

# Machine Learning I: Supervised Methods

B. Keith Jenkins

## Announcements

---



- Slido event code : 1345352
- Homework 5 is due Friday
- Sample exam problems
  - Discussion 7
  - Discussion 8
- Additional problems will be posted on D2L

## Today's lecture

---

- Lagrangian optimization (2)
  - Inequality constraints
- Kernels for nonlinear transformations
  - Examples of kernels
  - Kernel substitution
- Valid or Mercer kernels

# Lagrange Optimization with 1 Inequality Constraint

Extension of the equality constraint case.

Find min. of  $f(\underline{x})$  s.t.  $h(\underline{x}) \geq 0$

Let:

$$L(\underline{x}, \mu) = f(\underline{x}) - \mu h(\underline{x})$$

Find extremum of  $L$  over  $\underline{x}$  and  $\mu$ .

2 cases:

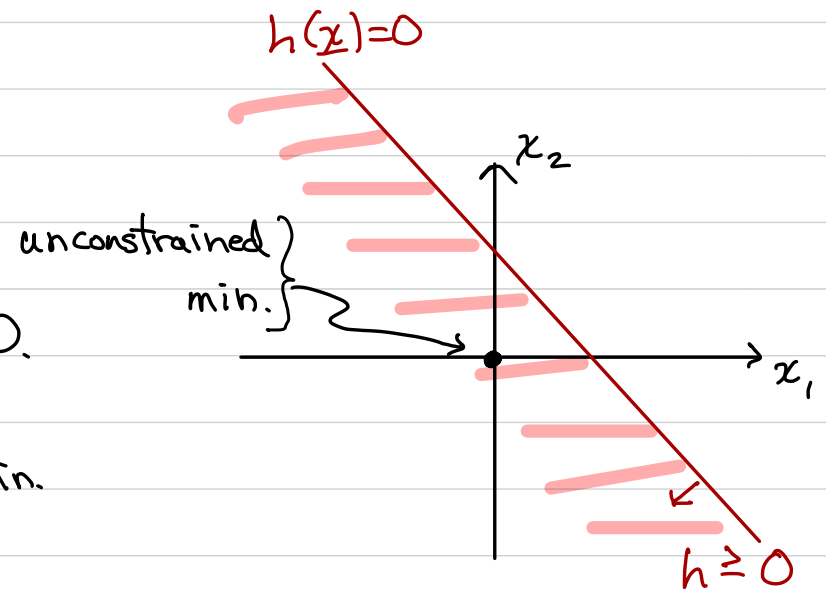
(a) unconstrained min. is in region  $h \geq 0$ .

$\Rightarrow$  constrained min. = unconstrained min.

Use  $\nabla_{\underline{x}} f(\underline{x}) = \underline{0}$

$\Rightarrow$  can solve using:

$L(\underline{x}, \mu)$  with  $\mu = 0$ .



(b) unconstrained min. is outside  $h \geq 0$  region.



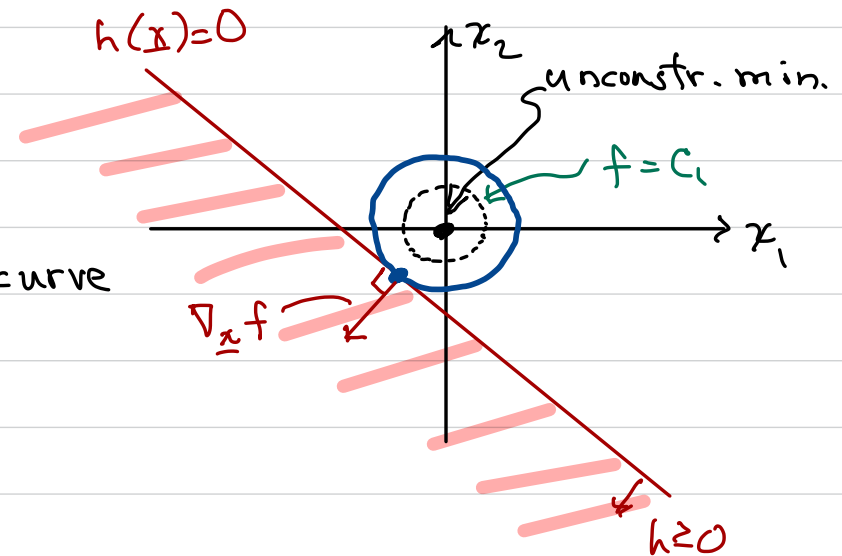
⑤

Constrained min. is on  $h(\underline{x})=0$  curve  
 $\Rightarrow$  same as equality constraint case

