# Variational Inference

**Setup**: Given some data $\mathcal{D}$, and a model with parameters $\vartheta \in \mathbb{R}^d$ and a likelihood $p(\mathcal{D}|\vartheta)$, and a prior $p(\vartheta)$.
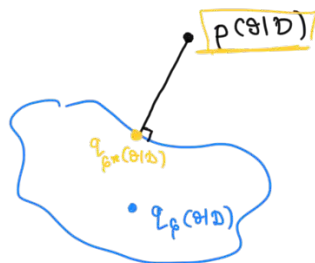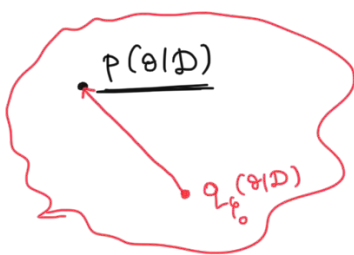
**Goal**: Approximate $\underset{\text{intractable posterior}}{p(\vartheta|\mathcal{D})} = \dfrac{p(\mathcal{D}|\vartheta)\, p(\vartheta)}{p(\mathcal{D}) \longrightarrow \int p(\mathcal{D}|\vartheta)\, p(\vartheta)\, d\vartheta}$

**Main idea**: Approx. $p(\vartheta|\mathcal{D})$ with a <u>family of distributions</u> that is easy to work with.

$$\vartheta = (\vartheta_1, \vartheta_2, \ldots, \vartheta_d)$$

e.g. <u>Mean-field family</u>: $p(\vartheta|\mathcal{D}) \approx q_\varphi(\vartheta|\mathcal{D}) = \prod_{i=1}^{d} \mathcal{N}(\vartheta_i \mid \mu_i, \sigma_i^2)$

$$\varphi := \{ \mu_1, \sigma_1^2, \mu_2, \sigma_2^2, \ldots, \mu_d, \sigma_d^2 \}$$

variational parameters

**Goal**: Find/estimate the variational parameters $\varphi$ such that $q_\varphi(\vartheta|\mathcal{D})$ is as close as possible to $p(\vartheta|\mathcal{D})$.



How to compute the "distance" between $q_\varphi(\vartheta|\mathcal{D})$ and $p(\vartheta|\mathcal{D})$?

In practice we use the Kullback-Leibler divergence as a way to compare $q_\varphi(\vartheta|\mathcal{D})$ to $p(\vartheta|\mathcal{D})$:

$$\left[ \int_\vartheta \frac{q_\varphi(\vartheta|\mathcal{D})}{} \right]

$$\text{KL}\left[q_\phi(\theta|D) \| p(\theta|D)\right] := \int \log \frac{1}{p(\theta|D)} \, q_\phi(\theta|D) \, d\theta$$

$$\underbrace{\phantom{\text{KL}\left[q_\phi(\theta|D) \| p(\theta|D)\right]}}_{\text{relative entropy}}$$

$$= \mathbb{E}_{\theta \sim q_\phi(\theta|D)} \left[ \log \frac{q_\phi(\theta|D)}{p(\theta|D)} \right]$$

Why use the KL-divergence?

... It is easy to work with (see below).

... but, it is **not** a distance, i.e. $\text{KL}[q_\phi \| p] \neq \underbrace{\text{KL}[p \| q_\phi]}_{\text{reverse KL}}$

however, $\text{KL}[q_\phi \| p] = 0 \iff q_\phi = p$.

We want want to find $\boxed{\phi^* = \arg\min_\phi \text{KL}\left[q_\phi(\theta|D) \| p(\theta|D)\right]}$

How to estimate the "optimal" variational parameters $\phi^*$?

