

CS109 – Data Science

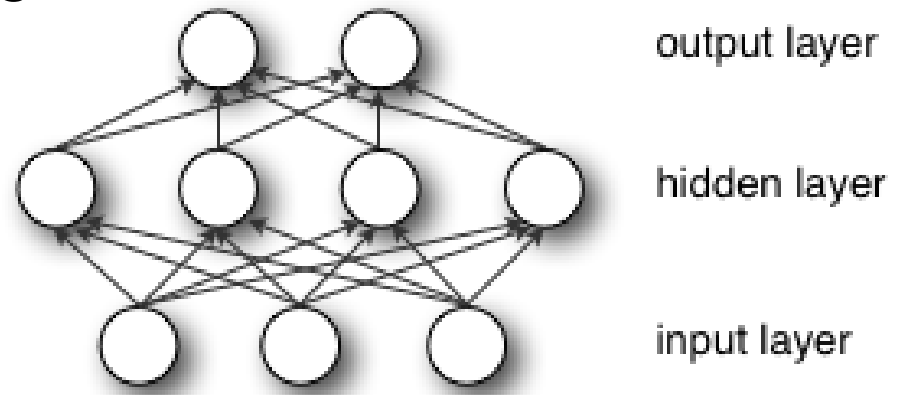
Deep Learning III - Tips

Hanspeter Pfister, Mark Glickman, Verena Kaynig-Fittkau

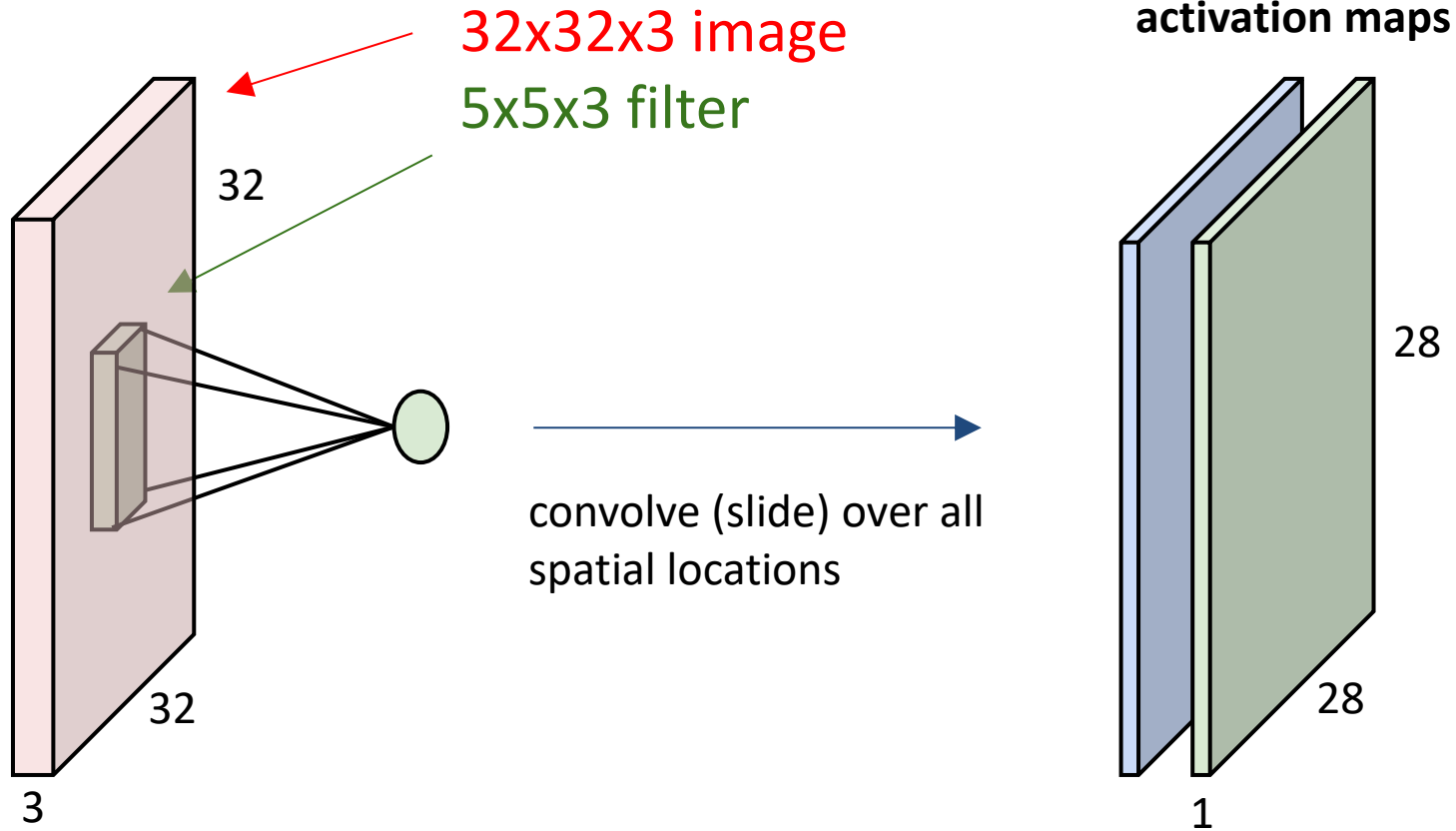


To Train a Simple Network We Need:

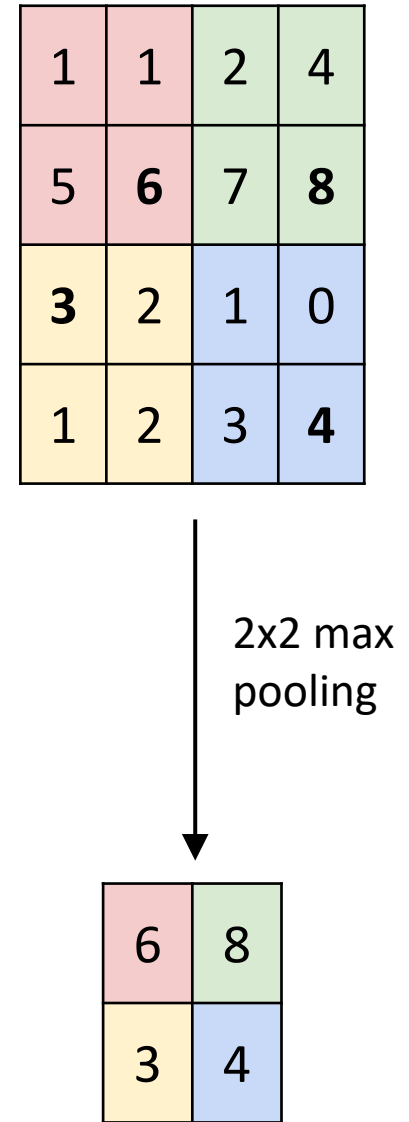
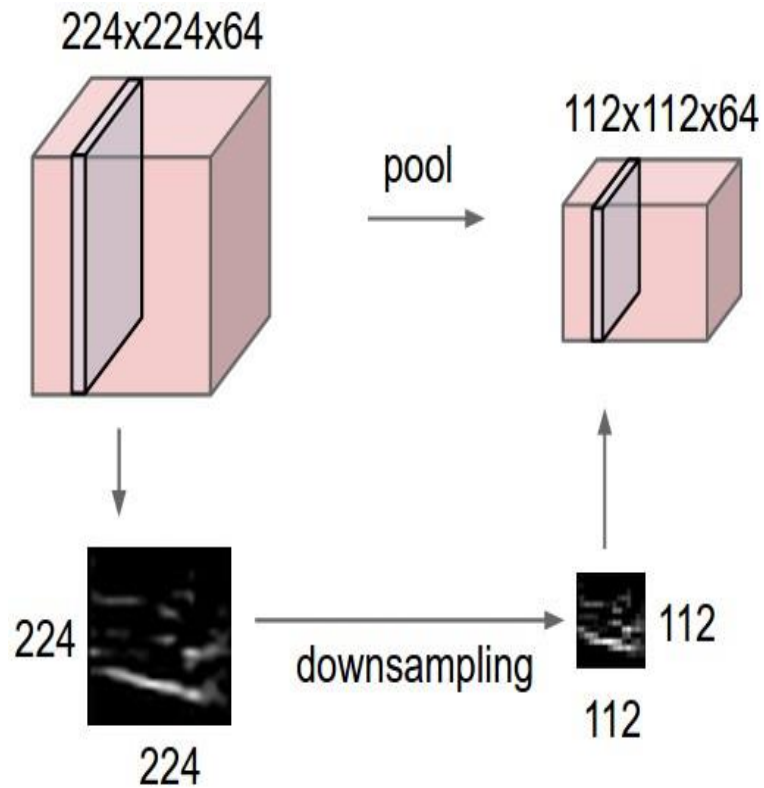
- Input layer size
 - Number of hidden layers
 - Sizes of hidden layers
 - Activation function
 - Number of output units
-
- Loss function
 - Optimization method



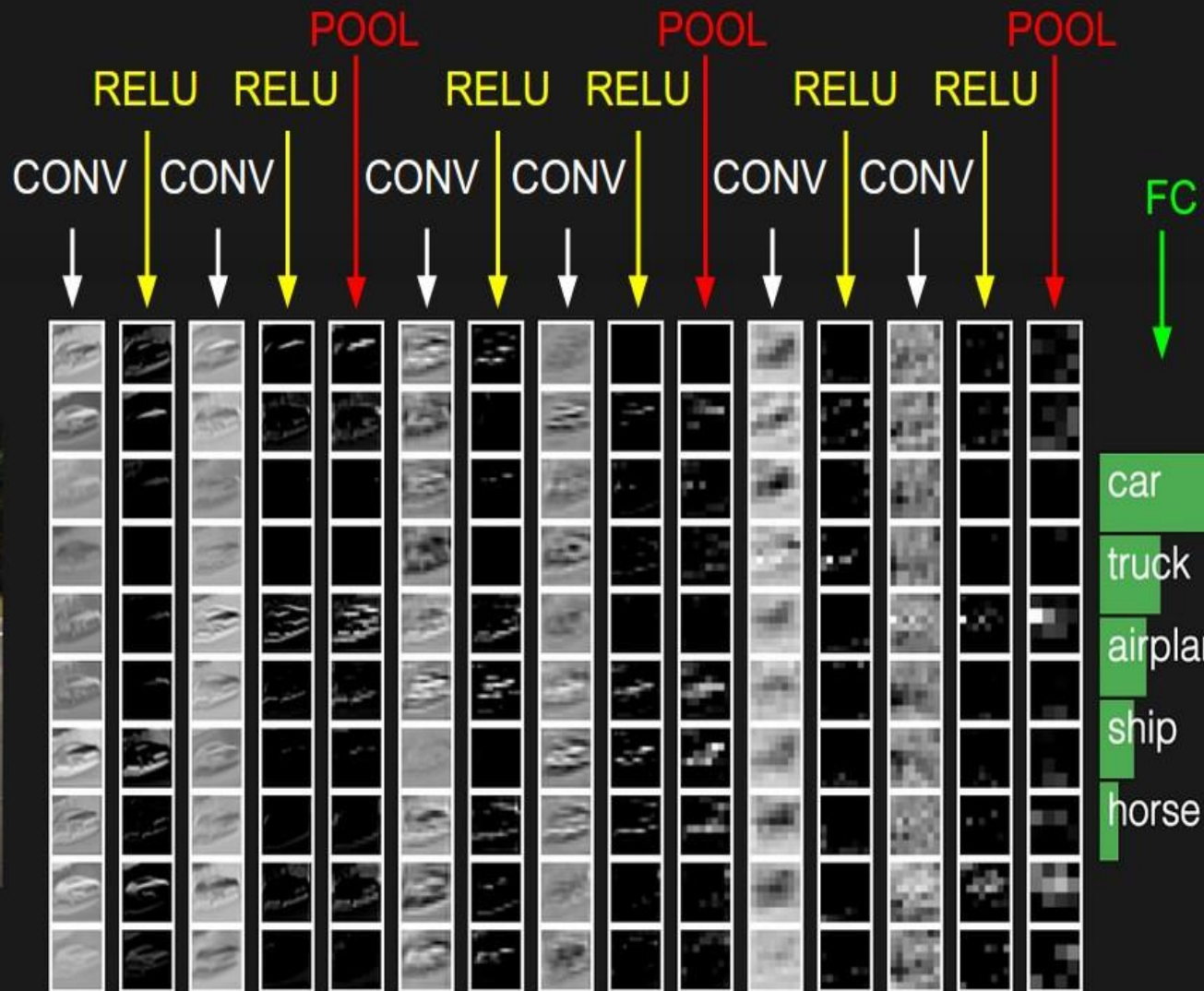
Convolution Layer



Pooling



Typical CNN



Keras Example CNN

Sanity Check

- Try to overfit small portion of training data
- Example:
 - take the first 20 examples from CIFAR-10
 - turn off regularization ($\text{reg} = 0.0$)
 - use simple vanilla 'sgd'
- If this doesn't work something is seriously wrong, as in you are using the library incorrectly.

Today's Lecture



Some Essential Things:

- When do you stop training?
- Strategies for learning rate updates
- Dropout and regularization

When to stop training

- Fixed number of epochs
- When the training converged
 - How do you measure convergence?
 - Optimization becomes too slow
 - Validation score doesn't improve
 - Combine with patience counter

Early Stopping in Keras

```
from keras.callbacks import EarlyStopping

early_stopping = EarlyStopping(monitor='val_loss',
                               patience=2)

model.fit(X, y, validation_split=0.2,
          callbacks=[early_stopping])
```

patience: number of epochs with no improvement after which training will be stopped.

More info: <https://keras.io/callbacks/>

Early Stopping in Keras

```
from keras.callbacks import EarlyStopping

early_stopping = EarlyStopping(monitor='val_loss',
                                patience=2)

model.fit(X, y, validation_split=0.2,
          callbacks=[early_stopping])
```

validation_split: fraction of training data used for validation. Data is split from the end of the data and remains the same during training.

More info: <https://keras.io/callbacks/>

Learning Rate

- So many options:
- Fixed learning rate
- Start with large LR for n epochs, then set to smaller value
- Slowly decay learning rate over time
- Start with large learning rate, when patience counter expires lower the rate, repeat.

Learning Rate Decay in Keras:

```
def step_decay(epoch):  
    lrate = 0.1  
    if epoch > 100:  
        lrate = 0.1 / (2. ** epoch)  
    return lrate  
  
lr = LearningRateScheduler(step_decay)  
callbacks_list = [lr]
```

schedule: a function that takes an epoch index as input (integer, indexed from 0) and returns a new learning rate as output (float).

