

构造更多质量好的特征! --Goal

特征

Pandas
Dataframe

公司利润	公司行业	市场情况			
100	冶金	好			
1	生物	不好			
2					

Numpy
Array

数据向量化 - 特征工程

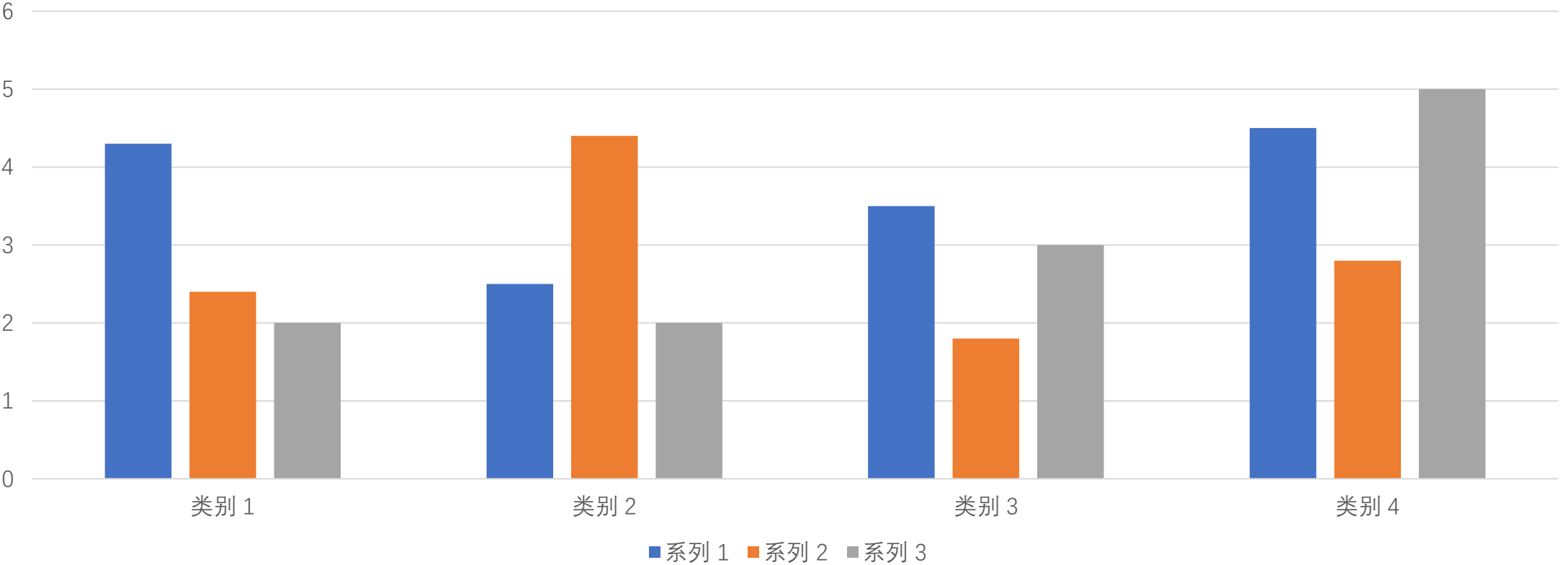


ML Model(KNN)

