

02477 – Bayesian Machine Learning: Lecture 4

Michael Riis Andersen

Technical University of Denmark,
DTU Compute, Department of Applied Math and Computer Science

Outline

- 1 Re-cap from last week and a few words about graphical models
- 2 Bayesian vs. classical statistics
- 3 Bayesian methods for classification
 - Generative modeling
 - Discriminative modelling
- 4 Bayesian logistic regression
- 5 Laplace approximations
- 6 The posterior predictive distribution

Re-cap from last week and a few words about graphical models

Bayesian linear regression model: the key equations

- Linear regression model with Gaussian noise and Gaussian priors

$$y_n = f(\phi(\mathbf{x}_n), \mathbf{w}) + e_n$$

- Given design matrix $\Phi \in \mathbb{R}^{N \times D}$ and observations $\mathbf{y} \in \mathbb{R}^N$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I}) \quad (\text{prior})$$

$$p(\mathbf{y} | \mathbf{w}) = \mathcal{N}(\mathbf{y} | \Phi \mathbf{w}, \beta^{-1} \mathbf{I}) \quad (\text{likelihood})$$

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mathbf{0}, \beta^{-1} \mathbf{I} + \alpha^{-1} \Phi \Phi^T) \quad (\text{marginal likelihood})$$

$$p(\mathbf{w} | \mathbf{y}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N) \quad (\text{posterior})$$

with *posterior parameters*

$$\mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{y} \quad (\text{posterior mean})$$

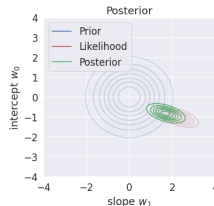
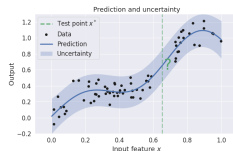
$$\mathbf{S}_N = \left(\alpha \mathbf{I} + \beta \Phi^T \Phi \right)^{-1} \quad (\text{posterior covariance})$$

- *Two hyperparameters*

α : prior precision of the regression weights

β : precision of the measurements

- *Lazy notation*: We should actually write $p(\mathbf{w} | \mathbf{y}, \alpha, \beta)$ etc., but we often suppress dependency of hyperparameter to ease notation



Posterior Predictive distributions

- The *posterior distribution* is $p(\mathbf{w}|\mathbf{y}) = \mathcal{N}(\mathbf{w}|\mathbf{m}, \mathbf{S})$ with

$$\mathbf{m} = \beta \mathbf{S} \Phi^T \mathbf{y}$$
$$\mathbf{S} = (\alpha \mathbf{I} + \beta \Phi^T \Phi)^{-1}$$

- When making predictions \mathbf{x}^* using Bayesian methods, we *average over all possible parameters values weighted by the posterior*

$$f(\mathbf{x}^*|\mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x}^*)$$
$$y(\mathbf{x}^*) = f(\mathbf{x}^*|\mathbf{w}) + \epsilon$$

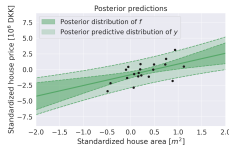
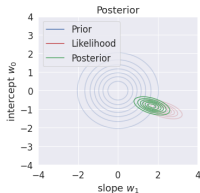
- First, we write the likelihood corresponding to the new input \mathbf{x}_*

$$p(y^*|\mathbf{x}^*, \mathbf{w}) = \mathcal{N}(y^*|\mathbf{w}^T \phi(\mathbf{x}^*), \beta^{-1})$$

- .. and marginalize with respect to the posterior distribution (sum rule)

$$p(y^*|\mathbf{y}) = \int p(y^*|\mathbf{x}^*, \mathbf{w}) p(\mathbf{w}|\mathbf{y}) d\mathbf{w}$$

- This is called *the posterior predictive distribution*



Posterior predictive distributions: how to make predictions?

