

Logistic Regression for Classification

Bernoulli trials

Bernoulli trials are experiments with binary (two) outcome

Probability of success = p, then probability of failure = q = 1-p

Transform binary classification problem into probabilistic framework.

Based upon the features, we assign objects to classes probabilistically.

Let $\mathbf{x}_i = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{in} \\ 1 \end{pmatrix}$ be the feature vector of i th object.

Corresponding class label be $y_i \in [-1, 1]$

Here actual feature vector is $\mathbf{x}_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{in} \end{pmatrix}$ $\mathbf{1}$ is appended at the end for the following reason

Why append 1 to the feature vector

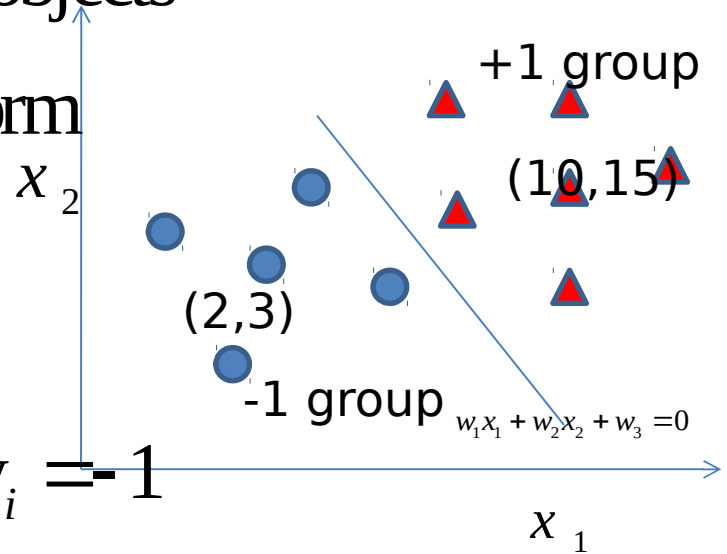
Let (x_1, x_2) be the two feature values of objects

A binary linear classifier will be of the form

$$f(x|w) = w_1x_1 + w_2x_2 + w_3$$

if $f(x_i | w) > 0$, class label is $y_i = 1$ else $y_i = -1$

or if $w_1x_1 + w_2x_2 + w_3 > 0$, class label is $y_i = 1$ else $y_i = -1$



If we append 1 with data vector, we can write the classifier

in a more compact form: if $w^T x > 0$, class label is $y_i = 1$ else $y_i = -1$

How the data look like

$$A = \begin{bmatrix} x_{11} & x_{12} & 1 \\ x_{21} & x_{22} & 1 \\ x_{31} & x_{32} & 1 \\ \vdots & \vdots & \vdots \\ x_{m1} & x_{m2} & 1 \end{bmatrix}; y = \begin{bmatrix} 1 \\ -1 \\ 1 \\ \vdots \\ -1 \end{bmatrix}$$

Note, here data vectors are in rows

Some time, data can be in this format

$$A_1 = \begin{bmatrix} x_{11} & x_{21} & x_{31} & \cdots & x_{m1} \\ x_{12} & x_{22} & x_{32} & \cdots & x_{m2} \\ 1 & 1 & 1 & \cdots & 1 \end{bmatrix};$$

