

# SPR Day 11

Goals: Review gradient + Hessian  
for linear & logistic regr.

Review 1<sup>st</sup> and 2<sup>nd</sup> order  
gradient descent methods

Show both LinRegr and LogisticRegr are CONVEX

Setup Bayesian Logistic Regression

Joint Model

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

↑ normal prior                      ↑ Bernoulli likelihood

Peek at Laplace Approximation

1<sup>st</sup> and 2<sup>nd</sup> order derivatives  
of Maximum Lik. objective for both Linear & Logistic Regr.

$\Phi$ :  $N \times M$  feature matrix

$t$ :  $N \times 1$  output vector

$$d = \sum_n \log p(t_n | x_n, w)$$

$$\hat{y} = \Phi w$$

predictions

actual labels

Linear Regr

gradient  $\nabla_w d = \Phi^T (\Phi w - t)$

Hessian  $\nabla_w \nabla_w d = \Phi^T \Phi$

always constant  
wrt to  $w$

will be positive definite

thus invertible if  $N \geq M$   
and full rank

Logistic Regression

gradient  $\nabla_w d = \Phi^T (\sigma(\Phi w) - t)$

apply sigmoid  
elementwise to vector  $\Phi w$

scalars vectors

Hessian  $\nabla_w \nabla_w d = \Phi^T R(w) \Phi = \sum_n \hat{y}_n(w) (1 - \hat{y}_n(w)) \phi_n \phi_n^T$

Let  $R(w)$  be  $N \times N$  diagonal matrix where entry  $n, n$   
is  $[R(w)]_{nn} = \hat{y}(w, x_n) (1 - \hat{y}(w, x_n))$

$$= \sigma(w^T \phi(x_n)) \sigma(-w^T \phi(x_n))$$

$\Phi^T \Phi$  is invertible

always  $> 0$

Because  $R$  has positive diagonal  
Log Reg. Hessian will be invertible if  
and pos. definite

# Gradients + Hessians of Linear + Logistic Regr.

$$\mathcal{L}(w) = - \sum_{n=1}^N \log p(t_n | x_n, w)$$

## Linear Regr.

gradient  $\nabla_w \mathcal{L} = \Phi^T (\Phi w - t)$  Mx1 vector

Hessian  $\nabla_w \nabla_w \mathcal{L} = \Phi^T \Phi$  MxM matrix

## Logistic Regr.

gradient  $\nabla_w \mathcal{L} = \Phi^T (\underbrace{\sigma(\Phi w)}_{\text{predicted mean}} - t)$

Hessian  $\nabla_w \nabla_w \mathcal{L} = \Phi^T R \Phi$  Bishop 4.98

↑ N x N diagonal matrix

each diagonal entry =  $\begin{bmatrix} \sigma(w^T \phi(x_1)) \sigma(-w^T \phi(x_1)) & & \\ & \ddots & \\ \sigma(w^T \phi(x_N)) \sigma(-w^T \phi(x_N)) & & \end{bmatrix}$

=  $y_n(1-y_n)$

Question: Is this objective function convex?  
Does it have unique global minima? ↔ same!

Punchline:  
Logistic Regression  
ML estimation  
is convex

Function is convex if Hessian is positive definite for all w.

Can show for logistic regression:  $u^T H u = (u^T \Phi^T R) (R^T \Phi u) = v^T v \geq 0$  thus, is / p.d.

# From 1<sup>st</sup> order to 2<sup>nd</sup> order gradient methods

For any loss  $\alpha(w)$  let  $g(w) = \nabla_w \alpha$  and  $H(w) = \nabla_w \nabla_w \alpha$

First order update:

$$w^{\text{new}} \leftarrow w^{\text{old}} - \epsilon g(w^{\text{old}}) \quad \begin{array}{l} \text{step size} \\ \epsilon > 0 \end{array}$$

Second order update (Newton-Raphson)

$$w^{\text{new}} \leftarrow w^{\text{old}} - H^{-1}(w^{\text{old}}) g(w^{\text{old}})$$

When loss is quadratic, Newton update gives optimal  $w^*$  in one step

Example: Linear Regression (use defn of  $g$  &  $H$  from prev. page)

$$\begin{aligned} w^{\text{new}} &\leftarrow w^{\text{old}} - H^{-1} g \\ &\leftarrow w^{\text{old}} - (\Phi^T \Phi)^{-1} (\Phi^T \Phi w^{\text{old}} - \Phi^T t) \\ &\leftarrow w^{\text{old}} - w^{\text{old}} + (\Phi^T \Phi)^{-1} \Phi^T t \\ &\leftarrow (\Phi^T \Phi)^{-1} \Phi^T t \quad \begin{array}{l} \text{"least squares"} \\ \text{optimal solution} \\ \text{(maximizes likelihood)} \end{array} \end{aligned}$$

