

# 最 优 化 模 型

## III 应 用

刘红英 编

北京航空航天大学数学与系统科学学院

2017 年 3 月

## 内容简介

本书是Optimization Models(G.C. Calafiore and L. El Ghaoui, Cambridge University Press. Oct. 2014)一书第III部分应用的中文翻译版.

## 前 言

本文档由北航本科生翻译并录入(详细分工如下),由刘红英统一校对并编辑排版。编者在此向他们的辛勤劳动表示衷心的感谢。

第十三章的协调组长为董昊,开头至13.4.2节结束由董昊翻译,13.4.3节至结束由王翔翻译。刘晗负责排版、插入图表,并校对公式。

第十四章的协调组长为周仕佶,开头至14.1节结束由周仕佶翻译,14.2节至14.3.1节结束由张林翻译,14.3.2节至14.4节结束由于浩翻译,14.5节至结束由徐亮翻译。

第十五章的协调组长为钞乾,开头至15.1节结束由钞乾翻译,15.2节由钱楚楚翻译,15.3节至结束由阮钰泽翻译,李冠巡负责排版校对。

第十六章的协调组长为沙铜,开头至16.1节结束由乐然翻译,16.2节和16.3节由沙铜翻译,16.4节和16.5节由周子易翻译,16.6节至结束由李博翻译,乐然负责处理和插入图表,沙铜负责排版校对。

由于这些翻译文稿尚未经过特别仔细的雕琢和严格的校对,仅供学习交流参考用。任何意见、建议或其它反馈都可以发送至liuhongying@buaa.edu.cn,在此深表感谢。

刘红英

2017.3 于北京



# 目 录

第十三章 基于数据的学习	1
13.1 监督学习概述	1
13.2 基于多项式模型的最小二乘预测	3
13.2.1 模型	3
13.2.2 性能评价	5
13.2.3 正则化和稀疏性	6
13.3 二分类	8
13.3.1 支持向量机	9
13.3.2 正则化与稀疏性	10
13.3.3 几何解释	10
13.3.4 鲁棒性	14
13.3.5 Logistic回归	15
13.3.6 Fisher判别	16
13.4 一种通用的监督学习问题	18
13.4.1 损失函数	19
13.4.2 惩罚和约束函数	19
13.4.3 Kernel方法	20
13.5 无监督学习	22
13.5.1 主成分分析 (Principal component analysis, PCA)	23
13.5.2 稀疏PCA	23
13.5.3 非负矩阵分解	26
13.5.4 鲁棒PCA	27
13.5.5 稀疏图的高斯模型	29
13.6 习题	33
第十四章 计算金融	39

14.1	单期投资组合最优化	39
14.1.1	均值-方差最优化	40
14.1.2	投资组合约束和交易成本	42
14.1.3	夏普比率(SR)最优化	43
14.1.4	风险价值(Value-at-Risk, VaR)最优化	45
14.2	鲁棒的投资组合最优化	46
14.2.1	鲁棒的均值-方差最优化	47
14.2.2	鲁棒的风险价值(VaR)最优化	48
14.3	多期投资组合最优化	50
14.3.1	投资组合动力学	51
14.3.2	最优的开环策略	53
14.3.3	带仿射价格的闭环分配	54
14.4	稀疏指标跟踪	56
14.5	习题	58
<b>第十五章 控制问题</b>		<b>67</b>
15.1	连续及离散时间模型	67
15.1.1	连续时间LTI 系统	67
15.1.2	离散时间LTI系统	69
15.2	基于最优化的控制综合问题	71
15.2.1	状态跟踪的控制命令的综合	71
15.2.2	轨迹追踪控制命令综合	75
15.2.3	模型预测控制	78
15.3	分析和控制器设计中的最优化	80
15.3.1	连续时间的Lyapunov 稳定性分析	80
15.3.2	镇定状态反馈设计	82
15.3.3	离散时间的Lyapunov稳定性和镇定	86
15.4	练习	87

---

<b>第十六章 工程设计</b>	<b>91</b>
16.1 数字滤波器设计 . . . . .	91
16.1.1 线性相位FIR滤波器 . . . . .	92
16.1.2 低通FIR设计规范 . . . . .	93
16.1.3 通过线性规划设计FIR . . . . .	94
16.1.4 逼近参考匹配的滤波器设计 . . . . .	96
16.2 天线阵列设计 . . . . .	99
16.2.1 天线方向图赋形 . . . . .	101
16.2.2 最小二乘设计法 . . . . .	103
16.2.3 SOCP设计法 . . . . .	105
16.3 数字电路设计 . . . . .	106
16.3.1 电路拓扑 . . . . .	107
16.3.2 设计变量 . . . . .	107
16.3.3 设计目标 . . . . .	108
16.3.4 一个电路设计问题 . . . . .	109
16.4 航天器设计 . . . . .	110
16.5 供应链管理 . . . . .	115
16.5.1 针对区间不确定需求的鲁棒法 . . . . .	116
16.5.2 区间不确定性下的仿射订购策略 . . . . .	117
16.5.3 用于一般的随机不确定性的场景法 . . . . .	121
16.6 练习 . . . . .	125

## 第十三章 基于数据的学习

如果事实与理论不符，改变那个事实。 阿尔伯特 爱因斯坦

本章以最优化的视角简单介绍了机器学习中出现的一些典型问题。我们首先研究所谓的监督学习问题，这里的基本问题是用给定的响应数据拟合出某一模型。然后我们研究无监督学习问题，这里也是为数据构建模型，与监督学习问题的区别之处是没有特定的响应。

在这章中，我们用  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m] \in \mathbb{R}^{n \times m}$  表示一般的数据矩阵，其中  $\mathbf{x}_i$  是第  $i$  个数据点，也称其为**样本**(example)。我们称数据点  $\mathbf{x}$  的一个特定维度为一个**特征**(feature)。常见的特征包括：给定单词在字典中出现的频率<sup>1</sup>；决定一个特定特征是否出现的Boolean 变量，例如一个特定演员是否在数据点代表的电影中出演；或者数值量，比如血压，温度，价格等等。

### §13.1 监督学习概述

在监督学习中，我们需要为未知的函数  $\mathbf{x} \rightarrow y(\mathbf{x})$  建立一个模型，这里  $\mathbf{x} \in \mathbb{R}^n$  是输入向量， $y(\mathbf{x}) \in \mathbb{R}$  是对应的输出。假设给了一些观测或样本的集合，即给定一些输入输出点对  $(\mathbf{x}_i, y_i), i = 1, \dots, m$ ，我们用这些样本来学习一个函数模型， $\mathbf{x} \rightarrow \hat{y}_{\mathbf{w}}(\mathbf{x})$ ，其中  $\mathbf{w}$  是模型中的参数向量。一旦得到这个模型，我们就可以预测：如果  $\mathbf{x}$  是一个新的输入点，并且也没有观察到它的输出，则我们预测其输出为  $\hat{y}(\mathbf{x}) = \hat{y}_{\mathbf{w}}(\mathbf{x})$ 。如果对预测的响应是任意实数，则我们可以假设输出是输入的线性函数(线性模型)，从而我们的模型形如  $\hat{y}_{\mathbf{w}} = \mathbf{w}^\top \mathbf{x}$ ，其中  $\mathbf{w} \in \mathbb{R}^n$  包含模型的所有系数。更一般地，我们假设  $\hat{y}_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$ ，其中  $\mathbf{x} \rightarrow \phi(\mathbf{x})$  是给定的非线性映射；置  $\phi(\mathbf{x}) = (1, \mathbf{x})$  可使这种一般性的表述涵盖仿射模型。

**例13.1. (库存管理中的需求量预测)** 一家大型商店需要预测来自特定地理区域的顾客对某种特定商品的需求量。假设需求量的对数是任意的实数，并且需求量的对数线性地依赖于某些输入(特征)：一年中的时刻，商品的类型，之前的销售量等等。这个问题是根据最近的刚刚过去的输入-需求对来预测未来的需求量。

有时候，需要预测二进制(binary)的响应，即对每个输入  $\mathbf{x}$ ， $y(\mathbf{x}) \in \{-1, 1\}$ ；我们将此称为一个**标签**(label)。这样，我们可以建立形如  $\hat{y}_{\mathbf{w}}(\mathbf{x}) =$

<sup>1</sup>参见2.1节中的文本bag-of-words表示。



$\text{sgn}(\mathbf{w}^\top \mathbf{x} + b)$ 的符号—线性模型，其中 $\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}$ 是模型的待定参数。类似地，我们可以用一种包含给定的非线性映射的更复杂的模型。这里不讨论分类问题中出现的其他响应类型，比如类别响应或者更复杂的对象(比如树)。

**例13.2. (医疗效果的分)** 二分类问题通常出现在医疗决策中，典型的场景是希望根据临床测量来预测治疗结果。威斯康新诊断乳腺癌数据(Wisconsin Diagnostic Breast Cancer, WDBC)集是公开数据库，该数据库包含了与 $m = 569$ 名病人有关的数据。对每个病人，我们有一个输入向量 $\mathbf{x}_i \in \mathbb{R}^{30}$ ，其由30个来自于FNA乳腺肿块的数字化映像的数值组成，其描述了数字影像所呈现的细胞核的特征，像半径，纹理，周长，面积，光滑度，紧凑性，中央凹陷，对称性和细胞核的分数维数。对于每个输入点，当诊断结果是恶性的时，对应的标签是 $+1$ ；当诊断为良性的时，标签是 $-1$ 。这里的目标是对已知根据进行训练以得到一个分类器，其可以对新来的病人进行诊断(恶性/良性诊断)。如果把新病人的各种临床测量结果表示为 $\mathbf{x} \in \mathbb{R}^{30}$ ，则我们可以根据这个新的数据 $\mathbf{x}$ 和通过学习得到的分类器的预测值对这个病人做出诊断。

**例13.3. (信用卡申请的二分类)** 一家信用卡公司收到上千份的信用卡申请单。每份申请包含了下列申请人信息：年龄，婚姻状况，每年的薪金，未偿还的债务等等。问题：决定是否通过一份申请，即将申请分成两类，通过和不通过。

**例13.4. (分类问题的其他例子)** 很多情况下都会出现二分类问题。在推荐系统中，比如说电影<sup>2</sup>，我们有特定用户的观影偏好信息；分类的目的是确定用户是否会喜欢一部给定的新电影。在垃圾邮件过滤系统中，我们可能会收到已确定是垃圾邮件的电子邮件子集，和另外一些已知是合法的电子邮件子集；分类的目的是确定新收到的一封邮件是垃圾邮件，还是合法邮件。在时间序列的预测中，我们试图基于过去观测到的价格来决定：相对与目前的价格而言，未来的价格(比如一支股票的价格)是增加还是减少。

**学习(或训练)**(learning(or training))问题指找到“最好的”模型的系数向量 $\mathbf{w}$ ，使得 $\hat{y}_{\mathbf{w}}(\mathbf{x}_i) \simeq y_i, i = 1, \dots, m$ 。因此，我们试图以模型变量 $\mathbf{w}$ 为决策变量，最小化预测向量 $\hat{\mathbf{y}}$ 和观测到的响应 $\mathbf{y}$ 的不一致性的某种度量。在实践中，我们希望对未看到的测试点 $\mathbf{x}$ 也能做出好的预测。后面会进一步讨论样本外性能(out-of-sample performance)的关键问题。

---

<sup>2</sup>参见11.4.1.4节。

如何把某种可靠性保证附加到预测输出  $\hat{y}$  中去是监督学习的另一个关键问题。这里的可靠性保证，换句话说，就是量化的新数据的**错误分类**(misclassification) 概率<sup>3</sup>。这个重要的问题是统计学习理论<sup>4</sup>的主题，其超出了这本书的范畴，将不再讨论。

我们首先考虑一个详细但是很基础的例子，其表明诸如正则化和样本外性能等基本概念是如何出现的。

### §13.2 基于多项式模型的最小二乘预测

考虑如图13.1所示的数据集，这里输入-输出对  $(x_i, y_i), x_i \in \mathbb{R}, i = 1, \dots, m$ 。我们的目标是预测与未看到的输入  $x$  对应的  $y$  的值。

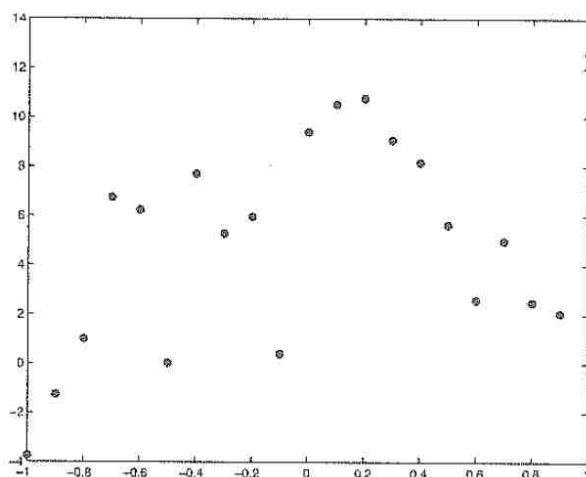


图 13.1: 输入-输出对为  $(x_i, y_i), i = 1, \dots, 20$ ，的数据集

#### §13.2.1 模型

**线性模型.** 我们先为这些数据建立基本的线性模型，即对输入  $x$ ，假定输出

$$y(x) = w_1 + xw_2 + e(x) = \mathbf{w}^\top \boldsymbol{\phi}(x) + e(x),$$

其中  $w_1, w_2$  是待定的权重， $\boldsymbol{\phi}(x) \doteq (1, x)$ ， $e(x)$  是误差项。为了确定权重向量  $\mathbf{w}$ ，我们用最小二乘法，此即导致所谓的**训练误差**(training error)最

<sup>3</sup>在例13.2的医疗背景下，错误分类概率对应于诊断误差。

<sup>4</sup>参加参考书，比如 T. Hastie, R. Tibshirani, and J. Friedman, The Elements of Statistical learning, Springer, 2008.

小化问题, 即

$$\min_{\mathbf{w}} \frac{1}{m} \sum_{i=1}^m (y_i - \phi_i^\top \mathbf{w})^2,$$

其中  $\phi_i \doteq \phi(x_i) = (1, x_i), i = 1, \dots, m$ . 上述问题用最小二乘可以表述为

$$\min_{\mathbf{w}} \|\Phi^\top \mathbf{w} - \mathbf{y}\|_2^2, \quad (13.1)$$

其中向量  $\mathbf{y} \doteq (y_1, \dots, y_m)$  是响应(response), 矩阵  $\Phi \in \mathbb{R}^{2 \times m}$  是以  $\phi_i (i = 1, \dots, m)$  为列向量的数据矩阵, 这里的特征包括置为1的常数特征( $\Phi$ 的第一行)和置为输入的第二特征( $\Phi$ 的第二行)。这个问题的目标涉及到所谓的平方损失(squared loss)函数

$$L(z, y) = \|z - y\|_2^2.$$

通常将这个使用已有数据来拟合一个特定模型的过程称为**训练**(training)阶段, 并且称问题(13.1)为训练(或学习)问题。一旦得到问题(13.1)的解 $\mathbf{w}^*$ , 我们就可以用向量 $\mathbf{w}^*$ 为任意的输入 $x$ 确定输出 $\hat{y}(x)$ , 即与 $x$ 对应的预测值, 即

$$\hat{y}(x) = \phi(x)^\top \mathbf{w}^* = w_1^* + w_2^* x.$$

如此选取 $\mathbf{w}$ 的合理性可以解释为:  $\mathbf{w}$ 已经使得已有数据的误差 $\mathbf{e} = \Phi^\top \mathbf{w} - \mathbf{y}$ 最小了; 对新的输入 $x$ , 我们希望上面的预测方法是准确的。在实践中, 除非能够对数据的生成方式做出一些其他的假定, 否则没法保证预测精度。

### 多项式模型

显然, 如果输入的数据 $x$ 接近于线性, 我们的预测算法会表现得很好。如果不是线性的, 我们可以尝试一种更复杂的模型, 比如形如

$$y(x) = w_1 + xw_2 + x^2w_3 + e(x) = \mathbf{w}^\top \phi(x) + e(x)$$

的二次模型. 这里  $\mathbf{w} \in \mathbb{R}^3$ ,  $\phi(x) = (1, x, x^2)$ 。同上, 求解最小二乘问题

$$\min_{\mathbf{w}} \|\Phi^\top \mathbf{w} - \mathbf{y}\|_2^2,$$

可以得到最合适的向量 $\mathbf{w}$ 。这里  $\mathbf{y} = (y_1, \dots, y_m)$ , 数据矩阵  $\Phi \in \mathbb{R}^{3 \times m}$  的列向量为  $\phi_i = (1, x_i, x_i^2), i = 1, \dots, m$ 。更一般地, 我们可以用 $k$ 次多项式模型(见图13.2), 即

$$y(x) = w_1 + xw_2 + x^2w_3 + \dots + x^kw_{k+1} + e(x).$$

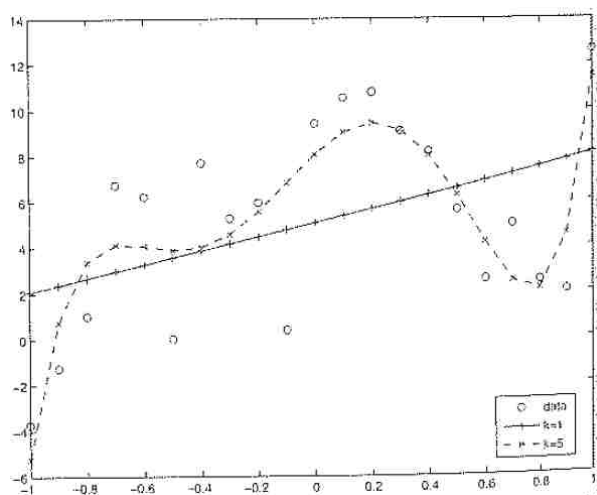


图 13.2: 两个不同的多项式模型，其中一个线性的，另一个是7次的

### §13.2.2 性能评价

我们的预测有多好？如何来比较这里的各种多项式模型？为了评价所得到的预测方法的性能，我们可以用一种称为**留一**(leave-one-out)的方法，即我们留出第 $j$ 个数据点 $(x_j, y_j)$ ，然后使用其余的数据对 $(x_i, y_i), i = 1, \dots, m, i \neq j$ ，来拟合模型，然后使用留出的第 $j$ 个数据点作为“测试点”（或者样本外），即用这个点来评价模型的性能。这样，我们把原始的数据点分成两部分，一部分是训练点的集合，其包含除了除第 $j$ 个点之外的所有数据点，剩下的是测试点集合，其只包含第 $j$ 个点。这种留一法可以推广到留出多个数据点的情况；另外一些方式将数据随机的剖分成训练集和测试集，后者通常包含 30% 的数据。这些不同的技术均可纳入**交叉验证**。

我们去除  $\mathbf{X}$  的最后一列和向量  $\mathbf{y}$  的最后一个分量，然后求解最小二乘问题。然后比较测试数据的实际值与对应的预测值，即可得到预测误差  $(\hat{y}(x_j) - y_j)^2$  这种留一误差与包含所有样本的误差可能有所不同。(13.1)的最优值度量了具体的不同。

我们让留一法中的指标  $j$  从 1 到  $m$  来重复上述过程，然后计算测试误差的平均值，即

$$\frac{1}{m} \sum_{j=1}^m (\hat{y}(x_j) - y_j)^2.$$

对每个  $j$ ，在预测算法中不会被看到对应的数据。因此，希望测试误差

变成零或者特别小都是不现实的。我们画出多项式模型的阶数所对应的平均测试和训练误差，得到图13.2。我们观察到：随着阶数  $k$  的增加，训练误差在减少，然而测试误差的一般是先减少，然后保持不变，然后再次增加。称这种现象为过拟合(over-fitting)。这种现象表明：高复杂度的模型来拟合训练数据时更好(即训练误差小)，但是一个太复杂的模型(比如一个高次多项式)对未知的数据拟合的不是很好。我们希望有这样一个“甜蜜点”，其对应的测试误差是最小的。

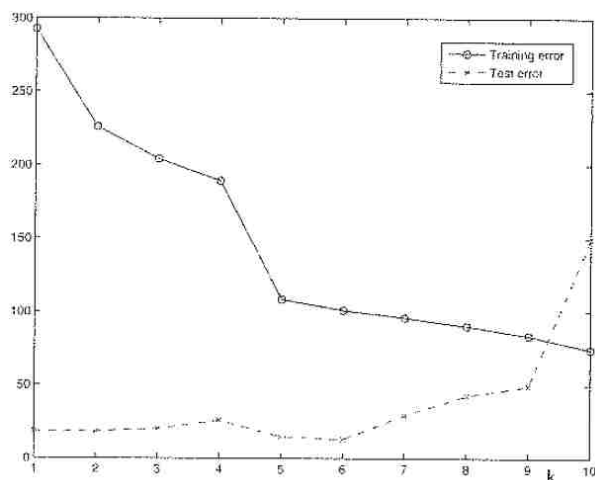


图 13.3: 训练误差和测试误差随多项式模型的阶数的变化

### §13.2.3 正则化和稀疏性

**正则化(Regularization).** 假设我们将多项式模型的阶数置为特定的  $k$ 。对于给定的阶数  $k$ ，所有  $k$  次多项式的变化不尽相同；有些变化非常剧烈，这对应着大的导数值。如果没有新输入点的精确值，多项式值的剧烈变化可能会导致严重的分类测试误差。

事实上，一个多项式的“复杂度”不仅依赖于阶数，还依赖于系数的大小(size)。更正式的表述<sup>5</sup>是：对多项式  $p_w(x) = w_1 + w_2x + \cdots + w_{k+1}x^k$ ，我们有

$$\forall x \in [-1, 1] : \left| \frac{d}{dx} p_w(x) \right| \leq k^{3/2} \|w\|_2. \quad (13.2)$$

这意味着，如果先验地知道数据是有界的(比如  $|x| \leq 1$ )，我们极小化系数向量  $w$  的欧氏范数，由此来控制这个给定阶数的多项式的导数。

<sup>5</sup>参见习题2.8.

这样，我们就可以给损失函数增加一个涉及  $\|\mathbf{w}\|_2$  的惩罚项，从而得到修正的最优化问题(13.1). 比如，我们用

$$\min_{\mathbf{w}} \|\Phi^T \mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

代替(13.1)，其中  $\lambda \geq 0$  是正则化参数。上述问题就是正则化的最小二乘问题. 利用6.7.3 节中的讨论，可以用线性代数方法给出这个问题的显式解。

选取好的参数  $\lambda$  的模式与确定模型的阶数的相同。事实上，为了同时选取两个参数(正则化参数  $\lambda$  和模型的阶数  $k$ )，需要在对应的二维空间上进行搜索。如图13.4所表明的，在很多情况下，额外的自由度确实会提高模型的测试误差。

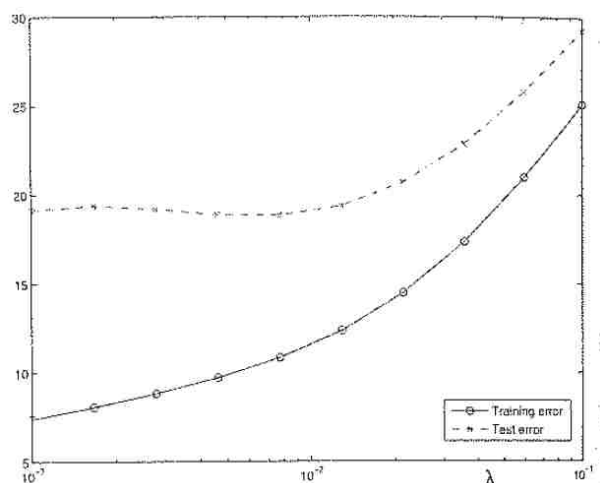


图 13.4: 固定阶数  $k = 10$  时，训练误差和测试误差随正则化参数  $\lambda$  的演变

欧氏范数的平方并不是唯一的用来控制多项式导数的方法。事实上，我们可以用(13.2)的许多变形；比如

$$\forall x \in [-1, 1] : \left| \frac{d}{dx} p_{\mathbf{w}}(u) \right| \leq k \|\mathbf{w}\|_1,$$

这对应于惩罚问题

$$\min_{\mathbf{w}} \|\Phi^T \mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_1. \quad (13.3)$$

上述问题是一个LASSO问题，已经在9.6.2节对其进行了广泛的讨论。

在训练问题中, 历史上一直偏好选欧氏范数的平方, 因为欧氏范数的平方是因为由此得到问题可以用直接用线性代数方法求解。与此相反, LASSO问题得不到不能直接用线性代数方法求解。随着凸优化方法的出现, LASSO 和相关的变形变得从计算上可以与正则化最小二乘法相媲美。使用非欧氏范数不仅得到可处理的问题, 下面还会看到它同时还提供了一些别的有用的选择。

**稀疏性(Sparsity).** 在有些情况下, 获得稀疏的模型是很理想的, 这里稀疏的模型指系数向量  $\mathbf{w}$  有许多零分量的模型。就像在例9.11中讨论的一样, 在惩罚项中使用  $\ell_1$ -范数能带来最优解向量  $\mathbf{w}$  的稀疏性。例如, 解决(13.3)的LASSO问题, 其中  $\lambda = 0.1$ , 产生了稀疏的多项式

$$p_{\mathbf{w}}(x) = 7.4950 + 10.7504x - 6.9644x^2 - 45.0750x^3 + 42.9250x^5 - 5.4516x^8 - 0.4539x^{15} + 9.2869x^{20},$$

图13.5画出了这个多项式和最小二乘确定的多项式的函数图形及数据点。由图可以看出, 使用  $\ell_1$ -范数作为惩罚项让我们找到了一种稀疏的多项式, 并且这个多项式在数据所在的区间  $[-1, 1]$  中的变化也不太剧烈。

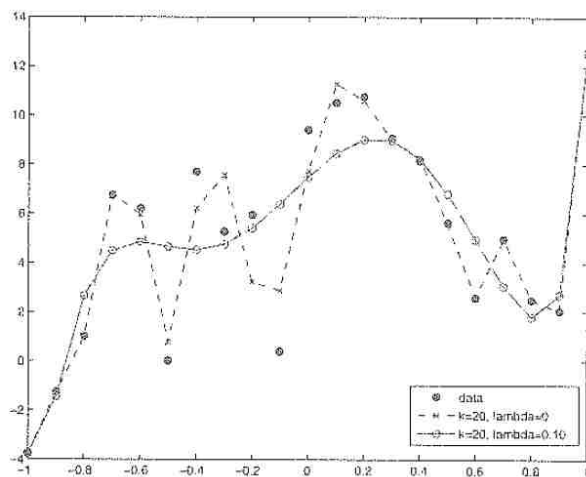


图 13.5: 在LASSO问题(13.3)中取  $k = 20, \lambda = 0.1$  时所得到的稀疏多项式模型

### §13.3 二分类

考虑二分类问题, 假设我们有  $m$  个数据点  $\mathbf{x}_i$  和对应的标签  $y_i \in$

$\{-1, 1\}, i = 1, \dots, m$ 。我们的目标是为未知的数据点  $\mathbf{x}$  预测对应的标签  $\hat{y}$ 。

### §13.3.1 支持向量机

在13.2节，我们需要预测一个实值变量，我们最初用了仿射函数  $x \rightarrow \mathbf{w}^\top \phi(x)$ ，其中  $\phi(x) = (1, x)$ 。当待预测的标签  $y \in \{-1, 1\}$ ，可能是最简单的修正是形如

$$\hat{y} = \text{sgn}(\mathbf{x}^\top \mathbf{w} + b)$$

的预测规则。一开始，我们可能试着去找使得训练误差的平均值取到最小值的  $\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}$ ，这里的训练误差可以用 0/1 函数正式地表述为

$$E(\alpha) \doteq \begin{cases} 1 & \text{如果 } \alpha < 0 \\ 0 & \text{否则.} \end{cases}$$

对给定的输入-输出对  $(\mathbf{x}, y)$ ，当且仅当  $y(\mathbf{w}^\top \mathbf{x} + b) < 0$  时，这个分类规则  $\mathbf{x} \rightarrow \text{sgn}(\mathbf{w}^\top \mathbf{x} + b)$  才会产生误差。因此，对给定的  $\mathbf{w}$  我们在训练集上观测到的误差的平均值是

$$\frac{1}{m} \sum_{i=1}^m E(y_i(\mathbf{w}^\top \mathbf{x}_i + b)).$$

由此得到了形如

$$\min_{\mathbf{w}, b} \frac{1}{m} \sum_{i=1}^m E(y_i(\mathbf{x}_i^\top \mathbf{w} + b)) \quad (13.4)$$

的训练问题。不幸的是，上面的问题不是凸的，从而算法可能收敛于局部极小点，这让人很不满意。

为了克服这个问题，我们观察到上面的函数  $E$  不大于凸的铰链(hinge)函数  $\alpha \rightarrow \max(0, 1 - \alpha)$ 。由此得到训练问题的凸近似(实际上是上界)，即

$$\min_{\mathbf{w}, b} \frac{1}{m} \sum_{i=1}^m \max(0, 1 - y_i(\mathbf{x}_i^\top \mathbf{w} + b)). \quad (13.5)$$

称这个目标函数是**铰链损失**(hinge loss)函数，它是凸函数，也是多面集的；因此可将上述问题重新表述成线性规划。一个可能的线性规划表述是

$$\begin{aligned} & \underset{\mathbf{w}, b, e}{\text{minimize}} && \frac{1}{m} \sum_{i=1}^m e_i \\ & \text{subject to} && e \geq 0, i = 1, \dots, m, \\ & && e_i \geq 1 - y_i(\mathbf{x}_i^\top \mathbf{w} + b), i = 1, \dots, m. \end{aligned} \quad (13.6)$$



求解这个模型, 设最优解是点对  $(\mathbf{w}^*, b^*)$ , 则对新数据点  $\mathbf{x} \in \mathbb{R}^n$  的预测标签为  $\hat{y} = \text{sgn}(\mathbf{w}^\top \mathbf{x} + b)$ 。(13.5) 是所谓的支持向量机模型的基本构件。支持向量集的命名源于 13.3.3 节讨论的几何解释。

### §13.3.2 正则化与稀疏性

像 13.2 节的例子那样, 我们对控制模型的复杂度感兴趣。一种方法是控制线性函数  $\mathbf{x} \rightarrow \mathbf{x}^\top \mathbf{w} + b$  在  $\mathbf{x}$  取遍观测到的输入时的变化量。易见该函数的梯度就是  $\mathbf{w}$ , 因此控制  $\mathbf{w}$  的大小就能实现我们的目标。

假设这些数据点都包含在一个以原点为圆心, 半径为  $R$  的欧氏球里, 因为

$$\max_{\mathbf{x}, \mathbf{x}': \|\mathbf{x}\|_2 \leq R, \|\mathbf{x}'\|_2 \leq R} |\mathbf{w}^\top (\mathbf{x} - \mathbf{x}')| \leq 2R \|\mathbf{w}\|_2,$$

所以我们用  $\ell_2$ -范数来正则化, 相应的学习问题就变成了

$$\min_{\mathbf{w}, b} \frac{1}{m} \sum_{i=1}^m \max(0, 1 - y_i(\mathbf{x}_i^\top \mathbf{w} + b)) + \lambda \|\mathbf{w}\|_2^2, \quad (13.7)$$

其中  $\lambda \geq 0$  是正则化参数。通常用留一法或者类似的方法选取该参数。

假设数据点都包含在大小为  $R$  的盒子里, 即  $\{\mathbf{x} : \|\mathbf{x}\|_\infty \leq R\}$ , 则相应的边界

$$\max_{\mathbf{x}, \mathbf{x}': \|\mathbf{x}\|_\infty \leq R, \|\mathbf{x}'\|_\infty \leq R} |\mathbf{w}^\top (\mathbf{x} - \mathbf{x}')| \leq 2R \|\mathbf{w}\|_1$$

则可将上述问题修正为

$$\min_{\mathbf{w}, b} \frac{1}{m} \sum_{i=1}^m \max(0, 1 - y_i(\mathbf{x}_i^\top \mathbf{w} + b)) + \lambda \|\mathbf{w}\|_1, \quad (13.8)$$

其中  $\lambda \geq 0$  仍然是正则化参数。 $\ell_1$ -范数在上述问题中使得最优的向量  $\mathbf{w}$  具有稀疏性。对这样一个向量, 标量  $\mathbf{w}^\top \mathbf{x}$  仅仅包含了  $\mathbf{x}$  的一部分分量。因此这个技巧能够来辨别一些有助于得到好的预测的关键特征( $\mathbf{x}$  的元素)。

### §13.3.3 几何解释

基本的支持向量机模型(13.5)和它的正则化变形有许多几何解释。

**线性可分数据**(linearly separable data). 假设问题(13.5)的最优值是零, 蕴含着训练误差为零是可达的。这表明存在向量  $\mathbf{w} \in \mathbb{R}^n$  和标量  $b$  使得

$$y_i(\mathbf{x}_i^\top \mathbf{w} + b) \geq 0, i = 1, \dots, m. \quad (13.9)$$

从几何的角度来说,这意味着超平面  $\{x : \mathbf{w}^\top \mathbf{x} + b = 0\}$  可以将数据**线性可分的**(linear separable). 这里线性可分指标签  $y_i = 1$  对应的数据点  $\mathbf{x}_i$  在这个超平面的一侧, 标签  $y_i = -1$  对应的数据点  $\mathbf{x}_i$  在另一侧, 如图13.6所示。这里决策规则  $\hat{y} = \text{sgn}(\mathbf{w}^\top \mathbf{x} + b)$  即对应着辨别新点  $\mathbf{x}$  落在超平面的哪一侧。

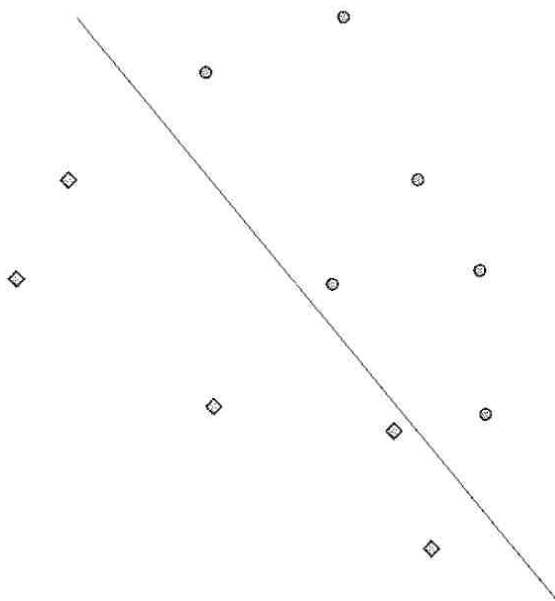


图 13.6: 严格线性可分的数据集

接下来假设是数据是**严格线性可分的**(strict linearly separable), 也就是说, 条件(13.9)中不等式严格成立。这种情况出现当且仅当存在正数  $\beta > 0$  使得

$$y_i(\mathbf{x}_i^\top \mathbf{w} + b) \geq \beta, i = 1, \dots, m.$$

在这种情形下, 我们可以给上述不等式两边同时除以  $\beta$  以规范化模型参数  $(\mathbf{w}, b)$ , 得出严格的线性可分性等价于存在  $(\mathbf{w}, b)$  使得

$$y_i(\mathbf{x}_i^\top \mathbf{w} + b) \geq 1, i = 1, \dots, m. \quad (13.10)$$

两个“临界”(critical)超平面  $\mathcal{H}_\pm = \{\mathbf{x} : \mathbf{w}^\top \mathbf{x} + b = \pm 1\}$  勾画出一个平板(slab). 没有点落在这个厚板内, 同时这个厚板也把标签分别为正的和负的数据点分开了, 如图13.7所示。

**最大间隔分类**(Maximum margin classifier). 在之前考虑的严格线性可分数据中, 存在无限多的  $(\mathbf{w}, b)$  满足条件(13.10), 即有无穷多个平

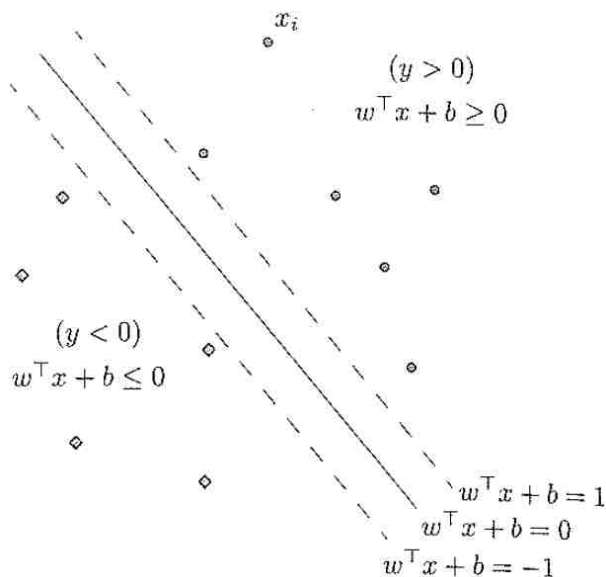


图 13.7: 严格线性可分数据集和一个分离平板

板可以分开数据集。解决这种解不惟一问题的合理作法是寻找分离间隔(separation margin)最大的分离平板, 这里分离间隔指两个“临界”超平面  $\mathcal{H}_{\pm} = \{x : w^T x + b = \pm 1\}$  之间的距离。最大化间隔的超平面无疑会具有好的性质: 对于一个非常接近正(或者负)标签数据的测试点  $x$  而言, 它将以更大的概率落在决策超平面相应的一侧。

易于验证两个临界超平面  $\mathcal{H}_+, \mathcal{H}_-$  之间的距离是  $2/\|w\|_2$ . 在满足标签约束的前提下, 为了最大化间隔, 我们就需要最小化  $\|w\|_2$ . 由此即得最优化问题

$$\begin{aligned} & \underset{w, b}{\text{minimize}} && \|w\|_2 \\ & \text{subject to} && y_i(x_i^T w + b) \geq 1, \quad i = 1, \dots, m. \end{aligned} \quad (13.11)$$

一般称此为**最大间隔SVM**(Maximum margin SVM), 见图13.8。**支持向量**(support vector)的名字源于这样的事实: 在实践中, 只有很少一些点落在临界边界  $\mathcal{H}_{\pm}$  上。粗略地说, 如果数据点是对某种连续分布的随机抽取得到的, 则超过  $n$  个特征点落在临界超平面上的概率为零。称这些落在临界超平面上的特定点为**支持向量**(support vectors)。

**不可分数据**(Non-separable data.) 当数据不可分(图13.9), 我们将在严格可分的条件(13.10)下引入“松弛变量”, 这些松弛变量表示对违返约

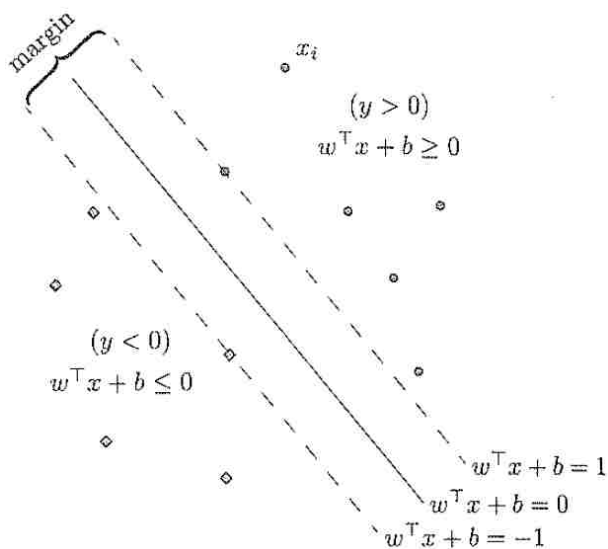


图 13.8: 最大间隔分类器

束的惩罚, 这时的变量为  $(w, b)$  和误差向量  $e$ , 与之对应的条件是

$$e \geq 0, \quad y_i(x_i^\top w + b) \geq 1 - e_i, \quad i = 1, \dots, m.$$

理想的作法是最小化  $e$  中非零元素的个数(number), 也就是向量  $e$  的基数。这将涉及非凸问题。然而, 我们可以近似这个问题, 用  $\ell_1$ -范数代替基数函数。因为  $e \geq 0$ ,  $\ell_1$ -范数的减少对应着  $e$  的元素之和的减少。通过这种方式, 我们精确地得到了(13.6)中表述的基本SVM问题。

**软间隔分类**(Classification with soft margin.) 在不可分的情况下, 我们需要找到一种间隔(两个临界超平面间的距离)和训练集上的分类误差之间的折衷。这就导致了形如(13.7)的问题。用松弛变量, 我们可将(13.7)重新表述成

$$\begin{aligned} & \underset{w, b, e}{\text{minimize}} && \frac{1}{m} \sum_{i=1}^m e_i + \lambda \|w\|_2^2 \\ & \text{subject to} && e \geq 0, e_i \geq 1 - y_i(x_i^\top w + b), \quad i = 1, \dots, m. \end{aligned} \quad (13.12)$$

这里的目标函数代表着分类间隔(我们希望大一些)和所有违反项的和(我们希望小一些)之间的折衷。针对这些原因, 称这种方法是软间隔分类。

由于当初的构造方式, 问题(13.12)总是可行的。如果  $(w^*, b^*, e^*)$  是一个最优解, 并且  $e^* = 0$ , 则  $(w^*, b^*)$  也是最大化间隔问题(13.11)的最

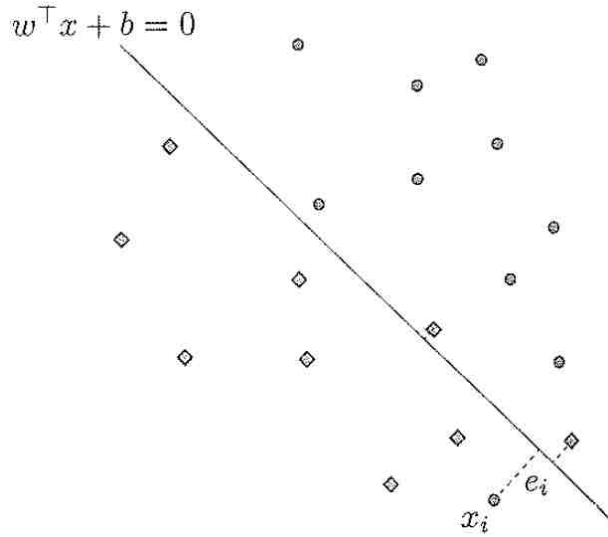


图 13.9: 不可分数据集的近似分离.

优解。如果  $e^* \geq 0, e^* \neq 0$ ，那么最大化间隔问题是不可行的，即数据是不可分的。在这种情况下，对  $e_i^*$  非零的所有  $i$ ，问题(13.12)中对应的约束在最优解处必须以等式成立，否则我们能够减小这些非最优的  $e_i^*$ 。如图13.10所示，这些  $e_i^*$  表示数据点  $x_i$  到超平面  $y_i(w^\top x + b) = 1$  的规范化距离。与  $e_i^* \in (0, 1)$  相应的点  $x_i$  落在间隔带内，但是仍然在决策边界正确的一侧；与  $e_i^* > 1$  对应的点落在决策边界错误的一侧，即被错误分类。

#### §13.3.4 鲁棒性

下面，我们从鲁棒性的角度解释惩罚式变形(13.7)和(13.8)。假设对每个  $i = 1, \dots, m$ ，仅知道数据点  $x_i$  落在以  $\hat{x}_i$  为中心，以  $\rho$  为半径的给定的球  $S_i$  内。我们首先来寻找数据点能被分离的条件。这里的分离指不管数据点在球上的精确位置，即要求数据点位于给定球内的任何位置时，都要求是可分的。然后我们尝试最大化鲁棒性水平，即球的半径。

对给定的模型参数向量  $(w, b)$ ，条件

$$y_i(x_i^\top w + b) \geq 0, \quad i = 1, \dots, m, \quad \forall x_i \in S_i$$

成立当且仅当

$$y_i(\hat{x}_i^\top w + b) \geq \rho \|w\|_2, \quad i = 1, \dots, m.$$

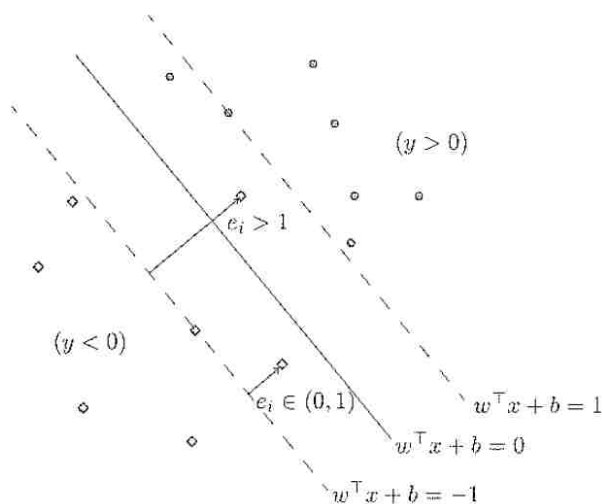


图 13.10: 软间隔分类.

由于上面的条件关于  $(w, b)$  是齐次的, 我们总是可以假设  $\rho \|w\|_2 = 1$ 。然后, 最大化  $\rho$  就相当于最小化  $\|w\|_2$ 。用  $x_i$  代替中心  $\hat{x}_i$ , 我们就得到了问题(13.11)。几何解释是我们分离了球而不是点。

类似的解释也可适用于  $\ell_1$ -范数, 此时假设数据未知但是有界, 即落在超立方体  $\{x : \|x - \hat{x}\|_\infty \leq \rho\}$  里。

### §13.3.5 Logistic回归

Logistic回归是支持向量机的一种变形, 即用光滑凸的对数函数  $\alpha \rightarrow \ln(1 + e^{-\alpha})$  代替SVM中的Hinge函数  $\alpha \rightarrow \max(0, 1 - \alpha)$ 。

**模型.** 与支持向量机的学习问题(13.5)相对照的Logistic回归是

$$\min_{w, b} \sum_{i=1}^m \ln(1 + e^{-y_i(\mathbf{x}_i^\top w + b)}). \quad (13.13)$$

当然, 在实践中通常使用惩罚变形, 即在目标中加上包含  $w$  的范数的惩罚项。因为目标中的log-sum-exp是凸的<sup>6</sup>。Logistic回归模型的一个优点是它们自然地推广到多分类问题, 此时标签的取值不是二进制的, 而是取有限个值。我们这里不讨论这种推广。

**概率解释.** SVMs中采用Logistic回归的主要优点是它与对应的有一个概率模型。这个模型允许为新数据点附加上属于正负类别的概率值。

<sup>6</sup>参见习题 2.14

这是它与SVMs的不同之处, 因为SVMs仅为新点的类别提供了“是或否”的答案。Logistic模型的解释如下. 假定一个点 $\mathbf{x}$ 的标签为  $y \in \{-1, 1\}$  的概率形如

$$\pi_{\mathbf{w}, b}(\mathbf{x}, y) = \frac{1}{1 + e^{-y(\mathbf{w}^\top \mathbf{x} + b)}}.$$

几何上可以清楚地解释这个概率值. 如果一个点 $\mathbf{x}$ 远离由  $\mathbf{w}^\top \mathbf{x} + b = 0$  定义的超平面, 且这个点所在区域对应于正标签, 则  $\mathbf{w}^\top \mathbf{x} + b$  的数值就会很大, 这时上面定义的概率就接近于 1. 反过来, 如果  $\mathbf{w}^\top \mathbf{x} + b$  是负数, 且绝对值很大, 这个点  $\mathbf{x}$  就落在超平面  $\mathbf{w}^\top \mathbf{x} + b = 0$  的另一侧, 且如上定义的概率值就接近于 0.