$$L(\underline{x}, \mu) \text{ with } \mu \neq 0$$

$$\nabla_{\underline{x}, \mu} L = \underline{0}$$

$$\Rightarrow \nabla_{\underline{x}} f = \mu \nabla_{\underline{x}} h, \quad \mu > 0$$



$$\therefore L(\underline{x}, \mu) = f(\underline{x}) - \mu h(\underline{x}), \quad h \geq 0 \text{ (constraint)}$$

In case (a):  $\mu = 0, h(\underline{x}^*) > 0$ .

In case (b):  $\mu > 0, h(\underline{x}^*) = 0$ .

$\Rightarrow$  in both cases,  $\mu h(\underline{x}^*) = 0$ .

## Summary of 1 ineq. constr.

$$\text{Min. } f(\underline{x}) \text{ s.t. } h(\underline{x}) \geq 0$$

$$L(\underline{x}, \mu) = f(\underline{x}) - \mu h(\underline{x})$$

$$\nabla_{\underline{x}, \mu} L(\underline{x}, \mu) = \underline{0}$$

$$\left. \begin{array}{l} h(\underline{x}^*) \geq 0 \\ \mu \geq 0 \\ \mu h(\underline{x}^*) = 0 \end{array} \right\} \begin{array}{l} \text{Karush-Kuhn-Tucker} \\ \text{(KKT) conditions} \end{array}$$

## General case: multiple equality and multiple inequality constraints

$$\text{Min. } f(\underline{x}) \text{ s.t. } g_i(\underline{x}) = 0 \quad \forall i, \quad h_j(\underline{x}) \geq 0 \quad \forall j$$

$$L(\underline{x}, \underline{\lambda}, \underline{\mu}) = f(\underline{x}) + \sum_{i=1}^R \lambda_i g_i(\underline{x}) - \sum_{j=1}^{R'} \mu_j h_j(\underline{x})$$

$$\nabla_{\underline{x}, \underline{\lambda}, \underline{\mu}} L(\underline{x}, \underline{\lambda}, \underline{\mu}) = \underline{0}$$

$$\text{Require: } \mu_j \geq 0 \quad \forall j$$

$$\mu_j h_j(\underline{x}^*) = 0 \quad \forall j$$

$$\text{Also given: } g_i(\underline{x}^*) = 0 \quad \forall i$$

$$h_j(\underline{x}^*) \geq 0 \quad \forall j$$

} KKT conditions

## Kernels

→ An alternate way of including nonlinear mappings  $\underline{u} = \underline{\phi}(\underline{x})$ .

Ex: Nearest-means classifier.

Stores class means  $\underline{\mu}_k$ ,  $k=1, 2, \dots, C$ .

Predicts  $\hat{y}(\underline{x})$  by taking  $\arg\min_k \|\underline{x} - \underline{\mu}_k\|_2^2$ .

$\underline{\mu}$  vs.  $\underline{u}$

Let  $\tilde{g}_k(\underline{x}) = -\|\underline{x} - \underline{\mu}_k\|_2^2 \leftarrow$  function that measures similarity between  $\underline{x}$  and  $\underline{\mu}_k$ .  
 Simplify  $\Rightarrow \tilde{g}_k(\underline{x}) = -(\underbrace{\underline{x}^T \underline{x}}_{\text{const. of } k} - 2 \underline{\mu}_k^T \underline{x} + \underline{\mu}_k^T \underline{\mu}_k)$

$$\Rightarrow \text{Let } g_k(\underline{x}) = 2 \underline{\mu}_k^T \underline{x} - \underline{\mu}_k^T \underline{\mu}_k$$

What if we want to implement nearest-means classifier in  $\underline{u}$ -space (after a nonlinear mapping  $\underline{u} = \underline{\phi}(\underline{x})$ )?

