

# Logistic/Softmax Classification, Laplace's Approximation, and Exponential Family

CS772A: Probabilistic Machine Learning

Piyush Rai

# Plan today

- Logistic and Softmax Classification, Generalized Linear Models
- Laplace's approximation: A method to approximate posterior for non-conjugate cases
- Exponential family



# Logistic Regression

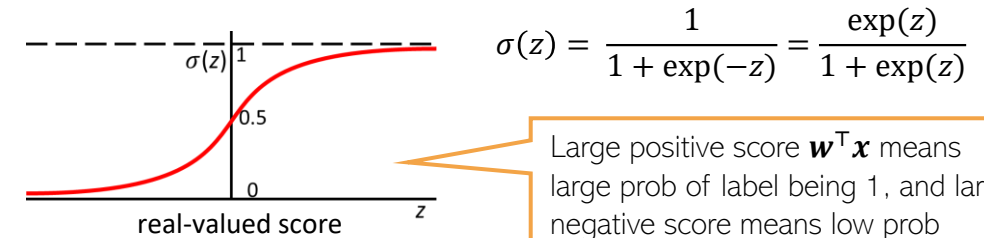
There are other ways too that can convert the score into a probability, such as a CDF:  $p(y = 1|\mathbf{x}, \mathbf{w}) = \mu = \Phi(\mathbf{w}^\top \mathbf{x})$  where  $\Phi$  is the CDF of  $\mathcal{N}(0,1)$ . This model is known as "Probit Regression".



- A discriminative model for binary classification ( $y \in \{0,1\}$ )
- A linear model with parameters  $\mathbf{w} \in \mathbb{R}^D$  computes a **score**  $\mathbf{w}^\top \mathbf{x}$  for input  $\mathbf{x}$
- A **sigmoid** function maps this real-valued score into probability of label being 1

Also used as a nonlinear "activation function" in deep neural networks

$$p(y = 1|\mathbf{x}, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x})$$



- Thus conditional distribution of label  $y \in \{0,1\}$  given  $\mathbf{x}$  is the following Bernoulli

Likelihood

$$p(y|\mathbf{x}, \mathbf{w}) = \text{Bernoulli}[y|\mu] = \mu^y (1 - \mu)^{1-y} = \left[ \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^y \left[ \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^{1-y}$$

- NLL is the **binary cross-entropy loss**:  $-[y_n \log \mu_n + (1 - y_n) \log (1 - \mu_n)]$
- NLL is convex in  $\mathbf{w}$ . Can also use a prior  $p(\mathbf{w}|\lambda) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1} \mathbf{I})$  if interested in MAP or full posterior on  $\mathbf{w}$



# Multiclass Logistic (a.k.a. Softmax) Regression

- Also called **multinoulli/multinomial regression**: Basically, LR for  $K > 2$  classes
- In this case,  $y_n \in \{1, 2, \dots, K\}$  and label probabilities are defined as

$$p(y_n = k | \mathbf{x}_n, \mathbf{W}) = \frac{\exp(\mathbf{w}_k^\top \mathbf{x}_n)}{\sum_{\ell=1}^K \exp(\mathbf{w}_\ell^\top \mathbf{x}_n)} = \mu_{nk}$$

Softmax function

Real-valued scores  $\mathbf{w}_k^\top \mathbf{x}_n$  are also known as “logits” (thus  $K$  logits for each input)

Also note that  $\sum_{\ell=1}^K \mu_{n\ell} = 1$  for any input  $\mathbf{x}_n$



- $K$  weight vecs  $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K$  (one per class), each  $D$ -dim, and  $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K]$
- Each likelihood  $p(y_n | \mathbf{x}_n, \mathbf{W})$  is a **multinoulli** distribution. Therefore total likelihood

$$p(\mathbf{y} | \mathbf{X}, \mathbf{W}) = \prod_{n=1}^N \prod_{\ell=1}^K \mu_{n\ell}^{y_{n\ell}}$$

Notation:  $y_{n\ell} = 1$  if true class of  $\mathbf{x}_n$  is  $\ell$  and  $y_{n\ell'} = 0 \forall \ell' \neq \ell$



# Generalized Linear Models

- (Probabilistic) Linear Regression: when response  $y$  is real-valued

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

- Logistic Regression: when response  $y$  is binary (0/1)

$$p(y|\mathbf{x}, \mathbf{w}) = \text{Bernoulli}[y|\sigma(\mathbf{w}^\top \mathbf{x})] = \left[ \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^y \left[ \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^{1-y}$$

- Both are examples of a Generalized Linear Model (GLM)

- The model depends on the inputs  $\mathbf{x}$  via a linear model  $\mathbf{w}^\top \mathbf{x}$

- GLM is defined using an exponential family distribution

$$p(y|\mathbf{x}, \mathbf{w}) = \text{ExpFam}[y|f(\mathbf{w}^\top \mathbf{x})]$$

MLE/MAP of  $\mathbf{w}$  is easy for GLMs (due to convex objective, thanks to exp-family). Posterior usually requires approximations if likelihood and prior are not conjugate pairs (Laplace approximation or other methods used)

- ExpFam can be any suitable distribution depending on the nature of outputs, e.g.,

- Gaussian for reals, Bernoulli for binary, Poisson for Count, gamma for positive reals

- ExpFam distributions are more generally useful in other contexts as well



# Logistic Regression: MAP and Posterior

- The posterior will be

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{w})p(\mathbf{y}|\mathbf{X}, \mathbf{w})}{p(\mathbf{y}|\mathbf{X})} = \frac{\overset{\text{Gaussian}}{p(\mathbf{w})} \prod_{n=1}^N \overset{\text{Bernoulli}}{p(y_n|\mathbf{w}, \mathbf{x}_n)}}{\int p(\mathbf{w}) \prod_{n=1}^N p(y_n|\mathbf{w}, \mathbf{x}_n) d\mathbf{w}}$$

