

# **Classification: (Bayesian) Logistic Regression (and our first tryst with non-conjugacy!)**

Piyush Rai

Topics in Probabilistic Modeling and Inference (CS698X)

Jan 25, 2018

# Probabilistic Models for Classification

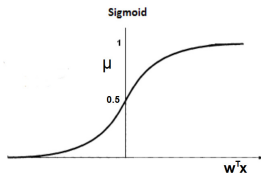
- The goal is to learn  $p(y|\mathbf{x})$ . Here  $p(y|\mathbf{x})$  will be a discrete distribution (e.g., Bernoulli, multinoulli)
- Usually two approaches to learn  $p(y|\mathbf{x})$ : Discriminative Classification and Generative Classification
- **Discriminative Classification:** Model and learn  $p(y|\mathbf{x})$  directly
  - This approach does not model the distribution of the inputs  $\mathbf{x}$
- **Generative Classification:** Model and learn  $p(y|\mathbf{x})$  “indirectly” as  $p(y|\mathbf{x}) = \frac{p(y)p(\mathbf{x}|y)}{p(\mathbf{x})}$ 
  - Called generative because, via  $p(\mathbf{x}|y)$ , we model how the inputs  $\mathbf{x}$  of each class are generated
  - The approach requires first learning **class-marginal**  $p(y)$  and **class-conditional** distributions  $p(\mathbf{x}|y)$
  - Usually harder to learn than discriminative but also has some advantages (more on this later)
- Both approaches can be given a non-Bayesian or Bayesian treatment
  - The Bayesian treatment won't rely on point estimates but infer the posterior over unknowns

# Discriminative Classification via Logistic Regression

- **Logistic Regression** (LR) is an example of discriminative **binary** classification, i.e.,  $y \in \{0, 1\}$
- Logistic Regression models  $\mathbf{x}$  to  $y$  relationship using the **sigmoid function**

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

where  $\mathbf{w} \in \mathbb{R}^D$  is the weight vector. Also note that  $p(y = 0|\mathbf{x}, \mathbf{w}) = 1 - \mu$



- A large positive (negative) “score”  $\mathbf{w}^\top \mathbf{x}$  means large probability of label being 1 (0)
- Is sigmoid the only way to convert the score into a probability?
  - No, while LR does that, there exist models that define  $\mu$  in other ways. E.g. **Probit Regression**

$$\mu = p(y = 1|\mathbf{x}, \mathbf{w}) = \Phi(\mathbf{w}^\top \mathbf{x}) \quad (\text{where } \Phi \text{ denotes the CDF of } \mathcal{N}(0, 1))$$

# Logistic Regression

- The LR classification rule is

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

$$p(y = 0|x, \mathbf{w}) = 1 - \mu = 1 - \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

- This implies a **Bernoulli likelihood** model for the labels

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

- Given  $N$  observations  $(\mathbf{X}, \mathbf{y}) = \{\mathbf{x}_n, y_n\}_{n=1}^N$ , we can do point estimation for  $\mathbf{w}$  by maximizing the log-likelihood (or minimizing the **negative log-likelihood**). This is basically MLE.

$$\mathbf{w}_{MLE} = \arg \max_{\mathbf{w}} \sum_{n=1}^N \log p(y_n|x_n, \mathbf{w}) = \arg \min_{\mathbf{w}} - \sum_{n=1}^N \log p(y_n|x_n, \mathbf{w}) = \arg \min_{\mathbf{w}} \text{NLL}(\mathbf{w})$$

- Convex loss function. Global minima. Both first order and second order methods widely used.

- Can also add a regularizer on  $\mathbf{w}$  to prevent overfitting. This corresponds to doing MAP estimation with a prior on  $\mathbf{w}$ , i.e.,  $\mathbf{w}_{MAP} = \arg \max_{\mathbf{w}} [\sum_{n=1}^N \log p(y_n|x_n, \mathbf{w}) + \log p(\mathbf{w})]$

# Bayesian Logistic Regression

- MLE/MAP only gives a point estimate. We would like to infer the full posterior over  $\mathbf{w}$
- Recall that the **likelihood model** is Bernoulli

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

- Just like the Bayesian linear regression case, let's use a Gaussian **prior** on  $\mathbf{w}$

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

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

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

- The denominator is intractable in general (logistic-Bernoulli and Gaussian are not conjugate)
  - Can't get a closed form expression for  $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$ . Must approximate it!
  - Several ways to do it, e.g., MCMC, variational inference, **Laplace approximation** (today)

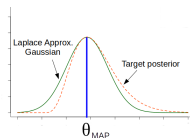
# Laplace Approximation

(Our first posterior approximation method)

# Laplace Approximation of Posterior Distribution

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

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

# Derivation of the Laplace Approximation

- Let's write the Bayes rule as

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

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

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

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

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



# Derivation of the Laplace Approximation

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

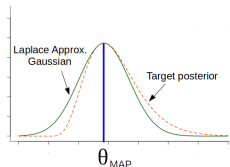
$$p(\theta|\mathcal{D}) = \frac{e^{\log p(\mathcal{D}, \theta)}}{\int e^{\log p(\mathcal{D}, \theta)} d\theta} \approx \frac{e^{\log p(\mathcal{D}, \theta_{MAP}) + \frac{1}{2}(\theta - \theta_{MAP})^\top \nabla^2 \log p(\mathcal{D}, \theta_{MAP})(\theta - \theta_{MAP})}}{\int e^{\log p(\mathcal{D}, \theta_{MAP})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

