# Training Neural Nets

# Training Neural Nets through SGD

Mini-batch SGD

Loop:

1. Sample a batch of data

2. Forward prop it through the graph (network), get loss

3. Backprop to calculate the gradients

4. Update the parameters using the gradient

# Convolution Layer - Stride=1, Pad=0, Channel=1

**Input** An image with height $n_h$ and width $n_w$, represented by an $n_h \times n_w$ matrix $X$, e.g., $n_h = n_w = 32$.

**Filter** Represented by an $F \times F$ matrix $W$, e.g., $F = 3$.

**Output** $Y = X * W$, where

$$Y_{ij} = \sum_{m=0}^{F-1} \sum_{n=0}^{F-1} W_{m,n} X_{i+m,j+n}$$

for $i = 0, \cdots, n_h - F + 1$ and $j = 0, \cdots, n_w - F + 1$.

For example, if the input is $32 \times 32$ and $F = 3$, then the output will be $30 \times 30$.

# Gradient Backpropagation in Convolution Layer - Stride=1, Pad=0, Channel=1

Forward $Y = X * W$, where $W \in \mathbb{R}^{F \times F}$ and $X \in \mathbb{R}^{n_h \times n_w}$

$$Y_{ij} = \sum_{m=0}^{F-1} \sum_{n=0}^{F-1} W_{mn} X_{i+m,j+n}$$

Gradient Backpropagation of gradients

$$\frac{\partial \mathcal{L}}{\partial X_{ij}} = \sum_{m=0}^{F-1} \sum_{n=0}^{F-1} W_{mn} \left( \frac{\partial \mathcal{L}}{\partial Y} \right)_{i-m,j-n}$$

$$= \sum_{m=0}^{F-1} \sum_{n=0}^{F-1} W_{F-1-m,F-1-n} \left( \frac{\partial \mathcal{L}}{\partial Y} \right)_{i-(F-1)+m,\ j-(F-1)+n}$$

$$= \left[ \left( \frac{\partial \mathcal{L}}{\partial Y} \right) * \hat{W} \right]_{i-(F-1),\ j-(F-1)}$$

where filter $\hat{W}$ is derived after rotating $W$ by 180 degrees.

# Gradient Backpropagation in Convolution Layer

Forward $Y = X * W$, where

$$Y_{ij} = \sum_{m=0}^{F-1} \sum_{n=0}^{F-1} W_{mn} X_{i+m,j+n}$$

Backward Gradients w.r.t. activities

$$\frac{\partial \mathcal{L}}{\partial X_{ij}} = \left[ \left( \frac{\partial \mathcal{L}}{\partial Y} \right) * \hat{W} \right]_{i-(F-1),\ j-(F-1)}$$

Gradients w.r.t. weights

$$\frac{\partial \mathcal{L}}{\partial W_{mn}} = \sum_i \sum_j \left( \frac{\partial \mathcal{L}}{\partial Y} \right)_{ij} X_{i+m,j+n}$$

$$\frac{\partial \mathcal{L}}{\partial W} = X * \left( \frac{\partial \mathcal{L}}{\partial Y} \right)$$

Both are convolution operations.

# Max Pooling - Stride=2

Input   An image with height $n_h$ and width $n_w$, represented by an $n_h \times n_w$ matrix $X$, e.g., $n_h = n_w = 32$.

Output   $Y = \mathrm{POOL(X)}$,

$$Y_{ij} = \max\{X_{2i,2j}, X_{2i+1,2j}, X_{2i,2j+1}, X_{2i+1,2j+1}\}$$

for $i = 0, \cdots, n_h/2$ and $j = 0, \cdots, n_w/2$.

For example, if the input is $32 \times 32$ and stride $S = 2$, then the output will be $16 \times 16$.

