

# **CS 4644-DL / 7643-A: LECTURE 10**

## **DANFEI XU**

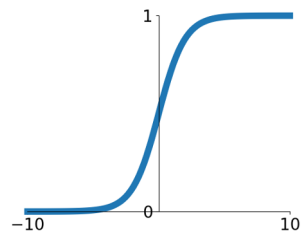
Topics:

- Training Neural Networks (Part 2)

# Activation Functions

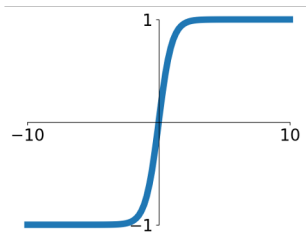
## Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



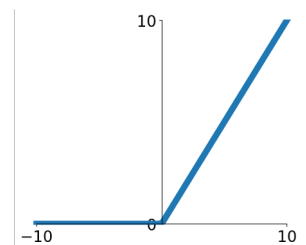
## tanh

$$\tanh(x)$$



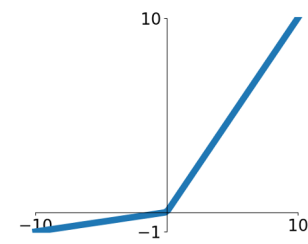
## ReLU

$$\max(0, x)$$



## Leaky ReLU

$$\max(0.1x, x)$$

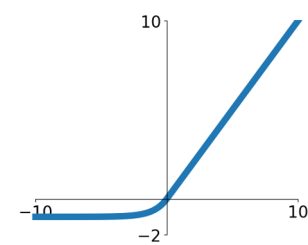


## Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

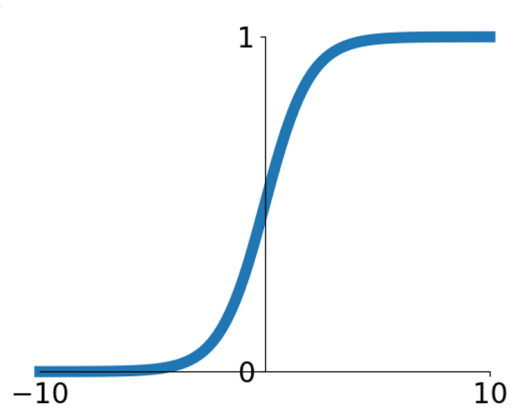
## ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



# Activation Functions

$$\sigma(x) = 1/(1 + e^{-x})$$



**Sigmoid**

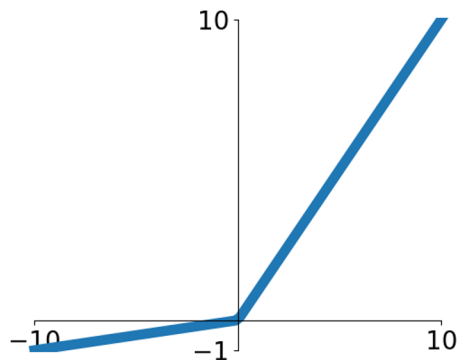
- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

3 problems:

- 1. Saturated neurons “kill” the gradients**
- 2. Sigmoid outputs are not zero-centered**
- 3.  $\exp()$  is a bit compute expensive**

# Activation Functions

[Mass et al., 2013]  
[He et al., 2015]



## Leaky ReLU

$$f(x) = \max(0.01x, x)$$

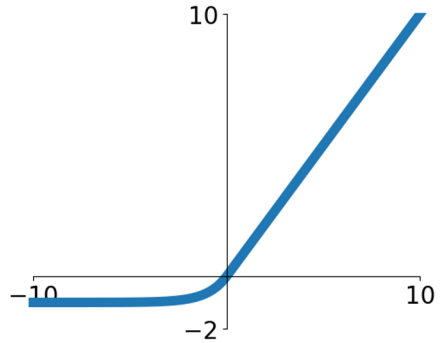
- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- **will not “die”.**

## Parametric Rectifier (PReLU)

$$f(x) = \max(\alpha x, x)$$

backprop into  $\alpha$   
(parameter)

## Exponential Linear Units (ELU)

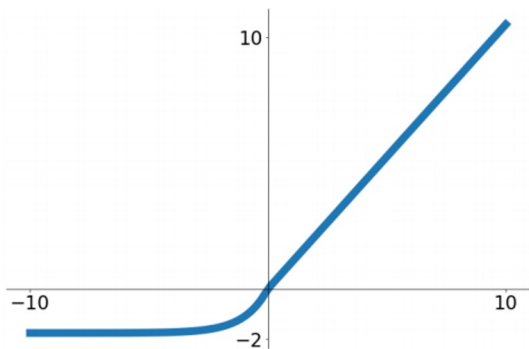


$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha (\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$

(Alpha default = 1)

- All benefits of ReLU
- Negative saturation encodes presence of features (all goes to  $-\alpha$ ), not magnitude
- Same in backprop
- Compared with Leaky ReLU: more robust to noise

## Scaled Exponential Linear Units (SELU)



$$f(x) = \begin{cases} \lambda x & \text{if } x > 0 \\ \lambda \alpha (e^x - 1) & \text{otherwise} \end{cases}$$

$$\alpha = 1.6732632423543772848170429916717$$
$$\lambda = 1.0507009873554804934193349852946$$

- Scaled version of ELU that works better for deep networks
- “Self-normalizing” property;
- Can train deep SELU networks without BatchNorm
  - (will discuss more later)

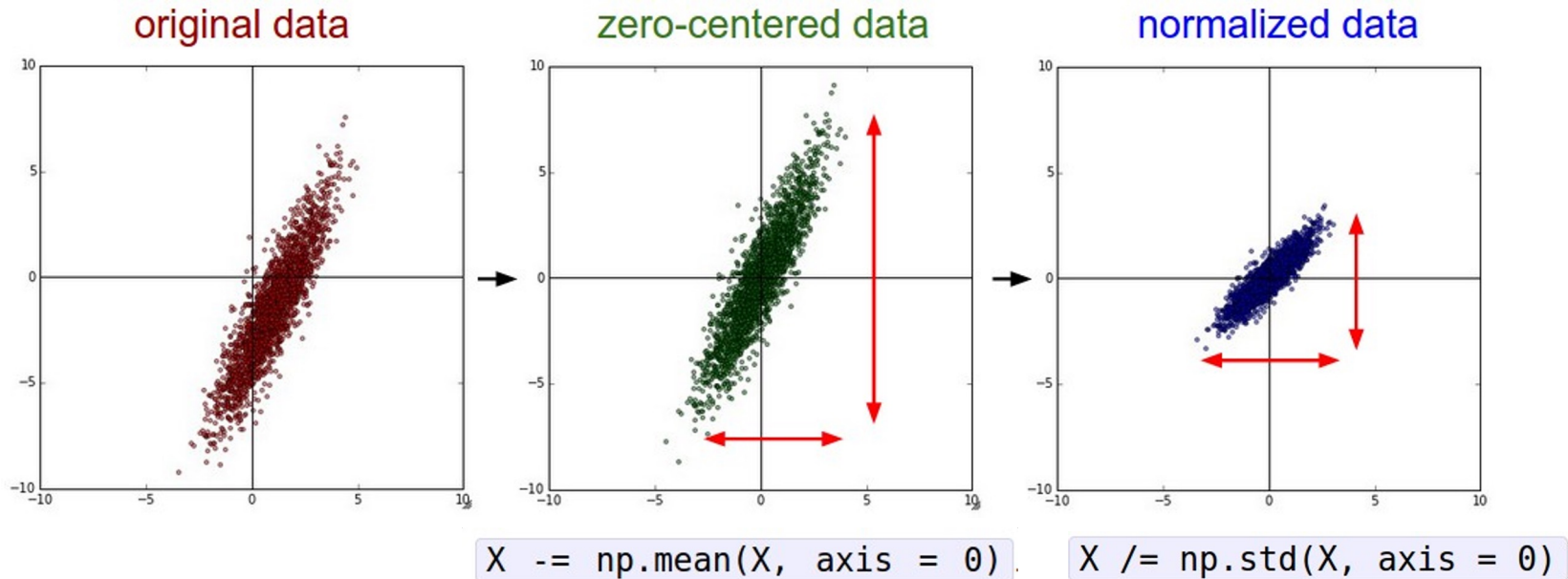
Derivation takes 91 pages of math in appendix...

(Klambauer et al, Self-Normalizing Neural Networks, ICLR 2017)

## TLDR: In practice:

- Many possible choices beyond what we've talked here, but ...
- Use **ReLU**. Be careful with your learning rates
- Try out **Leaky ReLU / ELU / SELU**
  - To squeeze out some marginal gains
- Don't use **sigmoid** or **tanh**

# Data Preprocessing

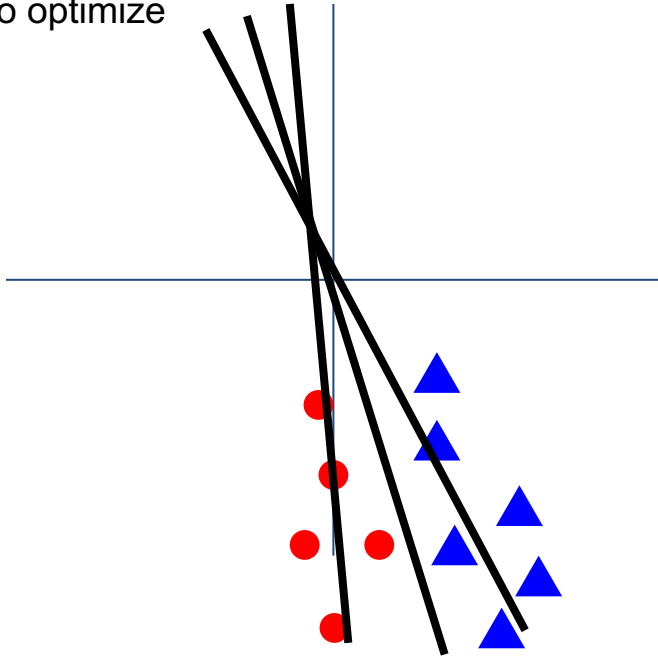


(Assume  $X$  [NxD] is data matrix,  
each example in a row)

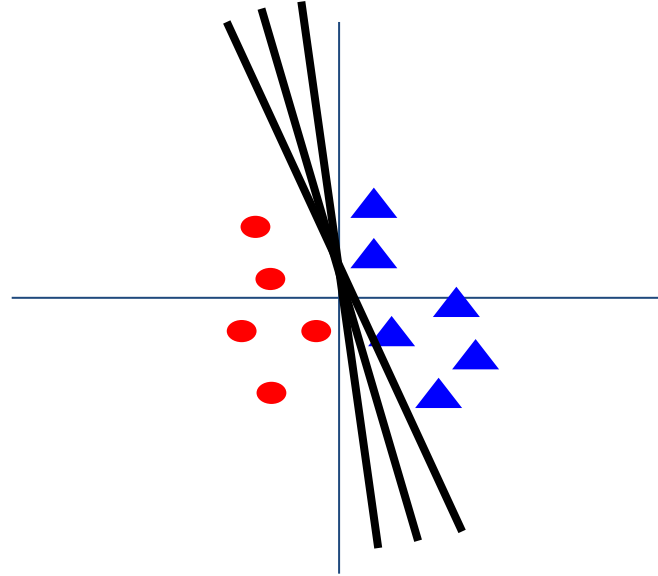


# Data Preprocessing

**Before normalization:** classification loss very sensitive to changes in weight matrix; hard to optimize



**After normalization:** less sensitive to small changes in weights; easier to optimize



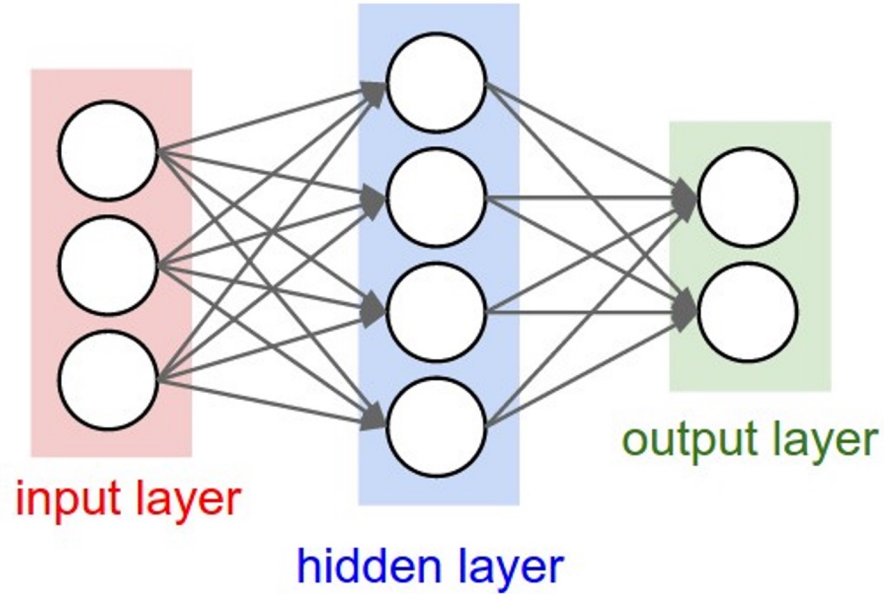
# This Time:

## **Training** Deep Neural Networks

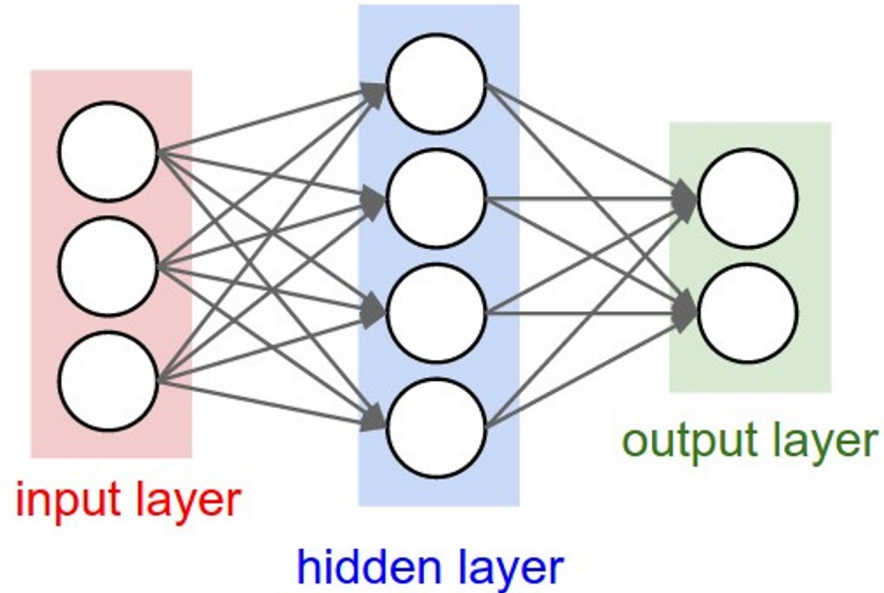
- Details of the non-linear activation functions
- Data normalization
- **Weight Initialization**
- **Batch Normalization**
- **Advanced Optimization**
- **Regularization**
- Data Augmentation
- Transfer learning
- Hyperparameter Tuning
- Model Ensemble

# Weight Initialization

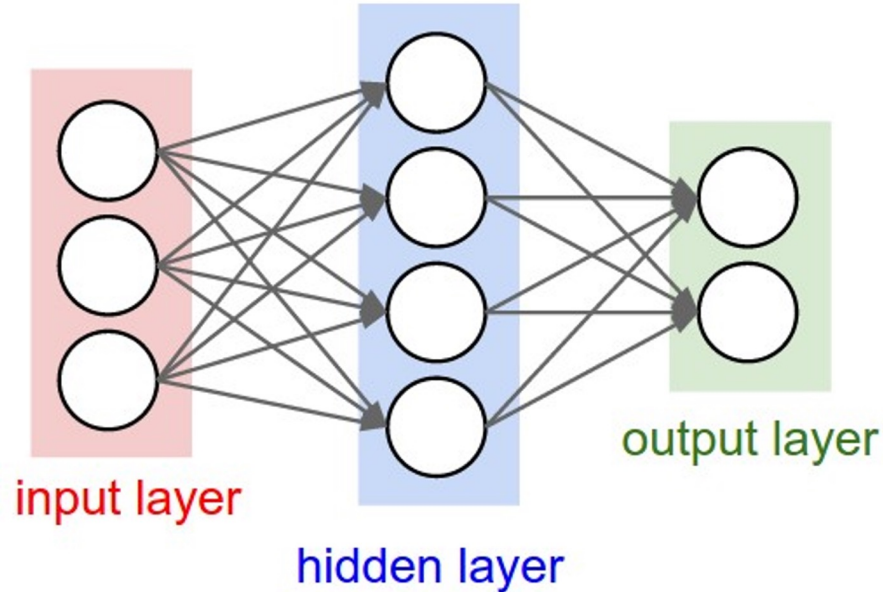
- Q: what happens when  $W$ =same initial value is used?



