

机器学习基础（二）

主讲：王皓 副研究员| 硕士导师

邮箱：wanghao3@ustc.edu.cn

主页：http://staff.ustc.edu.cn/~wanghao3

本章内容

➤ 线性模型

- 分类问题示例
- 线性分类模型

➤ 无监督学习

- 原型聚类
- 无监督特征学习
- 概率密度估计

线性模型

示例：图像分类

➤任务描述：其目标是根据图像信息中所反映的不同特征，把不同类别的图像区分开来，并从已知的类别标签集合中为给定的输入图片赋予一个类别标签。

➤数据集：CIFAR-10

- 60000张32x32色彩图像，共10类

- 每类6000张图像

airplane



automobile



bird



cat



deer



dog



frog



horse



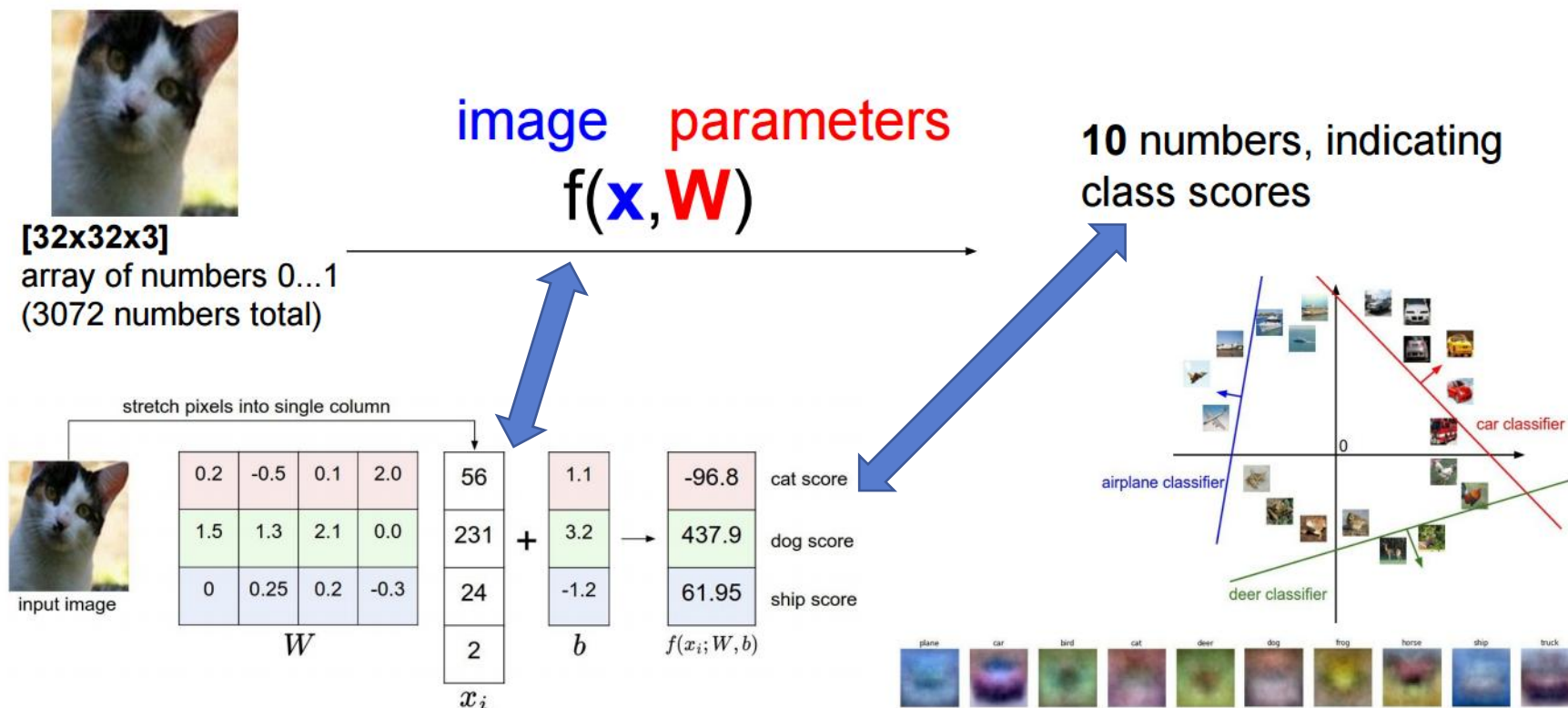
ship



truck

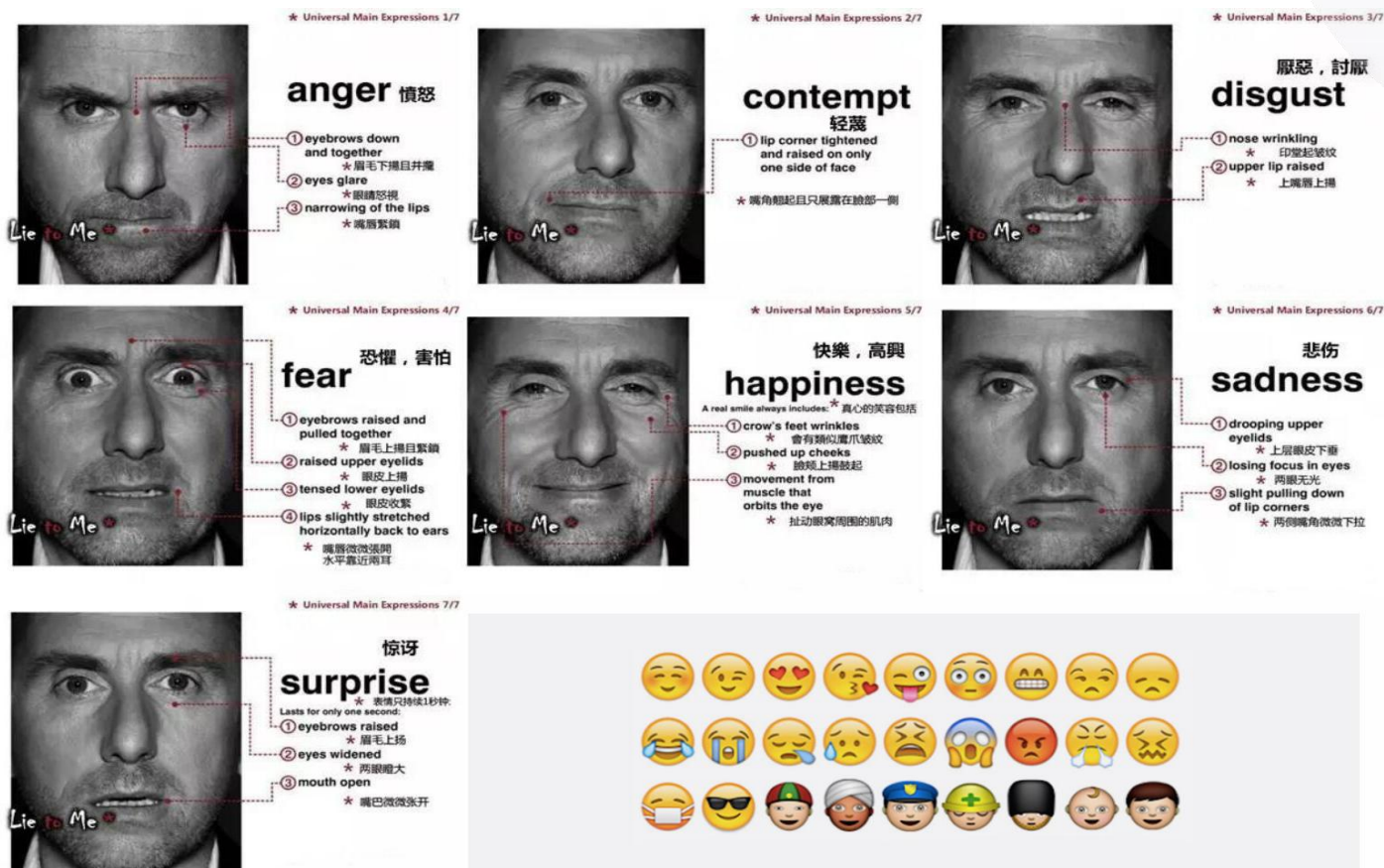


示例：图像分类



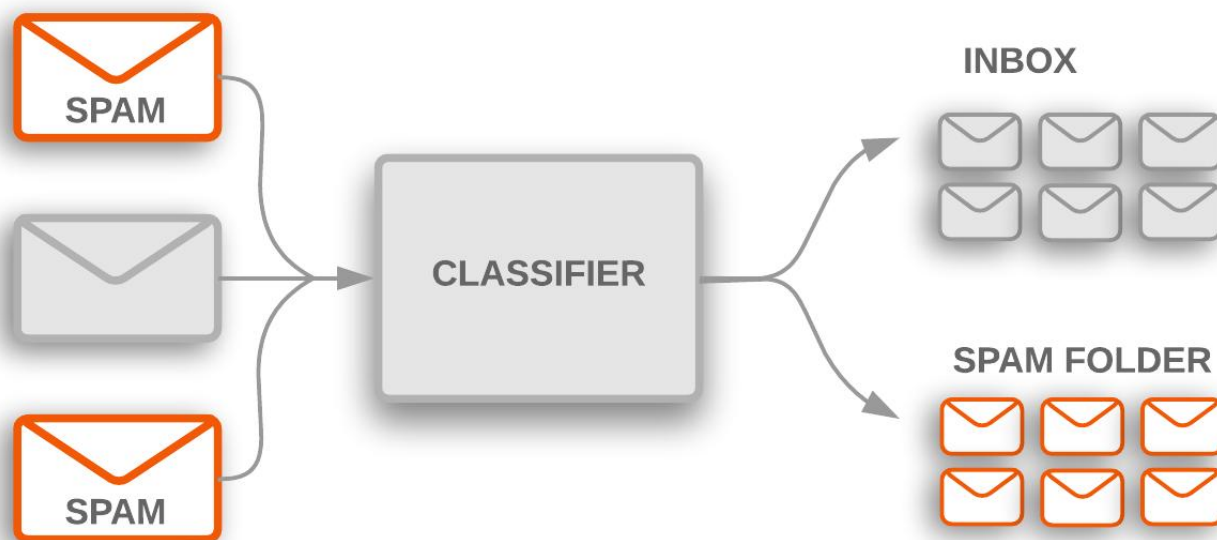
示例：情感分类

➤通过对面部特征的提取，判定人的表情对应的情感



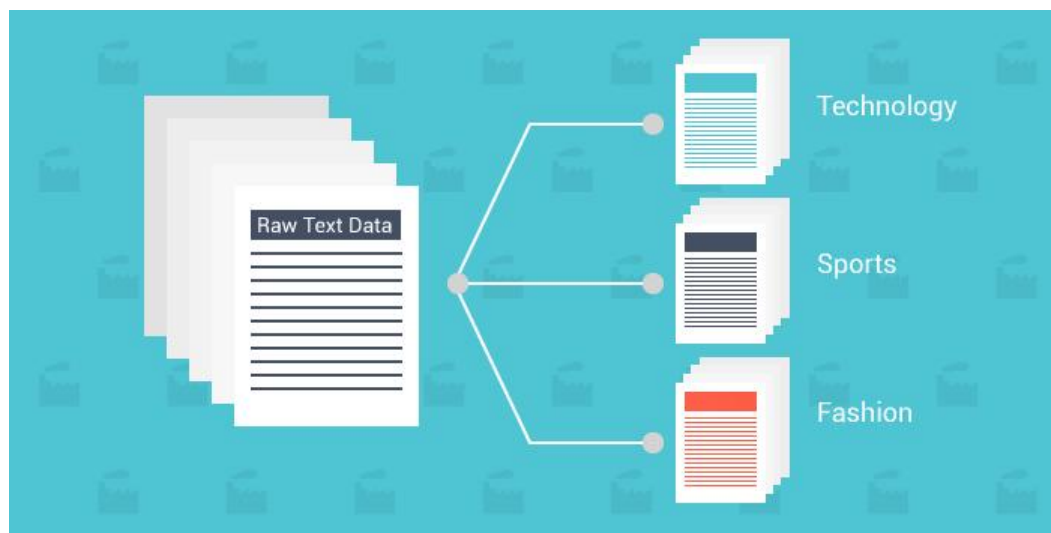
示例：垃圾邮件过滤

➤根据邮件标题、发件人等内容进行垃圾邮件过滤



示例：文档归类

➤ 为新闻、论文等文本进行归类



<https://towardsdatascience.com/automated-text-classification-using-machine-learning-3df4f9570b>

示例：文本分类

➤ 将样本 x 从文本形式转为向量形式

- 词袋模型 (Bag-of-Words, BoW) 模型

the dog is on the table

0	0	1	1	0	1	1	1
are	cat	dog	is	now	on	table	the

比如两个文本“我 喜欢 读书”和“我 讨厌 读书”中共有“我”、“喜欢”、“讨厌”、“读书”四个词，它们的BoW表示分别为

$$\mathbf{v}_1 = [1 \ 1 \ 0 \ 1]^T,$$

