

Logistic regression

cs434

Logistic regression

- Recall the problem of regression
 - Learns a linear mapping from input vector \mathbf{x} to a continuous output $y \in (-\infty, \infty)$
- Logistic regression can be viewed as extending linear regression to handle binary output y by warping the output of a linear function to the range between 0 and 1
- For convenience, we will assume $y \in \{0,1\}$ for logistic regression

Logistic regression (cont.)

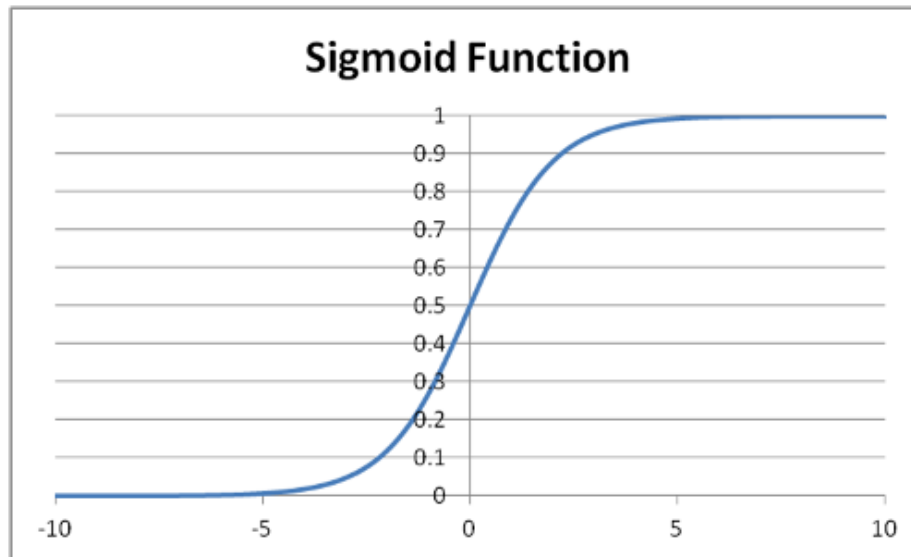
- Consider the linear regression function

$$\mathbf{w}^T \mathbf{x} = w_0 + w_1 x_1 + \cdots + w_m x_m$$

- We introduce a function σ :

$$\sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

Referred to as the sigmoid function or logistic function



Logistic Regression Makes Probabilistic Prediction

- We interpret the output of logistic regression as a probability:

$$P(y = 1|\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

i.e., probability of $y = 1$ given the input \mathbf{x} and the model's parameter is \mathbf{w}

$$P(y = 0|\mathbf{x}; \mathbf{w}) = 1 - \sigma(\mathbf{w}^T \mathbf{x}) = \frac{\exp(-\mathbf{w}^T \mathbf{x})}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

Logistic Regression learns a **linear** decision boundary

- We predict $y = 1$ if

$$P(y = 1|\mathbf{x}; \mathbf{w}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})} > P(y = 0|\mathbf{x}; \mathbf{w}) = \frac{\exp(-\mathbf{w}^T \mathbf{x})}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

$$\Rightarrow \frac{P(y = 1|\mathbf{x})}{P(y = 0|\mathbf{x})} = \exp(\mathbf{w}^T \mathbf{x}) > 1$$

$$\Rightarrow \log \frac{P(y = 1|\mathbf{x})}{P(y = 0|\mathbf{x})} = \mathbf{w}^T \mathbf{x} > 0$$

