We can revisit the multivariate interpolation problem in a higher-dimensional space.

PROBLEM : Given a set of $N$ different points
$$\{\underline{x}_i \in \mathbb{R}^{m_0}\} \; i = 1, 2, \ldots, N \quad \text{and a}$$
corresponding set of $N$ real nos $\{d_i \in \mathbb{R}\}_{i=1,\ldots,N}$
find a function $F : \mathbb{R}^N \to \mathbb{R}^1$/
$$F(\underline{x}_i) = d_i \; ; \; i = 1, \ldots, N$$

The idea of radial basis functions can help! (stemming from the Gaussian hidden units we saw in the XOR problem)

Suppose $F(\underline{x}) = \sum_{i=1}^{N} w_i \, \varphi\left(\|\underline{x} - \underline{x}_i\|\right)$,

$\|\cdot\|$ is the $L_2$-norm and $\varphi(\cdot)$ is a set of $N$ arbitrary 'smooth' non-linear functions. $\left(\begin{array}{l} L_2 - \text{norm is a} \\ \text{reason for the} \\ \text{radial symmetry} \end{array}\right)$

Let $\varphi_{ji} = \varphi\left(\|\underline{x}_j - \underline{x}_i\|\right)$, $j, i = 1, \ldots, N$

Let $\underline{d} = \begin{bmatrix} d_1 & \cdots & d_N \end{bmatrix}^T$ (desired)

$\underline{w} = \begin{bmatrix} w_1 & \cdots & w_N \end{bmatrix}^T$ (linear weight)

Let us form a matrix eqn under the
interpolation constraints.

$$\begin{bmatrix} \varphi_{11} & \varphi_{12} & \cdots & \varphi_{1N} \\ \vdots & & & \\ \varphi_{N1} & \varphi_{N2} & & \varphi_{NN} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_N \end{bmatrix} = \begin{bmatrix} d_1 \\ \vdots \\ d_N \end{bmatrix}$$

data point — basis

$$\phi := [\phi_{ji}] \quad j, i = 1, \ldots, N$$

$$\phi \, \underline{w} = \underline{d} \, ; \quad \underline{w} = \phi^{-1} \underline{d}$$

(existence of inverse?)

# Micchelli's Theorem :

Let $\{\underline{x}_i\}_{i=1}^{N}$ be a set of distinct points in $\mathbb{R}^{m_0}$. Then the $N \times N$ interpolation matrix $\Phi$ whose $j, i$ th element is $\varphi_{ji} = \varphi(\|\underline{x}_j - \underline{x}_i\|)$ is non singular.

Plenty of such functions $\varphi(\cdot)$ exist

## Examples:

1) **Multiquadrics:**
$$\varphi(r) = (r^2 + c^2)^{\frac{1}{2}} \qquad c > 0; \quad r \in \mathbb{R}$$

2) **Inverse multiquadrics**
$$\varphi(r) = (r^2 + c^2)^{-1/2} \qquad c > 0; \quad r \in \mathbb{R}$$

3) **Gaussian functions**
$$\varphi(r) = \exp\left(-\frac{r^2}{2\sigma^2}\right) \qquad \sigma > 0; \quad r \in \mathbb{R}$$

Inv. multiquadrics & Gaussians are __localized__ i.e., $\varphi(r) \xrightarrow[r \to \infty]{} 0$

Multiquadrics is __unbounded__ as $r \to \infty$

# Radial Basis Function Networks

## Ingredients

1) Input Layer: Consists of $m_0$ source nodes, $m_0$ is the dimensionality of the input vector $\underline{x}$.

2) Hidden layer: Consists of the same # of computational units as the size of the training samples $N$; each unit is mathematically described by a radial basis function

$$\varphi_j(\underline{x}) = \varphi\left(\|\underline{x} - \underline{x_j}\|\right); \quad j = 1, 2, \ldots, N$$

