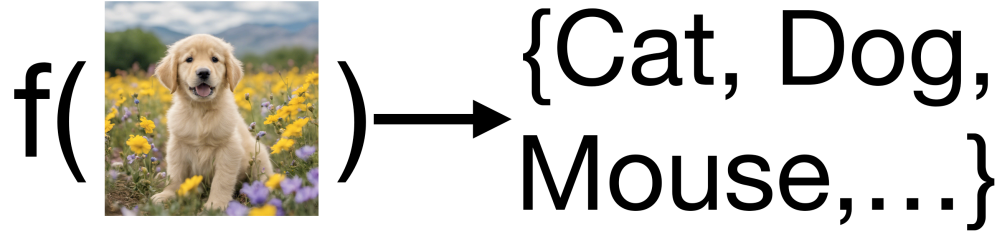


Multi-class classification



$$y = f(\mathbf{x}), \quad \text{Input: } \mathbf{x} \in \mathbb{R}^n \rightarrow \text{Output: } y \in \{1, 2, \dots, C\}$$

Ordering irrelevant!

1: Cat, 2: Dog, 3: Mouse

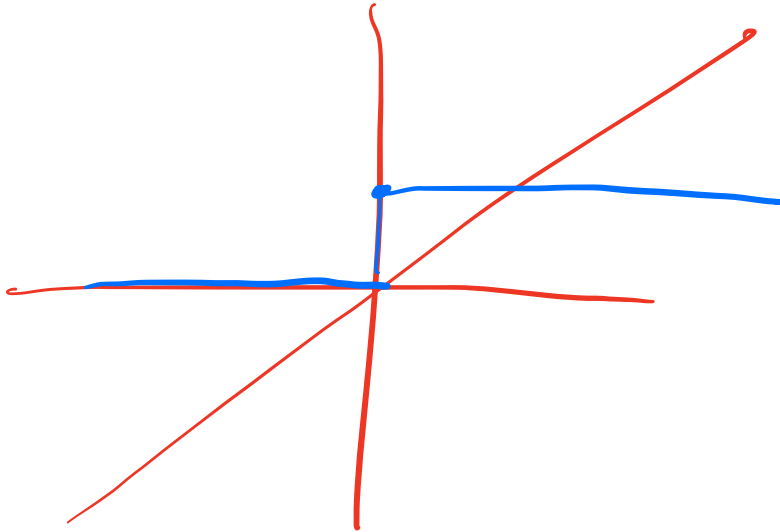
1: Dog, 2: Mouse, 3: Cat

Multi-class prediction functions

Binary thresholding

$$f(\mathbf{x}) = \mathbb{I}(\underbrace{\mathbf{x}^T \mathbf{w}}_{\text{linear function}} \geq 0) \rightarrow \{0, 1\}$$

Indicator



Multi-class thresholding

$$f(\mathbf{x}) = \operatorname{argmax}_{c \in \{1 \dots C\}} \underbrace{\mathbf{x}^T \mathbf{w}_c}_{\text{make } C \text{ different predictions}}$$

Multi-class prediction functions

Multi-class thresholding

$$f(\mathbf{x}) = \underset{c \in \{1 \dots C\}}{\operatorname{argmax}} \mathbf{x}^T \mathbf{w}_c$$

\hookrightarrow linear / ~~line~~ func.

$C \times d$ total params.
 \uparrow

num of weights

of classes / scalar

$$(\mathbf{x}^T \mathbf{w}^T)_c = \boxed{\mathbf{x}^T \mathbf{w}_c}$$

$(1 \times d)(d \times C) \rightarrow (1 \times C)$ vector

For convenience, we can also define a matrix that contains all C parameter vectors:

$$\mathbf{W} = \begin{bmatrix} \mathbf{w}_1^T \\ \mathbf{w}_2^T \\ \vdots \\ \mathbf{w}_C^T \end{bmatrix} = \begin{bmatrix} W_{11} & W_{12} & \dots & W_{1d} \\ W_{21} & W_{22} & \dots & W_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ W_{C1} & W_{C2} & \dots & W_{Cd} \end{bmatrix}$$

