

油门性能分布研究及异常检测

该部分分析基于北京朝阳艾瑞特2号的3-5月的日常行驶数据进行, 主要涉及:

- 从三至五月的日常行驶数据中提取出车辆连续平稳行驶的数据
- 对车辆平稳行驶时, 油门值的数据分布进行探索, 并检查是否存在time decay或系统上的前后分布差异
- 根据分布 $P(\text{油门值}|\text{车辆平稳行驶})$, 进一步构造车一个异常值监控的过程

1. 车辆平稳行驶数据

在这里, 我们对于车辆油门性能, 主要希望考究车辆保持连续平稳状态是, 发的油门值指令是否存在分布上的一致。倘若为了保持同一水平的平稳行驶, 车辆需要的油门值发生了较大的改变或者存在着明显的time decay, 则需要及时对器械进行相关调整, 使其性能维持在正确的标定水平。

基于这一想法, 我们首先需要提取出所有行驶数据中车辆处于连续平稳行驶状态的数据。对于北京朝阳的艾瑞特二号, 我们可以从**对可视化的观察中**和**对所有速度取值的分布探索**中发现, 车辆目标的平稳行驶速度大致在 1.4-1.5 m/s。平稳行驶时, 车辆的速度和油门值都各自处于某一水平, 实际的观测会以这一水平为中心上下波动。基于相关背景知识, 对此我们可以假设平稳行驶的速度(velocity)和油门值(command)都来自某一正态分布。

考虑到平稳行驶时, 车辆采用的PID控制方式, 使得车辆的油门值与速度存在很大的负相关性, 因此为了更加准确地刻画二者的联合分布, 我假设了车辆平稳行驶时, 速度与油门值来自某一个二维正态分布(加入了协方差信息), 即

$$(v, cmd|acc \approx 0)^T \sim \mathcal{N}((\mu_v, \mu_{cmd})^T, \Sigma_{v,cmd})$$

原本我希望通过对于速度和油门值取上下界来截取平稳行驶的数据, 即提取速度 $v_{lower} < v < v_{upper}$, $cmd_{lower} < cmd < cmd_{upper}$, 然后通过truncated normal利用截断的数据来对原来的二维正态分布进行估计, 但查过很多论文后, 我发现对于多维截断正态分布数据的估计并不是一个trivial的问题, 需要耗费大量的精力与计算资源。于是我退而求其次, 取了较大的截断范围(即值截去正态分布的非常末端的尾部数据), 然后直接利用截断的数据对于整体进行估计。虽然这样得到的估计显然是有偏的, 但由于截取的范围比较大, 实际的影响也很小。

基于上述的分析, 我们利用了如下逻辑来提取车辆平稳行驶数据。

- 用一个矩形去截取得到数据取值(v, cmd) 属于[1.3, 1.6] X[20, 30]内的数据
- 滤除加速度为0(传感器异常)和system_status!= 1(人工接管)
- 滤除方向盘状态值(steering_status)和加速度(acceleration)过大的数据
- 保留连续平稳行驶时长高于5s的数据

2. 数据分批

为了检查车辆的油门性能是否存在time decay或者前后差异, 我将数据依照时间顺序均分成100个批次(batch), 然后分别对每个批次用二维正态分布进行估计。比较100个分布的差异, 来判断是否存在time decay或者前后差异。

下图是其中某三个批次的数据散点图

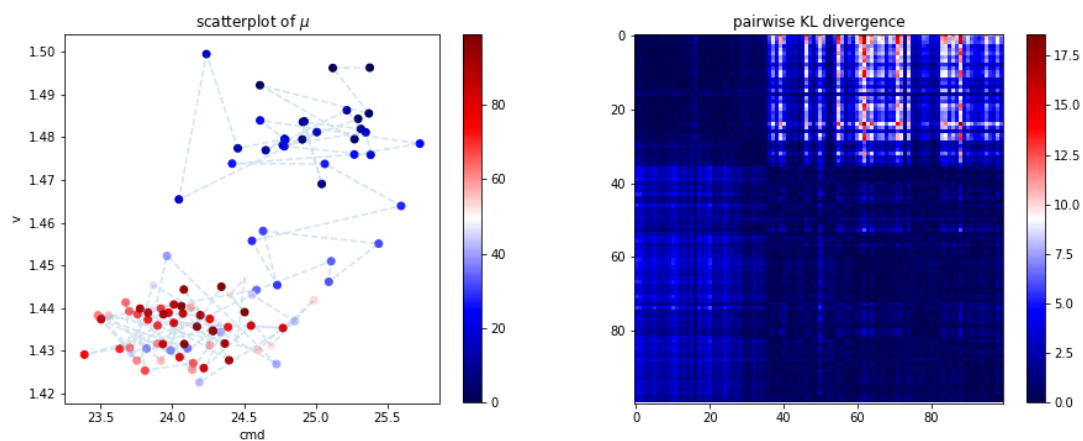


3. 联合分布

对每一个批次，我们都可以用样本均值，样本协方差来估计联合正态分布的均值和方差，得到其联合分布：

$$P_i(v, cmd | acc \approx 0) = N(\mu_i, \Sigma_i)$$

我们可以通过参数估计得到每个批次的均值 μ 和协方差矩阵 Σ ，然后我分别画了 μ 随时间变化的轨迹图，以及批次之间两两的KL散度（可以用以表示两个分布之间的差异）的热力图，其结果如下



