
CLASSIFICATION

Class 1. Logistic Regression Models

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

1

1

Outline

- Introduction
- Binary LR
- Softmax LR

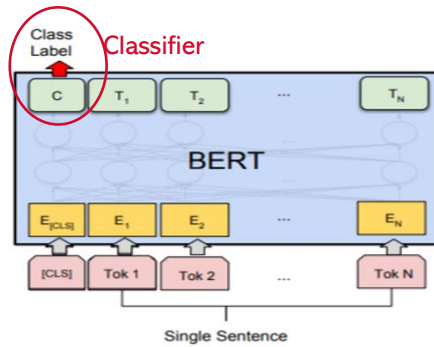
02/06/2024

@Yiming Yang, S24 Lecture on LR Models

2

2

BERT Fine Tuning for Classification



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

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

3

3

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 one and only one output label.
- **News article topic tagging (multi-label classification)**
 - Given article $x \in \mathbb{R}^d$, predict $y \subseteq \{0, 1\}^M$ where M is the number of category labels and y specifies the relevant subset given the input (e.g., China, Taiwan, economy, politics).

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

4

4

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

5

5

Classification via Label Scoring

- We usually learn scoring function $f_{\theta}(\mathbf{x}) \in \mathbb{R}^M$ via model training, where θ is the set of learnable parameters, and each element $f_{\theta}^{(j)}(\mathbf{x})$ is the predicted score with respect to label $j \in \{1, \dots, M\}$.
- We can obtain a *yes/no* decision on each category by assigning *yes* if $f_{\theta}^{(j)}(\mathbf{x}) \geq 0.5$ and *no* otherwise.
- Or we can use sort the labels based on their scores and assign *yes* to the top one and *no* to the rest as

$$\hat{y}(\mathbf{x}) = \operatorname{argmax}_j \{f_{\theta}^{(j)}(\mathbf{x})\}$$

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

6

6

Scoring Functions

- **Linear Function**

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

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

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

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

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

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

7

7

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 .

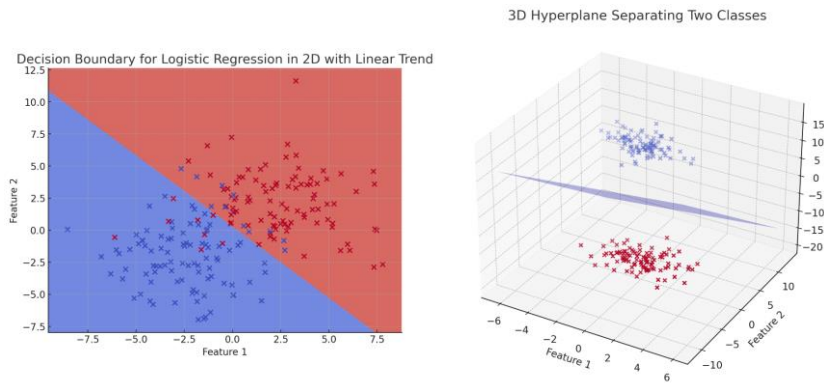
02/06/2024

@Yiming Yang, S24 Lecture on LR Models

8

8

Decision Boundaries of Binary LR (Generated by GPT4)



Binary LR has a linear decision boundary (e.g., in 2D and 3D).

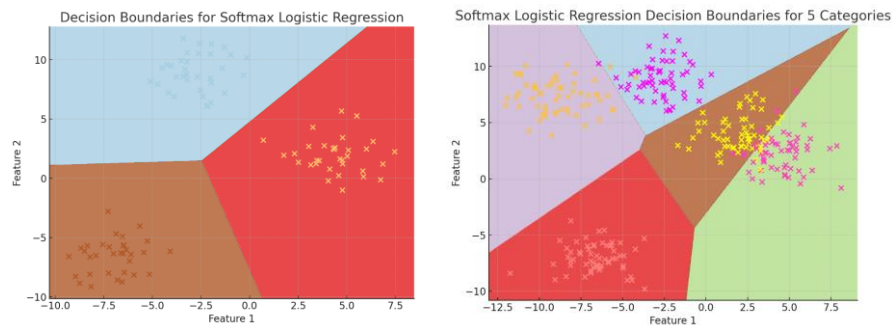
02/06/2024

@Yiming Yang, S24 Lecture on LR Models

9

9

Decision Boundaries of Softmax LR (Generated by GPT4)



The decision boundaries are piecewise linear but globally non-linear.