- MAP estimation is easy.  $-\log p(\mathbf{w}|\mathbf{X}, \mathbf{y})$  is convex for LR. Unique minima
  - Can use first or second order optimization with gradient and Hessian being

$$\begin{aligned} \mathbf{g} &= -\sum_{n=1}^N (y_n - \mu_n) \mathbf{x}_n + \lambda \mathbf{w} = \mathbf{X}^\top (\boldsymbol{\mu} - \mathbf{y}) + \lambda \mathbf{w} \quad (\text{a } D \times 1 \text{ vector}) \\ \mathbf{H} &= \sum_{n=1}^N \mu_n (1 - \mu_n) \mathbf{x}_n \mathbf{x}_n^\top + \lambda \mathbf{I} = \mathbf{X}^\top \mathbf{S} \mathbf{X} + \lambda \mathbf{I} \quad (\text{a } D \times D \text{ matrix}) \\ \mu_n &= \sigma(\mathbf{w}^\top \mathbf{x}_n) \end{aligned}$$

- Full posterior is intractable because of non-conjugacy
  - A popular option is to use the [Laplace's approximation](#) (other methods like MCMC and variational inference can also be used; will see them later)

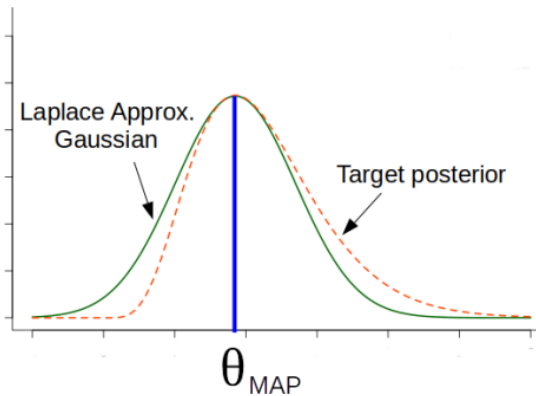


# Laplace's Approximation

- Consider a posterior distribution that is intractable to compute

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}, \theta)}{p(\mathcal{D})} = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})}$$

- Laplace approximation approximates the above using a **Gaussian** distribution



$$p(\theta|\mathcal{D}) \approx \mathcal{N}(\theta|\theta_{MAP}, \Lambda^{-1})$$

$$\theta_{MAP} = \operatorname{argmax}_{\theta} \log p(\theta|\mathcal{D})$$

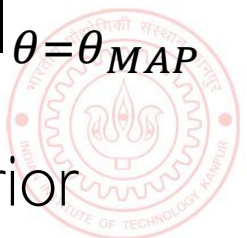
$$\Lambda = -\nabla_{\theta}^2 \log p(\theta|\mathcal{D}) \Big|_{\theta=\theta_{MAP}} = -\nabla_{\theta}^2 \log p(\mathcal{D}, \theta) \Big|_{\theta=\theta_{MAP}}$$

Tells us about the space (curvature) of the true posterior around  $\theta_{MAP}$

Related to the Fisher Information Matrix (FIM); will see shortly

Negative of the Hessian, i.e., the second derivative of the log joint, at  $\theta_{MAP}$

- Laplace's approx. is based on a second-order Taylor approx. of the posterior



# Derivation of the Laplace's Approximation

- Let's write the Bayes rule as

$$p(\mathcal{D}) \approx \exp(\log p(\mathcal{D}, \theta_{MAP})) \times (2\pi)^{D/2} \det(\Lambda)^{1/2}$$

We also get a Laplace approximation of the marginal likelihood (for free!)

Note: Sometimes marginal likelihood is also called model evidence

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}, \theta)}{p(\mathcal{D})} = \frac{p(\mathcal{D}, \theta)}{\int p(\mathcal{D}, \theta) d\theta} = \frac{\exp[\log p(\mathcal{D}, \theta)]}{\int \exp[\log p(\mathcal{D}, \theta)] d\theta}$$

- Consider second-order Taylor approximation of a function  $f(\theta)$  around some  $\theta_0$

$$f(\theta) \approx f(\theta_0) + (\theta - \theta_0)^\top \nabla_{\theta} f(\theta_0) + \frac{1}{2} (\theta - \theta_0)^\top \nabla_{\theta}^2 f(\theta_0) (\theta - \theta_0)$$

- Assuming  $f(\theta) = \log p(\mathcal{D}, \theta)$  and  $\theta_0 = \theta_{MAP}$

Constant w.r.t.  $\theta$

Same as  $\nabla^2 \log p(\theta_{MAP}|\mathcal{D})$

$$\log p(\mathcal{D}, \theta) \approx \log p(\mathcal{D}, \theta_{MAP}) + \frac{1}{2} (\theta - \theta_{MAP})^\top \nabla_{\theta}^2 \log p(\mathcal{D}, \theta_{MAP}) (\theta - \theta_{MAP})$$

$$\begin{aligned} p(\theta|\mathcal{D}) &\propto \exp \left[ -\frac{1}{2} (\theta - \theta_{MAP})^\top (-\nabla_{\theta}^2 \log p(\mathcal{D}, \theta_{MAP})) (\theta - \theta_{MAP}) \right] \\ &= \mathcal{N}(\theta | \theta_{MAP}, \Lambda^{-1}) \quad (\text{where } \Lambda = -\nabla_{\theta}^2 \log p(\mathcal{D}, \theta_{MAP}) = -\mathbf{H}) \end{aligned}$$



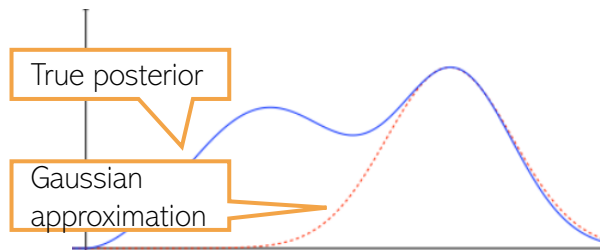


# Properties of Laplace's Approximation

- Straightforward if posterior's derivatives (first/second) can be computed easily
- Expensive if parameter  $\theta$  is very high dimensional
  - Reason: We need to compute and invert Hessian of size  $D \times D$  ( $D$  is the # of params)

E.g., a deep neural network, or even in simpler models (e.g., logistic reg with a very large number of features)

