

词向量

2018 10 16 by xzn

数据降维方法

- 当数据维度过大的情况下，模型的复杂度越高，训练模型所需的时间就会比较长，计算量就会越大
- 词汇表10w，one-hot或者tfidf都是10w维，过一个简单的线性变换，都是 $[10w, 10w]$ 参数级别。
- 可能需要用到降维的场景：
 - 数据维度之间存在一定的相关性
 - 从数据维度的含义判断可能为冗余属性
 - 训练数据集数量较多，想要缩短模型训练时间

数据降维方法

- 盲目减少属性可能会导致重要信息丢失，导致训练模型很不准确
- 目标：减少需要分析的数量，同时降低原数据信息损失量
 - 也就是利用减少后的属性可以比较大程度地代表原数据
- 常用方法：
 - PCA 主成分分析（非常简单介绍，可用）
 - LDA 线性判别分析（非常简单介绍，纯科普）
 - Word2vec 词向量（重点介绍）

PCA (Principal Component Analysis)

- 主成分分析
- 顾名思义：找到数据中的主成分，尝试用主成分信息代替原有数据，达到数据降维的目的
- 需要使用的知识：
 - 矩阵数值中心化
 - 矩阵的特征值与特征向量
 - 向量投影

PCA (Principal Component Analysis)

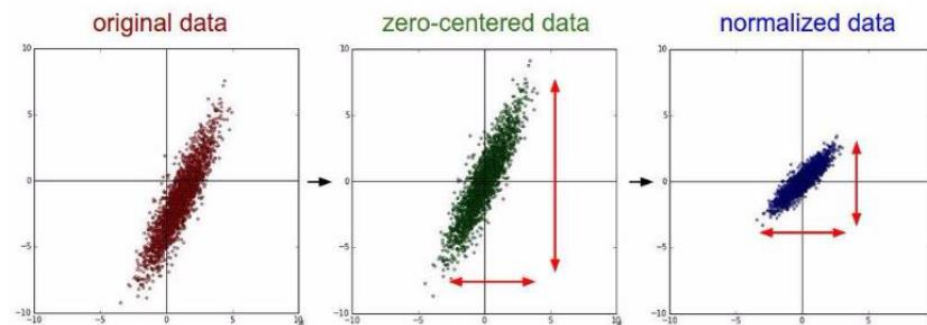
- 矩阵数值中心化
 - 每一列的各数值减去该列的均值

• 如：

3	2	3
2	3	4
4	4	2



0	-1	0
-1	0	1
1	1	-1



PCA (Principal Component Analysis)

- 计算数据矩阵的协方差矩阵

- 假设 D 维数据，协方差矩阵是 $D \times D$ 维，假设 $D = 3$ ，协方差矩阵如下：

- $$A = \begin{bmatrix} \text{cov}(d_1, d_1) & \text{cov}(d_1, d_2) & \text{cov}(d_1, d_3) \\ \text{cov}(d_2, d_1) & \text{cov}(d_2, d_2) & \text{cov}(d_2, d_3) \\ \text{cov}(d_3, d_1) & \text{cov}(d_3, d_2) & \text{cov}(d_3, d_3) \end{bmatrix}$$

- 当数据矩阵 A 是对称矩阵的时候，其奇异值等于其特征值，且存在正交矩阵 Q ，使得：
 - $Q^{-1} = Q^T$, $Q^T A Q = E$, E 为由特征值组成的对角矩阵
 - 对 A 进行奇异值分解，即可得到 A 的特征值以及特征向量

PCA (Principal Component Analysis)

- 计算数据矩阵的协方差矩阵
 - $Q^T A Q = E$, E 为由特征值组成的对角矩阵, Q 为特征向量矩阵, E 的每一列中的非 0 值为特征值 λ , 在 Q 的对应列即该特征值对应的特征向量 ev
 - $E = \begin{pmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_3 \end{pmatrix}$, $Q = [ev_1, ev_2, ev_3]$
- 此时判断 E 中的数值, 哪几个的加和占据了所有特征值之和的 95% 以上
- 这个阈值可以自由决定

PCA (Principal Component Analysis)

- 计算数据矩阵的协方差矩阵
 - 假设 $\lambda_1 + \lambda_2$ 占去了特征值总和的 95%，那么这两个特征足以代表原有的数据信息，所以我们将原数据矩阵 C (M 个输入向量) 映射到这两个特征向量上去，得到一个只有二维的数据矩阵
 - $C' = C * Q_{\lambda_1 + \lambda_2}, C' = C * [ev_1, ev_2]$
 - C 是 $M \times D$, $[ev_1, ev_2]$ 是 $D \times 2$, 得到的 C' 是 $M \times 2$
 - 以上就完成了 PCA 的流程
 - 对矩阵进行奇异值分解，求解特征值和特征向量的方法自行查阅

PCA (Principal Component Analysis)

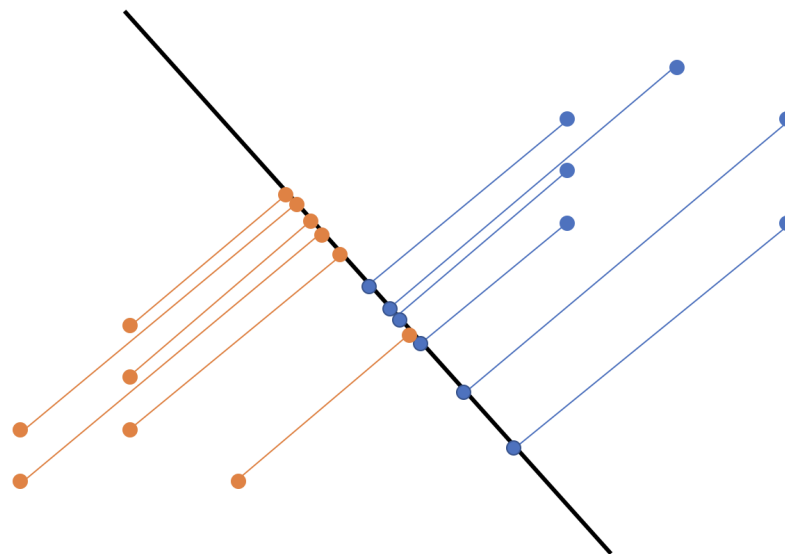
- 怎么用？
- 假设我们使用one-hot来编码，词汇表 10^w ，假设我们有 1^w 个文本作为输入，那么我们的语料库是 $[1^w, 10^w]$ ，这时候我们使用pca降维假设得到 $[1^w, 100]$ ，那么每个文本从原本的 10^w 降维到100。
- 这样后续的参数也会大幅度减少，并且将原本 10^w 的稀疏向量变成100维的高密度向量，对结果也有一定的提升。
- 允许调库sklearn.decomposition.PCA