Problem

Find a \mathbf{w} such that $\mathbf{w}^T \mathbf{x}_i > 0$ if $y_i = 1$

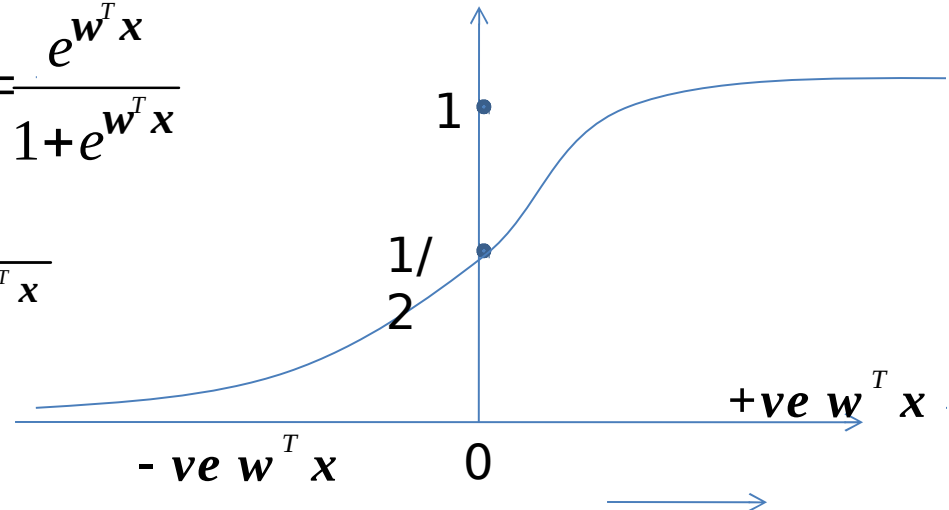
also $\mathbf{w}^T \mathbf{x}_j < 0$ if $y_j = -1$

\mathbf{w} should be such that for most **training** data,
the above relation is satisfied

Once we obtain \mathbf{w} , we give probabilistic statement for class / label prediction
using the formula

$$\text{Prob}(y = 1 | \mathbf{x}; \mathbf{w}) = p = f(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} = \frac{e^{\mathbf{w}^T \mathbf{x}}}{1 + e^{\mathbf{w}^T \mathbf{x}}}$$

$$\begin{aligned} \text{Prob}(y = -1 | \mathbf{x}; \mathbf{w}) &= 1 - p = 1 - \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} \\ &= \frac{e^{-\mathbf{w}^T \mathbf{x}}}{1 + e^{-\mathbf{w}^T \mathbf{x}}} = \frac{1}{1 + e^{\mathbf{w}^T \mathbf{x}}} \end{aligned}$$



Training the Classifier

We assume we are given labeled data of the form

$$A = \begin{bmatrix} x_{11} & x_{12} & 1 \\ x_{21} & x_{22} & 1 \\ x_{31} & x_{32} & 1 \\ \vdots & \vdots & \vdots \\ x_{m1} & x_{m2} & 1 \end{bmatrix}; y = \begin{bmatrix} 1 \\ -1 \\ 1 \\ \vdots \\ -1 \end{bmatrix}$$

*Training **involves** Finding a w such that $w^T x_i > 0$ if $y_i = 1$*

also $w^T x_j < 0$ if $y_j = -1$

w should be such that for most data the above relation is satisfied

Objective function

To obtain \mathbf{w} , we need an objective function that connects our data and \mathbf{w}

We use the concept of **likelihood function** from probability theory.

$$\text{Prob}(y=y_i=1 | \mathbf{x}_i; \mathbf{w}) = p_i = \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}_i}} = \frac{1}{1+e^{-y_i \mathbf{w}^T \mathbf{x}_i}}$$

$$\text{Prob}(y=y_i=-1 | \mathbf{x}_i; \mathbf{w}) = 1 - p_i = \frac{1}{1+e^{\mathbf{w}^T \mathbf{x}_i}} = \frac{1}{1+e^{-y_i \mathbf{w}^T \mathbf{x}_i}}$$

$$\therefore \text{Prob}(y=y_i | \mathbf{x}_i; \mathbf{w}) = \frac{1}{1+e^{-y_i \mathbf{w}^T \mathbf{x}_i}}$$

**This is for unknown data
assuming \mathbf{w} is known**

For training data, data_label is already known

$$\text{likelihood of the data } L(\mathbf{w}) = \prod_{i=1}^m \frac{1}{1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i}}$$

It is product of probabilities for the given data in terms of \mathbf{w}

To simplify we take log

$$\begin{aligned} \text{Log likelihood of the data } \text{Log} L(\mathbf{w}) &= \prod_{i=1}^m \frac{1}{1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i}} \\ &\quad - \sum_{i=1}^m \log(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i}) \end{aligned}$$

What is likelihood function

function that gives joint probability of occurrence of the data under the assumption that data is drawn independently from a given distribution.

