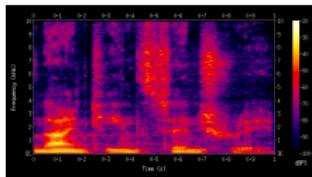


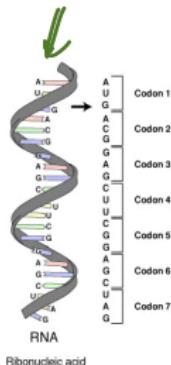
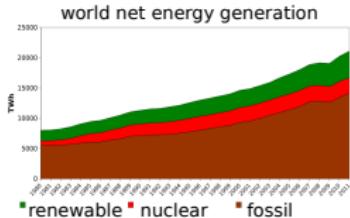
Sequence Modelling

Rich Turner

Sequence data



Some images taken from wikipedia

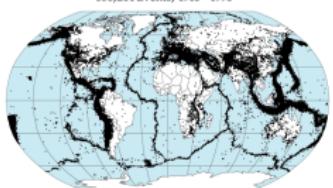


Ribonucleic acid



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



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

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

Remove noise from a sequence

$$p(\underbrace{y'_1, \dots, y'_t}_\text{clean} | \underbrace{y_1, \dots, y_t}_\text{mix})$$

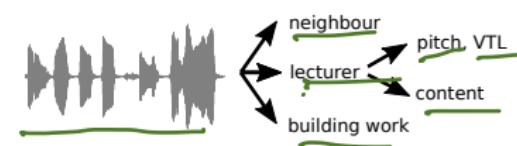
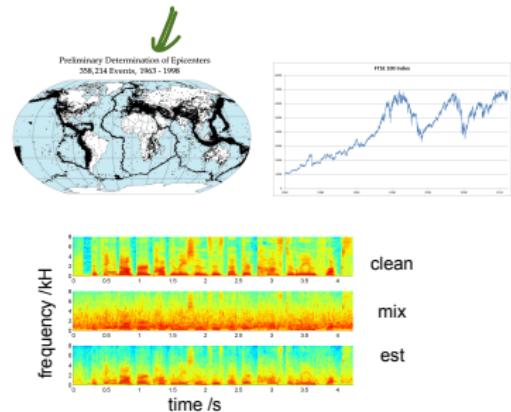
Predict one sequence from another

$$p(\underbrace{y'_1, \dots, y'_t}_\text{est} | \underbrace{y_1, \dots, y_t}_\text{mix})$$

Discover underlying latent variables

$$p(\underbrace{x_1, \dots, x_t}_\text{content} | \underbrace{y_1, \dots, y_t}_\text{building work})$$

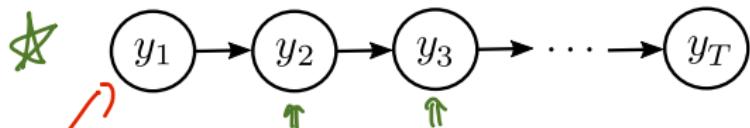
observing



Markov models

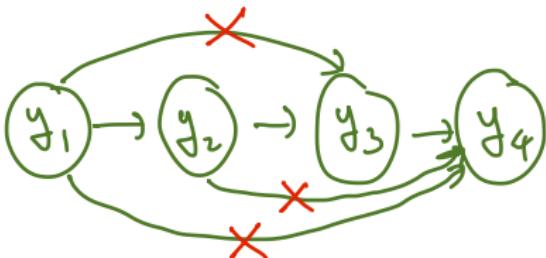
First order Markov

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



$T=4$

$$p(y_1, y_2, y_3, y_4) = p(y_1) p(y_2|y_1) p(y_3|y_2, \cancel{y_1}) p(y_4|y_3, \cancel{y_1}, \cancel{y_2})$$



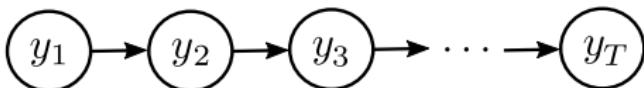
1st order Markov

$$y_t \perp\!\!\!\perp y_{t-2} | y_{t-1}$$

Markov models

First order Markov

$$p(y_1, y_2, y_3, \dots, y_T) = p(y_1) \underbrace{p(y_2|y_1)}_{\text{parameters tied}} \underbrace{p(y_3|y_2)}_{\text{parameters tied}} \dots p(y_T|y_{T-1})$$



$$\underset{\mathcal{N}}{\times} G \left[\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_T \end{bmatrix}; \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_T \end{bmatrix}, \Sigma_{1:T, 1:T}^{T \times T} \right]$$

Gaussian

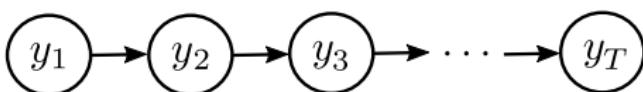
N

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
∞ number of variables
finite number of parameters

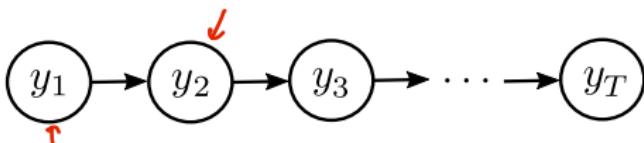
Markov model = conditional independence relationship + product rule

$$\text{future } \xrightarrow{\text{independent of past}} y_{t+1} \perp y_{1:t-1} | y_t \quad \xleftarrow{\text{given present}}$$
$$p(y_{1:T}) = \prod_{t=1}^T p(y_t | y_{1:t-1})$$

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

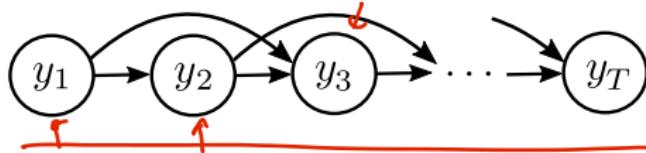
Markov model = conditional independence relationship + product rule

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

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

Second order Markov

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

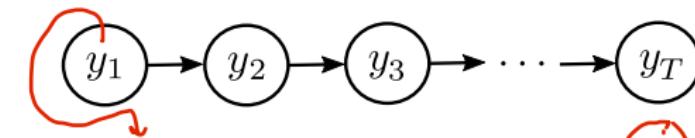


Markov models for discrete data: n-gram models

NLP

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

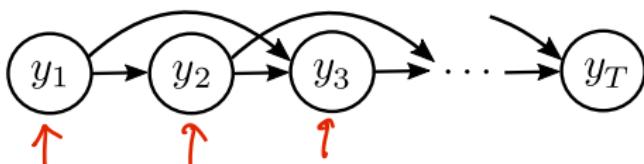
$$p(y_t = k | y_{t-1} = l) = T_{k,l}$$

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?

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

↑ ↑ ↓
 π_l^0 T_{kl}

sum product

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

$$p(y_2 = k) = \left[\begin{array}{c} \underline{\pi} \\ \underline{T} \end{array} \right]_k$$

3 F 3

$$\underline{S} = \underline{T}^T$$

$$\underline{\Pi}^0 = \underline{\pi}_0^T$$

$$P^0 \swarrow \searrow S_{z,L}^T$$

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?

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

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

$$p(y_\infty = k) = (\underline{\pi^\infty})_k$$

$$\underline{\pi^\infty} = \underline{\pi^\infty} \lambda \quad \lambda = 1$$

$$\begin{aligned} p(y_t = k) &= \sum_{l=1}^K p(y_t = k | y_{t-1} = l) p(y_{t-1} = l) \\ &= \sum_{l=1}^K T_{k,l} \underline{\pi^\infty}_l \end{aligned}$$

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?

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

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

Counts

to State

from state		
	A	B
A	2 6 2	0
B	3 4 1	1
C	0 2	2

= l) $\Sigma T_{k,l}^3$

First order Markov (bi-gram)

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

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

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

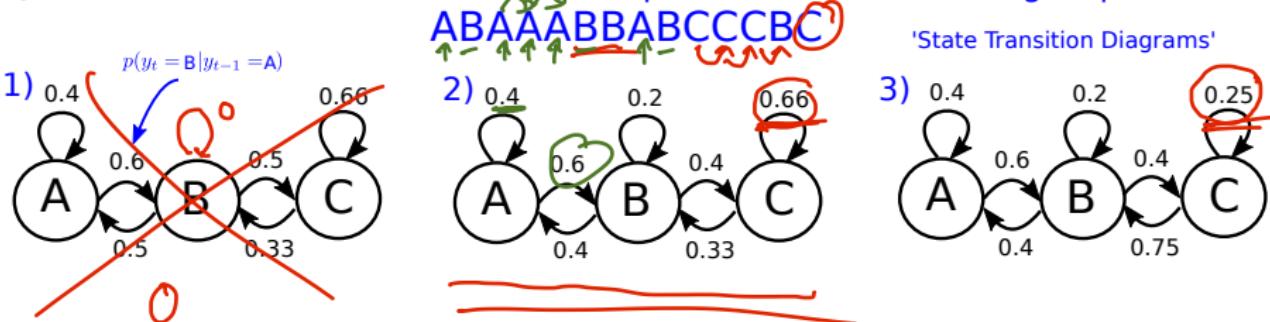
Q2. How can we compute the stationary distribution for the Markov chain?

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

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?



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?

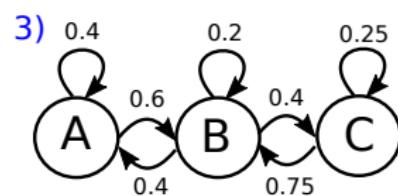
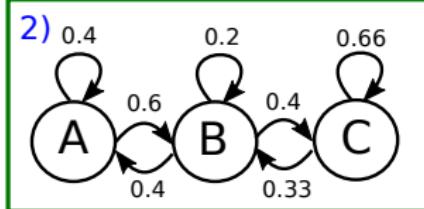
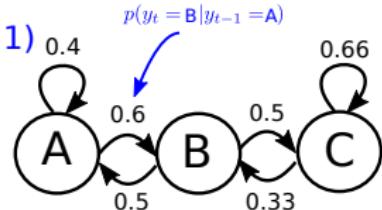
$$p(y_t = k) = \sum_{l=1}^K p(y_t = k | y_{t-1} = l) p(y_{t-1} = l)$$

eigenvectors of transition matrix with eigenvalue = 1

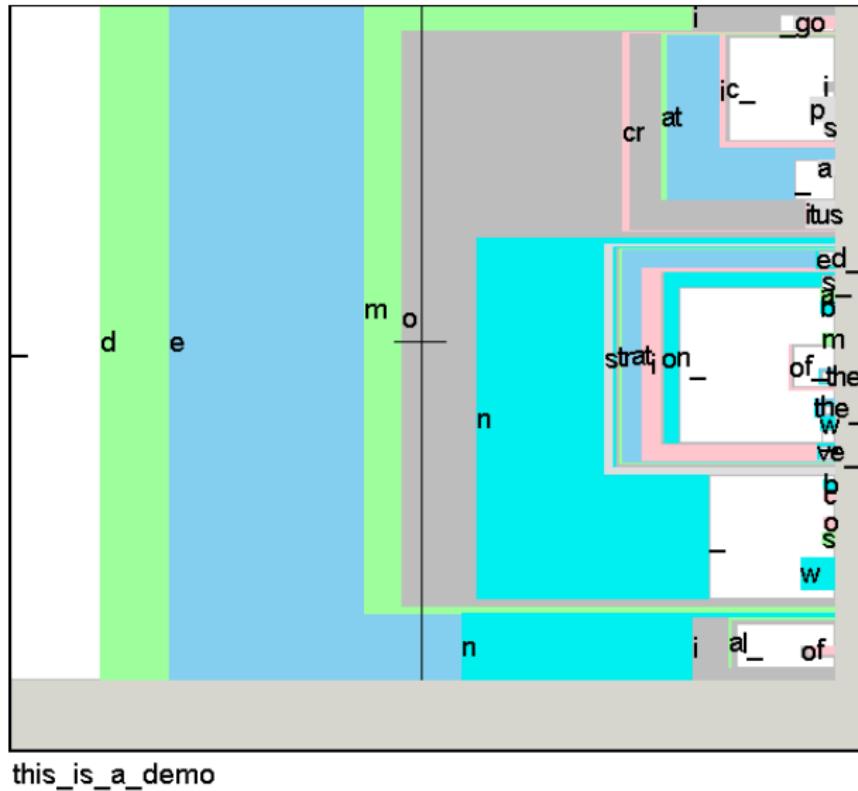
$$\pi_k^\infty = \sum_{l=1}^K T_{k,l} \pi_l^\infty$$

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

ABAAABBABCCCB



Example application of n-grams: text modelling for dasher



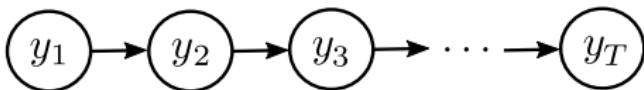
<http://www.inference.phy.cam.ac.uk/dasher/>

<https://www.youtube.com/watch?v=nr3s4613DX8>

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

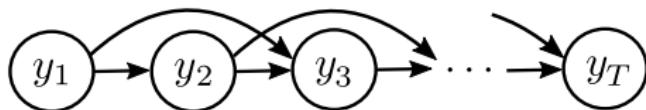
$$p(y_t = k | y_{t-1} = l) = T_{k,l}$$

transition probabilities
(stochastic matrix)

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

Second order Markov (tri-gram)

