

ECE276A: Sensing & Estimation in Robotics

Lecture 5: Logistic Regression

Instructor:

Nikolay Atanasov: nataanasov@ucsd.edu

Teaching Assistants:

Mo Shan: moshan@eng.ucsd.edu

Arash Asgharivaskasi: aasghari@eng.ucsd.edu



Supervised Learning

- ▶ **Data:** a set $D := \{\mathbf{x}_i, y_i\}_{i=1}^n$ of **iid** examples $\mathbf{x}_i \in \mathbb{R}^d$ with associated scalar labels y_i generated from an unknown joint pdf $p_*(y, \mathbf{x})$
- ▶ The training dataset is often also written in matrix notation, $D = (\mathbf{X}, \mathbf{y})$, with $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{y} \in \mathbb{R}^n$
- ▶ **Generative model:** choose model $p(y, \mathbf{x}; \omega)$ with parameters ω to approximate the unknown data-generating pdf
- ▶ **Discriminative model:** choose model $p(y|\mathbf{x}; \omega)$ with parameters ω to approximate the unknown label-generating pdf
- ▶ Optimize ω using $D = (\mathbf{X}, \mathbf{y})$:
 - ▶ **Maximum Likelihood Estimation (MLE):** maximize the likelihood of the data D given the parameters ω
 - ▶ **Maximum A Posteriori (MAP):** maximize the likelihood of the parameters ω given the data D
 - ▶ **Bayesian Inference:** estimate the whole distribution of the parameters ω given the data D

Parametric Learning

► Maximum Likelihood Estimation (MLE):

MLE	Discriminative Model	Generative Model
Training	$\omega^* \in \arg \max_{\omega} p(\mathbf{y} X, \omega)$	$\omega^* \in \arg \max_{\omega} p(\mathbf{y}, X \omega)$
Testing	$y_* \in \arg \max_y p(y \mathbf{x}_*, \omega^*)$	$y_* \in \arg \max_y p(y, \mathbf{x}_* \omega^*)$

► Maximum A Posteriori (MAP):

MAP	Discriminative Model	Generative Model
Training	$\omega^* \in \arg \max_{\omega} p(\omega \mathbf{y}, X)$ $= \arg \max_{\omega} p(\mathbf{y} X, \omega)p(\omega X)$	$\omega^* \in \arg \max_{\omega} p(\omega \mathbf{y}, X)$ $= \arg \max_{\omega} p(\mathbf{y}, X \omega)p(\omega)$
Testing	$y_* \in \arg \max_y p(y \mathbf{x}_*, \omega^*)$	$y_* \in \arg \max_y p(y, \mathbf{x}_* \omega^*)$

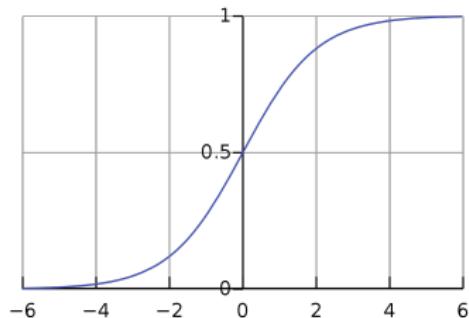
► Bayesian Inference:

BI	Discriminative Model	Generative Model
Training	$p(\omega \mathbf{y}, X) \propto p(\mathbf{y} X, \omega)p(\omega X)$	$p(\omega \mathbf{y}, X) \propto p(\mathbf{y}, X \omega)p(\omega)$
Testing	$p(y_* \mathbf{x}_*, \mathbf{y}, X) = \int p(y_* \mathbf{x}_*, \omega)p(\omega \mathbf{y}, X)d\omega$	$p(y_*, \mathbf{x}_* \mathbf{y}, X) = \int p(y_*, \mathbf{x}_* \omega)p(\omega \mathbf{y}, X)d\omega$

Logistic Sigmoid Function

- ▶ Useful for converting continuous (regression) preferences $z \in \mathbb{R}$ into a Bernoulli probability mass function, which can be used as a probabilistic model for **binary classification**, i.e., $y \in \{-1, 1\}$

$$\sigma(z) := \frac{1}{1 + \exp(-z)} = \frac{\exp(z)}{\exp(z) + \exp(0)}$$



- ▶ Properties:
 - ▶ $\sigma(z) = 1 - \sigma(-z) = \frac{1}{2} + \frac{1}{2} \tanh\left(\frac{z}{2}\right)$
 - ▶ $\sigma'(z) = \sigma(z)(1 - \sigma(z)) = \sigma(z)\sigma(-z)$
 - ▶ $\int \sigma(z)dz = \log(1 + \exp(z))$, which is known as the **softplus** function

