

Foundations of Machine Learning
African Masters in Machine Intelligence



AIMS | African Institute for
Mathematical Sciences
RWANDA

**Imperial College
London**

Logistic Regression

Marc Deisenroth

Quantum Leap Africa
African Institute for Mathematical
Sciences, Rwanda

Department of Computing
Imperial College London



@mpd37

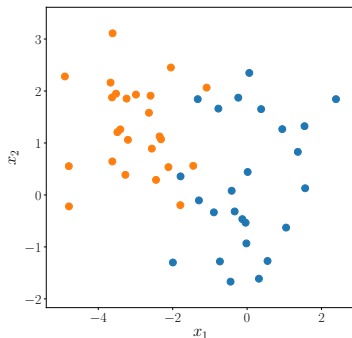
mdeisenroth@aimsammi.org

November 5, 2018

Learning Material

- ▶ Pattern Recognition and Machine Learning, Chapter 4 (Bishop, 2006)
- ▶ Machine Learning: A Probabilistic Perspective, Chapter 8 (Murphy, 2012)

Binary Classification



- ▶ Supervised learning setting with inputs $x_n \in \mathbb{R}^D$ and **binary** targets $y_n \in \{0, 1\}$ belonging to **classes** $\mathcal{C}_1, \mathcal{C}_2$.
- ▶ Objective: Find a decision boundary/surface that separates the two classes as well as possible

Class Posteriors

- ▶ Binary classification problem with two classes $\mathcal{C}_1, \mathcal{C}_2$.

Class Posteriors

- ▶ Binary classification problem with two classes $\mathcal{C}_1, \mathcal{C}_2$.
- ▶ Posterior class probability $p(y = 1|\mathbf{x}) = p(\mathcal{C}_1|\mathbf{x})$:

$$p(\mathcal{C}_1|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x})},$$

$$p(\mathbf{x}) = p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)$$

Class Posteriors

- ▶ Binary classification problem with two classes $\mathcal{C}_1, \mathcal{C}_2$.
- ▶ Posterior class probability $p(y = 1|\mathbf{x}) = p(\mathcal{C}_1|\mathbf{x})$:

$$p(\mathcal{C}_1|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x})},$$

$$p(\mathbf{x}) = p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)$$

- ▶ Define the **log-ratio of the posteriors (log-odds)**

$$a := \log \frac{p(\mathcal{C}_1|\mathbf{x})}{p(\mathcal{C}_2|\mathbf{x})} = \log \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}$$

Class Posteriors

- ▶ Binary classification problem with two classes $\mathcal{C}_1, \mathcal{C}_2$.
- ▶ Posterior class probability $p(y = 1|\mathbf{x}) = p(\mathcal{C}_1|\mathbf{x})$:

$$p(\mathcal{C}_1|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x})},$$

$$p(\mathbf{x}) = p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)$$

- ▶ Define the **log-ratio of the posteriors (log-odds)**

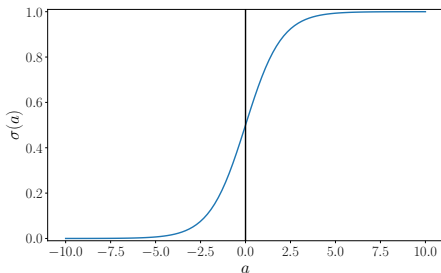
$$a := \log \frac{p(\mathcal{C}_1|\mathbf{x})}{p(\mathcal{C}_2|\mathbf{x})} = \log \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}$$

- ▶ Then

$$\underbrace{\sigma(a) := \frac{1}{1 + \exp(-a)}}_{\text{logistic sigmoid}} = ?$$

▶▶ **Discuss with your neighbors**

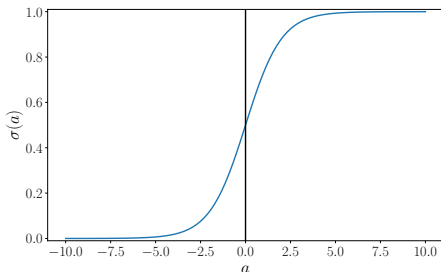
Logistic Sigmoid



$$a := \log \frac{p(\mathcal{C}_1|\mathbf{x})}{p(\mathcal{C}_2|\mathbf{x})} = \log \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}$$

$$\sigma(a) := \frac{1}{1 + \exp(-a)} = p(\mathcal{C}_1|\mathbf{x}) \quad \text{Logistic sigmoid}$$

