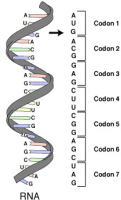
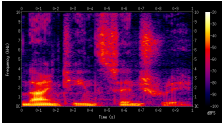


Sequence Modelling

Rich Turner and José Miguel Hernández-Lobato

Sequence data

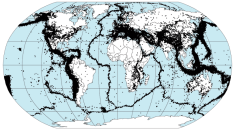


RNA
Ribonucleic acid

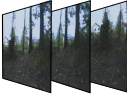
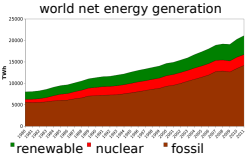
- Codon 1
A U G
- Codon 2
A C G
- Codon 3
A G G
- Codon 4
U C C
- Codon 5
U C G
- Codon 6
A G C
- Codon 7
U A G

*Good King Wenceslas looked out,
On the Feast of Stephen,
When the snow lay round about,
Deep and crisp and even;
Brightly shone the moon that night,
Though the frost was cruel,
When a poor man came in sight,
Gathering winter fuel.*

Preliminary Determination of Epicenters
358,214 Events, 1963 - 1998



Some images taken from wikipedia



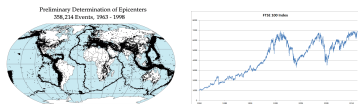
I believe that at the end of the century the use of words and general educated opinion will have altered so much that one will be able to speak of machines thinking without expecting to be contradicted.

A. Turing

Goals of sequence modelling

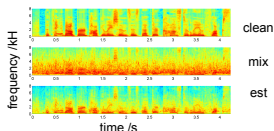
Predict future items in sequence

$$p(y_t | y_1, \dots, y_{t-1})$$



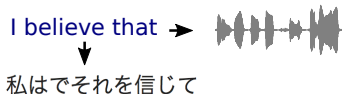
Remove noise from a sequence

$$p(y'_1, \dots, y'_t | y_1, \dots, y_t)$$



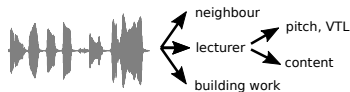
Predict one sequence from another

$$p(y'_1, \dots, y'_t | y_1, \dots, y_t)$$



Discover underlying latent variables

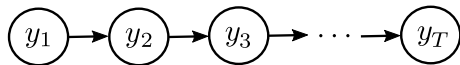
$$p(x_1, \dots, x_t | y_1, \dots, y_t)$$



Markov models

First order Markov

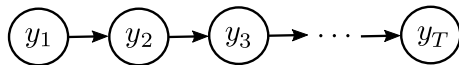
$$p(y_1, y_2, y_3, \dots, y_T) = p(y_1)p(y_2|y_1)p(y_3|y_2) \dots p(y_T|y_{T-1})$$



Markov models

First order Markov

$$p(y_1, y_2, y_3, \dots, y_T) = p(y_1)p(y_2|y_1)p(y_3|y_2) \dots p(y_T|y_{T-1})$$



parameters tied

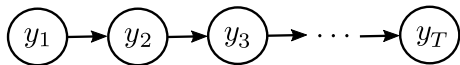


∞ number of variables
finite number of
parameters

Markov models

First order Markov

$$p(y_1, y_2, y_3, \dots, y_T) = p(y_1)p(y_2|y_1)p(y_3|y_2) \dots p(y_T|y_{T-1})$$



parameters tied \leftarrow ∞ number of variables
finite number of parameters

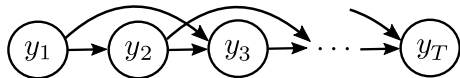
Markov model = conditional independence relationship + product rule

future $\rightarrow y_{t+1} \perp y_{1:t-1} | y_t$ \leftarrow independent of past
 \leftarrow given present

$$p(y_{1:T}) = \prod_{t=1}^T p(y_t | y_{1:t-1})$$

Second order Markov

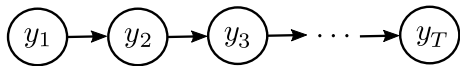
$$p(y_1, y_2, y_3, \dots, y_T) = p(y_1)p(y_2|y_1)p(y_3|y_2, y_1) \dots p(y_T|y_{T-1}, y_{T-2})$$



Markov models for discrete data: n-gram models

First order Markov (bi-gram)

$$p(y_1, y_2, y_3, \dots, y_T) = p(y_1)p(y_2|y_1)p(y_3|y_2) \dots p(y_T|y_{T-1})$$



$$y_t \in \{1, \dots, K\}$$

discrete states

$$p(y_1 = k) = \pi_k^0$$

initial state probabilities

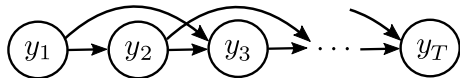
$$p(y_t = k | y_{t-1} = l) = T_{k,l}$$

transition probabilities
(stochastic matrix)

$$\sum_{k=1}^K T_{k,l} = 1$$

Second order Markov (tri-gram)

$$p(y_1, y_2, y_3, \dots, y_T) = p(y_1)p(y_2|y_1)p(y_3|y_2, y_1) \dots p(y_T|y_{T-1}, y_{T-2})$$



$$p(y_t = k | y_{t-1} = l, y_{t-2} = m) = T_{k,l,m}$$

n-grams require large
multidimensional arrays

Some questions about n-gram models

First order Markov (bi-gram)

$$y_t \in \{1, \dots, K\} \quad p(y_1 = k) = \pi_k^0 \quad p(y_t = k | y_{t-1} = l) = T_{k,l}$$

Q1. How can we compute the marginal distribution over the second state?

$$p(y_2 = l) = \sum_k p(y_2 = l | y_1 = k) p(y_1 = k) = \sum_k T_{lk} \pi_k^0$$

$$p(y_2) = \underline{T} \underline{\pi}^0$$

$$\left(\underline{\pi}^0\right)^T \underline{T}^T$$

Some questions about n-gram models

First order Markov (bi-gram)

$$y_t \in \{1, \dots, K\} \quad p(y_1 = k) = \pi_k^0 \quad p(y_t = k | y_{t-1} = l) = T_{k,l}$$

Q1. How can we compute the marginal distribution over the second state?

$$p(y_2 = k) = \sum_{l=1}^K p(y_2 = k | y_1 = l) p(y_1 = l) = \sum_{l=1}^K T_{k,l} \pi_l^0$$

Some questions about n-gram models

First order Markov (bi-gram)

$$y_t \in \{1, \dots, K\} \quad p(y_1 = k) = \pi_k^0 \quad p(y_t = k | y_{t-1} = l) = T_{k,l}$$

Q1. How can we compute the marginal distribution over the second state?

$$p(y_2 = k) = \sum_{l=1}^K p(y_2 = k | y_1 = l) p(y_1 = l) = \sum_{l=1}^K T_{k,l} \pi_l^0$$

Q2. How can we compute the stationary distribution for the Markov chain?

||
invariant distribution

eigenvalues of transition matrix

$$p(y_n = k) = \sum_l p(y_n = k | y_{n-1} = l) p(y_{n-1} = l)$$

$$p(y_n = k) = \sum_l T_{k,l} p(y_{n-1} = l)$$

$$\underline{1} \times \underline{p_\infty} = \underline{T} \underline{p_\infty} \quad \underline{T} \underline{e_\mu} = \lambda_\mu \underline{e_\mu} \quad \lambda_\mu = 1$$

Some questions about n-gram models

First order Markov (bi-gram)

$$y_t \in \{1, \dots, K\} \quad p(y_1 = k) = \pi_k^0 \quad p(y_t = k | y_{t-1} = l) = T_{k,l}$$

Q1. How can we compute the marginal distribution over the second state?

$$p(y_2 = k) = \sum_{l=1}^K p(y_2 = k | y_1 = l) p(y_1 = l) = \sum_{l=1}^K T_{k,l} \pi_l^0$$

Q2. How can we compute the stationary distribution for the Markov chain?

$$p(y_t = k) = \sum_{l=1}^K p(y_t = k | y_{t-1} = l) p(y_{t-1} = l)$$

Some questions about n-gram models

First order Markov (bi-gram)

$$y_t \in \{1, \dots, K\} \quad p(y_1 = k) = \pi_k^0 \quad p(y_t = k | y_{t-1} = l) = T_{k,l}$$

Q1. How can we compute the marginal distribution over the second state?

$$p(y_2 = k) = \sum_{l=1}^K p(y_2 = k | y_1 = l) p(y_1 = l) = \sum_{l=1}^K T_{k,l} \pi_l^0$$

Q2. How can we compute the stationary distribution for the Markov chain?

$$p(y_t = k) = \sum_{l=1}^K p(y_t = k | y_{t-1} = l) p(y_{t-1} = l) \quad \text{eigenvectors of transition matrix with eigenvalue} = 1$$
$$\pi_k^\infty = \sum_{l=1}^K T_{k,l} \pi_l^\infty$$

Some questions about n-gram models

First order Markov (bi-gram)

$$y_t \in \{1, \dots, K\} \quad p(y_1 = k) = \pi_k^0 \quad p(y_t = k | y_{t-1} = l) = T_{k,l}$$

Q1. How can we compute the marginal distribution over the second state?

$$p(y_2 = k) = \sum_{l=1}^K p(y_2 = k | y_1 = l) p(y_1 = l) = \sum_{l=1}^K T_{k,l} \pi_l^0$$

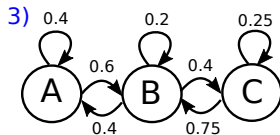
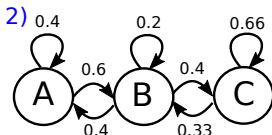
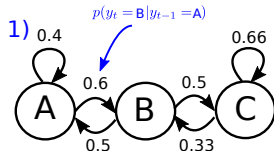
Q2. How can we compute the stationary distribution for the Markov chain?

$$p(y_t = k) = \sum_{l=1}^K p(y_t = k | y_{t-1} = l) p(y_{t-1} = l) \quad \text{eigenvectors of transition matrix with eigenvalue} = 1$$
$$\pi_k^\infty = \sum_{l=1}^K T_{k,l} \pi_l^\infty$$

Q3. Which transition matrix is most compatible with the following sequence?

ABAAABBABCCCB

'State Transition Diagrams'



Some questions about n-gram models

	to state			
	A	B	C	total
from state A	2/5	3/5	0	5
from state B	2/5	1/5	2/5	5
from state C	0	1/3	2/3	3

First order Markov (bi-gram)

$$y_t \in \{1, \dots, K\} \quad p(y_1 = k) = \pi_k^0 \quad p(y_t = k | y_{t-1} = l) = T_{k,l}$$

Q1. How can we compute the marginal distribution over the second state?

$$p(y_2 = k) = \sum_{l=1}^K p(y_2 = k | y_1 = l) p(y_1 = l) = \sum_{l=1}^K T_{k,l} \pi_l^0$$

Q2. How can we compute the stationary distribution for the Markov chain?

$$p(y_t = k) = \sum_{l=1}^K p(y_t = k | y_{t-1} = l) p(y_{t-1} = l)$$

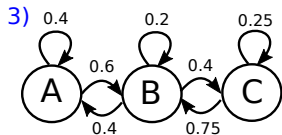
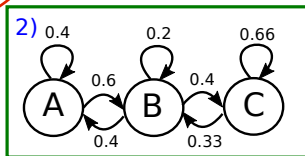
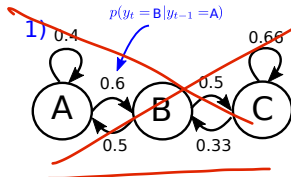
eigenvectors of transition matrix with eigenvalue = 1

$$\pi_k^\infty = \sum_{l=1}^K T_{k,l} \pi_l^\infty$$

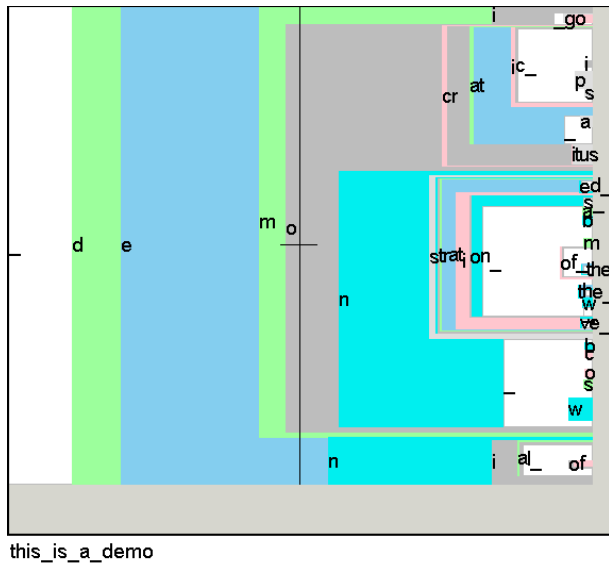
Q3. Which transition matrix is most compatible with the following sequence?

ABAAABBABCCCB

'State Transition Diagrams'



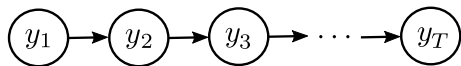
Example application of n-grams: text modelling for dasher



Markov models for discrete data: n-gram models

First order Markov (bi-gram)

$$p(y_1, y_2, y_3, \dots, y_T) = p(y_1)p(y_2|y_1)p(y_3|y_2) \dots p(y_T|y_{T-1})$$



$$y_t \in \{1, \dots, K\}$$

discrete states

$$p(y_1 = k) = \pi_k^0$$

initial state probabilities

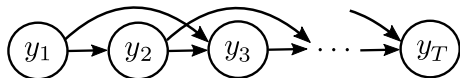
$$p(y_t = k | y_{t-1} = l) = T_{k,l}$$

transition probabilities
(stochastic matrix)

$$\sum_{k=1}^K T_{k,l} = 1$$

Second order Markov (tri-gram)

$$p(y_1, y_2, y_3, \dots, y_T) = p(y_1)p(y_2|y_1)p(y_3|y_2, y_1) \dots p(y_T|y_{T-1}, y_{T-2})$$



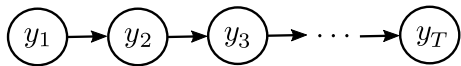
$$p(y_t = k | y_{t-1} = l, y_{t-2} = m) = T_{k,l,m}$$

n-grams require large
multidimensional arrays

Markov models for continuous data: Auto-Regressive (AR) Gaussian models

First order Markov (AR(1))

$$p(y_1, y_2, y_3, \dots, y_T) = p(y_1)p(y_2|y_1)p(y_3|y_2) \dots p(y_T|y_{T-1})$$



$$y_t \in \mathbb{R}^D$$

↑
continuous vector states

$$p(y_1) = \mathcal{G}(y_1; \mu_0, \Sigma_0)$$

↑
initial state density

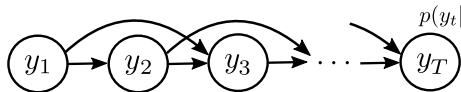
$$p(y_t|y_{t-1}) = \mathcal{G}(y_t; \Lambda y_{t-1}, \Sigma)$$

↑
transition density

$$\mathcal{G}(y; \mu, \Sigma) = \frac{1}{(2\pi)^{D/2} \det(\Sigma)^{1/2}} \exp \left\{ -\frac{1}{2} (y - \mu)^T \Sigma^{-1} (y - \mu) \right\}$$

Second order Markov (AR(2))

$$p(y_1, y_2, y_3, \dots, y_T) = p(y_1)p(y_2|y_1)p(y_3|y_2, y_1) \dots p(y_T|y_{T-1}, y_{T-2})$$



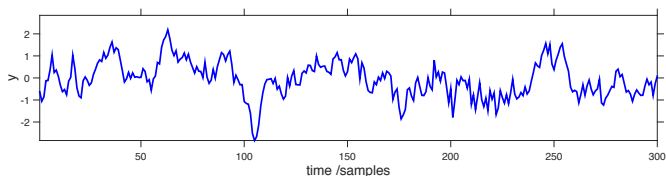
$$p(y_t|y_{t-1}, y_{t-2}) = \mathcal{G}(y_t; \Lambda_1 y_{t-1} + \Lambda_2 y_{t-2}, \Sigma)$$

joint distribution over all variables
is always multivariate Gaussian

Markov models for continuous data: Auto-Regressive (AR) Gaussian models

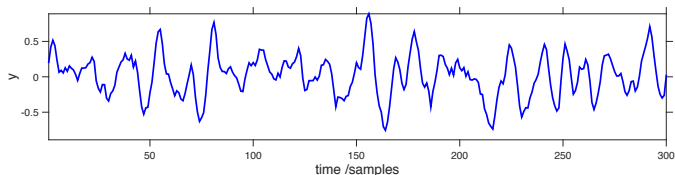
First order Markov (AR(1))

$$y_t \in \mathbb{R}^1 \quad p(y_t|y_{t-1}) = \mathcal{G}(y_t; \lambda y_{t-1}, \sigma^2) \quad \lambda = 0.9 \quad \sigma^2 = 0.01$$



Second order Markov (AR(2))

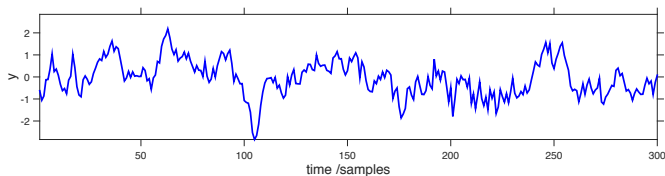
$$y_t \in \mathbb{R}^1 \quad p(y_t|y_{t-1}, y_{t-2}) = \mathcal{G}(y_t; \lambda_1 y_{t-1} + \lambda_2 y_{t-2}, \sigma^2)$$
$$[\lambda_1, \lambda_2] = [1.57, -0.78] \quad \sigma^2 = 0.01$$



Markov models for continuous data: Auto-Regressive (AR) Gaussian models

First order Markov (AR(1))

$$y_t \in \mathbb{R}^1 \quad p(y_t|y_{t-1}) = \mathcal{G}(y_t; \lambda y_{t-1}, \sigma^2) \quad \lambda = 0.9 \quad \sigma^2 = 0.01$$

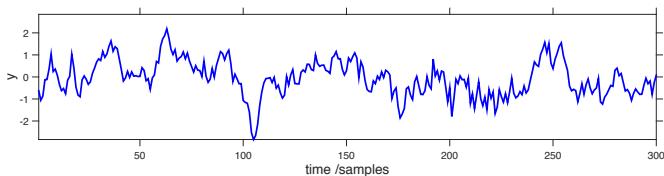


What is the stationary distribution of this process? $p(y_\infty) = ?$

Markov models for continuous data: Auto-Regressive (AR) Gaussian models

First order Markov (AR(1))

$$y_t \in \mathbb{R}^1 \quad p(y_t|y_{t-1}) = \mathcal{G}(y_t; \lambda y_{t-1}, \sigma^2) \quad \lambda = 0.9 \quad \sigma^2 = 0.01$$



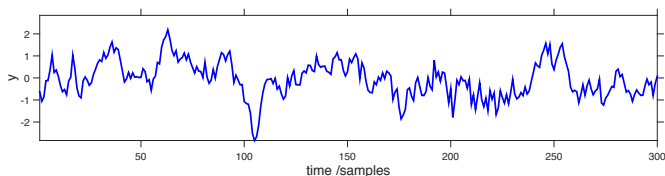
What is the stationary distribution of this process? $p(y_\infty) = ?$

Everything is linear Gaussian => must be Gaussian $p(y_\infty) = \mathcal{G}(y_\infty; \mu_\infty, \sigma_\infty^2)$

Markov models for continuous data: Auto-Regressive (AR) Gaussian models

First order Markov (AR(1))

$$y_t \in \mathbb{R}^1 \quad p(y_t|y_{t-1}) = \mathcal{G}(y_t; \lambda y_{t-1}, \sigma^2) \quad \lambda = 0.9 \quad \sigma^2 = 0.01$$



What is the stationary distribution of this process? $p(y_\infty) = ?$

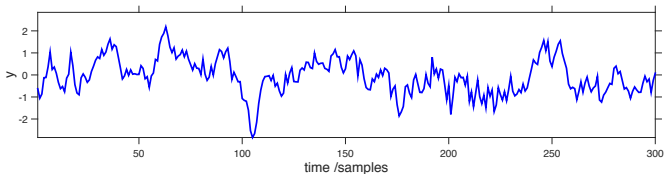
Everything is linear Gaussian => must be Gaussian $p(y_\infty) = \mathcal{G}(y_\infty; \mu_\infty, \sigma_\infty^2)$

$$y_t = \lambda y_{t-1} + \sigma \epsilon_t \quad \epsilon_t \sim \mathcal{G}(0, 1)$$

Markov models for continuous data: Auto-Regressive (AR) Gaussian models

First order Markov (AR(1))

$$y_t \in \mathbb{R}^1 \quad p(y_t|y_{t-1}) = \mathcal{G}(y_t; \lambda y_{t-1}, \sigma^2) \quad \lambda = 0.9 \quad \sigma^2 = 0.01$$



What is the stationary distribution of this process? $p(y_\infty) = ?$

Everything is linear Gaussian \Rightarrow must be Gaussian $p(y_\infty) = \mathcal{G}(y_\infty; \mu_\infty, \sigma_\infty^2)$

$$y_t = \lambda y_{t-1} + \sigma \epsilon_t \quad \epsilon_t \sim \mathcal{G}(0, 1)$$

Mean:

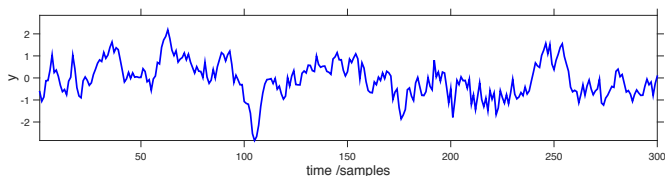
$$\begin{aligned} \langle y_t \rangle &= \langle \lambda y_{t-1} + \sigma \epsilon_t \rangle = \lambda \langle y_{t-1} \rangle + \sigma \langle \epsilon_t \rangle \\ \uparrow & \\ \mathbb{E}(y_t) &= \lambda \langle y_{t-1} \rangle + 0 \end{aligned}$$

$$\mu_\infty = \lambda \mu_\infty \Rightarrow \mu_\infty = 0$$

Markov models for continuous data: Auto-Regressive (AR) Gaussian models

First order Markov (AR(1))

$$y_t \in \mathbb{R}^1 \quad p(y_t|y_{t-1}) = \mathcal{G}(y_t; \lambda y_{t-1}, \sigma^2) \quad \lambda = 0.9 \quad \sigma^2 = 0.01$$



What is the stationary distribution of this process? $p(y_\infty) = ?$

Everything is linear Gaussian => must be Gaussian $p(y_\infty) = \mathcal{G}(y_\infty; \mu_\infty, \sigma_\infty^2)$

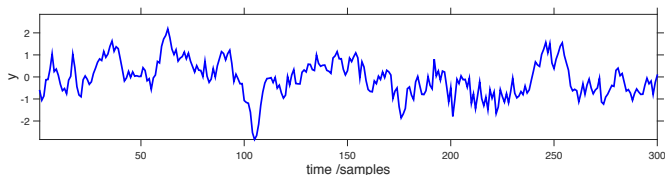
$$y_t = \lambda y_{t-1} + \sigma \epsilon_t \quad \epsilon_t \sim \mathcal{G}(0, 1)$$

Mean: $\langle y_t \rangle = \lambda \langle y_{t-1} \rangle + \sigma \langle \epsilon_t \rangle = 0$

Markov models for continuous data: Auto-Regressive (AR) Gaussian models

First order Markov (AR(1))

$$y_t \in \mathbb{R}^1 \quad p(y_t|y_{t-1}) = \mathcal{G}(y_t; \lambda y_{t-1}, \sigma^2) \quad \lambda = 0.9 \quad \sigma^2 = 0.01$$



What is the stationary distribution of this process? $p(y_\infty) = ?$

Everything is linear Gaussian => must be Gaussian $p(y_\infty) = \mathcal{G}(y_\infty; \mu_\infty, \sigma_\infty^2)$

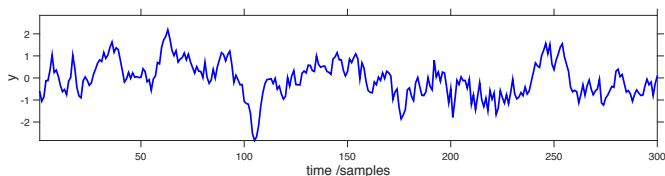
$$y_t = \lambda y_{t-1} + \sigma \epsilon_t \quad \epsilon_t \sim \mathcal{G}(0, 1)$$

Mean: $\langle y_t \rangle = \lambda \langle y_{t-1} \rangle + \sigma \langle \epsilon_t \rangle = 0 \quad \mu_\infty = 0$

Markov models for continuous data: Auto-Regressive (AR) Gaussian models

First order Markov (AR(1))

$$y_t \in \mathbb{R}^1 \quad p(y_t|y_{t-1}) = \mathcal{G}(y_t; \lambda y_{t-1}, \sigma^2) \quad \lambda = 0.9 \quad \sigma^2 = 0.01$$



What is the stationary distribution of this process? $p(y_\infty) = ?$

Everything is linear Gaussian => must be Gaussian $p(y_\infty) = \mathcal{G}(y_\infty; \mu_\infty, \sigma_\infty^2)$

$$y_t = \lambda y_{t-1} + \sigma \epsilon_t \quad \epsilon_t \sim \mathcal{G}(0, 1)$$

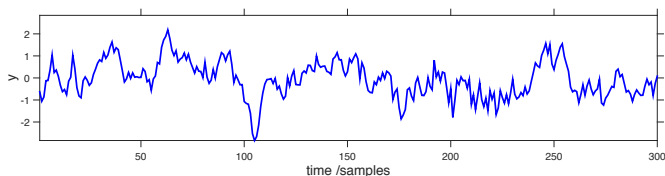
Mean: $\langle y_t \rangle = \lambda \langle y_{t-1} \rangle + \sigma \langle \epsilon_t \rangle = 0 \quad \mu_\infty = 0$

Variance: $\langle y_t^2 \rangle$

Markov models for continuous data: Auto-Regressive (AR) Gaussian models

First order Markov (AR(1))

$$y_t \in \mathbb{R}^1 \quad p(y_t|y_{t-1}) = \mathcal{G}(y_t; \lambda y_{t-1}, \sigma^2) \quad \lambda = 0.9 \quad \sigma^2 = 0.01$$



What is the stationary distribution of this process? $p(y_\infty) = ?$

Everything is linear Gaussian => must be Gaussian $p(y_\infty) = \mathcal{G}(y_\infty; \mu_\infty, \sigma_\infty^2)$

$$y_t = \lambda y_{t-1} + \sigma \epsilon_t \quad \epsilon_t \sim \mathcal{G}(0, 1)$$

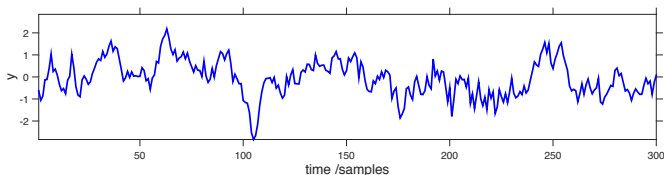
Mean: $\langle y_t \rangle = \lambda \langle y_{t-1} \rangle + \sigma \langle \epsilon_t \rangle = 0 \quad \mu_\infty = 0$

Variance: $\langle y_t^2 \rangle = \langle (\lambda y_{t-1} + \sigma \epsilon_t)^2 \rangle$

Markov models for continuous data: Auto-Regressive (AR) Gaussian models

First order Markov (AR(1))

$$y_t \in \mathbb{R}^1 \quad p(y_t|y_{t-1}) = \mathcal{G}(y_t; \lambda y_{t-1}, \sigma^2) \quad \lambda = 0.9 \quad \sigma^2 = 0.01$$



What is the stationary distribution of this process? $p(y_\infty) = ?$

Everything is linear Gaussian => must be Gaussian $p(y_\infty) = \mathcal{G}(y_\infty; \mu_\infty, \sigma_\infty^2)$

$$y_t = \lambda y_{t-1} + \sigma \epsilon_t \quad \epsilon_t \sim \mathcal{G}(0, 1)$$

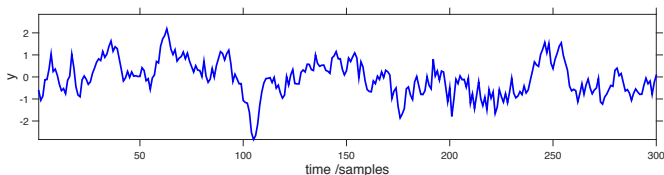
Mean: $\langle y_t \rangle = \lambda \langle y_{t-1} \rangle + \sigma \langle \epsilon_t \rangle = 0 \quad \mu_\infty = 0$

Variance: $\langle y_t^2 \rangle = \langle (\lambda y_{t-1} + \sigma \epsilon_t)^2 \rangle = \lambda^2 \langle y_{t-1}^2 \rangle + 2\lambda \sigma \langle y_{t-1} \epsilon_t \rangle + \sigma^2 \langle \epsilon_t^2 \rangle$

Markov models for continuous data: Auto-Regressive (AR) Gaussian models

First order Markov (AR(1))

$$y_t \in \mathbb{R}^1 \quad p(y_t|y_{t-1}) = \mathcal{G}(y_t; \lambda y_{t-1}, \sigma^2) \quad \lambda = 0.9 \quad \sigma^2 = 0.01$$



What is the stationary distribution of this process? $p(y_\infty) = ?$

Everything is linear Gaussian => must be Gaussian $p(y_\infty) = \mathcal{G}(y_\infty; \mu_\infty, \sigma_\infty^2)$

$$y_t = \lambda y_{t-1} + \sigma \epsilon_t \quad \epsilon_t \sim \mathcal{G}(0, 1)$$

Mean: $\langle y_t \rangle = \lambda \langle y_{t-1} \rangle + \sigma \langle \epsilon_t \rangle = 0 \quad \mu_\infty = 0$

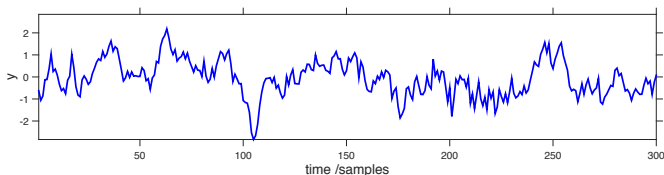
Variance: $\langle y_t^2 \rangle = \langle (\lambda y_{t-1} + \sigma \epsilon_t)^2 \rangle = \lambda^2 \langle y_{t-1}^2 \rangle + 2\lambda \sigma \langle y_{t-1} \epsilon_t \rangle + \sigma^2 \langle \epsilon_t^2 \rangle$

$$\langle y_t^2 \rangle = \lambda^2 \langle y_{t-1}^2 \rangle + \sigma^2$$

Markov models for continuous data: Auto-Regressive (AR) Gaussian models

First order Markov (AR(1))

$$y_t \in \mathbb{R}^1 \quad p(y_t|y_{t-1}) = \mathcal{G}(y_t; \lambda y_{t-1}, \sigma^2) \quad \lambda = 0.9 \quad \sigma^2 = 0.01$$



What is the stationary distribution of this process? $p(y_\infty) = ?$

Everything is linear Gaussian => must be Gaussian $p(y_\infty) = \mathcal{G}(y_\infty; \mu_\infty, \sigma_\infty^2)$

$$y_t = \lambda y_{t-1} + \sigma \epsilon_t \quad \epsilon_t \sim \mathcal{G}(0, 1)$$

Mean: $\langle y_t \rangle = \lambda \langle y_{t-1} \rangle + \sigma \langle \epsilon_t \rangle = 0 \quad \mu_\infty = 0$

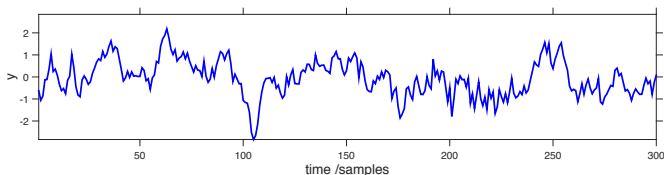
Variance: $\langle y_t^2 \rangle = \langle (\lambda y_{t-1} + \sigma \epsilon_t)^2 \rangle = \lambda^2 \langle y_{t-1}^2 \rangle + 2\lambda \sigma \langle y_{t-1} \epsilon_t \rangle + \sigma^2 \langle \epsilon_t^2 \rangle$

$$\langle y_t^2 \rangle = \lambda^2 \langle y_{t-1}^2 \rangle + \sigma^2 \quad \sigma_\infty^2 = \lambda^2 \sigma_\infty^2 + \sigma^2$$

Markov models for continuous data: Auto-Regressive (AR) Gaussian models

First order Markov (AR(1))

$$y_t \in \mathbb{R}^1 \quad p(y_t|y_{t-1}) = \mathcal{G}(y_t; \lambda y_{t-1}, \sigma^2) \quad \lambda = 0.9 \quad \sigma^2 = 0.01$$



What is the stationary distribution of this process? $p(y_\infty) = ?$

Everything is linear Gaussian => must be Gaussian $p(y_\infty) = \mathcal{G}(y_\infty; \mu_\infty, \sigma_\infty^2)$

$$y_t = \lambda y_{t-1} + \sigma \epsilon_t \quad \epsilon_t \sim \mathcal{G}(0, 1)$$

Mean: $\langle y_t \rangle = \lambda \langle y_{t-1} \rangle + \sigma \langle \epsilon_t \rangle = 0 \quad \mu_\infty = 0$

Variance: $\langle y_t^2 \rangle = \langle (\lambda y_{t-1} + \sigma \epsilon_t)^2 \rangle = \lambda^2 \langle y_{t-1}^2 \rangle + 2\lambda \sigma \langle y_{t-1} \epsilon_t \rangle + \sigma^2 \langle \epsilon_t^2 \rangle$

$$\langle y_t^2 \rangle = \lambda^2 \langle y_{t-1}^2 \rangle + \sigma^2 \quad \sigma_\infty^2 = \lambda^2 \sigma_\infty^2 + \sigma^2 \quad \sigma_\infty^2 = \frac{\sigma^2}{1-\lambda^2}$$

Markov Models

1st Order

$$p(y_{1:T}) = p(y_1) p(y_2 | y_1) \dots p(y_T | y_{T-1})$$

discrete $y \Rightarrow$ bigram models

$$p(y_1 = k) = \pi_k^0$$

$$p(y_t = k | y_{t-1} = l) = T_{kl}$$

\Downarrow

Stationary / invariant distribution

$$p(y_\infty = k) = \pi_k^\infty = \sum_l T_{kl} \pi_l^\infty$$

3FB Q&A Wed 9.30-11.00

CBL Seminar Room

see Moodle for details

continuous $y \Rightarrow$ autoregressive

$$p(\underline{y}_1) = \mathcal{G}(\underline{y}_1; \underline{\mu}_0, \underline{\Sigma}_0)$$

