# SPR day 10

Bayesian Logistic Regression
- Posterior
- Evidence
- Predictive Posterior

Reading: Sec 4.4 of Bishop PRML
Sec 4.5

Outline: (1) Overview of posterior + predictive

(2) Laplace approximation
in 1-dim and M-dim

(3) Laplace approximation for logistic regression

(4) Posterior predictive approximation

# Bayesian Logistic Regression

Model:

Prior on weights $w \in \mathbb{R}^M$

Usually $m_0 = \vec{0}$ all zero
$S_0 = \alpha^{-1} I_M$ with $\alpha > 0$

$$p(w) = N(w \mid \underset{\text{mean}}{m_0}, \underset{\text{covar}}{S_0})$$

Just like linear regr.

Likelihood of "outputs" $t_n \in \{0, 1\}$ (binary)

$$p(t \mid w) = \prod_{n=1}^{N} \text{Bern}(t_n \mid \sigma(w^T \phi(x_n)))$$

treat x as known, fixed

iid across examples

Goals are to estimate the posterior and predictive

Posterior: $p(w \mid t_{1:N})$

no closed form!
Not a Gaussian!

Predictive:
$$p(t_* \mid t_{1:N})$$

if we only denote random variables

$$= p(t_* \mid x_*, \{x_n, t_n\}_{n=1}^{N})$$

if we make known x vals explicit

$$= \int_w \underset{\text{likelihood}}{p(t_* \mid w, x_*)} \underset{\text{posterior}}{p(w \mid \{x_n, t_n\}_{n=1}^{N})} dw$$

Must be a Bernoulli
(r.v. $t_*$ is binary)
But no closed-form
for its parameter

# Laplace Approximation in 1D

Given: a random variable $w \in \mathbb{R}$
whose density $p(w)$ is known up to norm. const.

$$p(w) = \frac{1}{Z} f(w) \longleftrightarrow \log p(w) = \log f(w) + \text{const}$$

Here, $f(w) > 0$ is known and evaluatable and differentiable
but computing $Z = \int_W f(w)\, dw$ is hard

How can we estimate the distribution $p(w)$?

Idea: Approximate with a Gaussian: $q(w) = N(m, \beta^{-1})$    mean $m$   precision $\beta > 0$
- pick mean to match the <u>mode</u> of $p(w)$

$$M = \underset{w \in \mathbb{R}}{\text{argmax}}\ p(w) = \underset{w \in \mathbb{R}}{\text{argmax}}\ f(w)$$

     Can use Gradient Methods to solve this numerically

- pick precision to perform best possible 2nd-order
   Taylor approximation to $p(w)$ at the mode $w = m$

$$\beta = \frac{\partial}{\partial w}\frac{\partial}{\partial w}\left[-\log f(w)\right]\Big|_{w=m}$$

$$= -\ell''(m) \quad \text{where } \ell = \log f(w)$$

Advantages: Gives an approx distribution we can reason about!
     Second derivatives are often tractable

Limitations: bad if $p(w)$ multimodal
     bad if $p(w)$ has heavy tails, not symmetric about mode

# Derivation of Taylor approx to density $p(w)$ at $m = \underset{w}{\text{argmax}}\; p(w)$

$$\log p(w) = \log f(w) + \text{const wrt } w$$

by definition of $p(w)$

$$= \ell(w) + \text{const}_1$$

define $\ell(w) = \log f(w)$
note that $m$ is a mode of $\ell(w)$ too!

$$= \ell(m) + \ell'(m)(w-m) + \frac{\ell''(m)}{2}(w-m)^2 + \text{const}_1$$

(with $\ell'(m)(w-m)$ marked as $\to 0$)

2nd order Taylor approx. to func. $\ell$ at $w = m$

$$= -\frac{1}{2}\left[-\ell''(m)\right](w-m)^2 + \text{const}_2$$

$\ell(m)$ is const wrt $w$, so group w/ const

$\ell'(m) = 0$ bc. $m$ is a maximer of $f(w)$ & $\ell(w)$ so this term cancels

this is a Gaussian pdf
with mean $m = \underset{w}{\text{argmax}}\; \ell(w)$
and precision $\beta = -\ell''(m)$

# Laplace Approx in M-Dim.

Given: random variable $w \in \mathbb{R}^M$
whose PDF is known up to norm. const.

$$p(w) = \frac{1}{Z} f(w) \quad \longleftrightarrow \quad \log p(w) = \log f(w) + \text{const}$$

We know $f(w) > 0$ and its 1st/2nd derivatives
wrt vector $w$

Idea: Approximation: $q(w) = MVNorm\left(m, \Delta^{-1}\right)$

mean    precision: symmetric
        matrix    &
                  positive
                  definite!

