



2021年秋季 《机器学习概论》课程

第十一章：特征选择与稀疏学习

主讲：连德富 特任教授 | 博士生导师

邮箱： liandefu@ustc.edu.cn

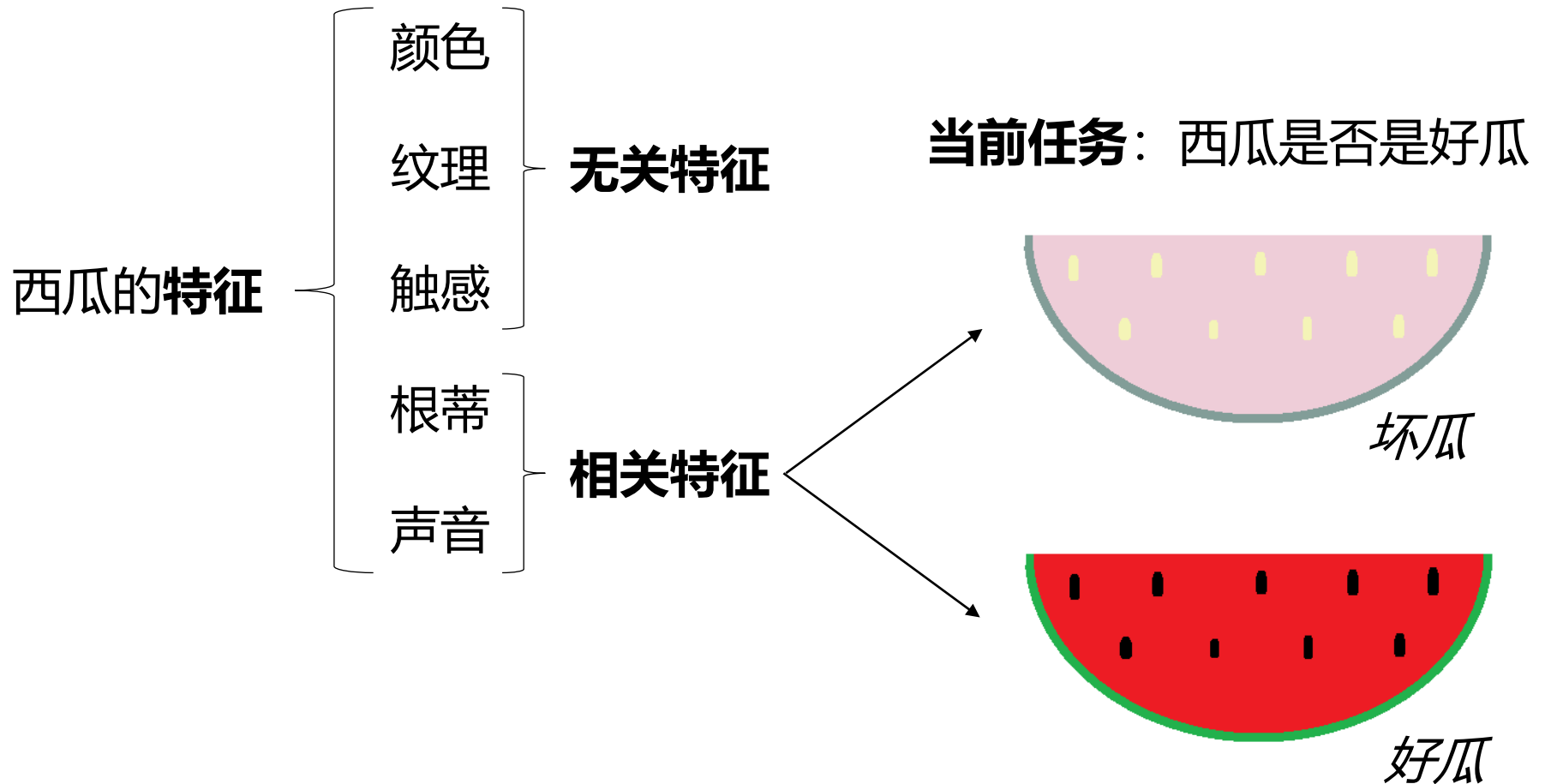
手机：13739227137

主页： <http://staff.ustc.edu.cn/~liandefu>

特征

- 特征
 - 描述物体的属性
- 特征的分类
 - 相关特征: 对**当前学习任务**有用的属性
 - 无关特征: 与**当前学习任务**无关的属性
 - 冗余特征: 其所包含信息能由其他特征推演出来