Reduce LR on Plateau

```
reduce_lr = ReduceLROnPlateau(monitor='val_loss',  
                               factor=0.2,  
                               patience=5,  
                               min_lr=0.001)  
  
model.fit(X_train, Y_train, callbacks=[reduce_lr])
```

factor: How much to reduce the learning rate

patience: How many epochs with no improvement until lr is updated

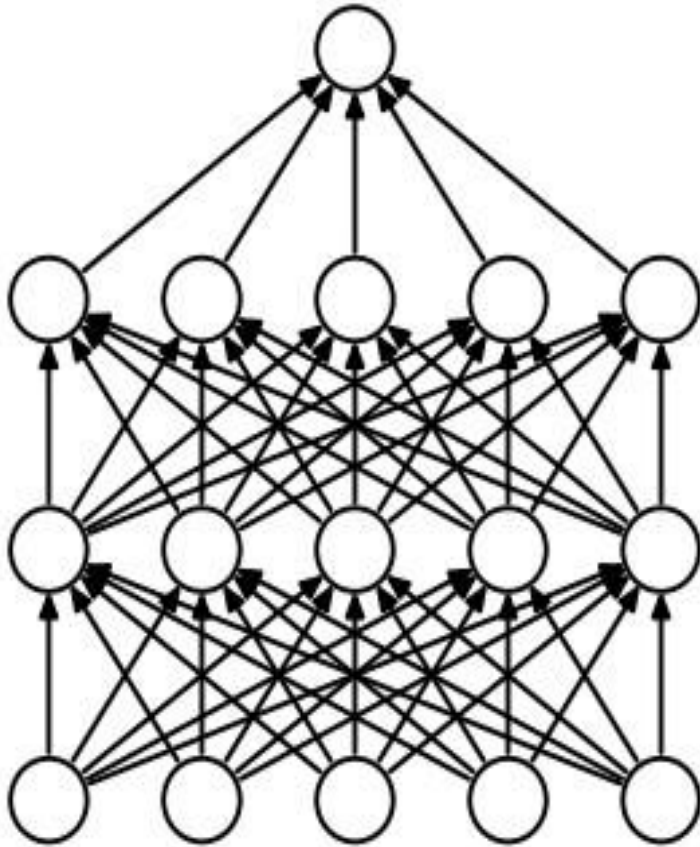
min_lr: When to completely stop

Regularization: dropout

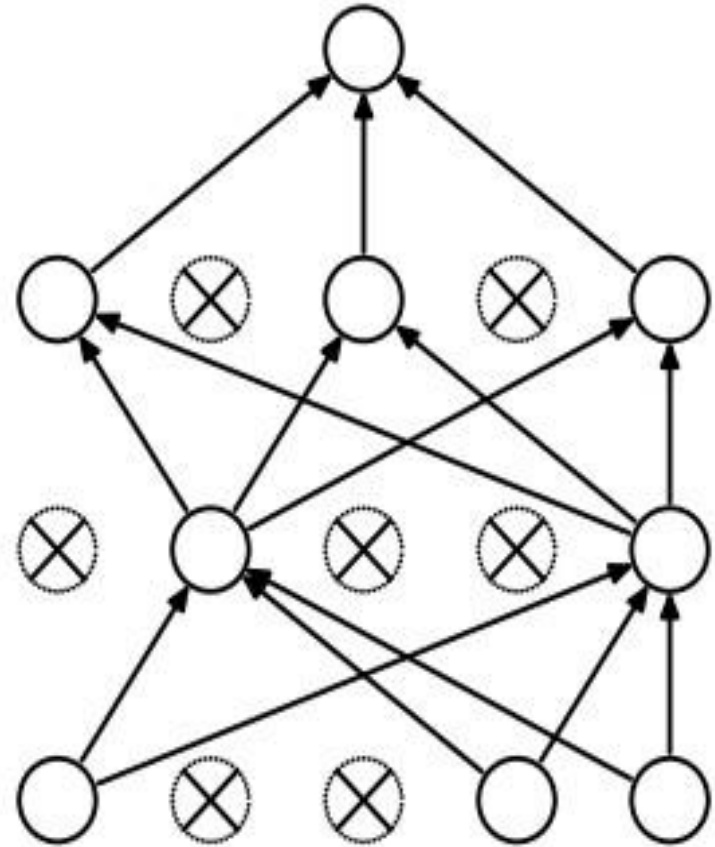
- We have seen that randomness can be used very effectively for regularization
- Remember Random forest
- Can we do something similar for deep learning?

Regularization: **Dropout**

“randomly set some neurons to zero”



(a) Standard Neural Net



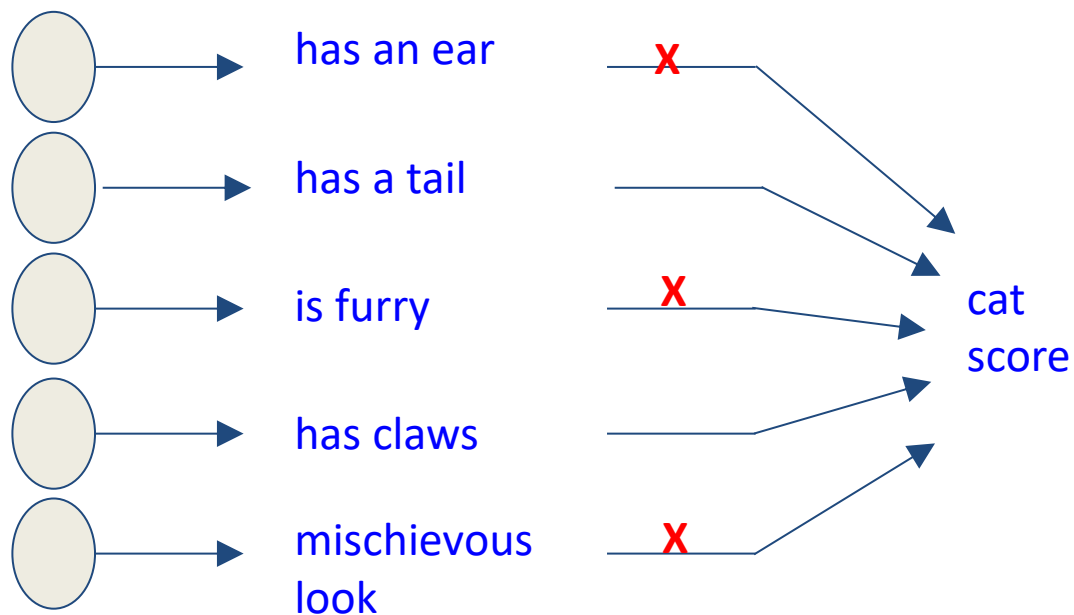
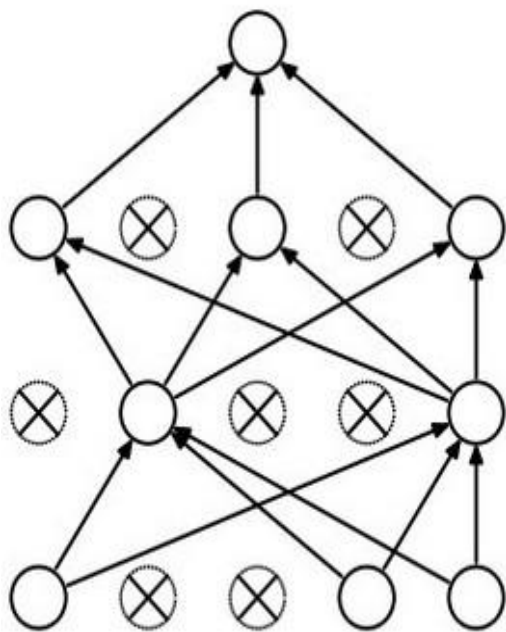
(b) After applying dropout.

[Srivastava et al., 2014]

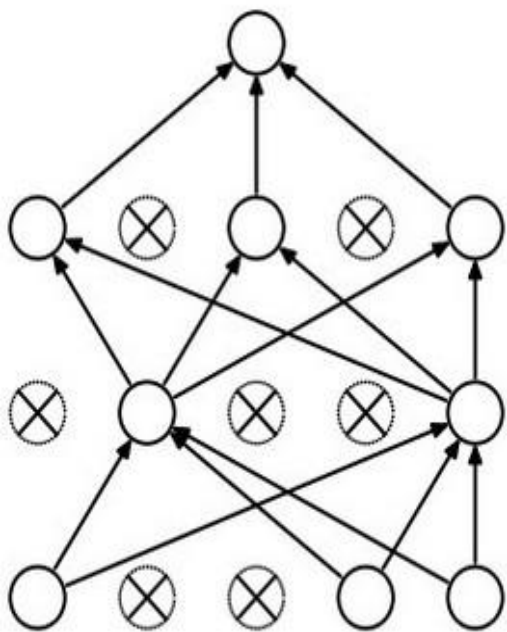
<http://cs231n.github.io/>

How could this possibly be a good idea?

Forces the network to have a redundant representation.



How could this possibly be a good idea?

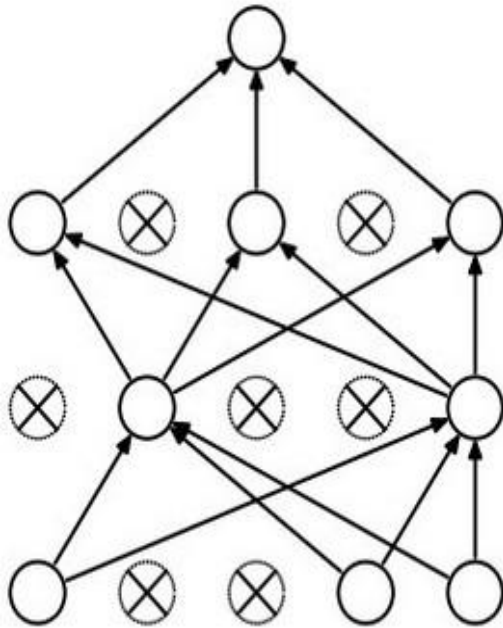


Another interpretation:

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

Each binary mask is one model, gets trained with only \sim one update.

At test time....



Ideally:

want to integrate out all the noise

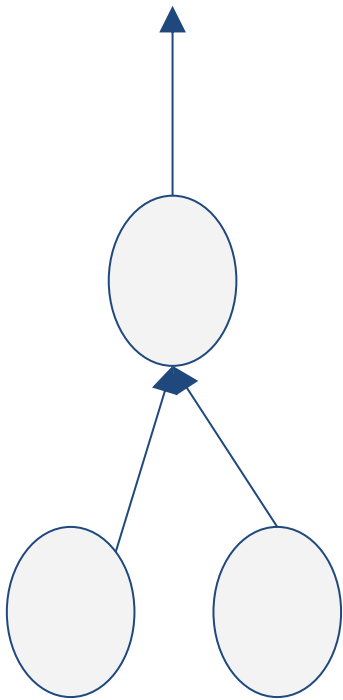
Monte Carlo approximation:

do many forward passes with different dropout masks, average all predictions

At test time....

Can in fact do this with a single forward pass! (approximately)

Leave all input neurons turned on (no dropout).

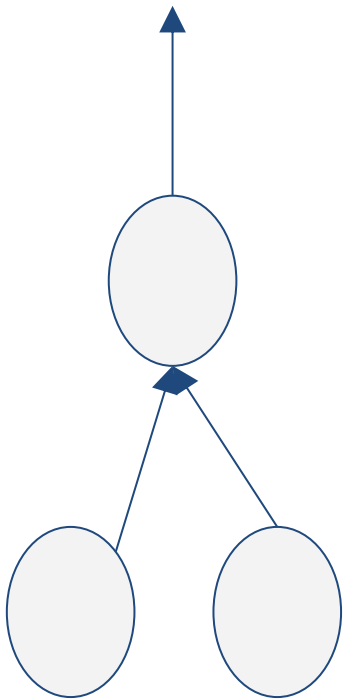


(this can be shown to be an approximation to evaluating the whole ensemble)

At test time....

Can in fact do this with a single forward pass! (approximately)

Leave all input neurons turned on (no dropout).



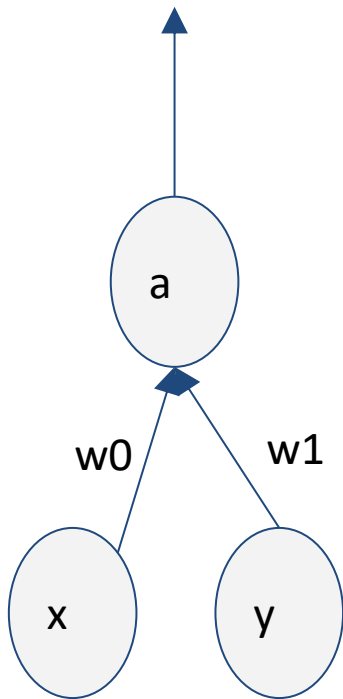
Q: Suppose that with all inputs present at test time the activation of this neuron is x .

What would its activation be during training time, in expectation? (e.g. if $p = 0.5$)

At test time....

