

# CLASSIFICATION

## CLS 1 & 2. LR Models

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

1

1

## 4 Lectures on Classification

CLS 1 & 2. Logistic Regression (LR) Models

CLS 3. Stochastic Gradient Descent & Evaluation Metrics

CLS 4. Neural Classifiers for Extremely Large Classification

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

2

2

## Outline on LR Models

- Introduction
- Decision boundaries
- Binary LR
- Optimization algorithms
- Convexity
- Regularization
- Softmax LR

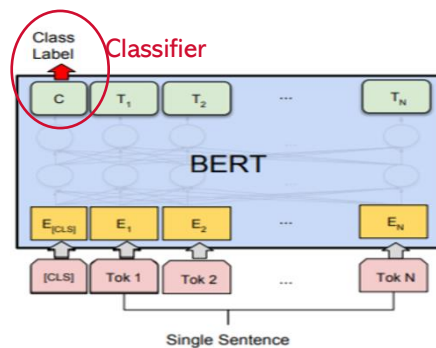
02/06/2024

@Yiming Yang, S24 Lecture on LR Models

3

3

## BERT Fine Tuning for Classification



(b) Single Sentence Classification Tasks:  
SST-2, CoLA

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

4

4

## Application Examples

- **Email spam detection (binary classification)**
  - Given message  $x \in \mathbb{R}^d$ , predict  $y \in \{yes, no\}$ .
- **Hand-written digit recognition (multi-class classification)**
  - Given image  $x \in \mathbb{R}^d$ , predict  $y \in \{0, 1, \dots, 9\}$ ;
  - Choosing 1 out of  $M > 2$  category labels.
- **Wikipedia page subject topics (multi-label classification)**
  - Given input text, predict the relevant labels;
  - Choosing 1 or more out of  $M > 2$  category labels.

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

5

5

## Mathematical Definition

- Find the mapping  $f: X \rightarrow Y$  for  $X \in \mathbb{R}^d$  and  $Y \in \{0, 1\}^M$ .
  - Practically, predict vector  $f(x) \in \mathbb{R}^M$  and apply a threshold to the elements of  $f(x)$  for yes/no decisions
    - Option 1. Assigning *yes* to the  $k$  top-ranking label and *no* to the rest where the  $k$  is a prespecified hyper-parameter;
    - Option 2. Assigning *yes* to label  $j$  is  $f_j(x) \geq 0.5$  for  $j = 1, \dots, M$ ;
    - Option 3. ...
- (see Y Yang, SIGIR 2001)

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

6

6

## Terminology

$X$	$Y$
Input Variables	Output Variables
Independent Variables	Dependent Variables
Predictors	Responses
Features	Categories or labels
Factors	Outcomes

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

7

7

## Scoring Functions

- **Linear Function** (e.g., linear regression or Naïve Bayes models)

$$f_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1 x_1 + \cdots + w_d x_d = \mathbf{w}^T \mathbf{x}$$

where  $\mathbf{x} = (1, x_1, \dots, x_d)$  is a data point, and

$\mathbf{w} = (w_0, w_1, \dots, w_d)$  are the model parameters.

- **Sigmoid Logistic Regression** (binary LR)

$$f_{\mathbf{w}}(\mathbf{x}) \equiv \widehat{P}_{\mathbf{w}}(Y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

8

8

## Scoring Functions (cont'd)

- **SoftMax LR:** for  $j \in \{1, 2, \dots, M\}$  and  $W = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_M)$

$$f_j(\mathbf{x}; W) \equiv \hat{P}(Y = j | \mathbf{x}; W) = \frac{\exp(\mathbf{w}_j^T \mathbf{x})}{\sum_{m=1}^M \exp(\mathbf{w}_m^T \mathbf{x})}$$

- **k-Nearest Neighbors (kNN)** (Non-parametric)

$$f_j(\mathbf{x} | D) = \frac{\sum_{x_i \in kNN(\mathbf{x})} \delta(y_i, j)}{k}, \quad \delta(y_i, j) = \begin{cases} 1 & \text{if } y_i = j \\ 0 & \text{otherwise} \end{cases}$$

$D = \{(x_i, y_i)\}_{i=1, \dots, N}$  is a labeled training set;

$kNN(\mathbf{x})$  is the set of k-nearest-neighbors of  $\mathbf{x}$  in  $D$ .

