

Lecture 1 Introduction

Generalized Cylinder 义圆柱

Pictorial Structure 图像结构

CS231 focuses on one of the most important problems of visual recognition - image classification

Year 2010 NEC-UIC

Year 2012 Super Vision

Year 2014 GoogleNet VGG

Year 2015 MSRA

Data → Labeled : PASCAL

Fei-Fei Li

Lecture 2 Image Classification

All pixel change when the camera moves

Challenges:

- | Deformation
- | Occlusion → robust
- | Background clutter
- | Intra-class variation
- ...

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

Distance Metric to compare images



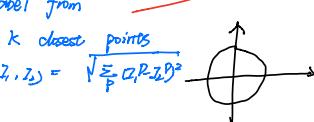
★ L_1 distance $d_1(z_1, z_2) = \sum_p |z_{1p} - z_{2p}|$

★ K -Nearest neighbors: Instead of copying label from nearest neighbor, take majority vote from K closest points

★ L_2 distance Euclidean distance $d_2(z_1, z_2) = \sqrt{\sum_p (z_{1p} - z_{2p})^2}$

problems: ① very slow at test time

② Distance metrics on pixels are not informative



Hyper parameters: choices best the algorithm that

we set rather than learn λ (regularization) accuracy and performance

train val test set → evaluate on test

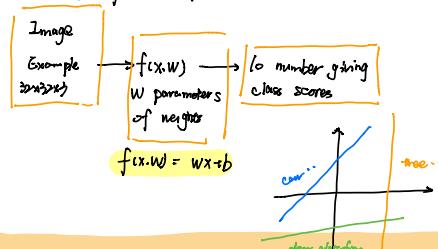
choose hyperparameters on val

Cross-validation Split data into folds by each fold

↓ as validation and average the results.

useful for small dataset but not frequently
in deep learning

Linear Classification ✓



Lecture 3 Loss Function and optimization

To Do: ① Define a loss function that quantifies our unhappiness with the scores across the training data

② Come up with a way of efficiently finding the parameters that minimize the loss function. (optimization)

$$\star \{ (x_i, y_i) \}_{i=1}^N \Rightarrow L = \sum_i L_i(f(x_i, w), y_i)$$

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

\Rightarrow if we're looping over all of the incorrect classes, so we're looping over, we will loss one.

↑
sum and mean have same effect

Loss function different choice. big loss not wrong then it's good.
square loss can enhance loss. ② You should try different loss function.

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq i} \max(0, f(x_i; w) - f(x_j; w) + 1) \Rightarrow \text{find } w \text{ for training data, resp. val}$$

↓ update
overfitting

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

Data loss = Model prediction should match training data \Rightarrow Regularization: Model should be "simple", so it works on test data.

Occam's Razor: Among competing hypotheses, the simplest is the best. — William of Ockham

$$\text{Regularization: } L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq i} \max(0, f(x_i; w) - f(x_j; w) + 1) + \lambda R(w)$$

L_2 regularization $R(w) = \sum_k \sum_l w_k w_l$

L_1 regularization $R(w) = \sum_k \sum_l |w_k|$

Elastic net ($L_1 + L_2$) $R(w) = \sum_k \sum_l (\beta w_k^2 + (1-\beta)w_k)$

Max norm regularization

Dropout

Fancier: batch normalization, stochastic depth.

$$\text{Softmax Classifier } P(Y=k | X=x_i) = \frac{e^{x_k}}{\sum_j e^{x_j}}$$

$$\text{Recap: } L_i = -\log \left(\frac{e^{y_i}}{\sum_j e^{y_j}} \right) \text{ softmax}$$

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

Optimization

follow the slope

SGD (Stochastic Gradient Descent)

$$L(w) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i; w) + \lambda R(w)$$

$$\nabla_w L(w) = \frac{1}{N} \sum_{i=1}^N \nabla_w L_i(x_i, y_i; w) + \lambda \nabla_w R(w)$$

Lecture 4: Backpropagation and Neural Networks

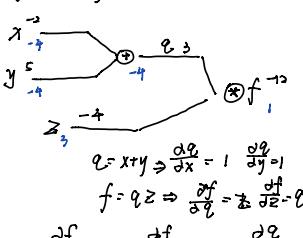
Gradient descent

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

Analytic gradients

Backpropagation: a simple example

$$f(x, y, z) = (x+y) \cdot z$$



$$\begin{aligned} f(x) &= e^x & \frac{df}{dx} &= e^x \\ f(y) &= \frac{1}{x} & \frac{df}{dy} &= -\frac{1}{x^2} \\ f(z) &= ax & \frac{df}{dz} &= a \\ f(g) &= c \cdot x & \frac{df}{dg} &= c \end{aligned}$$

local gradient \times upstream gradient

$$\begin{aligned} g(x) &= \frac{1}{1+e^{-x}} & \frac{dg}{dx} &= \frac{e^{-x}}{(1+e^{-x})^2} = \left(\frac{1+e^{-x}}{1+e^{-x}} \right) \cdot \left(\frac{1}{1+e^{-x}} \right) \\ \text{sigmoid function} & & &= (1-g(x)) \cdot g(x) \end{aligned}$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial y} \cdot \frac{\partial y}{\partial x} \quad \text{chain rule}$$

Gradients add at branches



$$\frac{\partial f}{\partial x} = \sum_i \frac{\partial f}{\partial y_i} \frac{\partial y_i}{\partial x}$$

Example: A vectorized example $f(x, w) = w^T x = \sum_{i=1}^n w_i x_i$

$$\begin{bmatrix} 0.1 \\ -0.2 \\ 0.8 \end{bmatrix} w \quad \begin{bmatrix} 0.2 \\ 0.1 \\ 0.4 \\ 0.5 \end{bmatrix} x \quad \begin{bmatrix} 0.2x_1 \\ 0.2x_2 \\ 0.8x_3 \\ 0.5x_4 \end{bmatrix} \quad \sum 0.1x_i = 1.00$$

$$q_i = w^T x = \begin{pmatrix} w_{1,1}x_1 + \dots + w_{1,n}x_n \\ \vdots \\ w_{k,1}x_1 + \dots + w_{k,n}x_n \end{pmatrix}$$

$$f(x) = \|q\|_2^2 = q_1^2 + q_2^2 + \dots + q_k^2$$

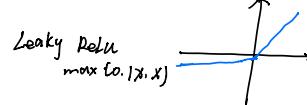
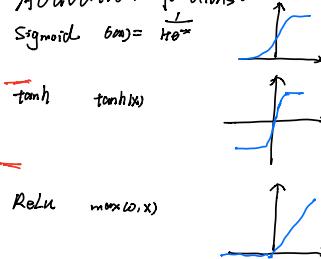
$$\frac{\partial q_k}{\partial w_{ij}} = 1 \quad k=i \quad j=j \\ \frac{\partial f}{\partial w_{ij}} = \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial w_{ij}} \\ = \sum_k 2q_k l / K = l \quad l = i, j \\ = 2q_i x_j$$

