

CS480/680: Introduction to Machine Learning

Lecture 4: Logistic Regression

Hongyang Zhang

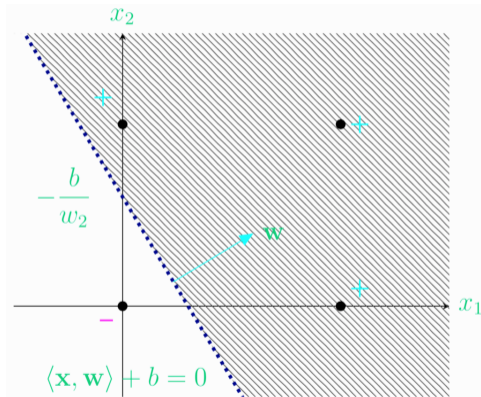


UNIVERSITY OF
WATERLOO

Jan 23, 2024

Predicting with Confidence

- **Caveat:** Though called “logistic regression”, it is for linear classification
- Recall that $\hat{y} = \text{sign}(\langle \mathbf{x}, \mathbf{w} \rangle)$ (assume using padding trick in this lecture)
- How confident we are about the prediction \hat{y} ?
- Use $|\langle \mathbf{x}, \mathbf{w} \rangle|$ (margin) as an indication
 - ▶ $\langle \mathbf{x}, \mathbf{w} \rangle$ is a.k.a. logit
 - ▶ in fact was used in perceptron
 - ▶ un-normalized w.r.t. the radius of data: hard to interpret
 - ▶ need ways to normalize it into $[0, 1]$
- **Better** idea: learn confidence directly



Max Likelihood Estimation (MLE)

- Let $\mathcal{Y} = \{0, 1\}$; Let's directly learn confidence $p(\mathbf{x}; \mathbf{w}) := \Pr(Y = 1|X = \mathbf{x})$
- Given $(\mathbf{x}_i, y_i), i = 1, \dots, n$, assume independence:

$$\Pr(Y_1 = y_1, \dots, Y_n = y_n | X_1 = \mathbf{x}_1, \dots, X_n = \mathbf{x}_n) = \prod_{i=1}^n \Pr(Y_i = y_i | X_i = \mathbf{x}_i)$$
$$\stackrel{\mathcal{Y}=\{0,1\}}{=} \prod_{i=1}^n [p(\mathbf{x}_i; \mathbf{w})]^{y_i} [1 - p(\mathbf{x}_i; \mathbf{w})]^{1-y_i}$$

- Maximizing the likelihood:

$$\max_{\mathbf{w}} \prod_{i=1}^n [p(\mathbf{x}_i; \mathbf{w})]^{y_i} [1 - p(\mathbf{x}_i; \mathbf{w})]^{1-y_i}$$

or minimizing the minus log-likelihood:

$$\min_{\mathbf{w}} \sum_{i=1}^n [-y_i \log p(\mathbf{x}_i; \mathbf{w}) - (1 - y_i) \log(1 - p(\mathbf{x}_i; \mathbf{w}))]$$

The Logit Transform

- $p(\mathbf{x}; \mathbf{w}) : \mathcal{X} \rightarrow [0, 1]$. How to parameterize $p(\mathbf{x}; \mathbf{w})$ using \mathbf{w} ?
- Odds Ratio = $\frac{\text{probability of event}}{\text{probability of non-event}}$

Assumption: the log of odds ratio is a linear function

Let's assume: $\log \frac{p(\mathbf{x}; \mathbf{w})}{1-p(\mathbf{x}; \mathbf{w})} = \langle \mathbf{x}, \mathbf{w} \rangle$.