- Decision boundary:

$$\mathbf{w}^T \mathbf{x} = 0$$

Learning \mathbf{w} for logistic regression

- Given a set of training data points (\mathbf{x}_i, y_i) , $i = 1, \dots, n$, where $y_i \in \{0, 1\}$
- The goal is to find a weight vector \mathbf{w} s.t. the probability of the correct y_i for \mathbf{x}_i is high for $i = 1, \dots, n$

$$\text{maximize } P(y = y_i | \mathbf{x}_i; \mathbf{w})$$

- Or equivalently

$$\text{minimize } - \sum_{i=1}^n \log P(y = y_i | \mathbf{x}_i; \mathbf{w})$$

Negative log likelihood, or cross-entropy loss
--

Learning \mathbf{w} for logistic regression

$$\text{minimize } - \sum_{i=1}^n \log P(y = y_i | \mathbf{x}_i; \mathbf{w})$$

$$P(y = y_i | \mathbf{x}_i; \mathbf{w}) = \begin{cases} \sigma(\mathbf{w}^T \mathbf{x}_i) & \text{if } y_i = 1 \\ 1 - \sigma(\mathbf{w}^T \mathbf{x}_i) & \text{if } y_i = 0 \end{cases}$$
$$= \sigma(\mathbf{w}^T \mathbf{x}_i)^{y_i} (1 - \sigma(\mathbf{w}^T \mathbf{x}_i))^{(1-y_i)}$$

Our final objective:

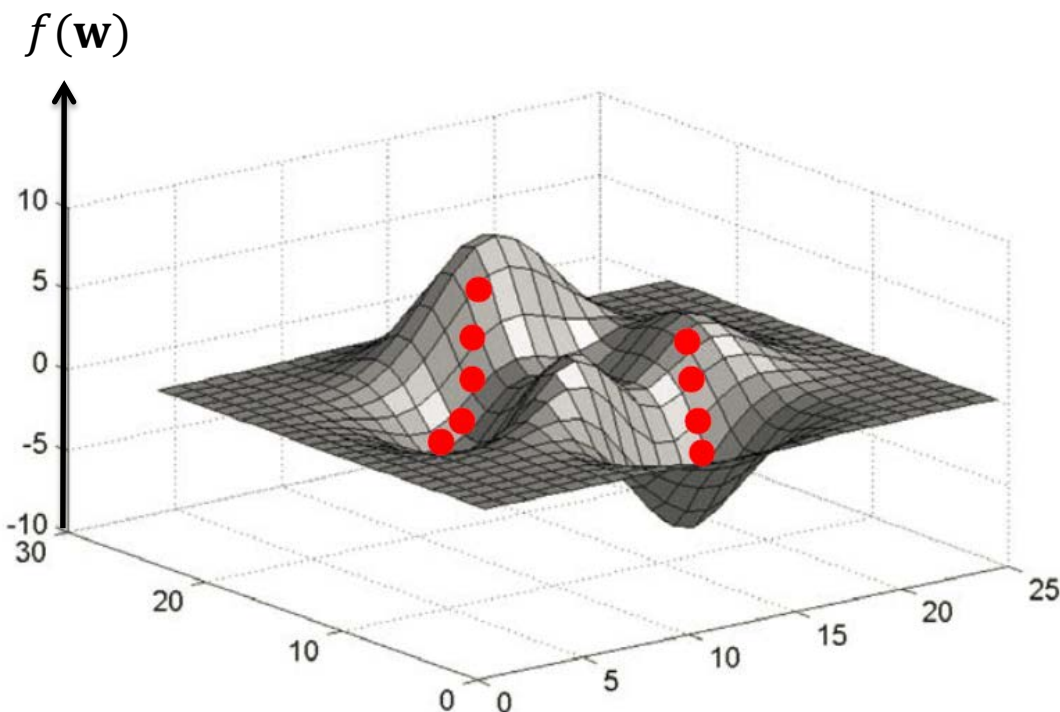
minimize

$$L(\mathbf{w}) = - \sum_{i=1}^n \left(y_i \log \sigma(\mathbf{w}^T \mathbf{x}_i) + (1 - y_i) \log (1 - \sigma(\mathbf{w}^T \mathbf{x}_i)) \right)$$

