

# **Machine Learning for Signal Processing**

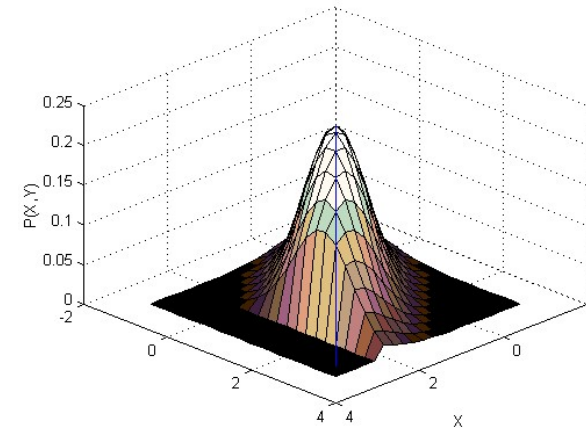
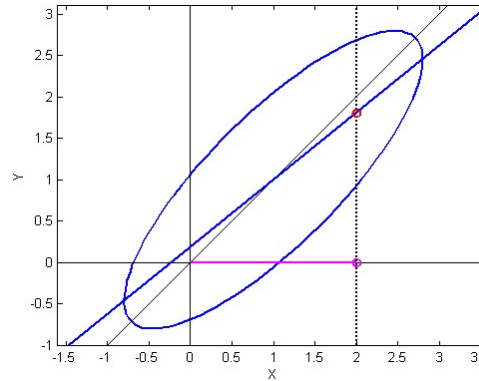
## **Predicting and Estimation from Time Series**

Bhiksha Raj

# Preliminaries : $P(y|x)$ for Gaussian

- If  $P(x,y)$  is Gaussian:

$$P(\mathbf{x}, \mathbf{y}) = N\left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} C_{xx} & C_{xy} \\ C_{yx} & C_{yy} \end{bmatrix}\right)$$



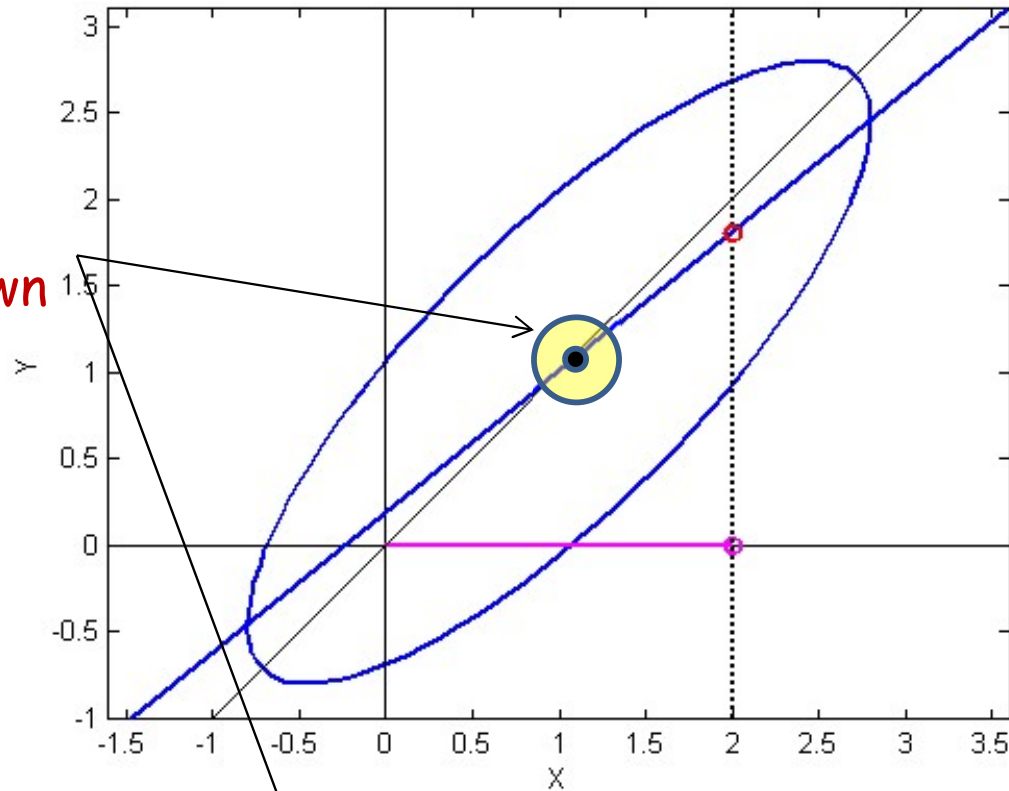
- The conditional probability of  $y$  given  $x$  is also Gaussian
  - The slice in the figure is Gaussian

$$P(y|x) = N(\mu_y + C_{yx}C_{xx}^{-1}(x - \mu_x), C_{yy} - C_{yx}C_{xx}^{-1}C_{xy})$$

- The mean of this Gaussian is a function of  $x$
- The variance of  $y$  reduces if  $x$  is known
  - Uncertainty is reduced

# Preliminaries : $P(y|x)$ for Gaussian

Best guess for  $Y$   
when  $X$  is not known



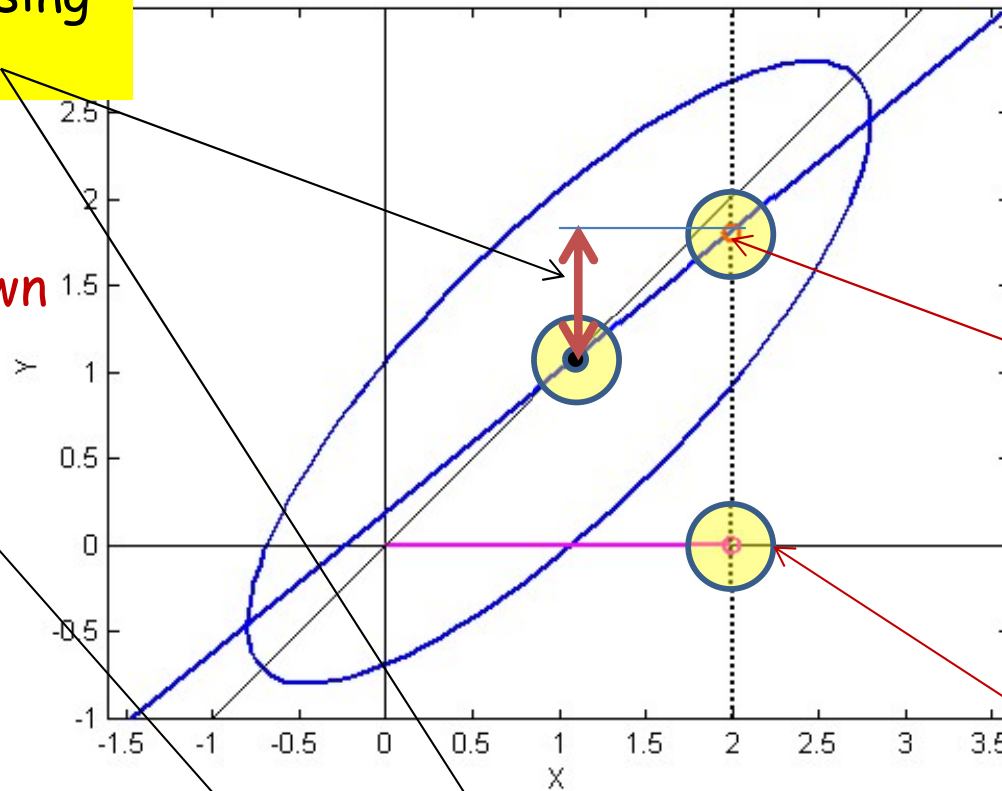
$$P(y|x) = N(\mu_y + C_{yx} C_{xx}^{-1} (x - \mu_x), C_{yy} - C_{yx} C_{xx}^{-1} C_{xy})$$

# Preliminaries : $P(y|x)$ for Gaussian

Update guess of  $Y$  based on information in  $X$   
Correction is 0 if  $X$  and  $Y$  are uncorrelated, i.e  $C_{yx} = 0$

Correction of  $Y$  using information in  $X$

Best guess for  $Y$  when  $X$  is not known



Mean of  $Y$  given  $X$

Given  $X$  value

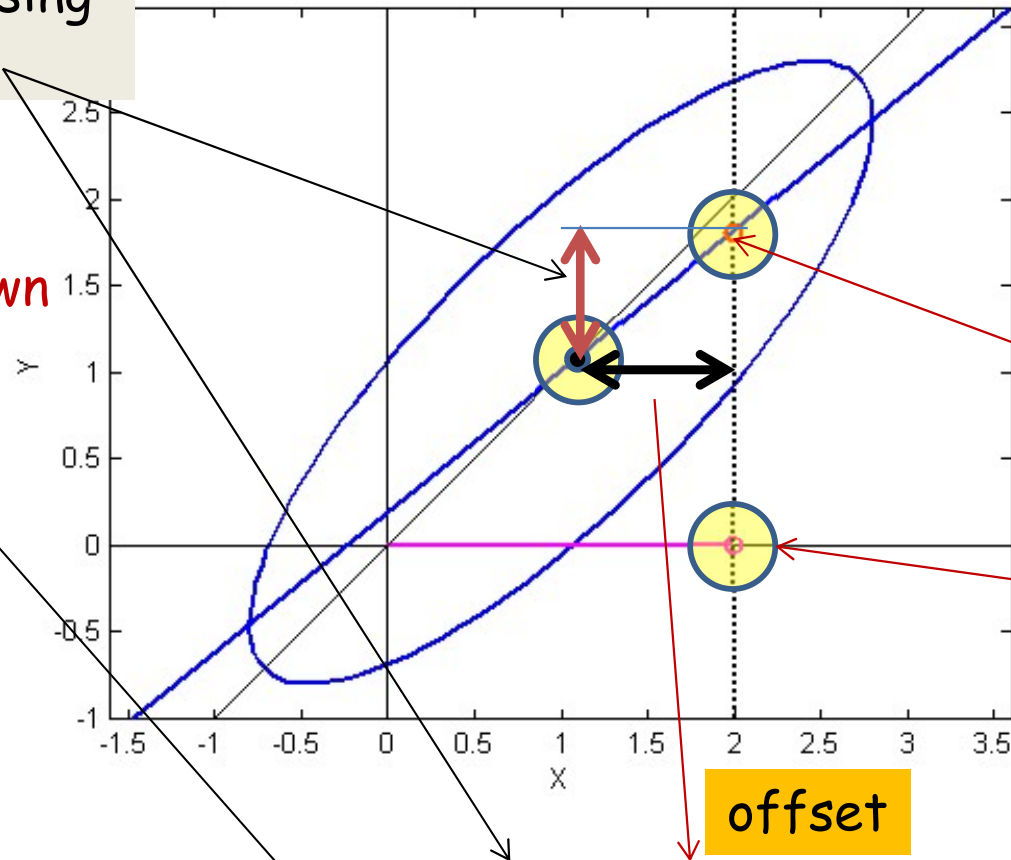
$$P(y|x) = N(\mu_y + C_{yx} C_{xx}^{-1} (x - \mu_x), C_{yy} - C_{yx} C_{xx}^{-1} C_{xy})$$

# Preliminaries : P(y|x) for Gaussian

Correction to Y = slope \* (offset of X from mean)

Correction of Y using information in X

Best guess for Y when X is not known



Mean of Y given X

Given X value

$$P(y | x) = N(\mu_y + C_{yx} C_{xx}^{-1} (x - \mu_x), C_{yy} - C_{yx} C_{xx}^{-1} C_{xy})$$

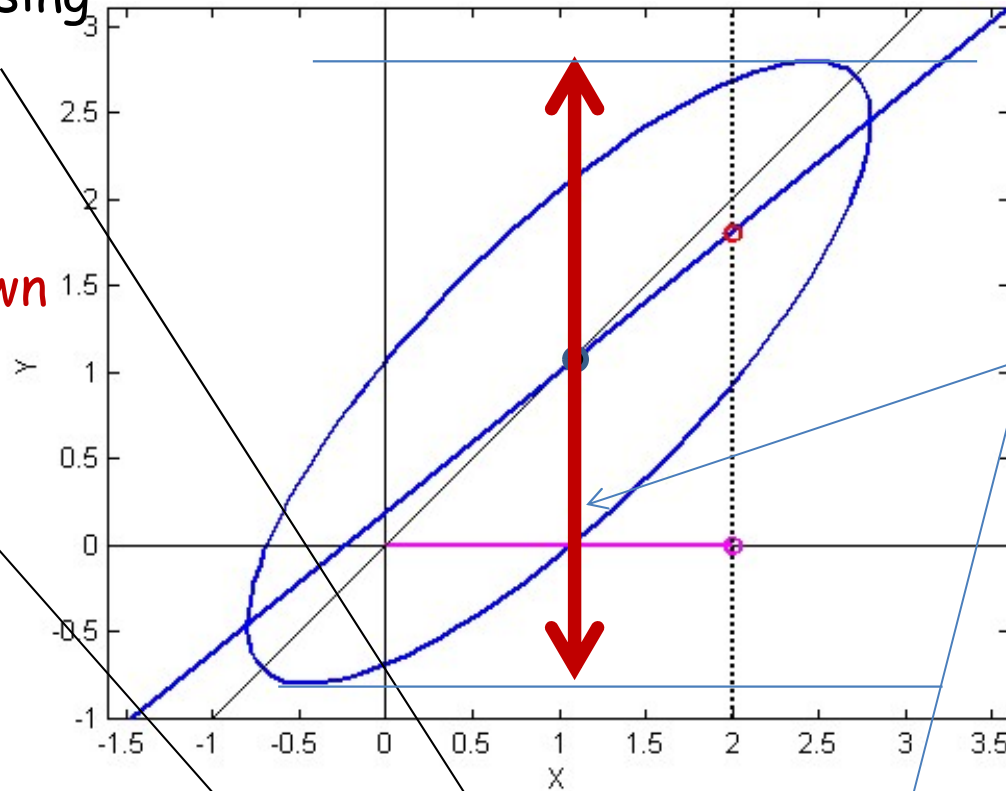
Slope

8797

# Preliminaries : $P(y|x)$ for Gaussian

Correction of  $Y$  using information in  $X$

Best guess for  $Y$  when  $X$  is not known



Uncertainty in  $Y$  when  $X$  is not known

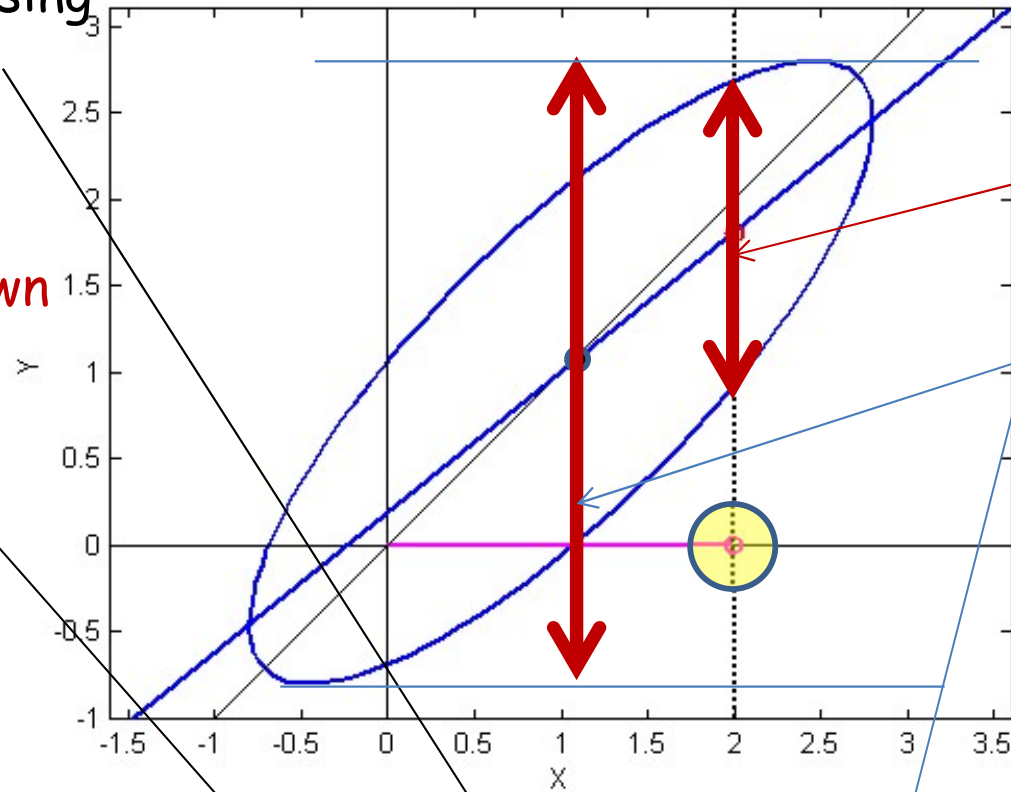
$$P(y|x) = N(\mu_y + C_{yx} C_{xx}^{-1} (x - \mu_x), C_{yy} - C_{yx} C_{xx}^{-1} C_{xy})$$

# Preliminaries : $P(y|x)$ for Gaussian

Shrinkage of variance is 0 if  $X$  and  $Y$  are uncorrelated, i.e  $C_{yx} = 0$

Correction of  $Y$  using information in  $X$

Best guess for  $Y$  when  $X$  is not known



Reduced uncertainty from knowing  $X$

Uncertainty in  $Y$  when  $X$  is not known

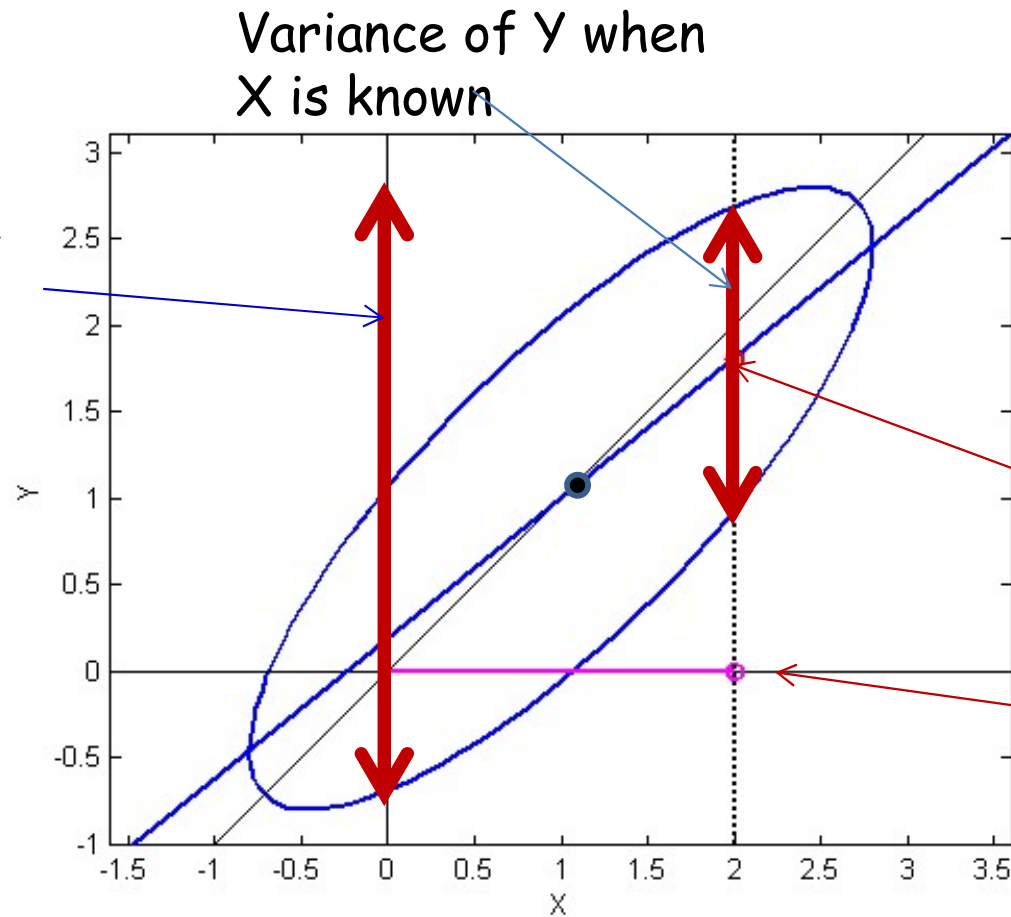
Shrinkage of uncertainty from knowing  $X$

$$P(y|x) = N(\mu_y + C_{yx} C_{xx}^{-1} (x - \mu_x), C_{yy} - C_{yx} C_{xx}^{-1} C_{xy})$$

