

# CS 231n Convolutional Neural Networks for Visual Recognition

## Lecture 1 Introduction

Computer vision → study of visual data.

2015 CISCO study estimated that by 2017, roughly 80% of all traffic on the Internet would be video. (from a pure number of bits perspective)  
Hubel & Wiesel, 1959.

Visual processing starts with simple structures of the visual world, oriented edges and as information moves along the visual processing pathway, the brain builds up the complexity of the visual information.  
David Marr 1970s.

Input image → Primal sketch → 2.5D sketch → 3D model representation.  
(Edge, image)

Brook & Buhfot 1979, Generalized Cylinder

Fischler and Elschlager, 1973. Pictorial Structure

basic idea → every object is composed of simple geometric primitives.

Sudeh Malik, 1997. Normalized Cut. → group pixels into meaningful areas.

Object Segmentation, image segmentation

Viola & Jones, 2001 Face Detection

AdaBoost algorithm → real-time face detection

first digital camera with face detection → 2006, Fuji Film

David Lowe, 1999. "SIFT" & Object Recognition.

feature based object recognition. (to match an entire object)  
invariant to changes. (light, viewpoint, etc)

Schmid & Ponce 2006, Spatial Pyramid Matching

SVM on feature descriptor.

Felzenszwalb, McAllester, Ramanan 2009, Deformable Part Model

Dalal & Triggs 2005, Histogram of Gradients.

PASCAL visual object challenge.

IMAGENET. Large Scale Visual Recognition Challenge.

物体识别精度高，但容易过拟合。训练样本数量不够时。

Convolutional Neural Network 虽早在90年代发明，但取决于计算能力  
和数据质量，直到2012年才取得突破。  
↓  
需要大量带标签的图像和像素

## Lecture 2 Image Classification Pipeline

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

e.g. 800x600x3 (3 channels RGB)

Challenges: → algorithm must be robust.

- viewpoint variation
- illumination.
- occlusion 遮挡.
- Deformation 形变.
- background clutter 背景杂乱
- Intra-class variation 类内差异.

- o hard-core algorithm (e.g. find edges, find corners...)  
↳ super brittle 不结实.  
not scalable, 对不同类别需要重新写对应算法

- o Data-Driven Approach.

1. collect a dataset of images and labels
2. Use machine learning to train a classifier
3. Evaluate the classifier on new images.

train → predict

## First classifier: Nearest Neighbor.

```

def train(images, labels): # Memorize all data and labels
    return model

def predict(model, test_images): # Predict the label of the most
    return test_labels
        similar training image.

```

with  $N$  examples, Train  $O(1)$ , Predict  $O(N)$ . ← Bad.

We want classifiers that are fast at prediction, slow for training is ok.

### Tips:

在解决 computer vision 问题时，在不同的角度反复思考验证很有用。

- idea of high dimensional points in the plane  
如将每张图片看作是高维度空间中的一点。
- concrete images (high dimensional vectors) of pixels.  
具体到图片像素。

## K-Nearest Neighbors

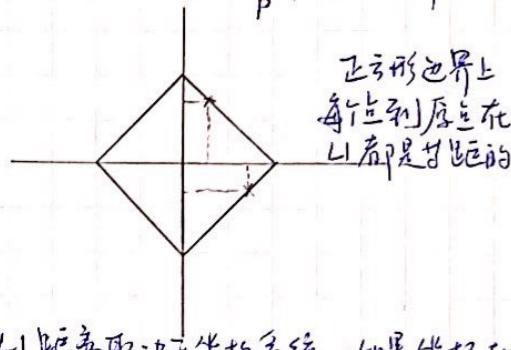
Instead of copying label from nearest neighbor, take majority vote from  $K$  closest points.

$K$  越大使得决策边界更平滑。

### Distance Metric.

$L_1$  (Manhattan) distance.

$$d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$



$L_1$  距离取决于坐标系，如果坐标轴转动， $L_1$  也会改变。

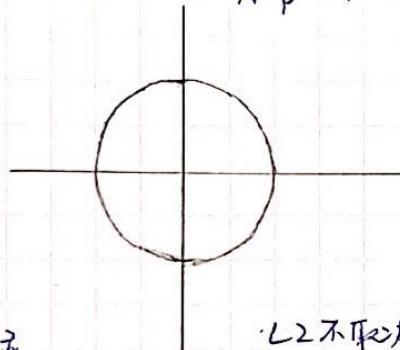
在输入 特征 的值具有重要意义（如用户年龄、收入等）， $L_1$  距离更合适。  
(features)

在 KNN 算法对  $L_1$  中， $L_1$  的结果显示决策边界倾向于沿着坐标轴。

而  $L_2$  结果显示不受坐标轴限制，边界比较圆润。

$L_2$  (Euclidean) distance

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



$L_2$  不取决于坐标系。

4

Hyperparameters, choices about the algorithm that we set rather than learn.



~~for~~, K-NN 中的 K, distance metric.

very problem-dependent, must try them all out and find the optimal.

找到最佳 hyperparameters 的训练模型的方式:

将数据集分为 train, validation, test.

In training set 上, 用不同的 hyperparameter 训练算法. In validation set 上评估找出表现最好的一组参数. 最后在 test set 上跑一次得到最终结果.

(test set 独立于训练数据, 检验模型泛化能力) 相当于“新数据”

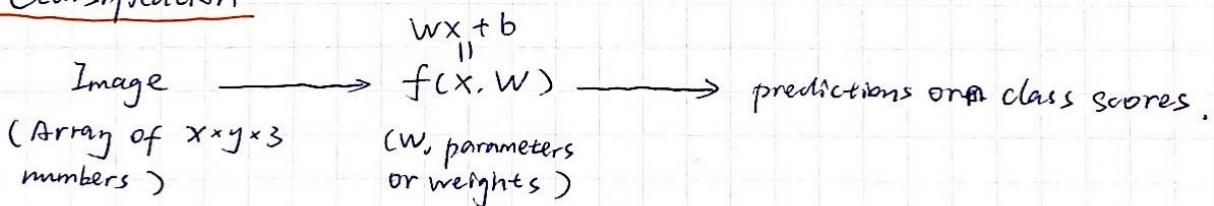
(cross-validation), split data into folds + test set, try each fold as validation and average the results)

实际应用中, K-NN on Images never used.

- very slow at test time
- Distance metrics on pixels are not informative.
- Curse of dimensionality.

随着维度增加, 计算量呈指数级增加.

## Linear Classification



通过训练得到  $w$  包含模型信息, 以后可以直接地应用该使用, 不需要原本的 data set  $x$ .

$x$  在图片中的 pixels.

弊端: - 有在线性不可分情况

- one single template per category ( $w$  中每一行对应一个类, 每个值对应一个 pixel)

换一个角度: 在高维空间中, 每个图片为一个点, Linear classifier 在决策边界上画一个线性分界面来划分一个类别.

## Further Reading

< A Few Useful Things to Know about Machine Learning >

"Recognizing and Learning Object Categories" [people.csail.mit.edu/torralba/shortCourseRLOC/index.html](http://people.csail.mit.edu/torralba/shortCourseRLOC/index.html)

## Lecture 3 Loss Functions and Optimization

1. Loss function quantifies our unhappiness with the scores across training data.
2. Optimization, a way of efficiently finding the parameters that minimize the loss function.

### Multiclass Support Vector Machine Loss

The SVM loss is setup so that the SVM "wants" the correct class for each image to have a score higher than the incorrect classes by some fixed margin  $\Delta$ .

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

data  $(x_i, y_i)$

$s_j$ : score for the  $j$ th class

$$s_j = f(x_i, w)_j$$

e.g. linear score function  $f(x_i, w) = w^T x_i$

$$L_i = \sum_{j \neq y_i} \max(0, w_j x_i - w_{y_i} x_i + \Delta)$$

$\max(0, -) \rightarrow$  hinge loss (max-margin loss)

Regularization.

$$L = \underbrace{\frac{1}{N} \sum_i L_i}_{\text{data loss}} + \underbrace{\lambda R(w)}_{\text{regularization loss}}$$

Full loss

flash-back,

$$\min \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i \quad \rightarrow \text{SVM with } L_1 \text{ regularization}$$

$$\text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1, i=1, \dots, m$$

$$\xi_i \geq 0, i=1, \dots, m$$

Occam's Razor: Among competing hypotheses, the simplest is the best.

Regularization methods:

- $L_1 \quad R(w) = \sum \|w\|_1$
- $L_2 \quad R(w) = \sum \|w\|_2^2$
- Elastic net ( $L_1 + L_2$ )  $R(w) = \sum (\beta \|w\|_2^2 + \|w\|_1)$
- Max norm regularization
- Dropout
- Batch normalization, stochastic depth ...

Interpretation for  $L_1$  and  $L_2$  regularization:

$L_1$  regularization  $\rightarrow$  encouraging sparsity

$L_2$  regularization  $\rightarrow$  for  $W$ , spread the influence across all the values in  $x$  to generate  $y_i$ .  
robust to the change of  $x_i$ .  
Sensitive to outlier.

Practical Considerations:

- The exact value of the margin between the scores (e.g.  $\Delta=1$  or,  $\Delta=100$ ) is meaningless because the weights can shrink or stretch the differences arbitrarily. (it's safe to set  $\Delta=1.0$  in all cases)
- In the case with Neural Networks in general, we will always work with the optimization objectives in their unconstrained primal form. (many of these objectives are technically not differentiable, but in practice, this is not a problem and it's common to use a subgradient).

### Softmax classifier

$$L_i = -\log \left( \frac{e^{f_{y_i}}}{\sum_j e^{f_j}} \right) = L_i = -f_{y_i} + \log \sum_j e^{f_j} \rightarrow \text{cross-entropy loss.}$$

↓  
softmax function

The cross-entropy between a "true" distribution  $p$  and an estimated distribution  $q$  is defined as:

$$H(p, q) = - \sum_x p(x) \log q(x)$$

The softmax classifier is hence minimizing the cross-entropy between the estimated class probabilities ( $q = \frac{e^{f_j}}{\sum_i e^{f_i}}$ ) and the "true" distribution ( $p = [0 \dots 1 \dots 0]$  contains a single 1 at the  $y_i$ -th position).

$$H(p, q) = H(p) + D_{KL}(p || q)$$

In this case  $H(p) = 0$ . Hence, softmax or this cross-entropy objective wants to minimize the KL divergence between the two distributions (a measure of distance). ( $D_{KL}(p || q) = \sum_x p(x) \cdot \log \frac{p(x)}{q(x)}$ )

Maximum Likelihood and Maximum a posteriori estimation.

(ML)

(MAP)

Estimate an unobserved population parameter  $\theta$  on the basis of observation  $x$ .

$f$  is the sampling distribution of  $x$ ,  $\rightarrow f(x|\theta)$

ML: maximum likelihood estimate of  $\theta \rightarrow \hat{\theta}_{MLE}(x) = \underset{\theta}{\operatorname{argmax}} f(x|\theta)$

MAP: assume a prior distribution  $g$  over  $\theta$  exists,

$$f(\theta|x) = \frac{f(x|\theta)g(\theta)}{\int_{\theta} f(x|\theta)g(\theta)d\theta} \quad P(AB) = P(A|B)P(B) = P(B|A)P(A)$$

$$\hat{\theta}_{MAP}(x) = \underset{\theta}{\operatorname{argmax}} f(\theta|x) = \underset{\theta}{\operatorname{argmax}} f(x|\theta)g(\theta)$$

For softmax classifier, we can interpret the  $R(w)$  in full loss function as coming from a gaussian prior over the weight matrix  $W$ , where instead of MLE, we perform MAP-estimation.

Practical issue:

Numeric stability,  $\frac{e^{f_{y_i}}}{\sum_j e^{f_j}} = \frac{c \cdot e^{f_{y_i}}}{c \cdot \sum_j e^{f_j}} = \frac{e^{f_{y_i} + \log c}}{\sum_j e^{f_j + \log c}}$  (common choice:  $\log c = -\max_j f_j$ )

Dividing large  $\sum_j e^{f_j}$  can be numerically unstable.

Softmax v.s. SVM

SVM gets correct score to be greater than a margin above the incorrect scores.

Softmax interprets the score as (unnormalized) log probabilities for each class and then encourages the (normalized) log probability of correct class to be high.

SVM 不在于单体 in score, 而是正负类 score 之差出错时在一定 margin 就满足, 但 softmax classifier 是 never fully happy with the scores, 因为正类 score 总会有更高的 probability (越近 1).

## Optimization

Visualizing the loss function.

1. generate random ~~W~~ (direction) -  $w_1, w_2, w$ .
2. evaluating  $L(w + aw_1 + bw_2)$  for different values of  $a$  and  $b$
3. plot the  $L$  with  $a, b$  as axis  $x$  and  $y$ .  
for 1 dimension, just  $a$ .

~~Note~~ Once we extend the score functions  $f$  to Neural Networks, the objective functions will become non-convex, and the visualizations will not feature bowls but complex bumpy terrains. (subgradient exits and commonly used.)

Strategy #1. Random search (bad idea)

Strategy #2. Random Local Search (wasteful and computationally expensive)

Start out with random  $W$ , generate random perturbation  $\delta W$  to find if  $(W + \delta W)$ 's loss is lower.

Strategy #3. Following the Gradient.

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \rightarrow \text{numerical gradient}$$

$$\nabla_w L_i \rightarrow \text{analytic gradient}$$

Practical : Always use analytic gradient, but check implementation with numerical gradient. (gradient check)

Compute the analytic gradient and compare it to the numerical to check the correctness.

Gradient Descent loop { 1. compute gradient

2. update the parameters with learning rate  $\times$  gradient }

Mini-batch gradient descent, compute the gradient over batches of training data.  
size: 32/64/128/256...

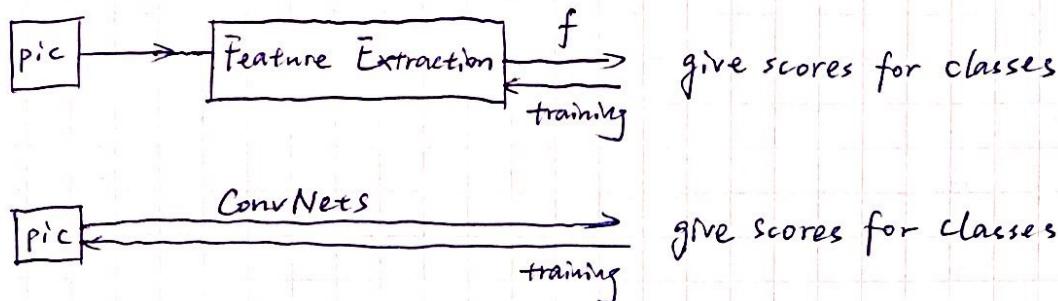
Stochastic gradient descent. One example each time to compute the gradient.

Other gradient descent methods : - gd with momentum.  
- Adam optimizer

Image Features → apply feature transform to make "points" separable by classifiers

- color histogram
- ~~histogram~~
- Histogram of Oriented Gradients (HOG)
- Bag of Words

Image features vs. ConvNets



Rather than writing down the features ahead of time, ConvNets learn the features directly from the data.

Further Reading

< Deep Learning using Linear Support Vector Machines >

### Some Definitions

$$\text{classification} \rightarrow L_i = \sum_{j \neq y_i} \max(0, f_j - f_{y_i} + 1), \quad L_i = -\log \left( \frac{e^{f_{y_i}}}{\sum_j e^{f_j}} \right)$$

Large number of classes → Hierarchical Softmax. 分层 labels 成一棵 tree. 每个 label 都是树的一个 path. 在节点上训练 softmax 作为分支.

Attribute classification →  $y_i$  是一个 binary vector. 每个 example 可能有多个 attributes

$$L_i = \sum_j \max(0, 1 - y_{ij} f_j), \quad L_i = \sum_j y_{ij} \log(\sigma(f_j)) + (1 - y_{ij}) \log(1 - \sigma(f_j))$$

$$\text{Regression} \rightarrow L_i = \|f - y_i\|_2^2, \quad L_i = \|f - y_i\|_1, \quad \frac{\partial L_i}{\partial f_i} = y_{ij} - \sigma(f_j)$$

L2 Loss 是更稳定的 loss (in softmax) 更难被优化. L2 Loss 需要非常脆弱和特定的网络 才能输出精确的正确值. for each input. 而 softmax 不同 每个 score is 很弱值没有那么重要. scores 相互之间较不重要. ② L2 Loss 对 outliers 敏感.

如下. 当面对 regression 任务时, 先考虑是否绝对必须. 因而代之的是更偏向于高散化输出. 然后 perform classification. (for AlexNet loss function 为 softmax 分类)

structured prediction → structured SVM loss (Labels are arbitrary structures. graphs, trees...)

## Lecture 4 Backpropagation and Neural Networks

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

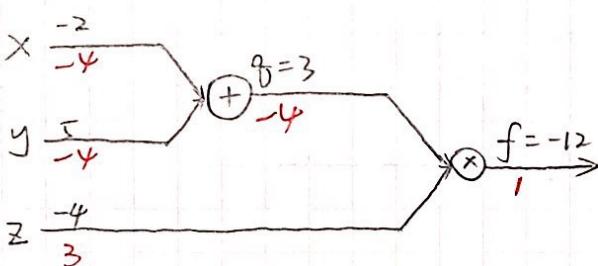
The derivative on each variable tells the sensitivity of the whole expression on its value.

### Chain rule

$$z = f(y) \quad y = g(x)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx} = f'(y)g'(x) = f'(g(x))g'(x)$$

exp:



### Backpropagation

Every gate in a circuit diagram gets some inputs and can right away compute two things :

1. its output value
2. the local gradient of its output with respect to its input value.

During backpropagation, we can get the gradient of output on the inputs of entire circuit by the chain rule.

The gates are relatively arbitrary. Any kind of differentiable function can act as a gate.

## Patterns in backward flow

add gate: gradient distributor.

takes the gradient on its output and distributes it equally to all of its inputs.

max gate: gradient router

distributes the gradient to exactly one of its inputs. (the input that had the highest value during forward pass)

multiply gate: gradient switcher.

## Gradients for vectorized operations

Tip: ① Explicitly write out a minimal vectorized example, derive the gradient on paper and then generalize the pattern to its efficient, vectorized form.

② The gradient with respect to a variable should have the same shape as the variable.

Because ~~the~~ each element of the gradient is quantifying how much that element is contributing to the final output.

## Lecture 5 Convolutional Neural Networks

### Differences with ordinary Neural Networks

- ConvNet architectures make the explicit assumption that the inputs are images, which allow us to encode certain properties into the architecture. (preserve spatial structure).
- forward function more efficient to implement.
- reduce the amount of parameters in the network.

A ConvNets is made up of layers. Every layer has a simple API: It transforms an input 3D volume to an output 3D volume with some differentiable function with or without parameters.

A simple ConvNet for CIFAR-10 classification:

INPUT - CONV - RELU - POOL - FC

### Convolutional Layer

#### Overview and Intuition

- The conv layer's parameters consists of a set of learnable filters.
- Every filter is small spatially (along width and height), but extends through the full depth of the input volume. (e.g.  $5 \times 5 \times 3$ )
- During forward pass, we slide (convolve) each filter across the width and height of the input volume and compute dot products between the entries of the filter and the input at any position.
- Several 2-dimensional activation maps is generated to give the responses of filter at every spatial position.
- The network will learn filters that activate when they see some type of visual features. (edge, blotch, honeycomb or wheel-like patterns)

## □ Local Connectivity.

- Each neuron will only connect to a local region of the input volume.
- The spatial extent of local connectivity is a hyperparameter called receptive field of the neuron (= filter size)
- The extend of the connectivity along the depth axis is always equal to the depth of the input volume.

## □ Spatial Arrangement.

Three hyperparameters control the size of the output volume: depth, stride and zero-padding.

- depth, the number of filters, each learning to look for  $s$ th different in the input depth column (fibre), a set of neurons that <sup>are</sup> looking at the same region of the input.
- stride, how many pixels will be skipped when we move the filters. 1 or 2 is commonly used. Larger stride will produce smaller output volumes spatially.
- zero-padding, pad the input volume with zeros around the border. exactly preserve the spatial size of the input volume.

Input volume size ( $W \times W$ ), the receptive field size ( $F$ ), the stride ( $S$ ), the amount of zero-padding ( $P$ ), number of filters ( $K$ )

$$\text{the size of output} = (W - F + 2P)/S + 1 \quad P = (F-1)/2 \quad S = 1$$

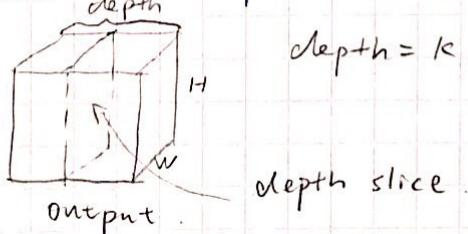
$\downarrow$   
input volume and output have  
the same size

The Conv Layer:

- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires four hyperparameters,  $K, F, S, P$
- produces a volume of size  $W_2 \times H_2 \times D_2$ , where  

$$W_2 = (W_1 - F + 2P)/S + 1, \quad H_2 = (H_1 - F + 2P)/S + 1, \quad D_2 = K.$$
- with parameter sharing, it introduces  $F \cdot F \cdot D_1$  weights per filter, for a total of  $(F \cdot F \cdot D_1)K$  weights and  $K$  biases.

