

数学的终极目标，是消除所有智力思考的需要。

— Alfred North Whitehead

## Genifer 4.1 理论笔记

YKY (甄景贤)

June 7, 2015

### 1 New insights

最近在香港认识了两个朋友，Dr 陈启良和 Dr 譚志斌 (?)，和他们谈过我的 AI 理论之后获益良多：

#### 1.1 命题空间的 dimension

在知识空间中，如果  $v_1, v_2, v_3$  是三个互不相关的命题，那么它们的线性相关  $a_1v_1 + a_2v_2 = a_3v_3$  是不容许的，否则会出现谬误。结论是：要么知识空间的 dimension 等於命题的总数（那会很大，可能无限大），要么放弃  $a_i$  是真假值的做法。

人的词汇 (vocabulary) 的个数大约在 3,000 至 10,000's 之间。暂时假设我们可以用 dimensionality reduction 把原始概念的个数缩小到 3000，那么概念空间  $C$  的维数就是 3000。

如果只考虑像  $a R b$  那样的关系，则命题空间是  $C \times C \times C$ ，维数是  $3000^3 = 9G$ 。现时的电脑记忆（例如 GPU）似乎可以处理。

## 1.2 「连续性」

第二个问题是「连续」指的是什么。例如命题  $P_1 = \text{「小明爱小娟」}$ ,  $P_2 = \text{「小强爱踢足球」}$ , 那么由  $P_1$  变化到  $P_2$  必然会有 discrete jump, 那是无可置疑的。除非我们重新定义命题是像  $k_1 P_1 + k_2 P_2$  那样的东西, 才可以有「连续变化」; 这一点需要再加以精确化。

我发觉「连续空间」不是一个严谨的术语, 因为在 discrete 空间中也可以定义 "continuous" 这概念, 例如在 computer science 的 functional programming 里有 domain theory、denotational semantics 等, 常常谈到离散空间中的连续函数。而且「可微分」的概念也有一些 generalizations, 例如 Fréchet derivatives (在泛函空间中)。

## 1.3 近似

那些算子或许可以合写成一个:

$$T_1 v + T_2 v + \dots + T_n v = \mathcal{T} v$$

然后合写后用 function approximation 来近似。

但这种近似必须有代价, 否则违反了两项原则:

1. Turing 已经很有远见地意识到, 逻辑推导的普适算法是不存在的, 因为它违反了停机原理 (halting problem)。后者是用 Cantor 的集合论中的 diagonal argument 证明的。如果我们的优化算法必然会给出答案, 而答案又可以转换回逻辑, 那是不可能的。所以在近似过程中必然会有某些误差。
2. 其二, 如果  $P \neq NP$ , 我们的算法也不可能总是在 polynomial time 之内回覆正确答案。

还有一个 "second-order" 的想法。在 first-order 我们的算子是这个形式:

$$T : V \rightarrow V$$

但 second-order 的做法是将所有物体都看成是算子, 算子可以作用在算子之上, 这似乎是一种 duality。暂时还未知道细节。

如果我们简单地将  $\mathcal{T}$  近似, 那些逻辑推导就会有错误。或者应用某个 duality 之后, 近似的做法会有较好结果?

## 2 Neural Tensor Network

Andrew Ng 是香港人, Coursera 的 machine learning 教授, 他在 Stanford, 最近加入了百度。他和合作者提出了 NTN 模型:

$$\top(a R b) = u^T \sigma(a^T U b + W \begin{pmatrix} a \\ b \end{pmatrix} + c)$$

$\top(a R b)$  表示  $a R b$  这个关系的强度。

$U$  是一个 tensor, 取两个向量  $a, b$  给出第三个向量  $a^T U b$  (它服从张量的双线性性质)。

$W$  是一个矩阵,  $W \begin{pmatrix} a \\ b \end{pmatrix}$  是一层传统的神经网络,  $\sigma$  是一个非线性函数。

