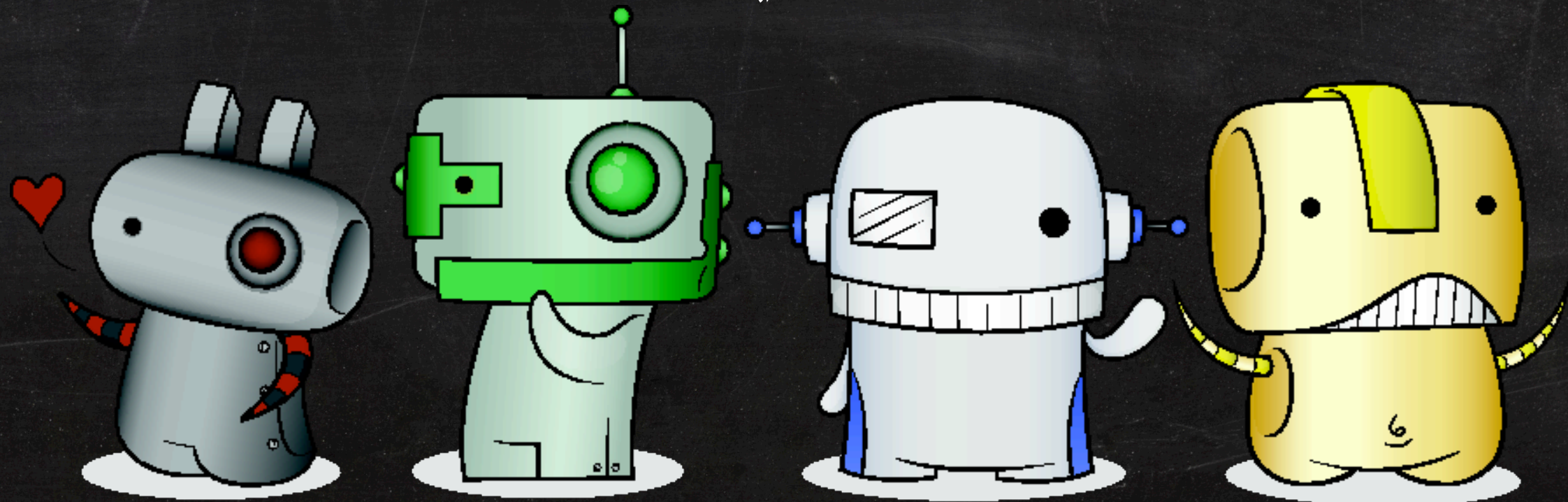


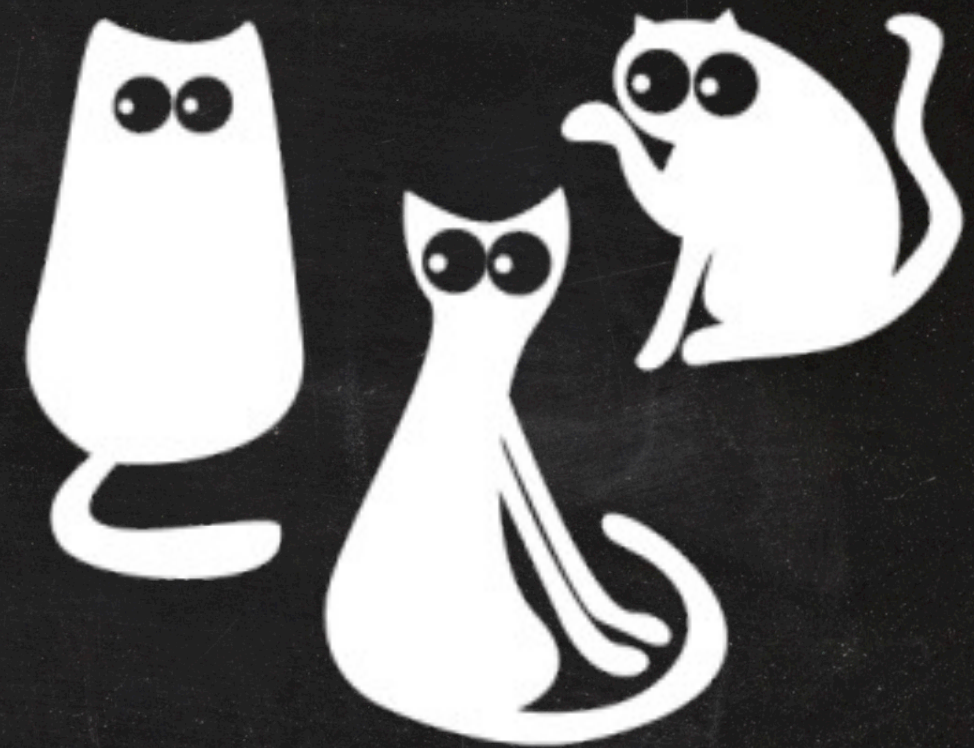
JAMES MCGIVERN

PROBABLY, DEFINITELY,  
MAYBE



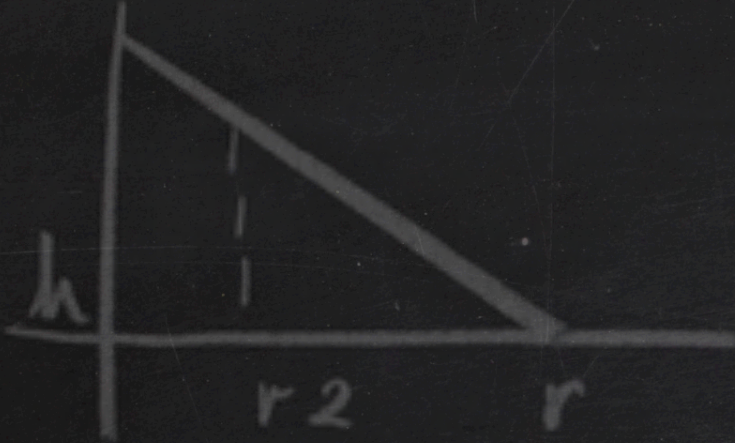
# ABOUT ME

 rockshore

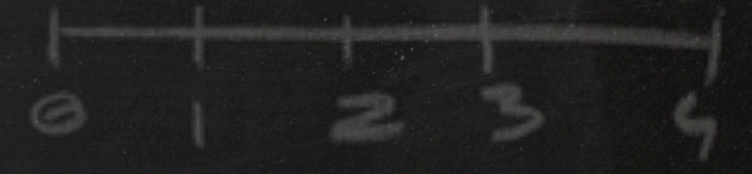


I Talk Fast!

# MATHS WARNING



$$x > \frac{y}{5} = \sqrt{\frac{582}{11}} \times \frac{18}{7}$$



$$c^2 = (AB)^2$$

$$f = \frac{dy}{dx} = x^4$$

$$M = \sqrt{\frac{2 \cdot 7 \cdot 10^2}{2 \cdot 14 \cdot 10^4}}$$



$$f(x) = x^2$$

$$2 \frac{1}{2} = \frac{5}{2}x$$



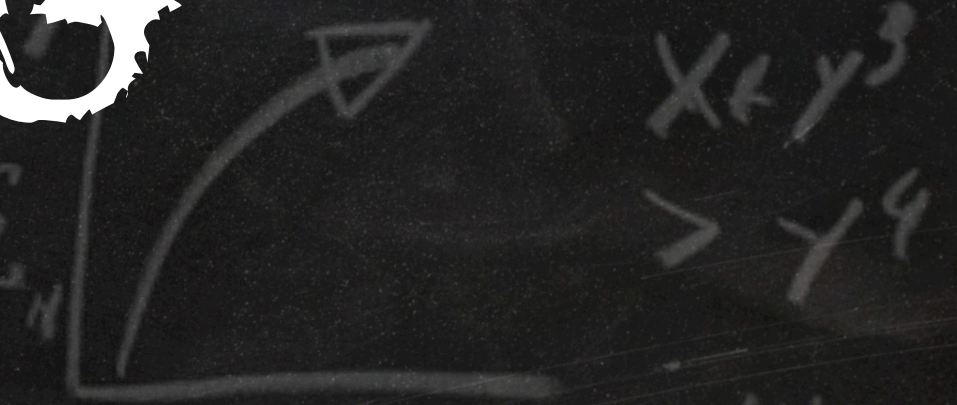
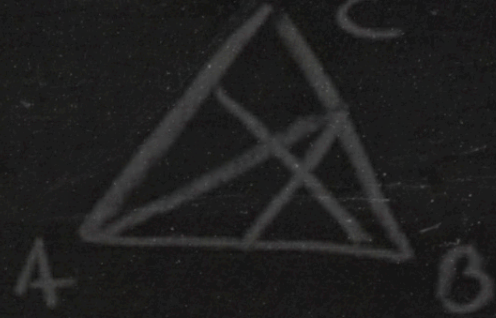
$$\frac{k}{x}$$

$$\frac{1}{5^2} + \frac{3}{2} \left( \frac{1}{5-2} - \frac{1}{5+2} \right) + \frac{1}{2} \sqrt{c}$$

$$+ 3 \frac{1}{4}$$

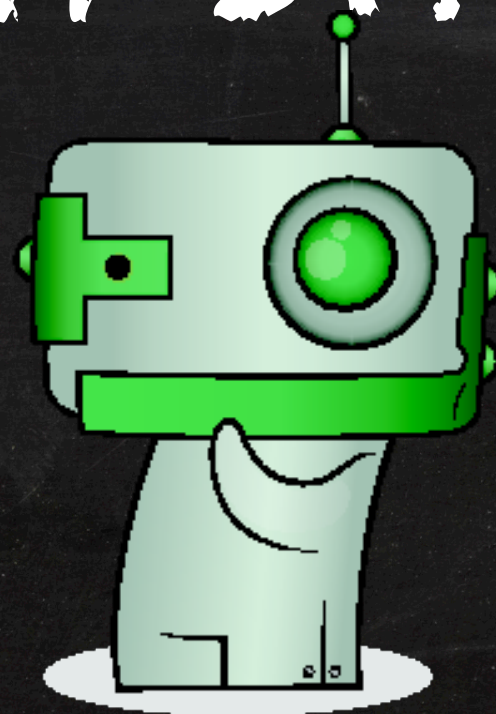
# MATHS WARNING

$$\begin{aligned} b^2 &= \cos^2 C - 2ab \cos C \\ &= a^2 (\cos^2 C + \sin^2 C) \\ &= a^2 + c^2 - 2ac \cos C \end{aligned}$$



$$A = \pi r^2 \times \frac{11}{25}$$

- CHAPTER 1 -  
BAYESIAN PROBABILITY  
&  
BAYESIAN STATISTICS



# THE MEAN

Given a set of variables

$$X = \{x_1, \dots, x_n\}$$

we define the mean (average)

$$\bar{X} = \frac{1}{n} \cdot \sum_{i=0}^n x_i$$

# EXAMPLE: AMAZON

- ▶ Users can rate a product by voting 1-5 stars
- ▶ product rating is the mean of the user votes

Q: how can we rank products with different number of votes?



# SIMPLE "BAYESIAN RANKING"

$$\text{rank} = \frac{Cm + Rv}{m + v}$$

- ▶ C - the mean vote across all items
- ▶ v - number of votes for a given item
- ▶ R - the mean rating of the item
- ▶ m - number of votes required to be in top n percentile

Book	Number of votes (v)	Average Rating (R)	Bayesian Rank
A	100	5	4.333333...
B	70	5	4.17
C	50	4	3.5
D	30	4	3.375
E	20	3.5	3.14
F	30	3	3
G	5	2	2.91