□ Parameter Sharing: the neurons in each depth slice use the same weights and biases



The parameter sharing assumption may not make sense when the input images to a ConvNet have some specific centered structure. (different features should be learned on one side of the image than another.)

One practical example  $\rightarrow$  faces centered in the image.

Instead of parameter sharing  $\rightarrow$  Locally-Connected Layer

$1 \times 1$  convolution  $\quad < \text{Network} \sim \text{In Network} >$

Dilated convolutions  $\quad < \text{Multi-Scale Context Aggregation by Dilated Convolution} >$

$$\text{dilation 0} \quad W[0] \times X[0] + W[1] \times X[1] + W[2] \times X[2]$$

$$\text{dilation 1} \quad W[0] \times X[0] + W[1] \times X[2] + W[2] \times X[4]$$

merge spatial information across the inputs much more aggressively with fewer layers.

## Pooling Layer

reduce the spatial size of the representation to reduce the amount of parameters and computation in the network, and hence to also control overfitting.

Typical form:  $2 \times 2$  filter max pooling, stride = 2

- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires two hyperparameters:  $F, S$
- Produces a volume of size  $W_2 \times H_2 \times D_2$

$$W_2 = (W_1 - F) / S + 1 \quad D_2 = D_1$$

$$H_2 = (H_1 - F) / S + 1$$

- Introduce zero parameters since it computes a fix function
- commonly no zero padding

Getting rid of pooling.  $\quad < \text{Striving for Simplicity: All Convolutional Net} >$

To reduce the size of representation using larger stride in Conv layer instead of pooling.

Discarding pooling layers has been found to be important in training good generative models (variational autoencoders, VAEs, generative adversarial networks, GANs)

## Converting FC layers to CONV layers

### Fully-connected

The only difference between FC and CONV layers is that the neurons in the CONV layer are connected only to a local region in the input, and that many of the neurons in a CONV volume share parameters. 连接共享性 ( $2 \times 2 \leftrightarrow 1 \times 1 \times 1 \times 1$ )  
Both layers compute dot products. 点积共享

- For any CONV layer there is an FC layer that implements the same forward function. The weight matrix would be a large matrix that is mostly zero except for at certain blocks (due to local connectivity) where the weights in many of the blocks are equal (due to parameter sharing)
- Any FC layer can be converted to a CONV layer. For example, an FC layer with  $K=4096$  that is looking at  $7 \times 7 \times 512$  input can be equivalently expressed as a CONV layer with  $F=7$ ,  $P=0$ ,  $S=1$ ,  $K=4096$ . The output will be  $1 \times 1 \times 4096$ . (filter size = input volume)

实际应用中将FC转变为CONV层更加有用。

假设一个卷积神经网络的输入是  $224 \times 224 \times 3$ . 第一层到第五层的激活层大小为  $7 \times 7 \times 512$ . (在 AlexNet 中是这样. 有 5 个池化层, 每次下采一半.  $224/32=7$ ). - 一个 AlexNet 用 3 个尺寸为 4096 的 FC, 最后一个有 1000 个神经元的 FC 用于计算分类评分。  
我们把这 3 个 FC 层中任意一个转化为 CONV 层:

$$FC_1. \quad 7 \times 7 \times 512 \rightarrow 1 \times 1 \times 4096. \quad \text{filter size } F=7 \quad K=4096$$

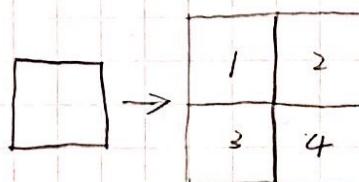
$$FC_2. \quad F=1 \quad K=4096$$

$$FC_3. \quad F=1 \quad K=1000.$$

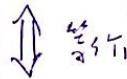
每次这样的变换需要把 FC 的 weight 重塑成 CONV layer 的滤波器。

将 FC 转化为 CONV 层的原理: 当让差值网络(带 FC)在一更更大的输入图片上

滑动(即把一张更大的图的不同区域分别输入到带有 FC 的差值网络. 滑动每个区域的得分), 得到多个输出。这样转化可以在单个转化后的差值网络向前传播中完成。



例如：上  $224 \times 224$  的滑窗以步长 32 在  $384 \times 384$  的图片滑动，把每个滑停的位置都带入有 FC 的 ConvNet，得到  $6 \times 6$  个位置的类别得分。



在  $384 \times 384$  大图上直接经过“把 FC 转换为 Conv layer”的卷积网络向前传播。

$$\begin{aligned} 224 \times 224 &\rightarrow 7 \times 7 \times 512 \quad (\text{5/3 池化}) \xrightarrow{\text{3FC}} 1 \times 1 \times 1000 \\ 384 \times 384 &\rightarrow 12 \times 12 \times 512 \quad (\text{5/3 池化}) \\ &\downarrow \text{3个卷积层} \\ &6 \times 6 \times 1000 \end{aligned}$$

### Summary

→ 相比“滑动”的方法，使用转化过的卷积网络 计算效率高得多。（只有 1 次向前传播，滑动需要计算 36 次，转化后相当于 36 次计算共享计算资源）。

→ 实践中，通常将一张图尺寸变大，然后使用变换后的 ConvNet 对空间不同位置进行并行计算分类得分，再求平均值。

当在大图上以  $\text{2}^{\text{pool}}$  进行滑动时，等同于 在大图上以步长为 1 直接卷积。

将原 FC 转变为 Conv layer.

→ 若使用 步长  $\text{2}^{\text{pool}}$  滑窗，需多次向前传播。  
先在原图用原步长向前传播一次，再分别沿宽、沿高 步长不变，把原图平移步长个像素，把平移后的图分别带入向前传播。（转化后的 ConvNet）

以原步长

$s = 32, 16, 8 \dots$

## ConvNet Architectures

最常见的卷积神经网络：

$$\text{Input} \rightarrow [[\text{conv} \rightarrow \text{RELU}] \times N \rightarrow \text{Pool?}] \times M \rightarrow [\text{FC} \rightarrow \text{RELU}] \times K \rightarrow \text{FC}$$

其中  $\times$  号表示重复次数，pool? 表示一个可选的汇聚层，其中  $N \geq 0$ ，通常  $N \leq 3$ ， $M \geq 0$ ， $K \geq 0$  通常  $K < 3$

常见结构模式：

- Input  $\rightarrow$  FC，实现一个线性分类器， $N = M = K = 0$
- Input  $\rightarrow$  CONV  $\rightarrow$  RELU  $\rightarrow$  FC
- Input  $\rightarrow$  [CONV  $\rightarrow$  RELU  $\rightarrow$  POOL]  $\times 2 \rightarrow$  FC  $\rightarrow$  RELU  $\rightarrow$  FC，每两个汇聚层间有一个卷积层
- Input  $\rightarrow$  [CONV  $\rightarrow$  RELU  $\rightarrow$  CONV  $\rightarrow$  RELU  $\rightarrow$  POOL]  $\times 3 \rightarrow$  [FC  $\rightarrow$  RELU]  $\times 2 \rightarrow$  FC  
每个汇聚层前有两个卷积层，这种思路适用于更复杂更深的网络。因为在执行具有破坏性的 POOL 前，多层的卷积层可以从输入数据中识别更多复杂特征。

## Intuition:

一些小滤波器卷积层的组合比一个大滤波器卷积层要好.

例如，有3个 $3 \times 3$ 的卷积层。第一层 Conv layer 的 Input 都有  $3 \times 3$  的 view，第二层 卷积层的神经元对第一层有  $3 \times 3$  的 view，即对输入有  $5 \times 1$  的 view，第三层的每个神经元 对于输入有  $7 \times 7$  的 view。这样与单独  $7 \times 7$  感受野的卷积层的感受野一样，但有大 filter 优点：① 多个卷积层与非线性激活层交替的结构，要比单一卷积层的结构更能提取出深层的更好特征。

② 假设所有数据有  $C$  个通道（深度为  $C$ ），单独的  $7 \times 7$  卷积层有  $C \times (7 \times 7 \times C) = 49C^2$  个参数，而 3 个  $3 \times 3$  卷积层的组合有  $C(3 \times 3 \times C) \times 3 = 27C^2$  个参数，小滤波器卷积层的组合参数更少。

大 filter 优点：进行反向传播时，若单用多个小滤波器的卷积层，可能会占用更多内存去存放中间层的结果。

Recent Departures：传统的将层按线性排列的方法度到了谷歌 Inception 和轻量级的残差网络（Residual Net）的挑战。

## Layer Sizing Patterns

- Input layer. 应该可以被 2 整除多次，如 CIFAR 的 32, 64, 96 (STL-10), 224, 384, 512。

- 卷积层. 应仅用小滤波器 ( $3 \times 3$ , 最多  $5 \times 5$ )，步长  $S=1$ 。关键是不对输入做 padding。  
当  $F=3$  且  $P=1$  保持输入尺寸，当  $F=5$ ;  $P=2$ 。对于一般的  $F$ ，当  $P=(F-1)/2$  时能保持输入尺寸。若必须使用更大的  $F$ ，如  $7 \times 7$ ，通常只用在第一个而非后续的卷积层上。

- Pooling layer: downsample the spatial dimensions of the input.  
通常用  $2 \times 2$  的感受野 ( $F=2$ ) 做最大值汇聚，步长为 2，会丢掉输入的 75%。

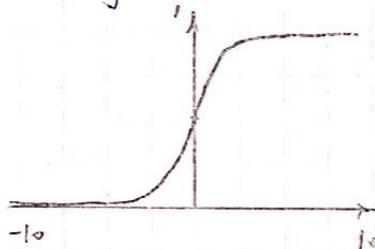
如果使用步长大于 1 的 filter 在卷积层且不做 padding，也能达到 down sampling 的效果，但是需要仔细监督根据通过卷积神经网络的过程，确保所有步长和 filter 尺寸都吻合。  
有时需要对因为内存限制的妥协。实际实践中倾向于在第一个卷积层以及，如 AlexNet 第一个卷积层使用步长为 4， $11 \times 11$  的滤波器。

Case studies : LeNet  
AlexNet  
ZF Net  
GoogLeNet  
VGG Net  
ResNet

## Lecture 6. Training Neural Networks I

### Activation Functions

- Sigmoid



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

缺点: 1. Saturated neurons "kill" the gradients

当  $x$  很大或很小时, 梯度趋近于 0, backprop 中无法得到梯度流的反馈.

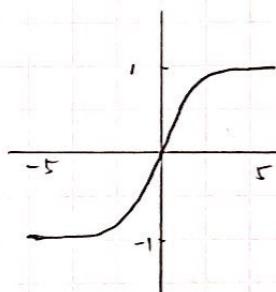
2. Sigmoid outputs are not zero-centered.

Sigmoid 输入到下一层 neuron 的输入恒为正, 则此 neuron weights 的 gradients 应为正或负. 在更新 weights 时会有 zigzagging 跳变现象.

当一个 batch 中 data 的 gradients 加起来对最终 w 的更新可以有不同的正负, 球解了不是 zero-centered 的问题, 但不太方便.

3.  $\exp()$  is a bit compute expensive.

-  $\tanh(x)$

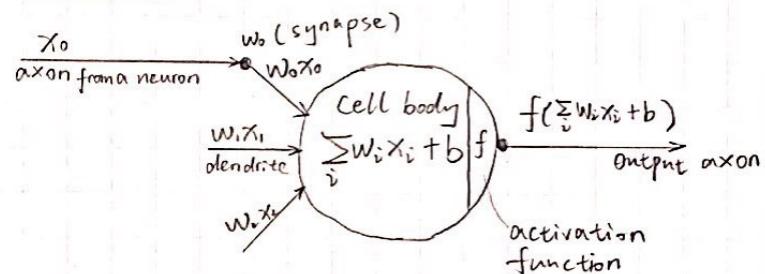


$$\tanh(x) = 2\sigma(2x) - 1$$

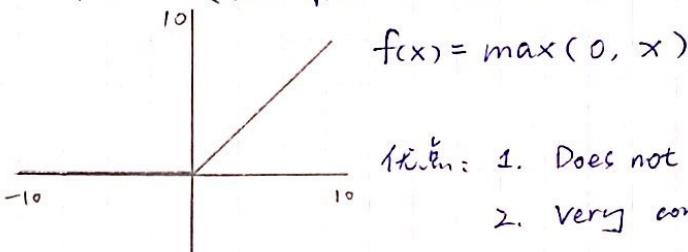
优点:  $\tanh$  是 zero-centered.

缺点: gradients saturation.

实际应用中  $\tanh$  优先于 sigmoid.



## - ReLU (Rectified Linear Unit)



- 优点: 1. Does not saturate in + region  
2. very computationally efficient

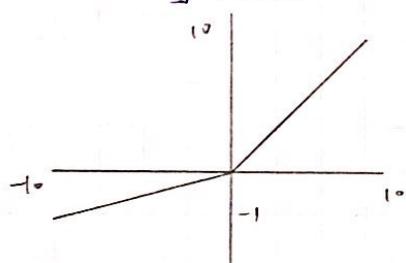
3. Converges much faster than sigmoid/tanh (e.g. 6x) AlexNet

缺点: 1. Not zero-centered output

2. 被置0的data将不会再更新. 可能的原因是 learning rate 过大 (40% of the network can be dead)

缓解方法: Initialize ReLU neurons with slightly positive biases (e.g. 0.01) 增加在初值时神经被激活的可能并得到更高的; proper setting of the Lr.

## - Leaky ReLU

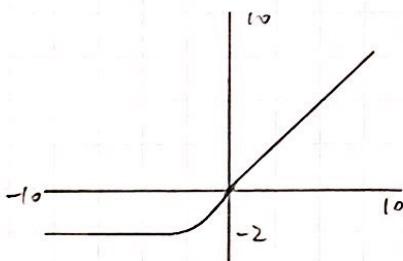


$\alpha$  is a small positive constant.  
 $\alpha$ 也可以由学习得到.

Leaky ReLU 拥有 ReLU 的优点. 并解决了被置0之后不会再更新的问题.

The consistency of the benefit across tasks is presently unclear.

## - ELU (Exponential Linear Units)



- 优点: 1. all benefits of ReLU  
2. closer to zero mean outputs  
3. Negative saturation regime compared with Leaky ReLU, adds some robustness to noise.

缺点: computation requires  $\exp()$

## — Maxout

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

优点: 1. generalizes ReLU and Leaky ReLU

有 ReLU 的优点, 也不会有 dying ReLU 的问题

2. Linear Regime. Does not saturate.

缺点: doubles the number of parameters/neuron.

## Summary

- 用 ReLU, 注意 learning rates 和网络中 "dead" units 的比例。
- 尝试 Leaky ReLU, Maxout 和 ELU, 防止 dead units 过多。
- 尝试 tanh, 但一般不会比 ReLU / Maxout 有效。
- 不要使用 sigmoid。

## Representational Power

Universality theorem: Neural networks with a single hidden layer can be used to approximate any continuous function to any desired precision.

实践中, 通常 3 层的神经网络比 2 层表现好, 4, 5, 6 层等更深的神经网络提升不大。

相反, 深度对于 CNN 非常重要。更深的 CNN 表现更好。可能的原因是图像包含了多种分层的结构信息(如脸由眼、嘴等构成, 它们都具有 edge)。

神经网络大小的选择: 不要因为害怕 overfitting 而选用小的神经网络。(更小的 NN 有更少的 local minima, 但是更容易收敛且 loss 更大, 训练一个小时的 NN, 最终 loss 会有很大的 variance, 即有的 loss 很小, 有的 loss 很大)

应该在计算资源允许的情况下选用更大的神经网络。然后使用 L2 regularization, dropout, input noise 等正则化方式来控制 overfitting。

(更大的 NN, 会找到更多不同的 local minima, 这些结果的 variance 很小, rely less on the luck of random initialization)

## Data Preprocessing

### — Mean subtraction

对每个数据减去其对应的特征的 mean. 具体上可以解释为将 data cloud 移到原点。  
 $X -= np.mean(X, axis=0)$

特别地对于图像，通常减去所有 pixels' mean. 或对 RGB channel 分别  
mean subtraction.  
 $X -= np.mean(X)$

### — Normalization.

归一化数据使其在大约相同的度量里。(范围)

一种方式是将经过 zero-centered 的数据除以标准差，在每个 dimension 上。

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

另一种方式是分别将最大值和最小值对应 1 和 -1.

归一化仅在不同输入的特征有不同 scales 才有意义。对于图像而言，pixels  
的相对 scales 大致相同 (0~255)，通常不必进行归一化。

### — PCA and Whitening

$$\frac{X^T X}{n} = V S^2 V^T$$

$$X -= np.mean(X, axis=0) \quad \# \text{zero-center } X (N \times D)$$

$$\text{cov} = np.dot(X.T, X) / X.shape[0] \quad \# \text{get the data covariance matrix}$$

$$U, S, V = np.linalg.svd(\text{cov}) \quad \# \text{SVD. } U \text{ 为 eigenvectors (因为 } X^T X \text{ 对称半正定)} \\ S \text{ 为 eigenvalues, 从大到小排列}$$

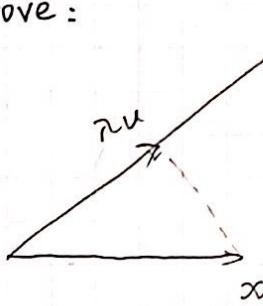
$$X_{\text{rot}} = np.dot(X, U) \quad \# \text{project the original zero-centered}$$

$$X_{\text{rot-reduced}} = np.dot(X, U[:, :100]) \quad \text{data into the eigenbasis.}$$

PCA 将原始数据在其 eigenbasis 上投影而达到降低维度和最大化数据间 variance

(尽可能地保留信息) 的目的。

Prove:



PCA maximizes the projection to receive information as much as possible, when  $\vec{u}$  is the eigenvector to guarantee the new  $x$  has the largest variance.

vector:  $(x - \lambda u)^T u = 0$   
 $\lambda = x^T u$

matrix:  
 $N \begin{bmatrix} D \\ -(x^T u)^T \\ X \end{bmatrix}$

$$(X^T - u\lambda)^T u = 0 \quad \text{当 } u \text{ 是 unit vec. } \lambda \text{ 为投影大小.}$$

$$\lambda = x^T u$$

最大化投影大小

$$\downarrow \quad \begin{aligned} & \text{maximize } (xu)^T xu \\ & = u^T x^T x u \end{aligned}$$

maximize the variance of the projections

$$\text{maximize } \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^{(i)\top} \mathbf{u})^2 = \frac{1}{m} \sum_{i=1}^m \mathbf{u}^\top \mathbf{x}^{(i)} \mathbf{x}^{(i)\top} \mathbf{u} \quad (m: \text{sample number})$$

$$= \mathbf{u}^\top \left( \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)} \mathbf{x}^{(i)\top} \right) \mathbf{u}$$

$$= \mathbf{u}^\top \Sigma \mathbf{u}$$

$$D \begin{bmatrix} N \\ \mathbf{x}^{(1)} & \dots & \mathbf{x}^{(N)} \end{bmatrix} \cdot \begin{bmatrix} D \\ \mathbf{x}^{(1)\top} & \dots & \mathbf{x}^{(N)\top} \end{bmatrix} = \sum_{i=0}^N \mathbf{x}^{(i)} \mathbf{x}^{(i)\top} = \mathbf{X}^\top \mathbf{X}$$

Lagrange:  $\max \mathbf{u}^\top \Sigma \mathbf{u}$   
s.t.  $\mathbf{u}^\top \mathbf{u} = 1$

$$\mathcal{L}(\mathbf{u}, \lambda) = \mathbf{u}^\top \Sigma \mathbf{u} - \lambda (\mathbf{u}^\top \mathbf{u} - 1)$$

$$d\mathcal{L} = (\mathbf{du})^\top \Sigma \mathbf{u} + \mathbf{u}^\top \Sigma \mathbf{du} - \lambda d\mathbf{u}^\top \mathbf{u} - \lambda u^\top du$$

$$\text{tr}(d\mathcal{L}) = 2 \text{tr}(\mathbf{u}^\top \Sigma - \lambda \mathbf{u}^\top) du$$

( $\Sigma$  is symmetric)

$$\frac{\partial \mathcal{L}}{\partial \mathbf{u}} = 2(\Sigma \mathbf{u} - \lambda \mathbf{u}) \xrightarrow{\text{Set}} 0$$

$$\Sigma \mathbf{u} = \lambda \mathbf{u}$$

So, if we wish to find a 1-dimensional subspace to approximate the data,  $\mathbf{u}$  should be the principal eigenvector of  $\Sigma$ .

