



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>

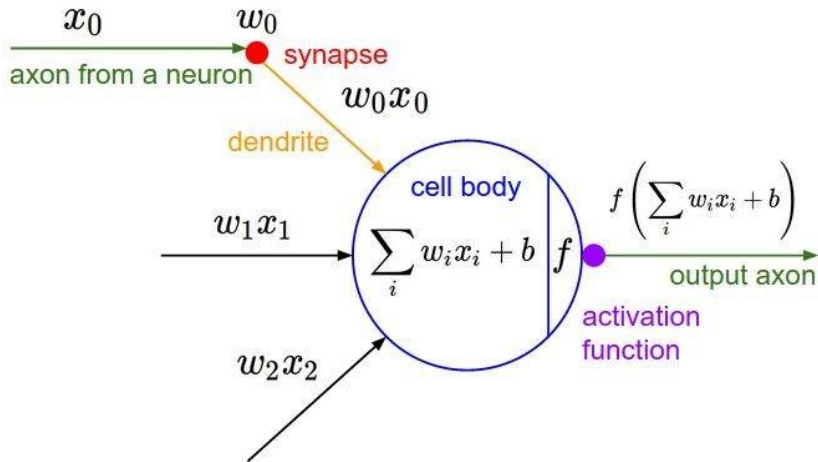
Lecture 4: Training Networks, Part 1

Part 1

- Activation Functions
- Data Preprocessing
- Weight Initialization
- Batch Normalization
- Babysitting the Learning Process
- Hyperparameter Optimization

Activation Functions

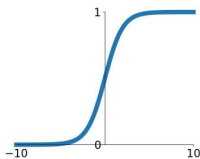
Activation Functions



Activation Functions

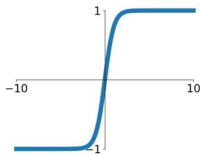
Sigmoid

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



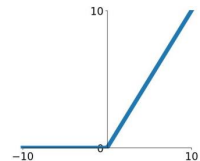
tanh

$$\tanh(x)$$



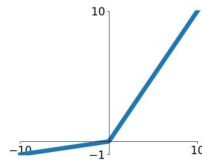
ReLU

$$\max(0, x)$$



Leaky ReLU

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

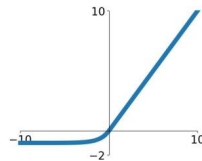


Maxout

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

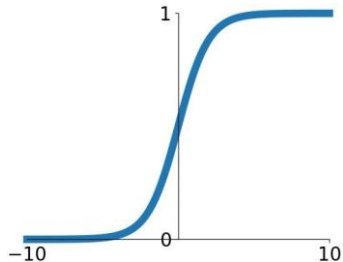
ELU

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



Activation Functions

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



Sigmoid

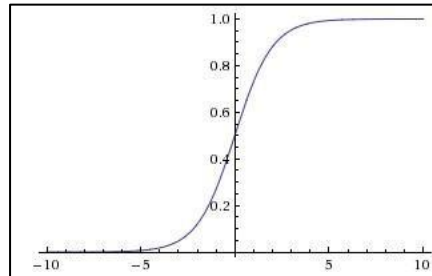
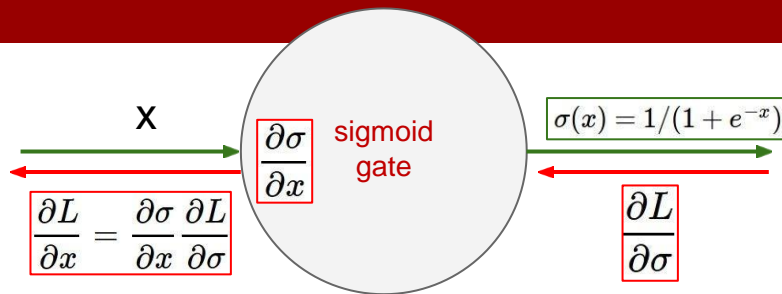
- Squashes numbers to range $[0,1]$
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

3 problems:

1. Saturated neurons “kill” the gradients (**Vanishing gradient**)

Your neuron is **saturated** if it outputs either 0 or 1, then the gradient will be killed. It'll just be multiplied by a very tiny number then gradients can't backpropagate through the network because they'll be stopped learning.

The gradients only flow if you're kind of in a safer zone and what we call an **active region** of a sigmoid



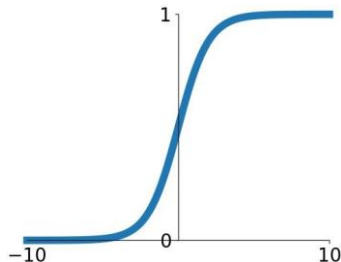
What happens when $x = -10$?

What happens when $x = 0$?

What happens when $x = 10$?

Activation Functions

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



Sigmoid
(logistic function)

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

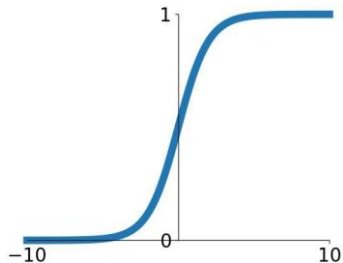
3 problems:

1. Saturated neurons “kill” the gradients
2. Sigmoid outputs are not zero-centered

It causes a zigzag path to reach the minima

Activation Functions

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



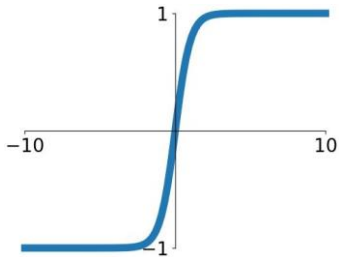
Sigmoid

- 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 computed expensively

Activation Functions



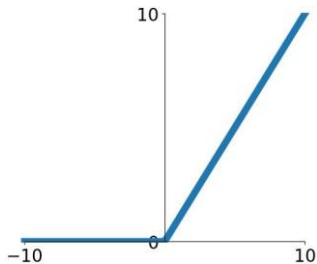
$\tanh(x)$

[LeCun et al., 1991]

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

Activation Functions

Computes $f(x) = \max(0, x)$



ReLU (Rectified Linear Unit)

[Krizhevsky et al., 2012]

- 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
 - Dying ReLU problem—A form of the vanishing gradient problem
hint: what is the gradient when $x < 0$?