我们采用极大似然法来拟合这个模型. 这里的似然是上述模型为每个数据点  $\mathbf{x}_i$  得到观测标签  $y_i$  所赋予的概率, 即

$$\prod_{i: y_i = +1} \pi_{\mathbf{w}, b}(\mathbf{x}_i, 1) \prod_{i: y_i = -1} \pi_{\mathbf{w}, b}(\mathbf{x}_i, -1).$$

取负对数, 并关于参数  $(\mathbf{w}, b)$  来极小化, 即得到问题(13.13)。

一旦观测到新数据点  $\mathbf{x}$ , 我们计算概率  $\pi_{\mathbf{w}, b}(\mathbf{x}, \pm 1)$ , 并以  $\pi_{\mathbf{w}, b}(\mathbf{x}, 1)$  和  $\pi_{\mathbf{w}, b}(\mathbf{x}, -1)$  中最大者对应的  $y$  作为  $\mathbf{x}$  的预测标签. 这类似于SVM中置  $\hat{y} = \text{sgn}(\mathbf{w}^\top \mathbf{x} + b)$ , 但Logistic回归有了这种额外的优点, 使得我们可以在我们的预测中附加上概率水平。

### §13.3.6 Fisher判别

目前解决二分类问题的另一种方法是Fisher判别(Fisher discrimination)。像我们将看到的那样, 用线性代数的方法即可解决这种方法中所涉及的计算问题. 费舍判别是要在数据空间找到一个方向  $\mathbf{u} \in \mathbb{R}^n$ , 并且使得分别属于正标签和负标签的两类数据各自在方向  $\mathbf{u}$  上的投影分得越开越好. 这里需要度量每类投影数据之间的距离相对于各自数据类的变化性。

为了便于表述, 我们分别定义与正负类别相对应的两个矩阵

$$\mathbf{A}_+ = [\mathbf{x}_i]_{i: y_i = +1} \in \mathbb{R}^{n \times m_+}, \quad \mathbf{A}_- = [\mathbf{x}_i]_{i: y_i = -1} \in \mathbb{R}^{n \times m_-},$$

其中  $m_{\pm} = \text{card}\{i : y_i = \pm 1\}$  表示每个类别的大小. 类似地, 我们用  $\bar{\mathbf{a}}_+, \bar{\mathbf{a}}_-$  表示两个类别的质心(centroid), 即

$$\bar{\mathbf{a}}_{\pm} = \frac{1}{m_{\pm}} \mathbf{A}_{\pm} \mathbf{1}.$$

最后, 令

$$\tilde{\mathbf{A}}_{\pm} = \mathbf{A}_{\pm} \left( \mathbf{I}_{\pm} - \frac{1}{m_{\pm}} \mathbf{1}\mathbf{1}^{\top} \right).$$

表示中心化的数据矩阵.

这里, 两个类别的质心在方向  $\mathbf{u}$  上的距离的平方可以表示为  $(\mathbf{u}^{\top}(\bar{\mathbf{a}}_+ - \bar{\mathbf{a}}_-))^2$ , 我们把这个量作为分离距离的一种度量, 并用数据的均方差来规范化. 数据集  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_p]$  关于平均值  $\bar{\mathbf{a}}$  在方向  $\mathbf{u}$  上的均方差即

$$\begin{aligned} s^2 &= \frac{1}{p} \sum_{i=1}^p (\mathbf{u}^{\top}(\mathbf{a}_i - \bar{\mathbf{a}}))^2 = \frac{1}{p} \sum_{i=1}^p \mathbf{u}^{\top} \tilde{\mathbf{a}}_i \tilde{\mathbf{a}}_i^{\top} \mathbf{u} \\ &= \frac{1}{p} \mathbf{u}^{\top} \left( \sum_{i=1}^p \tilde{\mathbf{a}}_i \tilde{\mathbf{a}}_i^{\top} \right) \mathbf{u} \\ &= \frac{1}{p} \mathbf{u}^{\top} \tilde{\mathbf{A}} \tilde{\mathbf{A}}^{\top} \mathbf{u}. \end{aligned}$$

将两个类别的质心沿方向  $\mathbf{u}$  的距离的平方进行规范化, 得判别准则

$$f_0(\mathbf{u}) = \frac{\mathbf{u}^{\top}(\bar{\mathbf{a}}_+ - \bar{\mathbf{a}}_-)(\bar{\mathbf{a}}_+ - \bar{\mathbf{a}}_-)^{\top} \mathbf{u}}{\mathbf{u}^{\top} \left( \frac{1}{m_+} \tilde{\mathbf{A}}_+ \tilde{\mathbf{A}}_+^{\top} + \frac{1}{m_-} \tilde{\mathbf{A}}_- \tilde{\mathbf{A}}_-^{\top} \right) \mathbf{u}} = \frac{(\mathbf{u}^{\top} \mathbf{c})^2}{\mathbf{u}^{\top} \mathbf{M} \mathbf{u}},$$

其中  $\mathbf{c} = \bar{\mathbf{a}}_+ - \bar{\mathbf{a}}_-$ ,  $\mathbf{M} = \frac{1}{m_+} \tilde{\mathbf{A}}_+ \tilde{\mathbf{A}}_+^{\top} + \frac{1}{m_-} \tilde{\mathbf{A}}_- \tilde{\mathbf{A}}_-^{\top}$ .

接下来, 我们需要以  $\mathbf{u}$  为变量来极大化判别准则  $f_0(\mathbf{u})$ , 这实质上是极大化两个二次型的比值, 即

$$\max_{\mathbf{u} \neq \mathbf{0}} \frac{(\mathbf{u}^{\top} \mathbf{c})^2}{\mathbf{u}^{\top} \mathbf{M} \mathbf{u}}.$$

由于问题的目标函数具有齐次性, 我们总可以重新调整  $\mathbf{u}$  使得  $\mathbf{u}^{\top} \mathbf{c} = 1$ , 于是得到最小化问题

$$\begin{aligned} &\underset{\mathbf{u}}{\text{minimize}} && \mathbf{u}^{\top} \mathbf{M} \mathbf{u} \\ &\text{subject to} && \mathbf{u}^{\top} \mathbf{c} = 1. \end{aligned}$$

这里对称矩阵  $\mathbf{M}$  是半正定的, 因此上面的问题是凸二次规划. 简单起见, 我们假设  $\mathbf{M} \succ \mathbf{0}$ . 为了求解上述的凸优化问题, 我们用拉格朗日乘子法, 即由Lagrange函数

$$L(\mathbf{u}, \mu) = \frac{1}{2} \mathbf{u}^{\top} \mathbf{M} \mathbf{u} + \mu(1 - \mathbf{u}^{\top} \mathbf{c}),$$

关于变量  $(\mathbf{u}, \mu)$  的稳定点条件

$$0 = \nabla_{\mathbf{u}} L(\mathbf{u}, \mu) = \mathbf{M} \mathbf{u} - \mu \mathbf{c}, \quad \mathbf{u}^{\top} \mathbf{c} = 1.$$



来表征解<sup>7</sup>。因为  $M$  是正定的，我们得到  $\mathbf{u} = \mu M^{-1}\mathbf{c}$ ，将其带入约束  $\mathbf{u}^\top \mathbf{c} = 1$ ，我们得到

$$\mathbf{u} = \frac{1}{\mathbf{c}^\top M^{-1} \mathbf{c}} M^{-1} \mathbf{c}.$$

图13.11的左半边展示了  $\mathbb{R}^2$  中的两组数据云，箭头代表着最优的Fisher判别方向  $\mathbf{u}$ 。图13.11的右半边表明：数据沿着方向  $\mathbf{u}$  投影后的直方图被很好地分开了。

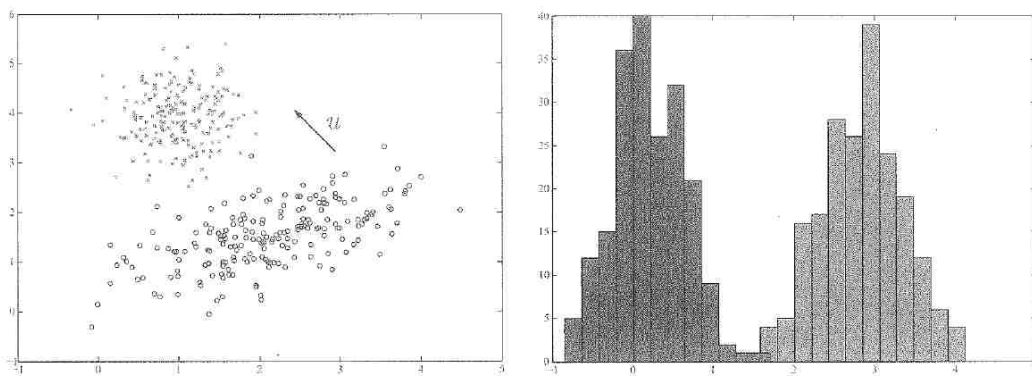


图 13.11: Fisher数据判别

### §13.4 一种通用的监督学习问题

在这一节，我们考虑一种形如

$$\min_{\mathbf{w}} L(\Phi^\top \mathbf{w}, \mathbf{y}) + \lambda p(\mathbf{w}) \quad (13.14)$$

的通用监督学习模型，其中  $\mathbf{w} \in \mathbb{R}^n$  包含了模型的参数； $p$  是惩罚函数(通常是一种范数或者范数的平方)； $\lambda$  是正则化参数；函数  $L$  代表损失(loss)函数，用来度量预测值和响应  $y$  之间的差异。在上述情况下，矩阵  $\Phi$  包含着某种变换之后的数据点；准确来说，每列  $\phi_i$  形如  $\phi_i = \phi(\mathbf{x}_i)$ ，其中  $\phi$  是某种由用户选定的映射(可能是非线性的)<sup>8</sup>。

这种通用的问题类还远不能涵盖所有的监督学习的方法，但它涵盖了最小二乘回归，SVM，Logistic回归和更多的方法。

与这类学习问题相关的预测规则取决于任务和模型。在回归中，输出的预测值是一个实数，并且预测规则形如  $\hat{y} = \phi(\mathbf{x})^\top \mathbf{w}$ 。在二分类中，

<sup>7</sup>参见8.5节。

<sup>8</sup>参见13.2节中关于映射  $\phi$  的说明

我们使用符号线性函数  $\hat{y} = \text{sgn}(\phi(\mathbf{x})^\top \mathbf{w})$ 。我们用符号  $\hat{y} = S(\phi(\mathbf{x})^\top \mathbf{w})$  来表示所有这些情况, 其中  $S$  是 1 或者是符号函数。

### §13.4.1 损失函数

损失函数  $L$  的形式是基于任务的(回归或者分类), 也就是说, 是基于待预测的输出值的属性, 还有其他因素, 比如对这些数据的先验假设。通常假设损失函数可以分解成求和的形式, 即

$$L(z, y) = \sum_{i=1}^m l(z_i, y_i),$$

其中  $l$  是某(通常是凸的)函数。表13.1给出了标准选择。

表 13.1: 标准的损失函数

损失函数	$l(z, y)$	任务
欧氏范数的平方	$\ z - y\ _2^2$	回归
$\ell_1$ -范数	$\ z - y\ _1$	回归
$\ell_\infty$ -范数	$\ z - y\ _\infty$	回归
Hinge函数	$\max(0, 1 - yz)$	二分类
Logistic函数	$\ln(1 + e^{-yz})$	分类

选择特定的损失函数需要基于任务和对数据的假设, 此外实践中还应考虑可用的软件(用其求解学习问题(13.14))等。

### §13.4.2 惩罚和约束函数

惩罚函数  $L$  有好多种, 具体的选择取决于任务和其它假设。历史上偏好选  $\ell_2$ -范数的平方是因为对应的最优化问题(13.14) 继承了诸如光滑性(当损失函数自身是光滑的时)等好的性质。

**选取一个惩罚(Choosing a penalty).** 如果我们认为最优的向量  $\mathbf{w}$  是稀疏的, 即只有很少一部分特征主导了对输出的预测, 则可用  $\ell_1$ -范数惩罚来提高稀疏性。如果向量  $\mathbf{w}$  由分块形式  $\mathbf{w}^{(i)}, i = 1, \dots, p$ , 如果我们认为许多块都应该是零向量, 则可以使用复合范数

$$\sum_{i=1}^p \|\mathbf{w}^{(i)}\|_2.$$

事实上,上述函数是由范数组成的向量  $(\|\mathbf{w}^{(1)}\|_2, \dots, \|\mathbf{w}^{(p)}\|_2)$  的  $\ell_1$ -范数。我们希望提高这个向量的稀疏性。

实践中很重要的一种考量是:设计的惩罚应该使得学习问题(13.14)有惟一解。为此,通常会添加  $\ell_2$ -范数的平方的小的倍数,这能确保目标函数是强凸的。

**约束版(Constrained Versions).** 注意到惩罚形式(13.14)涵盖了形如

$$\min_{\mathbf{w}} \{L(\Phi^\top \mathbf{w}, \mathbf{y}) : \mathbf{w} \in C\}$$

的约束式变形,这里  $C$  是一个(通常是凸的)集合,比如欧氏球。如果恰当地选取  $p$  为

$$p(\mathbf{w}) = \begin{cases} 1 & \text{若 } \mathbf{w} \in C \\ 0 & \text{否则} \end{cases}$$

就可将上述这个约束形式的学习问题表述成惩罚形式的问题。你也许好奇使用给定的惩罚,这个惩罚的平方,及与之相对应的约束版在实际中的差异<sup>9</sup>。为了说明这个问题,我们考虑LASSO和两种变形,即

$$\begin{aligned} p_{\text{lasso}}^* &\doteq \min_{\mathbf{w}} \|\Phi^\top \mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_1, \\ p_{\text{sqr}}^* &\doteq \min_{\mathbf{w}} \|\Phi^\top \mathbf{w} - \mathbf{y}\|_2 + \mu \|\mathbf{w}\|_1, \\ p_{\text{constr}}^* &\doteq \min_{\mathbf{w}} \{\|\Phi^\top \mathbf{w} - \mathbf{y}\|_2 : \|\mathbf{w}\|_1 \leq \alpha\}, \end{aligned}$$

其中  $\lambda \geq 0, \mu \geq 0, \alpha \geq 0$  是正则化参数。对一个固定的三元对  $(\lambda, \mu, \alpha)$ ,与上述问题相对应的解通常是不同的。然而,当参数  $\lambda, \mu, \alpha$  遍历所有非负实数时,则上述每个最优化问题的解所生成的路径对于同一个实际问题(即数据  $\Phi$  和  $\mathbf{y}$ )而言是相同的。

### §13.4.3 Kernel方法

倘若学习问题是欧氏范数的平方的惩罚式

$$\min_{\mathbf{w}} L(\Phi^\top \mathbf{w}, \mathbf{y}) + \lambda \|\mathbf{w}\|_2^2 \quad (13.15)$$

时,Kernel方法使得非线性法则的背景下也能够进行快速计算。我们发现第一项涉及损失,且仅在  $\Phi^\top \mathbf{w}$  中依赖于  $\mathbf{w}$ 。由线性代数基本定理(见3.2.4节),可将任一向量  $\mathbf{w} \in \mathbb{R}^n$  分解为  $\mathbf{w} = \Phi \mathbf{v} + \mathbf{r}$ , 其中  $\Phi^\top \mathbf{r} = \mathbf{0}$ 。从而可

<sup>9</sup>参见习题13.5.

将  $\mathbf{w}$  的欧氏范数的平方分解成  $\|\Phi\mathbf{v}\|_2^2 + \|\mathbf{r}\|_2^2$ 。将  $\mathbf{w}$  和  $\|\mathbf{w}\|_2^2$  的这种分解代入学习问题(13.15)，得

$$\min_{\mathbf{v}, \mathbf{r}} L(\Phi^\top \Phi \mathbf{v}, \mathbf{y}) + \lambda(\|\Phi\mathbf{v}\|_2^2 + \|\mathbf{r}\|_2^2).$$

显然最优的  $\mathbf{r}$  是零向量。从几何上看，这意味着最优的  $\mathbf{w}$  属于数据  $\Phi$  生成的空间，即存在  $\mathbf{v}$  使得最优解  $\mathbf{w} = \Phi\mathbf{v}$ 。此时，学习问题(13.15)就变成了

$$\min_{\mathbf{v}} L(\mathbf{K}\mathbf{v}, \mathbf{y}) + \lambda \mathbf{v}^\top \mathbf{K}\mathbf{v}, \quad (13.16)$$

我们称这里的  $\mathbf{K} = \Phi^\top \Phi \in \mathbb{R}^{m \times m}$  为核矩阵(kernel matrix)。因此，上述问题仅依赖于标量积  $K_{ij} = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$ ，其中  $1 \leq i, j \leq m$ 。

预测规则是相似的：如果  $\mathbf{x}$  是一个新数据点，则预测值

$$\hat{y} = S(\phi(\mathbf{x})^\top \mathbf{w}) = S(\phi(\mathbf{x})^\top \Phi \mathbf{v}).$$

现在的预测值也只涉及变换后的新点  $\phi(\mathbf{x})$  和变换后的数据  $\phi(\mathbf{x}_i)$  的标量积。

这意味着所有的计算取决于我们是否能快速求出任意一对变换后的点  $\phi(\mathbf{x}), \phi(\mathbf{x}')$  的标量积。考虑  $n = 2$  时二维问题的例子，我们寻找形如

$$\hat{y} = w_1 + \sqrt{2}w_2x_1 + \sqrt{2}w_3x_2 + w_4x_1^2 + \sqrt{2}w_5x_1x_2 + w_6x_2^2.$$

的二次分类规则。这个规则对应于非线性映射  $\mathbf{x} \rightarrow \mathbf{w}^\top \phi(\mathbf{x})$ ，其中

$$\phi(\mathbf{x}) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1x_2, x_2^2).$$

我们发现对任意一对点  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^n$ ，

$$\phi(\mathbf{x})^\top \phi(\mathbf{x}') = (1 + \mathbf{x}^\top \mathbf{x}')^2.$$

一般来说，对于形式是  $\mathbf{x}$  的  $d$  阶多项式的分类规则，我们可以作映射  $\phi(\mathbf{x}) \in \mathbb{R}^p$ ，其中<sup>10</sup>  $p = O(n^d)$ ，使得任意一对点  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^n$  有

$$\phi(\mathbf{x})^\top \phi(\mathbf{x}') = (1 + \mathbf{x}^\top \mathbf{x}')^d.$$

计算上述标量积先验上(a priori)需要  $O(p) = O(n^d)$ ，这是指数(关于特征数目  $n$ )时间的。上述公式可以将计算复杂度降到  $O(n)$ ，这样的计算复杂度是非常低的。

<sup>10</sup>这里  $O(q)$  表示是  $q$  的某线性函数。

表 13.2: 常用的有效核函数

名称	$K(\mathbf{x}, \mathbf{x}')$	参数
线性	$\mathbf{x}^\top \mathbf{x}'$	无
多项式	$(1 + \mathbf{x}^\top \mathbf{x}')^d$	阶数 $d$
高斯	$e^{\frac{1}{2\sigma^2} \ \mathbf{x} - \mathbf{x}'\ _2^2}$	宽度 $\sigma$

为了形成优化问题(13.16)以及给出预测值, 我们需要作  $O(m^2)$  个标量积。因为求解QP(13.16) 的复杂度为  $O(m^3)$ , kernel方法的总复杂度为  $O(m^3 + nm^2)$ , 这与多项式的阶数  $d$  无关。概括起来, 即我们求解多项式模型的且具有欧氏范数惩罚的学习问题所需的计算量与解决同一问题的线性模型所需的计算量一样多。

Kernel的思想也可用于非多项式模型。相对于多项式模型, 这里不必显式地定义非线性映射, 即可以绕过非线性映射  $\phi$  的显式定义, 只需确定核函数

$$K(\mathbf{x}, \mathbf{x}') \doteq \phi(\mathbf{x})^\top \phi(\mathbf{x}')$$

的值即可。当然, 并不是每个函数  $(\mathbf{x}, \mathbf{x}') \rightarrow K(\mathbf{x}, \mathbf{x}')$  都能定义一个有效的核。对形如  $\mathbf{x} \rightarrow \phi(\mathbf{x})$  的映射, 需要满足半正定条件。在实践中, 对核函数的取值附加一个条件将会使计算效率变得很低。表13.2 中给出了通常选取的有效核。

**例13.5. (乳腺癌的分类)** 回到关于乳腺癌的例13.2(基于威斯康星州的数据集)。为了使数据的软间隔分离的图像可视化, 我们先将数据降维, 即利用PCA方法(参见5.3.2 节)把数据从30维降到2 维。需要注意的是, 在实际中这不一定是一个好的想法, 因为将原空间中可分离的数据集投影到低维空间后, 所得投影数据不一定是可分离的。然而, 在这里我们让数据的维数急剧减少只是为了展示数据的二维图像。

对于投影后的二维数据点, 我们求解  $\lambda = 1$  的问题(13.16)以便训练出一个高斯核( $\sigma = 1$ )的SVM, 结果如图13.12所示, 其中  $+$  代表良性数据点,  $*$  代表恶性数据点, 圆圈标记的是支持向量, 实线表示分离面。

### §13.5 无监督学习

在无监督学习中, 仅有数据点  $\mathbf{x}_i \in \mathbb{R}^n, i = 1, \dots, m$ , 没有指定的标签

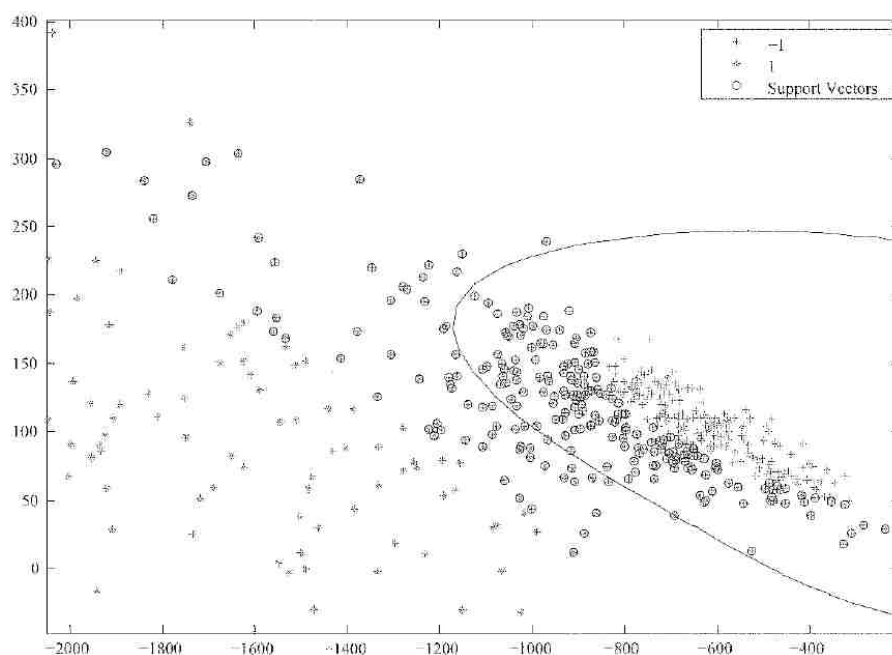


图 13.12: 乳腺癌数据集的二维投影数据的SVM分类

或反馈。我们的任务是从数据中得到一些信息或者结构。

### §13.5.1 主成分分析 (Principal component analysis, PCA)

在主成分分析中, 我们的目标是寻找数据集中最重要的(或就含有信息而言)方向, 即沿着该方向数据的变化最大。已经在5.3.2节中详细描述了PCA。数值上, PCA 问题的计算归结为计算数据矩阵的SVD。下面我们将描述基本PCA问题的一些变形。

### §13.5.2 稀疏PCA

PCA方法在许多领域中得到广泛应用, 且已发展出了多种变形。最近发展了一种称为稀疏PCA的变形, 其提升了所产生的主方向的可解释性, 具有更大的使用价值。在问题(5.18)中添加一个约束用以限制决策变量中非零元素的个数, 可得到稀疏PCA, 即

$$\begin{aligned} & \underset{\mathbf{z} \in \mathbb{R}^n}{\text{maximize}} && \mathbf{z}^\top (\tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top) \mathbf{z} \\ & \text{subject to} && \|\mathbf{z}\|_2 = 1, \text{card}(\mathbf{z}) \leq k, \end{aligned} \quad (13.16)$$

其中  $k \leq n$  是由用户指定的基数的上界。如同在注记5.3中讨论的那样, 原始PCA问题是寻找方向  $\mathbf{z}$ , 其能“解释”数据中的最大变化量; 在实践中, 这个方向  $\mathbf{z}$  通常不是稀疏的。相比之下, 由稀疏PCA 给出的方向如定义所述是稀疏的, 比如  $k \ll n$ , 而这种情况正是我们在实践中所感兴趣的。如果数据的每一维(数据矩阵的行)对应于一个可解释的量, 如价格、温度、体积等, 那么方向越稀疏, 表明它只涉及很少的(即  $k$  个)基本量, 从而它的可解释性就越强。

当数据的维数较小时, 稀疏PCA问题是易于求解的: 只需考虑将数据矩阵  $\tilde{\mathbf{X}}$  移走  $n - k$  维(行), 然后对剩下的行所组成的子矩阵求解一般的PCA 问题。对于维数较大的情况, 会出现行向量的组合爆炸问题, 即生成非常多的可能的子矩阵, 从而使得这种穷举法失效。在12.5.4节中描述了一个有效的解决稀疏PCA的方法, 13.7节的关于文本分析的问题使用了这种求解方法。

**例13.6. (市场数据的稀疏PCA)** 回到例5.3, 我们发现由矩阵  $\mathbf{U}$  的第一列

$$\mathbf{u}_1 = [-0.4143 \quad -0.4671 \quad -0.4075 \quad -0.5199 \quad -0.0019 \quad -0.4169]^\top$$

给出的方向对应的方差最大, 但是  $\mathbf{u}_1$  是一个稠密向量, 其中只有一个分量(第五个)接近于零. 与  $\mathbf{u}_1$  对应的最大奇异值  $\sigma_1 = 1.0765$ , 方差比

$$\eta_1 = \frac{\sigma_1^2}{\sigma_1^2 + \dots + \sigma_6^2} = 67.766\%,$$

这意味着大约 68% 的方差包含在一个由  $\mathbf{u}_1$  作为权重的特定指标中, 该指标几乎是所有指标的非平凡组合。

现在求解  $k = 2$  的稀疏PCA问题(13.16), 由此可得到一个基数小的方向, 沿着该方向方差也较大。我们可以通过穷尽原始六维空间的所有二维组合来搜索, 即我们要解决原始矩阵  $\tilde{\mathbf{X}}$  任意保留两行并去除其余行后所得到的所有可能的矩阵  $\tilde{\mathbf{X}}_2$  的PCA问题。我们发现与行指标 4 和 5 (称为指标PAC和BOT)对应的PCA的可解释方差在所有可能的情况中是最大的。这时对应的数据矩阵的最大奇异值  $\tilde{\sigma}_1 = 0.6826$ , 这与原始数据矩阵的最大奇异值  $\sigma_1 = 1.0765$  是可比的。依据解释方差, 这两个指标相对于总方差的比值为

$$\frac{\tilde{\sigma}_1^2}{\sigma_1^2 + \dots + \sigma_6^2} = 0.2725\%.$$

因此, 仅这两个指标就能捕获大约 28% 的总市场方差。

表 13.3: 与NYT数据的前5个稀疏主成分的非零元素对应的词语

第一PC	第二PC	第三PC	第四PC	第五PC
million	point	official	president	school
percent	play	government	Campaign	program
business	team	United_states	bush	children
company	season	u_s	administration	student
market	game	attack		
companies				

**例13.7. (主题发现中的稀疏PCA)** 当要从大量文档中发现主题时，可以应用稀疏PCA。UCI机器学习数据库中的**纽约时报**(The New York Times, NYT)文本集包含了300,000篇文章，和一个包含了102,660个独一无二的单词的字典，这使得文件大小达到1GB。这个给定的数据没有分类标签，且由于版权原因，文件名和文档级别元数据(如文章段落信息)均未知。文本通过“词袋”(bag-of-words)法转换为数值形式。

表13.3列出了排名前五的稀疏主成分的非零系数对应的非零特征(词语)。每个成分代表了一个确定的话题。的确，**纽约时报**的第一主成分是关于商业的，第二个是关于体育的，第三个是关于美国的，第四个是关于政治的，第五个是关于教育的。即使未提供元数据，稀疏PCA 仍然能清楚地识别，并完美地与**纽约时报**在其网站上对文章进行分类时所采用的主题相对应。

稀疏PCA的思想也可以推广到具有向量基数约束的矩阵的低秩逼近问题。例如，带基数约束的秩一逼近问题形如

$$\begin{aligned} & \underset{\mathbf{p} \in \mathbb{R}^n, \mathbf{q} \in \mathbb{R}^m}{\text{minimize}} && \|\mathbf{X} - \mathbf{p}\mathbf{q}^\top\|_F \\ & \text{subject to} && \text{card}(\mathbf{p}) \leq k, \text{card}(\mathbf{q}) \leq h, \end{aligned}$$

其中  $\mathbf{X} \in \mathbb{R}^{n \times m}$ ，且  $k \leq n, h \leq m$  是用于限制向量  $\mathbf{p}, \mathbf{q}$  中非零元数的数目的参数，由用户选定。当  $n, m$  较小时，类似于上述稀疏PCA问题那样，可用SVD解决这里的问题；否则，求解该问题将变得非常困难，但已经给出了基于幂迭代的有效算法，详见12.5.4节。



## §13.5.3 非负矩阵分解

给定非负的  $n \times m$  阶数据矩阵  $\mathbf{X}$ ，非负矩阵分解(Non-Negative Matrix Factorization, NNMF)指的是用非负低秩成分矩阵  $\mathbf{P}, \mathbf{Q}$  的乘积  $\mathbf{P}\mathbf{Q}^\top$  来逼近数据矩阵  $\mathbf{X}$ ，这里  $\mathbf{P} \in \mathbb{R}^{n \times k}, \mathbf{Q} \in \mathbb{R}^{m \times k}$  表示低秩非负矩阵，即它们的列数  $k \ll \min(n, m)$ 。因此，这是在5.3.2.4节中讨论过的低秩矩阵分解问题的简单改进。

为说明NNMF，我们先看数据矩阵  $\mathbf{X}$  的普通秩一逼近，即

$$\min_{\mathbf{p} \in \mathbb{R}^n, \mathbf{q} \in \mathbb{R}^m} \|\mathbf{X} - \mathbf{p}\mathbf{q}^\top\|_F.$$

如果存在  $\mathbf{p}, \mathbf{q}$  使得问题的目标值很小，即有  $\mathbf{X} \approx \mathbf{p}\mathbf{q}^\top$ ，我们就可以将  $\mathbf{p}$  解释为典型数据点， $\mathbf{q}$  解释为典型特征。如果  $\mathbf{X}$  非负，就可能存在一个秩一逼近  $\mathbf{X} \approx \mathbf{p}\mathbf{q}^\top$ ，但是这里  $\mathbf{p}, \mathbf{q}$  的某些分量是负的。对于这种情况，很难解释结果，比如当  $\mathbf{p}$  存在负分量时，我们就不能将  $\mathbf{p}$  解释成典型数据点。

NNMF方法可应用于非负数据矩阵。如果目标秩  $k = 1$ ，则NNMF问题即

$$\begin{aligned} & \underset{\mathbf{p} \in \mathbb{R}^n, \mathbf{q} \in \mathbb{R}^m}{\text{minimize}} && \|\mathbf{X} - \mathbf{p}\mathbf{q}^\top\|_F \\ & \text{subject to} && \mathbf{p} \geq \mathbf{0}, \mathbf{q} \geq \mathbf{0}. \end{aligned}$$

该结果的解释为：如果  $\mathbf{p}, \mathbf{q}$  为非负向量，且  $\mathbf{X} \approx \mathbf{p}\mathbf{q}^\top$ ，则  $\mathbf{X}$  的每列均与向量  $\mathbf{q}$  成正比例，向量  $\mathbf{p}$  的分量分别对应着这些比例系数。因此每个数据点都沿着单个“轮廓”  $\mathbf{q}$  运动，数据点的不同之处是与每个数据点对应一个特定的非负比例因子。更一般的，NNMF 问题可表述为

$$\begin{aligned} & \underset{\mathbf{P} \in \mathbb{R}^{n \times k}, \mathbf{Q} \in \mathbb{R}^{m \times k}}{\text{minimize}} && \|\tilde{\mathbf{X}} - \mathbf{P}\mathbf{Q}^\top\|_F \\ & \text{subject to} && \mathbf{P} \geq \mathbf{0}, \mathbf{Q} \geq \mathbf{0}. \end{aligned}$$

这里的解释是数据点是  $\mathbf{Q}$  的  $k$  个列向量确定的基本轮廓的线性组合。该问题是非凸的，难以精确求解。经常用到的方法是一种基于坐标块下降<sup>11</sup>的方法，即关于  $\mathbf{P}, \mathbf{Q}$  交替地进行极小化。每个子问题(关于  $\mathbf{P}$  或者  $\mathbf{Q}$  的极小化)都是凸问题，事实上是凸二次规划。

**例13.8.** (用于文本文档主题发现的NNMF) 利用“词袋”方法(见例2.1)，我们可以将包含  $m$  个文本文档的集合用矩阵  $\mathbf{X}$  来表示，其中  $x_{ij}$  表示术语  $i$  在文件  $j$  中出现的次数。对于这种矩阵，自然会想到用NNMF。秩

<sup>11</sup>参见12.5.3 节

一非负逼近  $\mathbf{X} \approx \mathbf{p}\mathbf{q}^\top$ , 其中  $\mathbf{p} \geq \mathbf{0}, \mathbf{q} \geq \mathbf{0}$  表示每个文档  $i$  ( $\mathbf{X}$  的列) 可用  $\mathbf{p}_i\mathbf{q}$  近似, 即秩一逼近下所有文档均相似于一个文档。为了将  $\mathbf{p}_i\mathbf{q}$  解释为一个文档, 这里  $\mathbf{p}, \mathbf{q}$  的非负性是至关重要的。秩  $k$  阶逼近表示所有文档是  $k$  个基本文档的“混合”。

**例13.9. (学生得分矩阵的NNMF)** 考虑关于四门考试和五位学生的二进制得分矩阵

$$\mathbf{X} = \begin{bmatrix} 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 & 1 \\ 1 & 1 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 & 1 \end{bmatrix},$$

其中当学生  $j$  通过考试  $i$  时,  $x_{ij}$  为1, 否则为 0。我们可以将其写作  $\mathbf{X} = \mathbf{P}\mathbf{Q}^\top$ , 其中  $\mathbf{P}, \mathbf{Q}$  均只有三列

$$\mathbf{P} = \begin{bmatrix} 0 & 1/2 & 0 \\ 1/2 & 0 & 0 \\ 0 & 1/2 & 1 \\ 1/2 & 0 & 0 \end{bmatrix}, \quad \mathbf{Q} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 2 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \\ 0 & 2 & 2 \end{bmatrix}.$$

我们可将  $\mathbf{P}, \mathbf{Q}$  的三列解释为“成功通过考试所需要的三项‘技能’”。矩阵  $\mathbf{P}$  代表着通过每门考试需要何种技能; 矩阵  $\mathbf{Q}$  表明每个学生各掌握了什么技能。以第一门考试和第一个学生为例, 我们从  $\mathbf{P}$  中可以看到通过考试 1 需要技能 2, 但从  $\mathbf{Q}$  中看到学生 1 仅掌握技能 1, 故学生 1 不能通过考试 1。事实上, 学生 1 只能通过考试 3, 因为其余考试都要求技能 2 或 3。

#### §13.5.4 鲁棒PCA

低秩矩阵逼近和与之密切相关的PCA均可解释为将给定的数据矩阵  $\mathbf{X}$  表示为矩阵之和  $\mathbf{X} = \mathbf{Y} + \mathbf{Z}$ , 其中  $\mathbf{Y}$  为低秩矩阵,  $\mathbf{Z}$  较小。

鲁棒PCA是PCA的另一种变形<sup>12</sup>, 我们将借此尝试将给定的  $n \times m$

<sup>12</sup>参见E.J. Candès, X. Li, Y. Ma, and J. Wright, Robust principal component analysis?, J. ACM, 2011.

矩阵  $\mathbf{X}$  分解为低秩成分和稀疏成分之和, 即形如

$$\begin{aligned} & \underset{\mathbf{Y}, \mathbf{Z} \in \mathbb{R}^{n \times m}}{\text{minimize}} && \text{rank}(\mathbf{Y}) + \lambda \text{card}(\mathbf{Z}) \\ & \text{subject to} && \mathbf{X} = \mathbf{Y} + \mathbf{Z}, \end{aligned}$$

其中  $\text{card}(\mathbf{Z})$  是  $n \times m$  矩阵  $\mathbf{Z}$  中非零元的个数,  $\lambda$  是一个参数, 其允许我们在  $\mathbf{X}$  的成分  $\mathbf{Y}$  的低秩性和另一成分  $\mathbf{Z}$  的稀疏性之间进行折衷。当认为测量数据具有低秩结构, 但是存在着由稀疏成分表现的“spike”(或outlier)时, 会出现这种问题。

鲁棒PCA问题通常很难求解。在11.4.1.4节中提到, 我们可以用矩阵的核范数作为秩的凸近似, 并用  $\ell_1$ -范数控制基数, 由此得到上述鲁棒PCA问题的一个凸逼近, 即

$$\begin{aligned} & \underset{\mathbf{Y}, \mathbf{Z} \in \mathbb{R}^{n \times m}}{\text{minimize}} && \|\mathbf{Y}\|_* + \lambda \|\mathbf{Z}\|_1 \\ & \text{subject to} && \mathbf{X} = \mathbf{Y} + \mathbf{Z}, \end{aligned}$$

这里  $\|\mathbf{Z}\|_1$  是矩阵  $\mathbf{Z}$  中元素的绝对值之和,  $\|\mathbf{Y}\|_*$  是矩阵  $\mathbf{Y}$  的核范数(即  $\mathbf{Y}$  的奇异值之和)。上述问题是一个SDP, 若问题规模适当, 则是可解的; 更大规模的实例需要用专门的算法来求解。

---

**例13.10. (视频中的背景建模)** 对视频数据而言, 由于图像间的相关性, 自然地可用低秩模型来对其进行建模。视频监控中最基本的算法功能之一是为场景中的背景变化建立好的模型。这个任务因为前景中出现的物体而变得复杂: 比如在繁忙的背景中, 每个图像中都可能出现一些异常现象。背景模型必须有足够强的适应性, 以便适应背景中的变化, 比如变化的光线。在这种情况下, 很自然地就可用低秩近似来建模背景变化。前景中的物体(如汽车和行人)通常只占据图像像素的一小部分, 因此可处理成稀疏误差。

我们考虑原始视频中的五幅图像, 如图13.13第一行所示, 其中包含一个过路人。每幅图的分辨率为  $176 \times 144$  像素。我们首先将其分成三个频段(RGB)。对于每个频段, 我们将每张图像排放成矩阵  $\mathbb{R}^{25344 \times 5}$  中的一列, 然后通过鲁棒PCA法将  $\mathbf{X}$  分解为一个低秩成分  $\mathbf{Y}$  和稀疏成分  $\mathbf{Z}$ 。之后我们将三个频段重新结合到一起, 形成两个图像集, 其中一个为低秩成分, 另一个为稀疏成分。如图13.13所示, 低秩成分正确恢复了背景, 同时稀疏成分正确地显示了正在移动的人。

---

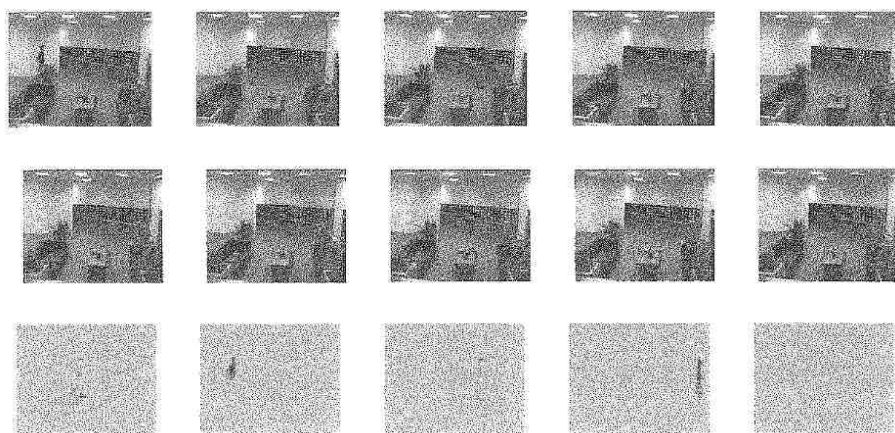


图 13.13: 视频帧集合: 顶部为原始集合, 中部为只显示背景的低秩成分, 底部为只显示前景的稀疏成分

### §13.5.5 稀疏图的高斯模型

稀疏图的高斯建模的目的是发现数据集中隐藏的图结构。基本的假设是数据由多维高斯分布生成, 而我们的任务是由数据得到该分布的参数, 这些参数能让模型显示出大量的条件独立性。为了让这些术语更清楚, 我们将给出一些背景知识。

**用高斯分布拟合数据.** 设取值为 $\mathbb{R}^n$ 中的向量的随机向量 $\boldsymbol{\xi}$ 服从高斯分布, 即它的概率密度函数形如

$$p(\boldsymbol{x}) = \frac{1}{\sqrt{(2\pi)^n \det(\boldsymbol{S})}} e^{-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^\top \boldsymbol{S}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})},$$

其中参数 $\boldsymbol{\mu} \in \mathbb{R}^n$ ,  $\boldsymbol{S}$  是一个  $n \times n$  的对称正定矩阵。这些参数都有自然的解释:  $\boldsymbol{\mu}$  是分布的平均值(期望),  $\boldsymbol{S}$  是协方差矩阵。事实上进行积分, 有

$$\mathbb{E}(\boldsymbol{\xi}) = \int \boldsymbol{x} p(\boldsymbol{x}) d\boldsymbol{x} = \boldsymbol{\mu},$$

$$\mathbb{E}\left((\boldsymbol{\xi} - \boldsymbol{\mu})(\boldsymbol{\xi} - \boldsymbol{\mu})^\top\right) = \int (\boldsymbol{x} - \boldsymbol{\mu})(\boldsymbol{x} - \boldsymbol{\mu})^\top p(\boldsymbol{x}) d\boldsymbol{x} = \boldsymbol{S}.$$

给定数据集  $\boldsymbol{X} = [\boldsymbol{x}_1, \dots, \boldsymbol{x}_m]$ , 定义**样本均值**(sample mean)

$$\bar{\boldsymbol{\mu}} \doteq \frac{1}{m} \sum_{i=1}^m \boldsymbol{x}_i$$

和样本协方差阵(sample covariance matrix)

$$\bar{\mathbf{S}} \doteq \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \bar{\boldsymbol{\mu}})(\mathbf{x}_i - \bar{\boldsymbol{\mu}})^\top.$$

下面在给定数据集  $\mathbf{X}$  后, 我们考虑用极大似然的思想来拟合高斯分布, 即寻找参数  $\boldsymbol{\mu}, \mathbf{S}$  的估计值  $\hat{\boldsymbol{\mu}}, \hat{\mathbf{S}}$ , 这些值使得似然函数取极大值. 为此, 需要极大化密度  $p(\mathbf{x}_i), i = 1, \dots, m$ , 的乘积, 对这个乘积取对数, 得问题

$$\begin{aligned} \underset{\boldsymbol{\mu}, \mathbf{S}}{\text{maximize}} \quad & -\frac{1}{2} \sum_{i=1}^m (\mathbf{x}_i - \boldsymbol{\mu})^\top \mathbf{S}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) + \log \det(\mathbf{S}) \\ \text{subject to} \quad & \mathbf{S} \succeq \mathbf{0}. \end{aligned}$$

记该问题的解为  $\boldsymbol{\mu}^*$  和  $\mathbf{S}^*$ , 易验证  $\boldsymbol{\mu}^* = \bar{\boldsymbol{\mu}}$ , 即样本均值  $\bar{\boldsymbol{\mu}}$  是总体均值  $\boldsymbol{\mu}$  的极大似然估计. 此外, 还有

$$\mathbf{S}^* = \arg \max_{\mathbf{S} \succ \mathbf{0}} \text{trace}(\mathbf{S}^{-1} \bar{\mathbf{S}}) - \log \det(\mathbf{S}), \quad (13.17)$$

其中  $\bar{\mathbf{S}}$  是样本协方差阵. 问题(13.17)本身的表述方式不是凸优化, 但通过变量替换  $\mathbf{S} \rightarrow \mathbf{P} = \mathbf{S}^{-1}$ , 可将其表述成关于  $\mathbf{P}$  的问题, 即

$$\max_{\mathbf{P} \succ \mathbf{0}} \text{trace}(\mathbf{P} \bar{\mathbf{S}}) + \log \det(\mathbf{P}). \quad (13.18)$$

该问题关于矩阵  $\mathbf{P}$  是凸的. 假定样本协方差矩阵  $\bar{\mathbf{S}}$  正定, 考虑该问题的最优性条件即可得到该问题的解, 即得到  $\hat{\mathbf{P}} = \bar{\mathbf{S}}^{-1}$ . 即我们得到总体协方差阵的极大似然估计就是样本协方差阵, 即  $\hat{\mathbf{S}} = \bar{\mathbf{S}}$ . 综上所述, 在该情况下, 极大似然方法得到的总体期望的估计是样本平均值, 总体协方差阵的估计是样本协方差阵.

