# CS396: Selected CS2
## (Deep Learning for visual recognition)

**Spring 2022**

**Dr. Wessam EL-Behaidy**
Associate Professor, Computer Science Department,
Faculty of Computers and Artificial Intelligence,
Helwan University.

**Lectures (Course slides) are based on** Stanford course :
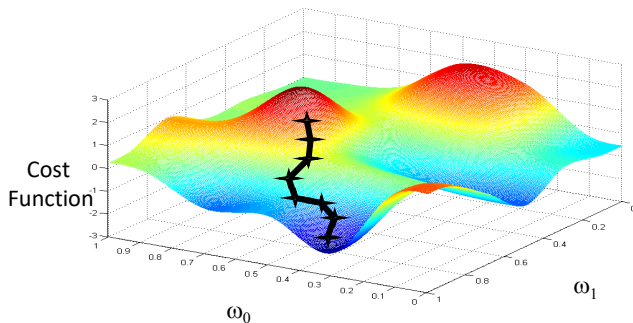Convolutional Neural Networks for Visual Recognition (CS231n):
http://cs231n.stanford.edu/index.html

Lecture 5: Training Networks, Part 2

- Fancier optimization
- Regularization

## **Following the Gradient**

- ⌃ Compute the best direction along which we should <u>change our parameter</u> (weight) vector that is mathematically guaranteed to be the direction of the **steepest descend.**
- ⌃ This direction will be related to the **gradient** of the <u>cost function.</u>



Cost Function, $\omega_0$, $\omega_1$

# Gradient Descent

- We make steps down the cost function in the direction with the steepest descent, and the size of each step is determined by the parameter **α**, which is called the **learning rate.**
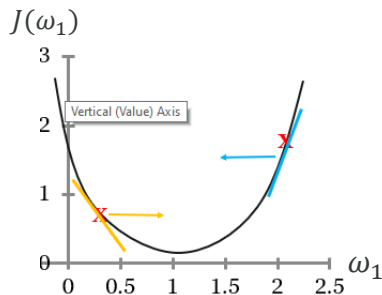- Type equation here.
- The gradient descent algorithm is:

repeat until **convergence** {

$$\omega_j := \omega_j - \alpha \frac{\partial}{\partial \omega_j} J(\omega_0, \omega_1)$$

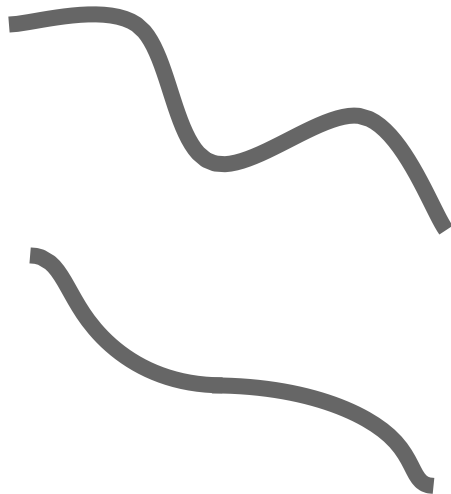$$(\text{for } j = 0 \text{ and } j = 1)$$

**Learning rate (step size)**

}



$J(\omega_1)$

Positive slope (positive number)→ $\omega_1$ will decrease
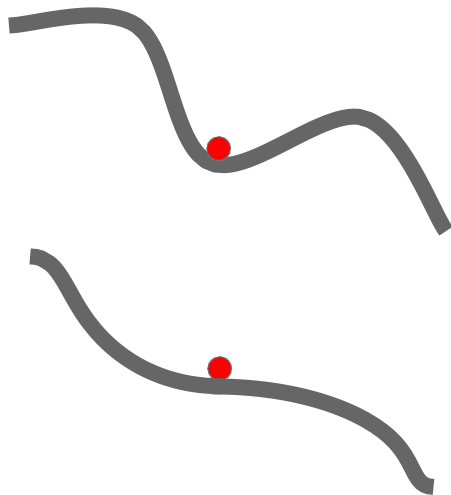Negative slope (negative number)→ $\omega_1$ will increase

Andrew Ng

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

# Optimization: Problems with SGD

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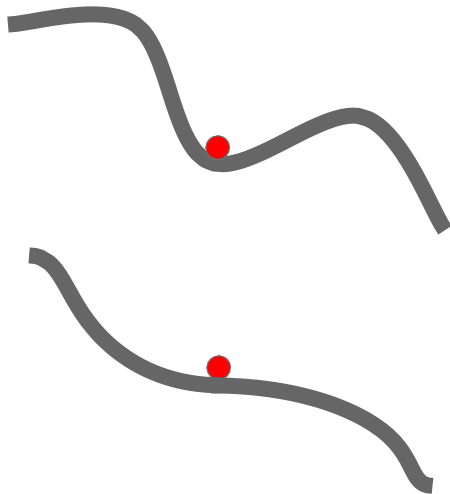function has a
**local minima** or
**saddle point**?

Zero gradient,
gradient descent
gets stuck

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

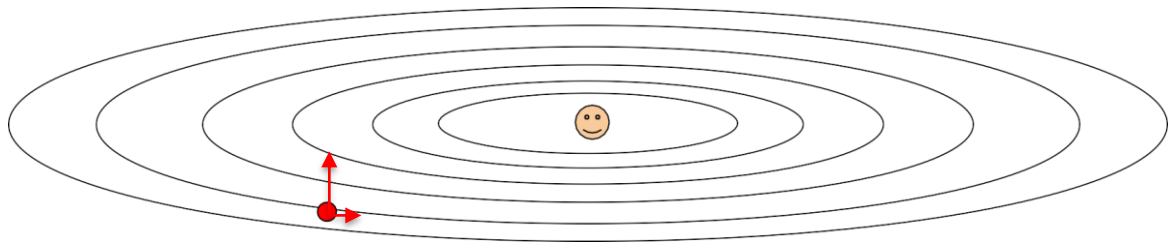Saddle points much more common in high dimension