Can in fact do this with a single forward pass! (approximately)

Leave all input neurons turned on (no dropout).



during test: $\mathbf{a} = \mathbf{w0} * \mathbf{x} + \mathbf{w1} * \mathbf{y}$

during train:

$$\begin{aligned} \mathbf{E}[\mathbf{a}] &= \frac{1}{4} * (\mathbf{w0} * 0 + \mathbf{w1} * 0 + \\ &\quad \mathbf{w0} * 0 + \mathbf{w1} * \mathbf{y} + \\ &\quad \mathbf{w0} * \mathbf{x} + \mathbf{w1} * 0 + \\ &\quad \mathbf{w0} * \mathbf{x} + \mathbf{w1} * \mathbf{y}) \\ &= \frac{1}{4} * (2\mathbf{w0} * \mathbf{x} + 2\mathbf{w1} * \mathbf{y}) \\ &= \frac{1}{2} * (\mathbf{w0} * \mathbf{x} + \mathbf{w1} * \mathbf{y}) \end{aligned}$$

=> Have to compensate by scaling
the activations back down by $\frac{1}{2}$

Dropout in Keras

```
model = Sequential()  
model.add(Conv2D(32, kernel_size=(3, 3),  
                 activation='relu',  
                 input_shape=input_shape))  
model.add(Conv2D(64, (3, 3), activation='relu'))  
model.add(MaxPooling2D(pool_size=(2, 2)))  
model.add(Dropout(0.25))  
  
model.add(Flatten())  
model.add(Dense(128, activation='relu'))  
model.add(Dropout(0.5))  
model.add(Dense(num_classes, activation='softmax'))
```


Dropout

- Dropout can be seen as ensemble averaging
- Each model in the ensemble is smaller than the original
- Reduces overfitting
- Introduces train and test mode
- Common settings are:
 - 0.2 dropout on input layer
 - 0.5 dropout on hidden layers

Good Old L1/L2 Regularization

- Keras layer parameters

```
keras.layers.core.Dense(units, activation=None,  
kernel_regularizer=None, bias_regularizer=None,  
activity_regularizer=None, kernel_constraint=None,  
bias_constraint=None)
```

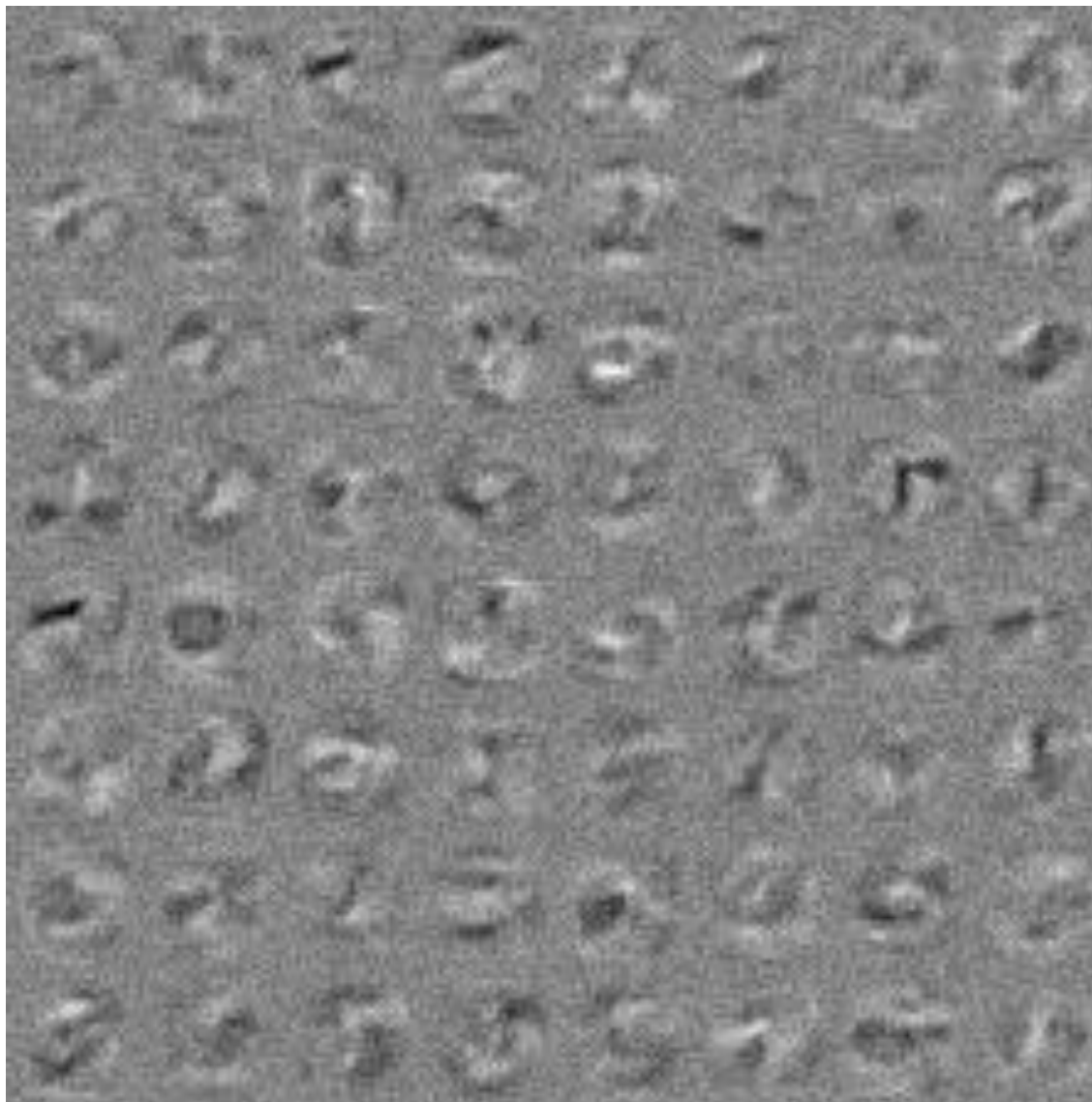
- **activation**: Default is linear!
- **regularizer**: This is our standard L1 or L2 option
- Example: <https://keras.io/regularizers/>

```
Dense(10, activation='relu',  
      kernel_regularizer=regularizers.l2(0.01))
```

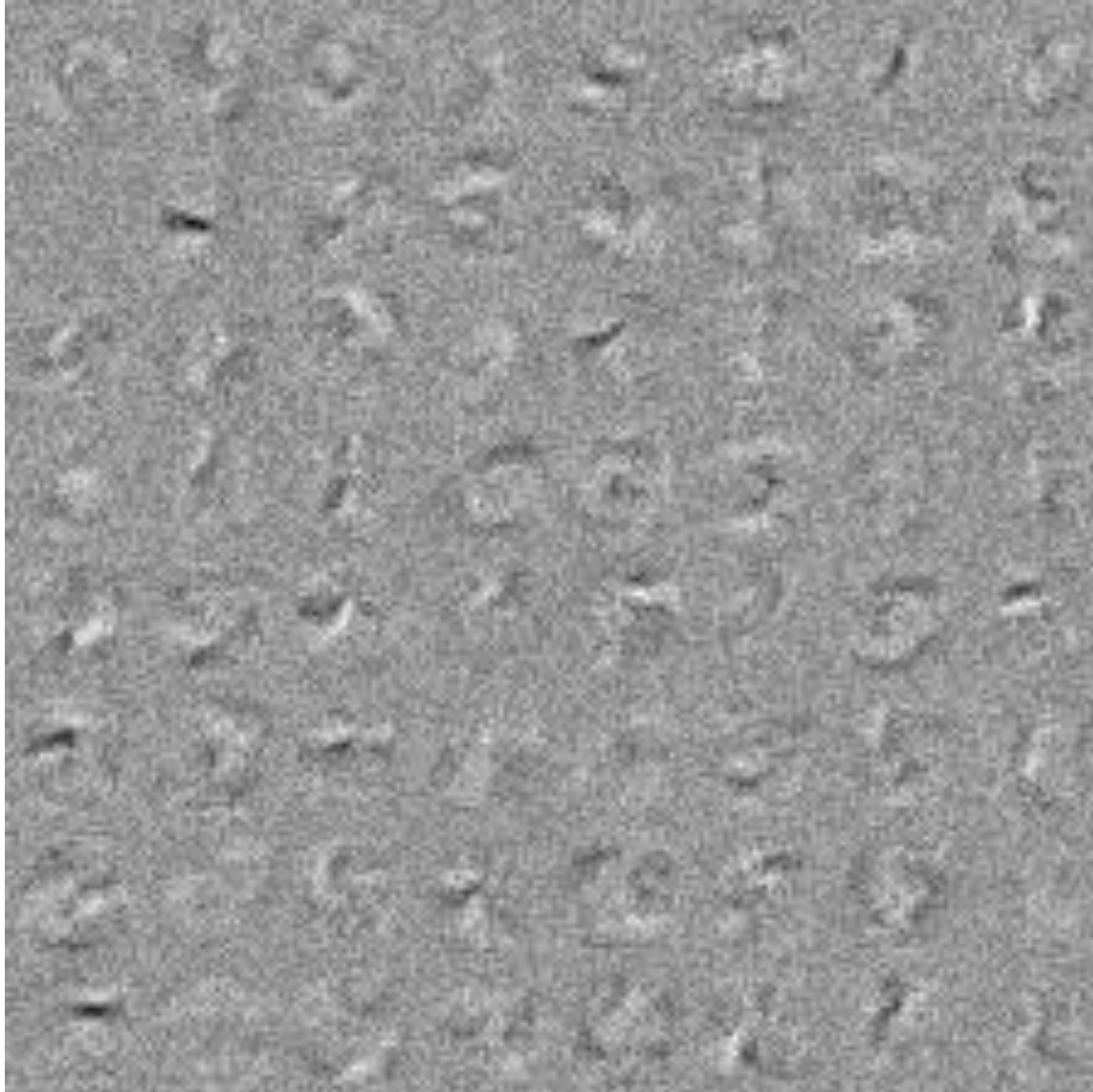
Filter visualizations

- original simple MLP network
- with regularization
- with dropout
- snapshots at 1, 10, 100, 200 epochs

