

Lecture 4

Machine Learning – III

Classification



Outline

- Brief review
- Binary Classification
- Multi-class classification

Multivariate Linear Regression

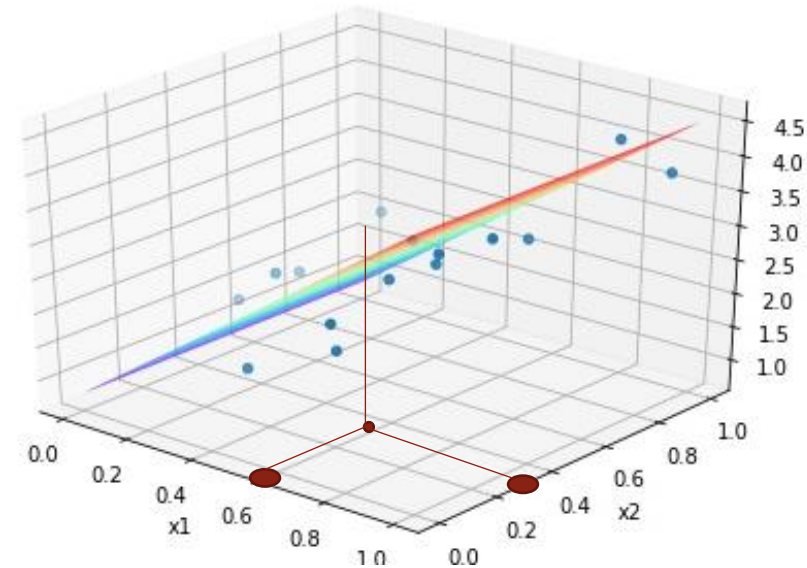
Size	0.488	0.681	0.655	0.088	0.139	0.721	0.999	0.775	0.881	0.007	0.14	0.817	0.681	0.851	0.436
Color	0.609	0.112	0.324	0.669	0.91	0.41	0.902	0.353	0.865	0.539	0.524	0.62	0.093	0.447	0.088
prize	2.737	2.237	2.869	1.626	2.22	3.038	3.929	3.351	4.308	1.304	1.917	3.189	2.656	3.517	1.663

- Data:

- Sample: $(\mathbf{x} \in R^n, y \in R)$, where $n=2$ is the *dimension* of the input vector \mathbf{x} .
- Dataset: $D = \{(\mathbf{x}^i \in R^2, y^i \in R) | i \in [1, m]\}$, where $m=15$ is the number of the samples.



$$\left(\mathbf{x}^i = \begin{bmatrix} x_1^i \\ x_2^i \end{bmatrix}, y^i \right)$$



Multivariate Linear Regression

Multivariate Linear Regression

- Data:

- Sample: $(\mathbf{x} \in R^n, y \in R)$, where $n=2$ is the *dimension* of the input vector \mathbf{x} .

- Dataset:

$D = \{(\mathbf{x}^i \in R^n, y^i \in R) | i \in [1, m]\}$,
where $m=15$ is the number of the samples.

- Linear Model:

$$a = \mathbf{w}\mathbf{x} + b = w_1x_1 + w_2x_2 + b,$$

$$\text{where } \mathbf{w} = [w_1, w_2], \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

- Object: $\operatorname{argmin}_{w,b} J(w, b) = \operatorname{argmin}_{w,b} \frac{1}{2m} \sum_{i=1}^m (a^i - y^i)^2$

Size	0.488	0.681	0.655	0.088
Color	0.609	0.112	0.324	0.669
bias	1	1	1	1
prize	2.737	2.237	2.869	1.626

$$\begin{aligned} a^i &= \mathbf{w}\mathbf{x}^i + b \\ &= w_1x_1^i + w_2x_2^i + b \cdot 1 \\ &= w_1x_1^i + w_2x_2^i + w_3 \cdot 1 \\ &= w_1x_1^i + w_2x_2^i + w_3 \cdot x_3^i \\ &= \sum_{j=1}^n w_j \cdot x_j^i \end{aligned}$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \longrightarrow \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix}$$

Multivariate Linear Regression

Multivariate Linear Regression

- Data:
 - Sample: $(\mathbf{x} \in R^n, y \in R)$, where $n=2$ is the *dimension* of the input vector \mathbf{x} .
 - Dataset: $D = \{(\mathbf{x}^i \in R^n, y^i \in R) | i \in [1, m]\}$, where $m=15$ is the number of the samples.
- Linear Model:
$$a = \mathbf{w}\mathbf{x} + b = w_1x_1 + w_2x_2 + b,$$
where $\mathbf{w} = [w_1, w_2]$, $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$
- Object:
$$\operatorname{argmin}_{w,b} J(w, b) = \operatorname{argmin}_{w,b} \frac{1}{2m} \sum_{i=1}^m (a^i - y^i)^2$$

Multivariate Linear Regression

- Data:
 - Sample: $(\mathbf{x} \in R^n, y \in R)$, where $n=3$ is the *dimension* of the input vector \mathbf{x} and $x_3=1$.
 - Dataset: $D = \{(\mathbf{x}^i \in R^n, y^i \in R) | i \in [1, m]\}$, where $m=15$ is the number of the samples.
- Linear Model:
$$a = \mathbf{w}\mathbf{x} = \sum_{j=1}^n w_j \cdot x_j^i$$
where $\mathbf{w} = [w_1, w_2, w_3]$, $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$
- Object:
$$\operatorname{argmin}_{\mathbf{w}} J(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w}} \frac{1}{2m} \sum_{i=1}^m (a^i - y^i)^2$$

Least Squares Approximations

Multivariate Linear Regression

- Data:
 - Sample: $(\mathbf{x} \in R^n, y \in R)$, where $n=3$ is the *dimension* of the input vector \mathbf{x} and $x_3=1$.
 - Dataset: $D = \{(\mathbf{x}^i \in R^n, y^i \in R) | i \in [1, m]\}$, where $m=15$ is the number of the samples.

- Linear Model:

$$a = \mathbf{w}\mathbf{x} = \sum_{j=1}^n w_j \cdot x_j^i$$

$$\text{where } \mathbf{w} = [w_1, w_2, w_3], \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

- Object:

$$\operatorname{argmin}_{\mathbf{w}} J(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w}} \frac{1}{2m} \sum_{i=1}^m (a^i - y^i)^2$$

For a differentiable function $J(\mathbf{w})$, if \mathbf{w}^* is a minimum point of J , then the following equation holds:

$$\left. \frac{\partial J}{\partial \mathbf{w}} \right|_{\mathbf{w}=\mathbf{w}^*} = 0$$



$$\begin{cases} \left. \frac{\partial J}{\partial w_1} \right|_{w_1=w_1^*} = 0 \\ \left. \frac{\partial J}{\partial w_2} \right|_{w_2=w_2^*} = 0 \\ \left. \frac{\partial J}{\partial w_3} \right|_{w_3=w_3^*} = 0 \end{cases}$$

Least Squares Approximations

$$\begin{aligned}
 \frac{1}{m} [(a^1 - y^1), \dots (a^i - y^i), \dots, (a^m - y^m)] & \begin{bmatrix} x_1^1 \\ x_1^2 \\ \vdots \\ x_1^i \\ \vdots \\ x_1^m \end{bmatrix} &= & \frac{\partial J}{\partial w_1} \\
 \frac{1}{m} [(a^1 - y^1), \dots (a^i - y^i), \dots, (a^m - y^m)] & \begin{bmatrix} x_2^1 \\ x_2^2 \\ \vdots \\ x_2^i \\ \vdots \\ x_2^m \end{bmatrix} &= & \frac{\partial J}{\partial w_2} \\
 \frac{1}{m} [(a^1 - y^1), \dots (a^i - y^i), \dots, (a^m - y^m)] & \begin{bmatrix} x_3^1 \\ x_3^2 \\ \vdots \\ x_3^i \\ \vdots \\ x_3^m \end{bmatrix} &= & \frac{\partial J}{\partial w_3}
 \end{aligned}$$

Least Squares Approximations

$$\frac{1}{m} [(a^1 - y^1), \dots, (a^i - y^i), \dots, (a^m - y^m)] X^T = \left[\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2}, \frac{\partial J}{\partial w_3} \right]$$

$$\frac{1}{m} (\mathbf{w}X - \mathbf{y})X^T = \left[\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2}, \frac{\partial J}{\partial w_3} \right]$$

$$(\mathbf{w}X - \mathbf{y})X^T = \mathbf{0}$$

$$\mathbf{w}XX^T = \mathbf{y}X^T$$

$$\mathbf{w} = \mathbf{y}X^T (XX^T)^{-1}$$

Steepest Gradient Descend Method

Multivariate Linear Regression

- Data:
 - Sample: $(\mathbf{x} \in R^n, y \in R)$, where $n=3$ is the *dimension* of the input vector \mathbf{x} and $x_3=1$.
 - Dataset: $D = \{(\mathbf{x}^i \in R^n, y^i \in R) | i \in [1, m]\}$, where $m=15$ is the number of the samples.

- Linear Model:

$$a = \mathbf{w}\mathbf{x} = \sum_{j=1}^n w_j \cdot x_j^i$$

$$\text{where } \mathbf{w} = [w_1, w_2, w_3], \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

- Object:

$$\operatorname{argmin}_{\mathbf{w}} J(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w}} \frac{1}{2m} \sum_{i=1}^m (a^i - y^i)^2$$

Vector form

$$\frac{\partial J}{\partial \mathbf{w}} = \frac{1}{m} (\mathbf{w}\mathbf{X} - \mathbf{y})\mathbf{X}^T$$

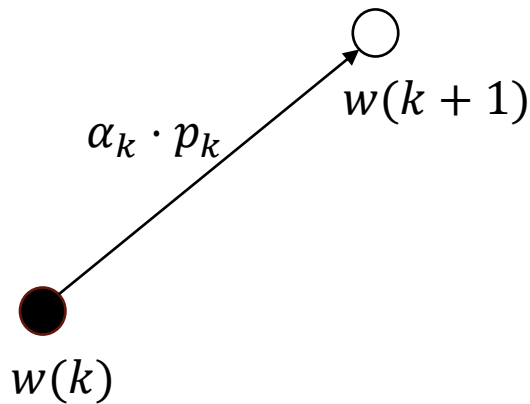
Component form

$$\frac{\partial J}{\partial w_j} = \frac{1}{m} \sum_{i=1}^m (a^i - y^i) \cdot x_j^i$$

Steepest Gradient Descend Method

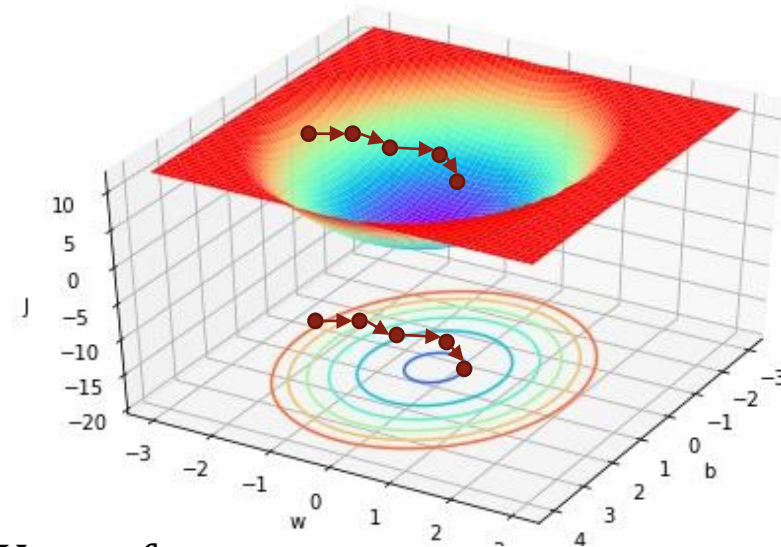
Finding a minimum point step by step

$$w(k+1) = w(k) + \alpha_k \cdot p_k^w$$



p_k , is called searching direction

α_k is learning rate at step k .



Vector form

$$\mathbf{w}(k+1) = \mathbf{w}(k) - \alpha_k \frac{1}{m} (\mathbf{w}X - \mathbf{y})X^T$$

Component form

$$w_j(k+1) = w_j(k) - \alpha_k \frac{1}{m} \sum_{i=1}^m (a^i - y^i) \cdot x_j^i$$

Steepest Gradient Descend Method

Multivariate Linear Regression

- Data:
 - Sample: $(\mathbf{x} \in R^n, y \in R)$, where $n=3$ is the *dimension* of the input vector \mathbf{x} and $x_3=1$.
 - Dataset: $D = \{(\mathbf{x}^i \in R^n, y^i \in R) | i \in [1, m]\}$, where $m=15$ is the number of the samples.

- Linear Model:

$$a = \mathbf{w}\mathbf{x} = \sum_{j=1}^n w_j \cdot x_j^i$$

$$\text{where } \mathbf{w} = [w_1, w_2, w_3], \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

- Object:

$$\operatorname{argmin}_{\mathbf{w}} J(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w}} \frac{1}{2m} \sum_{i=1}^m (a^i - y^i)^2$$

Steepest Descent Algorithm

Input: D, w

for k in $1, 2, \dots, K$:

{

for i in $1, 2, \dots, m$:

{

$$a^i \leftarrow \sum_{j=1}^n w_j x_j^i$$

}

for j in $1, 2, \dots, n$:

{

$$\frac{\partial J}{\partial w_j} \leftarrow \frac{1}{m} \sum_{i=1}^m (a^i - y^i) \cdot x_j^i$$

}

for j in $1, 2, \dots, n$:

{

$$w_j \leftarrow w_j - \alpha \frac{\partial J}{\partial w_j}$$

}

}

Steepest Gradient Descend Method

Multivariate Linear Regression

- Data:
 - Sample: $(\mathbf{x} \in R^n, y \in R)$, where $n=3$ is the *dimension* of the input vector \mathbf{x} and $x_3=1$.
 - Dataset: $D = \{(\mathbf{x}^i \in R^n, y^i \in R) | i \in [1, m]\}$, where $m=15$ is the number of the samples.

