# Homework on Variational Inference

This document provides detailed explanations for the models implemented in

- `logistic-regression/model_logreg_mvn.py` a Bayesian logistic regression / classification model for linearly separable data
- `vae/model_vae.py` an unsupervised model for handwritten digit generation

The implementations use the `pytorch-lightning` framework and the relevant parts of code are

- `logistic-regression/model_logreg_mvn.py`, functions `forward`, `accuracy` and `loss` as well as `logistic-regression/run.py`
- `vae/model_vae.py` functions `forward` and `loss` as well as `vae/run.py` and `vae/config_vae.yaml`

## Bayesian logistic regression / classification

The Bayesian logistic regression model is a classification model defined as

$$p(Y, w|X) = p(w) \prod_{i=1}^{n} p(y_i|x_i, w) \tag{1}$$

$$= \mathcal{N}(w; 0, \sigma_{\text{prior}} I_d) \prod_{i=1}^{n} \text{Bernoulli}(y_i; x_i w) \tag{2}$$

where

- $x_i \in \mathbb{R}^d$ are features $X = [x_1, x_2, \ldots, x_n]^T$
- $y_i \in \{0, 1\}$ are binary class labels $Y = [y_1, y_2, \ldots, y_n]^T$.
- $\mathcal{N}(w; 0, \sigma_{\text{prior}} I_d)$ denotes a multivariate normal distribution

Bayesian inference in this model means

- computing/approximating the posterior

$$q(w; Y, X) \approx p(w|Y, X) = \frac{p(Y, w|X)}{p(Y|X)} \tag{3}$$

- computing/approximating label predictions for new input features

$$p(y_*|x_*) = \int dw\, p(w|Y, X)\, p(y_*|x_*, w) \tag{4}$$

$$= \int dw\, q(w)\, p(y_*|x_*, w) \tag{5}$$

## Variational inference

Since the Bayesian posterior $p(w|Y, X)$ is analytically intractable we approximate it either by a multivariate Normal $q_\phi(w) = \mathcal{N}(w; \mu, \Sigma)$ with $\phi = \{\mu, \Sigma\}$ or diagonal/factorised multivariate Normal $q_\phi(w) = \mathcal{N}(w; \mu, \mathrm{diag}(\sigma^2))$ $\phi = \{\mu, \sigma^2\}$ by optimising the negative evidence lower bound (ELBO)

$$L(\phi; Y, X) = -\sum_{i=1}^{n} \mathbb{E}_{q_\phi(w)}[\log p(y_i|x_i, w)] + \mathrm{KL}[q_\phi(w)\|p(w)]. \tag{6}$$

$$\geq \log p(Y|X) \tag{7}$$

Here we need to numerically approximate $\mathbb{E}_{q_\phi(w)}[\log p(y_i|x_i w)]$ and the predictive distribution We can use

$$\int dw\, f(x^T w)\, \mathcal{N}(u; \mu, \Sigma) = \int du\, f(u)\, \mathcal{N}(u; x^T w, x^T \Sigma x) \tag{8}$$

$$= \int d\epsilon\, f(x^T w + \sqrt{x^T \Sigma x}\epsilon)\, \mathcal{N}(\epsilon; 0, 1) \tag{9}$$

$$\int du\, \sigma(u)\, \mathcal{N}(u; \mu, \sigma^2) \approx \sigma\left(\frac{1}{\sqrt{1 + \pi\sigma^2/8}}\mu\right) \tag{10}$$

$$\int du\, \sigma(u)\, \mathcal{N}(u; \mu, \sigma^2) \approx \sum_{k=1}^{K} \sigma(\mu + \sqrt{2}\sigma\tilde{u}_i)\frac{1}{\sqrt{\pi}}w_i, \tag{11}$$

where $\{\tilde{u}_k, w_k\}_{k=1}^{K}$ are the weights and nodes of the univariate Gauss-Hermite quadrature.

## Key variational inference concepts to learn from this model

The approximation of the EBLO objective in implemented in in the function.

```
1  def loss(self, features, labels):
```

please follow the implementation there.

