

Introduction to Machine Learning by Elhann Alpaydm

sample stat of a func computed from X_1, \dots, X_n random var.,
thus itself is also a random var. $g(X_1, \dots, X_n)$

→ Sample mean → sample variance

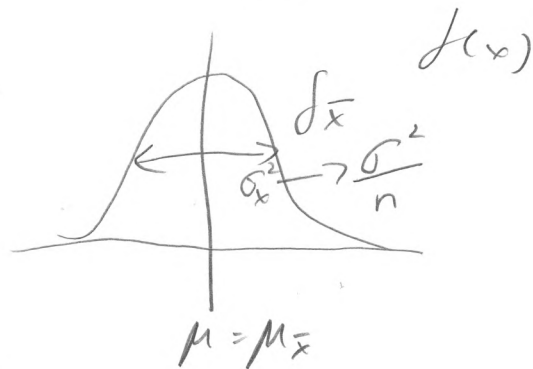
Assume X_1, \dots, X_n is iid (independent, identically distributed)

$$\bar{X} = \frac{X_1 + \dots + X_n}{n} \quad \text{Var}(\bar{X}) = E[(\bar{X} - \mu)^2] = E[\bar{X}^2] - (\mu)^2$$

mean of sampling distribution of means = $\mu_{\bar{X}}$

Thm. 1 $E(\bar{X}) = \mu_{\bar{X}} = \mu$

mean of stat. mean of pop.



Thm. 2 $E[(X - \mu)^2] = \sigma_X^2 = \frac{\sigma^2}{n}$

variance of sampling mean distribution

= variance of population

infinite pop. /
finite pop. w/
replacement

Thm. 3 $\sigma_{\bar{X}}^2 = \frac{\sigma^2}{n} \left(\frac{N-n}{N-1} \right)$ finite pop. w/o replacement
N sample size n

Thm. 4 If pop. is ~~uniform~~ normally dist. w/ μ, σ^2 ,
sample should also w/ μ & $\frac{\sigma^2}{n}$

Thm. 5 If pop. has μ, σ^2 but not normal,
standardized variable associated w/ \bar{X} , given by
 $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$ \bar{X} : sampling mean.

is asymptotically normal, $\lim_{n \rightarrow \infty} P(Z \leq z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-u^2/2} du$

Thm 5. \square from central limit then.

otherwise thm 5 will still be correct if replace $\frac{\sigma}{\sqrt{n}}$ by $\frac{\sigma^2}{\bar{x}}$ as given in thm 3

Sampling Binomial Dist

$$\mu = p \quad \sigma_p = \sqrt{\frac{p(1-p)}{n}}$$

$$\therefore \mu = p, \quad \sigma = \sqrt{pq}$$

give 0 to 1

$$E(x) = 0 \cdot \bar{q} + 1 \cdot p = p$$

$$E(x^2) = 0^2 \cdot \bar{q} + 1^2 \cdot p = p$$

$$p - p^2 = p(1-p) = pq$$

Sampling Variance

$$s^2 = \frac{(x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n}$$

unbiased estimator

$$\hat{s}^2 = \frac{n}{n-1} s^2 = \frac{\sum_{k=1}^n (x_k - \bar{x})^2}{n-1}$$

On Ave. \hat{s}^2 is target

Smaller Var gives more efficient estimator.

Bessel's correction: replace n by $(n-1)$

$\mu \pm n\sigma$ (68-95-99 rule) (confidence interval)

$\mu \pm 1.96\sigma$: 95% $\mu \pm 2.58\sigma$ (0.99) crit. values.
(confidence interval levels) z_c

$$n \geq 30 \left[\begin{array}{l} \bar{X} \pm \frac{1.96}{z_c} \sigma_{\bar{X}} \approx \pm 2.58 \sigma_{\bar{X}} \quad (\text{w/ replace / } \infty) \\ \bar{X} \pm z_c \frac{\sigma}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}} \quad (\text{w/o replacement}) \end{array} \right]$$