1. Old schoolers had to derive coordinate ascent rules for minimizing the KL. (see chapter 10 in Bishop).

2. New schoolers are using Automatic Differentiation Variational Infer..

ADVI: It is a black-box approach is agnostic to any details about $p(\theta$

Any model for which we can **evaluate** (and differentiate) the

**log-likelihood** and the **log-prior** works!

Recall : $\text{KL}\left[q_\phi(\theta|D) \| p(\theta\|D)\right] = \mathbb{E}_{\theta \sim q_\phi(\theta|D)}\left[\log q_\phi(\theta|D) - \log p(\theta|D)\right]$

$\underline{1^{st} \text{ term}}$ : $\mathbb{E}_{\theta \sim q_\phi(\theta|D)}\left[\log q_\phi(\theta|D)\right] = \int \log q_\phi(\theta|D) \cdot q_\phi(\theta|D) \, d\theta = $

$$\underbrace{\phantom{\int \log q_\phi(\theta|D) \cdot q_\phi(\theta|D) \, d\theta}}_{-H[q_\phi(\theta|D)] : \text{ negative entropy of } q_\phi(\theta|D)}$$

e.g. mean-field:

$$\mathbb{E}_{\vartheta \sim q_\phi(\vartheta \mid D)}\left[\log q_\phi(\vartheta \mid D)\right] = -\sum_{i=1}^{d} \log \sigma_i + \text{constant}$$

$\underline{2^{nd} \text{ term}}$ :

$$\mathbb{E}_{\vartheta \sim q_\phi(\vartheta \mid D)}\left[\log p(\vartheta \mid D)\right] =$$

$$\overset{\text{Bayes}}{=} \mathbb{E}_{\vartheta \sim q_\phi(\vartheta \mid D)}\left[\log p(D \mid \vartheta) + \underbrace{\log p(\vartheta)} - \underbrace{\log p(D)}^{\text{constant}}\right]$$

$$\Longrightarrow \mathbb{KL}\left[q_\phi(\vartheta \mid D) \,\|\, p(\vartheta \mid D)\right] = \overbrace{- H\left[q_\phi(\vartheta \mid D)\right]}^{\text{analytically}}$$

$$- \mathbb{E}_{\vartheta \sim q_\phi(\vartheta \mid D)}\left[\log p(D \mid \vartheta) + \log p(\vartheta)\right] +$$

Now, all terms can be evaluated.

But here's the catch :

$$\left\{ \begin{array}{c} \phi^* = \arg\min_\phi \mathcal{L}(\phi) \\[1em] \mathcal{L}(\phi) := - H\left[q_\phi(\vartheta \mid D)\right] - \mathbb{E}_{\vartheta \sim q_\phi(\vartheta \mid D)}\left[\log p(D \mid \vartheta) + \log p(\vartheta)\right] \end{array} \right\}$$

Minimization via gradient descent :

$$\phi_{n+1} = \phi_n - \eta \underline{\nabla_\phi \mathcal{L}(\phi)}$$

<u>Remarks</u> :

1.) All terms in $\mathcal{L}(\phi)$ can be evaluated. Sometimes this can
be done analytically (e.g. linear models with Gaussian likelihood and prior).

2.) We need to compute $\nabla_\phi \mathcal{L}(\phi)$:

$$\nabla_\phi \underbrace{\mathbb{E}_{\theta \sim q_\phi(\theta|D)} \left[ \log p(D|\theta) + \log p(\theta) \right]}_{=} = \nabla_\phi \int \left[ \log p(D|\theta) + \log p(\theta) \right] q_\phi(\theta|D)$$

We could sample $\theta_i \sim q_\phi(\theta|D)$ and use a Monte-Carlo estimator to compute the gradient:

$$\nabla_\phi \mathbb{E}_{\theta \sim q_\phi(\theta|D)} \left[ \log p(D|\theta) + \log p(\theta) \right] \approx \frac{1}{n} \sum_{i=1}^{n} \left[ \nabla_\phi \log p(D|\theta_i) + \nabla_\phi \log p(\theta_i) \right]$$

, where $\theta_i \overset{i.i.d}{\sim} q_\phi(\theta_i|D)$

However, notice that $q_\phi(\theta_i|D)$ depends on $\phi$, and as $\phi$ is changing during optimization, this M.C. estimator will exhibit very high variance. I.e. we will need a very large number of MC samples to get a reasonable approximat of the gradient.

## Reparametrization trick (next time):

Introduce a simple "change of variables" such that the required expectation can be computed with respect to distributions that do not depend on $\phi$.