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



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

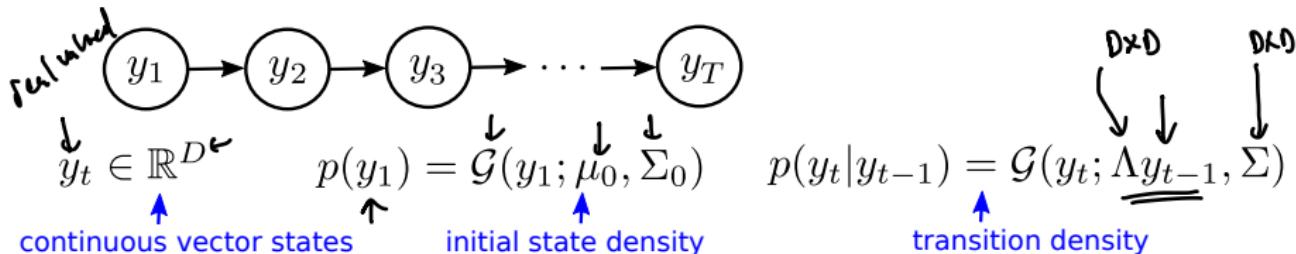
n-grams require large
multidimensional arrays

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

Markov order

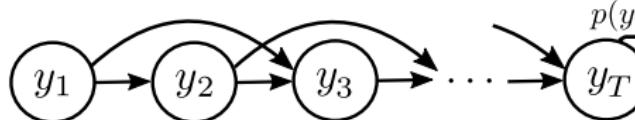
First order Markov (AR(1))

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



Second order Markov (AR(2))

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



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

joint distribution over all variables
is always multivariate Gaussian

$$p(y_{1:T}) = \mathcal{G}\left(y_{1:T}; \mu_{1:T}, \Sigma_{1:T, 1:T}\right)$$

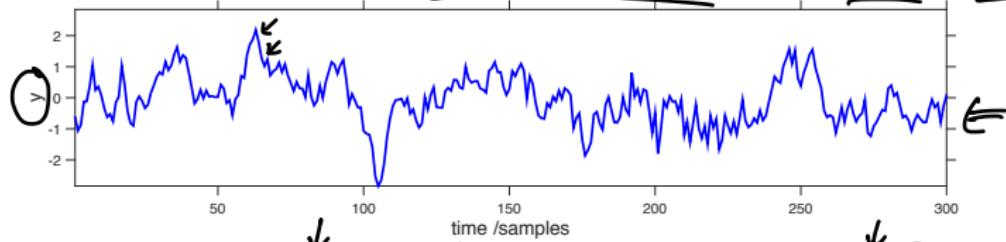
Markov models for continuous data: Auto-Regressive (AR) Gaussian models

$D=1$
First order Markov (AR(1))

$$y_t = \lambda y_{t-1} + \varepsilon_t \sigma \quad \varepsilon_t \sim N(0,1)$$

small

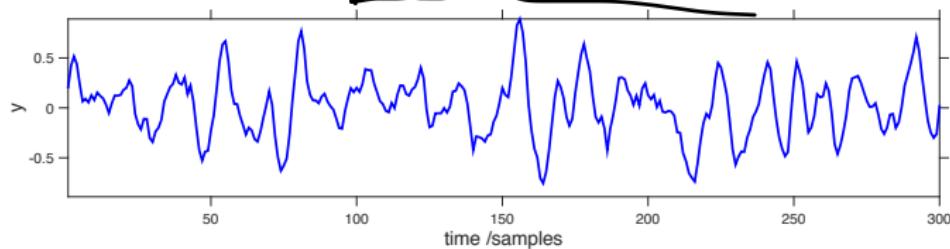
$$y_t \in \mathbb{R}^1 \quad p(y_t | y_{t-1}) = \mathcal{G}(y_t; \lambda y_{t-1}, \sigma^2) \quad \lambda = 0.9 \quad \sigma^2 = 0.01$$



Second order Markov (AR(2))

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

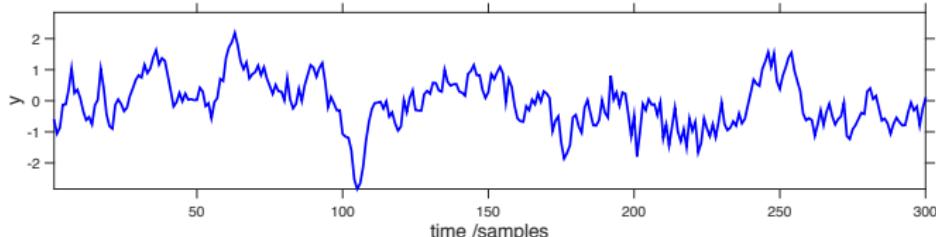
$$[\lambda_1, \lambda_2] = [1.57, -0.78] \quad \sigma^2 = 0.01$$



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

First order Markov (AR(1))

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



D
A
L

What is the stationary distribution of this process? $p(y_\infty) = ?$ $\mathcal{G}(y_\infty; \mu_\infty, \sigma_\infty^2)$

$$\mathbb{E}_{p(y_t)}(y_t) = \mathbb{E}_{p(y_{t-1}, \varepsilon_t)}(\lambda y_{t-1} + \varepsilon_t \sigma) \quad \underline{y_\infty \rightarrow 0}$$

$$\mu_t = \lambda \mu_{t-1} + 0.0 \Rightarrow \mu_t = \lambda \mu_{t-1}$$

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

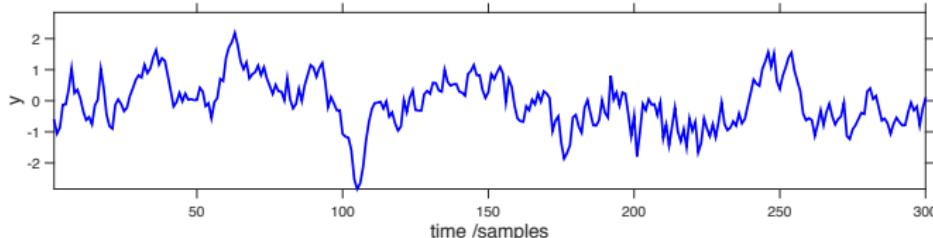
Markov models for continuous data: Auto-Regressive (AR) Gaussian models

$$y_t = \lambda y_{t-1} + \varepsilon_t \sigma$$

↑ ↑ ↑

First order Markov (AR(1))

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



$$\mathbb{E}(y_{t-1}) \mathbb{E}(\varepsilon_t)$$

↓ ↓

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

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

$$\mathbb{E}(y_t^2) = \mathbb{E}((\lambda y_{t-1} + \varepsilon_t \sigma)^2) = \lambda^2 \mathbb{E}(y_{t-1}^2) + 2\mathbb{E}(y_{t-1} \varepsilon_t) + \mathbb{E}(\varepsilon_t^2) \sigma^2$$

↓ ↓ ↓

$$= \lambda^2 \mathbb{E}(y_{t-1}^2) + \sigma^2$$

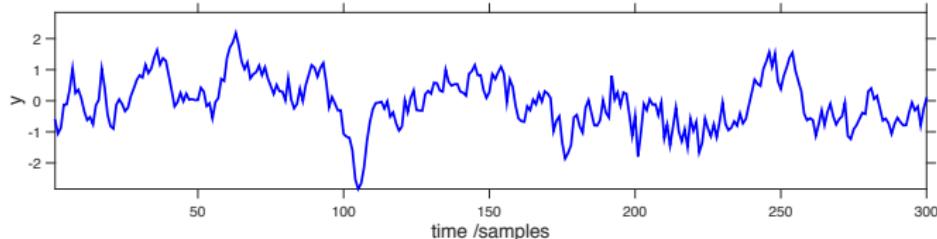
$$\sigma_\infty^2 = \lambda^2 \sigma_0^2 + \sigma^2$$

$$\sigma_0^2 = \frac{\sigma^2}{1-\lambda^2} = \frac{0.01}{1-0.9^2}$$

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? $p(y_\infty) = ?$

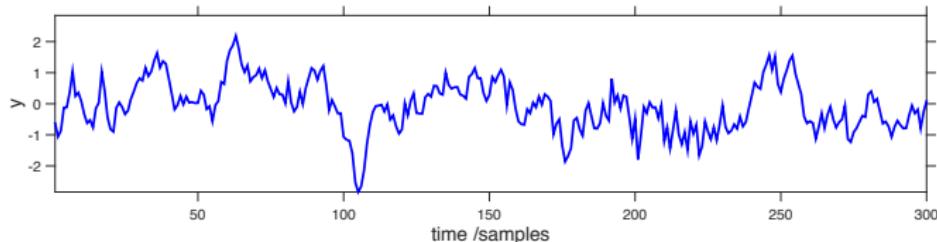
Everything is linear Gaussian \Rightarrow must be Gaussian $p(y_\infty) = \mathcal{G}(y_\infty; \mu_\infty, \sigma_\infty^2)$

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

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 $p(y_\infty) = \mathcal{G}(y_\infty; \mu_\infty, \sigma_\infty^2)$

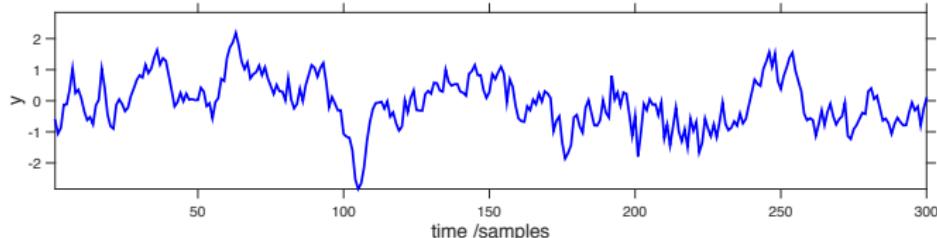
$$y_t = \lambda y_{t-1} + \sigma \epsilon_t \quad \epsilon_t \sim \mathcal{G}(0, 1)$$

Mean: $\langle y_t \rangle$

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

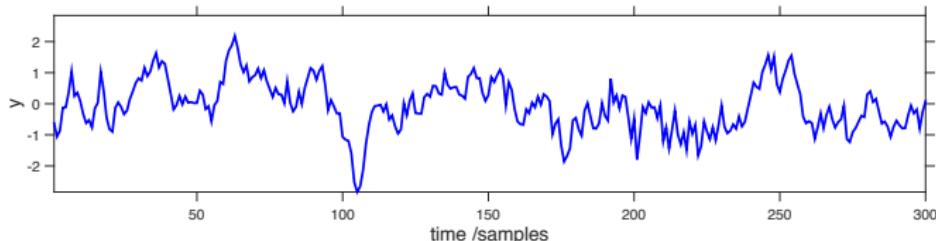
$$y_t = \lambda y_{t-1} + \sigma \epsilon_t \quad \epsilon_t \sim \mathcal{G}(0, 1)$$

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

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

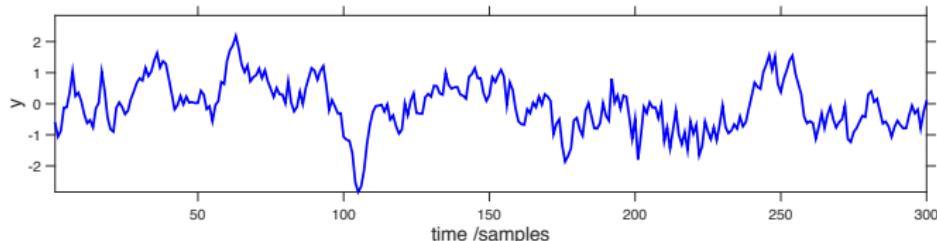
$$y_t = \lambda y_{t-1} + \sigma \epsilon_t \quad \epsilon_t \sim \mathcal{G}(0, 1)$$

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

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

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

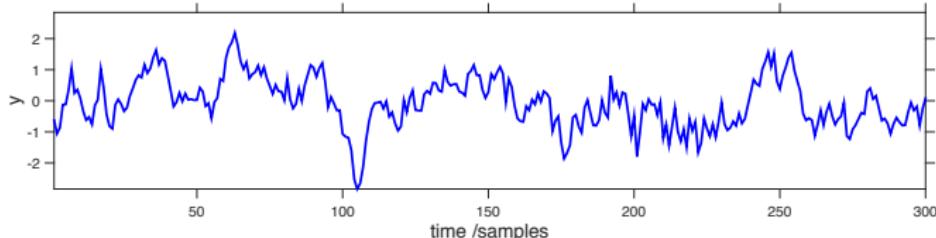
Mean: $\langle y_t \rangle = \lambda \langle y_{t-1} \rangle + \sigma \langle \epsilon_t \rangle = 0 \quad \mu_\infty = 0$

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

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

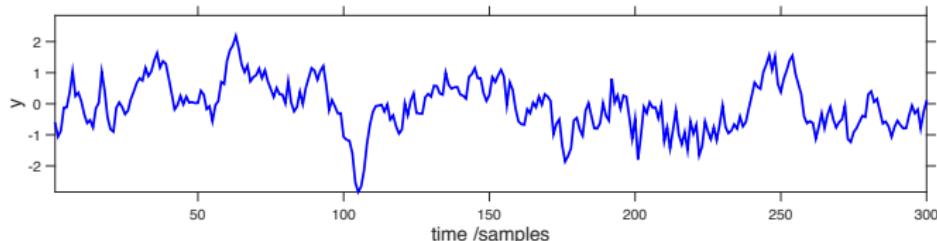
Mean: $\langle y_t \rangle = \lambda \langle y_{t-1} \rangle + \sigma \langle \epsilon_t \rangle = 0 \quad \mu_\infty = 0$

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

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

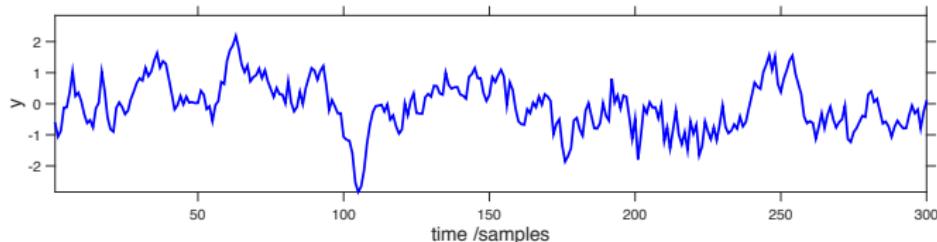
Mean: $\langle y_t \rangle = \lambda \langle y_{t-1} \rangle + \sigma \langle \epsilon_t \rangle = 0 \quad \mu_\infty = 0$

