

# Laplace Approximation (Contd) and Generalized Linear Models

CS772A: Probabilistic Machine Learning

Piyush Rai

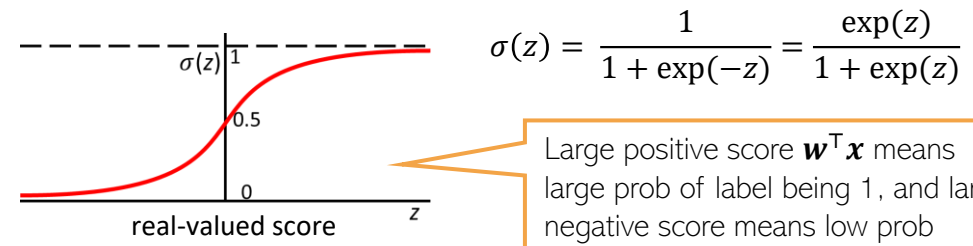
# Logistic Regression

There are other ways too that can convert the score into a probability, such as a CDF:  $p(y = 1|\mathbf{x}, \mathbf{w}) = \mu = \Phi(\mathbf{w}^T \mathbf{x})$  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  $\mathbf{w}^T \mathbf{x}$  for input  $\mathbf{x}$
- A sigmoid function maps this real-valued score into probability of label being 1

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



Large positive score  $\mathbf{w}^T \mathbf{x}$  means large prob of label being 1, and large negative score means low prob

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

Likelihood

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

- Can use a Gaussian prior on  $\mathbf{w}$ :  $p(\mathbf{w}|\lambda) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1} \mathbf{I})$
- Point estimation (MLE/MAP) for LR gives global optima (NLL is convex in  $\mathbf{w}$ )
- We will mainly focus on fully Bayesian inference (computing the posterior)

Can also use a sparsity-inducing prior, such as spike-and-slab or a scale-mixture of Gaussians



# Logistic Regression: The Posterior

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

- The posterior will be

Gaussian

Bernoulli

Hyperparam  $\lambda$  not shown

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

Unfortunately, Gaussian and Bernoulli are not conjugate with each other, so analytic expression for the posterior can't be obtained unlike prob. linear regression



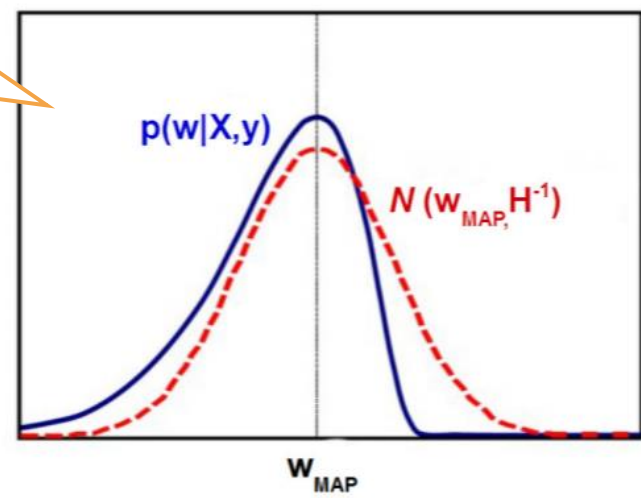
Other approx. inference methods, such as MCMC and VI later

- Need to approximate the posterior in this case

- For now, we will use a simple approximation called Laplace approximation

Laplace approx: Approximates the intractable posterior by a Gaussian whose mean is the MAP solution of the LR model

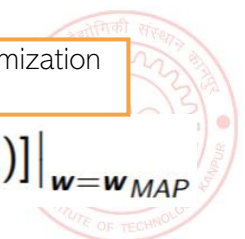
.. and the covariance matrix of this Gaussian is set to the inverse of the Hessian matrix (second derivative) of the model's negative log-joint of params and data, evaluated at the MAP solution



$$\begin{aligned}
 \mathbf{w}_{MAP} &= \arg \max_{\mathbf{w}} \log p(\mathbf{w} | \mathbf{y}, \mathbf{X}) \\
 &= \arg \max_{\mathbf{w}} \log p(\mathbf{y}, \mathbf{w} | \mathbf{X}) \\
 &= \arg \min_{\mathbf{w}} [-\log p(\mathbf{y}, \mathbf{w} | \mathbf{X})]
 \end{aligned}$$

