

Machine Learning



Lecture 0-1



机器学习
深度学习
[神经网络]

You write the program for learning

Looking for a function from Data

$f(\text{Information}) = \text{result}$

Framework supervised learning



f^* testing

steps:

1. define a set of functions

2. goodness of function

3. pick the best one

Learning Map:

Semi-supervised Learning: 少量有标注数据

Transfer Learning: 大量Data 和 大量无关 Data

Unsupervised Learning: 无标注 Data

Reinforcement Learning: 不固定输出, 只给输出打分

Supervised Learning: 大量 Data learn from critics

Regression 回归预测

Linear Model

Deep Learning SVM, decision tree, K-NN...
Non-linear Model

Structured Learning A lot more

Classification

Binary Classification: Yes/No

Multi-class classification: class1, class2...

Lecture 1 Regression

Cases:

stock Market Forecast
Self-driving Car → direction

Recommendation

x_i : 下标表示输入的特征量

$$f(x_1, x_2, \dots, x_n) = y$$

Step 1: Model

$$y = b + W^T x_p \rightarrow \text{set of function parameter}$$

linear Model

$$y = b + \sum w_i x_i \rightarrow x_i: \text{an attribute of input}, w_i: \text{weight feature}$$

Step 2: Goodness of function

x^i : 用上标表示一个完整的 object 编号

$$f(x^i), \hat{y}^i$$
: \hat{y}^i 表示真实值

Loss function L: $L(f) = \text{how bad}$ 可以是

$$L(w, b) \text{ 一般 } \leq (f^i - (b + w^T x^i))^2$$

draw a graph of $L(w, b)$

b : 偏移量 - f , 用梯度体现 $L(f)$

pick the best f Gradient Descent 可微

$$f^* = \arg \min L(f)$$

Step 3: Gradient Descent: ex. $L(w)$

Randomly pick w^0

$$w^1 = w^0 - \eta \frac{dL}{dw}|_{w=w^0} \leftarrow \frac{dL}{dw} \text{ 越大, 跨度越大}$$

η : learning rate

$$w^2 = \dots$$

$$w^3 = \dots$$

local optical $\frac{dL}{dw} = 0$ not global optical
多参数情况: $L(w, b)$

Random Pick (w^0, b^0)

$$w^1 = w^0 - \eta \frac{\partial L}{\partial w}|_{(w^0, b^0)} = (w^0, b^0), b^1 = b^0 - \eta \frac{\partial L}{\partial b}|_{(w^0, b^0)} = (w^0, b^0)$$

$\nabla L = \begin{bmatrix} \frac{\partial L}{\partial w} \\ \frac{\partial L}{\partial b} \end{bmatrix}$ Generalization

care about the error on new data
error过大 → 修改 model (更高泛化) (testing data)

可能会对 training data 效果很好但与 testing data 效果较差

$$\text{error} = \sum |\hat{y}^i - y^i|$$

某模型可以降低 training data error rate
但对 testing data 却不一定 (Overfitting)

有可能设计的 model (set of function) 里并没有正确的 f
 $f(x, y) \neq f(x)$, can't fit testing data → redesign the model

$$y = b_1 + \delta(x_{1,1})w_1x_{1,1} + b_2 + \delta(x_{1,2})w_2x_{1,2}$$

(Back to step 1)

$$y = b + \sum w_i x_i$$

Back to step 2: Regularization

redefine loss function

$$L = \frac{1}{n} \left(\sum_{i=1}^n (y_i - (b + \sum w_i x_i))^2 + \lambda \sum (w_i)^2 \right)$$

why? (无需考虑 bias)

The functions with smaller w_i are better smooth

使得单一变量对结果产生的影响更小

output 对输入的变化更不敏感

受 noises 的影响较小

Select λ to obtain the best model

Conclusion:

Gradient descent

Overfitting & Regularization

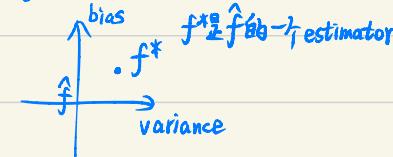
Validation

Lecture 2

error due to bias
variance

Estimator

$\hat{g} = \hat{f}(x)$ training data → \hat{f}^*



101: estimate \bar{x}

mean: M , variance: σ^2

estimate: N sample points: $\{x_1^n, x_2^n, \dots, x_N^n\}$

$$m = \frac{1}{N} \sum x^n = m + M \quad E(m) = E\left(\frac{1}{N} \sum x^n\right) = \frac{1}{N} \sum E(x^n) = M$$

$E(m) = M$: unbiased

$$\text{Var}[m] = \frac{\sigma^2}{N} \quad (\text{depends on the number of samples})$$

$$S^2 = \frac{1}{N} \sum (x_i - m)^2 \quad E(S^2) = \frac{N-1}{N} \sigma^2 + \sigma^2 \quad \text{biased estimator}$$

$$E(f^*) = \bar{f} \quad \bar{f} = \frac{1}{N} \sum f_i$$

每次试验选取
不同的 training data

Same model, different training data ⇒ diff f^*

Simple model → low variance & high bias vice versa
less influenced by the sampled data

但是 function set 更大, bias 偏大

variance 大: overfitting | bias 大: underfitting

能 fit training data | 不能 fit training data

不能 fit testing data | ↓
增加模型复杂度
(增加数据量并没有用)

增加 data (有效, 不实际)
(generate 更多的 training data)

自己制作 data

Regularization (smoothen)
相当于调整 function space, 可能增大 bias

Model Selection: trade-off between bias and variance
Should NOT do:

Training Set Testing Set

Training Set	Testing Set
Model 1	Err 0.9
Model 2	Err 0.7
Model 3	Err 0.5 pick X

Training Set ⊃ Validation Set

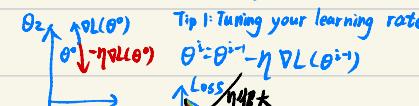
"kaggle"

N-fold Cross Validation 每次取出一部分作为 validation
将 training set 分割为 n 份

Lecture 3 Gradient Descent

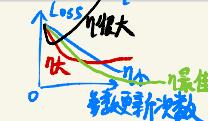
$$\theta^* = \arg \min_{\theta} L(\theta) \quad L: \text{loss function} \quad \theta: \text{parameter}$$

$$\nabla L(\theta) = \begin{bmatrix} \frac{\partial L}{\partial \theta_0} \\ \vdots \\ \frac{\partial L}{\partial \theta_n} \end{bmatrix} \quad \theta' = \theta^* - \eta \nabla L(\theta^*)$$



Tip 1: Tuning your learning rate

$$\theta^{(t+1)} = \theta^{(t)} - \eta \nabla L(\theta^{(t)})$$



Adaptive Learning Rates

Popular & simple Idea: $\downarrow \eta$ by some factor
开始: η 大, 一段时间后适当 $\downarrow \eta$

给不同的 parameter 设置不同的 η

Adagrad:

普通 Gradient Descent: $w^{t+1} \leftarrow w^t - \eta_t g^t$

Adagrad: $w^{t+1} \leftarrow w^t - \frac{\eta_t}{\sigma^t} g^t \quad \eta_t = \frac{\eta}{\sqrt{w^t}}$

$$w^t \in w^0 - \frac{\eta}{\sigma^0} g^0 \quad \sigma^0 = \sqrt{(g^0)^2}$$

$$w^2 \in w^1 - \frac{\eta^1}{\sigma^1} g^1 \quad \sigma^1 = \sqrt{\frac{1}{2}[(g^0)^2 + (g^1)^2]}$$

$$\vdots \quad \vdots$$

$$w^{t+1} \leftarrow w^t - \frac{\eta^t}{\sigma^t} g^t \quad \sigma^t = \sqrt{\frac{1}{t+1} \sum_{i=0}^t (g^i)^2}$$

