

CSC 4020 Fundamentals of Machine Learning: Support Vector Machine I

Baoyuan Wu
School of Data Science, CUHK-SZ

February 24, 2021

large

margin

Outline

- 1 Motivation
- 2 Derivation I: large margin
- 3 Derivation II: hinge loss

Classification

Binary classification:

Classification

Binary classification:

- Given training data set $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, and $\mathbf{x}_i \in \mathbb{R}^n$, $y_i \in \{-1, +1\}$

Classification

Binary classification:

- Given training data set $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$, and $\mathbf{x}_i \in \mathbb{R}^n, y_i \in \{-1, +1\}$
- We adopt the sign hypothesis function $y = \text{sign}(f_{\mathbf{w}}(\mathbf{x})) = \text{sign}(\mathbf{w}^\top \mathbf{x})$

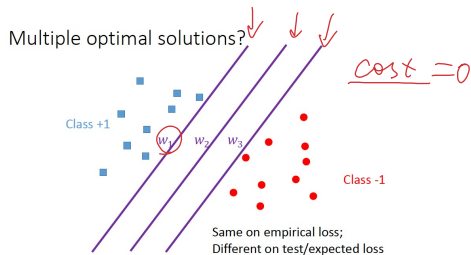
Classification

Binary classification:

- Given training data set $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$, and $\mathbf{x}_i \in \mathbb{R}^n, y_i \in \{-1, +1\}$
- We adopt the sign hypothesis function $y = \text{sign}(f_{\mathbf{w}}(\mathbf{x})) = \text{sign}(\mathbf{w}^\top \mathbf{x})$
- Then, we require that
 - If $y_i = +1$, then $\mathbf{w}^\top \mathbf{x} > 0$
 - If $y_i = -1$, then $\mathbf{w}^\top \mathbf{x} < 0$

$$\text{sign}(\mathbf{w}^\top \mathbf{x}) = 1$$
$$-1$$

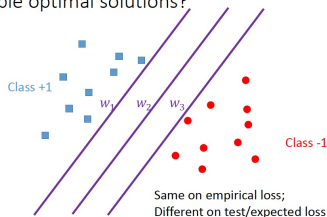
Classification



- There could be multiple decision boundaries to perfectly separate the above data.

Classification

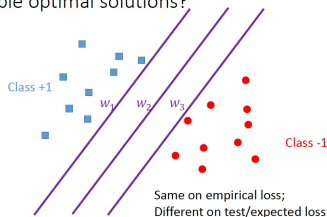
Multiple optimal solutions?



- There could be multiple decision boundaries to perfectly separate the above data. Why?

Classification

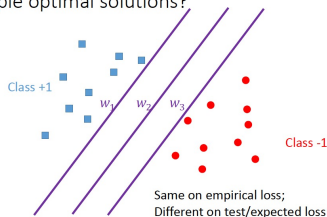
Multiple optimal solutions?



- There could be multiple decision boundaries to perfectly separate the above data. Why?
- For standard logistic regression, the objective function (*i.e.*, cross entropy loss) is convex, rather than strongly/strictly convex. Consequently, there are multiple values of parameters that can perfectly fit the training data.

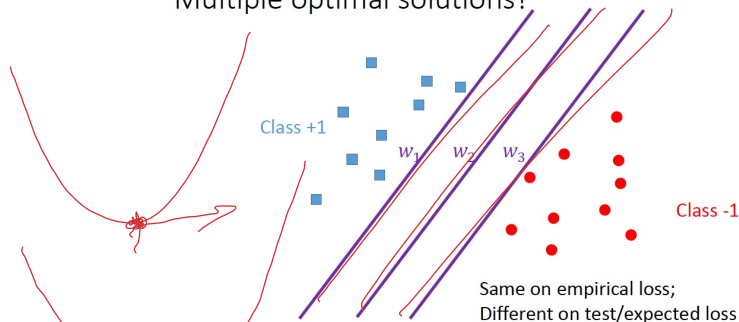
Classification

Multiple optimal solutions?



- There could be multiple decision boundaries to perfectly separate the above data. Why?
- For standard logistic regression, the objective function (*i.e.*, cross entropy loss) is convex, rather than strongly/strictly convex. Consequently, there are multiple values of parameters that can perfectly fit the training data.
- For regularized logistic regression, the objective function (*i.e.*, cross entropy loss + $\lambda \cdot \ell_2$ regularization) is strictly convex, which has the unique optimal solution. However, it depends on the trade-off hyper-parameter λ . For sure you can use cross-validation to use a suitable λ , but is there any more elegant approach?