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

9

9

## Example of a linear classifier

Elements of Statistical Learning ©Hastie, Tibshirani & Friedman 2001 Chapter 2

- A **linear decision boundary** is defined as a set of data points

$$h = \{\mathbf{x}: \mathbf{w}^T \mathbf{x} = b\}$$

which is

- a line in 2D
- a plane in 3D
- a hyperplane in  $\mathbb{R}^d$

- If the decision boundary by a classifier is linear, we call it a **linear classifier**.

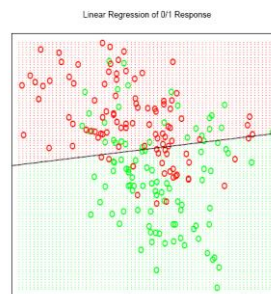


Figure 2.1: A classification example in two dimensions. The classes are coded as a binary variable—GREEN = 0, RED = 1—and then fit by linear regression. The line is the decision boundary defined by  $\mathbf{x}^T \hat{\beta} = 0.5$ .

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

10

10

## LDA (linear) vs. QDA (non-linear)

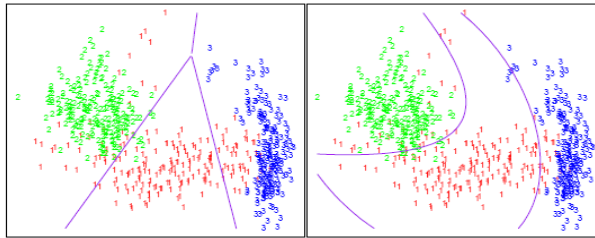


Figure 4.1: The left plot shows some data from three classes, with linear decision boundaries found by linear discriminant analysis. The right plot shows quadratic decision boundaries. These were obtained by finding linear boundaries in the five-dimensional space  $X_1, X_2, X_{12}, X_1^2, X_2^2$ . Linear inequalities in this space are quadratic inequalities in the original space.

(Hastie et al., ESL)

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

11

11

## Decision Boundaries of kNN (non-linear)

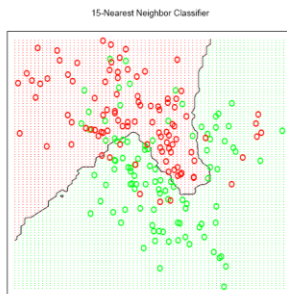


Figure 2.2: The same classification example in two dimensions as in Figure 2.1. The classes are coded as a binary variable (GREEN = 0, RED = 1) and then fit by 15-nearest-neighbor averaging as in (2.8). The predicted class is hence chosen by majority vote amongst the 15-nearest neighbors.

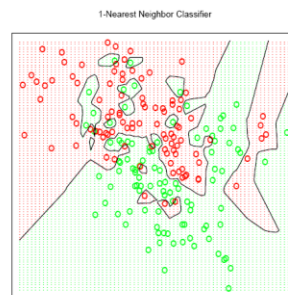


Figure 2.3: The same classification example in two dimensions as in Figure 2.1. The classes are coded as a binary variable (GREEN = 0, RED = 1), and then predicted by 1-nearest-neighbor classification.

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

12

12

## How to tell if a classifier is linear or not?

- We cannot tell by just looking at  $f_w$  (as being linear or non-linear).
- Instead, we must check if the decision boundary can be written as

$$h = \{x: \mathbf{w}^T x = b\}$$

- Let's take a look at a binary LR and Softmax classifiers.