Dauphin et al, "Identifying and attacking the saddle point problem in high-dimensional non-convex optimization", NIPS 2014

# Optimization: Problems with SGD

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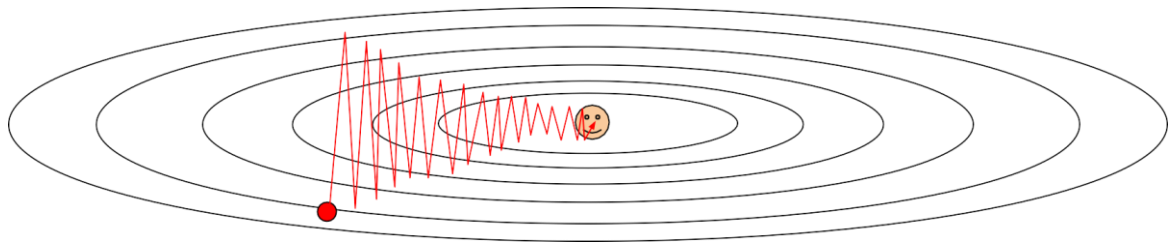What does gradient descent do?

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

<span style="color:red">Very slow progress along shallow dimension, jitter along steep direction</span>



So, we need to slow down changes on vertical dimension and increase it on the horizontal dimension.

## SGD

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

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

**Gradient**

**actual step**

**Velocity**

## 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" as a running <u>mean</u> 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

Local Minima

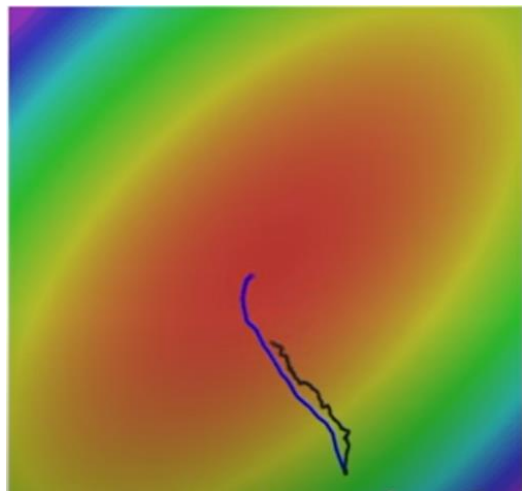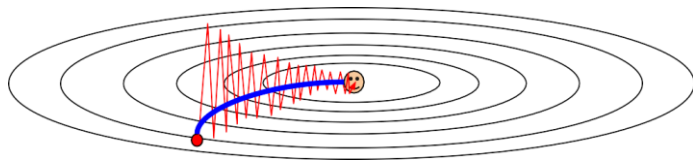Saddle points

Poor Conditioning
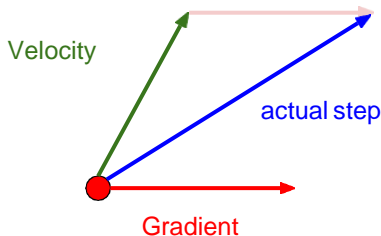
SGD
SGD+Momentum

## 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

## Momentum update:



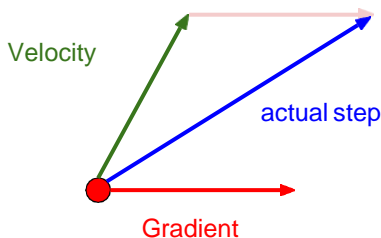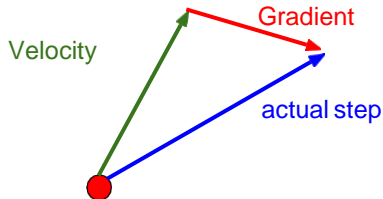Combine gradient at current point with velocity to get step used to update weights

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



"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 Momentum

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



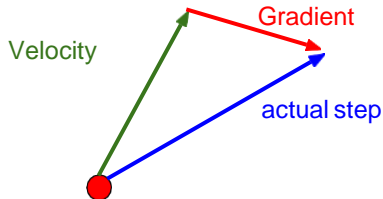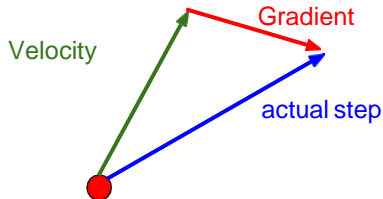Gradient

Velocity

actual step

"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 Momentum

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

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

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



"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 Momentum

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

Annoying, usually we want update in terms of $x_t, \nabla f(x_t)$
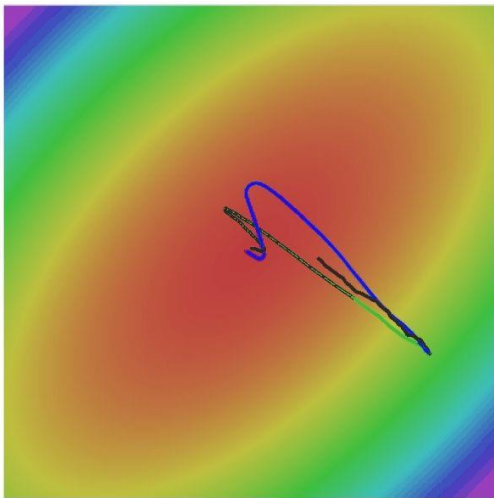
Change of variables $\tilde{x}_t = x_t + \rho v_t$ and rearrange:

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

The **first equation** looks exactly like updating the velocity in the SGD momentum case. At the **second equation**, we have our current point plus our current velocity plus **a weighted difference** between our current velocity and our previous velocity. Here, Nesterov momentum is kind of incorporating some kind of **error-correcting** term between your current velocity and your previous velocity.

