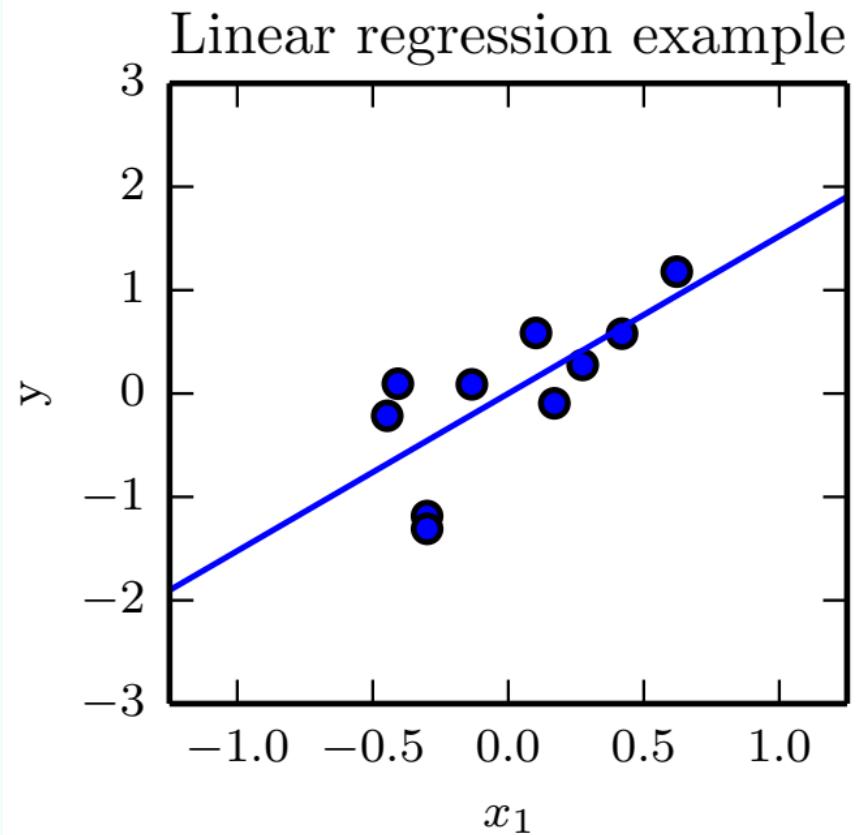


绪论

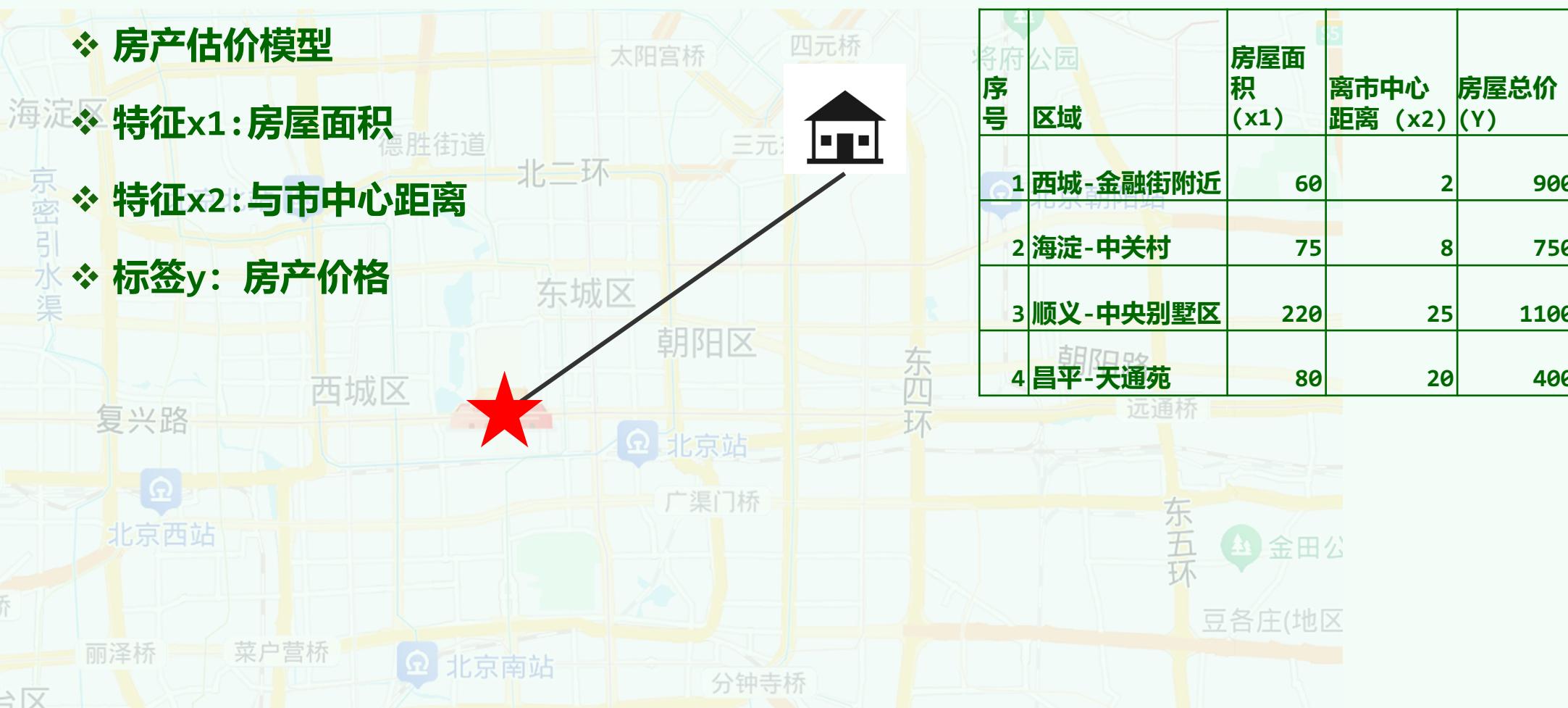
清华大学计算机系 陈键飞
《大模型计算》课程团队

函数拟合

- ◆ 给定训练数据集 $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$
- ◆ $x: D_X$ 维向量, $y: D_Y$ 维向量
- ◆ 找到函数 $f: \mathbb{R}^{D_X} \rightarrow \mathbb{R}^{D_Y}$
- ◆ 最小化拟合误差 $\mathcal{L} = \frac{1}{N} \sum_{i=1}^N \ell(y_i, f(x_i))$
- ◆ 函数包含某些未知参数 w
- ◆ 训练(training): 从数据集 \mathcal{D} 中学习得到模型 f (参数 w)
- ◆ 推理(inference): 根据模型 f 为新的 x 预测 $y = f(x)$



