



同濟大學
TONGJI UNIVERSITY

交通数据分析---期中作业

同济大学交通运输与工程学院

综合交通信息与控制工程系

1854084 徐晨龙

Part 01

数据说明

本数据取自 2017 年 2 月 6 日（星期一）至 2017 年 2 月 10 日（星期五）每天 0:00-12:00 的快速路微波车辆检测器数据。数据集内共包含 5 个检测器点位，检测范围内包括 2 个快速路入口匝道和 1 个快速路出口匝道。其具体布设位置如以下简图所示：

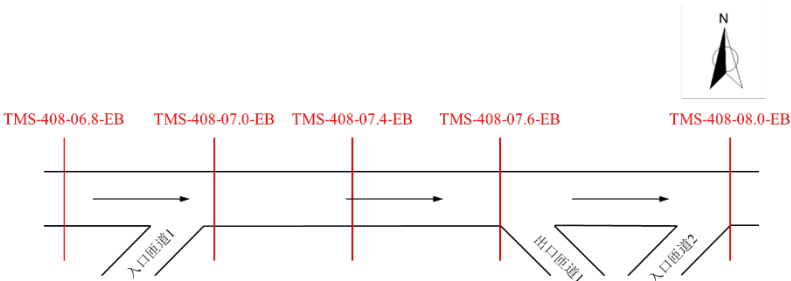


图 1：快速路微波的路面实际布置图

观察所给数据可以看出各检测线圈检测到的车道为：

- TMS – 408 – 06.8 – EB: line01 ,line02,line 03
- TMS – 408 – 07.0 – EB: line01 ,line02,line 03 ,line04
- TMS – 408 – 07.4 – EB: line01 ,line02,line 03
- TMS – 408 – 07.6 – EB: line01 ,line02,line 03 ,line04
- TMS – 408 – 08.0 – EB: line01 ,line02,line 03 ,Rmp – line04

理论上完整的采集数据量为：

$$N = 5 * (12 * 60 + 1) * 18 = 64890$$

处理要求

- (1) 分别采用最小值—最大值规范化法、 $z - score$ 规范化法（标准偏差）对检测数据速度进行标准化处理。并分别提供均值与标准差。
- (2) 分别运用阈值法的五分位数法、3 倍标准差法，剔除异常检测数据（速度、流量）。并给出异常数据的样本集
- (3) 分别使用时间序列法、或基于历史数据的修补方法、或基于空间位置的修补方法，修复缺失检测数据（包括运用阈值法剔除异常数据后产生的空缺数据）。给出填补的数据样本集

数据分析结果

- (1) 由于不同车道交通流参数分布的不同，因此首先将不同车道的数据分离，从数据和位置图之间的关系可以看出不同车道主要有：
lane1, lane2, lane3, ramp1, ramp2, ramp3
之后分别利用最小值-最大值规范化法、 $z - score$ 规范化法对不同车道的数据进

行处理，得到对应数据处理前后的均值和标准差如下表所示：

表格 1：处理后的不同车道的均值和标准差

	原始数据		最小值-最大值规范化		$z - score$ 规范化	
	<i>mean</i>	<i>sd</i>	<i>mean</i>	<i>sd</i>	<i>mean</i>	<i>sd</i>
<i>lane1</i>	67.8	6.95	0.7	0.05	0	1.6
<i>lane2</i>	63.09	5.92	0.3	0.03	0	1.59
<i>lane3</i>	61.6	6.02	0.3	0.038	0	1.54
<i>ramp1</i>	49.98	6.23	0.55	0.098	0	1.34
<i>ramp2</i>	53.18	5.72	0.46	0.096	0	1.37
<i>ramp3</i>	45.8	7.14	0.48	0.12	0	1.3

处理方式是首先提取所选取的六条车道的信息，之后使用函数 $data_standard = funciton(data)$ ，其作用对输入数据返回原始数据的均值和标准差、最小值-最大值规范后的均值和标准差、 $z - score$ 的均值和标准差。

(2) 本案例中异常数据可以分为缺失数据、异常数据、重复冗余数据。由于数据中存在大量的缺失数据、冗余数据、异常数据。因此我们首先生成完整的表格（理论上 64890 行）