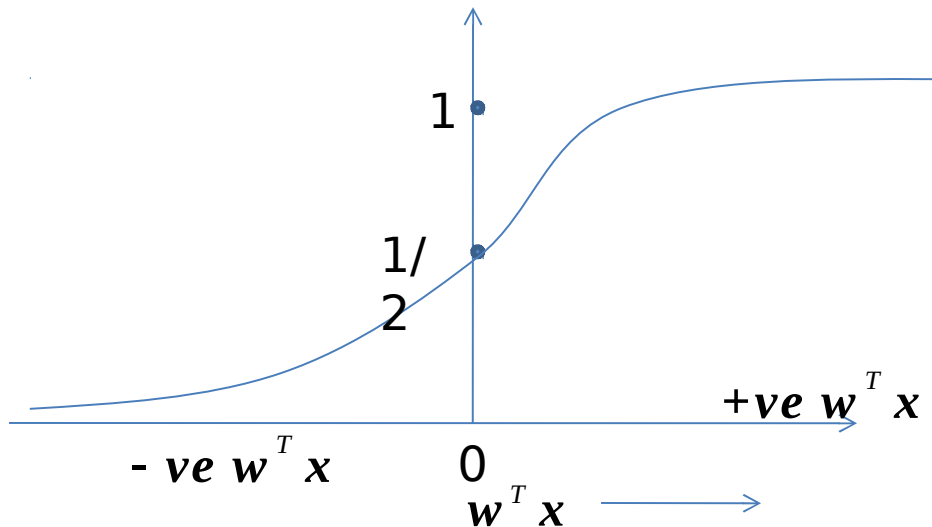
For a given \mathbf{w} , we can assign a probability to each data vector \mathbf{x}_i to belong to the given category (already given label).

To maximize the joint probability (for whole data) by finding a proper \mathbf{w} , for which, we use optimization theory

Logistic Regression

$$\text{Prob}(y = 1 | \mathbf{x}; \mathbf{w}) = p = f(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} = \frac{e^{\mathbf{w}^T \mathbf{x}}}{1 + e^{\mathbf{w}^T \mathbf{x}}}$$

$$\text{Prob}(y = -1 | \mathbf{x}; \mathbf{w}) = 1 - p = 1 - \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} = \frac{e^{-\mathbf{w}^T \mathbf{x}}}{1 + e^{-\mathbf{w}^T \mathbf{x}}} = \frac{1}{1 + e^{\mathbf{w}^T \mathbf{x}}}$$



By optimization theory we will find a \mathbf{w} such that

if $y_i = 1$ then $\mathbf{w}^T \mathbf{x}_i > 0$

else if $y_j = -1$ then $\mathbf{w}^T \mathbf{x}_j < 0$

formulation

class label be $y_i \in [1, -1]$

$$\text{Prob}(y=y_i=1 | \mathbf{x}_i; \mathbf{w}) = p_i = \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}_i}} = \frac{1}{1+e^{-y_i \mathbf{w}^T \mathbf{x}_i}}$$

$$\text{Prob}(y=y_i=-1 | \mathbf{x}_i; \mathbf{w}) = 1 - p_i = \frac{1}{1+e^{\mathbf{w}^T \mathbf{x}_i}} = \frac{1}{1+e^{-y_i \mathbf{w}^T \mathbf{x}_i}}$$

$$\therefore \text{Prob}(y=y_i | \mathbf{x}_i; \mathbf{w}) = \frac{1}{1+e^{-y_i \mathbf{w}^T \mathbf{x}_i}}$$

$$\text{likelihood of the data} = \prod_{i=1}^m \frac{1}{1+e^{-y_i \mathbf{w}^T \mathbf{x}_i}}$$

$$\text{Log likelihood} = \log \prod_{i=1}^m \frac{1}{1+e^{-y_i \mathbf{w}^T \mathbf{x}_i}} = - \sum_{i=1}^m \log(1+e^{-y_i \mathbf{w}^T \mathbf{x}_i})$$