# Overview of the training process

1.  One time setup

    activation functions, preprocessing, weight initialization, regularization, gradient checking

2.  Training dynamics

    babysitting the learning process, parameter updates, hyperparameter optimization
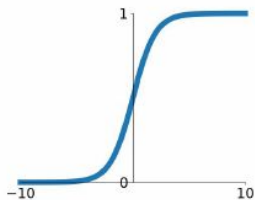
3.  Evaluation

    model ensembles
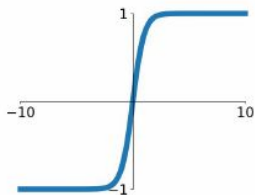
# Activation Function

# Activation Functions

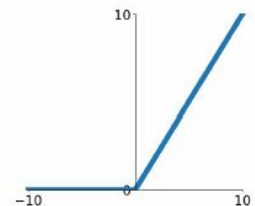**Sigmoid**
$$\sigma(x) = \frac{1}{1+e^{-x}}$$
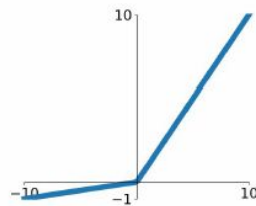
**tanh**
$$\tanh(x)$$

**ReLU**
$$\max(0, x)$$

**Leaky ReLU**
$$\max(0.1x, x)$$
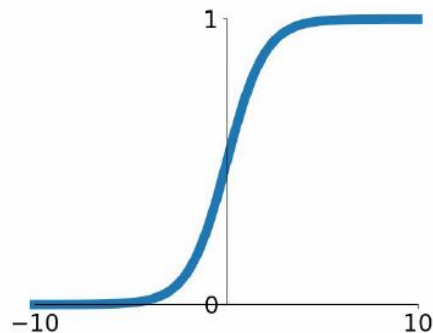
**Maxout**
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

**ELU**
$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$
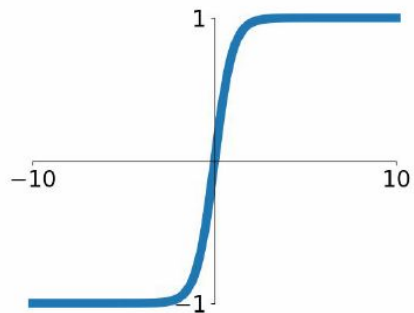
# Activation Functions

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



**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
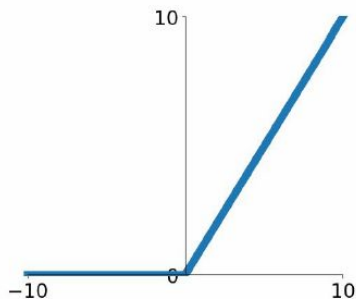2. Sigmoid outputs are not zero-centered
3. exp() is a bit compute expensive

# Activation Functions



**tanh(x)**

- Squashes numbers to range [-1,1]
- zero centered (nice)
- still kills gradients when saturated :(

# Activation Functions



**ReLU**
(Rectified Linear Unit)

- Computes **f(x) = max(0,x)**

- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than
  sigmoid/tanh in practice (e.g. 6x)
- Actually more biologically plausible
  than sigmoid

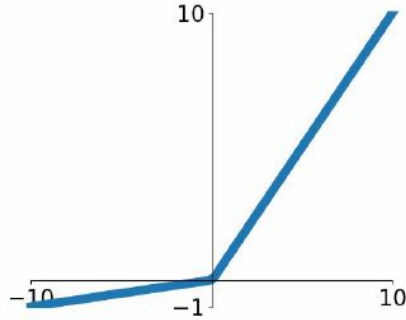- Not zero-centered output
- An annoyance:

hint: what is the gradient when x < 0?

Solution: Initialize ReLU neurons with slightly positive biases
(e.g. 0.01)

# Activation Functions

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

**Leaky ReLU**

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

**Parametric Rectifier (PReLU)**

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

backprop into \alpha
(parameter)

# Activation Functions

## Exponential Linear Units (ELU)



- All benefits of ReLU
- Closer to zero mean outputs
- Negative saturation regime compared with Leaky ReLU adds some robustness to noise

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

- Computation requires exp()

# Maxout "Neuron"

- Does not have the basic form of dot product -> nonlinearity
- Generalizes ReLU and Leaky ReLU
- Linear Regime! Does not saturate! Does not die!