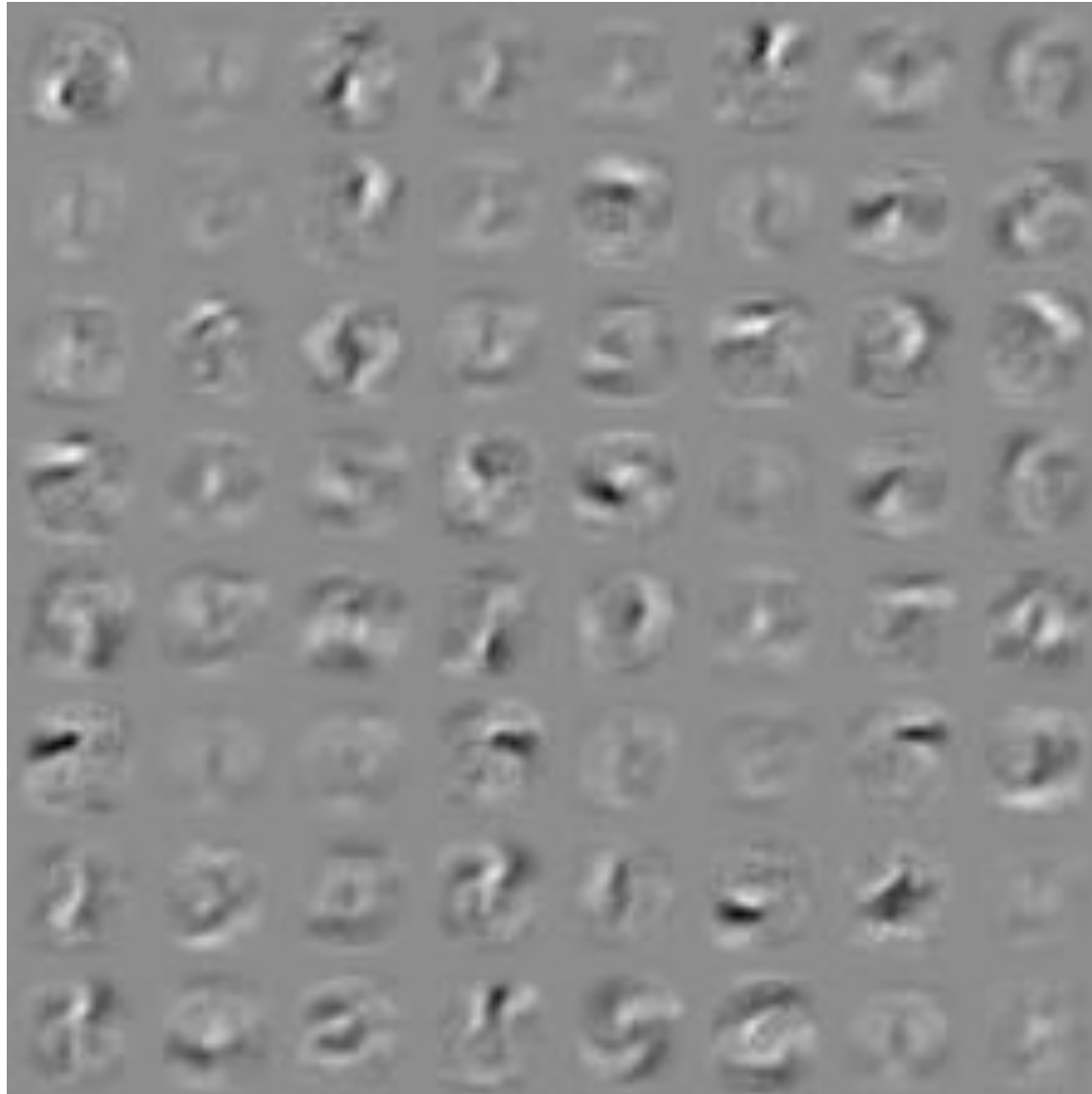
MLP with Sigmoid



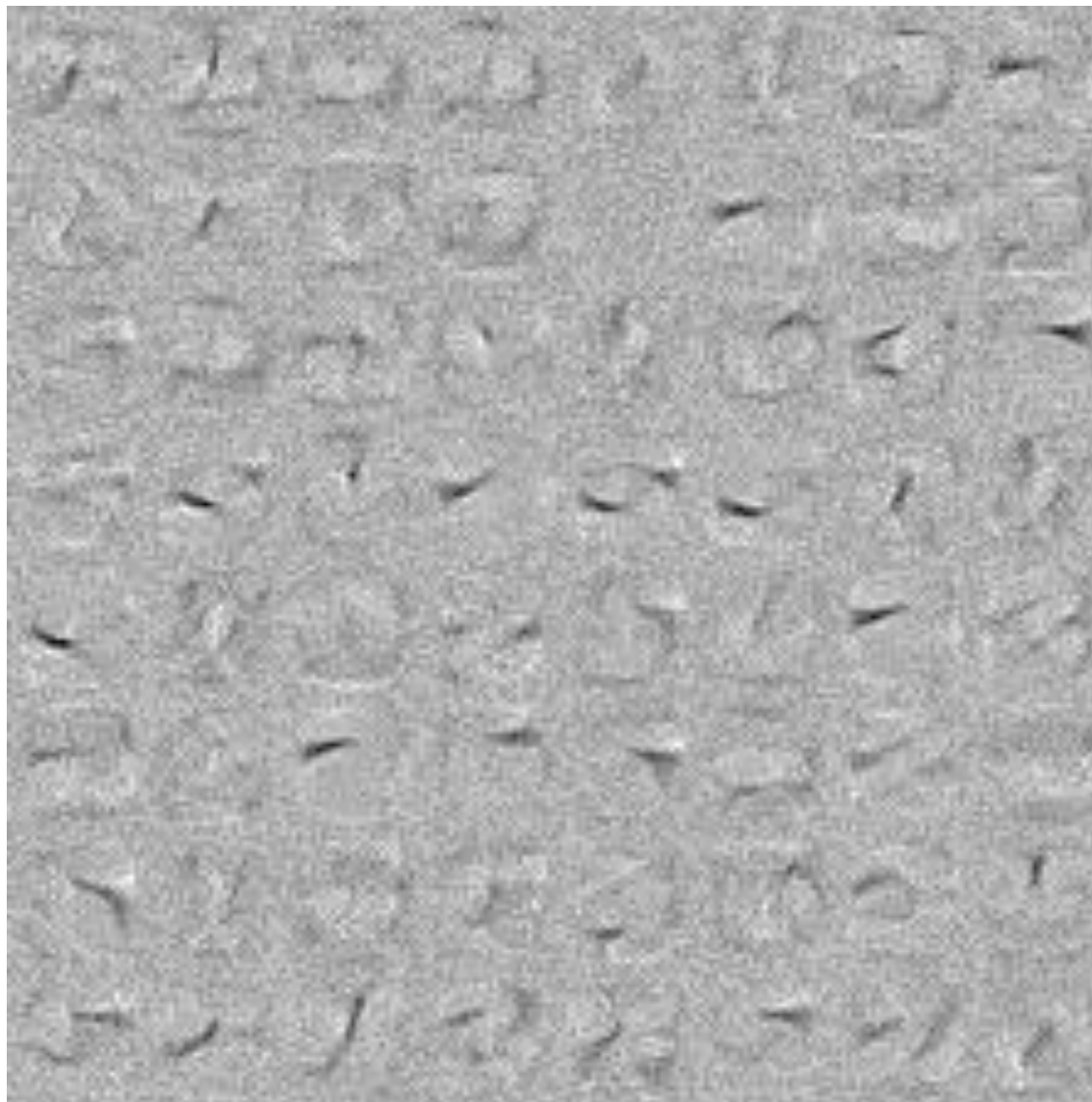
MLP with ReLU



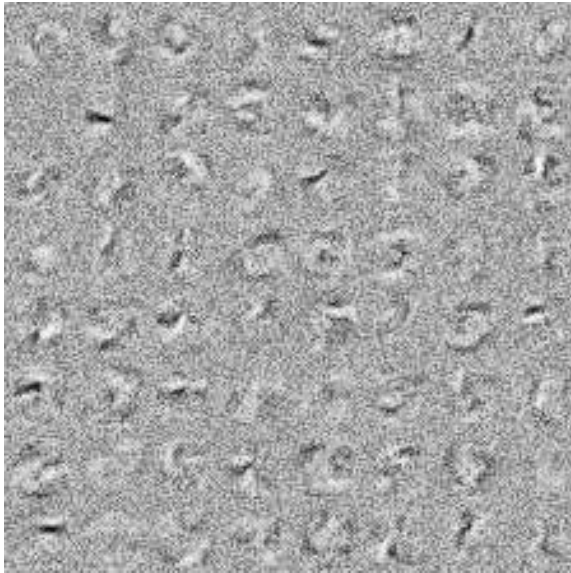
MLP ReLU + L2



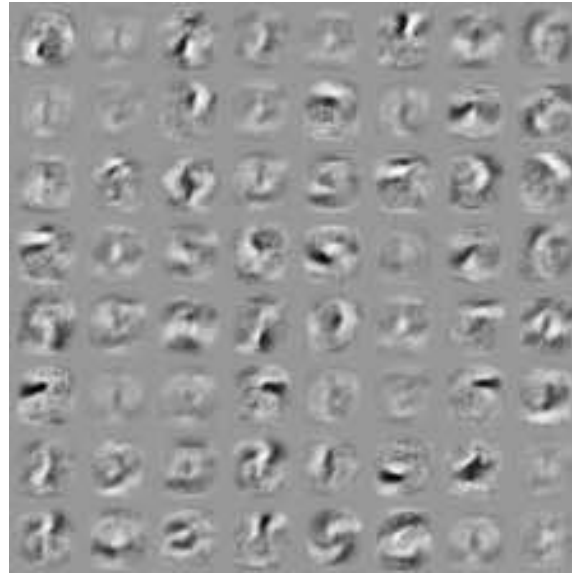
MLP ReLu + Dropout



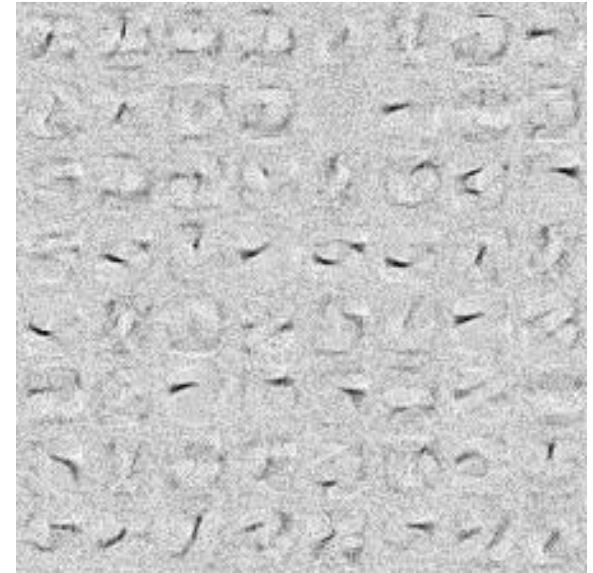
Overview



ReLU



ReLU + L2



ReLU + Dropout

Try at home: ReLU + L2 + Dropout

What About Accuracy?

200 epochs, lr = 0.001, batch_size = 100

MLP_ReLU

Test loss: 0.0887242713874

Test accuracy: 0.9757

MLP_ReLU + L2

Test loss: 0.266304204607

Test accuracy: 0.968

MLP_ReLU + Dropout

Test loss: 0.135238234081

Test accuracy: 0.9638

This is not a fair
comparison!

It is typical for
dropout to
increase
training time

Milestone!

- What we have so far:
 - CNN layer
 - MaxPool layer
 - FC layer
 - Activations:
 - **ReLU**
 - sigmoid or softmax for last layer
 - Cost function – cross entropy
 - **Training – SGD + Momentum**
 - Regularization – dropout + L2/L1
 - Convergence check

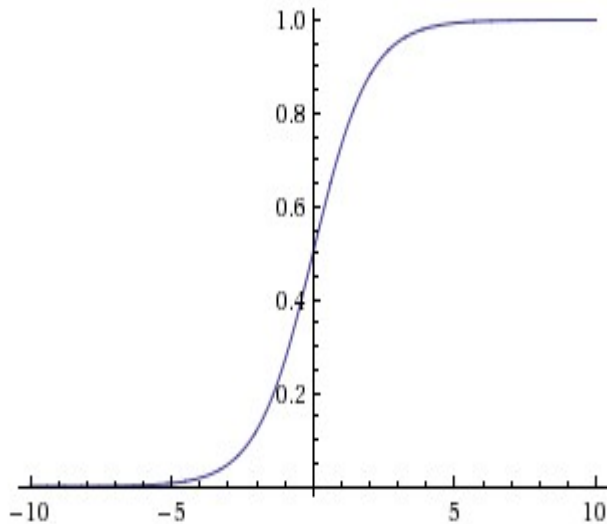
A few more notes on ReLU and Gradients

- One of the tricks that made deep learning work
- Need to be a bit careful
- We train our network with gradient information
- We take the gradient of our loss function
- And find the parameters that correspond to a good local minimum

Activation Functions

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

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



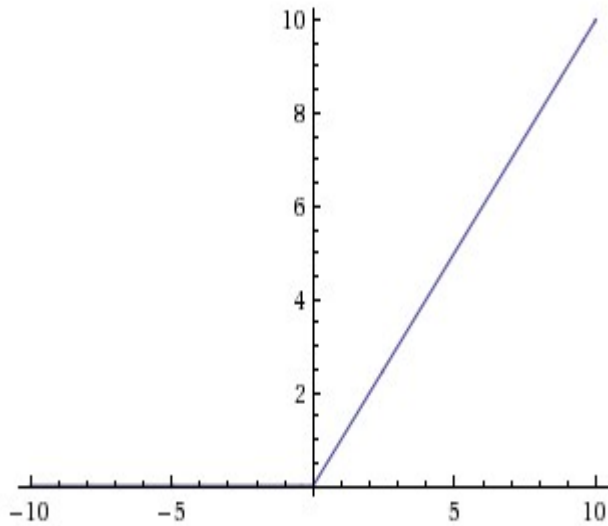
Sigmoid

problems:

1. Saturated neurons “kill” the gradients
2. $\exp()$ is a bit compute expensive

$$\frac{\partial L}{\partial x} = \frac{\partial \sigma}{\partial x} \frac{\partial L}{\partial \sigma}$$

Activation Functions

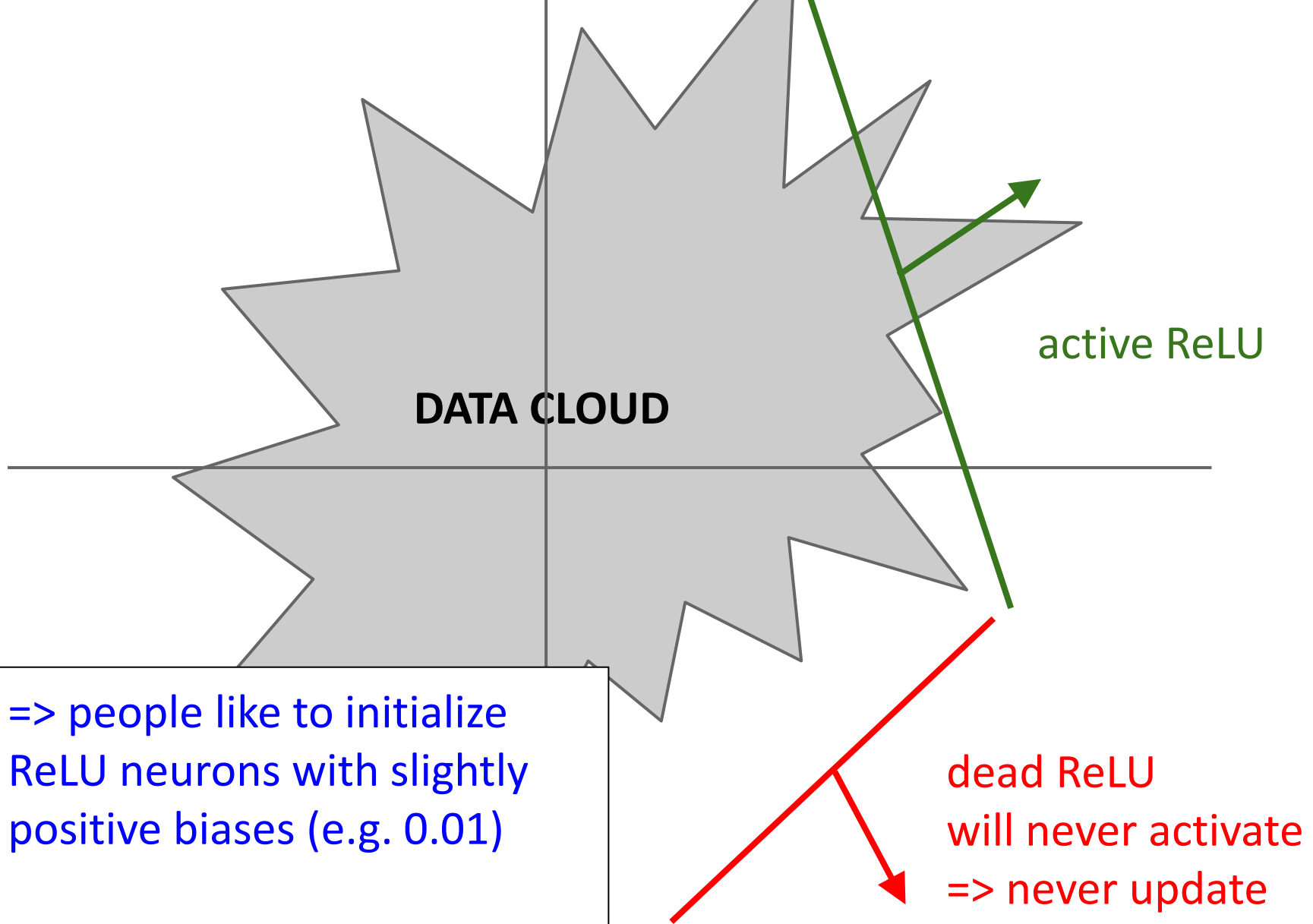


ReLU

(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)
- Not zero-centered output
- An annoyance:

hint: what is the gradient when $x < 0$?



ReLU in Practice

- ReLU is a good default to try
- Be careful with your learning rate
- Might benefit from small positive bias initialization

Weight Initialization

- Typically random if you train from scratch
- If weights are too large the network has a hard time learning from updates
- If weights are too small they might not break symmetry enough.
- This is actually an area of ongoing research.

Proper initialization is an active area of research...

Understanding the difficulty of training deep feedforward neural networks

by Glorot and Bengio, 2010

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

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

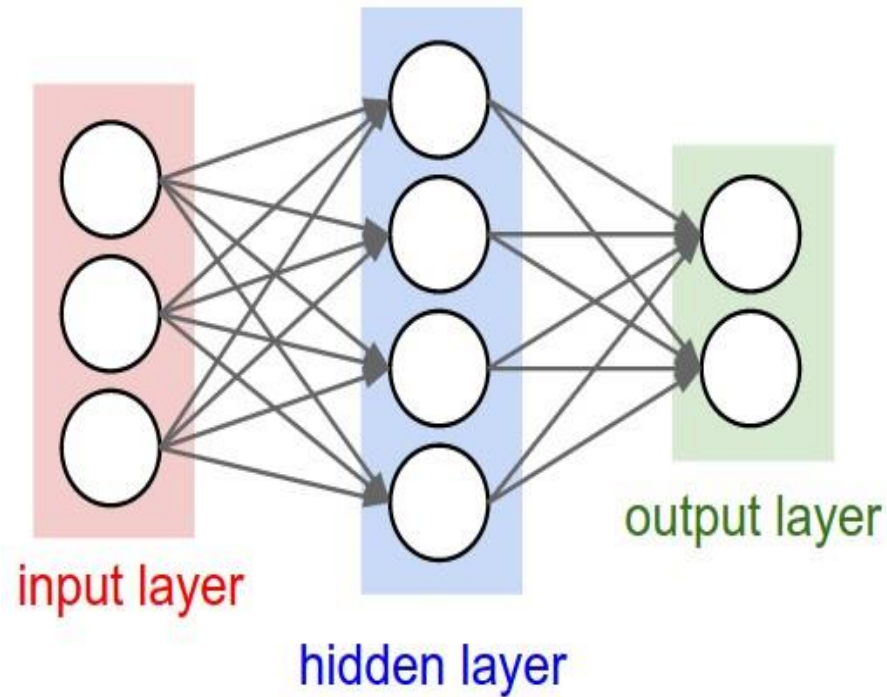
Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification by He et al., 2015

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

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

...