- \mathbf{x} has been padded with 1; \mathbf{w} has been padded with b
- Equivalently, the **Sigmoid** transformation: $p(\mathbf{x}; \mathbf{w}) = \frac{1}{1+\exp(-\langle \mathbf{x}, \mathbf{w} \rangle)}$

Logistic Regression

- The minus log-likelihood that we talked before:

$$\min_{\mathbf{w}} \sum_{i=1}^n -y_i \log[p(\mathbf{x}_i; \mathbf{w})] - (1 - y_i) \log[1 - p(\mathbf{x}_i; \mathbf{w})]$$

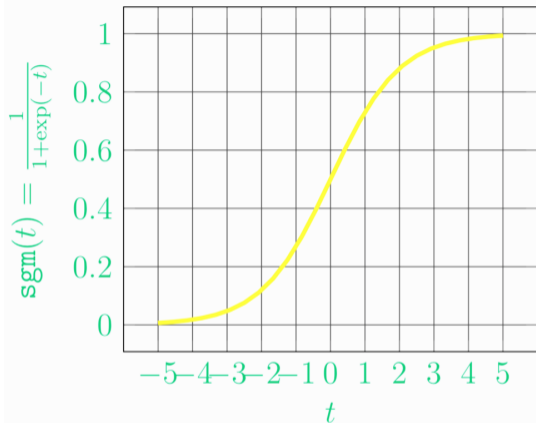
- Plug in the sigmoid parameterization $p(\mathbf{x}; \mathbf{w}) = \frac{1}{1 + \exp(-\langle \mathbf{x}, \mathbf{w} \rangle)}$

$$\min_{\mathbf{w}} \sum_{i=1}^n \boxed{\log[1 + \exp(-\langle \mathbf{x}_i, \mathbf{w} \rangle)] + (1 - y_i) \langle \mathbf{x}_i, \mathbf{w} \rangle}$$

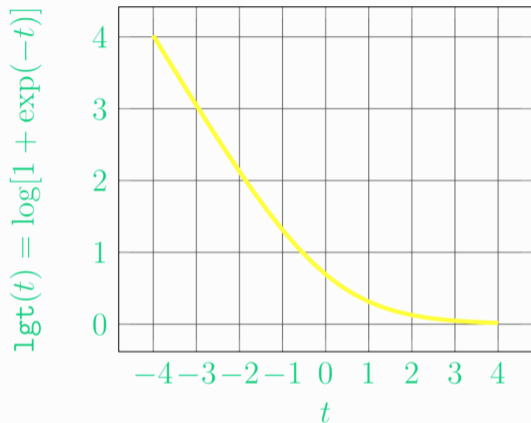
- Note the label encoding $y_i \in \{0, 1\}$ (let's call this y_i^{old}); if instead, $y_i^{new} \in \{\pm 1\}$, then

$$\min_{\mathbf{w}} \sum_{i=1}^n \underbrace{\boxed{\log[1 + \exp(-y_i^{new} \langle \mathbf{x}_i, \mathbf{w} \rangle)]}}_{\text{logistic loss}} \quad (\text{by transformation } y_i^{old} = \frac{y_i^{new} + 1}{2})$$

sigmoid function



logistic loss



D. R. Cox (1958). "The Regression Analysis of Binary Sequences". *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 20, no. 2, pp. 215–242.

Training: Solving Logistic Regression

$$\min_{\mathbf{w}} \sum_{i=1}^n \boxed{\log[1 + \exp(-y_i \langle \mathbf{x}_i, \mathbf{w} \rangle)]}$$

- Gradient descent algorithm:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \cdot \nabla_{\mathbf{w}} \text{Loss}(\mathbf{w})$$

- We will defer the introduction of gradient descent to **Lecture 8**.

Prediction

$$p(\mathbf{x}; \mathbf{w}) = \text{sigmoid}(\langle \mathbf{x}, \mathbf{w} \rangle) = \frac{1}{1 + \exp(-\langle \mathbf{x}, \mathbf{w} \rangle)}$$