- Linear Model:

$$a = \mathbf{w}\mathbf{x} = \sum_{j=1}^n w_j \cdot x_j^i$$

$$\text{where } \mathbf{w} = [w_1, w_2, w_3], \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

- Object:

$$\operatorname{argmin}_{\mathbf{w}} J(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w}} \frac{1}{2m} \sum_{i=1}^m (a^i - y^i)^2$$

Steepest Descent Algorithm

Input: D, \mathbf{w}

for k in $1, 2, \dots, K$:

{

$$\mathbf{a} \leftarrow \mathbf{w}\mathbf{X}$$

$$\frac{\partial J}{\partial \mathbf{w}} = \frac{1}{m} (\mathbf{w}\mathbf{X} - \mathbf{y})\mathbf{X}^T$$

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial J}{\partial \mathbf{w}}$$

}

$$\begin{aligned} J &= \frac{1}{2m} \sum_{i=1}^m (a^i - y^i)^2 \\ &= \frac{1}{2m} (\mathbf{a} - \mathbf{y})(\mathbf{a} - \mathbf{y})^T \end{aligned}$$

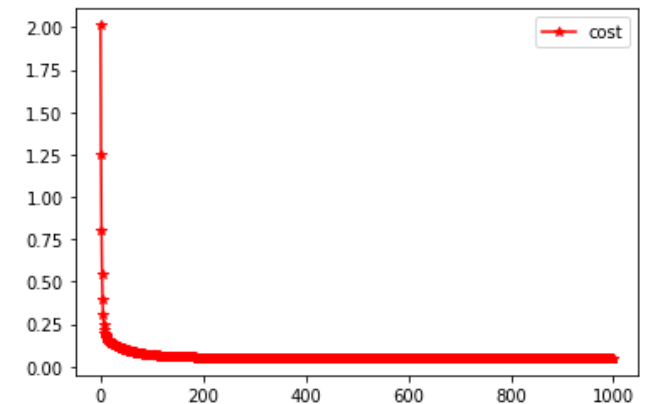
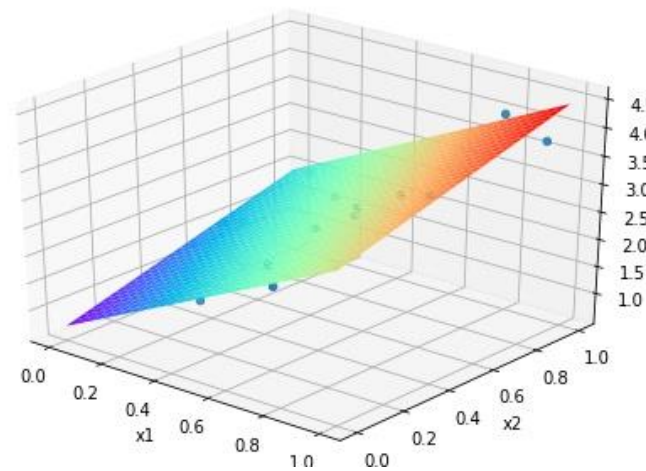
Steepest Gradient Descend Method

Size	0.488	0.681	0.655	0.088	0.139	0.721	0.999	0.775	0.881	0.007	0.14	0.817	0.681	0.851	0.436
Color	0.609	0.112	0.324	0.669	0.91	0.41	0.902	0.353	0.865	0.539	0.524	0.62	0.093	0.447	0.088
bias	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
prize	2.737	2.237	2.869	1.626	2.22	3.038	3.929	3.351	4.308	1.304	1.917	3.189	2.656	3.517	1.663

Please predict the price for an apple (size=0.6, color=0.3)

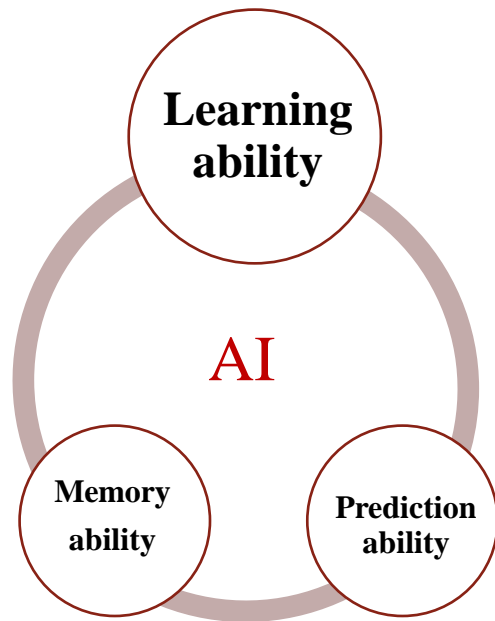
$$2.47 * 0.6 + 1.25 * 0.3 + 0.69 = 2.56$$

$$w = [2.47, 1.25, 0.69]$$



machine learning

- Supervised Learning
- Unsupervised Learning
- Reinforcement Learning



Supervised Learning

- The dataset contains the true answer(label):

$$D = \{ (x, y) \}$$

- Regression
- Classification



Unsupervised Learning



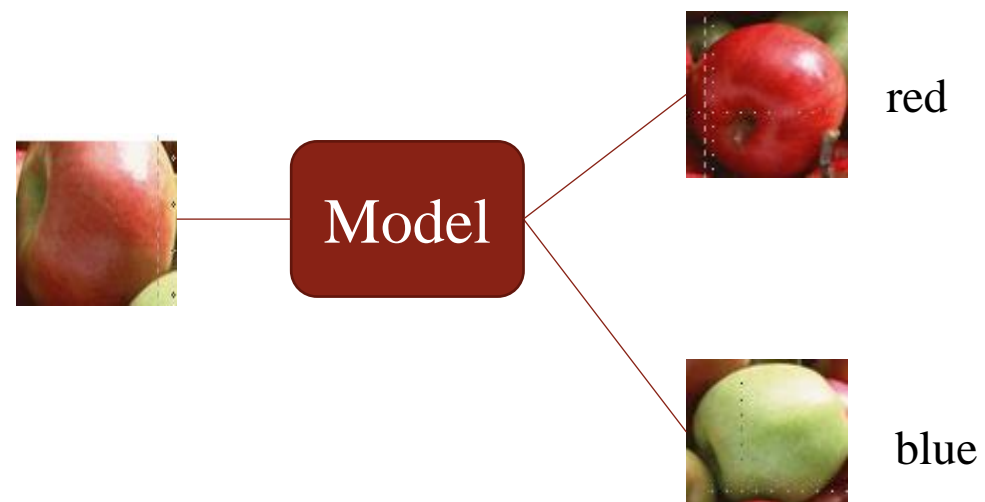
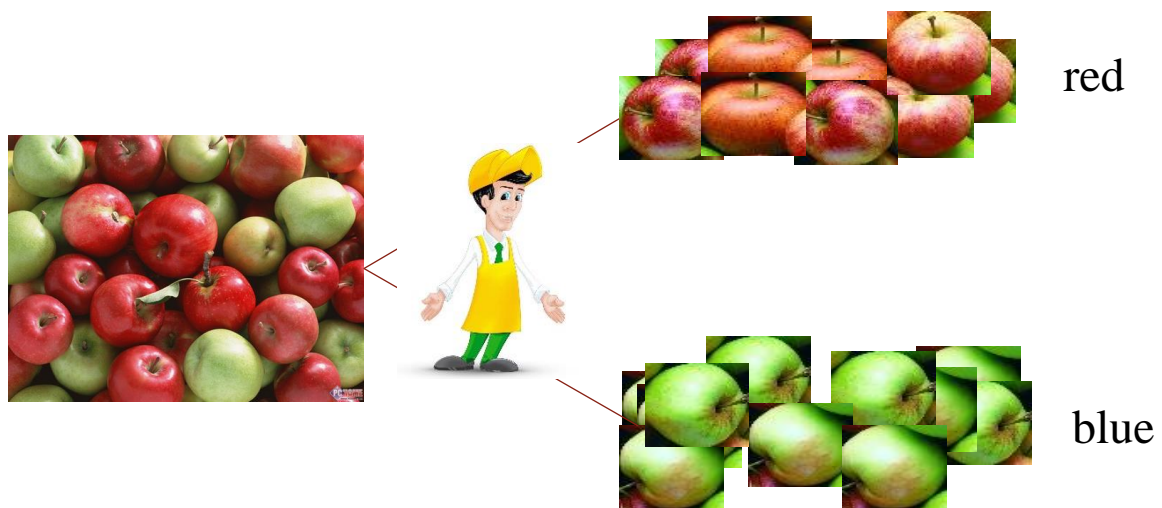
Reinforcement Learning

Outline

- Brief review
- **Binary Classification**
- Multi-class classification

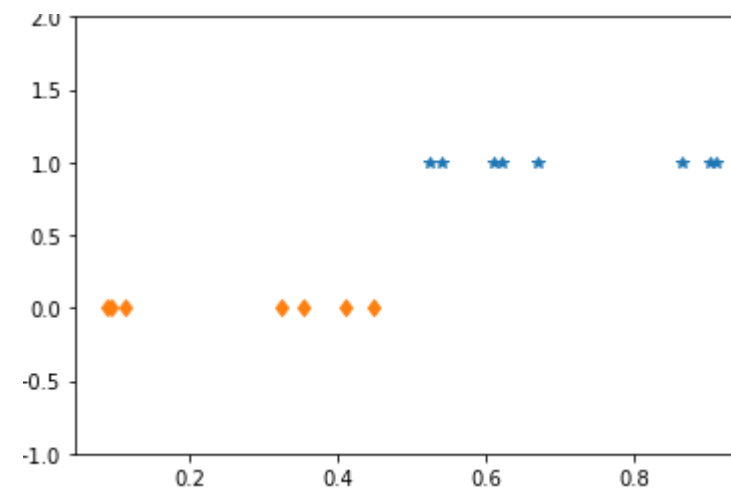
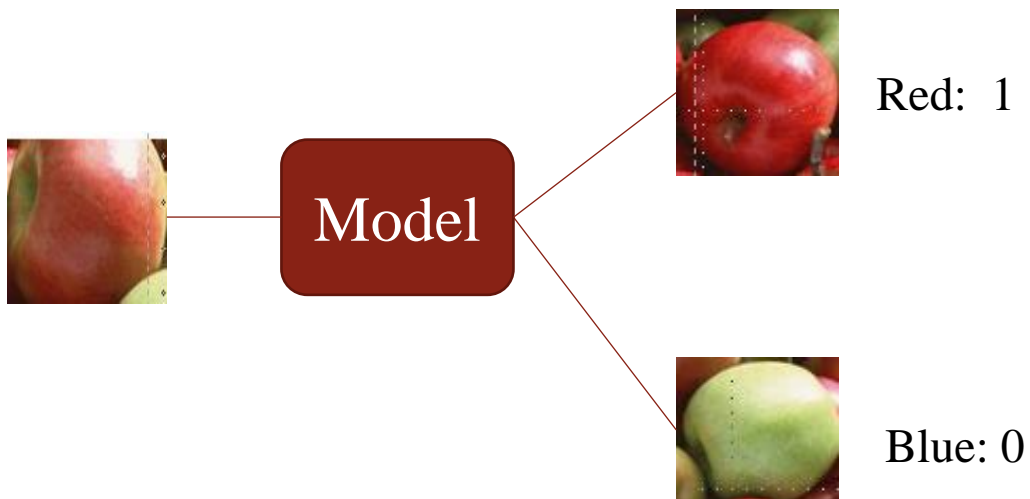
Examples

Color x	0.609	0.112	0.324	0.669	0.91	0.41	0.902	0.353	0.865	0.539	0.524	0.62	0.093	0.447	0.088
class	red	blue	blue	red	red	blue	red	blue	red	red	red	red	blue	blue	blue



Examples

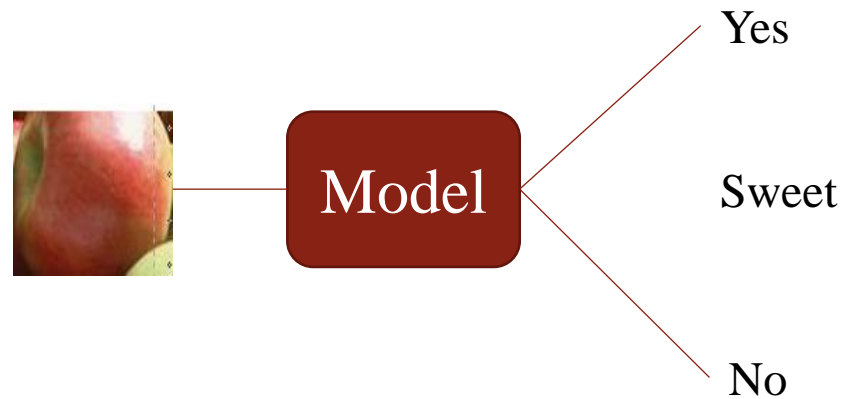
Color x	0.609	0.112	0.324	0.669	0.91	0.41	0.902	0.353	0.865	0.539	0.524	0.62	0.093	0.447	0.088
class	red	blue	blue	red	red	blue	red	blue	red	red	red	red	blue	blue	blue
label y	1	0	0	1	1	0	1	0	1	1	1	1	0	0	0



- Sample: $(\mathbf{x} \in \mathbb{R}^n, y \in \{0,1\})$
- Dataset: $D = \{(\mathbf{x}^i \in \mathbb{R}^n, y^i \in \{0,1\}) | i \in [1, m]\}$

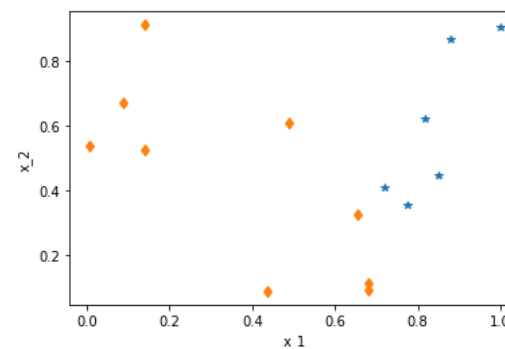
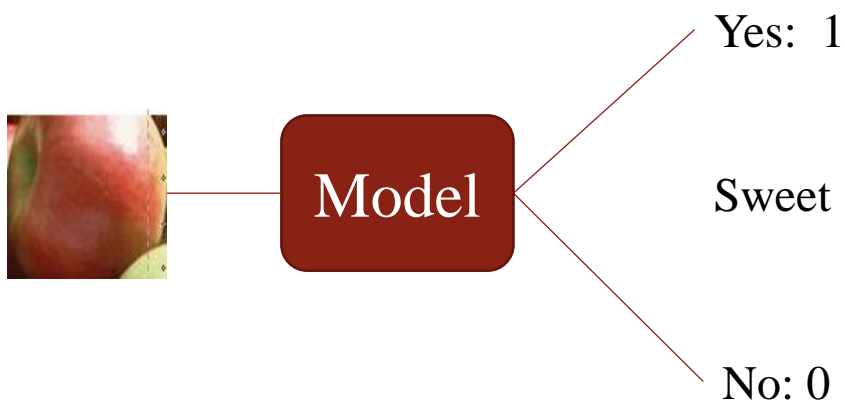
Examples

Size x_1	0.488	0.681	0.655	0.088	0.139	0.721	0.999	0.775	0.881	0.007	0.14	0.817	0.681	0.851	0.436
Color x_2	0.609	0.112	0.324	0.669	0.91	0.41	0.902	0.353	0.865	0.539	0.524	0.62	0.093	0.447	0.088
Class sweet	No	No	No	No	No	Yes	Yes	Yes	Yes	No	No	Yes	No	Yes	No



