

第一章 引言

人工智能的真正挑战在于解决那些对人来说很容易执行、但很难形式化描述的任务

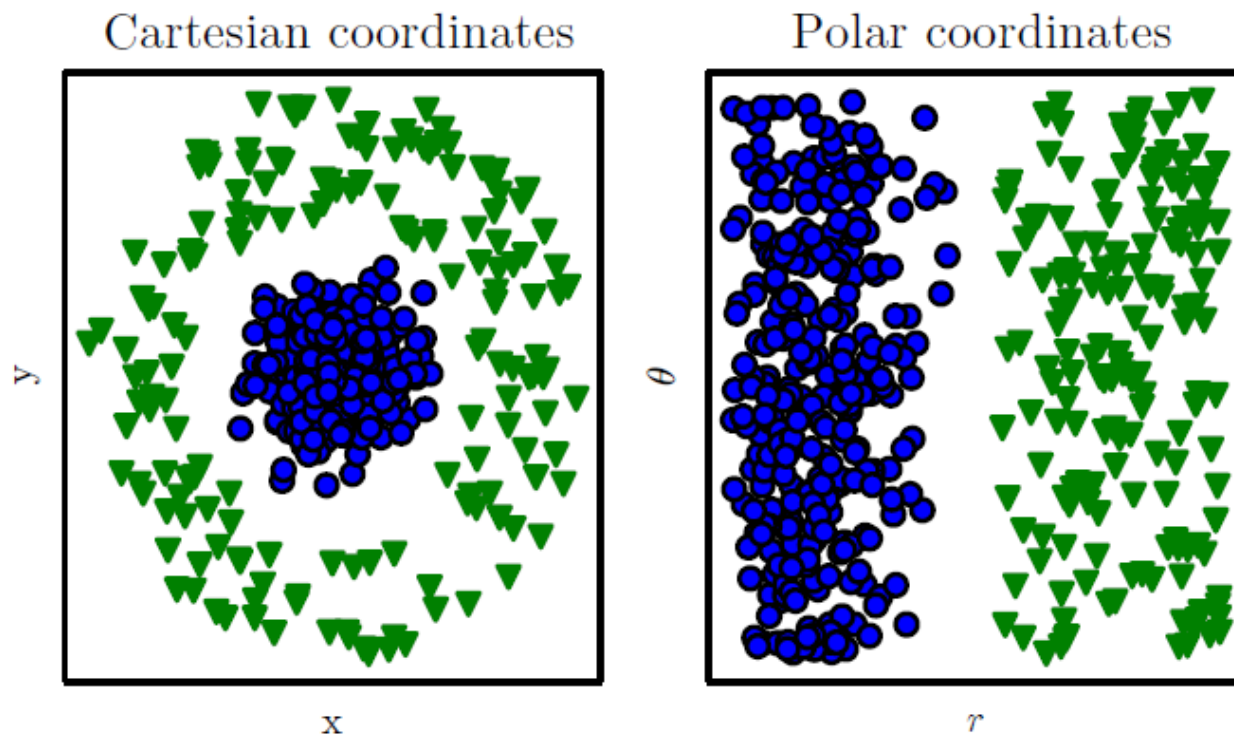
- 让计算机从经验中学习并根据层次化的概念体系来理解世界
- 层次化的概念让计算机构建较简单的概念来学习复杂概念

