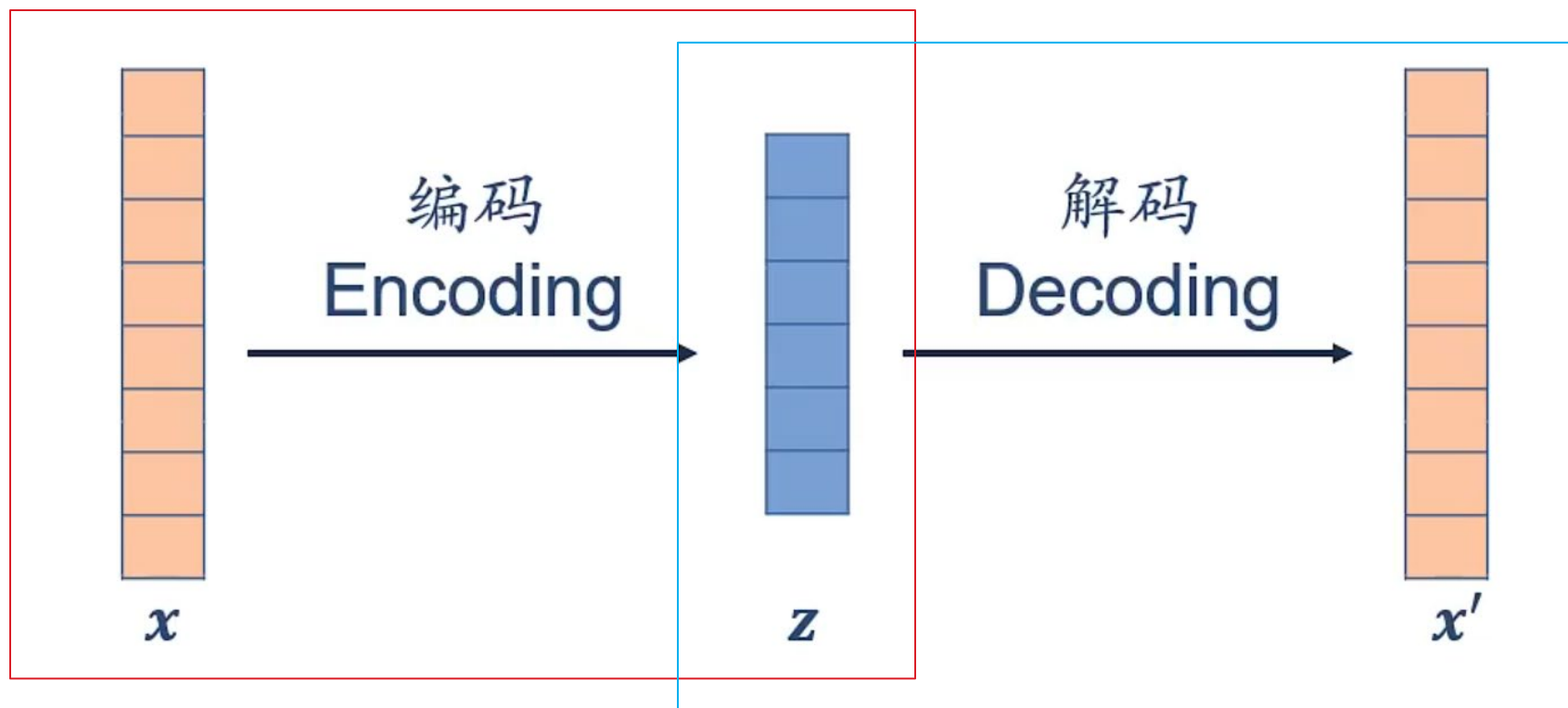




无监督特征学习

特征学习/编码

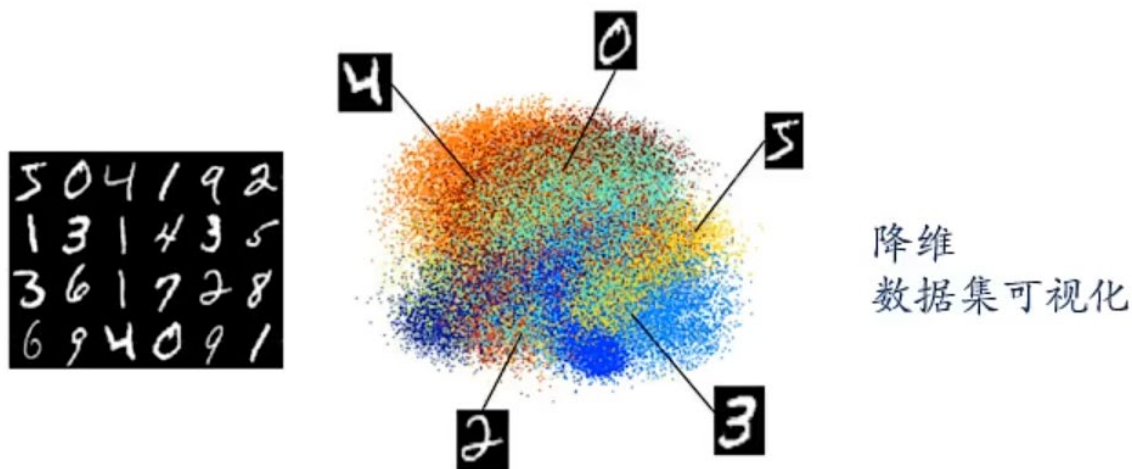


从无标签数据中学习有用的特征（表示）

特征学习的应用

► 目的：特征提取、去噪、降维、数据可视化

► 降维/可视化



► 稀疏编码

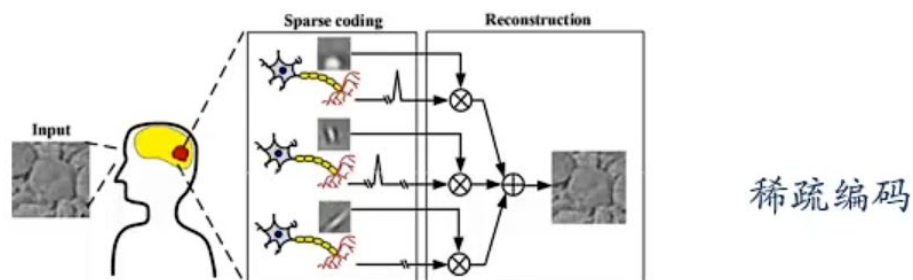


Fig. 1. Sparse coding mimicking sparse neural activities in primary visual cortex. The input is reconstructed by weighted sum of receptive fields of model neurons.

