

啟事式:

- What is my likelihood
- What should my model be? interaction?
- What are my parameter & hyperparameter?
- What kind of priors I should choose?
- Are there any clustering, time or spatial dependence?

EDA

可以先畫不同國家之間, household expenditure & density
food vs non-food expenditure density

Data 前處理:

check micro model 中, y 變數 & X 之間有沒有相關性 (畫散佈圖)
類別轉 dummy \rightarrow `pd.get_dummies(X, prefix_sep="_", drop_first=False)`
Training / Testing Split \rightarrow `X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.15)`
Data Specification: `cars_data = {"N": X_train.shape[0], "N_new": X_test.shape[0],
("在 python 中要用 dictionary") "K": X_train.shape[1], "y_obs": y_train.values.tolist(),
"X": np.array(X_train), "X_new": np.array(X_test)}`

Modeling:

Stan Model 分几功能塊要 specify (其中 data, parameter & model 是必要)

```
cars_code = """
data {
  int <lower = 1> N;
  int <lower = 0> K;
  matrix[N, K] X;
  vector[N] y_obs;
  int <lower = 1> N_new;
  matrix[N_new, K] X_new;
}
parameter {
  real <lower = 0> sigma;
  real alpha;
  vector[K] beta;
}
transformed parameter {
  vector[N] theta;
  theta = alpha + X * beta;
}
model {
  sigma ~ exponential(1);
  alpha ~ normal(0, 6);
  beta ~ multi_normal(rep_vector(0, K), diag_matrix(rep_vector(1, K)));
  y_obs ~ normal(theta, sigma);
  y_new ~ normal(alpha + X_new * beta, sigma);
}
"""
```

我們可以在 Stan 模型中
用類似此處 new 代號方式
建立預測值。

如果你有導性模型
放這
如果你很單純
只想討論參數
與前設, 就不需要這塊

prior 放這
likelihood
`~ normal(x, mu, sigma)`

v.s.
theta 可以想成 training data
fit 出模型, 參數估計值
此處是 mean, 為最好
data 最好, 代表
也可用本做預測

跑模型: `sm = pystan.StanModel(model_code = cars_code)`

`fit = sm.sampling(data = cars_data, iter = 6000, chains = 8)`

儲存模型: `with open("bayes-cars.pkl", "wb") as f:`

(註不用每週
一次建一次)

`pickle.dump(sm, f, protocol = pickle.HIGHEST_PROTOCOL)`

Result: ① `la = fit.extract(permutated = True)` → 印出每一受估參數在每次iteration的估計值
② `print(fit)` → 給你一了fit出來的summary table (含 mean, std error, credible interval, n-eff, Rhat)

Diagnosis: ① Chain mixing

Visualization

↳ `az.plot_trace(fit, var_names=["O", "O"])`
可省

[note]: 若想看density可用

`az_data = az.from_pystan(posterior=fit)`
`az.plot_density(az_data, var_names=["O"])`

② posterior credible intervals

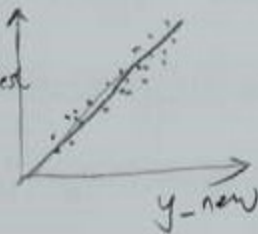
※ fit也可

↳ `az.plot_forest(az_data, kind="forestplot", var_names=[], combined=True)`

Prediction: $P(y_{new} | D_{old}, X_{new}) = \int P(y_{new} | \theta, D_{old}, X_{new}) P(\theta | D_{old}) d\theta$ 即為預測值

↳ ① 求 Bayes test MSE: `metrics.mean_squared_error(y_test, la["y_new"].mean())`

② 畫兩者 joint scatter plot: y_{test}



拿 y_{new} 和 y_{test} 做比較

③ 畫某受估對 y_{new} 的預測圖: `az.plot_hpd(某受估, la["y_new"], plot_kwargs={"ls": "--"})`
↑ 可省

	Train MSE	Test MSE
Bayesian	10.829	10.968
Frequentist / ML	10.747	10.558

model performance

更好, 衡量做法是用 bootstrap
做 train-test split ≥ 30 次
看 test MSE 95% confidence interval

[note]: 貝氏不見得比傳統 Maximum Likelihood (frequentist) 好, 特別備大量 dataset