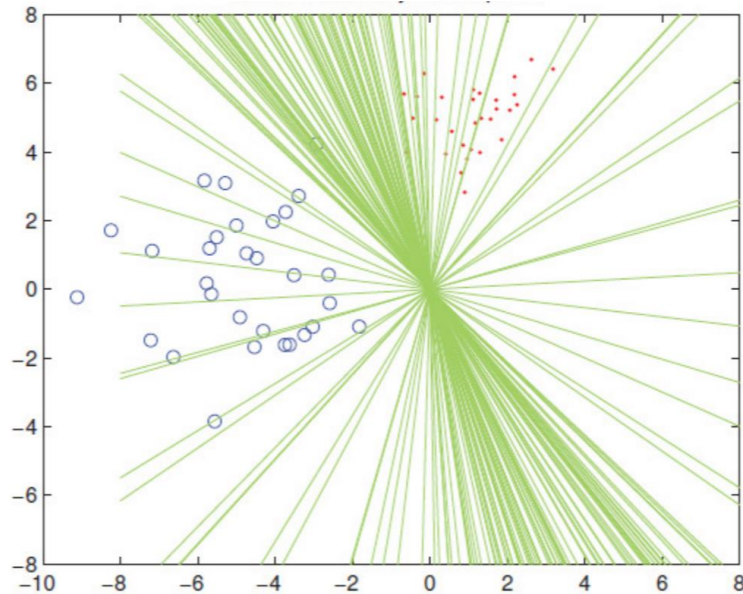
First or second-order optimization methods can be used

$$\mathbf{H} = \nabla^2 [-\log p(\mathbf{y}, \mathbf{w} | \mathbf{X})] \Big|_{\mathbf{w}=\mathbf{w}_{MAP}}$$



# LR Posterior: An Illustration

- Assuming the Gaussian approximation, some samples from the posterior of LR



Not all separators from from the posterior are equally good; their “goodness” will depends on their posterior probabilities  $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$



When making predictions, we can still use all of them but weighted by their importance based on their posterior probabilities

That's exactly what we do when computing the predictive distribution

- Each sample drawn from  $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$  will give a weight vector
- Each such  $\mathbf{w}$  corresponds to one of the separators in the above figure



# LR: Posterior Predictive Distribution

- 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

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

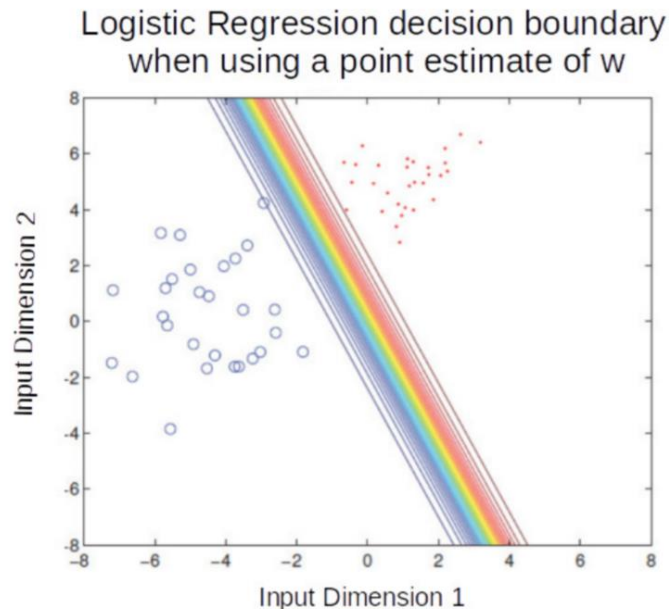
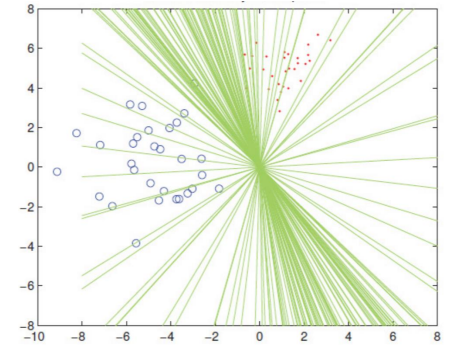
- In contrast, when using MLE/MAP solution  $\hat{\mathbf{w}}_{opt}$ , the plug-in pred. distribution

