# 14.332.435/16.332.530
# Introduction to Deep Learning
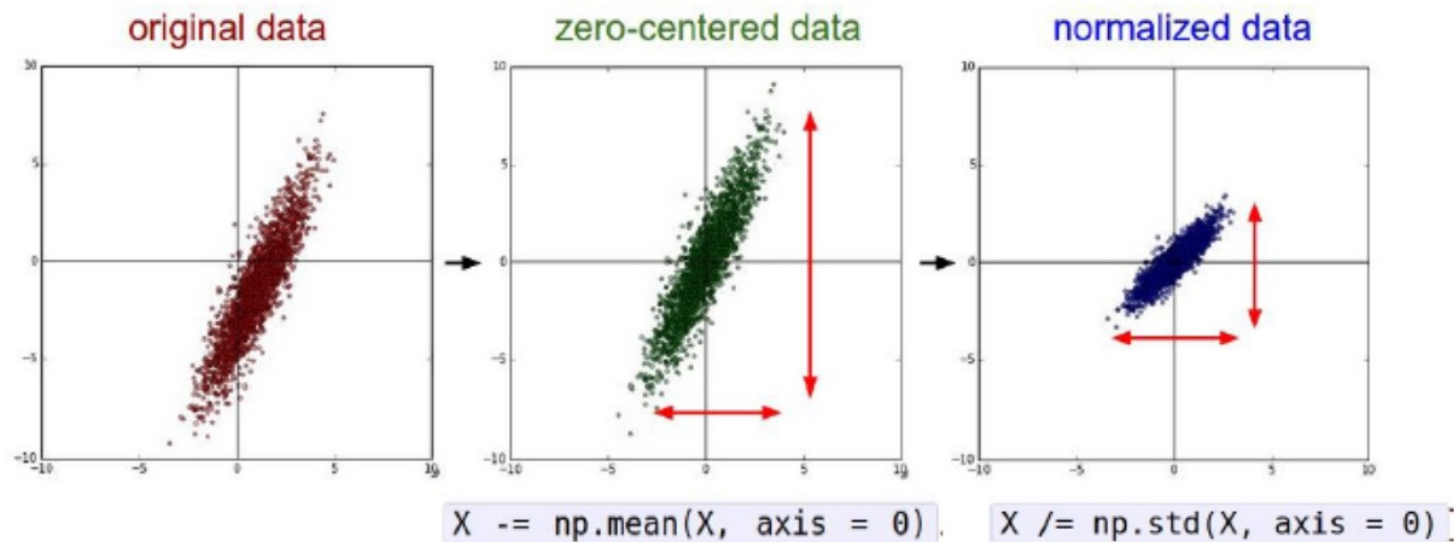
# Lecture 9
# Training Neural Network 2

## Yuqian Zhang

**Department of Electrical and Computer Engineering**
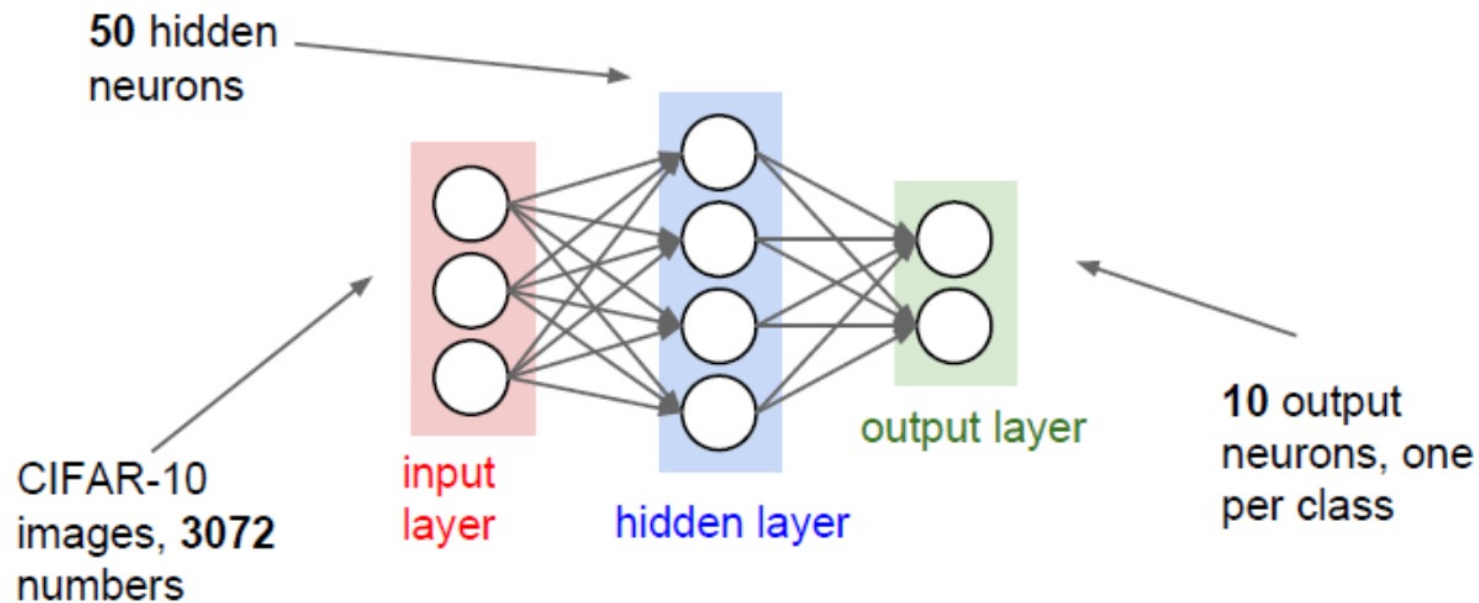
# Babysitting Learning Process