## Batch learning

For large datasets we cannot use all data in each training step due to memory restrictions, therefore, we use the approximation

$$\sum_{i=1}^{n} \mathbb{E}_{q_\phi(w)}[\log p(y_i|x_i, w)] \approx n \frac{1}{|S|} \sum_{s \in S} \mathbb{E}_{q_\phi(w)}[\log p(y_i|x_s, w)], \qquad (12)$$

that is, we approximate the objective by using only a random subset $S \subseteq \{1, \ldots, N\}$ to represent the dataset. This makes the objective stochastic w.r.t. sampling $S$ but with the right optimisation procedure (stochastic gradient descent with decaying learning rates) convergence can still be achieved.

**Implementation** This concept is implemented via the `DataModuleFromNPZ` in `run.py` which uses data batches of size `size_batch`, see

```
1  dm = DataModuleFromNPZ(
2      data_dir="data_logistic_regression_2d",
3      feature_labels=["inputs", "targets"],
4      batch_size=64,
5      num_workers=4,
6      shuffle_training=False
7  )
```

and the code line

```
1  logp_expct =
   self.size_data*torch.mean(p_labels.log_prob(labels.repeat((1,self.n_sample
   s_mc))))
```

where we take the batch average of the expected log likelihood and rescale it with the size of the training data set.

## Stochastic gradient learning

The expectations $\mathbb{E}_{q_\phi(w)}[\log p(y_i|x_s, w)]$ can rarely be computed exactly or approximated accurately via quadrature methods---quadrature work well only in 1D. For this reason we often use Monte-Carlo estimates

$$\mathbb{E}_{q_\phi(w)}[\log p(y_i|x_s, w)] \approx \frac{1}{R} \sum_{w^{(r)} \sim q(w)} \log p(y_i|x_s, w^{(r)}), \qquad (13)$$

that is, we sample $R$ samples $w^{(r)} \sim q(w)$ and average. This again makes the objective stochastic but we hope that with the right number of samples and the right optimisation procedure the optimisation can stlil converge (see above).

**Implementation** This is implemented in line

```
1  logp_expct = self.size_data*torch.mean(p_labels.log_prob(labels.repeat((1,
   self.n_samples_mc)))))
```

In our implementation we so something a cleverer detailed in the next two sections.

## Reparameterisation of stochastic variables

If a random variable can be represented as a deterministic (preferably differentiable) function of some other/basic random variable with fixed or no parameters, say,

$$z = f_\theta(\epsilon), \quad \epsilon \sim p_0(\epsilon), \qquad p_\theta(z) = \int d\epsilon \, p_0(\epsilon) \, \delta(z - f_\theta(\epsilon)) \qquad (14)$$

then we can rewrite expectations w.r.t this base distribution and make the source of stochasticity in $p_\theta$ independent of the parameters

$$E_{p(z)}[g(z)] = E_{p_0(\epsilon)}[g(f_\theta(\epsilon))] \approx \frac{1}{R} \sum_{\epsilon^{(r)} \sim p_0(\epsilon)} g(f_\theta(\epsilon^{(r)})). \qquad (15)$$

This makes the MC estimates of expectations easily computable and differentiable w.r.t. $\theta$, that is

$$\partial_\theta E_{p(z)}[g(z)] = E_{p_0(\epsilon)}[\partial g(f_\theta(\epsilon))\partial_\theta f_\theta(\epsilon)] \approx \frac{1}{R} \sum_{\epsilon_r \sim p_0(\epsilon)} \partial g(f_\theta(\epsilon))\partial_\theta f_\theta(\epsilon). \qquad (16)$$

In our model we have $w = \mu + L\epsilon, LL^T = \Sigma, \epsilon \sim \mathcal{N}(0, I)$. Hence we can use Monte-Carlo samples from $\epsilon$ to approximate the objective and easily differentiate the approximation.