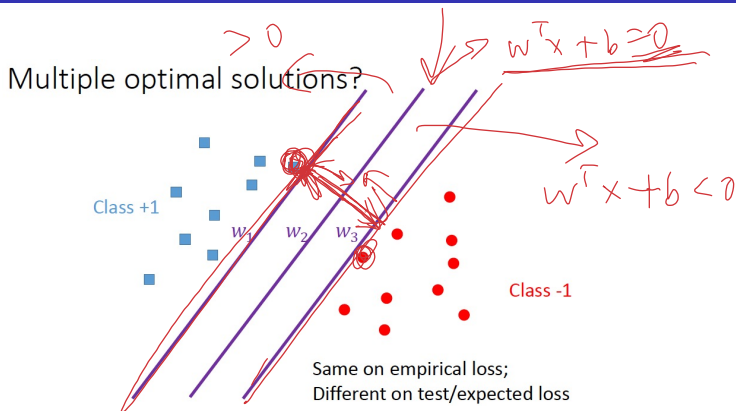
Multiple optimal solutions?



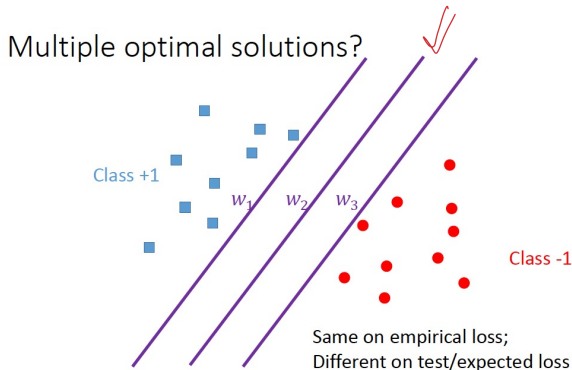
- Just following your intuition, which decision boundary do you prefer?

$$J(\theta) = \underbrace{CE(\theta)}_{\Delta} + \underbrace{\lambda \|\theta\|^2}_{\Delta} \quad [H_0 + \lambda I \geq 0]$$

Classification



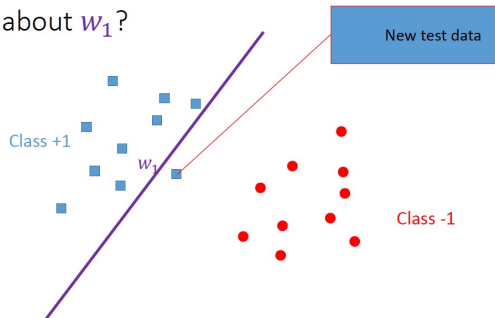
- Just following your intuition, which decision boundary do you prefer?
- The middle one (*i.e.*, $w_2^T x = 0$) seems better, as it is far from data of both positive and negative classes.



- Just following your intuition, which decision boundary do you prefer?
- The middle one (*i.e.*, $\mathbf{w}_2^\top \mathbf{x} = 0$) seems better, as it is far from data of both positive and negative classes.
- How to model such intuition?

Classification

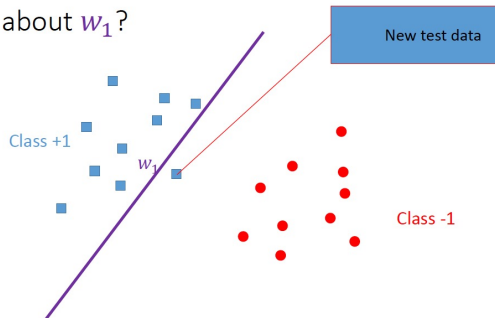
What about w_1 ?



- Just following your intuition, which decision boundary do you prefer?

Classification

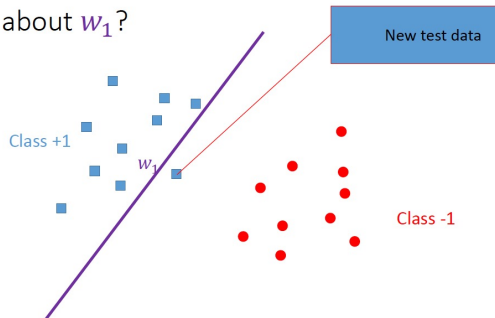
What about w_1 ?



- Just following your intuition, which decision boundary do you prefer?
- The middle one (*i.e.*, $w_2^\top x = 0$) seems better, as it is far from data of both positive and negative classes.

Classification

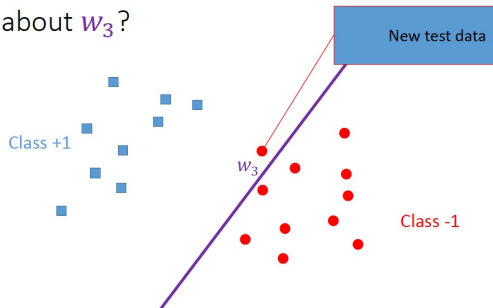
What about w_1 ?



- Just following your intuition, which decision boundary do you prefer?
- The middle one (*i.e.*, $w_2^\top x = 0$) seems better, as it is far from data of both positive and negative classes.
- How to model such intuition?