- Step 1: Data Preprocessing



original data     zero-centered data     normalized data

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

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

# Babysitting Learning Process

- Step 2: Choose Architecture

# Babysitting Learning Process

- Step 3a: Initial check the loss

```python
def init_two_layer_model(input_size, hidden_size, output_size):
    # initialize a model
    model = {}
    model['W1'] = 0.0001 * np.random.randn(input_size, hidden_size)
    model['b1'] = np.zeros(hidden_size)
    model['W2'] = 0.0001 * np.random.randn(hidden_size, output_size)
    model['b2'] = np.zeros(output_size)
    return model
```

```python
model = init_two_layer_model(32*32*3, 50, 10) # input size, hidden size, number of classes
loss, grad = two_layer_net(X_train, model, y_train, 0.0)
print loss
```

disable regularization

2.30261216167

loss ~2.3.
"correct" for
10 classes

returns the loss and the
gradient for all parameters

# Babysitting Learning Process

- Step 3b: Initial check the loss

```python
def init_two_layer_model(input_size, hidden_size, output_size):
    # initialize a model
    model = {}
    model['W1'] = 0.0001 * np.random.randn(input_size, hidden_size)
    model['b1'] = np.zeros(hidden_size)
    model['W2'] = 0.0001 * np.random.randn(hidden_size, output_size)
    model['b2'] = np.zeros(output_size)
    return model
```

```python
model = init_two_layer_model(32*32*3, 50, 10) # input_size, hidden size, number of classes
loss, grad = two_layer_net(X_train, model, y_train, 1e3)
print loss
```
crank up regularization

3.06859716482 ← loss went up, good. (sanity check)

# Babysitting Learning Process

- Step 4: Use very small data to exam train process

```
model = init_two_layer_model(32*32*3, 50, 10) # input size, hidden size, number of classes
trainer = ClassifierTrainer()
X_tiny = X_train[:20] # take 20 examples
y_tiny = y_train[:20]
best_model, stats = trainer.train(X_tiny, y_tiny, X_tiny, y_tiny,
                         model, two_layer_net,
                         num_epochs=200, reg=0.0,
                         update='sgd', learning_rate_decay=1,
                         sample_batches = False,
                         learning_rate=1e-3, verbose=True)
```

**Use very small data**

**Turn off regularization**

**Use simple optimizer**

```
Finished epoch 1 / 200: cost 2.302603, train: 0.400000, val 0.400000, lr 1.000000e-03
Finished epoch 2 / 200: cost 2.302258, train: 0.450000, val 0.450000, lr 1.000000e-03
Finished epoch 3 / 200: cost 2.301849, train: 0.600000, val 0.600000, lr 1.000000e-03
Finished epoch 4 / 200: cost 2.301196, train: 0.650000, val 0.650000, lr 1.000000e-03
Finished epoch 5 / 200: cost 2.300044, train: 0.650000, val 0.650000, lr 1.000000e-03
Finished epoch 6 / 200: cost 2.297864, train: 0.550000, val 0.550000, lr 1.000000c-03
Finished epoch 7 / 200: cost 2.293595, train: 0.600000, val 0.600000, lr 1.000000e-03
Finished epoch 8 / 200: cost 2.285096, train: 0.550000, val 0.550000, lr 1.000000e-03
Finished epoch 9 / 200: cost 2.268094, train: 0.550000, val 0.550000, lr 1.000000c-03
Finished epoch 10 / 200: cost 2.234787, train: 0.500000, val 0.500000, lr 1.000000e-03
Finished epoch 11 / 200: cost 2.173187, train: 0.500000, val 0.500000, lr 1.000000e-03
Finished epoch 12 / 200: cost 2.076862, train: 0.500000, val 0.500000, lr 1.000000e-03
Finished epoch 13 / 200: cost 1.974090, train: 0.400000, val 0.400000, lr 1.000000e-03
Finished epoch 14 / 200: cost 1.895885, train: 0.400000, val 0.400000, lr 1.000000e-03
Finished epoch 15 / 200: cost 1.820876, train: 0.450000, val 0.450000, lr 1.000000c-03
Finished epoch 16 / 200: cost 1.737430, train: 0.450000, val 0.450000, lr 1.000000e-03
Finished epoch 17 / 200: cost 1.642356, train: 0.500000, val 0.500000, lr 1.000000e-03
Finished epoch 18 / 200: cost 1.535239, train: 0.600000, val 0.600000, lr 1.000000e-03
Finished epoch 19 / 200: cost 1.421527, train: 0.600000, val 0.600000, lr 1.000000e-03

Finished epoch 195 / 200: cost 0.002694, train: 1.000000, val 1.000000, lr 1.000000e-03
Finished epoch 196 / 200: cost 0.002674, train: 1.000000, val 1.000000, lr 1.000000e-03
Finished epoch 197 / 200: cost 0.002655, train: 1.000000, val 1.000000, lr 1.000000e-03
Finished epoch 198 / 200: cost 0.002635, train: 1.000000, val 1.000000, lr 1.000000e-03
Finished epoch 199 / 200: cost 0.002617, train: 1.000000, val 1.000000, lr 1.000000e-03
```