Generally, if we wish to project our data into a  $k$ -dimensional subspace, we should choose  $\mathbf{u}_1, \dots, \mathbf{u}_k$  to be the top  $k$  eigenvectors of  $\Sigma$ .

$$\mathbf{X} \begin{bmatrix} \mathbf{x}^{(1)\top} \\ \vdots \\ \mathbf{x}^{(N)\top} \end{bmatrix} \cdot \begin{bmatrix} \mathbf{u}_1 & \dots & \mathbf{u}_k \end{bmatrix} = \begin{bmatrix} \mathbf{x}^{(1)\top} \mathbf{u}_1 & \dots & \mathbf{x}^{(1)\top} \mathbf{u}_k \\ \vdots & \vdots & \vdots \\ \mathbf{x}^{(N)\top} \mathbf{u}_1 & \dots & \mathbf{x}^{(N)\top} \mathbf{u}_k \end{bmatrix}_{N \times k}$$

Whitening operation takes the data in the eigenbasis and divides every dimension by the eigenvalue to normalize the scale.

$$\mathbf{X}_{\text{white}} = \mathbf{X}_{\text{rot}} / \underbrace{\text{np.sqrt}(S + 1e-5)}$$

prevent dividing by zero, but will exaggerate the noise in the data

mitigated by stronger smoothing (e.g. increasing  $1e-5$  to be a larger number)

## — Data preprocessing in practice

- 对于图像数据零 centering 比较主要，通常不做 PCA 和 whitening & normalization.
- 任何预处理的统计数字必须只在 training data 里计算，然后应用于 vali/test data.

## Weight Initialization

### — all zero initialization.

错误的做法。若所有  $w$  都从 0 为初始，每个神经元在输入数据的基础上有相同的操作，因为它们将输出相同的值，得到相同的梯度，从而以相同的方式更新。(no source of asymmetry between neurons)

### — small random numbers 不一定会 work strictly better

$W = 0.01 * np.random.randn(D, H) \rightarrow$  从  $N(0, 1)$  中生成随机数  
使用小的随机数 (zero mean) 初始化会使 activations 变得很快趋近于 0，从而在 backproping 梯度也很大，基本无更新，对于更深的网络影响更大。  
用更大的随机数可能有饱和的风险。

### — Calibrating the variances with $1/\sqrt{n}$

使用小随机数 ( $np.random.randn$ ) 的话是这样：从随机初始化的 neuron 的输出的 variance 随输出个数的增大而增大 →

Review:

$$\text{Var}(x) = E[(x - \mu)^2]$$

$$\text{Var}(x) = \text{cov}(x, x)$$

$$\text{Var}(x) = E(x^2) - E(x)^2$$

$\forall x, y$  are independent

$$\begin{aligned} \text{Var}(XY) &= [E(X)]^2 \text{Var}(Y) \\ &\quad + [E(Y)]^2 \text{Var}(X) \\ &\quad + \text{Var}(X)\text{Var}(Y) \\ &= E(X^2)E(Y^2) - [E(X)]^2[E(Y)]^2 \end{aligned}$$

$$\text{Var}(\alpha X) = \alpha^2 \text{Var}(X)$$

$$\begin{aligned} \text{内积 } s &= \sum_{i=0}^n w_i x_i \\ \text{Var}(s) &= \text{Var}\left(\sum_{i=0}^n w_i x_i\right) \quad \text{充分必要条件: } w_1 x_1, w_2 x_2, \dots, w_n x_n \text{ are independent} \\ &= \sum_{i=0}^n \text{Var}(w_i x_i) \\ &= \sum_{i=0}^n [E(w_i)]^2 \text{Var}(x_i) + [E(x_i)]^2 \text{Var}(w_i) + \text{Var}(x_i) \cdot \text{Var}(w_i) \\ &= \sum_{i=0}^n \text{Var}(x_i) \text{Var}(w_i) \quad \text{inputs and weights are zero-mean} \\ &= (n \text{Var}(w)) \text{Var}(x) \quad \text{assume } w_i \text{ and } x_i \text{ are identically distributed.} \end{aligned}$$

前向传播  
如果想保证  $S$  和输入  $X$  有相同的 Variance，因为  $\text{Var}(\alpha X) = \alpha^2 \text{Var}(X)$ ，那么  $\alpha$  的  
weights 需要满足  $\sqrt{\frac{1}{n}}$

$$w = np.random.randn(n) / \text{sqrt}(n)$$

类似初始化  $\rightarrow \text{Var}(w) = 2 / (\text{Nin} + \text{Nout})$  ↗ 同时考虑 backpropagation, 保证  
number of units in the previous layer and the next

Initialization for ReLU  $\rightarrow w = np.random.randn(n) * \text{sqrt}(2.0/n)$

### — Sparse initialization.

Another way to address the uncalibrated variances problem is to set all weight  
matrices to zero, but to break symmetry every neuron is randomly connected to  
a fixed number of neurons below it. (A typical number may be as small as 10)

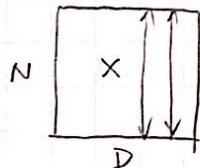
### — Initializing the biases

It is possible and common to initialize the biases to be zero, since the asymmetry  
breaking is provided by the small random numbers in the weights.

### — Batch Normalization

基本思想：“you want zero-mean unit-variance activations? just make them so”

位于 Fully Connected 或 Convolutional Layer 之后但在 nonlinearity 之前。



分别对每个 dimension 计算 empirical mean for variance.

$$\text{再做 } 1/2 - 1e. \quad \hat{x}^{(k)} = \frac{x^{(k)} - E[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

Input: Values of  $x$  over a mini-batch:  $B = \{x_1, \dots, x_m\}$

Parameters to be Learned:  $\gamma, \beta$ .

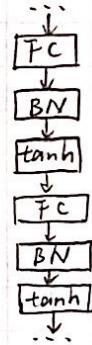
Output:  $\{y_i = BN_{\gamma, \beta}(x_i)\}$

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

$$\sigma_B^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \quad // \text{normalize}$$

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



Intuition:

- 提升了网络中的 gradient flow
- 允许更高的 learning rate
- 降低了对于初始化的依赖
- 其实意义上是 regularization，并稍微减少了对 dropout 的需求  
(因为每层输出都来自于输入  $x$  以及 batch 本身被平均到其它样本的  $\mu$  和  $\sigma^2$ )
- 引入  $\gamma$  和  $\beta$  的作用在于增加了解全局化和包含之即) 的 flexibility.
- Training 时  $\mu_B$  和  $\sigma_B^2$  由当前 batch 算出。  
而 Testing 时  $\mu_B$  和  $\sigma_B^2$  应使用 Training 时  
保存或类似经过处理的值
- BN 中所有操作都半透明，使得 Backprop  
可以有效运行并参与到相应的  $\gamma$  和  $\beta$

## Lecture 7. Training Neural Networks II

### Regularization

#### - L<sub>2</sub> regularization ( $\frac{1}{2}\lambda w^2$ )

L<sub>2</sub>正则化鼓励网络使用所有 input, a little. 而不是使用某些 input, a lot.

在参数更新时, 最终 L<sub>2</sub> 正则会将 weights 趋近于 0. ( $w \leftarrow -\text{lambda} * w$ ).

对  $x$  的改变不敏感, 但对 outlier 敏感. (L<sub>2</sub> 正则对于 NN 的意义可能不那么明确  
(对噪声的鲁棒性)  
有时会在 NN 中选择其它方案)

#### - L<sub>1</sub> regularization ( $\lambda |w|$ )

L<sub>1</sub> 正则鼓励稀疏 (sparse). 使用 L<sub>1</sub> 的网络最终只会使用稀疏的部分输入. 对  
噪声输入几乎不变.

实际上, 如果不关心明确特征的选择, L<sub>2</sub> 正则有望比 L<sub>1</sub> 正则更好.

#### - Elastic net regularization

combine the L<sub>1</sub> and L<sub>2</sub> regularization. ( $\lambda_1 |w| + \lambda_2 w^2$ )

通常 L<sub>1</sub> Lasso (L<sub>1</sub>) 表现好, 同时保留了相似的稀疏性. 鼓励 grouping effect (非常  
相关的特征趋向于一起在或不在模型中).

#### - Max norm constraints

对权重的大小设置上限. ( $\|w\|_2 < c$ ) 需要对  $w$  进行:  $w \leftarrow w \cdot \frac{c}{\|w\|_2}$

当 learning rate 太高时, 使用 max-norm 不会使网络 "explode". 因为上限是被限制的.

#### - Dropout

Dropout 指的是在 forward passing, 按一定概率  $P$ . 随机保留或丢弃 neurons.

Interpretation: 1. Forces the network to have a redundant representation.

Prevents co-adaptation of features.

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

Each binary mask is one model

In test or predicting 不做 dropout. 但要对 hidden layers 和 activations 从  
概率为  $P$  的缩放, 以保证 neurons 的输出和期望值一致.

#### o Inverted Dropout

为了保证 test 或 predict 更有效率. 不在 testing 做和  $P$  的乘法, 而在 training 时, 在  
每一次 dropout 除以  $P$ . 还有个好处在于可以选择部分做 dropout 而不是所有.

但因 dropout 通常训练会花更长 time), 因为每个 step 只更新网络的子部分

- o Dropout 对比 Batch Normalization.

BN中一个神经元可能出现在不同的 mini batches 中。对单个神经元来说，在训练中如何被正则化是有随机性的。Fw. BN在训练时引入了某种程度的随机性或噪声，又在 test 时平均(average)掉，也是一种正则化的方式。

当用 BN 在训练网络时，正则化效果有时是够好而不需要 dropout，但 dropout 的优点在于可调整正则化的强度，通过调整  $P$  来达到。而 BN 没有。

- o Drop Connect.

DropConnect sets a randomly selected of weights with the network to zero.

Each unit thus receives input from a random subset of units in the previous layer.

- o Data Augmentation

通过对有限的数据进行一定的变换生成新数据，防止过拟合。

变换方式包括 (Horizontal Flips, Random crops and scales, Color Jitter, PCA Jittering translation, rotation, stretching, shearing, lens distortions ...)

- o Fractional Max Pooling

随机部分 max pooling.

- o Stochastic Depth.

在训练过程中，随机丢弃部分 layers.

以上包括 Dropout 在内的正则化方法都得益于<sup>training</sup>引入了一定的随机性，而在 test 时。

随机性带来的 noise 会被解析地 (dropout 乘以  $p$ ) 或极简地 (随机丢弃变成多次 forward passes 最后取平均) 忽略，即归一化。

### Summary

通常不需要对 Bias 做正则化以及不需要对每层 layer 做不同强度的正则化。

实践中常见的是单一的全局 L2 正则并使用 cross-validation。常用于 dropout 一起使用。

$P=0.5$  是默认的推荐，可以在 validation tune.

- The recommended preprocessing is to center the data to have mean of zero, and normalize its scale to  $[-1, 1]$  along each feature.
- Initialize the weights by drawing them from gaussian distribution with standard deviation of  $\sqrt{\frac{2}{n}}$ , where  $n$  is the number of inputs to the neuron.
- use L2 regularization and dropout (the inverted version)
- use batch normalization.

## Learning

### Gradient Checks

In theory, performing a gradient check is as simple as comparing the analytic gradient to the numerical gradient.

- Use the centered formula. (numerical gradient)

$$\frac{df(x)}{dx} = \frac{f(x+h) - f(x)}{h} \quad (\text{bad, don't use})$$

实践中  $h \sim 10^{-5}$ , 更多使用 centered difference formula.

$$\frac{df(x)}{dx} = \frac{f(x+h) - f(x-h)}{2h} \quad (\text{use instead})$$

直观来讲, 这个式子评估了 loss function 2 次来检查每一个维度的梯度, 这个近似以更为准确. 使用 Taylor expansion 展开  $f(x+h)$  及  $f(x-h)$  发现第一个式子的误差为  $O(h)$ , 第二个式子误差为  $O(h^2)$ .

prove: 全  $f(x)$  在  $x=a$  处展开:  $f(x) = f(a) + f'(a)(x-a) + f''(a)\frac{(x-a)^2}{2} + f'''(a)\frac{(x-a)^3}{6} + O(x-a)^3$

$$\stackrel{\text{若 } x=a+h \text{ 代入}}{f(a+h) = f(a) + f'(a)h + f''(a)\frac{h^2}{2} + f'''(a)\frac{h^3}{6} + O(h^3)}$$

$$\stackrel{\text{若 } x=a-h \text{ 代入}}{f(a-h) = f(a) - f'(a)h + f''(a)\frac{h^2}{2} - f'''(a)\frac{h^3}{6} - O(h^3)}$$

$$\therefore \frac{f(a+h) - f(a-h)}{2h} = f'(a) + \frac{f''(a)\frac{h^2}{2} + f'''(a)\frac{h^2}{6} + O(h^2)}{h} = f'(a) + O(h)$$

$$\frac{f(a+h) - f(a-h)}{2h} = f'(a) + f''(a)\frac{h^2}{6} + O(h^2) = f'(a) + O(h^2)$$

$$R(x) = \frac{f(\xi)}{(n+1)!} (x-a)^{n+1} \quad (\alpha < \xi < x)$$

- Use relative error for the comparison.

比较 numerical gradient  $f'_n$  和 analytic gradient  $f'_a$  的差来追踪 gradient check is 正确性.

$$\text{relative error} = \frac{|f'_a - f'_n|}{\max(|f'_a|, |f'_n|)}$$

注意当  $f'_a$  和  $f'_n$  都为 0 时, 判断 pass gradient check.

实践中的判断: 相对误差  $> 10^{-2} \rightarrow$  gradient probably wrong

$10^{-2}$  相对误差  $> 10^{-4} \rightarrow$  should make you feel uncomfortable

$10^{-4} >$  相对误差  $\rightarrow$  Okay for objectives with kinks.

too high for no kinks (e.g. tanh and softmax)

$10^{-7} >$  相对误差  $\rightarrow$  happy.

越深的网络, 相对误差会越高 (errors build up on the way)

如果 10 层的网络误差为  $10^{-2}$ , 可能会 OK. 但是 1 层的误差  $10^{-2}$  肯定是不正确的.

aside:  
L2 error check  
 $\|f'_a - f'_n\|_2$   
 $\|f'_a\|_2 + \|f'_n\|_2$

- Use double precision

许多 Gradient check 可能使用双精度。因为经常会在梯度正确时因为仅用了单精度 float 而使 error 大到  $1e-2$ 。

- Stick around active range of floating point

如果梯度大小 ( $1e-10$ ) 可能需要临时 scale loss function up by a constant 来达到 nicer range where floats are more dense. (ideally on the order of 1.0)

- Kinks in the objective.

Kinks 指的是 non-differentiable parts of an objective function. (in ReLU, SVM Loss, Maxout neurons 等等)。

部分的不可导可能带来误差。在 ReLU 上,  $X = 1e-6$ , 此时梯度应该为 0. 但 numerical gradient 可能不为 0. 因为  $f(x+h) \neq f(x-h)$  提供了一个不为 0 的 contribution.

检测 kinks 的跨步的方法：评估  $f(x+h)$  和  $f(x-h)$ 。若  $\max(x, y)$  结果的选择发生了变化，则发生了 kinks cross. 由于 its numerical gradient 不准确。

- Use only few datapoints.

只用很少的数据点来做 gradient check 也可以减轻 kinks 的影响。

2, 3 个 data points 对于一个 batch 的 gradient check 已经足够 (faster and efficient)

- Be careful with the step size  $h$ .

$h$  不宜过小，过小的  $h$  会导致精度问题 ( $1e-6 \sim 1e-4$  可能比较合适)

- Gradcheck during or "characteristic" mode of operation.

梯度检验是在参数空间上指定一点 (randomly) 来进行的。不能立即确定这一点的正确性就是全局的。如在 SVM 上以很小的值初始化 weight。所有 loss score 会近 0. 则所有的 gradient 有相同的特征，这种情况不能 generalize 到更一般的梯度有大有小的情况。

所以安全的做法是给定一个  $\text{short}$  ~~loss~~  $\text{burn-in time}$ . 即等 ~~loss~~ 开始下降时开始做梯度检验。不要在第一次迭代时就做 gradient check.

- Don't let the regularization overwhelm the data.

通常 loss function 是 data loss 和 regularization loss 的和。而梯度的错误是正则化 loss 会覆盖 data loss. 推荐的做法是先关掉正则化，检查 data loss，再单独检查正则项。

单独检查正则项可以去掉 data loss 项，或加大 regularization strength. 来保证肯定它对 gradient check 有大的影响。即 an incorrect implementation would be spotted.

- Remember to turn off dropout / augmentations.

当执行梯度检验时，保证关闭所有网络中的不确定影响 (deterministic effect)。

不然这些因素会带来很多 error. 在做 numerical 的时候。

关闭这些影响的建议是在于无法验证 dropout 的梯度 (dropout 可能无法进行 backprop)

更好的方法是确定一个特定的 random seed. 再计算  $f(x+h)$  和  $f(x-h)$ ，再计算 analytic gradient.

- o Check only few dimensions

实践中，梯度可能有数百万个参数，比较实际的做法是检验某些 dimensions，并确保梯度检验到 ~~每一个~~ dimensions for 每一个单独的参数的一些 dimensions.

## Before Learning: sanity checks Tips/Tricks

- o Look for correct loss at chance performance

确保保在初始化后，得到正确的 loss. 最好先验证没有正则项，只有 data loss.

如 CIFAR-10 的初始 loss 应为  $-\ln(0.1) = 2.202$

Weston-Watkins SVM 初始 loss 为 9 (每个错误的 class vs margin  $\neq 1$ )

- o Second sanity check, 增加正则的 strength 应该增加 loss.

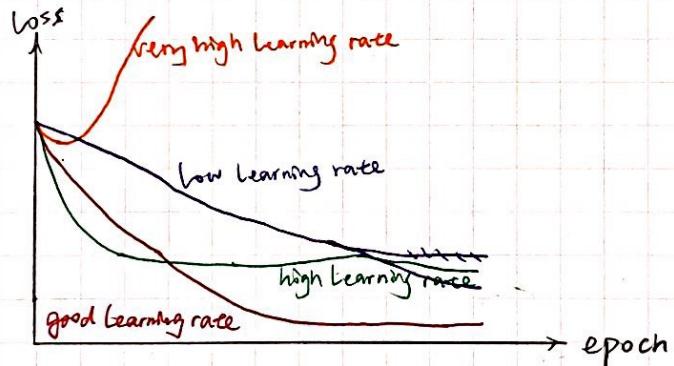
- o Overfit a tiny subset of data.

在训练所有数据之前，先拿出一小部分 (e.g. 20 examples) 训练，确保可以达到 0 cost. (需要先关闭正则项)

## Babysitting the learning process

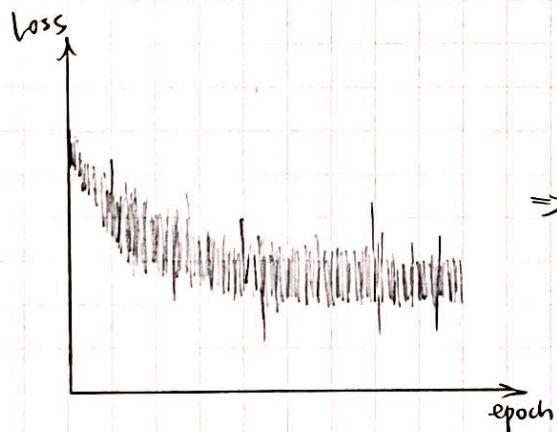
one epoch means that every example has been seen once. 更偏向于追踪 epochs 而不是 iterations. 因为迭代次数取决于 batch size.

### — Loss function.



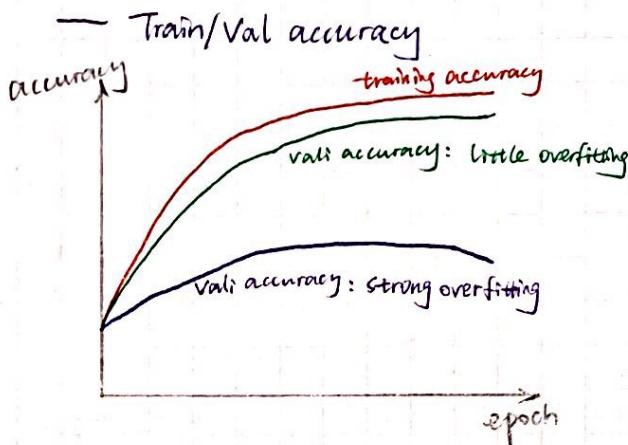
Low learning rates  $\rightarrow$  Improvements will be linear

High learning rates  $\rightarrow$  improvements look more exponential  
~~decay fast, but get stuck at worse value (due to parameters bouncing around chaotically)~~



$\Rightarrow \begin{cases} \text{loss function looks reasonable} \\ \text{maybe small learning rate} \\ \text{batch size small (noisy cost)} \end{cases}$

loss is "wiggly" ~~表示~~ batch size. ~~batch size~~ ~~越大~~, wiggles ~~越少~~.



training and validation accuracy gap 表明 overfitting 的程度

→ 当 validation accuracy 和 training accuracy 差值较大时为严重过拟合，需要增加正则化或得到更多数据。

→ 当 validation accuracy 跟随 training accuracy 相当好时（绿色），说明模型能力尚不够，需要增加模型的复杂度。

### — Ratio of weights: updates.

track weights 变更量和原值量的比值。（比值应在  $10^{-3}$  左右）如果过大，由 learning rate 太大。如果过大由 learning rate 过大。（而用 max, min 或 norm 等）

```
param_scale = np.linalg.norm(W.ravel())
update = -lr * dW # Simple SGD update
update_scale = np.linalg.norm(update.ravel())
W += update
print(update_scale / param_scale).
```

$W.ravel()$  返回矩阵  
 $\approx W.flatten()$   
 flatten 返回 copy  
 ravel 返回 view.

