

Inferensi Probabilistik

Kecerdasan Artifisial(CIF63310 / 2 sks)

Layout

- Inferensi Probabilistik
- Teori Bayesian
- Naïve Bayes
- Bayesian Belief Network



Building Up
Noble Future

Inferensi Probabilistik

Building Up
Noble Future



Inferensi



- Proses penarikan kesimpulan berdasarkan fakta yang dimiliki
- Inference system: sistem yang dapat melakukan proses reasoning/penalaran berdasarkan fakta/pengetahuan yang dimiliki
- Statistical inference system: proses penalaran yang menggunakan konsep-konsep dan teori statistika (frekuensi, peluang, rata-rata, distribusi, dll)

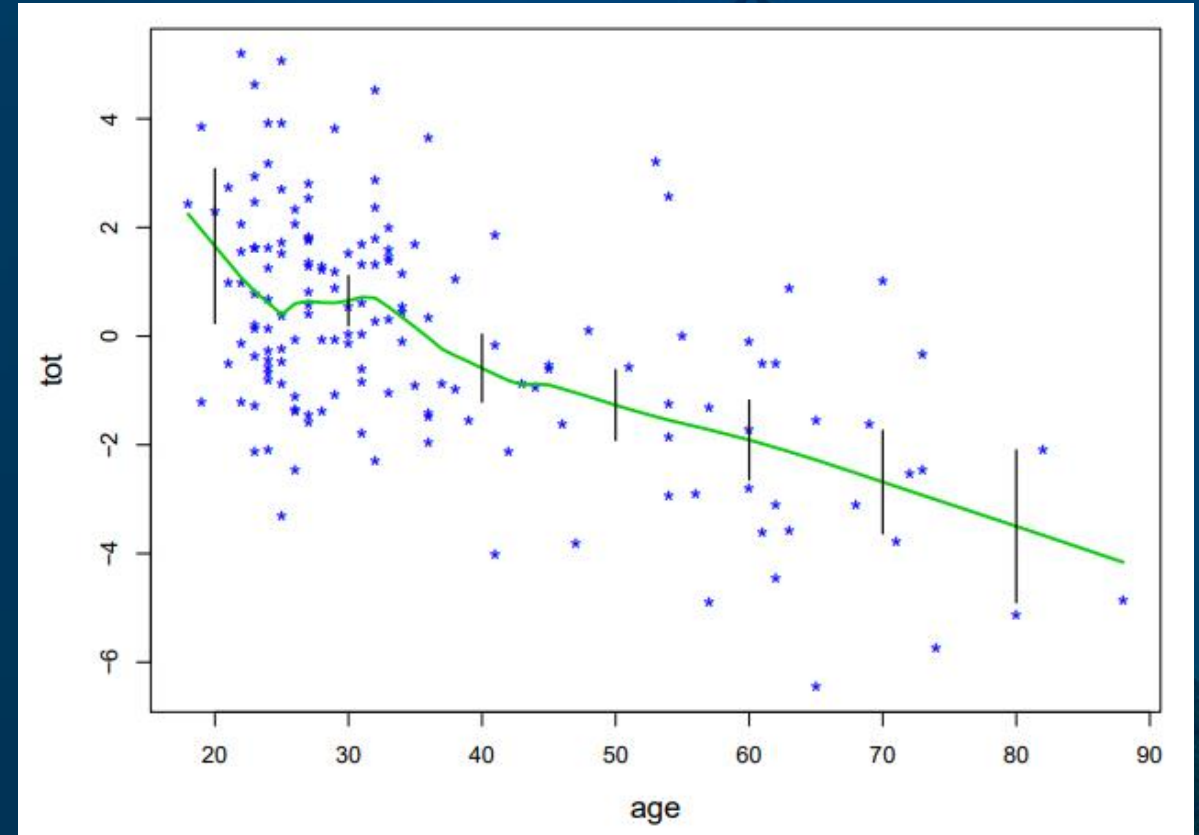
Statistical Inference system



- Frequentist Inference:
 - Kesimpulan didasarkan pada perhitungan frekuensi suatu kejadian yang dilakukan secara random/acak berulang-ulang dalam waktu yang lama.
 - Membutuhkan banyak random observation
- Bayesian Inference:
 - Kesimpulan didasarkan pada derajat kepercayaan terhadap sebuah kejadian

Frequentist Inference

- Contoh kasus: indeks kondisi kesehatan ginjal pasien (tot) berdasarkan umur (age).
- Prediksi tot pasien dapat dilakukan berdasarkan model yang dibangun dari dataset yang dikumpulkan dari banyak pasien sebelumnya.