$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sqrt{\sum_{i=0}^t (g^i)^2}} g^t$$

Adagrad 考虑 grad 的变化

adagrad 能自动选择最优路径

How surprise it is 反差

$\sqrt{\sum_i (g^i)^2} \leftarrow$ 造成反差的效果

$$\text{best step: } \frac{|y''|}{|y'|} \quad (\text{if } y = ax^2 + bx + c)$$

$$\sqrt{\frac{t}{t+1} (g^t)^2} \rightarrow \text{反映 } y'' \text{ 的大小}$$

Tip2: Stochastic Gradient Descent

选一个 example x^n

$$L^n = (y^n - (b + \sum m_i x_i^n))^2 \text{ 仅采用一个数计算 } L$$

Tip3: Feature Scaling

使不同 feature 的 Scale 是一样的

$$(\odot) \rightarrow (\odot)$$

对 $x_1^1, x_1^2, x_1^3, \dots$ 进行类似随机变量

标准化的操作: $\frac{x - \mu}{\sigma}$

Limitation of Gradient Descent

1. Local Min

2. Settle point (在驻点附近更新速率)

Lecture 4 Classification

$x \rightarrow f \rightarrow \text{class } n$

Classification as Regression? X

Ideal Alternatives

Function (Model): $x \rightarrow$

$f(x)$
$g_1(x) \rightarrow \text{class 1}$
$g_2(x) \rightarrow \text{class 2}$

Loss function: f 判断错误的次数

$$L(f) = \sum_n \delta(f(x^n) + \hat{y}^n)$$

Find the best function:

Perceptron, SVM...

贝叶斯公式 (从概率角度思考)

在已知抽到 x 的情况下, x 最可能来自哪个 class

Class 1	Class 2
<input type="checkbox"/>	<input type="checkbox"/>
$p(c_1 x)$	$p(c_2 x)$
$p(x c_1)$	$p(x c_2)$

$$p(c_1|x) = \frac{p(x|c_1)p(c_1)}{p(x|c_1)p(c_1) + p(x|c_2)p(c_2)}$$

$$\Rightarrow p(c_2|x) = \frac{p(x|c_2)p(c_2)}{p(x|c_1)p(c_1) + p(x|c_2)p(c_2)}$$

Generative Model \rightarrow generate x

$$p(x) = p(x|c_1)p(c_1) + p(x|c_2)p(c_2)$$

$p(c_1), p(c_2)$: prior 可直接由 training data 计算

$p(x|c_1) = ? \quad p(x|c_2) = ?$ 求 c_1, c_2 中数据的分布律

每一个 x 由一向量描述

Gaussian Distribution

输入 x, 输出取到 x 的概率 (密度)

假设所有点是从 Gaussian 分布中被选出来的 \rightarrow 求高斯分布
求 μ, Σ : Maximum Likelihood

$$\text{Likelihood: } L(\mu, \Sigma) = f_{\mu, \Sigma}(x^1) f_{\mu, \Sigma}(x^2) \dots f_{\mu, \Sigma}(x^n)$$

取 $L(\mu, \Sigma)$ 最大的 μ, Σ (μ^*, Σ^*)

$$\mu^* = \frac{1}{N} \sum x^i, \quad \Sigma^* = \frac{1}{N} \sum (x^i - \mu^*)^T$$

$$[] \quad [] \rightarrow []$$

为减少参数的个数, 让多个 class 共用一个 Σ

$$L(\mu^*, \mu^*, \Sigma) = f_{\mu^*, \Sigma}(x^1) f_{\mu^*, \Sigma}(x^2) \dots f_{\mu^*, \Sigma}(x^n) \\ f_{\mu^*, \Sigma}(x^{n+1}) \dots f_{\mu^*, \Sigma}(x^{140})$$

μ^* 和 Σ^* 计算方法相同, 根据数据量加权: $\Sigma = \frac{75}{140} \Sigma^1 + \frac{61}{140} \Sigma^2$