$j$th point defines the center of the radial basis function
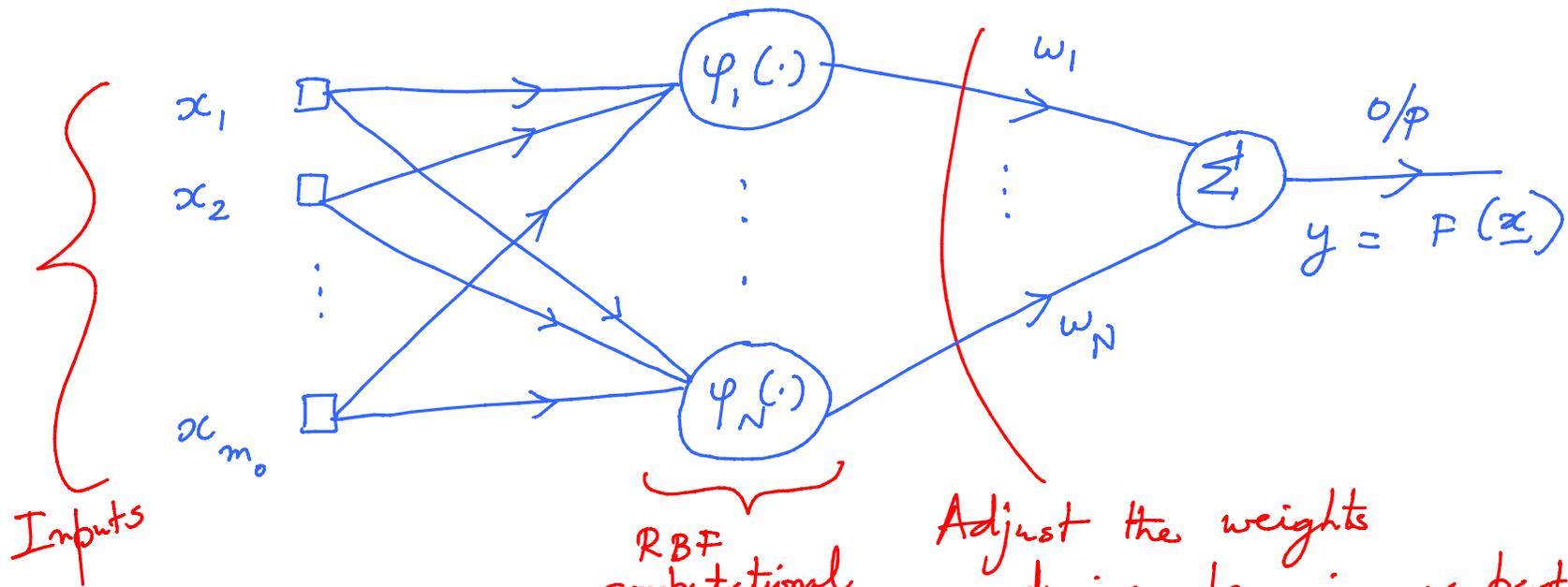
There are 'no' weights connecting source to hidden nodes

3) **O/p Layer**

This is a single computational unit. There is no restriction
(typically)
on the size of the o/p layer; o/p layer << hidden layer
size                      size

$$\varphi_j(\underline{x}) = \varphi\left(\|\underline{x} - \underline{x}_j\|\right)$$

$$= \exp\left(-\frac{1}{2\sigma_j^2} \|\underline{x} - \underline{x}_j\|^2\right); \quad j = 1, 2, \ldots, N$$

# Sketch of the RBF n/w architecture



Inputs

$x_1$

$x_2$

$\vdots$

$x_{m_0}$

$\varphi_1(\cdot)$

$\varphi_N(\cdot)$

$\Sigma$

$w_1$

$w_N$

o/p

$y = F(x)$

RBF computational units

Adjust the weights during learning as part of optimization

This is not like a traditional neuron with an activation function over a local receptive field as in a MLP

# Hybrid learning procedure for RBFs

## Idea:

1) Obtain the hidden layer elements i.e., the centers in each of the computational units through a clustering algorithm. (We do not need every data point to be a center of an RBF unit)

2) Solve for the optimal weight $\underline{w}$ linking the hidden layer and the o/p layer

# Sketch of the algo

**I/p layer:** The size of the i/p layer is determined by the dim. of the i/p vector $\underline{x}$, say $m_0$.

**Hidden Layer:** The size of the hidden layer $m_1$ is determined by the # clusters which is a tradeoff between performance & complexity

1) Using an algorithm such as the $K-$ means, obtain the cluster mean $\{\hat{\underline{\mu}}_j\}_{j=1}^{K}$ based on the inputs $\{\underline{x}_i\}_{i=1}^{N}$. (Compute the $\varphi_j(\cdot)$ using these <u>means</u> for all the data points. ($\varphi_j$ is based on the Gaussian)

Typically, the same $\sigma$ is applied to all Gaussians

2) $\sigma \sim \dfrac{d_{max}}{\sqrt{2K}}$ where $d_{max} := \max_{i,j} \|\hat{\underline{\mu}}_i - \hat{\underline{\mu}}_j\|$

(empirical rule)

The above choice of $\sigma$ ensures that individual Gaussians are not too `peaky` or `flat`. (Heuristic).