Classification

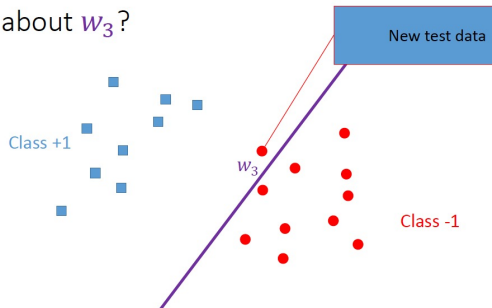
What about w_3 ?



- Just following your intuition, which decision boundary do you prefer?

Classification

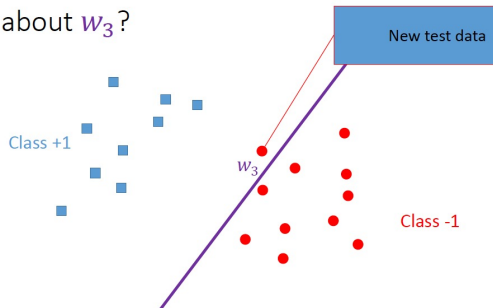
What about w_3 ?



- Just following your intuition, which decision boundary do you prefer?
- The middle one (*i.e.*, $w_2^\top x = 0$) seems better, as it is far from data of both positive and negative classes.

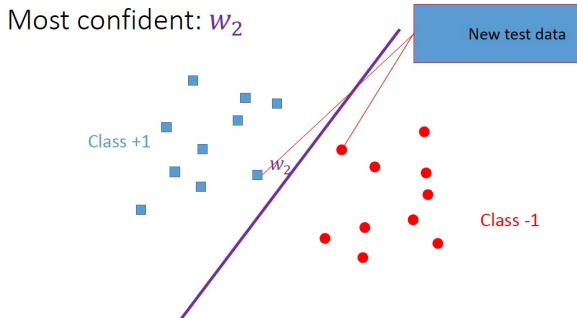
Classification

What about w_3 ?



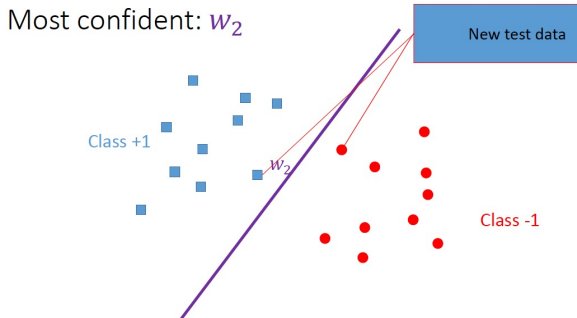
- Just following your intuition, which decision boundary do you prefer?
- The middle one (*i.e.*, $w_2^\top x = 0$) seems better, as it is far from data of both positive and negative classes.
- How to model such intuition?

Classification



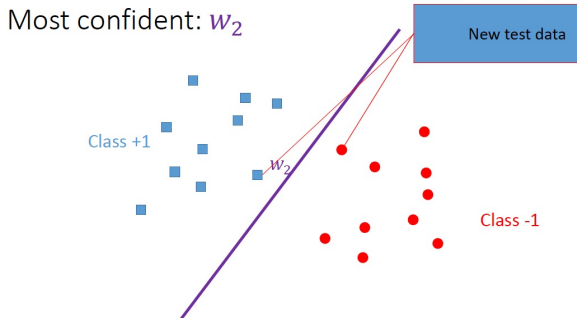
- Just following your intuition, which decision boundary do you prefer?

Classification



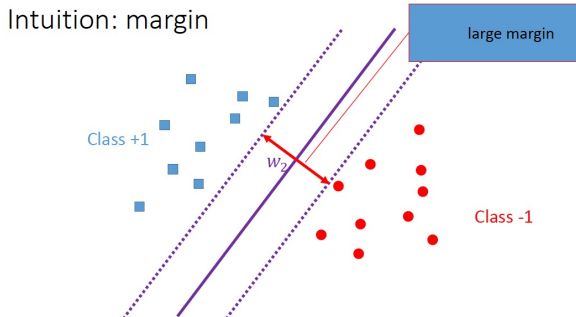
- Just following your intuition, which decision boundary do you prefer?
- The middle one (*i.e.*, $w_2^\top x = 0$) seems better, as it is far from data of both positive and negative classes.

Classification



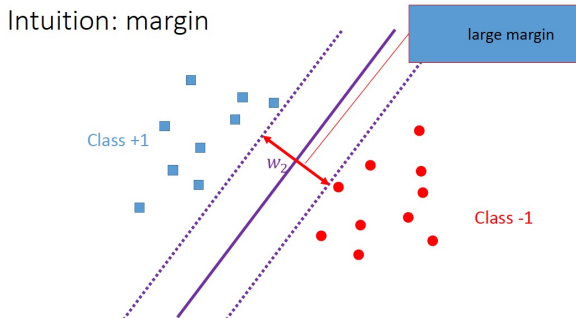
- Just following your intuition, which decision boundary do you prefer?
- The middle one (*i.e.*, $w_2^\top x = 0$) seems better, as it is far from data of both positive and negative classes.
- How to model such intuition?