LDA (Linear Discriminant Analysis)

- 线性判别分析
- 刚刚讲的 PCA 方法，是不涉及数据标签的无监督降维方法
- 如果某维度跟属性的划分没什么影响，那这个维度就可以删去
- 比如 身份ID 和 是否买电脑
- 但是不能够用人工判断的方式来删除属性
- 人的参与会直接影响算法得出的效果的有效性
- 需要一定的方法和指标来作为衡量的标准

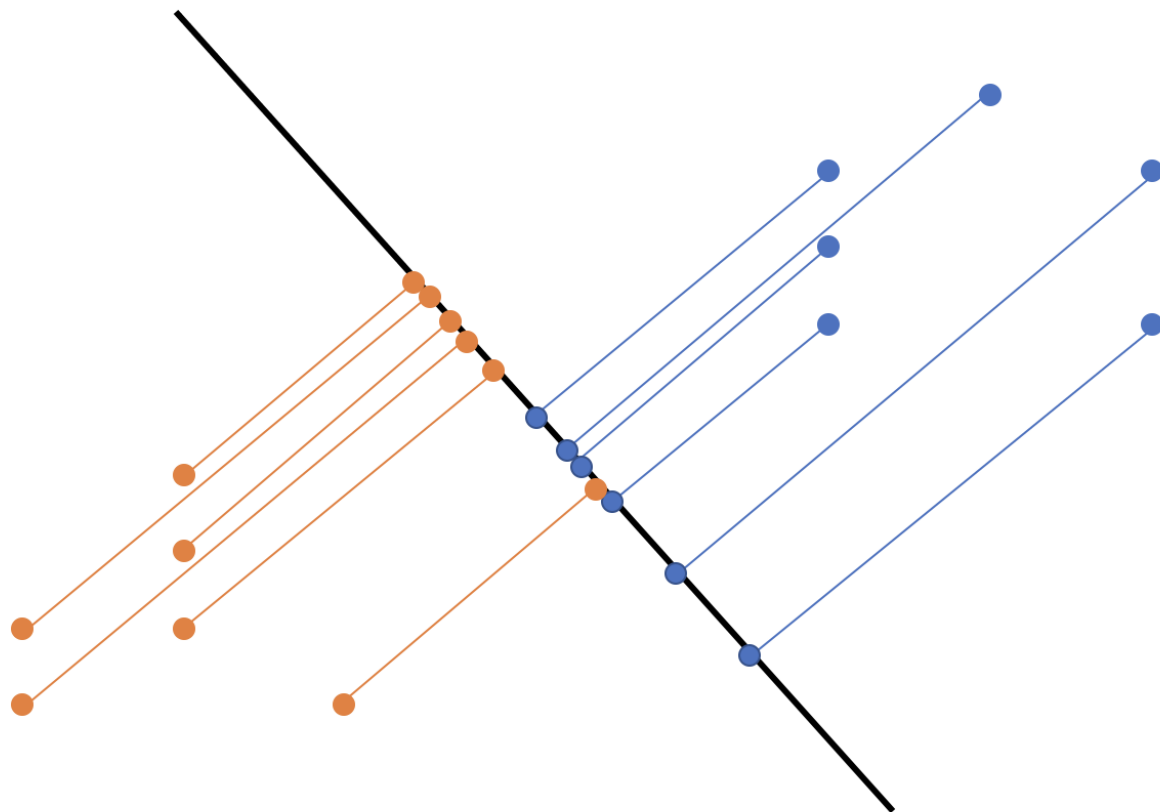
LDA (Linear Discriminant Analysis)

- 考虑二维特征，现在如果我们想要把数据的维度降到一维，但是又不想失去数据的类别信息，也就是我们想在一维的层面同样能够反映数据所属的类别
- 从二维到一维，其实就是找到一条直线，把所有的点投影到这条线上，就完成了降维的过程



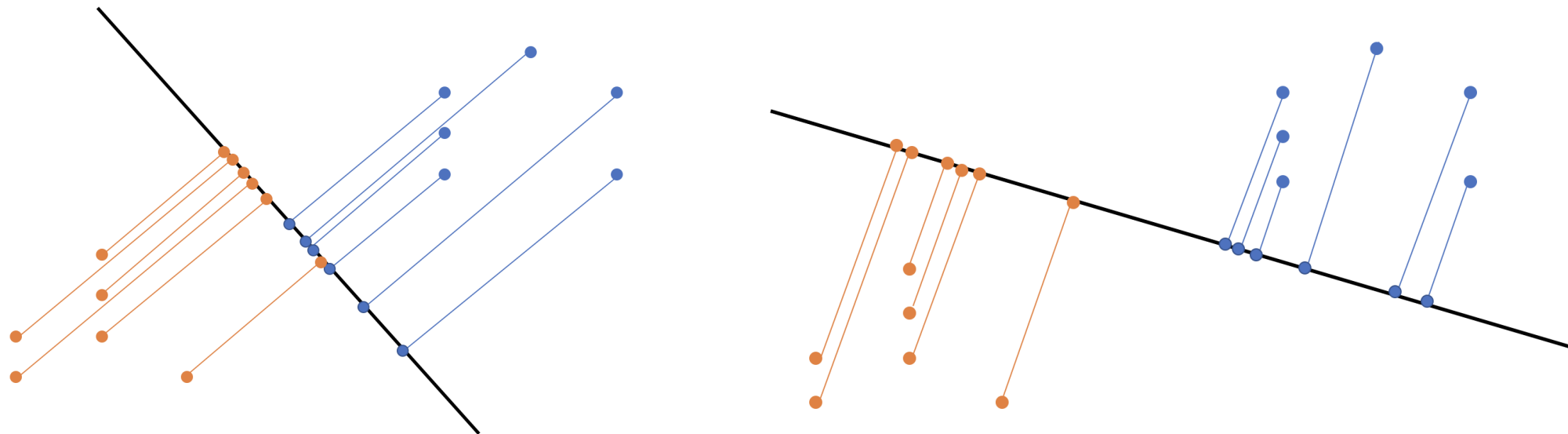
LDA (Linear Discriminant Analysis)

