



机器学习与数据挖掘

线性分类器



课程大纲

- 二分类
- 多分类

什么是分类

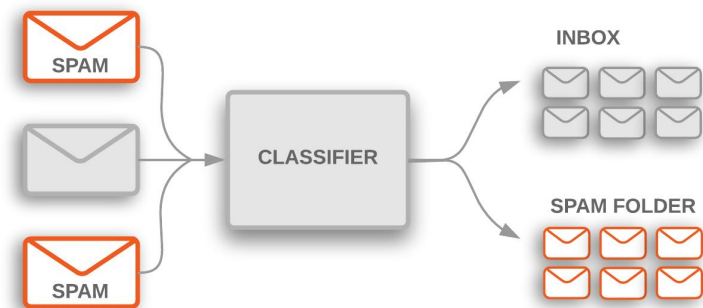
- 分类是监督学习的一个核心任务
- 目标是学习一个模型（称为分类器）
- 该模型能将输入数据（特征）映射到一个预定义的离散类别
- 与回归的区别： 回归预测的是连续值（如房价、温度），而分类预测的是类别标签（如“猫”或“狗”，“垃圾邮件”或“非垃圾邮件”）

示例

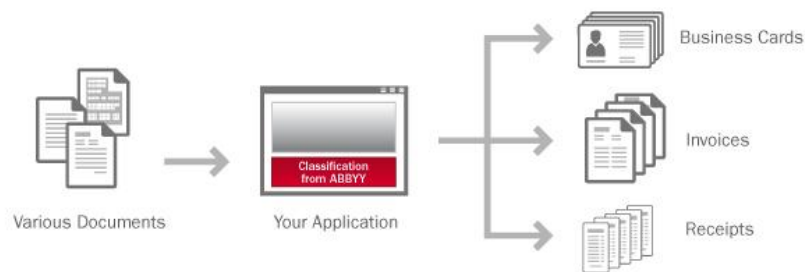
- 图像识别



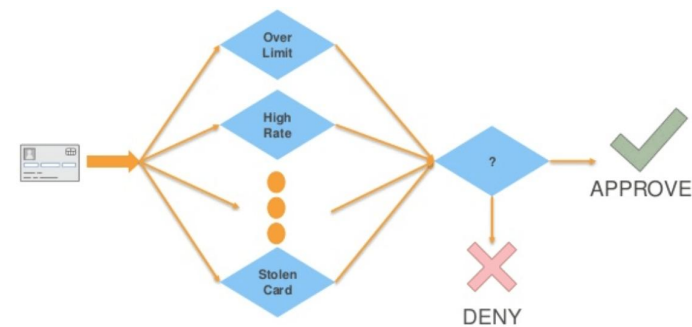
- 垃圾邮件检测



- 文档分类

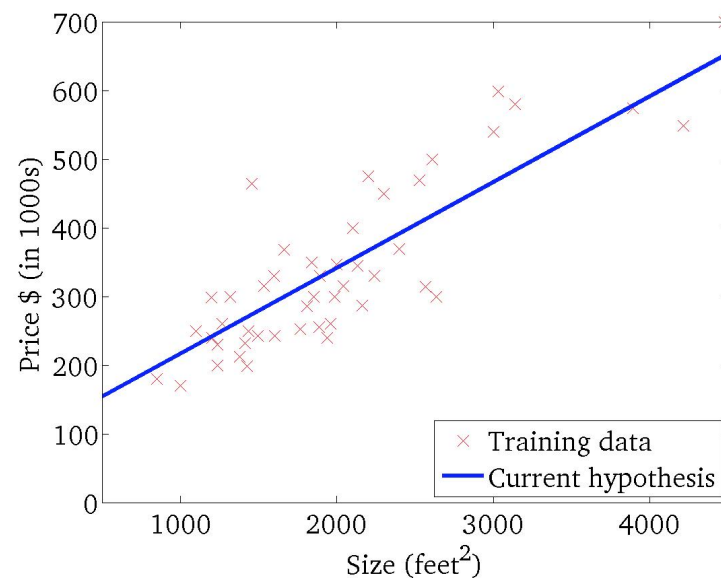


- 交易欺诈检测



逻辑回归

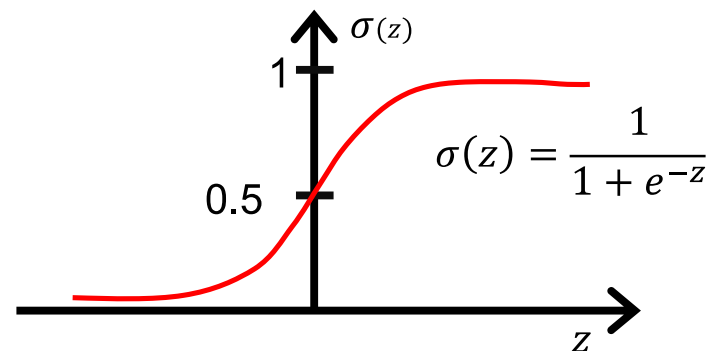
- 在分类问题中，目标变量 $y \in \{0, 1\}$
- 在线性回归中，输出 $f(\mathbf{x}) = \mathbf{xw}$ 的取值范围是 $[-\infty, +\infty]$