$C = 3$

$m = 50$

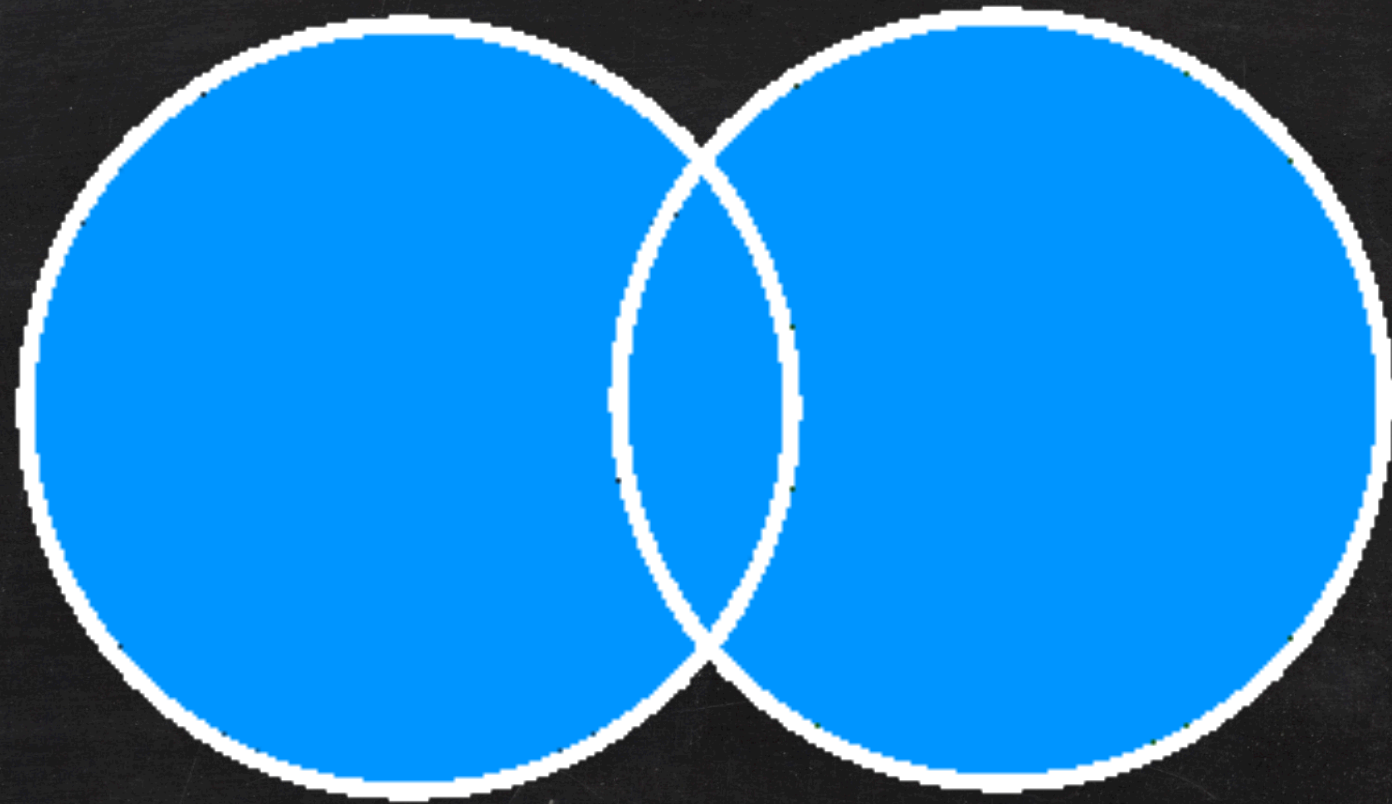


# A DETOUR IN TO PROBABILITY BASICS

# EVENTS

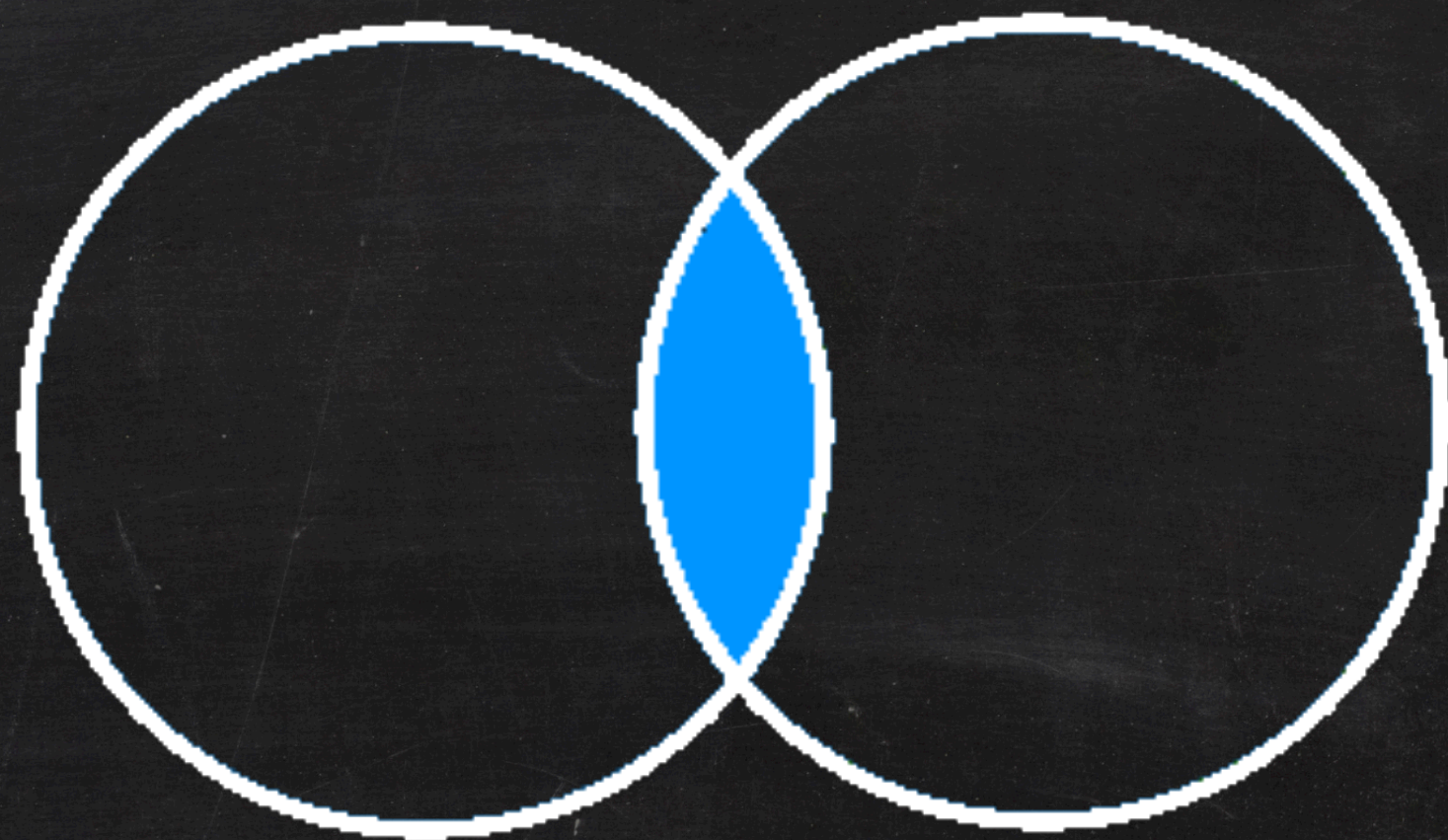
- ▶ Consider an experiment whose set of all possible outcomes  $\Omega$ , called the sample space, is  $\{x_1, \dots, x_n\}$
- ▶ We define an event  $E$  as a subset of  $\Omega$  and say that  $E$  occurs iff the experiment outcomes equal  $E$

# UNION



$$E_1 \cup E_2 \cup \dots \cup E_n = \bigcup_{i=1}^n E_i$$

# INTERSECTION



$$E_1 \cap E_2 \cap \dots \cap E_n = \bigcap_{i=1}^n E_i$$

# PROBABILITY AXIOMS: I

- ▶ We denote the probability of an event  $A$  by  $P(A)$
- ▶ For any event  $A$ ,  $0 \leq P(A) \leq 1$
- ▶ The certain event,  $\Omega$ , always occurs and  $P(\Omega)=1$
- ▶ The impossible event  $\emptyset$  never occurs and  $P(\emptyset)=0$

# PROBABILITY AXIOMS: 2

- ▶ We say that events  $A$  and  $B$  are disjoint if  $A \cap B = \emptyset$
- ▶ if  $A$  and  $B$  are disjoint then  $P(A \cup B) = P(A) + P(B)$
- ▶ for a set of disjoint events, the addition law gives us:

$$P\left(\bigcup_{i=1}^n E_i\right) = \sum_{i=1}^n P(E_i)$$

# PROBABILITY LEMMAS

▶ For any event  $E$ ,  $P(E^c) = P(\neg E) = 1 - P(E)$

▶  $P(A - B) = P(A) - P(A \cap B)$

▶ If  $A \subset B$  then  $P(A) \leq P(B)$

▶  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

$$P\left(\bigcup_{i=1}^n E_i\right) = \sum_{i=1}^n P(E_i) - \sum_{i < j} P(E_i \cap E_j)$$

$$+ \sum_{i < j < k} P(E_i \cap E_j \cap E_k) - \dots + (-1)^{n-1} P\left(\bigcap_{i=1}^n E_i\right)$$

# RANDOM VARIABLES

- ▶ Consider a random variable  $X$ , then  $\{X \leq x\}$  is the event that  $X$  has a value less than or equal to the real number  $x$ . Hence the probability that this event occurs is  $P(X \leq x)$
- ▶ If we allow  $x$  to vary we can define the distribution function

$$F(x) = P(X \leq x) \quad -\infty < x < \infty$$

- ▶ Note that:
  - ▶  $P(X > x) = 1 - F(x)$
  - ▶  $P(a < X < b) = F(b) - F(a)$



# PROBABILITY MASS FUNCTION

The probability mass function (PMF) of  $X$

$$f_x(x) = P(X = x) = P(\{s \in S : X(s) = x\})$$

is a probability measure of the possible values for the

random variable. Of course

$$\sum_{x \in A} f_x(x) = 1$$

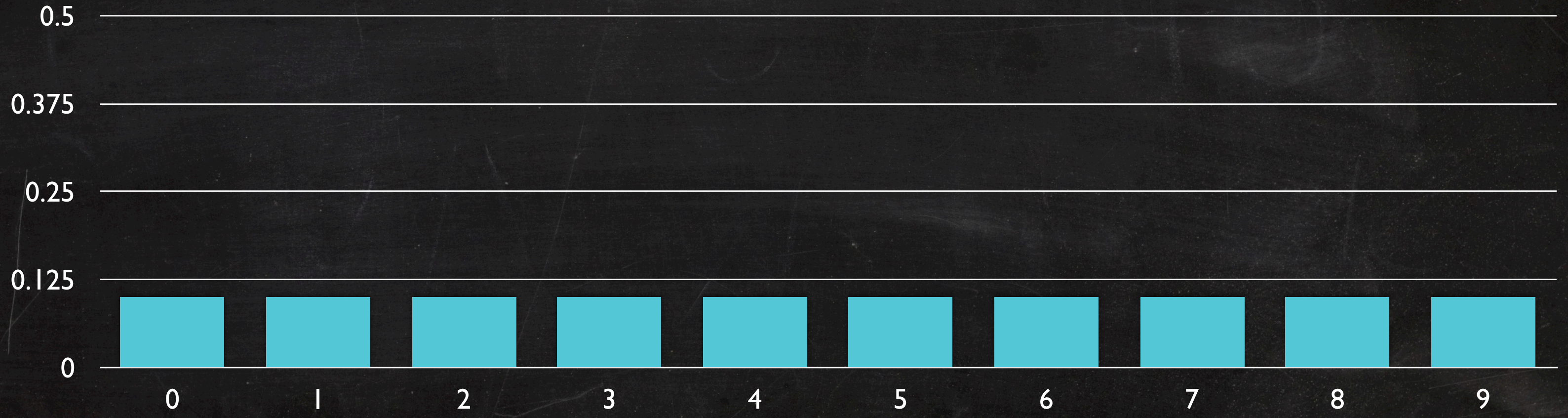
PMF for a fair  
6 sided dice

$$f_x = \begin{cases} \frac{1}{6} & x \in \{1, 2, 3, 4, 5, 6\} \\ 0 & \text{otherwise} \end{cases}$$

# EXAMPLE: PROOF OF 2 FAIR DIE

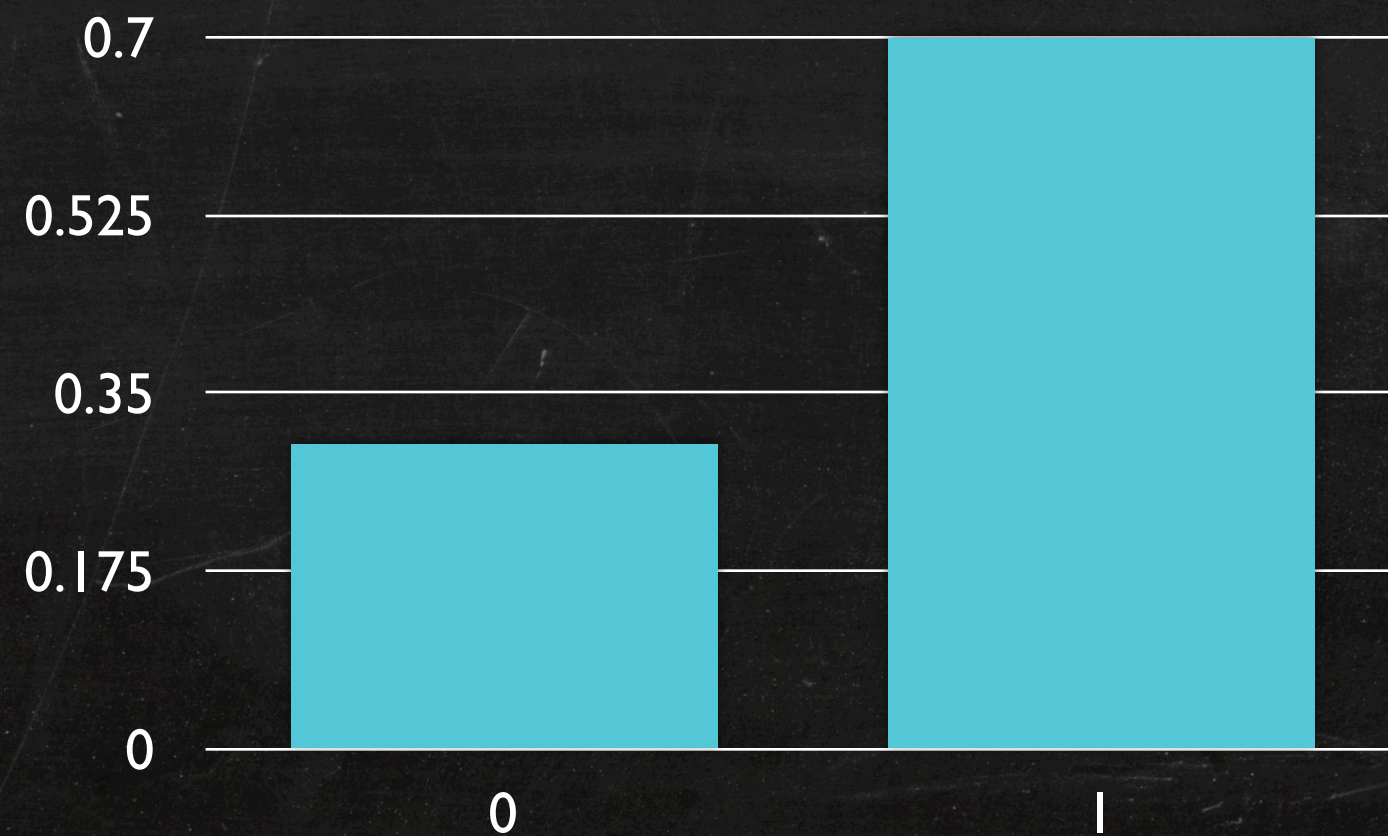
1/36	2/36	3/36	4/36	5/36	6/36	5/36	4/36	3/36	2/36	1/36
					1,6					
				1,5	2,5	2,6				
			1,4	2,4	3,4	3,5	3,6			
		1,3	2,3	3,3	4,3	4,4	4,5	4,6		
	1,2	2,2	3,2	4,2	5,2	5,3	5,4	5,5	5,6	
1,1	2,1	3,1	4,1	5,1	6,1	6,2	6,3	6,4	6,5	6,6
2	3	4	5	6	7	8	9	10	11	12

# UNIFORM DISTRIBUTION



# BERNOULLI DISTRIBUTION

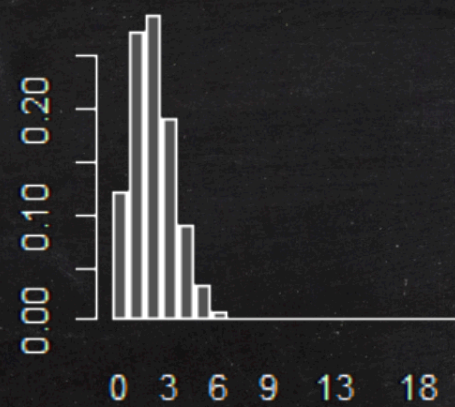
$$F(k, p) = \begin{cases} p & \text{if } k = 1 \\ 1 - p & \text{if } k = 0 \end{cases}$$