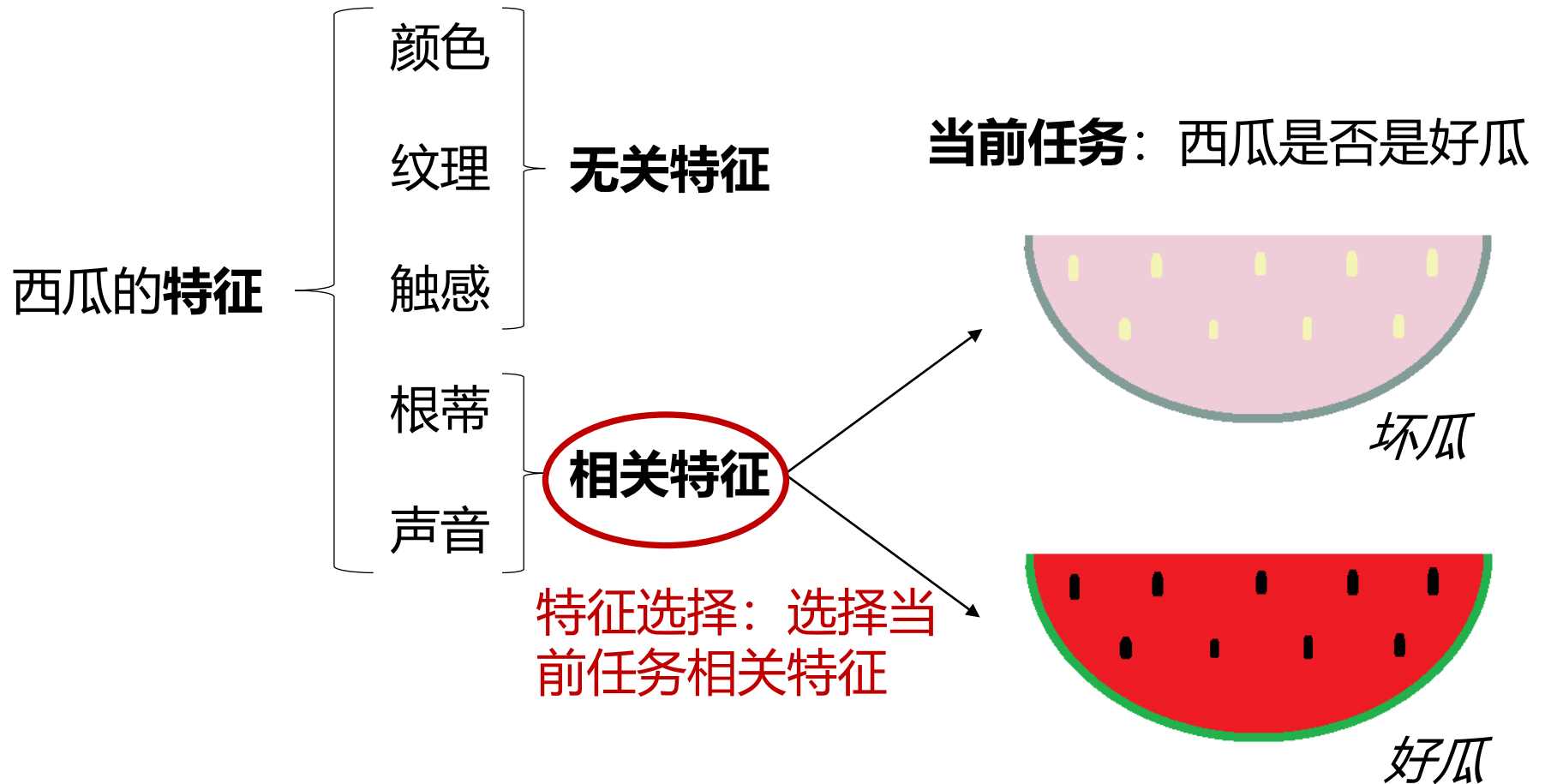
例子：西瓜的特征



特征选择

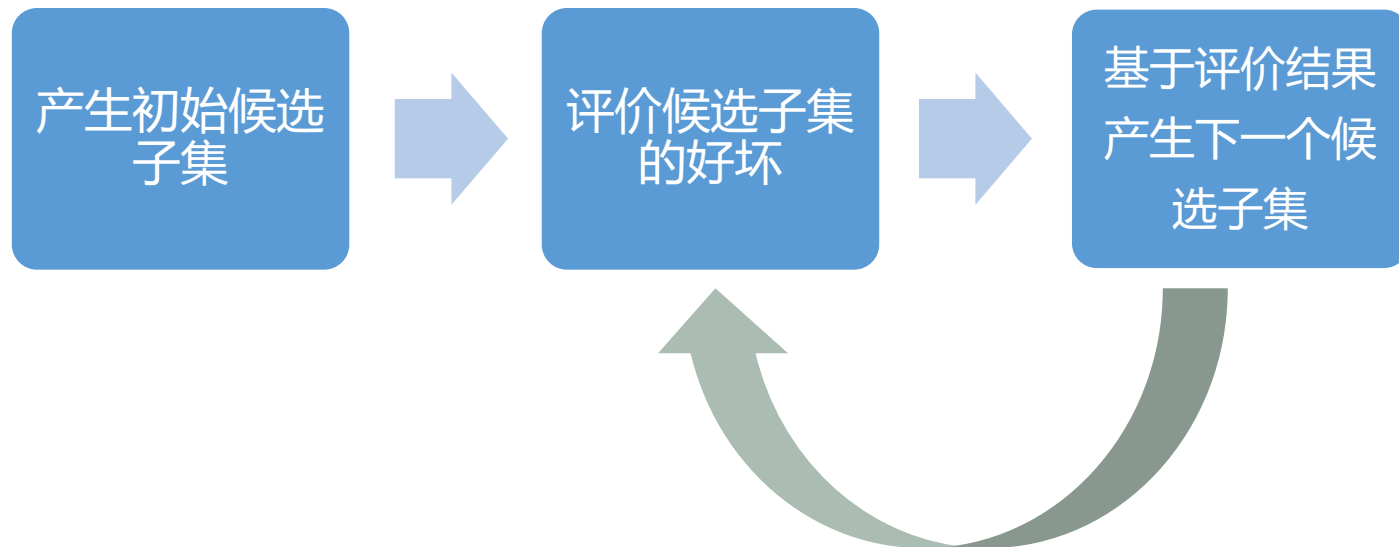
- 特征选择
 - 从给定的特征集合中选出**任务相关**特征子集
 - 必须确保不丢失重要特征
- 原因
 - 减轻维度灾难：在少量属性上构建模型
 - 降低学习难度：留下关键信息

例子：判断是否好瓜时的特征选择



特征选择的一般方法

- 遍历所有可能的子集
 - 计算上遭遇组合爆炸, **不可行**
- 可行方法



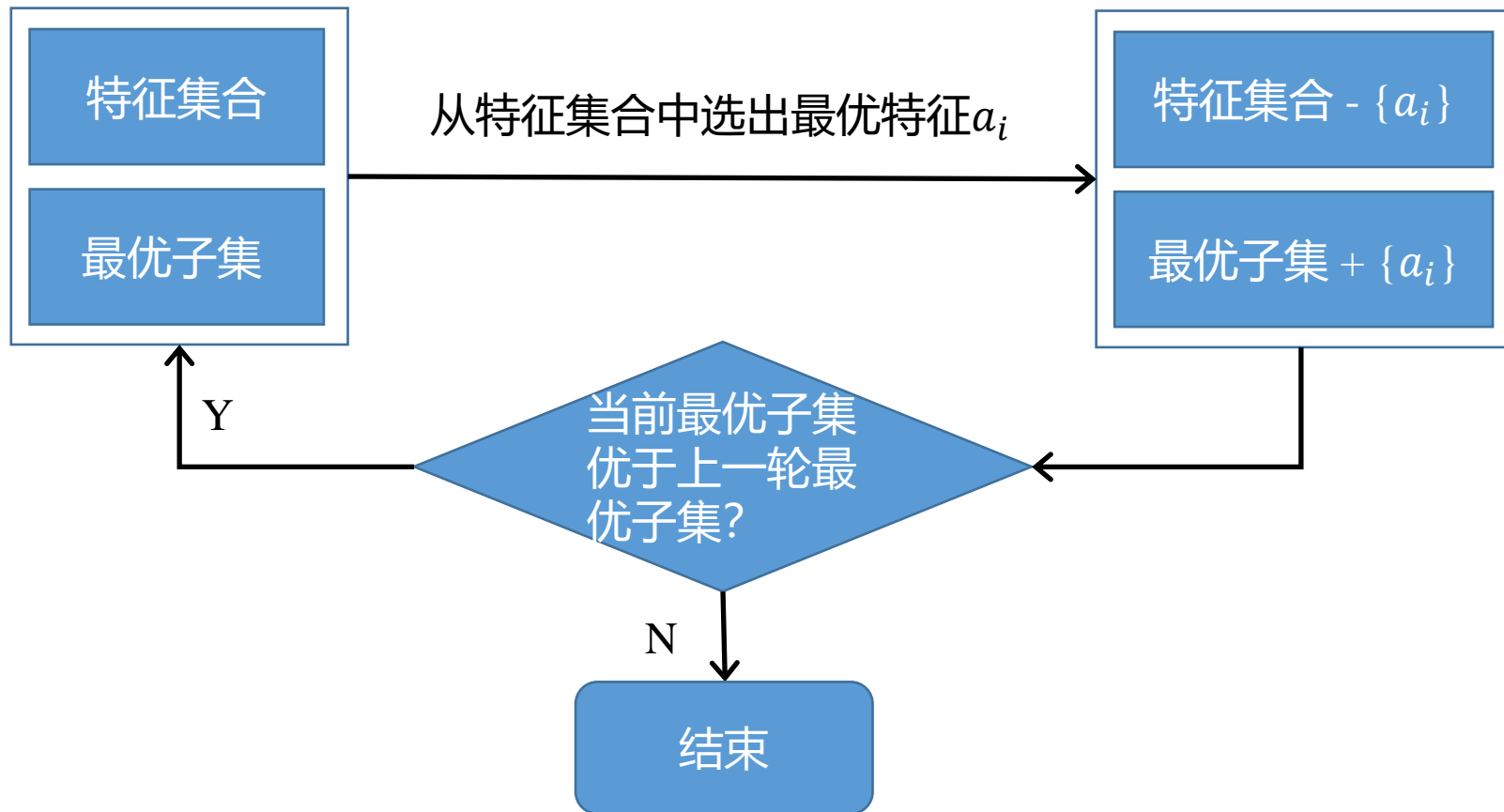
两个关键环节：子集搜索和子集评价

子集搜索

- 用贪心策略选择包含重要信息的特征子集
 - **前向搜索**: 逐渐增加相关特征
 - **后向搜索**: 从完整的特征集合开始, 逐渐减少特征
 - **双向搜索**: 每一轮逐渐增加相关特征, 同时减少无关特征

前向搜索

- 最优子集初始为空集，特征集合初始时包括所有给定特征



子集评价

- 特征子集确定了对数据集的一个划分
 - 每个划分区域对应着特征子集的某种取值
- 样本标记对应着对数据集的真实划分

通过估算这两个划分的差异，就能对特征子集进行评价；与样本标记对应的划分的差异越小，则说明当前特征子集越好

用信息熵进行子集评价

- 特征子集 A 确定了对数据集 D 的一个划分
 - A 上的取值将数据集 D 分为 V 份, 每一份用 D^v 表示
 - $\text{Ent}(D^v)$ 表示 D^v 上的信息熵
- 样本标记 Y 对应着对数据集 D 的真实划分
 - $\text{Ent}(D)$ 表示 D 上的信息熵

D 上的信息熵定义为

$$\text{Ent}(D) = - \sum_{k=1}^{|Y|} p_k \log_2 p_k$$

第 i 类样本所占比例为 p_i

特征子集 A 的信息增益为:

$$\text{Gain}(A) = \text{Ent}(D) - \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Ent}(D^v)$$

常见的特征选择方法

- 将特征子集搜索机制与子集评价机制相结合，即可得到特征选择方法

常见的特征选择方法大致分为如下三类：

- 过滤式
- 包裹式
- 嵌入式

过滤式选择

先用特征选择过程过滤原始数据，再用过滤后的特征来训练模型；
特征选择过程与后续学习器无关

- Relief (Relevant Features) 方法 [Kira and Rendell, 1992]
 - 为每个初始特征赋予一个“**相关统计量**”，度量特征的重要性
 - 特征子集的重要性由子集中每个特征所对应的相关统计量之和决定
 - 设计一个阈值，然后选择比阈值大的相关统计量分量所对应的特征
 - 或者指定欲选取的特征个数，然后选择相关统计量分量最大的指定个数特征

如何确定相关统计量？

Relief方法中相关统计量的确定

- 猜对近邻 (near-hit) : x_i 的同类样本中的最近邻 $x_{i,nh}$
- 猜错近邻 (near-miss) : x_i 的异类样本中的最近邻 $x_{i,nm}$

- 相关统计量对应于属性 j 的分量为

$$\delta^j = \sum_i -\text{diff}(x_i^j, x_{i,nh}^j)^2 + \text{diff}(x_i^j, x_{i,nm}^j)^2$$

若 j 为离散型, 则 $x_a^j = x_b^j$ 时 $\text{diff}(x_a^j, x_b^j) = 0$, 否则为 1 ;
若 j 为连续型, 则 $\text{diff}(x_a^j, x_b^j) = |x_a^j - x_b^j|$, 注意 x_a^j, x_b^j 已规范化到 $[0,1]$ 区间

- 相关统计量越大, 属性 j 上, 猜对近邻比猜错近邻越近, 即属性 j 对区分对错越有用
- Relief方法的时间开销随采样次数以及原始特征数线性增长, 运行效率很高

