

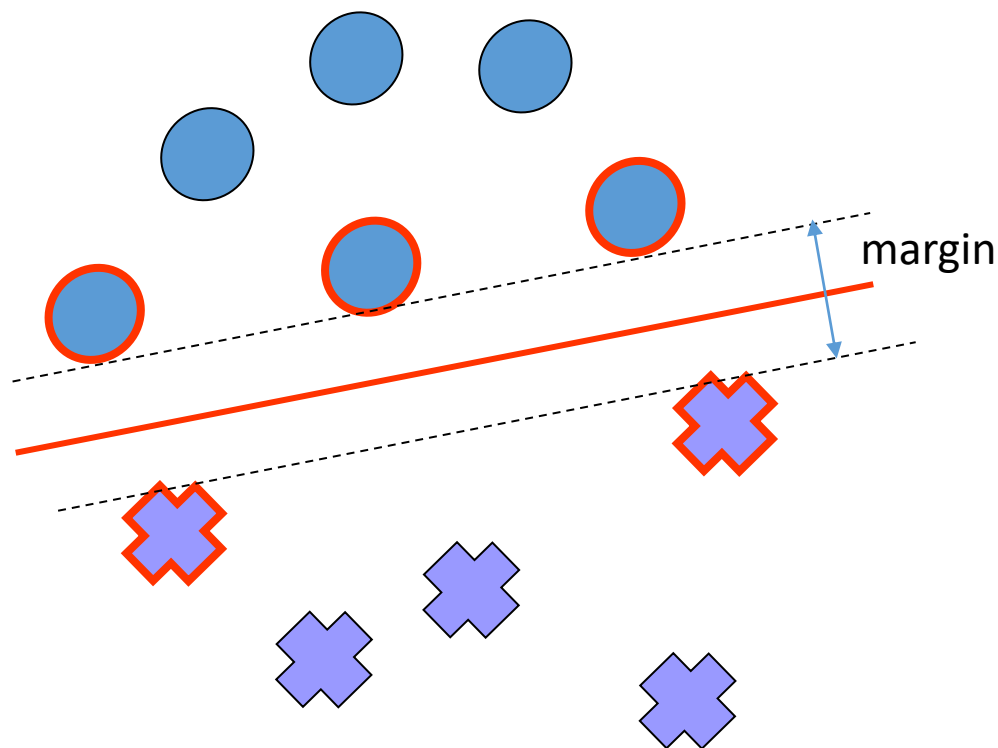
Support Vector Machine

SVM

Masaomi Kimura

SVM (Support vector machines)

- the method to find a border plane to classify two kinds of data



First order equations

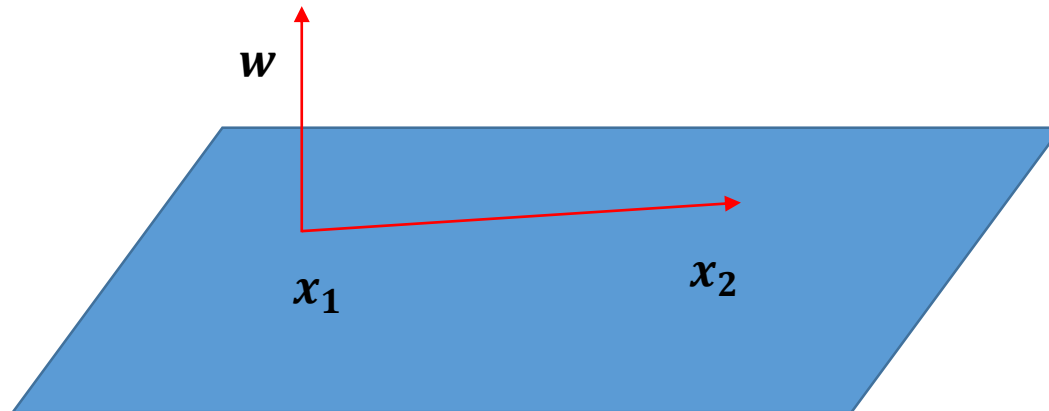
- First order equations express flat things!
 - Line: $2x + 3y = 1$
 - Plane: $2x + 3y + z = 1$
 - Hyper plane: $2x + 3y + z - w = 1$
- Using vectors and their inner product, we can express them easily:

$$\mathbf{w} \cdot \mathbf{x} - 1 (= 2x + 3y - 1) = 0,$$

$$\text{where } \mathbf{w} = \begin{pmatrix} 2 \\ 3 \end{pmatrix} \text{ and } \mathbf{x} = \begin{pmatrix} x \\ y \end{pmatrix}$$