Tips:  $\underline{\mu}_k$  is a constant for a given dataset

$\underline{x}_i, \underline{x}_\ell$  are datapoints (given with the dataset)

$\underline{x}$  is a (vector) variable in the original feature space

$\underline{u}, \underline{\phi}$  are (vector) variables in the expanded feature space

2 approaches to include a nonlinear mapping  $\underline{\phi}(\underline{x})$ :

1. Specify  $\underline{\phi}(\underline{x})$  explicitly. Nearest-means example: ( $D=2$ )

(1) Quadratic polynomial mapping  $\underline{\phi}(\underline{x}) = [1, x_1, x_2, x_1^2, x_1 x_2, x_2^2]^T$ .

(2) Then:  $g_k(\underline{x}) = -\underline{\mu}_k^T \underline{\mu}_k + 2 \underline{\mu}_k^T \underline{x}$

(3) 
$$= \frac{-1}{N_k} \left( \sum_{i=1}^{N_k} \underline{x}_i^{(k)T} \right) \frac{1}{N_k} \left( \sum_{l=1}^{N_k} \underline{x}_l^{(k)} \right) + 2 \frac{1}{N_k} \left( \sum_{i=1}^{N_k} \underline{x}_i^{(k)T} \right) \underline{x}$$

$\underline{\phi}(\underline{x}) \left( \underline{x} \rightarrow \underline{\phi}(\underline{x}), \quad \underline{x}_i^{(k)} \rightarrow \underline{\phi}(\underline{x}_i^{(k)}), \text{ etc.} \right.$

(4) 
$$g_k(\underline{x}) = -\frac{1}{N_k^2} \sum_{i=1}^{N_k} \sum_{l=1}^{N_k} \underline{\phi}^T(\underline{x}_i^{(k)}) \underline{\phi}(\underline{x}_l^{(k)}) + \frac{2}{N_k} \sum_{i=1}^{N_k} \underline{\phi}^T(\underline{x}_i^{(k)}) \underline{\phi}(\underline{x})$$

Can then plug in for  $\underline{\phi}$  everywhere and simplify, or compute numerically directly from (4) and (1).

For approach 2, first write  $g_k(\underline{x})$  (or  $J(\underline{w})$ ) as a function of inner products of  $\underline{x}$  only ( $\underline{x}^T \underline{x}_i, \underline{x}_i^T \underline{x}_l$ , etc.), if possible:

$$(5) \quad \text{From (3), } g_k(\underline{x}) = -\frac{1}{N_k^2} \sum_{i=1}^{N_k} \sum_{l=1}^{N_k} \underline{x}_i^{(k)T} \underline{x}_l^{(k)} + \frac{2}{N_k} \sum_{i=1}^{N_k} \underline{x}_i^{(k)T} \underline{x}$$

↗ (In Lagrange optimization problems, we will get this form by deriving a dual representation.)

Then use kernels.

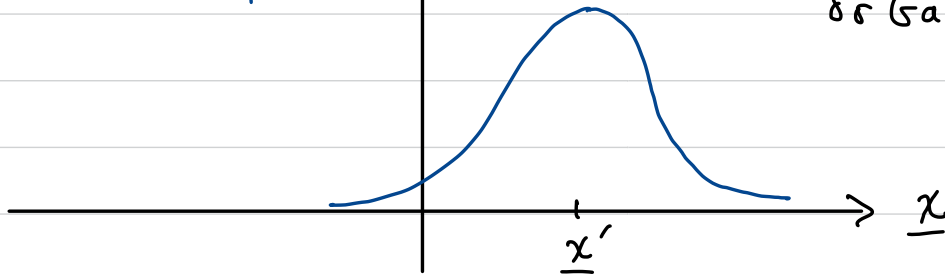
kernels are functions that help to measure similarity between 2 input vectors, call them  $\underline{x}$  and  $\underline{x}'$ ; and/or, add a nonlinear mapping to the system.

Ex:

$$(i) \quad k_{\text{linear}}(\underline{x}, \underline{x}') = \underline{x}^T \underline{x}' \quad (\text{linear kernel}).$$

