


Statistical Methods in AI

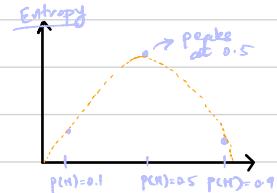
$$\text{Information } I(x) = \log\left(\frac{1}{p(x)}\right) \text{ when } p(x) \neq 0$$
$$= 0 \text{ when } p(x) = 0.$$

Average Information of RV \longleftrightarrow Expectation of Information

$$E(I(x)) = \sum_i p(x=i) \log\left(\frac{1}{p(x=i)}\right) = \sum_i p(x=i) \log(p(x=i))$$

This average information is called shannon entropy

$$\text{Entropy of 2-sided coin: } -\left(\frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2}\right) = 1.$$



$$p(H) = 0.9 \mid \text{Entropy} = 0.9 \log\left(\frac{10}{9}\right) + 0.1 \log(10) = 0.46894$$

Entropy $\uparrow \rightarrow$ Information $\uparrow \rightarrow$ uncertainty $\uparrow \rightarrow$ Impurity \uparrow

Step-1: Compute impurity score of training distribution (Starting Impurity)

Step-2: Compute impurity score for each unique value of candidate attribute.

Step-3: Compute impurity score for candidate attribute (Take weighted average)



$$\text{Information Gain: } \text{Gain}(S, A) = E(S) - I(S, A)$$

Step 6: Assign Root Node (with max Gain)

Recurse 8 Repeat step 1 to 6.

Properties of an impurity measure:

An Impurity measure is a function $i(V)$ s.t.....

→ Entropy

→ Gini index

→ Misclassification rate

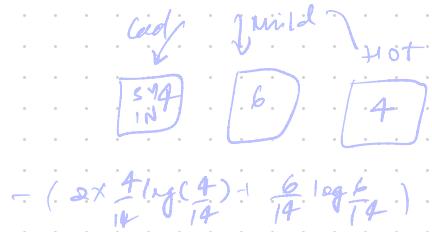
→ Intrinsic Information of split: \rightarrow more \leftarrow less useful

$$(Chu's attribute) I = -\sum \frac{|S_i|}{|S|} \log \frac{|S_i|}{|S|}$$

→ Gain Ratio = $\frac{\text{Information Gain}}{\text{Intrinsic Information}}$

$$GR(\text{Outlook}) = \frac{0.246}{1.5774} = 0.157$$

$$GR(\text{Day}) = 0.23$$



$$- \left(2 \times \frac{4}{14} \log \frac{4}{14} + \frac{6}{14} \log \frac{6}{14} \right)$$

$$GR(\text{Temp}) = \frac{0.0289}{1.5566} = 0.01856542927$$

$$GR(\text{Windy}) = 0.052$$

$$GR(\text{Humidity}) = 0.151$$

Handling Numerical attributes

Nominal can split once but not numeric attributes

Bayes Theorem

- Navresh Mawani

Disease	Not disease
99	999.00
99	

$$\text{population} = 100,000$$

$$\text{disease} = 100$$

$$\text{Bayes' Theorem: } P(H|E) = \frac{P(E|H) \cdot P(H)}{P(E|H) \cdot P(H) + P(E|\neg H) \cdot P(\neg H)}$$

$P(MA) = 0.4$	$P(MB) = 0.6$
$P(D MA) = 0.1$	$P(D MB) = 0.01$
$P(MA D) = ?$	(4/7)

Generalized Bayes Theorem:

If $\{E_i\}$ form a partition of sample space.