例：线性回归



例：线性回归

❖ 房产估价模型

❖ 特征 x_1 : 房屋面积

❖ 特征 x_2 : 与市中心距离

❖ 标签 y : 房产价格

❖ 模型 $y = w_1 \cdot x_1 + w_2 \cdot x_2 + c$

每平米单价 每公里惩罚 基础房价

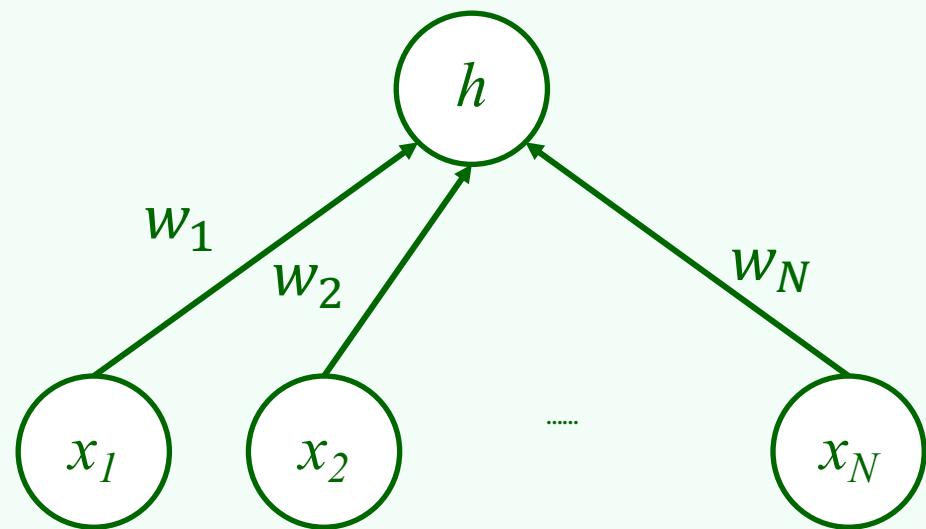
❖ $w_1 = 3.5$ $w_2 = -15$ $c = 700$

序号	区域	房屋面积 (x_1)	离市中心距离 (x_2)	房屋总价 (Y)	预测
1	西城-金融街附近	60	2	900	880
2	海淀-中关村	70	8	700	825
3	顺义-中央别墅区	220	25	1100	1095
4	昌平-天通苑	80	20	400	680

深度神经网络

清华大学计算机系 陈键飞
《大模型计算》课程团队

神经元



$$h = g\left(\sum_j w_j x_j + c\right)$$

向量形式

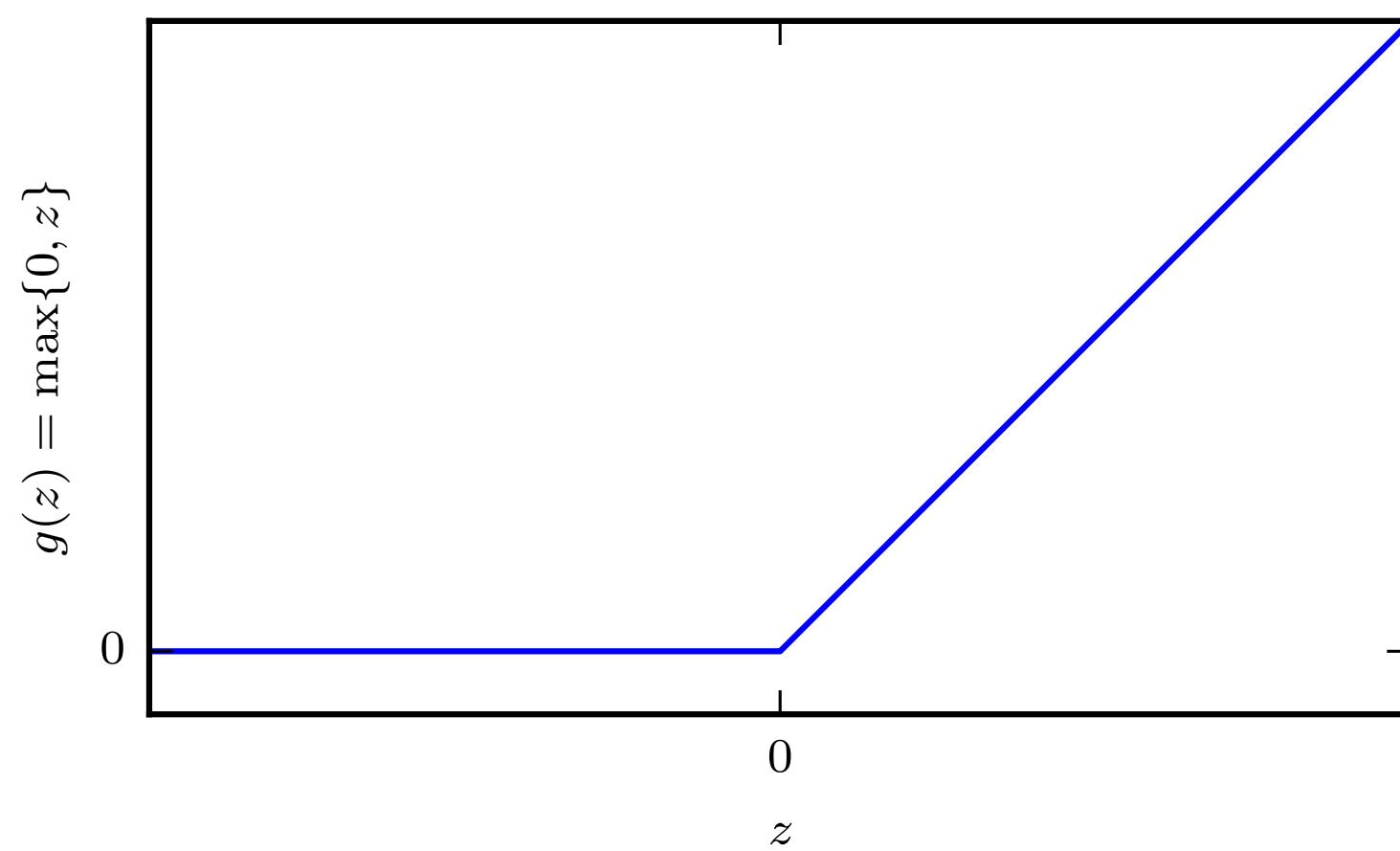
$$h = g(\mathbf{w}^\top \mathbf{x} + c)$$

