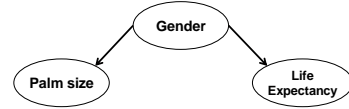## From Data To Understanding

- In machine learning, maintain critical perspective
  - Making predictions is only part of the story
  - Also try to get some <u>understanding</u> of the domain
- Example
  - True statement: palm size negatively correlates with life expectancy
    - The larger your palm size, the shorter your life (on average)
  - Why?
    - Women have smaller palms than men on average
    - Women live 5 years longer than men on average
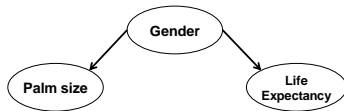  - Sometimes you need better model of your domain!

## Bayesian Networks

- Bayesian Network
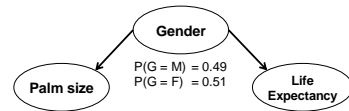  - Graphical representation of joint probability distribution



  - Node: random variable
  - Arc (X, Y): variable X has direct influence on variable Y
    - Call X a "parent" of Y
  - Each node X has conditional probability: P(X | parents(X))
  - Graph has no cycles (loops by following arcs)
    - Called "Directed Acyclic Graph" (DAG)

## Network Shows Conditional Independence



- Conditional independence encoded in network
  - Each node (variable) is conditionally independent of its non-descendants, given its parents
  - In network above Palm Size and Life Expectancy are conditionally independent, given Gender
    - Formally: P(PS, LE | G) = P(PS | G) P(LE | G)
- Network structure provides insight about domain
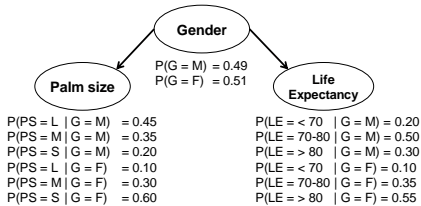
## Conditional Probability Tables



P(G = M) = 0.49
P(G = F) = 0.51

P(PS = L | G = M) = 0.45
P(PS = M | G = M) = 0.35
P(PS = S | G = M) = 0.20
P(PS = L | G = F) = 0.10
P(PS = M | G = F) = 0.30
P(PS = S | G = F) = 0.60

P(LE = < 70  | G = M) = 0.20
P(LE = 70-80 | G = M) = 0.50
P(LE = > 80  | G = M) = 0.30
P(LE = < 70  | G = F) = 0.10
P(LE = 70-80 | G = F) = 0.35
P(LE = > 80  | G = F) = 0.55

- Each node has conditional probability table (CPT)
  - For node X: P(X | Parents(X))
  - Conditional independence modularizes joint probability:

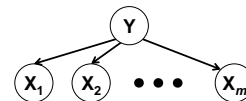$$P(X_1, X_2, ..., X_n) = \prod_{i=1}^{n} P(X_i | \mathrm{Parents}(X_i))$$

## Efficient Representation



P(G = M) = 0.49
P(G = F) = 0.51

P(PS = L | G = M) = 0.45
P(PS = M | G = M) = 0.35
P(PS = S | G = M) = 0.20
P(PS = L | G = F) = 0.10
P(PS = M | G = F) = 0.30
P(PS = S | G = F) = 0.60

P(LE = < 70  | G = M) = 0.20
P(LE = 70-80 | G = M) = 0.50
P(LE = > 80  | G = M) = 0.30
P(LE = < 70  | G = F) = 0.10
P(LE = 70-80 | G = F) = 0.35
P(LE = > 80  | G = F) = 0.55

- Each node has conditional probability table (CPT)
  - Reduces number of parameters needed in model
  - Normally, need 2 x 3 x 3 – 1 = 18 – 1 = 17 parameters
  - Here, need (2 – 1) + (6 – 2) + (6 – 2) = 9 parameters

## Bayesian Network for Naïve Bayes

- Welcome back, Naïve Bayes…
  - Now with new and improved "Bayesian Network" flavor!



  - Network structure encodes assumption:

$$P(X | Y) = P(X_1, X_2, ... X_m | Y) = \prod_{i=1}^{m} P(X_i | Y)$$
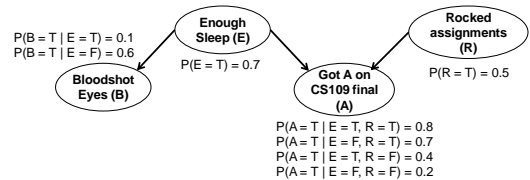
  - Full joint distribution can be computed as:

$$P(X, Y) = P(Y)P(X | Y) = P(Y) \prod_{i=1}^{m} P(X_i | Y)$$

## "Evidence" in Bayesian Networks

- In many machine learning examples:
  - We observe all $X_1, X_2, \ldots, X_m$ input variables and predict single output variable Y
- In general case of probabilistic inference:
  - Have a set of random variables $X_1, X_2, \ldots, X_m$
  - *Some* of the variables $X_1, X_2, \ldots, X_m$ are observed
    - Call observed variables $E_1, E_2, \ldots, E_k$ (E for "evidence")
  - Want to determine probability of some set of *unobserved* variables given the observed evidence
    - Call unobserved variables we care about $Y_1, Y_2, \ldots, Y_c$
  - Formally, want: $P(Y_1, Y_2, \ldots, Y_c \mid E_1, E_2, \ldots, E_k)$