### — Activation/Gradient distributions per layer

诊断 activation 和 gradient。plot histograms for all layers. 可以直观地看到 activation 和 gradient 的分布，判断是否异常。

### — First-layer visualization

plot out first layer features 可以直观地观察出训练的状态。

- plot out its weights 看上去模糊：有 noise，网络不收敛，learning rate 设置不好，很现代化。
- plot its weights 清晰，干净：大部分 features 在 plot 中是好的迹象，表明运行正常。

## Parameter updates.

### — SGD and bells and whistles

#### ◦ Vanilla update

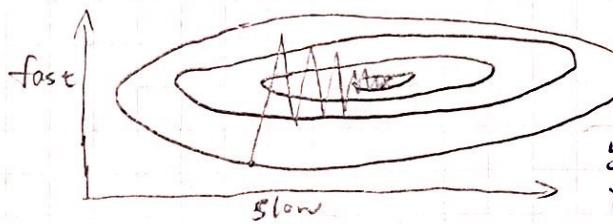
最简单的更新方式是沿着负梯度的方向更新（因为梯度表示着增长的方向，我们往往希望最小化 loss function）

$$x += -\text{learning\_rate} * \text{dx}.$$

## o Momentum update

单纯SGD有一些缺陷:

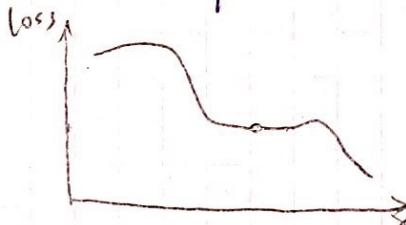
- taco shell problem



singular value  
Hessian matrix 中最大值和最小值之间差距较大.

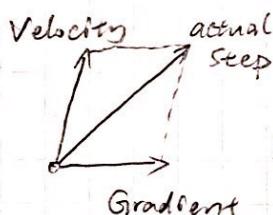
梯度在某个方向上下降很快，在某些方向上下降很慢，于是梯度下降时会出现这样 zig-zag 的情况，导致收敛速度慢。

- saddle point.



当遇到局部最小或鞍点时，梯度会变为0，从而SGD的梯度下降会停止。

在高维中，saddle point 的现象更为普遍，因为一个点上的许多维度的下降方向很可能被彼此抵消。



- noise

SGD的梯度来源于 mini-batch 的 estimation. 梯度更新时会显得 noisy。导致需要花更多的时间来达到 optimum。

解决缺陷的方法：

类比下山的过程，用随机数初始化相当于以0速度在某处放置一个质点。

$$\begin{aligned} \text{方法一} \quad V &= \rho \circ * V - \text{Learning-rate} * dx \\ X &+= V \end{aligned}$$

$\rho$  初始化为0。 $\rho = [0.5, 0.9, 0.95, 0.99]$

可以理解为磨擦力，如果没有  $\rho$ ，则质点永远不会停止运动。

类似于 learning rate 的“退火”机制，典型的设置是  $\rho$  开始为 0.5，在更新过程中逐渐增大到 0.99。

$$\begin{aligned} \text{方法二} \quad V &= \rho \circ * V + dx \\ X &-= \text{Learning-rate} * V \end{aligned}$$

with Momentum update, the parameter vector will build up velocity in any direction that has consistent gradient.

$$\begin{aligned} \text{方法三} \quad V &= \rho \circ * V + (1 - \rho) dx \rightarrow \text{和上面方法等效} \\ X &-= \text{Learning-rate} * V. \quad \text{只不过同时缩小了 } 1 - \rho \\ &\quad (\注意 V 从 0 为起始, 不是 -\rho) \end{aligned}$$

## ✓ Nesterov Momentum

在更新计算梯度时，可以把  $(x + \rho \circ * V)$  当作即将要达到点附近的预估。取  $(x + \rho \circ * V)$  点处的梯度作为更新。



$$\begin{aligned} x\_ahead &= x + \rho \circ * V \\ V &= \rho \circ * V - \text{Learning-rate} * dx\_ahead \\ X &+= V \end{aligned}$$

$$\begin{aligned} V_{t+1} &= \rho V_t - \alpha \nabla f(x_t + \rho V_t) \\ X_{t+1} &= x_t + V_{t+1} \end{aligned}$$

实践中我们不想更新  $x\_ahead$  处的梯度，我们要一直有“ahead version”。想让表达式类似 SGD in momentum。用  $\tilde{x}$  替代  $x\_ahead$ 。

$$\begin{aligned} V\_prev &= V \\ V &= \rho \circ * V - \text{Learning-rate} * dx \\ X &+= -\rho \circ * V - prev + (1 + \rho) * V \end{aligned}$$

$$\begin{aligned} \tilde{x}_t &= x_t + \rho V_t \\ \tilde{x}_{t+1} &= \tilde{x}_t - \rho V_t + (1 + \rho) V_{t+1} \\ \tilde{V}_{t+1} &= \rho V_t - \alpha \nabla f(\tilde{x}_t) \end{aligned}$$

## — Annealing the Learning rate

### ◦ Step decay

减少 learning rate by some factor every few epochs. (by half every 5 epochs, or 0.1 every 20 epochs)

One heuristic → 观察 validation error. 使用 fixed learning rate 训练，当 validation error 不再变小时减小 learning rate (to 0.1倍)

### ◦ Exponential decay.

$$\alpha = \alpha_0 e^{-kt}, \quad \alpha_0 \text{ and } K \text{ are hyperparameters. } t \text{ is iteration number in units of epochs}$$

### ◦ 1/t decay

$$\alpha = \alpha_0 / (1 + kt), \quad \alpha_0, k, t \text{ 同上.}$$

实践中 slightly prefer step decay. 其超参数对  $\alpha_0, k$  更强；如果计算资源允许可尝试使用更慢的 decay 并训练更长时间。

## — Second order methods.

Hessian matrix (a square matrix of second-order partial derivatives of the function)

基于牛顿法来更新参数： $x \leftarrow x - [H(x)]^{-1} \nabla f(x)$

牛顿法是求  $f(x) = 0$  的根。利用牛顿法来找到  $f'(x) = 0$  时的根即找到了最值或鞍点。这样的方法比梯度下降更加有效，原因之一 Hessian matrix 的逆矩阵可以使得在低曲率的方向上得到更激进的 steps，在高曲率的方向上得到小的步长。

另外这种方法没有 Learning rate 这样的超参 (unlike first-order methods).

缺点：

对于深度学习来说，牛顿法有些不切实际。

1. 计算 Hessian matrix 的逆十分 costly. 若有  $10^6$  个参数，Hessian 的大小为  $[10^6 \times 10^6]$ .

→ 有许多 quasi-Newton 方法来计算 Hessian matrix 逆的近似  
如最流行的 L-BFGS，使用梯度变化的趋势来做近似

2. L-BFGS 必须在整个训练集上计算。在 mini-batch 上计算 L-BFGS 目前仍是研究领域

实践中：

在大规模的 DL 中 second-order methods 不太常见。更多的使用的是 SGD 和一些 variants 基于 momentum.

## — Per-parameter adaptive learning rate methods

### o AdaGrad

```


$$dx = \text{compute\_gradient}(x)$$


$$\text{grad\_squared}_t = dx * dx \quad \# dx \text{ is a vector}$$


$$x \leftarrow \text{learning\_rate} * dx / (\text{np.sqrt}(\text{grad\_squared}) + \text{eps})$$


```

$\text{grad\_squared}$  和  $dx$  的 size 相同并且 keep track of 每个参数的梯度的平方和。

然后用  $\text{grad\_squared}$  来 normalize 参数更新项 (element-wise)，对高梯度的参数更新有在 learning rate 的削减作用，对低梯度的参数更新有加快的作用。

$\text{eps}$  在  $(1e-4 \sim 1e-8)$  用来防止除以 0。

(时间足够长，更新步长  $\rightarrow 0$ )

AdaGrad 的缺点在于，对于 DL，单次的字典通常被证明是过于激进和过早停止学习。

### o RMS prop.

Exponentially weighted average.

```


$$dx = \text{compute\_gradient}(x)$$


$$\text{grad\_squared} = \text{decay} * \text{grad\_squared} + (1 - \text{decay}) * dx * dx$$


$$x \leftarrow \text{learning\_rate} * dx / (\text{np.sqrt}(\text{grad\_squared}) + \text{eps})$$


```

RMS prop 引入了 AdaGrad，减少了 缓慢地降低的学习速率，但其不会单调降低。

$\text{decay}$  是一个超参  $[0.9, 0.99, 0.999]$  并且仔细选择了基于梯度大小调整学习速率。

### ✓ Adam (Adaptive Moment Estimation)

```


$$m = \text{beta1} * m + (1 - \text{beta1}) * dx$$


$$v = \text{beta2} * v + (1 - \text{beta2}) * (dx * dx)$$


$$x \leftarrow -\text{Learning\_rate} * m / (\text{np.sqrt}(v) + \text{eps})$$


```

recommendation:

$\text{esp} = 1e-8$

$\text{beta1} = 0.9$

$\text{beta2} = 0.999$

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

Adam 和 RMS prop 很像，只是多了对  $m$  的更新方式。

Adam 是现在推荐的算法。经常比 RMS prop 表现好一些。额外的，尝试一下 SGD + Nesterov Momentum。

完整的 Adam 还包含一个 bias correction 机制。刚开始时的  $m, v$  刚开始被初始化，都为 0。通过一个 bias correction 算法  $m$  和  $v$  都从 0 初始化为 0，使  $m$  和  $v$  更快

"warmup"

```


$$m = \text{beta1} * m + (1 - \text{beta1}) * dx \rightarrow \text{Momentum}$$


$$v = \text{beta2} * v + (1 - \text{beta2}) * (dx * dx)$$


$$mt = m / (1 - \text{beta1}^{** t})$$


$$vt = v / (1 - \text{beta2}^{** t}) \quad \} \text{Bias correction}$$


$$x \leftarrow -\text{Learning\_rate} * mt / (\text{np.sqrt}(vt) + \text{eps})$$


```

AdaGrad/RMSprop

$t \leftarrow 1$  (每 4 iteration + 1)

## Hyperparameter optimization

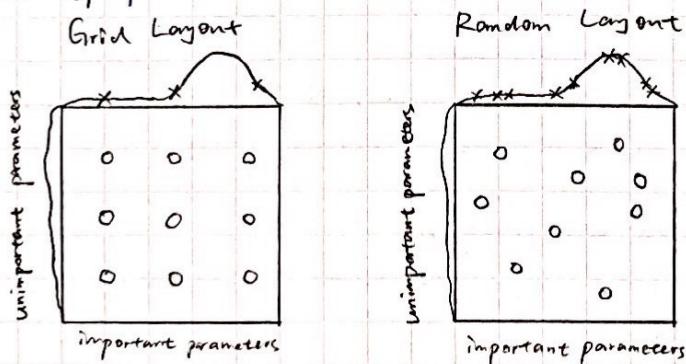
The most common hyperparameter in context of Neural Networks:

- the initial learning rate
- Learning rate decay schedule (such as the decay constant)
- regularization strength (L2 penalty, dropout strength).

- Implementation worker in master to guide the training worker, continuously samples random hyperparameters and performs the optimization. 在训练过程中，worker 持续地从 epoch 的 validation set 中挑选 model checkpoint (该 epoch 的 loss over time) 写进文件。 master, launches or kills workers across a computing cluster, and may additionally inspect the checkpoint written by workers and plot the statistics.

- Prefer one validation fold to cross-validation.  
大多数情况下一个 validation set 在网上可以找到，不需要 cross-validation.
- Hyperparameter ranges  
Search for hyperparameters on log scale. e.g.  $\text{learning\_rate} = 10^{** \text{uniform}(-6, 1)}$   
因为对于 training dynamics, learning rate 起到的是 multiplicative effects. (regularization strength 也相同)  
一些其它参数而在 uniform original scale 上搜索, 如  $\text{dropout} = \text{uniform}(0, 1)$

- Prefer random search to grid search.  
随机搜索 hyperparameter vs 网格搜索更有效，也更容易执行



- Careful with best values on border.  
当得到起始的边界后，需要确认结果是否在边界的边界上，若在边界则很有可能错过最优的超参数。

## — Stage your search from coarse to fine

阶梯式地搜索超参数。一步一步缩小搜索范围（先以粗略范围在 1 epoch 或更少的迭代次数里，因为很多超参数会导致根本不可行或很快 explode；再缩小范围在 5 个 epochs 中搜索；最终在更长的时间里做更细的搜索）。

## — Bayesian Hyperparameter Optimization.

是在更有效地探索参数空间的一系列算法。核心思想是合理地平衡 exploration 和 exploitation trade-off. when querying the performance at different hyperpara.

一些有名的 libraries: Spearmint, SMAC, Hyperopt.

但在实践中，对于 ConvNets 相对困难去仔细地 random search.

## Evaluation

### Model Ensembles

实践中，Ensemble 是一种简单的方法以提高 a few percent performance in its way. (通过许多独立模型，取 testing average 之后的结果)

A few approaches:

- Same model, different initializations.

用 cross-validation 来确定最佳的起参数，再同这组起参数来训练不同的模型。 (The danger: the only variety is due to initialization).

- Top models discovered during cross validation.

但用交叉验证来确定最佳起参数，选取 top few (the top 10) models 来形成 ensemble. 这种方式提高了 variety 但引入了 suboptimal models. 通常并不实施，因为不需要额外的部署。

- Different checkpoints of a single model.

如果不够昂贵，some people have had limited success in taking different checkpoints of a single network over time and using those to form an ensemble.

- Running average of parameters during training. (exponentially decaying sum)

在内存中保存另一组参数，并在最后几个 iterations 加权平均这些参数。最终使用“smoothed”参数。 (almost always achieves better performance)

Intuition: the objective is bowl-shaped and network is jumping around, so the average has a higher chance of being somewhere nearer the mode.

Ensemble in itself is not necessarily cheap prediction in my mind would be costly.

Graeff Hinton's idea: distill a good ensemble back to a single model by incorporating the ensemble log likelihoods into a modified objective.

## Summary

- Fit implementation 使用小 batch size Gradient check 并用 vs 7 例片.
- Sanity check, 确保初值 loss 是合理的. 并且可以很快地将损失降到 100% 以内 (过快)
- Train 1 epoch, 监视 loss, -training/validation accuracy, fancier  $\rightarrow$  updates 每 5 parameter value is 1e-3 (should be  $\sim 1e-3$ ), 处理 ConvNets 的 first-layer weights.
- 2 个推荐的 updates: SGD + Nesterov 或 Adam.
- 训练时学习率 learning rate. wide hyperpara ranges, train only for 1-5 epochs
- 使用随机搜索寻找起参数. 并 Stage the search from coarse  $\xrightarrow{\uparrow}$  to fine  $\downarrow$ .
- Form model ensemble for extra performance. narrower ranges, many more epochs

## Lecture 8 / Hardware and Software

### Deep Learning / Hardware

CPU. fewer cores, but each core is much faster and much more capable.  
great at sequential tasks.

GPU. More cores, but each core is much slower and "dumber"  
great for parallel tasks.

典型的最适合 GPU 的方法  $\rightarrow$  矩阵乘法.

Programming GPUs :

- CUDA (NVIDIA only)  
C-like code that runs directly on the GPU  
Optimized APIs: cuBLAS, cuFFT, cuDNN, etc.
- OpenCL  
similar to CUDA, but runs on anything.  
usually slower on NVIDIA hardware
- HIP.  
convert CUDA code to sth can run on AMD GPUs.

优化 CPU/GPU communication 需要注意不要被从内存中读取数据的速度限制.

- Solutions:
- Read all data into RAM.
  - Use SSD instead of HDD.
  - Use multiple CPU threads to prefetch data.

## Deep Learning Software

### — PyTorch.

- Tensor, like a numpy array, but can run on GPU
- Autograd, package for building computational graphs out of Tensors, and automatically computing gradients
- Module, a neural network layer; may store state or learnable weights.

### Dynamic Computation Graphs.

Building the graph and computing the graph happen at the same time.

### — TensorFlow.

- First define computational graph
- Then run the graph many times.

### Static Computation Graphs.

### — Static v.s. Dynamic.

#### Static vs. Optimization in Serialization.

- static graphs framework can optimize the graph before it runs.
- Once the graph is built, can serialize it and run it without code that built the graph.

#### Dynamic vs. Conditional for Loops.

- Pytorch 在于动态实现条件语句和循环。
- TensorFlow 将所有部分通过计算图实现，条件判断和循环动态实现。
- Dynamic Graph Applications: Recurrent networks, Recursive networks, Modular networks.

PyTorch Dynamic Graphs Static: ONNX, Caffe2	TensorFlow Static Graphs Dynamic: Eager
---	---

#### PyTorch,

- clean API,
- dynamic graphs make it very easy to develop and debug.
- Can build model in PyTorch then export to Caffe2 with ONNX for production/mobile.

#### Tensorflow,

- safe for most projects.
- Not perfect but has huge community, wide usage
- Same framework for research and production.
- Only choice to run on TPUs.

## Lecture 9 CNN Architectures

### AlexNet

- 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 :  $18.2\% \rightarrow 15.4\%$

### 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)

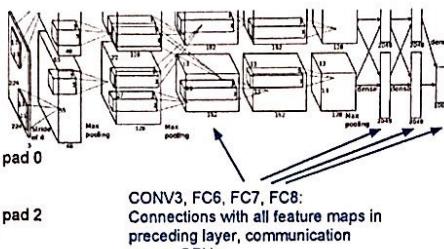


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 19 May 1, 2018

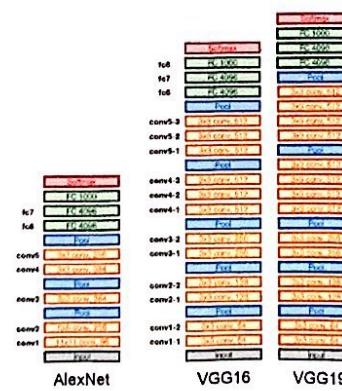
### VGGNet

### Case Study: VGGNet

[Simonyan and Zisserman, 2014]

#### Details:

- ILSVRC'14 2nd in classification, 1st in localization
- Similar training procedure as Krizhevsky 2012
- No Local Response Normalisation (LRN)
- Use VGG16 or VGG19 (VGG19 only slightly better, more memory)
- Use ensembles for best results
- FC7 features generalize well to other tasks



- Smaller filters, only  $3 \times 3$  conv stride 1, pad 1.  $2 \times 2$  MaxPool, stride 2.
- Deeper networks.

但 1 层的 filter 在第 1 层时：

$3 \times 3$  conv (stride 1)  $\Rightarrow$  第 1 层的 effective receptive field 是  $7 \times 7$  的。

- ① deeper, more non-linearities
- ② fewer parameters.

$$3 \times (3 \times 3 \times C) \times C < (7 \times 7 \times C) \times C$$

C for channels per longer.

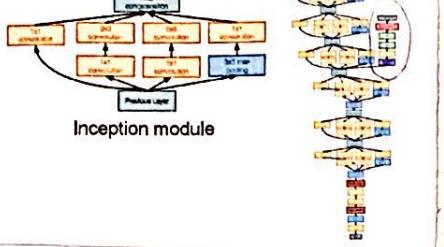
## GoogLeNet

### Case Study: GoogLeNet

[Szegedy et al., 2014]

Deeper networks, with computational efficiency

- 22 layers
- Efficient "Inception" module
- No FC layers
- Only 5 million parameters!
- 12x less than AlexNet
- ILSVRC'14 classification winner (6.7% top 5 error)



Auxiliary classification outputs

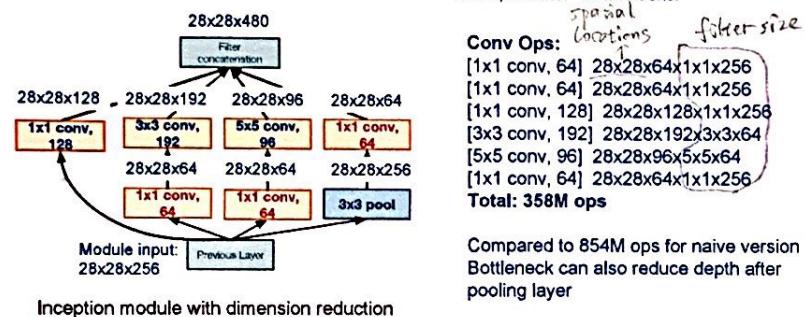
Inception module 是一个 "network within a network" 的思想。

将不同感受野的 filter 并行处理，包含  $1 \times 1$ ,  $3 \times 3$ ,  $5 \times 5$  的 convolution 和  $3 \times 3$  的 Pooling.

