$y(x, w)$     parametric model
       ↑
   learned    training data are thrown away after training

another class of PR techniques : nonparametric
e.g. "nearest neighbor"
    fast to train , slow at making prediction

___

Dual Representation
combinations of a kernel function evaluated at the training
data points

$$k(x, x') = \phi(x)^T \phi(x')$$

        `kernel trick'

1964   Aizerman et al.
1992   Boser et al.

Extensions of existing methods     scalar product $\longrightarrow$ kernel
`kernelize'
e.g.   Kernel PCA (1998) , Kernel Fisher Discriminant (1999)

different forms of kernel functions :
    stationary kernels    $k(x, x') = k(x - x')$
            invariant to translation in input space
    homogeneous kernels    $k(x, x') = k(\|x - x'\|)$
            radial basis functions (RBF)
        depend only on the magnitude of the distance

---

DUAL REPRESENTATIONS

Example :    Linear Regression , regularized sum-of-squares error

$$J(w) = \frac{1}{2} \sum_{n=1}^{N} \left\{ w^T \phi(x_n) - t_n \right\}^2 + \frac{\lambda}{2} w^T w$$

$$\lambda \geq 0$$

Solution for $w$ takes the form
$$w = -\frac{1}{\lambda} \sum_{n=1}^{N} \left\{ w^T \phi(x_n) - t_n \right\} \phi(x_n)$$

$$= \sum_{n=1}^{N} a_n \phi(x_n) = \overline{\Phi}^T a$$

$$a_n = -\frac{1}{\lambda} \left\{ w^T \phi(x_n) - t_n \right\} \qquad \overline{\Phi} : \begin{bmatrix} \vdots \\ \phi(x_n)^T \\ \vdots \end{bmatrix} \qquad a = (a_1, \cdots, a_N)^T$$

Dual Representation
$$J(a) = \frac{1}{2} a^T \overline{\Phi} \overline{\Phi}^T \overline{\Phi} \overline{\Phi}^T a - a^T \overline{\Phi} \overline{\Phi}^T t + \frac{1}{2} t^T t + \frac{\lambda}{2} a \overline{\Phi} \overline{\Phi}^T a$$

$$t = (t_1, \cdots, t_N)^T$$

Gram matrix    $K = \overline{\Phi} \overline{\Phi}^T$       symmetric

$$K_{nm} = \phi(x_n)^T \phi(x_m) = k(x_n, x_m)$$

$$J(a) = \frac{1}{2} a^T K K a - a K t + \frac{1}{2} t^T t + \frac{\lambda}{2} a K a$$

$$\nabla J(a) = 0$$
$$a = (K + \lambda I_N)^{-1} t \qquad coefficients$$

$$y(x) = w^T \phi(x) = a^T \overline{\Phi} \phi(x) = k(x)^T (K + \lambda I_N)^{-1} t$$

___

* avoid the explicit computation
  of feature mapping $\phi(x)$        $k(x) = \begin{pmatrix} \vdots \\ k(x_n, x) \\ \vdots \end{pmatrix}$  vector

valid  kernel          $k(x, x') = \phi(x)^T \phi(x')$

$$= \sum_{i=1}^{M} \phi_i(x) \phi_i(x')$$

$k(x, z) = (x^T z)^2$       a valid kernel?

2-dimensional case   $x = (x_1, x_2)^T$

$$k(x, z) = (x^T z)^2 = (x_1 z_1 + x_2 z_2)^2$$

$$= x_1^2 z_1^2 + 2 x_1 z_1 x_2 z_2 + x_2^2 z_2^2$$

$$= (x_1^2, \sqrt{2} x_1 x_2, x_2^2)(z_1^2, \sqrt{2} z_1 z_2, z_2^2)^T$$

$$= \phi(x)^T \phi(z)$$

so   $\phi(x) = (x_1^2, \sqrt{2} x_1 x_2, x_2^2)^T$

The  necessary  and  sufficient  condition  for  $k(x, x')$  to  be
a  valid  kernel :   Gram  matrix  $K = [k(x_n, x_m)]_{nm}$
should  be  positive  semidefinite  for  all  possible  choices  of
the  set  $\{x_n\}$

$K_{ij} = k(x_i, x_j)$     $\{x_1, \ldots, x_n\}$

K is symmetric

Let   $K = V \Lambda V^T$     $V$ is an orthogonal matrix      $V_t$

$\Lambda$  contains  the  eigenvalues      $\lambda_t$

assume $\lambda_t$  are  nonnegative.