3)  After we obtain the hidden layer,

obtain $\phi(\underline{x}_i) = \begin{bmatrix} \varphi(\underline{x}_i, \hat{\underline{\mu}}_1) \\ \vdots \\ \varphi(\underline{x}_i, \hat{\underline{\mu}}_k) \end{bmatrix}$  over all $i = 1 \text{ to } N$

From $\left\{ \left( \phi(\underline{x}_i), d_i \right) \right\}_{i=1}^{N}$ , obtain $\hat{\underline{w}} = \begin{bmatrix} w_1 \cdots w_k \end{bmatrix}^T$

by solving $\boxed{\phi \, \hat{\underline{w}} = \underline{d}}$   $\left( \begin{array}{l} \text{If } \phi \text{ is a rectangular} \\ \text{matrix, obtain the} \\ \text{pseudo-inverse} \end{array} \right)$

Typically, the direct solution of $\hat{\underline{w}}$ using $\phi^{-1} d$ is
avoided by using adaptive algorithms

Following our earlier notation,

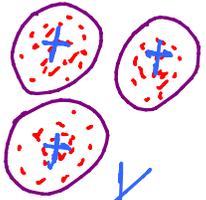$$\phi^T(\underline{x}_i) = \left[ \varphi(\underline{x}_i, \underline{\mu}_1), \cdots , \varphi(\underline{x}_i, \underline{\mu}_K) \right]$$

$$\therefore \quad \phi^T(\underline{x}_i) \, \underline{w} = d_i \; ; \text{ where } \underline{w} = \left[ w_1 \cdots w_K \right]^T \quad \textcircled{1}$$

Premultiply $\textcircled{1}$ by $\phi(\underline{x}_i)$ on both sides and
Sum up from $i = 1$ to $N$ $\qquad \textcircled{2}$

$$\underbrace{\sum_{i=1}^{N} \phi(\underline{x}_i) \, \phi^T(\underline{x}_i) \, \underline{w}}_{R} = \underbrace{\sum_{i=1}^{N} \phi(\underline{x}_i) \, d_i}_{\underline{r}}$$

$$R \, \underline{w} = \underline{r} \; ; \quad \underline{w} = R^{-1} \underline{r} \qquad \textcircled{3}$$

# K - means clustering

K - means is a simple idea for clustering. It is <u>unsupervised</u> in nature.

GOAL : Given a set of $N$ observations $\{\underline{x}_i\}_{i=1}^N$, find a Code book $C$ that assigns these observations to $K$ clusters in such a way that the <u>average measure of distortion</u> is minimized from the $K$ cluster mean

$$ J(C) = \sum_{j=1}^{K} \sum_{C(i)=j} \| \underline{x}_i - \hat{\underline{\mu}}_j \|^2 ; \quad i = 1, \ldots, N $$

Estimated mean for cluster 'j'

Except for a scaling factor of $N_j$ where $N_j$ is the # data points $\in$ cluster $'j'$,

$J(C)$ is a measure of overall cluster variance.

Now, how do we minimize $J(C)$

Approach: Use the familiar gradient descent approach

**Step 1:** For a given code book $C$, the total cluster variance is minimized w.r.t the assigned set of cluster means $\{\hat{\mu}_j\}_{j=1}^K$

i.e., $\min_{\{\hat{\mu}_j\}_{j=1}^K} \sum_{j=1}^K \sum_{\substack{C(i)=j \\ i=1,\ldots,N}} \|\underline{x}_i - \hat{\underline{\mu}}_j\|^2$

**Step 2:** Having computed $\{\hat{\underline{\mu}}_j\}_{j=1}^K$ in Step 1, optimize the encoder as

$$C(i) = \arg\min_{1 \leq j \leq K} \|\underline{x}_i - \hat{\underline{\mu}}_j\|^2$$

Iterate Steps 1 and 2 until <u>convergence</u>

# Recursive Least Squares Algo

Recall that $R \underline{w} = \underline{r}$

Suppose we are estimating $w$ i.e., $\hat{w}$ as a function of data points in an online manner

$$R(n) = \sum_{i=1}^{n} \phi(\underline{x}_i) \phi^T(\underline{x}_i)$$

where $\phi(\underline{x}_i) = \left[ \varphi(\underline{x}_i, \underline{\mu}_1), \quad \cdots \quad \varphi(\underline{x}_i, \underline{\mu}_K) \right]^T$

$$\varphi(\underline{x}_i, \underline{\mu}_j) = \exp\left( -\frac{1}{2\sigma_j^2} \| \underline{x}_i - \underline{\mu}_j \|^2 \right)$$

