# Reproducing Kernel Hilbert Space

Consider a Mercer kernel $K(\underline{x}, \cdot)$ where $\underline{x} \in \mathcal{H}$ and $\mathcal{F}$ be any vector space of all real valued functions of $\underline{x}$ generated by $K(\underline{x}, \cdot)$

Suppose we pick two functions $f(\cdot)$ and $g(\cdot)$ from $\mathcal{F}$

(I.P. space)

$$f(\cdot) = \sum_{i=1}^{\ell} a_i K(\underline{x}_i, \cdot)$$

$\uparrow$ $x$

$\uparrow$ $x$

$$g(\cdot) = \sum_{j=1}^{n} b_j K(\underline{\tilde{x}}_j, \cdot)$$

for all

$\underline{x}_i, \underline{\tilde{x}}_j \in \mathcal{H}$

Defn

Consider the bilinear form

$$\langle f, g \rangle = \sum_{i=1}^{\ell} \sum_{j=1}^{n} a_i \, b_j \, K(\underline{x}_i, \underline{\tilde{x}}_j)$$

$$= \underline{a}^T K \underline{b}$$

Gram matrix / Kernel matrix

$$\langle K(\underline{x}_i, \cdot), K(\underline{x}_j, \cdot) \rangle = K(\underline{x}_i, \underline{x}_j)$$

One element of the Gram matrix

We can rewrite $\langle f, g \rangle$ as

$$\langle f, g \rangle = \sum_{i=1}^{l} a_i \underbrace{\sum_{j=1}^{n} b_j \, K(\underline{x}_i, \underline{\tilde{x}}_j)}_{g(\underline{x}_i)} \qquad \left( \because K(\underline{x}_i, \underline{\tilde{x}}_j) = K(\underline{\tilde{x}}_j, \underline{x}_i) \right)$$

$$= \sum_{i=1}^{l} a_i \, g(\underline{x}_i)$$

$$\text{III}^{ly} \quad \langle f, g \rangle = \sum_{j=1}^{n} b_j \, f(\underline{\tilde{x}}_j)$$

# Properties

1) **Symmetry :** For all fns $f$ and $g \in \mathcal{F}$ the term $\langle f, g \rangle$ is symmetric

i.e., $\langle f, g \rangle = \langle g, f \rangle$

2) **Scaling and distribution**

For any pair of constants $c$ and $d$ and any set of functions $f, g$ and $h \in \mathcal{F}$

$$\langle (cf + dg), h \rangle = c \langle f, h \rangle + d \langle g, h \rangle$$

3) ## Squared norm

For any real valued fn $f \in \mathcal{F}$

$$\|f\|^2 = \langle f, f \rangle$$

$$= \underline{a}^T K \underline{a} \qquad \left( \begin{array}{c} \text{non negative} \\ \text{definite} \end{array} \right)$$

$$\|f\|^2 \geqslant 0$$

4)

Suppose $\quad g(\cdot) = K(\underline{x}, \cdot)$

$$\langle f, K(\underline{x}, \cdot) \rangle = \sum_{i=1}^{\ell} a_i K(\underline{x}, \underline{x}_i) \quad (\because \text{Symmetry})$$

$$K(\underline{x}, \underline{x}_i) = K(\underline{x}_i, \underline{x})$$

$$= \sum_{i=1}^{\ell} a_i K(\underline{x}_i, \underline{x})$$

$$= f(\underline{x}) \qquad \left( \begin{array}{c} \text{Reproducing} \\ \text{Kernel} \end{array} \right)$$

Mercer Kernel reproduces $f(\cdot)$

1) For every $x_i \in \mathcal{H}$, $K(\underline{x}, \underline{x_i})$ as a function of $\underline{x} \in \mathcal{F}$

Satisfies reproducing property

2)

Mercer kernel $\implies$ Reproducing kernel

Reproducing kernel space Complete $\implies$ Reproducing kernel Hilbert space

# Representer Theorem

**Theorem :** Any function defined in a RKHS can be represented as a linear combination of Mercer kernel functions.

**Proof :** Define a space $\mathcal{H}$ to represent RKHS induced by a Mercer kernel $K(\underline{x}, \cdot)$. Given any real valued fn $f(\cdot) \in \mathcal{H}$, we could decompose $f(\cdot)$ into 2 components lying in $\mathcal{H}$.

