



Practical Machine Learning

Lecture 5

Support Vector Machine (SVM) and Kernel trick

Dr. Suyong Eum



A question from the last class

□ How do we infer $p(Z|X)$ using $q(Z)$?

- Only data X are given
- $p(X,Z)$ is known

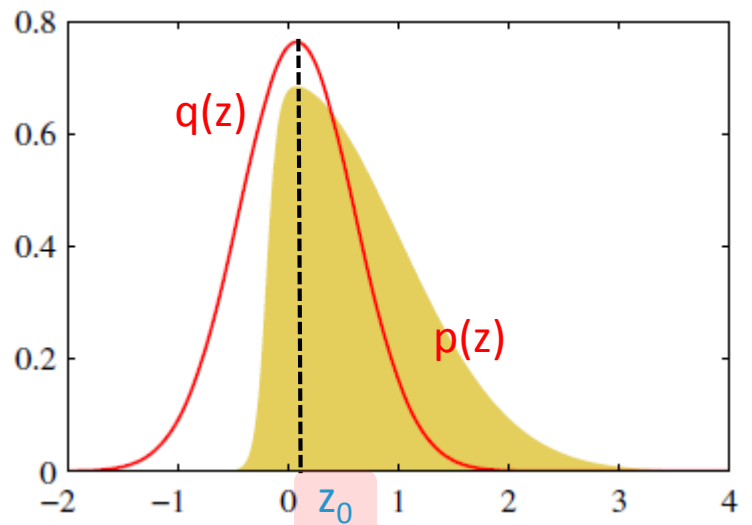
$$p(Z|X) = \frac{p(X,Z)}{p(X)} = \frac{p(X,Z)}{\sum_Z p(X,Z)}$$

- When the variable z is continuous
- When the number of variables z is many (?)

	Resource	Time	Accuracy
Laplace approach (Gaussian approximation)	Good	Good	Worse
Sampling (Numerical approach)	Worse	Worse	Good
Variational Inference (Analytical approach)	Medium	Medium	Medium

Laplace approach

- Finding the mode of the posterior distribution and then fitting a Gaussian centered at that mode.



Mode of $p(z)$: local maxima

$$p(z) \propto \exp(-z^2 / 2)(1 + \exp^{-20z-4})^{-1}$$

$$p(z) = \frac{1}{C} f(z) \quad C = \int f(z) dz \quad \text{Normalizing factor, which is unknown}$$

$$\left. \frac{df(z)}{dz} \right|_{z=z_0} = 0$$

It becomes mean of $q(z)$

$$A = - \left. \frac{d^2}{dz^2} \ln f(z) \right|_{z=z_0}$$

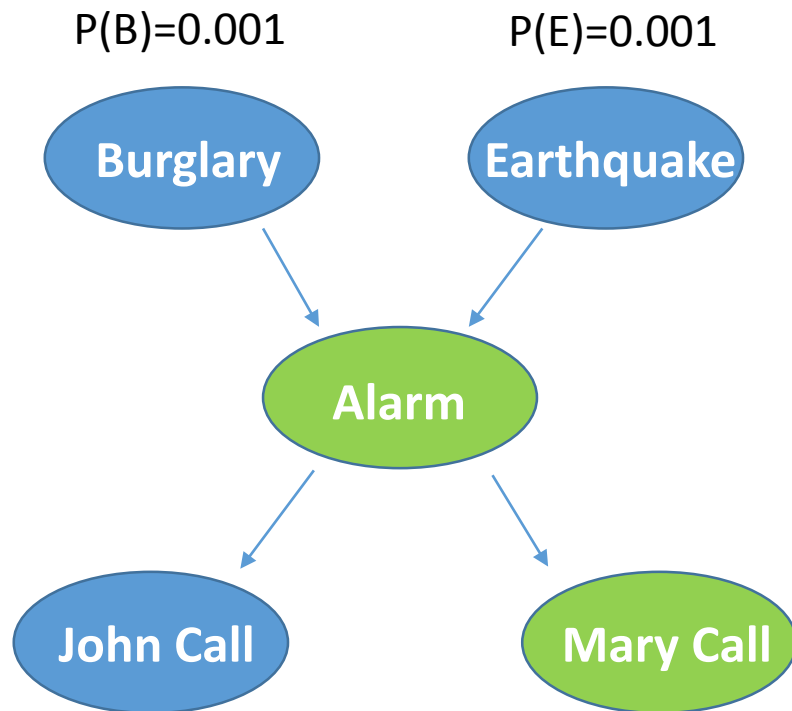
It becomes precision of $q(z)$

$$q(z) = \left(\frac{A}{2\pi} \right)^{1/2} \exp \left\{ -\frac{A}{2} (z - z_0)^2 \right\}$$

$$N(x | \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2} (x - \mu)^2 \right\}$$

$$\frac{d^2}{dx^2} \ln N(x | \mu, \sigma^2) = -\frac{1}{\sigma^2}$$

Sampling approach



$$p(B, E, A, J, M) = p(B)p(E)p(A | B, E)p(J | A)p(M | A)$$

$$p(B, E, J | A, M) \text{ ?}$$

$$p(E | B, A, J, M) = p(E | A, B)$$

$$p(B | E, A, J, M) = p(B | A, E)$$

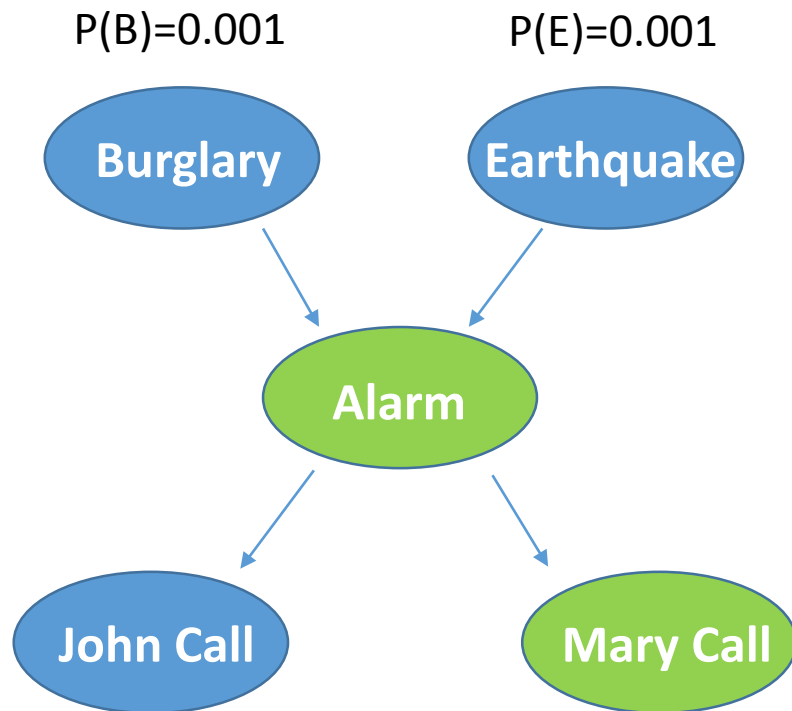
$$p(J | B, E, A, M) = p(J | A)$$

A	$P(J A)$
T	0.90
F	0.05

A	$P(M A)$
T	0.70
F	0.01

B	E	$P(A B,E)$
T	T	0.90
T	F	0.05
F	T	0.29
F	F	0.001

Sampling approach



$$p(B, E, A, J, M) = p(B)p(E)p(A | B, E)p(J | A)p(M | A)$$

$$p(B, E, J | A, M) \text{ ?}$$

$$p(E | B, A, J, M) = p(E | A, B)$$

$$p(E, A, B) = p(A | B, E)p(B)p(E)$$

$$p(E, A, B) = p(E | A, B)p(A, B)$$

$$= p(E | A, B)p(A | B)p(B)$$

$$p(A | B, E)p(B)p(E) = p(E | A, B)p(A | B)p(B)$$

$$p(E | A, B) = \frac{p(A | B, E)p(E)}{\sum_E p(A | B, E)}$$

A	P(J A)
T	0.90
F	0.05

A	P(M A)
T	0.70
F	0.01

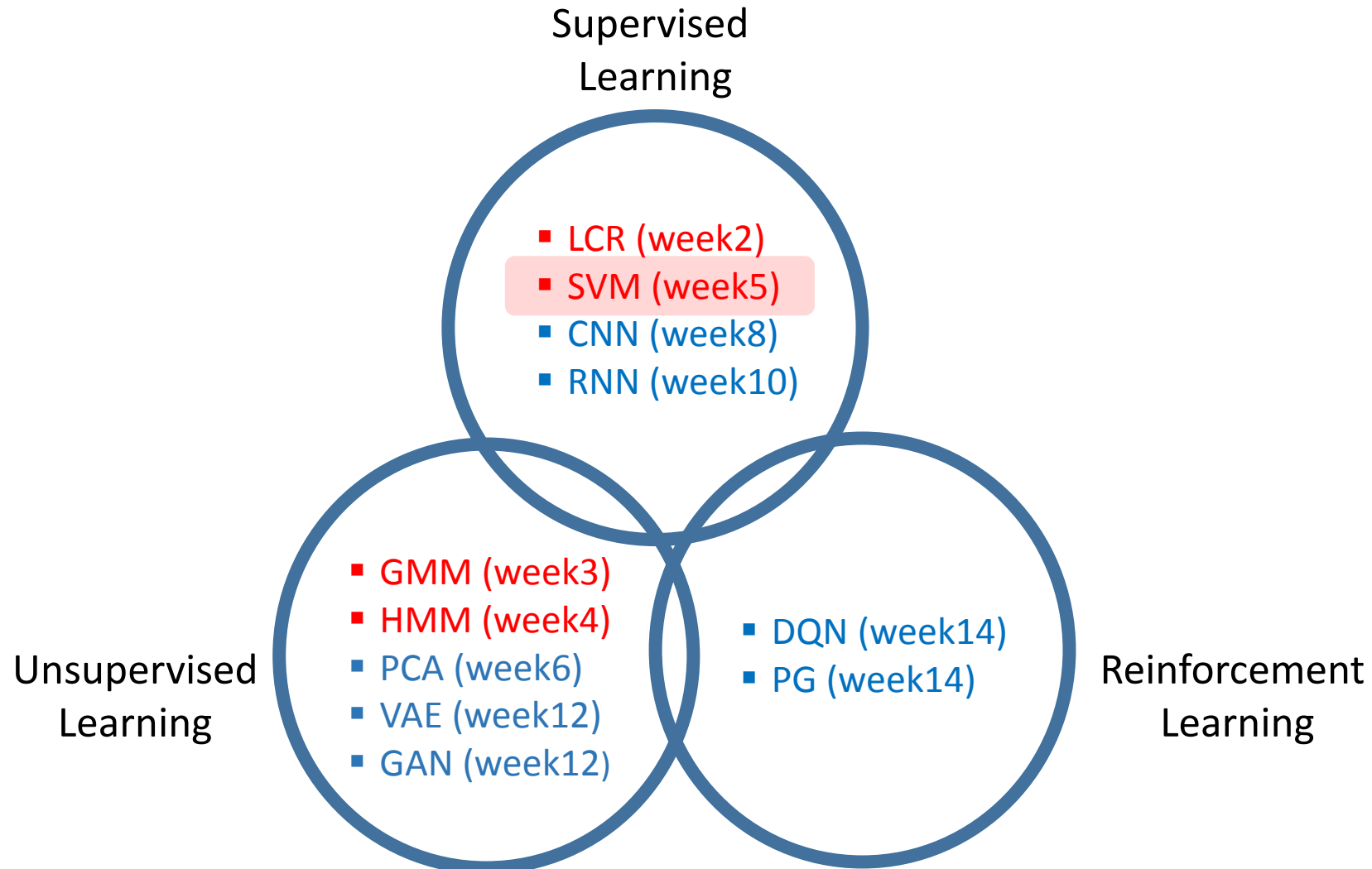
B	E	P(A B,E)
T	T	0.90
T	F	0.05
F	T	0.29
F	F	0.001