Logistic Sigmoid



$$a := \log \frac{p(\mathcal{C}_1|\mathbf{x})}{p(\mathcal{C}_2|\mathbf{x})} = \log \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}$$

$$\sigma(a) := \frac{1}{1 + \exp(-a)} = p(\mathcal{C}_1|\mathbf{x}) \quad \text{Logistic sigmoid}$$

- Assign the label for \mathcal{C}_1 to \mathbf{x} if $\sigma(a) = p(\mathcal{C}_1|\mathbf{x}) = p(y = 1|\mathbf{x}) \geq 0.5$

Generalization to the Multiclass Setting

- ▶ Assume we are given K classes. Then

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{\sum_{j=1}^K p(\mathbf{x}|\mathcal{C}_j)p(\mathcal{C}_j)}$$

is the generalization of the logistic sigmoid to K classes.

- ▶ **Softmax function, Boltzmann distribution, normalized exponential**

Implicit Modeling Assumptions

- ▶ Assume Gaussian class conditionals

$$p(\mathbf{x}|\mathcal{C}_k) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma})$$

where the covariance matrix $\boldsymbol{\Sigma}$ is shared across all K classes.

Implicit Modeling Assumptions

- ▶ Assume Gaussian class conditionals

$$p(\mathbf{x}|\mathcal{C}_k) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma})$$

where the covariance matrix $\boldsymbol{\Sigma}$ is shared across all K classes.

- ▶ For $K = 2$ we get (Bishop, 2006)

$$p(\mathcal{C}_1|\mathbf{x}) = \sigma(\boldsymbol{\theta}^\top \mathbf{x} + \theta_0),$$

$$\boldsymbol{\theta} := \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2), \quad \theta_0 := \frac{1}{2} \left(\boldsymbol{\mu}_2^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 - \boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 \right) + \log \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)}$$

Implicit Modeling Assumptions

- ▶ Assume Gaussian class conditionals

$$p(\mathbf{x}|\mathcal{C}_k) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma})$$

where the covariance matrix $\boldsymbol{\Sigma}$ is shared across all K classes.

- ▶ For $K = 2$ we get (Bishop, 2006)

$$p(\mathcal{C}_1|\mathbf{x}) = \sigma(\boldsymbol{\theta}^\top \mathbf{x} + \theta_0),$$

$$\boldsymbol{\theta} := \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2), \quad \theta_0 := \frac{1}{2} \left(\boldsymbol{\mu}_2^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 - \boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 \right) + \log \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)}$$

- ▶▶ Argument of the sigmoid is linear in \mathbf{x}

Implicit Modeling Assumptions

- ▶ Assume **Gaussian class conditionals**

$$p(\mathbf{x}|\mathcal{C}_k) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma})$$

where the **covariance matrix $\boldsymbol{\Sigma}$** is shared across all K classes.

- ▶ For $K = 2$ we get (Bishop, 2006)

$$p(\mathcal{C}_1|\mathbf{x}) = \sigma(\boldsymbol{\theta}^\top \mathbf{x} + \theta_0),$$

$$\boldsymbol{\theta} := \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2), \quad \theta_0 := \frac{1}{2} \left(\boldsymbol{\mu}_2^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 - \boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 \right) + \log \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)}$$

- ▶▶ Argument of the sigmoid is linear in \mathbf{x}
- ▶▶ Decision boundary is a surface along which the posterior class probabilities $p(\mathcal{C}_k|\mathbf{x})$ are constant
- ▶▶ **Decision boundary is a linear function of \mathbf{x}**

Implicit Modeling Assumptions

- ▶ Assume **Gaussian class conditionals**

$$p(\mathbf{x}|\mathcal{C}_k) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma})$$

where the **covariance matrix $\boldsymbol{\Sigma}$** is shared across all K classes.

- ▶ For $K = 2$ we get (Bishop, 2006)

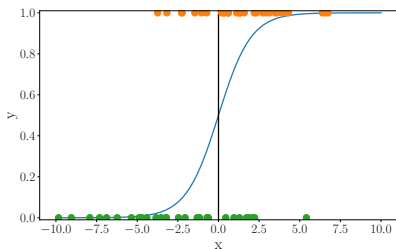
$$p(\mathcal{C}_1|\mathbf{x}) = \sigma(\boldsymbol{\theta}^\top \mathbf{x} + \theta_0),$$

$$\boldsymbol{\theta} := \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2), \quad \theta_0 := \frac{1}{2} \left(\boldsymbol{\mu}_2^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 - \boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 \right) + \log \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)}$$