- 线性回归的输出值与分类任务中的目标值不兼容

- Sigmoid/logistic 函数

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



- 逻辑回归

$$f(\mathbf{x}) = \sigma(\mathbf{x}\mathbf{w})$$

- 线性回归

$$f(\mathbf{x}) = \mathbf{x}\mathbf{w}$$

- 输出范围从 $[-\infty, +\infty]$ 转换为 $[0, 1]$

代价函数

- 目标
 - 如果真实标签 $y = 1$, 我们希望 $f(\mathbf{x}) = \sigma(\mathbf{x}\mathbf{w})$ 接近于1
 - 如果真实标签 $y = 0$, 我们希望 $f(\mathbf{x}) = \sigma(\mathbf{x}\mathbf{w})$ 接近于0
- 为了实现这个目标, 我们可以定义一个类似于回归中的代价函数:

$$L(\mathbf{w}) = (\sigma(\mathbf{x}\mathbf{w}) - y)^2$$

- 或者，我们也可以寻求最小化：

$$L(\mathbf{w}) = \begin{cases} -\log(\sigma(\mathbf{xw})) & \text{if } y = 1 \\ -\log(1 - \sigma(\mathbf{xw})) & \text{if } y = 0 \end{cases}$$

- 上面的目标可以等价地写成：

$$L(\mathbf{w}) = -y \log(\sigma(\mathbf{xw})) - (1 - y) \log(1 - \sigma(\mathbf{xw}))$$

如果 $y = 1$, $L(\mathbf{w})$ 简化为 $L(\mathbf{w}) = -\log(\sigma(\mathbf{xw}))$;

否则，如果 $y = 0$, $L(\mathbf{w})$ 简化为 $L(\mathbf{w}) = -\log(1 - \sigma(\mathbf{xw}))$

- 上面的损失函数被称为**交叉熵损失**

- 交叉熵损失：

$$L(\mathbf{w}) = -y \log(\sigma(\mathbf{xw})) - (1 - y) \log(1 - \sigma(\mathbf{xw}))$$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

- 二阶导数（海森矩阵）：

$$H = \frac{\partial^2 L(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^\top} = \mathbf{x} \mathbf{x}^\top \sigma(\mathbf{xw})(1 - \sigma(\mathbf{xw}))$$

- 判断半正定性：

$$\mathbf{v}^\top H \mathbf{v} = (\mathbf{x}^\top \mathbf{v})^2 \sigma(\mathbf{xw})(1 - \sigma(\mathbf{xw})) \geq 0$$

- 交叉熵损失函数的海森矩阵是半正定的，所以该函数是一个凸函数

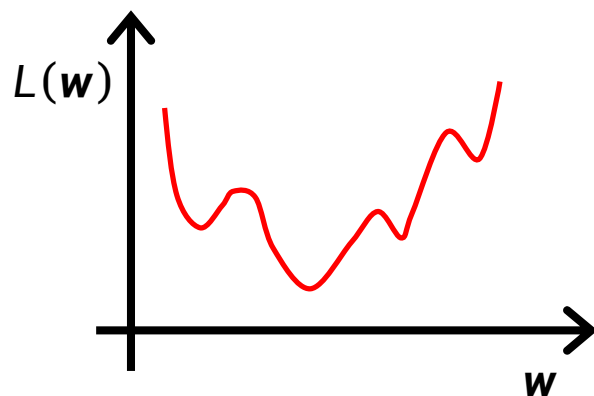
- 哪个代价函数更好？

平方误差: $L(\mathbf{w}) = (\sigma(\mathbf{xw}) - y)^2$

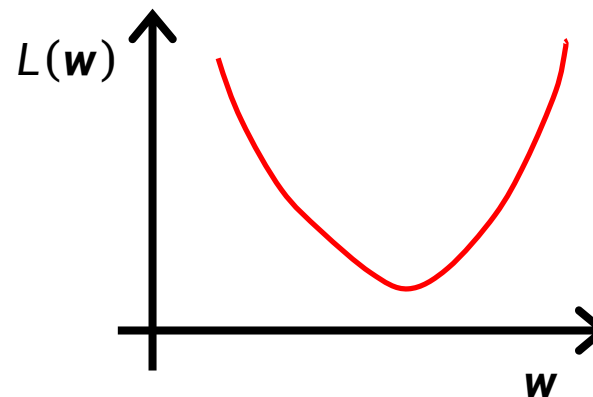
交叉熵: $L(\mathbf{w}) = -[y \log(\sigma(\mathbf{xw})) + (1 - y) \log(1 - \sigma(\mathbf{xw}))]$