Softmax Function

- ▶ Useful for converting continuous (regression) preferences $\mathbf{z} \in \mathbb{R}^K$ into a categorical probability mass function, which can be used as a probabilistic model for **classification**, i.e., $y \in \{1, \dots, K\}$

$$\mathbf{s}(\mathbf{z}) := \begin{bmatrix} \frac{\exp(z_1)}{\sum_j \exp(z_j)} & \cdots & \frac{\exp(z_K)}{\sum_j \exp(z_j)} \end{bmatrix} = \frac{e^{\mathbf{z}}}{\mathbf{1}^\top e^{\mathbf{z}}} \in \mathbb{R}^K$$

- ▶ Properties:
 - ▶ $\mathbf{s}(\mathbf{z}) = \mathbf{s}(\mathbf{z} - c\mathbf{1})$, where $c \in \mathbb{R}$ is any constant, e.g., $c = \max_i z_i$ is useful for numerical conditioning
 - ▶ $\log \frac{s_i(\mathbf{z})}{s_j(\mathbf{z})} = z_i - z_j$
 - ▶ $\frac{ds_i}{d\mathbf{z}_j} = \begin{cases} s_i(1 - s_i) & \text{if } i = j \\ -s_i s_j & \text{else} \end{cases}$
 - ▶ $\nabla_{\mathbf{z}} \log \left(\sum_{i=1}^K \exp(z_i) \right) = \mathbf{s}(\mathbf{z})$

Discriminative Classification via a Logistic Model

- Given an example $\mathbf{x} \in \mathbb{R}^d$, use a logistic sigmoid function to map the continuous value $\omega^\top \mathbf{x}$ to a Bernoulli probability mass function for a binary label $y \in \{-1, 1\}$:

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

- For $p(y | \mathbf{x}, \omega)$ to be a valid pmf, it needs to sum to 1 over $\{-1, 1\}$:

$$p(y = -1 | \mathbf{x}, \omega) = 1 - p(y = 1 | \mathbf{x}, \omega) = 1 - \sigma(\omega^\top \mathbf{x}) \xrightarrow[\text{properties}]{\text{sigmoid}} \sigma(-\omega^\top \mathbf{x})$$

- Combine the two cases to obtain the final discriminative model for the pmf of y conditioned on the example \mathbf{x} and the parameters ω :

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

Discriminative Classification via a Logistic Model

- ▶ **Logistic regression:** approximates the unknown label-generating pdf with a logistic sigmoid function:

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

- ▶ Since the data $D = (X, \mathbf{y})$ are **iid**, the joint data likelihood is:

$$p(\mathbf{y}|X, \omega) = \prod_{i=1}^n \sigma(y_i \mathbf{x}_i^\top \omega) = \prod_{i=1}^n \frac{1}{1 + \exp(-y_i \mathbf{x}_i^\top \omega)}$$

- ▶ Leads to these MLE and MAP (with $\omega \sim \mathcal{N}(0, \Lambda)$) estimates for ω :

$$\omega_{MLE} = \arg \max_{\omega} \log p(\mathbf{y} | X, \omega) = \arg \min_{\omega} \sum_{i=1}^n \log \left(1 + \exp(-y_i \mathbf{x}_i^\top \omega) \right)$$

$$\omega_{MAP} = \arg \max_{\omega} \log p(\mathbf{y} | X, \omega) + \log p(\omega)$$

$$= \arg \min_{\omega} \sum_{i=1}^n \log \left(1 + \exp(-y_i \mathbf{x}_i^\top \omega) \right) + \frac{1}{2} \omega^\top \Lambda^{-1} \omega$$