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

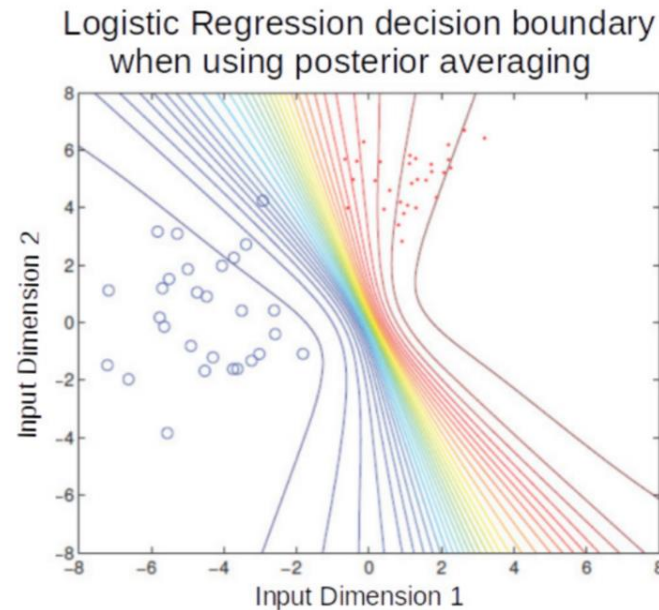


# LR: Plug-in Prediction vs Bayesian Averaging

- 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

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

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





# 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

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  $y_{n\ell'} = 0 \forall \ell' \neq \ell$

- Can do MLE/MAP/fully Bayesian estimation for  $\mathbf{W}$  similar to LR model



# Laplace Approximation of Posterior Distribution

- Consider a posterior distribution that is intractable to compute

$$p(\theta|\mathcal{D}) = \frac{\overset{\text{Unknowns of the model}}{p(\mathcal{D}|\theta)} \overset{\text{Data}}{p(\theta)}}{p(\mathcal{D})} = \frac{p(\mathcal{D}, \theta)}{p(\mathcal{D})}$$

- Laplace approximation approximates the above using a Gaussian distribution

$$\theta_{MAP} = \arg \max_{\theta} p(\theta|\mathcal{D}) = \arg \max_{\theta} p(\mathcal{D}, \theta)$$

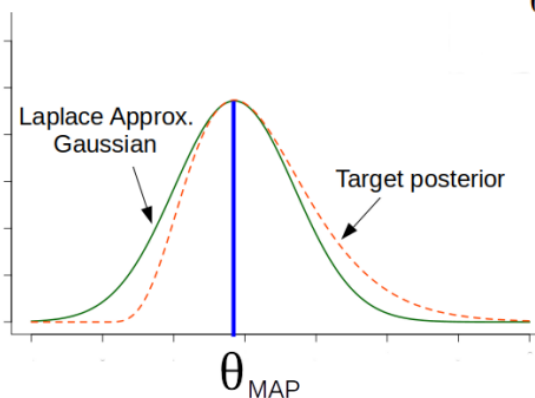
$$= \arg \max_{\theta} p(\mathcal{D}|\theta)p(\theta)$$

$$= \arg \max_{\theta} [\log p(\mathcal{D}|\theta) + \log p(\theta)]$$

$$\mathbf{H} = -\nabla^2 \log p(\theta|\mathcal{D})|_{\theta=\theta_{MAP}} = -\nabla^2 \log p(\mathcal{D}, \theta)|_{\theta=\theta_{MAP}}$$

$$= -\nabla^2 [\log p(\mathcal{D}|\theta) + \log p(\theta)]|_{\theta=\theta_{MAP}}$$

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



- Why is the above Gaussian a reasonable approximation to the posterior?





# Derivation of the Laplace Approximation

- Let's write the Bayes rule as

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}, \theta)}{p(\mathcal{D})} = \frac{p(\mathcal{D}, \theta)}{\int p(\mathcal{D}, \theta) d\theta} = \frac{e^{\log p(\mathcal{D}, \theta)}}{\int e^{\log p(\mathcal{D}, \theta)} d\theta}$$

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

$-\mathbf{H}$

Recall that Hessian is the second derivative of the negative of log-joint

Aha! This is a Gaussian!