- $\hat{y} = 1$ iff $p(\mathbf{x}; \mathbf{w}) = \Pr(Y = 1|X = \mathbf{x}) > \frac{1}{2}$ iff $\langle \mathbf{x}, \mathbf{w} \rangle > 0$
 - ▶ $\Pr(Y = 1|X = \mathbf{x}) > \frac{1}{2} \Leftrightarrow \Pr(Y = 1|X = \mathbf{x}) > \Pr(Y = -1|X = \mathbf{x})$
 - ▶ Therefore, $\hat{y} = \operatorname{argmax}_k \Pr(Y = k|X = \mathbf{x})$
- Decision boundary remains to be $H := \{\mathbf{x} : \langle \mathbf{x}, \mathbf{w} \rangle = 0\}$
- Can predict $\hat{y} = \operatorname{sign}(\langle \mathbf{x}, \mathbf{w} \rangle)$ as before, but now with confidence $p(\mathbf{x}; \mathbf{w})$
- An optional way to predict: check whether $p(\mathbf{x}; \mathbf{w}) > \frac{1}{2}$

More than a Classification Algorithm

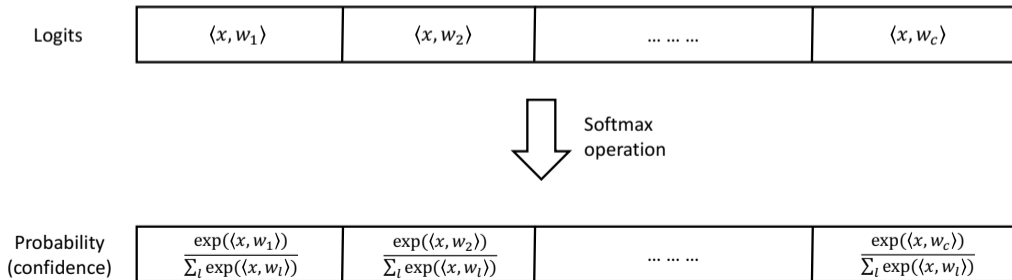
- Logistic regression estimates the posterior probability $\eta(\mathbf{x}) := \Pr(Y = 1|X = \mathbf{x})$ under the **linear log-odds ratio assumption**
 - ▶ Confidence is meaningless if the assumption does not hold true
- Classification itself only requires comparing $\eta(\mathbf{x})$ with $\frac{1}{2}$

Multi-Class Extension

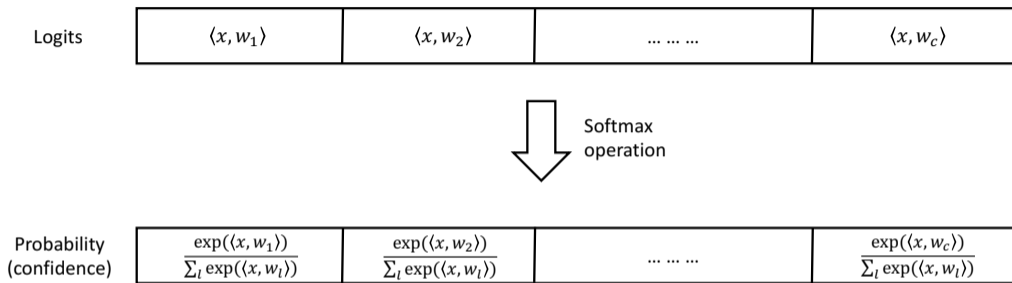
- Class $y \in \{1, \dots, c\}$; learn $\{\mathbf{w}_1, \dots, \mathbf{w}_c\}$ for each class
- **Sigmoid** \rightarrow **Softmax** function:

$$\Pr(Y = k | X = \mathbf{x}; \mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_c]) = \frac{\exp(\langle \mathbf{x}, \mathbf{w}_k \rangle)}{\sum_{l=1}^c \exp(\langle \mathbf{x}, \mathbf{w}_l \rangle)}$$

- ▶ map a real-valued vector to a probability vector
- ▶ non-negative and sum to 1



Multi-Class Extension



- Training: Again, maximum likelihood estimation

$$\min_{\mathbf{W}} \hat{\mathbb{E}} \left[-\log \frac{\exp(\langle \mathbf{X}, \mathbf{w}_Y \rangle)}{\sum_{l=1}^c \exp(\langle \mathbf{X}, \mathbf{w}_l \rangle)} \right]$$

- Prediction: $\hat{y} = \operatorname{argmax}_k \Pr(Y = k | X = \mathbf{x}; \mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_c])$

Questions

?

?

Answers

?