

首先, 理解清楚信息矩阵的概念, 根据[wiki](#)的解释,

In mathematical statistics, the Fisher information (sometimes simply called information[1]) is a way of measuring the amount of information that an **observable random variable X** carries about an **unknown parameter θ** of a distribution that models X.

Fisher Information 的公式表达

The **variance** of the score is defined to be the **Fisher information**.^[5]

$$\mathcal{I}(\theta) = \mathbb{E} \left[\left(\frac{\partial}{\partial \theta} \log f(X; \theta) \right)^2 \middle| \theta \right] = \int \left(\frac{\partial}{\partial \theta} \log f(x; \theta) \right)^2 f(x; \theta) dx,$$

关于这个公式, 我们可以看出来, 信息矩阵所针对的是未知参数 θ , 与 x 没有直接联系。此外, 可以看出来, 括号内部的部分期望为0。

$$\begin{aligned} & \mathbb{E} \left[\frac{\partial}{\partial \theta} \log f(X; \theta) \middle| \theta \right] \\ &= \int \frac{\frac{\partial}{\partial \theta} f(x; \theta)}{f(x; \theta)} f(x; \theta) dx \\ &= \frac{\partial}{\partial \theta} \int f(x; \theta) dx \\ &= \frac{\partial}{\partial \theta} 1 = 0. \end{aligned}$$

在[Agustinus](#)博客中, 它给出了向量的表达, 以及对应离散的表达

We can then see it as an information. The covariance of score function above is the definition of Fisher Information. As we assume θ is a vector, the Fisher Information is in a matrix form, called Fisher Information Matrix:

$$\mathbf{F} = \mathbb{E}_{p(x|\theta)} \left[\nabla \log p(x|\theta) \nabla \log p(x|\theta)^T \right].$$

However, usually our likelihood function is complicated and computing the expectation is intractable. We can approximate the expectation in \mathbf{F} using empirical distribution $\hat{q}(x)$, which is given by our training data $X = \{x_1, x_2, \dots, x_N\}$. In this form, \mathbf{F} is called Empirical Fisher:

$$\mathbf{F} = \frac{1}{N} \sum_{i=1}^N \nabla \log p(x_i|\theta) \nabla \log p(x_i|\theta)^T.$$

接下来我们来关注一下Hessian Matrix并寻找它和信息矩阵的关系。

首先关注Hessian matrix的定义

In mathematics, the Hessian matrix or Hessian is a square matrix of second-order partial derivatives of a scalar-valued function, or scalar field. It describes the local curvature of a function of many variables.

还是根据Agustinus的博客，我们首先得到其概率分布对应的Hessian矩阵的表达。

$$\begin{aligned}
 \mathbf{H}_{\log p(x|\theta)} &= \mathbf{J} \left(\frac{\nabla p(x|\theta)}{p(x|\theta)} \right) \\
 &= \frac{\mathbf{H}_{p(x|\theta)} p(x|\theta) - \nabla p(x|\theta) \nabla p(x|\theta)^T}{p(x|\theta) p(x|\theta)} \\
 &= \frac{\mathbf{H}_{p(x|\theta)} p(x|\theta)}{p(x|\theta) p(x|\theta)} - \frac{\nabla p(x|\theta) \nabla p(x|\theta)^T}{p(x|\theta) p(x|\theta)} \\
 &= \frac{\mathbf{H}_{p(x|\theta)}}{p(x|\theta)} - \left(\frac{\nabla p(x|\theta)}{p(x|\theta)} \right) \left(\frac{\nabla p(x|\theta)}{p(x|\theta)} \right)^T,
 \end{aligned}$$

然后，我们计算对应的期望

$$\begin{aligned}
 \mathbb{E}_{p(x|\theta)} [\mathbf{H}_{\log p(x|\theta)}] &= \mathbb{E}_{p(x|\theta)} \left[\frac{\mathbf{H}_{p(x|\theta)}}{p(x|\theta)} - \left(\frac{\nabla p(x|\theta)}{p(x|\theta)} \right) \left(\frac{\nabla p(x|\theta)}{p(x|\theta)} \right)^T \right] \\
 &= \mathbb{E}_{p(x|\theta)} \left[\frac{\mathbf{H}_{p(x|\theta)}}{p(x|\theta)} \right] - \mathbb{E}_{p(x|\theta)} \left[\left(\frac{\nabla p(x|\theta)}{p(x|\theta)} \right) \left(\frac{\nabla p(x|\theta)}{p(x|\theta)} \right)^T \right] \\
 &= \int \frac{\mathbf{H}_{p(x|\theta)}}{p(x|\theta)} p(x|\theta) dx - \mathbb{E}_{p(x|\theta)} [\nabla \log p(x|\theta) \nabla \log p(x|\theta)^T] \\
 &= \mathbf{H}_{\int p(x|\theta) dx} - \mathbf{F} \\
 &= \mathbf{H}_1 - \mathbf{F} \\
 &= -\mathbf{F}.
 \end{aligned}$$

请注意，H对 θ 的二阶导为0。

这样，得证，H的期望等于负的信息矩阵，而当H从期望变成固定的 θ 对应的信息矩阵时，E被去掉，H=-F。

Hessian和Covariance Matrix的关系

这个关系在提供的论文中，已经得到了证明，

Consider a Gaussian random vector θ with mean θ^* and covariance matrix Σ_θ so its joint probability density function (PDF) is given by:

$$p(\theta) = (2\pi)^{-\frac{N_\theta}{2}} |\Sigma_\theta|^{-\frac{1}{2}} \exp \left[-\frac{1}{2} (\theta - \theta^*)^T \Sigma_\theta^{-1} (\theta - \theta^*) \right] \quad (\text{A.1})$$

The objective function can be defined as its negative logarithm:

$$J(\theta) \equiv -\ln p(\theta) = \frac{N_\theta}{2} \ln 2\pi + \frac{1}{2} \ln |\Sigma_\theta| + \frac{1}{2} (\theta - \theta^*)^T \Sigma_\theta^{-1} (\theta - \theta^*) \quad (\text{A.2})$$

which is a quadratic function of the components in θ . By taking partial differentiations with respect to θ_l and $\theta_{l'}$, the (l, l') component of the Hessian matrix can be obtained:

$$\mathcal{H}^{(l,l')}(\theta^*) = \left. \frac{\partial^2 J(\theta)}{\partial \theta_l \partial \theta_{l'}} \right|_{\theta=\theta^*} = (\Sigma_\theta^{-1})^{(l,l')} \quad (\text{A.3})$$

so the Hessian matrix is equal to the inverse of the covariance matrix:

$$\mathcal{H}(\theta^*) = \Sigma_\theta^{-1} \quad (\text{A.4})$$

A.4给出了H和协方差矩阵的关系，从这里可以看出，信息矩阵 $F=-H=-\text{inv}(\text{cov})$