- First, we write up the *likelihood* corresponding for the new input $\phi_* = \phi(\mathbf{x}_*)$:

$$p(y_* | \mathbf{x}_*, \mathbf{w}) = \mathcal{N}(y_* | \mathbf{w}^T \phi_*, \beta^{-1})$$

- For MAP (and similar for MLE), we simply *plug in* the estimate of \mathbf{w}

$$p(y_* | \mathbf{y}, \mathbf{x}_*) \approx \mathcal{N}(y_* | \hat{\mathbf{w}}_{\text{MAP}}^T \phi_*, \beta^{-1})$$

- Bayesian: use the *sum rule* to *marginalize* wrt. the *posterior* distribution

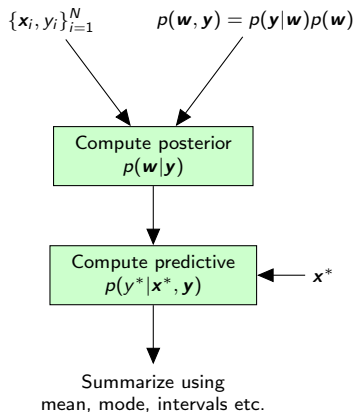
$$p(y_* | \mathbf{y}, \mathbf{x}_*) = \int p(y_* | \mathbf{x}_*, \mathbf{w}) p(\mathbf{w} | \mathbf{y}) d\mathbf{w} = \mathcal{N}(y_* | \hat{\mathbf{w}}_{\text{MAP}}^T \phi_*, \phi_*^T \mathbf{S} \phi_* + \beta^{-1})$$

- If posterior covariance \mathbf{S} is small, we get *approximately* the same result
- We can think of the *MAP* solution as an *approximate* posterior distribution

$$p(\mathbf{w} | \mathbf{y}) \approx \delta(\mathbf{w} - \mathbf{w}_{\text{MAP}}) = \begin{cases} \infty & \text{if } \mathbf{w} = \mathbf{w}_{\text{MAP}} \\ 0 & \text{otherwise} \end{cases}$$

- Such a distribution is called *Dirac's delta* distribution (mental picture: Gaussian with mean \mathbf{w}_{MAP} and variance going to zero)
- MAP is sometimes called *poor man's Bayes*, but can still be a useful tool!

Bayesian inference for supervised learning



- Same principles for linear regression, logistic regression, neural networks etc. etc.

Hyperparameters and the evidence approximation

- Posterior depends on the hyperparameters α and β (but often suppressed in notation)

$$p(\mathbf{w}|\mathbf{y}, \alpha, \beta) = \frac{p(\mathbf{y}|\mathbf{w}, \beta)p(\mathbf{w}|\alpha)}{p(\mathbf{y}|\alpha, \beta)}$$

- We could assign priors to α and β to get the posterior on α and β given the data

$$p(\alpha, \beta|\mathbf{y}) \propto p(\mathbf{y}|\alpha, \beta)p(\alpha)p(\beta)$$

- fully Bayesian solution: use the sum rule to marginalize over all unknowns ($\mathbf{w}, \alpha, \beta$)

$$p(y_*|\mathbf{y}, \mathbf{x}_*) = \int p(y_*|\mathbf{x}_*, \mathbf{w}_*, \beta)p(\mathbf{w}|\mathbf{y}, \alpha, \beta)p(\alpha, \beta|\mathbf{y}) d\mathbf{w}d\alpha d\beta$$

The evidence approximation

- We estimate $\hat{\alpha}, \hat{\beta}$ by *optimizing the marginal likelihood* $p(\mathbf{y}|\alpha, \beta)$

$$\hat{\alpha}, \hat{\beta} = \arg \max_{\alpha, \beta} \log p(\mathbf{y}|\alpha, \beta)$$

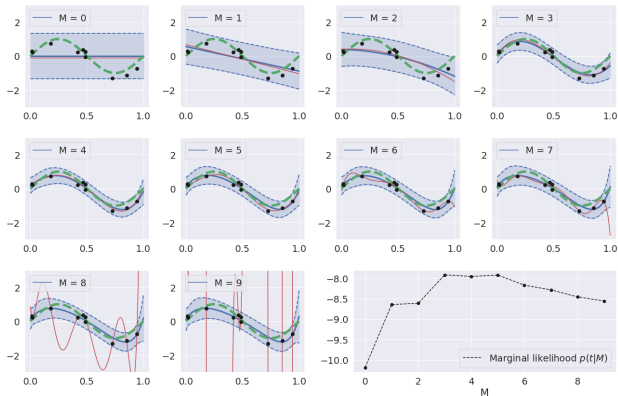
- Equivalent to maximizing the posterior $p(\alpha, \beta|\mathbf{y})$ assuming *flat priors* for α and β

$$p(\alpha, \beta|\mathbf{y}) \propto p(\mathbf{y}|\alpha, \beta)$$

- Equivalent to *poor man's Bayes* on hyperparameter level

Sinusoidal example revisited using the evidence approximation

- Also useful for model selection: $\alpha, \beta, M^* = \arg \max_{\alpha, \beta, M} p(\mathbf{y}|\alpha, \beta, M)$