---

## Evaluation of Evidence

- Consider the following Bayes Net:

$P(B = T \mid E = T) = 0.1$
$P(B = T \mid E = F) = 0.6$

**Enough Sleep (E)**
$P(E = T) = 0.7$

**Rocked assignments (R)**
$P(R = T) = 0.5$

**Bloodshot Eyes (B)**

**Got A on CS109 final (A)**

$P(A = T \mid E = T, R = T) = 0.8$
$P(A = T \mid E = F, R = T) = 0.7$
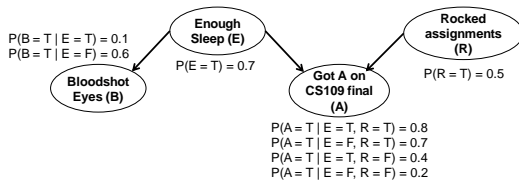$P(A = T \mid E = T, R = F) = 0.4$
$P(A = T \mid E = F, R = F) = 0.2$

- Determine $P(A = T \mid B = T, R = T)$
- Sum over unseen variables:

$$P(A=T \mid B=T, R=T) = \frac{P(A=T, B=T, R=T)}{P(B=T, R=T)} = \frac{\sum_{E=T,F} P(A=T, B=T, R=T, E)}{\sum_{E=T,F} \sum_{A=T,F} P(B=T, R=T, E, A)}$$

---

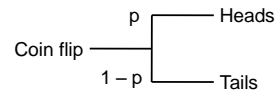## Evaluation of Evidence

- Consider the following Bayes Net:

$P(B = T \mid E = T) = 0.1$
$P(B = T \mid E = F) = 0.6$

**Enough Sleep (E)**
$P(E = T) = 0.7$

**Rocked assignments (R)**
$P(R = T) = 0.5$

**Bloodshot Eyes (B)**

**Got A on CS109 final (A)**

$P(A = T \mid E = T, R = T) = 0.8$
$P(A = T \mid E = F, R = T) = 0.7$
$P(A = T \mid E = T, R = F) = 0.4$
$P(A = T \mid E = F, R = F) = 0.2$

- Determine $P(A = T \mid B = T, R = T)$
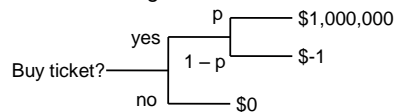- Note that joint probability decomposes as:
  $$P(A, B, E, R) = P(E)P(B \mid E)P(R)P(A \mid E, R)$$
- Plug in values from CPTs to compute joint probabilities

---

## Probability Tree
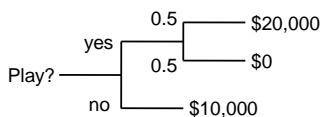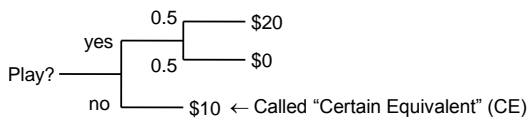
- Model outcomes of probabilistic events with tree

Coin flip
- p — Heads
- 1 − p — Tails

- Useful for modeling decisions

Buy ticket?
- yes
  - p — $1,000,000
  - 1 − p — $-1
- no — $0

  - Payoffs: yes = p(1000000) + (1 − p)(-1), no = 0

---

## Let's Play a Game

- Which choice would you make?

Play?
- yes
  - 0.5 — $20
  - 0.5 — $0
- no — $10  ← Called "Certain Equivalent" (CE)

Play?
- yes
  - 0.5 — $20,000
  - 0.5 — $0
- no — $10,000

  - Certain equivalent is how much game is worth to you

---

## Utility

- Utility U(x) is "value" you derive from x

Play?
- yes
  - 0.5 — $20,000
  - 0.5 — $0
- no — $10,000

Play?
- yes
  - 0.5 — U($20,000)
  - 0.5 — U($0)
- no — U($10,000)

  - Can be monetary, but often includes intangibles
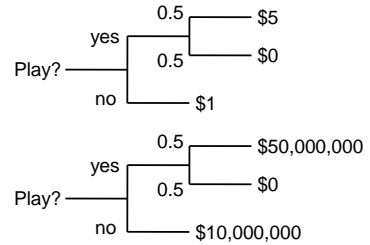    - E.g., quality of life, life expectancy, personal beliefs, etc.

