

STATS271/371: Applied Bayesian Statistics

Lap 2: Generalized Linear Models and the Laplace Approximation

Scott Linderman

April 12, 2021

Announcements

- ▶ Lecture slides and course content on Github:
<https://github.com/slinderman/stats271sp2021>
- ▶ Homework 2 due Friday (April 16).

Final Projects

- ▶ **Goal:** Go through Box's Loop on a real dataset.
 - ▶ Find an interesting real dataset, e.g. from your own research, a recent paper, or a public repository. (Synthetic datasets are not sufficient.)
 - ▶ Develop a model.
 - ▶ Think carefully about Bayesian inference algorithms and run one.
 - ▶ Evaluate your model, then revise it in light of your analysis.
 - ▶ Repeat...
 - ▶ Present your results and findings in the form of a short report.
- ▶ You may work individually or in a team of two.
- ▶ Proposals due late April/early May (TBD).

Survey Results

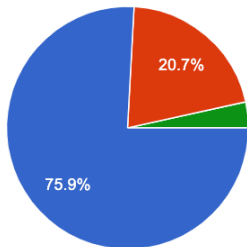
Where are you zooming in from?

- ▶ Campus or nearby: 24/29 (!!!)
- ▶ San Francisco
- ▶ India
- ▶ Senatobia, MS
- ▶ North Carolina
- ▶ Houston, TX

Survey Results II

What's your programming language of choice?

29 responses



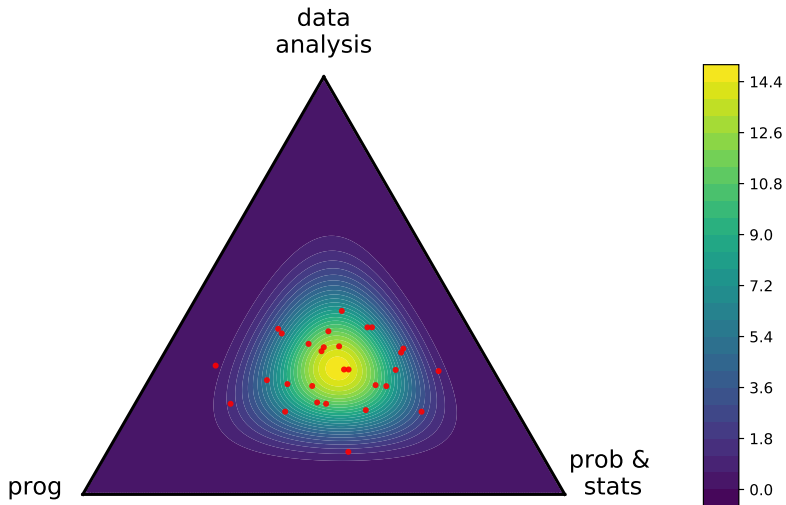
● Python

● R

● Julia

● I will likely use python, but I am a big fan of Mathematica which is used a lot in physics.

Survey Results III

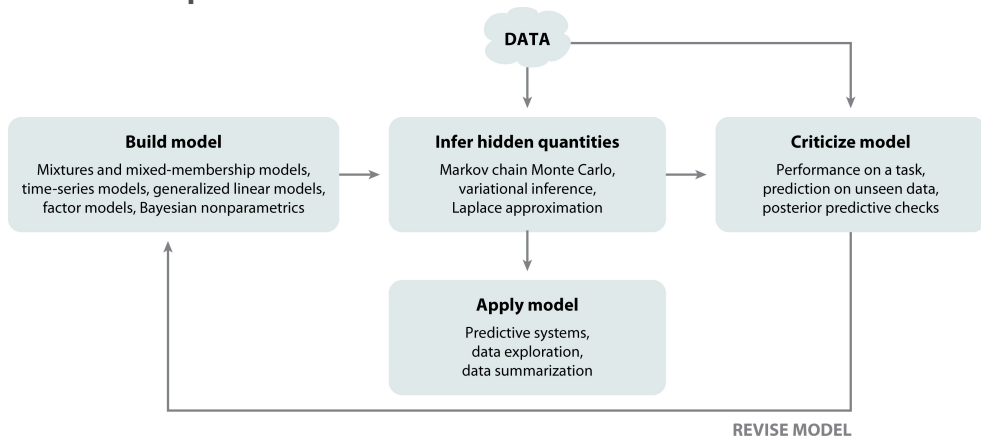


Survey Results IV

Any pandemic hobbies?

- ▶ Fingerstyle guitar
- ▶ Cycling
- ▶ Hiking
- ▶ Playing music
- ▶ Rock hard abs
- ▶ Golf
- ▶ Jouranling / Tik Tok
- ▶ Tennis (x3)
- ▶ Checking Covid numbers :(
- ▶ Cooking (x3)
- ▶ Reading (x2)
- ▶ Looking after my newborn :)
- ▶ Astrophotography
- ▶ Yoga/Pokemon Go
- ▶ Chess
- ▶ Learning Spanish
- ▶ Dominion
- ▶ Walking
- ▶ Piano

Lap 2 of Box's Loop

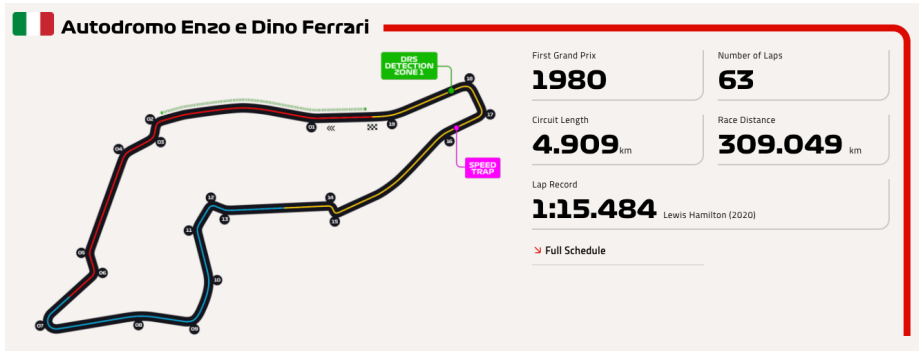


Blei DM. 2014.