Reasons of Dead ReLU

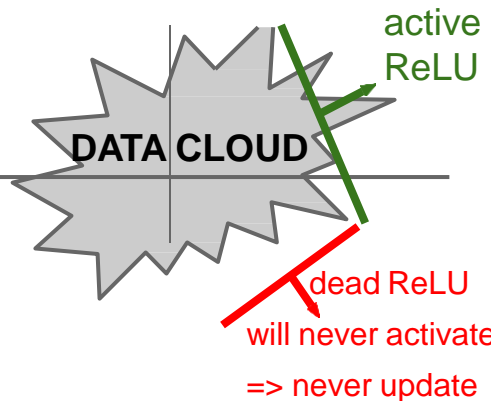
Dying ReLU problem—when inputs approach zero or are negative, the gradient of the function becomes zero, the network cannot perform backpropagation and cannot learn.

The issue can happen when:

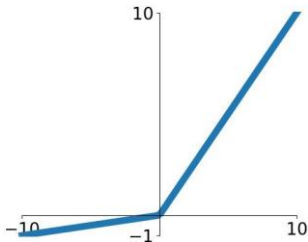
(1) Very **unlucky initialization** of your network may cause that the neurons only activate in the region outside of your data cloud then this dead ReLU you will never become activated and then it will never update.

=> people like to initialize ReLU neurons with slightly positive biases (e.g. 0.01)

2) **High learning rate** can cause saturated neurons.



Activation Functions



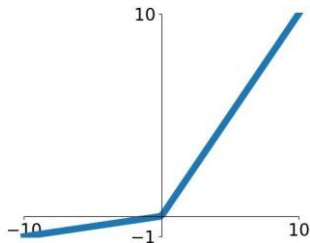
- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- **will not “die”.**

Leaky ReLU

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

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

Activation Functions



Leaky ReLU

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

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

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

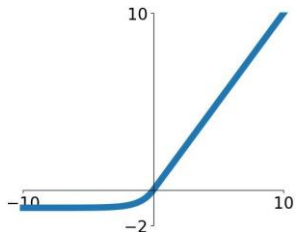
Parametric Rectifier (PReLU)

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

backprop into α
(parameter)

Activation Functions

Exponential Linear Units (ELU)



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

[Clevert et al., 2015]

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

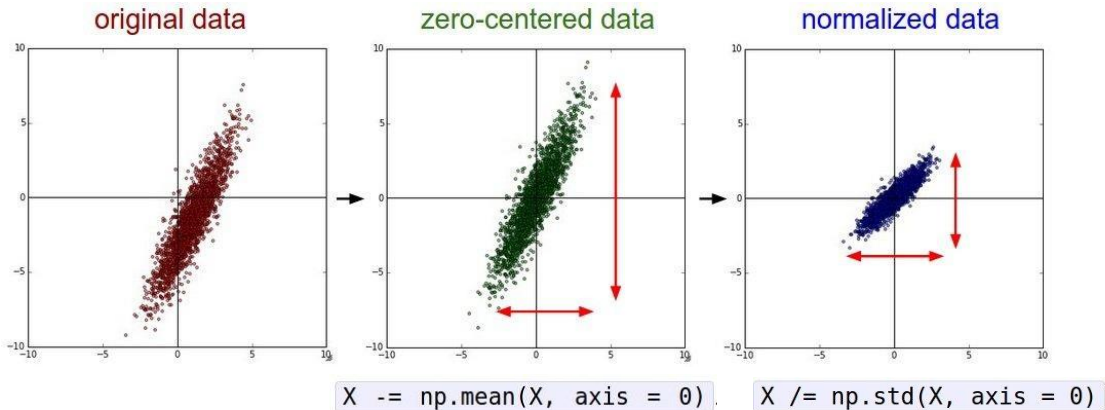
[Goodfellow et al., 2013]

TLDR: In practice:

- 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



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

TLDR: In practice for Images: center only

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

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

```
W = 0.01* np.random.randn(D,H)
```

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



generates an array of a shape (DxH), filled with random floats numbers sampled from a univariate standard normal (Gaussian) distribution.

Lets look at some activation statistics

E.g. 10-layer net with 500 neurons on each layer, using tanh non-linearities, and initializing as described in the last slide.

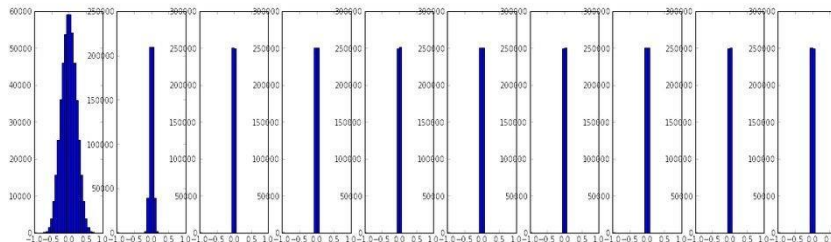
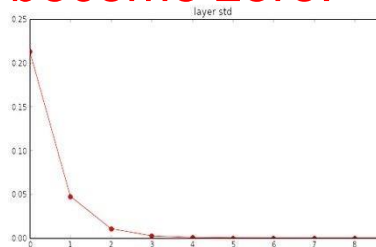
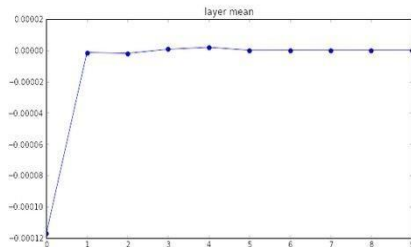
Q: think about the backward pass.
What do the gradients look like?

Hint: think about backward pass for a $W \cdot X$ gate.

```
input layer had mean 0.000927 and std 0.998388
hidden layer 1 had mean -0.000117 and std 0.213081
hidden layer 2 had mean -0.000001 and std 0.047551
hidden layer 3 had mean -0.000002 and std 0.010630
hidden layer 4 had mean 0.000001 and std 0.002378
hidden layer 5 had mean 0.000002 and std 0.000532
hidden layer 6 had mean -0.000000 and std 0.000119
hidden layer 7 had mean 0.000000 and std 0.000026
hidden layer 8 had mean -0.000000 and std 0.000006
hidden layer 9 had mean 0.000000 and std 0.000001
hidden layer 10 had mean -0.000000 and std 0.000000
```