Efron B. & Hastie T. Computer Age Statistical Inference, 2016, Cambridge University Press

Beberapa Metode Frequentist Inference



- Bootstrap:
 - Menambah data dengan cara melakukan pemilihan sample secara acak dari dataset yang ada secara berulang-ulang, dari pada mencari data baru yang membutuhkan waktu/sumber daya yang banyak.
- Monte Carlo Simulation
 - Menggunakan probability model (Probability Distribution Function-PDF, Cumulative Distribution Function-CDF, atau model lainnya) yang ada untuk men-generate sample baru.
 - Perlu memodelkan kasus (model matematis atau algoritma) berdasarakan dataset yang dimiliki.
 - Bootstrap adalah kasus khusus dari monte carlo simulation

Beberapa Metode Frequentist Inference (2)

- Analisis Regresi
 - Membangun model berdasarkan hasil eksperimen

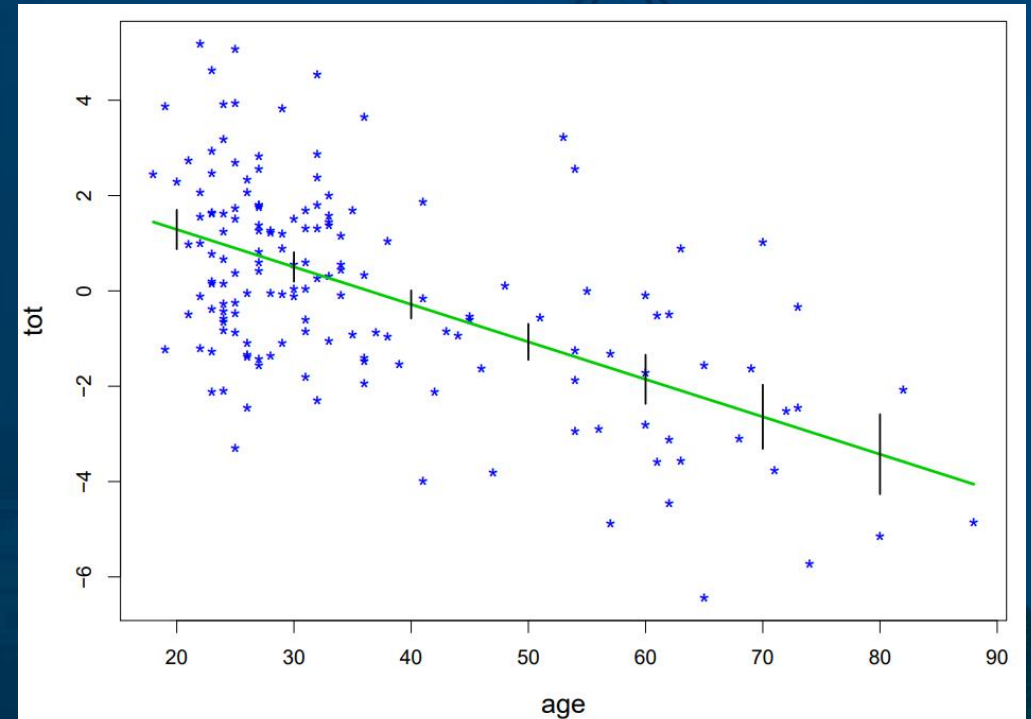


Figure 1.1 Kidney fitness **tot** vs **age** for 157 volunteers. The line is a linear regression fit, showing ± 2 standard errors at selected values of **age**.