# BINOMIAL DISTRIBUTION

$$F(k, n, p) = P_X(k) = \binom{n}{k} p^k (1-p)^{n-k}$$

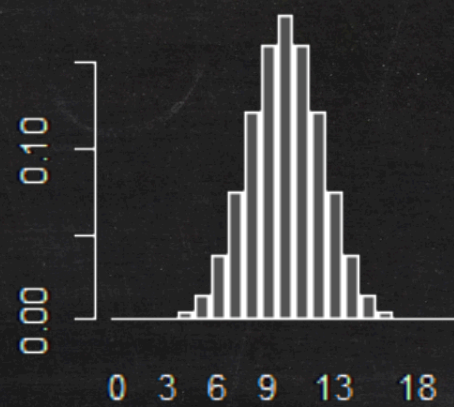
bin. dist. :20:0.1



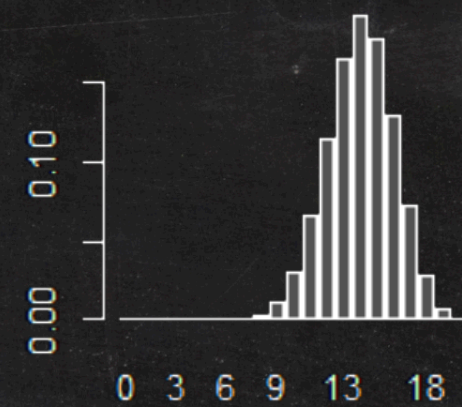
bin. dist. :20:0.3



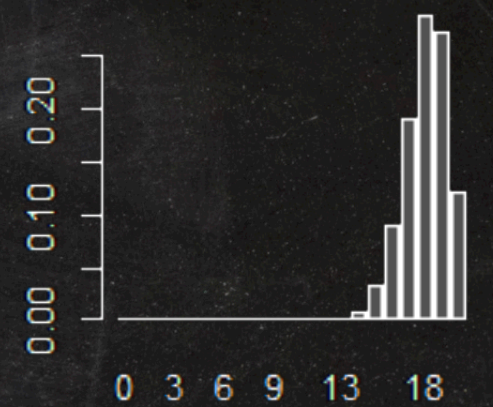
bin. dist. :20:0.5



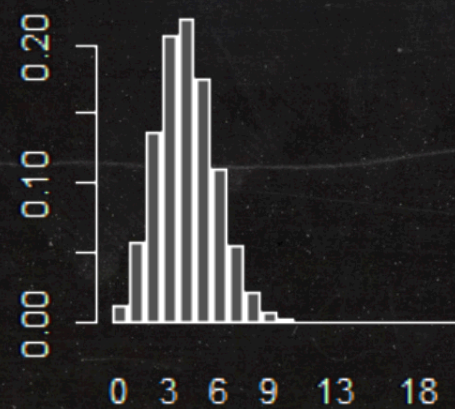
bin. dist. :20:0.7



bin. dist. :20:0.9



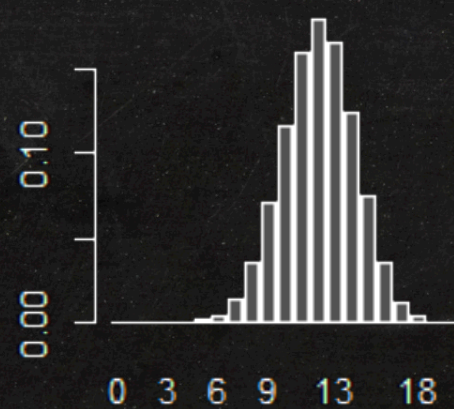
bin. dist. :20:0.2



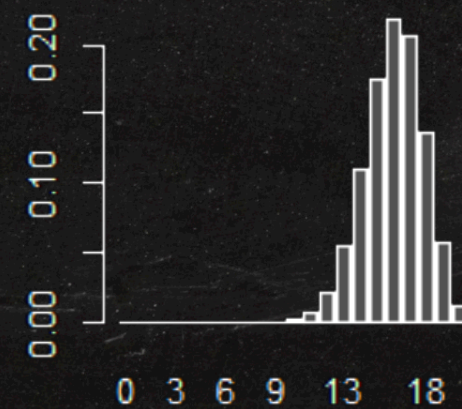
bin. dist. :20:0.4



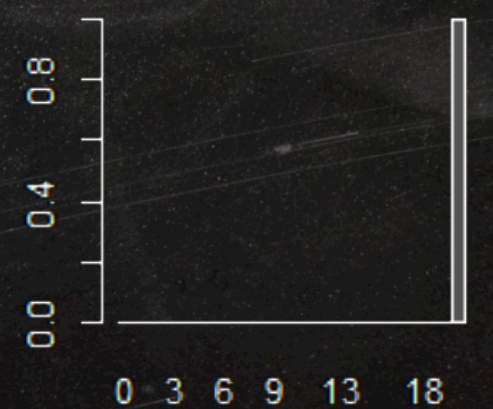
bin. dist. :20:0.6



bin. dist. :20:0.8



bin. dist. :20:1



# PROBABILITY AXIOMS: 3

- ▶ Given two events A and B we define the conditional probability  $P(A | B)$  by:

$$P(A \cap B) = P(A | B) P(B)$$

- ▶ Given two events A and B we say that they are independent iff:

$$P(A \cap B) = P(A) P(B)$$

# PRIOR & POSTERIOR DISTRIBUTIONS

- ▶  $P(\theta)$  the prior probability distribution of  $\theta$
- ▶  $P(\theta|X)$  is the posterior probability of  $\theta$  given  $X$
- ▶ The posterior probability can be written in the memorable form as:  
posterior probability \* likelihood \* prior probability
- ▶ If the posterior distributions  $P(\theta|X)$  are in the same family as the prior probability distribution  $P(\theta)$ , the prior and posterior are then called conjugate distributions, and the prior is called a conjugate prior

# BAYES' THEOREM

From the definition of conditional probability we know:

$$P(A|B)P(B) = P(A \cap B) = P(B|A)P(A)$$

Hence

$$P(B|A) = \frac{P(A|B)P(B)}{P(A)}$$

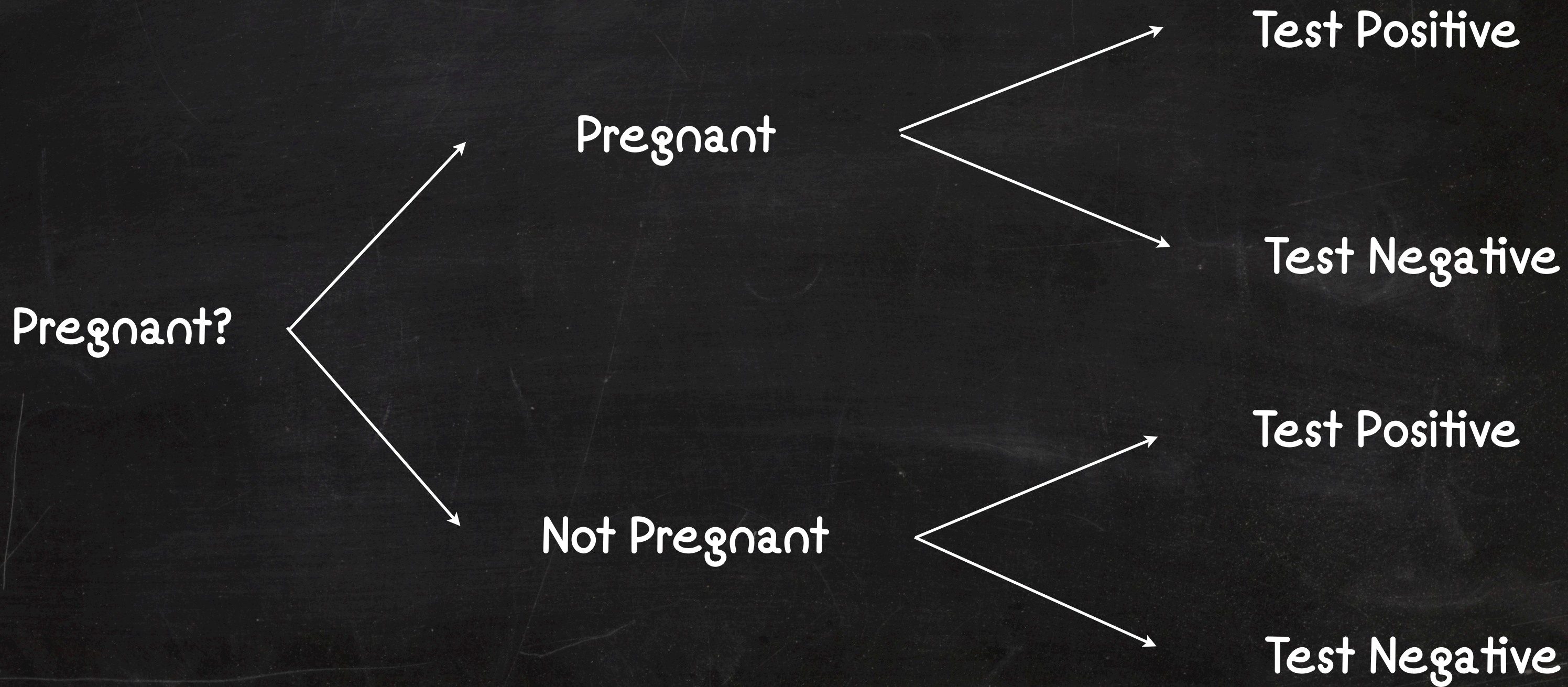
Or if  $B_1, \dots, B_n$  form a partition of the sample space

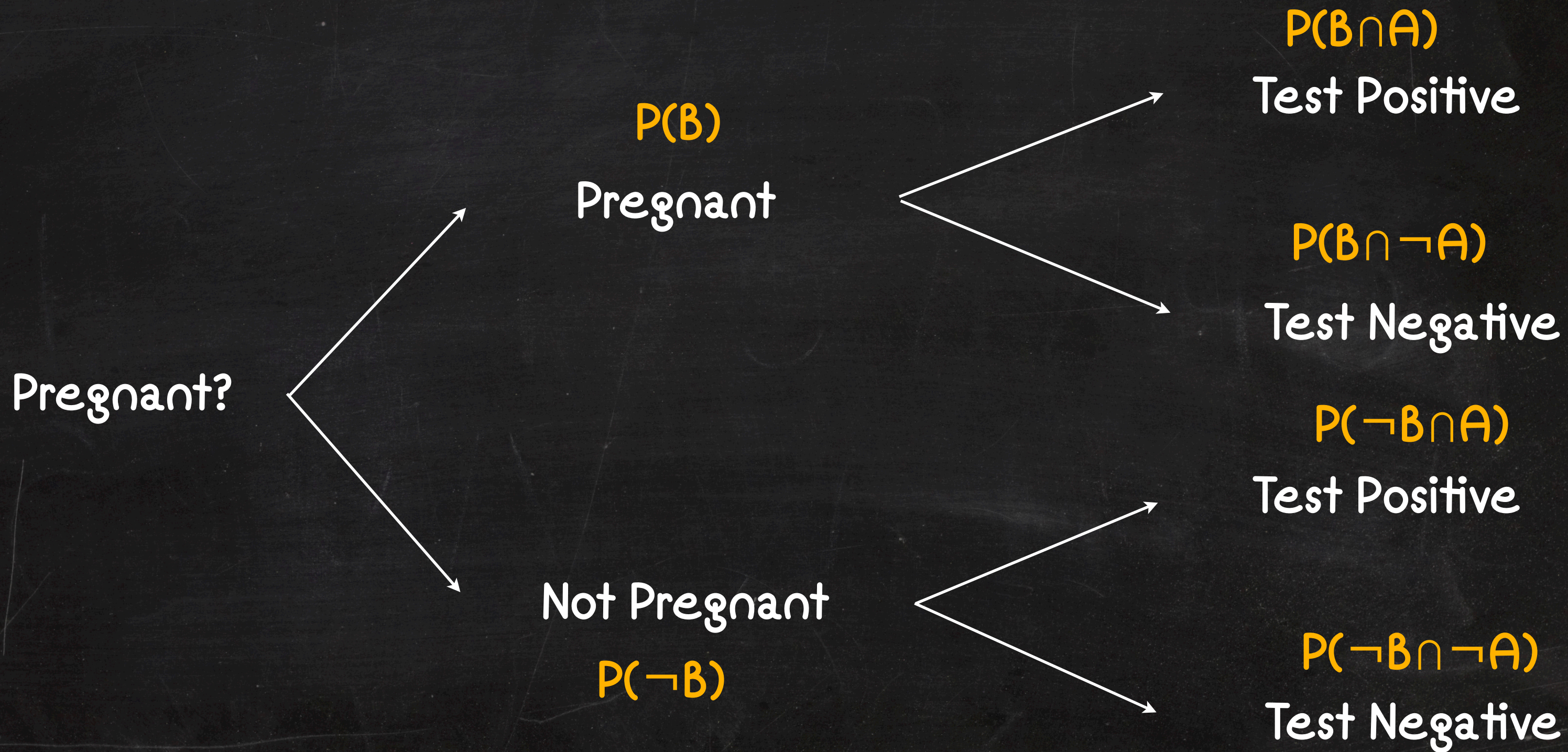
$$P(B_n|A) = \frac{P(A|B_n)P(B_n)}{\sum_i P(A|B_i)P(B_i)}$$