## Risk Premium

- A slightly different game:

```
                    0.5 ┌──── $20,000
            yes ────┤
                    0.5 └──── $0
Play? ──────┤
            no  └──── $8,000    Certain Equivalent
```
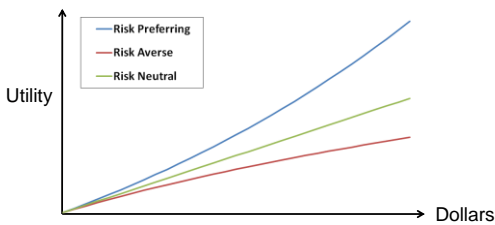
  - Expected monetary value (EMV) = expected dollar value of game (here = $10,000)
  - Risk premium = EMV – CE = $2,000
    - How much you would pay (give up) to avoid risk
    - This is what insurance is all about
    - It's also what the show "Deal or No Deal" is based on

## Non-Linear Utility of Money

- These two choices are different for most people

```
                    0.5 ┌──── $5
            yes ────┤
                    0.5 └──── $0
Play? ──────┤
            no  └──── $1
```

```
                    0.5 ┌──── $50,000,000
            yes ────┤
                    0.5 └──── $0
Play? ──────┤
            no  └──── $10,000,000
```

## Utility Curves



- Risk Preferring
- Risk Averse
- Risk Neutral

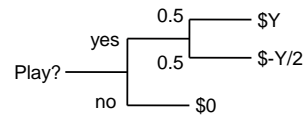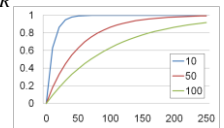Utility (vertical axis), Dollars (horizontal axis)

- Utility curve determines your "risk preference"
  - Can be different in different parts of the curve

## Exponential Utility Curves

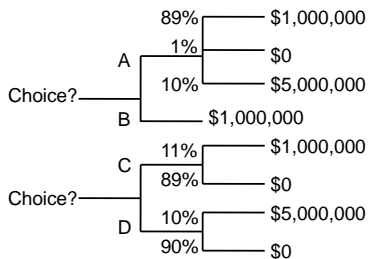- Many people have exponential utility curves

$$U(x) = 1 - e^{-x/R}$$



  - R is your "risk tolerance"
  - Larger R = less risk aversion
    - Makes utility function more linear
  - R ≈ highest value of Y for which you would play:

```
                    0.5 ┌──── $Y
            yes ────┤
                    0.5 └──── $-Y/2
Play? ──────┤
            no  └──── $0
```

## How Irrational Are You?

- Which option would you choose?

```
              89% ┌──── $1,000,000
          A ──┤1% ├──── $0
             10% └──── $5,000,000
Choice? ──┤
          B ──────── $1,000,000
```

```
              11% ┌──── $1,000,000
          C ──┤89%└──── $0
Choice? ──┤
             10% ┌──── $5,000,000
          D ──┤90%└──── $0
```

  - How many chose B and D?
    - You are inconsistent with utility theory (the Allais Paradox)

## Micromort

- A **micromort** is 1 in 1,000,000 chance of death
  - How much would you need to be paid to take on the risk of a micromort?
  - How much would you pay to avoid a micromort?
    - P(die in plane crash) ≈ 1 in 1,500,000
    - P(killed by lightning) ≈ 1 in 1,400,000
  - How much would you need to be paid to take on a decimort (1 in 10 chance of death)?
  - If you think this is morbid, companies actually do this
    - Car manufacturers
    - Insurance companies

## Let's Do a Real Test

- Game set-up
  - I will flip a fair coin
  - If "heads", you win $50. If "tails", you win $0
  - How much would you be willing to pay me to play?
    - $1 ?
    - $10 ?
    - $20 ?
    - $24.99 ?
    - $25.01 ?
    - $35 ?
  - Maximal value?
    - Come on down!
    - How did you determine that value?

## Just For Fun…

- Say we consider two batters in baseball
  - Batting averages of Player A and Player B for 2 years:

|  | Year 1 | Year 2 | Combined |
|---|---|---|---|
| Player A | .250 | .314 | .310 |
| Player B | .253 | .321 | .270 |

  - So is Player B the better player?
  - Is this possible?
  - In fact, it happened:

|  | 1995 | 1996 | Combined |
|---|---|---|---|
| Derek Jeter | 12/48 = .250 | 183/582 = .314 | 195/630 = .310 |
| David Justice | 104/411 = .253 | 45/140 = .321 | 149/551 = .270 |

  - This is known as Simpson's Paradox

## So Which Medicine Should You Choose?

- Consider medicine to treat a disease:
  - Success rates of two medicines on disease:

|  | Medicine A | Medicine B |
|---|---|---|
| Success rate | 273/350 = 78% | 289/350 = **83%** |

  - Seems reasonable to choose B
  - But now, let's consider your gender:

|  | Medicine A | Medicine B |
|---|---|---|
| Female success | 81/87 = **93%** | 234/270 = 87% |
| Male success | 192/263 = **73%** | 55/80 = 69% |
| Overall success | 273/350 = 78% | 289/350 = **83%** |

  - You're either male or female...
    - Results from gender preferences for different treatments