# Bayesian Logistic Regression, Bayesian Generative Classification

Piyush Rai

#### Topics in Probabilistic Modeling and Inference (CS698X)

Jan 23, 2019

Prob. Mod. & Inference - CS698X (Piyush Rai, IITK)

Bayesian Logistic Regression, Bayesian Generative Classification

イロト イポト イヨト イヨト

• Recall that the likelihood model for logistic regression is Bernoulli (since  $y \in \{0, 1\}$ )

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

• Recall that the likelihood model for logistic regression is Bernoulli (since  $y \in \{0, 1\}$ )

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

 $\, \bullet \,$  Just like the Bayesian linear regression case, let's use a Gausian prior on  $\, {\it w}$ 

$$p(\boldsymbol{w}) = \mathcal{N}(\boldsymbol{0}, \lambda^{-1} \boldsymbol{\mathsf{I}}_D) \propto \exp(-rac{\lambda}{2} \boldsymbol{w}^{ op} \boldsymbol{w})$$

• Recall that the likelihood model for logistic regression is Bernoulli (since  $y \in \{0, 1\}$ )

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

 $\,\circ\,$  Just like the Bayesian linear regression case, let's use a Gausian prior on  ${\it w}$ 

$$p(\boldsymbol{w}) = \mathcal{N}(0, \lambda^{-1} \boldsymbol{\mathsf{I}}_D) \propto \exp(-\frac{\lambda}{2} \boldsymbol{w}^\top \boldsymbol{w})$$

• Given N observations  $(\mathbf{X}, \mathbf{y}) = \{\mathbf{x}_n, y_n\}_{n=1}^N$ , where **X** is  $N \times D$  and **y** is  $N \times 1$ , the posterior over **w** 

$$p(w|\mathbf{X}, y) = \frac{p(y|\mathbf{X}, w)p(w)}{\int p(y|\mathbf{X}, w)p(w)dw}$$

《曰》 《曰》 《臣》 《臣》

• Recall that the likelihood model for logistic regression is Bernoulli (since  $y \in \{0,1\}$ )

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

 $\,\circ\,$  Just like the Bayesian linear regression case, let's use a Gausian prior on  ${\it w}$ 

$$p(\boldsymbol{w}) = \mathcal{N}(0, \lambda^{-1} \boldsymbol{\mathsf{I}}_D) \propto \exp(-rac{\lambda}{2} \boldsymbol{w}^{ op} \boldsymbol{w})$$

• Given N observations  $(\mathbf{X}, \mathbf{y}) = \{\mathbf{x}_n, y_n\}_{n=1}^N$ , where **X** is  $N \times D$  and **y** is  $N \times 1$ , the posterior over **w** 

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

《口》 《卽》 《言》 《言》

• Recall that the likelihood model for logistic regression is Bernoulli (since  $y \in \{0,1\}$ )

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

 $\,\circ\,$  Just like the Bayesian linear regression case, let's use a Gausian prior on  ${\it w}$ 

$$p(\boldsymbol{w}) = \mathcal{N}(0, \lambda^{-1} \boldsymbol{\mathsf{I}}_D) \propto \exp(-rac{\lambda}{2} \boldsymbol{w}^{ op} \boldsymbol{w})$$

• Given N observations  $(\mathbf{X}, \mathbf{y}) = \{\mathbf{x}_n, y_n\}_{n=1}^N$ , where **X** is  $N \times D$  and **y** is  $N \times 1$ , the posterior over **w** 

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

• The denominator is intractable in general (logistic-Bernoulli and Gaussian are not conjugate)

• Recall that the likelihood model for logistic regression is Bernoulli (since  $y \in \{0,1\}$ )

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

 $\,\circ\,$  Just like the Bayesian linear regression case, let's use a Gausian prior on  ${\it w}$ 

$$p(\boldsymbol{w}) = \mathcal{N}(0, \lambda^{-1} \boldsymbol{\mathsf{I}}_D) \propto \exp(-rac{\lambda}{2} \boldsymbol{w}^{ op} \boldsymbol{w})$$

• Given N observations  $(\mathbf{X}, \mathbf{y}) = \{\mathbf{x}_n, y_n\}_{n=1}^N$ , where **X** is  $N \times D$  and **y** is  $N \times 1$ , the posterior over **w** 

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

• The denominator is intractable in general (logistic-Bernoulli and Gaussian are not conjugate)

• Can't get a closed form expression for  $p(w|\mathbf{X}, y)$ . Must approximate it!

• Recall that the likelihood model for logistic regression is Bernoulli (since  $y \in \{0, 1\}$ )

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

 $\,\circ\,$  Just like the Bayesian linear regression case, let's use a Gausian prior on  ${\it w}$ 

$$p(\boldsymbol{w}) = \mathcal{N}(0, \lambda^{-1} \boldsymbol{\mathsf{I}}_D) \propto \exp(-rac{\lambda}{2} \boldsymbol{w}^{ op} \boldsymbol{w})$$

• Given N observations  $(\mathbf{X}, \mathbf{y}) = \{\mathbf{x}_n, y_n\}_{n=1}^N$ , where **X** is  $N \times D$  and **y** is  $N \times 1$ , the posterior over  $\mathbf{w}$ 

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

• The denominator is intractable in general (logistic-Bernoulli and Gaussian are not conjugate)

- Can't get a closed form expression for  $p(w|\mathbf{X}, y)$ . Must approximate it!

• Approximate the posterior distribution  $p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})} = \frac{p(\mathcal{D},\theta)}{p(\mathcal{D})}$  by the following Gaussian  $p(\theta|\mathcal{D}) \approx \mathcal{N}(\theta_{MAP}, \mathbf{H}^{-1})$ 



イロト イロト イヨト イヨト

• Approximate the posterior distribution  $p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})} = \frac{p(\mathcal{D},\theta)}{p(\mathcal{D})}$  by the following Gaussian  $p(\theta|\mathcal{D}) \approx \mathcal{N}(\theta_{MAP}, \mathbf{H}^{-1})$ 



• Note:  $\theta_{MAP}$  is the maximum-a-posteriori (MAP) estimate of  $\theta$ , i.e.,

$$heta_{\textit{MAP}} = rg\max_{ heta} p( heta | \mathcal{D}) = rg\max_{ heta} p(\mathcal{D}, heta)$$

• Approximate the posterior distribution  $p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})} = \frac{p(\mathcal{D},\theta)}{p(\mathcal{D})}$  by the following Gaussian  $p(\theta|\mathcal{D}) \approx \mathcal{N}(\theta_{MAP}, \mathbf{H}^{-1})$ 



• Note:  $\theta_{MAP}$  is the maximum-a-posteriori (MAP) estimate of  $\theta$ , i.e.,

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

• Approximate the posterior distribution  $p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})} = \frac{p(\mathcal{D},\theta)}{p(\mathcal{D})}$  by the following Gaussian  $p(\theta|\mathcal{D}) \approx \mathcal{N}(\theta_{MAP}, \mathbf{H}^{-1})$ 



• Note:  $\theta_{MAP}$  is the maximum-a-posteriori (MAP) estimate of  $\theta$ , i.e.,

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

• Approximate the posterior distribution  $p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})} = \frac{p(\mathcal{D},\theta)}{p(\mathcal{D})}$  by the following Gaussian  $p(\theta|\mathcal{D}) \approx \mathcal{N}(\theta_{MAP}, \mathbf{H}^{-1})$ 



• Note:  $\theta_{MAP}$  is the maximum-a-posteriori (MAP) estimate of  $\theta$ , i.e.,

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

• Usually  $\theta_{MAP}$  can be easily solved for (e.g., using first/second order iterative methods)

• Approximate the posterior distribution  $p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})} = \frac{p(\mathcal{D},\theta)}{p(\mathcal{D})}$  by the following Gaussian  $p(\theta|\mathcal{D}) \approx \mathcal{N}(\theta_{MAP}, \mathbf{H}^{-1})$ 



• Note:  $\theta_{MAP}$  is the maximum-a-posteriori (MAP) estimate of  $\theta$ , i.e.,

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

• Usually  $\theta_{MAP}$  can be easily solved for (e.g., using first/second order iterative methods)

• **H** is the Hessian matrix of the negative log-posterior (or negative log-joint-prob) at  $\theta_{MAP}$ 

$$\mathbf{H} = -
abla^2 \log oldsymbol{
ho}( heta | \mathcal{D}) igert_{ heta = heta_{ extsf{MAF}}}$$

• Approximate the posterior distribution  $p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})} = \frac{p(\mathcal{D},\theta)}{p(\mathcal{D})}$  by the following Gaussian  $p(\theta|\mathcal{D}) \approx \mathcal{N}(\theta_{MAP}, \mathbf{H}^{-1})$ 



• Note:  $\theta_{MAP}$  is the maximum-a-posteriori (MAP) estimate of  $\theta$ , i.e.,

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

• Usually  $\theta_{MAP}$  can be easily solved for (e.g., using first/second order iterative methods)

• **H** is the Hessian matrix of the negative log-posterior (or negative log-joint-prob) at  $\theta_{MAP}$ 

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

Prob. Mod. & Inference - CS698X (Piyush Rai, IITK)

• Approximate the posterior distribution  $p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})} = \frac{p(\mathcal{D},\theta)}{p(\mathcal{D})}$  by the following Gaussian  $p(\theta|\mathcal{D}) \approx \mathcal{N}(\theta_{MAP}, \mathbf{H}^{-1})$ 



• Note:  $\theta_{MAP}$  is the maximum-a-posteriori (MAP) estimate of  $\theta$ , i.e.,

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

• Usually  $\theta_{MAP}$  can be easily solved for (e.g., using first/second order iterative methods)

• **H** is the Hessian matrix of the negative log-posterior (or negative log-joint-prob) at  $\theta_{MAP}$ 

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

Prob. Mod. & Inference - CS698X (Piyush Rai, IITK)

• Let's write the Bayes rule as

$$p( heta | \mathcal{D}) = rac{p(\mathcal{D}, heta)}{p(\mathcal{D})}$$



メロト メロト メモト メモト

• Let's write the Bayes rule as

$$p( heta | \mathcal{D}) = rac{p(\mathcal{D}, heta)}{p(\mathcal{D})} = rac{p(\mathcal{D}, heta)}{\int p(\mathcal{D}, heta) d heta}$$



メロト メロト メモト メモト

• Let's write the Bayes rule as

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



メロト メロト メモト メモト

• Let's write the Bayes rule as

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

• Suppose log  $p(\mathcal{D}, \theta) = f(\theta)$ . Let's approximate  $f(\theta)$  using its 2nd order Taylor expansion

$$f( heta) pprox f( heta_0) + ( heta - heta_0)^ op 
abla f( heta_0) + rac{1}{2}( heta - heta_0)^ op 
abla^2 f( heta_0)( heta - heta_0)$$

where  $\theta_0$  is some arbitrarily chosen point in the domain of f

A B > A B > A B >
 A
 B > A B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A

• Let's write the Bayes rule as

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

• Suppose log  $p(\mathcal{D}, \theta) = f(\theta)$ . Let's approximate  $f(\theta)$  using its 2nd order Taylor expansion

$$f( heta) pprox f( heta_0) + ( heta - heta_0)^ op 
abla f( heta_0) + rac{1}{2}( heta - heta_0)^ op 
abla^2 f( heta_0)( heta - heta_0)$$

where  $\theta_0$  is some arbitrarily chosen point in the domain of f

• Let's choose  $\theta_0 = \theta_{MAP}$ .

• Let's write the Bayes rule as

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

• Suppose log  $p(\mathcal{D}, \theta) = f(\theta)$ . Let's approximate  $f(\theta)$  using its 2nd order Taylor expansion

$$f( heta) pprox f( heta_0) + ( heta - heta_0)^ op 
abla f( heta_0) + rac{1}{2}( heta - heta_0)^ op 
abla^2 f( heta_0)( heta - heta_0)$$

where  $\theta_0$  is some arbitrarily chosen point in the domain of f

• Let's choose  $\theta_0 = \theta_{MAP}$ . Note that  $\nabla f(\theta_{MAP}) = \nabla \log p(\mathcal{D}, \theta_{MAP}) = 0$ .

(日) (문) (문) (문)

• Let's write the Bayes rule as

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

• Suppose log  $p(\mathcal{D}, \theta) = f(\theta)$ . Let's approximate  $f(\theta)$  using its 2nd order Taylor expansion

$$f( heta) pprox f( heta_0) + ( heta - heta_0)^ op 
abla f( heta_0) + rac{1}{2}( heta - heta_0)^ op 
abla^2 f( heta_0)( heta - heta_0)$$

where  $\theta_0$  is some arbitrarily chosen point in the domain of f

• Let's choose  $\theta_0 = \theta_{MAP}$ . Note that  $\nabla f(\theta_{MAP}) = \nabla \log p(\mathcal{D}, \theta_{MAP}) = 0$ . Therefore

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

• Plugging in this 2nd order Taylor approximation for log  $p(\mathcal{D}, \theta)$ , we have

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



• Plugging in this 2nd order Taylor approximation for log  $p(\mathcal{D}, \theta)$ , we have

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



• Plugging in this 2nd order Taylor approximation for log  $p(\mathcal{D}, \theta)$ , we have

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

• Further simplifying, we have

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

A B > A B > A B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A

• Plugging in this 2nd order Taylor approximation for log  $p(\mathcal{D}, \theta)$ , we have

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

Further simplifying, we have  

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

• Therefore the Laplace approximation of the posterior  $p(\theta|\mathcal{D})$  is a Gaussian and is given by

 $p(\theta|\mathcal{D}) \approx \mathcal{N}(\theta|\theta_{MAP}, \mathbf{H}^{-1})$  where  $\mathbf{H} = -\nabla^2 \log p(\mathcal{D}, \theta_{MAP})$ 

• Plugging in this 2nd order Taylor approximation for log  $p(\mathcal{D}, \theta)$ , we have

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

Further simplifying, we have  

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

• Therefore the Laplace approximation of the posterior  $p(\theta|\mathcal{D})$  is a Gaussian and is given by

 $p(\theta|\mathcal{D}) \approx \mathcal{N}(\theta|\theta_{MAP}, \mathbf{H}^{-1})$  where  $\mathbf{H} = -\nabla^2 \log p(\mathcal{D}, \theta_{MAP})$ 



• Usually straightforward if derivatives (first and second) can be computed easily

イロト イロト イヨト イヨト

- Usually straightforward if derivatives (first and second) can be computed easily
- Expensive if the number of parameters is very large (due to Hessian computation and inversion)