由于 *date*、*tiemstamp*、*lane_id* 可以唯一确定一个样本数据，因此首先生成完整地三行对应数据，之后新建的 *New_midterm_data* 中 *speed*、*volume*、*occupancy* 三列向量对应的默认值为 -1，之后根据循环判断语句将对应的数据加到新建的数据框中

```
logic = (Midterm_assignment$date =
          = data_date[n]&Midterm_assignment$timestamp =
          = data_time[n]&Midterm_assignment$lane_id == data_lane[n])
```

寻找对应数据的逻辑语句如上式表示，其中重复数据值中逻辑下标值会存在多个，因此赋值是采用下式可以解决重复数据：

```
New_Midterm_data[n,c(5,6,7)]
= apply(Midterm_assignment[logic,c(5,6,7)],2,sum)/sum(logic)
```

这里仅采用“独立判断+联合判断+有效车长判断”来识别故障数据，对应函数 $judge1 = function(data)$ ：

a. 独立判断，即根据流量、平均速度和占有率的取值范围来判断异常值，根据交通经验设置采取的取值范围如下

$$0 \leq q \leq 51(pcu/min)$$

$$0 \leq v \leq 80 * 1.4(km/h)$$

$$0 \leq o \leq 95$$

b. 利用流量、速度、占有率之间的逻辑关系进行判别，可能的情况有：

- 情况 1. 流量、速度、占有率均不为 0
- 情况 2. 流量、速度、占有率均为 0

情况 3. 流量、密度为 0，占有率超过 95%

情况 4. 流量超过 5pcu/min，车速不为 0，占有率为 0

c. 根据联合车长判断，根据单位换算 $v(km/h)$ 、 $q(pcu/min)$

$$Avel = 1000 * V * O / Q = v * o / 6q$$

其中平均车长的合理范围(有部分数据车长约为 1.9m 也认为正常，因此将 2m 改为 1.5m)为：

$$1.5m < Avel < 22m$$

独立判断可以根据不同变量的阈值范围来单独判断是否出错，因此如果出错会单独设置对应的标记位。对于联合判断和有效车长判断的情况，由于其对应的三个之间对应的关系，因此如果出错会将三者的标志位均记为 1