- Q: what happens when  $W$ =same initial value is used?
- A: All output will be the same!  $w_1^T x = w_2^T x$  if  $w_1 = w_2$



- Q: what happens when  $W$ =same initial value is used?
- A: All output will be the same!  $w_1^T x = w_2^T x$  if  $w_1 = w_2$
- Want to **maintain variance** through the layers.



- First idea: **Small random numbers**  
(gaussian with zero mean and 1e-2 standard deviation)

```
W = 0.01 * np.random.randn(Din, Dout)
```

- First idea: **Small random numbers**  
(gaussian with zero mean and  $1e-2$  standard deviation)

```
W = 0.01 * np.random.randn(Din, Dout)
```

Works ~okay for small networks, but problems with deeper networks.



# Weight Initialization: Activation statistics

```
dims = [4096] * 7      Forward pass for a 6-layer
hs = []               net with hidden size 4096
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = 0.01 * np.random.randn(Din, Dout)
    x = np.tanh(x.dot(W))
    hs.append(x)
```

What will happen to the activations for the last layer?

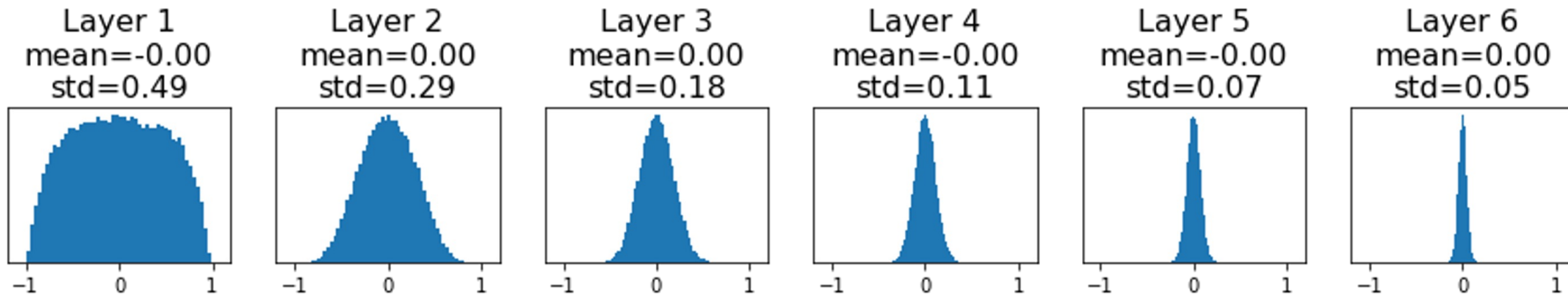
# Weight Initialization: Activation statistics

```
dims = [4096] * 7      Forward pass for a 6-layer
hs = []               net with hidden size 4096
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = 0.01 * np.random.randn(Din, Dout)
    x = np.tanh(x.dot(W))
    hs.append(x)
```

All activations tend to zero for deeper network layers

**Q:** What do the gradients  $dL/dW$  look like?

**Hint:**  $\frac{\partial L}{\partial w} = x^T \left( \frac{\partial L}{\partial y} \right)$



Visualize distribution of activations

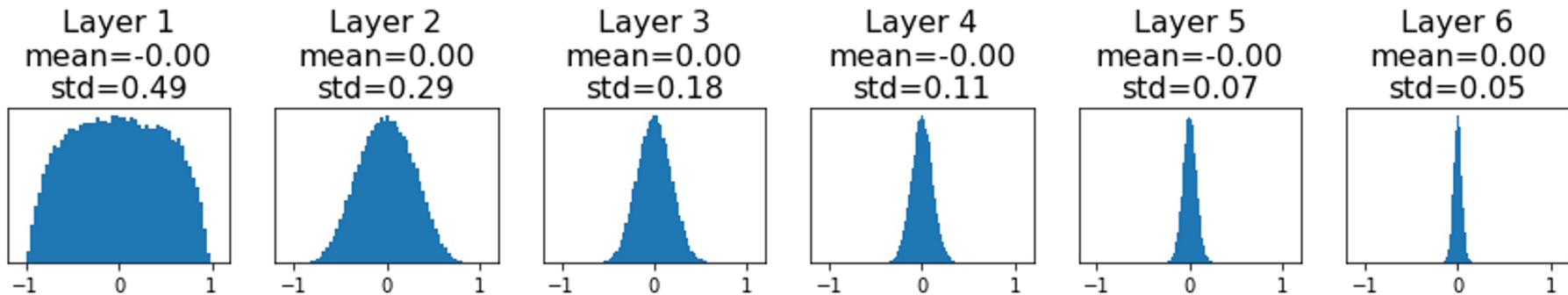
# Weight Initialization: Activation statistics

```
dims = [4096] * 7      Forward pass for a 6-layer
hs = []               net with hidden size 4096
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = 0.01 * np.random.randn(Din, Dout)
    x = np.tanh(x.dot(W))
    hs.append(x)
```

All activations tend to zero for deeper network layers

**Q:** What do the gradients  $dL/dW$  look like?

**A:** All zero, no learning =(



Visualize distribution of activations

# Weight Initialization: Activation statistics

```
dims = [4096] * 7    Increase std of initial
hs = []             weights from 0.01 to 0.05
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = 0.05 * np.random.randn(Din, Dout)
    x = np.tanh(x.dot(W))
    hs.append(x)
```

Initialize with higher values

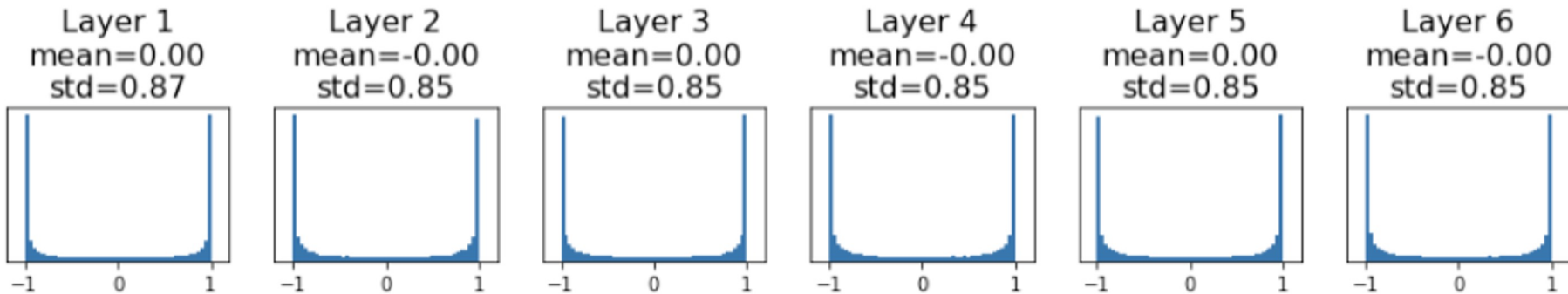
What will happen to the activations for the last layer?

# Weight Initialization: Activation statistics

```
dims = [4096] * 7      Increase std of initial
hs = []                weights from 0.01 to 0.05
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = 0.05 * np.random.randn(Din, Dout)
    x = np.tanh(x.dot(W))
    hs.append(x)
```

All activations saturate

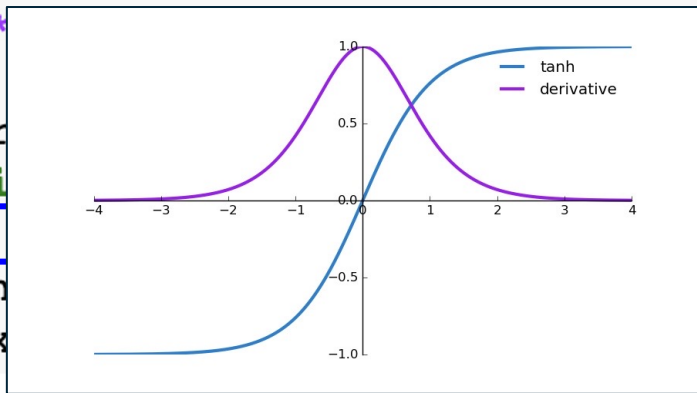
**Q:** What do the gradients look like?



Visualize distribution of activations

# Weight Initialization: Activation statistics

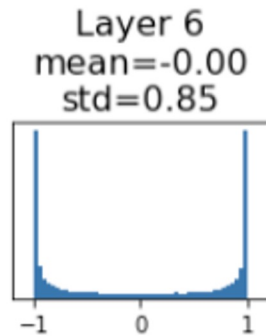
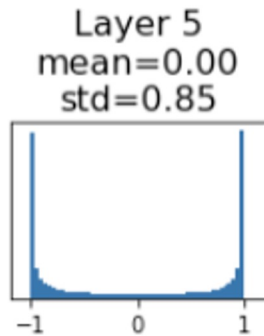
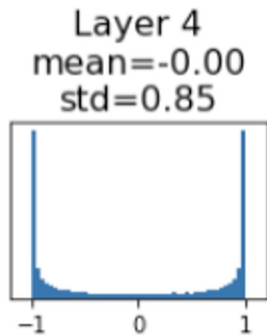
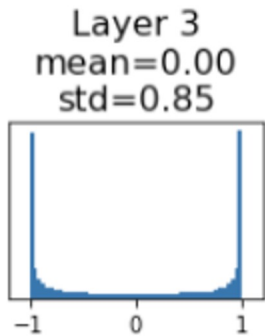
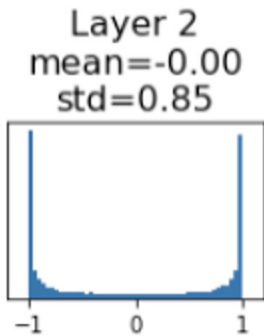
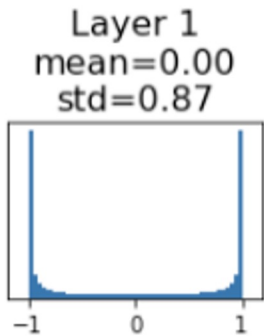
```
dims = [4096] *  
hs = []  
x = np.random.r  
for Din, Dout i  
    W = 0.05 *  
    x = np.tanh  
    hs.append(x
```



All activations saturate

**Q:** What do the gradients look like?

**A:** Local gradients all zero, no learning =(



Visualize distribution of activations

# Weight Initialization: Activation statistics

```
dims = [4096] * 7    Increase std of initial
hs = []             weights from 0.01 to 0.05
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = 0.05 * np.random.randn(Din, Dout)
    x = np.tanh(x.dot(W))
    hs.append(x)
```

All activations saturate

**Q:** What do the gradients look like?

More generally, *gradient explosion* (high  $w \rightarrow$  high output  $\rightarrow$  high gradient).

In general

- Small weights  $\rightarrow$  small output  $\rightarrow$  vanishing gradient in backpropagation.
- Large weights  $\rightarrow$  large output  $\rightarrow$  exploding gradient in backpropagation.

How do we initialize the weights “just right”?

# 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)
```

Glorot and Bengio, “Understanding the difficulty of training deep feedforward neural networks”, AISTAT 2010

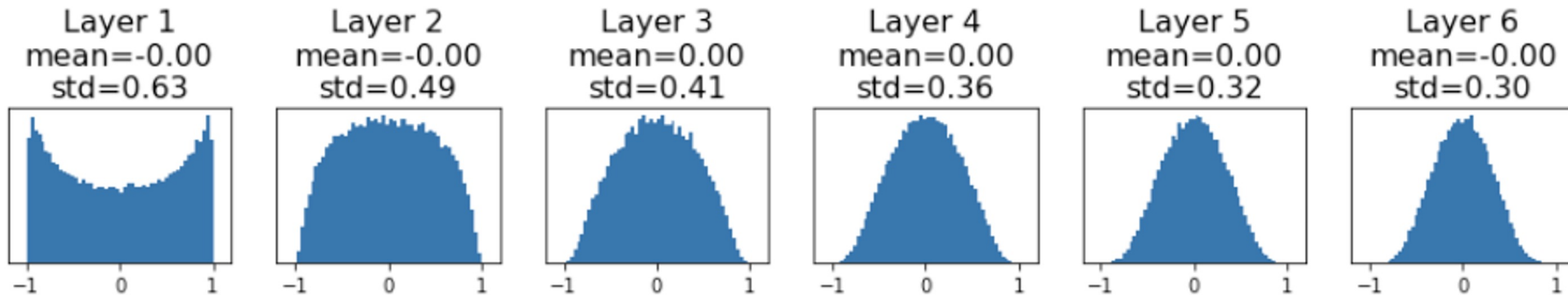


# Weight Initialization: “Xavier” Initialization

```
dims = [4096] * 7
hs = []
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)
```

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

“Just right”: Activations are nicely scaled for all layers!



Glorot and Bengio, “Understanding the difficulty of training deep feedforward neural networks”, AISTAT 2010

Visualize distribution of activations

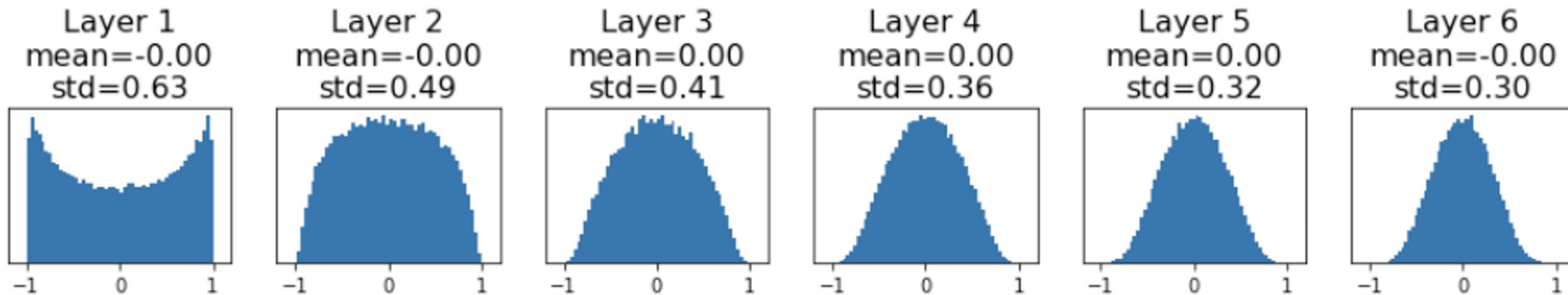
# Weight Initialization: “Xavier” Initialization

```
dims = [4096] * 7
hs = []
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)
```

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

“Just right”: Activations are nicely scaled for all layers!

For conv layers,  $D_{in}$  is  $filter\_size^2 * input\_channels$



Glorot and Bengio, “Understanding the difficulty of training deep feedforward neural networks”, AISTAT 2010

Visualize distribution of activations

# Weight Initialization: “Xavier” Initialization

```
dims = [4096] * 7
hs = []
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)
```

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

“Just right”: Activations are nicely scaled for all layers!

For conv layers,  $D_{in}$  is  $\text{filter\_size}^2 * \text{input\_channels}$

**Let:**  $y = x_1W_1 + x_2W_2 + \dots + x_{D_{in}}W_{D_{in}}$

# Weight Initialization: “Xavier” Initialization

```
dims = [4096] * 7
hs = []
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)
```

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

“Just right”: Activations are nicely scaled for all layers!

For conv layers,  $D_{in}$  is  $\text{filter\_size}^2 * \text{input\_channels}$

**Let:**  $y = x_1W_1 + x_2W_2 + \dots + x_{D_{in}}W_{D_{in}}$

**Assume:**  $\text{Var}(x_1) = \text{Var}(x_2) = \dots = \text{Var}(x_{D_{in}})$

# Weight Initialization: “Xavier” Initialization