3 steps:

Function Set (Model)

$$f: \mathbf{x} \rightarrow \begin{cases} C_1 & P(C_1|\mathbf{x}) > 0.5 \\ C_2 & \text{other} \end{cases}$$

Goodness of function:

μ, Σ 使得 likelihood 最大

Pick the best function

可任意选择概率分布(例子里是 Gaussian)

假设 $x_1, x_2, x_3, \dots, x_k$ 是相互独立的(有的时候不能这么假设)

$$P(\mathbf{x}|C_1) = p(x_1|C_1)p(x_2|C_1)\cdots p(x_k|C_1)$$

Naive Bayes Classifier

$$P(C_1|\mathbf{x}) = \frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_1)p(C_1) + p(\mathbf{x}|C_2)p(C_2)} = \frac{1}{1 + \frac{p(x_1|C_2)p(C_2)}{p(x_1|C_1)p(C_1)}}$$

$$z = \ln \frac{p(x|C_1)p(C_1)}{p(x|C_2)p(C_2)}, p(C_i|\mathbf{x}) = \frac{1}{1+e^{-z}} = \sigma(z)$$

$$z = \frac{(w^T - b)^T \Sigma^{-1} \mathbf{x} - b}{w^T} \Rightarrow p(C_1|\mathbf{x}) = \sigma(w^T \mathbf{x} + b), \text{ 只需找 } w \text{ 和 } b$$

Lecture 5 Logistic Regression

$$\text{找 } P_{w,b}(C_1|\mathbf{x}) = \sigma(z) \quad z = w^T \mathbf{x} + b = \sum_i w_i x_i + b$$

$$f_{w,b}(\mathbf{x}) = P_{w,b}(C_1|\mathbf{x})$$

$$\text{Step 1: } f_{w,b}(\mathbf{x}) = \sigma\left(\sum_i w_i x_i + b\right) \in (0,1)$$

Step 2: 假设数据是由 $f_{w,b}(\mathbf{x})$ 产生的

对于一组 w, b , 计算样本是由 $f_{w,b}$ 产生的可能性

$$L(w, b) = f_{w,b}(\mathbf{x}^1)f_{w,b}(\mathbf{x}^2)(1-f_{w,b}(\mathbf{x}^3))\dots$$

$$w^*, b^* = \arg \max_{w, b} L(w, b)$$

$$= w^*, b^* = \arg \min_{w, b} -\ln L(w, b)$$

$$\ln L(w, b) = \ln f_{w,b}(\mathbf{x}^1) + \ln f_{w,b}(\mathbf{x}^2) + \ln(1-f_{w,b}(\mathbf{x}^3)) + \dots$$

$$\mathbf{x}^1 \rightarrow \hat{y}^1 = 1 \quad \mathbf{x}^2 \rightarrow \hat{y}^2 = 1 \quad \mathbf{x}^3 \rightarrow \hat{y}^3 = 0$$

$$\ln f_{w,b}(\mathbf{x}^i) = \hat{y}^i \ln f(\mathbf{x}^i) + (1-\hat{y}^i) \ln(1-f(\mathbf{x}^i))$$

$$-\ln L(w, b) = \sum_n [\hat{y}^n \ln f(\mathbf{x}^n) + (1-\hat{y}^n) \ln(1-f(\mathbf{x}^n))]$$

$$L(f) = \sum_n C(f(\mathbf{x}^n), \hat{y}^n)$$

Loss, 与上面的 L 不同

$$C(f(\mathbf{x}^n), \hat{y}^n) = -[\hat{y}^n \ln f(\mathbf{x}^n) + (1-\hat{y}^n) \ln(1-f(\mathbf{x}^n))]$$

使 function 的 output 尽可能与 \hat{y} 接近

不能用线性回归里的方差

Step 3: find the best function

$$w_i \leftarrow w_i - \eta \sum_n (\hat{y}^n - f_{w,b}(\mathbf{x}^n)) \mathbf{x}_i^n$$

(与线性回归相同) $\hat{y} = \sigma(w^T \mathbf{x})$

Logistic Regression + Square Error?