Bayes Classifier

Bayes Classification

$$L_{0-1} = \begin{cases} 0 & \text{if } f(x) = y \\ 1 & \text{if } f(x) \neq y \end{cases}$$

$\{(x_i, y_i)\} \forall i \in 1 \dots N \Rightarrow$ Training set

Objective function: $\frac{1}{N} \sum_{i=1}^N L_{0-1}(f(x_i), y_i) \sim$ Average loss

IDEAL GOAL:

$$\min_{f(x)=y} E[L_{0-1}(f(x), y)]$$

def

$$RL(f) = E[L_{0-1}(f(x), y)]$$

Lecture #9: Regression

- Independent Variable / dependent variable

Eg. Price of crude oil / Retail price of crude oil

- does not depend on causality

- Linear regression

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \rightarrow \text{Random Error}$$

\downarrow
intercept \downarrow
slope

- Estimating Parameters: Least square Method

- "Best fit" means between Act Y & pred. Y are minimum. Since positive differences offset the -ve ones.
- While adding errors, +ves may cancel out with -ves, therefore we square the errors
- Therefore total error $E = \sum_{i=1}^n e_i^2$, where $e_i = y_i - \hat{y}_i$
 $= y_i - (\beta_0 + \beta_1 x_i)$

$$\therefore E = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 = f(\beta_0, \beta_1)$$

we want $f(\beta_0, \beta_1)$ minimized.

$$\frac{\partial f}{\partial \beta_0} = 0, \frac{\partial f}{\partial \beta_1} = 0$$

$$\Rightarrow \frac{\partial}{\partial \beta_0} \left(\sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 \right) = 0$$

$$\Rightarrow \sum_{i=1}^n \frac{\partial}{\partial \beta_0} (y_i - (\beta_0 + \beta_1 x_i))^2 = 0$$

$$\Rightarrow \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))(-1) = 0$$

$$\beta_0 = \frac{\sum y_i - \beta_1 \sum x_i}{n}$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

$$f(\beta_0, \beta_1) = \sum_i (y_i - (\beta_0 + \beta_1 x_i))^2$$

$$\frac{\partial f}{\partial \beta_1} = \sum_i -2(y_i - \beta_0 - \beta_1 x_i)x_i = 0 \\ = n\bar{y} - n\beta_0 - \beta_1 n\bar{x} = 0.$$

$$\Rightarrow \beta_1 = \frac{\bar{x} - \beta_0}{\bar{x}^2}$$

$$\Rightarrow \beta_1 = \frac{\bar{x} - (n\bar{y} - \beta_1 \bar{x})}{\bar{x}^2}$$

$$\beta_1 = \frac{\sum (y_i - \bar{y})x_i}{\sum (x_i - \bar{x})x_i}$$

$$\bar{y} = \beta_0 + \beta_1 \bar{x}$$

$$\bar{y} = \bar{y} + \beta_1 \bar{x}$$

$$\bar{y} - \beta_0 \bar{x} - \beta_1 \bar{x}^2 = 0$$

$$\bar{y} - (\bar{y} - \beta_1 \bar{x})\bar{x} - \beta_1 \bar{x}^2 = 0$$

$$\bar{y} - \bar{x}\bar{y} - (\bar{x}^2 + \bar{x}^2)\beta_1 = 0$$

$$\beta_1 = \frac{\bar{y} - \bar{x}\bar{y}}{\bar{x}^2 - \bar{x}^2}$$

$$= \frac{\sum n_i y_i - \sum \bar{x}_i \bar{y}}{\sum n_i^2 - \sum \bar{x}_i n_i} \Rightarrow \frac{\bar{y} - \bar{x}\bar{y}}{\bar{x}^2 - (\bar{x})^2}$$

$$\beta_1 =$$

$$MAE = \text{Mean Absolute Error} = \sum |y - \hat{y}| / n$$

$$MSE = \text{Mean squared error} = \sum (y - \hat{y})^2 / n / RMSE$$

look at the distribution error since MAE/MSE prone to outlier

$$MPE = \frac{100}{n} \times \sum \left(\frac{|y - \hat{y}|}{y} \right) \rightarrow \text{Mean Positive error} \propto \text{large MPE} \rightarrow \text{underestimating (MAPE)}$$

