

§ 1. Fisher Information Matrix.

Assume we have a distribution $p(x|\theta)$ parameterized by θ . We consider the gradient of its log likelihood

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

as the score function.

Proposition 1. The expectation of the score function w.r.t. our model is 0.

Proof.

$$\begin{aligned} & \mathbb{E}_{x \sim p(\cdot|\theta)} [\nabla_{\theta} \log p(x|\theta)] \\ &= \int_{\mathcal{X}} p(x|\theta) \cdot \nabla_{\theta} \log p(x|\theta) dx \\ &= \int_{\mathcal{X}} p(x|\theta) \cdot \frac{\nabla_{\theta} p(x|\theta)}{p(x|\theta)} \cdot dx \\ &= \int_{\mathcal{X}} \nabla_{\theta} p(x|\theta) dx \\ &= \nabla_{\theta} \int_{\mathcal{X}} p(x|\theta) dx \\ &= \nabla_{\theta} 1 \\ &= 0. \end{aligned}$$

□

Definition 2. The covariance of the score of our model is defined as the Fisher Information Matrix:

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

Remark. When it is untractable to compute the expectation,

F could be replaced by its empirical version,

$$\hat{F} = \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log p(x_i | \theta) \cdot \nabla_{\theta} \log p(x_i | \theta)^T,$$

which is given by our training data, $X = \{x_1, x_2, \dots, x_N\}$.

Proposition 3. The negative expected Hessian matrix of the log likelihood is the FIM F .

$$\text{Proof. } H \log p(x | \theta) = J(\nabla_{\theta} \log p(x | \theta))$$

$$\begin{aligned} &= J \left(\frac{\nabla_{\theta} p(x | \theta)}{p(x | \theta)} \right) \\ &= \underbrace{(\nabla_{\theta}^2 p(x | \theta)) p(x | \theta)}_{p^2(x | \theta)} - \nabla_{\theta} p(x | \theta) \cdot \nabla_{\theta} p(x | \theta)^T \\ &= \frac{H_{p(x | \theta)}}{p(x | \theta)} - \left(\frac{\nabla_{\theta} p(x | \theta)}{p(x | \theta)} \right) \left(\frac{\nabla_{\theta} p(x | \theta)}{p(x | \theta)} \right)^T. \end{aligned}$$

Taking expectation on both sides leads to

$$E_{x \sim p(\cdot | \theta)} [H \log p(x | \theta)]$$

$$= E_{x \sim p(\cdot | \theta)} \left[\frac{H_{p(x | \theta)}}{p(x | \theta)} - \left(\frac{\nabla_{\theta} p(x | \theta)}{p(x | \theta)} \right) \left(\frac{\nabla_{\theta} p(x | \theta)}{p(x | \theta)} \right)^T \right]$$

$$= \int_{\mathbb{X}} p(x | \theta) \cdot \frac{H_{p(x | \theta)}}{p(x | \theta)} dx - E_{x \sim p(\cdot | \theta)} [\nabla_{\theta} \log p(x | \theta) \cdot \nabla_{\theta} \log p(x | \theta)^T]$$

$$= H \cdot \int_{\mathbb{X}} p(x | \theta) dx - F$$

$$= H \cdot I - F$$

$$= -F.$$

III.

