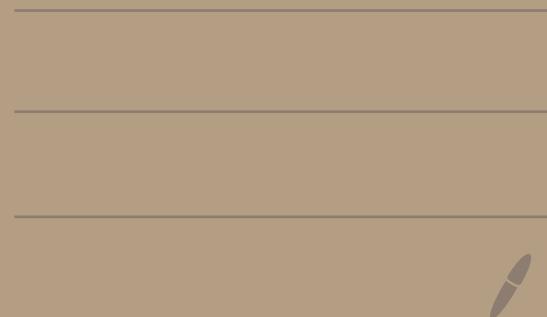


INF8953 CE - Fall 2020

# Machine Learning

- Sarath chandar

2. Statistical Decision Theory.



## 2. Statistical Decision Theory

Regression:-

$x \in \mathbb{R}^P$  — real valued random input vector.

$y \in \mathbb{R}$  — real valued random output variable.

$P_r(x, y)$  — joint distribution of  $x, y$ .

Goal: find a function  $f(x)$  for predicting  $y$   
given values of  $x$ .

$L(y, f(x))$  = loss function for penalizing errors  
in prediction.

Common loss function: Squared error loss.

$$L(y, f(x)) = (y - f(x))^2$$

Expected prediction error:-

$$EPE(f) = E[(y - f(x))^2]$$

$$= \int \int (y - f(x))^2 P(x, y) dx dy$$

$$= \int_n \left[ \int_y (y - f(x))^2 p(y|x) dy \right] p(n) dn$$

$$= \int_n E_{y|x} [(y - f)^2 | x=n] p(n) dn$$

$$= E_x E_{y|x} [(y - f)^2 | x=n]$$

It suffices to minimize EPE pointwise:

$$f^* = \underset{f}{\operatorname{argmin}} E_{y|x} [(y - f)^2 | x=n]$$

$$\frac{\partial}{\partial f} \left\{ E_{y|x} [(y - f)^2 | x=n] \right\} = 0$$

$$\frac{\partial}{\partial f} \left\{ E_{y|x} (y^2 | x=n) + f^2 - 2 E_{y|x} (y | x=n) f \right\} = 0$$

$$2f - 2 E_{y|x} (y | x=n) = 0$$

$$f^* = E_{y|x} (y | x=n)$$

↑  
regression function.

The best prediction of  $y$  at any point

$x = \bar{x}$  is the Conditional mean, when best is measured by average squared error.

What is Nearest Neighbor (N.N) doing?

$$\hat{f}(x) = \text{Ave} \left( y_i \mid x_i \in N_k(x) \right)$$

Two approximations to the regression function:

1. Expectation is approximated by averaging over sample data.
2. Conditioning at a point is relaxed to Conditioning on some region "close" to the target point.

Note 1: For large training sample size  $N$ , the points in the neighborhood are likely to be close to  $x$ .

Note 2: As ' $k$ ' gets larger, the average will get more stable.

Under mild regularity conditions on  $P(x,y)$ ,  
one can show that

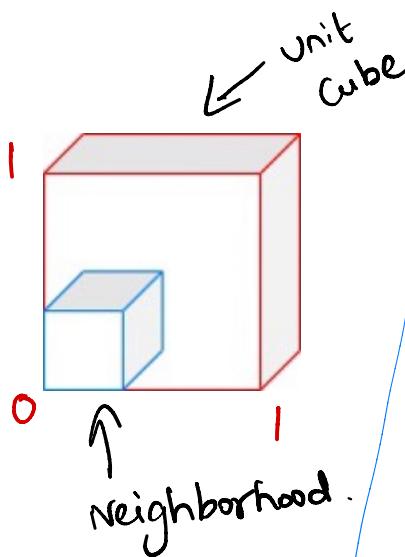
as  $N, k \rightarrow \infty$ , such that  $k/N \rightarrow 0$ ,

$$\hat{f}(x) \rightarrow E(y|x=x)$$

Looks like we have an universal function  
approximator!?

No! As the number of features increases, in  
high dimensions, we need very large  
samples.

### Local methods in high dimensions:-



Consider inputs uniformly distributed in a p-dimensional hypercube.

Consider a hypercubical neighborhood around a target point to capture a fraction 'r' of the observations.

This corresponds to a fraction 'r' of the unit volume.

$$\text{Expected edge length } e_p^{(r)} = r^{\frac{1}{p}}$$

$$e_{10}^{(0.01)} = 0.63$$

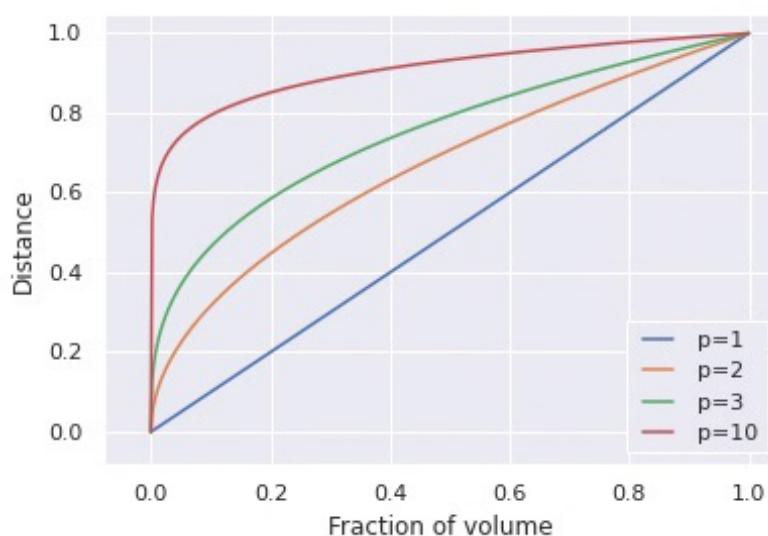
$$e_{10}^{(0.1)} = 0.80$$

Entire range for each input is only 1.0.