- ▶▶ Argument of the sigmoid is linear in \mathbf{x}
 - ▶▶ Decision boundary is a surface along which the posterior class probabilities $p(\mathcal{C}_k|\mathbf{x})$ are constant
 - ▶▶ **Decision boundary is a linear function of \mathbf{x}**
- ▶ If covariances are not shared: Quadratic decision boundaries

Model Specification (Logistic Regression)

likelihood



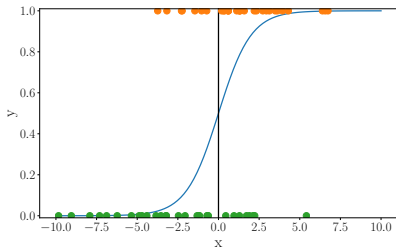
Model Specification (Logistic Regression)

- Bernoulli likelihood

$$y \in \{0, 1\}$$

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \text{Ber}(y|\mu(\mathbf{x})),$$

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



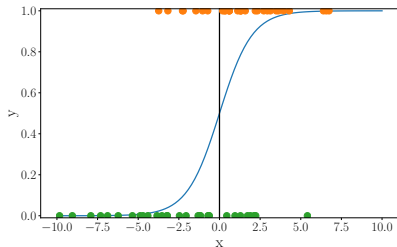
Model Specification (Logistic Regression)

- ▶ Bernoulli likelihood

$$y \in \{0, 1\}$$

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \text{Ber}(y|\mu(\mathbf{x})),$$

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



- ▶ Label y depends on input location \mathbf{x} , i.e., $\mu(\mathbf{x})$ needs to be a function of \mathbf{x}

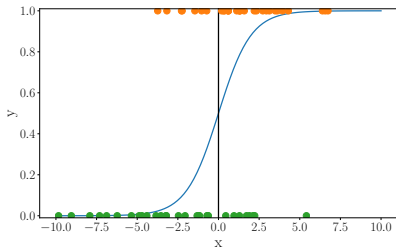
Model Specification (Logistic Regression)

- ▶ Bernoulli likelihood

$$y \in \{0, 1\}$$

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \text{Ber}(y|\mu(\mathbf{x})),$$

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



- ▶ Label y depends on input location \mathbf{x} , i.e., $\mu(\mathbf{x})$ needs to be a function of \mathbf{x}
- ▶ Idea: Linear model $\boldsymbol{\theta}^\top \mathbf{x}$ (as in linear regression)

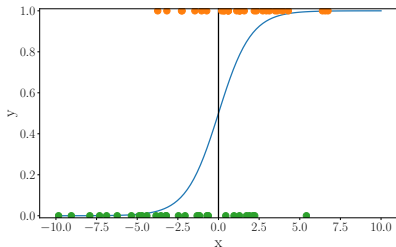
Model Specification (Logistic Regression)

- ▶ Bernoulli likelihood

$$y \in \{0, 1\}$$

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \text{Ber}(y|\mu(\mathbf{x})),$$

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



- ▶ Label y depends on input location \mathbf{x} , i.e., $\mu(\mathbf{x})$ needs to be a function of \mathbf{x}
- ▶ Idea: Linear model $\boldsymbol{\theta}^\top \mathbf{x}$ (as in linear regression)
- ▶ Ensure $0 \leq \mu(\mathbf{x}) \leq 1$

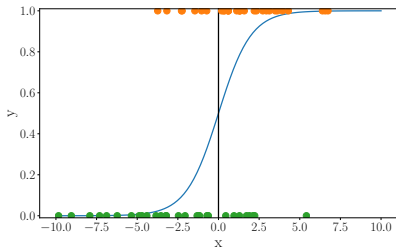
Model Specification (Logistic Regression)

- ▶ Bernoulli likelihood

$$y \in \{0, 1\}$$

$$p(y|x, \theta) = \text{Ber}(y|\mu(x)),$$

$$\mu(x) = p(y = 1|x) = \sigma(\theta^\top x)$$



- ▶ Label y depends on input location x , i.e., $\mu(x)$ needs to be a function of x
- ▶ Idea: Linear model $\theta^\top x$ (as in linear regression)
- ▶ Ensure $0 \leq \mu(x) \leq 1$
- ▶ Squash the linear combination through a function that guarantees this:

$$\mu(x) = \sigma(\theta^\top x)$$
$$\implies p(y|x, \theta) = \text{Ber}(y|\sigma(\theta^\top x))$$

Model Fitting

- ▶ Estimate model parameters θ (MLE or MAP)