We should maximize this by adjusting

$$\text{negative Log likelihood} = \log \prod_{i=1}^m \frac{1}{1+e^{-y_i \mathbf{w}^T \mathbf{x}_i}} = \sum_{i=1}^m \log(1+e^{-y_i \mathbf{w}^T \mathbf{x}_i})$$

gradient :

We should minimize this by adjusting \mathbf{w}

$$\nabla J(\mathbf{w}) = \sum_{i=1}^m \frac{-y_i e^{-y_i \mathbf{w}^T \mathbf{x}_i}}{(1+e^{-y_i \mathbf{w}^T \mathbf{x}_i})} \mathbf{x}_i$$

Programming

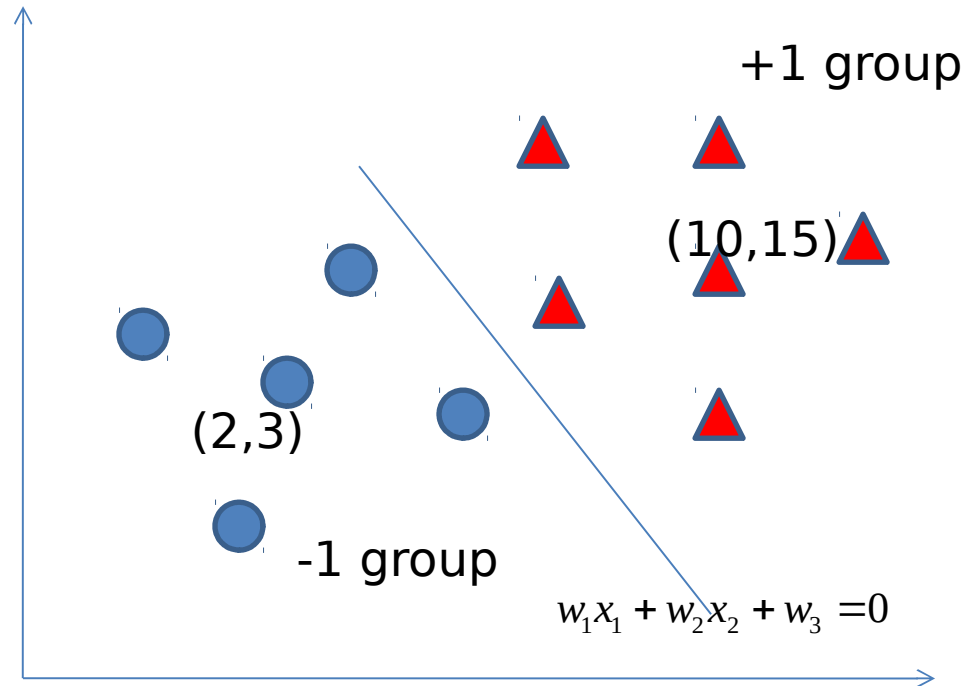
$$J(\mathbf{w}) = \sum_{i=1}^m \log(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i})$$

gradient :

$$\nabla J(\mathbf{w}) = \sum_{i=1}^m \frac{-y_i e^{-y_i \mathbf{w}^T \mathbf{x}_i}}{(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i})} \mathbf{x}_i$$

$$= \sum_i a_i \mathbf{x}_i = \text{linear combination of data_vectors}$$

$$A = \begin{bmatrix} x_{11} & x_{12} & 1 \\ x_{21} & x_{22} & 1 \\ x_{31} & x_{32} & 1 \\ \vdots & \vdots & \vdots \\ x_{m1} & x_{m2} & 1 \end{bmatrix}; y = \begin{bmatrix} 1 \\ -1 \\ 1 \\ \vdots \\ -1 \end{bmatrix}$$



