

线性分类器

鲁鹏

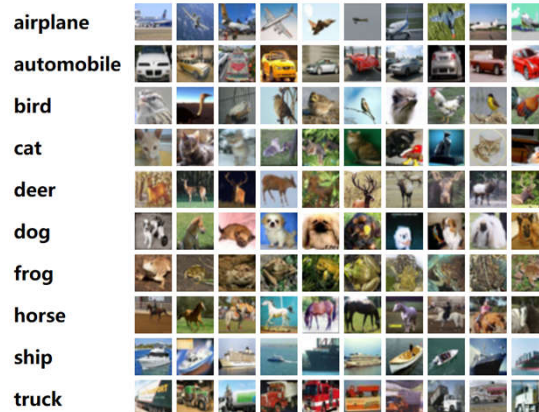
北京邮电大学 计算机学院 智能科学与技术中心

0

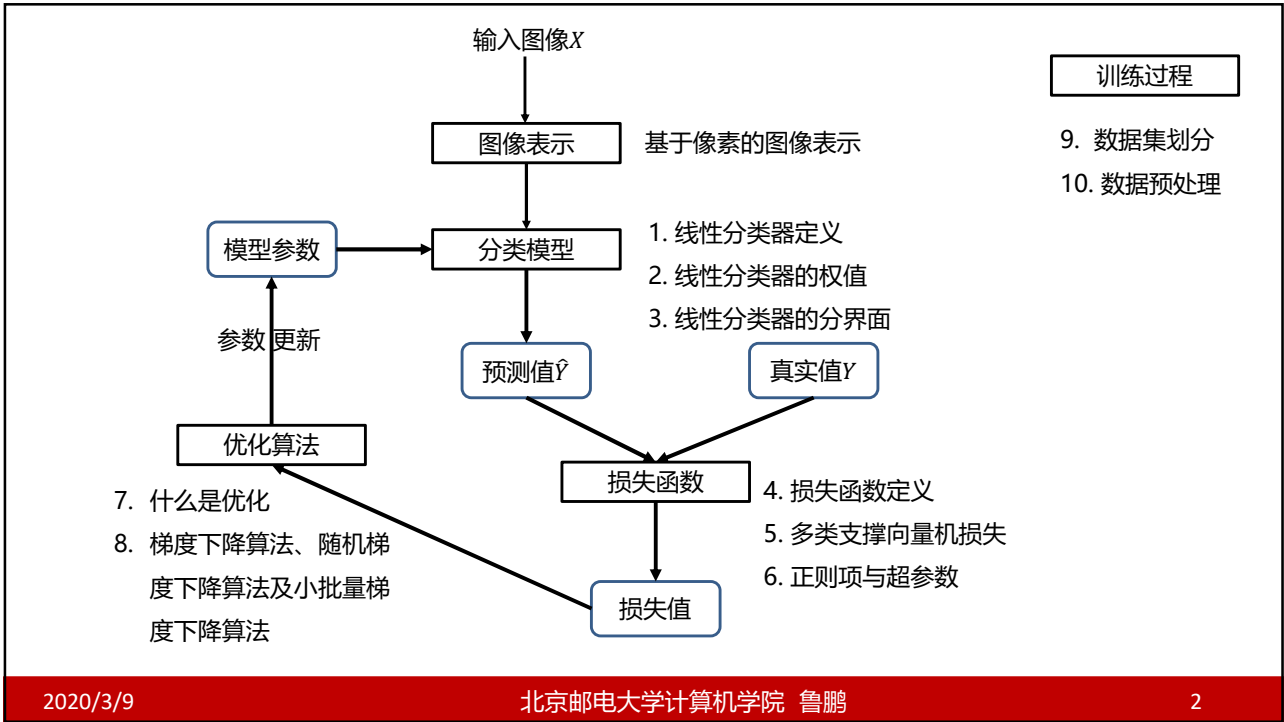
数据集介绍

➤ CIFAR10数据集：

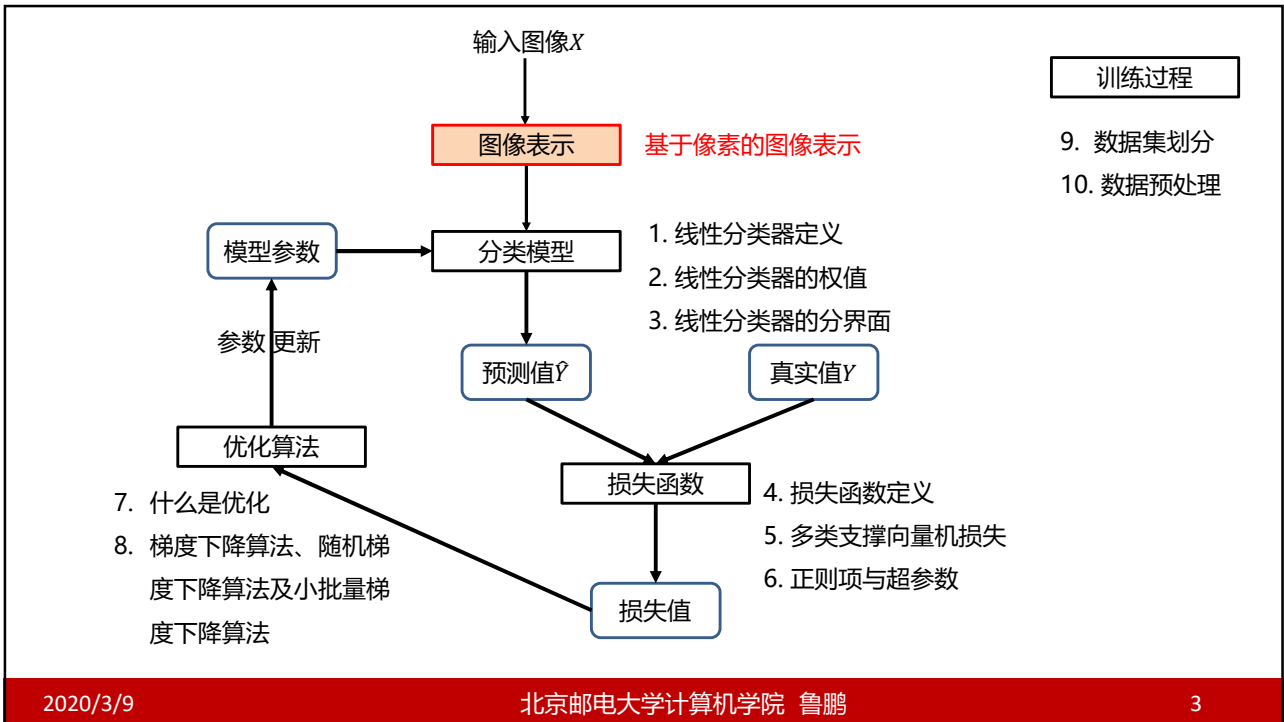
包含50000张训练样本、10000张
测试样本分为飞机、汽车、鸟、猫、
鹿、狗、蛙、马、船、卡车十个类
别图像为彩色图像，其大小为
32*32



1



2



3

图像类型

Binary



Gray Scale



Color



Source: Ulas Bagci

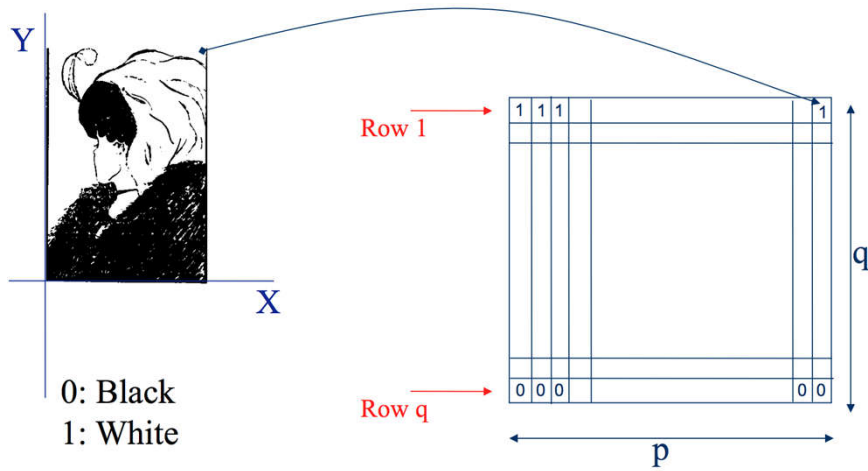
2020/3/9

北京邮电大学计算机学院 鲁鹏

4

4

二进制图像



Source: Ulas Bagci

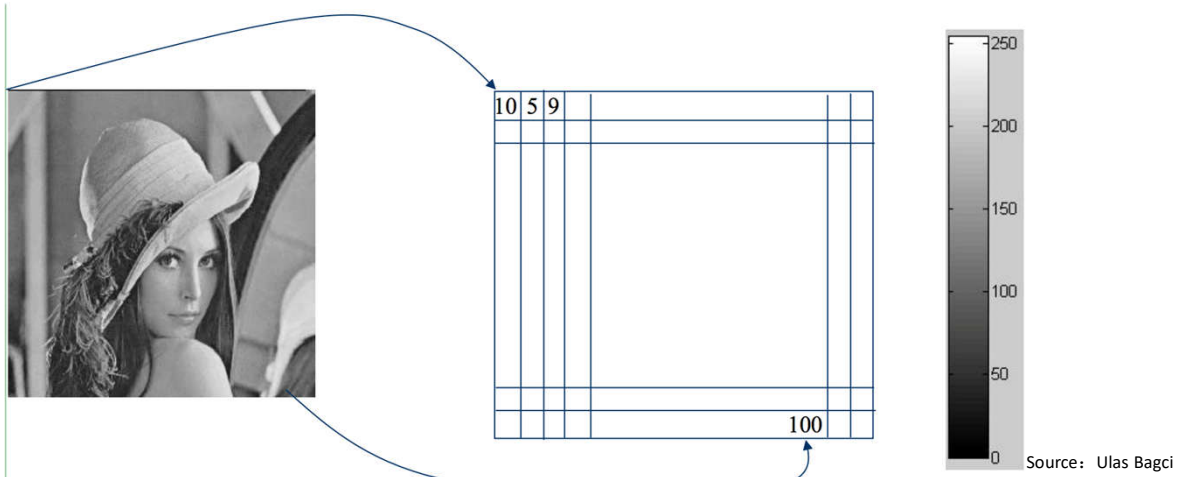
2020/3/9

北京邮电大学计算机学院 鲁鹏

5

5

灰度图像



2020/3/9

北京邮电大学计算机学院 鲁鹏

6

6

彩色图像



Source: Ulas Bagci

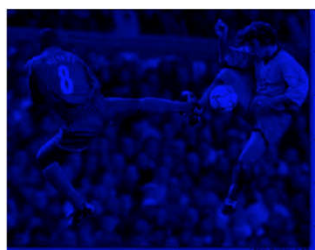
2020/3/9

北京邮电大学计算机学院 鲁鹏

7

7

彩色图像



Source: Ulas Bagci

2020/3/9

北京邮电大学计算机学院 鲁鹏

8

8

图像表示

大多数分类算法都要求**输入向量**!

2020/3/9

北京邮电大学计算机学院 鲁鹏

9

9

图像表示

将图像转换成向量的方法有很多，这里我们用一种最简单的方法，直接将图像矩阵转换成向量



图像

将矩阵转成列向量



$$x = \begin{bmatrix} r_1 \\ g_1 \\ b_1 \\ \vdots \\ r_n \\ g_n \\ b_n \end{bmatrix}$$

2020/3/9

北京邮电大学计算机学院 鲁鹏

10

10

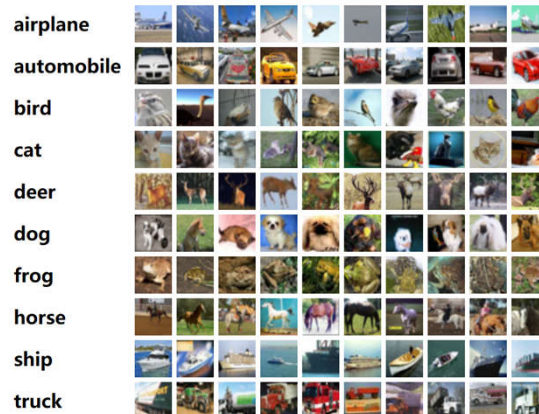
图像表示

CIFAR10中每一张图像转

换为向量是多少维？

答案：3072 (=32*32*3)

维列向量

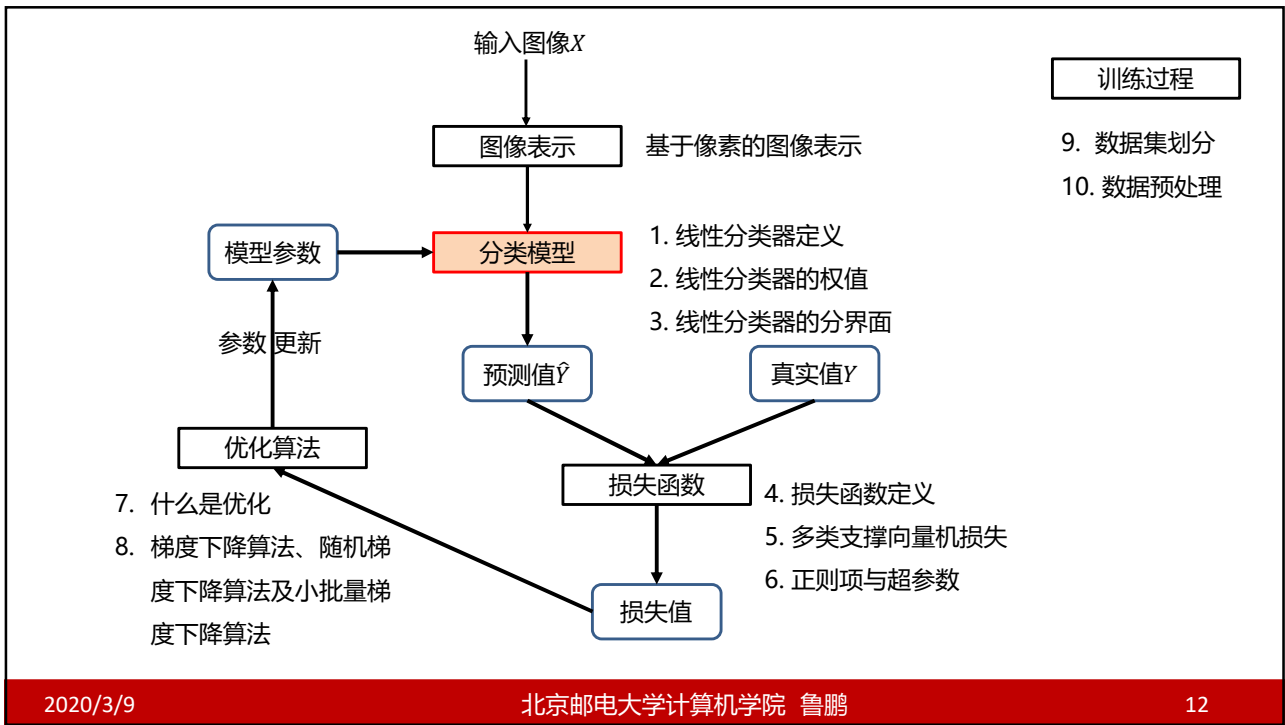


2020/3/9

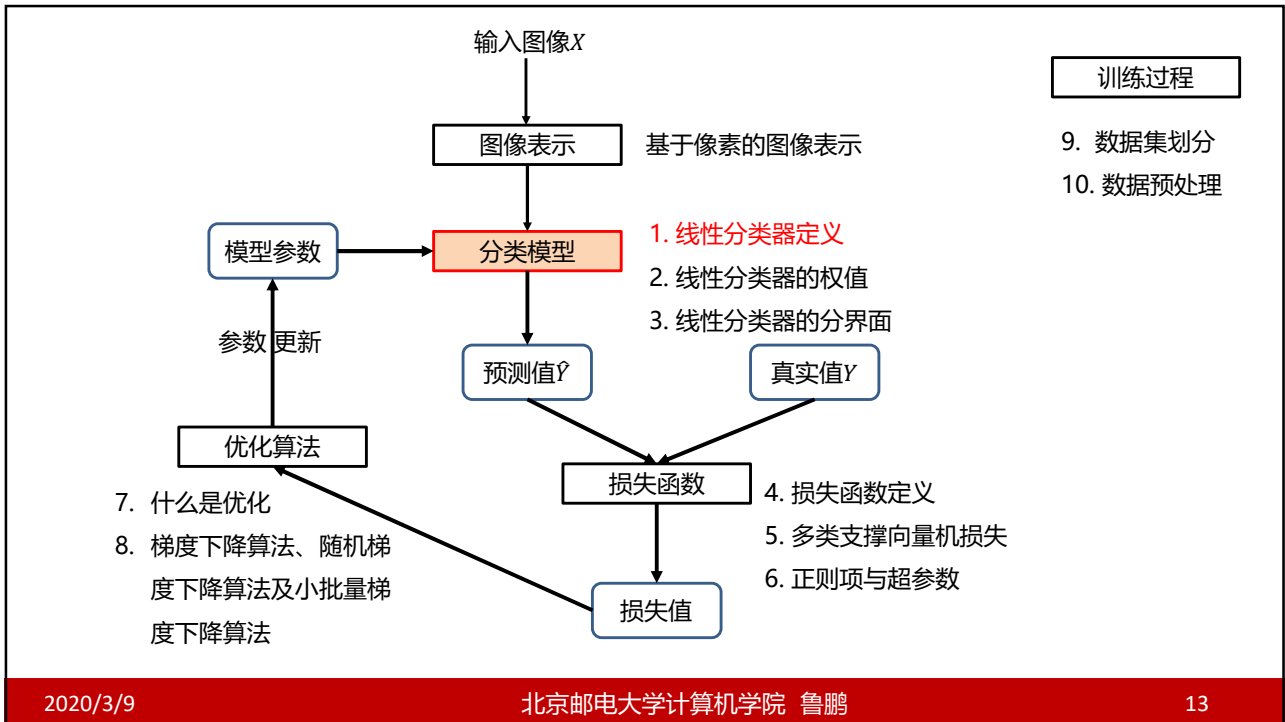
北京邮电大学计算机学院 鲁鹏

11

11



12



13

为什么从线性分类器开始?

