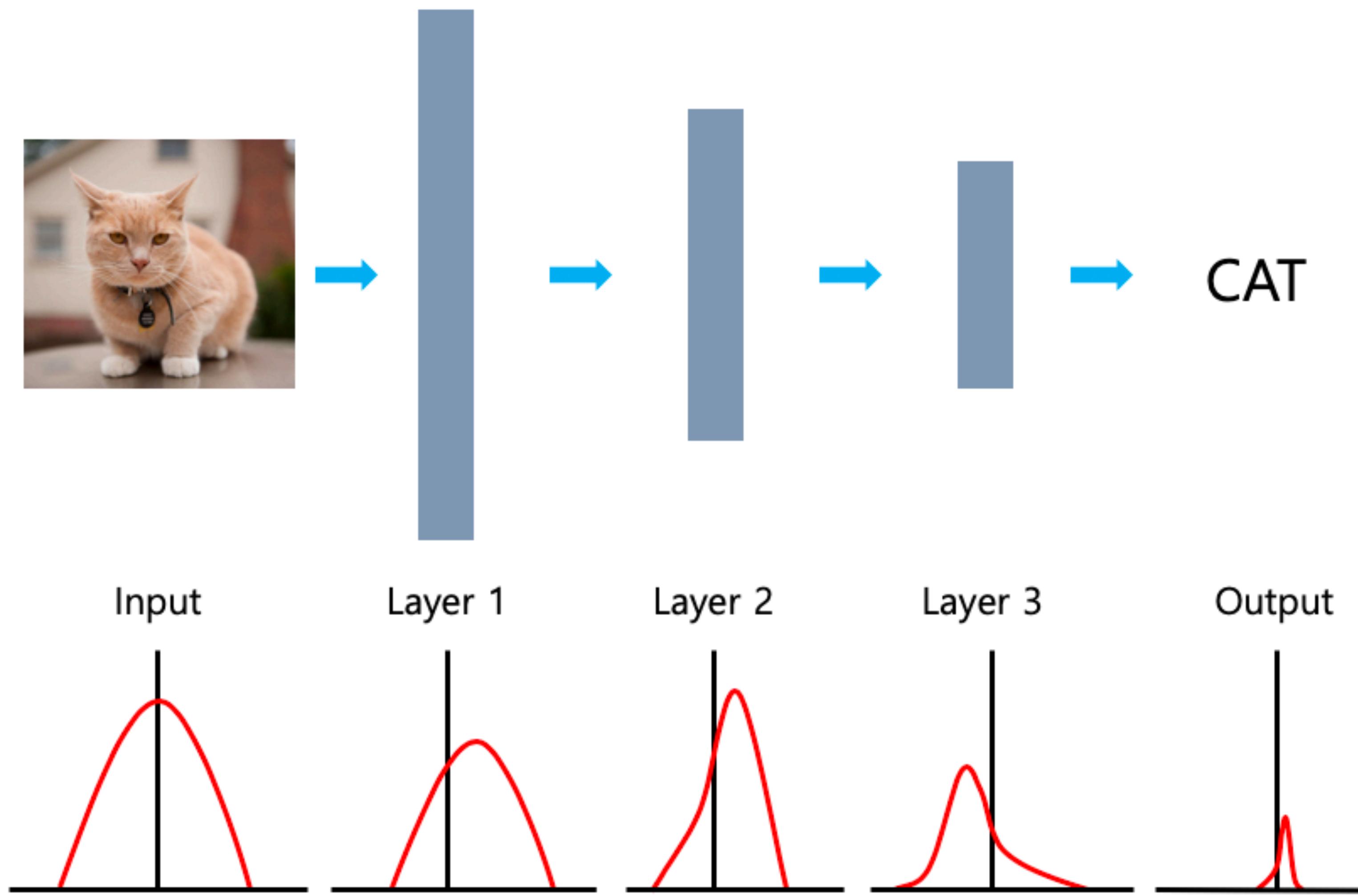


On the Reduction of Biases in Big Data Sets for the Detection of Irregular Power Usage

2. Generalization

Batch Normalization

Internal Covariate Shift



2. Generalization

Batch Normalization

Normalization via Mini-Batch Statistics

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots m\}$;

Parameters to be learned: γ, β

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.