$$\mathbf{v}_2 = [1 \ 0 \ 1 \ 1]^T.$$

示例：文本情感分类

- 根据文本内容来判断文本的相应类别

Review (X)

"This movie is fantastic! I really like it because it is so good!"

"Not to my taste, will skip and watch another movie"

"This movie really sucks! Can I get my money back please?"

Rating (Y)



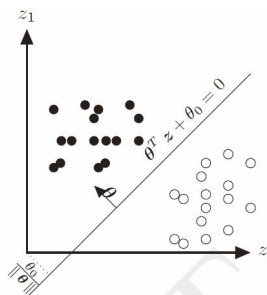
D_1 : “我喜欢读书”

D_2 : “我讨厌读书”

	我	喜欢	讨厌	读书
D_1	1	1	0	1
D_2	1	0	1	1

+

-



线性模型

- 线性模型 (Linear Model) 是机器学习中应用最广泛的模型, 指通过**样本特征的线性组合**来进行预测的模型。给定一个 D 维样本 $\mathbf{x} = [x_1, x_2, \dots, x_D]^T$, 其线性组合函数为

$$f(\mathbf{x}; \mathbf{w}) = w_1x_1 + w_2x_2 + \dots + w_Dx_D + b$$

$$= \mathbf{w}^T \mathbf{x} + b, \rightarrow \text{偏置}$$

$$\text{权重向量 } \mathbf{w} = [\omega_1, \omega_2, \dots, \omega_D]^T$$

- 回归问题: $y = f(\mathbf{x}; \mathbf{w})$, 值域为实数
- 分类问题: $y = g(f(\mathbf{x}; \mathbf{w}))$, $f(\mathbf{x}; \mathbf{w})$ 也称判别函数

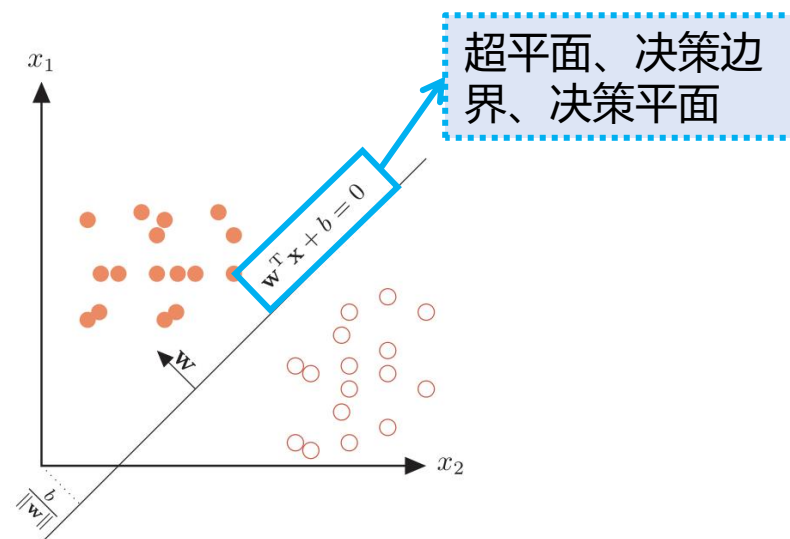
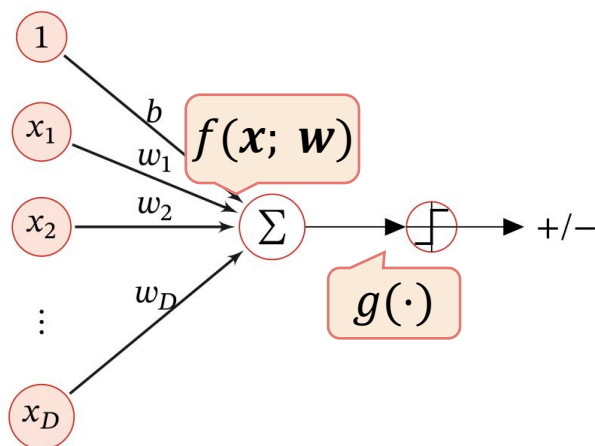
引入一个**非线性的决策函数 (Decision Function)** $g(\cdot)$ 来预测输出目标