- 形式简单、易于理解
- 通过层级结构（神经网络）或者高维映射（支撑向量机）可以形成功能强大的非线性模型

什么是线性分类器

线性分类器是一种线性映射，将输入的图像特征映射为类别分数。

线性分类器定义

x 代表输入的 d 维图像向量, c 为类别个数

线性分类器定义

第 i 个类的线性分类器:

$$f_i(\mathbf{x}, \mathbf{w}_i) = \mathbf{w}_i^T \mathbf{x} + b_i,$$

$$i = 1, \dots, c$$

x 代表输入的 d 维图像向量, c 为类别个数

线性分类器定义

第*i*个类的线性分类器:

$$f_i(\mathbf{x}, \mathbf{w}_i) = \mathbf{w}_i^T \mathbf{x} + b_i,$$

$$i = 1, \dots, c$$

\mathbf{x} 代表输入的 d 维图像向量, c 为类别个数

$\mathbf{w}_i = [w_{i1} \ \dots \ w_{id}]^T$ 为第 i 个类别的权值向量, b_i 为偏置

2020/3/9

北京邮电大学计算机学院 鲁鹏

18

18

线性分类器定义

第*i*个类的线性分类器: 每个类都有自己的参数 \mathbf{w} 和 b

$$f_i(\mathbf{x}, \mathbf{w}_i) = \mathbf{w}_i^T \mathbf{x} + b_i,$$

$$i = 1, \dots, c$$

\mathbf{x} 代表输入的 d 维图像向量, c 为类别个数

$\mathbf{w}_i = [w_{i1} \ \dots \ w_{id}]^T$ 为第 i 个类别的权值向量, b_i 为偏置

2020/3/9

北京邮电大学计算机学院 鲁鹏

19

19

线性分类器决策

第*i*个类的线性分类器: 每个类都有自己的参数 w 和 b

$$f_i(\mathbf{x}, \mathbf{w}_i) = \mathbf{w}_i^T \mathbf{x} + b_i,$$

$$i = 1, \dots, c$$

\mathbf{x} 代表输入的 d 维图像向量, c 为类别个数

$\mathbf{w}_i = [w_{i1} \ \dots \ w_{id}]^T$ 为第 i 个类别的权值向量, b_i 为偏置

决策规则:

如果 $f_i(\mathbf{x}) > f_j(\mathbf{x}), \forall j \neq i,$

则决策输入图像 \mathbf{x} 属于第 i 类

线性分类器示例

任务: 为图片分配类别标签 (汽车类、猫类、鸟类)



图片

线性分类器

$$\mathbf{w}_i^T \mathbf{x} + b_i = f_i, i = 1, 2, 3$$

?

线性分类器示例

图像表示 x

$$\begin{bmatrix} 56 \\ 231 \\ 24 \\ 2 \end{bmatrix}$$

线性分类器决策步骤:

1. 图像表示成向量



线性分类器示例

	权值 w_i^T	图像表示 x	偏移 b_i	得分 f_i					
汽车类	w_1^T <table border="1"><tr><td>0.2</td><td>-0.5</td><td>0.1</td><td>2.0</td></tr></table>	0.2	-0.5	0.1	2.0	$\begin{bmatrix} 56 \\ 231 \\ 24 \\ 2 \end{bmatrix}$	1.1 b_1	f_1 <table border="1"><tr><td>-97.9</td></tr></table>	-97.9
0.2	-0.5	0.1	2.0						
-97.9									
猫类	w_2^T <table border="1"><tr><td>1.5</td><td>1.3</td><td>2.1</td><td>0.0</td></tr></table>	1.5	1.3	2.1	0.0	3.2 b_2	f_2 <table border="1"><tr><td>434.7</td></tr></table>	434.7	
1.5	1.3	2.1	0.0						
434.7									
鸟类	w_3^T <table border="1"><tr><td>0</td><td>0.25</td><td>0.25</td><td>-0.3</td></tr></table>	0	0.25	0.25	-0.3	-1.2 b_3	f_3 <table border="1"><tr><td>63.15</td></tr></table>	63.15	
0	0.25	0.25	-0.3						
63.15									

线性分类器决策步骤:

1. 图像表示成向量

2. 计算当前图片每个类别的分数

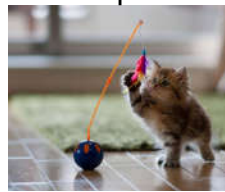


线性分类器示例

	权值 w_i^T	图像表示 x	偏移 b_i	得分 f_i
汽车类	w_1^T [0.2, -0.5, 0.1, 2.0]	[56, 231, 24, 2]	1.1 b_1	f_1 -97.9
猫类	w_2^T [1.5, 1.3, 2.1, 0.0]		3.2 b_2	f_2 434.7
鸟类	w_3^T [0, 0.25, 0.25, -0.3]		-1.2 b_3	f_3 63.15

线性分类器决策步骤:

1. 图像表示成向量
2. 计算当前图片每个类别的分数
3. 按类别得分判定当前图像



→ 猫类



线性分类器示例

	权值 w_i^T	图像表示 x	偏移 b_i	得分 f_i
汽车类	w_1^T [0.2, -0.5, 0.1, 2.0]	[56, 231, 24, 2]	1.1 b_1	f_1 -97.9
猫类	w_2^T [1.5, 1.3, 2.1, 0.0]		3.2 b_2	f_2 434.7
鸟类	w_3^T [0, 0.25, 0.25, -0.3]		-1.2 b_3	f_3 63.15

权值矩阵 W

偏移向量 b

得分向量 f



线性分类器的矩阵表示

	权值 w_i^T	图像表示 x	偏移 b_i	得分 f_i
汽车类	w_1^T	$\begin{bmatrix} 56 \\ 231 \\ 24 \\ 2 \end{bmatrix}$	1.1 b_1	f_1 -97.9
猫类	w_2^T		3.2 b_2	f_2 434.7
鸟类	w_3^T		-1.2 b_3	f_3 63.15
	权值矩阵 W		偏移向量 b	得分向量 f



线性分类器的矩阵表示:

$$f(x, W) = Wx + b$$

其中, x 代表输入图像, 其维度为 d ,

f 为分数向量, 其维度等于类别个数 c ,

$W = [w_1 \ \dots \ w_c]^T$ 为权值矩阵,

$w_i = [w_{i1} \ \dots \ w_{id}]^T$ 为第 i 个类别的权值向量,

$b = [b_1 \ \dots \ b_c]^T$ 为偏置向量, b_i 为第 i 个类别的偏置。

26

线性分类器的矩阵表示

问题: CIFAR10 数据集分类任务的分类器, W , x , b 的维度是多少?

线性分类器的矩阵表示:

$$f(x, W) = Wx + b$$

其中, x 代表输入图像, 其维度为 d ,

f 为分数向量, 其维度等于类别个数 c ,

$W = [w_1 \ \dots \ w_c]^T$ 为权值矩阵,

$w_i = [w_{i1} \ \dots \ w_{id}]^T$ 为第 i 个类别的权值向量,

$b = [b_1 \ \dots \ b_c]^T$ 为偏置向量, b_i 为第 i 个类别的偏置。

27

线性分类器的矩阵表示

问题：CIFAR10 数据集分类任务的分类器， W ， x ， b 的维度是多少？

回答：CIFAR10有10个类别且图像大小为32x32x3，因此：

x 是图像向量，其维度为3072维；

W 是权值矩阵，其维度为10x3072；

b 是偏置向量，其维度为10x1的向量；

f 是得分向量，其维度为10x1的向量；

线性分类器的矩阵表示：

$$f(x, W) = Wx + b$$

其中， x 代表输入图像，其维度为 d ，

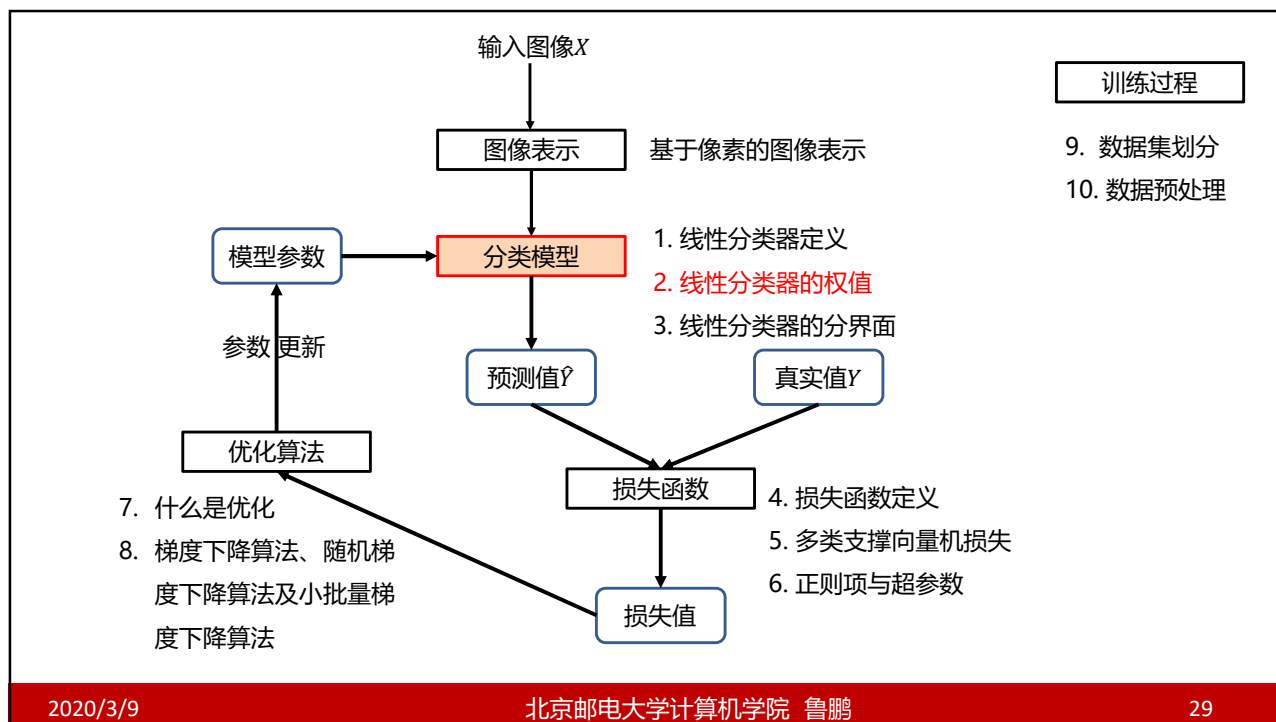
f 为分数向量，其维度等于类别个数 c ，

$W = [w_1 \ \dots \ w_c]^T$ 为权值矩阵，

$w_i = [w_{i1} \ \dots \ w_{id}]^T$ 为第 i 个类别的权值向量，

$b = [b_1 \ \dots \ b_c]^T$ 为偏置向量， b_i 为第 i 个类别的偏置。

28



2020/3/9

北京邮电大学计算机学院 鲁鹏

29

29

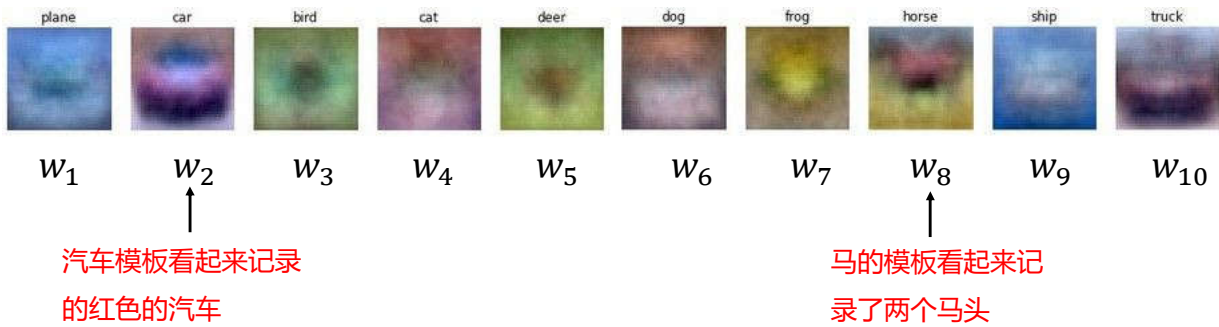
线性分类器的权值向量

第*i*个类的线性分类器:

$$f_i(x, w_i) = w_i^T x + b_i, \\ i = 1, \dots, c$$

x 代表输入的 d 维图像向量, c 为类别个数

$w_i = [w_{i1} \ \dots \ w_{id}]^T$ 为第 i 个类别的权值向量, b_i 为偏置



2020/3/9

北京邮电大学计算机学院 鲁鹏

30

30

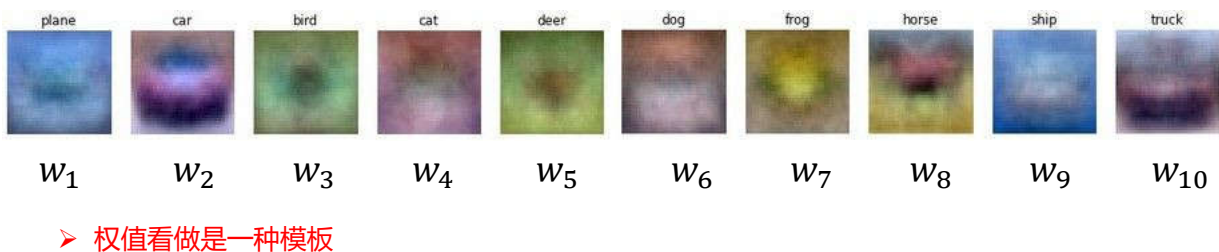
线性分类器的权值向量

第*i*个类的线性分类器:

$$f_i(x, w_i) = w_i^T x + b_i, \\ i = 1, \dots, c$$

x 代表输入的 d 维图像向量, c 为类别个数

$w_i = [w_{i1} \ \dots \ w_{id}]^T$ 为第 i 个类别的权值向量, b_i 为偏置



2020/3/9

北京邮电大学计算机学院 鲁鹏

31

31

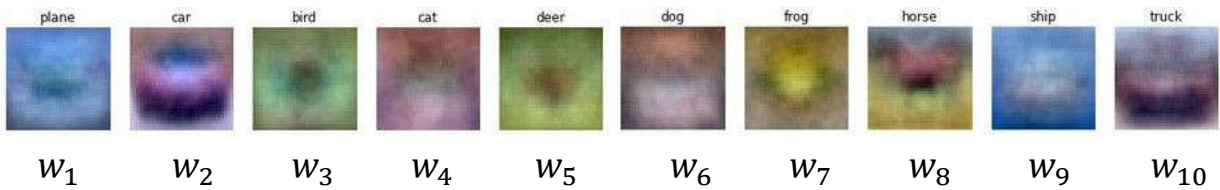
线性分类器的权值向量

第*i*个类的线性分类器:

$$f_i(x, w_i) = w_i^T x + b_i, \quad i = 1, \dots, c$$

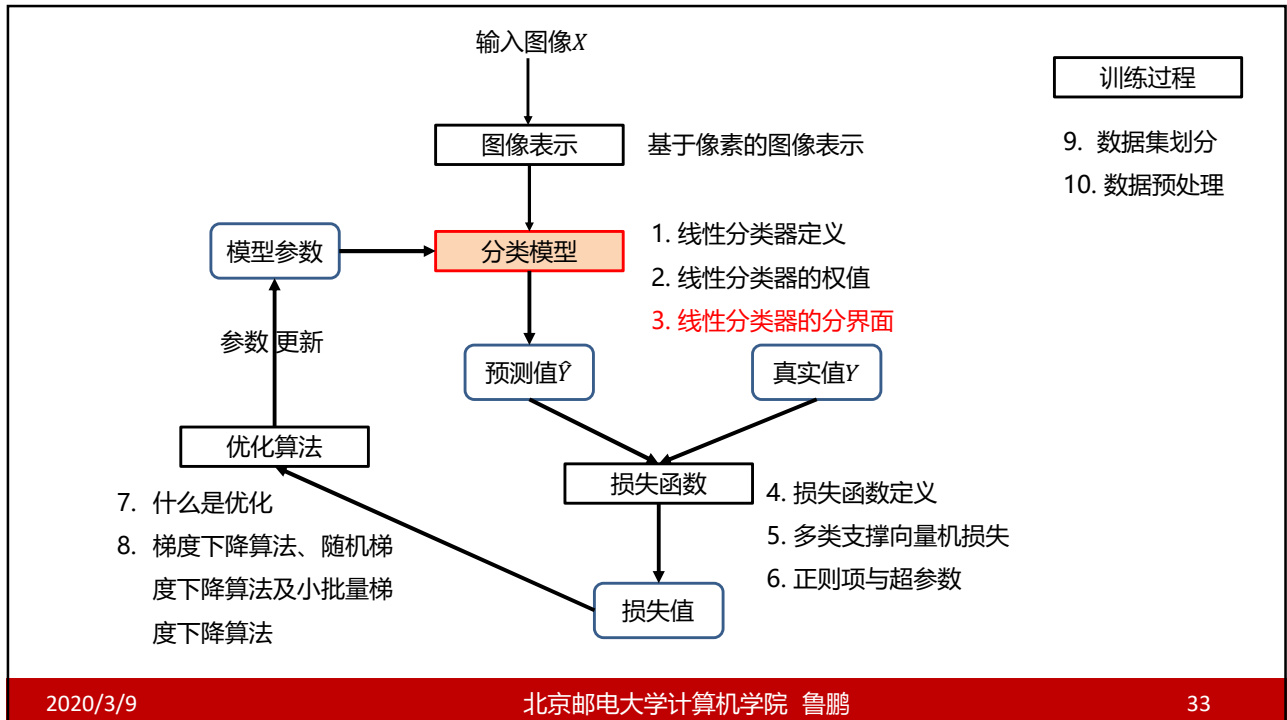
x 代表输入的 d 维图像向量, c 为类别个数

$w_i = [w_{i1} \ \dots \ w_{id}]^T$ 为第 i 个类别的权值向量, b_i 为偏置

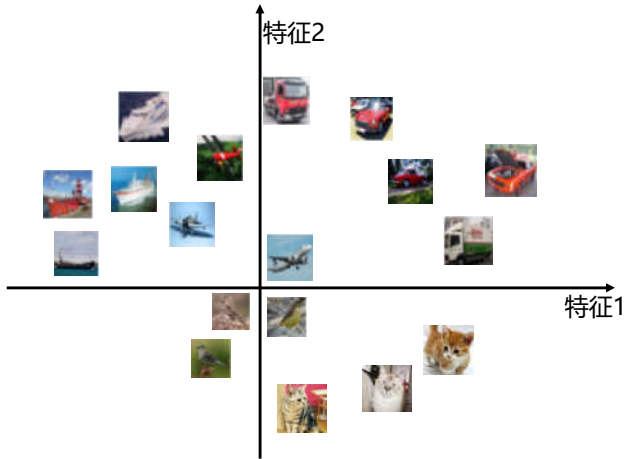


➤ 权值看做是一种模板

➤ 输入图像与评估模板的匹配程度越高, 分类器输出的分数就越高



线性分类器的决策边界



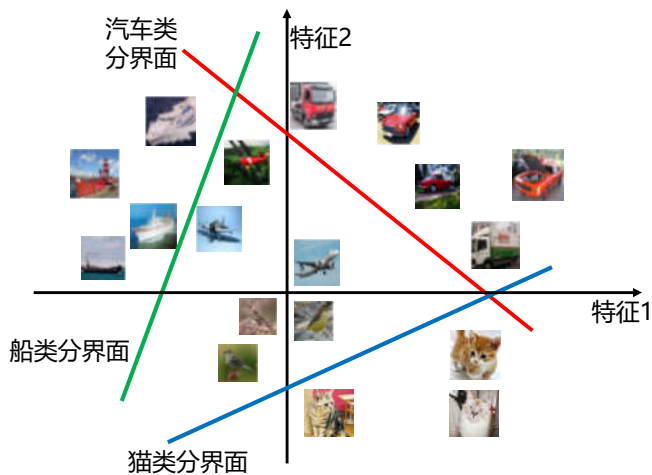
2020/3/9

北京邮电大学计算机学院 鲁鹏

34

34

线性分类器的决策边界



- 分数等于0的线就是决策面

$$w_i^T x + b_i = 0, i = 1, \dots, c$$

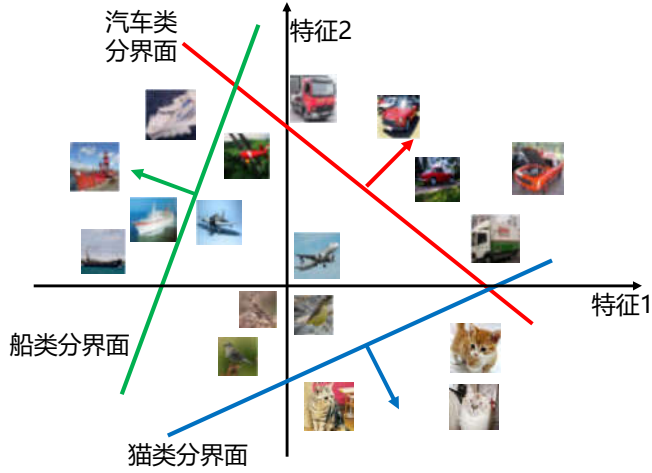
2020/3/9

北京邮电大学计算机学院 鲁鹏

35

35

线性分类器的决策边界



- 分数等于0的线就是决策面

$$w_i^T x + b_i = 0, i = 1, \dots, c$$

- w 控制着线的方向

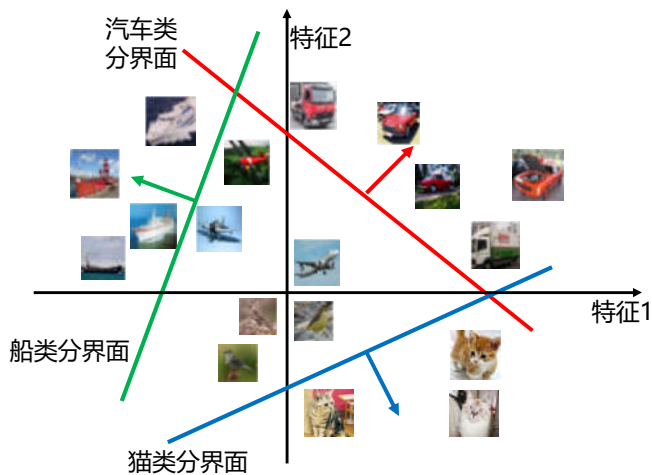
2020/3/9

北京邮电大学计算机学院 鲁鹏

36

36

线性分类器的决策边界



- 分数等于0的线就是决策面

$$w_i^T x + b_i = 0, i = 1, \dots, c$$

- w 控制着线的方向

- b 控制着线的偏移

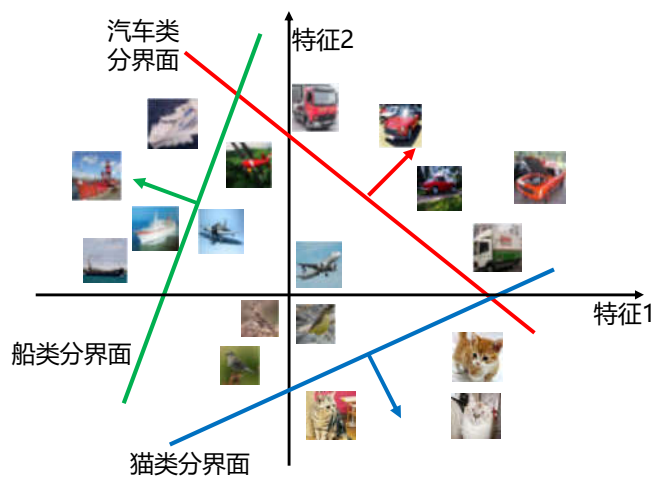
2020/3/9

北京邮电大学计算机学院 鲁鹏

37

37

线性分类器的决策边界



- 分数等于0的线就是决策面

$$w_i^T x + b_i = 0, i = 1, \dots, c$$
- w 控制着线的方向
- b 控制着线的偏移
- 箭头方向代表分类器的正方向, 沿着箭头方向距离决策面越远分数就越高。

2020/3/9

北京邮电大学计算机学院 鲁鹏

38

38

线性分类器小结

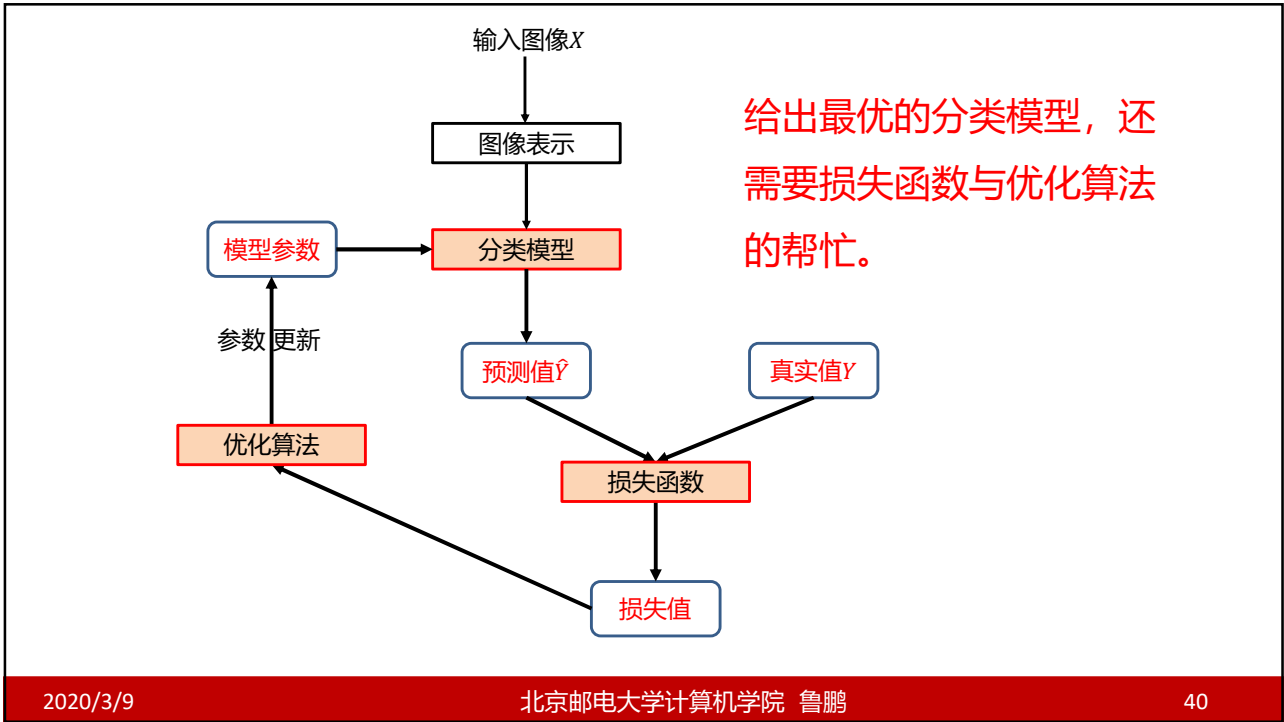
1. 线性分类器的定义
2. 线性分类器的决策
3. 线性分类器的矩阵表示
4. 线性分类器的权值向量
5. 线性分类器的决策边界

2020/3/9

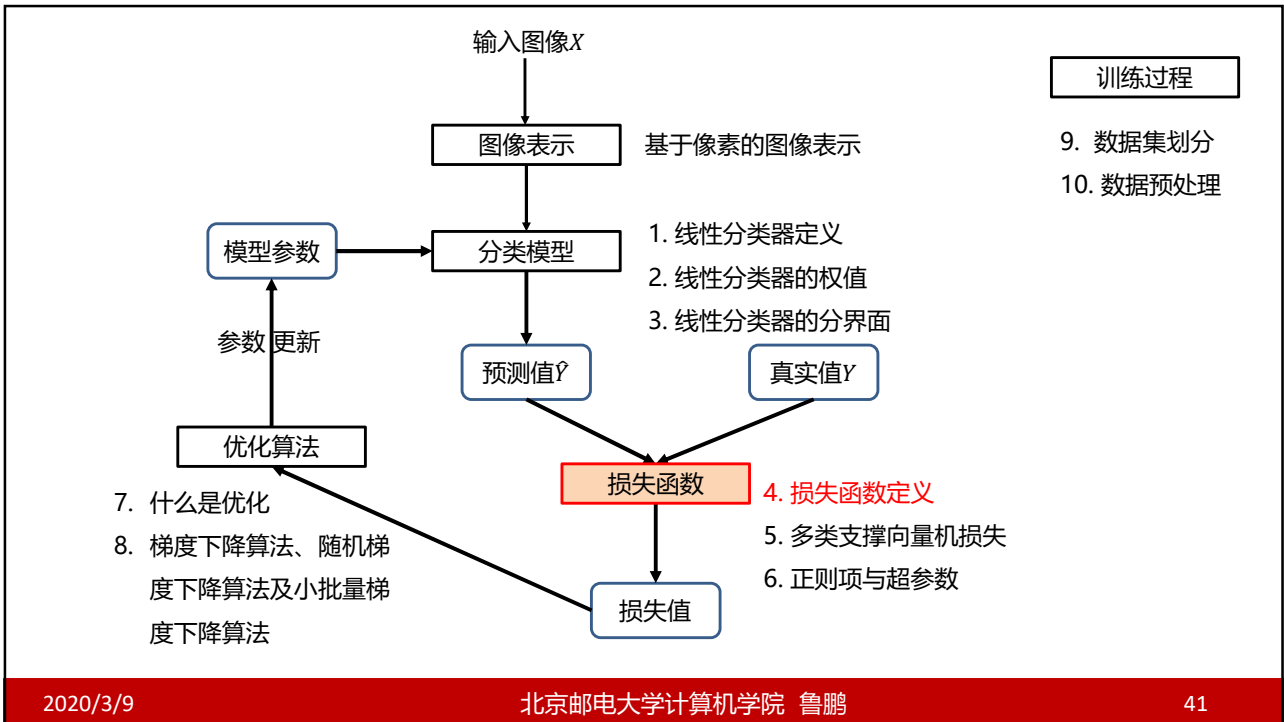
北京邮电大学计算机学院 鲁鹏

39

39



40



41

如何衡量分类器对当前样本的效果好坏？

分类器1

0.2	-0.5	0.1	2.0
1.5	1.3	2.1	0.0
0	0.25	0.25	-0.3


W1

$$\begin{bmatrix} 56 \\ 231 \\ 24 \\ 2 \end{bmatrix} + \begin{bmatrix} 1.1 \\ 3.2 \\ -1.2 \end{bmatrix} = \begin{bmatrix} -97.9 \\ 434.7 \\ 63.15 \end{bmatrix}$$

b1

f

- 汽车类分数: -97.9
- 猫类分数: 434.7 ✓
- 船类分数: 63.15



2020/3/9

北京邮电大学计算机学院 鲁鹏

42

42

如何衡量分类器对当前样本的效果好坏？

分类器2

0.2	-0.1	0.1	1.7
1.5	0.2	2.1	5.0
-1.2	1.5	0.25	0.5


W2

$$\begin{bmatrix} 56 \\ 231 \\ 24 \\ 2 \end{bmatrix} + \begin{bmatrix} 2.1 \\ -0.2 \\ 2.4 \end{bmatrix} = \begin{bmatrix} -6.1 \\ 190.6 \\ 286.3 \end{bmatrix}$$

b2

f

- 汽车类分数: -6.1
- 猫类分数: 190.6
- 船类分数: 286.3 ✗



2020/3/9

北京邮电大学计算机学院 鲁鹏

43

43

损失函数

对示例样本，分类器1与分类器2的分类谁的效果更好？

需要损失函数来帮忙

什么是损失函数？

损失函数搭建了模型性能与模型参数之间的桥梁，指导模型参数优化。

- **损失函数**是一个函数，用于度量给定分类器的预测值与真实值的不一致程度，其**输出**通常是一个**非负实值**。
- 其输出的非负实值可以作为**反馈信号**来对分类器参数进行调整，以**降低当前示例**对应的**损失值**，**提升**分类器的**分类效果**。

什么是损失函数?