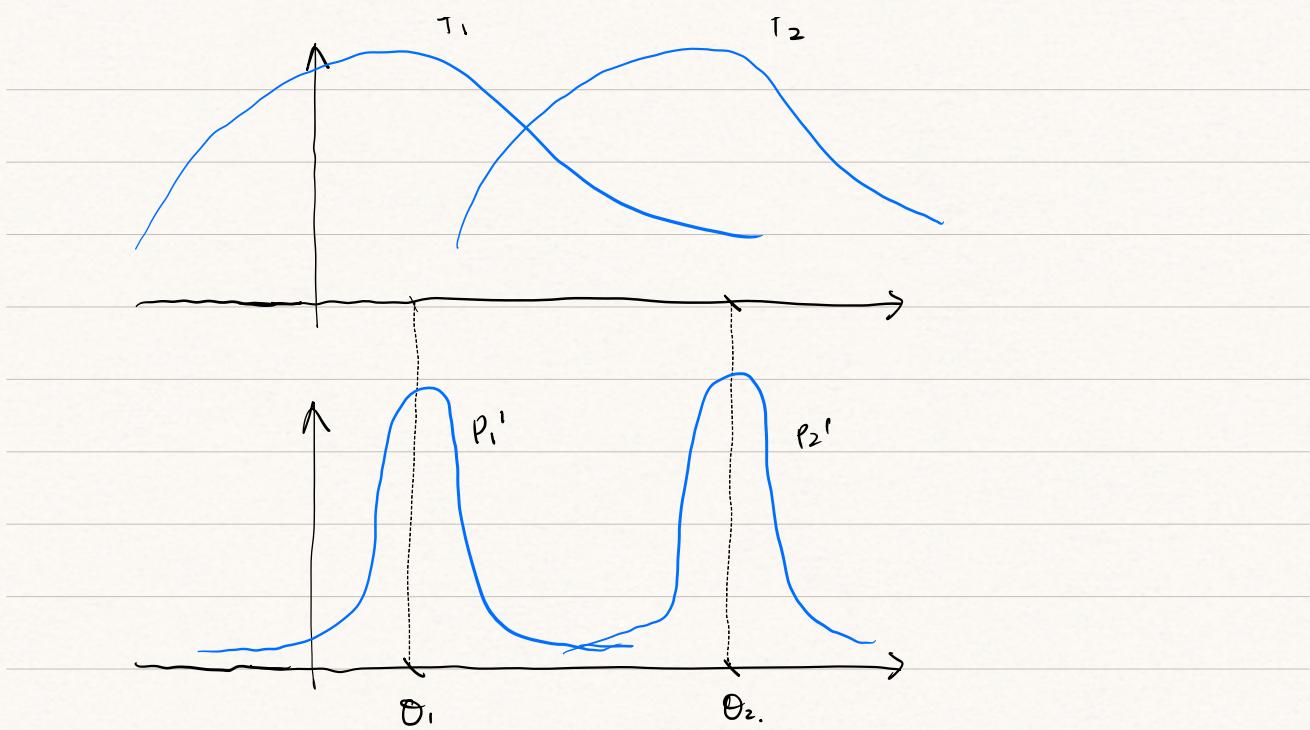
§2. Distribution Space

We would like to maximize the likelihood function to find the most likely parameter θ . Equivalent formulation would be to minimize the loss function $L(\theta)$, which is the negative log likelihood function.

Usual way to solve this optimization is to use gradient descent. In this case, we take a step in the steepest descent around the local neighbourhood of current value of parameter θ , in the parameter space. Formally, we have

$$\frac{-\nabla_{\theta} L(\theta)}{\|\nabla_{\theta} L(\theta)\|} = \lim_{\epsilon \rightarrow 0} \underset{d: \|d\| \leq \epsilon}{\operatorname{argmin}} L(\theta + d).$$

Thus, the optimization in gradient descent depends on the Euclidean geometry of the parameter space. Our purpose is correct, but our approach could be more appropriate, with the fact that two distributions might not be "close" to each other even though their parameters are "close" to each other in the parameter space. Consider the motivating example below:



Though the Euclidean distance between the mean of p_i and p_j and the Euclidean distance between the mean of p'_i and p'_j , are the same, it is clear that p_i and p_j are closer than p'_i and p'_j .

To this end we consider directly optimizing in the distribution space instead of the parameter space. A natural choice is to use KL-divergence as the metric in this distribution space, though KL-divergence is not symmetric and thus not a metric. Actually, it could be shown that as α goes to zero, KL-divergence is asymptotically symmetric.

Proposition 4. FIM F is the Hessian of KL-divergence between two distributions $p(X|\theta)$ and $p(X|\theta')$, w.r.t. θ' , evaluated at $\theta'=\theta$.

Proof. $KL(p(x|\theta), p(x|\theta')) = E_{p(x|\theta)} [\log p(x|\theta)] - E_{p(x|\theta')} [\log p(x|\theta')]$

$$\nabla_{\theta'} KL(p(x|\theta), p(x|\theta')) = - \int_{x \sim p(x|\theta)} p(x|\theta) \cdot \nabla_{\theta'} \log p(x|\theta') \cdot dx$$

$$\nabla_{\theta'}^2 KL(p(x|\theta), p(x|\theta'))|_{\theta'=\theta} = - \int_{x \sim p(x|\theta)} p(x|\theta) \cdot \nabla_{\theta'}^2 \log p(x|\theta')|_{\theta'=\theta} \cdot dx$$

$$= - \int_{x \sim p(x|\theta)} p(x|\theta) \cdot H[\log p(x|\theta)] \cdot dx.$$

$$= F.$$

III

Before we use the FIM to enhance the gradient descent, we need to derive the Taylor series expansion for KL-divergence around θ .

Proposition 5. Let $d \rightarrow 0$. The second order Taylor series expansion of KL-divergence is $KL(p(x|\theta), p(x|\theta+d)) \approx \frac{1}{2} \|d\|_F^2$.

Proof. $KL(p(x|\theta), p(x|\theta+d))$

$$= KL(p(x|\theta), p(x+\theta)) + (\nabla_{\theta'} KL(p(x|\theta), p(x|\theta'))|_{\theta'=\theta})^T d + \frac{d^T (\nabla_{\theta'}^2 KL(p(x|\theta), p(x|\theta'))|_{\theta'=\theta}) d}{2}$$

$$= 0 + 0 + \frac{1}{2} \|d\|_F^2.$$

$$= \frac{1}{2} \|d\|_F^2.$$

III

Now we would like to minimize the loss function $L(\theta)$ in distribution space with the KL-divergence as the metric.

Specifically, we do the minimization

$$d^* = \underset{d}{\operatorname{argmin}} \quad L(\theta + d). \quad (1)$$

$$\text{d} : \text{KL}(p(x|\theta), p(x|\theta+d)) = C$$

where C is a constant. The purpose of fixing the KL-divergence to some constant is to make the model more robust to the reparameterization of the model, i.e., the algorithm does not care how the model is parameterized, it only cares the distribution induced by the parameter.

If we write Eq.(1) in its Lagrangian form, with $L(\theta+d)$ approximated by the first order Taylor series expansion and KL-divergence approximated by the second order Taylor series expansion, we have

$$d^* = \underset{d}{\operatorname{argmin}} \quad L(\theta + d) + \lambda \cdot (\text{KL}(p(x|\theta), p(x|\theta+d)) - C)$$

$$\approx \underset{d}{\operatorname{argmin}} \quad L(\theta) + \nabla_{\theta} L(\theta)^T \cdot d + \frac{\lambda}{2} \|d\|_F^2 - \lambda C. \quad (2)$$

Setting the derivative of Eq. (2) w.r.t. θ to 0 leads to

$$d = -\frac{1}{\lambda} \cdot F^{-1} \cdot \nabla_{\theta} L(\theta),$$

where the constant factor could be absorbed into the learning rate.

Definition. 6. Natural gradient is defined as

$$\tilde{\nabla}_{\theta} L(\theta) = I^{-1} \nabla_{\theta} L(\theta).$$