Notice

- Vector \mathbf{w} is perpendicular
to the line/plane $\mathbf{w} \cdot \mathbf{x} + b = 0$



* If $\mathbf{w} \cdot \mathbf{a} = 0$, then \mathbf{w} and \mathbf{a} are perpendicular to each other.

Problem setting in SVM

- Classify data in two groups in advance, and assign labels $\{y_i | y_i = \pm 1\}$ to each of data

$$\mathbf{x}_i = \left(x_1^{(i)} \quad x_2^{(i)} \quad \dots \quad x_n^{(i)} \right)^T$$

- The data need to be linearly separable
 - SVM can be extended to allow margins or the cases with separable by other shapes

Inequality constraint in SVM

Points having $y_i = +1$ are expressed as \mathbf{x}_i^+

Points having $y_i = -1$ are expressed as \mathbf{x}_i^-

It is possible to find a small positive number ϵ to satisfy

$$\mathbf{w} \cdot \mathbf{x}_i^+ + b \geq \epsilon \quad \text{for } y_i = +1$$

$$\mathbf{w} \cdot \mathbf{x}_i^- + b \leq -\epsilon \quad \text{for } y_i = -1$$

Because rescaling coefficient does not change the plane, if we rescale $\mathbf{w} \rightarrow \epsilon \mathbf{w}$ and $b \rightarrow \epsilon b$, we have

$$\mathbf{w} \cdot \mathbf{x}_i^+ + b \geq 1 \quad \text{for } y_i = +1$$

$$\mathbf{w} \cdot \mathbf{x}_i^- + b \leq -1 \quad \text{for } y_i = -1$$



$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$$

Exercise

- Let us consider the case we have four points and a line:

$$\mathbf{x}_1^+ = (0 \quad 0.5)^T \quad \text{for } y_1 = +1$$

$$\mathbf{x}_2^+ = (3 \quad 1)^T \quad \text{for } y_2 = +1$$

$$\mathbf{x}_3^- = (-1 \quad 0.5)^T \quad \text{for } y_3 = -1$$

$$\mathbf{x}_4^- = (-2 \quad -4)^T \quad \text{for } y_4 = -1$$

$$2x + y + 1 = 0$$

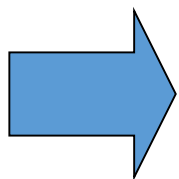
- Let's find a vector \mathbf{w} and a scalar b satisfying

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$$

SVM

- Find a hyperplane whose margins μ from both \mathbf{x}_i^+ and \mathbf{x}_i^- are maximized:

$$\begin{aligned}\mu &= \frac{|\mathbf{w} \cdot \mathbf{x}_i^+ + b|}{\|\mathbf{w}\|} + \frac{|\mathbf{w} \cdot \mathbf{x}_i^- + b|}{\|\mathbf{w}\|} \\ &= \frac{\mathbf{w} \cdot \mathbf{x}_i^+ + b}{\|\mathbf{w}\|} + \frac{-(\mathbf{w} \cdot \mathbf{x}_i^- + b)}{\|\mathbf{w}\|} \\ &\geq \frac{1}{\|\mathbf{w}\|} + \frac{-(-1)}{\|\mathbf{w}\|} = \frac{2}{\|\mathbf{w}\|} \quad (\text{should be maximized})\end{aligned}$$



Minimize $\frac{1}{2} \|\mathbf{w}\|^2$ subject to

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$$

Let's consider the following optimization problem...