最后再把输出连接 (concatenate) depth-wise.

### Case Study: GoogLeNet

[Szegedy et al., 2014]



Using same parallel layers as naive example, and adding "1x1 conv, 64 filter" bottlenecks:

Conv Ops: [spatial locations] filter size  
 [1x1 conv, 64] 28x28x64 1x1x256  
 [1x1 conv, 64] 28x28x64 1x1x256  
 [1x1 conv, 128] 28x28x128 1x1x256  
 [3x3 conv, 192] 28x28x192 3x3x64  
 [5x5 conv, 96] 28x28x96 5x5x64  
 [1x1 conv, 64] 28x28x64 1x1x256  
 Total: 358M ops

Compared to 854M ops for naive version  
 Bottleneck can also reduce depth after pooling layer

单纯的多个 conv layer 并行会带来计算复杂度过大的问题。

解决方法是加入 "bottlenecks".  
 ↪ 通过  $1 \times 1$  的 conv 缩小 depth.

GoogleNet 还有 2 个额外的辅助分类输出。  
 因为当网络很深时，backprop 时的一些梯度信息会被变小，并在靠近输入的时候丢失。  
 辅助分类输出提供了额外的梯度信号。

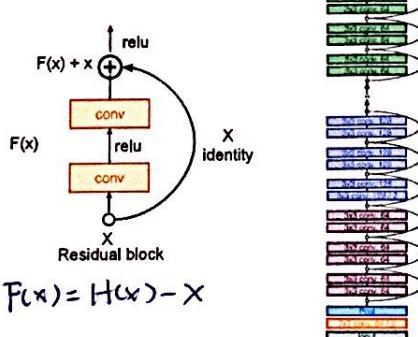
## ResNet

### Case Study: ResNet

[He et al., 2015]

Very deep networks using residual connections

- 152-layer model for ImageNet
- ILSVRC'15 classification winner (3.57% top 5 error)
- Swept all classification and detection competitions in ILSVRC'15 and COCO'15!



问题点：更深一直在 "plain" CNN 上堆叠的层。更深的网络表现会更差，并且不是由 overfitting 引起。  
 ↪ training and test error 都更差。

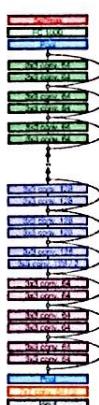
原因：hypothesis → 更深的网络更难被 optimized.

理论上，更深的网络应该至少和浅的网络表现一样好 (深网络 = 浅网络 + identity mapping layers)

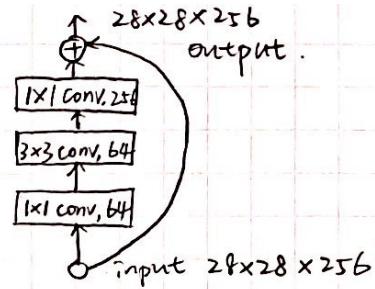
Solution: 使网络来 fit residual mapping 而不是直接 fit 一个 desired underlying mapping.  
 fit residual mapping, 例如只想学习到需要在输入 x 上添加或减去的部分。 $(H(x) - x)$

完整的结构：— stack residual blocks and each block has two  $3 \times 3$  conv layers.

- Periodically, double # of filters and downsample spatially using stride 2. ( $1/2$  in each dimension)
- Additional conv layer at beginning
- Only one "FC 1000" layer at the end.



对于更深的网络 (ResNet 50+), 使用 "bottleneck" 来提高效率 (similar to GooleNet)



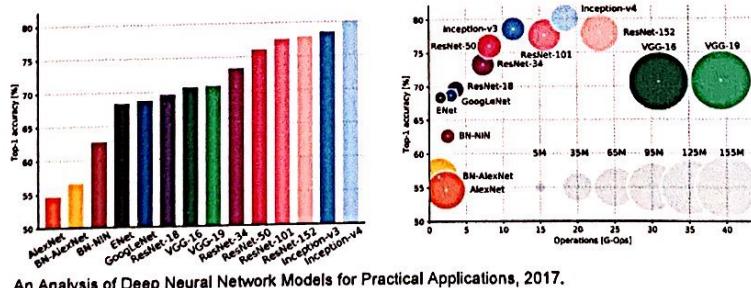
training ResNet in practice:

- Batch Normalization after every CONV layer
- Xavier 2/ initialization from He et al.
- SGD + Momentum 0.9.
- Learning rate 0.1, divided by 10 when validation error plateaus.
- Mini-batch size 256
- Weight decay of  $1e-5$
- No dropout used.

## Summary

- VGG, GooleNet, ResNet all in wide use, available in model zoos.
- ResNet current best default, also consider SENet when available.
- Trend towards extremely deep networks.
- Significant research centers around design of layer/skip connections and improving gradient flow.
- Efforts to investigate necessity of depth vs. width and residual connections.
- Even more recent trend towards meta-learning.

Comparing complexity...



An Analysis of Deep Neural Network Models for Practical Applications, 2017.

Figures copyright Alfredo Cantzani, Adam Pearce, Eugenio Culurciello. 2017. Reproduced with permission.

Inception-v4 : ResNet + Inception

VGG19 : highest memory, most operations

GooleNet : most efficient

AlexNet : smaller compute, memory heavy, lower accuracy

ResNet : moderate efficiency depending on model, highest accuracy.

## Lecture 10 Recurrent Neural Networks

Network input and output with application.

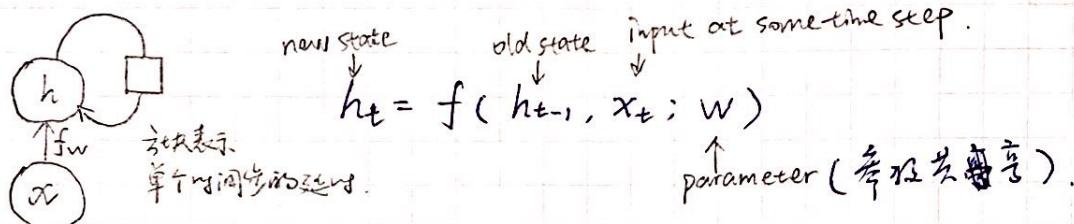
One to one : vanilla Neural Networks  
 (input)      (output)

One to many : Image Captioning. (image  $\rightarrow$  sequence of words).

many to one : Sentiment classification (sequence of words  $\rightarrow$  sentiment)

many to many : Machine Translation ; Video classification on frame level.

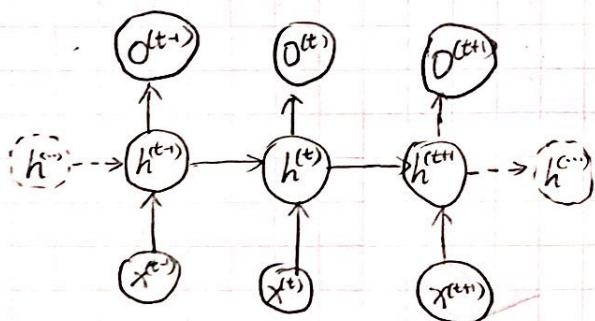
### Recurrent Neural Networks.



隐藏状态  $h_t$  作为过去序列 (直到  $t$ ) 与任务相关的有损摘要 (lossy summary).

有损是因为网络映射任意长度的序列 ( $x_1, x_2, x_3, \dots, x_t$ ) 到一固定长度的向量  $h_t$ .

根据不同的训练准则，摘要可能选择性精确保留过去序列的某些信息.

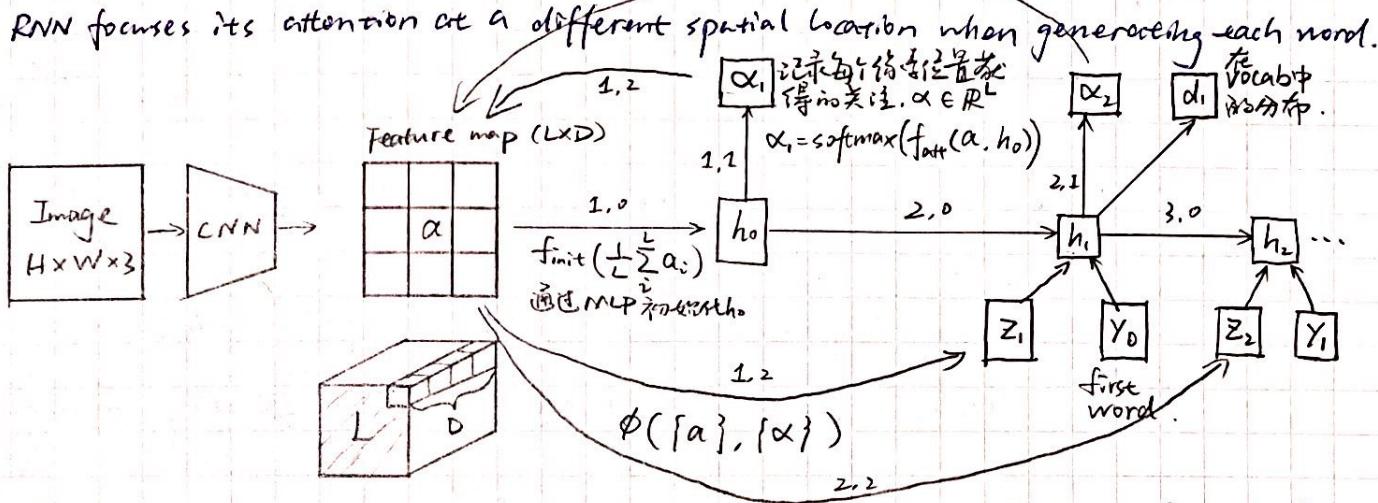


$$\begin{cases} h_t = \tanh(W_h h_{t-1} + W_{xh} x_t) \\ o_t = W_o h_t. \end{cases}$$

循环神经网络的设计模式包括：

- ① 每个时间步都有输出，并且在隐藏单元之间有循环连接。(如左图)
- ② 每个时间步都无输出，只有当前时刻输出到下一个时刻的隐藏单元之间有循环连接。
- ③ 选取整个序列后产生一个输出，且隐藏单元之间存在循环连接。

## Image Captioning with Attention



对于每个 location  $i, i \in L$ ,  $\alpha_i$  表示为用于产生下一个词, 模型应该聚焦位置  $i$  的概率。

$$\ell_{ti} = f_{\text{att}}(\alpha_i, h_{t-1})$$

$$\alpha_{ti} = \frac{\exp(\ell_{ti})}{\sum_{k=1}^L \exp(\ell_{tk})}$$

$f_{\text{att}}$  为 "attention model" (一个多层次感知器 MLP)。

$\hat{z}_t$  为 context vector,  $\hat{z}_t = \phi([\alpha_i], [\alpha_i])$ .

Intuition,

由 CNN 产生的 feature map 和隐藏单元  $h$  (前一状态的) 产生下一个单词所应该关注的区域信息  $\alpha$ . 然后  $\alpha$  再和 feature map 产生语境向量词一起作为下一个隐藏单元的输入.

输出下一个词的  $\alpha$  和  $d$  (即 vocabulary 中的分布).

## Long Short Term Memory (LSTM)

考虑vanilla RNN的梯度流，

$$\nabla_{h^{(t)}} L = W^T \text{diag}(1 - (h^{(t+1)})^2) \cdot \nabla_{h^{(t+1)}} L + V^T (\nabla_{o^{(t)}} L)$$

在计算  $\nabla_{h^{(t)}} L$  时，会包含后一步  $h^{(t+1)}$  的梯度，意味着越初的  $h^{(t)}$  的梯度会和越来越多的  $W$  相乘，导致：

对  $W$  而言 { 如果其最大的奇异值  $> 1$ , 则会梯度爆炸 (exploding gradients) }

{ 如果其最小的奇异值  $< 1$ , 则会梯度消失 (vanishing gradients) }

奇异值为  $W^T W$  或  $WW^T$  的特征值的开方。特征值的乘积为矩阵的体积值，其表示经过该矩阵线性变换的体积缩放因子 (volume scaling factor)

→ 当梯度爆炸时可以通过限制 norm 为 1. 来防止：

```
grad_norm = np.sum(grad * grad)
if grad_norm > threshold:
    grad *= (threshold / grad_norm)
```

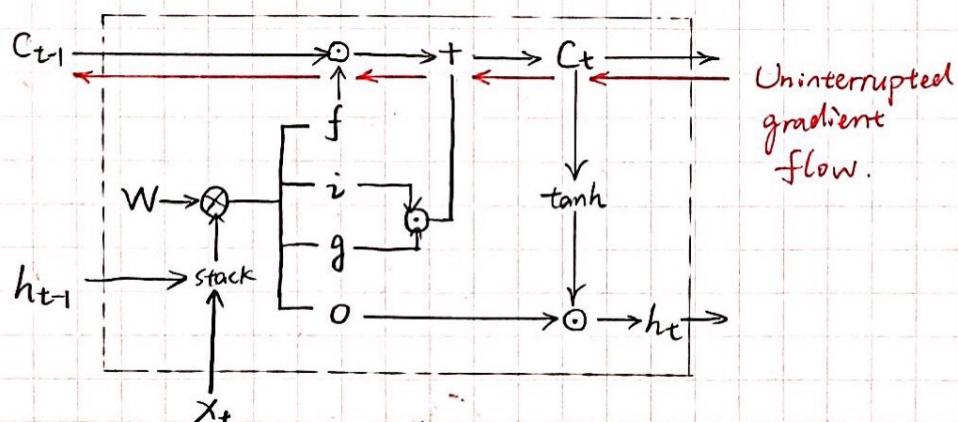
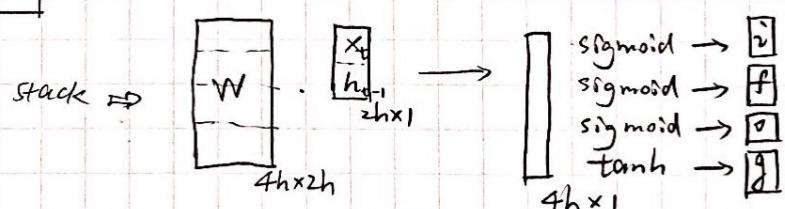
→ 当梯度消失时需要改变 RNN 的架构：LSTM.

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h^{t-1} \\ x_t \end{pmatrix}$$

$$C_t = f \odot C_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(C_t)$$

i. Input gate, whether to write to cell  
f. Forget gate, whether to erase cell  
o. Output gate, how much to reveal cell  
g. Grace gate, how much to write cell



## Summary

- Common to use LSTM or GRU, their additive interactions improve gradient flow
- Backward flow of gradients in RNN can explode or vanish  
Exploding is controlled with gradient clipping.  
Vanishing is controlled with additive interactions (LSTM)

$$\text{GRU: } r_t = \sigma(W_{xr}x_t + W_{hr}h_{t-1} + b_r) \rightarrow \text{复位门}$$

$$z_t = \sigma(W_{xz}x_t + W_{hz}h_{t-1} + b_z) \rightarrow \text{遗忘门}$$

$$\tilde{h}_t = \tanh(W_{xh}x_t + W_{hh}(r_t \odot h_{t-1}) + b_h)$$

$$h_t = z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t$$

复位门像条件语义累积器一样可以线性门控任意维度。

复位门控制当前状态中哪些部分用于计算下一目标状态。在过去状态和未来状态之间引入了附加的非线性效应。

## Lecture 11 Detection and Segmentation

几个CV task的分类：

1. No objects, just pixels: Semantic Segmentation. [语义分割]
2. Single object: Classification + Localization (目标分类并附上 bounding box)
3. Multiple objects: Object Detection (多个物体类似用框标出，并分类)  
Instance Segmentation (在OD基础上 segment 出个体) .
4. 3D object detection (object categories + 3D bounding boxes)

### Semantic Segmentation

- Label each pixel in the image with category label.  
Don't differentiate instances, only care about pixels. (任务描述)
- 一个比较直观的想法是，用一个 sliding window 在原图上滑动得到许多个 extract patch，之后将所有 patch 放入一个 CNN 中得到所有 patch 的一个唯一的 label。  
问题是这样做非常低效，没有利用上有密集 patches 之间共享的 features.
- Solution.
 

输入一个具有连接 Conv Layer 的网络识别每一个 pixel, at once.  
同时对原图做卷积可能会太 expensive. 因此  $\rightarrow$  % down sampling, 再在网络尾部 up sampling 回原大小。

  - Down sampling, pooling, strided convolution.

◦ Upsampling :

◦ Unpooling,

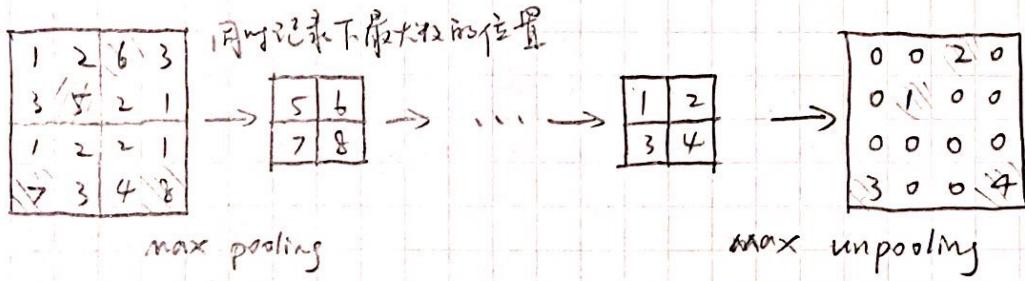
$\begin{array}{ c c } \hline 1 & 2 \\ \hline 3 & 4 \\ \hline \end{array}$	$\rightarrow$	$\begin{array}{ c c } \hline 1 & 2 & 2 \\ \hline 1 & 1 & 2 & 2 \\ \hline 3 & 3 & 4 & 4 \\ \hline 3 & 3 & 4 & 4 \\ \hline \end{array}$
---	---------------	---

Nearest Neighbor

$\begin{array}{ c c } \hline 1 & 2 \\ \hline 3 & 4 \\ \hline \end{array}$	$\rightarrow$	$\begin{array}{ c c } \hline 1 & 0 & 2 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline 3 & 0 & 4 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline \end{array}$
---	---------------	---

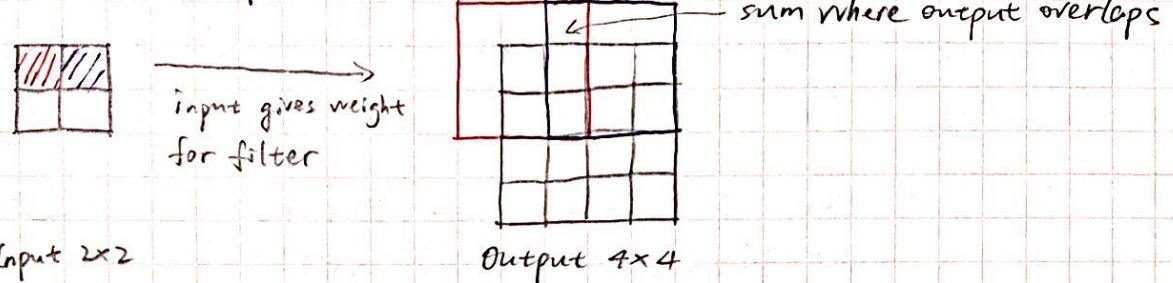
Bed of Nails.

△ Max unpooling.



△ Transpose Convolution,

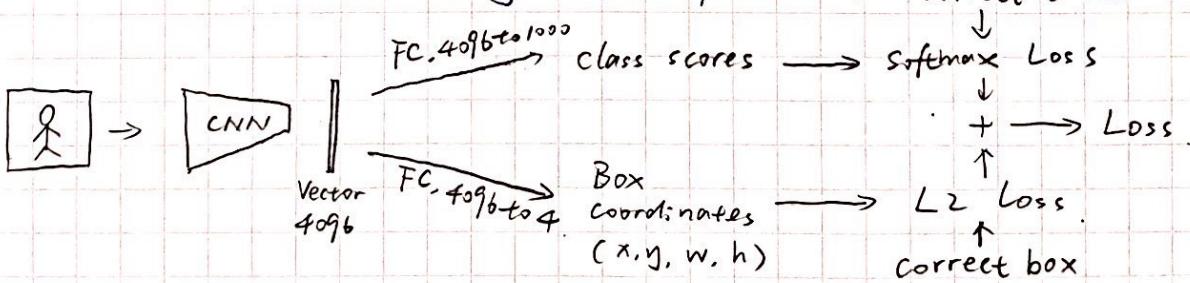
3x3 transpose convolution, stride 2 pad 1



在数学上, transpose convolution is forward pass 而普通的反向传播的 backward pass 是等价的。

## 2D Object Detection

对单个物体的 classification 和 Localization. 输入图像 -> CNN, 然后是 FC 层  
不仅给出 class scores, 还给出 bounding box 的四个坐标信息. correct label



这叫作 localization 也被称为 regression 问题。

但当遇到多个物体时, localization in classification 会有多组输出, 而且不同的图片物体数量不同, 网络的输出个数也不同. (OD as regression).