- Q: what happens when $W=0$ init is used?



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

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

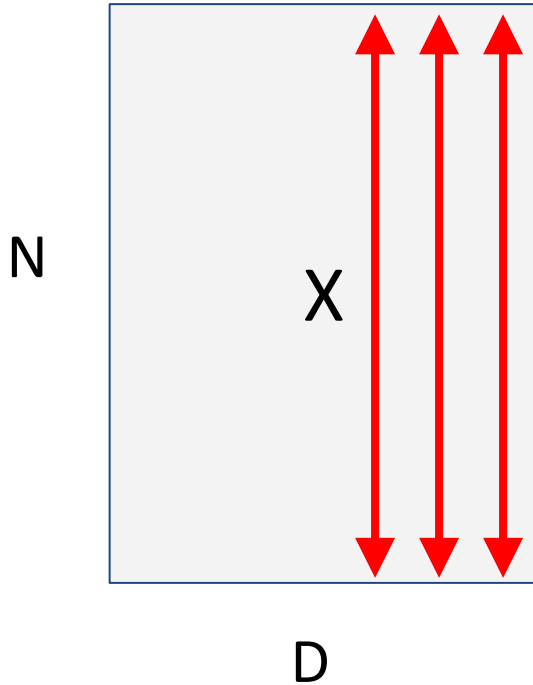
Works ~okay for small networks, but can lead to non-homogeneous distributions of activations across the layers of a network.

Weight Initialization Confusion

- glorot (also called Xavier)(for tanh)
 - $\text{stddev} = \sqrt{2 / (\text{fan_in} + \text{fan_out})}$ [Keras]
 - OR
 - $\text{stddev} = \sqrt{1 / \text{fan_in}}$ [Caffe]
- he (for ReLU)
 - $\text{stddev} = \sqrt{2 / \text{fan_in}}$ [Keras]

The Rescue: BatchNorm Layer

“you want unit gaussian activations?
just make them so.”



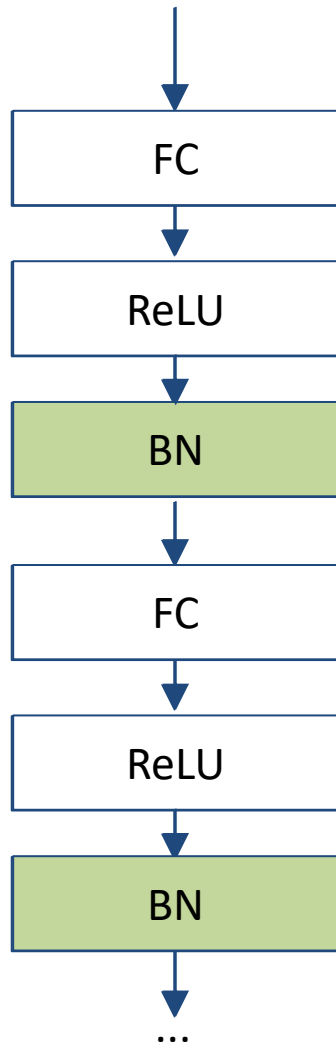
1. compute the empirical mean and
variance independently for each
dimension.

2. Normalize

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

Batch Normalization

[Ioffe and Szegedy, 2015]



Usually inserted after fully connected or convolutional layers, and before nonlinearity....
Or After....

- 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

Better Optimization

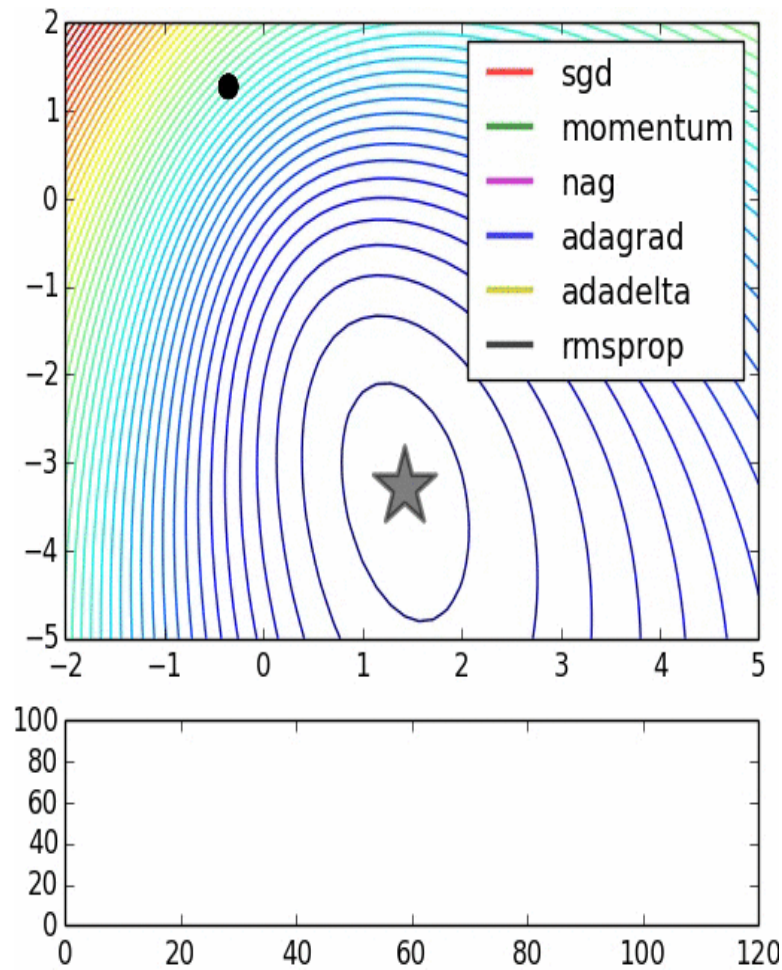


Image credits: Alec Radford

<http://cs231n.github.io/>

AdaGrad update

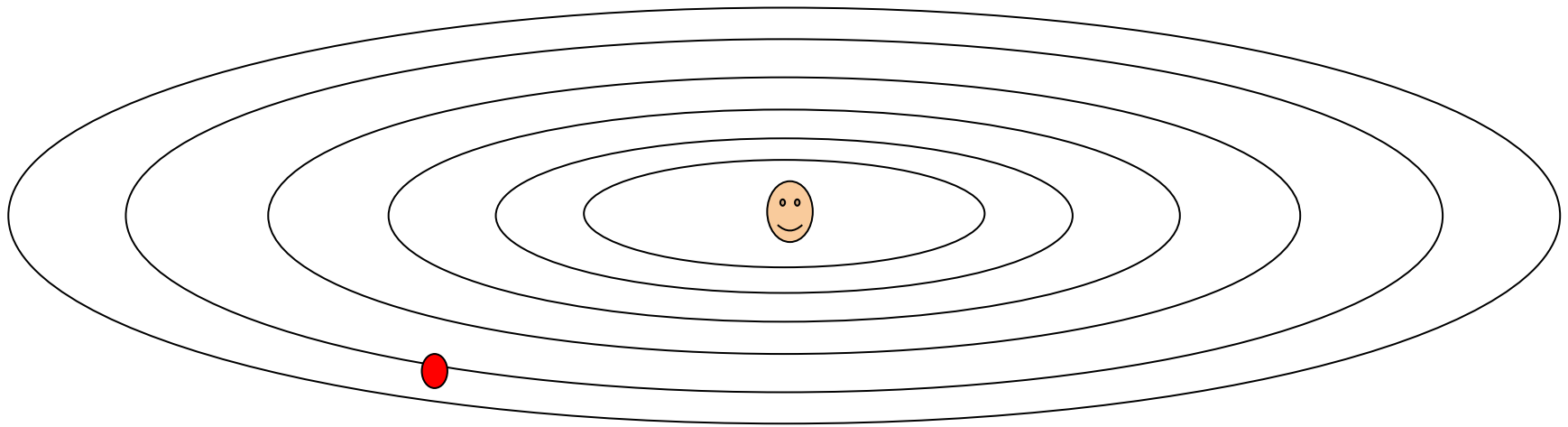
[Duchi et al., 2011]

```
# Adagrad update  
cache += dx**2  
x += - learning_rate * dx / (np.sqrt(cache) + 1e-7)
```


AdaGrad update

```
# Adagrad update  
cache += dx**2  
x += - learning_rate * dx / (np.sqrt(cache) + 1e-7)
```

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




Over time step size gets very small – training comes to a halt.

RMSProp update

Decay the cache over time

```
# Adagrad update  
cache += dx**2  
x += - learning_rate * dx / (np.sqrt(cache) + 1e-7)
```



```
# RMSProp  
cache = decay_rate * cache + (1 - decay_rate) * dx**2  
x += - learning_rate * dx / (np.sqrt(cache) + 1e-7)
```

Adam update

(incomplete, but close)

```
# Adam
m = beta1*m + (1-beta1)*dx # update first moment
v = beta2*v + (1-beta2)*(dx**2) # update second moment
x += - learning_rate * m / (np.sqrt(v) + 1e-7)
```

momentum

RMSProp-like

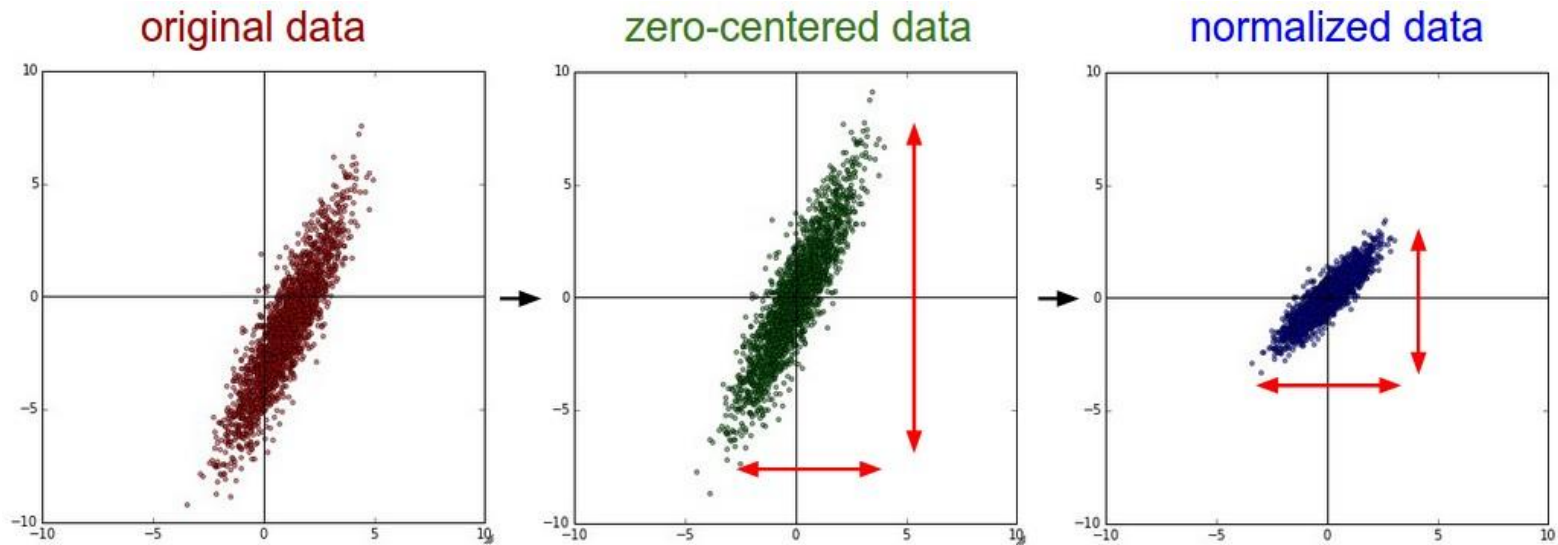
Looks a bit like RMSProp with momentum

```
# RMSProp
cache = decay_rate * cache + (1 - decay_rate) * dx**2
x += - learning_rate * dx / (np.sqrt(cache) + 1e-7)
```

Milestone Update

- What we have so far:
 - CNN layer
 - MaxPool layer
 - FC layer
 - **BatchNorm Layer**
 - Activations:
 - **ReLU**
 - sigmoid or softmax for last layer
 - Cost function – cross entropy
 - **Initialization**
 - **Training – Adam**
 - Regularization – dropout + L2/L1
 - Convergence check