权重 激活/输入 偏置
weight activation bias

单个神经元可以完成线性回归

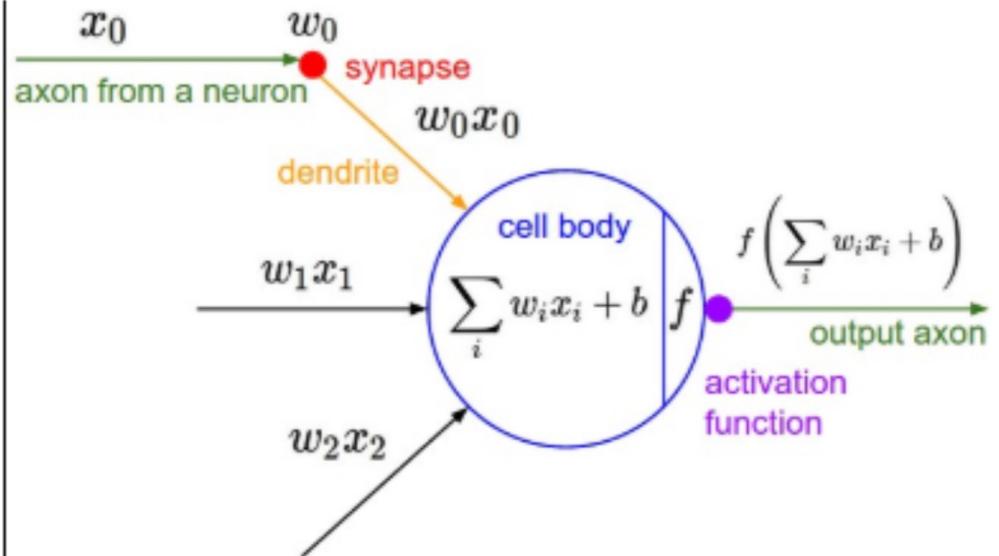
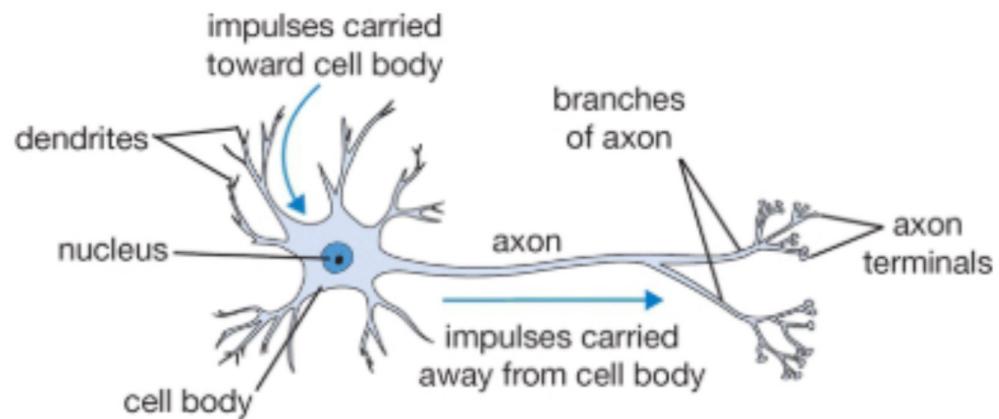
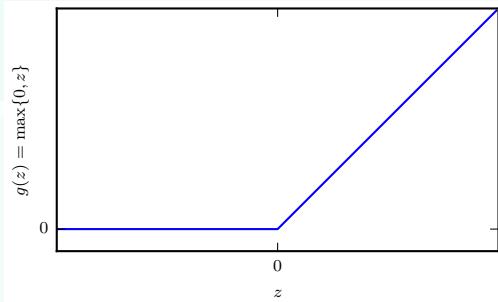
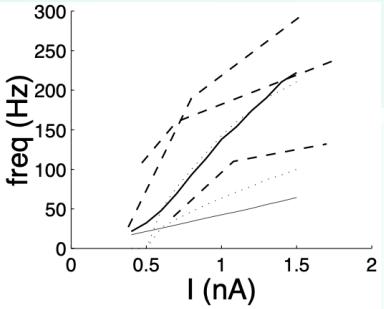
激活函数

整流线性单元 Rectifier Linear Unit (ReLU)
$$g(z) = \max\{0, z\}$$



生物神经元

Troyer T W, Miller K D. Integrate-and-fire neurons matched to physiological fI curves yield high input sensitivity and wide dynamic range[M]. Computational Neuroscience: Trends in Research, 1997. Boston, MA: Springer US, 1997: 197-201.