Model Fitting

- ▶ Estimate model parameters θ (MLE or MAP)
- ▶ Likelihood (training data \mathbf{X}, \mathbf{y}):

Model Fitting

- ▶ Estimate model parameters θ (MLE or MAP)
- ▶ Likelihood (training data \mathbf{X}, \mathbf{y}):

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}, \theta) &= \prod_{n=1}^N \text{Ber}(y_n | \sigma(\theta^\top \mathbf{x}_n)) = \prod_{n=1}^N (\sigma(\theta^\top \mathbf{x}_n))^{y_n} (1 - \sigma(\theta^\top \mathbf{x}_n))^{1-y_n} \\ &= \prod_{n=1}^N \mu_n^{y_n} (1 - \mu_n)^{1-y_n} \\ \mu_n &:= \sigma(\theta^\top \mathbf{x}_n) \end{aligned}$$

- ▶ **Negative log likelihood (cross-entropy):**

Model Fitting

- ▶ Estimate model parameters θ (MLE or MAP)
- ▶ Likelihood (training data \mathbf{X}, \mathbf{y}):

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}, \theta) &= \prod_{n=1}^N \text{Ber}(y_n | \sigma(\theta^\top \mathbf{x}_n)) = \prod_{n=1}^N (\sigma(\theta^\top \mathbf{x}_n))^{y_n} (1 - \sigma(\theta^\top \mathbf{x}_n))^{1-y_n} \\ &= \prod_{n=1}^N \mu_n^{y_n} (1 - \mu_n)^{1-y_n} \\ \mu_n &:= \sigma(\theta^\top \mathbf{x}_n) \end{aligned}$$

- ▶ **Negative log likelihood (cross-entropy):**

$$NLL = - \sum_{n=1}^N y_n \log \mu_n + (1 - y_n) \log(1 - \mu_n)$$

Model Fitting (2)

- ▶ Derivative of sigmoid w.r.t. its argument:

$$\sigma(z_n) = \frac{1}{1 + \exp(-z_n)}$$
$$\implies \frac{d\sigma(z_n)}{dz_n} =$$

Model Fitting (2)

- Derivative of sigmoid w.r.t. its argument:

$$\begin{aligned}\sigma(z_n) &= \frac{1}{1 + \exp(-z_n)} \\ \implies \frac{d\sigma(z_n)}{dz_n} &= \frac{\exp(-z_n)}{(1 + \exp(-z_n))^2} = \sigma(z_n)(1 - \sigma(z_n))\end{aligned}$$

Model Fitting (2)

- ▶ Derivative of sigmoid w.r.t. its argument:

$$\sigma(z_n) = \frac{1}{1 + \exp(-z_n)}$$
$$\implies \frac{d\sigma(z_n)}{dz_n} = \frac{\exp(-z_n)}{(1 + \exp(-z_n))^2} = \sigma(z_n)(1 - \sigma(z_n))$$

- ▶ Gradient of the negative log-likelihood:

$$\frac{dNLL}{d\theta} = - \sum_{n=1}^N \left(y_n \frac{1}{\mu_n} - (1 - y_n) \frac{1}{1 - \mu_n} \right) \frac{d\mu_n}{d\theta}$$
$$\frac{d\mu_n}{d\theta} =$$

Model Fitting (2)

- ▶ Derivative of sigmoid w.r.t. its argument:

$$\begin{aligned}\sigma(z_n) &= \frac{1}{1 + \exp(-z_n)} \\ \implies \frac{d\sigma(z_n)}{dz_n} &= \frac{\exp(-z_n)}{(1 + \exp(-z_n))^2} = \sigma(z_n)(1 - \sigma(z_n))\end{aligned}$$

- ▶ Gradient of the negative log-likelihood:

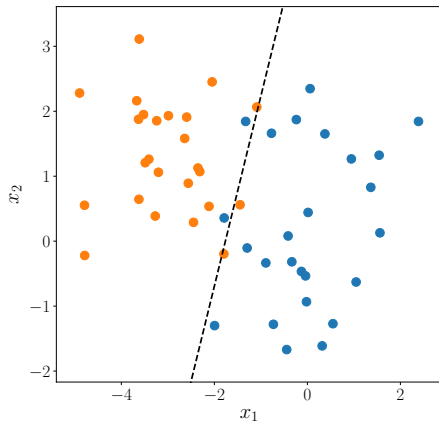
$$\begin{aligned}\frac{dNLL}{d\theta} &= - \sum_{n=1}^N \left(y_n \frac{1}{\mu_n} - (1 - y_n) \frac{1}{1 - \mu_n} \right) \frac{d\mu_n}{d\theta} \\ \frac{d\mu_n}{d\theta} &= \frac{d}{d\theta} \sigma(\underbrace{\theta^\top \mathbf{x}_n}_{z_n}) = \frac{d\sigma(z_n)}{dz_n} \frac{dz_n}{d\theta} = \sigma(z_n)(1 - \sigma(z_n)) \mathbf{x}_n^\top\end{aligned}$$