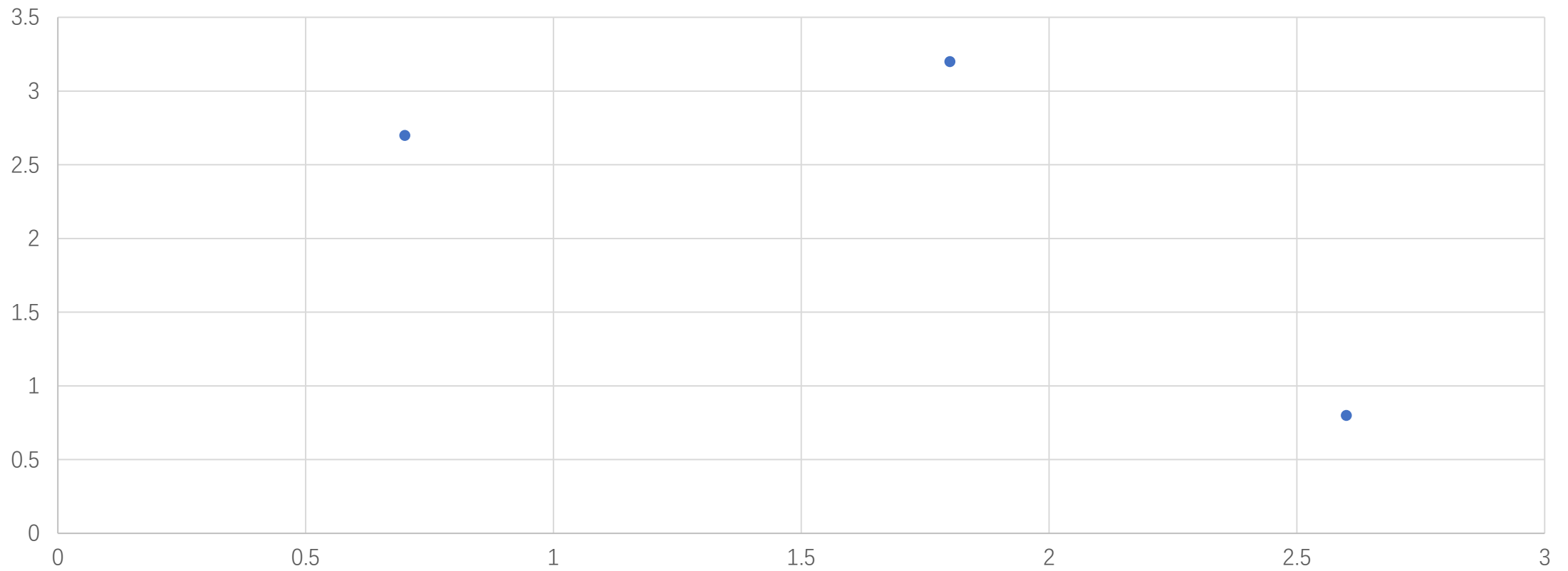


1

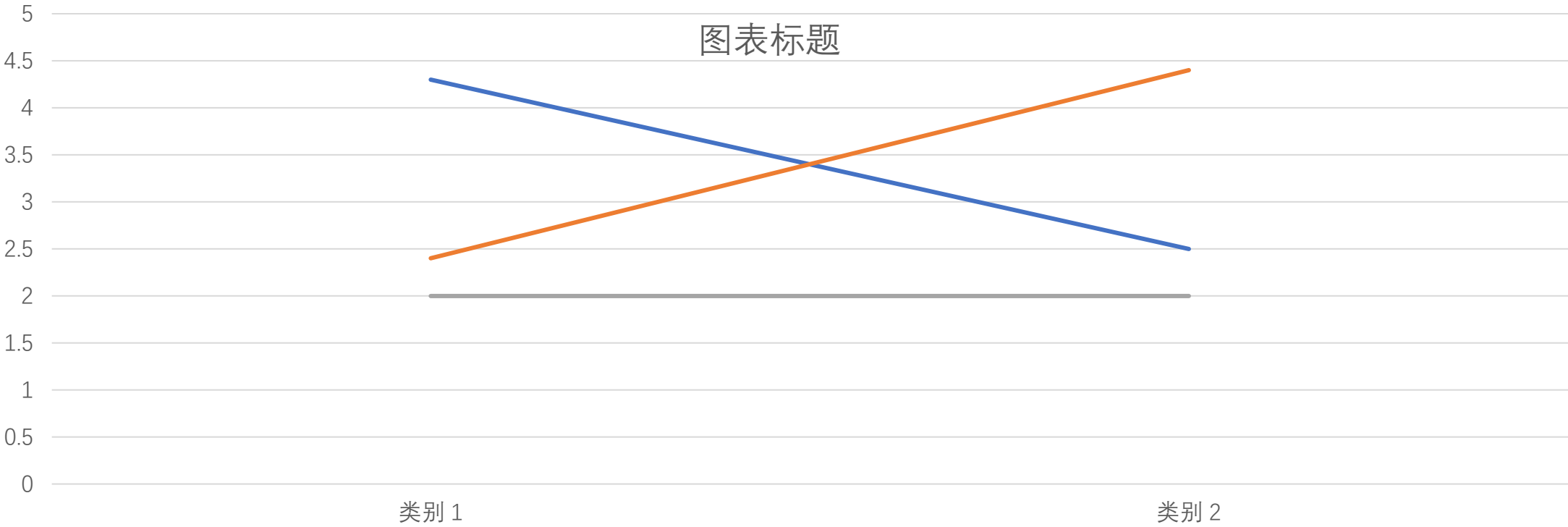
图表标题



Y 值

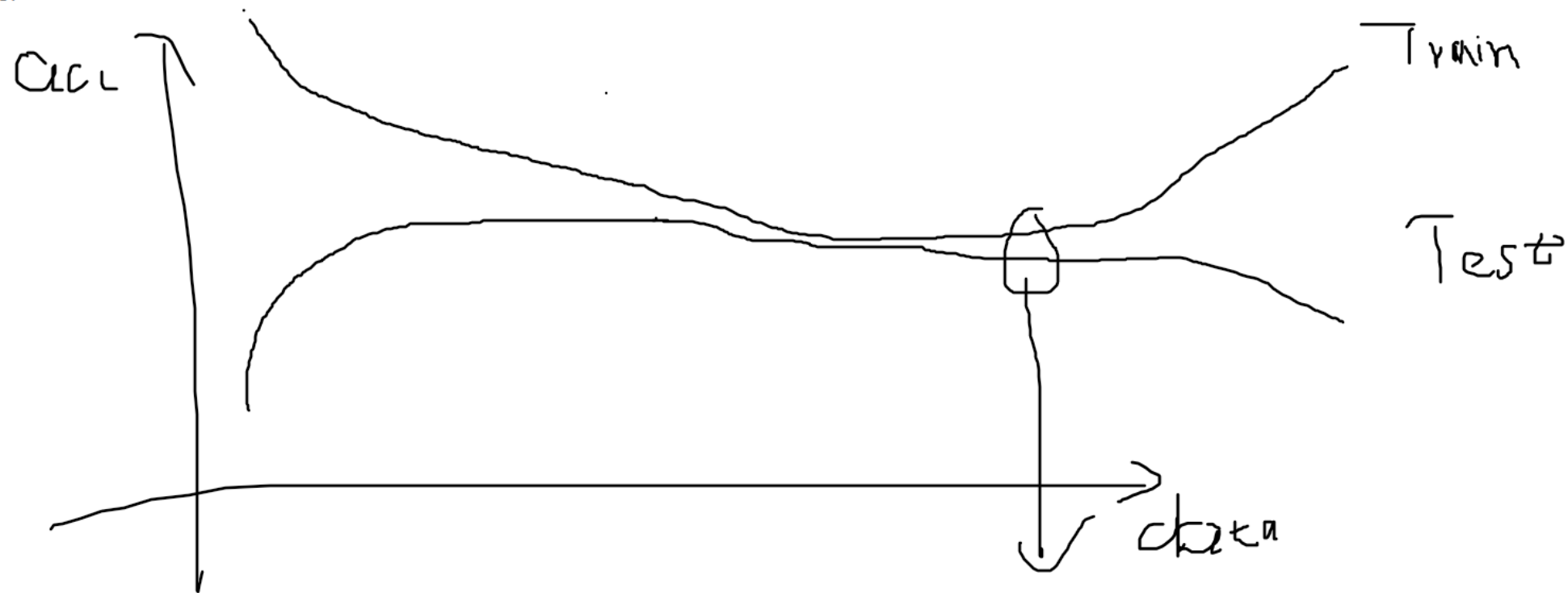


图表标题



— 系列 1 — 系列 2 — 系列 3