二分类问题

$g(\cdot)$ 可以是符号函数 (Sign Function) , 定义为

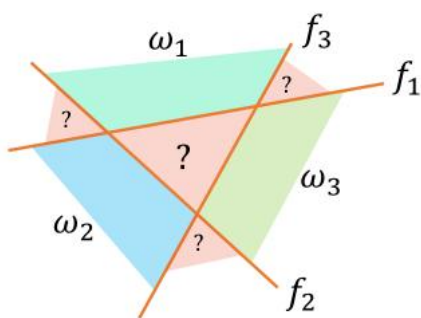
$$g(f(\mathbf{x}; \mathbf{w})) = \text{sgn}(f(\mathbf{x}; \mathbf{w}))$$

$$\triangleq \begin{cases} +1 & \text{if } f(\mathbf{x}; \mathbf{w}) > 0, \\ -1 & \text{if } f(\mathbf{x}; \mathbf{w}) < 0. \end{cases}$$

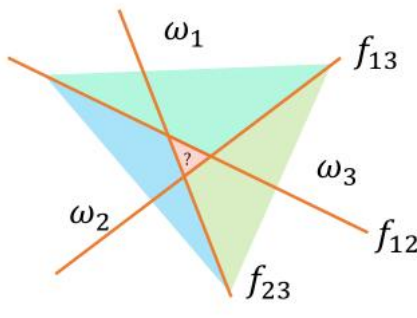


多分类问题

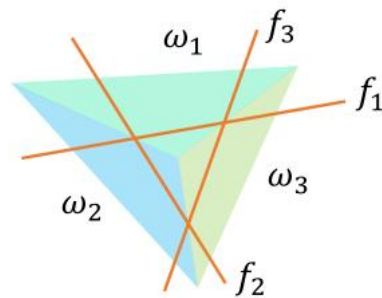
- 多分类 (Multi-class Classification) 问题是指分类的类别数 C 大于2, 因此一般需要多个线性判别函数, 但设计这些判别函数有很多方式, 常用的有以下三种:



(a) “一对其余”方式



(b) “一对一”方式



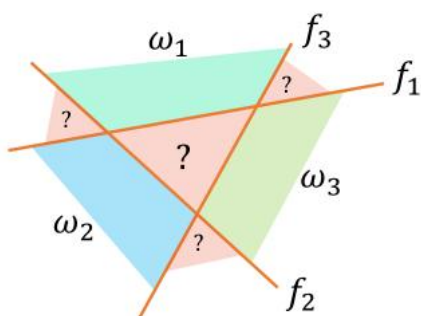
(c) “argmax”方式

- “一对其余”：把多分类问题转换为 C 个 “一对其余” 的二分类问题。这种方式共需要 C 个判别函数，每个判别函数的目的是将属于类别 c 和不属于类别 c 的样本区分开。
- “一对一”：把多分类问题转换为 $C(C-1)/2$ 个 “一对一” 的二分类问题。这种方式共需要 $C(C-1)/2$ 个判别函数，每个判别函数的目的是将属于类别 i 和属于类别 j 的样本区分开。

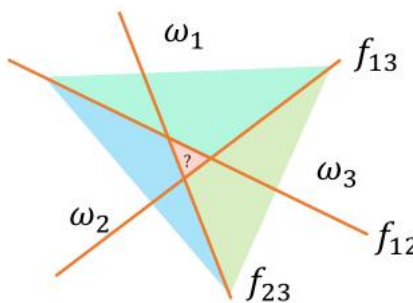
✓ 上述两种方式都存在一个缺陷：特征空间中存在会存在一些难以确定类别的区域

多分类问题

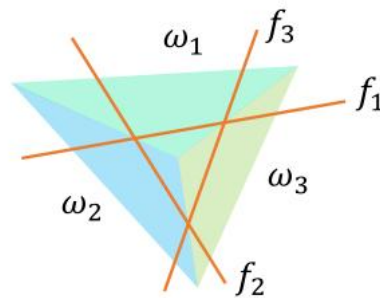
- 多分类 (Multi-class Classification) 问题是指分类的类别数 C 大于2, 因此一般需要**多个线性判别函数**, 但设计这些判别函数有很多方式, 常用的有以下三种:



(a) “一对其余” 方式



(b) “一对一” 方式



(c) “argmax” 方式

- “argmax” : 一种**改进的 “一对其余” 方式**, 共需要 C 个判别函数。

$$f_c(\mathbf{x}; \mathbf{w}_c) = \mathbf{w}_c^T \mathbf{x} + b_c, \quad c \in \{1, \dots, C\}$$

对于样本 \mathbf{x} , 如果存在一个类别 c , 相对于所有的其他类别 $\tilde{c} (\tilde{c} \neq c)$ 有 $f_c(\mathbf{x}; \mathbf{w}_c) > f_{\tilde{c}}(\mathbf{x}; \mathbf{w}_{\tilde{c}})$, 那么 \mathbf{x} 属于类别 c 。则 “argmax” 方式的预测函数定义为

$$y = \arg \max_{c=1}^C f_c(\mathbf{x}; \mathbf{w}_c).$$

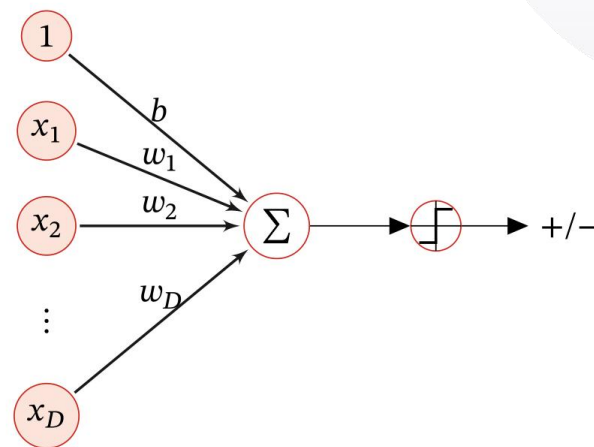
线性分类模型

- Logistic回归
- Softmax回归
- 支持向量机
-

Logistic Regression

➤ 模型 (本节中采用 $y \in \{0, 1\}$ 以符合Logistic回归的描述习惯)

$$g(f(\mathbf{x}; \mathbf{w})) = \begin{cases} 1 & \text{if } f(\mathbf{x}; \mathbf{w}) > 0 \\ 0 & \text{if } f(\mathbf{x}; \mathbf{w}) < 0 \end{cases}$$



分类问题

➤ 将分类问题看作条件概率估计问题

- 引入非线性函数 $g : \mathbb{R}^D \rightarrow (0, 1)$ 来预测类别标签的条件概率 $p(y = c | \mathbf{x})$ 。

- 以二分类为例：

$$p(y = 1 | \mathbf{x}) = g(f(\mathbf{x}; \mathbf{w}))$$

- 函数 f ：线性函数
- 函数 g ：称为**激活函数**，把线性函数的值域从实数区间“挤压”到了 $(0, 1)$ 之间，可以用来表示**概率**。

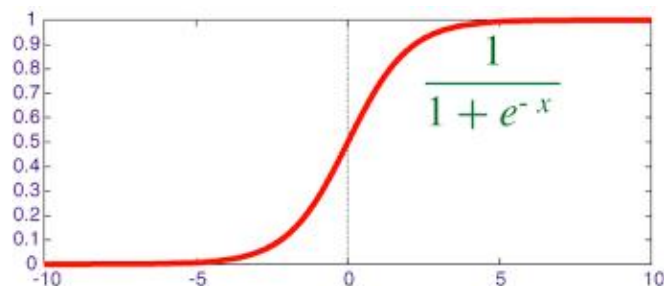


如何构建函数 g ?