http://www.jzlab.org/Knag15_SAILnet_Chip.pdf



无监督特征学习 主成分分析

主成分分析 (Principal Component Analysis, PCA)

▶ 数据原始特征的问题:

- ▶ 高维→维度灾难、过拟合
- ▶ 冗余性→学习效果差
- ▶ 解决方法: 降维

▶ 一种最常用的数据降维方法

▶ 线性投影

$$\mathbf{z} = \mathbf{W}^T \mathbf{x}$$

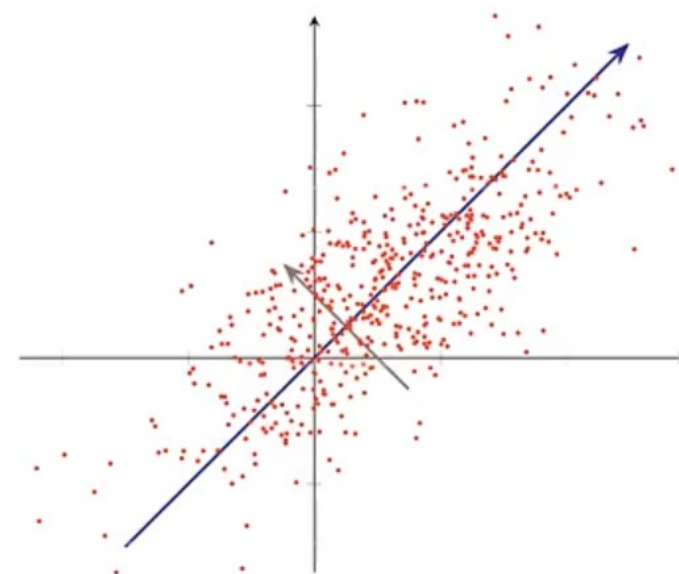
▶ 并满足

$$\mathbf{W}^T \mathbf{W} = \mathbf{I}$$

(保证尽可能小的冗余性)

▶ 优化准则

- ▶ 最大投影方差: 使得转换后的空间中数据方差最大, 尽可能多地保留原始数据信息
- ▶ 最小重构误差



主成分分析 (Principal Component Analysis, PCA)

▶ 目的：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

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

▶ 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

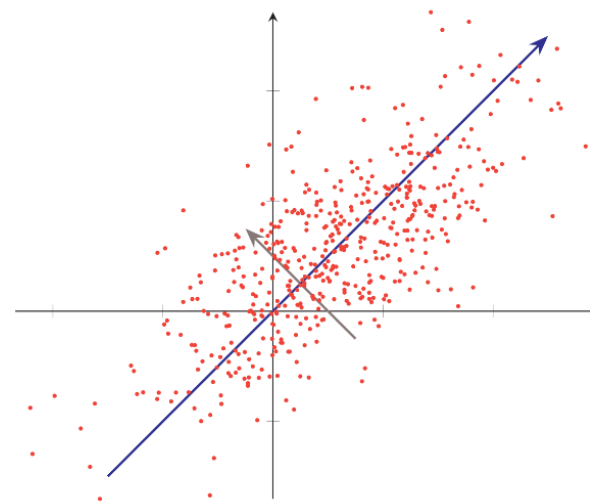
$$\max_{\mathbf{w}} \mathbf{w}^\top \Sigma \mathbf{w} + \lambda(1 - \mathbf{w}^\top \mathbf{w})$$

约束：内积=1
(拉格朗日乘子)

▶ 对目标函数求导并令导数等于 0，可得

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

(λ : 协方差矩阵的特征值)



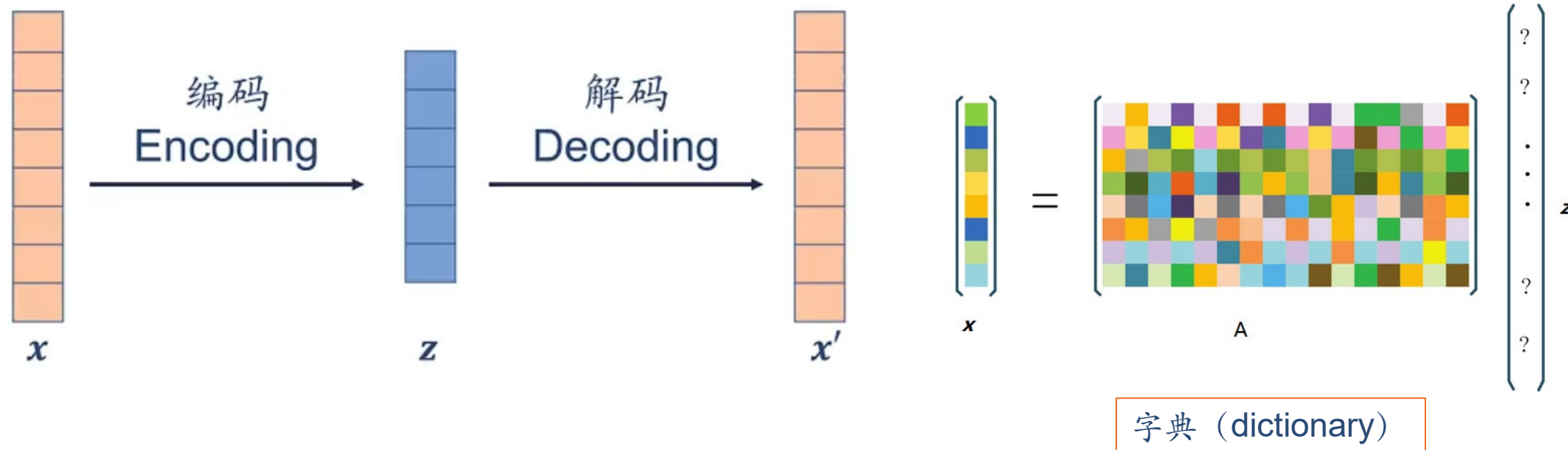


无监督特征学习 编码与稀疏编码

(线性) 编码

- ▶ 给定一组基向量 $A = [a_1, \dots, a_M]$ ，将输入样本 x 表示为这些基向量的线性组合

$$x = \sum_{m=1}^M z_m a_m$$
$$= Az,$$



稀疏编码 (Sparse Coding)