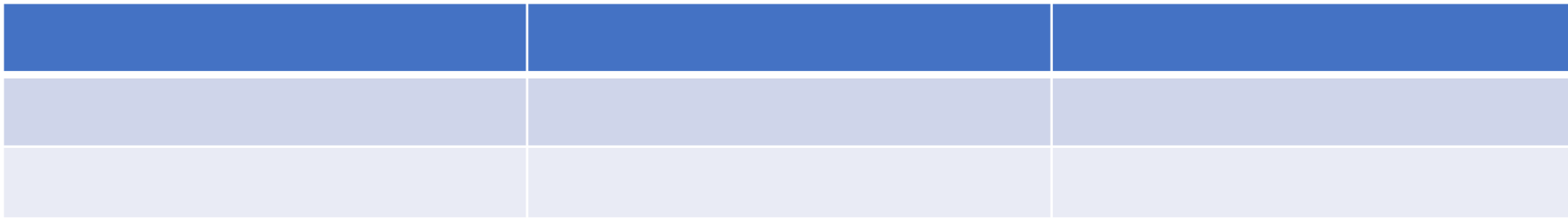
9:37



- $\text{Vec } a + \text{Vec } b = \text{Vec } c$
- For i in $\text{range}(10)$:
 - $c[i] = a[i] + b[i]$
- $\rightarrow c = a + b$