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

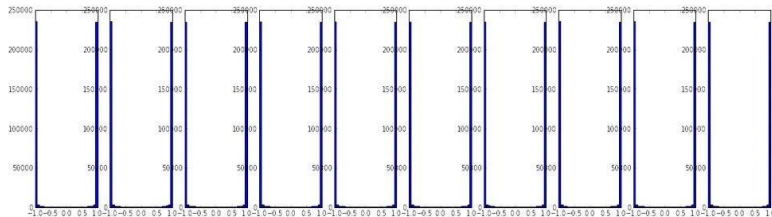
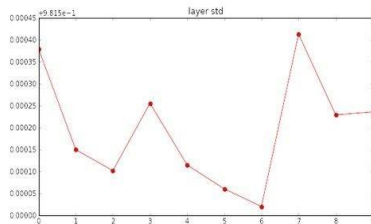
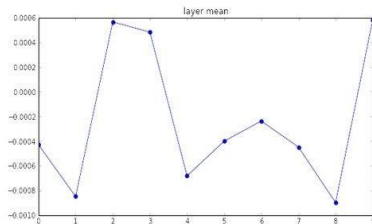
All activations become zero!




```
W = np.random.randn(fan_in, fan_out) * 1.0 # layer initialization
```

input layer had mean 0.001800 and std 1.001311
hidden layer 1 had mean -0.000430 and std 0.981879
hidden layer 2 had mean -0.000849 and std 0.981649
hidden layer 3 had mean 0.000566 and std 0.981601
hidden layer 4 had mean 0.000483 and std 0.981755
hidden layer 5 had mean -0.000682 and std 0.981614
hidden layer 6 had mean -0.000401 and std 0.981560
hidden layer 7 had mean -0.000237 and std 0.981520
hidden layer 8 had mean -0.000448 and std 0.981913
hidden layer 9 had mean -0.000899 and std 0.981728
hidden layer 10 had mean 0.000584 and std 0.981736

1.0 instead of *0.01



Almost all neurons are completely saturated, either -1 or 1.
Gradients will be all zero.

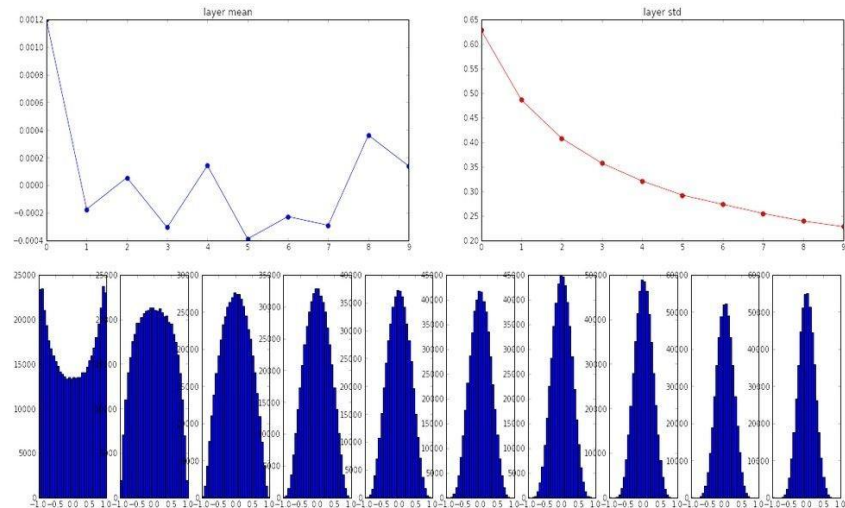
input layer had mean 0.001800 and std 1.001311
 hidden layer 1 had mean 0.001198 and std 0.627953
 hidden layer 2 had mean -0.000175 and std 0.486051
 hidden layer 3 had mean 0.000055 and std 0.407723
 hidden layer 4 had mean -0.000306 and std 0.357108
 hidden layer 5 had mean 0.000142 and std 0.320917
 hidden layer 6 had mean -0.000389 and std 0.292116
 hidden layer 7 had mean -0.000228 and std 0.273387
 hidden layer 8 had mean -0.000291 and std 0.254935
 hidden layer 9 had mean 0.000361 and std 0.239266
 hidden layer 10 had mean 0.000139 and std 0.228008

$W = \text{np.random.randn}(\text{fan_in}, \text{fan_out}) / \text{np.sqrt}(\text{fan_in})$ # layer initialization

“Xavier initialization”
[Glorot et al., 2010]

Reasonable initialization.
(Mathematical derivation assumes linear activations)

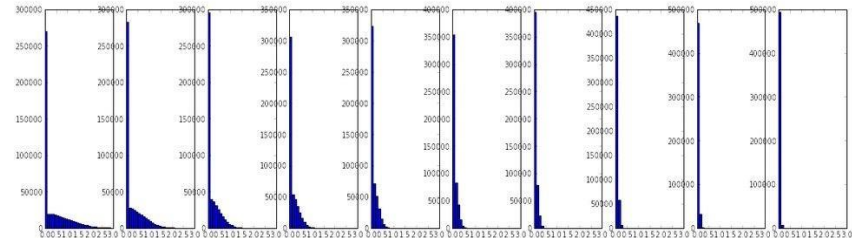
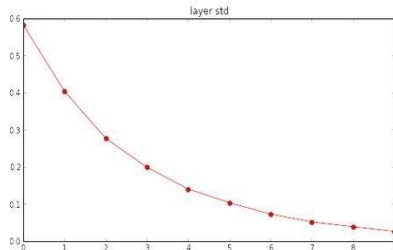
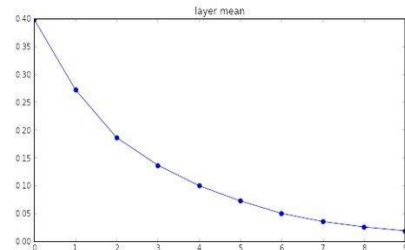
fan-in is the number of inputs to the neuron.
So, with large fan-in => lower weights
Small fan-in => large weights



input layer had mean 0.000501 and std 0.999444
 hidden layer 1 had mean 0.398623 and std 0.582273
 hidden layer 2 had mean 0.272352 and std 0.403795
 hidden layer 3 had mean 0.186076 and std 0.276912
 hidden layer 4 had mean 0.136442 and std 0.198685
 hidden layer 5 had mean 0.099568 and std 0.140299
 hidden layer 6 had mean 0.072234 and std 0.103280
 hidden layer 7 had mean 0.049775 and std 0.072748
 hidden layer 8 had mean 0.035138 and std 0.051572
 hidden layer 9 had mean 0.025404 and std 0.038583
 hidden layer 10 had mean 0.018408 and std 0.026076