# Preliminaries : $P(y|x)$ for Gaussian

Knowing  $X$  modifies the mean of  $Y$  and shrinks its variance

Overall variance of  $Y$  when  $X$  is unknown



Mean of  $Y$  given  $X$   
(MAP estimate of  $Y$ )

Given  $X$  value

$$P(y|x) = N(\mu_y + C_{yx} C_{xx}^{-1} (x - \mu_x), C_{yy} - C_{yx} C_{xx}^{-1} C_{xy})$$



# Background: Sum of Gaussian RVs

$$O = AS + \varepsilon$$

$$S \sim N(\mu_s, \Theta_s)$$

$$\varepsilon \sim N(\mu_\varepsilon, \Theta_\varepsilon)$$

- Consider a random variable  $O$  obtained as above
- The expected value of  $O$  is given by

$$E[O] = E[AS + \varepsilon] = A\mu_s + \mu_\varepsilon$$

- Notation:

$$E[O] = \mu_o$$

# Background: Sum of Gaussian RVs

$$\mathbf{O} = \mathbf{A}\mathbf{S} + \boldsymbol{\varepsilon}$$

$$\mathbf{S} \sim N(\boldsymbol{\mu}_s, \boldsymbol{\Theta}_s)$$

$$\boldsymbol{\varepsilon} \sim N(\boldsymbol{\mu}_\varepsilon, \boldsymbol{\Theta}_\varepsilon)$$

- The variance of  $\mathbf{O}$  is given by

$$\text{Var}(\mathbf{O}) = \boldsymbol{\Theta}_o = E[(\mathbf{O} - \boldsymbol{\mu}_o)(\mathbf{O} - \boldsymbol{\mu}_o)^T]$$

- This is just the sum of the variance of  $\mathbf{A}\mathbf{S}$  and the variance of  $\boldsymbol{\varepsilon}$

$$\boldsymbol{\Theta}_o = \mathbf{A}\boldsymbol{\Theta}_s\mathbf{A}^T + \boldsymbol{\Theta}_\varepsilon$$

# Background: Sum of Gaussian RVs

$$\mathbf{O} = \mathbf{A}\mathbf{S} + \boldsymbol{\varepsilon}$$

$$\mathbf{S} \sim N(\boldsymbol{\mu}_s, \boldsymbol{\Theta}_s)$$

$$\boldsymbol{\varepsilon} \sim N(\boldsymbol{\mu}_\varepsilon, \boldsymbol{\Theta}_\varepsilon)$$

- The conditional probability of  $\mathbf{O}$ :

$$P(\mathbf{O}|\mathbf{S}) = N(\mathbf{A}\mathbf{S} + \boldsymbol{\mu}_\varepsilon, \boldsymbol{\Theta}_\varepsilon)$$

- The overall probability of  $\mathbf{O}$ :

$$P(\mathbf{O}) = N(\mathbf{A}\boldsymbol{\mu}_s + \boldsymbol{\mu}_\varepsilon, \mathbf{A}\boldsymbol{\Theta}_s\mathbf{A}^T + \boldsymbol{\Theta}_\varepsilon)$$

# Background: Sum of Gaussian RVs

$$O = AS + \varepsilon$$

$$S \sim N(\mu_S, \Theta_S)$$

$$\varepsilon \sim N(\mu_\varepsilon, \Theta_\varepsilon)$$

- The *cross-correlation* between  $O$  and  $S$

$$\begin{aligned} \Theta_{OS} &= E[(O - \mu_O)(S - \mu_S)^T] = E[(A(S - \mu_S) + (\varepsilon - \mu_\varepsilon))(S - \mu_S)^T] \\ &= E[A(S - \mu_S)(S - \mu_S)^T + (\varepsilon - \mu_\varepsilon)(S - \mu_S)^T] \\ &= AE[(S - \mu_S)(S - \mu_S)^T] + E[(\varepsilon - \mu_\varepsilon)(S - \mu_S)^T] \\ &= AE[(S - \mu_S)(S - \mu_S)^T] \end{aligned}$$

- $= A \Theta_S$

- The cross-correlation between  $O$  and  $S$  is

$$\Theta_{OS} = A \Theta_S$$

$$\Theta_{SO} = \Theta_S A^T$$

# Background: Joint Prob. of $O$ and $S$

$$O = AS + \varepsilon$$

$$Z = \begin{bmatrix} O \\ S \end{bmatrix}$$

- The joint probability of  $O$  and  $S$  (i.e.  $P(Z)$ ) is also Gaussian

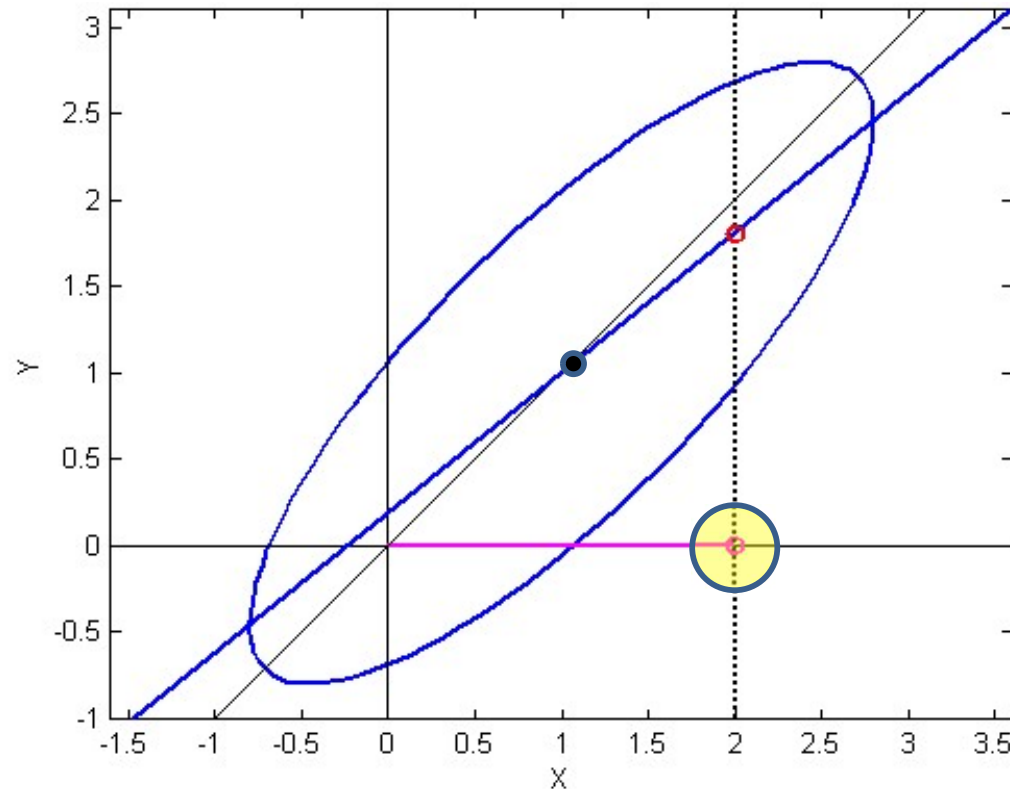
$$P(Z) = P(O, S) = N(\mu_Z, \Theta_Z)$$

- Where

$$\mu_Z = \begin{bmatrix} \mu_O \\ \mu_S \end{bmatrix} = \begin{bmatrix} A\mu_S + \mu_\varepsilon \\ \mu_S \end{bmatrix}$$

$$\Theta_Z = \begin{bmatrix} \Theta_O & \Theta_{OS} \\ \Theta_{SO} & \Theta_S \end{bmatrix} = \begin{bmatrix} A\Theta_S A^T + \Theta_\varepsilon & A\Theta_S \\ \Theta_S A^T & \Theta_S \end{bmatrix}$$

# Preliminaries : Conditional of S given O: $P(S|O)$



$$O = AS + \varepsilon$$

$$P(S|O) = N(\mu_S + \Theta_{SO} \Theta_O^{-1} (O - \mu_O), \Theta_S - \Theta_{SO} \Theta_O^{-1} \Theta_{OS})$$

$$P(S|O) = N(\mu_S + \Theta_S A^T (A \Theta_S A^T + \Theta_\varepsilon)^{-1} (O - A \mu_S - \mu_\varepsilon), \Theta_S - \Theta_S A^T (A \Theta_S A^T + \Theta_\varepsilon)^{-1} A \Theta_S)$$

# Poll 1

- X and Y are jointly Gaussian. Which of the following are true
  - Knowing X affects our expectation of Y, in all cases
  - Knowing X affects our expectation of Y if the two are correlated
  - Knowing X reduces the variance of the conditional distribution of Y by a value that depends on the observed X
  - Knowing X reduces the variance of Y by the same amount regardless of the observed X
- We are given that  $Y = AX + e$ , where X and e are Gaussian. Mark all that are true
  - Y and X are jointly Gaussian
  - The conditional distribution of X given Y is Gaussian
  - Knowing Y does not influence the variance of X, since Y is derived from X and not vice versa
  - Knowing Y does not influence the expected value of X since Y is derived from X and not vice versa

# Poll 1

- X and Y are jointly Gaussian. Which of the following are true
  - Knowing X affects our expectation of Y, in all cases
  - **Knowing X affects our expectation of Y if the two are correlated**
  - Knowing X reduces the variance of the conditional distribution of Y by a value that depends on the observed X
  - **Knowing X reduces the variance of Y by the same amount regardless of the observed X**
- We are given that  $Y = AX + e$ , where X and e are Gaussian. Mark all that are true
  - **Y and X are jointly Gaussian**
  - **The conditional distribution of X given Y is Gaussian**
  - Knowing Y does not influence the variance of X, since Y is derived from X and not vice versa
  - Knowing Y does not influence the expected value of X since Y is derived from X and not vice versa



# The little parable

You've been kidnapped

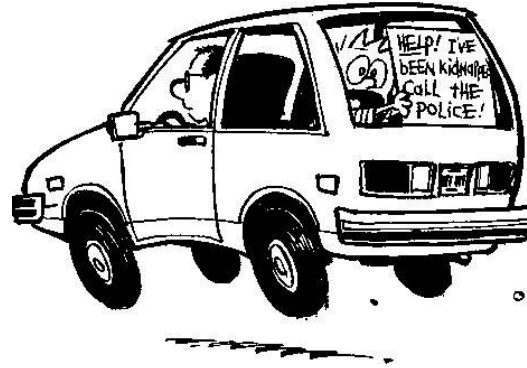


And blindfolded

You can only *hear* the car

You must find your way back home from wherever they drop you off

# Kidnapped!

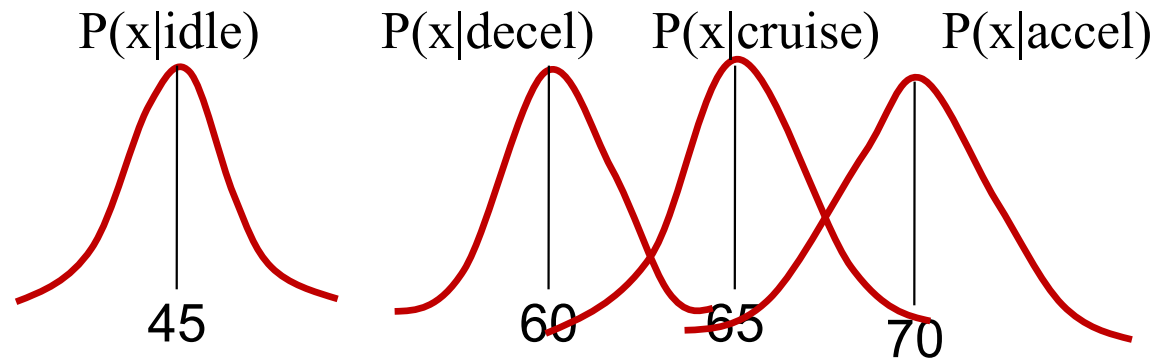


- Determine by only *listening* to a running automobile, if it is:
  - Idling; or
  - Travelling at constant velocity; or
  - Accelerating; or
  - Decelerating
- You only record energy level (SPL) in the sound
  - The SPL is measured once per second

# What we know

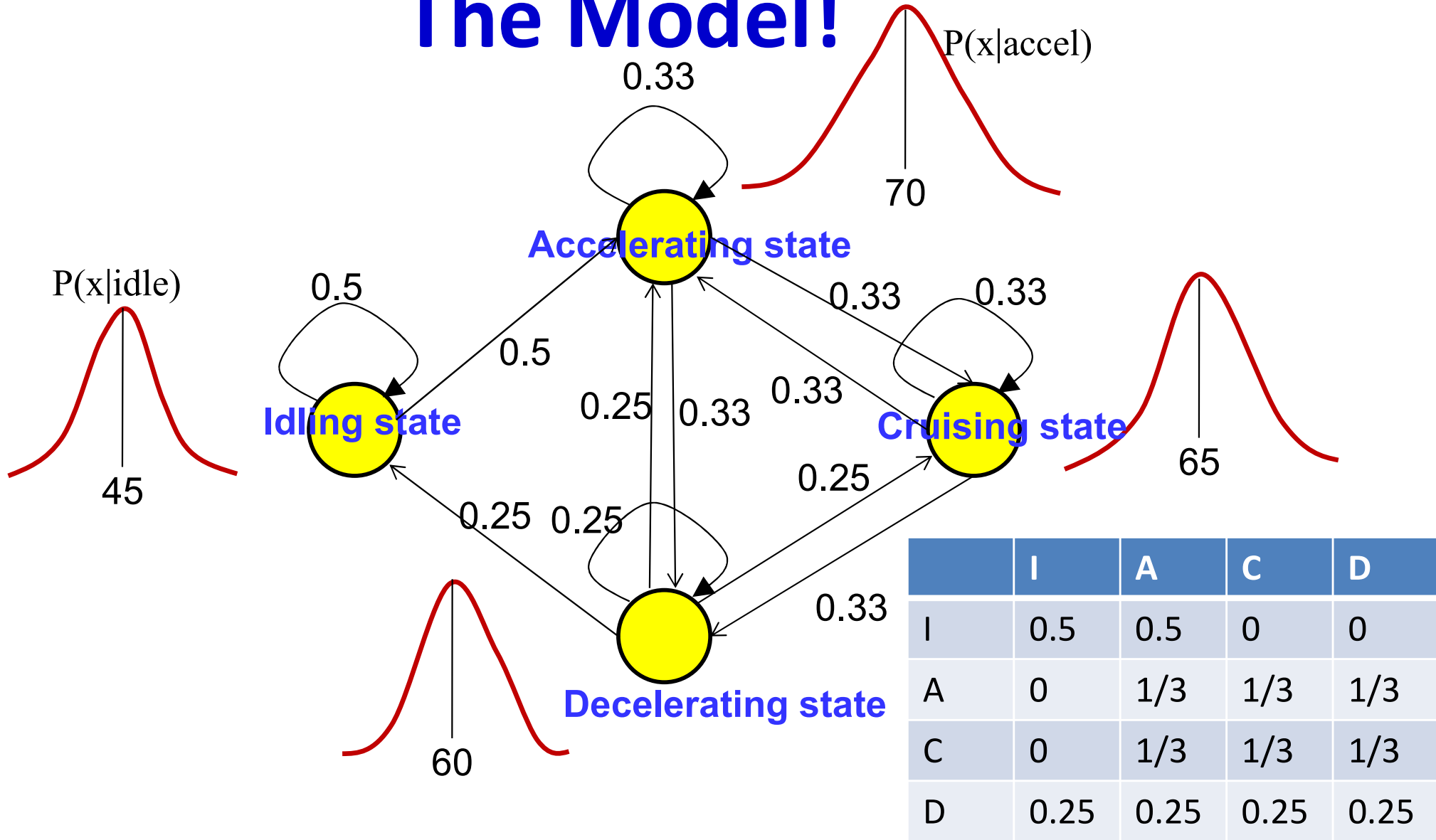
- An automobile that is at rest can accelerate, or continue to stay at rest
- An accelerating automobile can hit a steady-state velocity, continue to accelerate, or decelerate
- A decelerating automobile can continue to decelerate, come to rest, cruise, or accelerate
- A automobile at a steady-state velocity can stay in steady state, accelerate or decelerate

# What else we know



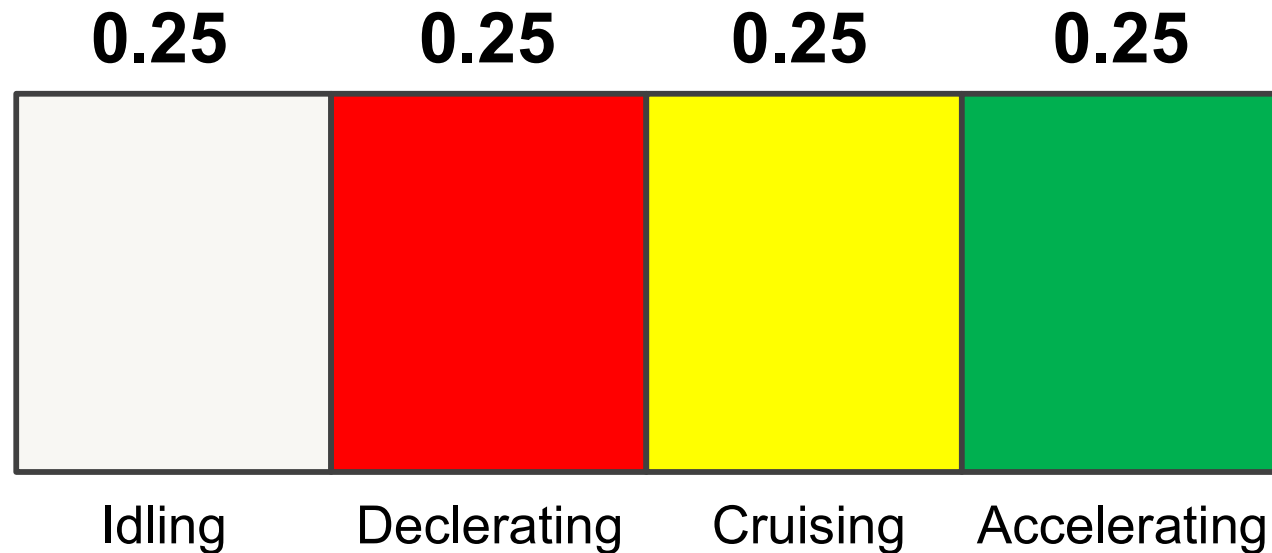
- The probability distribution of the SPL of the sound is different in the various conditions
  - As shown in figure
    - In reality, depends on the car
- The distributions for the different conditions overlap
  - Simply knowing the current sound level is not enough to know the state of the car

