

# Machine Learning - Supplement 1

Dom Hutchinson

January 11, 2020

## 1 Variational Bayes

*Variational Bayes* is a method for approximating an intractable integral,  $P$ , with a tractable one,  $Q$ . This is useful as *Posteriors* are integrals and generally intractable.

$$P(X|Y) \approx Q(X)$$

From the derivation below we get that

$$\ln P(Y) = KL(Q||P) - \mathbb{E}_X[\ln Q(X) - \ln P(X, Y)]$$

since  $\ln P(Y)$  is fixed wrt  $Q$  we can derive a likelihood function

$$\mathcal{L}(Q) = -\mathbb{E}_X[\ln Q(X) - \ln P(X, Y)]$$

If we maximise this likelihood function,  $\mathcal{L}$ , then we are minimising  $KL(Q||P)$  which means  $Q$  &  $P$  are becoming similar.

We have reduced the problem of approximation to just optimising  $\mathcal{L}(Q)$ .

By choosing a good form for  $Q$   $\mathcal{L}(Q)$  becomes tractable.

### 1.1 Kullback-Leibler Divergence

*Kullback-Leibler Divergence* is a similarity measure for two distributions,  $Q$  &  $P$ .

$$KL(Q||P) := \int Q(X) \ln \left( \frac{Q(X)}{P(X|Y)} \right) dX$$

The lower the value of  $KL(Q||P)$  the similar  $Q$  &  $P$  are.

*N.B.*  $KL(\cdot||\cdot) \geq 0$  and  $KL(Q||Q) = 0$ .

### 1.2 Intractibility

Here I show why the *Posterior* is intractable

$$\underbrace{P(X|Y)}_{\text{Intractable}} = \frac{\overbrace{P(Y|X)P(X)}^{\text{Tractable}}}{\underbrace{P(Y)}_{\text{Intractable}}} = \frac{P(Y|X)P(X)}{\underbrace{\int P(X, Y)dX}_{\text{Intractable}}}$$

$\int P(X, Y)dX$  is intractable since the space  $X$  is intractably large.

This makes the evidence, and thus posterior, intractable.

### 1.3 Derivation

$$\begin{aligned}
 KL(Q||P) &= \sum_X Q(X) \ln \left[ \frac{Q(X)}{P(X|Y)} \right] \\
 &= \sum_X Q(X) \ln \left[ \frac{Q(X)P(X)}{P(X,Y)} \right] \text{ by product rule} \\
 &= \sum_X Q(X) [\ln Q(X) + \ln P(X) - \ln P(X,Y)] \text{ by log rules} \\
 &= \mathbb{E}_X [\ln Q(X) - \ln P(X,Y)] + \ln P(Y) \text{ since } P(Y) \text{ is independent of } X \\
 \implies \ln P(Y) &= KL(Q||P) - \underbrace{\mathbb{E}_X [\ln Q(X) - \ln P(X,Y)]}_{\mathcal{L}(Q)}
 \end{aligned}$$

## 2 Predictive Gaussian Processes

*Gaussian Processes* are the class of *Stochastic Processes* st every finite linear combination of random variables is normally distributed.

Here I describe the process of make predictions for the value at a set of points, given training data  $(X, y)$ .

- i) Observed data points,  $(X, y)$ .
- ii) Define the set of points we wish to predict values at,  $(X^*)$ .
- iii) Define a kernel function to use as the covaraince funtion,  $k(\cdot, \cdot)$ .
- iv) Calculate  $\boldsymbol{\mu}^*, \Sigma^*$  for the points  $X^*$ , using the equations below.
- v) Draw samples from  $\text{Normal}(\boldsymbol{\mu}^*, \Sigma^*)$ . Each sample can be used to infer a function.

### 2.1 Equations

#### Without Noise

$$\begin{aligned}
 \boldsymbol{\mu}^* &= k(X^*, X^*)k(X, X)^{-1}y \\
 \Sigma^* &= k(X^*, X^*) - k(X^*, X)k(X, X)^{-1}k(X, X^*)^T
 \end{aligned}$$

#### With Noise

$$\begin{aligned}
 \boldsymbol{\mu}^* &= k(X^*, X^*)k(X, X+\mathbf{c})^{-1}y \\
 \Sigma^* &= k(X^*, X^*) - k(X^*, X)k(X, X+\mathbf{c})^{-1}k(X, X^*)^T
 \end{aligned}$$

The only difference is the  $+c$  in the  $K(X, X)$  terms.

N.B.  $\boldsymbol{\mu}^* \in \mathbb{R}^N$  &  $\Sigma^* \in \mathbb{R}^{N \times N}$  where  $N := |X^*|$ .

### 2.2 Kernels

Linear	$k(\mathbf{x}, \mathbf{y}) = \sigma^2 \times \mathbf{x}^T \mathbf{y}$
White	$k(\mathbf{x}, \mathbf{y}) = \sigma^2 I$
Periodic	$k(\mathbf{x}, \mathbf{y}) = \sigma^2 \exp \left\{ -\frac{2}{\ell^2} \sin^2 \left( \frac{\pi}{p} \ \mathbf{x} - \mathbf{y}\  \right) \right\}$
	where
	$\ell$ =length scale
	$p$ =period
Radial Basis Function	$k(\mathbf{x}, \mathbf{y}) = \exp \left\{ -\frac{1}{\ell^2} \ \mathbf{x} - \mathbf{y}\ ^2 \right\}$
	where
	$\ell$ =length scale

Vary  $\sigma^2$  depending on noise in readings.