**Make sure for very small portion we can overfit**

# Babysitting Learning Process

- Step 4: Begin to train via trying small regularization and learning rate to reduce loss

```
model = init_two_layer_model(32*32*3, 50, 10) # input size, hidden size, number of classes
trainer = ClassifierTrainer()
best_model, stats = trainer.train(X_train, y_train, X_val, y_val,
                        model, two_layer_net,
                        num_epochs=10, reg=0.000001,
                        update='sgd', learning_rate_decay=1,
                        sample_batches = True,
                        learning_rate=1e-6, verbose=True)
```
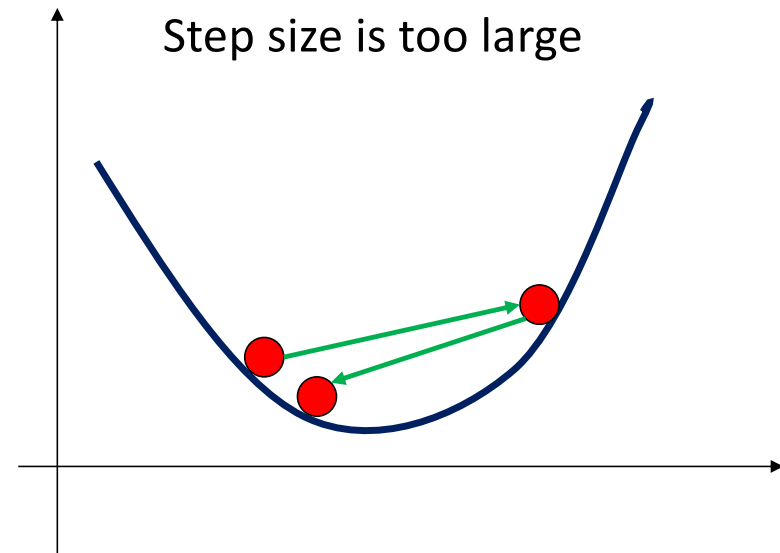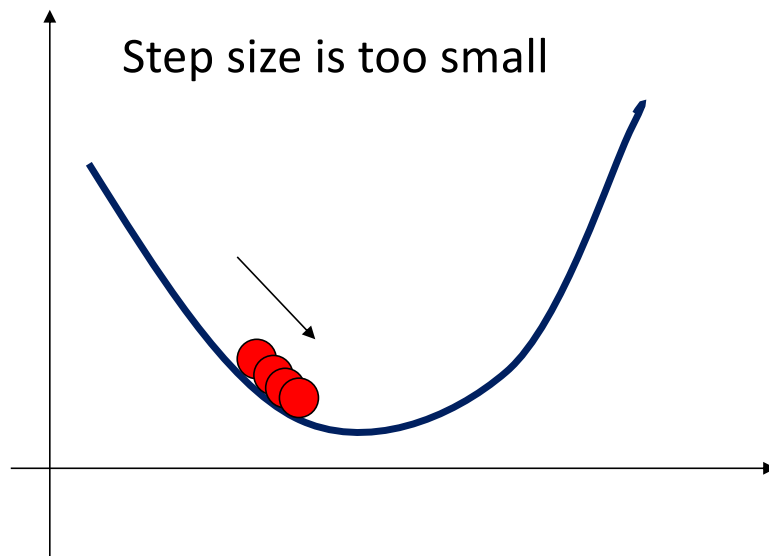
```
Finished epoch 1 / 10: cost 2.302576, train: 0.080000, val 0.103000, lr 1.000000e-06
Finished epoch 2 / 10: cost 2.302582, train: 0.121000, val 0.124000, lr 1.000000e-06
Finished epoch 3 / 10: cost 2.302558, train: 0.119000, val 0.138000, lr 1.000000e-06
Finished epoch 4 / 10: cost 2.302519, train: 0.127000, val 0.151000, lr 1.000000e-06
Finished epoch 5 / 10: cost 2.302517, train: 0.158000, val 0.171000, lr 1.000000e-06
Finished epoch 6 / 10: cost 2.302518, train: 0.179000, val 0.172000, lr 1.000000e-06
Finished epoch 7 / 10: cost 2.302466, train: 0.180000, val 0.176000, lr 1.000000e-06
Finished epoch 8 / 10: cost 2.302452, train: 0.175000, val 0.185000, lr 1.000000e-06
Finished epoch 9 / 10: cost 2.302459, train: 0.206000, val 0.192000, lr 1.000000e-06
Finished epoch 10 / 10: cost 2.302420, train: 0.190000, val 0.192000, lr 1.000000e-06
finished optimization. best validation accuracy: 0.192000
```

**Too small learning rate makes loss barely down**

# Babysitting Learning Process

- Step 4: Begin to train via trying small regularization and learning rate to reduce loss



```
model = init_two_layer_model(32*32*3, 50, 10) # input size, hidden size, number of classes
trainer = ClassifierTrainer()
best_model, stats = trainer.train(X_train, y_train, X_val, y_val,
                                  model, two_layer_net,
                                  num_epochs=10, reg=0.000001,
                                  update='sgd', learning_rate_decay=1,
                                  sample_batches = True,
                                  learning_rate=1e6, verbose=True)
```

```
/home/karpathy/cs231n/code/cs231n/classifiers/neural_net.py:50: RuntimeWarning: divide by zero en
countered in log
  data_loss = -np.sum(np.log(probs[range(N), y])) / N
/home/karpathy/cs231n/code/cs231n/classifiers/neural_net.py:48: RuntimeWarning: invalid value enc
ountered in subtract
  probs = np.exp(scores - np.max(scores, axis=1, keepdims=True))
Finished epoch 1 / 10: cost nan, train: 0.091000, val 0.087000, lr 1.000000e+06
Finished epoch 2 / 10: cost nan, train: 0.095000, val 0.087000, lr 1.000000e+06
Finished epoch 3 / 10: cost nan, train: 0.100000, val 0.087000, lr 1.000000e+06
```