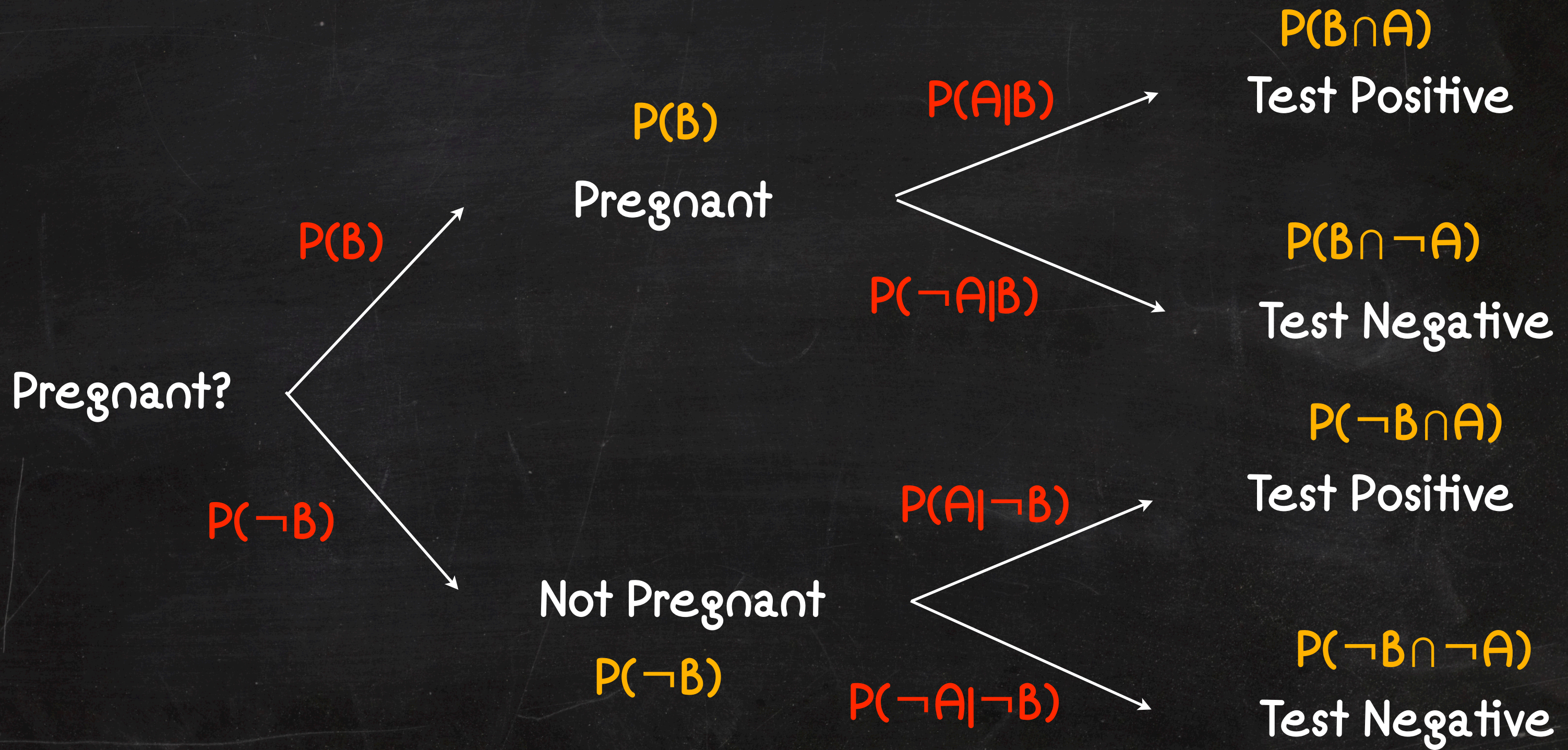


# EXAMPLE: PREGNANCY TESTS

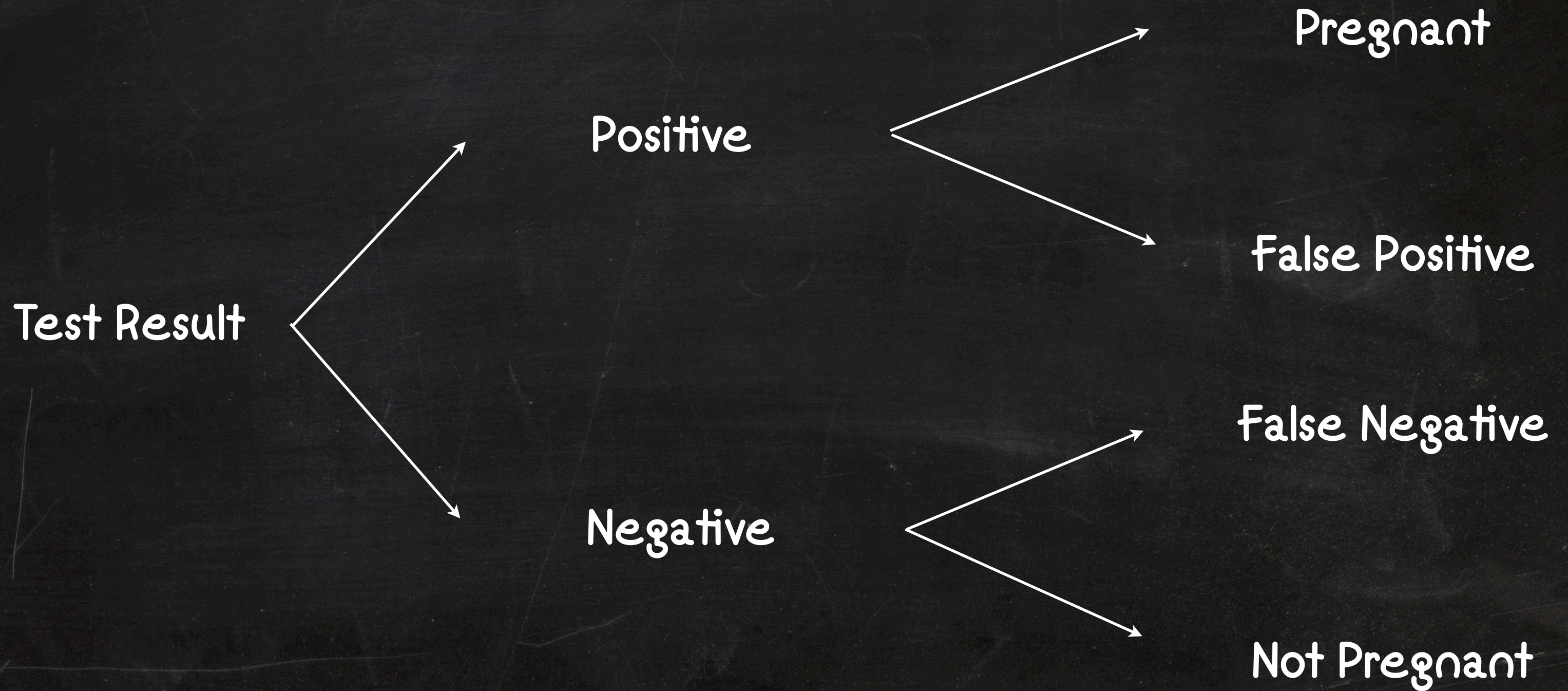
- ▶ Pregnancy tests detect the presence of hCG, or human chorionic gonadotropin, in the blood or urine
- ▶ A “false positive” is when the test incorrectly returns a positive result, and “false negative” when it incorrectly returns a false one.
- ▶ False positives in the hcg test include:
  - ▶ non-pregnant production of the hCG molecule
  - ▶ use of drugs containing the hCG molecule
  - ▶ Some medications cause a positive reaction in the tests
- ▶ The actual probability of being pregnant depends on many messy biological factors

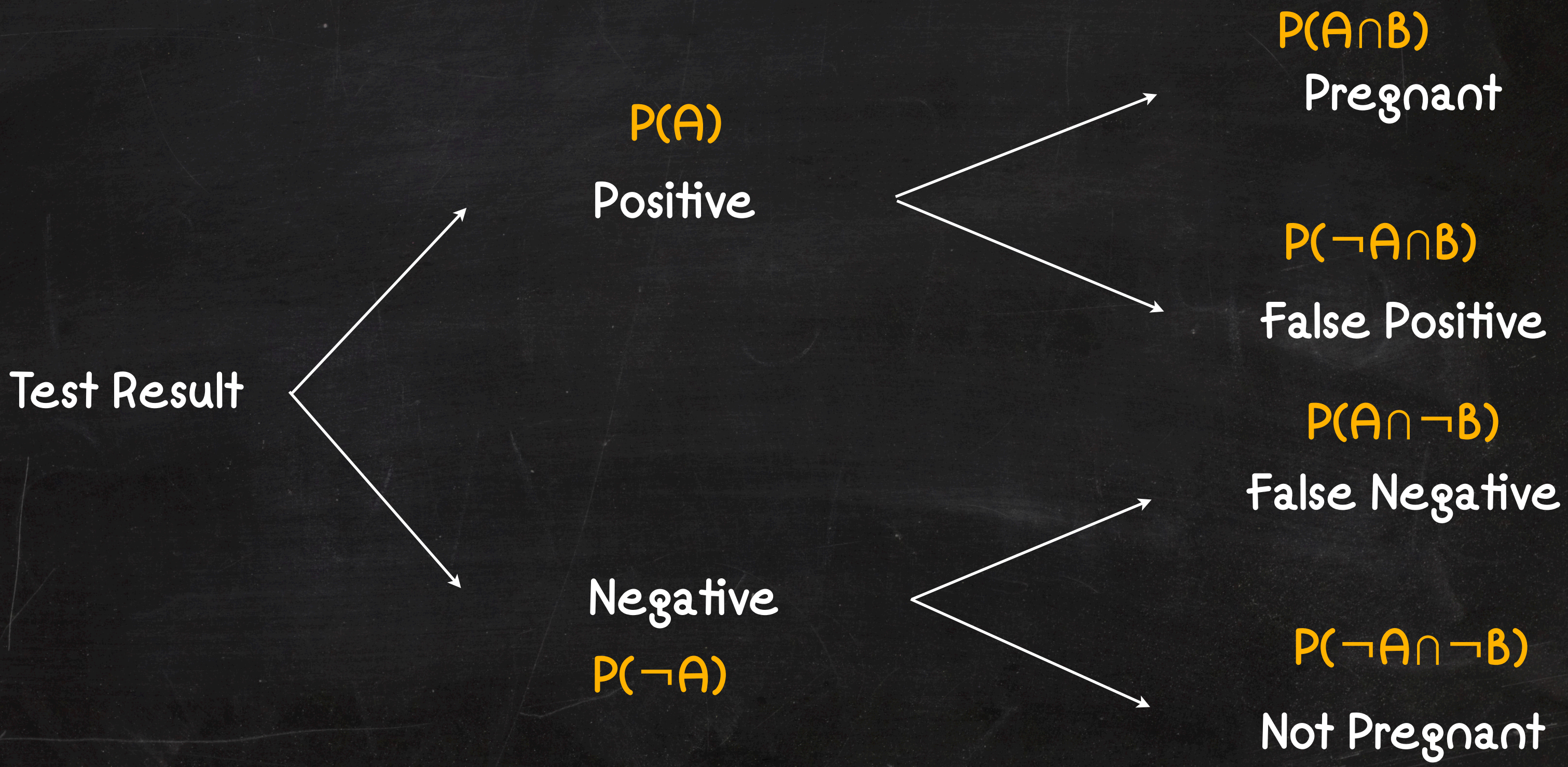


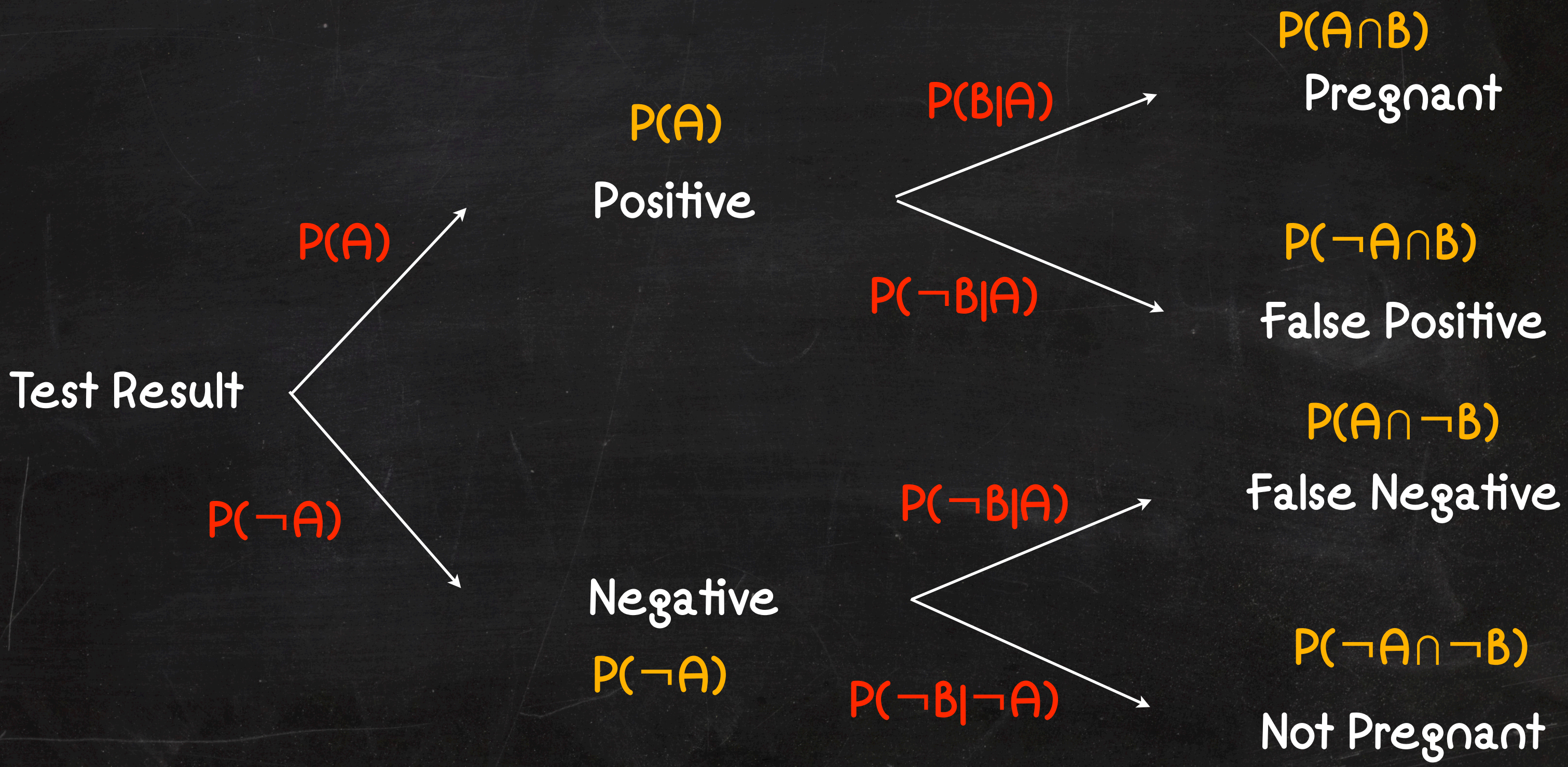




Q: Given the test is positive what is the probability that the subject is pregnant?









# EXAMPLE: DISEASE DIAGNOSIS

Q: Consider the set  $S = \{s_i : i = 0..N\}$  of all disease symptoms, and  $D_{s+} = \{s_i : s_i \text{ in } S\}$  are the diagnostically inclusive symptoms of Ebola, and  $D_{s-} = \{s_i : s_i \text{ in } S\}$  the exclusionary symptoms. Given a patient has some combination of symptoms  $P_s = \{s_i : s_i \text{ in } S\}$ , what is the probability they have Ebola?

