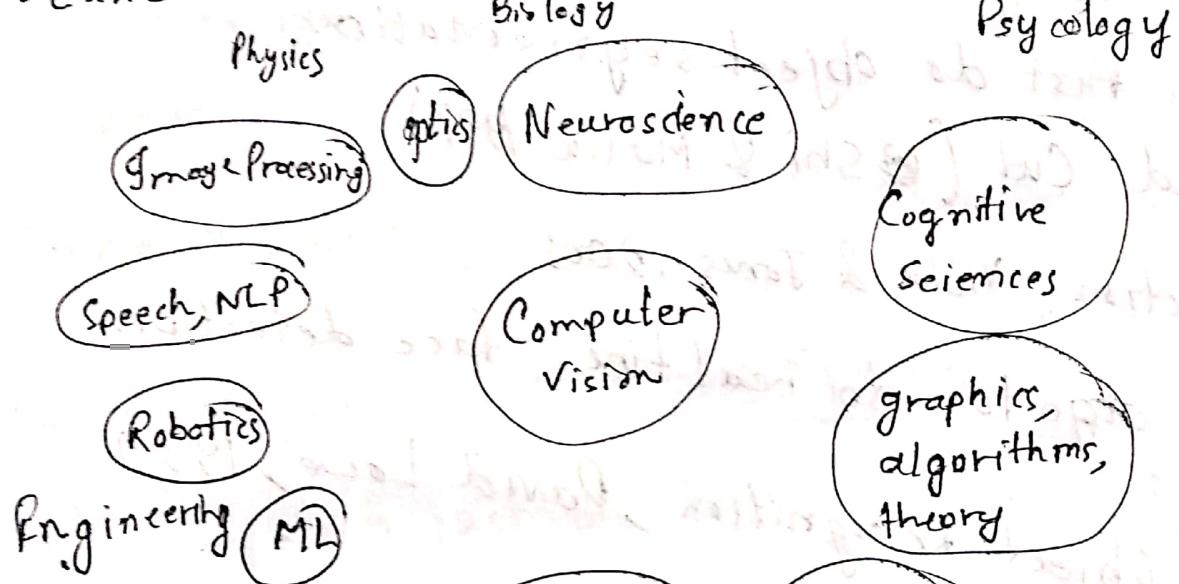


# CS231n: Convolutional Nets for Visual Recognition

Fei-Fei Li; Justin Johnson; Serena Yueng

Computer vision is the study of visual data

• called the dark matter of internet



Hubel & Wiesel, 1959

Electrical signal from brain  $\rightarrow$  stimulus

$\rightarrow$  mapping  $\leftarrow$

Block world [Larry Roberts, 1963]

1966 → Summer Vision Project

## Generalized Cylinder

1979

Pictorial Structure

1973

Nonna.

If object recognition is too hard, maybe we should first do object segmentation

Normalized Cut (shi & Malik 1997)

Face Detection, Viola & Jones, 2001

Adaboost algo to do real time face detection

'SIFT' & Object Recognition, David Lowe, 1990

HOG

Pascal Visual Object Challenge (10 categories)

Imagenet (Deng, Dong, Socher, Li, & Fei-Fei, 2009)  
Stanford & Princeton

(22K categories & 14M images)

Took 3 years (took help from Amazon Mechanical Turk)

② Imagenet Classification Challenge

There is a number of visual recognition problems that are related to image classification, such as object detection, image captioning

Lecun et al (conv-net)

Alexnet

Image  $\rightarrow$  Captions

## Lecture 2: Image Classification Pipeline

- An image is just a big grid of numbers between [0, 255]

Attempts: finding edges (not scalable)

## Data Driven approach

- 1) Collect data
- 2) Use ML to train
- 3) Evaluate the classifier

## Distance Metric

$$L_1 \text{ Distance} \Rightarrow d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$

