

Fisher Information Matrix

Suppose we have a model parameterized by parameter vector θ that models a distribution $p(x|\theta)$. In frequentist statistics, the way we learn θ is to maximize the likelihood $p(x|\theta)$ wrt. parameter θ . To assess the goodness of our estimate of θ we define a score function:

$$s(\theta) = \nabla_{\theta} \log p(x|\theta),$$

that is, score function is the gradient of log likelihood function. The result about score function below is important building block on our discussion.

Claim: The expected value of score wrt. our model is zero.

Proof. Below, the gradient is wrt. θ .

$$\begin{aligned} \mathbb{E}_{p(x|\theta)} [s(\theta)] &= \mathbb{E}_{p(x|\theta)} [\nabla \log p(x|\theta)] \\ &= \int \nabla \log p(x|\theta) p(x|\theta) \, dx \\ &= \int \frac{\nabla p(x|\theta)}{p(x|\theta)} p(x|\theta) \, dx \\ &= \int \nabla p(x|\theta) \, dx \\ &= \nabla \int p(x|\theta) \, dx \\ &= \nabla 1 \\ &= 0 \end{aligned}$$

□

But how certain are we to our estimate? We can define an uncertainty measure around the expected estimate. That is, we look at the covariance of score of our model. Taking the result from above:

$$\mathbb{E}_{p(x|\theta)} [(s(\theta) - 0)(s(\theta) - 0)^T].$$

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:

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

However, usually our likelihood function is complicated and computing the expectation is intractable. We can approximate the expectation in 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, F is called Empirical Fisher:

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

Fisher and Hessian

One property of F that is not obvious is that it has the interpretation of being the negative expected Hessian of our model's log likelihood.

Claim: The negative expected Hessian of log likelihood is equal to the Fisher Information Matrix F .

Proof. The Hessian of the log likelihood is given by the Jacobian of its gradient:

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

where the second line is a result of applying quotient rule of derivative. Taking expectation wrt. our model, we have:

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

Thus we have $\mathbf{F} = -\mathbb{E}_{p(x|\theta)} [\mathbf{H}_{\log p(x|\theta)}]$. □

Indeed knowing this result, we can see the role of \mathbf{F} as a measure of curvature of the log likelihood function.

Conclusion

Fisher Information Matrix is defined as the covariance of score function. It is a curvature matrix and has interpretation as the negative expected Hessian of log likelihood function. Thus the immediate application of \mathbf{F} is as drop-in replacement of \mathbf{H} in second order optimization methods.

One of the most exciting results of F is that it has connection to KL-divergence. This gives rise to natural gradient method, which we shall discuss further in the next article.

References

1. Martens, James. "New insights and perspectives on the natural gradient method." arXiv preprint arXiv:1412.1193 (2014).
2. Ly, Alexander, et al. "A tutorial on Fisher information." Journal of Mathematical Psychology 80 (2017): 40-55.

← **PREVIOUS POST** (</TECHBLOG/2017/12/23/ANNEALED-IMPORTANCE-SAMPLING/>)

NEXT POST → (</TECHBLOG/2018/03/14/NATURAL-GRADIENT/>)



(</feed.xml>)



(<https://github.com/wiseodd>)

Copyright © Agustinus Kristiadi's Blog 2021