Input: Network N with trainable parameters Θ ;
subset of activations $\{x^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference, $N_{\text{BN}}^{\text{inf}}$

- 1: $N_{\text{BN}}^{\text{tr}} \leftarrow N$ // Training BN network
- 2: **for** $k = 1 \dots K$ **do**
- 3: Add transformation $y^{(k)} = \text{BN}_{\gamma^{(k)}, \beta^{(k)}}(x^{(k)})$ to $N_{\text{BN}}^{\text{tr}}$ (Alg. 1)
- 4: Modify each layer in $N_{\text{BN}}^{\text{tr}}$ with input $x^{(k)}$ to take $y^{(k)}$ instead
- 5: **end for**
- 6: Train $N_{\text{BN}}^{\text{tr}}$ to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
- 7: $N_{\text{BN}}^{\text{inf}} \leftarrow N_{\text{BN}}^{\text{tr}}$ // Inference BN network with frozen
// parameters
- 8: **for** $k = 1 \dots K$ **do**
- 9: // For clarity, $x \equiv x^{(k)}, \gamma \equiv \gamma^{(k)}, \mu_{\mathcal{B}} \equiv \mu_{\mathcal{B}}^{(k)}$, etc.
- 10: Process multiple training mini-batches \mathcal{B} , each of size m , and average over them:

$$\text{E}[x] \leftarrow \text{E}_{\mathcal{B}}[\mu_{\mathcal{B}}]$$

$$\text{Var}[x] \leftarrow \frac{m}{m-1} \text{E}_{\mathcal{B}}[\sigma_{\mathcal{B}}^2]$$
- 11: In $N_{\text{BN}}^{\text{inf}}$, replace the transform $y = \text{BN}_{\gamma, \beta}(x)$ with

$$y = \frac{\gamma}{\sqrt{\text{Var}[x] + \epsilon}} \cdot x + \left(\beta - \frac{\gamma \text{E}[x]}{\sqrt{\text{Var}[x] + \epsilon}}\right)$$
- 12: **end for**

Algorithm 2: Training a Batch-Normalized Network

Batch Normalization

2. Generalization

Batch Normalization

tf.keras.layers.BatchNormalization

During inference (i.e. when using `evaluate()` or `predict()` or when calling the layer/model with the argument `training=False` (which is the default), the layer normalizes its output using a moving average of the mean and standard deviation of the batches it has seen during training. That is to say, it returns `gamma * (batch - self.moving_mean) / sqrt(self.moving_var + epsilon) + beta`.

`self.moving_mean` and `self.moving_var` are non-trainable variables that are updated each time the layer is called in training mode, as such:

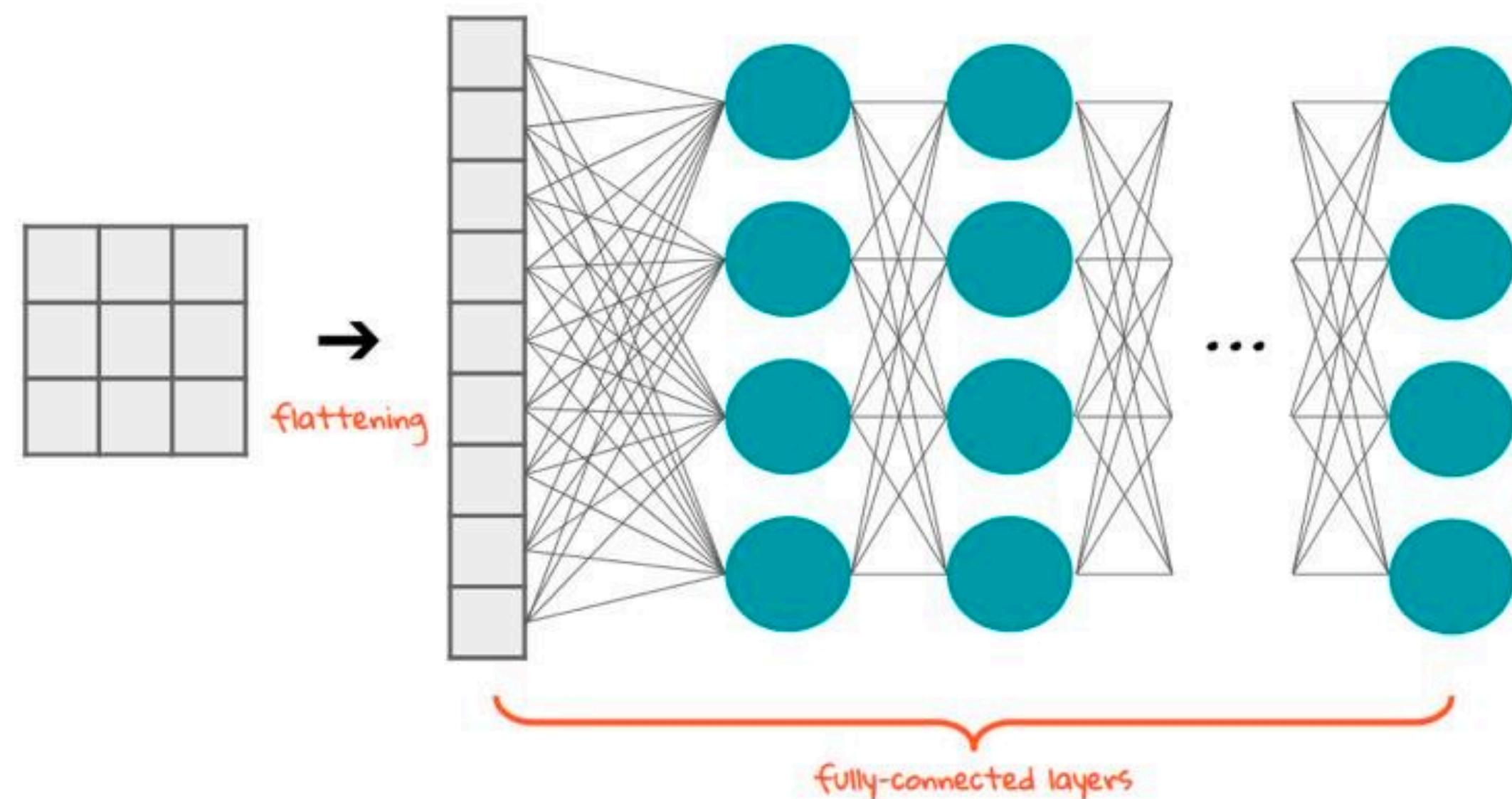
- `moving_mean = moving_mean * momentum + mean(batch) * (1 - momentum)`
- `moving_var = moving_var * momentum + var(batch) * (1 - momentum)`