Let us create 10 in the category +1 and another 10 in category -1

```
c1=[10;15];
c2=[2;3];
B=[randn(2,10)+repmat(c1,1,10)
  randn(2,10)+repmat(c2,1,10);ones(1,20) ]
A=B';
y=[ones(10,1); -1*ones(10,1)];
```

Computing J(w)

$$A = \begin{bmatrix} x_{11} & x_{12} & 1 \\ x_{21} & x_{22} & 1 \\ x_{31} & x_{32} & 1 \\ \vdots & \vdots & \vdots \\ x_{m1} & x_{m2} & 1 \end{bmatrix}; y = \begin{bmatrix} 1 \\ -1 \\ 1 \\ \vdots \\ -1 \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

$$J(\mathbf{w}) = \sum_{i=1}^m \log(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i})$$

$$J(\mathbf{w}) = \sum_{i=1}^m \log(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i}) = \text{sum}(\log(1 + \exp(-1 * y * A \mathbf{w})))$$

$$A * \mathbf{w} = \begin{bmatrix} x_{11} & x_{12} & 1 \\ x_{21} & x_{22} & 1 \\ x_{31} & x_{32} & 1 \\ \vdots & \vdots & \vdots \\ x_{m1} & x_{m2} & 1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} = \begin{bmatrix} \mathbf{w}^T \mathbf{x}_1 \\ \mathbf{w}^T \mathbf{x}_2 \\ \vdots \\ \mathbf{w}^T \mathbf{x}_{m-1} \\ \mathbf{w}^T \mathbf{x}_m \end{bmatrix}; \quad \exp(-1 * y * (A * \mathbf{w})) = \begin{bmatrix} e^{-y_1 \mathbf{w}^T \mathbf{x}_1} \\ e^{-y_2 \mathbf{w}^T \mathbf{x}_2} \\ \vdots \\ e^{-y_{m-1} \mathbf{w}^T \mathbf{x}_{m-1}} \\ e^{-y_m \mathbf{w}^T \mathbf{x}_m} \end{bmatrix}$$

$$1 + \exp(y * A \mathbf{w}) = \begin{bmatrix} 1 + e^{-y_1 \mathbf{w}^T \mathbf{x}_1} \\ 1 + e^{-y_2 \mathbf{w}^T \mathbf{x}_2} \\ \vdots \\ 1 + e^{-y_{m-1} \mathbf{w}^T \mathbf{x}_{m-1}} \\ 1 + e^{-y_m \mathbf{w}^T \mathbf{x}_m} \end{bmatrix} \quad \log(1 + y * A \mathbf{w}) = \log \begin{bmatrix} 1 + e^{-y_1 \mathbf{w}^T \mathbf{x}_1} \\ 1 + e^{-y_2 \mathbf{w}^T \mathbf{x}_2} \\ \vdots \\ 1 + e^{-y_{m-1} \mathbf{w}^T \mathbf{x}_{m-1}} \\ 1 + e^{-y_m \mathbf{w}^T \mathbf{x}_m} \end{bmatrix} = \begin{bmatrix} \log(1 + e^{-y_1 \mathbf{w}^T \mathbf{x}_1}) \\ \log(1 + e^{-y_2 \mathbf{w}^T \mathbf{x}_2}) \\ \vdots \\ \log(1 + e^{-y_{m-1} \mathbf{w}^T \mathbf{x}_{m-1}}) \\ \log(1 + e^{-y_m \mathbf{w}^T \mathbf{x}_m}) \end{bmatrix}$$

Computing $J(\mathbf{w})$

$$J(\mathbf{w}) = \sum_{i=1}^m \log(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i})$$

$$J(\mathbf{w}) = \sum_{i=1}^m \log(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i}) = \text{sum}(\log(1 + \exp(-1 * \mathbf{y} * \mathbf{A} * \mathbf{w})))$$

