

Best choice of pivot item = median.

## Sorting summary

	inplace?	stable?	worst	average	best	remarks
selection	✓		$N^2 / 2$	$N^2 / 2$	$N^2 / 2$	$N$ exchanges
insertion	✓	✓	$N^2 / 2$	$N^2 / 4$	$N$	use for small $N$ or partially ordered
shell	✓		?	?	$N$	tight code, subquadratic
merge		✓	$N \lg N$	$N \lg N$	$N \lg N$	$N \log N$ guarantee, stable
quick	✓		$N^2 / 2$	$2 N \ln N$	$N \lg N$	$N \log N$ probabilistic guarantee fastest in practice
3-way quick	✓		$N^2 / 2$	$2 N \ln N$	$N$	improves quicksort in presence of duplicate keys
???	✓	✓	$N \lg N$	$N \lg N$	$N$	holy sorting grail

1v1 CONV with 32 filters preserves spatial dimensions, reduces depth.

## Data Augmentation

- Horizontal or vertical Flips.
- Random crops and scales.
- Color Jitter.
- Random mix/combinations of :
  - Translation
  - Rotation
  - Stretching
  - Shearing
  - lens distortions

# Transfer Learning

	<b>very similar dataset</b>	<b>very different dataset</b>
<b>very little data</b>	Use Linear Classifier on top layer	You're in trouble... Try linear classifier from different stages
<b>quite a lot of data</b>	Finetune a few layers	Finetune a larger number of layers

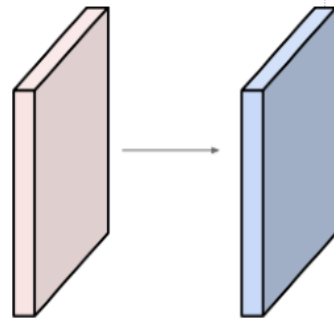
Why CNN?

- local connectivity
- parameter sharing
- pooling / subsampling hidden units

Examples time:

Input volume: **32x32x3**

**10** **5x5** filters with stride 1, pad 2



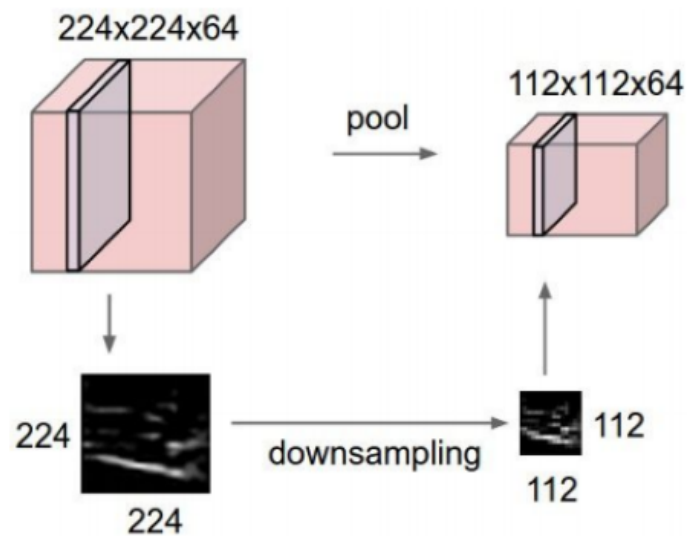
Number of parameters in this layer?

each filter has  $5*5*3 + 1 = 76$  params (+1 for bias)

=>  $76*10 = 760$

# Pooling layer

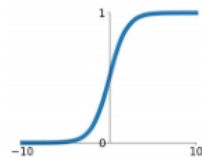
- Makes the representations smaller and more manageable.
- Operates over each activation map independently:



## Activation Functions

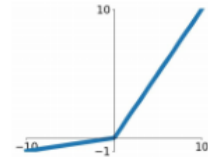
### Sigmoid

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



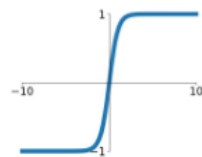
### Leaky ReLU

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



### tanh

$$\tanh(x)$$

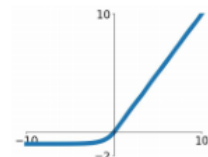


### Maxout

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

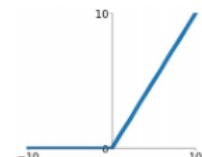
### ELU

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



### ReLU

$$\max(0, x)$$



# Batch Normalization

[Ioffe and Szegedy, 2015]

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_1 \dots x_m\}$ ;  
Parameters to be learned:  $\gamma, \beta$   
**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$
$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$
$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

- 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

Why RNN?