# The Model!



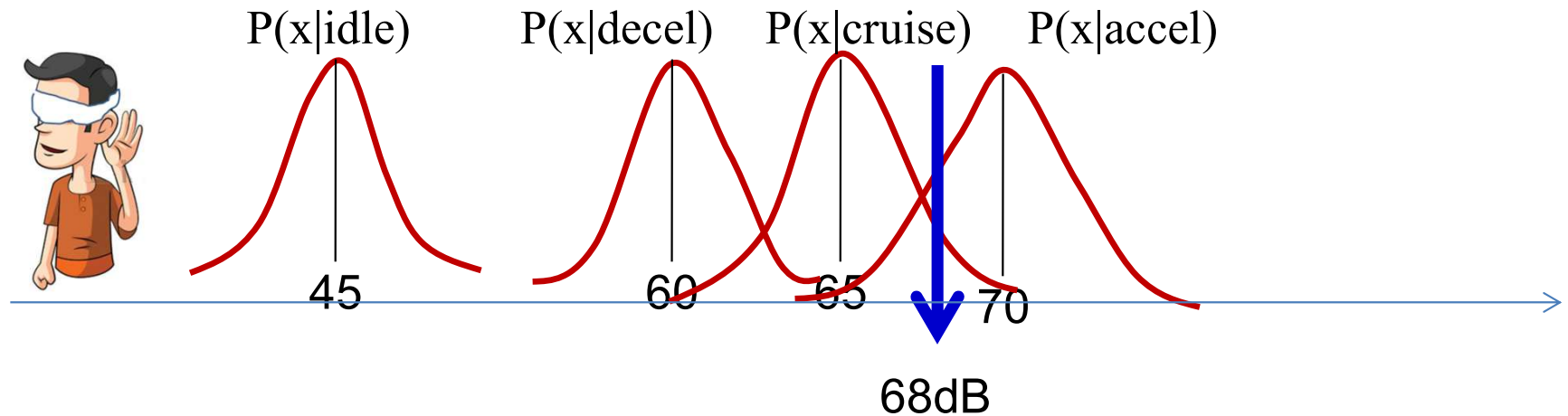
- The state-space model
  - Assuming all transitions from a state are equally probable
  - This is a Hidden Markov Model!

# Estimating the state at $T = 0$ -



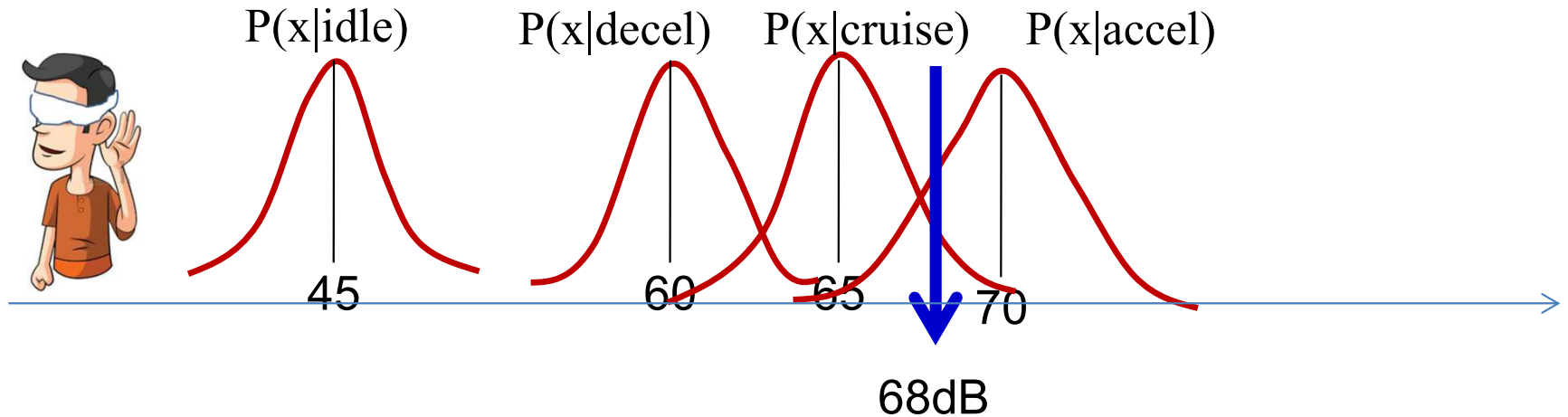
- A  $T=0$ , before the first observation, we know nothing of the state
  - Assume all states are equally likely

# The first observation: $T=0$



- At  $T=0$  you observe the sound level  $x_0 = 68\text{dB}$  SPL
- The observation modifies our belief in the state of the system

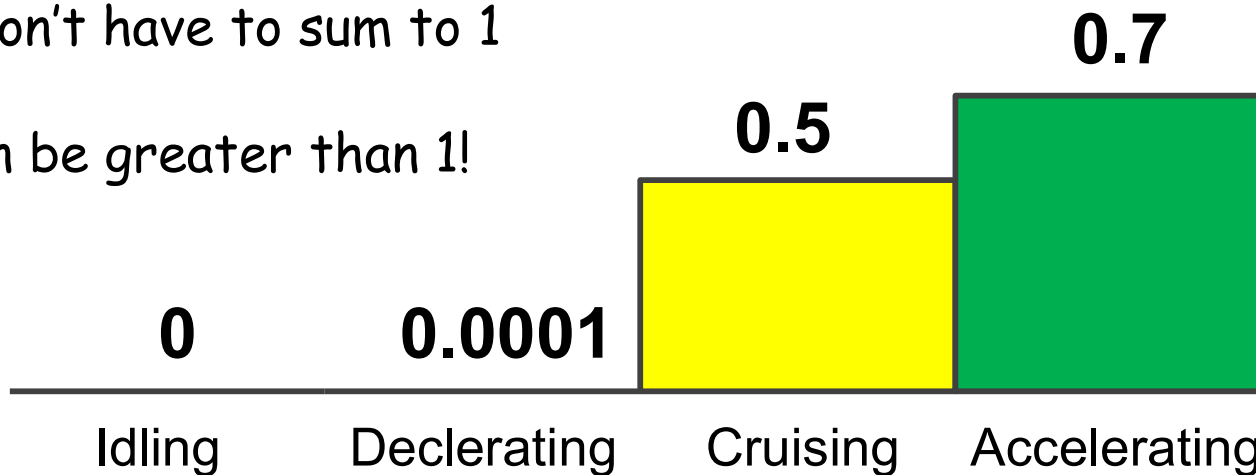
# The first observation: $T=0$



$P(x   \text{idle})$	$P(x   \text{deceleration})$	$P(x   \text{cruising})$	$P(x   \text{acceleration})$
0	0.0001	0.5	0.7

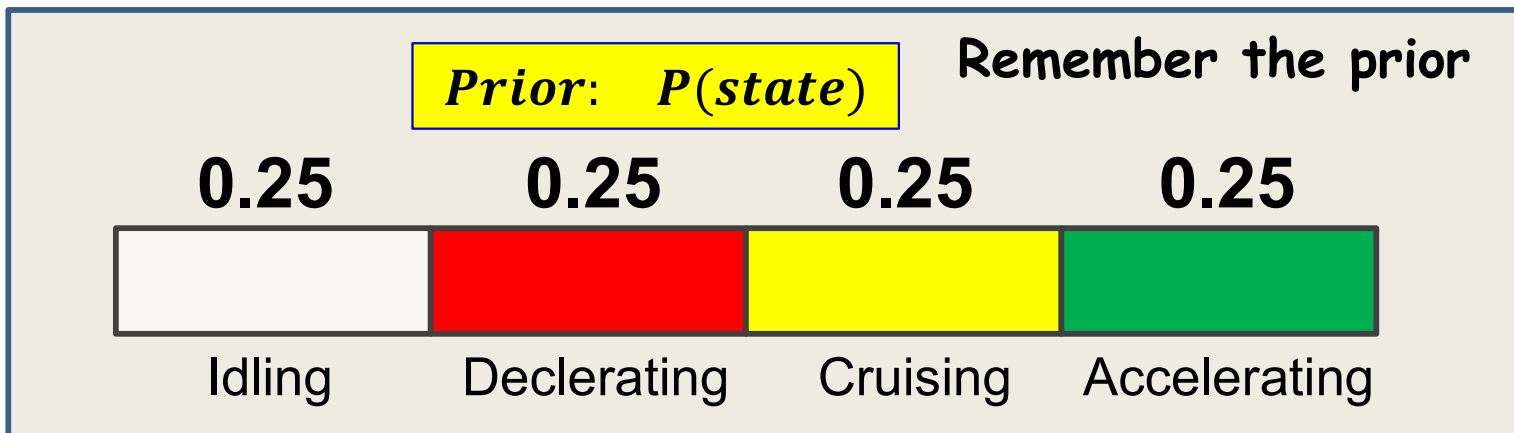
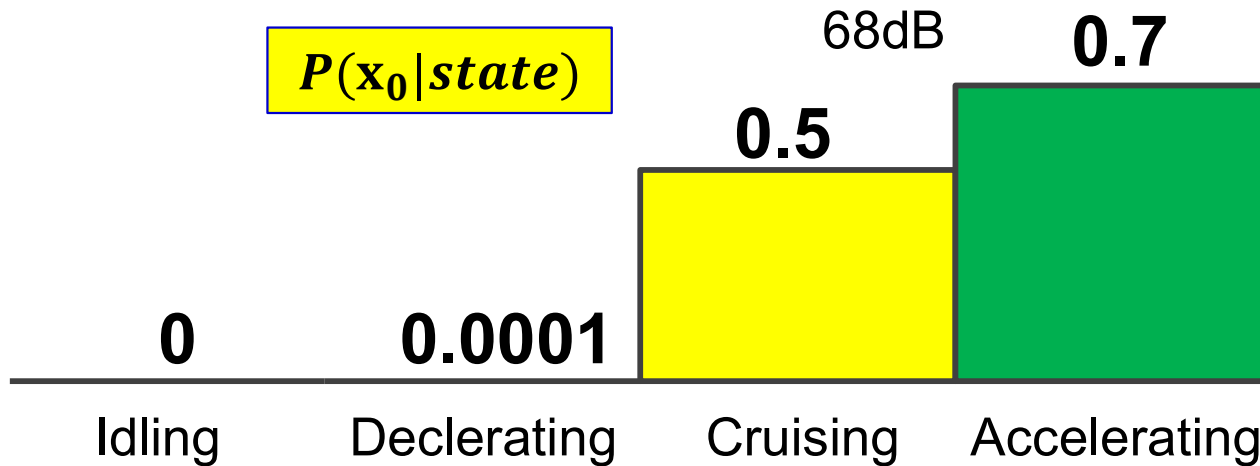
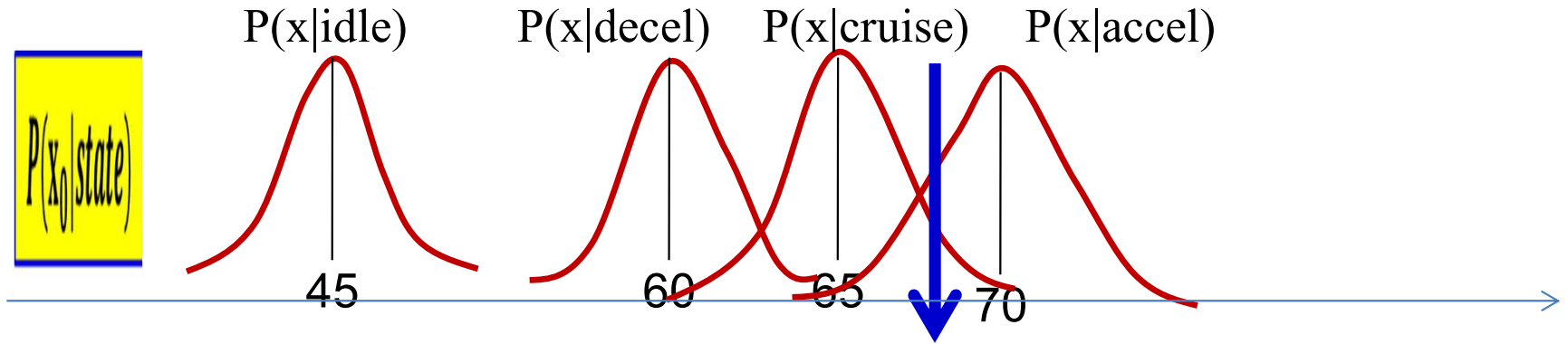
These don't have to sum to 1

Can even be greater than 1!





# The first observation: $T=0$

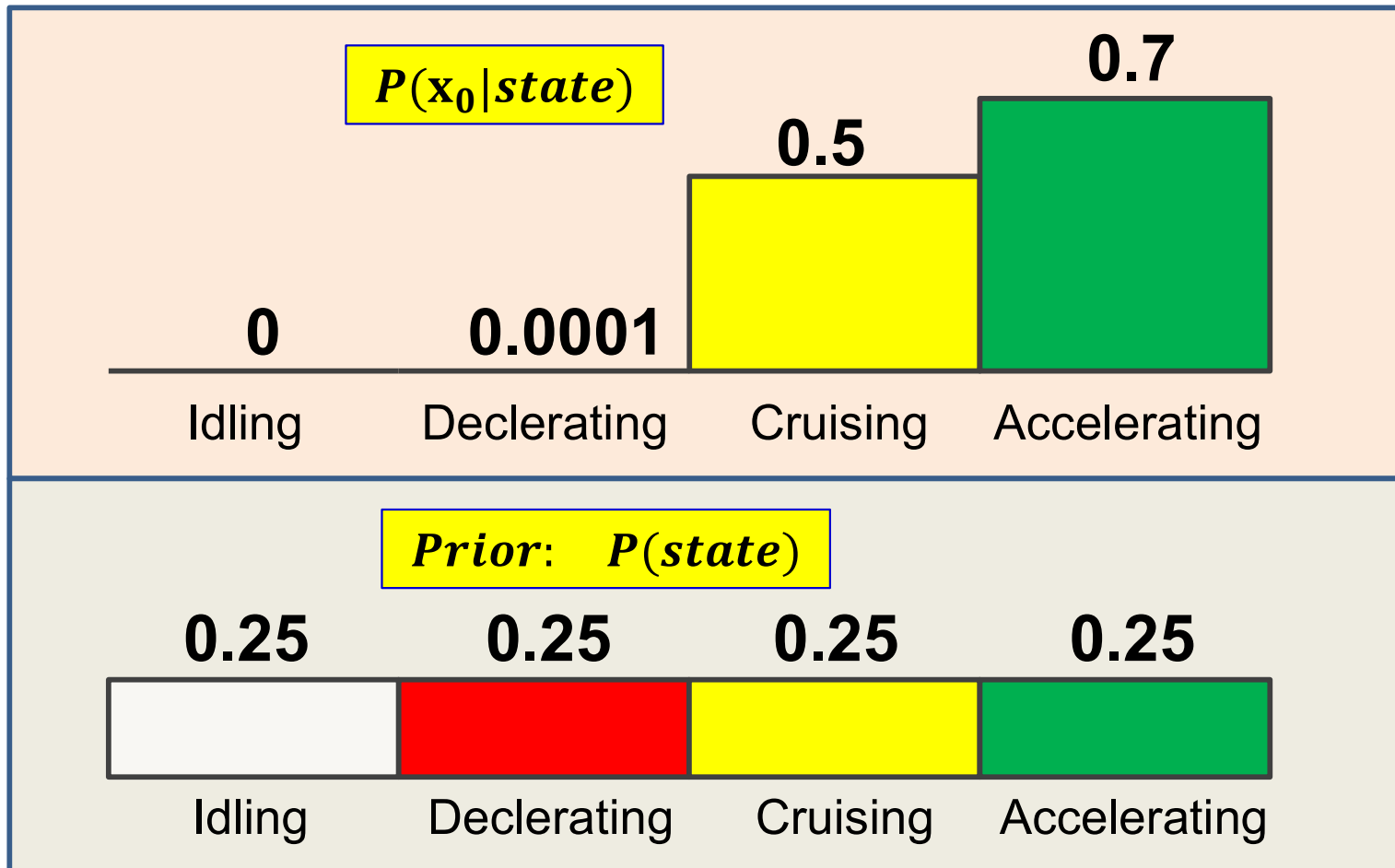


# Estimating state *after* at observing $\mathbf{x}_0$

- Combine prior information about state and evidence from observation
- We want  $P(\text{state}|\mathbf{x}_0)$
- We can compute it using Bayes rule as

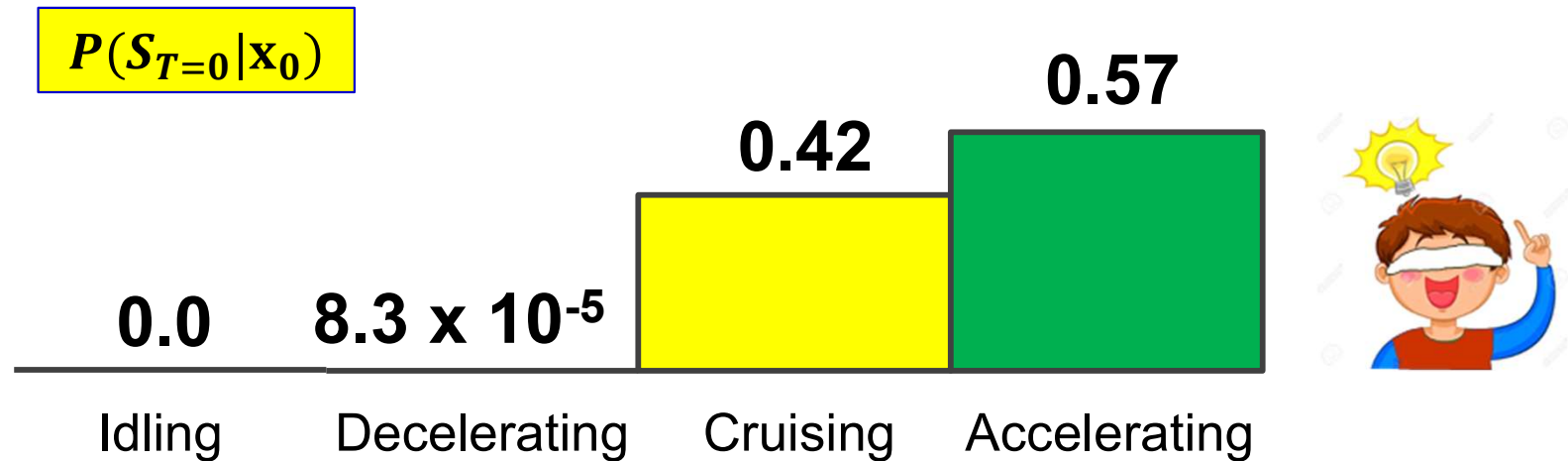
$$P(\text{state}|\mathbf{x}_0) = \frac{P(\text{state})P(\mathbf{x}_0|\text{state})}{\sum_{\text{state}'} P(\text{state}')P(\mathbf{x}_0|\text{state}')}$$