$$j = 1, \cdots, K$$

$K$: # clusters

The $K \times 1$ cross correlation vector is

$$\underline{r}(n) = \sum_{i=1}^{n} \phi(\underline{x}_i) \, d(i)$$

$\phi(\underline{x}_i)$ ← hidden response

$d(i)$ → desired response at the o/p of RBF n/w

$\underline{\hat{w}}(n)$ needs to be optimized in the <u>least square</u> sense

(Why: $R^{-1} \sim O(K^3)$ complex

It is computationally difficult for large $K$

i.e., We need an algo to overcome the inversion issue towards computational efficiency

$$\underline{r}(n) = \sum_{i=1}^{n-1} \phi(\underline{x}_i) d(i) + \phi(\underline{x}_n) d(n)$$

$$\underbrace{\qquad\qquad}_{\underline{r}(n-1)} + \phi(\underline{x}_n) d(n)$$

$$\underline{r}(n) =$$

$$= R(n-1) \hat{\underline{w}}(n-1) + \phi(\underline{x}_n) d(n)$$

Let us expand $R(n-1)$ further; Let us subsume $\underline{x}_n$ in $\phi(\underline{x}_n)$ as $\phi(n)$

$$\underline{r}(n) = \left[ R(n-1) + \phi(n) \underline{\phi}^T(n) \right] \hat{\underline{w}}(n-1)$$
$$\qquad\qquad + \phi(n) \left[ d(n) - \underline{\phi}^T(n) \hat{\underline{w}}(n-1) \right]$$

i.e., adding & subtracting $\phi(n) \phi^T(n)$ $\hat{\underline{w}}(n-1)$

$$R(n) = R(n-1) + \underline{\phi}(n) \underline{\phi}^T(n)$$

Let $\alpha(n) \overset{\Delta}{=} d(n) - \underline{\phi}^T_{(n)} \underline{\hat{w}}(n-1)$

$$= d(n) - \underline{\hat{w}}^T(n-1) \underline{\phi}(n)$$

$\alpha(n)$ is referred to as 'innovation'
Prior estimation error based on old $\hat{w}(n-1)$

$$\underline{r}(n) = R(n) \underline{\hat{w}}(n-1) + \underline{\phi}(n) \alpha(n)$$

$$R(n) \underline{\hat{w}}(n) = R(n) \underline{\hat{w}}(n-1) + \underline{\phi}(n) \alpha(n) \quad \textcircled{I}$$

$$\underline{\hat{w}}(n) = \underline{\hat{w}}(n-1) + R^{-1}(n) \underline{\phi}(n) \alpha(n)$$
Update rule

To solve for $R^{-1}$ appearing in the update rule, we invoke the matrix inversion lemma

Consider the matrix $A = B^{-1} + CDC^T$

$$A^{-1} = B - BC\left(D + C^T B C\right)^{-1} C^T B^T$$

For our set up

$A = R(n)$

$C = \underline{\phi}(n)$

$B^{-1} = R(n-1)$

$D = 1$

Plugging in,

$$R^{-1}(n) = R^{-1}(n-1) - \frac{R^{-1}(n-1)\,\phi(n)}{\left(1 + \phi^T(n)\,R^{-1}(n-1)\,\phi(n)\right)} \phi^T(n)\,R^{-1^T}(n-1)$$

$$= R^{-1}(n-1) - \frac{R^{-1}(n-1)\,\phi(n)\,\phi^T(n)\,R^{-1^T}(n-1)}{1 + \phi^T(n)\,R^{-1}(n-1)\,\phi(n)}$$

Let us define $P(n) \triangleq R^{-1}(n)$

$$R^T(n) = R(n) ; \quad R^{-1^T}(n-1) = R^{-1}(n-1) = P(n-1)$$

$$\text{Let} \qquad \underline{g}(n) \ \triangleq \ R^{-1}(n)\ \underline{\phi}(n)$$

gain vector

$$\underline{g}(n) \ = \ P(n)\ \underline{\phi}(n)$$

$$\underline{\hat{w}}(n) \ = \ \underline{\hat{w}}(n-1) \ + \ \underbrace{\underline{g}(n)\ \alpha(n)}_{\text{gain} \times \text{innovation}}$$

# Summary of the RLS Algo

Given $\{\phi(i), d(i)\}_{i=1}^{N}$, do the following

for $n = 1, \ldots, N$

$$P(n) = P(n-1) - \frac{P(n-1)\,\underline{\phi}(n)\,\underline{\phi}^T(n)\,P(n-1)}{1 + \underline{\phi}^T(n)\,P(n-1)\,\underline{\phi}(n)}$$

$$\underline{g}(n) = P(n)\,\underline{\phi}(n)$$

$$\alpha(n) = d(n) - \hat{\underline{w}}(n-1)\,\underline{\phi}(n)$$

$$\hat{\underline{w}}(n) = \hat{\underline{w}}(n-1) + \underline{g}(n)\,\underbrace{\alpha(n)}_{\text{scalar}}$$

To initialize,

Set $\quad \hat{w}(0) = \underline{0}$

$\overline{P(0)} = \chi^{-1} I \qquad$ small +ve constant.