Variational inference

□ The idea is

- Finding $p(Z|X)$ by **minimizing Kullback divergence** to $q(Z)$
- Minimizing KL between $p(Z|X)$ and $q(Z)$ is equivalent to maximizing a function where the conditional distribution **$p(Z|X)$ is replaced with** the joint distribution **$p(Z, X)$.**
- **Factorizing the joint distribution** on the assumption that the latent variables Z are independent.
- Developing the derivation in terms of one latent variable on the assumption of **the other latent variables are known.**
- Then, do some algebra..

□ Refer the backup slides which include the derivation



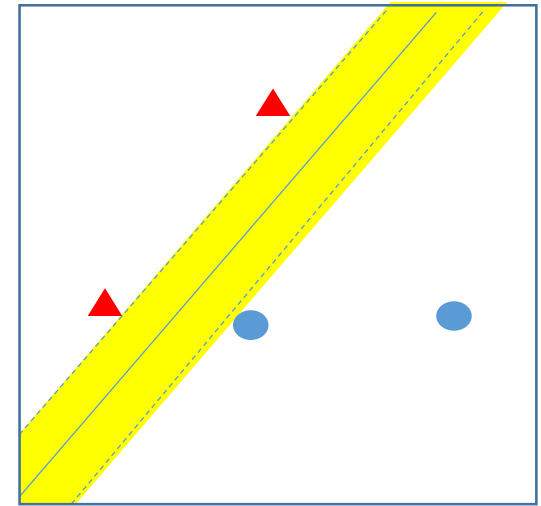
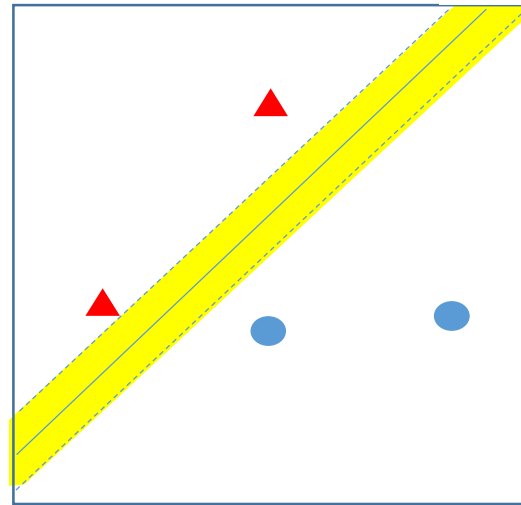
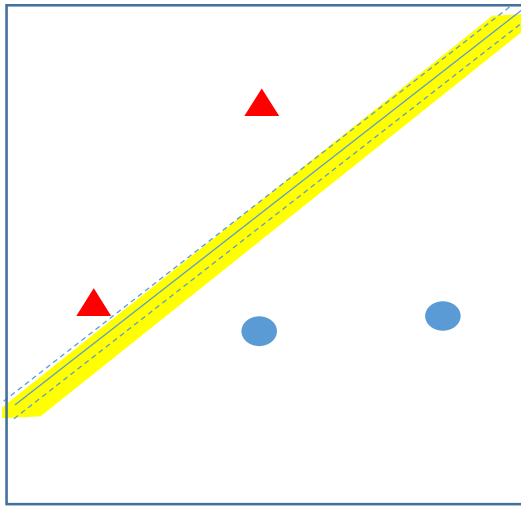
You are going to learn

- ❑ An idea of Support Vector Machine (SVM)
- ❑ Problem formulation of SVM
 - Linear classification: Hard Margin SVM
- ❑ Non-linear classification
 - Soft Margin SVM
 - Kernel trick

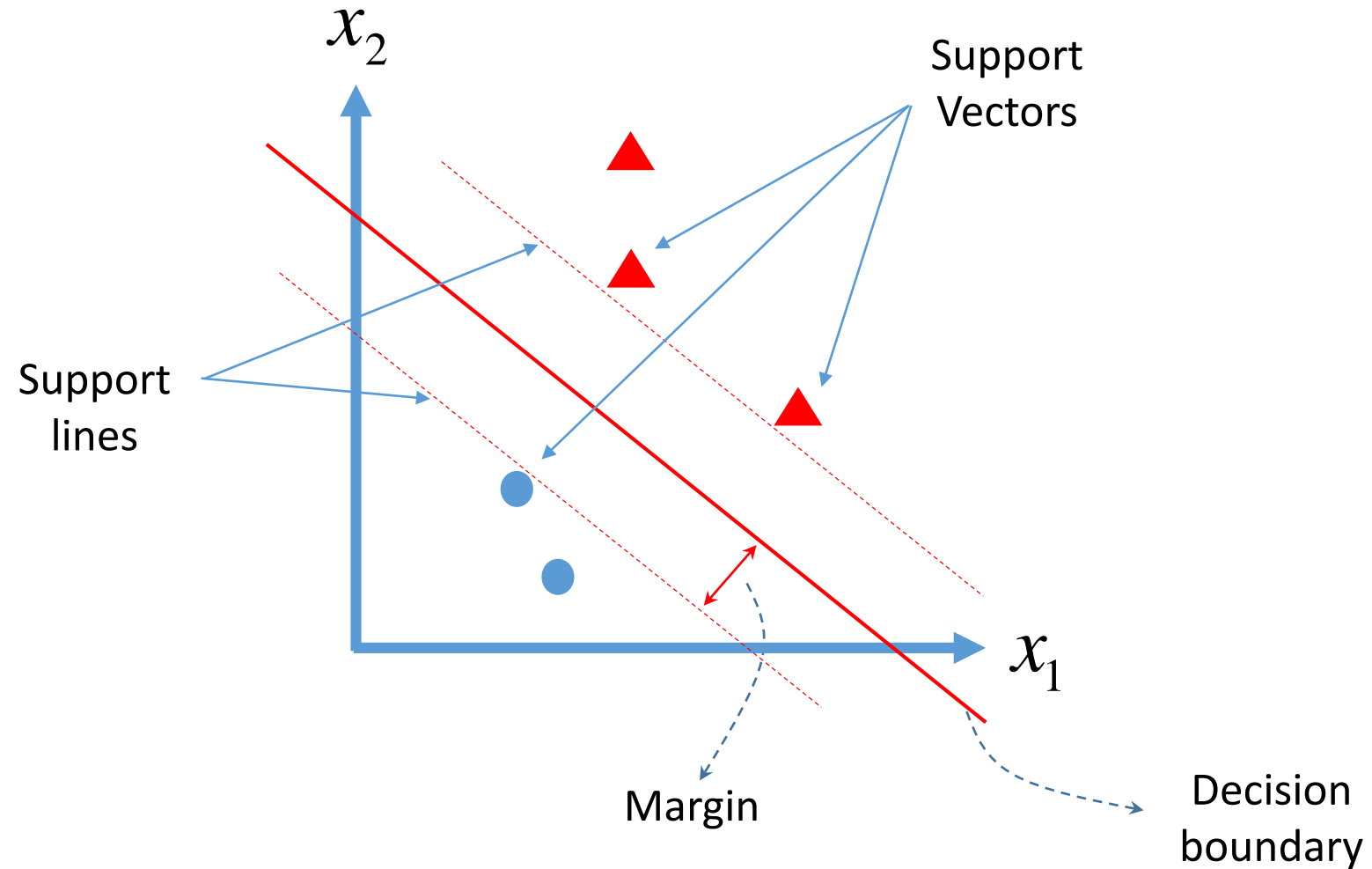
Why Support Vector Machine?

- ❑ Most widely used classification approach (practical)
 - Linearly separable data set
 - Linearly separable data set with a few violation
 - Non-linearly separable data set
- ❑ Supported by well defined mathematical theories
 - Geometry,
 - Optimization,
 - Quadratic programming,
 - Lagrange method,
 - Kernel, etc.

Which one is better for classification?



Terminology used in this lecture



Some geometry

$$y(\mathbf{x}) = w_2 x_2 + w_1 x_1 + w_0$$

$$y(\mathbf{x}^a) - y(\mathbf{x}^b) = w_2 x_2^a + w_1 x_1^a + w_0 - w_2 x_2^b - w_1 x_1^b - w_0$$