```
dims = [4096] * 7
hs = []
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)
```

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

“Just right”: Activations are nicely scaled for all layers!

For conv layers,  $D_{in}$  is  $\text{filter\_size}^2 * \text{input\_channels}$

**Let:**  $y = x_1W_1 + x_2W_2 + \dots + x_{D_{in}}W_{D_{in}}$

**Assume:**  $\text{Var}(x_1) = \text{Var}(x_2) = \dots = \text{Var}(x_{D_{in}})$

**We want:**  $\text{Var}(y) = \text{Var}(x_i)$

# Weight Initialization: “Xavier” Initialization

```
dims = [4096] * 7
hs = []
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)
```

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

“Just right”: Activations are nicely scaled for all layers!

For conv layers,  $D_{in}$  is  $\text{filter\_size}^2 * \text{input\_channels}$

**Let:**  $y = x_1w_1 + x_2w_2 + \dots + x_{D_{in}}w_{D_{in}}$

**Assume:**  $\text{Var}(x_1) = \text{Var}(x_2) = \dots = \text{Var}(x_{D_{in}})$

**We want:**  $\text{Var}(y) = \text{Var}(x_i)$

$\text{Var}(y) = \text{Var}(x_1w_1 + x_2w_2 + \dots + x_{D_{in}}w_{D_{in}})$   
[substituting value of  $y$ ]

# Weight Initialization: “Xavier” Initialization

```
dims = [4096] * 7
hs = []
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)
```

“Xavier” initialization:  
std = 1/sqrt(Din)

“Just right”: Activations are nicely scaled for all layers!

For conv layers, Din is filter\_size<sup>2</sup> \* input\_channels

**Let:**  $y = x_1w_1 + x_2w_2 + \dots + x_{Din}w_{Din}$

**Assume:**  $\text{Var}(x_1) = \text{Var}(x_2) = \dots = \text{Var}(x_{Din})$

**We want:**  $\text{Var}(y) = \text{Var}(x_i)$

$\text{Var}(y) = \text{Var}(x_1w_1 + x_2w_2 + \dots + x_{Din}w_{Din})$

$= \sum \text{Var}(x_iw_i) = \text{Din} \text{Var}(x_iw_i)$

[Assume all  $x_i, w_i$  are iid]  $\sigma_{X+Y}^2 = \sigma_X^2 + \sigma_Y^2$

# Weight Initialization: “Xavier” Initialization

```
dims = [4096] * 7
hs = []
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)
```

“Xavier” initialization:  
std = 1/sqrt(Din)

“Just right”: Activations are nicely scaled for all layers!

For conv layers, Din is filter\_size<sup>2</sup> \* input\_channels

**Let:**  $y = x_1w_1 + x_2w_2 + \dots + x_{Din}w_{Din}$

**Assume:**  $\text{Var}(x_1) = \text{Var}(x_2) = \dots = \text{Var}(x_{Din})$

**We want:**  $\text{Var}(y) = \text{Var}(x_i)$

$$\begin{aligned}\text{Var}(y) &= \text{Var}(x_1w_1 + x_2w_2 + \dots + x_{Din}w_{Din}) \\ &= \text{Din} \text{Var}(x_iw_i) \\ &= \text{Din} \text{Var}(x_i) \text{Var}(w_i)\end{aligned}$$

[Assume all  $x_i, w_i$  are zero mean]

$$\begin{aligned}\text{Var}(XY) &= E(X^2Y^2) - (E(XY))^2 = \text{Var}(X)\text{Var}(Y) + \text{Var}(X)(E(Y))^2 \\ &\quad + \text{Var}(Y)(E(X))^2\end{aligned}$$



# Weight Initialization: “Xavier” Initialization

```
dims = [4096] * 7
hs = []
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)
```

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

“Just right”: Activations are nicely scaled for all layers!

For conv layers,  $D_{in}$  is  $\text{filter\_size}^2 * \text{input\_channels}$

**Let:**  $y = x_1w_1 + x_2w_2 + \dots + x_{D_{in}}w_{D_{in}}$

**Assume:**  $\text{Var}(x_1) = \text{Var}(x_2) = \dots = \text{Var}(x_{D_{in}})$

**We want:**  $\text{Var}(y) = \text{Var}(x_i)$

$$\begin{aligned}\text{Var}(y) &= \text{Var}(x_1w_1 + x_2w_2 + \dots + x_{D_{in}}w_{D_{in}}) \\ &= D_{in} \text{Var}(x_iw_i) \\ &= D_{in} \text{Var}(x_i) \text{Var}(w_i) \\ &[\text{Assume all } x_i, w_i \text{ are iid}]\end{aligned}$$

So,  $\text{Var}(y) = \text{Var}(x_i)$  only when  $\text{Var}(w_i) = 1/D_{in}$

# Weight Initialization: “Xavier” Initialization

```
dims = [4096] * 7
hs = []
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)
```

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

“Just right”: Activations are nicely scaled for all layers!

For conv layers,  $D_{in}$  is  $\text{filter\_size}^2 * \text{input\_channels}$

**Let:**  $y = x_1w_1 + x_2w_2 + \dots + x_{D_{in}}w_{D_{in}}$

**Assume:**  $\text{Var}(x_1) = \text{Var}(x_2) = \dots = \text{Var}(x_{D_{in}})$

**We want:**  $\text{Var}(y) = \text{Var}(x_i)$

$$\begin{aligned}\text{Var}(y) &= \text{Var}(x_1w_1 + x_2w_2 + \dots + x_{D_{in}}w_{D_{in}}) \\ &= D_{in} \text{Var}(x_iw_i) \\ &= D_{in} \text{Var}(x_i) \text{Var}(w_i) \\ &[\text{Assume all } x_i, w_i \text{ are iid}]\end{aligned}$$

So,  $\text{Var}(y) = \text{Var}(x_i)$  only when  $\text{Var}(w_i) = 1/D_{in}$

Scaling a normal distribution (std=1) to have  $\text{Var}=1/D_{in}$  -> multiply by  $\sqrt{1/D_{in}}$

# Weight Initialization: “Xavier” Initialization

```
dims = [4096] * 7
hs = []
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)
```

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

“Just right”: Activations are nicely scaled for all layers!

For conv layers,  $D_{in}$  is  $\text{filter\_size}^2 * \text{input\_channels}$

**Let:**  $y = x_1w_1 + x_2w_2 + \dots + x_{D_{in}}w_{D_{in}}$

**Assume:**  $\text{Var}(x_1) = \text{Var}(x_2) = \dots = \text{Var}(x_{D_{in}})$

**We want:**  $\text{Var}(y) = \text{Var}(x_i)$

$$\begin{aligned}\text{Var}(y) &= \text{Var}(x_1w_1 + x_2w_2 + \dots + x_{D_{in}}w_{D_{in}}) \\ &= D_{in} \text{Var}(x_iw_i) \\ &= D_{in} \text{Var}(x_i) \text{Var}(w_i) \\ &[\text{Assume all } x_i, w_i \text{ are iid}]\end{aligned}$$

So,  $\text{Var}(y) = \text{Var}(x_i)$  only when  $\text{Var}(w_i) = 1/D_{in}$

In practice, use  $\text{Var} = 2 / (D_{in} + D_{out})$  to account for both forward and backward pass

# Weight Initialization: What about ReLU?

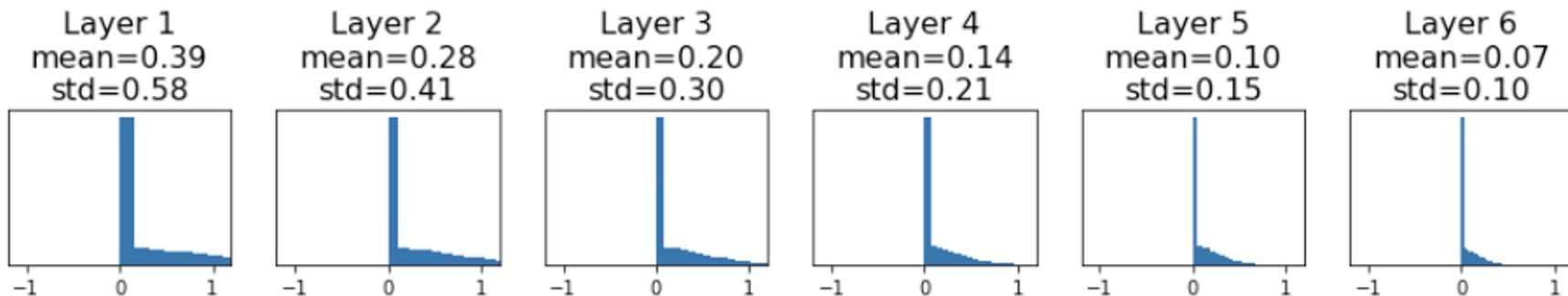
```
dims = [4096] * 7      Change from tanh to ReLU
hs = []
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = np.random.randn(Din, Dout) / np.sqrt(Din)
    x = np.maximum(0, x.dot(W))
    hs.append(x)
```

# Weight Initialization: What about ReLU?

```
dims = [4096] * 7      Change from tanh to ReLU
hs = []
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.maximum(0, x.dot(W))
    hs.append(x)
```

Xavier assumes zero centered activation function

Activations collapse to zero again, no learning =(



Visualize distribution of activations

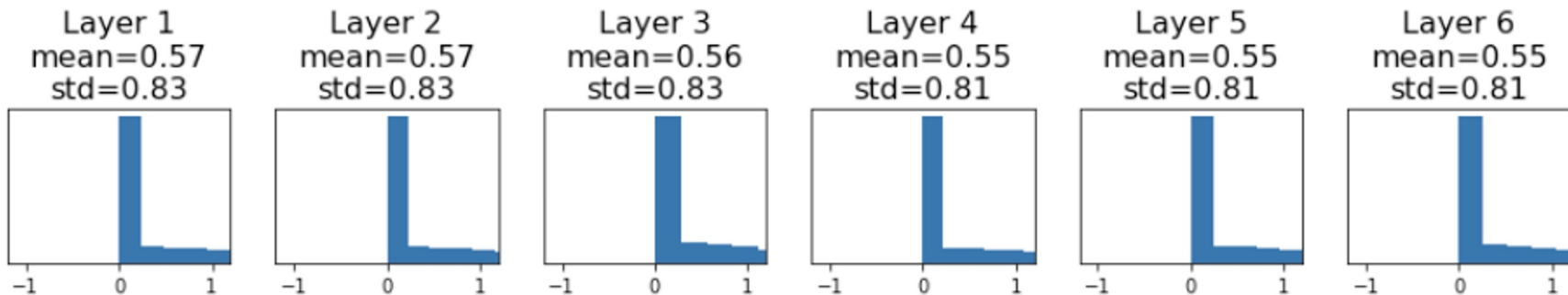
# Weight Initialization: Kaiming / MSRA Initialization

```
dims = [4096] * 7
hs = []
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = np.random.randn(Din, Dout) * np.sqrt(2/Din)
    x = np.maximum(0, x.dot(W))
    hs.append(x)
```

ReLU correction:  $\text{std} = \sqrt{2 / \text{Din}}$

Issue: Half of the activation get killed.

Solution: make the non-zero output variance twice as large as input



He et al, "Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification", ICCV 2015

Visualize distribution of activations

# Proper initialization is an active area of research...

***Understanding the difficulty of training deep feedforward neural networks***

by Glorot and Bengio, 2010

***Exact solutions to the nonlinear dynamics of learning in deep linear neural networks*** by Saxe et al, 2013

***Random walk initialization for training very deep feedforward networks*** by Sussillo and Abbott, 2014

***Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification*** by He et al., 2015

***Data-dependent Initializations of Convolutional Neural Networks*** by Krähenbühl et al., 2015

***All you need is a good init***, Mishkin and Matas, 2015

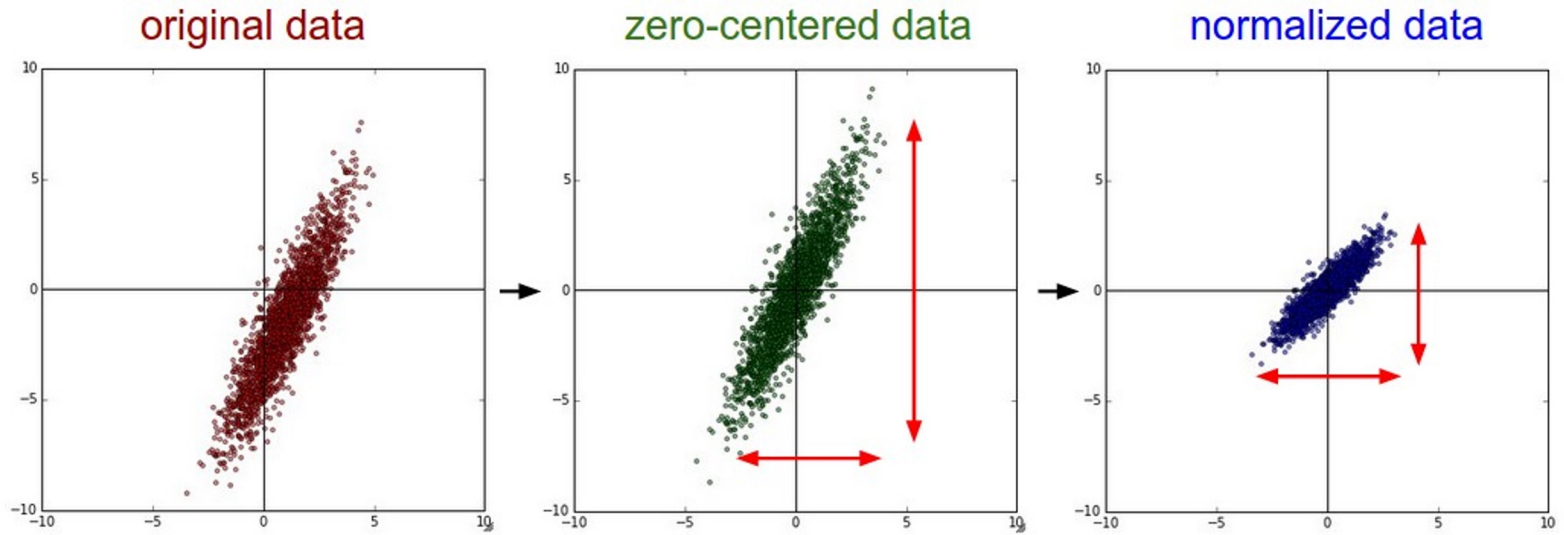
***Fixup Initialization: Residual Learning Without Normalization***, Zhang et al, 2019

***The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks***, Frankle and Carbin, 2019

# Batch Normalization



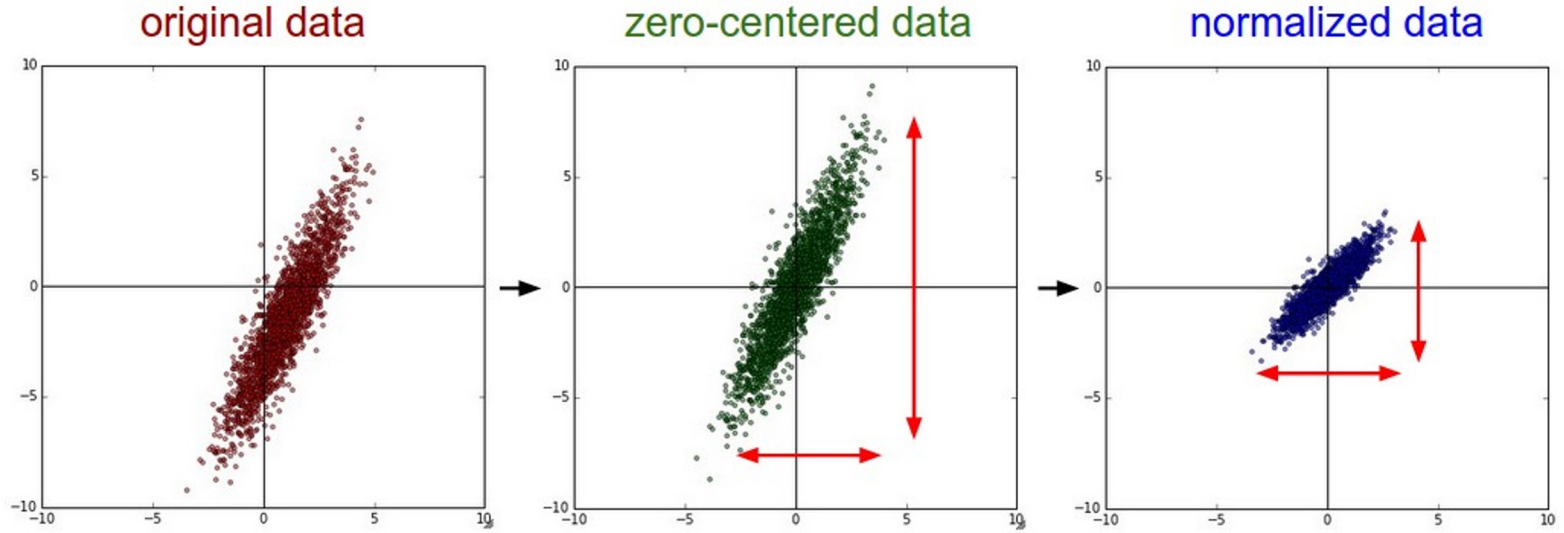
# Recall: Input Normalization