- 但是如果是对随意的直线进行投影，降维之后的类别可能会很混乱，可能本来可以线性划分的数据集反倒没办法线性划分了，比如刚刚的那幅图



LDA (Linear Discriminant Analysis)

- 降维后想要类别之间尽量分开

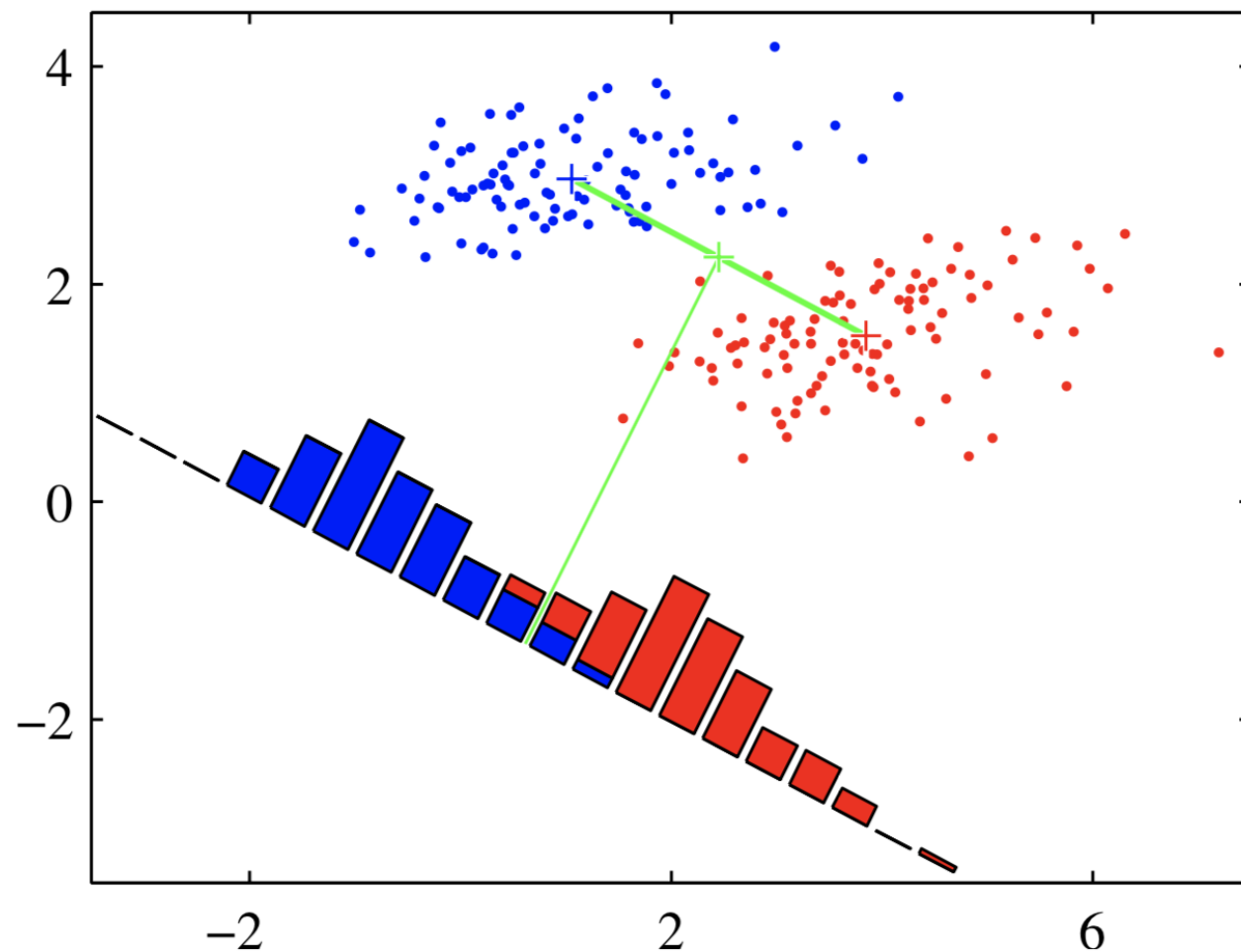


LDA (Linear Discriminant Analysis)

- 降维后想要类别之间尽量分开
- 如何实现这种目的呢？对两个类别 c_1, c_2 先求出中心点
- $m_1 = \frac{1}{N_1} \sum_{n \in c_1} x_n$ $m_2 = \frac{1}{N_2} \sum_{n \in c_2} x_n$
- 最简单的做法就是，让映射之后的中心点 m'_1, m'_2 离得越远越好
- $maximize (m'_1 - m'_2)^2$
- 那就是说选择用于投影的直线应该与两个中心点连成的直线平行

