training data

with targets    without targets

supervised      unsupervised

(probabilistic) models

parameters

training

find good
parameters

Predict

test data

Problem

classification
regression
clustering

algorithms

optimization
numerical methods

overfitting
regularization
reduce model complexity

Take a look at Appendix C

e.g.

$$(AB)^T = B^T A^T$$

$$AB\, B^{-1} A^{-1} = I \qquad (AB)^{-1} = B^{-1} A^{-1}$$

for square matrix

trace $\quad Tr(AB) = Tr(BA)$

$$\sum_i (AB)_{ii} = \sum_i \sum_j A_{ij} B_{ji} = \sum_j \sum_i B_{ji} A_{ij} = \sum_j (BA)_{jj}$$

$$Tr(ABC) = Tr(BCA) = Tr(CAB)$$

determinant

$$|A^{-1}| = \frac{1}{|A|}$$

calculus

$$\frac{\partial}{\partial A} Tr(AB) = {}_i\left[\begin{array}{c} {}^j \\ \vdots \\ \cdots \frac{\partial}{\partial A_{ij}} Tr(AB) \cdots \end{array}\right] = B^T$$

$$\frac{\partial}{\partial A_{ij}} Tr(AB) = \frac{\partial}{\partial A_{ij}} \sum_i \sum_j A_{ij} B_{ji} = B_{ji}$$

PRML
Ch 2.5

nonparametric methods    vs.    parametric methods

↓                    ↓

use data directly          assume a model

assume some distance measure      find parameters

---

Central limit theorem

sum of a set of uniformly distributed random variables

$\xi$

Gaussian

See Fig 2.6 or run Matlab

1D histogram

$x$: continuous variable

partition $x$ into distinct bins of width $\Delta_i$

count the number $n_i$ of observation of $x$ falling in bin $i$

a normalized probability density is obtained by

$$p_i = \frac{n_i}{N \Delta_i}$$

Simplified   $\Delta_i = \Delta$   fixed bin width

$$\lim_{\Delta \to 0} \sum_i p_i \Delta = \lim_{\Delta \to 0} \sum_i \frac{n_i}{N} = 1$$

$$\int p(x)\, dx = 1$$

See Fig 2.24 or Run Matlab

small $\Delta$ → spiky

large $\Delta$ → smooth

histograms
     ① quantized
     ② visualization of data for 1D or 2D
     ③ discontinuity on bin edge
     ④ curse of dimensionality

two insights
     ① local neighborhood for estimating the probability density

     "locality"    who are your neighbors?

     histogram: neighborhood is defined by the bins

     ② bin width $\Delta$: smoothing parameters
             too small or too large is not good
             related to regularization, model complexity

Density Estimation (high dimensional)

assume data from some unknown probability density $p(x)$ in D-dimensional Euclidean space

Probability mass with region $\mathcal{R}$

$$P = \int_{\mathcal{R}} p(x)\, dx$$

each data point has a probability $P$ of falling within $\mathcal{R}$

total number $K$ of points inside $\mathcal{R}$ is a binomial distribution

$$K \sim \text{Bin}(K|N,P) = \frac{N!}{K!\,(N-K)!}\, P^K (1-P)^{N-K}$$

$$E[K] = \sum_{K=1}^{N} K \cdot \text{Bin}(K|N,P) = NP$$

$E[K] = NP$

$$\sum_{K=1}^{N} \frac{N!}{K!\,(N-K)!}\, P^K (1-P)^{N-K} = (P+(1-P))^N = 1 \qquad \text{differentiate}$$

$$\sum_{K=1}^{N} \frac{N}{K!\,(N-K)!}\, P^K (1-P)^{N-K} \left\{ \frac{K}{P} - \frac{N-K}{1-P} \right\} = 0$$

multiply $P(1-P)$

$$\sum_{K=1}^{N} \frac{N!}{K!\,(N-K)!}\, P^K (1-P)^{N-K} \left\{ K(1-P) - P(N-K) \right\} = 0$$