```
X -= np.mean(X, axis = 0)
```

```
X /= np.std(X, axis = 0)
```

# Recall: Input Normalization



```
X -= np.mean(X, axis = 0)
```

```
X /= np.std(X, axis = 0)
```

Problem: Only for input to the first layer. Input for later layers are longer normalized!

But can't do dataset normalization for intermediate layers! Activation distribution changes as the training progresses.

# Batch Normalization

“you want zero-mean unit-variance activations? just make them so.”

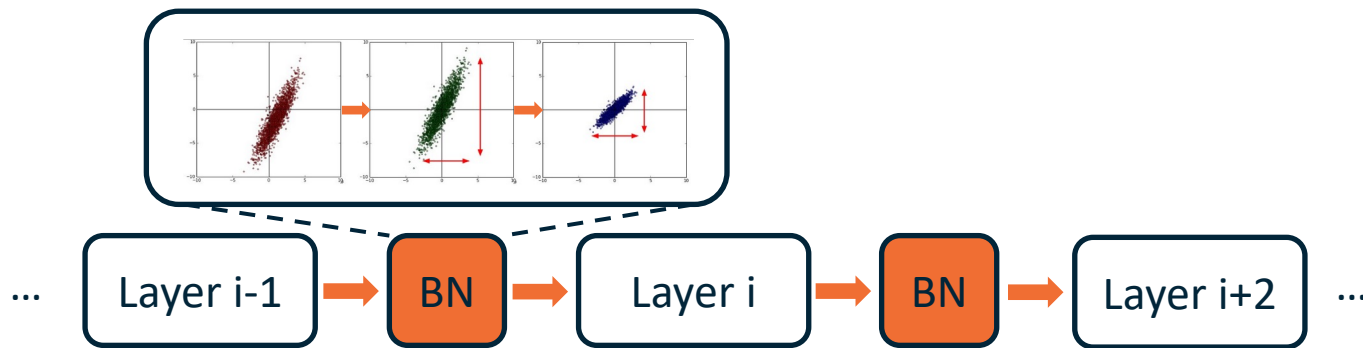
consider a **batch of activations**  $x$  at some layer. To make each dimension zero-mean unit-variance, apply:

$$\hat{x} = \frac{x - \mathbb{E}[x]}{\sqrt{\text{Var}[x]}}$$

this is a vanilla  
differentiable function...

# Batch Normalization

“you want zero-mean unit-variance activations? just make them so.”

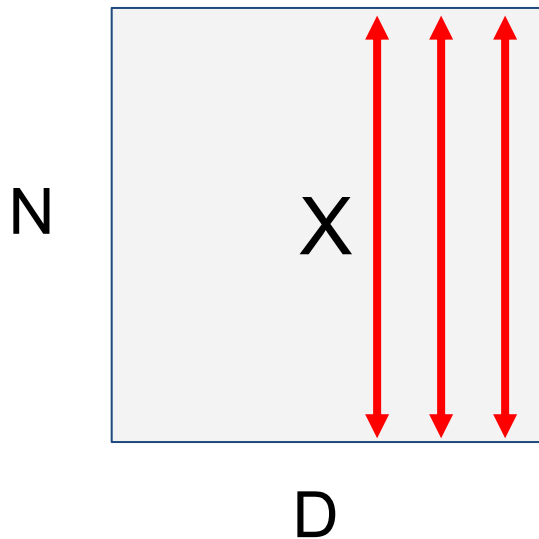


$$\hat{x} = \frac{x - \mathbb{E}[x]}{\sqrt{\text{Var}[x]}}$$

# Batch Normalization

[Ioffe and Szegedy, 2015]

**Input:**  $x : N \times D$



$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Per-channel mean,  
shape is D

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

Per-channel var,  
shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

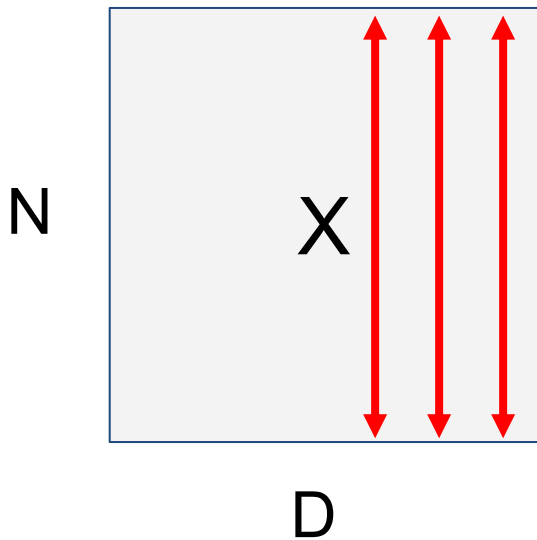
Normalized x,  
Shape is N x D

(Prevent div by 0 err)

# Batch Normalization

[Ioffe and Szegedy, 2015]

**Input:**  $x : N \times D$



$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Per-channel mean,  
shape is D

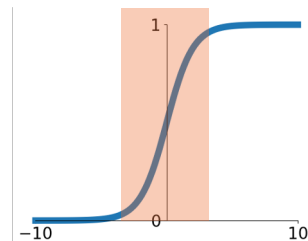
$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

Per-channel var,  
shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

Normalized x,  
Shape is N x D

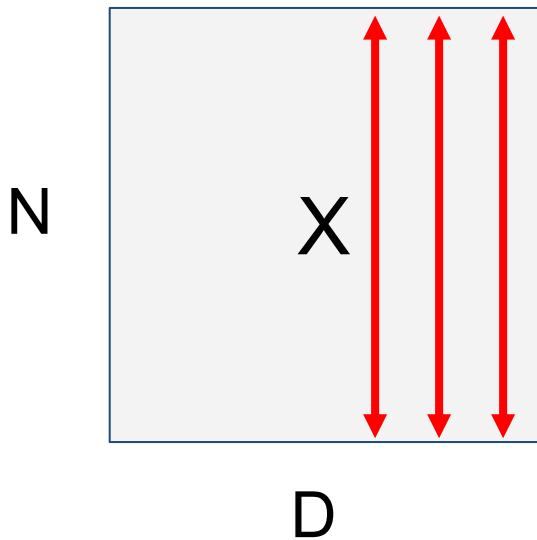
**Problem:** What if zero-mean, unit variance is too hard of a constraint?  
E.g., inserting a BN before sigmoid will constrain it to (mostly) linear regime



# Batch Normalization

[Ioffe and Szegedy, 2015]

**Input:**  $x : N \times D$



$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Per-channel mean,  
shape is D

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

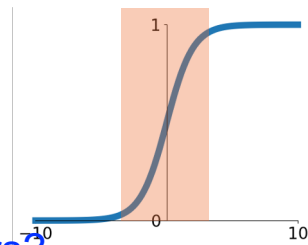
Per-channel var,  
shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

Normalized x,  
Shape is N x D

**Problem:** What if zero-mean, unit variance is too hard of a constraint?  
E.g., inserting a BN before sigmoid will constrain it to (mostly) linear regime

Can we learn the normalization parameters?



# Batch Normalization

[Ioffe and Szegedy, 2015]

**Input:**  $x : N \times D$   
**Learnable scale and shift parameters:**

$$\gamma, \beta : \mathbb{R}^D$$

We want to give the model a chance to **adjust batchnorm** if the default is not optimal.

Learning  $\gamma = \sigma$  and  $\beta = \mu$  will recover the identity function!

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Per-channel mean,  
shape is D

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

Per-channel var,  
shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

Normalized x,  
Shape is N x D

$$y_{i,j} = \underline{\gamma}_j \hat{x}_{i,j} + \underline{\beta}_j$$

Output,  
Shape is N x D



# Batch Normalization: Test-Time

Estimates depend on minibatch;  
can't do this at test-time!

**Input:**  $x : N \times D$   
**Learnable scale and  
shift parameters:**

$$\gamma, \beta: \mathbb{R}^D$$

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Per-channel mean,  
shape is D

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

Per-channel var,  
shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

Normalized x,  
Shape is N x D

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Output,  
Shape is N x D

# Batch Normalization: Test-Time

Estimates depend on minibatch;  
can't do this at test-time!

**Input:**  $x : N \times D$   
**Learnable scale and shift parameters:**

$$\gamma, \beta: \mathbb{R}^D$$

Activations become fixed after training. Can calculate training set-wide statistics for inference-time normalization.

At training time, do moving average to save compute.

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Per-channel mean,  
shape is D

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

Per-channel var,  
shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

Normalized x,  
Shape is N x D

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Output,  
Shape is N x D

# Batch Normalization: Test-Time

**Input:**  $x : N \times D$   
**Learnable scale and shift parameters:**

$$\gamma, \beta : \mathbb{R}^D$$

During testing batchnorm becomes a linear operator!  
Can be fused with the previous fully-connected or conv layer

$$\mu_j = \text{(Moving) average of values seen during training}$$

Per-channel mean, shape is D

$$\sigma_j^2 = \text{(Moving) average of values seen during training}$$

Per-channel var, shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

Normalized x, Shape is N x D

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Output, Shape is N x D

# Batch Normalization

[Ioffe and Szegedy, 2015]

Q: Should you put batchnorm before or after ReLU?

A: Topic of debate. Original paper says BN->ReLU. Now most commonly ReLU->BN. If BN-> ReLU and zero mean, ReLU kills half of the activations, but in practice makes insignificant differences.

Q: Should you normalize the **input** (e.g., images) with batchnorm?

A: No, you already have the fixed & correct dataset statistics, no need to do batchnorm.

Q: How many parameters does a batchnorm layer have?

A: Input dimension \* 4: beta, gamma, moving average mu, moving average sigma. Only beta and gamma are trainable parameters.

# Batch Normalization

[Ioffe and Szegedy, 2015]

- Makes deep networks **much** easier to train!
  - If you are interested in the theory, read <https://arxiv.org/abs/1805.11604>
  - TL;DR: makes optimization landscape smoother
- Allows higher learning rates, faster convergence
- More useful in deeper networks
- Networks become more robust to initialization
- More robust to range of input
- Zero overhead at test-time: can be fused with conv!
- Behaves differently during training and testing: this is a very common source of bugs!
- Needs large batch size to calculate accurate stats

# Batch Normalization for ConvNets

Batch Normalization for  
**fully-connected** networks

$$\mathbf{x}: \mathbf{N} \times \mathbf{D}$$

Normalize



$$\boldsymbol{\mu}, \boldsymbol{\sigma}: \mathbf{1} \times \mathbf{D}$$

$$\boldsymbol{\gamma}, \boldsymbol{\beta}: \mathbf{1} \times \mathbf{D}$$

$$\mathbf{y} = \boldsymbol{\gamma} (\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta}$$

Batch Normalization for  
**convolutional** networks  
(Spatial Batchnorm, BatchNorm2D)

$$\mathbf{x}: \mathbf{N} \times \mathbf{C} \times \mathbf{H} \times \mathbf{W}$$

Normalize



$$\boldsymbol{\mu}, \boldsymbol{\sigma}: \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1}$$

$$\boldsymbol{\gamma}, \boldsymbol{\beta}: \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1}$$

$$\mathbf{y} = \boldsymbol{\gamma} (\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta}$$

# Layer Normalization

**Batch Normalization** for fully-connected networks

$$\mathbf{x} : \mathbf{N} \times \mathbf{D}$$

Normalize



$$\boldsymbol{\mu}, \boldsymbol{\sigma} : \mathbf{1} \times \mathbf{D}$$

$$\boldsymbol{\gamma}, \boldsymbol{\beta} : \mathbf{1} \times \mathbf{D}$$

$$\mathbf{y} = \boldsymbol{\gamma} (\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta}$$

**Layer Normalization** for fully-connected networks  
Same behavior at train and test!

$$\mathbf{x} : \mathbf{N} \times \mathbf{D}$$

Normalize



$$\boldsymbol{\mu}, \boldsymbol{\sigma} : \mathbf{N} \times \mathbf{1}$$

$$\boldsymbol{\gamma}, \boldsymbol{\beta} : \mathbf{1} \times \mathbf{D}$$

$$\mathbf{y} = \boldsymbol{\gamma} (\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta}$$

Ba, Kiros, and Hinton, "Layer Normalization", arXiv 2016

More flexible (can use  $N = 1!$ ), works well with sequence models (RNN, Transformers)

# Instance Normalization

**Batch Normalization** for  
convolutional networks

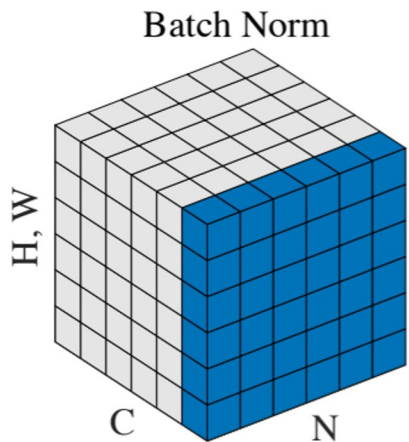
$$\begin{array}{l} \mathbf{x} : \mathbf{N} \times \mathbf{C} \times \mathbf{H} \times \mathbf{W} \\ \text{Normalize} \quad \downarrow \quad \downarrow \quad \downarrow \\ \boldsymbol{\mu}, \boldsymbol{\sigma} : \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1} \\ \boldsymbol{\gamma}, \boldsymbol{\beta}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

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

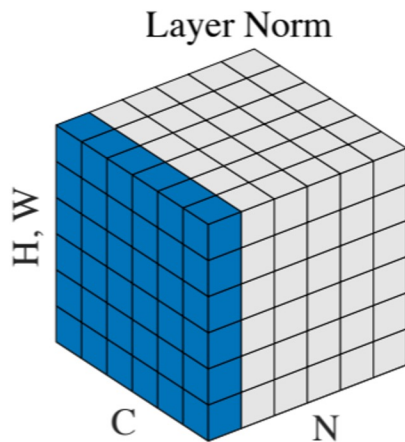
$$\begin{array}{l} \mathbf{x} : \mathbf{N} \times \mathbf{C} \times \mathbf{H} \times \mathbf{W} \\ \text{Normalize} \quad \quad \downarrow \quad \downarrow \\ \boldsymbol{\mu}, \boldsymbol{\sigma} : \mathbf{N} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1} \\ \boldsymbol{\gamma}, \boldsymbol{\beta}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



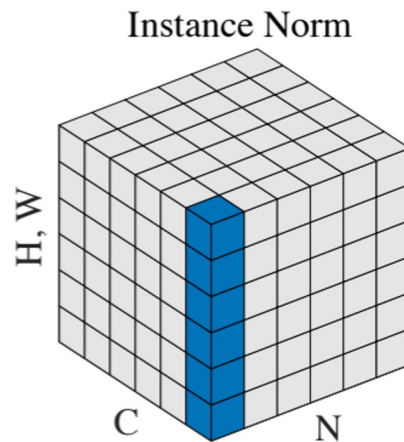
# Comparison of Normalization Layers



$N \times C \times H \times W \rightarrow$   
 $1 \times C \times 1 \times 1$

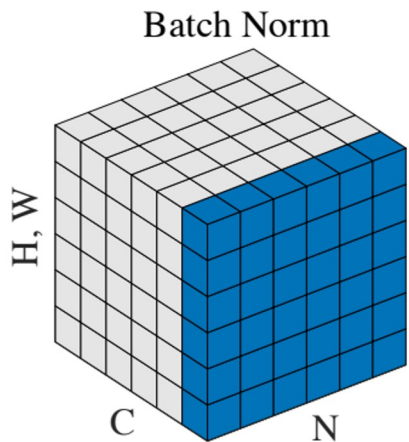


$N \times C \times H \times W \rightarrow$   
 $N \times 1 \times 1 \times 1$

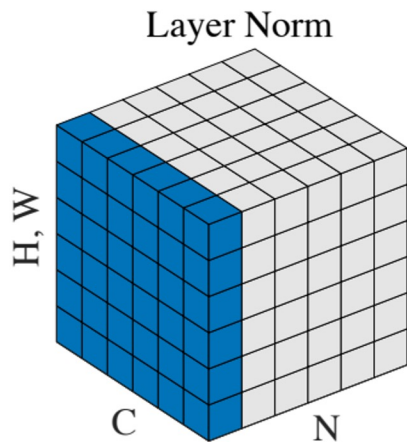


$N \times C \times H \times W \rightarrow$   
 $N \times C \times 1 \times 1$

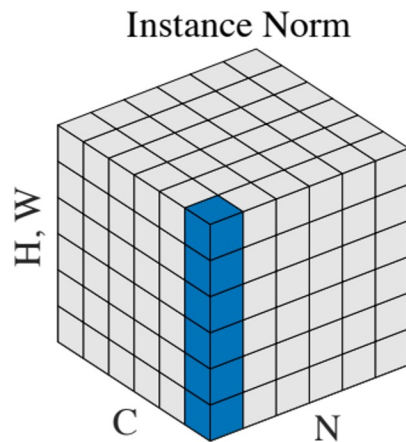
# Group Normalization



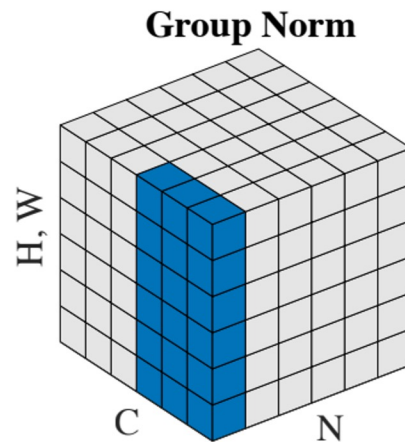
$$N \times C \times H \times W \rightarrow 1 \times C \times 1 \times 1$$



$$N \times C \times H \times W \rightarrow N \times 1 \times 1 \times 1$$



$$N \times C \times H \times W \rightarrow N \times C \times 1 \times 1$$



$$N \times C \times H \times W \rightarrow N \times C/G \times 1 \times 1$$

# (Fancier) Optimizers

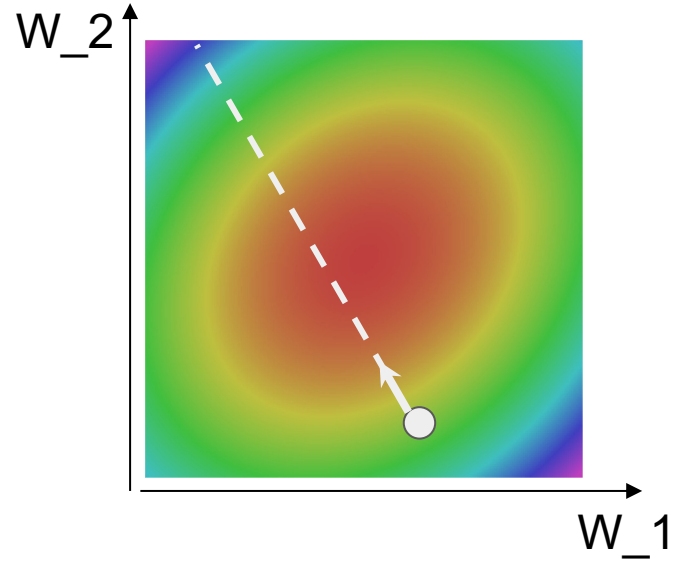
# Optimization