Neural Network

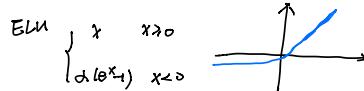
Linear score function $f = w^T x$

2-layer Neural Network $f = \max(0, w_1^T x)$

Activation functions:



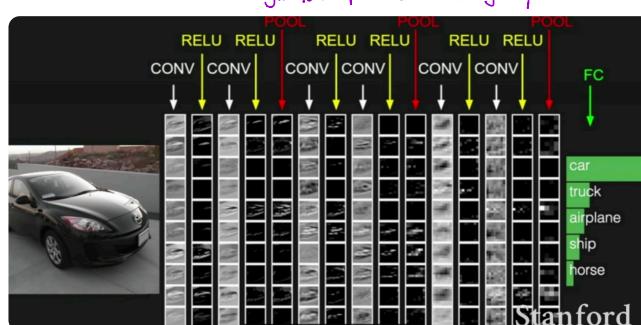
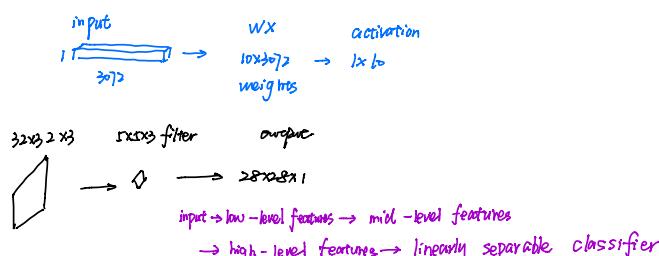
Max out
 $\max(w_1^T x + b_1, w_2^T x + b_2)$



Lecture 5 Convolutional Neural Networks

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

Topographical mapping in the cortex



stride one or two or three

input 7x7, filter 3x3

$$\text{stride } 1 \Rightarrow (7-3)/1+1 = 5$$

$$\text{stride } 2 \Rightarrow (7-3)/2+1 = 3$$

$$\text{stride } 3 \Rightarrow (7-3)/3+1 = 2$$

$$\boxed{(N-F)/\text{stride} + 1}$$

Common to zero pad the border

zero padding has influence !!

32x32x3 \rightarrow 10 x 3x3 filters with stride 1, pad 2
 \downarrow
output volume size $(32+2x2-5)/1+1 = 32$ spatially \cdot so $32x32x10$

Conclusion: Accepts a volume of size $W_1 \times H_1 \times D_1$

requires four hyperparameters $W_2 \times H_2 \times D_2$

$$\left\{ \begin{array}{l} W_2 = (W_1 - F + 2P) / S + 1 \\ H_2 = (H_1 - F + 2P) / S + 1 \\ D_2 = K \end{array} \right.$$

With parameter sharing, it introduces $F \cdot H_2 \cdot W_2 \cdot D_2$ weights per filter, for a total of $17 \cdot F \cdot H_2 \cdot W_2 \cdot K$ weights and K bias

1x1 convolution layers make perfect performance

in the output volume, the d -th depth slice (of size $H_2 \times W_2$) is the result of performing a valid convolution of the d -th filter over the input volume with a stride of S , and then offset by d -th bias

Pooling layer

- make the representations smaller and more manageable
- operates over each activation map independently

Max Pooling ✓

downsampling

Totally Connected Layer (FC layer)

Lecture 6 Training Neural Networks!

mini-batch SGD

a. One-time setup

b. Training dynamics

c. Evaluation

激活函数以0为中心的原因

权重更新方式是发生在反向更新阶段，过程：，对于神经元A来说， w_i 更新的方向和后几项都有关关系，先看超参数 η 他是人为规定的，默认值，所以学习率可以不考虑；再看最后两位乘积项 $\frac{\partial L}{\partial w_i}$ 这个对于神经元A来说，它是神经元A的误差表示，在某一次反向传播也是不变的，也不考虑，所以 w 值的更新方向只与值有关，这个值是上一层神经元的输入值，即经过了sigmoid函数激活过，所以肯定是正值，那么可以得出结论：在某一次反向传播时，对于神经元A来说， w_1 、 w_2 ...改变的方向是统一的，或正或负。如果你的最优值是需要 w_1 增加， w_2 减少，那么只能走Z字形才能实现。

Activation Functions

- squashes numbers to

range [0,1] **Sigmoid**

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



- historically popular since they have nice interpretation as a saturating "firing rate" of a neuron

tanh

$$\tanh(x)$$



ReLU

$$\max(0, x)$$



Leaky ReLU

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



Maxout

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



ELU

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



Stanford

Sigmoid

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

$f'(z)w_i x_i + b$

when $|x|$ is large, the gradient is killed

1. Saturated neurons "kill" the gradient



- 2. Sigmoid outputs are not zero-centered
- 3. $\exp(x)$ is a bit compute expensive

$$+ \xrightarrow{W_b} w_b$$

tanh

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

ReLU

$$\text{Computes } f(x) = \max(0, x) \quad x > 0 \quad \uparrow$$

- Does not saturate (in + region)

- very computationally efficient

- converges much faster than sigmoid/tanh in practice

- Actually more biologically plausible than sigmoid

dead ReLU will not activate, not update!

Leaky ReLU

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

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

- Does not saturate

- Computationally efficient

- Converges much faster than sigmoid/tanh in practice

- will not die

ELU (Exponential Linear Units)

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

- All benefits of ReLU

- Closer to zero mean outputs

- Negative saturation regime compared with leakyReLu adds some robustness to noise

* Computation require $\exp()$

1.2 Data pre process

Original data \Rightarrow zero-centered data

$$X = np.mean(X, axis=0)$$

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

Subtract the mean image or channel *

1.3 Weight Initialization

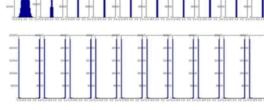
Idea 1: small random numbers

$$W = 0.01 * np.random(D, M)$$

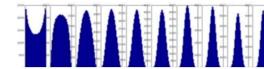


Proper initialization is hard.

Initialization too small:
Activations go to zero, gradients also zero,
No learning



Initialization too big:
Activations saturate (for tanh),
Gradients zero, no learning



Initialization just right:
Nice distribution of activations at all layers,
Learning proceeds nicely

1.4 Batch Normalization

$$\hat{x}^{(k)} = \frac{x^{(k)} - E[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}} \quad \text{this is a vanilly differentiable function}$$

usually inserted after Fully Connected or Convolutional layers,
and before non-linearity.