**Implementation** This is implemented in lines

```
1  # reparameterisation of stochastic variables
2  L      = self.weights_chol()
3  p_post = MultivariateNormal(loc=self.weights_loc.squeeze(), scale_tril=L)
```

via the helper function

```
1  def weights_chol(self):
2      return torch.tril(self.weights_scale_lower, -1) +
   torch.diag(torch.exp(self.weights_scale_logdiag))
```

and the parameterisation is defined the `__init__` function in lines

```
1  self.weights_loc          = nn.Parameter(torch.zeros((self.dim,1)),
   requires_grad=True)
2  self.weights_scale_logdiag = nn.Parameter(torch.zeros((self.dim)),
   requires_grad=True)
3  self.weights_scale_lower   = nn.Parameter(torch.zeros((self.dim,
   self.dim)), requires_grad=True)
```

## Local reparameterisation

We observe that the the likelihood terms $\mathrm{Bernoulli}(y_i; x_i w)$ depend only on $x_i w$ hence instead to sampling from $w \sim q(w) = \mathcal{N}(\mu, \Sigma)$ we can sample from $z_i = x_i w \sim= \mathcal{N}(x_i \mu, x_i \Sigma x_i^T)$, that is

$$\mathbb{E}_{q_\phi(w)}[\log \mathrm{Bernoulli}(y_i; x_i w)] = \mathbb{E}_{w \sim \mathcal{N}(\mu, \Sigma)}[\log \mathrm{Bernoulli}(y_i; x_i w)] \qquad (17)$$

$$= \mathbb{E}_{z_i \sim \mathcal{N}(x_i \mu, x_i \Sigma x_i^T)}[\log \mathrm{Bernoulli}(y_i; z_i)] \qquad (18)$$

$$\approx \frac{1}{R} \sum_{z_i^{(r)} \sim \mathcal{N}(x_i \mu, x_i \Sigma x_i^T)} \log \mathrm{Bernoulli}(y_i; z_i^{(r)}) \qquad (19)$$

thus significantly reducing the variance of the stochastic approximation of the objective.

**Implementation** This is implemented in lines

```
1  # local reparameterisation and MCsampling
2  z_loc     = torch.matmul(features, self.weights_loc).squeeze()
3  z_scale   = torch.sqrt(torch.sum(torch.matmul(features, L)**2, dim=-1,
   keepdim=True)).squeeze()
4  z_samples = Normal(loc=z_loc,
   scale=z_scale).rsample([self.n_samples_mc]).transpose(0,1)
5
6  # data distribution via MC samples
7  p_labels  = Bernoulli(logits=z_samples)
8  # computing the MC samples based expected log likelihood with batch
   learning correction
9  logp_expct = self.size_data*torch.mean(p_labels.log_prob(labels.repeat((1,
   self.n_samples_mc))))
```

## Questions and tasks (at home)

For the logistic regression model detailed above

- run the code with `python run.py` and check metrics with `tensorboard --logdir lightning_logs`
- change `batch_size`, `n_samples_mc`, `max_epochs`, what do you notice?
- try to implement the diagonal version of `class ModelLogisicRegressionMvn(LightningModule)`, what changes do you have to make?
- compare for $q_\phi(w; Y, X) = \mathcal{N}(w; \mu, \Sigma)$ and $q_\phi(w; Y, X) = \mathcal{N}(w; \mu, \mathrm{diag}(\sigma^2))$
  - what are the differences in terms of storage and computational complexity
  - compare the predictive results on a test set,
  - plot data and predictive class probability
  - plot the two distributions as functions of $w$ and compare to $p(w|Y, X)$, what can we learn?

# Unsupervised handwritten digit recognition /variational auto-encoder

Variational auto-encoders are unsupervised models that learn to embed and generate new data similar to one in a, i.i.d. dataset $X = [x_1, \ldots, x_n]^T$. They are Bayesian models where the distribution of the data is

$$p_\theta(X) = \prod_{i=1}^{n} p_\theta(x_i|z_i) \, p(z_i). \tag{20}$$

where generally we have

$$p_\theta(x_i|z_i) = \mathcal{N}(x_i; \mathrm{NN}^{\mathrm{dec}}_{\theta_\mu}(z_i), \theta_{\sigma^2} I_d) \quad \text{and} \quad p(z_i) = \mathcal{N}(0, I_d). \tag{21}$$

Training is done via maximum likelihood using variational Bayes with a posterior approximation

$$q_\phi(z_i; x) = \mathcal{N}(z_i; \mathrm{NN}^{\mathrm{enc-mean}}_{\phi_\mu}(x_i), \mathrm{diag}(\mathrm{NN}^{\mathrm{enc-var}}_{\theta_\sigma}(x_i))) \tag{22}$$

.Here $\mathrm{NN}$ denotes a neural network.