To capture 1% or 10% of data, to form a local average, we must cover 63% or 83% of the range of each input variable.

$\Rightarrow$  Such neighborhood are no longer "local".

$\Rightarrow$  if you reduce 'r', with fewer observations, variance will be high.



This is known as  
"Curse of dimensionality"

(Bellman, 1961)

Warning :- Our intuitions often break down in higher dimensions.

What is linear regression doing?

$$f(x) \approx x^T \omega$$

- This is a model-based approach. We are specifying a model for the regression function.

$$\omega = [E(x x^T)]^{-1} E(x y) \quad \text{--- } \star$$

Note : We have not conditioned on  $x$ .

rather we have used our knowledge of the functional relationship to pool over values of  $x$ .

Approximation :- Replace the expectation in  $\star$

by averages over the training data.

## Assumptions :-

- ① Linear regression assumes  $f(x)$  is well approximated by a globally linear function.
- ② k-nearest neighbor assumes  $f(x)$  is well approximated by a locally constant function.

These assumptions made by any algorithm is known as the inductive bias of the algorithm.

---

What happens if we replace  $L_2$  loss with  $L_1$  loss?

$$L_1 : E |y - f(x)|$$

Solution :  $\hat{f}(x) = \text{median}(y | x=x)$

which is a different measure of location, and its estimates are more robust than those for the Conditional mean.

Historically,  $L_2$  is more popular than  $L_1$  because of the discontinuities in the derivatives of  $L_1$ .

---

## Classification:-

Output : Categorical variable  $G_i$ .

Same paradigm works here.

We need a different loss fn. for penalizing prediction errors.

$G$  - set of possible classes.

$\hat{G}$  - estimate.

loss fn:  $k \times k$  matrix where  $k = \text{Card}(G)$

confusion mat

	1	2	..	$k$
1	0	0	0	0
2	0	0	0	0
:	0	0	0	0
$k$	0	0	0	0

Price paid for classifying an observation belonging to class 1 to class  $k$ .

$L(k, l)$  - price for classifying an observation belonging to class  $k$  to class  $l$ .

Common loss: 0/1 loss function.

0 for correct classification.

1 for misprediction.

$$EPE = E \left[ L(G_i, \hat{G}(x)) \right]$$

the output here is discrete, hence we use summazation.  
In the top the o/p is continous, hence we used integral

Expectation is w.r.t.  $\Pr(G_i | x)$

$$= E_x \sum_{k=1}^K L[G_k, \hat{G}(x)] \Pr(G_k | x)$$

again it suffices to minimize EPE point wise:

$$\hat{G}(x) = \operatorname{argmin}_{g \in \mathcal{G}} \sum_{k=1}^K L(G_k, g) \Pr(G_k | x=n)$$

For 0/1 loss:

$$\hat{G}(n) = \operatorname{argmin}_{g \in \mathcal{G}} [1 - \Pr(g | x=n)]$$

(or)

$$\boxed{\hat{G}(n) = \operatorname{argmax}_{g \in \mathcal{G}} [\Pr(g | x=n)]}$$

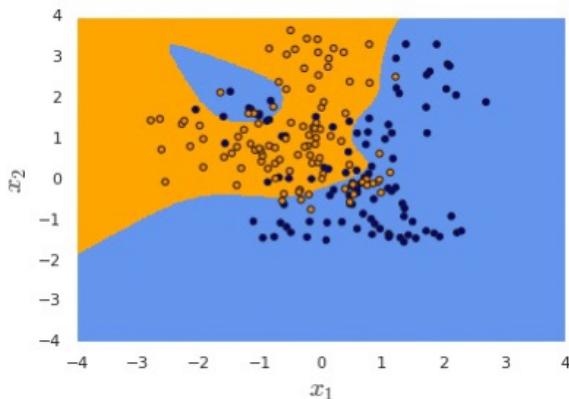
$\Rightarrow$  This solution is known as the Bayes classifier.

$\Rightarrow$  we classify to the most probable class, using  
the conditional distribution  $\Pr(G_i | x)$ .

$\Rightarrow$  The error rate of the Bayes classifier  
is called the Bayes error rate.

We cannot always calculate the bayes error, since we dont the distribution

Note: Bayes error rate for a particular classification problem is the lowest possible error rate any classifier can achieve.



The optimal Bayes decision boundary for the toy task.

K-Nearest neighbor directly approximates this solution.

Majority voting in the neighborhood.

- ① Conditional prob. at a point is relaxed to Conditional prob. within a neighborhood of a point.
- ② Probabilities are estimated by training same proportions.

You should know!

1. Expected prediction error
  2. Regression function
  3. Curse of dimensionality
  4. Inductive bias of an algorithm.
  5. Model based algorithms.
  6. 0/1 loss.
  7. Bayes classifier
  8. Bayes rate .
-