► 过完备

数学小知识 | 完备性

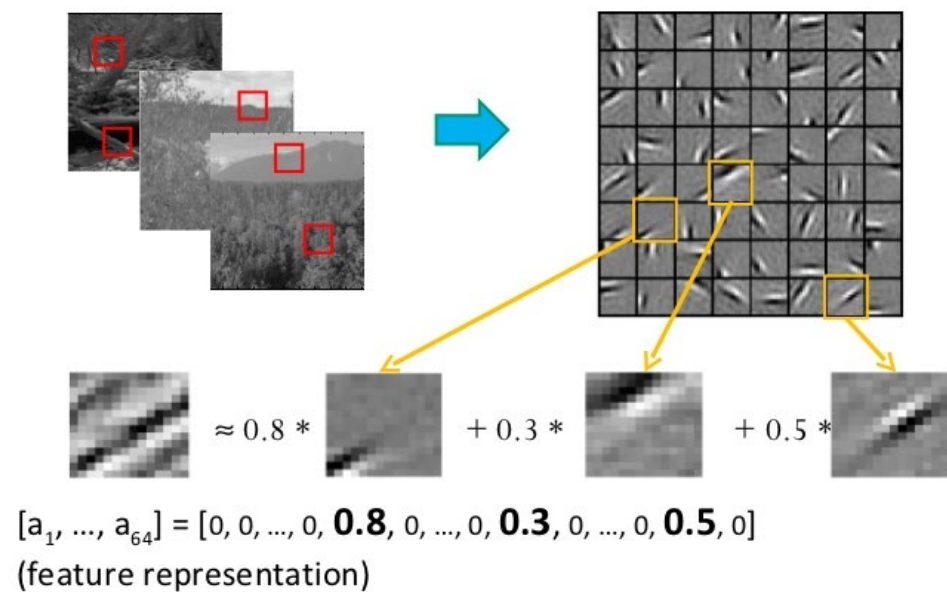
如果 M 个基向量刚好可以支撑 M 维的欧氏空间, 则这 M 个基向量是完备的. 如果 M 个基向量可以支撑 D 维的欧氏空间, 并且 $M > D$, 则这 M 个基向量是过完备的 (overcomplete)、冗余的.

“过完备”基向量是指基向量个数远远大于其支撑空间维度. 因此这些基向量一般不具备独立、正交等性质.

► 稀疏编码

► 找到一组“过完备”的基向量 (即 $M > D$) 来进行编码。

Sparse coding illustration



Slide credit: Andrew Ng

Compact & easily interpretable

稀疏编码 (Sparse Coding)

▶ 给定一组 N 个输入向量 $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$, 其稀疏编码的目标函数定义为

$$\mathcal{L}(\mathbf{A}, \mathbf{Z}) = \sum_{n=1}^N \left(\left\| \mathbf{x}^{(n)} - \mathbf{A}\mathbf{z}^{(n)} \right\|^2 + \eta \rho(\mathbf{z}^{(n)}) \right)$$

▶ $\rho(\cdot)$ 是一个稀疏性衡量函数, η 是一个超参数, 用来控制稀疏性的强度。

(不是连续可导)

$$\rho(\mathbf{z}) = \sum_{i=1}^p \mathbf{I}(|z_i| > 0)$$

$$\rho(\mathbf{z}) = \sum_{i=1}^p -\exp(-z_i^2)$$

$$\rho(\mathbf{z}) = \sum_{i=1}^p |z_i|$$

$$\rho(\mathbf{z}) = \sum_{i=1}^p \log(1 + z_i^2)$$

训练过程

► 稀疏编码的训练过程一般用交替优化的方法进行。

1) 固定基向量 \mathbf{A} , 对每个输入 $\mathbf{x}^{(n)}$, 计算其对应的最优编码

$$\min_{\mathbf{z}^{(n)}} \left\| \mathbf{x}^{(n)} - \mathbf{A}\mathbf{z}^{(n)} \right\|^2 + \eta \rho(\mathbf{z}^{(n)}), \quad \forall n \in [1, N].$$

2) 固定上一步得到的编码 $\{\mathbf{z}^{(n)}\}_{n=1}^N$, 计算其最优的基向量

$$\min_{\mathbf{A}} \sum_{n=1}^N \left(\left\| \mathbf{x}^{(n)} - \mathbf{A}\mathbf{z}^{(n)} \right\|^2 \right) + \lambda \frac{1}{2} \|\mathbf{A}\|^2,$$