Relief方法的多类拓展

- Relief方法是为二分类问题设计的，其扩展变体Relief-F [Kononenko, 1994]能处理多分类问题
- 数据集中的样本来自 $|y|$ 个类别，其中 x_i 属于第 k 类
- 猜中近邻：第 k 类中 x_i 的最近邻 $x_{i,nh}$
- 猜错近邻：第 k 类之外的每个类中找到一个 x_i 的最近邻作为猜错近邻，记为 $x_{i,l,nm} (l = 1, 2, \dots, |y|; l \neq k)$
- 相关统计量对应于属性的分量为

$$\delta^j = \sum_i \left(-\text{diff}(x_i^j, x_{i,nh}^j)^2 + \sum_{l \neq k} p_l \times \text{diff}(x_i^j, x_{i,l,nm}^j)^2 \right)$$

p_l 为第 l 类样本在数据集 D 中所占的比例

包裹式选择

- 包裹式选择直接把最终将要使用的学习器的性能作为特征子集的评价准则
 - 包裹式特征选择的目的是为给定学习器选择最有利于其性能、“量身定做”的特征子集
 - 包裹式选择方法直接针对给定学习器进行优化，因此从最终学习器性能来看，包裹式特征选择比过滤式特征选择更好
 - 包裹式特征选择过程中需多次训练学习器，计算开销通常比过滤式特征选择大得多

LVW包裹式特征选择方法

- LVW (Las Vegas Wrapper) [Liu and Setiono, 1996] 在拉斯维加斯方法框架下使用随机策略来进行子集搜索，并以最终分类器的误差作为特征子集评价准则
- 基本步骤
 - 在循环的每一轮随机产生一个特征子集
 - 在随机产生的特征子集上通过交叉验证推断当前特征子集的误差
 - 进行多次循环，在多个随机产生的特征子集中选择误差最小的特征子集作为最终解*

*若有运行时间限制，则该算法有可能给不出解

嵌入式选择

- 嵌入式特征选择是将特征选择过程与学习器训练过程融为一体，两者在同一个优化过程中完成，在学习器训练过程中自动地进行特征选择
- 考虑最简单的线性回归模型，以平方误差为损失函数，并引入 L_2 范数正则化项防止过拟合，则有

$$\min_{\mathbf{w}} \sum_{i=1}^m (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|_2^2$$

岭回归 (ridge regression)
[Tikhonov and Arsenin, 1977]

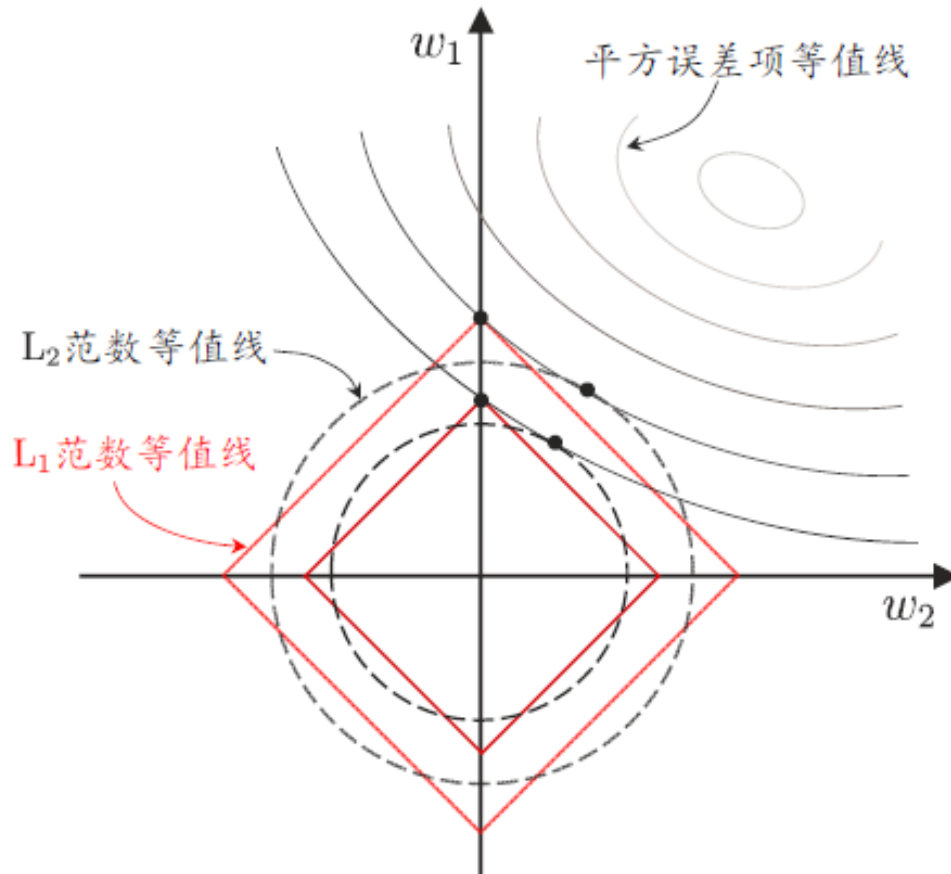
- 将 L_2 范数替换为 L_1 范数，则有**LASSO** [Tibshirani, 1996]

$$\min_{\mathbf{w}} \sum_{i=1}^m (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|_1$$

易获得稀疏解，是一种嵌入式特征选择方法

$$\|\mathbf{w}\|_1 = \sum_d |w_d|$$

使用 L_1 范数正则化易获得稀疏解



等值线即取值相同的点的连线

$$\|\mathbf{w}\|_1 = c$$

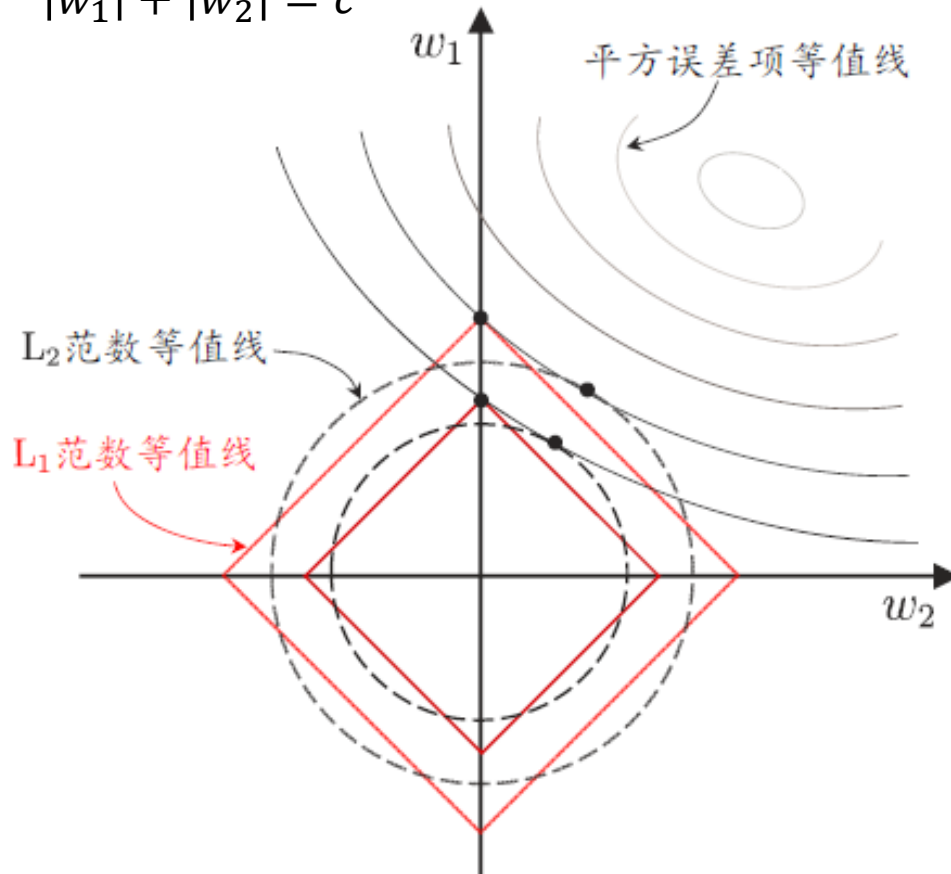


$$|w_1| + |w_2| = c$$

$$\begin{cases} w_1 + w_2 = c, & \text{if } w_1 \geq 0, w_2 \geq 0 \\ w_1 - w_2 = c, & \text{if } w_1 \geq 0, w_2 < 0 \\ -w_1 + w_2 = c, & \text{if } w_1 < 0, w_2 \geq 0 \\ -w_1 - w_2 = c, & \text{if } w_1 < 0, w_2 < 0 \end{cases}$$

使用 L_1 范数正则化易获得稀疏解

$$|w_1| + |w_2| = c$$



等值线即取值相同的点的连线

- 假设 x 仅有两个属性, 那么 w 有两个分量 w_1 和 w_2 . 那么目标优化的解要在平方误差项与正则化项之间折中, **即出现在图中平方误差项等值线与正则化等值线相交处.**

从图中看出, 采用 L_1 范数时交点常出现在坐标轴上, 即产生 w_1 或者 w_2 为0的稀疏解

L_1 正则化问题的求解(1)

- 可使用近端梯度下降(**Proximal Gradient Descend**, 简称**PGD**) 求解

[Boyd and Vandenberghe, 2004]

- 考虑更一般的问题

$$\min_{\mathbf{w}} f(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$$

- 假设 $f(\mathbf{w})$ 为凸函数, 且 $\nabla f(\mathbf{w})$ 满足L-Lipschitz条件, 即存在常数 $L > 0$, 使得

$$\forall \mathbf{w}, \mathbf{w}' \quad \|\nabla f(\mathbf{w}') - \nabla f(\mathbf{w})\| \leq L \|\mathbf{w}' - \mathbf{w}\|$$