Large margin intuition



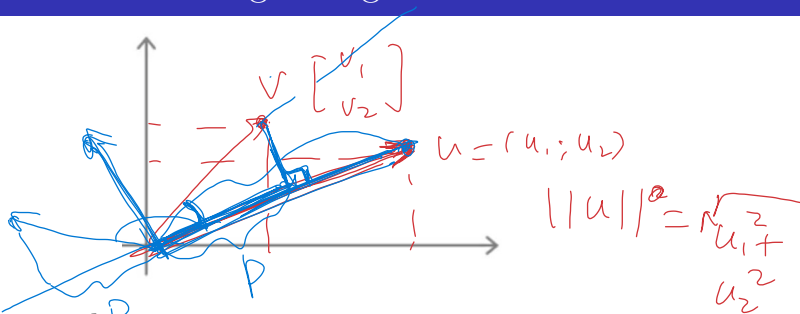
- We introduce the concept **margin**: the distance from the closest point of positive and negative classes to the decision boundary

Large margin intuition



- We introduce the concept **margin**: the distance from the closest point of positive and negative classes to the decision boundary
- The intuition is to choose the decision boundary with large margin, which is called **large margin classifier**, also called **support vector machine (SVM)**

Mathematics behind large margin classification



Inner vector product:

$$\bullet \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \nu = \begin{bmatrix} \nu_1 \\ \nu_2 \end{bmatrix}$$

$$u^T v = u_1 v_1 + u_2 v_2 \Rightarrow$$

$$= \frac{p}{>0} ||u||$$

Mathematics behind large margin classification



Inner vector product:

- $\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$, $\boldsymbol{\nu} = \begin{bmatrix} \nu_1 \\ \nu_2 \end{bmatrix}$
- $\|\boldsymbol{\mu}\| = \sqrt{\mu_1^2 + \mu_2^2}$, the length of $\boldsymbol{\mu}$

Mathematics behind large margin classification



Inner vector product:

- $\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \boldsymbol{\nu} = \begin{bmatrix} \nu_1 \\ \nu_2 \end{bmatrix}$
- $\|\boldsymbol{\mu}\| = \sqrt{\mu_1^2 + \mu_2^2}$, the length of $\boldsymbol{\mu}$
- $\boldsymbol{\mu}^\top \boldsymbol{\nu} = \mu_1 \nu_1 + \mu_2 \nu_2$.

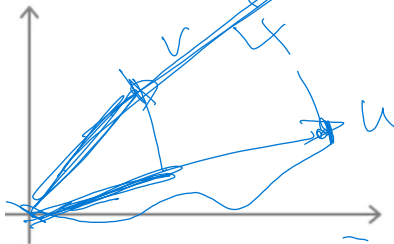
Mathematics behind large margin classification



Inner vector product:

- $\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$, $\boldsymbol{\nu} = \begin{bmatrix} \nu_1 \\ \nu_2 \end{bmatrix}$
- $\|\boldsymbol{\mu}\| = \sqrt{\mu_1^2 + \mu_2^2}$, the length of $\boldsymbol{\mu}$
- $\boldsymbol{\mu}^\top \boldsymbol{\nu} = \mu_1 \nu_1 + \mu_2 \nu_2$. How to represent it in the above plot?

Mathematics behind large margin classification



$$u^T v = p_1 \cdot \|u\| \\ = p_2 \cdot \|v\|$$

Inner vector product:

- $\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$, $\nu = \begin{bmatrix} \nu_1 \\ \nu_2 \end{bmatrix}$
- $\|\mu\| = \sqrt{\mu_1^2 + \mu_2^2}$, the length of μ
- $\mu^T \nu = \mu_1 \nu_1 + \mu_2 \nu_2$. How to represent it in the above plot?
- $\mu^T \nu = p \cdot \|\mu\|$, where p is the length of projection of ν on μ

Mathematics behind large margin classification



Inner vector product:

- $\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$, $\boldsymbol{\nu} = \begin{bmatrix} \nu_1 \\ \nu_2 \end{bmatrix}$
- $\|\boldsymbol{\mu}\| = \sqrt{\mu_1^2 + \mu_2^2}$, the length of $\boldsymbol{\mu}$
- $\boldsymbol{\mu}^\top \boldsymbol{\nu} = \mu_1 \nu_1 + \mu_2 \nu_2$. How to represent it in the above plot?
- $\boldsymbol{\mu}^\top \boldsymbol{\nu} = p \cdot \|\boldsymbol{\mu}\|$, where p is the length of projection of $\boldsymbol{\nu}$ on $\boldsymbol{\mu}$
- Note that if the angle between $\boldsymbol{\mu}$ and $\boldsymbol{\nu}$ is larger than 90° , then $p < 0$