$n < 30$: t-distribution $\pm t_{0.975}$ $\hat{\sigma}$: s.d.

$$95\%: -t_{0.975} < \frac{(\bar{X} - \mu)\sqrt{n}}{\hat{\sigma}} < t_{0.975}$$

Master Thm.
 $f(1) = O(1)$
 $f(n) = O(n^k) + f(\frac{n}{b})$
 $\log_b a < k = O(n^k)$
 $\log_b a = k = O(n^k \log n)$
 $\log_b a > k = O(n^{\log_b a})$

k-selection
 $\begin{matrix} x & y & z \\ x & y & z \\ x < y & y < z \\ x = k-1 & \text{return } y \end{matrix}$
 $x = k-1$ return y
 $<$ recurse left
 $>$ recurse right
choose v in
 $\leq O(n)$
 $x \geq \frac{n}{2} \quad z \geq \frac{n}{2}$

Dominance Counting
1. vert. ln. such that each side $\frac{n}{2}$ pts by k-selection $O(n)$
2. find pts dominated on left, right recursively
3. set left and right pts using y-axis
4. merge sort inversion

Strassen
 $\begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} E & F \\ G & H \end{bmatrix} = \begin{bmatrix} P_5 + P_4 - P_2 + P_6 & P_1 + P_2 \\ P_3 + P_4 & P_1 + P_3 - P_7 + P_7 \end{bmatrix}$
 $P_1 = A_{11}(B_{12} - B_{22}) \quad P_2 = (A_{11} + A_{22})B_{22}$
 $P_3 = (A_{21} + A_{22})B_{11} \quad P_4 = A_{22}(B_{21} - B_{11})$
 $P_5 = (A_{11} + A_{22})(B_{11} + B_{22}) \quad P_6 = (A_{21} - A_{22})(B_{11} + B_{22})$
 $P_7 = (A_{11} - A_{21})(B_{11} + B_{22})$

$$P_L \approx \sqrt{\frac{P_L P_R}{n}} \quad P_{\text{obs}} \quad P = \bar{p}$$

Error type: reject H_0 (false neg)

level of confidence $\alpha = 0.05$ or 0.01 use S.D.

complement of S.D.

True neg True pos.

95% confident if z-score between -1.96 and 1.96

If z-score outside only 0.05 prob.

1-tailed / 2-tailed

Level of significance

	0.10	0.05	0.01	0.005	
± 1.28		± 1.645	± 2.33	± 2.58	1-tailed
± 1.645		± 1.96	± 2.58	± 2.81	2-tailed

null: H_0 gives $\mu = \mu$

assume $\mu = 12$

H_1

$\mu > 12$

$\mu < 12$

$n = 36$

p values

$P(Z \geq 1.9) = 0.029$

$P(Z \leq -1.9)$

$\bar{X} = 12.98$

about $\frac{3}{100}$

$z = 1.9$

that $\bar{x} \geq 12.98$

$$\mu \neq \mu = P(Z \geq 1.9) + P(Z \leq -1.9)$$

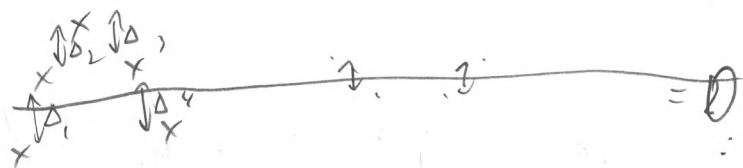
Small r values suggests reg. H_0

Method of least squares $\underbrace{d_1^2 + \dots + d_n^2}_{\text{minimum}}$ small

Linear approx.

$$\sum_{j=1}^n y_j = a n + b \sum_{j=1}^n x_j$$

$$\sum_{j=1}^n x_j y_j = a \sum_{j=1}^n x_j + b \sum_{j=1}^n x_j^2$$



$$a = \frac{\bar{y}(\sum x_j^2) - \sum x_j \bar{x} y_j}{n \sum x_j^2 - (\sum x_j)^2}$$

$$b = \frac{\sum (x_j - \bar{x})(y_j - \bar{y})}{\sum (x_j - \bar{x})^2}$$

3 curves

$$y, y^{-1}, (g(y))$$

$$\sum_{var} = \{x^e, v^e\}_{e=1}^N$$

data pt. no.