Nearest Neighbor  $O(1)$

Memoize data

$$L_2 \text{ dist} \Rightarrow d_2(I_1, I_2) = \sqrt{\sum_p (I_1^p - I_2^p)^2}$$

for each test image

find closest train image  $O(N)$

Predict label of nearest image

This is bad: we want classifiers that are fast at prediction; slow for training is ok

k-nearest neighbor (to get rid of noisy answer)

Hyperparameters: choices about the algo that we set rather than learn

Setting hyper params

K-NN on images never used

- very slow at testing

- distance metric not informative (not good to check similarity)

- curse of dimensionality

## Linear Classification

### Parametric Approach

$f(x, w) \xrightarrow{\text{Weights}} 10 \text{ numbers giving class scores}$

$$f(x, w) = \underbrace{w^T x + b}_{\substack{10 \times 1 \\ 10 \times 3072 \\ 3072 \times 1 \\ 10 \times 1}}$$

→ Linear classifier is learning only one template

## Lecture 3 - Loss Function & Optimization

### Challenges of recognition

i) Viewpoint

ii) Illumination

iii) Deformation

iv) Occlusion

v) Clutter

vi) Intraclass variation

- Define a loss function that quantifies our unhappiness with the scores across the training set.
- Come up with a way of efficiently finding parameters that minimize the loss function

$$L = \frac{1}{N} \sum_i L_i(f(x_i, w), y_i)$$

### Multiclass SVM loss

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_j > s_{y_i} + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

true value

$$= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \quad s = f(x_i, w)$$

$$\text{Then, } L = \frac{1}{N} \sum_{i=1}^N L_i$$

What happens to loss if a score is changed a bit?

→ Nothing happens because it still returns zero loss.

What is the min/max possible loss for SVM

→ min → 0  
max → ∞

Q3] At initialization  $\omega$  is small so all  $s_i > 0$ .

What is the loss?

(number of classes - L)

Q4] What if sum as ~~is~~ over all classes?

The loss increases by L

This is just for convention we omit the correct class so that our minimum loss is zero.

Q5] What if we used mean instead of sum?

→ answer would be same because we don't care about true scores

Q5] What if  $\sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)^2$

→ This would end up giving a different loss function that's not linear

Q6] If we find a  $\omega$  that gives  $L(\omega)$  is this  $\omega$  unique?  
No, there are many other  $\omega$ s. Like  $2\omega$

Regularization: Model should be simple so that it works on test data

$$L(\omega) = \frac{1}{N} \sum_{i=1}^N l_i(f(x_i, \omega), y_i) + \lambda R(\omega)$$

Pecan's Razor:

Among competing hypotheses, the simplest is the best

L2 regularization  $R(\omega) = \sum_k \sum_e \omega_{k,e}^2$

L1 regularization  $R(\omega) = \sum_k \sum_e |\omega_{k,e}|$

Elastic net ( $L1 + L2$ )  $R(\omega) = \sum_k \sum_e \beta \omega_{k,e}^2 + \gamma |\omega_{k,e}|$

Max norm regularization

Dropout

Batchnorm

Stochastic depth

If you're Bayesian: L2 regularization also corresponds

MAP inference using a Gaussian prior on  $\omega$

Softmax classifier (Multinomial Logistic Regression)

$$P(Y=k | X=x_i) = \frac{e^{s_i}}{\sum_j e^{s_j}} \quad \text{where } s = f(x_i; \omega)$$

$$L_i = -\log P(Y=y_i | X=x_i) \quad L_j = -\log \left( \frac{e^{s_{y_i}}}{\sum_j e^{s_j}} \right)$$

Want to maximize log likelihood, or to minimize the negative log likelihood of the ~~the~~ correct class

Q1: What is the min. and max loss?

$$\min = 0 \\ \max = \infty$$

Q2: Usually at initialization  $w$  is small so all  $s \approx 0$  what is the loss?

$$= -\log\left(\frac{1}{c}\right)$$

### Softmax vs SVM

Q: Suppose I take a datapoint and jiggle a bit (changing its score slightly). What happens to the loss in both cases?