- Usually straightforward if derivatives (first and second) can be computed easily
- Expensive if the number of parameters is very large (due to Hessian computation and inversion)
- Can do badly if the (true) posterior is multimodal



< 口 > < 同

- Usually straightforward if derivatives (first and second) can be computed easily
- Expensive if the number of parameters is very large (due to Hessian computation and inversion)
- Can do badly if the (true) posterior is multimodal



• Can actually apply it when working with any regularized loss function (not just probabilistic models) to get a Gaussian posterior distribution over the parameters

- Usually straightforward if derivatives (first and second) can be computed easily
- Expensive if the number of parameters is very large (due to Hessian computation and inversion)
- Can do badly if the (true) posterior is multimodal



- Can actually apply it when working with any regularized loss function (not just probabilistic models) to get a Gaussian posterior distribution over the parameters
  - ${\scriptstyle \circ }$  negative log-likelihood (NLL) = loss function, negative log-prior = regularizer

- Usually straightforward if derivatives (first and second) can be computed easily
- Expensive if the number of parameters is very large (due to Hessian computation and inversion)
- Can do badly if the (true) posterior is multimodal



- Can actually apply it when working with any regularized loss function (not just probabilistic models) to get a Gaussian posterior distribution over the parameters
  - negative log-likelihood (NLL) = loss function, negative log-prior = regularizer
  - Easy exercise: Try doing this for ℓ<sub>2</sub> regularized least squares regression (will get the same posterior as in Bayesian linear regression)

#### Laplace Approximation for Bayesian Logistic Regression

• Data  $\mathcal{D} = (\mathbf{X}, \mathbf{y})$  and parameter  $\theta = \mathbf{w}$ . The Laplace approximation of posterior will be

 $p(oldsymbol{w}|oldsymbol{X},oldsymbol{y}) pprox \mathcal{N}(oldsymbol{w}_{MAP},oldsymbol{H}^{-1})$ 



#### Laplace Approximation for Bayesian Logistic Regression

• Data  $\mathcal{D} = (\mathbf{X}, \mathbf{y})$  and parameter  $\theta = \mathbf{w}$ . The Laplace approximation of posterior will be

 $p(oldsymbol{w}|oldsymbol{X},oldsymbol{y}) pprox \mathcal{N}(oldsymbol{w}_{MAP},oldsymbol{H}^{-1})$ 

• The required quantities are defined as

$$w_{MAP} = \arg\max_{\boldsymbol{w}} \log p(\boldsymbol{w}|\boldsymbol{y}, \boldsymbol{X}) = \arg\max_{\boldsymbol{w}} \log p(\boldsymbol{y}, \boldsymbol{w}|\boldsymbol{X}) = \arg\min_{\boldsymbol{w}} [-\log p(\boldsymbol{y}, \boldsymbol{w}|\boldsymbol{X})]$$


• Data  $\mathcal{D} = (\mathbf{X}, \mathbf{y})$  and parameter  $\theta = \mathbf{w}$ . The Laplace approximation of posterior will be

 $p(oldsymbol{w}|oldsymbol{X},oldsymbol{y}) pprox \mathcal{N}(oldsymbol{w}_{MAP},oldsymbol{H}^{-1})$ 

• The required quantities are defined as

$$w_{MAP} = \arg \max_{\boldsymbol{w}} \log p(\boldsymbol{w}|\boldsymbol{y}, \boldsymbol{X}) = \arg \max_{\boldsymbol{w}} \log p(\boldsymbol{y}, \boldsymbol{w}|\boldsymbol{X}) = \arg \min_{\boldsymbol{w}} [-\log p(\boldsymbol{y}, \boldsymbol{w}|\boldsymbol{X})]$$
$$\boldsymbol{H} = \nabla^{2} [-\log p(\boldsymbol{y}, \boldsymbol{w}|\boldsymbol{X})]|_{\boldsymbol{w}=\boldsymbol{w}_{MAP}}$$

• Data  $\mathcal{D} = (\mathbf{X}, \mathbf{y})$  and parameter  $\theta = \mathbf{w}$ . The Laplace approximation of posterior will be

 $p(oldsymbol{w}|oldsymbol{X},oldsymbol{y}) pprox \mathcal{N}(oldsymbol{w}_{MAP},oldsymbol{H}^{-1})$ 

• The required quantities are defined as

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

• We can compute  $\boldsymbol{w}_{MAP}$  using iterative methods (gradient descent):

• Data  $\mathcal{D} = (\mathbf{X}, \mathbf{y})$  and parameter  $\theta = \mathbf{w}$ . The Laplace approximation of posterior will be

 $p(\boldsymbol{w}|\mathbf{X}, \boldsymbol{y}) \approx \mathcal{N}(\boldsymbol{w}_{MAP}, \mathbf{H}^{-1})$ 

• The required quantities are defined as

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

• We can compute  $\boldsymbol{w}_{MAP}$  using iterative methods (gradient descent):

• First-order (gradient) methods:  $\boldsymbol{w}_{t+1} = \boldsymbol{w}_t - \eta \boldsymbol{g}_t$ . Requires gradient  $\boldsymbol{g}$  of  $-\log p(\boldsymbol{y}, \boldsymbol{w} | \mathbf{X})$ 

$$m{g} = 
abla [-\log p(m{y},m{w}|m{X})]$$

• Data  $\mathcal{D} = (\mathbf{X}, \mathbf{y})$  and parameter  $\theta = \mathbf{w}$ . The Laplace approximation of posterior will be

 $p(\boldsymbol{w}|\mathbf{X}, \boldsymbol{y}) \approx \mathcal{N}(\boldsymbol{w}_{MAP}, \mathbf{H}^{-1})$ 

• The required quantities are defined as

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

• We can compute  $\boldsymbol{w}_{MAP}$  using iterative methods (gradient descent):

• First-order (gradient) methods:  $\boldsymbol{w}_{t+1} = \boldsymbol{w}_t - \eta \boldsymbol{g}_t$ . Requires gradient  $\boldsymbol{g}$  of  $-\log p(\boldsymbol{y}, \boldsymbol{w} | \mathbf{X})$ 

$$m{g} = 
abla [-\log p(m{y},m{w}|m{X})]$$

• Second-order methods.  $\boldsymbol{w}_{t+1} = \boldsymbol{w}_t - \boldsymbol{H}_t^{-1} \boldsymbol{g}_t$ . Requires both gradient and Hessian (defined above)

• Data  $\mathcal{D} = (\mathbf{X}, \mathbf{y})$  and parameter  $\theta = \mathbf{w}$ . The Laplace approximation of posterior will be

 $p(\boldsymbol{w}|\mathbf{X}, \boldsymbol{y}) \approx \mathcal{N}(\boldsymbol{w}_{MAP}, \mathbf{H}^{-1})$ 

• The required quantities are defined as

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

• We can compute  $\boldsymbol{w}_{MAP}$  using iterative methods (gradient descent):

• First-order (gradient) methods:  $\boldsymbol{w}_{t+1} = \boldsymbol{w}_t - \eta \boldsymbol{g}_t$ . Requires gradient  $\boldsymbol{g}$  of  $-\log p(\boldsymbol{y}, \boldsymbol{w} | \mathbf{X})$ 

$$m{g} = 
abla [-\log p(m{y},m{w}|m{X})]$$

• Second-order methods.  $\boldsymbol{w}_{t+1} = \boldsymbol{w}_t - \boldsymbol{H}_t^{-1} \boldsymbol{g}_t$ . Requires both gradient and Hessian (defined above)

 Note: When using second order methods for estimating *w<sub>MAP</sub>*, we anyway get the Hessian needed for the Laplace approximation of the posterior

• The LR objective function  $-\log p(\mathbf{y}, \mathbf{w} | \mathbf{X}) = -\log p(\mathbf{y} | \mathbf{X}, \mathbf{w}) - \log p(\mathbf{w})$  can be written as

$$-\log\prod_{n=1}^{N}p(y_n|x_n, w) - \log p(w)$$

• The LR objective function  $-\log p(\mathbf{y}, \mathbf{w} | \mathbf{X}) = -\log p(\mathbf{y} | \mathbf{X}, \mathbf{w}) - \log p(\mathbf{w})$  can be written as

$$-\log\prod_{n=1}^{N}p(y_n|\boldsymbol{x}_n,\boldsymbol{w}) - \log p(\boldsymbol{w}) = -\sum_{n=1}^{N}\log p(y_n|\boldsymbol{x}_n,\boldsymbol{w}) - \log p(\boldsymbol{w})$$



• The LR objective function  $-\log p(\mathbf{y}, \mathbf{w} | \mathbf{X}) = -\log p(\mathbf{y} | \mathbf{X}, \mathbf{w}) - \log p(\mathbf{w})$  can be written as

$$-\log\prod_{n=1}^{N}p(y_n|\boldsymbol{x}_n,\boldsymbol{w}) - \log p(\boldsymbol{w}) = -\sum_{n=1}^{N}\log p(y_n|\boldsymbol{x}_n,\boldsymbol{w}) - \log p(\boldsymbol{w})$$

• For the logistic regression model,  $p(y_n|\mathbf{x}_n, \mathbf{w}) = \mu_n^{y_n} (1-\mu_n)^{1-y_n}$  where  $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1+\exp(\mathbf{w}^\top \mathbf{x}_n)}$ 



• The LR objective function  $-\log p(\mathbf{y}, \mathbf{w} | \mathbf{X}) = -\log p(\mathbf{y} | \mathbf{X}, \mathbf{w}) - \log p(\mathbf{w})$  can be written as

$$-\log\prod_{n=1}^{N}p(y_n|oldsymbol{x}_n,oldsymbol{w}) - \log p(oldsymbol{w}) = -\sum_{n=1}^{N}\log p(y_n|oldsymbol{x}_n,oldsymbol{w}) - \log p(oldsymbol{w})$$

• For the logistic regression model,  $p(y_n|\mathbf{x}_n, \mathbf{w}) = \mu_n^{y_n} (1 - \mu_n)^{1-y_n}$  where  $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$ 

• With a Gaussian prior  $p(w) = \mathcal{N}(w|0, \lambda^{-1}\mathbf{I}) \propto \exp(-\lambda w^{\top}w)$ , the gradient and Hessian will be



• The LR objective function  $-\log p(\mathbf{y}, \mathbf{w} | \mathbf{X}) = -\log p(\mathbf{y} | \mathbf{X}, \mathbf{w}) - \log p(\mathbf{w})$  can be written as

$$-\log\prod_{n=1}^{N}p(y_n|\boldsymbol{x}_n, \boldsymbol{w}) - \log p(\boldsymbol{w}) = -\sum_{n=1}^{N}\log p(y_n|\boldsymbol{x}_n, \boldsymbol{w}) - \log p(\boldsymbol{w})$$

• For the logistic regression model,  $p(y_n|\mathbf{x}_n, \mathbf{w}) = \mu_n^{y_n} (1 - \mu_n)^{1-y_n}$  where  $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$ 

• With a Gaussian prior  $p(w) = \mathcal{N}(w|0, \lambda^{-1}\mathbf{I}) \propto \exp(-\lambda w^{\top}w)$ , the gradient and Hessian will be

$$\boldsymbol{g} = -\sum_{n=1}^{N} (y_n - \mu_n) \boldsymbol{x}_n + \lambda \mathbf{I} \boldsymbol{w}$$

• The LR objective function  $-\log p(\mathbf{y}, \mathbf{w} | \mathbf{X}) = -\log p(\mathbf{y} | \mathbf{X}, \mathbf{w}) - \log p(\mathbf{w})$  can be written as

$$-\log\prod_{n=1}^{N}p(y_n|\boldsymbol{x}_n, \boldsymbol{w}) - \log p(\boldsymbol{w}) = -\sum_{n=1}^{N}\log p(y_n|\boldsymbol{x}_n, \boldsymbol{w}) - \log p(\boldsymbol{w})$$

• For the logistic regression model,  $p(y_n|\mathbf{x}_n, \mathbf{w}) = \mu_n^{y_n} (1 - \mu_n)^{1-y_n}$  where  $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$ 

• With a Gaussian prior  $p(w) = \mathcal{N}(w|0, \lambda^{-1}\mathbf{I}) \propto \exp(-\lambda w^{\top}w)$ , the gradient and Hessian will be

$$\boldsymbol{g} = -\sum_{n=1}^{N} (\boldsymbol{y}_n - \boldsymbol{\mu}_n) \boldsymbol{x}_n + \lambda \mathbf{I} \boldsymbol{w} = \mathbf{X}^{\top} (\boldsymbol{\mu} - \boldsymbol{y}) + \lambda \boldsymbol{w}$$

• The LR objective function  $-\log p(\mathbf{y}, \mathbf{w} | \mathbf{X}) = -\log p(\mathbf{y} | \mathbf{X}, \mathbf{w}) - \log p(\mathbf{w})$  can be written as

$$-\log\prod_{n=1}^{N}p(y_n|\boldsymbol{x}_n, \boldsymbol{w}) - \log p(\boldsymbol{w}) = -\sum_{n=1}^{N}\log p(y_n|\boldsymbol{x}_n, \boldsymbol{w}) - \log p(\boldsymbol{w})$$

• For the logistic regression model,  $p(y_n|\mathbf{x}_n, \mathbf{w}) = \mu_n^{y_n} (1 - \mu_n)^{1-y_n}$  where  $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$ 

• With a Gaussian prior  $p(w) = \mathcal{N}(w|0, \lambda^{-1}\mathbf{I}) \propto \exp(-\lambda w^{\top}w)$ , the gradient and Hessian will be

$$\boldsymbol{g} = -\sum_{n=1}^{N} (y_n - \mu_n) \boldsymbol{x}_n + \lambda \boldsymbol{\mathsf{I}} \boldsymbol{w} = \boldsymbol{\mathsf{X}}^{ op} (\boldsymbol{\mu} - \boldsymbol{y}) + \lambda \boldsymbol{w}$$
 (a  $D \times 1$  vector)

• The LR objective function  $-\log p(\mathbf{y}, \mathbf{w} | \mathbf{X}) = -\log p(\mathbf{y} | \mathbf{X}, \mathbf{w}) - \log p(\mathbf{w})$  can be written as

$$-\log\prod_{n=1}^N p(y_n|oldsymbol{x}_n,oldsymbol{w}) - \log p(oldsymbol{w}) = -\sum_{n=1}^N \log p(y_n|oldsymbol{x}_n,oldsymbol{w}) - \log p(oldsymbol{w})$$

