

Non-Parametric Methods

3/30/2008

We do not assume a specific model  $g(x)$   
 Instead we assume that similar inputs have  
 similar outputs  
 instance-based or memory-based learning.  
 All (most) instances should be stored - requires  
 O(N) storage.

Nonparametric Density Estimation

Data  $\mathcal{X} = \{x^t\}_{t=1}^N$  i.i.d. from some unknown source  $p(x)$

$\rightarrow$  Cumulative distribution function (1-D)  
 nonparametric estimate  $\hat{F}(x) = \# \{x^t \leq x\}$

$\rightarrow$  Probability  $\hat{p}(x) = \frac{1}{h} \left[ \frac{\# \{x^t \leq x+h\} - \# \{x^t \leq x\}}{N} \right]$

$h$  - length of interval (close enough)

Histogram Estimator

Input space divided into equal sized intervals named bins. Origin  $x_0$  and bin width  $h$ .

$$\hat{p}(x) = \frac{1}{Nh} \sum_{x^t \text{ in same bin as } x} 1$$

choice of bin size  $h$  affects distribution.

(2)

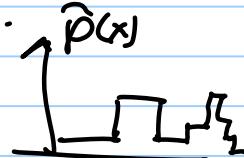
## Histogram Estimator (cont)

Naive estimator  
(no origin)

$$\hat{p}(x) = \frac{\# [x-h < x^t \leq x+h]}{2Nh}$$

This can be written as:

$$\hat{p}(x) = \frac{1}{Nh} \sum_{t=1}^N w\left(\frac{x-x^t}{h}\right)$$



weights defined by  $w(u) = 1$ ,  $|u| < 1$ ,  $= 0$ , otherwise  
|jumps at  $x^t \pm h$

## Kernel Estimator

To get a smooth estimate, we use a smooth weight function - kernel function.

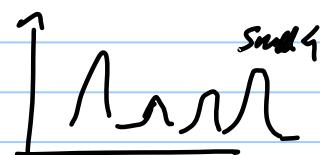
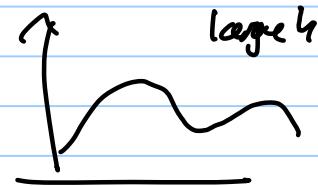
$$K(u) = \frac{1}{\sqrt{\pi}} e^{-u^2/2}$$

Kernel estimator / Parzen window

$$\hat{p}(x) = \frac{1}{Nh} \sum_{t=1}^N K\left(\frac{x-x^t}{h}\right)$$

(Can truncate  $K(u)$  for large  $|u|$ )

Again,  $h$  acts as scale.



(3)

## k-Nearest neighbour estimator

Adapt the amount of smoothing to the local density of data

Degree of smoothness is controlled by  $k$ , the number of neighbours

Distances  $|x - x_i|$  between example and

For each  $x \in X$ , define

$$d_1(x) \leq d_2(x) \leq \dots \leq d_N(x)$$

to be the distances from  $x$  to other sample points.

k-nearest neighbor is  $\hat{p}(x) = \frac{k}{2N d_k(x)}$

(Compared to Parzen -  $d_h(x) = h$  - but instead of fixing  $h$  and counting samples, we fix the no. of observations).

To get a smoother estimate, we use kernel

$$\hat{p}(x) = \frac{1}{N d_k(x)} \sum_{t=1}^N K\left(\frac{x - x_t}{d_k(x)}\right)$$

like kernel with adaptive smoothing  
 $h = d_k(x)$

(4)

## Generalization to Multivariate Data

$$X = \{x^t\}_{t=1}^N, \hat{p}(x) = \frac{1}{Nh^d} \sum_{i=1}^n K\left(\frac{x-x^i}{h}\right)$$

Typical Candidate

- multivariate Gaussian.

$$K(y) = \frac{1}{\sqrt{\pi^d}} e^{-\frac{1}{2} \|y\|_R^2}$$

$$\text{with } \int K(x) dx = 1.$$

But non-parametric methods are dangerous in high-dimensional spaces because of the curse of dimensionality.

If  $x$  is an 8-dim  
10 bins per dimension  
Then require  $10^8$  bins

For kernels, the data will typically have local preferred directions. So the Gaussian kernel should be of form  $K(u) = \frac{1}{(2\pi)^d |S|^{\frac{1}{2}}} e^{-\frac{1}{2} u^T S^{-1} u}$

(5)

## Nonparametric Classification

$$\hat{P}(\underline{x} | C_i) = \frac{1}{N_i h^d} \sum_{t=1}^N K\left(\frac{\underline{x} - \underline{x}^t}{h}\right) r_i^t$$

$r_i^t = 1, \text{ if } \underline{x}^t \in C_i$   
 $0, \text{ otherwise.}$

MLE  $\hat{P}(C_i) = N_i/N$ .  $N_i = \sum_t r_i^t$ .

Discriminant is  $g_i(\underline{x}) = \hat{P}(\underline{x} | C_i) \hat{P}(C_i)$

$$= \frac{1}{N h^d} \left\{ \sum_{t=1}^N K\left(\frac{\underline{x} - \underline{x}^t}{h}\right) r_i^t \right\}$$

For special case of  $k$ -NN estimator

$$\hat{P}(\underline{x} | C_i) = \frac{k_i}{N_i V^k(x)}$$

$V^k(x)$  volume containing the  $k$ -nearest neighbor

$$\hat{P}(C_i | \underline{x}) = \hat{P}(\underline{x} | C_i) \hat{P}(C_i) = \frac{k_i}{N}$$

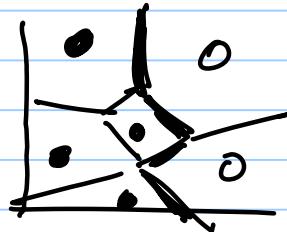
$k$ -nearest neighbor classifier assigns the input to the class having most examples among the  $k$ -nearest neighbors.

(6)

## Condensed Nearest Neighbors

decrease the number of stored instances without degrading performance

e.g. can remove those instances  $\bullet$  &  $\circ$  which are not at the boundaries



## Nonparametric Regression: Smoothing Models

regression  $X = \{x^t, r^t\}, r^t \in \mathcal{R}$ .

$$r^t = g(x^t) + \epsilon$$

Find neighbourhood of  $x$  and average the  $r$  values to calculate  $\hat{g}(x)$ .

### running mean smoother

$$\hat{g}(x) = \frac{\sum_{t=1}^n b(x, x^t) r^t}{\sum_{t=1}^n b(x, x^t)}$$

with  $b(x, x^t) = 1$ , if  $x^t$  in some bin as  $x$   
0, otherwise.

### weighted mean smoother

$$\hat{g}(x) = \frac{\sum_{t=1}^n w(\frac{x-x^t}{h}) r^t}{\sum_{t=1}^n w(\frac{x-x^t}{h})}$$

$$w(u) = 1 \quad \text{if } |u| < 1 \\ 0 \quad \text{otherwise.}$$

### Kernel smoother

$$\hat{g}(x) = \frac{\sum_t K(\frac{x-x^t}{h}) r^t}{\sum_t K(\frac{x-x^t}{h})}$$

(7)

## Choosing the Smoothing Parameter.

The choice of  $h$  (scale) or  $k$  (no. neighbor).  
Large  $h$  or  $k$  decreases variance but increases bias.

For smoothing splines

$$\sum_t [r^t - \hat{g}(x^t)]^2 + \lambda \int_a^b [\hat{g}''(x)]^2 dx$$

curvature.

Cross-validation is used to tune  $h$ ,  $k$  or  $\lambda$ .