`W = np.random.randn(fan_in, fan_out) / np.sqrt(fan_in) # layer initialization`

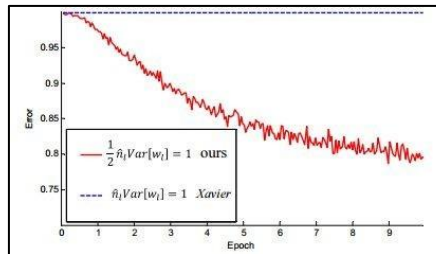
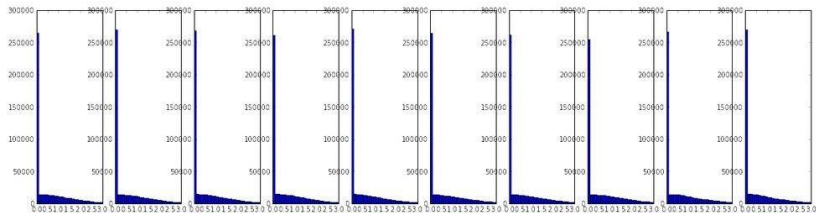
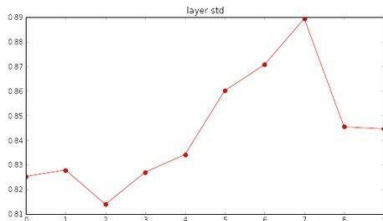
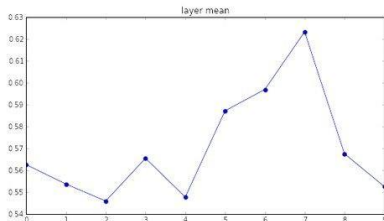
but when using the ReLU
nonlinearity it breaks.



input layer had mean 0.000501 and std 0.999444
 hidden layer 1 had mean 0.562488 and std 0.825232
 hidden layer 2 had mean 0.553614 and std 0.827835
 hidden layer 3 had mean 0.545867 and std 0.813855
 hidden layer 4 had mean 0.565396 and std 0.826902
 hidden layer 5 had mean 0.547678 and std 0.834092
 hidden layer 6 had mean 0.587103 and std 0.860035
 hidden layer 7 had mean 0.596867 and std 0.870610
 hidden layer 8 had mean 0.623214 and std 0.889348
 hidden layer 9 had mean 0.567498 and std 0.845357
 hidden layer 10 had mean 0.552531 and std 0.844523

```
W = np.random.randn(fan_in, fan_out) / np.sqrt(2/fan_in) # layer initialization
```

He et al., 2015
 (note additional 2/)



Batch Normalization (BN) and its variants

Batch Normalization (BN)

[Ioffe and Szegedy, 2015]

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

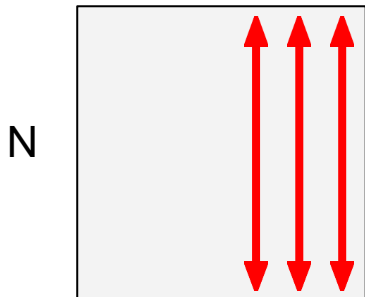
consider a batch of activations at some layer. To make each dimension zero-mean unit-variance, apply:

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

this is a vanilla
differentiable function...

Batch Normalization for fully-connected networks

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



1. compute the empirical mean and variance independently for each dimension.

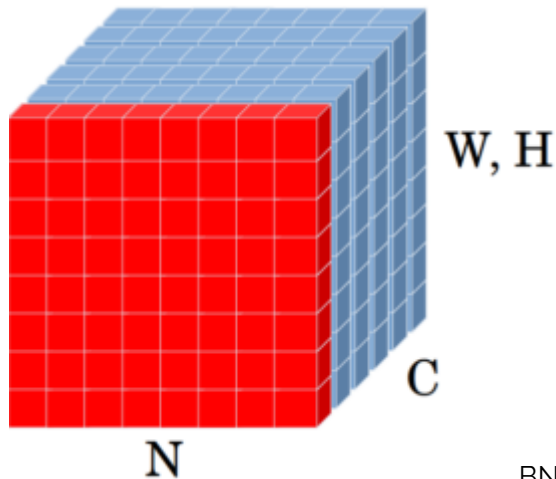
2. Normalize

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

$\mathbf{x} : \mathbf{N} \times \mathbf{D}$
↓
Normalize: $\mathbf{E}, \mathbf{Var} : \mathbf{1} \times \mathbf{D}$

[Ioffe and Szegedy, 2015]

Batch Normalization for ConvNets



Batch Normalization for
convolutional networks
(Spatial Batchnorm, BatchNorm2D)

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

Normalize ↓ ↓ ↓

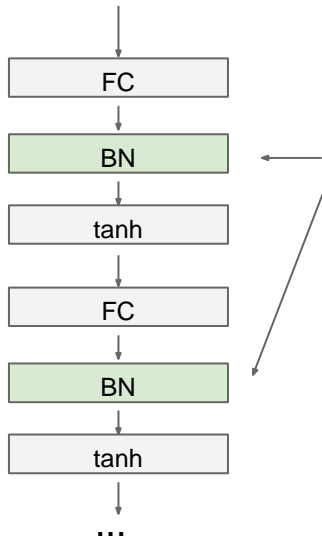
$\mathbf{E}, \mathbf{Var} : \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1}$

Batch of N images, each image has C feature maps with height H and weight W .

BN operates on one feature map over all the training samples in the mini-batch (specified in red):

Batch Normalization

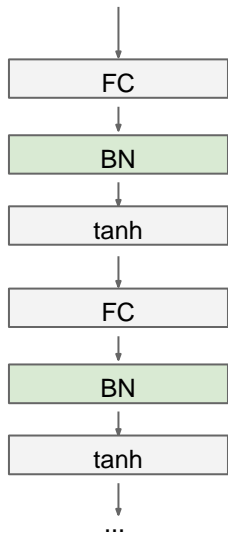
[Ioffe and Szegedy, 2015]



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

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

Batch Normalization



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

Problem: do we necessarily want a zero-mean unit-variance input?

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

Batch Normalization

[Ioffe and Szegedy, 2015]

Normalize:

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

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

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

Note, the network can learn:

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

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

to recover the identity mapping.