Batch Normalization

$$\mu_B = \frac{1}{m} \sum_{i=1}^m x_i$$

$$\sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$$

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$y_i = r\hat{x}_i + \beta = BN_{r, \beta}(x_i)$$

1.5 Babystepping the Learning Process

Double check the loss.

⇒ start with small training data

- loss not going down - learning rate too slow
- loss exploding - learning rate too fast

1.6 Hyperparameter Optimization

Cross-validation strategy
First stage: only a few epochs
Second stage: long time

1.7 monitor and visualize the loss curve

training and validation accuracy

- big gap: overfitting → increase regularization strength
- no gap: increase model capacity

Lecture 7 Training Neural Networks

Fancier optimization

saddle point → stop small, slow progress

Regularization

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

Transfer Learning

SGD: stochastic gradient descent

↓ come from mini batches so they can be noisy.

SGD + Momentum

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

handle local minima and saddle points → Nesterov Momentum



Nesterov Momentum

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

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

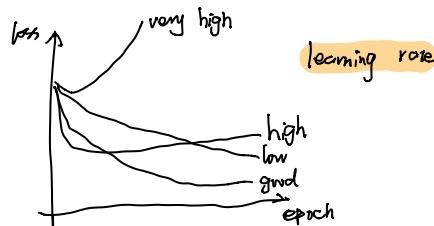
AdaGrad (Added element-wise scaling of the gradient based on the historical sum of squares in each dimension) think past

RMSProp

think past and now

Adam

Momentum AdaGrad / RMSProp Bias correction



13

Adam is good on most case. but it is batch training.

Regularization: Add term to loss

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

In common use:

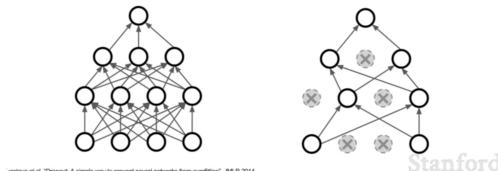
L2 regularization $R(W) = \sum_k \sum_l W_{k,l}^2$ (Weight decay)

L1 regularization $R(W) = \sum_k \sum_l$

$$\text{Elastic net (L1 + L2)} \quad R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$$

Regularization: Dropout

In each forward pass, randomly set some neurons to zero
Probability of dropping is a hyperparameter; 0.5 is common

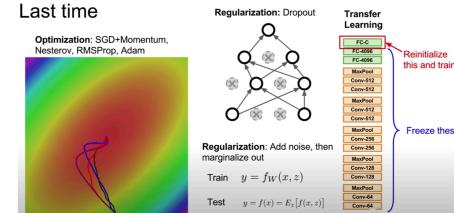
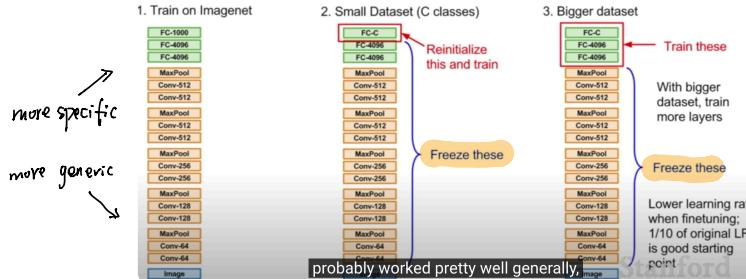


Dropout is training a large ensemble of models and make output random

Dropout, Batch Normalization, Data Augmentation
Drop Connect, Fractional Max Pooling, Stochastic Depth

Transfer Learning.

Transfer Learning with CNNs



Lecture 8 Deep Learning Software

CPU VS GPU Graphics Processing Units

Coffee Pytorch Tensorflow

Pytorch : three level of Abstraction
Tensor, Variable, Module

Lecture 9 CNN Architectures

Case Study : AlexNet

$$\text{InPut} = \text{d27} \times \text{d27} \times 3$$

First layer (conv1): 96 1×1 filters applied at stride 4 \Rightarrow Input $128 \times 112 \times 14 \times 1 = 55$

\downarrow $5 \times 5 \times 96$
 Second layer pooling 3×3 stride \Rightarrow output $27 \times 8 \times 96$
 no parameters

Parameters: $5 \times 5 \times 96 = 25K$
 ★ first use of ReLU
 use Norm Layers
 heavy Data augmentation
 dropout 0.5
 batch size 128

VGG

use smaller filter
 Stack of three 3×3 conv layers (stride 1)
 layers has same effective receptive field as 7×7
 $3 \times 3 \times C^2 < 7^2 \cdot C^2$

⇒ use ensembles for best results
 Similar training procedure as Krizhevsky
 FCN features generalize well to other tasks.

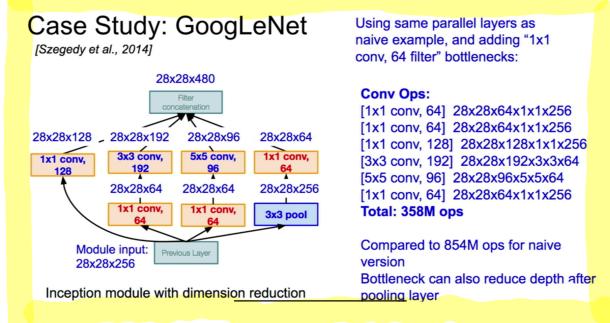
Google Net

22 layers
 Efficient "Inception" module
 No FC layers
 Only 5 million parameters

Inception module

Case Study: GoogLeNet

[Szegedy et al., 2014]



Naive Inception module

Apply parallel filter operations on the input from previous layer:

- Multiple receptive field sizes for convolution ($1 \times 1, 3 \times 3, 5 \times 5$)
- Pooling operation (3×3)

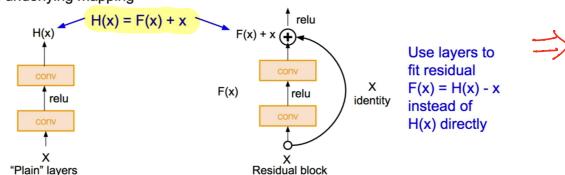
Concatenate all filter outputs together depth-wise

Case Study: ResNet very deep by using residual connection

Case Study: ResNet

[He et al., 2015]

Solution: Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping

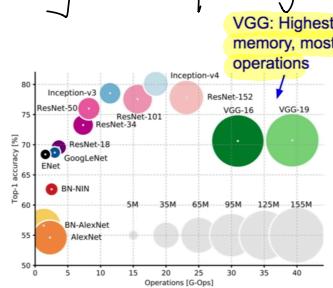
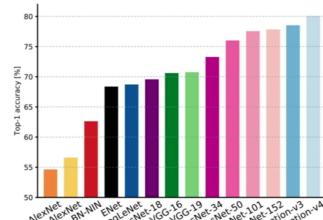


Training ResNet in practice:

- 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

Comparing complexity

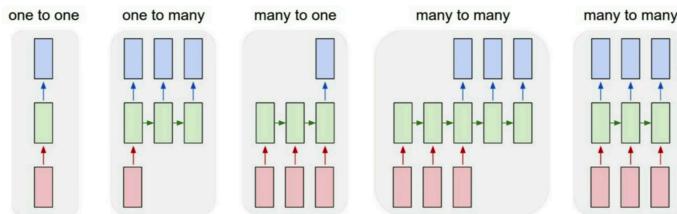
Comparing complexity...