```
# Vanilla Gradient Descent
```

```
while True:
```

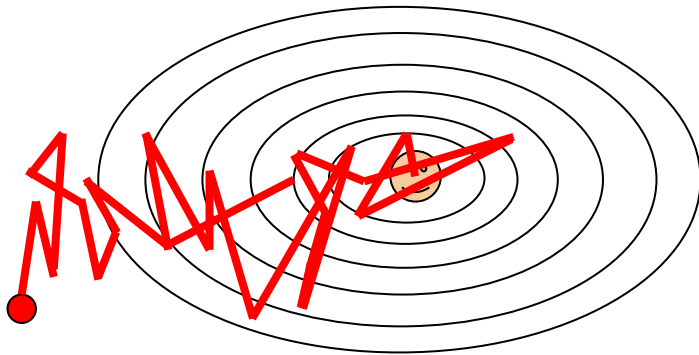
```
    weights_grad = evaluate_gradient(loss_fun, data, weights)
```

```
    weights += - step_size * weights_grad # perform parameter update
```



# Optimization: Problem #1 with SGD

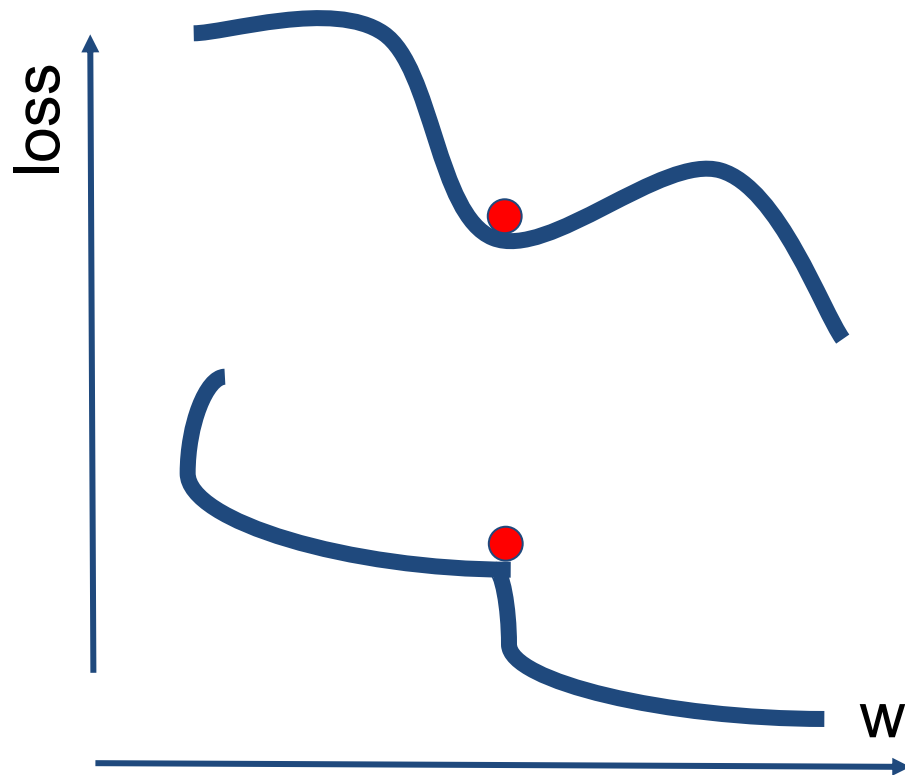
- Stochastic minibatch gives a noisy estimate of the true gradient direction. Very problematic when the batch size is small (e.g., due to compute resource limit).
- Poorly-selected learning rate makes the oscillation worse (overshoot)



[http://web.cs.ucla.edu/~chohsieh/teaching/CS260\\_Winter2019/lecture4.pdf](http://web.cs.ucla.edu/~chohsieh/teaching/CS260_Winter2019/lecture4.pdf)

# Optimization: Problem #2 with SGD

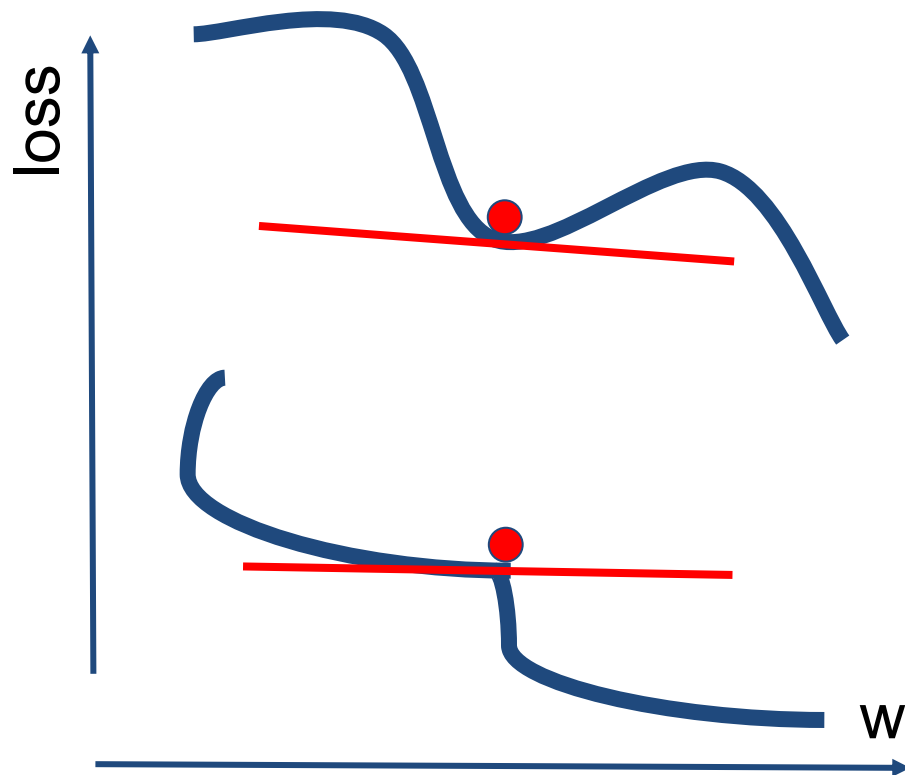
What if the loss function has a **local minima** or **saddle point**?



# Optimization: Problem #2 with SGD

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

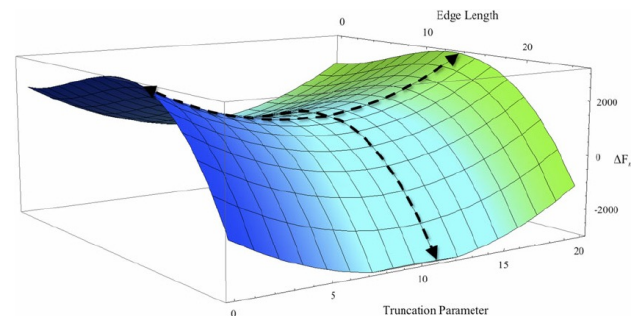
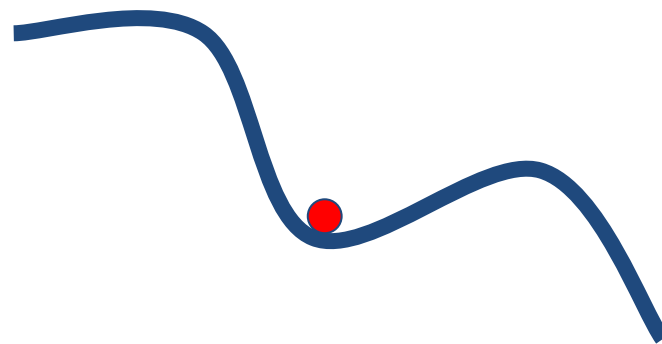
Zero gradient,  
gradient descent  
gets stuck



# Optimization: Problem #2 with SGD

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

Saddle points much more common in high dimension



<https://blog.paperspace.com/intro-to-optimization-in-deep-learning-gradient-descent/>



# SGD + Momentum

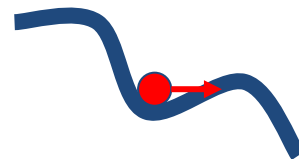
Intuitions:

- Think of a ball (set of parameters) moving in space (loss landscape), with momentum keeping it going in a direction.
- Individual gradient step may be noisy, the general trend accumulated over a few steps will point to the right direction.
- Momentum can “push” the ball over saddle points or local minima.

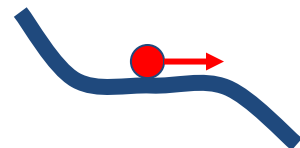
Noisy gradients



Local Minima



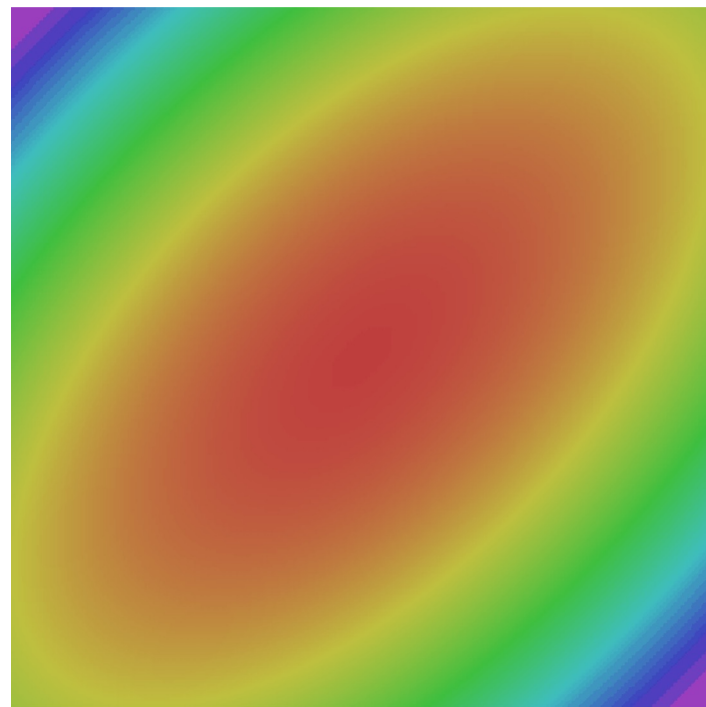
Saddle points



# SGD + Momentum

Intuitions:

- Think of a ball (set of parameters) moving in space (loss landscape), with momentum keeping it going in a direction.
- Individual gradient step may be noisy, the general trend accumulated over a few steps will point to the right direction.
- Momentum can “push” the ball over saddle points or local minima.



— SGD

— SGD+Momentum

# SGD: the simple two line update code

## SGD

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

```
while True:  
    dx = compute_gradient(x)  
    x -= learning_rate * dx
```

# SGD + Momentum:

continue moving in the general direction as the previous iterations

## SGD

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

```
while True:
    dx = compute_gradient(x)
    x -= learning_rate * dx
```

## SGD+Momentum

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

$$x_{t+1} = x_t - \alpha v_{t+1}$$

- Build up “velocity/momentum” 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

# SGD + Momentum:

continue moving in the general direction as the previous iterations

## SGD

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

```
while True:
    dx = compute_gradient(x)
    x -= learning_rate * dx
```

## SGD+Momentum

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

$$x_{t+1} = x_t - \alpha v_{t+1}$$

```
vx = 0
while True:
    dx = compute_gradient(x)
    vx = rho * vx + dx
    x -= learning_rate * vx
```

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

# SGD + Momentum:

alternative equivalent formulation

## SGD+Momentum

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

```
vx = 0
while True:
    dx = compute_gradient(x)
    vx = rho * vx - learning_rate * dx
    x += vx
```

## SGD+Momentum

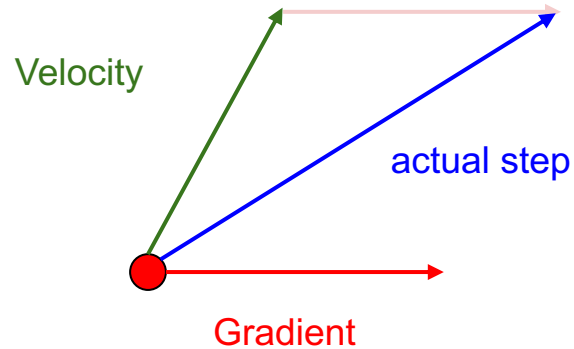
$$v_{t+1} = \rho v_t + \nabla f(x_t)$$
$$x_{t+1} = x_t - \alpha v_{t+1}$$