$\dim(\underline{y}_t) = D$ DxD matrices

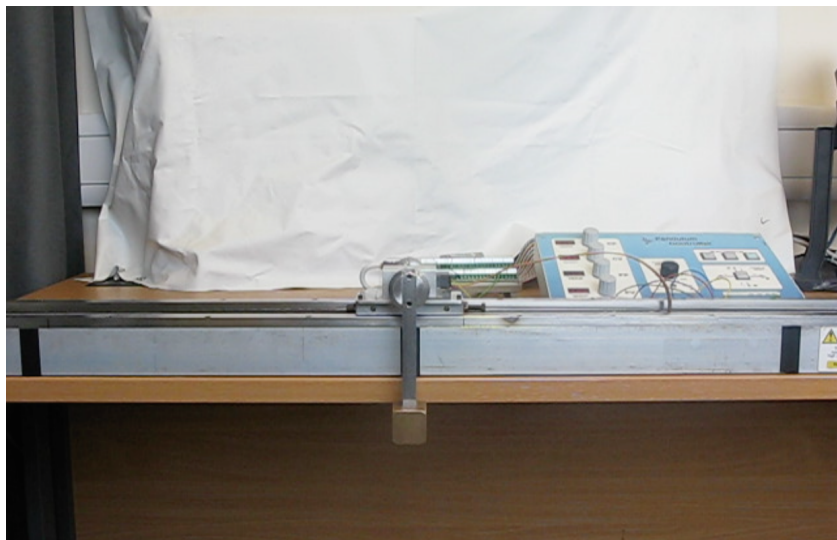
$$p(\underline{y}_t | \underline{y}_{t-1}) = \mathcal{G}(\underline{y}_t; \underline{\mu}_{t-1}, \underline{\Sigma}_{t-1})$$

\Downarrow

Stationary distribution

$$p(y_\infty) = \mathcal{G}(y_\infty; \mu_\infty = 0, \sigma_\infty^2 = \frac{\sigma^2}{1 - \lambda^2})$$

Example application of Markov Models: pendulum swing up control problem



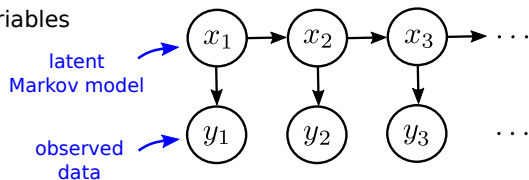
Hidden Markov models

Real data depend on latent variables

ASR

x phonemes/words

y waveform/feature



Computer Vision

x objects, pose, lighting

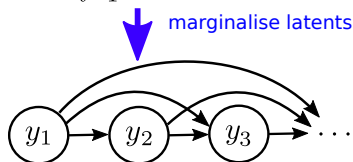
y image pixel intensities

$$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t | x_{t-1}) p(y_t | x_t)$$

Natural Language Processing

x topics

y words



Two prevalent Examples:

Hidden Markov Models (discrete x)

Linear Gaussian State Space Models (Gaussian x and y)

$p(y_{1:T})$
fully connected

Hidden Markov models: discrete hidden state

Discrete Hidden State

$$x_t \in \{1, \dots, K\}$$

$$p(x_t = k | x_{t-1} = l) = T_{k,l}$$

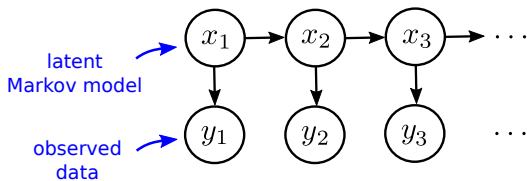
E.g. in examples below $K = 2$

$$T = \begin{bmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{bmatrix}$$

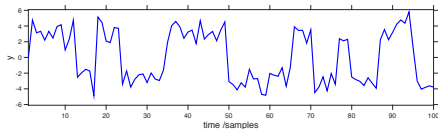
Continuous Observed State

$$p(y_t | x_t = k) = \mathcal{G}(y_t; \mu_k, \Sigma_k)$$

$$\mu_1 = 3 \quad \mu_2 = -3 \quad \sigma_1^2 = \sigma_2^2 = 1$$



$$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t | x_{t-1}) p(y_t | x_t)$$



Hidden Markov models: discrete hidden state

Discrete Hidden State

$$x_t \in \{1, \dots, K\}$$

$$p(x_t = k | x_{t-1} = l) = T_{k,l}$$

E.g. in examples below $K = 2$

$$T = \begin{bmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{bmatrix}$$

Continuous Observed State

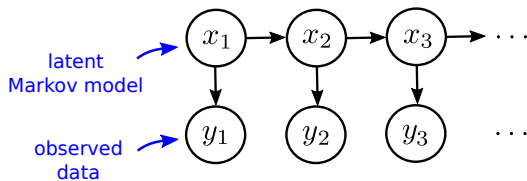
$$p(y_t | x_t = k) = \mathcal{G}(y_t; \mu_k, \Sigma_k)$$

$$\mu_1 = 3 \quad \mu_2 = -3 \quad \sigma_1^2 = \sigma_2^2 = 1$$

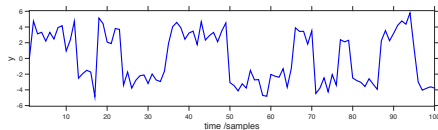
Discrete Observed State

$$p(y_t = l | x_t = k) = S_{l,k}$$

$$S = \begin{bmatrix} 0.5 & 0 \\ 0.5 & 0 \\ 0 & 1 \end{bmatrix}$$



$$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t | x_{t-1}) p(y_t | x_t)$$



ABBBBBAAABAACCCCB BBBB BCCCCCCCCC
AAABBBBAABAAB BCCCCCCCCCCCCC CBBA
AACCCCBABCCCCC CAABBAABABCCCC

Hidden Markov models: discrete hidden state

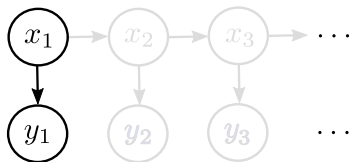
Discrete Hidden State, Continuous Observed State

$$x_t \in \{1, \dots, K\}$$

$$p(x_1 = k) = \pi_k^0$$

$$p(x_t = k | x_{t-1} = l) = T_{k,l}$$

$$p(y_t | x_t = k) = \mathcal{G}(y_t; \mu_k, \Sigma_k)$$



Consider $T = 1$

Q1: What type of distribution is $p(y_1)$?

$$p(y_1) = \sum_{k=1}^K p(y_1, x_1 = k) \quad \text{Sum rule}$$
$$= \sum_{k=1}^K p(x_1 = k) p(y_1 | x_1 = k) \quad \text{Product rule}$$

$$p(y_1) = \sum_{k=1}^K \pi_k^0 \mathcal{G}(y_1; \mu_k, \Sigma_k)$$

Hidden Markov models: discrete hidden state

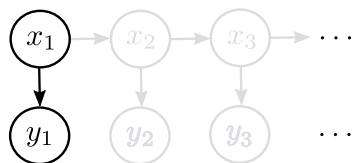
Discrete Hidden State, Continuous Observed State

$$x_t \in \{1, \dots, K\}$$

$$p(x_1 = k) = \pi_k^0$$

$$p(x_t = k | x_{t-1} = l) = T_{k,l}$$

$$p(y_t | x_t = k) = \mathcal{G}(y_t; \mu_k, \Sigma_k)$$



Consider $T = 1$

Q1: What type of distribution is $p(y_1)$?

$$p(y_1) = \sum_k p(y_1 | x_1 = k) p(x_1 = k)$$

Hidden Markov models: discrete hidden state

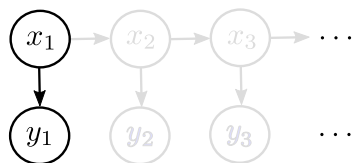
Discrete Hidden State, Continuous Observed State

$$x_t \in \{1, \dots, K\}$$

$$p(x_1 = k) = \pi_k^0$$

$$p(x_t = k | x_{t-1} = l) = T_{k,l}$$

$$p(y_t | x_t = k) = \mathcal{G}(y_t; \mu_k, \Sigma_k)$$



Consider T = 1

Q1: What type of distribution is $p(y_1)$?

$$p(y_1) = \sum_k p(y_1 | x_1 = k) p(x_1 = k) = \sum_k \pi_k^0 \mathcal{G}(y_1; \mu_k, \Sigma_k)$$

Hidden Markov models: discrete hidden state

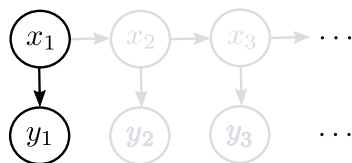
Discrete Hidden State, Continuous Observed State

$$x_t \in \{1, \dots, K\}$$

$$p(x_1 = k) = \pi_k^0$$

$$p(x_t = k | x_{t-1} = l) = T_{k,l}$$

$$p(y_t | x_t = k) = \mathcal{G}(y_t; \mu_k, \Sigma_k)$$



Q1: What type of distribution is $p(y_1)$?

$$p(y_1) = \sum_k p(y_1 | x_1 = k) p(x_1 = k) = \sum_k \pi_k^0 \mathcal{G}(y_1; \mu_k, \Sigma_k)$$

Q2: What distribution does $p(y_t)$ converge to after a long time?

Hidden Markov models: discrete hidden state

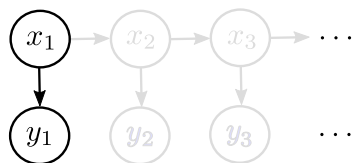
Discrete Hidden State, Continuous Observed State

$$x_t \in \{1, \dots, K\}$$

$$p(x_1 = k) = \pi_k^0$$

$$p(x_t = k | x_{t-1} = l) = T_{k,l}$$

$$p(y_t | x_t = k) = \mathcal{G}(y_t; \mu_k, \Sigma_k)$$



Consider $T = 1$

Q1: What type of distribution is $p(y_1)$?

$$p(y_1) = \sum_k p(y_1 | x_1 = k) p(x_1 = k) = \sum_k \pi_k^0 \mathcal{G}(y_1; \mu_k, \Sigma_k)$$

Q2: What distribution does $p(y_t)$ converge to after a long time?

stationary distribution of Markov chain satisfies $\pi_k^\infty = \sum_{l=1}^K T_{k,l} \pi_l^\infty$

Hidden Markov models: discrete hidden state

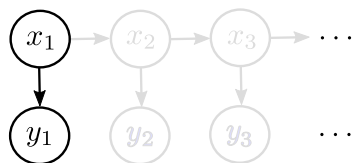
Discrete Hidden State, Continuous Observed State

$$x_t \in \{1, \dots, K\}$$

$$p(x_1 = k) = \pi_k^0$$

$$p(x_t = k | x_{t-1} = l) = T_{k,l}$$

$$p(y_t | x_t = k) = \mathcal{G}(y_t; \mu_k, \Sigma_k)$$



Consider $T = 1$

Q1: What type of distribution is $p(y_1)$?

$$p(y_1) = \sum_k p(y_1 | x_1 = k) p(x_1 = k) = \sum_k \pi_k^0 \mathcal{G}(y_1; \mu_k, \Sigma_k)$$

Q2: What distribution does $p(y_t)$ converge to after a long time?

stationary distribution of Markov chain satisfies $\pi_k^\infty = \sum_{l=1}^K T_{k,l} \pi_l^\infty$

$$p(y_t) = \sum_k p(y_t | x_t = k) p(x_t = k)$$

Hidden Markov models: discrete hidden state

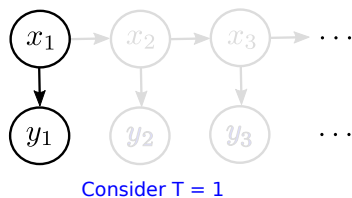
Discrete Hidden State, Continuous Observed State

$$x_t \in \{1, \dots, K\}$$

$$p(x_1 = k) = \pi_k^0$$

$$p(x_t = k | x_{t-1} = l) = T_{k,l}$$

$$p(y_t | x_t = k) = \mathcal{G}(y_t; \mu_k, \Sigma_k)$$



Q1: What type of distribution is $p(y_1)$?

$$p(y_1) = \sum_k p(y_1 | x_1 = k) p(x_1 = k) = \sum_k \pi_k^0 \mathcal{G}(y_1; \mu_k, \Sigma_k)$$

Q2: What distribution does $p(y_t)$ converge to after a long time?

stationary distribution of Markov chain satisfies $\pi_k^\infty = \sum_{l=1}^K T_{k,l} \pi_l^\infty$

$$p(y_t) = \sum_k p(y_t | x_t = k) p(x_t = k) \rightarrow \sum_k \pi_k^\infty \mathcal{G}(y_t; \mu_k, \Sigma_k)$$

Hidden Markov models: discrete hidden state

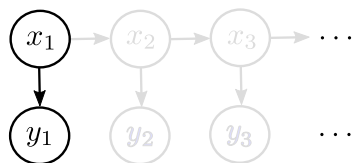
Discrete Hidden State, Continuous Observed State

$$x_t \in \{1, \dots, K\}$$

$$p(x_1 = k) = \pi_k^0$$

$$p(x_t = k | x_{t-1} = l) = T_{k,l}$$

$$p(y_t | x_t = k) = \mathcal{G}(y_t; \mu_k, \Sigma_k)$$



Consider $T = 1$

Q1: What type of distribution is $p(y_1)$?

$$p(y_1) = \sum_k p(y_1 | x_1 = k) p(x_1 = k) = \sum_k \pi_k^0 \mathcal{G}(y_1; \mu_k, \Sigma_k)$$

Q2: What distribution does $p(y_t)$ converge to after a long time?

stationary distribution of Markov chain satisfies $\pi_k^\infty = \sum_{l=1}^K T_{k,l} \pi_l^\infty$

$$p(y_t) = \sum_k p(y_t | x_t = k) p(x_t = k) \rightarrow \sum_k \pi_k^\infty \mathcal{G}(y_t; \mu_k, \Sigma_k)$$

this HMM = Mixture of Gaussian Models with dynamic cluster assignments

Hidden Markov models: continuous hidden state (LGSSMs)

Continuous Hidden State

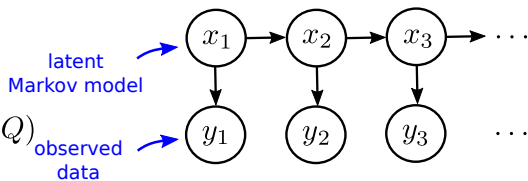
$$x_t \in \mathbb{R}^K$$

$$p(x_t|x_{t-1}) = \mathcal{G}(x_t; Ax_{t-1}, Q)$$

Continuous Observed State

$$y_t \in \mathbb{R}^D$$

$$p(y_t|x_t) = \mathcal{G}(y_t; Cx_t, R)$$



$$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t|x_{t-1})p(y_t|x_t)$$

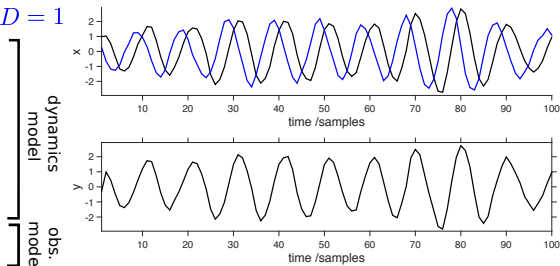
E.g. simple example $K = 2$ $D = 1$

$$A = \lambda \begin{bmatrix} \cos(\theta) & \sin(\theta) \\ -\sin(\theta) & \cos(\theta) \end{bmatrix}$$

$$\lambda = 0.99 \quad \theta = 2\pi/10$$

$$Q = (1 - \lambda^2) \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$C = [1, 0] \quad R = 0.01$$



$$\hat{z} \sim N(0, \sigma^2) \approx p^*(\theta) = \int p(\theta) d\theta \approx \int p^*(\theta) d\theta \approx p(z_{y_3} | x_1, x_2)$$

$$\underline{x}_t = \begin{bmatrix} x_{1t} \\ x_{2t} \end{bmatrix} = \underline{A} \underline{x}_{t-1} + \underline{Q}^{1/2} \underline{\varepsilon}_t \quad \underline{\varepsilon}_t \sim G(0, \underline{I})$$

matrix square root

$$\begin{bmatrix} x_{1t} \\ x_{2t} \end{bmatrix} = \lambda \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} x_{1t-1} \\ x_{2t-1} \end{bmatrix} + (1-\lambda^2)^{1/2} \underline{\varepsilon}_t$$



$$y_t = \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} x_{1t} \\ x_{2t} \end{bmatrix} + \text{noise} = x_{1t} + \text{noise}$$

Summary Sequence Modelling Lecture II

$$p(\underset{\substack{\uparrow \\ \text{observed}}}{y_{1:T}}, \underset{\substack{\uparrow \\ \text{latent}}}{x_{1:T}}) = \prod_{t=1}^T p(x_t | x_{t-1}) p(y_t | x_t)$$

Discrete Hidden State $x_t \in \{1 \dots K\}$ (also called HMMs!)

$$p(x_t = k | x_{t-1} = l) = T_{kl}$$

$$\rightarrow p(y_t | x_t = k) = G(y_t; \mu_k, \Sigma_k) \quad y_t \in \mathbb{R}^D$$

emission $\rightarrow p(y_t = l | x_t = k) = S_{lk} \quad y_t \in \{1 \dots D\}$

Continuous Hidden State $x_t \in \mathbb{R}^K$ (linear Gaussian state space models)

$$\rightarrow \underline{p(x_t | x_{t-1})} = G(\underline{x}_t; \underline{A} \underline{x}_{t-1}, \underline{Q}) \Leftrightarrow \underline{x}_t = \underline{A} \underline{x}_{t-1} + \underline{Q}^{1/2} \underline{\varepsilon}_t$$

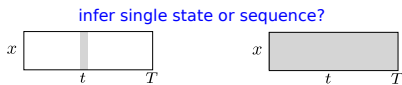
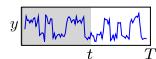
$$p(\underline{y}_t | \underline{x}_t) = G(\underline{y}_t; \underline{C} \underline{x}_t, \underline{R}) \quad \underline{y}_t \in \mathbb{R}^D$$

Today: Inference & Learning

Varieties of Inference

Distributional estimates

future data available?

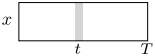
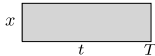
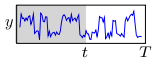



	marginal	joint
filter	$p(x_t y_{1:t})$	$p(x_{1:t} y_{1:t})$
smoother	$p(x_t y_{1:T})$	$p(x_{1:T} y_{1:T})$

$$\begin{aligned}
 1. \text{ LGSSM } \quad & p(x_{1:T} | y_{1:T}) = G_T(x_{1:T}; \mu_{1:T}, \Sigma_{1:T}) \\
 \Rightarrow & x'_{1:T} = \mu_{1:T} \quad \left. \begin{array}{l} \int dx_{\neq t} \\ \uparrow \\ \int dx_t \end{array} \right\} \Rightarrow x^*_{1:T} = x'_{1:T} \\
 & p(x_t | y_{1:T}) = G(x_t; \mu_t, \Sigma_{tt}) \\
 & x_t^* = \mu_t
 \end{aligned}$$

Varieties of Inference

Distributional estimates

		infer single state or sequence?	
			
		marginal	joint
future data available?	 filter	$p(x_t y_{1:t})$	$p(x_{1:t} y_{1:t})$
	 smoother	$p(x_t y_{1:T})$	$p(x_{1:T} y_{1:T})$

Point estimates

$$x_t^* = \arg \max_{x_t} p(x_t | y_{1:T}) \quad \leftarrow \text{most probable state @ } t$$

$$x'_{1:T} = \arg \max_{x_{1:T}} p(x_{1:T} | y_{1:T}) \quad \leftarrow \text{most probable sequence}$$

2. Discrete Hidden State HMM

T=2

x_1	x_2	$p(x_1, x_2 y_1, y_2)$
0	0	0.3
0	1	0.4
1	0	0.3
1	1	0

$x_{1:2}^* = [0, 0]$
 ~~$x_{1:2}^* = [0, 1]$~~
 $x_1^* = 0$
 $x_2^* = 0$

$p(x_1 | y_1, y_2) = [0.7, 0.3]$
 $p(x_2 | y_1, y_2) = [0.6, 0.4]$

Varieties of Inference

Distributional estimates

infer single state or sequence?

	marginal	joint
filter	$p(x_t y_{1:t})$	$p(x_{1:t} y_{1:t})$
smoother	$p(x_t y_{1:T})$	$p(x_{1:T} y_{1:T})$

future data available?

Point estimates

$$x_t^* = \arg \max_{x_t} p(x_t|y_{1:T}) \quad \text{most probable state @ } t$$
$$x'_{1:T} = \arg \max_{x_{1:T}} p(x_{1:T}|y_{1:T}) \quad \text{most probable sequence}$$


Question: are these estimates the same $x_{1:T}^* \stackrel{?}{=} x'_{1:T}$ for

1. Linear Gaussian State Space Models?
2. Discrete Hidden State HMMs?

Varieties of Inference

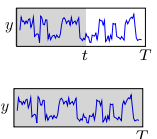
Distributional estimates

infer single state or sequence?



	marginal	joint
filter	$p(x_t y_{1:t})$	$p(x_{1:t} y_{1:t})$
smoother	$p(x_t y_{1:T})$	$p(x_{1:T} y_{1:T})$

future data available?



Point estimates

$$x_t^* = \arg \max_{x_t} p(x_t|y_{1:T}) \quad \text{most probable state @ } t$$
$$x'_{1:T} = \arg \max_{x_{1:T}} p(x_{1:T}|y_{1:T}) \quad \text{most probable sequence}$$


Question: are these estimates the same $x_{1:T}^* \stackrel{?}{=} x'_{1:T}$ for

1. Linear Gaussian State Space Models? $x_{1:T}^* = x'_{1:T}$ (Gaussian)
2. Discrete Hidden State HMMs?

Varieties of Inference

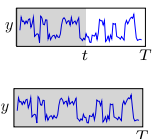
Distributional estimates

infer single state or sequence?



	marginal	joint
filter	$p(x_t y_{1:t})$	$p(x_{1:t} y_{1:t})$
smoother	$p(x_t y_{1:T})$	$p(x_{1:T} y_{1:T})$

future data available?



Point estimates

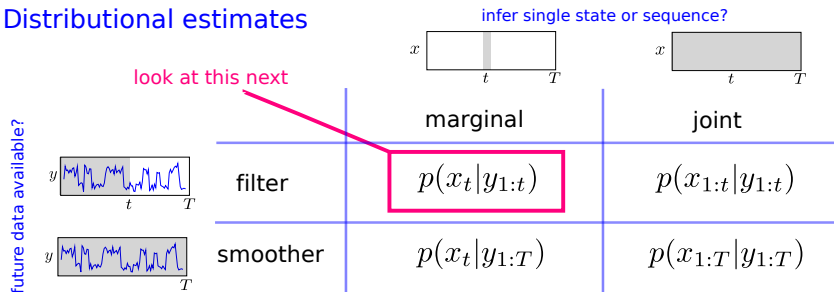
$$x_t^* = \arg \max_{x_t} p(x_t|y_{1:T}) \quad \text{most probable state @ } t$$
$$x'_{1:T} = \arg \max_{x_{1:T}} p(x_{1:T}|y_{1:T}) \quad \text{most probable sequence}$$

Question: are these estimates the same $x_{1:T}^* \stackrel{?}{=} x'_{1:T}$ for

1. Linear Gaussian State Space Models? $x_{1:T}^* = x'_{1:T}$ (Gaussian)
2. Discrete Hidden State HMMs? $x_{1:T}^* \neq x'_{1:T}$

Varieties of Inference

Distributional estimates



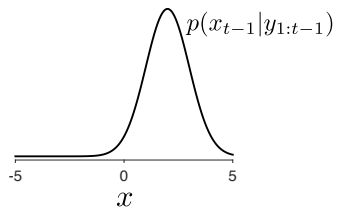
Point estimates

$$x_t^* = \arg \max_{x_t} p(x_t | y_{1:T}) \quad \text{most probable state @ } t$$
$$x'_{1:T} = \arg \max_{x_{1:T}} p(x_{1:T} | y_{1:T}) \quad \text{most probable sequence}$$

Question: are these estimates the same $x_{1:T}^* \stackrel{?}{=} x'_{1:T}$ for

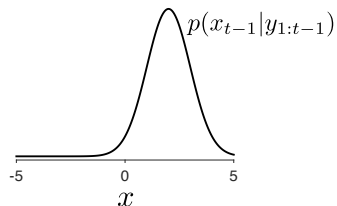
1. Linear Gaussian State Space Models? $x_{1:T}^* = x'_{1:T}$ (Gaussian)
2. Discrete Hidden State HMMs? $x_{1:T}^* \neq x'_{1:T}$

Inference: Kalman Filter



$$p(x_{t-1} | y_{1:t-1})$$

Inference: Kalman Filter

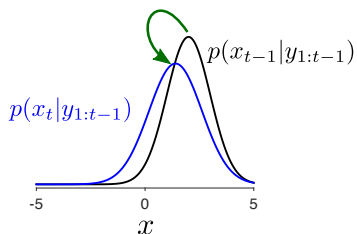


diffuse via dynamics

sum for discrete hidden state

$$p(x_t | y_{1:t-1}) = \int p(x_t | x_{t-1}) p(x_{t-1} | y_{1:t-1}) dx_{t-1}$$

Inference: Kalman Filter

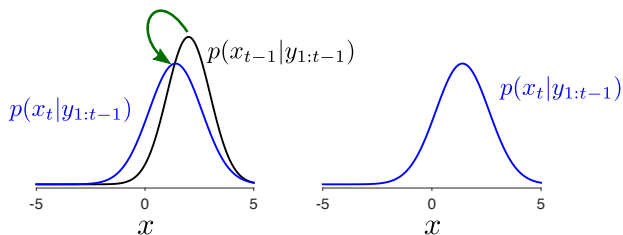


diffuse via dynamics

sum for discrete hidden state

$$p(x_t | y_{1:t-1}) = \int p(x_t | x_{t-1}) p(x_{t-1} | y_{1:t-1}) dx_{t-1}$$

Inference: Kalman Filter



diffuse via dynamics

combine with likelihood

sum for discrete hidden state

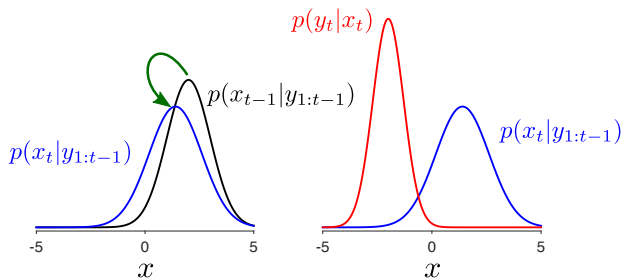
$$p(x_t|y_{1:t-1}) = \int p(x_t|x_{t-1})p(x_{t-1}|y_{1:t-1})dx_{t-1}$$

$p(x_t|y_{1:t}) \propto p(x_t|y_{1:t-1})p(y_t|x_t)$

prior likelihood

Bayes' Rule

Inference: Kalman Filter



diffuse via dynamics

combine with likelihood

$$p(x_t|y_{1:t-1}) = \int p(x_t|x_{t-1})p(x_{t-1}|y_{1:t-1})dx_{t-1}$$

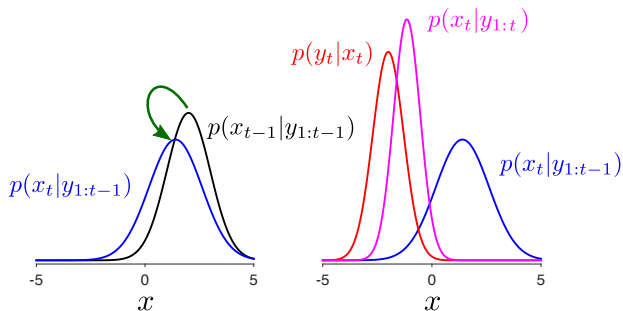
sum for discrete hidden state

$$p(x_t|y_{1:t}) \propto p(x_t|y_{1:t-1})p(y_t|x_t)$$

prior likelihood

Bayes' Rule

Inference: Kalman Filter



diffuse via dynamics

combine with likelihood

$$p(x_t|y_{1:t-1}) = \int p(x_t|x_{t-1})p(x_{t-1}|y_{1:t-1})dx_{t-1}$$

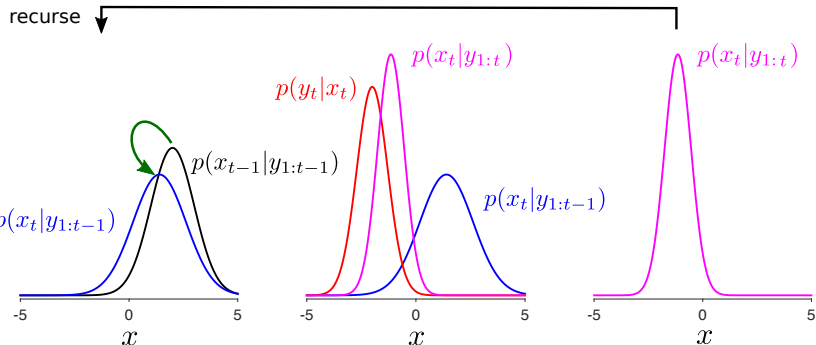
sum for discrete hidden state

$$p(x_t|y_{1:t}) \propto p(x_t|y_{1:t-1})p(y_t|x_t)$$

prior likelihood

Bayes' Rule

Inference: Kalman Filter



diffuse via dynamics

combine with likelihood

sum for discrete hidden state

$$p(x_t|y_{1:t-1}) = \int p(x_t|x_{t-1})p(x_{t-1}|y_{1:t-1})dx_{t-1}$$

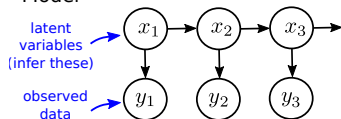
$$p(x_t|y_{1:t}) \propto p(x_t|y_{1:t-1})p(y_t|x_t)$$

prior likelihood

Bayes' Rule

Inference: Derivation of General Filtering Equations

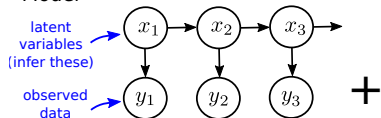
Model



$$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t | x_{t-1}) p(y_t | x_t)$$

Inference: Derivation of General Filtering Equations

Model



$$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t | x_{t-1}) p(y_t | x_t)$$

Rules of probability

product rule

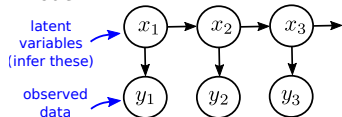
$$p(A|B, C) = \frac{1}{p(B|C)} p(B|A, C) p(A|C)$$

sum rule

$$p(A|C) = \sum_B p(A, B|C)$$

Inference: Derivation of General Filtering Equations

Model



$$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t | x_{t-1}) p(y_t | x_t)$$

Rules of probability

product rule

$$p(A|B, C) = \frac{1}{p(B|C)} p(B|A, C) p(A|C)$$

sum rule

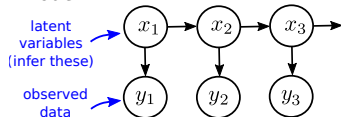
$$p(A|C) = \sum_B p(A, B|C)$$

Inference

= ?

Inference: Derivation of General Filtering Equations

Model



$$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t | x_{t-1}) p(y_t | x_t)$$

$$p(x_t | y_{1:t})$$

Rules of probability

product rule

$$p(A|B, C) = \frac{1}{p(B|C)} p(B|A, C) p(A|C)$$

sum rule

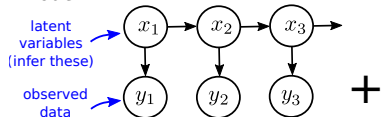
$$p(A|C) = \sum_B p(A, B|C)$$

Inference

= ?

Inference: Derivation of General Filtering Equations

Model



$$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t | x_{t-1}) p(y_t | x_t)$$

$$p(x_t | y_{1:t}) = p(x_t | y_t, y_{1:t-1})$$

Rules of probability

product rule

$$p(A|B, C) = \frac{1}{p(B|C)} p(B|A, C) p(A|C)$$

sum rule

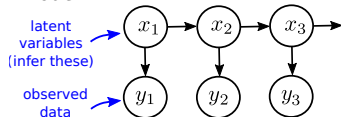
$$p(A|C) = \sum_B p(A, B|C)$$

Inference

= ?

Inference: Derivation of General Filtering Equations

Model



$$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t | x_{t-1}) p(y_t | x_t)$$

$$p(x_t | y_{1:t}) = p(x_t | y_t, y_{1:t-1})$$

$$= \frac{1}{p(y_t | y_{1:t-1})} p(y_t | x_t, y_{1:t-1}) p(x_t | y_{1:t-1})$$

product rule

$$A = x_t \quad B = y_t \quad C = y_{1:t-1}$$

Rules of probability

product rule

$$p(A|B, C) = \frac{1}{p(B|C)} p(B|A, C) p(A|C)$$

sum rule

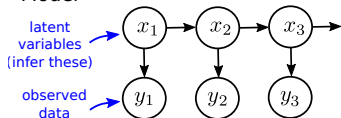
$$p(A|C) = \sum_B p(A, B|C)$$

Inference

= ?