**条件独立(conditional independence).** 假设协方差阵  $\mathbf{S}$  是可逆的. 称逆矩阵  $\mathbf{P} = \mathbf{S}^{-1}$  是精度(precision)矩阵. 就条件独立(conditional independence)而言, 精度矩阵的元素有一个很好的解释. 如果固定除  $\xi_k$  和  $\xi_l$  之外的其余随机变量时, 随机变量  $\xi_k$  和  $\xi_l$  独立, 我们称这对随机变量  $(\xi_k, \xi_l)$  条件独立. 这意味着对应的(条件)概率密度可被分解为两个标量分布的乘积, 这两个标量分布分别以  $x_k$  和  $x_l$  为变量. 现在假设存在随机变量对  $(\xi_k, \xi_l)$  使得精度矩阵中对应的  $P_{kl} = 0$ . 上面给定的密度函数  $p$  满足

$$-2 \ln p(\mathbf{x}) = (\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{S}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = \sum_{i,j} P_{ij} (x_i - \mu_i)(x_j - \mu_j).$$

显然, 当  $P_{kl} = 0$ , 上述二次函数中没有一项会同时包含  $x_k$  和  $x_l$ . 因此, 当固定  $\mathbf{x}$  中除  $x_k$  和  $x_l$  之外的其余变量后, 我们可将上式写作

$$-2 \ln p(\mathbf{x}) = p_k(x_k) + p_l(x_l).$$

意即概率密度函数  $p$  可分解成两个函数的乘积, 其中每个因子仅由  $x_k$  和  $x_l$  分别确定。

**例13.11.** 假设  $\mathbb{R}^3$  中的高斯分布的均值为  $\mathbf{0}$ , 它的  $3 \times 3$  精度矩阵  $\mathbf{P}$  中有  $P_{12} = 0$ 。让我们验证  $x_1, x_2$  是条件独立的。我们固定变量  $x_3$ , 观察到

$$\begin{aligned} -2 \ln p(\mathbf{x}) &= P_{11}x_1^2 + P_{22}x_2^2 + 2P_{12}x_1x_2 + 2P_{13}x_1x_3 + 2P_{23}x_2x_3 + P_{33}x_3^2 \\ &= (P_{11}x_1^2 + 2P_{13}x_1x_3) + (P_{22}x_2^2 + 2P_{23}x_2x_3) + \text{constant}. \end{aligned}$$

如上所述, 当  $x_3$  固定时, 密度函数可写为两项乘积, 每项分别只与  $x_1, x_2$  有关。

在仅基于观测数据来试图理解随机变量的关系时, 条件独立是一个自然的概念。在实际操作中, 众所周知的协方差的概念往往会失去作用。对于现实世界中的大部分数据集而言, 所有的变量都是相关的, 因此协方差阵是稠密的, 这样很难从协方差阵看出有什么明显的结构。与此相反, 许多变量对是条件独立的, 这表现为对应的精度矩阵是稀疏的, 该事实与其逆矩阵(协方差阵)稠密刚好相反。

上述事实促使我们去发现条件独立的结构, 例如不同随机变量  $\xi_1, \dots, \xi_n$  间的图, 在该图中任意一对条件独立的变量对  $\xi_i, \xi_j$  间没有边相连接。图越稀疏, 就越易读。

下面的例子可以说明这种思想, 即在揭示协方差矩阵无法呈现的结构时, 精度矩阵是如何起作用的。

**例13.12. (条件独立的一个例子)** 假设我们在大量人群中观察三个随机变量  $\xi_1, \xi_2, \xi_3$ , 分别对应鞋码, 是否白发, 以及年龄。假设我们经计算得到了协方差阵的一个很好的估计。不出所料, 这个协方差矩阵是稠密的, 这意味着所有变量都是相关的。特别的, 当年龄增长, 鞋码也增加, 白发出现的概率也会增加。

然而精度矩阵满足  $P_{12} = 0$ , 这表明当年龄给定时, 鞋码和是否有白发实际上是相互独立的。这说得通: 当我们将相似年纪的人分为一个群体时, 我们会发现对某一个年龄群体, 鞋码和是否有白发间并没有统计联系。

**稀疏的精度矩阵.** 为了在极大似然估计的模型中增大条件独立性, 我们修正凸问题(13.18), 即给其的目标函数增加一项关于  $\mathbf{P}$  的  $\ell_1$ -范数的惩罚, 得到变形问题

$$\max_{\mathbf{P} \succ \mathbf{0}} \text{trace}(\mathbf{P}\tilde{\mathbf{S}}) + \log \det(\mathbf{P}) - \lambda \|\mathbf{P}\|_1, \quad (13.19)$$

其中  $\lambda$  是参数, 它允许我们对模型的“拟合度”和精度矩阵的稀疏性进行折衷。与原问题(13.18)不同的是, 上述惩罚问题没有解析解。但是它是凸的, 存在算法能有效地求解它。

**例13.13. (利率数据图)** 在图13.14中, 我们呈现了从协方差矩阵的逆中得到的利率(采样超过一年多)的相关结构。每个节点表示一个特定的利率期限(interest rate maturity), 如果协方差矩阵的逆中对应的系数非零, 两个节点间有边相连, 即它们条件相关。我们将  $\lambda = 0$  (此时就是样本协方差矩阵的逆)和  $\lambda = 1$  时问题(13.19)的解相比较。在稀疏解中利率清晰地呈现出按期限聚类的现象。

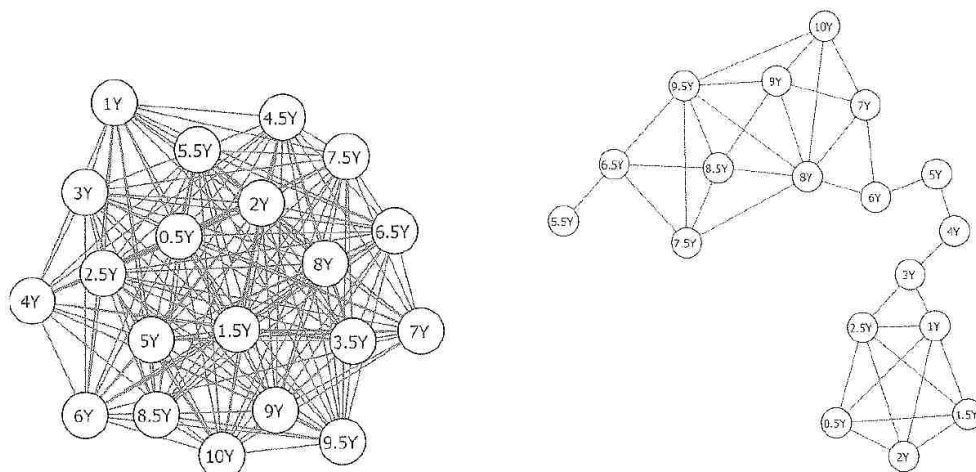


图 13.14: 利率数据的稀疏图模型

**例13.14. (议员的图)**(a graph of Senators) 回到例2.15中的参议院投票数据。图13.15显示了由数据得到的稀疏图, 其对应于  $\lambda = 0.1$  时问题(13.19)的

解。每当最优的稀疏协方差阵的逆  $\mathbf{P}$  中对应的元素为零，图13.15中与此元素对应的两个节点(参议员)之间就没有边相连。

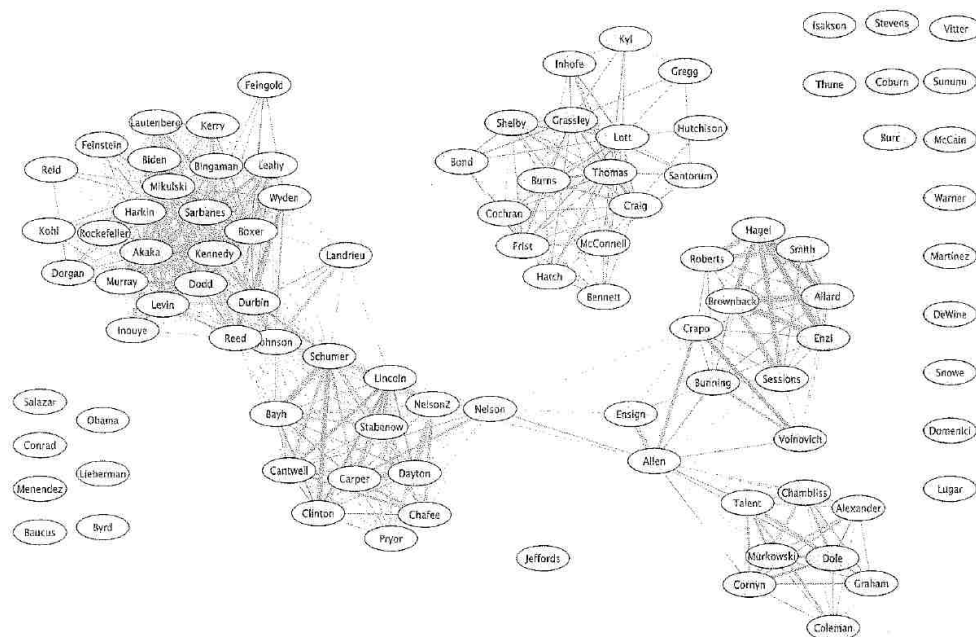


图 13.15: 参议院投票数据的稀疏图模型

### §13.6 习题

**习题13.1. (文本分析中的SVD)** 考虑对应于大量新闻文章所构成集合的数据集，即假设给定  $n \times m$  的术语文档(term-by-documents)矩阵  $\mathbf{X}$ 。确切地说， $\mathbf{X}$  的元素  $(i, j)$  表示单词  $i$  在文档  $j$  中出现的次数。我们想用二维图像对数据集进行可视化。请详述你将如何完成以下事情(用  $\mathbf{X}$  的适当中心化版本的SVD详述你的步骤)。

- 在单词空间内，将不同的新闻源作为点，画出方差最大的点。
- 在新闻源空间内，将不同的单词作为点，画出方差最大的点。

**习题13.2. (学习一个因素的模型)** 给定一个数据矩阵  $\mathbf{X} = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}]$ ，其中  $\mathbf{x}^{(i)} \in \mathbb{R}^n, i = 1, \dots, m$ 。假设已经将数据进行中心化处理，即  $\mathbf{x}^{(1)} +$



$\dots + \mathbf{x}^{(m)} = \mathbf{0}$  . 一种对协方差矩阵的(经验)估计<sup>13</sup>是

$$\mathbf{\Sigma} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)} \mathbf{x}^{(i)\top}.$$

在实践中, 我们通常会发现上述对协方差矩阵的估计有噪声。去除噪声的一种办法是将协方差矩阵近似为  $\mathbf{\Sigma} \approx \lambda \mathbf{I} + \mathbf{F} \mathbf{F}^\top$ , 其中  $\mathbf{F}$  是一个包含所谓“因素载荷”的  $n \times k$  矩阵,  $k \ll n$  是因素的个数,  $\lambda \geq 0$  为“异质噪声”的方差。对应于这种近似方法的随机模型为

$$\mathbf{x} = \mathbf{F} \mathbf{f} + \sigma \mathbf{e},$$

其中  $\mathbf{x}$  是中心化的观测(随机)向量;  $\mathbf{f}$  和  $\mathbf{e}$  均是随机向量, 二者的均值均为零, 协方差矩阵均是单位矩阵;  $\sigma = \sqrt{\lambda}$  为异质噪声成分  $\sigma \mathbf{e}$  的标准差。该随机模型可解释为观测是  $k$  个因素的线性组合(  $k$  较小)和噪声部分, 这个噪声对每个维度独立地产生影响。

为了用  $\mathbf{F}, \lambda$  拟合数据, 我们需要求解

$$\min_{\mathbf{F}, \lambda \geq 0} \|\mathbf{\Sigma} - \lambda \mathbf{I} - \mathbf{F} \mathbf{F}^\top\|_{\text{F}}. \quad (13.20)$$

- (a) 假设  $\lambda$  已知且小于  $\lambda_k$  (经验协方差矩阵  $\sigma$  中第  $k$  大的特征值)。将最优解  $\mathbf{F}$  表示为  $\lambda$  的函数, 记作  $\mathbf{F}(\lambda)$ 。意即: 你需要在固定  $\lambda$  的情况下解出  $\mathbf{F}$ 。
- (b) 定义误差  $E(\lambda) = \|\mathbf{\Sigma} - \lambda \mathbf{I} - \mathbf{F}(\lambda) \mathbf{F}(\lambda)^\top\|_{\text{F}}$ , 其中  $\mathbf{F}(\lambda)$  是在(a)中得到的矩阵。证明

$$E(\lambda)^2 = \sum_{i=k+1}^p (\lambda_i - \lambda)^2.$$

确定极小化误差的最优解  $\lambda$  的显式表示, 然后总结处理估计问题(13.20)的解决方案。

- (c) 假设我们想估计数据空间中关于特定方向的风险(由方差度量)。回忆例4.2, 当给定单位范数的  $n$ - 维向量  $\mathbf{w}$  时, 沿方向  $\mathbf{w}$  的方差为  $\mathbf{w}^\top \mathbf{\Sigma} \mathbf{w}$ 。证明与使用  $\mathbf{\Sigma}$  相比较, 使用对  $\mathbf{\Sigma}$  的  $k$  阶近似是定向风险的一个偏低估计。基于上述因素模型的近似将会怎样? 请进行讨论。

**习题13.3. (时间序列的运动预测)** 我们有一个历史数据集, 其包含时间序列  $r(1), \dots, r(T)$  的值。我们的目标是预测时间序列将会上升还是下降。

<sup>13</sup>参见例4.2.

基本想法是基于使用  $n$  个过往数据值(这里的  $n$  是固定的)的自回归模型输出的符号进行预测, 即对  $t$  时刻  $r(t+1) - r(t)$  的数值的符号预测为

$$\hat{y}_{\mathbf{w},b}(t) = \text{sgn}(w_1 r(t) + \cdots + w_n r(t-n+1) + b),$$

其中  $\mathbf{w} \in \mathbb{R}^n$  是我们的分类器的系数,  $b$  是偏移项,  $n \ll T$  决定了我们将用多久远的数据来做预测。

(a) 第一步, 我们将求解

$$\min_{\mathbf{w},b} \sum_{t=n}^{T-1} (\hat{y}_{\mathbf{w},b}(t) - y(t))^2.$$

其中  $y(t) = \text{sgn}(r(t+1) - r(t))$ 。换句话说, 我们尝试在最小二乘意义下将训练集的分类器作出的预测与观测到的事实进行匹配。上述问题是凸优化吗? 如果不是, 为什么?

(b) 如果要使用凸优化, 请说明如何表述问题及训练分类器。确保精确地定义了学习的步骤、产生的优化问题中的变量, 并说明如何求解优化问题及如何利用最优解来作出预测。

**习题13.4. (PCA的一种变形)** 回到例11.2中PCA的变形。使用你选定的(也可以是合成的)数据集来去比较经典PCA和这里所研究的变形, 尤其研究二者对于异常数据的敏感度。确保建立尽可能严谨的评价体系。讨论你的结果。

**习题13.5. (平方惩罚与非平方惩罚)** 考虑问题

$$\begin{aligned} P(\lambda) : p(\lambda) &\doteq \min_{\mathbf{x}} f(\mathbf{x}) + \lambda \|\mathbf{x}\|, \\ Q(\mu) : q(\mu) &\doteq \min_{\mathbf{x}} f(\mathbf{x}) + \frac{1}{2} \mu \|\mathbf{x}\|^2, \end{aligned}$$

其中  $f$  为凸函数,  $\|\cdot\|$  为任一向量范数, 参数  $\lambda > 0, \mu > 0$ 。假设对于给定的参数, 相应的问题均有唯一解。

通常, 对固定的  $\lambda$  和  $\mu$ , 上述两个问题的解是不重合的。这个练习将说明: 我们可以遍历第一个问题的所有解, 得到第二个问题的解集, 反之亦然。

(a) 证明:  $p, q$  均为凹函数, 且定义为  $\tilde{q}(\mu) = q(1/\mu)$  的函数在  $\mathbb{R}_+$  上是凸的。

(b) 证明

$$p(\lambda) = \min_{\mu > 0} q(\mu) + \frac{\lambda^2}{2\mu}, \quad q(\mu) = \max_{\lambda > 0} p(\lambda) - \frac{\lambda^2}{2\mu}.$$

对第二个式子, 你可假定  $f$  的定义域的内部非空。

(c) 从第一部分推出解的路径重合, 即对于每个  $\lambda > 0$ , 如果我们已经求解了第一个问题, 则对任一  $\mu > 0$ , 我们如上得到的最优点都将是第二个问题的最优点; 反之亦然。不妨将  $P(\lambda)(Q(\mu))$  的(惟一)解记作  $\mathbf{x}^*(\lambda)(\mathbf{z}^*(\mu))$ 。

(d) 关于第三个函数

$$r(\kappa) : r(\kappa) = \min\{f(\mathbf{x}) : \|\mathbf{x}\| \leq \kappa\}$$

陈述并证明相似的结论。

(e) 如果去掉惟一解的假设, 结果会怎样?

**习题13.6. (基数惩罚最小二乘)** 考虑问题

$$\phi(\lambda) \doteq \min_{\mathbf{w}} \|\mathbf{X}^\top \mathbf{w} - \mathbf{y}\|_2^2 + \rho^2 \|\mathbf{w}\|_2^2 + \lambda \text{card}(\mathbf{w}),$$

其中  $\mathbf{X} \in \mathbb{R}^{n \times m}$ ,  $\mathbf{y} \in \mathbb{R}^m$ ,  $\rho > 0$  为正则化参数,  $\lambda \geq 0$  允许我们控制解向量的基数(非零元素的个数)。解向量的基数越小, 对应的结果的可解释性越强。上述问题通常很难求解。在本练习中, 我们将  $\mathbf{X}$  的第  $i$  行记作  $\mathbf{a}_i^\top$ ,  $i = 1, \dots, n$ , 其对应于一个特定的“特征”(即变量  $\mathbf{w}$  的维数)。

(a) 首先假设没有基数惩罚, 即  $\lambda = 0$ 。证明

$$\phi(0) = \mathbf{y}^\top \left( \mathbf{I} + \frac{1}{\rho^2} \sum_{i=1}^n \mathbf{a}_i \mathbf{a}_i^\top \right)^{-1} \mathbf{y}.$$

(b) 现在考虑  $\lambda > 0$  的情况。证明

$$\phi(\lambda) = \min_{\mathbf{u} \in \{0,1\}^n} \mathbf{y}^\top \left( \mathbf{I} + \frac{1}{\rho^2} \sum_{i=1}^m u_i \mathbf{a}_i \mathbf{a}_i^\top \right)^{-1} + \lambda \sum_{i=1}^n u_i.$$

(c) 将约束  $\mathbf{u} \in \{0,1\}^n$  替换为区间约束  $\mathbf{u} \in [0,1]^n$ , 即可得到问题的一个自然的松弛。证明产生的下界  $\phi(\lambda) \geq \underline{\phi}(\lambda)$  是凸问题

$$\underline{\phi}(\lambda) = \max_{\mathbf{v}} 2\mathbf{y}^\top \mathbf{v} - \mathbf{v}^\top \mathbf{v} - \sum_{i=1}^n \left( \frac{(\mathbf{a}_i^\top \mathbf{v})^2}{\rho^2} - \lambda \right)_+$$

的最优值。如何从上述问题的解  $\mathbf{v}^*$  来恢复次优的稀疏模式?

(d) 将上述问题表示为一个SOCP。

(e) 写出该SOCP的对偶问题,并证明它可化简为

$$\underline{\phi}(\lambda) = \|\mathbf{X}^\top \mathbf{w} - \mathbf{y}\|_2^2 + 2\lambda \sum_{i=1}^n B\left(\frac{\rho x_i}{\sqrt{\lambda}}\right),$$

其中 $B$ 是(凸的)**Hübert函数的反函数**(inverse Hübér function): 对于  $\xi \in \mathbb{R}$ ,

$$B(\xi) \doteq \frac{1}{2} \min_{0 \leq z \leq 1} \left( z + \frac{\xi^2}{z} \right) = \begin{cases} |\xi|, & \text{如果 } |\xi| \leq 1, \\ \frac{\xi^2+1}{2}, & \text{否则} \end{cases}$$

再次说明如何由上述问题的解  $\mathbf{w}^*$  来恢复次优的稀疏模式?

(f) 一个处理基数惩罚的经典方法是将基数替换为  $\ell_1$ -范数,从而得到松弛的凸优化。将上述方法与  $\ell_1$ -范数松弛方法进行比较. 请对此进行讨论。



## 第十四章 计算金融

最优化在计算金融中扮演着越来越重要的角色. 因为H.Markowitz的开创性的工作<sup>1</sup>, 金融领域的从业人员和研究者一直在致力于为策略性投资决策寻找合适的模型的研究. 在本章中, 我们给出了基于凸优化模型的金融应用领域的一个简单的综述. 在这些金融应用案例中, 能效求解问题的重要工具是凸优化模型。

### §14.1 单期投资组合最优化

已经在例2.6、例4.3和例8.19中介绍了一个基本的金融投资组合模型。在那个问题的背景下, 向量 $\mathbf{r} \in \mathbb{R}^n$  包含了 $n$ 种资产在固定周期 $\Delta$  内的随机收益率,  $\mathbf{x} \in \mathbb{R}^n$ 描述了一个投资者给 $n$ 种资产中的每一种所分配的美元价值量。在这一节中, 我们假设投资者具有一个初始的投资组合 $\mathbf{x}(0) \in \mathbb{R}^n$ (所以这个投资者所拥有的初始总资产 $w(0) = \sum_{i=1}^n x_i(0) = \mathbf{1}^\top \mathbf{x}(0)$ ), 他希望进行一些市场交易, 从而更新了他的投资组合。交易量向量记为 $\mathbf{u} \in \mathbb{R}^n$ (如果给第 $i$ 种资产增加投资, 则 $u_i > 0$ , 反之小于0), 从而更新投资组合后, 投资组合变成

$$\mathbf{x} = \mathbf{x}(0) + \mathbf{u}.$$

我们进一步假设投资组合是自我筹资的(self financing), 意思是, 总资产是守恒的(没有新的资金进入或者输出), 这个可以用预算守恒(budget conservation) 条件 $\mathbf{1}^\top \mathbf{x} = \mathbf{1}^\top \mathbf{x}(0)$ 表示, 即

$$\mathbf{1}^\top \mathbf{u} = 0.$$

我们不失一般性地进一步假设初始资产是一个单位, 即( $w(0) = 1$ )。有时, 投资者的某种资产量可以是负的(即允许 $\mathbf{x}$ 的某些元素是负的)。在金融操作中称这种情况为卖空(short-selling), 这种情况在实践中意味着从银行或者经纪人处借来资产, 然后加以变卖。如果经纪人不允许该特征(或者我们不想去使用该操作), 我们给这个问题施加“无卖空”的限制条件 $\mathbf{x} \geq \mathbf{0}$ , 即

$$\mathbf{x}(0) + \mathbf{u} \geq \mathbf{0} \quad \text{无卖空.}$$

---

<sup>1</sup>由于他将金融折衷决策问题表述成一个凸优化问题的贡献, 他于1990年和M.Miller及W.Sharpe获得了诺贝尔经济学奖。

## §14.1.1 均值-方差最优化

在一种以诺贝尔经济学奖获得者H.Markowitz命名的经典方法中，人们假设收益向量的期望 $\hat{\mathbf{r}} = \mathbb{E}\{\mathbf{r}\}$ 和收益的协方差矩阵

$$\Sigma = \mathbb{E}\{(\mathbf{r} - \hat{\mathbf{r}})(\mathbf{r} - \hat{\mathbf{r}})^\top\}$$

是已知的，从而期末投资组合的收益就可以用随机变量

$$\rho(\mathbf{x}) = \mathbf{r}^\top \mathbf{x}$$

来描述，它的期望

$$\mathbb{E}\{\rho(\mathbf{x})\} = \hat{\mathbf{r}}^\top \mathbf{x},$$

方差

$$\text{var}\{\rho(\mathbf{x})\} = \mathbf{x}^\top \Sigma \mathbf{x}.$$

在Markowitz方法中，我们假设投资者偏好那些收益期望高的投资组合，同时他是规避风险(risk)的，这里的风险用投资组合的方差来度量(金融背景中一般叫做波动性)。为了得到最优的投资组合，或者以容许风险的上限为约束来极大化投资组合收益的期望，或者以收益的下限为约束来极小化风险。用 $\mu$ 表示想得到的期望收益的最小值， $\bar{\sigma}^2$ 是能接受的风险的最大值，则可将这两个问题表述为

$$\begin{aligned} \rho(\bar{\sigma}) = \underset{\mathbf{u}}{\text{maximize}} \quad & \hat{\mathbf{r}}^\top \mathbf{x} \\ \text{subject to} \quad & \mathbf{x}^\top \Sigma \mathbf{x} \leq \bar{\sigma}^2 \\ & \mathbf{x} \in \mathcal{X} \end{aligned}$$

和

$$\begin{aligned} \bar{\sigma}^2(\mu) = \underset{\mathbf{u}}{\text{minimize}} \quad & \mathbf{x}^\top \Sigma \mathbf{x} \\ \text{subject to} \quad & \hat{\mathbf{r}}^\top \mathbf{x} \geq \mu \\ & \mathbf{x} \in \mathcal{X}, \end{aligned}$$

其中

$$\mathcal{X} = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x} = \mathbf{x}(0) + \mathbf{u}, \mathbf{1}^\top \mathbf{x} = w(0), \text{ns} \cdot \mathbf{x} \geq \mathbf{0}\}$$

这里的ns是一个标识，当不允许卖空时ns = 1，否则，ns = 0。因为 $\Sigma \succeq \mathbf{0}$ ，这两种表述得到的都是凸规划。在实践中，通常对一个递增的关于 $\mu$ 的序列求解第二种表述；或者对一个递减的关于 $\bar{\sigma}$ 的序列求解第一种表述，由此得到风险-收益平面上的由最优值点组成的曲线，称该曲线是有效前沿(efficient frontier)。当投资组合 $\mathbf{x}$ 使得 $(\mathbf{x}^\top \Sigma \mathbf{x}, \hat{\mathbf{r}}^\top \mathbf{x})$ 在这条曲线上时，称 $\mathbf{x}$ 是有效的投资组合。

**例14.1. (交易所交易基金的分配)** 作为一个数值的例子, 考虑关于分配 $n = 7$ 时的投资组合分配:

1. SPDR道琼斯工业交易所交易基金, DIA;
2. iShare道琼斯交通交易所交易基金, IYT;
3. iShare道琼斯公共事业交易所交易基金, IDU;
4. 第一信托纳斯达克100前科技基金, QQXT;
5. SPDR欧元区斯托克50交易所交易基金, FEZ;
6. iShares巴克莱银行20+Yr Treas. 债券交易所交易基金, TLT;
7. iSharesiBoxx USD高犹太人公司债券交易所交易基金, HYG.

我们使用2013年4月16号之前所观察到的300组每日收益数据, 得到了收益向量的希望和这些资产的协方差矩阵的估计值<sup>2</sup>分别为

$$\hat{\boldsymbol{r}} = \begin{bmatrix} 5.996 & 4.584 & 6.202 & 7.374 & 3.397 & 1.667 & 3.798 \end{bmatrix} \times 10^{-4},$$

$$\hat{\boldsymbol{\Sigma}} = \begin{bmatrix} 0.5177 & 0.596 & 0.2712 & 0.5516 & 0.9104 & -0.3859 & 0.2032 \\ 0.596 & 1.22 & 0.3602 & 0.7671 & 1.095 & -0.4363 & 0.2469 \\ 0.2712 & 0.3602 & 0.3602 & 0.2866 & 0.4754 & -0.1721 & 0.1048 \\ 0.5516 & 0.7671 & 0.2866 & 0.8499 & 1.073 & -0.4363 & 0.2303 \\ 0.9104 & 1.095 & 0.4754 & 1.073 & 2.563 & -0.8142 & 0.4063 \\ -0.3859 & -0.4363 & -0.1721 & -0.4363 & -0.8142 & 0.7479 & -0.1681 \\ 0.2032 & 0.2469 & 0.1048 & 0.2303 & 0.4063 & -0.1681 & 0.1478 \end{bmatrix} \times 10^{-4}.$$

我们在 $\mu$ 所在的区间 $[3.54, 7.37] \times 10^{-4}$ 上以等间隔取 $N = 20$ 个点, 对于每个点对应的 $\mu$ 值, 我们求解 $n_s = 1$ (即无卖空)时的问题(14.1), 由此我们得到了有效前沿的离散点近似, 如图14.1 所示。比如, 有效前沿上与数字13对应的投资组合为

$$\boldsymbol{x} = [0.0145 \ 0.0014 \ 54.79 \ 31.68 \ 0.0008 \ 13.51 \ 0.0067] \times 0.01.$$

这个投资组合中含54.8% 的IDU, 31.7% 的QQXT, 13.5%的TLT; 对应的每日收益为0.0596%, 标准差为0.494%.

<sup>2</sup>估计收益的期望和协方差矩阵是证券组合中非常精细的一步. 这里我们利用向后回溯的这300天的数据, 将经验均值和经验协方差阵分别作为 $\hat{\boldsymbol{r}}$  和 $\hat{\boldsymbol{\Sigma}}$ 的估计值.



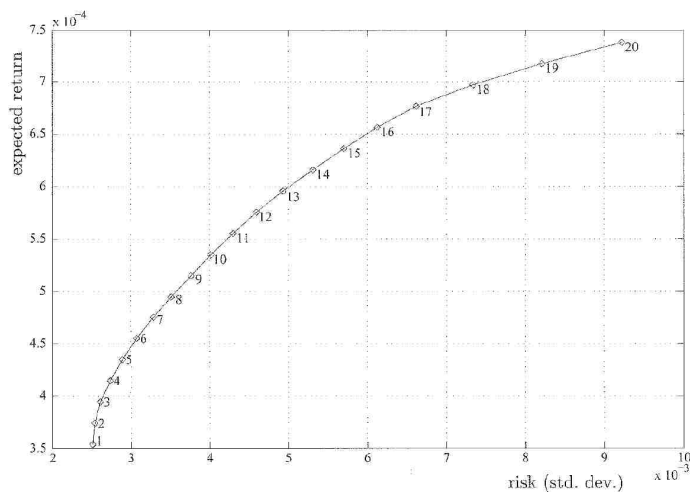


图 14.1: 有效前沿

### §14.1.2 投资组合约束和交易成本

下面讨论关于投资组合的进一步的可能提出的设计约束.

#### §14.1.2.1 板块限界

我们可以给 $\mathbf{x}$ 的某些元素施加上界, 以限制投资仅仅局限于单个的资产. 我们也可以限制在指定板块(比如说对应的前 $k$ 个资产)中的投资总量不要超出总投资的 $\alpha$ 倍. 给问题添加约束

$$\sum_{i=1}^k x_i \leq \alpha \mathbf{1}^\top \mathbf{x}$$

可以对此问题进行建模.

#### §14.1.2.2 多样性

我们还可以施加某些多样性约束. 比如, 我们可以施加对任何 $k(k < n)$ 种资产的投资均不超过投资总量的 $\eta$ 倍, 即

$$s_k(\mathbf{x}) \doteq \sum_{i=1}^k x_{[i]} \leq \eta \mathbf{1}^\top \mathbf{x},$$

其中 $x_{[i]}$ 指 $\mathbf{x}$ 的第 $i$ 大分量, 所以 $s_k(\mathbf{x})$ 是 $\mathbf{x}$ 中前 $k$ 大元素的总和, 参见例9.10. 利用(9.16)中的表示, 多样性约束可以表示为: 存在标量 $t$ 和 $n$ -维向量 $\mathbf{s}$ 使得

$$kt + \mathbf{1}^\top \mathbf{s} \leq \eta \mathbf{1}^\top \mathbf{x}, \quad \mathbf{s} \geq \mathbf{0}, \mathbf{s} \geq \mathbf{x} - t\mathbf{1}.$$

## §14.1.2.3 交易成本

我们在最优化模型中也可以考虑交易成本。特别的，我们易于包含形如

$$\phi(\mathbf{u}) = c\|\mathbf{u}\|_1$$

的正比例成本结构，其中  $c \geq 0$  是单位交易成本。如果存在交易成本，那么预算平衡条件就必须考虑到因交易而带来的开支： $\mathbf{1}^\top \mathbf{x}(0) - \phi(\mathbf{u}) = \mathbf{1}^\top \mathbf{x}$ ，即

$$\mathbf{1}^\top \mathbf{u} + \phi(\mathbf{u}) = 0$$

这个等式约束不是凸的。然而，如果  $\hat{\mathbf{r}} \geq 0$ ，并且我们正在考虑问题(14.1)，则可以将这个方程等价地松弛成不等式(参见例8.19)，因此得到

$$\begin{aligned} \rho(\bar{\sigma}) = \underset{\mathbf{u}}{\text{maximize}} \quad & \hat{\mathbf{r}}^\top \mathbf{x} \\ \text{subject to} \quad & \mathbf{x}^\top \Sigma \mathbf{x} \leq \bar{\sigma}^2 \\ & \mathbf{1}^\top \mathbf{u} + c\|\mathbf{u}\|_1 \leq 0, \\ & \mathbf{x} = \mathbf{x}(0) + \mathbf{u}, \\ & \text{ns} \cdot \mathbf{x} \geq \mathbf{0}. \end{aligned} \tag{14.1}$$

这是一个凸优化。

## §14.1.3 夏普比率(SR)最优化

夏普比(Sharpe ratio, SR)<sup>3</sup>是回报对风险一种度量，它定量地表示了每单位风险的(无风险利率之外)收益的期望值所对应的数量。投资组合的SR定义为

$$\text{SR}(\mathbf{x}) = \frac{\mathbb{E}\{\rho\} - r_f}{\sqrt{\text{var}\{\rho\}}} = \frac{\hat{\mathbf{r}}^\top \mathbf{x} - r_f}{\sqrt{\mathbf{x}^\top \Sigma \mathbf{x}}}$$

其中  $r_f \geq 0$  是**无风险**(risk-free)资产(例如把钱放在银行储蓄所获得的收益是0风险的)的收益，同时也和上面那样假设  $w(0) = \mathbf{1}^\top \mathbf{x}(0) = 1$ 。就收益/风险角度来说，SR越高，表明这个投资越好。假设我们的投资组合没有考虑无风险资产，且在  $\hat{\mathbf{r}}^\top \mathbf{x} > r_f$ ,  $\Sigma \succ \mathbf{0}$  的条件下，从几何上来讲，极大化SR最大的投资组合对应着有效前沿与过  $(0, r_f)$  点的直线(资金分配线, capital allocation line, CAL)的切点，如图14.2 所示。

<sup>3</sup>以William Sharpe命名，Nobel laureate(1990)

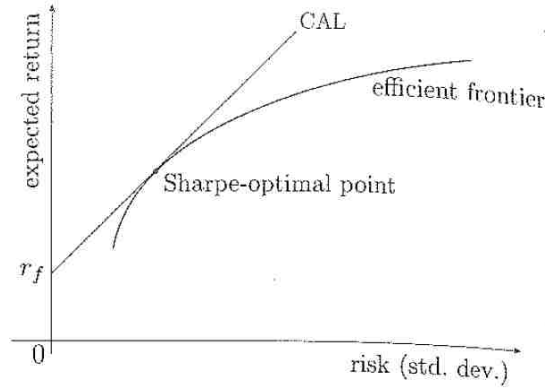


图 14.2: 资金分配线和Sharp最优点.

在如上的这些假设下, 寻求夏普最优的投资组合可通过求解如下问题, 即

$$\begin{aligned} & \underset{\mathbf{x}}{\text{maximize}} && \frac{\hat{\mathbf{r}}^\top \mathbf{x} - r_f}{\sqrt{\mathbf{x}^\top \boldsymbol{\Sigma} \mathbf{x}}} \\ & \text{subject to} && \hat{\mathbf{r}}^\top \mathbf{x} > r_f \\ & && \mathbf{x} \in \mathcal{X}. \end{aligned} \quad (14.2)$$

此种形式的问题是非凸的。然而, 我们可以进行如下操作: 给 $\text{SR}(\mathbf{x})$ 的分子和分母都乘上一个松弛变量 $\gamma > 0$ , 并令

$$\tilde{\mathbf{x}} \doteq \gamma \mathbf{x} = \gamma \mathbf{x}(0) + \tilde{\mathbf{u}}; \quad \tilde{\mathbf{u}} \doteq \gamma \mathbf{u}.$$

然后, 因为最大化 $\text{SR}(\mathbf{x})$ 等价于最小化 $1/\text{SR}(\mathbf{x})$ 是等价的(因为分子和分母都是正的), 我们就可以把问题重新写作

$$\begin{aligned} & \underset{\tilde{\mathbf{u}}, \gamma > 0}{\text{minimize}} && \frac{\sqrt{\tilde{\mathbf{x}}^\top \boldsymbol{\Sigma} \tilde{\mathbf{x}}}}{\hat{\mathbf{r}}^\top \tilde{\mathbf{x}} - \gamma r_f} \\ & \text{subject to} && \hat{\mathbf{r}}^\top \tilde{\mathbf{x}} > \gamma r_f, \\ & && \mathbf{1}^\top \tilde{\mathbf{u}} = 0, \\ & && \tilde{\mathbf{x}} = \gamma \mathbf{x}(0) + \tilde{\mathbf{u}}, \\ & && \text{ns} \cdot \tilde{\mathbf{x}} \geq 0. \end{aligned}$$

注意到这个问题关于变量 $(\tilde{\mathbf{u}}, \gamma)$ 是其次的, 即如果 $(\tilde{\mathbf{u}}, \gamma)$ 是一个解, 那么 $\alpha(\tilde{\mathbf{u}}, \gamma)$ 也是一个解 ( $\alpha > 0$ )。我们通过规范化目标函数的分母为 1 来解决齐次

性, 即我们施加  $\hat{\mathbf{r}}^\top \tilde{\mathbf{x}} - \gamma r_f = 1$ . 使用这种规范化, 上述问题就变成了

$$\begin{aligned} & \underset{\tilde{\mathbf{u}}, \gamma > 0}{\text{minimize}} && \sqrt{\tilde{\mathbf{x}}^\top \Sigma \tilde{\mathbf{x}}} \\ & \text{subject to} && \hat{\mathbf{r}}^\top \tilde{\mathbf{x}} - \gamma r_f = 1, \\ & && \mathbf{1}^\top \tilde{\mathbf{u}} = 0, \\ & && \tilde{\mathbf{x}} = \gamma \mathbf{x}(0) + \tilde{\mathbf{u}}, \\ & && \text{ns} \cdot \tilde{\mathbf{x}} \geq 0. \end{aligned}$$

因为  $\Sigma \succ 0$ , 我们可以分解为  $\Sigma = \Sigma^{1/2} \Sigma^{1/2}$  (参见§4.4.4), 因此  $\sqrt{\tilde{\mathbf{x}}^\top \Sigma \tilde{\mathbf{x}}} = \|\Sigma^{1/2} \tilde{\mathbf{x}}\|_2$ . 这样, SR的最优化问题就可以转化成形如

$$\begin{aligned} & \underset{\tilde{\mathbf{u}}, \gamma > 0}{\text{minimize}} && t \\ & \text{subject to} && \|\Sigma^{1/2} \tilde{\mathbf{x}}\|_2 \leq t, \\ & && \hat{\mathbf{r}}^\top \tilde{\mathbf{x}} - \gamma r_f = 1, \\ & && \mathbf{1}^\top \tilde{\mathbf{u}} = 0, \\ & && \tilde{\mathbf{x}} = \gamma \mathbf{x}(0) + \tilde{\mathbf{u}}, \\ & && \text{ns} \cdot \tilde{\mathbf{x}} \geq 0 \end{aligned} \tag{14.3}$$

的SOCP问题. 一旦由这个问题解出  $\tilde{\mathbf{u}}, \gamma$ , 我们就可以由  $\mathbf{u} = \gamma^{-1} \tilde{\mathbf{u}}$  恢复出初始变量  $\mathbf{u}$ . 作为一个数值例子, 在例14.1提供的数据之上求解问题(14.3), 这里假设无风险利率  $r_f = 1.5 \times 10^{-4}$ , 则我们得到夏普最优的投资组合

$$\mathbf{x} = [2.71 \ 0.00 \ 30.95 \ 14.08 \ 0.00 \ 23.82 \ 28.44]^\top \times 0.01,$$

对应的收益的期望和标准差分别是  $0.0460 * 0.01$  和  $0.3122 * 0.01$ . 这个投资组合在收益/风险有效前沿上, 是CAL 直线与有效前沿的切点, 如图14.3所示.

#### §14.1.4 风险价值(Value-at-Risk, VaR)最优化

假设收益向量  $\mathbf{r}$  服从正态分布, 且已知均值  $\hat{\mathbf{r}}$  和协方差阵  $\Sigma$ , 所谓的投资组合  $\mathbf{x}$  的水平为  $\alpha \in (0, 1)$  的**风险价值**(Value-at-Risk, VaR)定义为(参见例10.4)

$$\begin{aligned} \text{VaR}_\alpha(\mathbf{x}) &= -\sup\{\zeta : \text{Prob}\{\rho(\mathbf{x}) \leq \zeta\} \leq \alpha\} \\ &= -\inf\{\zeta : \text{Prob}\{\rho(\mathbf{x}) \leq \zeta\} \geq \alpha\} \\ &= -F_\rho^{-1}(\alpha) \end{aligned}$$

其中  $F_{\rho^{-1}(\alpha)}$  表示随机变量  $\rho(\mathbf{x})$  的累积分布函数的反函数, 这里  $\rho(\mathbf{x})$  表示投资组合的收益, 服从正态分布. 与经典的方差度量相比, 投资组合的VaR更

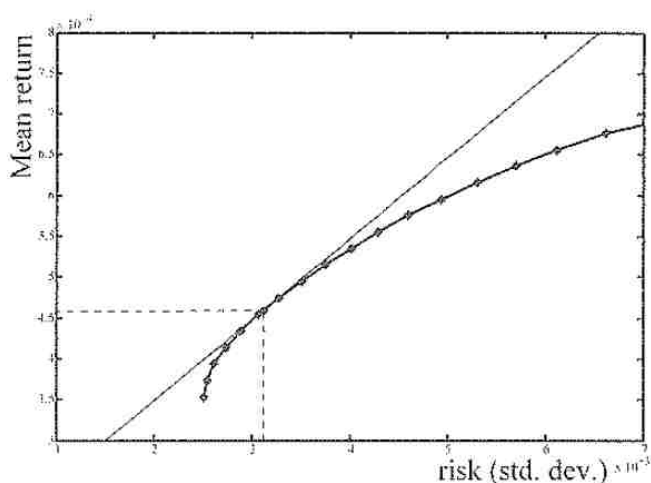


图 14.3: 当  $r_f = 1.5 \times 10^{-4}$  时 Sharp 最优的投资组合

受青睐, 是一个被广泛采用的风险度量。VaR $_{\alpha}$ 的意思是投资导致高于VaR $_{\alpha}$ 的损失(loss)<sup>4</sup>的概率不大于 $\alpha$ 。等价地说, 投资者的收益高于 $-VaR_{\alpha}$ 的概率不低于 $1 - \alpha$ 。在例10.4中我们已经看到了, 对于给定的 $\alpha \in (0, 0.5)$ ,

$$VaR_{\alpha} \leq \gamma \Leftrightarrow \Phi^{-1}(1 - \alpha) \|\Sigma^{1/2} \mathbf{x}\|_2 \leq \gamma + \hat{\mathbf{r}}^T \mathbf{x}, \quad (14.4)$$

其中 $\Phi$ 是标准正态累积分布函数。在保证关于VaR $_{\alpha}$ 的上界 $\gamma$ 的限制下, 最大化收益的期望即可得到最优的投资组合。求解形如

$$\begin{aligned} \rho(\gamma) = & \text{maximize}_{\mathbf{x}} \quad \hat{\mathbf{r}}^T \mathbf{x} \\ \text{subject to} \quad & \Phi^{-1}(1 - \alpha) \|\Sigma^{1/2} \mathbf{x}\|_2 \leq \hat{\mathbf{r}}^T \mathbf{x} + \gamma, \\ & \mathbf{x} \in \mathcal{X}. \end{aligned} \quad (14.5)$$

的SOCP可以得到这个最优的投资组合。

## §14.2 鲁棒的投资组合最优化

前面所说的投资组合优化模型是建立在随机收益向量的矩(例如: 期望和协方差矩阵)或整个分布是精确已知的假设条件下。然而, 在实践中, 这些量是不知道的, 可能需要结合历史数据和专家知识去估计。因此, 由于估计误差和/或在描述市场收益的随机过程性质的错误先验(priori)假设, 将导致这些估计量具有很高的不确定性。反过来, 由使用“标称”(nominal)数

<sup>4</sup>我们的损失(loss) $\ell$ 的意思是收益的相反数, 即 $\ell = -\rho$ 。这样, 用我们的记号, VaR $_{\alpha}$ 是一个典型的正值, 描述了投资者(以概率 $1 - \alpha$ )所面对的最大的损失。

据的最优化问题所得到的投资组合分配决策远远偏离实践中的最优策略，即最优的投资组合对输入数据非常敏感。克服这些问题的一种途径是一定程度上考虑这些数据存在的不确定性的先验(priori)信息，并寻求“鲁棒”(robust)的投资组合，其对于不确定性而言，在最坏情况下也表现良好。

### §14.2.1 鲁棒的均值-方差最优化

用鲁棒的均值方差投资组合优化方法中，我们假设 $\hat{\mathbf{r}}$ 和 $\Sigma$ 是不确定的，我们通过为这些参数确定隶属合集来建模这种不确定性，即我们假设 $(\hat{\mathbf{r}}, \Sigma)$ 属于某一给定的有界不确定集 $\mathcal{U}$ 。对于一个给定的投资组合 $\mathbf{x}$ ，最坏情况的投资组合的方差

$$\sigma_{\text{WC}}^2 = \sup_{(\hat{\mathbf{r}}, \Sigma) \in \mathcal{U}} \mathbf{x}^\top \Sigma \mathbf{x}.$$

最小方差投资组合设计问题(14.1)的鲁棒化版本是

$$\begin{aligned} \bar{\sigma}_{\text{WC}}^2(\mu) = & \underset{\mathbf{u}}{\text{minimize}} \quad \sup_{(\hat{\mathbf{r}}, \Sigma) \in \mathcal{U}} \mathbf{x}^\top \Sigma \mathbf{x} \\ & \text{subject to} \quad \inf_{(\hat{\mathbf{r}}, \Sigma) \in \mathcal{U}} \hat{\mathbf{r}}^\top \mathbf{x} \geq \mu, \\ & \mathbf{x} \in \mathcal{X}. \end{aligned} \quad (14.6)$$

只有在不确定集 $\mathcal{U}$ 足够“简单”时，才能高效准确地求解这个问题。比如 $\mathcal{U} = \{(\hat{\mathbf{r}}, \Sigma) : \hat{\mathbf{r}} \in \mathcal{U}_r, \Sigma \in \mathcal{U}_\Sigma\}$ ，其中 $\mathcal{U}_r$ 和 $\mathcal{U}_\Sigma$ 是区间集，即

$$\mathcal{U}_r = \{\hat{\mathbf{r}} : r_{\min} \leq \hat{\mathbf{r}} \leq r_{\max}\},$$

$$\mathcal{U}_\Sigma = \{\Sigma : \Sigma_{\min} \preceq \Sigma \preceq \Sigma_{\max}, \Sigma \succeq \mathbf{0}\}.$$

如果进一步假定 $\mathbf{x} \geq \mathbf{0}$ 和 $\Sigma_{\max} \succeq \mathbf{0}$ ，则问题(14.6)等价于

$$\begin{aligned} \bar{\sigma}_{\text{WC}}^2(\mu) = & \underset{\mathbf{u}}{\text{minimize}} \quad \mathbf{x}^\top \Sigma_{\max} \mathbf{x} \\ & \text{subject to} \quad \hat{\mathbf{r}}_{\min}^\top \mathbf{x} \geq \mu, \\ & \mathbf{1}^\top \mathbf{u} = 0, \\ & \mathbf{x} = \mathbf{x}(0) + \mathbf{u}, \\ & \mathbf{x} \geq \mathbf{0}. \end{aligned} \quad (14.7)$$

如果没有施加约束 $\mathbf{x} \geq \mathbf{0}$ ，可以用一个特定的鞍点(ad-hoc saddle-point)算法来高效求解问题(14.6)。该算法超出本节内容，不予详述<sup>5</sup>。

<sup>5</sup>感兴趣的读者可以参考Tütüncü和Koenig的文章: Robust asset allocation, Annals of Operations Research, 2004.