Batch Normalization

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1\dots m}\}$;

Parameters to be learned: γ, β

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

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

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

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

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

- 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

[Ioffe and Szegedy, 2015]

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1\dots m}\}$;

Parameters to be learned: γ, β

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

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

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

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

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

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)

Layer Normalization (LN)

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!
Can be used in recurrent networks

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

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

Instance Normalization (IN)

Batch Normalization for
convolutional networks

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

Normalize



$\mu, \sigma : 1 \times C \times 1 \times 1$

$\gamma, \beta : 1 \times C \times 1 \times 1$

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

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

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

Normalize



$\mu, \sigma : N \times C \times 1 \times 1$

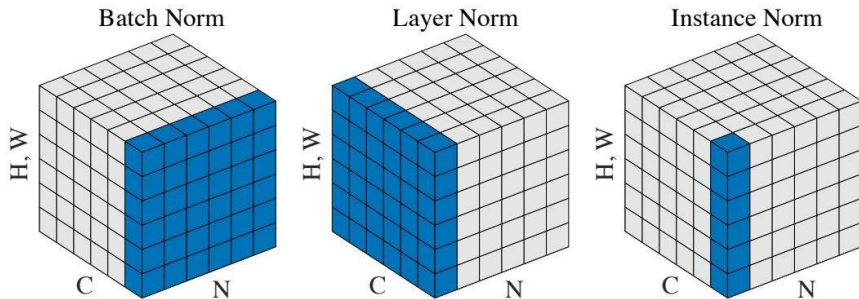
$\gamma, \beta : 1 \times C \times 1 \times 1$

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

Ulyanov et al, Improved Texture Networks: Maximizing Quality and Diversity in Feed-forward Stylization and Texture Synthesis, CVPR 2017

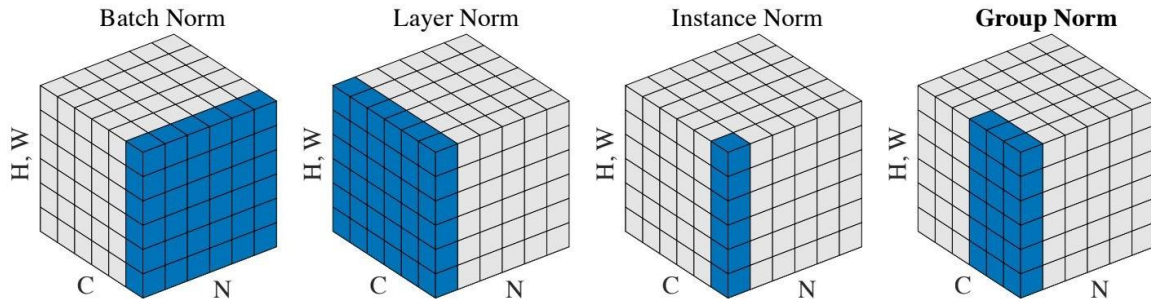
IN performs normalization across the width and height of a single feature map of a single example

Comparison of Normalization Layers



Wu and He, "Group Normalization", arXiv 2018

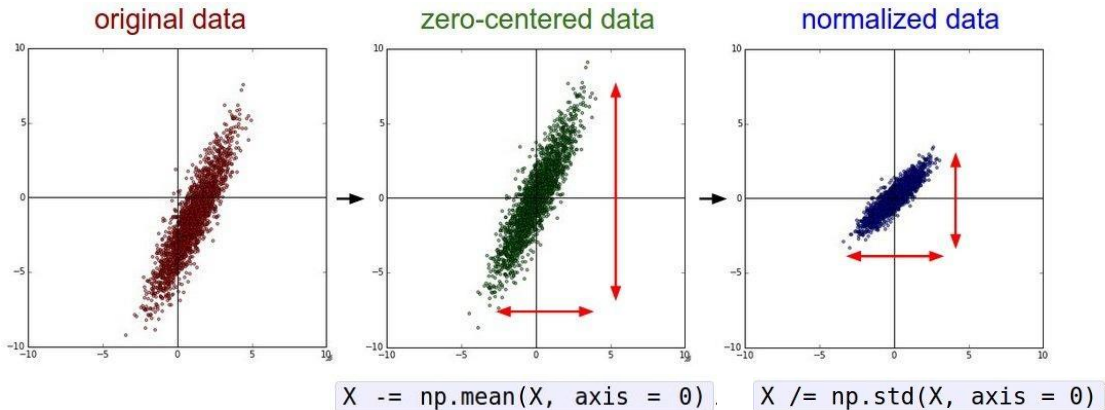
Group Normalization (GN)



Wu and He, "Group Normalization", arXiv 2018 (Appeared 3/22/2018)

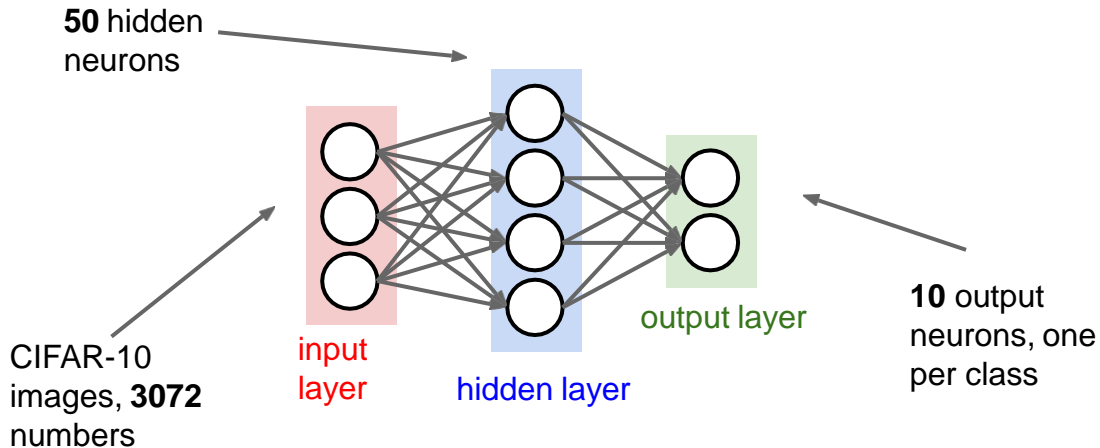
Babysitting the Learning Process

Step 1: Preprocess the data



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

Step 2: Choose the architecture:



Step 3: Network Sanity Checklist

1- Double check that **the loss** is reasonable:

a) **Disable the regularization parameter (0.0)** that is passed in the network and make sure that the **loss** comes out is **correct**.

b) **Add a big value for the regularization parameter ($1e3$)**, it is expected that the **loss go up** because we have this additional term in the objective function.

Lets try to train now...

```
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)
```

2- Take a **small portion of the training data**, and train the network. Make sure that you get a **very small loss (near zero)**, and **train accuracy of 1.00**. So, you fully overfit your training data because if you can't overfit a tiny piece of data then things are definitely broken.

The above code:

- take the first 20 examples from CIFAR-10
- turn off regularization (reg = 0.0)
- use simple vanilla 'sgd'

Lets try to train now...