Set $m$ to match mode of $f$

$$m = \arg\max_{w \in \mathbb{R}^M} \ell(w) \qquad \ell(w) = \log f(w)$$

Set $\Delta$ to negative Hessian at mode

$$\Delta = \nabla_w \nabla_w \left[ - \ell(w) \right] \Big|_{w=m}$$

Similar Pro/Con as in 1-dim case above

# Approximating the Posterior for Logistic Regression

True posterior intractable, but known to constant via Bayes

$$\log p(w \mid t_{1:N}) = \log p(w) + \log p(t_{1:N} \mid w) - \log p(t_{1:N})$$

$$= \underbrace{\log \text{MVNormPDF}(w \mid m_0, S_0) + \sum_{n=1}^{N} \log \text{BernPMF}\left(t_n \mid \sigma(w^T \phi(x_n))\right)}_{\ell(w)} \quad \underbrace{\log Z}_{\substack{\text{constant} \\ \text{wrt } w}}$$

Thus, we can apply Laplace Approx!

$$p\left(w \mid \{x_n, t_n\}_{n=1}^{N}\right) \approx N\left(w \mid m_{MAP}, S\right)$$

with mean vector $m_{MAP} = \underset{w \in \mathbb{R}^M}{\arg\max} \ \ell(w)$

Solved w/ Gradient descent

and precision matrix $S^{-1} = \nabla_w \nabla_w \left[-\ell(w)\right]\Big|_{w = m_{MAP}}$

uses plug-in formula

Using formulas for the Hessian of MAP objective, we know:

$$S^{-1} = S_0^{-1} + \Phi^T R(m_{MAP}) \Phi$$

$$= \alpha I_M + \Phi^T R(m_{MAP}) \Phi$$

positive scalar     $M \times M$ outer product

Recall that $\Phi^T R \Phi = \sum_{n=1}^{N} r(m_{MAP}, x_n) \phi(x_n) \phi(x_n)^T$

$$r = \sigma(m_{MAP}^T \phi_n)\left(1 - \sigma(m_{MAP}^T \phi_n)\right)$$

Mean precision

Punchline: Laplace parameters $M, S^{-1}$ possible to compute

# Posterior Predictive for Logistic Regr.

## Ideal (intractable) posterior predictive:

$$p(t_* | t_{1:N}) = \int_{w \in \mathbb{R}^M} p(t_* | w) \, p(w | t_{1:N}) \, dw$$

**Laplace approximation**

$$p(w | t_{1:N}) \approx N(m_{MAP}, S)$$

$$\approx \int_{w \in \mathbb{R}^M} p(t_* | w) \, N(w | m_{MAP}, S) \, dw$$

Still a tough integral over M dimensions
If M was 1 or 2, could use numerical
strategies like trapezoid approx.

### Option 1: Monte Carlo
Easy but need many samples **L**

$$p(t_* | t_{1:N}) = \mathbb{E}_{p(w | t_{1:N})} \left[ p(t_* | w) \right]$$

Average of L
samples from
approx. posterior

$$\approx \frac{1}{L} \sum_{\ell=1}^{L} p(t_* | w^\ell)$$

with $w^\ell \sim N(m_{MAP}, S)$

### Option 2: Probit approx.
Closed-form, hard to port to other models

Likelihood $p(t_* | w) = \begin{cases} \sigma(w^T \phi_*) & \text{if } t_* = 1 \\ \sigma(-w^T \phi_*) & \text{if } t_* = 0 \end{cases}$

See Bishop
Fig 4.9

$\sigma(a) \approx NormCDF(\sqrt{\frac{\pi}{8}} a)$

$\approx \begin{cases} NormCDF(\sqrt{\frac{\pi}{8}} w^T \phi_*) & \text{if } t_* = 1 \\ 1 - \text{above} & t_* = 0 \end{cases}$

Makes integral tractable when combined w/ Laplace!

(4.155): $p(t_* = 1 | t_{1:N}) = \sigma\left( m_{MAP}^T \phi_* \cdot \frac{1}{\sqrt{1 + \frac{\pi}{8} \phi_*^T S \phi_*}} \right)$