$C \times d$

With this notation, our prediction function becomes:

$$f(\mathbf{x}) = \underset{c \in \{1 \dots C\}}{\operatorname{argmax}} (\mathbf{x}^T \mathbf{W}^T)_c, \quad \mathbf{W} \in \mathbb{R}^{C \times d}$$

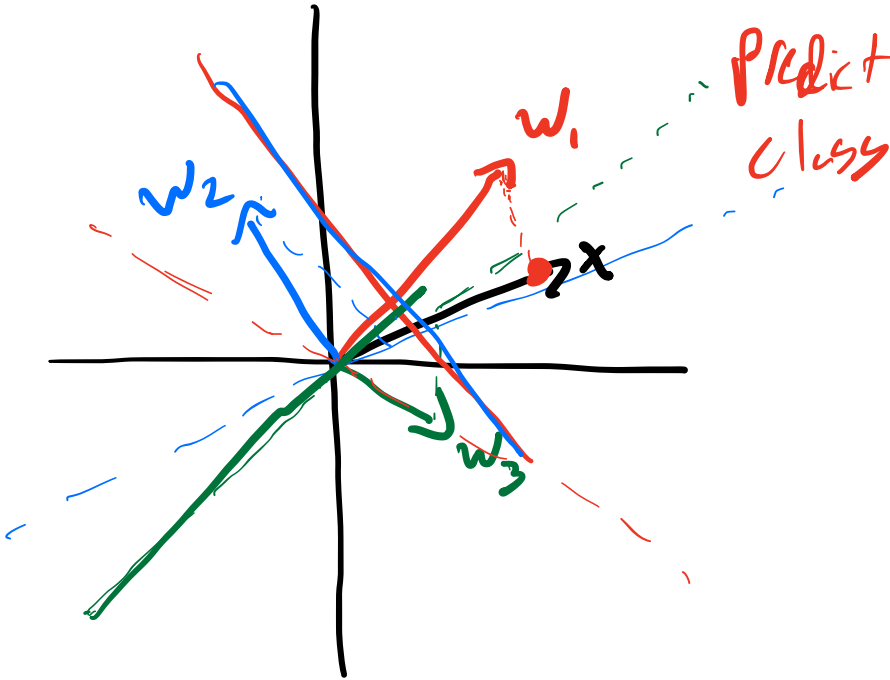
\mathbf{x} : $d \times 1$ column

\mathbf{x}^T : $1 \times d$ row vector

\mathbf{W} : $C \times d$ matrix, \mathbf{W}^T : $d \times C$ matrix

Multi-class decision boundaries

$$f(\mathbf{x}) = \underset{c \in \{0,1\}}{\operatorname{argmax}} (\mathbf{x}^T \mathbf{W}^T)_c = \mathbb{I}(\mathbf{x}^T \mathbf{w}_1 - \mathbf{x}^T \mathbf{w}_0 \geq 0)$$



$$C = 2$$

$$\mathbf{x}^T \mathbf{w}_0 \geq \mathbf{x}^T \mathbf{w}_1 \quad ? \begin{cases} \text{if } t_{x,0} = 1 \\ \text{if } t_{x,0} = 0 \end{cases}$$

$$\mathbf{x}^T \mathbf{w}_0 - \mathbf{x}^T \mathbf{w}_1 \geq 0$$

↓ 2 parameters

$$\mathbf{x}^T (\mathbf{w}_0 - \mathbf{w}_1) \geq 0$$

def. \mathbf{w}

logistic
Reg

$$= \mathbf{x}^T \mathbf{w} \geq 0$$

total probs
↓
prob. of outcome c
↓

Categorical distribution

$$p(y=c) = q_c, \quad y \in \{1 \dots C\}$$

$$\mathbf{q} \in \mathbb{R}^C \quad q_c \geq 0 \forall c \in \{1 \dots C\} \quad \sum_{c=1}^C q_c = 1$$