```
vx = 0
while True:
    dx = compute_gradient(x)
    vx = rho * vx + dx
    x -= learning_rate * vx
```

You may see SGD+Momentum formulated different ways,  
but they are equivalent - give same sequence of  $x$

# SGD+Momentum

Momentum update:



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

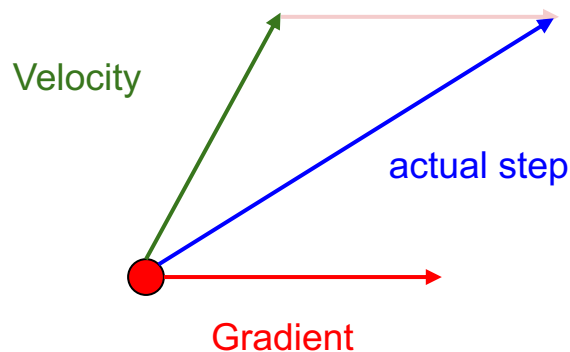
Nesterov, "A method of solving a convex programming problem with convergence rate  $O(1/k^2)$ ", 1983

Nesterov, "Introductory lectures on convex optimization: a basic course", 2004

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

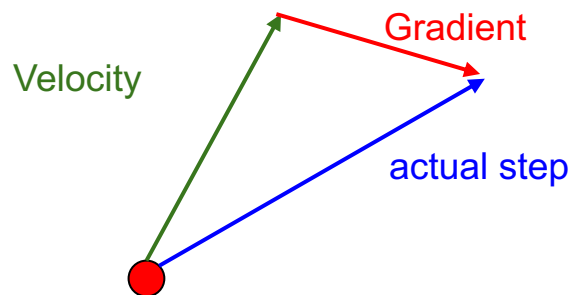
# Nesterov Momentum

Momentum update:



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

Nesterov Momentum

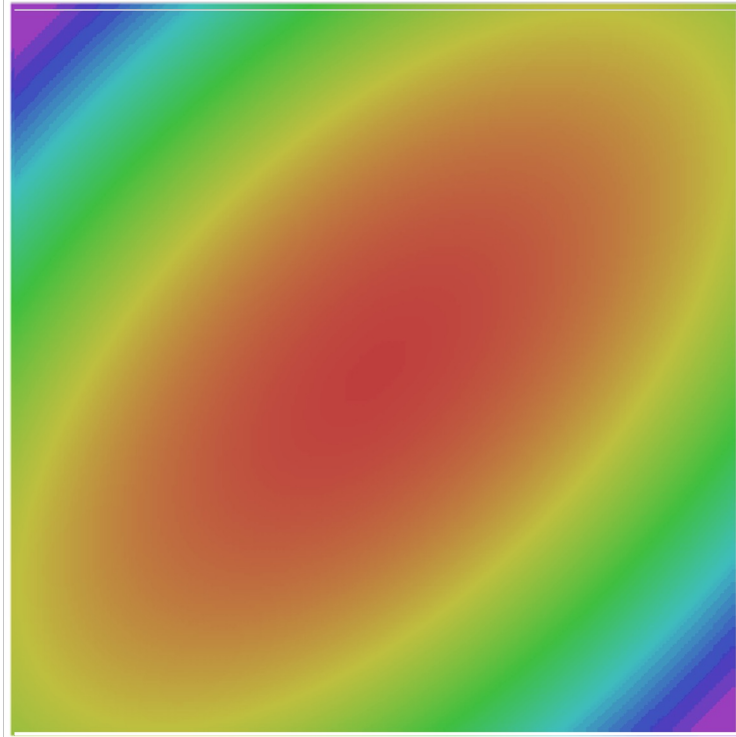


“Look ahead” to the point where updating using velocity would take us; compute gradient there and mix it with velocity to get actual update direction

Nesterov, “A method of solving a convex programming problem with convergence rate  $O(1/k^2)$ ”, 1983  
Nesterov, “Introductory lectures on convex optimization: a basic course”, 2004  
Sutskever et al, “On the importance of initialization and momentum in deep learning”, ICML 2013



# Nesterov Momentum

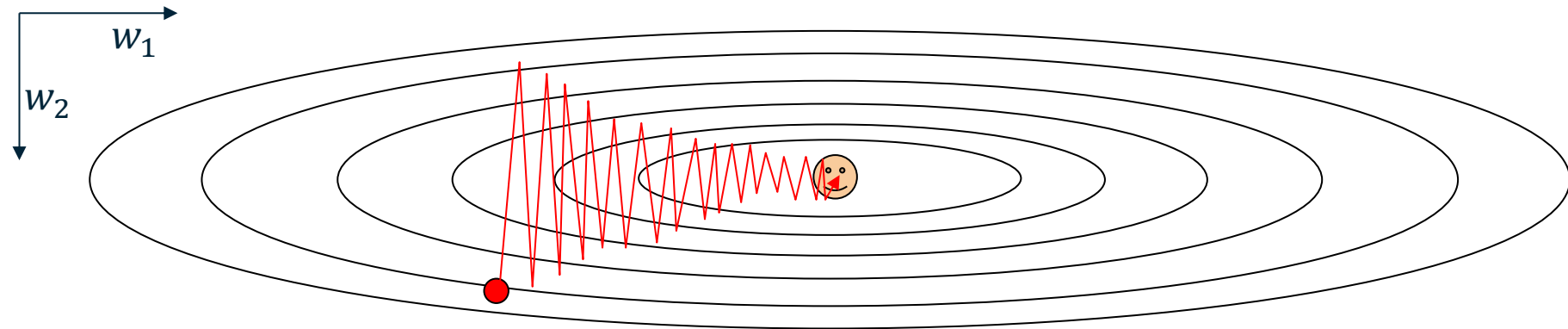


- SGD
- SGD+Momentum
- Nesterov

# Optimization: Problem #3 with SGD

What if loss changes quickly in one direction and slowly in another?  
What does gradient descent do?

Very slow progress along shallow dimension, jitter along steep direction



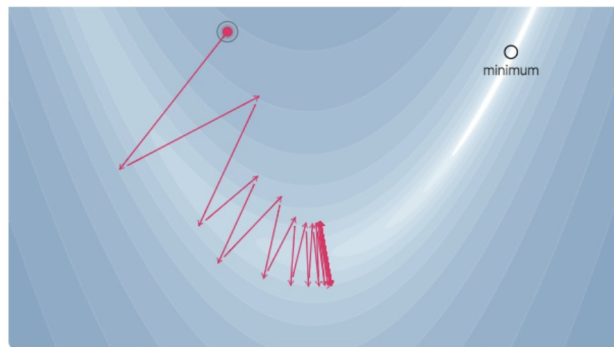
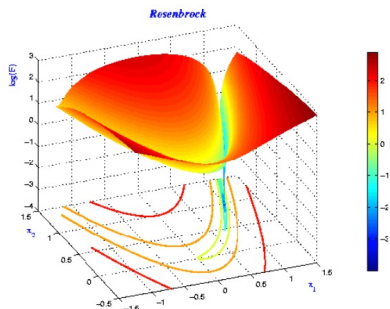
Assume each contour line has the same loss

# Optimization: Problem #3 with SGD

What if loss changes quickly in one direction and slowly in another?

Very slow progress along shallow dimension, jitter along steep direction

Long, narrow ravines:



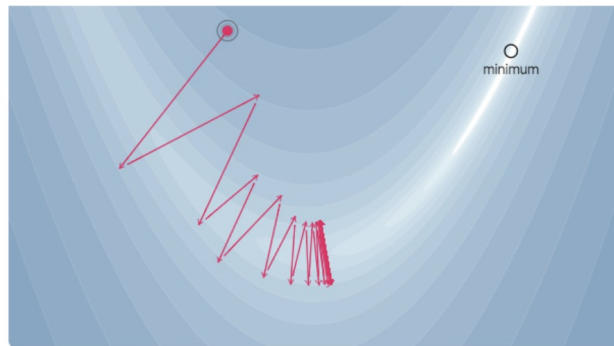
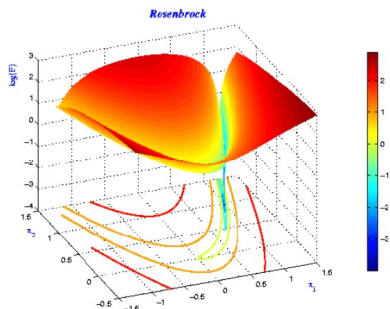
[https://www.cs.toronto.edu/~rgrosse/courses/csc421\\_2019/slides/lec07.pdf](https://www.cs.toronto.edu/~rgrosse/courses/csc421_2019/slides/lec07.pdf)

# Optimization: Problem #3 with SGD

What if loss changes quickly in one direction and slowly in another?

Very slow progress along shallow dimension, jitter along steep direction

Long, narrow ravines:



[https://www.cs.toronto.edu/~rgrosse/courses/csc421\\_2019/slides/lec07.pdf](https://www.cs.toronto.edu/~rgrosse/courses/csc421_2019/slides/lec07.pdf)

Loss function has high **condition number**: ratio of largest to smallest eigen value ( $\lambda_{max}/\lambda_{min}$ ) of the Hessian matrix of a loss function is large

Small condition number in loss Hessian -> circular contour

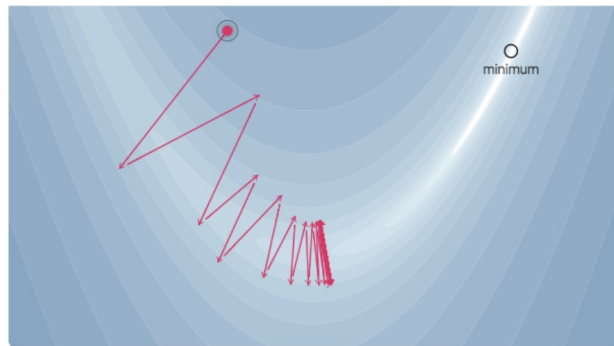
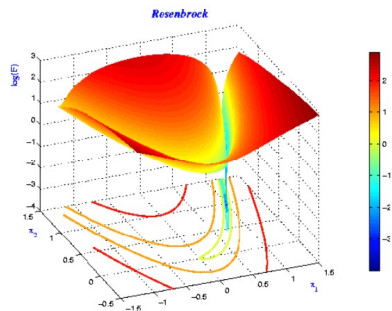
Large condition number in loss Hessian -> skewed contour

# Optimization: Problem #3 with SGD

What if loss changes quickly in one direction and slowly in another?

Very slow progress along shallow dimension, jitter along steep direction

Long, narrow ravines:



[https://www.cs.toronto.edu/~rgrosse/courses/csc421\\_2019/slides/lec07.pdf](https://www.cs.toronto.edu/~rgrosse/courses/csc421_2019/slides/lec07.pdf)

Loss function has high **condition number**: ratio of largest to smallest eigen value ( $\lambda_{max}/\lambda_{min}$ ) of the Hessian matrix of a loss function is large

Small condition number in loss Hessian -> circular contour

Large condition number in loss Hessian -> skewed contour

Can we enable SGD to adapt to this skew-ness?

# AdaGrad

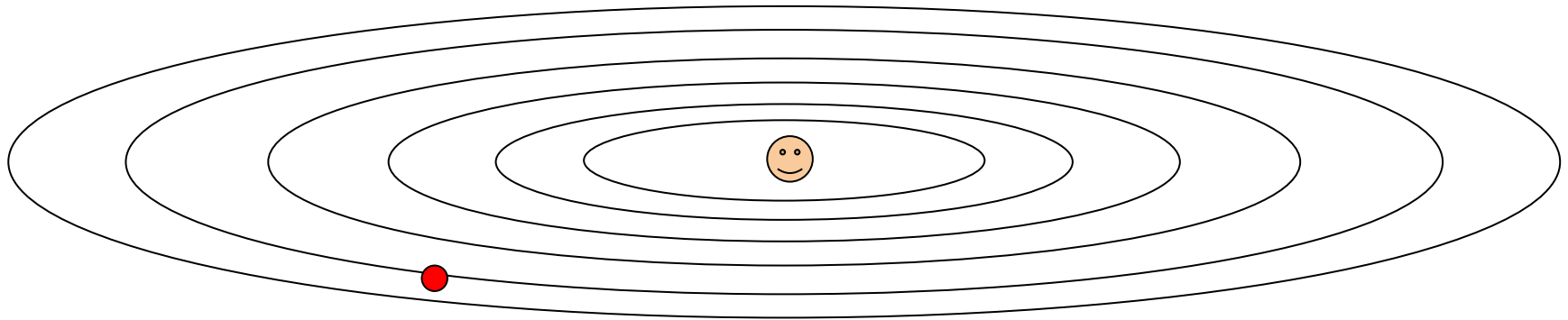
```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared += dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```

Added element-wise scaling of the gradient based on the historical sum of squares in each dimension

“Per-parameter learning rates”  
or “adaptive learning rates”

# AdaGrad

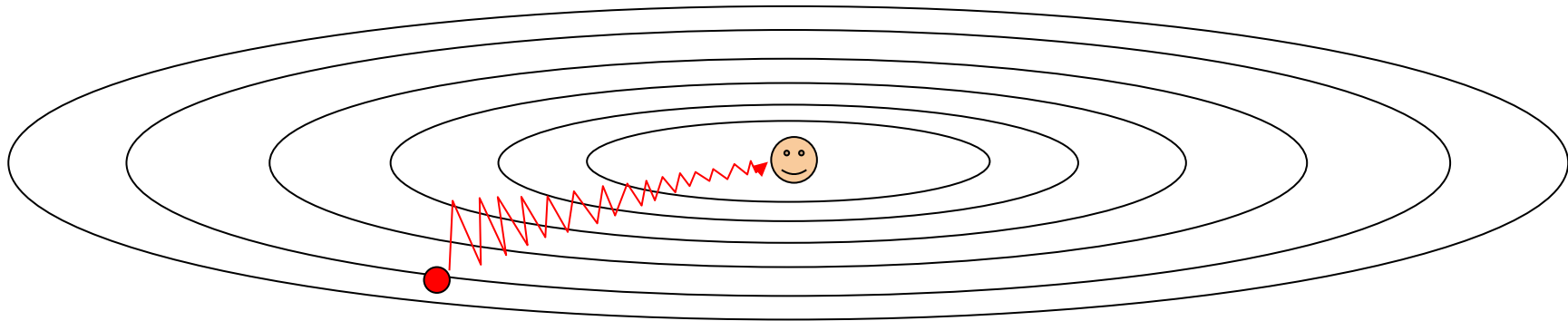
```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared += dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```



Q: What happens with AdaGrad?

# AdaGrad

```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared += dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```



Q: What happens with AdaGrad?

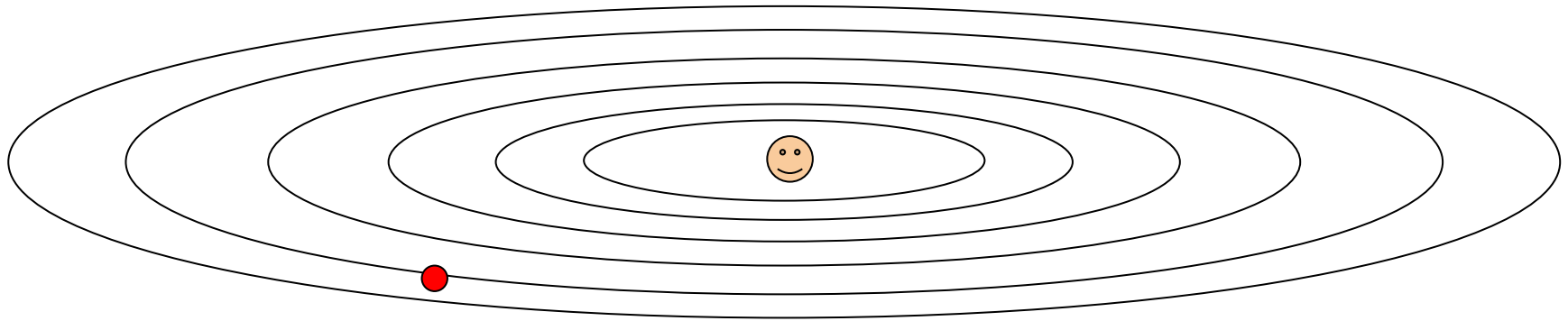
Progress along “steep” directions is damped;  
progress along “flat” directions is accelerated





# AdaGrad

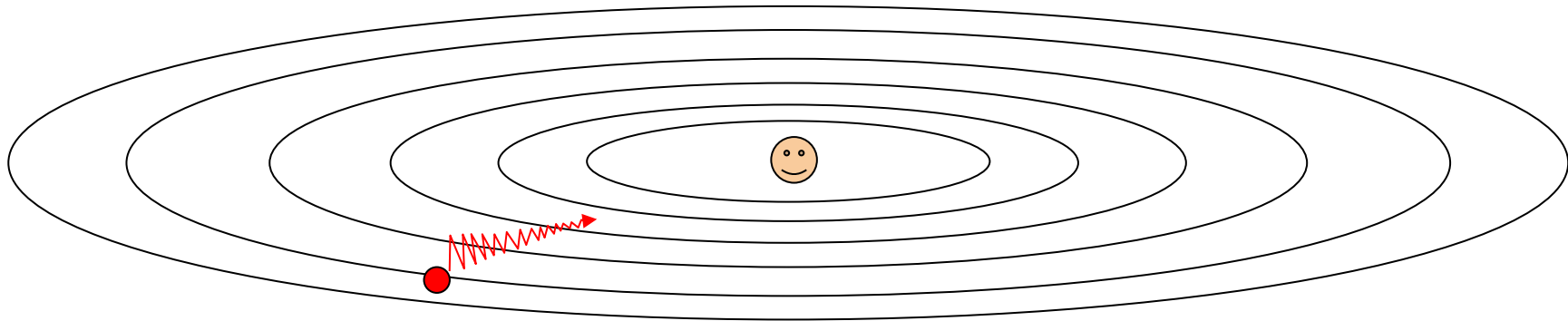
```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared += dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```



Q2: What happens to the step size over long time?

# AdaGrad

```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared += dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```



Q2: What happens to the step size over long time?