Summary: CNN Architectures

- VGG, GoogLeNet, ResNet all in wide use, available in model zoos
- ResNet current best default
- Trend towards extremely deep networks
- Significant research centers around design of layer / skip connections and improving gradient flow
- Even more recent trend towards examining necessity of depth vs. width and residual connections
- Next time: Recurrent neural networks

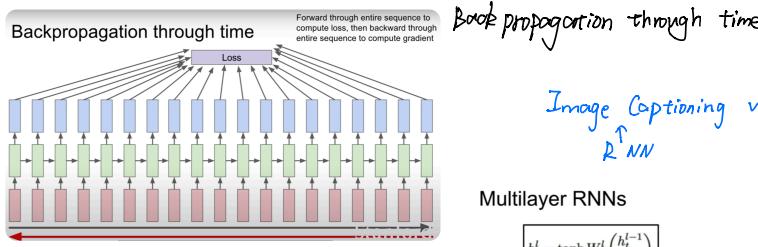
Lecture 10 Recurrent Neural Network → Process Sequences



$$h_t = f_w(h_{t-1}, x_t)$$

$$\begin{aligned} & y \\ \uparrow & \\ \text{RNN} & \geq \\ \uparrow & \\ x & \end{aligned}$$

$$\begin{aligned} h_t &= f_w(h_{t-1}, x_t) \\ h_t &= \tanh(W_h h_{t-1} + W_x x_t) \\ y_t &= W_y h_t \end{aligned}$$



Back propagation through time

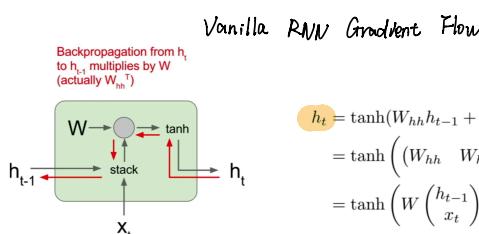
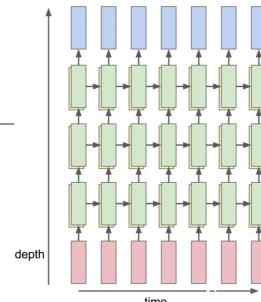
Image Captioning with Attention.
RNN

Multilayer RNNs

$$\begin{aligned} h_t^l &= \tanh W^l \begin{pmatrix} h_{t-1}^{l-1} \\ h_t^{l-1} \end{pmatrix} \\ h &\in \mathbb{R}^n. \quad W^l [n \times 2n] \end{aligned}$$

LSTM:

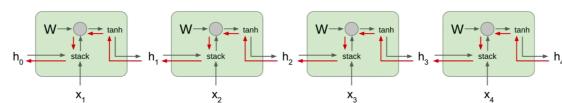
$$\begin{aligned} \begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} &= \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_{t-1} \\ h_t \end{pmatrix} \\ c_t^l &= f \odot c_{t-1}^l + i \odot g \\ h_t^l &= o \odot \tanh(c_t^l) \end{aligned}$$



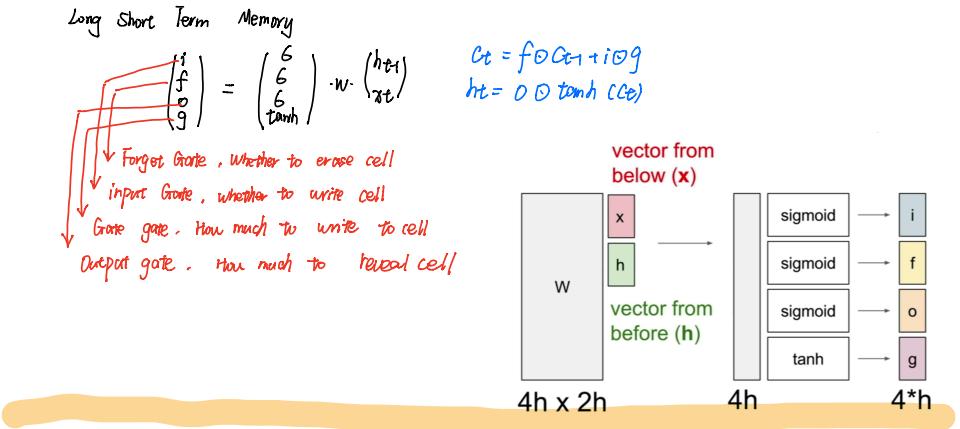
$$\begin{aligned} h_t &= \tanh(W_{hh} h_{t-1} + W_{xh} x_t) \\ &= \tanh((W_{hh} \quad W_{xh})(h_{t-1} \quad x_t)) \\ &= \tanh(W(h_{t-1} \quad x_t)) \end{aligned}$$

Vanilla RNN Gradient Flow

Bengio et al., "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al., "On the difficulty of training recurrent neural networks", ICML 2013



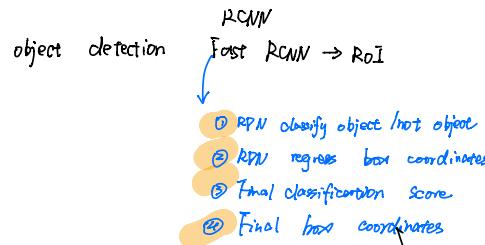
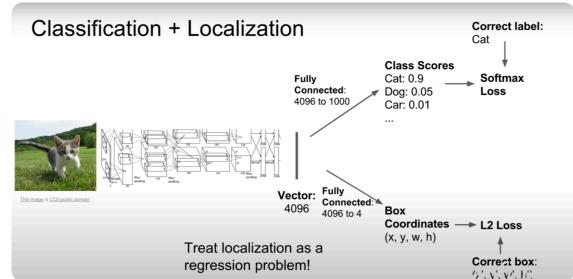
Computing gradient
of h_0 involves many
factors of W
(and repeated tanh)