会出现离最优解很远但 $L(f) \approx 0$ 的情况

Discriminative vs Generative

相同的 Model, 不同的假设

贝叶斯潜在问题: $P(C|x)$ 最大的 Class 并不一定是正确的 Class (可能是 training data 中该 class 占比很大)

Generative Model 做了更多的假设 (针对少量 data)
多 data 影响较小 (针对带 noise 的 data)

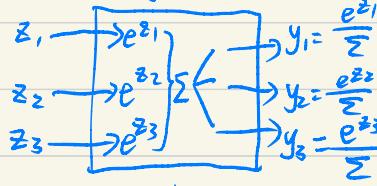
Multi-class Classification

$$C_1: w_1^T b_1 \quad z_1 = w_1^T x + b_1$$

$$C_2: w_2^T b_2 \quad z_2 = w_2^T x + b_2$$

$$C_3: w_3^T b_3 \quad z_3 = w_3^T x + b_3$$

Softmax



Normalization

对最大值做强化 $y_i = P(C_i|x)$

计算 y 与 \hat{y} 的 cross entropy $\sum -\hat{y}_i \ln y_i$

$$y = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

Logistic Regression 的局限性 (分界线是线性的)

feature transformation (很难找)

Cascading logistic regression models

$$x_1 \rightarrow z_1 \rightarrow x'_1$$

$$x_2 \rightarrow z_2 \rightarrow x'_2 \rightarrow \dots \rightarrow y$$

将 feature transformation 当作线性变换

通过 Logistic Regression 寻找线性变换

某一个 Logistic 的各个输入可以是其它 logistic 的输出, 其输出也可以是下一级 logistic 的输入

$$\begin{aligned} x_1 &\rightarrow \sigma(z_1) \\ x_2 &\rightarrow \sigma(z_2) \end{aligned}$$

$$\rightarrow \sigma(z_1) + \sigma(z_2) \rightarrow \sigma(z) \rightarrow y$$

neuron
Neural Network

Lecture 6 Deep Learning Introduction

Perceptron (Linear)

Perceptron Limits

Multi-layer perceptron (与现在的 DNN 差不多)

Back propagation

通常 3 个以上 hidden layer 没有太大用处

1 个 hidden layer "good enough"

RBM initialization (breakthrough) 实际上没作用

GPU 加速

Speech recognition

ILSVRC image competition

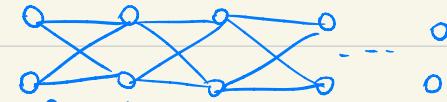
Step 1: Define function set (Neutral Network)

θ : 所有的 weight & bias

连接方式决定 Neutral Network 的结构

连接方式:

1. Fully Connect Feedforward Network



当 function 看, 输入向量, 输出向量

由 Neural Network 的结构确定一个 function set

Deep: Many hidden layers

用矩阵操作来表示 现在较少会用 sigmoid

$$\sigma \left(\begin{bmatrix} \cdot & \cdot & \cdot \end{bmatrix} \begin{bmatrix} \cdot & \cdot & \cdot \end{bmatrix} + \begin{bmatrix} \cdot & \cdot & \cdot \end{bmatrix} \right) = \begin{bmatrix} \cdot & \cdot & \cdot \end{bmatrix} \text{ 不一定是 sigmoid}$$

写成矩阵运算的形式由 GPU 计算

可进化神经网络? GA 来确定神经网络结构? 用的不多

2. Convolutional Neural Network (CNN)

一种神经网络的连接方式

$$\text{total loss } L = \sum_{n=1}^N C^n \quad C(y, \hat{y}) = -\sum_i y_i \ln \hat{y}_i$$

使 L 最小

Gradient Descent

Backpropagation: 在 NN 中有梯度下降法

Universality Theorem

连续函数 $f: \mathbb{R}^N \rightarrow \mathbb{R}^M$ (任意)

