

Homework on Variational Inference

This document provides detailed explanations about the models implemented in

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

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

- `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

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

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

where

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

Bayesian inference in this model means

- computing/approximating the posterior

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

- computing/approximating label predictions for new input features

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

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

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 a diagonal/factorised multivariate normal $q_\phi(w) = \mathcal{N}(w; \mu, \text{diag}(\sigma^2))$ $\phi = \{\mu, \sigma^2\}$ by optimising the negative evidence lower bound (ELBO)

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

$$\geq -\log p(Y|X) \quad (7)$$

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 the following tricks

$$\int dw f(x^T w) \mathcal{N}(w; \mu, \Sigma) = \int dz f(z) \mathcal{N}(z; x^T \mu, x^T \Sigma x) \quad (8)$$

$$\int dz \text{sigmoid}(z) \mathcal{N}(z; \mu, \sigma^2) \approx \text{sigmoid} \left(\frac{1}{\sqrt{1 + \pi \sigma^2 / 8}} \mu \right) \quad (9)$$

$$\int dz f(z) \mathcal{N}(z; \mu, \sigma^2) \approx \sum_{k=1}^K f(\mu + \sqrt{2} \sigma \tilde{z}_i) \frac{1}{\sqrt{\pi}} w_i, \quad (10)$$

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

`class LikelihoodLogistic(nn.Module)` for implementation. Alternatively, we can also use Monte-Carlo sampling as explained below.

Key variational inference concepts in this model/implementation

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

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

please follow the implementation there.

Batch learning

For large datasets it is often infeasible to use all training data in each training step, for example, due to memory restrictions. Therefore, we use the approximation

$$\sum_{i=1}^n \mathbb{E}_{q_\phi(w)} [\log p(y_i | x_i, w)] \approx n \frac{1}{|S|} \sum_{s \in S} \mathbb{E}_{q_\phi(w)} [\log p(y_i | x_s, w)], \quad (11)$$

that is, we approximate the objective by using only a random subset $S \subseteq \{1, \dots, N\}$ to represent the data. 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 or guaranteed (Robbins-Monro).

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 )
```

in `run.py` and the code line

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

in function `loss`. 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 works well only in 1D. For this reason we often use Monte-Carlo estimates

$$\mathbb{E}_{q_\phi(w)}[\log p(y_i|x_s, w)] \approx \frac{1}{R} \sum_{w^{(r)} \sim q(w)} \log p(y_i|x_s, w^{(r)}), \quad (12)$$

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

Implementation This is implemented in

```
1 # data distribution via MC samples  
2 p_labels = Bernoulli(logits=z_samples)  
3 # computing the MC samples based expected log likelihood with batch learning  
  correction  
4 logp_expct = self.size_data*torch.mean(p_labels.log_prob(labels.repeat((1,  
    self.n_samples_mc))))
```

As you can see, in our implementation we do something cleverer by using MC samples from $z = xw$, we explain this 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,

$$z = f_{\theta}(\epsilon), \quad \epsilon \sim p_0(\epsilon), \quad p_{\theta}(z) = \int d\epsilon p_0(\epsilon) \delta(z - f_{\theta}(\epsilon)) \quad (13)$$

then we can rewrite expectations w.r.t $p_{\theta}(s)$ as expectation w.r.t. the base distribution $p_0(\epsilon)$. For example, for a function $g()$ we can write and approximate the expectations as

$$\mathbb{E}_{p_{\theta}(z)}[g(z)] = \mathbb{E}_{p_0(\epsilon)}[g(f_{\theta}(\epsilon))] \approx \frac{1}{R} \sum_{\epsilon^{(r)} \sim p_0(\epsilon)} g(f_{\theta}(\epsilon^{(r)})). \quad (14)$$

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

$$\partial_{\theta} \mathbb{E}_{p(z)}[g(z)] = \mathbb{E}_{p_0(\epsilon)}[\partial g(f_{\theta}(\epsilon)) \partial_{\theta} f_{\theta}(\epsilon)] \approx \frac{1}{R} \sum_{\epsilon_r \sim p_0(\epsilon)} \partial g(f_{\theta}(\epsilon)) \partial_{\theta} f_{\theta}(\epsilon). \quad (15)$$

In our model we have $w = \mu + L\epsilon$, $LL^T = \Sigma$, $\epsilon \sim \mathcal{N}(0, I)$, therefore, we can use Monte-Carlo samples from ϵ to approximate the objective and easily differentiate the approximation. Since we need a positive definite Σ which requires an L with a positive diagonal, we use a parameterisation $L = L_{\text{lower}} + \text{diag}(e^{\gamma})$, where L_{lower} is a strictly lower triangular matrix and $\gamma \in \mathbb{R}^d$.

Implementation This is implemented in