- Implements "Occam's razor": choose the "simplest" model that explain the data
- Often works well for many models, but we should always assess the generalization error

A more general probabilistic perspective on supervised learning

Product rule

$$p(\mathbf{a}, \mathbf{b}) = p(\mathbf{b}|\mathbf{a})p(\mathbf{a})$$

Sum rule

$$p(\mathbf{b}) = \int p(\mathbf{a}, \mathbf{b})d\mathbf{a}$$

Conditional

$$p(\mathbf{a}|\mathbf{b}) = \frac{p(\mathbf{a}, \mathbf{b})}{p(\mathbf{b})}$$

Conditional independence

$$p(\mathbf{a}, \mathbf{b}|\mathbf{c}) = p(\mathbf{a}|\mathbf{c})p(\mathbf{b}|\mathbf{c})$$

Supervised learning: Given some data $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$, what can we say about a new test point $y^* = y(\mathbf{x}^*)$?

- Step 1: Formulate joint distribution for all variables of interests

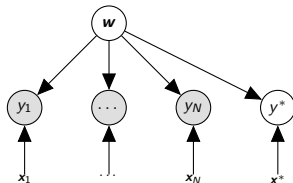
$$p(y^*, \mathbf{y}, \mathbf{w}) = p(y^*, \mathbf{y}|\mathbf{w})p(\mathbf{w}) = p(y^*|\mathbf{w})p(\mathbf{y}|\mathbf{w})p(\mathbf{w})$$

- Step 2: Conditioned on the observed data \mathbf{y}

$$p(y^*, \mathbf{w}|\mathbf{y}) = \frac{p(y^*|\mathbf{w})p(\mathbf{y}|\mathbf{w})p(\mathbf{w})}{p(\mathbf{y})}$$

- Step 3: Marginalize over all parameters \mathbf{w} using sum rule

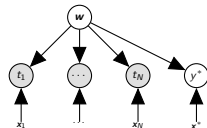
$$p(y^*|\mathbf{y}) = \int \frac{p(y^*|\mathbf{w})p(\mathbf{y}|\mathbf{w})p(\mathbf{w})}{p(\mathbf{y})}d\mathbf{w} = \int p(y^*|\mathbf{w})p(\mathbf{w}|\mathbf{y})d\mathbf{w}$$



A more complete graphical model

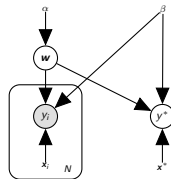
- Lazy, but common notation for Bayesian linear regression

$$p(y^*, \mathbf{y}, \mathbf{w}) = p(y^* | \mathbf{w}) p(\mathbf{y} | \mathbf{w}) p(\mathbf{w})$$



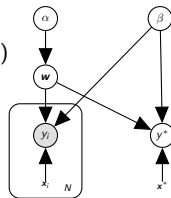
- More complete notation and corresponding graphical model

$$p(y^*, \mathbf{y}, \mathbf{w} | \alpha, \beta, \mathbf{X}, \mathbf{x}^*) = p(y^* | \mathbf{w}, \beta, \mathbf{x}^*) p(\mathbf{y} | \mathbf{w}, \beta, \mathbf{X}) p(\mathbf{w} | \alpha)$$



- Fully Bayesian inference on hyperparameter level

$$p(y^*, \mathbf{y}, \mathbf{w}, \alpha, \beta | \mathbf{X}, \mathbf{x}^*) = p(y^* | \mathbf{w}, \beta, \mathbf{x}^*) p(\mathbf{y} | \mathbf{w}, \beta, \mathbf{X}) p(\mathbf{w} | \alpha) p(\alpha) p(\beta)$$



Bayesian vs. classical statistics

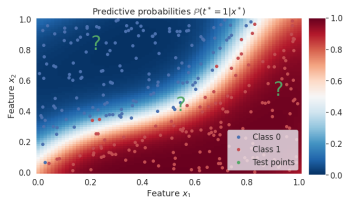
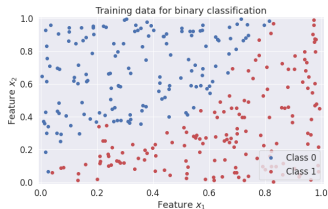
Bayesian vs. classical statistics