Model Fitting (3)

$$\frac{dNLL}{d\theta} = (\boldsymbol{\mu} - \mathbf{y})^\top \mathbf{X}$$
$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top$$

- ▶ No closed-form solution ➤ Gradient descent methods
- ▶ Unique global optimum exists

Example

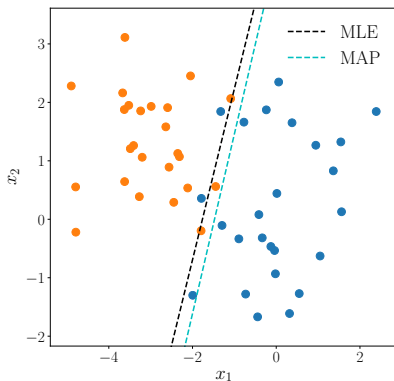


$$p(y|x, \theta) = \text{Ber}(\sigma(\theta_0 + \theta_1 x_1 + \theta_2 x_2))$$

Comments on Maximum Likelihood

- ▶ If the classes are linearly separable, the decision boundary is not unique and the likelihood will tend to infinity
- ▶ Overfitting is again a problem when we work with features $\phi(x)$ instead of x
- ▶ Maximum a posteriori estimation can address these issues to some degree

MAP Estimation



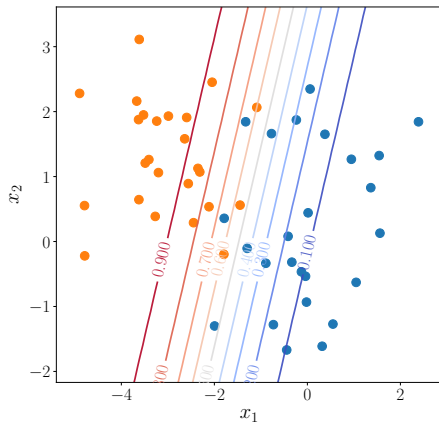
- ▶ Log-posterior:

$$\log p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{y}) = \log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) + \log p(\boldsymbol{\theta}) + \text{const}$$

- ▶ **No closed-form solution** for $\boldsymbol{\theta}_{\text{MAP}}$

▶▶ Numerical maximization of the log-posterior

Predictive Labels



$$p(y = 1 | \mathbf{x}, \boldsymbol{\theta}_{\text{MAP}}) = \text{Ber}(\sigma(\mathbf{x}^\top \boldsymbol{\theta}_{\text{MAP}}))$$

Bayesian Logistic Regression

Objective

For a given (i.i.d.) dataset $\mathcal{D} := \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ compute a posterior distribution on the parameters θ

Bayesian Logistic Regression

Objective

For a given (i.i.d.) dataset $\mathcal{D} := \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ compute a posterior distribution on the parameters $\boldsymbol{\theta}$

- ▶ Choose Gaussian prior $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} | \mathbf{0}, \mathbf{S}_0)$
- ▶ Posterior (via Bayes' theorem):

Bayesian Logistic Regression

Objective

For a given (i.i.d.) dataset $\mathcal{D} := \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ compute a posterior distribution on the parameters $\boldsymbol{\theta}$

- ▶ Choose Gaussian prior $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} | \mathbf{0}, \mathbf{S}_0)$
- ▶ Posterior (via Bayes' theorem):

$$p(\boldsymbol{\theta} | \mathcal{D}) = \frac{p(\boldsymbol{\theta}) p(\mathbf{y} | \boldsymbol{\theta}, \mathbf{X})}{p(\mathbf{y} | \mathbf{X})} = \frac{\mathcal{N}(\boldsymbol{\theta} | \mathbf{0}, \mathbf{S}_0) \prod_{n=1}^N \text{Ber}(\sigma(\mathbf{x}_n^\top \boldsymbol{\theta}))}{\int \mathcal{N}(\boldsymbol{\theta} | \mathbf{0}, \mathbf{S}_0) \prod_{n=1}^N \text{Ber}(\sigma(\mathbf{x}_n^\top \boldsymbol{\theta})) d\boldsymbol{\theta}}$$