# Nesterov Momentum
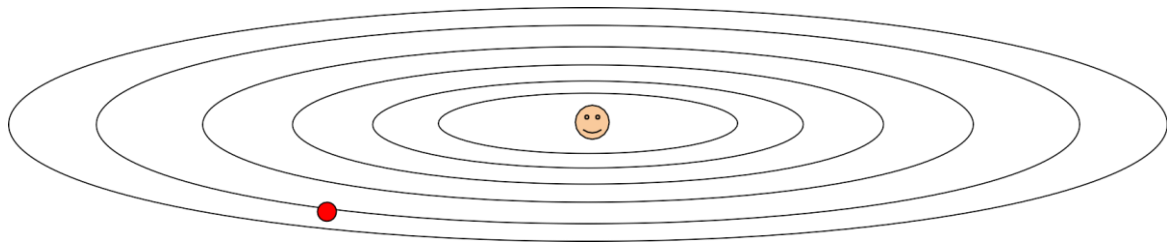


SGD

SGD+Momentum

Nesterov (nag)

# 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"

Duchi et al, "Adaptive subgradient methods for online learning and stochastic optimization", JMLR 2011
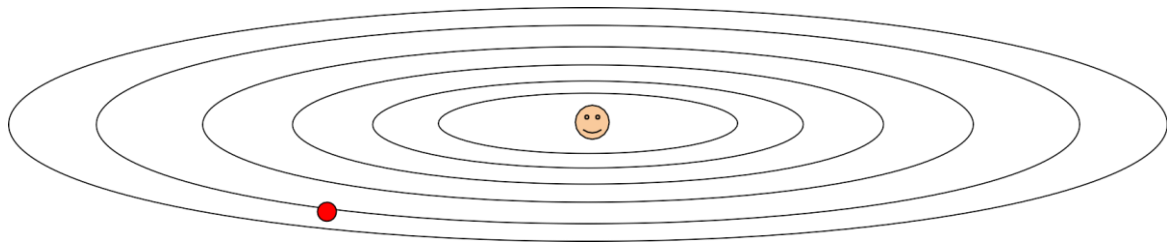
```
grad_squared = 0
while True:
  dx = compute_gradient(x)
  grad_squared += dx * dx
  x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```



Q: What happens with AdaGrad?

# AdaGrad

```python
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

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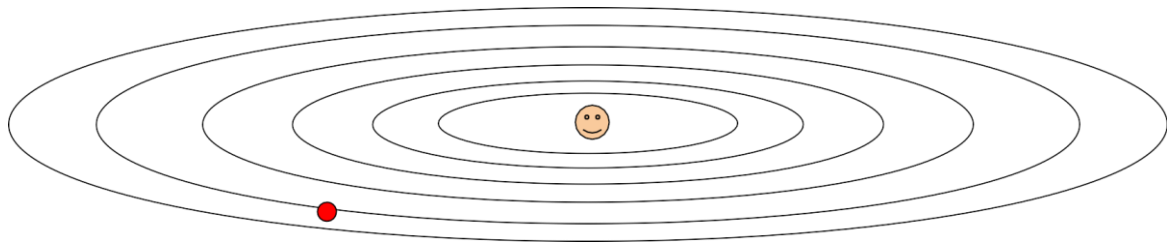


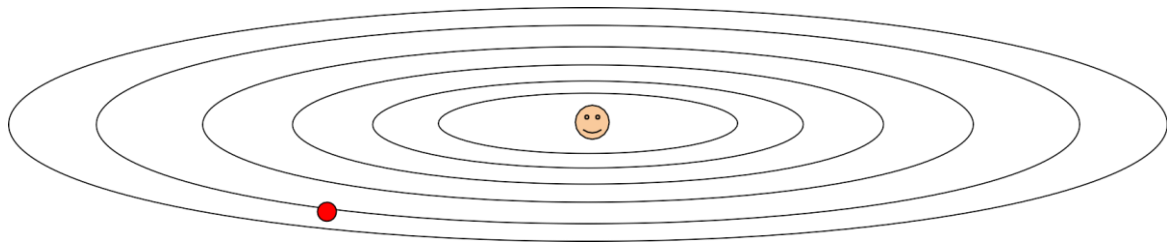Q2: What happens to the step size over long time?

```
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

AdaGrad

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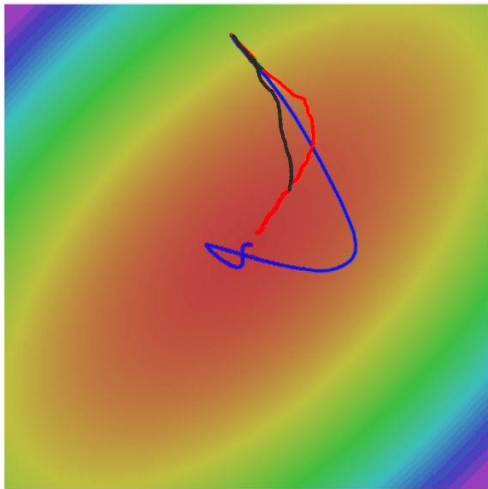
RMSProp

```
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)
```

Tieleman and Hinton, 2012          Typically, decay rate = 0.99

# RMSProp



SGD

SGD+Momentum

RMSProp

Note:
AdaGrad just gets stuck due to the problem of continually decaying learning rates. For that, it is not shown in this comparison.

# Adam (almost)

```python
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))
```

```
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

Sort of like RMSProp with momentum

Q: What happens at first timestep?