SVM doesn't care about the score. Softmax continuously try to ~~make~~ the datapoints like pushing the correct to  $\infty$  and pushing the incorrect to  $-\infty$

$$\text{Softmax}, t_i = \log\left(\frac{e^{s_i}}{\sum_j e^{s_j}}\right)$$

$$\text{SVM}, L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$\text{Full Loss } L = \frac{1}{N} \sum_{i=1}^N L_i + R(w)$$

## Optimization:

Strategy #1 Random search

Strategy #2 follow the slope

In 1 dimension, the derivative of a function

$$\frac{\delta f(x)}{\delta x} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \quad [\text{slope}]$$

In multiple dimensions, the gradient is the vector of partial derivatives along each direction

The slope in any direction is the dot product of the direction with the gradient

The direction of steepest descent is the negative gradient.

→ This is super slow and super bad.

We can use compute analytic gradient

Numerical gradient: approximate, slow, easy to write

Analytic gradients: exact, fast, error-prone

Gradient check: Using analytic gradient to find

grads but checking with numerical gradient to  
• This is an interesting debugging tool.

### Gradient descent

- i) find grads
- ii) weights  $\leftarrow \frac{\text{stepsize}}{\uparrow \text{hyperparameter}} \times \cancel{\text{grads}}$   
learning rate

### Stochastic Gradient Descent:

Full sum is expensive when it's large

Approx sum using minibatch of examples

32/64/128 common

### Image Features

Color Histogram

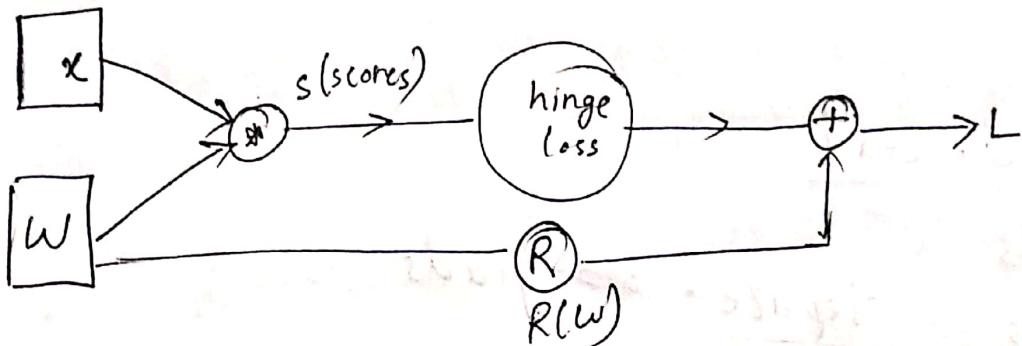
Histogram of Oriented Gradients

Bag of Words (Build Codebook, Encode Images)

## Image Features vs ConvNets

ConvNets (Krizhevsky 2012) AlexNet

## Lecture 4: Intro to Neural Nets



Leverage chain rule

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

$$\frac{d}{dx}(e^x) = e^x$$

$$\frac{d}{dx}(\alpha x) = \alpha$$

$$\frac{d}{dx}\left(\frac{1}{x}\right) = -\frac{1}{x^2}$$

$$f(x) = c+x = \alpha$$

$$\frac{d}{dx}(c+x) = 1$$

$$\begin{aligned} \frac{d \delta(x)}{dx} &= \frac{e^{-x}}{(1+e^{-x})^2} \\ &= \left( \frac{1+e^{-x}-1}{1+e^{-x}} \right) \left( \frac{1}{1+e^{-x}} \right) \\ &= [1-\delta(x)] \delta'(x) \end{aligned}$$

If any problem with gradients,  
break down to computational  
graph

add gate: gradient distributor

Q: What is a max gate?

the highest one back is gonna get the max value. other will be zeroed out

multip gate: gradient switcher

local gradient is the value of the other variable

Neural Nets:

2-layer Network  $f = w_2 \max(0, w_1 x)$

Biological Neurons are far more complicated

Activations

Sigmoid

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

Leaky ReLU

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

tanh

$$\tanh(x)$$

Maxout

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

ReLU

$$\max(0, x)$$

ELU

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

$$\nabla_w f = 2q \cdot w^T$$

$$\nabla_w f = 2w^T q$$

## Lecture 5: Convolution Neural Networks

- Try to maintain spatial structure

History of NN and CNN  
→ Mark I perception machine was the first implementation of the perceptron algorithm.  
\* no backprop

$$f(x) = \begin{cases} 1 & w \cdot x + b \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Frank Rosenblatt (1957)

→ Widrow and Hoff (1960) : Adaline/Madaline  
Multilayer perception  $\hookrightarrow$  no backprop

Rumelhart (1986)  $\rightarrow$  Backprop

Hinton and Salakhutdinov 2006 ]  $\rightarrow$  Reinigorated research in Deep Learning  
First strong results [2012] [Krizhevsky, 2012]

- ImageNet classification with deep CNN  
AlexNet

Neurocognition [Fukushima, 1980]

[Lecun, Bottou, Bengio, Haffner, 1998] -  
Gradient-based learning applied to  
document recognition

ConvNets for classification, retrieval, detection,  
segmentation, self driving car, pose estimation, game play  
Image Captioning

NVIDIA Tesla line [GPU] powers Neural Style Transfer

## Convolutional Neural Networks

Convolutional layers → preserve spatial structure

→ We call the layer convolutional because it is related to convolution of two signals

$$f[x, y] * g[x, y] = \sum_{n_1=-\infty}^{\infty} \sum_{n_2=-\infty}^{\infty} f[n_1, n_2] \cdot g[x-n_1, y-n_2]$$

elementwise multiplication  
and sum of a filter and  
the signed (image)

## Conv Layer

- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires four hyperparams:
  - $K \rightarrow \#$  of filters
  - $F \rightarrow \#$  their spatial extent
  - $S \rightarrow \#$  stride
  - $P \rightarrow$  Amount of zero padding
- \* Produces a volume of size  $W_2 \times H_2 \times D_2$ , where
  - $W_2 = (W_1 - F + 2P)/S + 1$
  - $H_2 = (H_1 - F + 2P)/S + 1$
  - $D_2 = K$
- with parameter sharing. It introduces  $F \cdot P \cdot D_1$  weights per filter for a total of  $(P \cdot F \cdot D_1) \cdot K$  weights and  $K$  biases.
- \* The output volume, the  $d$ -th depth slice (of size  $W_2 \times H_2$ ) is the result of performing a valid convolution of the  $d$ -th filter over the input volume with a stride of  $S$ , and then offset by  $d$ -th bias.

## Pooling layer

- makes the representation smaller and more manageable
- operates over each activation map independently

## Pooling layer:

Accepts a volume of size  $w_1 \times h_1 \times d_1$

- Requires three hyperparameters

→ their spatial extent  $f$

→ the stride  $s$

- Produces a volume of size  $w_2 \times h_2 \times d_2$  where:

$$\rightarrow w_2 = (w_1 - f) / s + 1$$

$$\rightarrow h_2 = (h_1 - f) / s + 1$$

$$\rightarrow d_2 = d_1$$

- Introduces zero parameters since it computes a fixed function of the input

- Note that it is not common to use zero-padding for pooling layers

## Fully connected Layer (FC-Layer)

→ Contains neurons that connect to the entire input volume, as in ordinary Neural Networks

## Typical architectures

$$[(\text{Conv-ReLU})^N \text{-POOL?}] * M - (\text{FC-ReLU}) * k, \text{ softmax}$$

- ResNet / GoogleNet challenge this paradigm

## Lecture 6: Training NNs I

### Minibatch SGD

Loop:

- 1) Sample a batch of data
- 2) forward prop it through the graph, get loss
- 3) Backprop to calculate the gradients
- 4) update the parameters using the gradient

P1

## Activation Function

### Data Preprocessing

### Weight Initialization

### Batch Normalization

### Babysitting the Learning Process

### Hyperparameter Optimization

#### Activation functions:

\* Sigmoid  $\Rightarrow$  Squashes number to range  $[0, 1]$

$\rightarrow$  historically popular [since they have nice interpretation as a saturating "firing rate" of a neuron]