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

10

10

Linear Decision Boundary Formally

- In 2D, a **line** can be written as $ax + by + c = 0$.
- In 3D, a **plane** can be written as $ax + by + cz + d = 0$.
- In d -dimensional vector space, a **hyperplane (h)** can be defined as

$$h = \{x: w^T x = 0\}$$

where $x = (1, x_1, x_2, \dots, x_d)$ and $w = (w_0, w_1, \dots, w_d)$.

- If the decision boundary of a classifier can be written as $w^T x = 0$, we call it a **linear classifier**.
- We cannot tell by looking at the scoring function $f_\theta(x)$ along; instead, we must take the **decision rule (thresholding strategy)** into account.

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

11

11

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 is set as

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

- The decision boundary is

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

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

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

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

12

12

How about softmax LR?

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

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

- Decision boundary between labels j and k

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

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

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

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

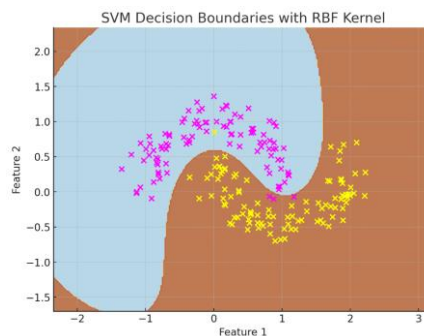
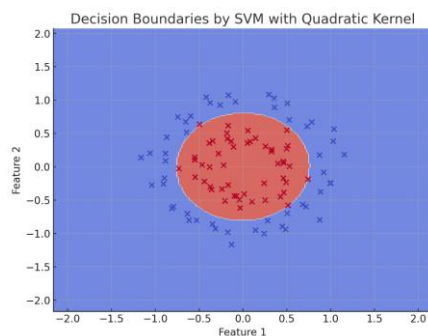
02/06/2024

@Yiming Yang, S24 Lecture on LR Models

13

13

Other Examples: Support Vector Machines



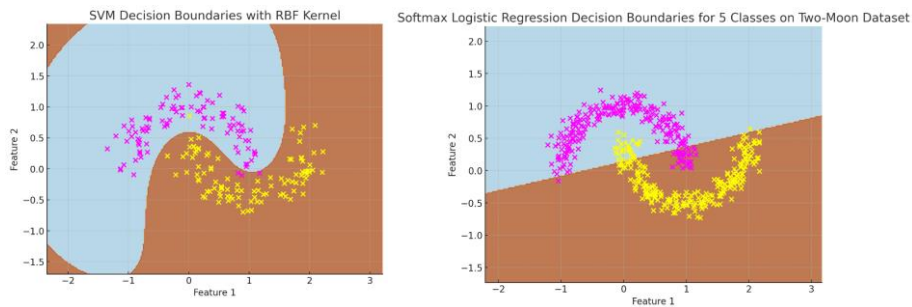
02/06/2024

@Yiming Yang, S24 Lecture on LR Models

14

14

SVM vs. Softmax (for 5 classes)



02/06/2024

@Yiming Yang, S24 Lecture on LR Models

15

15

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

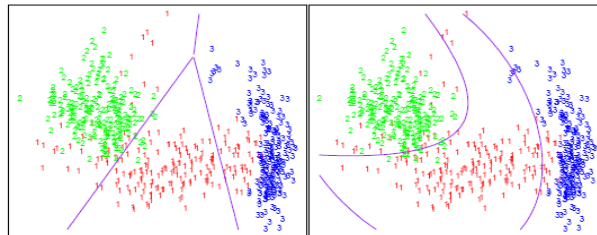


Figure 4.1: The left plot shows some data from three classes, with linear decision boundaries found by linear discriminant analysis. The right plot shows quadratic decision boundaries. These were obtained by finding linear boundaries in the five-dimensional space $X_1, X_2, X_{12}, X_1^2, X_2^2$. Linear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

16

16

Decision Boundaries of kNN (non-linear)

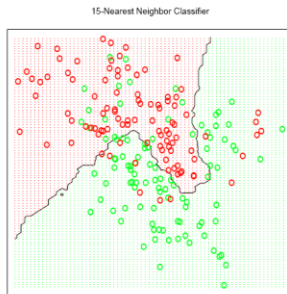


Figure 2.2: The same classification example in two dimensions as in Figure 2.1. The classes are coded as a binary variable (GREEN = 0, RED = 1) and then fit by 15-nearest-neighbor averaging as in (2.8). The predicted class is hence chosen by majority vote amongst the 15-nearest neighbors.

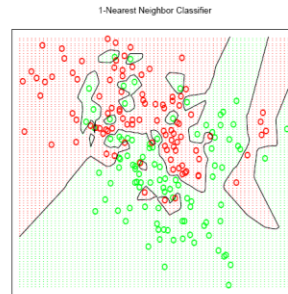


Figure 2.3: The same classification example in two dimensions as in Figure 2.1. The classes are coded as a binary variable (GREEN = 0, RED = 1), and then predicted by 1-nearest-neighbor classification.