# The Posterior



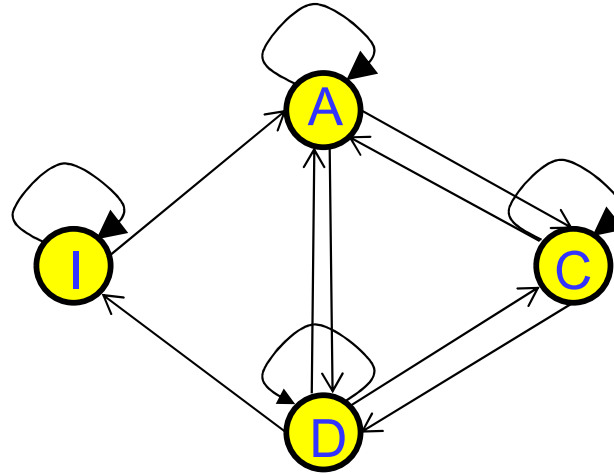
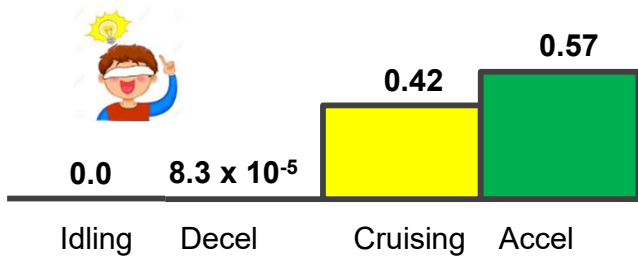
- Multiply the two, term by term, and normalize them so that they sum to 1.0

# Estimating the state at $T = 0+$



- At  $T=0$ , after the first observation  $\mathbf{x}_0$ , we update our belief about the states
  - The first observation provided some evidence about the state of the system
  - It modifies our belief in the state of the system

# Predicting the state at T=1

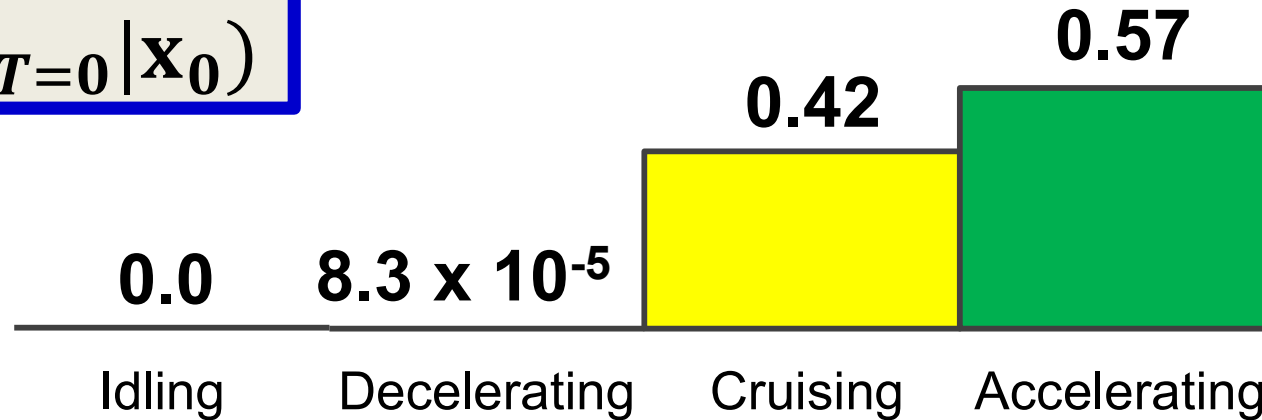


	I	A	C	D
I	0.5	0.5	0	0
A	0	1/3	1/3	1/3
C	0	1/3	1/3	1/3
D	0.25	0.25	0.25	0.25

- Predicting the probability of idling at T=1
  - $P(\text{idling} \mid \text{idling}) = 0.5$ ;
  - $P(\text{idling} \mid \text{deceleration}) = 0.25$
  - $P(\text{idling at } T=1 \mid \mathbf{x}_0) =$   
 $P(I_{T=0} \mid \mathbf{x}_0) P(I \mid I) + P(D_{T=0} \mid \mathbf{x}_0) P(I \mid D) = 2.1 \times 10^{-5}$
- In general, for any state S
  - $P(S_{T=1} \mid \mathbf{x}_0) = \sum_{S_{T=0}} P(S_{T=0} \mid \mathbf{x}_0) P(S_{T=1} \mid S_{T=0})$

# Predicting the state at T = 1

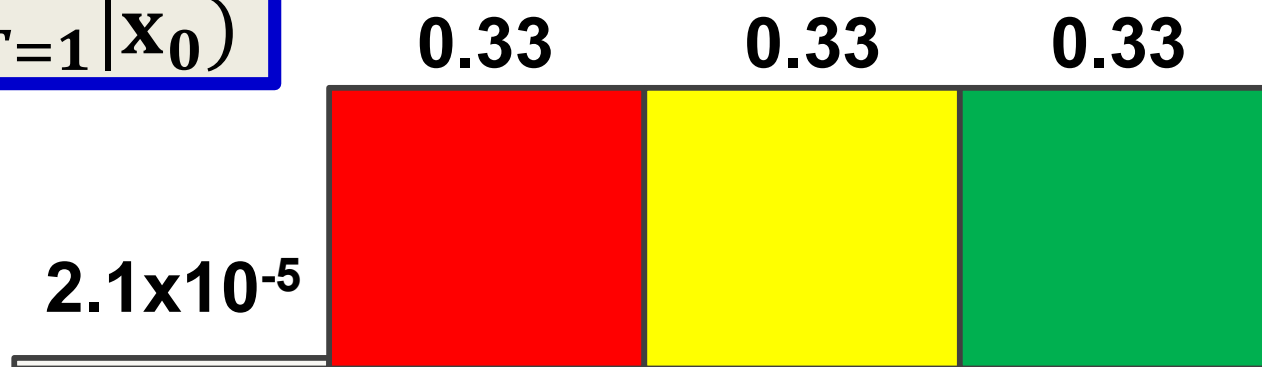
$$P(S_{T=0} | \mathbf{x}_0)$$



$$P(S_{T=1} | \mathbf{x}_0) = \sum_{S_{T=0}} P(S_{T=0} | \mathbf{x}_0) P(S_{T=1} | S_{T=0})$$

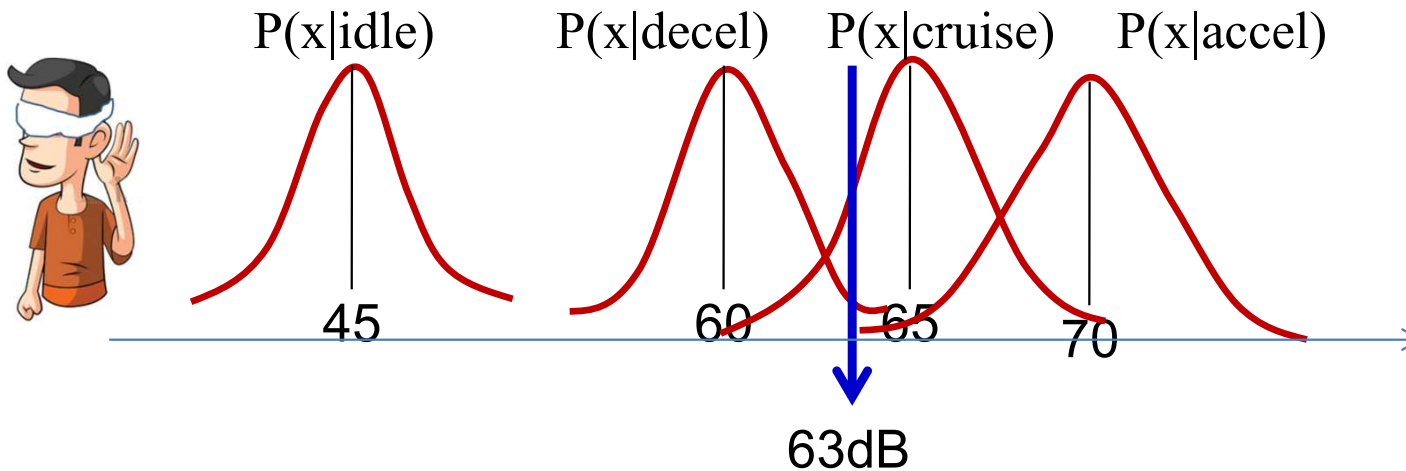


$$P(S_{T=1} | \mathbf{x}_0)$$



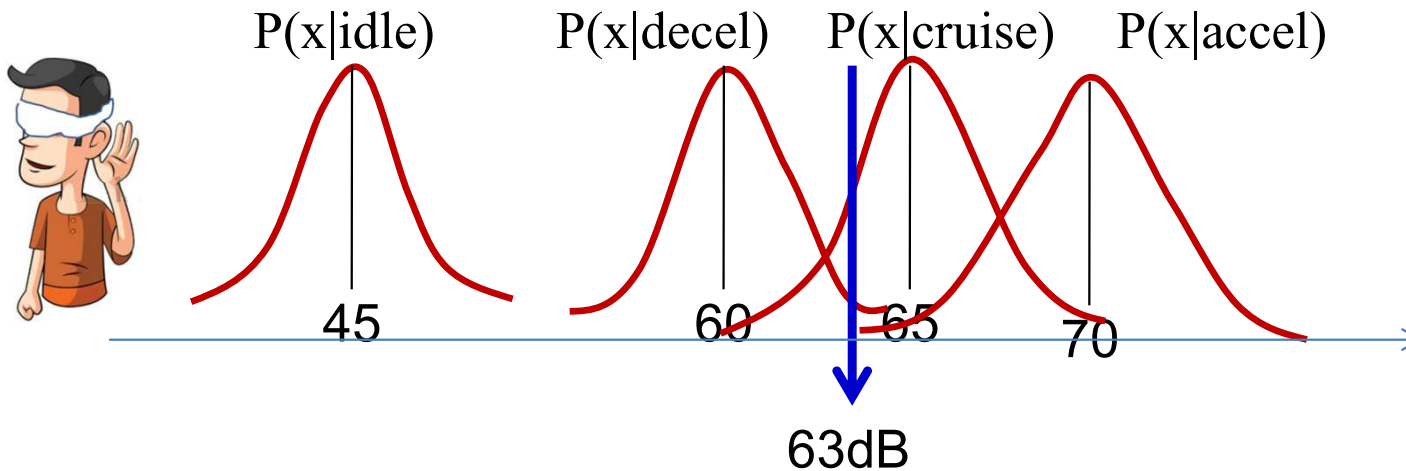
Rounded.  
In reality, they  
sum to 1.0

# Updating after the observation at $T=1$

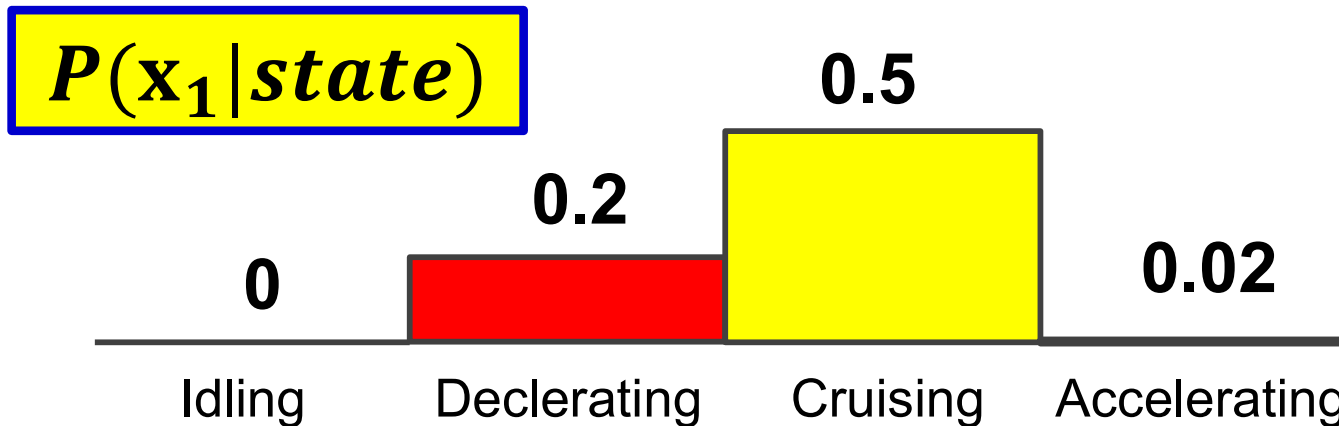


- At  $T=1$  we observe  $x_1 = 63\text{dB SPL}$

# Updating after the observation at T=1

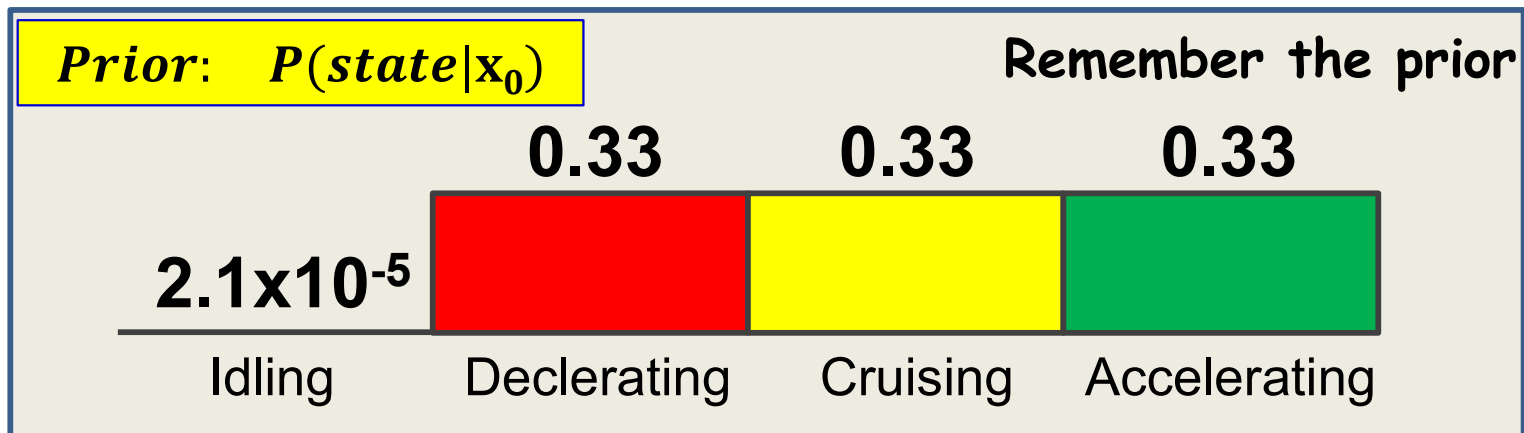
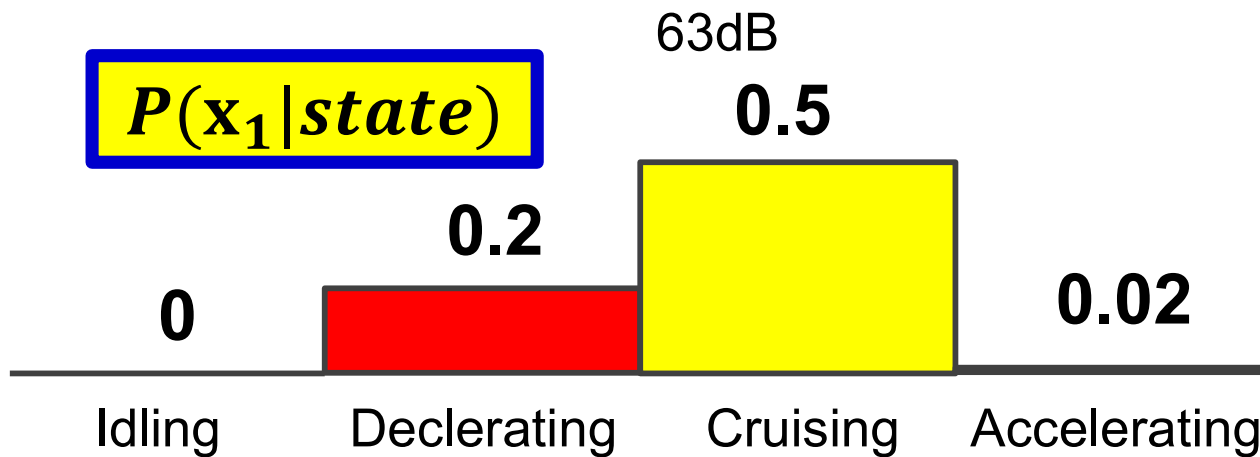
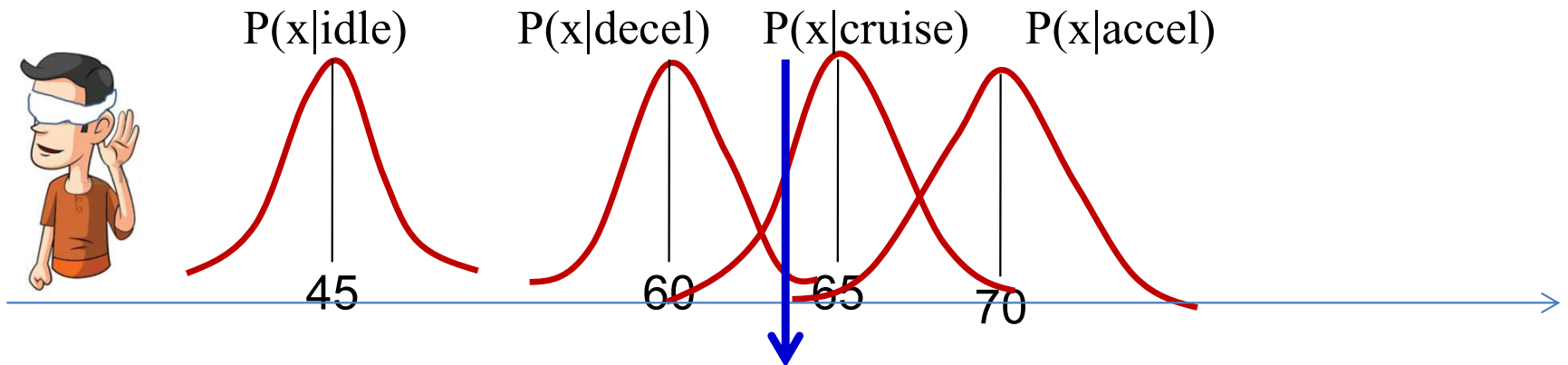