• For the logistic regression model,  $p(y_n|\mathbf{x}_n, \mathbf{w}) = \mu_n^{y_n} (1 - \mu_n)^{1-y_n}$  where  $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$ 

• With a Gaussian prior  $p(w) = \mathcal{N}(w|0, \lambda^{-1}\mathbf{I}) \propto \exp(-\lambda w^{\top}w)$ , the gradient and Hessian will be

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

Prob. Mod. & Inference - CS698X (Piyush Rai, IITK)

• The LR objective function  $-\log p(\mathbf{y}, \mathbf{w} | \mathbf{X}) = -\log p(\mathbf{y} | \mathbf{X}, \mathbf{w}) - \log p(\mathbf{w})$  can be written as

$$-\log\prod_{n=1}^{N}p(y_n|\boldsymbol{x}_n, \boldsymbol{w}) - \log p(\boldsymbol{w}) = -\sum_{n=1}^{N}\log p(y_n|\boldsymbol{x}_n, \boldsymbol{w}) - \log p(\boldsymbol{w})$$

• For the logistic regression model,  $p(y_n|\mathbf{x}_n, \mathbf{w}) = \mu_n^{y_n} (1 - \mu_n)^{1-y_n}$  where  $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$ 

• With a Gaussian prior  $p(w) = \mathcal{N}(w|0, \lambda^{-1}\mathbf{I}) \propto \exp(-\lambda w^{\top}w)$ , the gradient and Hessian will be

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

Prob. Mod. & Inference - CS698X (Piyush Rai, IITK)

• The LR objective function  $-\log p(\mathbf{y}, \mathbf{w} | \mathbf{X}) = -\log p(\mathbf{y} | \mathbf{X}, \mathbf{w}) - \log p(\mathbf{w})$  can be written as

$$-\log\prod_{n=1}^{N}p(y_n|\boldsymbol{x}_n, \boldsymbol{w}) - \log p(\boldsymbol{w}) = -\sum_{n=1}^{N}\log p(y_n|\boldsymbol{x}_n, \boldsymbol{w}) - \log p(\boldsymbol{w})$$

• For the logistic regression model,  $p(y_n|\mathbf{x}_n, \mathbf{w}) = \mu_n^{y_n} (1 - \mu_n)^{1-y_n}$  where  $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$ 

• With a Gaussian prior  $p(w) = \mathcal{N}(w|0, \lambda^{-1}\mathbf{I}) \propto \exp(-\lambda w^{\top}w)$ , the gradient and Hessian will be

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

Prob. Mod. & Inference - CS698X (Piyush Rai, IITK)

• The LR objective function  $-\log p(\mathbf{y}, \mathbf{w} | \mathbf{X}) = -\log p(\mathbf{y} | \mathbf{X}, \mathbf{w}) - \log p(\mathbf{w})$  can be written as

$$-\log\prod_{n=1}^{N}p(y_n|\boldsymbol{x}_n, \boldsymbol{w}) - \log p(\boldsymbol{w}) = -\sum_{n=1}^{N}\log p(y_n|\boldsymbol{x}_n, \boldsymbol{w}) - \log p(\boldsymbol{w})$$

• For the logistic regression model,  $p(y_n | \mathbf{x}_n, \mathbf{w}) = \mu_n^{y_n} (1 - \mu_n)^{1 - y_n}$  where  $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$ 

• With a Gaussian prior  $p(w) = \mathcal{N}(w|0, \lambda^{-1}\mathbf{I}) \propto \exp(-\lambda w^{\top}w)$ , the gradient and Hessian will be

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

•  $\boldsymbol{\mu} = [\mu_1, \dots, \mu_N]^{\top}$  is  $N \times 1$  and **S** is a  $N \times N$  diagonal matrix with  $S_{nn} = \mu_n(1 - \mu_n)$ 

990 중 4 4 4 4 4 4 4 4 4 4

• When using MLE, the predictive distribution will be

$$p(y_* = 1 | \boldsymbol{x}_*, \boldsymbol{w}_{MLE}) = \sigma(\boldsymbol{w}_{MLE}^{\top} \boldsymbol{x}_*)$$



Prob. Mod. & Inference - CS698X (Piyush Rai, IITK)

• When using MLE, the predictive distribution will be

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{w}_{MLE}) = \sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*)$$
  
$$p(y_* | \mathbf{x}_*, \mathbf{w}_{MLE}) = \text{Bernoulli}(\sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*))$$

• When using MLE, the predictive distribution will be

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{w}_{MLE}) = \sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*)$$
  
$$p(y_* | \mathbf{x}_*, \mathbf{w}_{MLE}) = \text{Bernoulli}(\sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*))$$

• When using MAP, the predictive distribution will be

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{w}_{MAP}) = \sigma(\mathbf{w}_{MAP}^{\top} \mathbf{x}_*)$$
  
 $p(y_* | \mathbf{x}_*, \mathbf{w}_{MAP}) = Bernoulli(\sigma(\mathbf{w}_{MAP}^{\top} \mathbf{x}_*))$ 

(日) (四) (日) (日) (日)

• When using MLE, the predictive distribution will be

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{w}_{MLE}) = \sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*)$$
  
$$p(y_* | \mathbf{x}_*, \mathbf{w}_{MLE}) = \text{Bernoulli}(\sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*))$$

• When using MAP, the predictive distribution will be

$$p(y_* = 1 | \boldsymbol{x}_*, \boldsymbol{w}_{MAP}) = \sigma(\boldsymbol{w}_{MAP}^\top \boldsymbol{x}_*)$$
$$p(y_* | \boldsymbol{x}_*, \boldsymbol{w}_{MAP}) = \text{Bernoulli}(\sigma(\boldsymbol{w}_{MAP}^\top \boldsymbol{x}_*))$$

• When using Bayesian inference, the posterior predictive distribution, based on posterior averaging

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

(日) (명) (분) (분)

• When using MLE, the predictive distribution will be

$$p(y_* = 1 | \boldsymbol{x}_*, \boldsymbol{w}_{MLE}) = \sigma(\boldsymbol{w}_{MLE}^\top \boldsymbol{x}_*)$$
  
$$p(y_* | \boldsymbol{x}_*, \boldsymbol{w}_{MLE}) = \text{Bernoulli}(\sigma(\boldsymbol{w}_{MLE}^\top \boldsymbol{x}_*))$$

• When using MAP, the predictive distribution will be

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{w}_{MAP}) = \sigma(\mathbf{w}_{MAP}^\top \mathbf{x}_*)$$
  
$$p(y_* | \mathbf{x}_*, \mathbf{w}_{MAP}) = \text{Bernoulli}(\sigma(\mathbf{w}_{MAP}^\top \mathbf{x}_*))$$

• When using Bayesian inference, the posterior predictive distribution, based on posterior averaging

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

(日) (명) (분) (분)

• When using MLE, the predictive distribution will be

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{w}_{MLE}) = \sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*)$$
  
$$p(y_* | \mathbf{x}_*, \mathbf{w}_{MLE}) = \text{Bernoulli}(\sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*))$$

• When using MAP, the predictive distribution will be

$$p(y_* = 1 | \boldsymbol{x}_*, \boldsymbol{w}_{MAP}) = \sigma(\boldsymbol{w}_{MAP}^\top \boldsymbol{x}_*)$$
  
$$p(y_* | \boldsymbol{x}_*, \boldsymbol{w}_{MAP}) = \text{Bernoulli}(\sigma(\boldsymbol{w}_{MAP}^\top \boldsymbol{x}_*))$$

• When using Bayesian inference, the posterior predictive distribution, based on posterior averaging

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

• Above is hard in general. :-( If using the Laplace approximation for  $p(w|\mathbf{X}, y)$ , it will be

• When using MLE, the predictive distribution will be

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{w}_{MLE}) = \sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*)$$
  
$$p(y_* | \mathbf{x}_*, \mathbf{w}_{MLE}) = \text{Bernoulli}(\sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*))$$

• When using MAP, the predictive distribution will be

$$p(y_* = 1 | \boldsymbol{x}_*, \boldsymbol{w}_{MAP}) = \sigma(\boldsymbol{w}_{MAP}^\top \boldsymbol{x}_*)$$
  
$$p(y_* | \boldsymbol{x}_*, \boldsymbol{w}_{MAP}) = \text{Bernoulli}(\sigma(\boldsymbol{w}_{MAP}^\top \boldsymbol{x}_*))$$

• When using Bayesian inference, the posterior predictive distribution, based on posterior averaging

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

• Above is hard in general. :-( If using the Laplace approximation for  $p(w|\mathbf{X}, y)$ , it will be

$$p(y_* = 1 | \boldsymbol{x}_*, \boldsymbol{X}, \boldsymbol{y}) \approx \int \sigma(\boldsymbol{w}^\top \boldsymbol{x}_*) \mathcal{N}(\boldsymbol{w} | \boldsymbol{w}_{MAP}, \boldsymbol{H}^{-1}) d\boldsymbol{w}$$

Prob. Mod. & Inference - CS698X (Piyush Rai, IITK)

( ) < </p>

• When using MLE, the predictive distribution will be

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{w}_{MLE}) = \sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*)$$
  
$$p(y_* | \mathbf{x}_*, \mathbf{w}_{MLE}) = \text{Bernoulli}(\sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*))$$

• When using MAP, the predictive distribution will be

$$p(y_* = 1 | \boldsymbol{x}_*, \boldsymbol{w}_{MAP}) = \sigma(\boldsymbol{w}_{MAP}^\top \boldsymbol{x}_*)$$
  
$$p(y_* | \boldsymbol{x}_*, \boldsymbol{w}_{MAP}) = \text{Bernoulli}(\sigma(\boldsymbol{w}_{MAP}^\top \boldsymbol{x}_*))$$

• When using Bayesian inference, the posterior predictive distribution, based on posterior averaging

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

• Above is hard in general. :-( If using the Laplace approximation for  $p(w|\mathbf{X}, y)$ , it will be

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) \quad \approx \quad \int \sigma(\mathbf{w}^\top \mathbf{x}_*) \mathcal{N}(\mathbf{w} | \mathbf{w}_{MAP}, \mathbf{H}^{-1}) d\mathbf{w}$$

• Even after Laplace approximation for  $p(\boldsymbol{w}|\mathbf{X}, \boldsymbol{y})$ , the above integral to compute posterior predictive is intractable. So we will need to also approximate the predictive posterior. :-)

#### **Posterior Predictive via Monte-Carlo Sampling**

• The posterior predictive is given by the following integral

$$p(y_* = 1 | \boldsymbol{x}_*, \boldsymbol{X}, \boldsymbol{y}) = \int \sigma(\boldsymbol{w}^\top \boldsymbol{x}_*) \mathcal{N}(\boldsymbol{w} | \boldsymbol{w}_{MAP}, \boldsymbol{H}^{-1}) d\boldsymbol{w}$$



イロト イロト イヨト イヨト

#### **Posterior Predictive via Monte-Carlo Sampling**

• The posterior predictive is given by the following integral

$$p(y_* = 1 | \boldsymbol{x}_*, \boldsymbol{X}, \boldsymbol{y}) = \int \sigma(\boldsymbol{w}^\top \boldsymbol{x}_*) \mathcal{N}(\boldsymbol{w} | \boldsymbol{w}_{MAP}, \boldsymbol{H}^{-1}) d\boldsymbol{w}$$

Monte-Carlo approximation: Draw several samples of *w* from *N*(*w*|*w*<sub>MAP</sub>, H<sup>-1</sup>) and replace the above integral by an empirical average of *σ*(*w*<sup>T</sup>*x*<sub>\*</sub>) computed using each of those samples

$$p(y_* = 1 | \boldsymbol{x}_*, \boldsymbol{X}, \boldsymbol{y}) \approx \frac{1}{S} \sum_{s=1}^{S} \sigma(\boldsymbol{w}_s^{\top} \boldsymbol{x}_*)$$

where  $\boldsymbol{w}_{s} \sim \mathcal{N}(\boldsymbol{w} | \boldsymbol{w}_{MAP}, \mathbf{H}^{-1})$ ,  $s = 1, \dots, S$ 

### **Posterior Predictive via Monte-Carlo Sampling**

• The posterior predictive is given by the following integral

$$p(y_* = 1 | \boldsymbol{x}_*, \boldsymbol{X}, \boldsymbol{y}) = \int \sigma(\boldsymbol{w}^\top \boldsymbol{x}_*) \mathcal{N}(\boldsymbol{w} | \boldsymbol{w}_{MAP}, \boldsymbol{H}^{-1}) d\boldsymbol{w}$$

Monte-Carlo approximation: Draw several samples of *w* from *N*(*w*|*w*<sub>MAP</sub>, H<sup>-1</sup>) and replace the above integral by an empirical average of *σ*(*w*<sup>T</sup>*x*<sub>\*</sub>) computed using each of those samples

$$p(y_* = 1 | \boldsymbol{x}_*, \boldsymbol{X}, \boldsymbol{y}) \approx \frac{1}{S} \sum_{s=1}^{S} \sigma(\boldsymbol{w}_s^{\top} \boldsymbol{x}_*)$$

where  $m{w}_{s} \sim \mathcal{N}(m{w} | m{w}_{MAP}, m{H}^{-1})$ ,  $s = 1, \dots, S$ 

More on Monte-Carlo methods when we discuss MCMC sampling

• The posterior predictive we wanted to compute was

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) \approx \int \sigma(\mathbf{w}^\top \mathbf{x}_*) \mathcal{N}(\mathbf{w} | \mathbf{w}_{MAP}, \mathbf{H}^{-1}) d\mathbf{w}$$



< ロ > < 同 > < 三 > < 三 >

• The posterior predictive we wanted to compute was

$$p(y_* = 1 | \boldsymbol{x}_*, \boldsymbol{X}, \boldsymbol{y}) \approx \int \sigma(\boldsymbol{w}^\top \boldsymbol{x}_*) \mathcal{N}(\boldsymbol{w} | \boldsymbol{w}_{MAP}, \boldsymbol{H}^{-1}) d\boldsymbol{w}$$

• In the above, let's replace the sigmoid  $\sigma(\boldsymbol{w}^{\top}\boldsymbol{x}_{*})$  by  $\Phi(\boldsymbol{w}^{\top}\boldsymbol{x}_{*})$ , i.e., CDF of standard normal

$$\Phi(z)=rac{1}{\sqrt{2\pi}}\int_{-\infty}^{z}e^{-t^{2}}dt$$
 (Note:  $z$  is a scalar and  $0\leq\Phi(z)\leq1)$ 