Annu. Rev. Stat. Appl. 1:203–32

Blei, *Ann. Rev. Stat. App.* 2014.

Lap 2 of Box's Loop



<https://www.formula1.com/en/racing/2021/EmiliaRomagna/Circuit.html>

Bayesian Generalized Linear Models

Lap 2 of Box's loop will introduce:

- ▶ **Model:** Bayesian generalized linear models
- ▶ **Algorithm:** Laplace approximation
- ▶ **Criticism:** Posterior predictive checks

Motivation

Linear regression assumed real-valued observations, but what if our data has some other support?

- ▶ $y_n \in \{0, 1\}$ binary observation; e.g. win/lose, healthy/sick.
- ▶ $y_n \in \{0, 1, 2, \dots\}$ count observation; e.g. number of spikes in 25ms bin.
- ▶ $y_n \in \{0, 1, \dots, M\}$ count observation (with max); e.g. number of races won out of M total
- ▶ $y_n \in \{1, \dots, K\}$ unordered, categorical observation; e.g.
 - ▶ Democrat/Republican/Independent
 - ▶ Hamilton/Madison/Jay (authors of the Federalist papers)
 - ▶ Hamilton/Verstappen/Bottas (only F1 drivers who win races)
 - ▶ Stanford/Berkeley/Not-Stanford-or-Berkeley (only institutions of higher learning)

We could still apply linear regression, but it would make weird predictions. *There will be -2 spikes in the next time bin? Bet on Hamilstappen to win?*

Notation

Let

- ▶ $y_n \in \mathcal{Y}$ denote the n -th *observation* (type to be specified later)
- ▶ $\mathbf{x}_n \in \mathbb{R}^P$ denote a the *covariates* (aka features) correspond the n -th datapoint
- ▶ $\mathbf{w} \in \mathbb{R}^P$ denote the *weights* of the model

Example: Bernoulli GLM

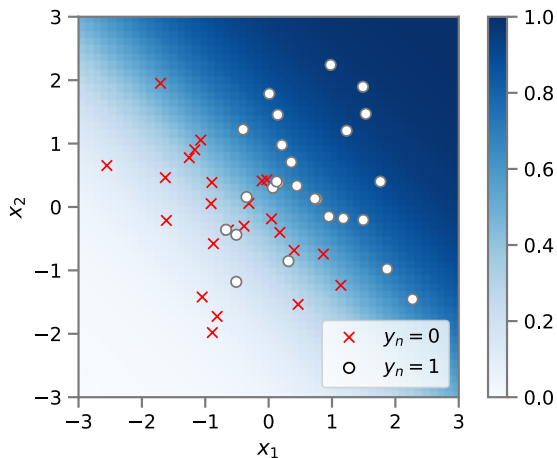
- For example, consider modeling binary observations $y_n \in \{0, 1\}$ given covariates $\mathbf{x}_n \in \mathbb{R}^P$.
- We model the observations as Bernoulli random variables,

$$y_n \mid \mathbf{x}_n \sim \text{Bern}(p_n) \quad (1)$$

where

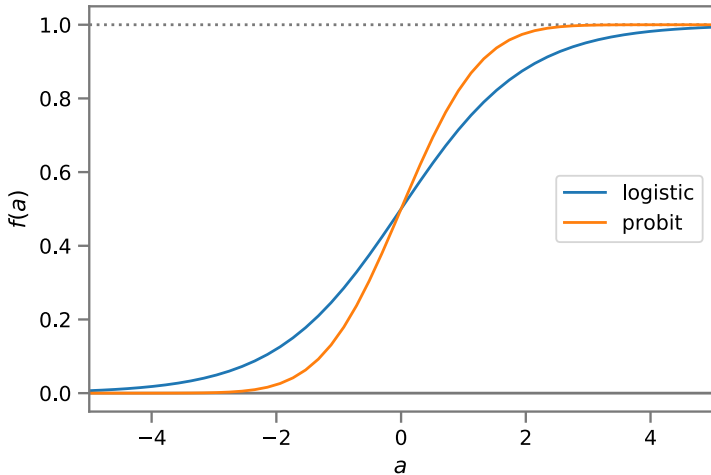
$$p_n = \Pr(y_n = 1 \mid \mathbf{x}_n) = \mathbb{E}[y_n \mid \mathbf{x}_n]. \quad (2)$$

- Then, model the conditional expectation as $\mathbb{E}[y_n \mid \mathbf{x}_n] = f(\mathbf{w}^\top \mathbf{x}_n)$ where $f: \mathbb{R} \rightarrow [0, 1]$ is the *mean function*.



Logistic and probit functions

Common choices for f include the logistic and probit functions, corresponding to the *logistic regression* and *probit regression* models, respectively.



Likelihood

The likelihood of N conditionally independent observations is,

$$p(\{y_n\}_{n=1}^N \mid \{\mathbf{x}_n\}_{n=1}^N, \mathbf{w}) = \prod_{n=1}^N \text{Bern}(y_n \mid f(\mathbf{w}^\top \mathbf{x}_n)) \quad (3)$$

$$= \prod_{n=1}^N f(\mathbf{w}^\top \mathbf{x}_n)^{y_n} (1 - f(\mathbf{w}^\top \mathbf{x}_n))^{1-y_n} \quad (4)$$

For example, consider the *logistic function* $f(a) = \frac{e^a}{1+e^a}$. Then,

$$p(\{y_n\}_{n=1}^N \mid \{\mathbf{x}_n\}_{n=1}^N, \mathbf{w}) = \prod_{n=1}^N \left[\frac{e^{\mathbf{w}^\top \mathbf{x}_n}}{1 + e^{\mathbf{w}^\top \mathbf{x}_n}} \right]^{y_n} \left[\frac{1}{1 + e^{\mathbf{w}^\top \mathbf{x}_n}} \right]^{1-y_n} \quad (5)$$

$$= \prod_{n=1}^N \frac{e^{\langle y_n \mathbf{x}_n, \mathbf{w} \rangle}}{1 + e^{\mathbf{w}^\top \mathbf{x}_n}} \quad (6)$$

The numerator can be written in terms of $\sum_n y_n \mathbf{x}_n$, but the denominator doesn't easily simplify.