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

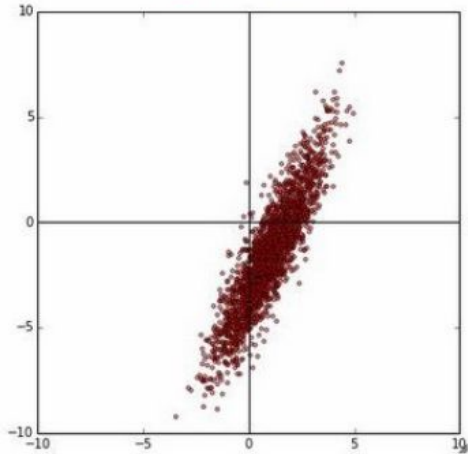Problem: doubles the number of parameters/neuron :(

# Empirical Suggestions

- ❏ Use ReLU. Be careful with your learning rates

- ❏ Try out Leaky ReLU / Maxout / ELU

- ❏ Try out tanh but don't expect much

- ❏ Don't use sigmoid

# Data Preprocessing

# Step 1: Preprocess the data



original data     zero-centered data     normalized data

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

# Step 1: Preprocess the data

In practice, you may also see **PCA** and **Whitening** of the data



| original data | decorrelated data | whitened data |
| --- | --- | --- |
| | (data has diagonal covariance matrix) | (covariance matrix is the identity matrix) |

# Mean shifting is common for image preprocessing

e.g. consider CIFAR-10 example with [32,32,3] images

- Subtract the mean image (e.g. AlexNet)
  (mean image = [32,32,3] array)
- Subtract per-channel mean (e.g. VGGNet)
  (mean along each channel = 3 numbers)

Not common to normalize variance, to do PCA or whitening

# Weight Initialization

# Weight Initialization

❏ All zero initialization (Wrong!  All parameters will undergo the same update and stay the same)

❏ Small random numbers (e.g., Gaussian with zero mean and 0.01 standard deviation)

$$\texttt{W = 0.01*np.random.randn(n,H)}$$

Problematic for deep networks as variances of neuron activities tend to become smaller and smaller

❏ Xavier Initialization:

$$\texttt{W = np.random.randn(fan\_in, fan\_out) / np.sqrt(fan\_in)}$$

❏ He Initialization

$$\texttt{W = np.random.randn(fan\_in, fan\_out) * np.sqrt(2.0/fan\_in)}$$

Recommended for ReLU units

# Batch Normalization

# Motivations on Batch Normalization

- ❏ Explicitly force the activations throughout a network to take on a unit gaussian distribution at the beginning of the training
- ❏ Possible because normalization is a simple differentiable operation
- ❏ Has become a very common practice to use Batch Normalization in neural networks. In practice, networks that use Batch Normalization are significantly more robust to bad initialization. Additionally, batch normalization can be interpreted as doing preprocessing at every layer of the network, but integrated into the network itself in a differentiable manner

# Batch Normalization

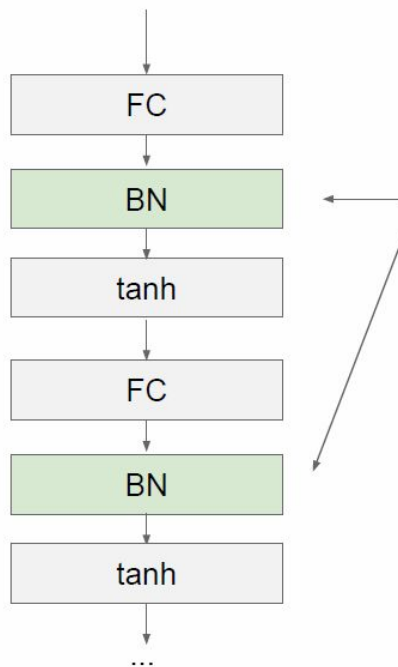"you want zero-mean unit-variance activations? just make them so."

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

$$\widehat{x}^{(k)} = \frac{x^{(k)} - \mathrm{E}[x^{(k)}]}{\sqrt{\mathrm{Var}[x^{(k)}]}}$$

this is a vanilla differentiable function...