Preprocess the data

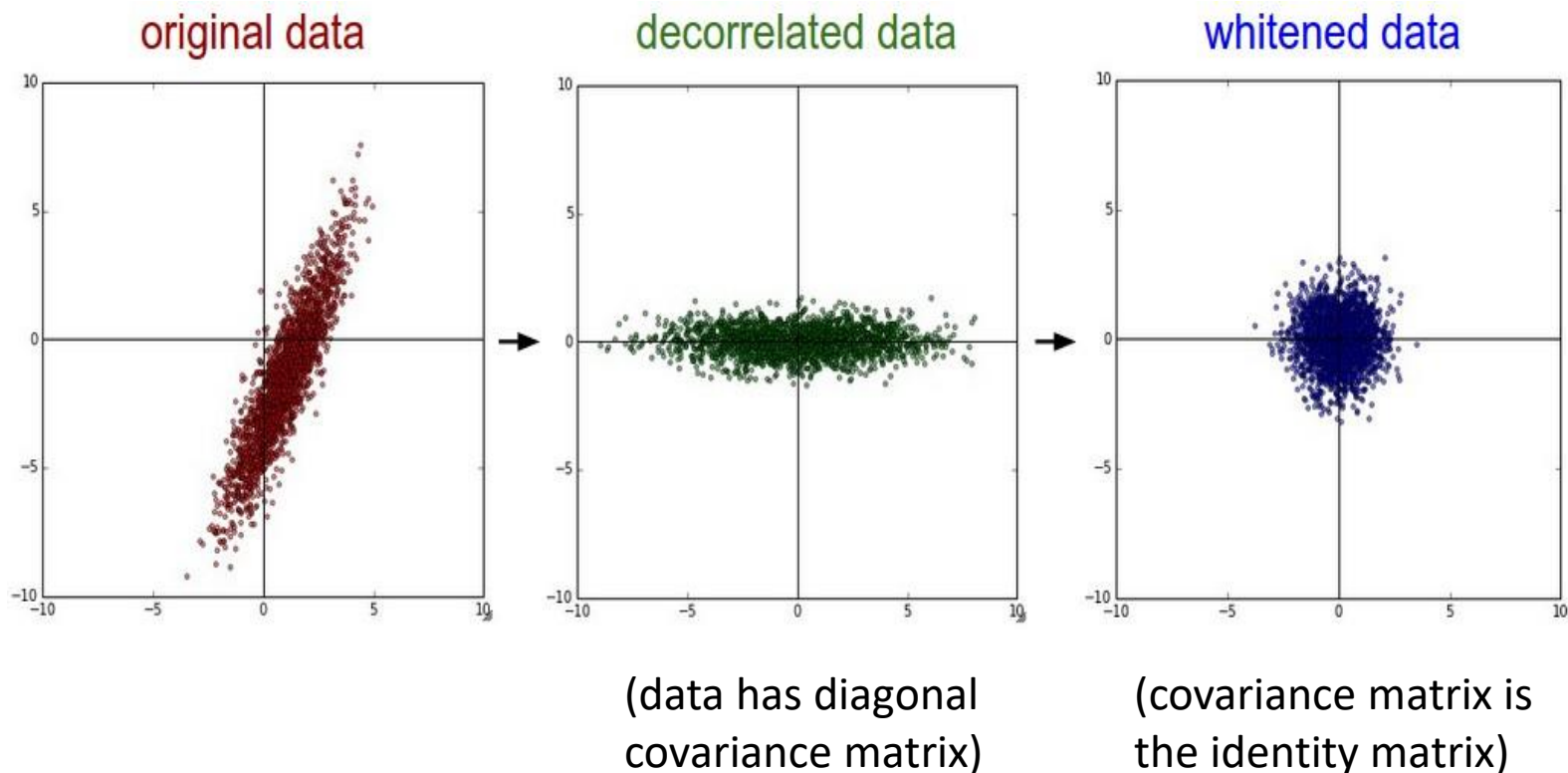


```
X -= np.mean(X, axis = 0)
```

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

Preprocess the data

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



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

Data Preprocessing in Keras

```
keras.preprocessing.image.ImageDataGenerator(  
    featurewise_center=False,  
    samplewise_center=False,  
    featurewise_std_normalization=False,  
    samplewise_std_normalization=False,  
    zca_whitening=False,  
    ...  
)
```

<https://keras.io/preprocessing/image/>

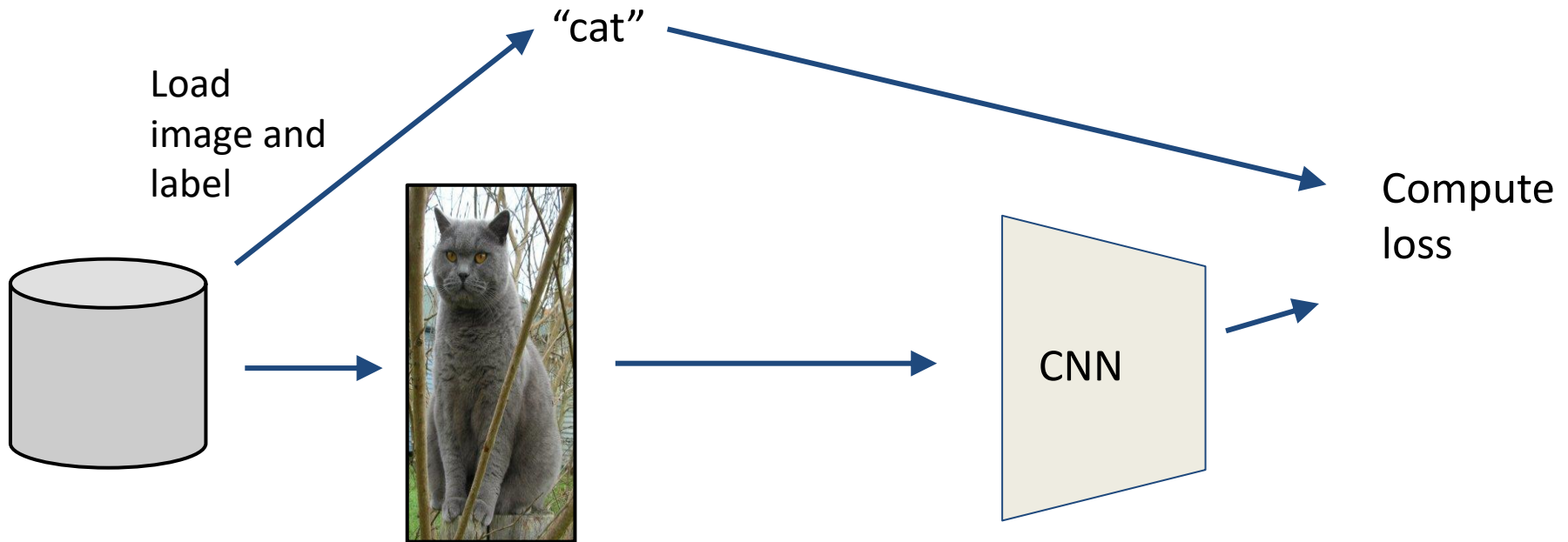
Remember

- You need to handle the test data in the same way!
- Save normalization parameters
- Apply them to test data
- In scikit learn terminology:

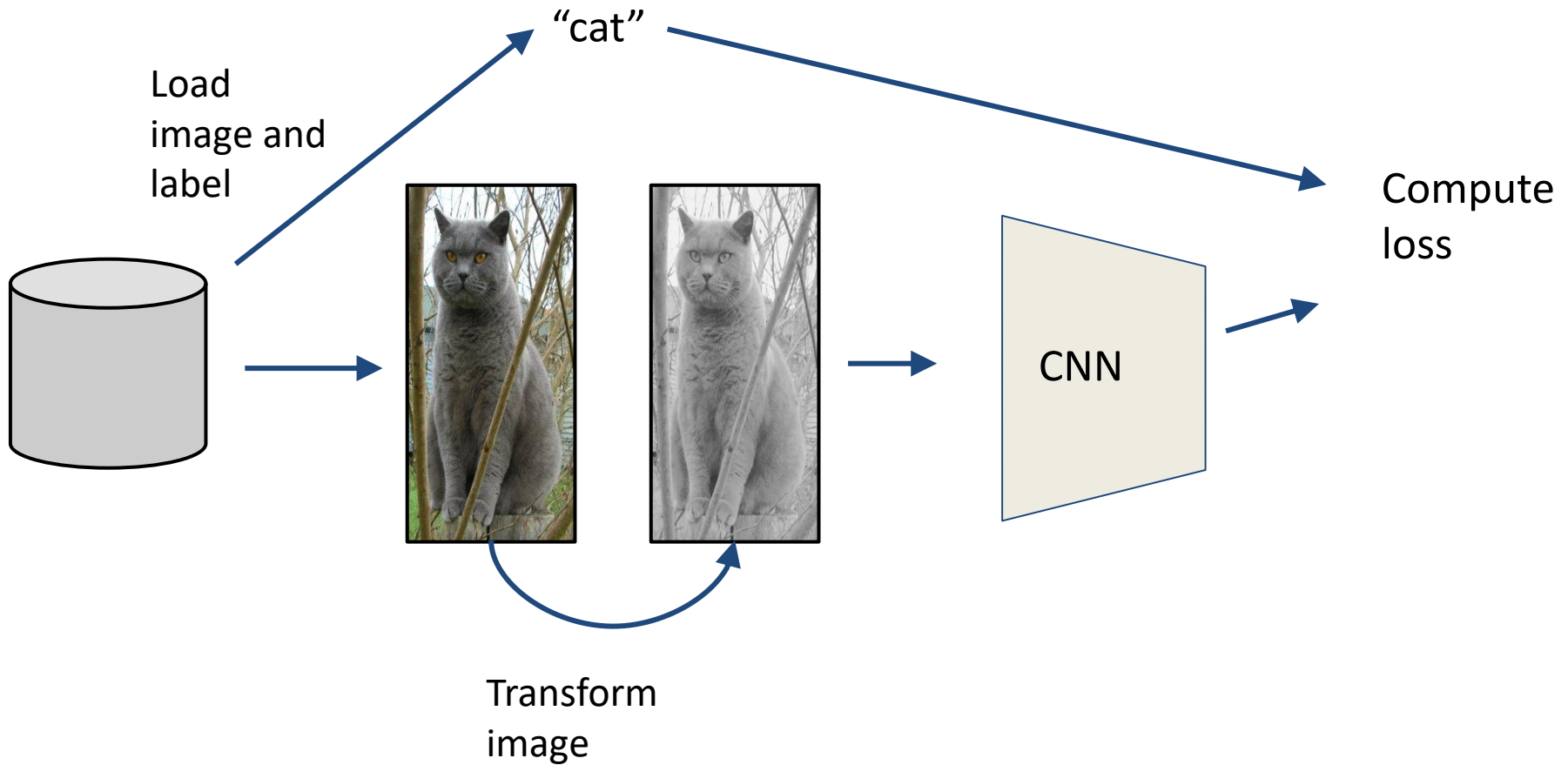
```
normalizer.fit_transform(train_data)  
normalizer.transform(test_data)
```

Data Augmentation

Data Augmentation



Data Augmentation



Data Augmentation

- Change the pixels without changing the label
- Train on transformed data
- VERY widely used

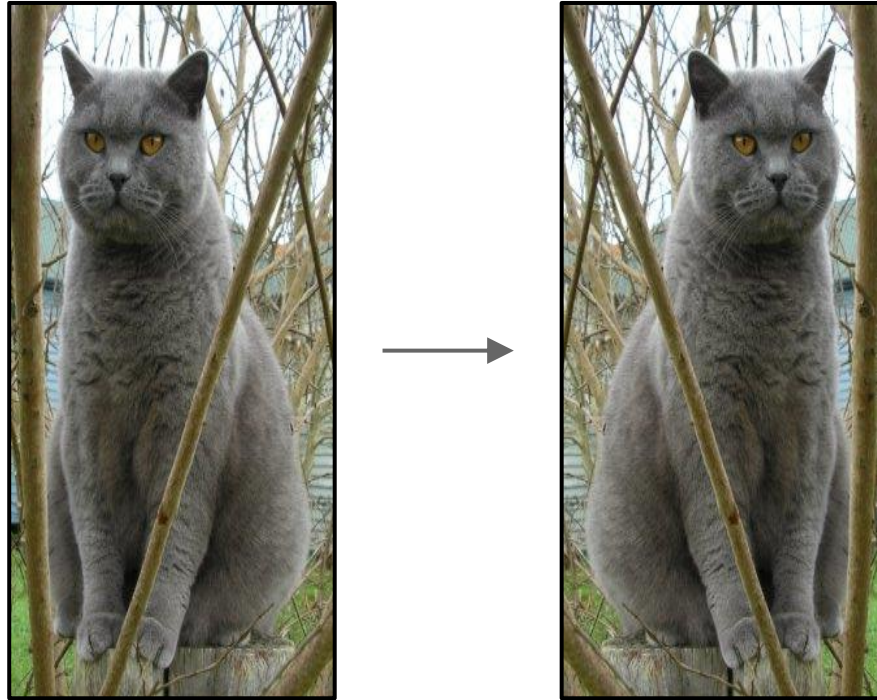


08	02	22	97	38	15	00	40	00	75	04	05	07	78	52	12	50	77	81	89
49	49	99	40	17	81	18	57	60	87	17	40	98	43	63	43	04	56	62	00
81	49	31	73	55	79	14	29	93	71	40	67	59	88	30	03	49	13	36	65
52	70	95	23	04	60	11	42	69	24	88	56	01	32	56	71	37	02	36	91
22	31	16	71	51	67	43	89	41	92	36	54	22	40	40	29	66	33	13	80
24	47	32	00	99	03	45	02	44	75	33	53	78	36	84	20	35	17	12	50
32	98	81	28	64	23	67	10	26	38	40	67	59	54	70	66	18	38	64	70
67	26	20	68	02	62	12	20	95	63	94	39	63	08	40	91	66	49	94	21
24	55	58	05	66	73	99	26	97	17	78	78	96	83	14	88	34	89	63	72
21	36	23	09	75	00	76	44	20	45	35	14	00	61	33	97	34	31	33	95
78	17	53	28	22	75	31	67	15	94	03	80	04	62	16	14	09	53	56	92
16	39	05	42	96	35	31	47	55	58	88	24	00	17	54	24	36	29	85	57
86	56	00	48	35	71	89	07	05	44	44	37	44	60	21	58	51	54	17	58
19	80	81	68	05	94	47	69	28	73	92	13	86	52	17	77	04	89	55	40
04	52	08	83	97	35	99	16	07	97	57	32	16	26	26	79	33	27	98	66
88	94	68	87	57	62	20	72	03	46	33	67	46	55	12	32	63	93	53	69
04	42	16	73	39	25	39	11	24	94	72	18	08	46	29	32	40	62	76	36
20	69	36	41	72	30	23	88	35	02	38	69	82	67	59	85	74	04	36	16
20	73	35	29	78	31	90	01	74	31	49	71	48	86	01	16	23	57	05	54
01	70	54	71	83	51	54	69	16	92	33	48	61	43	52	01	89	19	62	48