Lecture 11 Detection and Segmentation

Unpooling $2 \times 2 \rightarrow 4 \times 4$
 Max pooling \rightarrow Max Unpooling
 Remember which element was max
 Use positions from pooling layer

Transpose Convolution
 Strided convolution



Detection without Proposals: YOLO / SSD

with each grid cell

1. Regress from each of the B boxes to a final box with 5 numbers ($bx, by, bw, bh, confidence$)
2. Predict scores for each of C classes (including background as a class)

\Rightarrow output $H \times W \times (5 \times B + C)$

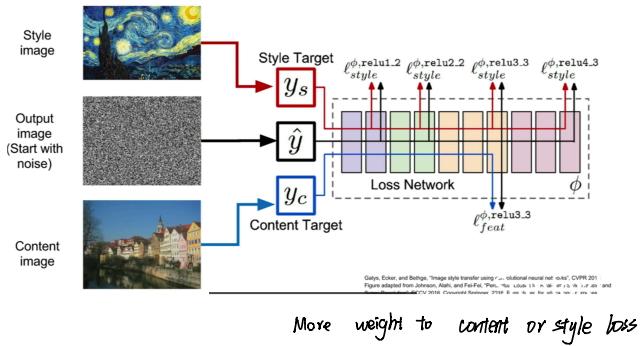
Object Detection: Lots of variables ...

Base Network	Object Detection architecture	Takeaways
VGG16	Faster R-CNN	Faster R-CNN is slower but more accurate
ResNet-101	R-FCN	
Inception V2	SSD	
Inception V3		SSD is much faster but not as accurate
Inception ResNet MobileNet	Image Size # Region Proposals	
...		

Lecture 12 Visualizing and Understanding

Visualizing CNN is an important task \Rightarrow Gradient Ascent

Neural Texture Synthesis



Style transfer very slow due to many forward/backward pass through VGG

Summary

Many methods for understanding CNN representations

Activations: Nearest neighbors, Dimensionality reduction, maximal patches, occlusion

Gradients: Saliency maps, class visualization, fooling images, feature inversion

Fun: DeepDream, Style Transfer.

Lecture 13 Generative Models

Supervised Learning

Data (x, y)

Goal: Learn a function to map $x \rightarrow y$

Example: Classification, regression, object detection, semantic segmentation, image captioning

Unsupervised Learning

Data: x

Goal: Learn some underlying hidden structures

Example: Clustering, dimensionality reduction, feature learning, density estimation

首先区分生成/判别方法和生成/判别模型。

有监督机器学习方法可以分为生成方法和判别方法（常见的生成方法有LDA主题模型、朴素贝叶斯算法和隐式马尔科夫模型等，常见的判别方法有SVM、LR等），生成方法学习出的是生成模型，判别方法学习出的是判别模型。

Generative Models

Given training data, generate new samples from same distribution



Training data $\sim p_{\text{data}}(x)$



Generated samples $\sim p_{\text{model}}(x)$

Want to learn $p_{\text{model}}(x)$ similar to $p_{\text{data}}(x)$

Addresses density estimation, a core problem in unsupervised learning \star

Several flavors:

- Explicit density estimation: explicitly define and solve for $p_{\text{model}}(x)$
- Implicit density estimation: learn model that can sample from $p_{\text{model}}(x)$ via explicitly defining it

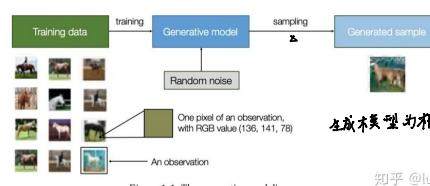
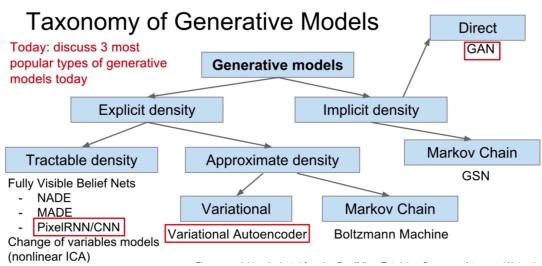


Figure 1-1. The generative modeling process

知乎 @luyan

判别模型: $p(y|x)$ 即给定观测 x 得到 y 的概率。

生成模型: $p(x)$ 即观测 x 出现的概率。



Generative Model: Pixel RNN, WaveNet, Video Pixel Networks

Fully visible belief network

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

likelihood of image x
 probability of i th pixel
 value given all previous pixels

Pixel RNN Generate image pixels starting from corner
 Dependency on previous pixels modeled using an RNN

Pixel CNN Generate image pixels starting from corner
 Dependency on previous pixels modeled using an CNN over context region
 Softmax loss at each pixel

PixelRNN and PixelCNN

Pros:

- Can explicitly compute likelihood $p(x)$
- Explicit likelihood of training data gives good evaluation metric
- Good samples

Con:

- Sequential generation => slow

Improving PixelCNN performance

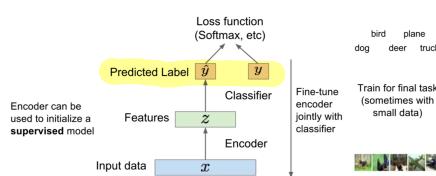
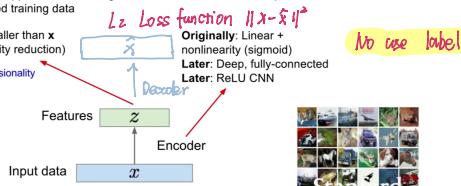
- Gated convolutional layers
- Short-cut connections
- Discretized logistic loss
- Multi-scale
- Training tricks
- Etc...

See

- Van der Oord et al. NIPS 2016
- Salimans et al. 2017 (PixelCNN++)

Some background first: Autoencoders

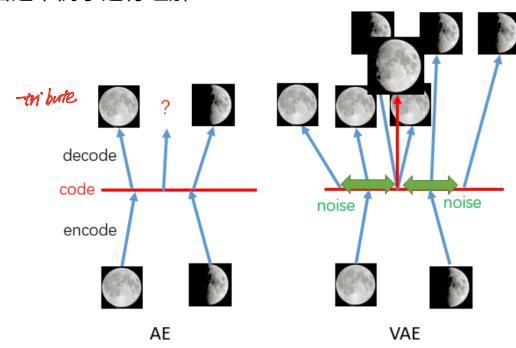
Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data



Variational Autoencoders

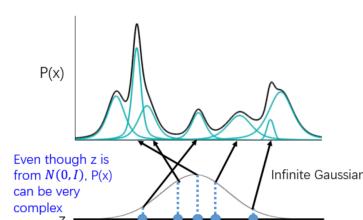
Data likelihood $P(x) = \int p(x|z) p(z) dz$
 VAE 是一个生成模型，其最终权生成概率分布 $P(x)$, 而不是直接数据。