可用带 1 个 hidden layer 的 NN 实现
(在 hidden layer 中用足够多的 neuron)

Lecture 8 "Hello World"

Keras MNIST

```

model = Sequential()
model.add(Dense(input_dim=28*28, output_dim=500))
model.add(Activation('sigmoid'))
model.add(Dense(output_dim=500))
model.add(Activation('sigmoid'))
model.add(Dense(output_dim=10))
model.add(Activation('softmax'))
model.compile(loss='categorical_crossentropy',
              optimizer='adam',
              metrics=['accuracy'])
model.fit(x_train, y_train, batch_size=100, nb_epoch=20)

```

↑ ↑
numpy array numpy array

$[10000, 784]$ $[10000, 10]$

score = model.evaluate(x_test, y_test)

score[0]: total loss score[1]: accuracy

result = model.predict(x_test)

Lecture 7 Backpropagation

为了高效计算 Gradient (很多维)

本质上是 Chain Rule

$$L(\theta) = \sum_{n=1}^N C_n(\theta) \quad \frac{\partial L(\theta)}{\partial \theta} = \sum_{n=1}^N \frac{\partial C_n(\theta)}{\partial \theta}$$

$$\frac{\partial C}{\partial w} = \frac{\partial C}{\partial z} \frac{\partial z}{\partial w}$$

$\frac{\partial z}{\partial w}$: forward pass 好求

$\frac{\partial C}{\partial z}$: backward pass 不好求

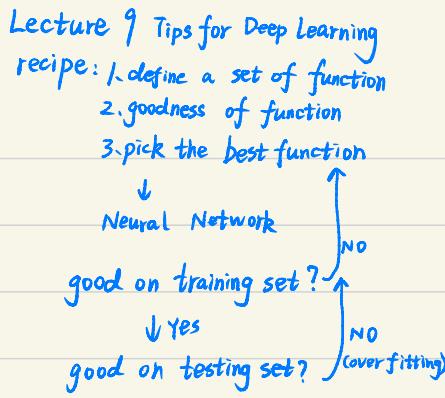
forward pass: $\frac{\partial z}{\partial w} = x_i (z = z_w; x_i + b)$

backward pass: $\frac{\partial C}{\partial z}$

$$\frac{\partial C}{\partial z} = \frac{\partial C}{\partial z} \frac{\partial z}{\partial z} + \frac{\partial C}{\partial z''} \frac{\partial z''}{\partial z} \dots$$

$$= \sum w \frac{\partial C}{\partial z} \text{ 再向下一层...}$$

实际上是从 output 层开始算, 反向传播
建立反向的 NN 用于 backward pass



一般第一个问题是 over fitting

training 结果不好：

1. New activation function
2. Adaptive Learning Rate

testing 结果不好：

1. Early Stopping
2. Regularization
3. Dropout

Vanishing Gradient Problem:

在靠近 output 层梯度更新快，靠近 input 层梯度慢

ΔW 的变化在由 Sigmoid 处理后变得很小 \rightarrow 改用 ReLU

Leaky ReLU Parametric ReLU

Maxout (自动学习 Activation)

RMSProp

$$w' \leftarrow w - \frac{\eta}{\sigma^2} g^2 \quad \sigma^2 = g^2$$

$$w^2 \leftarrow w^2 - \frac{\eta}{\sigma^2} g^2 \quad \sigma^2 = \sqrt{\alpha(\sigma^2)^2 + (1-\alpha)(g^2)^2}$$

$$w^3 \leftarrow w^2 - \frac{\eta}{\sigma^2} g^2 \quad \sigma^2 = \sqrt{\alpha(\sigma^2)^2 + (1-\alpha)(g^2)^2}$$

$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sigma^2} g^t \quad \sigma^2 = \sqrt{\alpha(\sigma^2)^2 + (1-\alpha)(g^t)^2}$$