人工智能的一个关键挑战就是如何将这些非形式化的知识传达给计算机。

- 依靠硬编码的知识体系不能解决
- AI系统需要具备自己获取知识的能力——机器学习

第一章 引言

- 机器学习算法的性能在很大程度上依赖于给定数据表示
- 表示的选择会对机器学习算法的性能产生巨大的影响



第一章 引言

很多AI任务的一个通用解决方案：

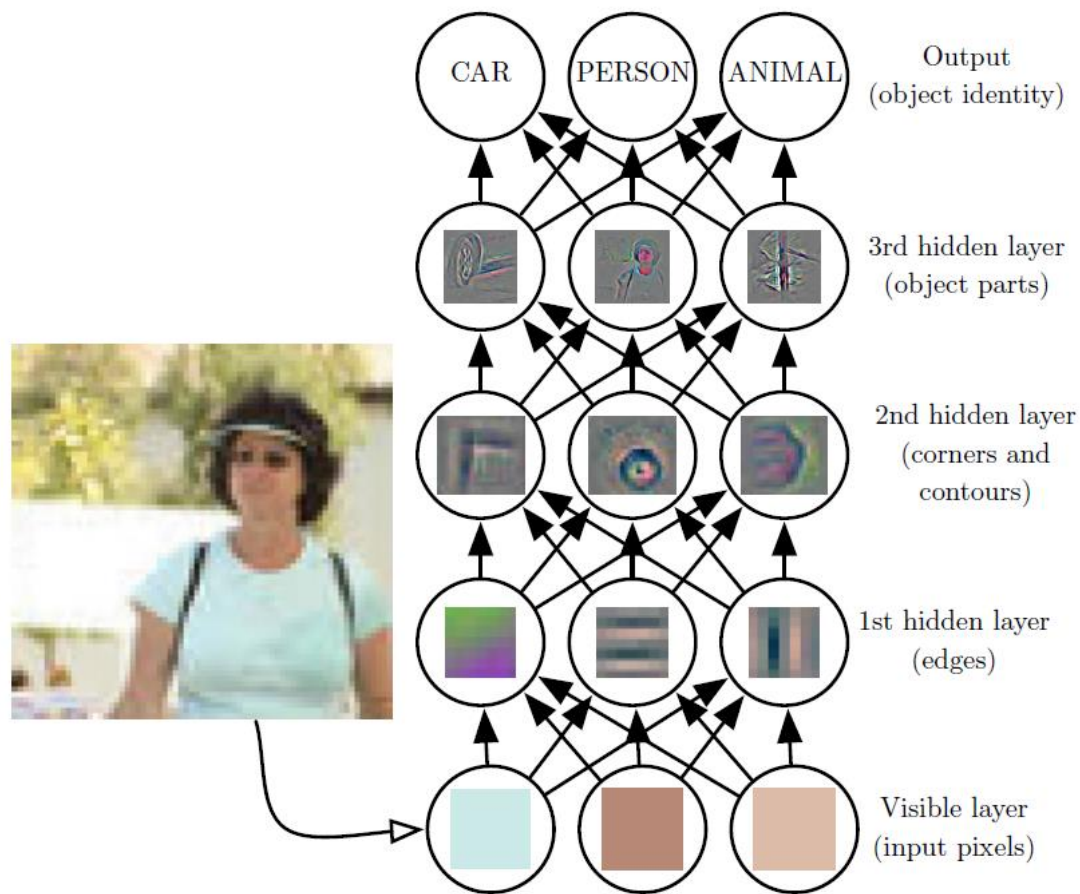
- 提取一个合适的特征集
- 将这些特征提供给简单的机器学习算法

然而，有很多任务很难知道应该提取哪些特征：

- 表示学习——使用机器学习来发掘表示本身
- 深度学习——通过其他较简单的表示来表达复杂表示

第一章 引言

深度学习模型示意图：



输出层：
对象识别

隐藏层3：
第三层检测特定对象的整个部分

隐藏层2：
第二层搜索可识别为角和扩展轮廓的边的集合

隐藏层1：
第一层识别边缘

可见层（输入层）：
输入像素

第一章 引言

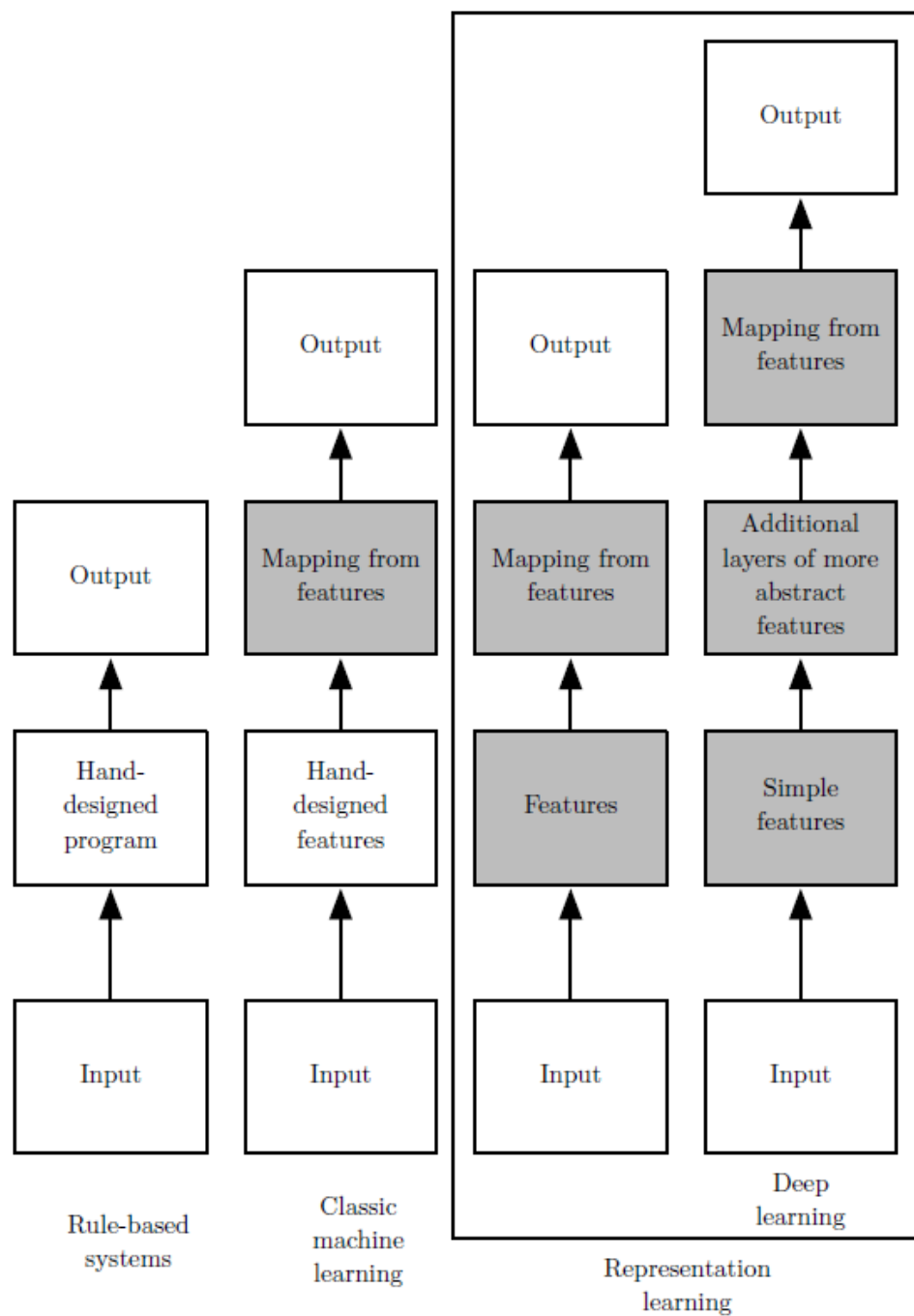
解释深度学习：

- 学习数据的正确表示
- 深度促使计算机学习一个多步骤的计算机程序

两种度量模型深度的方式：

- 基于评估架构所需执行的顺序指令的数
- 将描述概念彼此如何关联的图的深度视为模型深度