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

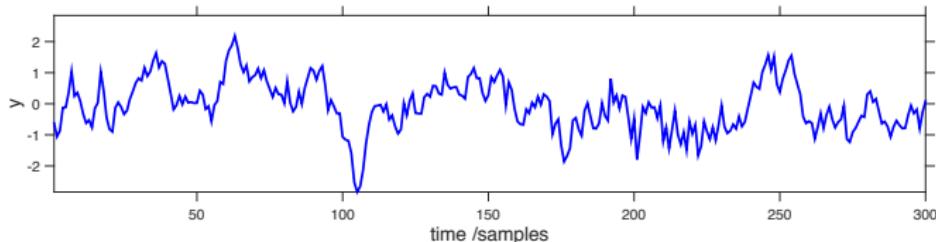
Variance: $\langle y_t^2 \rangle = \langle (\lambda y_{t-1} + \sigma \epsilon_t)^2 \rangle = \lambda^2 \langle y_{t-1}^2 \rangle + 2\lambda \sigma \langle y_{t-1} \epsilon_t \rangle + \sigma^2 \langle \epsilon_t^2 \rangle$

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

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

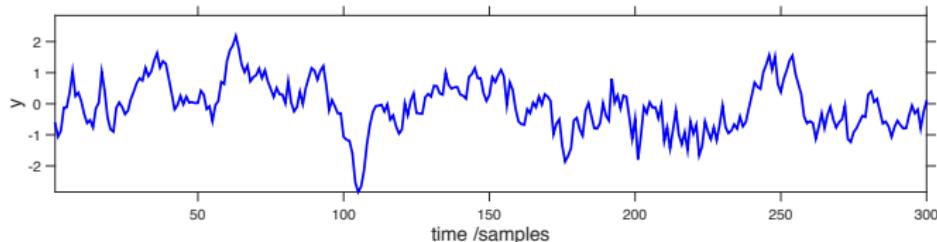
Variance: $\langle y_t^2 \rangle = \langle (\lambda y_{t-1} + \sigma \epsilon_t)^2 \rangle = \lambda^2 \langle y_{t-1}^2 \rangle + 2\lambda \sigma \langle y_{t-1} \epsilon_t \rangle + \sigma^2 \langle \epsilon_t^2 \rangle$

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

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

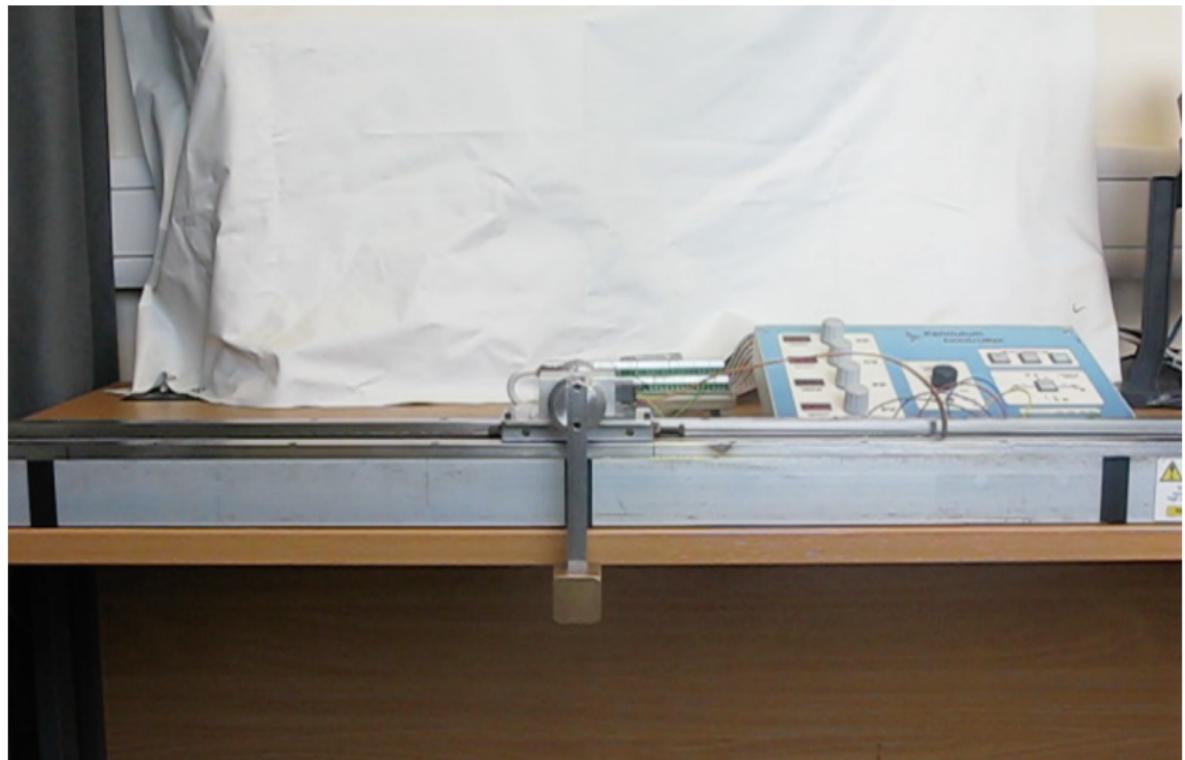
$$y_t = \lambda y_{t-1} + \sigma \epsilon_t \quad \epsilon_t \sim \mathcal{G}(0, 1)$$

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

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

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

Example application of Markov Models: pendulum swing up control problem



SUMMARY OF MARKOV MODELS

1ST ORDER

$$p(y_{1:T}) = p(y_1) \underbrace{p(y_2|y_1)}_{\text{1st-order transition}} \underbrace{p(y_3|y_2)}_{\text{1st-order transition}} \underbrace{p(y_4|y_3)}_{\text{1st-order transition}} \cdots \underbrace{p(y_T|y_{T-1})}_{\text{1st-order transition}}$$

Discrete y \Rightarrow BIGRAM Models

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

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



STATIONARY DISTRIBUTION

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

Continuous $y \Rightarrow$ AUTOREGRESSIVE Models

$$p(y_1) = G(y_1; \mu_0, \Sigma_0) \star$$

$$p(y_t | y_{t-1}) = G(y_t; \underline{\mu}_{t-1}, \underline{\Sigma}_{t-1}) \star$$

$$\dim(y_t) = 1 \quad \downarrow \quad \text{STATIONARY DISTRIBUTION}$$

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



$$\begin{aligned} \lambda &\rightarrow 1 \\ \text{Var}(y_\infty) &\rightarrow \infty \end{aligned}$$

Hidden Markov models

Real data depend on latent variables

ASR

x phonemes/words

y waveform/feature

Computer Vision

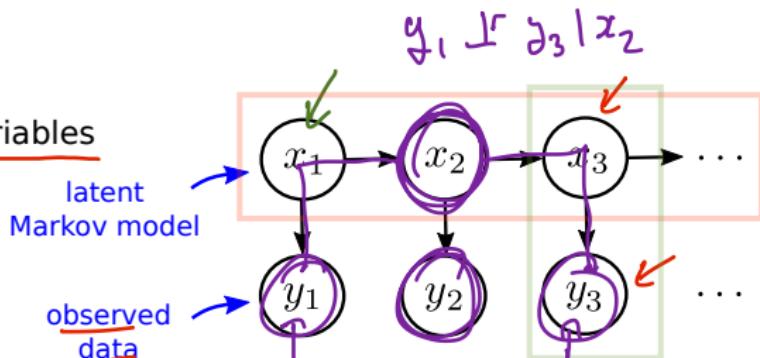
x objects, pose, lighting

y image pixel intensities

Natural Language Processing ↵

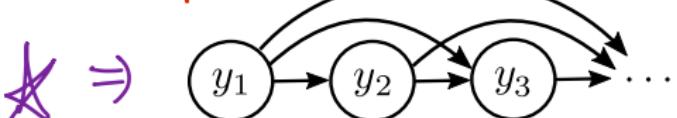
x topics

y words



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

$p(x_1|x_0) = p(x_1) = \text{initial state distribution}$



Two prevalent Examples:

⇒ Hidden Markov Models (discrete x)

⇒ Linear Gaussian State Space Models (Gaussian x and y)

$$p(y_{1:T}) = \int p(y_{1:T}, x_{1:T}) dx_{1:T}$$

Hidden Markov models: discrete hidden state

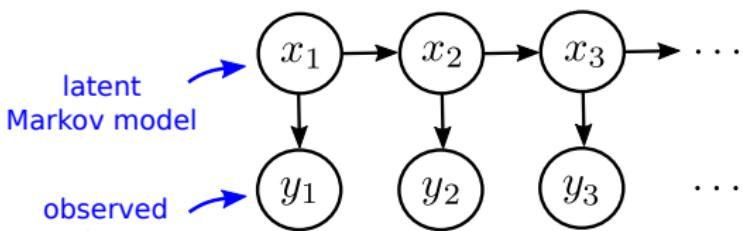
Discrete Hidden State

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

$$p(x_t = k | x_{t-1} = l) = \underline{\underline{T}_{k,l}}$$

E.g. in examples below

$$T = \begin{bmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{bmatrix} \in \begin{array}{l} x=1 \\ x=2 \end{array}$$

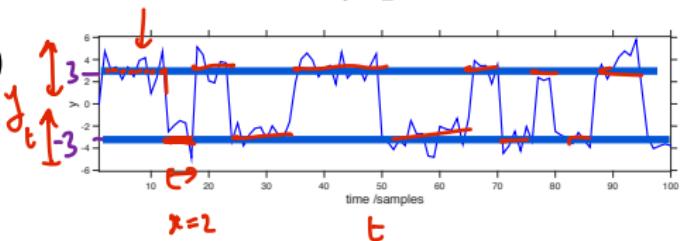


$$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t | x_{t-1}) p(y_t | x_t)$$

Continuous Observed State

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

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



Hidden Markov models: discrete hidden state

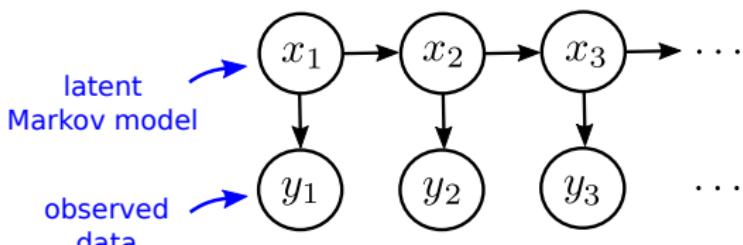
Discrete Hidden State

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

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

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

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

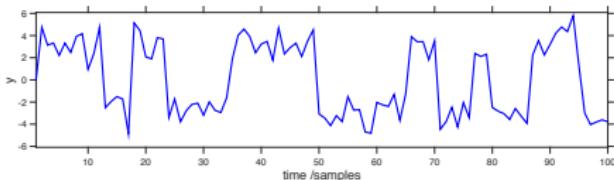


$$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t | x_{t-1}) p(y_t | x_t)$$

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

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

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

↑
emission distribution $x_t = 1$ $x_t = 2$

$$y_t \in \{A, B, C\}$$

↑ ABBBBAAABAAACCCCCCBBBCCCCCCCCCCCC
AAABBBBAABAAABBCCCCCCCCCCCCCCCCCBBA
AACCCCCCBABCBBBBBAAABBAABABCCCCC

Hidden Markov models: discrete hidden state

Discrete Hidden State, Continuous Observed State

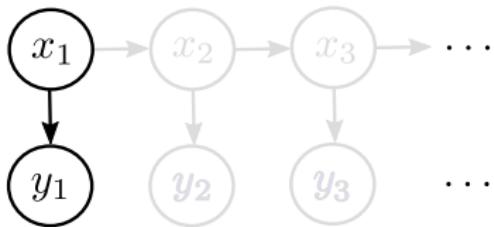
$$x_t \in \{1, \dots, K\}$$

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

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

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

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



Consider $T = 1$

$$\begin{aligned} p(y_1) &= \sum_{k=1}^K p(y_1 | x_1 = k) p(x_1 = k) \\ &= \sum_{k=1}^K \mathcal{G}(y_1; \mu_k, \Sigma_k) \pi_k^0 \end{aligned}$$

Hidden Markov models: discrete hidden state

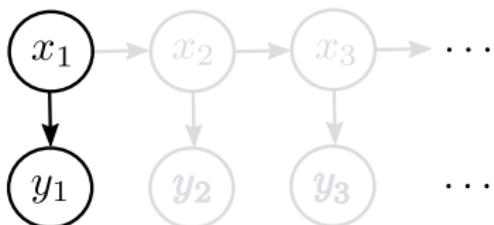
Discrete Hidden State, Continuous Observed State

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

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

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

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



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

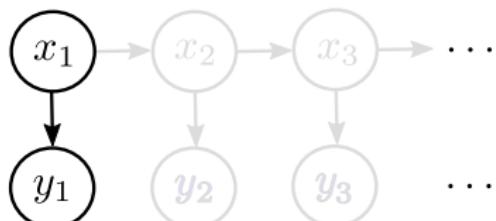
Discrete Hidden State, Continuous Observed State

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

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

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

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



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

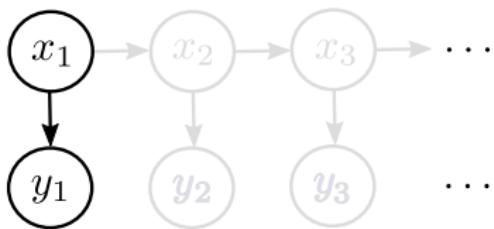
Discrete Hidden State, Continuous Observed State

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

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

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

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



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?