损失函数的一般定义:

$$L = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i)$$

2020/3/9

北京邮电大学计算机学院 鲁鹏

46

46

什么是损失函数?

损失函数的一般定义:

x_i 表示数据集中第 i 张图片;

$$L = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i)$$

2020/3/9

北京邮电大学计算机学院 鲁鹏

47

47

什么是损失函数?

损失函数的一般定义:

$$L = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i)$$

x_i 表示数据集中第 i 张图片;

$f(x_i, W)$ 为分类器对 x_i 的类别预测;

什么是损失函数?

损失函数的一般定义:

$$L = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i)$$

x_i 表示数据集中第 i 张图片;

$f(x_i, W)$ 为分类器对 x_i 的类别预测;

y_i 为样本 i 真实类别标签 (整数);

什么是损失函数?

损失函数的一般定义:

$$L = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i)$$

x_i 表示数据集中第 i 张图片;

$f(x_i, W)$ 为分类器对 x_i 的类别预测;

y_i 为样本 i 真实类别标签 (整数);

L_i 为第 i 个样本的损失当预测值;

什么是损失函数?

损失函数的一般定义:

$$L = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i)$$

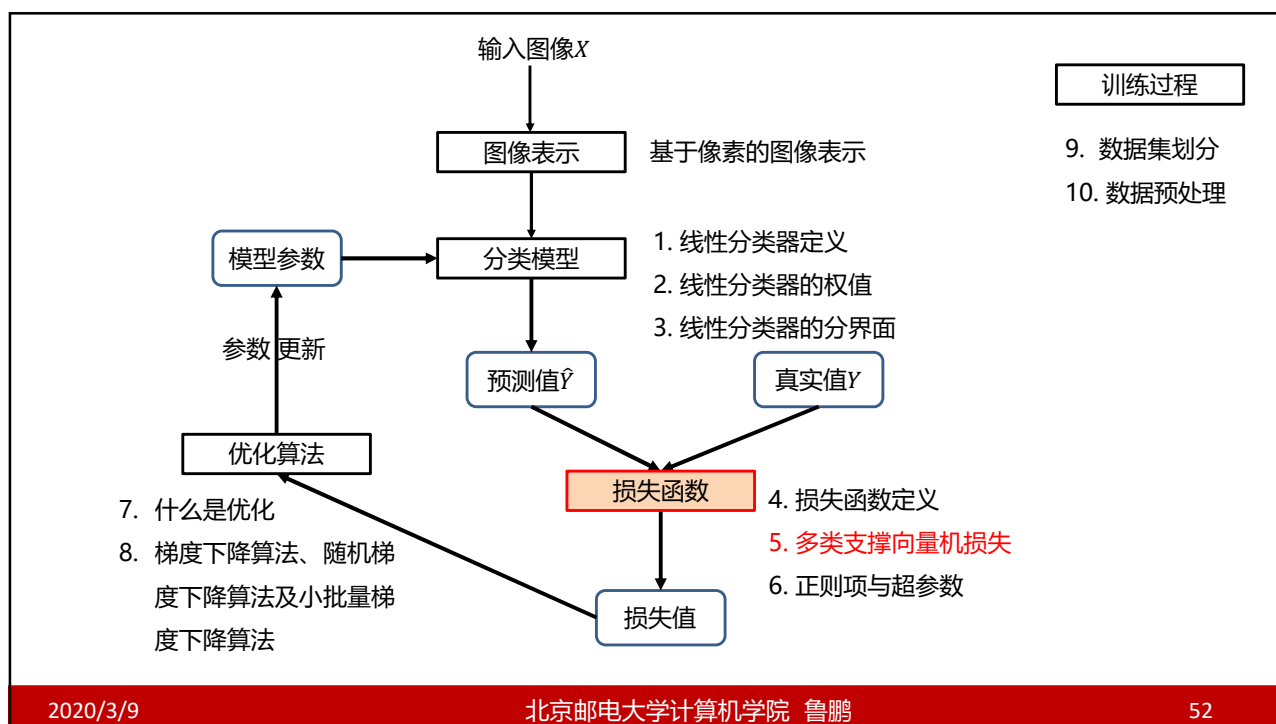
x_i 表示数据集中第 i 张图片;

$f(x_i, W)$ 为分类器对 x_i 的类别预测;

y_i 为样本 i 真实类别标签 (整数);

L_i 为第 i 个样本的损失当预测值;

L 为数据集损失, 它是数据集中所有样本损失的平均。



52

多类支撑向量机损失

2020/3/9

北京邮电大学计算机学院 鲁鹏

53

53

多类支撑向量机损失

$$s_{ij} = f_j(x_i, w_j, b_j) = w_j^T x_i + b_j$$

j —— 类别标签, 取值范围{1, 2, ..., c};

多类支撑向量机损失

$$s_{ij} = f_j(x_i, w_j, b_j) = w_j^T x_i + b_j$$

j —— 类别标签, 取值范围{1, 2, ..., c};

w_j, b_j —— 第 j 个类别分类器的参数;

多类支撑向量机损失

$$s_{ij} = f_j(x_i, w_j, b_j) = w_j^T x_i + b_j$$

j —— 类别标签, 取值范围 $\{1, 2, \dots, c\}$;

w_j, b_j —— 第 j 个类别分类器的参数;

x_i —— 表示数据集中的第 i 个样本

多类支撑向量机损失

$$s_{ij} = f_j(x_i, w_j, b_j) = w_j^T x_i + b_j$$

j —— 类别标签, 取值范围 $\{1, 2, \dots, c\}$;

w_j, b_j —— 第 j 个类别分类器的参数;

x_i —— 表示数据集中的第 i 个样本

s_{ij} —— 第 i 个样本第 j 类别的预测分数

多类支撑向量机损失

$$s_{ij} = f_j(x_i, w_j, b_j) = w_j^T x_i + b_j$$

j —— 类别标签, 取值范围 $\{1, 2, \dots, c\}$;

w_j, b_j —— 第 j 个类别分类器的参数;

x_i —— 表示数据集中的第 i 个样本

s_{ij} —— 第 i 个样本第 j 类别的预测分数

第 i 个样本的多类支撑向量机损失定义如下:

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_{ij} + 1 \\ s_{ij} - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

$$= \sum_{j \neq y_i} \max(0, s_{ij} - s_{y_i} + 1)$$

多类支撑向量机损失

$$s_{ij} = f_j(x_i, w_j, b_j) = w_j^T x_i + b_j$$

j —— 类别标签, 取值范围 $\{1, 2, \dots, c\}$;

w_j, b_j —— 第 j 个类别分类器的参数;

x_i —— 表示数据集中的第 i 个样本

s_{ij} —— 第 i 个样本第 j 类别的预测分数

s_{y_i} —— 第 i 个样本真实类别的预测分数

第 i 个样本的多类支撑向量机损失定义如下:

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_{ij} + 1 \\ s_{ij} - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

$$= \sum_{j \neq y_i} \max(0, s_{ij} - s_{y_i} + 1)$$

多类支撑向量机损失

$$s_{ij} = f_j(x_i, w_j, b_j) = w_j^T x_i + b_j$$

第*i*个样本的多类支撑向量机损失定义如下:

j —— 类别标签, 取值范围 $\{1, 2, \dots, c\}$;

w_j, b_j —— 第*j*个类别分类器的参数;

x_i —— 表示数据集中的第 i 个样本

s_{ij} —— 第*i*个样本第*j*类别的预测分数

s_{y_i} —— 第*i*个样本真实类别的预测分数

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_{ij} + 1 \\ s_{ij} - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

$$= \sum_{j \neq y_i} \max(0, s_{ij} - s_{y_i} + 1)$$

◆ 正确类别的得分比不正确类别的得分高出1分, 就没有损失

2020/3/9

北京邮电大学计算机学院 鲁鹏

60

60

多类支撑向量机损失

$$s_{ij} = f_j(x_i, w_j, b_j) = w_j^T x_i + b_j$$

第*i*个样本的多类支撑向量机损失定义如下:

j —— 类别标签, 取值范围 $\{1, 2, \dots, c\}$;

w_j, b_j —— 第*j*个类别分类器的参数;

x_i —— 表示数据集中的第 i 个样本

s_{ij} —— 第*i*个样本第*j*类别的预测分数

s_{y_i} —— 第*i*个样本真实类别的预测分数

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_{ij} + 1 \\ s_{ij} - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

$$= \sum_{j \neq y_i} \max(0, s_{ij} - s_{y_i} + 1)$$

◆ 正确类别的得分比不正确类别的得分高出1分, 就没有损失

◆ 否则, 就会产生损失

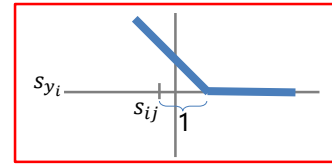
2020/3/9

北京邮电大学计算机学院 鲁鹏

61

61

多类支撑向量机损失



$$s_{ij} = f_j(x_i, w_j, b_j) = w_j^T x_i + b_j$$

第*i*个样本的多类支撑向量机损失定义如下:

j —— 类别标签, 取值范围 $\{1, 2, \dots, c\}$;

w_j, b_j —— 第*j*个类别分类器的参数;

x_i —— 表示数据集中的第 i 个样本

s_{ij} —— 第*i*个样本第*j*类别的预测分数

s_{y_i} —— 第*i*个样本真实类别的预测分数

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_{ij} + 1 \\ s_{ij} - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

$$= \sum_{j \neq y_i} \max(0, s_{ij} - s_{y_i} + 1)$$

◆ 正确类别的得分比不正确类别的得分高出1分, 就没有损失

◆ 否则, 就会产生损失

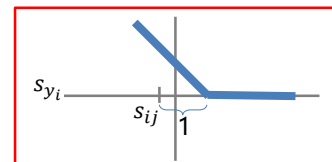
2020/3/9

北京邮电大学计算机学院 鲁鹏

62

62

多类支撑向量机损失



$$s_{ij} = f_j(x_i, w_j, b_j) = w_j^T x_i + b_j$$

第*i*个样本的多类支撑向量机损失定义如下:

j —— 类别标签, 取值范围 $\{1, 2, \dots, c\}$;

w_j, b_j —— 第*j*个类别分类器的参数;

x_i —— 表示数据集中的第 i 个样本

s_{ij} —— 第*i*个样本第*j*类别的预测分数

s_{y_i} —— 第*i*个样本真实类别的预测分数

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_{ij} + 1 \\ s_{ij} - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

$$= \sum_{j \neq y_i} \max(0, s_{ij} - s_{y_i} + 1)$$

◆ 正确类别的得分比不正确类别的得分高出1分, 就没有损失

◆ 否则, 就会产生损失

$\max(0, \cdot)$ 损失——常被称为折页损失 (hinge loss)




2020/3/9

北京邮电大学计算机学院 鲁鹏

63

63

示例：假设有3个类别的训练样本各一张，分类器是线性分类器 $f(x, W) = Wx + b$ ，其中权重 W ， b 已知，分类器对三个样本的打分如下：

	bird	cat	car	loss
	0.6	-2.3	1.9	
	1.7	2.9	2.3	
	3.1	-2.6	4.3	




2020/3/9

北京邮电大学计算机学院 鲁鹏

64

64

示例：假设有3个类别的训练样本各一张，分类器是线性分类器 $f(x, W) = Wx + b$ ，其中权重 W ， b 已知，分类器对三个样本的打分如下：

	bird	cat	car	loss
	0.6	-2.3	1.9	2.3
	1.7	2.9	2.3	0.4
	3.1	-2.6	4.3	0

当前分类器对于鸟这张图像的损失：

$$\begin{aligned}
 L_i &= \sum_{j \neq y_i} \max(0, s_{ij} - s_{y_i} + 1) \\
 &= \max(0, -2.3 - 0.6 + 1) \\
 &\quad + \max(0, 1.9 - 0.6 + 1) \\
 &= \max(0, -1.9) + \max(0, 2.3) \\
 &= 0 + 2.3 \\
 &= 2.3
 \end{aligned}$$




2020/3/9

北京邮电大学计算机学院 鲁鹏

65

65

示例：假设有3个类别的训练样本各一张，分类器是线性分类器 $f(x, W) = Wx + b$ ，其中权重 W ， b 已知，分类器对三个样本的打分如下：

	bird	cat	car	loss
	0.6	-2.3	1.9	2.3
	1.7	2.9	2.3	0.4
	3.1	-2.6	4.3	0

当前分类器对于鸟这张图像的损失：

$$\begin{aligned}
 L_i &= \sum_{j \neq y_i} \max(0, s_{ij} - s_{y_i} + 1) \\
 &= \max(0, 1.7 - 2.9 + 1) \\
 &\quad + \max(0, 2.3 - 2.9 + 1) \\
 &= \max(0, -0.2) + \max(0, 0.4) \\
 &= 0 + 0.4 \\
 &= 0.4
 \end{aligned}$$




2020/3/9

北京邮电大学计算机学院 鲁鹏

66

66

示例：假设有3个类别的训练样本各一张，分类器是线性分类器 $f(x, W) = Wx + b$ ，其中权重 W ， b 已知，分类器对三个样本的打分如下：

	bird	cat	car	loss
	0.6	-2.3	1.9	2.3
	1.7	2.9	2.3	0.4
	3.1	-2.6	4.3	0

当前分类器对于鸟这张图像的损失：

$$\begin{aligned}
 L_i &= \sum_{j \neq y_i} \max(0, s_{ij} - s_{y_i} + 1) \\
 &= \max(0, 3.1 - 4.3 + 1) \\
 &\quad + \max(0, -2.6 - 4.3 + 1) \\
 &= 0
 \end{aligned}$$