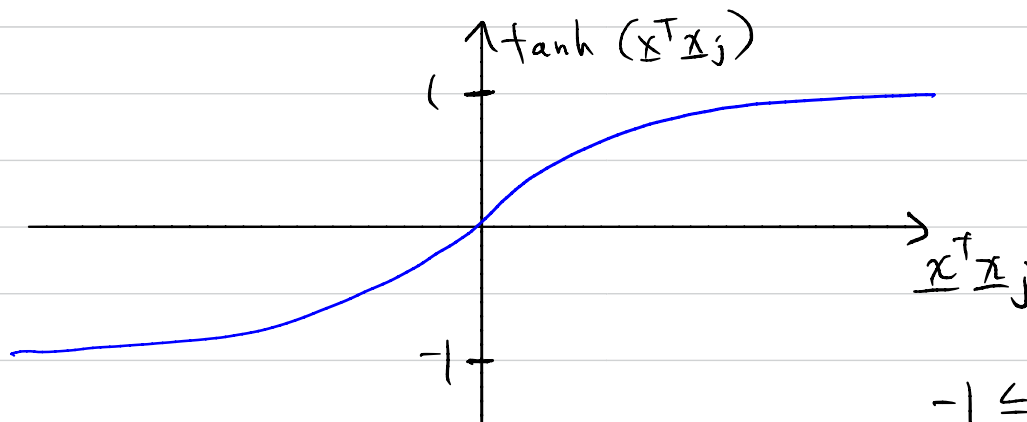
$$(ii) \quad k_{\text{poly}}(\underline{x}, \underline{x}') = (1 + \underline{x}^T \underline{x}')^d \quad (\text{polynomial kernel})$$

$$(iii) \quad k_{\text{RBF}}(\underline{x}, \underline{x}') = \exp \{ -\gamma \|\underline{x} - \underline{x}'\|_2^2 \}, \quad \gamma > 0$$

$k_{\text{rbf}}(\underline{x}, \underline{x}')$ (Radial Basis Function (RBF)  
or Gaussian kernel)

(iv)  $k_{\text{sigmoid}}(\underline{x}, \underline{x}_j) = \tanh[(\underline{x}^T \underline{x}_j) \alpha + c]$  (Sigmoid kernel)

Interpret: Let  $\alpha = 1, c = 0 \Rightarrow k_{\text{sigm}} = \tanh(\underline{x}^T \underline{x}_j)$



$$-1 \leq k_{\text{sigmoid}} \leq +1 \text{ always.}$$

$k_{\text{sigmoid}}$  is also a similarity measure.

Often kernels are chosen so that  $k(\underline{x}, \underline{x}') = \underline{\phi}^T(\underline{x}) \underline{\phi}(\underline{x}')$   
( $k$  is an inner product in  $\underline{u}$ -space).



2. Specify  $k(\underline{x}, \underline{x}')$  explicitly, which defines the nonlinear mapping implicitly.

Ex: RBF kernel:  $k_{\text{RBF}}(\underline{x}, \underline{x}') = \exp\{-\gamma \|\underline{x} - \underline{x}'\|_2^2\}$ .

Once  $g_k(\underline{x})$  or  $J(\underline{w})$  is written as a function of  $\underline{x}^T \underline{x}'$  (e.g.,  $\underline{x}^T \underline{x}_n$ ),  
substitute:

$$\underline{x}^T \underline{x}' \rightarrow k(\underline{x}, \underline{x}')$$

to implement the nonlinear mapping.

Nearest-means example: We had

$$(5) \quad g_k(\underline{x}) = -\frac{1}{N_k^2} \sum_{i=1}^{N_k} \sum_{l=1}^{N_k} \underline{x}_i^{(k)T} \underline{x}_l^{(k)} + \frac{2}{N_k} \sum_{i=1}^{N_k} \underline{x}_i^{(k)T} \underline{x}$$