Sum to 1

Fair die: $\{1, 2, 3, 4, 5, 6\}$

$$p(y=c) = \frac{1}{6} \quad \frac{1}{6} \quad \frac{1}{6} \quad \frac{1}{6} \quad \frac{1}{6} \quad \frac{1}{6}$$

A probabilistic model for multi-class classification

Model $p(y|x)$ = Categorical

$y_i \sim \text{Categorical}(\mathbf{q}=?),$
Output / label

$$\mathbf{q} = \mathbf{x}_i^T \mathbf{W}^T?$$

Vector length C

$$\mathbf{x}^T \mathbf{W}^T \in \mathbb{R}^C, \quad q_c \geq 0 \quad \forall c \in \{1 \dots C\}, \quad \sum_{c=1}^C q_c = 1$$

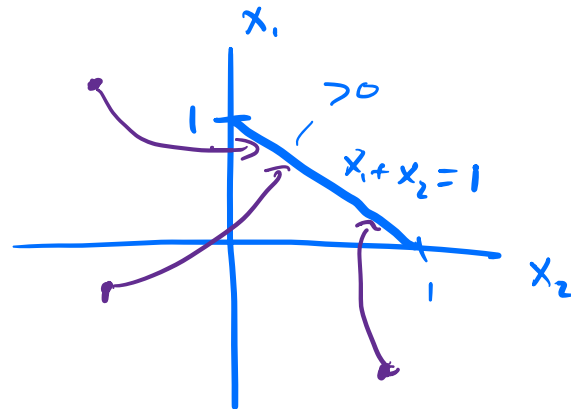
$q=?$

$\begin{bmatrix} \\ \\ \end{bmatrix}$

Need $f(\mathbf{x}) : \mathbb{R}^C \rightarrow [0, \infty)^C, \quad \sum_{i=1}^C f(\mathbf{x})_c = 1$

logistic Reg. 

$\mathbf{x}^T \mathbf{W} \in \mathbb{R} \xrightarrow{\sigma} [0, 1]$
Sigmoid func.



Categorical distribution

$$p(y = c) = q_c, \quad y \in \{1 \dots C\}$$

Defn.

$$P(y=1) = q$$

$$P(y) = (q)^y (1-q)^{(1-y)}$$

$$\mathbf{q} \in \mathbb{R}^C \quad q_c \geq 0 \quad \forall c \in \{1 \dots C\} \quad \sum_{c=1}^C q_c = 1$$

total # of classes or categories

index of the summation

$$p(y) = \prod_{c=1}^C q_c^{\mathbb{I}(y=c)} = q_y$$

q_1 if $y=1$

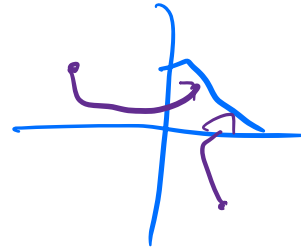
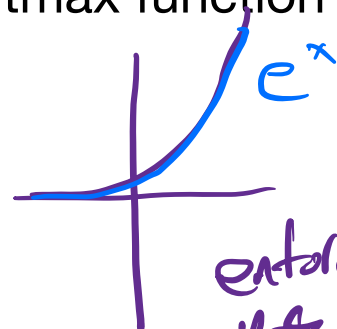
q_2 if $y=2$

$$\log p(y) = \log \prod_c q_c^{\mathbb{I}(y=c)} = \sum_{c=1}^C \mathbb{I}(y=c) \log q_c$$

$$\mathbb{R}^C \rightarrow \mathbb{R}^C$$

Softmax function

$$\text{softmax}(\mathbf{x})_c = \frac{e^{x_c}}{\sum_{j=1}^C e^{x_j}}$$



$$\text{softmax}(\mathbf{x}) = \begin{bmatrix} \frac{e^{x_1}}{\sum_{j=1}^C e^{x_j}} \\ \frac{e^{x_2}}{\sum_{j=1}^C e^{x_j}} \\ \vdots \\ \frac{e^{x_C}}{\sum_{j=1}^C e^{x_j}} \end{bmatrix}$$