Minimizing $L(\mathbf{w})$

- Unfortunately this does not have a close form solution
 - You take the derivative, set it to zero, but no closed form solution
- Instead, we iteratively search for the optimal \mathbf{w}
- Start with a random \mathbf{w} , iteratively improve \mathbf{w} (similar to Perceptron) by taking the negative gradient
- This is referred to as gradient descent

Gradient descent to minimize $L(\mathbf{w})$



1. Start from some initial guess \mathbf{w}^0
2. Find the direction of steepest descent – opposite of the gradient direction $\nabla f(\mathbf{w})$
3. Take a step toward that direction
$$\mathbf{w}^{t+1} = \mathbf{w}^t - \nabla f(\mathbf{w}^t)$$
4. Repeat until no local improvement is possible
($|\nabla f(\mathbf{w}^t)| \leq \epsilon$)

Different starting points may lead to different local minimum if objective is not convex (like what's shown in the above figure)

Gradient of $L(\mathbf{w})$

Focus on the loss of a particular example (\mathbf{x}_i, y_i) :

$$l_i = -y_i \log \sigma(\mathbf{w}^T \mathbf{x}_i) - (1 - y_i) \log (1 - \sigma(\mathbf{w}^T \mathbf{x}_i))$$

Useful fact: $\sigma'(t) = \sigma(1 - \sigma)$

$$\nabla l_i = \frac{-y_i}{\sigma(\mathbf{w}^T \mathbf{x}_i)} \nabla \sigma(\mathbf{w}^T \mathbf{x}_i) + \frac{1 - y_i}{1 - \sigma(\mathbf{w}^T \mathbf{x}_i)} \nabla \sigma(\mathbf{w}^T \mathbf{x}_i)$$

$$\nabla \sigma(\mathbf{w}^T \mathbf{x}_i) = \sigma(\mathbf{w}^T \mathbf{x}_i) (1 - \sigma(\mathbf{w}^T \mathbf{x}_i)) \mathbf{x}_i$$

$$\nabla l_i = (\sigma(\mathbf{w}^T \mathbf{x}_i) - y_i) \mathbf{x}_i$$

Online gradient descent for Logistic Regression

Recall: y takes 0/1 here, not 1/-1

Given: training examples $(\mathbf{x}_i, y_i), i = 1, \dots, n$

Let $\mathbf{w} \leftarrow (0, \dots, 0)$ // initialization

Repeat until convergence //

for $i = 1, \dots, n$

$$\hat{y}_i \leftarrow \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}}$$

$$\mathbf{w} \leftarrow \mathbf{w} - \eta(\hat{y}_i - y_i)\mathbf{x}_i$$

Note the striking similarity between LR and perceptron

Both learn a linear decision boundary

The iterative algorithm takes very similar form

Batch Learning for Logistic Regression

Note: y takes 0/1 here, not 1/-1

Given: training examples $(\mathbf{x}_i, y_i), i = 1, \dots, n$

Let $\mathbf{w} \leftarrow (0, \dots, 0)$ //initialization

Repeat

$$\nabla = (0, \dots, 0)$$

for $i = 1, \dots, n$

$$\hat{y}_i \leftarrow \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}}$$

$$\nabla \leftarrow \nabla + (\hat{y}_i - y_i) \mathbf{x}_i$$

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla$$

Until $|\nabla| \leq \epsilon$

Logistic regression: Summary

- LR uses σ to warp the output of a linear function to $(0,1)$, which is interpreted as $P(y = 1|\mathbf{x}; \mathbf{w})$
- Learn a vector \mathbf{w} s.t. for input \mathbf{x}_i the desired correct output y_i has high probability, and other incorrect outputs have low probabilities
 - Maximize the log likelihood or
 - Minimize the negative log likelihood
- Learn \mathbf{w} iteratively using gradient descent
- Strong similarity with Perceptron
- Logistic regression learns a linear decision boundaries
 - By introducing nonlinear features (i.e., $x_1^2, x_2^2, x_1x_2, \dots$), can be extended to nonlinear boundary.