$$\parallel p(x_\infty = k) = \pi_k^\infty$$

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

Hidden Markov models: discrete hidden state

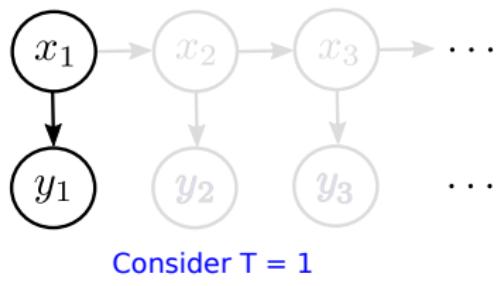
Discrete Hidden State, Continuous Observed State

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

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

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

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



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

Hidden Markov models: discrete hidden state

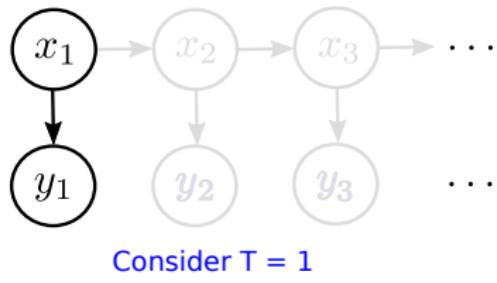
Discrete Hidden State, Continuous Observed State

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

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

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

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



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

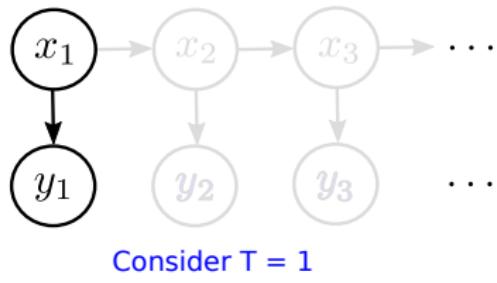
Discrete Hidden State, Continuous Observed State

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

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

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

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



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

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

Hidden Markov models: discrete hidden state

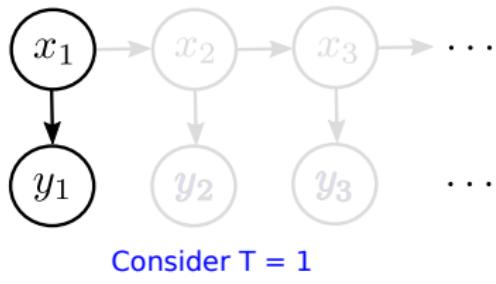
Discrete Hidden State, Continuous Observed State

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

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

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

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



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

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

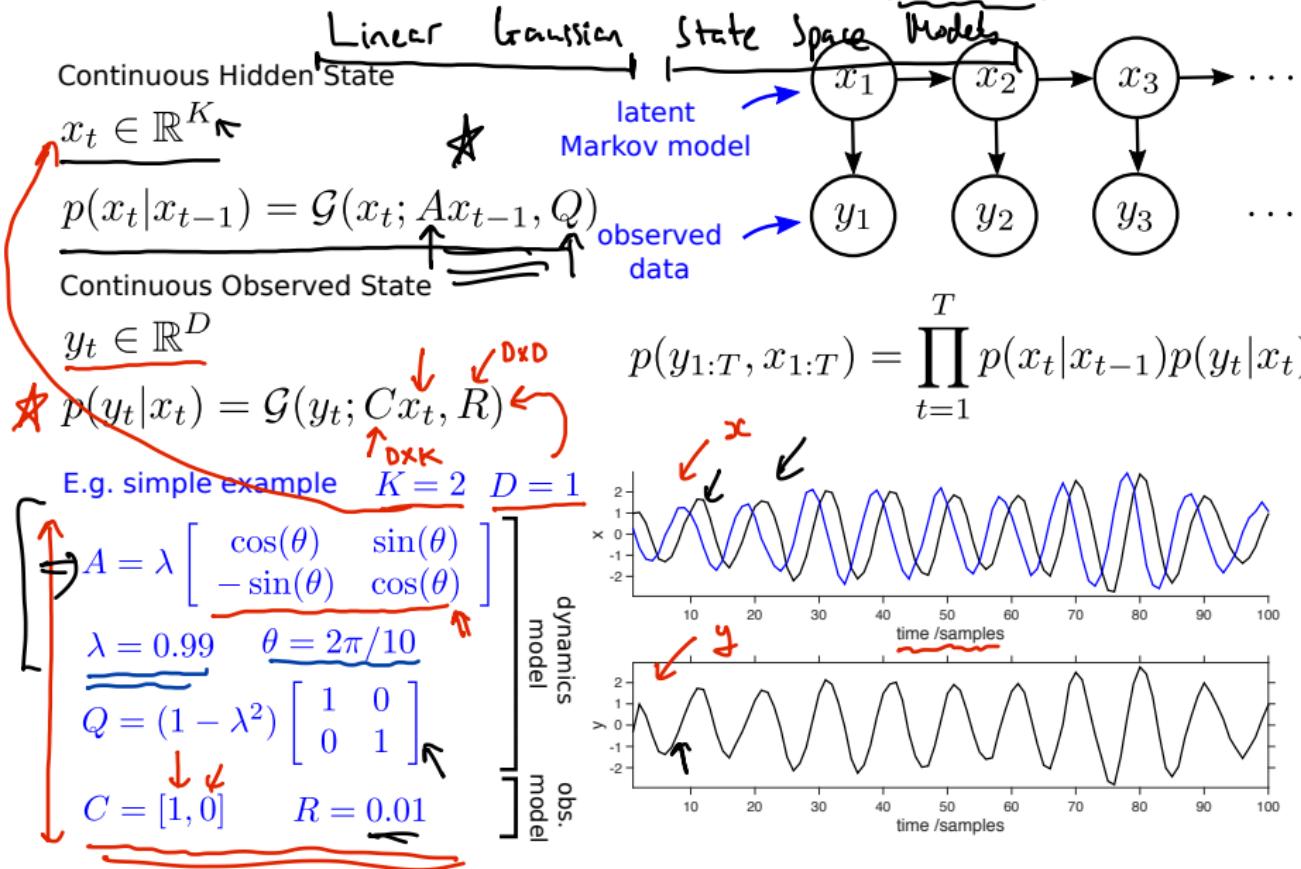
Q2: What distribution does $p(y_t)$ converge to after a long time?

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

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

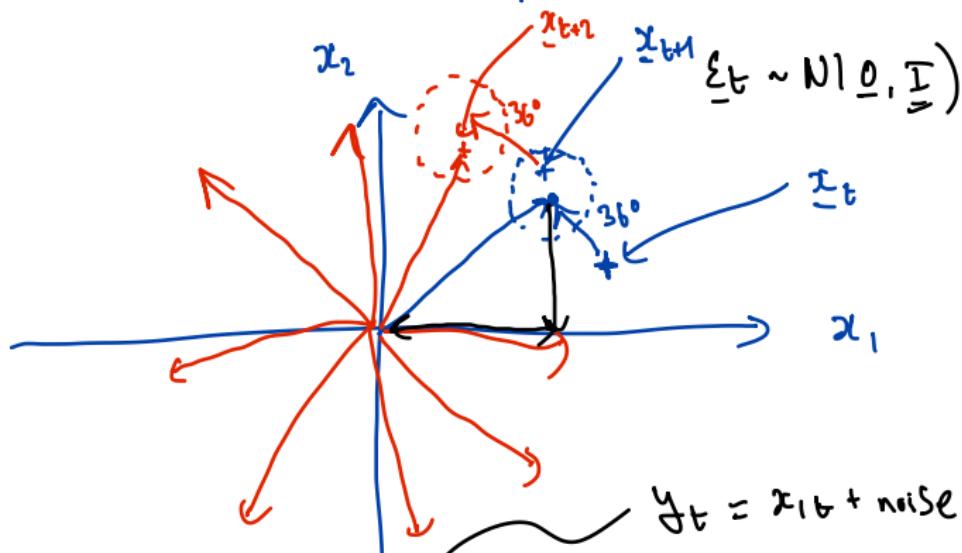
this HMM = Mixture of Gaussian Models with dynamic cluster assignments

Hidden Markov models: continuous hidden state (LGSSMs)



$$\lambda = 0.99$$

$$\begin{bmatrix} x_{1t} \\ x_{2t} \end{bmatrix} = \lambda \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} x_{1,t-1} \\ x_{2,t-1} \end{bmatrix} + (1-\lambda^2)^{\frac{1}{2}} \varepsilon_t$$

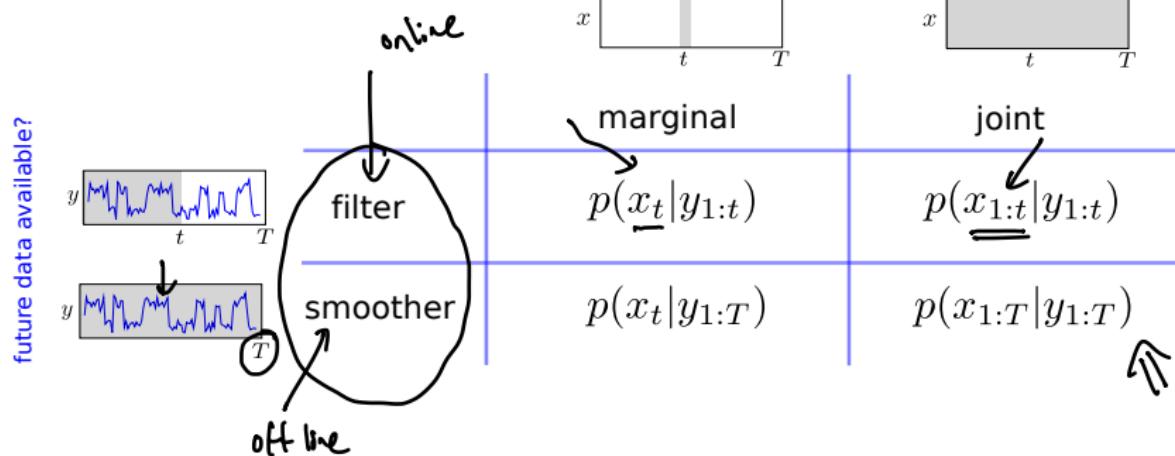


$$y_t = x_{1t} + \text{noise}$$

$$y_t = [1, 0] \begin{bmatrix} x_{1,t} \\ x_{2,t} \end{bmatrix} + \sqrt{0.01} n_t \quad n_t \sim N(0, 1)$$

Varieties of Inference

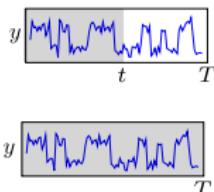
Distributional estimates



Varieties of Inference

Distributional estimates

future data available?



infer single state or sequence?

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

* Point estimates

SUMMARY SEQUENCE MODELLING LECTURE III

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

observed latent Markov ↑ ↑

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

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

$$p(y_t | x_t = k) = G(y_t; \mu_k, \Sigma_k) \leftarrow$$

$$p(x_t = l | x_t = k) = S_{lk} \leftarrow$$

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

$$y_t \in \{1 \dots D\}$$

* Continuous Hidden State $x_t \in \mathbb{R}^K$

$$\Rightarrow p(x_t | x_{t-1}) = G(x_t; A x_{t-1}^\top, Q) \leftarrow$$

$$\Rightarrow p(y_t | x_t) = G(y_t; j \leq x_t, R) \leftarrow$$

linear Gaussian state space models

$$\Leftrightarrow \begin{cases} x_t = A x_{t-1} + Q^{1/2} \varepsilon_t \\ y_t = j \leq x_t + R^{-1/2} \eta_t \end{cases}$$

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

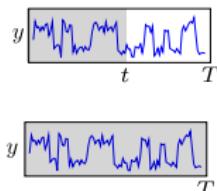
Today : Inference & Learning

$p(x | y)$

Varieties of Inference

Distributional estimates

future data available?



		infer single state or sequence?	
		x [] t T	x [] t T
		marginal	joint
filter	future data available?	$p(x_t y_{1:t})$	$p(x_{1:t} y_{1:t})$
	no future data available?	$p(x_t y_{1:T})$	$p(x_{1:T} y_{1:T})$

Point estimates

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

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? ↗

① LGSSM

$$p(x_{1:T} | y_{1:T}) = G(x_{1:T}; \underline{\mu}_{1:T}, \underline{\Sigma}_{1:T, 1:T})$$

$$\underline{x'_{1:T}} = \underset{x_{1:T}}{\operatorname{arg\,max}} p(x_{1:T} | y_{1:T}) = \underline{\mu_{1:T}}$$

$$p(x_t | y_{1:T}) = G(x_t; \underline{\mu}_t, \underline{\Sigma}_{t,t})$$

$$x_t^* = \underset{x_t}{\operatorname{arg\,max}} p(x_t | y_{1:T}) = \underline{\mu_t}$$

$$x_{1:T}^* = \underline{\mu_{1:T}}$$

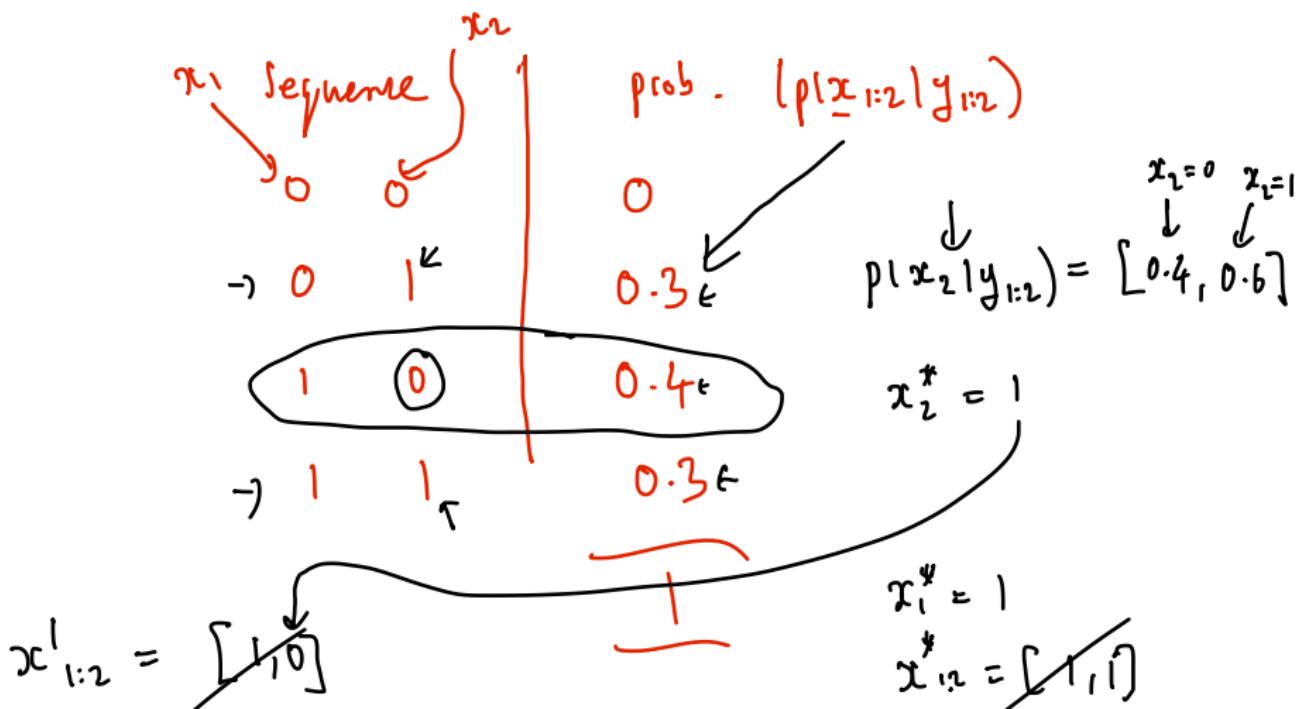


1

Discrete Hidden State HMMs

Counter example

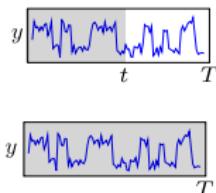
$$T = 2, K = 2$$



Varieties of Inference

Distributional estimates

future data available?



infer single state or sequence?