L_1 正则化问题的求解(1)

- 假设 $f(\mathbf{x})$ 为凸函数, 且 $\nabla f(\mathbf{x})$ 满足L-Lipschitz条件, 即存在常数 $L > 0$, 使得

$$\forall \mathbf{w}, \mathbf{w}' \quad \|\nabla f(\mathbf{w}') - \nabla f(\mathbf{w})\| \leq L \|\mathbf{w}' - \mathbf{w}\|$$

- 等价于 $g(\mathbf{w}) = \frac{L}{2} \|\mathbf{w}\|^2 - f(\mathbf{w})$ 为凸函数

$$f(\mathbf{w}') \geq f(\mathbf{w}) + \nabla f(\mathbf{w})^\top (\mathbf{w}' - \mathbf{w}) \quad + \quad f(\mathbf{w}) \geq f(\mathbf{w}') + \nabla f(\mathbf{w}')^\top (\mathbf{w} - \mathbf{w}')$$

$$\Rightarrow (\nabla f(\mathbf{w}) - \nabla f(\mathbf{w}'))^\top (\mathbf{w} - \mathbf{w}') \geq 0$$

$$\Rightarrow (\nabla f(\mathbf{w}) - \nabla f(\mathbf{w}'))^\top (\mathbf{w} - \mathbf{w}') \leq L \|\mathbf{w} - \mathbf{w}'\|^2$$

$$\begin{aligned} (\nabla g(\mathbf{w}') - \nabla g(\mathbf{w}))^\top (\mathbf{w}' - \mathbf{w}) &= \left(L(\mathbf{w}' - \mathbf{w}) - (\nabla f(\mathbf{w}') - \nabla f(\mathbf{w})) \right)^\top (\mathbf{w}' - \mathbf{w}) \\ &\geq 0 \end{aligned}$$

L_1 正则化问题的求解(1)

- 假设 $f(\mathbf{x})$ 为凸函数, 且 $\nabla f(\mathbf{x})$ 满足L-Lipschitz条件, 即存在常数 $L > 0$, 使得

$$\forall \mathbf{w}, \mathbf{w}' \quad \|\nabla f(\mathbf{w}') - \nabla f(\mathbf{w})\| \leq L \|\mathbf{w}' - \mathbf{w}\|$$

- 等价于 $g(\mathbf{w}) = \frac{L}{2} \|\mathbf{w}\|^2 - f(\mathbf{w})$ 为凸函数
- 也等价于 $f(\mathbf{w}') \leq f(\mathbf{w}) + \nabla f(\mathbf{w})(\mathbf{w}' - \mathbf{w}) + \frac{L}{2} \|\mathbf{w}' - \mathbf{w}\|^2$

$$g(\mathbf{w}') \geq g(\mathbf{w}) + \nabla g(\mathbf{w})(\mathbf{w}' - \mathbf{w})$$

$$\frac{L}{2} \|\mathbf{w}'\|^2 - f(\mathbf{w}') \geq \frac{L}{2} \|\mathbf{w}\|^2 - f(\mathbf{w}) + (L\mathbf{w} - \nabla f(\mathbf{w}))^\top (\mathbf{w}' - \mathbf{w})$$

$$f(\mathbf{w}') \leq f(\mathbf{w}) + \nabla f(\mathbf{w})^\top (\mathbf{w}' - \mathbf{w}) - L\mathbf{w}(\mathbf{w}' - \mathbf{w}) + \frac{L}{2} \|\mathbf{w}'\|^2 - \frac{L}{2} \|\mathbf{w}\|^2$$

$$f(\mathbf{w}') \leq f(\mathbf{w}) + \nabla f(\mathbf{w})(\mathbf{w}' - \mathbf{w}) + \frac{L}{2} \|\mathbf{w}' - \mathbf{w}\|^2$$

L_1 正则化问题的求解(1)

- 假设 $f(\mathbf{x})$ 为凸函数, 且 $\nabla f(\mathbf{x})$ 满足L-Lipschitz条件, 即存在常数 $L > 0$, 使得

$$\forall \mathbf{w}, \mathbf{w}' \quad \|\nabla f(\mathbf{w}') - \nabla f(\mathbf{w})\| \leq L \|\mathbf{w}' - \mathbf{w}\|$$

- 等价于 $g(\mathbf{w}) = \frac{L}{2} \|\mathbf{w}\|^2 - f(\mathbf{w})$ 为凸函数
- 也等价于 $f(\mathbf{w}') \leq f(\mathbf{w}) + \nabla f(\mathbf{w})(\mathbf{w}' - \mathbf{w}) + \frac{L}{2} \|\mathbf{w}' - \mathbf{w}\|^2$

$$= \frac{L}{2} \left\| \mathbf{w}' - \left(\mathbf{w} - \frac{1}{L} \nabla f(\mathbf{w}) \right) \right\|_2^2 + \text{const}$$

$$f(\mathbf{w}) + \lambda \|\mathbf{w}\|_1 \leq \frac{L}{2} \left\| \mathbf{w} - \left(\mathbf{w}^k - \frac{1}{L} \nabla f(\mathbf{w}^k) \right) \right\|_2^2 + \lambda \|\mathbf{w}\|_1 + \text{const}$$

L_1 正则化问题的求解(1)

$$\bullet f(\mathbf{w}) + \lambda \|\mathbf{w}\|_1 \leq \frac{L}{2} \left\| \mathbf{w} - \left(\mathbf{w}^k - \frac{1}{L} \nabla f(\mathbf{w}^k) \right) \right\|_2^2 + \lambda \|\mathbf{w}\|_1 + \text{const}$$

$$\mathbf{w}^{k+1} = \operatorname{argmin}_{\mathbf{w}} \frac{L}{2} \left\| \mathbf{w} - \left(\mathbf{w}^k - \frac{1}{L} \nabla f(\mathbf{w}^k) \right) \right\|_2^2 + \lambda \|\mathbf{w}\|_1$$

$$\mathbf{w}^{k+1} = \operatorname{argmin}_{\mathbf{w}} \frac{L}{2} \left\| \mathbf{w} - \mathbf{z}^k \right\|_2^2 + \lambda \|\mathbf{w}\|_1$$

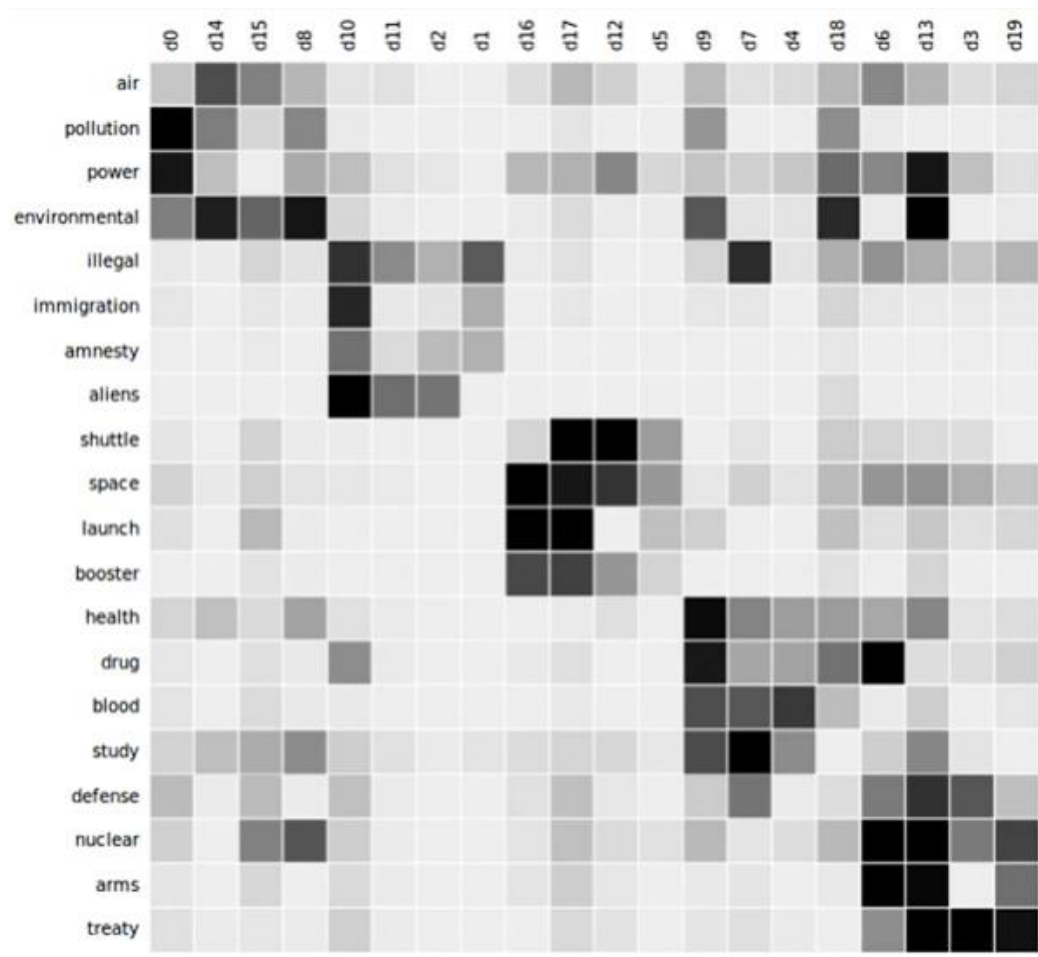
$$[\mathbf{w}^{k+1}]_i = \operatorname{argmin}_{w_i} \frac{L}{2} (w_i - z_i^k)^2 + \lambda |w_i|$$

$$\begin{aligned} & \frac{L}{2} (w_i - z_i^k)^2 + \lambda |w_i| \\ = & \begin{cases} \frac{L}{2} \left(w_i - z_i^k + \frac{\lambda}{L} \right)^2 + \text{const}, & \text{if } w_i \geq 0 \\ \frac{L}{2} \left(w_i - z_i^k - \frac{\lambda}{L} \right)^2 + \text{const}, & \text{if } w_i < 0 \end{cases} \end{aligned}$$

$$w_i = \begin{cases} z_i^k - \frac{\lambda}{L}, & \text{if } \frac{\lambda}{L} < z_i^k \\ 0, & \text{if } \frac{\lambda}{L} \geq |z_i^k| \\ z_i^k + \frac{\lambda}{L}, & \text{if } z_i^k < -\frac{\lambda}{L} \end{cases}$$