Examples

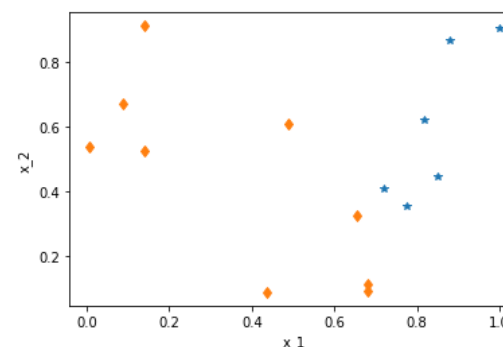
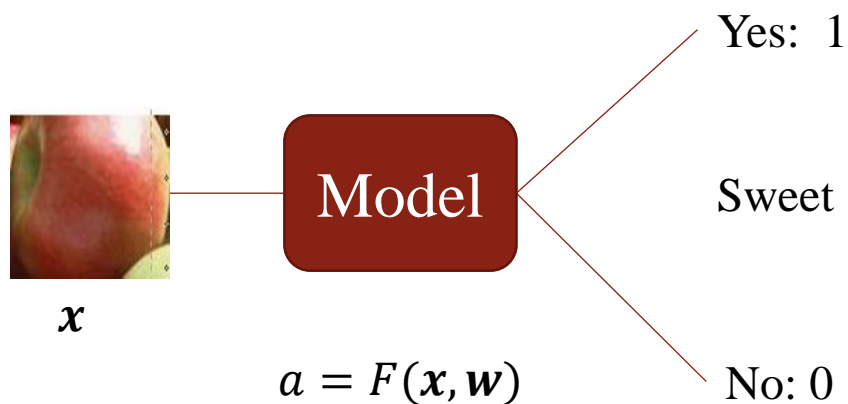
Size x_1	0.488	0.681	0.655	0.088	0.139	0.721	0.999	0.775	0.881	0.007	0.14	0.817	0.681	0.851	0.436
Color x_2	0.609	0.112	0.324	0.669	0.91	0.41	0.902	0.353	0.865	0.539	0.524	0.62	0.093	0.447	0.088
Class sweet	No	No	No	No	No	Yes	Yes	Yes	Yes	No	No	Yes	No	Yes	No
Label y	0	0	0	0	0	1	1	1	1	0	0	1	0	1	0



- Sample: $(\mathbf{x} \in R^n, y \in \{0,1\})$
- Dataset: $D = \{(\mathbf{x}^i \in R^n, y^i \in \{0,1\}) | i \in [1, m]\}$

Examples

Size x_1	0.488	0.681	0.655	0.088	0.139	0.721	0.999	0.775	0.881	0.007	0.14	0.817	0.681	0.851	0.436
Color x_2	0.609	0.112	0.324	0.669	0.91	0.41	0.902	0.353	0.865	0.539	0.524	0.62	0.093	0.447	0.088
Class sweet	No	No	No	No	No	Yes	Yes	Yes	Yes	No	No	Yes	No	Yes	No
Label y	0	0	0	0	0	1	1	1	1	0	0	1	0	1	0



- Sample: $(\mathbf{x} \in R^n, y \in \{0,1\})$
- Dataset: $D = \{(\mathbf{x}^i \in R^n, y^i \in \{0,1\}) | i \in [1, m]\}$

Binary classification

- Data:

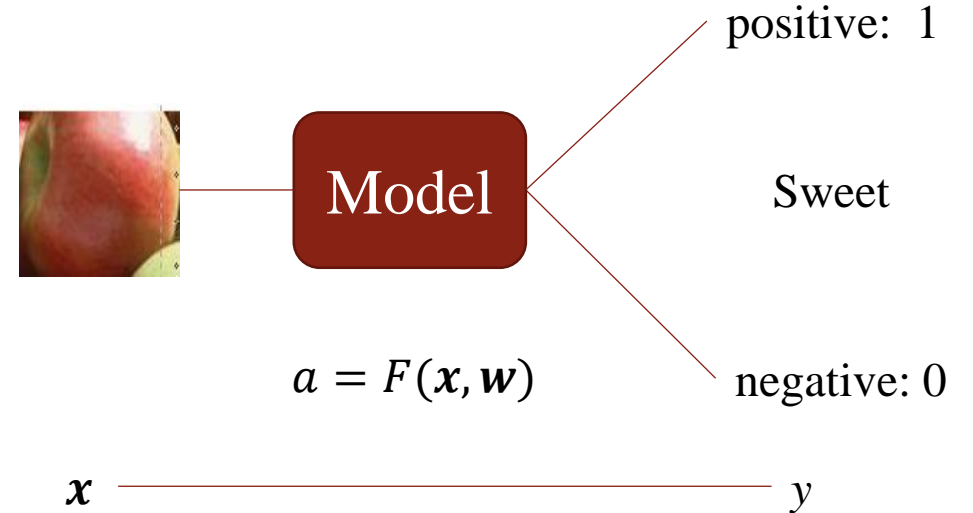
- Sample: $(\mathbf{x} \in R^n, y \in \{0,1\})$, where n is the *dimension* of the input vector \mathbf{x} .
- Dataset: $D = \{(\mathbf{x}^i \in R^n, y^i \in \{0,1\}) | i \in [1, m]\}$, where m is the number of the samples.

- Model:

$$a = F(\mathbf{x}, \mathbf{w})$$

- Object:

$\forall (\mathbf{x}, y) \in D$, the output a is close to y



How to define $F(\mathbf{x}, \mathbf{w})$

Binary classification

- Data:

- Sample: $(\mathbf{x} \in R^n, y \in \{0,1\})$, where n is the *dimension* of the input vector \mathbf{x} .
- Dataset: $D = \{(\mathbf{x}^i \in R^n, y^i \in \{0,1\}) | i \in [1, m]\}$, where m is the number of the samples.

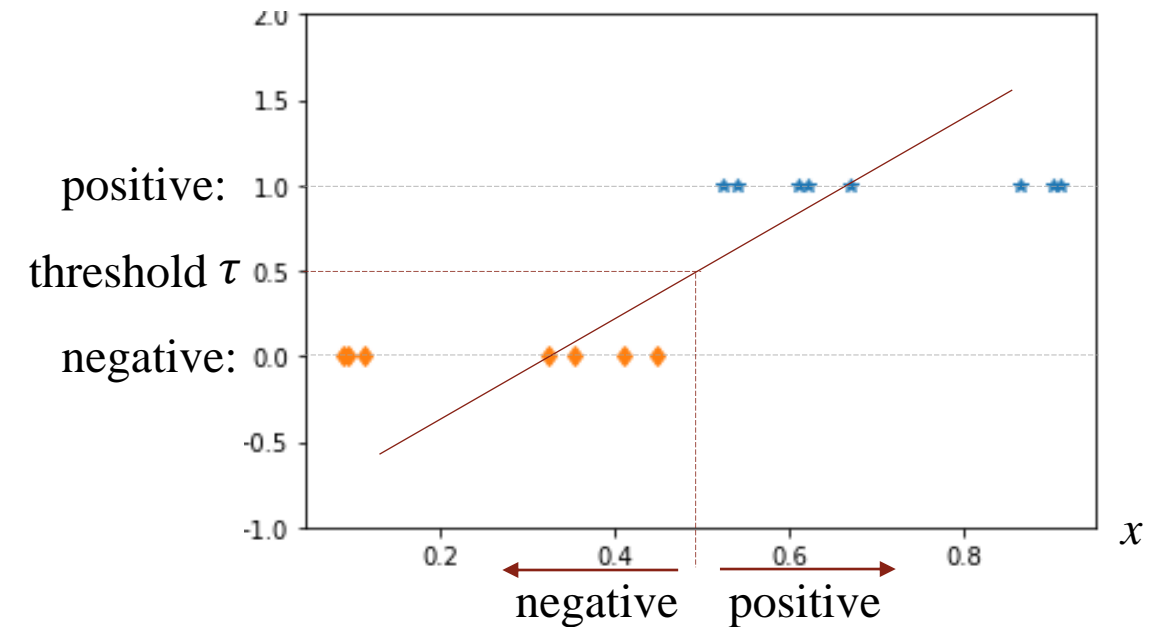
- Linear Regression Model :

$$a = \mathbf{w}\mathbf{x} = \sum_{j=1}^n w_j \cdot x_j^i$$

- Object:

$\forall (\mathbf{x}, y) \in D$, the output a is close to y

Color x	0.609	0.112	0.324	0.669
label y	1	0	0	1



Threshold classifier output a at τ :

if $a \geq \tau$, predict “y=1”

if $a < \tau$, predict “y=0”.

For example $\tau = 0.5$

Binary classification

- Data:

- Sample: $(\mathbf{x} \in \mathbb{R}^n, y \in \{0,1\})$, where n is the *dimension* of the input vector \mathbf{x} .
- Dataset: $D = \{(\mathbf{x}^i \in \mathbb{R}^n, y^i \in \{0,1\}) | i \in [1, m]\}$, where m is the number of the samples.

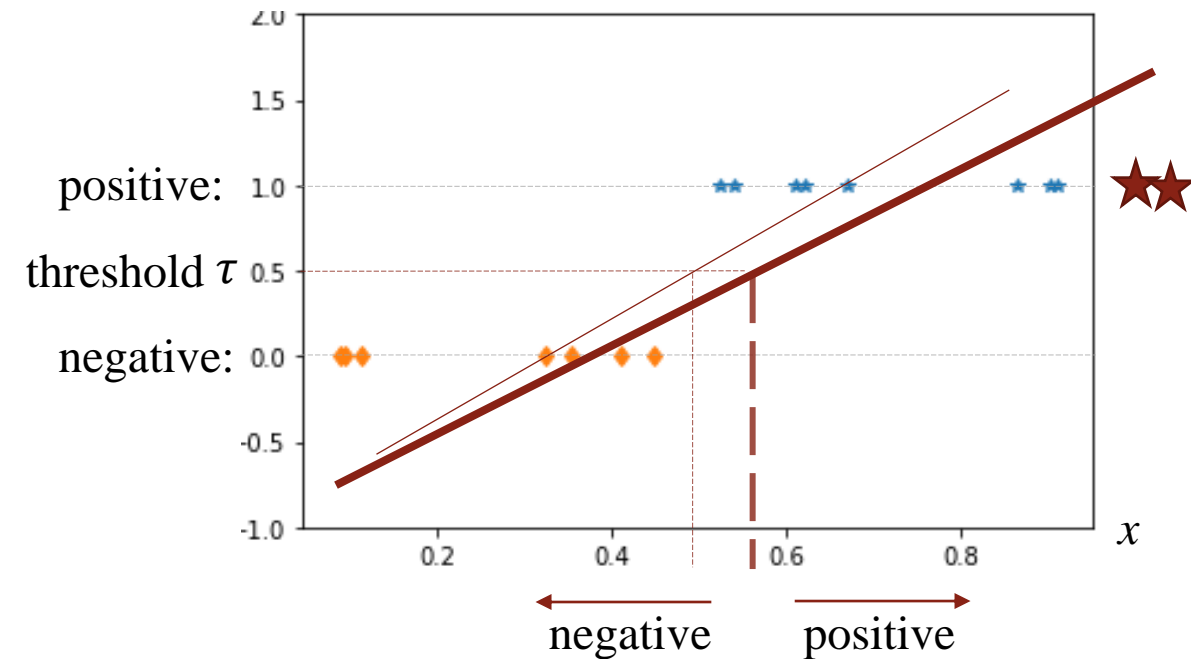
- Linear Regression Model:

$$a = \mathbf{w}\mathbf{x} = \sum_{j=1}^n w_j \cdot x_j^i$$

- Object:

$\forall (\mathbf{x}, y) \in D$, the output a is close to y

Color x	0.609	0.112	0.324	0.669
label y	1	0	0	1



label $y \in \{0,1\}$

output $a \in [-\infty, +\infty]$

Binary classification

- Data:

- Sample: $(\mathbf{x} \in R^n, y \in \{0,1\})$, where n is the *dimension* of the input vector \mathbf{x} .
- Dataset: $D = \{(\mathbf{x}^i \in R^n, y^i \in \{0,1\}) | i \in [1, m]\}$, where m is the number of the samples.

- Logistic Regression Model:

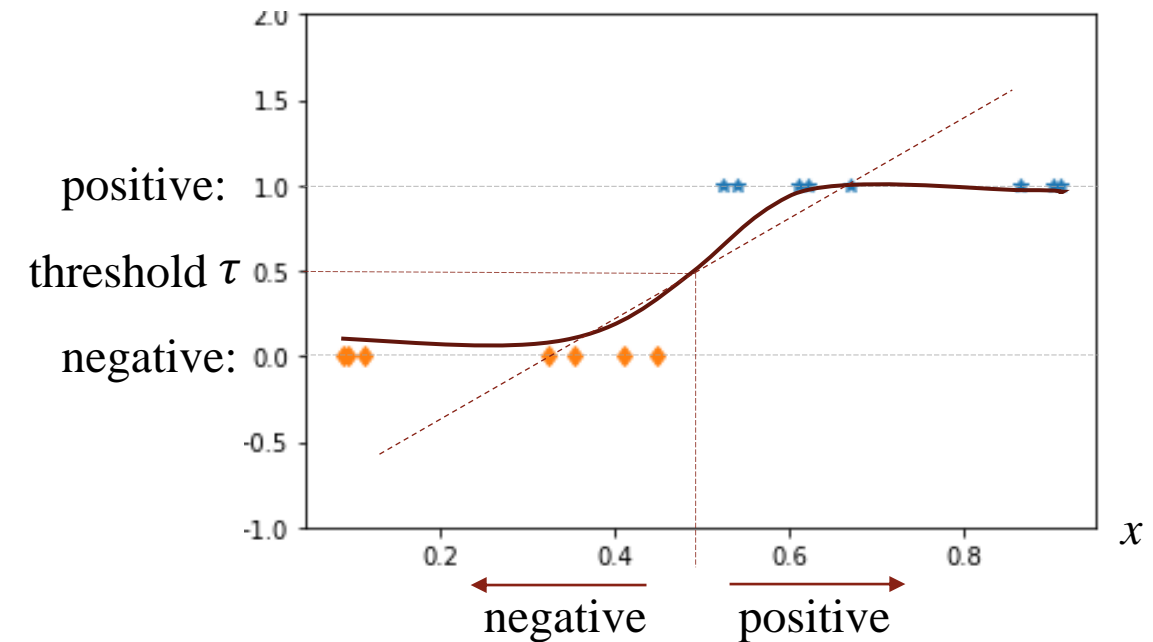
$$a = \sigma(\mathbf{w}\mathbf{x}) = \sigma\left(\sum_{j=1}^n w_j \cdot x_j^i\right)$$

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

- Object:

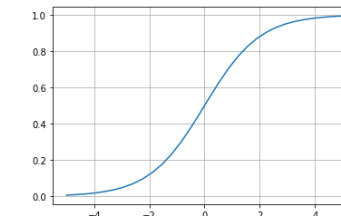
$\forall (\mathbf{x}, y) \in D$, the output a is close to y

Color x	0.609	0.112	0.324	0.669
label y	1	0	0	1



label $y \in \{0,1\}$

output $a \in (0,1)$



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

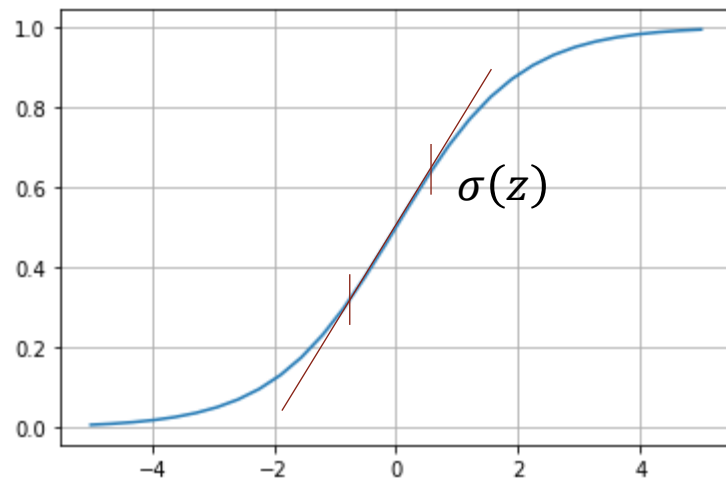
Sigmoid function

Sigmoid Function

A sigmoid function is a bounded, monotonic, differentiable function that is defined for all real input values and has a non-negative derivative at each point.