	Frequentist/classical	Bayesian
Probability interpretation	Long run frequencies	Degrees of belief
Parameters	Deterministic, but unknown Cannot make probabilistic statement about parameters	Random variables Probabilistic reasoning at levels: models, parameters and observations
Intepretation of intervals	Point estimates <i>Confidence intervals</i> If the experiment is repeated infinitely many times, 95% of the intervals will contain the true population value	Probability distributions <i>Credibility intervals</i> The interval will contain the population value with 95% probability given the data
Sources of information	Data only	Data & prior knowledge
Computation	Often less computationally expensive	Often more computationally expensive

Bayesian methods for classification

Probabilistic approaches for classification

- Dataset $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$
 - Input features: $\mathbf{x}_i \in \mathbb{R}^D$
 - Targets: $y_i \in \{0, 1\}$
- How to predict label for test point $\mathbf{x}^* \in \mathbb{R}^D$?
- *Predictive distributions*: what is the probability that \mathbf{x}^* belong to class 1, i.e. $p(y^* = 1 | \mathcal{D}, \mathbf{x}^*)$?
- Two probabilistic approaches
 1. *Discriminative methods*
 2. *Generative methods*



Discriminative vs generative methods

- The *generative approach* models the *joint distribution* $p(\mathbf{x}_n, y_n)$, e.g. via Bayes rule

$$p(y_n = k | \mathbf{x}_n) = \frac{p(\mathbf{x}_n | y_n = k)p(y_n = k)}{p(\mathbf{x}_n)}$$

- Joint distribution over inputs \mathbf{x}_n and labels y_n allow sampling (*generating*) from the model

$$\mathbf{x}^{(i)}, y^{(i)} \sim p(\mathbf{x}, y)$$

- Pros and cons for generative models

- + Optimal if the assumptions are correct
- + Can easily handle *missing data*
- + Can reason about input data
- Assumptions are often hard to get correct

- The *discriminative approach* models the *conditional distribution* $p(y_n | \mathbf{x}_n)$ directly by assuming some parametric form for the posterior

$$p(y_n | \mathbf{x}_n) = f(\mathbf{x}_n | \mathbf{w})$$

- The function $f(\mathbf{x}_n | \mathbf{w})$ can be based on a linear model, a neural network etc.

- Pros and cons for discriminative models

- + Often superior when the assumptions for generative models are wrong
- + Often better calibrated (compared to e.g. generative methods like Naïve Bayes etc)
- + Easy to make flexible
- Difficult to handle missing data
- Cannot reason about input data

Bayesian methods for classification: Generative modeling

The generative approach I

- Binary classification $y_n \in \{0, 1\}$

$$p(y_n = 1 | \mathbf{x}_n) = \frac{p(\mathbf{x}_n | y_n = 1)p(y_n = 1)}{p(\mathbf{x}_n)}$$

- Terminology

- *Class-conditional distribution* $p(\mathbf{x}_n | y_n)$

- *Prior probabilities* $p(y_n = k) = \pi_k$

- *Marginal data density* $p(\mathbf{x}_n)$

- The marginal density of \mathbf{x}_n is a *mixture distribution* and is obtained using the *sum rule*

$$p(\mathbf{x}_n) = \sum_{k \in \{0,1\}} p(\mathbf{x}_n | y_n = k)p(y_n = k) = \pi_0 p(\mathbf{x}_n | y_n = 0) + \pi_1 p(\mathbf{x}_n | y_n = 1)$$

- Let's plug the result into Bayes' rule

The generative approach II

- The posterior of y_n given the input \mathbf{x}_n

$$p(y_n = 1|\mathbf{x}_n) = \frac{\pi_1 p(\mathbf{x}_n|y_n = 1)}{\pi_0 p(\mathbf{x}_n|y_n = 0) + \pi_1 p(\mathbf{x}_n|y_n = 1)}$$

- Divide by numerator

$$p(y_n = 1|\mathbf{x}_n) = \frac{1}{1 + \frac{\pi_0 p(\mathbf{x}_n|y_n=0)}{\pi_1 p(\mathbf{x}_n|y_n=1)}}$$

- Define

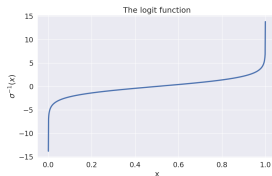
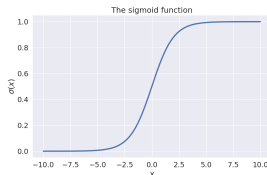
$$a = \ln \frac{\pi_1 p(\mathbf{x}_n|y_n = 1)}{\pi_0 p(\mathbf{x}_n|y_n = 0)}$$