2020/3/9

北京邮电大学计算机学院 鲁鹏

67

67

示例：假设有3个类别的训练样本各一张，分类器是线性分类器 $f(x, W) = Wx + b$ ，其中权重 W ， b 已知，分类器对三个样本的打分如下：

	bird	cat	car	loss
	0.6	-2.3	1.9	2.3
	1.7	2.9	2.3	0.4
	3.1	-2.6	4.3	0

当前分类器对于整个数据集图像的损失：

$$L = \frac{1}{N} \sum_{i=1}^N L_i \quad (\text{整个数据集损失的平均值})$$

$$L = (2.3 + 0.4 + 0)/3 = 0.9$$

问题抢答

➤ 损失函数：

$$L = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i)$$

➤ 单样本的多类支撑向量机损失：

$$L_i = \sum_{j \neq y_i} \max(0, s_{ij} - s_{y_i} + 1)$$

➤ 线性分类器：

$$s_{ij} = w_j^T x_i + b_j$$

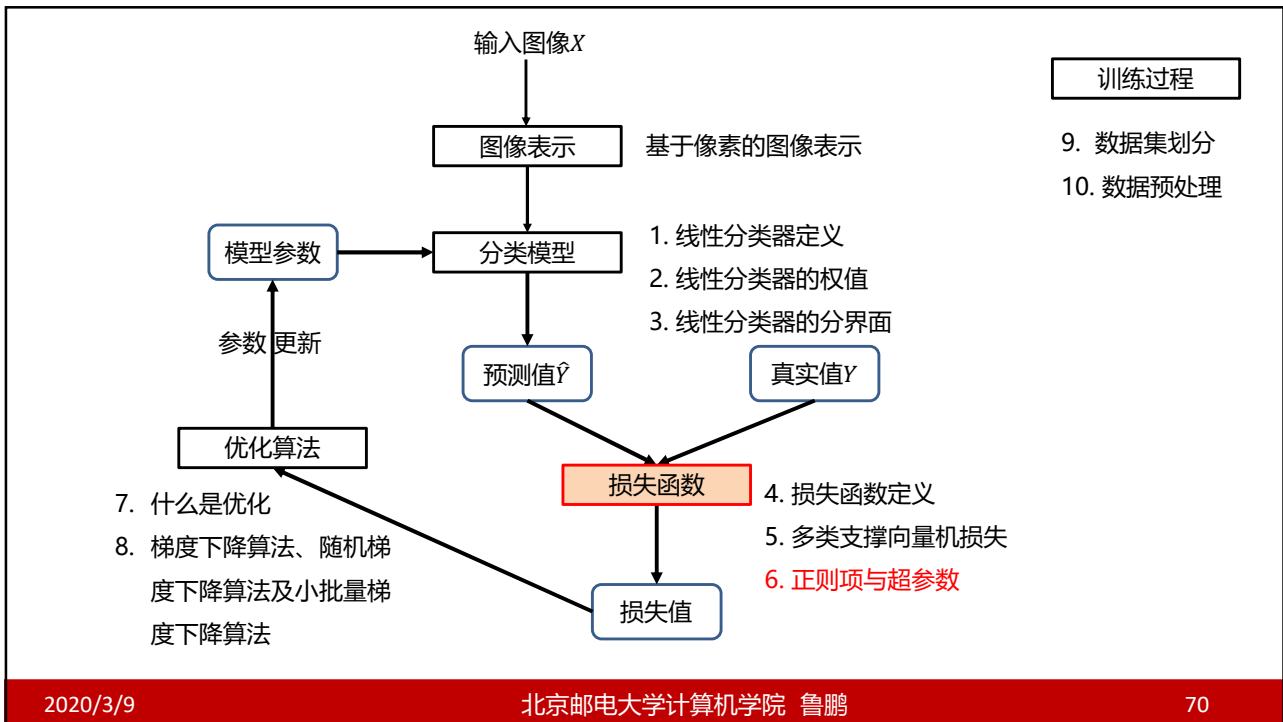
1: 多类支撑向量机损失 L_i 的最大/最小值会是多少？

2: 如果初始化时 w 和 b 很小，损失 L 会是多少？

3: 考虑所有类别 (包括 $j = y_i$)，损失 L_i 会有什么变化？

4: 在总损失 L 计算时，如果用求和代替平均？

5: 如果使用 $L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)^2$



70

再谈损失函数

$$L(W) = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i)$$

问题：假设存在一个W使损失函数L=0，这个W是唯一的吗？

71

再谈损失函数

示例：假设两个线性器分类 $f_1(x, W_1) = W_1 x$, $f_2(x, W_2) = W_2 x$, 其中, $W_2 = 2 W_1$ 。对于下面图像, 已知分类器1的打分结果 (如表所示) , 请计算分类器2的打分结果, 以及两个分类器对当前样本的多类支撑向量机损失:



	鸟	猫	汽车	损失
分类器1	3.1	-2.6	4.3	?
分类器2	?	?	?	?

第*i*个样本的多类支撑向量机损失:

$$L_i = \sum_{j \neq y_i} \max(0, s_{ij} - s_{y_i} + 1)$$

2020/3/9

北京邮电大学计算机学院 鲁鹏

72

72

再谈损失函数

示例：假设两个线性器分类 $f_1(x, W_1) = W_1 x$, $f_2(x, W_2) = W_2 x$, 其中, $W_2 = 2 W_1$ 。对于下面图像, 已知分类器1的打分结果 (如表所示) , 请计算分类器2的打分结果, 以及两个分类器对当前样本的多类支撑向量机损失:



	鸟	猫	汽车	损失
分类器1	3.1	-2.6	4.3	?
分类器2	6.2	-5.2	8.6	?

第*i*个样本的多类支撑向量机损失:

$$L_i = \sum_{j \neq y_i} \max(0, s_{ij} - s_{y_i} + 1)$$

2020/3/9

北京邮电大学计算机学院 鲁鹏

73

73

再谈损失函数

示例：假设两个线性器分类 $f_1(x, W_1) = W_1 x$, $f_2(x, W_2) = W_2 x$, 其中, $W_2 = 2 W_1$ 。对于下面图像, 已知分类器1的打分结果 (如表所示), 请计算分类器2的打分结果, 以及两个分类器对当前样本的多类支撑向量机损失:



	鸟	猫	汽车	损失
分类器1	3.1	-2.6	4.3	0
分类器2	6.2	-5.2	8.6	?

第 i 个样本的多类支撑向量机损失:

$$L_i = \sum_{j \neq y_i} \max(0, s_{ij} - s_{y_i} + 1)$$

分类器1损失:

$$\begin{aligned} &= \max(0, 3.1 - 4.3 + 1) \\ &+ \max(0, -2.6 - 4.3 + 1) \\ &= \max(0, -0.2) + \max(0, -5.9) \\ &= 0 \end{aligned}$$

2020/3/9

北京邮电大学计算机学院 鲁鹏

74

74

再谈损失函数

示例：假设两个线性器分类 $f_1(x, W_1) = W_1 x$, $f_2(x, W_2) = W_2 x$, 其中, $W_2 = 2 W_1$ 。对于下面图像, 已知分类器1的打分结果 (如表所示), 请计算分类器2的打分结果, 以及两个分类器对当前样本的多类支撑向量机损失:



	鸟	猫	汽车	损失
分类器1	3.1	-2.6	4.3	0
分类器2	6.2	-5.2	8.6	0

第 i 个样本的多类支撑向量机损失:

$$L_i = \sum_{j \neq y_i} \max(0, s_{ij} - s_{y_i} + 1)$$

分类器1损失:

$$\begin{aligned} &= \max(0, 3.1 - 4.3 + 1) \\ &+ \max(0, -2.6 - 4.3 + 1) \\ &= \max(0, -0.2) + \max(0, -5.9) \\ &= 0 \end{aligned}$$

分类器2损失:

$$\begin{aligned} &= \max(0, 6.2 - 8.6 + 1) \\ &+ \max(0, -5.2 - 8.6 + 1) \\ &= \max(0, -1.4) + \max(0, -12.8) \\ &= 0 \end{aligned}$$

2020/3/9

北京邮电大学计算机学院 鲁鹏

75

75

再谈损失函数

$$L(W) = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i)$$

问题：假设存在一个W使损失函数L=0，这个W是唯一的吗？

答：不唯一，因为 W_2 同样有 $L = 0$

应该如何在 W_1 和 W_2 之间做出选择？

正则项

$$L(W) = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i) + \lambda R(W)$$

数据损失

$$L(W) = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i) + \lambda R(W)$$

数据损失: 模型预测需
要和训练集相匹配

正则项损失

$$L(W) = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i) + \lambda R(W)$$

数据损失: 模型预测需
要和训练集相匹配

正则损失: 防止模型在训
练集上学习得“太好”。

正则项损失

$$L(W) = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i) + \lambda R(W)$$

数据损失: 模型预测需
要和训练集相匹配

正则损失: 防止模型在训
练集上学习得“太好”。

◆ $R(W)$ 是一个与权值有
关, 跟图像数据无关的
函数

2020/3/9

北京邮电大学计算机学院 鲁鹏

80

80

正则项损失

$$L(W) = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i) + \lambda R(W)$$

数据损失: 模型预测需
要和训练集相匹配

正则损失: 防止模型在训
练集上学习得“太好”。

◆ $R(W)$ 是一个与权值有
关, 跟图像数据无关的
函数

◆ λ 是一个超参数控制着
正则损失在总损失中所
占的比重

2020/3/9

北京邮电大学计算机学院 鲁鹏

81

81

什么是超参数?

- 在开始学习过程之前设置值的参数，而不是学习得到。
- 超参数一般都会对模型性能有着重要的影响。

$$L(W) = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i) + \lambda R(W)$$

数据损失: 模型预测需
要和训练集相匹配

正则损失: 防止模型在训
练集上学习得“太好”。

- $\lambda = 0$ 优化结果仅与数据损失相关
- $\lambda = \infty$ 优化结果与数据损失无关，仅考虑权重损失。此时，系统最优解为 $W = 0$ 。

如何设置一个好的超参数呢，这个问题我们会在后面的课程中专门探讨。

2020/3/9

北京邮电大学计算机学院 鲁鹏

82

82

L2正则项

$$L(W) = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i) + \lambda R(W)$$

数据损失: 模型预测需
要和训练集相匹配

正则损失: 防止
模型在训练集上
学习得“太好”。

L2正则项 $R(W) = \sum_k \sum_l W_{k,l}^2$

2020/3/9

北京邮电大学计算机学院 鲁鹏

83

83

L2正则项

L2损失示例

$$L(W) = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i) + \lambda R(W)$$

数据损失: 模型预测需
要和训练集相匹配

正则损失: 防止
模型在训练集上
学习得“太好”。

L2正则项 $R(W) = \sum_k \sum_l W_{k,l}^2$

样本: $x = [1,1,1,1]$

分类器1: $w_1 = [1,0,0,0]$

分类器2: $w_2 = [0.25,0.25,0.25,0.25]$

2020/3/9

北京邮电大学计算机学院 鲁鹏

84

84

L2正则项

L2损失示例

$$L(W) = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i) + \lambda R(W)$$

数据损失: 模型预测需
要和训练集相匹配

正则损失: 防止
模型在训练集上
学习得“太好”。

L2正则项 $R(W) = \sum_k \sum_l W_{k,l}^2$

样本: $x = [1,1,1,1]$

分类器1: $w_1 = [1,0,0,0]$

分类器2: $w_2 = [0.25,0.25,0.25,0.25]$

2020/3/9

北京邮电大学计算机学院 鲁鹏

85

85

L2正则项

L2损失示例

$$L(W) = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i) + \lambda R(W)$$

数据损失: 模型预测需
要和训练集相匹配

正则损失: 防止
模型在训练集上
学习得“太好”。

L2正则项 $R(W) = \sum_k \sum_l W_{k,l}^2$

样本: $x = [1,1,1,1]$

分类器1: $w_1 = [1,0,0,0]$

分类器2: $w_2 = [0.25,0.25,0.25,0.25]$

分类器输出: $w_1^T x = w_2^T x = 1$

2020/3/9

北京邮电大学计算机学院 鲁鹏

86

86

L2正则项

L2损失示例

$$L(W) = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i) + \lambda R(W)$$

数据损失: 模型预测需
要和训练集相匹配

正则损失: 防止
模型在训练集上
学习得“太好”。

L2正则项 $R(W) = \sum_k \sum_l W_{k,l}^2$

样本: $x = [1,1,1,1]$

分类器1: $w_1 = [1,0,0,0]$

分类器2: $w_2 = [0.25,0.25,0.25,0.25]$

分类器输出: $w_1^T x = w_2^T x = 1$

正则损失:
 $R(w_1) = 1$ $R(w_2) = 0.25$

2020/3/9

北京邮电大学计算机学院 鲁鹏

87

87

L2正则项

L2损失示例

$$L(W) = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i) + \lambda R(W)$$

数据损失: 模型预测需
要和训练集相匹配

正则损失: 防止
模型在训练集上
学习得“太好”。

L2正则项 $R(W) = \sum_k \sum_l W_{k,l}^2$

样本: $x = [1,1,1,1]$

分类器1: $w_1 = [1,0,0,0]$

分类器2: $w_2 = [0.25,0.25,0.25,0.25]$

分类器输出: $w_1^T x = w_2^T x = 1$

正则损失:

$$R(w_1) = 1 \quad R(w_2) = 0.25$$

w_2 总损失小

2020/3/9

北京邮电大学计算机学院 鲁鹏

88

88

L2正则项

L2损失示例

$$L(W) = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i) + \lambda R(W)$$

数据损失: 模型预测需
要和训练集相匹配

正则损失: 防止
模型在训练集上
学习得“太好”。

L2正则项 $R(W) = \sum_k \sum_l W_{k,l}^2$

样本: $x = [1,1,1,1]$

分类器1: $w_1 = [1,0,0,0]$

分类器2: $w_2 = [0.25,0.25,0.25,0.25]$

分类器输出: $w_1^T x = w_2^T x = 1$

正则损失:

$$R(w_1) = 1 \quad R(w_2) = 0.25$$

w_2 总损失小

2020/3/9

北京邮电大学计算机学院 鲁鹏

89

89

L2正则项

L2损失示例

$$L(W) = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i) + \lambda R(W)$$

数据损失: 模型预测需
要和训练集相匹配

正则损失: 防止
模型在训练集上
学习得“太好”。

$$\text{L2正则项} \quad R(W) = \sum_k \sum_l W_{k,l}^2$$

L2正则损失对大数值权值进行惩罚，喜欢分散权值，鼓励分类器将所有维度的特征都用起来，而不是强烈的依赖其中少数几维特征

样本: $x = [1,1,1,1]$

分类器1: $w_1 = [1,0,0,0]$

分类器2: $w_2 = [0.25,0.25,0.25,0.25]$

分类器输出: $w_1^T x = w_2^T x = 1$

正则损失:

$$R(w_1) = 1 \quad R(w_2) = 0.25$$

w_2 总损失小

2020/3/9

北京邮电大学计算机学院 鲁鹏

90

90

L2正则项

L2损失示例

$$L(W) = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i) + \lambda R(W)$$

数据损失: 模型预测需
要和训练集相匹配

正则损失: 防止
模型在训练集上
学习得“太好”。

$$\text{L2正则项} \quad R(W) = \sum_k \sum_l W_{k,l}^2$$

L2正则损失对大数值权值进行惩罚，喜欢分散权值，鼓励分类器将所有维度的特征都用起来，而不是强烈的依赖其中少数几维特征

样本: $x = [1,1,1,1]$

分类器1: $w_1 = [1,0,0,0]$

分类器2: $w_2 = [0.25,0.25,0.25,0.25]$

分类器输出: $w_1^T x = w_2^T x = 1$

正则损失:

$$R(w_1) = 1 \quad R(w_2) = 0.25$$

w_2 总损失小

正则项让模型有了偏好!!!