- Consider the problem we need to find (x^*, y^*) :

$$(x^*, y^*) = \arg \min_{x, y} f(x, y)$$

under the condition

$$g(x, y) \geq 0.$$

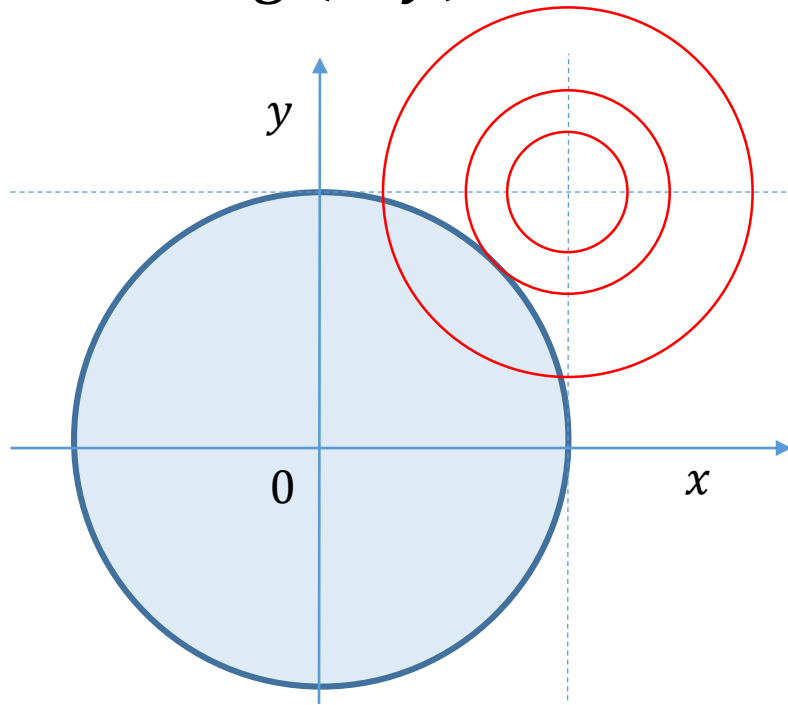
Example

- Consider a problem to minimize

$$f(x, y) = (x - 1)^2 + (y - 1)^2 = k$$

where x and y satisfy

$$g(x, y) = 1 - x^2 - y^2 \geq 0$$



Which circle has the least k ?