Comparing with a Gaussian PDF  
Mean =  $\theta_{MAP}$   
Cov. Matrix =  $\mathbf{H}^{-1}$

- Approximating  $\log p(\mathcal{D}, \theta)$  by a quadratic function of  $\theta$  will make it a Gaussian
- Consider the second-order Taylor approx of a function  $f(\theta)$  around some  $\theta_0$

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

- Assuming  $f(\theta) = \log p(\mathcal{D}, \theta)$  and  $\theta_0 = \theta_{MAP}$ ,  $\nabla f(\theta_{MAP}) = \nabla \log p(\mathcal{D}, \theta_{MAP}) = 0$

Constant w.r.t.  $\theta$

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

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

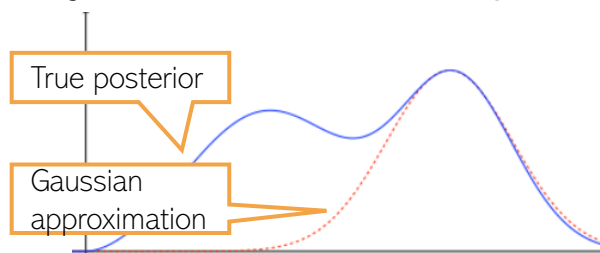


# Properties of Laplace Approximation

- Usually straightforward if derivatives (first and second) can be computed easily
- Expensive if parameter  $\theta$  is very high dimensional
  - Reason: We need to **invert the Hessian** whose size is  $D \times D$  ( $D$  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 approximation\*

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)

- Applicable only when  $\theta$  is real-valued (won't if, say, it is positive, binary etc)
- Note: Even if we have a non-probabilistic model (loss function + regularization), we can obtain an approx “posterior” for that model using the Laplace approximation
  - Optima of the regularized loss function will be Gaussian's mean
  - Second derivative of the regularized loss function will be the Hessian



# Laplace Approx. for High-Dimensional Problems

- When  $\theta$  is very high dim, one option is to approximate the Hessian itself
- One such approx. of the Hessian is a diagonal approximation

Fisher Information Matrix (FIM)

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

FIM is easily computable in auto-diff frameworks used in deep learning

$$\begin{aligned} \mathbf{F} &= \mathbb{E}_{p_{\theta}(\mathbf{x}, \mathbf{y})} [\nabla \log p(\mathbf{y}|\mathbf{x}, \theta) \nabla \log p(\mathbf{y}|\mathbf{x}, \theta)^{\top}] \\ &\approx \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x}, \mathbf{y})} [\nabla \log p(\mathbf{y}|\mathbf{x}, \theta) \nabla \log p(\mathbf{y}|\mathbf{x}, \theta)^{\top}] \\ &= \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \nabla \log p(\mathbf{y}|\mathbf{x}, \theta) \nabla \log p(\mathbf{y}|\mathbf{x}, \theta)^{\top} \end{aligned}$$

Assuming a discriminative model with parameters  $\theta$

Example: A Bayesian neural net for regression/classification ( $\theta$  denotes the weights of the network)

- The diagonal approx. of Hessian may be too crude 😞
  - Ignores covariances among params and treats them as being independent of each other
- A **block-diagonal approx.** proposed recently (in the context of deep neural nets)
  - Treats params across layers to be independent but correlated within the same layer
  - The approach known as Kronecker-Product Factored (KFAC) Laplace approximation



# Generalized Linear Models

- (Probabilistic) Linear Regression: when response  $\mathbf{y}$  is real-valued

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

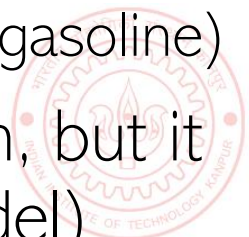
- Logistic Regression: when response  $\mathbf{y}$  is binary (0/1)

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

$$\sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

- In both, the model depends on the inputs  $\mathbf{x}$  via a linear model  $\mathbf{w}^\top \mathbf{x}$
- **Generalized Linear Models** (GLM) allow modeling other types of responses, e.g.,
  - Counts (e.g., predicting the hourly hits on a website)
  - Positive reals (e.g., predicting depth of different pixels in a scene, or stock prices)
  - Fractions between 0 and 1 (e.g., predicting proportion of crude oil convertible to gasoline)
- Note: Can convert responses to real values and apply standard regression, but it is better to model them directly (e.g., for better interpretability of the model)

Note: Probabilistic Linear Regression and Logistic Regression are also GLMs



# Generalized Linear Models: Formally

The reason why GLMs can model a wide variety of responses

- GLMs model the response using an exponential family distribution

Response  $y$  assumed univariate but vector GLMs also exist

Scalar natural param (depends on input  $x$ )