Data Collecting Amount:

Probably approximately correct (PAC)

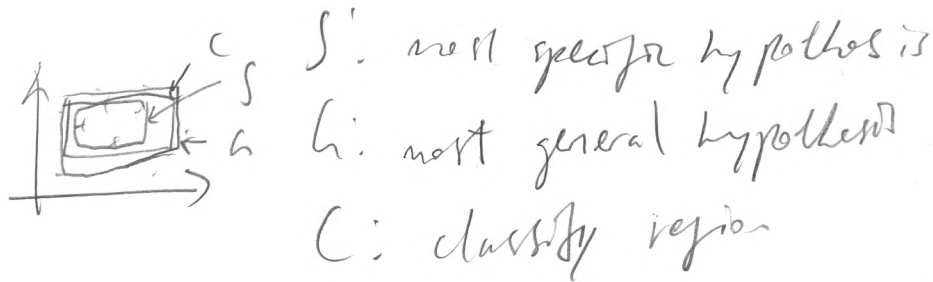
VC Dimension (relate learner complexity to errors)

class C : class to classify

$$\bar{E}(h|X) = \sum_{e=1}^N 1(h(x^e) \neq v^e)$$

Error (Total Err)

$$h(x) = \begin{cases} 1 & \text{if classifies as pos} \\ 0 & \text{if classifies as neg} \end{cases}$$



Let H , between S & h is consistent & make up the version space.

RAC

How many N training examples such that w/ prob $(1-\delta)$, h has error at most ϵ ?

$$P\{C(h) \leq \epsilon\} \geq 1 - \delta \quad (C(h), \text{ region between } C \text{ and } h)$$

VC Dim Capacity of learning machine
 N pts can be labeled in 2^N for shatter

If $h \in H$ separates, H shatters N
 max no. of N pts h shatters.

look for test err w/ train err

\uparrow cap \rightarrow \downarrow performance over-fitting

$$\text{test error} \leq \text{train error} + \sqrt{\frac{C \left(\log\left(\frac{2N}{\epsilon} + 1\right) - \log\left(\frac{\eta}{\epsilon}\right) \right)}{N}}$$

N : num. of samples

C : VC dim

Reality (complexity)

VC Dim for linear

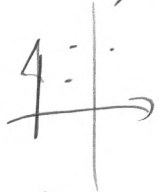
n -dim feature space set of n : $(n \geq n)$

it is ~~shatter~~ general dimension $1 \leq n$ no subset of $(n+1)$ pts lies on $(n-1)$ hyperplane

\therefore no 3 pts. lie on 1D line.

H shatters m pts in n -dim of all possible combinations of m pts. in n -dim space exactly classified by H

H , not h , shatters, classifier, not actual hypothesis



st. lin. shatters 3 pts in 2 dim

VC dim \Rightarrow cardinality of largest set of pts of H can shatter

VC dim linear: $(n+1)$ n : dim

Shattering: $\exists h \in H : E(h, X) = 0$

Neur, Latent Var

$$\hat{r}_{ij} = \sum_{k=1}^K p_{ik} q_{kj} \quad \text{for } i=1 \dots |U| \text{ \& } j=1 \dots |D|$$

$$e_{ij}^2 = (r_{ij} - \hat{r}_{ij})^2 = \left(r_{ij} - \sum_{k=1}^K p_{ik} q_{kj} \right)^2$$

$$\frac{\partial e_{ij}^2}{\partial p_{ik}} \quad \frac{\partial e_{ij}^2}{\partial q_{kj}}$$