Prior

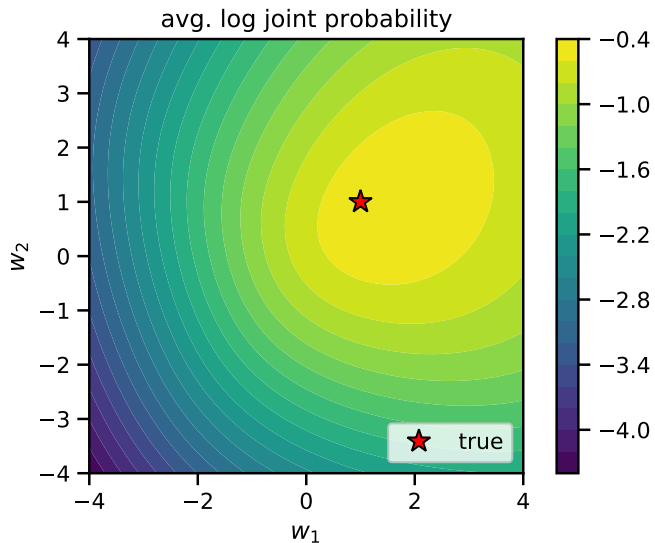
- Here, we'll use a spherical multivariate normal prior,

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} \mid \mathbf{0}, \sigma^2 \mathbf{I}), \quad (7)$$

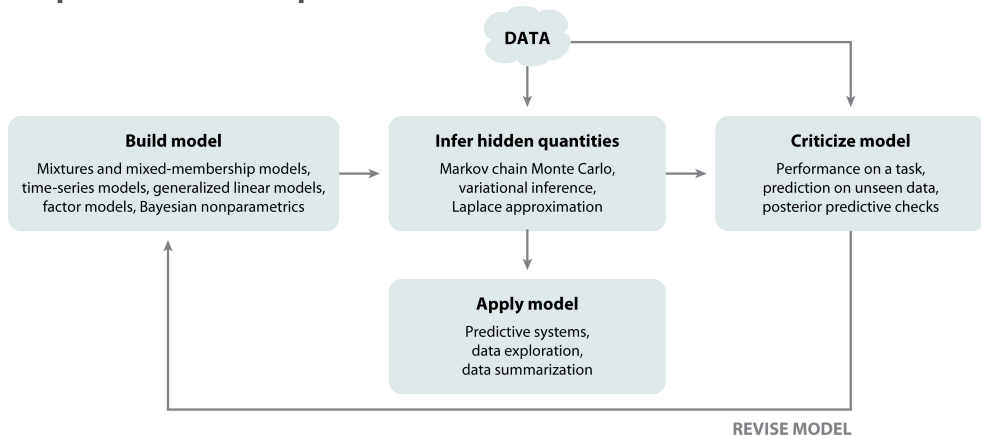
where $\mathbf{0} \in \mathbb{R}^P$ and $\sigma^2 \in \mathbb{R}_+$ are the mean and variance, respectively.

- This is the simplest choice, but we have lots of flexibility if a stronger prior is warranted.
- For example, it's easy to generalize to arbitrary mean and covariance, $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$.

Log probability contours



Box's Loop: Infer hidden quantities



Blei DM. 2014.

Annu. Rev. Stat. Appl. 1:203–32

Blei, *Ann. Rev. Stat. App.* 2014.

Maximum a posteriori (MAP) estimation

Before going to full Bayesian inference, let's just look for the posterior mode.

The log posterior probability is equal to the log joint minus the log marginal likelihood, which is constant with respect to the parameters \mathbf{w} ,

$$\mathcal{L}(\mathbf{w}) = \log p(\{y_n\}_{n=1}^N \mid \{\mathbf{x}_n\}_{n=1}^N, \mathbf{w}) + \log p(\mathbf{w}) - \underbrace{\log p(\{y_n\}_{n=1}^N \mid \{\mathbf{x}_n\}_{n=1}^N)}_{\text{const. w.r.t. } \mathbf{w}} \quad (8)$$

$$= \sum_{n=1}^N \log \frac{e^{\langle y_n \mathbf{x}_n, \mathbf{w} \rangle}}{1 + e^{\mathbf{w}^\top \mathbf{x}_n}} - \frac{1}{2\sigma^2} \mathbf{w}^\top \mathbf{w} + c \quad (9)$$

$$= \left\langle \sum_{n=1}^N y_n \mathbf{x}_n, \mathbf{w} \right\rangle - \sum_{n=1}^N \log(1 + \exp\{\mathbf{x}_n^\top \mathbf{w}\}) - \frac{1}{2\sigma^2} \mathbf{w}^\top \mathbf{w} + c \quad (10)$$

Gradient of the log probability

Taking the gradient and setting it to zero,

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \sum_{n=1}^N y_n \mathbf{x}_n - \sum_{n=1}^N \frac{\exp\{\mathbf{x}_n^\top \mathbf{w}\}}{1 + \exp\{\mathbf{x}_n^\top \mathbf{w}\}} \mathbf{x}_n - \frac{1}{\sigma^2} \mathbf{w} \quad (11)$$

$$= \sum_{n=1}^N y_n \mathbf{x}_n - \sum_{n=1}^N f(\mathbf{x}_n^\top \mathbf{w}) \mathbf{x}_n - \frac{1}{\sigma^2} \mathbf{w} \quad (12)$$