Scalar suff-stats  $\phi(y) = y$

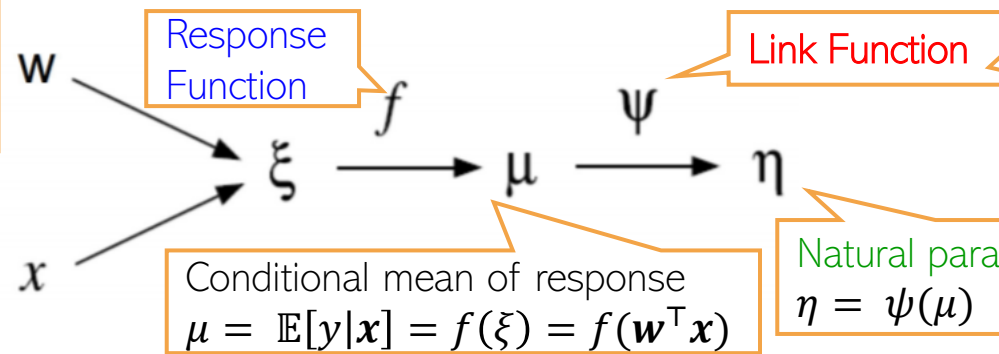
$$p(y|\eta) = h(y) \exp(\eta y - A(\eta))$$

- The inputs  $x$  only appear via a linear model  $\xi = w^T x$  and the overall pipeline is



For prob. linear regression with Gaussian lik,  $f$  is **identity** since mean  $\mu = \mathbb{E}[y|x] = w^T x$

For logistic regression,  $f$  is **sigmoid** since mean  $\mu = \mathbb{E}[y|x] = \sigma(w^T x)$



For GLM with **Canonical Response Function** (next slide),  $\psi = f^{-1}$  and thus nat. param.  $\eta = \xi = w^T x$

Natural parameter  $\eta = \psi(\mu)$

$f$  known as "**inverse link function**" in this case

- Note: Some GLM are represented via exponential **dispersion** family given by

If  $\sigma^2$  is fixed, it is the standard exponential family

$$p(y|\eta, \sigma^2) = h(y, \sigma^2) \exp \left[ \frac{\eta y - A(\eta)}{\sigma^2} \right]$$

Called the "dispersion parameter"

Examples: Gaussian GLM, Gamma GLM

Recall cumulant results of exp-fam

$$\mathbb{E}[y] = A'(\eta)$$

$$\text{var}[y] = A''(\eta)\sigma^2$$



# Generalized Linear Models: Examples

- Consider the overdispersed GLMs

$$p(y|\eta, \sigma^2) = h(y, \sigma^2) \exp \left[ \frac{\eta y - A(\eta)}{\sigma^2} \right] = \exp \left[ \frac{\eta y - A(\eta)}{\sigma^2} + \log h(y, \sigma^2) \right]$$

Note that here we expressed the Gaussian in the overdispersed GLM form unlike how we did it earlier when discussing exp-family

- Consider a linear regression model with Gaussian likelihood

$$p(y|\mathbf{x}, \mathbf{w}, \sigma^2) \propto \exp \left[ -\frac{(y - \mathbf{w}^\top \mathbf{x})^2}{2\sigma^2} \right] = \exp \left[ -\frac{y^2 + (\mathbf{w}^\top \mathbf{x})^2 - 2y\mathbf{w}^\top \mathbf{x}}{2\sigma^2} \right] = \exp \left[ \frac{y\mathbf{w}^\top \mathbf{x} - (\mathbf{w}^\top \mathbf{x})^2/2}{\sigma^2} - \frac{y^2}{2\sigma^2} \right]$$

- Comparing the expressions,  $\eta = \mathbf{w}^\top \mathbf{x}$ ,  $A(\eta) = \frac{\eta^2}{2}$ ,  $\log h(y, \sigma^2) = -y^2/2\sigma^2$
- Can likewise express other models for exp-family distributions  $p(\mathbf{y}|\mathbf{x})$ 
  - Regardless of the form, all will have  $\eta = \mathbf{w}^\top \mathbf{x}$





# GLM with Canonical Response Function

- For GLM with Canon Resp Func (a.k.a., canonical GLM)