3. DNN with TF Fashion MNIST

실습

4. CNN

기초 이론



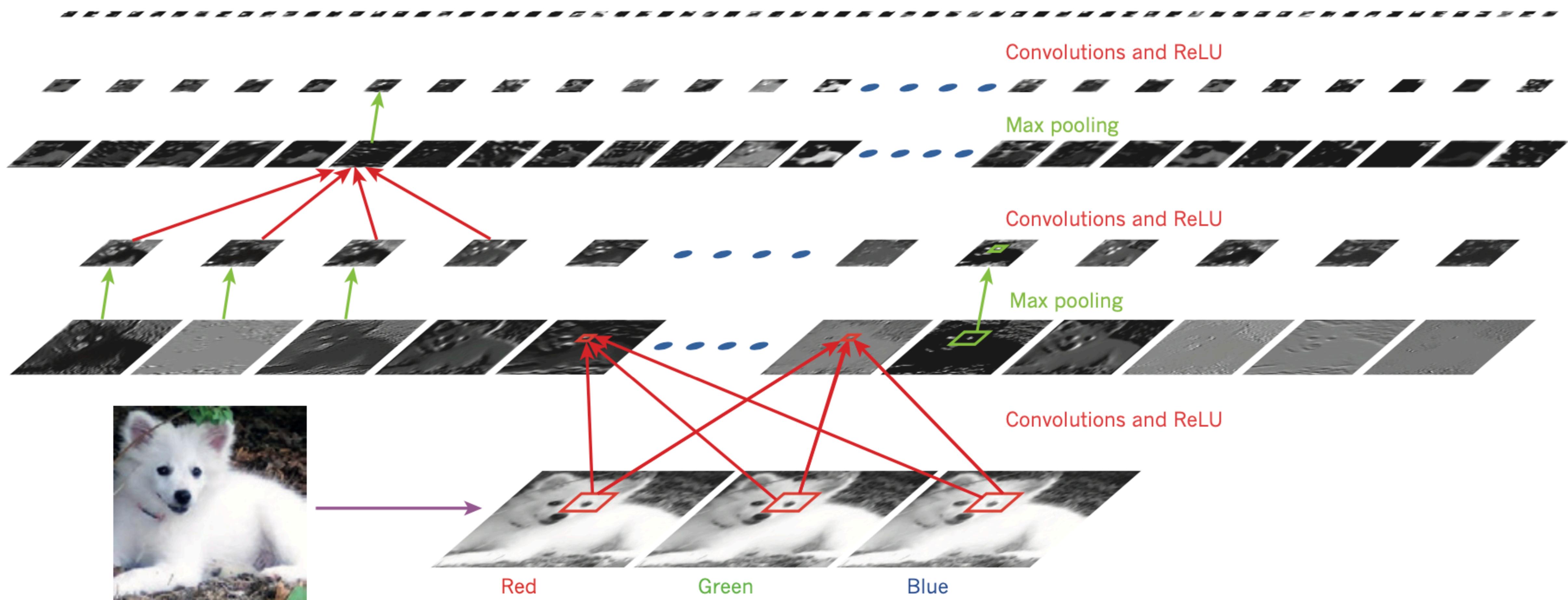
1. Image의 공간적 특성을 학습에 이용하지 못 함
2. data의 크기가 커짐에 따라 parameters의 수가 매우 커짐