Start with small regularization and find learning rate that makes the loss go down.

```
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)
```

Lets try to train now...

Start with small regularization and find learning rate that makes the loss go down.

```
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 | | |

Loss barely changing

Lets try to train now...

Start with small regularization and find learning rate that makes the loss go down.

loss not going down:
learning rate too low

```
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)
```

| Epoch | Cost | Train Loss | Val Loss | Learning Rate |
|---------|----------|------------|----------|---------------|
| 1 / 10 | 2.302576 | 0.080000 | 0.103000 | 1.000000e-06 |
| 2 / 10 | 2.302582 | 0.121000 | 0.124000 | 1.000000e-06 |
| 3 / 10 | 2.302558 | 0.119000 | 0.138000 | 1.000000e-06 |
| 4 / 10 | 2.302519 | 0.127000 | 0.151000 | 1.000000e-06 |
| 5 / 10 | 2.302517 | 0.158000 | 0.171000 | 1.000000e-06 |
| 6 / 10 | 2.302518 | 0.179000 | 0.172000 | 1.000000e-06 |
| 7 / 10 | 2.302466 | 0.180000 | 0.176000 | 1.000000e-06 |
| 8 / 10 | 2.302452 | 0.175000 | 0.185000 | 1.000000e-06 |
| 9 / 10 | 2.302459 | 0.206000 | 0.192000 | 1.000000e-06 |
| 10 / 10 | 2.302420 | 0.190000 | 0.192000 | 1.000000e-06 |

finished optimization. best validation accuracy: 0.192000

Loss barely changing: Learning rate is probably too low

Lets try to train now...

Start with small regularization and find learning rate that makes the loss go down.

```
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=3e-3, verbose=True)
```

```
Finished epoch 1 / 10: cost 2.186654, train: 0.308000, val 0.306000, lr 3.000000e-03
Finished epoch 2 / 10: cost 2.176230, train: 0.330000, val 0.350000, lr 3.000000e-03
Finished epoch 3 / 10: cost 1.942257, train: 0.376000, val 0.352000, lr 3.000000e-03
Finished epoch 4 / 10: cost 1.827868, train: 0.329000, val 0.310000, lr 3.000000e-03
Finished epoch 5 / 10: cost inf, train: 0.128000, val 0.128000, lr 3.000000e-03
Finished epoch 6 / 10: cost inf, train: 0.144000, val 0.147000, lr 3.000000e-03
```

3e-3 is too high. Cost explodes....

loss not going down:
learning rate too low
loss exploding:
learning rate too high

=> Rough range for learning rate we should be cross-validating is somewhere [1e-3 ... 1e-5]

Hyperparameter Optimization

Cross-validation strategy

coarse -> fine cross-validation in stages

First stage: only a few epochs to get rough idea of what params work

Second stage: longer running time, finer search

... (repeat as necessary)

Tip for detecting explosions in the solver:

If the cost is ever $> 3 \times$ original cost, break out early

For example: run coarse search for 5 epochs

```
max_count = 100
for count in xrange(max_count):
    reg = 10**uniform(-5, 5)
    lr = 10**uniform(-3, -6)

    trainer = ClassifierTrainer()
    model = init_two_layer_model(32*32*3, 50, 10) # input size, hidden size, number of classes
    trainer = ClassifierTrainer()
    best_model_local, stats = trainer.train(X_train, y_train, X_val, y_val,
                                           model, two_layer_net,
                                           num_epochs=5, reg=reg,
                                           update='momentum', learning_rate_decay=0.9,
                                           sample_batches = True, batch_size = 100,
                                           learning_rate=lr, verbose=False)
```

note it's best to optimize
in log space!

```
val_acc: 0.412000, lr: 1.405206e-04, reg: 4.793564e-01, (1 / 100)
val_acc: 0.214000, lr: 7.231888e-06, reg: 2.321281e-04, (2 / 100)
val_acc: 0.208000, lr: 2.119571e-06, reg: 8.011857e+01, (3 / 100)
val_acc: 0.196000, lr: 1.551131e-05, reg: 4.374936e-05, (4 / 100)
val_acc: 0.079000, lr: 1.753300e-05, reg: 1.200424e+03, (5 / 100)
val_acc: 0.223000, lr: 4.215128e-05, reg: 4.196174e+01, (6 / 100)
val_acc: 0.441000, lr: 1.750259e-04, reg: 2.110807e-04, (7 / 100)
val_acc: 0.241000, lr: 6.749231e-05, reg: 4.226413e+01, (8 / 100)
val_acc: 0.482000, lr: 4.296863e-04, reg: 6.642555e-01, (9 / 100)
val_acc: 0.079000, lr: 5.401602e-06, reg: 1.599828e+04, (10 / 100)
val_acc: 0.154000, lr: 1.618508e-06, reg: 4.925252e-01, (11 / 100)
```

nice

Now run finer search...

```
max_count = 100
for count in xrange(max_count):
    reg = 10**uniform(-5, 5)
    lr = 10**uniform(-3, -6)
```

adjust range



```
max_count = 100
for count in xrange(max_count):
    reg = 10**uniform(-4, 0)
    lr = 10**uniform(-3, -4)
```

```
val_acc: 0.527000, lr: 5.340517e-04, reg: 4.097824e-01, (0 / 100)
val_acc: 0.492000, lr: 2.279484e-04, reg: 9.991345e-04, (1 / 100)
val_acc: 0.512000, lr: 8.680827e-04, reg: 1.349727e-02, (2 / 100)
val_acc: 0.461000, lr: 1.028377e-04, reg: 1.220193e-02, (3 / 100)
val_acc: 0.460000, lr: 1.113730e-04, reg: 5.244309e-02, (4 / 100)
val_acc: 0.498000, lr: 9.477776e-04, reg: 2.001293e-03, (5 / 100)
val_acc: 0.469000, lr: 1.484369e-04, reg: 4.328313e-01, (6 / 100)
val_acc: 0.522000, lr: 5.586261e-04, reg: 2.312685e-04, (7 / 100)
val_acc: 0.530000, lr: 5.808183e-04, reg: 8.259964e-02, (8 / 100)
val_acc: 0.489000, lr: 1.979168e-04, reg: 1.010889e-04, (9 / 100)
val_acc: 0.490000, lr: 2.036031e-04, reg: 2.406271e-03, (10 / 100)
val_acc: 0.475000, lr: 2.021162e-04, reg: 2.287807e-01, (11 / 100)
val_acc: 0.460000, lr: 1.135527e-04, reg: 3.905040e-02, (12 / 100)
val_acc: 0.515000, lr: 6.947668e-04, reg: 1.562808e-02, (13 / 100)
val_acc: 0.531000, lr: 9.471549e-04, reg: 1.433895e-03, (14 / 100)
val_acc: 0.509000, lr: 3.140888e-04, reg: 2.857518e-01, (15 / 100)
val_acc: 0.514000, lr: 6.438349e-04, reg: 3.033781e-01, (16 / 100)
val_acc: 0.502000, lr: 3.921784e-04, reg: 2.707126e-04, (17 / 100)
val_acc: 0.509000, lr: 9.752279e-04, reg: 2.850865e-03, (18 / 100)
val_acc: 0.500000, lr: 2.412048e-04, reg: 4.997821e-04, (19 / 100)
val_acc: 0.466000, lr: 1.319314e-04, reg: 1.189915e-02, (20 / 100)
val_acc: 0.516000, lr: 8.039527e-04, reg: 1.528291e-02, (21 / 100)
```

53% - relatively good
for a 2-layer neural net
with 50 hidden neurons.

Now run finer search...