- Can do badly if the (true) posterior is multimodal



For multimodal posteriors, can use a mixture of Laplace approximations\*

Useful for deep learning models

If  $K$  local modes, then define the approx. posterior as a mixture of  $K$  Gaussians

$$p(\theta|D) \approx \sum_{k=1}^K \pi^{(k)} \mathcal{N}(\theta | \theta_{MAP}^{(k)}, H^{(k)-1})$$

(see paper cited below for details)



- Used only when  $\theta$  is a real-valued vector (because of Gaussian approximation)
- Note: Even if we have a non-probabilistic model (loss function + regularization), we can obtain an approx “posterior” for that model using the Laplace's approximation
  - Optima of the regularized loss function will be Gaussian's mean
  - Inverse of the second derivative of the regularized loss function will be covariance matrix



# Detour: Hessian and Fisher Information Matrix

- Hessian is related to the Fisher Information Matrix (FIM)
- Gradient of the log likelihood is also called score function:  $s(\theta) = \nabla_{\theta} \log p(y|\theta)$ 
  - Note: At some places (some generative models)  $\nabla_y \log p(y|\theta)$  also called score function
- Expectation of score function is zero:  $\mathbb{E}_{p(y|\theta)}[s(\theta)] = 0$  (exercise)
- **Fisher Information Matrix (FIM)** is covariance matrix of score function
 
$$\mathbf{F} = \mathbb{E}_{p(y|\theta)}[(s(\theta) - 0)(s(\theta) - 0)^{\top}] = \mathbb{E}_{p(y|\theta)}[\nabla_{\theta} \log p(y|\theta) \nabla_{\theta} \log p(y|\theta)^{\top}]$$

Note: If we have a prior  $p(\theta)$  too, then also add the second derivative of  $\log p(\theta)$
- $\mathbf{F} = - \mathbb{E}_{p(y|\theta)} [\nabla_{\theta}^2 \log p(y|\theta)]$ , i.e., negative of expected **Hessian** (exercise)
- Each entry  $F_{ij}$  tells us how “sensitive” the model is w.r.t. the pair  $(\theta_i, \theta_j)$ 
  - Each diagonal entry  $F_{ii} = (\nabla_{\theta_i} \log p(y|\theta))^2$  tells “important”  $\theta_i$  is by itself
- Can compute **empirical FIM** using data:  $\hat{\mathbf{F}} = \frac{1}{N} \sum_{n=1}^N [\nabla_{\theta} \log p(y_n|\theta) \nabla_{\theta} \log p(y_n|\theta)^{\top}]$



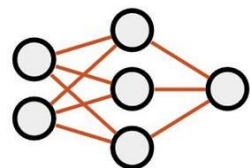
# Laplace Approx. for High-Dimensional Problems

- For high-dim  $\theta$ , Laplace's approx  $p(\theta|\mathcal{D}) \approx \mathcal{N}(\theta|\theta_{MAP}, \Lambda^{-1})$  can be expensive
- Many methods to address this, e.g.,
  - Use a diagonal of (empirical) Fisher as the precision

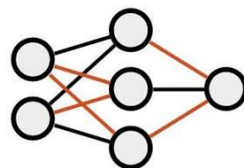
$$\Lambda \approx \text{diag}(\mathbf{F})$$

Diagonal approximation assumes that the weights are all independent whereas block-diagonal assumes that the weights within each block may have correlations

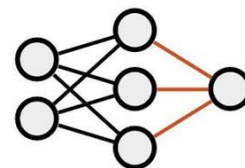
- Use a **block-diagonal** approximation\* of  $\Lambda$  (better than diagonal approx)
- For deep nets, use LA only for some weights + point estimates for others
  - Option 1: Use LA only for last layer weights - “last layer Laplace’s approximation” (LLLA)
  - Option 2: Use LA for weights from an identified “subnetwork”



(a) All



(b) Subnetwork



(c) Last-Layer

- See the “[Laplace Redux](#)” paper for more options and discussion on scalability of LA



# PPD when using Laplace's Approximation

- The PPD when using the Laplace's approximation of the posterior

$$p(y_*|\mathbf{x}_*, \mathcal{D}) = \int p(y_*|\mathbf{x}_*, \theta) p(\theta|\mathcal{D}) d\theta$$

This PPD is an approximation because we are using an approximation of the posterior

$$\approx \int p(y_*|\mathbf{x}_*, \theta) \mathcal{N}(\theta|\theta_{MAP}, \Lambda^{-1}) d\theta$$

- PPD may be intractable depending on the form of  $p(y_*|\mathbf{x}_*, \theta) = p(y_*|f(\mathbf{x}_*, \theta))$
- We can use further approximations if the integral is intractable. Two options:
  - Generate  $M$  samples  $\{\theta^{(i)}\}_{i=1}^M$  from  $\mathcal{N}(\theta|\theta_{MAP}, \Lambda^{-1})$  and compute a Monte Carlo approx.

$$\int p(y_*|\mathbf{x}_*, \theta) \mathcal{N}(\theta|\theta_{MAP}, \Lambda^{-1}) d\theta \approx \frac{1}{M} \sum_{i=1}^M p(y_*|\mathbf{x}_*, \theta^{(i)})$$

Generalized Gauss-Newton method

Using MC approximation is the general purpose option when computing intractable PPD

- Use the [GGN approximation](#) of LA. Equivalent to using a “linearized” model for  $p(y_*|\mathbf{x}_*, \theta)$ , using which we can easily compute PPD using linear Gaussian model results