# Comparison of RBFs and MLPs

## Similarities

1) Both are non-linear layered feed forward networks

2) Both are universal approximators

## Differences

1) RBF in basic form has a <u>single hidden layer</u>, where as, an MLP has <u>more than 1 hidden layer</u>. The neuronal model is the same in the MLP.

2) For RBF, each unit can have a <u>different $\mu$</u>. Computations in the hidden layer in a MLP require local gradients. This is not so in a RBF

3) For RBF, hidden layer is non-linear, but o/p layer is linear. However, both hidden and o/p layers of MLP are non-linear

4) The argument of activation fn in MLP, involves inner product i.e., $\varphi(\underline{w}^T\underline{x} + b)$. In case of RBF, one looks into the Euclidean norm i.e., $\varphi(||\underline{x}_i - \underline{\mu}_j||)$ for the $j^{th}$ RBF unit

5) Given the same level of n/w complexity, MLP (?) Could provide better accuracy than RBF RBF is faster than MLP.

# Kernel Regression

**Motivation** : Can we link RBFs to solve the regression problem?

Let us revisit the kernel regression idea built on density estimation

$$y_i = f(\underline{x}_i) + \varepsilon_i \; ; \; i = 1, \ldots, N$$
$$f(\cdot) \text{ is unknown}$$

**Qn:** What is a reasonable estimate of $f(\cdot)$?

If we look into the mean of the observables around a point $\underline{x}$ i.e., confine the observations in a small neighborhood around $\underline{x}$, we can form an estimate for $f(\underline{x})$

$$f(\underline{x}) = E\left(y \mid \underline{x}\right) \qquad \text{(conditional mean)}$$

$$= \int_{-\infty}^{\infty} y \; P_Y\left(y \mid \underline{x}\right) \, dy$$

$$P_Y\left(y \mid \underline{x}\right) = \frac{P_{\underline{X}Y}\left(\underline{x}, y\right)}{P_{\underline{X}}\left(\underline{x}\right)}$$

joint p.d.f of $\underline{x}, Y$

marginal of $\underline{x}$

regression function

$$f(\underline{x}) = \frac{\int_{-\infty}^{\infty} y \; P_{\underline{X}Y}(\underline{x}, y) \, dy}{P_{\underline{X}}(\underline{x})} \qquad \text{I}$$

From $P_Y(y | \underline{x})$

A few points to note

1) Joint density $P_{\underline{X}Y}(\underline{x}, y)$ is unknown

2) We may need a non parametric estimate

Typically a kernel defined by $K(\underline{x})$ has properties similar to a prob. density function (pdf)

1) Kernel $K(\underline{x})$ is continuous, bounded; and a real function of $\underline{x}$ symmetric about the origin where it attains a max. value e.g., Gaussian kernel

2) Volume under the kernel is unity

$$\int_{\mathbb{R}^m} K(\underline{x})\, d\underline{x} = 1$$

(Normalization)

Assuming $\underline{x}_1, \underline{x}_2, \cdots, \underline{x}_N$ are independent and identically distributed random vectors, the PARZEN ROZENBLATT density estimate of $P_{\underline{x}}(\underline{x})$ is

$$\hat{P}_{\underline{x}}(\underline{x}) = \frac{1}{N h^{m_o}} \sum_{i=1}^{N} K\left(\frac{\underline{x} - \underline{x}_i}{h}\right)$$

estimate

data point

$\underline{x} \in \mathbb{R}^{m_o}$

Controls the size (bandwidth)

$\boxed{\underline{\mathbb{I}}}$

# PROPERTY (BIAS)

If $h = h(N)$ is a function such that

$$\lim_{N \to \infty} h(N) = 0,$$ then

$$\lim_{N \to \infty} E\left[\hat{P}_{\underline{x}}(\underline{x})\right] = P_{\underline{x}}(\underline{x}) \quad \left(\text{Asymptotically unbiased}\right)$$