The first component $f_{||}(\cdot)$ is contained in the span of the kernel fns $K(\underline{x}_1, \cdot), K(\underline{x}_2, \cdot) \ldots$

$$f_{||}(\cdot) = \sum_{i=1}^{\ell} a_i K(\underline{x}_i, \cdot) \quad \underline{\qquad} \quad \text{①}$$

The second component is orthogonal to the span of the kernel fns; $f_{\perp}(\cdot)$

$$f(\cdot) = f_{||}(\cdot) + f_{\perp}(\cdot) \quad \underline{\qquad} \quad \text{②}$$

$$f(\cdot) = \sum_{i=1}^{l} a_i \, K(x_i, \cdot) + \underline{f_\perp(\cdot)} \qquad ③$$

From the <u>reproducing property</u>

$$f(x_j) = \langle f(\cdot), K(x_j, \cdot) \rangle_{\mathcal{H}} \qquad\qquad ④$$

Plug in ③ into ④

$$f(x_j) = \langle \left[ \sum_{i=1}^{l} a_i \, K(x_i, \cdot) + f_\perp(\cdot) \right] K(x_j, \cdot) \rangle$$

$$f(\underline{x}_j) = \left\langle \sum_{i=1}^{\ell} a_i K(\underline{x}_i, \cdot), K(\underline{x}_j, \cdot) \right\rangle + \left\langle f_\perp(\cdot), K(\underline{x}_j, \cdot) \right\rangle$$

$$= \sum_{i=1}^{\ell} a_i K(\underline{x}_i, \underline{x}_j)$$

$$\left( \because K(\underline{x}_i, \underline{x}_j) = \langle K(\underline{x}_i, \cdot), K(\underline{x}_j, \cdot) \rangle \right)$$

Mercer Kernel functions

# Generalized Applicability

**Theorem:** $f(\underline{x}_j) = \sum\limits_{i=1}^{l} a_i \, K(\underline{x}_i, \underline{x}_j)$ is the minimizer of the <u>regularized</u> empirical risk

$$\mathcal{E}(f) = \frac{1}{2N} \sum_{i=1}^{N} \left( d(n) - f(x(n)) \right)^2 + \Omega \left( \|f\|_{\mathcal{H}} \right)$$

Std error

regularizing fn.

$\left( \begin{array}{l} \Omega(\cdot) \text{ must be} \\ \text{a non decreasing fn.} \end{array} \right)$

$f(\cdot)$ unknown

$(x(n), d(n))$ Data pairs

$n = 1, \ldots, N$

## Proof :

**Step 1 :** Let $f_\perp$ denote the orthogonal complement to the span of the kernel fns $\{K(x_i, \cdot)\}_{i=1}^{\ell}$

Now, every fn can be expressed as a Kernel expansion on the training $+ f_\perp$

$$\Omega\left(\|f\|_{\mathcal{H}}\right) = \Omega\left(\left\|\sum_{i=1}^{\ell} a_i K(x_i, \cdot) + f_\perp(\cdot)\right\|_{\mathcal{H}}\right)$$

Introduce

$$\tilde{\Omega}\left(\|f\|_{\mathcal{H}}^2\right) = \Omega\left(\|f\|_{\mathcal{H}}\right)$$

$$\tilde{\Omega}\left(\|f\|_{\mathcal{H}}^2\right) = \tilde{\Omega}\left(\left\|\sum_{i=1}^{\ell} a_i\, k(x_i, \cdot) + f_1(\cdot)\right\|_{\mathcal{H}}^2\right)$$

**Step 2 :**  Apply Pythagorus theorem

$$\tilde{\Omega}\left(\| f \|_{\mathcal{H}}^2\right) = \tilde{\Omega}\left(\left\| \sum_{i=1}^{\ell} a_i K(x_i, \cdot) \right\|_{\mathcal{H}}^2 + \| f_\perp(\cdot) \|_{\mathcal{H}}^2\right)$$

$$\geq \tilde{\Omega}\left(\left\| \sum_{i=1}^{\ell} a_i K(x, \cdot) \right\|_{\mathcal{H}}^2\right)$$