稀疏表示与字典学习

- 将数据集考虑成一个矩阵，每行对应一个样本，每列对应一个特征
- 矩阵中有很多零元素，且非整行整列出现
- 稀疏表达的优势：
 - 文本数据线性可分
 - 存储高效



能否将稠密表示的数据集转化为“稀疏表示”，
使其享受稀疏表达的优势？

稀疏表示与字典学习

- 一般的学习任务中，并没有字典可用，需学习出这样一个字典
- 为普通稠密表达的样本找到合适的字典，将样本转化为稀疏表示，这一过程称为字典学习
- 给定数据集 $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}, \mathbf{x}_i \in \mathbb{R}^d$
- 最简单的字典学习的优化形式为

$$\min_{\mathbf{B}, \alpha_i} \sum_{i=1}^m \|\mathbf{x}_i - \mathbf{B}\alpha_i\|_2^2 + \lambda \sum_{i=1}^m \|\alpha_i\|_1$$

采用交替迭代优化进行求解

$\mathbf{B} \in \mathbb{R}^{d \times k}$ 为字典矩阵、 $\alpha_i \in \mathbb{R}^k$ 为样本的稀疏表示

k 称为字典的词汇量，通常由用户指定

字典学习的解法(1)

- 固定字典 \mathbf{B} , 参考LASSO的方法求解

$$\min_{\alpha_i} \sum_{i=1}^m \|\mathbf{x}_i - \mathbf{B}\alpha_i\|_2^2 + \lambda \sum_{i=1}^m \|\alpha_i\|_1$$

- 以 α_i 为初值更新字典 \mathbf{B}

$$\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m] \in \mathbb{R}^{d \times m}$$

$$\min_{\mathbf{B}} L = \sum_{i=1}^m \|\mathbf{x}_i - \mathbf{B}\alpha_i\|_2^2 = \|\mathbf{X} - \mathbf{B}\mathbf{A}\|_F^2$$

$$\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m] \in \mathbb{R}^{k \times m}$$

$$\nabla_{\mathbf{B}} L = -2(\mathbf{X} - \mathbf{B}\mathbf{A})\mathbf{A}^\top = 0 \quad \Rightarrow \quad \mathbf{B} = \mathbf{X}\mathbf{A}^\top(\mathbf{A}\mathbf{A}^\top)^{-1}$$

字典学习的解法(1)

- 当k很大的时候，求逆计算代价高，可以通过逐列更新

$$L = \|X - BA\|_F^2 = \left\| X - \sum_{i=1}^k b_i a_{(i)}^\top \right\|_F^2$$

b_i 表示字典 B 的第*i*列
 $a_{(i)}$ 表示矩阵 A 的第*i*行

$$= \left\| \left(X - \sum_{j \neq i} b_j a_{(j)}^\top \right) - b_i a_{(i)}^\top \right\|_F^2$$

$$= \|E_i - b_i a_{(i)}^\top\|_F^2$$

$$\nabla_{b_i} L = -2(E_i - b_i a_{(i)}^\top) a_{(i)} = 0 \quad \Rightarrow \quad b_i = \frac{E_i a_{(i)}}{a_{(i)}^\top a_{(i)}}$$

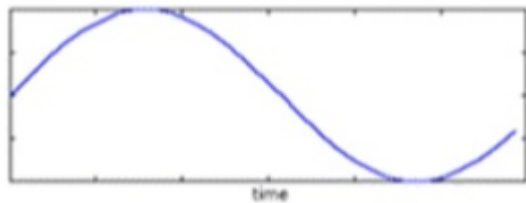
反复迭代以获得最优解

压缩感知

- 能否利用部分数据恢复全部数据？
- 数据传输中，能否利用接收到的压缩、丢包后的数字信号，精确重构出原信号？
- 压缩感知 (compressive sensing) [Cándes et al., 2006, Donoho, 2006] 为解决此类问题提供了新的思路。

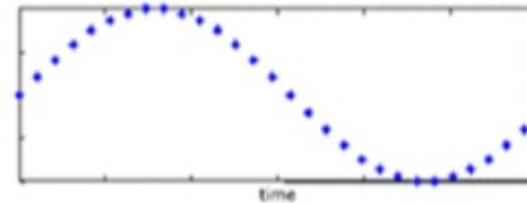
针对信号采样的技术，它通过一些手段，实现了“压缩的采样”，准确说是在采样过程中完成了数据压缩的过程。

模拟信号采样

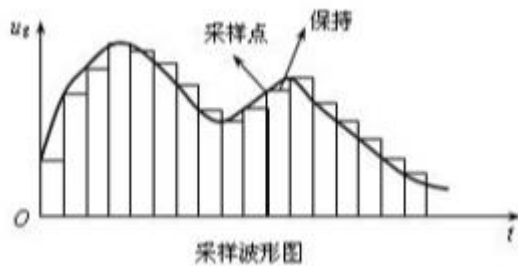


模拟信号

采样

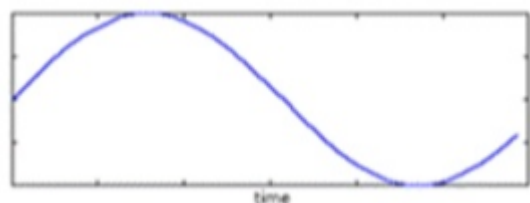


数字信号



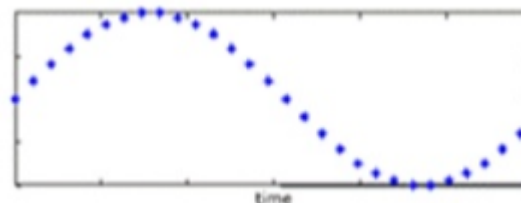
用多大的采样频率，才能完全恢复模拟信号？

模拟信号采样

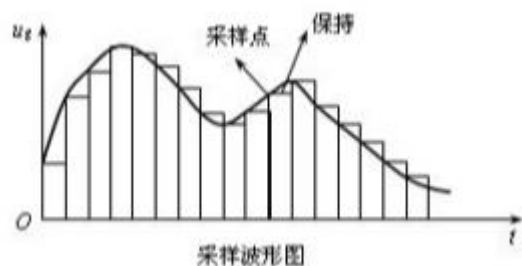


模拟信号

采样



数字信号



采样波形图

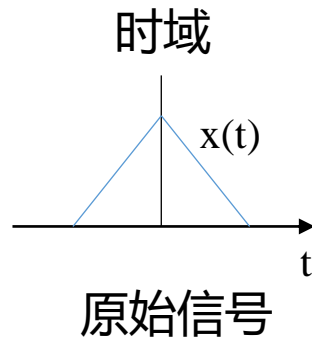
奈奎斯特采样定理

为了不失真地恢复模拟信号，
采样频率应该大于模拟信号频谱中最高频率的2倍。

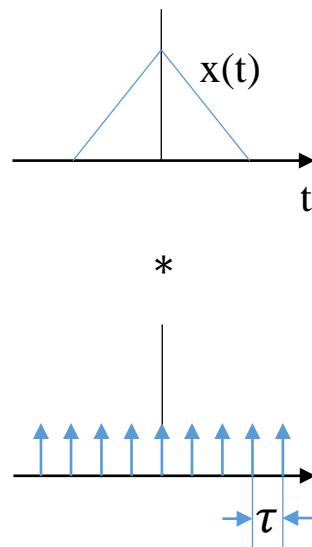
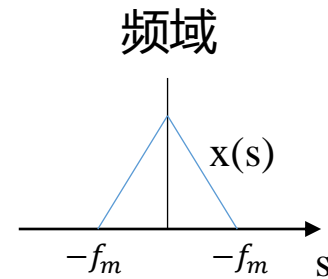