Logistic函数与回归

- Logistic函数

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$



- Logistic回归

$$p(y = 1|\mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x}) \rightarrow \begin{array}{l} \text{增广特征向量 } \mathbf{x} = [x_1, x_2, \dots, x_D, 1]^\top \\ \text{增广权重向量 } \mathbf{w} = [\omega_1, \omega_2, \dots, \omega_D, b]^\top \end{array}$$
$$\triangleq \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})}$$

学习准则

➤ 模型预测条件概率 $p_{\theta}(y|\mathbf{x})$

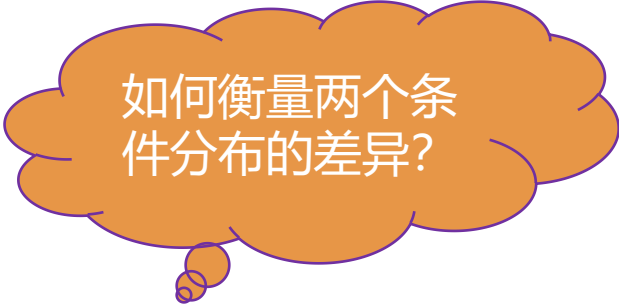
$$p_{\theta}(y = 1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$$

➤ 真实条件概率 $p_r(y|\mathbf{x})$

• 对于一个样本 (\mathbf{x}, y^*) , $y^* \in \{0, 1\}$, 其真实条件概率为

$$p_r(y = 1|\mathbf{x}) = y^*$$

$$p_r(y = 0|\mathbf{x}) = 1 - y^*$$



如何衡量两个条件分布的差异?

熵 (Entropy)

➤在信息论中，熵用来衡量一个随机事件的不确定性。

- 自信息 (Self Information)

$$I(x) = -\log(p(x))$$

- 熵 $H(X) = \mathbb{E}_X[I(x)] = \mathbb{E}_X[-\log p(x)] = -\sum_{x \in \mathcal{X}} p(x) \log p(x)$

✓熵越高，则事件的不确定性越大，随机变量的信息越多；

✓熵越低，则事件的不确定性越小，随机变量的信息越少。

- 在对某一特定概率分布 $q(y)$ 进行编码时，熵 $H(q)$ 也是理论上最优的平均编码长度，这种编码方式称为熵编码 (Entropy Encoding)

交叉熵 (Cross Entropy)

- 交叉熵是按照概率分布 q 的最优编码对真实分布为 p 的信息进行编码的长度。

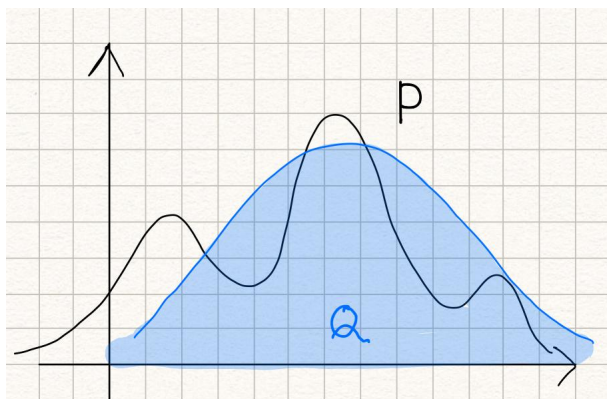
$$\begin{aligned} H(p, q) &= \mathbb{E}_p[-\log q(x)] \\ &= -\sum_x p(x) \log q(x) \end{aligned}$$

在给定分布 q 的情况下

- 如果 p 和 q 越接近, 交叉熵越小;
- 如果 p 和 q 越远, 交叉熵就越大。

KL散度 (Kullback-Leibler Divergence)

- KL散度是用概率分布 q 来近似 p 时所造成的**信息损失量**。
- KL散度是按照概率分布 q 的最优编码对真实分布为 p 的信息进行编码，其平均编码长度（即交叉熵） $H(p, q)$ 和 p 的最优平均编码长度（即熵） $H(p)$ 之间的差异。



交叉熵 — 熵

$$\begin{aligned} \text{KL}(p, q) &= H(p, q) - H(p) \\ &= \sum_x p(x) \log \frac{p(x)}{q(x)} \end{aligned}$$

越小效果越好

交叉熵损失

$$D_{KL}(p_r(y|x) || p_\theta(y|x)) = \sum_{y=0}^1 p_r(y|x) \log \frac{p_r(y|x)}{p_\theta(y|x)} \quad \text{KL散度}$$

$$\propto - \sum_{y=0}^1 p_r(y|x) \log p_\theta(y|x) \quad \text{交叉熵损失}$$

$$\begin{aligned} & \boxed{y^* \in \{0, 1\} \text{为} x \text{的真实标签}} - I(y^* = 1) \log p_\theta(y = 1|x) - I(y^* = 0) \log p_\theta(y = 0|x) \\ &= - y^* \log p_\theta(y = 1|x) - (1 - y^*) \log p_\theta(y = 0|x) \end{aligned}$$

$$= - \log p_\theta(y^*|x) \quad \text{负对数似然}$$

梯度下降法

- 给定 N 个训练样本 $\{(x^{(n)}, y^{(n)})\}_{n=1}^N$ ，用Logistic回归模型对每个样本 $x^{(n)}$ 进行预测，输出其标签为1的后验概率，记为 $\hat{y}^{(n)}$ 。

$$\hat{y}^{(n)} = \sigma(\mathbf{w}^T \mathbf{x}^{(n)}), 1 \leq n \leq N$$

- 由于 $y^{(n)} \in \{0, 1\}$ ，样本 $(\mathbf{x}^{(n)}, y^{(n)})$ 的真实条件概率可以表示为

$$p_r(y^{(n)} = 1 | \mathbf{x}^{(n)}) = y^{(n)}$$

$$p_r(y^{(n)} = 0 | \mathbf{x}^{(n)}) = 1 - y^{(n)}$$

- 故使用交叉熵损失函数，模型在训练集的风险函数为

$$\mathcal{R}(\mathbf{w}) = -\frac{1}{N} \sum_{n=1}^N (p_r(y^{(n)} = 1 | \mathbf{x}^{(n)}) \log \hat{y}^{(n)} + p_r(y^{(n)} = 0 | \mathbf{x}^{(n)}) \log(1 - \hat{y}^{(n)}))$$

$$\text{化简后得: } \mathcal{R}(\mathbf{w}) = -\frac{1}{N} \sum_{n=1}^N (y^{(n)} \log \hat{y}^{(n)} + (1 - y^{(n)}) \log(1 - \hat{y}^{(n)}))$$

梯度下降法

- 风险函数 $\mathcal{R}(\mathbf{w})$ 关于参数 \mathbf{w} 的偏导数为

$$\begin{aligned}\frac{\partial \mathcal{R}(\mathbf{w})}{\partial \mathbf{w}} &= -\frac{1}{N} \sum_{n=1}^N \left(y^{(n)} \frac{\hat{y}^{(n)}(1 - \hat{y}^{(n)})}{\hat{y}^{(n)}} \mathbf{x}^{(n)} - (1 - y^{(n)}) \frac{\hat{y}^{(n)}(1 - \hat{y}^{(n)})}{1 - \hat{y}^{(n)}} \mathbf{x}^{(n)} \right) \\ &= -\frac{1}{N} \sum_{n=1}^N \left(y^{(n)}(1 - \hat{y}^{(n)}) \mathbf{x}^{(n)} - (1 - y^{(n)}) \hat{y}^{(n)} \mathbf{x}^{(n)} \right) \\ &= -\frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)} (y^{(n)} - \hat{y}^{(n)}).\end{aligned}$$