第一章 引言



第一章 引言

本书组织为三个部分：

- 第一部分介绍基本的数学工具和机器学习的概念
 - 第二到第五章
- 第二部分介绍最成熟的深度学习算法
 - 第六到第十二章
- 第三部分讨论某些具有展望性的想法
 - 第十三到第二十章

第二章 线性代数

标量、向量、矩阵和张量：

- 标量：一个标量就是一个单独的数
- 向量：一个向量是一列数
- 矩阵：矩阵是一个二维数组， \mathbf{A}
- 张量：超过两维的数组， \mathbf{A}

转置：矩阵的转置是以主对角线为轴的镜像， \mathbf{A}^T

矩阵和向量相加：

- 向量 \mathbf{b} 和矩阵 \mathbf{A} 的每一行相加
- 这种隐式地复制向量 \mathbf{b} 到很多位置的方式，被称为广播

第二章 线性代数

矩阵和向量相乘：

- 两个矩阵**A**和**B**的矩阵乘积是第三个矩阵**C**
- 乘法操作定义为

$$C_{i,j} = \sum_k A_{i,k} B_{k,j}.$$

- 两个维度相同的向量**x**和**y**的点积可看作是矩阵乘积 $\mathbf{x}^T \mathbf{y}$
- 计算**C**_{*i,j*}的步骤可以看作计算**A**的第*i*行和**B**的第*j*列的点积