稀疏编码的优点

▶降低后续计算量

▶稀疏性带来的最大好处就是可以极大地降低计算量。

▶可解释性强

▶因为稀疏编码只有少数的非零元素，相当于将一个输入样本表示为少数几个相关的特征。这样我们可以更好地描述其特征，并易于理解。

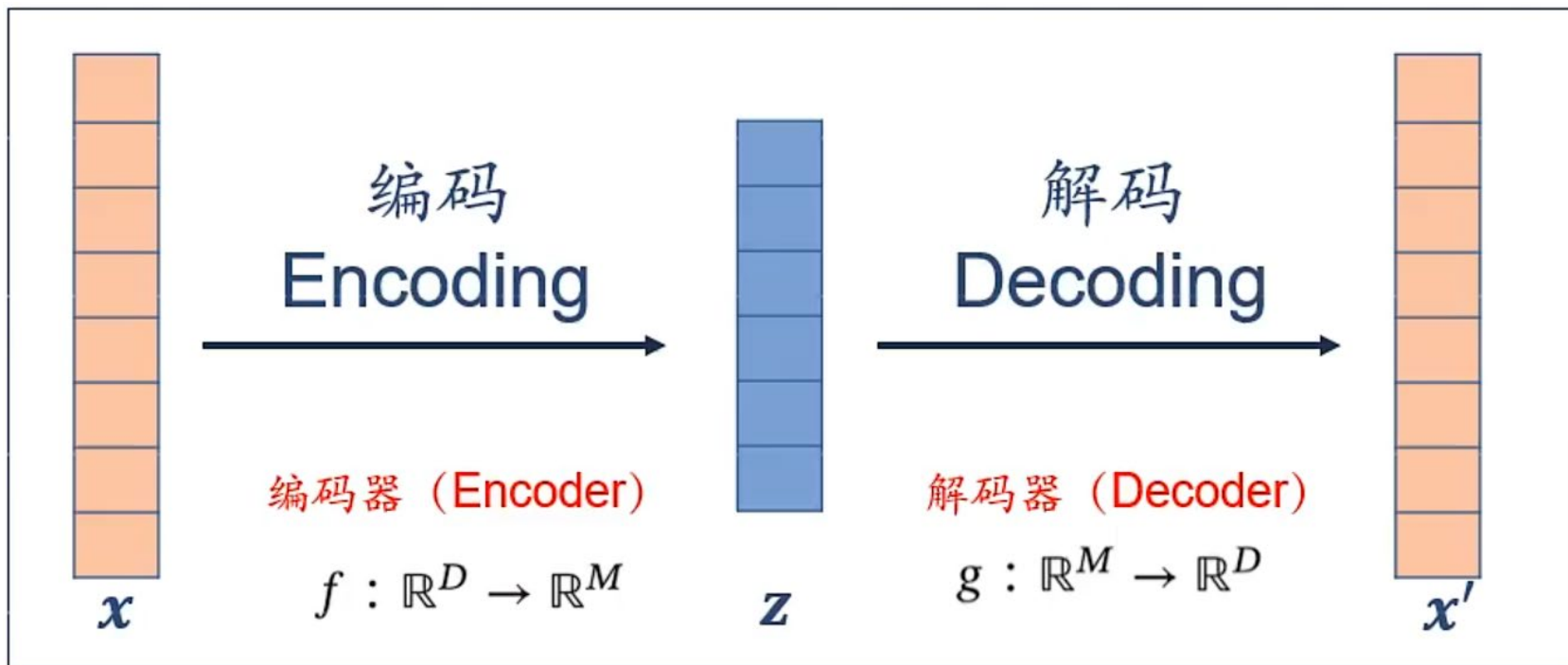
▶实现隐式的特征选择

▶稀疏性带来的另外一个好处是可以实现特征的自动选择，只选择和输入样本相关的最少特征，从而可以更好地表示输入样本，降低噪声并减轻过拟合。



无监督特征学习 自编码器

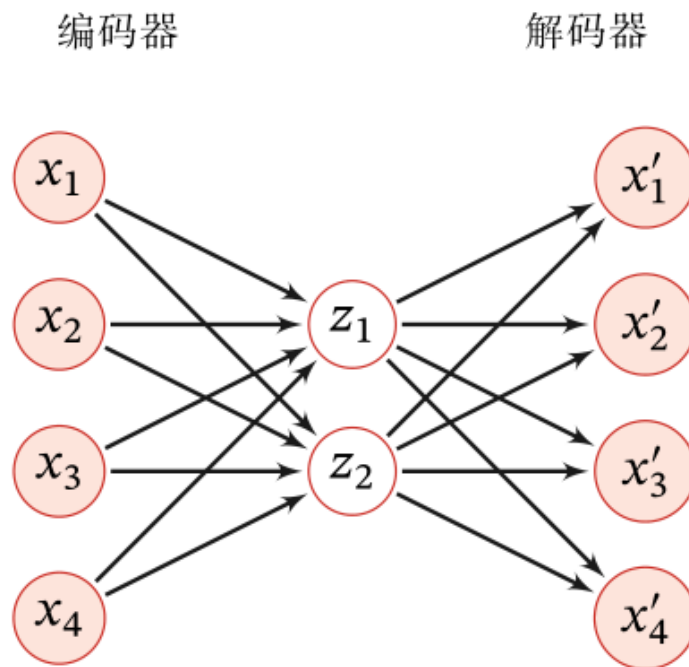
自编码器 (Auto-Encoder)



目标函数：重构错误

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

自编码器 (Auto-Encoder)



目标函数：重构错误

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

稀疏自编码器

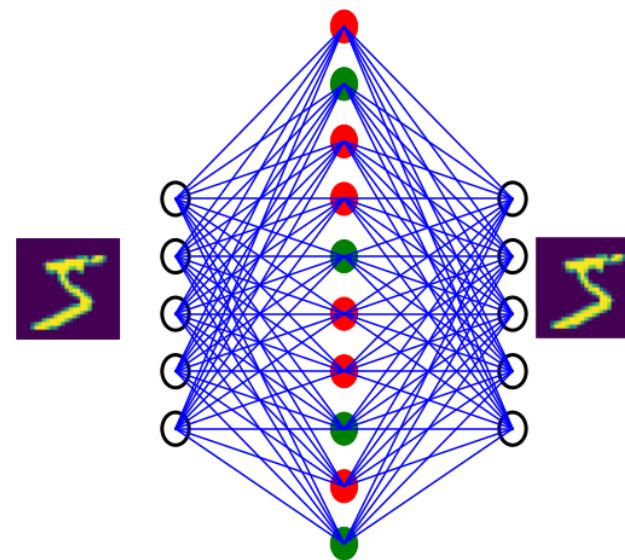
▶ 通过给自编码器中隐藏层单元 z 加上稀疏性限制，自编码器可以学习到数据中一些有用的结构。

▶ 目标函数

$$\mathcal{L} = \sum_{n=1}^N \|\mathbf{x}^{(n)} - \mathbf{x}'^{(n)}\|^2 + \eta \rho(\mathbf{Z}) + \lambda \|\mathbf{W}\|^2$$