reshape

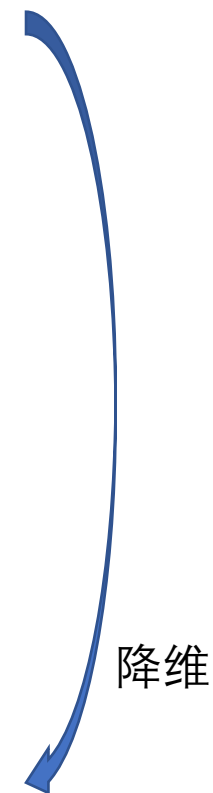


a	b	c	e					
---	---	---	---	--	--	--	--	--



特征选择

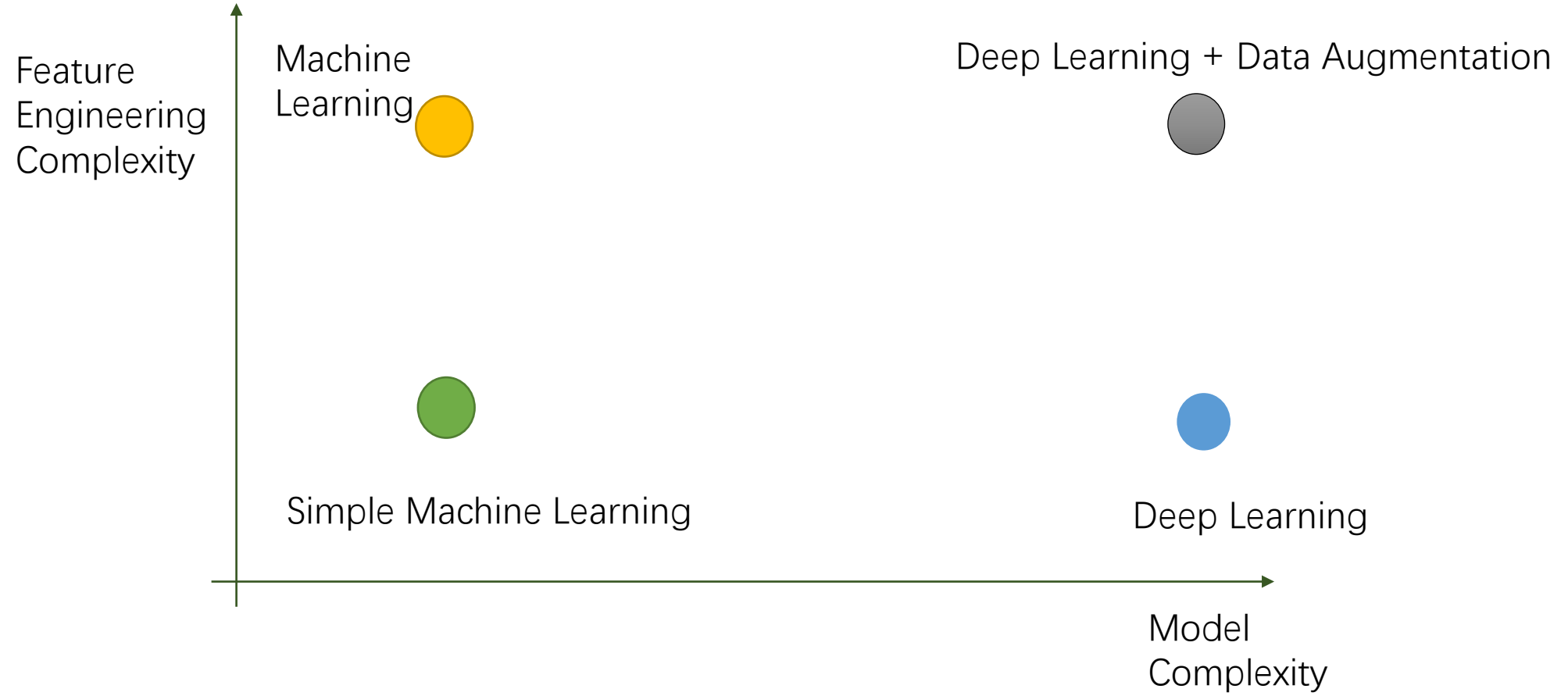