(更倾向于相信新的 g , 前面的 g 的影响减弱)

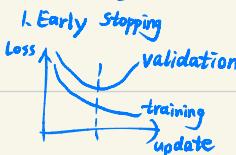
Momentum: 加上前一次移动的方向

$$V^0 = 0, V^t = \lambda V^0 - \eta \nabla L(\theta^0)$$

$$V^t = \lambda V^t - \eta \nabla L(\theta^t) \dots$$

adam \rightarrow momentum + RMSprop

在 testing 上结果不好:



1. Early Stopping

Weight Decay (不常用)

3. dropout

每次参数 update 之前:

每个 neuron 有 $p\%$ 的几率被丢掉

\hookrightarrow network 结构调整

用新的 network 训练

testing 时不 dropout:

如果 training 的 dropout rate 是 $p\%$, 则所有 weight 增上 $(1-p)\%$

Keras Demo

`model.add(Dropout(0.7))`

Lecture 10 CNN

CNN：人为简化后的NN

大部分的pattern比整个image小，不同区域会有相同图像（Conv）

Connecting to small region with less parameters

subsampling：拿掉image的奇数行和偶数列（MaxPooling）

→使image变小，减少参数

CNN：图像 \rightarrow convolution \rightarrow max pooling

重复多次

flatten
↓
DNN

用filter来识别图像的一部分



一个filter生成一个更小的pattern

用Maxpooling进行subsampling

多filter \rightarrow feature map

采样：RGB \rightarrow 对应的filter也有3层

feature map \rightarrow MaxPooling \rightarrow 取每块最大值

weight sharing 得到一个新的image

Flatten：将feature map拉直

输入三维(RGB)的vector(tensor)

RGB:3

keras: model.add(Convolution2D(25, 3, 3, input_shape=(1, 28, 28)))

model.add(MaxPooling2D((2, 2)))

model.add(Convolution2D(50, 3, 3))

model.add(MaxPooling2D((2, 2)))

model.add(Flatten())

第k层

分析CNN：第k层filter被激活的程度： $\Sigma \sum a_{ij}^k$

找一个image使第k层的activation最大

$x^* = \arg \max x^k$ (gradient ascent)

让classifier最active的可能完全不像那个class

Deep Dream Deep Style

什么时候使用CNN：

1. image中的pattern很小

2. pattern在image的不同区域出现

3. subsampling后不会改变object

Lecture 11 Why deep?



or



Deep \rightarrow Modularization

module: detect attribute

Speech

Phoneme 发音受前后的影响

Tri-phone

State: Tied-state

End-to-end Learning

Production line

Lecture 12 Semi-Supervised Learning

training data 有部分 unlabeled (一般有很多)

unlabeled data 的分布能帮助学习(带有一定假设)

Generative Model

Low-density Separation :不同的 class 间有明显的分界

Self-training

labeled data $\rightarrow f^*$ \rightarrow 给未标记的数据标记

\rightarrow Pseudo-label (选一些 label \rightarrow training)

Entropy-based Regularization

假设 x 的分布集中在某些区域



Cluster and then label

Graph-based smoothness

$$S = \frac{1}{2} \sum w_{ij} (y_i - y_j)^2 = y^T L y \quad L: \text{Graph Laplacian}$$

$$L = D - W$$

Better Representation

Lecture 13 Unsupervised Learning

Reduction 化繁为简

Generation 无中生有

Clustering

k-means

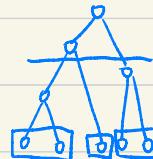
$$X = \{x^1, x^2, \dots, x^N\} \rightarrow k \text{ cluster}$$

找 k 个 center

\rightarrow 求其它数据与 k center 的哪个最近
update center
归类

Hierarchical Agglomerative Clustering (HAC)

建树: 两两计算相似度, 取最相似的时



Distributed Representation

用向量描述 attribute
分量

Dimension Reduction

Feature selection

Principle component analysis (PCA)