Bayesian Logistic Regression

Objective

For a given (i.i.d.) dataset $\mathcal{D} := \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ compute a posterior distribution on the parameters $\boldsymbol{\theta}$

- ▶ Choose Gaussian prior $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} | \mathbf{0}, \mathbf{S}_0)$
- ▶ Posterior (via Bayes' theorem):

$$p(\boldsymbol{\theta} | \mathcal{D}) = \frac{p(\boldsymbol{\theta}) p(\mathbf{y} | \boldsymbol{\theta}, \mathbf{X})}{p(\mathbf{y} | \mathbf{X})} = \frac{\mathcal{N}(\boldsymbol{\theta} | \mathbf{0}, \mathbf{S}_0) \prod_{n=1}^N \text{Ber}(\sigma(\mathbf{x}_n^\top \boldsymbol{\theta}))}{\int \mathcal{N}(\boldsymbol{\theta} | \mathbf{0}, \mathbf{S}_0) \prod_{n=1}^N \text{Ber}(\sigma(\mathbf{x}_n^\top \boldsymbol{\theta})) d\boldsymbol{\theta}}$$

- ▶ **No analytic solution**
 - ▶ Approximations necessary

Laplace Approximation

- ▶ Objective: Approximate an unknown distribution

$$p(\mathbf{x}) \propto \exp(-E(\mathbf{x})) =: \tilde{p}(\mathbf{x})$$

with a Gaussian distribution $q(\mathbf{x})$.

Laplace Approximation

- ▶ Objective: Approximate an unknown distribution

$$p(\mathbf{x}) \propto \exp(-E(\mathbf{x})) =: \tilde{p}(\mathbf{x})$$

with a Gaussian distribution $q(\mathbf{x})$.

- ▶ Idea: Taylor-series expansion of $-\log \tilde{p}(\mathbf{x}) = E(\mathbf{x})$ around the mode \mathbf{x}^* (MAP estimate)

Laplace Approximation

- ▶ Objective: Approximate an unknown distribution

$$p(\mathbf{x}) \propto \exp(-E(\mathbf{x})) =: \tilde{p}(\mathbf{x})$$

with a Gaussian distribution $q(\mathbf{x})$.

- ▶ Idea: Taylor-series expansion of $-\log \tilde{p}(\mathbf{x}) = E(\mathbf{x})$ around the mode \mathbf{x}^* (MAP estimate)

$$-\log \tilde{p}(\mathbf{x}) \approx E(\mathbf{x}^*) + \mathbf{J}(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_*)^\top \mathbf{H}(\mathbf{x}_*)(\mathbf{x} - \mathbf{x}^*),$$

\mathbf{J} : Jacobian, \mathbf{H} : Hessian

Laplace Approximation

- ▶ Objective: Approximate an unknown distribution

$$p(\mathbf{x}) \propto \exp(-E(\mathbf{x})) =: \tilde{p}(\mathbf{x})$$

with a Gaussian distribution $q(\mathbf{x})$.

- ▶ Idea: Taylor-series expansion of $-\log \tilde{p}(\mathbf{x}) = E(\mathbf{x})$ around the mode \mathbf{x}^* (MAP estimate)

$$-\log \tilde{p}(\mathbf{x}) \approx E(\mathbf{x}^*) + J(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_*)^\top \mathbf{H}(\mathbf{x}_*)(\mathbf{x} - \mathbf{x}^*),$$

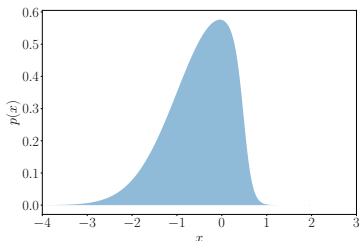
J : Jacobian, \mathbf{H} : Hessian

- ▶ $J(\mathbf{x}^*) = \mathbf{0}^\top$ because \mathbf{x}^* is a stationary point (mode) of $\log \tilde{p}$

$$\tilde{p}(\mathbf{x}) \approx \exp(-E(\mathbf{x}^*)) \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{x}_*)^\top \mathbf{H}(\mathbf{x}_*)(\mathbf{x} - \mathbf{x}^*)\right)$$

$$\propto \mathcal{N}(\mathbf{x} | \mathbf{x}^*, \mathbf{H}^{-1}) =: q(\mathbf{x})$$