- ▶ The presence or absence of some symptoms can completely rule out the diagnosis
- ▶ By updating the model based on real outcomes it is possible to provide more and more accurate predictions

# EXPECTED VALUES & MOMENTS

- ▶ Suppose random variable  $X$  can take value  $x_1$  with probability  $p_1$ , value  $x_2$  with probability  $p_2$ , and so on. Then the expectation of this random variable  $X$  is defined as:

$$E[X] = p_1x_1 + p_2x_2 + \dots + p_kx_k$$

- ▶ The variance of a random variable  $X$  is its second central moment, the expected value of the squared deviation from the mean  $\mu = E[X]$ :

$$\text{Var}(X) = E[(X - \mu)^2]$$

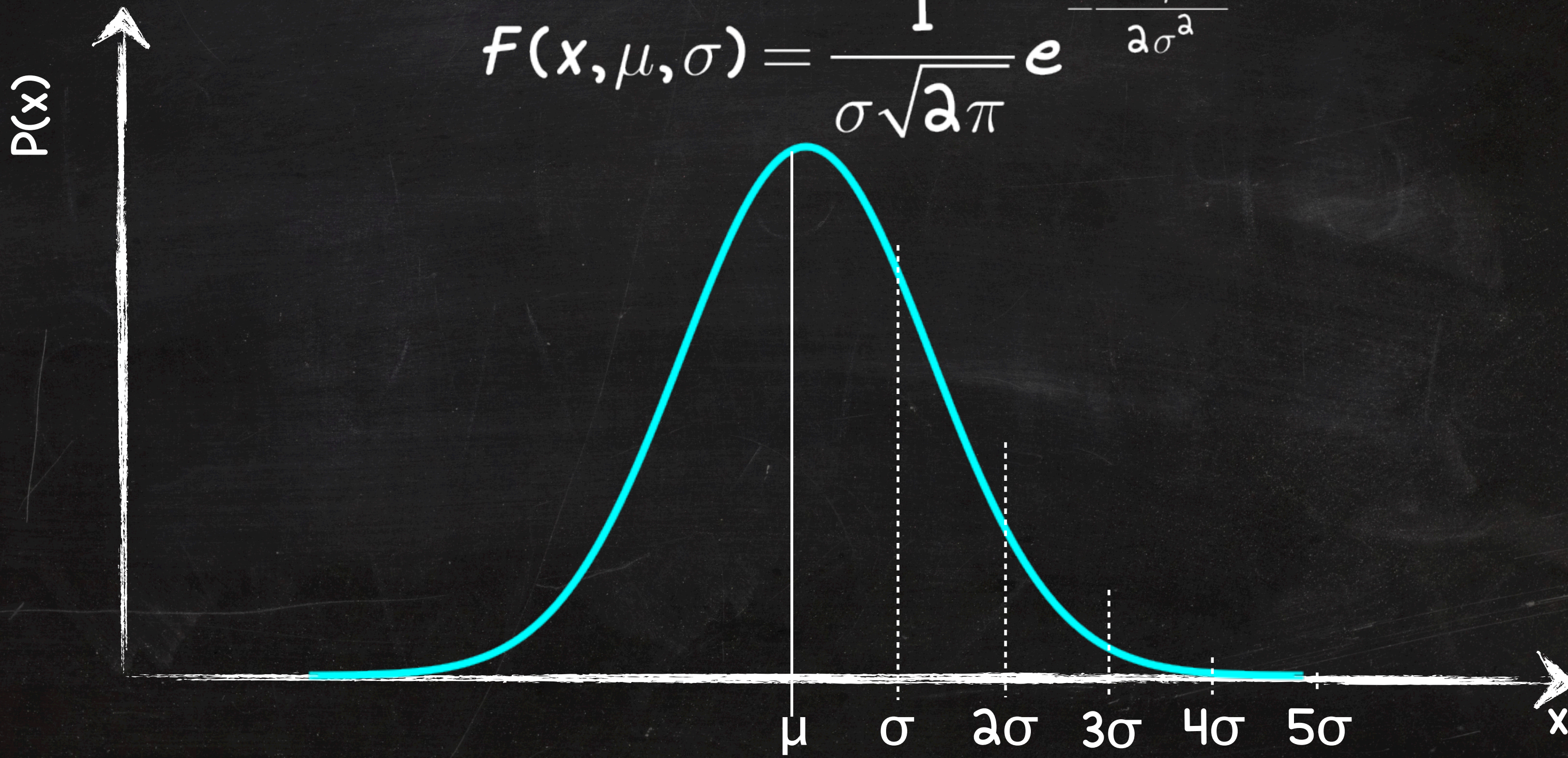
# VARIANCE & COVARIANCE

- ▶ Variance is a measure of how far a set of numbers differs from the mean of those numbers. The square root of the variance is the standard deviation  $\sigma$ 
  - ▶ CERN uses the 5-sigma rule to rule out statistical anomalies in sensor readings, i.e. is the value NOT the expected value of noise
- ▶ The covariance between two jointly distributed random variables X and Y with finite second moments is defined as:

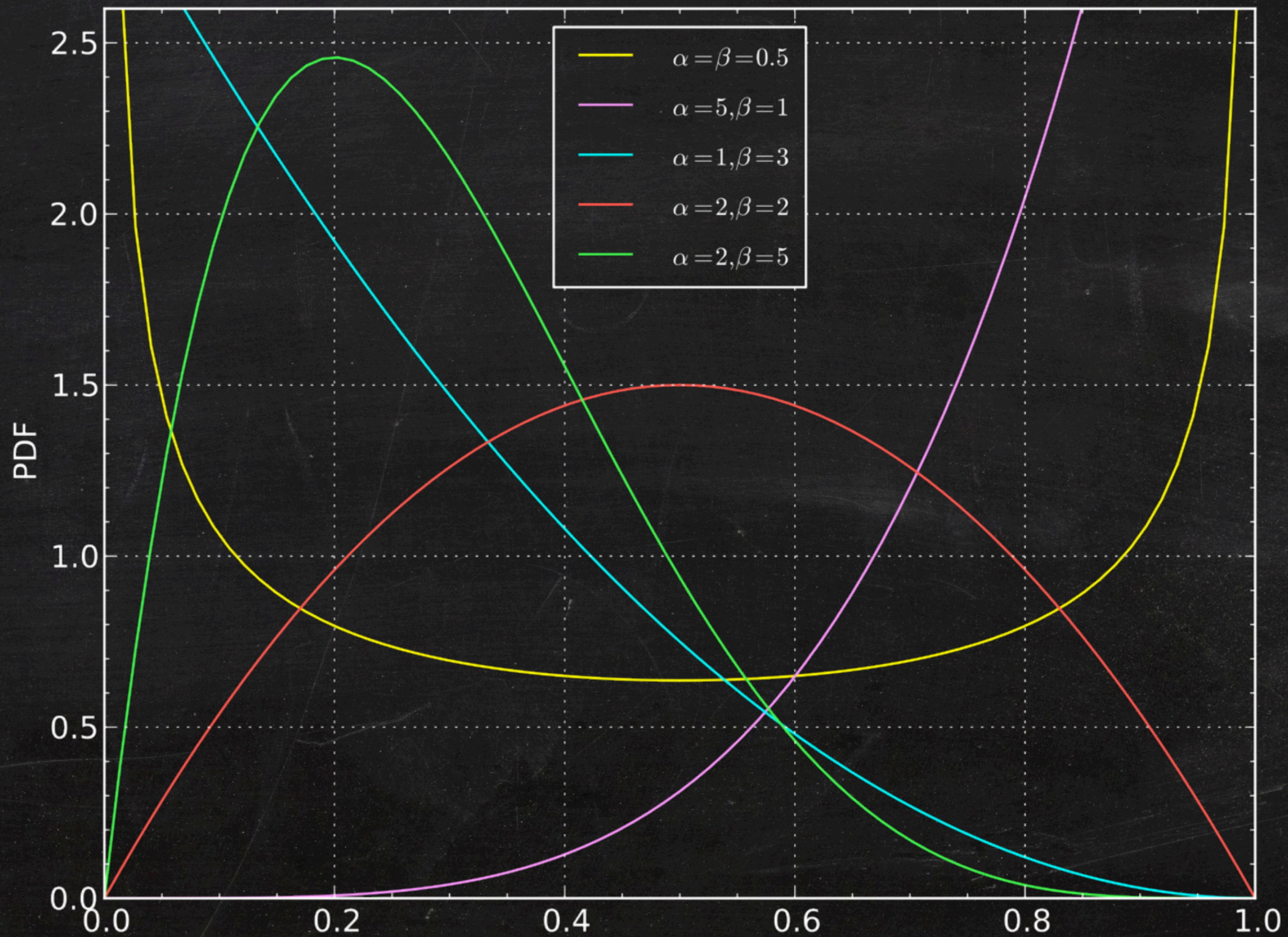
$$\sigma(X, Y) = E[(X - E[X]) \cdot (Y - E[Y])]$$

# GAUSSIAN DISTRIBUTION

$$f(x, \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



# BETA DISTRIBUTION



END OF DETOUR

# EXAMPLE: AMAZON (REVISITED)

- ▶ Users can rate a product by voting 1-5 stars
- ▶ product rating is the mean of the user votes

Q: how can we rank products with different number of votes?



# SIMPLE "BAYESIAN RANKING"

$$\overline{\text{rank}} = \frac{Cm + Rv}{m + v}$$

Assume the vote posterior distribution is a Normal, then the prior is also a Normal\*, with mean

prior mean

prior precision

$$\frac{\tau_0 \mu_0 + \tau \sum_{i=1}^n x_i}{\tau_0 + n\tau}$$

precision of vote distribution

(\*) [http://en.wikipedia.org/wiki/Conjugate\\_prior](http://en.wikipedia.org/wiki/Conjugate_prior)



# EXAMPLE: YOUTUBE

- ▶ Users can rate a clip by voting +/-1
  - ▶ clip rating is the mean of the user votes
  - ▶ clips also record the number of views
- Q: how can we compare clip ratings with different number of votes and views?
- Q: how can we make results more relevant? e.g. take in to account how old ratings are, author provenance, cost of incorrect ranking promotion

# "BAYESIAN RANK"

- ▶ Since this is a +/- Bernoulli Trial we can model the prior belief distribution by a beta function\*

$$f(x, \alpha, \beta) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}$$

- ▶ Let  $\alpha$  = upvote bias + number of up votes
- ▶ Let  $\beta$  = downvote bias + number of down votes + 1
- ▶ Every time we receive a new vote we just recalculate the distribution

- ▶ To map between a belief and a sorting criterion we make a decision using a loss function  $L$
- ▶ Since the value of  $L$  depends on the value of a random variable we use instead the expected value of  $L$
- ▶ Consider a multilinear loss function:

$$L_k(x, X) = \begin{cases} k(X - x) & : x < X \\ x - X & : x \geq X \end{cases}$$

since we want to minimise the loss we have:

$$\begin{aligned} \min(E[L_k(x, X)]) &= I_x(U + 1, D + 1) \\ &= \frac{1}{1 + k} \end{aligned}$$

# EXTENDING THE LOSS FUNCTION

- ▶ Suppose in addition to the vote counts we also record the timestamp of the votes
- ▶ Items becomes less relevant in the rank the longer it is since the last vote following a pattern of exponential decay
- ▶ Hence the current up or down vote count is now determined by
- ▶ Thus we derive an updated rank function from

$$v' = v \times a^{-t/\lambda} + 1$$

$$I_x(U \times a^{-t/\lambda}, D \times a^{-t/\lambda}) = \frac{1}{1+k}$$

# THE HOUSE THAT SKYNET BUILT

- ▶ SkyNet SmartHome™ is a system designed to manage the state of various household resources, e.g. heating, lighting, media-centre, etc
- ▶ it communicates with users smart phones to identify their location
- ▶ its aims to are to maximise the comfort (e.g. room temperature, hot water) of the users and minimise on waste (e.g. power consumption)
- ▶ users can provide feedback to correct inappropriate behaviour, e.g turning the heating up too high

# PREDICTING BEHAVIOUR

Q: Given the user is leaving the office what actions should Skynet smarthome take?

▶ the variables could include:

▶ day of the week

▶ work day or holiday

▶ weather

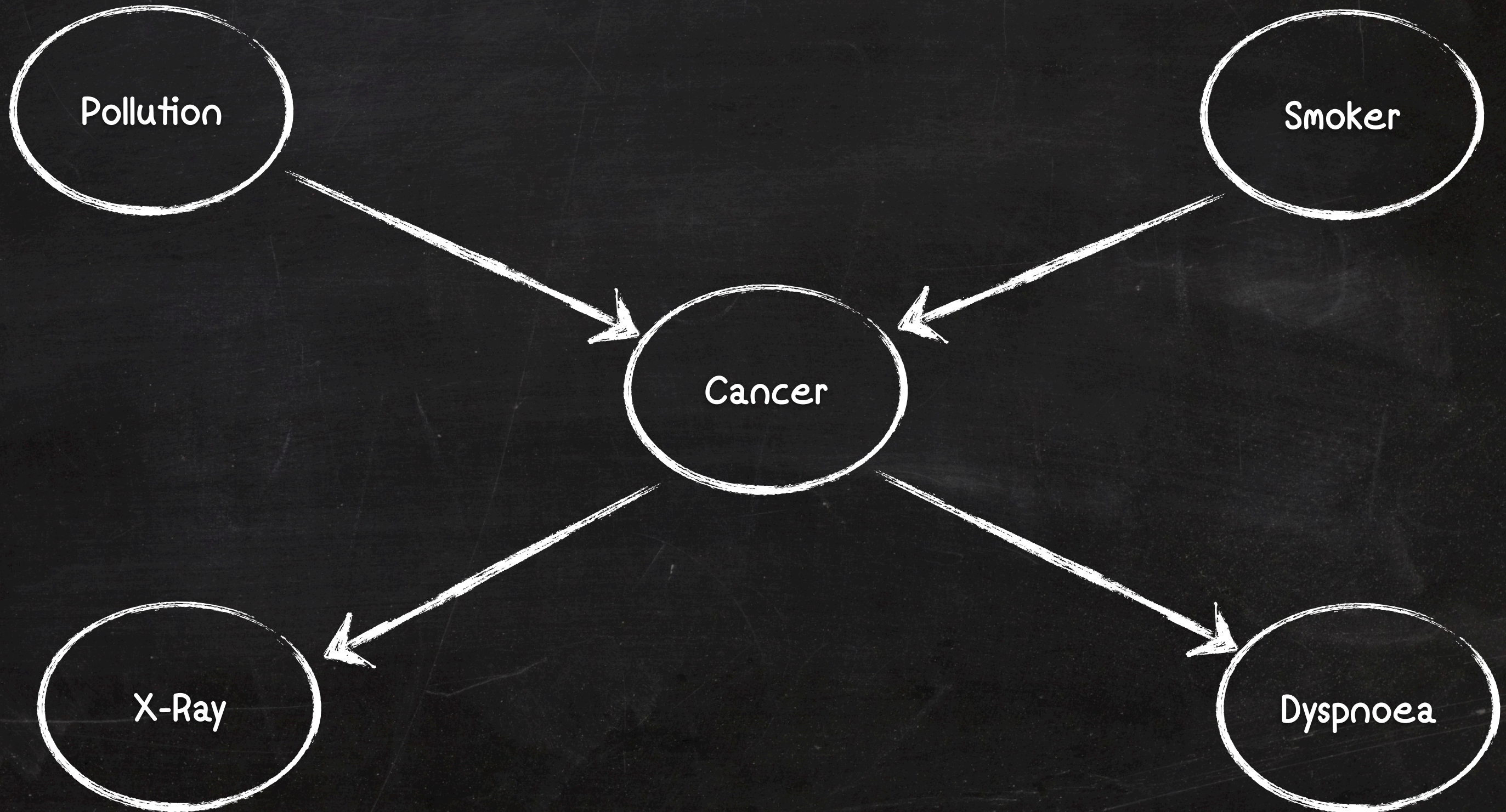
▶ season

▶ calendar events

# BAYESIAN NETWORKS

- ▶ A BN is a probabilistic directed acyclic graphical model that represents a set of random variables and their conditional dependencies
- ▶ Vertices may be observable quantities, latent variables, unknown parameters or hypotheses
- ▶ Vertices that are not connected represent variables that are conditionally independent of each other