Inference: Derivation of General Filtering Equations

Model



+

$$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t | x_{t-1}) p(y_t | x_t)$$

Rules of probability

product rule

$$p(A|B, C) = \frac{1}{p(B|C)} p(B|A, C) p(A|C)$$

sum rule

$$p(A|C) = \sum_B p(A, B|C)$$

Inference

= ?

$$p(x_t | y_{1:t}) = p(x_t | y_t, y_{1:t-1})$$

$$= \frac{1}{p(y_t | y_{1:t-1})} p(y_t | x_t, y_{1:t-1}) p(x_t | y_{1:t-1})$$

product rule

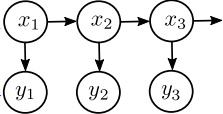
$$A = x_t \quad B = y_t \quad C = y_{1:t-1}$$

$$= \frac{1}{p(y_t | y_{1:t-1})} p(y_t | x_t) p(x_t | y_{1:t-1})$$

conditional independence from model

$$y_t \perp y_{1:t-1} | x_t$$

Inference: Derivation of General Filtering Equations

Model	Rules of probability	Inference
 <p>latent variables (infer these)</p> <p>observed data</p>	<p>product rule</p> $p(A B, C) = \frac{1}{p(B C)} p(B A, C)p(A C)$ <p>sum rule</p> $p(A C) = \sum_B p(A, B C)$	= ?
$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t x_{t-1})p(y_t x_t)$		
$\begin{aligned} p(x_t y_{1:t}) &= p(x_t y_t, y_{1:t-1}) \\ &= \frac{1}{p(y_t y_{1:t-1})} p(y_t x_t, y_{1:t-1})p(x_t y_{1:t-1}) \\ &= \frac{1}{p(y_t y_{1:t-1})} p(y_t x_t)p(x_t y_{1:t-1}) \\ &\propto p(y_t x_t)p(x_t y_{1:t-1}) \end{aligned}$	<p>product rule</p> $A = x_t \quad B = y_t \quad C = y_{1:t-1}$ <p>conditional independence from model</p> $y_t \perp y_{1:t-1} x_t$ <p>constant of proportionality $p(y_t y_{1:t-1})$ (see learning)</p>	

Inference: Derivation of General Filtering Equations

Model

latent variables (infer these)

observed data

Rules of probability

product rule

$$p(A|B, C) = \frac{1}{p(B|C)} p(B|A, C)p(A|C)$$

sum rule

$$p(A|C) = \sum_B p(A, B|C)$$

Inference

$$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t|x_{t-1})p(y_t|x_t) \quad + \quad = ?$$
$$p(x_t|y_{1:t}) = p(x_t|y_t, y_{1:t-1})$$
$$= \frac{1}{p(y_t|y_{1:t-1})} p(y_t|x_t, y_{1:t-1})p(x_t|y_{1:t-1})$$

product rule $A = x_t \ B = y_t \ C = y_{1:t-1}$

$$= \frac{1}{p(y_t|y_{1:t-1})} p(y_t|x_t)p(x_t|y_{1:t-1})$$

conditional independence from model $y_t \perp y_{1:t-1} | x_t$

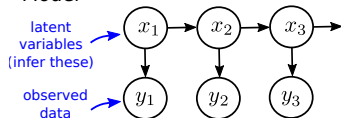
$$\propto p(y_t|x_t)p(x_t|y_{1:t-1})$$

constant of proportionality $p(y_t|y_{1:t-1})$ (see learning)

$$p(x_t|y_{1:t-1})$$

Inference: Derivation of General Filtering Equations

Model



$$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t | x_{t-1}) p(y_t | x_t)$$

Rules of probability

product rule

$$p(A|B, C) = \frac{1}{p(B|C)} p(B|A, C) p(A|C)$$

sum rule

$$p(A|C) = \sum_B p(A, B|C)$$

Inference

= ?

$$p(x_t | y_{1:t}) = p(x_t | y_t, y_{1:t-1})$$

$$= \frac{1}{p(y_t | y_{1:t-1})} p(y_t | x_t, y_{1:t-1}) p(x_t | y_{1:t-1})$$

product rule
 $A = x_t \quad B = y_t \quad C = y_{1:t-1}$

$$= \frac{1}{p(y_t | y_{1:t-1})} p(y_t | x_t) p(x_t | y_{1:t-1})$$

conditional independence from model
 $y_t \perp y_{1:t-1} | x_t$

$$\propto p(y_t | x_t) p(x_t | y_{1:t-1})$$

constant of proportionality $p(y_t | y_{1:t-1})$ (see learning)

$$p(x_t | y_{1:t-1}) = \int p(x_t, x_{t-1} | y_{1:t-1}) dx_{t-1}$$

sum rule

$$A = x_t \quad B = x_{t-1} \quad C = y_{1:t-1}$$

Inference: Derivation of General Filtering Equations

Model	Rules of probability	Inference
	<p>product rule</p>	$= ?$
$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t x_{t-1})p(y_t x_t)$	<p>sum rule</p> $p(A C) = \sum_B p(A, B C)$	
$p(x_t y_{1:t}) = p(x_t y_t, y_{1:t-1})$	<p>product rule</p>	$A = x_t \quad B = y_t \quad C = y_{1:t-1}$
$= \frac{1}{p(y_t y_{1:t-1})} p(y_t x_t, y_{1:t-1})p(x_t y_{1:t-1})$	<p>conditional independence from model</p>	$y_t \perp y_{1:t-1} x_t$
$= \frac{1}{p(y_t y_{1:t-1})} p(y_t x_t)p(x_t y_{1:t-1})$	<p>constant of proportionality $p(y_t y_{1:t-1})$ (see learning)</p>	
$p(x_t y_{1:t-1}) = \int p(x_t, x_{t-1} y_{1:t-1})dx_{t-1}$	<p>sum rule</p>	$A = x_t \quad B = x_{t-1} \quad C = y_{1:t-1}$
$= \int p(x_t x_{t-1}, y_{1:t-1})p(x_{t-1} y_{1:t-1})dx_{t-1}$	<p>product rule</p>	

Inference: Derivation of General Filtering Equations

Model	Rules of probability	Inference
	<p>product rule</p> $p(A B, C) = \frac{1}{p(B C)} p(B A, C)p(A C)$	$= ?$
$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t x_{t-1})p(y_t x_t)$	<p>sum rule</p> $p(A C) = \sum_B p(A, B C)$	
$p(x_t y_{1:t}) = p(x_t y_t, y_{1:t-1})$	<p>product rule</p> $A = x_t \quad B = y_t \quad C = y_{1:t-1}$	
$= \frac{1}{p(y_t y_{1:t-1})} p(y_t x_t, y_{1:t-1})p(x_t y_{1:t-1})$	<p>conditional independence from model</p> $y_t \perp y_{1:t-1} x_t$	
$= \frac{1}{p(y_t y_{1:t-1})} p(y_t x_t)p(x_t y_{1:t-1})$	<p>constant of proportionality $p(y_t y_{1:t-1})$ (see learning)</p>	
$p(x_t y_{1:t-1}) = \int p(x_t, x_{t-1} y_{1:t-1})dx_{t-1}$	<p>sum rule</p> $A = x_t \quad B = x_{t-1} \quad C = y_{1:t-1}$	
$= \int p(x_t x_{t-1}, y_{1:t-1})p(x_{t-1} y_{1:t-1})dx_{t-1}$	<p>product rule</p>	
$= \int p(x_t x_{t-1})p(x_{t-1} y_{1:t-1})dx_{t-1}$	<p>conditional independence from model</p>	

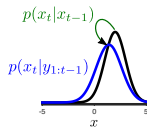
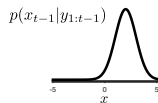
Inference: Kalman Filter

$$p(x_{t-1} | y_{1:t-1})$$

diffuse via
dynamics



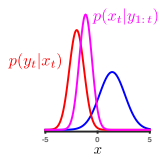
$$p(x_t | y_{1:t-1}) = \int p(x_t | x_{t-1}) p(x_{t-1} | y_{1:t-1}) dx_{t-1}$$



combine
with
likelihood



$$p(x_t | y_{1:t}) \propto \underbrace{p(x_t | y_{1:t-1})}_{\text{prior}} \underbrace{p(y_t | x_t)}_{\text{likelihood}}$$



Inference: Kalman Filter

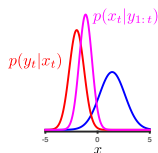
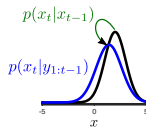
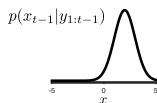
$$p(x_{t-1}|y_{1:t-1}) = \mathcal{G}(x_{t-1}; \mu_{t-1}^{t-1}, V_{t-1}^{t-1})$$

← most recent data used in prediction

← variable being predicted

diffuse via dynamics

$$p(x_t|y_{1:t-1}) = \int p(x_t|x_{t-1})p(x_{t-1}|y_{1:t-1})dx_{t-1}$$



combine with likelihood

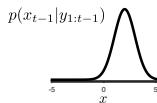
$$p(x_t|y_{1:t}) \propto \underbrace{p(x_t|y_{1:t-1})}_{\text{prior}} \underbrace{p(y_t|x_t)}_{\text{likelihood}}$$

Inference: Kalman Filter

$$p(x_{t-1}|y_{1:t-1}) = \mathcal{G}(x_{t-1}; \mu_{t-1}^{t-1}, V_{t-1}^{t-1})$$

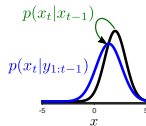
← most recent data used in prediction

← variable being predicted



diffuse via dynamics

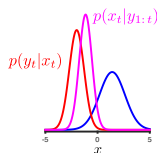
$$p(x_t|y_{1:t-1}) = \int p(x_t|x_{t-1})p(x_{t-1}|y_{1:t-1})dx_{t-1}$$



$$p(x_t|y_{1:t-1}) = \mathcal{G}(x_t; \mu_t^{t-1}, V_t^{t-1})$$

combine with likelihood

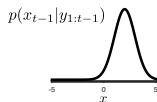
$$p(x_t|y_{1:t}) \propto \underbrace{p(x_t|y_{1:t-1})}_{\text{prior}} \underbrace{p(y_t|x_t)}_{\text{likelihood}}$$



Inference: Kalman Filter

$$p(x_{t-1}|y_{1:t-1}) = \mathcal{G}(x_{t-1}; \mu_{t-1}^{t-1}, V_{t-1}^{t-1})$$

← most recent data used in prediction
 ← variable being predicted

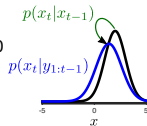


diffuse via dynamics



$$p(x_t|y_{1:t-1}) = \int p(x_t|x_{t-1})p(x_{t-1}|y_{1:t-1})dx_{t-1}$$

diffuses toward 0



$$p(x_t|y_{1:t-1}) = \mathcal{G}(x_t; \mu_t^{t-1}, V_t^{t-1}) \quad \mu_t^{t-1} = A\mu_{t-1}^{t-1}$$

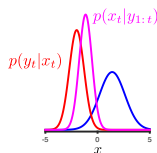
$$V_t^{t-1} = AV_{t-1}^{t-1}A^T + Q$$

combine with likelihood



$$p(x_t|y_{1:t}) \propto \underbrace{p(x_t|y_{1:t-1})}_{\text{prior}} \underbrace{p(y_t|x_t)}_{\text{likelihood}}$$

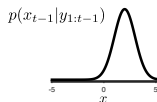
variance inflates



Inference: Kalman Filter

$$p(x_{t-1}|y_{1:t-1}) = \mathcal{G}(x_{t-1}; \mu_{t-1}^{t-1}, V_{t-1}^{t-1})$$

← most recent data used in prediction
 ← variable being predicted

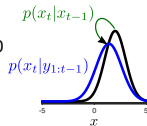


diffuse via dynamics



$$p(x_t|y_{1:t-1}) = \int p(x_t|x_{t-1})p(x_{t-1}|y_{1:t-1})dx_{t-1}$$

diffuses toward 0



$$p(x_t|y_{1:t-1}) = \mathcal{G}(x_t; \mu_t^{t-1}, V_t^{t-1}) \quad \mu_t^{t-1} = A\mu_{t-1}^{t-1}$$

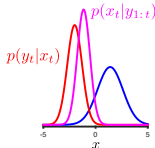
$$V_t^{t-1} = AV_{t-1}^{t-1}A^T + Q$$

combine with likelihood



$$p(x_t|y_{1:t}) \propto \underbrace{p(x_t|y_{1:t-1})}_{\text{prior}} \underbrace{p(y_t|x_t)}_{\text{likelihood}}$$

variance inflates

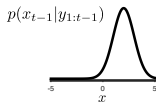


$$p(x_t|y_{1:t}) = \mathcal{G}(x_t; \mu_t^t, V_t^t)$$

Inference: Kalman Filter

$$p(x_{t-1}|y_{1:t-1}) = \mathcal{G}(x_{t-1}; \mu_{t-1}^{t-1}, V_{t-1}^{t-1})$$

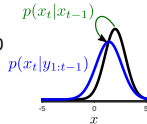
← most recent data used in prediction
 ← variable being predicted



diffuse via dynamics

$$p(x_t|y_{1:t-1}) = \int p(x_t|x_{t-1})p(x_{t-1}|y_{1:t-1})dx_{t-1}$$

diffuses toward 0



$$p(x_t|y_{1:t-1}) = \mathcal{G}(x_t; \mu_t^{t-1}, V_t^{t-1}) \quad \mu_t^{t-1} = A\mu_{t-1}^{t-1}$$

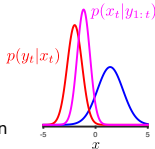
$$V_t^{t-1} = AV_{t-1}^{t-1}A^T + Q$$

combine with likelihood

$$p(x_t|y_{1:t}) \propto p(x_t|y_{1:t-1})p(y_t|x_t)$$

prior
likelihood

variance inflates



$$p(x_t|y_{1:t}) = \mathcal{G}(x_t; \mu_t^t, V_t^t)$$

$$\mu_t^t = \mu_t^{t-1} + K_t(y_t - C\mu_t^{t-1})$$

$$V_t^t = V_t^{t-1} - K_tCV_t^{t-1}$$

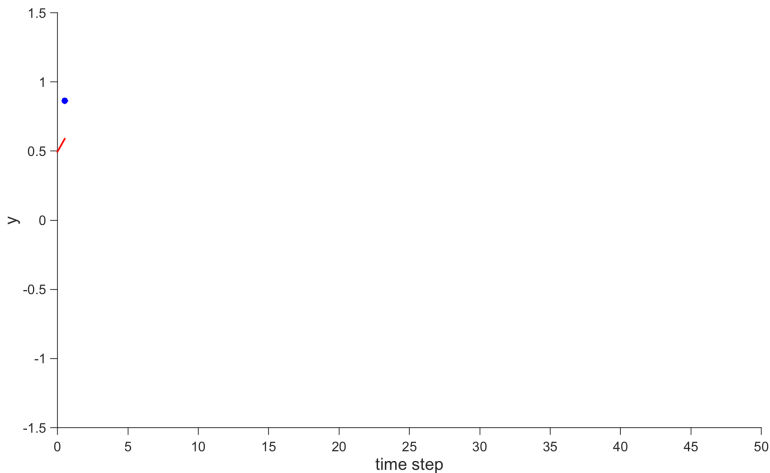
$$\text{Kalman gain} \rightarrow K_t = V_t^{t-1}C^T(CV_t^{t-1}C^T + R)^{-1}$$

Kalman Filter Demo

- ▶ data: $y_t = \sin(\omega t) + \sigma_y \epsilon_t$ where $\sigma_y^2 = 0.1$
- ▶ model: $x_t = \lambda x_{t-1} + \sigma \eta$ and $y_t = x_t + \sigma_y \eta'_t$
where $\lambda = 0.99$ and $\sigma^2 = 1 - \lambda^2$
- ▶ demo shows how the Kalman filter processes the data to form estimates of the hidden state at each time point $p(x_t | y_{1:t})$

Kalman Filter Demo

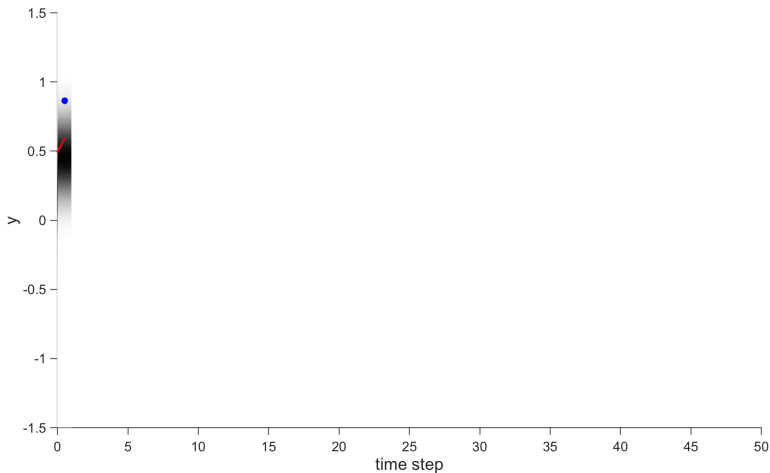
observed noisy data y_t , ground truth sinusoid



observe first data point y_1

Kalman Filter Demo

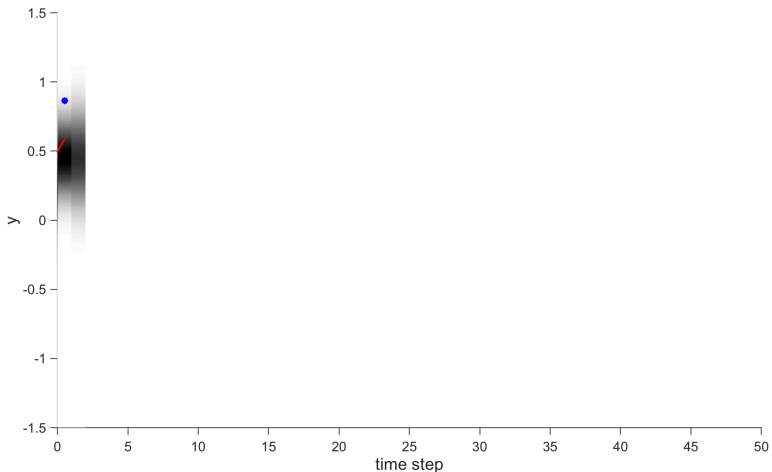
observed noisy data y_t , ground truth sinusoid



posterior over first latent variable $p(x_1|y_1)$

Kalman Filter Demo

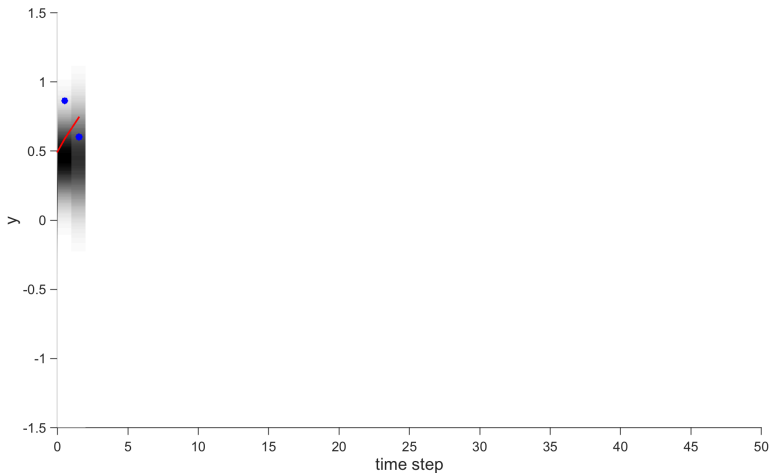
observed noisy data y_t , ground truth sinusoid



prediction for second latent variable $p(x_2|y_1)$

Kalman Filter Demo

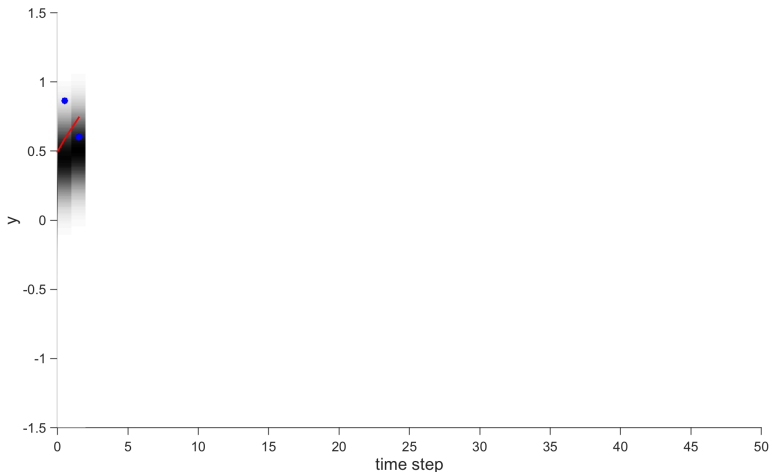
observed noisy data y_t , ground truth sinusoid



observe next data point y_2

Kalman Filter Demo

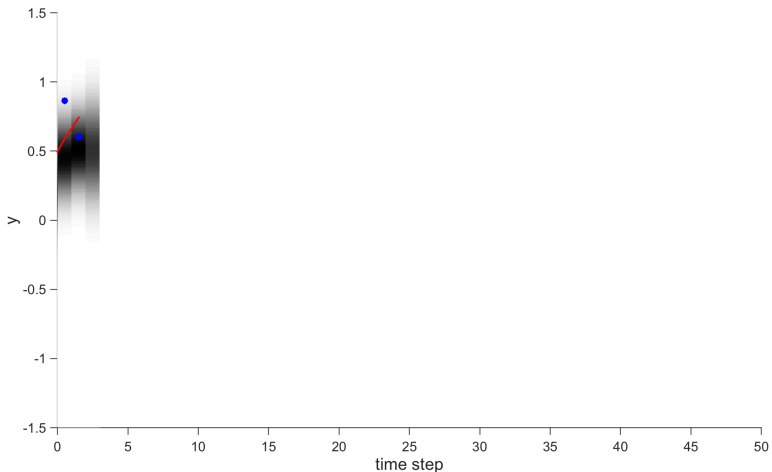
observed noisy data y_t , ground truth sinusoid



form posterior over second latent variable $p(x_2|y_1, y_2)$

Kalman Filter Demo

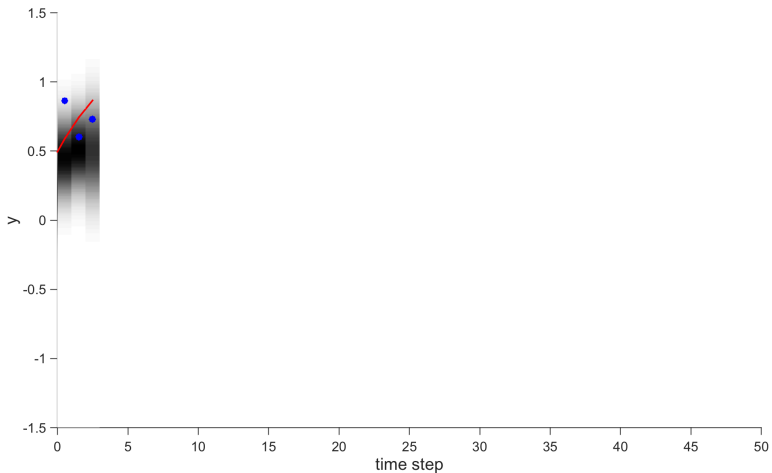
observed noisy data y_t , ground truth sinusoid



prediction for third latent variable $p(x_3|y_1, y_2)$

Kalman Filter Demo

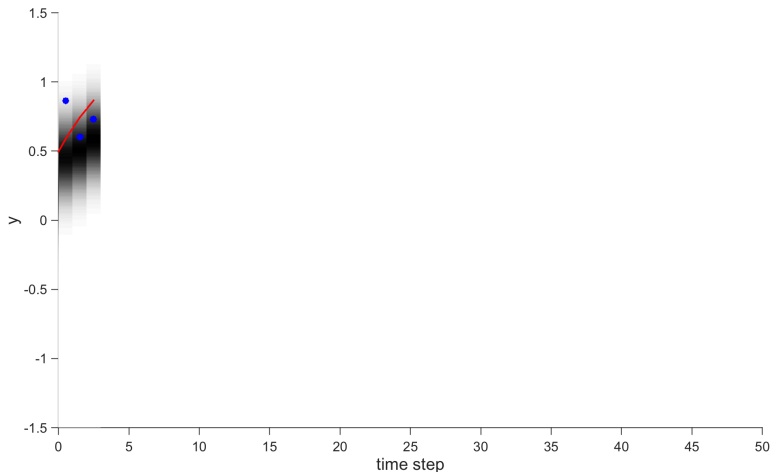
observed noisy data y_t , ground truth sinusoid



observe next data point y_3

Kalman Filter Demo

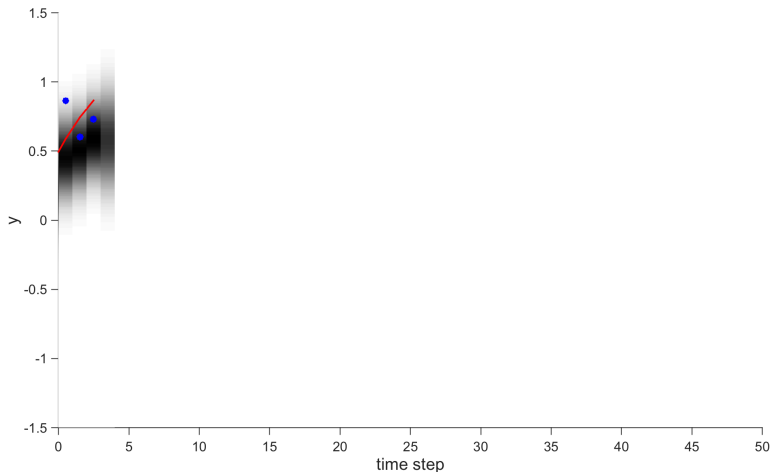
observed noisy data y_t , ground truth sinusoid



form posterior over third latent variable $p(x_3 | y_1, y_2, y_3)$

Kalman Filter Demo

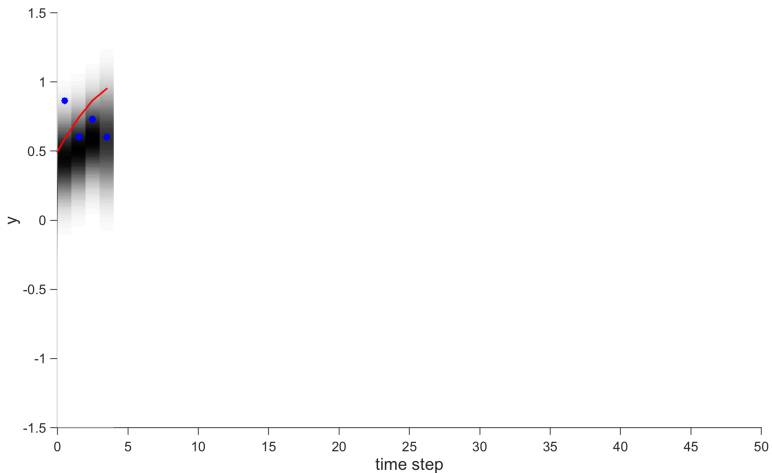
observed noisy data y_t , ground truth sinusoid



prediction for fourth latent variable $p(x_4|y_{1:3})$

Kalman Filter Demo

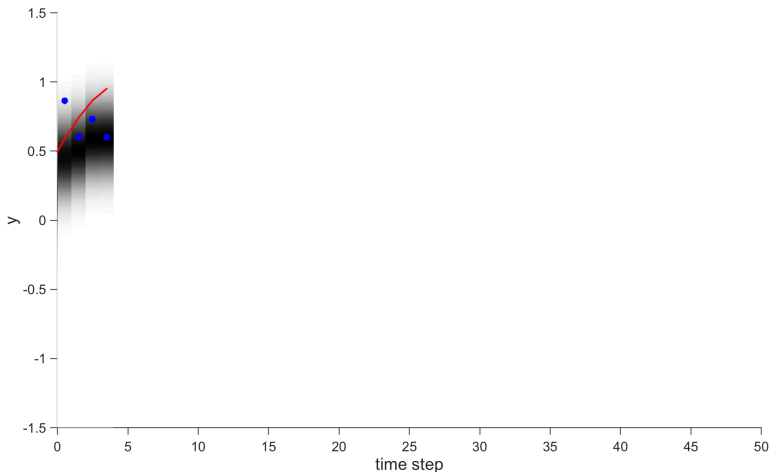
observed noisy data y_t , ground truth sinusoid



observe next data point y_4

Kalman Filter Demo

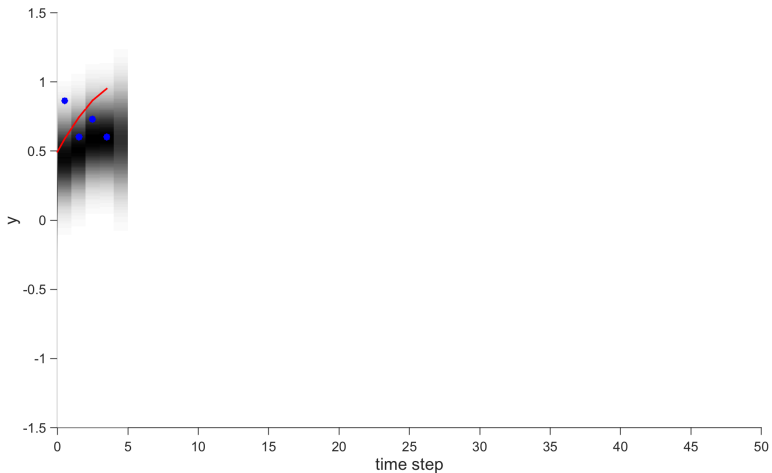
observed noisy data y_t , ground truth sinusoid



form posterior over fourth latent variable $p(x_4|y_{1:4})$

Kalman Filter Demo

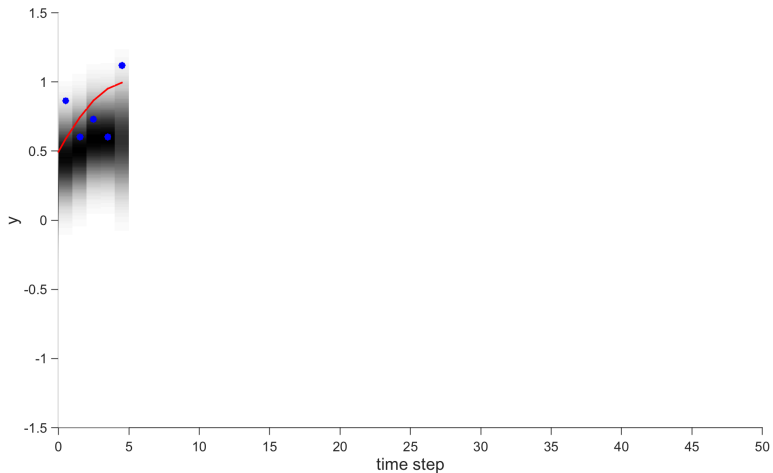
observed noisy data y_t , ground truth sinusoid



prediction for fifth latent variable $p(x_5 | y_{1:4})$

Kalman Filter Demo

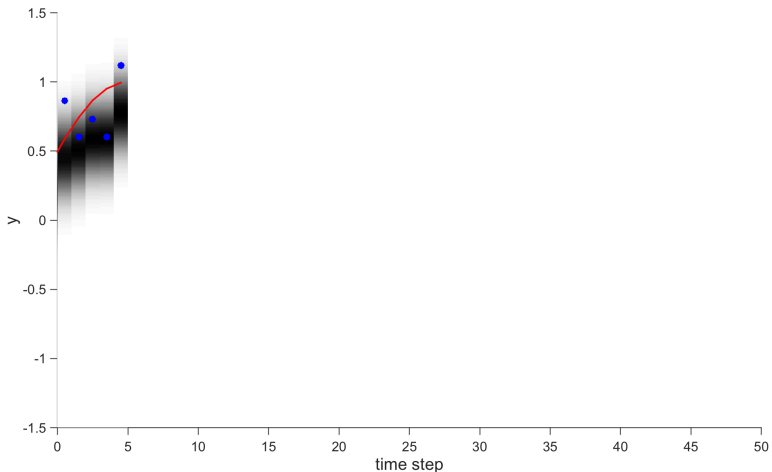
observed noisy data y_t , ground truth sinusoid



observe next data point y_5

Kalman Filter Demo

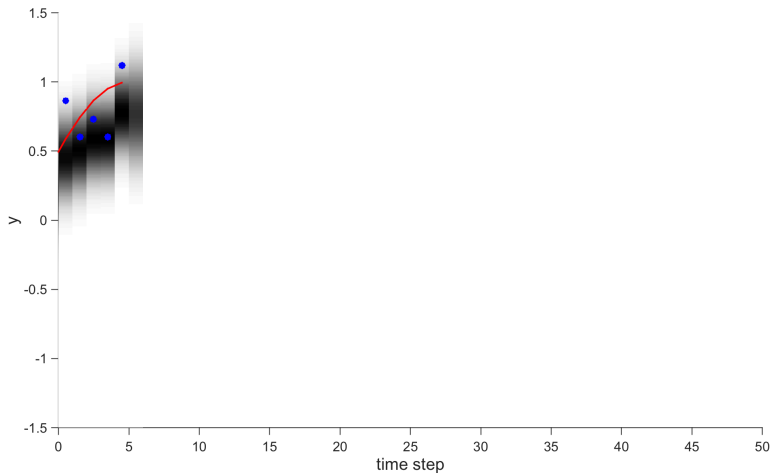
observed noisy data y_t , ground truth sinusoid



form posterior over fifth latent variable $p(x_5 | y_{1:5})$

Kalman Filter Demo

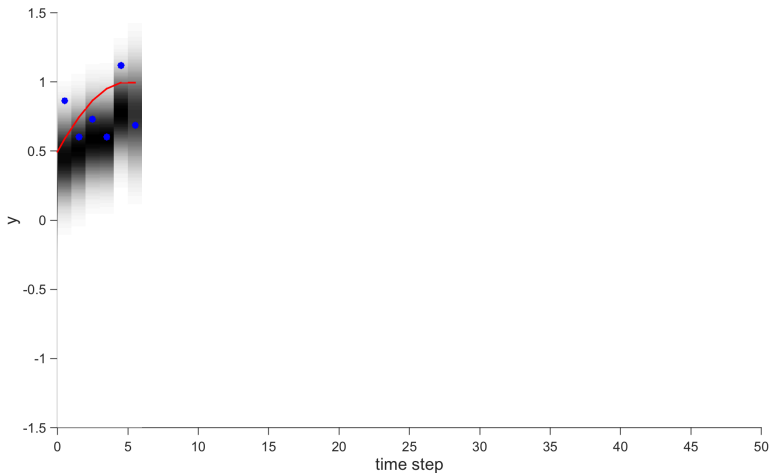
observed noisy data y_t , ground truth sinusoid



prediction for sixth latent variable $p(x_6|y_{1:5})$

Kalman Filter Demo

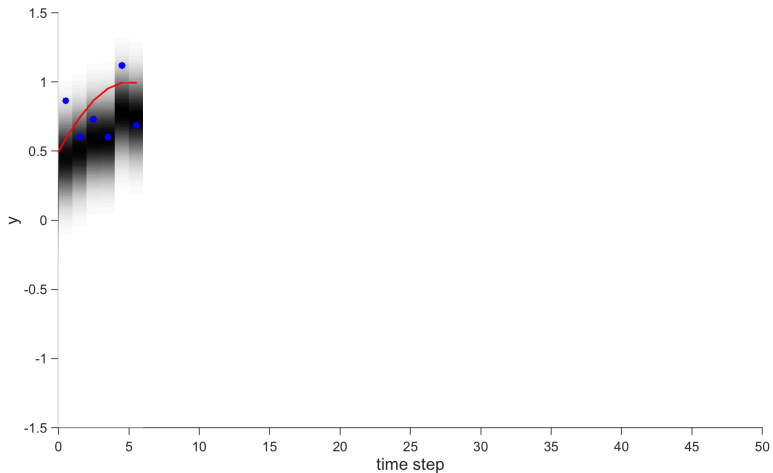
observed noisy data y_t , ground truth sinusoid



observe next data point y_6

Kalman Filter Demo

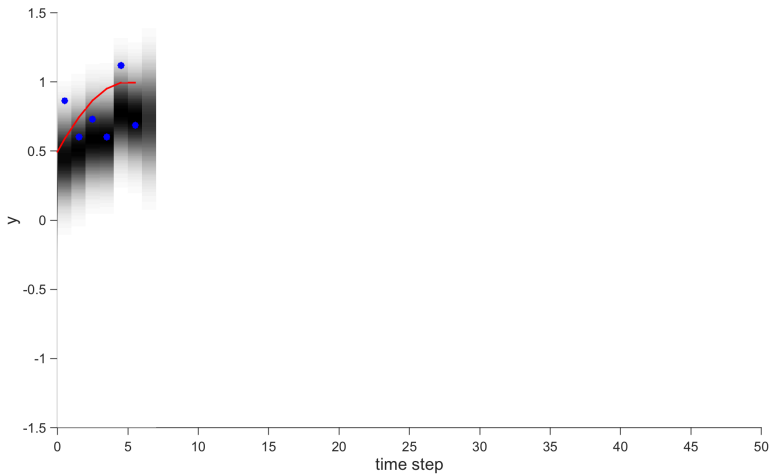
observed noisy data y_t , ground truth sinusoid



form posterior over sixth latent variable $p(x_6|y_{1:6})$

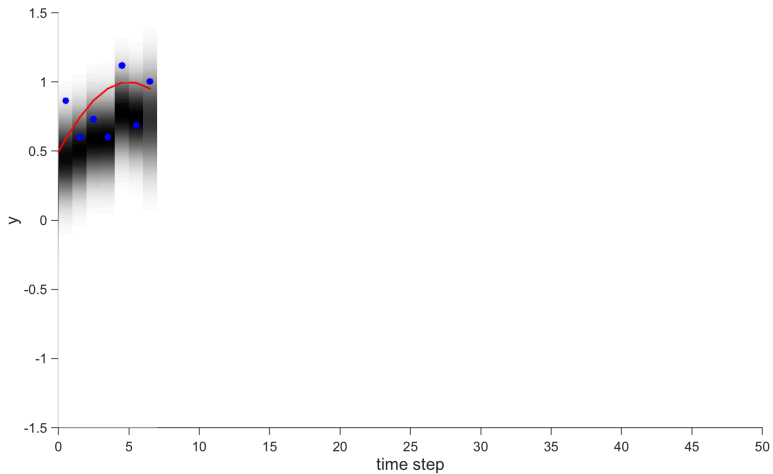
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