Discriminative Classification via a Logistic Model

- ▶ $\nabla_{\omega} (-\log p(\mathbf{y} | X, \omega)) = 0$ does not have a closed-form solution so we need to use an iterative optimization algorithm like gradient descent
- ▶ The negative log-likelihood $-\log p(\mathbf{y} | X, \omega)$ is **convex** in ω :
 - ▶ The composition of an affine function $f_i(\omega) := -y_i \mathbf{x}_i^\top \omega$ and a convex function $g(z) := \log(1 + \exp^z)$ is convex
 - ▶ The sum of convex functions $\sum_{i=1}^n g(f_i(\omega))$ is convex
- ▶ The negative log-likelihood can be minimized iteratively to obtain a **global** minimum:

$$\begin{aligned}\omega_{MLE}^{(t+1)} &= \omega_{MLE}^{(t)} - \alpha \nabla_{\omega} (-\log p(\mathbf{y}|X, \omega)) \Big|_{\omega=\omega_{MLE}^{(t)}} \\ &= \omega_{MLE}^{(t)} - \alpha \sum_{i=1}^n \frac{1}{1 + \exp(-y_i \mathbf{x}_i^\top \omega_{MLE}^{(t)})} \exp(-y_i \mathbf{x}_i^\top \omega_{MLE}^{(t)}) (-y_i \mathbf{x}_i) \\ &= \omega_{MLE}^{(t)} + \alpha \sum_{i=1}^n y_i \mathbf{x}_i (1 - \sigma(y_i \mathbf{x}_i^\top \omega_{MLE}^{(t)}))\end{aligned}$$

Logistic Regression Summary

- ▶ **Discriminative model:** $p(\mathbf{y}|X, \omega)$ for binary labels $\mathbf{y} \in \{-1, 1\}^n$:

$$p(\mathbf{y}|X, \omega) = \prod_{i=1}^n \sigma(y_i \mathbf{x}_i^\top \omega) = \prod_{i=1}^n \frac{1}{1 + \exp(-y_i \mathbf{x}_i^\top \omega)}$$

- ▶ **Training:** given data $D = (X, \mathbf{y})$, optimize the model parameters:

- ▶ MLE: $\omega_{MLE}^{(t+1)} = \omega_{MLE}^{(t)} + \alpha \sum_{i=1}^n y_i \mathbf{x}_i (1 - \sigma(y_i \mathbf{x}_i^\top \omega_{MLE}^{(t)}))$

- ▶ MAP: $\omega_{MAP}^{(t+1)} = \omega_{MAP}^{(t)} + \alpha \left(\sum_{i=1}^n y_i \mathbf{x}_i (1 - \sigma(y_i \mathbf{x}_i^\top \omega_{MAP}^{(t)})) - \Lambda^{-1} \omega_{MAP}^{(t)} \right)$

- ▶ **Testing:** given a test example $\mathbf{x}_* \in \mathbb{R}^d$, use the optimized parameters ω^* to predict the label:

$$y_* = \begin{cases} 1 & \mathbf{x}_*^\top \omega^* \geq 0 \\ -1 & \mathbf{x}_*^\top \omega^* < 0 \end{cases}$$

- ▶ Logistic regression generates a **linear decision boundary**, i.e., the values of \mathbf{x}_* for which the predictions $y_* = 1$ and $y_* = -1$ are equally likely are:

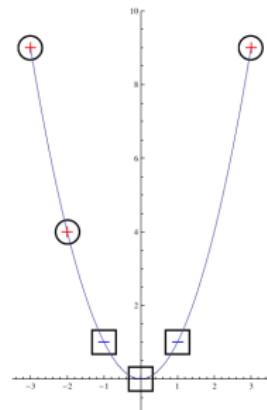
$$1 = \frac{p(y = 1 | \mathbf{x}_*, \omega^*)}{p(y = -1 | \mathbf{x}_*, \omega^*)} = \frac{1 + \exp(\mathbf{x}_*^\top \omega^*)}{1 + \exp(-\mathbf{x}_*^\top \omega^*)} \quad \Leftrightarrow \quad \mathbf{x}_*^\top \omega^* = 0$$

Logistic Regression Example

- ▶ Consider the same data as before:

$$X = \begin{bmatrix} -3 & 9 & 1 \\ -2 & 4 & 1 \\ -1 & 1 & 1 \\ 0 & 0 & 1 \\ 1 & 1 & 1 \\ 3 & 9 & 1 \end{bmatrix} \in \mathbb{R}^{n \times d}$$

$$\mathbf{y} = \begin{bmatrix} +1 \\ +1 \\ -1 \\ -1 \\ -1 \\ +1 \end{bmatrix} \in \mathbb{R}^n$$