# EXAMPLE: CANCER





# EXAMPLE: CANCER

P	P(P)
Low	0.9
High	0.1

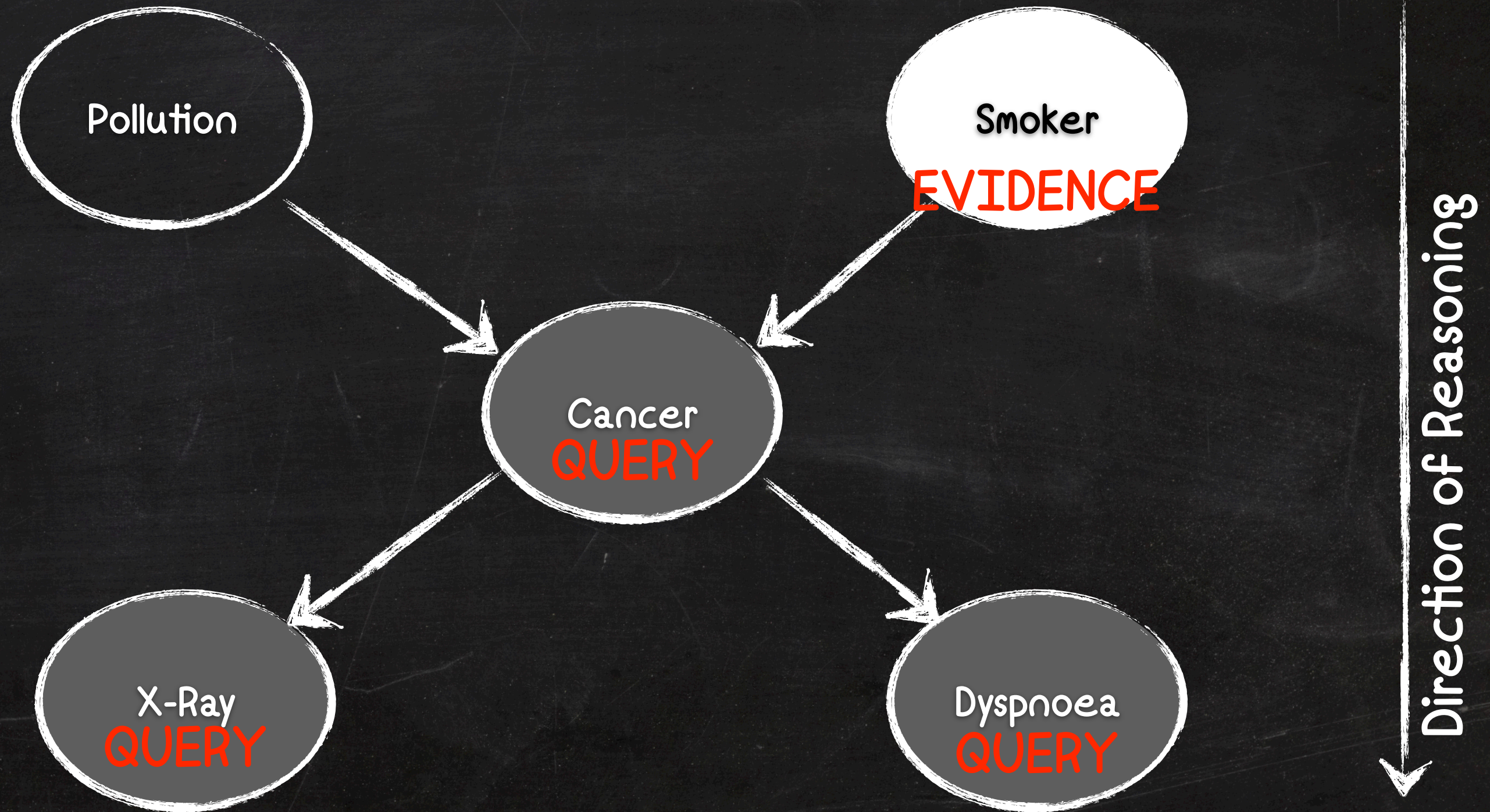
S	P(S)
TRUE	0.3
FALSE	0.7

P	S	P(C P∩S)
Low	T	0.03
Low	F	0.001
High	T	0.05
High	F	0.02

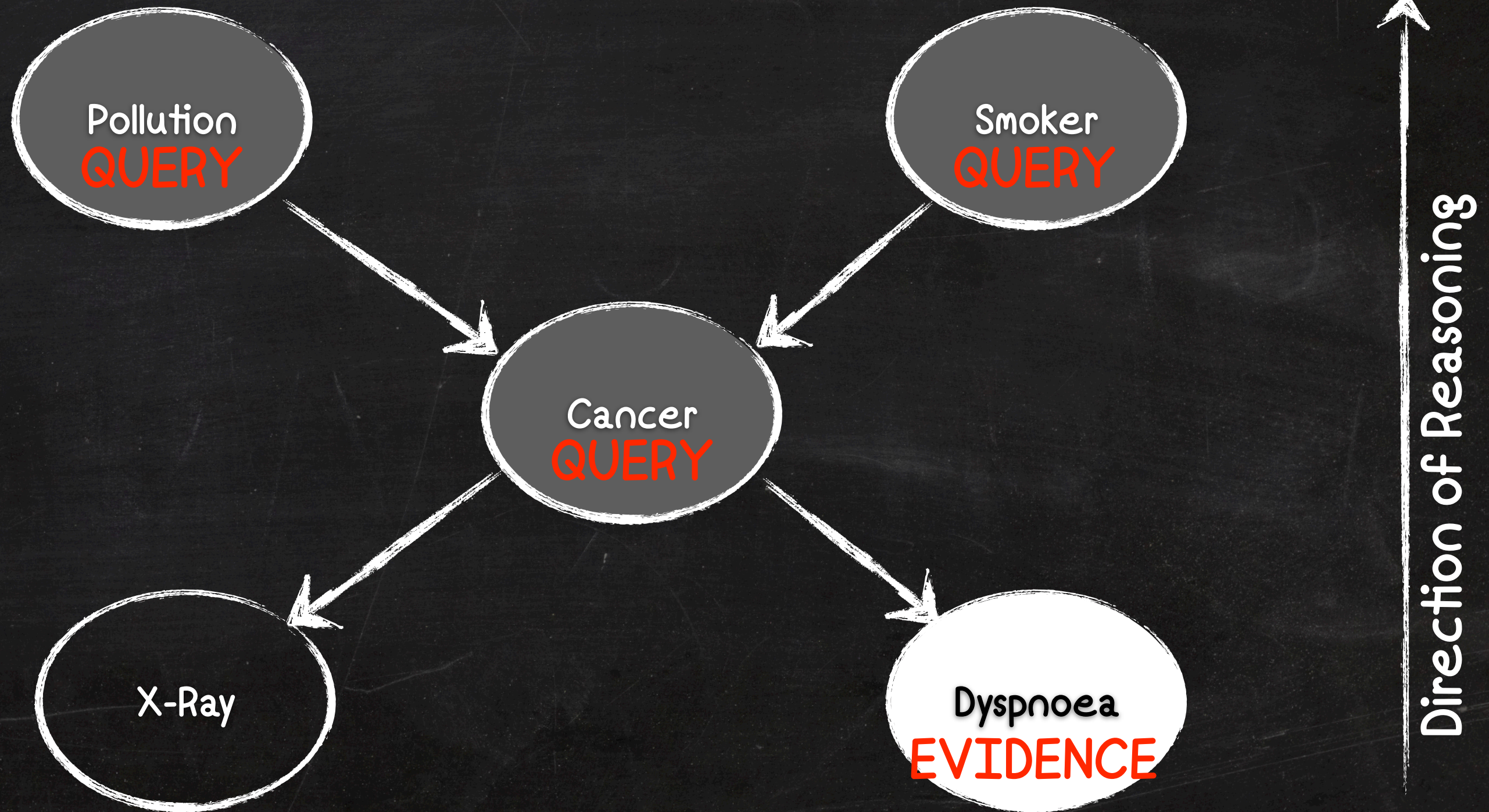
C	P(XRay+ C)
T	0.9
F	0.2

C	P(D+ C)
T	0.65
F	0.3

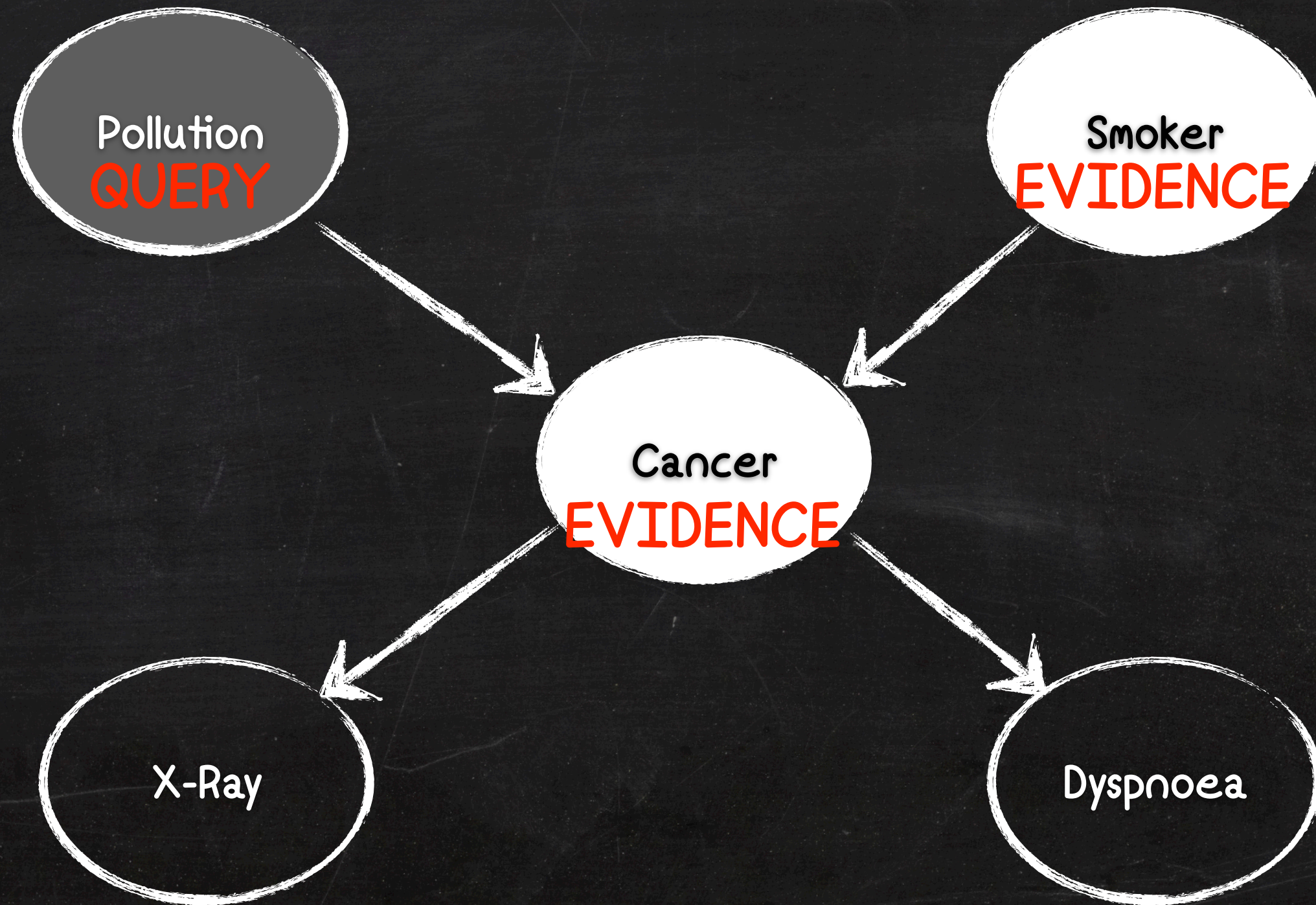
# PREDICTIVE REASONING



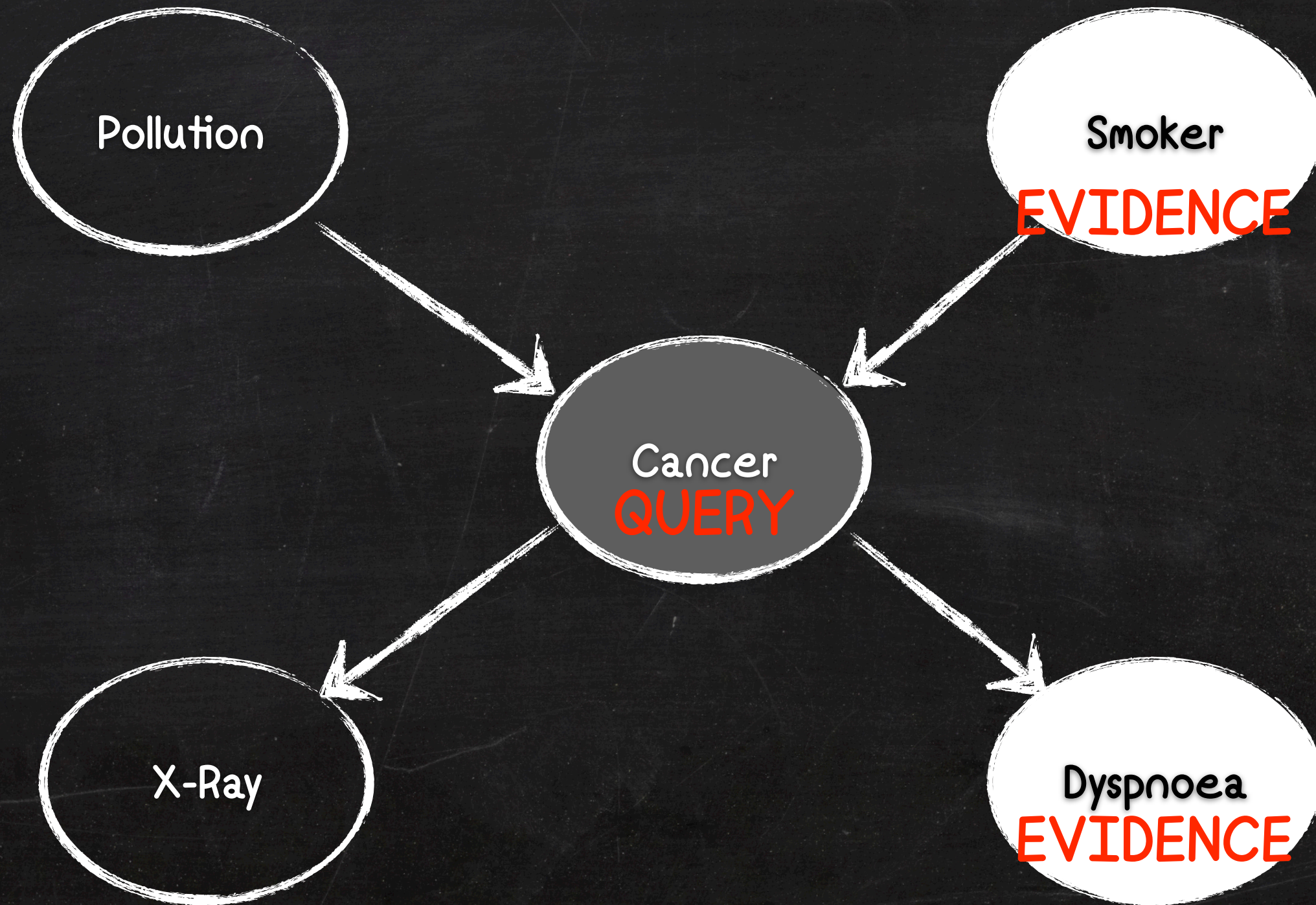
# DIAGNOSTIC REASONING



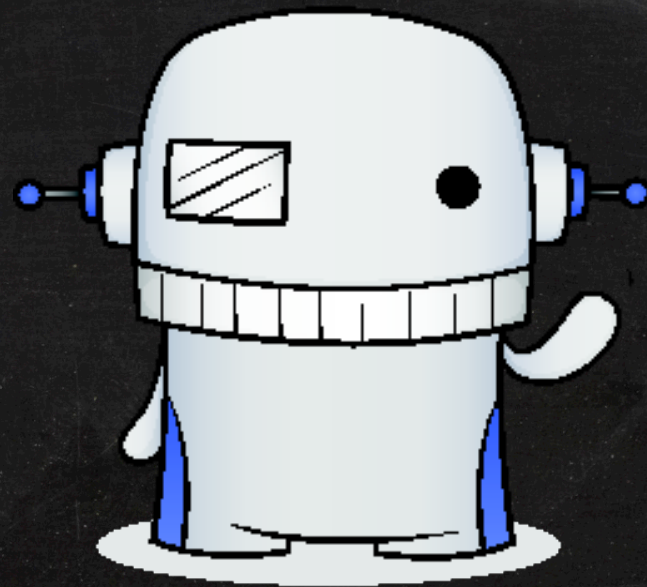
# INTERCAUSAL REASONING



# COMBINED REASONING



# - CHAPTER II - MARKOV MODELS



# STOCHASTIC PROCESS

- ▶ A stochastic process, or random process, is a collection of random variables representing the evolution of some system over time
- ▶ Examples: stock market value and exchange rate fluctuations, audio and video signals, EKG & EEG readings
- ▶ They can be classified as:
  - ▶ Discrete time & discrete space
  - ▶ Discrete time & continuous space
  - ▶ Continuous time & discrete space
  - ▶ Continuous time & continuous space