▶ \mathbf{W} 表示自编码器中的参数

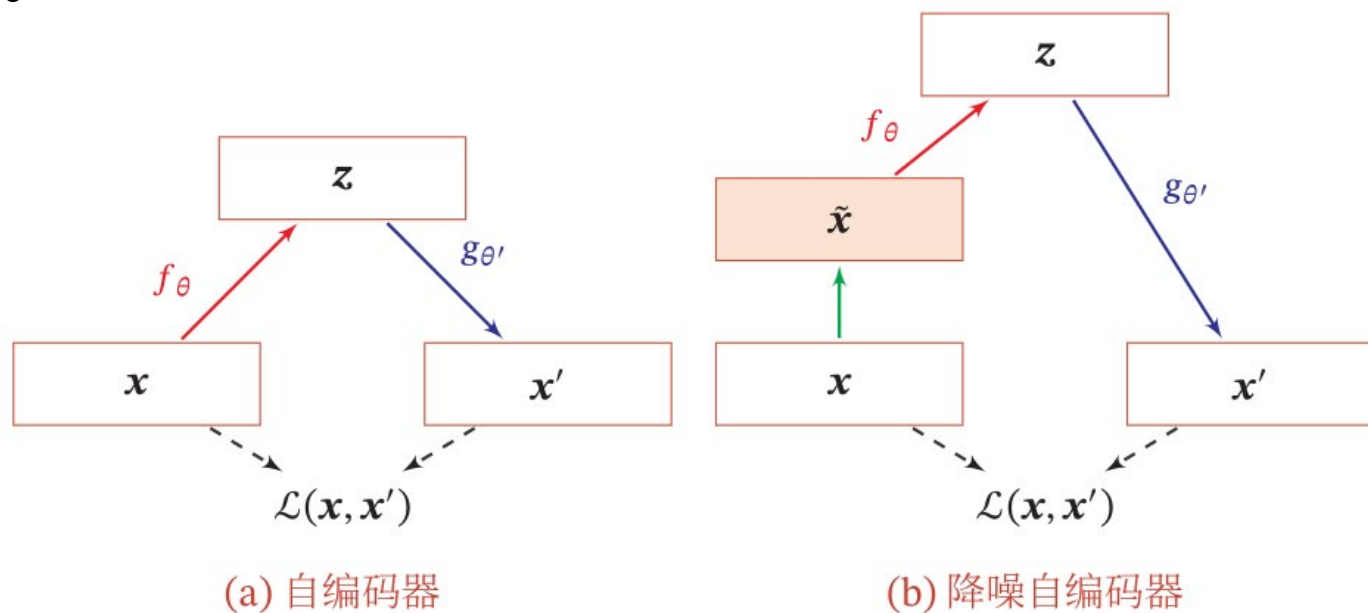
▶ 和稀疏编码一样，稀疏自编码器的优点是有很高的可解释性，并同时进行了隐式的特征选择。



降噪自编码器

▶ 通过引入噪声来增加编码鲁棒性的自编码器

- ▶ 对于一个向量 \mathbf{x} ，我们首先根据一个比例 μ 随机将 \mathbf{x} 的一些维度的值设置为0，得到一个被损坏的向量 $\tilde{\mathbf{x}}$ 。
- ▶ 然后将被损坏的向量 $\tilde{\mathbf{x}}$ 输入给自编码器得到编码 \mathbf{z} ，并重构出原始的无损输入 \mathbf{x} 。



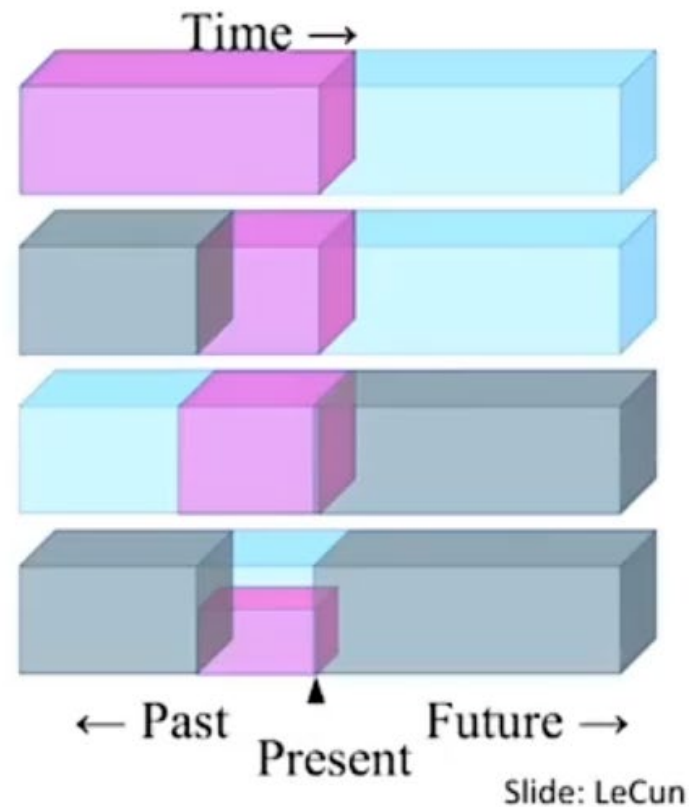


自监督学习

自监督学习 (Self-supervised Learning)

▶ 不再只是以输入重构作为目标，在无标签样本自身寻找更多样的“目标”