Decays to zero 😞

# RMSProp: “Leaky AdaGrad”

AdaGrad

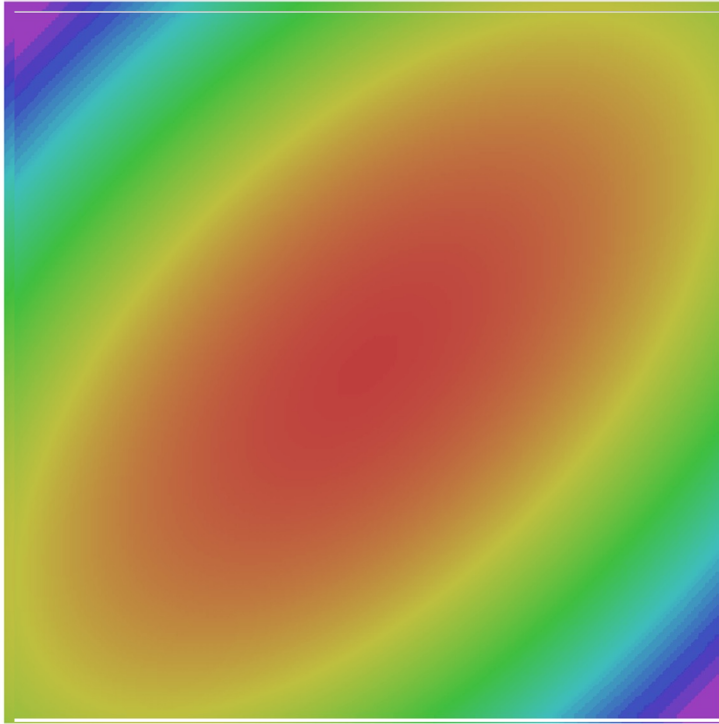
```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared += dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```



RMSProp

```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared = decay_rate * grad_squared + (1 - decay_rate) * dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```

# RMSProp



- SGD
- SGD+Momentum
- RMSProp
- AdaGrad  
(stuck due to decaying lr)

# Adam (almost)

```
first_moment = 0
second_moment = 0
while True:
    dx = compute_gradient(x)
    first_moment = beta1 * first_moment + (1 - beta1) * dx
    second_moment = beta2 * second_moment + (1 - beta2) * dx * dx
    x -= learning_rate * first_moment / (np.sqrt(second_moment) + 1e-7))
```

Typical hyperparams: beta1=0.9, beta2=0.999

# Adam (almost)

```
first_moment = 0
second_moment = 0
while True:
    dx = compute_gradient(x)
    first_moment = beta1 * first_moment + (1 - beta1) * dx
    second_moment = beta2 * second_moment + (1 - beta2) * dx * dx
    x -= learning_rate * first_moment / (np.sqrt(second_moment) + 1e-7))
```

Momentum

AdaGrad / RMSProp

Typical hyperparams: beta1=0.9, beta2=0.999

Sort of like RMSProp with momentum

Q: What happens at first timestep?

# Adam (almost)

```
first_moment = 0
second_moment = 0
while True:
    dx = compute_gradient(x)
    first_moment = beta1 * first_moment + (1 - beta1) * dx
    second_moment = beta2 * second_moment + (1 - beta2) * dx * dx
    x -= learning_rate * first_moment / (np.sqrt(second_moment) + 1e-7)
```

Momentum

AdaGrad / RMSProp

Typical hyperparams: beta1=0.9, beta2=0.999

Small -> divide by small number -> bad initial step

Q: What happens at first timestep?

# Adam (full form)

```
first_moment = 0
second_moment = 0
for t in range(1, num_iterations):
    dx = compute_gradient(x)
    first_moment = beta1 * first_moment + (1 - beta1) * dx
    second_moment = beta2 * second_moment + (1 - beta2) * dx * dx
    first_unbias = first_moment / (1 - beta1 ** t)
    second_unbias = second_moment / (1 - beta2 ** t)
    x -= learning_rate * first_unbias / (np.sqrt(second_unbias) + 1e-7)
```

Momentum

Bias correction

AdaGrad / RMSProp

Typical hyperparams: beta1=0.9, beta2=0.999

Bias correction for the fact that  
first and second moment  
estimates start at zero



# Adam (full form)

```
first_moment = 0
second_moment = 0
for t in range(1, num_iterations):
    dx = compute_gradient(x)
    first_moment = beta1 * first_moment + (1 - beta1) * dx
    second_moment = beta2 * second_moment + (1 - beta2) * dx * dx
    first_unbias = first_moment / (1 - beta1 ** t)
    second_unbias = second_moment / (1 - beta2 ** t)
    x -= learning_rate * first_unbias / (np.sqrt(second_unbias) + 1e-7))
```

Momentum

Bias correction

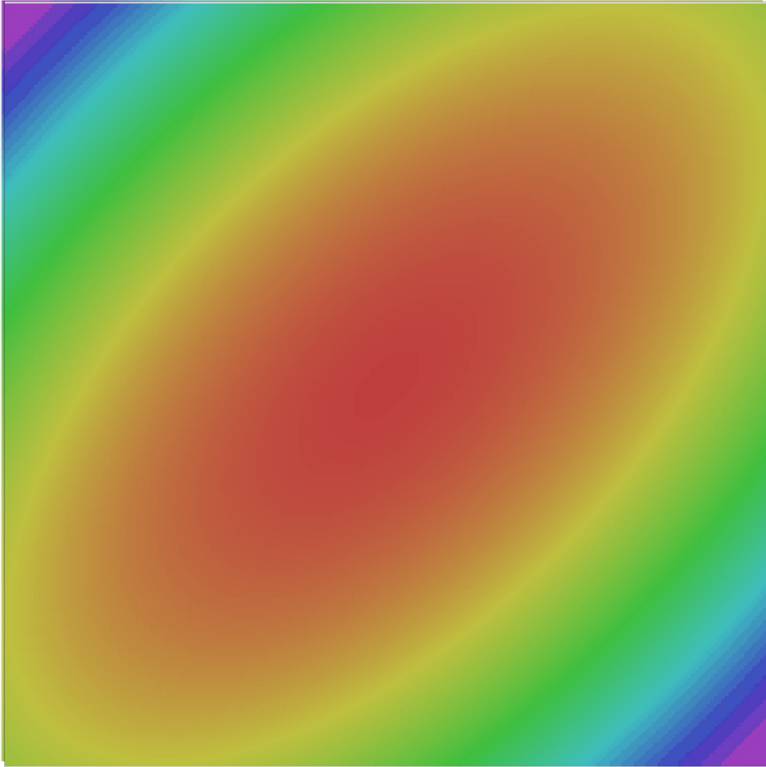
AdaGrad / RMSProp

Typical hyperparams: beta1=0.9, beta2=0.999

Bias correction for the fact that first and second moment estimates start at zero

Adam with **beta1 = 0.9**, **beta2 = 0.999**, and **learning\_rate = 1e-3 or 5e-4** is a great starting point for many models!

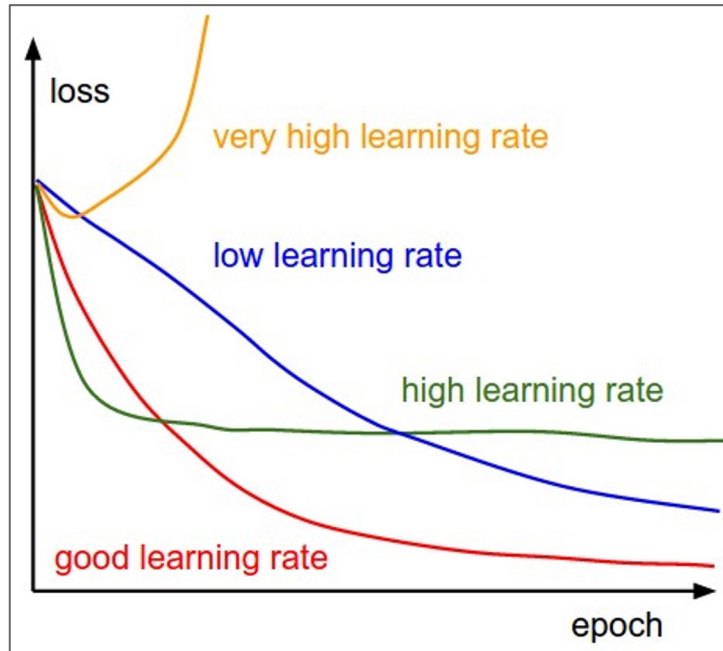
# Adam



- SGD
- SGD+Momentum
- RMSProp
- Adam

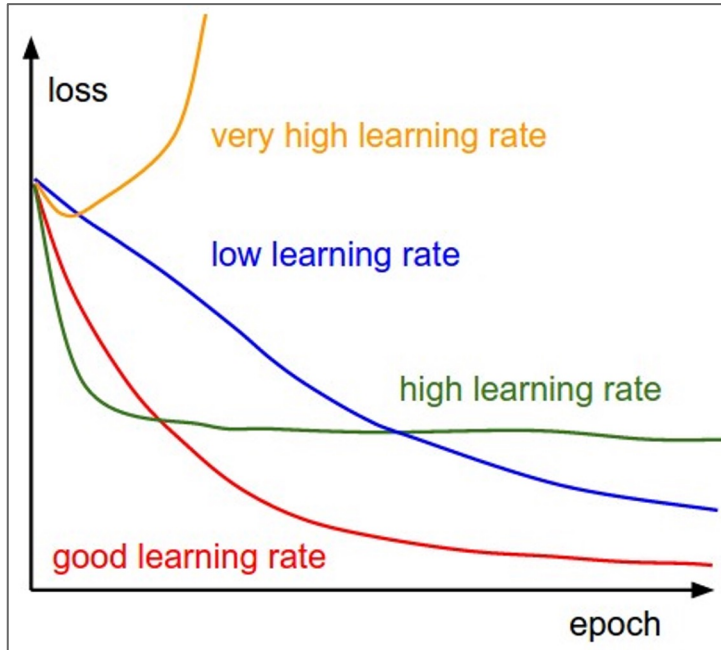
# Learning rate schedules

SGD, SGD+Momentum, Adagrad, RMSProp, Adam all have **learning rate** as a hyperparameter.



Q: Which one of these learning rates is best to use?

SGD, SGD+Momentum, Adagrad, RMSProp, Adam all have **learning rate** as a hyperparameter.

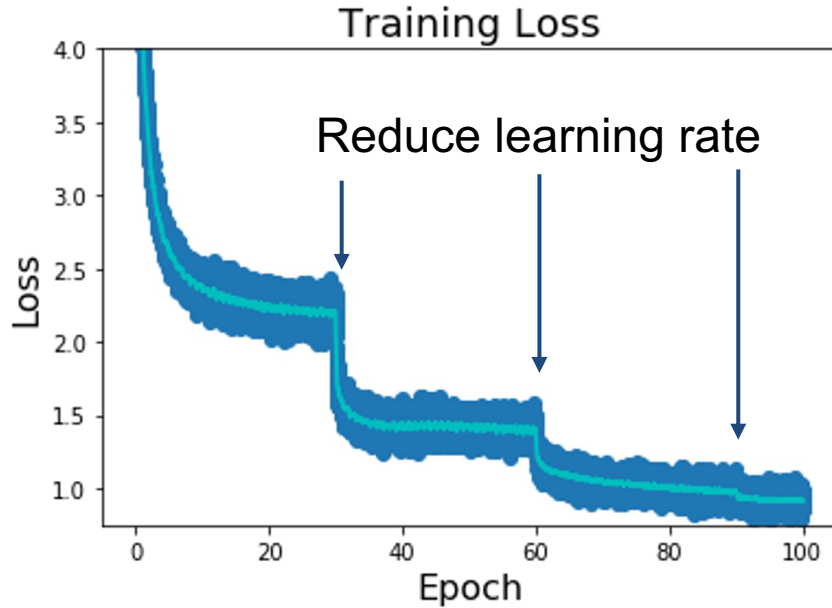


Q: Which one of these learning rates is best to use?

A: In reality, all of these are good learning rates.

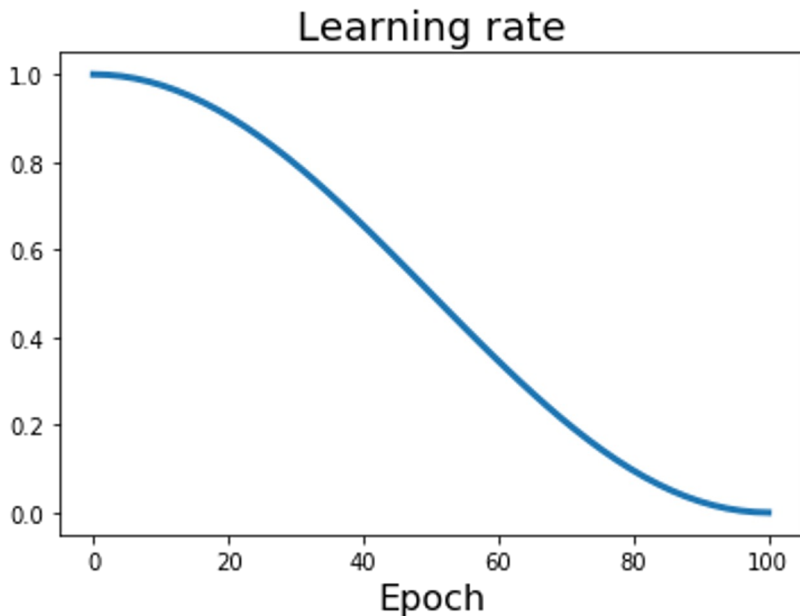
Need finer adjustment closer to convergence, so we want to reduce learning rate over time to keep making progress.