Laplace Approximation: Example

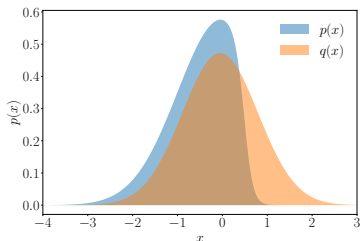
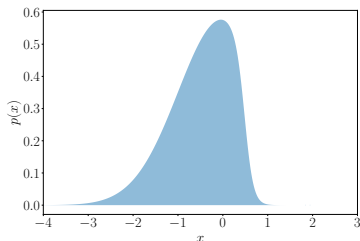


- ▶ Unnormalized distribution:

$$\tilde{p}(x) = \exp\left(-\frac{1}{2}x^2\right)\sigma(ax + b)$$

▶▶ Discuss with your neighbors

Laplace Approximation: Example



- Unnormalized distribution:

$$\tilde{p}(x) = \exp\left(-\frac{1}{2}x^2\right)\sigma(ax + b)$$

$$q(x) = \mathcal{N}\left(x \mid x^*, (1 + a^2\mu_*(1 - \mu_*))^{-1}\right), \quad \mu_* := \sigma(ax_* + b)$$

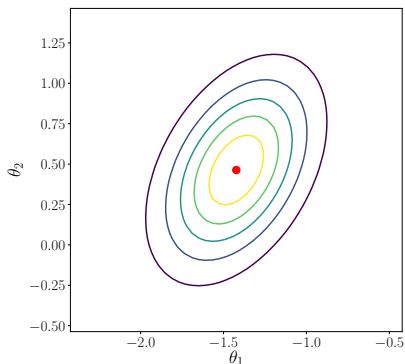
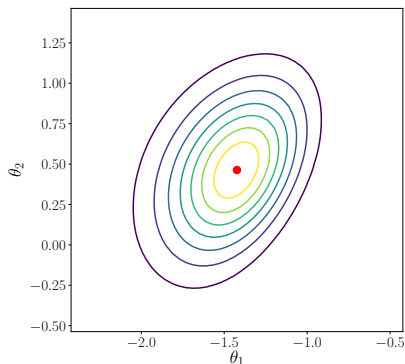
Laplace Approximation: Properties

- ▶ Only need to know the **unnormalized distribution** \tilde{p}
- ▶ Finding the mode: numerical methods (optimization problem)
- ▶ **Captures only local properties** of the distribution
- ▶ Multimodal distributions: Approximation will be different depending on which mode we are in (**not unique**)

Laplace Approximation: Properties

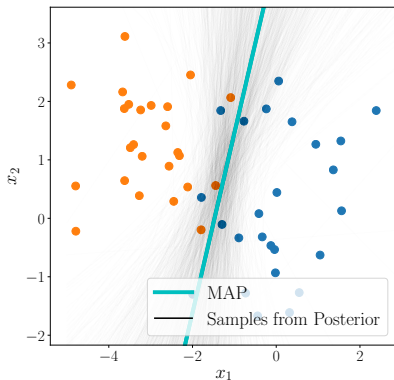
- ▶ Only need to know the **unnormalized distribution** \tilde{p}
- ▶ Finding the mode: numerical methods (optimization problem)
- ▶ **Captures only local properties** of the distribution
- ▶ Multimodal distributions: Approximation will be different depending on which mode we are in (**not unique**)
- ▶ For large datasets, we would expect the posterior to converge to a Gaussian (central limit theorem)
 - ▶▶ Laplace approximation should work well in this case

Posterior Approximation



- ▶ Left: true parameter posterior
- ▶ Right: Laplace approximation

Posterior Decision Boundary



- ▶ Parameter samples θ_i drawn from Laplace approximation $q(\theta)$ of posterior $p(\theta|\mathbf{X})$
- ▶ Decision boundary drawn for each θ_i

Predictions

Assume a Gaussian distribution $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ on the parameters (e.g., Laplace approximation of the posterior). Then:

$$\begin{aligned} p(\mathbf{y}|\mathbf{x}) &= \int p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta} \\ &= \int \text{Ber}(\sigma(\boldsymbol{\theta}^\top \mathbf{x})) \mathcal{N}(\boldsymbol{\theta} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\boldsymbol{\theta} \\ &= \mathbb{E}_{\boldsymbol{\theta}}[\text{Ber}(\sigma(\boldsymbol{\theta}^\top \mathbf{x}))] \end{aligned}$$