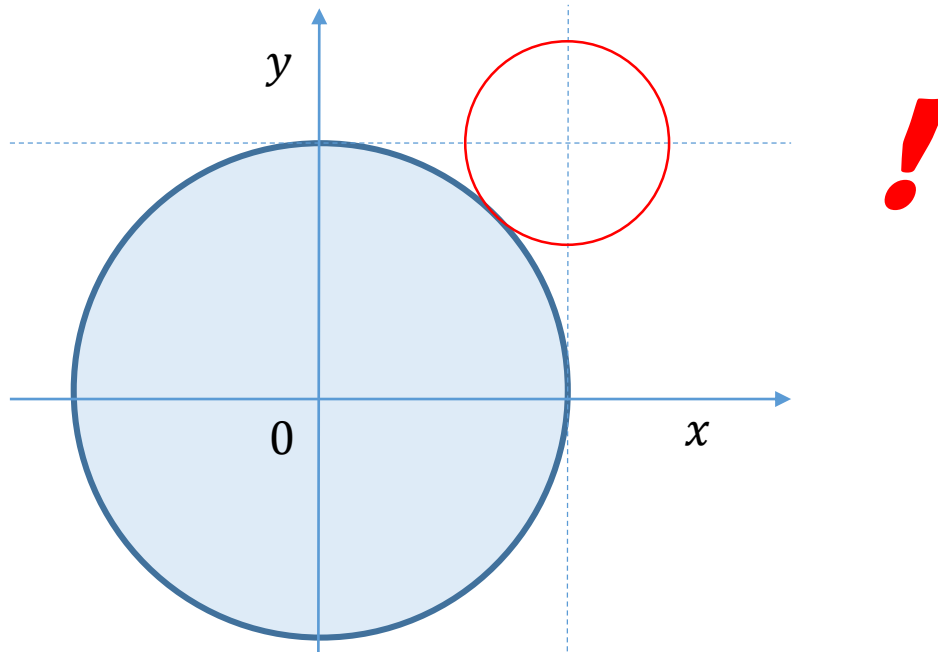


Answer

- The circle tangent to the border

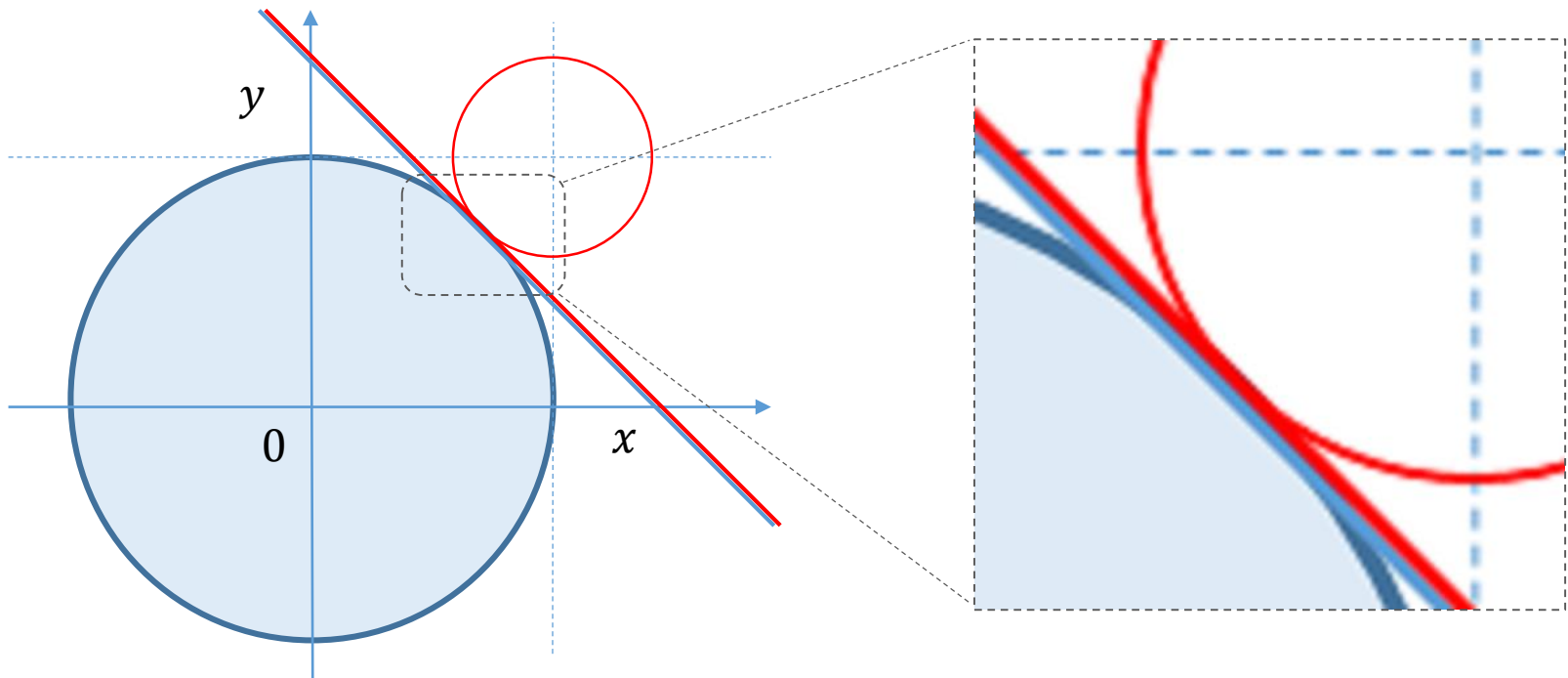
$$g(x, y) = 1 - x^2 - y^2 = 0$$

has the least k



The important thing is...

- The points where two objects (circle-line/circle-circle) were tangent to each other were important.
- At those points, they share the same tangent lines.



Tangent lines

- Around the tangent point (x_0, y_0) ,

$$g(x_0, y_0) = 0.$$

- These two tangent lines needs to express the same line:

$$\left(\frac{\partial f}{\partial x}(x_0, y_0)\right)(x - x_0) + \left(\frac{\partial f}{\partial y}(x_0, y_0)\right)(y - y_0) = 0$$

$$\left(\frac{\partial g}{\partial x}(x_0, y_0)\right)(x - x_0) + \left(\frac{\partial g}{\partial y}(x_0, y_0)\right)(y - y_0) = 0$$