- Temporal dependencies
- Variable sequence length

## Keras Simple RNN

```
inp = Input(shape = (10,8))
x = SimpleRNN(4)(inp)
model = Model(inp, x)
```

```
model.summary()
```

Model: "model\_2"

Layer (type)	Output Shape	Param #
input_5 (InputLayer)	[(None, 10, 8)]	0
simple_rnn_3 (SimpleRNN)	(None, 4)	52

=====  
Total params: 52  
Trainable params: 52  
Non-trainable params: 0  
=====

## Vanishing and Exploding gradient solutions

- $\tanh(a)$
- $\max(a, 0)$

## Vanishing / exploding gradient

### 1. Exploding gradients

- Truncated BPTT
- Clip gradients at threshold
- RMSprop to adjust learning rate

### 2. Vanishing gradients

- Harder to detect
- Weight initialization
- ReLu activation functions
- RMSprop
- LSTM, GRUs

## LSTM (Long term short memory) networks

- They are capable of learning long-term dependencies.
- LSTMs are explicitly designed to avoid the long-term dependency problem.
- Remembering information for long periods of time is practically their default behavior, not something they struggle to learn.

# Keras LSTM

```
inp = Input(shape=(28,28))
x = LSTM(16)(inp)
x = Dense(10,activation='softmax')(x)
model = Model(inp, x)
model.summary()
```

Model: "model\_6"

Layer (type)	Output Shape	Param #
=====		
input_9 (InputLayer)	[(None, 28, 28)]	0
lstm (LSTM)	(None, 16)	2880
dense_2 (Dense)	(None, 10)	170
=====		
Total params: 3,050		
Trainable params: 3,050		
Non-trainable params: 0		

$28 \times 16 \times 4 + 16 \times 4 + 16 \times 16 \times 4$

2880

# Gated Recurrent Units (GRU)

- As RNNs and particularly the LSTM architecture rapidly gained popularity during the 2010s, a number of papers began to experiment with simplified architectures in hopes of retaining the key idea of incorporating an internal state and multiplicative gating mechanisms but with the aim of speeding up computation.
- The gated recurrent unit (GRU) (Cho et al., 2014) offered a streamlined version of the LSTM memory cell that often achieves comparable performance but with the advantage of being faster to compute (Chung et al., 2014).

## Keras GRU

```
inp = Input(shape=(28,28))
x = GRU(16,reset_after=False)(inp)
x = Dense(10,activation='softmax')(x)
model = Model(inp, x)
model.summary()
```

Model: "model\_9"

Layer (type)	Output Shape	Param #
=====		
input_12 (InputLayer)	[(None, 28, 28)]	0
gru_2 (GRU)	(None, 16)	2160
dense_5 (Dense)	(None, 10)	170

=====  
Total params: 2,330  
Trainable params: 2,330  
Non-trainable params: 0

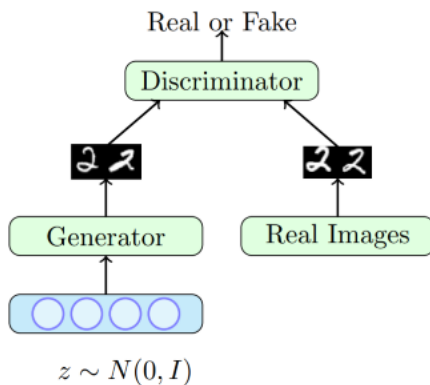
```
28*16*3+16*16*3+16*3
```

```
2160
```

# Properties of Autoencoders

- **Data-specific:** Autoencoders are only able to compress data similar to what they have been trained on.
- **Lossy:** The decompressed outputs will be degraded compared to the original inputs.
- **Learned automatically from examples:** It is easy to train specialized instances of the algorithm that will perform well on a specific type of input.