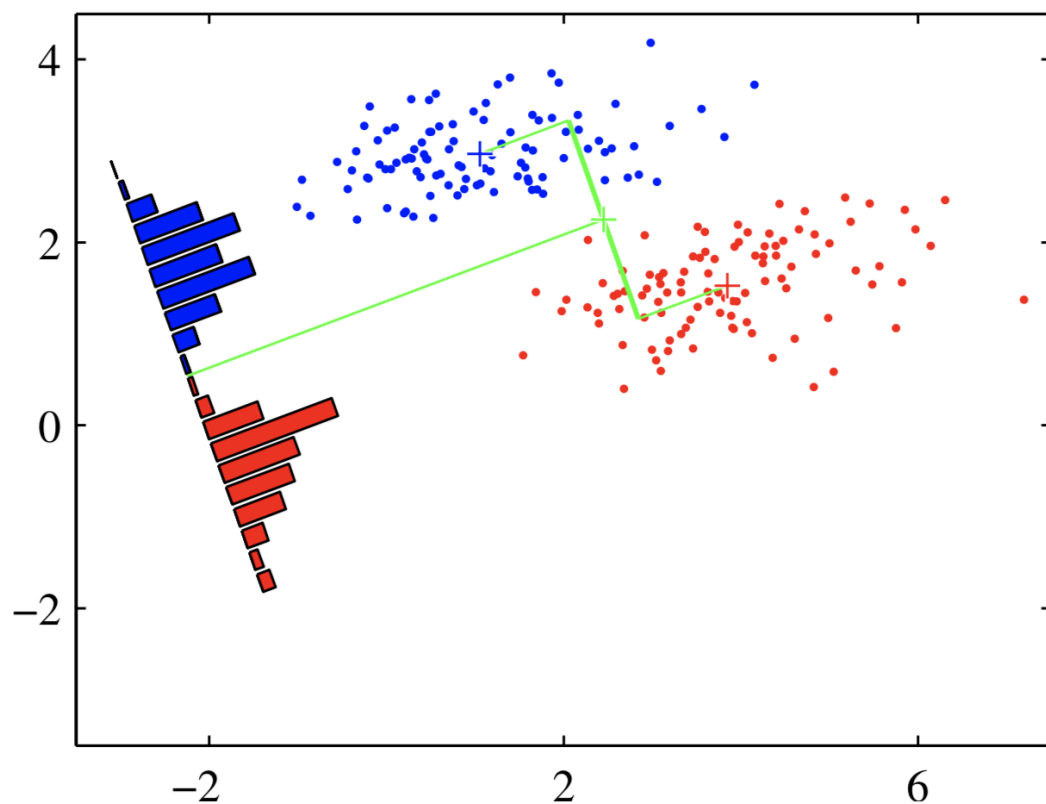
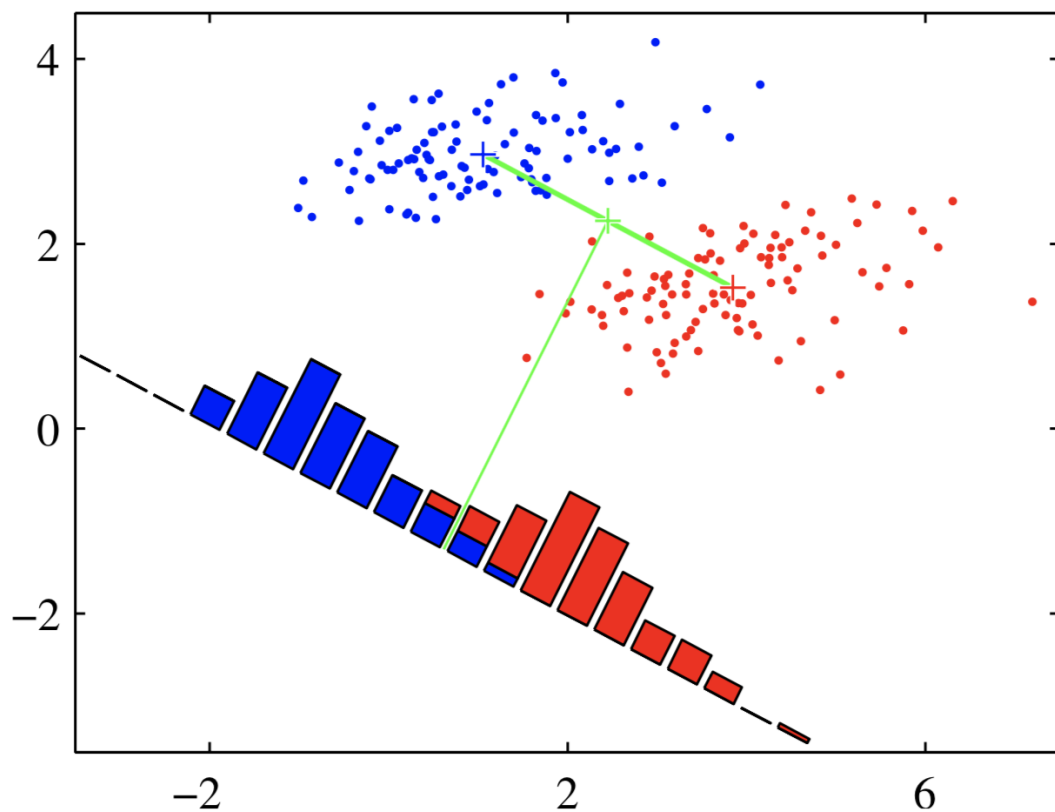
LDA (Linear Discriminant Analysis)

- 用于投影的直线应该与两个中心点连成的直线平行



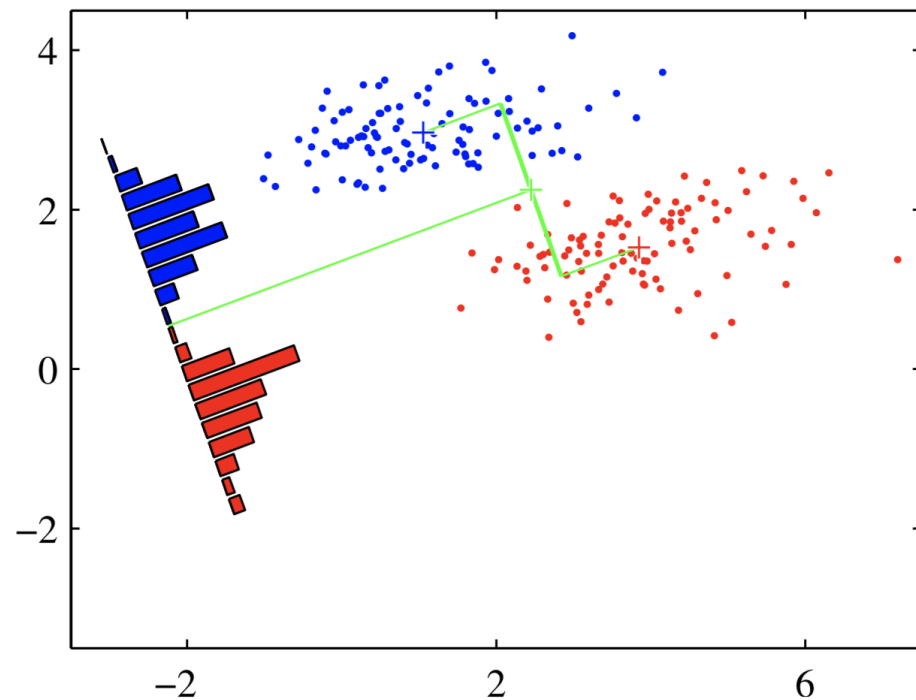
LDA (Linear Discriminant Analysis)