Motivation

1. Sparse interactions
2. Parameter sharing
3. Equivariant representations

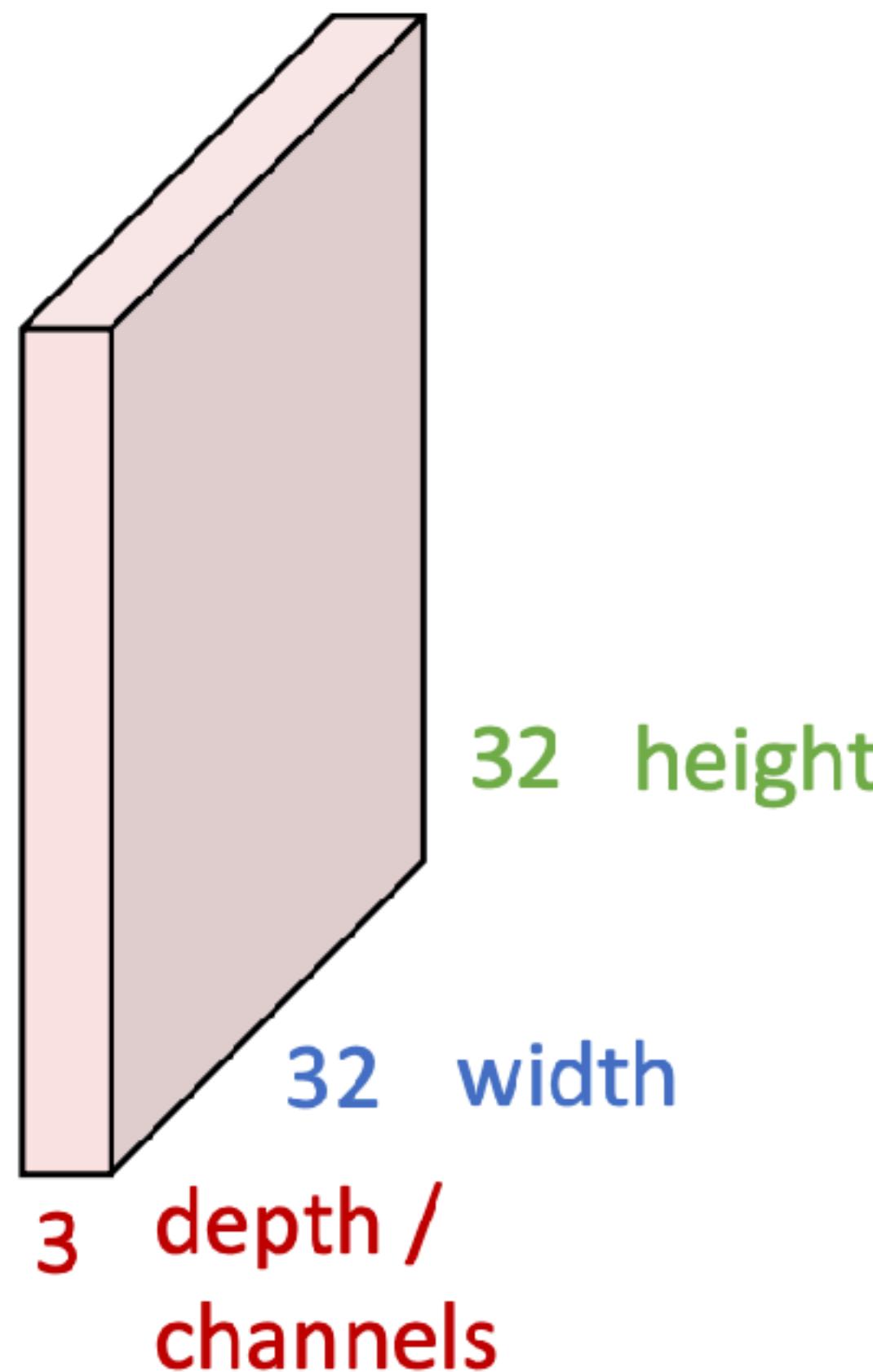
4. CNN 기초 이론

Samoyed (16); Papillon (5.7); Pomeranian (2.7); Arctic fox (1.0); Eskimo dog (0.6); white wolf (0.4); Siberian husky (0.4)



