

Nevertheless, the distributions in these two plots are different, and thus we should find a metric that mirrors this in its definition.

Properties of the KL divergence: Observe that the KL divergence, by definition, satisfies:

$$\begin{aligned} D_{\text{KL}}(p_1(\cdot) || p_2(\cdot)) &= \mathbb{E}_{p_1(\cdot)} \left[\log \frac{p_1(\cdot)}{p_2(\cdot)} \right] \\ &= \mathbb{E}_{p_1(\cdot)} [\log p_1(\cdot)] - \mathbb{E}_{p_1(\cdot)} [\log p_2(\cdot)]. \end{aligned}$$

To help our discussion, and make connections with our problem so far, we use $p_1(\cdot) = p(x|\theta)$ and $p_2(\cdot) = p(x|\theta')$. Then, the gradient of KL with respect to θ' satisfies:

$$\begin{aligned} \nabla_{\theta'} (D_{\text{KL}}(p(x|\theta) || p(x|\theta'))) &= \nabla_{\theta'} \mathbb{E}_{p(x|\theta)} [\log p(x|\theta)] - \nabla_{\theta'} \mathbb{E}_{p(x|\theta)} [\log p(x|\theta')] \\ &= -\mathbb{E}_{p(x|\theta)} [\nabla_{\theta'} \log p(x|\theta')] \\ &= -\int p(x|\theta) \nabla_{\theta'} \log p(x|\theta') dx. \end{aligned}$$

and the second derivative satisfies:

$$\nabla_{\theta'}^2 (D_{\text{KL}}(p(x|\theta) || p(x|\theta'))) = -\int p(x|\theta) \nabla_{\theta'}^2 \log p(x|\theta') dx$$

Then, the Hessian evaluated at $\theta' = \theta$ is:

$$\begin{aligned} H_{D_{\text{KL}}(p(x|\theta) || p(x|\theta'))} &= -\int p(x|\theta) \nabla_{\theta'}^2 \log p(x|\theta')|_{\theta'=\theta} dx \\ &= -\int p(x|\theta) H_{\log p(x|\theta)} dx \\ &= -\mathbb{E}_{p(x|\theta)} [H_{\log p(x|\theta)}] \\ &= F, \end{aligned}$$

which is what we have shown above; i.e., the expected Hessian of the log function is the Fisher information matrix.

2nd-order Taylor expansion of KL divergence and natural gradient: Let us now connect the dots. Following similar reasoning to classical optimization, given an objective function, we can approximate locally the objective with its second-order Taylor approximation (which involves both the gradient and the Hessian information), and then locally minimize that approximation; then, we iterate.

The second-order approximation of the KL divergence metric satisfies:

$$\begin{aligned} D_{\text{KL}}(p(x|\theta) || p(x|\theta + d)) &\approx D_{\text{KL}}(p(x|\theta) || p(x|\theta)) + \langle \nabla D_{\text{KL}}(p(x|\theta) || p(x|\theta)), d \rangle + \frac{1}{2} \langle Fd, d \rangle \\ &= \frac{1}{2} \langle Fd, d \rangle. \end{aligned}$$

Similar to the Euclidean case, we seek for an update vector d that minimizes the loss function $\mathcal{L}(\theta)$ in the *distribution space*. Analogously to steepest descent:

$$d^* = \operatorname{argmin}_{D_{\text{KL}}(p(x|\theta) || p(x|\theta+d))=c} \mathcal{L}(\theta + d),$$

where c is some constant. Compare this with the Euclidean case where:

$$d^* = \operatorname{argmin}_{\|d\|_2 \leq \epsilon} \mathcal{L}(\theta + d).$$

The purpose of fixing the KL-divergence to some constant is to make sure that we move along the space of distributions with constant speed, regardless the curvature.

How do we solve this part? If we write the above minimization in Lagrangian form (*we have not talked about Lagrange*

multipliers yet, but accept this for the moment), we get:

$$\begin{aligned} d^* &= \arg \min_d \{ \mathcal{L}(\theta + d) + \lambda \cdot (D_{\text{KL}}(p(x|\theta) || p(x|\theta + d)) - c) \} \\ &\approx \arg \min_d \left\{ \mathcal{L}(\theta) + \langle \nabla \mathcal{L}(\theta), d \rangle + \frac{1}{2} \lambda \langle Fd, d \rangle - \lambda c \right\}. \end{aligned}$$

To solve this minimization, we set its derivative with respect to d to zero; this will lead to the solution:

$$d = -\frac{1}{\lambda} F^{-1} \nabla \mathcal{L}(\theta).$$

(*The λ constant can be absorbed in the definition of F .*)

The above lead to the definition of the *natural gradient* descent method:

• **Repeat:**

1. Compute the gradient $\nabla \mathcal{L}(\theta_t)$.
2. Compute the Fisher information matrix F_t .
3. Compute the natural gradient direction: $d_t = F_t^{-1} \nabla \mathcal{L}(\theta_t)$.
4. For a step size η , compute $\theta_{t+1} = \theta_t - \eta d_t$.

Take-away messages:

- The natural gradient descent is a generalization of the Newton's method, as there are cases where from the natural gradient descent we can obtain the Newton's iteration.
- Remember that, afterall, the Fisher information matrix is computed per iteration, and inverted; and it turns out that on expectation it corresponds to the expected Hessian of the objective.
- How do we implement the natural gradient method in reality? Remember the definition of the Fisher information:

$$F = -\mathbb{E}_{p(x|\theta)} [H_{\log p(x|\theta)}].$$

In realistic scenarios, we do not have access to the distribution $p(x|\theta)$, but rather have data that come from that distribution. In that case, we refer to the *empirical* Fisher information matrix, defined as:

$$F = \frac{1}{n} \sum_{i=1}^n \nabla \log p(x_i|\theta) \cdot \nabla \log p(x_i|\theta)^\top$$

where $\{x_i\}_{i=1}^n$ denote the training set of examples. In that case, the main recursion of natural gradient descent becomes:

$$\theta_{t+1} = \theta_t - \eta \left(\frac{1}{n} \sum_{i=1}^n \nabla \log p(x_i|\theta_t) \cdot \nabla \log p(x_i|\theta_t)^\top \right)^{-1} \cdot \nabla \hat{\mathcal{L}}(\theta_t),$$

where also the objective and its gradient are evaluated in their empirical form:

$$\nabla \hat{\mathcal{L}}(\theta_t) := \frac{1}{n} \sum_{i=1}^n \nabla \log p(x_i|\theta_t).$$

- The main reason we studied natural gradient descent is to motivate our discussion later on, regarding algorithms in training neural networks. The empirical Fisher information matrix appears in almost all modern algorithms in ML, usually further approximated to be easily computed. E.g., one way to get around computing the exact empirical Fisher information matrix is to constrain it to be diagonal matrix; a technique that is heavily used in algorithms such as AdaGrad, AdaDelta, RMSprop, Adam, AMSGrad, Yogi, etc. In other words, in the mostly used algorithms in neural network training.