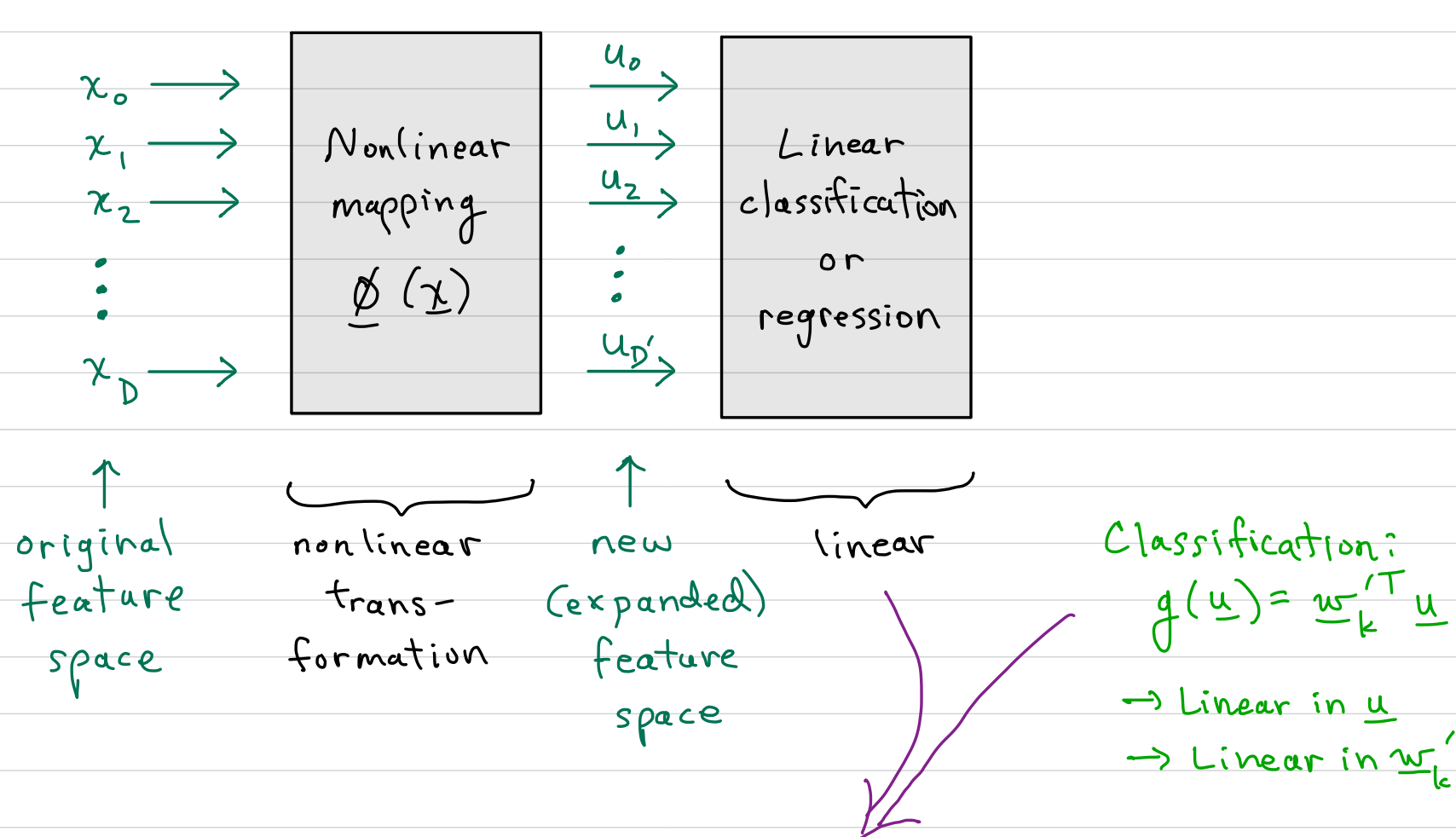
$\rightarrow$  Nonlinear mapping  $\underline{u} = \underline{\phi}(\underline{x}) \rightarrow$

$$g_k(\underline{u}) = -\frac{1}{N_k^2} \sum_{i=1}^{N_k} \sum_{l=1}^{N_k} \underbrace{k(\underline{x}_i^{(k)}, \underline{x}_l^{(k)})}_{\text{kernel}} + \frac{2}{N_k} \sum_{i=1}^{N_k} \underbrace{k(\underline{x}_i^{(k)}, \underline{x})}_{\text{kernel}}$$

For RBF kernel use  $k_{\text{RBF}}(\underline{x}, \underline{x}') = \exp\{-\gamma \|\underline{x} - \underline{x}'\|_2^2\}$

$$(6) \quad g_k(\underline{u}) = -\frac{1}{N_k^2} \sum_{i=1}^{N_k} \sum_{l=1}^{N_k} \exp\{-\gamma \|\underline{x}_i^{(k)} - \underline{x}_l^{(k)}\|_2^2\} + \frac{2}{N_k} \sum_{i=1}^{N_k} \exp\{-\gamma \|\underline{x} - \underline{x}_i^{(k)}\|_2^2\}$$

But our nonlinear-transformation approach was [Lecture 11, p.3]:



The linear classifier or regressor should be the same as original classifier or regressor, except with  $\underline{x} \rightarrow \underline{u}$  or  $\underline{x} \rightarrow \underline{\phi}(\underline{x})$ .