What the computer
sees

Data Augmentation

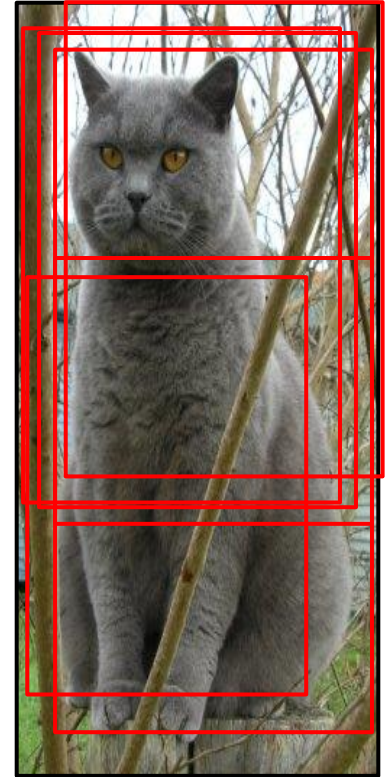
1. Horizontal flips



Data Augmentation

2. Random crops/scales

Training: sample random crops / scales



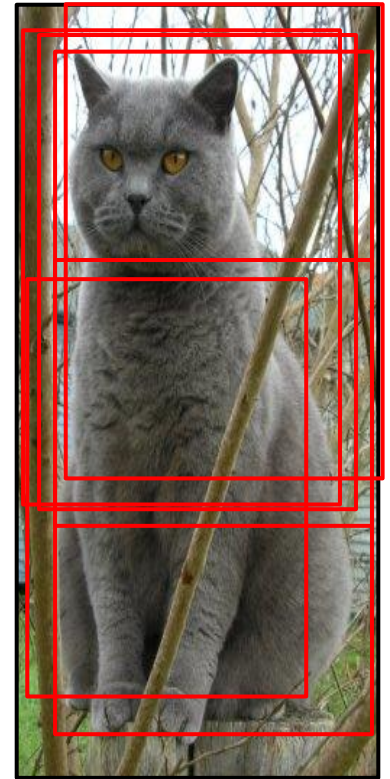
Data Augmentation

2. Random crops/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 224×224 patch



Data Augmentation

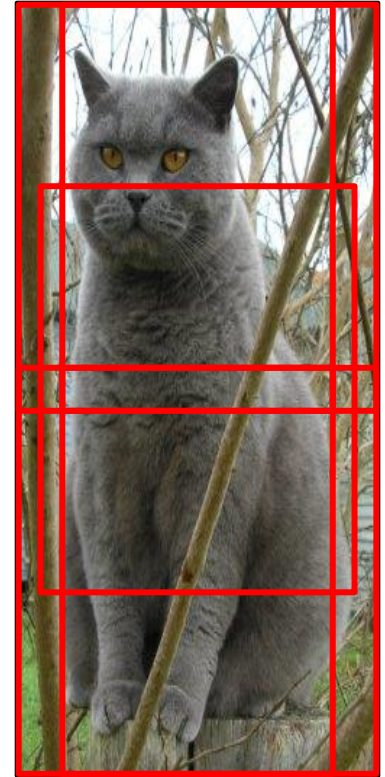
2. Random crops/scales

Training: sample random crops / scales

ResNet:

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

Testing: average a fixed set of crops



Data Augmentation

2. Random crops/scales

Training: sample random crops / scales

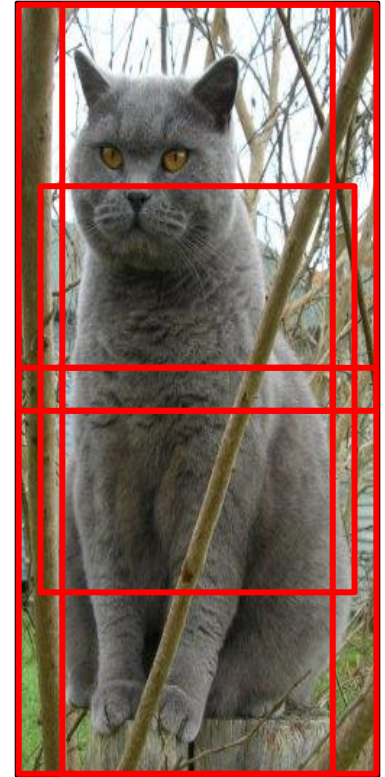
ResNet:

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

Testing: average a fixed set of crops

ResNet:

1. Resize image at 5 scales: $\{224, 256, 384, 480, 640\}$
2. For each size, use 10 224×224 crops: 4 corners + center, + flips

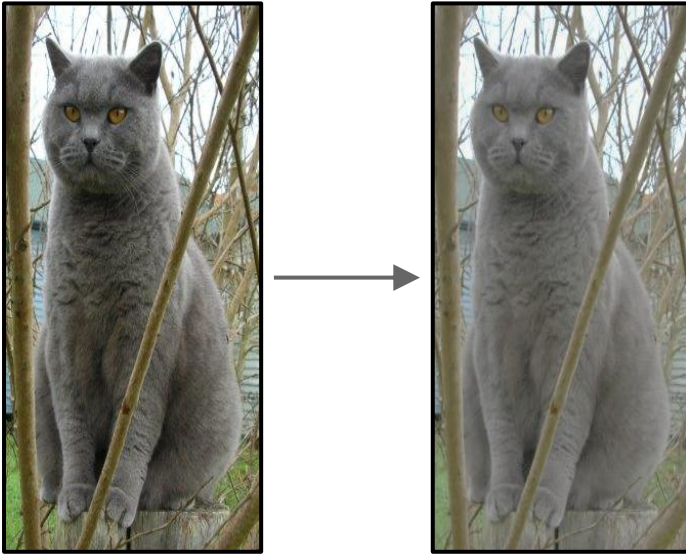


Data Augmentation

3. Color jitter

Simple:

Randomly jitter contrast

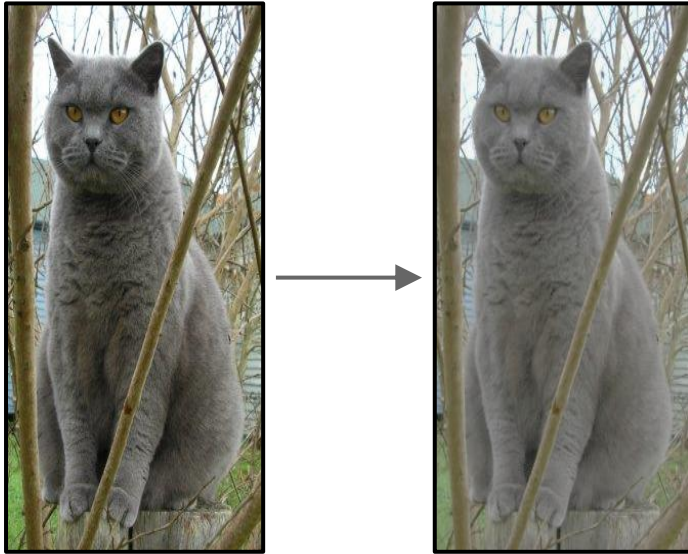


Data Augmentation

3. Color jitter

Simple:

Randomly jitter contrast



Complex:

1. Apply PCA to all [R, G, B] pixels in training set
2. Sample a “color offset” along principal component directions
1. Add offset to all pixels of a training image

(As seen in [Krizhevsky et al. 2012], ResNet, etc)

Data Augmentation

4. Get creative!

Random mix/combinations of :

- translation
- rotation
- stretching
- shearing,
- lens distortions, ... (go crazy)

Data Augmentation: Takeaway

- Simple to implement, use it
- Especially useful for small datasets
- Fits into framework of noise / marginalization

Data Augmentation in Keras

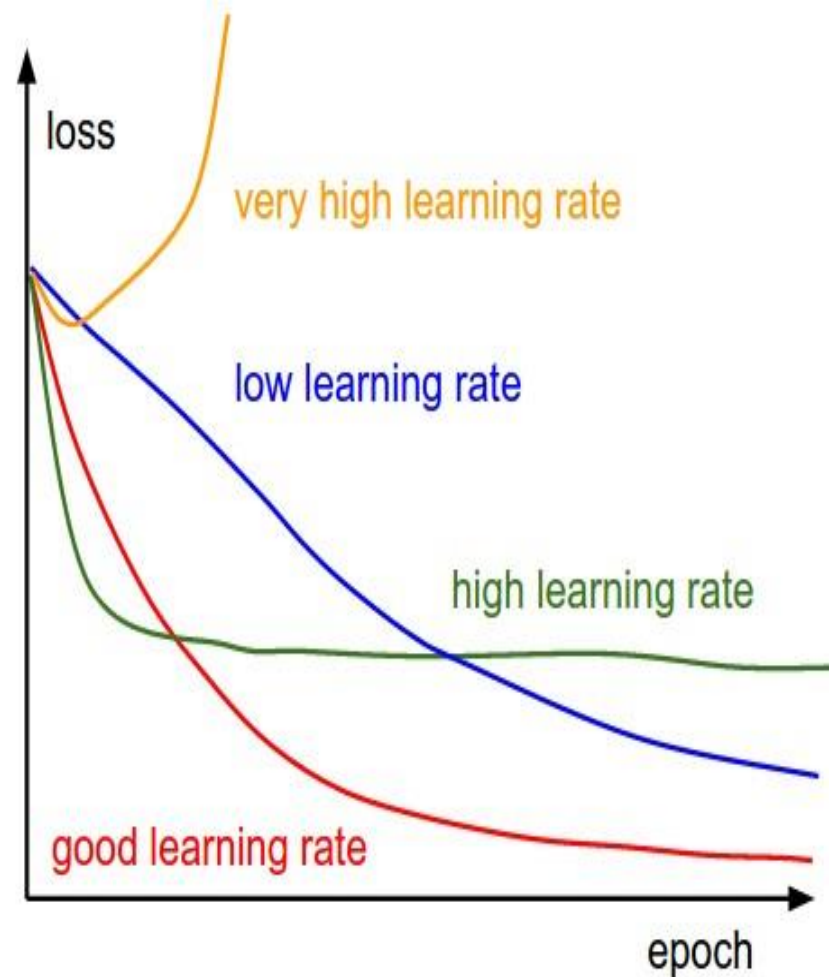
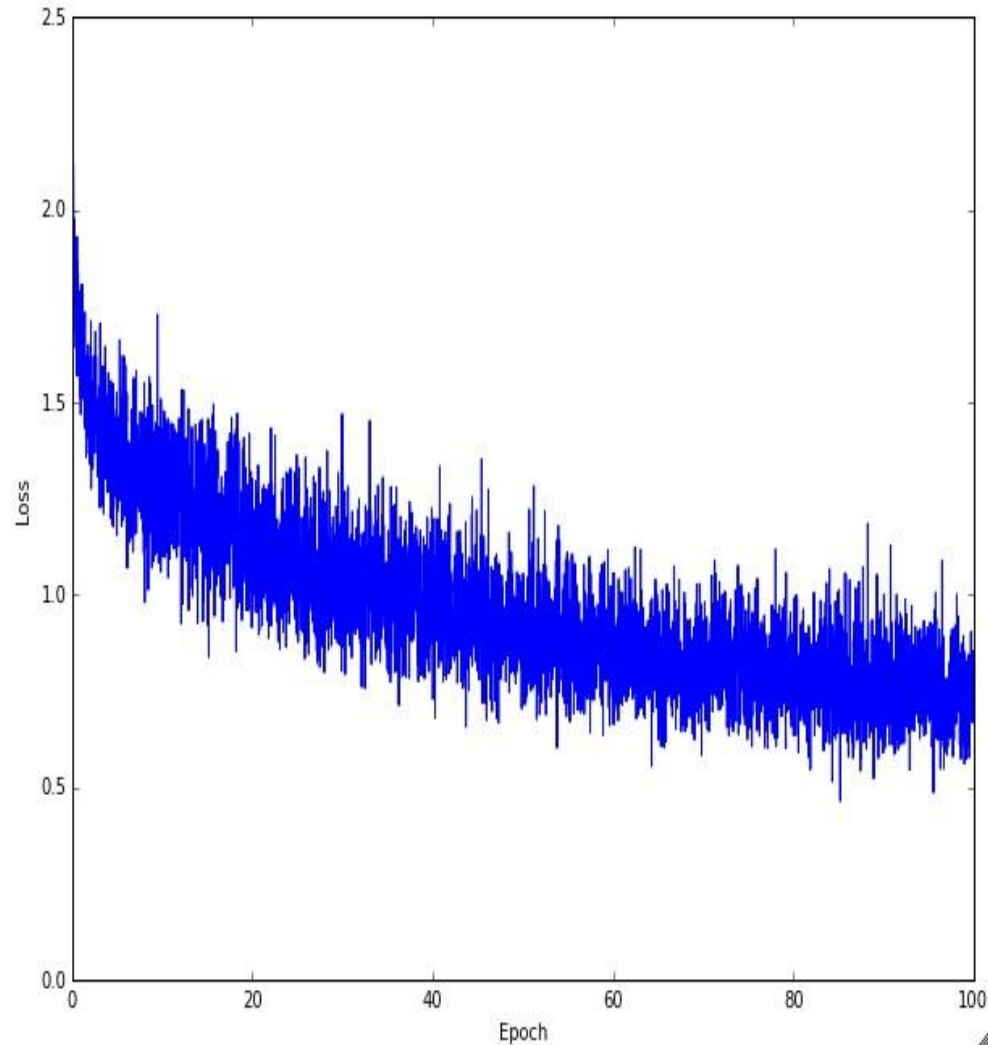
```
keras.preprocessing.image.ImageDataGenerator(  
    ...  
    rotation_range=0., width_shift_range=0.,  
    height_shift_range=0., shear_range=0.,  
    zoom_range=0., channel_shift_range=0.,  
    fill_mode='nearest', cval=0.,  
    horizontal_flip=False, vertical_flip=False,  
    rescale=None, preprocessing_function=None,  
    data_format=K.image_data_format())  
)
```

<https://keras.io/preprocessing/image/>

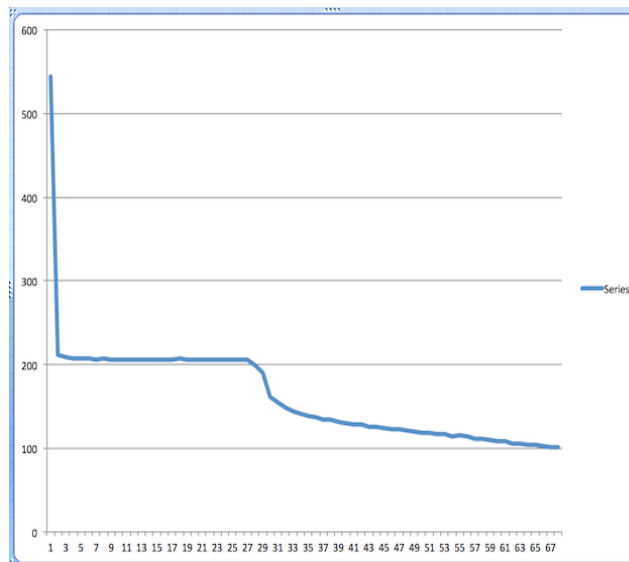
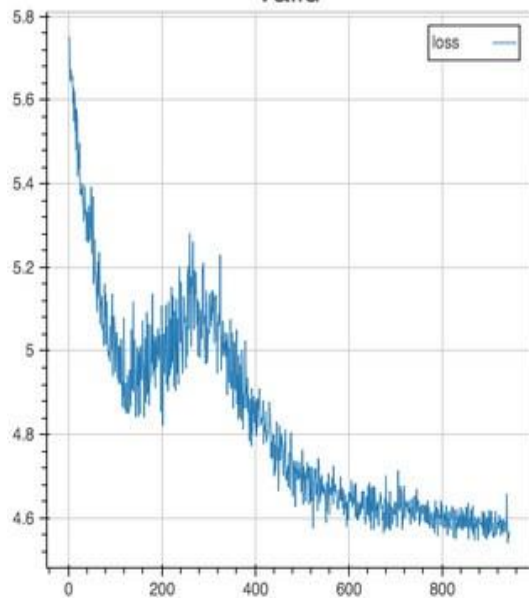
Learning Rate (once again)