- 采用梯度下降法，Logistic回归的训练过程为：

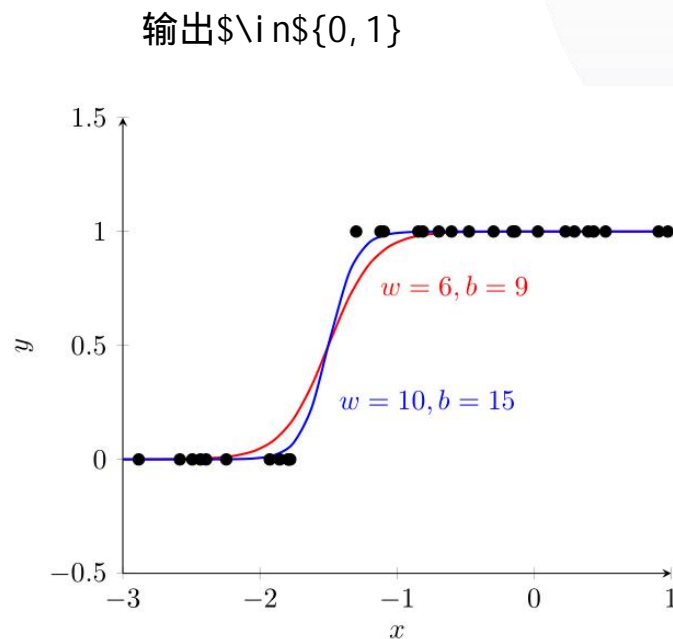
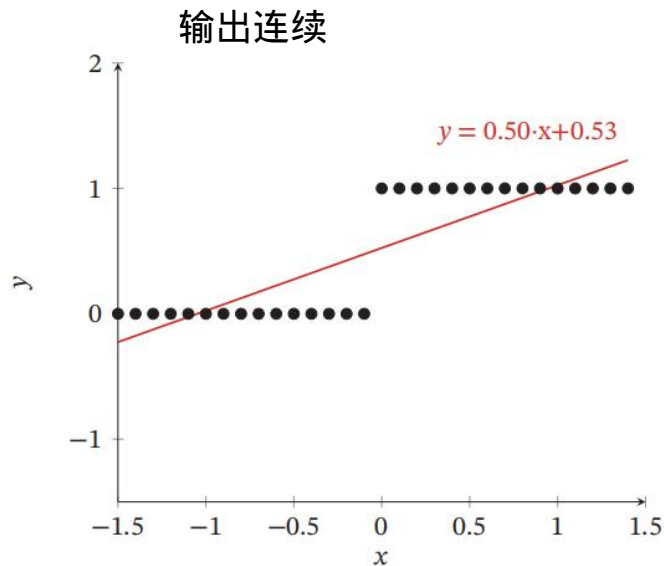
- ✓ 初始化 $\mathbf{w}_0 \leftarrow \mathbf{0}$

- ✓ 迭代更新参数： $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \alpha \frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)} (y^{(n)} - \hat{y}_{\mathbf{w}_t}^{(n)})$

Logistic回归

本质上是线性判别 $y=wx+b \geq 0 \Rightarrow 1:0$

➤ 一维数据的二分类问题示例



Softmax回归

没有3元线性

➤ Softmax 回归，也称为**多项或多类的Logistic回归**，是 Logistic回归在多分类问题上的推广。

✓ Softmax函数

$$\text{softmax}(x_k) = \frac{\exp(x_k)}{\sum_{i=1}^K \exp(x_i)}$$

Logistic回归的作用

The diagram illustrates the Softmax function's role in Logistic regression. It shows a 2x1 column vector $\begin{bmatrix} 1.2 \\ 0.9 \end{bmatrix}$ as input to a box labeled 'Softmax'. The output is a 3x1 column vector $\begin{bmatrix} 0.46 \\ 0.34 \\ 0.20 \end{bmatrix}$. The text 'Logistic回归的作用' is positioned below the input vector.

✓ 对于多类问题，类别标签 $y \in \{1, 2, \dots, C\}$ 可以有 C 个取值。给定一个样本 \mathbf{x} ，Softmax回归预测的属于类别 c 的条件概率为：

$$\begin{aligned} p(y = c | \mathbf{x}) &= \text{softmax}(\mathbf{w}_c^\top \mathbf{x}) \\ &= \frac{\exp(\mathbf{w}_c^\top \mathbf{x})}{\sum_{c'=1}^C \exp(\mathbf{w}_{c'}^\top \mathbf{x})}, \end{aligned}$$

Softmax回归

➤ 因此Softmax回归的**决策函数**可以表示为：

$$\begin{aligned}\hat{y} &= \arg \max_{c=1}^C p(y = c | x) \\ &= \arg \max_{c=1}^C \frac{\exp(\omega_c^T x)}{\sum_{c'}^C \exp(\omega_{c'}^T x)} \\ &= \arg \max_{c=1}^C \omega_c^T x\end{aligned}$$

参数学习

➤ 给定 N 个训练样本 $\{(\mathbf{x}^{(n)}, \mathbf{y}^{(n)})\}_{n=1}^N$, $\mathbf{y}^{(n)}$ 用 C 维 one-hot 向量表示

- 学习准则：使用交叉熵损失函数，Softmax 回归模型的风险函数为

$$\begin{aligned}\mathcal{R}(\mathbf{W}) &= -\frac{1}{N} \sum_{n=1}^N \sum_{c=1}^C \mathbf{y}_c^{(n)} \log \hat{\mathbf{y}}_c^{(n)} \\ &= -\frac{1}{N} \sum_{n=1}^N (\mathbf{y}^{(n)})^\top \log \hat{\mathbf{y}}^{(n)},\end{aligned}$$

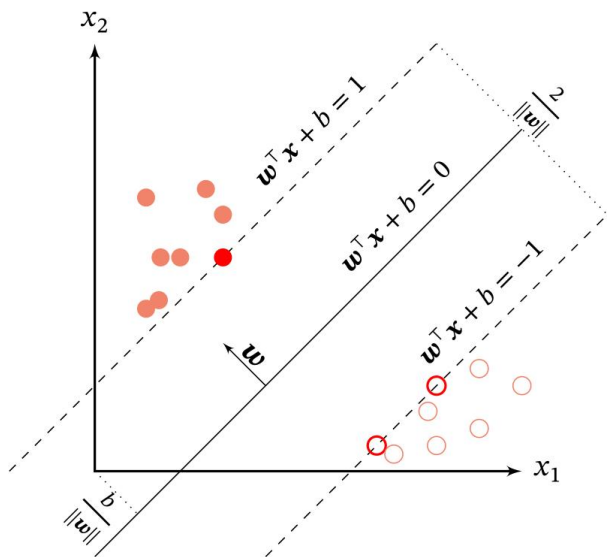
$\hat{\mathbf{y}}^{(n)} = \text{softmax}(\mathbf{W}^T \mathbf{x}^{(n)})$
为样本 $\mathbf{x}^{(n)}$ 在每个类别的
后验概率

- 优化：使用梯度下降法

$$\frac{\partial \mathcal{R}(\mathbf{W})}{\partial \mathbf{W}} = -\frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)} (\mathbf{y}^{(n)} - \hat{\mathbf{y}}^{(n)})^\top.$$

- 参数更新： $\mathbf{W}_{t+1} \leftarrow \mathbf{W}_t + \alpha \left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)} (\mathbf{y}^{(n)} - \hat{\mathbf{y}}^{(n)} \mathbf{W}_t)^\top \right)$

支持向量机 (Support Vector Machine, SVM)



数据集中所有满足 $y^{(n)}(\mathbf{w}^T \mathbf{x}^{(n)} + b) = 1$ 的样本点，都称为“支持向量”

- 数据集 \mathcal{D} 中每个样本 $\mathbf{x}^{(n)}$ 到分割超平面 $\mathbf{w}^T \mathbf{x} + b = 0$ 的距离为：

$$\gamma^{(n)} = \frac{|\mathbf{w}^T \mathbf{x}^{(n)} + b|}{\|\mathbf{w}\|} = \frac{y^{(n)}(\mathbf{w}^T \mathbf{x}^{(n)} + b)}{\|\mathbf{w}\|}.$$