另一种思路, 在图上任意取 crops 放入 CNN 而得到这些 crops 的分类结果。

这个方法也有问题, 因为 crops 实在太多, 计算上很昂贵。 (OD as classification)

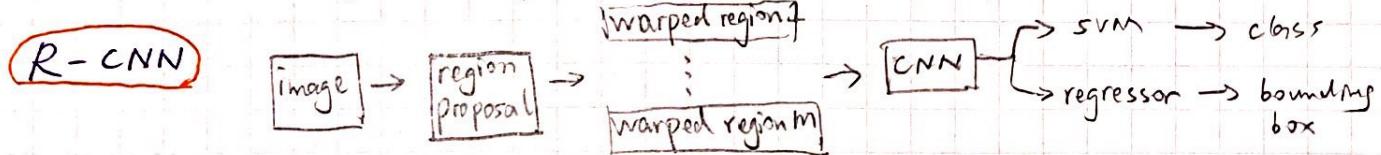
跟着 OP 做 classification 的思路，可以这样减少 crops 的数量：

"Selective Search" 算法输出大约 2000 个 region proposals 在 CPU 上运行时间内。

Step 1, Generate initial sub-segmentation. 通过粗略的分割方法保证每个 region 都属于一个物体

Step 2, Recursively combine similar regions into larger ones.

Step 3, use the generated regions to produce candidate object locations.



step 1, 在原图中选取 Regions of Interest (RoI). 通过一些 这部分是 Selective search

step 2, 将 RoI 转形成 规范化的形状

step 3, 将每个变形后的 RoI 输入 CNN. 输出分类结果以及 bounding box.

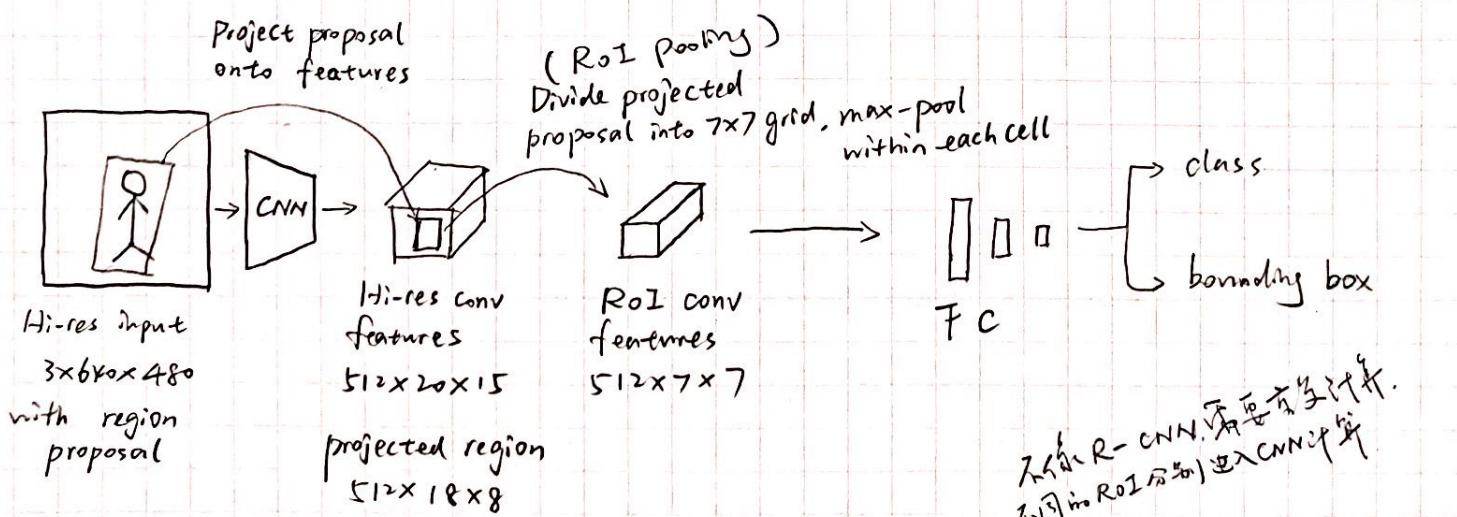
对分类结果做 classification 优化, 对 bounding box 做 linear regression.

- problems :
1. ROI 是相对于训练来说是固定的, 通过单抽选取方法, 而不是通过学习得到。
  2. 训练十分缓慢并且所占的存储空间也很大
  3. 测试时间也很慢 (47s per image with VGG16)  
(Inference) 需要在 2.3s 完成一次预测, 而生成 RoIs 的时间就有 2s

## Fast R-CNN

与 R-CNN 不同的是, Fast R-CNN 将整张图片为 CNN 的输入. 得到一个 feature map.

在 Feature map 上映射 ROI, 将不同的 crops 通过 FC 输出 分类结果及 bounding box.



为什么 Fast R-CNN 快? 因为当选择 ROI 时, Fast R-CNN 对于所有 ROI 只有单次计算。

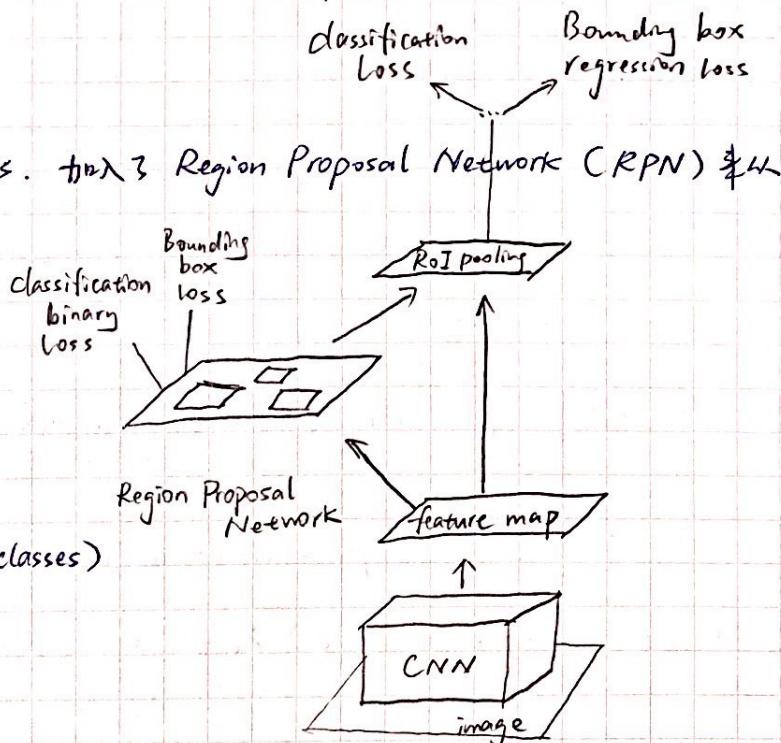
但是 Fast R-CNN 的速度还是被 ROI 的选取速度限制了 (run time).

## Faster R-CNN

让 CNN 自己来找 region proposals. 通过 Region Proposal Network (RPN) 得到 features of 提出的 proposals.

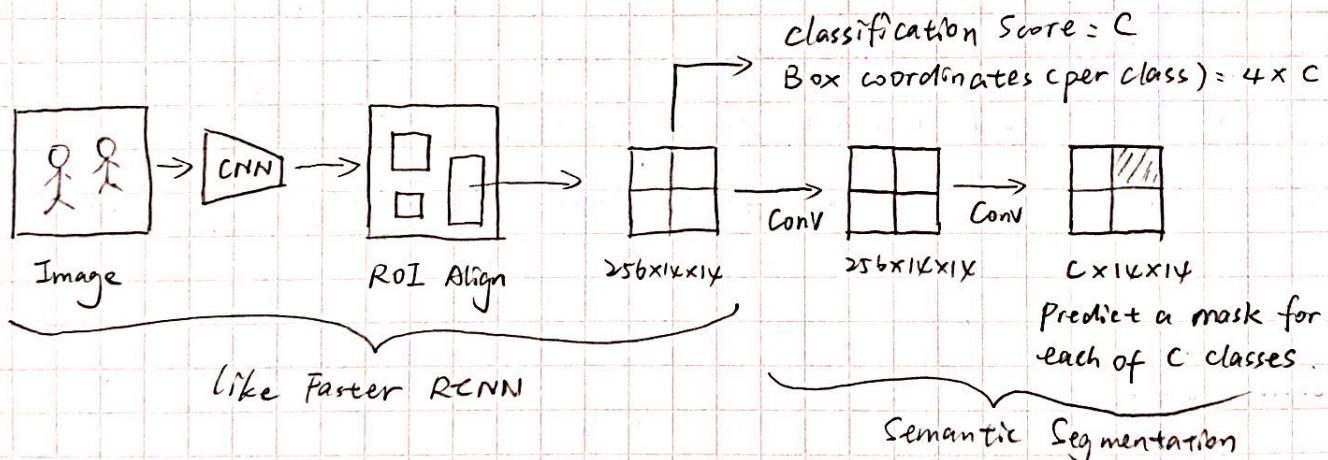
Jointly train with 4 losses

1. RPN classify object / not object
2. RPN regress box coordinates
3. Final classification score (object classes)
4. Final box coordinates



## Mask R-CNN

for Instance Segmentation. and also does poses. estimation.



Mask R-CNN 和 Faster R-CNN 相同，只不过产生的是一个 mask 放在所选区域中对像素进行分类，不确定这是不是属于某一个物体

## Detection without Proposals

YOLO / SSD ( You Only Look Once / Single-shot MultiBox Detector )

网络首先将图片分成  $7 \times 7$  的 grid. 然后将一组 base boxes 分别放在 grid 的单元格中

对于每个 grid cell:

- Regress from each of the  $B$  base boxes to a final box with  $5$  numbers.  
( $\text{dx}$ ,  $\text{dy}$ ,  $\text{dh}$ ,  $\text{dw}$ , confidence)
- Predict scores for each of  $C$  classes (including background as a class)

Output:  $7 \times 7 \times (5 \times B + C)$  的 tensor.

## Summary

Object Detection has a lot of variables. (Base network, Object detection architecture, Imagesize of region Proposals).

Faster R-CNN is slower but more accurate.

SSD is much faster but not as accurate.

paper: Speed/accuracy trade-offs for modern convolutional object detection.

## [Pseudo-code]

- R-CNN.

```

RoIs = region_proposal(image)
for ROI in RoIs:
    patch = get_patch(image, ROI)
    results = detector(patch)
  
```

run CNN

## — Fast R-CNN.

run CNN. (现代卷积网络 for FPN-513)  
从RPN中直接  
得到RoIs

```

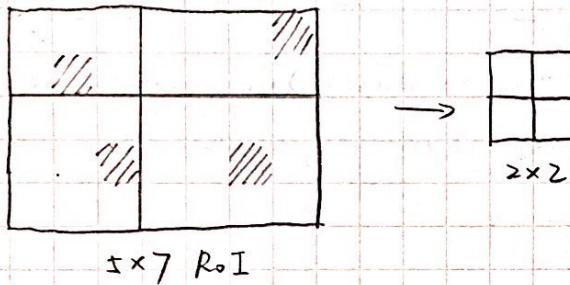
feature-maps = process(image)
RoIs = region-proposal(image) ← Expensive!
for ROI in RoIs:
    patch = roi-pooling(feature-maps, ROI)
    results = detector2(patch)
    ↓ fully connected layer.

```

ROI pooling:

因为 Fast R-CNN 使用的是 FC 层，输入以不同 size 的 RoI 需要 map 成相同大小。

做法是，将 RoI 划分成若干 predefined size，再取 max。



## — Faster R-CNN

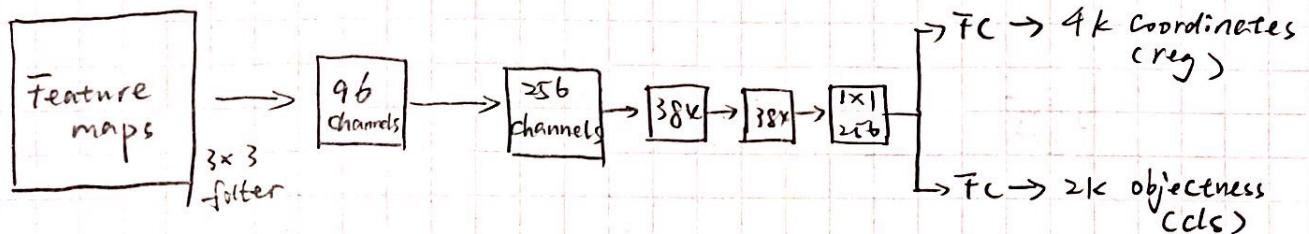
```

feature-maps = process(image)
RoIs = region-proposal-network(feature-maps)
for ROI in RoIs:
    patch = roi-pooling(feature-maps, ROI)
    class-scores, box = detector(patch)
    class-probabilities = softmax(class-scores)

```

与 Fast R-CNN 类似，但不同点在于引入了 RPN，并且 RoIs 的生成也从 513 变为 3 feature maps。从而取代了 Fast R-CNN 中从 513 产生 RoIs 的 region proposal methods。

## Region proposal network (RPN)



- RPN 用第一阶段的 CNN 来输出 feature maps 作为输入。
- 再使用  $3 \times 3$  filter 的 CNN (ZF network) 提取特征，(也可以使用 VGG 或 ResNet)。
- ZF network 输出的 256 个值分别进入 2 个 FC。一个预测 bounding box，一个产生 2 个 objectness scores。

Faster R-CNN 中将其作为 classification 问题处理，即有物体类和无物体（背景）类。  
也可以简化为 regression 及 objectness score。

- 对于 feature maps 上的每个位置 ( $3 \times 3$  input) , RPN 都会产生  $K$  个预测。  
其中  $K$  是指  $K$  个固定形状 (carefully pre-selected) 的 anchor boxes.  
Anchor 也被称作 priors 或 default boundary boxes.

Faster R-CNN 中使用了 9 个 anchors。

## — Region-based Fully Convolutional Networks (R-FCN)

```

feature-maps = process(image)
ROIs = region-proposal(feature-maps)
score-maps = compute-score-map(feature-maps)
for ROI in ROIs
    V = region-roi-pool(score-maps, ROI)
    class-scores, box = average(V)
    class-probabilities = softmax(class-scores)
  
```

Faster R-CNN 需要对每一个 ROI 通过 FC 计算结果，对于很多的 ROI，计算量会很大。  
R-FCN 通过减少每个 ROI 需要的计算量来加快速度。

## Lecture 12 Generative Models.

### Supervised Learning v.s. Unsupervised Learning

	Supervised	Unsupervised
Data	$(x, y)$	$x \leftarrow \text{training data}$ $\uparrow \text{cheap}$
Goal	Learn a function to map $x \rightarrow y$	Learn some underlying hidden structure of the data.
Examples	Classification, regression, object detection, semantic segmentation, image captioning...	Clustering, dimensionality reduction, feature learning, density estimation.

### Generative Models

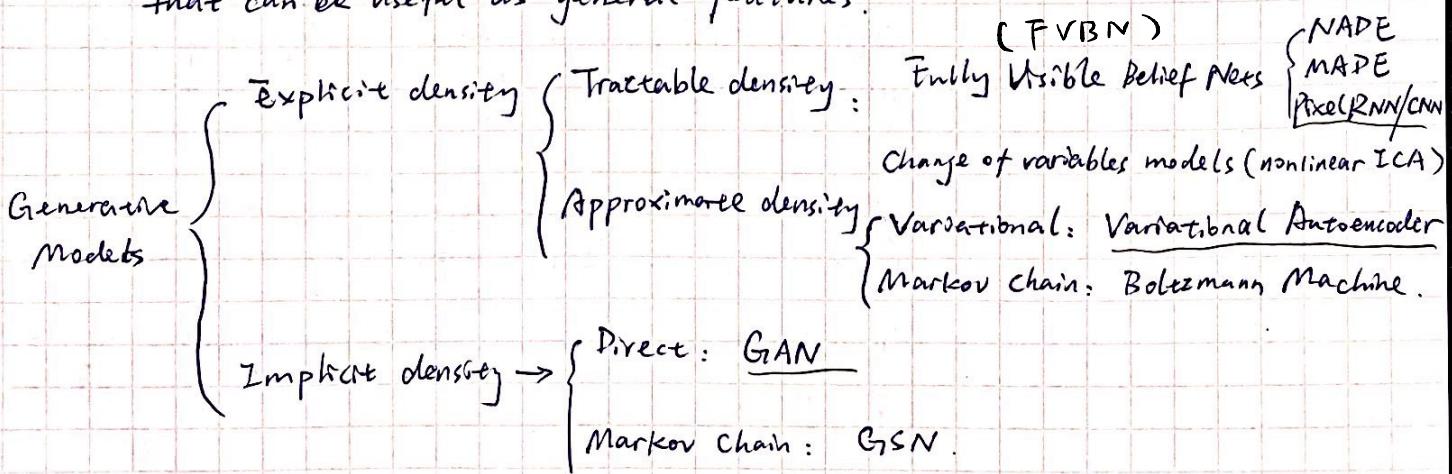
— Give training data, generate new images(samples) from same distribution.

Generative models address density estimation, which is a core problem in unsupervised learning. (Learn  $P_{\text{model}}(x)$  similar to  $P_{\text{data}}(x)$ ).

- Explicit density estimation  
explicitly define and solve for  $P_{\text{model}}(x)$
- Implicit density estimation.  
Learn model that can sample from  $P_{\text{model}}(x)$  without explicitly defining it.

— Why generative models?

1. Realistic samples for artwork, super-resolution, colorization, etc.
2. Generative models of time-series data can be used for simulation and planning (reinforcement learning applications)
3. Training generative models can also enable inference of latent representations that can be useful as general features.

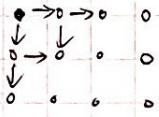


## Pixel RNN and Pixel CNN.

Use chain rule to decompose likelihood of an image  $x$  into product of

1-d distributions:

$$P(x) = \prod_{i=1}^n P(x_i | x_1, x_2, \dots, x_{i-1})$$



Pixel RNN, generate image pixels starting from corner dependency on previous pixels modeled using an RNN (LSTM)

Drawback: sequential generation is slow.

Pixel CNN, still generate image pixels starting from corner.

dependency on previous pixels now modeled using a CNN over context region.

Training is faster than Pixel RNN (can parallelize convolutions since context region values known from training set).

Generation still sequential (slow).

## VAE (Variational Autoencoders)

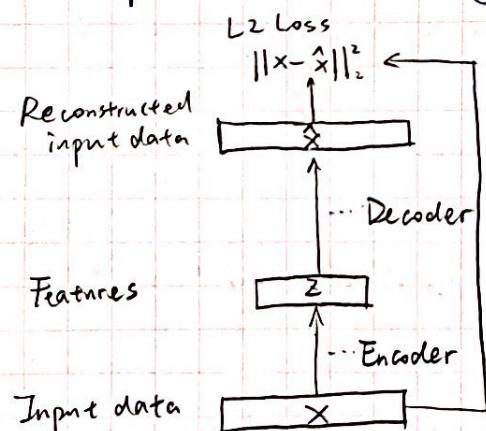
VAEs define intractable density function with latent  $z$ :

$$P_\theta(x) = \int P_\theta(z) P_\theta(x|z) dz$$

Cannot optimize directly, derive and optimize lower bound on likelihood instead.

### — Auto-encoders.

Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data.



L2 Loss  
 $\|x - \hat{x}\|_2^2$   
 通常 Z 不仅仅是 raw 像素，因为我们需要的是提取有意义的特征 /  
 变化的因素。

Encoder  $\rightarrow$  CNN      Decoder  $\rightarrow$  CNN (upconv)  
 $\rightarrow$  maximize the likelihood of training data.  
 训练这个模型时，这个 autoencoder 是不可用的。因为它需要 decoder，所以  
 它 encoder 里的 supervised model no initialization.

这样做的目的是让模型可以用许多元特征来训练，而不是只用  
 特定的普遍特征 (general feature representation)

用这些特征来进一步优化一个监督学习问题。通过将特征  
 分成不同的分支。

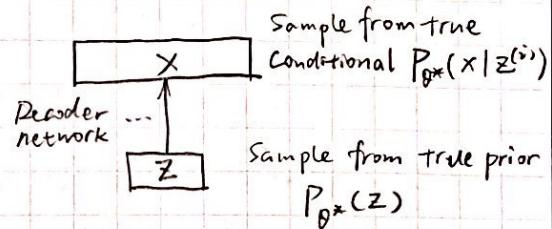
## Variational Autoencoders

Autoencoders 可以重建(reconstruct)输入图片并将其到 Features  $Z$ .

Probabilistic spin on auto-encoders will let us sample from the model to generate data.

假设训练集和后  $\{x^{(i)}\}_{i=1}^N$  是从隐含表示  $Z$  中生成的。

我们想要的是未估计的 true parameters  $\theta^*$ .



• 如何表示模型?

1. 选择最简单的先验  $P(z)$ , for Gaussian.