**情景法.** 对于一般的不确定集 $\mathcal{U}_r, \mathcal{U}_\Sigma$ 及与之相对应的问题, 我们无法找到一个有效的方式来表述鲁棒性, 如§10.3.4 描述的那样, 我们可以用情景法来逼近鲁棒性。在当前的背景下, 我们应该假设定义在 $\mathcal{U}_r, \mathcal{U}_\Sigma$  上的概率分布, 并收集了 $N$ 组独立同分布样本 $(\hat{\mathbf{r}}^{(i)}, \Sigma^{(i)}), i = 1, \dots, N$ , 然后我们求解情景问题

$$\begin{aligned} & \underset{\mathbf{x}, t}{\text{minimize}} && t \\ & \text{subject to} && \mathbf{x}^\top \Sigma^{(i)} \mathbf{x} \leq t, i = 1, \dots, N, \\ & && \hat{\mathbf{r}}^{(i)\top} \mathbf{x} \geq \mu, \quad i = 1, \dots, N, \\ & && \mathbf{x} \in \mathcal{X}. \end{aligned} \tag{14.8}$$

如果这个情景问题中的 $N$ 与式(10.30)所描述的相匹配, 则求解这个问题可以得到一个具有理想水平的概率鲁棒的证券组合。这种方法的优点是适用于一般的不确定性结构, 并且从类别上来说不会增加求解问题的难度, 即“鲁棒化”问题(14.8)仍然是一个凸优化问题。

### §14.2.2 鲁棒的风险价值(VaR)最优化

我们在§14.1.4已经看到, 当假设收益向量 $\mathbf{r}$ 服从均值和协方差矩阵分别为 $\hat{\mathbf{r}}$ 和 $\Sigma$ 的正态分布时, 式(14.4)给出了投资组合 $\mathbf{x}$ 的水平为 $\alpha \in (0, 0.5)$ 的风险价值。然而, 在实践中, 我们可能会遇到两种不确定性来源, 其妨碍我们直接使用这个公式。第一种是实际收益的分布不是精确的正态分布(分布不确定性); 第二种是参数 $\hat{\mathbf{r}}, \Sigma$  的值可能不准确(参数不确定性)。接下来我们讨论在这两种不确定性下如何处理VaR投资组合最优化问题。

#### §14.2.2.1 分布不确定性下的鲁棒VaR

首先, 我们假设收益分布的参数 $\hat{\mathbf{r}}, \Sigma$ 是精确已知的, 但分布本身未知。然后我们考虑一类均值为 $\hat{\mathbf{r}}$ , 协方差阵为 $\Sigma$  的概率分布 $\mathcal{P}$ 。最坏的情况下, 损失大于某水平 $\zeta$ 的概率是

$$\text{Prob}_{\text{WC}}\{-\rho(\mathbf{x}) \geq \zeta\} \doteq \sup_{\mathcal{P}} \text{Prob}\{-\rho(\mathbf{x}) \geq \zeta\}, \tag{14.9}$$

其中上确界是关于 $\mathcal{P}$ 中的所有分布来计算的。最坏情况的 $\alpha - \text{VaR}$ 定义为

$$\text{wc-VaR}_\alpha(\mathbf{x}) \doteq \sup\{\zeta : \text{Prob}_{\text{WC}}\{-\rho(\mathbf{x}) \geq \zeta\} \leq \alpha\},$$

显然有

$$\text{wc-VaR}_\alpha(\mathbf{x}) \leq \gamma \quad \Leftrightarrow \quad \text{Prob}_{\text{WC}}\{-\rho(\mathbf{x}) \geq \gamma\} \leq \alpha. \tag{14.10}$$

而Chebyshev-Cantelli不等式表明, 对于任意方差有限的随机变量 $z$ 及 $t > 0$ , 有

$$\sup \text{Prob}\{z - \mathbb{E}\{z\} \geq t\} = \frac{\text{var}(z)}{\text{var}(z) + t^2},$$

其中的上确界是对所有具有相同的均值和相同的协方差的分布来计算的。现在, 考虑(14.9)中的随机变量 $\rho(\mathbf{x})$ , 它的期望为 $\hat{\mathbf{r}}^\top \mathbf{x}$ , 方差为 $\|\Sigma^{1/2} \mathbf{x}\|_2^2$ , 将Chebyshev-Cantelli 不等式应用于 $\rho(\mathbf{x})$ , 我们有

$$\begin{aligned} \text{Prob}_{\text{WC}}\{-\rho(\mathbf{x}) \geq \zeta\} &= \text{Prob}_{\text{WC}}\{-\rho(\mathbf{x}) + \hat{\mathbf{r}}^\top \mathbf{x} \geq \zeta + \hat{\mathbf{r}}^\top \mathbf{x}\} \\ &= \frac{\|\Sigma^{1/2} \mathbf{x}\|_2^2}{\|\Sigma^{1/2} \mathbf{x}\|_2^2 + (\zeta + \hat{\mathbf{r}}^\top \mathbf{x})^2}. \end{aligned}$$

因此, 由(14.10)可知, 对于 $\alpha \in (0, 0.5)$ ,

$$\text{WC-VaR}_\alpha(\mathbf{x}) \leq \gamma \Leftrightarrow \kappa(\alpha) \left\| \Sigma^{1/2} \mathbf{x} \right\|_2 \leq \hat{\mathbf{r}}^\top \mathbf{x} + \gamma,$$

其中 $\kappa(\alpha) \doteq \sqrt{\frac{1-\alpha}{\alpha}}$ . 如上, 在VaR优化问题(14.5)中, 将SOC 约束的系数 $\Phi^{-1}(1-\alpha)$ 替换成 $\kappa(\alpha)$ , 即得对应的分布不确定性下的鲁棒VaR优化问题。

#### §14.2.2.2 矩不确定性下的鲁棒VaR

除了分布不确定性外, 如果参数 $\hat{\mathbf{r}}$ 和 $\Sigma$ 的取值也具有不确定性, 那么我们可以用如下方法<sup>6</sup> 来确定鲁棒的VaR 证券组合。假设描述 $\hat{\mathbf{r}}, \Sigma$ 的不确定性的集合是区间型的, 即

$$\begin{aligned} \mathcal{U}_{\hat{\mathbf{r}}} &= \{\hat{\mathbf{r}} : \mathbf{r}_{\min} \leq \hat{\mathbf{r}} \leq \mathbf{r}_{\max}\}, \\ \mathcal{U}_{\Sigma} &= \{\Sigma : \Sigma_{\min} \preceq \Sigma \preceq \Sigma_{\max}\}, \end{aligned}$$

并且存在一个 $\Sigma \in \mathcal{U}_{\Sigma}$ 使得 $\Sigma \succ \mathbf{0}$ 。然后, 定义 $\text{WC-VaR}_\alpha(\mathbf{x})$  是 $\text{VaR}_\alpha(\mathbf{x})$ 关于所有 $\hat{\mathbf{r}} \in \mathcal{U}_{\hat{\mathbf{r}}}$ ,  $\Sigma \in \mathcal{U}_{\Sigma}$ , 以及所有均值为 $\hat{\mathbf{r}}$ , 协方差阵为 $\Sigma$  的分布取上确界, 可以证明 $\text{WC-VaR}_\alpha(\mathbf{x}) \leq \gamma$  当且仅当存在对称矩阵 $\Lambda_+, \Lambda_- \in \mathbb{S}_+^n$ , 向量 $\lambda_+, \lambda_- \in \mathbb{R}^n$ 和标量 $v \in \mathbb{R}$ , 使得

$$\begin{aligned} \text{trace}(\Lambda_+ \Sigma_{\max}) - \text{trace}(\Lambda_- \Sigma_{\min}) + k^2(\alpha)v + \lambda_+^\top \mathbf{r}_{\max} - \lambda_-^\top \mathbf{r}_{\min} &\leq \gamma, \\ \begin{bmatrix} \Lambda_+ - \Lambda_- & \mathbf{x}/2 \\ \mathbf{x}^\top/2 & v \end{bmatrix} &\succeq \mathbf{0}, \\ \lambda_+ &\geq \mathbf{0}, \lambda_- \geq \mathbf{0}, \Lambda_+ \succeq \mathbf{0}, \Lambda_- \succeq \mathbf{0}, \\ \mathbf{x} &= \lambda_- - \lambda_+. \end{aligned} \tag{14.11}$$

<sup>6</sup>关于该方法的描述详见论文: EI Ghaoui, Oks and Oustry, Worst-case value-at-risk and robust portfolio optimization: a conic programming approach, Operations Research, 2003.



在分布和矩的双重不确定性下,与VaR优化问题(14.5)对应的鲁棒VaR优化可以显式地表述成如下的SDP,即

$$\begin{aligned} & \underset{\mathbf{u}, \mathbf{\Lambda}_+, \mathbf{\Lambda}_-, \boldsymbol{\lambda}_+, \boldsymbol{\lambda}_-, v}{\text{minimize}} && \hat{\mathbf{r}}^\top \mathbf{x} \\ & \text{subject to} && (14.11), \\ & && \mathbf{x} \in \mathcal{X}. \end{aligned}$$

### §14.3 多期投资组合最优化

前两节中的投资分配模型仅仅考虑一个投资期,因此关注的是投资组合在投资期末的某种业绩。然而,在实践中,一名投资者在本投资期结束时,将再次面临如何在下一个周期重新分配他的资产的问题。当然,他可以再次求解单一投资期的分配问题,然后依次重复下去。然而,从长远目标来看,这种基于按期策略的决策有可能是“目光短浅的”(myopic),即如果投资者提前知道他的投资目标是设定在自此向前的 $T(\geq 1)$ 个投资期,投资者一开始就考虑自己投资问题的多周期性有利于得到更有远见的策略。

在接下来的几节,我们将讨论两种在多期投资组合上用以优化投资决策的**数据驱动**(data-driven)技术。第一种技术我们称之为**开环**(open-loop)方法,其目标是在时刻 $k=0$ 确定整个未来的投资组合调整量,这包括从时刻0到最终时刻 $T$ 上的投资组合调整量序列。第二种技术我们称之为**闭环**(closed-loop)策略,它增加了未来决策的灵活性,从而可看成第一种策略的提升。具体地,它将未来决策作为观测收益的函数(所谓的方法),由此以减少未来不确定性的影响。两种方法都是数据驱动的,即它们是基于收益流的 $N$ 种仿真情景的可得性。

我们从一些记号和准备性工作开始:资产集合记为 $a_1, \dots, a_n$ ,时刻 $k\Delta$ 时资产 $a_i$ 的市场价格记为 $p_i(k)$ ,其中 $k$ 是一个带符号的整数, $\Delta$ 是一个固定的时间段。对资产 $i$ 的投资在第 $k$ 期(从 $(k-1)\Delta$ 到 $k\Delta$ )的简单收益是

$$r_i(k) \doteq \frac{p_i(k) - p_i(k-1)}{p_i(k-1)}, i = 1, \dots, n; k = 1, 2, \dots,$$

对应的**增益**(gain)定义为

$$g_i(k) \doteq 1 + r_i(k), i = 1, \dots, n; k = 1, 2, \dots$$

第 $k$ 期资产的收益向量记为 $\mathbf{r}(k) \doteq [r_1(k) \cdots r_n(k)]^\top$ ,  $\mathbf{g}(k)$ 是对应的增益向量。符号 $\mathbf{G}(k) = \text{diag}(\mathbf{g}(k))$ 表示对角线元素是 $\mathbf{g}(k)$ 的分量的对角矩阵。假

定收益和增益向量是随机向量。我们用  $T \geq 1$  表示需要做分配决策的正向周期数目，并假设有一个场景生成先哲(scenario-generating oracle)，它能产生正向收益流  $\{\mathbf{r}(1), \dots, \mathbf{r}(T)\}$  的  $N$  个独立同分布的样本。

### §14.3.1 投资组合动力学

我们考虑  $T$  个时段(或阶段(stage))的决策问题。在每个时段我们有机会重新平衡我们的投资组合，在保证满足每个时段的的投资组合约束的前提下，使得最后一个时段的某种恰当的成本函数(稍后讨论)取最小值。如图14.4所示，第  $k$  个决策时段从时刻  $k-1$  开始，到时刻  $k$  结束。

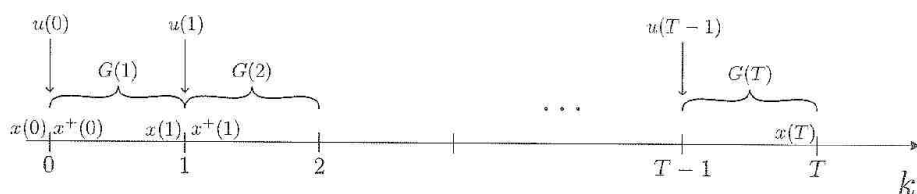


图 14.4: 投资决策水平和时段

我们用  $x_i(k)$  表示投资者时刻  $k$  在证券  $a_i$  上的投资占总资产的比例。时刻  $k$  的投资组合向量

$$\mathbf{x}(k) \doteq [x_1(k) \quad \dots \quad x_n(k)]^\top.$$

时刻  $k$  投资者的总资产

$$w(k) \doteq \sum_{i=1}^n x_i(k) = \mathbf{1}^\top \mathbf{x}(k).$$

令  $\mathbf{x}(0)$  是时刻  $k=0$  时给定的初始投资组合布局(例如，可以假设  $\mathbf{x}(0)$  仅一个分量非零，其代表初始的可用资金数目)。在时刻  $k=0$  时，我们有机会在市场中进行交易，并且因此可以增加或减少对每项资产的投资数目，即调整投资组合。在交易过后，调整后的投资组合是  $\mathbf{x}^+ = \mathbf{x}(0) + \mathbf{u}(0)$ ，如果我们增加了对第  $i$  个资产的投资，则  $u_i > 0$ ；若减少了，则  $u_i < 0$ ；若不变，则  $u_i = 0$ 。现在假设投资组合在第一时段  $\Delta$  内固定不变。在这个时段结束时，投资组合布局是

$$\mathbf{x}(1) = \mathbf{G}(1)\mathbf{x}^+(0) = \mathbf{G}(1)\mathbf{x}(0) + \mathbf{G}(1)\mathbf{u}(0),$$

其中  $\mathbf{G}(1) = \text{diag}(g_1(1), \dots, g_n(1))$  是对角矩阵，对角线元素是时刻 0 到时刻 1 这个时段的资产增益。在时刻  $k=1$ ，我们再次对投资组合调整  $\mathbf{u}(1)$ ，

即  $\mathbf{x}^+(1) = \mathbf{x}(1) + \mathbf{u}(1)$ , 然后在紧接着的这个时段  $\Delta$  内持有这个更新后的投资组合。因此, 时刻  $k = 2$  的投资组合布局是

$$\mathbf{x}(2) = \mathbf{G}(2)\mathbf{x}^+(1) = \mathbf{G}(2)\mathbf{x}(1) + \mathbf{G}(2)\mathbf{u}(1).$$

对  $k = 0, 1, 2, \dots$ , 依次按这种方式作用, 我们得到了第  $(k+1)$  时段末的投资组合布局的迭代动力学方程

$$\mathbf{x}(k+1) = \mathbf{G}(k+1)\mathbf{x}(k) + \mathbf{G}(k+1)\mathbf{u}(k), k = 0, \dots, T-1 \quad (14.12)$$

和刚经过第  $(k+1)$  次交易后的投资组合布局方程(参见图14.4)

$$\mathbf{x}^+(k) = \mathbf{x}(k) + \mathbf{u}(k). \quad (14.13)$$

由(14.12)可得时刻  $k = 1, \dots, T$  的(随机的)投资组合布局是

$$\mathbf{x}(k) = \Phi(1, k)\mathbf{x}(0) + \sum_{j=1}^k \Phi(j, k)\mathbf{u}(j-1), \quad (14.14)$$

其中我们定义  $\Phi(v, k) (v \leq k)$  是从时段  $v$  开始到时段  $k$  结束的**复合增益**(compounded gain) 矩阵, 即

$$\Phi(v, k) \doteq \mathbf{G}(k)\mathbf{G}(k-1)\cdots\mathbf{G}(v), \quad \Phi(k, k) \doteq \mathbf{G}(k).$$

投资组合的表达式可以重新紧促地重写为

$$\mathbf{x}(k) = \Phi(1, k)\mathbf{x}(0) + \Omega_k \mathbf{u},$$

其中

$$\begin{aligned} \mathbf{u} &\doteq [\mathbf{u}(0)^\top \cdots \mathbf{u}(T-2)^\top \mathbf{u}(T-1)^\top]^\top, \\ \Omega_k &\doteq [\Phi(1, k) \cdots \Phi(k-1, k) \quad \Phi(k, k) | 0 \cdots 0]. \end{aligned}$$

因此, 我们的总资产

$$w(k) = \mathbf{1}^\top \mathbf{x}(k) = \phi(1, k)^\top \mathbf{x}(0) + \mathbf{w}_k^\top \mathbf{u}$$

其中

$$\begin{aligned} \phi(v, k)^\top &\doteq \mathbf{1}^\top \Phi(v, k), \\ \mathbf{w}_k^\top &\doteq \mathbf{1}^\top \Omega_k = [\phi(1, k)^\top \cdots \phi(k-1, k)^\top \quad \phi(k, k)^\top | 0 \cdots 0]. \end{aligned}$$

我们考虑自筹资金的投资组合，即

$$\sum_{i=1}^n u_i(k) = 0, \quad k = 0, 1, \dots, T-1,$$

并且通过强迫更新后的投资组合  $\mathbf{x}^+(k)$  属于多面体  $\mathcal{X}(k)$  允许我们在模型中包含一般的线性约束。投资在整体范围内的累积总收益是

$$\rho(\mathbf{u}) \doteq \frac{w(T)}{w(0)} = \frac{\mathbf{1}^\top \mathbf{x}(T)}{\mathbf{1}^\top \mathbf{x}(0)} = \frac{\phi(1, T)^\top \mathbf{x}(0)}{\mathbf{1}^\top \mathbf{x}(0)} + \frac{1}{\mathbf{1}^\top \mathbf{x}(0)} \mathbf{w}_T^\top \mathbf{u}.$$

我们看到  $\rho(\mathbf{u})$  是决策变量  $\mathbf{u}$  的仿射函数，系数是随机向量  $\mathbf{w}_T$ ，其依赖于  $T$  个时段上的随机增益。

### §14.3.2 最优的开环策略

现在假设可以从场景生成先哲那儿得到各时段增益的  $N$  个独立同分布样本(场景)  $\{\mathbf{G}^{(i)}(k), k = 1, \dots, T\}$ ,  $i = 1, \dots, N$ 。这些样本因此对每个矩阵  $\mathbf{\Omega}_k$ ,  $k = 1, \dots, T$ ，进而对  $\omega(k)$  和  $\phi(1, k)$  等向量依次产生了  $N$  个场景。我们分别用  $\mathbf{\Omega}_k^{(i)}, \omega_k^{(i)}, \phi^{(i)}(1, k)$  表示这些场景，用  $\mathbf{x}^{(i)}(k), \mathbf{w}^{(i)}(k), \rho^{(i)}(\mathbf{u})$  分别表示第  $i$  个场景中，时刻  $k$  的投资组合布局，时刻  $k$  的总资产以及累积的期末收益。

利用采样的场景，我们可以构造几个可能的经验风险测度，将其作为我们最小化的目标。常用的目标如下：令  $\gamma$  是给定的期末收益水平(通常但不是必须地，将  $\gamma$  设置成累积期末收益的平均值，即  $\gamma = \frac{1}{N} \sum_{i=1}^N \rho^{(i)}$ )：

$$J_1 \doteq \frac{1}{N} \sum_{i=1}^N |\gamma - \rho^{(i)}|, \quad (14.15)$$

$$J_2 \doteq \frac{1}{N} \sum_{i=1}^N (\gamma - \rho^{(i)})^2, \quad (14.16)$$

$$J_{p1} \doteq \frac{1}{N} \sum_{i=1}^N \max(0, \gamma - \rho^{(i)}), \quad (14.17)$$

$$J_{p2} \doteq \frac{1}{N} \sum_{i=1}^N (\max(0, \gamma - \rho^{(i)}))^2. \quad (14.18)$$

这些目标表示收益与  $\gamma$  间偏差的经验平均值。特别地，(14.15) 正比例于收益  $\rho^{(i)}$  与  $\gamma$  的偏差的  $\ell_1$  范数；(14.16) 正比例于这些偏差的  $\ell_2$  范数。(14.17) 和 (14.18) 是非对称的测度(有时称之为下偏距(lower partial moments))： $J_{p1}$  度量收益变量  $\rho^{(i)}$  低于水平  $\gamma$  的经验平均值； $J_{p2}$  度量这种偏差的平方的平均

值。选择一阶还是二阶成本测度取决于投资者的**风险规避**(risk aversion)水平。因为二阶成本是残量的平方，因此更看重成本，所以高阶测度反映了对风险规避的水平更高。

我们的开环多阶段分配策略即寻求投资组合调整参数  $\mathbf{u} = (\mathbf{u}(0), \dots, \mathbf{u}(T-1))$ ，它最小化上面描述的某个成本测度，同时受限制于整个范围内每个时段所给定的投资组合约束，即考虑

$$\begin{aligned} J^*(\gamma) = & \underset{\mathbf{u}}{\text{minimize}} \quad J(\mathbf{u}) \\ \text{subject to} \quad & \mathbf{x}^{(i)+}(k) \in \mathcal{X}(k), k = 0, \dots, T-1; i = 1, \dots, N, \\ & \mathbf{1}^\top \mathbf{u}(k) = 0, k = 0, \dots, T-1, \end{aligned}$$

其中  $J(\mathbf{u})$  是刚才提到的一个成本， $\mathbf{x}^{(i)+}(k)$  是第  $i$  个采样场景中由 (14.13) 和 (14.14) 确定的。求解一个线性规划问题或者一个凸二次规划问题是确定这些最优分配数值解的有效方法。

### §14.3.3 带仿射价格的闭环分配

因为所有的调整量  $\mathbf{u}(0), \dots, \mathbf{u}(T-1)$  是在时刻  $k=0$  计算出来的，所以上节研究的开环策略在实践中实施时可能是次优的。尽管第一个决策  $\mathbf{u}(0)$  必定立即实施(现时变量)，但是实际上可以等到获得向前推的一个时段的收益的实际结果后，再对未来的决策做出判断，这些观察结果可以减小不确定性。例如，在时刻  $k \geq 1$ ，当我们需要实施  $\mathbf{u}(k)$  时，我们已经**观察到**(observed)了从时段1到时段  $k$  的资产收益的一个实现。因此，我们可以利用这个信息来考虑该分配决策  $\mathbf{u}(k)$  的**条件**(conditional)分配决策，由此对之前时段的收益做出反应。这蕴含着，我们不应专注于固定的决策，而是确定一些适当的方法(policy)，其描述了基于时刻 1 至  $k$  的观测收益流以后如何作出实际决策。当确定投资方法的结构时，我们应该对所得到的最优化问题的普适性和可计算性进行折衷。一些研究工作<sup>7</sup>发现，线性或者仿射方法是一种有效的折衷方法，因为用凸优化技术可以有效地计算反应策略。在本节中，我们遵循该思路，考虑如下形式的仿射方法，即

$$\mathbf{u}(k) = \bar{\mathbf{u}}(k) + \Theta(k)(\mathbf{g}(k) - \bar{\mathbf{g}}(k)), k = 1, \dots, T-1, \quad (14.19)$$

且  $\mathbf{u}(0) = \bar{\mathbf{u}}(0)$ ，这里  $\bar{\mathbf{u}}(k) \in \mathbb{R}^n (k = 0, \dots, T-1)$  是“名义上的”分配决策变量， $\mathbf{g}(k)$  是第  $k$  个时段的增益向量， $\bar{\mathbf{g}}(k)$  是给定的对  $\mathbf{g}(k)$  的期望值的一

<sup>7</sup> 比如参见：G. Calafiore, Multi-period portfolio optimization with linear control policies, *Automatica*, 2008.

个估计值,  $\Theta(k) \in \mathbb{R}^{n \times n} (k = 0, \dots, T-1)$  是策略“反应矩阵”, 它的作用是对名义分配的调整量和增益  $g(k)$  与它的期望值的偏差成比例。因为预算平衡约束  $\mathbf{1}^\top \mathbf{u}(k) = 1$  要对任意的增益的实现都满足, 所以我们施加了限制

$$\mathbf{1}^\top \bar{\mathbf{u}}(k) = 0, \mathbf{1}^\top \Theta(k) = 0, k = 1, \dots, T-1.$$

### §14.3.3.1 仿射方法下的投资组合动力学

应用调整方法(14.19), 即将(14.19)代入投资组合动力学方程(14.12)和(14.13), 我们得到

$$\mathbf{x}^+(k) = \mathbf{x}(k) + \bar{\mathbf{u}}(k) + \Theta(k)(\mathbf{g}(k) - \bar{\mathbf{g}}(k)), \quad (14.20)$$

$$\mathbf{x}(k+1) = \mathbf{G}(k+1)\mathbf{x}^+(k), \quad k = 1, \dots, T-1, \quad (14.21)$$

令  $\Theta(0) \doteq \mathbf{0}$ 。反复应用公式(14.20)和(14.21), 我们得到瞬时时刻  $k = 1, \dots, T-1$  处投资组合的表达式:

$$\mathbf{x}(k) = \Phi(1, k)\mathbf{x}(0) + \Omega_k \bar{\mathbf{u}} + \sum_{t=1}^k \Phi(t, k)\Theta(t-1)\tilde{\mathbf{g}}(t-1), \quad (14.22)$$

其中  $\Theta(0) = \mathbf{0}$ ,

$$\bar{\mathbf{u}} \doteq [\bar{\mathbf{u}}(0)^\top \cdots \bar{\mathbf{u}}(T-2)^\top \bar{\mathbf{u}}(T-1)^\top]^\top,$$

$$\tilde{\mathbf{g}}(k) \doteq \mathbf{g}(k) - \bar{\mathbf{g}}(k), \quad k = 1, \dots, T.$$

一个关键的观察所得是  $\mathbf{x}(k)$  是决策变量  $\bar{\mathbf{u}}(k)$  和  $\Theta(k)$ ,  $k = 1, \dots, T-1$  的仿射函数。在整个过程中, 投资收益的累积总收益是

$$\begin{aligned} \rho(\bar{\mathbf{u}}, \Theta) &= \frac{w(T)}{w(0)} = \frac{\mathbf{1}^\top \mathbf{x}(T)}{\mathbf{1}^\top \mathbf{x}(0)} \\ &= \frac{1}{\mathbf{1}^\top \mathbf{x}(0)} (\phi(1, T)^\top \mathbf{x}(0) + \omega_T^\top \bar{\mathbf{u}} + \sum_{t=1}^T \Phi(t, T)\Theta(t-1)\tilde{\mathbf{g}}(t-1)), \end{aligned}$$

这也是一个关于变量  $\bar{\mathbf{u}}$  和  $\Theta \doteq [\Theta(1) \cdots \Theta(T-1)]$  的仿射函数。

### §14.3.3.2 使用仿射方法的最佳策略

给定从场景生成先哲那儿得到的各时段增益  $\{\mathbf{G}(k), k = 1, \dots, T\}$  的  $N$  个独立同分布样本(场景)  $\mathbf{G}^{(i)}(k), k = 1, \dots, T, i = 1, \dots, T$ 。则我们可以

将(14.15)-(14.18)最为最优化的目标函数，即求解如下凸优化问题以求确定最优政策：

$$\begin{aligned} J_{\text{cl}}^*(\gamma) = & \underset{\bar{\mathbf{u}}, \Theta}{\text{minimize}} && J(\mathbf{u}, \Theta) \\ \text{subject to} &&& \mathbf{x}^{(i)+}(k) \in \mathcal{X}(k), k = 0, \dots, T-1; i = 1, \dots, N, \\ &&& \mathbf{1}^\top \bar{\mathbf{u}}(k) = 0, k = 0, \dots, T-1, \\ &&& \mathbf{1}^\top \Theta(k) = 0, k = 1, \dots, T-1, \end{aligned}$$

其中 $\mathbf{x}^{(i)+}(k)$ 由(14.20)确定，但是其中的 $\mathbf{x}(k)$ 应该换成 $\mathbf{x}^{(i)}(k)$ ，即将第 $i$ 个采样场景代入(14.22)得到的与第 $i$ 个采样场景相匹配的投资组合。

#### §14.4 稀疏指标跟踪

我们再次考虑在例9.13中介绍的replicating(跟踪)指标问题。这个问题可以视为有约束的最小二乘问题，即

$$\begin{aligned} & \underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} && \|\mathbf{R}\mathbf{x} - \mathbf{y}\|_2^2 \\ \text{subject to} &&& \mathbf{1}^\top \mathbf{x} = 1, \\ &&& \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

其中 $\mathbf{R} \in \mathbb{R}^{T \times n}$ ，它的第 $i$  ( $i = 1, \dots, n$ )列数据表示组合中的资产 $i$ 在 $T$  ( $> n$ )个时期内的历史收益， $\mathbf{y} \in \mathbb{R}^T$ ，它包含了 $T$ 个时期投资者的参考收益。我们的目标是找到资产组合的比例向量 $\mathbf{x}$ ，使其对应的收益尽可能地匹配收益流的参考指标。然而，在实践中， $n$ 可能很大，并且用户更喜欢用组分资产中的一个小(small)的子集来replicating指标，因为资产越少越便于管理，并且需要的交易成本更少。因此，客户愿意放弃一些跟踪精度来换取解 $\mathbf{x}$ 的稀疏性。我们在前面的章节中提到，一种通常的提高解的稀疏性的作法是向目标中添加一个 $\ell_1$ 范数正则项，即考虑形如 $\|\mathbf{R}\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{x}\|_1$ 的目标。遗憾的是，这种方法在我们的背景下不能工作，因为这里的决策变量属于标准单纯形 $\mathcal{X} = \{\mathbf{x} : \mathbf{1}^\top \mathbf{x} = 1, \mathbf{x} \geq \mathbf{0}\}$ ，因此所有可行的 $\mathbf{x}$ 的 $\ell_1$ 范数是常值，均等于1。我们接下来介绍一个当变量被限制在单形 $\mathcal{X}$ 中取值时，可以获得稀疏解的另一个有效的方法。首先考虑理想的问题，即我们要求解

$$p^* = \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) + \lambda \|\mathbf{x}\|_0, \quad (14.23)$$

其中 $f(\mathbf{x}) \doteq \|\mathbf{R}\mathbf{x} - \mathbf{y}\|_2^2$ ， $\|\mathbf{x}\|_0$ 是 $\mathbf{x}$ 的基数， $\lambda \geq 0$ 是一个折衷参数。因为我们知道这是一个非凸的问题，难以求解，所以我们要寻找一种可以有效

求解的松弛方法。为此，我们研究

$$\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i| \leq \|\mathbf{x}\|_0 \max_{i=1,\dots,n} |x_i| \leq \|\mathbf{x}\|_0 \|\mathbf{x}\|_\infty.$$

对  $\mathbf{x} \in \mathcal{X}$ ，上面表达式左边部分等于1。因此我们得到

$$\mathbf{x} \in \mathcal{X} \Rightarrow \|\mathbf{x}\|_0 \geq \frac{1}{\|\mathbf{x}\|_\infty} = \frac{1}{\max_{i=1,\dots,n} x_i}.$$

所以，对  $\mathbf{x} \in \mathcal{X}$ ，

$$f(\mathbf{x}) + \lambda \|\mathbf{x}\|_0 \geq f(\mathbf{x}) + \lambda \frac{1}{\max_{i=1,\dots,n} x_i}.$$

因此，解决问题

$$p_\infty^* = \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) + \lambda \frac{1}{\max_{i=1,\dots,n} x_i} \quad (14.24)$$

后，我们就得到了原始最优目标值的下界  $p_\infty^*$ ，即  $p_\infty^* \leq p^*$ 。进一步，记(14.24)的最优解为  $\mathbf{x}_\infty^*$ ，则我们立即得到

$$f(\mathbf{x}_\infty^*) + \lambda \|\mathbf{x}_\infty^*\|_0 \geq p^* \geq p_\infty^*,$$

其中第一个不等式源于(14.23)。这个关系允许我们进行如下操作：如果我们获得了一个解  $\mathbf{x}_\infty^*$ ，那么去检查它的后验性(a posteriori)，即它的次优性水平，因为该关系表明“真正”的最优值  $p^*$  一定包含于区间  $[p_\infty^*, f(\mathbf{x}_\infty^*) + \lambda \|\mathbf{x}_\infty^*\|_0]$ 。

尚未解决的问题是如何得到问题(14.24)的解，因为这个问题关于  $\mathbf{x}$  是非凸的(函数  $\frac{1}{\max_i x_i}$  在  $\mathcal{X}$  上是凹函数)，所以它的解并不是唾手可得的。然而我们可以进行如下推导<sup>8</sup>

$$\begin{aligned} p_\infty^* &= \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) + \min_{i=1,\dots,n} \frac{\lambda}{x_i} \\ &= \min_{i=1,\dots,n} \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) + \frac{\lambda}{x_i} \\ &= \min_{i=1,\dots,n} \min\{f(\mathbf{x}) + t : \mathbf{x} \in \mathcal{X}, t \geq 0, tx_i \geq \lambda\}. \end{aligned}$$

<sup>8</sup>参见Pilanci, El Ghaoui, and Chandrasekaran, Recovery of sparse probability measures via convex programming, Proc. Conference on Neural Information Processing Systems, 2012.



进一步, 可以利用(10.2)将双曲型约束 $tx_i \geq \lambda$ 表示为SOC约束, 即

$$tx_i \geq \lambda, x_i \geq 0, t \geq 0 \Leftrightarrow \left\| \begin{bmatrix} \sqrt{\lambda} \\ x_i - t \end{bmatrix} \right\|_2 \leq x_i + t.$$

上述推导表明, 我们可以通过求解 $n$ 个凸优化问题, 即对 $i = 1, 2, \dots, n$ , 求解

$$\begin{aligned} & \underset{\mathbf{x}, t}{\text{minimize}} && f(\mathbf{x}) + t \\ & \text{subject to} && \mathbf{x} \in \mathcal{X}, \\ & && t \geq 0, \\ & && \left\| \begin{bmatrix} \sqrt{\lambda} \\ x_i - t \end{bmatrix} \right\|_2 \leq x_i + t. \end{aligned}$$

然后比较最优值, 取最优值最小的一个问题对应的解, 由此即获得非凸问题(14.24)的解. 下面用一个数值例子来说明, 我们再次使用例9.13中的数据, 它涉及到使用 $n = 5$ 的组分指标来跟踪摩根士丹利世界(MSCI WORLD)指标. 求解 $\lambda = 0.1$ 时的问题(14.24), 我们获得只有两个非零元素的投资组合

$$\mathbf{x}_\infty^* = [95.19, 0, 4.81, 0, 0]^\top \times 0.01,$$

与此对应的跟踪误差 $\|\mathbf{R}\mathbf{x}_\infty^* - \mathbf{y}\|_2^2 = 0.0206$ .

**注记14.1** 我们这里讨论的稀疏性方法, 除了应用到指标跟踪问题外, 显然可以被应用于本文提到的其它决策变量被限制在标准单纯形中的投资组合分配问题. 例如, 一旦施加了非卖空条件, 就可以用它来得到14.1.1节提到的均值方差最优化问题的稀疏解. 相反的, 当组合变量没有施加非负限制时(允许卖空), 可以使用标准的 $\ell_1$ 范数法进行松弛, 由此提升投资组合的稀疏性.

## §14.5 习题

**习题14.1. (投资组合最优化问题)** 我们考虑单时段的涉及 $n$ 种资产的最优化问题, 其中决策向量 $\mathbf{x} \in \mathbb{R}^n$ 包含了我们对每种资产的投资数量, 确定以下哪些目标或约束可以使用凸优化进行建模。

- (a) 风险水平(由投资组合的方差测度)等于一个给定的目标 $t$ (假设协方差矩阵是已知的)。
- (b) 风险水平(由投资组合的方差测度)给给定的目标 $t$ 之下。

- (c) 夏普比率(定义为投资组合收益与投资组合标准差的比值)在目标 $t \geq 0$ 之上。这里假设收益向量的期望和协方差矩阵都是已知的。
- (d) 假设收益向量服从已知的高斯分布, 确保投资组合收益低于目标 $t$ 的概率小于3%。
- (e) 假设收益向量 $\mathbf{r} \in \mathbb{R}^n$ 有三种可能取值 $\mathbf{r}^{(i)}, i = 1, 2, 3$ . 施加如下约束: 这三个场景下投资收益的最小值要在目标水平 $t$ 之上。
- (f) 在类似于(e)中的假设下: 两个最小投资组合的收益的平均值均在目标水平 $t$ 之上。提示: 利用新变量 $s_i = \mathbf{x}^\top \mathbf{r}^{(i)} (i = 1, 2, 3)$  并且考虑函数 $\mathbf{s} \rightarrow s_{[2]} + s_{[3]}$ , 其中 $s_{[k]}$ 表示 $\mathbf{s}$ 的第 $k$ 大元素,  $k = 1, 2, 3$ 。
- (g) 交易成本(在线性交易成本及初始投资组合 $\mathbf{x}(0) = \mathbf{0}$ 的模型中)低于某个目标。
- (h) 从初始投资组合 $\mathbf{x}(0) = \mathbf{0}$ 到最优投资组合 $\mathbf{x}$ 的交易数目低于某个目标。
- (i) 投资组合收益的期望值与目标收益 $t$ 之差的绝对值小于任意一个给定的正数 $\epsilon$ (在这里, 我们假设收益向量的期望 $\bar{\mathbf{r}}$ 是已知的)。
- (j) 投资组合收益的期望值或者在某一特定值 $t_{\text{up}}$ 之上, 或者在另一个特定值 $t_{\text{low}}$ 之下。

**习题14.2. (中位风险, Median risk)** 我们考虑一个单时段关于 $n$ 种资产的投资组合最优化问题。我们利用过去的样本, 即单时段收益向量 $\mathbf{r}_1, \dots, \mathbf{r}_N$ , 其中 $\mathbf{r}_t \in \mathbb{R}^n$  代表着对资产从时段 $t-1$ 到时段 $t$ 的投资组合。我们用

$$\hat{\mathbf{r}} \doteq \left(\frac{1}{N}\right)(\mathbf{r}_1 + \dots + \mathbf{r}_N)$$

来表示样本平均值向量, 这是一个基于过去的样本对期望收益进行估计的估计量。

我们用以下量来度量风险。我们用 $\rho_t(\mathbf{x})$ 来表示时刻 $t$ 的收益(如果 $t$ 时刻的投资组合是 $\mathbf{x}$ )。我们的风险测度是

$$\mathcal{R}_1(\mathbf{x}) \doteq \frac{1}{N} \sum_{t=1}^n |\rho_t(\mathbf{x}) - \hat{\rho}(\mathbf{x})|$$

其中 $\hat{\rho}(\mathbf{x})$ 是投资组合的样本平均收益。

- (a) 证明 $\mathcal{R}_1(\mathbf{x}) = \|\mathbf{R}^\top \mathbf{x}\|_1$ , 你需要确定出这里的 $\mathbf{R}$ , 它是一个 $n \times N$  矩阵。这个风险测度 $\mathcal{R}_1$ 是凸的吗?

- (b) 说明在投资组合收益率的样本平均值大于目标 $\mu$ 的条件下, 极小化风险测度 $R_1$ 可以用线性规划来表述. 将问题化成标准形式的线性规划, 并且精确地定义变量和约束。
- (c) 定量评价这种风险度量方式所得投资组合和使用更经典的, 即基于方差

$$R_2(\mathbf{x}) \doteq \frac{1}{N} \sum_{t=1}^n (\rho_t(\mathbf{x}) - \hat{\rho}(\mathbf{x}))^2$$

度量风险所得投资组合的差异。

#### 习题14.3. (带因素模型的投资组合优化-1)

- (a) 考虑下面的投资组合优化问题:

$$\begin{aligned} p^* = \underset{\mathbf{x}}{\text{minimize}} \quad & \mathbf{x}^\top \boldsymbol{\Sigma} \mathbf{x} \\ \text{subject to} \quad & \hat{\mathbf{r}}^\top \mathbf{x} \geq \mu, \end{aligned}$$

其中 $\hat{\mathbf{r}} \in \mathbb{R}^n$ 是收益向量的期望,  $\boldsymbol{\Sigma} \in \mathbb{S}^n$ ,  $\boldsymbol{\Sigma} \succeq \mathbf{0}$ 是收益的协方差矩阵, 并且 $\mu$ 是投资组合预期收益的一个目标值。假定随机收益向量 $\mathbf{r}$ 服从简化因素模型, 即

$$\mathbf{r} = \mathbf{F}(\mathbf{f} + \hat{\mathbf{f}}), \quad \hat{\mathbf{r}} \doteq \mathbf{F}\hat{\mathbf{f}},$$

其中 $\mathbf{F} \in \mathbb{R}^{n \times k}$  ( $k \ll n$ ) 是一个因素载荷矩阵,  $\hat{\mathbf{f}} \in \mathbb{R}^k$ 是给定的, 并且 $\mathbf{f} \in \mathbb{R}^k$ 满足 $\mathbb{E}\{\mathbf{f}\} = \mathbf{0}$  和  $\mathbb{E}\{\mathbf{f}\mathbf{f}^\top\} = \mathbf{I}$ . 上述优化问题是一个包含 $n$ 个决策变量的凸二次优化问题。解释如何将这个问题转化为等价的, 且只包含 $k$ 个决策变量的凸二次规划问题。从几何上解释既约后的问题. 给出问题的一个具有封闭形式的解。

- (b) 考虑(a)中问题的变形, 即

$$\begin{aligned} p^* = \underset{\mathbf{x}}{\text{minimize}} \quad & \mathbf{x}^\top \boldsymbol{\Sigma} \mathbf{x} - \gamma \hat{\mathbf{r}}^\top \mathbf{x} \\ \text{subject to} \quad & \mathbf{x} \geq \mathbf{0} \end{aligned}$$

其中 $\gamma > 0$ 是一个折衷参数, 用来权衡目标中的风险项和收益项的相关性。由于存在约束 $\mathbf{x} \geq \mathbf{0}$ , 从而这个问题没有封闭形式的通解。

假定 $\mathbf{r}$ 由如下形式的因素模型指定, 即

$$\mathbf{r} = \mathbf{F}(\mathbf{f} + \hat{\mathbf{f}}) + \mathbf{e},$$

其中  $F, f$  以及  $\hat{f}$  如(a)中所定义,  $e$  是一个与  $f$  不相关(即  $\mathbb{E}\{fe^\top\} = 0$ ) 的噪声项, 并且满足  $\mathbb{E}\{e\} = 0$  和  $\mathbb{E}\{ee^\top\} = D^2 \doteq \text{diag}([d_1^2, \dots, d_n^2]) \succ 0$ . 假设我们期望去用12.3.1节中讨论的对数障碍法来求解这个问题, 解释如何利用收益的因素结构提高算法的数值性能. **提示:** 添加恰当的松弛变量后, (加上障碍项的)目标的海森矩阵将是对角的.

**习题14.4. (带因素模型的投资组合优化-2)** 再次考虑练习14.3中的(b). 令  $z \doteq F^\top x$ , 并且验证(b)中的问题可以被重写成

$$\begin{aligned} p^* = \underset{x, z}{\text{minimize}} \quad & x^\top D^2 x + z^\top z - \gamma \hat{r}^\top x \\ \text{subject to} \quad & F^\top x = z, x \geq 0. \end{aligned}$$

考虑拉格朗日函数

$$L(x, z, \lambda) = x^\top D^2 x + z^\top z - \gamma \hat{r}^\top x + \lambda^\top (z - F^\top x)$$

以及对偶函数

$$g(\lambda) \doteq \min_{x \geq 0, z} L(x, z, \lambda).$$

因为原始问题是凸的并且是严格可行的, 所以强对偶成立, 因此  $p^* = d^* = \max_{\lambda} g(\lambda)$ .

(a) 找到对偶函数  $g(\lambda)$  的闭合形式的表达式。

(b) 用对偶最优变量  $\lambda^*$  来表示原始最优解  $x^*$ 。

(c) 给出  $-g(\lambda)$  的一个次梯度。

**习题14.5. (Kelly的投注策略)** 一个赌徒的原始的资本是  $W_0$ , 并且在一次赌局中将自己所有可得的资本下注, 赌局中赢赌注的概率  $p \in [0, 1]$ , 输掉赌注的概率是  $1 - p$ . 赌局重复进行. 在  $k$  次投注后, 他的资本  $W_k$  是一个随机变量:

$$W_k = \begin{cases} 2^k W_0 & \text{依概率 } p^k, \\ 0 & \text{以概率 } 1 - p^k. \end{cases}$$

(a) 确定  $k$  次投注后该赌徒资产的期望. 确定赌徒在某时刻  $k$  破产的概率。

(b) (a)中的结论描述的是毁灭性赌博策略. 假设现在赌徒变得跟谨慎, 即在每一步决将资本的一部分作为赌注. 用  $w$  和  $\ell$  分别表示赌徒赢和输的赌局的次数, 二者均是随机的. 则赌徒在  $k$  次投注后的资本

$$W_k = (1 + x)^w (1 - x)^\ell W_0,$$

其中  $x \in [0, 1]$  是投注比例, 并且  $w + \ell = k$ 。定义赌徒资本的指数增长率为

$$G = \lim_{k \rightarrow \infty} \frac{1}{k} \log_2 \frac{W_k}{W_0}.$$

- (i) 将指数增长率  $G$  表示成  $x$  的函数。这个函数是凹的吗?
- (ii) 求  $x \in [0, 1]$  使指数增长率  $G$  取到最大值。称这个  $x$  是最优的凯利投注策略。
- (c) 考虑一个更加一般的情况, 其中一个投资者可以在一个投资机会上面仅投资他的资金中的一部分, 并且该投资以不同的概率具有不同的收益。具体地, 如果投资了  $W_0 x$  美金, 那么投资导致资产  $W = (1 + rx)W_0$ , 其中  $r$  表示投资的收益, 是一个离散随机变量, 它以概率  $p_1, \dots, p_m$  取  $r_1, \dots, r_m$  ( $p_i \geq 0, r_i \geq -1, i = 1, \dots, m, \sum_{i=1}^m p_i = 1$ )。

对于这种情况, (b) 中引入的指数增长率就是投资的对数增益的期望值

$$G = \mathbb{E}\{\log(W/W_0)\} = \mathbb{E}\{\log(1 + rx)\}.$$

(b) 中考虑的问题是这种推广问题的特殊情况, 对应于  $m = 2, r_1 = 1, r_2 = -1, p_1 = p, p_2 = 1 - p$ 。

- (i) 给出指数增长函数  $G$  作为  $x \in [0, 1]$  的函数的显式表达式。
- (ii) 设计一个简单的计算方案以找到使  $G$  取最大值的最优的投资比例  $x$ 。

**习题14.6. (多期投资)** 我们在  $n$  个时段上考虑一个多时段单一资产的投资决策问题。对任意给定的时段  $i = 1, \dots, n$ , 我们用  $y_i$  来表示预期收益,  $\sigma_i$  表示对应的方差,  $u_i$  表示投资的资金数量。假定我们的初始投资是  $u_0 = w$ , 投资问题是

$$\phi(w) \doteq \max_{\mathbf{u}} \left\{ \sum_{i=1}^{n+1} (y_i u_i - \lambda \sigma_i^2 u_i^2 - c |u_i - u_{i-1}|) : u_0 = w, u_{n+1} = 0 \right\},$$

其中第一项代表利润, 第二项代表风险, 第三项是交易成本的近似。在这里,  $c > 0$  是单位交易成本,  $\lambda > 0$  是风险-收益折衷参数。(不失一般性我们假定  $\lambda = 1$ .)

- (a) 确定这个问题的对偶问题。
- (b) 证明函数  $\phi$  是凹的, 并且给出  $-\phi$  在  $w$  处的一个次梯度。如果函数  $\phi$  在  $w$  处可微, 它在  $w$  处的梯度是什么?