# Batch Normalization

[Ioffe and Szegedy, 2015]

```
   │
   ▼
┌──────────┐
│    FC    │
└──────────┘
   │
   ▼
┌──────────┐
│    BN    │  ◄─────
└──────────┘
   │
   ▼
┌──────────┐
│   tanh   │
└──────────┘
   │
   ▼
┌──────────┐
│    FC    │
└──────────┘
   │
   ▼
┌──────────┐
│    BN    │  ◄─────
└──────────┘
   │
   ▼
┌──────────┐
│   tanh   │
└──────────┘
   │
   ▼
  ...
```

Usually inserted after Fully Connected or Convolutional layers, and before nonlinearity.

$$\widehat{x}^{(k)} = \frac{x^{(k)} - \mathrm{E}[x^{(k)}]}{\sqrt{\mathrm{Var}[x^{(k)}]}}$$

# Batch Normalization

Normalize:

$$\widehat{x}^{(k)} = \frac{x^{(k)} - \mathrm{E}[x^{(k)}]}{\sqrt{\mathrm{Var}[x^{(k)}]}}$$

And then allow the network to squash the range if it wants to:

$$y^{(k)} = \gamma^{(k)}\widehat{x}^{(k)} + \beta^{(k)}$$

Note, the network can learn:

$$\gamma^{(k)} = \sqrt{\mathrm{Var}[x^{(k)}]}$$

$$\beta^{(k)} = \mathrm{E}[x^{(k)}]$$

to recover the identity mapping.

# Batch Normalization

**Input:** Values of $x$ over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$;
Parameters to be learned: $\gamma$, $\beta$

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

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

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

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

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

- Improves gradient flow through the network
- Allows higher learning rates
- Reduces the strong dependence on initialization
- Acts as a form of regularization in a funny way, and slightly reduces the need for dropout, maybe

# Batch Normalization

**Input:** Values of $x$ over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$;
Parameters to be learned: $\gamma$, $\beta$

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

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

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

$$\widehat{x}_i \leftarrow \frac{x_i - \mu_\mathcal{B}}{\sqrt{\sigma^2_\mathcal{B} + \epsilon}} \qquad \text{// normalize}$$

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

**Note: at test time BatchNorm layer functions differently:**

The mean/std are not computed based on the batch. Instead, a single fixed empirical mean of activations during training is used.

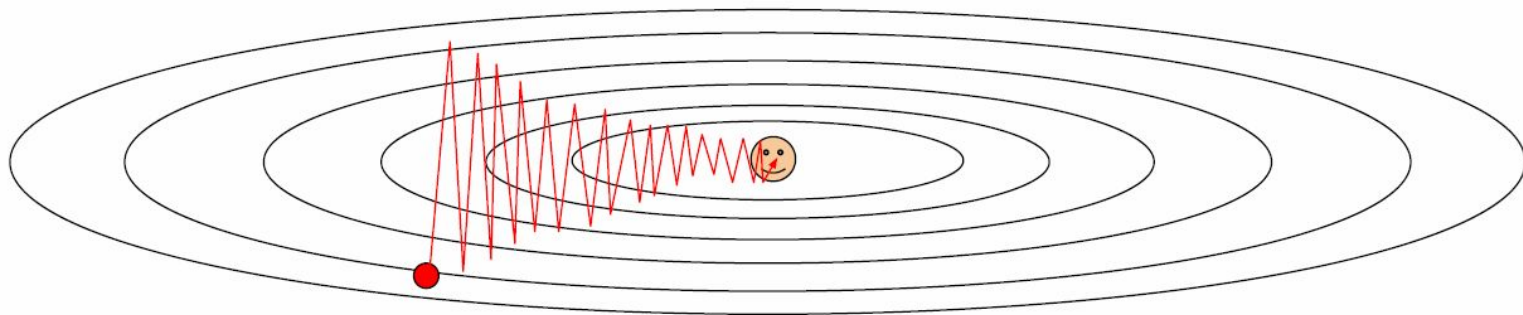(e.g. can be estimated during training with running averages)

# Training

# Optimization: Problems 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



Loss function has high **condition number**: ratio of largest to smallest singular value of the Hessian matrix is large

# SGD + Momentum

### SGD

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

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

### SGD+Momentum

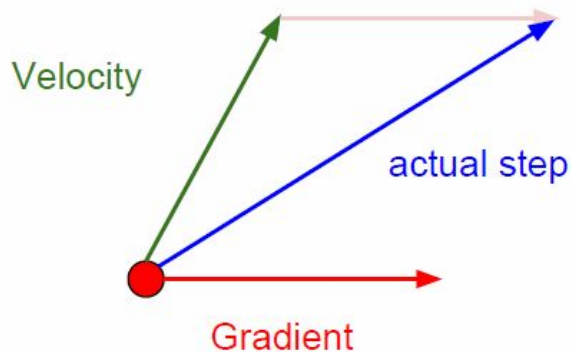$$v_{t+1} = \rho v_t + \nabla f(x_t)$$
$$x_{t+1} = x_t - \alpha v_{t+1}$$