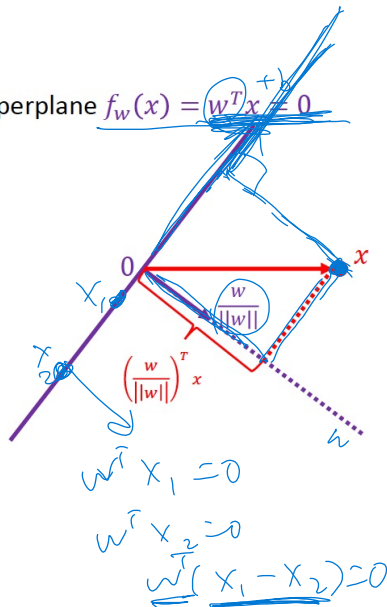
Mathematics behind large margin classification

- Lemma 1: x has distance $\frac{|f_w(x)|}{\|w\|}$ to the hyperplane $f_w(x) = w^T x = 0$

Proof:

- w is orthogonal to the hyperplane
- The unit direction is $\frac{w}{\|w\|}$
- The projection of x is $\left(\frac{w}{\|w\|}\right)^T x = \frac{f_w(x)}{\|w\|}$

$$\frac{w^T x}{\|w\|} = \text{p. } \|w\|$$



Mathematics behind large margin classification

- Claim 1: w is orthogonal to the hyperplane $f_{w,b}(x) = w^T x + b = 0$

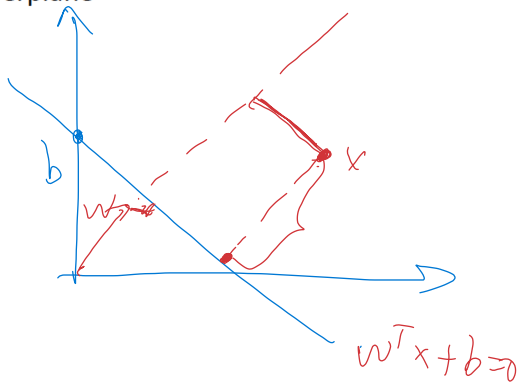
Proof:

- pick any x_1 and x_2 on the hyperplane

- $w^T x_1 + b = 0$

- $w^T x_2 + b = 0$

- So $w^T (x_1 - x_2) = 0$

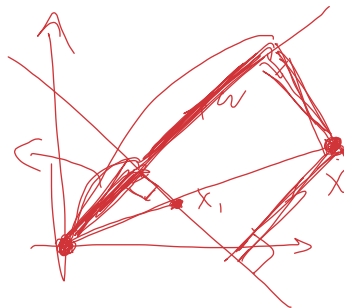


Mathematics behind large margin classification

- Claim 2: 0 has distance $\frac{-b}{||w||}$ to the hyperplane $w^T x + b = 0$

Proof:

- pick any x_1 the hyperplane
- Project x_1 to the unit direction $\frac{w}{||w||}$ to get the distance
- $\left(\frac{w}{||w||}\right)^T x_1 = \frac{-b}{||w||}$ since $w^T x_1 + b = 0$



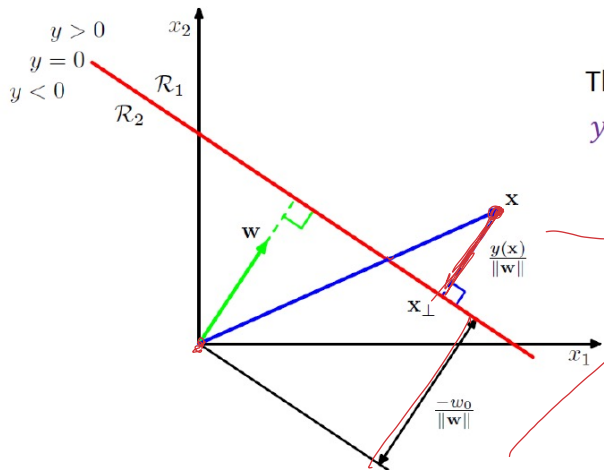
Mathematics behind large margin classification

- Lemma 2: x has distance $\frac{|f_{w,b}(x)|}{||w||}$ to the hyperplane $f_{w,b}(x) = w^T x + b = 0$

Proof:

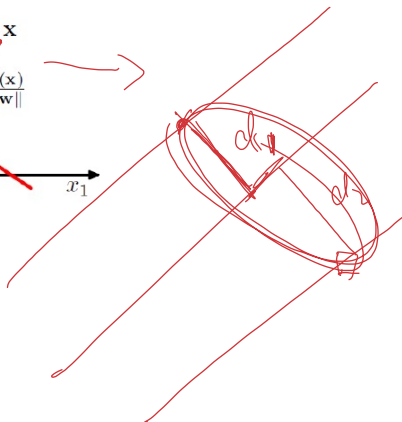
- Let $x = x_{\perp} + r \frac{w}{||w||}$, then $|r|$ is the distance
- Multiply both sides by w^T and add b
- Left hand side: $w^T x + b = f_{w,b}(x)$
- Right hand side: $w^T x_{\perp} + r \frac{w^T w}{||w||} + b = 0 + r ||w||$