# A DETOUR INTO MATRIX ALGEBRA



# VECTORS & MATRICES

$$A = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix}$$

$$B = \begin{bmatrix} b_{1,1} & b_{1,2} & b_{1,3} \\ b_{2,1} & b_{2,2} & b_{2,3} \end{bmatrix}$$

1x3 matrix  
or vector

3x2 matrix

# MATRIX ADDITION

If A and B are two m by n matrices then addition is defined by:

$$A + B$$

$$= C, \text{ where } C_{ij} = A_{ij} + B_{ij}$$

$$= \begin{bmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & a_{ij} & \vdots \\ a_{m1} & \cdots & a_{mn} \end{bmatrix} + \begin{bmatrix} b_{11} & \cdots & b_{1n} \\ \vdots & b_{ij} & \vdots \\ b_{m1} & \cdots & b_{mn} \end{bmatrix}$$

$$= \begin{bmatrix} a_{11} + b_{11} & \cdots & a_{1n} + b_{1n} \\ \vdots & a_{ij} + b_{ij} & \vdots \\ a_{m1} + b_{m1} & \cdots & a_{mn} + b_{mn} \end{bmatrix}$$

# MATRIX MULTIPLICATION

If  $A$  is  $n \times m$  matrix and  $B$  is an  $m \times p$  matrix then multiplication is defined by:

$$A \cdot B = AB = \begin{bmatrix} (AB)_{11} & \cdots & (AB)_{1p} \\ \vdots & \ddots & \vdots \\ (AB)_{n1} & \cdots & (AB)_{np} \end{bmatrix}$$

where

$$(AB)_{i,j} = \sum_{k=1}^m A_{i,k} B_{k,j}$$

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \bullet \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix} =$$

$$\begin{bmatrix} a_{11}b_{11} + a_{12}b_{21} & a_{11}b_{12} + a_{12}b_{22} \\ a_{21}b_{11} + a_{22}b_{21} & a_{21}b_{12} + a_{22}b_{22} \end{bmatrix}$$

$$\begin{bmatrix} a_{11} & a_{1a} \\ a_{a1} & a_{aa} \end{bmatrix} \cdot \begin{bmatrix} b_{11} & b_{1a} \\ b_{a1} & b_{aa} \end{bmatrix} =$$

$$\begin{bmatrix} a_{11}b_{11} + a_{1a}b_{a1} & a_{11}b_{1a} + a_{1a}b_{aa} \\ a_{a1}b_{11} + a_{aa}b_{a1} & a_{a1}b_{1a} + a_{aa}b_{aa} \end{bmatrix}$$

$$\begin{bmatrix} a_{11} & a_{1a} \\ a_{a1} & a_{aa} \end{bmatrix} \bullet \begin{bmatrix} b_{11} & b_{1a} \\ b_{a1} & b_{aa} \end{bmatrix} =$$

$$\begin{bmatrix} a_{11}b_{11} + a_{1a}b_{a1} & a_{11}b_{1a} + a_{1a}b_{aa} \\ a_{a1}b_{11} + a_{aa}b_{a1} & a_{a1}b_{1a} + a_{aa}b_{aa} \end{bmatrix}$$

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \cdot \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix} =$$

$$\begin{bmatrix} a_{11}b_{11} + a_{12}b_{21} & a_{11}b_{12} + a_{12}b_{22} \\ a_{21}b_{11} + a_{22}b_{21} & a_{21}b_{12} + a_{22}b_{22} \end{bmatrix}$$

$$\begin{bmatrix} a_{11} & a_{1a} \\ a_{a1} & a_{aa} \end{bmatrix} \cdot \begin{bmatrix} b_{11} & b_{1a} \\ b_{a1} & b_{aa} \end{bmatrix} =$$

$$\begin{bmatrix} a_{11}b_{11} + a_{1a}b_{a1} & a_{11}b_{1a} + a_{1a}b_{aa} \\ a_{a1}b_{11} + a_{aa}b_{a1} & a_{a1}b_{1a} + a_{aa}b_{aa} \end{bmatrix}$$



# THE IDENTITY MATRIX

If  $A$  is  $n \times m$  matrix and then the identity matrix  $I$  is an  $m \times n$  matrix such that  $A I = A$ .  
 $I$  is defined by:

$$A_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

So the  $3 \times 3$  identity matrix is:

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

# INVERSE MATRIX

If  $A$  is  $n \times n$  matrix and then the inverse of  $A$ ,  $A^{-1}$ , is an  $n \times n$  matrix such that  $AA^{-1} = I$ . Non-square matrices do not have inverses

$$\begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$a_{11}b_{11} + a_{12}b_{21} = 1, \quad a_{11}b_{12} + a_{12}b_{22} = 0$$

$$a_{21}b_{11} + a_{22}b_{21} = 0, \quad a_{21}b_{12} + a_{22}b_{22} = 1$$

# MATRIX TRANSPOSITION

If  $A$  is  $n \times m$  matrix and then the transpose of  $A$ ,  $A^T$ , is an  $m \times n$  matrix defined by:

$$A^T_{ij} = A_{ji}$$

For example, consider a  $2 \times 2$  matrix:

$$\begin{pmatrix} a & b \\ d & c \end{pmatrix}^T = \begin{pmatrix} a & d \\ b & c \end{pmatrix}$$

END OF DETOUR

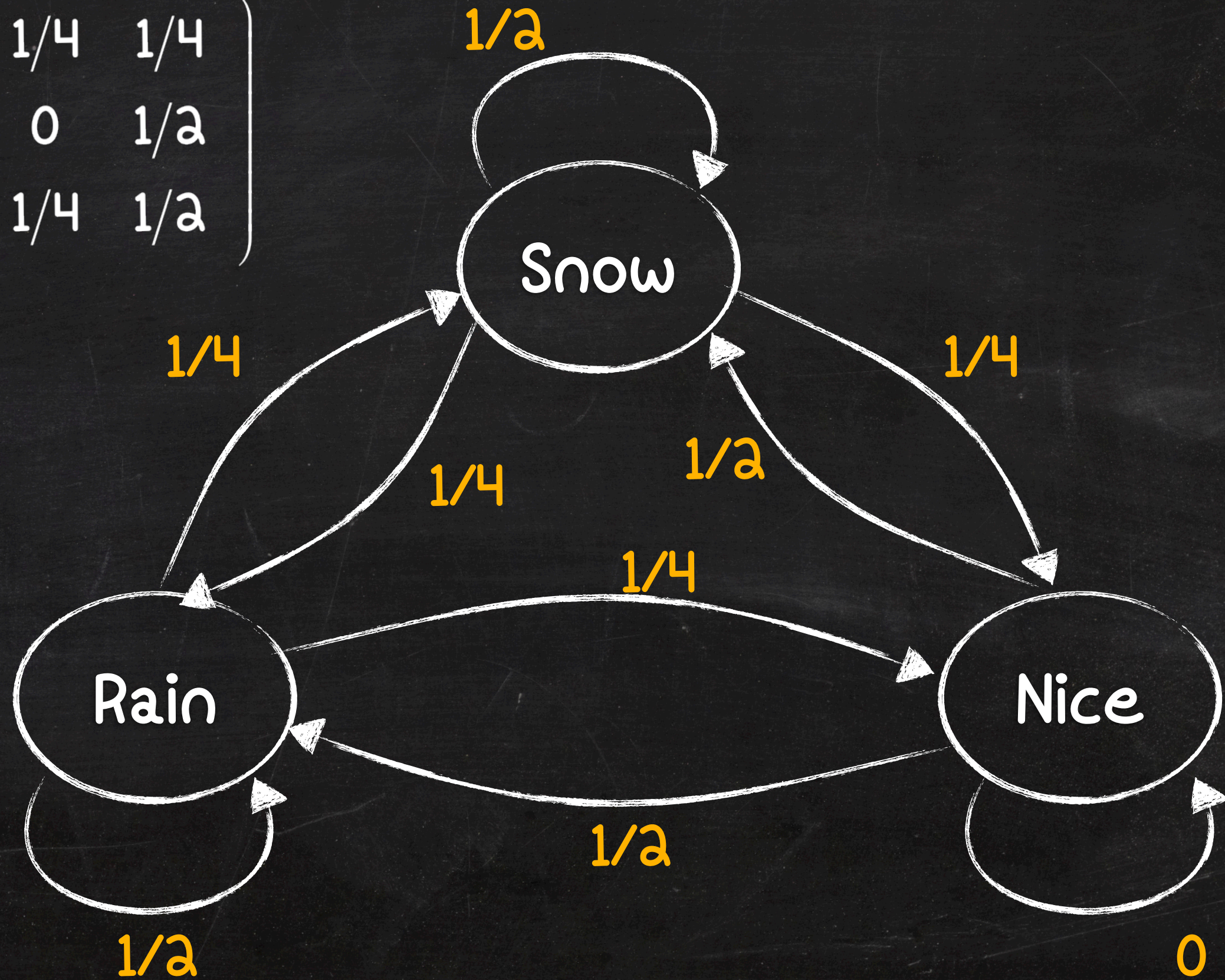
- ▶ A Markov chain is a directed graph whose vertices represent states and the edges the probability of transition between the two states
- ▶ Frequently we use an adjacency matrix representation of the graph  $T$  called the transition matrix
- ▶ Hence for a chain with  $N$  vertices the transition matrix is  $N \times N$
- ▶ The initial state of the system,  $S_0$ , is also an  $N \times N$  matrix
- ▶ The state evolves according to  $S_{n+1} = S_n T = S T^n$
- ▶ The value of  $S_{n+1}(i,j)$  is the probability of being at that state in step  $n+1$

# EXAMPLE: BRITISH WEATHER

- ▶ England is the land of rain. We never have two nice days in a row. In fact if a nice day is always followed by either rain or snow. If there is a change from rain or snow only half the time is this a change to sunny weather

$$T = \begin{array}{l} \text{Rain} \\ \text{Nice} \\ \text{Snow} \end{array} \begin{pmatrix} \text{Rain} & \text{Nice} & \text{Snow} \\ 1/2 & 1/4 & 1/4 \\ 1/2 & 0 & 1/2 \\ 1/4 & 1/4 & 1/2 \end{pmatrix}$$

$$T = \begin{matrix} & \begin{matrix} R & N & S \end{matrix} \\ \begin{matrix} R \\ N \\ S \end{matrix} & \begin{pmatrix} 1/2 & 1/4 & 1/4 \\ 1/2 & 0 & 1/2 \\ 1/4 & 1/4 & 1/2 \end{pmatrix} \end{matrix}$$



- ▶ We can determine the long-term state of the process by calculating  $T^n$  as  $n$  increases towards infinity. The system will converge on a stationary value.
- ▶ For our weather example we find:

	Rain	Nice	Snow
Rain	4/10	2/10	4/10
Nice	4/10	2/10	4/10
Snow	4/10	2/10	4/10



# ABSORBING MARKOV CHAINS

- ▶ A state  $s_i$  of a Markov chain is called absorbing if it is impossible to leave it (i.e.,  $T_{i,i} = 1$ )
- ▶ A Markov chain is absorbing if:
  - ▶ it has at least one absorbing state
  - ▶ every state is connected to at least one absorbing state
- ▶ In an absorbing Markov chain, the probability that the process will be absorbed is 1 (i.e.,  $Q^n \rightarrow 0$  as  $n \rightarrow \infty$ )

# EXAMPLE: THE WANDERING DRUNK

- ▶ Consider a city divided up in some some grid, e.g. square blocks.
- ▶ The drunk can move 1 block per turn, each direction has equal probability
- ▶ If the drunk reaches Home or the Bar they will stay there
- ▶ Questions we can answer:
  - ▶ What is the expect time until an absorbing state is reached?
  - ▶ How many times does the drunk visit each intersection?

EXAMPLE: PREDICTIVE TEXTING

# HIDDEN MARKOV MODELS

- ▶ The name is misleading! Nothing is unknown...
- ▶ In a basic Markov model the states of the system are visible. E.g. we can see if it is snowing
- ▶ In a hidden Markov Model the (entire) state is not directly visible but some outputs dependent on the state are observable
- ▶ However we still need to know all the transition probability values!