- 支持向量机的目标是寻找一个超平面 (\mathbf{w}^*, b^*) 使得 γ 最大，即

$$\max_{\mathbf{w}, b} \quad \gamma$$

$$\text{s.t.} \quad \frac{y^{(n)}(\mathbf{w}^T \mathbf{x}^{(n)} + b)}{\|\mathbf{w}\|} \geq \gamma, \forall n \in \{1, \dots, N\}.$$



$$\begin{aligned} \max_{\mathbf{w}, b} \quad & \frac{1}{\|\mathbf{w}\|^2} \\ \text{s.t.} \quad & y^{(n)}(\mathbf{w}^T \mathbf{x}^{(n)} + b) \geq 1, \forall n \in \{1, \dots, N\}. \end{aligned}$$

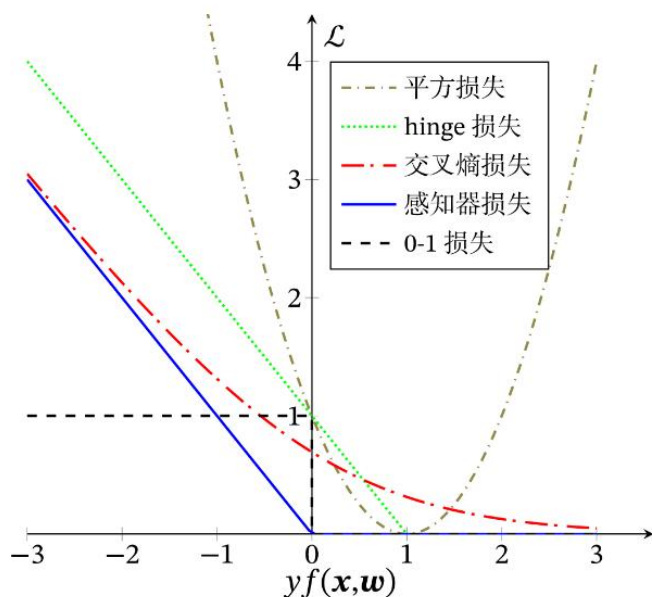
线性分类模型小结

线性模型	激活函数	损失函数	优化方法
线性回归	-	$(y - \mathbf{w}^\top \mathbf{x})^2$	最小二乘、梯度下降
Logistic 回归	$\sigma(\mathbf{w}^\top \mathbf{x})$	$\mathbf{y} \log \sigma(\mathbf{w}^\top \mathbf{x})$	梯度下降
Softmax 回归	$\text{softmax}(\mathbf{W}^\top \mathbf{x})$	$\mathbf{y} \log \text{softmax}(\mathbf{W}^\top \mathbf{x})$	梯度下降
感知器	$\text{sgn}(\mathbf{w}^\top \mathbf{x})$	$\max(0, -y\mathbf{w}^\top \mathbf{x})$	随机梯度下降
支持向量机	$\text{sgn}(\mathbf{w}^\top \mathbf{x})$	$\max(0, 1 - y\mathbf{w}^\top \mathbf{x})$	二次规划、SMO 等

不同损失函数的对比

损失函数measure样本的输出与实际tag的误差
perception中为 $\max(0, -yf(x; w))$

- 为方便比较, 定义类别标签 $y \in \{+1, -1\}$, 并定义 $f(x; w) = w^T x + b$.
- 对于样本 (x, y) , 若 $yf(x; w) > 0$, 则分类正确, 并且 $yf(x; w)$ 越大, 模型预测越正确; 反之若 $yf(x; w) < 0$, 则分类错误, 且越小越错误。
- 一个好的损失函数应该随着 $yf(x; w)$ 的增大而减小。



$$\mathcal{L}_{LR} = \log(1 + \exp(-yf(x; w))). \quad \text{LR回归}$$

$$\mathcal{L}_p = \max(0, -yf(x; w)). \quad \text{感知机}$$

$$\mathcal{L}_{hinge} = \max(0, 1 - yf(x; w)). \quad \text{SVM}$$

$$\mathcal{L}_{squared} = (1 - yf(x; w))^2. \quad \text{平方损失}$$

无监督学习

无监督学习 (Unsupervised Learning)

➤ 监督学习

- 建立输入与输出之间的映射关系 $f: x \rightarrow y$

➤ 无监督学习

- 指从**无标签的数据**中学习出一些**有用的模式**。一般直接从原始数据中学习，不借助于任何人工给出标签或者反馈等指导信息。
- 发现隐藏的数据中的有价值信息，包括有效的**特征、类别、结构以及概率分布**等

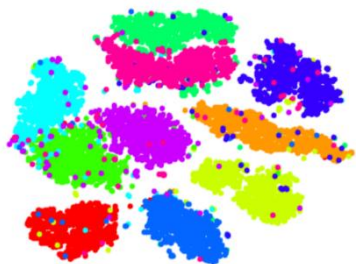
为什么要无监督学习？

大脑有大约 10^{14} 个突触, 我们只能活大约 10^9 秒. 所以我们有比数据更多的参数. 这启发了我们必须进行大量无监督学习的想法, 因为感知输入 (包括本体感受) 是我们可以获得每秒 10^5 维约束的唯一途径.

——杰弗里·辛顿 (Geoffrey Hinton)

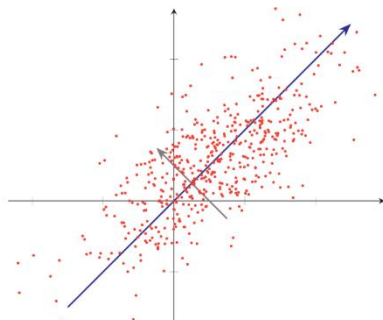
2018 年图灵奖获得者

典型的无监督学习问题



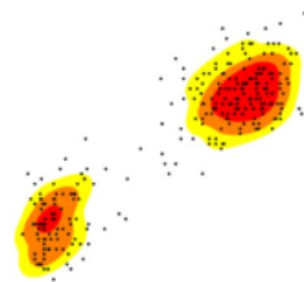
聚类

聚类 (Clustering) 是将一组样本根据一定的准则划分到不同的组 (也称为簇 Cluster)。一个比较通用的准则是组内样本的相似性要高于组间样本的相似性



无监督特征学习

无监督特征学习 (Unsupervised Feature Learning) 是从无标签训练数据中挖掘有效的特征或表示。无监督特征学习一般用来进行降维、数据可视化或监督学习前期的数据预处理



概率密度估计

概率密度估计 (Probabilistic Density Estimation) 简称密度估计, 是根据一组训练样本来估计样本空间的概率密度。密度估计可以分为参数密度估计和非参数密度估计