VAE 变分自动编码器作为 AE 的变体，它主要的变动是对编码(code)的生成上。编码(code)不再像 AE 中是唯一映射的，而是具有某种分布，使得编码(code)在某范围内波动时都可产生对应输出。借助下面这个例子进行理解：

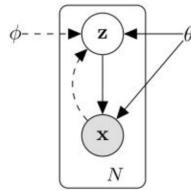


VAE 通过高斯混合模型 (Gaussian Mixture Model) 来生成 $P(x)$ ，也就是说 $P(x)$ 由一系列高斯分布叠加而成，每一个高斯分布都有自己的参数 μ 和 σ

$\{z \sim N(0, I)\}$ z is a vector from normal distribution
 $z|B \sim N(\mu_B, \sigma_B)$ Each dimension of z represents an attribute

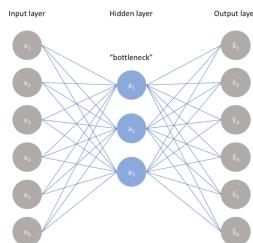


$$\begin{aligned}
x &= \{x_1, x_2, \dots, x_n\} \\
p(x) &= \prod_x p(x_i) \cdot p(x) \\
&= \prod_x N(0, I) \prod_x p(x) = N(0, I) \\
&= N(0, I) \prod_x p(x) = N(0, I) \\
&\quad \text{VAE 框架图} \\
&\quad \text{VAE 是分布的映射关系 即 } D_x \rightarrow D_z \\
&\quad \text{A 是单道映射关系} \\
k &\in [N(\mu_k), \Sigma_k] \quad \text{from } N(0, I) \\
&\quad \text{Encoder } Q(z|x) \\
&\quad \text{Decoder } P(x|z) \\
&\quad \text{sample } z \text{ from } N(0, I) \\
&\quad \mathcal{L} = -\log(p(x)) = \sum_k q(z|x) \log(p(x)) \\
&\quad = \sum_k q(z|x) \log\left(\frac{p(x|z)}{p(z)}\right)
\end{aligned}$$



VAE的图模型。我们能观测到的数据是 $p(x)$ ，而 $p(x)$ 由隐变量 z 产生，由 $p(x|z)$ 是生成模型 $p(x|z)$ ，从自编码器（auto-encoder）的角度来看，就是解码器；而由 $q(z|x)$ 是识别模型（recognition model） $q(z|x)$ ，类似于自编码器的编码器。

Auto Encoder (异常检测, 数据来源, 数据降维, 图像修复, 信息检索)



(1) 压缩编码的数据维度一定要比元时输入数据更少。

(2) 不管是编码器，还是解码器，本质上都是神经网络层，神经网络层一定要具有一定的“容量”。

最小重构误差 $\mathcal{L}(x, \hat{x})$

常见的 AE

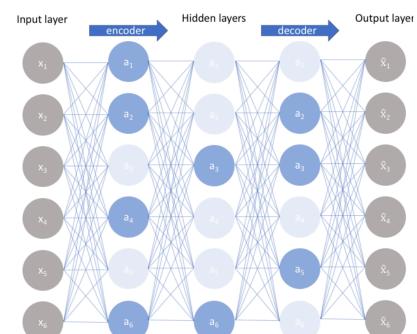
① 欠完备的自编码器 — 非通一般自编码器

欠完备自编码器没有明确的正则化项，只是根据重构损失来训练我们的模型。因此，确保模型不记忆输入数据的唯一方法就是确保我们已经充分限制了隐藏层中的节点数量。

② 稀疏自编码器 — 隐层选择性的激活

稀疏自编码器为我们提供了一种不需要减少我们隐藏层的节点数量，就可以引入信息瓶颈的方法。相反，我们将构造我们的损失函数，以惩罚层中的激活。对于任何给定的观察，我们都会鼓励我们的网络学习只依赖激活少量神经元进行编码和解码。值得注意的是，这是一种比较特殊的正则化实现方法，因为我们通常调整网络的权重，而不是激活。

Two approaches to enhance 稀疏性约束



L1正则化：我们可以添加一个对损失函数的正则化项，在 h 层中为观察 i 惩罚激活 a 的向量值的绝对值，使用微调参数 λ 进行缩放。 $\mathcal{L}(x, \hat{x}) + \lambda \sum_i |a_i^{(h)}|$

KL -散度（相对熵）：本质上，KL散度是两个概率分布差异的度量。我们可以定义一个参数 ρ 稀疏，它表示一个神经元在样本集合上的平均激活。这种期望可以计算为

$$\hat{\rho}_j = \frac{1}{m} \sum_i [a_i^{(h)}(x)]$$

下标 j 表示表示层 h 中特定的神经元，对 m 个训练观察的表征 x 的激活求和。本质上，通过限制一个神经元在样本集合上的平均激活，我们鼓励神经元只对观测的一个子集进行激活。我们可以将 ρ 描述为一个伯努利随机变量分布，我们可以利用KL散度（下展开）来比较理想的分布在所有隐藏层节点上的观察分布。

$$\mathcal{L}(x, \hat{x}) + \sum_j KL(p||\hat{\rho}_j)$$

③ 降噪自编码器 \Rightarrow 多数据添加噪音

采用这种方法，我们的模型不能简单地开发一个记忆训练数据的映射，因为我们的输入和目标输出不再相同。更确切的说，该模型学习矢量场以将输入数据映射到较低维流形；如果这个流形精确地描述了自然数据，我们就有效地“消除”了多余的噪音。

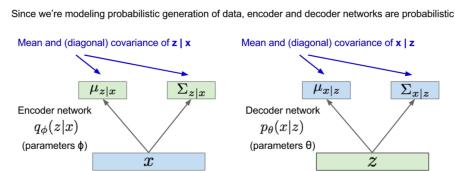
④ 压缩自编码器

人们会期望对于非常相似的输入，学习的编码也会非常相似。我们可以为此训练我们的模型，以便通过要求隐藏层激活的导数相对于输入而言很小。换句话说，对于输入比较小的改动，我们仍然应该保持一个非常类似的编码状态。这与降噪自编码器相似，因为输入的小扰动本质上被认为是噪音，并且我们希望我们的模型对噪音具有很强的鲁棒性。“降噪自编码器使重构函数（解码器）抵抗输入有限小的扰动，而压缩自编码器使特征提取函数（编码器）抵抗输入无限小的扰动。”

$$VAE - Probabilistic spin on autoencoders$$

$$\text{Data likelihood: } P_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$$

Variational Autoencoders



GAN: Generator DISCriminator

$$\min_D \max_G V(D, G) = \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log (1 - D(G(z)))]$$