Mathematics behind large margin classification



The notation here is:

$$y(x) = w^T x + w_0$$



Large margin classification

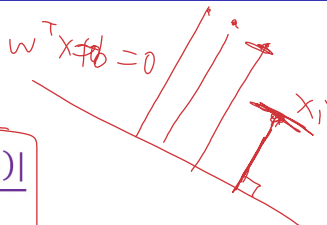
- Margin over all training data points:

$$\gamma = \min_i \frac{|f_{w,b}(x_i)|}{||w||}$$

- Since only want correct $f_{w,b}$, and recall $y_i \in \{+1, -1\}$, we have

$$\gamma = \min_i \frac{y_i f_{w,b}(x_i)}{||w||}$$

- If $f_{w,b}$ incorrect on some x_i , the margin is negative



Large margin classification

- Maximize margin over all training data points:

$$\max_{w,b} \gamma = \max_{w,b} \min_i \frac{\gamma_i f_{w,b}(x_i)}{\|w\|} = \max_{w,b} \min_i \frac{y_i (w^T x_i + b)}{\|w\|}$$

- A bit complicated ...

min max

$$\max_{w,b} \frac{1}{\min \|w\|^2}$$

Large margin classification

- Observation: when (w, b) scaled by a factor c , the margin unchanged

$$\frac{y_i(cw^T x_i + cb)}{\|cw\|} = \frac{y_i(w^T x_i + b)}{\|w\|}$$

$c \|w\|$

- Let's consider a fixed scale such that

$$y_{i^*}(w^T x_{i^*} + b) = 1$$

where x_{i^*} is the point closest to the hyperplane



Large margin classification

- Let's consider a fixed scale such that

$$y_{i^*}(w^T x_{i^*} + b) = 1$$

where x_{i^*} is the point closet to the hyperplane

- Now we have for all data

$$y_i(w^T x_i + b) \geq 1$$

and at least for one i the equality holds

- Then the margin is $\frac{1}{\|w\|}$

Large margin classification

- Optimization simplified to

$$\min_{w,b} \frac{1}{2} \|w\|^2$$
$$y_i(w^T x_i + b) \geq 1, \forall i$$

Alternative view of logistic regression

- Hypothesis function:

$$h_{\mathbf{w}}(\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = g(z)$$

where $z = \mathbf{w}^\top \mathbf{x}$

Alternative view of logistic regression

- Hypothesis function:

$$h_{\mathbf{w}}(\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = g(z)$$

where $z = \mathbf{w}^\top \mathbf{x}$

- If $y = 1$, we want $h_{\mathbf{w}}(\mathbf{x}) \approx 1$, *i.e.*, $\mathbf{w}^\top \mathbf{x} \gg 0$

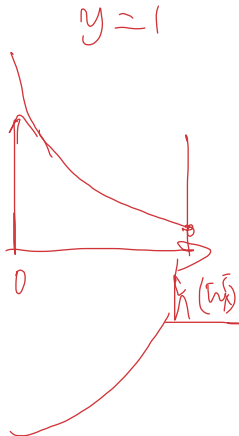
Alternative view of logistic regression

- Hypothesis function:

$$h_{\mathbf{w}}(\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = g(z)$$

where $z = \mathbf{w}^\top \mathbf{x}$

- If $y = 1$, we want $h_{\mathbf{w}}(\mathbf{x}) \approx 1$, i.e., $\mathbf{w}^\top \mathbf{x} \gg 0$
- If $y = -1$, we want $h_{\mathbf{w}}(\mathbf{x}) \approx 0$, i.e., $\mathbf{w}^\top \mathbf{x} \ll 0$



Alternative view of logistic regression

- Hypothesis function:

$$h_{\mathbf{w}}(\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = g(z)$$

where $z = \mathbf{w}^\top \mathbf{x}$

- If $y = 1$, we want $h_{\mathbf{w}}(\mathbf{x}) \approx 1$, i.e., $\mathbf{w}^\top \mathbf{x} \gg 0$
- If $y = -1$, we want $h_{\mathbf{w}}(\mathbf{x}) \approx 0$, i.e., $\mathbf{w}^\top \mathbf{x} \ll 0$
- Objective function of logistic regression

$$J(\mathbf{w}) = -\delta_{y=1} \log(h_{\mathbf{w}}(\mathbf{x})) - \delta_{y=-1} \log(1 - h_{\mathbf{w}}(\mathbf{x})), \quad (1)$$

where $\delta_a = 1$ if a is true, otherwise 0.