# Detour: Gradient and Hessian

- For LA (and for optimization general), we need  $\nabla_{\theta} \log p(\mathcal{D}, \theta)$  and  $\nabla_{\theta}^2 \log p(\mathcal{D}, \theta)$
- These depend on the likelihood function,  $p(\mathbf{y}|\mathbf{x}, \theta) = p(\mathbf{y}|f(\mathbf{x}, \theta))$
- The form of the function  $f$  depends on the likelihood model. Some examples:

$$p(\mathbf{y}|\mathbf{x}, \theta) = \mathcal{N}(\mathbf{y}|\theta^{\top} \mathbf{x}, \sigma^2) \quad p(\mathbf{y}|\mathbf{x}, \theta) = \text{multinoulli}(\mathbf{y}|\text{softmax}(\theta^{\top} \mathbf{x}))$$

$$p(\mathbf{y}|\mathbf{x}, \theta) = \mathcal{N}(\mathbf{y}|\text{NN}(\mathbf{x}, \theta), \sigma^2) \quad p(\mathbf{y}|\mathbf{x}, \theta) = \text{multinoulli}(\mathbf{y}|\text{softmax}(\text{NN}(\mathbf{x}, \theta)))$$

- Assume  $\mathbf{y}$  and  $\mathbf{f} = f(\mathbf{x}, \theta)$  both to be vectors of size  $C$ ,  $\theta \in \mathbb{R}^P$  and define

Jacobian of size  
 $C \times P$  with  
 $[J_{\theta}(\mathbf{x})]_{ci} =$   
 $\nabla_{\theta_i} f_c(\mathbf{x}, \theta)$

$$J_{\theta}(\mathbf{x}) = \nabla_{\theta} f$$

Hessian of size  
 $C \times P \times P$  with  
 $[\mathcal{H}_{\theta}(\mathbf{x})]_{cij} =$   
 $\nabla_{\theta_i} \nabla_{\theta_j} f_c(\mathbf{x}, \theta)$

$$\mathcal{H}_{\theta}(\mathbf{x}) = \nabla_{\theta}^2 f$$

$$\mathbf{r}(\mathbf{y}; \mathbf{f}) = \nabla_{\mathbf{f}} \log p(\mathbf{y}|\mathbf{f})$$

$$\mathbf{L}(\mathbf{y}; \mathbf{f}) = -\nabla_{\mathbf{f}}^2 \log p(\mathbf{y}|\mathbf{f})$$

$$\nabla_{\theta} \log p(\mathbf{y}|f(\mathbf{x}, \theta)) = J_{\theta}(\mathbf{x})^{\top} \mathbf{r}(\mathbf{y}; \mathbf{f})$$

$$\nabla_{\theta}^2 \log p(\mathbf{y}|f(\mathbf{x}, \theta)) = \mathcal{H}_{\theta}(\mathbf{x})^{\top} \mathbf{r}(\mathbf{y}; \mathbf{f}) - J_{\theta}(\mathbf{x})^{\top} \mathbf{L}(\mathbf{y}; \mathbf{f}) J_{\theta}(\mathbf{x})$$



# Generalized Gauss-Newton (GGN) Approximation<sup>14</sup>

- The Hessian of the log-likelihood turned out to be

$$\nabla_{\theta}^2 \log p(\mathbf{y}|f(\mathbf{x}, \theta)) = \mathcal{H}_{\theta}(\mathbf{x})^{\top} \mathbf{r}(\mathbf{y}; \mathbf{f}) - \mathcal{J}_{\theta}(\mathbf{x})^{\top} \mathbf{L}(\mathbf{y}; \mathbf{f}) \mathcal{J}_{\theta}(\mathbf{x})$$

- Ignoring the term involving  $\mathcal{H}_{\theta}(\mathbf{x}) = \nabla_{\theta}^2 f$ , we have an approximation

$$\nabla_{\theta}^2 \log p(\mathbf{y}|f(\mathbf{x}, \theta)) \approx -\mathcal{J}_{\theta}(\mathbf{x})^{\top} \mathbf{L}(\mathbf{y}; \mathbf{f}) \mathcal{J}_{\theta}(\mathbf{x})$$

This approximation of the Hessian is guaranteed to be positive semi-definite unlike the original Hessian because  $-\log p(\mathbf{y}|f(\mathbf{x}, \theta))$  may not be convex in  $\theta$

- This is called the Generalized Gauss-Newton (GGN) approximation\* of the precision matrix used in Laplace Approximation

- We can further apply diagonal or block-diagonal approximations for efficiency\*

Reason:  $\mathcal{H}_{\theta}(\mathbf{x})$  will be 0 for a linear function

- GGN is also equivalent to approximating  $\mathbf{f} = f(\mathbf{x}, \theta)$  by a linear function of  $\theta$

Gradient vector acting as “features” in this linear model

Gradient of  $f$  at  $\theta = \theta_{MAP}$

Not the gradient of the (log) likelihood – that gradient is zero at  $\theta_{MAP}$

A linear function of  $\theta$

Makes PPD easy to compute when using Laplace approximation

A nonlinear function, e.g., a neural net, approximated by a linear function

$$f(\mathbf{x}_*, \theta) \approx f(\mathbf{x}_*, \theta_{MAP}) + \nabla_{\theta_{MAP}} f^{\top} (\theta - \theta_{MAP}) = f_{lin}(\mathbf{x}_*, \theta)$$

\*Improving predictions of Bayesian neural nets via local linearization (Immer et al, 2021)



# PPD with GGN/Linearized Laplace's Approximation<sup>15</sup>

- Assuming  $p(\mathbf{y}|\mathbf{x}, \theta) = p(\mathbf{y}|f(\mathbf{x}, \theta))$ , LA based PPD is

$$p(\mathbf{y}_*|\mathbf{x}_*, \mathcal{D}) = \int p(\mathbf{y}_*|f(\mathbf{x}_*, \theta))p(\theta|\mathcal{D})d\theta \approx \int p(\mathbf{y}_*|f(\mathbf{x}_*, \theta))\mathcal{N}(\theta|\theta_{MAP}, \Lambda^{-1})d\theta$$

- We can use GGN and Linearized Laplace idea in two ways for the above PPD
- Use  $f(\mathbf{x}_*, \theta)$  but use  $\mathcal{N}(\theta|\theta_{MAP}, \Lambda_{GGN}^{-1})$  as approx post instead  $\mathcal{N}(\theta|\theta_{MAP}, \Lambda^{-1})$ 
  - May require Monte Carlo integration if PPD integral is intractable (e.g., if  $f$  is a neural net or non-lin func)
  - Less commonly used and is less accurate\*
- Use  $f_{lin}(\mathbf{x}_*, \theta)$  instead of  $f(\mathbf{x}_*, \theta)$  and also use  $\mathcal{N}(\theta|\theta_{MAP}, \Lambda_{GGN}^{-1})$  as approx. post.
  - Assuming  $p(\mathbf{y}_*|f(\mathbf{x}_*, \theta)) = \mathcal{N}(y_*|f_{lin}(\mathbf{x}_*, \theta), \beta^{-1})$  for scalar-valued regression