Variational Autoencoders

Probabilistic spin to traditional autoencoders => allows generating data
Defines an intractable density => derive and optimize a (variational) lower bound

Pros:

- Principled approach to generative models
- Allows inference of $q(z|x)$, can be useful feature representation for other tasks

Cons:

- Maximizes lower bound of likelihood: okay, but not as good evaluation as PixelRNN/PixelCNN
- Samples blurrier and lower quality compared to state-of-the-art (GANs)

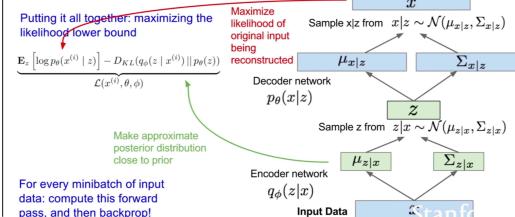
Active areas of research:

- More flexible approximations, e.g. richer approximate posterior instead of diagonal Gaussian
- Incorporating structure in latent variables

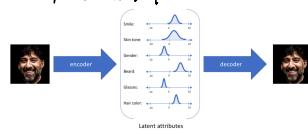
Training GANs

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

Variational Autoencoders



VAE 将每个给定的隐属性视为概率分布，当从隐状态解码时，我们将从每个隐状态分布中随机采样来生成向量化的解码器的输入。





Discriminator network: try to distinguish between real and fake images.



整个式子由两项构成。 x 表示真实图片， z 表示输入G网络的噪声，而 $G(z)$ 表示G网络生成的图片。

Train jointly in minimax game

$$\min_{\theta_D} \left[\mathbb{E}_{x \sim p_{data}} \log D(x) + \mathbb{E}_{z \sim p_{z}} \log (1 - D(G(z))) \right]$$

Gradient on discriminator ↴

Gradient on generator ↑

$D(x)$ 表示D网络判断真实图片是否真实的概率

(因为 x 就是真实的，所以对于D来说，这个值越接近1越好)。而 $D(G(z))$ 是D网络判断G生成的图片的是否真实的概率。

G的目的：上面提到过， $D(G(z))$ 是D网络判断G生成的图片是否真实的概率，G应该希望自己生成的图片“越接近真实越好”。也就是说，G希望 $D(G(z))$ 尽可能得大，这时 $V(D, G)$ 会变小。

因此我们看到式子的最前面的记号是 \min_G 。

D的目的：D的能力越强， $D(x)$ 应该越大。

$D(G(z))$ 应该越小。这时 $V(D, G)$ 会变大。因此式子对于D来说是求最大(max_D)

Lecture 14 Reinforcement Learning

Markov Decision Process

- Mathematical formulation of the RL problem
- **Markov property:** Current state completely characterises the state of the world

Defined by: $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbb{P}, \gamma)$

\mathcal{S} : set of possible states

\mathcal{A} : set of possible actions

\mathcal{R} : distribution of reward given (state, action) pair

\mathbb{P} : transition probability i.e. distribution over next state given (state, action) pair

γ : discount factor

Markov Decision Process

- At time step $t=0$, environment samples initial state $s_0 \sim p(s_0)$
- Then, for $t=0$ until done:
 - Agent selects action a_t
 - Environment samples reward $r_t \sim R(\cdot | s_t, a_t)$
 - Environment samples next state $s_{t+1} \sim P(\cdot | s_t, a_t)$
 - Agent receives reward r_t and next state s_{t+1}
- A policy π is a function from S to A that specifies what action to take in each state
- **Objective:** find policy π^* that maximizes cumulative discounted reward: $\sum_{t \geq 0} \gamma^t r_t$

The optimal policy π^*

We want to find optimal policy π^* that maximizes the sum of rewards.

How do we handle the randomness (initial state, transition probability...)?
Maximize the **expected sum of rewards!**

$$\text{Formally: } \pi^* = \arg \max_{\pi} \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t | \pi \right] \text{ with } s_0 \sim p(s_0), a_t \sim \pi(\cdot | s_t), s_{t+1} \sim p(\cdot | s_t, a_t)$$

Definitions: Value function and Q-value function

$$V^\pi(s) = \mathbb{E} \left[\sum_{t \geq 0} r^t | s_0 = s, \pi \right] \quad (\text{How good is a state})$$

$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{t \geq 0} r^t | s_0 = s, a_0 = a, \pi \right] \quad (\text{How good is a state-action pair})$$

Reinforce algorithm

$$J(\theta) = \mathbb{E}_{\pi_\theta} [r_t | \theta] = \int r(t) p(t; \theta) dt$$

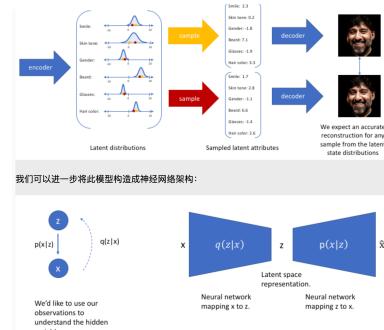
Q-value

policy Q-function

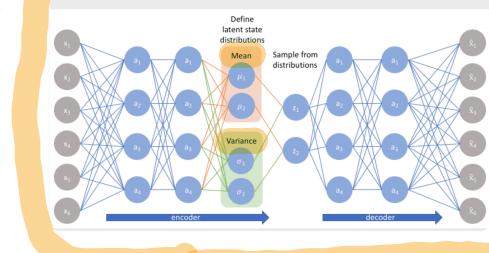
Action - critic Algorithm

定义：收获 [公式] 为在 t 个马尔科夫奖励链上从 t 时刻开始往后所有的奖励的有衰减的总和。也有翻译成“收益”或“回报”。公式如下：

通过构造我们的编码器来输出一系列可能的值（统计分布），然后随机采样该值作为解码器的输入，我们能够学习到一个连续，平滑的隐空间。因此，在隐空间中彼此相邻的值应该与非常类似的重建相对应。而从隐分布中采样到的任何样本，我们都希望解码器理解，并准确重构出来。



下图是VAE的结构图：



马尔科夫性 Markov Property

某一状态信息包含了所有相关的历史，只要当前状态可知，所有的历史信息都不再需要，当前状态就可以决定未来，则认为该状态具有马尔科夫性。

可以用下面的状态转移概率公式来描述马尔科夫性：

$$P_{ss'} = P[S_{t+1}=s' | S_t=s]$$

状态转移矩阵

$$P = \begin{bmatrix} p_{11} & \dots & p_{1n} \\ \vdots & \ddots & \vdots \\ p_{n1} & \dots & p_{nn} \end{bmatrix}$$