For logistic regression, 2<sup>nd</sup> order methods are gold standard (but still require many iterations)

$$\begin{aligned} w^{\text{new}} &\leftarrow w^{\text{old}} - (\Phi^T R \Phi)^{-1} \Phi^T (\hat{y} - t) \quad \text{remember } R(w) \text{ and } \hat{y}(w) \text{ depend on } \underline{w} \\ &\leftarrow (\Phi^T R \Phi)^{-1} (\Phi^T R \Phi w^{\text{old}} - \Phi^T \hat{y} + \Phi^T t) \\ &\leftarrow (\Phi^T R \Phi)^{-1} \Phi^T R z \quad \text{where } z = \Phi w^{\text{old}} - R^{-1}(\hat{y} - t) \\ &\text{looks like least squares with weights } R_{nn} = \hat{y}_n(1 - \hat{y}_n) \end{aligned}$$

Algorithm: 2<sup>nd</sup> order GD using Hessian (aka Newton method) is known as "Iteratively Reweighted Least Squares" (IRLS)

↪ update  $\Gamma_n$  given  $w$  ↪ update  $w$  given  $\Gamma_n$  ↪ until converged!

# Bayesian Logistic Regression

Model:

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

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

Just like linear regr.

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

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

iid across examples

Goals are to estimate the posterior and predictive

Posterior:  $p(w \mid \{x_n, t_n\}_{n=1}^N)$  no closed form!  
Not a Gaussian!

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

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

likelihood

posterior

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.

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

Here,  $f(w) > 0$  is known and evaluable 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:  $q(w) = \mathcal{N}(m, \beta^{-1})$   
- pick mean to match the mode 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 2<sup>nd</sup>-order Taylor approximation to  $p(w)$  at the mode  $w=m$

$$\begin{aligned} \beta &= \left. \frac{\partial}{\partial w} \frac{\partial}{\partial w} \left[ -\log f(w) \right] \right|_{w=m} \\ &= -l''(m) \quad \text{where } l = \log f(w) \end{aligned}$$

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  $m = \operatorname{argmax}_w p(w)$

$$\log p(w) = \log f(w) + \text{const w.r.t } w$$

by definition  
of  $p(w)$

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

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

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

2<sup>nd</sup> order Taylor  
approx. to func.  $l$   
at  $w=m$

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

$l(m)$  is const  
w.r.t  $w$ ,  
so group w/ const

$$l'(m) = 0 \text{ bc.}$$

this is a Gaussian pdf

with mean  $m = \operatorname{argmax}_w l(w)$

and precision  $\beta = -l''(m)$

$m$  is a maximizer  
of  $f(w)$  &  $l(w)$   
so this term  
cancels

# Approximating the Posterior for Logistic Regression

True posterior intractable:

$$P(w | \{x_n, t_n\}_{n=1}^N) = \frac{1}{Z} \cdot \prod_{n=1}^N \text{Bern}(t_n | \sigma(w^T \phi(x_n)))$$

• Normal( $w | m_0, S_0$ )

But we can apply Laplace Approx!

$$P(w | \{x_n, t_n\}_{n=1}^N) \approx \mathcal{N}(w | m_{\text{MAP}}, S^{-1})$$

where  $m_{\text{MAP}} = \underset{w \in \mathbb{R}^M}{\text{argmax}} \log P(w | \{x_n, t_n\}_{n=1}^N)$

solved w/ Gradient descent

and precision  $S^{-1} = -l''(m_{\text{MAP}})$

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

where  $l'' = \nabla^2 \log p(w|x,t)$   
= Hessian of log likelihood  
+ Hessian of log prior

Hessian of log prior =  $S_0^{-1}$

Hessian of log likelihood =  $\Phi^T R(m_{\text{MAP}}) \Phi$

$$= \sum_{n=1}^N r(m_{\text{MAP}}, x_n) \phi(x_n) \phi(x_n)^T$$

$$\uparrow = \sigma(m_{\text{MAP}}^T \phi_n) (1 - \sigma(m_{\text{MAP}}^T \phi_n))$$