- Further simplifying, we have

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

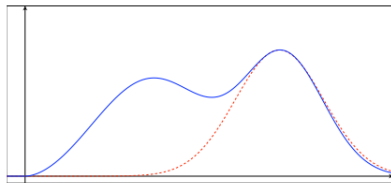
- Therefore the Laplace approximation of the posterior  $p(\theta|\mathcal{D})$  is a Gaussian and is given by

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



# Properties of Laplace Approximation

- 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 something like a posterior distribution over the parameters
  - negative log-likelihood (NLL) = loss function, negative log-prior = regularizer

# 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

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

- The required quantities are defined as

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

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

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

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

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

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

- Note: When using second order methods for estimating  $\mathbf{w}_{MAP}$ , we anyway get the Hessian needed for the Laplace approximation of the posterior

# An Aside: Gradient and Hessian for Logistic Regression

- 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

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

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

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

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

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

# Logistic Regression: Predictive Distributions

- When using MLE, the predictive distribution will be

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

- When using MAP, the predictive distribution will be

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

- 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.  $\therefore$  ( If using the Laplace approximation for  $p(\mathbf{w} | \mathbf{X}, \mathbf{y})$ , it will be

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

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

# Posterior Predictive via Monte-Carlo Sampling

- The posterior predictive is given by the following integral

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

- **Monte-Carlo approximation:** Draw several samples of  $\mathbf{w}$  from  $\mathcal{N}(\mathbf{w} | \mathbf{w}_{MAP}, \mathbf{H}^{-1})$  and replace the above integral by an empirical average of  $\sigma(\mathbf{w}^\top \mathbf{x}_*)$  computed using each of those samples

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

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

- More on Monte-Carlo methods when we discuss MCMC sampling

# Predictive Posterior via Probit Approximation

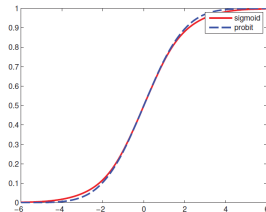
- The posterior predictive we wanted to compute was

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

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

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

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



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

# Predictive Posterior via Probit Approximation

- With this approximation, the predictive posterior will be

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

- Since  $a = \mathbf{w}^\top \mathbf{x}_* = \mathbf{x}_*^\top \mathbf{w}$ , 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

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

- 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.  $t^2 = \pi/8$

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



# Bayesian Logistic Regression: Posterior over Linear Classifiers!

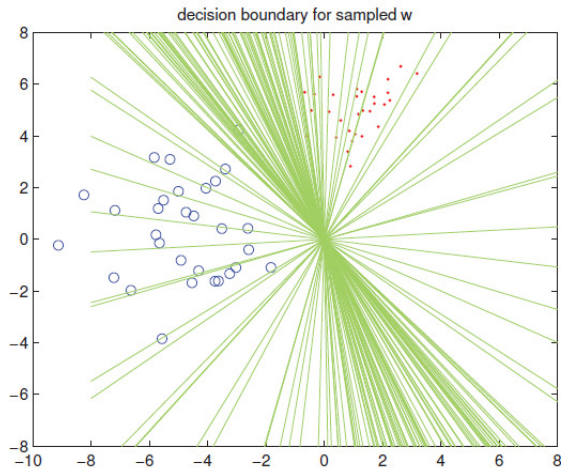
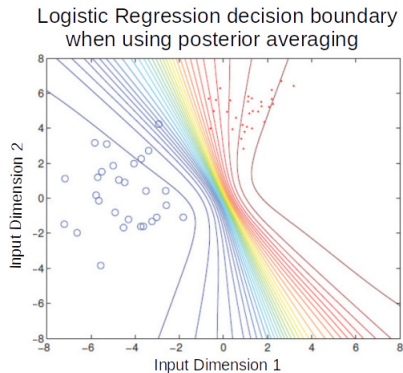
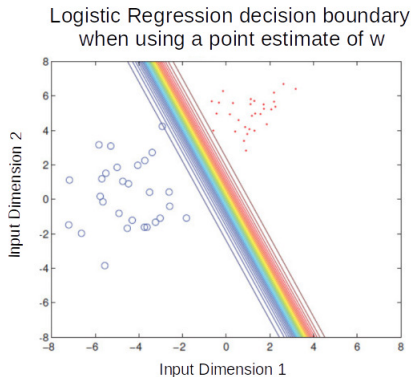


Figure courtesy: MLAPP (Murphy)

# Logistic Regression: Plug-in Prediction vs Bayesian Averaging

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



# Some Comments

- 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  $p(\mathbf{w}) = \mathcal{N}(0, \lambda^{-1}\mathbf{I})$ ) 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