经过计算可得，故障数据的样本集为：

	date	timestamp	detector_id	lane_id	type	speed	s_ab	volume	v_ab	occupancy	o_ab	error
1	2017-02-10	1899-12-31 05:51:00	TMS-408-08.0-EB	TMS-408-08.0-EB-Rmp-Lane04	Ramp	37	1	1	1	4	1	3
2	2017-02-10	1899-12-31 09:01:00	TMS-408-08.0-EB	TMS-408-08.0-EB-Rmp-Lane04	Ramp	45	1	2	1	6	1	3
3	2017-02-09	1899-12-31 11:01:00	TMS-408-07.0-EB	TMS-408-07.0-EB-Lane01	Mainline	58	1	1	1	6	1	3
4	2017-02-09	1899-12-31 11:01:00	TMS-408-07.0-EB	TMS-408-07.0-EB-Lane03	Mainline	51	1	5	1	30	1	3
5	2017-02-08	1899-12-31 00:38:00	TMS-408-07.0-EB	TMS-408-07.0-EB-Lane04	Ramp	53	1	1	1	3	1	3
6	2017-02-08	1899-12-31 00:39:00	TMS-408-07.0-EB	TMS-408-07.0-EB-Lane03	Mainline	59	1	1	1	4	1	3
7	2017-02-08	1899-12-31 00:42:00	TMS-408-07.6-EB	TMS-408-07.6-EB-Lane04	Ramp	43	1	1	1	4	1	3
8	2017-02-08	1899-12-31 00:47:00	TMS-408-06.8-EB	TMS-408-06.8-EB-Lane02	Mainline	53	1	2	1	9	1	3
9	2017-02-08	1899-12-31 00:47:00	TMS-408-07.0-EB	TMS-408-07.0-EB-Lane02	Mainline	55	1	2	1	9	1	3
10	2017-02-08	1899-12-31 01:25:00	TMS-408-08.0-EB	TMS-408-08.0-EB-Lane03	Mainline	72	1	1	1	2	1	3
11	2017-02-08	1899-12-31 07:50:00	TMS-408-06.8-EB	TMS-408-06.8-EB-Lane01	Mainline	0	1	30	1	17	1	3
12	2017-02-06	1899-12-31 05:11:00	TMS-408-07.4-EB	TMS-408-07.4-EB-Lane03	Mainline	70	1	1	1	2	1	3
13	2017-02-10	1899-12-31 07:22:00	TMS-408-08.0-EB	TMS-408-08.0-EB-Lane01	Mainline	64	0	30	0	114	1	1
14	2017-02-09	1899-12-31 08:18:00	TMS-408-07.6-EB	TMS-408-07.6-EB-Lane03	Mainline	170	1	29	0	12	0	1
15	2017-02-09	1899-12-31 08:21:00	TMS-408-06.8-EB	TMS-408-06.8-EB-Lane02	Mainline	180	1	29	0	14	0	1
16	2017-02-08	1899-12-31 07:05:00	TMS-408-08.0-EB	TMS-408-08.0-EB-Lane01	Mainline	65	0	30	0	120	1	1
17	2017-02-08	1899-12-31 07:13:00	TMS-408-08.0-EB	TMS-408-08.0-EB-Lane01	Mainline	63	0	30	0	-30	1	1
18	2017-02-08	1899-12-31 08:35:00	TMS-408-07.4-EB	TMS-408-07.4-EB-Lane02	Mainline	65	0	200	1	15	0	1
19	2017-02-07	1899-12-31 07:14:00	TMS-408-06.8-EB	TMS-408-06.8-EB-Lane01	Mainline	-30	1	30	0	14	0	1