Set $f_\perp(\cdot) = 0$ for optimality

$$\tilde{\Omega}\left(\| f \|_{\mathcal{H}}^2\right) = \tilde{\Omega}\left(\left\| \sum_{i=1}^{\ell} a_i K(x, \cdot) \right\|_{\mathcal{H}}^2\right)$$

Step 3 : In light of monotonicity

$$\Omega\left(\|f\|_{\mathcal{H}}\right) = \Omega\left(\left\|\sum_{i=1}^{l} a_i K(x_i, \cdot)\right\|\right)$$

$\implies$ For fixed $a_i \in \mathbb{R}$, the representer theorem is also a minimizer of the regularizing fn $\Omega\left(\|f\|_{\mathcal{H}}\right)$ provided monotonicity is satisfied !

# MOTIVATION TO REGULARIZATION THEORY

Often, in machine learning problems, we encounter situations where problems are not well-posed.

For example, when the # of data points in the training samples >> # of degrees of freedom, the problem is over determined.
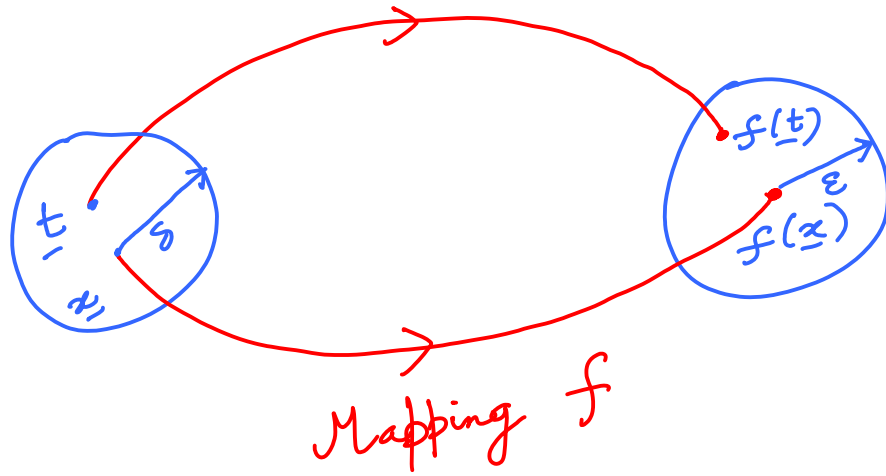
One may fit misleading variations in the data!

Learning is a sort of multi-D mapping (f), and can be viewed as a problem of hyper surface reconstruction given a set of sparse points

Now, given $X$ (domain) and $Y$ (range) that are metric spaces, related by a fixed but unknown mapping

$$f : X \rightarrow Y$$

The problem of reconstructing $f$ is well-posed if it satisfies the following:

a) Existence: For every input vector $\underline{x} \in X$, $\exists$ a $\underline{y} = f(\underline{x})$, $\underline{y} \in Y$

b) Uniqueness: For any pair of input vectors $\underline{x}, \underline{t} \in X$

$$f(\underline{x}) = f(\underline{t}) \text{ iff } \underline{x} = \underline{t}$$

c) Continuity: For any $\varepsilon > 0$, $\exists \; \delta = \delta(\varepsilon)$ /

$$d(\underline{x}, \underline{t}) < \delta \implies d\left(f(\underline{x}), f(\underline{t})\right) < \varepsilon$$

Mapping f

How can one make an ill-posed problem, well-posed?

SOLN :     Regularization     ( Tikhonov)

Consider the following problem

Input signal : $\quad \underline{x}_i \in \mathbb{R}^{m_0} \qquad i = 1, \ldots, N$

Desired signal : $\quad d_i \in \mathbb{R} \qquad i = 1, \ldots, N$

Data

$\{\underline{x}_i, d_i\}_{i=1}^{N}$

Let the approximating function be $F(\underline{x})$

$$\mathcal{E}_s(F) = \frac{1}{2} \sum_{i=1}^{N} \left( d_i - F(\underline{x}_i) \right)^2$$

( Approximation error)

Introduce the regularization term that
depends on the geometry of the problem

$$\mathcal{E}_c^{(Reg)}(F) = \frac{1}{2} \|DF\|^2$$

Linear differential operator $D$

Choice of '$D$' is problem dependent !