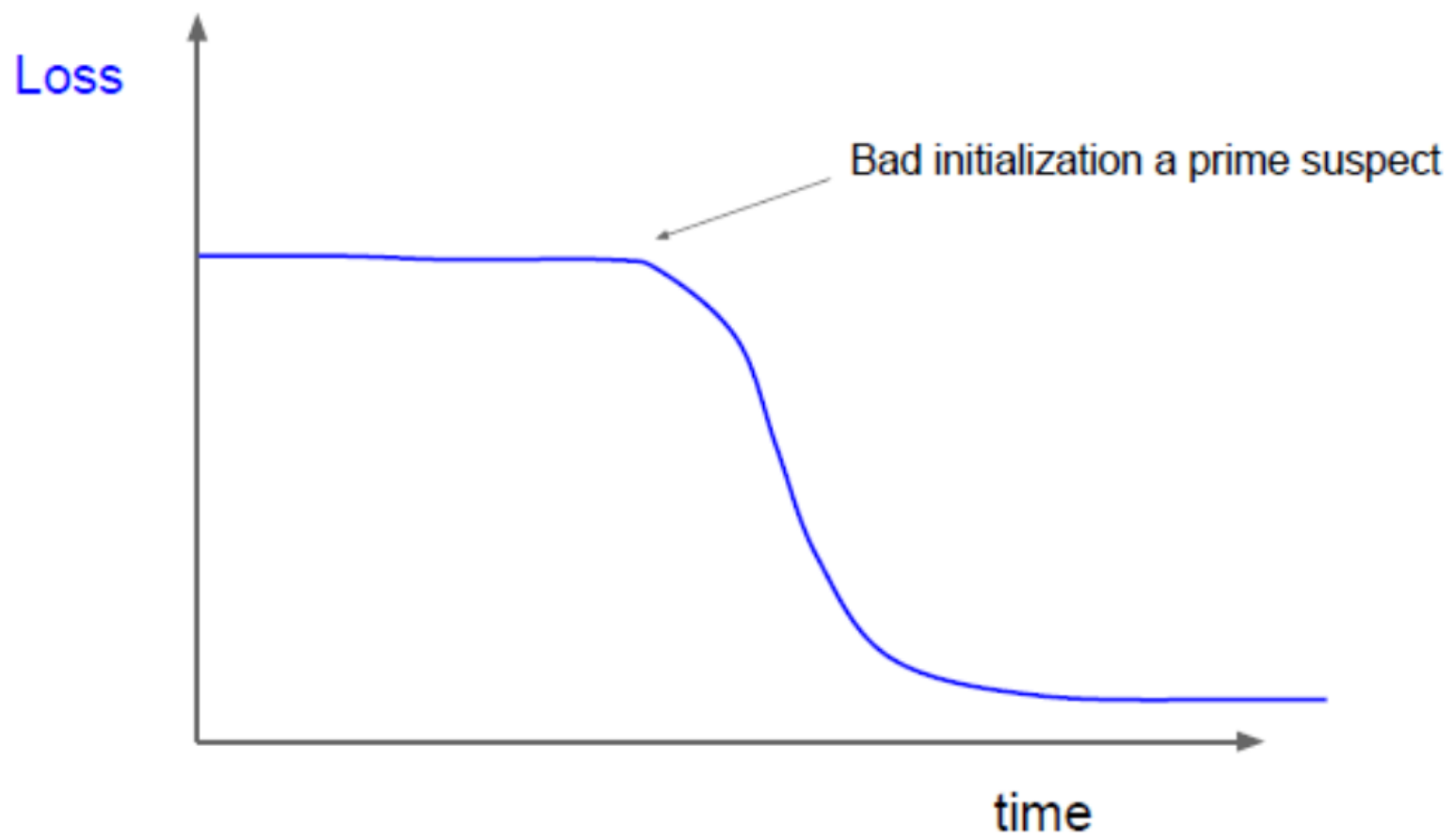**Too high learning rate makes loss explode**

# Proper Step Size is Important

- Too small: Long training time
- Too big: Skip optimal point (hard to converge)
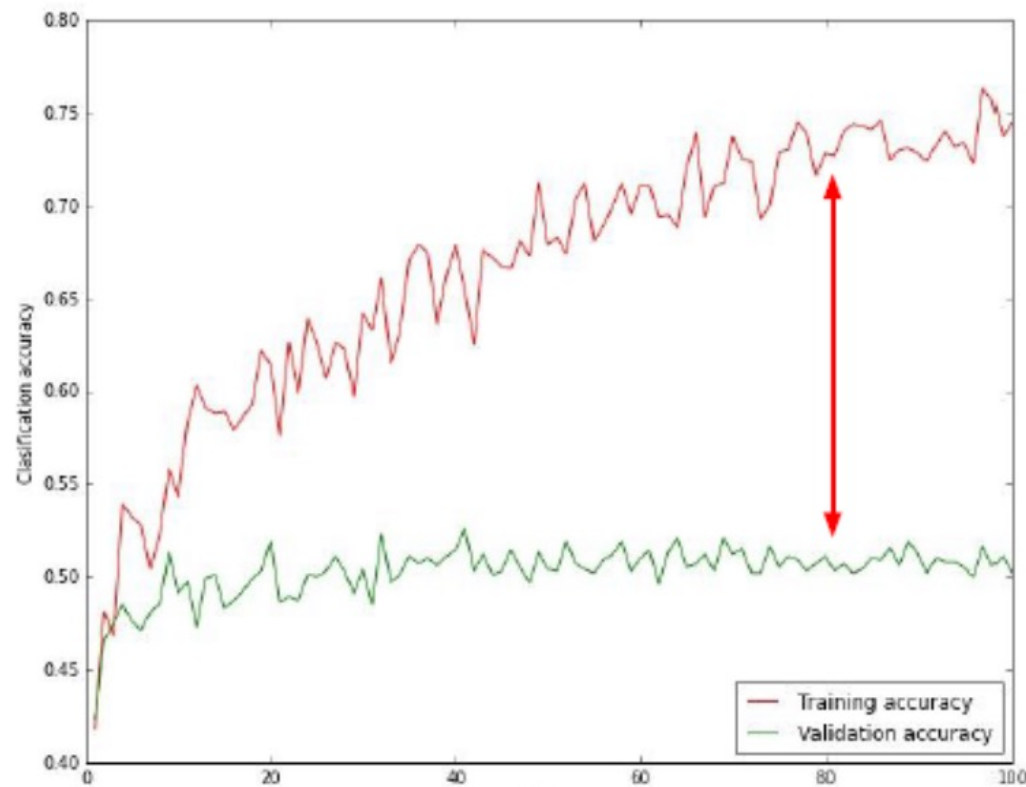
Step size is too small

Step size is too large

# Notes on Training

- Typical rang of learning rate

  -- within the range [1e-3 … 1e-5]

- How to determine the end of training

  -- set maximum iteration

  -- observe convergence of loss

  -- observe validation error rate

# Overfit and Underfit
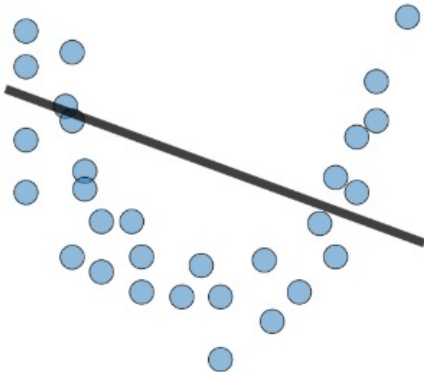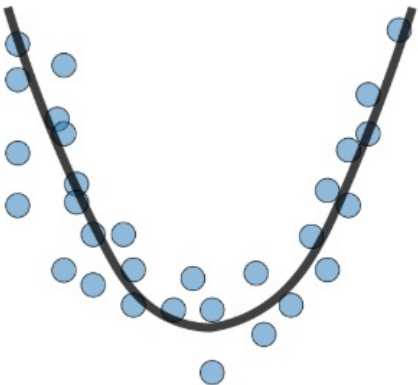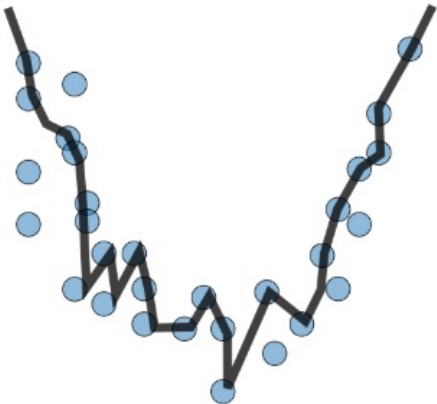