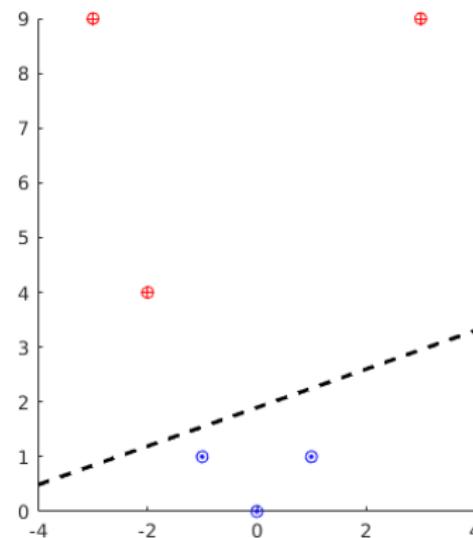


Logistic Regression Example

- ▶ **Training:** start with $\omega_{MLE}^{(0)} = \mathbf{0} \in \mathbb{R}^3$ and iterate:

$$\omega_{MLE}^{(t+1)} = \omega_{MLE}^{(t)} + \alpha \sum_{i=1}^n y_i \mathbf{x}_i (1 - \sigma(y_i \mathbf{x}_i^\top \omega_{MLE}^{(t)}))$$

- ▶ After 10 iterations with $\alpha = 0.1$, we have: $\omega_{MLE}^{(10)} = \begin{bmatrix} 0.2115 \\ -0.6015 \\ 1.1408 \end{bmatrix}$
- ▶ **Testing:** the decision boundary is a line with equation $0 = \mathbf{x}^\top \omega$:



Logistic Regression vs Gaussian Naïve Bayes

- ▶ Logistic regression generates a **linear decision boundary**: $\omega^\top \mathbf{x} = 0$.
- ▶ Gaussian Naïve Bayes generates a **quadratic decision boundary**. It looks like an ellipse, parabola, or hyperbola in 2-D:

$$\log \frac{\theta_-^2}{\theta_+^2} + \sum_{I=1}^d \log \frac{\sigma_{+,I}^2}{\sigma_{-,I}^2} + \frac{(x_I - \mu_{+,I})^2}{\sigma_{+,I}^2} - \frac{(x_I - \mu_{-,I})^2}{\sigma_{-,I}^2} = 0$$

- ▶ When the variance is shared among the classes, i.e., $\sigma_{-,I} = \sigma_{+,I} = \sigma_I$, Gaussian Naïve Bayes generates a **linear decision boundary**:

$$\log \frac{\theta_-^2}{\theta_+^2} + \sum_{I=1}^d \frac{2(\mu_{-,I} - \mu_{+,I})}{\sigma_I^2} x_I + \frac{(\mu_{+,I}^2 - \mu_{-,I}^2)}{\sigma_I^2} = 0$$

- ▶ Logistic regression has **lower bias** but **higher variance** than Gaussian Naïve Bayes.

K-ary Logistic Regression

- **Logistic regression with K -classes:** approximates the unknown label-generating pdf for $\mathbf{y} \in \{1, \dots, K\}^n$ with a softmax function with parameters $W \in \mathbb{R}^{K \times d}$:

$$p(\mathbf{y}|X, W) = \prod_{i=1}^n \mathbf{e}_{y_i}^\top \mathbf{s}(W\mathbf{x}_i) := \prod_{i=1}^n \mathbf{e}_{y_i}^\top \frac{\exp(W\mathbf{x}_i)}{\mathbf{1}^\top \exp(W\mathbf{x}_i)}$$

where \mathbf{e}_j is the j -th standard basis vector, e.g., $\mathbf{e}_1 := [1 \ 0 \ \cdots \ 0]^\top$

- To optimize the parameters W via MLE, we need to compute the gradient of the data log-likelihood:

$$\begin{aligned} W_{MLE}^{(t+1)} &= W_{MLE}^{(t)} + \alpha \left(\nabla_W [\log p(\mathbf{y} | X, W)] \Big|_{W=W_{MLE}^{(t)}} \right) \\ &= W_{MLE}^{(t)} + \alpha \left(\sum_{i=1}^n \left(\mathbf{e}_{y_i} - \mathbf{s}(W_{MLE}^{(t)} \mathbf{x}_i) \right) \mathbf{x}_i^\top \right) \end{aligned}$$