$\|\cdot\|$ is the norm over which the function space belongs to.

$$\mathcal{E}(F) = \mathcal{E}_S(F) + \lambda\, \mathcal{E}_c(F)$$

$$= \frac{1}{2} \sum_{i=1}^{N} \left( d_i - F(\underline{x}_i) \right)^2 + \frac{1}{2}\lambda \, \| D F \|^2$$

<span style="color:red">approx. error</span>

<span style="color:red">regularization term</span>

$\mathcal{E}(F)$ is also called the <span style="color:red">"Tikhonov functional"</span>

$\lambda \longrightarrow 0$ ; unconstrained

$\lambda \longrightarrow \infty$ ; $x_i$'s are unreliable

Choose $\lambda$ in between $(0, \infty)$

Normalize / $\lambda \in (0, 1)$ i.e., a fraction

Now $\quad F_\lambda(x) = \min_{\lambda, \underline{\omega}} \; \mathcal{E}(F)$ $\left(\begin{array}{l}\text{min. Tikhonov}\\\text{functional}\end{array}\right)$

$\longleftarrow$ parameter in $F(\cdot)$

Consider the standard error term differential

$$d\,\mathcal{E}_s(F, h) = \left[\frac{d}{d\beta}\,\mathcal{E}_s\left(F + \beta\, h\right)\right]_{\beta = 0}$$

$h(\underline{x})$ is a fixed function of '$x$'

$$d\left(\mathcal{E}(F,h)\right) = d\left(\mathcal{E}_s(F,h)\right) + \lambda d\left(\mathcal{E}_c(F,h)\right) = 0$$

$$d\left(\mathcal{E}_s(F,h)\right) = \frac{1}{2}\frac{d}{d\beta}\sum_{i=1}^{N}\left[d_i - F(\underline{x}_i) - \beta h(\underline{x}_i)\right]^2$$

$$= -\sum_{i=1}^{N}\left[d_i - F(\underline{x}_i) - \beta h(\underline{x}_i)\right]h(\underline{x}_i)\Bigg|_{\beta=0}$$

$$= -\sum_{i=1}^{N}\left(d_i - F(x_i)\right)h(\underline{x}_i)$$

$$= -\left\langle h, \left(\underline{d} - F(\underline{x})\right)\delta(\underline{x}-x_i)\right\rangle$$

Illy doing it over the regularization terms

$$d(\mathcal{E}_c(F, h)) = \frac{d}{d\beta} \mathcal{E}_c(F + \beta h)\Big|_{\beta = 0}$$

$$= \frac{1}{2} \frac{d}{d\beta} \int_{\mathbb{R}^{m_0}} \left(D(F + \beta h)\right)^2 d\underline{x}\Big|_{\beta = 0}$$

$$= \int_{\mathbb{R}^{m_0}} D(F + \beta h) \cdot D h \, d\underline{x}\Big|_{\beta = 0}$$

$$= \int_{\mathbb{R}^{m_0}} DF \cdot D h \, d\underline{x} \qquad = \langle DF, D h \rangle_{\mathcal{H}}$$

# Euler - Lagrange equation

Given a linear differential operator $D$, we can find a uniquely determined adjoint operator by $\widetilde{D}$ for any pair of functions $u(\underline{x})$ and $v(\underline{x})$ that are sufficiently differentiable ( upto a certain degree) & satisfy proper boundary conditions

$$\int_{\mathbb{R}^m} u(\underline{x}) \, D \, V(\underline{x}) \, d\underline{x} = \int_{\mathbb{R}^m} v(\underline{x}) \, \widetilde{D} \, u(\underline{x}) \, d\underline{x}$$

$D$ is a matrix.