Kingma and Ba, "Adam: A method for stochastic optimization", ICLR 2015

# 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))
```
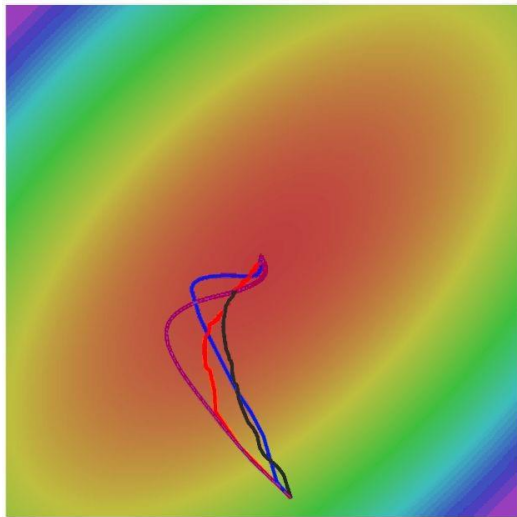
Momentum

Bias correction

AdaGrad / RMSProp

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

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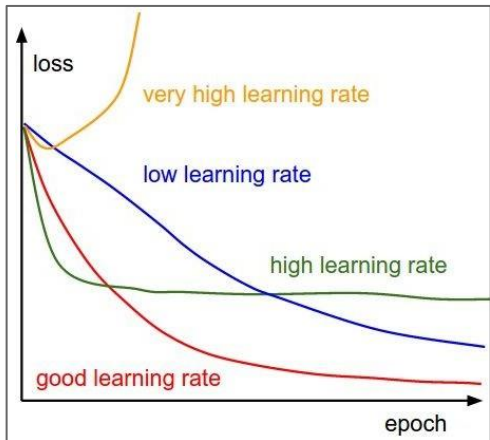beta2 = 0.999, and learning_rate = 1e-3 or 5e-4
is a great starting point for many models!

Kingma and Ba, "Adam: A method for stochastic optimization", ICLR 2015
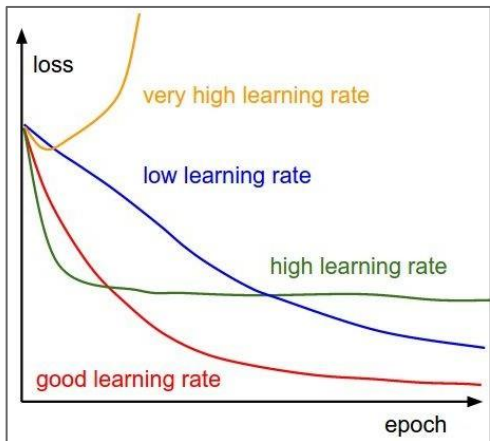
SGD

SGD+Momentum

RMSProp

Adam

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.



**=> Learning rate decay over time!**

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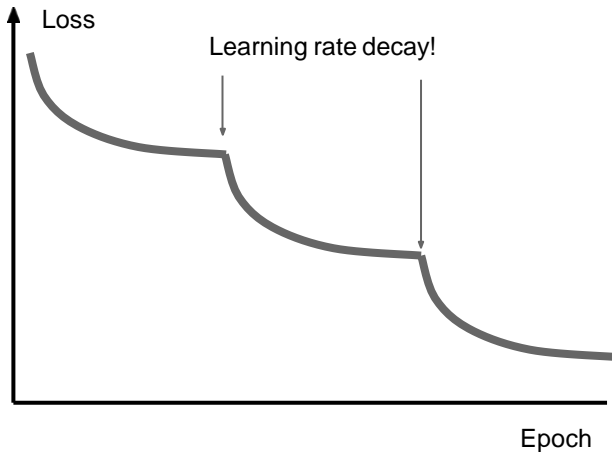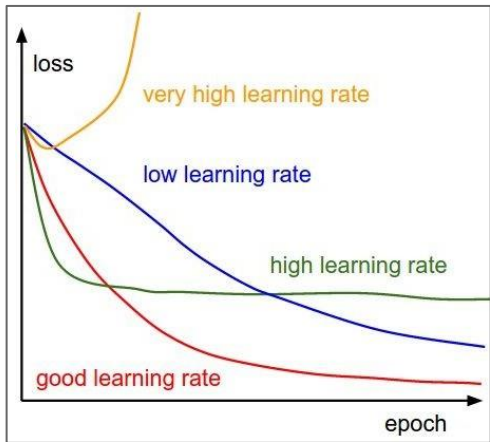e.g. decay learning rate by half every few epochs.

**exponential decay:**

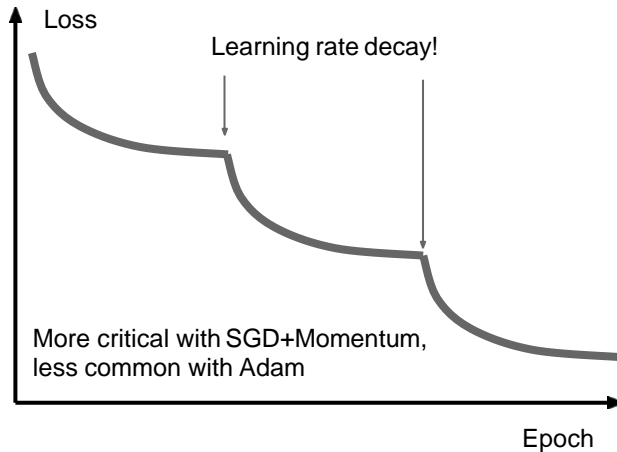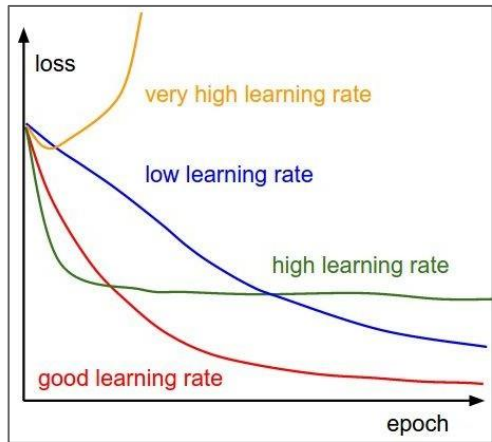$$\alpha = \alpha_0 e^{-kt}$$

**1/t decay:**

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

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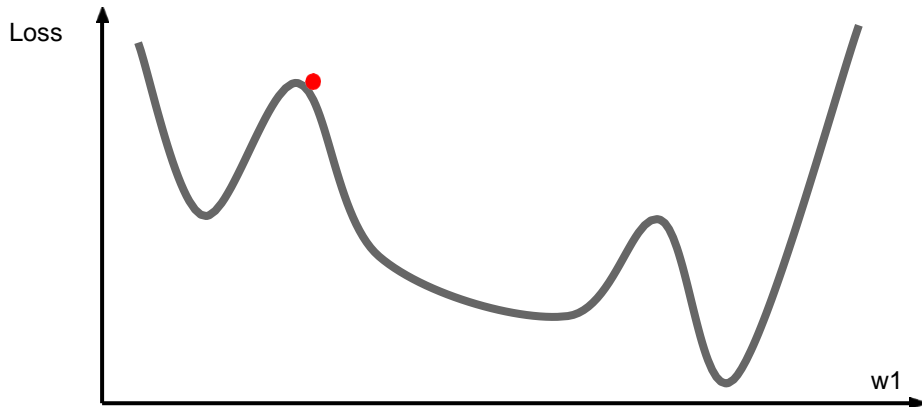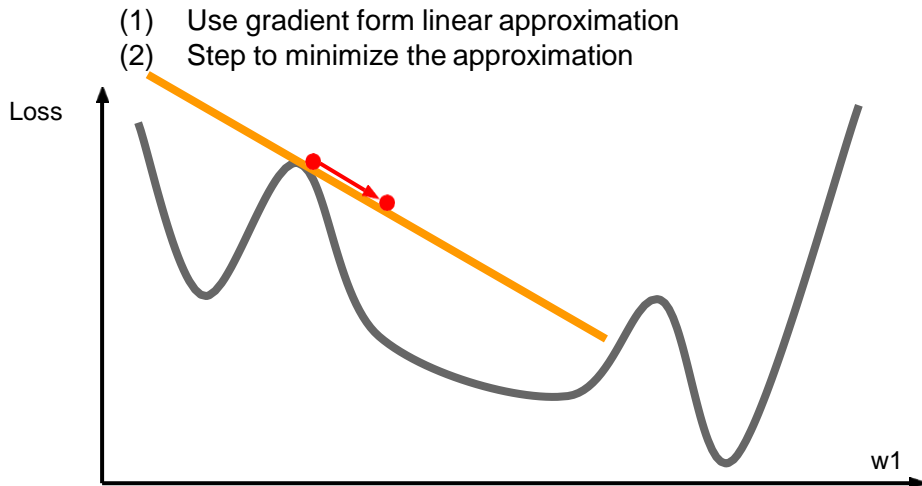
# SGD, SGD+Momentum, Adagrad, RMSProp, Adam all have **learning rate** as a hyperparameter.

# First-Order Optimization

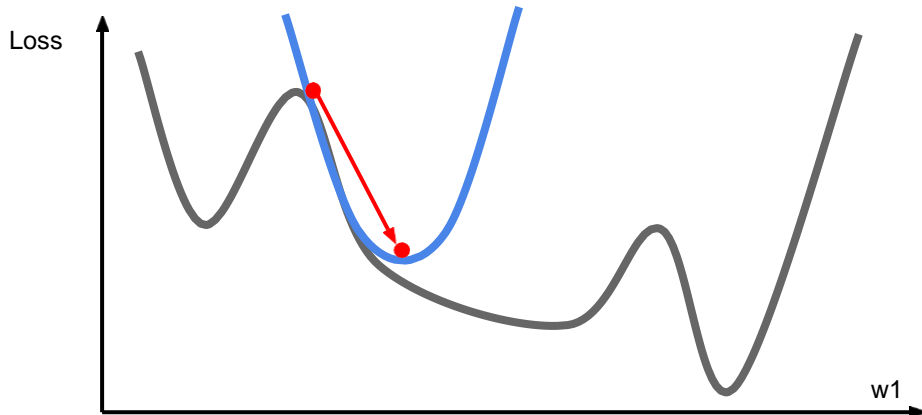# First-Order Optimization

(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

# Advantage and Disadvantage of Second-Order Optimization

Advantage:
- No hyperparameters
- No learning rate

Disadvantage:
- High complexity which is bad for deep learning

# L-BFGS (Limited-memory 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,
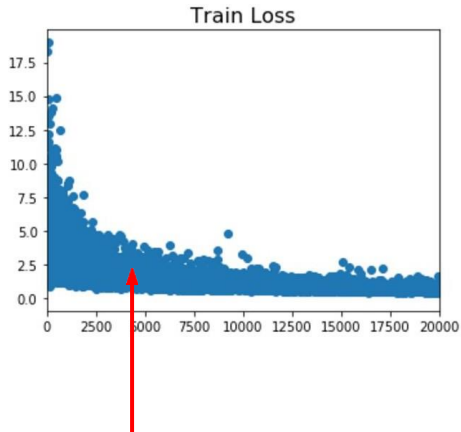  stochastic setting is an active area of research.

Le et al, "On optimization methods for deep learning, ICML 2011"
Ba et al, "Distributed second-order optimization using Kronecker-factored approximations", ICLR 2017
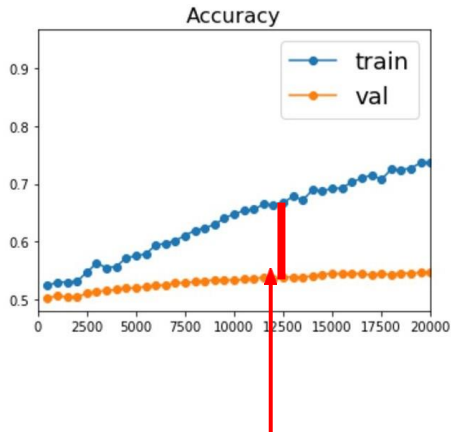