原型聚类

➤ 原型聚类

- 也称为“基于原型的聚类” (prototype-based clustering), 此类算法假设聚类结构能通过一组原型刻画。

➤ 算法过程：

- 通常情况下，算法先对原型进行初始化，再对原型进行迭代更新求解。
- 采用不同的原型表示、不同的求解方式，将产生不同的算法。

原型聚类：K均值

➤ 优化问题： $\min_{T, \mu} E(T, \mu) = \|X - T\mu\|_F^2$

算法流程（迭代优化）：

初始化每个簇的均值向量（共k个簇）

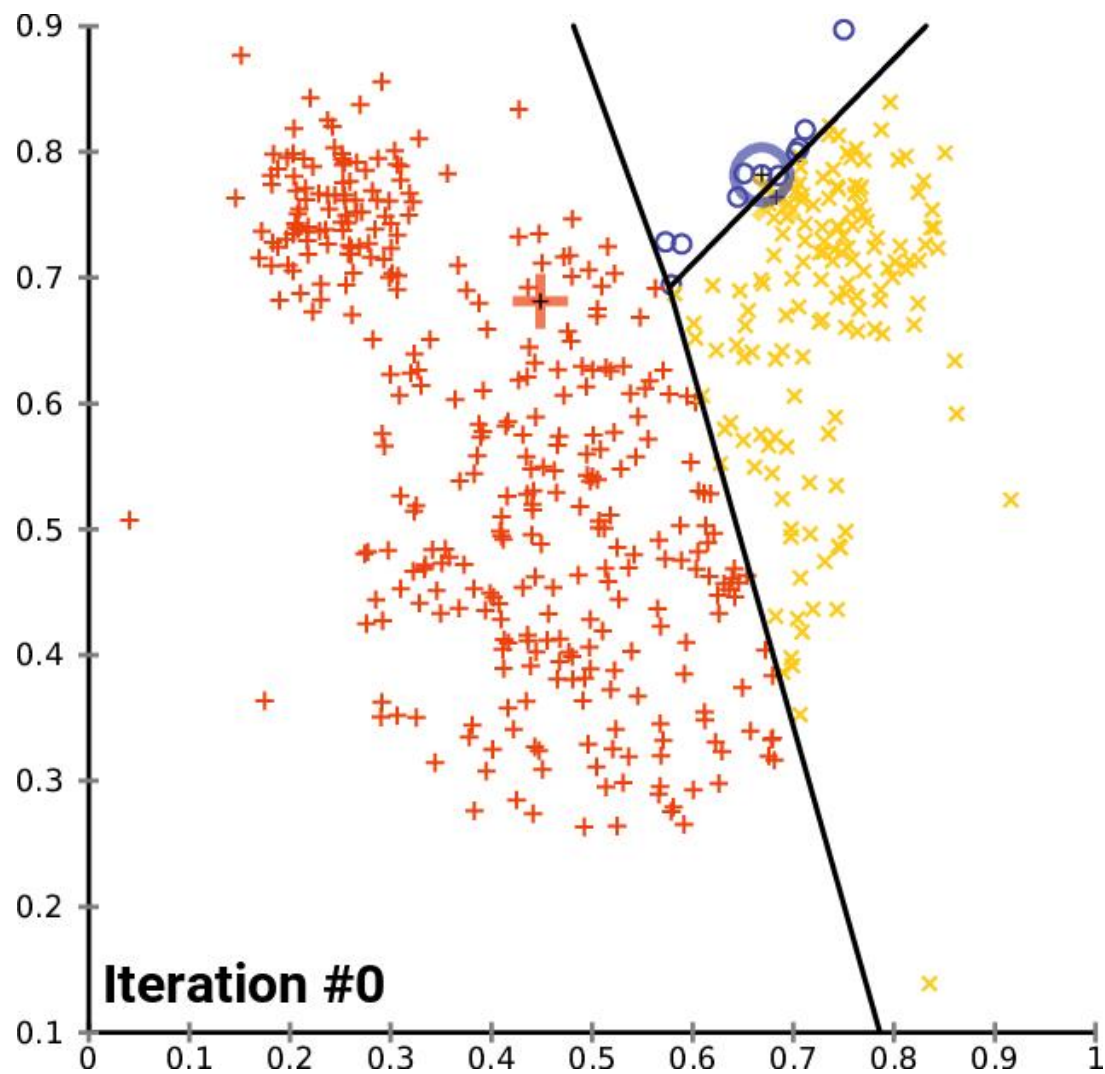
repeat

1. 将每个样本分配给最近的簇；
2. 计算每个簇的均值向量

until 当前均值向量均未更新

$$T^{(t)} \leftarrow \min_T E(T, \mu^{(t-1)})$$
$$\mu^{(t)} \leftarrow \min_{\mu} E(T^{(t)}, \mu)$$

原型聚类：K均值



主成份分析PCA

➤ 一种最常用的数据降维方法，使得在转换后的空间中数据的方差最大。

- 样本点 $\mathbf{x}^{(n)}$ 投影之后的表示为

$$z^{(n)} = \mathbf{w}^\top \mathbf{x}^{(n)}$$

- 所有样本投影后的方差为

$$\begin{aligned}\sigma(\mathbf{X}; \mathbf{w}) &= \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}^{(n)} - \mathbf{w}^\top \bar{\mathbf{x}})^2 \\ &= \frac{1}{N} (\mathbf{w}^\top \mathbf{X} - \mathbf{w}^\top \bar{\mathbf{X}})(\mathbf{w}^\top \mathbf{X} - \mathbf{w}^\top \bar{\mathbf{X}})^\top \\ &= \mathbf{w}^\top \Sigma \mathbf{w},\end{aligned}$$

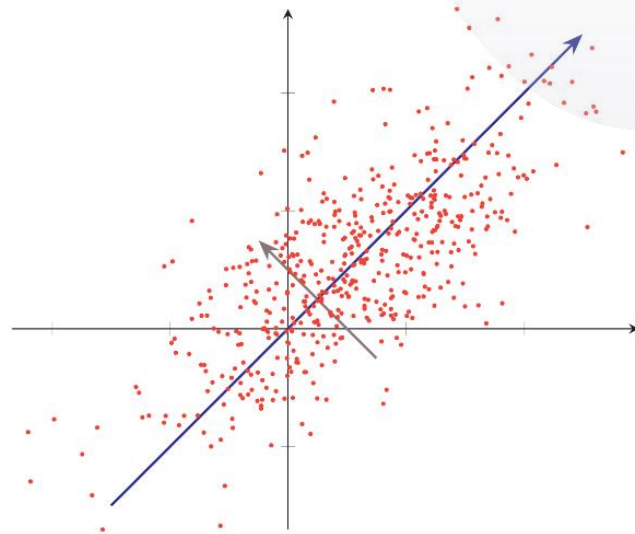
- 目标函数

$$\max_{\mathbf{w}} \mathbf{w}^\top \Sigma \mathbf{w} + \lambda(1 - \mathbf{w}^\top \mathbf{w})$$



- 对目标函数求导并令导数等于 0，可得

$$\Sigma \mathbf{w} = \lambda \mathbf{w}$$



(线性) 编码

- 给定一组基向量 $A = [\mathbf{a}_1, \dots, \mathbf{a}_M]$, 将输入样本 \mathbf{x} 表示为这些基向量的线性组合

$$\mathbf{x} = \sum_{m=1}^M z_m \mathbf{a}_m$$

$$= \mathbf{A}\mathbf{z},$$

字典 (dictionary)

编码 (encoding)

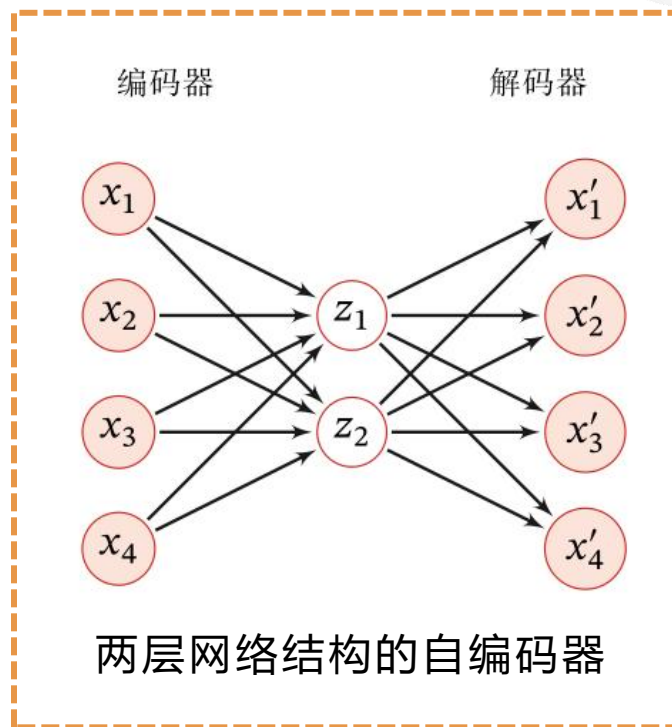
$$\mathbf{x} = \mathbf{A}\mathbf{z}$$

自编码器 (Auto-Encoder)

- 编码器 (Encoder) $f : \mathbb{R}^D \rightarrow \mathbb{R}^M$
- 解码器 (Decoder) $g : \mathbb{R}^M \rightarrow \mathbb{R}^D$

- 学习目标：最小化**重构错误**

$$\begin{aligned}\mathcal{L} &= \sum_{n=1}^N \|\mathbf{x}^{(n)} - g(f(\mathbf{x}^{(n)}))\|^2 \\ &= \sum_{n=1}^N \|\mathbf{x}^{(n)} - f \circ g(\mathbf{x}^{(n)})\|^2.\end{aligned}$$



概率密度估计

➤ 参数密度估计 (Parametric Density Estimation)

- 根据先验知识假设随机变量服从某种分布，然后通过训练样本来估计 **分布的参数**。
- 估计方法：**最大似然估计**

$$\log p(\mathcal{D}; \theta) = \sum_{n=1}^N \log p(\mathbf{x}^{(n)}; \theta).$$

对数似然函数

$$\theta^{ML} = \arg \max_{\theta} \sum_{n=1}^N \log p(\mathbf{x}^{(n)}; \theta).$$

参数估计问题转化为最优化问题

➤ 非参数密度估计 (Nonparametric Density Estimation)

- 不假设数据服从某种分布，通过**将样本空间划分为不同的区域**并估计每个区域的概率来近似数据的概率密度函数。

参数密度估计

- 正态分布

假设样本 $\mathbf{x} \in \mathbb{R}^D$ 服从正态分布

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right),$$

其中 $\boldsymbol{\mu}$ 和 $\boldsymbol{\Sigma}$ 分别为正态分布的均值和方差.