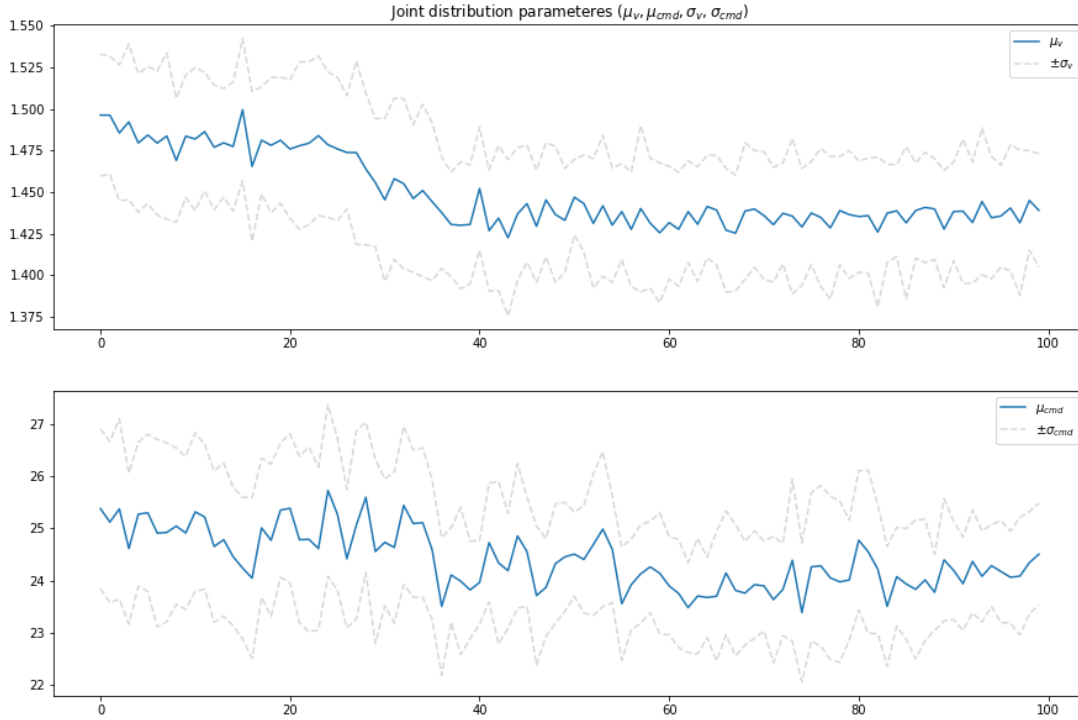
1. 由左图我们可以看到，联合分布的均值位置随着时间变化(批次序号增大，颜色由蓝至红，逐渐从右上角迁移到了左下角)，呈现出较为明显的2个总体。
2. 右图的横纵坐标都为批次序号，颜色表示分布之间的KL散度的取值大小，从图中我们也能很明显能感受到2个聚类总体（0-36，37-99）。

综上，通过对联合分布的观察，我们不难发现，表示油门性能的参数的分布似乎在第36个批次(4月16日)前后，有较大的改变。

注：对于正态分布的KL散度计算，有如下公式：

$$D_{KL}(P_1||P_2) = \frac{1}{2} \left[\log \frac{|\Sigma_2|}{|\Sigma_1|} - d + \text{tr}(\Sigma_2^{-1}\Sigma_1) + (\mu_2 - \mu_1)^T \Sigma_2^{-1}(\mu_2 - \mu_1) \right]$$

联合分布的参数按批次(随时间)变化的趋势图如($\mu - 1 * \sigma$, $\mu + 1 * \sigma$)下



4. 条件分布

一个较为明显的问题是，如果我们相比较油门性能，我们应当看当车辆的速度恒定在某一特定值 v_{target} 时，油门指令 cmd 的分布变化，因此我们需要用到条件分布，即得到

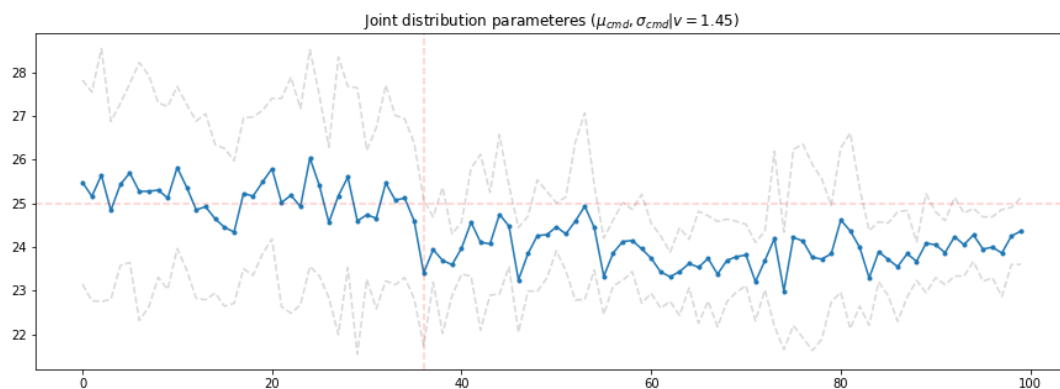
$$P(cmd|v = v_{target}, acc \approx 0)$$

由多维正态分布的性质，我们能够很容易地由下面的公式得到条件分布：

$$(X_1^T, X_2^T)^T \sim N((\mu_1^T, \mu_2^T)^T, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix})$$

$$(X_1|X_2) \sim N(\mu_1 - \Sigma_{11}^{-1}\Sigma_{12}(x_2 - \mu_2), \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21})$$

各个批次得到的条件分布参数如下图



5. 异常检测

对于异常值的检测，我们可以选用某最近一段使其的正常形式数据，得到上面的条件分布

$$P(cmd|v = v_{target}, acc \approx 0)$$

然后对于每一个给定的 v_{target} , 可以得到一个相对应的置信区间，如果观测落到置信区间之外，即表明：

这个观测与之前一段时间的观测可能来自于不同的总体，即可能是异常。