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

If $z_1 > z_2$ then $\sigma(z_1) > \sigma(z_2)$

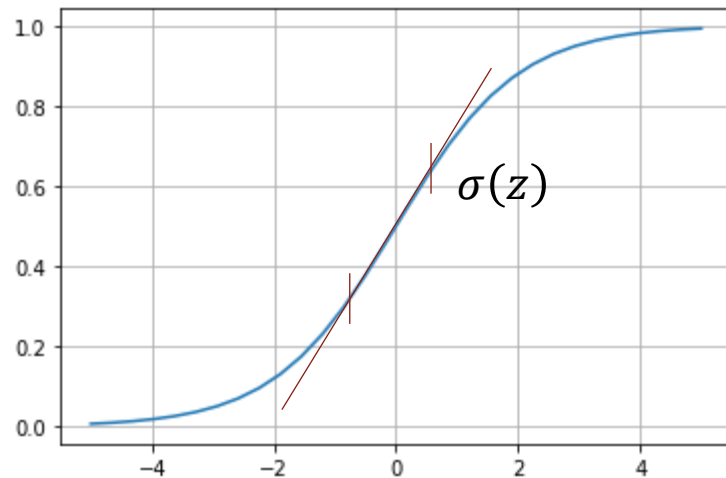


Sigmoid Function

A sigmoid function is a bounded, monotonic, differentiable function that is defined for all real input values and has a non-negative derivative at each point.

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

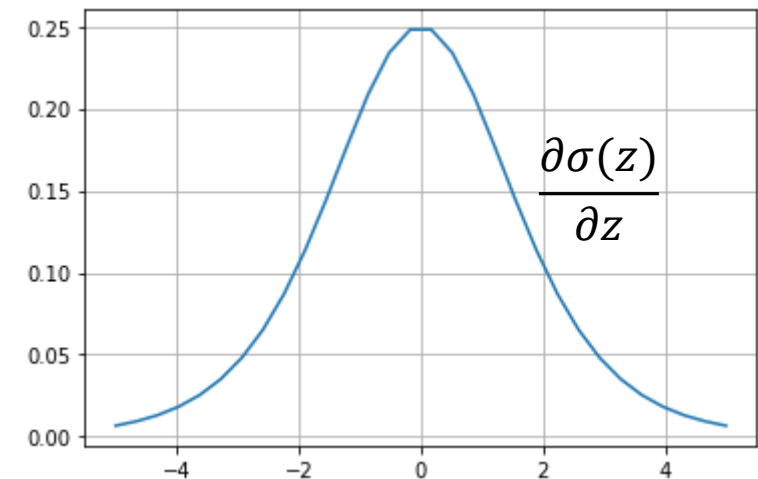
If $z_1 > z_2$ then $\sigma(z_1) > \sigma(z_2)$



$$\frac{\partial \sigma(z)}{\partial z} = -\frac{-e^{-z}}{(1 + e^{-z})^2}$$

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

$$\frac{\partial \sigma(z)}{\partial z} = \sigma(z) \cdot (1 - \sigma(z))$$



Interpretation of output

- Data:

- Sample: $(\mathbf{x} \in R^n, y \in \{0,1\})$, where n is the *dimension* of the input vector \mathbf{x} .
- Dataset: $D = \{(\mathbf{x}^i \in R^n, y^i \in \{0,1\}) | i \in [1, m]\}$, where m is the number of the samples.

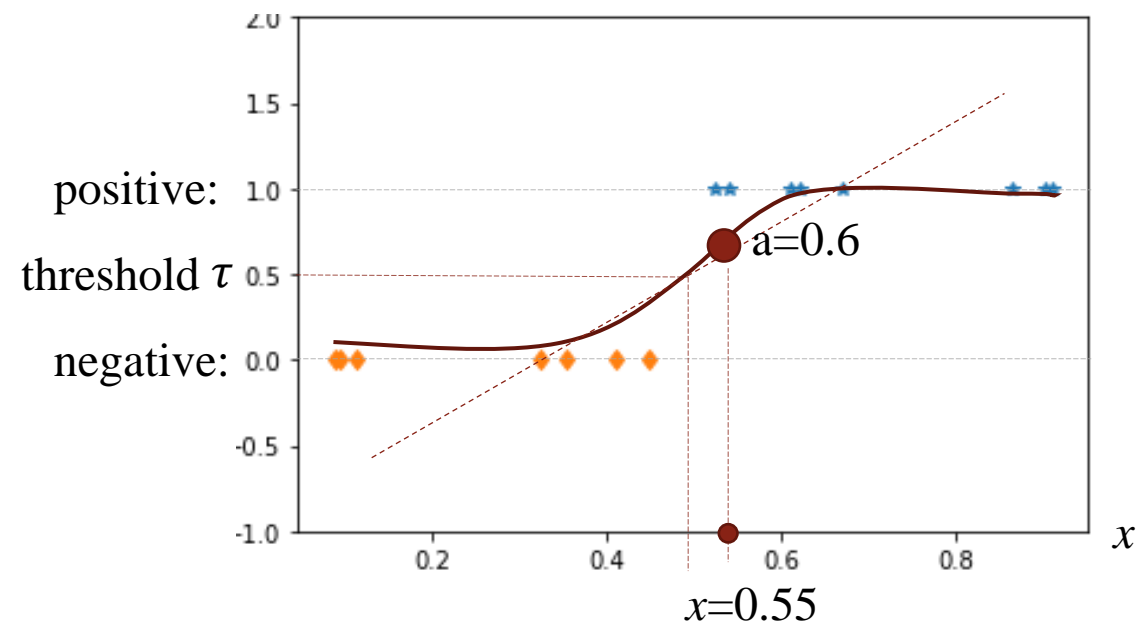
- Logistic Regression Model:

$$a = \sigma(\mathbf{w}\mathbf{x}) = \sigma\left(\sum_{j=1}^n w_j \cdot x_j^i\right)$$

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

- Object:

$\forall (\mathbf{x}, y) \in D$, the output a is close to y



Given an input x ,
the output a estimates the probability that “ $y=1$ ”

$$p(y = 1 | \mathbf{x}, \mathbf{w}) = a$$

$$p(y = 0 | \mathbf{x}, \mathbf{w}) = 1 - a$$

Interpretation of output

- Data:

- Sample: $(\mathbf{x} \in R^n, y \in \{0,1\})$, where n is the *dimension* of the input vector \mathbf{x} .
- Dataset: $D = \{(\mathbf{x}^i \in R^n, y^i \in \{0,1\}) | i \in [1, m]\}$, where m is the number of the samples.

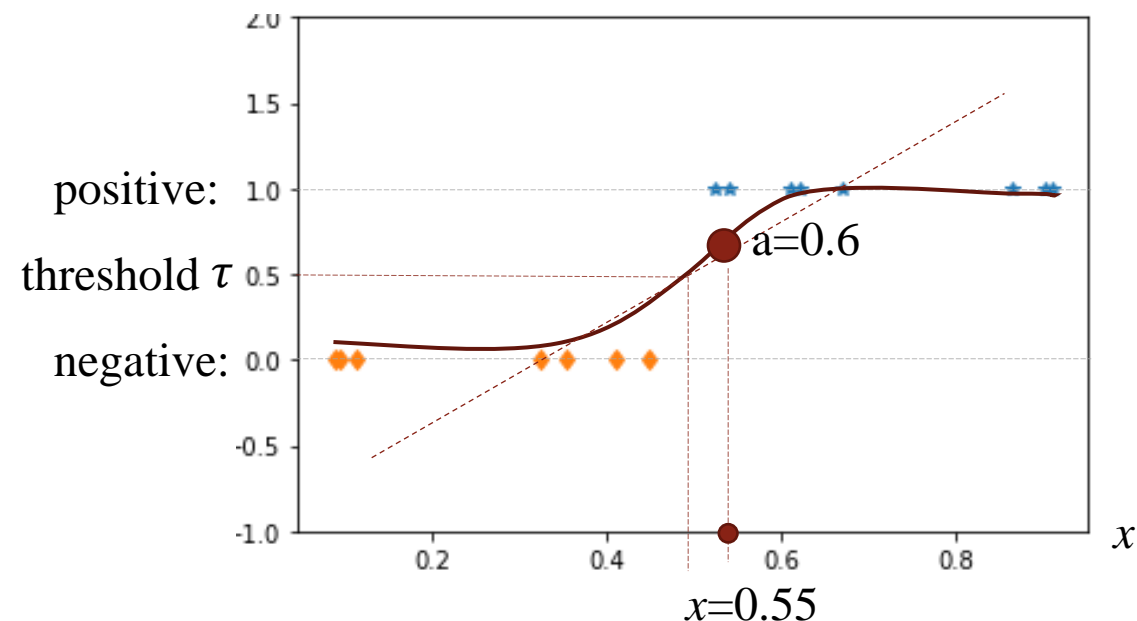
- Logistic Regression Model:

$$a = \sigma(\mathbf{w}\mathbf{x}) = \sigma\left(\sum_{j=1}^n w_j \cdot x_j^i\right)$$

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

- Object:

$\forall (\mathbf{x}, y) \in D$, the output a is close to y



$$p(y = 1 | \mathbf{x}, \mathbf{w}) = a$$

$$p(y = 0 | \mathbf{x}, \mathbf{w}) = 1 - a$$

How to find \mathbf{w} ?



$$p(y | \mathbf{x}, \mathbf{w}) = a^y \cdot (1 - a)^{1-y}$$

Cost function

Joint probability for sampled dataset D

$$P(Y|X, \mathbf{w}) = \prod_{i=1}^m p(y^i | \mathbf{x}^i, \mathbf{w})$$

$$\mathbf{X} = \begin{bmatrix} x_1^1, x_1^2, \dots, x_1^i, \dots, x_1^m \\ x_2^1, x_2^2, \dots, x_2^i, \dots, x_2^m \\ x_3^1, x_3^2, \dots, x_3^i, \dots, x_3^m \end{bmatrix} \quad Y = [y^1, \dots, y^i, \dots, y^m]$$

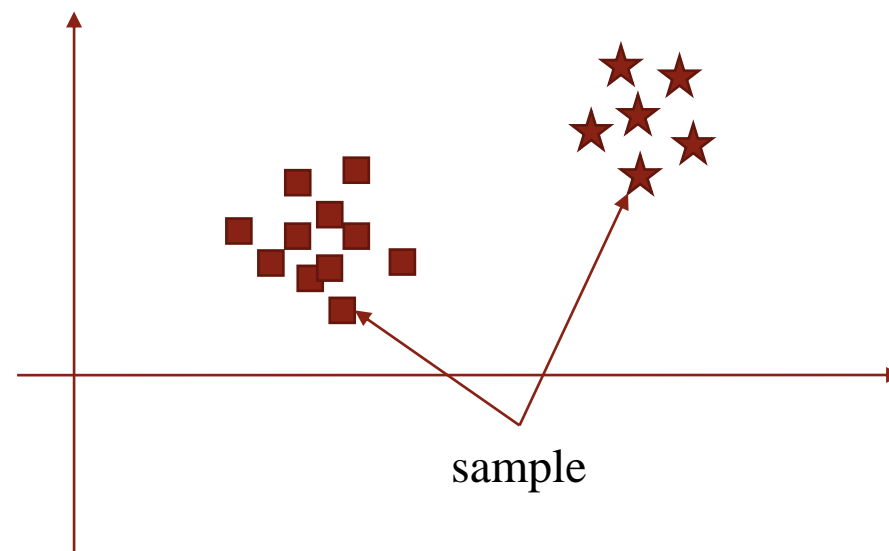
we would like to maximize the joint probability $P(Y|X, \mathbf{w})$

↓

$$\underset{\mathbf{w}}{\operatorname{argmax}} P(Y|X, \mathbf{w})$$

$$p(y|\mathbf{x}, \mathbf{w}) = a^y \cdot (1 - a)^{1-y}$$

Size \mathbf{x}_1	0.088	0.999	0.007	0.14	0.817	0.681	0.851	0.436
Color \mathbf{x}_2	0.669	0.902	0.539	0.524	0.62	0.093	0.447	0.088
Label \mathbf{y}	0	1	0	0	1	0	1	0

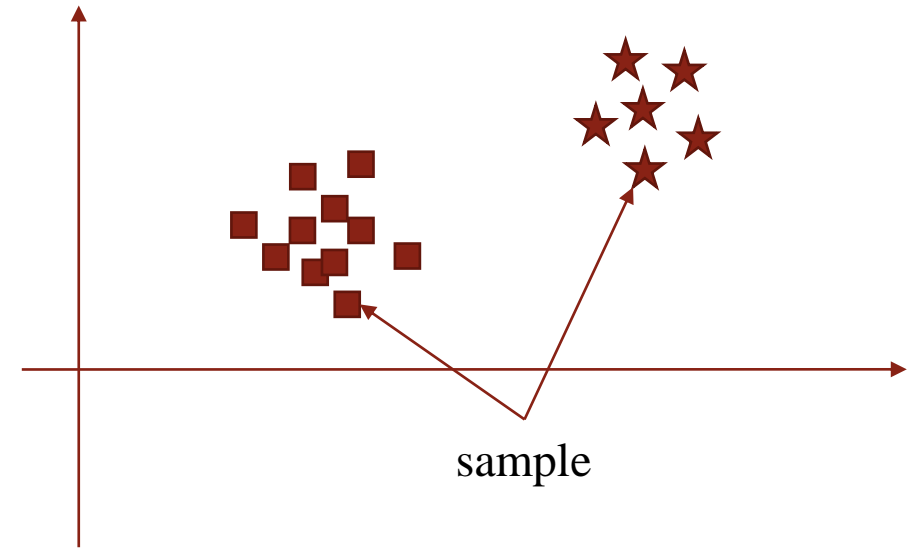


Cost function

$$\operatorname{argmax}_{\mathbf{w}} P(Y|X, \mathbf{w}) = \operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^m p(y^i | \mathbf{x}^i, \mathbf{w})$$



$$\begin{aligned} \operatorname{argmax}_{\mathbf{w}} \log P(Y|X, \mathbf{w}) &= \operatorname{argmax}_{\mathbf{w}} \log \prod_{i=1}^m p(y^i | \mathbf{x}^i, \mathbf{w}) \\ &= \operatorname{argmax}_{\mathbf{w}} \sum_{i=1}^m \log(p(y^i | \mathbf{x}^i, \mathbf{w})) \\ &= \operatorname{argmax}_{\mathbf{w}} \sum_{i=1}^m \log(a^{y^i} \cdot (1 - a^i)^{1-y^i}) \\ &= \operatorname{argmax}_{\mathbf{w}} \sum_{i=1}^m [y^i \cdot \log(a^i) + (1 - y^i) \cdot \log(1 - a^i)] \end{aligned}$$



$$p(y|x, \mathbf{w}) = a^y \cdot (1 - a)^{1-y}$$

Cost function

$$\operatorname{argmax}_{\mathbf{w}} P(Y|\mathbf{X}, \mathbf{w}) = \operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^m p(y^i | \mathbf{x}^i, \mathbf{w})$$



$$\operatorname{argmax}_{\mathbf{w}} \log P(Y|\mathbf{X}, \mathbf{w}) = \operatorname{argmax}_{\mathbf{w}} \sum_{i=1}^m [y^i \cdot \log(a^i) + (1 - y^i) \cdot \log(1 - a^i)]$$