## In practice:

- **Adam** is a good default choice in many cases
- **SGD+Momentum** with learning rate decay often outperforms Adam by a bit, but requires more tuning

- If you can afford to do full batch updates then try out **L-BFGS** (and don't forget to disable all sources of noise)

Better optimization algorithms help reduce training loss

But we really care about error on new data - how to reduce the gap?

1. Train multiple independent models  or
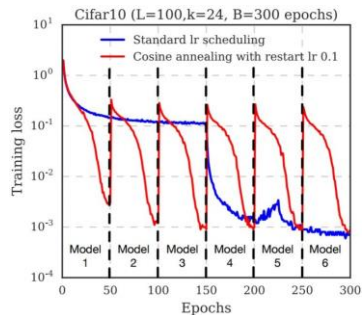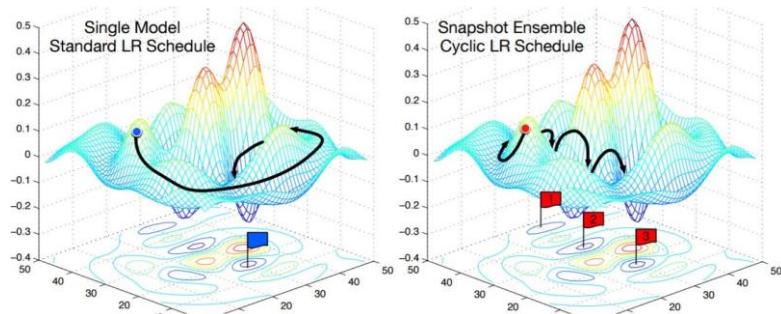2. Use multiple snapshots of a single model during training

At test time average their results

(Take average of predicted probability distributions, then choose argmax)

Enjoy 2% extra performance
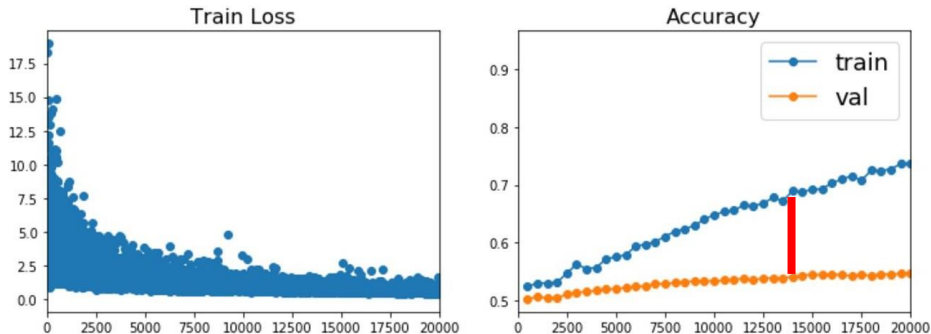
# Model Ensembles: Tips and Tricks

Instead of training independent models, use multiple snapshots of a single model during training!



Loshchilov and Hutter, "SGDR: Stochastic gradient descent with restarts", arXiv 2016
Huang et al, "Snapshot ensembles: train 1, get M for free", ICLR 2017
Figures copyright Yixuan Li and Geoff Pleiss, 2017. Reproduced with permission.

Cyclic learning rate schedules can make this work even better!

# How to improve single-model performance?



Regularization

# What is missing in this Loss function?

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

**Data loss**: Model predictions
should match training data ( encourage overfitting)

# What is missing in this Loss function?

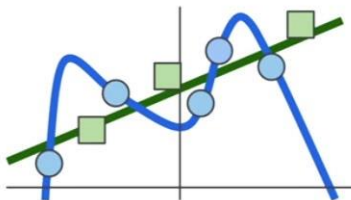$$L(W) = \frac{1}{N} \sum_{i=1}^{N} L_i(f(x_i, W), y_i)$$

**Data loss**: Model predictions
should match training data ( encourage overfitting)

# 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**: Technique to discourage the complexity of the model (i.e, express preferences over weights). It does this by penalizing the loss function. This helps to solve the overfitting problem.

# Regularization parameter λ

Regularization works on assumption that smaller weights generate simpler model and thus helps avoid **overfitting**.