- (c) 函数 $\phi$ 关于初始投资 $w$ 的灵敏度问题是什么?确切地,对于固定的 $y, \sigma, c$ 和任意的 $\epsilon > 0$ , 为 $|\phi(w + \epsilon) - \phi(w)|$ 提供一个紧的上界. 你可以假设函数 $\phi$ 对任意的 $u \in [w, w + \epsilon]$ 是可微的。

**习题14.7. (个人理财问题)** 考虑如下的个人理财问题。在接下来的六个月内你将从事咨询工作, 你因此而得到总共 $C = 30000$ 美元的酬金。你打算用这笔报酬来清还信用卡在过去的债务, 总共 $D = 7000$ 美元。信用卡的APR(annual interest rate, 年利率) $r_1 = 15.95\%$ 。你要考虑下面的一些因素:

- 在每个月的一开始, 你可以将信用卡上债务的一部分转移到另一张年利率较低的信用卡上, 该张信用卡的年利率 $r_2 = 2.9\%$ 。这种转移需要付交易费, 交易费是总共转移数量的 $r_3 = 0.2\%$ 。并且只能将债务从卡 1 转到卡 2。
- 雇主允许你自己选择付款时间表: 你可以在至多六个月中分期付款。由于资金流动的原因, 雇主限制每个月的付款量不超过 $\frac{4}{3} \times (\frac{C}{6})$ 。
- 你每年的底薪 $B = 70000$ 美金。你不能用你的底薪来还信用卡上的债务。然而它会影响你要缴纳的税金(如下所述)。
- 这六个月中的前三个月是本财年的最后三个月, 后三个月是下一个财年的前三个月。所以, 你要是想在本财年(咨询的前三个月)支取的酬金多一些, 你就要交的税就很多。如果你在不同时期支取酬金的话, 你要交的税将会少一些。具体的税费取决于你的年收入总额 $G$ , 也就是你的底薪加上其它收入。边际税率表见表14.1。

表 14.1: 边际税收率表

所有总收入(\$)	边际税率	所有的税
$0 \leq G \leq 80,000$	10%	$10\% \times G$
$80,000 \leq G$	28%	$28\% \times (G - 80,000) + 8000$

- 无风险利率(储蓄利率)是0。
- 事件的时间列表: 所有事件发生在每个月的开始。也就是说, 在每个月的开始, 雇主支付给你想要的数量的酬金, 同时你要决定还信用卡上债务的数量, 并且决定从信用卡 1 转到信用卡 2 的债务的数量, 并进行转移。任何未偿还债务的利率都将累积到该月底记账。

- 你的目标是在两个财政年结束时，还清信用卡上的所有债务，并且使得总资产最大。

- (a) 将决策问题表述成一个优化问题。要保证精确地定义变量和约束。为了描述税收，我们使用下面的约束：

$$T_i = 0.1 \min(G_i, \alpha) + 0.28 \max(G_i - \alpha, 0) \quad (14.25)$$

其中 $T_i$ 和 $G_i$ 分别是第 $i$ 个财年要缴的税金和总收入， $i = 1, 2$ ， $\alpha = 80000$ 是税收阈值参数。

- (b) 这个问题是一个线性规划问题吗？请解释。
- (c) 当 $\alpha$ 和 $G_i$ 满足什么条件时，税收约束可以由下面的一组约束所代替？这是我们所要考虑的情况吗？在你的问题中可以用(14.26)来代替(14.25)吗？请解释。

$$\begin{aligned} T_i &= 0.1d_{1,i} + 0.28d_{2,i}, \\ d_{2,i} &\geq G_i - \alpha, \\ d_{2,i} &\geq 0, \\ d_{1,i} &\geq G_i - d_{2,i}, \\ d_{1,i} &\geq d_{2,i} - \alpha. \end{aligned} \quad (14.26)$$

- (d) 这个问题的含(14.26)的新表述是凸的吗？解释你的答案。
- (e) 用你喜爱的求解器求解这个问题。写出最优策略：收到酬金的方式(时间和数量)，支付/转移信用卡债务的时间表(时间和数量)，以及两年末最优总资产。最大资产 $W$ 是多少？
- (f) 对于 $\alpha \in [70k, 90k]$ ，计算最优的 $W$ ，并且画出 $W$ 作为 $\alpha$ 的函数在这个范围的图形。你能解释这个图形吗？

**习题14.8. (交易成本和市场影响)** 我们考虑下面的投资组合优化问题：

$$\max_{\mathbf{x}} \{ \hat{\mathbf{r}}^\top \mathbf{x} - \lambda \mathbf{x}^\top \mathbf{C} \mathbf{x} - c \cdot T(\mathbf{x} - \mathbf{x}^0) : \mathbf{x} \geq \mathbf{0}, \mathbf{x} \in \mathcal{X} \} \quad (14.27)$$

其中 $\mathbf{C}$ 是经验协方差矩阵， $\lambda > 0$ 是风险参数， $\hat{\mathbf{r}}$ 表示在给定时段上每项资产收益的时间-平均。这里的约束集 $\mathcal{X}$ 由下面的条件所决定。

- 不允许卖空。

- 存在预算约束  $x_1 + \cdots + x_n = 1$ 。

在上面,  $T$  是函数, 代表着交易费以及市场影响,  $c \geq 0$  是控制这些费用大小的参数,  $\mathbf{x}^0 \in \mathbb{R}^n$  是表示初始投资组合的向量。函数  $T$  形如

$$T(\mathbf{x}) = \sum_{i=1}^n B_M(x_i),$$

其中函数  $B_M$  对于较小的  $x$  来说是分段线性的, 对于较大的  $x$  来说是二次的; 这样, 我们寻求捕获事实: 对于小宗交易(较小的  $x$ ), 交易成本是占主导地位的; 同时对于那些大宗交易, 市场影响占主导地位。确切地说, 我们定义一个称之为“reverse Hüber”的截止参数为  $M$  的函数: 对于一个标量  $z$ , 函数值

$$B_M(z) \doteq \begin{cases} |z| & \text{当 } |z| \leq M, \\ \frac{z^2 + M^2}{2M} & \text{否则.} \end{cases}$$

标量  $M > 0$  描述直线形状的惩罚转变成二次形状的惩罚时的转变点。

- (a) 说明  $B_M$  可以被表述为如下优化问题的解, 即

$$B_M(z) = \min_{v, w \in \mathbb{R}} \left\{ v + w + \frac{w^2}{2M} : |z| \leq v + w, v \leq M, w \geq 0 \right\}.$$

上面的表述可以证明  $B_M$  是凸的. 请给出理由?

- (b) 证明对于给定的  $\mathbf{x} \in \mathbb{R}^n$ :

$$T(\mathbf{x}) = \min_{\mathbf{w}, \mathbf{v} \in \mathbb{R}^n} \left\{ \mathbf{1}^\top (\mathbf{v} + \mathbf{w}) + \frac{1}{2M} \mathbf{w}^\top \mathbf{w} : \mathbf{v} \leq M\mathbf{1}, \mathbf{w} \geq \mathbf{0}, |\mathbf{x} - \mathbf{x}^0| \leq \mathbf{v} + \mathbf{w} \right\},$$

其中  $\mathbf{w}, \mathbf{v}$  现在表示  $n$  维向量,  $\mathbf{1}$  是分量全为 1 的向量, 并且其中的  $|\cdot|$  和不等式都是逐分量意义的。

- (c) 把最优化问题(14.27)表述成凸优化的形式。转化后的问题隶属于我们熟悉的那一类(LP, QP, SOCP 等)?





## 第十五章 控制问题

**动力系统**是时下热门的物理问题。典型动力系统的数学模型是常微分方程，其中指令与干扰信号是输入量，它们与已有的初始条件一起，决定了内在变量（状态）以及输出信号随时间如何演化。这种模型泛见于工程学，例如，用来描述飞机的飞行状态，内燃机和机械臂各自的运行情况，以及导弹的弹道等。

宽泛地讲，一个动力系统的**控制问题**就是通过确定适当的输入信号使系统向我们想要的方向演进，比如，遵循一个既定的轨道或在外界扰动下保持稳定等等。动力系统控制问题的入门介绍都需要一整本教科书。这里，我们仅考虑一类受限动力系统（即有限维线性定常系统）的几个具体方面。

我们首先介绍连续时间模型及其对应的离散模型。对离散时间模型来说，我们着重于在有限阶段下输入输出量的表现与由线性方程组描述的静态线性映射之间的联系。我们将看到，在这种情况下会自然而然地会涉及到最优化问题，并讨论这些问题在控制背景下的诠释。本章的第一部分处理所谓的基于最优化的控制问题(在这种问题中，通过求解一个有限阶段上的最优化问题来直接确定控制输入)及其在滑动阶段的迭代实现(此即所谓的模型预测控制(MPC)的范例)。本章的第二部分比较简略，我们反过来讨论一种更经典的控制方法，其基于诸如稳定性等这样的基于无限时间范畴的概念。在这里，最优化方法(尤其是SDP)作为稳定性分析或设计稳定反馈控制器的工具，它是以间接的方式出现在问题中的。

### §15.1 连续及离散时间模型

#### §15.1.1 连续时间LTI 系统

先来考虑一个简单的例子。

**例15.1. (轨上小车)**如图15.1所示，考虑沿着一条水平轨道移动的质量为 $m$ 的小车，这里小车受阻于阻尼系数为 $\beta$ 的粘滞阻尼(即阻尼与速度成正比例)。记小车质心的位置为 $p(t)$ ，记质心受到的力为 $u(t)$ ，其中 $t$ 是时间。

根据牛顿第二定律，有

$$u(t) - \beta \dot{p}(t) = m\ddot{p}(t),$$

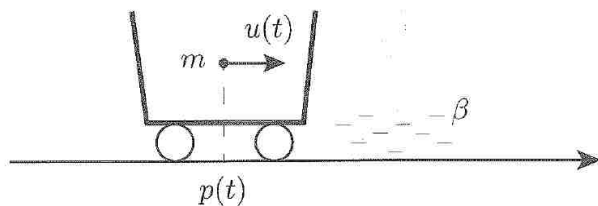


图 15.1: 轨上小车

这是一个二阶微分方程，它描述了这个系统的动力学。如果我们引入(状态)变量  $x_1(t) = p(t)$ ,  $x_2(t) = \dot{p}(t)$ ，则我们可以把牛顿方程重写为一阶常微分方程组：

$$\begin{aligned}\dot{x}_1(t) &= x_2(t), \\ \dot{x}_2(t) &= \alpha x_2(t) + bu(t),\end{aligned}$$

其中  $\alpha \doteq -\frac{\beta}{m}$ ,  $b \doteq \frac{1}{m}$ . 这个系统可以用更紧凑的矩阵形式表示为

$$\dot{\mathbf{x}}(t) = \mathbf{A}_c \mathbf{x}(t) + \mathbf{B}_c u(t), \quad (15.1)$$

其中

$$\mathbf{A}_c = \begin{bmatrix} 0 & 1 \\ 0 & \alpha \end{bmatrix}, \quad \mathbf{B}_c = \begin{bmatrix} 0 \\ b \end{bmatrix}.$$

我们也可以给这个系统联立一个输出方程

$$y(t) = \mathbf{C} \mathbf{x}(t), \quad \mathbf{C} = \begin{bmatrix} 1 & 0 \end{bmatrix} \quad (15.2)$$

以表示信号  $y$  是我们特别感兴趣的，这里表示小车自己的位置是我们特别关注的。

方程(15.1)和(15.2)实际上表述了一类相当有趣的动力系统，即所谓的(严格真(strictly proper)，有限维)连续时间的线性时不变(Linear Time-Invariant, LTI) 系统。如此的一个LTI 系统可以用  $(\mathbf{A}, \mathbf{B}, \mathbf{C})$  这样的三元矩阵对，以状态方程

$$\dot{\mathbf{x}} = \mathbf{A} \mathbf{x} + \mathbf{B} \mathbf{u} \quad (15.3)$$

$$\mathbf{y} = \mathbf{C} \mathbf{x} \quad (15.4)$$

来定义，这里  $\mathbf{x} \in \mathbb{R}^n$  是状态向量， $\mathbf{u} \in \mathbb{R}^m$  是输入向量， $\mathbf{y} \in \mathbb{R}^p$  是输出向量( $\mathbf{x}, \mathbf{u}, \mathbf{y}$  都是时间  $t \in \mathbb{R}$  的函数，但为了简洁，我们省略对  $t$  的显式

依赖性)。给定 $t_0$ 时刻的状态的值, 以及 $t \geq t_0$ 时的输入向量 $\mathbf{u}(t)$ , 则系统(15.3)的状态随时间的演化可以显式表示为(一般称之为拉格朗日公式)

$$\mathbf{x}(t) = e^{\mathbf{A}(t-t_0)}\mathbf{x}(t_0) + \int_{t_0}^t e^{\mathbf{A}(t-\tau)}\mathbf{B}\mathbf{u}(\tau)d\tau, \quad t \geq t_0. \quad (15.5)$$

### §15.1.2 离散时间LTI系统

对于分析特定的动力学现象, 例如人口动力学或经济学, 在离散的时间间隔上描述系统可能更加符合实际情况。例如, 如果我们想要对一个国家的经济进行建模, 当我们考虑要纳入模型的相关量(比如国内生产总值, 或者失业率)时, 我们会将离散的时间间隔 $\Delta$  ( $\Delta$ 为一周, 或一个月)处的数据, 而非连续时间下的相关数据。通常用一个带符号的整型变量 $k, k = \dots, -1, 0, 1, 2, \dots$ , 来表示“离散”的时间, 这里 $k$ 代表时刻 $t = k\Delta$ 。(严格真, 有限维的)离散时间线性时不变(LTI)系统用一阶差分方程组表示为

$$\mathbf{x}(k+1) = \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k), \quad (15.6)$$

$$\mathbf{y}(k) = \mathbf{C}\mathbf{x}(k). \quad (15.7)$$

给定系统在时刻 $k_0$ 的状态, 以及 $k \geq k_0$ 的输入 $\mathbf{u}(k)$ , 通过递归应用方程(15.6), 易验证有

$$\mathbf{x}(k) = \mathbf{A}^{k-k_0}\mathbf{x}(k_0) + \sum_{i=k_0}^{k-1} \mathbf{A}^{k-i-1}\mathbf{B}\mathbf{u}(i), \quad k \geq k_0. \quad (15.8)$$

#### §15.1.2.1 离散化

在实践中, 常常会出现一个连续时间动力系统必须被数字器件(如DSP和计算机)分析或控制的情况。数字器件本身是离散时间对象, 只能在离散时刻 $k\Delta$ 与外界交互, 这里的 $\Delta$ 是一个很小的时间间隔, 代表I/O的时钟频率。因此很自然地, 我们需要把一个形如(15.3)和(15.4)的连续时间系统“转化”为对应的离散时间的形式, 通过在 $t = k\Delta$ 这些时刻“给系统拍照”, 这里的 $\Delta$ 称为采样间隔(sampling interval), 并假定输入信号 $\mathbf{u}(t)$ 在两个相继的取样时刻之间是不变的, 即

$$\mathbf{u}(t) = \mathbf{u}(k\Delta), \quad \forall t \in [k\Delta, (k+1)\Delta).$$

这样的离散时间转换可以根据如下步骤完成: 给定系统(15.3)在时刻 $t = k\Delta$  (这里我们把 $\mathbf{x}(k\Delta)$ 记为 $\mathbf{x}(k)$ )的状态, 我们可以用方程(15.5)来计算

该状态在时刻 $(k+1)\Delta$ 的值:

$$\begin{aligned}\mathbf{x}(k+1) &= e^{\mathbf{A}\Delta}\mathbf{x}(k) + \int_{k\Delta}^{(k+1)\Delta} e^{\mathbf{A}((k+1)\Delta-\tau)}\mathbf{B}\mathbf{u}(\tau)d\tau \\ &= e^{\mathbf{A}\Delta}\mathbf{x}(k) + \int_{k\Delta}^{(k+1)\Delta} e^{\mathbf{A}((k+1)\Delta-\tau)}\mathbf{B}d\tau\mathbf{u}(k) \\ &= e^{\mathbf{A}\Delta}\mathbf{x}(k) + \int_0^\Delta e^{\mathbf{A}\tau}\mathbf{B}d\tau\mathbf{u}(k)\end{aligned}$$

对于连续时间系统(15.3)和(15.4)的采样版本, 这意味着根据形如(15.6), (15.7)的离散时间递归时, 有

$$\begin{aligned}\mathbf{x}(k+1) &= \mathbf{A}_\Delta\mathbf{x}(k) + \mathbf{B}_\Delta\mathbf{u}(k), \\ \mathbf{y}(k) &= \mathbf{C}\mathbf{x}(k),\end{aligned}$$

这里

$$\mathbf{A}_\Delta = e^{\mathbf{A}\Delta}, \mathbf{B}_\Delta = \int_0^\Delta e^{\mathbf{A}\tau}\mathbf{B}d\tau. \quad (15.9)$$

**例15.2.** 考虑例15.1中轨上小车问题的连续时间模型, 并假定 $m = 1\text{kg}$ ,  $\beta = 0.1\text{Ns/m}$ . 采样间隔 $\Delta = 0.1\text{s}$ , 对系统(15.1)离散化, 我们得到离散时间系统

$$\mathbf{x}(k+1) = \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k), \quad (15.10)$$

$$\mathbf{y}(k) = \mathbf{C}\mathbf{x}(k), \quad (15.11)$$

这里 $\mathbf{A}, \mathbf{B}$ 由(15.9)按如下方式计算得到. 通过直接计算, 我们首先观察到

$$\mathbf{A}_c^2 = \mathbf{A}_c\mathbf{A}_c = \begin{bmatrix} 0 & \alpha \\ 0 & \alpha^2 \end{bmatrix}, \mathbf{A}_c^3 = \mathbf{A}_c^2\mathbf{A}_c = \begin{bmatrix} 0 & \alpha^2 \\ 0 & \alpha^3 \end{bmatrix}, \dots,$$

$$\mathbf{A}_c^k = \mathbf{A}_c^{k-1}\mathbf{A}_c = \begin{bmatrix} 0 & \alpha^{k-1} \\ 0 & \alpha^k \end{bmatrix},$$

因此

$$\begin{aligned}e^{\mathbf{A}_c t} &= \mathbf{I} + \sum_{k=1}^{\infty} \frac{1}{k!} \mathbf{A}_c^k t^k = \mathbf{I} + \sum_{k=1}^{\infty} \frac{1}{k!} \begin{bmatrix} 0 & \alpha^{k-1} \\ 0 & \alpha^k \end{bmatrix} t^k \\ &= \begin{bmatrix} 1 & \sum_{k=1}^{\infty} \frac{1}{k!} \alpha^{k-1} t^k \\ 0 & 1 + \sum_{k=1}^{\infty} \frac{1}{k!} \alpha^k t^k \end{bmatrix} \\ &= \begin{bmatrix} 1 & \frac{1}{\alpha}(e^{\alpha t} - 1) \\ 0 & e^{\alpha t} \end{bmatrix}\end{aligned}$$

故并且

$$\mathbf{A} = e^{\mathbf{A}_c \Delta} = \begin{bmatrix} 1 & \frac{1}{\alpha}(e^{\alpha\Delta} - 1) \\ 0 & e^{\alpha\Delta} \end{bmatrix} = \begin{bmatrix} 1 & 0.0995017 \\ 0 & 0.9900498 \end{bmatrix}.$$

类似地, 对 $\mathbf{B}$ 我们有

$$\begin{aligned} \mathbf{B} &= \int_0^\Delta e^{\mathbf{A}_c \tau} \mathbf{B}_c d\tau = \int_0^\Delta \begin{bmatrix} 1 & \frac{1}{\alpha}(e^{\alpha\tau} - 1) \\ 0 & e^{\alpha\tau} \end{bmatrix} \begin{bmatrix} 0 \\ b \end{bmatrix} d\tau \\ &= b \int_0^\Delta \begin{bmatrix} \frac{1}{\alpha}(e^{\alpha\tau} - 1) \\ e^{\alpha\tau} \end{bmatrix} d\tau \\ &= \frac{b}{\alpha} \begin{bmatrix} \frac{1}{\alpha}(e^{\alpha\Delta} - 1) - \Delta \\ e^{\alpha\Delta} - 1 \end{bmatrix} = \begin{bmatrix} 0.0049834 \\ 0.0995016 \end{bmatrix} \end{aligned}$$

## §15.2 基于最优化的控制综合问题

### §15.2.1 状态跟踪的控制命令的综合

在这个章节中, 我们集中讨论形如(15.6)和(15.7)的一个通用的离散时间系统, 其有一个标量输入 $u(k)$ 和一个标量输出 $y(k)$ (称这样的系统为SISO, 其代表单输入单输出(Single-Input Single-Output)). 这样的模型可能来自某系统的一个原始的离散时间表示, 也可能是原始的连续时间系统的一个离散版本。我们讨论以下的控制问题: 给定一个初始状态 $\mathbf{x}(0) = \mathbf{x}_0$ , 一个目标状态 $x_T$ 和一个目标整数时刻 $T > 0$ , 确定对系统的一系列控制动作 $u(k)$ ,  $k = 0, \dots, T-1$ 使得在 $T$ 时刻时系统状态为 $x_T$ 。也就是说, 我们寻找一个可以使系统状态在任意设定的时刻达到目标状态的控制命令列。

利用式(15.8)来表述这个问题, 得到

$$\begin{aligned}
 \mathbf{x}(T) &= \mathbf{A}^\top \mathbf{x}_0 + \sum_{i=0}^{T-1} \mathbf{A}^{T-i-1} \mathbf{B} u(i) \\
 &= \mathbf{A}^\top \mathbf{x}_0 + \begin{bmatrix} \mathbf{A}^{T-1} \mathbf{B} & \mathbf{A}^{T-2} \mathbf{B} & \dots & \mathbf{A} \mathbf{B} & \mathbf{B} \end{bmatrix} \begin{bmatrix} u(0) \\ u(1) \\ \vdots \\ u(T-2) \\ u(T-1) \end{bmatrix} \\
 &= \mathbf{A}^\top \mathbf{x}_0 + \mathbf{R}_T \boldsymbol{\mu}_T,
 \end{aligned}$$

其中我们定义了 $T$ -可达矩阵( $T$ -reachability matrix), 即 $\mathbf{R}_T \doteq \begin{bmatrix} \mathbf{A}^{T-1} \mathbf{B} & \dots & \mathbf{A} \mathbf{B} & \mathbf{B} \end{bmatrix}$ , 和包含了控制作用列的向量 $\boldsymbol{\mu}_T$ 。因此原控制问题等价于寻找一个向量 $\boldsymbol{\mu}_T \in \mathbb{R}^T$ 满足

$$\mathbf{R}_T \boldsymbol{\mu}_T = \boldsymbol{\xi}_T, \quad \boldsymbol{\xi}_T \doteq \mathbf{x}_T - \mathbf{A}^\top \mathbf{x}_0, \quad (15.12)$$

其中 $\mathbf{R}_T \in \mathbb{R}^{n \times T}$ 是系统的 $T$ -可达矩阵。式(15.12)非常有趣, 它告诉我们离散时间LTI系统在 $T$ 时刻的状态是一个时间从 0 到 $T-1$ 的输入命令列的线性方程, 如图15.2所示。

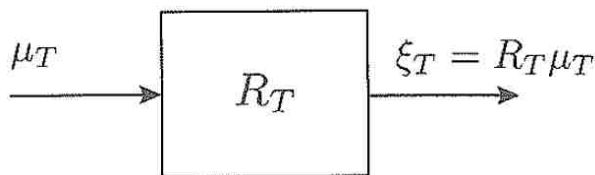


图 15.2: 从命令序列到状态的线性映射

决定一个使系统达到目标状态的输入序列相当于解一个关于 $\boldsymbol{\mu}_T$ 的形如(15.12)的线性方程组。我们假设 $T \geq n$ , 且 $\mathbf{R}_T$ 是行满秩的(用控制工程语言描述, 这个系统是完全 $T$ -可达的)。然后, 式(15.12)总是有一个解; 同时, 如果 $T > n$ , 它会有无穷多的可行解。这意味着这里有无限个可行的输入序列使得系统能达到指定的目标状态。在所有的可能性中, 很自然地会选择一个用“最小的气力”达到目标状态的控制序列, 其中“气力”(effort)可以定义为包含这些输入序列的向量 $\boldsymbol{\mu}_T$ 的某种范数。

## §15.2.1.1 最小能量控制

举例来说, 有限时间间隔 $\{0, \dots, T-1\}$ 上的输入信号 $\{u(0), \dots, u(T-1)\}$ 的能量(energy)定义为

$$\|\boldsymbol{\mu}_T\|_2^2 = \sum_{k=0}^{T-1} |u(k)|^2.$$

因此找能达到理想目标的能量最小的命令序列, 就相当于寻找一个线性方程组的2-范数最小的解(见§6.3.2), 即

$$\begin{aligned} & \underset{\boldsymbol{\mu}_T}{\text{minimize}} && \|\boldsymbol{\mu}_T\|_2^2 \\ & \text{subject to} && \mathbf{R}_T \boldsymbol{\mu}_T = \boldsymbol{\xi}_T. \end{aligned}$$

此系统有显式解(在完全可达性的假设下)

$$\boldsymbol{\mu}_T^* = \mathbf{R}_T^+ \boldsymbol{\xi}_T = \mathbf{R}_T^\top (\mathbf{R}_T \mathbf{R}_T^\top)^{-1} \boldsymbol{\xi}_T, \quad (15.13)$$

称其中的矩阵 $\mathbf{G}_T \doteq \mathbf{R}_T \mathbf{R}_T^\top$ 为系统的 $T$ 可达Gram矩阵。

与最小能量控制问题相关的另一个问题是确定用单位能量的输入序列时, 时刻 $T$ 可以达到的所有可能状态的集合, 即设法描述集合

$$\mathcal{E}_T = \{\mathbf{x} = \mathbf{R}_T \boldsymbol{\mu}_T + \mathbf{A}^\top \mathbf{x}_0 : \|\boldsymbol{\mu}_T\|_2 \leq 1\}.$$

上述集合是单位球 $\{\boldsymbol{\mu}_T : \|\boldsymbol{\mu}_T\|_2 \leq 1\}$ 在由矩阵 $\mathbf{R}_T$ 描述的线性映射下的像集加上一个常数偏差项为 $\mathbf{A}^\top \mathbf{x}_0$ 的平移。因此, 根据引理6.4, 我们推出单位能量的输入的可达状态集是以 $\mathbf{c} = \mathbf{A}^\top \mathbf{x}_0$ 为中心, 以 $\mathbf{G}_T = \mathbf{R}_T \mathbf{R}_T^\top$ 为构型矩阵的椭圆, 即

$$\mathcal{E}_T = \{\mathbf{x} : (\mathbf{x} - \mathbf{c})^\top \mathbf{G}_T^{-1} (\mathbf{x} - \mathbf{c}) \leq 1\}.$$

## §15.2.1.2 最小燃料控制

确定控制序列的另一种方法是最小化命令序列的1-范数而不是2-范数。 $\boldsymbol{\mu}_T$ 的1-范数与产生输入命令的所需的“燃料”消耗成正比例。比如, 在航天领域的应用中, 控制的输入通常是指推进器通过喷射压缩气体而产生的推力, 因此 $\|\boldsymbol{\mu}_T\|_1$ 与控制驱动所需气体的总质量成正比例。于是问题现在变成

$$\begin{aligned} & \underset{\boldsymbol{\mu}_T}{\text{minimize}} && \|\boldsymbol{\mu}_T\|_1 \\ & \text{subject to} && \mathbf{R}_T \boldsymbol{\mu}_T = \boldsymbol{\xi}_T, \end{aligned} \quad (15.14)$$



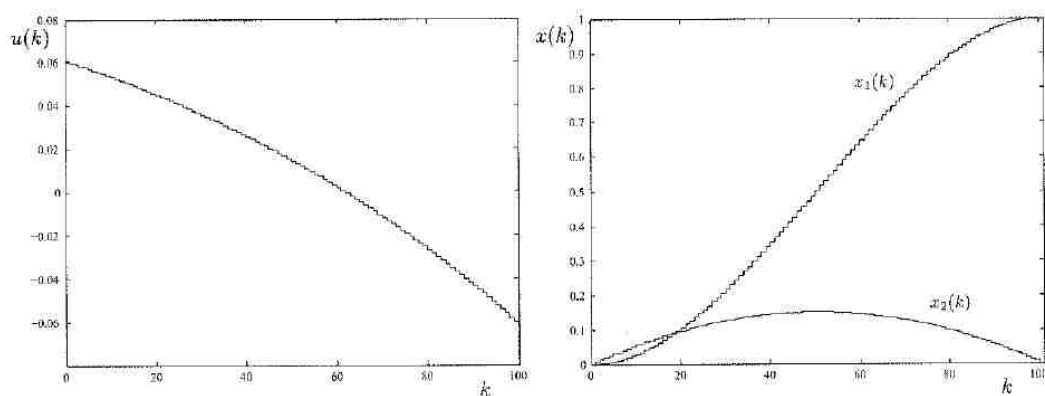


图 15.3: 最小能量控制信号(左)和相应的状态轨道(右)

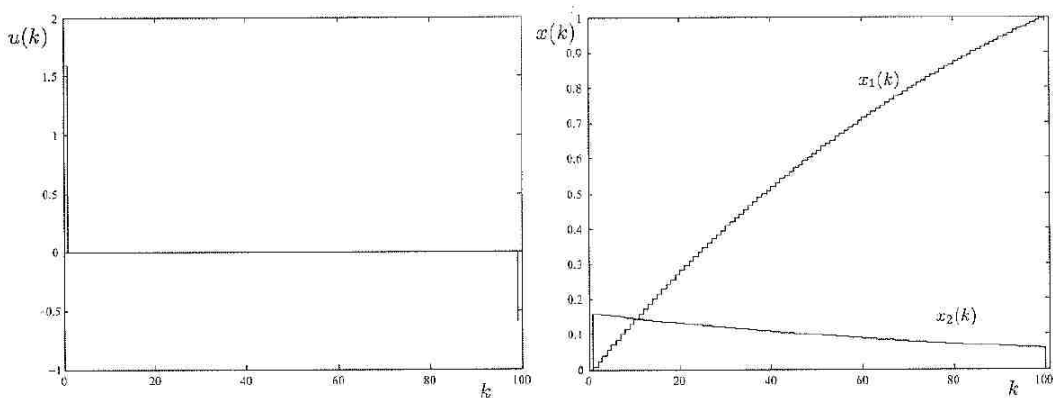


图 15.4: 最小燃料控制信号(左)和相应的状态轨道(右)

这可被重写成一个等价的线性规划

$$\begin{aligned}
 & \underset{\mu_T, s}{\text{minimize}} && \sum_{k=0}^{T-1} s_k \\
 & \text{subject to} && |u(k)| \leq s_k, k = 0, \dots, T-1, \\
 & && \mathbf{R}_T \mu_T = \xi_T.
 \end{aligned}$$

**例15.3.** 考虑例15.2中轨上小车的离散模型。给定初始条件  $\mathbf{x}(0) = [0 \ 0]^\top$ ，我们寻找最小能量和最小燃料的输入序列使得小车在  $t = 10\text{s}$  时到达位置  $p(t) = 1$ ，同时速度为0。因为采样时间是  $\Delta = 0.1\text{s}$ ，从而知道最终的整

型目标时刻是 $T = 100s$ ，因此我们的输入序列有100个未知量，即

$$\boldsymbol{\mu}_T = [u(0), \dots, u(T-1)]^\top.$$

由式(15.13)可以找到最小能量控制序列，其结果如图15.3所示。最小燃料控制序列可以通过求解与(15.14)等价的LP问题得到，其结果如图15.4所示。从定性的角度来看，我们观察到从最小能量方法和最小燃料方法得到的输入信号的形状非常不同。尤其是最小燃料方法的解是**稀疏**(sparse)的，意味着在这种情况下，除了开始和结束的时刻，控制作用处处为0（在控制术语中，这叫做**开关**(bang-bang)控制序列）。

### §15.2.2 轨迹追踪控制命令综合

在§15.2.1中我们讨论了如何确定控制序列以达到一个目标状态(target state)的问题，但是我们并没有特别地关注初始状态和目标状态之间的状态的运动轨迹。在这里，我们研究如何确定控制序列，其使得一个离散时间LTI系统的输出 $y(k)$ 在给定的有限时间阶段 $k \in \{1, \dots, T\}$ 上尽可能地接近指定的参考轨道 $y_{\text{ref}}(k)$ 。我们接下来假设这个系统是SISO，且 $\boldsymbol{x}(0) = \mathbf{0}, y_{\text{ref}}(0) = 0$ 。

根据方程(15.8)和 $\boldsymbol{x}_0 = \mathbf{0}$ ，我们得到

$$\boldsymbol{x}(k) = \boldsymbol{A}^{k-1} \boldsymbol{B} u(0) + \dots + \boldsymbol{A} \boldsymbol{B} u(k-2) + \boldsymbol{B} u(k-1).$$

然后，考虑输出方程 $y(k) = \boldsymbol{C} \boldsymbol{x}(k)$ 我们得到

$$y(k) = \boldsymbol{C} \boldsymbol{A}^{k-1} \boldsymbol{B} u(0) + \dots + \boldsymbol{C} \boldsymbol{A} \boldsymbol{B} u(k-2) + \boldsymbol{C} \boldsymbol{B} u(k-1), k = 1, \dots, T.$$

将上式重写为矩阵形式，我们得到

$$\begin{bmatrix} y(1) \\ y(2) \\ \vdots \\ y(T) \end{bmatrix} = \begin{bmatrix} \boldsymbol{C} \boldsymbol{B} & 0 & \dots & 0 \\ \boldsymbol{C} \boldsymbol{A} \boldsymbol{B} & \boldsymbol{C} \boldsymbol{B} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ \boldsymbol{C} \boldsymbol{A}^{T-1} \boldsymbol{B} & \dots & \boldsymbol{C} \boldsymbol{A} \boldsymbol{B} & \boldsymbol{C} \boldsymbol{B} \end{bmatrix} \begin{bmatrix} u(0) \\ u(1) \\ \vdots \\ u(T-1) \end{bmatrix}$$

由此，我们显然可以定义输出序列 $\boldsymbol{y}_T$ 和“变换矩阵” $\boldsymbol{\Phi}_T \in \mathbb{R}^{T \times T}$ 满足

$$\boldsymbol{y}_T = \boldsymbol{\Phi}_T \boldsymbol{\mu}_T. \quad (15.15)$$

矩阵 $\boldsymbol{\Phi}_T$ 具有特殊结构(对角线元素是常数)，称为Toeplitz结构；此外， $\boldsymbol{\Phi}_T$ 是下三角矩阵，因此当 $\boldsymbol{C} \boldsymbol{B} \neq 0$ 时， $\boldsymbol{\Phi}_T$ 是可逆的。更进一步，称 $\boldsymbol{\Phi}_T$ 的第一

列的元素为系统的“脉冲响应”，因为它们代表了当输入信号是 $u(0) = 1$ 和 $u(k) = 0$ (所有 $k \geq 1$ )这样的离散脉冲时，对应的系统的输出值序列。因此我们可以观察到，在一个固定的有限的时域中，一个离散时间LTI系统的输入输出行为可以用变换矩阵 $\Phi_T$ 决定的线性映射(15.15)来描述。现在，如果想要得到的输出序列被指定为某一参考序列 $y_{\text{ref}}(k)$ ，其中 $k \in \{1, \dots, T\}$ ，我们可以由式(15.15)解出输入序列 $\mu_T$ ，其能产生目标输出序列(当 $BC \neq 0$ 时，这样的解存在且唯一)。定义 $\mathbf{y}_{\text{ref}} = [y_{\text{ref}}(1) \cdots y_{\text{ref}}(T)]^\top$ ，这也等价于找到 $\mu_T$ 满足

$$\min_{\mu_T} \|\Phi_T \mu_T - \mathbf{y}_{\text{ref}}\|_2^2. \quad (15.16)$$

这个最小二乘表述的好处是当(15.15)是奇异的时，也可以为我们提供一个输入序列。在这种情况下( $\Phi_T$ 奇异)下，我们找到这样一个输入序列 $\mu_T$ 使得对应的输出 $\mathbf{y}_T$ 是“封闭”的，虽然在最小二乘意义上与 $\mathbf{y}_{\text{ref}}$ 不完全相同。

把问题(15.16)从数学语言转换成一个控制设计“规范”(specification)，即要找到这样一个控制输入，它控制的输出与给定的参考输出 $y_{\text{ref}}(k)$ 的跟踪误差(tracking error)最小。然而，对于控制系统想要找到仅有这样的规范的控制律是很难的。比如说注意到问题(15.16)所有的关注点都在跟踪误差上，然而我们完全忽视了命令的活动力(command activity)，也即输入序列 $\mu_T$ 的行为表现。事实上，有些问题必须考虑，比如说最大输入幅度的限制，转换速率约束或者至少让输入信号的能量或燃料消耗保持在可控范围中。举例来说，即使当 $\Phi_T$ 是可逆的，考虑下列正则化的轨迹追踪问题也是很好的实践，即

$$\min_{\mu_T} \|\Phi_T \mu_T - \mathbf{y}_{\text{ref}}\|_2^2 + \gamma \|\mu_T\|_2^2. \quad (15.17)$$

在这个问题里面我们引入了一个与输入信号能量成正比比例的惩罚项。这样的表述保证了追踪精度和输入能量间的一个权衡：典型的情况是为了达到一个小的跟踪精度需要一个高能量的输入信号。如果计算所得到的输入(当 $\gamma = 0$ 时得到的)的能量过高，我们就需要尝试更大的 $\gamma > 0$ 并再次求解这个问题，以此类推，直到找到一个令人满意的折中方案。

当然，以上问题里面可以有很多的变化。比如说，若想将1-范数纳入考量，那么对应的正则化问题是

$$\min_{\mu_T} \|\Phi_T \mu_T - \mathbf{y}_{\text{ref}}\|_2^2 + \gamma \|\mu_T\|_1,$$

或者想要给(15.17)添加关于控制信号的振幅的显式约束, 即

$$\begin{aligned} & \underset{\boldsymbol{\mu}_T}{\text{minimize}} \quad \|\boldsymbol{\Phi}_T \boldsymbol{\mu}_T - \mathbf{y}_{\text{ref}}\|_2^2 + \gamma \|\boldsymbol{\mu}_T\|_2^2 \\ & \text{subject to} \quad \|\boldsymbol{\mu}_T\|_\infty \leq u_{\max}, \end{aligned}$$

得到的这个问题是凸二次规划。也易于处理关于信号的瞬时速率的改变(转换速率)的约束, 通过观察可得

$$\begin{bmatrix} u(1) - u(0) \\ u(2) - u(1) \\ \vdots \\ u(T-1) - u(T-2) \end{bmatrix} = \mathbf{D} \boldsymbol{\mu}_T, \quad \mathbf{D} \doteq \begin{bmatrix} -1 & 1 & 0 & \cdots & 0 \\ 0 & -1 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & -1 & 1 \end{bmatrix}$$

如果用 $s_{\max}$ 表示对输入的转换速率的上限, 则具有振幅限制和输入信号转换速率限制的正则化问题为

$$\begin{aligned} & \underset{\boldsymbol{\mu}_T}{\text{minimize}} \quad \|\boldsymbol{\Phi}_T \boldsymbol{\mu}_T - \mathbf{y}_{\text{ref}}\|_2^2 + \gamma \|\boldsymbol{\mu}_T\|_2^2 \\ & \text{subject to} \quad \|\boldsymbol{\mu}_T\|_\infty \leq u_{\max}, \quad \|\mathbf{D} \boldsymbol{\mu}_T\|_\infty \leq s_{\max}. \end{aligned}$$

这是一个凸二次规划问题。

**例15.4.** 在此考虑例15.2中轨上小车的离散模型。给定给定初始状态 $\mathbf{x}(0) = [0 \ 0]^\top$ 和参考输出轨道

$$y_{\text{ref}}(k) = \sin(\omega k \Delta), \omega = \frac{2\pi}{10}, k = 1, \dots, 100,$$

我们寻找使得系统(15.10)的输出尽可能地接近 $y_{\text{ref}}$ 的输入信号 $u(k)$ ,  $k \in \{1, \dots, 100\}$ 。为此, 我们考虑式(15.17)给出的正则化最小二乘表述。我们首先求解 $\gamma = 0$ 时对应的问题。在这种情况下, 因为 $\mathbf{CB} = 0.005 \neq 0$ , 此时从原理上讲, 可以可得输入控制其使得输出信号精确地等于参考信号, 即可以做到精确跟踪, 结果如图15.5所示: 我们观察到跟踪误差的确是数值上的0(注意右图中 $e(k)$ 的量级是 $10^{-15}$ )。然而, 实现这个结果的代价是输入信号 $u(k)$ 剧烈震荡, 且振幅很大, 见图15.5中的左图。产生这种现象的原因之一是: 当 $\mathbf{CB} \neq 0$ 时, 尽管 $\boldsymbol{\Phi}_T$ 是可逆的, 但是 $\boldsymbol{\Phi}_T$ 的条件数随 $T$ 迅速上升(在这个问题中,  $\boldsymbol{\Phi}_T$ 的条件数大约为 $10^6$ )。

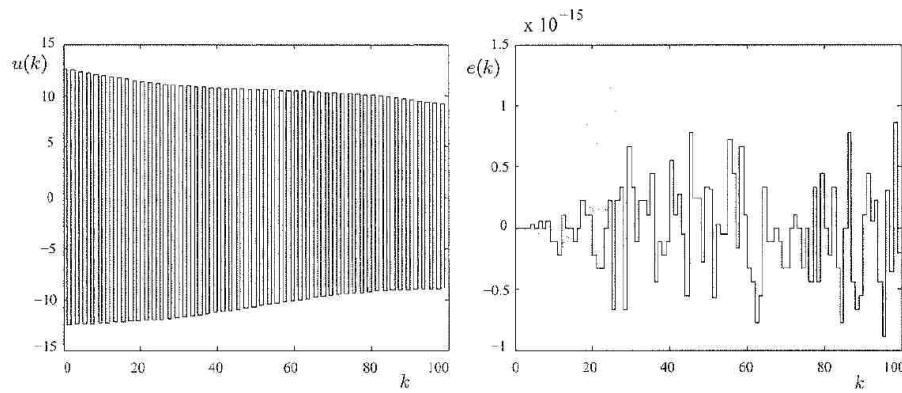


图 15.5: 输入信号 $u(k)$ (左)和 $\gamma = 0$ 时的跟踪误差 $e(k) = y(k) - y_{\text{ref}}(k)$ (右)

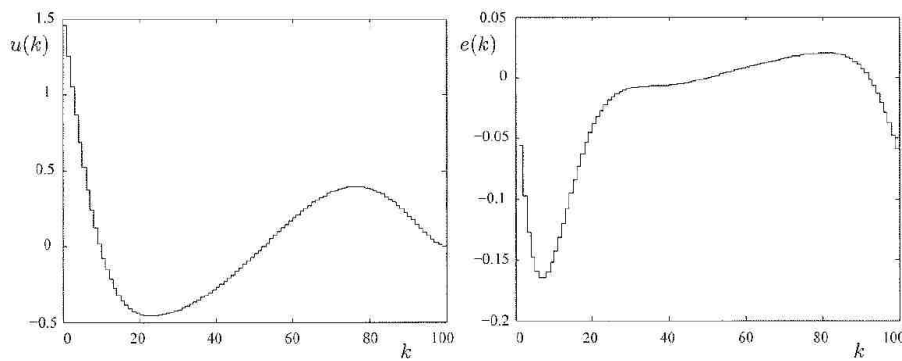


图 15.6: 输入信号 $u(k)$ (左)和 $\gamma = 0.1$ 时的跟踪误差 $e(k) = y(k) - y_{\text{ref}}(k)$ (右)

如果考虑输入信号的能量非零惩罚, 则可以改变(减小)最小二乘问题的法方程的系数矩阵的条件数(见§6.7.3.1), 从而大幅提升输入信号的性能。举例来说, 对于上例, 取 $\gamma = 0.1$ 并求解正则化的最小二乘问题(15.17), 所得结果如图15.6所示。我们观察到这里的跟踪误差比之前的高(但是可能仍然在可接受范围内, 这个需要根据具体问题具体分析), 然而与 $\gamma = 0$ 时的情况相比, 输入信号振幅变小了, 也平缓得多。

### §15.2.3 模型预测控制

模型预测控制(Model Predict Control, MPC)是一种非常有效的, 且被广泛应用于动力系统的控制方法, 它是基于最优化技术的。我们在这里简单地阐述LTI 离散时间系统中的MPC方法。前几小节中讨论的控制设计方法的共同特点是: 需要事先(也就是在时刻 0 时)计算定义在时长

为 $T$ 的给定时刻上的**整个**控制序列 $u(0), \dots, u(T-1)$ , 然后再将其应用到系统中。在实际操作的过程中, 因为将计算所得的输入应用于真实的系统时, 这个系统因干扰可能随时间发生变化, 从而有可能出现某些问题。如果时间阶段 $T$ 很大, 即使我们计算得到了控制序列, 然后靠后的瞬时时刻的系统状态因为干扰可能完全不同于我们在时刻 $0$  (即计算控制序列)时所事先预测的那样。像这种需要在多个时间段上进行决策(这里就是确定控制的输入)是一个很普遍的问题。

直觉表明以下作法更合理. 我们在时刻 $k=0$  计算序列

$$\mathbf{u}^{(0)} \doteq (u_{0|0}, \dots, u_{T-1|0}),$$

而后仅仅把这些控制中的第一个应用于 $k=0$ 时的系统, 即 $u(0) = u_{0|0}$ 。然后, 我们**静观其变**(wait and see), 即我们观察真实系统对第一个输入如何响应(即我们测量实际的状态 $\mathbf{x}(1)$ ), 然后重新计算向前平移时域一步后的新序列(记作 $\mathbf{u}^{(1)} \doteq (u_{0|1}, \dots, u_{T-1|1})$ )。因此, 在时刻 $k=1$ 时, 我们给系统施加这个序列的第一个值, 即 $u(1) = u_{0|1}$ , 然后不停地重复这个过程(即我们得到结果 $\mathbf{x}(2)$ , 计算新的序列 $\mathbf{u}^{(2)} = (u_{0|2}, \dots, u_{T-1|2})$ , 将 $u(2) = u_{0|2}$ 施加给系统, 以此类推)。称这种方法是MPC 控制率的**滑动时域**(sliding-horizon)实现。

下面考虑确定预测控制的一种基本模式. 对给定的时刻 $k$ , 首先测量真实系统在时刻 $k$ 的状态变量值, 记为 $\mathbf{x}(k)$ 。然后以包含了需要计算的时刻 $k$ 的预测序列 $\mathbf{u}_{0|k}$ 的变量

$$\mathbf{u}^{(k)} \doteq (\mathbf{u}_{0|k}, \dots, \mathbf{u}_{T|k})$$

作为决策变量, 并用 $\mathbf{x}_{j|k}, j=0, \dots, T$ , 表示预测的状态序列, 则其需要满足系统的动力学方程

$$\mathbf{x}_{j+1|k} = \mathbf{A}\mathbf{x}_{j|k} + \mathbf{B}\mathbf{u}_{j|k}, j=0, \dots, T-1. \quad (15.18a)$$

$$\mathbf{x}_{0|k} = \mathbf{x}(k). \quad (15.18b)$$

最后, MPC的滑动时域实现求解如下形式的最优化问题来得到 $\mathbf{u}^{(k)}$ , 即

$$\begin{aligned} & \underset{\mathbf{u}^{(k)}}{\text{minimize}} & J &= \sum_{j=0}^{T-1} (\mathbf{x}_{j|k}^\top \mathbf{Q} \mathbf{x}_{j|k} + \mathbf{u}_{j|k}^\top \mathbf{R} \mathbf{u}_{j|k}) + \mathbf{x}_{T|k}^\top \mathbf{S} \mathbf{x}_{T|k} \\ & \text{subject to} & \mathbf{u}_{1b}^{(j)} &\leq \mathbf{u}_{j|k} \leq \mathbf{u}_{ub}^{(j)}, j=0, \dots, T-1, \\ & & \mathbf{F}^{(j)} \mathbf{x}_{j|k} &\leq \mathbf{g}^{(j)}, j=1, \dots, T, \\ & & \mathbf{x}_{j+1|k} &= \mathbf{A}\mathbf{x}_{j|k} + \mathbf{B}\mathbf{u}_{j|k}, j=0, \dots, T-1 \\ & & \mathbf{x}_{0|k} &= \mathbf{x}(k), \end{aligned} \quad (15.19)$$