奈奎斯特
1889-1976

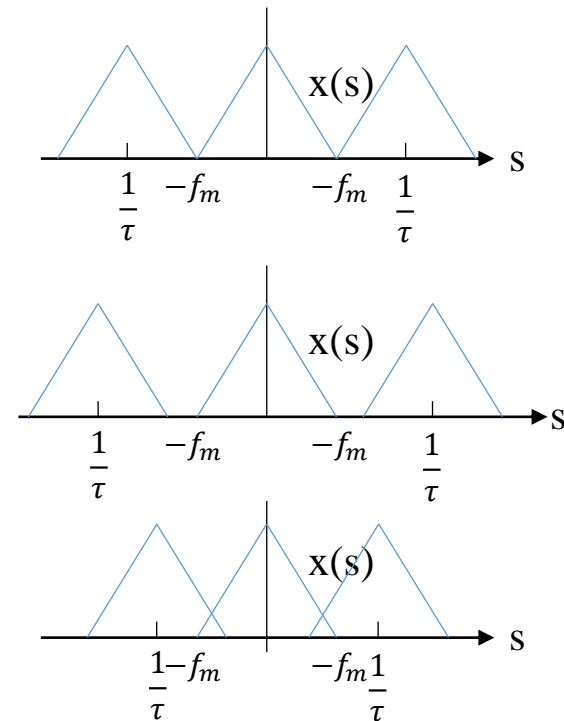
模拟信号采样



傅里叶变换



傅里叶变换



混叠

时域以 τ 为间隔进行采样，频域会以 $1/\tau$ 为周期发生周期延拓

压缩感知

- 2004年，陶哲轩等人证明，如果信号是稀疏的，那么可以由远低于采样定理要求的采样点重建恢复，并与2007年正式提出压缩感知这个概念



陶哲轩



Emmanuel Candes



David Donoho

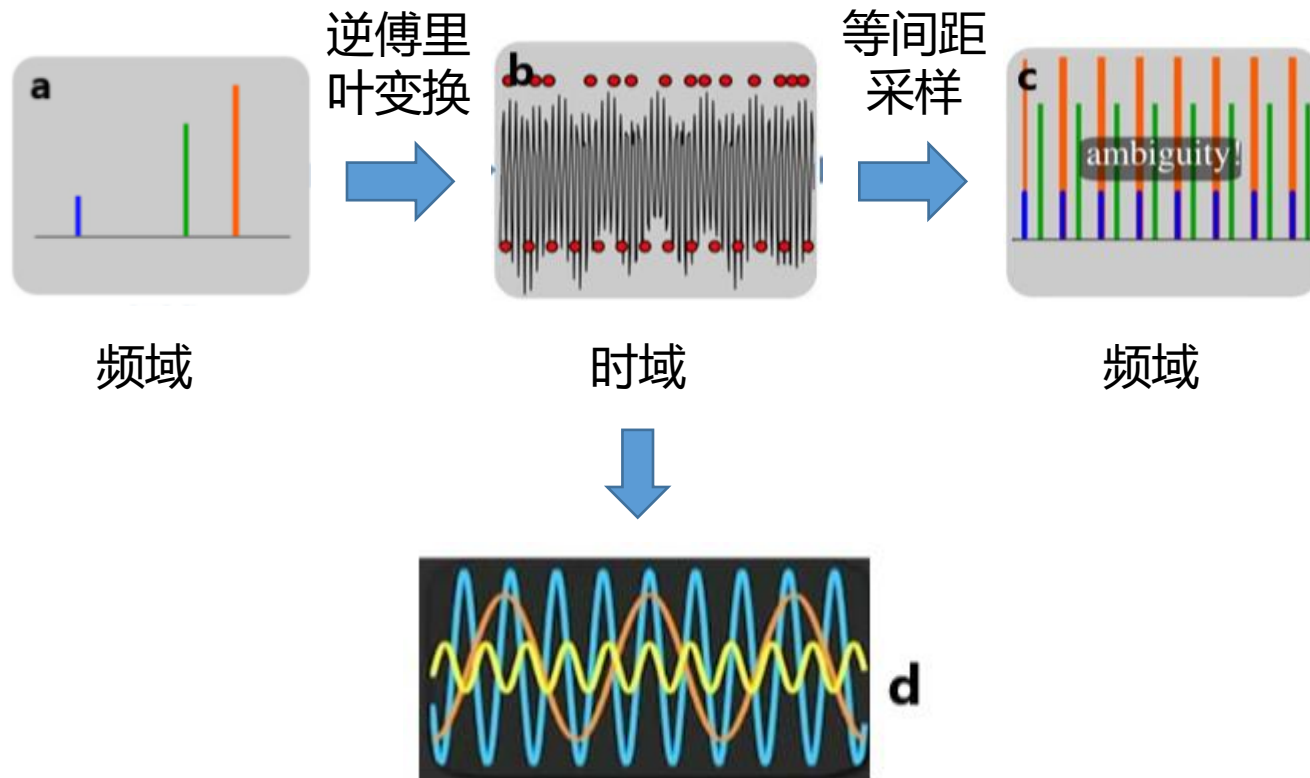
压缩感知

- 采样频率应该大于模拟信号频谱中最高频率的2倍

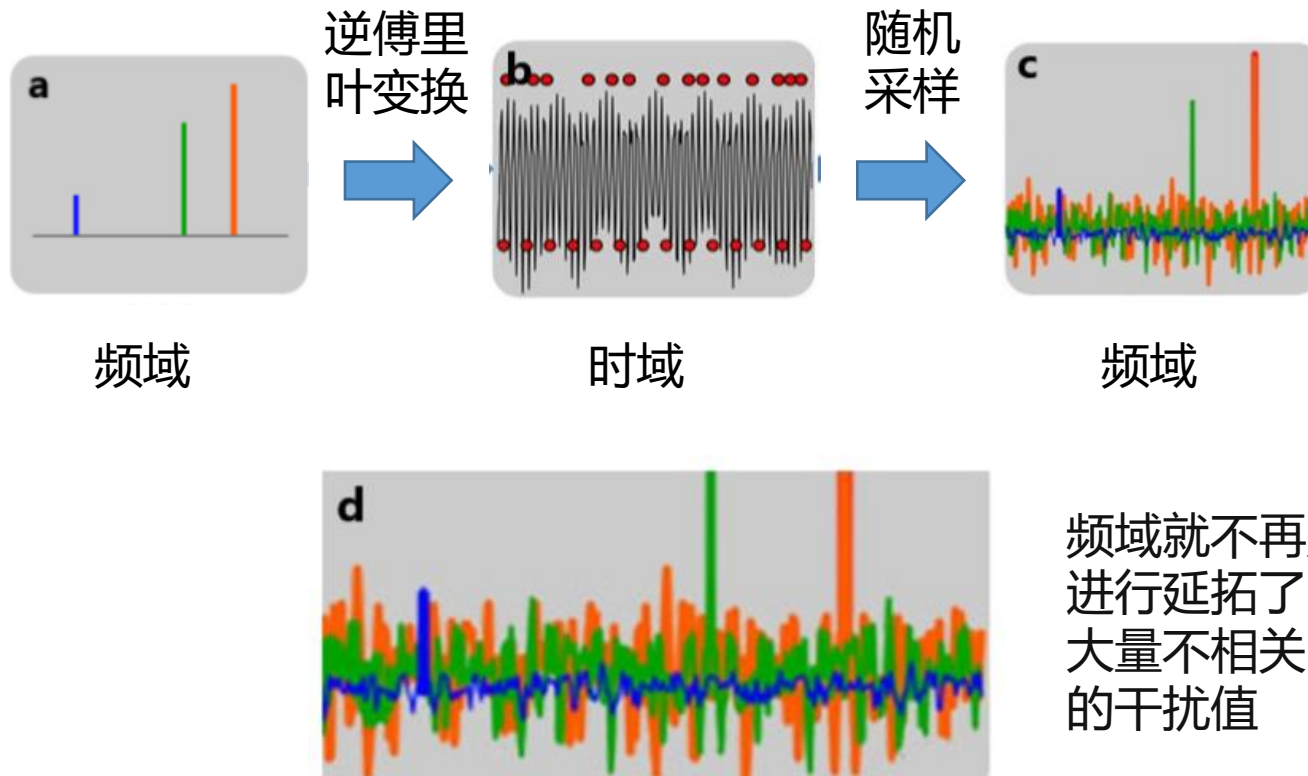
奈奎斯特采样定理

- 给定采样频率采样，意味着等间距采样
- 等间距采样，频域将以以 $1/\tau$ 为周期延拓，采样频率低的时候会一起混叠
- 如果是不等间距采样或者随机采样呢？