With $\quad u(\underline{x}) \doteq DF(\underline{x})$

and $\quad \underline{v}(\underline{x}) \doteq h(\underline{x})$

$$d\mathcal{E}_c(F, h) = \int_{\mathbb{R}^m} \underbrace{h(\underline{x})}_{v(\underline{x})} \underbrace{\widetilde{D} D F(\underline{x})}_{u(\underline{x})} d\underline{x}$$

$$= \langle h(\underline{x}), \widetilde{D} D F \rangle_{\mathcal{H}}$$

Interpret this as an inner product

With the inclusion of a regularization parameter,

$$d\mathcal{E}(F,h) = \left\langle h, \left[ \tilde{D}DF - \frac{1}{\lambda}\sum_{i=1}^{N}(d_i - F)\,\delta_{x_i} \right] \right\rangle_{\mathcal{H}}$$

regulatory

Fréchet differential

desired value

approximating fn

$\lambda \in (0, \infty)$

$d\mathcal{E}(F,h)$ is zero for every $h(\underline{x})$ in $\mathcal{H}$ space

$$\text{iff} \quad \tilde{D}DF - \frac{1}{\lambda}\sum_{i=1}^{N}(d_i - F)\,\delta_{x_i} = 0$$

$$\text{i.e.,} \quad \tilde{D}DF_{\lambda}(\underline{x}) = \frac{1}{\lambda}\sum_{i=1}^{N}\left(d_i - F(\underline{x}_i)\right)\delta(\underline{x}-\underline{x}_i) \qquad (A)$$

# Green's Function

Eqn (A) represents a partial differential eqn in the approximating function $F$, whose solution is of interest.

Let $G(\underline{x}, \underline{\xi})$ be a function of $\underline{x}$ and $\underline{\xi}$.

(Green's function)

↑ Some argument

satisfying certain properties.

For a given linear differential operator $L$, $G(\underline{x}, \underline{\xi})$ satisfies the following properties: (Courant & Hilbert)

1) For a fixed $\underline{\xi}$, $G(\underline{x}, \underline{\xi})$ is a function of $\underline{x}$ satisfying the boundary Conditions

2) Except @ $\underline{x} = \underline{\xi}$, the derivatives of $G(\underline{x}, \underline{\xi})$ w.r.t $\underline{x}$ are all Continuous; the # of derivatives is determined by $L$

3) $\quad L \, G(\underline{x}, \underline{\xi}) = 0 \quad$ everywhere $\quad$ except

@ $\underline{x} = \underline{\xi}$, where it is singular.

$$L \, G(\underline{x}, \underline{\xi}) = \delta(\underline{x} - \underline{\xi}) \quad @ \quad \underline{x} = \underline{\xi} \quad \text{exists}$$

The function $G(\underline{x}, \underline{\xi})$ is called the Green's

function of operator $L$.

$\quad$ ( Similar to the inverse of a matrix eqn! )

Let $\varphi(\underline{x})$ be a continuous/piecewise continuous function of $\underline{x} \in \mathbb{R}^m$, then

**Claim:** $f(\underline{x}) = \int_{\mathbb{R}^m} G(\underline{x}, \underline{\xi}) \varphi(\underline{\xi}) d\underline{\xi}$ is a solution to $L f(\underline{x}) = \varphi(\underline{x})$

Let us verify the validity!

$$L F(\underline{x}) = L \int_{\mathbb{R}^m} G(\underline{x}, \underline{\xi}) \varphi(\underline{\xi}) \, d\underline{\xi}$$

$$\underbrace{\qquad\qquad\qquad}_{F(\underline{x})}$$

$$= \int_{\mathbb{R}^m} L \, G(\underline{x}, \underline{\xi}) \varphi(\underline{\xi}) \, d\underline{\xi}$$

From the property of Green's functions.

$$= \int_{\mathbb{R}^m} \delta(\underline{x} - \underline{\xi}) \varphi(\underline{\xi}) \, d\underline{\xi} = \varphi(\underline{x})$$

$$L F(\underline{x}) = \varphi(\underline{x})$$

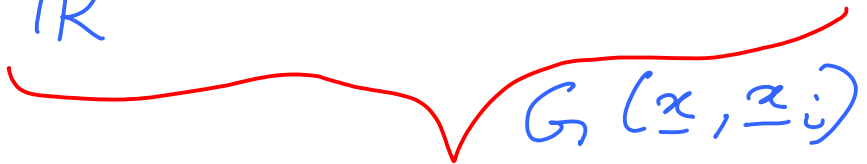Let us look into the regularization problem

$$L = \tilde{D} D$$

$$\psi(\underline{\xi}) = \frac{1}{\lambda} \sum_{i=1}^{N} \left( d_i - F(\underline{x}_i) \right) \delta(\underline{x}_i - \underline{\xi})$$