$P(x   \text{idle})$	$P(x   \text{deceleration})$	$P(x   \text{cruising})$	$P(x   \text{acceleration})$
0	0.2	0.5	0.01





# The second observation: T=1

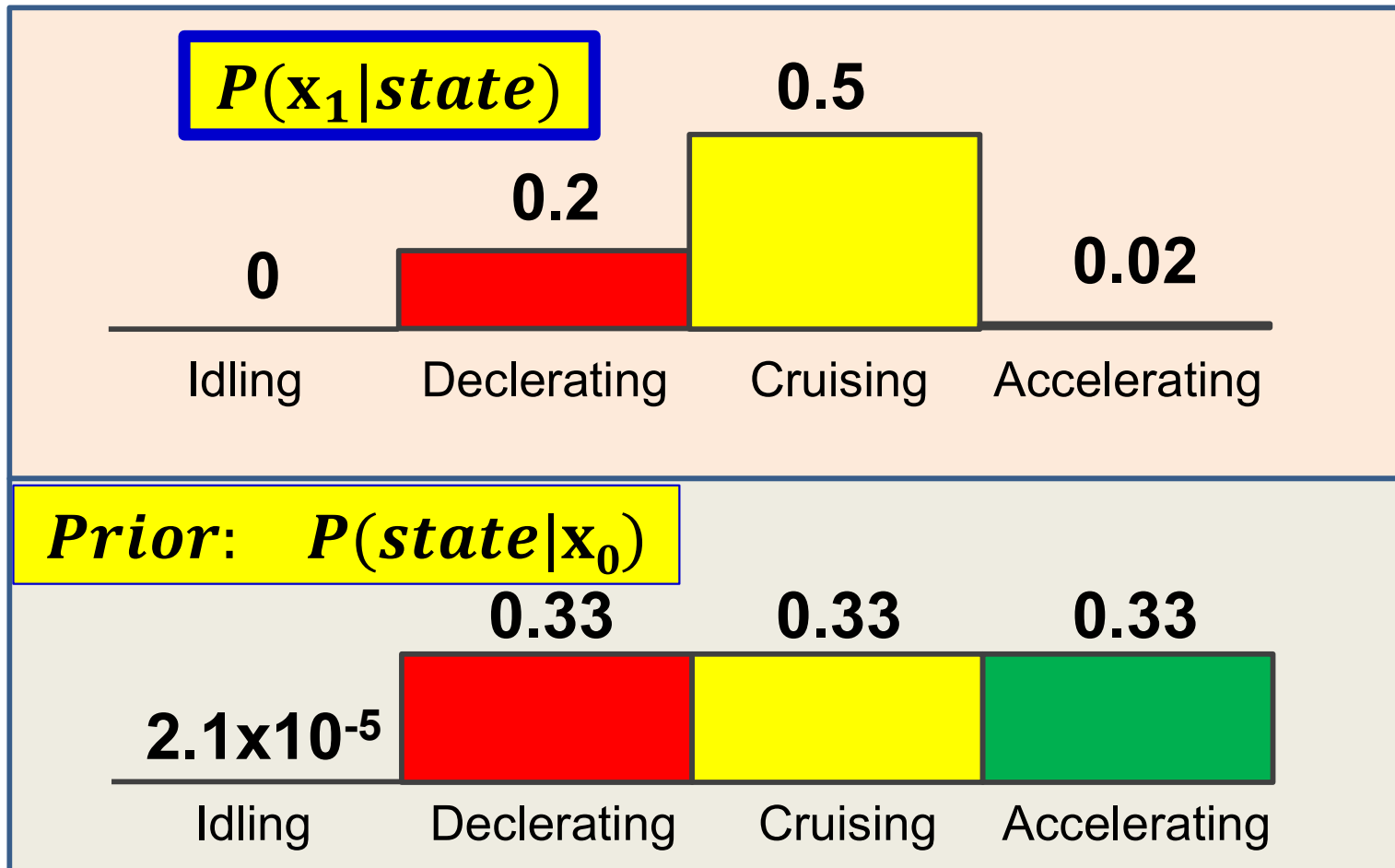


# Estimating state *after* at observing $\mathbf{x}_1$

- Combine prior information from the observation at time  $T=0$ , AND evidence from observation at  $T=1$  to estimate **state** at  $T=1$
- We want  $P(\text{state}|\mathbf{x}_0, \mathbf{x}_1)$
- We can compute it using Bayes rule as

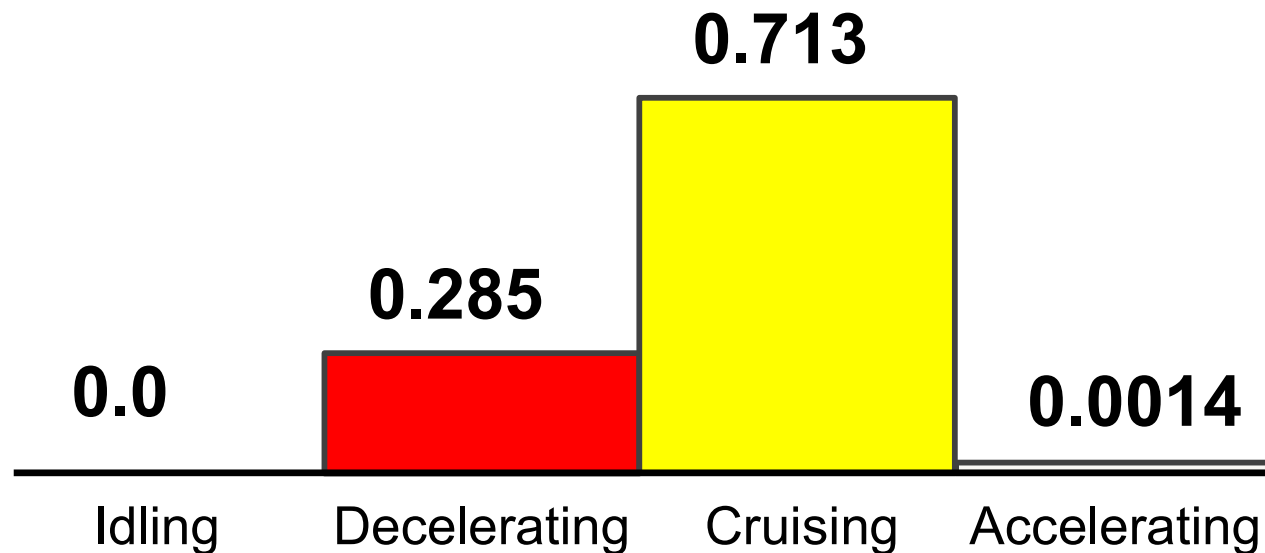
$$P(\text{state}|\mathbf{x}_0, \mathbf{x}_1) = \frac{P(\text{state}|\mathbf{x}_0)P(\mathbf{x}_1|\text{state})}{\sum_{\text{state}'} P(\text{state}'|\mathbf{x}_0)P(\mathbf{x}_1|\text{state}')}$$

# The Posterior at $T = 1$



- Multiply the two, term by term, and normalize them so that they sum to 1.0

# Estimating the state at $T = 1+$



- The updated probability at  $T=1$  incorporates information from both  $x_0$  and  $x_1$ 
  - It is NOT a local decision based on  $x_1$  alone
  - Because of the Markov nature of the process, the state at  $T=0$  affects the state at  $T=1$ 
    - $x_0$  provides evidence for the state at  $T=1$

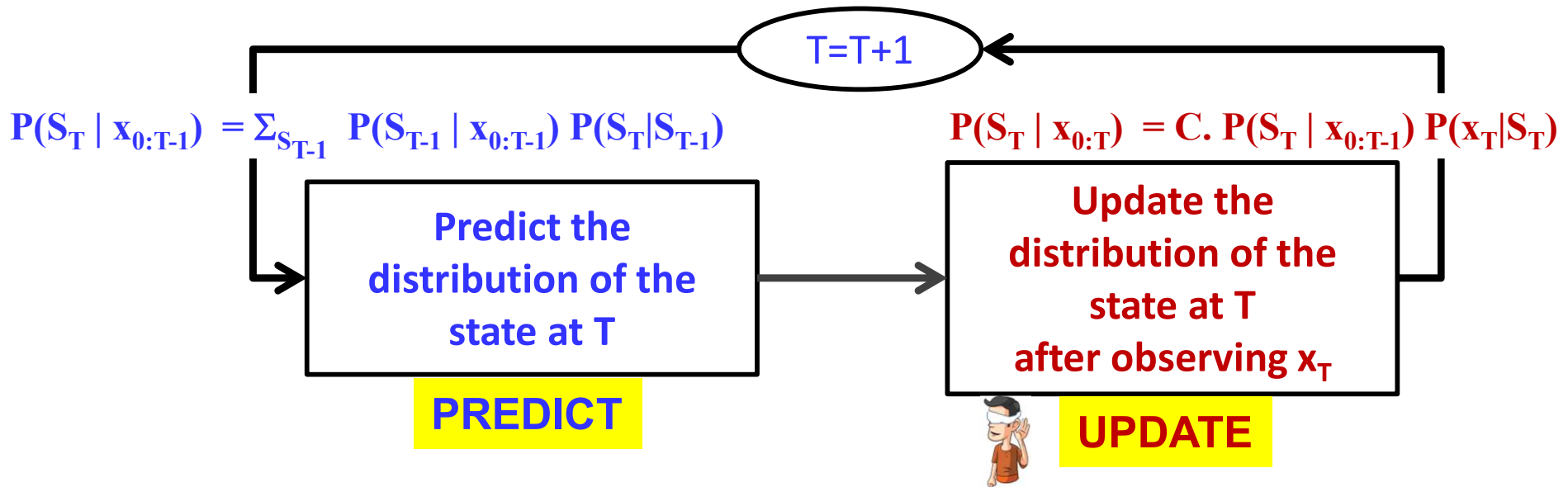
# Overall Process

## Time

## Computation

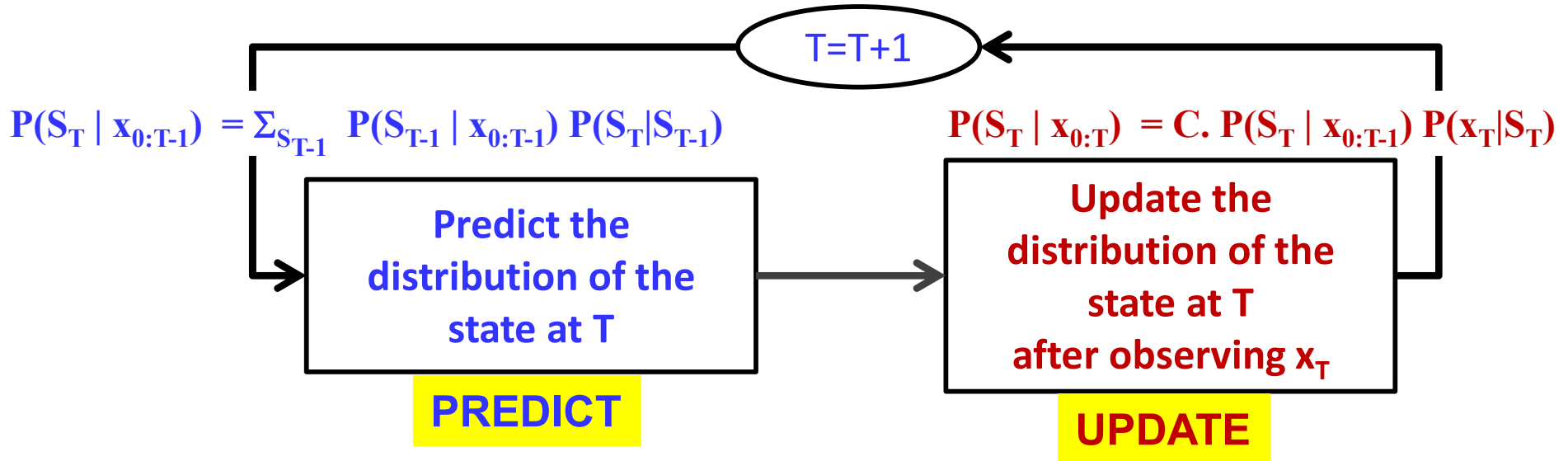
- |   |   |
|---|---|
| <ul style="list-style-type: none"> <li>• T=0- : A priori probability</li> <li>• T = 0+: Update after <math>X_0</math></li> </ul>                | <ul style="list-style-type: none"> <li>• <math>P(S_0) = P(S)</math></li> <li>• <math>P(S_0 X_0) = C \cdot P(S_0)P(X_0 S_0)</math></li> </ul>  |
| <ul style="list-style-type: none"> <li>• T=1- (Prediction before <math>X_1</math>)</li> <li>• T = 1+: Update after <math>X_1</math></li> </ul>  | <ul style="list-style-type: none"> <li>• <math>P(S_1 X_0) = \sum_{S_0} P(S_1 S_0)P(S_0 X_0)</math></li> <li>• <math>P(S_1 X_{0:1}) = C \cdot P(S_1 X_0)P(X_1 S_1)</math></li> </ul>                               |
| <ul style="list-style-type: none"> <li>• T=2- (Prediction before <math>X_2</math>)</li> <li>• T = 2+: Update after <math>X_2</math></li> </ul>  | <ul style="list-style-type: none"> <li>• <math>P(S_2 X_{0:1}) = \sum_{S_1} P(S_2 S_1)P(S_1 X_{0:1})</math></li> <li>• <math>P(S_2 X_{0:2}) = C \cdot P(S_2 X_{0:1})P(X_2 S_2)</math></li> </ul>                   |
| <ul style="list-style-type: none"> <li>• ...</li> </ul>   | <ul style="list-style-type: none"> <li>• ...</li> </ul>   |
| <ul style="list-style-type: none"> <li>• T= t- (Prediction before <math>X_t</math>)</li> <li>• T = t+: Update after <math>X_t</math></li> </ul> | <ul style="list-style-type: none"> <li>• <math>P(S_t X_{0:t-1}) = \sum_{S_{t-1}} P(S_t S_{t-1})P(S_{t-1} X_{0:t-1})</math></li> <li>• <math>P(S_t X_{0:t}) = C \cdot P(S_t X_{0:t-1})P(X_t S_t)</math></li> </ul> |

# Overall procedure



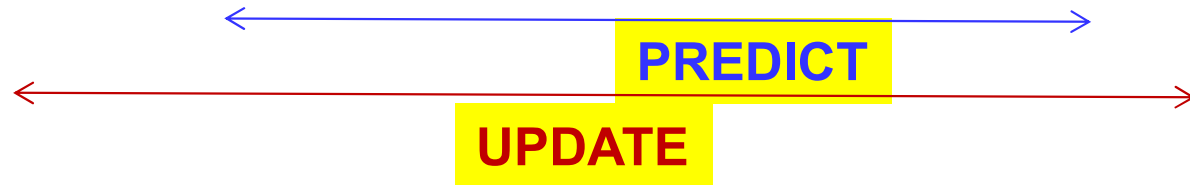
- At  $T=0$  the predicted state distribution is the initial state probability
- At each time  $T$ , the current estimate of the distribution over states considers *all* observations  $x_0 \dots x_T$ 
  - A natural outcome of the Markov nature of the model
- The prediction+update is identical to the forward computation for HMMs to within a normalizing constant

# Comparison to Forward Algorithm



- Forward Algorithm:

- $P(x_{0:T}, S_T) = P(x_T | S_T) \sum_{S_{T-1}} P(x_{0:T-1}, S_{T-1}) P(S_T | S_{T-1})$



- Normalized:

- $P(S_T | x_{0:T}) = (\sum_{S'_T} P(x_{0:T}, S'_T))^{-1} P(x_{0:T}, S_T) = C P(x_{0:T}, S_T)$

# Decomposing the Algorithm

$$P(S_t, X_{0:t}) = P(X_t|S_t) \sum_{S_{t-1}} P(S_t|S_{t-1})P(S_{t-1}, X_{0:t-1})$$



Predict:  $P(S_t|X_{0:t-1}) = \sum_{S_{t-1}} P(S_t|S_{t-1})P(S_{t-1}|X_{0:t-1})$

Update:  $P(S_t|X_{0:t}) = \frac{P(S_t|X_{0:t-1})P(X_t|S_t)}{\sum_S P(S|X_{0:t-1})P(X_t|S)}$

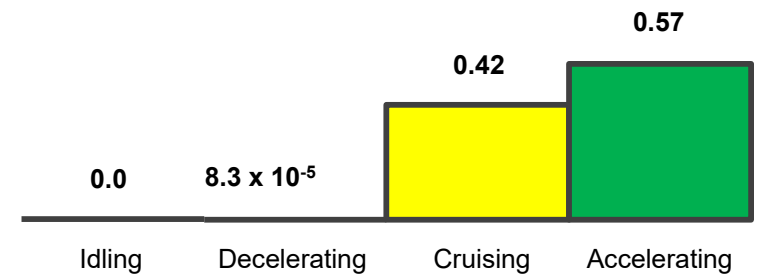




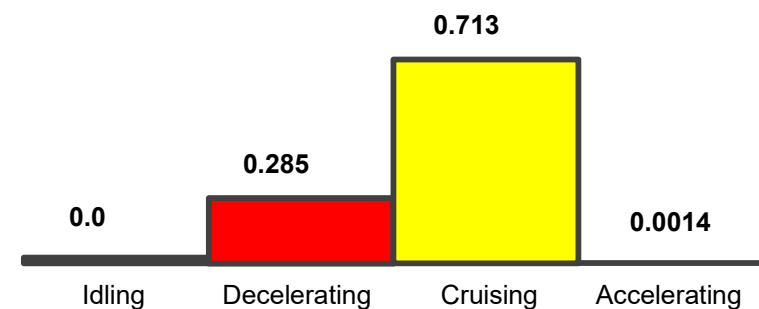
# Estimating a Unique state

- What we have estimated is a *distribution* over the states
- If we had to guess **a** state, we would pick the most likely state from the distributions