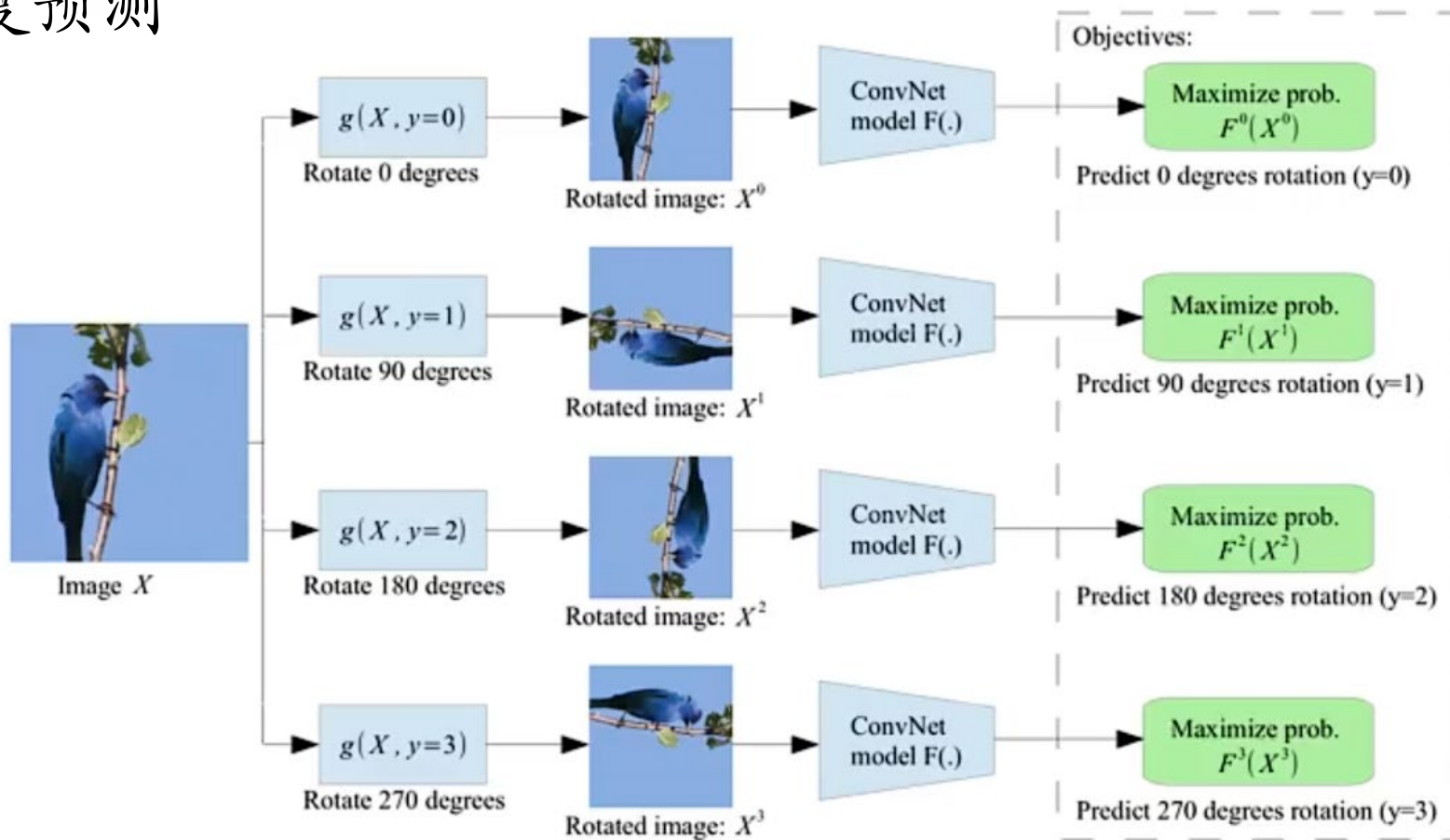
- ▶ Predict any part of the input from any other part.
- ▶ Predict the **future** from the **past**.
- ▶ Predict the **future** from the **recent past**.
- ▶ Predict the **past** from the **present**.
- ▶ Predict the **top** from the **bottom**.
- ▶ Predict the occluded from the visible
- ▶ Pretend there is a part of the input you don't know and predict that.



自监督学习 (Self-supervised Learning)

► 图像任务中的自监督学习

► 旋转角度预测



自监督学习 (Self-supervised Learning)

► 文本任务中的自监督学习

► 掩码语言模型 (Masked Language Model)

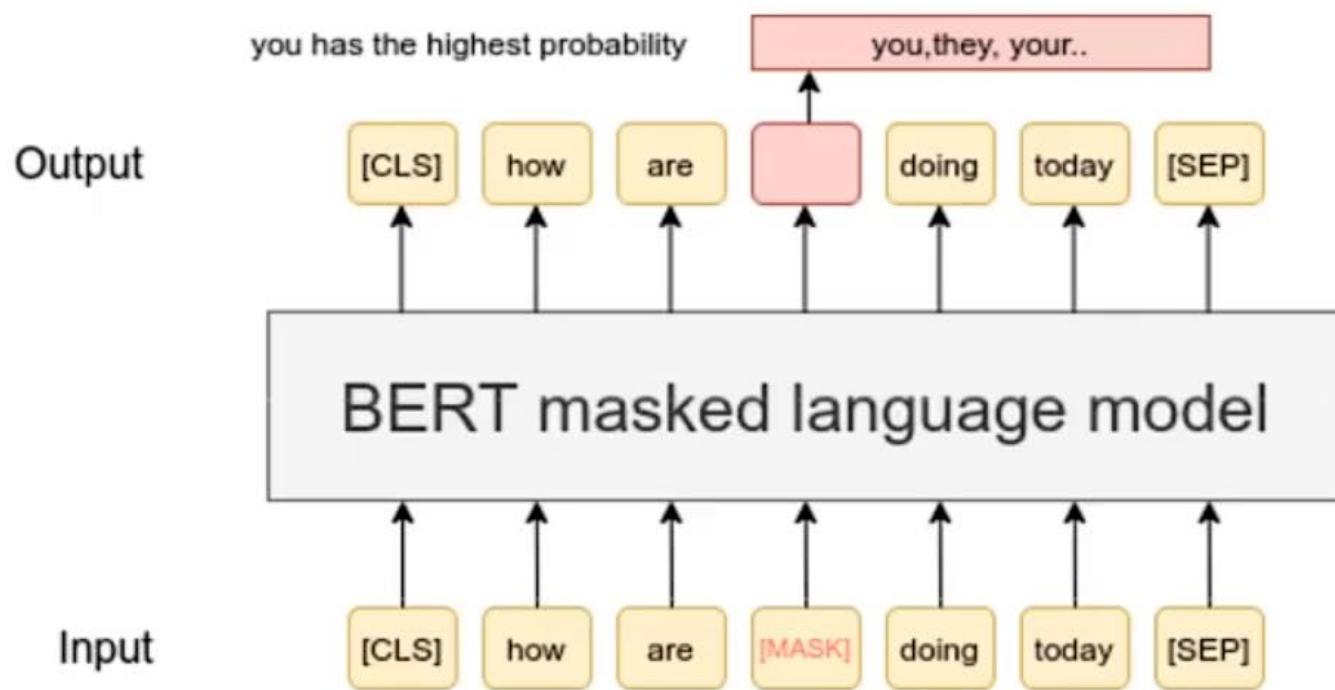


Fig. 1. BERT-Original sentence 'how are you doing today'



自监督学习 (Self-supervised Learning)



Yann LeCun

SSL is the future!

How Much Information is the Machine Given during Learning?

Y. LeCun

- ▶ **“Pure” Reinforcement Learning (cherry)**
 - ▶ The machine predicts a scalar reward given once in a while.
 - ▶ **A few bits for some samples**
- ▶ **Supervised Learning (icing)**
 - ▶ The machine predicts a category or a few numbers for each input
 - ▶ Predicting human-supplied data
 - ▶ **10→10,000 bits per sample**
- ▶ **Self-Supervised Learning (cake génoise)**
 - ▶ The machine predicts any part of its input for any observed part.
 - ▶ Predicts future frames in videos
 - ▶ **Millions of bits per sample**



© 2019 IEEE International Solid-State Circuits Conference

1.1: Deep Learning Hardware: Past, Present, & Future

59



概率密度估计

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

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

▶ 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.