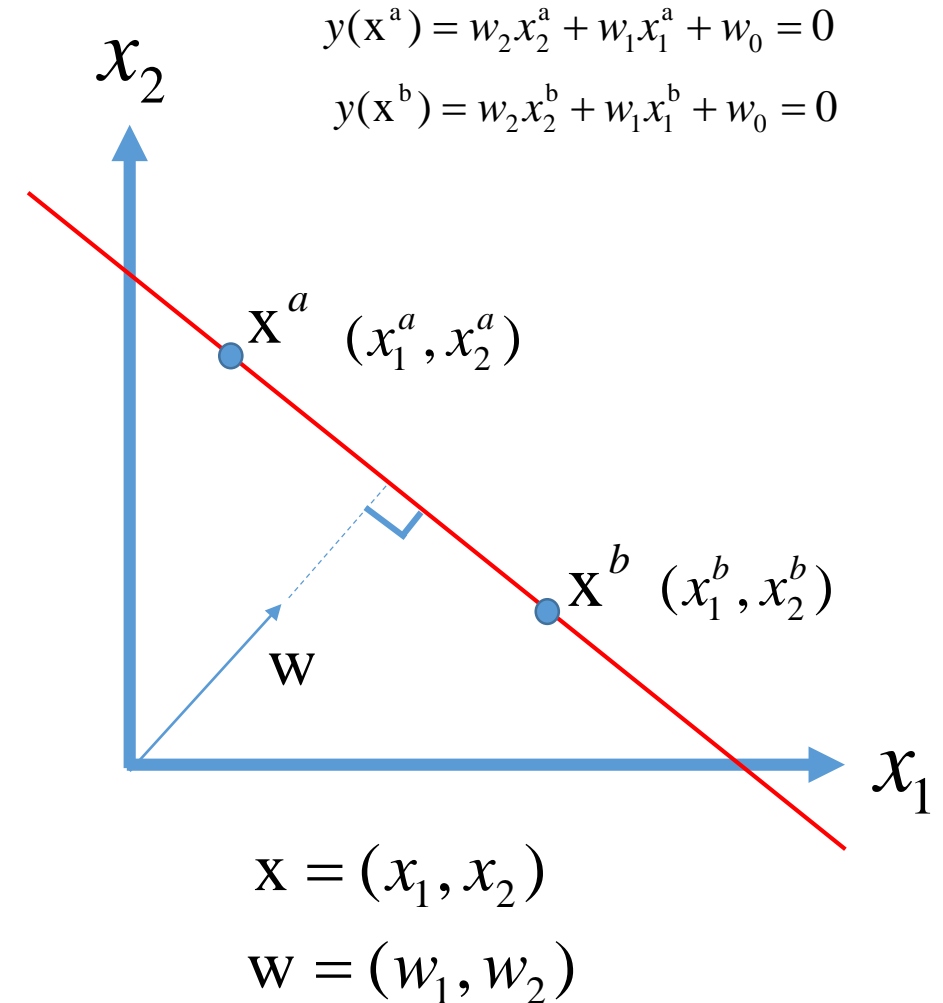
$$= w_2 (x_2^a - x_2^b) + w_1 (x_1^a - x_1^b)$$

$$= [w_1, w_2] \begin{bmatrix} x_1^a - x_1^b \\ x_2^a - x_2^b \end{bmatrix} \quad (1 \times 2)(2 \times 1) = (1 \times 1)$$

$$0 = \mathbf{w}^T (\mathbf{x}^a - \mathbf{x}^b)$$

$$\mathbf{w}^T \perp (\mathbf{x}^a - \mathbf{x}^b)$$

Vector on the decision boundary



□ Inner product

$$(x_1^a, x_2^a) \cdot (w_1, w_2) = \| (x_1^a, x_2^a) \| \| (w_1, w_2) \| \cos \theta$$

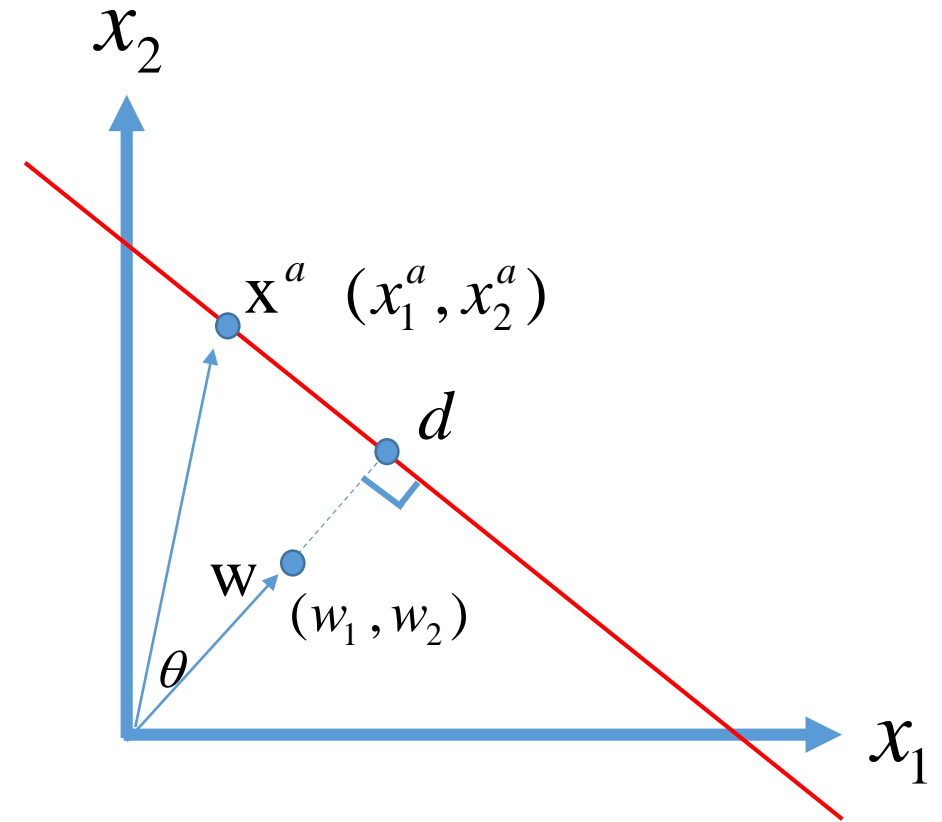
□ $\cos \theta$ definition

$$\cos \theta = \frac{\| d \|}{\| (x_1^a, x_2^b) \|} \Rightarrow \| d \| = \| (x_1^a, x_2^b) \| \cos \theta$$

$$\| (w_1, w_2) \| \| d \| = (x_1^a, x_2^b) \cdot (w_1, w_2)$$

$$\| d \| = \frac{w_2 x_2^a + w_1 x_1^a}{\| (w_1, w_2) \|} = \frac{-w_0}{\| (w_1, w_2) \|}$$

$$\| d \| = \frac{-w_0}{\| \mathbf{w} \|}$$



Margin distance

Size of the vector
($x^b \rightarrow x^c$)

$$x^c = x^b + \left\| r \right\| \frac{w}{\left\| w \right\|}$$

Unit vector showing
the direction only

□ Let's multiply w^T and add w_0 in both sides.

$$w^T x^c + w_0 = w^T x^b + w_0 + w^T \left\| r \right\| \frac{w}{\left\| w \right\|}$$

$$y(x^c) = w^T \left\| r \right\| \frac{w}{\left\| w \right\|}$$

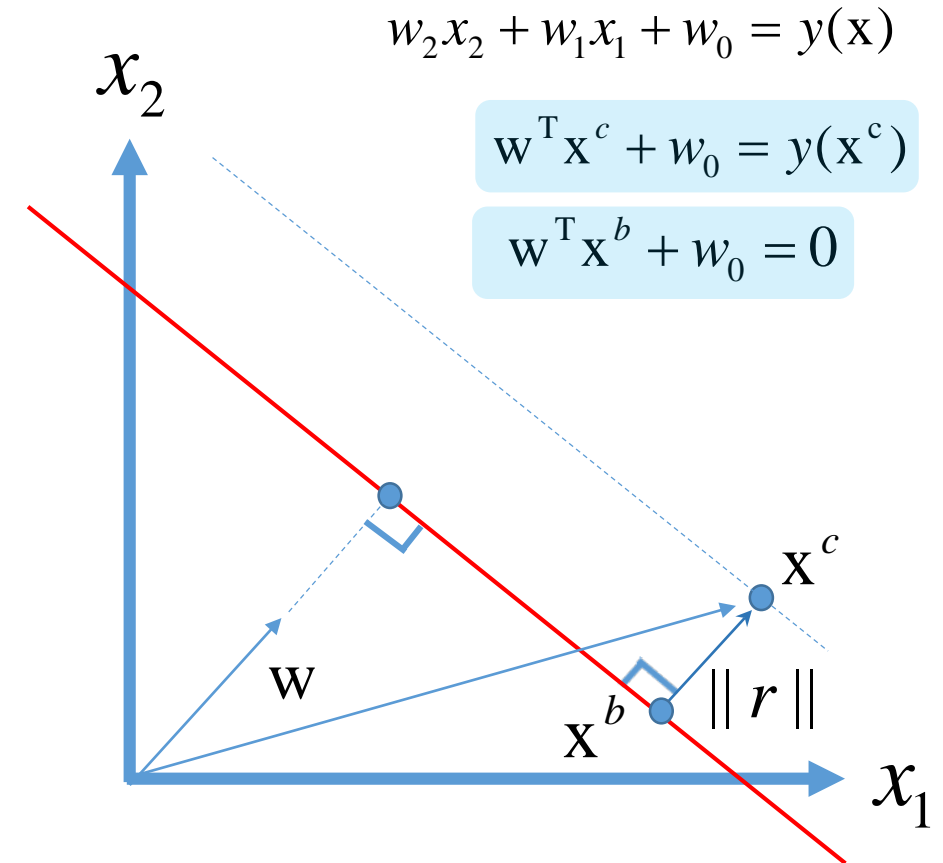
$$\left\| r \right\| = \frac{y(x^c)}{\left\| w \right\|}$$

$$\left\| r \right\| = \frac{1}{\left\| w \right\|}$$

Let's say

$$|y(x^c)| = 1$$

To make the margin become one... ([see later](#))



Problem formulation

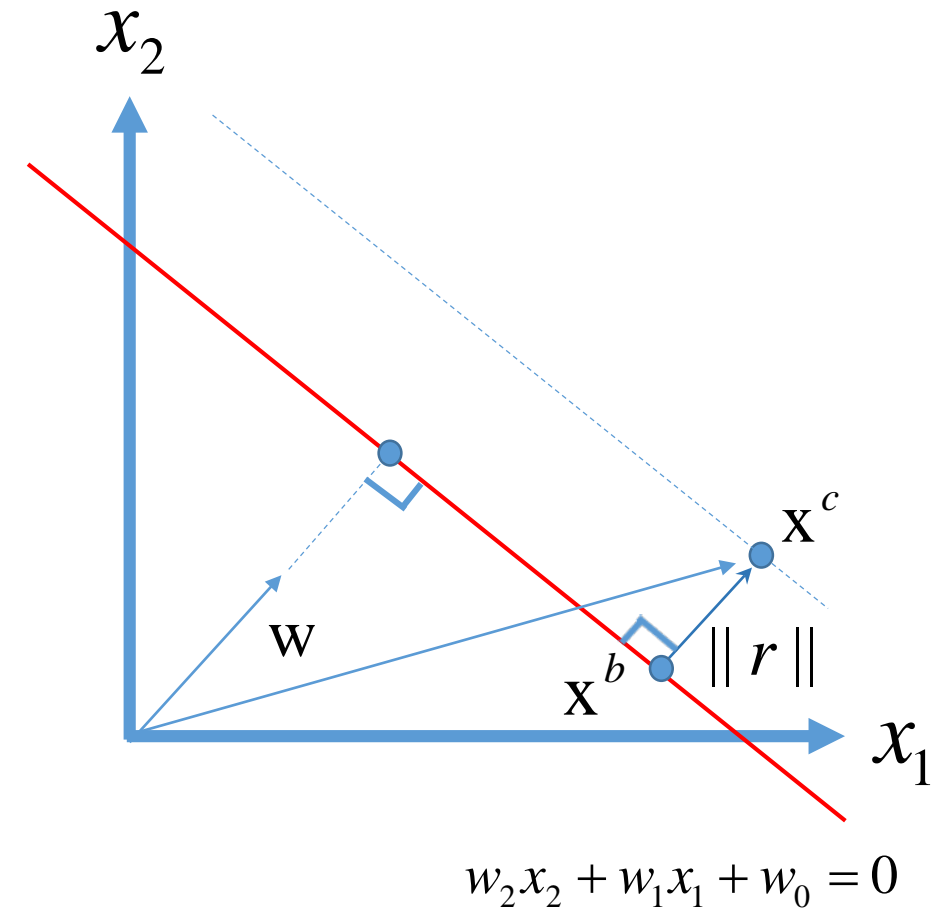
- Finding a decision boundary which maximizes the margin.