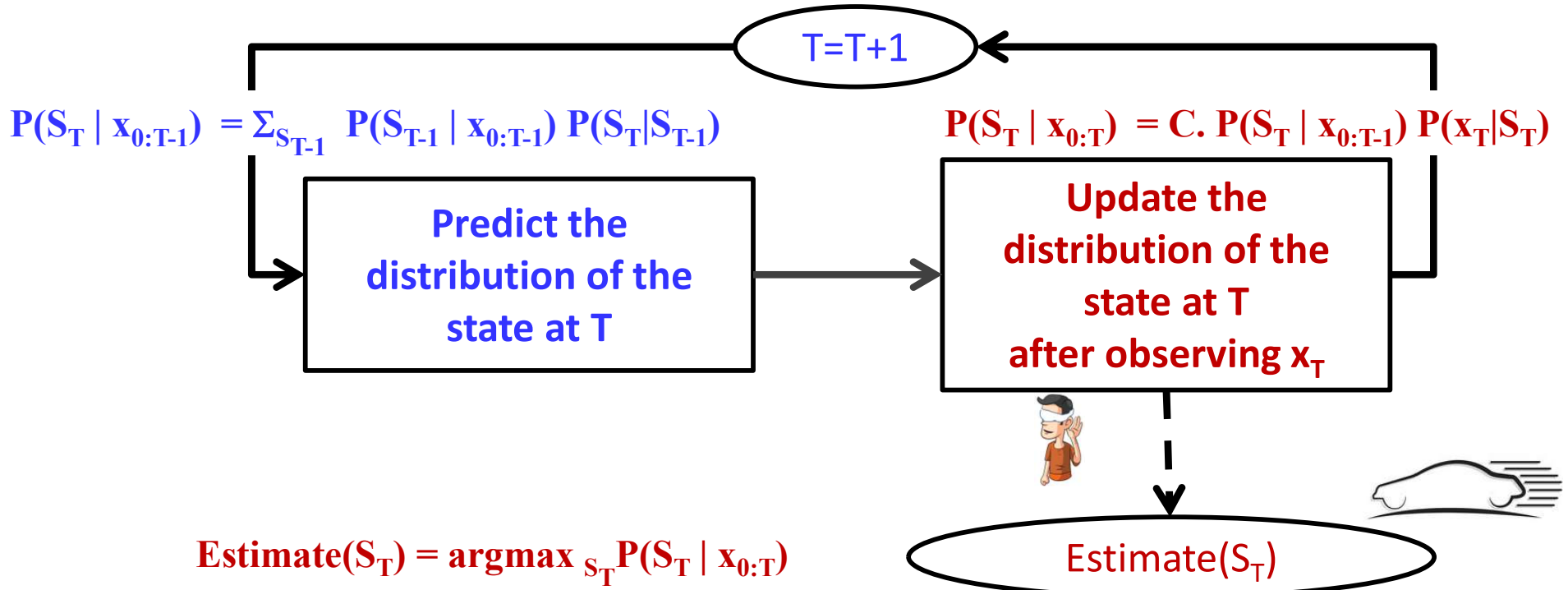
- State( $T=0$ ) = Accelerating



- State( $T=1$ ) = Cruising

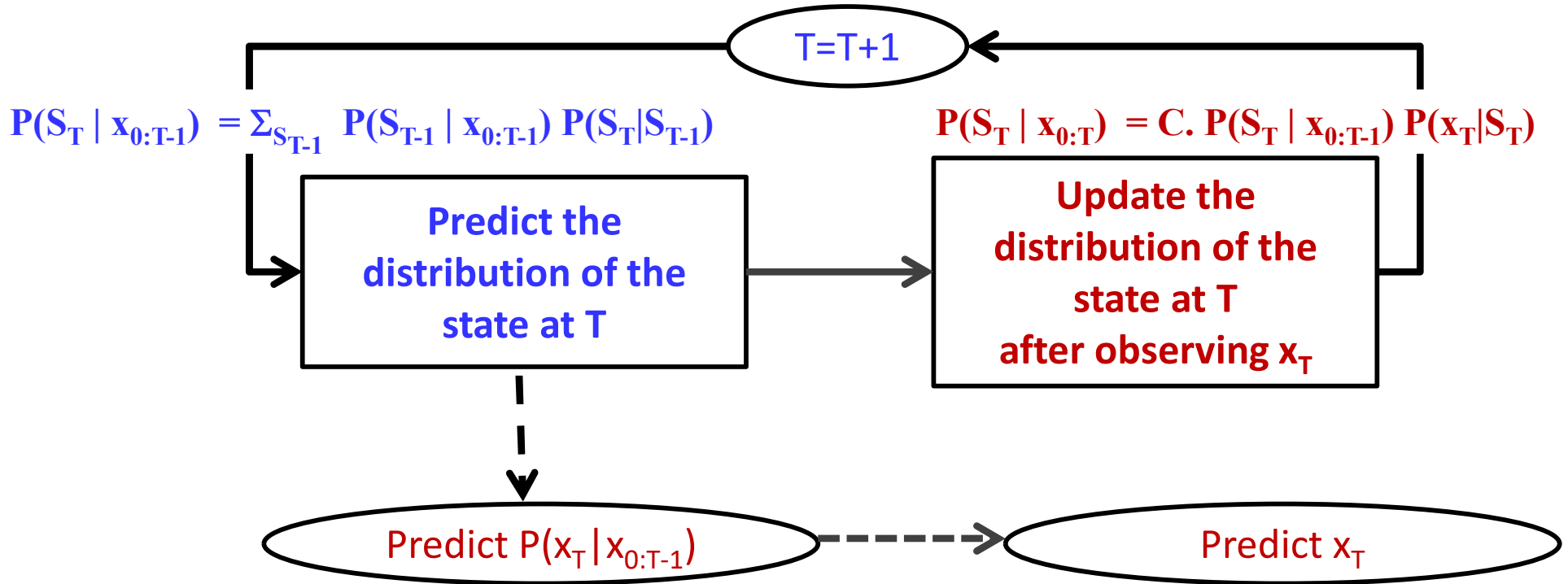


# Estimating the *state*



- The state is estimated from the updated distribution
  - The updated distribution is propagated into time, not the state

# Predicting the *next observation*



- The probability distribution for the observations at the next time is a mixture:
- $P(X_t | X_{0:t-1}) = \sum_{S_t} P(X_t | S_t) P(S_t | X_{0:t-1})$
- The actual observation can be predicted from  $P(x_T | x_{0:T-1})$

# Predicting the next observation

- Can use any of the various estimators of  $x_T$  from  $P(x_T|x_{0:T-1})$
- MAP estimate:
  - $\operatorname{argmax}_{x_T} P(x_T|x_{0:T-1})$
- MMSE estimate:
  - $\operatorname{Expectation}(x_T|x_{0:T-1})$

# Difference from Viterbi decoding

- Estimating only the *current* state at any time
  - Not the state sequence
  - Although we are considering all past observations
- The most likely state at  $T$  and  $T+1$  may be such that there is no valid transition between  $S_T$  and  $S_{T+1}$

# Poll 2

- To find your way back home...
  - At each time  $t$  you \*predict\* your beliefs about what your state will be at the *next time  $t+1$*  based on all you have observed until now (time  $t$ )
  - At each time  $t$ , you \*update\* your beliefs about the state at  $t$ , that you made when still at  $t-1$ , based on the latest observation  $O(t)$
  - At each time  $t$  you predict your belief at the state at  $t+1$ , and then update your belief after observing  $O(t+1)$
  - At each time you predict the distribution of the state at  $t+1$ , and then update your predicted distribution based on  $O(t+1)$
  - Your guess for the actual state must be derived from the estimated distribution for the state

# Poll 2

- To find your way back home...
  - At each time  $t$  you **\*predict\*** your beliefs about what your state will be at the *next time  $t+1$*  based on all you have observed until now (time  $t$ )
  - At each time  $t$ , you **\*update\*** your beliefs about the state at  $t$ , that you made when still at  $t-1$ , based on the latest observation  $O(t)$
  - At each time  $t$  you predict the actual state at  $t+1$ , and then update your guess for the state after observing  $O(t+1)$
  - **At each time you predict the distribution of the state at  $t+1$ , and then update your predicted distribution based on  $O(t+1)$**
  - **Your guess for the actual state must be derived from the estimated distribution for the state**

# *A continuous state model*

- HMM assumes a very coarsely quantized state space
  - Idling / accelerating / cruising / decelerating
- Actual state can be finer
  - Idling, accelerating at various rates, decelerating at various rates, cruising at various speeds
- Solution: Many more states (one for each acceleration /deceleration rate, cruising speed)?
- Solution: *A continuous* valued state



# Tracking and Prediction: The wind and the target

- Aim: measure wind velocity
- Using a noisy wind speed sensor
  - E.g. arrows shot at a target



- **State:** Wind speed at time  $t$  depends on speed at time  $t-1$

$$S_t = S_{t-1} + \epsilon_t$$



- **Observation:** Arrow position at time  $t$  depends on wind speed at time  $t$

$$Y_t = AS_t + \gamma_t$$



# The real-valued state model

- A state equation describing the dynamics of the system

$$s_t = f(s_{t-1}, \varepsilon_t)$$

- $s_t$  is the state of the system at time  $t$
  - $\varepsilon_t$  is a driving function, which is assumed to be random
- The state of the system at any time depends only on the state at the previous time instant and the driving term at the current time
- An observation equation relating state to observation

$$o_t = g(s_t, \gamma_t)$$

- $o_t$  is the observation at time  $t$
  - $\gamma_t$  is the noise affecting the observation (also random)
- The observation at any time depends only on the current state of the system and the noise

# States are still “hidden”



$$s_t = f(s_{t-1}, \epsilon_t)$$

$$o_t = g(s_t, \gamma_t)$$

- The state is a continuous valued parameter that is not directly seen
  - The state is the position of the automobile or the star
- The observations are dependent on the state and are the only way of knowing about the state
  - Sensor readings (for the automobile) or recorded image (for the telescope)

# Statistical Prediction and Estimation

- Given an *a priori* probability distribution for the state
  - $P_0(s)$ : Our belief in the state of the system before we observe any data
    - Probability of state of navlab
    - Probability of state of stars
- Given a sequence of observations  $o_0..o_t$
- Estimate state at time  $t$

# Prediction and update at $t = 0$

- Prediction
  - Initial probability distribution for state
  - $P(s_0) = P_0(s_0)$
- Update:
  - Then we observe  $o_0$
  - We must update our belief in the state

$$P(s_0 | o_0) = \frac{P(s_0)P(o_0 | s)}{P(o_0)} = \frac{P_0(s_0)P(o_0 | s_0)}{P(o_0)}$$

- $P(s_0 | o_0) = C.P_0(s_0)P(o_0 | s_0)$

# Prediction and update at $t = 0$

- Prediction
  - Initial probability distribution for state
  - $P(s_0) = P_0(s_0)$
- Update:
  - Then we observe  $o_0$
  - We must update our belief in the state

$$P(s_0 | o_0) = \frac{P(s_0)P(o_0 | s)}{P(o_0)} = \frac{P_0(s_0)P(o_0 | s_0)}{P(o_0)}$$

- $P(s_0 | o_0) = C.P_0(s_0)P(o_0 | s_0)$

# The observation probability: $P(o | s)$

- $o_t = g(s_t, \gamma_t)$ 
  - This is a (possibly many-to-one) stochastic function of state  $s_t$  and noise  $\gamma_t$
  - Noise  $\gamma_t$  is random. Assume it is the same dimensionality as  $o_t$
- Let  $P_\gamma(\gamma_t)$  be the probability distribution of  $\gamma_t$
- Let  $\{\gamma: g(s_t, \gamma) = o_t\}$  be all  $\gamma$  that result in  $o_t$

$$P(o_t | s_t) = \sum_{\gamma: g(s_t, \gamma) = o_t} \frac{P_\gamma(\gamma)}{|J_\gamma(g(s_t, \gamma))|}$$

# The observation probability

- $P(o|s) = ?$        $o_t = g(s_t, \gamma_t)$

$$P(o_t | s_t) = \sum_{\gamma: g(s_t, \gamma) = o_t} \frac{P_\gamma(\gamma)}{|J_\gamma(g(s_t, \gamma))|}$$

- The  $J$  is a Jacobian

$$|J_\gamma(g(s_t, \gamma))| = \begin{vmatrix} \frac{\partial o_t(1)}{\partial \gamma(1)} & \dots & \frac{\partial o_t(1)}{\partial \gamma(n)} \\ \vdots & \ddots & \vdots \\ \frac{\partial o_t(n)}{\partial \gamma(1)} & \dots & \frac{\partial o_t(n)}{\partial \gamma(n)} \end{vmatrix}$$

- For scalar functions of scalar variables, it is simply a derivative:

$$|J_\gamma(g(s_t, \gamma))| = \left| \frac{\partial o_t}{\partial \gamma} \right|$$



# Predicting the next state at t=1

- Given  $P(s_0 | o_0)$ , what is the probability of the state at t=1

$$P(s_1 | o_0) = \int_{\{s_0\}} P(s_1, s_0 | o_0) ds_0 = \int_{\{s_0\}} P(s_1 | s_0) P(s_0 | o_0) ds_0$$

- State progression function:

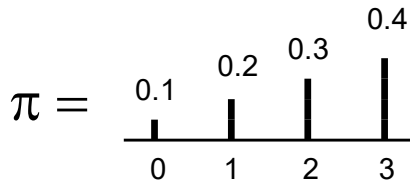
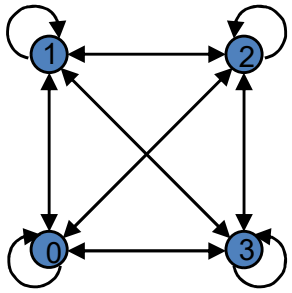
$$s_t = f(s_{t-1}, \varepsilon_t)$$

- $\varepsilon_t$  is a driving term with probability distribution  $P_\varepsilon(\varepsilon_t)$
- $P(s_t | s_{t-1})$  can be computed similarly to  $P(o | s)$ 
  - $P(s_1 | s_0)$  is an instance of this

# And moving on

- $P(s_1 | o_0)$  is the predicted state distribution for  $t=1$
- Then we observe  $o_1$ 
  - We must update the probability distribution for  $s_1$
  - $P(s_1 | o_{0:1}) = CP(s_1 | o_0)P(o_1 | s_1)$
- We can continue on

# Discrete vs. Continuous state systems



Prediction at time 0:

$$P(S_0) = \pi(S_0)$$

Update after  $O_0$ :

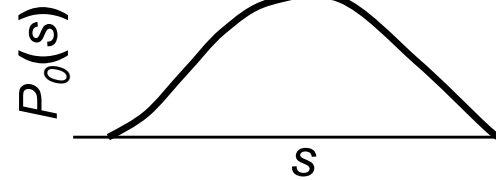
$$P(S_0|O_0) = C \cdot \pi(S_0)P(O_0|S_0)$$

Prediction at time 1:

$$P(S_1|O_0) = \sum_{S_0} P(S_0|O_0)P(S_1|S_0)$$

Update after  $O_1$ :

$$P(S_1|O_{0:1}) = C \cdot P(S_1|O_0)P(O_1|S_1)$$



$$s_t = f(s_{t-1}, \epsilon_t)$$

$$o_t = g(s_t, \gamma_t)$$

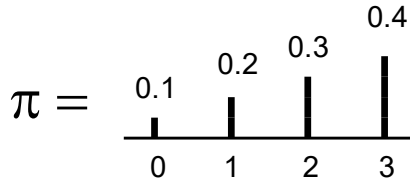
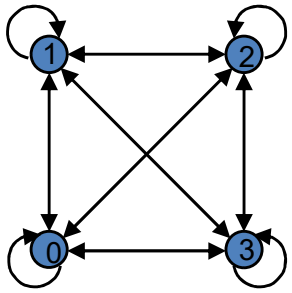
$$P(S_0) = P_0(S_0)$$

$$P(S_0|O_0) = C \cdot P(S_0)P(O_0|S_0)$$

$$P(S_1|O_0) = \int_{-\infty}^{\infty} P(S_0|O_0)P(S_1|S_0)dS_0$$

$$P(S_1|O_{0:1}) = C \cdot P(S_1|O_0)P(O_1|S_1)$$

# Discrete vs. Continuous State Systems



$$s_t = f(s_{t-1}, \varepsilon_t)$$

$$o_t = g(s_t, \gamma_t)$$

## Prediction at time t:

$$P(S_t | O_{0:t-1}) = \sum_{S_{t-1}} P(S_{t-1} | O_{0:t-1}) P(S_t | S_{t-1})$$

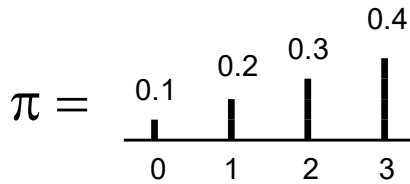
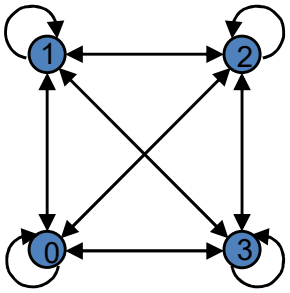
$$P(S_t | O_{0:t-1}) = \int_{-\infty}^{\infty} P(S_{t-1} | O_{0:t-1}) P(S_t | S_{t-1}) dS_{t-1}$$

## Update after observing $O_t$ :

$$P(S_t | O_{0:t}) = C \cdot P(S_t | O_{0:t-1}) P(O_t | S_t)$$

$$P(S_t | O_{0:t}) = C \cdot P(S_t | O_{0:t-1}) P(O_t | S_t)$$

# Discrete vs. Continuous State Systems



$$s_t = f(s_{t-1}, \epsilon_t)$$

$$o_t = g(s_t, \gamma_t)$$

## Parameters

Initial state prob.

$$\pi$$

$$P(s)$$

Transition prob

$$P(s_t = j | s_{t-1} = i)$$

$$P(s_t | s_{t-1})$$

Observation prob

$$P(O|s)$$

$$P(O|s)$$

# Special case: Linear Gaussian model



$$s_t = A_t s_{t-1} + \varepsilon_t$$

$$P(\varepsilon) = \frac{1}{\sqrt{(2\pi)^d |\Theta_\varepsilon|}} \exp\left(-0.5(\varepsilon - \mu_\varepsilon)^T \Theta_\varepsilon^{-1} (\varepsilon - \mu_\varepsilon)\right)$$



$$o_t = B_t s_t + \gamma_t$$

$$P(\gamma) = \frac{1}{\sqrt{(2\pi)^d |\Theta_\gamma|}} \exp\left(-0.5(\gamma - \mu_\gamma)^T \Theta_\gamma^{-1} (\gamma - \mu_\gamma)\right)$$

- A *linear* state dynamics equation
  - Probability of state driving term  $\varepsilon$  is Gaussian
  - Sometimes viewed as a driving term  $\mu_\varepsilon$  and additive zero-mean noise
- A *linear* observation equation
  - Probability of observation noise  $\gamma$  is Gaussian
- $A_t$ ,  $B_t$  and Gaussian parameters assumed known
  - May vary with time