consider  the  feature  mapping    $\phi: x_i \mapsto (\sqrt{\lambda_t} v_{ti})_{t=1}^n \in \mathbb{R}^n$

$i = 1, \ldots, n$

$\phi(x_i)^T \phi(x_j) = \sum_{t=1}^{n} \lambda_t v_{ti} v_{tj} = (V \Lambda V^T)_{ij} = K_{ij} = k(x_i, x_j)$

therefore  we  find  a  feature  mapping  $\phi$  for  the  kernel

requirement  of  $\lambda$  being nonnegative  is  necessary:  if $\lambda_s < 0$, eigenvector $v_s$

we have $z = \sum_{i=1}^{n} v_{si} \phi(x_i) = \sqrt{\Lambda} V^T v_s$   $\Rightarrow$   $\|z\|^2 = v_s^T V \Lambda V^T v_s = v_s^T K v_s = \lambda_s < 0$

---

assume $k_1$ is a valid kernel function

P. 296      $k(x, x') = f(x) k_1(x, x') f(x')$   (6.14)

$k(x, x') = \exp(k_1(x, x'))$      (6.16)

Proof of  (6.14):   for  all  $\alpha_i, \alpha_j \in \mathbb{R}$

$$\sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j k(x_i, x_j) = \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j f(x_i) k_1(x_i, x_j) f(x_j)$$

$$= \sum_{i=1}^{n} \alpha_i f(x_i) \phi^T(x_i) \sum_{j=1}^{n} \alpha_j \phi(x_j) f(x_j)$$

$$= \Psi^T \Psi \geq 0        \text{positive  semidefinite}$$

$$\Psi = \sum_{i=1}^{n} \alpha_i f(x_i) \phi(x_i)$$

other  popular  kernels :

polynomial  kernels     $k(x, x') = (x^T x' + c)^M$

Gaussian  kernels     $k(x, x') = \exp(-\|x - x'\|^2 / 2\sigma^2)$

$k(x, x') = \exp(-\|x - x'\|^2 / 2\sigma^2)$

$$= \exp(-x^T x / 2\sigma^2) \exp(x^T x' / \sigma^2) \exp(-(x')^T x' / 2\sigma^2)$$

so  by  (6.14) (6.16)  $\Rightarrow$   valid  kernels

another  property  of  kernels

$$k(x, x')^2 = (\phi(x)^T \phi(x'))^2 \leq \|\phi(x)\|^2 \|\phi(x')\|^2$$

$$= k(x, x) k(x', x')$$

## Other types of Kernels

generalized Gaussian kernels :

$$k(x, x') = \exp\left\{-\frac{1}{2\sigma^2}\left(K(x,x) + K(x',x') - 2K(x,x')\right)\right\}$$

$K(x, x')$ a nonlinear kernel

set kernels :

$A_1, A_2$ are two subsets $\quad A_1, A_2 \subseteq A$

$$k(A_1, A_2) = 2^{|A_1 \cap A_2|} \quad \text{is a valid kernel}$$

encoding :

$$\phi(A_1) = \begin{pmatrix} 1 \\ 1 \\ 0 \\ 0 \\ \vdots \end{pmatrix}$$

$2^{|A_1|}$ bits that index the subsets of A

$|A_1|$ is the number of elements
$2^{|A_1|}$ is the number of subsets

### Generative + Discriminative

① probability kernels :  $\quad k(x, x') = p(x) p(x')$

or $\quad k(x, x') = \sum_i p(x|i) p(x'|i) p(i)$

or $\quad k(x, x') = \int p(x|\mathbb{Z}) p(x'|\mathbb{Z}) p(\mathbb{Z}) \, d\mathbb{Z}$

↑
latent
variable

② Fisher kernel

Fisher score  $g(\theta, x) = \nabla_\theta \ln p(x|\theta)$

Fisher kernel  $k(x, x') = g(\theta, x)^T F^{-1} g(\theta, x')$

$$F = E_x\left[g(\theta, x) g(\theta, x)^T\right] \quad \text{is the Fisher information}$$