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

13

13

## Is binary LR a linear classifier?

- Scoring function given  $x$  is sigmoid (**non-linear**)

$$\sigma_w(x) = (1 + e^{-w^T x})^{-1} \quad (1)$$

- A popular threshold for a binary decision is set as

$$\sigma_w(x) = 0.5 \quad (2)$$

- Denoting the decision boundary as  $h$  we have

$$h = \{x: (1 + e^{-w^T x})^{-1} = 0.5\} \quad (3)$$

$$\Rightarrow 1 + e^{-w^T x} = 2 \Rightarrow e^{-w^T x} = 1 \Rightarrow w^T x = 0$$

$$\Rightarrow h = \{x: w^T x = 0\} \quad (4)$$

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

14

14

## Is softmax LR a linear classifier?

- Scoring function for  $k = 1, 2, \dots, K$

$$\Pr(y = k|x) = \frac{\exp(\mathbf{w}_k^T x)}{\sum_{k'=1}^K \exp(\mathbf{w}_{k'}^T x)} \equiv \hat{p}_k(x) \quad (5)$$

- Decision boundary between labels  $j$  and  $k$

$$h_{jk} = \{x: \hat{p}_j(x) = \hat{p}_k(x)\} \quad (6)$$

$$\Rightarrow \frac{\exp(\mathbf{w}_j^T x)}{\sum_{k'=1}^K \exp(\mathbf{w}_{k'}^T x)} = \frac{\exp(\mathbf{w}_k^T x)}{\sum_{k'=1}^K \exp(\mathbf{w}_{k'}^T x)} \Rightarrow \mathbf{w}_j^T x = \mathbf{w}_k^T x$$

$$\Rightarrow (\mathbf{w}_j^T - \mathbf{w}_k^T)x = 0$$

$$\Rightarrow h_{jk} = \{x: \underbrace{(\mathbf{w}_j - \mathbf{w}_k)^T}_{\mathbf{w}} x = 0\} \quad (7)$$

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

15

15

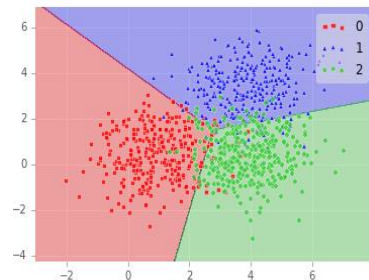
## Is softmax LR a linear classifier?

- Answer

- Locally (pairwise) linear but globally nonlinear

- Thresholding strategy

$$\hat{y}(\mathbf{x}) = \operatorname{argmax}_k \{\hat{p}_k(\mathbf{x})\}$$



02/06/2024

@Yiming Yang, S24 Lecture on LR Models

16

16



## Outline

- ✓ Introduction
- ✓ Decision boundaries
- Binary LR
- Optimization algorithms
- Convexity
- Regularization
- Softmax LR

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

17

17

## LR for Binary Classification

- Label probabilities estimated using a sigmoid function

$$P_{\mathbf{w}}(y = 1|\mathbf{x}) = \sigma_{\mathbf{w},\mathbf{x}} = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

$$P_{\mathbf{w}}(y = 0|\mathbf{x}) = 1 - \sigma_{\mathbf{w},\mathbf{x}} = \frac{\exp(-\mathbf{w}^T \mathbf{x})}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

- Compact formula

$$P_{\mathbf{w}}(y|\mathbf{x}) = (\sigma_z)^y (1 - \sigma_z)^{(1-y)} \text{ with } z = \mathbf{w}^T \mathbf{x}$$

$$\log P_{\mathbf{w}}(y|\mathbf{x}) = y \log \sigma_z + (1 - y) \log (1 - \sigma_z)$$

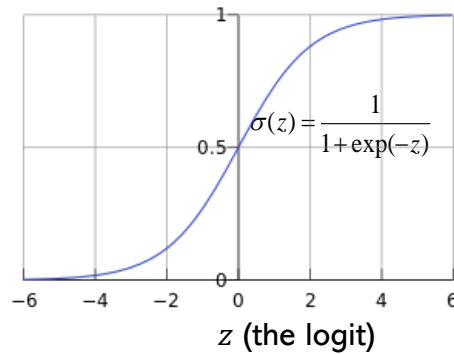
02/06/2024

@Yiming Yang, S24 Lecture on LR Models

18

18

## Sigmoid Function



$$z \in (-\infty, \infty), \quad \sigma(z) \in (0, 1), \quad \sigma(0) = 0.5, \quad \sigma(z) = (1 - \sigma(-z))$$