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

Point estimates

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

most probable state @ t most probable sequence

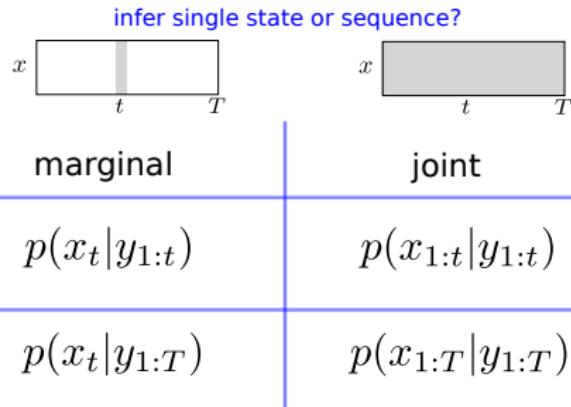
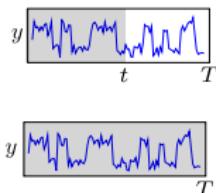
Question: are these estimates the same $x_{1:T}^* \stackrel{?}{=} x'_{1:T}$ for

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

Varieties of Inference

Distributional estimates

future data available?



Point estimates

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

most probable state @ t most probable sequence

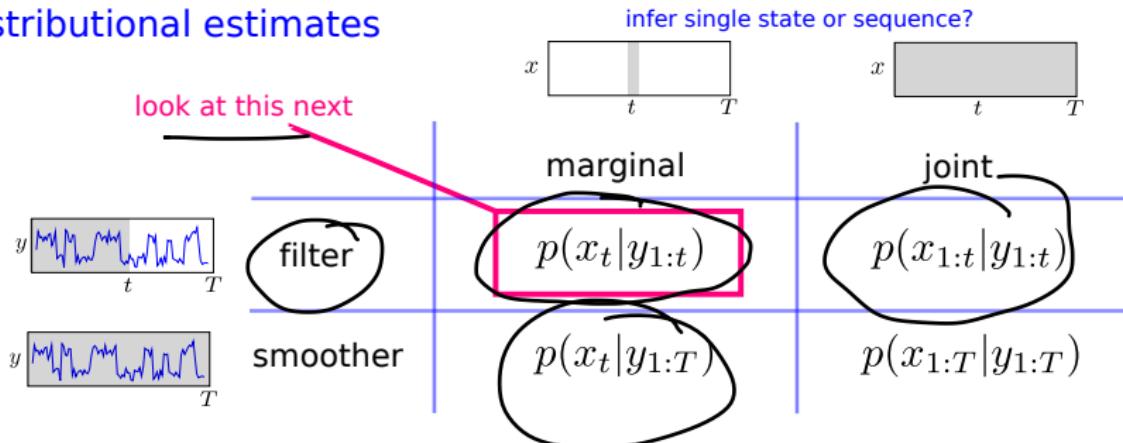
Question: are these estimates the same $x_{1:T}^* \stackrel{?}{=} x'_{1:T}$ for

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

future data available?



Point estimates

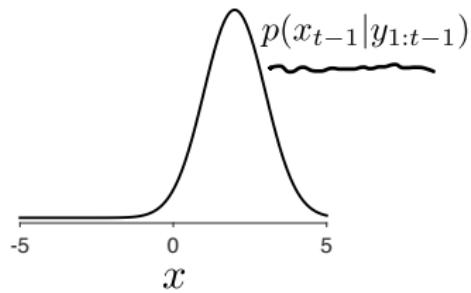
$$x_t^* = \arg \max_{x_t} p(x_t|y_{1:T}) \quad \text{most probable state @ t}$$
$$x'_{1:T} = \arg \max_{x_{1:T}} p(x_{1:T}|y_{1:T}) \quad \text{most probable sequence}$$

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? $x_{1:T}^* \neq x'_{1:T}$

Inference: Kalman Filter

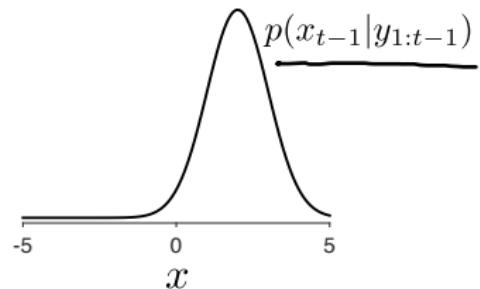
$$p(x_t | y_{1:t}) \in$$



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



Inference: Kalman Filter

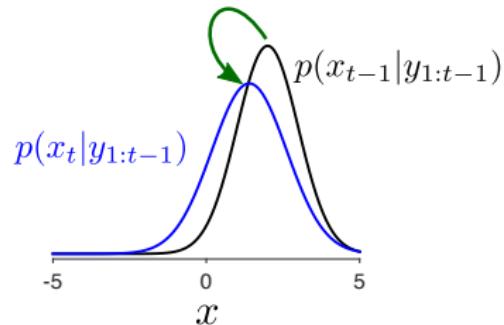


diffuse via dynamics

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

sum for discrete hidden state

Inference: Kalman Filter



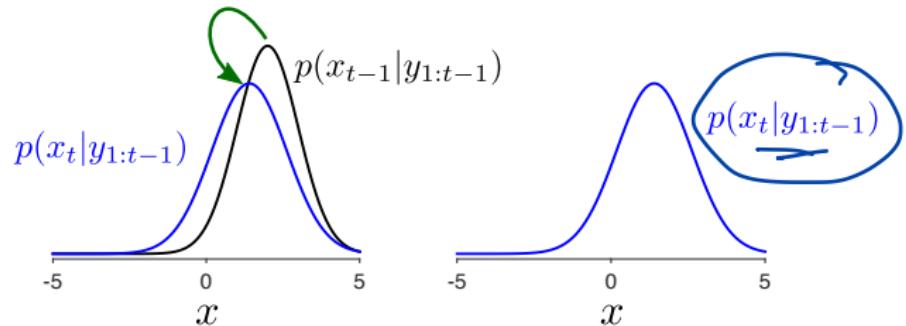
diffuse via dynamics

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

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

sum for discrete hidden state

Inference: Kalman Filter



diffuse via dynamics

combine with likelihood

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

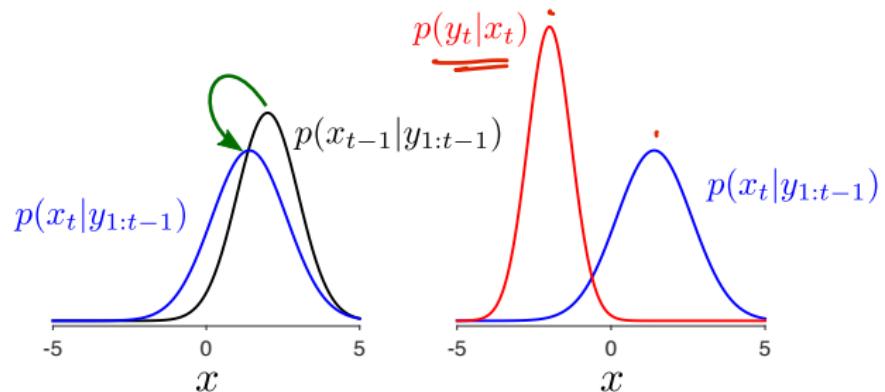
sum for discrete hidden state

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

Bayes' Rule

prior likelihood

Inference: Kalman Filter



diffuse via dynamics
combine with likelihood

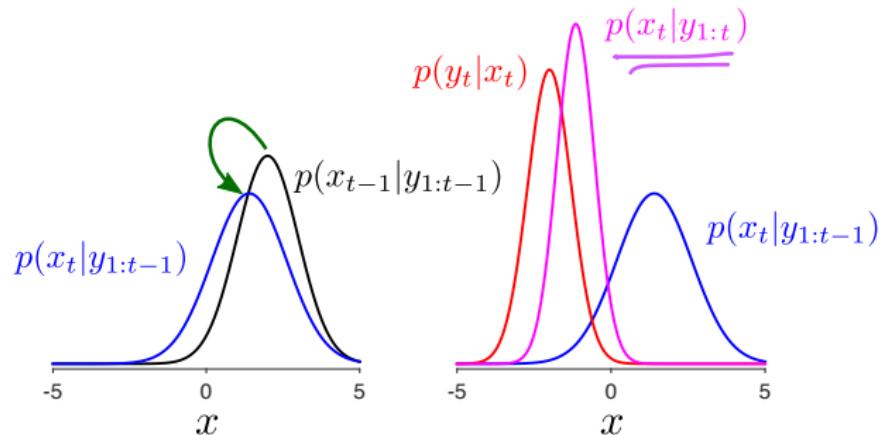
sum for discrete hidden state

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

prior likelihood

Bayes' Rule

Inference: Kalman Filter



diffuse via dynamics
combine with likelihood

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

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

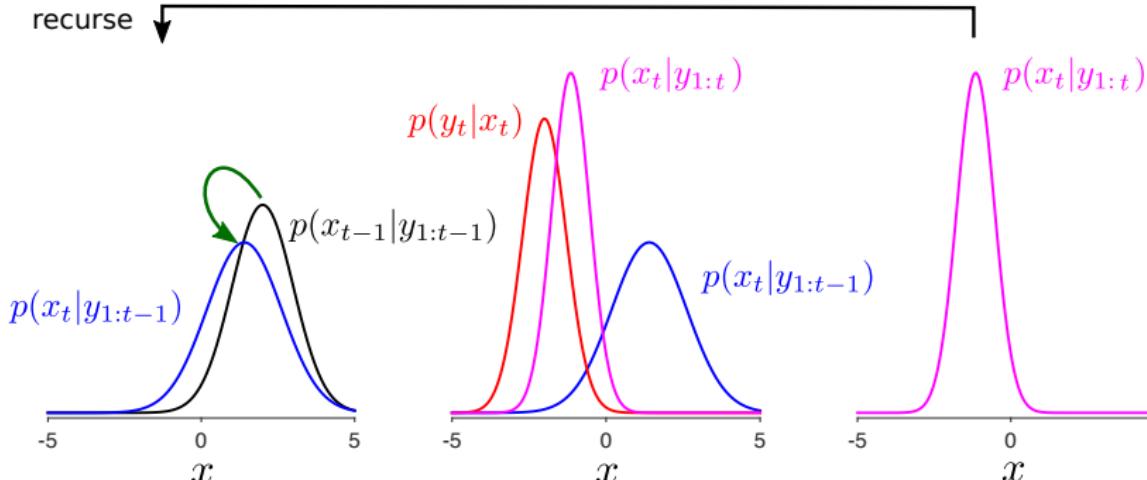
$p(x_t|y_{1:t}) \propto p(x_t|y_{1:t-1})p(y_t|x_t)$

sum for discrete hidden state

Bayes' Rule

prior likelihood

Inference: Kalman Filter



diffuse via dynamics

combine with likelihood

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

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

sum for discrete hidden state

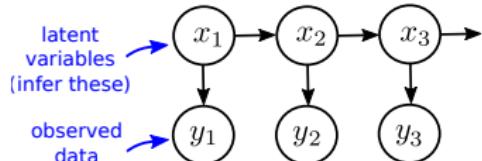
$p(x_t|y_{1:t}) \propto p(x_t|y_{1:t-1})p(y_t|x_t)$

Bayes' Rule

prior likelihood

Inference: Derivation of General Filtering Equations

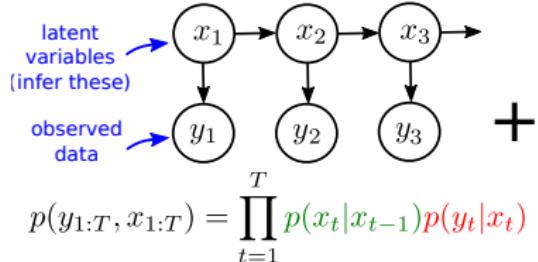
Model



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

Inference: Derivation of General Filtering Equations

Model



Rules of probability

product rule

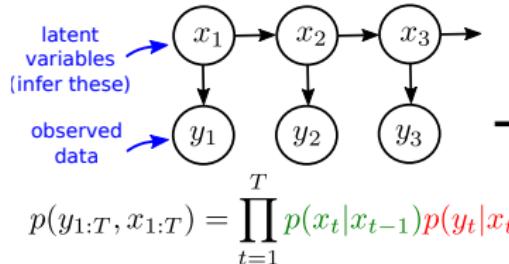
$$p(A|B, C) = \frac{1}{p(B|C)} p(B|A, C)p(A|C)$$

sum rule

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

Inference: Derivation of General Filtering Equations

Model



Rules of probability

product rule

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

sum rule

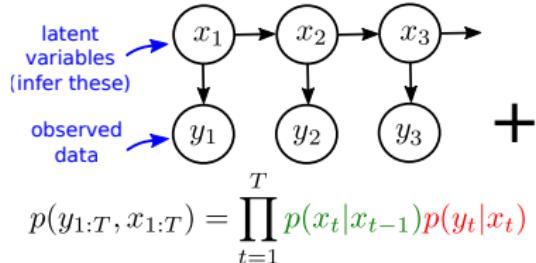
$$p(A|C) = \sum_B p(A, B|C)$$

Inference

= ?