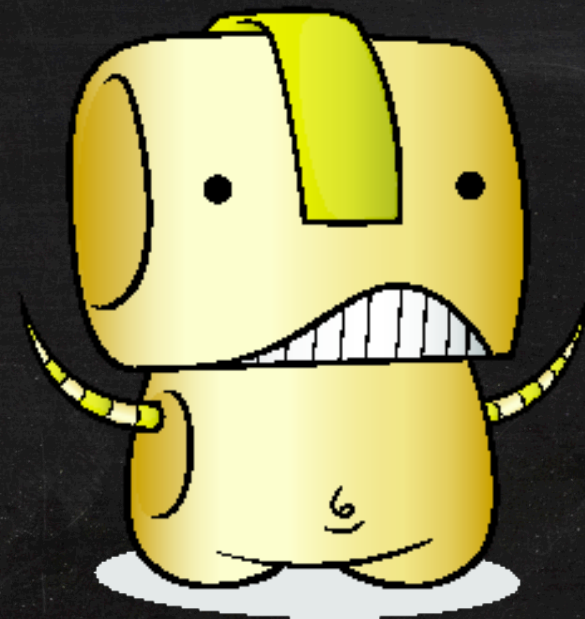
# MARKOV DECISION PROCESSES

- ▶ Markov decision processes are an extension of Markov chains; the difference is the addition of actions (allowing choice) and rewards (giving motivation)
- ▶ A Markov decision process is a 5-tuple  $(S, A, P_A, R_A, L)$
- ▶ The core problem is to choose an action  $\pi$  that will maximise some cumulative function, e.g.

$$\sum_{t=0}^{\infty} \gamma^t R_{A_t}(s_t, s_{t+1}), \text{ where } A_t = \pi(s)$$

- ▶ MDPs can be easily solved by linear (e.g. Simplex method) or dynamic programming (e.g. Map-Reduce)

# - CHAPTER III - KALMAN FILTERS



# A LINEAR DYNAMIC SYSTEM

- ▶ Continuous time definition:

$$\frac{\partial}{\partial t} x(t) = A \cdot x(t)$$

- ▶ Discrete time definition:

$$x_{n+1} = A \cdot x_n$$

- ▶ The systems are called linear since given any two solutions  $x(t)$  &  $y(t)$  then any linear combination of these solutions is also a solution

$$z(t) = \alpha x(t) + \beta y(t)$$

# EXAMPLE: CLIMATE CONTROL

- ▶ SkyNet SmartHome™ is able to monitor the temperature of rooms in the house and effect heating/AC to regulate the temperature
- ▶ The temperature sensors contain noise
- ▶ heating and AC are either ON or OFF
- ▶ similar systems exist for humidity



CONSIDER...

$$X_n = AX_{n-1} + Bu_k + \omega_{k-1}$$

# KALMAN FILTERS

- ▶ Developed ~1960 by Rudolf E. Kálmán, Peter Swerling, and Richard S. Bucy.
- ▶ First implemented in NASA as part of the Apollo navigation computer
- ▶ Still used in many aeronautic and military applications, e.g. submarines, cruise missiles, NASA Space Shuttle, ISS
- ▶ There are generalisation of the basic Kalman filters for continuous time systems as well as non-linear systems

- ▶ Kalman Filters work by making a prediction of the future, getting a measurement from reality, comparing the two, moderating this difference, and adjusting its estimate with this moderated value.
- ▶ Kalman filters are:
  - ▶ discrete
  - ▶ recursive
  - ▶ Extremely accurate if you have a good model

# OVERVIEW

current  
estimate

measured  
value

$$\hat{X}_k = K_k \cdot Z_k + (1 - K_k) \cdot \hat{X}_{k-1}$$

Kalman gain

previous  
estimate

# BASIC KALMAN FILTERS

- ▶ Modelled on a Markov chain built on linear operators perturbed by errors that may include Gaussian noise
- ▶ The state of the system is represented by a vector of real numbers
- ▶ The filter is recursive, i.e. only the estimated state from the previous time step and the current measurement are needed to compute the estimate for the current state
- ▶ typically we describe the algorithm in two phases: predict & update

# 1. PREDICT: STATE ESTIMATION

State transition  
matrix

Control matrix

$$\hat{X}_n = A \cdot \hat{X}_{n-1} + B \cdot U_{n-1}$$

Predicted state

Previous estimate of  
state

Control vector

# 2. PREDICT: ERROR ESTIMATION

State transition  
matrix

$$\hat{P}_n = A \cdot \hat{P}_{n-1} \cdot A^T + Q$$

Covariance prediction

Previous  
covariance  
estimate

Estimated process  
error covariance

### 3: UPDATE: INNOVATION COVARIANCE

Measurement Vector

Observation Matrix

$$\tilde{y} = z_n - H \cdot \hat{x}_n$$

Covariance Innovation

Predicted state (step 1)

The diagram shows the equation  $\tilde{y} = z_n - H \cdot \hat{x}_n$  on a chalkboard background. The term  $\tilde{y}$  is labeled 'Covariance Innovation' with an arrow pointing to it. The term  $z_n$  is labeled 'Measurement Vector' with an arrow pointing to it. The term  $H$  is labeled 'Observation Matrix' with an arrow pointing to it. The term  $\hat{x}_n$  is labeled 'Predicted state (step 1)' with an arrow pointing to it. The minus sign and the dot are not labeled.



# 4: UPDATE: INNOVATION COVARIANCE

Observation Matrix

$$S = H \cdot \hat{P}_n \cdot H^T + R$$

Covariance  
Innovation

Covariance prediction  
(step 2)

Estimated measurement error  
covariance

# 5: UPDATE: KALMAN GAIN

Covariance prediction  
(step 2)

Observation Matrix

$$K = \hat{P}_{\hat{n}} \cdot H^T \cdot S^{-1}$$

Kalman Gain

Covariance Innovation  
(step 3)

# 6. UPDATE STATE

Predicted state (step 1)

Covariance Innovation  
(step 3)

$$\hat{X}_n = \hat{X}_{\hat{n}} + K \tilde{y}$$

New state estimate

Kalman Gain (step 5)

# 7: UPDATE: COVARIANCE

Kalman Gain (step 5)

$$\hat{P}_n = (I - K \cdot H) \hat{P}_{\hat{n}}$$

New estimate of error

Observation Matrix

Covariance prediction (step 2)

# EXAMPLE: VOLTMETER

- ▶ Consider a voltmeter measuring a constant DC voltage via a sensor with noise.
- ▶ The system can be described by:

$$V_n = V_{n-1} + \omega_n$$

- ▶ Since the voltage is constant using a Kalman filter allows us to filter out the noise  $\omega_n$
- ▶ Also since this is a single state example all matrices are of size 1\*1

- ▶ A: State transition - since the previous state should equal the current state  
 $A=1$
- ▶ H: Observation transform - since we're taking direct measurements from the sensor  $H=1$
- ▶ B: Control matrix - we have no controls so  $B=0$
- ▶ Q: Process covariance - since we know the model very accurately  $Q=0.00001$
- ▶ R: Measurement covariance - we don't trust the sensor too much so  $R=0.1$
- ▶  $\hat{X}$ : Initial state estimate = any number
- ▶  $\hat{P}$ : Initial covariance estimate = 1 (because)

$$\hat{X}_{\hat{n}} = \hat{X}_{n-1}$$

$$\hat{P}_{\hat{n}} = \hat{P}_{n-1} + 0.000001$$

$$\tilde{y} = z_n - \hat{X}_{\hat{n}}$$

$$S = \hat{P}_{\hat{n}} + 0.1$$

$$K = \hat{P}_{\hat{n}} \cdot S^{-1}$$

$$\hat{X}_n = \hat{X}_{\hat{n}} + K\tilde{y}$$

$$\hat{P}_n = (I - K)\hat{P}_{\hat{n}}$$

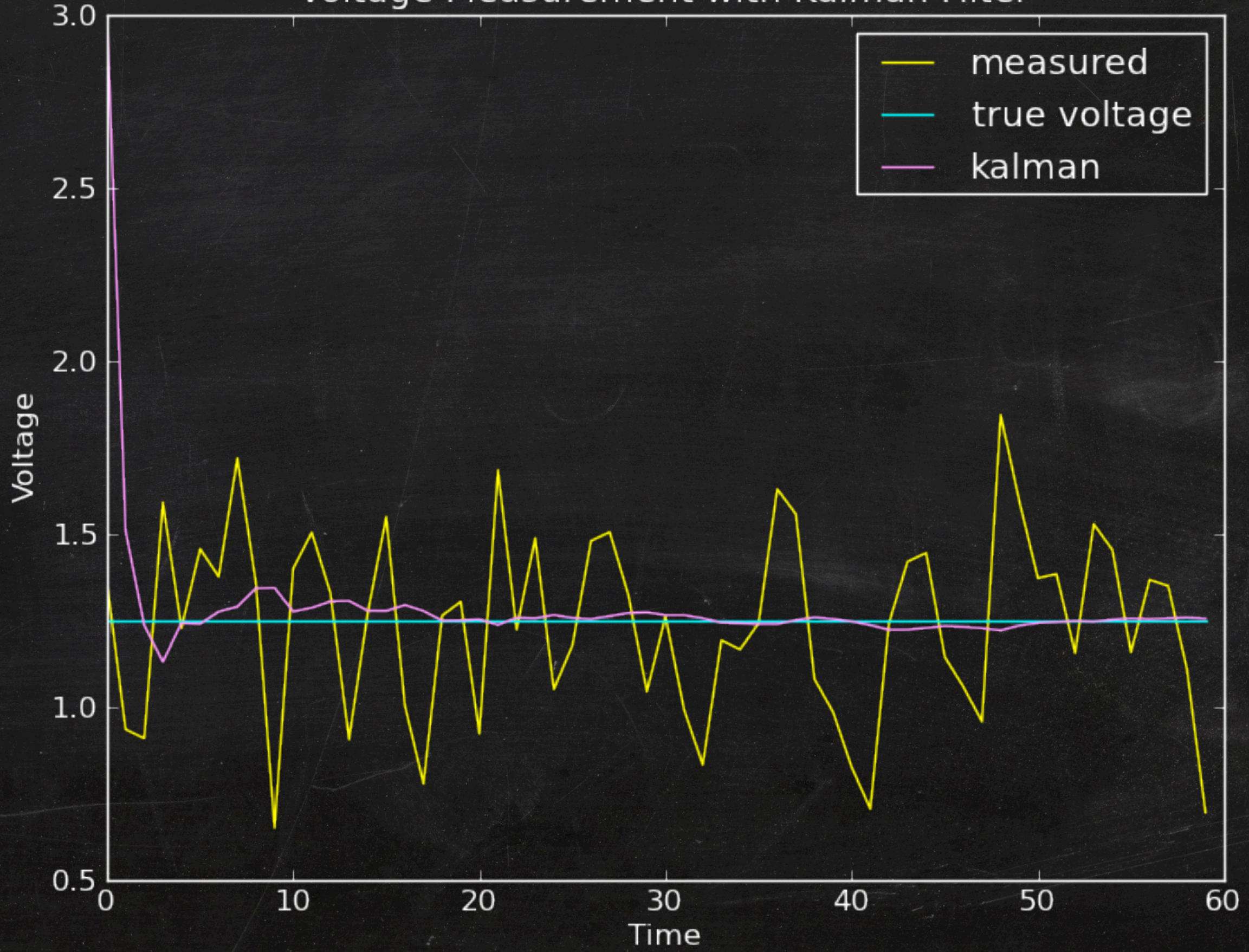
$$K = \frac{\hat{P}_{n-1} + 0.00001}{(\hat{P}_{n-1} + 0.00001) + 0.1}$$

$$\hat{X}_n = Kz_n - (K - 1)\hat{X}_{n-1}$$

$$\hat{P}_n = (1 - K)(\hat{P}_{n-1} + 0.00001)$$



# Voltage Measurement with Kalman Filter



AS A PROGRAMMER YOUR  
CHALLENGE IS TO FIND THE RIGHT  
FILTER MODEL AND DETERMINE  
THE VALUES OF THE MATRICES

# EXAMPLE: ROBO-COPTER

FREEZIN  
ESKIMO



XCELL TEMPEST  
HELICOPTER



# RL HELICOPTER

- ▶ [http://library.rl-community.org/wiki/Helicopter\\_\(Java\)](http://library.rl-community.org/wiki/Helicopter_(Java))
- ▶ Sensors to determine:
  - ▶ bearing
  - ▶ acceleration (velocity)
  - ▶ position (GPS)
  - ▶ rotational rates
  - ▶ inertial measurement unit
  - ▶ and more...

SUMMARY



# EVENTS & RANDOM VARIABLES

EVENTS & RANDOM VARIABLES

CONDITIONAL  
PROBABILITY



EVENTS & RANDOM VARIABLES  
BAYESIAN  
NETWORKS

CONDITIONAL  
PROBABILITY

EVENTS & RANDOM VARIABLES  
MARKOV CHAINS  
BAYESIAN NETWORKS

CONDITIONAL  
PROBABILITY

EVENTS & RANDOM VARIABLES  
MARKOV CHAINS  
BAYESIAN NETWORKS  
RANDOM WALKS  
CONDITIONAL PROBABILITY

EVENTS & RANDOM VARIABLES  
MARKOV CHAINS  
BAYESIAN NETWORKS  
RANDOM

MARKOV DECISION PROCESSES  
CONDITIONAL PROBABILITY

EVENTS & RANDOM VARIABLES  
MARKOV CHAINS    KALMAN FILTERS    BAYESIAN NETWORKS  
RANDOM WALKS    RANDOM PROCESSES  
MARKOV DECISION PROCESSES    CONDITIONAL PROBABILITY

EVENTS & RANDOM VARIABLES  
MARKOV CHAINS    KALMAN FILTERS    BAYESIAN NETWORKS  
RANDOM PROCESSES  
MARKOV DECISION PROCESSES    CONDITIONAL PROBABILITY

# BOOKS

- ▶ Introduction to Probability - Grinstead & Snell  
[http://www.dartmouth.edu/~chance/teaching\\_aids/books\\_articles/probability\\_book/book.html](http://www.dartmouth.edu/~chance/teaching_aids/books_articles/probability_book/book.html)
- ▶ Bayesian Artificial Intelligence - Kevin B. Korb & Ann E. Nicholson
- ▶ An Introduction to Stochastic Modelling - Mark A Pinsky & Samuel Karlin
- ▶ Stochastic Processes and Filtering Theory - Andrew H. Jazwinski
- ▶ Artificial Intelligence: A Modern Approach - Stuart Russell and Peter Norvig

# JAVA LIBRARIES

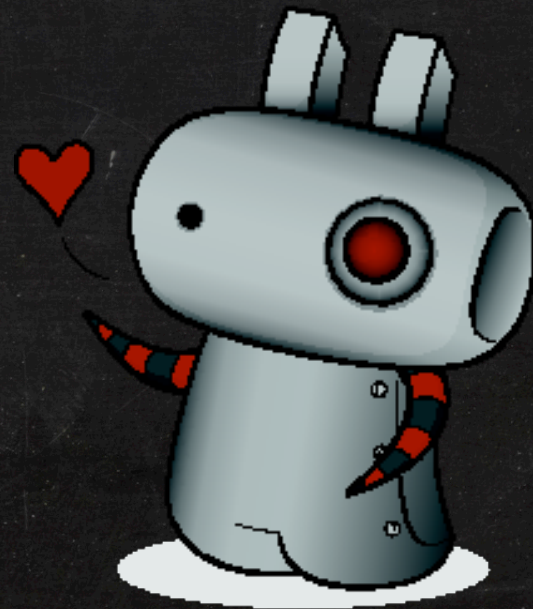
- ▶ Apache Commons Math: <http://commons.apache.org/proper/commons-math/>
- ▶ Colt - high performance data structures and algorithms: <http://dst.lbl.gov/ACSSoftware/colt/>
- ▶ Parallel Colt: <https://sites.google.com/site/piotrwendykier/software/parallelcolt>
- ▶ JBlas - high performance Java API for native libraries LAPACK, BLAS, & ATLAS: <http://mikiobraun.github.io/jblas/>
- ▶ The rest... <http://code.google.com/p/java-matrix-benchmark/>
- ▶ Jayes - A Java framework for Bayesian Networks



# OTHER RESOURCES

- ▶ <http://www.probabilitycourse.com>
- ▶ <http://masanjin.net/blog/bayesian-average> - detailed derivation of bayesian averaging via normal distributions
- ▶ [http://fulmicoton.com/posts/bayesian\\_rating/](http://fulmicoton.com/posts/bayesian_rating/) - an alternative derivation of bayesian "averaging"
- ▶ <http://www.tina-vision.net/docs/memos/1996-002.pdf> - a beautifully simple derivation of Kalman filters
- ▶ <http://www.intechopen.com/books/kalman-filter> - articles on applications of Kalman filters

THANK YOU



JAMES MCGIVERN

PROBABLY, DEFINITELY,  
MAYBE