1) enforce $\sum_{j=1}^C e^{x_j} = 1$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_C \end{bmatrix} \rightarrow \begin{bmatrix} e^{x_1} \\ e^{x_2} \\ \vdots \\ e^{x_C} \end{bmatrix}$$

2) enforce sum to 1

$$T = \sum_{j=1}^C e^{x_j}$$

$$\frac{1}{T} \sum_{j=1}^C e^{x_j}$$

$$\begin{bmatrix} \frac{e^{x_1}}{T} \\ \vdots \\ \frac{e^{x_C}}{T} \end{bmatrix} \rightarrow \begin{bmatrix} \frac{e^{x_1}}{\sum_{j=1}^C e^{x_j}} \\ \vdots \\ \frac{e^{x_C}}{\sum_{j=1}^C e^{x_j}} \end{bmatrix}$$

Softmax function

$$\text{softmax}(\mathbf{x})_c = \frac{e^{x_c}}{\sum_{j=1}^C e^{x_j}}$$

$$\sum_{i=1}^C \frac{e^{x_i}}{\sum_{j=1}^C e^{x_j}} = \frac{\sum_{i=1}^C e^{x_i}}{\sum_{j=1}^C e^{x_j}} = 1$$

$$\operatorname{argmax}_{c \in \{1, \dots, C\}} \mathbf{x}_c = \operatorname{argmax}_{c \in \{1, \dots, C\}} \text{softmax}(\mathbf{x})_c$$

A probabilistic model for multi-class classification: Multinomial Logistic regression

$y_i \sim \text{Categorical}(\text{softmax}(\mathbf{x}^T \mathbf{W}))$

linear func \mathcal{L}

$$p(y_i = c \mid \mathbf{x}, \mathbf{W}) = \text{softmax}(\mathbf{x}^T \mathbf{W})_c = \frac{e^{\mathbf{x}^T \mathbf{w}_c}}{\sum_{j=1}^C e^{\mathbf{x}^T \mathbf{w}_j}}$$

input parameters index \mathcal{L} \leftarrow removed all linear func. exp.

Maximum likelihood estimation for Multinomial Logistic regression

Model

Loss

$$y_i \sim \text{Categorical}(\text{softmax}(\mathbf{x}^T \mathbf{W}))$$

$$p(y_i = c | \mathbf{x}, \mathbf{W}) = \text{softmax}(\mathbf{x}^T \mathbf{W})_c = \frac{e^{\mathbf{x}^T \mathbf{w}_c}}{\sum_{j=1}^C e^{\mathbf{x}^T \mathbf{w}_j}}$$

$$\text{Loss}(\mathbf{W}) = \text{NLL}(\mathbf{W}, \mathbf{X}, \mathbf{y}) = - \sum_{i=1}^N \log p(y_i | \mathbf{x}_i, \mathbf{W})$$

output Input

$$\text{NLL}(\mathbf{W}, \mathbf{X}, \mathbf{y}) = - \sum_{i=1}^N \log \left[\frac{e^{\mathbf{x}_i^T \mathbf{w}_{y_i}}}{\sum_{j=1}^C e^{\mathbf{x}_i^T \mathbf{w}_j}} \right]$$

$$= - \sum_{i=1}^N \log e^{\mathbf{x}_i^T \mathbf{w}_{y_i}} - \log \sum_{j=1}^C e^{\mathbf{x}_i^T \mathbf{w}_j}$$

$$= - \sum_{i=1}^N \mathbf{x}_i^T \mathbf{w}_{y_i} - \log \sum_{j=1}^C e^{\mathbf{x}_i^T \mathbf{w}_j}$$