- 问题：然而这样也不代表着投影后有最好的划分效果，比如右图应该比较好的划分效果，让中心点映射之后最远的想法不完全正确



LDA (Linear Discriminant Analysis)

- 要保证映射之后中心点尽可能远, $maximize (x_1 - x_2)^2$
- 同时要保证映射之后每个类别的点尽可能地集中, 也就是说每个类别的方差要尽可能小
- $s_1^2 = \sum_{n \in c_1} (y_n - m'_1)^2$
- $s_2^2 = \sum_{n \in c_2} (y_n - m'_2)^2$
- 要保证 $s_1^2 + s_2^2$ 尽可能的小



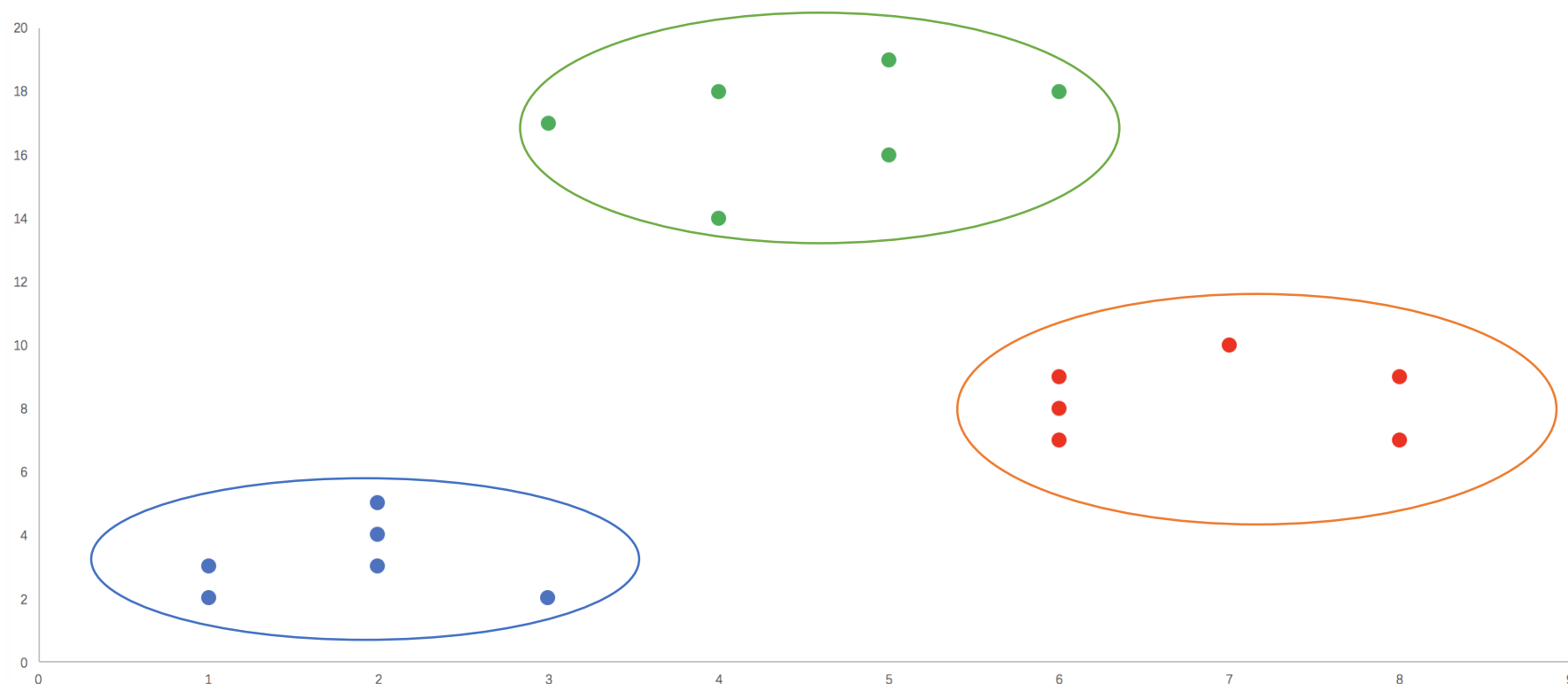
LDA (Linear Discriminant Analysis)

- 目标函数:

$$\textit{maximize} \quad \frac{(m'_1 - m'_2)^2}{s_1^2 + s_2^2}$$

LDA (Linear Discriminant Analysis)

- 多个类别的时候，我们一样先从 2 维开始分析，假设有三个类别：



- 还是让两个类别的中心点越远越好吗？目的是想让各个类别分开

Word2Vec

- Word Embedding：指的是一种单词映射的做法，将单词从原始的数据空间映射到新的多维空间上去
- 可以用在多元分类的数据集处理上
- Skip-Gram 模型是给定 input 单词预测上下文
- CBOW 模型是给定上下文预测 input 单词

government debt problems turning into banking crises as has happened in
saying that Europe needs unified banking regulation to replace the hodgepodge