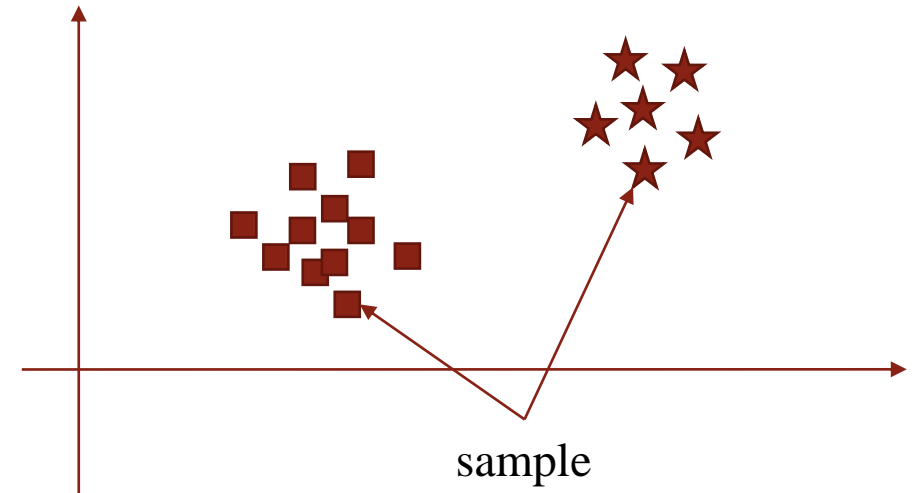


$$\operatorname{argmax}_{\mathbf{w}} \frac{1}{m} \log P(Y|\mathbf{X}, \mathbf{w}) = \operatorname{argmax}_{\mathbf{w}} \frac{1}{m} \sum_{i=1}^m [y^i \cdot \log(a^i) + (1 - y^i) \cdot \log(1 - a^i)]$$



$$\operatorname{argmin}_{\mathbf{w}} -\frac{1}{m} \log P(Y|\mathbf{X}, \mathbf{w}) = \operatorname{argmin}_{\mathbf{w}} -\frac{1}{m} \sum_{i=1}^m [y^i \cdot \log(a^i) + (1 - y^i) \cdot \log(1 - a^i)]$$

$$J(\mathbf{w}) = -\frac{1}{m} \sum_{i=1}^m [y^i \cdot \log(a^i) + (1 - y^i) \cdot \log(1 - a^i)]$$



Cost function

$$\operatorname{argmax}_{\mathbf{w}} P(Y|\mathbf{X}, \mathbf{w}) = \operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^m p(y^i | \mathbf{x}^i, \mathbf{w})$$

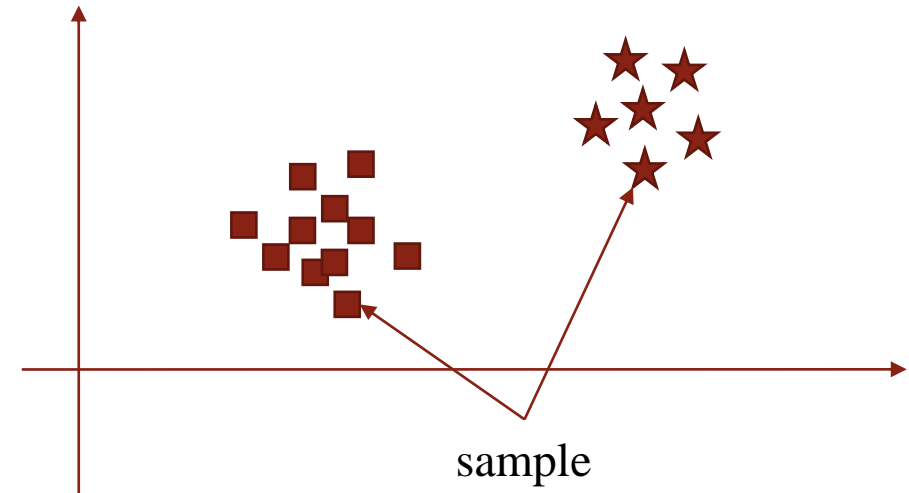
$$\operatorname{argmax}_{\mathbf{w}} \log P(Y|\mathbf{X}, \mathbf{w}) = \operatorname{argmax}_{\mathbf{w}} \sum_{i=1}^m [y^i \cdot \log(a^i) + (1 - y^i) \cdot \log(1 - a^i)]$$

$$\operatorname{argmax}_{\mathbf{w}} \frac{1}{m} \log P(Y|\mathbf{X}, \mathbf{w}) = \operatorname{argmax}_{\mathbf{w}} \frac{1}{m} \sum_{i=1}^m [y^i \cdot \log(a^i) + (1 - y^i) \cdot \log(1 - a^i)]$$

$$\operatorname{argmin}_{\mathbf{w}} -\frac{1}{m} \log P(Y|\mathbf{X}, \mathbf{w}) = \operatorname{argmin}_{\mathbf{w}} -\frac{1}{m} \sum_{i=1}^m [y^i \cdot \log(a^i) + (1 - y^i) \cdot \log(1 - a^i)]$$

$$J(\mathbf{w}) = -\frac{1}{m} \sum_{i=1}^m [y^i \cdot \log(a^i) + (1 - y^i) \cdot \log(1 - a^i)]$$

Cross-entropy cost function



Binary classification

- Data:

- Sample: $(\mathbf{x} \in R^n, y \in \{0,1\})$, where n is the *dimension* of the input vector \mathbf{x} .
- Dataset: $D = \{(\mathbf{x}^i \in R^n, y^i \in \{0,1\}) | i \in [1, m]\}$, where m is the number of the samples.

- Logistic Regression Model:

$$a = \sigma(\mathbf{w}\mathbf{x}) = \sigma\left(\sum_{j=1}^n w_j \cdot x_j^i\right)$$

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

- Object:

$$\operatorname{argmin}_{\mathbf{w}} J(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w}} \frac{1}{m} \sum_{i=1}^m -[y^i \cdot \log(a^i) + (1 - y^i) \cdot \log(1 - a^i)]$$

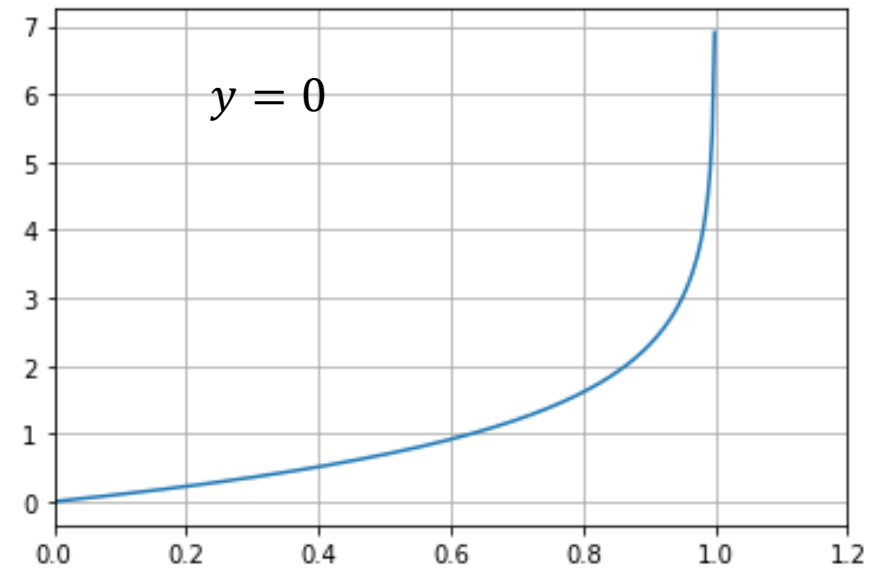
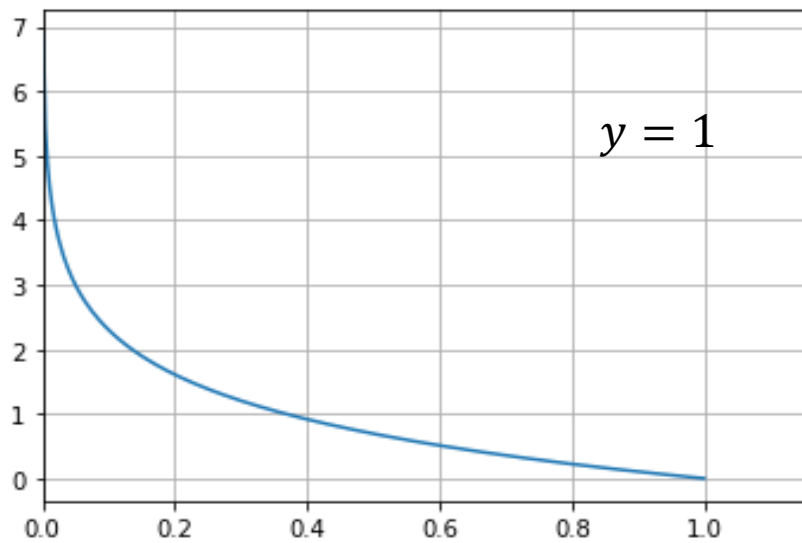
$e(a^i, y^i)$

$$e(a, y) = \begin{cases} -\log(a), & y = 1 \\ -\log(1 - a) & y = 0 \end{cases}$$

$$\begin{cases} e(a, y) \rightarrow \infty, & \text{if } a \text{ is not close to } y \\ e(a, y) \rightarrow 0, & \text{if } a \text{ is close to } y \end{cases}$$

Binary classification

$$e(a, y) = \begin{cases} -\log(a), & y = 1 \\ -\log(1 - a) & y = 0 \end{cases} \quad \begin{cases} e(a, y) \rightarrow \infty, & \text{if } a \text{ is not close to } y \\ e(a, y) \rightarrow 0, & \text{if } a \text{ is close to } y \end{cases}$$



Steepest Gradient Descend Method

- Data:

- Sample: $(\mathbf{x} \in \mathbb{R}^n, y \in \{0,1\})$, where n is the *dimension* of the input vector \mathbf{x} .
- Dataset: $D = \{(\mathbf{x}^i \in \mathbb{R}^n, y^i \in \{0,1\}) \mid i \in [1, m]\}$, where m is the number of the samples.

- Logistic Regression Model:

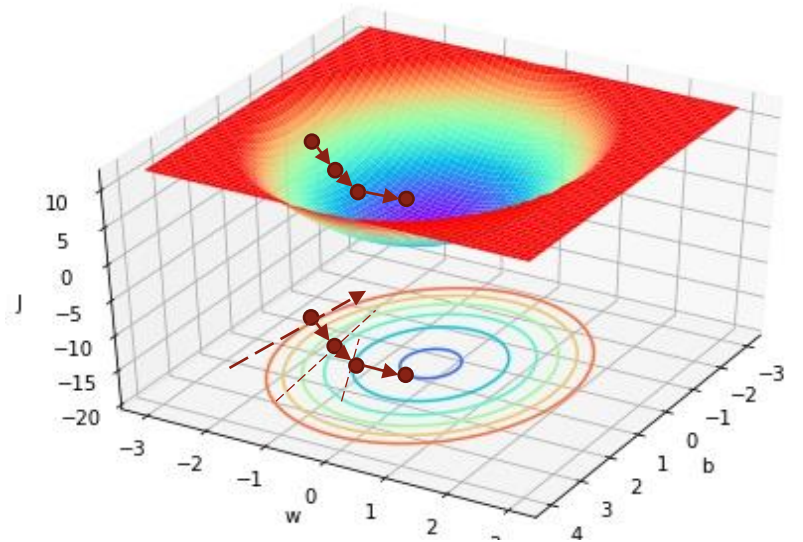
$$a = \sigma(\mathbf{w}\mathbf{x}) = \sigma\left(\sum_{j=1}^n w_j \cdot x_j^i\right)$$

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

- Object:

$$\operatorname{argmin}_{\mathbf{w}} J(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w}} \frac{1}{m} \sum_{i=1}^m -[y^i \cdot \log(a^i) + (1 - y^i) \cdot \log(1 - a^i)]$$

$e(a^i, y^i)$



Vector form

$$\mathbf{w}(k+1) = \mathbf{w}(k) - \alpha_k \frac{\partial J}{\partial \mathbf{w}}$$

Component form

$$w_j(k+1) = w_j(k) - \alpha_k \frac{\partial J}{\partial w_j}$$

Steepest Gradient Descend Method

$$\frac{\partial J}{\partial w_j} = -\frac{1}{m} \sum_{i=1}^m \frac{\partial [y^i \cdot \log a^i + (1 - y^i) \cdot \log(1 - a^i)]}{\partial a^i} \cdot \frac{\partial a^i}{\partial w_j}$$

$$\frac{\partial J}{\partial w_j} = -\frac{1}{m} \sum_{i=1}^m \left[y^i \cdot \frac{1}{a^i} + (1 - y^i) \cdot \frac{-1}{1 - a^i} \right] \cdot \frac{\partial a^i}{\partial w_j}$$

$$\frac{\partial J}{\partial w_j} = -\frac{1}{m} \sum_{i=1}^m \frac{y^i - a^i}{a^i(1 - a^i)} \cdot \frac{\partial a^i}{\partial w_j}$$

$$a = \sigma \left(\sum_j^n w_j x_j \right)$$

Steepest Gradient Descend Method

$$\frac{\partial J}{\partial w_j} = -\frac{1}{m} \sum_{i=1}^m \frac{y^i - a^i}{a^i(1 - a^i)} \cdot \frac{\partial a^i}{\partial w_j}$$

$$\frac{\partial J}{\partial w_j} = -\frac{1}{m} \sum_{i=1}^m \frac{y^i - a^i}{a^i(1 - a^i)} \cdot \frac{\partial a^i}{\partial \sum_l^n w_l x_l^i} \cdot \frac{\partial \sum_l^n w_l x_l^i}{\partial w_j}$$

$$\frac{\partial J}{\partial w_j} = -\frac{1}{m} \sum_{i=1}^m \frac{y^i - a^i}{a^i(1 - a^i)} \cdot [a^i \cdot (1 - a^i)] \cdot x_j^i$$

$$\frac{\partial J}{\partial w_j} = -\frac{1}{m} \sum_{i=1}^m (y^i - a^i) \cdot x_j^i = \frac{1}{m} \sum_{i=1}^m (a^i - y^i) \cdot x_j^i$$

$$a = \sigma \left(\sum_j^n w_j x_j \right)$$

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

$$\frac{\partial \sigma(z)}{\partial z} = \sigma(z) \cdot (1 - \sigma(z))$$

Steepest Gradient Descend Method

$$\frac{\partial J}{\partial w_j} = \frac{1}{m} \sum_{i=1}^m (a^i - y^i) \cdot x_j^i = \frac{1}{m} [(a^1 - y^1), (a^2 - y^2), \dots, (a^i - y^i) \dots (a^m - y^m)] \begin{bmatrix} x_j^1 \\ x_j^2 \\ \vdots \\ x_j^i \\ \vdots \\ x_j^m \end{bmatrix}$$

$$\left[\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2}, \dots, \frac{\partial J}{\partial w_n} \right] = \frac{1}{m} [(a^1 - y^1), \dots (a^i - y^i), \dots, (a^m - y^m)] \begin{bmatrix} x_1^1 \\ x_1^2 \\ \vdots \\ x_1^i \\ \vdots \\ x_1^m \end{bmatrix} \begin{bmatrix} x_2^1 \\ x_2^2 \\ \vdots \\ x_2^i \\ \vdots \\ x_2^m \end{bmatrix} \dots \begin{bmatrix} x_n^1 \\ x_n^2 \\ \vdots \\ x_n^i \\ \vdots \\ x_n^m \end{bmatrix}$$

$$\frac{\partial J}{\partial \mathbf{w}} = \frac{1}{m} [\mathbf{a} - \mathbf{Y}] \mathbf{X}^T$$

Steepest Gradient Descend Method

- Data:

- Sample: $(\mathbf{x} \in \mathbb{R}^n, y \in \{0,1\})$, where n is the *dimension* of the input vector \mathbf{x} .
- Dataset: $D = \{(\mathbf{x}^i \in \mathbb{R}^n, y^i \in \{0,1\}) \mid i \in [1, m]\}$, where m is the number of the samples.