$$\max \|r\| = \frac{1}{\|w\|}$$

s.t.

$$t_n y(x_n) > 0 \quad \longrightarrow \quad \text{Every data points are classified correctly.}$$

$$\begin{cases} t_n = +1, & y(x_n) > 0 \\ t_n = -1, & y(x_n) < 0 \end{cases}$$



Problem formulation

- ❑ Let's make it a quadratic programming problem.

$$\max \frac{1}{\|w\|}$$

$$s.t. \quad t_n y(x_n) > 0, \quad \forall n$$

- ❑ Do you remember?

$$\max \frac{1}{\|w\|}$$

$$s.t. \quad t_n y(x_n) \geq 1, \quad \forall n$$

Let's say

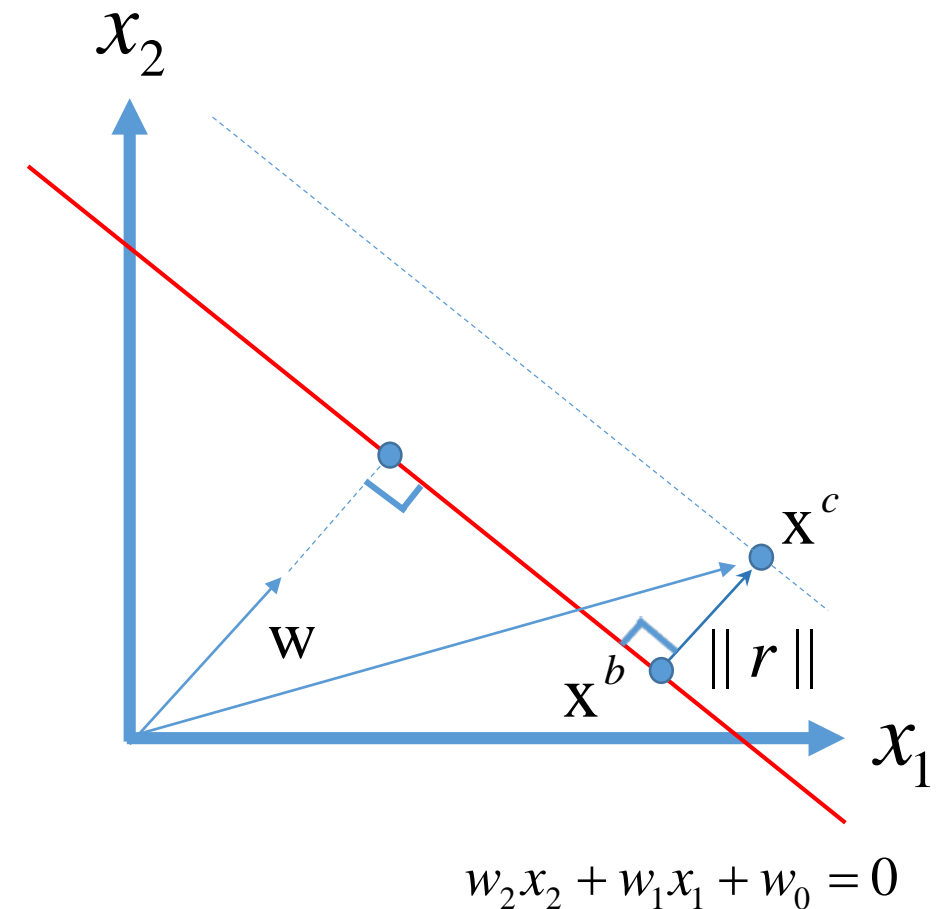
$$|y(x^c)| = 1$$

meaning that any data point is away from the decision boundary at least 1

- ❑ Finally

$$\min \frac{1}{2} \|w\|^2$$

$$s.t. \quad t_n (w^T x_n + w_0) \geq 1, \quad \forall n$$

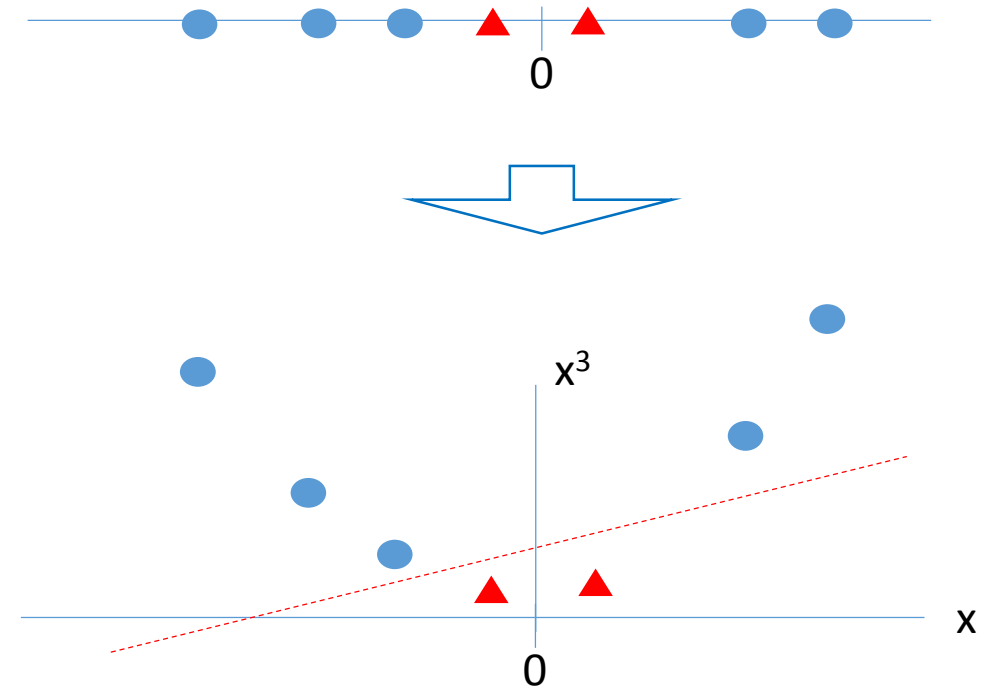
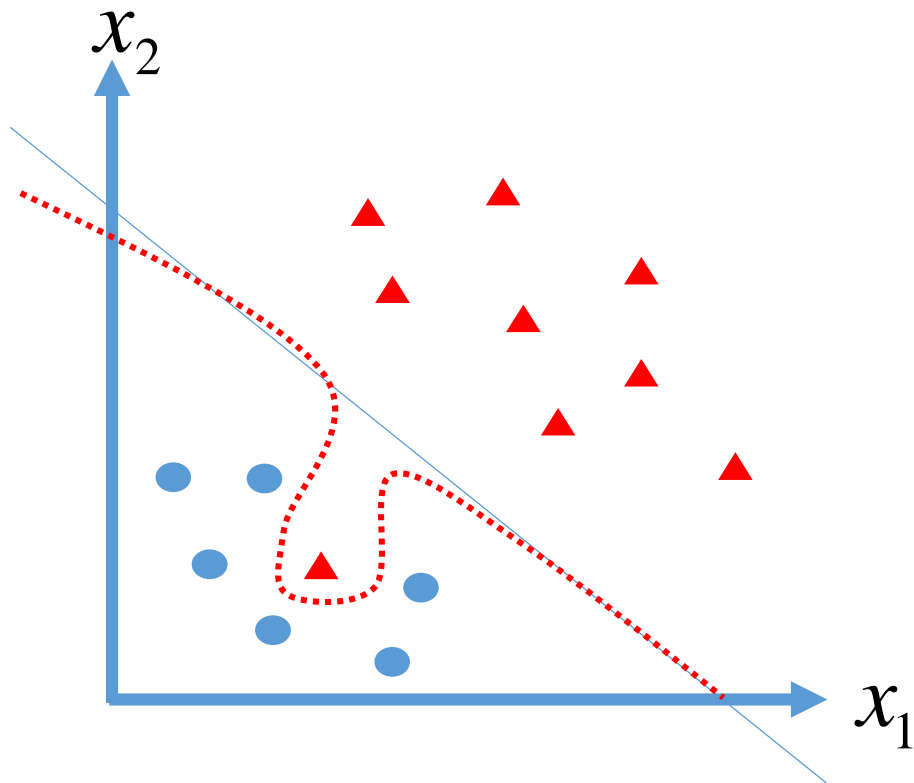


Quadratic programming

How about non-linearly separable case?

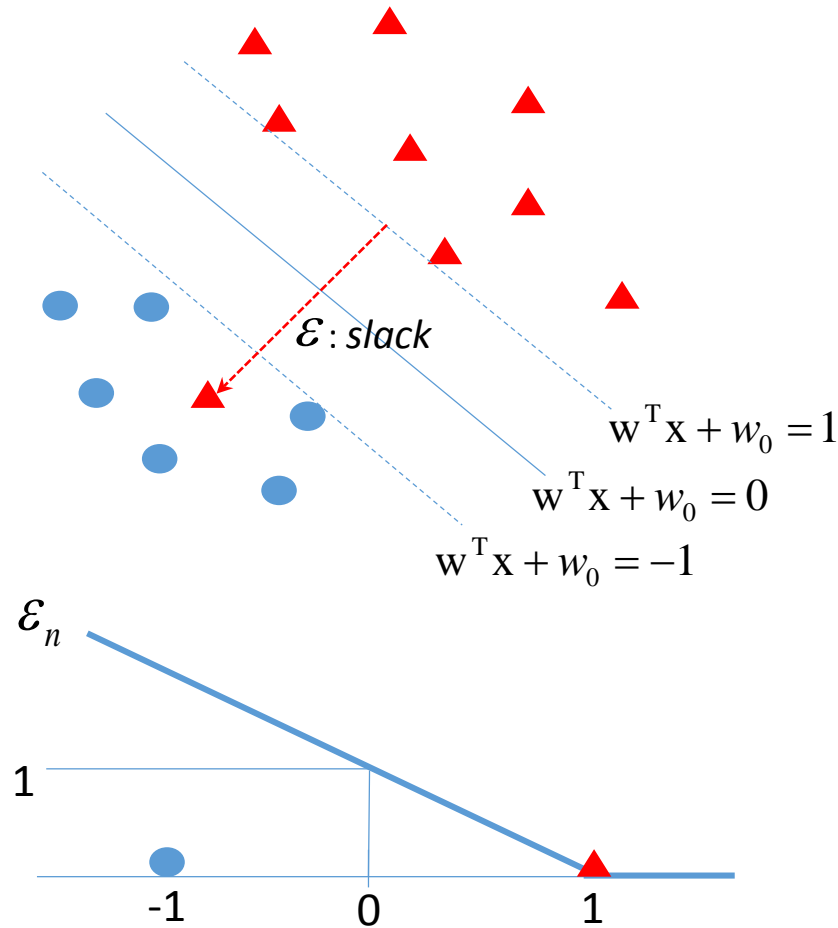
$$\begin{aligned} \min \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & t_n(w^T x_n + w_0) \geq 1, \quad \forall n \end{aligned}$$

	Approaches
Option 1	Soft margin SVM
Option 2	Kernel trick



Soft margin SVM

Option 1: soft margin SVM



- Remember the constraint below?