big gap = overfitting
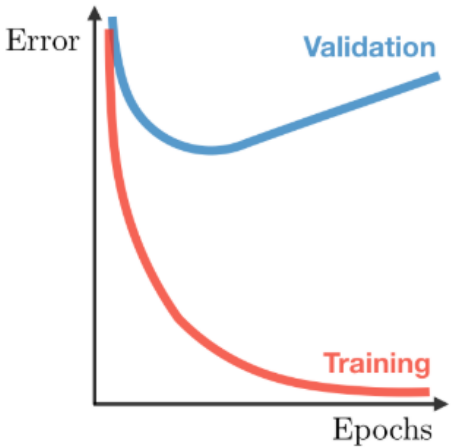=> increase regularization strength?

no gap
=> increase model capacity?

# Bias and Variance

- The *bias* is an error from erroneous assumptions in the learning algorithm. High bias can cause an algorithm to miss the relevant relations between features and target outputs (underfitting).
- The *variance* is an error from sensitivity to small fluctuations in the training set. High variance can cause an algorithm to model the random noise in the training data, rather than the intended outputs (overfitting).

In statistics and machine learning, the **bias–variance tradeoff** is the property of a set of predictive models whereby models with a lower bias in parameter estimation have a higher variance of the parameter estimates across samples, and vice versa. The **bias–variance**
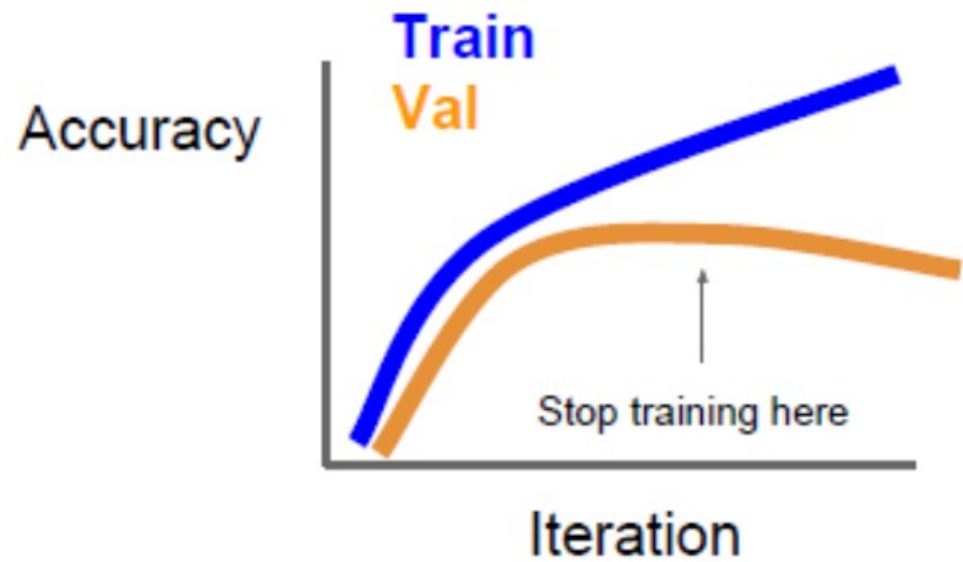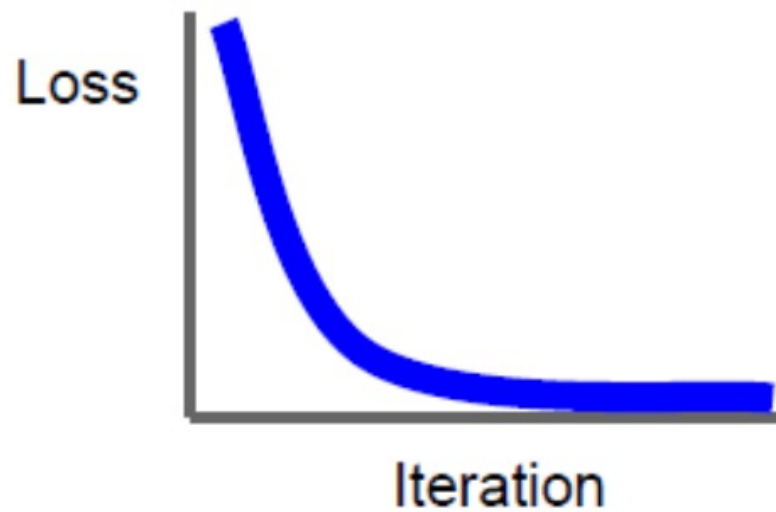
From wiki

| | Underfitting | Just right | Overfitting |
|---|---|---|---|
| **Symptoms** | • High training error<br>• Training error close to test error<br>• High bias | • Training error slightly lower than test error | • Very low training error<br>• Training error much lower than test error<br>• High variance |
| **Regression illustration** |  |  |  |