Let us formulate the Parzen-Rosenblatt density estimate for the joint pdf $P_{\underline{X}Y}(\underline{x}, y)$; assuming $(\underline{x}, y)$ pairs are independent and identically distributed

$$\hat{P}_{\underline{X}y}(\underline{x}, y) = \frac{1}{N h^{m_o+1}} \sum_{i=1}^{N} K\left(\frac{\underline{x} - \underline{x}_i}{h}\right) K\left(\frac{y - y_i}{h}\right)$$

$$\underline{x} \in \mathbb{R}^{m_o}$$
$$y \in \mathbb{R}$$

estimate of the joint density

Consider the numerator of $\boxed{I}$
which can be simplified as

$$\underbrace{\int_{-\infty}^{\infty} y \, \widehat{P}_{\underline{x}, y}(\underline{x}, y) \, dy}_{} = \frac{1}{N h^{m_o + 1}} \sum_{i=1}^{N} K\left(\frac{\underline{x} - \underline{x}_i}{h}\right) \underbrace{\int_{-\infty}^{\infty} y \, K\left(\frac{y - y_i}{h}\right) dy}_{}$$

<span style="color:red">Let us compute the integral</span>

We need to compute the integral carefully

Consider $\displaystyle\int_{-\infty}^{\infty} y\, K\left(\frac{y-y_i}{h}\right) dy$

Let $z = (y - y_i)/h$     (Change of variable)

$$y = y_i + zh \;\; ; \;\; dy = h\, dz$$

$$\therefore \int_{-\infty}^{\infty} \left(y_i + zh\right) K(z)\, h\, dz = h\left[\underbrace{\int_{-\infty}^{\infty} y_i\, K(z)\, dz}_{\text{Term 1}} + \underbrace{\int_{-\infty}^{\infty} zh\, K(z)\, dz}_{\text{Term 2}}\right]$$

Term 1 evaluates to $y_i$

since $\int_{-\infty}^{\infty} K(z) = 1$ <span style="color:red">( Normalization )</span>

Term 2 evaluate to 0 since $\int_{-\infty}^{\infty} z\, K(z)\, dz = 0$

$\underbrace{\qquad\qquad}_{\text{Zero mean by assumption}}$

$\therefore \int_{-\infty}^{\infty} y\, \hat{P}_{\underline{x}y}(\underline{x}, y)\, dy = \dfrac{1}{N \underbrace{h^{m_0+1}}_{h^{m_0}}} \cdot h \sum_{i=1}^{N} y_i\, K\left(\dfrac{\underline{x} - \underline{x}_i}{h}\right) \quad \boxed{111}$

$$\therefore \quad \hat{f}_{reg} = E\left(y \mid \underline{x}\right) = \frac{\int_{-\infty}^{\infty} y \, \hat{P}_{\underline{x}\, y}\left(\underline{x}, y\right) dy}{P_x\left(x\right)}$$

<span style="color:red">Recall!</span> Ⓘ

Using Ⅱ and Ⅲ in Ⓘ , we have the <span style="color:red">Kernel reg. estimator</span>

$$\hat{f}_{reg}\left(\underline{x}\right) = \frac{\sum_{i=1}^{N} y_i \, K\left(\frac{\underline{x} - \underline{x}_i}{h}\right)}{\sum_{i=1}^{N} K\left(\frac{\underline{x} - \underline{x}_i}{h}\right)}$$

NOTE:
Denominator is "not" zero.
Pander why?

<span style="color:red">↖ compact form</span>

# Nadaraya Watson Regression Estimator

Let us define the normalized weighting function

$$W_{N,i}(\underline{x}) = \frac{K\left(\frac{x - x_i}{h}\right)}{\sum_{j=1}^{N} K\left(\frac{x - x_j}{h}\right)}$$

$\left(\begin{array}{c} \text{weight to} \\ \text{the 'i' th} \\ \text{data point } \underline{x}_i \end{array}\right)$

$$\sum_{i=1}^{N} W_{N,i}(\underline{x}) = 1 \qquad \forall\ \underline{x}$$

$$\hat{f}_{reg}(\underline{x}) = \sum_{i=1}^{N} w_{N,i}(\underline{x}) \, y_i$$

Regression fn

observable

weight depends on $\underline{x}_i$

Weighted average of $y$- observables.

# Link to RBF n/w

Since we assume spherical symmetry for the kernel in RBFs ( Gaussian case )

$$K\left(\frac{\underline{x} - \underline{x}_i}{h}\right) = K\left(\frac{\|\underline{x} - \underline{x}_i\|}{h}\right) \quad \forall i$$

$L_2$ norm (·) radial symmetry

Define $\psi_N(\underline{x}, \underline{x}_i) = \dfrac{K\left(\|\frac{\underline{x} - \underline{x}_i}{h}\|\right)}{\sum\limits_{j=1}^{N} K\left(\frac{\|\underline{x} - \underline{x}_j\|}{h}\right)}$ ; $i = 1, \ldots, N$

weighting fn

Again $\sum_{i=1}^{N} \psi_N (\underline{x}, \underline{x}_i) = 1 \quad \forall \underline{x}$