其中  $\mathbf{Q} \in \mathbb{S}_+^n, \mathbf{S} \in \mathbb{S}_+^n, \mathbf{R} \in \mathbb{S}_+^m$  是给定的权重矩阵,  $\mathbf{u}_{ub}^{(j)}, \mathbf{u}_{lb}^{(j)} \in \mathbb{R}^m$  是关于预测控制输入序列的给定的上界和下界向量,  $\mathbf{F}^{(j)} \in \mathbb{R}^{q \times n}, \mathbf{g}^{(j)} \in \mathbb{R}^q$  描述了可能的关于系统状态的多面集约束。

如果我们从  $i = 0$  到  $j$  递归使用方程(15.18a), 并注意到(15.18b), 得

$$\mathbf{x}_{j|k} = \mathbf{A}^j \mathbf{x}(k) + \mathbf{R}_j \mathbf{u}^{(k)}, \quad \mathbf{R}_j = [\mathbf{A}^{j-1} \mathbf{B} \cdots \mathbf{A} \mathbf{B} \quad \mathbf{B} \quad 0 \cdots 0].$$

这表明  $\mathbf{x}_{j|k}$  是决策变量  $\mathbf{u}^{(k)}$  的线性函数。将这个表述代入问题(15.19)的目标和前两组约束, 易于看到所得问题是关于变量  $\mathbf{u}^{(k)}$  的凸二次规划。

因此, 在每一个时刻  $k = 0, 1, \dots$ , 我们观测  $\mathbf{x}(k)$  (这需要状态变量可测量), 然后通过求解(15.19)来找到最佳的预测序列  $\mathbf{u}^{(k)} = (\mathbf{u}_{0|k}, \dots, \mathbf{u}_{T|k})$ , 最后将控制输入  $\mathbf{u}(k) = \mathbf{u}_{0|k}$  施加给真实的系统, 以此类推, 进行迭代。

问题(15.19)的目标函数量化了状态所蕴含能量和控制信号所蕴含能量的折衷 ( $\mathbf{x}_{j|k}^\top \mathbf{Q}_{j|k} \mathbf{x}_{j|k}$  和  $\mathbf{u}_{j|K}^\top \mathbf{R} \mathbf{u}_{j|K}$  分别代表着状态向量和控制向量的加权2-范数), 外加一个与此相似的关于最终状态的惩罚项  $\mathbf{x}_{T|k}^\top \mathbf{S} \mathbf{x}_{T|k}$ 。这种控制策略的目的是驱使(调节)系统的状态渐近地逼近 0。MPC的一个特色是它允许我们直接对控制信号施加限制(比如上下界等), 也可以直接对容许状态的轨道施加限制。

### §15.3 分析和控制器设计中的最优化

在本节中我们将简要概要讨论利用最优化分析系统的性质(比如稳定性), 或者为系统设计更合适的控制器。这些方法通常着重考虑系统的渐进性质。与上一节我们所讨论的利用最优化直接确定控制序列不同, 这些技术是间接的, 即最优化不是用来直接确定控制序列, 而是用来决定一些参数的值(例如控制器的增益), 这些参数与连接在系统上(例如一个反馈装置)以便对系统进行控制的理想器件(控制器)有关。以控制系统的分析和设计的最优化方法(尤其是基于SDP)为主题的文献现在已经非常多了。这里将通过考虑一些关于LTI系统的稳定性和反馈稳定性等非常基本的问题, 让大家体会这一类结论的特点。

#### §15.3.1 连续时间的Lyapunov 稳定性分析

称连续时间LTI系统

$$\dot{\mathbf{x}}(t) = \mathbf{A} \mathbf{x}(t) + \mathbf{B} \mathbf{u}(t) \quad (15.20)$$

是(渐进)稳定的, 如果对  $\mathbf{u}(t) = 0$ , 系统对所有的初态  $\mathbf{x}(0) = \mathbf{x}_0$  都有  $\lim_{t \rightarrow \infty} \mathbf{x}(t) = 0$  成立。用文字描述, 系统的自由响应(free response) 渐进地趋于 0, 并且

与初值无关。从Lyapunov稳定性理论, 我们知道一个系统稳定的充分必要条件是存在一个矩阵  $\mathbf{P} \succ \mathbf{0}$  使得

$$\mathbf{A}^\top \mathbf{P} + \mathbf{P} \mathbf{A} \prec \mathbf{0}, \quad (15.21)$$

这是一个针对矩阵变量  $\mathbf{P}$  的(严格)LMI条件。存在满足条件的矩阵  $\mathbf{P}$  为形如二次Lyapunov 函数  $V(\mathbf{x}) = \mathbf{x}^\top \mathbf{P} \mathbf{x}$  的系统稳定性提供了判据(certificate)。注意到  $\mathbf{P} \succ \mathbf{0}$  与(15.21)这两个条件是齐次的。这意味着如果存在某个  $\bar{\mathbf{P}}$  满足条件, 即

$$\bar{\mathbf{P}} \succ \mathbf{0}, \quad \mathbf{A}^\top \bar{\mathbf{P}} + \bar{\mathbf{P}} \mathbf{A} = -\bar{\mathbf{Q}}, \quad \bar{\mathbf{Q}} \succ \mathbf{0}$$

那么形如  $\mathbf{P} = \alpha \bar{\mathbf{P}}, \alpha > 0$  的这些矩阵同样满足这些条件。特别的, 如果我们选取

$$\alpha = \max\{\lambda_{\min}^{-1}(\bar{\mathbf{P}}), \lambda_{\min}^{-1}(\bar{\mathbf{Q}})\},$$

那么

$$\mathbf{P} = \alpha \bar{\mathbf{P}} \succeq \mathbf{I}, \quad \mathbf{A}^\top \mathbf{P} + \mathbf{P} \mathbf{A} \preceq -\mathbf{I}.$$

因此, 稳定性等价于满足下面两个关于  $\mathbf{P}$  的非严格LMI条件:

$$\mathbf{P} \succeq \mathbf{I}, \quad \mathbf{A}^\top \mathbf{P} + \mathbf{P} \mathbf{A} \preceq -\mathbf{I}.$$

这样, 寻找一个Lyapunov稳定性判据  $\mathbf{P}$  的问题可以表述成形如

$$\begin{aligned} & \underset{\mathbf{P}, \nu}{\text{minimize}} && \nu \\ & \text{subject to} && \mathbf{A}^\top \mathbf{P} + \mathbf{P} \mathbf{A} \preceq -\mathbf{I}, \\ & && \mathbf{I} \preceq \mathbf{P} \preceq \nu \mathbf{I}, \end{aligned}$$

的SDP问题。该表述的最后一个LMI约束表示  $\mathbf{P}$  的条件数不超过  $\nu$ , 因此目标函数表明我们要寻找一个条件数最小的Lyapunov矩阵。

一种等价的方法是用逆Lyapunov矩阵  $\mathbf{W} = \mathbf{P}^{-1}$  进行另一种表述。事实上, 给(15.21)的两边同时左乘和右乘  $\mathbf{P}^{-1}$ , 我们得到另一个稳定性条件:  $\exists \mathbf{W} \succ \mathbf{0}$ , 使得  $\mathbf{W} \mathbf{A}^\top + \mathbf{A} \mathbf{W} \prec \mathbf{0}$ , 与上面关于  $\mathbf{P}$  的基于齐次性的推理类似, 这些条件等价于

$$\exists \mathbf{W} \succeq \mathbf{I} : \mathbf{W} \mathbf{A}^\top + \mathbf{A} \mathbf{W} \preceq -\mathbf{I}. \quad (15.22)$$

从而, 求解SDP问题

$$\begin{aligned} & \underset{\mathbf{W}, \nu}{\text{minimize}} && \nu \\ & \text{subject to} && \mathbf{W} \mathbf{A}^\top + \mathbf{A} \mathbf{W} \preceq -\mathbf{I}, \\ & && \mathbf{I} \preceq \mathbf{W} \preceq \nu \mathbf{I} \end{aligned}$$

可以找到一个最优的(最小条件数)逆Lyapunov判据。



### §15.3.2 镇定状态反馈设计

能够将之前的方法从稳定性分析推广到反馈式镇定控制律设计。假设控制输入是状态反馈(state-feedback)式的, 即

$$\mathbf{u}(t) = \mathbf{K}\mathbf{x}(t), \quad (15.23)$$

其中  $\mathbf{K} \in \mathbb{R}^{m \times n}$  是状态反馈增益(gain)矩阵(或者控制器), 需要设计  $\mathbf{K}$  使得被控系统是稳定的。将(15.23)代入状态方程(15.20), 我们得到被控系统

$$\dot{\mathbf{x}} = (\mathbf{A} + \mathbf{BK})\mathbf{x}.$$

这个系统是稳定的当且仅当闭环系统的矩阵  $\mathbf{A}_{cc} = \mathbf{A} + \mathbf{BK}$  满足条件(15.22), 即当且仅当

$$\exists \mathbf{W} \succeq \mathbf{I}: \quad \mathbf{W}\mathbf{A}^\top + \mathbf{A}\mathbf{W} + (\mathbf{K}\mathbf{W})^\top \mathbf{B}^\top + \mathbf{B}(\mathbf{K}\mathbf{W}) \preceq -\mathbf{I}.$$

引入新变量  $\mathbf{Y} \doteq \mathbf{K}\mathbf{W}$ , 则可将这些条件改写为

$$\exists \mathbf{W} \succeq \mathbf{I}: \quad \mathbf{W}\mathbf{A}^\top + \mathbf{A}\mathbf{W} + \mathbf{Y}^\top \mathbf{B}^\top + \mathbf{B}\mathbf{Y} \preceq -\mathbf{I}.$$

这是关于变量  $\mathbf{W}$  和  $\mathbf{Y}$  的LMI条件。我们先求解一个以上述LMI为约束的凸优化问题, 然后由  $\mathbf{K} = \mathbf{Y}\mathbf{W}^{-1}$  得到镇定反馈增益  $\mathbf{K}$ 。例如, 我们可以将  $\mathbf{W}$  的条件数和  $\mathbf{Y}$  的范数的折衷作为最优化问题的目标, 即考虑

$$\begin{aligned} & \underset{\mathbf{W}, \mathbf{Y}, \nu}{\text{minimize}} && \nu + \eta \|\mathbf{Y}\|_2 \\ & \text{subject to} && \mathbf{W}\mathbf{A}^\top + \mathbf{A}\mathbf{W} + \mathbf{Y}^\top \mathbf{B}^\top + \mathbf{B}\mathbf{Y} \preceq -\mathbf{I}, \\ & && \mathbf{I} \preceq \mathbf{W} \preceq \nu \mathbf{I}, \end{aligned}$$

其中  $\eta \geq 0$  是折衷参数。

#### §15.3.2.1 鲁棒的反馈设计

上面描述的基于LMI的反馈稳定化是相当灵活的, 因为可以将其推广到不确定(uncertain)系统的稳定化。在这里, 我们只讨论不确定性影响系统时最简单的情况, 即所谓的场景(scenario)或者多面体不确定性。考虑系统

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t), \quad (15.24)$$

其中系统矩阵  $\mathbf{A}, \mathbf{B}$  可以在由给定的矩阵  $(\mathbf{A}_i, \mathbf{B}_i), i = 1, \dots, N$ , 张成的多面体内变化, 即

$$(\mathbf{A}, \mathbf{B}) \in \text{co}\{(\mathbf{A}_1, \mathbf{B}_1), \dots, (\mathbf{A}_N, \mathbf{B}_N)\}$$

矩阵对 $(\mathbf{A}_i, \mathbf{B}_i)$ 是多面体的顶点, 并且可以将它们解释为场景, 或者一个不确定的设备的不同的可能实现, 也可以将它们解释为在不同的操作条件或者时间点下的测量得到。可以证明不确定系统(15.24)稳定的充分条件 是对所有顶点系统存在一个公共的二次Lyapunov 函数, 即

$$\exists \mathbf{W} \succeq \mathbf{I} : \mathbf{W} \mathbf{A}_i^\top + \mathbf{A}_i \mathbf{W} \preceq -\mathbf{I}, i = 1, \dots, N.$$

基于这个充分条件, 我们可以设计形如(15.23)的反馈控制律以便能鲁棒地镇定不确定系统。可将待控系统描述为

$$\dot{\mathbf{x}}(t) = (\mathbf{A} + \mathbf{B}\mathbf{K})\mathbf{x}(t).$$

如果

$$\exists \mathbf{W} \succeq \mathbf{I} : \mathbf{W}(\mathbf{A}_i + \mathbf{B}_i \mathbf{K})^\top + (\mathbf{A}_i + \mathbf{B}_i \mathbf{K})\mathbf{W} \preceq -\mathbf{I}, i = 1, \dots, N,$$

则待控系统是稳定的。引入新变量 $\mathbf{Y} = \mathbf{K}\mathbf{W}$ , 则求解如下优化问题可得镇定控制器的反馈增益矩阵 $\mathbf{K}$ 。为此, 先求解

$$\begin{aligned} & \underset{\mathbf{W}, \mathbf{Y}, \nu}{\text{minimize}} && \nu + \eta \|\mathbf{Y}\|_2 \\ & \text{subject to} && \mathbf{W} \mathbf{A}_i^\top + \mathbf{A}_i \mathbf{W} + \mathbf{Y}^\top \mathbf{B}_i^\top + \mathbf{B}_i \mathbf{Y} \preceq -\mathbf{I}, i = 1, \dots, N, \\ & && \mathbf{I} \preceq \mathbf{W} \preceq \nu \mathbf{I}, \end{aligned} \quad (15.25)$$

然后令 $\mathbf{K} = \mathbf{Y}\mathbf{W}^{-1}$ 。

**例15.5. (倒立摆的鲁棒镇定)** 考虑一个简化的倒立摆模型, 如图15.7所示, 系统由长为 $\ell$ 的刚性杆和顶端的质量为 $m$ 的点组成, 其中 $\theta(t)$ 表示杆关于竖直坐标轴的角位移,  $u(t)$ 是施加在连接点的输入扭矩,  $g$  是重力加速度。针对该系统牛顿平衡方程是

$$m\ell^2\ddot{\theta}(t) = mg\ell \sin \theta(t) + u(t).$$

这个二阶非线性微分方程描述了系统的动力学。然而, 如果角位移 $\theta(t)$ 持续很小时, 即 $|\theta(t)| \simeq 0$ 时, 我们可以利用 $\sin \theta(t) \simeq \theta(t)$ 这个事实, 用一个线性微分方程

$$m\ell^2\ddot{\theta}(t) = mg\ell\theta(t) + u(t)$$

来近似这个非线性方程。引入状态变量 $x_1(t) = \theta(t), x_2 = \dot{\theta}(t)$ , 这个连续时间LTI系统为

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}u, \mathbf{A} = \begin{bmatrix} 0 & 1 \\ \alpha & 0 \end{bmatrix}, \mathbf{B} = \begin{bmatrix} 0 \\ \beta \end{bmatrix},$$

其中  $\alpha \doteq \frac{g}{\ell}$ ,  $\beta \doteq \frac{1}{m\ell^2}$ . 对于参数任何取值(正的), 这个系统都是不稳定的。假设给定了如下的数值:

$$g = 9.8\text{m/s}^2, \quad m = 0.5\text{kg}, \quad \ell \in [0.9, 1.1]\text{m},$$

即我们不能确定杆的长度。我们希望设计一个状态控制反馈律  $u(t) = Kx(t)$  能鲁棒地镇定该系统。为此, 我们观察到系统矩阵在由两个顶点系统定义的线段上变化, 具体的顶点系统

$$A_1 = \begin{bmatrix} 0 & 1 \\ 10.8889 & 0 \end{bmatrix}, \quad B_1 = \begin{bmatrix} 0 \\ 2.4691 \end{bmatrix};$$

$$A_2 = \begin{bmatrix} 0 & 1 \\ 8.9091 & 0 \end{bmatrix}, \quad B_2 = \begin{bmatrix} 0 \\ 1.6529 \end{bmatrix}.$$

因此, 由多面体系统的关于稳定性的充分条件, 取  $N = 2, \eta = 1$  并求解凸规划问题(15.25)得到控制增益  $K$ 。利用 CVX(一种求解凸优化的软件)求得数值解, 得到鲁棒镇定控制器  $K = -[6.7289, 2.3151]$  及 Lyapunov 判据

$$P = W^{-1} = \begin{bmatrix} 0.8796 & 0.2399 \\ 0.2399 & 0.5218 \end{bmatrix}.$$

**注记15.1 频域分析.** 前面几小节的分析和设计方法都是基于系统的时域(time-domain)表示。另一种可选的经典方法是在变换域上分析系统, 通常称这个域为频域(frequency domain)。一个特殊的映射将时域信号变换成复变量的函数, 从而完成信号从时域到频域的转变, 称这个映射为(单边)拉普拉斯变换(Laplace transform), 其将时域信号  $w(t) \in \mathbb{C}^n$  映射为函数

$$W(s) \doteq \int_0^\infty w(t)e^{-st}dt,$$

其中  $s \in \mathbb{C}$ . 如果对所有  $s$ , 上面的积分均收敛到一个有限值, 则  $W(s)$  是有定义的。将  $W(s)$  在  $s$  取值为纯虚数时的限制记为

$$\hat{W}(\omega) \doteq W(s)|_{s=j\omega} = \int_0^\infty w(t)e^{-j\omega t}dt,$$

而且当  $w(t)$  是一个因果信号<sup>1</sup>, 那么该限制与  $w(t)$  的傅立叶变换是一致的。如果  $w$  是实数, 那么  $\hat{W}(-\omega)^\top = \hat{W}(\omega)^*$ , 这里  $*$  代表转置共轭。由

<sup>1</sup>即当  $t < 0$  时取值全为零的信号。通常假定控制中的信号是因果的, 因为习惯上假定瞬态  $t = 0$  为系统启动的时间。

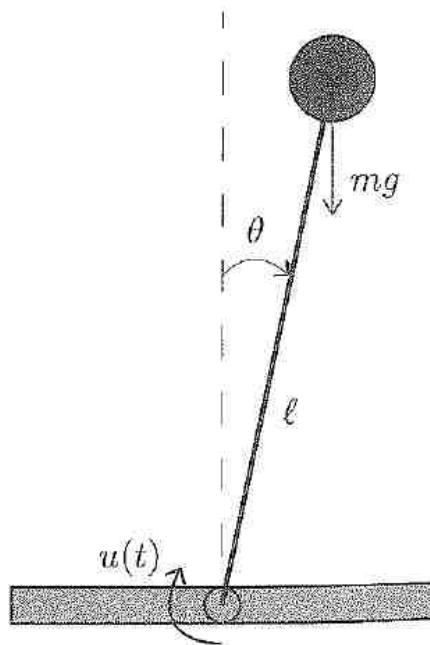


图 15.7: 一个倒立摆

此知道，对一个实信号，它的范数的平方  $\|\hat{\mathbf{W}}(\omega)\|_2^2 = \hat{\mathbf{W}}(-\omega)^* \hat{\mathbf{W}}(\omega) = \hat{\mathbf{W}}(-\omega)^\top \hat{\mathbf{W}}(\omega)$  关于  $\omega$  是中心对称的。

对于一个连续时间LTI系统

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t), \mathbf{x}(0) = \mathbf{0}, \quad (15.26)$$

$$\mathbf{y}(t) = \mathbf{C}\mathbf{x}(t) + \mathbf{D}\mathbf{u}(t), \quad (15.27)$$

其中  $\mathbf{x}(t) \in \mathbb{R}^n$  是状态， $\mathbf{u}(t) \in \mathbb{R}^p$  是输入， $\mathbf{y}(t) \in \mathbb{R}^m$  是输出，并且系统矩阵  $(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$  的维数是相容的。定义系统的传递矩阵(transfer matrix)为

$$\mathbf{H}(s) \doteq \mathbf{C}(s\mathbf{I} - \mathbf{A})^{-1}\mathbf{B} + \mathbf{D}.$$

在Laplace域中，可以将输入和输出的关系间接地表述为

$$\mathbf{Y}(s) = \mathbf{H}(s)\mathbf{U}(s),$$

其中  $\mathbf{U}(s), \mathbf{Y}(s)$  分别是  $\mathbf{u}(t), \mathbf{y}(t)$  的Laplace变换。从而传递矩阵能够简单地描述系统在Laplace域的某些方面的行为。矩阵  $\mathbf{H}(s)$  的元素  $H_{ij}(s)$  是实信号  $h_{ij}(t)$  的Laplace变换，它表示着当输入向量的第  $j$  个元素  $u_j(t)$  是一个脉冲(狄拉克  $\delta$  函数)，其它输入  $u_k(t), k \neq j$ ，均为 0 时，输出向量的第  $i$  个

元 $y_i(t)$ 。当然, 如果输入信号的傅里叶变换是 $U(\omega)$ , 并且系统是稳定的, 那么输出信号的傅立叶变换是

$$\hat{Y}(\omega) = H(j\omega)\hat{U}(\omega).$$

这里的 $\hat{H}(\omega) \doteq H(j\omega)$ ,  $\omega \geq 0$ , 就是所谓的系统的频率响应(frequency response)。

### §15.3.3 离散时间的Lyapunov稳定性和镇定

和连续时间的情况类似, 称一个离散时间 LTI 系统

$$\mathbf{x}(k+1) = \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k) \quad (15.28)$$

是(渐进)稳定的, 如果它的自由响应趋于零, 并与初值无关。系统(15.28)是稳定的当且仅当存在  $\mathbf{P} \succ \mathbf{0}$  满足

$$\mathbf{A}^\top \mathbf{P} \mathbf{A} - \mathbf{P} \prec \mathbf{0}, \quad (15.29)$$

或者, 由齐次性, 等价地有

$$\exists \mathbf{P} \succeq \mathbf{I} : \mathbf{A}^\top \mathbf{P} \mathbf{A} - \mathbf{P} \preceq -\mathbf{I}.$$

同样, 定义  $\mathbf{W} = \mathbf{P}^{-1}$  并且给方程(15.29)左乘  $\mathbf{W}$ , 右乘  $\mathbf{W}$ , 我们得到条件

$$\exists \mathbf{W} \succ \mathbf{0} : \mathbf{W} - \mathbf{W}\mathbf{A}^\top \mathbf{W}^{-1} \mathbf{A}\mathbf{W} \succ \mathbf{0}.$$

并将其去齐次化, 得

$$\exists \mathbf{W} \succeq \mathbf{I} : \mathbf{W} - \mathbf{W}\mathbf{A}^\top \mathbf{W}^{-1} \mathbf{A}\mathbf{W} \succeq \mathbf{I}.$$

根据Schur补的性质, 可将上面的条件写成仅仅是关于  $\mathbf{W}$  的LMI条件

$$\begin{bmatrix} \mathbf{W} - \mathbf{I} & \mathbf{W}\mathbf{A}^\top \\ \mathbf{A}\mathbf{W} & \mathbf{W} \end{bmatrix} \succeq \mathbf{0}. \quad (15.30)$$

该表述在反馈控制设计中很有用: 假设我们希望为系统(15.28)设计一个状态反馈律  $\mathbf{u}(k) = \mathbf{K}\mathbf{x}(k)$ , 那么受控系统是

$$\mathbf{x}(k+1) = (\mathbf{A} + \mathbf{B}\mathbf{K})\mathbf{x}(k),$$

该系统是稳定的当且仅当(应用(15.30))

$$\begin{bmatrix} \mathbf{W} - \mathbf{I} & \mathbf{W}(\mathbf{A} + \mathbf{B}\mathbf{K})^\top \\ (\mathbf{A} + \mathbf{B}\mathbf{K})\mathbf{W} & \mathbf{W} \end{bmatrix} \succeq \mathbf{0}.$$

定义新变量  $\mathbf{Y} = \mathbf{K}\mathbf{W}$ , 上式就变成一个关于  $\mathbf{W}$  和  $\mathbf{Y}$  的 LMI 条件。从而, 求解关于  $\mathbf{W}, \mathbf{Y}$  的最优化问题, 即

$$\begin{aligned} & \underset{\mathbf{W}, \mathbf{Y}, \nu}{\text{minimize}} && \nu + \eta \|\mathbf{Y}\|_2 \\ & \text{subject to} && \begin{bmatrix} \mathbf{W} - \mathbf{I} & \mathbf{W}\mathbf{A}^\top + \mathbf{Y}^\top \mathbf{B}^\top \\ \mathbf{A}\mathbf{W} + \mathbf{B}\mathbf{Y} & \mathbf{W} \end{bmatrix} \succeq \mathbf{0}, \\ & && \mathbf{W} \preceq \nu \mathbf{I}, \end{aligned} \quad (15.31)$$

并且令  $\mathbf{K} = \mathbf{Y}\mathbf{W}^{-1}$ , 即得镇定控制增益  $\mathbf{K}$ 。用和之前连续时间系统基本相似的思路, 也可以将这种方法推广以解决不确定的多面体系统的鲁棒镇定问题。

### §15.4 练习

**习题15.1. (稳定性和特征值)** 证明连续时间系统(15.20)是渐进稳定的(或者简称为稳定的)当且仅当  $\mathbf{A}$  矩阵的所有特征值  $\lambda_i(\mathbf{A})$  的实部是(严格)负的,  $i = 1, \dots, n$ 。

证明离散时间系统(15.28)是稳定的当且仅当矩阵  $\mathbf{A}$  的所有特征值  $\lambda_i(\mathbf{A})$  的模(严格)小于1,  $i = 1, \dots, n$ 。

**提示:** 利用表达式  $\mathbf{x}(t) = e^{\mathbf{A}t}\mathbf{x}_0$  和  $\mathbf{x}(k) = \mathbf{A}^k\mathbf{x}_0$  分别表示连续时间系统和离散时间系统的自由响应。你可以在  $\mathbf{A}$  是可对角化的假设下推导你的证明。

**习题15.2. (信号范数)** 一个连续时间信号  $\mathbf{w}(t)$  映射  $t \in \mathbb{R}$  到  $\mathbf{w}(t) \in \mathbb{C}^m$  或者  $\mathbb{R}^m$ 。信号的**能量**(energy)定义为

$$E(\mathbf{w}) \doteq \|\mathbf{w}\|_2^2 = \int_{-\infty}^{\infty} \|\mathbf{w}(t)\|_2^2 dt,$$

其中  $\|\mathbf{w}\|_2$  是信号的2范数。能量有限的信号类指上面的2范数有限的信号全体。

周期信号是典型的能量无限的信号。对一个周期为  $T$  的信号, 我们定义它的**功率**(power)为

$$P(\mathbf{w}) \doteq \frac{1}{T} \int_{t_0}^{t_0+T} \|\mathbf{w}(t)\|_2^2 dt.$$

- (a) 计算谐波信号  $\mathbf{w}(t) = \mathbf{v}e^{j\omega t}$ ,  $\mathbf{v} \in \mathbb{R}^m$ , 和因果指数信号  $\mathbf{w}(t) = \mathbf{v}e^{at}$ ,  $a < 0$ ,  $t \geq 0$ ,  $\mathbf{w}(t) = 0$ ,  $t < 0$ , 的能量。

(b) 计算谐波信号  $\mathbf{w}(t) = \mathbf{v}e^{j\omega t}$  和正弦信号  $\mathbf{w}(t) = \mathbf{v} \sin(\omega t)$  的功率.

**习题15.3. (系统的状态演化的能量上界)** 考虑一个连续时间LTI系统  $\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t), t \geq 0$ , 没有输入(称这种系统是自治的(autonomous)), 输出为  $\mathbf{y}(t) = \mathbf{C}\mathbf{x}$ . 我们希望求出系统的输出信号所包含的能量, 其由指标

$$J(\mathbf{x}_0) \triangleq \int_0^\infty \mathbf{y}(t)^\top \mathbf{y}(t) dt = \int_0^\infty \mathbf{x}(t)^\top \mathbf{Q} \mathbf{x}(t) dt,$$

所度量, 这里  $\mathbf{Q} \triangleq \mathbf{C}^\top \mathbf{C} \succeq \mathbf{0}$ .

(a) 证明: 如果系统稳定, 那么对所有的  $\mathbf{x}_0$ ,  $J(\mathbf{x}_0) < \infty$ .

(b) 证明: 如果系统稳定, 并且存在一个矩阵  $\mathbf{P} \succeq \mathbf{0}$ , 使得

$$\mathbf{A}^\top \mathbf{P} + \mathbf{P} \mathbf{A} + \mathbf{Q} \preceq \mathbf{0},$$

那么  $J(\mathbf{x}_0) \leq \mathbf{x}_0^\top \mathbf{P} \mathbf{x}_0$  成立. 提示: 考虑二次型  $V(\mathbf{x}(t)) = \mathbf{x}(t)^\top \mathbf{P} \mathbf{x}(t)$ , 计算它对时间的导数.

(c) 对于给定的初始条件, 解释如何计算一个系统的状态能量的最小上界.

**习题15.4. (系统增益)** 一个系统的增益(gain)是指从输入信号到输出信号的最小的能量放大率. 任意具有有限能量的输入信号  $\mathbf{u}(t)$  通过稳定系统后, 会被映射成一个同样具有有限能量的输出信号  $\mathbf{y}(t)$ . 帕塞瓦恒等式(Parseval identity)给出了信号  $\mathbf{w}(t)$  在时域的能量和信号在频域的能量(参考注记15.1)的关系, 即

$$E(\mathbf{w}) \triangleq \|\mathbf{w}\|_2^2 = \int_{-\infty}^\infty \|\mathbf{w}(t)\|_2^2 dt = \frac{1}{2\pi} \int_{-\infty}^\infty \|\hat{\mathbf{W}}(\omega)\|_2^2 d\omega \triangleq \|\hat{\mathbf{W}}\|_2^2.$$

系统(15.26)的能量增益(energy gain)定义为

$$\text{EG} \triangleq \sup_{\mathbf{u}(t): \|\mathbf{u}\|_2 < \infty, \mathbf{u} \neq \mathbf{0}} \frac{\|\mathbf{y}\|_2^2}{\|\mathbf{u}\|_2^2},$$

(a) 利用上面的信息证明: 对于一个稳定的系统,

$$\text{EG} \leq \sup_{\omega \geq 0} \|\mathbf{H}(j\omega)\|_2^2,$$

其中  $\|\mathbf{H}(j\omega)\|_2$  是系统(15.26)的传递矩阵在  $s = j\omega$  处的谱范数. 称系统的能量增益的平方根为  $\mathcal{H}_\infty$  范数, 记作  $\|\mathbf{H}\|_\infty$

提示: 运用帕塞瓦恒等式并证明某积分有上界. 注意等式实际上蕴含在前面的公式中, 但是并没有要求你去证明该事实.

(b) 假设系统(15.26)是稳定的,  $\mathbf{x}(0) = \mathbf{0}, \mathbf{D} = \mathbf{0}$ 。证明: 如果存在  $\mathbf{P} \succeq \mathbf{0}$  使得

$$\begin{bmatrix} \mathbf{A}^\top \mathbf{P} + \mathbf{P} \mathbf{A} + \mathbf{C}^\top \mathbf{C} & \mathbf{P} \mathbf{B} \\ \mathbf{B}^\top \mathbf{P} & -\gamma^2 \mathbf{I} \end{bmatrix} \preceq \mathbf{0}, \quad (15.32)$$

那么

$$\|\mathbf{H}\|_\infty \leq \gamma.$$

设计一种计算格式, 其能得到系统能量增益的最小可能的上界  $\gamma^*$ 。

提示: 定义二次函数  $V(\mathbf{x}) = \mathbf{x}^\top \mathbf{P} \mathbf{x}$ , 并观察  $V$  对时间  $t$  的导数. 沿着系统(15.26)的轨道, 这个导数是

$$\frac{dV(\mathbf{x})}{dt} = \mathbf{x}^\top \mathbf{P} \dot{\mathbf{x}} + \dot{\mathbf{x}}^\top \mathbf{P} \mathbf{x}.$$

然后说明LMI条件(15.32)等价于条件:

$$\frac{dV(\mathbf{x})}{dt} + \|\mathbf{y}\|^2 - \gamma^2 \|\mathbf{u}\|^2 \leq 0, \forall \mathbf{x}, \mathbf{u}, \text{ satisfying (15.26),}$$

并且这反过来蕴含着  $\|\mathbf{H}\|_\infty \leq \gamma$ 。

**习题15.5. (扩展超稳定矩阵)** 设矩阵  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , 如果存在  $\mathbf{d} \in \mathbb{R}^n$  使得

$$\sum_{j \neq i} |a_{ij}| d_j < -a_{ii} d_i, d_i > 0, i = 1, \dots, n,$$

则称矩阵  $\mathbf{A}$  是连续时间**扩展超稳定的**(extended wuperstable)<sup>2</sup>(记为  $\mathbf{A} \in E_c$ ). 类似的, 如果存在  $\mathbf{d} \in \mathbb{R}^n$  使得

$$\sum_{j=1}^n |a_{ij}| d_j < d_i, d_i > 0, i = 1, \dots, n,$$

则称矩阵  $\mathbf{A}$  是离散时间扩展超稳定的(记为  $\mathbf{A} \in E_d$ ).

如果  $\mathbf{A} \in E_c$ , 那么它所有的特征值的实部都小于0, 因此相应的连续时间LTI系统  $\dot{\mathbf{x}} = \mathbf{A} \mathbf{x}$  是稳定的。类似的, 如果  $\mathbf{A} \in E_d$ , 那么  $\mathbf{A}$  的所有特征值的模小于1, 因此相应的离散时间LTI系统  $\mathbf{x}(k+1) = \mathbf{A} \mathbf{x}(k)$  是稳定的。扩展超稳定提供了稳定性的充分条件, 它的优势在于仅需检验一系列线性不等式的可行性。

<sup>2</sup>见B.T. Polyak, Extended superstability in control theory, *Automation and Remote Control*, 2004.



- (a) 给定一个连续时间系统  $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}$ ,  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{u} \in \mathbb{R}^m$ , 描述一种有效方法, 其设计的形为  $\mathbf{u} = -\mathbf{K}\mathbf{x}$  的状态反馈控制律使得受控系统是扩展超稳定的。
- (b) 给定离散时间系统  $\mathbf{x}(k+1) = \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k)$ , 假设矩阵  $\mathbf{A}$  具有区间不确定性, 即

$$a_{ij} = \hat{a}_{ij} + \delta_{ij}, i, j = 1, \dots, n$$

其中  $\hat{a}_{ij}$  是名义上的元素,  $\delta_{ij}$  是一个不确定项. 对给定的  $r_{ij} \geq 0$ , 仅知道  $|\delta_{ij}| \leq \rho r_{ij}$ . 定义扩展超稳定的半径是使得  $\mathbf{A}$  对  $\rho \geq 0$  时所有可容许的不确定量都是扩展超稳定的那个最大的, 记为  $\rho^*$ . 给出一种确定  $\rho^*$  的计算方法。

## 第十六章 工程设计

凸模型和凸优化在过去几十年对工程设计产生巨大的影响. 随着可靠的最优化技术(起先的线性规划, linear programming, 后来的二次规划, quadratic programming, 和锥规划)的出现, 工程师开始回顾各种分析和设计问题, 发现可以用凸模型来表述它们, 因此可以被有效求解. 整个领域, 诸如自动控制、电路设计, 在二十世纪九十年代因为引入了凸规划的方法论(特别是SDP)而发生了革命性的变化. 今天, 凸模型常被用来解决相关问题中的一部分, 例如结构力学、识别、过程控制、滤波器设计、电子电路的宏建模、物流和管理、网络设计等等. 在这一章中, 将详述从这些应用的一部分.

### §16.1 数字滤波器设计

一个单输入和单输出的数字滤波器是一个动力系统, 具有标量输入信号 $u(t)$ 和标量输出信号 $y(t)$ , 这里 $t$ 代表了(离散的)时间变量. 有限冲击响应(finite-impulse response, FIR)滤波器是一种特殊的滤波器, 形如

$$y(t) = \sum_{i=0}^{n-1} h_i u(t-i), \quad t \in \mathbb{Z},$$

这里 $h_0, \dots, h_{n-1}$ 叫做滤波器的**脉冲响应**(impulse response), 该滤波器得名源于如下事实: 滤波器对离散时间脉冲

$$u(t) = \begin{cases} 1 & t = 0, \\ 0 & \text{其它}, \end{cases}$$

的时间响应刚好是有限支集信号

$$y(t) = \begin{cases} h_t & 0 \leq t \leq n-1, \\ 0 & \text{其它}. \end{cases}$$

脉冲响应的离散傅里叶变换是一个复值函数 $H: \mathbb{R} \rightarrow \mathbb{C}$ , 其值

$$H(\omega) = \sum_{t=0}^{n-1} h_t e^{-j\omega t}, \quad \omega \in [-\pi, \pi].$$

该函数非常重要, 因为它表明了滤波器如何对周期信号做出响应. 准确地说, 如果输入是复指数 $u(t) = e^{j\omega_0 t}$ , 则输出信号将被放缩成复指数 $y(t) =$

$H(\omega_0)e^{j\omega_0 t}$ . 因为 $H(\omega)$ 以 $2\pi$ 为周期并且 $H^*(\omega) = H(-\omega)$ , 我们可以将分析限定在规范区间 $[0, \pi]$ 上. 一个简单的FIR滤波器是滑动平均滤波器: 长度为2的滑动平均滤波形如

$$y(t) = \frac{1}{2}(u(t) + u(t-1)), \quad t \in \mathbb{Z}.$$

### §16.1.1 线性相位FIR滤波器

**I型线性相位滤波器**(Type I linear-phase filter)是一类特殊的FIR滤波器, 它的特征是项数为奇数, 即 $n = 2N + 1$ , 并且脉冲响应关于中点对称, 即

$$h_t = h_{n-1-t}, \quad t = 0, \dots, n-1.$$

这里的修饰语“线性相位”源于如下事实: 对这类滤波器, 它的频率响应形如

$$H(\omega) = e^{-j\omega N} \tilde{H}(\omega), \quad \omega \in [0, \pi],$$

其中 $\tilde{H}(\omega)$ 是实值函数

$$\tilde{H}(\omega) = h_N + 2 \sum_{t=0}^{N-1} h_t \cos((N-t)\omega).$$

称函数 $\tilde{H}(\omega)$ 为滤波器的**振幅响应**(amplitude response). 需要注意的是 $H(\omega)$ 的实际相位也许是不连续的(因为实际相位为 $-\text{sgn}(\tilde{H}(\omega))\omega N$ , 并且 $H(\omega)$ 的模可以表示为

$$|H(\omega)| = |\tilde{H}(\omega)|.$$

有些情况下, 为了方便工程中采用“连续线性相位” $\theta(\omega) = -\omega N$ 和振幅响应函数 $\tilde{H}(\omega)$ (它是实的, 但是能够取正值或者负值)来表示 $H(\omega)$ , 而不是直接用 $H(\omega)$ 实际的相位和模来表示. 例如, 我们例子中的 $\tilde{H}(\omega)$ 是设计参数 $\mathbf{h}$ 的一个简单线性函数, 即, 我们可以写为

$$\tilde{H}(\omega) = \mathbf{a}^\top(\omega) \mathbf{h}, \quad \mathbf{h} = \begin{bmatrix} h_0 \\ h_1 \\ \vdots \\ h_{N-1} \\ h_N \end{bmatrix}, \quad \mathbf{a}(\omega) = \begin{bmatrix} 2 \cos \omega N \\ 2 \cos \omega(N-1) \\ \vdots \\ 2 \cos \omega \\ 1 \end{bmatrix}.$$

我们进一步发现, 不失一般性, 可选取脉冲响应向量 $\mathbf{h}$ , 使得 $\tilde{H}(0) > 0$ .

## §16.1.2 低通FIR设计规范

与FIR滤波器有关的设计问题通常涉及选取滤波器的脉冲响应 $h$ ，以便获得具有理想形状的振幅响应。例如，要确保滤波器阻止高频信号，但是允许低频信号通过(低通滤波器)。可以将这些要求对应到下面规定的滤波器的振幅响应，也可以参考图16.1所示的振幅设计约束。

- 阻带约束：

$$-\delta_s \leq \tilde{H}(\omega) \leq \delta_s, \quad \omega \in [\Omega_s, \pi],$$

其中 $\Omega_s$ 是阻带边沿频率， $\delta_s > 0$ 对应着我们对高频需要获取的衰减水平。

- 通带约束：

$$1 - \delta_p \leq \tilde{H}(\omega) \leq 1 + \delta_p, \quad \omega \in [0, \Omega_p],$$

其中 $\Omega_p$ 是通带截止频率， $\delta_p > 0$ 界定了低频通带波纹。

注意到，对极小的 $\delta_s$ 和 $\delta_p$ ，上述约束近似等同于下面的两个关于 $H(\omega)$ 的模取对数(以10为底)意义下的约束：

$$\begin{aligned} \log |H(\omega)| &\leq \log \delta_s, & \omega \in [\Omega_s, \pi], \\ -\log \delta_p &\leq \log |H(\omega)| \leq \log \delta_p, & \omega \in [0, \Omega_p]. \end{aligned}$$

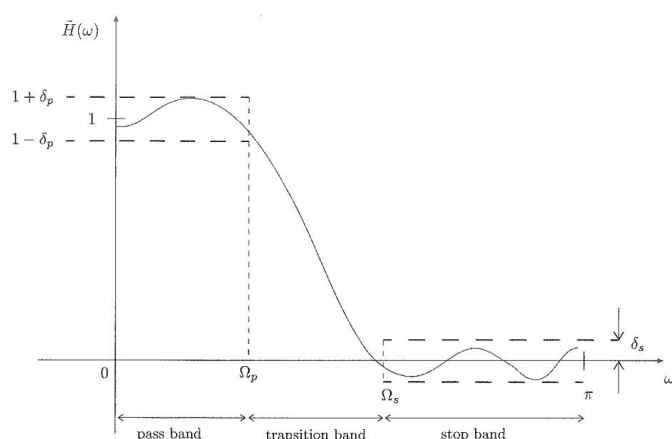


图 16.1: 低通FIR设计的模板

## §16.1.3 通过线性规划设计FIR

注意到上述设计约束涉及到各自区间上关于区间中每一个频率 $\omega$ 处的线性不等式约束集. 因此, 它们实际上涉及到无限多个线性约束. 为了克服该问题, 这里简单地将频率区间进行离散处理. 做法就是在对应的区间上取有限个不同的频率点, 对这些点处的频率要求约束满足, 而不是强迫区间上的每个频率都满足约束. 将在高频区间 $[\Omega_s, \pi]$ 取的点记为 $\omega_i$ ,  $i = 1, \dots, N_s$ , 并由有限个关于 $\mathbf{h}$ 的线性不等式

$$-\delta_s \leq \mathbf{a}^\top(\omega_i)\mathbf{h} \leq \delta_s, \quad i = 1, \dots, N_s$$

来近似阻带约束. 同样的, 将在低频区间 $[0, \Omega_p]$ 上取的点记为 $\omega_i$ ,  $i = N_s + 1, \dots, N_s + N_p$ , 得到通带约束的近似:

$$1 - \delta_p \leq \mathbf{a}^\top(\omega_i)\mathbf{h} \leq 1 + \delta_p, \quad i = N_s + 1, \dots, N_s + N_p.$$

在前文提出的假设条件下, 低通FIR滤波器设计问题可以通过各种方式表述成线性规划. 例如, 我们限定要求的阻带衰减水平 $\delta_s > 0$ , 寻找 $\mathbf{h}$ 和 $\delta_s$ 来最小化通带波纹, 即

$$\begin{aligned} & \underset{\mathbf{h} \in \mathbb{R}^{N+1}, \delta_p \in \mathbb{R}}{\text{minimize}} && \delta_p \\ & \text{subject to} && -\delta_s \leq \mathbf{a}^\top(\omega_i)\mathbf{h} \leq \delta_s, && i = 1, \dots, N_s, \\ & && 1 - \delta_p \leq \mathbf{a}^\top(\omega_i)\mathbf{h} \leq 1 + \delta_p, && i = N_s + 1, \dots, N_s + N_p. \end{aligned}$$

同样地, 还可以限定通带波纹 $\delta_p > 0$ , 来最小化阻带衰减水平 $\delta_s$ , 得到的LP问题形如

$$\begin{aligned} & \underset{\mathbf{h} \in \mathbb{R}^{N+1}, \delta_s \in \mathbb{R}}{\text{minimize}} && \delta_s \\ & \text{subject to} && -\delta_s \leq \mathbf{a}^\top(\omega_i)\mathbf{h} \leq \delta_s, && i = 1, \dots, N_s, \\ & && 1 - \delta_p \leq \mathbf{a}^\top(\omega_i)\mathbf{h} \leq 1 + \delta_p, && i = N_s + 1, \dots, N_s + N_p. \end{aligned}$$

同样地, 还可以对参数 $\mu > 0$ 取不同值, 来最小化 $\delta_p$ 与 $\delta_s$ 的加权值, 即

$$\begin{aligned} & \underset{\mathbf{h} \in \mathbb{R}^{N+1}, \delta_s \in \mathbb{R}, \delta_p \in \mathbb{R}}{\text{minimize}} && \delta_s + \mu\delta_p \\ & \text{subject to} && -\delta_s \leq \mathbf{a}^\top(\omega_i)\mathbf{h} \leq \delta_s, && i = 1, \dots, N_s, \\ & && 1 - \delta_p \leq \mathbf{a}^\top(\omega_i)\mathbf{h} \leq 1 + \delta_p, && i = N_s + 1, \dots, N_s + N_p. \end{aligned}$$

这里向读者介绍一个受约束于给定通带波纹 $\delta_p$ , 最小化阻带衰减水平 $\delta_s$ 的数值实例. 这里选取参数 $N = 10$ (因此滤波器的抽头数 $n = 2N + 1 =$

21),  $\Omega_p = 0.35\pi$ ,  $\Omega_s = 0.5\pi$ ,  $\delta_p = 0.02$ (即 $-33.98$  dB), 且以 $N_s = N_p = 100$ 将频率区间离散化, 线性等距取点. 求解相应的LP问题, 得到的最优阻带衰减 $\delta_s = 0.0285$ (即 $-30.9$  dB). 振幅响应曲线见图16.2, 相应的 $|H(\omega)|$ 的对数见图16.3, 滤波器系数(脉冲响应)见图16.4.