In matrix notation,

$$\mathbf{X} = \begin{bmatrix} - & \mathbf{x}_1^\top & - \\ & \vdots & \\ - & \mathbf{x}_N^\top & - \end{bmatrix}, \quad \text{and} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}, \quad (13)$$

and letting $\hat{\mathbf{y}} = f(\mathbf{X}\mathbf{w})$, where f is applied elementwise, we can write the gradient as,

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = (\mathbf{y} - \hat{\mathbf{y}})^\top \mathbf{X} - \frac{1}{\sigma^2} \mathbf{w}. \quad (14)$$

Gradient of the log probability II

- ▶ In other words, the *gradient of the log likelihood is a weighted sum of the covariates* where the weights are given by the error $y_n - \hat{y}_n$.
- ▶ Note that the gradient is still a *nonlinear function of \mathbf{w}* since $\hat{y} = f(\mathbf{x}_n^\top \mathbf{w})$.
- ▶ The nonlinear equation $\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = 0$ does not have a closed form solution.

Hessian of the log probability

- ▶ The Hessian—i.e. the matrix of second partial derivatives, or the Jacobian of the gradient—is,

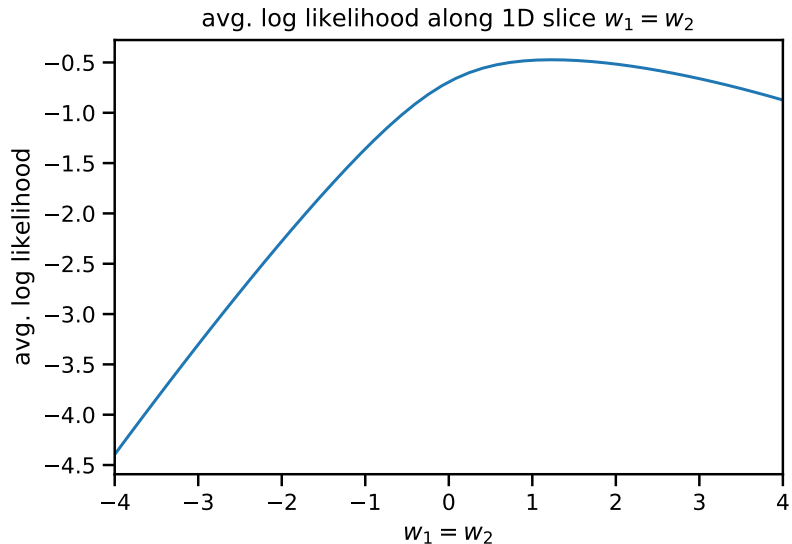
$$\nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{w}) = \nabla_{\mathbf{w}} \sum_{n=1}^N \left(y_n - f(\mathbf{x}_n^\top \mathbf{w}) \right) \mathbf{x}_n - \frac{1}{\sigma^2} \mathbf{w}. \quad (15)$$

$$= - \left(\sum_{n=1}^N f'(\mathbf{x}_n^\top \mathbf{w}) \mathbf{x}_n \mathbf{x}_n^\top + \frac{1}{\sigma^2} \mathbf{I} \right) \quad (16)$$

where $f' : \mathbb{R} \rightarrow \mathbb{R}$ is the derivative of f .

- ▶ Since f is monotonically increasing, $f' \geq 0$.
- ▶ The negative Hessian is a weighted sum of outer products with positive weights $f'(\mathbf{x}_n^\top \mathbf{w})$. Thus, the Hessian is negative definite.
- ▶ That means the *log probability is concave* and has a global optimum.

Concavity of the log probability



Exercise: The derivative of the logistic function

Show that

$$f'(a) = \left[\frac{d}{d\alpha} \frac{e^\alpha}{1 + e^\alpha} \right]_{\alpha=a} = f(a)(1 - f(a)) = f(a)f(-a) \geq 0. \quad (17)$$

Algorithm

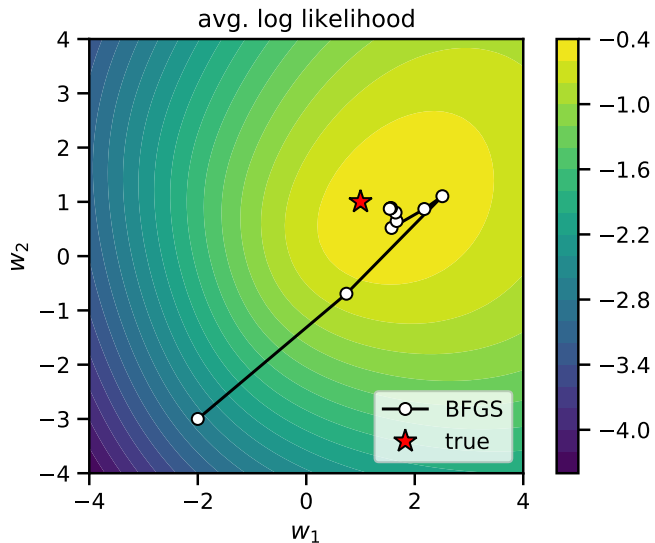
- ▶ Though the MAP estimate doesn't have a closed form solution, we can find it via optimization.
- ▶ Since the log probability is concave, standard optimization methods will find the global optimum.
- ▶ For example, (damped) *Newton's method*,

$$\mathbf{w}^{(i+1)} \leftarrow \mathbf{w}^{(i)} - \gamma [\nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{w}^{(i)})]^{-1} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}^{(i)}) \quad (18)$$

(with sufficiently small step size $\gamma \in (0, 1]$) achieves second-order convergence to the optimum, \mathbf{w}_{MAP} .