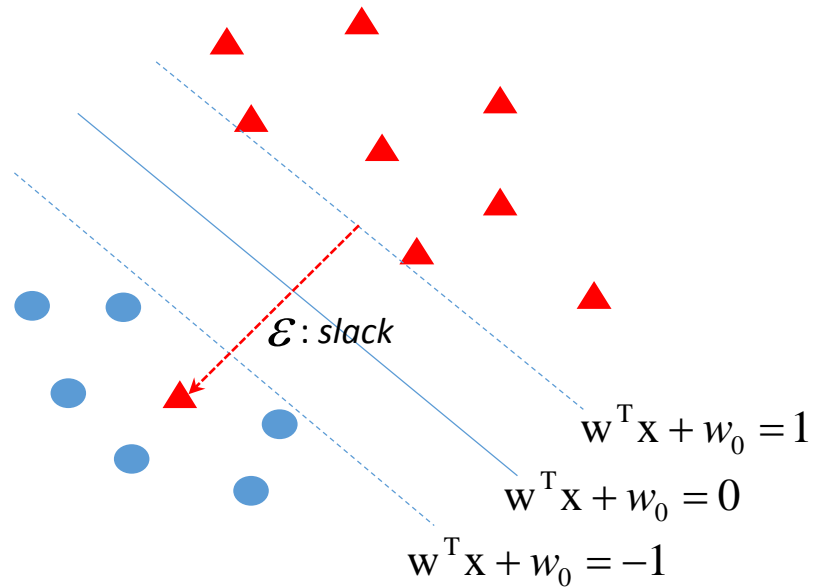
$$t_n(w^T x_n + w_0) \geq 1, \quad \forall n$$

- For the data points which are non-separable, we **relax** the constraint:

$$t_n(w^T x_n + w_0) \geq 1 - \varepsilon_n, \quad \forall n \quad \varepsilon_n \geq 0$$

- It says that **the distance between a data point and the decision boundary is allowed to be less than 1.**
- ε_n is called slack variables.
- Question. Where is a data point when $\varepsilon_n = 1$?

Option 1: soft margin SVM



- So we have the constraint below. How about the objective function?

$$t_n(w^T x_n + w_0) \geq 1 - \varepsilon_n, \quad \forall n \quad \varepsilon_n \geq 0$$

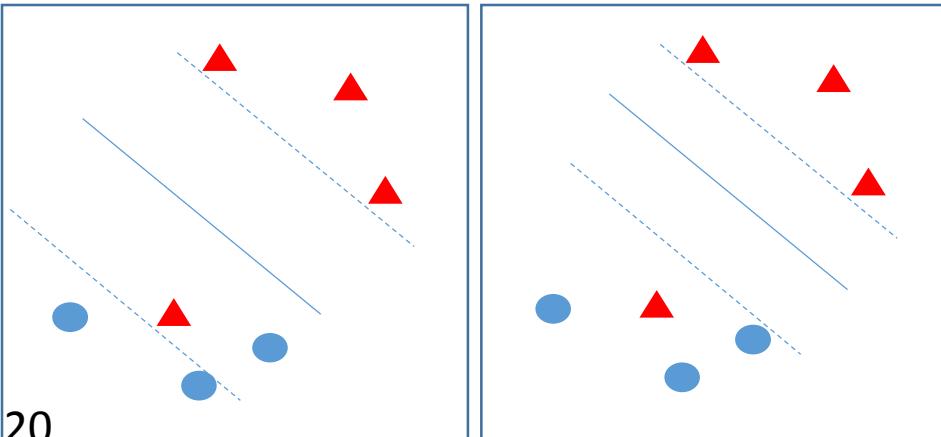
- We want to minimize the slack.

$$\min \frac{1}{2} \|w\|^2 + C \sum_n \varepsilon_n$$

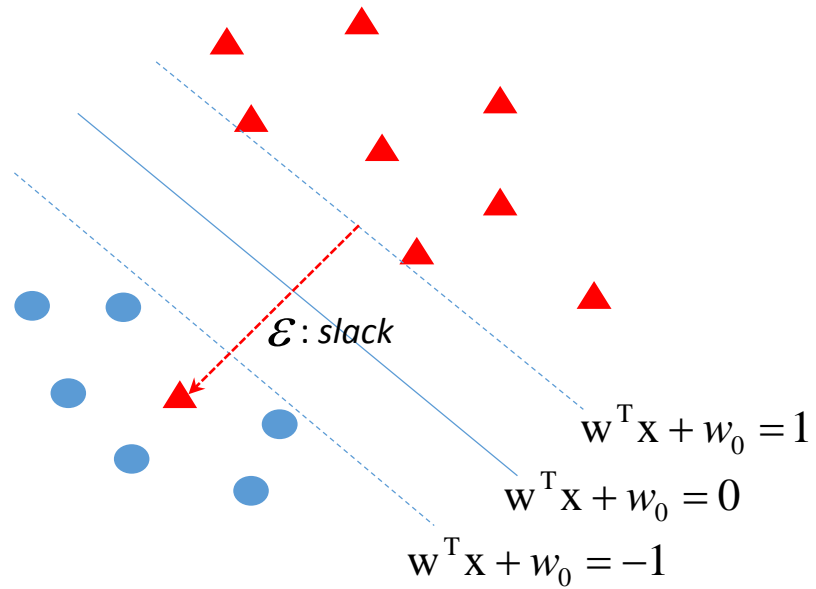
- If “C” is small, the slack contributes more
 - 1) Prefer large margin
 - 2) May cause large # of misclassified data points.
- If “C” is large, the slack contributes less
 - 1) Prefer less # of misclassified data points.
 - 2) May cause small margin.

“C” is small

“C” is large



Option 1: soft margin SVM



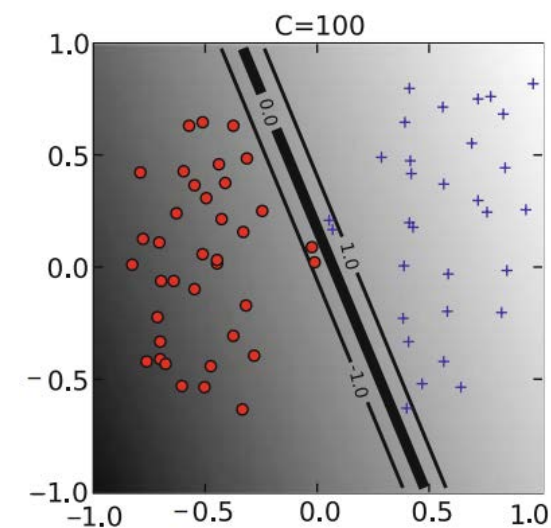
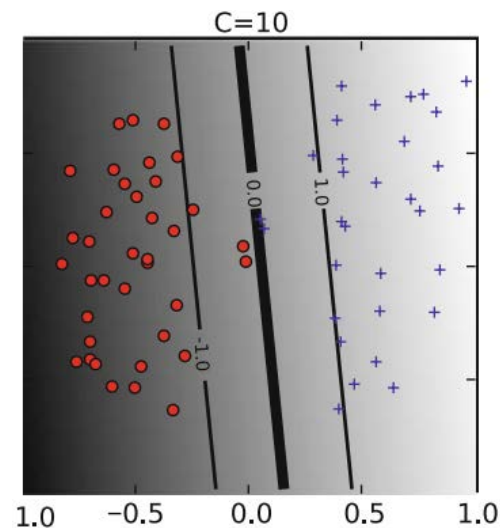
□ The formulation finally becomes

$$\min \frac{1}{2} \|w\|^2 + C \sum_n \varepsilon_n$$

s.t.

$$t_n(w^T x_n + w_0) \geq 1 - \varepsilon_n, \forall n$$

$$\varepsilon_n \geq 0$$

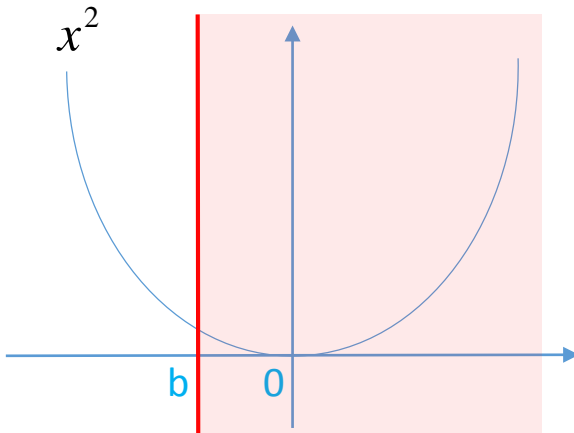


Kernel trick

Lagrange method for an optimization problem with inequality constraints

$$\begin{array}{ll} \min & x^2 \\ \text{s.t.} & x \geq b \end{array}$$

$$\begin{array}{ll} \min_x \max_{\lambda} & x^2 - \lambda(x - b) \\ \text{s.t.} & \lambda \geq 0 \end{array}$$



- Minima is zero when $b \leq 0$
- Minima is " b^2 " when $b > 0$
- It means at optima: $\lambda(x-b) = 0$ (complementary slackness)
- Maximizing λ results in minimizing the objective value
 - $\lambda \geq 0$ (it should be because $x-b \geq 0$)

Convert the quadratic problem in SVM to Lagrange optimization problem

Primal
problem



$$\begin{aligned} \min_{\mathbf{w}} \max_{\lambda} & \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{n=1}^n \lambda_n (t_n (\mathbf{w}^T \mathbf{x}_n + w_0) - 1) \\ \text{s.t.} \quad & \lambda_n \geq 0 \end{aligned}$$

KKT conditions

1) Stationarity condition

$$\frac{\partial}{\partial \mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} - \frac{\partial}{\partial \mathbf{w}} \sum_{n=1}^n \lambda_n (t_n (\mathbf{w}^T \mathbf{x}_n + w_0) - 1) = 0$$

2) Complementary slackness condition

$$\lambda_n (t_n (\mathbf{w}^T \mathbf{x}_n + w_0) - 1) = 0$$

3) Duality feasibility condition

$$\lambda_n \geq 0$$

Dual
problem



$$\begin{aligned} \max_{\lambda} \min_{\mathbf{w}} & \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{n=1}^n \lambda_n (t_n (\mathbf{w}^T \mathbf{x}_n + w_0) - 1) \\ \text{s.t.} \quad & \lambda_n \geq 0 \end{aligned}$$

- We would like to convert again the optimization problem above into another form, which provides same results.
- Because we want to solve the optimization problem in term of “lagrange multiplier (λ_n)”.