The regression estimate is a <u>weighted sum</u> of '$N$' <u>basis functions</u> $\psi_N (\underline{x}, \underline{x}_i)$

Let $y_i = w_i$ for $i = 1, 2, \dots, N$ (weights are simply observables)

$$\hat{f}_{reg}(\underline{x}) = \sum_{i=1}^{N} w_i \psi_N (\underline{x}, \underline{x}_i)$$

weight

R B basis functions

(A)

NOTE :

1) $(A)$ denotes the input/output mapping of a normalized RBF with $\quad 0 \leq \psi_N(\underline{x}, \underline{x}_i) \leq 1 \qquad \forall \underline{x}, \underline{x}_i$

2) $\psi_N(\underline{x}, \underline{x}_i)$ is interpreted as the prob. of an event described by input vector $\underline{x}$ conditioned on $\underline{x}_i$

3) Density est. can be ill-posed & can be made well-posed by regularization.

Kernel functions can be of various forms

E.g. $K(\underline{x}) = \dfrac{1}{\sqrt{(2\pi)}^{m_0}} \exp\left(-\dfrac{\|\underline{x}\|^2}{2}\right)$ ← $\sigma^2 = 1$

Gaussian kernel

With a spread parameter $\sigma$,

$K\left(\dfrac{\underline{x} - \underline{x}_i}{h}\right) = \dfrac{1}{(2\pi\sigma^2)^{m_0/2}} \exp\left(-\dfrac{\left\|\dfrac{\underline{x} - \underline{x}_i}{h}\right\|^2}{2\sigma^2}\right)$

$i = 1, \ldots, N$

With $h = 1$, following NWRE

$$\hat{f}_{reg}(\underline{x}) = \frac{\sum\limits_{i=1}^{N} y_i \exp\left(-\frac{\|\underline{z} - \underline{x}_i\|^2}{2\sigma^2}\right)}{\sum\limits_{j=1}^{N} \exp\left(-\frac{\|\underline{x} - \underline{x}_j\|^2}{2\sigma^2}\right)}$$

Final form using Gaussian Fns.

# Basics of constrained optimization

**Motivation:** We are interested in minimizing/maximizing functions subject to constraints over the variable

**Example:**
1) Optimize the path from point A to point B, subject to traffic conditions over all connecting roads.

2) Optimize the square error in the MLP subject to a sparsity constraint in the network connectivity across layers etc

We have to deal with minimizing functions subject to equality and inequality constraints

Formulation

$$\min_{x \in \mathbb{R}^n} f(x) \quad s.t. \quad \begin{cases} c_i(x) = 0 \\ \quad i \in \text{Equality constraints 'E'} \\ c_i(x) \geq 0 \\ \quad i \in \text{Inequality constraints 'I'} \end{cases}$$

objective function

$c_i(x) \leq 0$

NOTE : The ' $\leq$ ' in the constraint inequality can be suitably transformed to ' $\geq$ ' by a sign flip

# Assumptions / terminology

1) We assume $f$ and $c_i$'s to be <u>smooth</u> and <u>real</u> <u>valued</u> <u>operating</u> on a subset of $\mathbb{R}^n$.

2) $f$ is the objective function and $c_i \in E$, $c_j \in I$

Let $\Omega = \left\{ \underline{x} \;\middle|\; c_i(\underline{x}) = 0, \; c_i \in E, \; c_j(\underline{x}) \geqslant 0, j \in I \right\}$

<span style="color:red">(overall constraint set)</span>

Written Succinctly,

$$\min_{\underline{x} \in \Omega} f(\underline{x}) \quad \text{——————} \quad \boxed{I}$$

Recall : For an optimal solution $\underline{x}^*$ (Minimization)

$$\nabla f(\underline{x}^*) = \underline{0}$$

and $\quad \nabla^2 f(\underline{x}^*) \geq 0 \qquad$ ( Positive Semi definite Property )

If $\quad \nabla^2 f(\underline{x}^*) > 0$

Strict inequality ( +ve definite )