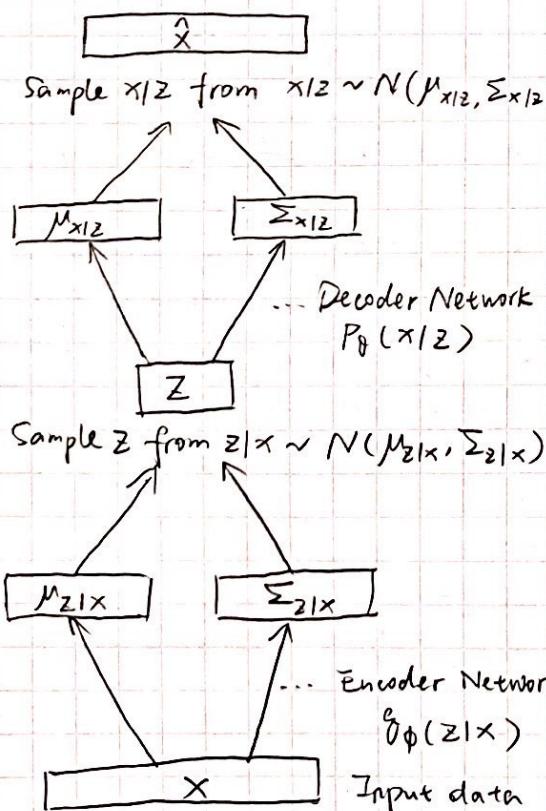
2.  $P(x|z)$  较复杂, 用 neural network 建立.

• 如何训练模型? Fully Visible Belief Networks

In FVBN 中, 通过最大似然 likelihood of training data 来训练模型. 但是这里  $P_\theta(x) = \int P_\theta(x|z) P_\theta(z) dz$  是无法直接优化的. 因为这个积分是 intractable 的.

Intractable to compute  $P(x|z)$  for every  $z$ . 同时后验  $P(z|x) = P(x|z)P(z)/P(x)$  也是 intractable, 因为不知道  $P(x)$ .

Solution: 在 decoder network  $P_\theta(x|z)$  之外, 增加 encoder  $\theta_\phi(z|x)$  逼近  $P_\theta(z|x)$ ,



$$\begin{aligned}
 \log P_\theta(x^{(i)}) &= E_{z \sim \theta_\phi(z|x^{(i)})} [\log P_\theta(x^{(i)})] \quad \text{可以这样构造} \\
 &= E_z \left[ \log \frac{P_\theta(x^{(i)}|z) P_\theta(z)}{P_\theta(z|x^{(i)})} \right] \quad \text{原因在于 } P_\theta(x^{(i)}) \text{ 不依赖于 } z; \\
 &= E_z \left[ \log \frac{P_\theta(x^{(i)}|z) P_\theta(z)}{P_\theta(z|x^{(i)})} \cdot \frac{\theta_\phi(z|x^{(i)})}{\theta_\phi(z|x^{(i)})} \right] \quad \text{这样才好处理} \\
 &= E_z [\log P_\theta(x^{(i)}|z)] - E_z [\log \frac{\theta_\phi(z|x^{(i)})}{P_\theta(z|x^{(i)})}] + E_z [\log \frac{\theta_\phi(z|x^{(i)})}{P_\theta(z|x^{(i)})}] \\
 &= E_z [\log P_\theta(x^{(i)}|z)] - \underbrace{D_{KL}(\theta_\phi(z|x^{(i)}) || P_\theta(z))}_{\textcircled{1}} - \underbrace{D_{KL}(\theta_\phi(z|x^{(i)}) || P_\theta(z|x^{(i)}))}_{\textcircled{2}} + \underbrace{\log \theta_\phi(z|x^{(i)})}_{\geq 0 \text{ intractable}}
 \end{aligned}$$

$$\log P_\theta(x^{(i)}) \geq L(x^{(i)}, \theta, \phi)$$

$$\text{Training: } \theta^*, \phi^* = \arg \max_{\theta, \phi} \sum_{i=1}^N L(x^{(i)}, \theta, \phi)$$

maximize the lower bound.

① Maximize likelihood of original input being reconstructed

② Make approximate posterior distribution close to prior

① 通过 sampling 得到 fit (sampling differentiable through reparam. "trick")

②  $\theta_\phi$  是 encoder 的参数分布  $\theta_\phi(z|x)$  和 prior  $P_\theta(z|x)$  的 KL 距离. 是 closed-form.

## • Synthetic data?

只使用 decoder 来生成图片的生成。Z 和  $p_{\phi}(z|x)$  不同于训练时从  $p_{\phi}(z|x)$  中采样，而是仅用我们指定的 prior，如  $Z \sim N(0, I)$ 。

Z 表示的其实仅仅是隐含的特征。因为 Z 的 prior 是 diagonaliz，所以表示互相独立的隐含变量。不同 in Z 的值及编码码了可解释的变化因素。（维度也是人为规定）

类似于 Autoencoder，Z 同样是属于特征表示 (feature representation)，也可以认为 encoder 部分输入图片生成 Z 用来进行 supervised learning。

## • 总结：

Probabilistic spin to traditional auto-encoders  $\Rightarrow$  allows generating data.

Defines an intractable density  $\Rightarrow$  derive and optimize a (variational) lower bound.

## Pros:

Principled approach to generative models.

Allows inference of  $q_{\phi}(z|x)$ , can be useful feature representation.

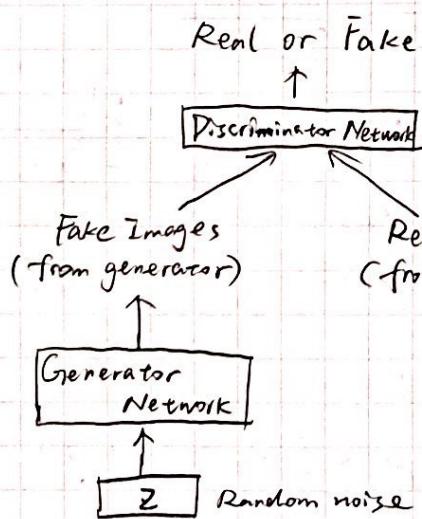
## Cons:

Maximize lower bound of likelihood, not as good evaluation as PixelRNN/CNN

Samples blurrier and lower quality compared to state-of-the-art (GANs)

## Generative Adversarial Networks (GAN) unsupervised learning with a supervised component

GANs don't work with any explicit density function. (没有显式概率生成函数)。



Generator network: try to fool the discriminator by generating real-looking images.

Discriminator network: try to distinguish between real and fake images  
Outputs likelihood in  $(0, 1)$  of real image

Minimax objective function.

$$\min_{\theta_g} \max_{\theta_d} \left[ E_{x \sim P_{data}} \log D_{\theta_d}(x) + E_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Discriminator output for real data(x)      Discriminator output for generated data  $G(z)$

Alternate between:

→ 1. Gradient ascent on discriminator

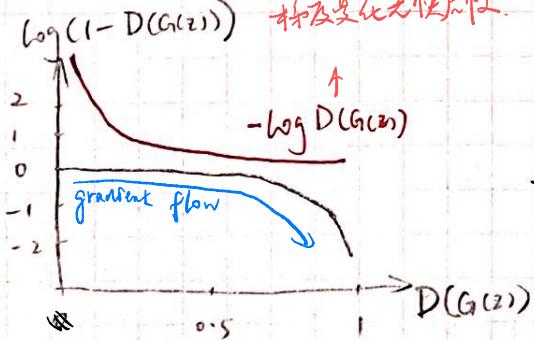
$$\max_{\theta_d} \left[ E_{x \sim P_{data}} \log D_{\theta_d}(x) + E_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

→ 2. Gradient descent on generator.

$$\min_{\theta_g} E_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z)))$$

真实的 discriminator  
输出正确的决策概率

虚假的 discriminator  
输出错误决策概率



$\log(1 - D(G(z)))$  is 1 if  $z \in \mathcal{M}$  当生成 its sample. 0 if  $z \in \mathcal{D}$  (false)  
即，梯度小，训练慢. 相反当生成 in sample 时梯度大  
梯度大，容易振荡，甚至发散.

Jointly training two networks is challenging, can be unstable.  
Choosing objectives with better loss landscapes helps training.

优化  $\rightarrow$  Alternate between:

1. Gradient ascent on discriminator

$$\max_{\theta_d} \left[ E_{x \sim p_{\text{data}}} \log D_{\theta_d}(x) + E_{z \sim p_g(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. Gradient ascent on generator.  $\rightarrow$  使 generator, 使 discriminator

$$\max_{\theta_g} E_{z \sim p_g(z)} \log (D_{\theta_d}(G_{\theta_g}(z)))$$

outperform

- GAN training algorithm.

for number of training iterations do

```

    for k steps do
        |   for K steps do
        |   |   • Sample minibatch of m noise samples {z(1)...z(m)} from noise prior Pg(z)
        |   |   • Sample minibatch of m examples {x(1)...x(m)} from data generating distribution Pdata(x)
        |   |   • Update the discriminator by ascending its stochastic gradient:
        |   |    $\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m [\log D_{\theta_d}(x^{(i)}) + \log (1 - D_{\theta_d}(G_{\theta_g}(z^{(i)})))]$ 
        |   end for
        |   • Sample minibatch of m noise samples {z(1)...z(m)} from Pg(z)
        |   • Update the generator by ascending its stochastic gradient:
        |    $\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (D_{\theta_d}(G_{\theta_g}(z^{(i)})))$ 
    end for.
  
```

Wasserstein GAN balance the training times between 2 networks.

After training, use generator network to generate new images.

## — Deep Convolutional GANs.

Generator is an up sampling network with fractionally-strided convolutions.

Discriminator is a convolutional network.

Architecture guidelines for stable Deep Conv GANs =

- Replace any pooling layers with stride conv in discriminator and fractional-strided convolutions in generator.
- Use batchnorm in both discriminator and generator
- Remove FC for deeper architectures
- Use ReLU in generator for all layers except for output, which uses tanh.
- Use Leaky ReLU in discriminator for all layers.

## Lecture 13 Visualizing and Understanding

# Lecture 14 Deep Reinforcement Learning

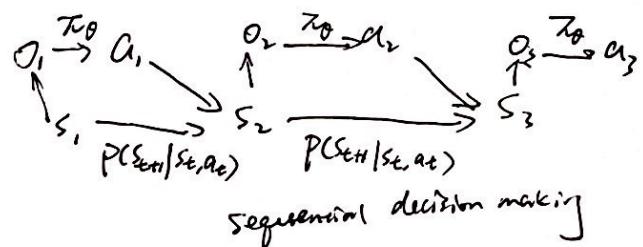
Markov

→ Definition:

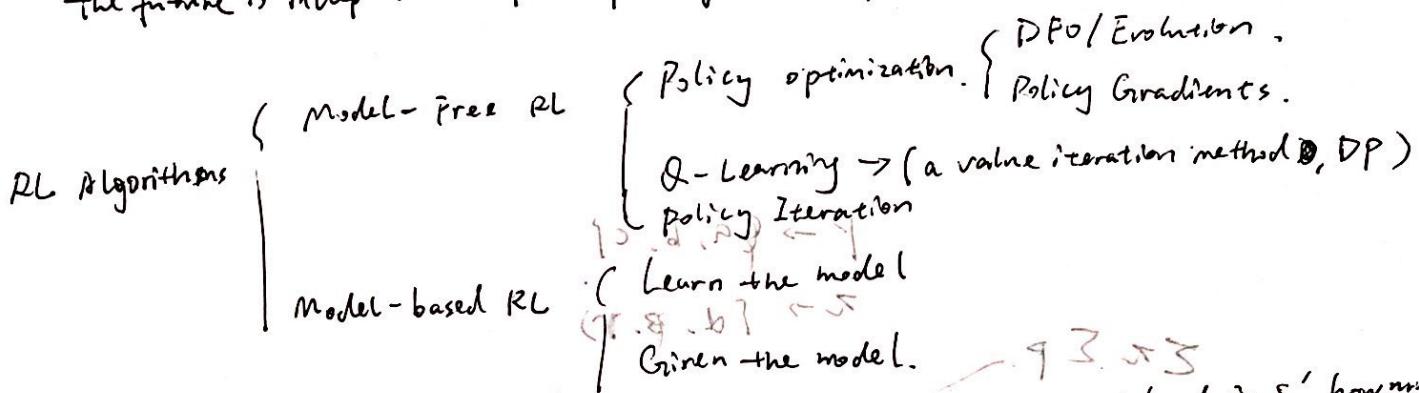
A state  $s_t$  is Markov if and only if:

$$P(s_{t+1}|s_t) = P(s_{t+1}|s_1, \dots, s_t)$$

"The future is independent of the past given the present."



sequential decision making



Markov Decision Process (MDP)

States:  $S$  Actions:  $A$  Transition function:  $P(s'|s, a)$  How long it will be  
Reward function:  $R(s, a, s')$  Discount factor:  $\gamma$  in state  $s$  take action  $a$  and land in  $s'$ , how much reward get for that.

Goal:  $\max_{\pi} E \left[ \sum_{t=0}^H \gamma^t R(s_t, a_t, s_{t+1}) | \pi \right]$

Deterministic policy:  $a = \pi(s)$

Stochastic policy:  $\pi(a|s) = P(A_t=a|s_t=s)$

Map from state to distributions of actions.

→ Value Iteration

Optimal Value Function:

(Optimal) State-value function:

$$V^*(s) = \max_{\pi} E \left[ \sum_{t=0}^H \gamma^t R(s_t, a_t, s_{t+1}) | \pi, s_0 = s \right]$$

↓  
sum of discounted rewards when starting from state  $s$  if use the best possible policy

$$V_k^*(s) = \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_{k-1}^*(s'))$$

↓  
optimal value for state  $s$  when  $H = k$

$$V_H^* \xrightarrow{H \rightarrow \infty} V_\infty^* = V^*$$

Value Iteration Algorithm:

Start with  $V_0^*(s) = 0$  for all  $s$

For  $k = 1 \dots H$ :

For all states  $s$  in  $S$ :

$$V_k^*(s) \leftarrow \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_{k-1}^*(s'))$$

$$\pi_k^*(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_{k-1}^*(s'))$$

action-value function  $Q$ -Values.

$Q^*(s, a) \rightarrow$  expected utility starting in  $s$ , taking action  $a$ , and thereafter acting optimally (a doesn't have to be optimal action)

$$Q^*(s, a) = \sum_{s'} p(s'|s, a) (R(s, a, s') + \gamma \max_{a'} Q^*(s', a')) \xrightarrow{\text{Optimally}} V^*(s)$$

$$= E[R_{t+1} + \gamma V^*(s') | S_t = s, A_t = a]$$

$$= E[R_{t+1} + \gamma \max_{a'} Q^*(s', a') | S_t = s, A_t = a]$$

$$V^*(s) = \max_{a'} Q^*(s, a')$$

~~Max~~

$Q$ -Value Iteration.

$$Q_{k+1}^*(s, a) \leftarrow \sum_{s'} p(s'|s, a) (R(s, a, s') + \gamma \max_{a'} Q_k^*(s', a'))$$

$Q$ -Value iteration requires many steps to converge, where it's better to start with initial values

— Policy Iteration.

Policy Evaluation.

$$V_k^*(s) = \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) (R(s, a, s') + \gamma V_k^*(s'))$$

Policy Iteration:

(1) policy evaluation for current policy  $\pi_k$ :

$$\text{iterate until convergence } V_{k+1}^*(s) \leftarrow \sum_{s'} p(s'|s, \pi_k(s)) [R(s, \pi_k(s), s') + \gamma V_k^*(s')]$$

(2) policy improvement. (find the best action according to one-step lookahead)

$$\pi_{k+1}(s) \leftarrow \arg \max_a \sum_{s'} p(s'|s, a) [R(s, a, s') + \gamma V_k^*(s')]$$

In Value iteration  $\Rightarrow$  policy iteration is policy evaluation in Policy iteration  $\Rightarrow$  max.

Value iteration  $\Rightarrow$  policy iteration in改進版. (using combination). (2) of Policy iteration in改進版 policy evaluation, is "one-step lookahead" protracted iterative computation requiring multiple sweeps through the state set. Value iteration in改進版 has only one sweep.

$$V_{k+1}^*(s) \leftarrow \max_a \sum_{s'} p(s'|s, a) [R(s, a, s') + \gamma V_k^*(s')]$$

sweep: one update for each state.

## — Limitations.

- (1) iteration over storage for all states and actions: requires small, discrete state-action space  $\rightarrow Q/V$  function fitting
- (2) Update equations requires access to dynamics model ( $P(s'|s, a)$ )  
 $\searrow$  sampling-based approximations  
 $\uparrow$  unknown for model-free method.

## — Sampling-based approximation:

$$Q_{k+1} \leftarrow E_{\substack{s' \\ s \sim P(s', s)}} [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')]$$

$$Q\text{-value Iteration: } Q_{k+1}(s, a) \leftarrow \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma \max_{a'} Q_k(s', a'))$$

Tabular Q-Learning: replace expectation by samples. (off-policy learning)

Algorithm:

Start with  $Q_0(s, a)$  for all  $s, a$ .

Get initial state  $s$ .

For  $k = 1, 2, \dots$  till convergence

sample action  $a$ , get next state  $s'$

If  $s'$  is terminal:

$\text{target} = R(s, a, s') \leftarrow$  immediate reward.

Sample new initial state  $s'$

else:

$\text{target} = R(s, a, s') + \gamma \max_{a'} Q_k(s', a')$

$Q_{k+1}(s, a) \leftarrow (1 - \alpha) Q_k(s, a) + \alpha \cdot \text{target}$ .

$s \leftarrow s'$

← iterative step  
off-policy update  
 $\hat{\pi}$  is policy used  
arg max Q function.

on-policy  
↓  
start with initial  
policy and collect  
samples according to  
it strictly.

How to sample actions?

1. choose random actions?  $\rightarrow$  over all states and actions.
2. choose action that maximizes  $Q_k(s, a)$  (greedily).  $\rightarrow$  choose actions greedily according to current estimate, which is almost perfect, so that chose good actions most of the time.

✓ 3. trade-off 1 & 2.

(exploring  $\leftrightarrow$  exploring)

$\epsilon$ -greedy: random action with prob  $\epsilon$ .  
otherwise choose action greedily

$\downarrow$   
may lead to insufficient exploration if it initializes during initial phase of learning.

Challenges:

- have to explore enough
- have to eventually make the learning rate  $\alpha$  very small

value Iteration 不适合 sampling-based approximation. 因为不能同时最大化 draw samples.  
for Policy Evaluation 也不适合 sampling-based approx → TD Learning (Temporal difference)

Policy Improvement 也不适合

因为  $\pi^*$  是一个常数 (constant) 简化问题 (simplifies the problem) tabular method is not scalable.

其实在连续环境 (continuous environment) 中也一样。

## - Approximate Q-Learning

假设 Q-Learning 的实现方法是 Q-table. Instead of a table, we have a parametrized Q-function:  $Q_\theta(s, a) \rightarrow \begin{cases} \text{linear function} \\ \text{complicated neural net.} \end{cases}$  可以根据不同的 Q 值，处理 states 210 generalise.

Learning rule:

$$\textcircled{1} \text{ Remember: } \text{target}(s') = R(s, a, s') + \gamma \max_{a'} Q_K(s', a')$$

\textcircled{2} update:

$$Q_{K+1} \leftarrow Q_K - \alpha \nabla_{Q_K} \left[ \frac{1}{2} (Q_K(s, a) - \text{target}(s'))^2 \right] \Big|_{Q=Q_K}$$

Tabular Q-Learning is Approximate Q-Learning in 233 种情况 当  $Q_\theta(s, a) = \theta_{sa}$ .

假设  $\theta \in \mathbb{R}^{181 \times 181}$ ,  $Q_\theta(s, a) = \theta_{sa}$

$$\nabla_{\theta_{sa}} \left[ \frac{1}{2} (Q_\theta(s, a) - \text{target}(s'))^2 \right] = \nabla_{\theta_{sa}} \left[ \frac{1}{2} (\theta_{sa} - \text{target}(s'))^2 \right] \\ = \theta_{sa} - \text{target}(s')$$

$$\theta_{sa} \leftarrow \theta_{sa} - \alpha (\theta_{sa} - \text{target}(s')) \\ = (1-\alpha) \theta_{sa} + \alpha \text{target}(s')$$

## Deep Q-Networks

deep neural net T in approximate Q-learning task 为什么?

\textcircled{1} 因为 network 处理 states 之间具有依赖性, 处理着长期的时序信息。因此必须用 deep learning. 因为 target 是 non-stationary

\textcircled{2} 随机梯度下降 (SGD) 无法处理 data 是 i.i.d.s. 因为 RL 中存在强相关性 (strong correlations).

~~SGD~~

DQN.

Two main ideas for stabilizing Q-Learning : Experience replay

breaks local correlations from online updates  $\hookrightarrow$  ① Apply Q-updates on batches of past experience instead of online makes the data distribution more stationary

② Use an older set of weights to compute the targets (target network)  
- keeps the target function from changing too quickly

$$\downarrow \quad L_i(\theta_i) = E_{s,a,s',r,\pi} \left( r + \max_{a'} Q(s'; a'; \theta^-_i) - Q(s; a; \theta_i) \right)^2$$

alleviate correlations between Q-value and target.