Beberapa Video Tutorial



- Bootstrap:
 - <https://www.youtube.com/watch?v=Xz0x-8-cgaQ&t=38s>
- Montecarlo:
 - <https://www.youtube.com/watch?v=7ESK5SaP-bc>
 - <https://www.youtube.com/watch?v=EaR3C4e600k>

Bayesian Inference



- Penalaran dengan menggunakan prior probability, likelihood, dan evidence
- Prior probability: Tingkat kepercayaan awal terhadap peluang terjadinya sebuah kejadian
- Likelihood: Kemungkinan sebuah kejadian data sample muncul di dalam sebuah populasi
- Evidence: fakta yang diketahui saat ini



Peluang kondisional (conditional probability):

$$P(Y|X) = \frac{P(X, Y)}{P(X)}$$

$$P(X|Y) = \frac{P(X, Y)}{P(Y)}$$



Teorema Bayes

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

- Posterior Probability: $P(Y|X)$
- Likelihood: $P(X|Y)$
- Prior Probability: $P(Y)$
- Evidence: $P(X)$

Contoh:



- Given:
 - A doctor knows that meningitis causes stiff neck 50% of the time
 - Prior probability of any patient having meningitis is 1/50,000
 - Prior probability of any patient having a stiff neck is 1/20
- If a patient has a stiff neck, what's the probability he/she has meningitis?

$$P(M|S) = \frac{P(S|M)P(M)}{P(S)} = \frac{0.5 \times 1/50000}{1/20} = 0.0002$$



Naïve Bayes

- Naïve Bayes merupakan bentuk khusus dari teorma bayes
- Dalam teorema bayes, sulit untuk menghitung likelihood pada data sample dengan banyak parameter.
- $P(X|Y) = P(x_1, x_2, \dots, x_n|Y)$
- Naïve Bayes menganggap setiap kejadian adalah kejadian yang saling lepas, sehingga likelihoodnya menjadi
- $P(X|Y) = P(x_1|Y) P(x_2|Y) \dots P(x_n|Y)$

Naïve Bayes Classification



- Misalkan ada sebuah data sampel (record) dengan atribut (x_1, x_2, \dots, x_d)
- Goal: ingin memprediksi kelas untuk record tersebut
- Kelas Y : $\{y_1, y_2, \dots, y_n\}$
- Mencari nilai tertinggi dari peluang $P(y_i|x_1, x_2, \dots, x_d)$

Contoh Kasus



Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- Diketahui record yang ingin diklasifikasi adalah $X = \{\text{refund} = \text{No}, \text{deforced}, \text{income} = 120\text{K}\}$
 - x_1 : refund=No
 - x_2 : deforced
 - x_3 : income=120K
- Nilai $Y = \{\text{Evade} = \text{yes}, \text{Evade} = \text{No}\}$
- Hitung posterior probability:
 - $P(\text{Evade} = \text{yes} \mid X) = \dots$
 - $P(\text{Evade} = \text{No} \mid X) = \dots$

Likelihood X terhadap yes



Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- $P(yes|X) = \frac{P(X|yes)P(yes)}{P(X)}$
- Menghitung likelihood $P(X | yes)$:
 - $= P(x_1, x_2, x_3 | yes)$
 - $= P(x_1 | yes) P(x_2 | yes) P(x_3 | yes)$
 - $= P(\text{refund=no} | yes) * P(\text{divorced} | yes) * P(\text{income} = 120K | yes)$
 - $= 3/3 * 1/3 * 0/3$
 - $= 0$

Likelihood X terhadap no



Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- $P(no|X) = \frac{P(X|no)P(no)}{P(X)}$
- Menghitung likelihood $P(X | no)$:
 - $= P(x_1, x_2, x_3 | no)$
 - $= P(x_1 | no) P(x_2 | no) P(x_3 | no)$
 - $= P(\text{refund=no} | no) * P(\text{divorced} | no) * P(\text{income} = 120K | no)$
 - $= 4/7 * 1/7 * 1/7$
 - $= 0.012$

Prior yes dan no

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	no Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- $P(yes) = \frac{3}{10}$