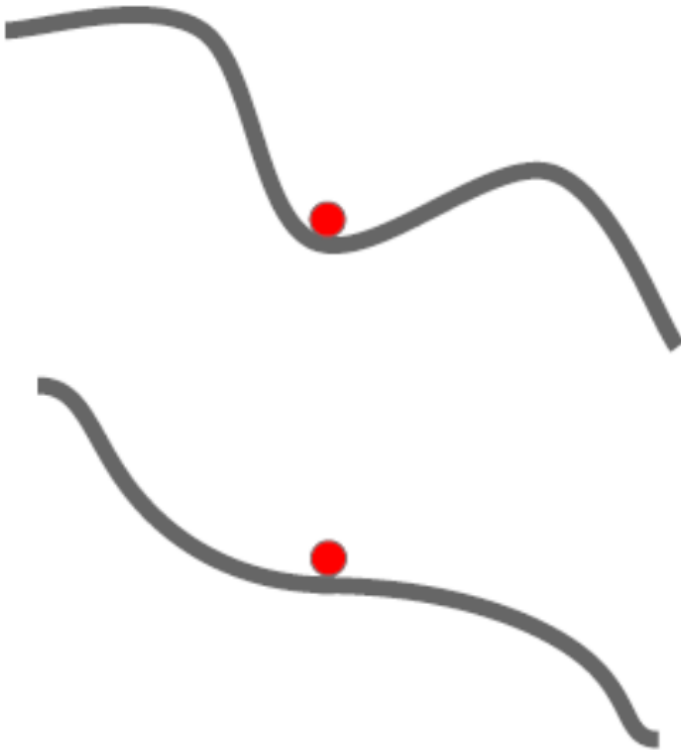| | Underfitting | Just right | Overfitting |
|---|---|---|---|
| **Symptoms** | • High training error<br>• Training error close to test error<br>• High bias | • Training error slightly lower than test error | • Very low training error<br>• Training error much lower than test error<br>• High variance |
| **Classification illustration** | | | |

| | Underfitting | Just right | Overfitting |
|---|---|---|---|
| **Symptoms** | • High training error<br>• Training error close to test error<br>• High bias | • Training error slightly lower than test error | • Very low training error<br>• Training error much lower than test error<br>• High variance |
| **Deep learning illustration** |  |  |  |
| **Possible remedies** | • Complexify model<br>• Add more features<br>• Train longer | | • Perform regularization<br>• Get more data |

https://stanford.edu/~shervine/teaching/cs-229/cheatsheet-machine-learning-tips-and-tricks#

# Early Stopping

# Problem of SGD



**Local Minima**

**Saddle Point**

Zero gradient, gradient descent gets stuck

# SGD with Momentum

Think an analogy

- Interpret loss as the height ($h$) of a hilly terrain
- Interpret weight initialization as setting a particle with zero initial velocity ($v$) at some location.
- Interpret optimization process as simulating the particle rolling on the landscape.
- Interpret the gradient as the force felt by the particle

    -- Because $U = mgh$    $F = -\nabla U$

# SGD with Momentum

Now consider $F = ma$  $U = mgh$  $F = -\nabla U$

- Gradient can be viewed propositional the acceleration ($a$) of the particle.

- Conventional SGD gradient directly integrates the position.

- This analogy implies gradient only directly influences the velocity, which in turn has an effect on the position

# SGD with Momentum

## Conventional SGD

```
# Vanilla update
x += - learning_rate * dx
```

## SGD with Momentum

```
# Momentum update
v = mu * v - learning_rate * dx # integrate velocity
x += v # integrate position
```

v initialized as 0          mu: momentum, hyperparameter, usually 0.9
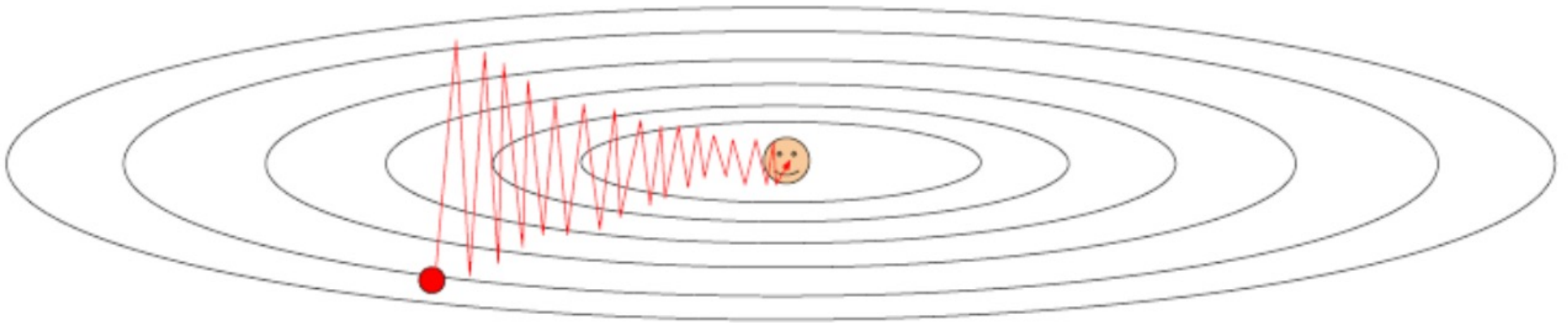
# After using SGD+Momentum