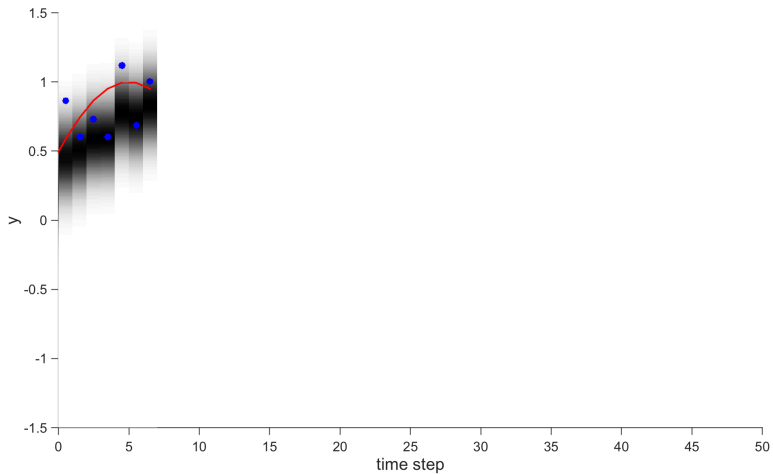
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



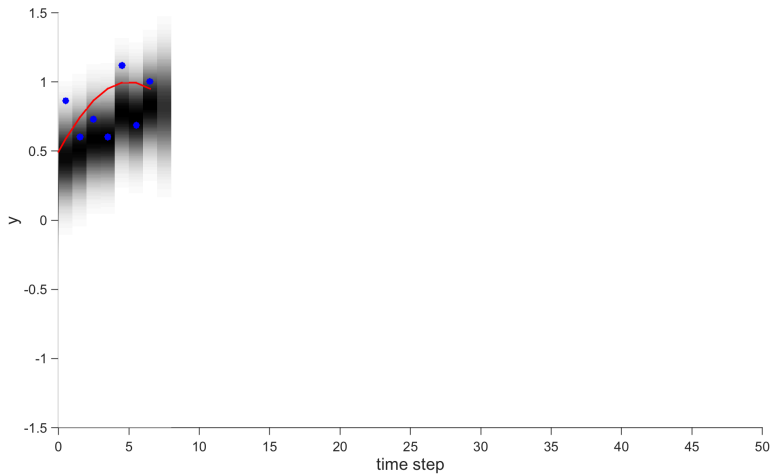
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



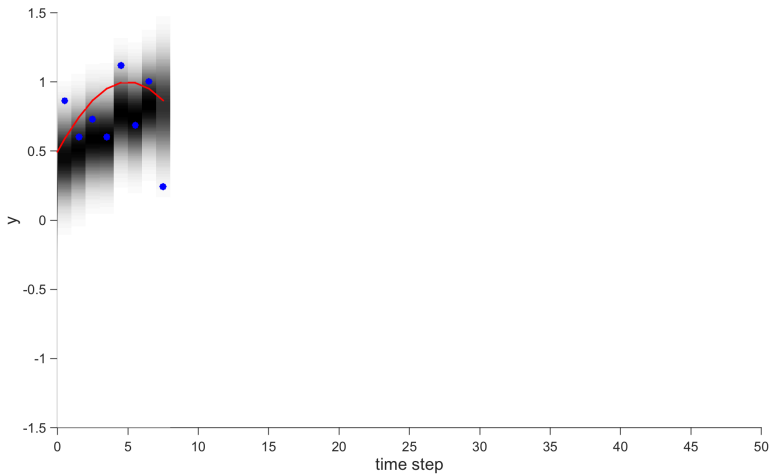
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



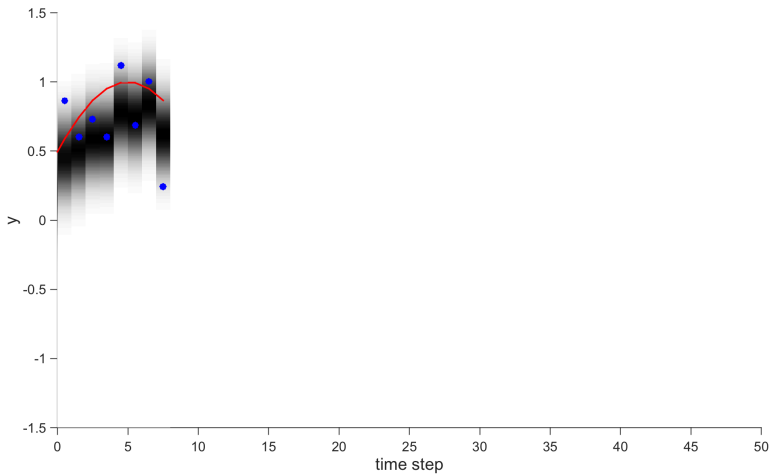
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



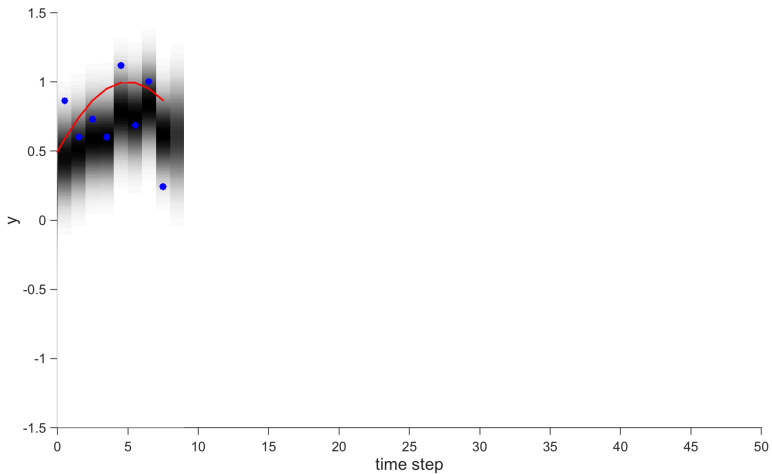
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



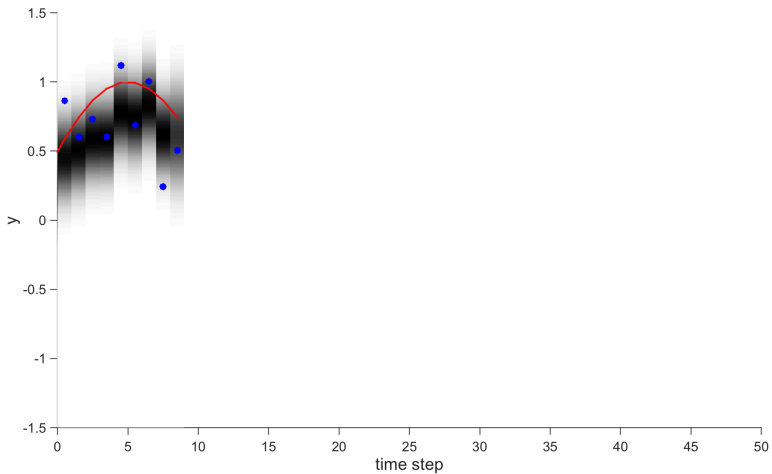
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



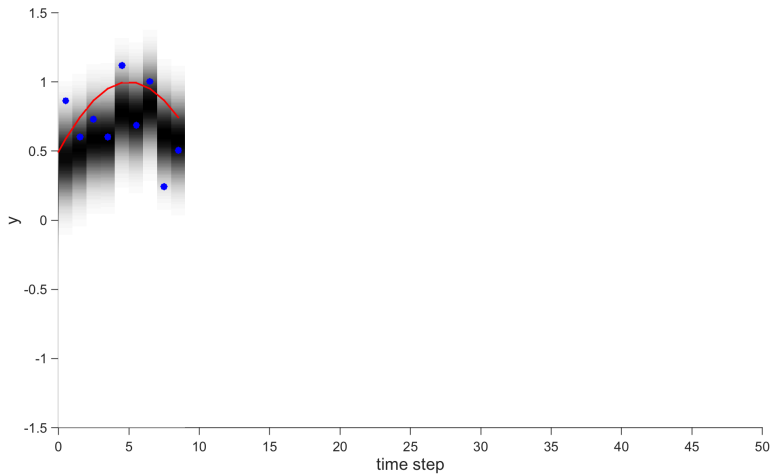
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



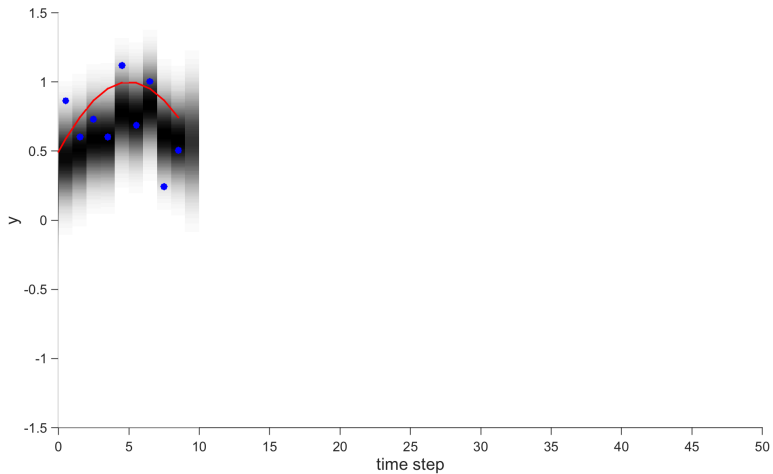
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



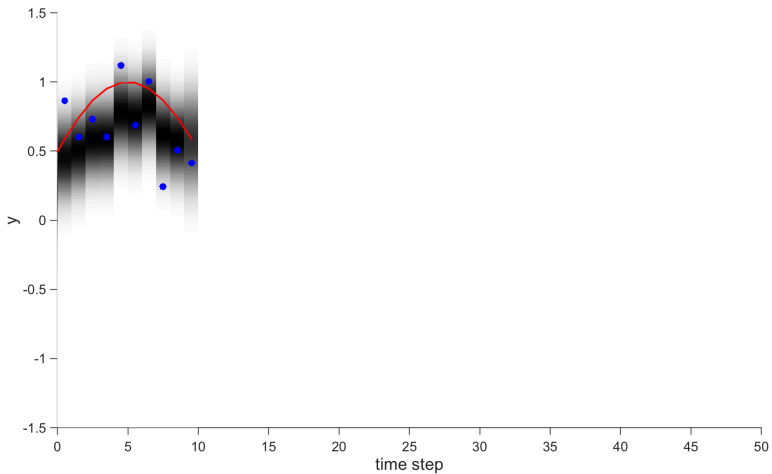
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



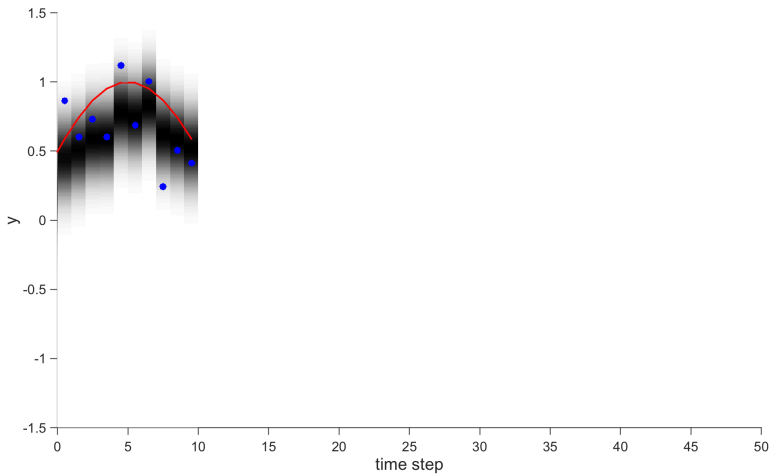
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



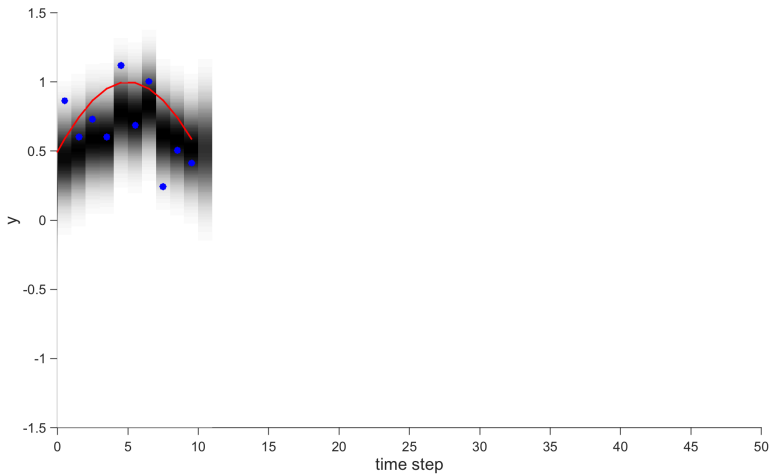
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



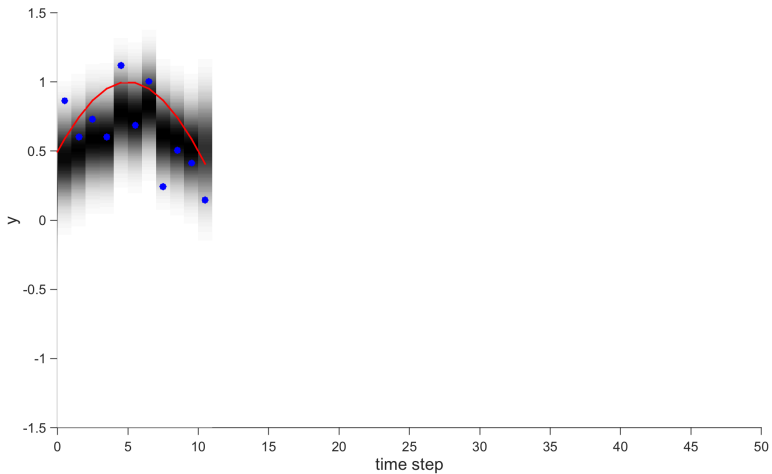
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



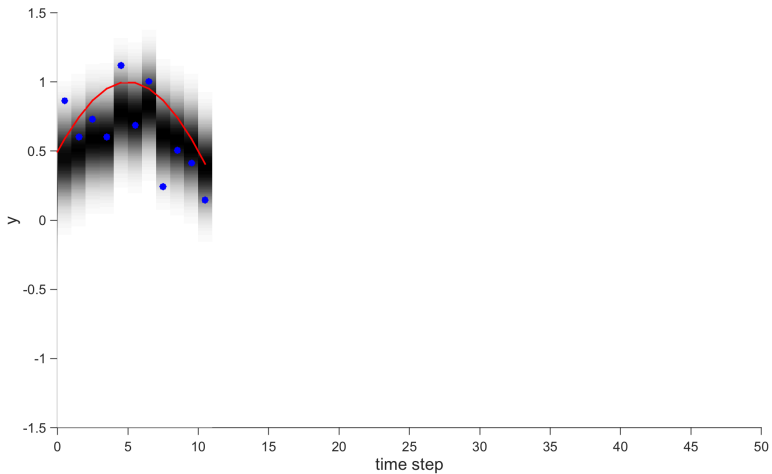
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



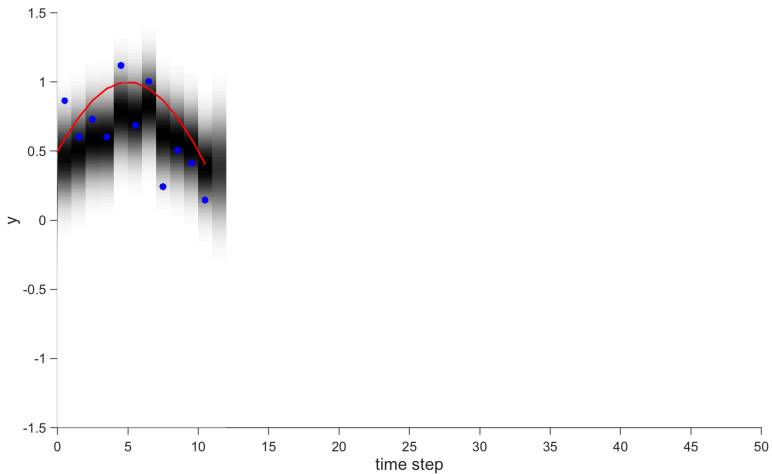
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



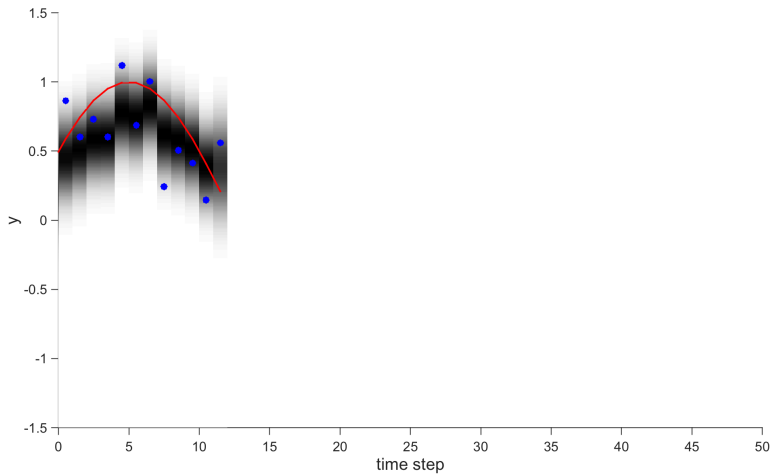
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



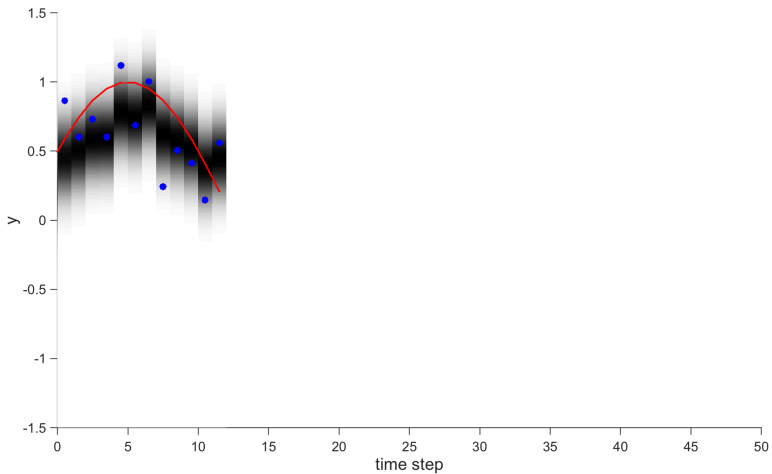
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



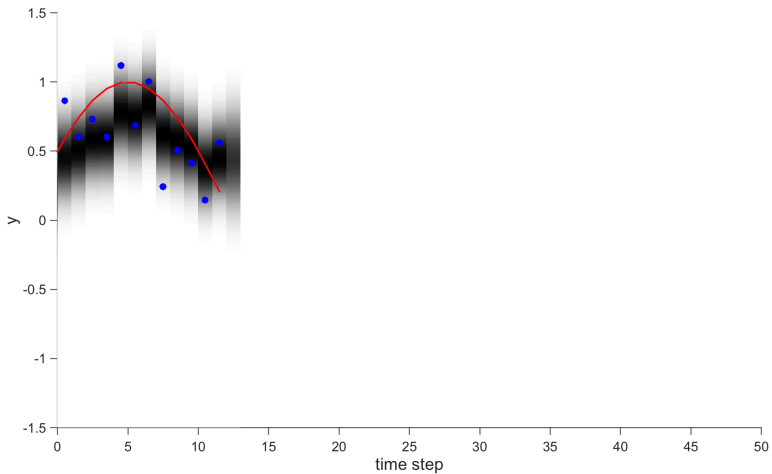
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



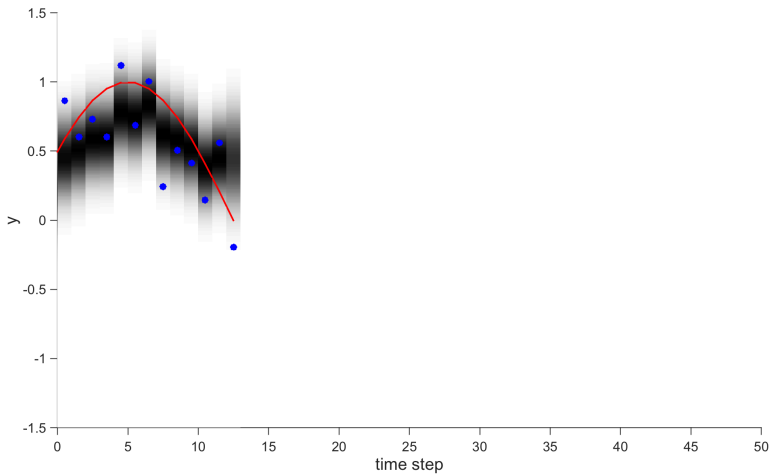
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



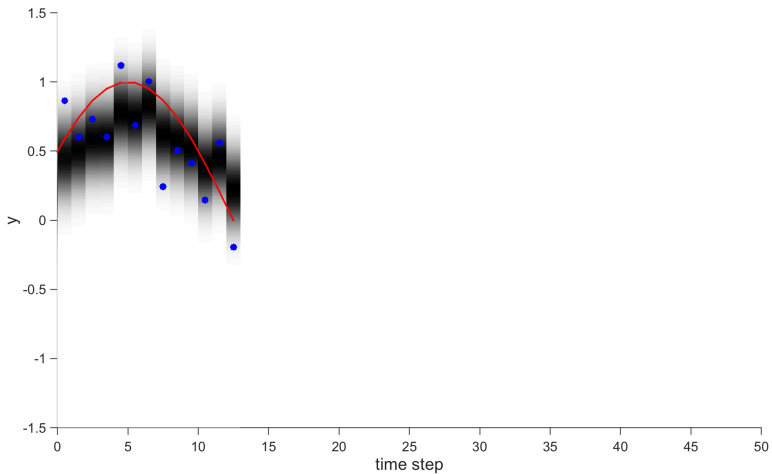
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



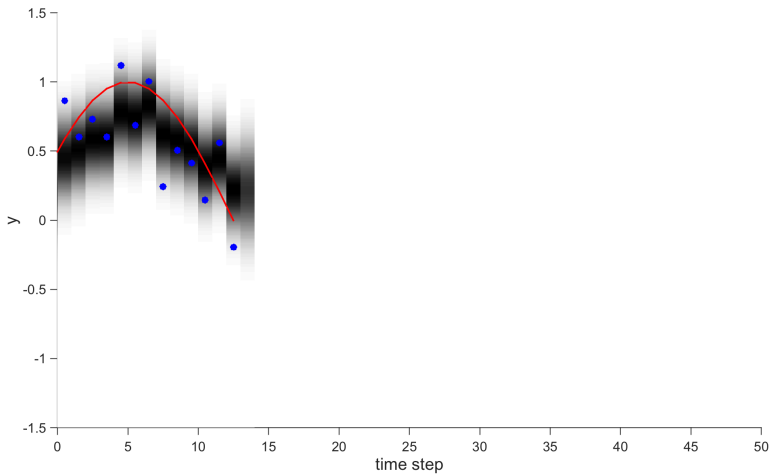
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



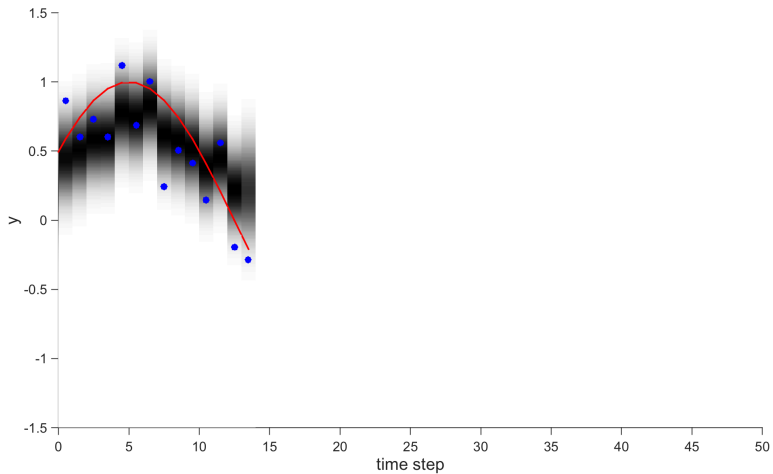
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



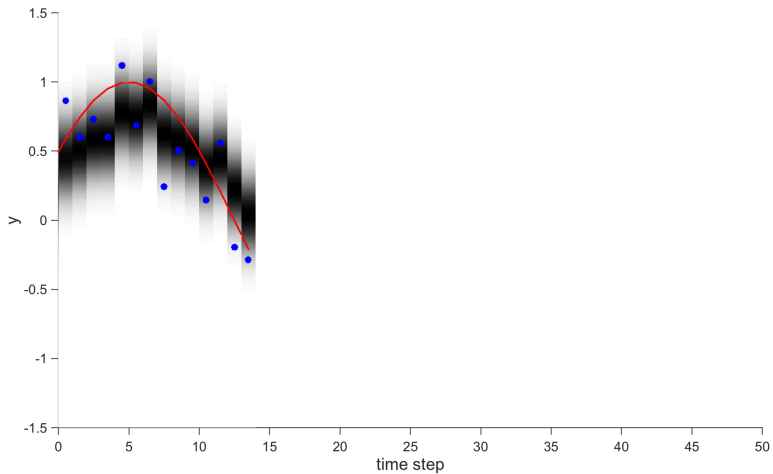
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



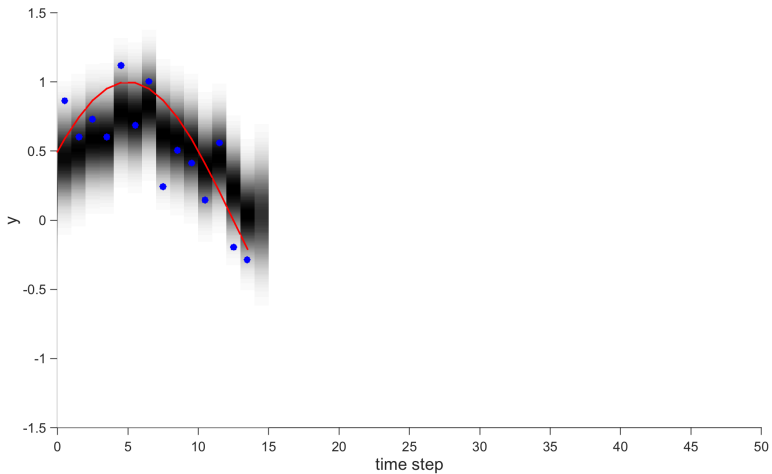
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



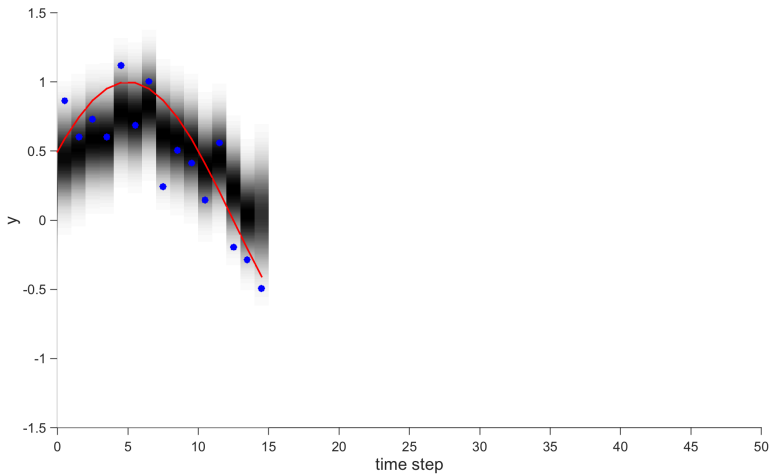
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



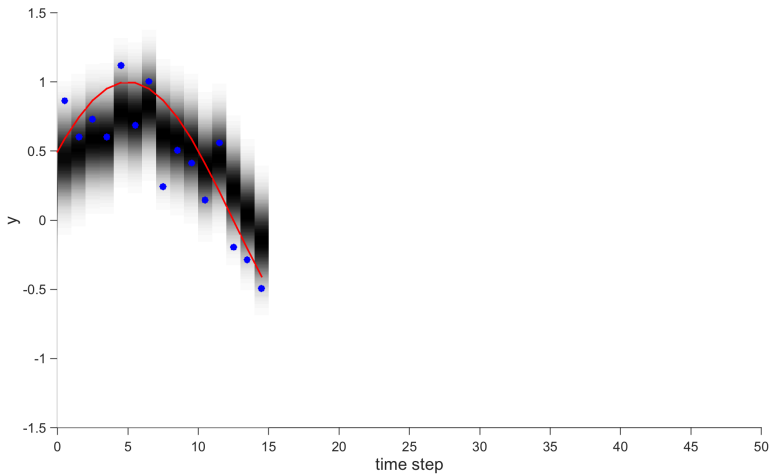
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



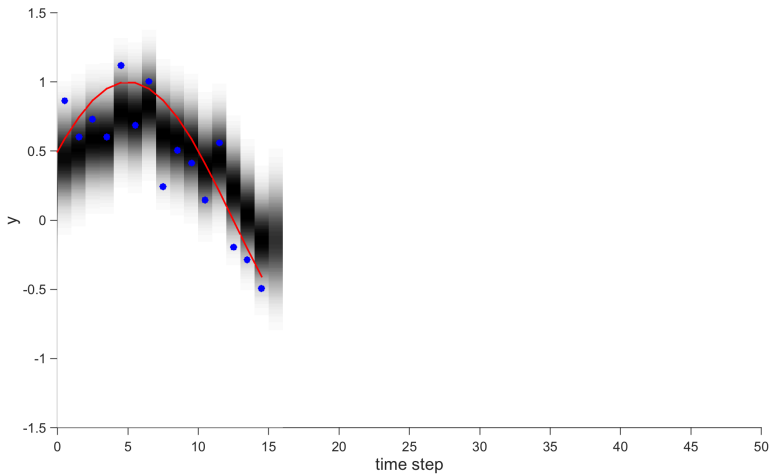
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



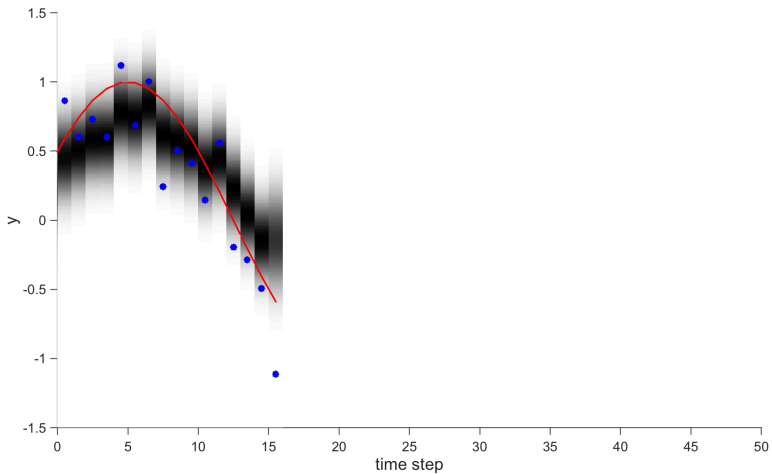
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



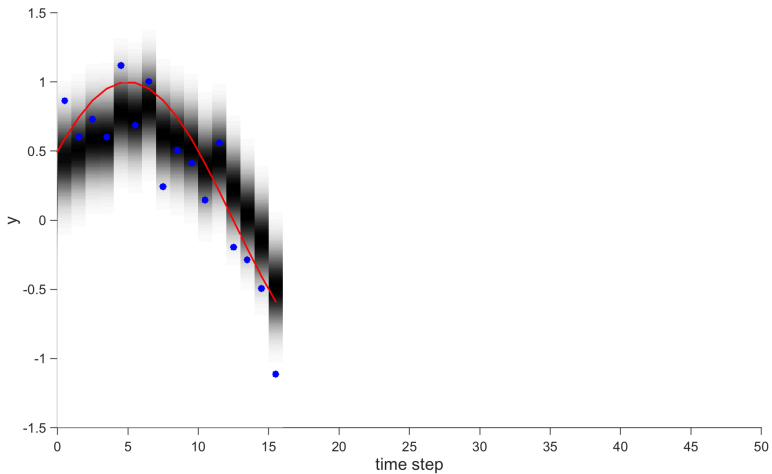
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