↓ Plug in

$$F_\lambda(\underline{x}) = \int_{\mathbb{R}^m} G(\underline{x}, \underline{\xi}) \, \psi(\underline{\xi}) \, d\underline{\xi}$$

$$F_\lambda(\underline{x}) = \int_{\mathbb{R}^m} G(\underline{x}, \underline{\xi}) \left\{ \frac{1}{\lambda} \sum_{i=1}^{N} [d_i - F(\underline{x}_i)] \, \delta(\underline{x}_i - \underline{\xi}) \right\} d\underline{\xi}$$

$$= \frac{1}{\lambda} \sum_{i=1}^{N} [d_i - F(\underline{x}_i)] \underbrace{\int_{\mathbb{R}^m} G(\underline{x}, \underline{\xi}) \, \delta(\underline{x}_i - \underline{\xi}) \, d\underline{\xi}}_{G(\underline{x}, \underline{x}_i)}$$

$$F_\lambda(\underline{x}) = \frac{1}{\lambda} \sum_{i=1}^{N} (d_i - F(\underline{x}_i)) \, G(\underline{x}, \underline{x}_i)$$

The minimizing function to the regularization problem
is a linear superposition of $N$ - green functions.
The points $\underline{x}_i$ represent the <u>centers of the</u>
<u>expansion</u> and $\left(d_i - F(\underline{x}_i)\right)/\lambda$ represent the
<u>weights of the expansion</u>

$\left\{ G(\underline{x}, \underline{x}_i) \right\}_{i=1}^{N}$ centered @ $\underline{x} = \underline{x}_i$ constitute
the basis of a subspace of smooth
fn. where the soln to the regularization problem
lies

# How do we determine the Coeffts ($w_i$)?

Let $w_i \stackrel{\Delta}{=} \dfrac{1}{\lambda}\left[d_i - F(\underline{x}_i)\right]$; $i = 1, \ldots, N$

Continuous

$$F_\lambda(\underline{x}) = \sum_{i=1}^{N} w_i\, G(\underline{x}, \underline{x}_i) \underline{\qquad} \textcircled{1}$$

($\#$ of Green's functions $= \#$ of data points)

Evaluate $\textcircled{1}$ @ $\underline{x}_j$; $j = 1, \ldots, N$

data points

Let $F_\lambda \triangleq \left[ F_\lambda(\underline{x}_1) \cdots F_\lambda(\underline{x}_N) \right]^T$

$d \triangleq \left[ d_1 \cdots d_N \right]^T ; \quad \underline{w} = \left[ w_1 \cdots w_N \right]^T$

$$G \triangleq \begin{bmatrix} G(\underline{x}_1, \underline{x}_1) & \cdots & G(\underline{x}_1, \underline{x}_N) \\ \vdots & \ddots & \\ G(\underline{x}_N, \underline{x}_1) & & G(\underline{x}_N, \underline{x}_N) \end{bmatrix}$$

(Gram matrix)

$N \times N$

Writing in matrix form,

$$\underline{w} = \frac{1}{\lambda}\left[\underline{d} - \underline{F}_\lambda\right] \implies \underline{F}_\lambda = \underline{d} - \lambda\underline{w}$$

$$\underline{F}_\lambda = G\,\underline{w}$$

$$G := \left[G(\underline{x}_i, \underline{x}_j)\right]$$

$$\therefore \left(G + \lambda I\right)\underline{w} = \underline{d}$$

But, the adjoint of the linear differential operator $L$

$$\tilde{L} = L \implies \text{Green's fns are symmetric!}$$
$$G(\underline{x}_i, \underline{x}_j) = G(\underline{x}_j, \underline{x}_i)$$

However, all functions in the null space
of $D$ are invisible to the smoothing term /
regulatory constraints $\|DF\|^2$
and is problem dependent.