数据集 $\mathcal{D} = \{\mathbf{x}^{(n)}\}_{n=1}^N$ 的对数似然函数为

$$\log p(\mathcal{D}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{N}{2} \log\left((2\pi)^D |\boldsymbol{\Sigma}|\right) - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}^{(n)} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}^{(n)} - \boldsymbol{\mu}).$$

分别求上式关于 $\boldsymbol{\mu}, \boldsymbol{\Sigma}$ 的偏导数, 并令其等于 0, 可得,

$$\begin{aligned}\boldsymbol{\mu}^{ML} &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)}, \\ \boldsymbol{\Sigma}^{ML} &= \frac{1}{N} \sum_{n=1}^N (\mathbf{x}^{(n)} - \boldsymbol{\mu}^{ML})(\mathbf{x}^{(n)} - \boldsymbol{\mu}^{ML})^\top.\end{aligned}$$

参数密度估计

- 多项分布

数据集 $\mathcal{D} = \{\mathbf{x}^{(n)}\}_{n=1}^N$ 的对数似然函数为

$$\log p(\mathcal{D}|\boldsymbol{\mu}) = \sum_{n=1}^N \sum_{k=1}^K x_k^{(n)} \log(\mu_k). \quad (9.34)$$

多项分布的参数估计为约束优化问题. 引入拉格朗日乘子 λ , 将原问题转换为无约束优化问题.

$$\max_{\boldsymbol{\mu}, \lambda} \sum_{n=1}^N \sum_{k=1}^K x_k^{(n)} \log(\mu_k) + \lambda \left(\sum_{k=1}^K \mu_k - 1 \right). \quad (9.35)$$

分别求上式关于 μ_k, λ 的偏导数, 并令其等于 0. 可得,

$$\mu_k^{ML} = \frac{m_k}{N}, \quad 1 \leq k \leq K \quad (9.36)$$

其中 $m_k = \sum_{n=1}^N x_k^{(n)}$ 为数据集中取值为第 k 个状态的样本数量.

参数密度估计一般存在以下问题

➤ 模型选择问题

- 如何选择数据分布的密度函数？
- 实际数据的分布往往是**非常复杂**的，而不是简单的正态分布或多项分布。

➤ 不可观测变量问题

- 即我们用来训练的样本只包含**部分的可观测变量**，还有一些非常关键的变量是无法观测的，这导致我们很难准确估计数据的真实分布。

➤ 维度灾难问题

- **高维数据**的参数估计十分困难
- 随着维度的增加，估计参数所需要的样本数量指数增加。在样本不足时会出现**过拟合**。

非参数密度估计

- 对于高维空间中的一个随机向量 \mathbf{x} ，假设其服从一个未知分布 $p(\mathbf{x})$ ，则 \mathbf{x} 落入空间中的小区域 \mathcal{R} 的概率为

$$P = \int_{\mathcal{R}} p(\mathbf{x}) d\mathbf{x}.$$

- 给定 N 个训练样本 $D = \{\mathbf{x}^{(n)}\}_{n=1}^N$ ，落入区域 \mathcal{R} 的样本数量 K 服从二项分布

$$P_K = \binom{N}{K} P^K (1-P)^{1-K},$$

- 当 N 非常大时，我们可以近似认为

$$P \approx p(\mathbf{x})V$$

- 假设区域 \mathcal{R} 足够小，其内部的概率密度是相同的，则有

$$P \approx \frac{K}{N}$$

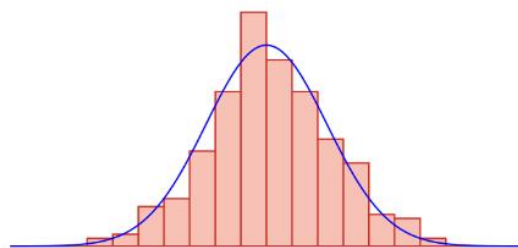
- 结合上述两个公式，得到

$$p(\mathbf{x}) \approx \frac{K}{NV}$$

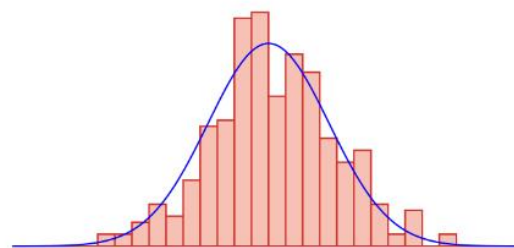
直方图方法 (Histogram Method)

- 一种非常直观的估计连续变量密度函数的方法，可以表示为一种柱状图。
 - 以一维随机变量为例，首先将其取值范围分成 M 个连续的、不重叠的区间，每个区间的宽度为 Δ_m 。给定 N 个训练样本 $\mathcal{D} = \{x^{(n)}\}_{n=1}^N$ ，我们统计这些样本落入每个区间的数量 K_m ，然后将它们归一化为密度函数。

$$p_m = \frac{K_m}{N\Delta_m}, \quad 1 \leq m \leq M$$



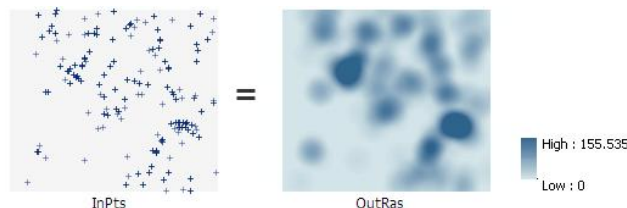
(a) 10 个区间 (bin)



(b) 30 个区间 (bin)

核密度估计 (Kernel Density Estimation)

- 核密度估计，也叫Parzen窗方法，是一种直方图方法的改进。



- 假设 \mathcal{R} 为 d 维空间中的一个以点 \mathbf{x} 为中心的“超立方体”，并定义核函数来表示一个样本 \mathbf{z} 是否落入该超立方体中。其中 H 为超立方体的边长，也称为核函数的宽度

$$\phi\left(\frac{\mathbf{z} - \mathbf{x}}{H}\right) = \begin{cases} 1 & \text{if } |z_i - x_i| < \frac{H}{2}, 1 \leq i \leq D \\ 0 & \text{else} \end{cases}$$

- 点 \mathbf{x} 的密度估计为

$$p(\mathbf{x}) = \frac{K}{NH^D} = \frac{1}{NH^D} \sum_{n=1}^N \phi\left(\frac{\mathbf{x}^{(n)} - \mathbf{x}}{H}\right),$$

K近邻方法

- 核密度估计方法中的核宽度 H 是固定的，因此同一个宽度可能对高密度的区域过大，而对低密度区域过小。
- 一种更灵活的方式是设置一种可变宽度的区域，并使得落入每个区域中样本数量为固定的 K 。
- 要估计点 x 的密度，首先找到一个以 x 为中心的球体，使得落入球体的样本数量为 K ，就可以计算出点 x 的密度。
- 因为落入球体的样本也是离 x 最近的 K 个样本，所以这种方法称为 **K近邻方法**。