- Since there is a redundancy of coefficients in an equation to express a line,

$$\frac{\partial f}{\partial x}(x_0, y_0) = \lambda \frac{\partial g}{\partial x}(x_0, y_0)$$

$$\frac{\partial f}{\partial y}(x_0, y_0) = \lambda \frac{\partial g}{\partial y}(x_0, y_0)$$

The constraint of λ

- Let the region G^+ be a set of points satisfying the condition $g(x, y) \geq 0$.
- If a point $(x', y') \in G^+$ is close to the tangent point (x_0, y_0) ,
$$g(x', y') - g(x_0, y_0) \cong \left(\frac{\partial g}{\partial x}(x_0, y_0) \right) (x - x_0) + \left(\frac{\partial g}{\partial y}(x_0, y_0) \right) (y - y_0) \geq 0,$$
since $g(x', y') \geq 0$ and $g(x_0, y_0) = 0$.
- Since $f(x_0, y_0)$ is less than any $f(x', y')$,
$$f(x', y') - f(x_0, y_0) \cong \left(\frac{\partial f}{\partial x}(x_0, y_0) \right) (x - x_0) + \left(\frac{\partial f}{\partial y}(x_0, y_0) \right) (y - y_0) \geq 0.$$

This leads λ to be positive.

In summary...

- To solve a problem to minimize $f(x, y)$ where x and y satisfy $g(x, y) \geq 0$, we should solve

$$\frac{\partial f}{\partial x}(x, y) - \lambda \frac{\partial g}{\partial x}(x, y) = 0$$

$$\frac{\partial f}{\partial y}(x, y) - \lambda \frac{\partial g}{\partial y}(x, y) = 0$$

$$\lambda > 0$$

$$g(x, y) = 0.$$

- Notice that the point (x, y) minimizing $f(x, y)$ without conditions sits outside G^+ .

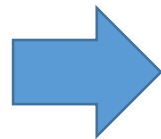
Another case we need to consider

- Consider the case that the point (x, y) minimizing $f(x, y)$ without conditions sits inside G^+ .
- It's super easy!! We can use a USUAL differential technique.

$$\frac{\partial f}{\partial x}(x, y) = 0$$

$$\frac{\partial f}{\partial y}(x, y) = 0$$

$$g(x, y) > 0$$



$$\frac{\partial f}{\partial x}(x, y) - \lambda \frac{\partial g}{\partial x}(x, y) = 0$$

$$\frac{\partial f}{\partial y}(x, y) - \lambda \frac{\partial g}{\partial y}(x, y) = 0$$

$$\lambda = 0$$

$$g(x, y) > 0$$

Kuhn-Tucker conditions

- In the case we need to find (x^*, y^*) :

$$(x^*, y^*) = \arg \min_{x, y} f(x, y)$$

under the condition $g(x, y) \geq 0$, we find it by solving

$$\frac{\partial f}{\partial x}(x^*, y^*) - \lambda \frac{\partial g}{\partial x}(x^*, y^*) = 0$$

$$\frac{\partial f}{\partial y}(x^*, y^*) - \lambda \frac{\partial g}{\partial y}(x^*, y^*) = 0$$

$$\lambda \geq 0$$

$$\lambda \cdot g(x^*, y^*) = 0.$$

Or...

- Define a function $L(x, y, \lambda) = f(x, y) - \lambda g(x, y)$, then the equations to be solved are:

$$\frac{\partial L}{\partial x}(x^*, y^*, \lambda^*) = 0$$

$$\frac{\partial L}{\partial y}(x^*, y^*, \lambda^*) = 0$$

$$\lambda^* \cdot \frac{\partial L}{\partial \lambda}(x^*, y^*, \lambda^*) = 0$$

$$\lambda^* \geq 0$$

More generally...

- Using a vector \mathbf{x} , define a function

$$L(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) - \sum_i \lambda_i g_i(\mathbf{x}),$$

then the equations to be solved are:

$$\frac{\partial L}{\partial \mathbf{x}}(\mathbf{x}^*, \boldsymbol{\lambda}^*) = 0$$

$$\lambda_i^* \cdot \frac{\partial L}{\partial \lambda_i}(\mathbf{x}^*, \boldsymbol{\lambda}^*) = 0$$

$$\lambda_i^* \geq 0$$

Then we can extend the discussion to a space whose dimension is more than two.