$$= -2e_{ij} q_{kj} \quad = 2e_{ij} p_{ik}$$

Gradient descent $\eta \approx 0.002$

Minimize for observed ratings only $\rightarrow \hat{R}$

$$\sum_T e_{ij}^2 = \sum_{i,j \in T} \left(r_{ij} - \sum_{k=1}^K p_{ik} q_{kj} \right)^2$$

$$e_{ij}^2 = \left(r_{ij} - \sum_{k=1}^K p_{ik} q_{kj} \right)^2 + \frac{\beta}{2} \left(\sum_{a=1}^{|U|} \sum_{b=1}^K (p_{ab})^2 + \sum_{a=1}^K \sum_{b=1}^{|D|} (q_{ab})^2 \right)$$

Matrix Factorization

- make recommendation under partial information
- Let A be $m \times n$ mat of rank r

$$A = XY \quad X: m \times r \quad Y: r \times n$$

\therefore storage saving fr by

$$A = \begin{bmatrix} 1 & 2 & 3 & 5 \\ 2 & 4 & 6 & 10 \\ 3 & 6 & 9 & 15 \end{bmatrix} \quad r=2$$

$$A_1 = 1 \cdot A_1 + 0 \cdot A_3 \quad A_3 = 0 \cdot A_1 + 1 \cdot A_3$$

$$A_2 = 2 \cdot A_1 + 0 \cdot A_3 \quad A_4 = 2 \cdot A_1 + 1 \cdot A_3$$

$$\therefore X = \begin{bmatrix} 1 & 3 \\ 2 & 6 \\ 3 & 9 \end{bmatrix} \quad Y = \begin{bmatrix} 1 & 2 & 0 & 2 \\ 0 & 0 & 1 & 1 \end{bmatrix}$$

e.g. Netflix

$$U = |4| \text{ users} \quad D = |10| \text{ movies} \quad \therefore R = \overset{|4| \times |10|}{\cancel{A \times D}} (X)$$

$$R \text{ by } R (|4| \text{ by } |10|)$$

$$\therefore R = P Q \quad \text{relates each user w/ some latent features}$$

$$R \approx P \times Q = \hat{R}$$

M-step: $\bar{Z}^{l+1} = \underset{\bar{Z}}{\operatorname{argmax}} Q(\bar{Z} | \bar{Z}^l) = \sum_t \sum_i h_i^t [\log \pi_i + \log p_i(x^t | \bar{Z}^l)]$

$\pi_i = \frac{\sum_t h_i^t}{N}$

$$\nabla_{\pi_i} \sum_t \sum_i h_i^t \log \pi_i - \lambda \sum_i \pi_i - 1 = 0$$

$$\nabla_{\bar{Z}} \sum_t \sum_i h_i^t \log p_i(x^t | \bar{Z}) = 0$$

If Gaussian $\hat{p}(x^t | \bar{Z}) \sim \mathcal{N}(m_i, S_i)$ M-step \rightarrow

$$m_i^{l+1} = \frac{\sum_t h_i^t x^t}{\sum_t h_i^t} \quad S_i^{l+1} = \frac{\sum_t h_i^t (x^t - m_i^{l+1})(x^t - m_i^{l+1})^T}{\sum_t h_i^t}$$

$$h_i^t = \frac{\pi_i |S_i|^{-\frac{1}{2}} \exp(-\frac{1}{2} (x^t - m_i)^T S_i^{-1} (x^t - m_i))}{\sum_j \pi_j |S_j|^{-\frac{1}{2}} \exp(-\frac{1}{2} (x^t - m_j)^T S_j^{-1} (x^t - m_j))}$$

Hierarchical Clustering \rightarrow Agglomerative ($N \rightarrow k$ Group) / Divisive ($1 \rightarrow k$)

Minkowski $d_m(x^i, x^j) = \left(\sum_j (x_j^i - x_j^j)^p \right)^{1/p}$ end to end for $p \geq 2$

City-block $d_{cb}(x^i, x^j) = \sum_j |x_j^i - x_j^j|$ $p=1 \rightarrow$ Manhattan

single-link = $\min d()$ complete link = $\max d()$ Ave-link = ave $d()$ centroid $x^i \in C_i, x^j \in C_j$

Non-param est.