• The posterior predictive we wanted to compute was

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) \approx \int \sigma(\mathbf{w}^\top \mathbf{x}_*) \mathcal{N}(\mathbf{w} | \mathbf{w}_{MAP}, \mathbf{H}^{-1}) d\mathbf{w}$$

• In the above, let's replace the sigmoid  $\sigma(\boldsymbol{w}^{\top}\boldsymbol{x}_{*})$  by  $\Phi(\boldsymbol{w}^{\top}\boldsymbol{x}_{*})$ , i.e., CDF of standard normal

$$\Phi(z)=rac{1}{\sqrt{2\pi}}\int_{-\infty}^{z}e^{-t^{2}}dt$$
 (Note:  $z$  is a scalar and  $0\leq\Phi(z)\leq1$ )

• Note:  $\Phi(z)$  is also called the probit function



• The posterior predictive we wanted to compute was

$$p(y_* = 1 | \boldsymbol{x}_*, \boldsymbol{X}, \boldsymbol{y}) \approx \int \sigma(\boldsymbol{w}^\top \boldsymbol{x}_*) \mathcal{N}(\boldsymbol{w} | \boldsymbol{w}_{MAP}, \boldsymbol{H}^{-1}) d\boldsymbol{w}$$

• In the above, let's replace the sigmoid  $\sigma(\boldsymbol{w}^{\top}\boldsymbol{x}_{*})$  by  $\Phi(\boldsymbol{w}^{\top}\boldsymbol{x}_{*})$ , i.e., CDF of standard normal

$$\Phi(z)=rac{1}{\sqrt{2\pi}}\int_{-\infty}^{z}e^{-t^{2}}dt$$
 (Note:  $z$  is a scalar and  $0\leq\Phi(z)\leq1)$ 

• Note:  $\Phi(z)$  is also called the probit function



< ∃ > < ∃

• The posterior predictive we wanted to compute was

$$p(y_* = 1 | \boldsymbol{x}_*, \boldsymbol{X}, \boldsymbol{y}) \approx \int \sigma(\boldsymbol{w}^\top \boldsymbol{x}_*) \mathcal{N}(\boldsymbol{w} | \boldsymbol{w}_{MAP}, \boldsymbol{H}^{-1}) d\boldsymbol{w}$$

• In the above, let's replace the sigmoid  $\sigma(\boldsymbol{w}^{\top}\boldsymbol{x}_{*})$  by  $\Phi(\boldsymbol{w}^{\top}\boldsymbol{x}_{*})$ , i.e., CDF of standard normal

$$\Phi(z)=rac{1}{\sqrt{2\pi}}\int_{-\infty}^{z}e^{-t^{2}}dt$$
 (Note:  $z$  is a scalar and  $0\leq\Phi(z)\leq1)$ 

• Note:  $\Phi(z)$  is also called the probit function



• This approach relies on numerical approximation (as we will see)

-

• With this approximation, the predictive posterior will be

$$p(y_* = 1 | \boldsymbol{x}_*, \boldsymbol{X}, \boldsymbol{y}) = \int \Phi(\boldsymbol{w}^\top \boldsymbol{x}_*) \mathcal{N}(\boldsymbol{w} | \boldsymbol{w}_{MAP}, \boldsymbol{\mathsf{H}}^{-1}) d\boldsymbol{w}$$
 (an expectation)



イロト イボト イモト イモン

• With this approximation, the predictive posterior will be

$$\begin{split} p(y_* = 1 | \boldsymbol{x}_*, \boldsymbol{X}, \boldsymbol{y}) &= \int \Phi(\boldsymbol{w}^\top \boldsymbol{x}_*) \mathcal{N}(\boldsymbol{w} | \boldsymbol{w}_{MAP}, \boldsymbol{H}^{-1}) d\boldsymbol{w} \quad \text{(an expectation)} \\ &= \int_{-\infty}^{\infty} \Phi(a) p(a | \mu_a, \sigma_a^2) da \quad \text{(an equivalent expectation)} \end{split}$$



イロト イロト イヨト イヨト

• With this approximation, the predictive posterior will be

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int \Phi(\mathbf{w}^\top \mathbf{x}_*) \mathcal{N}(\mathbf{w} | \mathbf{w}_{MAP}, \mathbf{H}^{-1}) d\mathbf{w} \quad \text{(an expectation)}$$
$$= \int_{-\infty}^{\infty} \Phi(a) p(a | \mu_a, \sigma_a^2) da \quad \text{(an equivalent expectation)}$$

• Since  $\mathbf{a} = \mathbf{w}^{\top} \mathbf{x}_{*} = \mathbf{x}_{*}^{\top} \mathbf{w}$ , and  $\mathbf{w}$  is normally distributed,  $p(a|\mu_{a},\sigma_{a}^{2}) = \mathcal{N}(a|\mu_{a},\sigma_{a}^{2})$ , with  $\mu_{a} = \mathbf{w}_{MAP}^{\top} \mathbf{x}_{*}$  and  $\sigma_{a}^{2} = \mathbf{x}_{*}^{\top} \mathbf{H}^{-1} \mathbf{x}_{*}$  (follows from the linear trans. property of random vars)



• With this approximation, the predictive posterior will be

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int \Phi(\mathbf{w}^\top \mathbf{x}_*) \mathcal{N}(\mathbf{w} | \mathbf{w}_{MAP}, \mathbf{H}^{-1}) d\mathbf{w} \quad \text{(an expectation)}$$
$$= \int_{-\infty}^{\infty} \Phi(a) p(a | \mu_a, \sigma_a^2) da \quad \text{(an equivalent expectation)}$$

• Since  $\mathbf{a} = \mathbf{w}^{\top} \mathbf{x}_{*} = \mathbf{x}_{*}^{\top} \mathbf{w}$ , and  $\mathbf{w}$  is normally distributed,  $p(a|\mu_{a}, \sigma_{a}^{2}) = \mathcal{N}(a|\mu_{a}, \sigma_{a}^{2})$ , with  $\mu_{a} = \mathbf{w}_{MAP}^{\top} \mathbf{x}_{*}$  and  $\sigma_{a}^{2} = \mathbf{x}_{*}^{\top} \mathbf{H}^{-1} \mathbf{x}_{*}$  (follows from the linear trans. property of random vars)

• Given  $\mu_a = \mathbf{w}_{MAP}^{\top} \mathbf{x}_*$  and  $\sigma_a^2 = \mathbf{x}_*^{\top} \mathbf{H}^{-1} \mathbf{x}_*$ , the predictive posterior will be

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) \approx \int_{-\infty}^{\infty} \Phi(a) \mathcal{N}(a | \mu_a, \sigma_a^2) da = \Phi\left(\frac{\mu_a}{\sqrt{1 + \sigma_a^2}}\right)$$

Prob. Mod. & Inference - CS698X (Piyush Rai, IITK)

A B > A B > A B >
 A
 B >
 A
 B > A B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
#### **Predictive Posterior via Probit Approximation**

• With this approximation, the predictive posterior will be

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int \Phi(\mathbf{w}^\top \mathbf{x}_*) \mathcal{N}(\mathbf{w} | \mathbf{w}_{MAP}, \mathbf{H}^{-1}) d\mathbf{w} \quad \text{(an expectation)}$$
$$= \int_{-\infty}^{\infty} \Phi(a) p(a | \mu_a, \sigma_a^2) da \quad \text{(an equivalent expectation)}$$

• Since  $\mathbf{a} = \mathbf{w}^{\top} \mathbf{x}_{*} = \mathbf{x}_{*}^{\top} \mathbf{w}$ , and  $\mathbf{w}$  is normally distributed,  $p(a|\mu_{a}, \sigma_{a}^{2}) = \mathcal{N}(a|\mu_{a}, \sigma_{a}^{2})$ , with  $\mu_{a} = \mathbf{w}_{MAP}^{\top} \mathbf{x}_{*}$  and  $\sigma_{a}^{2} = \mathbf{x}_{*}^{\top} \mathbf{H}^{-1} \mathbf{x}_{*}$  (follows from the linear trans. property of random vars)

• Given  $\mu_a = \mathbf{w}_{MAP}^{\top} \mathbf{x}_*$  and  $\sigma_a^2 = \mathbf{x}_*^{\top} \mathbf{H}^{-1} \mathbf{x}_*$ , the predictive posterior will be

$$p(y_* = 1 | \boldsymbol{x}_*, \boldsymbol{X}, \boldsymbol{y}) pprox \int_{-\infty}^{\infty} \Phi(\boldsymbol{a}) \mathcal{N}(\boldsymbol{a} | \mu_{\boldsymbol{a}}, \sigma_{\boldsymbol{a}}^2) d\boldsymbol{a} = \Phi\left(rac{\mu_{\boldsymbol{a}}}{\sqrt{1 + \sigma_{\boldsymbol{a}}^2}}
ight)$$

• Note that the variance  $\sigma_a^2$  also "moderates" the probability of  $y_n$  being 1 (MAP would give  $\Phi(\mu_a)$ )

#### **Predictive Posterior via Probit Approximation**

• With this approximation, the predictive posterior will be

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int \Phi(\mathbf{w}^\top \mathbf{x}_*) \mathcal{N}(\mathbf{w} | \mathbf{w}_{MAP}, \mathbf{H}^{-1}) d\mathbf{w} \quad \text{(an expectation)}$$
$$= \int_{-\infty}^{\infty} \Phi(a) p(a | \mu_a, \sigma_a^2) da \quad \text{(an equivalent expectation)}$$

• Since  $\mathbf{a} = \mathbf{w}^{\top} \mathbf{x}_{*} = \mathbf{x}_{*}^{\top} \mathbf{w}$ , and  $\mathbf{w}$  is normally distributed,  $p(\mathbf{a}|\mu_{a}, \sigma_{a}^{2}) = \mathcal{N}(\mathbf{a}|\mu_{a}, \sigma_{a}^{2})$ , with  $\mu_{a} = \mathbf{w}_{MAP}^{\top} \mathbf{x}_{*}$  and  $\sigma_{a}^{2} = \mathbf{x}_{*}^{\top} \mathbf{H}^{-1} \mathbf{x}_{*}$  (follows from the linear trans. property of random vars)

• Given  $\mu_a = \mathbf{w}_{MAP}^{\top} \mathbf{x}_*$  and  $\sigma_a^2 = \mathbf{x}_*^{\top} \mathbf{H}^{-1} \mathbf{x}_*$ , the predictive posterior will be

$$p(y_* = 1 | \boldsymbol{x}_*, \boldsymbol{X}, \boldsymbol{y}) pprox \int_{-\infty}^{\infty} \Phi(\boldsymbol{a}) \mathcal{N}(\boldsymbol{a} | \mu_{\boldsymbol{a}}, \sigma_{\boldsymbol{a}}^2) d\boldsymbol{a} = \Phi\left(rac{\mu_{\boldsymbol{a}}}{\sqrt{1 + \sigma_{\boldsymbol{a}}^2}}
ight)$$

• Note that the variance  $\sigma_a^2$  also "moderates" the probability of  $y_n$  being 1 (MAP would give  $\Phi(\mu_a)$ )

• Since logistic and probit aren't exactly identical, we usually scale *a* by a scalar *t* s.t.  $t^2 = \pi/8$  $p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int_{-\infty}^{\infty} \Phi(ta) \mathcal{N}(a | \mu_a, \sigma_a^2) da = \Phi\left(\frac{\mu_a}{\sqrt{t^{-2} + \sigma_a^2}}\right) da = \Phi\left(\frac{\mu_a}{\sqrt{t^{-2} + \sigma_a^2}}\right) da$ 

#### **Bayesian Logistic Regression: Posterior over Linear Classifiers!**



Figure courtesy: MLAPP (Murphy)

DOC E

## Logistic Regression: Plug-in Prediction vs Bayesian Averaging

- (Left) Predictive distribution when using a point estimate uses only a single linear hyperplane  $\boldsymbol{w}$
- (Right) Posterior predictive distribution averages over many linear hyperplanes  ${m w}$



• We saw basic logistic regression model and some ways to perform Bayesian inference for this model

<ロト < 四ト < 三ト < 三ト

- We saw basic logistic regression model and some ways to perform Bayesian inference for this model
  - We assumed the hyperparameters (e.g., precision/variance of  $p(w) = \mathcal{N}(0, \lambda^{-1}I)$ ) to be fixed. However, these can also be learned if desired

- We saw basic logistic regression model and some ways to perform Bayesian inference for this model
  - We assumed the hyperparameters (e.g., precision/variance of p(w) = N(0, λ<sup>-1</sup>I)) to be fixed. However, these can also be learned if desired
  - LR is a linear classification model. Can be extended to nonlinear classification (more on this later)

- We saw basic logistic regression model and some ways to perform Bayesian inference for this model
  - We assumed the hyperparameters (e.g., precision/variance of p(w) = N(0, λ<sup>-1</sup>I)) to be fixed. However, these can also be learned if desired
  - LR is a linear classification model. Can be extended to nonlinear classification (more on this later)
- Logistic Regression (and its Bayesian version) is widely used in probabilistic classification

- We saw basic logistic regression model and some ways to perform Bayesian inference for this model
  - We assumed the hyperparameters (e.g., precision/variance of p(w) = N(0, λ<sup>-1</sup>I)) to be fixed. However, these can also be learned if desired
  - LR is a linear classification model. Can be extended to nonlinear classification (more on this later)
- Logistic Regression (and its Bayesian version) is widely used in probabilistic classification
- Its multiclass extension is softmax regression (which again can be treated in a Bayesian manner)

- We saw basic logistic regression model and some ways to perform Bayesian inference for this model
  - We assumed the hyperparameters (e.g., precision/variance of p(w) = N(0, λ<sup>-1</sup>I)) to be fixed. However, these can also be learned if desired
  - LR is a linear classification model. Can be extended to nonlinear classification (more on this later)
- Logistic Regression (and its Bayesian version) is widely used in probabilistic classification
- Its multiclass extension is softmax regression (which again can be treated in a Bayesian manner)
- LR and softmax some of the simplest models for discriminative classification but non-conjugate