函数 $f(z) = \log(\sigma(z))$ 的形状是什么样的？

➤ 平方损失是非凸的



➤ 交叉熵是凸的



凸函数更易于优化

- 在下一次课中，我们将从精确建模的角度，展示使用交叉熵损失的另一个优点

解析解

- 交叉熵损失

$$L(\mathbf{w}) = -y \log(\sigma(\mathbf{xw})) - (1 - y) \log(1 - \sigma(\mathbf{xw}))$$

- 交叉熵损失的梯度

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{N} \sum_{\ell=1}^N [\sigma(\mathbf{x}^{(\ell)} \mathbf{w}) - y^{(\ell)}] \mathbf{x}^{(\ell)T}$$

- 求解 $\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = 0$ 可以得到最优的 \mathbf{w}^* 。将 Sigmoid 函数代入，我们得到：

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{N} \sum_{\ell=1}^N \left[\frac{1}{1 + e^{-\mathbf{x}^{(\ell)} \mathbf{w}}} - y^{(\ell)} \right] \mathbf{x}^{(\ell)T} = 0$$

- 这个方程是关于 \mathbf{w} 的一个非线性方程，无法直接求解这个方程，不能得到一个解析解

梯度下降

- 交叉熵损失的梯度

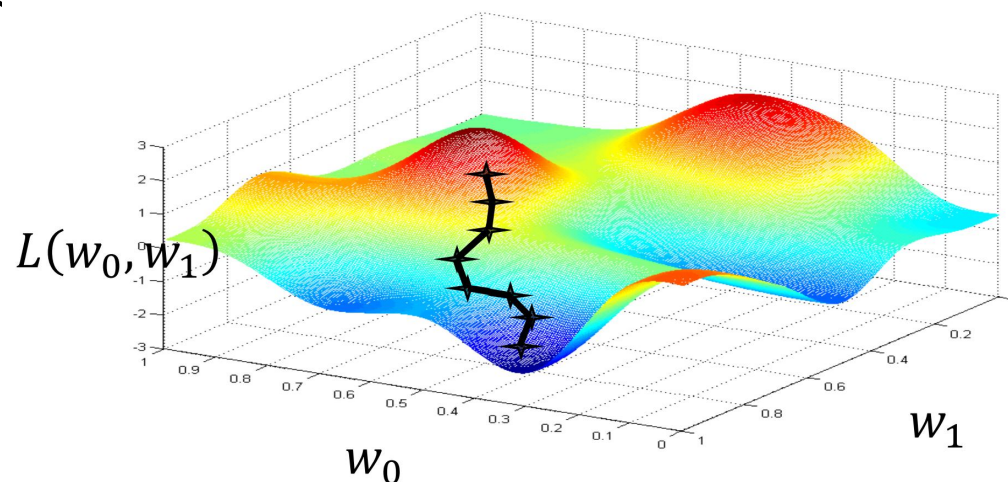
$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{N} \sum_{\ell=1}^N [\sigma(\mathbf{x}^{(\ell)} \mathbf{w}) - y^{(\ell)}] \mathbf{x}^{(\ell)T}$$

- 求解 $\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = 0$ 可以得到最优的 \mathbf{w}^* 。但在这里，*解析解不存在*
- 因此，我们只能求助于数值方法
 - 梯度下降
 - 牛顿法
 - 坐标下降
 -

- 梯度下降法
- 从一个初始的 \mathbf{w}_0 开始，沿着梯度的反方向迭代更新参数，以逐步逼近最小值点

$$\mathbf{w}_{t+1} = \mathbf{w}_t - r \cdot \frac{\partial L(\mathbf{w})}{\partial \mathbf{w}}$$

- r 是学习率，控制每一步更新的幅度
- 收敛性： 由于交叉熵损失是凸函数 梯度下降法保证能收敛到最优值 \mathbf{w}^*



- 梯度的直观理解

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{N} \sum_{\ell=1}^N \left[\underbrace{\sigma(\mathbf{x}^{(\ell)} \mathbf{w}) - y^{(\ell)}}_{\text{预测误差}} \right] \mathbf{x}^{(\ell)T}$$