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

19

19

## Logit = logarithm of the odds

$$p = \frac{1}{1 + \exp(-z)} \quad \text{Start with the sigmoid}$$

$$p(1 + \exp(-z)) = 1 \quad \text{Multiply the denominator on both sides}$$

$$\exp(-z) = \frac{1-p}{p} \quad \text{Arrange } p \text{ to the RHS}$$

$$\exp(z) = \frac{p}{1-p} \quad \text{Flip over}$$

$$\text{logit } z = \log \frac{p}{1-p} \quad \text{odds of } p \quad \text{Take the log on both sides}$$

02/06/2024

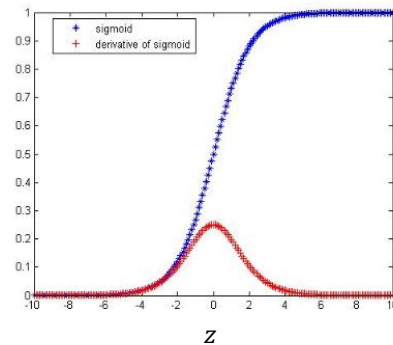
@Yiming Yang, S24 Lecture on LR Models

20

20

## Derivative of Sigmoid

$$\begin{aligned}\frac{d\sigma(z)}{dz} &= \frac{d}{dz} \left( \frac{1}{1 + \exp(-z)} \right) \\ &= (-1)(-1) \frac{\exp(-z)}{(1 + \exp(-z))^2} \\ &= \frac{1}{(1 + \exp(-z))} \frac{\exp(-z)}{(1 + \exp(-z))} \\ &= \sigma(z)(1 - \sigma(z))\end{aligned}$$



02/06/2024

@Yiming Yang, S24 Lecture on LR Models

21

21

## A good online lecture on LR by Andrew Ng

- [Andrew Ng on LR decision boundary](#)

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

22

22

## Outline

- ✓ Introduction
- ✓ Decision boundaries
- ✓ Binary Logistic Regression (LR)
  - Optimization algorithms
  - Convexity
  - Regularization
  - Softmax LR

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

23

23

## Training a Binary Classifier

- Labeled Training Data

$$D = \{(x_i, y_i)\}_{i=1}^n \text{ with } x_i \in \mathbb{R}^d \text{ and } y_i \in \{-1, 1\}$$

- Model Training

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \{ \operatorname{Loss}(D; \mathbf{w}) + C \|\mathbf{w}\|^2 \}$$

- $\operatorname{Loss}(D; \mathbf{w})$  is the **training-set loss**, measuring how well the model fits the labeled data;
- $\|\mathbf{w}\|^2$  is the **regularization term**, controlling the model complexity to avoid overfitting.

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

24

24

## Parameter Optimization

$$\hat{\mathbf{w}} = \operatorname{argmax}_{\mathbf{w}} \{ \log P(D|\mathbf{w}) \} \quad (\text{the log-likelihood})$$

$$= \operatorname{argmin}_{\mathbf{w}} \{ -\log P(D|\mathbf{w}) \} \quad (\text{the negative log-likelihood})$$

$$= \operatorname{argmin}_{\mathbf{w}} \underbrace{- \sum_{i=1}^n \{ y_i \ln \sigma_i + (1 - y_i) \ln (1 - \sigma_i) \}}_{l(\mathbf{w})}$$

the cross-entropy loss

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

25

25

## Popular Algorithms

### ■ Gradient Descent

- Use the first-order derivative of  $l(\boldsymbol{\beta})$
- Need to pre-specify the "learning rate" (step size)
- Fast to compute in each step but may take many steps

### ■ Newton-Raphson

- Use the first-order and second-order derivatives of  $l(\boldsymbol{\beta})$
- Automatically find the optimal step size for each iteration
- Converge faster but may be too costly in each step