- We saw basic logistic regression model and some ways to perform Bayesian inference for this model
  - We assumed the hyperparameters (e.g., precision/variance of p(w) = N(0, λ<sup>-1</sup>I)) to be fixed. However, these can also be learned if desired
  - LR is a linear classification model. Can be extended to nonlinear classification (more on this later)
- Logistic Regression (and its Bayesian version) is widely used in probabilistic classification
- Its multiclass extension is softmax regression (which again can be treated in a Bayesian manner)
- LR and softmax some of the simplest models for discriminative classification but non-conjugate
- The Laplace approximation is one of the simplest approximations to handle non-conjugacy

- We saw basic logistic regression model and some ways to perform Bayesian inference for this model
  - We assumed the hyperparameters (e.g., precision/variance of p(w) = N(0, λ<sup>-1</sup>I)) to be fixed. However, these can also be learned if desired
  - LR is a linear classification model. Can be extended to nonlinear classification (more on this later)
- Logistic Regression (and its Bayesian version) is widely used in probabilistic classification
- Its multiclass extension is softmax regression (which again can be treated in a Bayesian manner)
- LR and softmax some of the simplest models for discriminative classification but non-conjugate
- The Laplace approximation is one of the simplest approximations to handle non-conjugacy
- A variety of other approximate inference algorithms exist for these models

E Dac

イロト 不得 トイヨト イヨト

- We saw basic logistic regression model and some ways to perform Bayesian inference for this model
  - We assumed the hyperparameters (e.g., precision/variance of p(w) = N(0, λ<sup>-1</sup>I)) to be fixed. However, these can also be learned if desired
  - LR is a linear classification model. Can be extended to nonlinear classification (more on this later)
- Logistic Regression (and its Bayesian version) is widely used in probabilistic classification
- Its multiclass extension is softmax regression (which again can be treated in a Bayesian manner)
- LR and softmax some of the simplest models for discriminative classification but non-conjugate
- The Laplace approximation is one of the simplest approximations to handle non-conjugacy
- A variety of other approximate inference algorithms exist for these models
  - We will revisit LR when discussing such approximate inference methods

E DOO

・ロト ・日下・ ・日下・

# Bayesian Generative Classification

• Consider N labeled examples  $\{(\mathbf{x}_i, y_i)\}_{n=1}^N$ . Assume binary labels, i.e.,  $y_i \in \{0, 1\}$ 



A B > A B > A B >
 A
 B >
 A
 B > A B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A

• Consider N labeled examples  $\{(\mathbf{x}_i, y_i)\}_{n=1}^N$ . Assume binary labels, i.e.,  $y_i \in \{0, 1\}$ 

• Goal: Classify a new example  $\boldsymbol{x}$  by assigning a label  $y \in \{0,1\}$  to it

<ロト < 四ト < 三ト < 三ト

- Consider N labeled examples  $\{(\mathbf{x}_i, y_i)\}_{n=1}^N$ . Assume binary labels, i.e.,  $y_i \in \{0, 1\}$
- ${\ }$  Goal: Classify a new example  ${\it {\it x}}$  by assigning a label  ${\it y}\in\{0,1\}$  to it
- We will assume a Generative Model for both labels y and and features x



< □ > < □ > < Ξ > < Ξ > < Ξ >

- Consider N labeled examples  $\{(\mathbf{x}_i, y_i)\}_{n=1}^N$ . Assume binary labels, i.e.,  $y_i \in \{0, 1\}$
- $\, \circ \,$  Goal: Classify a new example  ${\it x}$  by assigning a label  $y \in \{0,1\}$  to it
- We will assume a Generative Model for both labels y and and features x
  - What it means: We will have (probabilistic) observation models for both y as well as x

< □ > < □ > < Ξ > < Ξ > < Ξ >

- Consider N labeled examples  $\{(\mathbf{x}_i, y_i)\}_{n=1}^N$ . Assume binary labels, i.e.,  $y_i \in \{0, 1\}$
- $\, \circ \,$  Goal: Classify a new example  ${\it x}$  by assigning a label  $y \in \{0,1\}$  to it
- We will assume a Generative Model for both labels y and and features x
  - What it means: We will have (probabilistic) observation models for both y as well as x
  - In contrast, in Bayesian linear regression model (and Bayesian logistic regression model), we didn't model x (there, we simply conditioned y on x, treating x as "fixed")

- Consider N labeled examples  $\{(\mathbf{x}_i, y_i)\}_{n=1}^N$ . Assume binary labels, i.e.,  $y_i \in \{0, 1\}$
- $\, \circ \,$  Goal: Classify a new example  ${\it x}$  by assigning a label  $y \in \{0,1\}$  to it
- We will assume a Generative Model for both labels y and and features x
  - What it means: We will have (probabilistic) observation models for both y as well as x
  - In contrast, in Bayesian linear regression model (and Bayesian logistic regression model), we didn't model x (there, we simply conditioned y on x, treating x as "fixed")
  - When we don't model x and simply model y as a function of x: Discriminative Model

- Consider N labeled examples  $\{(\mathbf{x}_i, y_i)\}_{n=1}^N$ . Assume binary labels, i.e.,  $y_i \in \{0, 1\}$
- $\, \circ \,$  Goal: Classify a new example  ${\it x}$  by assigning a label  $y \in \{0,1\}$  to it
- We will assume a Generative Model for both labels y and and features x
  - What it means: We will have (probabilistic) observation models for both y as well as x
  - In contrast, in Bayesian linear regression model (and Bayesian logistic regression model), we didn't model x (there, we simply conditioned y on x, treating x as "fixed")
  - When we don't model x and simply model y as a function of x: Discriminative Model
- Generative classification models have many benefits. E.g.,

- Consider N labeled examples  $\{(\mathbf{x}_i, y_i)\}_{n=1}^N$ . Assume binary labels, i.e.,  $y_i \in \{0, 1\}$
- $\, \circ \,$  Goal: Classify a new example  ${\it x}$  by assigning a label  $y \in \{0,1\}$  to it
- We will assume a Generative Model for both labels y and and features x
  - What it means: We will have (probabilistic) observation models for both y as well as x
  - In contrast, in Bayesian linear regression model (and Bayesian logistic regression model), we didn't model x (there, we simply conditioned y on x, treating x as "fixed")
  - When we don't model x and simply model y as a function of x: Discriminative Model
- Generative classification models have many benefits. E.g.,
  - Can also utilize unlabeled examples (semi-supervised learning)

- Consider N labeled examples  $\{(\mathbf{x}_i, y_i)\}_{n=1}^N$ . Assume binary labels, i.e.,  $y_i \in \{0, 1\}$
- $\, \circ \,$  Goal: Classify a new example  ${\it x}$  by assigning a label  $y \in \{0,1\}$  to it
- We will assume a Generative Model for both labels y and and features x
  - What it means: We will have (probabilistic) observation models for both y as well as x
  - In contrast, in Bayesian linear regression model (and Bayesian logistic regression model), we didn't model x (there, we simply conditioned y on x, treating x as "fixed")
  - When we don't model x and simply model y as a function of x: Discriminative Model
- Generative classification models have many benefits. E.g.,
  - Can also utilize unlabeled examples (semi-supervised learning)
  - Can handle missing/corrupted features in  $\boldsymbol{x}$

・ロト ・四ト ・ヨト ・ヨト

- Consider N labeled examples  $\{(\mathbf{x}_i, y_i)\}_{n=1}^N$ . Assume binary labels, i.e.,  $y_i \in \{0, 1\}$
- $\, \circ \,$  Goal: Classify a new example  ${\it x}$  by assigning a label  $y \in \{0,1\}$  to it
- We will assume a Generative Model for both labels y and and features x
  - What it means: We will have (probabilistic) observation models for both y as well as x
  - In contrast, in Bayesian linear regression model (and Bayesian logistic regression model), we didn't model x (there, we simply conditioned y on x, treating x as "fixed")
  - When we don't model x and simply model y as a function of x: Discriminative Model
- Generative classification models have many benefits. E.g.,
  - Can also utilize unlabeled examples (semi-supervised learning)
  - Can handle missing/corrupted features in x
  - Can properly handle cases when features in x could be of mixed type (e.g., real, binary, count)

- Consider N labeled examples  $\{(\mathbf{x}_i, y_i)\}_{n=1}^N$ . Assume binary labels, i.e.,  $y_i \in \{0, 1\}$
- $\, \circ \,$  Goal: Classify a new example  ${\it x}$  by assigning a label  $y \in \{0,1\}$  to it
- We will assume a Generative Model for both labels y and and features x
  - What it means: We will have (probabilistic) observation models for both y as well as x
  - In contrast, in Bayesian linear regression model (and Bayesian logistic regression model), we didn't model x (there, we simply conditioned y on x, treating x as "fixed")
  - When we don't model x and simply model y as a function of x: Discriminative Model
- Generative classification models have many benefits. E.g.,
  - Can also utilize unlabeled examples (semi-supervised learning)
  - Can handle missing/corrupted features in x
  - Can properly handle cases when features in x could be of mixed type (e.g., real, binary, count)
  - And many others (more on this later during the semester)

• Basic idea: Each  $x_i$  is assumed generated conditioned on the value of corresponding label  $y_i$ 

・ロト ・四ト ・モト ・モト

- Basic idea: Each  $x_i$  is assumed generated conditioned on the value of corresponding label  $y_i$
- The associated generative story is as follows



- Basic idea: Each  $x_i$  is assumed generated conditioned on the value of corresponding label  $y_i$
- The associated generative story is as follows
  - First draw ("generate") a binary label  $y_i \in \{0,1\}$

 $y_i \sim \mathsf{Bernoulli}(\pi)$ 

・ロト ・四ト ・モト ・モト

- Basic idea: Each  $x_i$  is assumed generated conditioned on the value of corresponding label  $y_i$
- The associated generative story is as follows
  - First draw ("generate") a binary label  $y_i \in \{0, 1\}$

 $y_i \sim \mathsf{Bernoulli}(\pi)$ 

• Now draw ("generate") the feature vector x from a distribution specific to the value  $y_i$  takes

 $oldsymbol{x}_i|y_i \sim p(oldsymbol{x}| heta_{y_i})$ 



- Basic idea: Each  $x_i$  is assumed generated conditioned on the value of corresponding label  $y_i$
- The associated generative story is as follows
  - First draw ("generate") a binary label  $y_i \in \{0,1\}$

 $y_i \sim \mathsf{Bernoulli}(\pi)$ 

• Now draw ("generate") the feature vector x from a distribution specific to the value  $y_i$  takes

$$|\mathbf{x}_i|y_i \sim p(\mathbf{x}| heta_{y_i})$$

• The above generative model shown in "plate notation" (shaded = observed)



• Our generative model for classification is

$$y_i \sim \mathsf{Bernoulli}(\pi), \qquad oldsymbol{x}_i | y_i \sim p(oldsymbol{x} | heta_{y_i})$$

• Note: We have two distributions  $p(x|\theta_0)$  and  $p(x|\theta_1)$  for feature vector x (depending on its label)

• Our generative model for classification is

$$y_i \sim \mathsf{Bernoulli}(\pi), \qquad oldsymbol{x}_i | y_i \sim p(oldsymbol{x} | heta_{y_i})$$

• Note: We have two distributions  $p(x|\theta_0)$  and  $p(x|\theta_1)$  for feature vector x (depending on its label)

• These distributions are also known as "class-conditional distributions"

• Our generative model for classification is

$$y_i \sim \text{Bernoulli}(\pi), \qquad oldsymbol{x}_i | y_i \sim p(oldsymbol{x} | heta_{y_i})$$

• Note: We have two distributions  $p(\mathbf{x}|\theta_0)$  and  $p(\mathbf{x}|\theta_1)$  for feature vector  $\mathbf{x}$  (depending on its label)

- These distributions are also known as "class-conditional distributions"
- For now, we will not assume any specific form for the distriutions  $p(\mathbf{x}|\theta_0)$  and  $p(\mathbf{x}|\theta_1)$

<ロト < 四ト < 三ト < 三ト

• Our generative model for classification is

$$y_i \sim \mathsf{Bernoulli}(\pi), \qquad oldsymbol{x}_i | y_i \sim p(oldsymbol{x} | heta_{y_i})$$

• Note: We have two distributions  $p(\mathbf{x}|\theta_0)$  and  $p(\mathbf{x}|\theta_1)$  for feature vector  $\mathbf{x}$  (depending on its label)

- These distributions are also known as "class-conditional distributions"
- For now, we will not assume any specific form for the distriutions  $p(\mathbf{x}|\theta_0)$  and  $p(\mathbf{x}|\theta_1)$ 
  - Depends on nature of x (real-valued vectors? binary vectors? count vectors?)