↩ These words will represent *banking* ↗

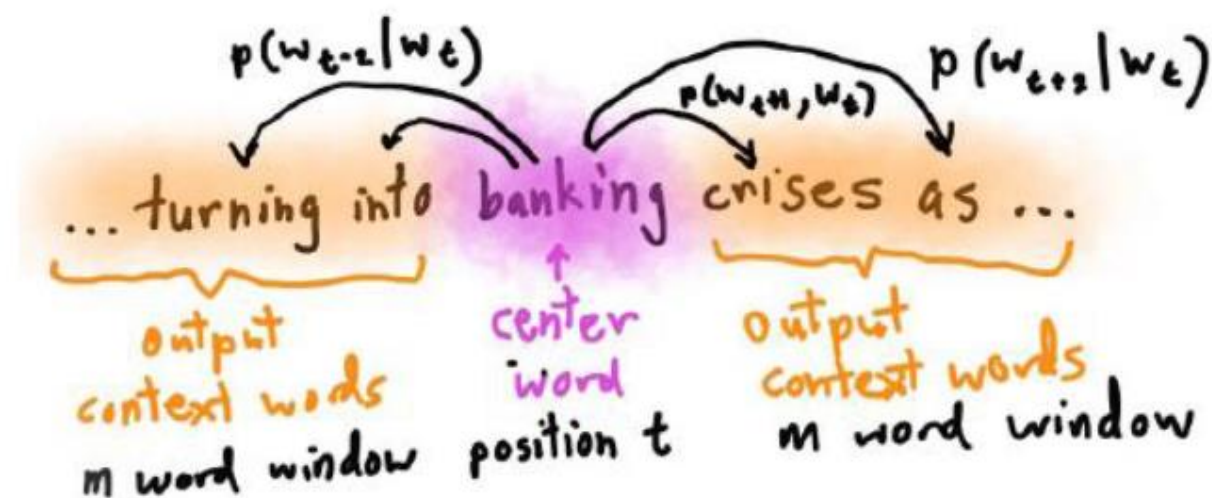
我们可以通过两种手段来训练我们的向量，一种是根据词的上下文单词来预测这个词（CBOW 模型），一种是根据词来预测该词上下文的单词（SG 模型）。我们根据这两种模型会得到一些 dense vector（与 one-hot 比较而言是 dense 的），长相可能如图 2.4 所示。

$$\textit{linguistics} = \begin{pmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \end{pmatrix}$$

表达式

$$p(\text{context}|w_t) \text{ or } p(w_t|\text{context})$$

前者是根据中心词计算上下文词的概率，后者是根据上下文词计算中心词的概率。



Word2Vec – Skip-Gram

- 举例文本：I live in China and I love China. ($V = 6$)
- 选择一个**输入词**，假设为 live
- 定义一个 skip-window 的参数，假设为 2，这个参数代表着从这个输入词的左侧或右侧选择的单词数量
- 那么选到的词就是 I live in China

Word2Vec – Skip-Gram

- I live in China
- 定义一个 num-skip 的参数，也假设为 2，这个参数代表着从根据 skip-window 选出的单词中，选择多少个不同的词作为**输出词**
- 假设选到的词是 I 和 in
- 那么我们就得到了两组训练数据，(live, I) (live, in)

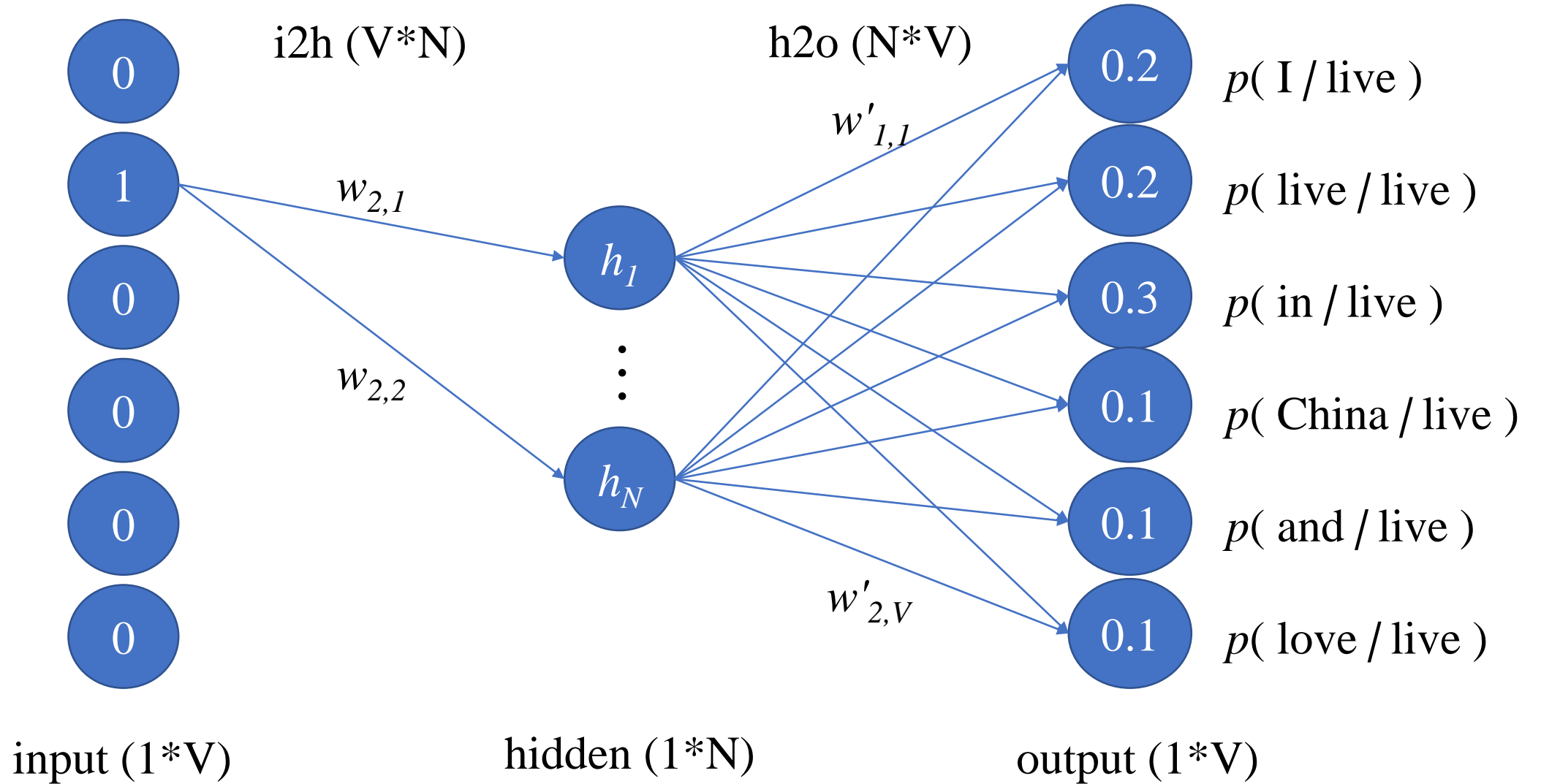
Word2Vec – Skip-Gram

- 两组训练数据, (live, l) (live, in)
- 我们想要通过模型得出的是, 由于输出词是在以输入词为中心的窗口中选择出来的, 那么这些词应该跟输入词比较相近, 也就是说**输入词有更大的几率得到这些输出词**
- 那么不是这些输出词的其他词, **输入词得到他们的概率就应该比较小**