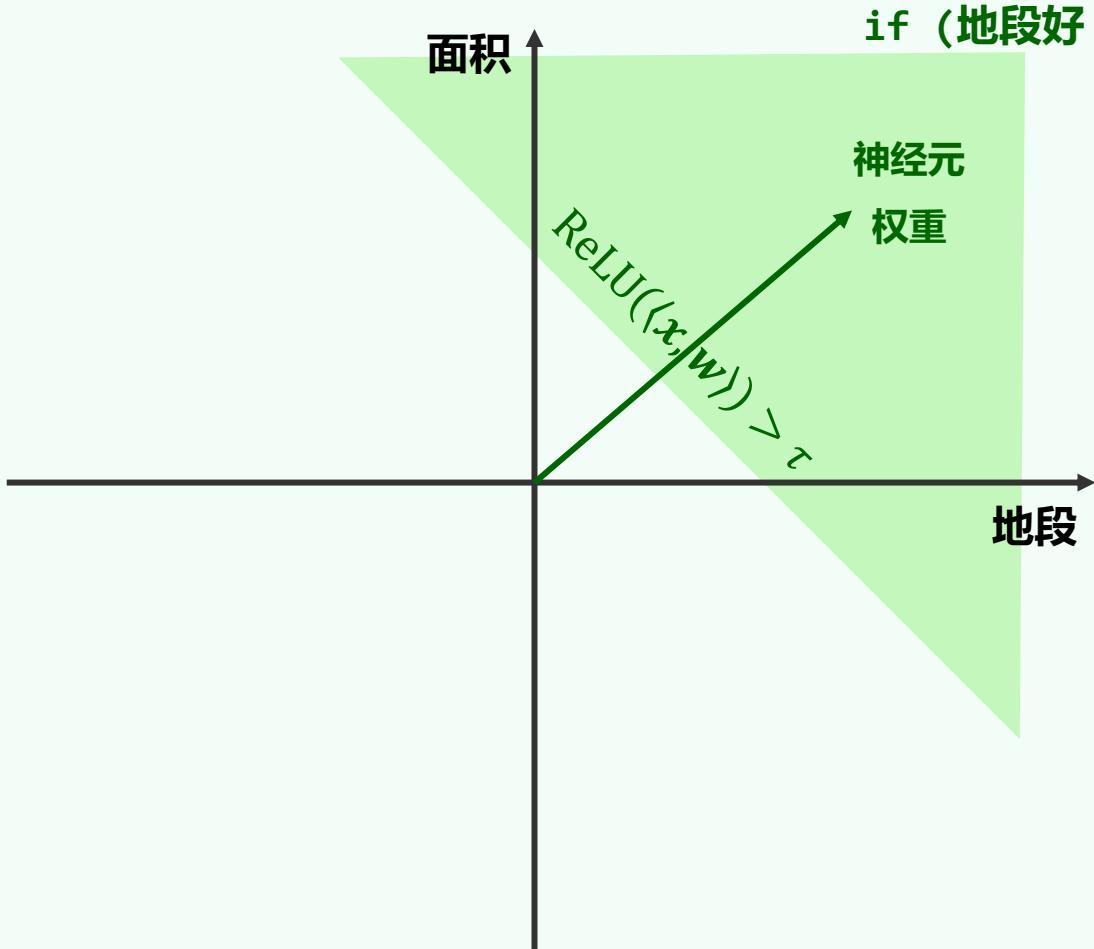
A cartoon drawing of a biological neuron (left) and its mathematical model (right).

生物神经元

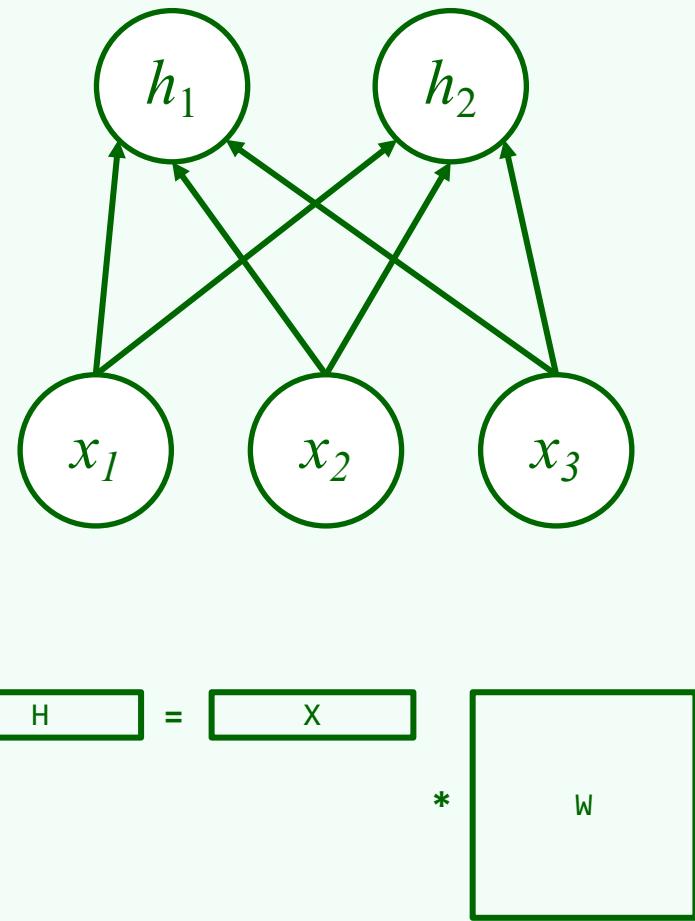
人工神经元

神经元 = 组合特征捕捉器

捕捉 “地段好 **且** 面积大” 概念



全连接层（线性层）



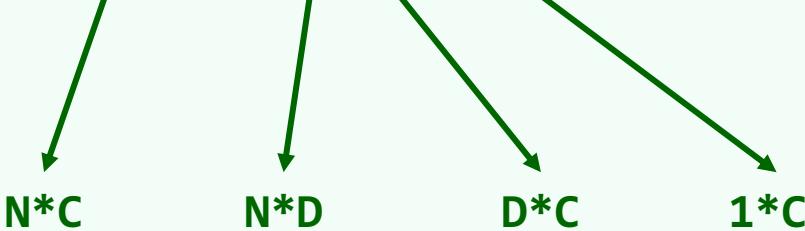
$$h_i = g\left(\sum_j w_{ij}x_j + c_i\right)$$