The RBF happens to be a special case of Green's function that is <u>translationally</u> and <u>rotationally invariant</u> i.e., if $G(\underline{x}, \underline{x}_i) = G(||\underline{x} - \underline{x}_i||)$

For RBF,

$$F_\lambda(\underline{x}) = \sum_{i=1}^{N} w_i \, G(||\underline{x} - \underline{x}_i||) \quad \left(\text{Linear function space}\right) \left(\begin{matrix}\text{depends} \\ \text{on} \\ \text{data!}\end{matrix}\right)$$
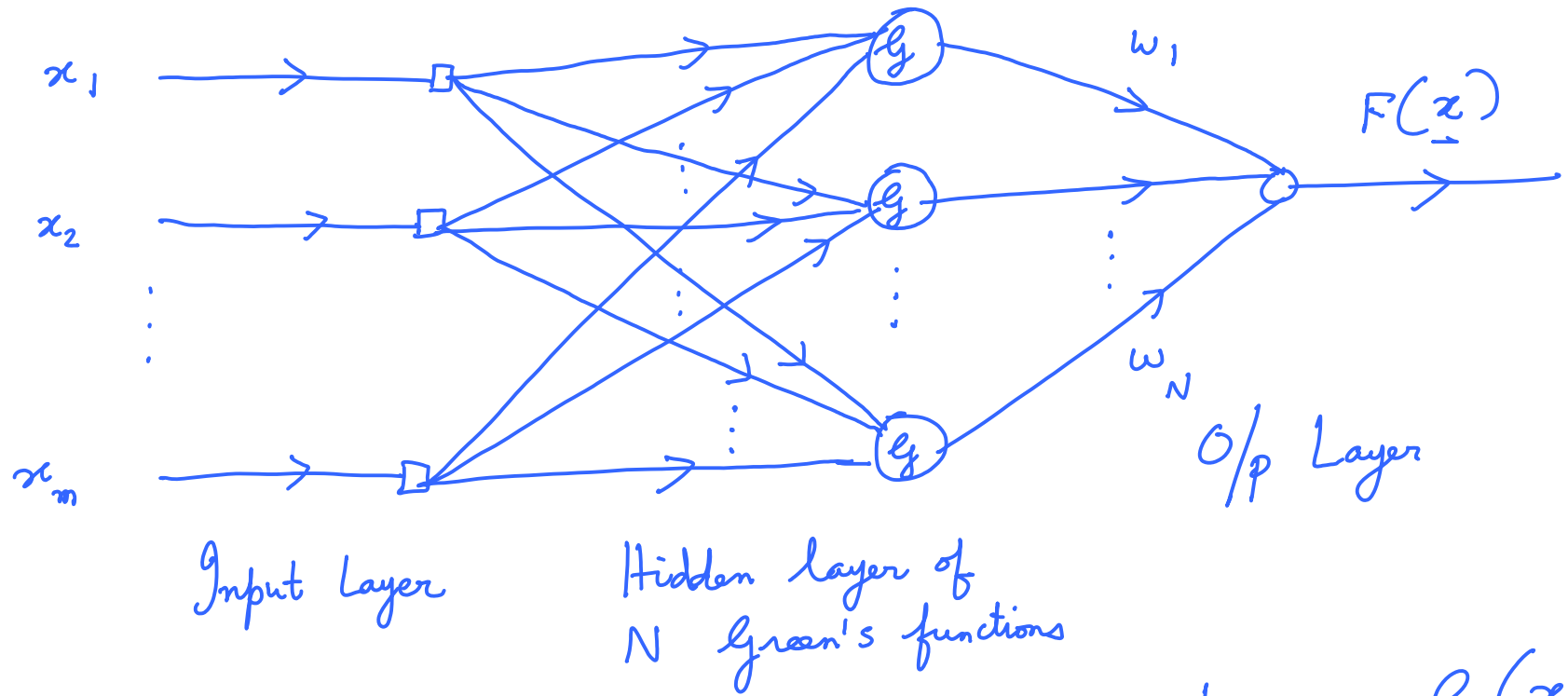
Assuming       Gaussian units,