- ▶ With many covariates (when P is large), Newton's method incurs an $O(P^3)$ complexity.
- ▶ *Quasi-Newton* methods, like the BFGS algorithm, build an approximation to approximate Newton's method at lower computational cost.

BFGS Convergence



Posterior Distribution

What if we want more than a point estimate?

The posterior distribution is given by,

$$p(\mathbf{w} \mid \{\mathbf{x}_n, y_n\}_{n=1}^N) = \frac{p(\mathbf{w}) p(\{y_n\}_{n=1}^N \mid \mathbf{w}, \{\mathbf{x}_n\}_{n=1}^N)}{p(\{y_n\}_{n=1}^N \mid \{\mathbf{x}_n\}_{n=1}^N)} \quad (19)$$

(20)

Unfortunately, the marginal likelihood is hard to compute since

$$p(\{y_n\}_{n=1}^N \mid \{\mathbf{x}_n\}_{n=1}^N) = \int \prod_{n=1}^N \text{Bern}(y_n \mid f(\mathbf{w}^\top \mathbf{x}_n)) \mathcal{N}(\mathbf{w} \mid 0, \sigma^2 I) d\mathbf{w} \quad (21)$$

does not have a closed form solution.

Laplace Approximation

Idea: *approximate the posterior with a multivariate normal distribution centered on the mode.*

To motivate this, consider a second-order Taylor approximation to the log posterior,

$$\mathcal{L}(\mathbf{w}) \approx \mathcal{L}(\mathbf{w}_{\text{MAP}}) + (\mathbf{w} - \mathbf{w}_{\text{MAP}})^{\top} \underbrace{\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_{\text{MAP}})}_{\mathbf{0} \text{ at the mode}} + \frac{1}{2} (\mathbf{w} - \mathbf{w}_{\text{MAP}})^{\top} \nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{w}_{\text{MAP}}) (\mathbf{w} - \mathbf{w}_{\text{MAP}}) \quad (22)$$

$$= -\frac{1}{2} (\mathbf{w} - \mathbf{w}_{\text{MAP}})^{\top} \Sigma^{-1} (\mathbf{w} - \mathbf{w}_{\text{MAP}}) + c \quad (23)$$

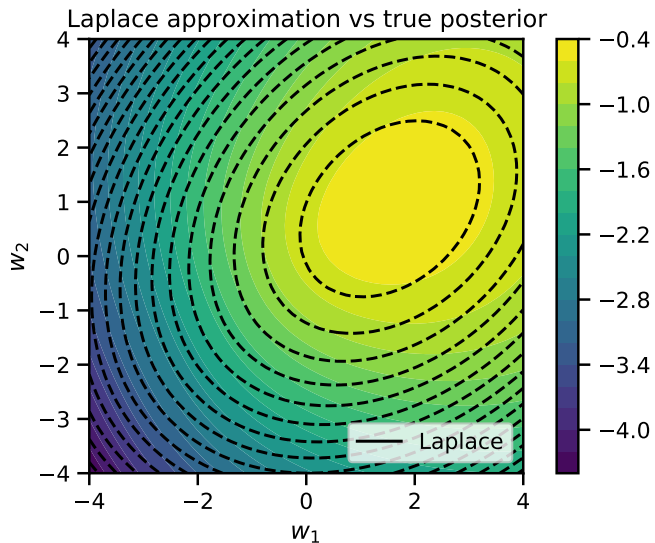
$$= \log \mathcal{N}(\mathbf{w} \mid \mathbf{w}_{\text{MAP}}, \Sigma) \quad (24)$$

where $\Sigma = -[\nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{w}_{\text{MAP}})]^{-1}$

In other words, the posterior is approximately Gaussian with covariance given by the (negative) inverse Hessian at the mode.

Since the Hessian is *negative* definite, the covariance is *positive* definite, as required.

Laplace Approximation II



Bernstein-von Mises Theorem

In the large data limit (as $N \rightarrow \infty$), the posterior is asymptotically normal, justifying the Laplace approximation in this regime.

Consider a simpler setting in which we have data $\{y_n\}_{n=1}^N \stackrel{\text{iid}}{\sim} p(y \mid \theta_{\text{true}})$.

Under some conditions (e.g. θ_{true} not on the boundary of Θ and θ_{true} has nonzero prior probability), then the MAP estimate is consistent. As $N \rightarrow \infty$, $\theta_{\text{MAP}} \rightarrow \theta_{\text{true}}$.

Likewise,

$$p(\theta \mid \{y_n\}_{n=1}^N) \rightarrow \mathcal{N}(\theta \mid \theta_{\text{true}}, \frac{1}{N} [J(\theta_{\text{true}})]^{-1}) \quad (25)$$

where

$$J(\theta) = -\mathbb{E}_{p(y|\theta)} \left[\frac{d^2}{d\theta^2} \log p(y \mid \theta) \right] \quad (26)$$

is the *Fisher information* of parameter θ .

Approximating the marginal likelihood

The Laplace approximation also offers an approximation of the intractable marginal likelihood,

$$\mathcal{L}(\mathbf{w}_{\text{MAP}}) = \log p(\mathbf{w}_{\text{MAP}}) + \log p(\{y_n\}_{n=1}^N \mid \mathbf{w}_{\text{MAP}}, \{\mathbf{x}_n\}_{n=1}^N) - \log p(\{y_n\}_{n=1}^N \mid \{\mathbf{x}_n\}_{n=1}^N) \quad (27)$$

$$\approx \log \mathcal{N}(\mathbf{w}_{\text{MAP}} \mid \mathbf{w}_{\text{MAP}}, \Sigma) \quad (28)$$

$$= -\frac{P}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma| \quad (29)$$

Rearranging terms,

$$\begin{aligned} \log p(\{y_n\}_{n=1}^N \mid \{\mathbf{x}_n\}_{n=1}^N) &\approx \log p(\mathbf{w}_{\text{MAP}}) + \log p(\{y_n\}_{n=1}^N \mid \mathbf{w}_{\text{MAP}}, \{\mathbf{x}_n\}_{n=1}^N) \\ &\quad + \frac{P}{2} \log(2\pi) + \frac{1}{2} \log |\Sigma| \end{aligned} \quad (30)$$

Note: combine this with $\Sigma \approx \frac{1}{N} [J(\mathbf{w}_{\text{MAP}})]^{-1}$ and $\frac{1}{2} \log |\Sigma| \approx \frac{P}{2} \log N + O(1)$ to derive the *Bayesian information criterion (BIC)*, a technique for penalized maximum likelihood estimation.