矩阵形式（输入为向量）

$$h = g(Wx + c)$$

矩阵形式（输入为矩阵）

$$H = g(XW^T + c)$$

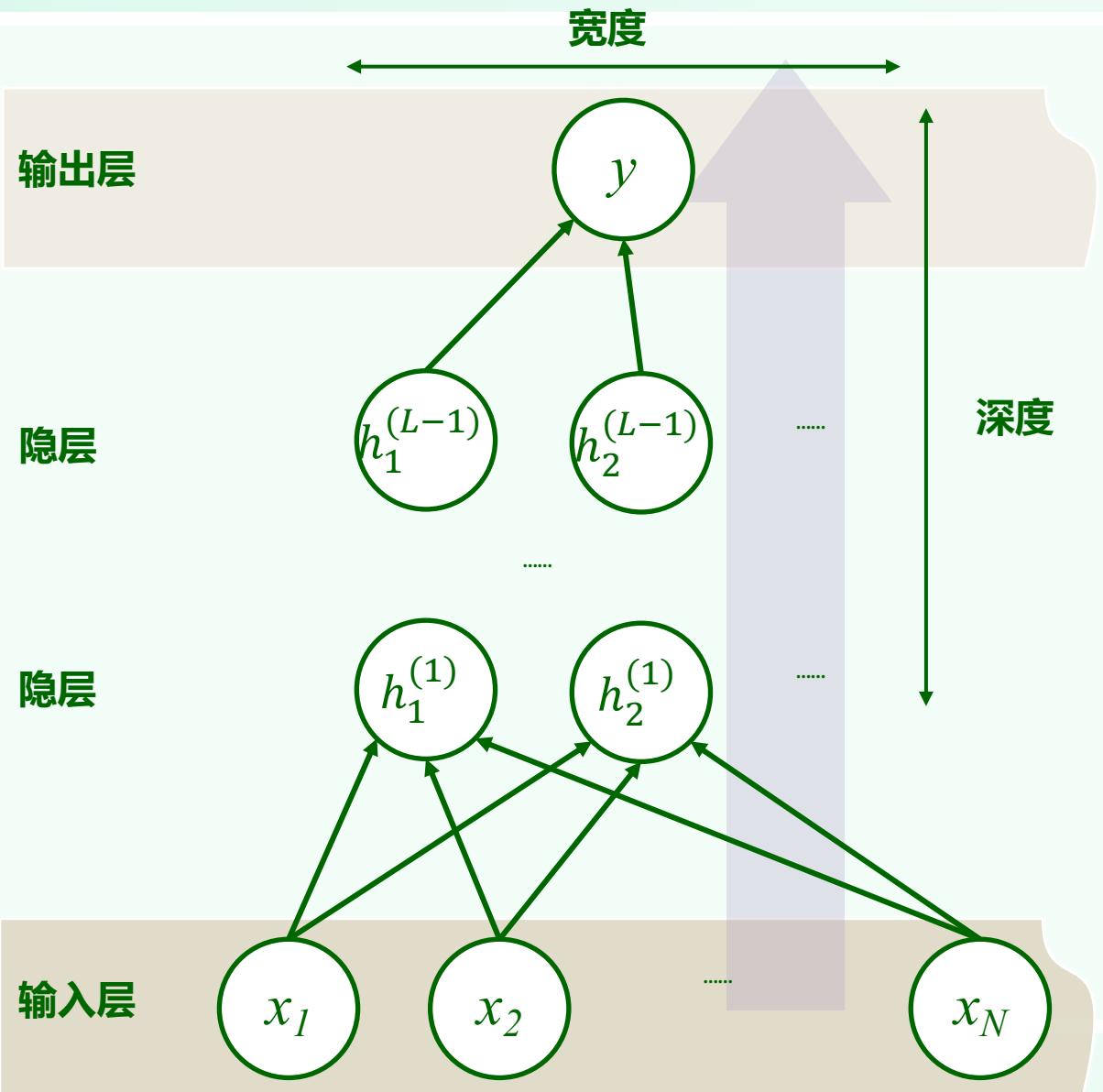


N: 数据

C: 输出维度

D: 输入维度

神经网络



```
h1 = linear(x, w1)  
h2 = linear(h1, w2)  
...  
y = linear(h_{L-1}, w_L)
```

$$f(x; W)$$

数据

模型

全连接神经网络 (Fully-connected network)
多层感知机 (Multi-layer perceptron, MLP)
前馈网络 (Feed-forward network, FFN)

例：线性回归

- ❖ 房产估价模型
- ❖ 特征 x_1 : 房屋面积
- ❖ 特征 x_2 : 与市中心距离
- ❖ 标签 y : 房产价格

$$\text{❖ 模型 } y = w_1 \cdot x_1 + w_2 \cdot x_2 + c$$

每平米单价 每公里惩罚 基础房价

$$\text{❖ } w_1 = 3.5 \quad w_2 = -15 \quad c = 700$$

序号	区域	房屋面积 (x_1)	离市中心距离 (x_2)	房屋总价 (Y)	预测
1	西城-金融街附近	60	2	900	880
2	海淀-中关村	70	8	700	825
3	顺义-中央别墅区	220	25	1100	1095
4	昌平-天通苑	80	20	400	680

例：线性回归

❖ 房产估价模型

❖ 特征 x_1 : 房屋面积

❖ 特征 x_2 : 与市中心距离

❖ 标签 y : 房产价格

❖ 模型 $y = w_1 \cdot x_1 + w_2 \cdot x_2 + c$

每平米单价 每公里惩罚 基础房价

❖ $w_1 = 3.5$ $w_2 = -15$ $c = 700$

序号	区域	房屋面积 (x_1)	离市中心距离 (x_2)	房屋总价 (Y)	预测
1	西城-金融街附近	60	2	900	880
2	海淀-中关村	70	8	700	825
3	顺义-中央别墅区	220	25	1100	1095
4	昌平-天通苑	80	20	400	680
5	西城-金融街附近	200	2	4000	1370

例：线性回归

❖ 房产估价模型

❖ 特征 x_1 : 房屋面积

❖ 特征 x_2 : 与市中心距离

❖ 标签 y : 房产价格

❖ 模型 $y = w_1 \cdot x_1 + w_2 \cdot x_2 + c$

每平米单价 每公里惩罚 基础房价

❖ $w_1 = 3.5$ $w_2 = -15$ $c = 700$

序号	区域	房屋面积 (x_1)	离市中心距离 (x_2)	单价	房屋总价 (Y)	预测
1	西城-金融街附近	60	2	15	900	880
2	海淀-中关村	70	8	10	700	825
3	顺义-中央别墅区	220	25	5	1100	1095
4	昌平-天通苑	80	20	5	400	680

单价受地段影响很大!

例：线性回归

- ❖ 房产估价模型
- ❖ 特征 x_1 : 房屋面积
- ❖ 特征 x_2 : 与市中心距离

❖ 标签 y : 房产价格

❖ 基础房价

$$h_1 = 5 * x_1$$

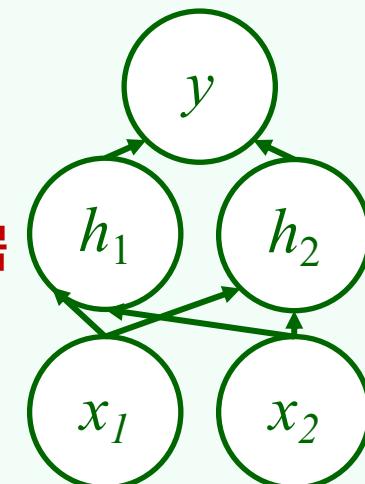
❖ 核心城区溢价

$$h_2 = \text{ReLU}(12 * x_1 - 62 * x_2)$$
 有指向性的筛选出部分数据

❖

$$y = h_1 + h_2$$

序号	区域	房屋面积 (x_1)	离市中心距离 (x_2)	单价	房屋总价 (Y)	h_1	h_2	预测
1	西城-金融街附近	60	2	15	900	300	596	896
2	海淀-中关村	70	8	10	700	350	344	694
3	顺义-中央别墅区	220	25	5	1100	1100	0	1100
4	昌平-天通苑	80	20	5	400	400	0	680