# Linear model example

## The wind and the target



- **State:** Wind speed at time  $t$  depends on speed at time  $t-1$

$$S_t = S_{t-1} + \epsilon_t$$



- **Observation:** Arrow position at time  $t$  depends on wind speed at time  $t$

$$O_t = BS_t + \gamma_t$$



# Model Parameters:

## The initial state probability

$$P_0(s) = \frac{1}{\sqrt{(2\pi)^d |R|}} \exp\left(-0.5(s - \bar{s})R^{-1}(s - \bar{s})^T\right)$$

$$P_0(s) = \text{Gaussian}(s; \bar{s}, R)$$

- We also assume the *initial* state distribution to be Gaussian
  - Often assumed zero mean

$$s_t = A_t s_{t-1} + \varepsilon_t$$

$$o_t = B_t s_t + \gamma_t$$



# Model Parameters:

## The observation probability

$$o_t = B_t s_t + \gamma_t$$

$$P(\gamma) = \text{Gaussian}(\gamma; \mu_\gamma, \Theta_\gamma)$$

$$P(o_t | s_t) = \text{Gaussian}(o_t; \mu_\gamma + B_t s_t, \Theta_\gamma)$$

- The probability of the observation, given the state, is simply the probability of the noise, with the mean shifted
  - Since the only uncertainty is from the noise
- The new mean is the mean of the distribution of the noise + the value of the observation in the absence of noise

# Model Parameters:

## State transition probability

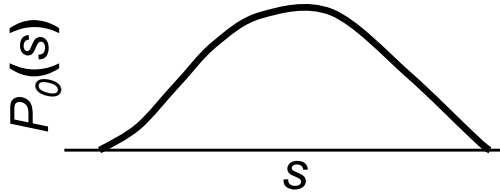
$$s_{t+1} = A_t s_t + \varepsilon_t$$

$$P(\varepsilon) = \text{Gaussian}(\varepsilon; \mu_\varepsilon, \Theta_\varepsilon)$$

$$P(s_{t+1} | s_t) = \text{Gaussian}(s_t; \mu_\varepsilon + A_t s_t, \Theta_\varepsilon)$$

- The probability of the state at time t, given the state at t-1, is simply the probability of the driving term, with the mean shifted

# Continuous state systems



$$s_{t+1} = A_t s_t + \varepsilon_t$$

$$o_t = B_t s_t + \gamma_t$$

Prediction at time 0:

$$P(S_0) = P_0(S_0)$$

---

Update after  $O_0$ :

$$P(S_0|O_0) = C \cdot P(S_0)P(O_0|S_0)$$

---

Prediction at time 1:

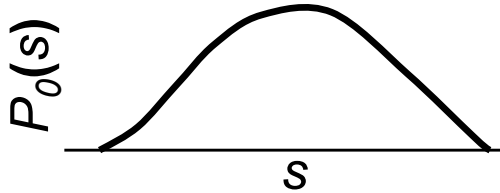
$$P(S_1|O_0) = \int_{-\infty}^{\infty} P(S_0|O_0)P(S_1|S_0)dS_0$$

---

Update after  $O_1$ :

$$P(S_1|O_{0:1}) = C \cdot P(S_1|O_0)P(O_1|S_1)$$

# Continuous state systems



$$s_{t+1} = A_t s_t + \varepsilon_t$$

$$o_t = B_t s_t + \gamma_t$$

Prediction at time 0:

$$P(S_0) = P_0(S_0)$$

---

Update after  $O_0$ :

$$P(S_0|O_0) = C \cdot P(S_0)P(O_0|S_0)$$

---

Prediction at time 1:

$$P(S_1|O_0) = \int_{-\infty}^{\infty} P(S_0|O_0)P(S_1|S_0)dS_0$$

---

Update after  $O_1$ :

$$P(S_1|O_{0:1}) = C \cdot P(S_1|O_0)P(O_1|S_1)$$

# Model Parameters:

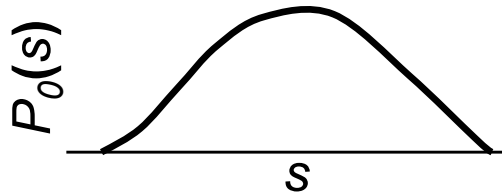
## The initial state probability

$$P_0(s) = \frac{1}{\sqrt{(2\pi)^d |R_0|}} \exp\left(-0.5(s - \bar{s}_0)R_0^{-1}(s - \bar{s}_0)^T\right)$$

$$P_0(s) = \text{Gaussian}(s; \bar{s}_0, R_0)$$

- We assume the *initial* state distribution to be Gaussian
  - Often assumed zero mean

# Continuous state systems



$$s_{t+1} = A_t s_t + \varepsilon_t$$

$$o_t = B_t s_t + \gamma_t$$

*a priori* probability  
distribution of state  $s$

Prediction at time 0:

$$P(S_0) = P_0(S_0)$$

$$= N(\bar{s}_0, R_0)$$

Update after  $O_0$ :

$$P(S_0|O_0) = C \cdot P(S_0)P(O_0|S_0)$$

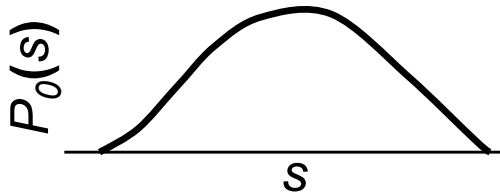
Prediction at time 1:

$$P(S_1|O_0) = \int_{-\infty}^{\infty} P(S_0|O_0)P(S_1|S_0)dS_0$$

Update after  $O_1$ :

$$P(S_1|O_{0:1}) = C \cdot P(S_1|O_0)P(O_1|S_1)$$

# Continuous state systems



$$s_{t+1} = A_t s_t + \varepsilon_t$$

$$o_t = B_t s_t + \gamma_t$$

Prediction at time 0:

$$P(S_0) = N(\bar{s}_0, R_0)$$

---

Update after  $O_0$ :

$$P(S_0|O_0) = C \cdot P(S_0)P(O_0|S_0)$$

---

Prediction at time 1:

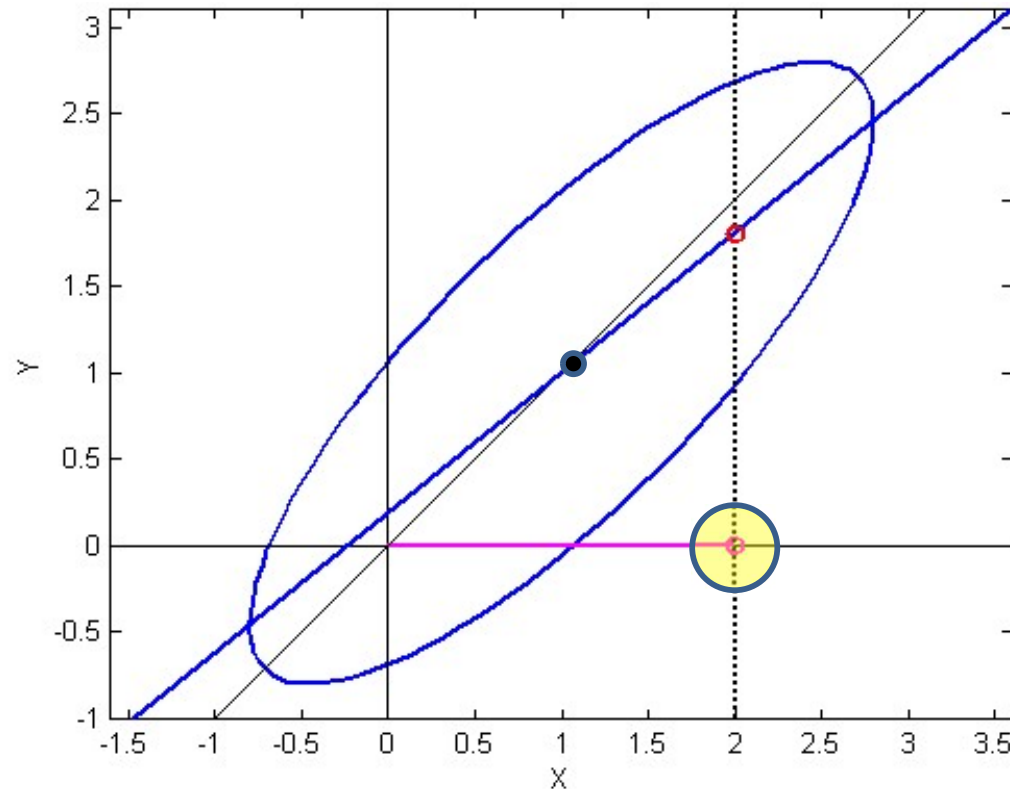
$$P(S_1|O_0) = \int_{-\infty}^{\infty} P(S_0|O_0)P(S_1|S_0)dS_0$$

---

Update after  $O_1$ :

$$P(S_1|O_{0:1}) = C \cdot P(S_1|O_0)P(O_1|S_1)$$

# Recap: Conditional of S given O: $P(S|O)$ for Gaussian RVs

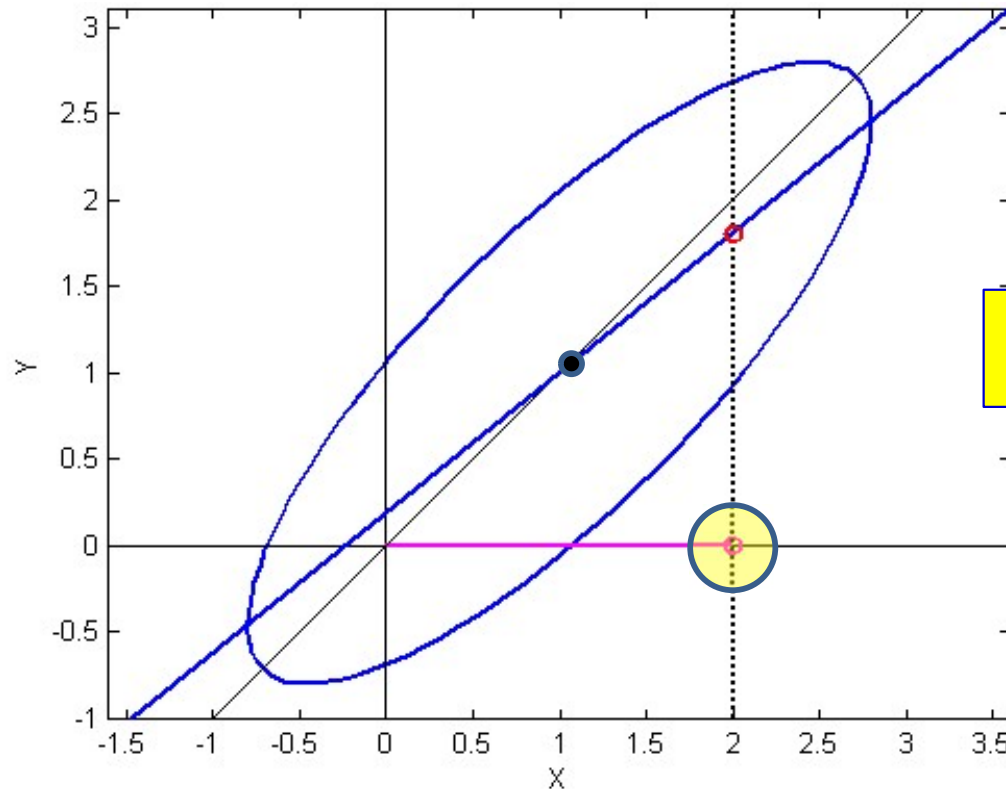


$$O = BS + \gamma$$

$$P(S|O) = N(\mu_S + \Theta_{SO}\Theta_O^{-1}(O - \mu_O), \Theta_S - \Theta_{SO}\Theta_O^{-1}\Theta_{OS})$$



# Recap: Conditional of S given O: $P(S|O)$ for Gaussian RVs



$$O = BS + \gamma$$

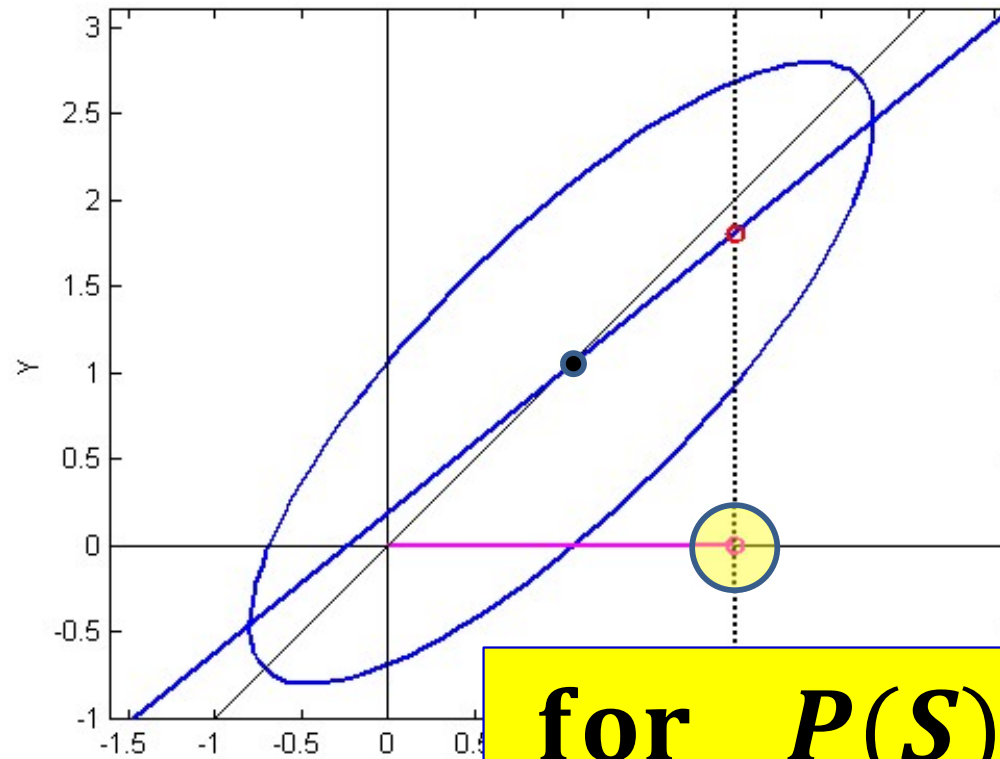
$$\Theta_{SO} = \Theta_S B^T$$

$$\Theta_O = B\Theta_S B^T + \Theta_\gamma$$

$$P(S|O) = N(\mu_S + \Theta_{SO}\Theta_O^{-1}(O - \mu_O), \Theta_S - \Theta_{SO}\Theta_O^{-1}\Theta_{OS})$$

$$P(S|O) = N(\mu_S + \Theta_S B^T (B\Theta_S B^T + \Theta_\gamma)^{-1} (O - B\mu_S - \mu_\gamma), \Theta_S - \Theta_S B^T (B\Theta_S B^T + \Theta_\gamma)^{-1} B\Theta_S)$$

# Recap: Conditional of S given O: $P(S|O)$ for Gaussian RVs

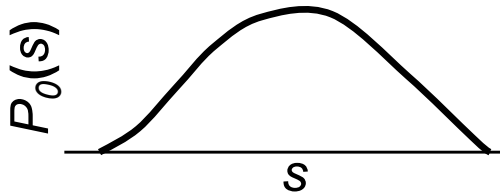


$$O = BS + \varepsilon$$

$$\text{for } P(S) = N(\bar{s}_0, R_0)$$

$$P(S_0|O_0) = N(\bar{s}_0 + R_0 B^T (BR_0 B^T + \Theta_\gamma)^{-1} (O_0 - B\bar{s}_0 - \mu_\gamma), \\ R_0 - R_0 B^T (BR_0 B^T + \Theta_\gamma)^{-1} BR_0)$$

# Continuous state systems



$$s_{t+1} = A_t s_t + \varepsilon_t$$

$$o_t = B_t s_t + \gamma_t$$

Prediction at time 0:

$$P(S_0) = N(\bar{s}_0, R_0)$$

---

Update after  $O_0$ :

$$P(S_0|O_0) = C \cdot P(S_0)P(O_0|S_0)$$

$$P(S_0|O_0) = N(\hat{s}_0, \hat{R}_0)$$

---

Prediction at time 1:

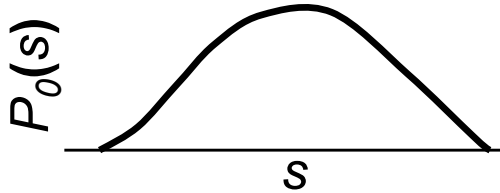
$$P(S_1|O_0) = \int_{-\infty}^{\infty} P(S_0|O_0)P(S_1|S_0)dS_0$$

---

Update after  $O_1$ :

$$P(S_1|O_{0:1}) = C \cdot P(S_1|O_0)P(O_1|S_1)$$

# Continuous state systems



$$s_{t+1} = A_t s_t + \varepsilon_t$$

$$o_t = B_t s_t + \gamma_t$$

Prediction at time 0:

$$P(S_0) = N(\bar{s}_0, R_0)$$

Update after  $O_0$ :

$$P(S_0|O_0) = N(\hat{s}_0, \hat{R}_0)$$

$$K_0 = R_0 B^T (B R_0 B^T + \Theta_\gamma)^{-1}$$

$$\hat{s}_0 = \bar{s}_0 + K_0 (O_0 - B \bar{s}_0 - \mu_\gamma)$$

$$\hat{R}_0 = (I - K_0) R_0$$

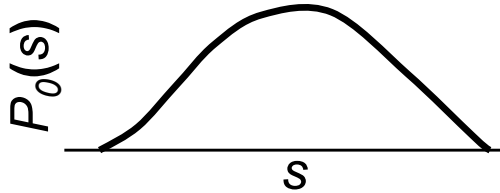
Prediction at time 1:

$$P(S_1|O_0) = \int_{-\infty}^{\infty} P(S_0|O_0) P(S_1|S_0) dS_0$$

Update after  $O_1$ :

$$P(S_1|O_{0:1}) = C \cdot P(S_1|O_0) P(O_1|S_1)$$

# Continuous state systems



$$s_{t+1} = A_t s_t + \varepsilon_t$$

$$o_t = B_t s_t + \gamma_t$$

Prediction at time 0:

$$P(S_0) = N(\bar{s}_0, R_0)$$

Update after  $O_0$ :

$$P(S_0|O_0) = C \cdot P(S_0)P(O_0|S_0)$$

$$= N(\bar{s}_0 + R_0 B^T (B R_0 B^T + \Theta_\gamma)^{-1} (O_0 - B \bar{s}_0 - \mu_\gamma), \\ R_0 - R_0 B^T (B R_0 B^T + \Theta_\gamma)^{-1} B R_0)$$

Prediction at time 1:

$$P(S_1|O_0) = \int_{-\infty}^{\infty} P(S_0|O_0)P(S_1|S_0)dS_0$$

Update after  $O_1$ :

$$P(S_1|O_{0:1}) = C \cdot P(S_1|O_0)P(O_1|S_1)$$

# Introducing shorthand notation

$$P(S_0|O_0) = N(\bar{s}_0 + R_0 B^T (B R_0 B^T + \Theta_\gamma)^{-1} (O_0 - B \bar{s}_0 - \mu_\gamma), \\ R_0 - R_0 B^T (B R_0 B^T + \Theta_\gamma)^{-1} B R_0)$$

$$\hat{s}_0 = \bar{s}_0 + R_0 B^T (B R_0 B^T + \Theta_\gamma)^{-1} (O - B \bar{s}_0 - \mu_\gamma)$$

$$\hat{R}_0 = R_0 - R_0 B^T (B R_0 B^T + \Theta_\gamma)^{-1} B R_0$$

$$P(S_0|O_0) = N(\hat{s}_0, \hat{R}_0)$$

# Introducing shorthand notation

$$P(S_0|O_0) = N(\bar{s}_0 + R_0 B^T (B R_0 B^T + \Theta_\gamma)^{-1} (O_0 - B \bar{s}_0 - \mu_\gamma), \\ R_0 - R_0 B^T (B R_0 B^T + \Theta_\gamma)^{-1} B R_0)$$