**λ** is the **penalty term** or **regularization parameter** which determines how much to penalizes the weights.

- **λ = 0,** then the regularization term becomes zero (back to the original Loss function).

- **λ is large**, the weights become close to zero (i.e. a very simple model have **underfitting**).

  **λ** is a _hyperparameter_ between 0 and a large value.

**Simple Examples:**

**L2 regularization:** $R(W) = \sum_k \sum_l W_{k,l}^2$

L1 regularization: $R(W) = \sum_k \sum_l |W_{k,l}|$

Elastic net (L1 + L2): $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$

**More complex:**

Dropout

Batch normalization

Stochastic depth, fractional pooling, etc

# L1 regularization

L1 regularization: $R(W) = \sum_k \sum_l |W_{k,l}|$

**Built-in feature selection** : L1 regularization does feature selection. It does this by **assigning insignificant input features with zero weight** and useful features with a non zero weight.

$$w_1 = [1, 0, 0, 0]$$

**Sparsity :** In L1 regularization we shrink the parameters to zero. When input features have weights closer to zero that leads to sparse L1 norm. In Sparse solution majority of the input features have zero weights and very few features have non zero weights.

# L2 regularization

<u>L2 regularization:</u> $R(W) = \sum_k \sum_l W_{k,l}^2$

L2 regularization <u>forces the weights to be small but does not make them zero</u> and does **non-sparse solution**.

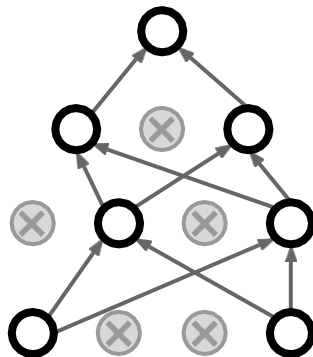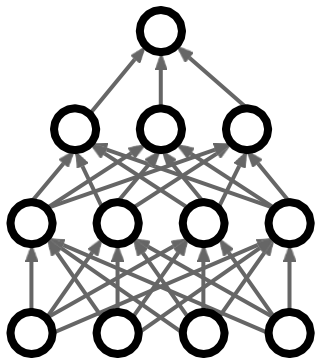$$w_2 = [0.25, 0.25, 0.25, 0.25]$$

L2 regularization likes to "spread out" the weights

*L2 has **no feature selection**, it gives better prediction when output variable is a function of all input features*

*L2 regularization is able to learn complex data patterns*

# Regularization: Dropout

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



Srivastava et al, "Dropout: A simple way to prevent neural networks from overfitting", JMLR 2014

# Regularization: Dropout