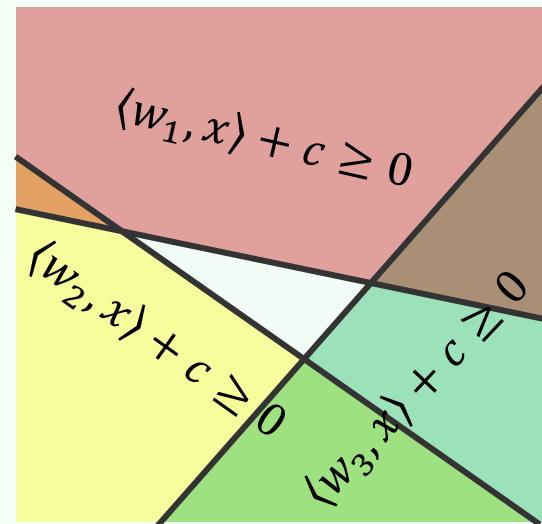
例：线性回归



万有近似定理 (Universal Approximation Theorem)

给定任意数据分布 $p(x)$ 及连续函数 $f(x)$, 存有一族单隐层神经网络, 当网络宽度趋于无穷时, 对 $(x, f(x))$ 的拟合误差趋于 0。

```
if input is x1 then output y1  
if input is x2 then output y2  
if input is x3 then output y3  
.....
```



Cybenko G. Approximation by superpositions of a sigmoidal function[J]. Mathematics of control, signals and systems, 1989, 2(4): 303-314.

PyTorch

清华大学计算机系 陈键飞
《大模型计算》课程团队

张量

❖ 数值组成的网格

❖ 形状Shape [D_1, D_2, ..., D_R]



R阶张量

R个轴 (axis)

❖ 标量

[]

❖ 3维向量

[3]

❖ 5*3矩阵

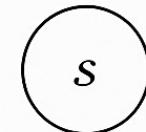
[5, 3]

❖ Transformer隐层

[128, 512, 768]

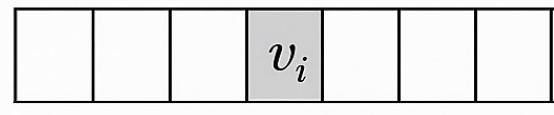
(a) 标量

Scalar



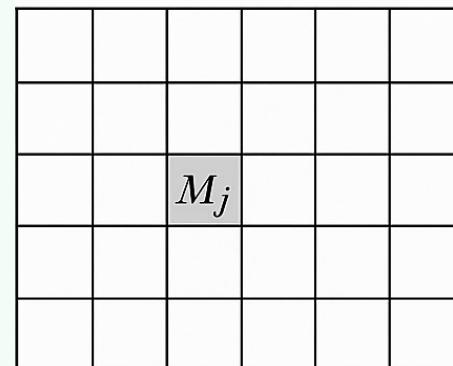
(b) 向量

Vector



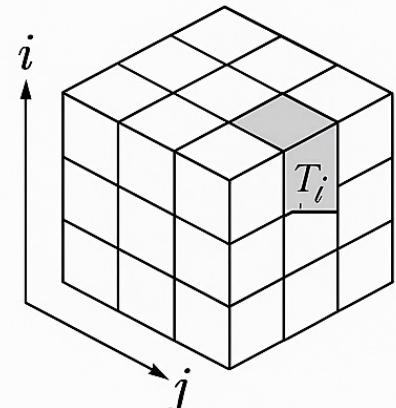
(c) 矩阵

Matrix



(d) 三阶张量

3rd-order Tensor

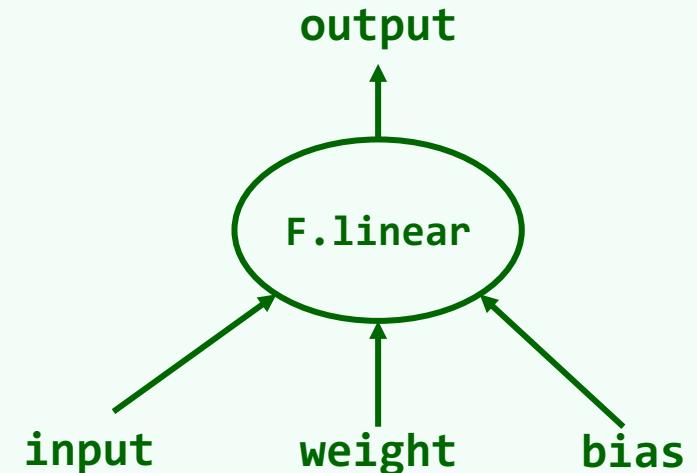


算子

❖ 多张量输入，多张量输出的函数

$$y_1, y_2, \dots, y_m = f(x_1, x_2, \dots, x_n)$$

- ❖ 逐元素加法 `torch.add(x, y)`
- ❖ 矩阵乘法 `torch.dot(x, y)`
- ❖ 转置 `x.t()`
- ❖ 奇异值分解 `U, s, V = torch.svd(X)`
- ❖ 全连接层 `h = F.linear(input, weight, bias)`



模块

❖ 神经网络模型的抽象

- 参数 Parameter
- 前向传播 forward

```
class Linear:  
    def __init__(self, in_features: int, out_features: int, bias: bool = True) :  
        self.weight = Parameter(torch.empty((out_features, in_features)))  
        if bias:  
            self.bias = Parameter(torch.empty(out_features,))  
        self.reset_parameters()                                # 初始化参数  
  
    def forward(self, input: Tensor):  
        return F.linear(input, self.weight, self.bias)          # input @ self.weight.t() + self.bias
```