d	a	e

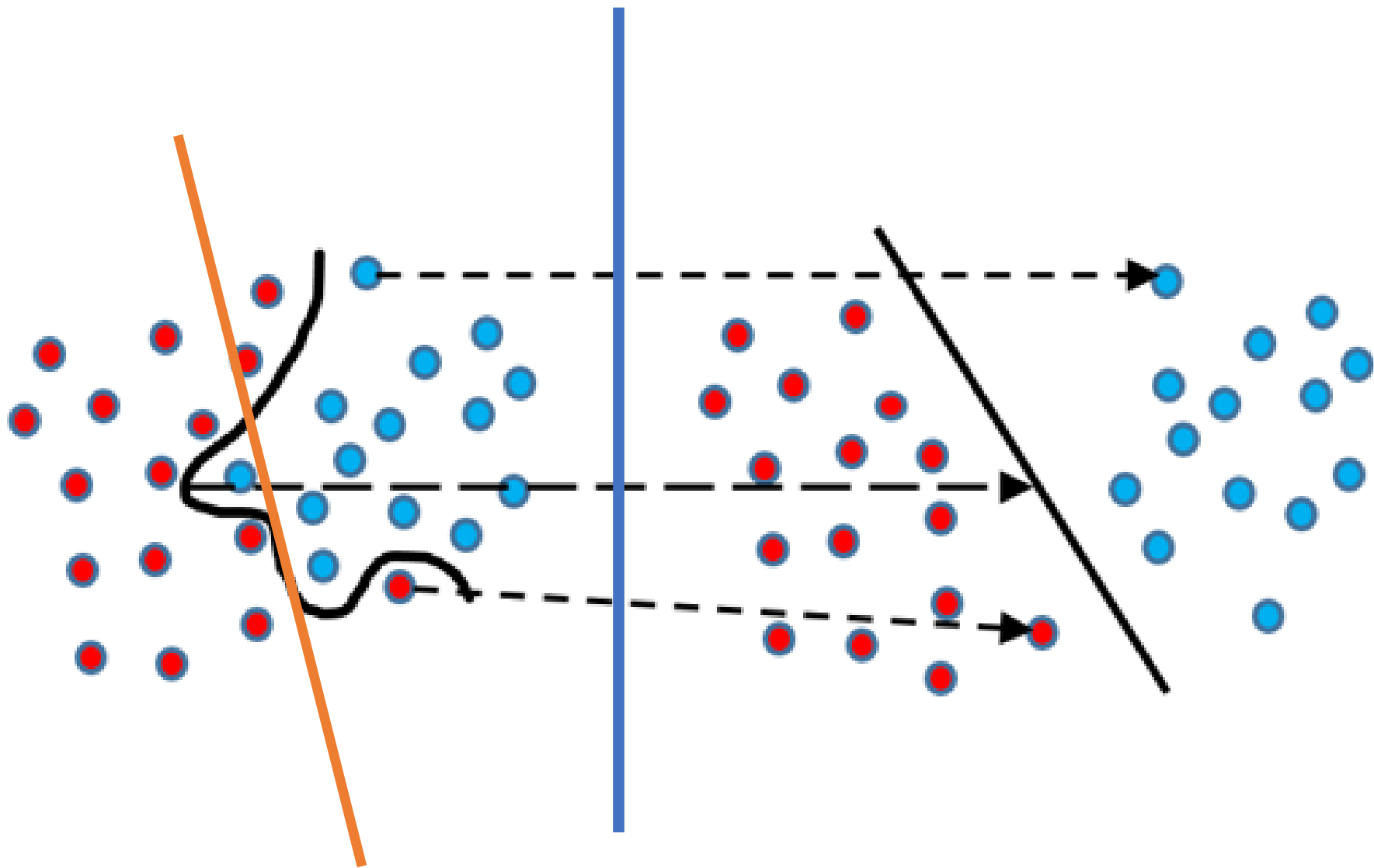


降维

- $Y = wX + b$

- 模型被保存的文件 = 权重 W, b + meta (函数形式)