Inference: Derivation of General Filtering Equations

Model



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

Rules of probability

product rule

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

sum rule

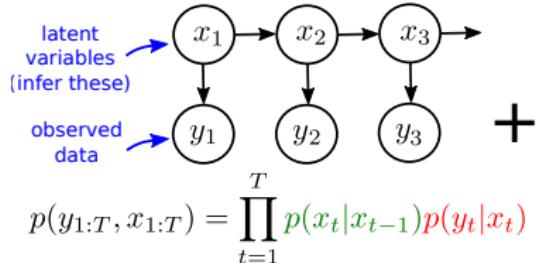
$$p(A|C) = \sum_B p(A, B|C)$$

Inference

= ?

Inference: Derivation of General Filtering Equations

Model



Rules of probability

product rule

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

sum rule

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

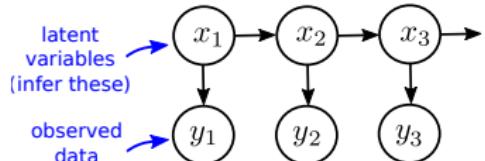
Inference

= ?

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

Inference: Derivation of General Filtering Equations

Model



Rules of probability

product rule

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

Inference

= ?

sum rule

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

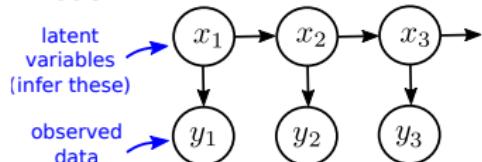
$$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t|x_{t-1})p(y_t|x_t)$$

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

$$= \frac{1}{p(y_t|y_{1:t-1})} p(y_t|x_t, y_{1:t-1})p(x_t|y_{1:t-1}) \quad \begin{matrix} \text{product rule} \\ A = x_t \ B = y_t \ C = y_{1:t-1} \end{matrix}$$

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

= ?

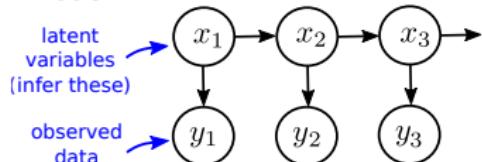
$$p(x_t|y_{1:t}) = p(x_t|y_t, y_{1:t-1})$$

$$= \frac{1}{p(y_t|y_{1:t-1})} p(y_t|x_t, y_{1:t-1})p(x_t|y_{1:t-1}) \quad \begin{matrix} \text{product rule} \\ A = x_t \ B = y_t \ C = y_{1:t-1} \end{matrix}$$

$$= \frac{1}{p(y_t|y_{1:t-1})} p(y_t|x_t)p(x_t|y_{1:t-1}) \quad \begin{matrix} \text{conditional independence from model} \\ y_t \perp y_{1:t-1}|x_t \end{matrix}$$

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

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

Inference

= ?

sum rule

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

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

$$= \frac{1}{p(y_t|y_{1:t-1})} p(y_t|x_t, y_{1:t-1})p(x_t|y_{1:t-1}) \quad \begin{matrix} \text{product rule} \\ A = x_t \ B = y_t \ C = y_{1:t-1} \end{matrix}$$

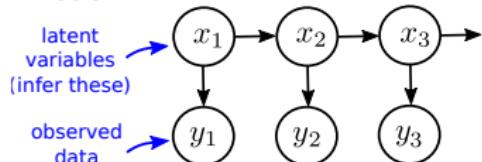
$$= \frac{1}{p(y_t|y_{1:t-1})} p(y_t|x_t)p(x_t|y_{1:t-1}) \quad \begin{matrix} \text{conditional independence from model} \\ y_t \perp y_{1:t-1}|x_t \end{matrix}$$

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

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

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

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

Inference

= ?

sum rule

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

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

$$= \frac{1}{p(y_t|y_{1:t-1})} p(y_t|x_t, y_{1:t-1})p(x_t|y_{1:t-1}) \quad \begin{matrix} \text{product rule} \\ A = x_t \ B = y_t \ C = y_{1:t-1} \end{matrix}$$

$$= \frac{1}{p(y_t|y_{1:t-1})} p(y_t|x_t)p(x_t|y_{1:t-1}) \quad \begin{matrix} \text{conditional independence from model} \\ y_t \perp y_{1:t-1}|x_t \end{matrix}$$

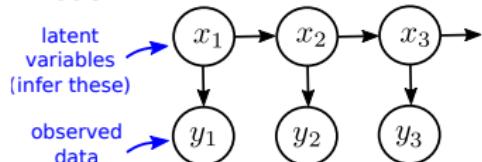
$$\propto p(y_t|x_t)p(x_t|y_{1:t-1})$$

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

= ?

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

$$= \frac{1}{p(y_t|y_{1:t-1})} p(y_t|x_t, y_{1:t-1})p(x_t|y_{1:t-1}) \quad \begin{matrix} \text{product rule} \\ A = x_t \ B = y_t \ C = y_{1:t-1} \end{matrix}$$

$$= \frac{1}{p(y_t|y_{1:t-1})} p(y_t|x_t)p(x_t|y_{1:t-1}) \quad \begin{matrix} \text{conditional independence from model} \\ y_t \perp y_{1:t-1}|x_t \end{matrix}$$

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

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

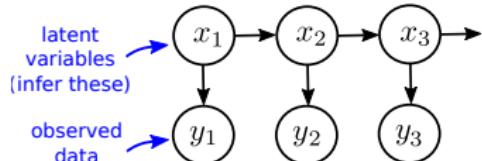
$$p(x_t|y_{1:t-1}) = \int p(x_t, x_{t-1}|y_{1:t-1})dx_{t-1}$$

sum rule

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

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

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

Inference

= ?

sum rule

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

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

$$= \frac{1}{p(y_t|y_{1:t-1})} p(y_t|x_t, y_{1:t-1})p(x_t|y_{1:t-1}) \quad \begin{matrix} \text{product rule} \\ A = x_t \ B = y_t \ C = y_{1:t-1} \end{matrix}$$

$$= \frac{1}{p(y_t|y_{1:t-1})} p(y_t|x_t)p(x_t|y_{1:t-1}) \quad \begin{matrix} \text{conditional independence from model} \\ y_t \perp y_{1:t-1}|x_t \end{matrix}$$

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

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

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

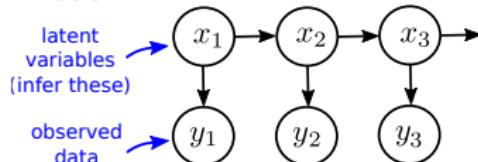
sum rule

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

$$= \int p(x_t|x_{t-1}, y_{1:t-1})p(x_{t-1}|y_{1:t-1})dx_{t-1} \quad \begin{matrix} \text{product rule} \\ \end{matrix}$$

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

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

Inference

= ?

sum rule

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

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

product rule

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

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

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

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

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

$$\propto \underbrace{p(y_t|x_t)}_{\text{constant of proportionality}} \underbrace{p(x_t|y_{1:t-1})}_{p(y_t|y_{1:t-1}) \text{ (see learning)}}$$

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

sum rule

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

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

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

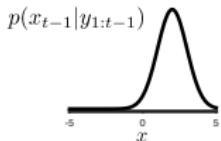
conditional independence from model

Inference: Kalman Filter

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

diffuse via
dynamics

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

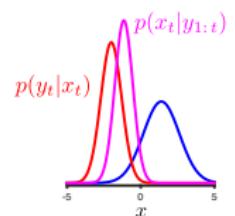
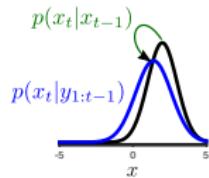


combine
with
likelihood

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

prior

likelihood



Inference: Kalman Filter

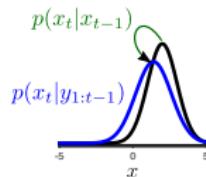
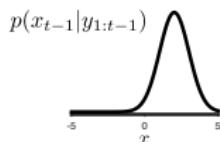
$$p(x_{t-1}|y_{1:t-1}) = \mathcal{G}(x_{t-1}; \mu_{t-1}^{t-1}, V_{t-1}^{t-1})$$

most recent data used
in prediction

variable being predicted

diffuse via dynamics

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

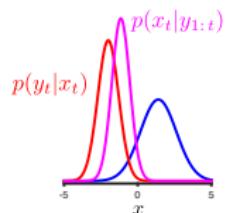


combine with likelihood

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

prior 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

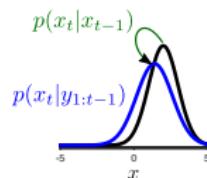
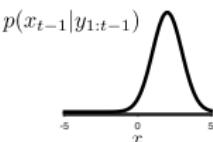
variable being predicted

diffuse via dynamics

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

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

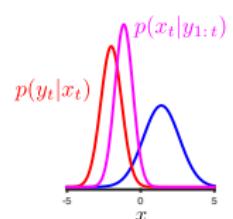
↖ ↘



combine
with
likelihood

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

prior likelihood



Inference: Kalman Filter

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

- 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

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

- $\mu_t^{t-1} = A\mu_{t-1}^{t-1}$
- $V_t^{t-1} = AV_{t-1}^{t-1}A^\top + Q$
- variance inflates

combine with likelihood

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

prior likelihood

$$p(x_t|x_{t-1}) = \mathcal{G}(x_t; \underbrace{\mu_{t-1}}_A, \underbrace{V_{t-1}}_Q)$$

$$E(x_t) = E(a x_t + \underbrace{e_t}_{\sim N(0, Q)})$$

$$p(x_{t-1}) = \mathcal{G}(x_{t-1}; \underbrace{\mu_{t-1}^{t-1}}_A, \underbrace{V_{t-1}^{t-1}}_Q)$$

$$E(x_t^2) = E[(a x_t + e_t)^2]$$

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

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

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

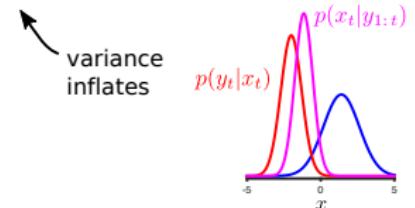
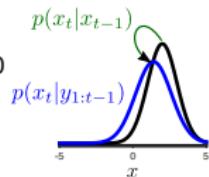
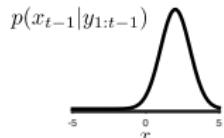
combine
with
likelihood

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

prior likelihood

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

~~prior~~



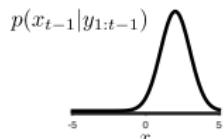
Inference: Kalman Filter

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

most recent data used
in prediction

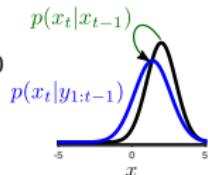
variable being predicted

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



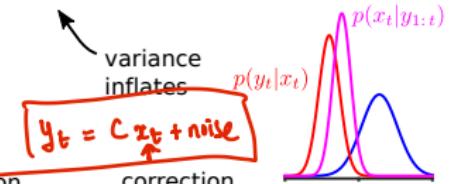
||* $p(x_t|y_{1:t-1}) = \mathcal{G}(x_t; \mu_t^{t-1}, V_t^{t-1}) \quad \mu_t^{t-1} = A\mu_{t-1}^{t-1}$

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

combine
with
likelihood

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

prior likelihood



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

Kalman gain

$$\begin{aligned} \mu_t^t &= \underline{\mu}_{t-1}^{t-1} + K_t(\underline{y}_t - C\underline{\mu}_{t-1}^{t-1}) \\ V_t^t &= \underline{V}_{t-1}^{t-1} - K_t C \underline{V}_{t-1}^{t-1} C^\top + R \\ K_t &= V_t^{t-1} C^\top (C V_t^{t-1} C^\top + R)^{-1} \end{aligned}$$