统计机器学习

清华大学计算机系 陈键飞
《大模型计算》课程团队

离散概率分布

- ❖ **V面骰子** $Y \sim \text{Categorical}(p)$, 概率 $P(Y = v) = p_v$
- ❖ 观测到了N次骰子投掷的结果 $\mathcal{D} = y_1, \dots, y_N$, 其中 $y_i \in \{1, \dots, V\}$
- ❖ 问题：如何估计骰子自身的参数 p ？
- ❖ 极大似然估计：给定参数 p , 出现结果 \mathcal{D} 的概率为

$$P(\mathcal{D}; p) = \prod_{i=1}^N p_{y_i} = \prod_{v=1}^V p_v^{s_v}$$



离散概率分布

- ❖ **V面骰子** $Y \sim \text{Categorical}(p)$, 概率 $P(Y = v) = p_v$
- ❖ 观测到了N次骰子投掷的结果 $\mathcal{D} = y_1, \dots, y_N$, 其中 $y_i \in \{1, \dots, V\}$
- ❖ 问题：如何估计骰子自身的参数 p ？
- ❖ 极大似然估计：给定参数 p , 出现结果 \mathcal{D} 的概率为

$$P(\mathcal{D}; \mathbf{p}) = \prod_{i=1}^N p_{y_i} = \prod_{v=1}^V p_v^{s_v}$$

- ❖ 负对数似然度 (negative log-likelihood, NLL)

$$-\log P(\mathcal{D}; \mathbf{p}) = -\sum_v s_v \log p_v$$

❖ 独热 (one-hot) 编码

$y = 1$	$y = (1, 0, 0, 0)$
$y = 2$	$y = (0, 1, 0, 0)$
$y = 3$	$y = (0, 0, 1, 0)$
$y = 4$	$y = (0, 0, 0, 1)$

整数表示

坐标向量表示

❖ 充分统计量

各面计数

$$\mathbf{s} = \prod_{i=1}^N y_i$$

极大似然估计

❖ 极大似然估计：从所有可能的参数 p 中找出现结果 \mathcal{D} 的概率最大者

$$\begin{aligned} \min_p -\log P(\mathcal{D}; p) &= -\sum_v s_v \log p_v \\ \text{s. t. } \sum_{v=1}^V p_v &= 1 \end{aligned}$$

❖ 拉格朗日乘子

$$\mathcal{L} = -\sum_v s_v \log p_v + \lambda \left(\sum_v p_v - 1 \right)$$

❖ 求解

$$\begin{aligned} \mathcal{L}'_{p_v} &= \frac{s_v}{p_v} + \lambda = 0 \quad \Rightarrow \quad p_v = -\frac{s_v}{\lambda} \\ \sum_{v=1}^V p_v &= 1 \quad \Rightarrow \quad -\sum_v \frac{s_v}{\lambda} = 1 \quad \Rightarrow \quad \lambda = -\sum_v s_v \\ &\Rightarrow p_v = \frac{s_v}{\sum_w s_w} \end{aligned}$$

分类问题

❖ 问题：骰子的概率取决于条件 x ，要估计的是条件概率 $P(y|x)$

❖ 例：天气预测



$$\begin{aligned}P(y = \text{晴} | x) &= 0.9 \\P(y = \text{雨} | x) &= 0.1\end{aligned}$$



$$P(y = \text{多云} | x) = 0.7$$



$$P(y = \text{雨} | x) = 0.9$$

❖ 挑战： x 和 y 之间关系可能很复杂

❖ 思路：利用神经网络建模

分类问题

- ❖ 模型: $P(y|x) \sim \text{Categorical}(f(x; w))$
- ❖ f : 从 x (D维) 到 p (V维)
- ❖ 模型参数由概率向量 p 变成网络权重 w
- ❖ 负对数似然度 $-\log P(\mathcal{D}; w) = -\sum_v s_v \log f(x; w)_v$
- ❖ 等等.....哪里不太对?

分类问题

- ❖ 模型: $P(y|x) \sim \text{Categorical}(f(x; w))$
- ❖ f : 从 x (D维) 到 p (V维)
- ❖ 模型参数由概率向量 p 变成网络权重 w

- ❖ 负对数似然度 $-\log P(\mathcal{D}; w) = -\sum_v s_v \log f(x; w)_v$
- ❖ 骰子  的概率 p 应该是归一化的, 但神经网络的输出无法保证这点.....
- ❖ 解决方法: 手工归一化

$$f(x; w) = \text{softmax}(\underbrace{g(x; w)}_{\substack{\text{概率} \\ \text{probability}}})$$

←

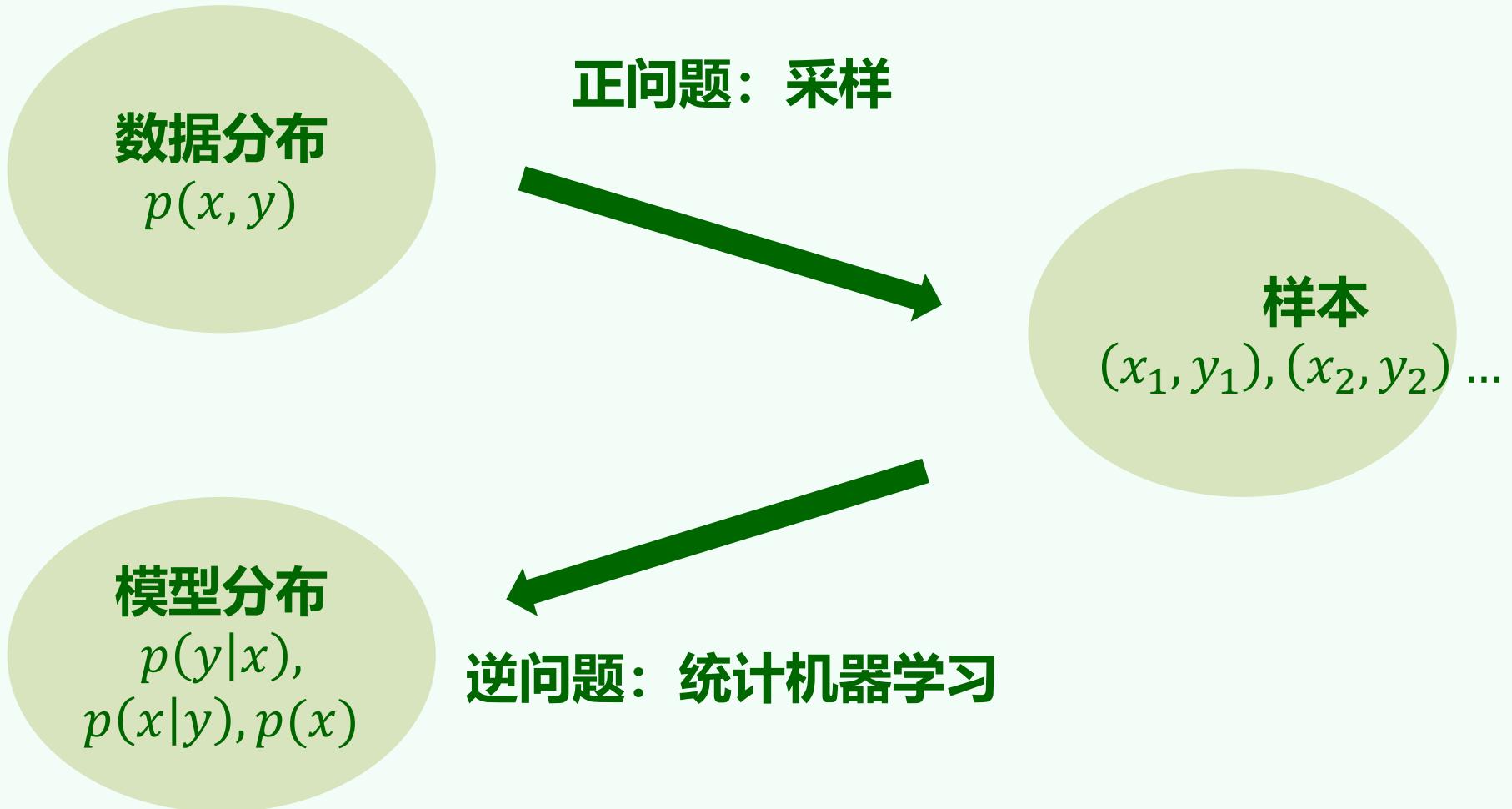
$$\text{softmax}(\mathbf{g})_v = \exp(g_v) / \sum_w \exp(g_w)$$