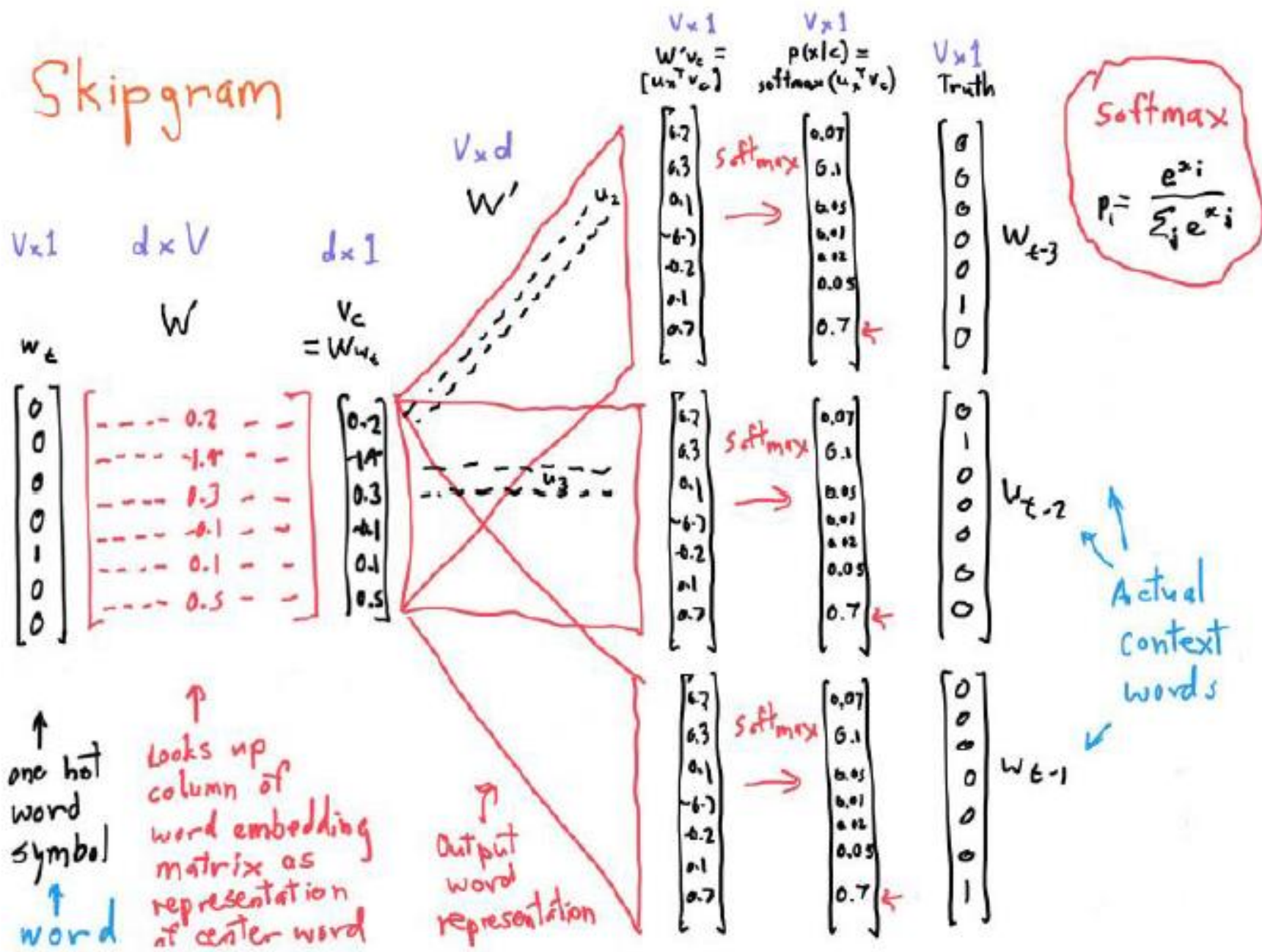
Word2Vec – Skip-Gram

- $V = 6$, 那么 live 用 onehot 向量可以表示成 $(0,1,0,0,0,0)$
- 输入 $(0,1,0,0,0,0)$ 假设输出是一个概率分布 $(0.3, 0.2, 0.1, 0.1, 0.2, 0.2)$
- 0.3最大, 对应词汇表第一个词

Word2Vec – Skip-Gram



Skipgram



Word2Vec – Skip-Gram

- softmax 函数

$$p(y_n = i | x_n, W) = \frac{e^{W_i^T x_n}}{\sum_{j=1}^M e^{(W_j^T x_n)}}$$

- 假设一共有 M 个类别，那么对应每个类别都会有一个权重向量 w ，上述公式的含义就是，利用 M 个类别的 w ，可以算出 x_n 分别属于 M 个类别的概率，做一个让概率和为 1 的归一化即可得到 x_n 属于每一个类别的概率

Word2Vec – CBOW

- CBOW 和 Skip-Gram 不一样的地方是，对于一个词来说，组成的（输入词，输出词）词对只有一个了，并且**选择的这个词是输出词**
- 我们要求解的其实就是，**依据这个词附近的这些词，有多大的概率可以得到这个词**