$$F_\lambda(\underline{x}) = \sum_{i=1}^{N} w_i \, \exp\left(-\frac{1}{2\sigma_i^2} \left\| \underline{x} - \underline{x}_i \right\|^2\right)$$
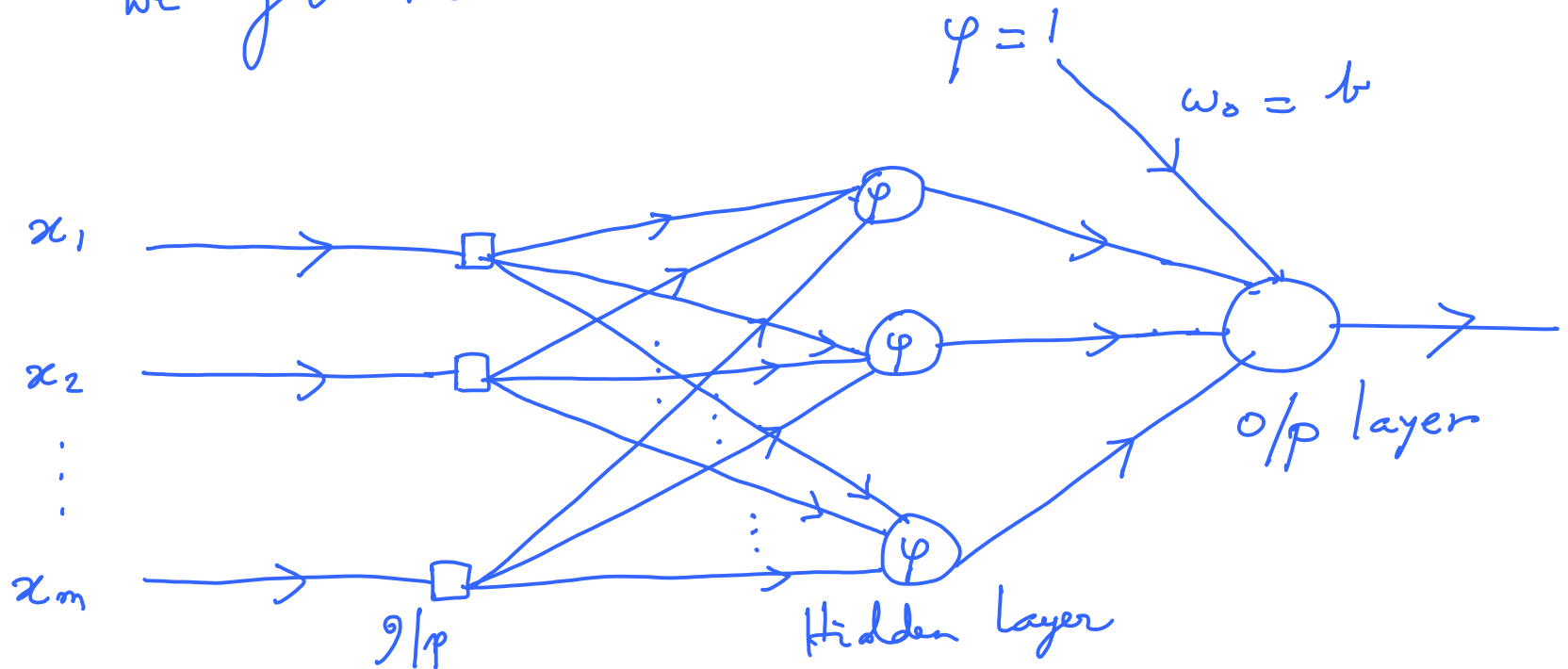
usual weight

# Regularization Networks

The idea of Green's $\text{fns}$ $G(\underline{x}, \underline{x}_i)$ centered @ $\underline{x}_i$ gives us a feel of the $n/w$ structure.

1) One hidden unit for each data point $\underline{x}_i$ $i = 1, \dots, N$. The o/p of the hidden unit is $G(\underline{x}, \underline{x}_i)$.

2) The o/p of the $n/w$ is $F(\underline{x})$ by combining the Green's functions

Input Layer

Hidden layer of
N Green's functions

O/p Layer

For the $i$th hidden unit, the O/p is $G(\underline{x}, \underline{x}_i)$

By imposing certain constraints such as ( +ve definite property )
and making $G(\cdot)$ to be rotationally invariant,
we get the Gaussian form used in RBF n/ws.

3 desirable properties for regularization n/ws from approximation theory perspective

1) It is a <u>universal approximator</u>; approx. any multivariate continuous fn very well.

2) Since the approx. scheme is derived from regularization theory & <u>linear in</u> the unknown coeffts, the unknown non-linear function can be always be approx. through an appropriate choice of the coeffts.

3) The soln. computed by a regularization n/w is optimal, and based on minimizing a functional