```
max_count = 100
for count in xrange(max_count):
    reg = 10**uniform(-5, 5)
    lr = 10**uniform(-3, -6)
```

adjust range

```
max_count = 100
for count in xrange(max_count):
    reg = 10**uniform(-4, 0)
    lr = 10**uniform(-3, -4)
```

| |
|--|
| val_acc: 0.527000, lr: 5.340517e-04, reg: 4.097824e-01, (0 / 100) |
| val_acc: 0.492000, lr: 2.279484e-04, reg: 9.991345e-04, (1 / 100) |
| val_acc: 0.512000, lr: 8.680827e-04, reg: 1.349727e-02, (2 / 100) |
| val_acc: 0.461000, lr: 1.028377e-04, reg: 1.220193e-02, (3 / 100) |
| val_acc: 0.460000, lr: 1.113730e-04, reg: 5.244309e-02, (4 / 100) |
| val_acc: 0.498000, lr: 9.477776e-04, reg: 2.001293e-03, (5 / 100) |
| val_acc: 0.469000, lr: 1.484369e-04, reg: 4.328313e-01, (6 / 100) |
| val_acc: 0.522000, lr: 5.586261e-04, reg: 2.312685e-04, (7 / 100) |
| val_acc: 0.530000, lr: 5.808183e-04, reg: 8.259964e-02, (8 / 100) |
| val_acc: 0.489000, lr: 1.979168e-04, reg: 1.010889e-04, (9 / 100) |
| val_acc: 0.490000, lr: 2.036031e-04, reg: 2.406271e-03, (10 / 100) |
| val_acc: 0.475000, lr: 2.021162e-04, reg: 2.287807e-01, (11 / 100) |
| val_acc: 0.460000, lr: 1.135527e-04, reg: 3.905040e-02, (12 / 100) |
| val_acc: 0.515000, lr: 6.947668e-04, reg: 1.562808e-02, (13 / 100) |
| val_acc: 0.531000, lr: 9.471549e-04, reg: 1.433895e-03, (14 / 100) |
| val_acc: 0.509000, lr: 3.140888e-04, reg: 2.857518e-01, (15 / 100) |
| val_acc: 0.514000, lr: 6.438349e-04, reg: 3.033781e-01, (16 / 100) |
| val_acc: 0.502000, lr: 3.921784e-04, reg: 2.707126e-04, (17 / 100) |
| val_acc: 0.509000, lr: 9.752279e-04, reg: 2.850865e-03, (18 / 100) |
| val_acc: 0.500000, lr: 2.412048e-04, reg: 4.997821e-04, (19 / 100) |
| val_acc: 0.466000, lr: 1.319314e-04, reg: 1.189915e-02, (20 / 100) |
| val_acc: 0.516000, lr: 8.039527e-04, reg: 1.528291e-02, (21 / 100) |

53% - relatively good for a 2-layer neural net with 50 hidden neurons.

But this best cross-validation result is worrying. Why?

Random Search vs. Grid Search

*Random Search for
Hyper-Parameter Optimization
Bergstra and Bengio, 2012*

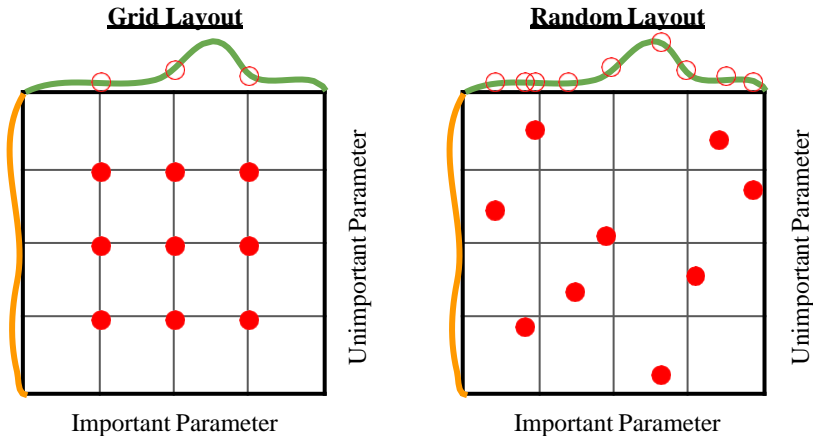
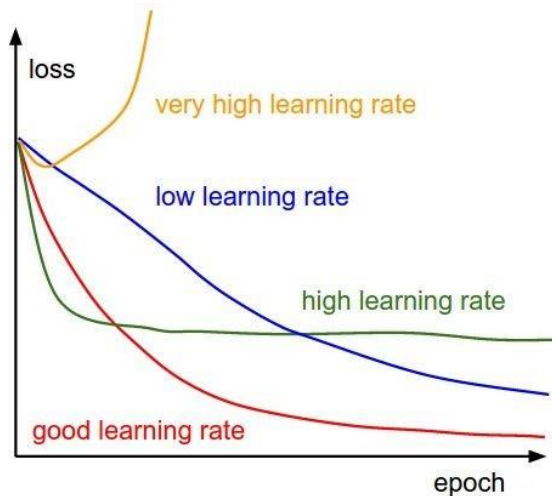
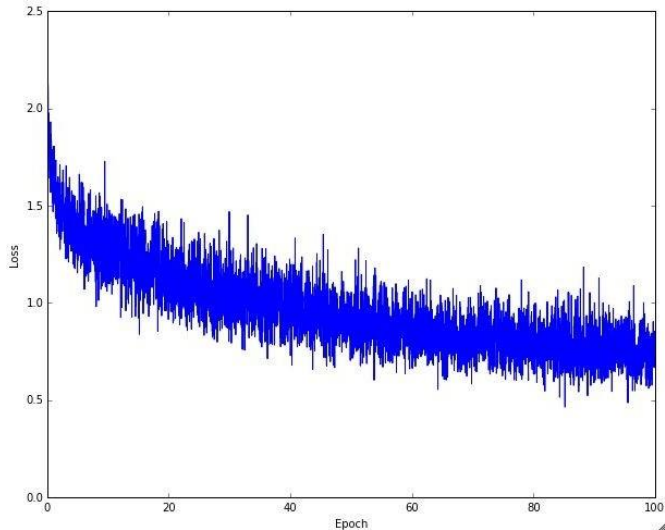


Illustration of Bergstra et al., 2012 by Shayne
Longpre, copyright CS231n 2017

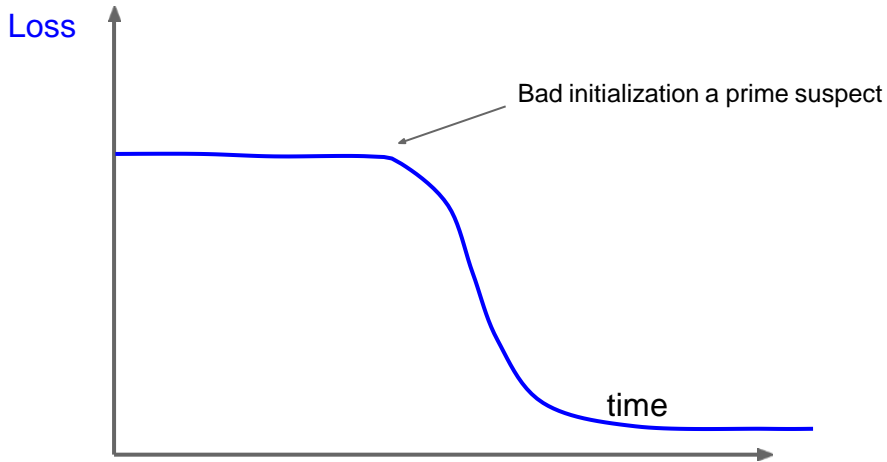
Hyperparameters to play with:

- Network architecture
- learning rate, its decay schedule, update type
- regularization (L2/Dropout strength)

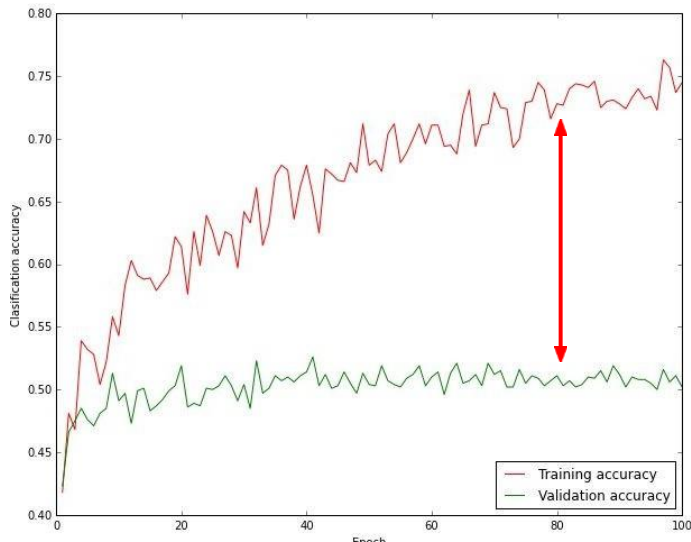
Monitor and visualize the loss curve



What is wrong with this curve?



Monitor and visualize the accuracy:



big gap = overfitting

=> increase regularization strength?

no gap

=> increase model capacity?

Track the ratio of weight updates/ weight magnitudes:

Tracking the difference between:

(a) the scale of your parameters and (b) the scale of updates to those parameters

You don't want those updates to be much larger than the weights and you don't want them to be tiny. Ex: you don't want your updates to be on the order of $1e-7$ when your weights are on the order of $1e-2$.

Let's look at the parameters' norm and compare it to the updated scale of your parameters. It is usually a good “**rule of thumb**” is this should be roughly $1e-3$.

→ So, you're not making huge updates or very small updates.

If this ratio is **too high**, you **decrease** the learning rate, and if it is **too low**, you **increase** the learning rate.

Note: The norm of a matrix is **a measure of how large its elements are**. It is a real number that is a measure of the magnitude of the matrix.

Track the ratio of weight updates/ weight magnitudes:

```
# assume parameter vector W and its gradient vector dW
param_scale = np.linalg.norm(W.ravel())
update = -learning_rate*dW # simple SGD update
update_scale = np.linalg.norm(update.ravel())
W += update # the actual update
print update_scale / param_scale # want ~1e-3
```

ratio between the updates and values: $\sim 0.0002 / 0.02 = 0.01$ (about okay)
want this to be somewhere around 0.001 or so.

Final Grading Distribution

| | | |
|---------------------|-----|-------|
| • Midterm | 20% | |
| • Lab participation | 3% | Bonus |
| • Project | 20% | |
| • Final Exam | 60% | |

CS396: Selected Topics CS-2 Projects_Spring 2021-2022

Team **Members: 6-8**

Description:

1. As a team, you will choose a recent paper (2018-2022) on one of these topics:
 - a. Multi-class classification using CNN or pre-trained models.
 - b. Image-to-image translation using GAN implementation.
 - c. Object detection using deep techniques.
2. Then, choose another image dataset of your selection (not in the paper you had chosen), and apply the same algorithms on the chosen dataset.

Grading:

Total: 20 marks

- 5 marks: Presentation on the paper. ✓
- 10 marks: Implementation ✓
- 5 marks: Individual Assessment.

Project Bonus:

Dataset Bonus: As a team, you will be awarded **+2** marks as a bonus if you used unique datasets (not being used by another team)

Optimization of model: As a team, you will be awarded **+2** marks as a bonus if you optimize your model to reach an accuracy greater than 97%. But you should include in your report the hyperparameters and all results details before and after optimization.

We looked in detail at:

- Activation Functions (use ReLU)
- Data Preprocessing (images: subtract mean)
- Weight Initialization (use Xavier/He init)
- Batch Normalization (use)
- Babysitting the Learning process
- Hyperparameter Optimization
(random sample hyperparams, in log space when appropriate)