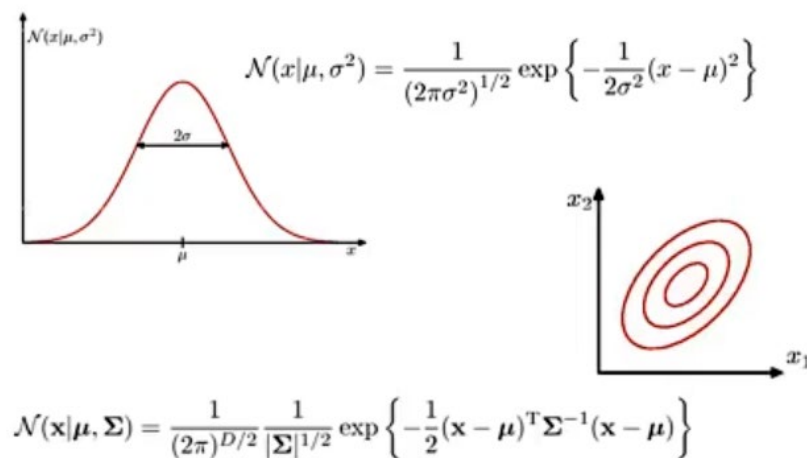
数据集 $\mathcal{D} = \{\mathbf{x}^{(n)}\}_{n=1}^N$ 的对数似然函数为

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

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

$$\boldsymbol{\mu}^{ML} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)},$$

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



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

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

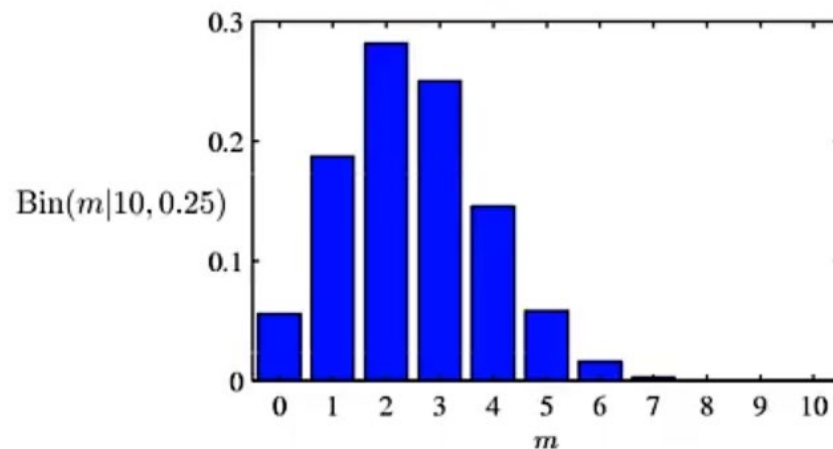
多项分布的参数估计为约束优化问题. 引入拉格朗日乘子 λ , 将原问题转换为无约束优化问题.

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

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

- ▶ 样本可能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

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

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

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

- ▶ 结合上述两个公式，得到

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

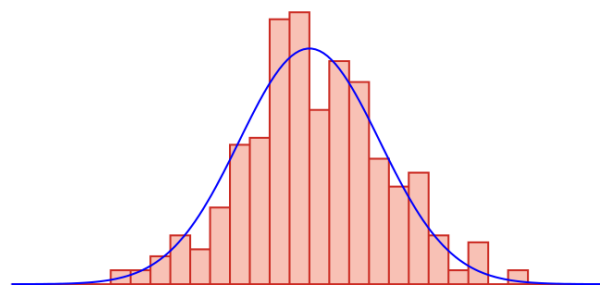
直方图方法 (Histogram Method)

- ▶ 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 (bin), 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.



(a) 10 个区间 (bin)

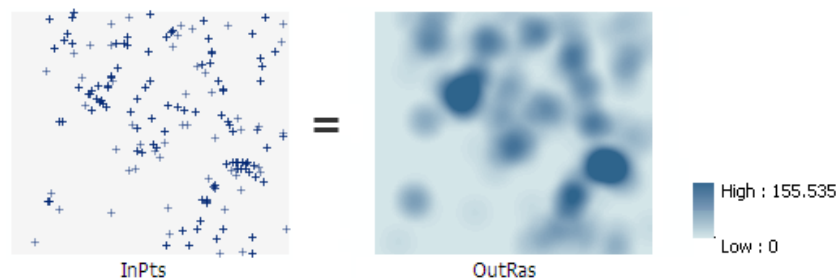


(b) 30 个区间 (bin)

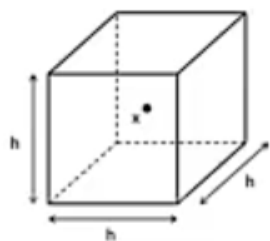
核密度估计 (Kernel Density Estimation)