Outline

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

LR for Binary Classification

- Using sigmoid function to estimate label probabilities

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

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

- Merged formula

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

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

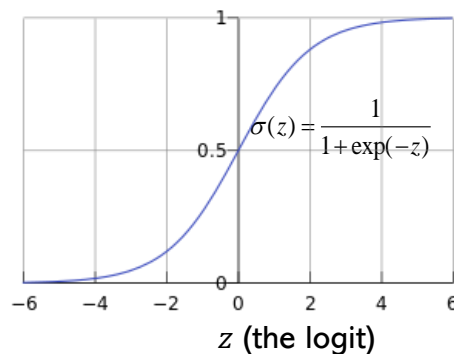
02/06/2024

@Yiming Yang, S24 Lecture on LR Models

19

19

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

20

20

Logit = logarithm of the odds

$$p = \frac{1}{1 + \exp(-z)}$$

Start with the sigmoid

$$p(1 + \exp(-z)) = 1$$

Multiply the denominator on both sides

$$\exp(-z) = \frac{1-p}{p}$$

Arrange p to the RHS

$$\exp(z) = \frac{p}{1-p}$$

Flip over

$$\text{logit } z = \log \frac{p}{1-p} \text{ odds of } p$$

Take the log on both sides

02/06/2024

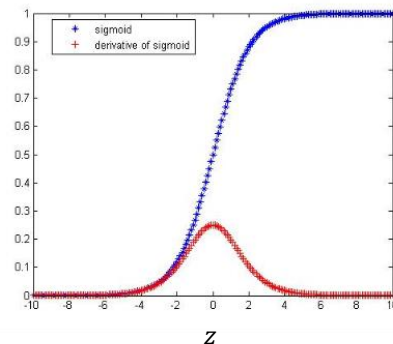
@Yiming Yang, S24 Lecture on LR Models

21

21

Derivative of Sigmoid

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



02/06/2024

@Yiming Yang, S24 Lecture on LR Models

22

22

Training Binary LR

- Labeled Training Data

$$D = \{(x_i, y_i)\}_{i=1}^n \text{ with } x_i \in R^d \text{ and } y_i \in \{0,1\}$$

- 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 data;
- $\|\mathbf{w}\|^2$ is the **regularization term (on model parameters only)**, controlling the model complexity to avoid overfitting.

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

23

23

Parameter Optimization

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

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

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

cross-entropy loss

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

24

24

Popular Algorithms

■ Gradient Descent

- Use the first-order derivative of $l(\beta)$
- 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 $l(\beta)$
- 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

25

25

Gradient on a single training pair

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

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

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

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

26

26

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

02/06/2024

@Yiming Yang, S24 Lecture on LR Models 27

27

Newton-Raphson Method

(in the case of one-dimensional w)

- Given current w , we want move it with the optimal step size (ε) 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 ε)

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

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

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

28

28

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

29

29

Newton-Raphson Method (cont'd)

First order derivative (as shown in slide #11):

$$\frac{\partial}{\partial w_j} l(\mathbf{w}) = \sum_{i=1}^n (y_i - \sigma(\mathbf{w}^T \mathbf{x}_i)) x_{ij}$$

Second order derivative:

$$\begin{aligned} \frac{\partial^2}{\partial w_j \partial w_{j'}} l(\mathbf{w}) &= \frac{\partial}{\partial w_{j'}} \left(\frac{\partial}{\partial w_j} l(\mathbf{w}) \right) = \sum_{i=1}^n \frac{\partial}{\partial w_{j'}} (y_i - \sigma(\mathbf{w}^T \mathbf{x}_i)) x_{ij} \\ &= - \sum_{i=1}^n \left(\frac{d\sigma}{dz_i} \frac{\partial z_i}{\partial w_{j'}} \right) x_{ij} = - \sum_{i=1}^n \sigma_i (1 - \sigma_i) x_{ij} x_{ij'} \end{aligned}$$

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

30

30

Newton-Raphson Method (cont'd)

The gradient (compact notion) :

$$\nabla l(\mathbf{w}) = \sum_{i=1}^n (y_i - \sigma_i) \mathbf{x}_i = \mathbf{X}^T (\mathbf{y} - \boldsymbol{\sigma})$$

$\nabla l(\mathbf{w})$ is the weighted sum of the training documents;

\mathbf{X}^T is $(m+1) \times n$, whose columns (\mathbf{x}_i 's) are the training documents;

$\mathbf{y} = (y_1, y_2, \dots, y_n)^T$ is the vector of true labels of n training doc's;

$\boldsymbol{\sigma} = (\sigma_1, \sigma_2, \dots, \sigma_n)^T$ is the vector of predicted probabilities $\sigma_i = \sigma(\mathbf{w}^T \mathbf{x}_i)$.