矩阵乘积满足分配率和结合律，不满足交换律，向量点积满足交换律

第二章 线性代数

单位矩阵和逆矩阵：

- 保持 n 维向量不变的单位矩阵记作 \mathbf{I}_n
- 形式上， $\mathbf{I}_n \in \mathbb{R}^{n \times n}$

$$\forall \mathbf{x} \in \mathbb{R}^n, \mathbf{I}_n \mathbf{x} = \mathbf{x}.$$

- 单位矩阵的结构很简单：所有沿主对角线的元素都是1，而所有其他位置的元素都是0

矩阵 \mathbf{A} 的矩阵逆矩阵（matrix inversion）记作 \mathbf{A}^{-1} ，

$$\mathbf{A}^{-1} \mathbf{A} = \mathbf{I}_n.$$

第二章 线性代数

线性相关和生成子空间：

- 一组向量的线性组合，指的是每个向量乘以对应标量系数之后的和
$$\sum_i c_i \mathbf{v}^{(i)}.$$
- 如果一组向量中的任意一个向量都不能表示成其他向量的线性组合，那么这组向量称为线性无关（linearly independent）
- 反之，这组向量线性相关
- 一组向量的生成子空间（span）是原始向量线性组合后所能抵达的点的集合

第二章 线性代数

范数（norm）：

- L^p 范数定义如下：

$$\|\mathbf{x}\|_p = \left(\sum_i |x_i|^p \right)^{\frac{1}{p}} \quad p \in \mathbb{R}, p \geq 1$$

- 范数是满足下列性质的任意函数：

- $f(\mathbf{x}) = 0 \Rightarrow \mathbf{x} = \mathbf{0}$
- $f(\mathbf{x} + \mathbf{y}) \leq f(\mathbf{x}) + f(\mathbf{y})$ （三角不等式（triangle inequality））
- $\forall \alpha \in \mathbb{R}, f(\alpha \mathbf{x}) = |\alpha| f(\mathbf{x})$

- L^2 范数被称为欧几里得范数，常简化表示为 $\|\mathbf{x}\|$ ，可以简单地通过点积 $\mathbf{x}^\top \mathbf{x}$ 计算

第二章 线性代数

范数（norm）：

- L^1 范数常作为表示非零元素数目的替代函数
- L^∞ 范数，也被称为最大范数，表示向量中具有最大幅值的元素的绝对值：

$$\|x\|_\infty = \max_i |x_i|.$$

- **Frobenius** 范数，衡量矩阵的大小：

$$\|A\|_F = \sqrt{\sum_{i,j} A_{i,j}^2},$$

- 两个向量的点积（dot product）可以用范数来表示

$$x^\top y = \|x\|_2 \|y\|_2 \cos \theta$$

第二章 线性代数

特殊类型的矩阵和向量：

- 对角矩阵：只在主对角线上含有非零元素，其他位置都是零
- 对称矩阵：矩阵是转置和自己相等的矩阵， $\mathbf{A} = \mathbf{A}^\top$.
- 单位向量：具有单位范数的向量， $\|\mathbf{x}\|_2 = 1$.
- 正交： $\mathbf{x}^\top \mathbf{y} = 0$ ， 向量 \mathbf{x} 和向量 \mathbf{y} 互相正交
- 如果这些向量不仅互相正交，并且范数都为1，那么我们称它们标准正交
- 正交矩阵：指行向量和列向量是分别标准正交的方阵

$$\mathbf{A}^\top \mathbf{A} = \mathbf{A} \mathbf{A}^\top = \mathbf{I}. \quad \mathbf{A}^{-1} = \mathbf{A}^\top.$$