20	2017-02-07	1899-12-31 07:30:00	TMS-408-06.8-EB	TMS-408-06.8-EB-Lane01	Mainline	60	0	-20	1	17	0	1
21	2017-02-06	1899-12-31 07:25:00	TMS-408-06.8-EB	TMS-408-06.8-EB-Lane01	Mainline	34	0	60	1	32	0	1

图 2：异常数据样本集

date	timestamp	detector_id	lane_id	speed	s_ab	volume	v_ab	occupancy	o_ab	error
8	2017-02-10	1899-12-31 00:00:00	TMS-408-07.4-EB-Lane01	-1	1	-1	1	-1	1	3
18	2017-02-10	1899-12-31 00:00:00	TMS-408-08.0-EB	-1	1	-1	1	-1	1	3
32	2017-02-10	1899-12-31 00:01:00	TMS-408-07.6-EB	-1	1	-1	1	-1	1	3
37	2017-02-10	1899-12-31 00:02:00	TMS-408-06.8-EB	-1	1	-1	1	-1	1	3
40	2017-02-10	1899-12-31 00:02:00	TMS-408-07.0-EB	-1	1	-1	1	-1	1	3
44	2017-02-10	1899-12-31 00:02:00	TMS-408-07.4-EB	-1	1	-1	1	-1	1	3
54	2017-02-10	1899-12-31 00:02:00	TMS-408-08.0-EB-Rmp-Lane04	-1	1	-1	1	-1	1	3
55	2017-02-10	1899-12-31 00:03:00	TMS-408-06.8-EB	-1	1	-1	1	-1	1	3
61	2017-02-10	1899-12-31 00:03:00	TMS-408-07.0-EB	-1	1	-1	1	-1	1	3
62	2017-02-10	1899-12-31 00:03:00	TMS-408-07.4-EB	-1	1	-1	1	-1	1	3
65	2017-02-10	1899-12-31 00:03:00	TMS-408-07.6-EB	-1	1	-1	1	-1	1	3
73	2017-02-10	1899-12-31 00:04:00	TMS-408-06.8-EB	-1	1	-1	1	-1	1	3
74	2017-02-10	1899-12-31 00:04:00	TMS-408-06.8-EB	-1	1	-1	1	-1	1	3
76	2017-02-10	1899-12-31 00:04:00	TMS-408-07.0-EB	-1	1	-1	1	-1	1	3
79	2017-02-10	1899-12-31 00:04:00	TMS-408-07.0-EB	-1	1	-1	1	-1	1	3
80	2017-02-10	1899-12-31 00:04:00	TMS-408-07.4-EB	-1	1	-1	1	-1	1	3
83	2017-02-10	1899-12-31 00:04:00	TMS-408-07.6-EB	-1	1	-1	1	-1	1	3
86	2017-02-10	1899-12-31 00:04:00	TMS-408-07.6-EB	-1	1	-1	1	-1	1	3
87	2017-02-10	1899-12-31 00:04:00	TMS-408-08.0-EB	-1	1	-1	1	-1	1	3
93	2017-02-10	1899-12-31 00:05:00	TMS-408-06.8-EB	-1	1	-1	1	-1	1	3
94	2017-02-10	1899-12-31 00:05:00	TMS-408-07.0-EB	-1	1	-1	1	-1	1	3
96	2017-02-10	1899-12-31 00:05:00	TMS-408-07.0-EB	-1	1	-1	1	-1	1	3
97	2017-02-10	1899-12-31 00:05:00	TMS-408-07.0-EB	-1	1	-1	1	-1	1	3
98	2017-02-10	1899-12-31 00:05:00	TMS-408-07.4-EB	-1	1	-1	1	-1	1	3