Predictions

Assume a Gaussian distribution $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ on the parameters (e.g., Laplace approximation of the posterior). Then:

$$\begin{aligned} p(\mathbf{y}|\mathbf{x}) &= \int p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta} \\ &= \int \text{Ber}(\sigma(\boldsymbol{\theta}^\top \mathbf{x})) \mathcal{N}(\boldsymbol{\theta} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\boldsymbol{\theta} \\ &= \mathbb{E}_{\boldsymbol{\theta}}[\text{Ber}(\sigma(\boldsymbol{\theta}^\top \mathbf{x}))] \end{aligned}$$

▶▶ **Integral intractable**

Predictions

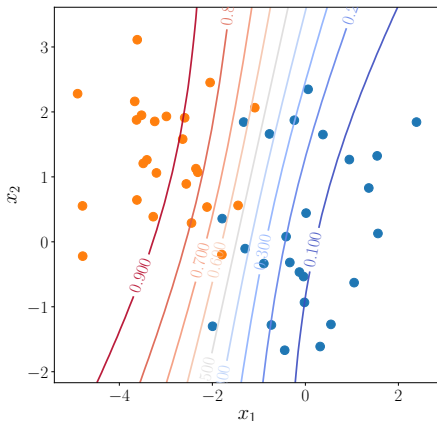
Assume a Gaussian distribution $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ on the parameters (e.g., Laplace approximation of the posterior). Then:

$$\begin{aligned} p(\mathbf{y}|\mathbf{x}) &= \int p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta} \\ &= \int \text{Ber}(\sigma(\boldsymbol{\theta}^\top \mathbf{x})) \mathcal{N}(\boldsymbol{\theta} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\boldsymbol{\theta} \\ &= \mathbb{E}_{\boldsymbol{\theta}}[\text{Ber}(\sigma(\boldsymbol{\theta}^\top \mathbf{x}))] \end{aligned}$$

▶▶ Integral intractable

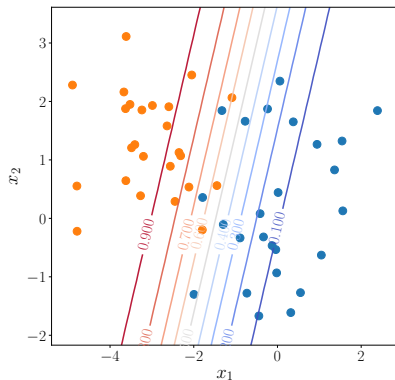
- ▶ “Plug-in approximation”: use posterior mean (MAP estimate)
 $\mathbb{E}[\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}]$
- ▶ Monte Carlo estimate (sampling from $p(\boldsymbol{\theta})$ is easy)

Predictions (2)

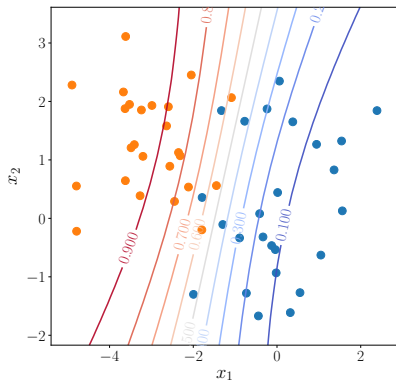


1. Samples from Laplace approximation of the posterior
2. Monte-Carlo estimate of label prediction

Comparison with MAP Predictions



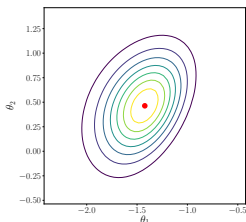
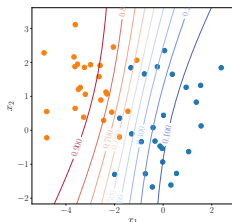
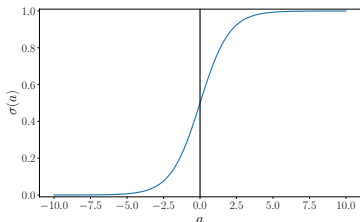
(a) MAP



(b) Bayesian Logistic Regression

► Predictive labels

Summary



- ▶ Binary classification problems
- ▶ Linear model with non-Gaussian likelihood
- ▶ Implicit modeling assumptions
- ▶ Parameter estimation (MLE, MAP) no longer in closed form
- ▶ Bayesian logistic regression with Laplace approximation of the posterior

References I

- [1] C. M. Bishop. *Pattern Recognition and Machine Learning*. Information Science and Statistics. Springer-Verlag, 2006.
- [2] K. P. Murphy. *Machine Learning: A Probabilistic Perspective*. MIT Press, Cambridge, MA, USA, 2012.