- $P(no) = \frac{7}{10}$

Evidence



Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- $P(X) = P(x_1, x_2, x_3)$
- $= P(x_1)P(x_2)P(x_3)$
- $= P(\text{refund=no}) * P(\text{divorced}) * P(\text{income}=120K)$
- $= 7/10 * 2/10 * 1/10$
- $= 0.014$

asumsi variable saling lepas

Posteriror Probability



- $P(yes|X) = \frac{P(X|yes)P(yes)}{P(X)} = \frac{0*0.3}{0.014} = 0$
- $P(no|X) = \frac{P(X|no)P(no)}{P(X)} = \frac{0.012*0.7}{0.014} = 0.58$
- Karena $P(no|x) > P(yes|X)$, maka kelas yang sesuai untuk data tersebut adalah evade = no
- Perhatikan nilai evidence $P(X)$ pada kedua perhitungan di atas. Karena yang diperlukan untuk keputusan adalah posterior terbesar, dan nilai $P(X)$ sama antara kedua nilai posterior, maka perhitungan evidence boleh diabaikan

Continue Probability



- Perhatikan nilai likelihood untuk kelas yes
 - $= P(\text{refund=no} \mid \text{yes}) * P(\text{deforced} \mid \text{yes}) * P(\text{income} = 120\text{K} \mid \text{yes})$
 - $= 3/3 * 1/3 * 0/3$
 - $= 0$
- Nilainya = 0 karena kolom income tidak ada yang tepat 120K
- Nilai income dalam dunia nyata selalu spesifik untuk tiap orang dan tidak bisa dihitung secara diskrit.
- Ada 2 teknik penyelesaiannya:
 - Discretization: buat interval pada nilai income
 - Menghitung peluang kontiny dengan probability density estimation

Probability Density Estimation



- Distribusi Normal (Normal Distribution)

$$P(X_i|Y_j) = \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} e^{-\frac{(X_i - \mu_{ij})^2}{2\sigma_{ij}^2}}$$

σ_{ij}^2 : sample variance
 μ_{ij} : sample mean

- Untuk (income, class = no):

- Sample mean = 110
- Sample variance = 2975

$$P(\text{Income} = 120|\text{No}) = \frac{1}{\sqrt{2\pi}(54.54)} e^{-\frac{(120-110)^2}{2(2975)}}$$
$$= 0.0072$$

- Untuk (income, class = yes):

- Sample mean = 90
- Sample variance = 25

$$P(\text{Income} = 120|\text{Yes}) = \frac{1}{\sqrt{2\pi}(5)} e^{-\frac{(120-90)^2}{2(25)}}$$
$$= 1.2 \times 10^{-9}$$

Likelihood X terhadap yes (hitung ulang)



Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- $P(yes|X) = \frac{P(X|yes)P(yes)}{P(X)}$
- Menghitung likelihood $P(X | yes)$:
 - $= P(x_1, x_2, x_3 | yes)$
 - $= P(x_1 | yes) P(x_2 | yes) P(x_3 | yes)$
 - $= P(\text{refund=no} | yes) * P(\text{deforced} | yes) * P(\text{income} = 120K | yes)$
 - $= 3/3 * 1/3 * 1.2 \times 10^{-9}$
 - $= 4 \times 10^{-10}$

Likelihood X terhadap no



Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- $P(no|X) = \frac{P(X|no)P(no)}{P(X)}$

- Menghitung likelihood $P(X | no)$:

- $= P(x1, x2, x3 | no)$

- $= P(x1 | no) P(x2 | no) P(x3 | no)$

- $= P(\text{refund=no} | no) * P(\text{deforced} | no) * P(\text{income} = 120K | no)$