Exponential family distributions

Now we'll consider GLMs with *exponential family* observations,

$$p(y_n | \eta_n) = h(y_n) \exp \{ \langle t(y_n), \eta_n \rangle - A(\eta_n) \}, \quad (31)$$

where

- ▶ $h(y_n) : \mathcal{Y} \rightarrow \mathbb{R}_+$ is the *base measure*,
- ▶ $t(y_n) \in \mathbb{R}^T$ are the *sufficient statistics*,
- ▶ $\eta_n \in \mathbb{R}^T$ are the *natural parameters*, and
- ▶ $A(\eta_n) : \mathbb{R}^T \rightarrow \mathbb{R}$ is the *log normalizing function* (aka the *partition function*).

Example: Bernoulli distribution in exponential family form

The Bernoulli distribution is a special case,

$$\text{Bern}(y_n | p_n) = p_n^{y_n} (1 - p_n)^{1-y_n} \quad (32)$$

$$= \exp \{ y_n \log p_n + (1 - y_n) \log(1 - p_n) \} \quad (33)$$

$$= \exp \left\{ y_n \log \frac{p_n}{1 - p_n} + \log(1 - p_n) \right\} \quad (34)$$

$$= h(y_n) \exp \{ y_n \eta_n - A(\eta_n) \} \quad (35)$$

where

$$h(y_n) = 1 \quad t(y_n) = y_n \quad \eta_n = \log \frac{p_n}{1 - p_n} \quad A(\eta_n) = -\log(1 - p_n) \quad (36)$$

$$= -\log \left(1 - \frac{e^{\eta_n}}{1 + e^{\eta_n}} \right) \quad (37)$$

$$= \log(1 + e^{\eta_n}). \quad (38)$$

Gradient of the log normalizer yields the expected sufficient statistics

By definition,

$$A(\eta_n) = \log \int h(y_n) \exp \{ \langle t(y_n), \eta_n \rangle \} dy_n \quad (39)$$

so

$$\nabla_{\eta_n} A(\eta_n) = \nabla_{\eta_n} \log \int h(y_n) \exp \{ \langle t(y_n), \eta_n \rangle \} dy_n \quad (40)$$

$$= \frac{\int h(y_n) \exp \{ \langle t(y_n), \eta_n \rangle \} t(y_n) dy_n}{\int h(y_n) \exp \{ \langle t(y_n), \eta_n \rangle \} dy_n} \quad (41)$$

$$= \int p(y_n | \eta_n) t(y_n) dy_n \quad (42)$$

$$= \mathbb{E}_{p(y_n | \eta_n)} [t(y_n)] \quad (43)$$

Hessian of the log normalizer yields the covariance of the sufficient statistics

Likewise,

$$\nabla_{\eta_n}^2 A(\eta_n) = \nabla_{\eta_n} \int p(y_n | \eta_n) t(y_n) dy_n \quad (44)$$

$$= \int p(y_n | \eta_n) t(y_n) (t(y_n) - \nabla_{\eta_n} A(\eta_n))^{\top} dy_n \quad (45)$$

$$= \mathbb{E}_{p(y_n | \eta_n)} [t(y_n) t(y_n)^{\top}] - \mathbb{E}_{p(y_n | \eta_n)} [t(y_n)] \mathbb{E}_{p(y_n | \eta_n)} [t(y_n)]^{\top} \quad (46)$$

$$= \text{Cov}_{p(y_n | \eta_n)} [t(y_n)] \quad (47)$$

Covariances are positive semi-definite, so the *log normalizer is a convex function*.

Exponential family GLMs

- To construct a generalized linear model with exponential family observations, we set

$$\mathbb{E}[y_n | \mathbf{x}_n] = f(\mathbf{w}^\top \mathbf{x}_n). \quad (48)$$

- From above, this implies,

$$\nabla_{\eta_n} A(\eta_n) = f(\mathbf{w}^\top \mathbf{x}_n) \quad (49)$$

$$\Rightarrow \eta_n = [\nabla A]^{-1}(f(\mathbf{w}^\top \mathbf{x}_n)), \quad (50)$$

when $\nabla A(\cdot)$ is invertible. (In this case, the exponential family is said to be *minimal*).

- The *canonical mean function* is $f(\cdot) = \nabla A(\cdot)$ so that $\eta_n = \mathbf{w}^\top \mathbf{x}_n$.
- The (canonical) *link function* is the inverse of the (canonical) mean function.

Logistic regression revisited

Consider the Bernoulli distribution once more. The gradient of the log normalizer is,

$$\nabla_{\eta} A(\eta) = \nabla_{\eta} \log(1 + e^{\eta}) = \frac{e^{\eta}}{1 + e^{\eta}} \quad (51)$$

This is the logistic function!

Thus, logistic regression with is a Bernoulli GLM with the canonical mean function.

Why care about canonical mean functions?

(**Spoiler:** It's 2021 and we have automatic differentiation so it's not such a big deal anymore.)