Dual problem of the quadratic problem: applying stationarity condition

$$\max_{\lambda} \min_{w, w_0} L(w, w_0, \lambda) = \frac{1}{2} w^T w - \sum_{n=1}^N \lambda_n (t_n (w^T x_n + w_0) - 1)$$

$$w = \sum_{n=1}^N \lambda_n t_n x_n \quad \sum_{n=1}^N \lambda_n t_n = 0$$

$$L(\lambda) = \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N t_n t_m \lambda_n \lambda_m x_n^T x_m - \sum_{n=1}^N \sum_{m=1}^N t_n t_m \lambda_n \lambda_m x_n^T x_m - \sum_{n=1}^N \lambda_n t_n w_0 + \sum_{n=1}^N \lambda_n$$


$$\max_{\lambda} L(\lambda) = \sum_{n=1}^N \lambda_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N t_n t_m \lambda_n \lambda_m x_n^T x_m$$

$$\lambda_n \geq 0$$

$$\sum_{n=1}^N \lambda_n t_n = 0$$

$$w = \sum_{n=1}^N \lambda_n t_n x_n$$

“w” does not appear in the equation, and so we do not use this constraint anymore



$$\begin{aligned}
 \max_{\lambda} L(\lambda) &= \sum_{n=1}^N \lambda_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m \\
 s.t. \quad \lambda_n &\geq 0, \quad \sum_{n=1}^N \lambda_n t_n = 0
 \end{aligned}$$

$$\begin{aligned}
 \min_{\lambda} L(\lambda) &= \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m - \sum_{n=1}^N \lambda_n \\
 s.t. \quad \lambda_n &\geq 0, \quad \sum_{n=1}^N \lambda_n t_n = 0
 \end{aligned}$$

□ Again, the optimization problem becomes a quadratic programming problem.

Let's summarize

$$\begin{aligned} \min_{\lambda} L(\lambda) &= \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m - \sum_{n=1}^N \lambda_n \\ \text{s.t. } \lambda &\geq 0, \quad t^T \lambda = 0 \end{aligned}$$

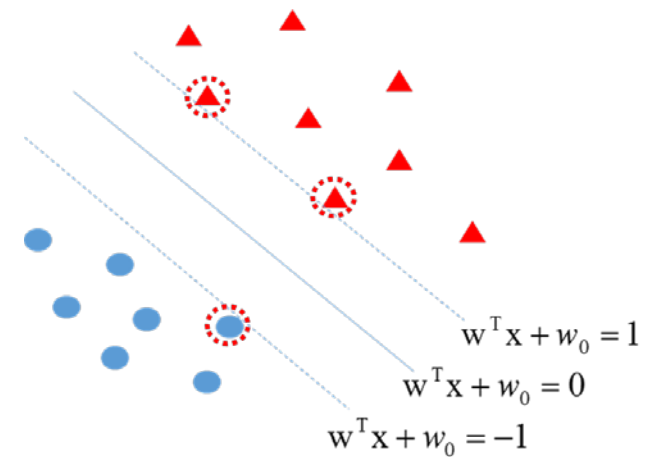
- ❑ The solution from the quadratic programming is “lagrange multipliers” (λ_n)
- ❑ Many of the solutions (lagrange multipliers) are zero
- ❑ Complementary slackness (one of KKT conditions) should be satisfied.

$$\lambda_n (t_n (w^T x_n + w_0) - 1) = 0$$

- ❑ In other words, if λ_n are not zero, $(t_n (w^T x_n + w_0) - 1)$ should be zero where corresponding data points should be support vectors.
- ❑ With the non-zero λ_n , w and w_0 can be calculated using $t_n (w^T x_n + w_0) = 1$

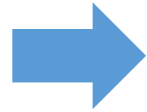
$$w = \sum_{n=1}^N \lambda_n t_n x_n$$

$$w_0 = t_n - \sum_{n=1}^N \lambda_n t_n x_n x_n$$



$$\min \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

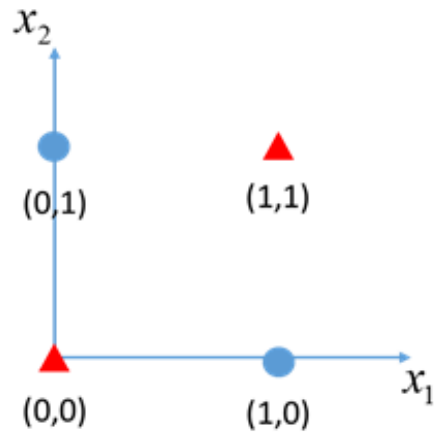
$$s.t. \quad t_n (\mathbf{w}^T \mathbf{x}_n + w_0) \geq 1$$



$$\min_{\lambda} L(\lambda) = \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m - \sum_{n=1}^N \lambda_n$$

$$s.t. \quad \lambda \geq 0, \quad t^T \lambda = 0$$

□ If data \mathbf{x}_n are not linearly separable, what should we do?

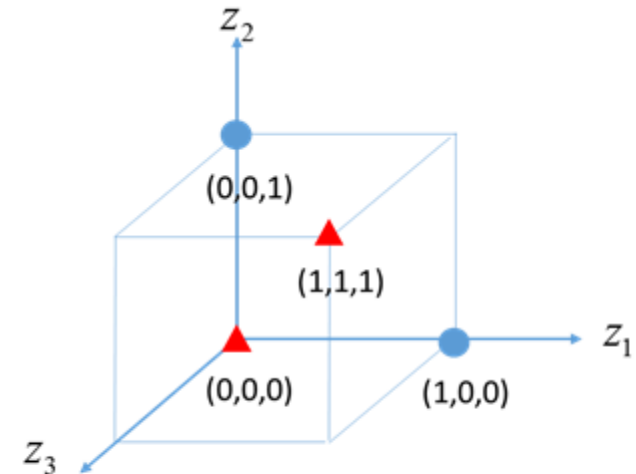


Space X

$$\phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ x_1 x_2 \\ x_2^2 \end{pmatrix}$$

$$\phi(\mathbf{x})$$

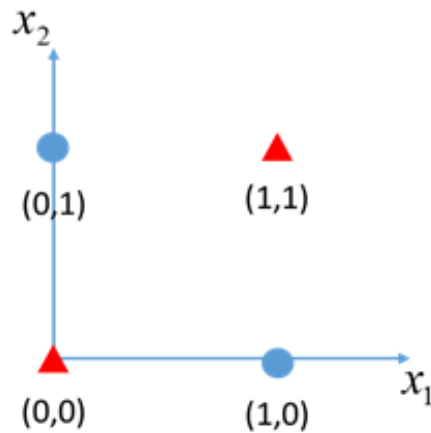
$$X \rightarrow Z$$



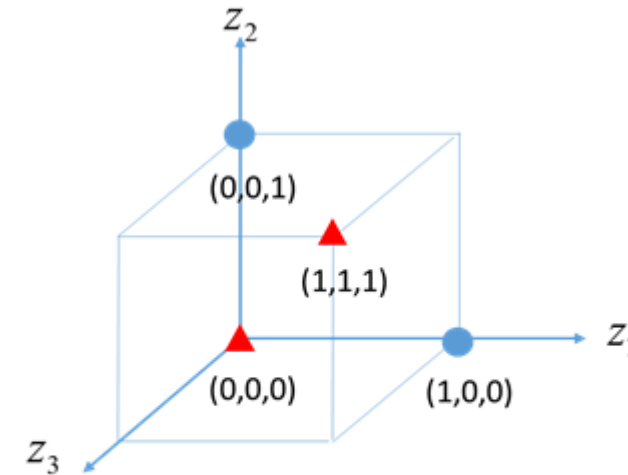
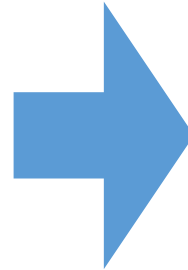
Space Z

Kernel trick

- The idea of Kernel trick begins from here: to find the scalar values (the inner product of two vectors: \mathbf{z}_n and \mathbf{z}_m) and so we can formulate the quadratic problem which can be linearly separable.



Space X



Space Z

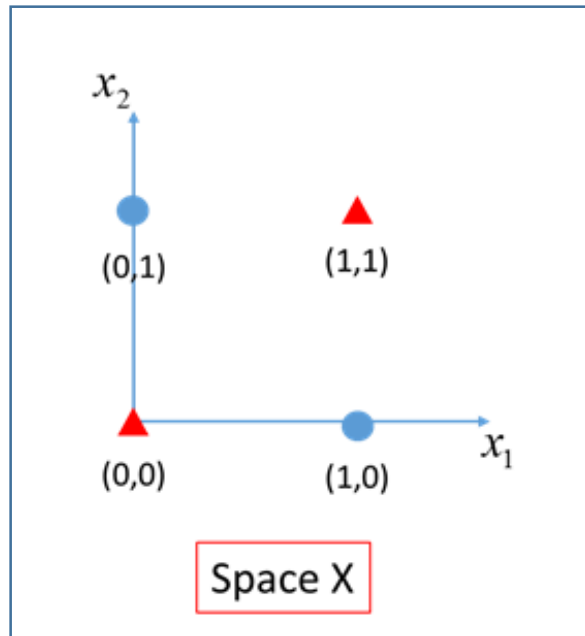
$$\begin{aligned} \min_{\lambda} \quad L(\lambda) &= \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m - \sum_{n=1}^N \lambda_n \\ \text{s.t.} \quad \lambda &\geq 0, \quad t^T \lambda = 0 \end{aligned}$$

$$\begin{aligned} \min_{\lambda} \quad L(\lambda) &= \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N t_n t_m \lambda_n \lambda_m \mathbf{z}_n^T \mathbf{z}_m - \sum_{n=1}^N \lambda_n \\ \text{s.t.} \quad \lambda &\geq 0, \quad t^T \lambda = 0 \end{aligned}$$