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

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

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

# Nesterov Momentum

## Momentum update:

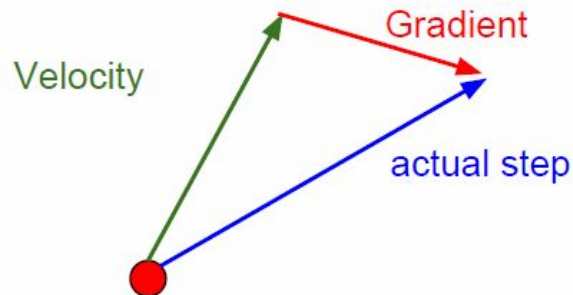

Velocity

actual step

Gradient

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



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

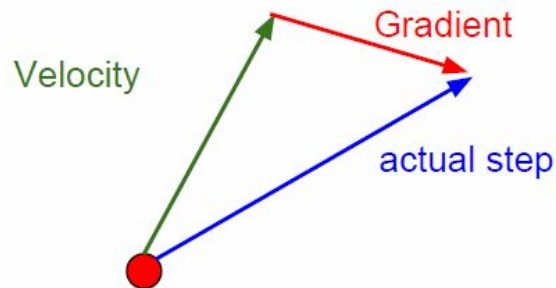$$v_{t+1} = \rho v_t - \alpha \nabla f(\boxed{x_t + \rho v_t})$$

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

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

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

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

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

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



"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

# 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)
```

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"

# RMSProp

**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)
```