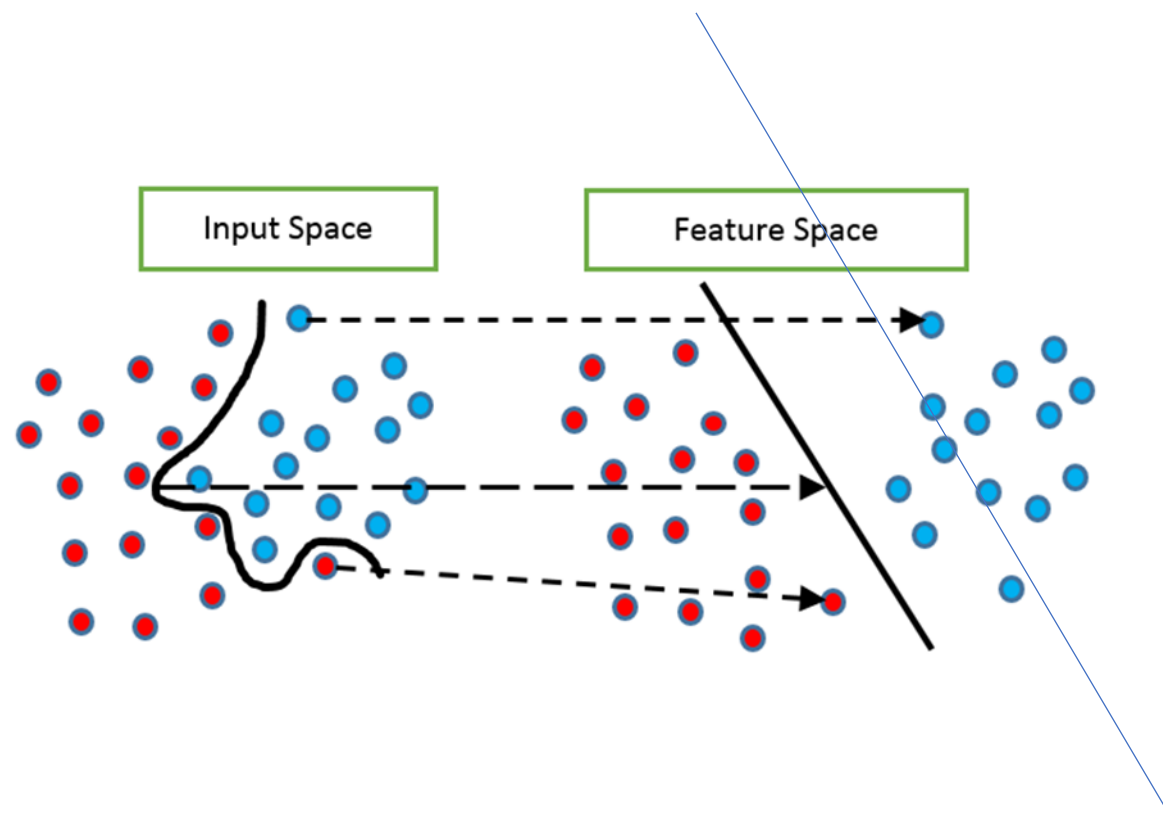
Data Augmentation



空值处理

- 缺失较多，可以丢弃
- 填充
 - 填充固定数值
 - 离散数值 – 众数或新数值
 - 连续数值 – 均值
- ML预测

- Label Encoder
 - Company_type: big, middle, small.
 - 0,1,2
- One encoder
 - Company_type -> c_t_a, c_t_b, c_t_c
 - [1,0,0], [0,1,0], [0,0,1]





0 100 500

Embedded feature selection

- $Y = WX + B$
 - W 反应特征重要性
- Decision tree



AutoML – HPO