- $= 4/7 * 1/7 * 0.0072$

- $= 0.006$

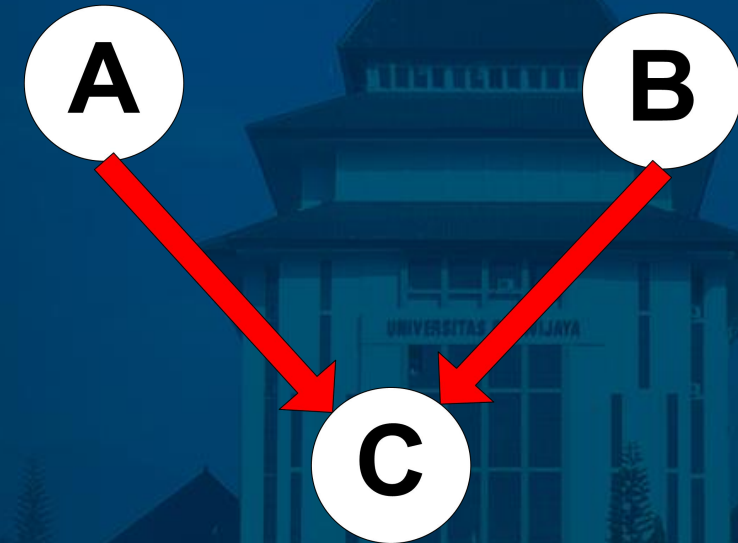
Posteriror Probability (hitung ulang)

- $P(yes|X) = \frac{P(X|yes)P(yes)}{P(X)} = \frac{4 \times 10^{-10} * 0.3}{0.014} = 8.5 \times 10^{-9}$
- $P(no|X) = \frac{P(X|no)P(no)}{P(X)} = \frac{0.006 * 0.7}{0.014} = 0.3$
- Karena $P(no|x) > P(yes|X)$, maka kelas yang sesuai untuk data tersebut adalah evade = no
- Perhatikan nilai evidence $P(X)$ pada kedua perhitungan di atas. Karena yang diperlukan untuk keputusan adalah posterior terbesar, dan nilai $P(X)$ sama antara kedua nilai posterior, maka perhitungan evidence boleh diabaikan

Bayesian Belief Network



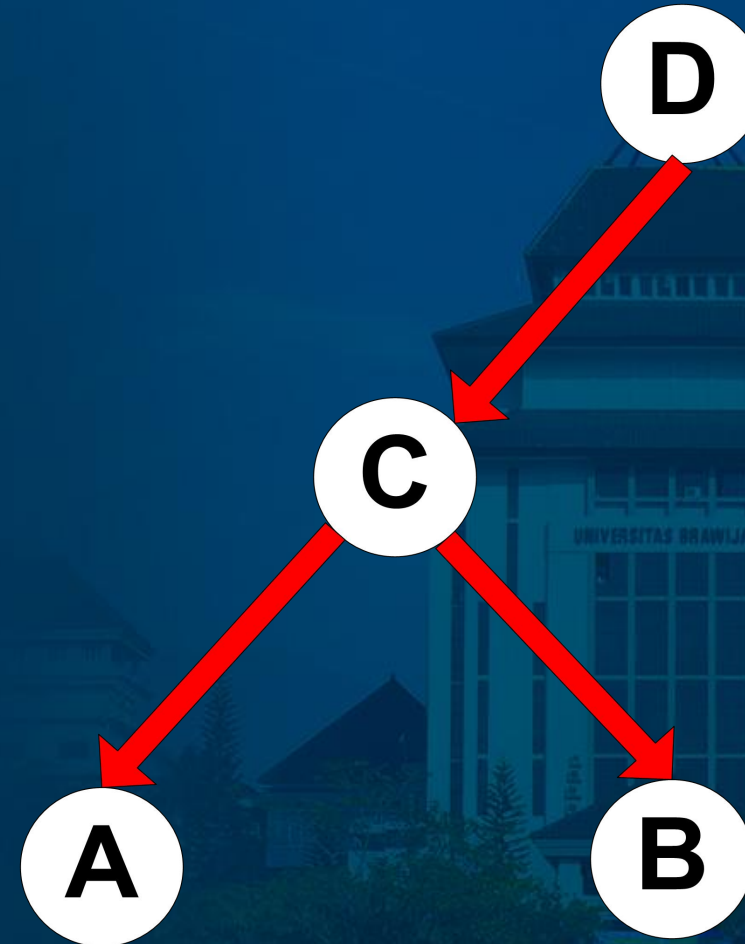
- Representasi grafis dari hubungan probabilitas antar random variable
- Terdiri dari:
 - Directed acyclic graph
 - Node adalah variable
 - Edge adalah hubungan antar variable
 - Table probabilitas (probability table) yang menghubungkan tiap node dengan parentnya