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

#### Problems:

- 1) Saturated neurons "kill" the gradients
- 2) Sigmoid outputs are not zero centered
- 3)  $\exp(\cdot)$  is a bit computationally expensive

- \*  $\tanh(x)$ 
  - squashes numbers to range  $[-1, 1]$
  - zero-centered
  - still kills gradients when saturated  
[Krizhevsky et. al. 2012]

- \* ReLU (Rectified Linear Unit)  $f(x) = \max(0, x)$ 
  - Does not saturate (in + region)
  - very computationally efficient
  - converges much faster than sigmoid/tanh in practice  $[6x]$
  - actually more biologically plausible than sigmoid

cons → not zero centered output

An annoyance:

what is the gradient when  $x < 0$ ?

- \* Leaky ReLU  $\rightarrow$  to activate ~~the~~ activate dying ReLU
  - Does not ~~s~~
  - will not die  $f(x) = \max(0.01x, x)$

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

## Exponential Linear Unit (ELU) [Clevert et. al. 2015]

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

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

- computation requires exp()

## Maxout Neuron

[Goodfellow et. al. 2013]

- Doesn't have the basic form of dot products
- nonlinearity
- generalizes ReLU and Leaky ReLU
- Linear Regime! Doesn't saturate! Doesn't die!

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

problem: doubles the number of param. neurons

## 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

zero center and normalization so that all of them can contribute PCA and whitened data (for high dimensional ones)

In practice for Images (center only)

→ Not common to normalize variance, to do PCA or whitening

## Weight Initialization

What happens when  $\omega = 0$  init is used?

→ all the neuron start learning the same thing

First idea: small random numbers

(gaussian with zero mean and  $1e-2$  std deviation)

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

→ Works ~okay for small nets but problems with deeper networks

All activation become zero for bigger nets.

What do gradients look like?

things get smaller and smaller and later'

it's the same

Almost all neurons completely saturate and either -1 and 1  
and grads will be zero

Xavier initialization [Glorot et al. 2010]

Reasonable → mathematical derivation assumes linear  
activations

- but when using the ReLU non-linearity it breaks

He et al. 2015 (note initialization/2)

- proper initialization is an active area of research

Batch Normalization:

[Toffo and Segey, 2015]

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

- Consider a batch of activations at some layer.

Usually inserted after Fully connected or conv layers and before non-linearity

- do we necessarily want to a unit gaussian input to an activation layer

$$\rightarrow \hat{y}^{(k)} = \underbrace{\gamma^{(k)} \hat{x}^{(k)}}_{\text{scale}} + \underbrace{\beta^{(k)}}_{\text{shift}} \quad \# \text{ scale and shift}$$

network can squash the range if it wants to (learnable)

- 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

Note: at test time BatchNorm layer are not computed based on batch. Instead, a single fixed empirical mean of activations during training is used.

## Babysitting Training Process

check loss

Overfit on small portion

take the first 20 examples from CIFAR-10

- turn off regularization

- use simple vanilla 'sgd'

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

→ loss not going down → learning rate <sup>too</sup> ~~low~~

increase learning rate ~~and~~ <sup>too high</sup> → loss exploding  
Loss reaches NaN

## Hyperparam Optimization

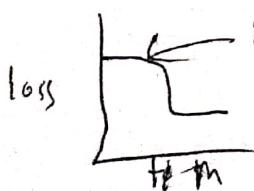
Cross-validation strategy

coarse → fine

run coarse search for some epochs

Random search, Grid Search  
Bengio 2011

Bad initialization: prime suspect



## Lecture 7: Training NNs II

Fancier optimization

Regularization

Transfer learning

### Optimization!

Vanilla gradient descent

#### Problems with SGD

What if loss changes quickly in one dimension and slowly in another? what does gradient do?

