

Regression: how much
Classification: which one.

3-1 One-hot encoding (独热编码)

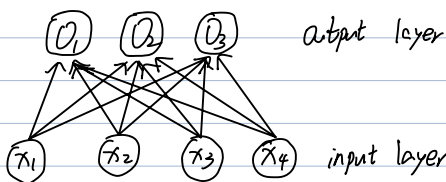
One-hot encoding 是一个 vector, 其 component 的数量就是类别数.

比如: 我们希望把一张图分入三类中: 猫, 鸡, 狗

那么 $y \in \{ \overset{\text{猫}}{\langle 1, 0, 0 \rangle}, \overset{\text{鸡}}{\langle 0, 1, 0 \rangle}, \overset{\text{狗}}{\langle 0, 0, 1 \rangle} \}$

如果我们有 4 个 feature 即每个 sample $\vec{x} = \langle x_1, x_2, x_3, x_4 \rangle$

那么我们的 NN 是:



(此为 single layer NN)

$$\begin{cases} O_1 = \vec{w}_1^T \vec{x} + b_1 \\ O_2 = \vec{w}_2^T \vec{x} + b_2 \\ O_3 = \vec{w}_3^T \vec{x} + b_3 \end{cases} \Rightarrow \vec{O} = W \cdot \vec{x} + \vec{b}$$

$$\hat{y} = \text{softmax}(\vec{O}), \text{ where } \hat{y}_i = \frac{\exp(O_i)}{\sum_j \exp(O_j)}$$

$$\text{即 } \begin{bmatrix} O_1 \\ O_2 \\ O_3 \end{bmatrix} \mapsto \begin{bmatrix} \frac{e^{O_1}}{e^{O_1} + e^{O_2} + e^{O_3}} \\ \frac{e^{O_2}}{e^{O_1} + e^{O_2} + e^{O_3}} \\ \frac{e^{O_3}}{e^{O_1} + e^{O_2} + e^{O_3}} \end{bmatrix}$$

input: n 个 sample, d 个 features

$$X = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_d^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \dots & x_d^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(n)} & x_2^{(n)} & \dots & x_d^{(n)} \end{bmatrix} \in \mathbb{R}^{n \times d}$$

$$q \text{ 个类别 } (\hat{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_q \end{bmatrix})$$

$$W = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1d} \\ w_{21} & w_{22} & \dots & w_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ w_{d1} & w_{d2} & \dots & w_{dq} \end{bmatrix} \in \mathbb{R}^{d \times q} \quad \vec{b} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_q \end{bmatrix} \in \mathbb{R}^{1 \times q}$$

(广播后为 $n \times q, n$ 个相同 \vec{b})

这里用 3 个 affine function 得到 3 个 weights vector 和 3 个 bias. (每个有 4 个 component)

合成一个 3×4 W/matrix 和一个 3×1 \vec{b} vector.

$$O_i = w_{i1}x_1 + w_{i2}x_2 + w_{i3}x_3 + w_{i4}x_4 + b_i$$

$$\vec{O} = \begin{bmatrix} O_1 \\ O_2 \\ O_3 \end{bmatrix}$$

这里由 one-hot wde 得出的是 3 个 logit (未经规范化的预测)

我们希望得到的是: $\langle 1, 0, 0 \rangle, \langle 0, 1, 0 \rangle, \langle 0, 0, 1 \rangle$

这样的 hard labels.

但我们首先得到由 prob 形成的 soft label.

比如 $\langle 0, 0.5, 0.5 \rangle$ 表示 50% 为鸡, 50% 为狗

然后再把 soft labels 转化 hard labels.

然而这里 \vec{O} 甚至不是一个 soft label, 因为它 components 和不为 1, 并且 components 值可以为负.

因而我们需要先把 \vec{O} 转为 soft label, 再转化为 hard label.