2020/3/9

北京邮电大学计算机学院 鲁鹏

91

91

常用的正则项损失

$$L(W) = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i) + \lambda R(W)$$

数据损失: 模型预测需要和训练集相匹配
正则损失: 防止模型在训练集上学习得“太好”。

$$\text{L1正则项: } R(W) = \sum_k \sum_l |W_{k,l}|$$

$$\text{L2正则项: } R(W) = \sum_k \sum_l W_{k,l}^2$$

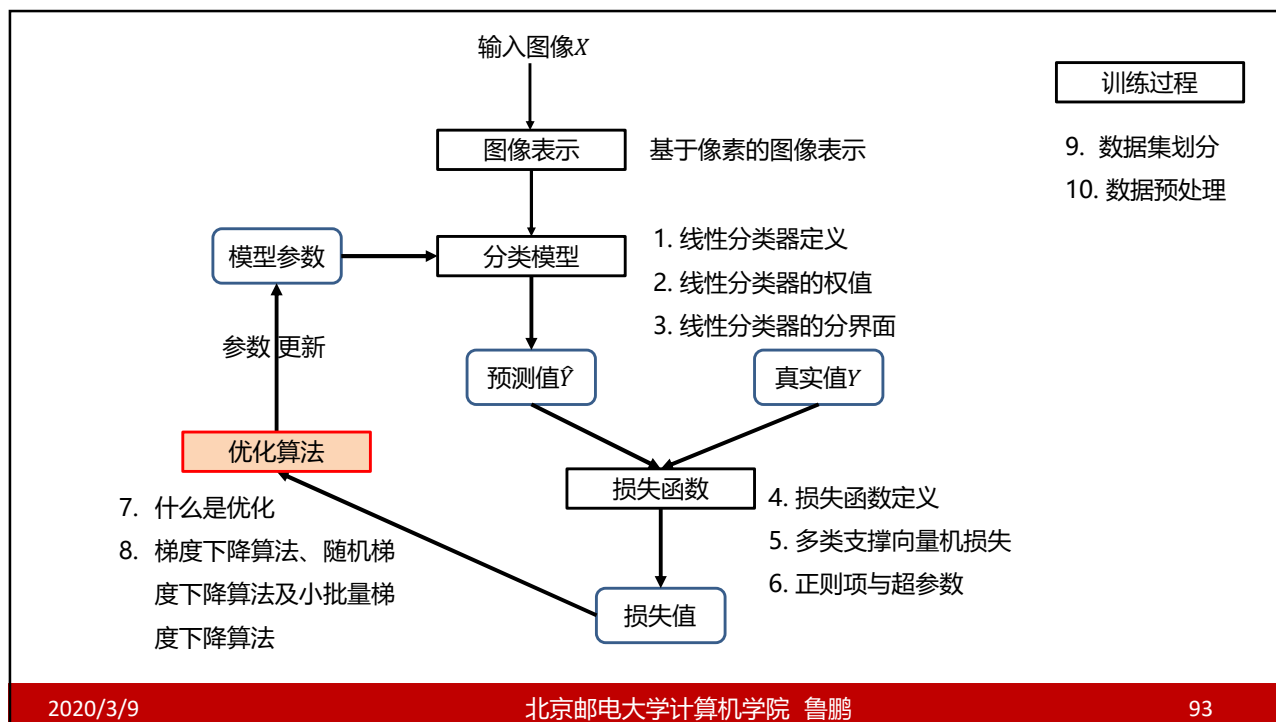
$$\text{Elastic net(L1+L2): } R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$$

2020/3/9

北京邮电大学计算机学院 鲁鹏

92

92

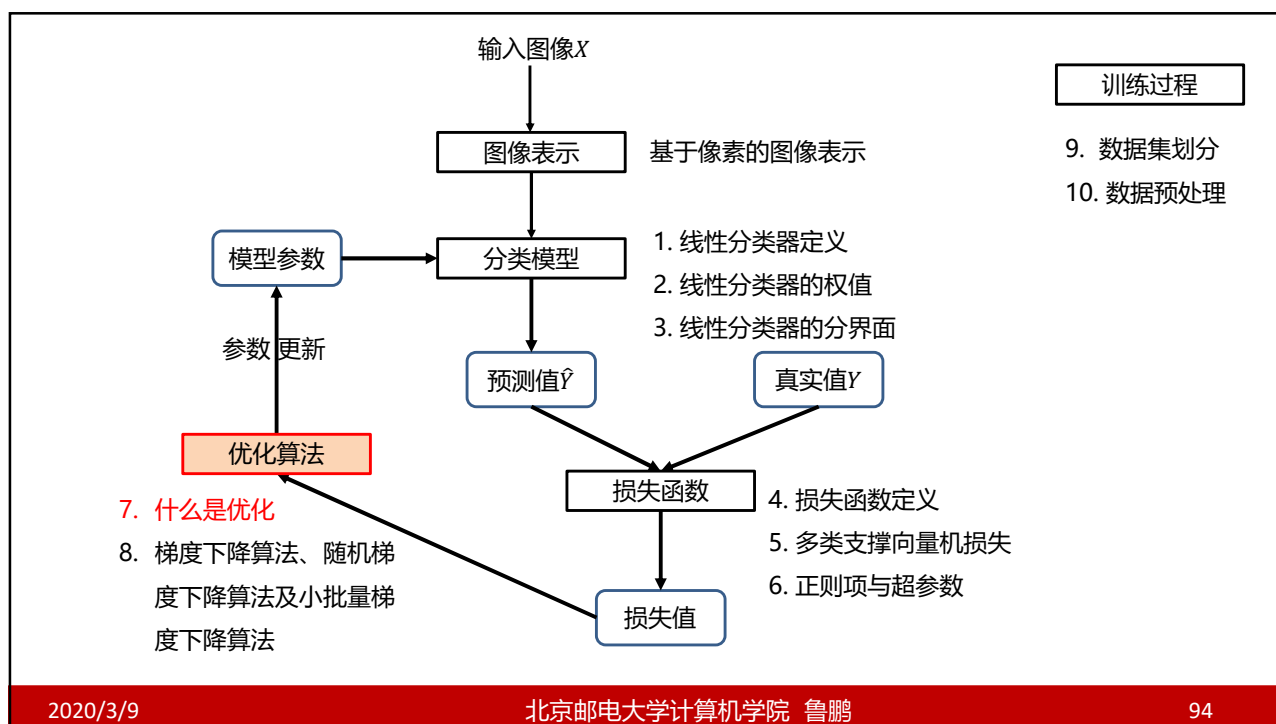


2020/3/9

北京邮电大学计算机学院 鲁鹏

93

93



94

什么是参数优化?

参数优化是机器学习的核心步骤之一，它**利用损失函数的输出值作为反馈信号来调整分类器参数**，以提升分类器对训练样本的**预测性能**。

2020/3/9

北京邮电大学计算机学院 鲁鹏

95

95

优化算法目标

$$\text{损失函数 } L = \frac{1}{N} \sum_{i=1}^N L_i + \lambda R(W)$$

损失函数 L 是一个与参数 W 有关的函数，优化的目标就是找到使损失函数 L 达到最优的那组参数 W 。

直接方法：

$$\frac{\partial L}{\partial W} = 0$$

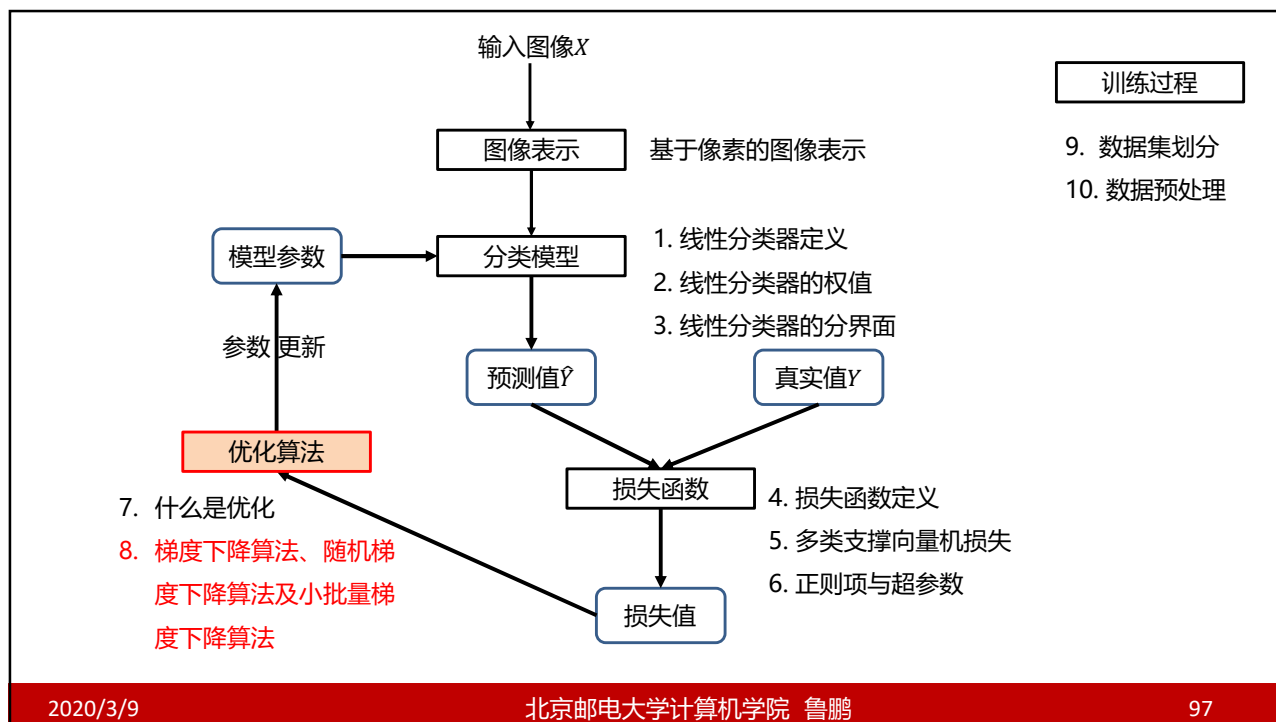
通常， L 形式比较复杂，很难从这个等式直接求解出 W ！

2020/3/9

北京邮电大学计算机学院 鲁鹏

96

96



2020/3/9

北京邮电大学计算机学院 鲁鹏

97

97

梯度下降算法

一种简单而高效的迭代优化方法!

2020/3/9

北京邮电大学计算机学院 鲁鹏

98

98

梯度下降算法

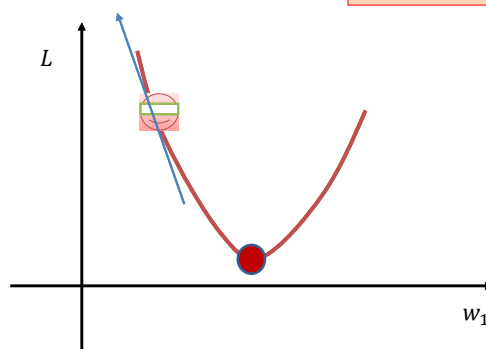
$$L(W) = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i) + \lambda R(W)$$

➤ 往哪儿走?

答: 负梯度方法

➤ 走多远?

答: 步长来决定



2020/3/9

北京邮电大学计算机学院 鲁鹏

99

99

梯度下降算法

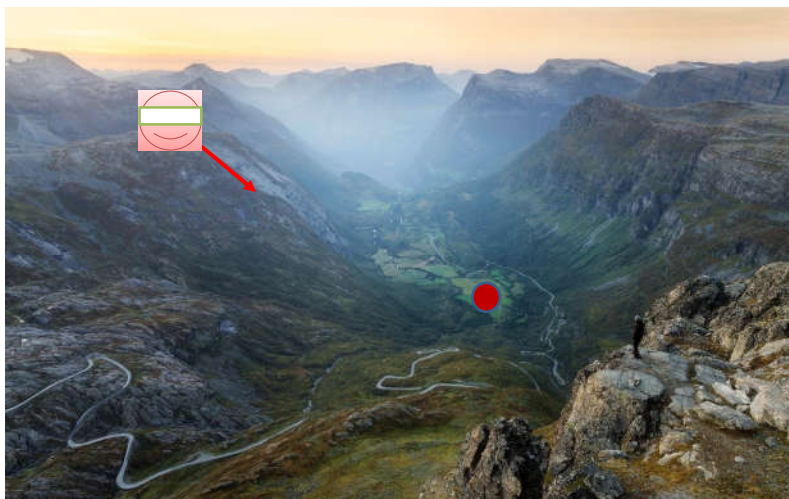
$$L(W) = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i) + \lambda R(W)$$

➤ 往哪儿走?

答: 负梯度方法

➤ 走多远?

答: 步长来决定



2020/3/9

北京邮电大学计算机学院 鲁鹏

100

100

梯度下降算法

梯度下降: 利用所有样本计算损失并更新梯度

```
while True
```

```
  权值的梯度 ← 计算梯度(损失, 训练样本, 权值)
```

```
  权值 ← 权值 - 学习率 * 权值的梯度
```

2020/3/9

北京邮电大学计算机学院 鲁鹏

101

101

梯度下降算法

梯度下降: 利用所有样本计算损失并更新梯度

```
while True
    权值的梯度 ← 计算梯度(损失, 训练样本, 权值)
    权值 ← 权值 - 学习率 * 权值的梯度
```

2020/3/9

北京邮电大学计算机学院 鲁鹏

102

102

梯度下降算法

梯度下降: 利用所有样本计算损失并更新梯度

```
while True
    权值的梯度 ← 计算梯度(损失, 训练样本, 权值)
    权值 ← 权值 - 学习率 * 权值的梯度
```

2020/3/9

北京邮电大学计算机学院 鲁鹏

103

103

梯度下降算法

梯度下降: 利用所有样本计算损失并更新梯度

```
while True
  权值的梯度 ← 计算梯度(损失, 训练样本, 权值)
  权值 ← 权值 - 学习率 * 权值的梯度
```

更新后
的权值

当前的
权值

2020/3/9

北京邮电大学计算机学院 鲁鹏

104

104

梯度下降算法

梯度下降: 利用所有样本计算损失并更新梯度

```
while True
  权值的梯度 ← 计算梯度(损失, 训练样本, 权值)
  权值 ← 权值 - 学习率 * 权值的梯度
```

更新后
的权值

当前的
权值

2020/3/9

北京邮电大学计算机学院 鲁鹏

105

105

梯度下降算法

梯度下降: 利用所有样本计算损失并更新梯度

```
while True
    权值的梯度 ← 计算梯度(损失, 训练样本, 权值) ?
    权值 ← 权值 - 学习率 * 权值的梯度
```

2020/3/9

北京邮电大学计算机学院 鲁鹏

106

106

梯度计算

1.数值法 计算量大, 不精确!

一维变量, 函数求导:

$$\frac{dL(w)}{dw} = \lim_{h \rightarrow 0} \frac{L(w+h) - L(w)}{h}$$

示例: 损失函数 $L(w) = w^2$ 求 $w=1$ 点处的梯度

$$\frac{dL(w)}{dw} = \lim_{h \rightarrow 0} \frac{L(w+h) - L(w)}{h} \approx \frac{L(1+0.0001) - L(1)}{0.0001} = 2.0001$$

2020/3/9

北京邮电大学计算机学院 鲁鹏

107

107

梯度计算

2.解析法 精确，速度快，导数函数推导易错！

示例：损失函数 $L(w) = w^2$ 求 $w=1$ 点处的梯度

$$\nabla L(w) = 2w$$

$$\nabla_{w=1} L(w) = 2$$



2020/3/9

北京邮电大学计算机学院 鲁鹏

108

108

梯度计算总结

- 数值梯度: 近似, 慢, 易写
- 解析梯度: 精确, 快, 易错

数值梯度有什么作用?