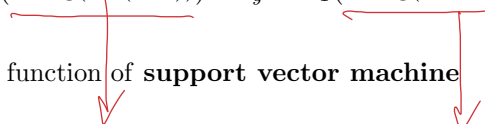
Objective of SVM

- Objective function of logistic regression

$$\frac{1}{m} \sum_i^m \left[\delta_{y^{(i)}=1} \left(-\log(h_{\mathbf{w}}(\mathbf{x}^{(i)})) \right) + \delta_{y^{(i)}=-1} \left(-\log(1 - h_{\mathbf{w}}(\mathbf{x}^{(i)})) \right) \right] + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

Objective of SVM

- Objective function of logistic regression

$$\frac{1}{m} \sum_i^m \left[\delta_{y^{(i)}=1} (- \log(h_{\mathbf{w}}(\mathbf{x}^{(i)}))) + \delta_{y^{(i)}=-1} (- \log(1 - h_{\mathbf{w}}(\mathbf{x}^{(i)}))) \right] + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$


- Objective function of **support vector machine**

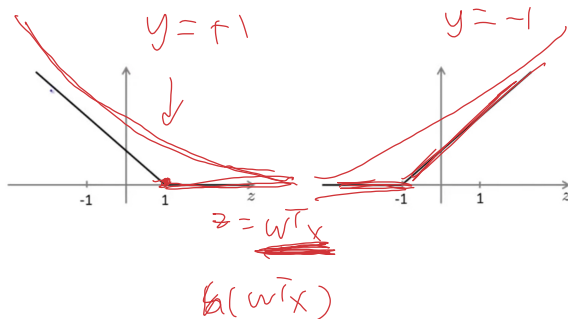
$$\frac{1}{m} \sum_i^m \left[\delta_{y^{(i)}=1} \text{cost}_1(\mathbf{w}^\top \mathbf{x}^{(i)} + b) + \delta_{y^{(i)}=-1} \text{cost}_{-1}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) \right] + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$
$$\equiv C \sum_i^m \left[\delta_{y^{(i)}=1} \text{cost}_1(\mathbf{w}^\top \mathbf{x}^{(i)} + b) + \delta_{y^{(i)}=-1} \text{cost}_{-1}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) \right] + \frac{1}{2} \sum_{j=1}^n w_j^2$$

Objective of SVM

Objective of SVM

- Objective function of support vector machine

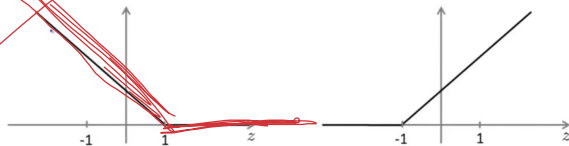
$$C \sum_i^m \left[\delta_{y^{(i)}=1} \text{cost}_1(\mathbf{w}^\top \mathbf{x}^{(i)} + b) + \delta_{y^{(i)}=-1} \text{cost}_{-1}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) \right] + \frac{1}{2} \sum_{j=1}^n w_j^2$$



Objective of SVM

- Objective function of support vector machine

$$C \sum_i^m \left[\delta_{y^{(i)}=1} \underbrace{\text{cost}_1(\mathbf{w}^\top \mathbf{x}^{(i)} + b)}_{y \approx 1} + \delta_{y^{(i)}=-1} \text{cost}_{-1}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) \right] + \frac{1}{2} \sum_{j=1}^n w_j^2$$



- If $y = +1$, we require that $\mathbf{w}^\top \mathbf{x}^{(i)} + b \geq 1$. In other words, $\text{cost}_1(\mathbf{w}^\top \mathbf{x}^{(i)} + b) = 0$ if $\mathbf{w}^\top \mathbf{x}^{(i)} + b \geq 1$

Objective of SVM

- Objective function of support vector machine

$$C \sum_i^m \left[\delta_{y^{(i)=1}} \underbrace{\text{cost}_1(\mathbf{w}^\top \mathbf{x}^{(i)} + b)} + \delta_{y^{(i)=-1}} \text{cost}_{-1}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) \right] + \frac{1}{2} \sum_{j=1}^n w_j^2$$



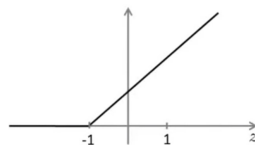
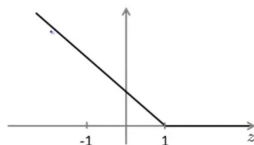
- If $y = +1$, we require that $\mathbf{w}^\top \mathbf{x}^{(i)} + b \geq 1$. In other words, $\text{cost}_1(\mathbf{w}^\top \mathbf{x}^{(i)} + b) = 0$ if $\mathbf{w}^\top \mathbf{x}^{(i)} + b \geq 1$
- If $y = -1$, we require that $\mathbf{w}^\top \mathbf{x}^{(i)} + b \leq -1$. In other words, $\text{cost}_{-1}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) = 0$ if $\mathbf{w}^\top \mathbf{x}^{(i)} + b \leq -1$

$$\max(0, 1 - y \cdot z)$$

Objective of SVM

- Objective function of support vector machine

$$C \sum_i^m \left[\underbrace{\delta_{y^{(i)=1} \text{cost}_1(\mathbf{w}^\top \mathbf{x}^{(i)} + b)} + \delta_{y^{(i)=-1} \text{cost}_{-1}(\mathbf{w}^\top \mathbf{x}^{(i)} + b)} \right] + \frac{1}{2} \sum_{j=1}^n w_j^2$$



