# Reading Notes for Mackey ch28 Model Comparison and Occam's Razor

Xiang Pan

September 21, 2021

## 1 Occam's Razor

Occam's razor: Accept the simplest explanation that fits the data. Reason:

- aesthetic ('A theory with mathematical beauty is more likely to be correct than an ugly one that fits some experimental data')

- past empirical success of Occam's razor

It is essential to use proper priors – otherwise the evidences and the Occam factors are not meaningful.

$$\frac{P\left(\mathcal{H}_{1} \mid D\right)}{P\left(\mathcal{H}_{2} \mid D\right)}=\frac{P\left(\mathcal{H}_{1}\right)}{P\left(\mathcal{H}_{2}\right)} \frac{P\left(D \mid \mathcal{H}_{1}\right)}{P\left(D \mid \mathcal{H}_{2}\right)}$$

Occam's Razor gives a favor of simplicity to the model.

**Model fitting**

$$P\left(\mathbf{w} \mid D, \mathcal{H}_{i}\right)=\frac{P\left(D \mid \mathbf{w}, \mathcal{H}_{i}\right) P\left(\mathbf{w} \mid \mathcal{H}_{i}\right)}{P\left(D \mid \mathcal{H}_{i}\right)}$$

$$\text{Posterior }=\frac{\text{Likelihood } \times \text{ Prior}}{\text{Evidence}} \tag{1}$$

**Model Comparison**

Models $H_i$ are ranked by evaluating the evidence,

$$P\left(D \mid \mathcal{H}_{i}\right) \simeq \underbrace{P\left(D \mid \mathbf{w}_{\mathrm{MP}}, \mathcal{H}_{i}\right)}_{} \times \underbrace{P\left(\mathbf{w}_{\mathrm{MP}} \mid \mathcal{H}_{i}\right) \sigma_{w \mid D}}_{}. \tag{2}$$

$$\text{Evidence } \simeq \text{ Best fit likelihood } \times \text{ Occam factor}$$

$$\text{Occam factor }=\frac{\sigma_{w \mid D}}{\sigma_{w}}$$

Occam factor is equal to the ratio of the posterior accessible volume of $H_i$'s parameter space to the prior accessible volume.

**Occam factor for several parameters** If the posterior is well approximated by a Gaussian, then the Occam factor is obtained from the determinant of the corresponding covariance matrix.

$$P\left(D \mid \mathcal{H}_{i}\right) \simeq P\left(D \mid \mathbf{w}_{\mathrm{MP}}, H_{i}\right) \times P\left(\mathbf{w}_{\mathrm{MP}} \mid \mathcal{H}_{i}\right) \operatorname{det}^{-\frac{1}{2}}(\mathbf{A} / 2 \pi) \tag{3}$$

$$\mathbf{A}=-\nabla \nabla \ln P\left(\mathbf{w} \mid D, \mathcal{H}_{i}\right) \tag{4}$$

**On-line learning and cross-validation.**

$$\log P(D \mid \mathcal{H})=\log P\left(\mathbf{t}^{(1)} \mid \mathcal{H}\right)+\log P\left(\mathbf{t}^{(2)} \mid \mathbf{t}^{(1)}, \mathcal{H}\right)$$

$$+\log P\left(\mathbf{t}^{(3)} \mid \mathbf{t}^{(1)}, \mathbf{t}^{(2)}, \mathcal{H}\right)+\cdots+\log P\left(\mathbf{t}^{(N)} \mid \mathbf{t}^{(1)} \ldots \mathbf{t}^{(N-1)}, \mathcal{H}\right) \tag{5}$$

Cross-validation examines the average value of just the last term. The evidence, on the other hand, sums up how well the model predicted all the data, starting from scratch.