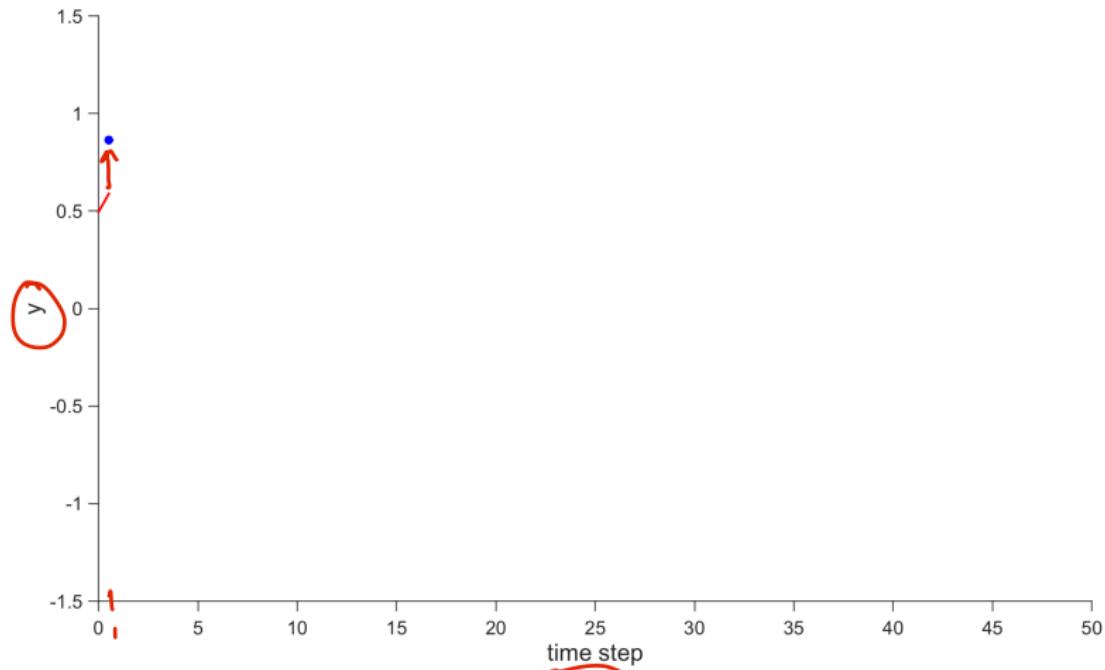
ε

Kalman Filter Demo

- ▶ data: $y_t = \underline{\sin(\omega t)} + \sigma_y \epsilon_t$ 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

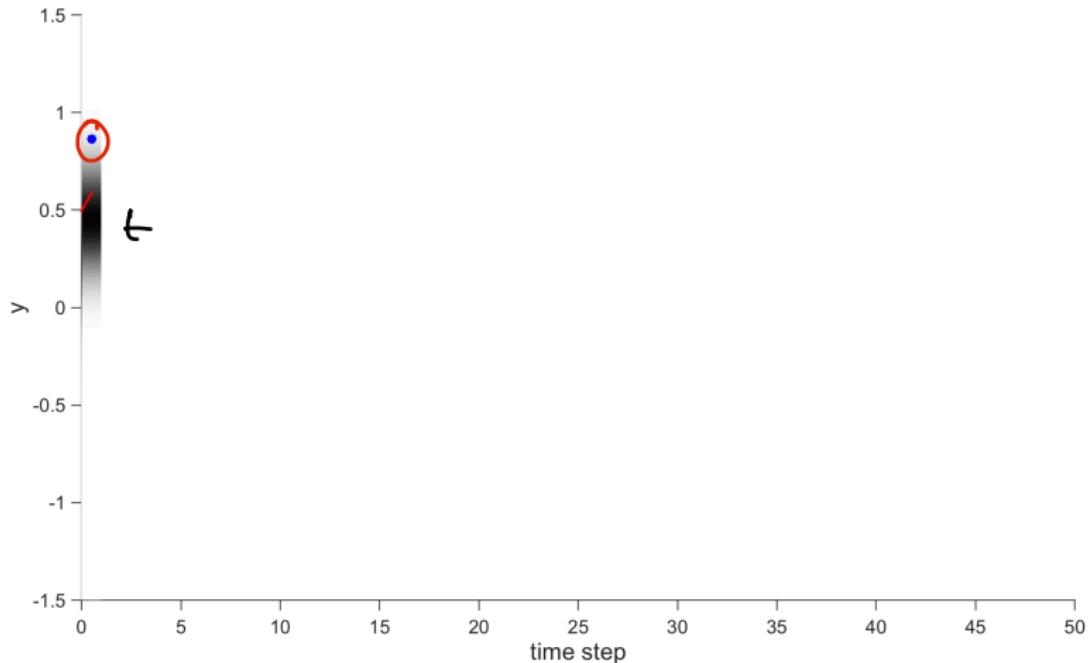
observed noisy data y_t , ground truth sinusoid



observe first data point y_1

Kalman Filter Demo

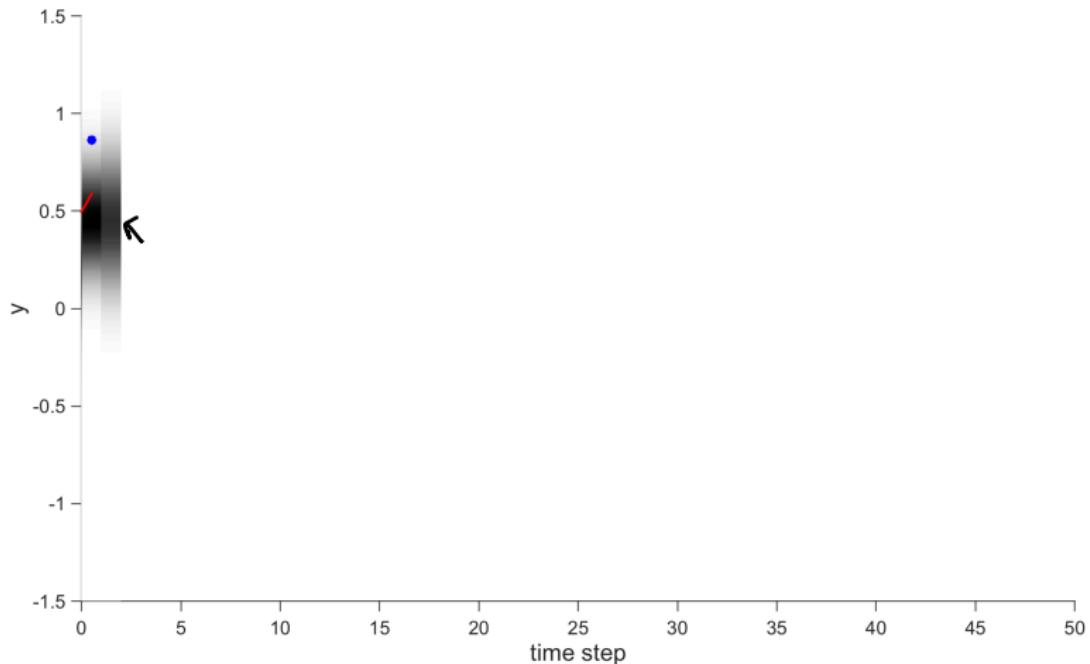
observed noisy data y_t , ground truth sinusoid



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

Kalman Filter Demo

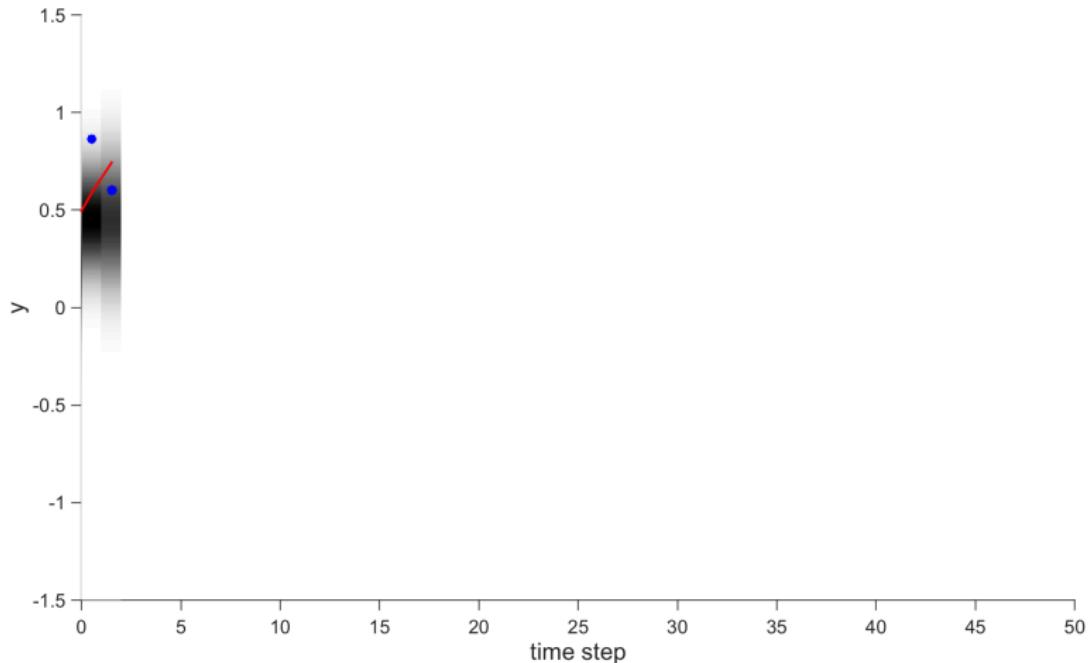
observed noisy data y_t , ground truth sinusoid



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

Kalman Filter Demo

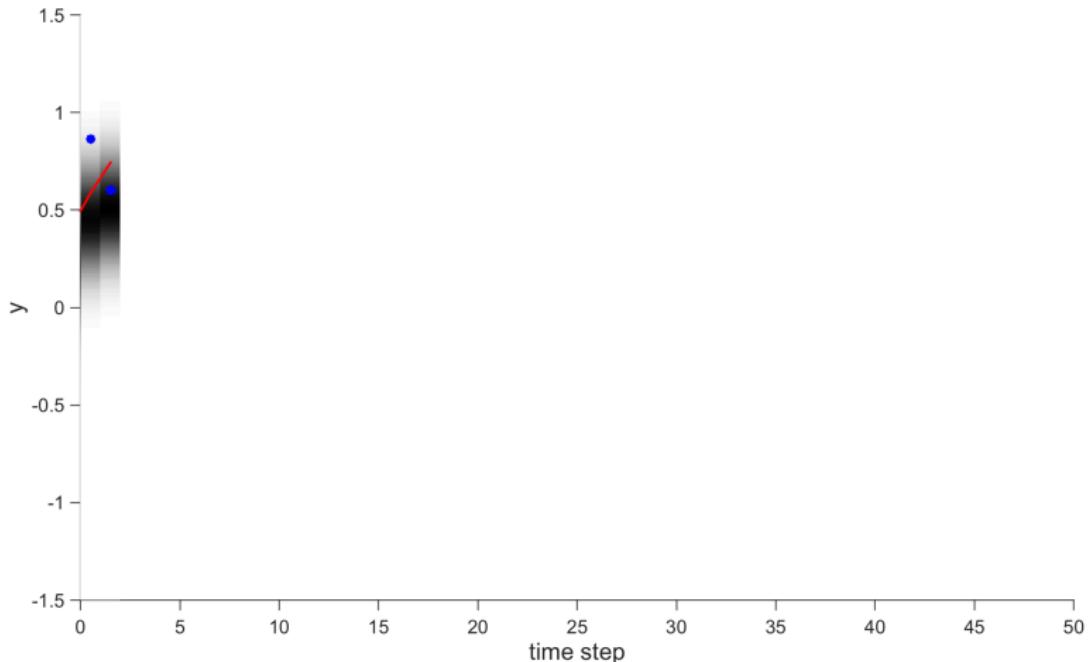
observed noisy data y_t , ground truth sinusoid



observe next data point y_2

Kalman Filter Demo

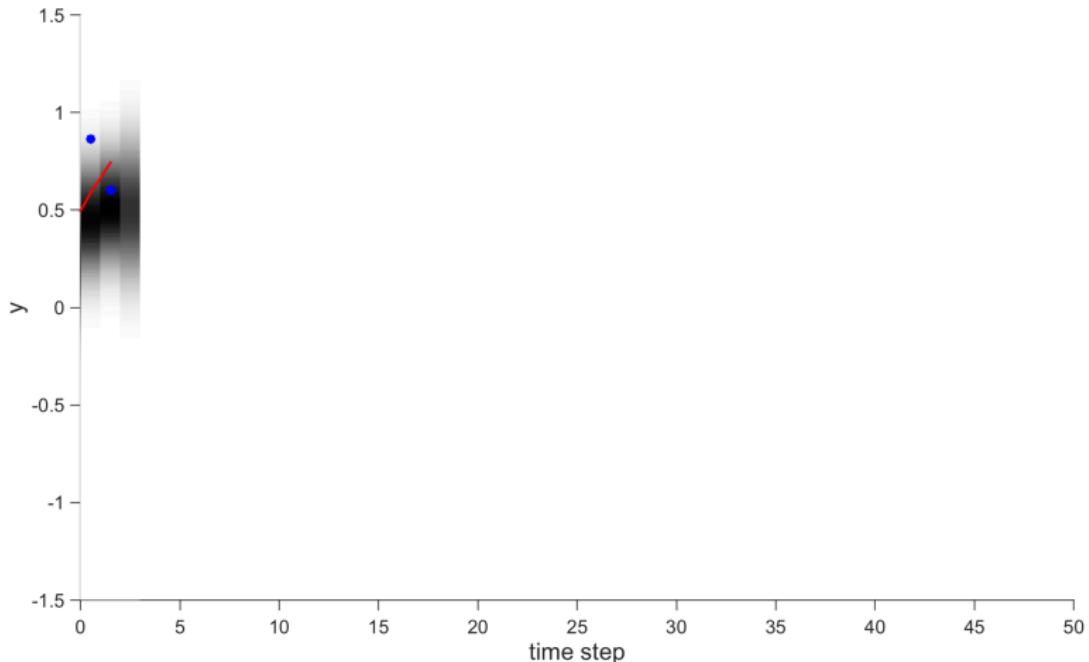
observed noisy data y_t , ground truth sinusoid