The Hessian (compact notion)

$$\mathbf{H}(\mathbf{w}) = - \sum_{i=1}^n \sigma_i (1 - \sigma_i) \mathbf{x}_i (\mathbf{x}_i)^T = -\mathbf{X}^T \Lambda \mathbf{X}$$

$$\Lambda = \text{diag}(\sigma_1(1 - \sigma_1), \sigma_2(1 - \sigma_2), \dots, \sigma_n(1 - \sigma_n))$$

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

31

31

Newton-Raphson Method (cont'd)

- Update rule in 1-dimensional LR

$$w := w^{\text{old}} - \frac{l'(w^{\text{old}})}{l''(w^{\text{old}})}$$

- Update rule in high-dimensional LR

$$\mathbf{w} := \mathbf{w}^{\text{old}} - H(\mathbf{w}^{\text{old}})^{-1} \nabla l(\mathbf{w}^{\text{old}})$$

$$:= \mathbf{w}^{\text{old}} + (\mathbf{X}^T \Lambda \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \boldsymbol{\sigma})$$

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

32

32

Outline

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

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

33

33

Globally optimal solution guaranteed?

- Check the convexity of the objective function
 - If it is **convex**, then there is a single global minimum.
 - If it is **concave**, then there is a single global maximum.
 - If it is **neither**, then the global optimal is not guaranteed.

02/06/2024

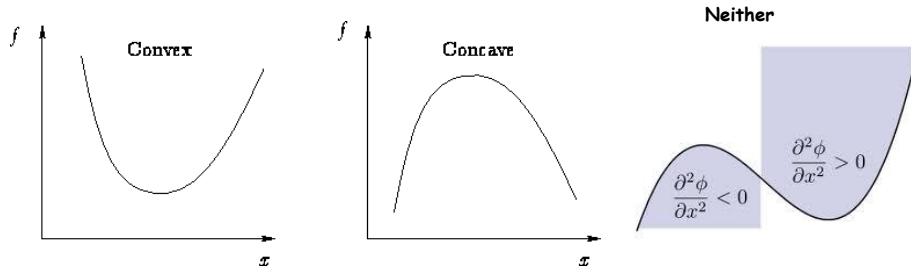
@Yiming Yang, S24 Lecture on LR Models

34

34

Convexity

Examples of 1-dimensional functions



02/06/2024

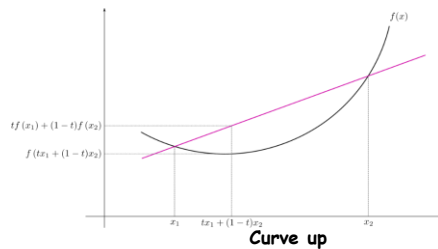
@Yiming Yang, S24 Lecture on LR Models

35

35

Convex Function

https://en.wikipedia.org/wiki/Convex_function



- f is called **convex** if:

$$\forall x_1, x_2 \in X, \forall t \in [0, 1]: \quad f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2).$$

- f is called **strictly convex** if:

$$\forall x_1 \neq x_2 \in X, \forall t \in (0, 1): \quad f(tx_1 + (1-t)x_2) < tf(x_1) + (1-t)f(x_2).$$

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

36

36

Convexity of One-dimensional Function

- If $f(x)$ has a second derivative in $[a, b]$, then a necessary and sufficient condition for it to be **convex** on that interval is that the second derivative $f''(x) \geq 0$ for all $x \in [a, b]$.
- We call it **strictly convex** if $f''(x) > 0$ for all $x \in [a, b]$.
- We call it **strongly convex** if $f''(x) > m$ with a positive m for all $x \in [a, b]$.

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

37

37

Examples (One-dimensional Case)

(Garret Thomas, <https://gwrthomas.github.io/docs/math4ml.pdf>)

- Consider real $x \in \{-\infty, \infty\}$
 - 1) $f(x) = 2x + 3$: $f' = 2$, $f'' = 0$ ← Convex but not strictly nor strongly
 - 2) $f(x) = x^2 + 2x + 1$: $f' = 2x$, $f'' = 2$ ← Strongly (and strictly)
 - 3) $f(x) = e^x$: $f' = f'' = e^x$ ← Convex? Strictly? Strongly?
 - 4) $f(x) = x^4$ ← ???

$$f' = 4x^3, f'' = 12x^2, f'''(0) = 0$$

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

38

38

Convexity of High-dimensional Function

We can its Hessian matrix $H \in \mathbb{R}^{d \times d}$.