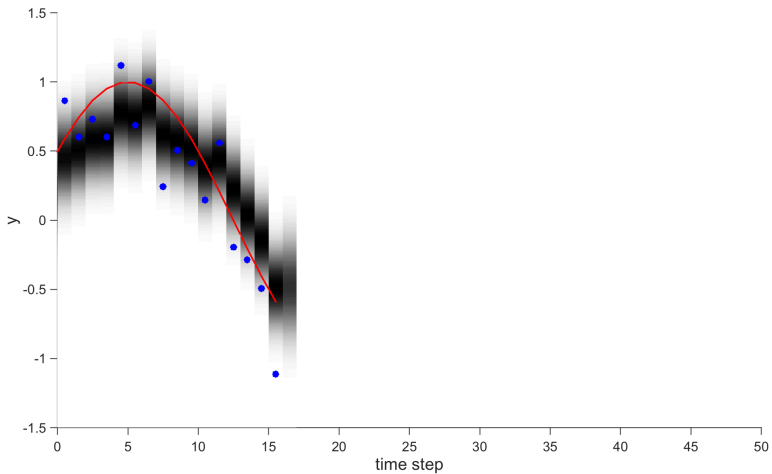
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



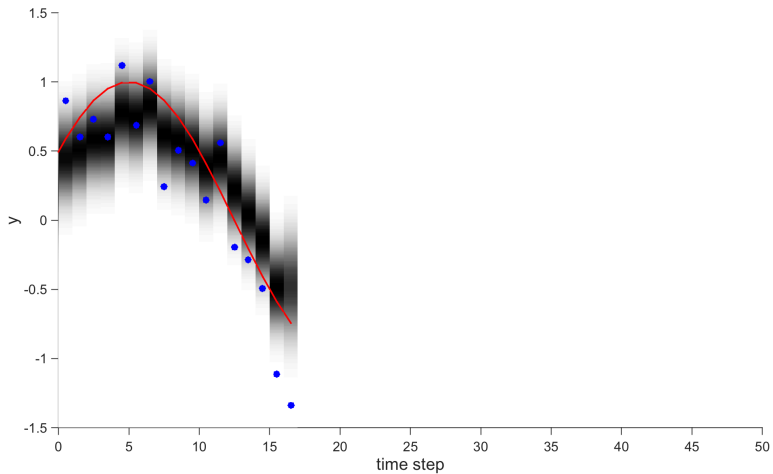
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



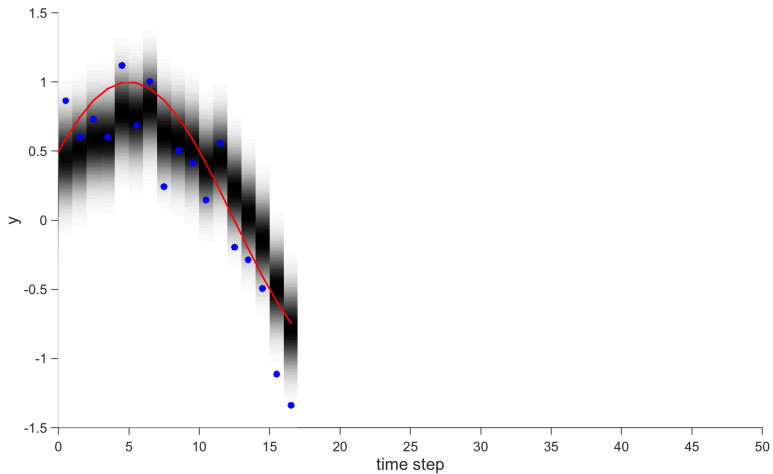
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



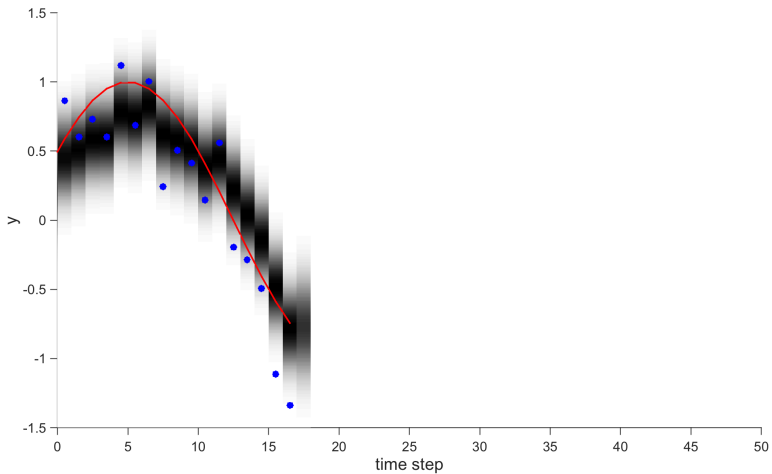
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



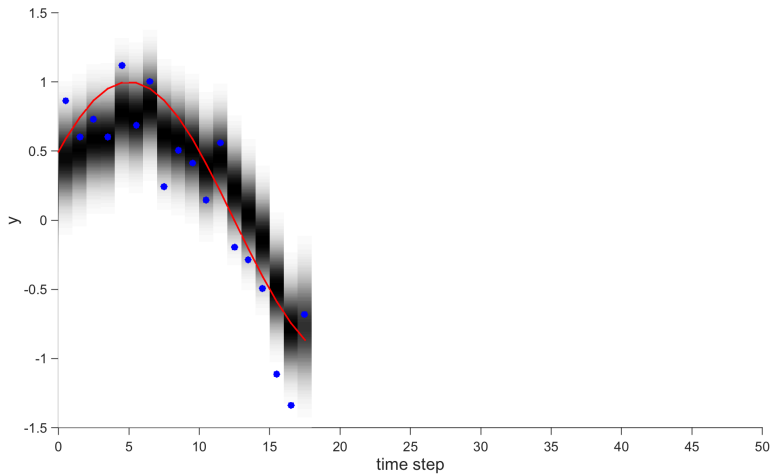
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



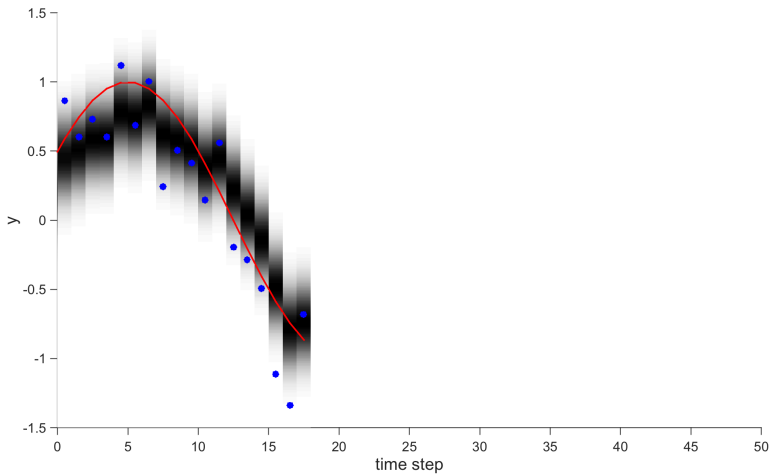
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



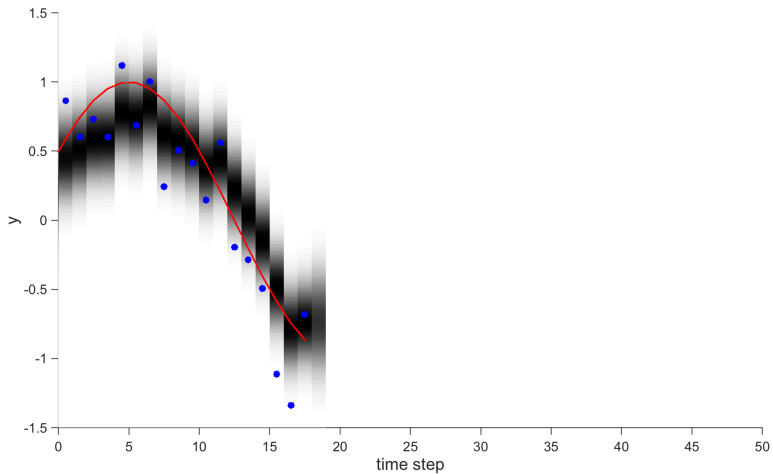
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



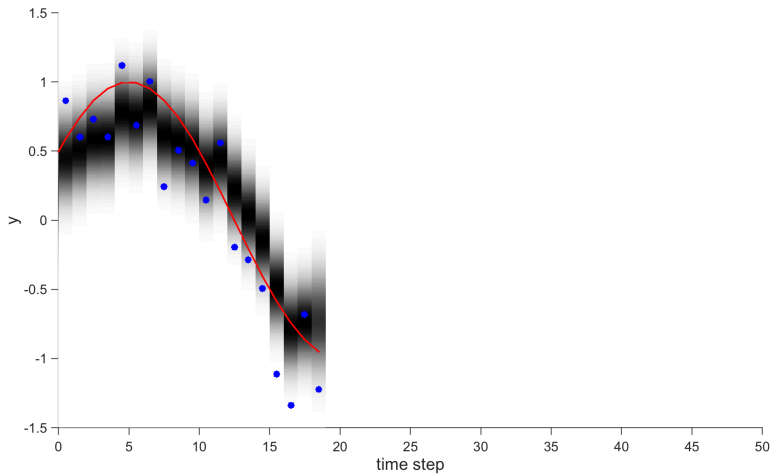
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



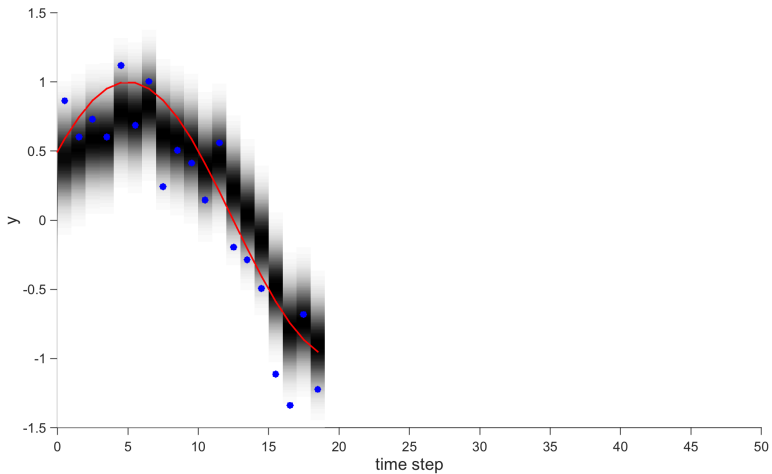
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



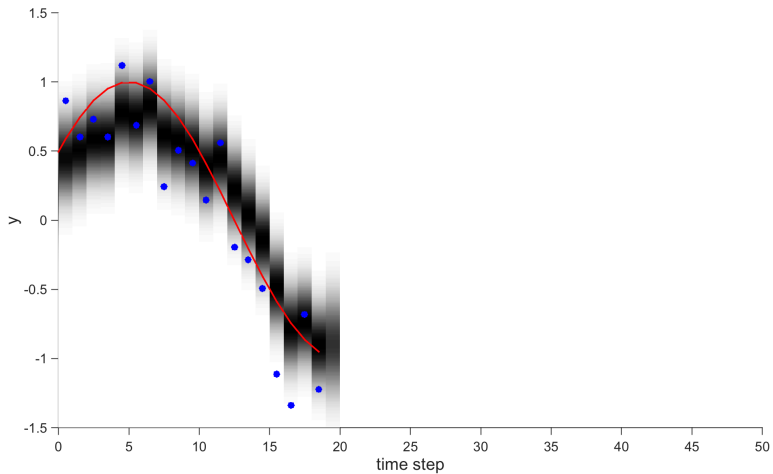
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



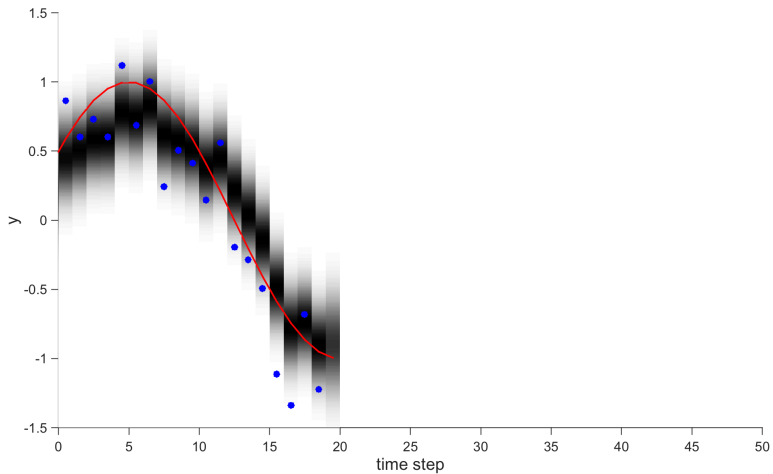
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



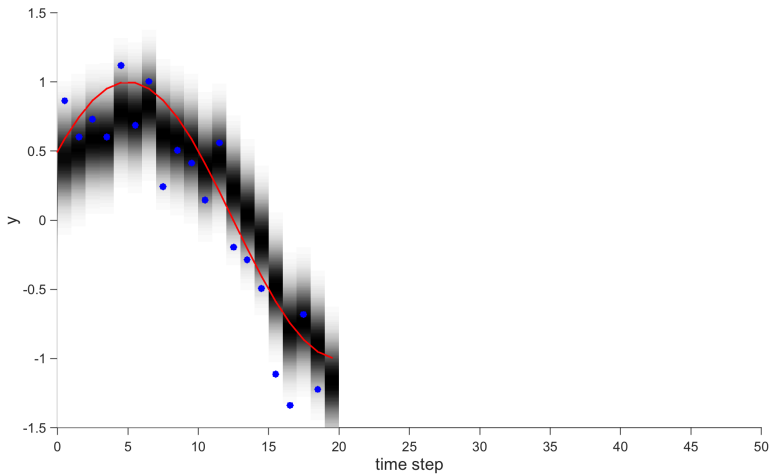
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



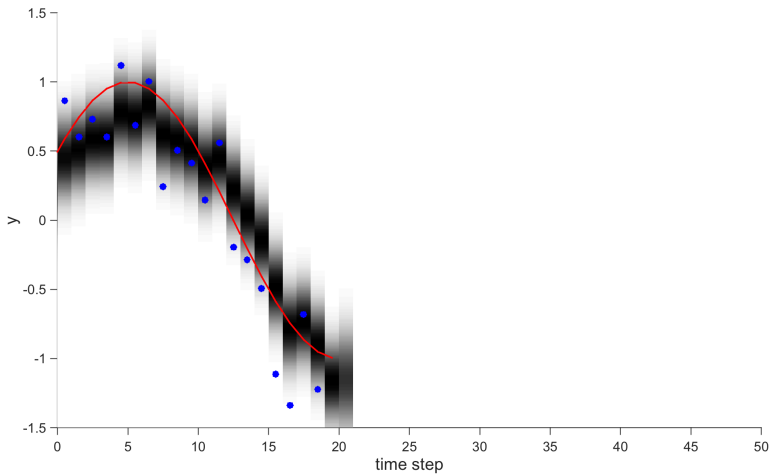
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



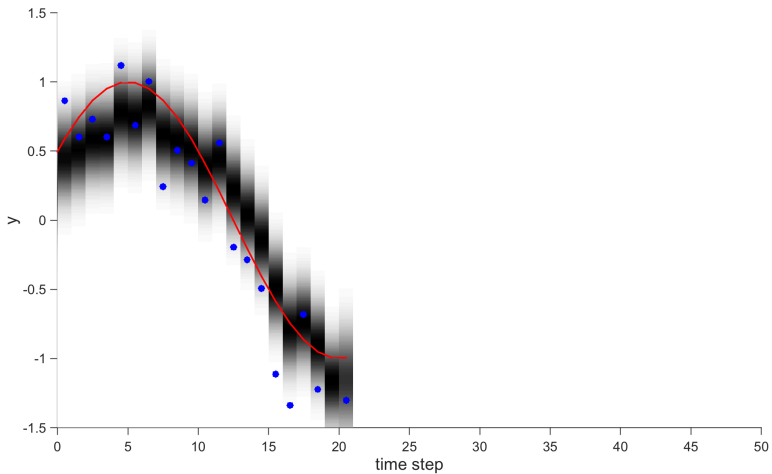
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



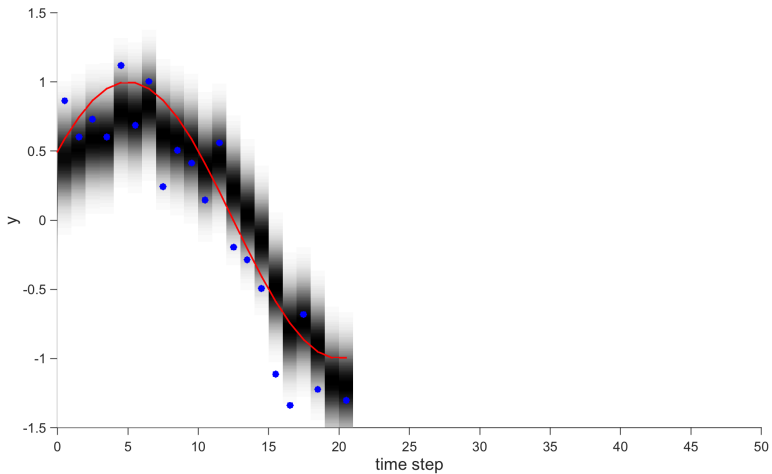
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



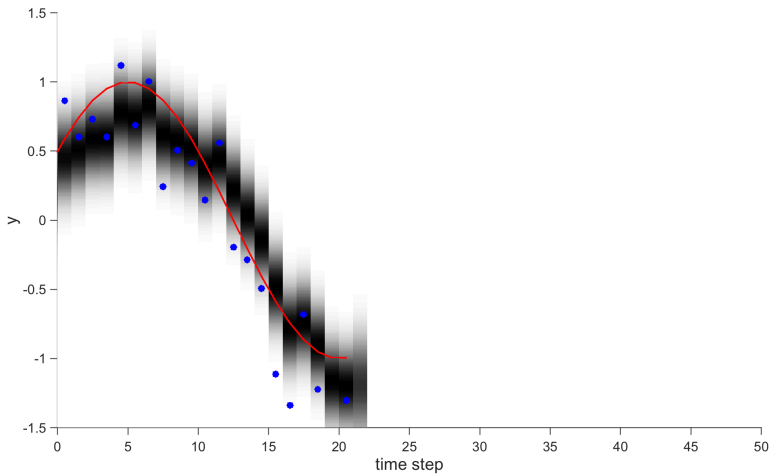
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



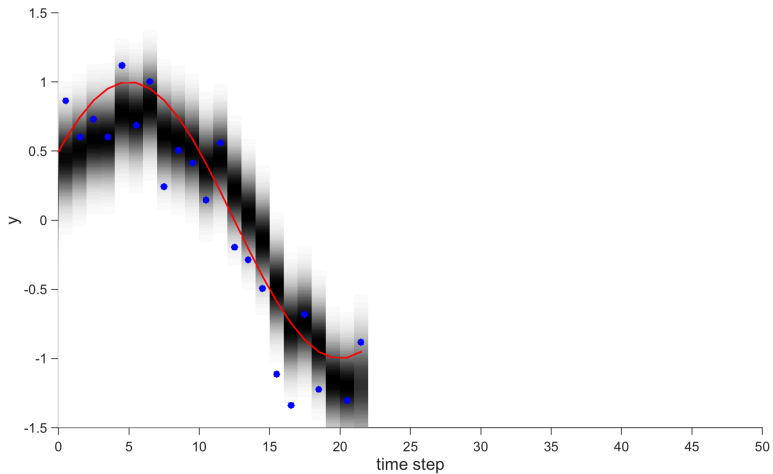
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



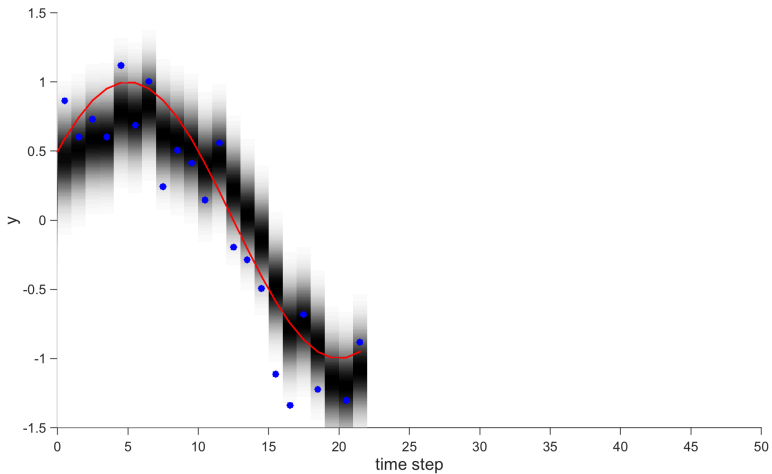
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



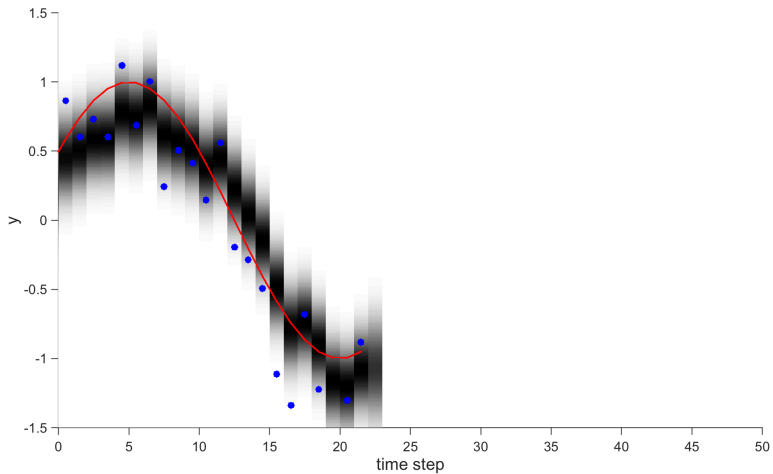
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



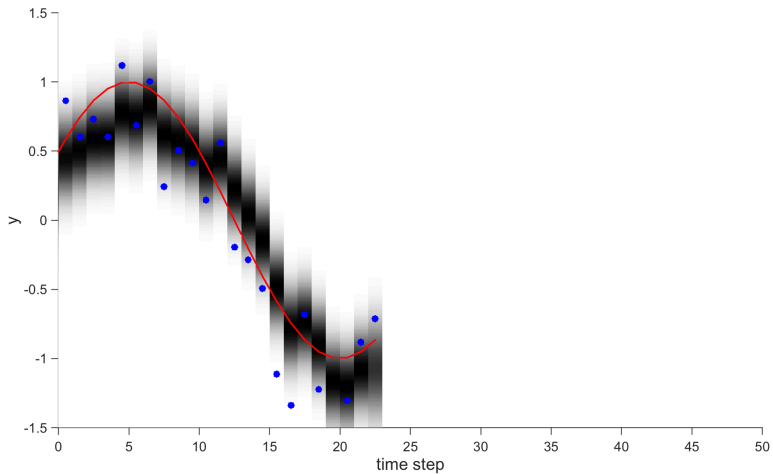
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



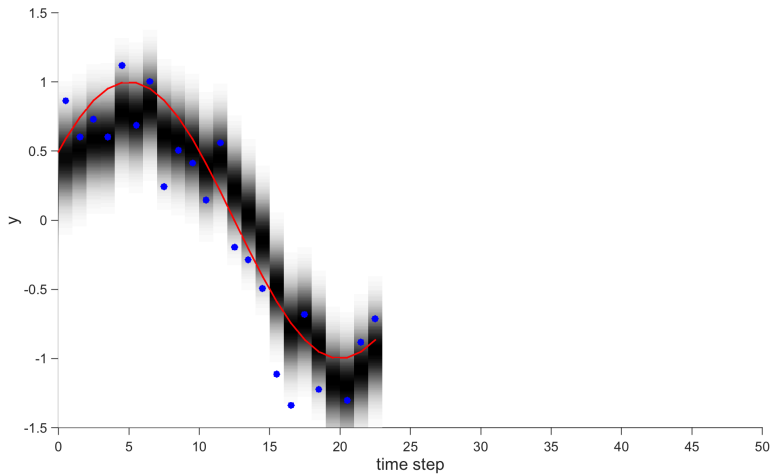
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



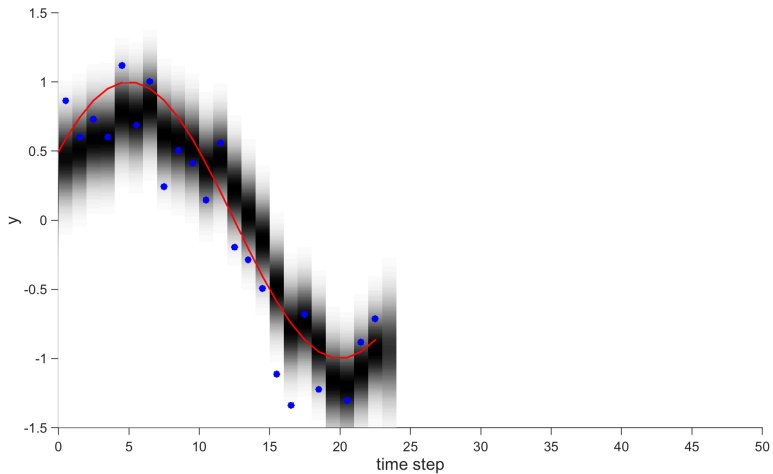
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



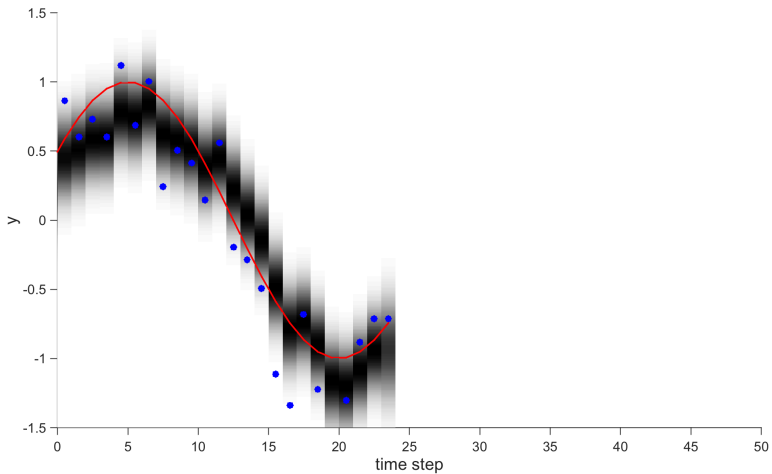
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



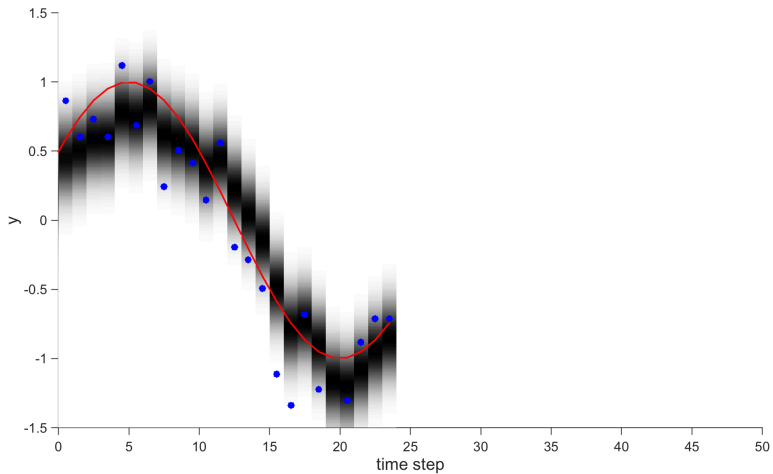
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



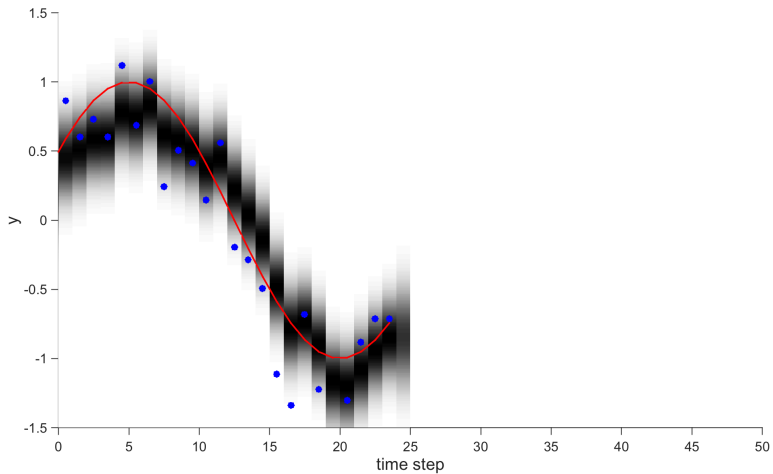
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



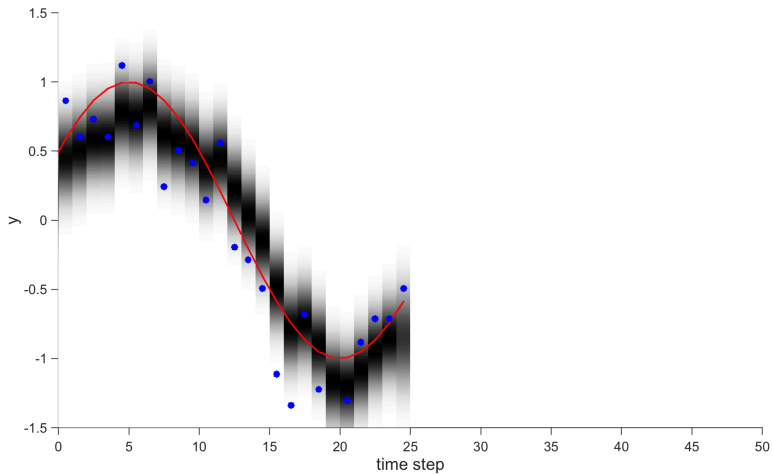
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



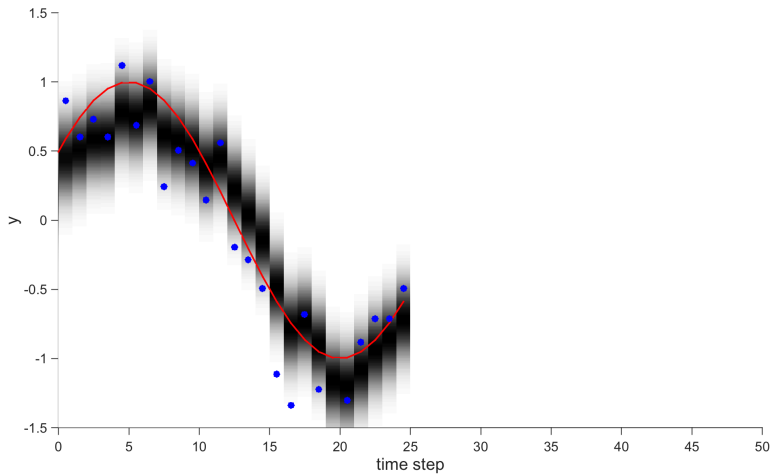
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



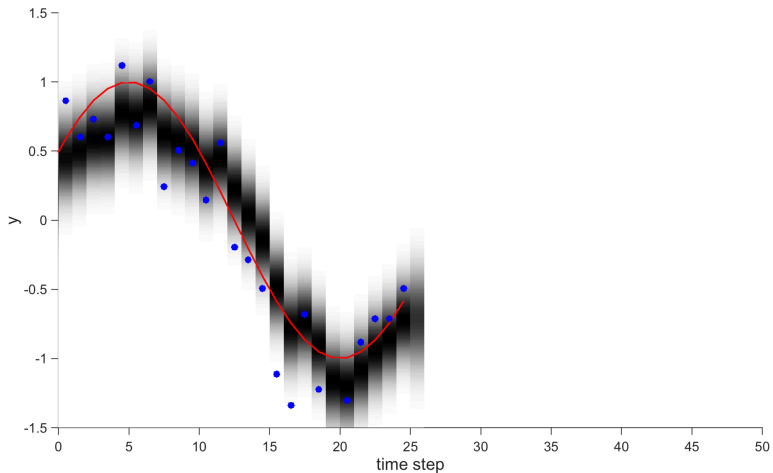
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



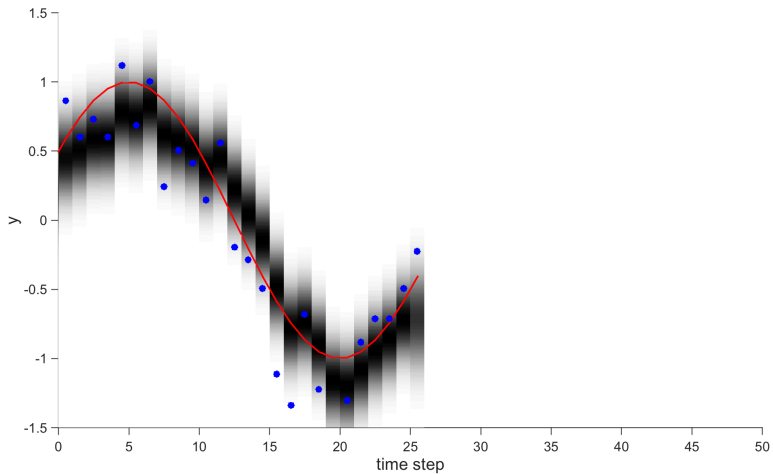
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



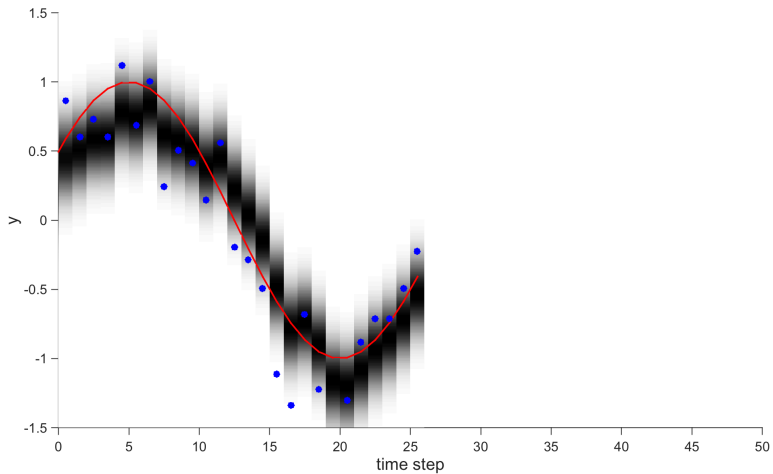
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