Histogram = $\frac{1}{N} \frac{\text{bin}}{\text{range}}$

Naive = $\frac{1}{N} \frac{x^t \text{ in range}}{h}$

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

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

\hookrightarrow kNN: $\hat{p}(x) = \frac{k}{2N d_k(x)}$

Generalization on multivariate

Gaussian: $K(u) = \left(\frac{1}{\sqrt{2\pi}}\right)^d e^{(-\frac{|u|^2}{2})}$

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

curse of dimensionality

Classification

$\hat{p}(x | C_i) = \frac{1}{Nh^d} \sum_{t=1}^N K\left(\frac{x - x^t}{h}\right) h_i^t \rightarrow g_i(x) = \hat{p}(x | C_i) \hat{p}(C_i)$

can $\hat{p}(x | C_i) = \frac{k_i}{N_i V^k(x)} \leftarrow$ no. of neighbors $\in C_i$ in k th $\rightarrow P(C_i | x) = \frac{k_i}{k}$
 $N_i V^k(x) \leftarrow$ volume of k th x^k from x

Condensed nearest neighbour $\rightarrow 1-n$

For all $x \in X$

$$\text{Find } x' \in Z \text{ s.t. } |x - x'| = \min_{x'' \in Z} |x - x''|$$

If $\text{class}(x) \neq \text{class}(x')$ add to Z

Until Z does not change (end)

Distance based

$$D(x, m_i) = \min_{j=1}^k (x, m_j)$$

$$D(x, m_i) = |x - m_i| \sqrt{(x - m_i)^T S_i^{-1} (x - m_i)}$$

or Mahalanobis

$$D(x, x^e | M) = (x - x^e)^T M (x - x^e)$$

$$M = L^T L \quad \begin{matrix} n \times d \\ k \times d \end{matrix} \rightarrow$$

$$= \cancel{|x - x^e|^2} (Lx - Lx^e)^T (Lx - Lx^e)$$

$$= |x - x^e|^2$$

Outlier detection

$$LOF(x) = \frac{d_k(x)}{\sum_{s \in N(x)} d_k(s) / |N(x)|} \leftarrow \text{number of samples} \quad \gg 1 \rightarrow \text{outlier}$$

Non-param regression

same as classification, where r^e no longer classifier, but result

regression $\rightarrow \hat{f}(x) = \frac{\sum b(x, x^e) r^e}{\sum b(x, x^e)} \quad b(x, x^e) \text{ if } x^e \text{ is in bin}$

running-mean

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

$$w(u) = \begin{cases} 1 & \text{if } |u| < 1 \\ 0 & \text{otherwise} \end{cases} \quad \hat{f}(x)$$

reg kernel

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

k, h small $\rightarrow \uparrow$ complexity
 \uparrow var \downarrow bias

k, h big $\rightarrow \downarrow$ complexity
 \downarrow var \uparrow bias

decision tree
greedy algorithm
recursive split

leave: result
internal decision nodes: Numerical split / Discrete split (binary)

Univariate

$$\hat{p}(C_i | x, m) = p_m^i = \frac{N_m^i}{N_m} \quad N_m = \text{Number of instances at node } m$$

node m pure if $p_m^i = 0$ or 1 , impurity $I_m = -\sum_i p_m^i \log_2 p_m^i$ (entropy)

$$\bar{I} = -p \log_2 p - (1-p) \log_2 (1-p) \text{ (entropy)} \quad (0,1] \Rightarrow 0 \rightarrow \text{pure}$$

$$\bar{I}(p) = 2p(1-p) \text{ (Gini index)} \quad \bar{I}(p, 1-p) = 1 - \max(p, 1-p) \text{ (misclassification err.)}$$