The function to optimise is the negative evidence lower bound

$$L(\theta, \phi; X) = -\sum_i E_{q_\phi(z_i; x_i)}[\log p_\theta(x_i|x_i)] + \mathrm{KL}[q_\phi(z_i; x_i)\|p(z_i)]. \qquad (23)$$

**Implementation** The parameters of the distributions are implemented via a pair of neural networks. The decoder implements $p_\theta(x_i|x_i)$ while the encoder implements $q_\phi(z_i; x_i)$ .

```
1  class Encoder(nn.Module):
2      def __init__(self, config):
3          super().__init__()
4          ...
5          self.model  = nn.Sequential(...)
6      def forward(self, data):
7          input  = data
8          output = self.model(input)
9          loc, scale_isp = torch.split(output, [self.d_state,
   self.d_state], dim=-1)
10          return loc, STD_MIN + torch.nn.functional.softplus(scale_isp)
11
12   class Decoder(nn.Module):
13      def __init__(self, config):
14          super().__init__()
15          self.config = config
16          ...
17          self.loc  = nn.Sequential(...)
18          self.scale_isp = nn.Parameter(0.5*torch.ones(1),
   requires_grad=True)
19
20      def forward(self, input:
21          loc   = self.loc(input)
22          scale = STD_MIN + torch.nn.functional.softplus(self.scale_isp) *
   torch.ones(loc.shape, device=loc.device)
23          return loc, scale
```

# Key variational inference concepts to learn from this model

The loss function is implemented in

```
1  def loss(self, imgs):
```

## Maximum likelihood via expectation maximisation

In this model we not only have to approximate $p_\theta(z_i|x_i)$ but we also need to maximise the lower bound with respext to the model parameters $\theta$., that is we jointly maximise the w.r.t. the model parameters and the posterior approximation $q_\phi(z_i; x_i)$.

**Implementation** joint learning is implemented via a single optimiser

```
1  def configure_optimizers(self):
2      opt = torch.optim.Adam(itertools.chain(self.encoder.parameters(),
3                                             self.decoder.parameters()),
4                             lr=self.hparams.lr, betas=(self.hparams.b1,
   self.hparams.b2))
5      return opt
```

## Amortised variational inference

If we would proceed according to the logistic regression model we would have to approximate each $p_\theta(z_i|x_i)$ in a separate inner loop for each new $\theta$ value. Instead we learn $q_\phi(z_i; x_i) \approx p_\theta(z_i|x_i)$ thus replacing the variational inference algorithm with learning the parameter mappings $\text{NN}^{\text{enc-mean}}_{\phi_\mu}(x_i)$ and $\text{NN}^{\text{enc-var}}_{\theta_\mu}(x_i)$.

```
1  q_zx =
   torch.distributions.normal.Normal(*self.encoder(imgs.view([size_batch] +
   self.encoder.shape_input)))
2  z_samples = q_zx.rsample((self.hparams.n_samples,))
3  p_xz =
   torch.distributions.normal.Normal(*self.decoder(z_samples.view([self.hpara
   ms.n_samples, size_batch] + self.decoder.shape_input)))
4  lik_eval = p_xz.log_prob(imgs.view([size_batch] +
   self.encoder.shape_input))
5  logp = torch.sum(torch.mean(torch.mean(lik_eval, dim=0, keepdim=True),
   dim=1, keepdim=True))
```

## Optional questions and tasks

- run the code with `python run.py --config config_vae.yaml` and check metrics with `tensorboard --logdir lightning_logs`

- if you have time, read the original paper https://arxiv.org/pdf/1312.6114.pdf

- try to answer the questions

- how is this model different from the logistic regression one
    - in terms of latent variables?
    - in terms of likelihood model?
- what are the differences in term of approximate inference
    - in term of latent variables?
    - in terms of parameterisation, what does amortisation mean?
    - can we use amortisation for the logistic regression model?