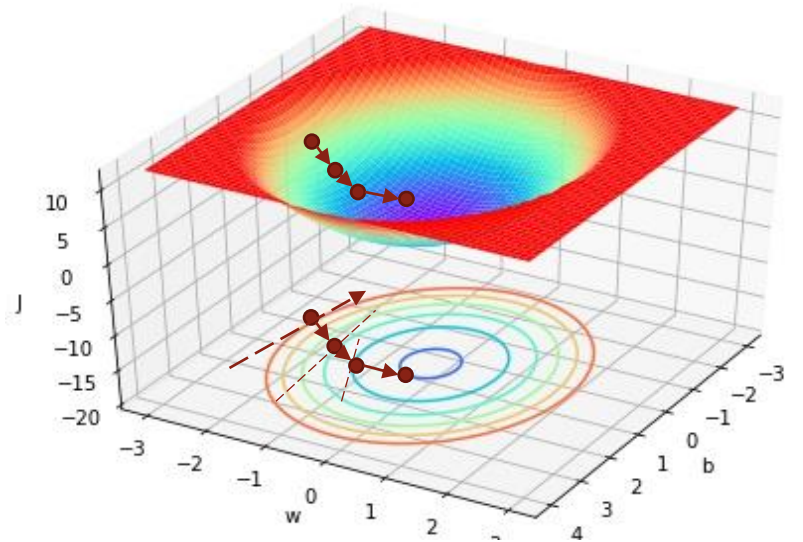
- Logistic Regression Model:

$$a = \sigma(\mathbf{w}\mathbf{x}) = \sigma\left(\sum_{j=1}^n w_j \cdot x_j^i\right)$$

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

- Object:

$$\operatorname{argmin}_{\mathbf{w}} J(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w}} -\frac{1}{m} \sum_{i=1}^m [y^i \cdot \log(a^i) + (1 - y^i) \cdot \log(1 - a^i)]$$



Vector form

$$\mathbf{w}(k+1) = \mathbf{w}(k) - \alpha_k \frac{1}{m} [\mathbf{a} - \mathbf{Y}] \mathbf{X}^T$$

Component form

$$w_j(k+1) = w_j(k) - \alpha_k \frac{1}{m} \sum_{i=1}^m (a^i - y^i) \cdot x_j^i$$

Steepest Gradient Descend Method

- Data:

- Sample: $(\mathbf{x} \in R^n, y \in \{0,1\})$, where n is the *dimension* of the input vector \mathbf{x} .
- Dataset: $D = \{(\mathbf{x}^i \in R^n, y^i \in \{0,1\}) | i \in [1, m]\}$, where m is the number of the samples.

- Logistic Regression Model:

$$a = \sigma(\mathbf{w}\mathbf{x}) = \sigma\left(\sum_{j=1}^n w_j \cdot x_j^i\right)$$

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

- Object:

$$\operatorname{argmin}_{\mathbf{w}} J(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w}} -\frac{1}{m} \sum_{i=1}^m [y^i \cdot \log(a^i) + (1 - y^i) \cdot \log(1 - a^i)]$$

Steepest Descent Algorithm

Input: D, \mathbf{w}, α

for k in $1, 2, \dots, K$:

{

for i in $1, 2, \dots, m$:

{

$$a^i \leftarrow \sigma\left(\sum_{j=1}^n w_j \cdot x_j^i\right)$$

}

for j in $1, 2, \dots, n$:

{

$$\frac{\partial J}{\partial w_j} \leftarrow \frac{1}{m} \sum_{i=1}^m (a^i - y^i) \cdot x_j^i$$

}

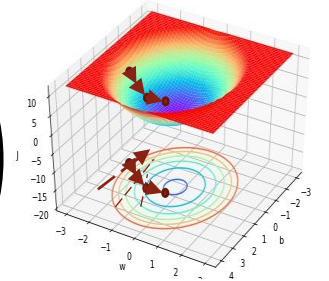
for j in $1, 2, \dots, n$:

{

$$w_j \leftarrow w_j - \alpha \frac{\partial J}{\partial w_j}$$

}

}



Steepest Gradient Descend Method

- Data:

- Sample: $(\mathbf{x} \in R^n, y \in \{0,1\})$, where n is the *dimension* of the input vector \mathbf{x} .
- Dataset: $D = \{(\mathbf{x}^i \in R^n, y^i \in \{0,1\}) | i \in [1, m]\}$, where m is the number of the samples.

- Logistic Regression Model:

$$a = \sigma(\mathbf{w}\mathbf{x}) = \sigma\left(\sum_{j=1}^n w_j \cdot x_j^i\right)$$

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

- Object:

$$\operatorname{argmin}_{\mathbf{w}} J(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w}} -\frac{1}{m} \sum_{i=1}^m [y^i \cdot \log(a^i) + (1 - y^i) \cdot \log(1 - a^i)]$$

Steepest Descent Algorithm

Input: D, \mathbf{w}, α

for k in $1, 2, \dots, K$:

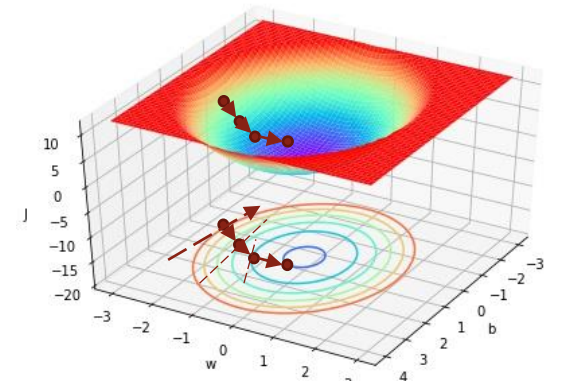
{

$$\mathbf{a} \leftarrow \sigma(\mathbf{w}\mathbf{X})$$

$$\frac{\partial J}{\partial \mathbf{w}} = \frac{1}{m} [\mathbf{a} - \mathbf{Y}] \mathbf{X}^T$$

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial J}{\partial \mathbf{w}}$$

}



Steepest Gradient Descend Method

- Data:

- Dataset: $D = \{(\mathbf{x}^i \in R^n, y^i \in \{0,1\}) | i \in [1, m]\}$, where m is the number of the samples.

- Logistic Regression Model:

$$a = \sigma(\mathbf{w}\mathbf{x}) = \sigma\left(\sum_{j=1}^n w_j \cdot x_j^i\right)$$

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

Vector form

$$\mathbf{w}(\mathbf{k} + 1) = \mathbf{w}(\mathbf{k}) - \alpha_k \frac{1}{m} [\mathbf{a} - Y] X^T$$

Component form

$$w_j(k + 1) = w_j(k) - \alpha_k \frac{1}{m} \sum_{i=1}^m (a^i - y^i) \cdot x_j^i$$

- Data:

- Dataset: $D = \{(\mathbf{x}^i \in R^n, y^i \in R) | i \in [1, m]\}$, where m is the number of the samples.

- Linear Regression Model:

$$a = \mathbf{w}\mathbf{x} = \sum_{j=1}^n w_j \cdot x_j^i$$

Vector form

$$\mathbf{w}(\mathbf{k} + 1) = \mathbf{w}(\mathbf{k}) - \alpha_k \frac{1}{m} [\mathbf{a} - Y] X^T$$

Component form

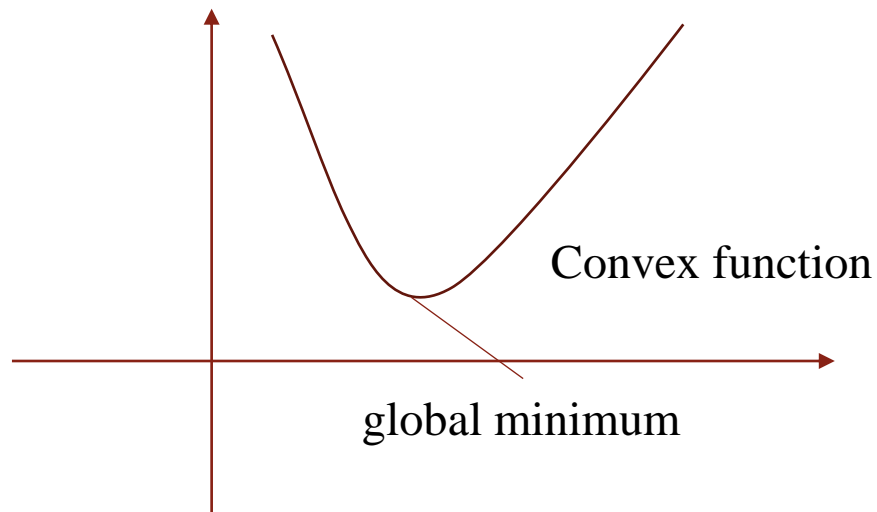
$$w_j(k + 1) = w_j(k) - \alpha_k \frac{1}{m} \sum_{i=1}^m (a^i - y^i) \cdot x_j^i$$

Cross-entropy vs Mean Square Error

$$a = \sigma(\mathbf{w}x) = \sigma\left(\sum_{j=1}^n w_j \cdot x_j^i\right) \quad \text{where } \sigma(z) = \frac{1}{1+e^{-z}}$$

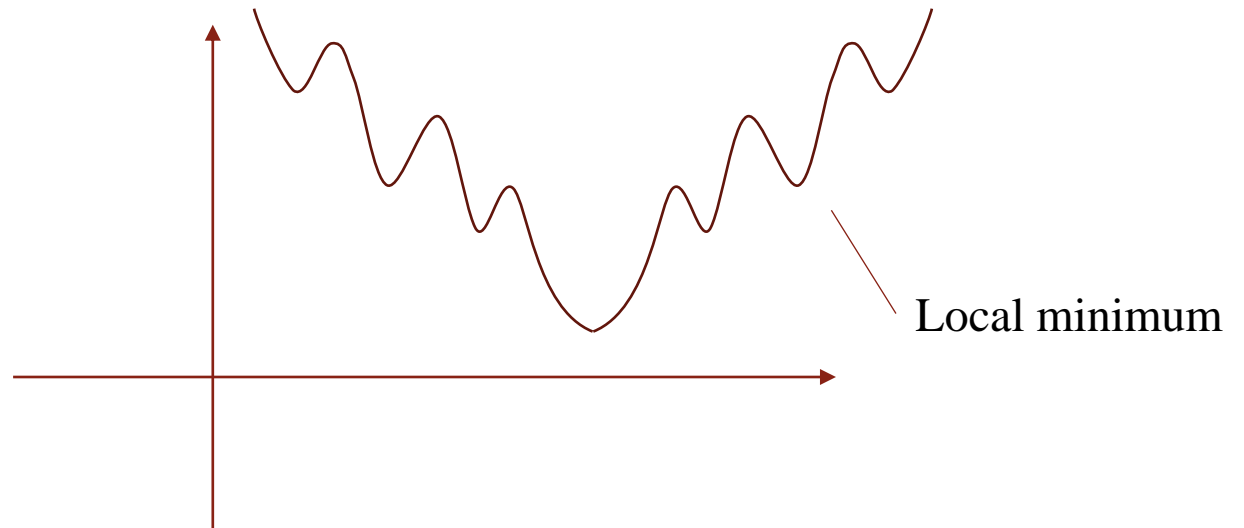
Cross-entropy

$$J(\mathbf{w}) = -\frac{1}{m} \sum_{i=1}^m [y^i \cdot \log(a^i) + (1 - y^i) \cdot \log(1 - a^i)]$$



Mean Square Error

$$J(\mathbf{w}) = -\frac{1}{m} \sum_{i=1}^m (a^i - y^i)^2$$



Cross-entropy cost function

If Hessian Matrix $\nabla^2 J(w)$ is positive determinate matrix, then, $J(w)$ is a convex function.

$$\nabla J(w) = \left[\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2}, \dots, \frac{\partial J}{\partial w_n} \right]$$

$$H = \begin{bmatrix} \frac{\partial^2 J}{\partial w_1 \partial w_1} & \frac{\partial^2 J}{\partial w_1 \partial w_2} & \dots & \frac{\partial^2 J}{\partial w_1 \partial w_q} & \dots & \frac{\partial^2 J}{\partial w_1 \partial w_n} \\ \frac{\partial^2 J}{\partial w_2 \partial w_1} & \frac{\partial^2 J}{\partial w_2 \partial w_2} & \dots & \frac{\partial^2 J}{\partial w_2 \partial w_q} & \dots & \frac{\partial^2 J}{\partial w_2 \partial w_n} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \frac{\partial^2 J}{\partial w_p \partial w_1} & \frac{\partial^2 J}{\partial w_p \partial w_2} & \dots & \frac{\partial^2 J}{\partial w_p \partial w_q} & \dots & \frac{\partial^2 J}{\partial w_p \partial w_n} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \frac{\partial^2 J}{\partial w_n \partial w_1} & \frac{\partial^2 J}{\partial w_n \partial w_2} & \dots & \frac{\partial^2 J}{\partial w_n \partial w_q} & \dots & \frac{\partial^2 J}{\partial w_n \partial w_n} \end{bmatrix}_{n \times n}$$

$$\frac{\partial^2 J}{\partial w_p \partial w_q} = \frac{\partial}{\partial w_q} \frac{\partial J}{\partial w_p}$$

Cross-entropy cost function

$$\frac{\partial^2 J}{\partial w_p \partial w_q} = \frac{\partial}{\partial w_q} \frac{\partial J}{\partial w_p}$$

$$\frac{\partial^2 J}{\partial w_p \partial w_q} = \frac{1}{m} \frac{\partial \sum_{i=1}^m (a^i - y^i) \cdot x_p^i}{\partial w_q}$$

$$\frac{\partial^2 J}{\partial w_p \partial w_q} = \frac{1}{m} \sum_{i=1}^m \frac{\partial a^i}{\partial w_q} \cdot x_p^i$$

$$\frac{\partial^2 J}{\partial w_p \partial w_q} = \frac{1}{m} \sum_{i=1}^m \frac{\partial a^i}{\partial \sum_{l=1}^n w_l \cdot x_l^i} \cdot \frac{\partial \sum_{l=1}^n w_l \cdot x_l^i}{\partial w_q} x_p^i$$

$$\frac{\partial^2 J}{\partial w_p \partial w_q} = \frac{1}{m} \sum_{i=1}^m a^i (1 - a^i) \cdot x_q^i \cdot x_p^i$$

$$\frac{\partial^2 J}{\partial w_p \partial w_q} = \frac{1}{m} \sum_{i=1}^m a^i (1 - a^i) \cdot x_p^i \cdot x_q^i$$

$$a = \sigma(\mathbf{w}\mathbf{x}) = \sigma\left(\sum_{j=1}^n w_j \cdot x_j^i\right)$$

$$\frac{\partial \sigma(z)}{\partial z} = \sigma(z) \cdot (1 - \sigma(z))$$

$$\frac{\partial J}{\partial w_j} = \frac{1}{m} \sum_{i=1}^m (a^i - y^i) \cdot x_j^i$$