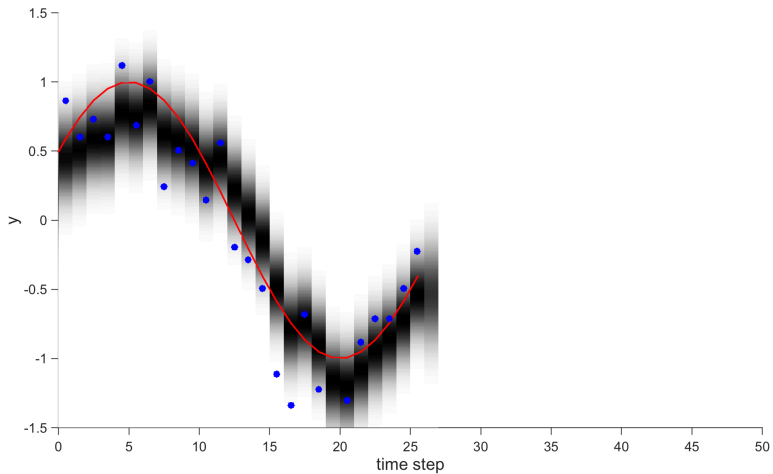
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



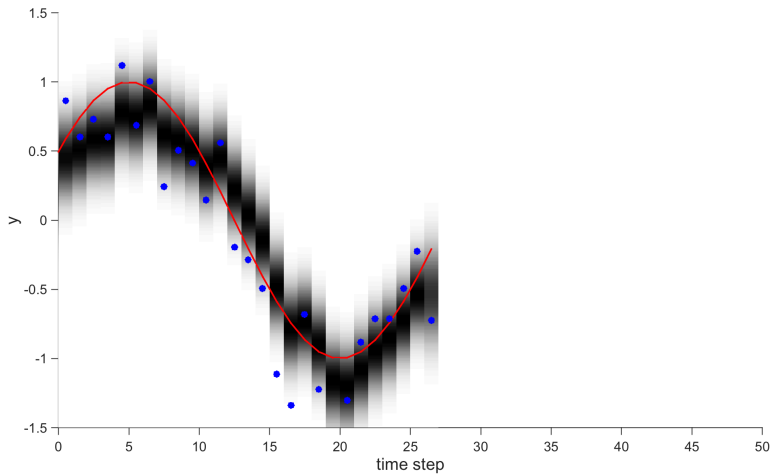
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



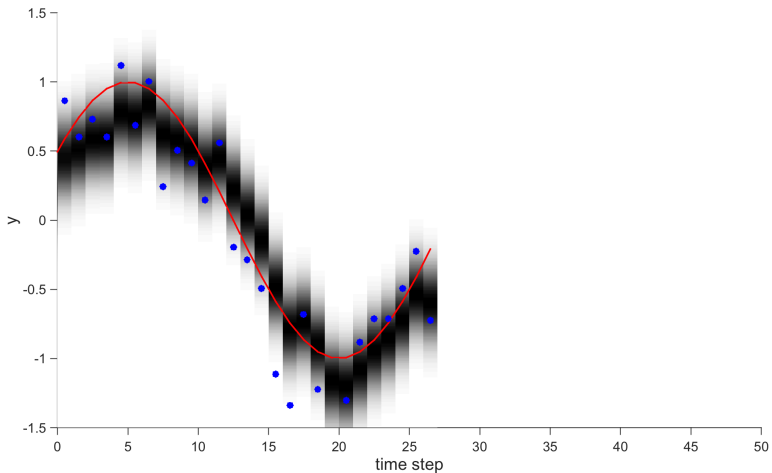
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



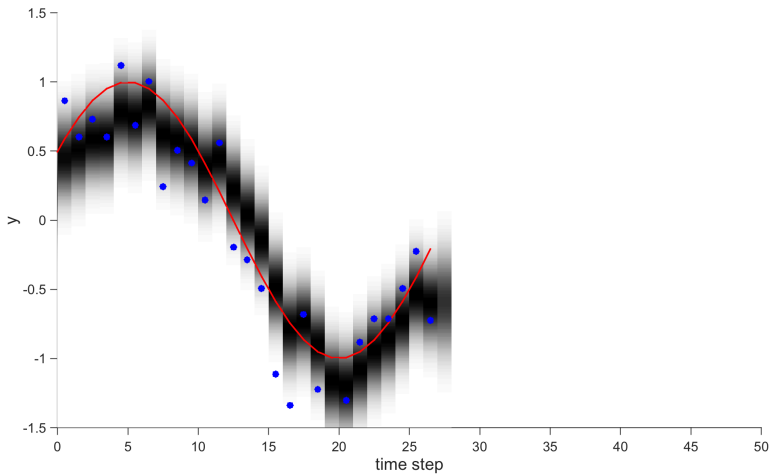
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



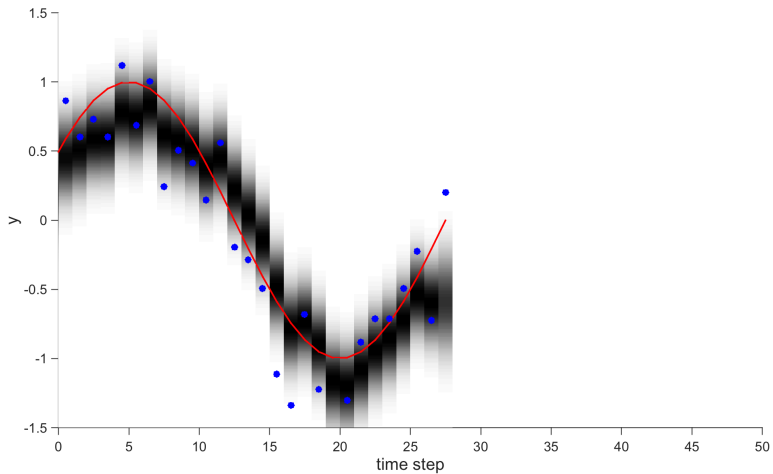
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



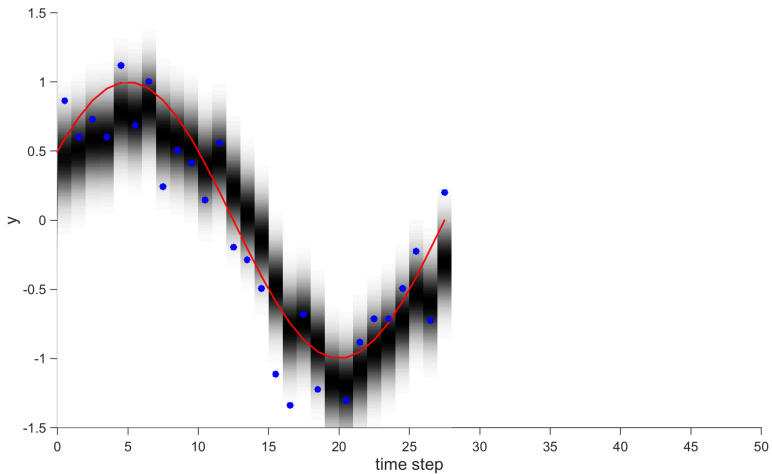
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



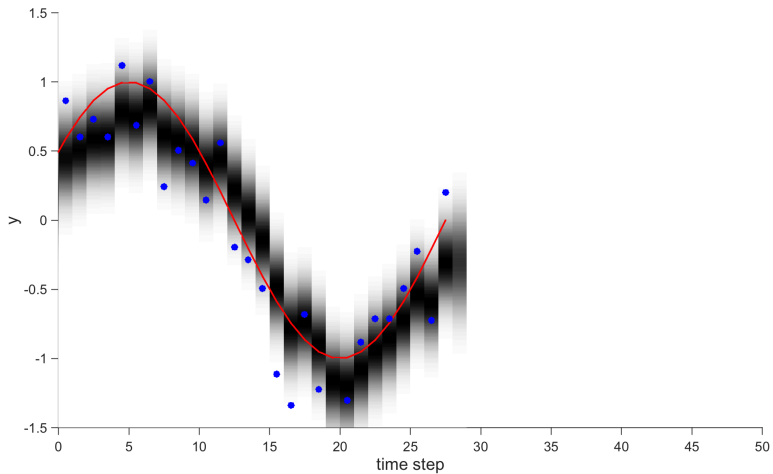
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



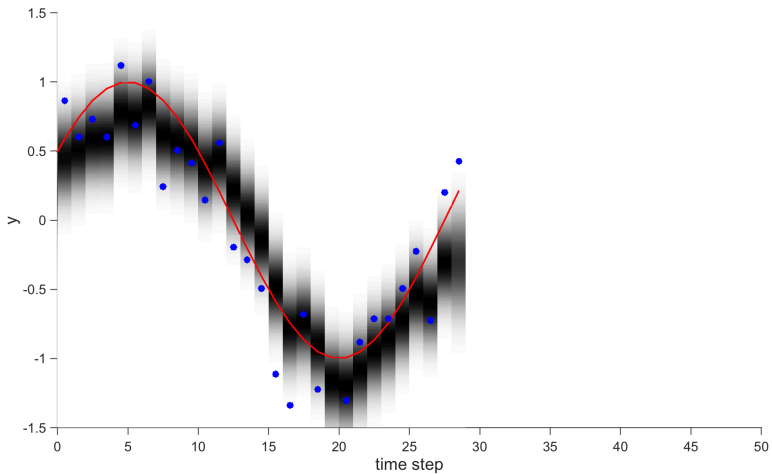
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



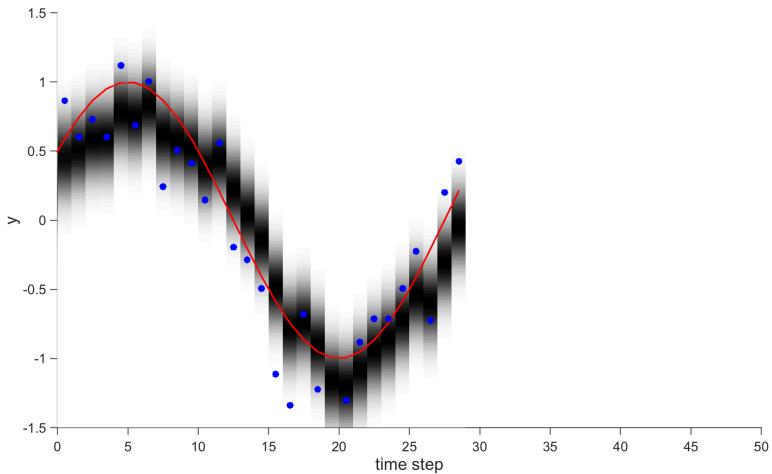
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



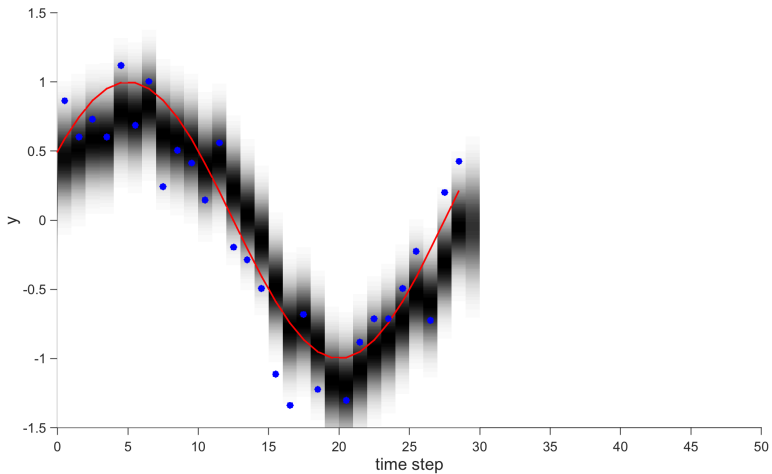
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



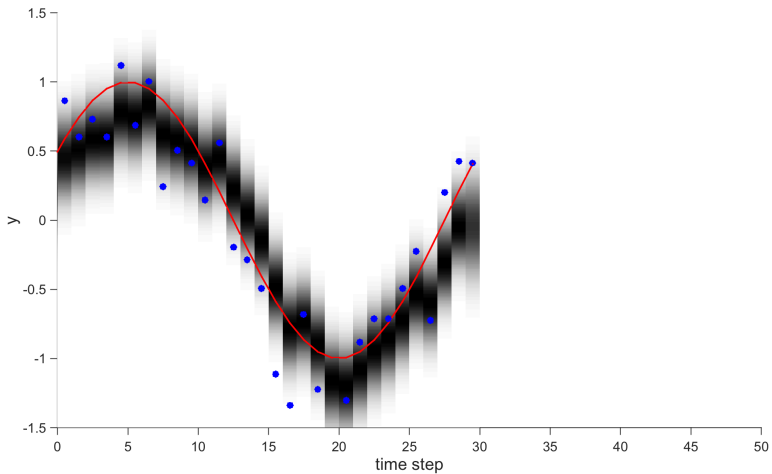
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



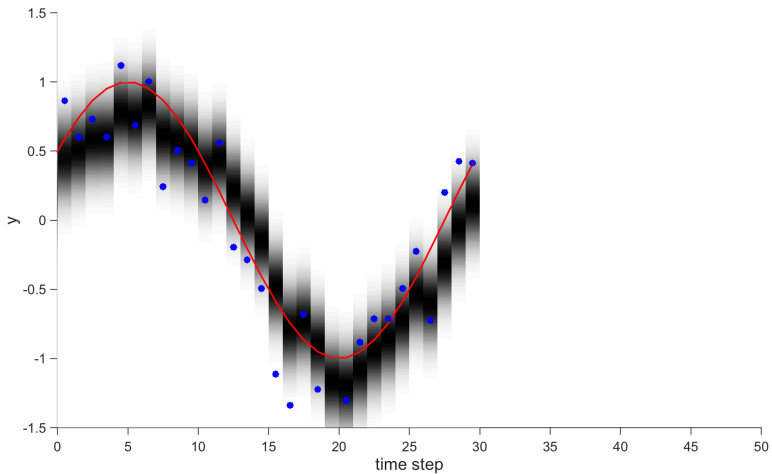
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



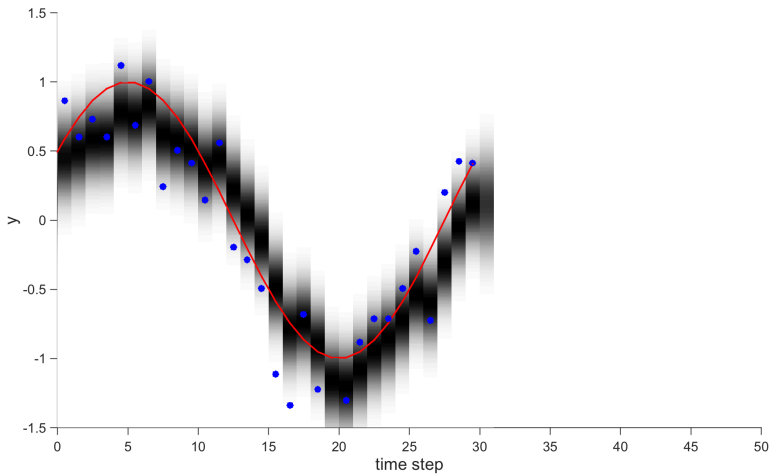
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



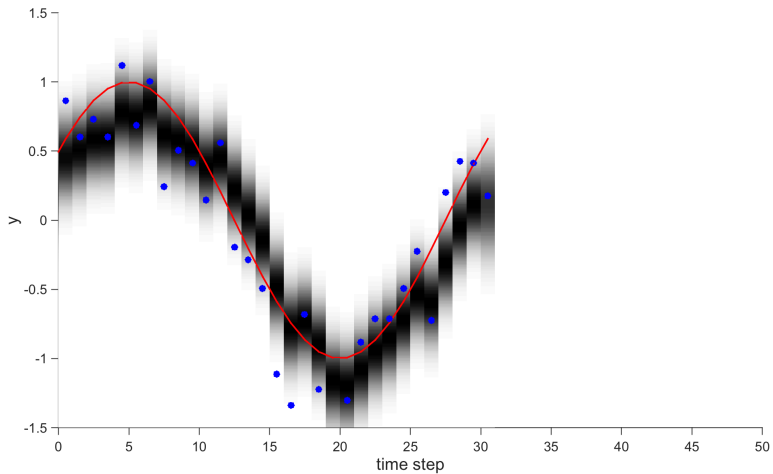
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



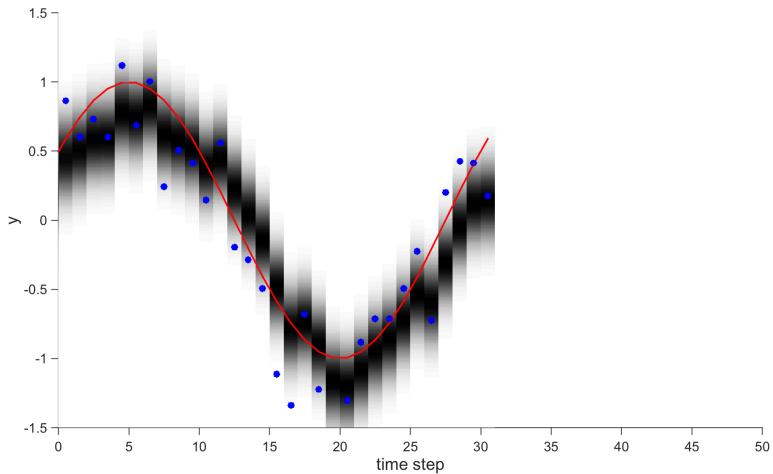
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



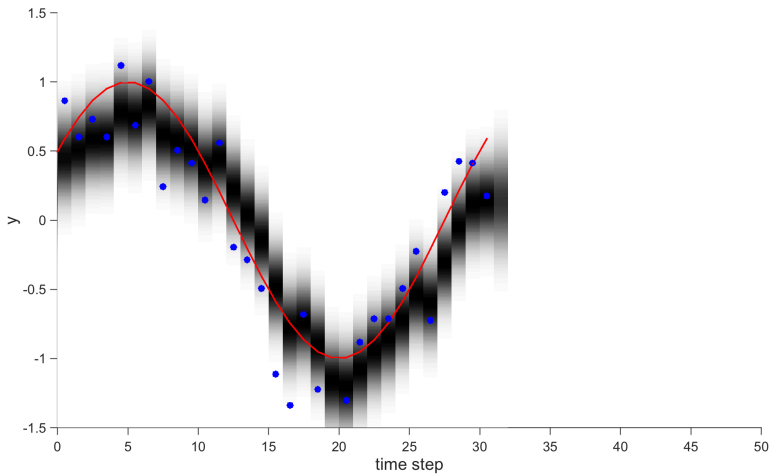
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



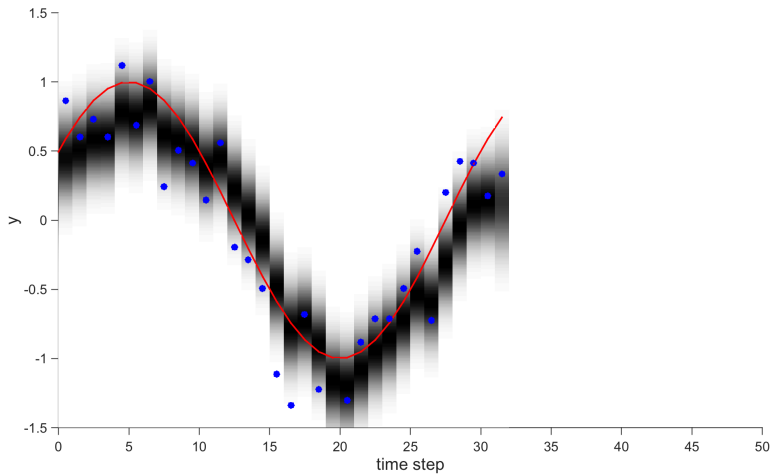
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



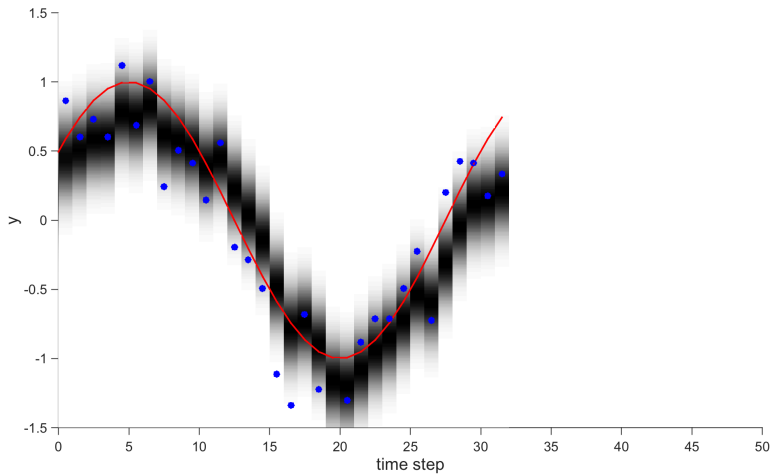
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



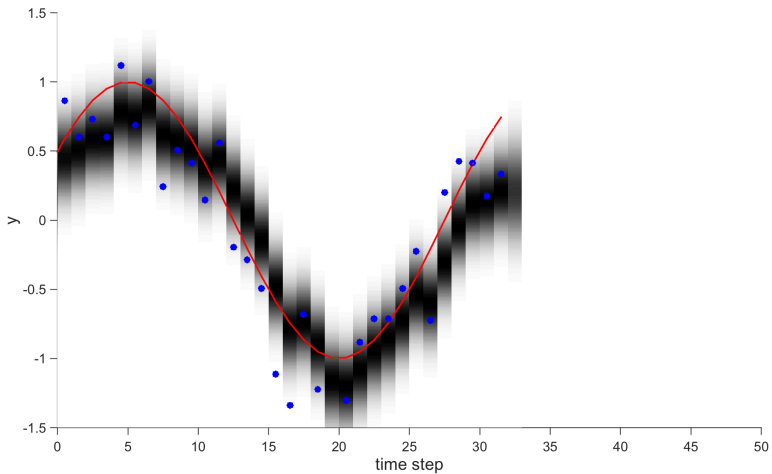
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



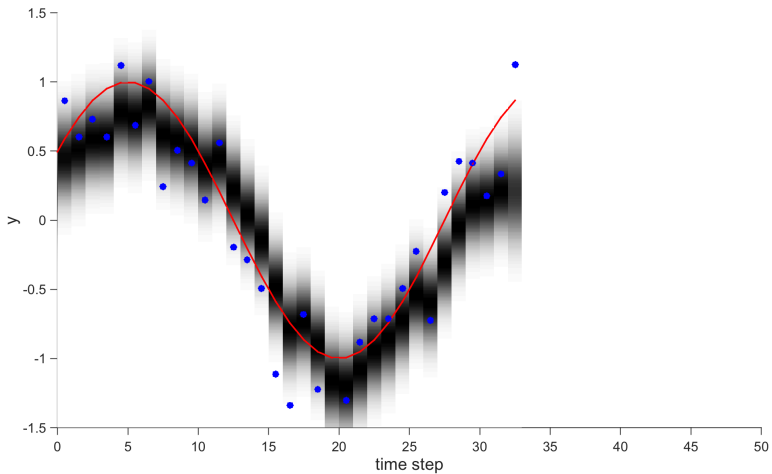
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



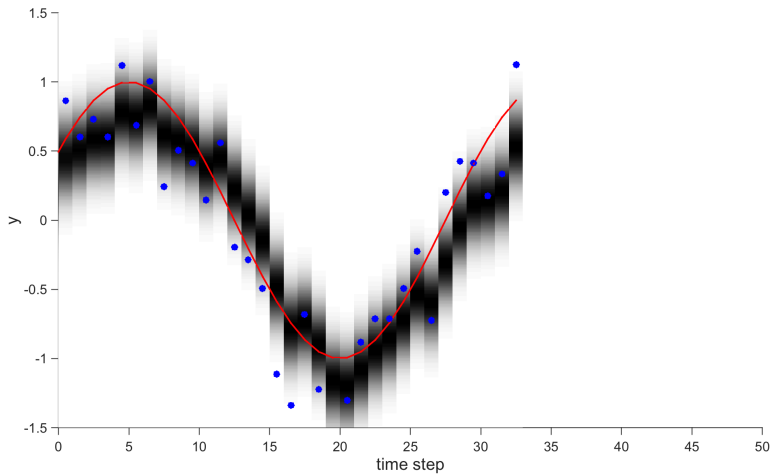
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



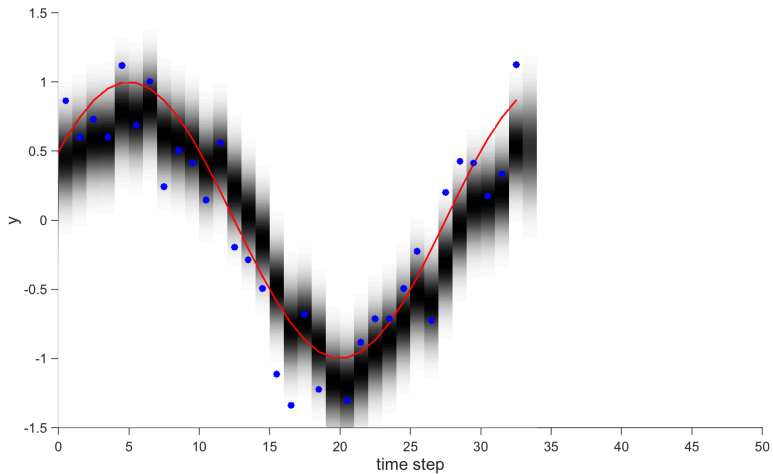
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



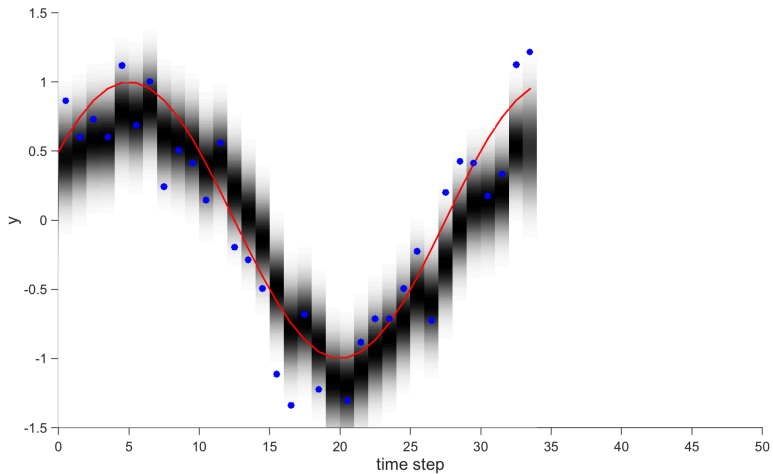
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



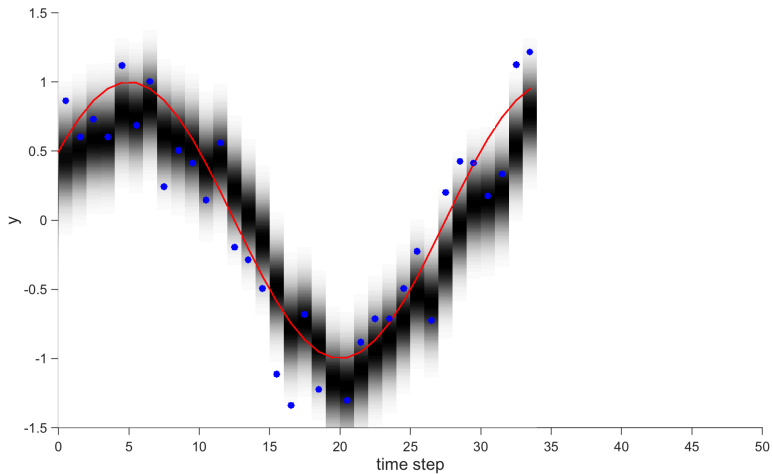
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



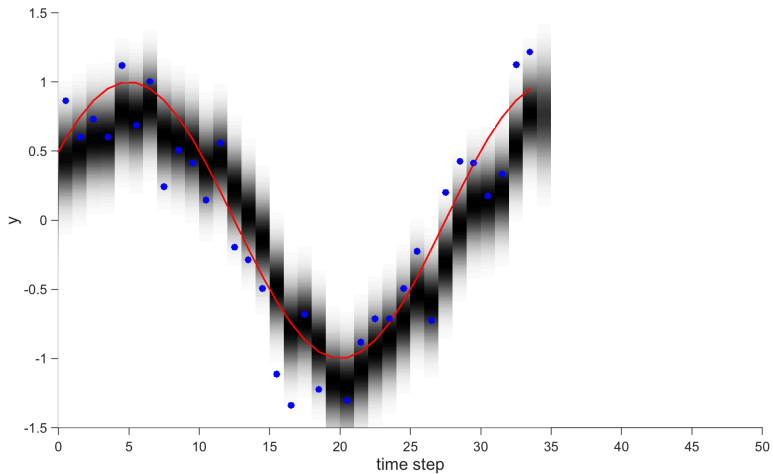
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



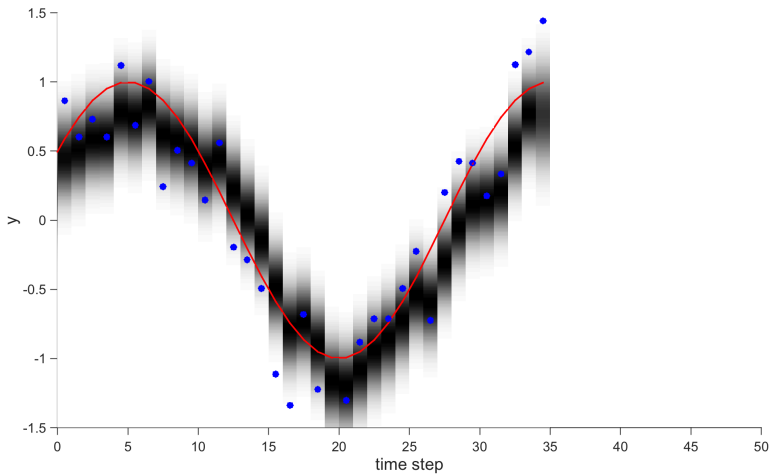
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



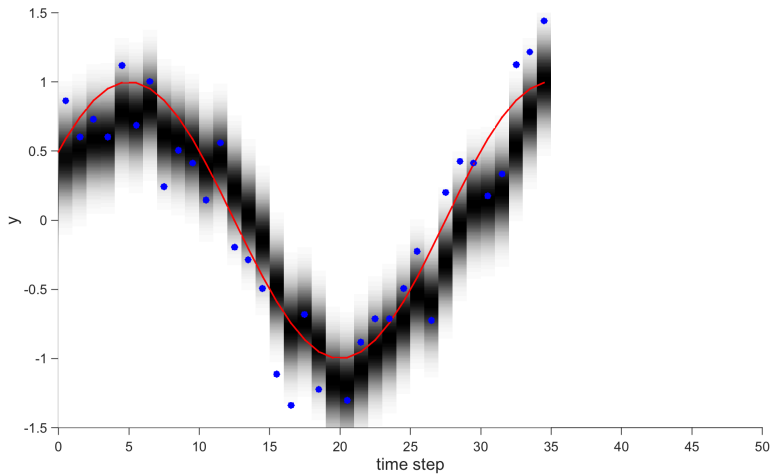
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



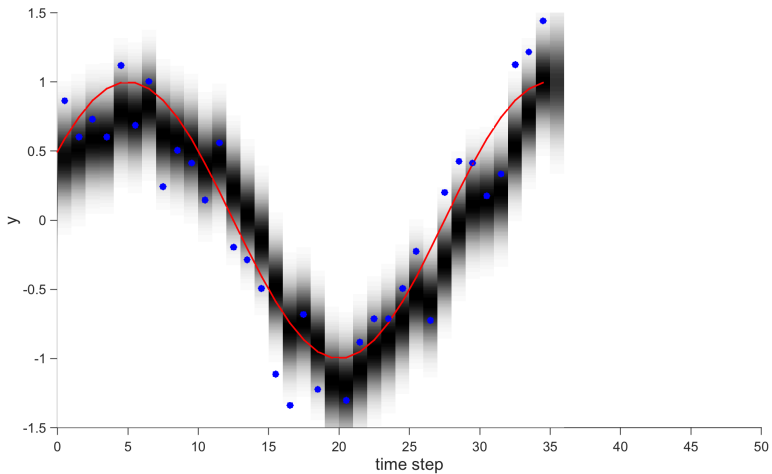
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



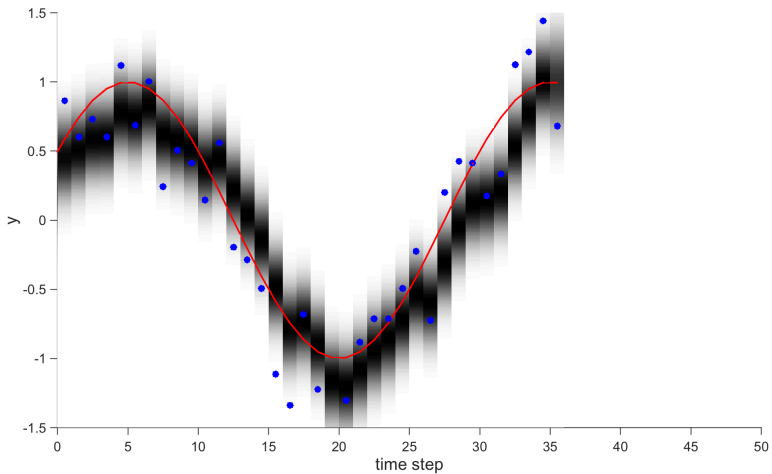
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



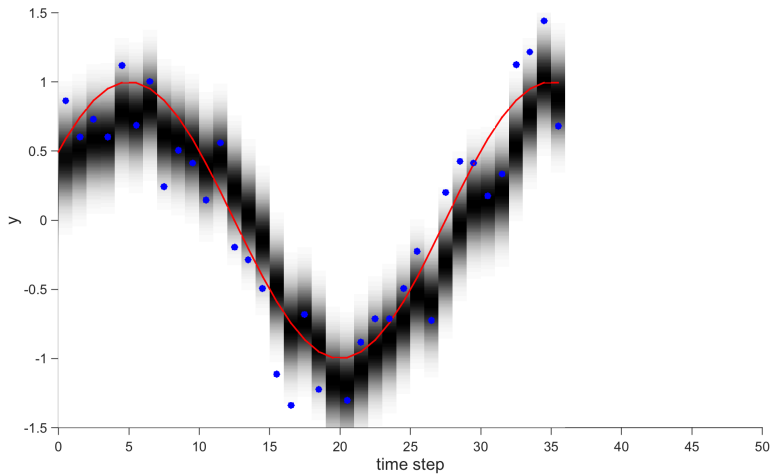
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