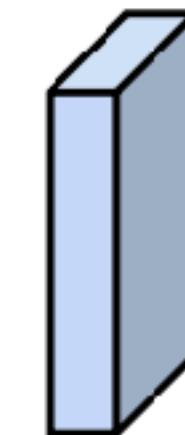
Convolution Layer

$3 \times 32 \times 32$ image



Filters always extend the full depth of the input volume

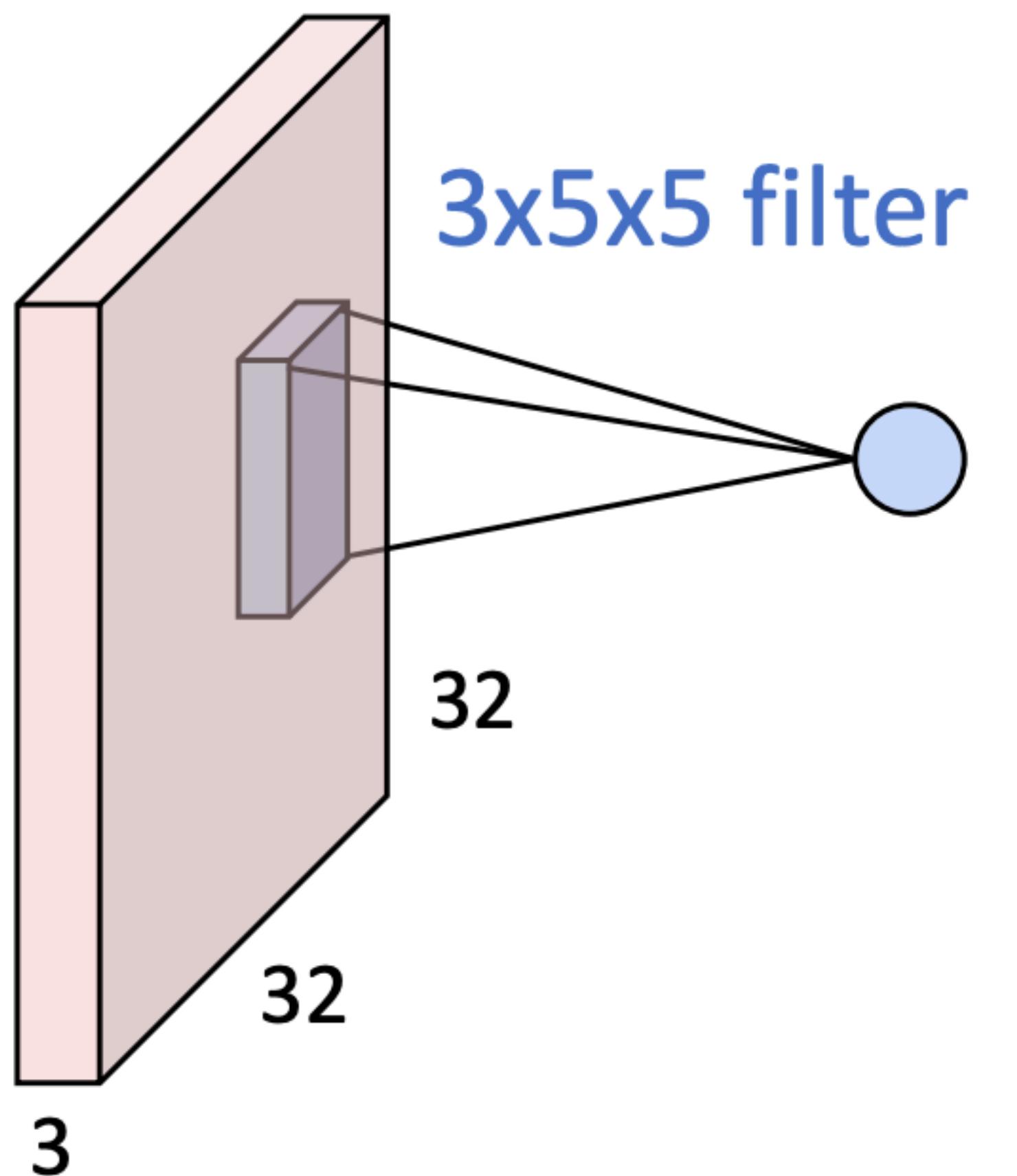
$3 \times 5 \times 5$ filter



Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