3-2 Softmax function.

softmax 可以把 \vec{O} 变为总和为 1 的 $[0, 1]$ 之间的数.

$$= \begin{bmatrix} p_1 \\ \vdots \\ p_q \end{bmatrix} \text{ 的 components}$$

$$O = XW + \vec{b} = \begin{bmatrix} \vec{O}_1 & \vec{O}_2 & \dots & \vec{O}_q \end{bmatrix} \quad \text{(广播)} \quad \begin{matrix} n \text{ 个 samples.} \\ \in \mathbb{R}^{n \times q} \end{matrix}$$

$$\vec{O}_1 = \begin{bmatrix} w_{11}x_1^{(1)} + w_{12}x_2^{(1)} + \dots + w_{1d}x_d^{(1)} + b_1 \\ w_{11}x_1^{(2)} + w_{12}x_2^{(2)} + \dots + w_{1d}x_d^{(2)} + b_1 \\ \vdots \\ w_{11}x_1^{(n)} + w_{12}x_2^{(n)} + \dots + w_{1d}x_d^{(n)} + b_1 \end{bmatrix}$$

$$\vec{O}_q = \begin{bmatrix} w_{q1}x_1^{(1)} + w_{q2}x_2^{(1)} + \dots + w_{qd}x_d^{(1)} + b_q \\ w_{q1}x_1^{(2)} + w_{q2}x_2^{(2)} + \dots + w_{qd}x_d^{(2)} + b_q \\ \vdots \\ w_{q1}x_1^{(n)} + w_{q2}x_2^{(n)} + \dots + w_{qd}x_d^{(n)} + b_q \end{bmatrix}$$

即: 由 design matrix X predict 得到

Consider: $\{X, Y\}$ 有 n 个 samples.

- 1 sample i 由一个 $\vec{x}^{(i)} \in \mathbb{R}^d$ 和 $\vec{y}^{(i)} \in \mathbb{R}^2$
(d features) (q classes)
 $[x_1, x_2, \dots, x_d]^T$ $[y_1, y_2, \dots, y_q]^T$

那么 $p(Y|X) = \prod_{i=1}^n p(\vec{y}^{(i)} | \vec{x}^{(i)})$

3-2 Cross-entropy loss.

注意到 $p(\vec{y}^{(i)} | \vec{x}^{(i)}) = \prod_{j=1}^q (\hat{y}_j)^{y_j}$

为什么呢? 因为 $\vec{y}^{(i)}$ 有一个 component 为 1, 其余都是 0.

所以 $\vec{y}^{(i)}$ 中, $y_j = 0$ 的 component 上, $(\hat{y}_j)^{y_j} = 1$

$y_j = 1$ 的 component 上, $(\hat{y}_j)^{y_j} = \hat{y}_j$

因而 $\prod_{j=1}^q (\hat{y}_j)^{y_j} = 1 \cdot 1 \cdot \dots \cdot \hat{y}_{\text{correct}} \cdot 1 \cdot 1 \cdot \dots \cdot 1 = \hat{y}_{\text{correct}}$

因为 softmax 算法, 所有 \hat{y}_j 的和为 1

所以只会保留预测到的结果正确为 $\vec{y}^{(i)}$ 的概率

也就是在给定 \vec{x} , 下列模型正确从 \vec{y} 推出 \vec{y} 的概率即 $p(\vec{y} | \vec{x}^{(i)})$

那么同样, 求 negative log likelihood 使其 minimize:

$$-\ln(p(Y|X)) = \sum_{i=1}^n -\ln(p(\vec{y}^{(i)} | \vec{x}^{(i)}))$$

我们已说到 $p(\vec{y}^{(i)} | \vec{x}^{(i)}) = \prod_{j=1}^q (\hat{y}_j)^{y_j}$

因而 $= \sum_{i=1}^n -\ln\left(\prod_{j=1}^q (\hat{y}_j)^{y_j}\right)$

$= \sum_{i=1}^n \left(-\sum_{j=1}^q y_j \ln \hat{y}_j \right)$

我们定义这个式为 cross-entropy loss.

$$L(\vec{y}, \hat{\vec{y}}) = -\sum_{j=1}^q y_j \log \hat{y}_j$$

如果进一步拆解: $L(\vec{y}, \hat{\vec{y}}) = -\sum_{j=1}^q y_j \ln \frac{\exp(o_j)}{\sum_{k=1}^q \exp(o_k)}$

$= \ln \sum_{k=1}^q \exp(o_k) - \sum_{j=1}^q y_j o_j$

$\Rightarrow \partial_{o_j} L(\vec{y}, \hat{\vec{y}}) = \frac{\exp(o_j)}{\sum_{k=1}^q \exp(o_k)} - y_j = \text{softmax}(\vec{o})_j - y_j$