图 3：缺失数据样本集

(3) 对于缺失值的填补

表格中出现缺失值的情况主要有两种原因：

a. 第二问中，由于检测器异常导致删除的异常值，造成数据的缺失。
 b. 由于检测自身的误差导致数据缺失，存在的异常值
 修补方式为采用相同车道在时间或者空间上最接近的两点的正常值之和的平均值来进行修复。设置函数`data_correct`与`data_find`
`data_correct`:寻找靠近错误数据的最近的两个正常值的平均值来代替错误数据
`data_find`:用于进行修正
 其中`data_correct`的逻辑判断过程为：

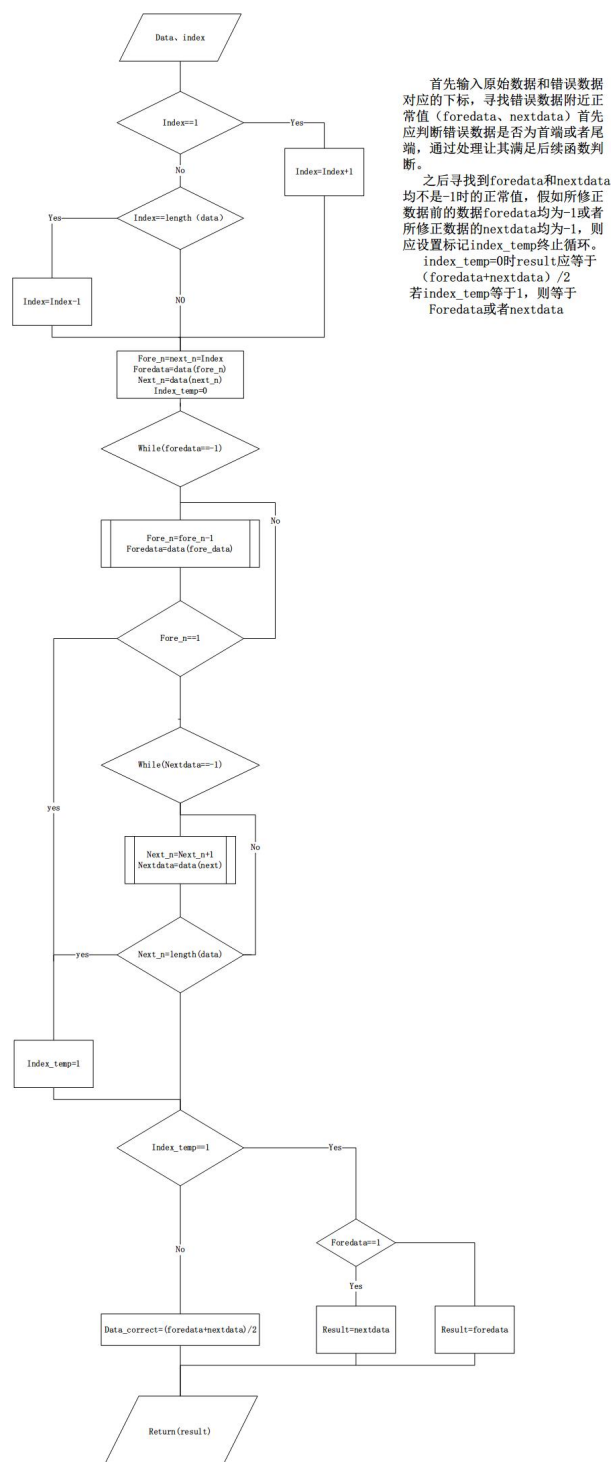


图 4：数据修正中`data_correct`函数思路

	date	timestamp	detector_i	lane_id	type	speed	volume	occupanc
1	2017/2/10	1899-12-31 07:22:00	TMS-408	TMS-408	Mainline	64	30	11.5
2	2017/2/9	1899-12-31 11:01:00	TMS-408	TMS-408	Mainline	71	7.5	2.5
3	2017/2/8	1899-12-31 07:05:00	TMS-408	TMS-408	Mainline	65	30	10
4	2017/2/8	1899-12-31 07:13:00	TMS-408	TMS-408	Mainline	63	30	14
5	2017/2/8	1899-12-31 07:50:00	TMS-408	TMS-408	Mainline	56.5	30	16
6	2017/2/7	1899-12-31 07:14:00	TMS-408	TMS-408	Mainline	64	30	14
7	2017/2/7	1899-12-31 07:30:00	TMS-408	TMS-408	Mainline	60	30	17
8	2017/2/6	1899-12-31 07:25:00	TMS-408	TMS-408	Mainline	34	28	32
9	2017/2/9	1899-12-31 08:21:00	TMS-408	TMS-408	Mainline	46	29	14
10	2017/2/8	1899-12-31 00:47:00	TMS-408	TMS-408	Mainline	58	2	5.5
11	2017/2/8	1899-12-31 00:47:00	TMS-408	TMS-408	Mainline	56.5	2.5	5
12	2017/2/8	1899-12-31 08:35:00	TMS-408	TMS-408	Mainline	65	27	15
13	2017/2/9	1899-12-31 08:18:00	TMS-408	TMS-408	Mainline	54	29	12
14	2017/2/9	1899-12-31 11:01:00	TMS-408	TMS-408	Mainline	61.5	14	8
15	2017/2/8	1899-12-31 00:39:00	TMS-408	TMS-408	Mainline	54.5	1	1
16	2017/2/8	1899-12-31 01:25:00	TMS-408	TMS-408	Mainline	61.5	1.5	0
17	2017/2/6	1899-12-31 05:11:00	TMS-408	TMS-408	Mainline	61	2.5	1.5
18	2017/2/8	1899-12-31 00:38:00	TMS-408	TMS-408	Ramp	36	2	2
19	2017/2/8	1899-12-31 00:42:00	TMS-408	TMS-408	Ramp	45	1.5	2.5
20	2017/2/10	1899-12-31 05:51:00	TMS-408	TMS-408	Ramp	52	1.5	0.5
21	2017/2/10	1899-12-31 09:01:00	TMS-408	TMS-408	Ramp	44	8.5	7.5

图 5:a 类数据修正

B 类数据由于数据量过多不在此展示，两者错误数据修正的结果写入 *Midterm_data_error.csv* 中

修正之后的完整样本集输出 *Midterm_data_correct.csv* 中

Part 02

数据说明

本数据取自一段快速路 3 年的事故相关数据, 包含 case(事件是否是事故, 1=crash, 0=non-crash), rain(事件发生时是否下雨, 1=rain, 0=non-rain), speed(事件的速度)。

处理要求

完成事故-下雨-速度, 三者之间的两两相关性分析, 给出具体的过程, 计算检验值, 并判断相关性是否显著。不建议采用程序自动计算的方法, 可以通过具体的公式, 一步一步计算。