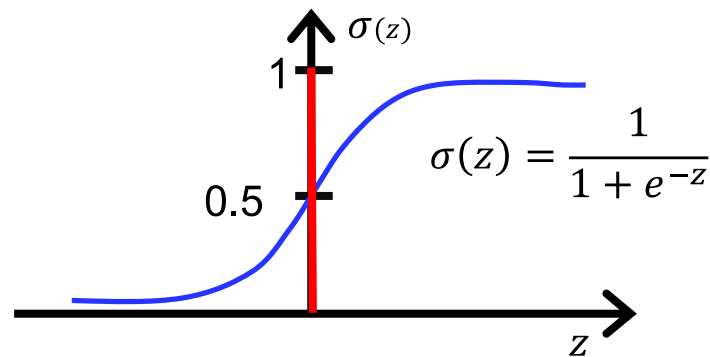
可以发现梯度下降法（GD）总是试图减小预测误差

- 如果 $y^{(\ell)} = 1$ ，算法会驱使 $\sigma(\mathbf{x}^{(\ell)} \mathbf{w})$ 趋近于 1
- 如果 $y^{(\ell)} = 0$ ，算法会驱使 $\sigma(\mathbf{x}^{(\ell)} \mathbf{w})$ 趋近于 0

决策边界

- 决策边界是特征空间中的一个“边界”，它将不同类别的点分隔开
- 样本按以下方式被分为 1 和 0 两类：

$$\hat{y} = \begin{cases} 1, & \text{if } \sigma(\mathbf{x}\mathbf{w}) \geq 0.5 \\ 0, & \text{if } \sigma(\mathbf{x}\mathbf{w}) < 0.5 \end{cases}$$

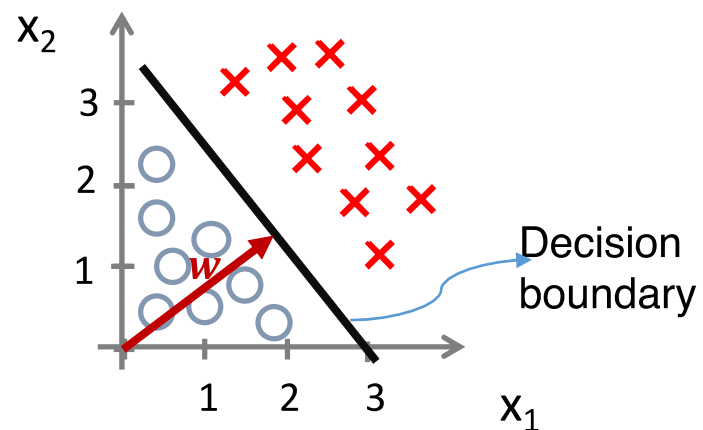


- 这等价于按以下方式对样本进行分类：

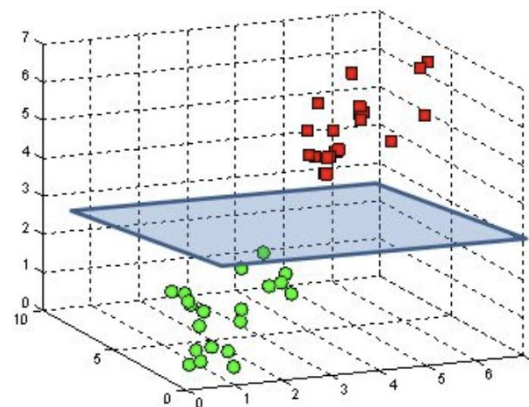
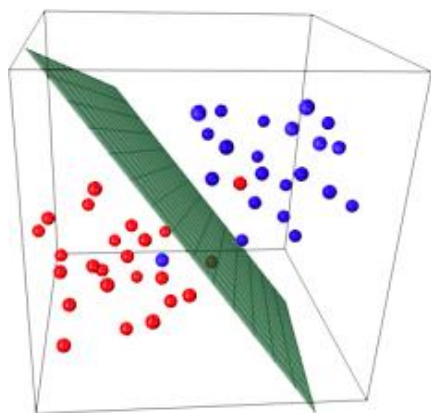
$$\hat{y} = \begin{cases} 1, & \text{if } \mathbf{x}\mathbf{w} \geq 0 \\ 0, & \text{if } \mathbf{x}\mathbf{w} < 0 \end{cases}$$