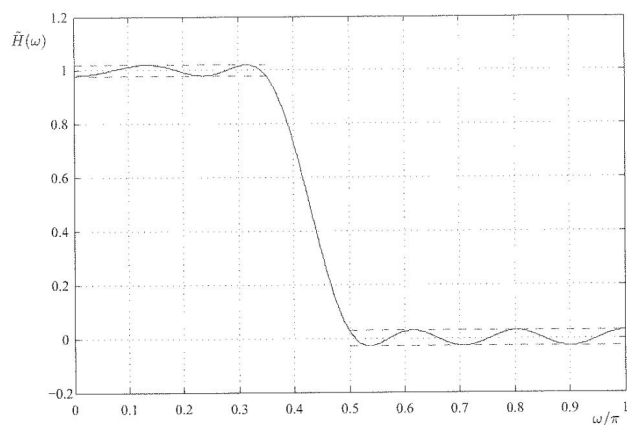


图 16.2: FIR滤波器振幅响应

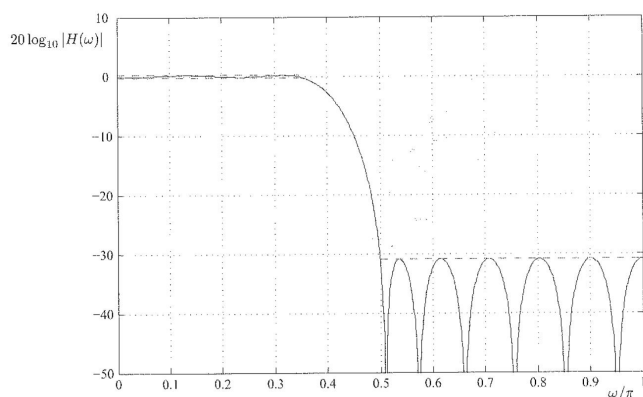


图 16.3: FIR滤波器的模作为(归一化)频率的函数

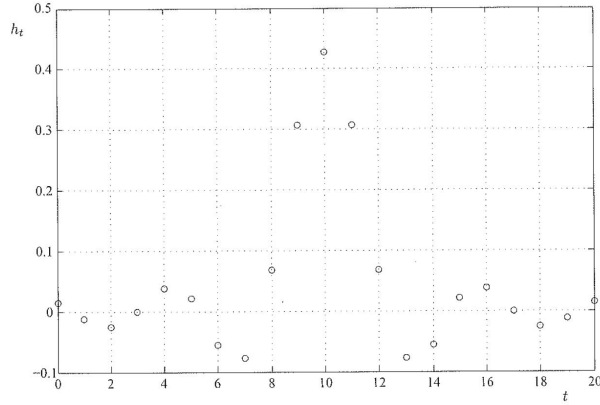


图 16.4: FIR滤波器的脉冲响应

#### §16.1.4 逼近参考匹配的滤波器设计

设计线性相位滤波器的一个可行方案, 是提供一个“参考的”理想幅度响应, 然后尝试寻找滤波器参数 $\mathbf{h}$ 使得滤波器响应与参考响应在所有频率处尽可能地接近.

##### §16.1.4.1 最小二乘设计

设 $\tilde{H}_{\text{ref}}(\omega)$ 是给定的振幅响应, 并且将频率区间 $[0, \pi]$ 离散化为 $\omega_i$ ,  $i = 1, \dots, M$ . 一种可能的方法是寻找滤波器系数 $\mathbf{h}$ 来极小化错配度量

$$\underset{\mathbf{h}}{\text{minimize}} \sum_{i=1}^M (\tilde{H}(\omega_i) - \tilde{H}_{\text{ref}}(\omega_i))^2.$$

因为 $\tilde{H}(\omega) = \mathbf{a}^\top(\omega)\mathbf{h}$ , 因此显然是一个最小二乘(least-square, LS)问题:

$$\underset{\mathbf{h}}{\text{minimize}} \|\mathbf{A}\mathbf{h} - \mathbf{b}\|_2^2,$$

这里

$$\mathbf{A} = \begin{bmatrix} \mathbf{a}^\top(\omega_1) \\ \vdots \\ \mathbf{a}^\top(\omega_M) \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} \tilde{H}_{\text{ref}}(\omega_1) \\ \vdots \\ \tilde{H}_{\text{ref}}(\omega_M) \end{bmatrix}.$$

这种方法的一个变形是对不同频率处的错配度量引入权重: 选定一组权重 $w_i \geq 0$ ,  $i = 1, \dots, M$ , 并且求解修正问题

$$\underset{\mathbf{h}}{\text{minimize}} \sum_{i=1}^M w_i^2 (\tilde{H}(\omega_i) - \tilde{H}_{\text{ref}}(\omega_i))^2.$$

在这种设置下,  $\omega_i$  的一个相对高的值意味着减小频率  $\omega_i$  处的错配是很重要的. 相反,  $\omega_i$  的一个相对小值则蕴含着  $\omega_i$  处的错配是不太重要的. 由这种方法得到加权LS问题

$$\underset{\mathbf{h}}{\text{minimize}} \quad \|\mathbf{W}(\mathbf{A}\mathbf{h} - \mathbf{b})\|_2^2, \quad \mathbf{W} = \text{diag}(w_1, \dots, w_M).$$

用一个实例来说明这种方法. 考虑一个  $N = 10$ ,  $M = 200$  的线性等距离散化频率的滤波器, 当频率小于通带截止频率  $\Omega_p = 0.35\pi$  时, 参考响应等于1; 当  $\omega > \Omega_s = 0.5\pi$  时参考响应等于0. 在参考幅度在过渡带  $[\omega_p, \omega_s]$  从1线性递减到0. 利用常数(单位)频率权重, 由LS 最优化问题得到的结果的振幅响应见图16.5, 相应的模图见图16.6.

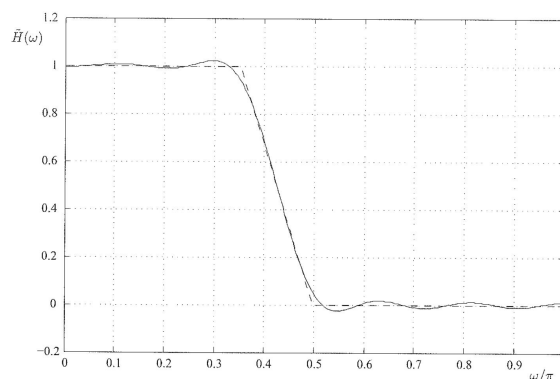


图 16.5: 由LS匹配得到的FIR滤波器的振幅响应, 虚线表示的是参考响应  $\tilde{H}_{\text{ref}}(\omega)$

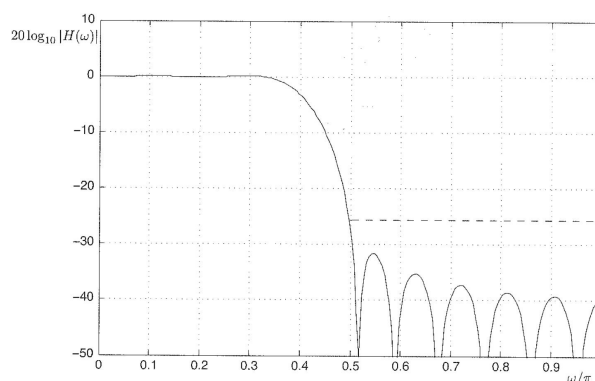


图 16.6: 由LS匹配得到的FIR滤波器的模, 虚线显示了阻带边沿  $\Omega_s = 0.5$  处的模水平, 其等于0.0518(即-25.71 dB)



## §16.1.4.2 切比雪夫设计

易于计算LS设计的解. 然而, 这种方法中我们没办法控制逐点的理想(参考)响应与滤波器实际的响应之间错配误差的最大值. 切比雪夫型设计便是将最小化最大加权错配作为目标, 即

$$\underset{\mathbf{h}}{\text{minimize}} \quad \max_{i=1,\dots,M} w_i |\tilde{H}(\omega_i) - \tilde{H}_{\text{ref}}(\omega_i)|.$$

这一问题可以表述成如下LP:

$$\begin{aligned} & \underset{\mathbf{h}, \gamma}{\text{minimize}} \quad \gamma \\ & \text{subject to} \quad w_i (\mathbf{a}^\top(\omega_i) \mathbf{h} - \tilde{H}_{\text{ref}}(\omega_i)) \leq \gamma, \quad i = 1, \dots, M, \\ & \quad \quad \quad w_i (\mathbf{a}^\top(\omega_i) \mathbf{h} - \tilde{H}_{\text{ref}}(\omega_i)) \geq -\gamma, \quad i = 1, \dots, M. \end{aligned}$$

将这一设计方法用于上文实例中的数据, 得到的振幅响应见图16.7, 对应的模图见图16.8. 由此得到的实际幅度响应与设定的振幅响应之间偏差的绝对值的最大者为  $\gamma = 0.0313$  (即-30.1 dB).

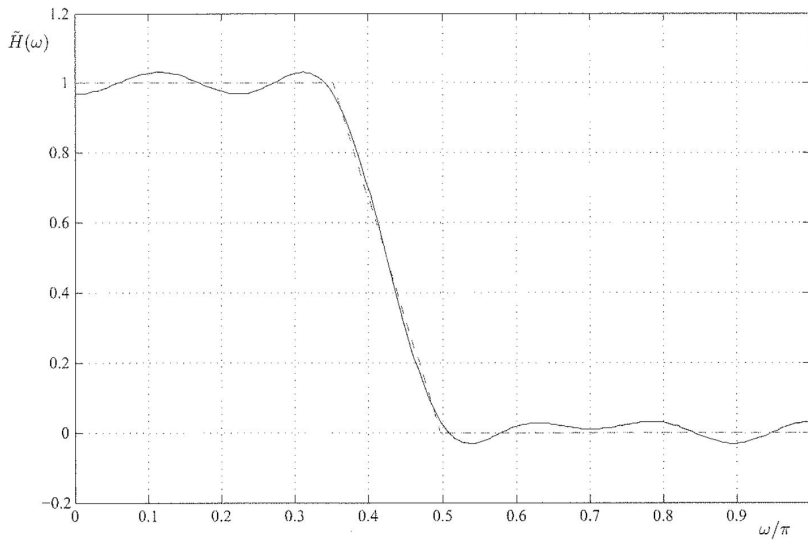


图 16.7: 由切比雪夫匹配得到的FIR滤波器的振幅响应, 虚线是参考响应  $\tilde{H}_{\text{ref}}(\omega)$

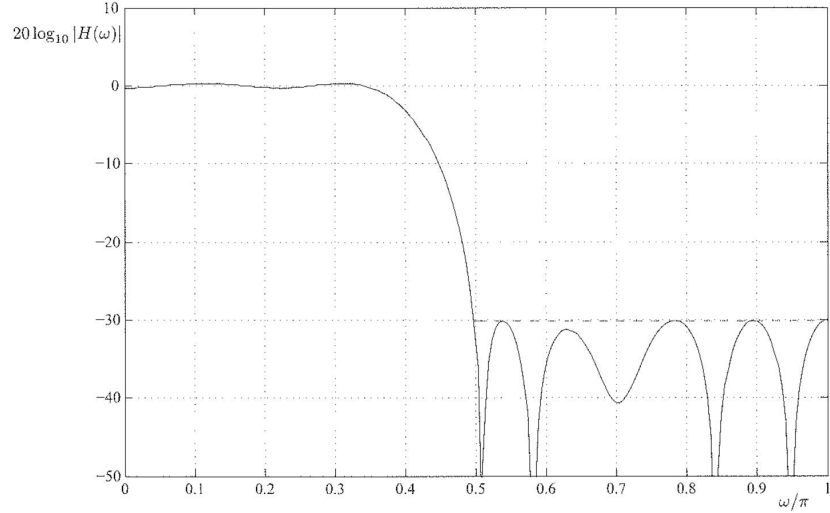


图 16.8: 由切比雪夫匹配得到的FIR滤波器的模, 虚线表示的是阻带边沿 $\Omega_s = 0.5$ 处的模水平, 其等于0.0313(即-30.1 dB).

## §16.2 天线阵列设计

在一个天线阵列中, 对若干发射天线单元的输出进行线性组合, 由此产生一个复合的阵列输出. 阵列输出具有方向图, 其取决于在组合过程中使用的相对权重或比例因子. 权重设计的目标是去选取权重, 以实现期望的方向图.

发射天线中的基本单元元件是各向同性的谐波振荡器, 其发射波长为 $\lambda$ , 频率为 $\omega$ 的球形单色波. 振荡器会产生电磁场, 其在离天线距离为 $d$ 的特定点 $\mathbf{p}$ 处的电分量为

$$\frac{1}{d} \operatorname{Re} \left( \mathbf{z} \cdot \exp \left( j \left( \omega t - \frac{2\pi d}{\lambda} \right) \right) \right),$$

其中 $\mathbf{z} \in \mathbb{C}$ 是一个设计参数, 用来缩放和改变电场的相位. 我们将这个复数称为天线单元的**权重**(weight). 现在, 我们放置 $n$ 个这样的振荡器, 分别在点 $\mathbf{p}_k \in \mathbb{R}^3$ ,  $k = 1, \dots, n$ , 处, 复值权重分别为 $\mathbf{z}_k \in \mathbb{C}$ ,  $k = 1, \dots, n$ . 这样, 在点 $\mathbf{p} \in \mathbb{R}^3$ 处的总电场强度由以下加权和给出:

$$E = \operatorname{Re} \left( \exp(j\omega t) \cdot \sum_{k=1}^n \frac{1}{d_k} \mathbf{z}_k \cdot \exp \left( \frac{-2\pi j d_k}{\lambda} \right) \right),$$

其中  $d_k = \|\mathbf{p} - \mathbf{p}_k\|_2$  是点  $\mathbf{p}$  到点  $\mathbf{p}_k$  的距离,  $k = 1, \dots, n$ .

**线性阵列的远场近似.** 前面的公式在以下两个假设下可以近似简化:

(a) 振荡器形成线性阵列, 即它们被放置在某一直线上的等距点上, 如果在这条直线上建立  $x$  轴, 那么这  $n$  个振荡器所处的网格点放置在  $x$ -轴上, 即点可分别表示为  $\mathbf{p}_k = \ell k \mathbf{e}_1$ ,  $k = 1, \dots, n$ , 其中  $\mathbf{e}_1 = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix}$  (即  $\mathbb{R}^3$  的第一标准单位基向量);

(b) 所考虑的点  $\mathbf{p}$  离原点足够远: 设  $\mathbf{p} = r\mathbf{u}$ , 其中  $\mathbf{u} \in \mathbb{R}^3$  是给定方向的单位范数的向量,  $r$  是点  $\mathbf{p}$  离原点的距离, 假定  $r$  足够大, 以至于可以忽略天线阵列的几何尺寸, 即  $r \gg n\ell$ , 参见图16.9.

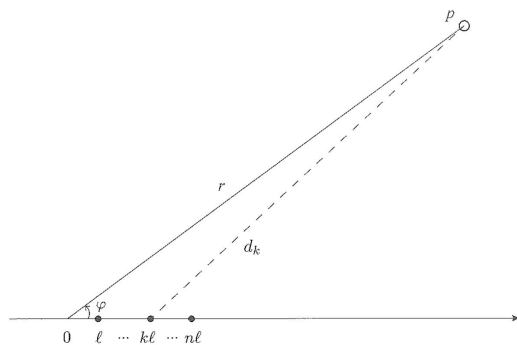


图 16.9: 线性天线阵列

对于线性阵列, 电场强度  $E$  仅由阵列所在直线和上述假设(b)中的远点  $\mathbf{p}$  与原点连线之间的夹角  $\varphi$  近似决定. 因此, 对足够小的  $\frac{k\ell}{r}$ , 我们有

$$d_k = r \sqrt{1 + \left(\frac{k\ell}{r}\right)^2 + 2\left(\frac{k\ell}{r}\right) \cos \varphi} \simeq r + k\ell \cos \varphi.$$

这样, 我们就得到了一个  $E$  的好的近似形式

$$E \simeq \frac{1}{r} \operatorname{Re} \left( \exp \left( j\omega t - \frac{2\pi jr}{\lambda} \right) \cdot D_{\mathbf{z}}(\varphi) \right),$$

其中函数  $D_{\mathbf{z}} : [0, 2\pi] \rightarrow \mathbb{C}$  被称为天线的**方向图**(diagram):

$$D_{\mathbf{z}}(\varphi) \doteq \sum_{k=1}^n \mathbf{z}_k \cdot \exp \left( \frac{-2\pi j k \ell \cos \varphi}{\lambda} \right). \quad (16.1)$$

我们使用下标“ $\mathbf{z}$ ”是为了强调方向图由复值权重向量  $\mathbf{z} = [\mathbf{z}_1, \dots, \mathbf{z}_n]$  决定.

类似的结果适用于**接收器**(receiver)天线的线性阵列, 如图16.10: 具有频率 $\omega$ 和波长 $\lambda$ 的平面谐波以方向角 $\varphi$ 入射并且通过阵列传播, 信号的输出被转换成基带(复数), 由权重 $\mathbf{z}_k$ 进行加权, 并求和, 再次给出线性阵列波束方向图 $D_{\mathbf{z}}(\varphi)$ .

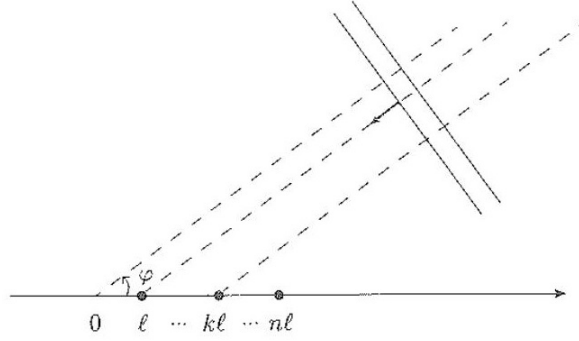


图 16.10: 一个接收器天线的线性阵列

### §16.2.1 天线方向图赋形

可证明, 天线方向图的模平方 $|D_{\mathbf{z}}(\varphi)|^2$ 与天线发送的电磁场能量的方向密度成比例. 因此, 我们感兴趣的是给天线方向图的模 $|D_{\mathbf{z}}(\cdot)|$ “赋形”(通过选择 $\mathbf{z}$ 来实现), 以便满足一些方向性要求. 注意到, 当 $\varphi$ 固定时,  $D_{\mathbf{z}}(\varphi)$ 是权重向量 $\mathbf{z}$ 的实部与虚部的线性函数. 特别地, 由式(16.1), 我们有

$$D_{\mathbf{z}}(\varphi) = \mathbf{a}^{\top}(\varphi) \mathbf{z}, \quad \mathbf{a}^{\top}(\varphi) = [a_1(\varphi) \cdots a_n(\varphi)],$$

$$a_k(\varphi) = \exp\left(\frac{-2\pi j k \ell \cos \varphi}{\lambda}\right), \quad k = 1, \dots, n.$$

进一步注意到,  $D_{\mathbf{z}}(\varphi)$ 的函数值是一个复数, 因此, 定义

$$\mathbf{a}_R(\varphi) = \text{Re}(\mathbf{a}(\varphi)), \mathbf{a}_I(\varphi) = \text{Im}(\mathbf{a}(\varphi)), \boldsymbol{\zeta} = \begin{bmatrix} \text{Re}(\mathbf{z}) \\ \text{Im}(\mathbf{z}) \end{bmatrix},$$

$$\mathbf{C}(\varphi) = \begin{bmatrix} \mathbf{a}_R^{\top}(\varphi) & -\mathbf{a}_I^{\top}(\varphi) \\ \mathbf{a}_I^{\top}(\varphi) & \mathbf{a}_R^{\top}(\varphi) \end{bmatrix},$$

我们得到

$$D_{\mathbf{z}}(\varphi) = \begin{bmatrix} \mathbf{a}_R^{\top}(\varphi) & -\mathbf{a}_I^{\top}(\varphi) \end{bmatrix} \boldsymbol{\zeta} + j \begin{bmatrix} \mathbf{a}_I^{\top}(\varphi) & \mathbf{a}_R^{\top}(\varphi) \end{bmatrix} \boldsymbol{\zeta},$$

$$|D_{\mathbf{z}}(\varphi)| = \|\mathbf{C}(\varphi) \boldsymbol{\zeta}\|_2.$$

一个典型的要求是天线沿着期望的方向（在给定角度上或其附近）很好地发射（或接收），而沿其它角度效果不佳. 以这种方式，发出的能量会集中在给定的“目标”方向上，比如说 $\varphi_{\text{target}} = 0^\circ$ ，并且在偏离该方向稍大的方向上能量较小. 另一类要求涉及由天线产生的热噪声功率.

**归一化.** 首先，我们归一化沿目标方向的能量. 当将所有权重乘以同一个非零常复数时，我们不改变能量的方向分布. 因此，通过归一化权重，即使得

$$D_{\mathbf{z}}(0) = 1,$$

我们没有任何损失. 该约束等价于决策变量 $\mathbf{z} \in \mathbb{C}^n$ 的实部和虚部中的两个线性等式约束：

$$\begin{bmatrix} \mathbf{a}_R^\top(0) & -\mathbf{a}_I^\top(0) \end{bmatrix} \boldsymbol{\zeta} = 1, \quad \begin{bmatrix} \mathbf{a}_I^\top(0) & \mathbf{a}_R^\top(0) \end{bmatrix} \boldsymbol{\zeta} = 0.$$

**旁瓣电平约束.** 我们接下来定义“通带” $[-\phi, \phi]$ ，其中 $\phi > 0$ 给定，在这个区间内，能量被集中；相应的“阻带”是该区间的外部. 为了满足能量被集中的需求，我们要求

$$|D_{\mathbf{z}}(\varphi)| \leq \delta, \quad \forall \varphi : |\varphi| \geq \phi,$$

其中 $\delta$ 是阻带上的期望衰减电平（有时也被称为旁瓣电平）. 事实上，旁瓣电平约束需要无穷多个约束. 处理这种连续无限多约束的一个实际可行方法是简单地离散化它们，即施加

$$|D_{\mathbf{z}}(\varphi_i)| \leq \delta, \quad i = 1, \dots, N.$$

其中 $\varphi_1, \dots, \varphi_N$ 是在阻带中取值的 $N$ 个间隔均匀的离散角度. 这是 $N$ 个关于 $\mathbf{z}$ 的实部和虚部的二次锥(second-order cone, SOC) 约束：

$$\|\mathbf{C}(\varphi_i) \boldsymbol{\zeta}\|_2 \leq \delta, \quad i = 1, \dots, N.$$

举个例子，如图16.11所示，天线方向图的模必须通过右侧 $\phi = 0^\circ$ 上的点，否则将被包含在白色区域中. 在阻带（阴影区域）中，至少在离散点处，天线方向图的模必须保持低于 $\delta$ .

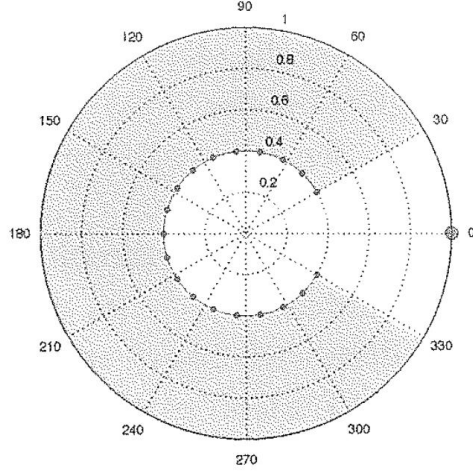


图 16.11: 关于天线方向图的约束

**热噪声功率约束.** 通常我们还希望控制由发射天线产生的热噪声功率. 结果证明, 该功率与 (复) 向量  $\mathbf{z}$  的欧几里得范数的平方成正比, 即:

$$\text{热噪声功率} = \alpha \|\mathbf{z}\|_2^2 = \alpha \sum_{i=1}^n |z_i|^2.$$

### §16.2.2 最小二乘设计法

第一种简化方法是最小二乘法, 其通过考虑旁瓣电平衰减和热噪声功率之间的权衡来解决天线设计问题. 为了取代在每个离散角度  $\varphi_i$  处对阻带水平施加的约束, 我们在目标函数中添加一个总阻带电平平方和的惩罚项, 即考虑问题

$$\underset{\mathbf{z}}{\text{minimize}} \quad \|\mathbf{z}\|_2^2 + \mu \sum_{i=1}^N |D_{\mathbf{z}}(\varphi_i)|^2, \quad \text{subject to } D_{\mathbf{z}}(0) = 1,$$

其中  $\mu \geq 0$  是权衡参数. 这是一个有等式约束的LS问题. 利用含  $\mathbf{z}$  的实部与虚部的变量  $\boldsymbol{\zeta}$ , 可以把该问题更明晰地表示为:

$$\underset{\boldsymbol{\zeta}}{\text{minimize}} \quad \|\boldsymbol{\zeta}\|_2^2 + \mu \boldsymbol{\zeta}^\top \mathbf{A} \boldsymbol{\zeta}, \quad \text{subject to } \mathbf{C}(0) \boldsymbol{\zeta} = \begin{bmatrix} 1 \\ 0 \end{bmatrix},$$

其中我们定义

$$\mathbf{A} = \sum_{i=1}^N \mathbf{C}^\top(\varphi_i) \mathbf{C}(\varphi_i).$$

由于 $\mathbf{A}$ 是半正定的, 所以该问题是一个凸二次规划问题. 通过分解 $\mathbf{A} = \mathbf{F}^\top \mathbf{F}$ , 我们进一步可以把该问题写成形如

$$\min_{\boldsymbol{\zeta}} \left\| \begin{bmatrix} \mathbf{I} \\ \sqrt{\mu} \mathbf{F} \end{bmatrix} \boldsymbol{\zeta} \right\|_2^2, \quad \text{subject to } \mathbf{C}(0) \boldsymbol{\zeta} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

的最小二乘形式. 注意, 这种惩罚的方法不需要知道阻带电平的一个先验(a priori)期望阈值 $\delta$ . 我们仅希望, 当 $\mu$ 足够大时, 平方和中的所有项将小于期望的临界值 $\delta$ . 因此, 设计所达到的阻带衰减电平只能通过后验检查确定.

考虑一个数值例子, 我们设定以下参数: 天线数目 $n = 16$ , 波长 $\lambda = 8$ , 通带宽度 $\phi = \frac{\pi}{6}$ , 天线的间距 $\ell = 1$ , 离散的角度数目 $N = 100$ , 权衡参数 $\mu = 0.5$ . 使用CVX求解所得到的解见图16.12.

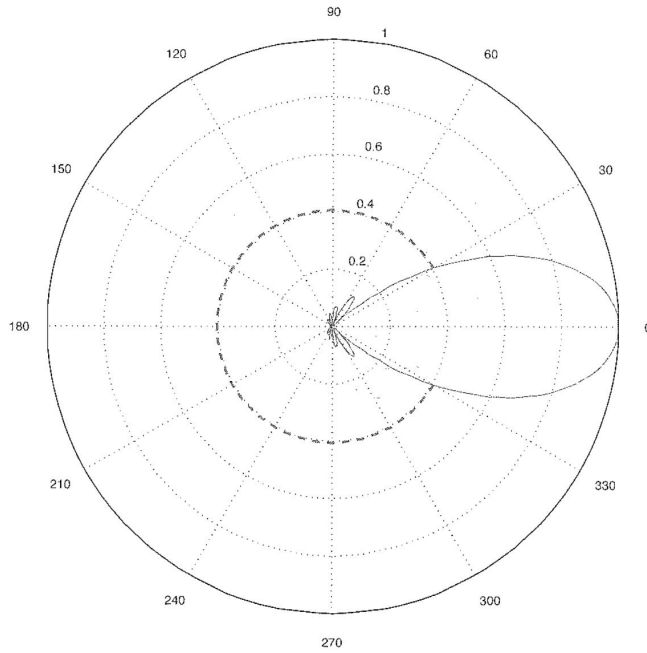
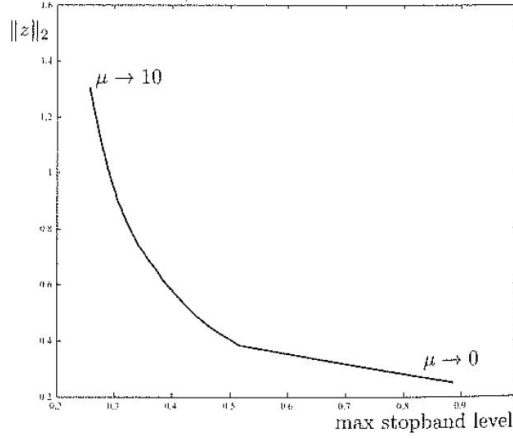


图 16.12: 由最小二乘的解得到的天线方向图

热噪声功率(平方根)是 $\|\mathbf{z}\|_2 = 0.5671$ , 最大阻带电平是0.4050 (在图16.12中用虚线圆弧来突出显示). 此外, 我们多次增加 $\mu$ 的值并重复求解该问题, 绘制出了热噪声功率的平方根 $\|\mathbf{z}\|_2$ 与阻带衰减电平之间的相应的权衡曲线, 见图16.13.

图 16.13:  $\mu$  在区间  $[0, 10]$  上变动时对应的权衡曲线

### §16.2.3 SOCP设计法

实际的天线设计问题，对离散角度上的阻带电平有明确的约束，可以直接表述为二次锥规划(second-order cone programming, SOCP). 一种可能的方法是在旁瓣电平约束下最小化热噪声功率，由此得到以含有  $z$  的实部和虚部的向量  $\zeta$  为变量的显式SOCP:

$$\begin{aligned}
 & \underset{\zeta \in \mathbb{R}^{2n}, \gamma}{\text{minimize}} && \gamma \\
 & \text{subject to} && C(0) \zeta = \begin{bmatrix} 1 & 0 \end{bmatrix}^T, \\
 & && \|C(\varphi_i) \zeta\|_2 \leq \delta, \quad i = 1, \dots, N, \\
 & && \|\zeta\|_2 \leq \gamma.
 \end{aligned} \tag{16.2}$$

或者，可以通过给定热噪声功率的阈值  $\gamma$  来最小化旁瓣电平衰减  $\delta$ . 更进一步，我们可以通过考虑利用如下SOCP来求出阻带衰减和噪声之间的最佳权衡:

$$\begin{aligned}
 & \underset{\zeta \in \mathbb{R}^{2n}, \gamma, \delta}{\text{minimize}} && \delta + w\gamma \\
 & \text{subject to} && C(0) \zeta = \begin{bmatrix} 1 & 0 \end{bmatrix}^T, \\
 & && \|C(\varphi_i) \zeta\|_2 \leq \delta, \quad i = 1, \dots, N, \\
 & && \|\zeta\|_2 \leq \gamma.
 \end{aligned} \tag{16.3}$$

其中  $w \geq 0$  是给定的权衡参数.

作为数值示例，求解  $\delta = 0.35$  时的问题(16.2)，得到  $\|z\|_2 = 0.435$ ，相应的方向图见图16.14.



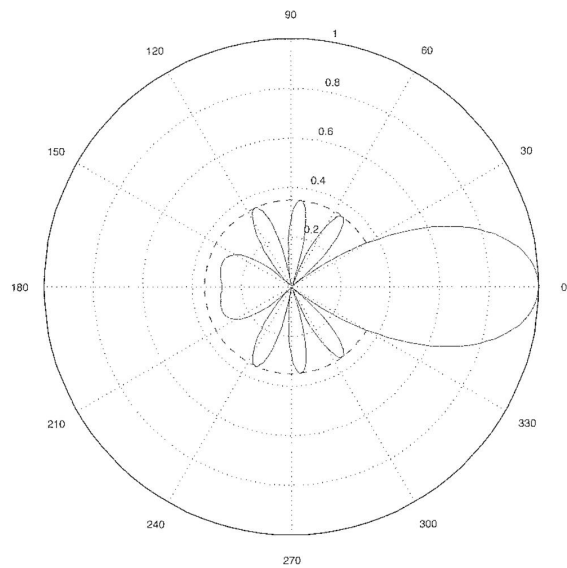


图 16.14:  $\delta = 0.35$ 时求解问题(16.2)所的结果的天线方向图

同样，我们多次增加 $w$ 的值并重复求解问题(16.3)，得到权衡曲线. 图16.15显示了 $w \in [0.2, 10]$ 的权衡曲线.

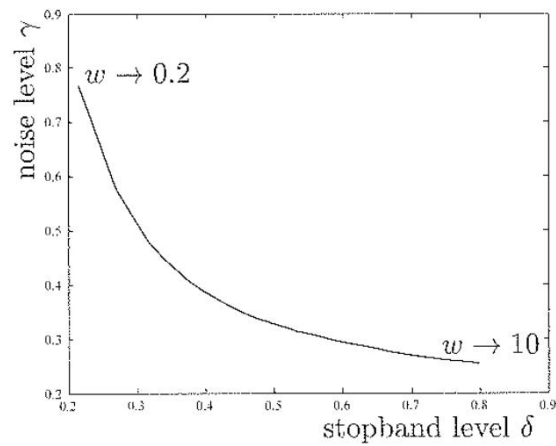


图 16.15:  $w$ 在区间 $[0.2, 10]$ 上变动时的权衡曲线

§16.3 数字电路设计

我们把数字电路设计问题考虑为一个组合逻辑模块. 在这样一个电

路中, 基本构建模块是门: 一个执行着一些简单布尔函数的电路, 例如**逆变器**(inverter), 它表现为逻辑非, 或者一些更为复杂的函数, 如“与非(not-AND或NAND)”, 它把两个布尔输入值 $A, B$ 转换为“ $A$ 且 $B$ ”整体的非. 数字电路设计的基本思想是设计门(确定门的大小), 使得电路运行快而且占据的面积小. 设计问题中还涉及其他因素, 比如功率, 但在这里我们不考虑这些. 设计变量是比例因子, 它决定了每一个门的大小, 同时也是门的基本电学参数. 这些参数连同电路的拓扑(这里固定)一同影响电路的速度. 这一节, 我们将用几何规划(geometric programming, GP)模型来处理电路设计问题. 电路设计中GP在有一个很长的历史.<sup>1</sup>

### §16.3.1 电路拓扑

组合电路由有主输入和输出的连通门组成. 我们假设在相应的电路拓扑图中没有环. 对于每一个门, 我们可以定义**扇入**(fan-in), 即电路图中该门的前驱的集合, 以及**扇出**(fan-out), 它是该门的后继的集合. 电路由 $n$ 个“门”(由一些输入和一个输出的逻辑模块)组成, 这些门连接着主输入(在图16.16中标记为 $\{8, 9, 10\}$ )和主输出(在图中标记为 $\{11, 12\}$ ). 每一个门都用一个指定其类型的符号代表, 例如, 标记为 $\{1, 3, 6\}$ 的门是逆变器. 对于这个电路, 门4的扇入和扇出分别是

$$FI(4) = \{1, 2\}, \quad FO(4) = \{6, 7\}.$$

根据定义, 主输入扇入为空, 主输出扇出为空.

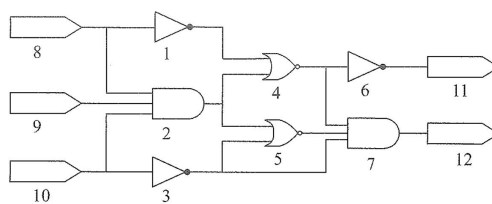


图 16.16: 一个数字电路的例子

### §16.3.2 设计变量

我们模型中的设计变量是**比例因子** $x_i$ ,  $i = 1, \dots, n$ , 它们粗略的决定了每一个门的尺寸. 这些比例因子满足 $x_i \geq 1$ ,  $i = 1, \dots, n$ , 其中 $x_i = 1$ 表

<sup>1</sup>这里给出的介绍以及示例来自论文Digital circuit optimization via geometric programming(S. Boyd, S.-J. Kim, D. Patil, and M. Horowitz, *Operations Research*, 2005), 其中也包含了许多参考文献.

示对应门的尺寸最小, 当一个比例因子  $x_i = 16$  时表示门中所有装置的总宽度是最小门的16倍. 比例因子决定了门的大小和各种电学特征, 例如门的电阻和电导. 它们之间的关系可以被很好的近似如下.

- 门  $i$  的面积  $A_i(\mathbf{x})$  与比例因子  $x_i$  成正比, 即  $A_i(\mathbf{x}) = a_i x_i$  对某个  $a_i > 0$  成立.

- 门  $i$  的固有电容形如

$$C_i^{\text{intr}}(\mathbf{x}) = C_i^{\text{intr}} x_i,$$

其中  $C_i^{\text{intr}}$  是正系数.

- 门  $i$  的负载电容是关于门  $i$  的扇出中门的比例因子的线性函数:

$$C_i(\mathbf{x}) = \sum_{j \in \text{FO}(i)} C_j x_j,$$

其中  $C_j$  是正系数.

- 每一个门都有电阻, 它与比例因子成反比(门越大, 通过它的电流越大):

$$R_i(\mathbf{x}) = \frac{r_i}{x_i},$$

其中  $r_i$  是正系数.

- 门延迟是衡量门执行它应该执行的逻辑操作的速度的量, 这个延迟可以近似为

$$D_i(\mathbf{x}) = 0.7 R_i(\mathbf{x}) (C_i^{\text{intr}}(\mathbf{x}) + C_i(\mathbf{x})).$$

我们可以看到, 上述所有量都是关于(正的)设计向量  $\mathbf{x}$  的正多项式函数.

### §16.3.3 设计目标

一个可能的设计目标是最小化电路的总延迟  $D$ . 总延迟可表示为

$$D = \max_{1 \leq i \leq n} T_i, \quad (16.4)$$

其中  $T_i$  表示门  $i$  直到成功输出一共所消耗的时间, 并假设主输入信号在  $t = 0$  时传递. 即,  $T_i$  是在所有始于主输入并终于门  $i$  的所有路上延迟的最大值. 我们可以递归地表示  $T_i$  为

$$T_i = \max_{j \in \text{FI}(i)} (T_j + D_i). \quad (16.5)$$

$D$ 的计算过程仅包括加法和逐点取最大. 由于每一个 $D_i$ 是 $\mathbf{x}$ 的正多项式, 我们可以将总延迟表示为一个 $\mathbf{x}$ 的广义正多项式. 对于图16.16中的电路, 总延迟 $D$ 可表示为

$$\begin{aligned} T_i &= D_i, \quad i = 1, 2, 3; \\ T_4 &= \max(T_1, T_2) + D_4; \\ T_5 &= \max(T_2, T_3) + D_5; \\ T_6 &= T_4 + D_6; \\ T_7 &= \max(T_3, T_4, T_5) + D_7; \\ D &= \max(T_6, T_7). \end{aligned}$$

#### §16.3.4 一个电路设计问题

我们现在考虑问题: 在面积约束条件下, 选择比例因子 $x_i$ 使得总延迟最小, 即

$$\underset{\mathbf{x}}{\text{minimize}} D(\mathbf{x}) \quad \text{subject to} \quad A_i(\mathbf{x}) \leq A_{\max}, \quad x_i \geq 1, \quad i = 1, \dots, n,$$

其中 $A_{\max}$ 是每个门的面积上限. 因为 $D$ 是 $\mathbf{x}$ 的广义正多项式, 所以上述问题可以被看做是一个GP问题. 为了找到一个紧凑且明确的GP表示, 我们可以使用在总延迟的定义中出现的中间变量 $T_i$ . 通过观察可知, 由不等式关系

$$T_i \geq \max_{j \in \text{FI}(i)} (T_j + D_i) \quad (16.6)$$

代替等式关系(16.5)不改变优化问题的最优解. 其原因为: 由于式(16.4)中的目标 $D$ 相对于 $T_i$ 是非递减的, 并且 $T_i$ 随着 $T_j, j \in \text{FI}(i)$ 的增加而增加, 所以最优解将选择所有尽可能小的 $T_i$ , 因此(16.6)是以等式成立的. 更进一步, 上述不等式等价于

$$T_i \geq T_j + D_i, \quad \forall j \in \text{FI}(i).$$

之后, 我们得到设计问题的如下显式GP表述:

$$\begin{aligned} &\underset{\mathbf{x}, T_i > 0, D}{\text{minimize}} && D \\ &\text{subject to} && A_i(x) \leq A_{\max}, \quad x_i \geq 1, \quad i = 1, \dots, n; \\ & && D \geq T_i, \quad i = 1, \dots, n; \\ & && T_i \geq T_j + D_i(x), \quad \forall j \in \text{FI}(i), \quad i = 1, \dots, n. \end{aligned}$$

我们考虑图16.16中的电路给出的数值实例，其中 $A_{\max} = 16$ ，其他数据见表16.1.

门	$C_i$	$C_i^{\text{intr}}$	$r_i$	$a_i$
1,3,6	3	3	0.48	3
2,7	4	6	0.48	8
4,5	5	6	0.48	10

表 16.1: 图16.16中电路的数据

利用CVX求得最优比例因子

$$\mathbf{x}^* = \begin{bmatrix} 4.6375 \\ 2.0000 \\ 4.8084 \\ 1.6000 \\ 1.0000 \\ 1.0000 \\ 1.0000 \end{bmatrix},$$

对应的最小延迟 $D^* = 7.686$ .

#### §16.4 航天器设计

近年来，我们已经注意到，一些关于航空器结构与运行方面的问题能够表述为几何规划的形式，从而，通过凸优化被高效的求解. 我们在这里展示一个简单的机翼设计问题的例子.<sup>2</sup>

我们的设计变量是总面积 $S$ ，机翼展弦比 $A = b^2/S$ ，其中 $b$ 为翼展，和巡航速度 $V(\text{m/s})$ ，目标是最小化阻力

$$D = \frac{1}{2}\rho V^2 C_D S, \quad (16.7)$$

其中 $\rho$ 是空气密度， $C_D$ 是阻力系数，见图16.17.

<sup>2</sup>我们的陈述来自论文Geometric programming for aircraft design optimization, by W. Hoburg and P. Abbeel, in proc. *Structures, Structural Dynamics and Materials Conference*, 2012, 读者可以阅读该文以获得更多细节.

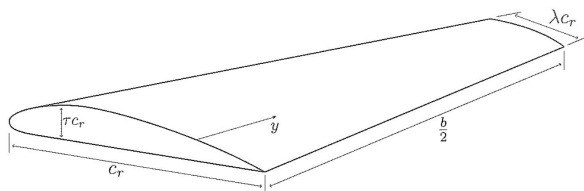


图 16.17: 单锥形机翼几何:  $c_r$  是根弦,  $b/2$  是半翼展,  $\lambda$  是弦锥因子,  $\tau$  是翼型厚弦比.

当飞机稳定飞行时, 它必须满足两个基本的平衡条件, 正像图16.18所展示的: 升力 $L$ 必须抵消飞机的重力 $W$ , 并且推力 $T$ 与阻力 $D$ 必需是平衡的, 即

$$L = W, \quad T = D.$$

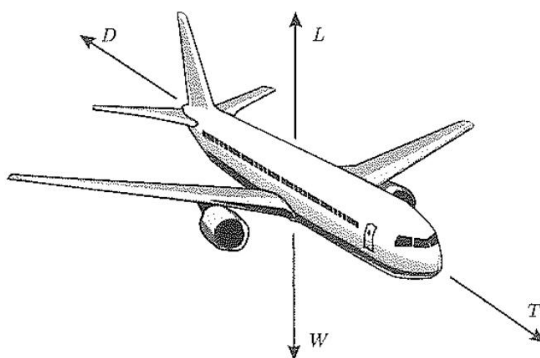


图 16.18: 在稳态飞行下, 升力 $L$ 等于重力 $W$ , 推力 $T$ 等于阻力 $D$

用 $C_L$ 表示升力系数, 升力可以表示为

$$L = \frac{1}{2} \rho V^2 C_L S.$$

式(16.7)中的阻力系数可建模为三项的和: 机身寄生阻力, 机翼寄生阻力, 诱导阻力, 即

$$C_D = \frac{CDA_0}{S} + \kappa C_f \frac{S_{\text{wet}}}{S} + \frac{C_L^2}{\pi A e}, \quad (16.8)$$

其中 $CDA_0$ 是机身阻力面积,  $\kappa$ 是导致压差阻力的形式因子,  $S_{\text{wet}}$ 是受潮面积(即机翼曲面的实际面积, 而 $S$ 是二维投影面积或是机翼阴影面积),  $C_L$ 是

升力系数,  $e$  是Oswald效率因子. 对于一个完全湍流边界层, 式(16.8)中的表面摩擦系数  $C_f$  可以粗略地近似为

$$C_f = \frac{0.074}{\text{Re}^{0.2}},$$

其中

$$\text{Re} = \frac{\rho V}{\mu} \sqrt{S/A}$$

是平均翼弦<sup>3</sup>  $c = \sqrt{S/A}$  处的雷诺数,  $\mu$  是空气粘度. 航空器的总重量  $W$  被建模为一个固定重量  $W_0$  与机翼重量  $W_w$  之和, 后者被建模为

$$W_w = \kappa_s S + \kappa_l \frac{N_{\text{ult}} b^3 \sqrt{W_0 W}}{S \tau},$$

这里  $\kappa_s, \kappa_l$  是合适的常数,  $N_{\text{ult}}$  是结构尺寸决定的极限载荷因子,  $\tau$  是翼型厚弦比(注意到  $W_w$  是  $W$  自己的函数). 在稳态飞行条件下, 约束将重力方程耦合到阻力方程中, 升力必须要等于重力, 即  $L = W$ . 因此, 必须满足

$$\begin{aligned} \frac{1}{2} \rho V^2 C_L S &= W, \\ W_0 + W_w &= W. \end{aligned}$$

最后, 飞机在落地时, 必须没有拖延地能以最小速度  $V_{\min}$  飞行. 这个要求由下面的限制得以满足:

$$\frac{1}{2} \rho V_{\min}^2 C_L^{\max} S \geq W,$$

这里,  $C_L^{\max}$  是着陆(放下副翼)时的最大升力系数.

我们必须选择  $S, A$  和  $V$  的值, 使得阻力最小化且满足上述**所有**关系. 表16.2给定了常量参数. 这个问题看起来是一个很难的优化问题, 涉及非线性耦合变量. 然而, 这个问题中的一个关键点是: 我们能够将正项多项式**等式**约束松弛成正项多项式**不等式**约束, 且保持问题的最优目标值不变. 比如, 如果  $C_D$  没有出现在其他任何单项等式限制中, 且目标与不等式限制

$$C_D \geq \frac{C D A_0}{S} + k C_f \frac{S_{\text{wet}}}{S} + \frac{C_L^2}{\pi A e} \quad (16.9)$$

关于  $C_D$  都是单调递增的(或者是常数), 则我们可用不等式关系(16.9)来等价替换等式(16.8). 在这些条件下, 如果等式关系(16.8)在最优解处不成立, 显然, 我们能减小  $C_D$  直到等式满足, 而不会增大目标值或是将解移

<sup>3</sup>平均翼弦是区间  $[\lambda c_r, c_r]$  中且满足  $S = bc$  的值  $c$ .

表 16.2: 航天器设计实例中的常数

量	值	单位	描述
$CDA_0$	0.0306	$m^2$	机身阻力面积
$\rho$	1.23	$Kg/m^3$	空气密度
$\mu$	$1.78 \times 10^{-5}$	$Kg/ms$	空气粘度
$S_{wet}/S$	2.05		受潮面积比
$\kappa$	1.2		形式因子
$e$	0.96		Oswald效率因子
$W_0$	4940	N	除机翼外的机身重量
$N_{ult}$	2.5		极限载荷因子
$\tau$	0.12		翼型厚弦比
$V_{min}$	22	m/s	着陆速度
$C_L^{max}$	2.0		(放下副翼)时 $C_L$ 的最大值
$\kappa_s$	45.42	$N/m^2$	
$\kappa_l$	$8.71 \times 10^{-5}$	$m^{-1}$	

出可行集. 因为在目前的背景下这些条件都是满足的, 所以我们能够将设计问题写成如下显示的GP形式:

$$\begin{aligned}
 & \underset{S, A, V, C_D, C_L, C_f, Re, W, W_w}{\text{minimize}} && \frac{1}{2} \rho V^2 C_D S \\
 & \text{subject to} && \frac{0.074}{C_f Re^{0.2}} = 1, \quad \frac{2W}{\rho V^2 C_L S} = 1, \\
 & && \frac{2W}{\rho V_{min}^2 C_L^{max} S} \leq 1, \quad \frac{\rho V}{\mu Re} \sqrt{S/A} = 1, \\
 & && \frac{CDA_0}{C_D S} + \kappa \frac{C_f}{C_D} \frac{S_{wet}}{S} + \frac{C_L^2}{C_D \pi A e} \leq 1, \\
 & && \kappa_s \frac{S}{W_w} + \kappa_l \frac{N_{ult} A^{3/2} \sqrt{W_0 W S}}{W_w \tau} \leq 1, \\
 & && \frac{W_0}{W} + \frac{W_w}{W} \leq 1.
 \end{aligned} \tag{16.10}$$

常量参数见表16.2. 求解GP(16.10)得到的最优设计结果见表16.3. 值得指出的是, 这里的 $S, A, V$ 是独立设计变量, 其余的变量都是由这三个独立变量和常量参数决定的.