The Solutions to opt. problems can have

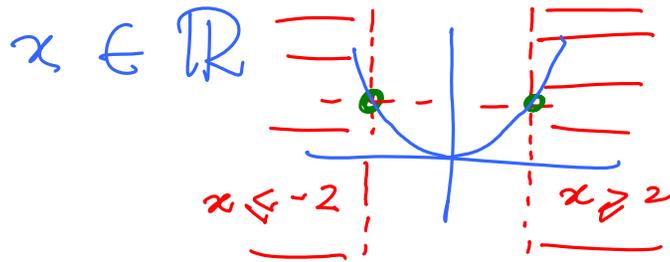(a) local Solution

(b) global Solution
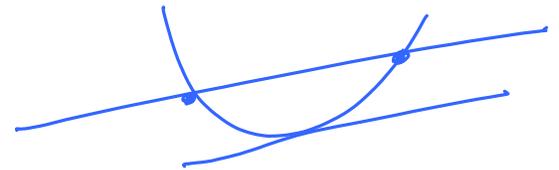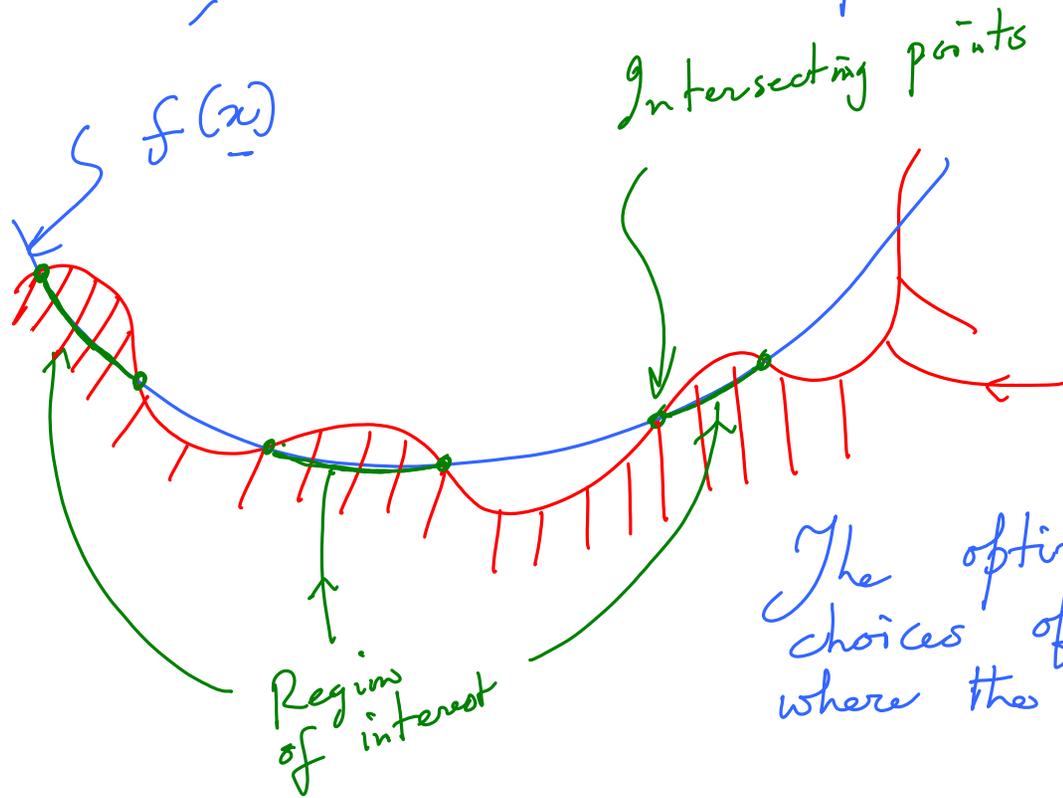

global

Example:

min $x^2$

$x \in \mathbb{R}$

s.t. $|x| \geq 2$
(constraint)

(There are $\underline{2}$ solns & $\underline{\text{not unique}}$)

$x \leq -2$   $x \geq 2$

No Constraint $\Rightarrow$ $x = 0$   trivially unique

Instead,    Consider the problem

Intersecting points

$f(x)$

Constraint
$c_i(x) \leq 0$

Region of interest

The optimization is over the choices of $x$ over the objective where the constraints are satisfied

1) A vector $x^*$ is a _local solution_ to problem $\boxed{I}$ if $x^* \in \Omega$ and there is a neighborhood $N$ of $x^*$ such that $f(x) \geq f(x^*) \quad \forall \, x \in N \cap \Omega$

2) A vector $x^*$ is a _strict local solution_ if $x^* \in \Omega$ and there is a neighborhood $N$ of $x^*$ / $f(x) > f(x^*)$ $\forall \, x \in N \cap \Omega$ with $x \neq x^*$

3) A point $x^*$ is an _isolated local solution_ if $x^* \in \Omega$ and there is a neighborhood $N$ of $x^*$ / $x^*$ is the only minimizer in $N \cap \Omega$

# Smoothness

Smoothness of objective fns & constraints can help algorithms to make better choices during gradient search towards the soln.
(local soln)

# Active/Inactive Constraints

At a feasible point $x$, the inequality constraint $c_i \in I$ is <u>active</u> if $c_i(x) = 0$ & <u>inactive</u> if $c_i(x) > 0$

$c_i(x) \geq 0$