Is that true here? Compare (6) with (2), (4) for nearest-means example.

Because we substituted  $\underline{x}^T \underline{x}' \rightarrow k(\underline{x}, \underline{x}')$  to map  $\underline{x} \rightarrow \underline{u} = \underline{\phi}(\underline{x})$ , we want  $k(\underline{x}, \underline{x}') = \underline{\phi}^T(\underline{x}) \underline{\phi}(\underline{x}')$  for some N.L. mapping  $\underline{\phi}$ .

How do we know if our choice of  $k(\underline{x}, \underline{x}')$  can be expressed this way?

→ If  $k(\underline{x}, \underline{x}')$  is a valid kernel, then yes, it can.

Def: Gram matrix  $\underline{K}$ :

$$\underline{K} = \begin{bmatrix} k(\underline{x}_1, \underline{x}_1) & \cdots & k(\underline{x}_1, \underline{x}_N) \\ \vdots & \ddots & \vdots \\ k(\underline{x}_N, \underline{x}_1) & \cdots & k(\underline{x}_N, \underline{x}_N) \end{bmatrix}$$

If the Gram matrix  $\underline{K}$  is positive definite for all sets of input vectors  $\{\underline{x}_i\}_{i=1}^N$ , then  $k(\underline{x}, \underline{x}')$  is a valid kernel

(also called Mercer kernel or positive-definite kernel).

How to prove  $k(\underline{x}, \underline{x}')$  is a valid kernel, for a given  $k(\underline{x}, \underline{x}')$ ?

More commonly, for a given (possible) kernel  $k(\underline{x}, \underline{x}')$ , rather than proving  $\underline{K}$  is positive definite, one can use known valid kernels  $k'(\underline{x}, \underline{x}')$ , and properties to build up a new valid kernel from known valid kernels  $k'(\underline{x}, \underline{x}')$  [Bishop 6.2, "Constructing Kernels"].

Ex:

The RBF (Gaussian) kernel  $k(\underline{x}, \underline{x}') = \exp\{-\gamma \|\underline{x} - \underline{x}'\|_2^2\}$  can be shown to be a valid kernel, by building it out of  $k'(\underline{x}, \underline{x}') = \underline{x}^T \underline{x}'$  using known properties [Bishop 6.2].

So, in our earlier ( $\underline{u}$ -space) equation:

$$(6) \quad g_k(\underline{a}) = \underbrace{-\frac{1}{N_k^2} \sum_{i=1}^{N_k} \sum_{l=1}^{N_k} \exp\{-\gamma \|\underline{x}_i^{(k)} - \underline{x}_l^{(k)}\|_2^2\}}_{k(\underline{x}_i^{(k)}, \underline{x}_l^{(k)})} + \underbrace{\frac{2}{N_k} \sum_{i=1}^{N_k} \exp\{-\gamma \|\underline{x} - \underline{x}_i^{(k)}\|_2^2\}}_{k(\underline{x}, \underline{x}_i^{(k)})}$$

$$= \underbrace{\underline{\phi}^T(\underline{x}_i^{(k)}) \underline{\phi}(\underline{x}_l^{(k)})}_{k(\underline{x}_i^{(k)}, \underline{x}_l^{(k)})} = \underbrace{\underline{\phi}^T(\underline{x}) \underline{\phi}(\underline{x}_i^{(k)})}_{k(\underline{x}, \underline{x}_i^{(k)})}$$

for some mapping  $\underline{\phi}(\underline{x})$  that exists but hasn't been stated.

## Comments:

1. Implementing a N.L. mapping  $\phi(\underline{x})$  by substituting:

$$\underline{x}^T \underline{x}' \rightarrow k(\underline{x}, \underline{x}')$$

is known as the kernel trick or kernel substitution.

2. Why bother?

(i) Sometimes we know a good  $k(\underline{x}, \underline{x}')$ , but don't know the  $\phi(\underline{x})$ .

Ex:  $k_{\text{rbf}}(\underline{x}, \underline{x}') = \exp\{-\gamma \|\underline{x} - \underline{x}'\|_2^2\}$  can be a good similarity fcn.

(ii) Sometimes  $\phi(\underline{x})$  is difficult to work with.

Ex:  $\phi(\underline{x})$  corresponding to  $k_{\text{rbf}}(\underline{x}, \underline{x}')$  is infinite dimensional!