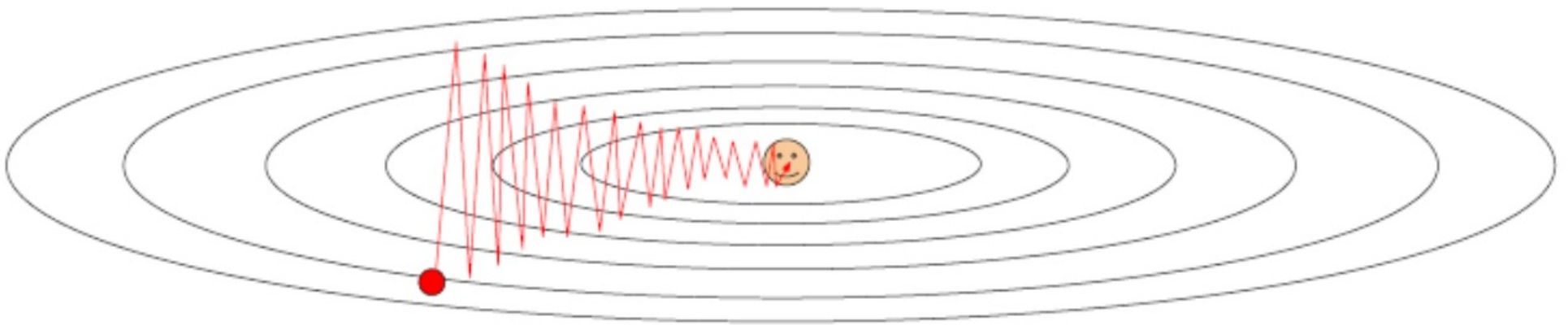Local Minima        Saddle points

# Problem of SGD

- When loss changes quickly in one direction and slowly in another

  -- Very slow progress along shallow dimension, jitter along steep direction

# Adjust Learning Rate Adaptively

- Avoid globally equal learning rate can alleviate

$$w = w - lamda * dw$$

# AdaGrad

- Element-wise adaptively adjust *effective* learning rate
  -- for weights that receive high gradients, reduce
  -- for weights that receive small or infrequent update, increase

```python
# Assume the gradient dx and parameter vector x
cache += dx**2
x += - learning_rate * dx / (np.sqrt(cache) + eps)
```

eps is between 1e-4 and 1e-8

# RMSprop (Root Mean Squared Propagation)

- Cons of AdaGrad
  - -- Learning rate monotonically decreases
  - -- Stops learning too early
- RMSprop avoid this via using moving average

```
cache = decay_rate * cache + (1 - decay_rate) * dx**2
x += - learning_rate * dx / (np.sqrt(cache) + eps)
```

decay rate is typically 0.9, 0.99, 0.999...

# Note about Learning Rate

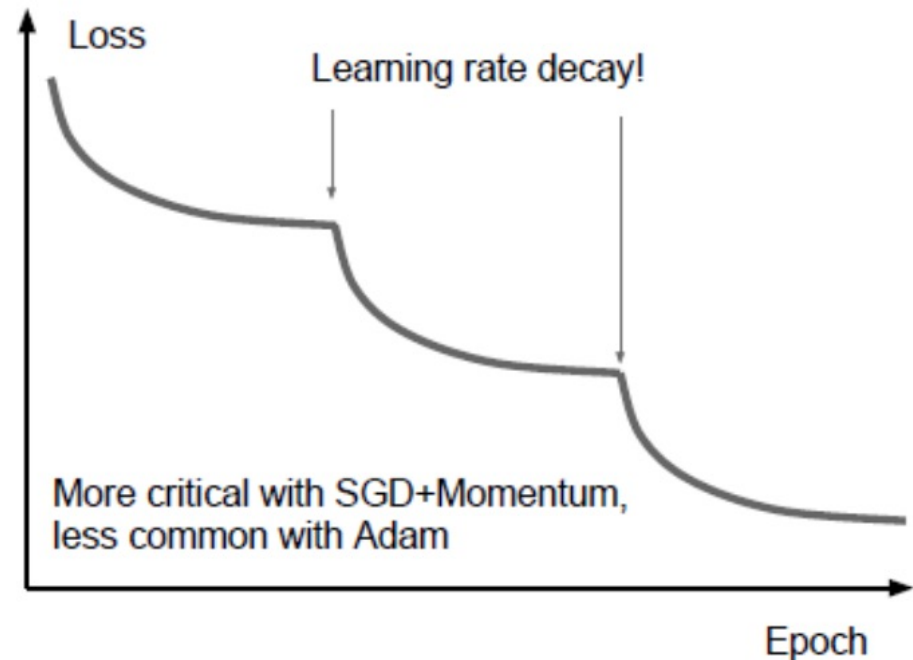- (Global) Learning rate should also be decayed gradually

**step decay:**
e.g. decay learning rate by half every few epochs.

**exponential decay:**

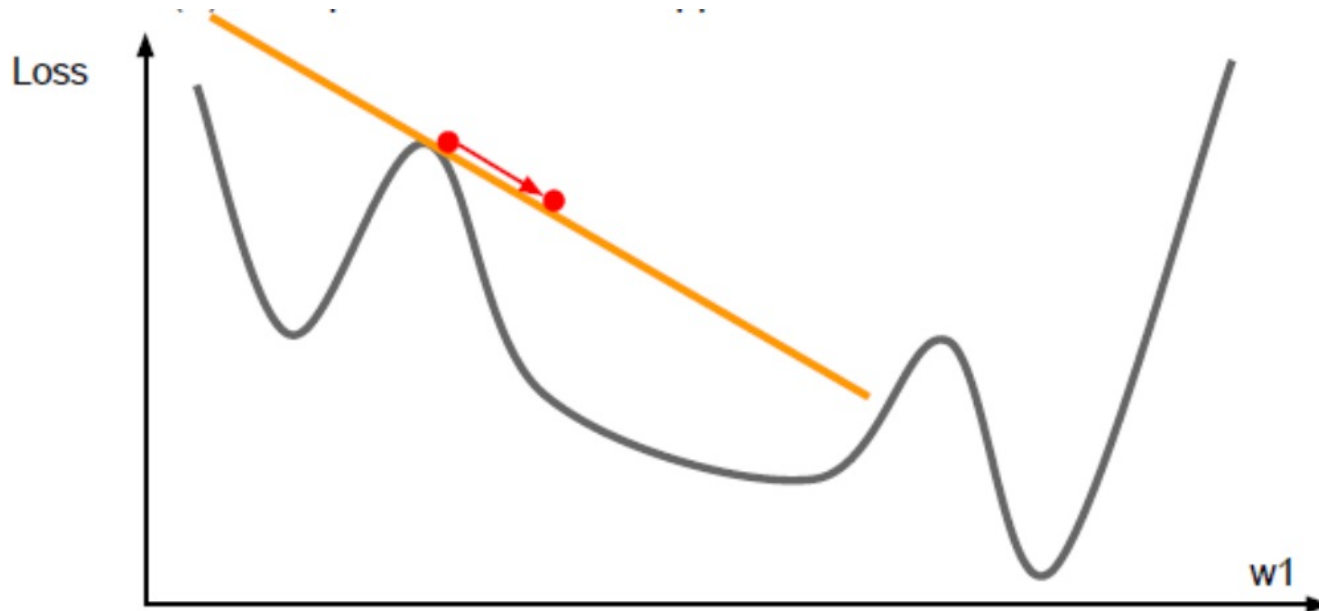$$\alpha = \alpha_0 e^{-kt}$$

**1/t decay:**

$$\alpha = \alpha_0 / (1 + kt)$$

# Pytorch Learning Rate Scheduler

Decreases the learning rate by gamma every step_size epochs.