The simple form of canonical GLM (nat. param just a linear function  $\mathbf{w}^\top \mathbf{x}$ ) makes parameter estimation via MLE/MAP easy since gradient and Hessian have simple expressions (though the Hessian may be expensive to compute/invert)

$$p(y|\eta) = h(y) \exp(\eta y - A(\eta)) = h(y) \exp(y \mathbf{w}^\top \mathbf{x} - A(\eta))$$

- Consider doing MLE (assuming  $N$  i.i.d. responses). The log likelihood

$$L(\eta) = \log p(Y|\eta) = \log \prod_{n=1}^N h(y_n) \exp(y_n \mathbf{w}^\top \mathbf{x}_n - A(\eta_n)) = \sum_{n=1}^N \log h(y_n) + \mathbf{w}^\top \sum_{n=1}^N y_n \mathbf{x}_n - \sum_{n=1}^N A(\eta_n)$$

- Convexity of  $A(\eta)$  guarantees a global optima. Gradient of log-lik w.r.t.  $\mathbf{w}$

$$\mathbf{g} = \sum_{n=1}^N \left( y_n \mathbf{x}_n - A'(\eta_n) \frac{d\eta_n}{d\mathbf{w}} \right) = \sum_{n=1}^N (y_n \mathbf{x}_n - \mu_n \mathbf{x}_n) = \sum_{n=1}^N (y_n - \mu_n) \mathbf{x}_n$$

Exp of suff-stats  $\mathbb{E}[y_n]$       Corrective updates for  $\mathbf{w}$

The Hessian can also be shown to be

$$\mathbf{H} = -\nabla \mathbf{g} = \sum_{n=1}^N f'(\eta_n) \mathbf{x}_n \mathbf{x}_n^\top$$



- Note  $\mu_n = f(\xi_n) = f(\mathbf{w}^\top \mathbf{x}_n)$  and  $f = \psi^{-1}$  (“inverse link”) depends on the model

- Real-valued  $y$  (linear regression):  $f$  is identity, i.e.,  $\mu_n = \mathbf{w}^\top \mathbf{x}_n$
- Binary  $y$  (logistic regression):  $f$  is sigmoid function, i.e.,  $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$
- Count-valued  $y$  (Poisson regression):  $f$  is exp, i.e.,  $\mu_n = \exp(\mathbf{w}^\top \mathbf{x}_n)$
- Non-negative  $y$  (gamma regression):  $f$  is inverse negative i.e.,  $\mu_n = -1/(\mathbf{w}^\top \mathbf{x}_n)$

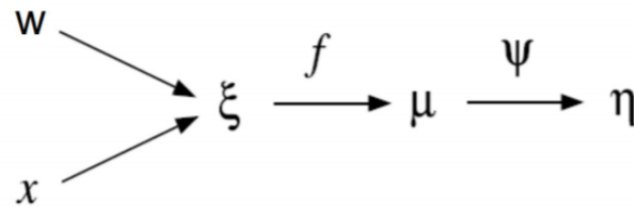


# Fully Bayesian Inference for GLMs

- Most GLMs, except linear regression with Gaussian lik. and Gaussian prior, do not have conjugate pairs of likelihood and priors (recall logistic regression)
- Posterior over the weight vector  $\mathbf{w}$  is intractable
- Approximate inference methods needed, e.g.,
  - Laplace approximation (have already seen): Easily applicable since derivatives (first and second) can be easily computed (note that we need  $\mathbf{w}_{MAP}$  and Hessian)
  - MCMC or variational inference (will see later)



# Various Types of GLMs



Type of response	Type of GLM	Link Function $\Psi$	Response Function $f$ (Inv Link Func if canon. GLM) (Operates on $\xi = w^T x$ )	Activation
Real	Gaussian	Identity	Identity	Linear
Binary	Logistic	Log-odds: $\log \frac{\mu}{1-\mu}$	Sigmoid	Sigmoid
Binary	Probit	Inv CDF: $\Phi^{-1}(\mu)$	$\Phi$ (CDF of $N(0,1)$ )	Probit
Categorical	Multinoulli	Log-odds: $\log \frac{\mu_k}{1-\mu_k}$	Softmax	Softmax
Count	Poisson	$\log \mu$	exp	
Non-negative real	gamma	Negative of inverse	Negative of inverse	
Binary	Gumbel	Gumbel Inv CDF: $\log(-\log())$	Gumbel CDF: $\exp(-\exp(-))$	

.. and several others (exponential, inverse Gaussian, Binomial, Tweedie, etc)