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

26

26

## Gradient on a single training pair

$$l_i(\mathbf{w}) = y_i \ln \frac{\sigma(z_i)}{\sigma_i} + (1 - y_i) \ln(1 - \sigma(z_i)) \quad z_i = \mathbf{w}^T \mathbf{x}_i = w_0 + \sum_{j=1}^m w_j x_{ij}$$

$$\begin{aligned} \frac{\partial}{\partial w_j} l_i(\mathbf{w}) &= \frac{dl_i}{d\sigma_i} \frac{d\sigma_i}{dz_i} \frac{\partial z_i}{\partial w_j} \\ &= \left( y_i \frac{1}{\sigma_i} - (1 - y_i) \frac{1}{1 - \sigma_i} \right) \sigma_i(1 - \sigma_i) x_{ij} = (y_i - \sigma_i) x_{ij} \end{aligned}$$

$$\nabla l_i(\mathbf{w}) = \left( \frac{\partial}{\partial w_0} l_i(\mathbf{w}), \frac{\partial}{\partial w_1} l_i(\mathbf{w}), \dots, \frac{\partial}{\partial w_m} l_i(\mathbf{w}) \right)^T$$

02/06/2024

@Yiming Yang, S24 Lecture on  
LR Models

27

27

## Gradient ascent on a training set

The single-instance version:  $D = \{(\mathbf{x}^{(i)}, y^{(i)})\}$

Loop until convergence {

for  $i = 1$  to  $|D|$  {

$\mathbf{w} := \mathbf{w} + \eta \nabla l_i(\mathbf{w})$

( $\eta > 0$  is prespecified or adapted  
via backtracking line search)

}

The batch version:

Loop until convergence {

$\mathbf{w} := \mathbf{w} + \eta \sum_{i=1}^{|D|} \nabla l_i(\mathbf{w})$

}

Guaranteed:  $l(\mathbf{w}^{(0)}) \geq l(\mathbf{w}^{(1)}) \geq l(\mathbf{w}^{(2)}) \dots$

02/06/2024

@Yiming Yang, S24 Lecture on LR Models 28

28

## Newton-Raphson Method

(in the case of one-dimensional  $w$ )

- Given current  $w$ , we want move it with the optimal step size ( $\varepsilon$ ) in the right direction.
- Taylor series:

$$l(w + \varepsilon) = l(w) + \frac{l'(w)}{1!} \varepsilon + \frac{l''(w)}{2!} \varepsilon^2 + \dots$$

- At the mode (with respect to  $\varepsilon$ )

$$0 = \frac{d}{d\varepsilon} l(w + \varepsilon) \approx l'(w) + l''(w)\varepsilon \Rightarrow \varepsilon = -\frac{l'(w)}{l''(w)}$$

- Update rule:
- $$w := w - \frac{l'(w)}{l''(w)}$$

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

29

29

## Newton-Raphson Method

(in the case of multi-dimensional  $\mathbf{w}$ )

- Taylor series:

$$l(\mathbf{w} + \boldsymbol{\varepsilon}) = l(\mathbf{w}) + \nabla l(\mathbf{w})\boldsymbol{\varepsilon} + \boldsymbol{\varepsilon}^T \frac{\nabla \nabla l(\mathbf{w})}{2!} \boldsymbol{\varepsilon} + \dots$$

- Update rule:

$$\mathbf{w} := \mathbf{w} - \underbrace{(\nabla \nabla l(\mathbf{w}))^{-1}}_{\mathbf{H}(\mathbf{w})} \underbrace{\nabla l(\mathbf{w})}_{\text{the gradient}}$$

$$\nabla l(\mathbf{w}) = \left( \frac{\partial}{\partial w_0} l(\mathbf{w}), \frac{\partial}{\partial w_1} l(\mathbf{w}), \dots, \frac{\partial}{\partial w_m} l(\mathbf{w}) \right)^T$$

$$\nabla \nabla l(\mathbf{w}) \equiv \mathbf{H}(\mathbf{w}) = (H_{jj'}) , \quad H_{jj'} = \frac{\partial^2}{\partial w_j \partial w_{j'}} l(\mathbf{w})$$

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

30

30