处理过程

2.1 判断是否发生事故和是否下雨之间的相关性分析

H0: 假设发生事故和事件发生时下雨相互独立

H1: 假设发生事故和事件发生时下雨有关联

利用卡方检验来分析：

是否发生事故	事件是否下雨		合计
	1 (是)	0 (否)	
1 (是)	30	95	125
0 (否)	45	567	612
	75	662	737

此时可以考虑利用卡方检验，公式如下

$$\chi^2 = \sum_i^j \frac{(E_{ij} - T_{ij})^2}{E_{ij}}$$

是否发生事故	事件是否下雨		合计
	1 (是)	0 (否)	
1 (是)	12.72	112.28	125
0 (否)	62.27	549.72	612
	75	662	737

计算可得

$$\chi^2 = 31.4670$$

查表可得， $p_{value} < 0.00001$ 时， $\chi^2 = 29.6743$ ，因此有理由相信发生事故和发生事件下雨存在相关性关系。

2.2 判断是否发生事故和速度之间是否存在相关性

Value	Speed		Speed	
	Cash	Noncash	Rain	Non_rain
Mean	23.1	32.12	28.6	30.82
SD	10.76	13.87	12.55	13.94
Min	1.73	5.73	5.73	1.73
Max	67	85	70	85
t-value	8.09			

$$N_{cash} = 125 \quad N_{noncash} = 612$$

先做假设：

H_0 : 发生事故和未发生事故的速度相互独立

H_1 : 发生事故和未发生事故的速度相互关联

首先需要进行方差齐性检验：

$$H_0: \sigma_1^2 = \sigma_2^2$$

$$H_1: \sigma_1^2 \neq \sigma_2^2$$

$$F = \frac{\sigma_1^2}{\sigma_2^2} = 1.6616$$

$$F(611, 124)_{threshold} = 1.2712$$

(利用matlab: `finv(0.95, 611, 124)`)

说明两组数据存在显著性差异，即方差不相等.之后进行 t 检验：

$$t_{value} = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\left(\frac{s_1^2}{N_1} + \frac{s_2^2}{N_2}\right)}} = 8.09$$

其自由度为：

$$dt = \frac{\left(\frac{s_1^2}{N_1} + \frac{s_2^2}{N_2}\right)^2}{\frac{\left(\frac{s_1^2}{N_1}\right)^2}{N_1 - 1} + \frac{\left(\frac{s_2^2}{N_2}\right)^2}{N_2 - 1}} = 217$$

查表可得双侧置信度为 95%时

$$t(217)_{threshold} = 1.971$$

(利用matlab: `tinv(0.975,217)`)

因此可以认为发生事故和速度存在关系。

2.3 判断是否下雨和事件的速度是否有关系

$$N_{rain} = 75 \quad N_{nonrain} = 662$$

同 2.2 先做方差齐性检验

$$F' = 1.23 < F(661,74)_{threshold} = 1.3575$$

(利用matlab: `finc(0.95,661,74)`)

说明两组数据并不存在显著性差异，即满足方差相等

$$S_p = 190.616 \quad t_{value} = 1.31977$$

其自由度

$$N = N_1 + N_2 - 2 = 735$$

查表可得双侧置信度为 95%时

$$t(735)_{threshold} = 1.9632$$

(利用matlab: `tinv(0.975,735)`)