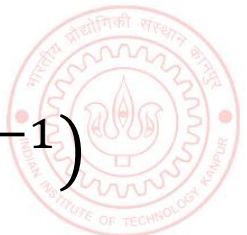
Linear transformation of  $\theta$  with  
 $p(\theta|\mathcal{D}) = \mathcal{N}(\theta|\theta_{MAP}, \Lambda_{GGN}^{-1})$  and  
Gaussian noise  $\epsilon \sim \mathcal{N}(0, \beta^{-1})$

$$y_* \approx f_{lin}(\mathbf{x}_*, \theta) + \epsilon$$

Even though  $f(\mathbf{x}_*, \theta)$  is a complex function like neural net, using linearized Laplace approx, we get PPD in closed form

$$p(y_*|\mathbf{x}_*, \mathcal{D}) \approx \mathcal{N}(y_*|f(\mathbf{x}_*, \theta_{MAP}), \nabla_{\theta_{MAP}} f^\top \Lambda_{GGN}^{-1} \nabla_{\theta_{MAP}} f + \beta^{-1})$$

• \*'In-Between' Uncertainty in Bayesian Neural Networks (Foong et al, 2019),  
• \*Improving predictions of Bayesian neural nets via local linearization (Immer et al, 2021)



# Standard Laplace vs Linearized Laplace

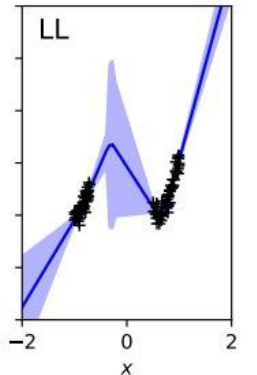
- Standard LA based PPD is usually computed using Monte Carlo sampling

$$p(y_*|\mathbf{x}_*, \mathcal{D}) \approx \int p(y_*|f(\mathbf{x}_*, \theta)) \mathcal{N}(\theta|\theta_{MAP}, \Lambda^{-1}) d\theta \approx \frac{1}{M} \sum_{i=1}^M p(y_*|f(\mathbf{x}_*, \theta^{(i)}))$$

- If the samples  $\theta^{(i)}$  don't come from high-prob regions of the posterior, the above PPD may have poor accuracy (often happens for high-dim posteriors)

- Linearized Laplace based PPD is computed as

$$p(y_*|\mathbf{x}_*, \mathcal{D}) \approx \mathcal{N}(y_*|f(\mathbf{x}_*, \theta_{MAP}), \nabla_{\theta_{MAP}} f^\top \Lambda_{\text{GGN}}^{-1} \nabla_{\theta_{MAP}} f + \beta^{-1})$$



- Linearized Laplace based PPD typically is reasonably accurate and sometimes even more accurate than standard LA with PPD computed using MC sampling\*





# Logistic Regression PPD using Monte Carlo

- The posterior predictive distribution can be computed as

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int p(y_* = 1 | \mathbf{w}, \mathbf{x}_*) p(\mathbf{w} | \mathbf{X}, \mathbf{y}) d\mathbf{w}$$

Integral not tractable and must be approximated

sigmoid

Gaussian (if using Laplace approx.)

- Monte-Carlo approximation of this integral is one possible way

- Draw  $M$  samples  $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_M$ , from the approx. of posterior
  - Approximate the PPD as follows

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) \approx \frac{1}{M} \sum_{m=1}^M p(y_* = 1 | \mathbf{w}_m, \mathbf{x}_*) = \frac{1}{M} \sum_{m=1}^M \sigma(\mathbf{w}_m^\top \mathbf{x}_*)$$

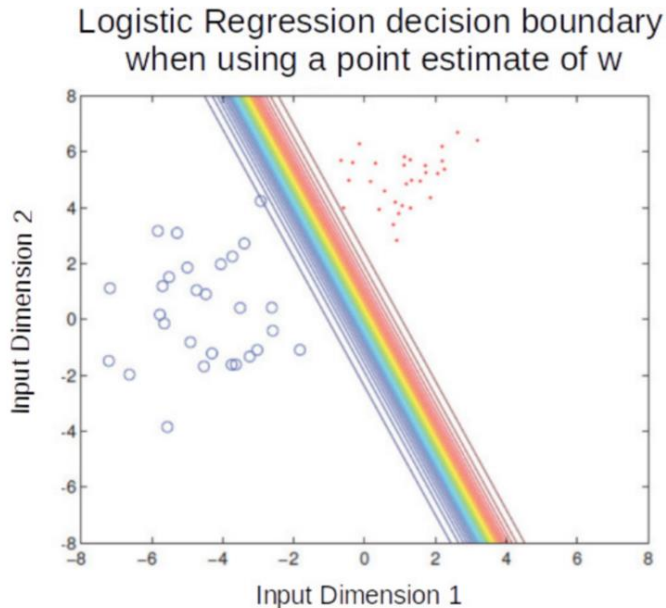
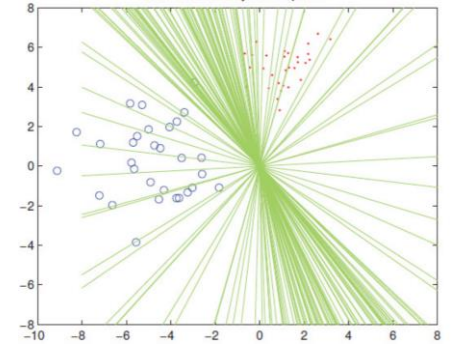
- In contrast, when using MLE/MAP solution  $\hat{\mathbf{w}}_{opt}$ , the plug-in pred. distribution

$$\begin{aligned} p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \int p(y_* = 1 | \mathbf{w}, \mathbf{x}_*) p(\mathbf{w} | \mathbf{X}, \mathbf{y}) d\mathbf{w} \\ &\approx p(y_* = 1 | \hat{\mathbf{w}}_{opt}, \mathbf{x}_*) = \sigma(\hat{\mathbf{w}}_{opt}^\top \mathbf{x}_*) \end{aligned}$$