- If $y = +1$, we require that $\mathbf{w}^\top \mathbf{x}^{(i)} + b \geq 1$. In other words, $\text{cost}_1(\mathbf{w}^\top \mathbf{x}^{(i)} + b) = 0$ if $\mathbf{w}^\top \mathbf{x}^{(i)} + b \geq 1$
- If $y = -1$, we require that $\mathbf{w}^\top \mathbf{x}^{(i)} + b \leq -1$. In other words, $\text{cost}_{-1}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) = 0$ if $\mathbf{w}^\top \mathbf{x}^{(i)} + b \leq -1$
- Hinge loss:**

$$\max(0, 1 - y(\mathbf{w}^\top \mathbf{x}^{(i)} + b)) \quad (2)$$

Mathematics behind large margin classification

- However, hinge loss is non-smooth. We transform the objective function of support vector machine to the following

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \sum_{j=1}^n w_j^2 \\ \text{s.t.} \quad & \mathbf{w}^\top \mathbf{x}^{(i)} + b \geq 1, \text{ if } y^{(i)} = 1; \quad \mathbf{w}^\top \mathbf{x}^{(i)} + b < -1, \text{ if } y^{(i)} = -1. \end{aligned} \tag{3}$$

Mathematics behind large margin classification

- However, hinge loss is non-smooth. We transform the objective function of support vector machine to the following

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \sum_{j=1}^n w_j^2 \\ \text{s.t.} \quad & \mathbf{w}^\top \mathbf{x}^{(i)} + b \geq 1, \text{ if } y^{(i)} = 1; \quad \mathbf{w}^\top \mathbf{x}^{(i)} + b < -1, \text{ if } y^{(i)} = -1. \end{aligned} \tag{3}$$

- It can be simplified as follows

$$\left\{ \begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & y^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) \geq 1, \forall i \end{aligned} \right. \tag{4}$$

Mathematics behind large margin classification

- However, hinge loss is non-smooth. We transform the objective function of support vector machine to the following

$$\min_{\mathbf{w}} \frac{1}{2} \sum_{j=1}^n w_j^2 \quad (3)$$

s.t. $\mathbf{w}^\top \mathbf{x}^{(i)} + b \geq 1$, if $y^{(i)} = 1$; $\mathbf{w}^\top \mathbf{x}^{(i)} + b < -1$, if $y^{(i)} = -1$.

- It can be simplified as follows

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 \quad (4)$$

s.t. $y^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) \geq 1, \forall i$

- Utilizing $p = \frac{\mathbf{w}^\top \mathbf{x} + b}{\|\mathbf{w}\|}$, we have

$$\mathbf{w}^\top \mathbf{x}^{(i)} + b = p^{(i)} \cdot \|\mathbf{w}\| \quad (5)$$

Mathematics behind large margin classification

- The objective function of support vector machine is transformed to

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & y^{(i)} \cdot p^{(i)} \cdot \|\mathbf{w}\| \geq 1, \forall i \end{aligned} \tag{6}$$

where $p^{(i)}$ indicates the projection length of $\mathbf{x}^{(i)}$ on \mathbf{w} .

Mathematics behind large margin classification

- The objective function of support vector machine is transformed to

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & y^{(i)} \cdot \underbrace{p^{(i)}}_{\text{projection length}} \cdot \underbrace{\|\mathbf{w}\|}_{\text{norm}} \geq 1, \forall i \end{aligned} \tag{6}$$

where $p^{(i)}$ indicates the projection length of $\mathbf{x}^{(i)}$ on \mathbf{w} .

- Let's see the following two decision boundaries (plot below)

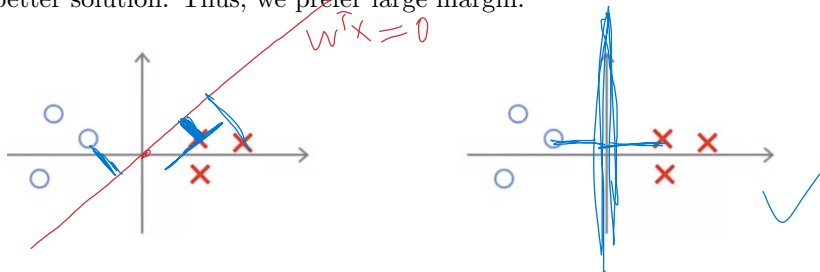
Mathematics behind large margin classification

- The objective function of support vector machine is transformed to

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & y^{(i)} \cdot p^{(i)} \cdot \|\mathbf{w}\| \geq 1, \forall i \end{aligned} \quad (6)$$

where $p^{(i)}$ indicates the projection length of $\mathbf{x}^{(i)}$ on \mathbf{w} .

- Let's see the following two decision boundaries (plot below)
- If the projection length p is larger, then $\|\mathbf{w}\|$ could be smaller, leading to better solution. Thus, we prefer large margin.



Reading material

Reading materials:

- Andrew Ng's note on SVM:
<https://see.stanford.edu/materials/aimlcs229/cs229-notes3.pdf>
- Chapter 7.1 of Bishop's book