第二章 线性代数

特征分解：

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}.$$

\mathbf{v} 是特征向量，标量 λ 被称为这个特征向量对应的特征值

\mathbf{A} 的特征分解： $\mathbf{A} = \mathbf{V}\text{diag}(\boldsymbol{\lambda}) \mathbf{V}^{-1}.$

$$\mathbf{V} = [\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(n)}]$$

$$\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_n]^\top$$

每个实对称矩阵都可以分解成实特征向量和实特征值：

$$\mathbf{A} = \mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^\top.$$

其中 \mathbf{Q} 是 \mathbf{A} 的特征向量组成的正交矩阵， $\boldsymbol{\Lambda}$ 是对角矩阵

第二章 线性代数

特征分解：

- 正定：所有特征值都是正数的矩阵被称为正定
- 半正定：所有特征值都是非负数的矩阵被称为半正定
- 负定：所有特征值都是负数的矩阵被称为负定
- 半负定：所有特征值都是非正数的矩阵被称为半负定

第二章 线性代数

奇异值分解：

$$\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^{\top}.$$

\mathbf{A} 是 $m \times n$ 的矩阵， \mathbf{U} 是 $m \times m$ 的矩阵

\mathbf{D} 是 $m \times n$ 的矩阵， \mathbf{V} 是 $n \times n$ 的矩阵

这些矩阵中的每一个经定义后都拥有特殊的结构：

- \mathbf{U} 和 \mathbf{V} 都定义为正交矩阵
- 矩阵 \mathbf{D} 定义为对角矩阵，注意， \mathbf{D} 不一定是方阵

第二章 线性代数

Moore-Penrose 伪逆:

矩阵 \mathbf{A} 的伪逆定义为:

$$\mathbf{A}^+ = \lim_{\alpha \searrow 0} (\mathbf{A}^\top \mathbf{A} + \alpha \mathbf{I})^{-1} \mathbf{A}^\top.$$

计算伪逆的实际算法是使用下面的公式:

$$\mathbf{A}^+ = \mathbf{V} \mathbf{D}^+ \mathbf{U}^\top.$$

矩阵 \mathbf{U} , \mathbf{D} 和 \mathbf{V} 是矩阵 \mathbf{A} 奇异值分解后得到的矩阵

对角矩阵 \mathbf{D} 的伪逆 \mathbf{D}^+ 是其非零元素取倒数之后再转置得到的

第二章 线性代数

Moore-Penrose 伪逆:

矩阵 \mathbf{A} 不是方阵:

$$\mathbf{A}\mathbf{x} = \mathbf{y}$$

如果矩阵 \mathbf{A} 的行数小于列数, 那么上述方程可能有多个解

如果矩阵 \mathbf{A} 的行数大于列数, 那么上述方程可能没有解

如果矩阵 \mathbf{A} 的行数小于列数, $\mathbf{x} = \mathbf{A}^+\mathbf{y}$ 是方程所有可行解中欧几里得范数最小的一个

如果矩阵 \mathbf{A} 的行数大于列数, 可能没有解, 在这种情况下, 通过伪逆得到的 \mathbf{x} 使得 \mathbf{Ax} 和 \mathbf{y} 的欧几里得距离 $\|\mathbf{Ax} - \mathbf{y}\|_2$ 最小

第二章 线性代数

迹运算：

迹运算返回的是矩阵主对角线上的元素的和：

$$\text{Tr}(\mathbf{A}) = \sum_i \mathbf{A}_{i,i}.$$

迹运算提供了另一种描述矩阵Frobenius范数的方式：

$$\|\mathbf{A}\|_F = \sqrt{\text{Tr}(\mathbf{A}\mathbf{A}^\top)}.$$

多个矩阵相乘得到的方阵的迹，和将这些矩阵中的最后一个挪到最前面之后相乘的迹是相同的

$$\text{Tr}\left(\prod_{i=1}^n \mathbf{F}^{(i)}\right) = \text{Tr}\left(\mathbf{F}^{(n)} \prod_{i=1}^{n-1} \mathbf{F}^{(i)}\right).$$

标量在迹运算后仍然是它自己

实例 主成分分析 (PCA)

假设在 \mathbb{R}^n 空间中有 m 个点 $\{x^{(1)}, \dots, x^{(m)}\}$:

对这些点进行有损压缩, 使用更少的内存, 但损失一些精度, 我们希望损失的精度尽可能少。

方法:

$$x^{(i)} \in \mathbb{R}^n \rightarrow \text{编码向量 } c^{(i)} \in \mathbb{R}^l, \quad l < n$$

编码函数, 根据输入返回编码: $f(x) = c$

解码函数, 给定编码的重构输入: $x \approx g(f(x))$

简化解码器: 令 $g(c) = Dc$, $D \in \mathbb{R}^{n \times l}$

限制1: D 的列向量都有单位范数 (为了获得唯一解)

限制2: D 的列向量彼此正交 (简化解码器的最优编码问题)

实例 主成分分析 (PCA)

解最优编码问题:

思路: 最小化原始输入向量 \mathbf{x} 和重构向量 $g(\mathbf{c})$ 之间的距离
可以使用 L^2 范数:

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} \|\mathbf{x} - g(\mathbf{c})\|_2.$$

因为 L^2 范数非负, 可以使用平方 L^2 范数:

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} \|\mathbf{x} - g(\mathbf{c})\|_2^2.$$

根据定义:

$$\begin{aligned} \|\mathbf{x} - g(\mathbf{c})\|_2^2 &= (\mathbf{x} - g(\mathbf{c}))^T (\mathbf{x} - g(\mathbf{c})) \\ &= \mathbf{x}^T \mathbf{x} - 2\mathbf{x}^T g(\mathbf{c}) + g(\mathbf{c})^T g(\mathbf{c}) \end{aligned}$$

实例 主成分分析 (PCA)

解最优编码问题:

$\mathbf{x}^T \mathbf{x}$ 不依赖与 \mathbf{c} , 可以直接忽略:

$$\begin{aligned}\mathbf{c}^* &= \arg \min_{\mathbf{c}} -2\mathbf{x}^T g(\mathbf{c}) + g(\mathbf{c})^T g(\mathbf{c}). \\ &= \arg \min_{\mathbf{c}} -2\mathbf{x}^T D\mathbf{c} + \mathbf{c}^T D^T D\mathbf{c} \\ &= \arg \min_{\mathbf{c}} -2\mathbf{x}^T D\mathbf{c} + \mathbf{c}^T I_l \mathbf{c} \\ &= \arg \min_{\mathbf{c}} -2\mathbf{x}^T D\mathbf{c} + \mathbf{c}^T \mathbf{c}\end{aligned}$$

微积分求解: $\nabla_{\mathbf{c}}(-2\mathbf{x}^T D\mathbf{c} + \mathbf{c}^T \mathbf{c}) = 0$

$$-2D^T \mathbf{x} + 2\mathbf{c} = 0$$

于是: $\mathbf{c} = D^T \mathbf{x}.$

$$f(\mathbf{x}) = D^T \mathbf{x}. \quad r(\mathbf{x}) = g(f(\mathbf{x})) = DD^T \mathbf{x}.$$

实例 主成分分析 (PCA)

挑选编码矩阵 \mathbf{D} :

回顾问题: 最小化原始输入向量 \mathbf{x} 和重构向量 $g(\mathbf{c})$ 之间的距离

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} \|\mathbf{x} - g(\mathbf{c})\|_2.$$

此时, 必须最小化所有维度和所有点上的误差矩阵, 所以使用 Frobenius 范数:

$$\mathbf{D}^* = \arg \min_{\mathbf{D}} \sqrt{\sum_{i,j} \left(\mathbf{x}_j^{(i)} - r(\mathbf{x}^{(i)})_j \right)^2} \text{ subject to } \mathbf{D}^\top \mathbf{D} = \mathbf{I}_l.$$

我们首先考虑 $l = 1$ 的情况, 此时, \mathbf{D} 是一个向量 \mathbf{d} :

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \sum_i \left\| \mathbf{x}^{(i)} - \mathbf{d} \mathbf{d}^\top \mathbf{x}^{(i)} \right\|_2^2 \text{ subject to } \|\mathbf{d}\|_2 = 1.$$

实例 主成分分析 (PCA)

挑选编码矩阵 \mathbf{d} :

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \sum_i \left\| \mathbf{x}^{(i)} - \mathbf{d} \mathbf{d}^\top \mathbf{x}^{(i)} \right\|_2^2 \text{ subject to } \|\mathbf{d}\|_2 = 1. \quad \mathbf{d}^\top \mathbf{x}^{(i)} \text{ 是标量}$$