---

continued.  Fisher kernel

consider a Gaussian distribution $N(x|\mu, S)$

$$g(\mu, x) = \nabla_\mu \ln N(x|\mu S) = S^{-1}(x-\mu)$$

$$F = E_x\left[g(\mu, x) g(\mu, x)^T\right] = S^{-1} E_x\left[(x-\mu)(x-\mu)^T\right] S^{-T}$$
$$= S^{-1} S S^{-T}$$
$$= S^{-1}$$

$$k(x, x') = g(\theta, x)^T F^{-1} g(\theta, x)$$
$$= (x-\mu)^T S^{-T} \left(S^{-1}\right)^{-1} S^{-1}(x'-\mu)$$
$$= (x-\mu)^T S^{-1} (x'-\mu) \quad \text{(Mahalanobis)}$$

---

### Radial Basis Function

$$f(x) = \sum_{n=1}^{N} w_n h(\|x - x_n\|) \quad \text{centered on each data point}$$

Derivation :

consider noisy input , the error function is defined by
$$E = \frac{1}{2} \sum_{n=1}^{N} \int \{y(x_n+\xi) - t_n\}^2 \nu(\xi) \, d\xi$$

$\nu$ is isotropic

to find  $y(x)$  that minimizes  $E$
we need to apply calculus of variations
perturb  $y(x) \quad \Rightarrow \quad y(x) + \epsilon \eta(x)$

$$E(y+\epsilon\eta) = \frac{1}{2} \sum_{n=1}^{N} \int \{y(x_n+\xi) + \epsilon\eta(x+\xi) - t_n\}^2 \nu(\xi) \, d\xi$$
$$= E(y) + \epsilon \underbrace{\sum_{n=1}^{N} \int \{y(x_n+\xi) - t_n\} \nu(\xi) \eta(x+\xi) \, d\xi}_{= 0} + O(\epsilon^2)$$

$\epsilon \to 0$

At the optimum, $E$ should be stationary for arbitrary $\eta$

choose $\eta(x) = \delta(x - z)$

$$\sum_{n=1}^{N} \int \{ y(x_n + \xi) - t_n \} \, \delta(x_n + \xi - z) \, \nu(\xi) \, d\xi$$

$$= \sum_{n=1}^{N} \{ y(z) - t_n \} \, \nu(z - x_n) = 0$$

$$\left| \begin{array}{l} x_n + \xi - z = 0 \\ \xi = z - x_n \end{array} \right.$$

$$y(z) = \frac{t_n \, \nu(z - x_n)}{\sum_{n=1}^{N} \nu(z - x_n)}$$

$$\Rightarrow \qquad y(x) = \sum_{n=1}^{N} t_n \, h(x - x_n)$$

$$h(x - x_n) = \frac{\nu(x - x_n)}{\sum_{n=1}^{N} \nu(x - x_n)}$$

---

Nadaraya $-$ Watson model

$$\{ x_n, t_n \} \qquad p(x, t) = \frac{1}{N} \sum_{n=1}^{N} f(x - x_n, t - t_n)$$

$f$: component density function

$$y(x) = E[t \,|\, x] = \int_{-\infty}^{\infty} t \, p(t \,|\, x) \, dt$$

$$= \frac{\int t \, p(x, t) \, dt}{\int p(x, t) \, dt}$$

$$\left| \begin{array}{l} \text{from conditional} \\ \text{to joint} \end{array} \right.$$

$$= \frac{\sum_{n} \int t \, f(x - x_n, t - t_n) \, dt}{\sum_{m} f(x - x_m, t - t_m) \, dt}$$

assume $\displaystyle\int_{-\infty}^{\infty} f(x, t) \, t \, dt = 0 \qquad$ zero mean

$$\int t \, f(x - x_n, t - t_n) \, dt$$

$$= \underbrace{\int (t - t_n) \, f(x - x_n, t - t_n) \, dt}_{\| \, 0} + \int t_n \, f(x - x_n, t - t_n) \, dt$$

so after "change of variable"

$$y(x) = \frac{\sum_{n} g(x - x_n) \, t_n}{\sum_{m} g(x - x_n)}$$

$$= \sum_{n} k(x, x_n) \, t_n$$

$$g(x) = \int_{-\alpha}^{\infty} f(x, t) \, dt$$

$$p(t \,|\, x) = \frac{p(t, x)}{\int p(t, x) \, dt} = \frac{\sum_{n} f(x - x_n, t - t_n)}{\sum_{m} \int f(x - x_m, t - t_m) \, dt}$$

$f(x, t)$ zero-mean isotropic Gaussian