$$\sum_{K=1}^{N} \frac{N!}{K!\,(N-K)!}\, P^K (1-P)^{N-K} \left\{ K - PN \right\} = 0$$

so $E\left[\dfrac{K}{N}\right] = P$

similarly, variance: $\text{Var}\left[\dfrac{K}{N}\right] = \dfrac{P(1-P)}{N}$ (differentiate again)

for large $N$, we get a distribution sharply peaked around the mean, so

$$K \simeq NP$$

assume $\mathcal{R}$ is small that the probability density $p(x)$ is roughly constant over the region

$$P \simeq p(x)\, V \qquad V \text{ is the volume of } \mathcal{R}$$

we are interested in

$$p(x) = \frac{K}{NV}$$

we have two contradictory assumptions

① $\mathcal{R}$ should be sufficiently small that the density in $\mathcal{R}$ is constant

② $\mathcal{R}$ should be sufficiently large so that the number $K$ of points is sufficient for the binomial distribution to be sharply peaked

$$p(x) = \frac{K}{NV} \qquad \text{two different approaches}$$
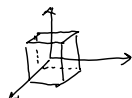
fix $K$ : k-nearest-neighbor

fix $V$ : kernel approaches

both converge to true probability density
in the limit $N \to \infty$, provided $V\downarrow$ with $N\uparrow$, $K\uparrow$ with $N\uparrow$

---

## Kernel Density Estimation (K. D. E.)

consider the kernel function

$$k(u) = \begin{cases} 1, & |u_i| \leq \frac{1}{2}, \ i=1,\cdots,D \\ 0, & \text{otherwise} \end{cases}$$

i.e. $R$ is a unit cube
centered at the origin

The total number $K$ of data points inside a cube of side $h$
centered on $x$

$$K = \sum_{n=1}^{N} k\left(\frac{x-x_n}{h}\right)$$

Therefore

$$p(x) = \frac{K}{NV} = \frac{1}{N h^D} \sum_{n=1}^{N} k\left(\frac{x-x_n}{h}\right) \qquad \Big| V = h^D$$

artificial discontinuities across the cube boundary

usually we use a Gaussian-like kernel function

$$p(x) = \frac{1}{N} \sum_{n=1}^{N} \frac{1}{(2\pi h^2)^{\frac{D}{2}}} \exp\left\{ -\frac{\|x-x_n\|^2}{2h^2} \right\}$$

$h$ is like the standard deviation in Gaussian

other kernel functions are allowable as long as

$$k(u) \geq 0$$
$$\int k(u) \, du = 1$$

$\Rightarrow$ no "training" for kernel density estimation
but the computational cost for testing is hig

Try Fig 25 MATLAB

Nearest Neighbor Methods    (KNN)

In K.D.E., optimal choice for $h$ may be dependent on location

(there is an issue called "bandwidth selection" in K.D.E.)

we may fix $K$ and use the data to find an
appropriate value for $V$

e.g. $K=5$

Small $V$        large $V$

KNN can be easily applied to multiclass classification problems

(Homework!)

Consider $N_m$ points in class $C_m$

$$\sum_{m=1}^{M} N_m = N$$

conditional    $P(x \mid C_m) = \dfrac{K_m}{N_m V}$        (likelihood)
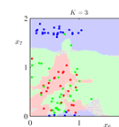
$$P(x) = \frac{K}{NV}$$

class priors are    $P(C_m) = \dfrac{N_m}{N}$

by Baye's theorem

the posterior is    $P(C_m \mid x) = \dfrac{P(x|C_m) \, P(C_m)}{P(x)} = \dfrac{K_m}{K}$

So the decision criterion is very simple:

To classify a new point, we find the $K$ nearest neighbor points
from the training data and assigned the new point to the
class having the largest number of representatives among the
$K$ nearest neighbors (the largest $K_m$)



$N \to \infty$, the error rate is never more than
twice minimum achievable error rate
of an optimal classifier, i.e., one that uses
the true class distributions.

① require entire training set to be stored
② the computational cost may be reduced by approximate
nearest neighbor techniques