K-ary Logistic Regression MLE

- ▶ Taking the gradient of $\log p(\mathbf{y} | \mathbf{X}, \mathbf{W})$ with respect to \mathbf{W} gives:

$$\nabla_{\mathbf{W}} \left[\sum_{i=1}^n \log \mathbf{e}_{y_i}^\top \mathbf{s}(\mathbf{W}\mathbf{x}_i) \right] = \sum_{i=1}^n \frac{1}{\mathbf{e}_{y_i}^\top \mathbf{s}(\mathbf{W}\mathbf{x}_i)} \nabla_{\mathbf{W}} \left[\mathbf{e}_{y_i}^\top \mathbf{s}(\mathbf{W}\mathbf{x}_i) \right]$$

- ▶ Let $\mathbf{z}_i := \mathbf{W}\mathbf{x}_i$. The derivative of the last term via the chain rule is:

$$\begin{aligned} \frac{\partial}{\partial \mathbf{z}_i} \mathbf{e}_{y_i}^\top \mathbf{s}(\mathbf{z}_i) \frac{\partial}{\partial \mathbf{W}} \mathbf{W}\mathbf{x}_i &= \frac{\partial}{\partial \mathbf{W}} \mathbf{e}_{y_i}^\top \frac{\partial \mathbf{s}(\mathbf{z}_i)}{\partial \mathbf{z}_i} \mathbf{W}\mathbf{x}_i = \frac{\partial}{\partial \mathbf{W}} \text{tr} \left(\mathbf{x}_i \mathbf{e}_{y_i}^\top \frac{\partial \mathbf{s}(\mathbf{z}_i)}{\partial \mathbf{z}_i} \mathbf{W} \right) \\ &= \mathbf{x}_i \mathbf{e}_{y_i}^\top \frac{\partial \mathbf{s}(\mathbf{z}_i)}{\partial \mathbf{z}_i} = \mathbf{x}_i \mathbf{e}_{y_i}^\top \frac{\partial \mathbf{s}}{\partial \mathbf{z}_i} (\mathbf{W}\mathbf{x}_i) \end{aligned}$$

- ▶ Arrange the derivatives of the softmax in matrix form:

$$\frac{d\mathbf{s}(\mathbf{z})}{d\mathbf{z}} = \begin{bmatrix} s_1 - s_1^2 & -s_1 s_2 & \cdots & -s_1 s_d \\ -s_2 s_1 & s_2 - s_2^2 & \cdots & -s_2 s_d \\ \vdots & & \ddots & \vdots \\ -s_d s_1 & -s_d s_2 & \cdots & -s_d s_d \end{bmatrix} = \mathbf{diag}(\mathbf{s}(\mathbf{z})) - \mathbf{s}(\mathbf{z})\mathbf{s}(\mathbf{z})^\top$$

K-ary Logistic Regression MLE

- The gradient is the transpose of the derivative:

$$\begin{aligned}\nabla_W [\log p(\mathbf{y} | \mathbf{X}, \mathbf{W})] &= \sum_{i=1}^n \frac{1}{\mathbf{e}_{y_i}^\top \mathbf{s}(W\mathbf{x}_i)} (\text{diag}(\mathbf{s}(W\mathbf{x}_i)) - \mathbf{s}(W\mathbf{x}_i)\mathbf{s}(W\mathbf{x}_i)^\top) \mathbf{e}_{y_i} \mathbf{x}_i^\top \\ &= \sum_{i=1}^n (\mathbf{e}_{y_i} - \mathbf{s}(W\mathbf{x}_i)) \mathbf{x}_i^\top \in \mathbb{R}^{K \times d}\end{aligned}$$

- MLE leads to gradient ascent of the form:

$$W_{MLE}^{(t+1)} = W_{MLE}^{(t)} + \alpha \left(\sum_{i=1}^n (\mathbf{e}_{y_i} - \mathbf{s}(W_{MLE}^{(t)} \mathbf{x}_i)) \mathbf{x}_i^\top \right)$$

- MAP with $\mathbf{W} \sim \mathcal{N}(0, \lambda I_{Kd \times Kd})$ leads to gradient ascent of the form:

$$W_{MAP}^{(t+1)} = W_{MAP}^{(t)} + \alpha \left(\sum_{i=1}^n (\mathbf{e}_{y_i} - \mathbf{s}(W_{MAP}^{(t)} \mathbf{x}_i)) \mathbf{x}_i^\top \right) - \alpha \lambda^{-1} W_{MAP}^{(t)}$$