- from torch.optim.lr_scheduler import StepLR

- optimizer = torch.optim.SGD(model.parameters(), lr=0.1)
- scheduler = StepLR(optimizer, step_size=30, gamma=0.1)

- for epoch in range(num_epochs):
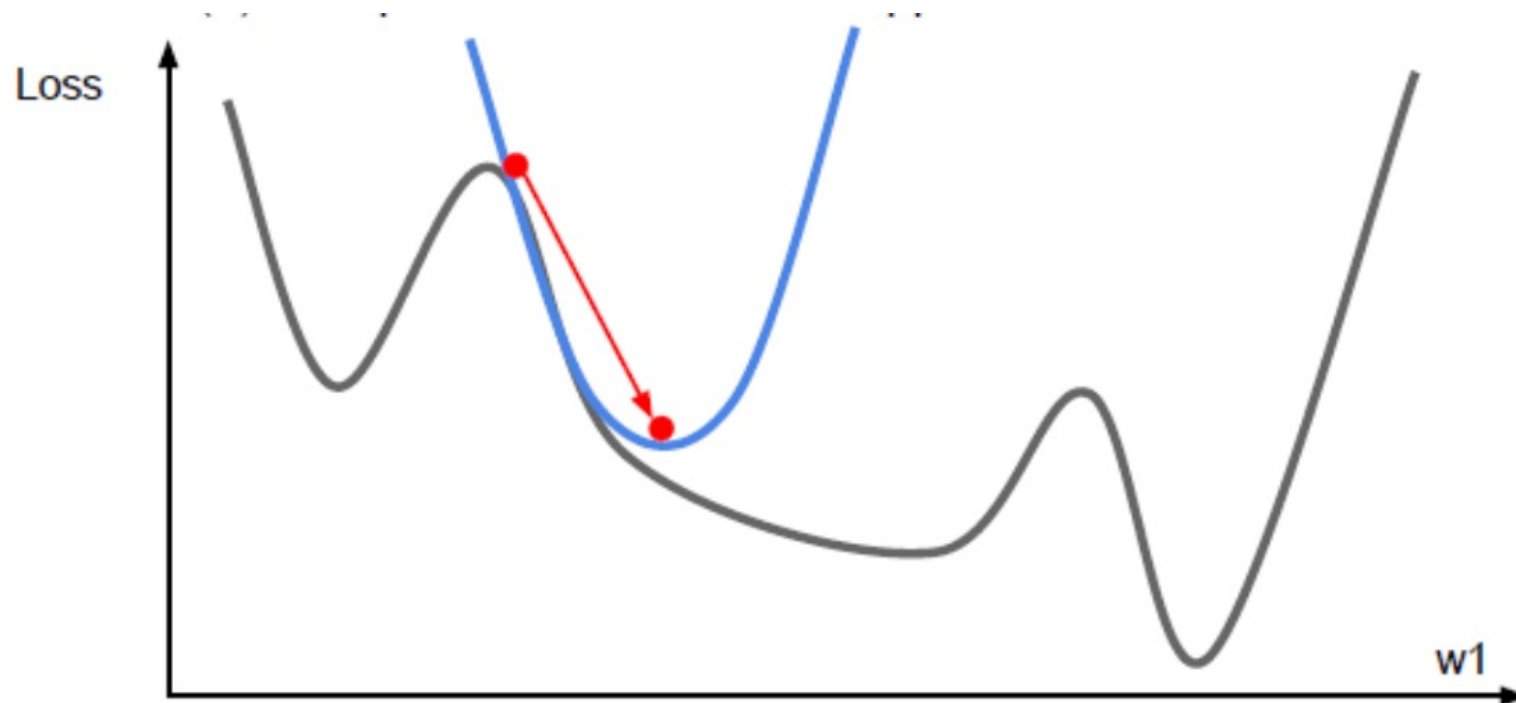-     train(...)
-     scheduler.step()

# First-Order Optimization

- Gradient descent belongs to first-order optimization
- Convergence rate is relatively slow

# Second-Order Optimization

- Second-order optimization has fast convergence rate

# Pros of Second-Order Optimization

second-order Taylor expansion:

Hessian matrix

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

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

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

No hyperparameters!
No learning rate!
(Though you might use one in practice)

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a function taking as input a vector $\mathbf{x} \in \mathbb{R}^n$ and outputting a scalar $f(\mathbf{x}) \in \mathbb{R}$; if all second partial derivatives of $f$ exist and are continuous over the domain of the function, then the Hessian matrix $\mathbf{H}$ of $f$ is a square $n \times n$ matrix, usually defined and arranged as follows:

$$\mathbf{H} = \begin{bmatrix} \dfrac{\partial^2 f}{\partial x_1^2} & \dfrac{\partial^2 f}{\partial x_1\, \partial x_2} & \cdots & \dfrac{\partial^2 f}{\partial x_1\, \partial x_n} \\[2ex] \dfrac{\partial^2 f}{\partial x_2\, \partial x_1} & \dfrac{\partial^2 f}{\partial x_2^2} & \cdots & \dfrac{\partial^2 f}{\partial x_2\, \partial x_n} \\[2ex] \vdots & \vdots & \ddots & \vdots \\[2ex] \dfrac{\partial^2 f}{\partial x_n\, \partial x_1} & \dfrac{\partial^2 f}{\partial x_n\, \partial x_2} & \cdots & \dfrac{\partial^2 f}{\partial x_n^2} \end{bmatrix}.$$

or, by stating an equation for the coefficients using indices i and j:

$$\mathbf{H}_{i,j} = \frac{\partial^2 f}{\partial x_i\, \partial x_j}.$$

From wiki

# Cons of Second-Order Optimization

second-order Taylor expansion:

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

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

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

Hessian has O(N^2) elements
Inverting takes O(N^3)
N = (Tens or Hundreds of) Millions

# Acknowledgement

Many materials of the slides of this course are adopted and re-produced from several deep learning courses and tutorials.