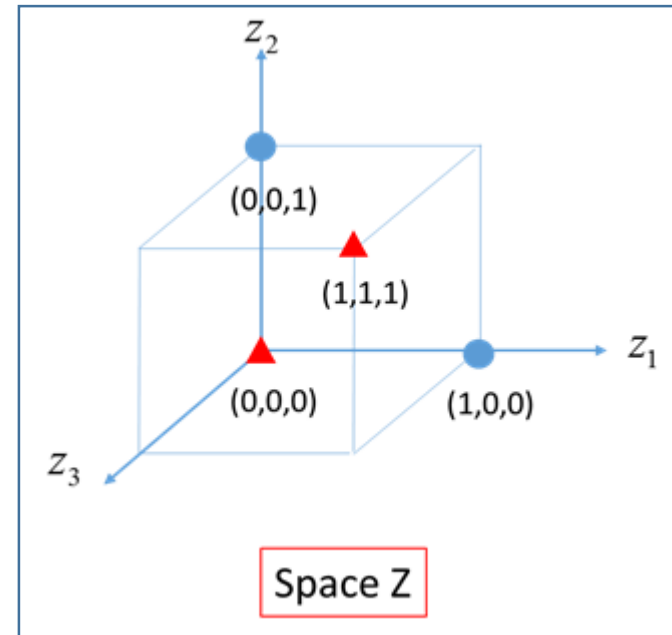
Kernel trick

- Kernel function $K()$ is a function which returns the scalar values (the inner product of two vectors: \mathbf{z}_n and \mathbf{z}_m in Z space) when the data points (\mathbf{x}_n and \mathbf{x}_m in X space) are given.

$$K(\mathbf{x}_n^T, \mathbf{x}_m) = \phi(\mathbf{x}_n^T)\phi(\mathbf{x}_m) = \mathbf{z}_n^T \mathbf{z}_m$$



$$\mathbf{z}_n^T = \phi(\mathbf{x}_n^T)$$
$$\mathbf{z}_m = \phi(\mathbf{x}_m)$$



Finally finally...

- With the Kernel function defined previously, we want to change the quadratic problem as follows:
- Because the Kernel function is a function of data points (\mathbf{x}_n and \mathbf{x}_m) which we already have.

$$\begin{aligned} \min_{\lambda} L(\lambda) &= \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N t_n t_m \lambda_n \lambda_m \mathbf{z}_n^T \mathbf{z}_m - \sum_{n=1}^N \lambda_n \\ \text{s.t. } \lambda &\geq 0, \quad t^T \lambda = 0 \end{aligned} \quad \Rightarrow \quad \begin{aligned} \min_{\lambda} L(\lambda) &= \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N t_n t_m \lambda_n \lambda_m K(\mathbf{x}_n^T \mathbf{x}_m) - \sum_{n=1}^N \lambda_n \\ \text{s.t. } \lambda &\geq 0, \quad t^T \lambda = 0 \end{aligned}$$

$$\min_{\lambda} L(\lambda) = \frac{1}{2} \lambda^T \begin{bmatrix} t_1 t_1 K(\mathbf{x}_1, \mathbf{x}_1) & t_1 t_2 K(\mathbf{x}_1^T, \mathbf{x}_2) & \cdots & t_1 t_N K(\mathbf{x}_1^T, \mathbf{x}_N) \\ t_2 t_1 K(\mathbf{x}_2, \mathbf{x}_1) & t_2 t_2 K(\mathbf{x}_2^T, \mathbf{x}_2) & \cdots & t_2 t_N K(\mathbf{x}_2^T, \mathbf{x}_N) \\ \cdots & \cdots & \cdots & \cdots \\ t_N t_1 K(\mathbf{x}_N, \mathbf{x}_1) & t_N t_2 K(\mathbf{x}_N^T, \mathbf{x}_2) & \cdots & t_N t_N K(\mathbf{x}_N^T, \mathbf{x}_N) \end{bmatrix} \lambda + (-1^T) \lambda$$

Finally finally...

$$\min_{\lambda} L(\lambda) = \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N t_n t_m \lambda_n \lambda_m \mathbf{K}(\mathbf{x}_n^T \mathbf{x}_m) - \sum_{n=1}^N \lambda_n$$
$$s.t. \quad \lambda \geq 0, \quad t^T \lambda = 0$$

$$\mathbf{w} = \sum_{z_n \in SV} \lambda_n t_n \mathbf{z}_n \quad w_0 = t_n - \sum_{z_n \in SV} \lambda_n t_n z_n z_n = t_n - \sum_{z_n \in SV} \lambda_n t_n K(\mathbf{x}_n, \mathbf{x}_n)$$

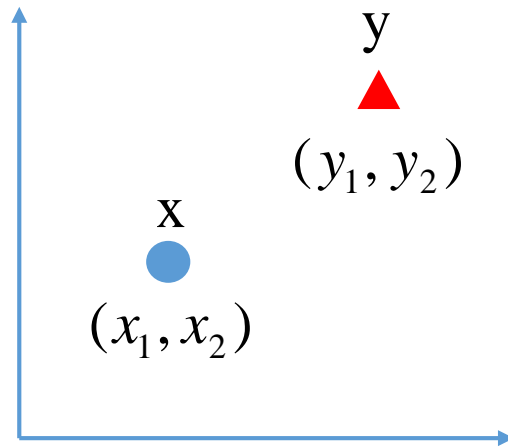
$$\text{sign}(\mathbf{w}^T \mathbf{z} + w_0)$$

$$\text{sign} \left(\sum \lambda_n t_n \mathbf{z}_n \mathbf{z} + t_n - \sum_{z_n \in SV} \lambda_n t_n K(\mathbf{x}_n, \mathbf{x}_n) \right)$$

$$\text{sign} \left(\sum \lambda_n t_n K(x_n, x) + t_n - \sum \lambda_n t_n K(x_n, x_n) \right)$$

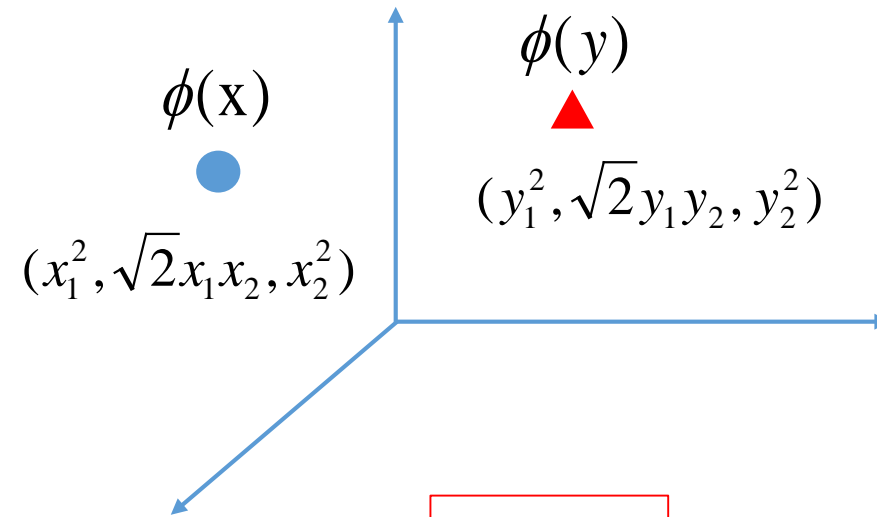
- Now you have a function, which classifies a data point in **z space** without mapping the data point to **z space** at all.
- Do you see why it is called a trick?

Polynomial kernel of degree 2



Space X

$$\begin{aligned} K(x, y) &= (xy)^2 \\ &= ((x_1, x_2) \cdot (y_1, y_2))^2 \\ &= (x_1 y_1 + x_2 y_2)^2 \\ &= x_1^2 y_1^2 + 2x_1 x_2 y_1 y_2 + x_2^2 y_2^2 \end{aligned}$$



Space Z

$$\begin{aligned} \phi(x)\phi(y) &= (x_1^2, \sqrt{2}x_1x_2, x_2^2) \cdot (y_1^2, \sqrt{2}y_1y_2, y_2^2) \\ &= x_1^2 y_1^2 + 2x_1 x_2 y_1 y_2 + x_2^2 y_2^2 \end{aligned}$$

Mapping to 3-dimension

Gaussian Kernel: derivation (inner product in the infinite z space)

$$K(\mathbf{x}_n, \mathbf{x}_m) = \exp(-\alpha \|\mathbf{x}_n - \mathbf{x}_m\|^2)$$

$$= \exp(-\alpha \mathbf{x}_n^2) \exp(-\alpha \mathbf{x}_m^2) \exp(2\alpha \mathbf{x}_n \mathbf{x}_m)$$

Taylor series expansion of
an exponential function

$$\exp(x) = \frac{x^0}{0!} + \frac{x}{1!} + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots$$

$$= \exp(-\alpha \mathbf{x}_n^2) \exp(-\alpha \mathbf{x}_m^2) \sum_{k=0}^{\infty} \frac{(2\alpha)^k (\mathbf{x}_n)^k (\mathbf{x}_m)^k}{k!}$$

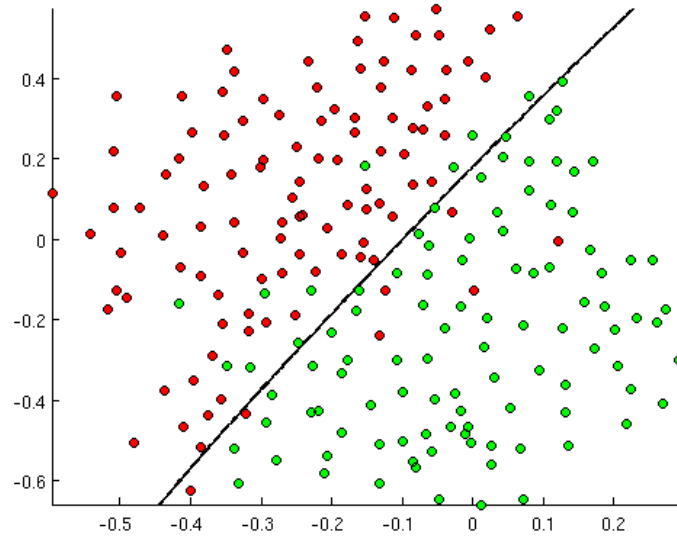
$$= \sum_{k=0}^{\infty} \sqrt{\frac{(2\alpha)^k}{k!}} \exp(-\alpha \mathbf{x}_n^2) (\mathbf{x}_n)^k \sqrt{\frac{(2\alpha)^k}{k!}} \exp(-\alpha \mathbf{x}_m^2) (\mathbf{x}_m)^k$$

$$= \phi(\mathbf{x}_n) \phi(\mathbf{x}_m)$$

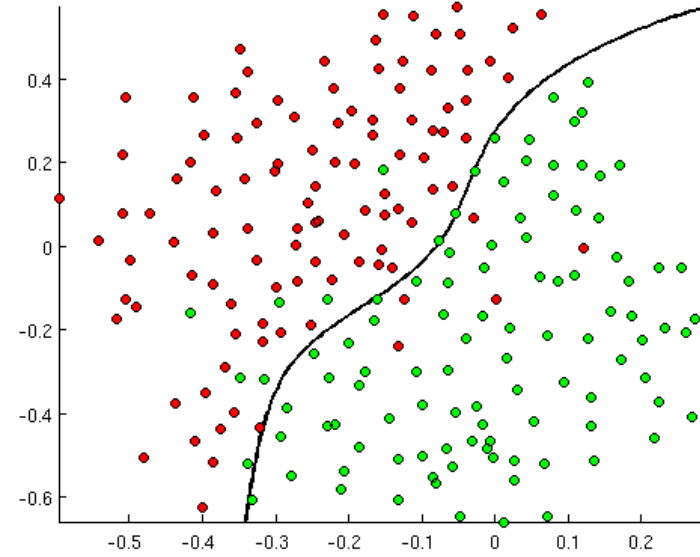
Mapping to infinite-dimension !