Deep Q-Learning with experience replay Algorithm:

Initialize replay memory  $\mathcal{D}$  to capacity  $N$ .

Initialize action-value function  $Q$  with random weights  $\theta$

Initialize target action-value function  $\hat{Q}$  with weights  $\theta^- = \theta$ .

For episode  $= 1, \dots, M$  do

  Initialize sequence  $s_i = \{x_i\}$  and preprocessed sequence  $\phi_i = \phi(s_i)$

  For  $t=1, \dots, T$  do

    with probability  $\epsilon$  select a random action  $a_t$ .  
    otherwise select  $a_t = \arg\max_a Q(\phi(s_t), a; \theta)$

    Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$

    Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$

    Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $\mathcal{D}$ .

    Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $\mathcal{D}$ .

    Set  $y_j = \begin{cases} r_j & \text{if episode terminate at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$

    Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  w.r.t the network param  $\theta$

    Every C steps reset  $\hat{\theta} = Q$

  End For

End For

Other details : ① uses Huber loss instead of squared loss on Bellman error.

$$L_\delta(a) = \begin{cases} \frac{1}{2}a^2 & |a| \leq \delta \\ \delta(|a| - \frac{1}{2}\delta) & \text{otherwise} \end{cases}$$

② use RMSProp

$$\text{③ } \sum_{\text{Start at 1}}^1 \frac{1}{\text{million frame}} \rightarrow 0.1 \text{ or } 0.05$$

## — Policy Gradients.

- Indirect policy  $\pi_\theta(u|s) \rightarrow$  1. policy parameterized by parameter vector  $\theta$
- $\max_{\theta} E \left[ \sum_{t=0}^H R(s_t) | \pi_\theta \right]$
2. Stochastic policy class (smooth out the problem)  
 $\pi_\theta(u|s)$ : probability of action  $u$  in state  $s$ .  
 Action distribution  
 Action selection based on the distribution of actions.  
 Action selection based on the distribution of actions, which is determined by gradient descent type.

Why policy optimization?

- often  $\pi$  can be simpler than  $Q$  or  $V$
- ~~Value function~~, doesn't prescribe actions (would need dynamics model.)
- $\pi$ : need to be able to efficiently solve  $\arg \max_u Q_\pi(s, u)$   
(challenge for continuous / high-dimensional action spaces)  $\rightarrow$   $\begin{cases} \text{HAF} \\ \text{Input Convex NNs} \\ \text{Zero Energy } \pi. \end{cases}$

	Policy Optimization	Dynamic Programming	
Conceptually	Optimize what you care about	Indirect, exploit the problem structure self-consistency	
Empirically	More compatible with rich architectures (e.g. recurrence) More versatile More compatible with auxiliary objectives	More compatible with exploration and off-policy learning More sample-efficient when they work	

Likelihood Ratio Policy Gradient.

$\tau$  denotes a state-action sequence  $s_0, u_0, \dots, s_H, u_H$ .

$$R(\tau) = \sum_{t=0}^H R(s_t, u_t)$$

$$U(\theta) = E \left[ \sum_{t=0}^H R(s_t, u_t) ; \pi_\theta \right] = \sum_{\tau} P(\tau; \theta) R(\tau)$$

The goal is to find  $\theta$ :

$$\max_{\theta} U(\theta) = \max_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau)$$

$$U(\theta) = \sum_{\tau} p(\tau; \theta) R(\tau)$$

$$\nabla_{\theta} U(\theta) = \nabla_{\theta} \sum_{\tau} p(\tau; \theta) R(\tau)$$

$$= \sum_{\tau} \nabla_{\theta} p(\tau; \theta) R(\tau)$$

$$= \sum_{\tau} \frac{p(\tau; \theta)}{p(\tau; \theta)} \nabla_{\theta} p(\tau; \theta) R(\tau)$$

$$= \sum_{\tau} p(\tau; \theta) \frac{\nabla_{\theta} p(\tau; \theta)}{p(\tau; \theta)} R(\tau)$$

$$= \sum_{\tau} p(\tau; \theta) \nabla_{\theta} \log p(\tau; \theta) R(\tau)$$

$$\log'(x) = \frac{1}{x}$$

Approximate with the empirical estimate for  $m$  sample paths under policy  $\pi_{\theta}$ ,

$$\nabla_{\theta} U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log p(\tau^{(i)}; \theta) R(\tau^{(i)})$$

unbiased estimate of the  $\nabla_{\theta} U(\theta)$

Valid even when:  $R$  is discontinuous and/or unknown  
Sample space (of paths) is a discrete set.

The gradient tries to increase probability of paths with positive  $R$  and decrease the prob of paths with negative  $R$ .

! Likelihood ratio changes probabilities of experienced paths, does not try to change the paths ( $\leftrightarrow$  Path Derivative)

→ Decompose:

$$\begin{aligned} \nabla_{\theta} \log p(\tau^{(i)}; \theta) &= \nabla_{\theta} \log \left[ \prod_{t=0}^H p(s_t^{(i)} | s_t^{(i)}, u_t^{(i)}) \cdot \pi_{\theta}(u_t^{(i)} | s_t^{(i)}) \right] \\ &= \nabla_{\theta} \left[ \sum_{t=0}^H \log p(s_t^{(i)} | s_t^{(i)}, u_t^{(i)}) + \sum_{t=0}^H \log \pi_{\theta}(u_t^{(i)} | s_t^{(i)}) \right] \\ &= \nabla_{\theta} \sum_{t=0}^H \log \pi_{\theta}(u_t^{(i)} | s_t^{(i)}) \\ &= \sum_{t=0}^H \underbrace{\nabla_{\theta} \log \pi_{\theta}(u_t^{(i)} | s_t^{(i)})}_{\text{no dynamics model required.}} \end{aligned}$$

this is the neural net, which encodes the distribution of actions given states.

→ Derivation from Importance Sampling.

$$U(\theta) = \mathbb{E}_{\tau \sim \theta_{\text{old}}} \left[ \frac{P(\tau|\theta)}{P(\tau|\theta_{\text{old}})} R(\tau) \right]$$

$P(\tau|\theta_{\text{old}})$

Current policy  
 $P(\tau|\theta)$

Importance sampling means sample from the old policy to derive new ~~policy~~.

The ratio  $\frac{P(\tau|\theta)}{P(\tau|\theta_{\text{old}})}$  gives the correction.

$$\nabla_{\theta} U(\theta) = \mathbb{E}_{\tau \sim \theta_{\text{old}}} \left[ \frac{\nabla_{\theta} P(\tau|\theta)}{P(\tau|\theta_{\text{old}})} R(\tau) \right]$$

$$\nabla_{\theta} U(\theta)|_{\theta=\theta_{\text{old}}} = \mathbb{E}_{\tau \sim \theta_{\text{old}}} \left[ \frac{\nabla_{\theta} P(\tau|\theta)|_{\theta_{\text{old}}}}{P(\tau|\theta_{\text{old}})} R(\tau) \right]$$

$$= \mathbb{E}_{\tau \sim \theta_{\text{old}}} [\nabla_{\theta} \log P(\tau|\theta)|_{\theta_{\text{old}}} R(\tau)] \rightarrow \text{the same as former derivation}$$

suggests we can also look at more than just gradient. E.g. can use importance sampled objectives as "surrogate" loss. (change your policy after your gradient update, use this ~~equation~~ equation to estimate how good that update was.)

→ So far, unbiased but very noisy

~~Fixes~~ that lead to real-world practicality:

Increase gradient of paths that better than average and decrease worse than average.

2 Baseline.

$$\nabla_{\theta} U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log P(\tau^{(i)}; \theta) (R(\tau^{(i)}) - b)$$

Remove terms that don't depend on current action can lower variance:

$$\frac{1}{m} \sum_{i=1}^m \sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \left( \sum_{k=t}^{H-1} R(s_k^{(i)}, a_k^{(i)}) - b(s_t^{(i)}) \right)$$

→ choices of  $b$

• constant baseline:  $b = \mathbb{E}[R(\tau)] \approx \frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})$

• optimal constant baseline:  $b = \frac{\sum_i (\nabla_{\theta} \log P(\tau^{(i)}; \theta))^2 R(\tau^{(i)})}{\sum_i (\nabla_{\theta} \log P(\tau^{(i)}; \theta))^2}$

• time based:  $b_t = \frac{1}{m} \sum_{i=1}^m \sum_{k=t}^{H-1} R(s_k^{(i)}, a_k^{(i)})$

• State-dependent:

$$b(s_t) = \mathbb{E}[r_t + r_{t+1} + \dots + r_{H-1}] = V^{\pi}(s_t)$$

Increase the logprob of action proportionally to how much its returns are better than the expected return under the current policy

1. Baseline
2. Temporal structure
3. KL-divergence trust region/natural gradient.

↓  
General trick, equally applicable to perturbation analysis and finite differences.

estimation of  $V^\pi$ :

Init  $V_{\phi_0}^\pi$

Collect trajectories  $\tau_1, \dots, \tau_m$

Regress against empirical return:  $\phi_{i+1} \leftarrow \arg \min_{\phi} \frac{1}{m} \sum_{v=1}^m \sum_{t=0}^{H-1} (V_{\phi}^\pi(s_t^{(v)}) - (\sum_{k=t}^{H-1} R(s_k^{(v)}, u_k^{(v)}))^2)$

Monte-Carlo.



~~ID~~ to estimate of  $V^\pi$ :

Bellman Equation.  $V^\pi(s) = \sum_u \pi(u|s) \sum_{s'} P(s'|s, u) [R(s, u, s') + \gamma V^\pi(s')]$

Init  $V_{\phi_0}^\pi$ ,

Collect data  $\{s, u, s', r\}$

Fitted V iteration:  $\phi_{i+1} \leftarrow \min_{\phi} \sum_{(s, u, s', r)} \|r + V_{\phi_i}^\pi(s') - V_{\phi}(s)\|^2 + \lambda \|\phi - \phi_i\|_2^2$

→ Recall the likelihood Ratio PG Estimator

$$\frac{1}{m} \sum_{i=1}^m \sum_{t=0}^{H-1} \nabla_\theta \log \pi_\theta(u_t^{(i)} | s_t^{(i)}) \cdot \underbrace{\left( \sum_{k=t}^{H-1} R(s_k^{(i)}, u_k^{(i)}) - V^\pi(s_k^{(i)}) \right)}$$

Estimation of  $\hat{Q}$  from single roll-out

$$\hat{Q}^\pi(s, u) = E[r_0 + \gamma r_1 + \gamma^2 r_2 + \dots | s_0 = s, u_0 = u]$$

high variance per sample based / no generalization

→ reduce variance by discounting  $E[r_0 + \gamma r_1 + \gamma^2 r_2 + \dots | s_0 = s, u_0 = u]$   
by function approximation → critic

Reducing Variance by Function Approximation.

$$\begin{aligned} Q^{\pi, \gamma}(s, u) &= E[r_0 + \gamma r_1 + \gamma^2 r_2 + \dots | s_0 = s, u_0 = u] && (1-\gamma) \\ &= E[r_0 + \gamma V^\pi(s_1) | s_0 = s, u_0 = u] && (1-\gamma)\gamma \\ &= E[r_0 + \gamma r_1 + \gamma^2 \cancel{V^\pi}(s_2) | s_0 = s, u_0 = u] && (1-\gamma)\gamma^2 \\ &= \dots && (1-\gamma)\gamma^3 \end{aligned}$$

Async Advantage Actor Critic (A3C):  $\hat{Q} = \text{one of the above choices (e.g. } k=5 \text{ step look ahead)}$   
Generalized Advantage Estimation (GAE):  $\hat{Q} = \text{lambda exponentially weighted average of all the above.}$

Actor-Critic with A3C or GAE.

init  $\pi_{\theta_0}$ ,  $V_{\phi_0}^{\pi}$

collect roll-outs  $\{s, u, s', r\}$  and  $\hat{Q}_i(s, u)$

$$\text{update: } \phi_{i+1} \leftarrow \min_{\phi} \sum_{\{(s, u, s', r)\}} \|\hat{Q}_i(s, u) - V_{\phi}^{\pi}(s)\|_2^2 + \gamma \|(\phi - \phi_i)\|_2^2$$

$$\phi_{i+1} \leftarrow \phi_i + \alpha \frac{1}{m} \sum_{k=1}^m \sum_{t=0}^{H-1} \nabla_{\phi} \log \pi_{\phi_i}(u_t^{(k)} | s_t^{(k)}) (\hat{Q}_i(s_t^{(k)}, u_t^{(k)}) - V_{\phi_i}^{\pi}(s_t^{(k)}))$$

TD



Initialize policy param  $\theta$ . critic param  $\phi$ .

For iteration = 1, 2 ... do

Sample m trajectories under the current policy

$$\Delta \theta \leftarrow 0$$

For i=1 ... m do

For t=1 ... T do

$$A_t = \sum_{t' \geq t} \gamma^{t'-t} \cdot r_t^{(i)} - V_{\phi}(s_t^{(i)})$$

$$\Delta \theta \leftarrow \Delta \theta + A_t \cdot \nabla_{\theta} \log \pi_{\theta}$$

$$\Delta \phi \leftarrow \sum_i \sum_t \nabla_{\phi} \|A_t^{(i)}\|^2$$

Monte-Carlo

$$\theta \leftarrow \alpha \Delta \theta$$

$$\phi \leftarrow \beta \Delta \phi$$

End For

## Lecture 16 Adversarial Examples and Adversarial Training

### — What are adversarial examples?

- 1) 选择任一图片
- 2) 选择任一类别 (错误类别)
- 3) 修改图片并最大化该类别
- 4) Repeat until network is fooled.

最终生成的 adversarial example 对于 network 来说是不同类别，但人脸识别不出且 network 对错误的类别有较高的 confidence (e.g. 0.999)

Not just for neural nets. → Linear models (Logistic regression, softmax regression, SVMs)  
→ Decision Trees  
→ Nearest neighbors.

### — Why do they happen?

最初怀疑是由于 overfitting，但是因为 overfitting，导致 its mismatching 会有一定的随机性。而实际发生的现象是不同的模型会 misclassify 相同的 adversarial examples，并会 assign 相同的 class 给它们。另一个现象是将 input in adversarial example 之后加 offset vector 加到干净的图片上，这个图片也会变成 adversarial example. (almost always) 可以说这是 a systematic effect 而不是 random effect.

↓

真正的原因是 underfitting. → modern deep nets are very piecewise linear

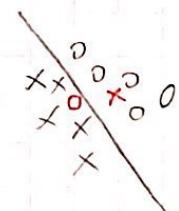
激活函数如 ReLU, Maxout 都是 piecewise linear 的。

Carefully tuned sigmoid 也在某些情况下是线性的。

甚至 LSTM 使用加法来累计和记录信息 over time. 也是线性的。

这里的线性是指从输入到输出的映射是线性的。从参数到输出的映射是非线性的，因为 weight matrices 是相乘的。这也使得训练 neural network 变得很难。  
不容易

从输入到输出的 mapping 更加线性，意味着 optimization problems that aim to optimize the input to the model are much more easier than optimization problems that aims to optimize the parameters.



## The Fast Gradient Sign Method.

$$J(\tilde{x}, \theta) \approx J(x, \theta) + (\tilde{x} - x)^T \nabla_x J(x)$$

cost w.r.t. input  $x$ .

$$\text{Maximize } J(x, \theta) + (\tilde{x} - x)^T \nabla_x J(x)$$

$$\text{Subject to } \|\tilde{x} - x\|_\infty \leq \epsilon \rightarrow \text{保证 pixels 的变化对于人眼不可分割}$$

$$\Rightarrow \tilde{x} = x + \epsilon \cdot \text{sign}(\nabla_x J(x))$$

其它生成 adversarial examples 的方法，Nicholas Carlini's attack based on multiple steps of the Adam optimizer.

FGSM 对于无 defenses 的普通神经网络有超过 99% 的 attack success rate.

FGSM 表明只需要找到与 gradient 在梯度上升的方向，并在该方向上移动很小距离就可以欺骗模型。

另外需要注意的是 ~~Adversarial examples are not noise.~~ Adversarial examples are not noise. 特别在于，在高维的空間里，random vector 平均有 zero dot product with reference vector. (参考 FGSM 中  $\epsilon \cdot \text{sign}(\nabla_x J(x))$ ，adversarial examples 选择最大的 cost).

Adversarial region (subspace) 指在 gradient 向量所在空间的正交子空间。

MNIST 有平均 25% 的 adversarial dimensions. adversarial dimensions 通常很大。意：未着起成容易产生 adversarial examples，且不同模型生成容易 intersect 而生成相同的 adversarial examples.

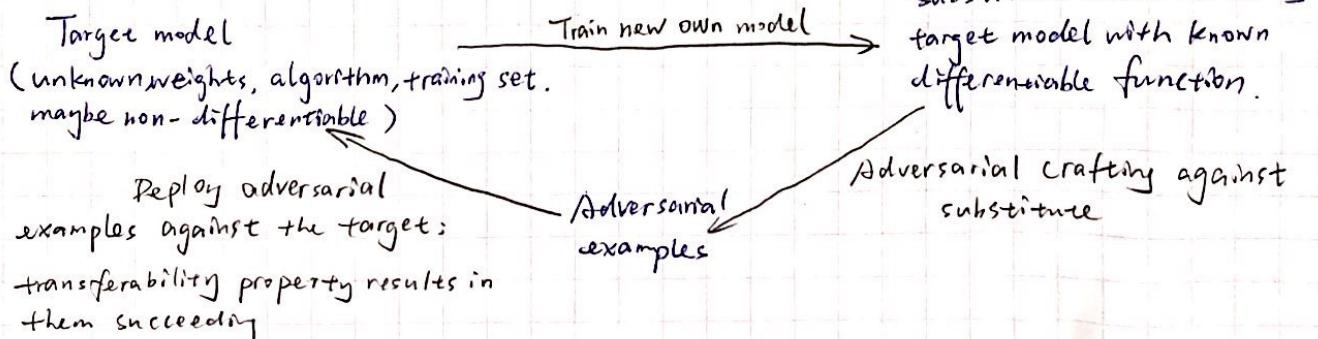
random noise.

— How can they be used to compromise machine learning systems?

由于 adversarial examples 具有 transferability，意味着对于不同 dataset，甚至不同 machine learning technique (DNN, LR, SVM, DT, KNN)，Adversary 具有可转移性。

具有可转移性的原因在于，模型即使不同，dataset <sup>即使</sup> 不同，但最终学习到的 function 是类似的，其对 adversarial examples 都是相同的。

Transferability Attack :



substitute model 与真实用多种模型的 ensemble 加强 transferability.

## — What are the defenses?

Generative Modeling is not sufficient to solve the problem.

如果已知真实正确  $P(x)$  及  $P(y|x)$  可以判断出哪些是 adversarial example. 但现实是  $P(y|x)$  是 posterior  $P(y=1|x)$  也不等于  $P(y=1|x)$  不等于  $P(x)$ .

对于 posterior  $P(y=1|x)$  不能更正. 如果一个够好的生成模型可以得到真实的 image 以及足够好的 posterior 分布. 则这种方法是可行的. 但如何训练却是很难得的, 通常得到的只是  $\approx$  linear, 但得 posterior distribution to extrapolate linearly again.

Universal approximator theorem: Neural nets can represent (approximate) continuous functions on compact subset of  $R^n$ .

现在我们问: 把 NN 训练出具有 linear decision function, 如果我们想有 step function?

↓

Training on Adversarial Examples.

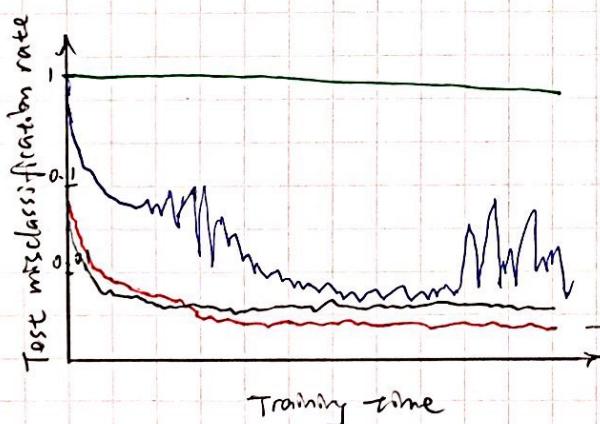
↓

对所有输入  $x$  都加上 adversarial example 并映射 (map to right label)

↓

这样可以在很大程度上避免对抗性攻击.

empirical



— Train = Clean Test = clean  
— Train = Clean Test = Adv  
— Train = Adv Test = clean  
— Train = Adv Test = Adv } Test on the same kind of adv examples that trained on.

红色比蓝色还坏, 说明这里 training on adv 是个好 regularizer

如果原模型 overfitting, training on adv 可以帮助其减少 overfitting.

Adversarial trained neural nets have the best empirical success rate on adversarial examples of any machine learning models.

↓

Deep NN can be trained to be nonlinear, so it can be a path to a solution.

## — Conclusion.

Attacking is easy.

Defending is difficult.

Adversarial training provides regularization and semi-supervised learning.

The one-of-domain input problem is a bottleneck for model-based optimization generally.