Canonical mean functions lead to nice math. Consider the log joint probability,

$$\mathcal{L}(\mathbf{w}) = \log p(\mathbf{w}) + \sum_{n=1}^N \langle t(y_n), \eta_n \rangle - A(\eta_n) + c \quad (52)$$

$$= -\frac{1}{2\sigma^2} \mathbf{w}^\top \mathbf{w} + \sum_{n=1}^N \langle t(y_n), \mathbf{w}^\top \mathbf{x}_n \rangle - A(\mathbf{w}^\top \mathbf{x}_n) + c, \quad (53)$$

where we have assumed $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} \mid \mathbf{0}, \sigma^2 \mathbf{I})$ and canonical mean function so $\eta_n = \mathbf{w}^\top \mathbf{x}_n$.

The gradient is,

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = -\frac{1}{\sigma^2} \mathbf{w} + \sum_{n=1}^N \langle t(y_n), \mathbf{x}_n \rangle - \langle \nabla A(\mathbf{w}^\top \mathbf{x}_n), \mathbf{x}_n \rangle \quad (54)$$

$$= -\frac{1}{\sigma^2} \mathbf{w} + \sum_{n=1}^N \langle t(y_n) - \mathbb{E}_{p(y_n | \mathbf{w}^\top \mathbf{x}_n)}[t(y_n)], \mathbf{x}_n \rangle \quad (55)$$

Why care about canonical mean functions? II

In many cases, $t(y_n) = y_n \in \mathbb{R}$ so

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = -\frac{1}{\sigma^2} \mathbf{w} + \sum_{n=1}^N (y_n - \hat{y}_n) \mathbf{x}_n. \quad (56)$$

And in that case the Hessian is

$$\nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{w}) = -\frac{1}{\sigma^2} I - \sum_{n=1}^N \nabla^2 A(\mathbf{w}^\top \mathbf{x}_n) \mathbf{x}_n \mathbf{x}_n^\top \quad (57)$$

$$= -\left(\frac{1}{\sigma^2} I + \sum_{n=1}^N \text{Var}_{p(y_n | \mathbf{w}^\top \mathbf{x}_n)}[y_n] \mathbf{x}_n \mathbf{x}_n^\top \right) \quad (58)$$

Why care about canonical mean functions? III

Now recall the Newton's method updates, written here in terms of the change in weights,

$$\Delta \mathbf{w} = -[\nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{w})]^{-1} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) \quad (59)$$

$$= \left[\frac{1}{\sigma^2} \mathbf{I} + \sum_{n=1}^N \text{Var}_{p(y_n | \mathbf{w}^\top \mathbf{x}_n)}[y_n] \mathbf{x}_n \mathbf{x}_n^\top \right]^{-1} \left[-\frac{1}{\sigma^2} \mathbf{w} + \sum_{n=1}^N (y_n - \hat{y}_n) \mathbf{x}_n \right] \quad (60)$$

Letting $\hat{v}_n = \text{Var}_{p(y_n | \mathbf{w}^\top \mathbf{x}_n)}[y_n]$ and taking the uninformative limit of $\sigma^2 \rightarrow \infty$,

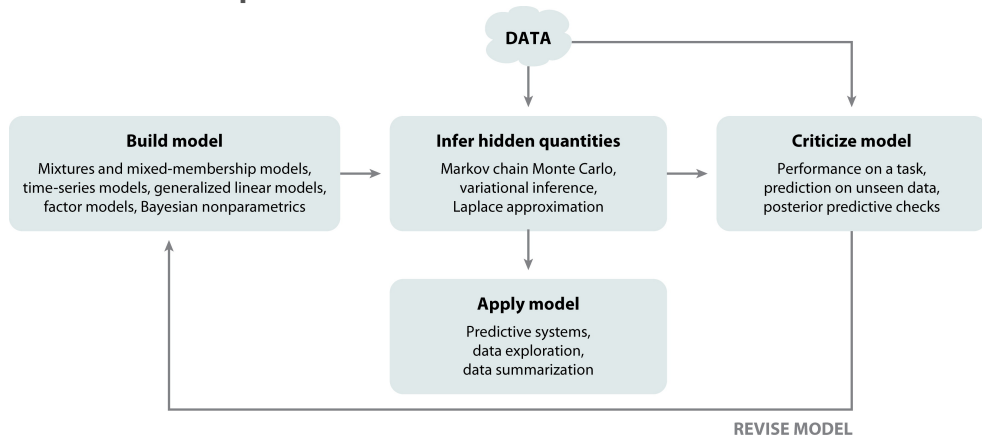
$$\Delta \mathbf{w} = \left[\sum_{n=1}^N \hat{v}_n \mathbf{x}_n \mathbf{x}_n^\top \right]^{-1} \left[\sum_{n=1}^N (y_n - \hat{y}_n) \mathbf{x}_n \right] \quad (61)$$

$$= (\mathbf{X}^\top \hat{\mathbf{V}} \mathbf{X})^{-1} [\mathbf{X}^\top \hat{\mathbf{V}} \hat{\mathbf{z}}] \quad (62)$$

where $\hat{\mathbf{V}} = \text{diag}([\hat{v}_1, \dots, \hat{v}_N])$ and $\hat{\mathbf{z}} = \hat{\mathbf{V}}^{-1}(\mathbf{y} - \hat{\mathbf{y}})$.

This is *iteratively reweighted least squares* (IRLS) with weights \hat{v}_n and targets $\hat{z}_n = \frac{y_n - \hat{y}_n}{\hat{v}_n}$, both of which are functions of the current weights \mathbf{w} .

Model evaluation via predictive likelihood



Blei DM. 2014.

Annu. Rev. Stat. Appl. 1:203–32

Blei, *Ann. Rev. Stat. App.* 2014.

Measures of predictive accuracy

One way to evaluate a model is through the accuracy of its predictions.