►核密度估计是一种直方图方法的改进。

►也叫Parzen窗方法



►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



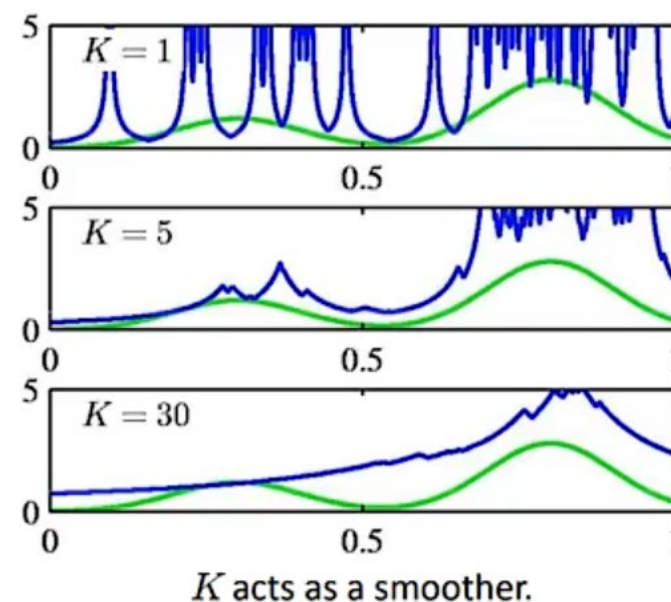
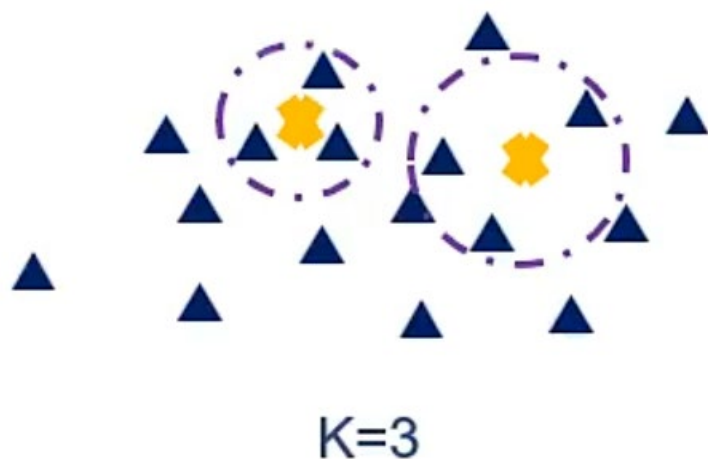
$$\phi\left(\frac{\mathbf{z} - \mathbf{x}}{h}\right) = \begin{cases} 1 & \text{if } |z_i - x_i| < \frac{H}{2}, 1 \leq i \leq D \\ 0 & \text{else} \end{cases}$$

►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

- ▶ 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。



非参数密度估计

- ▶ 非参数密度估计(直方图方法除外)需要保留整个训练集。
- ▶ 而参数密度估计不需要保留整个训练集，因此在存储和计算上更加高效。

思考题

- ▶ 估计以下三种密度函数的区别： $p(x)$, $p(y|x)$ and $p(x|y)$?
 - ▶ 生成模型 Generative model
 - ▶ 判别式模型 Discriminative model
- ▶ 最大似然函数的本质是什么？
 - ▶ 真实数据分布 $p_r(x)$
 - ▶ 模型数据分布 $p_\theta(x)$



半监督学习

半监督学习 (Semi-supervised Learning)

▶ 监督学习

- ▶ 提供任务相关的标签
- ▶ 打标签费时费力，数量有限

▶ 无监督学习

- ▶ 不用打标签，数量充足
- ▶ 仅限特定任务

同时利用少量有标签数据和大量无标签数据？



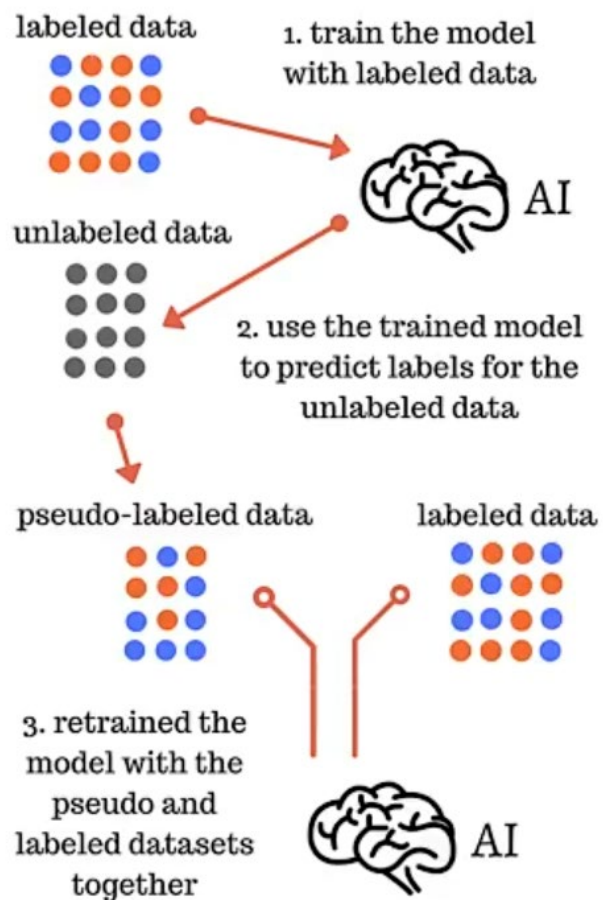
但是无标签数据缺失了很多重要的信息！

但是它或许仍然能够提供一些有用的信息。

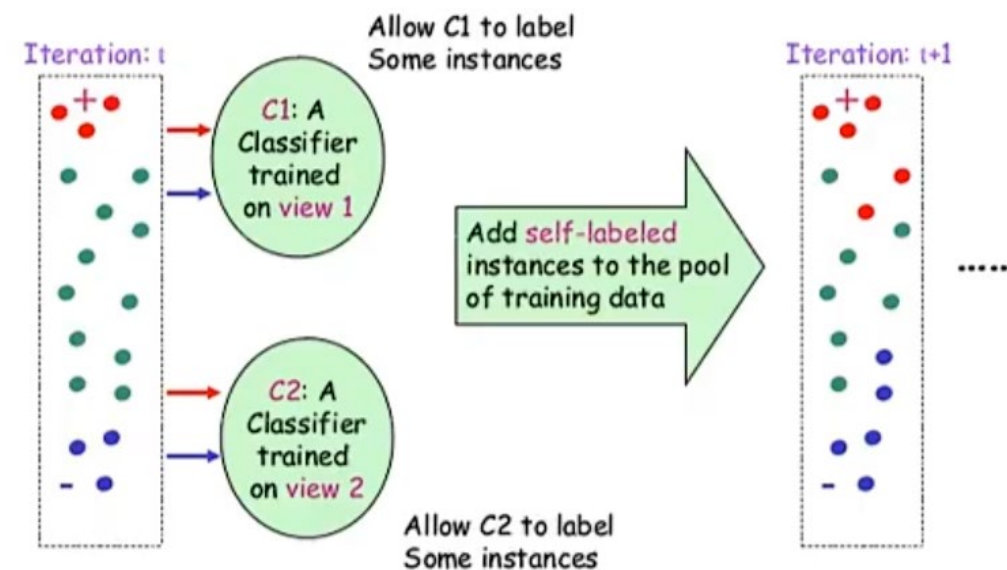


半监督学习 (Semi-supervised Learning)

► 自训练 (Self-training)



► 协同训练 (Co-training)



谢 谢