N.B.  $\|\mathbf{x}\| := \sqrt{\sum x_i^2} \implies \|\mathbf{x} - \mathbf{y}\| = \sqrt{\sum (x_i - y_i)^2}$ . This is the *Euclidean Distance*.

### 3 Dirichlet Processes

*Dirichlet Processes* are the class of *Stochastic Processes* whose realisations are probability distributions.

*Dirichlet Processes* take a base distribution,  $f(\cdot)$ , and a concentration parameter,  $\alpha \in \mathbb{R}$ . Realisations become more continuous the greater the value  $\lim \alpha$  tends to.

The following algorithm is used to construct a realisation

- i) With probability  $\frac{\alpha}{\alpha + n - 1}$  draw  $X_n$  from  $f(\cdot)$ .
- ii) With probability  $\frac{n_x}{\alpha + n - 1}$  set  $X_n = x$   
 where  $n_x := |\{j < n : X_j = x\}|$  (*i.e.* the number of previous observations of  $x$ ).

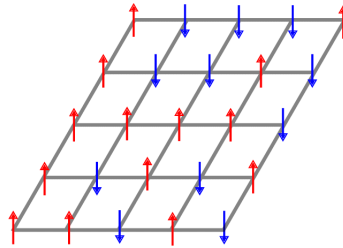
These  $X_1, X_2, \dots$  represent the relative frequencies with which each value should occur in the distribution. In practice we cannot produce full distributions as that would require infinite iterations of the algorithm, instead we are just approximating it.

*N.B.*  $X_1, X_2, \dots$  are not independent since they depend on the previously generated results.

*N.B.* The base distribution,  $f(\cdot)$ , is the expected result.

### 4 Ising Model

An *Ising Model* is one where each latent variable takes one of two states and only has dependency on neighbours in such a way that a grid is formed. Thus dependencies are undirected, forming a *Markov Random Field*.



#### 4.1 Ising Prior

$$p(\mathbf{x}) = \frac{1}{Z_0} e^{\sum_{i \in \mathbf{x}} \sum_{j \in N_i} w_{ij} x_i x_j}$$

where  $Z_0$  is a normalising term,  $N_i$  is the neighbourhood of  $x_i$  &  $w_{ij}$  is the weighting of the relationship between variables  $x_i$  &  $x_j$ .

*N.B.*  $x_i x_j = 1$  iff  $x_i = x_j$ , otherwise  $x_i x_j = -1$ . Thus this term only increases when  $x_i$  has the same value as many of its neighbours.

#### 4.2 Iterative Conditional Modes

*Iterative Conditional Modes* is a technique for inferring latent variable values in an *Ising Model*.

- i) Randomly initialise  $\mathbf{x}$ .
- ii) For each  $x_i \in \mathbf{x}$ .
  - (a) Assume all latent values are fixed except for  $x_i$ .
  - (b) Assign  $x_i$  to the most likely value given the other values.
- iii) Repeat ii) until time is up

### 4.3 Gibbs Sampling

*Gibbs Sampling* is an implementation of *Markov Chain Monte Carlo*.

The general idea is that given a distribution  $p(\mathbf{x})$  which we wish to sample from we shall draw samples from one dimension at a time, using the other dimensions as conditions,  $p(x_i|\mathbf{x}_{-i})$ .

If we now just consider the *Ising Model* we can perform some derivations

$$\begin{aligned}
 p(x_i|\mathbf{x}_{-i}, \mathbf{y}) &= \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x}_{-i}, \mathbf{y})} \\
 &= \frac{p(\mathbf{x}, \mathbf{y})}{\int p(\mathbf{x}, \mathbf{y}) dx_i} \\
 &= \frac{p(\mathbf{x}, \mathbf{y})}{\sum_{x_i \in \{1, -1\}} p(\mathbf{x}, \mathbf{y})} \\
 &= \frac{p(\mathbf{x}, \mathbf{y})}{p(x_i = 1, \mathbf{x}_{-i}, \mathbf{y}) + p(x_i = -1, \mathbf{x}_{-i}, \mathbf{y})} \\
 \implies p(x_i = 1|\mathbf{x}_{-i}, \mathbf{y}) &= \frac{p(x_i = 1, \mathbf{x}_{-i}, \mathbf{y})}{p(x_i = 1, \mathbf{x}_{-i}, \mathbf{y}) + p(x_i = -1, \mathbf{x}_{-i}, \mathbf{y})}
 \end{aligned}$$

This is a tractable expression.

Here is an algorithm for *Gibbs Sampling*

- i) Randomly initialise  $\mathbf{x}$ .
- ii) For each  $x_i \in \mathbf{x}$ .
  - (a) Calculate  $p(x_i = 1|\mathbf{x}_{-i}, \mathbf{y})$ .
  - (b) Draw  $u$  from Uniform[0, 1].
  - (c) If  $p(x_i = 1|\mathbf{x}_{-i}, \mathbf{y}) > u$  set  $x_i = 1$ , otherwise set  $x_i = -1$ .
- iii) Repeat ii) until time is up