- If $H \succcurlyeq 0$ (positive semi-definite), meaning that $\forall u, u^T H u \geq 0$, then f is **convex**.
- If $H \preccurlyeq 0$ (negative semidefinite), meaning that $\forall u, u^T H u \leq 0$, then f is **concave**.
- If none of the above is true, we can only reach a local optimal.

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

39

39

Convexity of High-dimensional Function

- **Proposition.** Suppose $f(x)$ is twice differentiable in its domain, and denote its Hessian as $H \stackrel{\text{def}}{=} \nabla^2 f(x) \in \mathbb{R}^{d \times d}$
 - 1) f is **convex** if and only if $H \succcurlyeq 0$ for every $x \in \text{dom } f$.
 - 2) f is **strictly convex** if and only if $H \succ 0$ for every $x \in \text{dom } f$.
 - 3) f is **m-strongly convex** if and only if $H \succcurlyeq mI$ ($m > 0$) for every $x \in \text{dom } f$, where $A \succcurlyeq B$ means that $A - B$ is positive semi-definite.

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

40

40

Positive (Semi-)Definite Matrices

■ Definitions

- A symmetric real matrix is **positive-definite** if the real number $\mathbf{z}^T \mathbf{M} \mathbf{z}$ is **positive** for every nonzero vector $\mathbf{z} \in \mathbb{R}^n$.
- A symmetric real matrix is **positive semi-definite** if the real number $\mathbf{z}^T \mathbf{M} \mathbf{z}$ is **non-negative** for every nonzero vector $\mathbf{z} \in \mathbb{R}^n$.

■ Examples

- Identity matrix \mathbf{I} ($n \times n$) is **positive-definite** as $\mathbf{z}^T \mathbf{I} \mathbf{z} > 0$ for every nonzero vector \mathbf{z} .
- Zero matrix $\mathbf{0}$ ($n \times n$) is **positive semi-definite** as $\mathbf{z}^T \mathbf{0} \mathbf{z} = 0$ for every nonzero vector \mathbf{z} .
- What about real matrix $\mathbf{M} = \mathbf{A}^T \mathbf{A}$ where $\mathbf{A} \in \mathbb{R}^{n \times m}$ is rectangular?

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

41

41

Let's check if $\mathbf{M} = \mathbf{A}^T \mathbf{A}$ is positive(semi-)definite

- For every nonzero $\mathbf{z} \in \mathbb{R}^d$, we have

$$\mathbf{z}^T \mathbf{M} \mathbf{z} = \mathbf{z}^T \mathbf{A}^T \mathbf{A} \mathbf{z} = \|\mathbf{A} \mathbf{z}\|^2 \geq 0$$

- If \mathbf{M} has a full rank (but we did not assume that), then $\mathbf{M} \succcurlyeq m \mathbf{I}$ and m is the smallest eigenvalue (positive) of the matrix (we omit the explanation for now).

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

42

42

Examples (High Dimensional Cases)

(Garret Thomas, <https://gwthomas.github.io/docs/math4ml.pdf>)

- $f(\mathbf{x}) = \mathbf{a}^T \mathbf{x} - b$ for $\mathbf{x} \in \mathbb{R}^d, \mathbf{a} \in \mathbb{R}^d, b \in \mathbb{R}$

$$\nabla_{\mathbf{x}} f = \mathbf{a}, \quad \nabla_{\mathbf{x}}^2 = 0_{d \times d} \quad \leftarrow \text{Convex but not strictly nor strongly}$$

- $f(\mathbf{x}) = \|\mathbf{x}\|_1 = \sum_{i=1}^d |x_i|$ for $\mathbf{x} \in \mathbb{R}^d$

$$\nabla_{\mathbf{x}} f = (\text{sign}(x_1), \text{sign}(x_2), \dots, \text{sign}(x_d)), \quad \nabla_{\mathbf{x}}^2 = ?$$

- $f(\mathbf{x}) = \|\mathbf{x}\|_2^2 = \sum_{i=1}^d x_i^2$ for $\mathbf{x} \in \mathbb{R}^d$

$$\nabla_{\mathbf{x}} f = (2x_1, 2x_2, \dots, 2x_d), \quad \nabla_{\mathbf{x}}^2 = 2\mathbf{I} \quad \leftarrow \text{Strongly convex}$$