```
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) +
3     torch.diag(torch.exp(self.weights_scale_logdiag))
```

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

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

Local reparameterisation

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

$$\mathbb{E}_{q_\phi(w)}[\log \text{Bernoulli}(y_i; x_i w)] = \mathbb{E}_{w \sim \mathcal{N}(\mu, \Sigma)}[\log \text{Bernoulli}(y_i; x_i w)] \quad (16)$$

$$= \mathbb{E}_{z_i \sim \mathcal{N}(x_i^T \mu, x_i^T \Sigma x_i)}[\log \text{Bernoulli}(y_i; z_i)] \quad (17)$$

$$\approx \frac{1}{R} \sum_{z_i^{(r)} \sim \mathcal{N}(x_i \mu, x_i \Sigma x_i^T)} \log \text{Bernoulli}(y_i; z_i^{(r)}) \quad (18)$$

$$= \frac{1}{R} \sum_{\epsilon_i^{(r)} \sim \mathcal{N}(0, 1)} \log \text{Bernoulli}(y_i; x_i^T \mu + \sqrt{x_i^T \Sigma x_i} \epsilon_i^{(r)}) \quad (19)$$

thus significantly reducing the variance of the stochastic approximation of the objective, the latter step is implemented in `pytorch` in `Normal(loc=z_loc, scale=z_scale).rsample([self.n_samples_mc])`.

Implementation Approximating the expectations and using all the approximation steps detailed above we obtain

```
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,
4                                   keepdim=True)).squeeze()
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
9 logp_expct = self.size_data*torch.mean(p_labels.log_prob(labels.repeat((1,
10                                   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`, `lr` and the optimisation algorithm, what do you notice?
- Try to implement the diagonal version of `class ModelLogisticRegressionMvn(LightningModule)`. What changes do you have to make? What advantages does this method have when the number of features d is large? Hint: you need to compute `z_scale` in a different way and use `Normal` instead of `MultivariateNormal` to compute the KL divergence.
- How would you change the model to do classification on linearly non-separable data?

- Compare for $q_\phi(w; Y, X) = \mathcal{N}(w; \mu, \Sigma)$ and $q_\phi(w; Y, X) = \mathcal{N}(w; \mu, \text{diag}(\sigma^2))$.
 - 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 an i.i.d. dataset $X = [x_1, \dots, x_n]^T$. They are Bayesian models where the distribution of the data is defined as

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

and learning is done by maximising the likelihood $p_\theta(X)$ w.r.t. the model parameters θ . After successful training the following properties should hold

- samples $z \sim p(z)$, $x \sim p(x|z)$ generate data similar to those in X
- the Bayesian posteriors $p(z|x)$ embed the data in a space where components of z or the new dataset $z \sim p(z_i|x_i)$ can tell us more about the data, say, the embedding cluster nicely, relevant variations in features align along the axes, etc

Where generally choose

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

Here NN denotes a neural network and θ_μ denote the weights.

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

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

Please check the lecture slides/video for more details.

The function to optimise the negative ELBO for this model is

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

which now also has to be optimised w.r.t. the model parameters θ . Note that we did not do this in the logistic regression model because we had a very limited number of model parameters.

Implementation The parameters of the distributions are implemented via a pair of neural networks. The decoder implements $p_{\theta}(x_i|z_i)$ 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) *
23         torch.ones(loc.shape, device=loc.device)
24         return loc, scale
```

Key variational inference concepts to learn from this model

The loss function is implemented in

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

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 respect to the model parameters θ , that is we jointly maximise the ELBO w.r.t. the model parameters θ and the posterior approximation's parameters ϕ .

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,
5 self.hparams.b2))
6     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 θ value. Instead we learn $q_{\phi}(z_i;x_i) \approx p_{\theta}(z_i|x_i)$ thus replacing the variational inference algorithm to learn $q_{\phi}(z_i;x_i)$ with learning the parameter mappings $\text{NN}_{\phi_{\mu}}^{\text{enc-mean}}(x_i)$ and $\text{NN}_{\theta_{\mu}}^{\text{enc-var}}(x_i)$.

```
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.hparams.n_sampl
    es, 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?