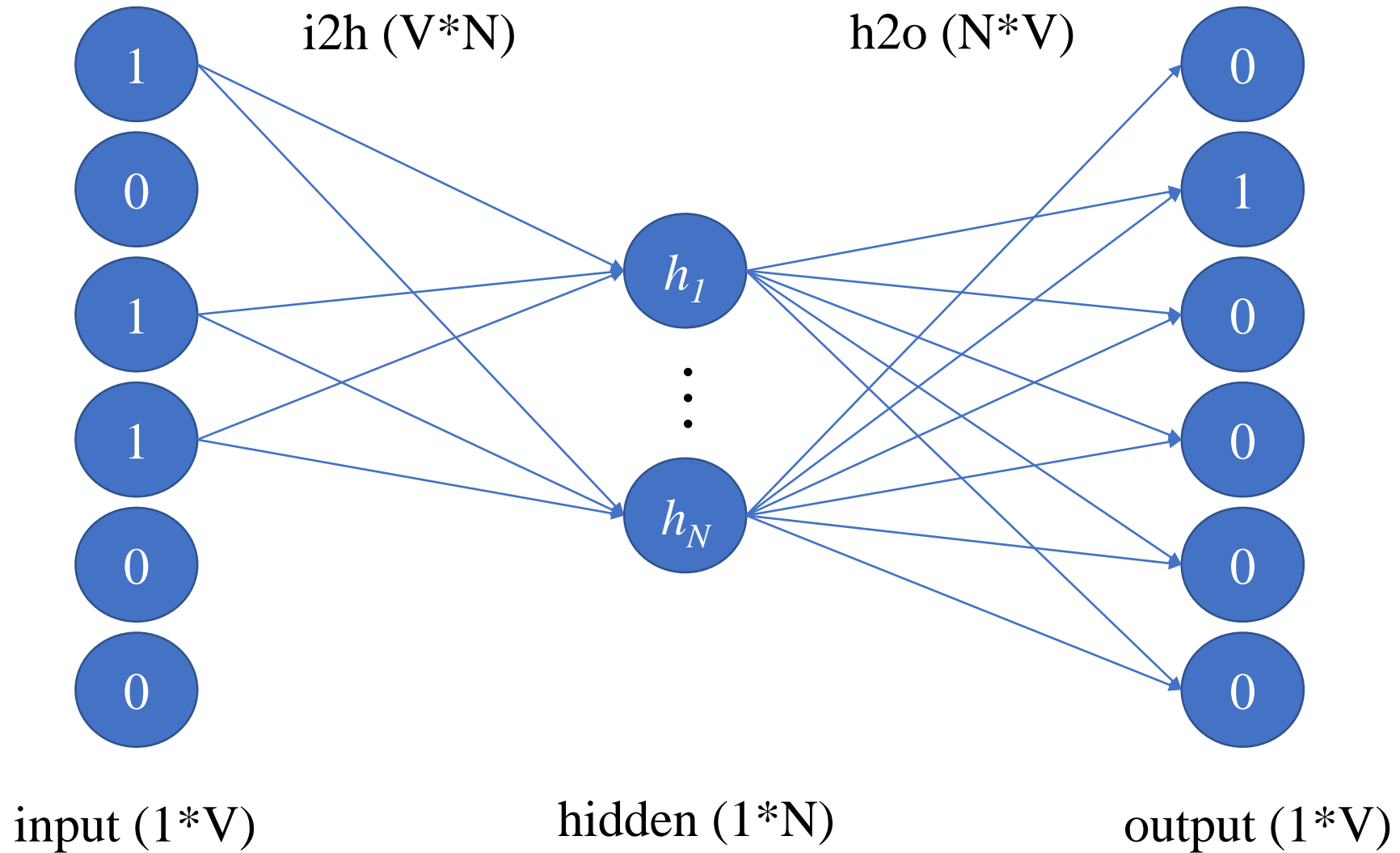
Word2Vec – CBOW

- 举例文本：I live in China and I love China. ($V = 6$)
- 选择一个**输出词**，假设为 live
- 定义一个 skip-window 的参数，假设为 2，这个参数代表着从这个输入词的左侧或右侧选择的单词数量
- 那么选到的词就是 I live in China

Word2Vec – CBOW

- I live in China
- 那么其输入的词就是 I in China
- 同样我们转为 onehot 向量，就是 (1,0,1,1,0,0)
- 那么同样的，我们传入到神经网络里面去训练
- 标准的输出是 live，也就是 (0,1,0,0,0,0)

Word2Vec – CBOW



Word2Vec – CBOW

- 把输入词的 onehot 向量当成输入向量，输出向量的第 i 维是这个输入词可以得到词汇表第 i 个单词的概率
- 同样利用的是 softmax 函数作为输出层的激活，这部分跟 Skip-Gram是一样的

Word2Vec

- 目前我们使用的主要是GLOVE
- 下载<https://nlp.stanford.edu/projects/glove/>

GloVe: Global Vectors for Word Representation

Jeffrey Pennington, Richard Socher, Christopher D. Manning

Introduction

GloVe is an unsupervised learning algorithm for obtaining vector representations for words. Training is performed on aggregated global word-word co-occurrence statistics from a corpus, and the resulting representations showcase interesting linear substructures of the word vector space.

Getting started (Code download)

- Download the [code](#) (licensed under the [Apache License, Version 2.0](#))
- Unpack the files: `unzip GloVe-1.2.zip`
- Compile the source: `cd GloVe-1.2 && make`
- Run the demo script: `./demo.sh`
- Consult the included README for further usage details, or ask a [question](#)
- The code is also available [on GitHub](#)

Download pre-trained word vectors

- Pre-trained word vectors. This data is made available under the [Public Domain Dedication and License](#) v1.0 whose full text can be found at: <http://www.opendatacommons.org/licenses/pddl/1.0/>.
 - [Wikipedia 2014](#) + [Gigaword 5](#) (6B tokens, 400K vocab, uncased, 50d, 100d, 200d, & 300d vectors, 822 MB download): [glove.6B.zip](#)
 - Common Crawl (42B tokens, 1.9M vocab, uncased, 300d vectors, 1.75 GB download): [glove.42B.300d.zip](#)
 - Common Crawl (840B tokens, 2.2M vocab, cased, 300d vectors, 2.03 GB download): [glove.840B.300d.zip](#)
 - Twitter (2B tweets, 27B tokens, 1.2M vocab, uncased, 25d, 50d, 100d, & 200d vectors, 1.42 GB download): [glove.twitter.27B.zip](#)
- Ruby [script](#) for preprocessing Twitter data

Citing GloVe

Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. [GloVe: Global Vectors for Word Representation](#). [pdf] [bib]

Highlights

1. Nearest neighbors