Cross-entropy cost function

$$\frac{\partial^2 J}{\partial w_p \partial w_q} = \frac{1}{m} \sum_{i=1}^m a^i (1 - a^i) \cdot x_p^i \cdot x_q^i$$

$$H_{pq} = \frac{1}{m} \sum_{i=1}^m a^i (1 - a^i) \cdot x_p^i \cdot x_q^i$$



$$H = \frac{1}{m} \sum_{i=1}^m a^i (1 - a^i) \cdot \mathbf{x}^i (\mathbf{x}^i)^T$$



$\mathbf{x}^i (\mathbf{x}^i)^T$ is positive determinate matrix

Hessian Matrix H is positive determinate matrix,
then, $J(w)$ is a convex function.

$$H = \begin{bmatrix} \frac{\partial^2 J}{\partial w_1 \partial w_1} & \frac{\partial^2 J}{\partial w_1 \partial w_2} & \cdots & \frac{\partial^2 J}{\partial w_1 \partial w_q} & \cdots & \frac{\partial^2 J}{\partial w_1 \partial w_n} \\ \frac{\partial^2 J}{\partial w_2 \partial w_1} & \frac{\partial^2 J}{\partial w_2 \partial w_2} & \cdots & \frac{\partial^2 J}{\partial w_2 \partial w_q} & \cdots & \frac{\partial^2 J}{\partial w_2 \partial w_n} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \frac{\partial^2 J}{\partial w_p \partial w_1} & \frac{\partial^2 J}{\partial w_p \partial w_2} & \cdots & \frac{\partial^2 J}{\partial w_p \partial w_q} & \cdots & \frac{\partial^2 J}{\partial w_p \partial w_n} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \frac{\partial^2 J}{\partial w_n \partial w_1} & \frac{\partial^2 J}{\partial w_n \partial w_2} & \cdots & \frac{\partial^2 J}{\partial w_n \partial w_q} & \cdots & \frac{\partial^2 J}{\partial w_n \partial w_n} \end{bmatrix}_{n \times n}$$

$$\begin{bmatrix} x_1^i \\ x_2^i \\ \vdots \\ x_p^i \\ \vdots \\ x_n^i \end{bmatrix}$$

$$[x_1^i, x_2^i, \cdots, x_q^i, \cdots, x_n^i]$$

Steepest Gradient Descend Method

- Data:

- Sample: $(\mathbf{x} \in R^n, y \in \{0,1\})$, where n is the *dimension* of the input vector \mathbf{x} .
- Dataset: $D = \{(\mathbf{x}^i \in R^n, y^i \in \{0,1\}) | i \in [1, m]\}$, where m is the number of the samples.

- Logistic Regression Model:

$$a = \sigma(\mathbf{w}\mathbf{x}) = \sigma\left(\sum_{j=1}^n w_j \cdot x_j^i\right)$$

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

- Object:

$$\operatorname{argmin}_{\mathbf{w}} J(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w}} -\frac{1}{m} \sum_{i=1}^m [y^i \cdot \log(a^i) + (1 - y^i) \cdot \log(1 - a^i)]$$

Steepest Descent Algorithm

Input: D, w

for k in $1, 2, \dots, K$:

{

for i in $1, 2, \dots, m$:

$$\left\{ \begin{aligned} a^i &\leftarrow \sigma\left(\sum_{j=1}^n w_j \cdot x_j^i\right) \end{aligned} \right.$$

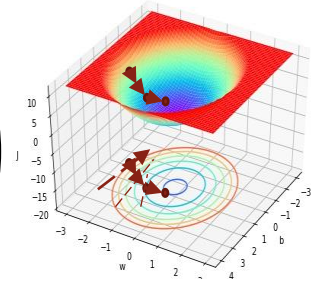
for j in $1, 2, \dots, n$:

$$\left\{ \begin{aligned} \frac{\partial J}{\partial w_j} &\leftarrow \frac{1}{m} \sum_{i=1}^m (a^i - y^i) \cdot x_j^i \end{aligned} \right.$$

for j in $1, 2, \dots, n$:

$$\left\{ \begin{aligned} w_j &\leftarrow w_j - \alpha \frac{\partial J}{\partial w_j} \end{aligned} \right.$$

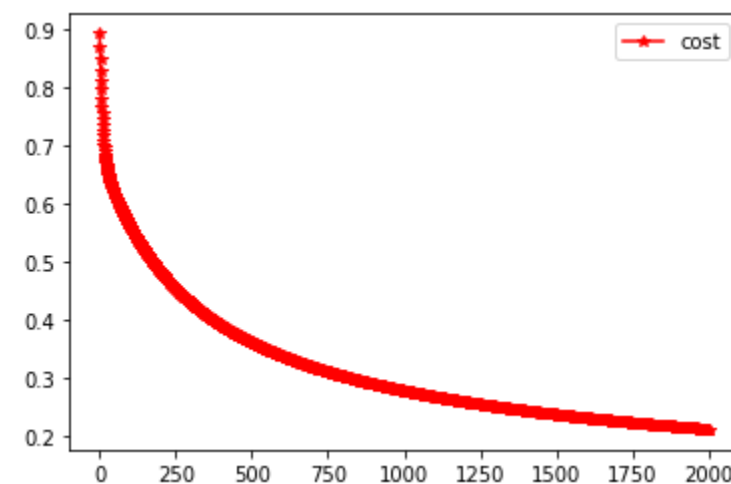
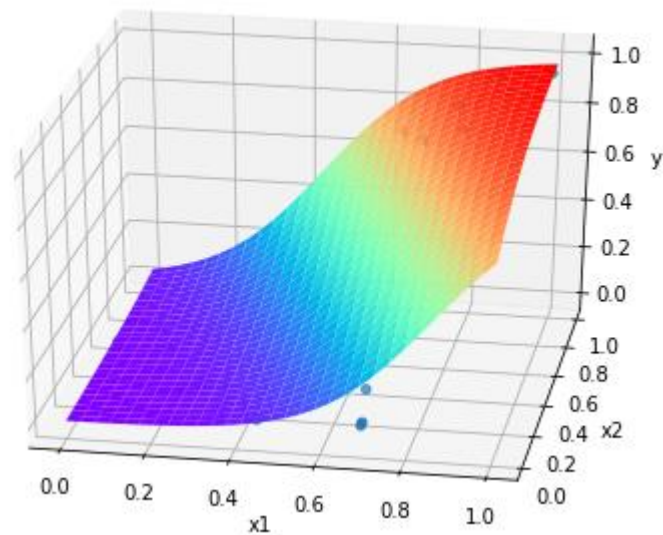
}



Examples

Size x_1	0.488	0.681	0.655	0.088	0.139	0.721	0.999	0.775	0.881	0.007	0.14	0.817	0.681	0.851	0.436
Color x_2	0.609	0.112	0.324	0.669	0.91	0.41	0.902	0.353	0.865	0.539	0.524	0.62	0.093	0.447	0.088
Label y	0	0	0	0	0	1	1	1	1	0	0	1	0	1	0

$$w = [7.52, 3.19, -6.46]$$

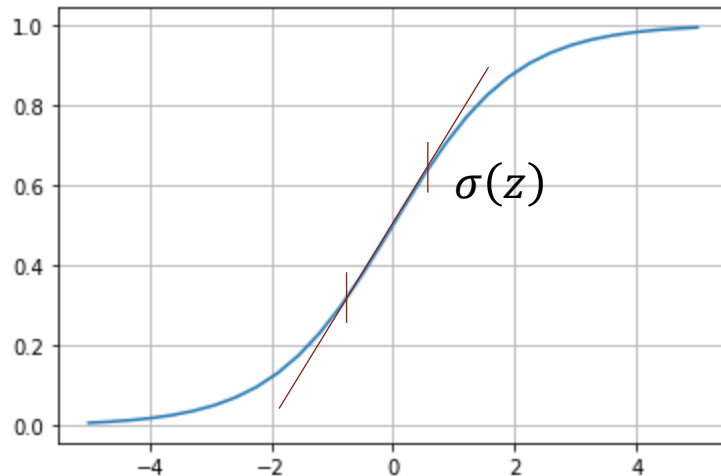


Decision Line

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

If $z_1 > z_2$ then $\sigma(z_1) > \sigma(z_2)$

If $\sigma(z_1) > \sigma(z_2)$ then $z_1 > z_2$



Threshold classifier output a at τ :

if $a \geq \tau$, predict “y=1”

if $a < \tau$, predict “y=0”.

For example $\tau = 0.5$

$$a = \sigma(\mathbf{w}\mathbf{x}) = \sigma\left(\sum_{j=1}^n w_j \cdot x_j^i\right)$$

Threshold classifier output a at 0.5 :

if $a \geq 0.5$, predict “y=1”

if $a < 0.5$, predict “y=0”.



Threshold classifier output a at 0.5 :

if $\mathbf{w}\mathbf{x} \geq 0$, predict “y=1”

if $\mathbf{w}\mathbf{x} < 0$, predict “y=0”.

$\mathbf{w}\mathbf{x} = 0$ is called the decision line

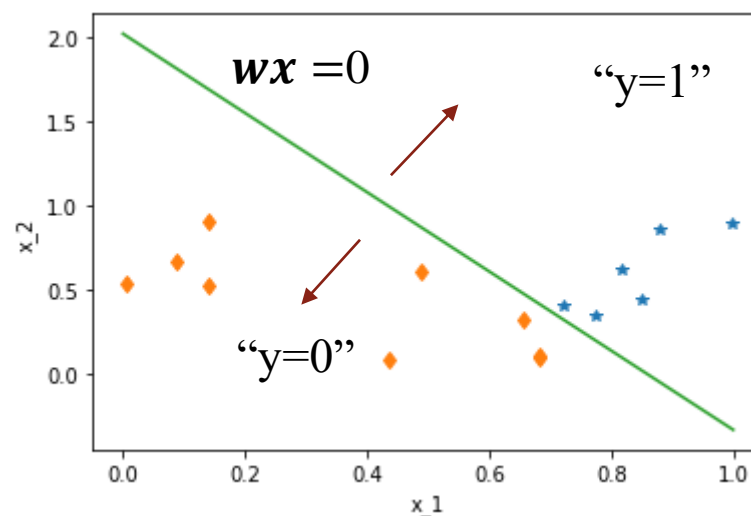
Decision Line

Size x_1	0.488	0.681	0.655	0.088	0.139	0.721	0.999	0.775	0.881	0.007	0.14	0.817	0.681	0.851	0.436
Color x_2	0.609	0.112	0.324	0.669	0.91	0.41	0.902	0.353	0.865	0.539	0.524	0.62	0.093	0.447	0.088
Label y	0	0	0	0	0	1	1	1	1	0	0	1	0	1	0

$w = [7.52, 3.19, -6.46]$

Threshold classifier output a at τ :
 if $a \geq \tau$, predict “y=1”
 if $a < \tau$, predict “y=0”.

For example $\tau = 0.5$



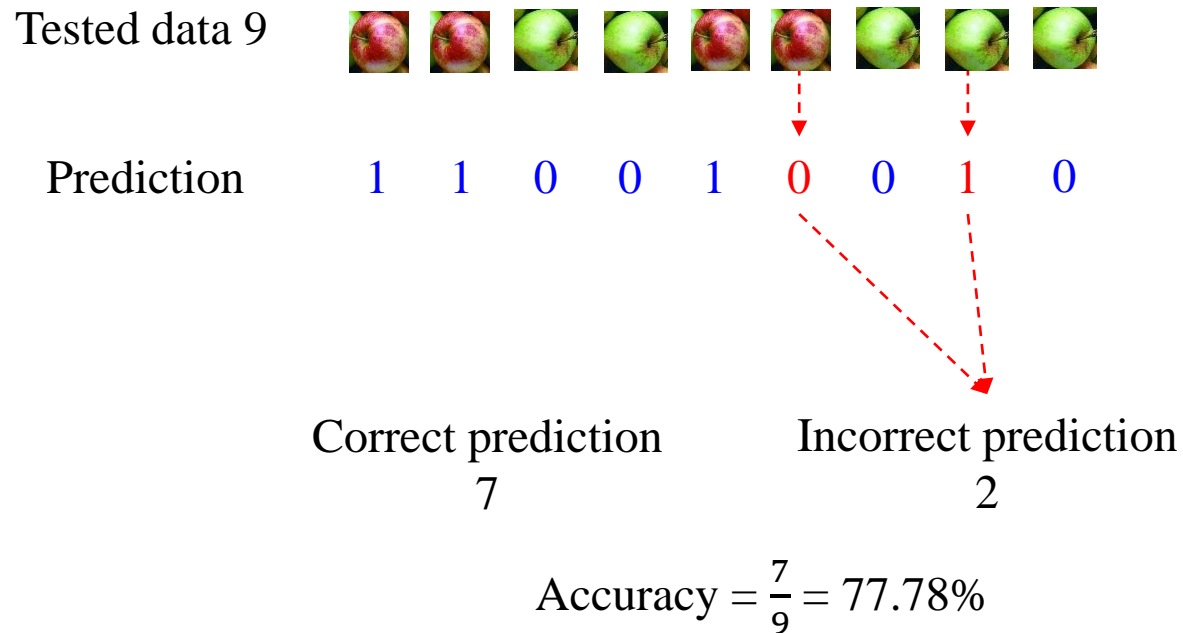
$$a = \sigma(w\mathbf{x}) = \sigma\left(\sum_{j=1}^n w_j \cdot x_j^i\right)$$

Threshold classifier output a at 0.5 :
 if $w\mathbf{x} \geq 0$, predict “y=1”
 if $w\mathbf{x} < 0$, predict “y=0”.

Accuracy

$$\text{Accuracy} = \frac{\text{number of correct prediction}}{\text{number of samples}}$$

An example



Threshold classifier output a at τ :

- if $a \geq \tau$, predict “y=1”
- if $a < \tau$, predict “y=0”.

For example $\tau = 0.5$

Test on **training** set:

- Reflect the progress of training.
- Evaluate the ability of the model to fit given data.