**RMSProp**

```python
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)
```

# 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))
```
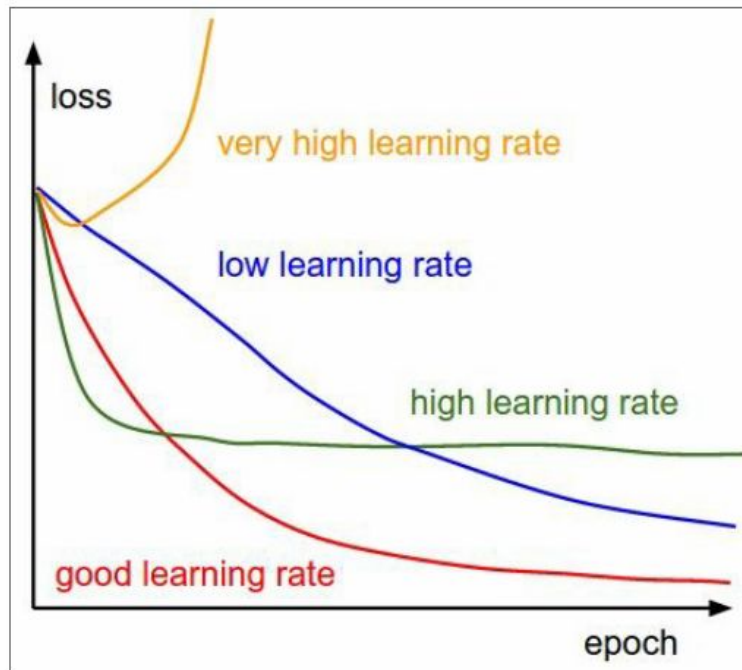
Momentum

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



loss
very high learning rate
low learning rate
high learning rate
good learning rate
epoch

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

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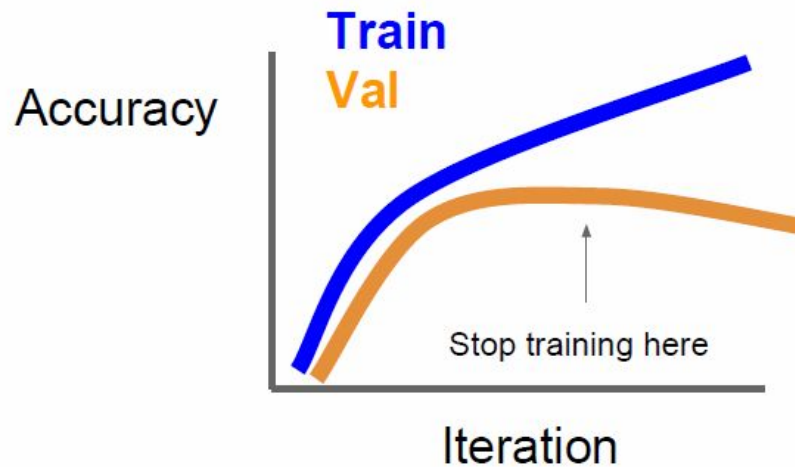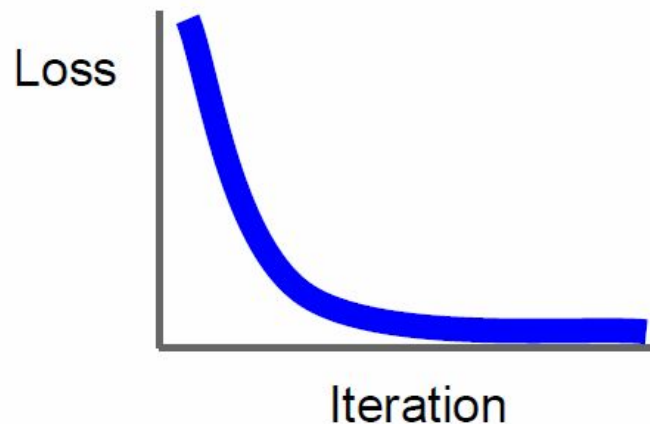e.g. decay learning rate by half every few epochs.

**exponential decay:**

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

**1/t decay:**

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

# Early Stopping



Loss / Iteration

Accuracy / Iteration — **Train** **Val** — Stop training here

Stop training the model when accuracy on the validation set decreases
Or train for a long time, but always keep track of the model snapshot that worked best on val

# Model Ensembles

1.  Train multiple independent models
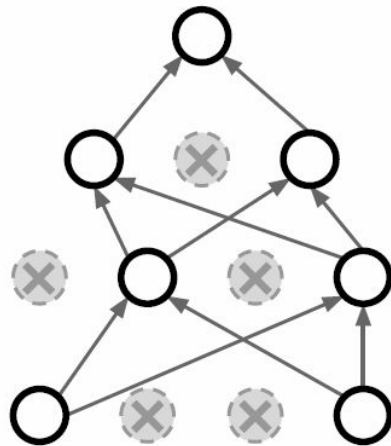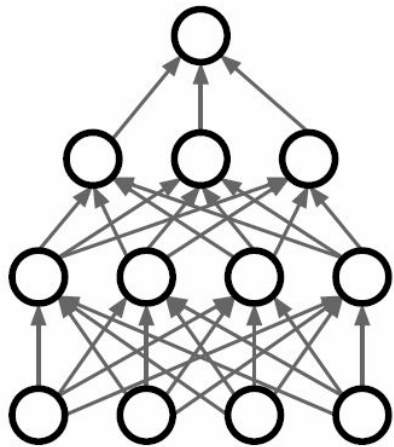2.  At test time average their results
    (Take average of predicted probability distributions, then choose argmax)

Enjoy 2% extra performance

# Dropout

# Regularization: Dropout

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

# Regularization: Dropout

Example forward pass with a 3-layer network using 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)
```
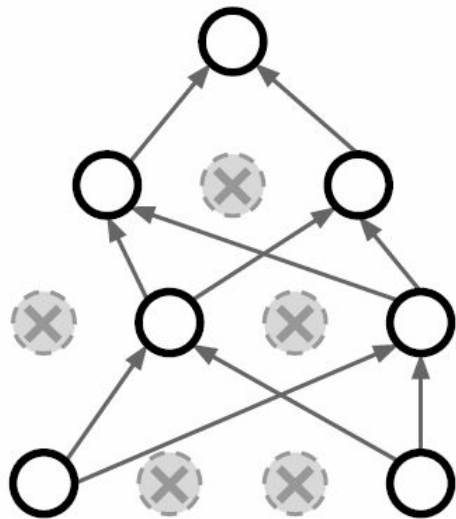
# Regularization: Dropout

How can this possibly be a good idea?



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

has an ear ✗

has a tail

is furry ✗

has claws

mischievous look ✗

cat score

# Regularization: Dropout

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!
Only $\sim 10^{82}$ atoms in the universe...

# 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!