- then

$$p(y_n = 1|\mathbf{x}_n) = \frac{1}{1 + \exp(-a)} = \sigma(a)$$

- Recall $\sigma(a)$ is the *logistic sigmoid* function and its inverse is called the *logit* function

$$a = \ln \left(\frac{\sigma}{1-\sigma} \right)$$



The generative approach III: multi-class problems and softmax

- Assume we have K different classes, where $k = 1, \dots, K$

- Define a_k

$$a_k = \ln p(\mathbf{x}_n | y_n = k) p(y_n = k)$$

- Using similar line of reasoning for K classes

$$p(y_n = k | \mathbf{x}_n) = \frac{p(\mathbf{x}_n | y_n = k) p(y_n = k)}{\sum_{i=1}^K p(\mathbf{x}_n | y_n = i) p(y_n = i)} = \frac{\exp(a_k)}{\sum_{j=1}^K \exp(a_j)}$$

- The *normalized exponentials* is the known as the *softmax* function

Example: Gaussian class conditionals

- Binary classification, normal class conditionals with common variance in 1D

$$p(x_n|y_n = 0) = \mathcal{N}(x_n|\mu_0, \sigma^2)$$

$$p(x_n|y_n = 1) = \mathcal{N}(x_n|\mu_1, \sigma^2)$$

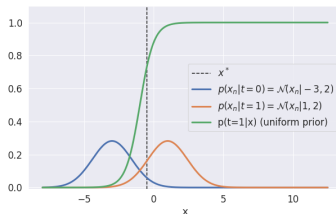
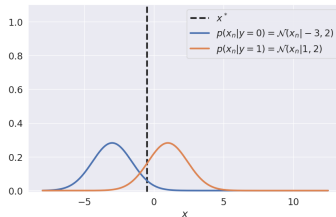
- We know

$$p(y_n = 1|\mathbf{x}_n) = \frac{1}{1 + \exp(-a)} = \sigma(a) = \sigma(w_0 + w_1 x_n)$$

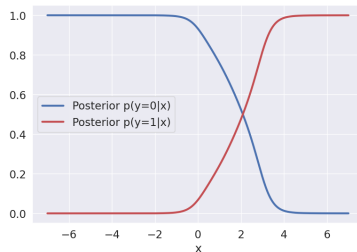
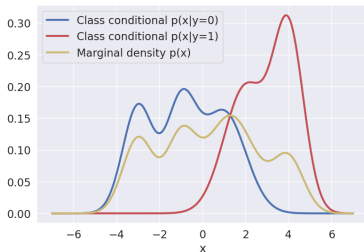
- The quantity a has a simple expression

$$\begin{aligned} a &= \ln \frac{\pi_1 p(\mathbf{x}_n|y_n = 1)}{\pi_0 p(\mathbf{x}_n|y_n = 0)} = \ln \frac{\pi_1 \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{(x_n - \mu_1)^2}{2\sigma^2})}{\pi_0 \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{(x_n - \mu_0)^2}{2\sigma^2})} \\ &= \ln \frac{\pi_1}{\pi_0} - \frac{\mu_1^2 - \mu_0^2}{2\sigma^2} + \frac{\mu_1 - \mu_0}{\sigma^2} x_n \\ &= w_0 + w_1 x_n \end{aligned}$$

$$\text{where } w_0 = \ln \frac{\pi_1}{\pi_0} - \frac{\mu_1^2 - \mu_0^2}{2\sigma^2} \text{ and } w_1 = \frac{\mu_1 - \mu_0}{\sigma^2}$$



Example 2: More complex distributions



Bayesian methods for classification: Discriminative modelling

Discriminative modelling for binary classification

- In the generative model, we defined *priors* $p(y_n = 1) = \pi_1$ and a set of *class-conditionals* $p(\mathbf{x}_n | y_n = 1)$, applied Bayes rule and ended up with

$$a(\mathbf{x}_n) = \ln \frac{\pi_1 p(\mathbf{x}_n | y_n = 1)}{\pi_0 p(\mathbf{x}_n | y_n = 0)} = \ln \frac{p(y_n = 1 | \mathbf{x}_n)}{p(y_n = 0 | \mathbf{x}_n)}$$

and

$$p(y_n = 1 | \mathbf{x}_n) = \frac{1}{1 + \exp(-a)} = \sigma(a(\mathbf{x}_n))$$

- The distributional assumptions gave the specific functional form for $a(\mathbf{x})$, but in *discriminative modelling*, we directly assume a functional form for $a(\mathbf{x})$
- Example: logistic regression

$$p(y_n = 1 | \mathbf{x}_n, \mathbf{w}) = \frac{1}{1 + \exp(-a)} = \sigma(\phi(\mathbf{x}_n)^T \mathbf{w})$$