→ lots of underestimating / overestimating

Linear Regression - Matrix Form

$$Y = \beta X + \epsilon$$

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \quad X = \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} \quad \epsilon = \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

$$E = e^T e$$

$$\min_e e^T e = \min_\beta (Y - X\beta)^T (Y - X\beta) = f(\beta)$$

$$\begin{aligned} \nabla_\beta f(\beta) = 0 &\Rightarrow \nabla_\beta (Y^T - \beta^T X^T)(Y - X\beta) \\ &\Rightarrow \nabla_\beta (Y^T Y - Y^T X\beta - \beta^T X^T Y + \beta^T X^T X\beta) = 0 \end{aligned}$$

Dimensions: $X = n \times p+1$

$$Y = n \times 1$$

$$\beta = (p+1) \times 1$$

$$\nabla_\beta (Y^T Y - 2\beta^T X^T Y + \beta^T X^T X\beta)$$

$$\Rightarrow 0 - 2X^T Y + 2X^T X\beta = 0$$

$$\Rightarrow X^T X\beta = X^T Y$$

$$\Rightarrow \boxed{\beta = (X^T X)^{-1} X^T Y} \Rightarrow H = (X^T X)^{-1} X^T = \text{influence matrix}$$

$$\beta = HY$$

$n \approx p \rightarrow \beta \text{ comp stable}$

$n > p \rightarrow \text{more samples than unknowns}$

$n \prec p \rightarrow \text{less samples than unknowns}$

Computing inverse $\rightarrow O(n^3)$

PCA - Eigen analysis

$$\begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} = 1+2^2 = 5$$

$$\begin{bmatrix} 1 \\ 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix}$$

$$V^T V = S$$

$$V^T V$$

$$V V^T$$

Gradient descent

Linear Regression: linear in coefficients (parameters) and NOT Variables

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + e \rightarrow \text{we can do linear regression}$$

Correlation \neq causation

Regression (k-NN) (8 DT)

Unsupervised Learning

learning the density of the data

(Given $x \in X$; learn $f(x)$)

Expectation Maximization

K-medoids

→ use L1 instead of L2 & instead of mean, medians in KNN

Maximum Likelihood Estimation

Mid-sum will be proof based



for a concave function f and a random variable X ,
 $f(E(X)) \geq E(f(X)) \rightarrow$ Jensen's Inequality

$$L(\theta) = \prod_{i=1}^N p(x_i; \theta) \quad \ell(\theta) = \sum_{i=1}^N \log(p(x_i; \theta)) \quad \text{can be anything}$$

$$\theta^* = \underset{\theta}{\operatorname{argmax}} L(\theta) \quad \ell(\theta) = \sum_{i=1}^N \log(p(x_i; \theta))$$

$$\alpha^* = \underset{\alpha}{\operatorname{argmax}} \ell(\alpha)$$

To prove $\alpha^* = \theta^*$

Assume not, then by contradiction it can be proved

$$\alpha^* \neq \theta^*$$

$$\mu(\theta) = \sum_i Q_i(z_i) \log \left(\frac{P(x_i; z_i, \theta)}{Q_i(z_i)} \right) \rightarrow \text{Evidence Lower Bound (ELBO)}$$

EM is co-ordinate ascent on ELBO

$$\max_a f(a, b)$$

$$\begin{aligned} \hat{b} &= \max_b f_b(\hat{a}, b) \\ \hat{a} &= \max_a f_a(a, \hat{b}) \end{aligned} \quad \text{iteratively}$$

KL-Divergence: Maximizing ELBO is minimizing $D_{KL}(Q_i(z_i) || P(z_i|y_i))$

$$KL\text{-div} \triangleq \sum_{i=1}^B p_i \log \left(\frac{p_i}{q_i} \right)$$

ELBO maximized when $Q_i(z_i) \approx P(z_i|y_i)$

proof of convergence for EM:

$$l(\theta^{(t+1)}) \geq l(\theta^{(t)})$$

$$l(\theta) \geq \sum_i \sum_{z_i} Q_i(z_i) \log \left(\frac{P(x_i; z_i; \theta)}{Q_i(z_i)} \right) \quad \text{--- ①}$$