Markov Chain $\langle S, P \rangle$

衰减系数

Markov Reward Process $\langle S, P, R, \gamma \rangle$

奖励系数

Discount Factor 衰减系数

$\gamma \in [0, 1]$

Return

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

Actor-Critic Algorithm

```

Initialize policy parameters  $\theta$ , critic parameters  $\phi$ 
For iteration=1, 2, ... do
    Sample m trajectories under the current policy
     $\Delta\theta \leftarrow 0$ 
    For  $i=1, \dots, m$  do
        For  $t=1, \dots, T$  do
             $A_t = \sum_{t' \geq t} \gamma^{t'-t} r_t^i - V_\phi(s_t^i)$ 
             $\Delta\theta \leftarrow \Delta\theta + A_t \nabla_\theta \log(a_t^i | s_t^i)$ 
         $\Delta\phi \leftarrow \sum_t \sum_i \nabla_\phi ||A_t^i||^2$ 
         $\theta \leftarrow \alpha \Delta\theta$ 
         $\phi \leftarrow \beta \Delta\phi$ 
    End for

```

价值函数给出了某一状态或某一行为的长期价值。

$r \rightarrow 0$ 近期价值
 $r \rightarrow 1$ 远期价值

定义：一个马尔科夫奖励过程中某一状态的价值函数为从该状态开始的马尔可夫链收获的期望

$$Value\ Function\ 价值函数$$

$$V(s) = E[A_t | S_t=s]$$

收获是针对一个马尔科夫链中的某一个状态来说的。

Bellman Function - MRP

$$\begin{aligned} v(s) &= E[G_t | S_t=s] \\ &= E[R_{t+1} + \gamma R_{t+2} + \dots | S_t=s] \\ &= E[R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \dots) | S_t=s] \\ &= E[R_{t+1} + \gamma G_{t+1} | S_t=s] \\ &= E[R_{t+1} + \gamma V(S_{t+1}) | S_t=s] \end{aligned}$$

$V(s)$

通过方程可以看出 [公式] 由两部分组成，一是该状态的即时奖励期望，即时奖励期望等于即时奖励，因为根据即时奖励的定义，它与下一个状态无关；另一个是下一时刻状态的价值期望，可以根据下一时刻状态的概率分布得到其期望。如果用 s' 表示 s 状态下一时刻任一可能的状态，那么 Bellman 方程可以写成 $V(s) = R_s + \gamma \sum_{s' \in S} P_{ss'} V(s')$

$$V = R + \gamma PV$$

$$\begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix} = \begin{bmatrix} R_1 \\ \vdots \\ R_n \end{bmatrix} + \gamma \begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & \ddots & \vdots \\ P_{n1} & \dots & P_{nn} \end{bmatrix} \begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix}$$

$$V = R + \gamma PV$$

$$(I - \gamma P)V = R$$

$$V = C(I - \gamma P)^{-1} R$$

Markov Decision Process

行为集合 $A < s, a, p, r, \pi >$

$$P_{ss'}^\alpha = P[S_{t+1}=s' | S_t=s, A_t=a]$$

$$R_s^\alpha = E[R_{t+1} | S_t=s, A_t=a]$$

策略 π 是本概率集成或分布， $\pi(a|s)$ 为又过过程中某一状态采取可能行为 a 的概率

$$\pi(a|s) = P[A_t=a | S_t=s]$$

S_1, S_2, \dots

当给定一个MDP [公式] 和一个策略 π ，那么状态序列 [公式] 是一个马尔科夫过程 [公式]，同样，状态和奖励序列 [公式] 是一个马尔科夫奖励过程 [公式]，并且在这个奖励过程中满足下面两个方程：

$$P_{s,s'}^\pi = \sum_{a \in A} \pi(a|s) P_{ss'}^\pi$$

用文字描述是这样的，在执行策略 π 时，状态从 s 转移至 s' 的概率等于一系列概率的和，这一系列概率指的是在执行当前策略时，执行某一个行为的概率与该行为能使状态从 s 转移至 s' 的概率的乘积。

奖励函数表示如下：

$$R_s^\pi = \sum_{a \in A} \pi(a|s) R_s^\pi$$

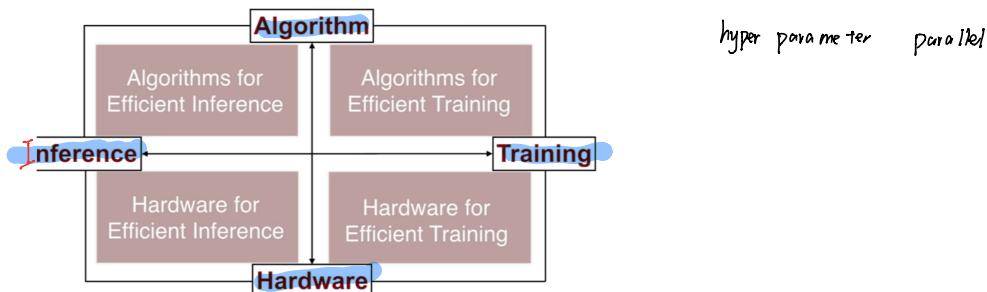
用文字表述是这样的：当前状态 s 下执行某一指定策略得到的即时奖励是该策略下所有可能行为得到的奖励与该行为发生的概率的乘积的和。

基于策略 π 的价值策略

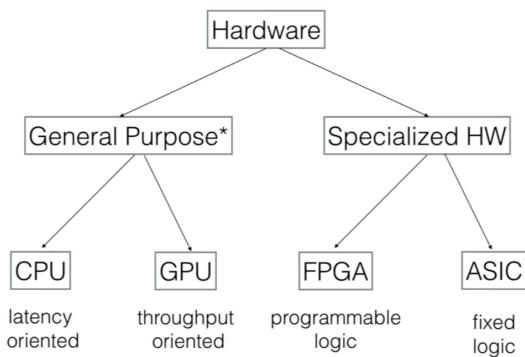
$$V_\pi(s) = E_\pi[G_t | S_t=s]$$

$$行为价值函数 Q_\pi(s,a) = E_\pi[G_t | S_t=s, A_t=a]$$

Agenda



Hardware 101: the Family



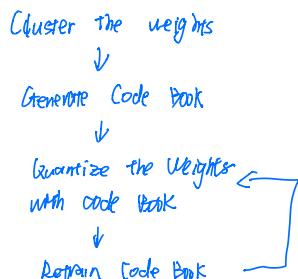
Part 1 Algorithm for Efficient Inference

① Pruning (剪枝)

pruning and retraining

② weight sharing

③ Q quantization (量化)



Lecture 1b Adversarial Examples on Adversarial Training

Adversarial Examples

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

Noise \star (sigmoid)

Human Brain



Map of Adversarial Cross Section

Attack easy, defense difficult

Ch23). END.

2020.07.04

老, 又, 好