- We model each observation with a *Bernoulli* distribution with probability $\sigma(\phi(\mathbf{x}_n)^T \mathbf{w})$
- We *estimate* \mathbf{w} using maximum likelihood, MAP or Bayesian inference with the likelihood function

$$p(\mathbf{y} | \mathbf{w}, \mathbf{X}) = \prod_{n=1}^N \sigma(\phi(\mathbf{x}_n)^T \mathbf{w})^{y_n} (1 - \sigma(\phi(\mathbf{x}_n)^T \mathbf{w}))^{1-y_n}$$

Maximum likelihood estimator for logistic regression: Quiz

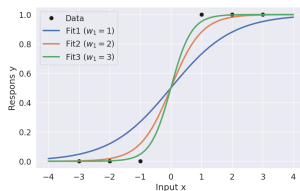
Set-up

- Consider a simple dataset with $N = 6$, $\phi(x) = x$

- Logistic regression likelihood

$$p(\mathbf{y}|\mathbf{w}_1) = \prod_{n=1}^N \sigma(w_1 x_n)^{y_n} (1 - \sigma(w_1 x_n))^{1-y_n}$$

- One parameter: w_1



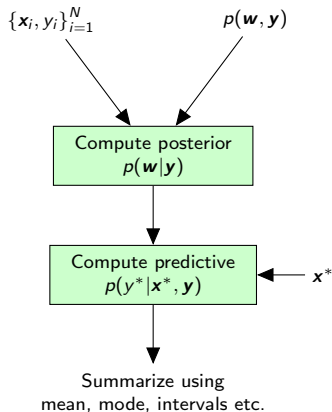
Questions (5mins)

- Spend 5 minutes DTU Learn quiz: "Lecture 4: Logistic regression"

Bayesian logistic regression

Bayesian supervised learning

- For conjugate models, both the posterior and predictive distributions can be computed analytically
- The real strength of the Bayesian framework lies in the modelling flexibility
- ... but for even rather simple models like *Bayesian logistic regression*, we cannot compute
 1. the posterior distribution
 2. the predictive distribution
- Today we will see how to use the *Laplace approximation* to approximate the posterior
- ... and we will discuss different strategies to evaluate the predictive distribution



Bayesian logistic regression

- Likelihood for logistic regression

$$p(\mathbf{y}|\mathbf{w}) = \prod_{n=1}^N \sigma(\phi(\mathbf{x}_n)^T \mathbf{w})^{y_n} (1 - \sigma(\phi(\mathbf{x}_n)^T \mathbf{w}))^{1-y_n}$$

- Let's impose a prior distribution on the weights \mathbf{w} assuming the individual weights w_i are *independent and identically distributed (i.i.d)* a priori

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1} \mathbf{I}) = \prod_{i=1}^D \mathcal{N}(w_i|\mathbf{0}, \alpha^{-1})$$

- The posterior follows from Bayes' theorem

$$p(\mathbf{w}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{w})p(\mathbf{w})}{p(\mathbf{y})}$$

- Let's calculate the posterior mean

$$\mathbb{E}_{p(\mathbf{w}|\mathbf{y})} [\mathbf{w}] = \int \mathbf{w} p(\mathbf{w}|\mathbf{y}) d\mathbf{w} = \frac{1}{p(\mathbf{y})} \int \mathbf{w} p(\mathbf{y}|\mathbf{w}) p(\mathbf{w}) d\mathbf{w}$$

- Clearly, we need $p(\mathbf{y})$ as well

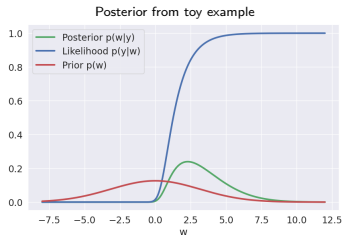
$$p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{w}) p(\mathbf{w}) d\mathbf{w}$$

Bayesian logistic regression II

- General problem: we *cannot* compute the posterior mean analytically

$$\mathbb{E}_{p(\mathbf{w}|\mathbf{y})}[\mathbf{w}] = \frac{1}{p(\mathbf{y})} \int \mathbf{w} p(\mathbf{y}|\mathbf{w}) p(\mathbf{w}) d\mathbf{w}$$

- This *roadblock* occurs for almost any other interesting posterior summary
- The posterior distribution and the marginal likelihood for most Bayesian models is *analytically intractable*
- Posterior distribution from our toy example is asymmetric, but almost resembles a Gaussian
- *Bernstein von Mises theorem*
Assuming certain regularity conditions, the posterior distribution of a parametric model becomes more and more Gaussian as N increases
- Let's approximate $p(\mathbf{w}|\mathbf{y})$ with a Gaussian!