压缩感知

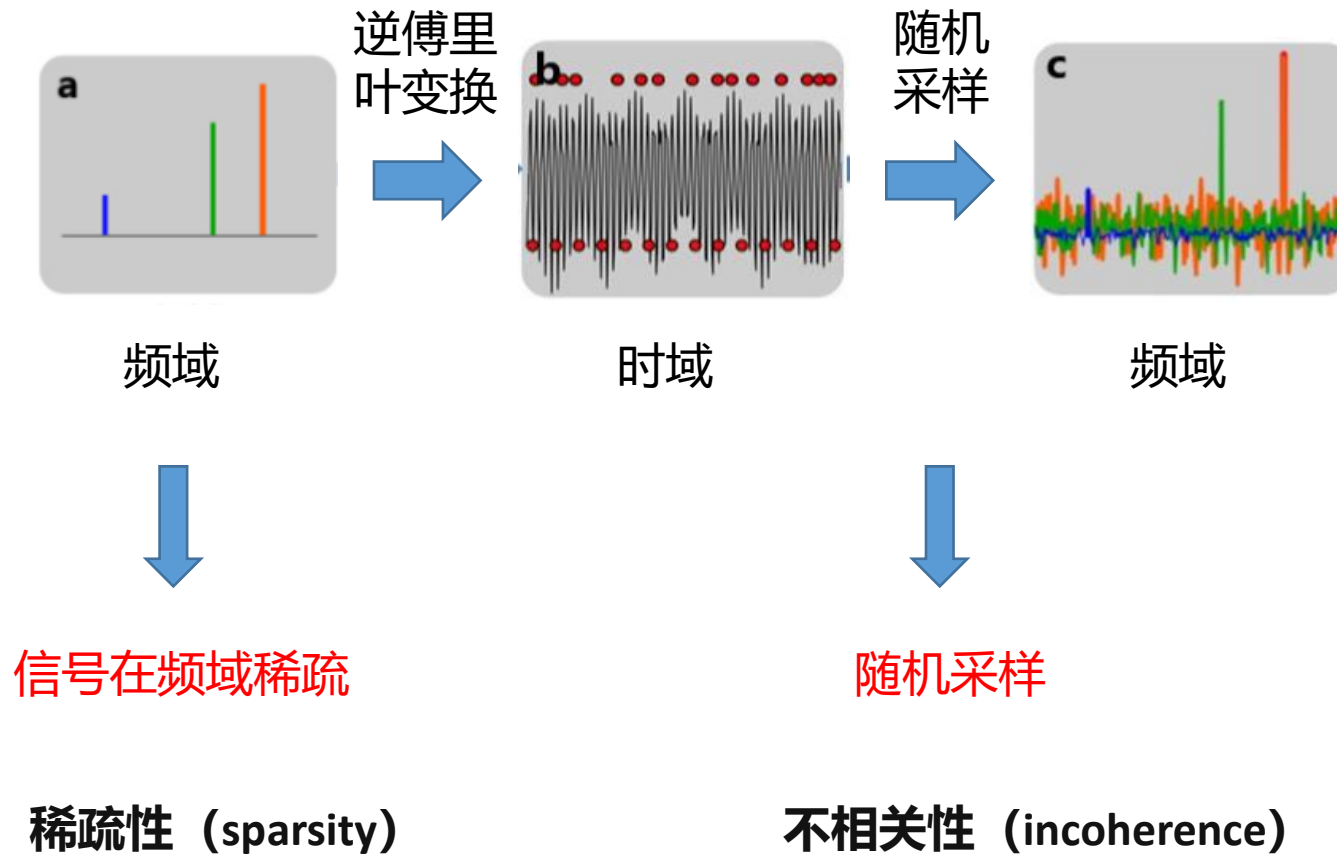


压缩感知



随机采样使得频谱不再是整齐地搬移，而是一小部分一小部分胡乱地搬移

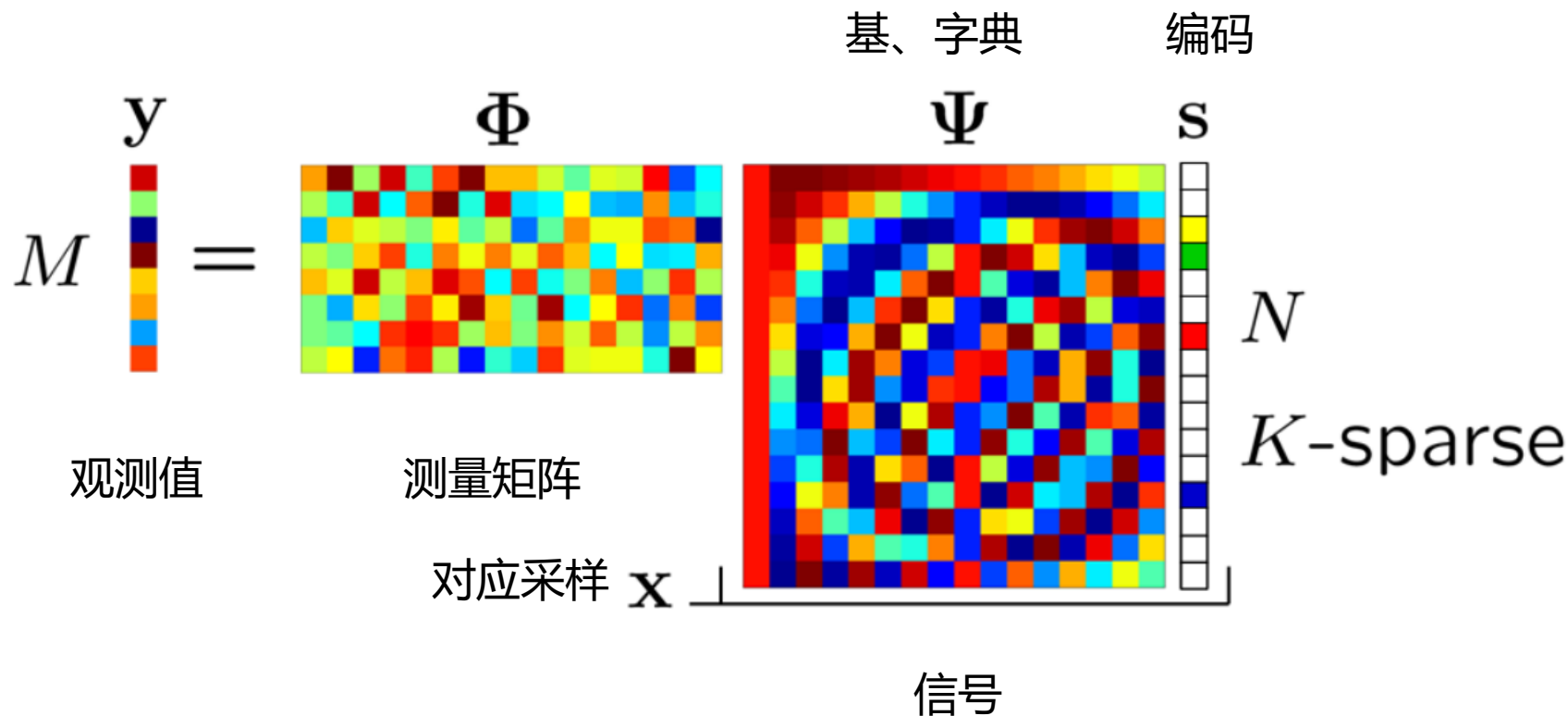
压缩感知—前提条件



压缩感知

一般的自然信号 x 本身并不是稀疏的，需要在某种稀疏基上进行稀疏表示

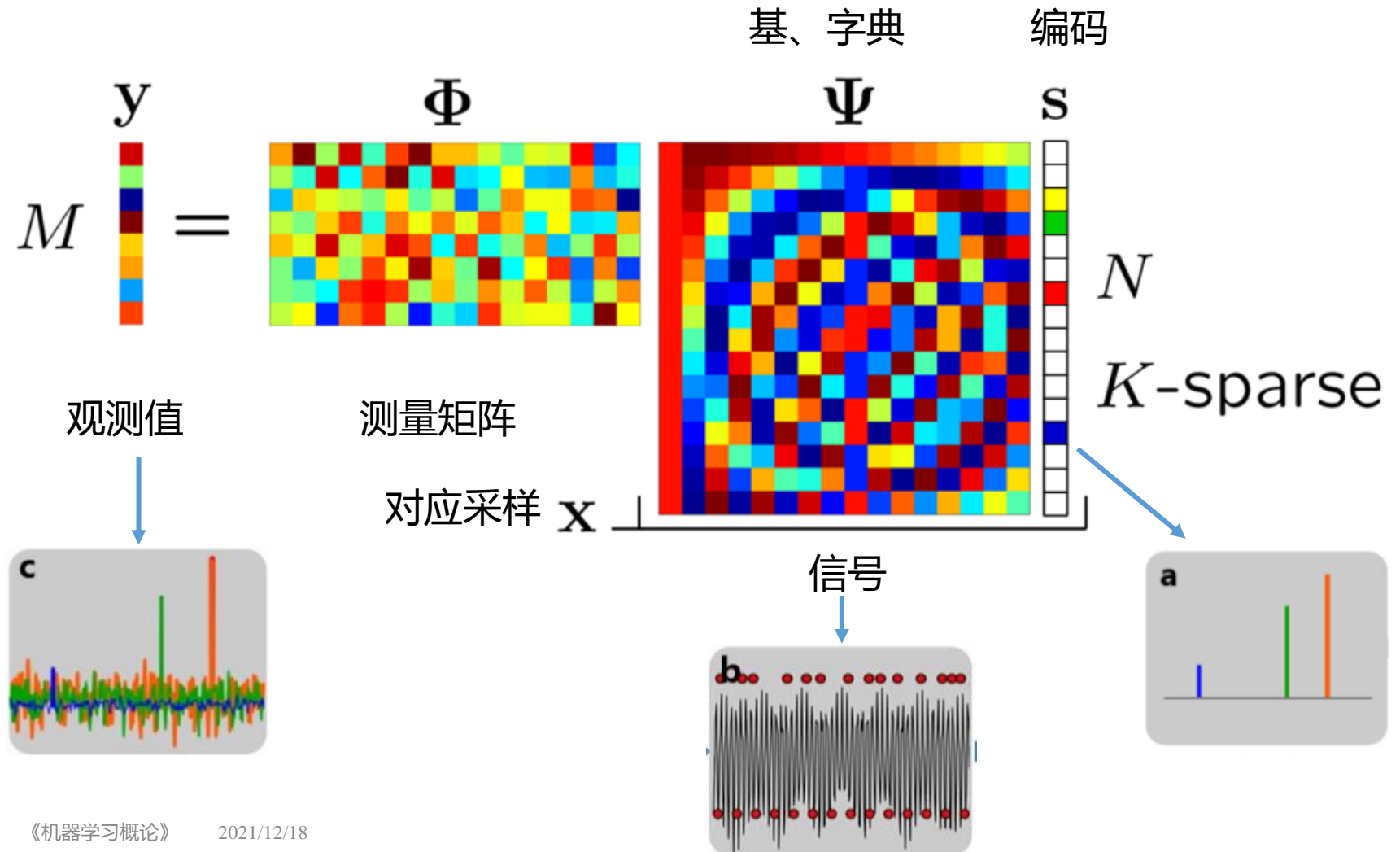
对 x 在基 Ψ 上进行稀疏表示， $x = \Psi s$



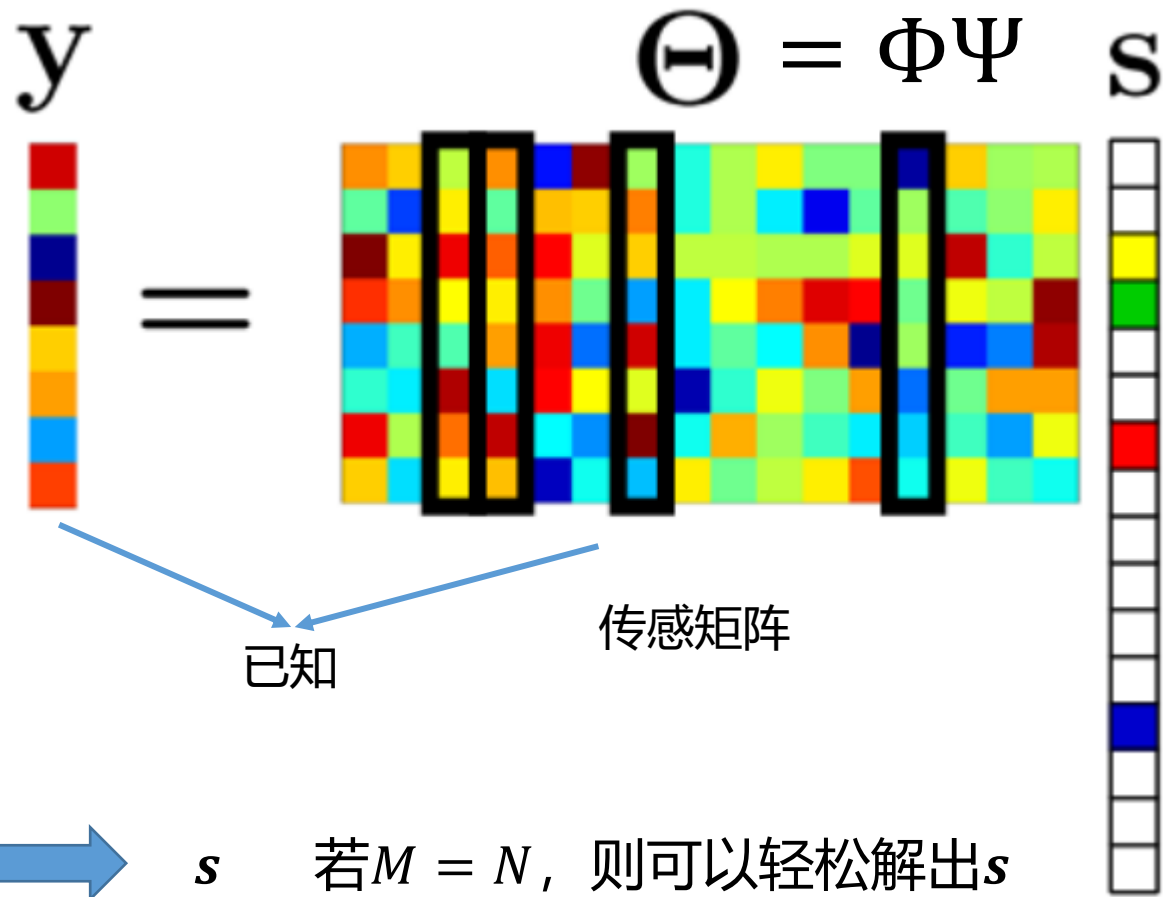
测量矩阵将高维信号投影到低维空间

压缩感知问题就是在已知测量值 y 和测量矩阵 Φ 的基础上，求解欠定方程组 $y = \Phi x$ 得到原信号 x

压缩感知



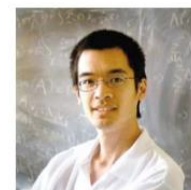
压缩感知



RIP特性— k -RIP

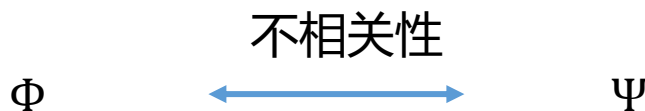
- 存在某个 $\epsilon > 0$, 对任意含有 k 个非零元素的向量 v ,

$$\text{满足 } 1 - \epsilon \leq \frac{\|\Theta v\|_2}{\|v\|_2} \leq 1 + \epsilon,$$



矩阵 Θ 必须维持向量 v 的长度只发生小量变化

- Baraniuk证明: RIP的等价条件是测量矩阵和稀疏基的不相关性



陶哲轩和Candès又证明:

独立同分布的高斯随机测量矩阵可以成为普适的压缩感知测量矩阵

压缩感知的优化目标和解法

- 若 \mathbf{A} 满足 k 限定等距性, 则可通过下面的优化问题近乎完美地从 \mathbf{y} 中恢复出稀疏信号 \mathbf{s} , 进而恢复出 \mathbf{x} :

$$\min_{\mathbf{s}} \|\mathbf{s}\|_0 \quad s.t. \mathbf{y} = \mathbf{\Theta} \mathbf{s}$$

- L_0 范数最小化是NP难问题。不过, L_1 范数最小化在一定条件下与 L_0 范数最小化问题共解 [Cándes et al., 2006]: 将上式转化为共解的 L_1 范数最小化问题

$$\min_{\mathbf{s}} \|\mathbf{s}\|_1 \quad s.t. \mathbf{y} = \mathbf{\Theta} \mathbf{s}$$

- 转化为LASSO的等价形式再通过近端梯度下降法求解, 即使用“基追踪去噪” (Basis Pursuit De-Noising) [Chen et al., 1998]

矩阵补全

客户对书籍的喜好程度的评分

	《笑傲江湖》	《万历十五年》	《人间词话》	《云海玉弓缘》	《人类的故事》
赵大	5	?	?	3	2
钱二	?	5	3	?	5
孙三	5	3	?	?	?
李四	3	?	5	4	?

- 能否将表中已经通过读者评价得到的数据当作部分信号，基于压缩感知的思想恢复出完整信号从而进行书籍推荐呢？从题材、作者、装帧等角度看（相似题材的书籍有相似的读者），表中反映的信号是稀疏的，能通过类似压缩感知的思想加以处理。

“矩阵补全” 技术解决此类问题

矩阵补全的优化问题和解法

- 矩阵补全 (matrix completion) 技术的优化形式为

$$\begin{aligned} \min_{\mathbf{X}} \text{rank}(\mathbf{X}) \\ \text{s. t. } X_{ij} = A_{ij}, (i, j) \in \Omega \end{aligned}$$