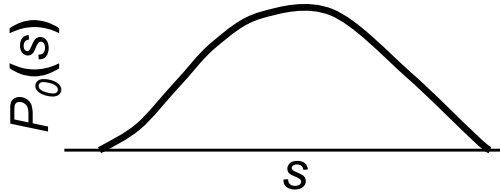
$$K_0 = R_0 B^T (B R_0 B^T + \Theta_\gamma)^{-1}$$

$$\hat{s}_0 = \bar{s}_0 + K_0 (O - B \bar{s}_0 - \mu_\gamma)$$

$$\hat{R}_0 = (I - K_0 B) R_0$$

$$P(S_0|O_0) = N(\hat{s}_0, \hat{R}_0)$$

# Continuous state systems



$$s_{t+1} = A_t s_t + \varepsilon_t$$

$$o_t = B_t s_t + \gamma_t$$

Prediction at time 0:

$$P(S_0) = N(\bar{s}_0, R_0)$$

Update after  $O_0$ :

$$P(S_0|O_0) = N(\hat{s}_0, \hat{R}_0)$$

$$K_0 = R_0 B^T (B R_0 B^T + \Theta_\gamma)^{-1}$$

$$\hat{s}_0 = \bar{s}_0 + K_0 (O_0 - B \bar{s}_0 - \mu_\gamma)$$

$$\hat{R}_0 = (I - K_0) R_0$$

Prediction at time 1:

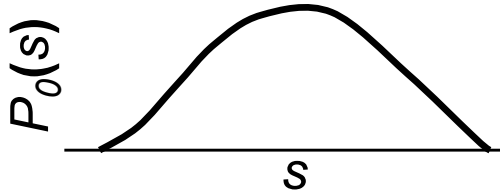
$$P(S_1|O_0) = \int_{-\infty}^{\infty} P(S_0|O_0) P(S_1|S_0) dS_0$$

Update after  $O_1$ :

$$P(S_1|O_{0:1}) = C \cdot P(S_1|O_0) P(O_1|S_1)$$



# Continuous state systems



$$s_{t+1} = A_t s_t + \varepsilon_t$$

$$o_t = B_t s_t + \gamma_t$$

Prediction at time 0:

$$P(S_0) = N(\bar{s}_0, R_0)$$

Update after  $O_0$ :

$$P(S_0|O_0) = N(\hat{s}_0, \hat{R}_0)$$

$$K_0 = R_0 B^T (B R_0 B^T + \Theta_\gamma)^{-1}$$

$$\hat{s}_0 = \bar{s}_0 + K_0 (O_0 - B \bar{s}_0 - \mu_\gamma)$$

$$\hat{R}_0 = (I - K_0 B) R_0$$

Prediction at time 1:

$$P(S_1|O_0) = \int_{-\infty}^{\infty} P(S_0|O_0) P(S_1|S_0) dS_0$$

Update after  $O_1$ :

$$P(S_1|O_{0:1}) = C \cdot P(S_1|O_0) P(O_1|S_1)$$

# The prediction equation

$$P(S_1|O_0) = \int_{-\infty}^{\infty} P(S_0|O_0)P(S_1|S_0)dS_0$$

$$P(S_0|O_0) = N(\hat{s}_0, \hat{R}_0)$$

$$P(S_1|S_0) = N(AS_0 + \mu_\varepsilon, \Theta_\varepsilon)$$

$$P(\varepsilon) = N(\mu_\varepsilon, \Theta_\varepsilon)$$

$$s_{t+1} = A_t s_t + \varepsilon_t$$

- The integral of the product of two Gaussians

$$P(S_1|O_0) = \int_{-\infty}^{\infty} \text{Gaussian}(S_0; \hat{s}_0, \hat{R}_0) \text{Gaussian}(S_1; AS_0, \Theta_\varepsilon) dS_0$$

# The Prediction Equation

- The integral of the product of two Gaussians is Gaussian!

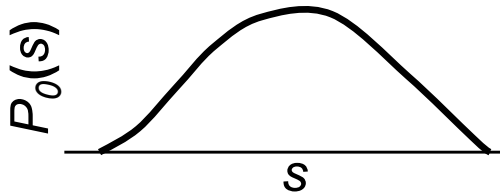
$$P(S_1|O_0) = \int_{-\infty}^{\infty} \text{Gaussian}(S_0; \hat{s}_0, \hat{R}_0) \text{Gaussian}(S_1; AS_0 + \mu_\varepsilon, \Theta_\varepsilon) dS_0$$

$$= \int_{-\infty}^{\infty} c_1 \exp(-0.5(S_0 - \hat{s}_0)\hat{R}_0^{-1}(S_0 - \hat{s}_0)^T) \cdot c_2 \exp(-0.5(S_1 - AS_0 - \mu_\varepsilon)\Theta_\varepsilon^{-1}(S_1 - AS_0 - \mu_\varepsilon)^T) dS_0$$

$$= \text{Gaussian}(S_1; A\hat{s}_0 + \mu_\varepsilon, \Theta_\varepsilon + A\hat{R}_0A^T)$$

$$P(S_1|O_0) = N(A\hat{s}_0 + \mu_\varepsilon, \Theta_\varepsilon + A\hat{R}_0A^T)$$

# Continuous state systems



$$s_{t+1} = A_t s_t + \varepsilon_t$$

$$o_t = B_t s_t + \gamma_t$$

Prediction at time 0:

$$P(S_0) = N(\bar{s}_0, R_0)$$

Update after  $O_0$ :

$$P(S_0|O_0) = N(\hat{s}_0, \hat{R}_0)$$

$$K_0 = R_0 B^T (B R_0 B^T + \Theta_\gamma)^{-1}$$

$$\hat{s}_0 = \bar{s}_0 + K_0 (O_0 - B \bar{s}_0 - \mu_\gamma)$$

$$\hat{R}_0 = (I - K_0 B) R_0$$

Prediction at time 1:

$$P(S_1|O_0) = \int_{-\infty}^{\infty} P(S_0|O_0) P(S_1|S_0) dS_0$$

$$= N(A \hat{s}_0 + \mu_\varepsilon, \Theta_\varepsilon + A \hat{R}_0 A^T)$$

Update after  $O_1$ :

$$P(S_1|O_{0:1}) = C \cdot P(S_1|O_0) P(O_1|S_1)$$

# More shorthand notation

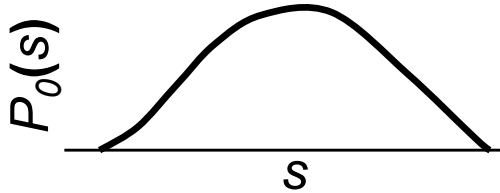
$$P(S_1|O_0) = N(A\hat{S}_0 + \mu_\varepsilon, \Theta_\varepsilon + A\hat{R}_0A^T)$$

$$\bar{S}_1 = A\hat{S}_0 + \mu_\varepsilon$$

$$R_1 = \Theta_\varepsilon + A\hat{R}_0A^T$$

$$P(S_1|O_0) = N(\bar{S}_1, R_1)$$

# Continuous state systems



$$s_{t+1} = A_t s_t + \varepsilon_t$$

$$o_t = B_t s_t + \gamma_t$$

Prediction at time 0:

$$P(S_0) = N(\bar{s}_0, R_0)$$

Update after  $O_0$ :

$$P(S_0|O_0) = N(\hat{s}_0, \hat{R}_0)$$

$$K_0 = R_0 B^T (B R_0 B^T + \Theta_\gamma)^{-1}$$

$$\hat{s}_0 = \bar{s}_0 + K_0 (O_0 - B \bar{s}_0 - \mu_\gamma)$$

$$\hat{R}_0 = (I - K_0 B) R_0$$

Prediction at time 1:

$$P(S_1|O_0) = N(\bar{s}_1, R_1)$$

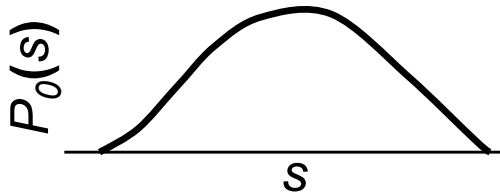
$$\bar{s}_1 = A \hat{s}_0 + \mu_\varepsilon$$

$$R_1 = \Theta_\varepsilon + A \hat{R}_0 A^T$$

Update after  $O_1$ :

$$P(S_1|O_{0:1}) = C \cdot P(S_1|O_0) P(O_1|S_1)$$

# Continuous state systems



$$s_{t+1} = A_t s_t + \varepsilon_t$$

$$o_t = B_t s_t + \gamma_t$$

Prediction at time 0:

$$P(S_0) = N(\bar{s}_0, R_0)$$

Update after  $O_0$ :

$$P(S_0|O_0) = N(\hat{s}_0, \hat{R}_0)$$

$$K_0 = R_0 B^T (B R_0 B^T + \Theta_\gamma)^{-1}$$

$$\hat{s}_0 = \bar{s}_0 + K_0 (O_0 - B \bar{s}_0 - \mu_\gamma)$$

$$\hat{R}_0 = (I - K_0 B) R_0$$

Prediction at time 1:

$$P(S_1|O_0) = N(\bar{s}_1, R_1)$$

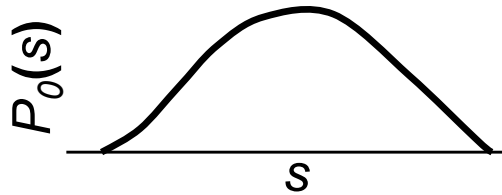
$$\bar{s}_1 = A \hat{s}_0 + \mu_\varepsilon$$

$$R_1 = \Theta_\varepsilon + A \hat{R}_0 A^T$$

Update after  $O_1$ :

$$P(S_1|O_{0:1}) = C \cdot P(S_1|O_0) P(O_1|S_1)$$

# Continuous state systems



$$s_{t+1} = A_t s_t + \varepsilon_t$$

$$o_t = B_t s_t + \gamma_t$$

Prediction at time 0:

$$P(S_0) = N(\bar{s}_0, R_0)$$

Update after  $O_0$ :

$$P(S_0|O_0) = N(\hat{s}_0, \hat{R}_0)$$

$$K_0 = R_0 B^T (B R_0 B^T + \Theta_\gamma)^{-1}$$

$$\hat{s}_0 = \bar{s}_0 + K_0 (O_0 - B \bar{s}_0 - \mu_\gamma)$$

$$\hat{R}_0 = (I - K_0 B) R_0$$

Prediction at time 1:

$$P(S_1|O_0) = N(\bar{s}_1, R_1)$$

$$\bar{s}_1 = A \hat{s}_0 + \mu_\varepsilon$$

$$R_1 = \Theta_\varepsilon + A \hat{R}_0 A^T$$

Update after  $O_1$ :

$$P(S_1|O_{0:1}) = C \cdot P(S_1|O_0) P(O_1|S_1) = N(\hat{s}_1, \hat{R}_1)$$

$$K_1 = R_1 B^T (B R_1 B^T + \Theta_\gamma)^{-1}$$

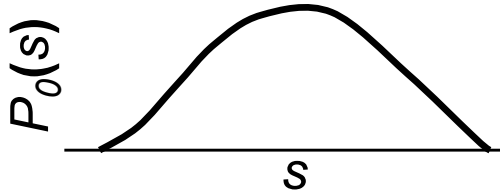
$$\hat{s}_1 = \bar{s}_1 + K_1 (O_1 - B \bar{s}_1 - \mu_\gamma)$$

$$\hat{R}_1 = (I - K_1 B) R_1$$





# Continuous state systems



$$s_{t+1} = A_t s_t + \varepsilon_t$$

$$o_t = B_t s_t + \gamma_t$$

Prediction at time 0:

$$P(S_0) = N(\bar{s}_0, R_0)$$

Update after  $O_0$ :

$$P(S_0|O_0) = N(\hat{s}_0, \hat{R}_0)$$

$$K_0 = R_0 B^T (B R_0 B^T + \Theta_\gamma)^{-1}$$

$$\hat{s}_0 = \bar{s}_0 + K_0 (O_0 - B \bar{s}_0 - \mu_\gamma)$$

$$\hat{R}_0 = (I - K_0 B) R_0$$

Prediction at time 1:

$$P(S_1|O_0) = N(\bar{s}_1, R_1)$$

$$\bar{s}_1 = A \hat{s}_0 + \mu_\varepsilon$$

$$R_1 = \Theta_\varepsilon + A \hat{R}_0 A^T$$

Update after  $O_1$ :

$$P(S_1|O_{0:1}) = N(\hat{s}_1, \hat{R}_1)$$

$$K_1 = R_1 B^T (B R_1 B^T + \Theta_\gamma)^{-1}$$

$$\hat{s}_1 = \bar{s}_1 + K_1 (O_1 - B \bar{s}_1 - \mu_\gamma)$$

$$\hat{R}_1 = (I - K_1 B) R_1$$

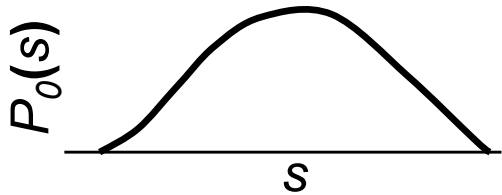
# Poll 3

- Tracking state with a continuous-state system is strictly analogous to doing so with an HMM
  - True
  - False
- When the state and observation relations are given by equations between continuous variables, rather than probabilistic dependencies, state estimation becomes a deterministic procedure
  - True
  - False
- In a linear Gaussian model, where the initial state distribution is Gaussian and state and observation equations are affine, the predicted and updated state probability distributions are:
  - Always Gaussian
  - Predicted distributions are Gaussian, but updated distributions may not be
  - Neither is assured to be Gaussian

# Poll 3

- Tracking state with a continuous-state system is strictly analogous to doing so with an HMM
  - **True**
  - False
- When the state and observation relations are given by equations between continuous variables, rather than probabilistic dependencies, state estimation becomes a deterministic procedure
  - True
  - **False**
- In a linear Gaussian model, where the initial state distribution is Gaussian and state and observation equations are affine, the predicted and updated state probability distributions are:
  - **Always Gaussian**
  - Predicted distributions are Gaussian, but updated distributions may not be
  - Neither is assured to be Gaussian

# Gaussian Continuous State Linear Systems



$$s_{t+1} = A_t s_t + \varepsilon_t$$

$$o_t = B_t s_t + \gamma_t$$



**Prediction at time t:**

$$P(S_t | O_{0:t-1}) = \int_{-\infty}^{\infty} P(S_{t-1} | O_{0:t-1}) P(S_t | S_{t-1}) dS_{t-1}$$

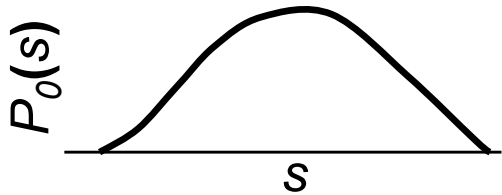


**Update after observing  $O_t$ :**

$$P(S_t | O_{0:t}) = C \cdot P(S_t | O_{0:t-1}) P(O_t | S_t)$$



# Gaussian Continuous State Linear Systems



$$s_{t+1} = A_t s_t + \varepsilon_t$$

$$o_t = B_t s_t + \gamma_t$$



## Prediction at time t:

$$P(S_t | O_{0:t-1}) = N(\bar{s}_t, R_t)$$

$$\bar{s}_t = A \hat{s}_{t-1} + \mu_\varepsilon$$

$$R_t = \Theta_\varepsilon + A \hat{R}_{t-1} A^T$$

## Update after observing $O_t$ :

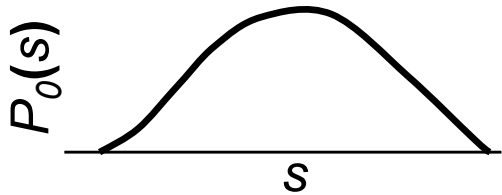
$$P(S_t | O_{0:t}) = N(\hat{s}_t, \hat{R}_t)$$

$$K_t = R_t B^T (B R_t B^T + \Theta_\gamma)^{-1}$$

$$\hat{s}_t = \bar{s}_t + K_t (O_t - B \bar{s}_t - \mu_\gamma)$$

$$\hat{R}_t = (I - K_t B) R_t$$

# Gaussian Continuous State Linear Systems



$$s_{t+1} = A_t s_t + \varepsilon_t$$

$$o_t = B_t s_t + \gamma_t$$



**Prediction at time t:**

$$P(S_t | O_{0:t-1}) = N(\bar{s}_t, R_t)$$

**Update after observing  $O_t$ :**

$$P(S_t | O_{0:t}) = N(\hat{s}_t, \hat{R}_t)$$

## KALMAN FILTER

$$\bar{s}_t = A \hat{s}_{t-1} + \mu_\varepsilon$$

$$R_t = \Theta_\varepsilon + A \hat{R}_{t-1} A^T$$

$$K_t = R_t B^T (B R_t B^T + \Theta_\gamma)^{-1}$$

$$\hat{s}_t = \bar{s}_t + K_t (O_t - B \bar{s}_t - \mu_\gamma)$$

$$\hat{R}_t = (I - K_t B) R_t$$

# The Kalman filter

- Prediction (based on state equation)

$$\bar{s}_t = A_t \hat{s}_{t-1} + \mu_\varepsilon$$

$$s_t = A_t s_{t-1} + \varepsilon_t$$

$$R_t = \Theta_\varepsilon + A_t \hat{R}_{t-1} A_t^T$$

- Update (using observation and observation equation)

$$K_t = R_t B_t^T (B_t R_t B_t^T + \Theta_\gamma)^{-1}$$

$$o_t = B_t s_t + \gamma_t$$

$$\hat{s}_t = \bar{s}_t + K_t (o_t - B_t \bar{s}_t - \mu_\gamma)$$

$$\hat{R}_t = (I - K_t B_t) R_t$$

# Explaining the Kalman Filter

- Prediction

$$s_t = A_t s_{t-1} + \varepsilon_t$$

$$\bar{s}_t = A_t \hat{s}_{t-1} + \mu_\varepsilon$$

$$o_t = B_t s_t + \gamma_t$$

$$R_t = \Theta_\varepsilon + A_t \hat{R}_{t-1} A_t^T$$

- The Kalman filter can be explained intuitively without working through the math

$$\hat{s}_t = \bar{s}_t + K_t (o_t - B_t \bar{s}_t - \mu_\gamma)$$

$$\hat{R}_t = (I - K_t B_t) R_t$$

NEXT CLASS!