```python
p = 0.5 # probability of keeping a unit active. higher = less dropout

def train_step(X):
  """ X contains the data """

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

  # backward pass: compute gradients... (not shown)
  # perform parameter update... (not shown)
```
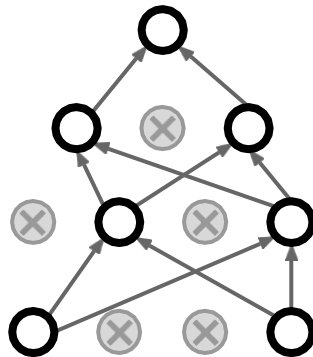
# Regularization: Dropout

How can this possibly be a good idea?



Forces the network to have a redundant representation;
Prevents co-adaptation of features

has an ear **X**

has a tail

is furry **X**

has claws

mischievous look **X**

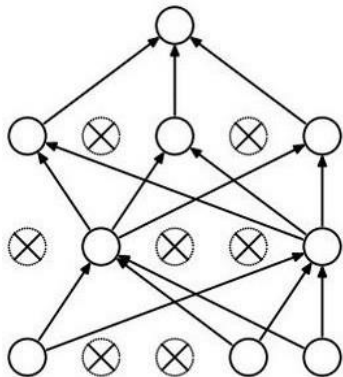cat score

How can this possibly be a good idea?



Another interpretation:

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

Each binary mask is one model

An FC layer with 4096 units has $2^{4096} \sim 10^{1233}$ possible masks!
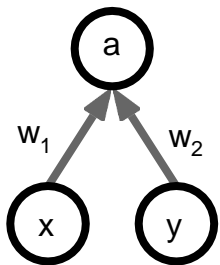
# At test time….



**Ideally**:
want to integrate out all the noise

**Monte Carlo approximation:**
do many forward passes with different dropout masks, average all predictions

Can in fact do this with a single forward pass! (approximately)
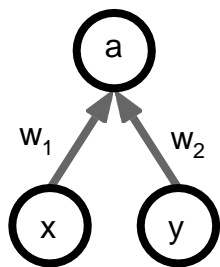Leave all input neurons turned on (no dropout).

$a$

$w_1$    $w_2$

$x$    $y$

(this can be shown to be an
approximation to evaluating the
whole ensemble)

Can in fact do this with a single forward pass! (approximately)
Leave all input neurons turned on (no dropout).

At test time we have: $E\big[a\big] = w_1 x + w_2 y$

During training we have:

With p=0.5, using all inputs in the forward pass would inflate the activations by 2x from what the network was "used to" during training!

=> Have to compensate by scaling the activations back down by ½

$$E[a] = \frac{1}{4}(w_1 x + w_2 y) + \frac{1}{4}(w_1 x + 0y)$$
$$+ \frac{1}{4}(0x + 0y) + \frac{1}{4}(0x + w_2 y)$$
$$= \frac{1}{2}(w_1 x + w_2 y)$$

At test time, **multiply** by dropout probability

a

$w_1$   $w_2$

x    y

# Dropout: Test time

```python
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
```

At test time all neurons are active always
=> We must scale the activations so that for each neuron:
<u>output at test time</u> = <u>expected output at training time</u>

# Dropout Summary

```python
""" Vanilla Dropout: Not recommended implementation (see notes below) """

p = 0.5 # probability of keeping a unit active. higher = less dropout

def train_step(X):
  """ X contains the data """

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

  # backward pass: compute gradients... (not shown)
  # perform parameter update... (not shown)

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
```

drop in forward pass

scale at test time

# More common: "Inverted dropout"

```python
p = 0.5 # probability of keeping a unit active. higher = less 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