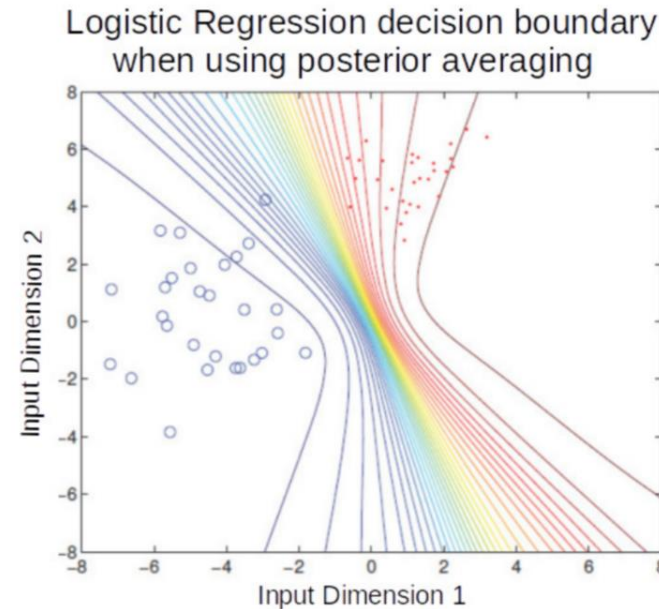


# LR: Plug-in Prediction vs Bayesian Averaging

- Plug-in prediction uses a single  $\mathbf{w}$  (point est) to make prediction
- PPD does an averaging using all possible  $\mathbf{w}$ 's from the posterior



Color transitions (red to blue) in both plots denote how the probability of an input changes from belonging to red class to belonging to blue class. All inputs on a line (or curve on RHS plot) have the same probability of belonging to the red/blue class



Posterior averaging is like using an ensemble of models. In this example, each model is a linear classifier but the ensemble-like effect resulted in nonlinear boundaries

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) \approx \sigma(\hat{\mathbf{w}}_{opt}^T \mathbf{x}_n)$$

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) \approx \frac{1}{M} \sum_{m=1}^M \sigma(\mathbf{w}_m^T \mathbf{x}_n)$$



# Exp. Family (Pitman, Darmois, Koopman, 1930s)

- Defines a class of distributions. An Exponential Family distribution is of the form

$$p(\mathbf{x}|\theta) = \frac{1}{Z(\theta)} h(\mathbf{x}) \exp[\theta^\top \phi(\mathbf{x})] = h(\mathbf{x}) \exp[\theta^\top \phi(\mathbf{x}) - A(\theta)]$$

- $\mathbf{x} \in \mathcal{X}^m$  is the r.v. being modeled ( $\mathcal{X}$  denotes some space, e.g.,  $\mathbb{R}$  or  $\{0,1\}$ )
- $\theta \in \mathbb{R}^d$  : Natural parameters or canonical parameters defining the distribution
- $\phi(\mathbf{x}) \in \mathbb{R}^d$  : Sufficient statistics (another random variable)
  - Knowing this quantity suffices to estimate parameter  $\theta$  from  $\mathbf{x}$
- $Z(\theta) = \int h(\mathbf{x}) \exp[\theta^\top \phi(\mathbf{x})] d\mathbf{x}$ : Partition Function
- $A(\theta) = \log Z(\theta)$ : Log-partition function (also called cumulant function)
- $h(\mathbf{x})$ : A constant (doesn't depend on  $\theta$ )



# Expressing a Distribution in Exp. Family Form

- Recall the form of exp-fam distribution  $p(x|\theta) = h(x)\exp[\theta^\top \phi(x) - A(\theta)]$
- To write any exp-fam dist  $p()$  in the above form, write it as  $\exp(\log p())$

$$\begin{aligned} \exp(\log \text{Binomial}(x|N, \mu)) &= \exp\left(\log \binom{N}{x} \mu^x (1 - \mu)^{N-x}\right) \\ &= \exp\left(\log \binom{N}{x} + x \log \mu + (N - x) \log(1 - \mu)\right) \\ &= \binom{N}{x} \exp\left(x \log \frac{\mu}{1 - \mu} - N \log(1 - \mu)\right) \end{aligned}$$

- Now compare the resulting expression with the exponential family form

$$p(x|\theta) = h(x)\exp[\theta^\top \phi(x) - A(\theta)]$$

.. to identify the natural parameters, sufficient statistics, log-partition function, etc.



# (Univariate) Gaussian as Exponential Family

- Let's try to write a univariate Gaussian in the exponential family form

$$p(\mathbf{x}|\theta) = h(\mathbf{x}) \exp[\theta^\top \phi(\mathbf{x}) - A(\theta)]$$

- Recall the PDF of a univar Gaussian (already has exp, so less work needed :))

$$\begin{aligned} \mathcal{N}(x|\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right] = \frac{1}{\sqrt{2\pi}} \exp\left[\frac{\mu}{\sigma^2}x - \frac{1}{2\sigma^2}x^2 - \frac{\mu^2}{2\sigma^2} - \log \sigma\right] \\ &= \frac{1}{\sqrt{2\pi}} \exp\left[\begin{bmatrix} \frac{\mu}{\sigma^2} \\ -\frac{1}{2\sigma^2} \end{bmatrix}^\top \begin{bmatrix} x \\ x^2 \end{bmatrix} - \left(\frac{\mu^2}{2\sigma^2} + \log \sigma\right)\right] \end{aligned}$$

$$\theta = \begin{bmatrix} \frac{\mu}{\sigma^2} \\ -\frac{1}{2\sigma^2} \end{bmatrix} = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \quad \phi(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix}, \text{ and } \begin{bmatrix} \mu \\ \sigma^2 \end{bmatrix} = \begin{bmatrix} -\frac{\theta_1}{2\theta_2} \\ -\frac{1}{2\theta_2} \end{bmatrix}$$

$$h(x) = \frac{1}{\sqrt{2\pi}} \quad A(\theta) = \frac{\mu^2}{2\sigma^2} + \log \sigma = \frac{-\theta_1^2}{4\theta_2} - \frac{1}{2} \log(-2\theta_2) - \frac{1}{2} \log(2\pi)$$