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \sum_i \left\| \mathbf{x}^{(i)} - \mathbf{d}^\top \mathbf{x}^{(i)} \mathbf{d} \right\|_2^2 \text{ subject to } \|\mathbf{d}\|_2 = 1, \quad \text{写到左边}$$

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \sum_i \left\| \mathbf{x}^{(i)} - \mathbf{x}^{(i)\top} \mathbf{d} \mathbf{d}^\top \right\|_2^2 \text{ subject to } \|\mathbf{d}\|_2 = 1. \quad \text{转置}$$

将表示各点的向量按行堆叠成一个矩阵 \mathbf{X} :

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \left\| \mathbf{X} - \mathbf{X} \mathbf{d} \mathbf{d}^\top \right\|_F^2 \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1.$$

实例 主成分分析 (PCA)

挑选编码矩阵 \mathbf{D} :

将表示各点的向量按行堆叠成一个矩阵 \mathbf{X} :

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \left\| \mathbf{X} - \mathbf{X} \mathbf{d} \mathbf{d}^\top \right\|_F^2 \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1.$$

不考虑约束, 简化Frobenius 范数:

$$\begin{aligned} \arg \min_{\mathbf{d}} \left\| \mathbf{X} - \mathbf{X} \mathbf{d} \mathbf{d}^\top \right\|_F^2 &= \arg \min_{\mathbf{d}} \text{Tr} \left(\left(\mathbf{X} - \mathbf{X} \mathbf{d} \mathbf{d}^\top \right)^\top \left(\mathbf{X} - \mathbf{X} \mathbf{d} \mathbf{d}^\top \right) \right) \\ &= \arg \min_{\mathbf{d}} \text{Tr} \left(\mathbf{X}^\top \mathbf{X} - \mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top - \mathbf{d} \mathbf{d}^\top \mathbf{X}^\top \mathbf{X} + \mathbf{d} \mathbf{d}^\top \mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top \right) \\ &= \arg \min_{\mathbf{d}} \text{Tr}(\mathbf{X}^\top \mathbf{X}) - \text{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) - \text{Tr}(\mathbf{d} \mathbf{d}^\top \mathbf{X}^\top \mathbf{X}) + \text{Tr}(\mathbf{d} \mathbf{d}^\top \mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) \\ &= \arg \min_{\mathbf{d}} - \text{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) - \text{Tr}(\mathbf{d} \mathbf{d}^\top \mathbf{X}^\top \mathbf{X}) + \text{Tr}(\mathbf{d} \mathbf{d}^\top \mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) \\ &= \arg \min_{\mathbf{d}} - 2\text{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) + \text{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top \mathbf{d} \mathbf{d}^\top) \end{aligned}$$

实例 主成分分析 (PCA)

挑选编码矩阵 \mathbf{D} :

加上约束:

$$\arg \min_{\mathbf{d}} -2\text{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) + \text{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top \mathbf{d} \mathbf{d}^\top) \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1$$

$$= \arg \min_{\mathbf{d}} -2\text{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) + \text{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1$$

$$= \arg \min_{\mathbf{d}} -\text{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1$$

$$= \arg \max_{\mathbf{d}} \text{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1$$

$$= \arg \max_{\mathbf{d}} \text{Tr}(\mathbf{d}^\top \mathbf{X}^\top \mathbf{X} \mathbf{d}) \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1.$$

实例 主成分分析（PCA）

挑选编码矩阵 **D** :

$$\arg \max_d \text{Tr}(\mathbf{d}^\top \mathbf{X}^\top \mathbf{X} \mathbf{d}) \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1.$$

这个优化问题可以通过特征分解来求解。具体来讲，最优的 **d** 是 **$\mathbf{X}^\top \mathbf{X}$** 最大特征值对应的特征向量

以上推导特定于 **$l = 1$** 的情况，仅得到了第一个主成分。

当我们希望得到主成分的基时，矩阵 **D** 是由 **$\mathbf{X}^\top \mathbf{X}$** 的前 **$l$** 个最大的特征值对应的特征向量组成的。