答：求梯度时一般使用解析梯度，而数值梯度主要用于解析梯度的正确性校验（梯度检查）。

2020/3/9

北京邮电大学计算机学院 鲁鹏

109

109

作业：梯度计算

如何计算多类支撑向量机损失的导数函数？

$$L_i = \sum_{j \neq y_i} \max(0, s_{ij} - s_{y_i} + 1)$$

$$s_{ij} = w_j^T x_i + b_j$$

$$L_i = \sum_{j \neq y_i} \max(0, w_j^T x_i + b_j - (w_{y_i}^T x_i + b_{y_i}) + 1)$$

2020/3/9

北京邮电大学计算机学院 鲁鹏

110

110

梯度下降算法的计算效率

梯度下降：利用所有样本计算损失并更新梯度

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W) + \lambda R(W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

```
while True
```

```
  权值的梯度 ← 计算梯度(损失, 训练样本, 权值)
```

```
  权值 ← 权值 - 学习率 * 权值的梯度
```



当N很大时，权值的梯度计算量很大！

2020/3/9

北京邮电大学计算机学院 鲁鹏

111

111

随机梯度下降算法

随机梯度下降: 每次随机选择一个样本 x_i , 计算损失并更新梯度

$$L(W) = L_i(x_i, y_i, W) + \lambda R(W)$$

$$\nabla_W L(W) = \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

while True

数据 \leftarrow 从训练数据采样(训练数据, 1)

权值的梯度 \leftarrow 计算梯度(损失, 数据, 权值)

权值 \leftarrow 权值 - 学习率 * 权值的梯度



单个样本的训练可能会带来很多噪声, 不是每次迭代都向着整体最优化方向,

2020/3/9

北京邮电大学计算机学院 鲁鹏

112

112

小批量梯度下降算法

小批量随机梯度下降: 每次随机选择 m (批量的大小) 个样本, 计算损失并更新梯度

超参数

$$L(W) = \frac{1}{m} \sum_{i=1}^m L_i(x_i, y_i, W) + \lambda R(W)$$

$$\nabla_W L(W) = \frac{1}{m} \sum_{i=1}^m \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

while True

数据 \leftarrow 从训练数据采样(训练数据, 批量大小)

权值的梯度 \leftarrow 计算梯度(损失, 数据, 权值)

权值 \leftarrow 权值 - 学习率 * 权值的梯度

- > iteration: 表示1次迭代, 每次迭代更新1次网络结构的参数;
- > batch-size: 1次迭代所使用的样本量;
- > epoch: 1个epoch表示过了1遍训练集中的所有样本。



tip: 通常使用2的幂数作为批量大小, 比如每次选取32或64或128个样本

2020/3/9

北京邮电大学计算机学院 鲁鹏

113

113

总结

梯度下降方法

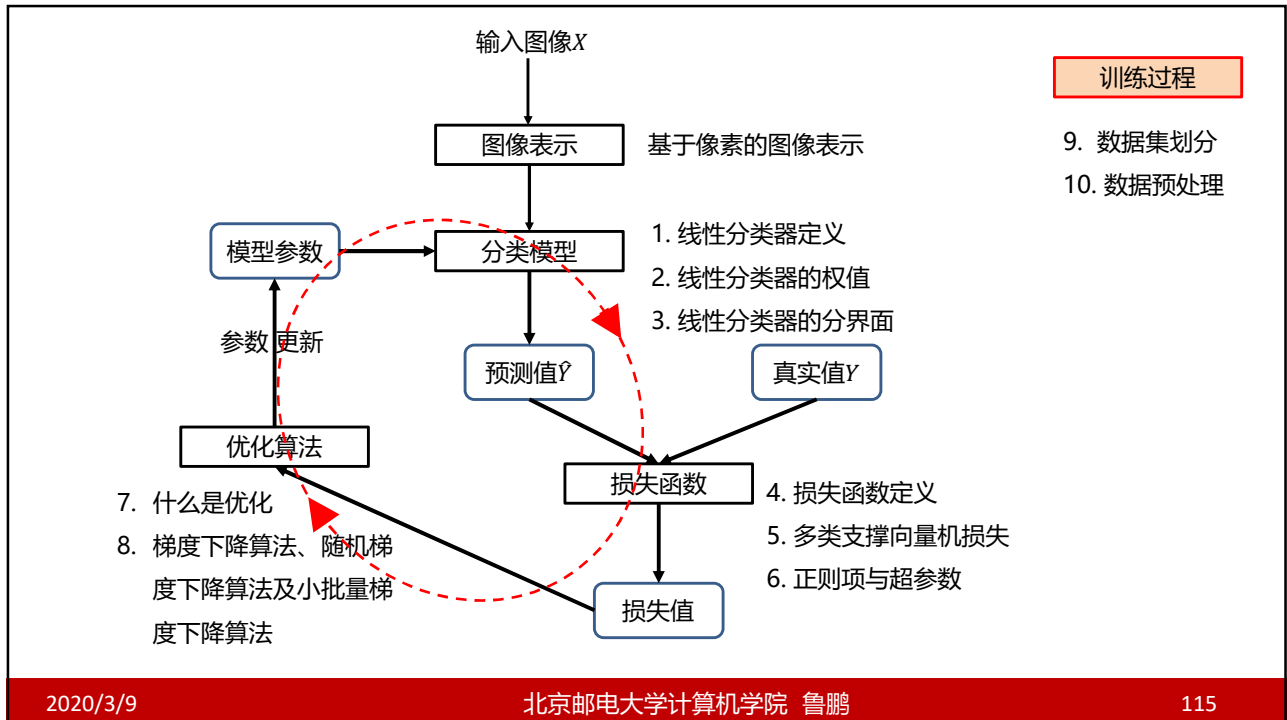
```
while True
    权值的梯度 ← 计算梯度(损失, 训练数据, 权值)
    权值 ← 权值 - 学习率 * 权值的梯度
```

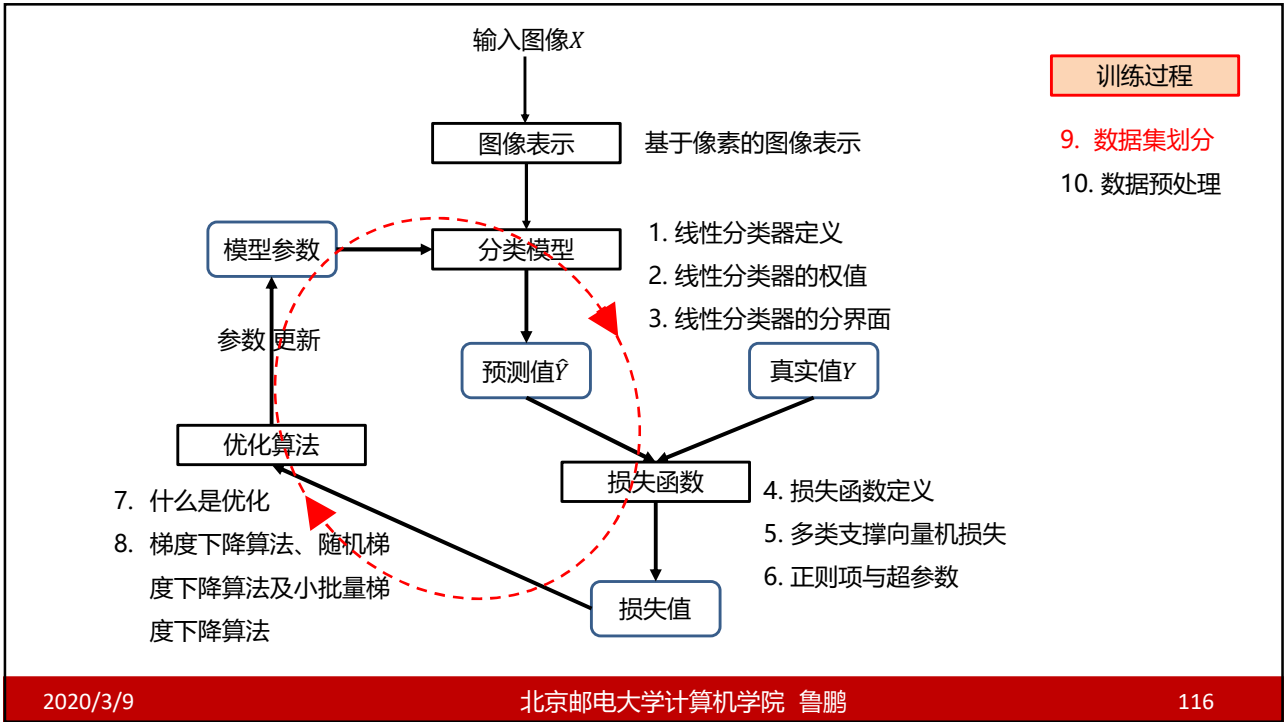
随机梯度下降算法

```
while True
    数据 ← 从训练数据采样(训练数据, 1)
    权值的梯度 ← 计算梯度(损失, 数据, 权值)
    权值 ← 权值 - 学习率 * 权值的梯度
```

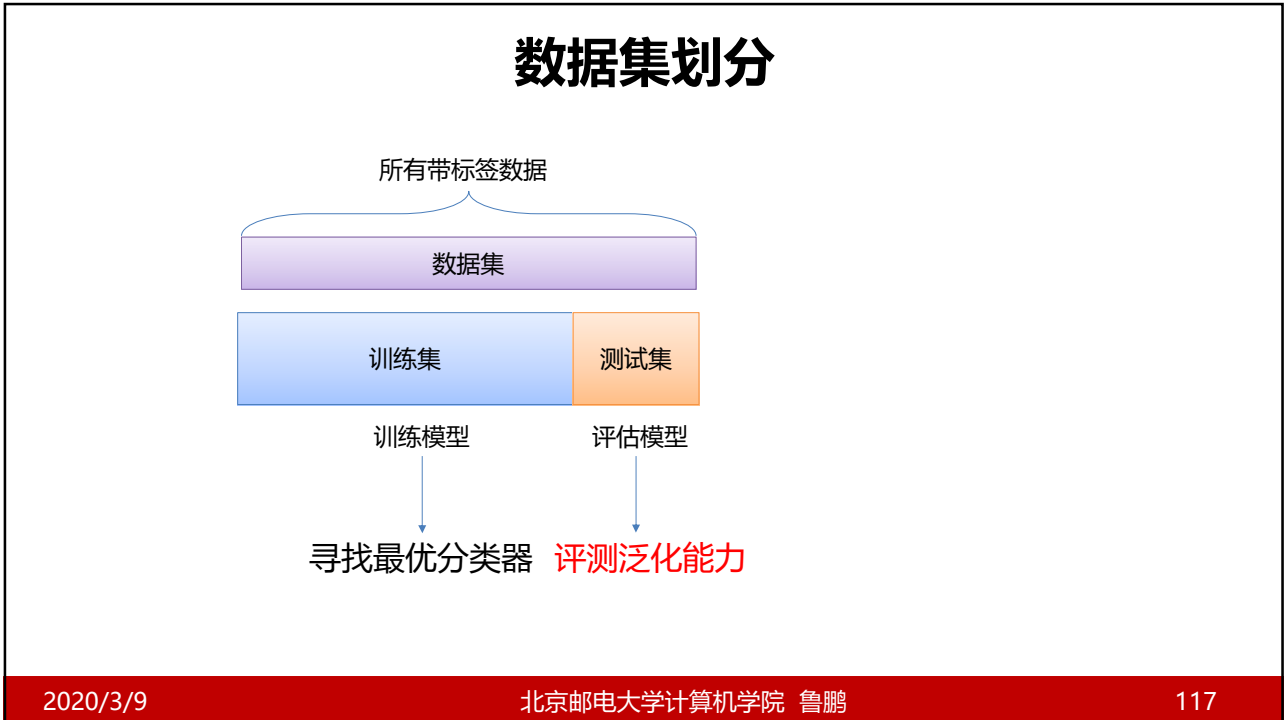
小批量梯度下降算法

```
while True
    小批量数据 ← 从训练数据采样(训练数据, 批量大小)
    权值的梯度 ← 计算梯度(损失, 小批量数据, 权值)
    权值 ← 权值 - 学习率 * 权值的梯度
```