Conditional independence

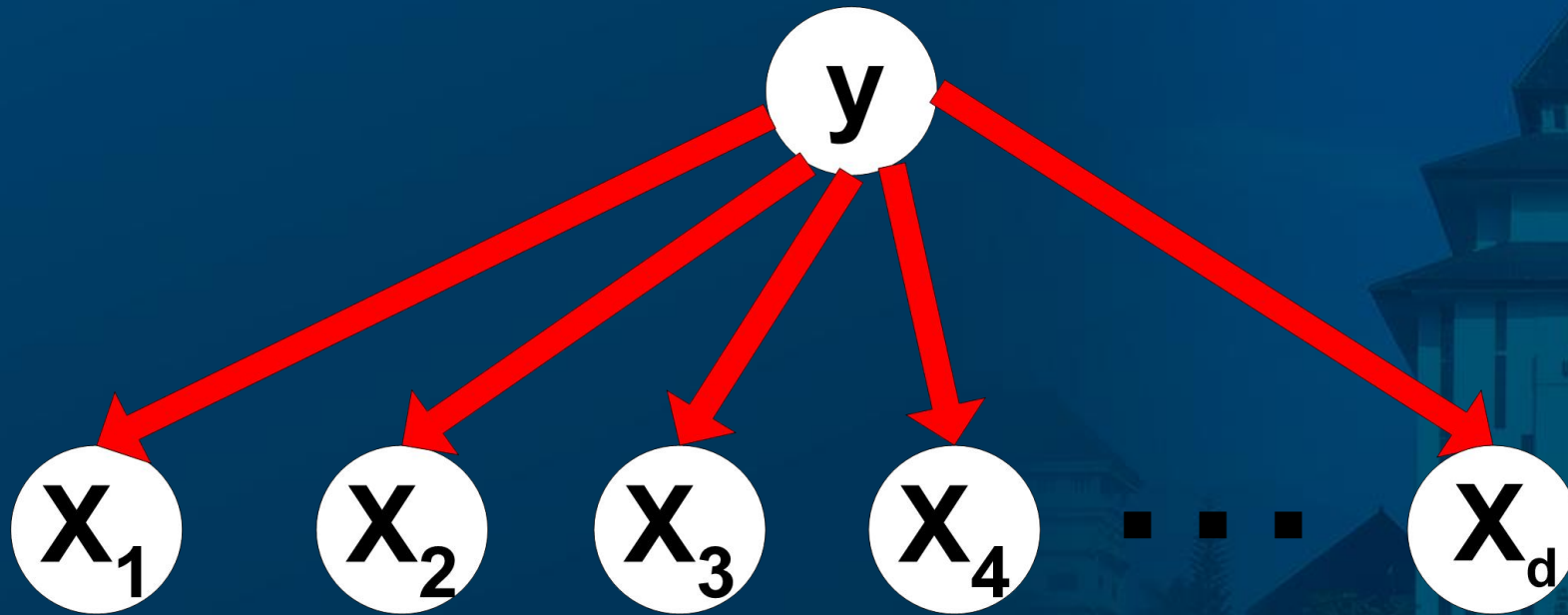


- D is parent of C
 - A is child of C
 - B is descendant of D
 - D is ancestor of A
-
- A node in a Bayesian network is conditionally independent of all of its nondescendants, if its parents are known