- 决策边界由满足 $\mathbf{x}\mathbf{w} = 0$ 的 \mathbf{x} 组成

- 由于 w 是一个向量，所有满足 $xw = 0$ 的 x 构成一个与 w 正交的空间



- 在二维情况下，这个空间是一条直线
- 在三维情况下，这个空间是一个平面

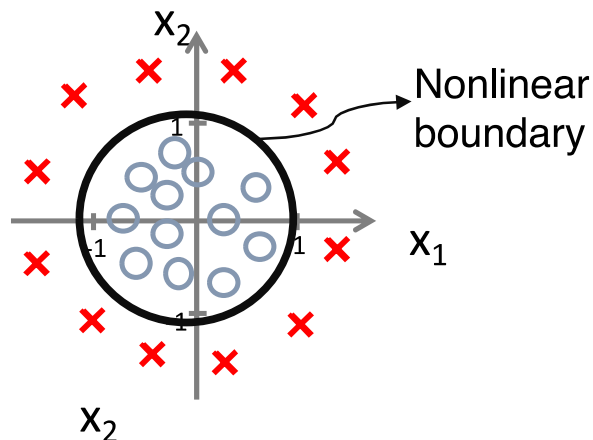


- 对于一个固定的向量 $\mathbf{w} \in \mathbb{R}^K$, 点集

$$\mathbf{x} \in \{\mathbf{x} | \mathbf{x}\mathbf{w} = 0\}$$

构成一个 $(K - 1)$ 维的超平面

- 超平面永远无法表示非线性的决策边界, 例如:

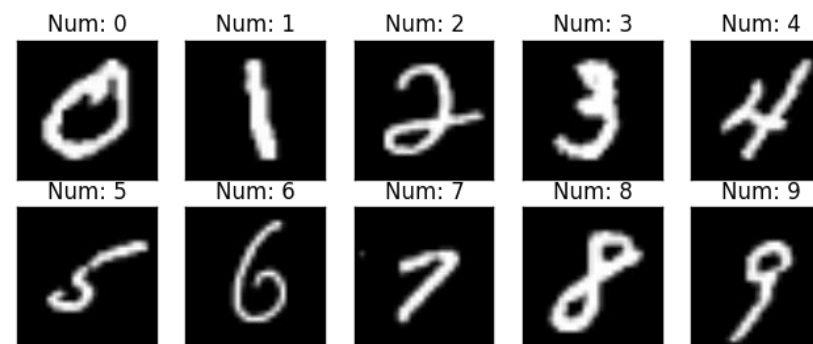


这就是为什么逻辑回归被称为线性分类器的原因

课程大纲

- 二分类
- 多分类

- 许多应用包含多于 2 个类别



- 处理多类别分类的两种方法:

- One-vs-All

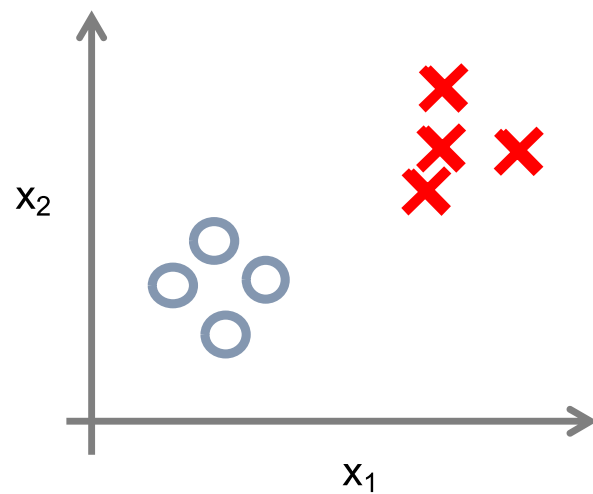
将多类别问题转化为多个二元分类问题

- Softmax 函数

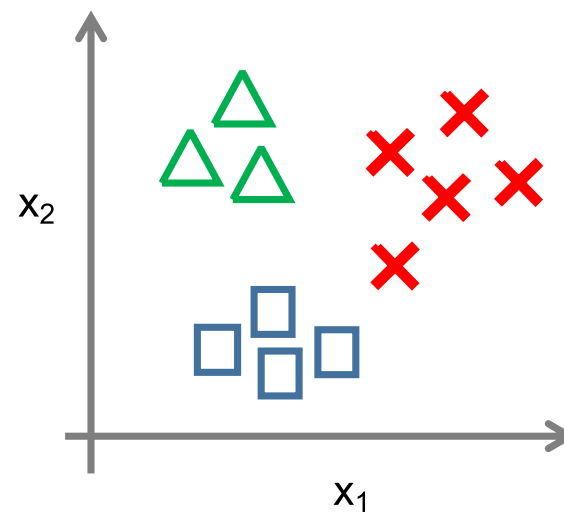
直接将样本分到其中一个类别

One-vs-All

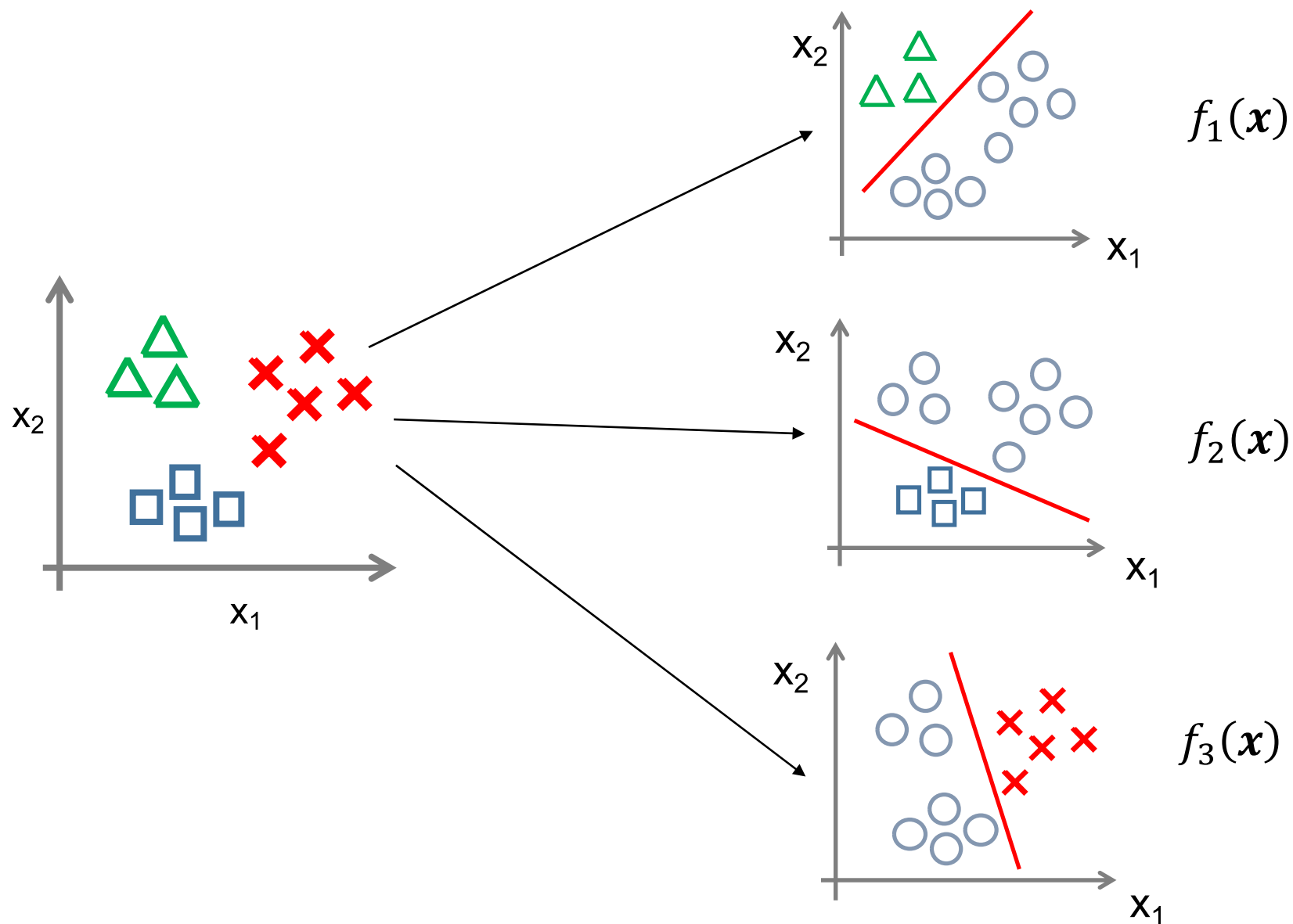
二分类



多分类

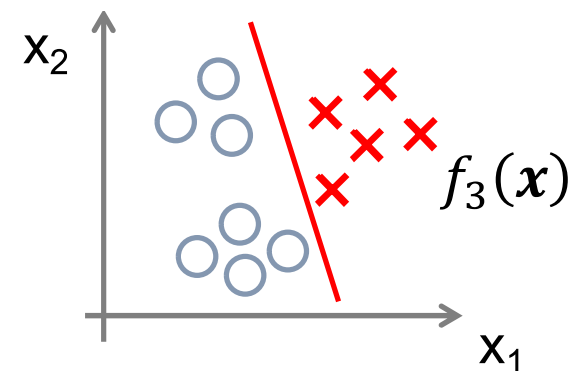
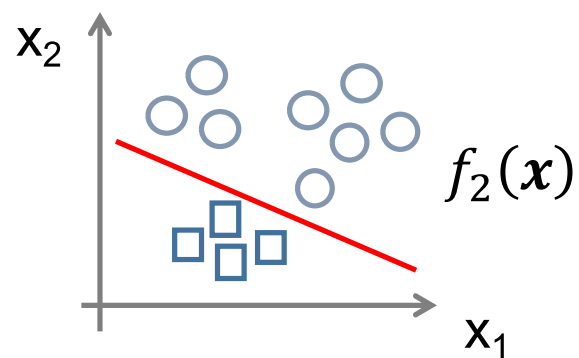
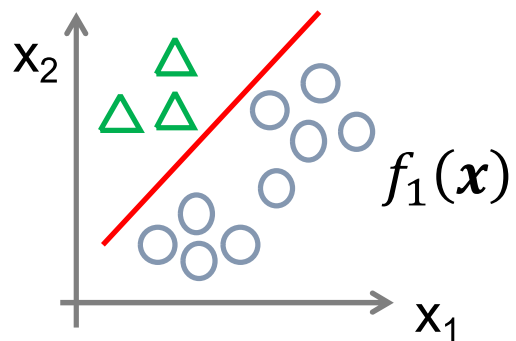


- 为每个类别训练一个分类器



- 为了预测新样本 \mathbf{x} 的类别，选择满足以下条件的类别：

$$k = \arg \max_i f_i(\mathbf{x})$$