Dual theory

- To solve it,

$$\frac{\partial L}{\partial \mathbf{x}}(\mathbf{x}, \boldsymbol{\lambda}) = 0 \Rightarrow \mathbf{x} = \mathbf{h}(\boldsymbol{\lambda})$$

- If we consider the case $\lambda_i > 0$, we solve $g_i(\mathbf{h}(\boldsymbol{\lambda})) = 0$ to find $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$.
- Dual function is defined as

$$\theta(\boldsymbol{\lambda}) = L(\mathbf{h}(\boldsymbol{\lambda}), \boldsymbol{\lambda}) = f(\mathbf{h}(\boldsymbol{\lambda})) - \sum_i \lambda_i g_i(\mathbf{h}(\boldsymbol{\lambda}))$$

$$\frac{\partial \theta(\boldsymbol{\lambda})}{\partial \lambda_i} = \frac{\partial L}{\partial \mathbf{x}} \bigg|_{\mathbf{x}=\mathbf{h}} \frac{d\mathbf{h}}{d\lambda_i} - g_i(\mathbf{h}(\boldsymbol{\lambda})) = -g_i(\mathbf{h}(\boldsymbol{\lambda})) \leq \mathbf{0}$$

Dual theory (cont')

- Let λ^* be a solution of $\frac{\partial \theta(\lambda)}{\partial \lambda} = 0$. Then $\theta(\lambda^*)$ takes a maximal value:
 - $\theta(\lambda^*) = \theta(\lambda) + \frac{\partial \theta(\lambda)}{\partial \lambda_i} (\lambda_i^* - \lambda_i) \geq \theta(\lambda)$
- Since $f(\mathbf{h}(\lambda)) \geq f(\mathbf{h}(\lambda)) - \sum_i \lambda_i g_i(\mathbf{h}(\lambda)) = \theta(\lambda)$, $f(\mathbf{h}(\lambda^*)) = \theta(\lambda^*)$. The $f(\mathbf{h}(\lambda^*))$ takes a minimum value, because

$$\left. \frac{\partial L}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{h}(\lambda^*)} = \mathbf{0} \quad \text{and} \quad \left. \frac{\partial \theta(\lambda)}{\partial \lambda_i} \right|_{\lambda=\lambda^*} = -g_i(\mathbf{h}(\lambda^*)) = \mathbf{0}$$

Differential by a vector

- Let a vector denote $\mathbf{x} = (x_1 \quad x_2 \quad \dots \quad x_n)^T$
 - \mathbf{x}^T is a transposed vector (matrix) of \mathbf{x} .

- Definition:

$$\frac{\partial f}{\partial \mathbf{x}} = \left(\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \quad \dots \quad \frac{\partial f}{\partial x_n} \right)^T$$

- Examples

$$\begin{aligned} \frac{\partial}{\partial \mathbf{w}} (\mathbf{w} \cdot \mathbf{w}) &= \left(\frac{\partial}{\partial w_1} \sum_i w_i^2 \quad \frac{\partial}{\partial w_2} \sum_i w_i^2 \quad \dots \quad \frac{\partial}{\partial w_n} \sum_i w_i^2 \right)^T \\ &= (2w_1 \quad 2w_2 \quad \dots \quad 2w_n)^T = 2\mathbf{w} \end{aligned}$$

$$\frac{\partial}{\partial \mathbf{w}} (\mathbf{x} \cdot \mathbf{w}) = \mathbf{x}$$

Application of Kuhn-Tucker cond.

Define $L = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_i \alpha_i \{y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1\}$, then
KT condition is:

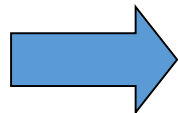
$$\begin{aligned}\frac{\partial L}{\partial \mathbf{w}} &= \mathbf{w} - \sum_i \alpha_i y_i \mathbf{x}_i = 0 \\ \frac{\partial L}{\partial b} &= - \sum_i \alpha_i y_i = 0 \\ \alpha_i &\geq 0 \\ \alpha_i \{y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1\} &= 0\end{aligned}$$