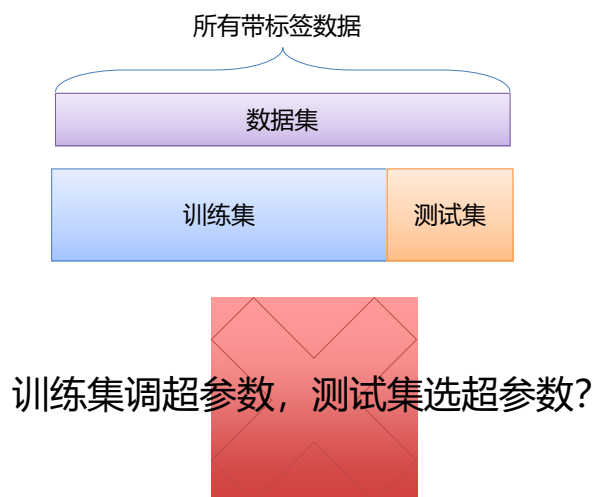


116



117

数据集划分



问题：如果模型含有超参数（比如正则化强度），如何找到泛化能力最好的超参数？

回答：使用验证集

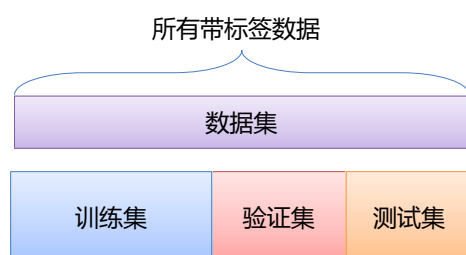
2020/3/9

北京邮电大学计算机学院 鲁鹏

118

118

数据集划分



- 训练集用于给定的超参数时分类器参数的学习；
- 验证集用于选择超参数；
- 测试集评估泛化能力；

2020/3/9

北京邮电大学计算机学院 鲁鹏

119

119

K折交叉验证

问题：如果数据很少，那么可能验证集包含的样本就太少，从而无法在统计上代表数据。

这个问题很容易发现：如果在划分数据前进行不同的随机打乱，最终得到的模型性能差别很大，那么就存在这个问题。

接下来会介绍K折验证与重复的K折验证，它们是解决这一问题的两种方法。

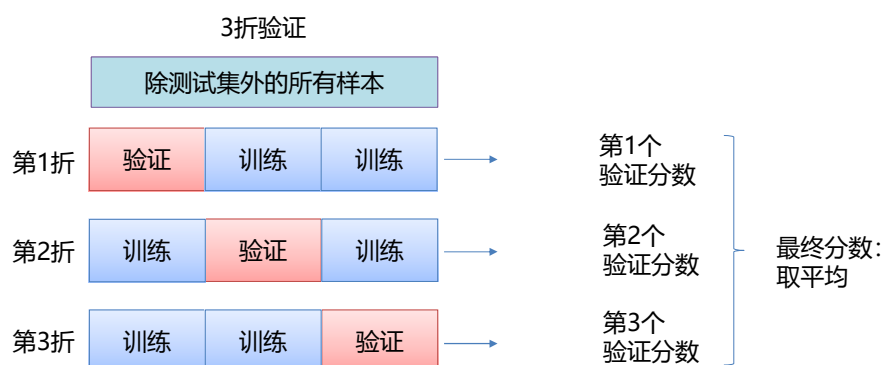
2020/3/9

北京邮电大学计算机学院 鲁鹏

120

120

K折交叉验证



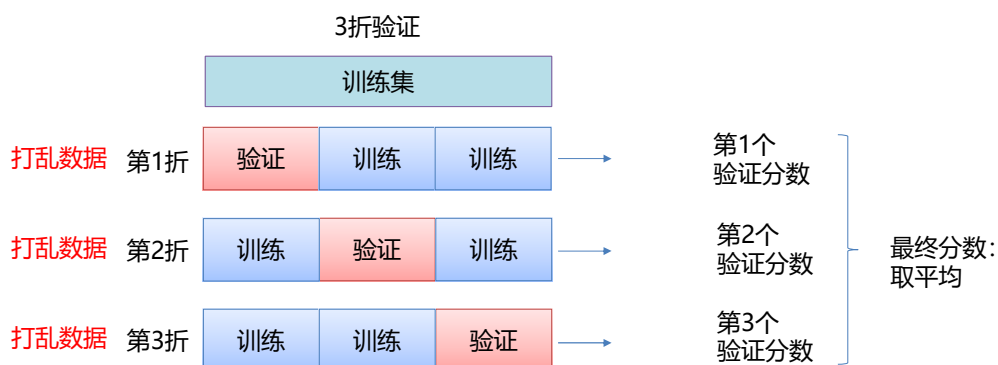
2020/3/9

北京邮电大学计算机学院 鲁鹏

121

121

带有打乱数据的重复K折验证

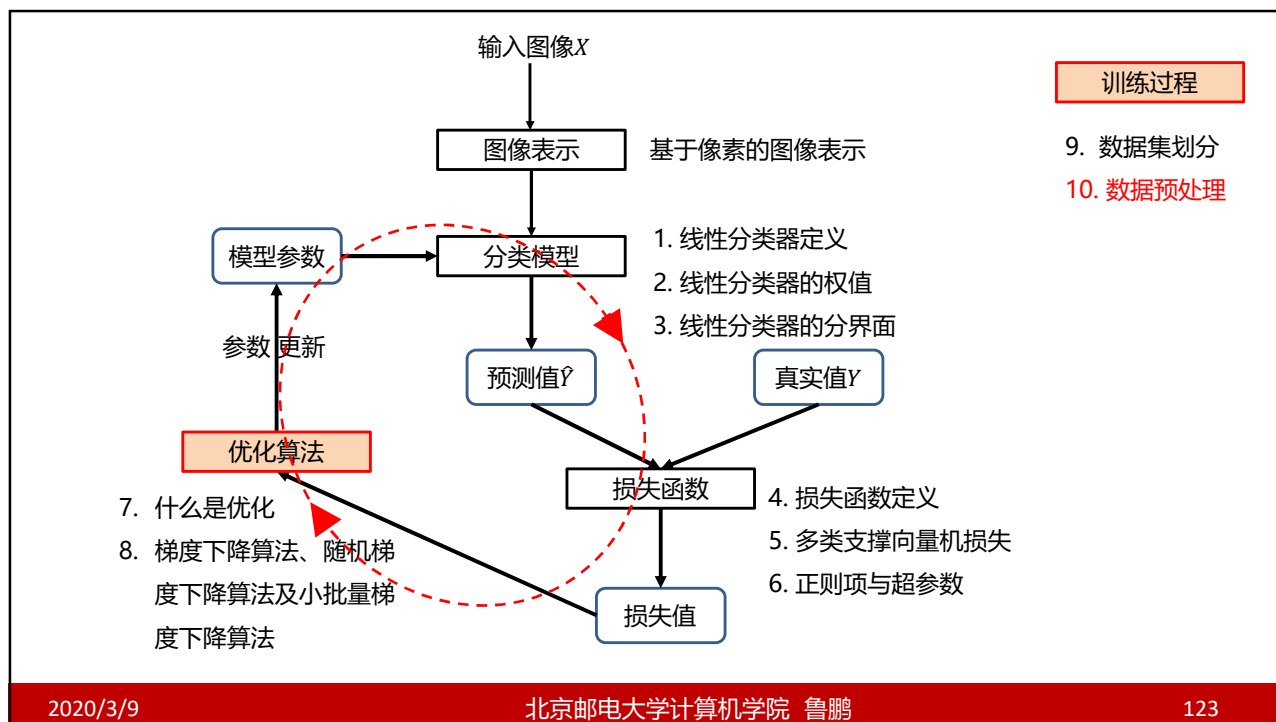


2020/3/9

北京邮电大学计算机学院 鲁鹏

122

122



2020/3/9

北京邮电大学计算机学院 鲁鹏

123

123

数据预处理

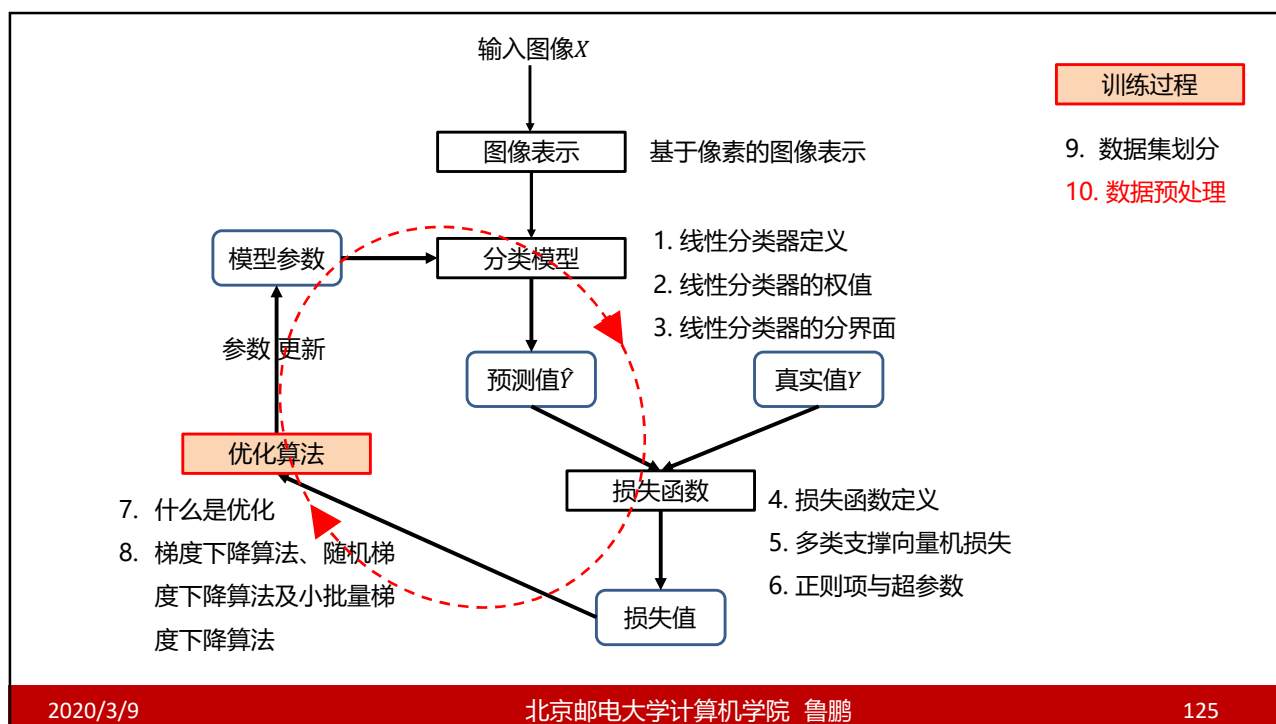
数据能否直接使用？有哪些处理方式？

2020/3/9

北京邮电大学计算机学院 鲁鹏

124

124



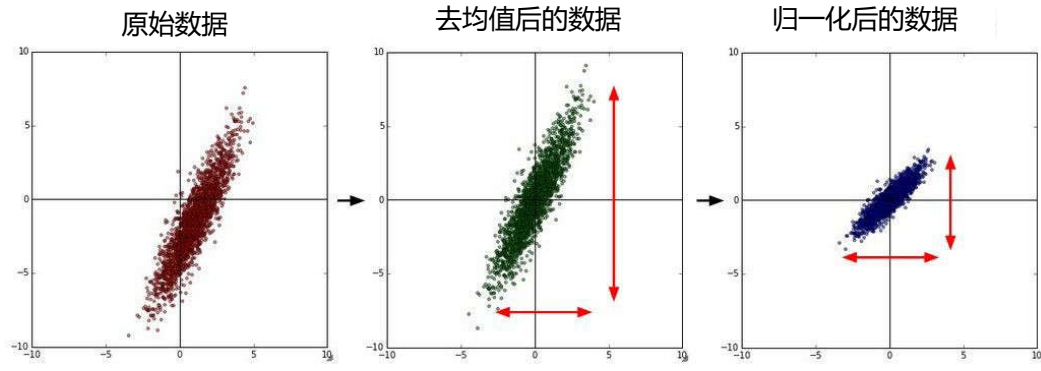
2020/3/9

北京邮电大学计算机学院 鲁鹏

125

125

数据预处理-1



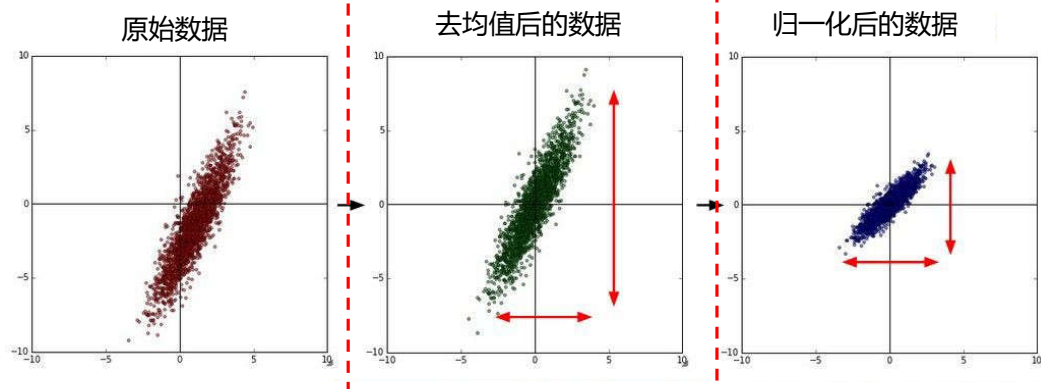
2020/3/9

北京邮电大学计算机学院 鲁鹏

126

126

数据预处理-1



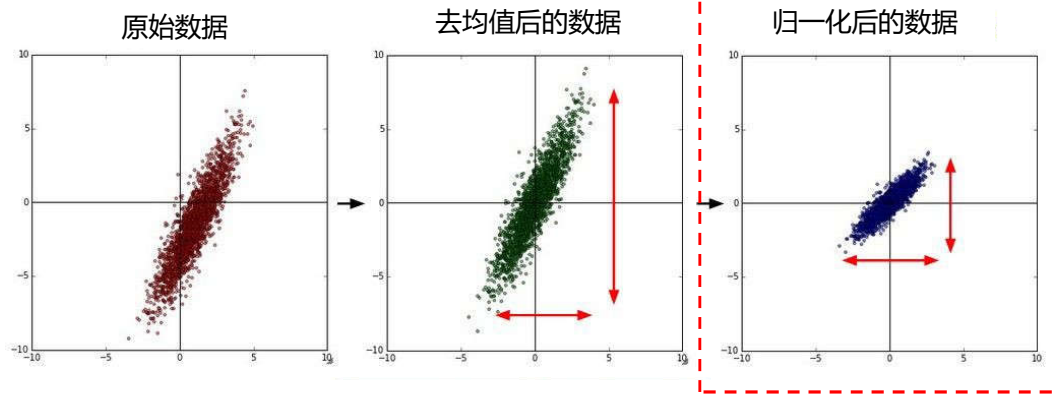
2020/3/9

北京邮电大学计算机学院 鲁鹏

127

127

数据预处理-1



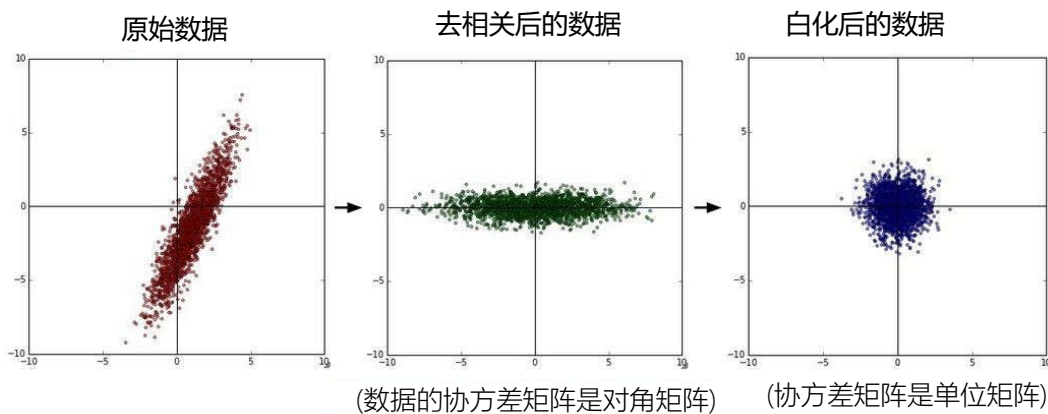
2020/3/9

北京邮电大学计算机学院 鲁鹏

128

128

数据预处理-2



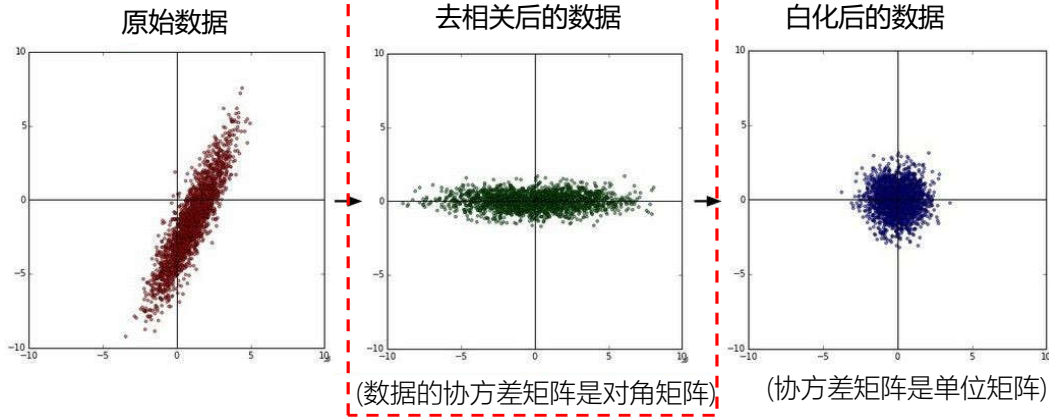
2020/3/9

北京邮电大学计算机学院 鲁鹏

129

129

数据预处理-2



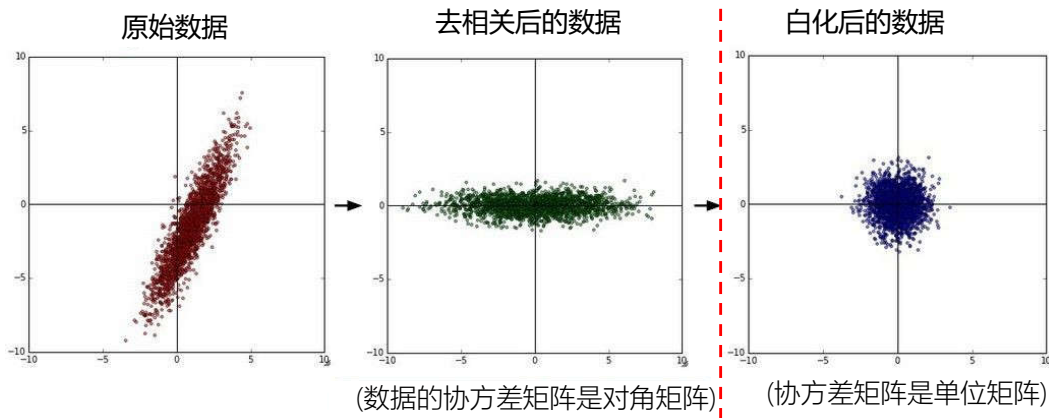
2020/3/9

北京邮电大学计算机学院 鲁鹏

130

130

数据预处理-2

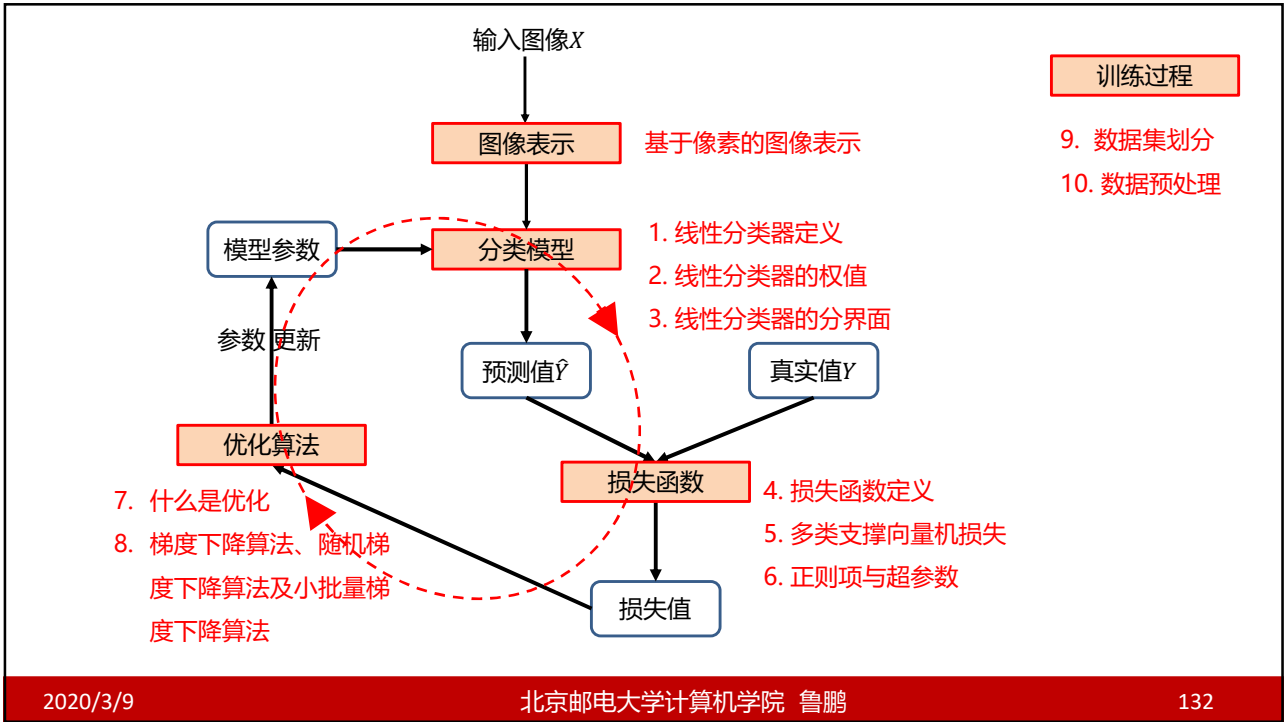


2020/3/9

北京邮电大学计算机学院 鲁鹏

131

131



132

线性分类器体验

<http://vision.stanford.edu/teaching/cs231n-demos/linear-classify/>

2020/3/9 北京邮电大学计算机学院 鲁鹏 133

133