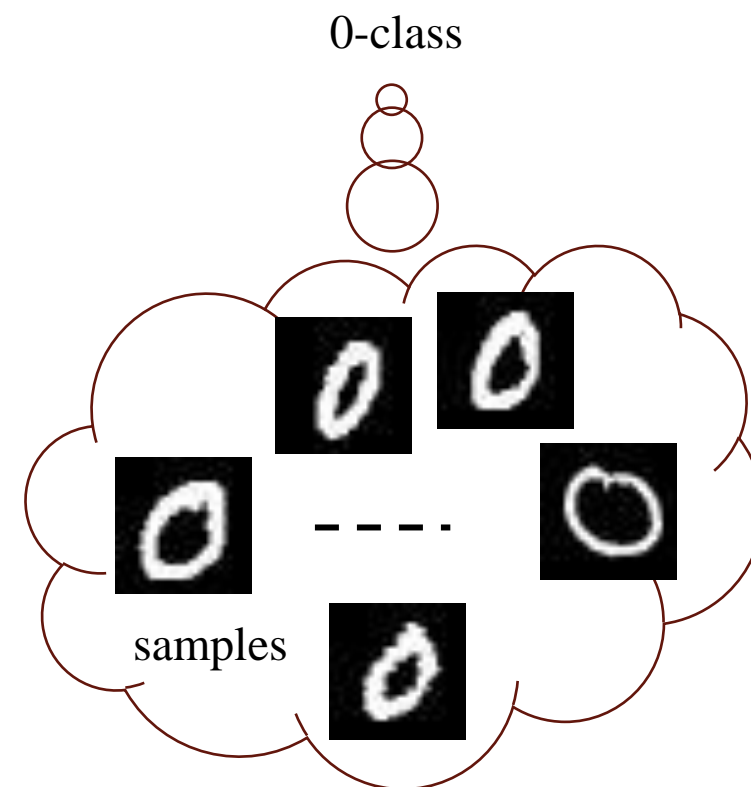
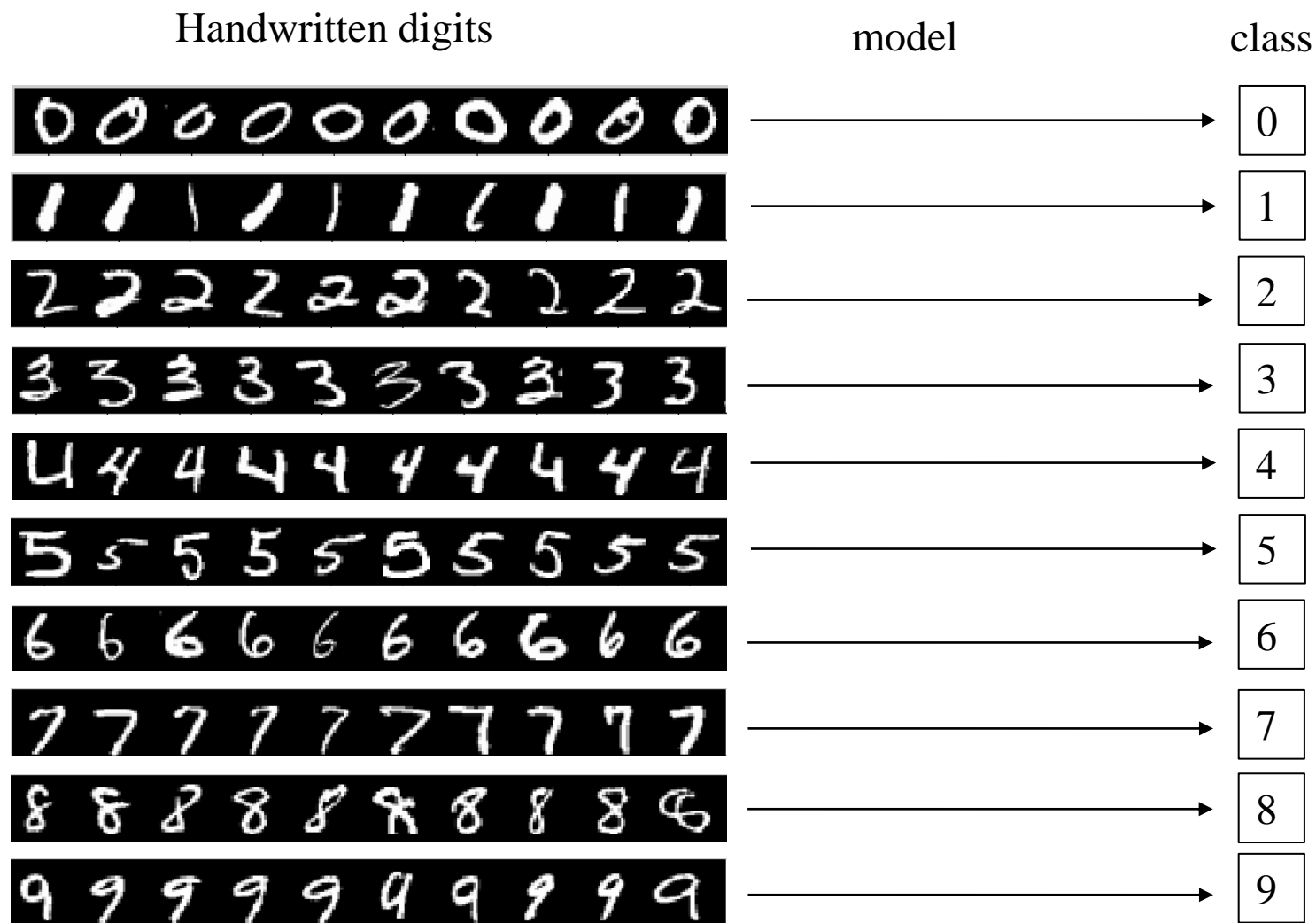
Test on **testing** set:

- Evaluate the ability of the model to generalize the knowledge.

Outline

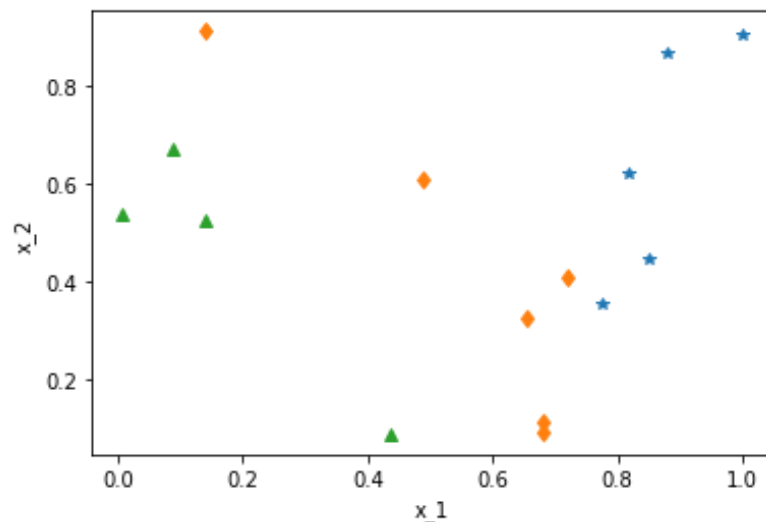
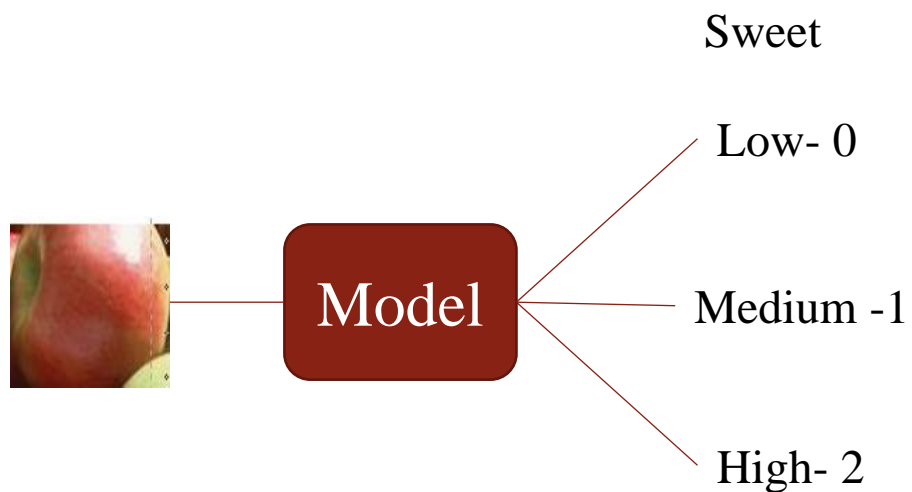
- Brief review
- Binary Classification
- Multi-class classification

Examples



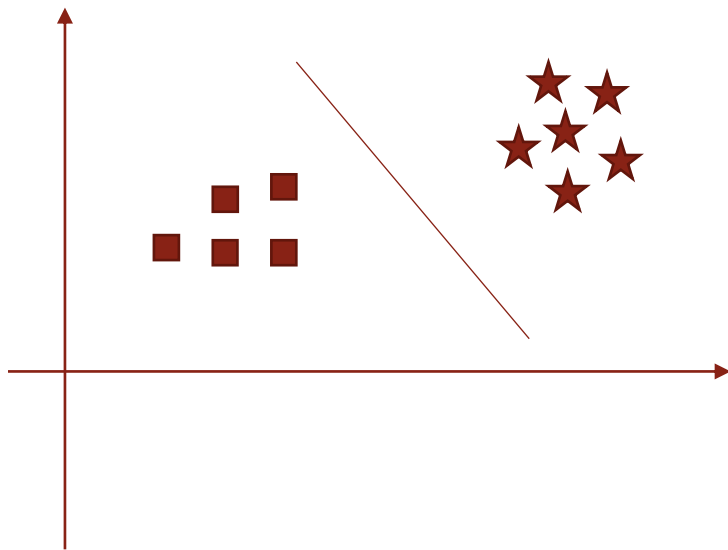
Examples

Size x_1	0.488	0.681	0.655	0.088	0.139	0.721	0.999	0.775	0.881	0.007	0.14	0.817	0.681	0.851	0.436
Color x_2	0.609	0.112	0.324	0.669	0.91	0.41	0.902	0.353	0.865	0.539	0.524	0.62	0.093	0.447	0.088
Class sweet	medium 1	medium 1	medium 1	low 0	medium 1	medium 1	high 2	high 2	high 2	low 0	low 0	high 2	medium 1	high 2	low 0



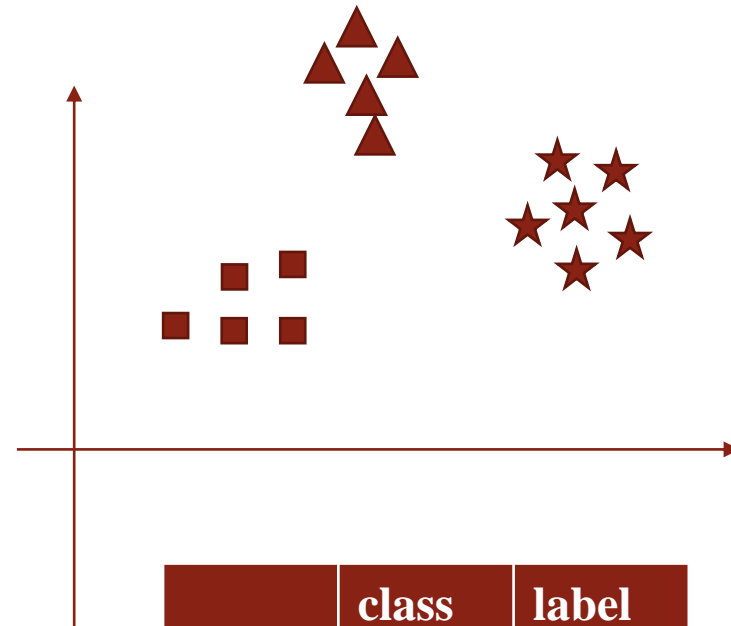
Examples

Binary classification



	class	label
★	1	0
■	2	1

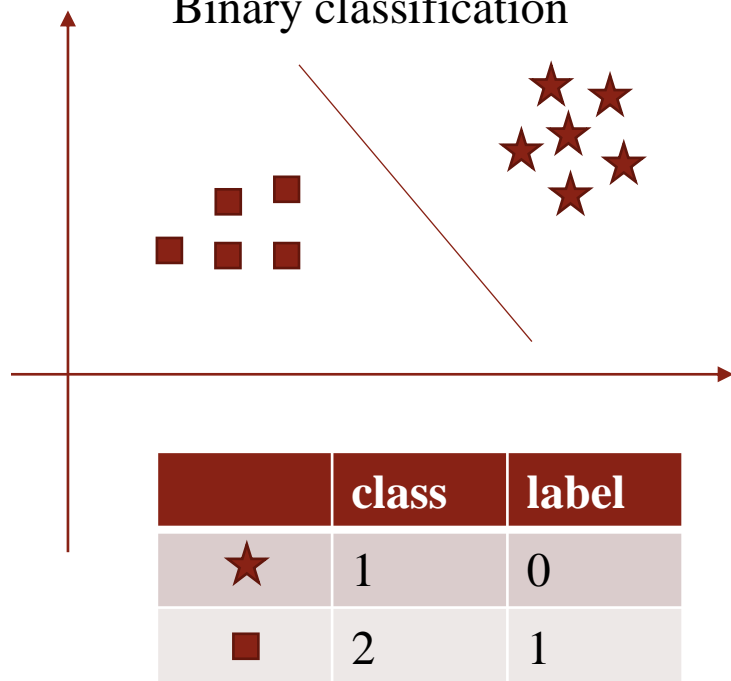
Multi-class classification



	class	label
★	1	?
■	2	?
▲	3	?

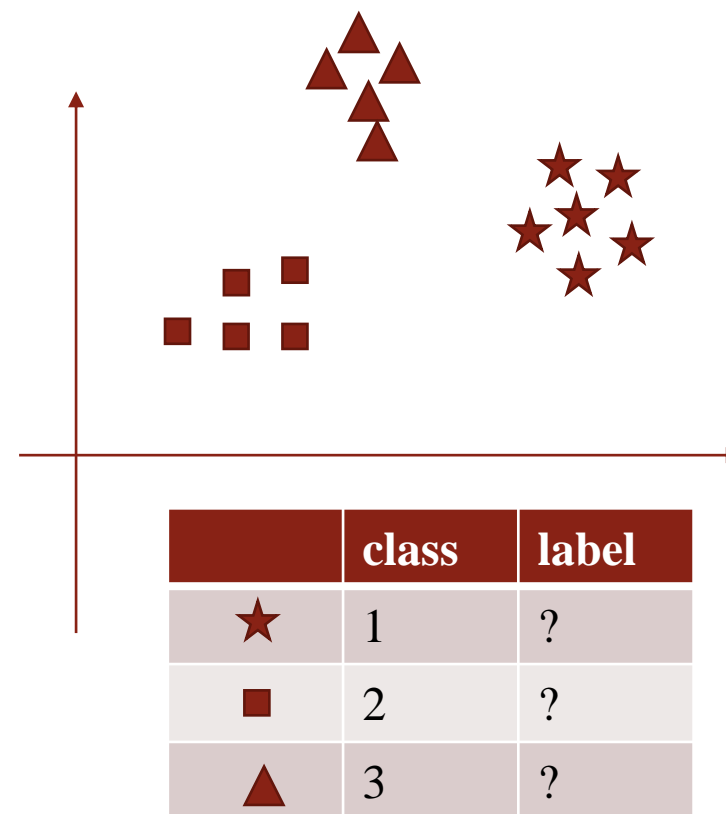
Examples

Binary classification



Threshold classifier output a at τ :
if $a \geq \tau$, predict “y=1”
if $a < \tau$, predict “y=0”.

Multi-class classification



How to interpret the output a ?

Thanks!!!