・ロト ・四ト ・モト ・モト

• Our generative model for classification is

$$y_i \sim \mathsf{Bernoulli}(\pi), \qquad oldsymbol{x}_i | y_i \sim p(oldsymbol{x} | heta_{y_i})$$

• Note: We have two distributions  $p(\mathbf{x}|\theta_0)$  and  $p(\mathbf{x}|\theta_1)$  for feature vector  $\mathbf{x}$  (depending on its label)

- These distributions are also known as "class-conditional distributions"
- For now, we will not assume any specific form for the distriutions  $p(\mathbf{x}|\theta_0)$  and  $p(\mathbf{x}|\theta_1)$ 
  - Depends on nature of x (real-valued vectors? binary vectors? count vectors?)
- Model parameters to be learned here:  $(\pi, \theta_0, \theta_1)$

• Our generative model for classification is

$$y_i \sim \text{Bernoulli}(\pi), \qquad \mathbf{x}_i | y_i \sim p(\mathbf{x} | \theta_{y_i})$$

• Note: We have two distributions  $p(\mathbf{x}|\theta_0)$  and  $p(\mathbf{x}|\theta_1)$  for feature vector  $\mathbf{x}$  (depending on its label)

- These distributions are also known as "class-conditional distributions"
- For now, we will not assume any specific form for the distributions  $p(\mathbf{x}|\theta_0)$  and  $p(\mathbf{x}|\theta_1)$ 
  - Depends on nature of x (real-valued vectors? binary vectors? count vectors?)
- Model parameters to be learned here:  $(\pi, \theta_0, \theta_1)$
- Note: Can extend to more than 2 classes (e.g., by replacing the Bernoulli on y by multinoulli)

・ロト ・日 ・ ・ ヨ ・ ・ ヨ ・
## A Generative Model for Classification

• Our generative model for classification is

$$y_i \sim \text{Bernoulli}(\pi), \qquad \mathbf{x}_i | y_i \sim p(\mathbf{x} | \theta_{y_i})$$

• Note: We have two distributions  $p(\mathbf{x}|\theta_0)$  and  $p(\mathbf{x}|\theta_1)$  for feature vector  $\mathbf{x}$  (depending on its label)

- These distributions are also known as "class-conditional distributions"
- For now, we will not assume any specific form for the distributions  $p(\mathbf{x}|\theta_0)$  and  $p(\mathbf{x}|\theta_1)$ 
  - Depends on nature of x (real-valued vectors? binary vectors? count vectors?)
- Model parameters to be learned here:  $(\pi, \theta_0, \theta_1)$
- Note: Can extend to more than 2 classes (e.g., by replacing the Bernoulli on y by multinoulli)
- Note: When  $y_i$  for each  $x_i$  is a hidden variable, we can think of it as the cluster id of x

イロト イロト イヨト イヨト

## A Generative Model for Classification

• Our generative model for classification is

$$y_i \sim \text{Bernoulli}(\pi), \qquad \mathbf{x}_i | y_i \sim p(\mathbf{x} | \theta_{y_i})$$

• Note: We have two distributions  $p(\mathbf{x}|\theta_0)$  and  $p(\mathbf{x}|\theta_1)$  for feature vector  $\mathbf{x}$  (depending on its label)

- These distributions are also known as "class-conditional distributions"
- For now, we will not assume any specific form for the distributions  $p(\mathbf{x}|\theta_0)$  and  $p(\mathbf{x}|\theta_1)$ 
  - Depends on nature of x (real-valued vectors? binary vectors? count vectors?)
- Model parameters to be learned here:  $(\pi, \theta_0, \theta_1)$
- Note: Can extend to more than 2 classes (e.g., by replacing the Bernoulli on y by multinoulli)
- Note: When  $y_i$  for each  $x_i$  is a hidden variable, we can think of it as the cluster id of x
  - It then becomes a mixture model based data clustering problem (unsupervised learning)

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・

• Note: The generative model only defines  $p(y|\pi)$  and  $p(\mathbf{x}|\theta_y)$ . Doesn't define  $p(y|\mathbf{x})$ 



- Note: The generative model only defines  $p(y|\pi)$  and  $p(\mathbf{x}|\theta_y)$ . Doesn't define  $p(y|\mathbf{x})$
- We combine these using Bayes rule to get p(y|x)

$$p(y|\mathbf{x}) = \frac{p(y|\pi)p(\mathbf{x}|\theta_y)}{p(\mathbf{x})} = \frac{p(y|\pi)p(\mathbf{x}|\theta_y)}{\sum_y p(y|\pi)p(\mathbf{x}|\theta_y)}$$

- Note: The generative model only defines  $p(y|\pi)$  and  $p(x|\theta_y)$ . Doesn't define p(y|x)
- We combine these using Bayes rule to get  $p(y|\mathbf{x})$

$$p(y|\mathbf{x}) = \frac{p(y|\pi)p(\mathbf{x}|\theta_y)}{p(\mathbf{x})} = \frac{p(y|\pi)p(\mathbf{x}|\theta_y)}{\sum_y p(y|\pi)p(\mathbf{x}|\theta_y)}$$

• Parameters of distributions  $p(y|\pi)$  and  $p(\mathbf{x}|\theta_y)$  are estimated from training data using point estimation methods (MLE or MAP) or using fully Bayesian inference (discussed today)

- Note: The generative model only defines  $p(y|\pi)$  and  $p(x|\theta_y)$ . Doesn't define p(y|x)
- We combine these using Bayes rule to get  $p(y|\mathbf{x})$

$$p(y|\mathbf{x}) = \frac{p(y|\pi)p(\mathbf{x}|\theta_y)}{p(\mathbf{x})} = \frac{p(y|\pi)p(\mathbf{x}|\theta_y)}{\sum_y p(y|\pi)p(\mathbf{x}|\theta_y)}$$

- Parameters of distributions  $p(y|\pi)$  and  $p(x|\theta_y)$  are estimated from training data using point estimation methods (MLE or MAP) or using fully Bayesian inference (discussed today)
- Once these parameters  $\pi$  and  $\theta_y$  are estimated (point estimates, or full posterior if doing Bayesian inference), the above Bayes rule can be applied to a new input  $\hat{x}$  to compute  $p(\hat{y}|\hat{x})$

- Note: The generative model only defines  $p(y|\pi)$  and  $p(x|\theta_y)$ . Doesn't define p(y|x)
- We combine these using Bayes rule to get  $p(y|\mathbf{x})$

$$p(y|\mathbf{x}) = \frac{p(y|\pi)p(\mathbf{x}|\theta_y)}{p(\mathbf{x})} = \frac{p(y|\pi)p(\mathbf{x}|\theta_y)}{\sum_y p(y|\pi)p(\mathbf{x}|\theta_y)}$$

- Parameters of distributions  $p(y|\pi)$  and  $p(x|\theta_y)$  are estimated from training data using point estimation methods (MLE or MAP) or using fully Bayesian inference (discussed today)
- Once these parameters  $\pi$  and  $\theta_y$  are estimated (point estimates, or full posterior if doing Bayesian inference), the above Bayes rule can be applied to a new input  $\hat{x}$  to compute  $p(\hat{y}|\hat{x})$
- Let's now set up the parameter estimation for  $\pi$  and  $\theta_y$  as a Bayesian inference problem

- Note: The generative model only defines  $p(y|\pi)$  and  $p(x|\theta_y)$ . Doesn't define p(y|x)
- We combine these using Bayes rule to get  $p(y|\mathbf{x})$

$$p(y|\mathbf{x}) = \frac{p(y|\pi)p(\mathbf{x}|\theta_y)}{p(\mathbf{x})} = \frac{p(y|\pi)p(\mathbf{x}|\theta_y)}{\sum_y p(y|\pi)p(\mathbf{x}|\theta_y)}$$

- Parameters of distributions  $p(y|\pi)$  and  $p(x|\theta_y)$  are estimated from training data using point estimation methods (MLE or MAP) or using fully Bayesian inference (discussed today)
- Once these parameters  $\pi$  and  $\theta_y$  are estimated (point estimates, or full posterior if doing Bayesian inference), the above Bayes rule can be applied to a new input  $\hat{x}$  to compute  $p(\hat{y}|\hat{x})$
- Let's now set up the parameter estimation for  $\pi$  and  $\theta_y$  as a Bayesian inference problem
  - Note: As we will see in the end, in this approach, computing  $p(\hat{y}|\hat{x})$  for a new input  $\hat{x}$  will NOT use a point estimate of the parameters  $\pi, \theta_y$  but would use posterior averaging

イロト イロト イヨト イヨト (星) りんの

• Let us focus on the supervised, binary classification setting for now

道公

• Let us focus on the supervised, binary classification setting for now

 ${\, \bullet \, }$  In this case, we have three parameters to be learned:  $\pi, \, \theta_0,$  and  $\theta_1$ 

・ロト ・四ト ・モト ・モト

• Let us focus on the supervised, binary classification setting for now

- ${}_{\odot}$  In this case, we have three parameters to be learned:  $\pi,\,\theta_0,$  and  $\theta_1$ 
  - $\, \circ \,$  Probability  $\pi \in (0,1)$  of the Bernoulli. Can assume the following Beta prior

 $\pi \sim \mathsf{Beta}(a, b)$ 

イロト イロト イヨト イヨト

• Let us focus on the supervised, binary classification setting for now

- $\, \bullet \,$  In this case, we have three parameters to be learned:  $\pi, \, \theta_0,$  and  $\theta_1$ 
  - $\, \circ \,$  Probability  $\pi \in (0,1)$  of the Bernoulli. Can assume the following Beta prior

 $\pi \sim \mathsf{Beta}(a, b)$ 

 $\circ\,$  Parameters  $\theta_0,$  and  $\theta_1$  of the class-conditional distributions. Will assume the same prior on both

 $heta_0, heta_1\sim \pmb{p}( heta)$ 

「「「「」

• Let us focus on the supervised, binary classification setting for now

- $\, \bullet \,$  In this case, we have three parameters to be learned:  $\pi, \, \theta_0,$  and  $\theta_1$ 
  - $\, \circ \,$  Probability  $\pi \in (0,1)$  of the Bernoulli. Can assume the following Beta prior

 $\pi \sim \mathsf{Beta}(a, b)$ 

• Parameters  $\theta_0$ , and  $\theta_1$  of the class-conditional distributions. Will assume the same prior on both

 $heta_0, heta_1\sim p( heta)$ 

• Note: The actual form of  $p(\theta)$  will depend on what the class conditional distributions  $p(\mathbf{x}|\theta_0)$  and  $p(\mathbf{x}|\theta_1)$  are (e.g., if these are Gaussians and if we want to learn both mean and covariance matrix of these Gaussians, then  $p(\theta)$  will be some distribution over mean and covariance matrix, e.g., a Normal-inverse Wishart distribution)

• Let us focus on the supervised, binary classification setting for now

- $\, \bullet \,$  In this case, we have three parameters to be learned:  $\pi, \, \theta_0,$  and  $\theta_1$ 
  - $\, \circ \,$  Probability  $\pi \in (0,1)$  of the Bernoulli. Can assume the following Beta prior

 $\pi \sim \mathsf{Beta}(a, b)$ 

• Parameters  $\theta_0$ , and  $\theta_1$  of the class-conditional distributions. Will assume the same prior on both

 $heta_0, heta_1\sim p( heta)$ 

• Note: The actual form of  $p(\theta)$  will depend on what the class conditional distributions  $p(\mathbf{x}|\theta_0)$  and  $p(\mathbf{x}|\theta_1)$  are (e.g., if these are Gaussians and if we want to learn both mean and covariance matrix of these Gaussians, then  $p(\theta)$  will be some distribution over mean and covariance matrix, e.g., a Normal-inverse Wishart distribution)

• We will jointly denote the prior on  $\pi$ ,  $\theta_0$ , and  $\theta_1$  as  $p(\pi, \theta_0, \theta_1) = p(\pi)p(\theta_0)p(\theta_1)$ 

### The Likelihood

• Denote the  $N \times D$  feature matrix by X and the  $N \times 1$  label vector by  $\boldsymbol{y}$ 

重心

・ロト ・四ト ・モト ・モト

### The Likelihood

- Denote the  $N \times D$  feature matrix by X and the  $N \times 1$  label vector by  $\boldsymbol{y}$
- Since both X and y are being modeled here, the likelihood function will be

#### The Likelihood

- Denote the  $N \times D$  feature matrix by X and the  $N \times 1$  label vector by  $\boldsymbol{y}$
- Since both X and y are being modeled here, the likelihood function will be

$$p(X, \vec{y}|\pi, \theta_1, \theta_0) = \prod_{i=1}^N p(x_i, y_i|\pi, \theta_1, \theta_0)$$
  
= 
$$\prod_{i=1}^N p(x_i|y_i, \pi, \theta_1, \theta_0) p(y_i|\pi, \theta_1, \theta_0)$$
  
= 
$$\prod_{i=1}^N p(x_i|\theta_{y_i}) p(y_i|\pi)$$

イロト イロト イヨト イヨト

• We need to infer the following posterior distribution

$$p(\pi,\theta_1,\theta_0|\vec{y},X) = \frac{p(X,\vec{y}|\pi,\theta_1,\theta_0)p(\pi,\theta_1,\theta_0)}{\int_{\Omega_\theta}\int_0^1 p(X,\vec{y}|\pi,\theta_1,\theta_0)p(\pi,\theta_1,\theta_0)d\pi d\theta_1 d\theta_0}$$

• Note:  $\Omega_{\theta}$  denotes the domain of  $\theta$ 

• We need to infer the following posterior distribution

$$p(\pi,\theta_1,\theta_0|\vec{y},X) = \frac{p(X,\vec{y}|\pi,\theta_1,\theta_0)p(\pi,\theta_1,\theta_0)}{\int_{\Omega_\theta}\int_0^1 p(X,\vec{y}|\pi,\theta_1,\theta_0)p(\pi,\theta_1,\theta_0)d\pi d\theta_1 d\theta_0}$$

- Note:  $\Omega_{\theta}$  denotes the domain of  $\theta$
- Might look scary at first but it isn't actually

イロト イロト イヨト イヨト

• We need to infer the following posterior distribution

$$p(\pi,\theta_1,\theta_0|\vec{y},X) = \frac{p(X,\vec{y}|\pi,\theta_1,\theta_0)p(\pi,\theta_1,\theta_0)}{\int_{\Omega_\theta}\int_0^1 p(X,\vec{y}|\pi,\theta_1,\theta_0)p(\pi,\theta_1,\theta_0)d\pi d\theta_1 d\theta_0}$$

- Note:  $\Omega_{\theta}$  denotes the domain of  $\theta$
- Might look scary at first but it isn't actually
- Recall the prior  $p(\pi, \theta_0, \theta_1) = p(\pi)p(\theta_0)p(\theta_1)$ .