约束表明, 恢复出的矩阵中 X_{ij} 应当与已观测到的对应元素相同

- \mathbf{X} : 需要恢复的稀疏信号
- $\text{rank}(\mathbf{X})$: \mathbf{X} 的秩
- \mathbf{A} : 已观测信号
- Ω : \mathbf{A} 中已观测到的元素的位置下标的集合

- NP难问题. 将 $\text{rank}(\mathbf{X})$ 转化为其凸包 “核范数” (nuclear norm)

$$\begin{aligned} \min_{\mathbf{X}} \|\mathbf{X}\|_* \\ \text{s. t. } X_{ij} = A_{ij}, (i, j) \in \Omega \end{aligned} \quad \|\mathbf{X}\|_* = \sum_{j=1}^{\min\{m,n\}} \sigma_j(\mathbf{X})$$

$\sigma_j(\mathbf{X})$ 为 \mathbf{X} 的奇异值
 $\|\mathbf{X}\|_*$ 矩阵的核范数为矩阵的奇异值之和

- 凸优化问题, 通过半正定规划求解 (SDP, Semi-Definite Programming)

满足一定条件时, 只需观察到 $O(mr \log^2 m)$ 个元素就能完美恢复 \mathbf{A}
[Recht, 2011]

作业

- 11.5
- 11.7
- PPT 20页：证明回归和对率回归的损失函数的梯度是否满足L-Lipschitz条件，并求出L

阅读材料

- Lee, D. D., & Seung, H. S. (1999). Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755), 788-791.
- Hinton, G. E., & Salakhutdinov, R. R. (2006). Reducing the dimensionality of data with neural networks. *science*, 313(5786), 504-507.
- Rodriguez, A., & Laio, A. (2014). Clustering by fast search and find of density peaks. *Science*, 344(6191), 1492-1496.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., ... & Dieleman, S. (2016). Mastering the game of Go with deep neural networks and tree search. *nature*, 529(7587), 484-489.
- Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., ... & Chen, Y. (2017). Mastering the game of go without human knowledge. *nature*, 550(7676), 354-359.
- Zhu, J., Wen, C., Zhu, J., Zhang, H., & Wang, X. (2020). A polynomial algorithm for best-subset selection problem. *Proceedings of the National Academy of Sciences*.