# Learning rate decays over time



**Step:** Reduce learning rate at a few fixed points. E.g. for ResNets, multiply LR by 0.1 after epochs 30, 60, and 90.

# Learning Rate Decay



**Step:** Reduce learning rate at a few fixed points. E.g. for ResNets, multiply LR by 0.1 after epochs 30, 60, and 90.

**Cosine:**  $\alpha_t = \frac{1}{2}\alpha_0 (1 + \cos(t\pi/T))$

$\alpha_0$  : Initial learning rate

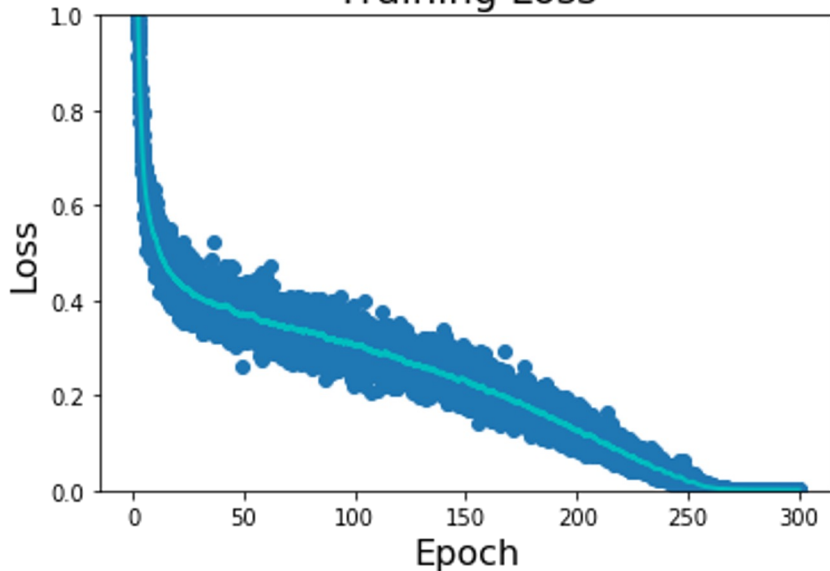
$\alpha_t$  : Learning rate at epoch  $t$

$T$  : Total number of epochs

Loshchilov and Hutter, “SGDR: Stochastic Gradient Descent with Warm Restarts”, ICLR 2017  
Radford et al, “Improving Language Understanding by Generative Pre-Training”, 2018  
Feichtenhofer et al, “SlowFast Networks for Video Recognition”, arXiv 2018  
Child et al, “Generating Long Sequences with Sparse Transformers”, arXiv 2019

# Learning Rate Decay

Training Loss



**Step:** Reduce learning rate at a few fixed points. E.g. for ResNets, multiply LR by 0.1 after epochs 30, 60, and 90.

**Cosine:**  $\alpha_t = \frac{1}{2}\alpha_0 (1 + \cos(t\pi/T))$

$\alpha_0$  : Initial learning rate

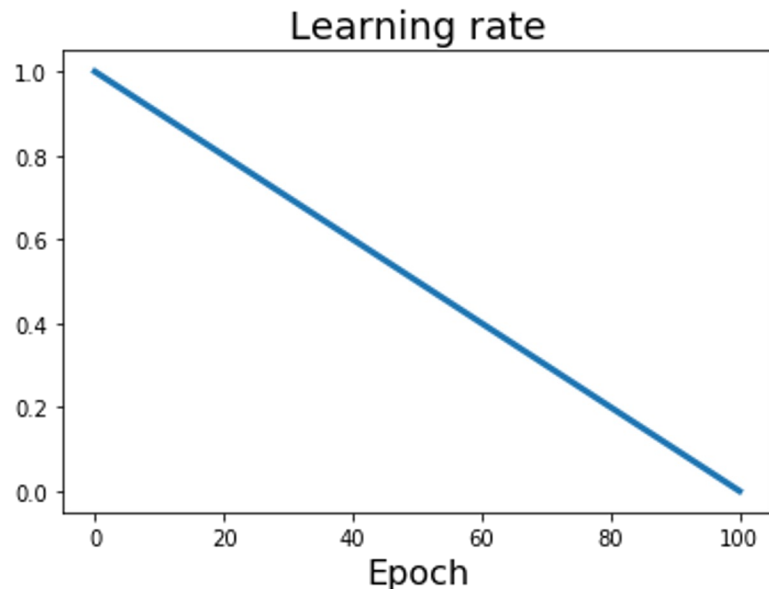
$\alpha_t$  : Learning rate at epoch  $t$

$T$  : Total number of epochs

Loshchilov and Hutter, "SGDR: Stochastic Gradient Descent with Warm Restarts", ICLR 2017  
Radford et al, "Improving Language Understanding by Generative Pre-Training", 2018  
Feichtenhofer et al, "SlowFast Networks for Video Recognition", arXiv 2018  
Child et al, "Generating Long Sequences with Sparse Transformers", arXiv 2019



# Learning Rate Decay



**Step:** Reduce learning rate at a few fixed points. E.g. for ResNets, multiply LR by 0.1 after epochs 30, 60, and 90.

**Cosine:**  $\alpha_t = \frac{1}{2}\alpha_0 (1 + \cos(t\pi/T))$

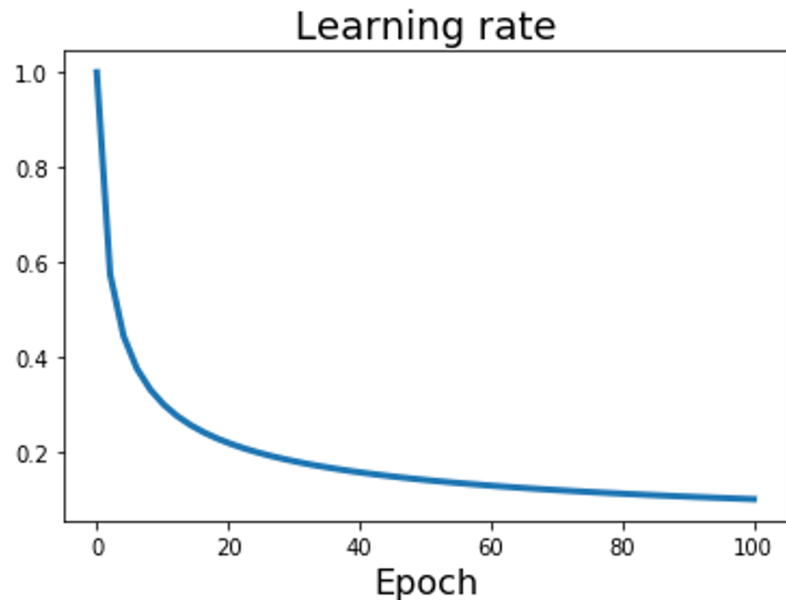
**Linear:**  $\alpha_t = \alpha_0(1 - t/T)$

$\alpha_0$  : Initial learning rate

$\alpha_t$  : Learning rate at epoch  $t$

$T$  : Total number of epochs

# Learning Rate Decay



**Step:** Reduce learning rate at a few fixed points. E.g. for ResNets, multiply LR by 0.1 after epochs 30, 60, and 90.

**Cosine:**  $\alpha_t = \frac{1}{2}\alpha_0 (1 + \cos(t\pi/T))$

**Linear:**  $\alpha_t = \alpha_0(1 - t/T)$

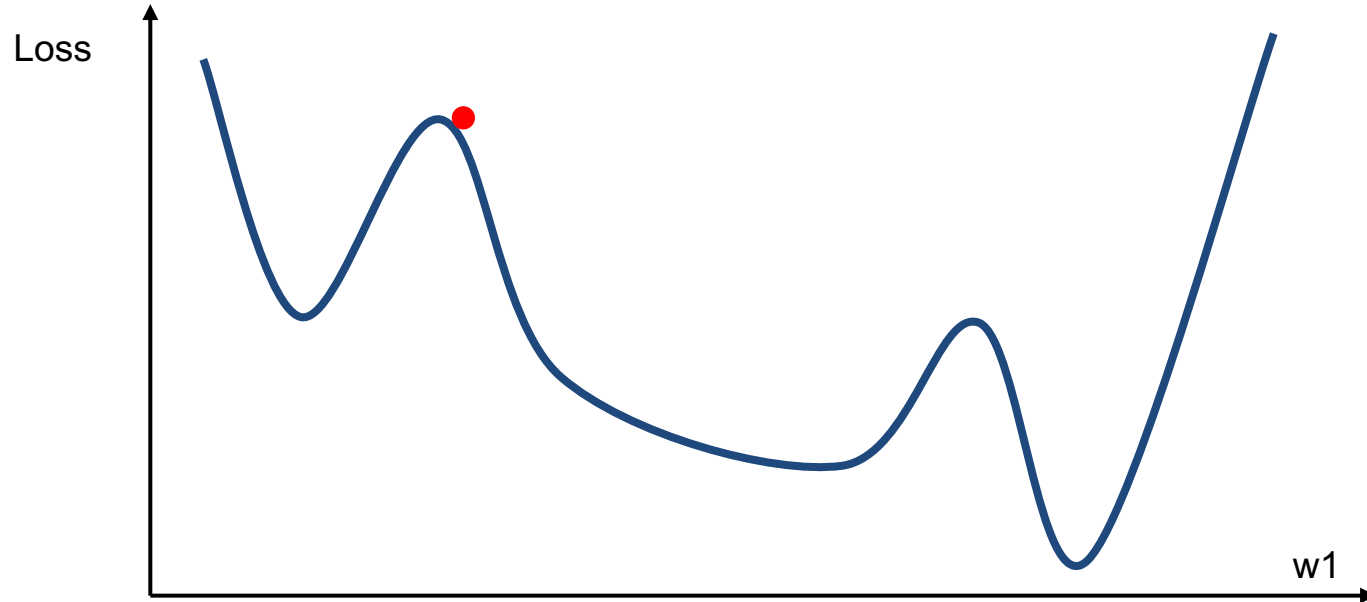
**Inverse sqrt:**  $\alpha_t = \alpha_0/\sqrt{t}$

$\alpha_0$  : Initial learning rate

$\alpha_t$  : Learning rate at epoch  $t$

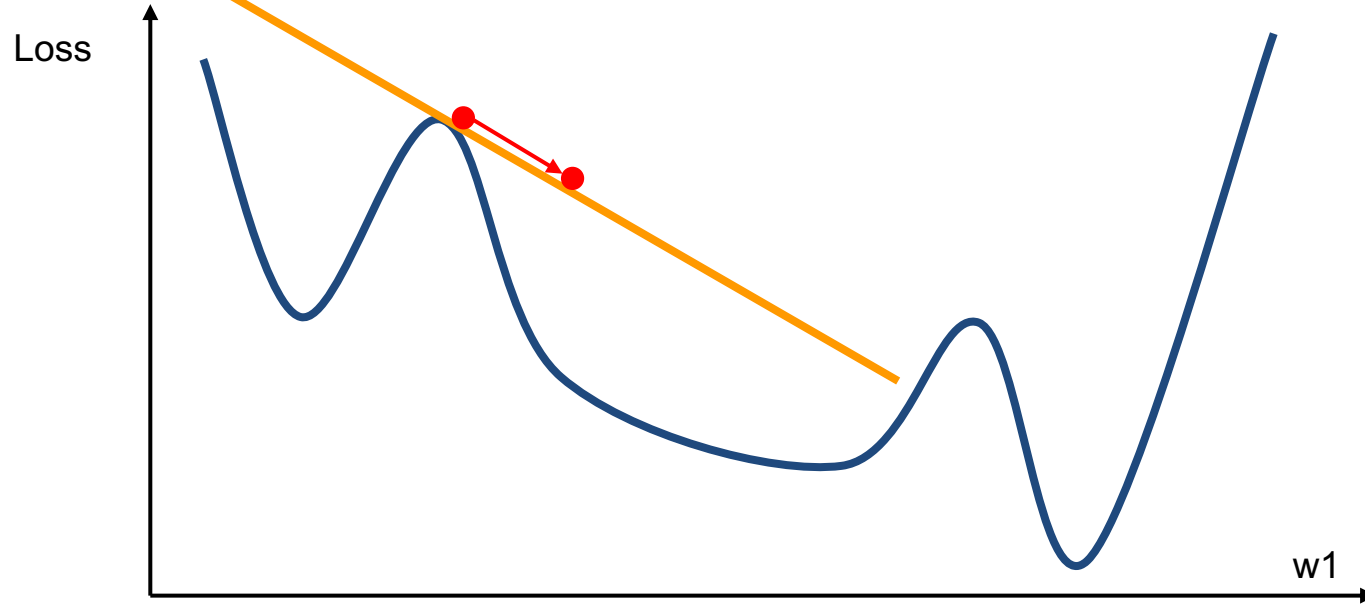
$T$  : Total number of epochs

# First-Order Optimization



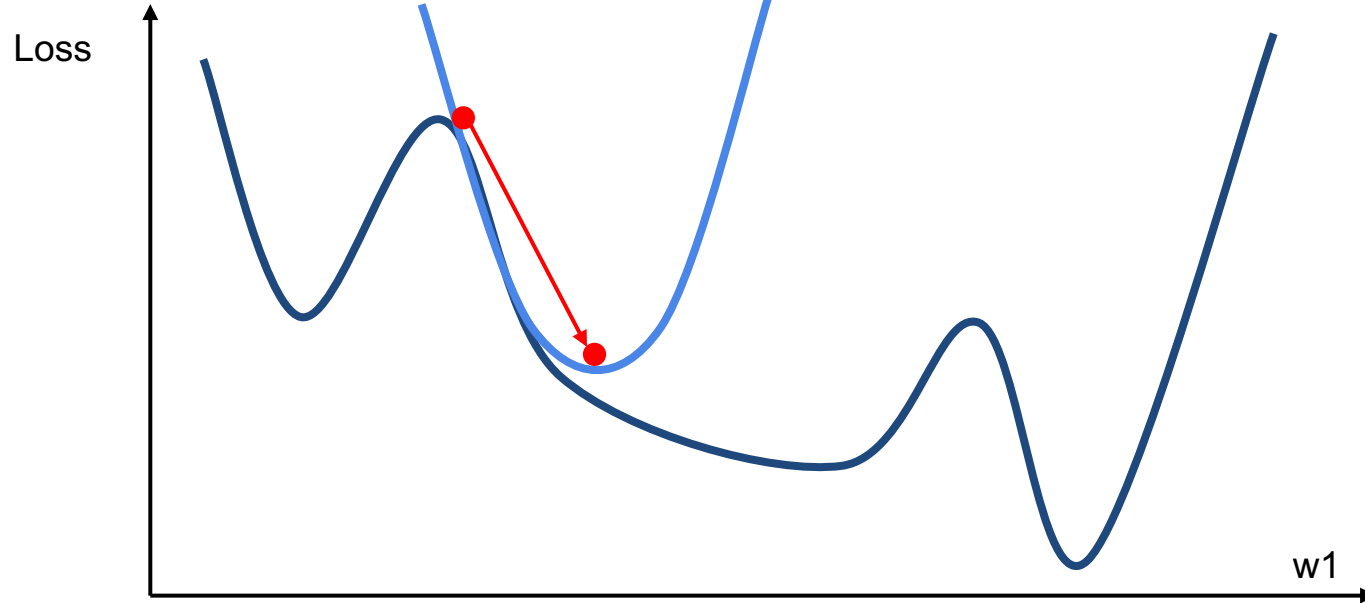
# First-Order Optimization

- (1) Use gradient form linear approximation
- (2) Step to minimize the approximation



# Second-Order Optimization

- (1) Use gradient **and Hessian** to form **quadratic** approximation
- (2) Step to the **minima** of the approximation



# Second-Order Optimization

second-order Taylor expansion:

$$J(\boldsymbol{\theta}) \approx J(\boldsymbol{\theta}_0) + (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_0) + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \mathbf{H} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)$$

Solving for the critical point we obtain the Newton parameter update:

$$\boldsymbol{\theta}^* = \boldsymbol{\theta}_0 - \mathbf{H}^{-1} \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_0)$$

Q: Why is this bad for deep learning?

# Second-Order Optimization

second-order Taylor expansion:

$$J(\boldsymbol{\theta}) \approx J(\boldsymbol{\theta}_0) + (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_0) + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \mathbf{H} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)$$

Solving for the critical point we obtain the Newton parameter update:

$$\boldsymbol{\theta}^* = \boldsymbol{\theta}_0 - \mathbf{H}^{-1} \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_0)$$

Hessian has  $O(N^2)$  elements  
Inverting takes  $O(N^3)$   
 $N = \text{Millions}$

Q: Why is this bad for deep learning?

$$\mathbf{H}_f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix},$$

# Second-Order Optimization

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0)$$

- Quasi-Newton methods (**BGFS** most popular):  
*instead of inverting the Hessian ( $O(n^3)$ ), approximate inverse Hessian with rank 1 updates over time ( $O(n^2)$  each).*
- **L-BFGS** (Limited memory BFGS):  
*Does not form/store the full inverse Hessian.*



# L-BFGS

- **Usually works very well in full batch, deterministic mode** i.e. if you have a single, deterministic  $f(x)$  then L-BFGS will probably work very nicely
- **Does not transfer very well to mini-batch setting.** Gives bad results. Adapting second-order methods to large-scale, stochastic setting is an active area of research.

Le et al, "On optimization methods for deep learning, ICML 2011"

Ba et al, "Distributed second-order optimization using Kronecker-factored approximations", ICLR 2017

# In practice:

- **Adam** is a good default choice in many cases; it often works ok even with constant learning rate
- **SGD+Momentum** can outperform Adam but may require more tuning of LR and schedule
  - Try cosine schedule, very few hyperparameters!
- If you can afford to do full batch updates (very rare for deep learning applications) then try out **L-BFGS** (and don't forget to disable all sources of noise)

# Next Time:

## **Training** Deep Neural Networks

- Details of the non-linear activation functions
- Data normalization
- Weight Initialization
- Batch Normalization
- Advanced Optimization
- Regularization
- Data Augmentation
- Transfer learning
- Hyperparameter Tuning
- Model Ensemble