# Stochastic Process

Lect 1:

# Introduction

## Frequentist and bayesian.

### Data Matrix

$$X = \begin{bmatrix} -x^{(1)}- \\ -x^{(2)}- \\ -x^{(N)}- \end{bmatrix} \in R^{N \times D}$$

where $x^{(i)} = [x_1^{(i)}, \cdots x_D^{(i)}] \in R^{D}$

we assume each $x^{(i)}$ is sampled from $P(D|\theta)$ in $\boxed{\text{i.i.d}}$ manner.

**Frequentist.:** assume $\theta$ is a constant.

Then the probability to observe $N$ data points in i.i.d manner is

$$P(D|\theta) = \prod_{i=1}^{N} P(x^{(i)}|\theta)$$

To calculate $\theta$, we can use MLE (Maximum likelihood estimator)

$$\hat{\theta}_{MLE} = \underset{\theta}{\arg\max} \, P(D|\theta)$$

$$= \underset{\theta}{\arg\max} \, \log P(D|\theta)$$

$$= \arg\max \sum_{i=1}^{N} \log P(x^{(i)}|\theta)$$

---

**Bayesian:** Assume $\theta$ is not a constant.

$\theta \sim P(\theta)$ ⟸ Preset prior distribution.

By Baye's Rule: the posterior distribution.

( the prob of $\theta$ given the evidence X )

$$P(\theta|D) = \frac{P(D|\theta) \, P(\theta)}{P(D)} \longrightarrow \text{Marginal likelihood}$$
$$\text{or evidence.}$$

$$\boxed{= \frac{P(D|\theta) \cdot P(\theta)}{\int_{\theta} P(D|\theta) \cdot P(\theta) d\theta.}}$$

To get $\theta$, we will maximize the posterior distribution.

$$\hat{\theta}_{MAP} = \underset{\theta}{\arg\max} \, P(\theta|D)$$

This is not fully "Bayesian" since This is not fully "Bayesian" since $\hat{\theta}_{MAP}$ is a point estimate.

$$\overset{\text{prevent overfitting}}{=} \underset{\theta}{\arg\max} \, P(D|\theta) \cdot P(\theta)$$

because $P(D) = \int_{\theta} P(D|\theta) \cdot P(\theta) \, d\theta$

is not a function of $\theta$.

$$= \underset{\theta}{\arg\max} \, \log P(D|\theta) + \log P(\theta).$$

posterior distribution. $P(\hat{\theta}_{MAP}|D)$

$$= \frac{\boxed{P(D|\hat{\theta}_{MAP}) \cdot P(\hat{\theta}_{MAP})}}{P(D)}$$

It is a function of $D$.

$\rightarrow$ likelihood function.

predictive uncertainty:

the uncertainty in the prediction

induced by uncertainty in the parameter

Compute posterior predictive distribution

$$P(y \mid x, D) \swarrow$$

$$= \int_{\theta} P(y \mid x, \underline{\theta}) \, P(\underline{\theta} \mid D) \, d\theta.$$

$\hookrightarrow$ Marginalizing out the parameter.
reduce the overfitting //