说明不可以认为是是否下雨和速度存在关系

附录：

Midterm_assignment.xlsx 原始数据

New_Midterm_data.csv 包含错误和缺失的原始数据

Midterm_data_errort.csv 错误和缺失数据样本集

Midterm_data_correct.csv 改正后的完整样本集

Midterm_correct.R 第一问相关代码

```
1. library(readxl)
2. library(stringr)
3. Midterm_assignment <- read_excel("Midterm_assignment.xlsx")
4. Midterm_assignment=Midterm_assignment[, -5]#type 无用，为了后续处理方便删除
   type
5. unique(Midterm_assignment$lane_id)#观察有多少车道
6.
7. ##part01
8. #1.1 读取数据，分*_Lane01\*_Lane02\*_Lane03\*07.0-EB-Lane04\*07.6-EB-
   Lane04\*08.0-EB-Rmp-Lane04
9. data_lane01=Midterm_assignment[endsWith(Midterm_assignment$lane_id, 'Lane01')
,]
```



```
10. data_lane02=Midterm_assignment[endsWith(Midterm_assignment$lane_id,'Lane02')
,]
11. data_lane03=Midterm_assignment[endsWith(Midterm_assignment$lane_id,'Lane03')
,]
12. data_Ramp1=Midterm_assignment[endsWith(Midterm_assignment$lane_id,'07.0-EB-
Lane04'),]
13. data_Ramp2=Midterm_assignment[endsWith(Midterm_assignment$lane_id,'07.6-EB-
Lane04'),]
14. data_Ramp3=Midterm_assignment[endsWith(Midterm_assignment$lane_id,'08.0-EB-
Rmp-Lane04'),]
15. data_standard=function(data){
16.   standard1_data=(data-min(data))/(max(data)-min(data))
17.
18.   sigma_data=sum(abs(data-mean(data)))/length(data)
19.   standard2_data=(data-mean(data))/sigma_data
20.
21.   result=vector()
22.   result[1]=mean(data)
23.   result[2]=sd(data)
24.   result[3]=mean(standard1_data)
25.   result[4]=sd(standard1_data)
26.   result[5]=mean(standard2_data)
27.   result[6]=sd(standard2_data)
28.   return(result)
29. }
30. #计算均值
31. data_standard(data_lane01$speed)
32. data_standard(data_lane02$speed)
33. data_standard(data_lane03$speed)
34. data_standard(data_Ramp1$speed)
35. data_standard(data_Ramp2$speed)
36. data_standard(data_Ramp3$speed)
37.
38. #1.2 异常值的处理（对不同车道实行阈值法）
39.
40. #生成完整的没有缺失数据的数据框
41. data_date=rep(unique(Midterm_assignment$date),each=721*18)
42.
43. temp_time=unique(Midterm_assignment$timestamp)
44. temptime=temp_time[c(1:100,721,101:720)]
45. data_time=rep(temptime,times=5,each=18)
46.
47. temp_lane=unique(Midterm_assignment$lane_id)
```

```
48. temp_lane=temp_lane[c(1,2,3,4,5,6,7,17,8,9,10,11,12,13,14,15,16,18)]#调整顺序
49. data_lane=rep(temp_lane,times=721*5)
50.
51. data_detector=str_sub(data_lane,1,15)
52. data_speed=rep(-1,times=5*721*18)
53. data_volume=rep(-1,times=5*721*18)
54. data_occupancy=rep(-1,times=5*721*18)
55.
56. New_Midterm_data=data.frame(date=data_date,timestamp=data_time,detector_id=data_detector,
57.                               lane_id=data_lane,speed=data_speed,volume=data_volume,occupancy=data_occupancy)
58.
59. for(n in 1:(5*721*18)){
60.   #将现有的数据加到新建的数据框中
61.   logic=(Midterm_assignment$date==data_date[n]&Midterm_assignment$timestamp==data_time[n]&Midterm_assignment$lane_id==data_lane[n])
62.   if(sum(logic)==1){
63.     New_Midterm_data[n,]=Midterm_assignment[logic,]
64.   }
65.   print(n)
66. }
67.
68. judge1=function(data){
69.   #五分位数和三标准差法判断范围过大，采用单独的判断、联合判断和有效车长来进行判断，只有第一种可以区分三个参数中唯一出错的
70.   s_ab=(data$speed<0|data$speed>112)
71.   v_ab=(data$volume<0|data$volume>51)
72.   o_ab=(data$occupancy<0|data$occupancy>95)
73.
74.   temp1=data$speed==0
75.   temp2=data$volume==0
76.   temp3=data$occupancy==0
77.   temp=temp1+temp2+temp3
78.
79.   add_01=((temp==3)|(temp==0)|data$occupancy>95|(data$volume<5&data$occupancy==0))|((s_ab+v_ab+o_ab)!=0))==0