- $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x}$ for $\mathbf{x} \in \mathbb{R}^d$ and $\mathbf{A} \in \mathbb{R}^{d \times k}$

$$\nabla_{\mathbf{x}} f = \mathbf{A}^T \mathbf{A} \mathbf{x}, \quad \nabla \nabla_{\mathbf{x}} f = \mathbf{A}^T \mathbf{A} \quad \leftarrow \text{Convexity?}$$

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

43

43

The Hessian of LR

- We have shown $H^{(LR)} = -\mathbf{X}^T \mathbf{A} \mathbf{X}$

- It can be re-written as

$$\mathbf{H}^{(LR)} = -\mathbf{X}^T \mathbf{A} \mathbf{X} = -\underbrace{\mathbf{X}^T \mathbf{A}^{1/2}}_{\mathbf{A}^T} \underbrace{\mathbf{A}^{1/2} \mathbf{X}}_{\mathbf{A}}$$

$$\text{where } \mathbf{A} = \text{diag}(\sigma_i(1 - \sigma_i))_{i=1 \dots n}$$

$$\forall \mathbf{u} \in \mathbb{R}^{m+1}, \quad \mathbf{u}^T \mathbf{H}^{(LR)} \mathbf{u} = -\mathbf{u}^T \mathbf{A}^T \mathbf{A} \mathbf{u} = -\|\mathbf{A} \mathbf{u}\|_2^2 \leq 0$$

\mathbf{H} is **negative semi-definite**. LR has a **concave** objective function.

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

44

44

Outline

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

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

45

45

Regularized Logistic Regression (RLR)

- So far, we have focus on the MLE objective as:

$$\hat{\mathbf{w}}^{LR} = \operatorname{argmax}_{\mathbf{w}, w_0} \{ l_D(\mathbf{w}) \}$$
$$l_D(\mathbf{w}) = \sum_{i=1}^n \{ y_i \ln \sigma(\mathbf{w}^T \mathbf{x}_i) + (1 - y_i) \ln(1 - \sigma(\mathbf{w}^T \mathbf{x}_i)) \}$$

- Now we add a regularization term as:

$$\hat{\mathbf{w}}^{RLR} = \operatorname{argmax}_{\mathbf{w}} \left\{ l_D(\mathbf{w}) - \frac{1}{2} C \|\mathbf{w}\|^2 \right\}$$

- Equivalent to adding a Bayesian prior for \mathbf{w} (next slide)

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

46

46

Maximum A Posterior (MAP) Solution

- Bayesian Prior (Assumption)

$$\mathbf{w} \sim N(0, \sigma^2 I) \quad P(\mathbf{w}) = \frac{1}{Z_0} \exp\left(-\frac{\mathbf{w}^T \mathbf{w}}{2\sigma^2}\right)$$

where Z_0 is some constant (normalization factor).

- Posterior Probability

$$P(\mathbf{w}|D) = \frac{P(D|\mathbf{w})P(\mathbf{w})}{P(D)} \propto P(D|\mathbf{w})P(\mathbf{w})$$

- Objective

$$\begin{aligned} \hat{\mathbf{w}}^{RLR} &= \operatorname{argmax}_{\mathbf{w}} P(\mathbf{w}|D) = \operatorname{argmax}_{\mathbf{w}} P(D|\mathbf{w})P(\mathbf{w}) \\ &= \operatorname{argmax}_{\mathbf{w}} \left\{ \underbrace{\log P(D|\mathbf{w})}_{l(\mathbf{w})} + \log P(\mathbf{w}) \right\} \\ &= \operatorname{argmax}_{\mathbf{w}} \{ l(\mathbf{w}) - \lambda \mathbf{w}^T \mathbf{w} + \text{some constant} \} \\ &\quad \text{where } \lambda = \frac{1}{2\sigma^2} \end{aligned}$$

- MAP solution for RLR assumes a “non-informative” Gaussian prior of \mathbf{w} .

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

47

47

L1 vs L2 regularization

(figure from Elements of Stat. Learn., Hastie et al.)

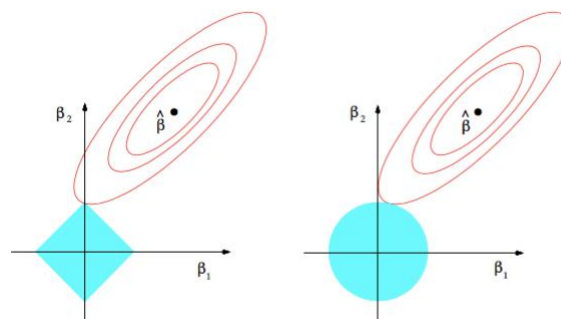


FIGURE 3.11. Estimation picture for the lasso (left) and ridge regression (right). Shown are contours of the error and constraint functions. The solid blue areas are the constraint regions $|\beta_1| + |\beta_2| \leq t$ and $\beta_1^2 + \beta_2^2 \leq t^2$, respectively, while the red ellipses are the contours of the least squares error function.