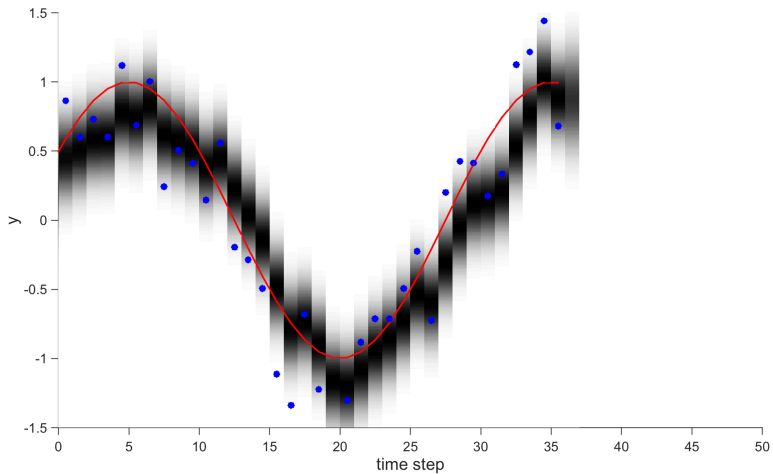
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



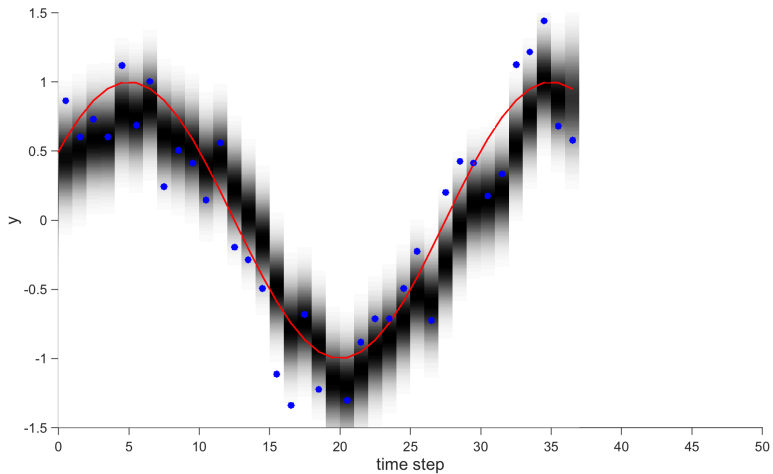
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



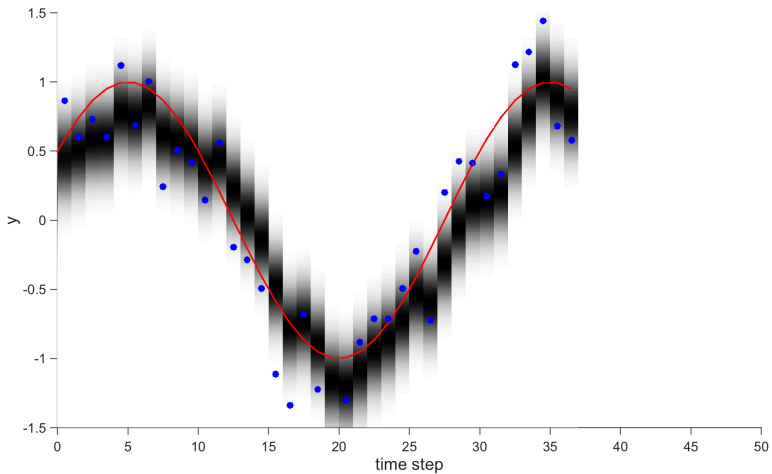
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



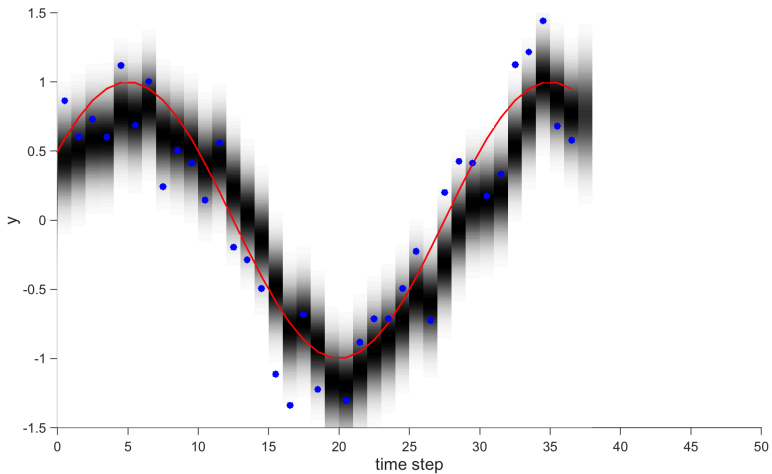
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



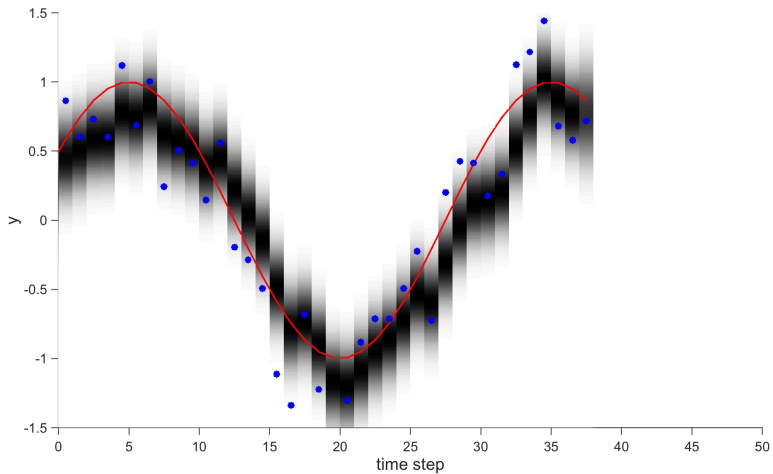
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



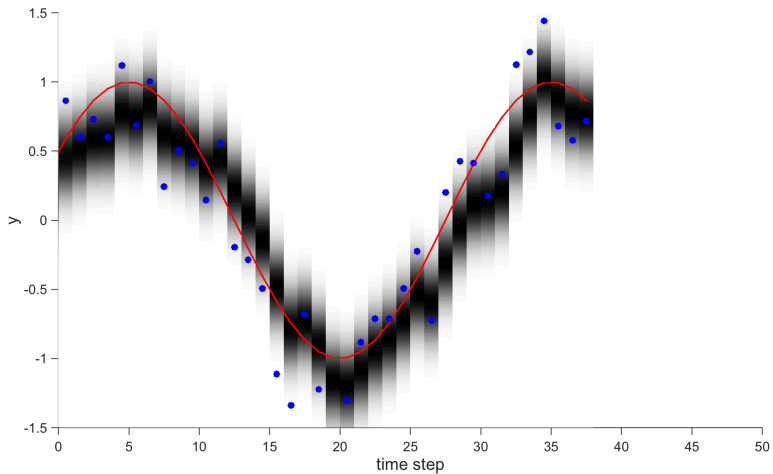
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



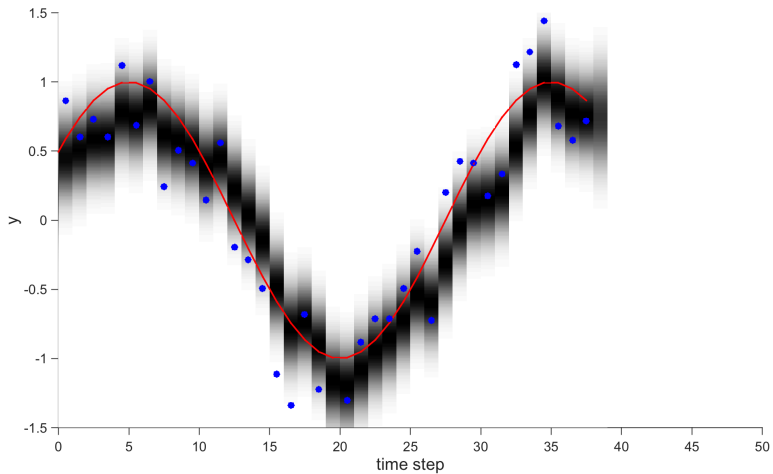
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



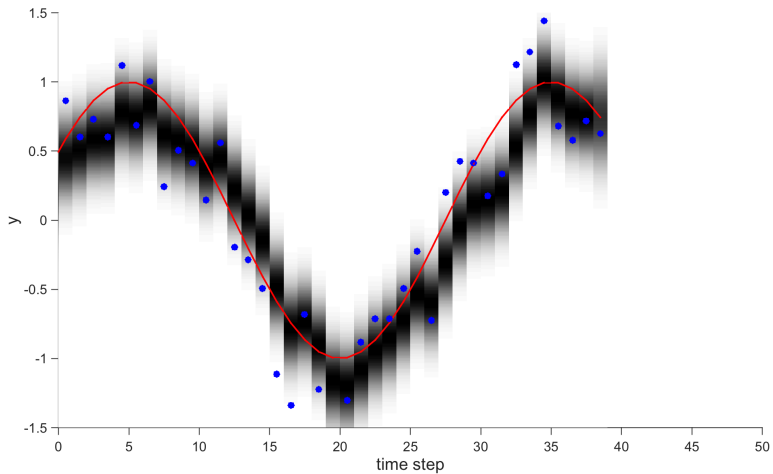
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



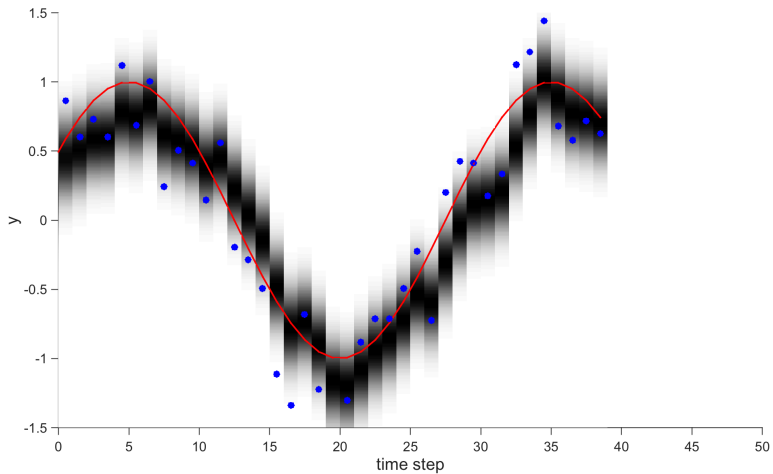
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



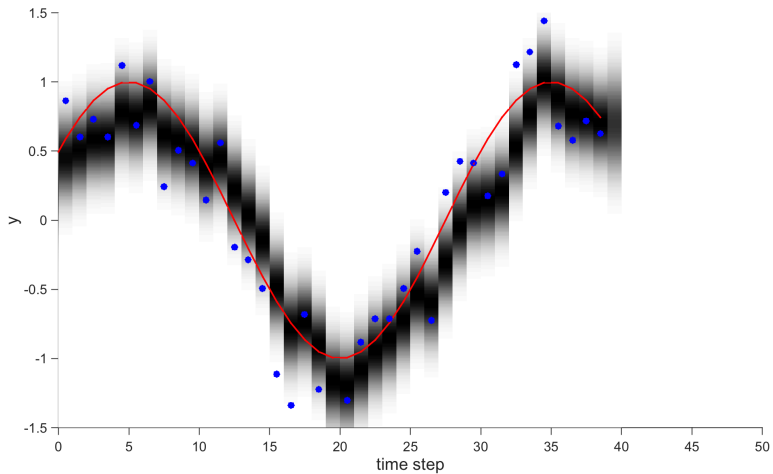
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



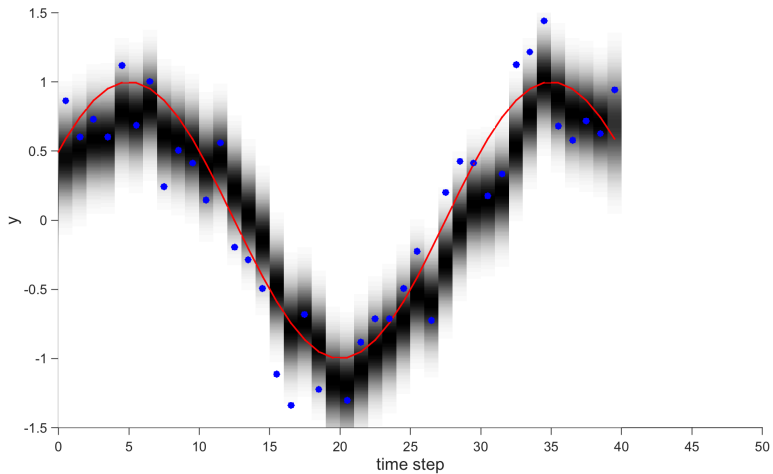
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



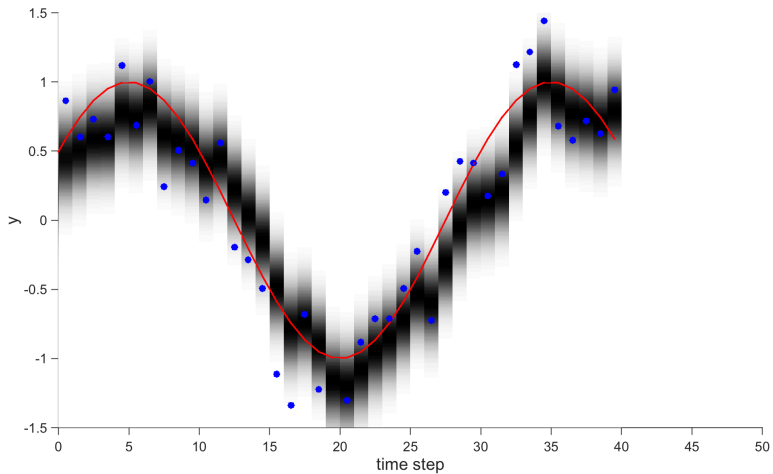
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



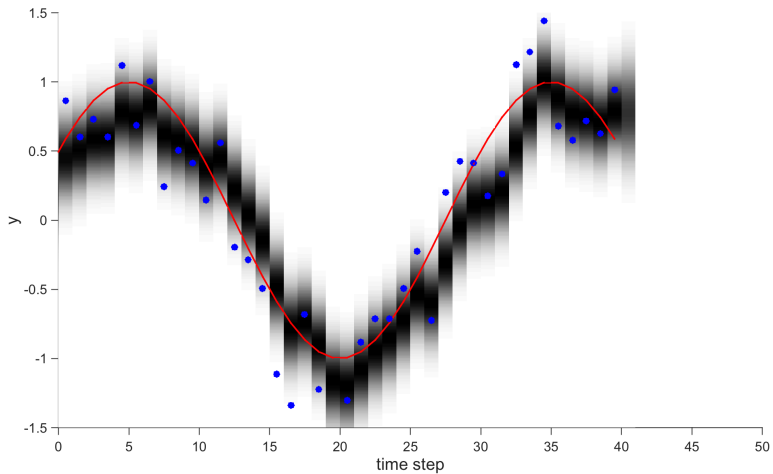
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



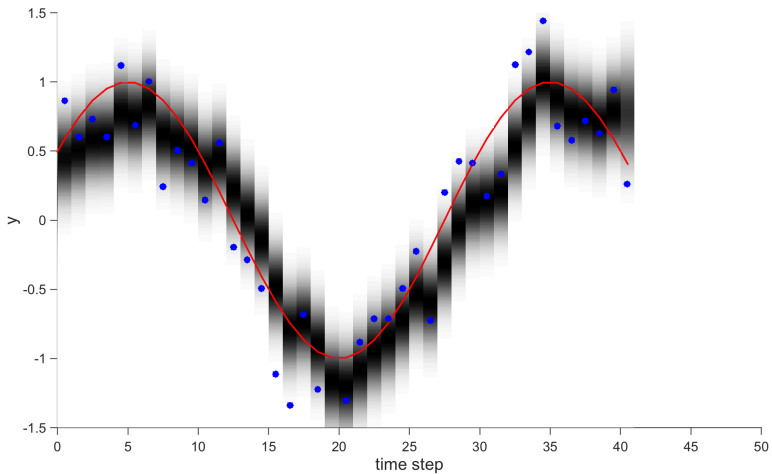
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



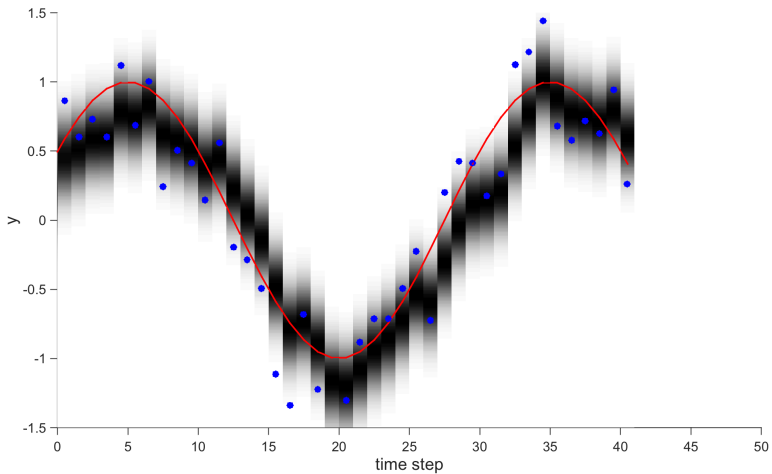
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



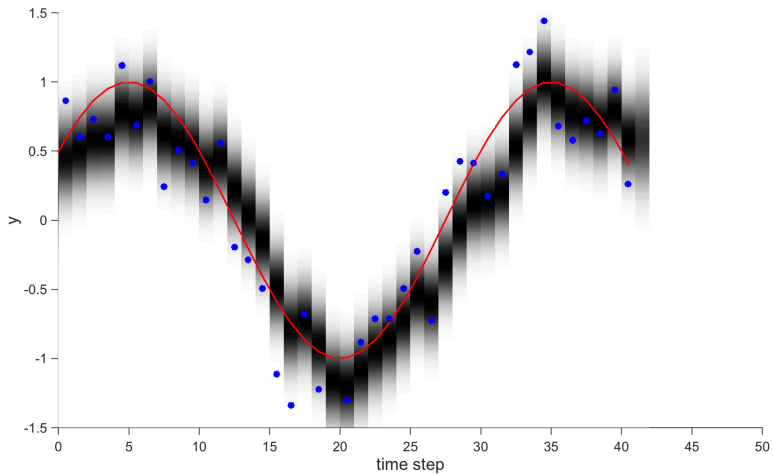
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



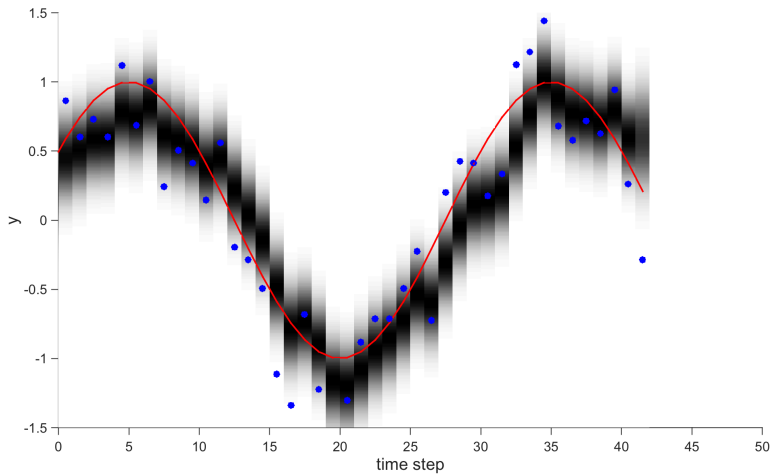
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



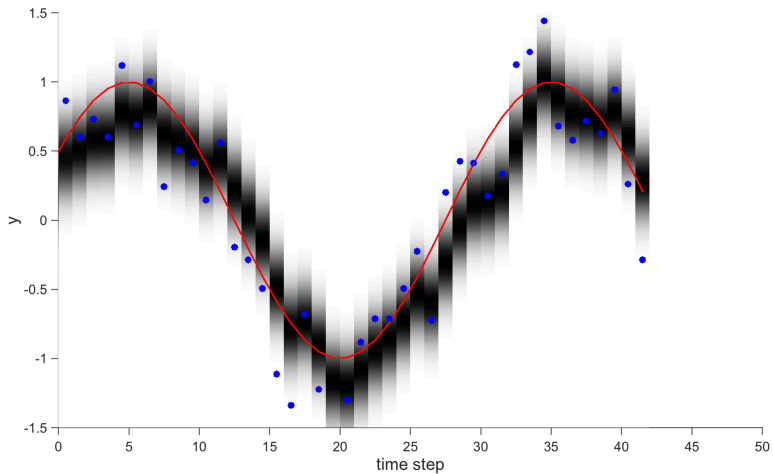
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



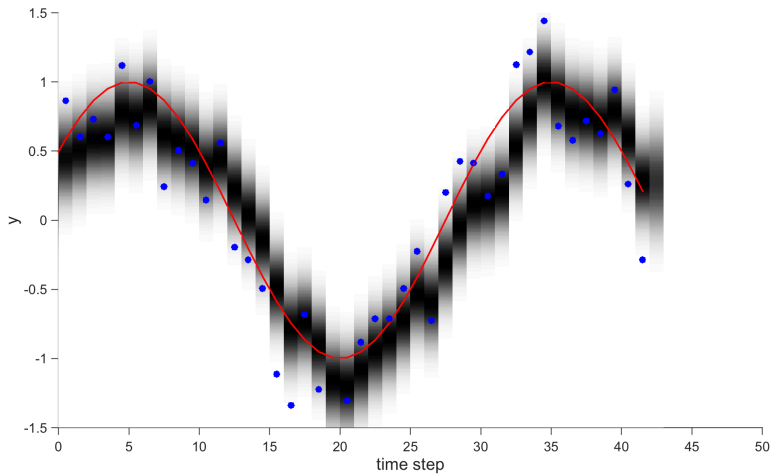
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



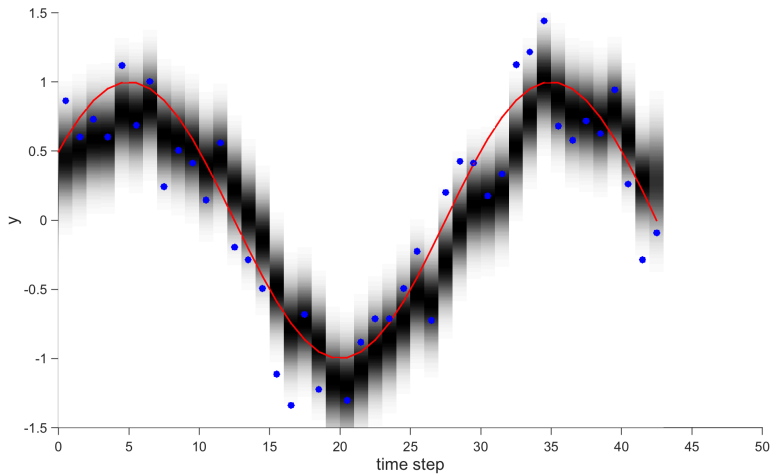
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



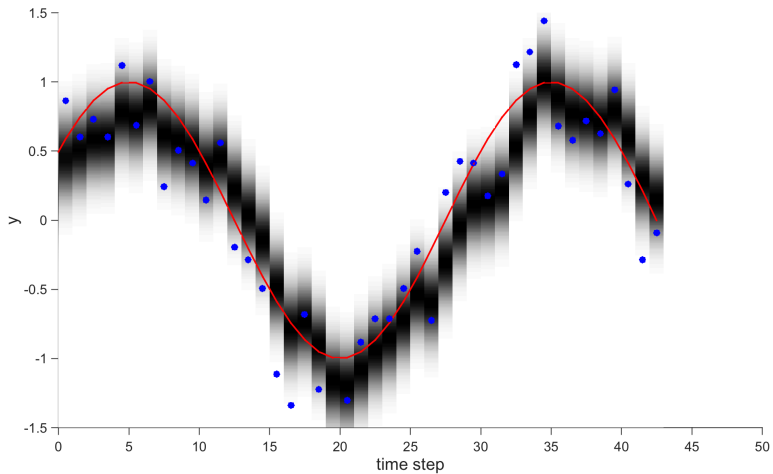
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



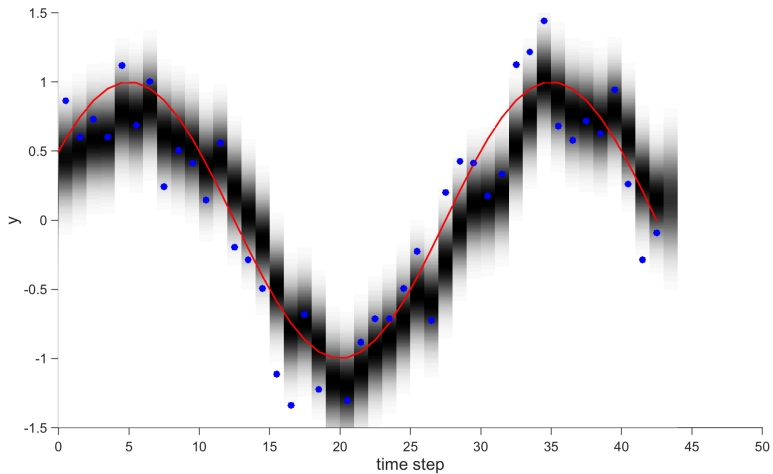
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



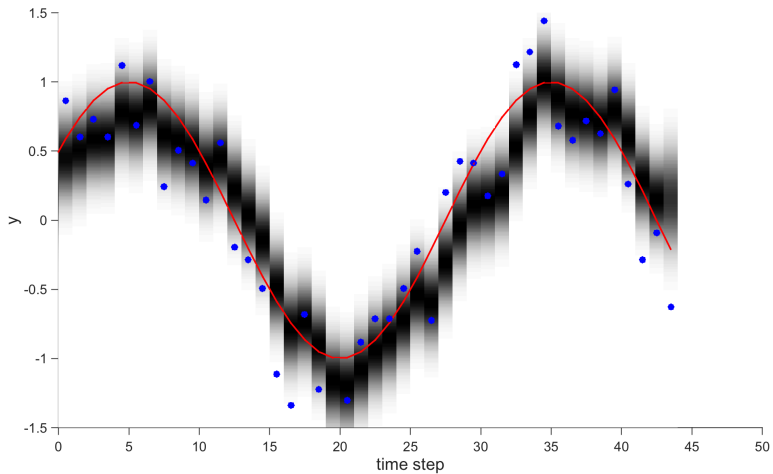
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



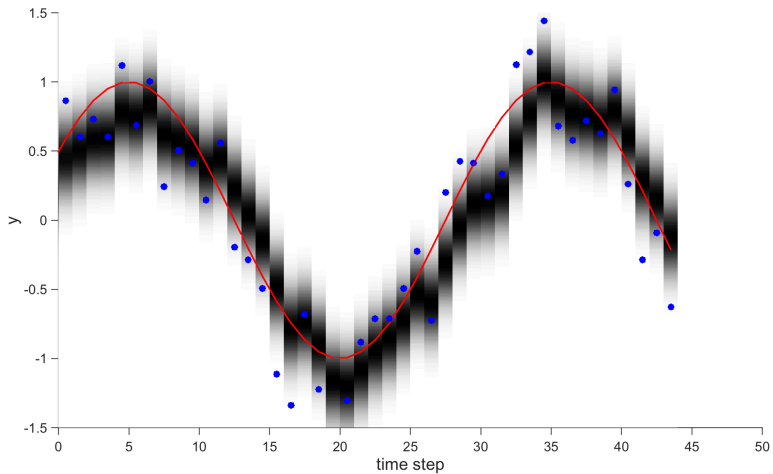
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



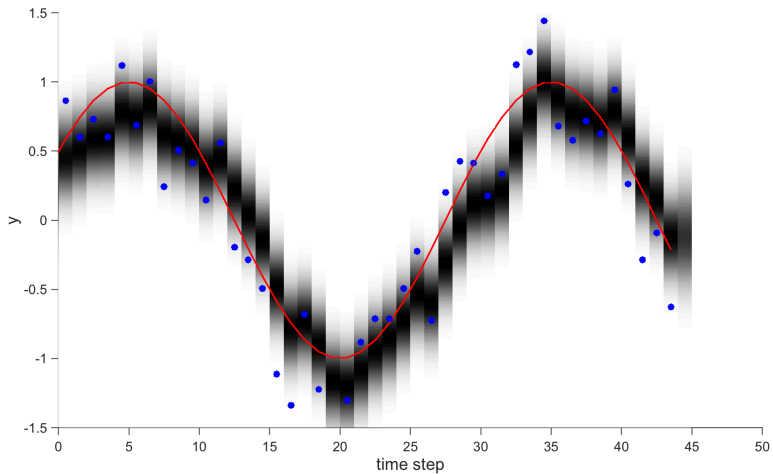
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



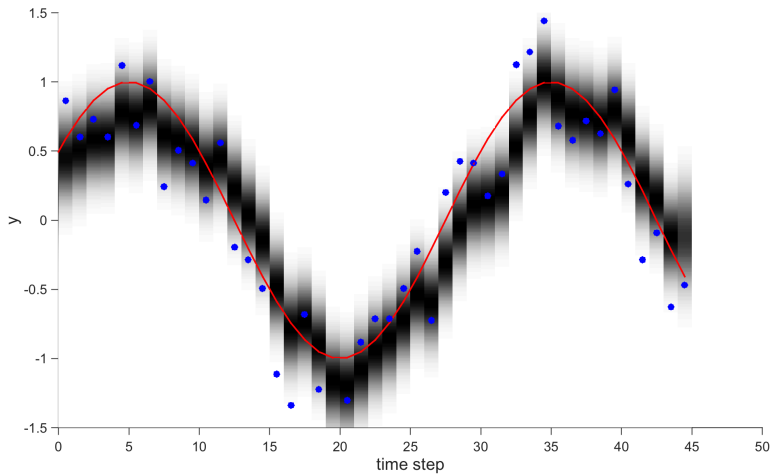
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



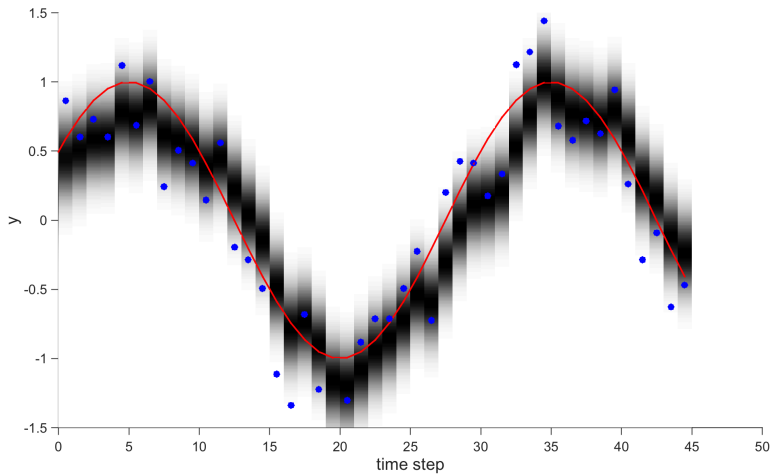
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



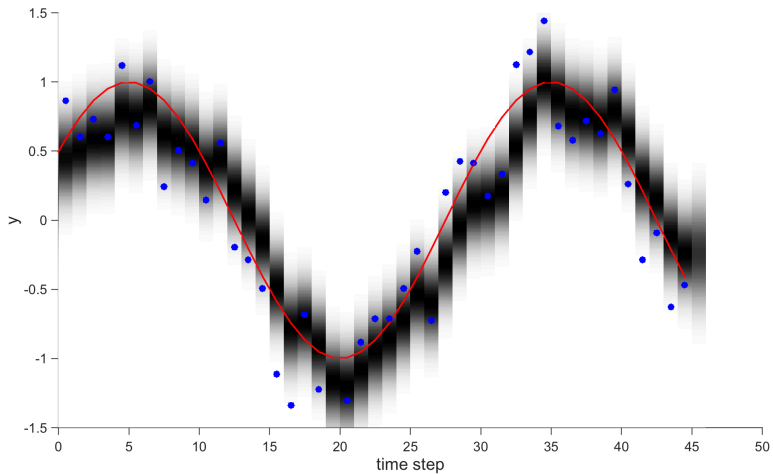
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



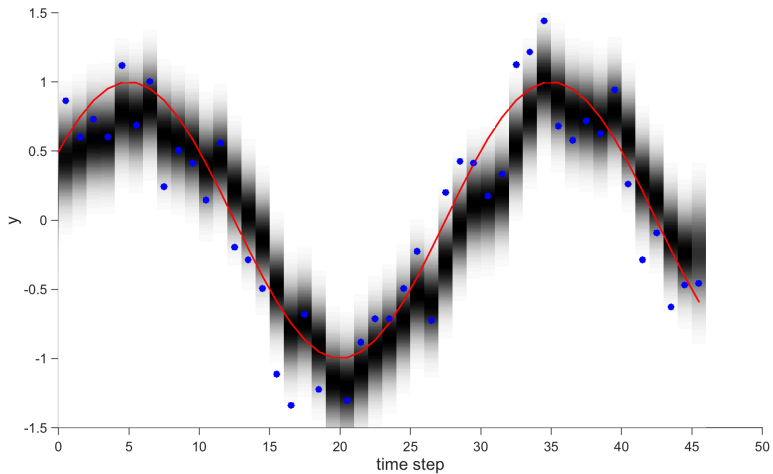
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



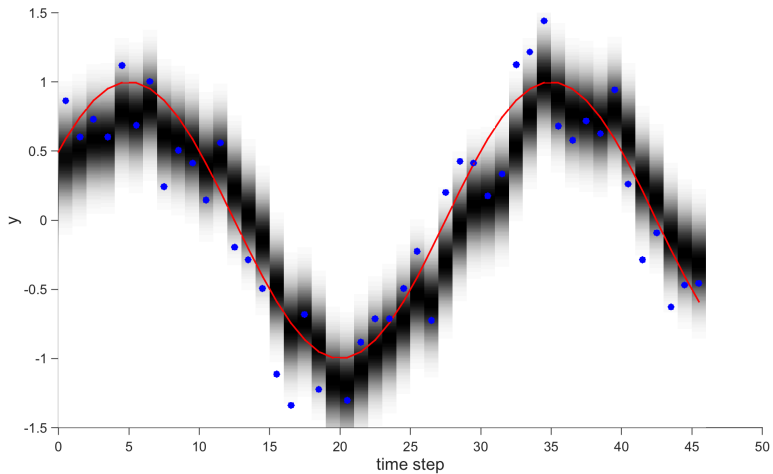
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