If node is pure, generate leaf and stop, else split recursively

Impurity after split: N_{mj} of N_m branch j

$$\hat{p}(C_i | x, m, j) = p_{mj}^i = \frac{N_{mj}^i}{N_{mj}} \quad I_m = -\sum_{j=1}^n \frac{N_{mj}}{N_m} \sum_{i=1}^K p_{mj}^i \log_2 p_{mj}^i$$

To generate tree:

if entropy $< \theta_1$

create leaf of majority

return

else

split

for all branches,

generate tree(branch)

Split: All Attributes $i, 1 \dots d$

if x_i is discrete of n vals. (discrete)

$e = \text{split entropy}(X_1, \dots, X_n)$

if $e < \min$, best $\leftarrow i$, $\min \leftarrow e$

else (Numerical)

split X into X_1, X_2 on x_i

$e \leftarrow \text{split entropy}$

if $e < \min$, $\min \leftarrow e$, best $\leftarrow i$

return best

Univariate regression $x_{mj} \subset x_m$ of table branch j

$$b_{mj}(x) = \begin{cases} 1 & \text{if } x \in X_{mj} \\ 0 & \text{otherwise} \end{cases}$$

$$g_{mj} = \frac{\sum_t b_{mj}(x^t) r^t}{\sum_t b_{mj}(x^t)}$$

$$E'_m = \frac{1}{N_m} \sum_j \sum_t (r^t - g_{mj})^2 b_{mj}(x^t)$$

minimize error E'_m , replace entropy with E'_m , θ_1 , $\frac{1}{2}$, θ_r

Pre-pruning: early stop Post-pruning: grow tree then prune what overfitted on pruning set

Rule learning

tree induction (BFS) rule induction (DFS)

1. rule set contains rules

2. remove covered (true) samples of rules in ruleset when added

Sequential, adds 1 rule at a time till all evaluated.

Multivariate,

$$\text{Node} = f_m(x) = w_m^T x + w_{m0} > 0 \quad \text{linear hyperplane}$$

For discrete attributes shld be 0/1 dummy numeric divide until defined by polyhedra in input space

$$\text{or } f_m(x) = x^T W_m x + w_m^T x + w_{m0} > 0 \rightarrow \text{quadratic}$$

$$\text{or } f_m(x) = |x - c_m| \leq \alpha_m \quad \text{sphere node}$$

Discriminant Assume model for $g_i(x|\Phi_i)$ and find boundary

Linear: $g_i(x|w_i, w_{i0}) = w_i^T x + w_{i0} = \sum_{j=1}^d w_{ij} x_j + w_{i0} \quad O(d)$

when shared cov., and almost linearly separable

Quad: $g_i(x|W_i, w_i, w_{i0}) = x^T W_i x + w_i^T x + w_{i0} \quad O(d)$

or convert to linear

$g_i(x) = \sum_{j=1}^k w_j \phi_{ij}(x) \rightarrow$ basis func: $\sin(x), \exp(-(x-m)^2/c), (g(x), 1(x>0), \dots)$

$K=2: g(x) = g_1(x) - g_2(x) = w^T x - w_0$

$r = \frac{g(x)}{|w|}$

$r_0 = \frac{w_0}{|w|}$

find w using LDA

$K \geq 2 \quad g_i(x|w_i, w_{i0}) = w_i^T x + w_{i0} = \begin{cases} > 0 & \text{if } i \in C_i \\ \leq 0 & \text{otherwise} \end{cases}$
 $|g_i(x)| / |w_i|$ is dist
 $\max_j g_j(x)$

Pairwise Separation Use $K(K-1)/2$ linear discriminants, $g_{ij}(x)$
 where $g_{ij}(x|w_{ij}, w_{ij0}) = w_{ij}^T x + w_{ij0}$

$P(C_i|x) = \text{sigmoid}(w^T x + w_0) = \frac{1}{1 + \exp(-(w^T x + w_0))}$