80.
81.   avel=data$speed*data$occupancy/(6*data$volume)
82.   add_02=(avel<1.5|avel>22)&((data$volume<5&data$occupancy==0)==0)&((s_ab+v_ab+o_ab+add_01)==0)
83.
84.   s_ab=s_ab+add_01+add_02
```

```
85. v_ab=v_ab+add_01+add_02
86. o_ab=o_ab+add_01+add_02
87.
88. data['s_ab']=as.integer(s_ab!=0)
89. data['v_ab']=as.integer(v_ab!=0)
90. data['o_ab']=as.integer(o_ab!=0)
91. data['error']=s_ab+v_ab+o_ab
92.
93. return(data)
94. }
95. New_Midterm_data=judge1(New_Midterm_data)
96. New_Midterm_data=New_Midterm_data[c(1,2,3,4,5,8,6,9,7,10,11)]#调整参数
97. write.csv(New_Midterm_data,file='New_Midterm_data.csv')
98.
99. #1.3 缺失值的填补
100. data_find=function(data,n){
101.     #用来寻找异常值或者缺失值时空最接近点上正常的数
102.
103.     #首先确定 n 是否为第一位或者最后以为
104.     if(n==1){
105.         n=n+1
106.     }
107.     if(n==length(data)){
108.         n=n-1
109.     }
110.     #初始化前值和后值，以及初始下标和标记位
111.     foredata=data[n-1]
112.     nextdata=data[n+1]
113.     fore_n=n
114.     next_n=n
115.     temp_fore=0
116.     temp_next=0
117.     #函数主体循环，说明
118.     while(foredata!=-1){
119.         fore_n=fore_n -1
120.         #说明前面没有正确的值
121.         if(length(fore_n)==0){
122.             temp_fore=1
123.             break
124.         }
125.         foredata=data[fore_n]
126.         if(length(foredata)==0){
127.             foredata=0.01
128.         }
```

```
129. }
130. while(nextdata==-1){
131.     next_n=next_n+1
132.     if(nextdata>length(data)){
133.         temp_next=1
134.         break
135.     }
136.     nextdata=data[next_n]
137.     if(length(nextdata)==0){
138.         nextdata=0.01
139.     }
140. }
141. if((temp_fore!=1)&(temp_next!=1)){
142.     result=(foredata+nextdata)/2
143. }else{
144.     if((temp_fore==1)&(temp_next!=1)){
145.         result=nextdata
146.     }else{
147.         if((temp_fore!=1)&(temp_next==1)){
148.             result=foredata
149.         }
150.     }
151. }
152. return(result)
153. }
154.
155. data_correct=function(data){
156.     #不同进行修正
157.     index_speed=which(data$s_ab==1)
158.     index_volume=which(data$v_ab==1)
159.     index_occupancy=which(data$o_ab==1)
160.     for(index in index_speed)
161.         data$speed[index]=data_find(data$speed,index)
162.     for(index in index_volume)
163.         data$volume[index]=data_find(data$volume,index)
164.     for(index in index_occupancy)
165.         data$occupancy[index]=data_find(data$occupancy,index)
166.     return(data)
167. }
168.
169. data_lane01=data_correct(New_Midterm_data[endsWith(New_Midterm_data$lane_id
    , 'Lane01'),])
170. data_lane02=data_correct(New_Midterm_data[endsWith(New_Midterm_data$lane_id
    , 'Lane02'),])
```

```
171. data_lane03=data_correct(New_Midterm_data[endsWith(New_Midterm_data$lane_id
, 'Lane03'),])
172. data_Ramp1=data_correct(New_Midterm_data[endsWith(New_Midterm_data$lane_id,
'07.0-EB-Lane04'),])
173. data_Ramp2=data_correct(New_Midterm_data[endsWith(New_Midterm_data$lane_id,
'07.6-EB-Lane04'),])
174. data_Ramp3=data_correct(New_Midterm_data[endsWith(New_Midterm_data$lane_id,
'08.0-EB-Rmp-Lane04'),])
175.
176. data_after_correct=rbind(data_lane01,data_lane02,data_lane03,data_Ramp1,dat
a_Ramp2,data_Ramp3)
177. data_error=data_after_correct[data_after_correct$error!=0,]
    178. write.csv(data_error,file='Midterm_data_error.csv')
```