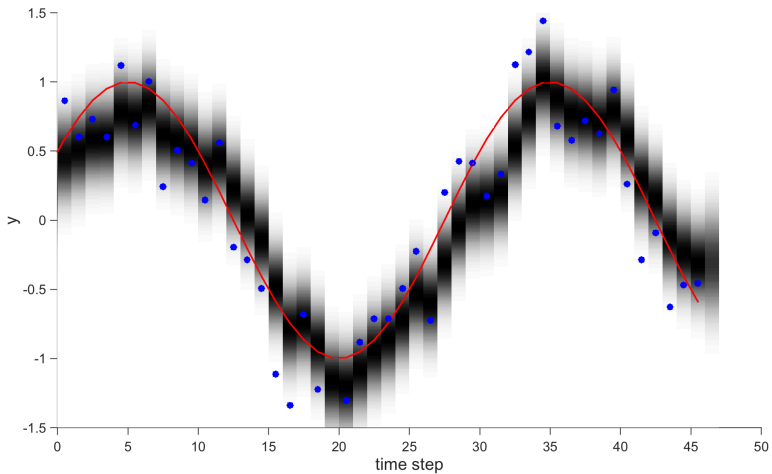
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



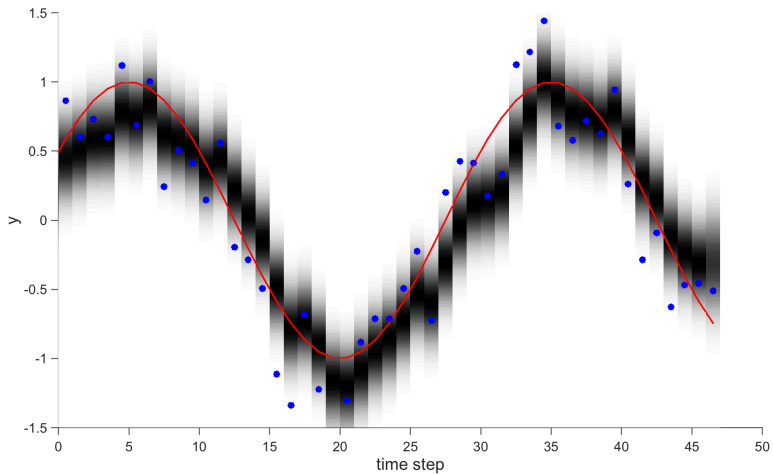
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



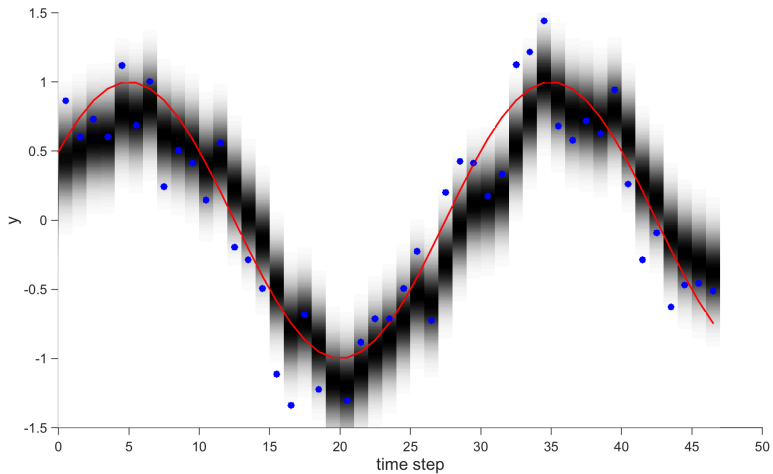
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



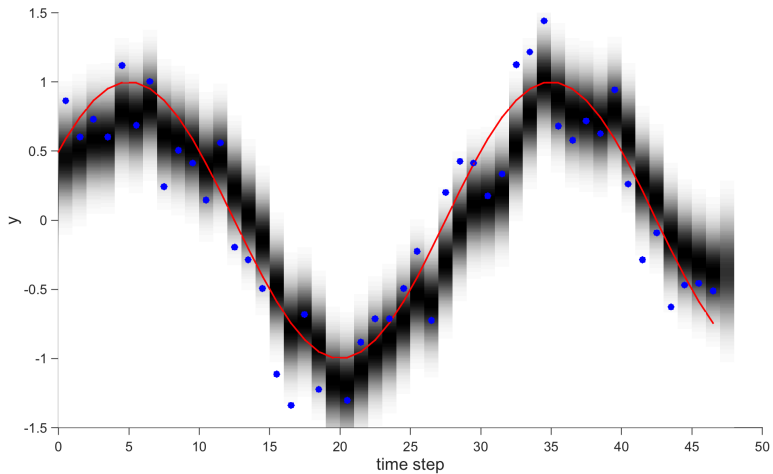
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



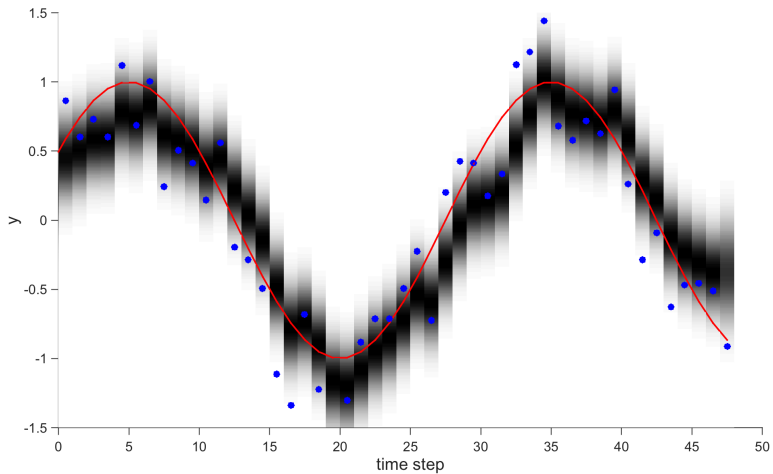
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



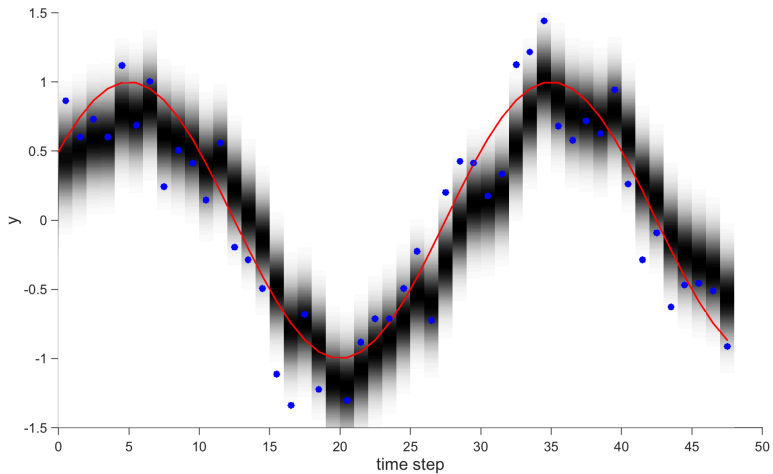
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



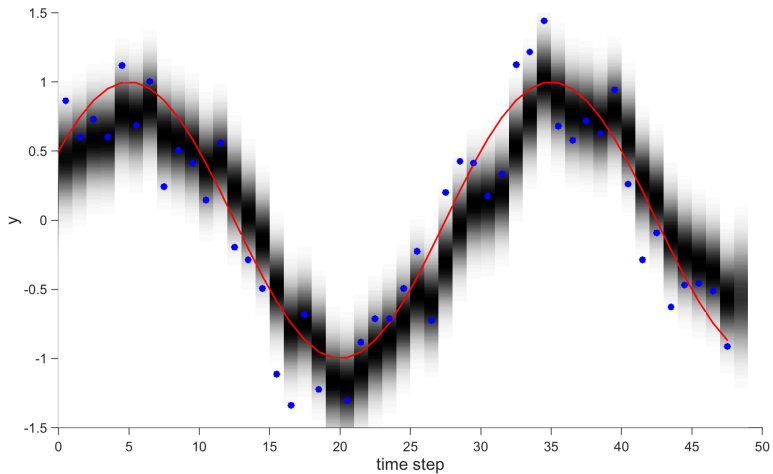
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



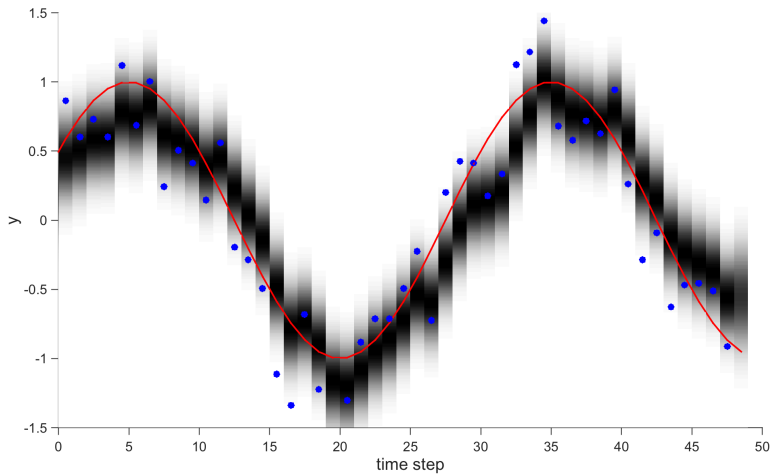
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



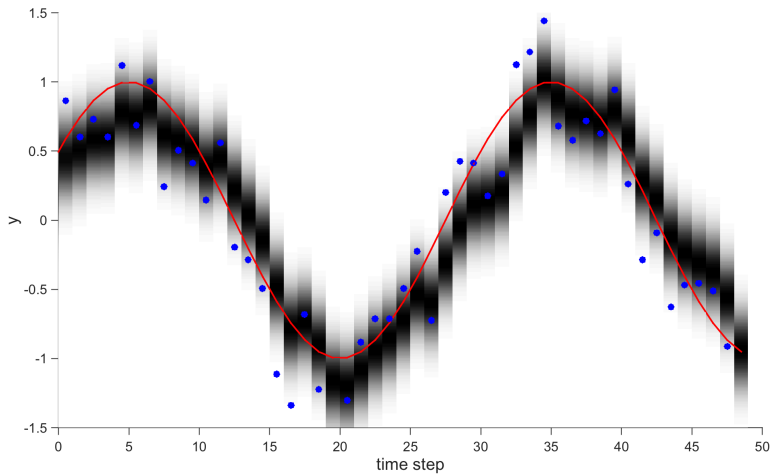
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



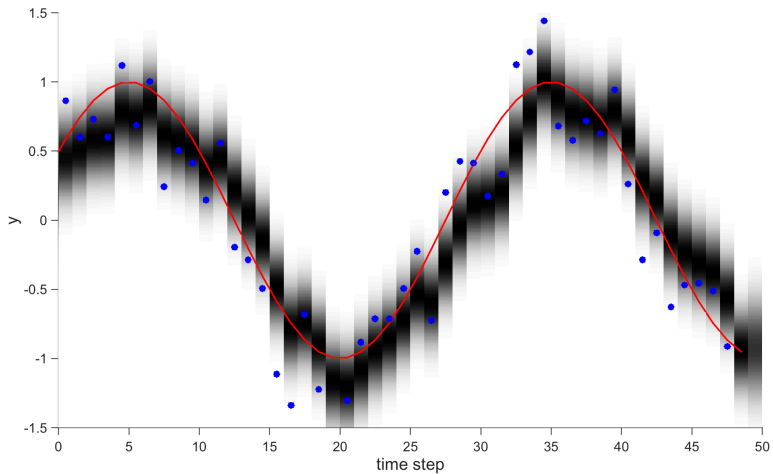
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



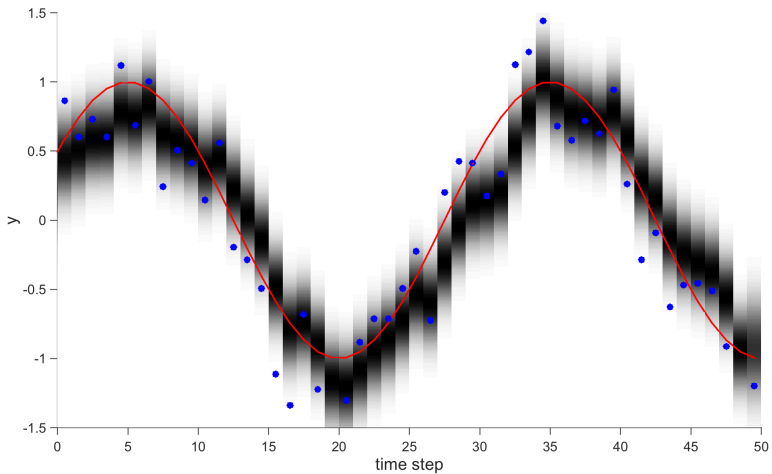
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



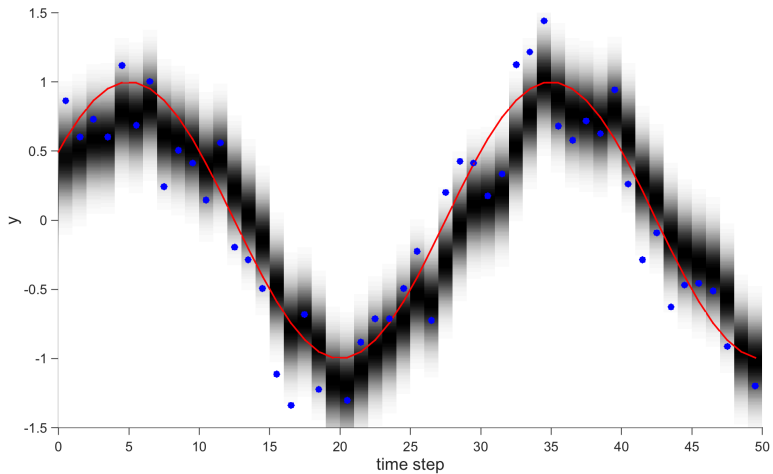
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid





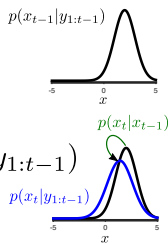
Inference: Forward Algorithm

$$p(x_{t-1} = k | y_{1:t-1})$$

diffuse via
dynamics



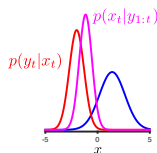
$$p(x_t = k | y_{1:t-1}) = \sum_{l=1}^K p(x_t = k | x_{t-1} = l) p(x_{t-1} = l | y_{1:t-1})$$



combine
with
likelihood



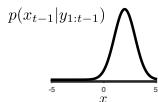
$$p(x_t = k | y_{1:t}) \propto \underbrace{p(x_t = k | y_{1:t-1})}_{\text{prior}} \underbrace{p(y_t | x_t = k)}_{\text{likelihood}}$$



Inference: Forward Algorithm

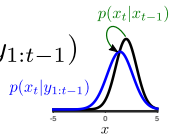
$$p(x_{t-1} = k | y_{1:t-1}) = \rho_{t-1}^{t-1}(k)$$

← most recent data used in prediction
← variable being predicted



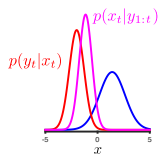
diffuse via dynamics

$$p(x_t = k | y_{1:t-1}) = \sum_{l=1}^K p(x_t = k | x_{t-1} = l) p(x_{t-1} = l | y_{1:t-1})$$



combine with likelihood

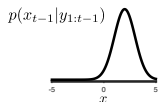
$$p(x_t = k | y_{1:t}) \propto \underbrace{p(x_t = k | y_{1:t-1})}_{\text{prior}} \underbrace{p(y_t | x_t = k)}_{\text{likelihood}}$$



Inference: Forward Algorithm

$$p(x_{t-1} = k | y_{1:t-1}) = \rho_{t-1}^{t-1}(k)$$

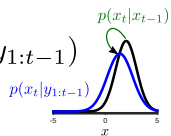
← most recent data used in prediction
← variable being predicted



diffuse via dynamics

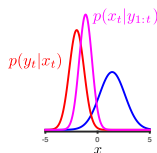
$$p(x_t = k | y_{1:t-1}) = \sum_{l=1}^K p(x_t = k | x_{t-1} = l) p(x_{t-1} = l | y_{1:t-1})$$

$$\rho_t^{t-1}(k) = \sum_{l=1}^K T(k, l) \rho_{t-1}^{t-1}(l)$$



combine with likelihood

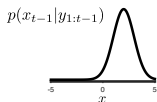
$$p(x_t = k | y_{1:t}) \propto \underbrace{p(x_t = k | y_{1:t-1})}_{\text{prior}} \underbrace{p(y_t | x_t = k)}_{\text{likelihood}}$$



Inference: Forward Algorithm

$$p(x_{t-1} = k | y_{1:t-1}) = \rho_{t-1}^{t-1}(k)$$

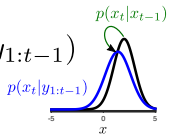
← most recent data used in prediction
← variable being predicted



diffuse via dynamics

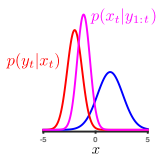
$$p(x_t = k | y_{1:t-1}) = \sum_{l=1}^K p(x_t = k | x_{t-1} = l) p(x_{t-1} = l | y_{1:t-1})$$

$$\rho_t^{t-1}(k) = \sum_{l=1}^K T(k, l) \rho_{t-1}^{t-1}(l)$$



combine with likelihood

$$p(x_t = k | y_{1:t}) \propto \underbrace{p(x_t = k | y_{1:t-1})}_{\text{prior}} \underbrace{p(y_t | x_t = k)}_{\text{likelihood}}$$

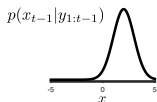


$$\rho_t^t(k) \propto \rho_{t-1}^{t-1}(k) p(y_t | x_t = k)$$

Inference: Forward Algorithm

$$p(x_{t-1} = k | y_{1:t-1}) = \rho_{t-1}^{t-1}(k)$$

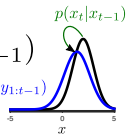
← most recent data used in prediction
← variable being predicted



diffuse via dynamics

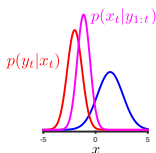
$$p(x_t = k | y_{1:t-1}) = \sum_{l=1}^K p(x_t = k | x_{t-1} = l) p(x_{t-1} = l | y_{1:t-1})$$

$$\rho_t^{t-1}(k) = \sum_{l=1}^K T(k, l) \rho_{t-1}^{t-1}(l)$$



combine with likelihood

$$p(x_t = k | y_{1:t}) \propto \underbrace{p(x_t = k | y_{1:t-1})}_{\text{prior}} \underbrace{p(y_t | x_t = k)}_{\text{likelihood}}$$



$$\rho_t^t(k) \propto \rho_{t-1}^{t-1}(k) p(y_t | x_t = k)$$

When implementing, take care with numerical underflow/overflow.

Computing the likelihood

How can we compute the likelihood efficiently?

Computing the likelihood

How can we compute the likelihood efficiently?

$$p(y_{1:T}) = \prod_{t=1}^T p(y_t | y_{1:t-1})$$

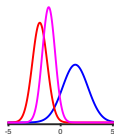
Computing the likelihood

How can we compute the likelihood efficiently?

$$p(y_{1:T}) = \prod_{t=1}^T p(y_t | y_{1:t-1})$$

already returned by Kalman Filter/Forward algorithm

$$p(x_t | y_{1:t}) = \frac{1}{p(y_t | y_{1:t-1})} p(y_t | x_t) p(x_t | y_{1:t-1})$$
$$\propto p(y_t | x_t) p(x_t | y_{1:t-1})$$



Computing the likelihood

How can we compute the likelihood efficiently?

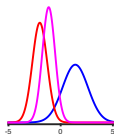
$$p(y_{1:T}) = \prod_{t=1}^T p(y_t | y_{1:t-1})$$

already returned by Kalman Filter/Forward algorithm

$$p(x_t | y_{1:t}) = \frac{1}{p(y_t | y_{1:t-1})} p(y_t | x_t) p(x_t | y_{1:t-1})$$

$$\propto p(y_t | x_t) p(x_t | y_{1:t-1})$$

$p(y_t | y_{1:t-1})$ is normaliser of filter/forward algorithm update



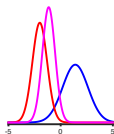
Computing the likelihood

How can we compute the likelihood efficiently?

$$p(y_{1:T}) = \prod_{t=1}^T p(y_t | y_{1:t-1})$$

already returned by Kalman Filter/Forward algorithm

$$p(x_t | y_{1:t}) = \frac{1}{p(y_t | y_{1:t-1})} p(y_t | x_t) p(x_t | y_{1:t-1})$$
$$\propto p(y_t | x_t) p(x_t | y_{1:t-1})$$



$p(y_t | y_{1:t-1})$ is normaliser of filter/forward algorithm update

How can we compute the smoothing estimate?

$$p(x_t | y_{1:T})$$

LGSSM: Kalman Smoother

HMM: Forward-Backward= Algorithm

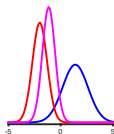
Computing the likelihood

How can we compute the likelihood efficiently?

$$p(y_{1:T}) = \prod_{t=1}^T p(y_t | y_{1:t-1})$$

already returned by Kalman Filter/Forward algorithm

$$p(x_t | y_{1:t}) = \frac{1}{p(y_t | y_{1:t-1})} p(y_t | x_t) p(x_t | y_{1:t-1})$$
$$\propto p(y_t | x_t) p(x_t | y_{1:t-1})$$



$p(y_t | y_{1:t-1})$ is normaliser of filter/forward algorithm update

How can we compute the smoothing estimate?

$$p(x_t | y_{1:T})$$

LGSSM: Kalman Smoother
HMM: Forward-Backward= Algorithm

How can we compute the most probable sequence?

$$x'_{1:T} = \arg \max_{x_{1:T}} p(x_{1:T} | y_{1:T})$$

LGSSM: Kalman Smoother
HMM: Viterbi Decoding

The magic of the Forward Algorithm: Dynamic Programming

What's going on here?

In discrete case, likelihood involves sum over all sequences: $x_{1:T}^{(k)}$

$$p(y_{1:T}) = \sum_{\text{all sequences } k} p(y_{1:T}, x_{1:T}^{(k)})$$

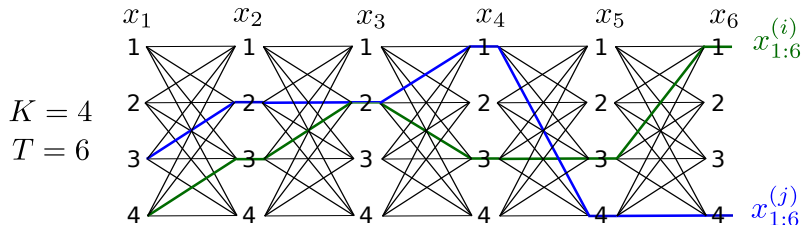
The magic of the Forward Algorithm: Dynamic Programming

What's going on here?

In discrete case, likelihood involves sum over all sequences: $x_{1:T}^{(k)}$

$$p(y_{1:T}) = \sum_{\text{all sequences } k} p(y_{1:T}, x_{1:T}^{(k)})$$

Trellis diagram represents possible sequences:



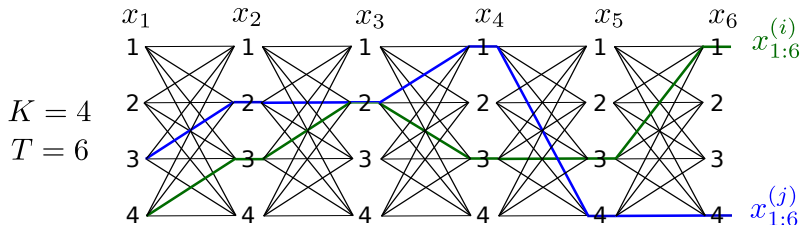
The magic of the Forward Algorithm: Dynamic Programming

What's going on here?

In discrete case, likelihood involves sum over all sequences: $x_{1:T}^{(k)}$

$$p(y_{1:T}) = \sum_{\text{all sequences } k} p(y_{1:T}, x_{1:T}^{(k)})$$

Trellis diagram represents possible sequences:



Exponential number of sequences: K^T

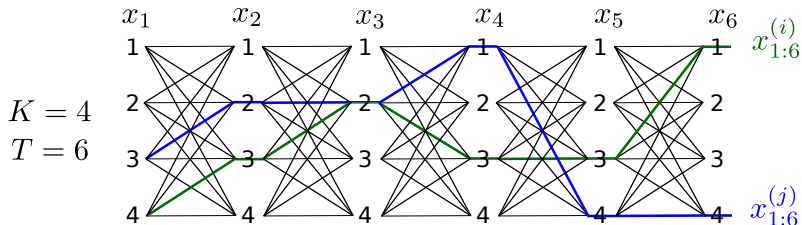
The magic of the Forward Algorithm: Dynamic Programming

What's going on here?

In discrete case, likelihood involves sum over all sequences: $x_{1:T}^{(k)}$

$$p(y_{1:T}) = \sum_{\text{all sequences } k} p(y_{1:T}, x_{1:T}^{(k)})$$

Trellis diagram represents possible sequences:



Exponential number of sequences: K^T

But Forward algorithm had linear complexity in time (loop over t)

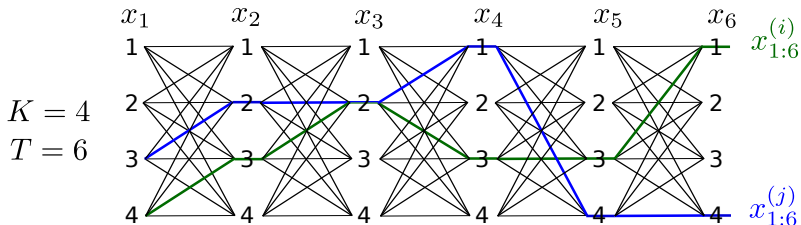
The magic of the Forward Algorithm: Dynamic Programming

What's going on here?

In discrete case, likelihood involves sum over all sequences: $x_{1:T}^{(k)}$

$$p(y_{1:T}) = \sum_{\text{all sequences } k} p(y_{1:T}, x_{1:T}^{(k)})$$

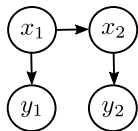
Trellis diagram represents possible sequences:



Exponential number of sequences: K^T

But Forward algorithm had linear complexity in time (loop over t)

Markov property means we can forget history of previous states:
just remember last one (dynamic programming/belief propagation)



Maximum Likelihood Learning of HMMs: simple once inference is solved

log-likelihood:
$$\log p(y_{1:T}|\theta) = \log \int p(y_{1:T}, x_{1:T}|\theta) dx_{1:T}$$

Maximum Likelihood Learning of HMMs: simple once inference is solved

log-likelihood: $\log p(y_{1:T}|\theta) = \log \int p(y_{1:T}, x_{1:T}|\theta) dx_{1:T}$

gradient of
log-likelihood: $\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \frac{1}{p(y_{1:T}|\theta)} \int \frac{d}{d\theta} p(y_{1:T}, x_{1:T}|\theta) dx_{1:T}$

Maximum Likelihood Learning of HMMs: simple once inference is solved

log-likelihood: $\log p(y_{1:T}|\theta) = \log \int p(y_{1:T}, x_{1:T}|\theta) dx_{1:T}$

gradient of
log-likelihood: $\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \frac{1}{p(y_{1:T}|\theta)} \int \frac{d}{d\theta} p(y_{1:T}, x_{1:T}|\theta) dx_{1:T}$

show gradient depends
on simple moments
of posterior:

$$\frac{d}{d\theta} \log p(y_{1:T}|\theta)$$

Maximum Likelihood Learning of HMMs: simple once inference is solved

log-likelihood: $\log p(y_{1:T}|\theta) = \log \int p(y_{1:T}, x_{1:T}|\theta) dx_{1:T}$

gradient of
log-likelihood: $\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \frac{1}{p(y_{1:T}|\theta)} \int \frac{d}{d\theta} p(y_{1:T}, x_{1:T}|\theta) dx_{1:T}$

show gradient depends
on simple moments
of posterior:

$$\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \frac{1}{p(y_{1:T}|\theta)} \int \frac{d}{d\theta} \exp(\log p(y_{1:T}, x_{1:T}|\theta)) dx_{1:T}$$

Maximum Likelihood Learning of HMMs: simple once inference is solved

log-likelihood: $\log p(y_{1:T}|\theta) = \log \int p(y_{1:T}, x_{1:T}|\theta) dx_{1:T}$

gradient of
log-likelihood: $\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \frac{1}{p(y_{1:T}|\theta)} \int \frac{d}{d\theta} p(y_{1:T}, x_{1:T}|\theta) dx_{1:T}$

show gradient depends
on simple moments
of posterior:

$$\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \frac{1}{p(y_{1:T}|\theta)} \int \frac{d}{d\theta} \exp(\overbrace{\log p(y_{1:T}, x_{1:T}|\theta)}^{E(\theta; x_{1:T}, y_{1:T})}) dx_{1:T}$$

Maximum Likelihood Learning of HMMs: simple once inference is solved

log-likelihood: $\log p(y_{1:T}|\theta) = \log \int p(y_{1:T}, x_{1:T}|\theta) dx_{1:T}$

gradient of
log-likelihood: $\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \frac{1}{p(y_{1:T}|\theta)} \int \frac{d}{d\theta} p(y_{1:T}, x_{1:T}|\theta) dx_{1:T}$

show gradient depends
on simple moments
of posterior:

$$E(\theta; x_{1:T}, y_{1:T}) = \sum_t [\log p(y_t|x_t, \theta) + \log p(x_t|x_{t-1}, \theta)]$$

$$E(\theta; x_{1:T}, y_{1:T})$$

$$\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \frac{1}{p(y_{1:T}|\theta)} \int \frac{d}{d\theta} \exp(\overbrace{\log p(y_{1:T}, x_{1:T}|\theta)}) dx_{1:T}$$

Maximum Likelihood Learning of HMMs: simple once inference is solved

log-likelihood: $\log p(y_{1:T}|\theta) = \log \int p(y_{1:T}, x_{1:T}|\theta) dx_{1:T}$

gradient of log-likelihood: $\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \frac{1}{p(y_{1:T}|\theta)} \int \frac{d}{d\theta} p(y_{1:T}, x_{1:T}|\theta) dx_{1:T}$

simple form: e.g. quadratic in x for LGSSMs

$$E(\theta; x_{1:T}, y_{1:T}) = \sum_t [\log p(y_t|x_t, \theta) + \log p(x_t|x_{t-1}, \theta)]$$

show gradient depends on simple moments of posterior:

$$E(\theta; x_{1:T}, y_{1:T})$$

$$\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \frac{1}{p(y_{1:T}|\theta)} \int \frac{d}{d\theta} \exp(\overbrace{\log p(y_{1:T}, x_{1:T}|\theta)}) dx_{1:T}$$

Maximum Likelihood Learning of HMMs: simple once inference is solved

log-likelihood: $\log p(y_{1:T}|\theta) = \log \int p(y_{1:T}, x_{1:T}|\theta) dx_{1:T}$

gradient of log-likelihood: $\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \frac{1}{p(y_{1:T}|\theta)} \int \frac{d}{d\theta} p(y_{1:T}, x_{1:T}|\theta) dx_{1:T}$

simple form: e.g. quadratic in x for LGSSMs

$$E(\theta; x_{1:T}, y_{1:T}) = \sum_t [\log p(y_t|x_t, \theta) + \log p(x_t|x_{t-1}, \theta)]$$

show gradient depends on simple moments of posterior:

$$E(\theta; x_{1:T}, y_{1:T})$$

$$\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \frac{1}{p(y_{1:T}|\theta)} \int \frac{d}{d\theta} \exp(\overbrace{\log p(y_{1:T}, x_{1:T}|\theta)}) dx_{1:T}$$

$$\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \frac{1}{p(y_{1:T}|\theta)} \int p(y_{1:T}, x_{1:T}|\theta) \frac{d}{d\theta} E(\theta; x_{1:T}, y_{1:T}) dx_{1:T}$$

Maximum Likelihood Learning of HMMs: simple once inference is solved

log-likelihood: $\log p(y_{1:T}|\theta) = \log \int p(y_{1:T}, x_{1:T}|\theta) dx_{1:T}$

gradient of
log-likelihood: $\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \frac{1}{p(y_{1:T}|\theta)} \int \frac{d}{d\theta} p(y_{1:T}, x_{1:T}|\theta) dx_{1:T}$

simple form: e.g. quadratic in x for LGSSMs

$$E(\theta; x_{1:T}, y_{1:T}) = \sum_t [\log p(y_t|x_t, \theta) + \log p(x_t|x_{t-1}, \theta)]$$

show gradient depends
on simple moments
of posterior:

$$E(\theta; x_{1:T}, y_{1:T})$$

$$\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \frac{1}{p(y_{1:T}|\theta)} \int \frac{d}{d\theta} \exp(\log p(y_{1:T}, x_{1:T}|\theta)) dx_{1:T}$$

$$\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \frac{1}{p(y_{1:T}|\theta)} \int p(y_{1:T}, x_{1:T}|\theta) \frac{d}{d\theta} E(\theta; x_{1:T}, y_{1:T}) dx_{1:T}$$

$$\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \int p(x_{1:T}|y_{1:T}, \theta) \frac{d}{d\theta} E(\theta; x_{1:T}, y_{1:T}) dx_{1:T}$$

Maximum Likelihood Learning of HMMs: simple once inference is solved

log-likelihood: $\log p(y_{1:T}|\theta) = \log \int p(y_{1:T}, x_{1:T}|\theta) dx_{1:T}$

gradient of log-likelihood: $\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \frac{1}{p(y_{1:T}|\theta)} \int \frac{d}{d\theta} p(y_{1:T}, x_{1:T}|\theta) dx_{1:T}$

simple form: e.g. quadratic in x for LGSSMs

$$E(\theta; x_{1:T}, y_{1:T}) = \sum_t [\log p(y_t|x_t, \theta) + \log p(x_t|x_{t-1}, \theta)]$$

show gradient depends
on simple moments
of posterior:

$$E(\theta; x_{1:T}, y_{1:T})$$

$$\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \frac{1}{p(y_{1:T}|\theta)} \int \frac{d}{d\theta} \exp(\log p(y_{1:T}, x_{1:T}|\theta)) dx_{1:T}$$

$$\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \frac{1}{p(y_{1:T}|\theta)} \int p(y_{1:T}, x_{1:T}|\theta) \frac{d}{d\theta} E(\theta; x_{1:T}, y_{1:T}) dx_{1:T}$$

$$\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \int p(x_{1:T}|y_{1:T}, \theta) \frac{d}{d\theta} E(\theta; x_{1:T}, y_{1:T}) dx_{1:T}$$

$$\frac{d}{d\theta} \log p(y_{1:T}|\theta) = \left\langle \frac{d}{d\theta} E(\theta; x_{1:T}, y_{1:T}) \right\rangle_{p(x_{1:T}|y_{1:T}, \theta)}$$

requires posterior moments: marginals and pairwise marginals

Course Survey: please complete this!