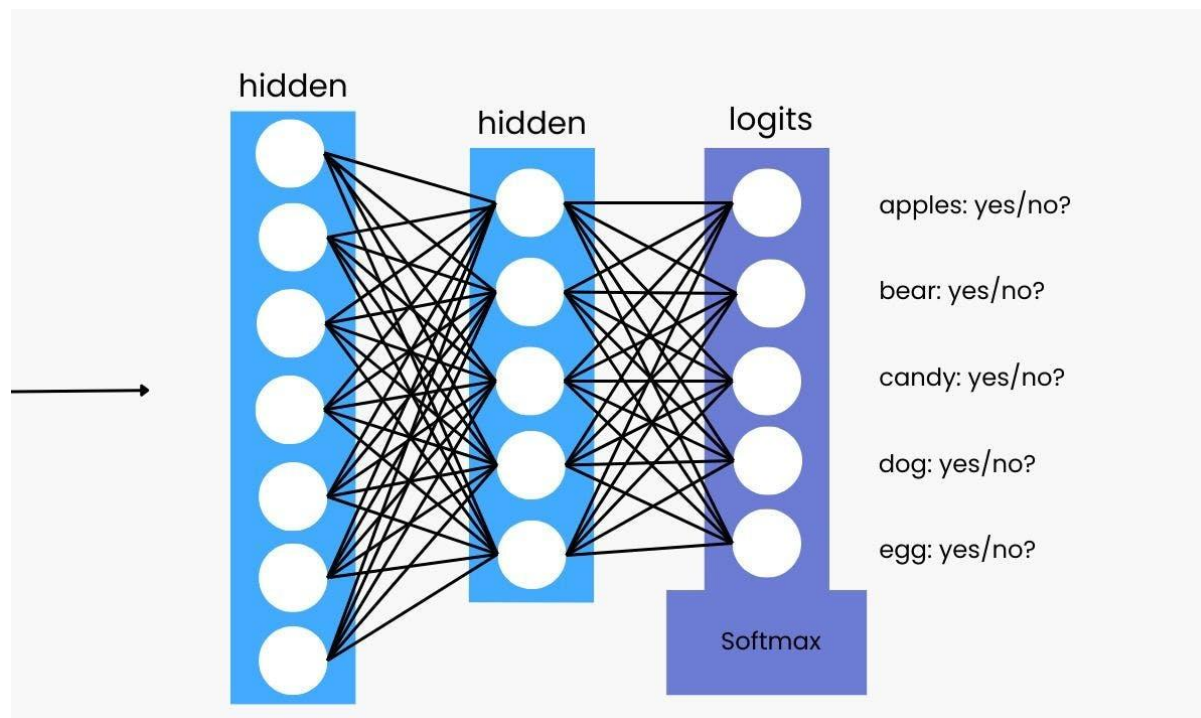


Softmax 函数

- Softmax 函数

$$\text{softmax}_i(\mathbf{z}) = \frac{e^{z_i}}{\sum_{k=1}^K e^{z_k}}$$

可以看出, $\sum_{i=1}^K \text{softmax}_i(\mathbf{z}) = 1$



$$\text{softmax}_i(\mathbf{z}) = \frac{e^{z_i}}{\sum_{k=1}^K e^{z_k}}$$

- 数据 \mathbf{x} 被分类到第 i 个类别的概率是

$$f_i(\mathbf{x}) = \text{softmax}_i(\mathbf{x}\mathbf{W}) = \frac{e^{\mathbf{x}\mathbf{w}_i}}{\sum_{k=1}^K e^{\mathbf{x}\mathbf{w}_k}}$$