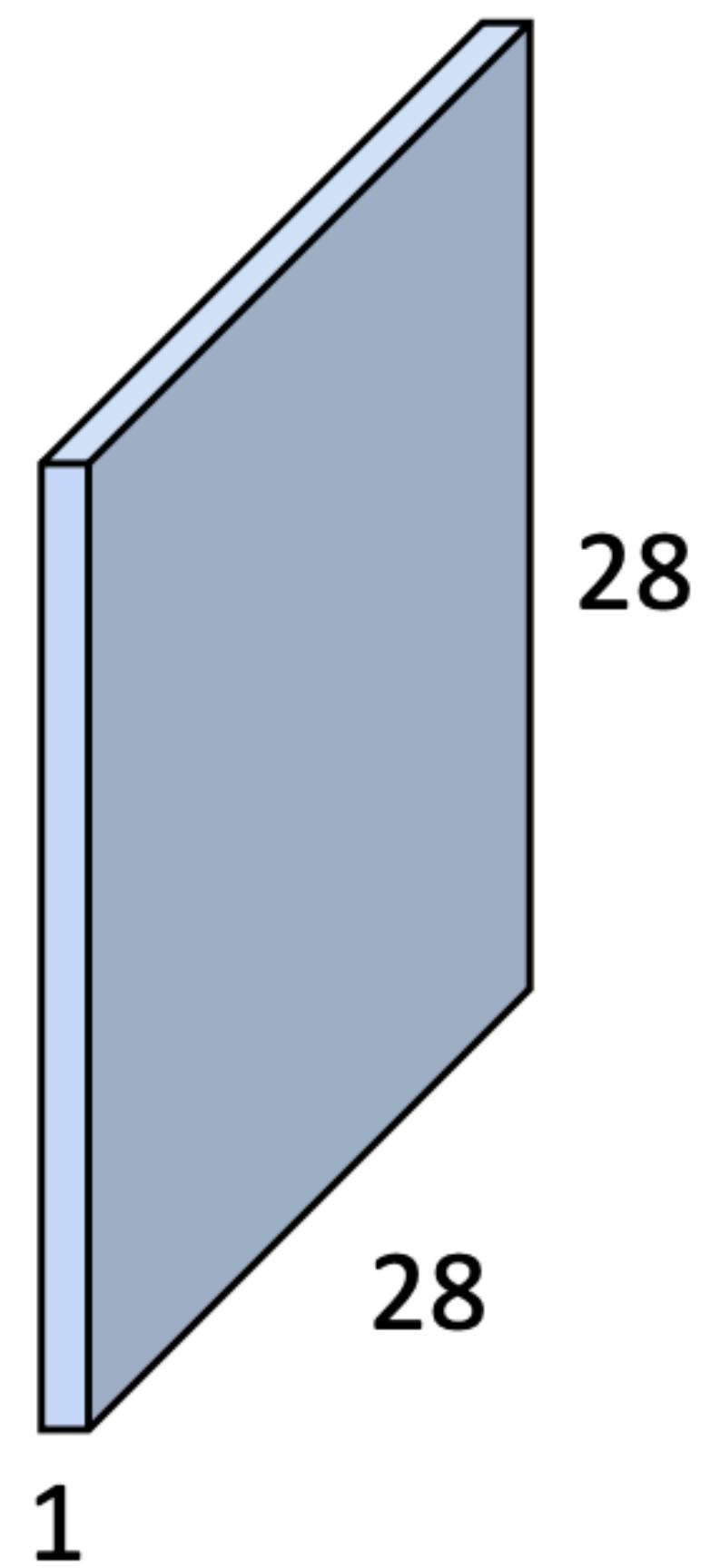
Convolution Layer

3x32x32 image

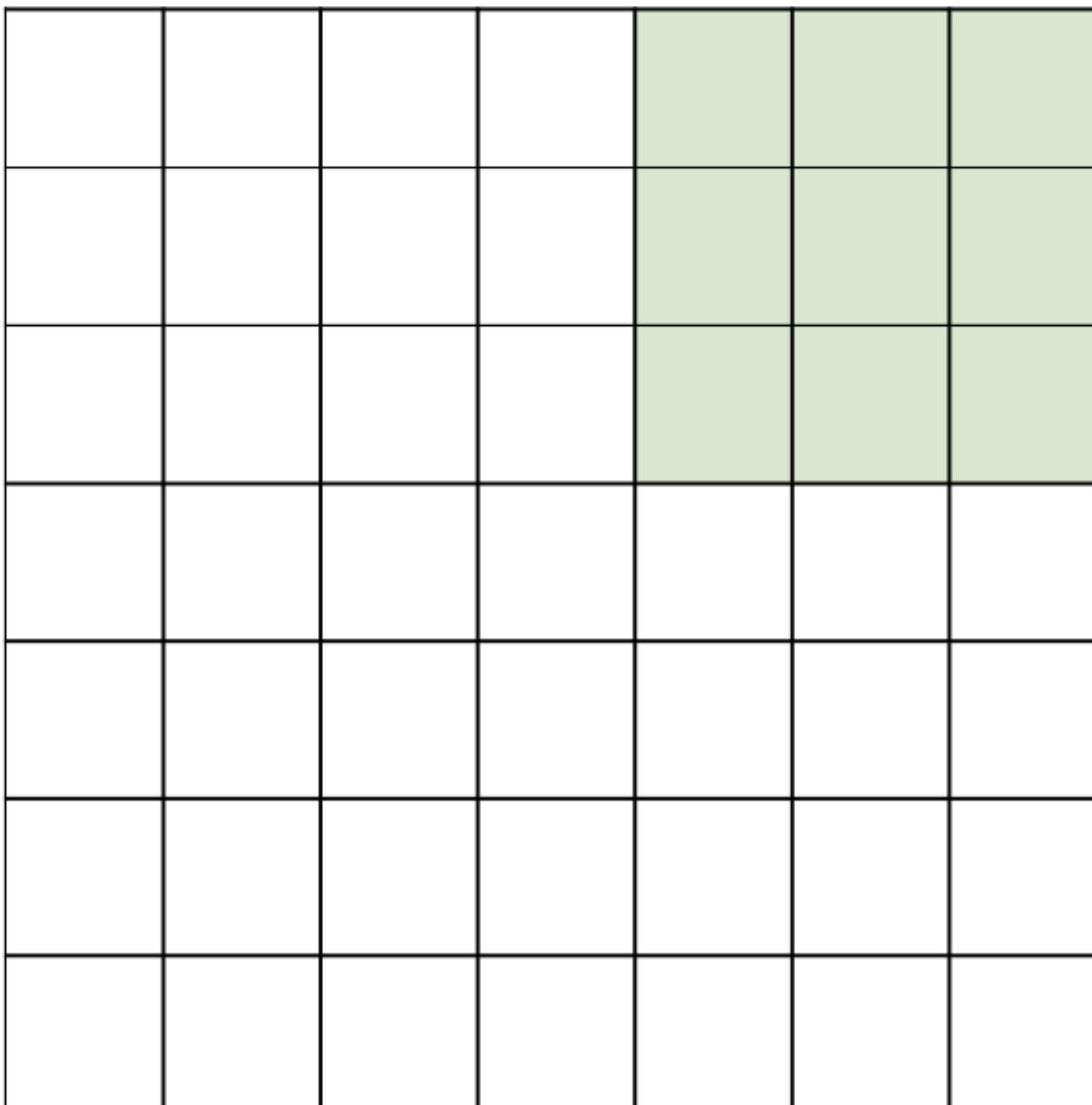


convolve (slide) over
all spatial locations

1x28x28
activation map



Stride



Input: 7x7

Filter: 3x3

Stride: 2

Output: 3x3

In general:

Input: W

Filter: K

Padding: P

Stride: S

Output: $(W - K + 2P) / S + 1$

Padding

A closer look at spatial dimensions

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input: 7x7

Filter: 3x3

Output: 5x5

In general: Problem: Feature

Input: W

Filter: K

Output: $W - K + 1$

maps “shrink”

with each layer!