$u^T$  是一粒神经元, 作用是把各个分量加起来, 得出一个数。

## 3 Paul Smolensky's tensor representation

Paul Smolensky 的书是《The harmonic mind》(2006), 第一卷是 AI 理论, 第二卷是语言学理论。

他提出用张量来表示关系:

$$\begin{aligned} father(john, pete) &\Leftrightarrow father \otimes john \otimes pete \\ Var_1 : val_1, Var_2 : val_2 &\Leftrightarrow Var_1 \otimes val_1 + Var_2 \otimes val_2 \end{aligned}$$

第二句的意思是, variable 1 的值是 value 1, variable 2 的值是 value 2。

我在博客<sup>1</sup>上解释过, tensor 是所有 bi-linear forms 的 universal form。如果有向量空间  $U$  和  $V$ ,

$$\begin{aligned} T : U \otimes V &\rightarrow W \\ t : u \otimes v &\mapsto w \end{aligned}$$

那么在  $U, V$  中线性无关的  $\{u_1, \dots, u_m\}$  和  $\{v_1, \dots, v_n\}$ , 它们的乘积  $\{u_i \otimes v_j\}$  在  $W$  中仍然是线性无关的。换句话说:

$$\text{线性无关集} \otimes \text{线性无关集} \mapsto \text{线性无关集}$$

而这正是我们需要的, 因为根据「scalar = 逻辑命题真假值」的诠释, 这正是命题之间不「相撞」的条件。

Antony Browne & Ron Sun 的较早的论文《Connectionist inference models, 2001》有讲述更多 neural-symbolic integration 的做法。

---

<sup>1</sup>blog post

## 4 Metric embedding

最近有一本新书: [Mikhail Ostrovskii 2013] *Metric embeddings: bi-Lipschitz and coarse embeddings into Banach spaces*, 它似乎正讲述了我们如何将 logic 嵌入到「连续空间」的问题。

What is the bi-Lipschitz condition? A map  $f : X \rightarrow Y$  is called a  $C$ -bi-Lipschitz embedding if there exists  $r > 0$  such that

$$\forall u, v \in X, \quad r d_X(u, v) \leq d_Y(f(u), f(v)) \leq rC d_X(u, v)$$

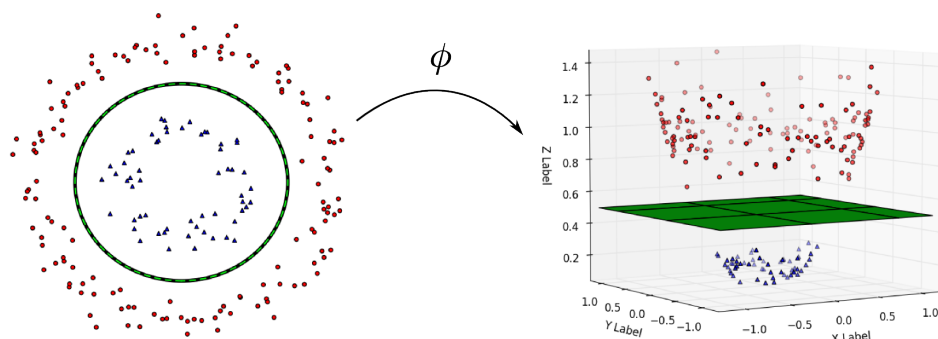
for some  $C < \infty$ . The smallest constant  $C$  for which there exists  $r > 0$  such that the above is satisfied is called the *distortion* of the map  $f$ . 这似乎定义了「远的保持远、近的保持近」的意思。

## 5 SVMs and the kernel trick

印象中 **SVM** (support vector machine) 似乎可以将非凸的问题转化成凸优化的问题; 可不可以将这个想法应用到我们的问题上呢?

SVM 的算法, 首先用 kernel trick 将 data points 转换到高维空间, 然后在高维空间中进行线性的 hyperplane 分割。

所谓 **kernel trick** 是指: 给定一个 inner product, 它暗含了一个到高维空间的转换  $\phi$  (但这个转换不需要真的计出来)。



而, 在高维空间中用线性分割, 那 error term 是一些正交距离, 所以是 **quadratic form**, 所以这问题是凸优化。

我们的问题是性质不同的, 但两者似乎有些相似...