神经网络
logits

交叉熵损失

- ❖ 模型: $P(y|x) \sim \text{Categorical}(f(x; w))$
- ❖ f : 从 x (D维) 到 p (V维)
- ❖ 模型参数由概率向量 p 变成网络权重 w
- ❖ 负对数似然度 $-\log P(\mathcal{D}; w) = -\sum_v s_v \log f(x; w)_v = -\sum_v s_v g_v(x; w) + N \log \sum_w \exp(g_w(x; w))$
- ❖ 无解析最小值, 需要用数值方法 (AdamW) 迭代求解.....
- ❖ 模型有什么用?
- ❖ 采样: 给定 x , 产生 $y \sim p(y|x)$
- ❖ 概率计算: 给定 x 和 y , 计算 $p(y|x)$

小结



大语言模型

清华大学计算机系 陈键飞
《大模型计算》课程团队

语言模型

- ❖ 学习概率分布 P (文档)
- ❖ P 是 文档 $\rightarrow [0,1]$ 的映射，给每篇可能的文档一个概率值

我的家乡有一条清澈的小河，鱼儿在水中嬉戏，河边柳树轻轻摇曳。我和小伙伴们常在这里玩耍，快乐无比。我爱我的家乡！

$P=1e-100$

我滴家多有一条清彻见底的小盒，鱼仔在水中自游自在地嬉戏，河边柳树轻摇业。我和小伙半常在这里晚要，快乐比无。我哎我家多

$P=1e-120$

荇葵蕡美蔻蔻葩葵菌萎葉莘
殖莖萼蘋莖蘋莖蘋莖蘋莖蘋
造淺莖陪許莖蘋莖蘋莖蘋莖
噴莖蘋莖蘋莖蘋莖蘋莖蘋莖
葵雍菽葵蔬冀莖斯酷薑擅薑

$P=1e-500$

- ❖ 如何把文档转换成模型能理解的形式？（最好是向量）

分词 (tokenize)

❖ 如何把文档表示成向量?

It's a good day today.

↓ 文字转整数

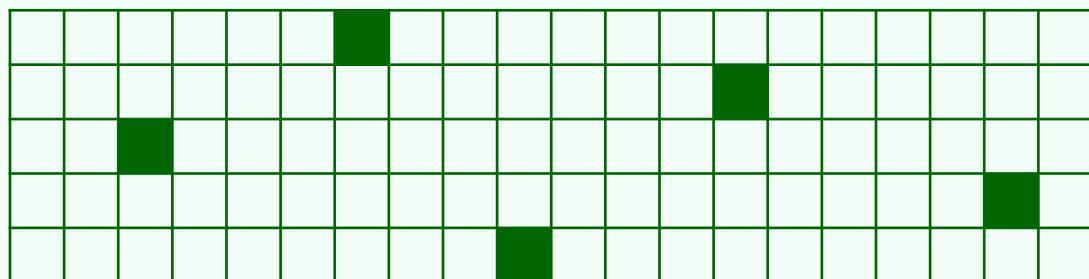
(15000, 1, 10000, 3000, 50000, 80000)

Token
词元

↓ 整数转独热

v列 (词汇表大小)

T行
(文档长度)



注：该过程并不精确，实际分词的算法比此处复杂。

❖ 词汇表

Id	词
1	a
2	abandon
3000	day
10000	good
15000	It's
50000	today
80000	.
100000	<EOS>

序列建模

- ❖ 如何建模长度可变，且最长可能可达128K，甚至1M的词元序列？
- ❖ 长度可变：设计特殊词元<EOS>

It's a good day today. <EOS> <EOS> <EOS> <EOS> ...

固定长度 (如8192)

- ❖ 数据变为定长的词元序列 (x_1, \dots, x_T)
- ❖ 直接利用统计的手段建模 $P(x_1, \dots, x_T)$ 需要存储 V^T 个参数.....
- ❖ $V = 10^6, T = 8192$, 则 $V^T \approx 10^{50000}$

自回归模型 (autoregressive model)

- ❖ 思路：每次只预测一个词
- ❖ $p(x_1, \dots, x_T) = p(x_1)p(x_2|x_1) \dots p(x_T|x_1, \dots, x_{T-1})$
- ❖ 把 T 维分布的建模问题转换为 T 个1维分布的建模问题
- ❖ $p(x_i|x_{<i}) = \text{softmax}(g(x_{<i}))$
- It's a good day _____
- $P(x_5 | \text{"It's a good day"}) = \dots$
- ❖ 采样

Id	词	概率
1	a	1e-7
2	abandon	1e-7
3000	day	0.01
10000	good	0.0001
15000	It's	1e-7
50000	today	1e-7
80000	.	1e-7
100000	<EOS>	1e-7

小知识：从离散分布中采样



Input: probability p

```
s = cumsum(p)          #  $s_i = \sum_{j \leq i} p_j$ 
pos = rand() * s[-1]    #  $[0, \sum_{j=1}^N p_j]$ 
i = upper_bound(s, pos) #  $\min_i \text{ s.t. } s_i > pos$ 
```