Laplace approximations

Laplace approximations I

- The *Laplace approximation* is a method for approximating intractable probability densities

- Assume we have a posterior distribution of interest $p(\mathbf{w}|\mathbf{y})$

$$p(\mathbf{w}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{w})p(\mathbf{w})}{p(\mathbf{y})} = \frac{1}{Z} f(\mathbf{w}) \approx \mathcal{N}(\mathbf{w}|\mathbf{m}, \mathbf{S})$$

- The log density for Gaussians is quadratic wrt. \mathbf{w}

$$\ln \mathcal{N}(\mathbf{w}|\mathbf{m}, \mathbf{S}) = -\frac{D}{2} \ln(2\pi) - \frac{1}{2} \log |\mathbf{S}| - \frac{1}{2} (\mathbf{w} - \mathbf{m})^T \mathbf{S}^{-1} (\mathbf{w} - \mathbf{m})$$

- Let's make a second order Taylor expansion of $f(\mathbf{w})$ around the mode \mathbf{w}_{MAP}

$$\ln f(\mathbf{w}) \approx \ln f(\mathbf{w}_{\text{MAP}}) - \frac{1}{2} (\mathbf{w} - \mathbf{w}_{\text{MAP}})^T \mathbf{A} (\mathbf{w} - \mathbf{w}_{\text{MAP}}),$$

where \mathbf{A} is the Hessian at the mode, i.e. $\mathbf{A} = -\nabla \nabla \ln f(\mathbf{w})|_{\mathbf{w}=\mathbf{w}_{\text{MAP}}}$

- The Laplace approximation is defined

$$q(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{w}_{\text{MAP}}, \mathbf{A}^{-1})$$

- That is, we approximate the posterior mean using the MAP and the posterior covariance using the curvature at the MAP solution.

Laplace approximation II

- Suppose we want to approximate $p(\mathbf{w}|\mathbf{y})$ using the Laplace approximation

$$p(\mathbf{w}|\mathbf{y}) \approx q(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{w}_{\text{MAP}}, \mathbf{A}^{-1})$$

- Computational steps

1. Locate the mode of $p(\mathbf{w}|\mathbf{y})$

$$\mathbf{w}_{\text{MAP}} = \arg \max_{\mathbf{w}} p(\mathbf{w}|\mathbf{y}) = \arg \max_{\mathbf{w}} p(\mathbf{y}|\mathbf{w})p(\mathbf{w})$$

2. Evaluate the Hessian at \mathbf{w}_{MAP}

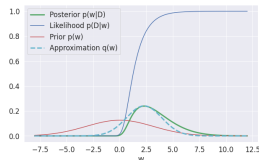
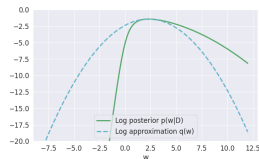
$$\mathbf{A} = -\nabla \nabla \ln p(\mathbf{y}|\mathbf{w})p(\mathbf{w}) \Big|_{\mathbf{w}=\mathbf{w}_{\text{MAP}}}$$

- Advantages

1. Simple and well-understood
2. Very fast to compute
3. Gives good results for many problems

- Limitations

1. Only applies to continuous parameters
2. Gaussian (symmetric distribution, thin tails)
3. Only capture local properties of $p(\mathbf{w}|\mathbf{y})$ near \mathbf{w}_{MAP}
4. Does not work for hierarchical models in general



Laplace approximations III: Approximating the marginal likelihood

- Our second order Taylor approximation for $\ln f(\mathbf{w})$

$$\ln f(\mathbf{w}) \approx \ln f(\mathbf{w}_{\text{MAP}}) - \frac{1}{2}(\mathbf{w} - \mathbf{w}_{\text{MAP}})^T \mathbf{A}(\mathbf{w} - \mathbf{w}_{\text{MAP}}),$$

- By assumption

$$p(\mathbf{w}|\mathbf{y}) = \frac{1}{Z} f(\mathbf{w}) \quad \Rightarrow \quad Z = \int f(\mathbf{w}) d\mathbf{w}$$