For example, given a *point estimate* $\hat{\theta}$, the *mean squared error* summarizes predictive accuracy as,

$$\frac{1}{N_{\text{test}}} \sum_{n=1}^{N_{\text{test}}} (y_n - \mathbb{E}[y_n \mid \hat{\theta}])^2 \quad (63)$$

In *probabilistic prediction* (aka *probabilistic forecasting*), we report a distribution over parameters θ . In Bayesian data analysis, this is the posterior distribution $p(\theta \mid \{y_n\}_{n=1}^N)$.

The *posterior predictive distribution* averages over the posterior to make a prediction,

$$p(y_{n'} \mid \{y_n\}_{n=1}^N) = \int p(y_{n'} \mid \theta) p(\theta \mid \{y_n\}_{n=1}^N) d\theta, \quad (64)$$

weighting parameters by their posterior probability.

Scoring rules

A *scoring rule* measures the accuracy of a probabilistic prediction.

The *logarithmic scoring rule* is unique in that it is both *proper* and *local*.

The *log predictive likelihood* of new datapoint $y_{n'}$ is,

$$\log p(y_{n'} \mid \{y_n\}_{n=1}^N) = \log \int p(y_{n'} \mid \theta) p(\theta \mid \{y_n\}_{n=1}^N) d\theta. \quad (65)$$

Averaging over the distribution of future data

Ideally, we would measure the expected out-of-sample predictive performance,

$$\mathbb{E}_{f(y)}[\log p(y \mid \{y_n\}_{n=1}^N)] = \int \log p(y \mid \{y_n\}_{n=1}^N) f(y) dy \quad (66)$$

where $f(y)$ is the true, but unknown data-generating distribution.

BDA3 calls this the *expected out-of-sample log predictive density* (elpd).

But we don't know f ! Instead, we approximate it with,

$$\frac{1}{N_{\text{test}}} \sum_{n'=1}^{N_{\text{test}}} \log p(y_{n'} \mid \{y_n\}_{n=1}^N) \quad (67)$$

for some *test dataset* $\{y_{n'}\}_{n'=1}^{N_{\text{test}}}$.

Estimating out-of-sample predictive accuracy with available data

- ▶ Naively, the *within-sample* predictive accuracy takes the test set to be the training set. As such, it over-estimates the the predictive accuracy.
- ▶ *Information criteria* like the AIC, DIC, and WAIC try to correct for this bias by penalizing the within-sample estimate.
- ▶ *Cross-validation* splits the data into train and test to approximate the predictive likelihood.
- ▶ Hardcore Bayesians don't like to withhold data when computing the posterior. Nevertheless, cross-validated predictive log likelihood is the current standard.

Akaike Information Criterion (AIC)

The *Akaike information criterion* (AIC) is based on an asymptotic normal approximation to the posterior.

It corrects for the bias of the within-sample log likelihood by accounting for the number of parameters (an approximation to the amount of overfitting),

$$-\frac{1}{2}\text{AIC} = \sum_{n=1}^N \log p(y_n \mid \theta_{\text{MLE}}) - K, \quad (68)$$

where K is the number of parameters.

The AIC makes some sense for linear models with uniform priors, but not so much for more complicated models.

Deviance Information Criterion (DIC)

The *deviance information criterion* (DIC) makes two changes,

1. Replace the MLE with the posterior mean, $\bar{\theta} = \mathbb{E}_{p(\theta|y)}[\theta]$.
2. Replace K with a bias correction based on the data.

Then,

$$-\frac{1}{2}\text{DIC} = \sum_{n=1}^N \log p(y_n | \bar{\theta}) - K_{\text{DIC}} \quad (69)$$

where

$$K_{\text{DIC}} = 2 \left(\log \sum_{n=1}^N p(y_n | \bar{\theta}) - \mathbb{E}_{p(\theta|y)} \left[\log \sum_{n=1}^N p(y_n | \theta) \right] \right) \quad (70)$$

An alternative definition uses,

$$K_{\text{DIC}} = 2 \text{Var}_{p(\theta|y)} \left[\log \sum_{n=1}^N p(y_n | \bar{\theta}) \right] \quad (71)$$

Watanabe-Akaike Information Criterion (WAIC)

The *Watanabe-Akaike information criterion* (WAIC) is similar to the DIC with,

$$K_{\text{WAIC}} = 2 \left(\log \sum_{n=1}^N \mathbb{E}_{p(\theta|y)} [p(y_n | \theta)] - \mathbb{E}_{p(\theta|y)} \left[\log \sum_{n=1}^N p(y_n | \theta) \right] \right) \quad (72)$$

An alternative definition uses,

$$K_{\text{WAIC}} = \sum_{n=1}^N \text{Var}_{p(\theta|y)} [\log p(y_n | \bar{\theta})], \quad (73)$$

Then, the WAIC is,

$$-\frac{1}{2} \text{WAIC} = \sum_{n=1}^N \log \int p(y_n | \theta) p(\theta | \{y_n\}_{n=1}^N) d\theta - K_{\text{WAIC}} \quad (74)$$

Leave-one-out cross validation

Idea: approximate out-of-sample predictive accuracy by randomly splitting the training data.

Leave-one-out cross validation (LOOCV) withholds one datapoint out at a time, computes the posterior given the other $N - 1$, and then evaluates predictive likelihood on the held out datapoint.

$$\mathbb{E}_{f(y)}[\log p(y \mid \{y_n\}_{n=1}^N)] \approx \frac{1}{N} \sum_{n=1}^N \log p(y_n \mid \{y_m\}_{m \neq n}). \quad (75)$$

For small N (or when using a small number of folds), a bias-correction can be used. (See pg. 175-176 of the book.)

Cross-validated predictive log likelihood estimates are similar to the *jackknife* resampling method in the sense that it is estimating an expectation wrt the unknown data-generating distribution f by resampling the given data.

Thus, it measures the frequentist properties of a Bayesian procedure.