イロト イポト イヨト イヨト

• We need to infer the following posterior distribution

$$p(\pi,\theta_1,\theta_0|\vec{y},X) = \frac{p(X,\vec{y}|\pi,\theta_1,\theta_0)p(\pi,\theta_1,\theta_0)}{\int_{\Omega_\theta}\int_0^1 p(X,\vec{y}|\pi,\theta_1,\theta_0)p(\pi,\theta_1,\theta_0)d\pi d\theta_1 d\theta_0}$$

- Note:  $\Omega_{\theta}$  denotes the domain of  $\theta$
- Might look scary at first but it isn't actually

• Recall the prior  $p(\pi, \theta_0, \theta_1) = p(\pi)p(\theta_0)p(\theta_1)$ . The likelihood also factorized over data points, i.e.,

$$p(X, \mathbf{y}|\pi, \theta_1, \theta_0) = \prod_{i=1}^{N} p(x_i|\theta_{y_i}) p(y_i|\pi)$$

イロト イ伊ト イヨト イヨン

• We need to infer the following posterior distribution

$$p(\pi,\theta_1,\theta_0|\vec{y},X) = \frac{p(X,\vec{y}|\pi,\theta_1,\theta_0)p(\pi,\theta_1,\theta_0)}{\int_{\Omega_\theta}\int_0^1 p(X,\vec{y}|\pi,\theta_1,\theta_0)p(\pi,\theta_1,\theta_0)d\pi d\theta_1 d\theta_0}$$

- Note:  $\Omega_{ heta}$  denotes the domain of heta
- Might look scary at first but it isn't actually
- Recall the prior  $p(\pi, \theta_0, \theta_1) = p(\pi)p(\theta_0)p(\theta_1)$ . The likelihood also factorized over data points, i.e.,

$$p(X, \mathbf{y}|\pi, \theta_1, \theta_0) = \prod_{i=1}^{N} p(x_i|\theta_{y_i}) p(y_i|\pi)$$

• Thus, the posterior will be

$$p(\pi, \theta_1, \theta_0 | \vec{y}, X) \propto \left[ \prod_{i:y_i=1} p(x_i | \theta_1) p(\theta_1) \right] \left[ \prod_{i:y_i=0} p(x_i | \theta_0) p(\theta_0) \right] \left[ \prod_{i=1}^N p(y_i | \pi) p(\pi) \right]$$

• We need to infer the following posterior distribution

$$p(\pi,\theta_1,\theta_0|\vec{y},X) = \frac{p(X,\vec{y}|\pi,\theta_1,\theta_0)p(\pi,\theta_1,\theta_0)}{\int_{\Omega_\theta}\int_0^1 p(X,\vec{y}|\pi,\theta_1,\theta_0)p(\pi,\theta_1,\theta_0)d\pi d\theta_1 d\theta_0}$$

- Note:  $\Omega_{\theta}$  denotes the domain of  $\theta$
- Might look scary at first but it isn't actually
- Recall the prior  $p(\pi, \theta_0, \theta_1) = p(\pi)p(\theta_0)p(\theta_1)$ . The likelihood also factorized over data points, i.e.,

$$p(X, \mathbf{y}|\pi, \theta_1, \theta_0) = \prod_{i=1}^{N} p(x_i|\theta_{y_i}) p(y_i|\pi)$$

• Thus, the posterior will be

$$p(\pi, \theta_1, \theta_0 | \vec{y}, X) \propto \left[ \prod_{i:y_i=1} p(x_i | \theta_1) p(\theta_1) \right] \left[ \prod_{i:y_i=0} p(x_i | \theta_0) p(\theta_0) \right] \left[ \prod_{i=1}^N p(y_i | \pi) p(\pi) \right]$$

• But what about the normalization constant in the denominator?

Prob. Mod. & Inference - CS698X (Piyush Rai, IITK)

< d> < d> < d> < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d > < d

• Luckily, in this case, the same factorization structure simplies the denominator as well

$$p(\pi, \theta_1, \theta_0 | \vec{y}, X) = \frac{\prod_{i:y_i=1} p(x_i | \theta_1) p(\theta_1)}{\int \prod_{i:y_i=1} p(x_i | \theta_1) p(\theta_1) d\theta_1} \cdot \frac{\prod_{i:y_i=0} p(x_i | \theta_0) p(\theta_0)}{\int \prod_{i:y_i=0} p(x_i | \theta_0) p(\theta_0) d\theta_0} \cdot \frac{\prod_{i=1}^N p(y_i | \pi) p(\pi)}{\int \prod_{i=1}^N p(y_i | \pi) p(\pi) d\pi}$$

1

• Luckily, in this case, the same factorization structure simplies the denominator as well

$$p(\pi, \theta_1, \theta_0 | \vec{y}, X) = \frac{\prod_{i:y_i=1} p(x_i | \theta_1) p(\theta_1)}{\int \prod_{i:y_i=1} p(x_i | \theta_1) p(\theta_1) d\theta_1} \cdot \frac{\prod_{i:y_i=0} p(x_i | \theta_0) p(\theta_0)}{\int \prod_{i:y_i=0} p(x_i | \theta_0) p(\theta_0) d\theta_0} \cdot \frac{\prod_{i=1}^N p(y_i | \pi) p(\pi)}{\int \prod_{i=1}^N p(y_i | \pi) p(\pi) d\pi}$$

• The above is just a product of three posterior distributions !

$$p(\pi, \theta_1, \theta_0 | \vec{y}, X) = p(\theta_1 | \{ x_i : y_i = 1 \}) p(\theta_0 | \{ x_i : y_i = 0 \}) p(\pi | \vec{y})$$

(日) (명) (분) (분)

• Luckily, in this case, the same factorization structure simplies the denominator as well

$$p(\pi, \theta_1, \theta_0 | \vec{y}, X) = \frac{\prod_{i:y_i=1} p(x_i | \theta_1) p(\theta_1)}{\int \prod_{i:y_i=1} p(x_i | \theta_1) p(\theta_1) d\theta_1} \cdot \frac{\prod_{i:y_i=0} p(x_i | \theta_0) p(\theta_0)}{\int \prod_{i:y_i=0} p(x_i | \theta_0) p(\theta_0) d\theta_0} \cdot \frac{\prod_{i=1}^N p(y_i | \pi) p(\pi)}{\int \prod_{i=1}^N p(y_i | \pi) p(\pi) d\pi}$$

• The above is just a product of three posterior distributions !

$$p(\pi, \theta_1, \theta_0 | \vec{y}, X) = p(\theta_1 | \{ x_i : y_i = 1 \}) p(\theta_0 | \{ x_i : y_i = 0 \}) p(\pi | \vec{y})$$

• We also know what  $p(\pi|\mathbf{y})$  will be (recall the coin-toss example)

$$p(\pi|\vec{y}) \propto \prod_{i=1}^{N} p(y_i|\pi) p(\pi) \longrightarrow p(\pi|\vec{y}) = \text{Beta}(a + \sum_i y_i, b + N - \sum_i y_i)$$

Prob. Mod. & Inference - CS698X (Piyush Rai, IITK)

(日) (四) (三) (三) (三)

• Luckily, in this case, the same factorization structure simplies the denominator as well

$$p(\pi, \theta_1, \theta_0 | \vec{y}, X) = \frac{\prod_{i:y_i=1} p(x_i | \theta_1) p(\theta_1)}{\int \prod_{i:y_i=1} p(x_i | \theta_1) p(\theta_1) d\theta_1} \cdot \frac{\prod_{i:y_i=0} p(x_i | \theta_0) p(\theta_0)}{\int \prod_{i:y_i=0} p(x_i | \theta_0) p(\theta_0) d\theta_0} \cdot \frac{\prod_{i=1}^N p(y_i | \pi) p(\pi)}{\int \prod_{i=1}^N p(y_i | \pi) p(\pi) d\pi}$$

• The above is just a product of three posterior distributions !

$$p(\pi, \theta_1, \theta_0 | \vec{y}, X) = p(\theta_1 | \{ x_i : y_i = 1 \}) p(\theta_0 | \{ x_i : y_i = 0 \}) p(\pi | \vec{y})$$

• We also know what  $p(\pi|\mathbf{y})$  will be (recall the coin-toss example)

$$p(\pi|\vec{y}) \propto \prod_{i=1}^{N} p(y_i|\pi) p(\pi) \longrightarrow p(\pi|\vec{y}) = \text{Beta}(a + \sum_i y_i, b + N - \sum_i y_i)$$

• Form of posteriors on  $\theta_1$  and  $\theta_2$  will depend on  $p(\mathbf{x}|\theta_1)$  and  $p(\theta_1)$ , and  $p(\mathbf{x}|\theta_0)$  and  $p(\theta_0)$ , resp.

• We have already seen how to compute the parameter posterior  $p(\pi, \theta_1, \theta_0 | \boldsymbol{y}, X)$  for this model

• We have already seen how to compute the parameter posterior  $p(\pi, \theta_1, \theta_0 | \boldsymbol{y}, X)$  for this model

• Original goal is classification. We thus also want the predictive posterior for label of a new input, i.e.,  $p(\hat{y}|\hat{x})$ , for which the more "complete" notation in this Bayesian setting would be  $p(\hat{y}|\hat{x}, X, y)$ 

$$p(\hat{y}|\hat{x}, X, \vec{y}) = \int_{\Omega_{\theta}} \int_{\Omega_{\theta}} \int_{0}^{1} p(\hat{y}|\hat{x}, \theta_{1}, \theta_{0}, \pi) p(\theta_{1}, \theta_{0}, \pi | X, \vec{y}) d\pi d\theta_{1} d\theta_{0}$$

• We have already seen how to compute the parameter posterior  $p(\pi, \theta_1, \theta_0 | \boldsymbol{y}, X)$  for this model

• Original goal is classification. We thus also want the predictive posterior for label of a new input, i.e.,  $p(\hat{y}|\hat{x})$ , for which the more "complete" notation in this Bayesian setting would be  $p(\hat{y}|\hat{x}, X, y)$ 

$$p(\hat{y}|\hat{x}, X, \vec{y}) = \int_{\Omega_{\theta}} \int_{\Omega_{\theta}} \int_{0}^{1} p(\hat{y}|\hat{x}, \theta_{1}, \theta_{0}, \pi) p(\theta_{1}, \theta_{0}, \pi | X, \vec{y}) d\pi d\theta_{1} d\theta_{0}$$

• Luckily, in this case, this too has a rather simple form. Using Bayes rule, we have  $p(\hat{y}|\hat{x}, X, \vec{y}) = \frac{p(\hat{x}|\hat{y}, X, \vec{y})p(\hat{y}|X, \vec{y})}{p(\hat{x}|\hat{y} = 1, X, \vec{y})p(\hat{y} = 1|X, \vec{y}) + p(\hat{x}|\hat{y} = 0, X, \vec{y})p(\hat{y} = 0|X, \vec{y})}$ 

• We have already seen how to compute the parameter posterior  $p(\pi, \theta_1, \theta_0 | \boldsymbol{y}, X)$  for this model

• Original goal is classification. We thus also want the predictive posterior for label of a new input, i.e.,  $p(\hat{y}|\hat{x})$ , for which the more "complete" notation in this Bayesian setting would be  $p(\hat{y}|\hat{x}, X, y)$ 

$$p(\hat{y}|\hat{x}, X, \vec{y}) = \int_{\Omega_{\theta}} \int_{\Omega_{\theta}} \int_{0}^{1} p(\hat{y}|\hat{x}, \theta_{1}, \theta_{0}, \pi) p(\theta_{1}, \theta_{0}, \pi | X, \vec{y}) d\pi d\theta_{1} d\theta_{0}$$

• Luckily, in this case, this too has a rather simple form. Using Bayes rule, we have  $p(\hat{y}|\hat{x}, X, \vec{y}) = \frac{p(\hat{x}|\hat{y}, X, \vec{y})p(\hat{y}|X, \vec{y})}{p(\hat{x}|\hat{y} = 1, X, \vec{y})p(\hat{y} = 1|X, \vec{y}) + p(\hat{x}|\hat{y} = 0, X, \vec{y})p(\hat{y} = 0|X, \vec{y})}$   $= \frac{p(\hat{x}|\hat{y}, X, \vec{y})p(\hat{y}|\vec{y})}{p(\hat{x}|\hat{y} = 1, X, \vec{y})p(\hat{y} = 1|\vec{y}) + p(\hat{x}|\hat{y} = 0, X, \vec{y})p(\hat{y} = 0|\vec{y})}$ 

• We have already seen how to compute the parameter posterior  $p(\pi, \theta_1, \theta_0 | \boldsymbol{y}, X)$  for this model

• Original goal is classification. We thus also want the predictive posterior for label of a new input, i.e.,  $p(\hat{y}|\hat{x})$ , for which the more "complete" notation in this Bayesian setting would be  $p(\hat{y}|\hat{x}, X, y)$ 

$$p(\hat{y}|\hat{x}, X, \vec{y}) = \int_{\Omega_{\theta}} \int_{\Omega_{\theta}} \int_{0}^{1} p(\hat{y}|\hat{x}, \theta_{1}, \theta_{0}, \pi) p(\theta_{1}, \theta_{0}, \pi | X, \vec{y}) d\pi d\theta_{1} d\theta_{0}$$

• Luckily, in this case, this too has a rather simple form. Using Bayes rule, we have  $\begin{array}{rcl} p(\hat{y}|\hat{x}, X, \vec{y}) &=& \frac{p(\hat{x}|\hat{y}, X, \vec{y})p(\hat{y}|X, \vec{y})}{p(\hat{x}|\hat{y}=1, X, \vec{y})p(\hat{y}=1|X, \vec{y}) + p(\hat{x}|\hat{y}=0, X, \vec{y})p(\hat{y}=0|X, \vec{y})} \\ &=& \frac{p(\hat{x}|\hat{y}, X, \vec{y})p(\hat{y}|\vec{y})}{p(\hat{x}|\hat{y}=1, X, \vec{y})p(\hat{y}=1|\vec{y}) + p(\hat{x}|\hat{y}=0, X, \vec{y})p(\hat{y}=0|\vec{y})} \end{array}$ 

• In order to compute this, we need  $p(\hat{x}|\hat{y}, X, \boldsymbol{y})$  and  $p(\hat{y}|\boldsymbol{y})$ 

• We have already seen how to compute the parameter posterior  $p(\pi, \theta_1, \theta_0 | \boldsymbol{y}, X)$  for this model

• Original goal is classification. We thus also want the predictive posterior for label of a new input, i.e.,  $p(\hat{y}|\hat{x})$ , for which the more "complete" notation in this Bayesian setting would be  $p(\hat{y}|\hat{x}, X, y)$ 

$$p(\hat{y}|\hat{x}, X, \vec{y}) = \int_{\Omega_{\theta}} \int_{\Omega_{\theta}} \int_{0}^{1} p(\hat{y}|\hat{x}, \theta_{1}, \theta_{0}, \pi) p(\theta_{1}, \theta_{0}, \pi | X, \vec{y}) d\pi d\theta_{1} d\theta_{0}$$

• Luckily, in this case, this too has a rather simple form. Using Bayes rule, we have  $\begin{array}{rcl} p(\hat{y}|\hat{x}, X, \vec{y}) &=& \frac{p(\hat{x}|\hat{y}, X, \vec{y})p(\hat{y}|X, \vec{y})}{p(\hat{x}|\hat{y}=1, X, \vec{y})p(\hat{y}=1|X, \vec{y}) + p(\hat{x}|\hat{y}=0, X, \vec{y})p(\hat{y}=0|X, \vec{y})} \\ &=& \frac{p(\hat{x}|\hat{y}, X, \vec{y})p(\hat{y}|\vec{y})}{p(\hat{x}|\hat{y}=1, X, \vec{y})p(\hat{y}=1|\vec{y}) + p(\hat{x}|\hat{y}=0, X, \vec{y})p(\hat{y}=0|\vec{y})} \end{array}$ 

• In order to compute this, we need  $p(\hat{x}|\hat{y},X,\textbf{\textit{y}})$  and  $p(\hat{y}|\textbf{\textit{y}})$ 