Conditional independence

- Asumsi dalam Naïve Bayes

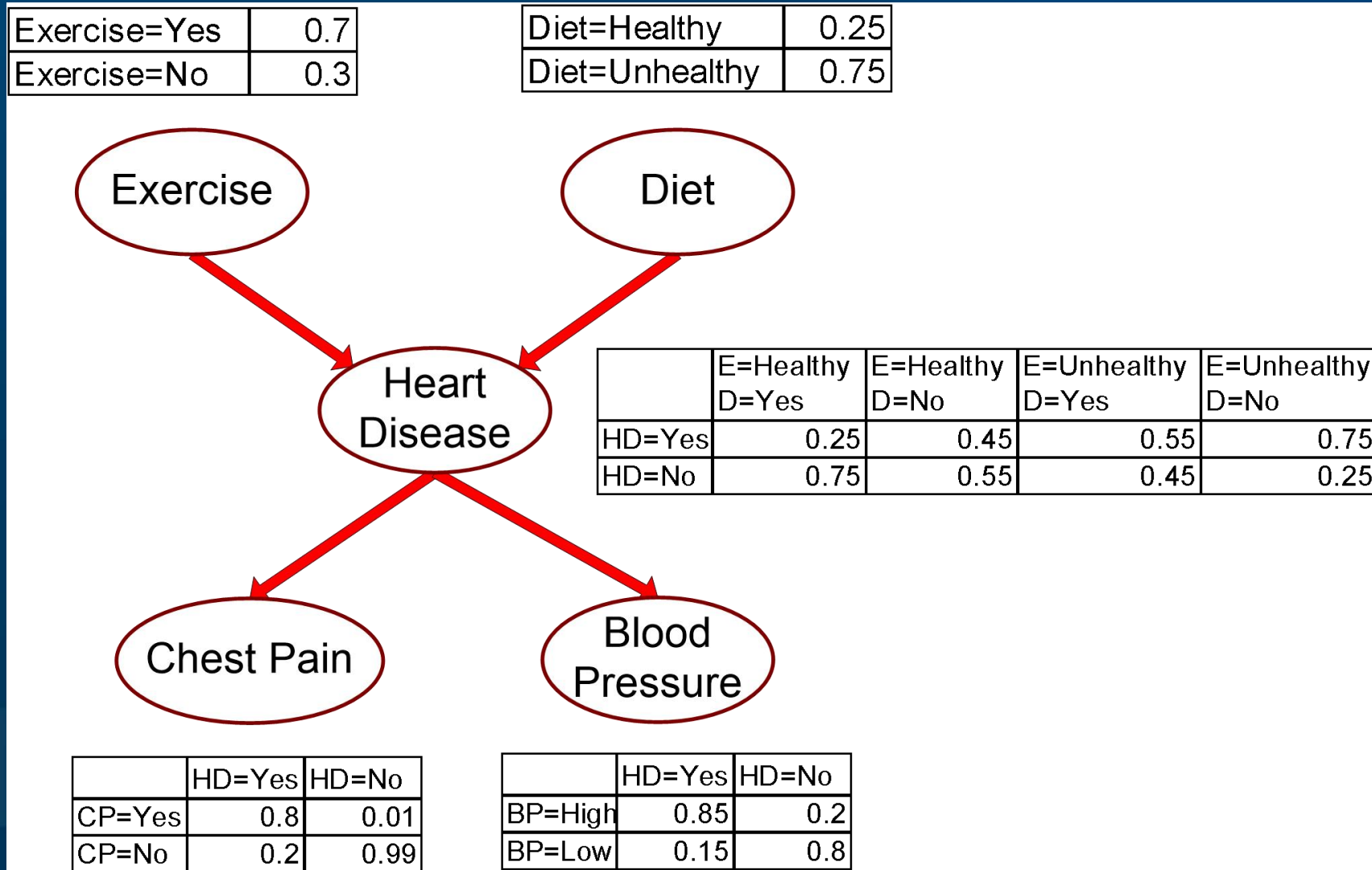


Probability Tables



- Jika sebuah node tidak memiliki parent, table probability-nya diisi dengan prior probability $P(x)$
- Jika variable x hanya memiliki sebuah parent (y), maka table probability-nya diisi dengan conditional probability $P(x|y)$
- Dan jika variable x memiliki banyak parents (y_1, y_2, \dots, y_k), maka table probability-nya diisi dengan conditional probability $P(x|y_1, y_2, \dots, y_k)$

Contoh Bayesian Belief Network



Terima Kasih

Building Up
Noble Future