## The intuition



- What can we use for such a complex transformation?  
A Neural Network
- How do you train such a neural network?  
Using a two player game
- There are two players in the game: a **generator** and a **discriminator**
- The job of the generator is to produce images which look so natural that the discriminator thinks that the images came from the real data distribution
- The job of the discriminator is to get better and better at distinguishing between true images and generated (fake) images



# VAE Encoder

- The encoder takes input and returns parameters for a probability density (e.g.,  $q_{\theta}(z | x)$ ): i.e., gives the mean and co-variance matrix.
  - We can sample from this distribution to get random values of the lower-dimensional representation  $z$ .
  - Implemented via a neural network: each input  $x$  gives a vector mean and diagonal covariance matrix that determine the Gaussian density
  - Parameters  $\theta$  for the NN need to be learned – need to set up a loss function.
- 
- The discriminator wants to maximize the second term whereas the generator wants to minimize it (hence it is a two-player game)
  - Architecture guidelines for stable Deep Convolutional GANs
    - Replace any pooling layers with strided convolutions (discriminator) and fractional-strided convolutions (generator).
    - Use batchnorm in both the generator and the discriminator.
    - Remove fully connected hidden layers for deeper architectures.
    - Use ReLU or LeakyReLU activation in generator for all layers except for the output, which uses tanh/sigmoid.
    - Use LeakyReLU activation in the discriminator for all layers
- 
- Advantages of WGAN:
    - Better training stability.
    - Loss has meaning so it helps in termination criteria
  - 
  - Disadvantages of WGAN
    - Longer to train.

# The full algorithm

---

**Algorithm 1** WGAN, our proposed algorithm. All experiments in the paper used the default values  $\alpha = 0.00005$ ,  $c = 0.01$ ,  $m = 64$ ,  $n_{\text{critic}} = 5$ .

---

**Require:** :  $\alpha$ , the learning rate.  $c$ , the clipping parameter.  $m$ , the batch size.  $n_{\text{critic}}$ , the number of iterations of the critic per generator iteration.

**Require:** :  $w_0$ , initial critic parameters.  $\theta_0$ , initial generator's parameters.

```
1: while  $\theta$  has not converged do
2:   for  $t = 0, \dots, n_{\text{critic}}$  do
3:     Sample  $\{x^{(i)}\}_{i=1}^m \sim \mathbb{P}_r$  a batch from the real data.
4:     Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
5:      $g_w \leftarrow \nabla_w \left[ \frac{1}{m} \sum_{i=1}^m f_w(x^{(i)}) - \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)})) \right]$ 
6:      $w \leftarrow w + \alpha \cdot \text{RMSPProp}(w, g_w)$ 
7:      $w \leftarrow \text{clip}(w, -c, c)$ 
8:   end for
9:   Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
10:   $g_\theta \leftarrow -\nabla_\theta \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)}))$ 
11:   $\theta \leftarrow \theta - \alpha \cdot \text{RMSPProp}(\theta, g_\theta)$ 
12: end while
```

---

## Training GAN

```
1.  procedure GAN Training
2.    for number of training iterations do
3.      for k steps do
4.        Sample minibatch of m noise samples  $\{z^{(1)}, \dots, z^{(m)}\}$  from noise prior  $p_g(z)$ 
5.        Sample minibatch of m examples  $\{x^{(1)}, \dots, x^{(m)}\}$  from data generating distribution  $p_{\text{data}}(x)$ 
6.        Update the discriminator by ascending its stochastic gradient:

7.          
$$\nabla_\theta \frac{1}{m} \sum_{i=1}^m \left[ \log D_\theta(x^{(i)}) + \log(1 - D_\theta(G_\phi(z^{(i)})) \right]$$

8.        end for
9.        Sample minibatch of m noise samples  $\{z^{(1)}, \dots, z^{(m)}\}$  from noise prior  $p_g(z)$ 
10.       Update the generator by ascending its stochastic gradient

11.       
$$\nabla_\phi \frac{1}{m} \sum_{i=1}^m \left[ \log(D_\theta(G_\phi(z^{(i)}))) \right]$$

12.     end for
13.  end procedure
```