# Other Examples

- Many other distribution belong to the exponential family
  - Bernoulli
  - Beta
  - Gamma
  - Multinoulli/Multinomial
  - Dirichlet
  - Multivariate Gaussian
  - .. and many more ( [https://en.wikipedia.org/wiki/Exponential\\_family](https://en.wikipedia.org/wiki/Exponential_family) )
- Note: Not all distributions belong to the exponential family, e.g.,
  - Uniform distribution ( $x \sim \text{Unif}(a, b)$ )
  - Student-t distribution
  - Mixture distributions (e.g., mixture of Gaussians)



# Log-Partition Function

- The log-partition function is  $A(\theta) = \log Z(\theta) = \log \int h(\mathbf{x}) \exp[\theta^\top \phi(\mathbf{x})] d\mathbf{x}$
- $A(\theta)$  is also called the **cumulant function**
- Derivatives of  $A(\theta)$  can be used to generate the cumulants of the sufficient statistics
- Exercise: Assume  $\theta$  to be a scalar (thus  $\phi(\mathbf{x})$  is also scalar). Show that the first and the second derivatives of  $A(\theta)$  are

$$\begin{aligned}\frac{dA}{d\theta} &= \mathbb{E}_{p(\mathbf{x}|\theta)}[\phi(\mathbf{x})] \\ \frac{d^2A}{d\theta^2} &= \mathbb{E}_{p(\mathbf{x}|\theta)}[\phi^2(\mathbf{x})] - [\mathbb{E}_{p(\mathbf{x}|\theta)}[\phi(\mathbf{x})]]^2 = \text{var}[\phi(\mathbf{x})]\end{aligned}$$

- Above result also holds when  $\theta$  and  $\phi(\mathbf{x})$  are vector-valued (the “var” will be “covar”)
- Important:  $A(\theta)$  is a convex function of  $\theta$ . Why?





# MLE for Exponential Family Distributions

- Assume data  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  drawn i.i.d. from an exp. family distribution

$$p(x|\theta) = h(x)\exp[\theta^\top \phi(x) - A(\theta)]$$

- To do MLE, we need the overall likelihood -- a product of the individual likelihoods

$$p(\mathcal{D}|\theta) = \prod_{i=1}^N p(\mathbf{x}_i|\theta) = \left[ \prod_{i=1}^N h(\mathbf{x}_i) \right] \exp \left[ \theta^\top \sum_{i=1}^N \phi(\mathbf{x}_i) - NA(\theta) \right] = \left[ \prod_{i=1}^N h(\mathbf{x}_i) \right] \exp \left[ \theta^\top \phi(\mathcal{D}) - NA(\theta) \right]$$

- To estimate  $\theta$  (as we'll see shortly), we only need  $\phi(\mathcal{D}) = \sum_{i=1}^N \phi(\mathbf{x}_i)$  and  $N$
- Size of  $\phi(\mathcal{D}) = \sum_{i=1}^N \phi(\mathbf{x}_i)$  does not grow with  $N$  (same as the size of each  $\phi(\mathbf{x}_i)$ )
- Only exponential family distributions have finite-sized sufficient statistics
  - No need to store all the data; can simply update the sufficient statistics as data comes
  - Useful in probabilistic inference with large-scale data sets and “online” parameter estimation





# Bayesian Inference for Expon. Family Distributions<sup>25</sup>

- Already saw that the total **likelihood** given  $N$  i.i.d. observations  $\mathcal{D} = \{x_1, \dots, x_N\}$

$$p(\mathcal{D}|\theta) \propto \exp \left[ \theta^\top \phi(\mathcal{D}) - NA(\theta) \right] \quad \text{where} \quad \phi(\mathcal{D}) = \sum_{i=1}^N \phi(x_i)$$

- Let's choose the following **prior** (note: looks similar in terms of  $\theta$  within exp)

$$p(\theta|\nu_0, \tau_0) = h(\theta) \exp \left[ \theta^\top \tau_0 - \nu_0 A(\theta) - A_c(\nu_0, \tau_0) \right]$$

- Ignoring the prior's log-partition function  $A_c(\nu_0, \tau_0) = \log \int_{\theta} h(\theta) \exp \left[ \theta^\top \tau_0 - \nu_0 A(\theta) \right] d\theta$

$$p(\theta|\nu_0, \tau_0) \propto h(\theta) \exp \left[ \theta^\top \tau_0 - \nu_0 A(\theta) \right]$$

- Comparing the prior's form with the likelihood, note that
  - $\nu_0$  is like the number of “pseudo-observations” coming from the prior
  - $\tau_0$  is the total sufficient statistics of the pseudo-observations ( $\tau_0 / \nu_0$  per pseudo-obs)



# The Posterior

- The likelihood and prior were

$$p(\mathcal{D}|\theta) \propto \exp \left[ \theta^\top \phi(\mathcal{D}) - NA(\theta) \right] \quad \text{where} \quad \phi(\mathcal{D}) = \sum_{i=1}^N \phi(\mathbf{x}_i)$$

Assume its log partition function denoted as  $A_c(\nu_0, \tau_0)$

$$p(\theta|\nu_0, \tau_0) \propto h(\theta) \exp \left[ \theta^\top \tau_0 - \nu_0 A(\theta) \right]$$

Posterior is also from the same family as the prior

Happens when the prior is conjugate to the likelihood

- The posterior  $p(\theta|\mathcal{D}) \propto p(\theta)p(\mathcal{D}|\theta)$  therefore will be