- Need to find good range for learning rate
 - Look at the loss during training
 - If it is NaN => learning rate too high
 - If it decreases too slow => learning rate too low
- Also remember batch size and gradient estimation influence
 - Large batch – large learning rate – fewer updates
 - Small batch – smaller learning rate – more updates

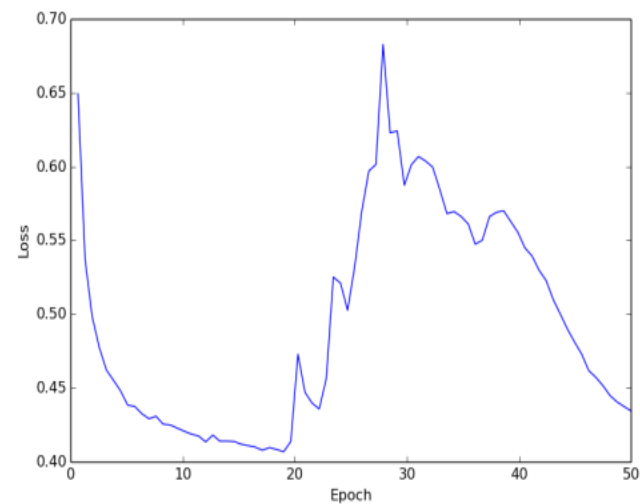
Monitor and visualize the loss curve



valid



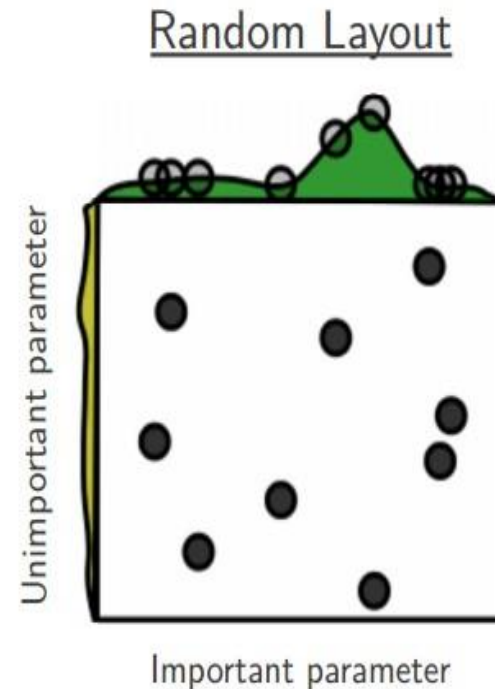
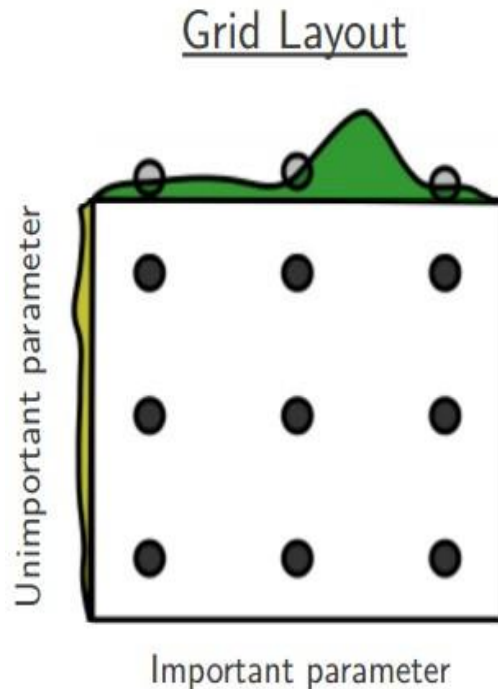
Training Loss



Cross Validation for Parameter Tuning

- How to fine tune the learning rate and regularization?
- Just like any classifier you can also use CV for deep learning
- Costs for trying a parameter setting are high
- People tend to use coarse grids
- This can be bad

Random Search vs. Grid Search



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

Architecture Design

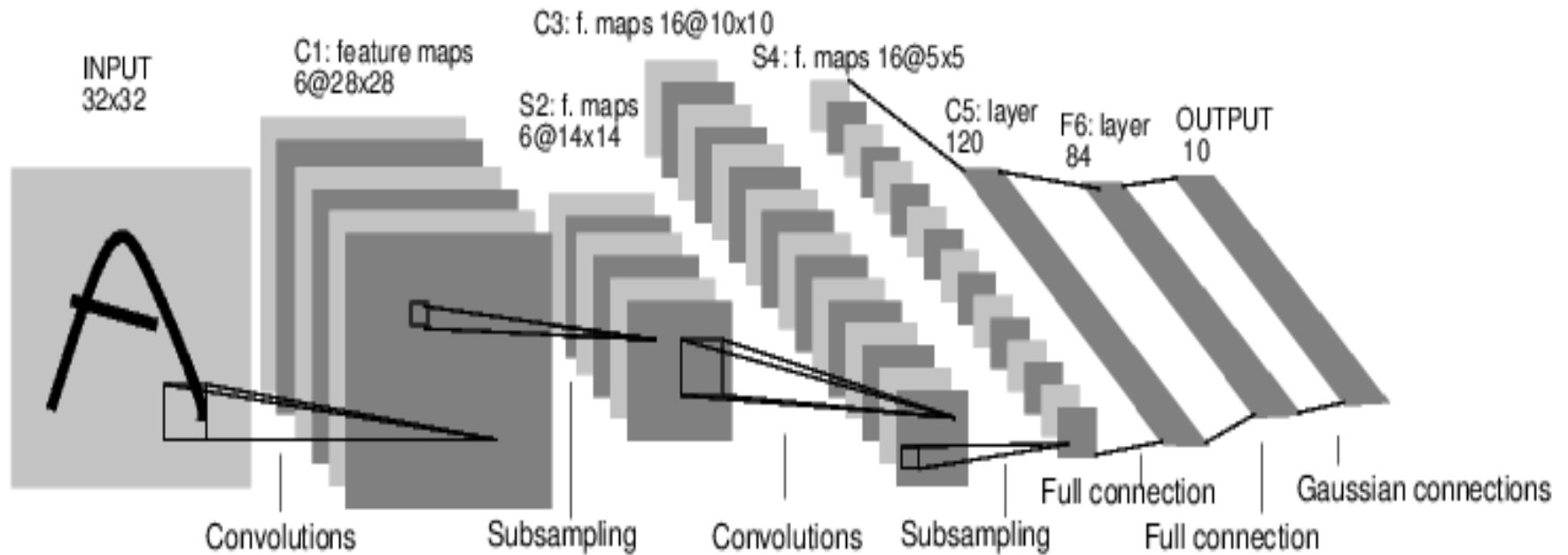
- Things to consider:
 - Previous work done on similar problems
 - Computational resources
- Wider, and deeper means more degrees of freedom
 - Better classification performance
 - IF you can tune it right

Case studies

- Don't go from scratch when you learn deep learning
- Look at example models
- Use tricks that have been shown to work

Case Study: LeNet-5

[LeCun et al., 1998]



Conv filters were 5x5, applied at stride 1

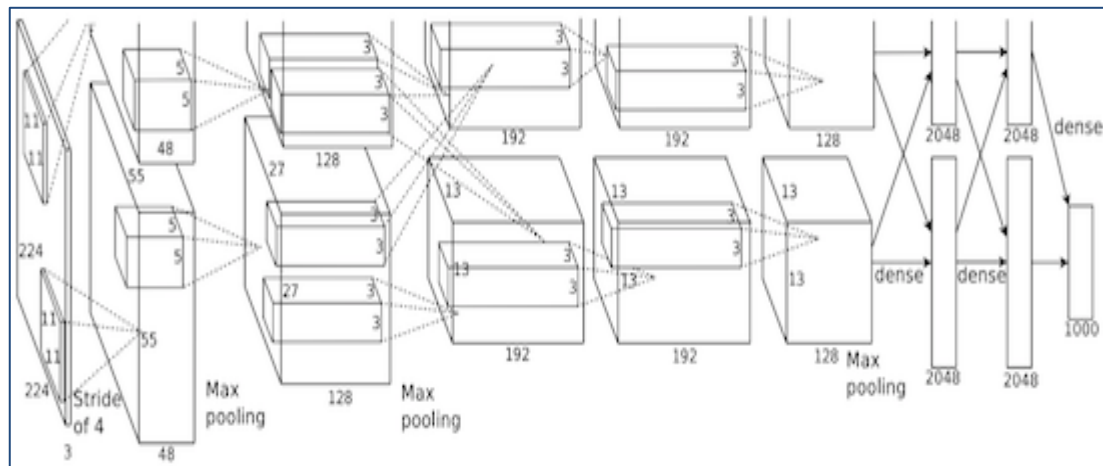
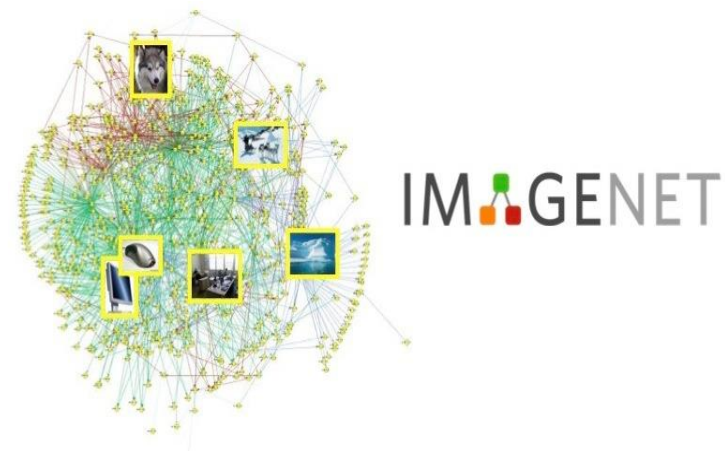
Subsampling (Pooling) layers were 2x2 applied at stride 2

i.e. architecture is [CONV-POOL-CONV-POOL-CONV-FC]

A bit of history:

ImageNet Classification with Deep Convolutional Neural Networks

[Krizhevsky, Sutskever, Hinton, 2012]



“AlexNet”

Case Study: AlexNet

[Krizhevsky et al. 2012]

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] **CONV1**: 96 11x11 filters at stride 4, pad 0

[27x27x96] **MAX POOL1**: 3x3 filters at stride 2

[27x27x96] **NORM1**: Normalization layer

[27x27x256] **CONV2**: 256 5x5 filters at stride 1, pad 2

[13x13x256] **MAX POOL2**: 3x3 filters at stride 2

[13x13x256] **NORM2**: Normalization layer

[13x13x384] **CONV3**: 384 3x3 filters at stride 1, pad 1

[13x13x384] **CONV4**: 384 3x3 filters at stride 1, pad 1

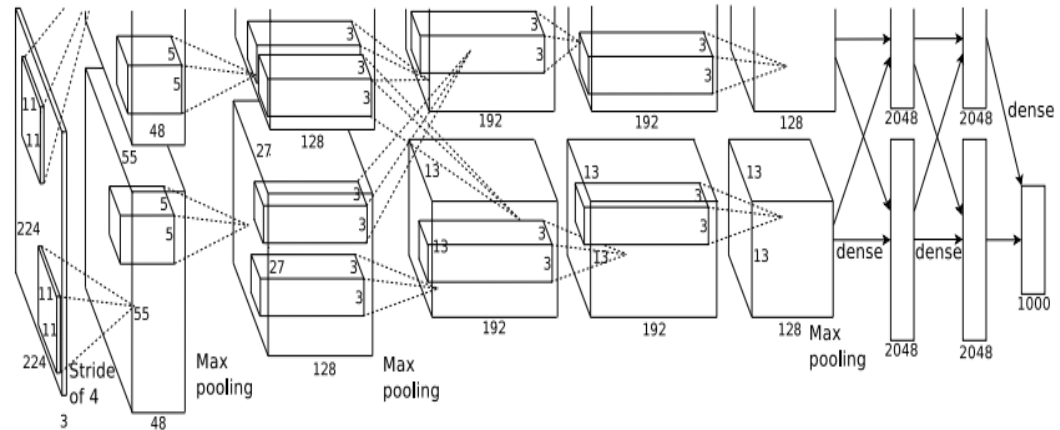
[13x13x256] **CONV5**: 256 3x3 filters at stride 1, pad 1

[6x6x256] **MAX POOL3**: 3x3 filters at stride 2

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

[1000] **FC8**: 1000 neurons (class scores)



Details/Retrospectives:

- first use of ReLU
- used Norm layers (not common anymore)
- heavy data augmentation
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
- L2 weight decay 5e-4
- **7 CNN ensemble: 18.2% -> 15.4%**

Case Study: VGG-16

- Simonyan & Zisserman won the ImageNet ILSVRC-2014 challenge.
- They published a paper called “Very Deep Convolutional Networks for Large-Scale Image Recognition”
- They compare networks of up to 19 layers.
- They also **published** their network
- It is one of the most popular pre-trained networks and available for a lot of deep learning libraries, including Keras.
- <https://keras.io/applications/>

VGG-16 Architecture

input: 224x244 RGB image

64 Conv 3x3

64 Conv 3x3

maxpool

128 Conv 3x3

128 Conv 3x3

maxpool

256 conv 3x3

256 conv 3x3

256 conv 3x3

maxpool

512 conv 3x3

512 conv 3x3

512 conv 3x3

maxpool

512 conv 3x3

512 conv 3x3

512 conv 3x3

maxpool

4096 FC

4096 FC

1000 FC

softmax

All layers using ReLU

Trained with SGD and momentum

VGG-16 - Training

- SGD + momentum
- With Gaussian initialization training needed tricks:
 - First trained smaller network (less layers)
 - Used the smaller network as initialization for deeper network
- With Glorot initialization it could be trained from scratch!
- See paper for further details:
<https://arxiv.org/pdf/1409.1556.pdf>