• zigzag

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

Zero gradient, gradient descent gets stuck

- saddle point much common in high dimension
- Our gradients come from minibatches so they can be noisy.

- Full batch gradient descent doesn't solve these problems.

### SGD + momentum

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

$$x_{t+1} = x_t - \alpha v_{t+1}$$

- Build up velocity as a running mean of gradients
- $\rho$  gives friction: typically  $\rho = 0.9$  or  $0.99$

This solves all the problem is SGD

### Nesterov momentum:

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



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

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

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

$$= x_t + v_{t+1} + \rho(v_{t+1} - v_t)$$

## Adagrad (don't choose in NN)

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

Step size gets smaller and smaller.

## RMSProp

step size is managed to ~~stay as~~ shrink slowly

## Adam (almost)

Momentum + Adagrad/RMSProp

what happen at the first step

## Full form

Default

$\beta_1 = 0.9$

$\beta_2 = 0.999$

$\text{lr} = 1e-3 \text{ or } 5e-4$

### Momentum

- Bias correction  $\rightarrow$  for the fact that first and second momentum estimates start at zero.

Which one of those learning rates is best to use?

step decay  $\rightarrow$  decaying by half after each epoch

exponential decay:  $d = d_0 e^{-kt}$

$1/t$  decay  $\Rightarrow d = d_0 / (1 + kt)$

Linear decay is more common with SGD momentum

common with Adam

### First order optimization

- (1) use gradient form linear approx
- (2) use step to minimize approx

### Second order approximation:

- (1) Use gradient and Hessian to form quadratic approx

- (2) Step to the minima of the approx

$$J(\theta) = J(\theta_0) + (\theta - \theta_0)^T \nabla_{\theta} J(\theta_0) + \frac{1}{2} (\theta - \theta_0)^T H(\theta - \theta_0)$$

$$\theta^* = \theta_0 - t^{-1} H_{\theta_0}^{-1} J(\theta_0) \quad [\text{No learning rate}]$$

Practically bad for deep learning [impractical]  
Hessian has  $O(N^2)$  and inverting takes  $O(N^3)$

$N$  is large

## In practice

- Adam is a good default
- L-BPGS can be tried out if we can afford full batch updates (can be good for style transfer)

## Beyond training error

### Model Ensembles

- Train multiple independent models
  - At train time average their results
- Enjoy 2x extra performance
- Instead of training independent models, use multiple snapshots of a single model during training
  - keep a moving average of the parameter vector and use that at test time (Polyak averaging)

Improve single model performance

Regularization

$$L = \frac{1}{N} \sum_{i=1}^N \sum_j \max(0, f(x_i; w)_j - f(x_i; w)y_i + 1) + \lambda R(w)$$

$R(w)$

L2 Reg

L1  
Elastic Net ( $L1 + L2$ )

Dropout

Randomly set some neurons to zero. 0.5 is common

Forces the network to have a redundant representation;

Prevents co-adaption of features

Dropout: Test time

Dropout makes our output random

Want to 'average out' the randomness at

test time

$$y = f(x) = \mathbb{E}_z [f(x, z)] = \int p(z) f(x, z) dz$$

$$E[a] = w_1x + w_2y$$

During training we have  $E[a] = \frac{1}{4}(w_1x + w_2y) + \frac{1}{2}(w_1x + w_2y)$

At test time, multiply  
by dropout prob

$$+ \frac{1}{2}(0x + 0y) + \frac{1}{4}(0x + w_2y)$$
$$= \frac{1}{2}(w_1x + w_2y)$$

### Data Augmentation:

Horizontal flip

Crop

Occlusion

Contrast and brightness

translation

rotation

stretching

shearing

Lens distortion

(go crazy)

### Drop Connect

We zero out the value of weight matrices

almost like dropout

### Fractional Max Pooling

Stochastic Depth Randomly drop network

## Transfer Learning

Reduce overfitting

Reduce need for huge data

### Small dataset

freeze the taken ones and train.

Later try to train on the whole network

for different dataset fine-tuning  
we need to try out ~~more~~ training more layers