

Logistic regression (part II)

Prof. Marco Alvarez, Computer Science
University of Rhode Island

Solving logistic regression

Maximum likelihood estimation (MLE)

► MLE principle

- statistical method used to estimate the parameters of a probability distribution (e.g., mean and standard deviation for normal distributions)
 - if we observe \mathcal{D} , choose the parameters that make \mathcal{D} most probable
- many machine learning algorithms follow this principle

► The **likelihood** function

- probability of observing the data given some parameters:

$$\mathcal{L}(\mathbf{w}) = P(X; \mathbf{w})$$

- for independent and identically distributed observations:

$$\mathcal{L}(\mathbf{w}) = P(\mathbf{x}_1; \mathbf{w}) \cdot P(\mathbf{x}_2; \mathbf{w}) \cdot \dots \cdot P(\mathbf{x}_n; \mathbf{w})$$

MLE and logistic regression

► MLE objective

- find \mathbf{w} that **maximizes** the likelihood function

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \mathcal{L}(\mathbf{w})$$

► Logistic regression

- conditional data likelihood:

$$\mathcal{L}(\mathbf{w}) = \prod_{i=1}^n P(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w})$$

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \prod_{i=1}^n P(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w})$$

Solving logistic regression

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \prod_{i=1}^n P(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w})$$

$$= \arg \max_{\mathbf{w}} \log \left(\prod_{i=1}^n \frac{1}{1 + e^{-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)}}} \right)$$

$$= \arg \max_{\mathbf{w}} - \sum_{i=1}^n \log \left(1 + e^{-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)}} \right)$$

$$= \arg \max_{\mathbf{w}} - \sum_{i=1}^n \log \left(1 + e^{-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)}} \right)$$

$$\text{negative log likelihood} = \arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)}} \right) \text{error (loss)} e(h(\mathbf{x}), y)$$

Objective function

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)}} \right) \text{cross-entropy loss (over a dataset)}$$

► the loss function (objective) is convex

- however, no closed-form solution
- can use **gradient descent** or second-order methods
- coming soon ...

How to classify new data?

► Once the final hypothesis $h(\mathbf{x})$ is known ...

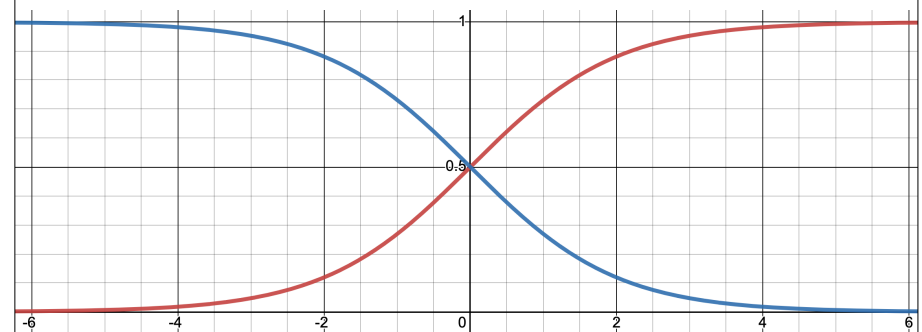
$$h(\mathbf{x}) = P(y = +1 | \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$

- predict label **+1** to input instance \mathbf{x}
 - if $p(+1 | \mathbf{x}) \geq 0.5$
- predict label **-1** to input instance \mathbf{x}
 - if $p(+1 | \mathbf{x}) < 0.5$

Decision boundary

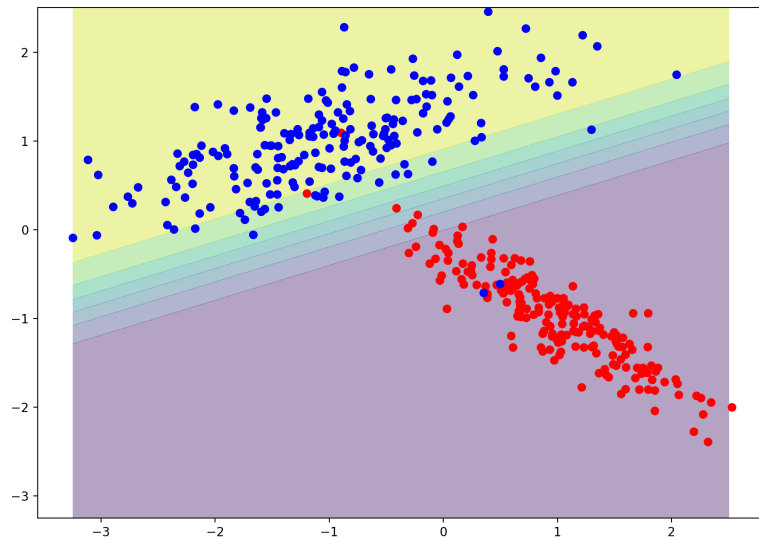
$$P(y = +1 | \mathbf{x}) = P(y = -1 | \mathbf{x}) = 0.5$$

$\sigma(\mathbf{w}^T \mathbf{x})$ $\sigma(-\mathbf{w}^T \mathbf{x})$



Logistic regression finds a linear decision boundary with $\mathbf{w}^T \mathbf{x} = 0$

Decision boundary



Final remarks

- ▶ Simple classifier with **probabilistic outputs**
- ▶ Loss function is convex
 - guaranteed global minimum
- ▶ Robust to overfitting
 - use regularization (coming soon)
- ▶ Offers interpretability to weights
 - feature importance
- ▶ Decision boundary is still linear