  # backward pass: compute gradients... (not shown)
  # perform parameter update... (not shown)

def predict(X):
  # ensembled forward pass
  H1 = np.maximum(0, np.dot(W1, X) + b1) # no scaling necessary
  H2 = np.maximum(0, np.dot(W2, H1) + b2)
  out = np.dot(W3, H2) + b3
```
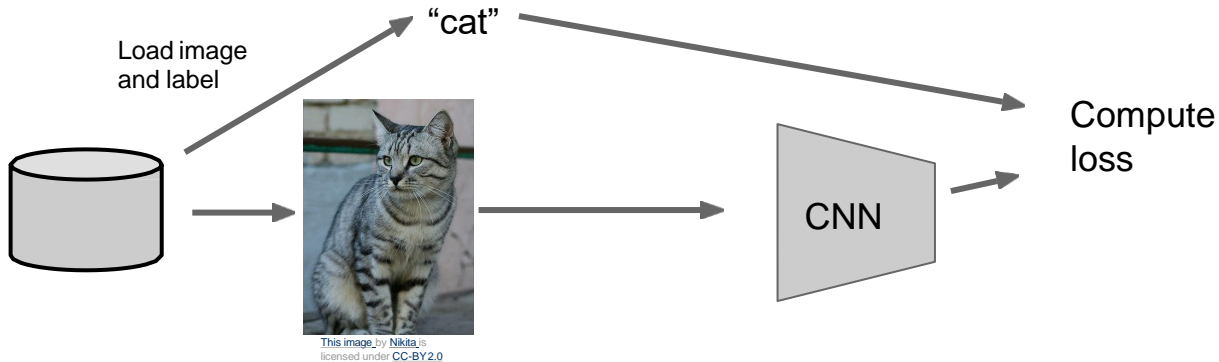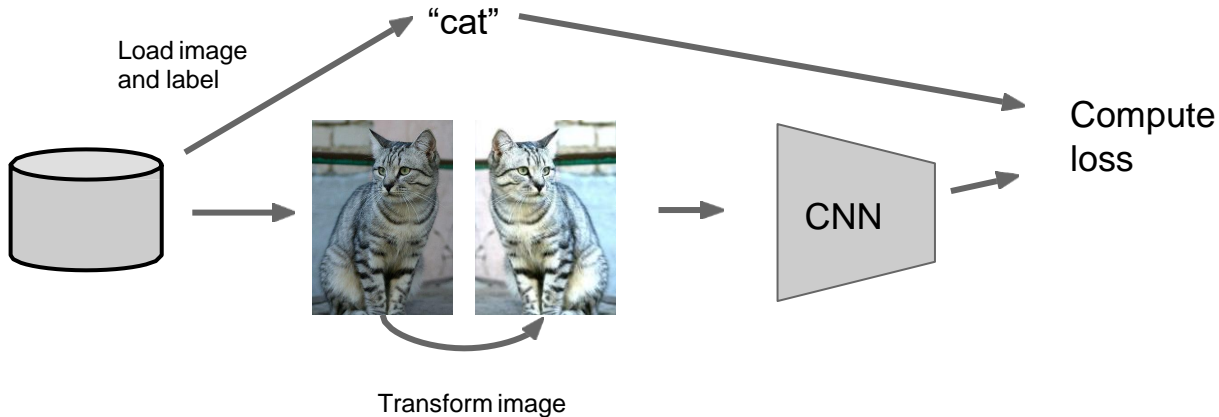
test time is unchanged!

# Regularization: Data Augmentation



Load image
and label

"cat"

Compute
loss

CNN

This image by Nikita is
licensed under CC-BY 2.0

Load image
and label

"cat"

Transform image

CNN

Compute
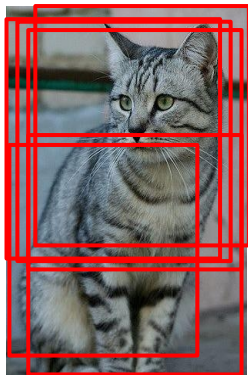loss

# Random crops and scales

**Training**: sample random crops / scales

ResNet:
1. Pick random L in range [256, 480]
2. Resize training image, short side = L
3. Sample random 224 x 224 patch

# Color Jitter

Simple: Randomize
contrast and brightness

Get creative for your problem!

Random mix/combinations of :
- translation
- rotation
- stretching
- shearing,
- …

# Regularization: A common pattern

**Training**: Add random noise
**Testing**: Marginalize over the noise

**Examples**:
Dropout
Batch Normalization
Data Augmentation

**Training**: Add random noise
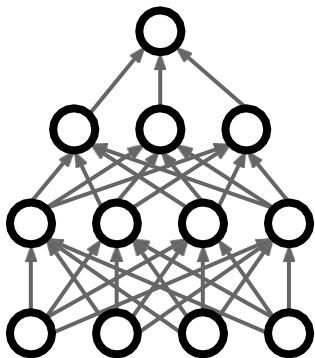**Testing**: Marginalize over the noise

**Examples**:
Dropout
Batch Normalization
Data Augmentation
DropConnect



Wan et al, "Regularization of Neural Networks using DropConnect", ICML 2013

# Regularization: A common pattern

**Training**: Add random noise
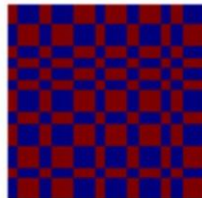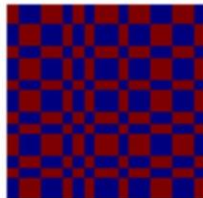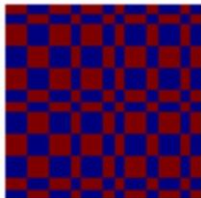**Testing**: Marginalize over the noise

**Examples**:
Dropout
Batch Normalization
Data Augmentation
DropConnect
Fractional Max Pooling



Graham, "Fractional Max Pooling", arXiv 2014

# Regularization: A common pattern

**Training**: Add random noise
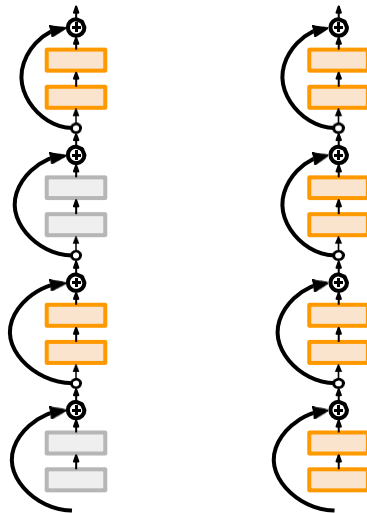**Testing**: Marginalize over the noise

**Examples**:
Dropout
Batch Normalization
Data Augmentation
DropConnect
Fractional Max Pooling
Stochastic Depth

Huang et al, "Deep Networks with Stochastic Depth", ECCV 2016

- Optimization
  - Momentum, RMSProp, Adam, etc
- Regularization
  - Dropout, etc