•  $p(\hat{x}|\hat{y}, X, y)$ : Marginal class-conditional distribution of the new input vector  $\hat{x}$ 

• We have already seen how to compute the parameter posterior  $p(\pi, \theta_1, \theta_0 | \boldsymbol{y}, X)$  for this model

• Original goal is classification. We thus also want the predictive posterior for label of a new input, i.e.,  $p(\hat{y}|\hat{x})$ , for which the more "complete" notation in this Bayesian setting would be  $p(\hat{y}|\hat{x}, X, y)$ 

$$p(\hat{y}|\hat{x}, X, \vec{y}) = \int_{\Omega_{\theta}} \int_{\Omega_{\theta}} \int_{0}^{1} p(\hat{y}|\hat{x}, \theta_{1}, \theta_{0}, \pi) p(\theta_{1}, \theta_{0}, \pi | X, \vec{y}) d\pi d\theta_{1} d\theta_{0}$$

• Luckily, in this case, this too has a rather simple form. Using Bayes rule, we have  $\begin{array}{rcl} p(\hat{y}|\hat{x}, X, \vec{y}) &=& \frac{p(\hat{x}|\hat{y}, X, \vec{y})p(\hat{y}|X, \vec{y})}{p(\hat{x}|\hat{y}=1, X, \vec{y})p(\hat{y}=1|X, \vec{y}) + p(\hat{x}|\hat{y}=0, X, \vec{y})p(\hat{y}=0|X, \vec{y})} \\ &=& \frac{p(\hat{x}|\hat{y}, X, \vec{y})p(\hat{y}|\vec{y})}{p(\hat{x}|\hat{y}=1, X, \vec{y})p(\hat{y}=1|\vec{y}) + p(\hat{x}|\hat{y}=0, X, \vec{y})p(\hat{y}=0|\vec{y})} \end{array}$ 

• In order to compute this, we need  $p(\hat{x}|\hat{y},X,\textbf{y})$  and  $p(\hat{y}|\textbf{y})$ 

- $p(\hat{x}|\hat{y}, X, y)$ : Marginal class-conditional distribution of the new input vector  $\hat{x}$
- $p(\hat{y}|y)$ : Marginal probability of its label  $\hat{y}$  given the labels of training data

# The Predictive Posterior Distribution (Contd.)

• Predictive posterior requires computing  $p(\hat{x}|\hat{y}, X, y)$  and  $p(\hat{y}|y)$ 

# The Predictive Posterior Distribution (Contd.)

- Predictive posterior requires computing  $p(\hat{x}|\hat{y}, X, y)$  and  $p(\hat{y}|y)$
- The marginal likelihood  $p(\hat{x}|\hat{y}, X, y)$  of  $\hat{x}$  can be computed as

$$p(\hat{x}|\hat{y}, X, \vec{y}) = \int_{\Omega_{\theta}} \int_{\Omega_{\theta}} p(\hat{x}|\hat{y}, \theta_1, \theta_0) p(\theta_1, \theta_0 | X, \vec{y}) d\theta_1 d\theta_0$$
$$= \int_{\Omega_{\theta}} p(\hat{x}|\theta_{\hat{y}}) p(\theta_{\hat{y}}|\{x_i : y_i = \hat{y}\}) d\theta_{\hat{y}}$$

イロト イロト イヨト イヨト
# The Predictive Posterior Distribution (Contd.)

- Predictive posterior requires computing  $p(\hat{x}|\hat{y}, X, y)$  and  $p(\hat{y}|y)$
- The marginal likelihood  $p(\hat{x}|\hat{y}, X, y)$  of  $\hat{x}$  can be computed as

$$p(\hat{x}|\hat{y}, X, \vec{y}) = \int_{\Omega_{\theta}} \int_{\Omega_{\theta}} p(\hat{x}|\hat{y}, \theta_1, \theta_0) p(\theta_1, \theta_0 | X, \vec{y}) d\theta_1 d\theta_0$$
$$= \int_{\Omega_{\theta}} p(\hat{x}|\theta_{\hat{y}}) p(\theta_{\hat{y}}|\{x_i : y_i = \hat{y}\}) d\theta_{\hat{y}}$$

 The above is simply the posterior predictive distribution of class ŷ. The final expression will depend on the forms of p(x̂|θ<sub>ŷ</sub>) and p(θ<sub>ŷ</sub>|.). If exp-family, we will have closed form expression!

## The Predictive Posterior Distribution (Contd.)

- Predictive posterior requires computing  $p(\hat{x}|\hat{y}, X, y)$  and  $p(\hat{y}|y)$
- The marginal likelihood  $p(\hat{x}|\hat{y}, X, y)$  of  $\hat{x}$  can be computed as

$$p(\hat{x}|\hat{y}, X, \vec{y}) = \int_{\Omega_{\theta}} \int_{\Omega_{\theta}} p(\hat{x}|\hat{y}, \theta_1, \theta_0) p(\theta_1, \theta_0 | X, \vec{y}) d\theta_1 d\theta_0$$
$$= \int_{\Omega_{\theta}} p(\hat{x}|\theta_{\hat{y}}) p(\theta_{\hat{y}}|\{x_i : y_i = \hat{y}\}) d\theta_{\hat{y}}$$

- The above is simply the posterior predictive distribution of class ŷ. The final expression will depend on the forms of p(x̂|θ<sub>ŷ</sub>) and p(θ<sub>ŷ</sub>|.). If exp-family, we will have closed form expression!
- The marginal likelihood  $p(\hat{y}|y)$  is something we have already seen (recall Bernoulli coin-toss)

$$p(\hat{y} = 1|m{y}) = \int p(\hat{y} = 1|\pi) p(\pi|m{y}) d\pi = \int \pi p(\pi|m{y}) d\pi = rac{a + \sum_{i=1}^{N} y_i}{a + b + N}$$

## The Predictive Posterior Distribution (Contd.)

- Predictive posterior requires computing  $p(\hat{x}|\hat{y}, X, y)$  and  $p(\hat{y}|y)$
- The marginal likelihood  $p(\hat{x}|\hat{y}, X, y)$  of  $\hat{x}$  can be computed as

$$p(\hat{x}|\hat{y}, X, \vec{y}) = \int_{\Omega_{\theta}} \int_{\Omega_{\theta}} p(\hat{x}|\hat{y}, \theta_1, \theta_0) p(\theta_1, \theta_0 | X, \vec{y}) d\theta_1 d\theta_0$$
$$= \int_{\Omega_{\theta}} p(\hat{x}|\theta_{\hat{y}}) p(\theta_{\hat{y}}|\{x_i : y_i = \hat{y}\}) d\theta_{\hat{y}}$$

- The above is simply the posterior predictive distribution of class ŷ. The final expression will depend on the forms of p(x̂|θ<sub>ŷ</sub>) and p(θ<sub>ŷ</sub>|.). If exp-family, we will have closed form expression!
- The marginal likelihood  $p(\hat{y}|y)$  is something we have already seen (recall Bernoulli coin-toss)

$$p(\hat{y}=1|oldsymbol{y})=\int p(\hat{y}=1|\pi)p(\pi|oldsymbol{y})d\pi=\int \pi p(\pi|oldsymbol{y})d\pi=rac{oldsymbol{a}+\sum_{i=1}^N y_i}{oldsymbol{a}+b+N}$$

• .. and  $p(\hat{y}=0|\boldsymbol{y})=1-p(\hat{y}=1|\boldsymbol{y})=rac{b+N-\sum_{i=1}^{N}y_i}{a+b+N}$ 

• Usually the most critical choice in generative classification is that of class conditional  $p(\mathbf{x}|\theta_y)$ 

- Usually the most critical choice in generative classification is that of class conditional  $p(\mathbf{x}|\theta_y)$
- Very complex  $p(\mathbf{x}|\theta_y)$  with lots of parameters may make estimation difficult

- Usually the most critical choice in generative classification is that of class conditional  $p(\mathbf{x}|\theta_y)$
- Very complex  $p(\mathbf{x}|\theta_y)$  with lots of parameters may make estimation difficult
- Often however we can choose simple forms of  $p(\mathbf{x}|\theta_y)$  to make estimation easier

- Usually the most critical choice in generative classification is that of class conditional  $p(\mathbf{x}|\theta_y)$
- Very complex  $p(\mathbf{x}|\theta_y)$  with lots of parameters may make estimation difficult
- Often however we can choose simple forms of  $p(\mathbf{x}|\theta_y)$  to make estimation easier
- The naïve Bayes assumption: The conditional distribution  $p(\mathbf{x}|\theta_y)$  factorizes over individual features (or over groups of features)

< ロ > < 回 > < 三 > < 三 >

- Usually the most critical choice in generative classification is that of class conditional  $p(\mathbf{x}|\theta_y)$
- Very complex  $p(\mathbf{x}|\theta_y)$  with lots of parameters may make estimation difficult
- Often however we can choose simple forms of  $p(\mathbf{x}|\theta_y)$  to make estimation easier
- The naïve Bayes assumption: The conditional distribution  $p(\mathbf{x}|\theta_y)$  factorizes over individual features (or over groups of features)
  - Suppose the features of  $\hat{x}$  can be partitioned into v groups  $\hat{x} = \{\hat{x}(j)\}_{j=1}^{v}$

- Usually the most critical choice in generative classification is that of class conditional  $p(\mathbf{x}|\theta_y)$
- Very complex  $p(\mathbf{x}|\theta_y)$  with lots of parameters may make estimation difficult
- Often however we can choose simple forms of  $p(\mathbf{x}|\theta_y)$  to make estimation easier
- The naïve Bayes assumption: The conditional distribution  $p(\mathbf{x}|\theta_y)$  factorizes over individual features (or over groups of features)
  - Suppose the features of  $\hat{x}$  can be partitioned into v groups  $\hat{x} = \{\hat{x}(j)\}_{j=1}^{v}$
  - $\,\circ\,$  Can also assume a similar partitioning for the parameters  $\theta_{\hat{y}}$

- Usually the most critical choice in generative classification is that of class conditional  $p(\mathbf{x}|\theta_y)$
- Very complex  $p(\mathbf{x}|\theta_y)$  with lots of parameters may make estimation difficult
- Often however we can choose simple forms of  $p(\mathbf{x}|\theta_y)$  to make estimation easier
- The naïve Bayes assumption: The conditional distribution  $p(\mathbf{x}|\theta_y)$  factorizes over individual features (or over groups of features)
  - Suppose the features of  $\hat{x}$  can be partitioned into v groups  $\hat{x} = \{\hat{x}(j)\}_{j=1}^{v}$
  - $\circ~$  Can also assume a similar partitioning for the parameters  $\theta_{\hat{y}}$
  - This further simplifies calculation of marginal likelihood  $p(\hat{x}|\hat{y}, X, y)$

$$\begin{split} p(\hat{x}|\hat{y}, X, \vec{y}) &= \int_{\Omega_{\theta}} \prod_{j=1}^{v} p(\hat{x}(j)|\theta_{\hat{y}}(j)) p(\theta_{\hat{y}}(j)|\{x_{i}(j) : y_{i} = \hat{y}\}) d\theta_{\hat{y}} \\ &= \prod_{j=1}^{v} \int p(\hat{x}(j)|\theta_{\hat{y}}(j)) p(\theta_{\hat{y}}(j)|\{x_{i}(j) : y_{i} = \hat{y}\}) d\theta_{\hat{y}}(j) \end{split}$$

- Usually the most critical choice in generative classification is that of class conditional  $p(\mathbf{x}|\theta_y)$
- Very complex  $p(\mathbf{x}|\theta_y)$  with lots of parameters may make estimation difficult
- Often however we can choose simple forms of  $p(\mathbf{x}|\theta_y)$  to make estimation easier
- The naïve Bayes assumption: The conditional distribution  $p(\mathbf{x}|\theta_y)$  factorizes over individual features (or over groups of features)
  - Suppose the features of  $\hat{x}$  can be partitioned into v groups  $\hat{x} = \{\hat{x}(j)\}_{j=1}^{v}$
  - $\circ~$  Can also assume a similar partitioning for the parameters  $\theta_{\hat{y}}$
  - This further simplifies calculation of marginal likelihood  $p(\hat{x}|\hat{y}, X, y)$

$$\begin{aligned} p(\hat{x}|\hat{y}, X, \vec{y}) &= \int_{\Omega_{\theta}} \prod_{j=1}^{v} p(\hat{x}(j)|\theta_{\hat{y}}(j)) p(\theta_{\hat{y}}(j)|\{x_{i}(j) : y_{i} = \hat{y}\}) d\theta_{\hat{y}} \\ &= \prod_{j=1}^{v} \int p(\hat{x}(j)|\theta_{\hat{y}}(j)) p(\theta_{\hat{y}}(j)|\{x_{i}(j) : y_{i} = \hat{y}\}) d\theta_{\hat{y}}(j) \end{aligned}$$

• This modeling choice in a Bayesian setting gives rise to a "Bayesian naïve Bayes" model

• In the Bayesian naïve Bayes model, we can still choose different types of class conditional  $p(\mathbf{x}|\theta_y)$ 

- In the Bayesian naïve Bayes model, we can still choose different types of class conditional  $p(\mathbf{x}|\theta_{y})$ 
  - Gaussian naïve Bayes: if x is modeled using a multivariate Gaussian (assumed factorized as per the naïve Bayes assumption)

イロト イロト イヨト イヨト

- In the Bayesian naïve Bayes model, we can still choose different types of class conditional  $p(\mathbf{x}|\theta_y)$ 
  - Gaussian naïve Bayes: if x is modeled using a multivariate Gaussian (assumed factorized as per the naïve Bayes assumption)
  - Multivariate Bernoulli naïve Bayes: if x is modeled using a multivariate Bernoulli (assumed factorized as per the naïve Bayes assumption)

イロト イロト イヨト イヨト

- In the Bayesian naïve Bayes model, we can still choose different types of class conditional  $p(\mathbf{x}|\theta_y)$ 
  - Gaussian naïve Bayes: if x is modeled using a multivariate Gaussian (assumed factorized as per the naïve Bayes assumption)
  - Multivariate Bernoulli naïve Bayes: if x is modeled using a multivariate Bernoulli (assumed factorized as per the naïve Bayes assumption)
- MLAPP (Murphy) Section 3.5.1.2 and 3.5.5 contains an example of Multivariate Bernoulli case

イロト イロト イヨト イヨト