GP模型使我们很容易地得到设计参数的全局最优值. 但是, 在一个现实的设计背景下, 设计者需要考虑在相互冲突的目标间考虑一系列可能的折衷解(比如增加着陆速度和(或)巡航速度). 对一定范围内的一系列 $V_{min}$ 和 $V$ 的值, 求解GPGP(16.10)很易于得到这些数值折衷解. 并且, 最



表 16.3: 航天器实例中的最有设计

参数	结果	单位	描述
$A$	12.7		机翼展弦比
$S$	12.08	$\text{m}^2$	机翼表面积
$V$	38.55	$\text{m/s}$	巡航速度
$C_D$	0.0231		阻力系数
$C_L$	0.6513		升力系数
$C_f$	0.0039		表面摩擦系数
$\text{Re}$	$2.5978 \times 10^6$		雷诺数
$W$	7189.1	N	总重量
$W_w$	2249.1	N	机翼重量

优化模型中也包含了几其他方面. 例如由模型可以计算出燃料的质量, 即通过

$$W = (W_0 + W_w)(1 + \theta_{\text{fuel}}),$$

其中 $\theta_{\text{fuel}}$ 是燃料的质量比. 进一步, Brequet的航程方程

$$R = \frac{h_{\text{fuel}}}{g} \eta_0 \frac{L}{D} \log(1 + \theta_{\text{fuel}}) \quad (16.11)$$

描述了燃料质量比与航空器的最大航程之间的关系, 这里假定升阻比 $L/D$ 是常数,  $\eta_0$ 是全体燃料能量对推力能量的效率因子. 等式 $P_{\text{fuel}} = \dot{m}_{\text{fuel}} h_{\text{fuel}}$ 描述了 $h_{\text{fuel}}$ 与燃料质量流速率与燃料能量的关系. 在稳态飞行条件下, 必须满足

$$TV \leq \eta_0 P_{\text{fuel}}.$$

虽然Brequet的航程方程(16.11)并不是直接表述成正的多项式, 但可以用泰勒级数展开来很好地逼近它. 由(16.11)解得

$$1 + \theta_{\text{fuel}} = \exp\left(\frac{gRD}{h_{\text{fuel}}\eta_0 L}\right),$$

且指数函数的级数展式具有正项式结构. 因此, 利用

$$z = \frac{gRD}{h_{\text{fuel}}\eta_0 L},$$

$$\theta_{\text{fuel}} \geq z + \frac{z^2}{2!} + \frac{z^3}{3!} + \cdots$$

可将Brequet的航程方程近似地引入到GP模型中.

## §16.5 供应链管理

本节讨论一个生产工程中出现的问题,即在不确定的需求下如何控制库存量.<sup>4</sup>这个问题由 $T$ 个时段的订购、库存和存储决策组成,目标是在满足需求的前提下最小化成本.总成本由实际采购成本、持有与缺货带来的成本组成.考虑一个单一商品,此商品在时段 $k$ 的库存量记作 $x(k)$ .那么,商品库存量随着(离散)时刻的基本演化可以写作:

$$x(k+1) = x(k) + u(k) - w(k), \quad k = 0, 1, \dots, T-1,$$

这里 $x(0) = x_0$ 是一个给定的初始库存量, $u(k)$ 是时段 $k$ 的采购量, $w(k)$ 是时段 $k$ 到 $k+1$ 的需求量.我们假设一个单位的仓储价格是 $h$ ,一个单位的短缺成本是 $p$ ,购买单位商品的价格是 $c$ .那么时段 $k$ 的成本由

$$cu(k) + \max(hx(k+1), -px(k+1))$$

确定.进一步假设任一次订货规模的上界是 $M$ ,我们可以将 $T$ 个时段的库存控制问题写作:

$$\begin{aligned} & \underset{u(0), \dots, u(T-1)}{\text{minimize}} && \sum_{k=0}^{T-1} cu(k) + \max(hx(k+1), -px(k+1)) \\ & \text{subject to} && 0 \leq u(k) \leq M, \quad k = 0, \dots, T-1, \end{aligned}$$

这里

$$x(k) = x_0 + \sum_{i=0}^{k-1} (u(i) - w(i)), \quad k = 1, \dots, T.$$

引入松弛变量 $y(0), \dots, y(T-1)$ ,此问题可表述成如下的线性规划:

$$\begin{aligned} & \underset{u(0), \dots, u(T-1), y(0), \dots, y(T-1)}{\text{minimize}} && \sum_{k=0}^{T-1} y(k) \\ & \text{subject to} && cu(k) + hx(k+1) \leq y(k), \quad k = 0, \dots, T-1, \\ & && cu(k) - px(k+1) \leq y(k), \quad k = 0, \dots, T-1, \\ & && x(k) = x_0 + \sum_{i=0}^{k-1} (u(i) - w(i)), \quad k = 1, \dots, T, \\ & && 0 \leq u(k) \leq M, \quad k = 0, \dots, T-1. \end{aligned}$$

接着,我们定义向量 $\mathbf{u} = (u(0), \dots, u(T-1))$ ,  $\mathbf{y} = (y(0), \dots, y(T-1))$ ,  $\mathbf{w} =$

<sup>4</sup>这里展示的处理方法所采用的设置与记号源于Bertsimas 和Thiele的论文: A robust optimization approach to supply chain management, *Operation Research*, 2006.

$(w(0), \dots, w(T-1)), \mathbf{x} = (x(0), \dots, x(T))$ . 可用紧凑的记号将问题重写为:

$$\begin{aligned} & \underset{\mathbf{u}, \mathbf{y}}{\text{minimize}} && \mathbf{1}^\top \mathbf{y} \\ & \text{subject to} && c\mathbf{u} + h\mathbf{x} \leq \mathbf{y}, \quad c\mathbf{u} - p\mathbf{x} \leq \mathbf{y}, \\ & && \mathbf{x} = x_0\mathbf{1} + \mathbf{U}\mathbf{u} - \mathbf{U}\mathbf{w}, \\ & && \mathbf{0} \leq \mathbf{u} \leq M\mathbf{1}, \end{aligned} \quad (16.12)$$

其中

$$\mathbf{U} \doteq \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 1 & 1 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 1 & 1 & 1 & \cdots & 1 \end{bmatrix}.$$

在这种设置下, 我们假设在所有时段  $k = 0, \dots, T-1$ , 需求  $w(k)$  是确定且已知的, 并且时段  $k = 0$  的所有决策都是已知的. 在下一节中, 我们将讨论如何放松这两个假设.

### §16.5.1 针对区间不确定需求的鲁棒法

这里我们假设可能精准地知道需求  $\mathbf{w}$ . 特别地, 考虑需求在名义上的预期需求  $\hat{\mathbf{w}} \geq \mathbf{0}$  周围有  $\rho\%$  的波动的情况, 即

$$\mathbf{w} \in W \doteq \{\mathbf{w} : \mathbf{w}_{\text{lb}} \leq \mathbf{w} \leq \mathbf{w}_{\text{ub}}\}, \quad (16.13)$$

其中

$$\mathbf{w}_{\text{lb}} \doteq (1 - \rho/100)\hat{\mathbf{w}}, \quad \mathbf{w}_{\text{ub}} \doteq (1 + \rho/100)\hat{\mathbf{w}}.$$

然后, 我们寻求一个订货序列  $\mathbf{u}$ , 使得在所有可能的需求下的最坏情况的成本最小, 即

$$\begin{aligned} & \underset{\mathbf{u}, \mathbf{y}}{\text{minimize}} && \mathbf{1}^\top \mathbf{y} \\ & \text{subject to} && c\mathbf{u} + h\mathbf{x} \leq \mathbf{y}, \quad \forall \mathbf{w} \in W, \\ & && c\mathbf{u} - p\mathbf{x} \leq \mathbf{y}, \quad \forall \mathbf{w} \in W, \\ & && \mathbf{x} = x_0\mathbf{1} + \mathbf{U}\mathbf{u} - \mathbf{U}\mathbf{w}, \quad \forall \mathbf{w} \in W, \\ & && \mathbf{0} \leq \mathbf{u} \leq M\mathbf{1}. \end{aligned} \quad (16.14)$$

因为  $\mathbf{x} = x_0\mathbf{1} + \mathbf{U}\mathbf{u} - \mathbf{U}\mathbf{w}$ , 其中  $\mathbf{U}$  的元素都是非负的, 并且  $\mathbf{w}$  中的每个元素是属于一个已知区间的, 约束  $c\mathbf{u} + h\mathbf{x} \leq \mathbf{y}, \forall \mathbf{w} \in W$  成立当且仅当

$$c\mathbf{u} + h(x_0\mathbf{1} + \mathbf{U}\mathbf{u} - \mathbf{U}\mathbf{w}_{\text{lb}}) \leq \mathbf{y},$$

约束  $cu - px \leq y, \forall w \in W$  成立当且仅当

$$cu - p(x_0\mathbf{1} + Uu - Uw_{ub}) \leq y.$$

因此可以通过求解LP

$$\begin{aligned} & \underset{u, y}{\text{minimize}} && \mathbf{1}^\top y \\ & \text{subject to} && cu + h(x_0\mathbf{1} + Uu - Uw_{lb}) \leq y, \\ & && cu - p(x_0\mathbf{1} + Uu - Uw_{ub}) \leq y, \\ & && \mathbf{0} \leq u \leq M\mathbf{1} \end{aligned} \quad (16.15)$$

得到最坏情况下的最优决策.

### §16.5.2 区间不确定性下的仿射订购策略

在以上的表述中, 是在时段  $k = 0$  且以一种“开环”的方式来计算所有向前的订货决策  $u(0), \dots, u(T-1)$ . 但在实际中, 人们只会执行第一次决策  $u(0)$  (所谓的“此时此地”决策), 然后等待并观察接下来发生的事, 直到下一个决定时段. 在时段  $k = 1$ , 可以观察到不确定需求的实际实现. 因此, 当确定决策  $u(1)$  时可以获取这个信息. 显然, 这些信息会使新的决策  $u(1)$  受益, 因此原先的忽略这些信息的方法是次优的.

一般地, 在时段  $k$  将执行的决策  $u(k)$  将受益于从时段 0 到  $k-1$  之间观察到的需求信息 (这些需求在时间  $k$  时已经是确定的了), 也就是说  $u(k)$  是过去需求的泛函, 即  $u(k) = \varphi_k(w(0), \dots, w(k-1))$ , 在这样的反映式或者“闭环”的方法下, 这个优化问题相当于是寻找最优函数  $\varphi_k$  (通常称为政策), 使得最坏情况成本最小. 但是寻找所有泛函中搜索使得问题变得极为困难 (需要在一个由所有可能函数  $\varphi_k$  组成的“无限维”集合中进行搜寻). 因此, 一种常用方法是固定  $\varphi_k$  中的参数, 然后优化有限维的参数集合 (参考 14.3.3 节描述了一个金融领域中的多期优化问题的类似方法). 举例来说, 有一种有效的方法是, 考虑仿射式参数化, 也就是考虑  $k = 1, \dots, T-1$  时决策  $\varphi_k$  形如

$$u(k) = \varphi_k(w(0), \dots, w(k-1)) = \bar{u}(k) + \sum_{i=0}^{k-1} \alpha_{k,i}(w(i) - \hat{w}(i)),$$

这里  $\bar{u}(k)$  是“名义”决策,  $u(0) = \bar{u}(0)$ , 对  $k \geq 1$ ,  $\alpha_{k,i}$  是我们修正名义决策使其和需求  $w(i)$  与名义需求  $\hat{w}(i)$  之间偏差成正比例的系数. 使用矩阵记号,

我们写成

$$\mathbf{u} = \bar{\mathbf{u}} + \mathbf{A}(\mathbf{w} - \hat{\mathbf{w}}), \quad \mathbf{A} \doteq \begin{bmatrix} 0 & 0 & \cdots & 0 \\ \alpha_{1,0} & 0 & \cdots & 0 \\ \alpha_{2,0} & \alpha_{2,1} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ \alpha_{T-1,0} & \cdots & \alpha_{T-1,T-2} & 0 \end{bmatrix},$$

其中  $\bar{\mathbf{u}} = (\bar{u}(0), \dots, \bar{u}(T-1))$ . 特别的, 在区间不确定模型(16.13)中, 需求波动  $\tilde{\mathbf{w}} \doteq \mathbf{w} - \hat{\mathbf{w}}$  属于对称向量区间

$$-\bar{\mathbf{w}} \leq \tilde{\mathbf{w}} \leq \bar{\mathbf{w}}; \quad \bar{\mathbf{w}} \doteq \frac{\rho}{100} \hat{\mathbf{w}} \geq \mathbf{0}.$$

将这里的  $\mathbf{u}$  的表达式代入问题(16.14), 我们得到

$$\begin{aligned} & \underset{\bar{\mathbf{u}}, \mathbf{y}, \mathbf{A}}{\text{minimize}} && \mathbf{1}^\top \mathbf{y} \\ & \text{subject to} && (c\bar{\mathbf{u}} + c\mathbf{A}\tilde{\mathbf{w}}) + h(x_0\mathbf{1} + \mathbf{U}(\bar{\mathbf{u}} + \mathbf{A}\tilde{\mathbf{w}}) - \mathbf{U}\mathbf{w}) \leq \mathbf{y}, \quad \forall \mathbf{w} \in W \\ & && (c\bar{\mathbf{u}} + c\mathbf{A}\tilde{\mathbf{w}}) - p(x_0\mathbf{1} + \mathbf{U}(\bar{\mathbf{u}} + \mathbf{A}\tilde{\mathbf{w}}) - \mathbf{U}\mathbf{w}) \leq \mathbf{y}, \quad \forall \mathbf{w} \in W \\ & && \mathbf{0} \leq \mathbf{u} + \mathbf{A}\tilde{\mathbf{w}} \leq M\mathbf{1}, \quad \forall \mathbf{w} \in W. \end{aligned} \tag{16.16}$$

如果  $\mathbf{v}$  是一个向量并且  $\bar{\mathbf{w}} \geq \mathbf{0}$ , 易于验证

$$\begin{aligned} \underset{-\bar{\mathbf{w}} \leq \tilde{\mathbf{w}} \leq \bar{\mathbf{w}}}{\text{maximize}} & \mathbf{v}^\top \tilde{\mathbf{w}} = |\mathbf{v}|^\top \bar{\mathbf{w}}, \\ \underset{-\bar{\mathbf{w}} \leq \tilde{\mathbf{w}} \leq \bar{\mathbf{w}}}{\text{minimize}} & \mathbf{v}^\top \tilde{\mathbf{w}} = -|\mathbf{v}|^\top \bar{\mathbf{w}}, \end{aligned}$$

这里  $|\mathbf{v}|$  是由  $\mathbf{v}$  分量的绝对值组成的向量. 针对问题(16.16)的约束逐行应用这个规则, 即可得到这个鲁棒问题等价于

$$\begin{aligned} & \underset{\bar{\mathbf{u}}, \mathbf{y}, \mathbf{A}}{\text{minimize}} && \mathbf{1}^\top \mathbf{y} \\ & \text{subject to} && c\bar{\mathbf{u}} + h\mathbf{U}\bar{\mathbf{u}} + hx_0\mathbf{1} - h\mathbf{U}\hat{\mathbf{w}} + |c\mathbf{A} + h\mathbf{U}\mathbf{A} - h\mathbf{U}|\bar{\mathbf{w}} \leq \mathbf{y}, \\ & && c\bar{\mathbf{u}} - p\mathbf{U}\bar{\mathbf{u}} - px_0\mathbf{1} + p\mathbf{U}\hat{\mathbf{w}} + |c\mathbf{A} - p\mathbf{U}\mathbf{A} + p\mathbf{U}|\bar{\mathbf{w}} \leq \mathbf{y}, \\ & && \bar{\mathbf{u}} + |\mathbf{A}|\bar{\mathbf{w}} \leq M\mathbf{1}, \\ & && \bar{\mathbf{u}} - |\mathbf{A}|\bar{\mathbf{w}} \geq \mathbf{0}. \end{aligned}$$

引入三个由松弛变量组成的下三角矩阵  $\mathbf{Z}_1, \mathbf{Z}_2, \mathbf{Z}_3$ , 我们将此问题重新表

述成LP

$$\begin{aligned}
 & \underset{\bar{\mathbf{u}}, \mathbf{y}, \mathbf{A}, \mathbf{Z}_1, \mathbf{Z}_2, \mathbf{Z}_3}{\text{minimize}} && \mathbf{1}^\top \mathbf{y} \\
 & \text{subject to} && c\bar{\mathbf{u}} + h\mathbf{U}\bar{\mathbf{u}} + hx_0\mathbf{1} - h\mathbf{U}\hat{\mathbf{w}} + \mathbf{Z}_1\bar{\mathbf{w}} \leq \mathbf{y}, \\
 & && c\bar{\mathbf{u}} - p\mathbf{U}\bar{\mathbf{u}} - px_0\mathbf{1} + p\mathbf{U}\hat{\mathbf{w}} + \mathbf{Z}_2\bar{\mathbf{w}} \leq \mathbf{y}, \\
 & && \bar{\mathbf{u}} + \mathbf{Z}_3\bar{\mathbf{w}} \leq M\mathbf{1}, \\
 & && \bar{\mathbf{u}} - \mathbf{Z}_3\bar{\mathbf{w}} \geq 0, \\
 & && |c\mathbf{A} + h\mathbf{U}\mathbf{A} - h\mathbf{U}| \leq \mathbf{Z}_1, \\
 & && |c\mathbf{A} - p\mathbf{U}\mathbf{A} + p\mathbf{U}| \leq \mathbf{Z}_2, \\
 & && |\mathbf{A}| \leq \mathbf{Z}_3.
 \end{aligned} \tag{16.17}$$

一旦求解了这个问题，即可获得订货策略和不确定的库存量

$$\begin{aligned}
 \mathbf{u} &= \bar{\mathbf{u}} + \mathbf{A}\tilde{\mathbf{w}}, \\
 \mathbf{x} &= x_0\mathbf{1} + \mathbf{U}(\bar{\mathbf{u}} - \hat{\mathbf{w}}) + (\mathbf{U}\mathbf{A} - \mathbf{U})\tilde{\mathbf{w}}.
 \end{aligned} \tag{16.18}$$

然后，我们能获得订购和库存量的上下极限，如下

$$\begin{aligned}
 \mathbf{u}_{\text{lb}} &= \bar{\mathbf{u}} - |\mathbf{A}|\bar{\mathbf{w}}, \\
 \mathbf{u}_{\text{ub}} &= \bar{\mathbf{u}} + |\mathbf{A}|\bar{\mathbf{w}}, \\
 \mathbf{x}_{\text{lb}} &= x_0\mathbf{1} + \mathbf{U}(\bar{\mathbf{u}} - \hat{\mathbf{w}}) - |\mathbf{U}\mathbf{A} - \mathbf{U}|\bar{\mathbf{w}}, \\
 \mathbf{x}_{\text{ub}} &= x_0\mathbf{1} + \mathbf{U}(\bar{\mathbf{u}} - \hat{\mathbf{w}}) + |\mathbf{U}\mathbf{A} - \mathbf{U}|\bar{\mathbf{w}}.
 \end{aligned}$$

实践中，将根据策略(16.18)在线计算时段 $k$ 时的实际订购量 $u(k)$ 。随着时间的演进，将观察到 $\tilde{w}(i)$ ,  $i = 0, \dots, k-1$ , 的实际值。我们已知的先验信息是：当 $k = 0, \dots, T-1$ 时，值 $u(k)$ 将包含于区间 $[u_{\text{lb}}(k), u_{\text{ub}}(k)]$ 。

下面用一个例子来说明这些想法。考虑有 $T = 12$ (比如月)个时段的决策，相应的数据为：购买成本 $c = 5$ ，存储成本 $h = 4$ ，短缺成本 $p = 6$ ，订购上界 $M = 110$ 。进一步，假设名义需求服从正弦曲线

$$\hat{w}(k) = 100 + 20 \sin\left(2\pi \frac{k}{T-1}\right), \quad k = 0, \dots, T-1, \tag{16.19}$$

并且初始库存量 $x_0 = 100$ 。在没有不确定性的情况下(需求量 $\mathbf{w}$ 等于名义需求量 $\hat{\mathbf{w}}$ )，求解问题(16.12)得到的最优成本是5,721.54，并且订购与存储量的图形见图16.19。

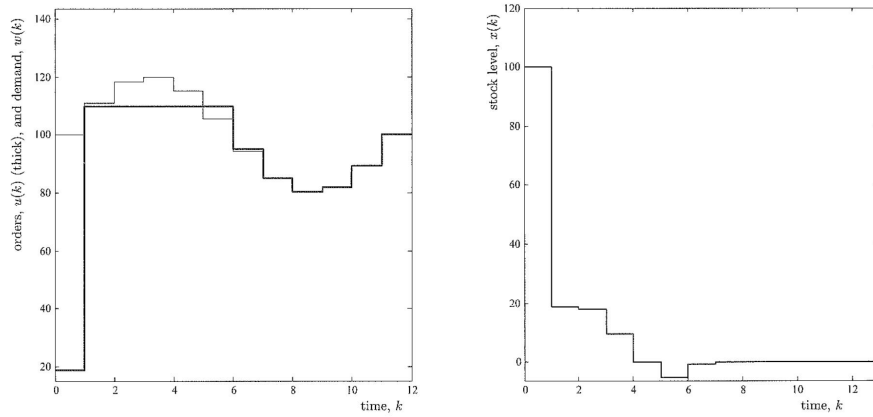


图 16.19: 名义需求下的最优订购策略. 左图: 需求和订购曲线. 右图: 库存量曲线. 最优的名义成本为5,721.54.

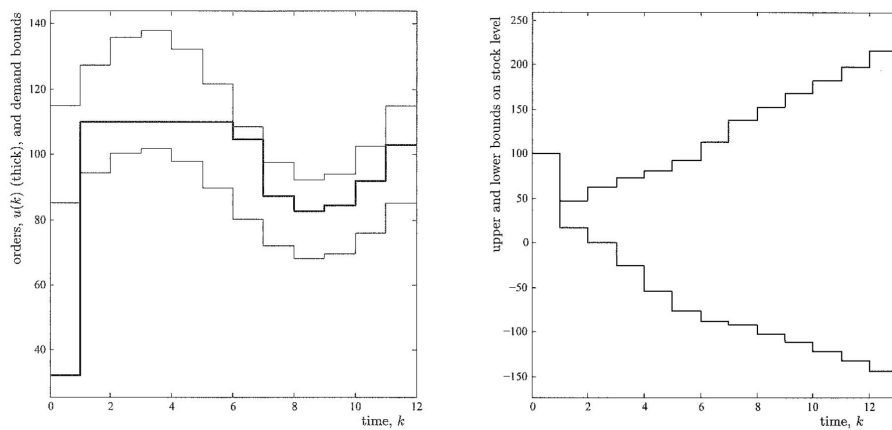


图 16.20: 在 $\rho = 15\%$ 不确定性下的最优“开环”订购. 左图: 需求和订购曲线, 右图: 库存量曲线. 最坏情况成本为11,932.50.

接下来考虑需求量存在 $\rho = 15\%$ 的不确定性的情况. 我们求解鲁棒的“开环”问题(16.15), 得到的最坏情况的成本是11,932.50. 订购与库存量曲线如图16.20所示.

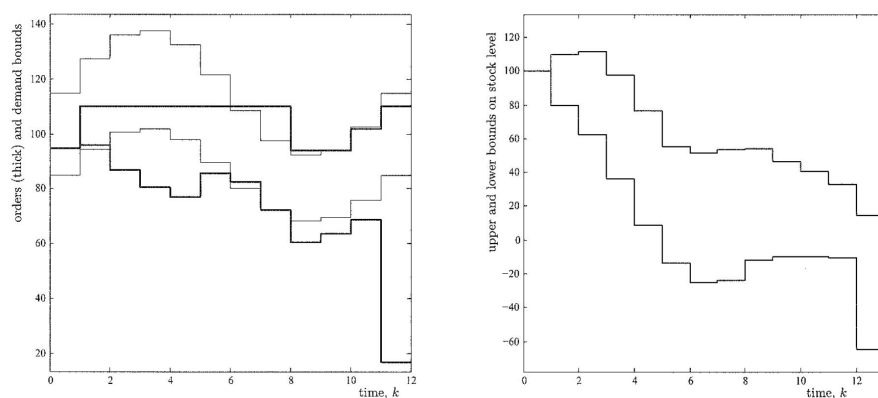


图 16.21: 不确定性为  $\rho = 15\%$  时的最优“闭环”订购. 左图, 需求和订购曲线; 右图: 库存量曲线. 最坏情况的成本为 8,328.64

在需求量同样有 15% 的不确定性的情况下, 我们求解鲁棒的“闭环”问题(16.17), 得到最坏情况下的成本是 8,328.64. 订购与库存量的曲线如图 16.21 所示. 进一步, 由(16.18)得到的订购策略的最优参数为

$$\bar{\mathbf{u}}^T = \begin{bmatrix} 94.69 & 103.0 & 98.36 & 95.26 & 93.49 & 97.83 & 96.26 & 91.25 & 77.33 & 78.68 & 85.2 & 63.4 \end{bmatrix},$$

$$\mathbf{A} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0.4697 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0.2506 & 0.4743 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0.1332 & 0.2518 & 0.4828 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0.07069 & 0.1324 & 0.254 & 0.4868 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0.007817 & 0.01524 & 0.04958 & 0.1856 & 0.4395 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0.009959 & 0.0148 & 0.03499 & 0.0915 & 0.2258 & 0.4529 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0.02248 & 0.05136 & 0.08535 & 0.114 & 0.1627 & 0.2676 & 0.4913 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0.01838 & 0.03096 & 0.04785 & 0.06251 & 0.0879 & 0.1427 & 0.2594 & 0.5104 & 0 & 0 & 0 & 0 & 0 \\ 0.009341 & 0.01578 & 0.02445 & 0.03202 & 0.04514 & 0.07356 & 0.1343 & 0.2649 & 0.5463 & 0 & 0 & 0 & 0 \\ 0.004745 & 0.008067 & 0.01258 & 0.01654 & 0.02343 & 0.03848 & 0.07088 & 0.1412 & 0.2945 & 0.6839 & 0 & 0 & 0 \\ 0.01105 & 0.02035 & 0.03276 & 0.04322 & 0.06086 & 0.09775 & 0.1717 & 0.3123 & 0.547 & 0.8799 & 1.372 & 0 & 0 \end{bmatrix}.$$

通过使用“反应式的”(或“闭环”)决策取代静态的“开环”法, 就“最坏情况下的成本”而言, 获得的改进超过 30%.

### §16.5.3 用于一般的随机不确定性的场景法

当需求量  $w(k)$  的不确定性是随机的, 且不能建模为前面的区间模型时, 我们可以使用一种虽然是近似但却简单而有效的方法, 即基于不确定性的场景抽样. 该法中, 我们假设需求量  $\mathbf{w}$  (包含需求  $w(0), \dots, w(T-1)$ ) 是一个随机向量, 它的期望  $\hat{\mathbf{w}} = \mathbb{E}\{\mathbf{w}\}$  是已知的, 它的概率分布也是已知的, 或者至少得到了它的  $N$  个独立同分布的样本  $\mathbf{w}^{(1)}, \dots, \mathbf{w}^{(N)}$ . 在这种设



定下,  $w(k)$ 不必有界, 且基于所设定的分布, 有可能是时间相关的.

这种基于场景的方法<sup>5</sup>的简单思想假定: 如果 $N$ 足够大, 则 $N$ 个生成的场景的全体将会合理地呈现出不确定性. 这样, 这里的场景—鲁棒问题相当于求解形如(16.16)的优化问题, 但是要用 $\forall \mathbf{w} \in \{\mathbf{w}^{(1)}, \dots, \mathbf{w}^{(N)}\}$ 代替 $\forall \mathbf{w} \in W$ . 也就是说, 我们的目的不是对所有可能的不确定性的实现均要求约束成立, 而是对样本场景值要求这些约束成立. 显然, 这样所获得的解在绝对的最坏情况下将不再是鲁棒的, 但在非严格的概率情况下是鲁棒的. 我们可以通过选择足够大的 $N$ 以获得一个好的鲁棒性. 对给定的 $\alpha \in (0, 1)$ , 由等式(10.30), 如果已经确定了场景数量, 那么

$$N \geq \frac{2}{1-\alpha}(n+10), \quad (16.20)$$

其中 $n$ 是优化问题中决策变量的总数目, 那么这个场景解将是水平为 $\alpha$ 的“概率鲁棒的”(然而, 通常对于即使较小的情境数目 $N$ , 也能获得好的结果).

按以上描述, 我们需要求解一个如下的线性优化问题:

$$\begin{aligned} & \underset{\bar{\mathbf{u}}, \mathbf{y}, \mathbf{A}}{\text{minimize}} && \mathbf{1}^\top \mathbf{y} \\ & \text{subject to} && \mathbf{c}\mathbf{u}^{(i)} + h\mathbf{x}^{(i)} \leq \mathbf{y}, \quad i = 1, \dots, N, \\ & && \mathbf{c}\mathbf{u}^{(i)} - p\mathbf{x}^{(i)} \leq \mathbf{y}, \quad i = 1, \dots, N, \\ & && \mathbf{u}^{(i)} = \bar{\mathbf{u}} + \mathbf{A}(\mathbf{w}^{(i)} - \hat{\mathbf{w}}), \quad i = 1, \dots, N, \\ & && \mathbf{x}^{(i)} = x_0 \mathbf{1} + \mathbf{U}\mathbf{u}^{(i)} - \mathbf{U}\mathbf{w}^{(i)}, \quad i = 1, \dots, N, \\ & && 0 \leq \mathbf{u}^{(i)} \leq M\mathbf{1}, \quad i = 1, \dots, N. \end{aligned} \quad (16.21)$$

作为一个数值例子, 我们再次考虑16.5.2节中的数据. 此时, 问题(16.21)的决策变量数目

$$n = T + T + \frac{T(T-1)}{2} = 90.$$

考虑鲁棒性水平 $\alpha = 0.9$ , 由公式(16.20)计算得知, 在问题中需要使用 $N = 2,000$ 个情境. 我们进一步假设期望需求由(16.19)给定, 向量 $\mathbf{w}$ 满足正态分布, 其协方差矩阵

$$\Sigma = \text{diag}(\sigma_0^2, \dots, \sigma_{T-1}^2),$$

其中方差随时间演进的方式(刻画不确定性随时间增加的情况)为

$$\sigma_k^2 = (1+k)\bar{\sigma}^2; \quad k = 0, \dots, T-1, \quad \text{其中 } \bar{\sigma}^2 = 1.$$

<sup>5</sup>见G.Calafiore的Random convex programs, *SIAM J. Opt.*, 2010.

将这些参数代入场境(16.21), 求解所得的场景问题, 得到最坏情况下的最优成本是6,296.80, 最优决策如下

$$\bar{\mathbf{u}}^T = \begin{bmatrix} 24.4939 & 108.77 & 108.603 & 108.808 & 110.0 & 110.0 & 98.7917 & 83.7308 & 80.6833 & 81.5909 & 89.8084 & 94.6564 \end{bmatrix},$$

$$\mathbf{A} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0.3774 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0.02567 & 0.2845 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ -0.1184 & 0.2025 & 0.1344 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0.4868 & 0.03163 & 0.496 & 0.4467 & 0.4006 & 0.5332 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0.3849 & -0.1964 & 0.254 & 0.2024 & 0.4895 & 0.2111 & 0.6726 & 0 & 0 & 0 & 0 & 0 \\ 0.135 & 0.072 & -0.1547 & 0.1603 & -0.03367 & 0.02405 & 0.04367 & 0.5894 & 0 & 0 & 0 & 0 \\ -0.2695 & 0.6738 & 0.1186 & -0.2117 & 0.08057 & 0.001069 & 0.1795 & 0.2256 & 0.6008 & 0 & 0 & 0 \\ 0.2583 & 0.1248 & 0.04462 & 0.07893 & 0.05989 & 0.09069 & 0.007373 & 0.3262 & 0.2858 & 0.5088 & 0 & 0 \\ 0.4861 & 0.6754 & 0.009505 & 0.3034 & 0.3853 & 0.08153 & 0.1048 & 0.4221 & 0.6108 & 0.3069 & 1.086 & 0 \end{bmatrix}.$$

然后, 我们可以通过Monte Carlo法, 也就是生成新的场境(与优化中的不同), 来仿真这个决策在随随机需求下的后验行为. 使用1,500 个新的需求场境, 我们得到了图16.22中的结果. 由各个情境产生的成本的直方图见图16.23. 这个仿真得到的最大成本是6,190.22, 而场境法中的最坏情况成本为6,296.80. 显然前者要低于后者.

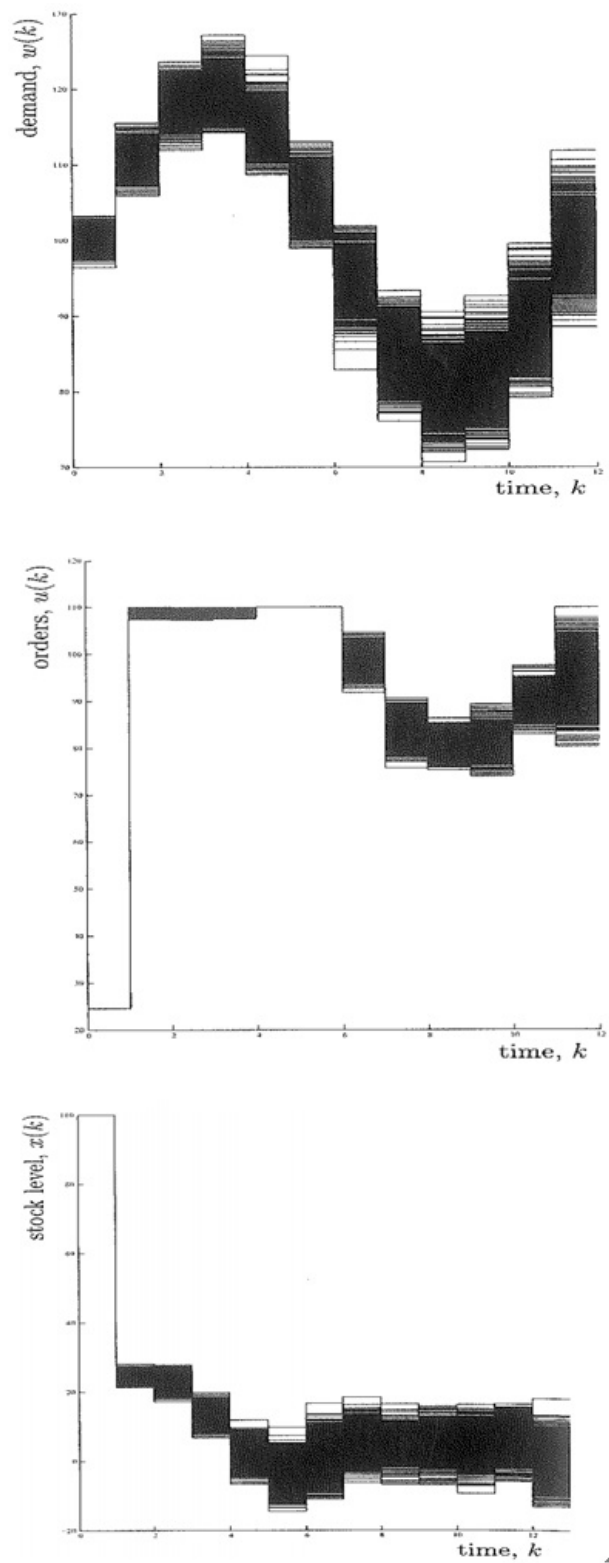


图 16.22: 1,500个随机生成的需求下的最优场境策略的Monte Carlo仿真. 上图: 随机需求. 中图: 订购曲线. 下图: 库存量曲线.

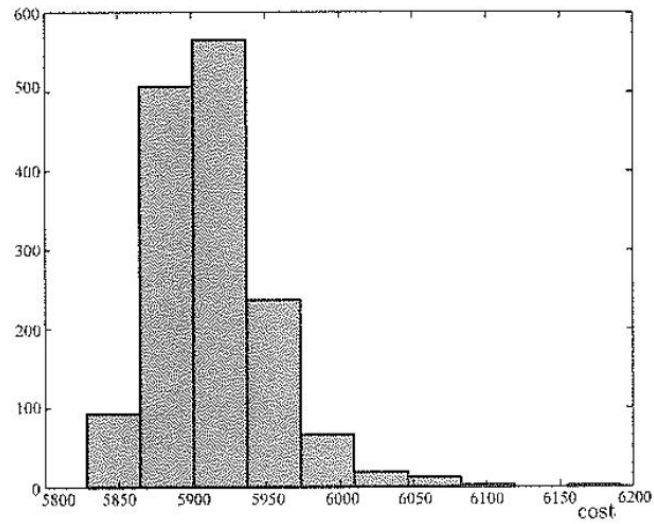


图 16.23: 由1,500个随机生成需求下最优场景策略的Monte Carlo仿真所产生的成本的直方图.

### §16.6 练习

#### 习题16.1. (网络拥塞控制)

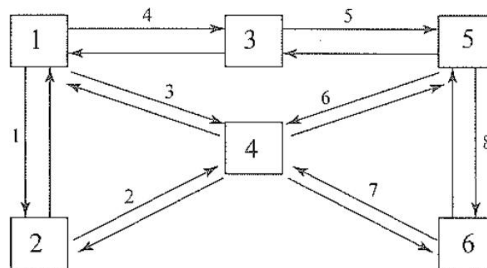


图 16.24: 一个小型网络

图16.24显示了由 $n = 6$ 个计算机组成的P2P网络. 每一台计算机可以在图中所示链接以特定的速率上传和下载数据. 令 $\mathbf{b}^+ \in \mathbb{R}^8$ 为一个向量, 表示图中标号对应的连接上的数据传输速率, 令 $\mathbf{b}^- \in \mathbb{R}^8$ 为一向量, 表示反向连接的数据传输速率, 这里必须有 $\mathbf{b}^+ \geq \mathbf{0}, \mathbf{b}^- \geq \mathbf{0}$ .

为此网络定义节点弧关联矩阵

$$\mathbf{A} \doteq \begin{bmatrix} 1 & 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ -1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 & 1 & 0 & 0 & 0 \\ 0 & -1 & -1 & 0 & 0 & -1 & -1 & 0 \\ 0 & 0 & 0 & 0 & -1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & -1 \end{bmatrix},$$

令  $\mathbf{A}_+ \doteq \max(\mathbf{A}, 0)$  ( $\mathbf{A}$  的正部),  $\mathbf{A}_- \doteq \min(\mathbf{A}, 0)$  ( $\mathbf{A}$  的负部). 由此得到节点的总输出(上传)速率  $\mathbf{v}_{\text{upl}} = \mathbf{A}_+ \mathbf{b}^+ - \mathbf{A}_- \mathbf{b}^-$ , 总输入(下载)速率  $\mathbf{v}_{\text{dwl}} = \mathbf{A}_+ \mathbf{b}^- - \mathbf{A}_- \mathbf{b}^+$ . 因此节点上网络的出流可以表示为

$$\mathbf{v}_{\text{net}} = \mathbf{v}_{\text{upl}} - \mathbf{v}_{\text{dwl}} = \mathbf{A} \mathbf{b}^+ - \mathbf{A} \mathbf{b}^-,$$

此外还需要满足流平衡方程  $[v_{\text{net}}]_i = f_i$ , 如果计算机  $i$  既没有产生数据也没有接收数据(它仅仅是传递接收到的数据, 即它的表现像一个中继站), 则  $f_i = 0$ ; 如果计算机  $i$  产生包, 则  $f_i > 0$ ; 若它接收包, 则  $f_i < 0$ , 其中  $f_i$  是已分配的速率.

每一台计算机下载数据的最大速率  $\bar{v}_{\text{dwl}} = 20$  Mbit/s, 上传数据的最大速率  $\bar{v}_{\text{upl}} = 10$  Mbit/s (是对每台计算机所有上传或下载连接的速率之和的限制). 每一条连接的拥堵水平定义为

$$c_j = \max(0, (b_j^+ + b_j^- - 4)), \quad j = 1, \dots, 8.$$

假设点1必须以  $f_1 = 9$  Mbit/s 的速率传输数据到点5, 点2必须以  $f_2 = 8$  Mbit/s 的速率传输数据到点6. 求出网络平均拥挤水平最小时所有连接的数据传输速率.

**习题16.2. (水库设计问题)** 我们需要设计一个用于蓄水和蓄能的水库. 如图16.25所示.

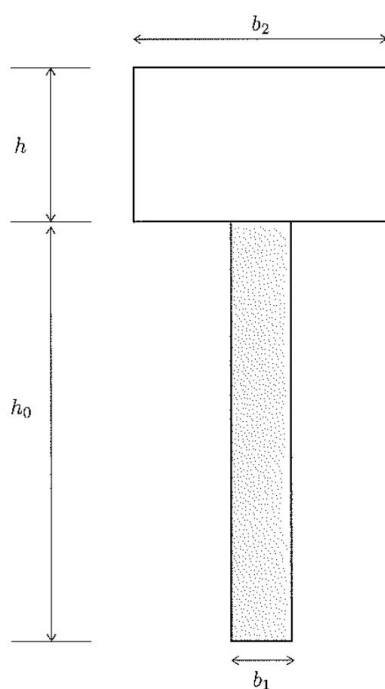


图 16.25: 混凝土基座上的蓄水池

混凝土基座截面为长方形, 长为 $b_1$ , 高为 $h_0$ , 水库本身截面也为长方形, 长为 $b_2$ , 高为 $h$ . 表格16.4中给出了一些需要用到的数据.

表 16.4: 水库问题的数据

参数	结果	单位	描述
$g$	9.8	$\text{m/s}^2$	重力加速度
$E$	$30 \times 10^9$	$\text{N/m}^2$	基座弹性模量
$\rho_w$	$10 \times 10^3$	$\text{N/m}^3$	水的密度
$\rho_b$	$25 \times 10^3$	$\text{N/m}^3$	基座密度
$J$	$b_1^4/12$	$\text{m}^4$	基座惯性矩
$N_{cr}$	$\pi^2 JE / (2h_0)^2$	N	基座的临界载荷限制

基座的临界负载限制为 $N_{cr}$ , 它至少为水的重量的两倍, 结构要满足规范 $h_0/b_1^2 \leq 35$ . 蓄水池的形状应满足 $1 \leq b_2/h \leq 2$ . 结构的总高度不能够超过30 m. 结构的总重量(基座和装满水的蓄水池)不应超过 $9.8 \times 10^5$  N. 问

题是找到维度  $b_1, b_2, h_0, h$  使得蓄水池中水的势能  $P_w$  最大(假设  $P_w = (\rho_w h b_2^2) h_0$ ). 请问这个问题能否被建模为凸优化问题? 如果可以的话, 请给出建立的模型和解释, 并求解所得模型给出最优设计.

**习题16.3. (电路中的导线尺寸问题)** 现代电子芯片间的互连可以被建模为置于基质之上的导电表面积. 因此可以将一根“导线”看成图16.26所示的矩形段序列.



图 16.26: 将导线表示成基质上的矩形表面. 每个片段的长度  $\ell_i$  是固定的, 宽度  $x_i$  是决策变量. 该例中有三个导线段.

我们假设每个矩形段的长度是固定的, 而宽度需要依据下文说明的准则来确定. 一个通用的方法是把导线模型化为数个RC段的级联, 其中对每个段而言,  $S_i = 1/R_i$  和  $C_i$  分别是第  $i$  个矩形段的电导和电容, 参见图16.27.

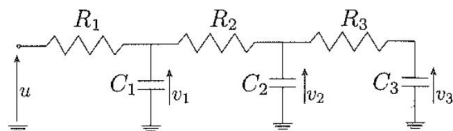


图 16.27: 三分段的导线RC模型

$S_i$  和  $C_i$  的值与导线段的表面积成比例, 由于已知长度  $\ell_i$  是定值, 因此电容和电导是宽度的仿射函数, 即

$$S_i = S_i(x_i) = \sigma_i^{(0)} + \sigma_i x_i, \quad C_i = C_i(x_i) = c_i^{(0)} + c_i x_i,$$

这里  $\sigma_i^{(0)}, \sigma_i, c_i^{(0)}, c_i$  是给定的正常数. 对于图16.27中的三分段导线模型, 可以写出如下的动力方程组来描述节点间电压  $v_i(t)$ ,  $i = 1, \dots, 3$ , 随时间的变化, 即

$$\begin{bmatrix} C_1 & C_2 & C_3 \\ 0 & C_2 & C_3 \\ 0 & 0 & C_3 \end{bmatrix} \dot{\mathbf{v}}(t) = - \begin{bmatrix} S_1 & 0 & 0 \\ -S_2 & S_2 & 0 \\ 0 & -S_3 & S_3 \end{bmatrix} \mathbf{v}(t) + \begin{bmatrix} S_1 \\ 0 \\ 0 \end{bmatrix} u(t),$$

其中 $u(t)$ 是输入.

如果我们进行换元

$$\mathbf{v}(t) = \mathbf{Q}\mathbf{z}(t), \quad \mathbf{Q} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix},$$

则上述方程组就可表述成更有用的形式. 具体地, 利用这里引入的新变量 $\mathbf{z}(t)$ , 上述方程变成

$$\mathbf{C}(\mathbf{x})\dot{\mathbf{z}}(t) = -\mathbf{S}(\mathbf{x})\mathbf{z}(t) + \begin{bmatrix} S_1 \\ 0 \\ 0 \end{bmatrix} u(t),$$

其中

$$\mathbf{C}(\mathbf{x}) \doteq \begin{bmatrix} C_1 + C_2 + C_3 & C_2 + C_3 & C_3 \\ C_2 + C_3 & C_2 + C_3 & C_3 \\ C_3 & C_3 & C_3 \end{bmatrix}, \quad \mathbf{S}(\mathbf{x}) \doteq \text{diag}(S_1, S_2, S_3).$$

$\mathbf{C}(\mathbf{x})$ ,  $\mathbf{S}(\mathbf{x})$ 显然是对称矩阵, 他们的每一个分量是决策变量 $\mathbf{x} = (x_1, x_2, x_3)$ 的仿射函数. 此外我们可以看出: 只要 $\mathbf{x} \geq \mathbf{0}$ ,  $\mathbf{C}(\mathbf{x})$ 就是非奇异的(我们问题中的物理情形就是这样的), 因此 $\mathbf{z}(t)$ 的演化可以被表示为(我们接下来假设 $u(t) = 0$ , 也就是说我们仅考虑系统自由响应的时间演化)

$$\dot{\mathbf{z}}(t) = -\mathbf{C}(\mathbf{x})^{-1}\mathbf{S}(\mathbf{x})\mathbf{z}(t).$$

定义电路的主时间常数(dominant time constant)为

$$\tau = \frac{1}{\lambda_{\min}(\mathbf{C}(\mathbf{x})^{-1}\mathbf{S}(\mathbf{x}))'},$$

这是电路“速度”的一种量测( $\tau$ 越小, 电路的响应越快).

在保证主时间常数不超过一个指定水平 $\eta > 0$ 的前提下, 描述一个可以有效确定导线宽度 $\mathbf{x}$ 的方法, 其使得导线所占总面积最小.