Solution: padding

Add zeros around the input

Pooling

8	7	5	3
12	9	5	7
13	2	10	3
9	4	5	14

2x2 pooling,
stride 2

Max pooling

12	7
13	14

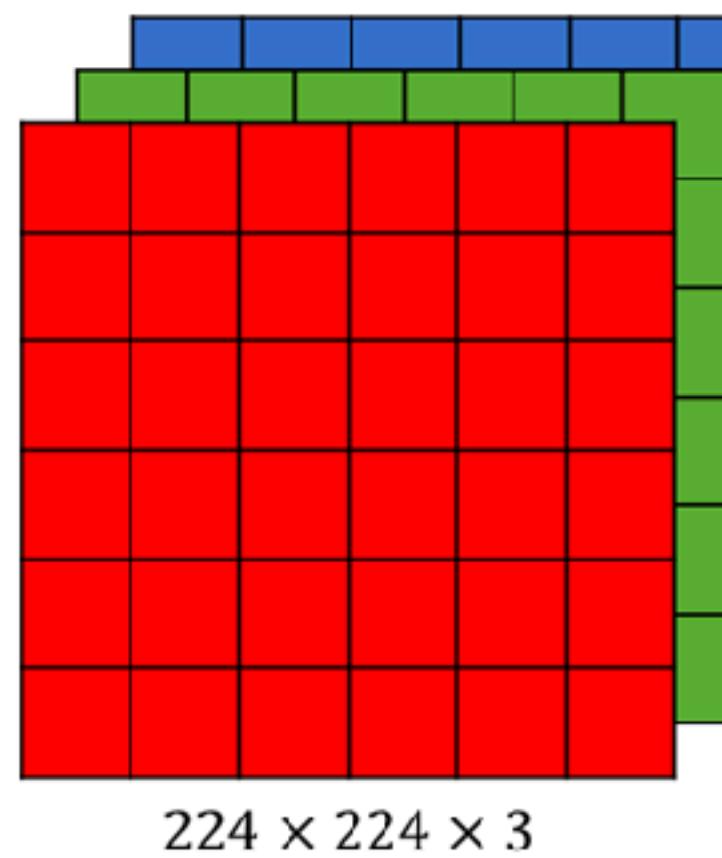
Average pooling

9	5
7	8

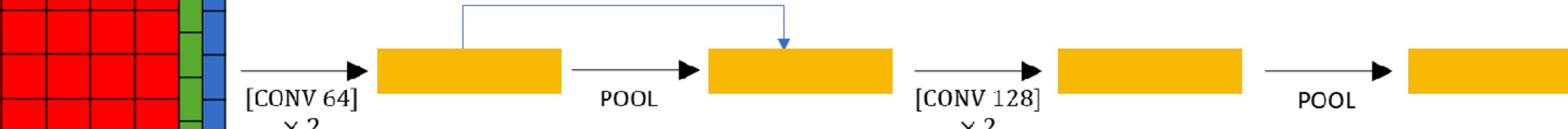
4. CNN

기초 이론

CONV = 3×3 filter, s=1, same

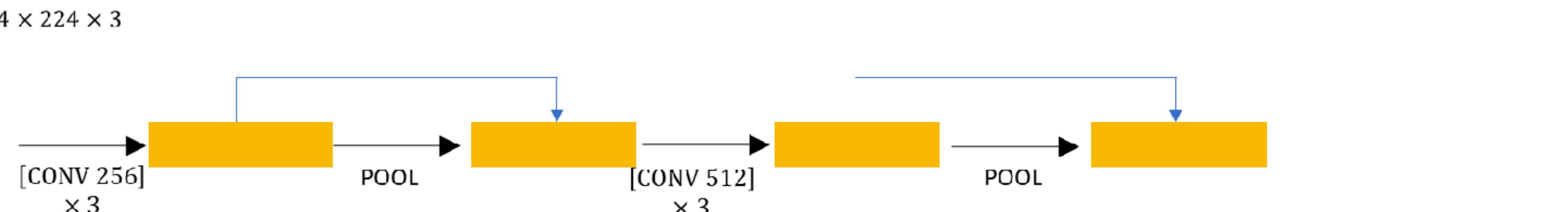


MAX-POOL = 2×2 , s=2



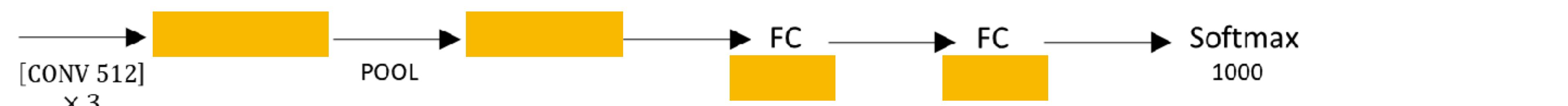
[CONV 128] $\times 2$

POOL



[CONV 512] $\times 3$

POOL



[CONV 512] $\times 3$

POOL

FC

FC

Softmax
1000

$n_h, n_w \downarrow$

n_c, \uparrow

$\sim 138M$ parameters

5. CNN with TF Fashion MNIST

실습