SVD autoencoder deep autoencoder

PCA 弱点 $\left\{ \begin{array}{l} \text{unsupervised} \\ \text{Linear} \end{array} \right.$

non-negative matrix factorization

LSA, PLSA, LDA

Lecture 15 Neighbor Embedding

LLE -Locally Linear Embedding

Laplacian Eigenmaps

t-SNE → 常用于 visualization

Lecture 16 Auto-encoder

→ NN encoder → code

→ NN decoder →



Vector Space Model

pre-training DNN fine tune

De-noising auto-encoder

auto-encoder

→ → code →

Contractive auto-encoder

auto-encoder for CNN

depooling de convolution (实际上就是 convolution)

可以利用decoder再从 code 生成 image

Lecture 17 Deep Generative Model 1

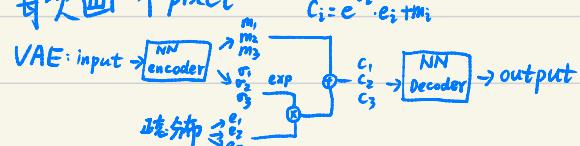
PixelRNN

Variational Autoencoder (VAE)

Generative Adversarial Network (GAN)

Pixel Recurrent Neural Network

每次画一个 pixel



minimize: $\sum (l_i \sigma_i - (m_i)^2 - e^{\sigma_i})$ for reconstruction error

Lecture 18 Deep generative model 2

VAE: 加上 noise m_i : 原来的 code

c_i : 加上 noise 后的 code

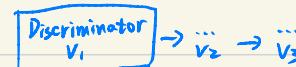
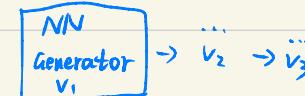
Gaussian Mixture Model

KL divergence

Conditional VAE

VAE 的问题: VAE 只产生与 training data 相近的图像
而不是以假乱真的图像

GAN since sliced bread 有史以来



generator 不看 real image, 由 discriminator 看 real image 来评价 generator 的好坏

image → discriminator →
 1. real
 0. fake

vectors from a distribution → generator → image

用 real image 和 fake image 为 label 训练 discriminator

谓 generator 让 discriminator 识别是 real image

将 generator 和 discriminator 当作一个 NN 来训练

LSTM: 4 input | output

固定 discriminator 的参数
difficult to optimise

Lecture 19 Transfer Learning

存在一些不相干的 data

source data target data

用 target data fine-tune model

用 source data 得到的 model 作为初值
(容易 over fitting)

Conservative training

新旧 model 差距不要太大

Layer Transfer

Image: 通常 copy 前面几层，只 train 后面几层

Multitask Learning

Progressive Neural Networks

Lecture 20 SVM

Hinge loss
+
kernel Method } Support Vector (SVM) Machine

convex ↴
Kernel Trick

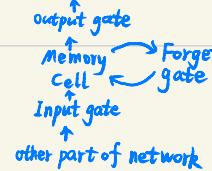
Lecture 21-1 RNN 1

Slot Filling

Elman Network & Jordan Network

Bidirectional RNN

Long Short-term Memory (LSTM)



Lecture 21-2 RNN 2

BPTT

LSTM 可以解决 gradient vanishing

Gate Recurrent Unit (GRU)

Attention-based Model

Neural Turing Machine

HMM

deep and structured

Lecture 22 Ensemble

A mixture of all approaches

Bagging: data → set1 → f1 → test → y₁ → average/voting → y
→ set2 → f2 → test → y₂
→ set3 → f3 → test → y₃

Model 很复杂，可能 over fitting 时用 (high variance)

Decision Tree → bagging → random forest

Boosting: improve weak classifiers

f₁ → f₂ → f₃ ...

AdaBoost: training f_t(x) on the new training set that fails f_{t-1}(x)

Gradient Boosting: g_t ← g_{t-1} - η $\frac{\partial L(g)}{\partial g}$ |_{g=g_{t-1}} g(x) 是函数

Stacking