Its log partition function will be  $A_c(\nu_0 + N, \tau_0 + \phi(\mathcal{D}))$

$$p(\theta|\mathcal{D}) \propto h(\theta) \exp \left[ \theta^\top (\tau_0 + \phi(\mathcal{D})) - (\nu_0 + N)A(\theta) \right]$$

- Every exp family likelihood has a conjugate prior having the form above
- Posterior's hyperparams  $\tau'_0, \nu'_0$  obtained by adding "stuff" to prior's hyperparams

Number of pseudo-observations plus number of actual observations

$$\nu'_0 \leftarrow \nu_0 + N$$

Suff-stats of pseudo-observations plus suff-stats of actual observations

$$\tau'_0 \leftarrow \tau_0 + \phi(\mathcal{D})$$

Another equivalent form

$$p(\theta|\mathcal{D}) \propto h(\theta) \exp \left[ \theta^\top (\nu_0 + N) \frac{\nu_0 \bar{\tau}_0 + \phi(\mathcal{D})}{\nu_0 + N} - (\nu_0 + N)A(\theta) \right]$$

$$\bar{\tau}_0 = \tau_0 / \nu_0$$

$$\begin{aligned} \nu'_0 &\leftarrow \nu_0 + N \\ \bar{\tau}'_0 &\leftarrow \frac{\nu_0 \bar{\tau}_0 + N \bar{\phi}}{\nu_0 + N} \end{aligned}$$

$$\bar{\phi} = \frac{\phi(\mathcal{D})}{N}$$

Convex comb of avg suff-stats of pseudo obs and actual obs



# Posterior Predictive Distribution

- Assume some training data  $\mathcal{D} = \{x_1, \dots, x_N\}$  from some exp-fam distribution
- Assume some test data  $\mathcal{D}' = \{\tilde{x}_1, \dots, \tilde{x}_{N'}\}$  from the same distribution
- The posterior pred. distr. of  $\mathcal{D}'$

$$\begin{aligned}
 p(\mathcal{D}' | \mathcal{D}) &= \int \underbrace{p(\mathcal{D}' | \theta)}_{\text{Exp. Fam. likelihood w.r.t. test data}} \underbrace{p(\theta | \mathcal{D})}_{\text{Posterior (same form as the prior due to conjugacy)}} d\theta \\
 &= \int \underbrace{\left[ \prod_{i=1}^{N'} h(\tilde{x}_i) \right]}_{\text{constant w.r.t. } \theta} \exp \left[ \theta^\top \phi(\mathcal{D}') - N' A(\theta) \right] h(\theta) \exp \left[ \theta^\top (\tau_0 + \phi(\mathcal{D})) - (\nu_0 + N) A(\theta) - \underbrace{A_c(\nu_0 + N, \tau_0 + \phi(\mathcal{D}))}_{\text{constant w.r.t. } \theta} \right] d\theta
 \end{aligned}$$

- This gets further simplified into

$$\begin{aligned}
 p(\mathcal{D}' | \mathcal{D}) &= \left[ \prod_{i=1}^{N'} h(\tilde{x}_i) \right] \frac{\int h(\theta) \exp \left[ \theta^\top (\tau_0 + \phi(\mathcal{D}) + \phi(\mathcal{D}')) - (\nu_0 + N + N') A(\theta) \right] d\theta}{\exp [A_c(\nu_0 + N, \tau_0 + \phi(\mathcal{D}))]} \\
 &= \left[ \prod_{i=1}^{N'} h(\tilde{x}_i) \right] \frac{Z_c(\nu_0 + N + N', \tau_0 + \phi(\mathcal{D}) + \phi(\mathcal{D}'))}{\exp [A_c(\nu_0 + N, \tau_0 + \phi(\mathcal{D}))]}
 \end{aligned}$$



# Posterior Predictive Distribution

- Since  $A_c = \log Z_c$  or  $Z_c = \exp(A_c)$ , we can write the PPD as

$$\begin{aligned} p(\mathcal{D}'|\mathcal{D}) &= \left[ \prod_{i=1}^{N'} h(\tilde{\mathbf{x}}_i) \right] \frac{Z_c(\boldsymbol{\nu}_0 + N + N', \boldsymbol{\tau}_0 + \phi(\mathcal{D}) + \phi(\mathcal{D}'))}{Z_c(\boldsymbol{\nu}_0 + N, \boldsymbol{\tau}_0 + \phi(\mathcal{D}))} \\ &= \left[ \prod_{i=1}^{N'} h(\tilde{\mathbf{x}}_i) \right] \exp [A_c(\boldsymbol{\nu}_0 + N + N', \boldsymbol{\tau}_0 + \phi(\mathcal{D}) + \phi(\mathcal{D}')) - A_c(\boldsymbol{\nu}_0 + N, \boldsymbol{\tau}_0 + \phi(\mathcal{D}))] \end{aligned}$$

Thus PPD as well as marginal likelihood has closed form expression when working with exp-family distributions



- Therefore the **posterior predictive** is proportional to
  - Ratio of two partition functions of two “posterior distributions” (one with  $N + N'$  examples and the other with  $N$  examples)
  - Exponential of the difference of the corresponding log-partition functions
- Note that the form of  $Z_c$  (and  $A_c$ ) will simply depend on the chosen conjugate prior
- Very useful result. Also holds for  $N = 0$ 
  - In this case  $p(\mathcal{D}') = \int p(\mathcal{D}'|\theta)p(\theta)d\theta$  is simply the **marginal likelihood** of test data  $\mathcal{D}'$



# Summary

- Exp. family distributions are very useful for modeling diverse types of data/parameters
- Conjugate priors to exp. family distributions make parameter updates very simple
- Other quantities such as posterior predictive can be computed in closed form
- Useful in designing generative classification models. Choosing class-conditional from exponential family with conjugate priors helps in parameter estimation
- Useful in designing generative models for unsupervised learning
- Used in designing **Generalized Linear Models**: Model  $p(\mathbf{y}|\mathbf{x})$  using exp. fam distribution
  - Linear regression (with Gaussian likelihood) and logistic regression are GLMs
- Will see several use cases when we discuss approx inference algorithms (e.g., Gibbs sampling, and especially variational inference)