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

48

48

Outline

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

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

49

49

Softmax LR

- Let $\mathbf{x} \in R^d$ be the input vector, $k = 1, 2, \dots, K$ be the index of category labels, and $Y \in \{1, 2, \dots, K\}$ be the output variable.
- Assume a categorical (multinomial) distribution of labels given \mathbf{x} , with $p_k(\mathbf{x}) \stackrel{\text{def}}{=} P(Y = k | \mathbf{x})$ and $\sum_{k=1}^K p_k(\mathbf{x}) = 1$,
$$Y | \mathbf{x} \sim \text{Cat}(p_1(\mathbf{x}), p_2(\mathbf{x}), \dots, p_K(\mathbf{x}))$$
- The system predicts the probability of each category as

$$\hat{p}_k(\mathbf{x}) = \frac{\exp(\mathbf{w}_k^T \mathbf{x})}{\sum_{k'=1}^K \exp(\mathbf{w}_{k'}^T \mathbf{x})}, \quad k = 1, 2, \dots, K$$

where $\mathbf{w}_k^T \in R^d$ is a category-specific vector of model parameters .

- Decision Rule: $\hat{Y} | \mathbf{x} = \text{argmax}_k \{ \hat{p}_k(\mathbf{x}) \}$

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

50

50

Log-likelihood of Softmax LR

$$P(D|W) = \prod_{i=1}^N \prod_{k=1}^K p_k(x_i)^{y_{ik}} \quad y_{ik} \in \{0,1\}$$

$$\log P(D|W) = \sum_{i=1}^N \sum_{k=1}^K y_{ik} \log p_k(x_i)$$

$$= \sum_{i=1}^N \sum_{k=1}^K y_{ik} \log \left(\frac{\exp(w_k^T x_i)}{\sum_{k'=1}^K \exp(w_{k'}^T x_i)} \right)$$

$$= \sum_{i=1}^N \sum_{k=1}^K y_{ik} w_k^T x_i - \sum_{i=1}^N \log \sum_{k'=1}^K \exp(w_{k'}^T x_i)$$

where $D = \{(x_i, y_i)\}$, $W = \{w_k\}_{k=1}^K$ and $y_{ik} = I(y(x_i) = k)$

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

51

51

Loss Function of Regularized LR

$$J(W; D) = \frac{\lambda}{2} \sum_{k=1}^K \|w_k\|^2 - \log P(D|W) \quad (\lambda > 0)$$

$$= \frac{\lambda}{2} \sum_{k=1}^K \|w_k\|^2 - \sum_{i=1}^N \sum_{k=1}^K y_{ik} w_k^T x_i + \sum_{i=1}^N \log \sum_{k'=1}^K \exp(w_{k'}^T x_i)$$

- Is Softmax convex? (Yes, but tricks are needed for proving)
- How to minimize the loss function?
- Where is the computational bottleneck, and how to alleviate it?

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

52

52

Convexity of Softmax LR

$$J(W; D) = \frac{\lambda}{2} \sum_{k=1}^K \|\mathbf{w}_k\|^2 - \log P(D|W)$$

$$= \frac{\lambda}{2} \sum_{k=1}^K \|\mathbf{w}_k\|^2 - \sum_{i=1}^N \sum_{k=1}^K y_{ik} \mathbf{w}_k^T \mathbf{x}_i + \sum_{i=1}^N \log \sum_{k'=1}^K \exp(\mathbf{w}_{k'}^T \mathbf{x}_i)$$

- The first term (regularization) is convex because the non-negatively weighted sum of convex function is convex.
- The 2nd term is a linear function (convex and concave) of model parameters.
- The 3rd term is about the convexity of the **log-sum-exp (LSE) function**.

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

53

53

Convexity of log-sum-exponential

- <https://math.stackexchange.com/questions/2418554/why-is-log-of-sum-of-exponentials-convex>
 - Proving based on the convexity definition
- <https://math.stackexchange.com/questions/2416837/the-second-derivative-of-log-left-sum-limits-i-1-n-e-x-i-right-seems-negative>
 - By showing the Hessian of Softmax LR to be positive semi-definite

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

54

54

Optimization of Softmax LR

$$\begin{aligned} \min_W J(W; D) &= \\ &= \min_W \frac{\lambda}{2} \sum_{k=1}^K \|\mathbf{w}_k\|^2 - \sum_{i=1}^N \sum_{k=1}^K y_{ik} \mathbf{w}_k^T \mathbf{x}_i + \sum_{i=1}^N \log \sum_{k'=1}^K \exp(\mathbf{w}_{k'}^T \mathbf{x}_i) \end{aligned}$$