其中 $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K]$

- 如果 \mathbf{x} 属于第 i 个类别，模型应该促使 $f_i(\mathbf{x})$ 尽可能地大

- 当 $K = 2$ 时, Softmax 函数等价于 logistic 函数

➤ 在二分类情况下, 我们有:

$$\text{softmax}_1(\mathbf{x}\mathbf{W}) = \frac{e^{xw_1}}{e^{xw_1} + e^{xw_2}}$$

$$= \frac{1}{1 + e^{-x(w_1 - w_2)}}$$

$$\text{softmax}_2(\mathbf{x}\mathbf{W}) = \frac{e^{xw_2}}{e^{xw_1} + e^{xw_2}}$$

$$= \frac{e^{-x(w_1 - w_2)}}{1 + e^{-x(w_1 - w_2)}}$$

➤ 可以看出:

$$\text{softmax}_1(\mathbf{x}\mathbf{W}) = \sigma(\mathbf{x}(\mathbf{w}_1 - \mathbf{w}_2))$$

$$\text{softmax}_2(\mathbf{x}\mathbf{W}) = 1 - \sigma(\mathbf{x}(\mathbf{w}_1 - \mathbf{w}_2))$$

二分类 Softmax 分类等价于逻辑回归, 其模型参数为 $\mathbf{w}_1 - \mathbf{w}_2$

代价函数

- 对于一个有 K 个类别的训练数据集，其标签 y 用一个独热向量 (one-hot vector)

$$\begin{bmatrix} 1, 0, 0, \dots, 0, \\ 0, 1, 0, \dots, 0, \\ \vdots \\ 0, 0, 0, \dots, 1 \end{bmatrix}$$

[1 , 4 , 2 , 0 , 3]


$$\begin{bmatrix} [0, 1, 0, 0, 0] \\ [0, 0, 0, 0, 1] \\ [0, 0, 1, 0, 0] \\ [1, 0, 0, 0, 0] \\ [0, 0, 0, 1, 0] \end{bmatrix}$$

Normal array

One hot encoding

- 目标是最大化相应的概率 $f_i(\mathbf{x})$ 。因此，代价函数可以写为

$$L(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K) = -\frac{1}{N} \sum_{\ell=1}^N \sum_{k=1}^K y_k^{(\ell)} \log [\text{softmax}_k(\mathbf{x}^{(\ell)} \mathbf{w})]$$

— $y_k^{(\ell)}$ 是 $\mathbf{y}^{(\ell)}$ 的第 k 个元素

Cross-entropy 损失

梯度下降

- 关于 \mathbf{w}_j 的梯度是

$$\frac{\partial L(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K)}{\partial \mathbf{w}_j} = \frac{1}{N} \sum_{\ell=1}^N \left(\underbrace{\text{softmax}_j(\mathbf{x}^{(\ell)} \mathbf{W}) - y_j^{(\ell)}}_{\text{预测误差}} \right) \mathbf{x}^{(\ell)T}$$

注意，对于 $j = 1, \dots, K$ ，所有的 \mathbf{w}_j 都应该同时更新

$$\frac{\partial L(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K)}{\partial \mathbf{w}_j} = \frac{1}{N} \sum_{\ell=1}^N \left(\underbrace{\text{softmax}_j(\mathbf{x}^{(\ell)} \mathbf{W}) - y_j^{(\ell)}}_{\text{预测误差}} \right) \mathbf{x}^{(\ell)T}$$

- 将 \mathbf{W} 表示为 $[\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K]$, 我们有:

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}} = \frac{1}{N} \sum_{\ell=1}^N \mathbf{x}^{(\ell)T} (\text{softmax}(\mathbf{x}^{(\ell)} \mathbf{W}) - \mathbf{y}^{(\ell)})$$

- $\text{softmax}(\mathbf{x}^{(\ell)} \mathbf{W}) = [\text{softmax}_1(\mathbf{x}^{(\ell)} \mathbf{W}), \dots, \text{softmax}_K(\mathbf{x}^{(\ell)} \mathbf{W})]$ 是一个行向量

- 更新: $\mathbf{W}_{t+1} = \mathbf{W}_t - r \cdot \frac{\partial L(\mathbf{W})}{\partial \mathbf{W}} \Big|_{\mathbf{W}=\mathbf{W}_t}$