Gaussian Kernel

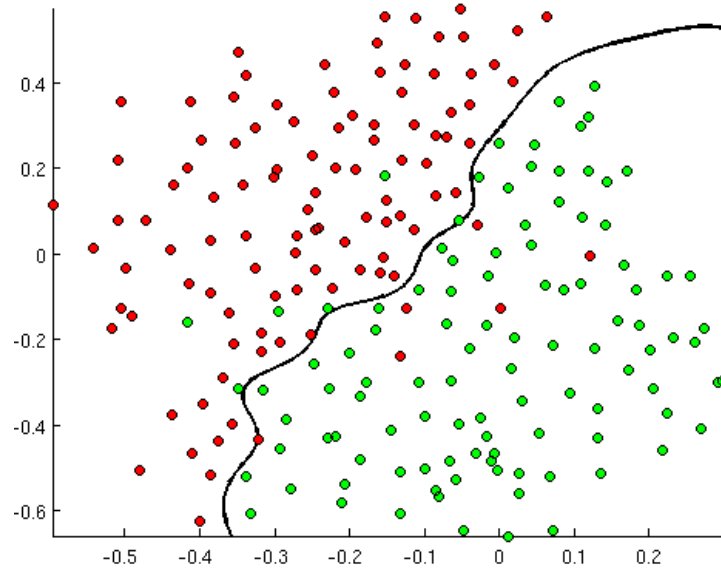
$\alpha = 1$



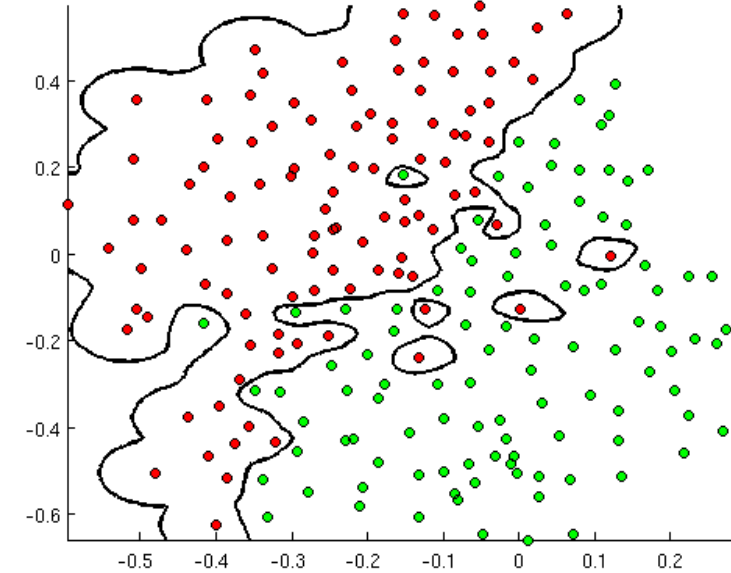
$\alpha = 10$



$\alpha = 100$



$\alpha = 1000$



Backup Slides

$$KL(q(Z) \parallel p(Z \mid X)) = -\sum_Z q(Z) \log \frac{p(Z \mid X)}{q(Z)} \longrightarrow \text{Finding } q(Z) \text{ which minimizes the Kullback divergence}$$

$$= -\sum_Z q(Z) \log \frac{p(X, Z)}{q(Z) p(X)}$$

$$KL = -\sum_Z q(Z) \log \frac{p(X, Z)}{q(Z)} + \sum_Z q(Z) \log p(X)$$

$$KL + \sum_Z q(Z) \log \frac{p(X, Z)}{q(Z)} = \sum_Z q(Z) \log p(X)$$

Variational inference: derivation

$$KL + \sum_Z q(Z) \log \frac{p(X, Z)}{q(Z)} = \sum_Z q(Z) \log p(X)$$

- 4-8=-1
- 1-2=-1

- Always positive
- Always **Negative**
- Lower bound (**L**)
- Always negative
- It is a **fixed** value

- ❑ KL divergence and lower bound are a function of “q(Z)”
- ❑ Minimizing KL divergence is equivalent to maximizing the lower bound (**L**).

↓

$$KL = - \sum_Z q(Z) \log \frac{p(Z | X)}{q(Z)}$$

- We do not have this conditional distribution

↑

$$L = \sum_Z q(Z) \log \frac{p(X, Z)}{q(Z)}$$

- We do have the joint distribution

Variational inference: derivation

$$\max L = \sum_Z q(Z) \log \frac{p(X, Z)}{q(Z)}$$

$$= \sum_{z_1} \sum_{z_2} q(z_1) q(z_2) \log \frac{p(x_1, x_2, z_1, z_2)}{q(z_1) q(z_2)}$$

Assuming that z_1 and z_2 are independent

$$= \sum_{z_1} \sum_{z_2} q(z_1) q(z_2) [\log p(x_1, x_2, z_1, z_2) - \log q(z_1) q(z_2)]$$

$$= \sum_{z_1} \sum_{z_2} q(z_1) q(z_2) [\log p(x_1, x_2, z_1, z_2) - \log q(z_1) - \log q(z_2)]$$

$$= \sum_{z_1} \sum_{z_2} q(z_1) q(z_2) \log p(x_1, x_2, z_1, z_2) - \sum_{z_1} \sum_{z_2} q(z_1) q(z_2) \log q(z_1) - \sum_{z_1} \sum_{z_2} q(z_1) q(z_2) \log q(z_2)$$

Assuming that $q(z_2)$ is known, and so we just look for $q(z_1)$

Variational inference: derivation

$$L = \sum_{z_1} \sum_{z_2} q(z_1)q(z_2) \log p(x_1, x_2, z_1, z_2) - \sum_{z_1} \sum_{z_2} q(z_1)q(z_2) \log q(z_1) - \sum_{z_1} \sum_{z_2} q(z_1)q(z_2) \log q(z_2)$$

Assuming that $q(z_2)$ is known, and so we just look for $q(z_1)$

$$= \sum_{z_1} \sum_{z_2} q(z_1)q(z_2) \log p(x_1, x_2, z_1, z_2) - \sum_{z_1} q(z_1) \log q(z_1) \sum_{z_2} q(z_2) - \sum_{z_1} q(z_1) \sum_{z_2} q(z_2) \log q(z_2)$$

$$= \sum_{z_1} q(z_1) \sum_{z_2} q(z_2) \log p(x_1, x_2, z_1, z_2) - \sum_{z_1} q(z_1) \log q(z_1) \sum_{z_2} q(z_2) - \sum_{z_1} q(z_1) \sum_{z_2} q(z_2) \log q(z_2)$$

$$= \sum_{z_1} q(z_1) E_{z_2} [\log p(x_1, x_2, z_1, z_2)] - \sum_{z_1} q(z_1) \log q(z_1) - K \sum_{z_1} q(z_1)$$

It is one but we keep it for a while

$$= \sum_{z_1} q(z_1) [E_{z_2} [\log p(x_1, x_2, z_1, z_2)] - K] - \sum_{z_1} q(z_1) \log q(z_1)$$

Variational inference: derivation

$$L = \sum_{z_1} q(z_1) [E_{z_2} [\log p(X, Z)] - K_1 - K_2] - \sum_{z_1} q(z_1) \log q(z_1)$$

$$\log f(X, Z) = E_{z_2} [\log p(X, Z)] - K_1$$

$$f(X, Z) = e^{E_{z_2} [\log p(X, Z)] - K_1} = e^{-K_1} e^{E_{z_2} [\log p(X, Z)]} = C e^{E_{z_2} [\log p(X, Z)]}$$

If we choose “C” carefully, $f(X, Z)$ can be a probability distribution.

$$\iint C e^{E_{z_2} [\log p(X, Z)]} dX dZ = 1$$

$$L = \sum_{z_1} q(z_1) [\log f(X, Z) - K_2] - \sum_{z_1} q(z_1) \log q(z_1)$$

$$= \sum_{z_1} q(z_1) \log f(X, Z) - \sum_{z_1} q(z_1) K_2 - \sum_{z_1} q(z_1) \log q(z_1)$$

$$= \sum_{z_1} q(z_1) \frac{\log f(X, Z)}{\log q(z_1)} - \sum_{z_1} q(z_1) K_2 = \sum_{z_1} q(z_1) \frac{\log f(X, Z)}{\log q(z_1)} + C'$$

Variational inference: derivation

$$L = \sum_{z_1} q(z_1) \frac{\log f(X, Z)}{\log q(z_1)} - \sum_{z_1} q(z_1) K_2 = \sum_{z_1} q(z_1) \frac{\log f(X, Z)}{\log q(z_1)} + C'$$

$$\log f(X, Z) = E_{z_2} [\log p(X, Z)] - K_1 \quad \text{We defined it previously}$$

$$f(X, Z) = e^{E_{z_2} [\log p(X, Z)] - K_1} = e^{-K_1} e^{E_{z_2} [\log p(X, Z)]} = C e^{E_{z_2} [\log p(X, Z)]}$$

Lower bound (L) is maximized when $\log q(z_1)$ and $\log p(X, Z)$ are equal because it is a negative KL.
Thus, ...

$$\log q(z_1) = \log f(X, Z)$$

$$q(z_1) = f(X, Z) = C_1 e^{E_{z_2} [\log p(X, Z)]} = C_1 e^{\sum q(z_2) \log p(X, Z)}$$
$$q(z_2) = f(X, Z) = C_2 e^{E_{z_1} [\log p(X, Z)]} = C_2 e^{\sum q(z_1) \log p(X, Z)}$$