- It cannot be solved analytically.
- Instead, it should be solved with GD or SGD iteratively.
- Exercise by yourself in HW2: derive $\nabla_{\mathbf{w}_k} J(W; D)$

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

55

55

Computing the gradients

$$J(W; D) = \underbrace{\frac{\lambda}{2} \sum_{k=1}^K \|\mathbf{w}_k\|^2}_{f(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K)} - \underbrace{\sum_{i=1}^N \sum_{k=1}^K y_{ik} \mathbf{w}_k^T \mathbf{x}_i}_{g(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K)} + \underbrace{\sum_{i=1}^N \log \sum_{k'=1}^K \exp(\mathbf{w}_{k'}^T \mathbf{x}_i)}_{\varphi(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K)}$$

$$\frac{\partial f}{\partial \mathbf{w}_k} = \frac{\lambda}{2} \frac{\partial}{\partial \mathbf{w}_k} \|\mathbf{w}_k\|^2 = \lambda \mathbf{w}_k \quad \begin{array}{l} \text{input is a vector; output is a scalar} \\ \rightarrow \text{the gradient is a vector.} \end{array}$$

$$\frac{\partial g}{\partial \mathbf{w}_k} = \sum_{i=1}^N y_{ik} \frac{\partial}{\partial \mathbf{w}_k} \mathbf{w}_k^T \mathbf{x}_i = \sum_{i=1}^N y_{ik} \mathbf{x}_i$$

$$\frac{\partial \varphi}{\partial \mathbf{w}_k} = \sum_{i=1}^N \frac{\partial}{\partial \mathbf{w}_k} LSE = \sum_{i=1}^N \frac{\partial(\log v)}{\partial v} \dots \frac{\partial z}{\partial \mathbf{w}_k} \quad (v = \sum_{k'=1}^K \exp(\mathbf{w}_{k'}^T \mathbf{x}_i))$$

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

56

56

Parallel Computing Bottleneck

$$\min_w \frac{\lambda}{2} \sum_{k=1}^K \|w_k\|^2 - \sum_{i=1}^N \sum_{k=1}^K y_{ik} w_k^T x_i + \sum_{i=1}^N \log \sum_{k'=1}^K \exp(w_{k'}^T x_i)$$

- Easy to decouple the updates of w_k 's updates for first two terms
- But how can we decouple for 3rd term (on different processors)?

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

57

57

Distributed training of regularized LR [Gopal & Yang, ICML 2013]

- Objective

$$L(W) = \frac{\lambda}{2} \sum_{k=1}^K \|w_k\|^2 - \sum_{i=1}^N \sum_{k=1}^K y_{ik} w_k^T x_i + \sum_{i=1}^N \log \sum_{k'=1}^K \exp(w_{k'}^T x_i)$$

- **Log-concavity bound** (the 1st order concavity property) of log function

$$\log(v) \leq \alpha v - \log(\alpha) - 1 \quad \forall v, \quad \alpha > 0$$

- **Log-partition term** for each instance i in LR is bounded as **substitute**

$$\log \left(\sum_{k=1}^K \exp(w_k^T x_i) \right) \leq \underbrace{\alpha_i}_{\text{variational parameter}} \sum_{k=1}^K \exp(w_k^T x_i) - \log(\alpha_i) - 1 \quad \alpha_i > 0$$

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

58

58

The Modified Objective Function

$$\min_{\alpha > 0, \mathbf{w}} F(\mathbf{W}, \alpha)$$
$$F(\mathbf{W}, \alpha) = \frac{\lambda}{2} \sum_{k=1}^K \|w_k\|^2 - \sum_{i=1}^N \sum_{k=1}^K (y_{ik} w_k^T x_i - \alpha_i \exp(w_k^T x_i) - \log \alpha_i - 1)$$

Convergence-related Properties (proof in Gopal's ICML 2013 paper)

- 1) It does not preserve the convexity of the original objective function (because we have α_i 's as the additional variables).
- 2) However, it has exactly one stationary point that is the same stationary point of the original convex function of softmax.
- 3) A block co-ordinate descent procedure guarantees to converge to the stationary point which is the optimal solution of softmax.

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

59

59

Concluding Remarks about LR

- Explicit probabilistic reasoning
- Easy to extend with regularization terms (e.g., L1 or L2 norm of the parameter vector)
- Commonly used as building blocks in neural nets
- Expressive enough for complicated decision boundaries?

02/06/2024

@Yiming Yang, S24 Lecture on LR Models

60

60

A good online lecture on LR by Andrew Ng

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