Computing $\nabla J(\mathbf{w})$

gradient :

$$\nabla J(\mathbf{w}) = \sum_{i=1}^m \frac{-y_i e^{-y_i \mathbf{w}^T \mathbf{x}_i}}{(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i})} \mathbf{x}_i$$

$$1 + \exp(-\mathbf{y} \cdot (\mathbf{A} \cdot \mathbf{w})) = \begin{bmatrix} 1 + e^{-y_1 \mathbf{w}^T \mathbf{x}_1} \\ 1 + e^{-y_2 \mathbf{w}^T \mathbf{x}_2} \\ \vdots \\ 1 + e^{-y_{m-1} \mathbf{w}^T \mathbf{x}_{m-1}} \\ 1 + e^{-y_m \mathbf{w}^T \mathbf{x}_m} \end{bmatrix} \quad -\mathbf{y} \cdot \exp(\mathbf{y} \cdot (\mathbf{A} \cdot \mathbf{w})) = \begin{bmatrix} -y_1 e^{-y_1 \mathbf{w}^T \mathbf{x}_1} \\ -y_2 e^{-y_2 \mathbf{w}^T \mathbf{x}_2} \\ \vdots \\ -y_{m-1} e^{-y_{m-1} \mathbf{w}^T \mathbf{x}_{m-1}} \\ -y_m e^{-y_m \mathbf{w}^T \mathbf{x}_m} \end{bmatrix}$$

Computing $\nabla J(\mathbf{w})$

$$(-y \cdot \exp(y \cdot (A \cdot w))) ./ (1 + y \cdot A \cdot w) =$$

$$\begin{bmatrix} \frac{-y_1 e^{-y_1 \mathbf{w}^T \mathbf{x}_1}}{(1 + e^{-y_1 \mathbf{w}^T \mathbf{x}_1})} \\ \frac{-y_2 e^{-y_2 \mathbf{w}^T \mathbf{x}_2}}{(1 + e^{-y_2 \mathbf{w}^T \mathbf{x}_2})} \\ \vdots \\ \frac{-y_{m-1} e^{-y_{m-1} \mathbf{w}^T \mathbf{x}_{m-1}}}{(1 + e^{-y_{m-1} \mathbf{w}^T \mathbf{x}_{m-1}})} \\ \frac{-y_m e^{-y_m \mathbf{w}^T \mathbf{x}_m}}{(1 + e^{-y_m \mathbf{w}^T \mathbf{x}_m})} \end{bmatrix}$$

gradient :

$$\nabla J(\mathbf{w}) = \sum_{i=1}^m \frac{-y_i e^{-y_i \mathbf{w}^T \mathbf{x}_i}}{(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i})} \mathbf{x}_i$$

Computing $\nabla J(\mathbf{w})$

$$\nabla J(\mathbf{w}) = \nabla J(\mathbf{w}) = \sum_{i=1}^m \frac{-y_i e^{-y_i \mathbf{w}^T \mathbf{x}_i}}{(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i})} \mathbf{x}_i$$

= *linear combination of* \mathbf{x}_i

= *linear combination of columns of* A^T

$$= A^T * (-y .* \exp(-y .* (A * w))) ./ (1 + \exp(-y .* A * w))$$

Matlab Function to evaluate

$$J(\mathbf{w}); \nabla J(\mathbf{w})$$

```
[Jfun, Jgrad] = Mylogistic(A, y,  
w)  
com=exp(-1*y.*A*w;  
Jfun=sum(log(1+com));  
GradJ=A'*(-y.*com)./(1+com);  
end
```

Now apply gradient method to find the solution

$$A^T * (-y.*\exp(-y.*(A*w)))./(1+\exp(-y.*A*w))$$

New formulation

L2 - regularized logistic regression

$$\min_{\mathbf{w}} J(\mathbf{w}) = \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^n \log(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i})$$

gradient :

$$\nabla J(\mathbf{w}) = \lambda \mathbf{w} + \sum_{i=1}^n \frac{-y_i e^{-y_i \mathbf{w}^T \mathbf{x}_i}}{(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i})} \mathbf{x}_i$$