- Plugging in the approximation for $\ln f(\mathbf{w})$

$$\begin{aligned} Z &= \int f(\mathbf{w}) d\mathbf{w} \\ &\approx f(\mathbf{w}_{\text{MAP}}) \int \exp\left(-\frac{1}{2}(\mathbf{w} - \mathbf{w}_{\text{MAP}})^T \mathbf{A}(\mathbf{w} - \mathbf{w}_{\text{MAP}})\right) \\ &= f(\mathbf{w}_{\text{MAP}}) \frac{(2\pi)^{\frac{D}{2}}}{|\mathbf{A}|^{\frac{1}{2}}} \end{aligned}$$

- Using $f(\mathbf{w}) = p(\mathbf{y}|\mathbf{w})p(\mathbf{w})$, our approximation for $p(\mathbf{y})$ becomes

$$\ln p(\mathbf{y}) \approx \ln p(\mathbf{y}|\mathbf{w}_{\text{MAP}}) + \ln p(\mathbf{w}_{\text{MAP}}) + \frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{A}|$$

- Very useful for model selection, parameter tuning etc

The posterior predictive distribution

How to make predictions?

- For classification, we need the predictive distribution for a new input \mathbf{x}^* . The likelihood for a input data point \mathbf{x}^* is

$$p(y^* = 1 | \mathbf{w}, \mathbf{x}^*) = \sigma(\phi(\mathbf{x}^*)^T \mathbf{w})$$

- As always, we want to take the posterior uncertainty into account using the sum rule

$$\begin{aligned} p(y^* = 1 | \mathbf{y}, \mathbf{x}^*) &= \int p(y^* = 1 | \mathbf{x}^*, \mathbf{w}) p(\mathbf{w} | \mathbf{y}) d\mathbf{w} \\ &\approx \int p(y^* = 1 | \mathbf{x}^*, \mathbf{w}) q(\mathbf{w}) d\mathbf{w} \\ &= \int \sigma(\phi(\mathbf{x}^*)^T \mathbf{w}) \mathcal{N}(\mathbf{w} | \mathbf{m}, \mathbf{S}) d\mathbf{w} \\ &= \int \sigma(f) \mathcal{N}(f | \mu, \sigma^2) df \end{aligned}$$

where

$$\mu = \phi(\mathbf{x}^*)^T \mathbf{m} \qquad \sigma^2 = \phi(\mathbf{x}^*)^T \mathbf{S} \phi(\mathbf{x}^*)$$

- The good news: we only have to calculate 1D integrals to make predictions
- The bad news: the integral does not have analytical solution

Evaluating predictive distributions for logistic regression

How does uncertainty in f affect the distribution of $\sigma(f)$?

$$p(y^* = 1 | \mathbf{y}, \mathbf{x}^*) = \int \sigma(f) \mathcal{N}(f | \mu, \sigma^2) df$$

■ General strategies for evaluating this integral

1. Monte Carlo methods (sampling)

$$p(y^* = 1 | \mathbf{y}, \mathbf{x}^*) \approx \frac{1}{S} \sum_{i=1}^S \sigma(f^{(i)}) \quad \text{for} \quad f^{(i)} \sim \mathcal{N}(f | \mu, \sigma^2)$$

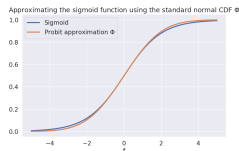
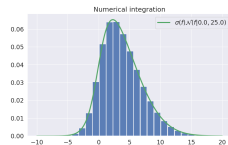
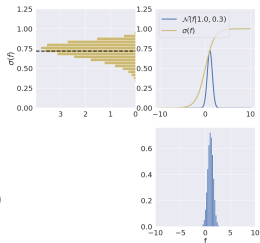
2. Numerical integration (Gauss-Hermite integration)

$$p(y^* = 1 | \mathbf{y}, \mathbf{x}^*) \approx \frac{1}{\sqrt{\pi}} \sum_{i=1}^S w_i h(\sqrt{2}\sigma x_i + \mu)$$

3. Probit approximation

$$\sigma(y) \approx \Phi\left(y \sqrt{\frac{\pi}{8}}\right)$$

where Φ is the CDF of the standard normal



Let's zoom out and summarize

- We introduced *logistic regression* as a *discriminative* approach for binary classification
- We saw to use the *Laplace approximation* to approximate the *posterior* of the weights
- We briefly discussed three strategies to compute *the predictive distribution*
 1. Sampling
 2. Numerical integration
 3. Probit approximation