① hold for any $g(Q_i(z_i))$ and θ

In particular, it holds for $Q_i = Q_i^{(t)}(z_i)$ and $\theta = \theta^{(t+1)}$

$$l(\theta^{(t+1)}) \geq \sum_i \sum_{z_i} Q_i^{(t+1)}(z_i) \log \left(\frac{P(x_i; z_i; \theta^{(t+1)})}{Q_i^{(t)}(z_i)} \right)$$

$$l(\theta^{(t)}) \geq \sum_i \sum_{z_i} Q_i^{(t)}(z_i) \log \left(\frac{P(x_i; z_i; \theta^{(t+1)})}{Q_i^{(t)}(z_i)} \right)$$

what is $l(\theta^{(t+1)}) - l(\theta^{(t)})$?

GMM advantages \rightarrow flexible, density estimation (outlier)

useful as a generative model

1) pick cluster id by sampling from

$$\pi \rightarrow K$$

2) sample from $N(\mu_K, \Sigma_K)$

How to choose K ?

GMM does not work in discontinuous dataset

higher dimension / Numerical instability

In sufficient data (N v/s #params)

Restricting Σ

EM is a general-purpose algorithm

PCA is about find a good coordinate system such that you are able to represent your data better.

Generalized Linear Model (GLM)

Neural Networks

Neuron:

Input

$$\begin{matrix} x_1 \\ x_2 \\ \vdots \\ x_n \\ 1 \end{matrix}$$



Activation

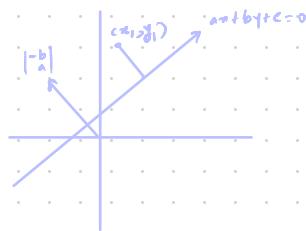
$$s = \sum_{i=1}^n w_i x_i + b \rightarrow \phi(s) \rightarrow \text{Output } f$$

Perceptron Rule : Reduces the number of misclassification (can be proved)

Linear classifier : Represent decision boundary by hyperplane

Batch perceptron

Difference : online learning / Batch learning



d5:title|0:Bittorrent|q:subject|3:ACN|e:keywords|2:ratio:bitmagnet

$$bx - ay + c' = 0$$

$$bx - ay + ay_1 - bx_1 = 0$$

$$b(x - x_1) - a(y - y_1) = 0$$

$$\frac{x - x_1}{a} = \frac{y - y_1}{b} = t$$

$$x = at + x_1, y = bt + y_1$$

$$a(at + x_1) + b(bt + y_1) + c' = 0$$

$$-(a^2 + b^2)t = ax_1 + by_1 + c$$

$$t = -\frac{(ax_1 + by_1 + c)}{a^2 + b^2}$$

$$\sqrt{a^2 t^2 + b^2 t^2} = |t| \sqrt{a^2 + b^2}$$

$$= \frac{|ax_1 + by_1 + c|}{\sqrt{a^2 + b^2}}$$

$$\frac{\sinh(n)}{\cosh(n)} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$e^{inx} = \cos n + i \sin n$$

$$e^{-inx} = \cos n - i \sin n$$

$$\cos n = \frac{(e^{inx} + e^{-inx})}{2}$$

$$\sin n = \frac{(e^{inx} - e^{-inx})}{2i}$$

$$\tanh(n) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$= \frac{1 - e^{-2n}}{1 + e^{-2n}}$$

$$= \frac{e^{2n} - 1}{e^{2n} + 1}$$

$$x \rightarrow \infty, x \rightarrow -\infty$$

$$t \rightarrow 1, t \rightarrow -1$$



Convolutional Neural Network

filters look for patterns

Head + body + Tail + leg = sparrow

each filter will
identify individuals (maybe)

CNN will learn filters + weights + biases

filter \rightarrow pattern detectors

3-d filter : $3 \times 3 \times 3$ filter

Most assume depth (3 in this case)

stride = jump

which filters to construct the network will decide / learn

dropout \rightarrow zero out some % of the image randomly