Dual expression of SVM

- Inserting $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$ into L , we get a dual expression based on Kuhn-Tucker condition, which should be maximized is:

$$\theta(\boldsymbol{\alpha}) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i y_i \alpha_j y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

$$\begin{aligned} \alpha_i &\geq 0 \\ \alpha_i \{y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1\} &= 0 \\ \sum_i \alpha_i y_i &= 0 \end{aligned}$$



$\alpha_i \neq 0$ requires \mathbf{x}_i to satisfy

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 = 0$$

(support vectors)

After obtaining a solution

- We can calculate coefficients of the hyperplane:

$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$$
$$b = -\frac{\min(\mathbf{w} \cdot \mathbf{x}_i^+) + \max(\mathbf{w} \cdot \mathbf{x}_i^-)}{2}$$

- A class for new data \mathbf{x} is given by the value

$$H(x) = \text{sign} \left(\sum_i \alpha_i y_i (\mathbf{x}_i \cdot \mathbf{x}) + b \right)$$

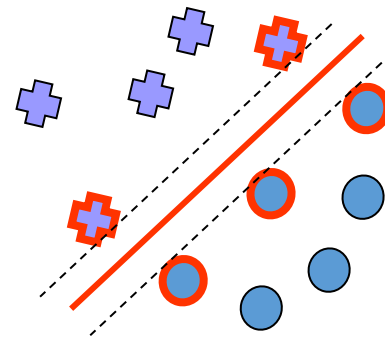
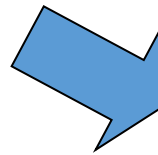
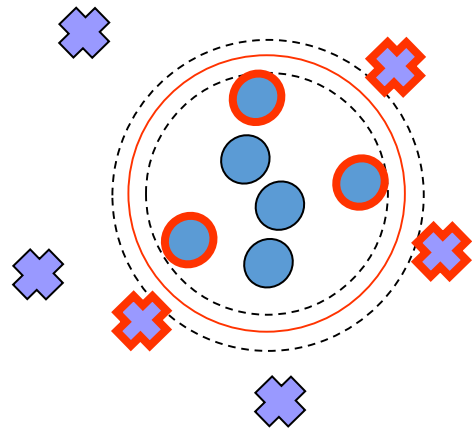
After obtaining a solution

- We can calculate coefficients of the hyperplane and get \mathbf{w} and b .
- A class for new data \mathbf{x} is given by the value

$$H(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$$

$$\text{sign}(x) = \begin{cases} 1 & (x \geq 0) \\ -1 & (x < 0) \end{cases}$$

Extension to a high dimensional space



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

$$(x-a)^2 + (y-b)^2 = r^2$$

$$\phi(x) = (X, Y, Z, V)$$

$$X + Y + aZ + bV + c = 0$$

$$(X = x^2, Y = y^2, Z = -2x, V = -2y, c = a^2 + b^2 - r^2)$$

Generalization

- If we can map \mathbf{x} to a higher dimensional space

$$\mathbf{X} = \boldsymbol{\phi}(\mathbf{x})$$

and data in two classes are separable by a hyper plane, the problem is to find the hyper plane:

$$\mathbf{W} \cdot \boldsymbol{\phi}(\mathbf{x}) + b = 0,$$

which maximizes margin between the two classes.

- By the detail discussion (see Appendix), \mathbf{W} is given as $\mathbf{W} = \sum_i \alpha_i y_i \boldsymbol{\phi}(\mathbf{x}_i)$.

$$H(\mathbf{x}) = \text{sign} \left(\sum_i \alpha_i y_i \boldsymbol{\phi}(\mathbf{x}_i) \cdot \boldsymbol{\phi}(\mathbf{x}) + b \right)$$

RBF kernel

- The properties of a kernel $K(\mathbf{x}, \mathbf{y}) = \boldsymbol{\phi}(\mathbf{x}) \cdot \boldsymbol{\phi}(\mathbf{y})$ needs to have are:
 - $K(\mathbf{x}, \mathbf{x}) \geq 0$
 - $K(\mathbf{x}, \mathbf{y}) = K(\mathbf{y}, \mathbf{x})$
- The most used kernel is RBF (radial based function) kernel defined as
$$K(\mathbf{x}, \mathbf{y}) = \exp(-\gamma \|\mathbf{x} - \mathbf{y}\|^2)$$