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

Kalman Filter Demo

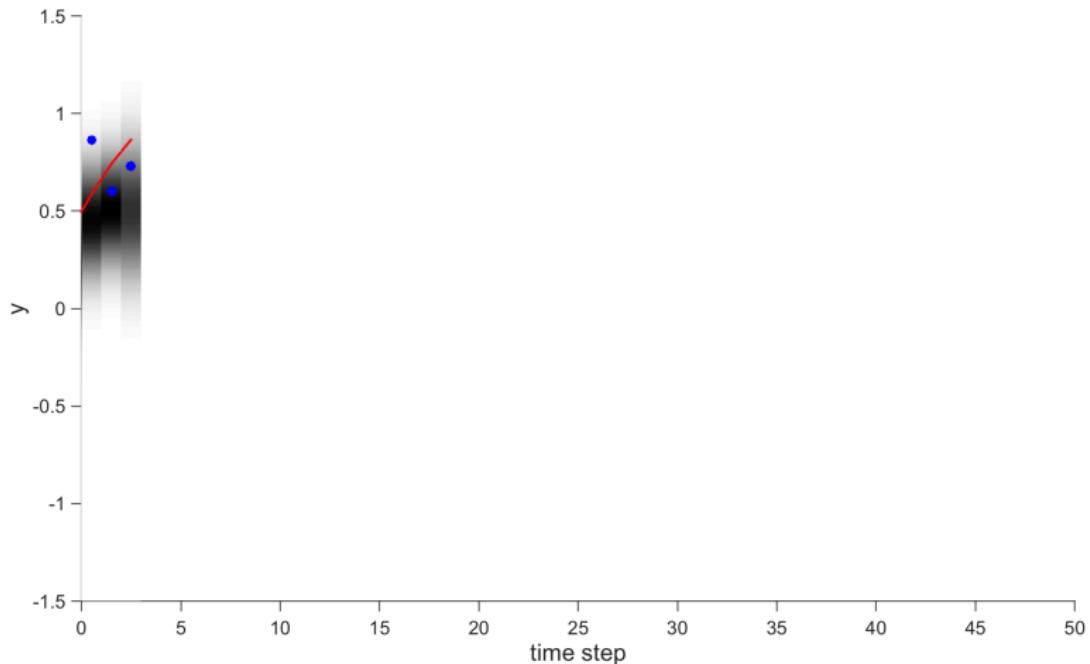
observed noisy data y_t , ground truth sinusoid



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

Kalman Filter Demo

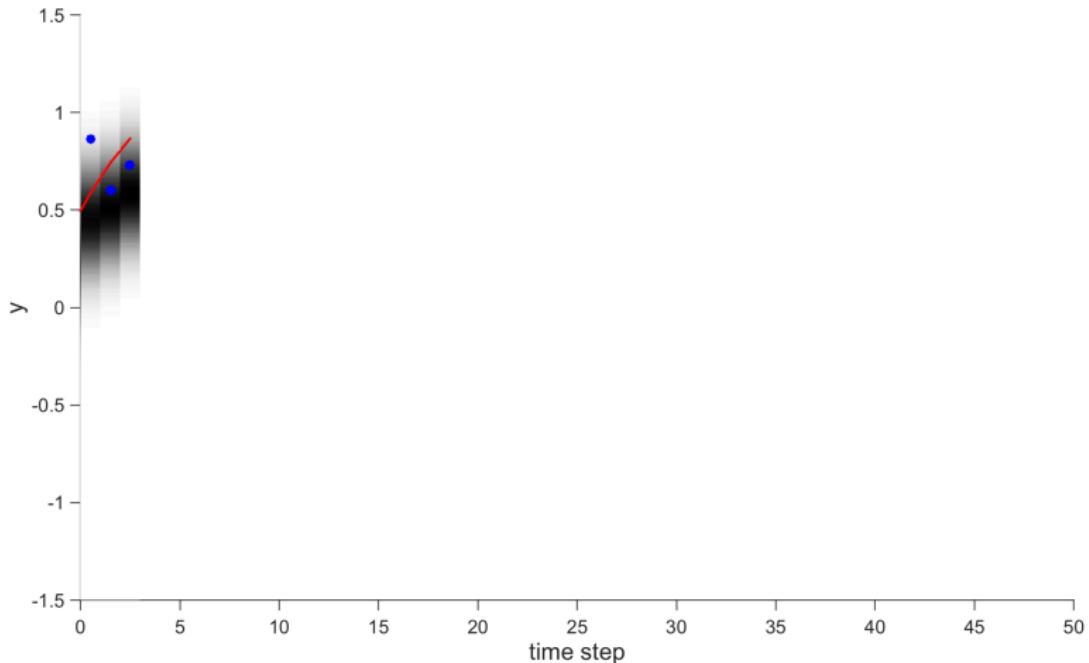
observed noisy data y_t , ground truth sinusoid



observe next data point y_3

Kalman Filter Demo

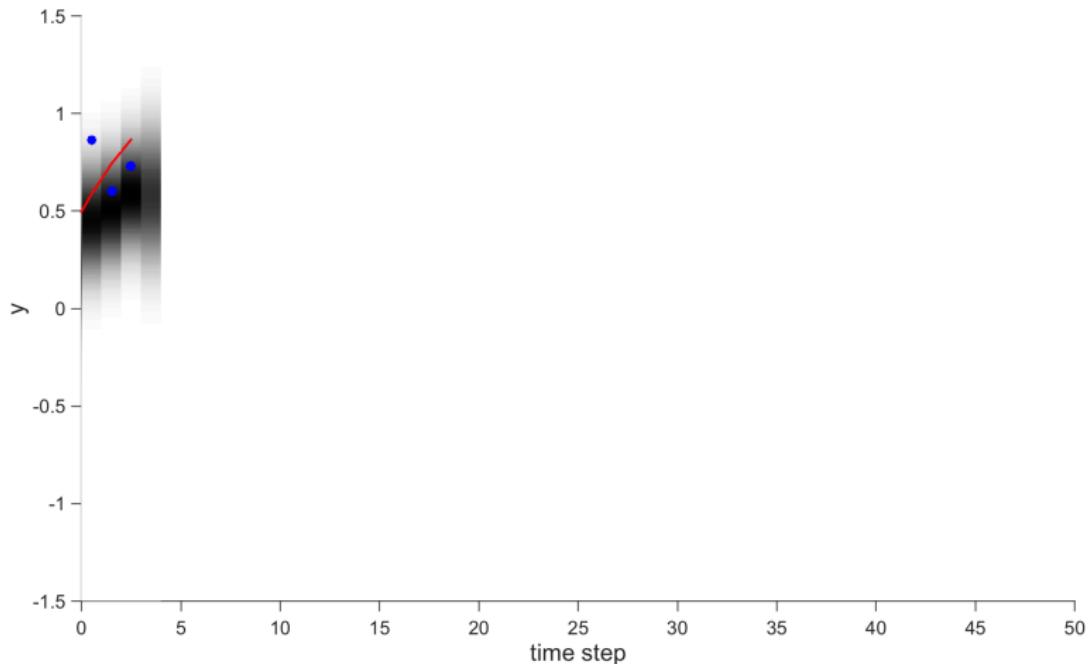
observed noisy data y_t , ground truth sinusoid



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

Kalman Filter Demo

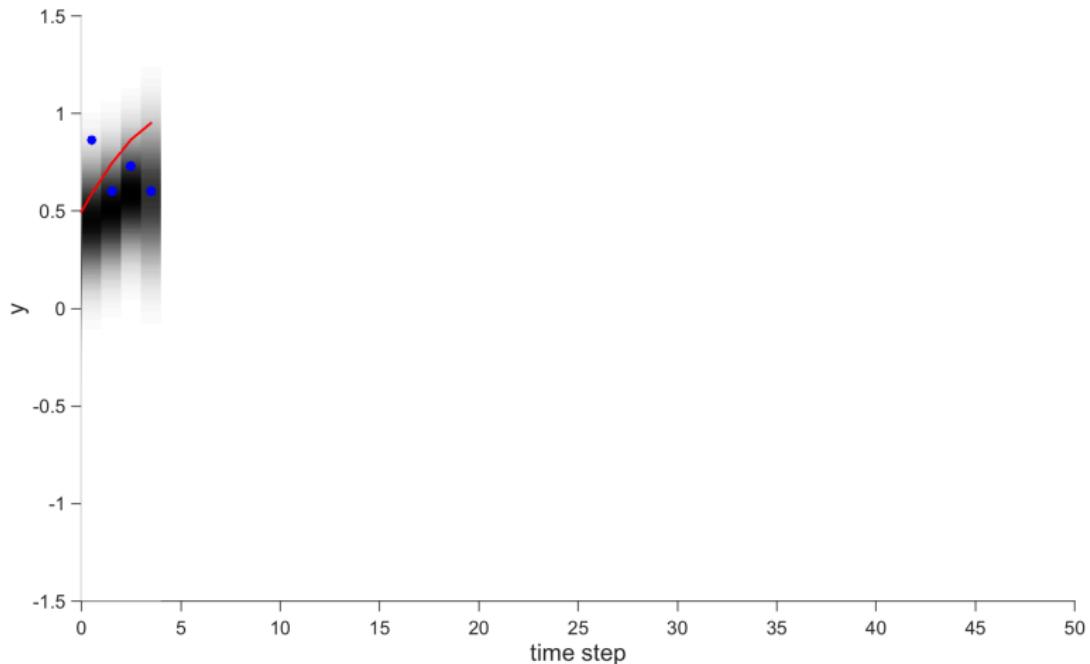
observed noisy data y_t , ground truth sinusoid



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

Kalman Filter Demo

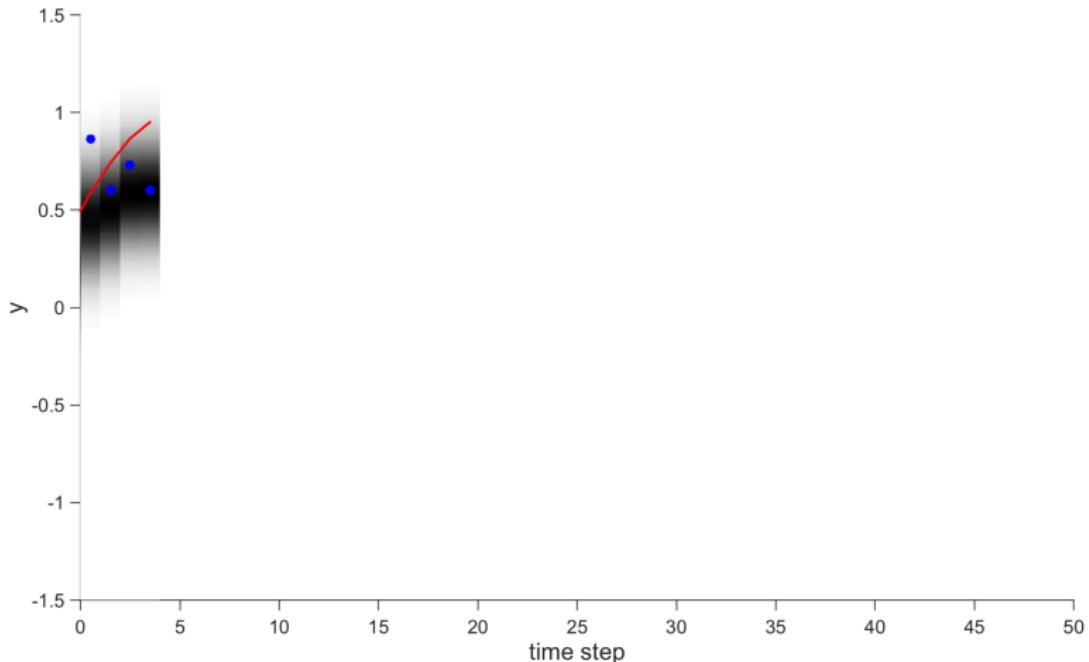
observed noisy data y_t , ground truth sinusoid



observe next data point y_4

Kalman Filter Demo

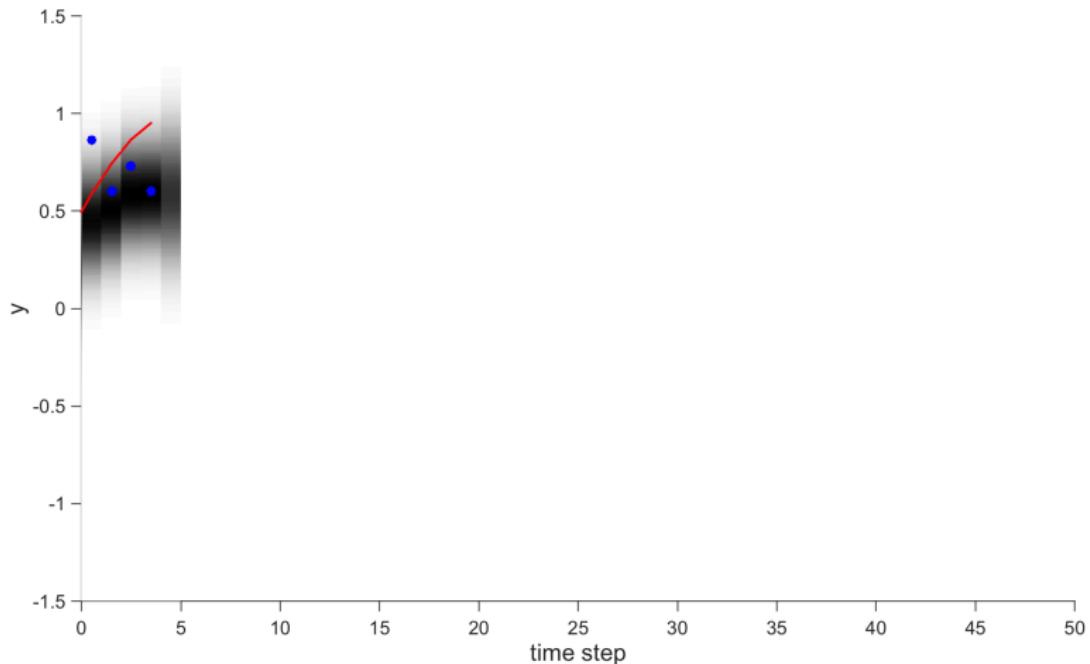
observed noisy data y_t , ground truth sinusoid



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

Kalman Filter Demo

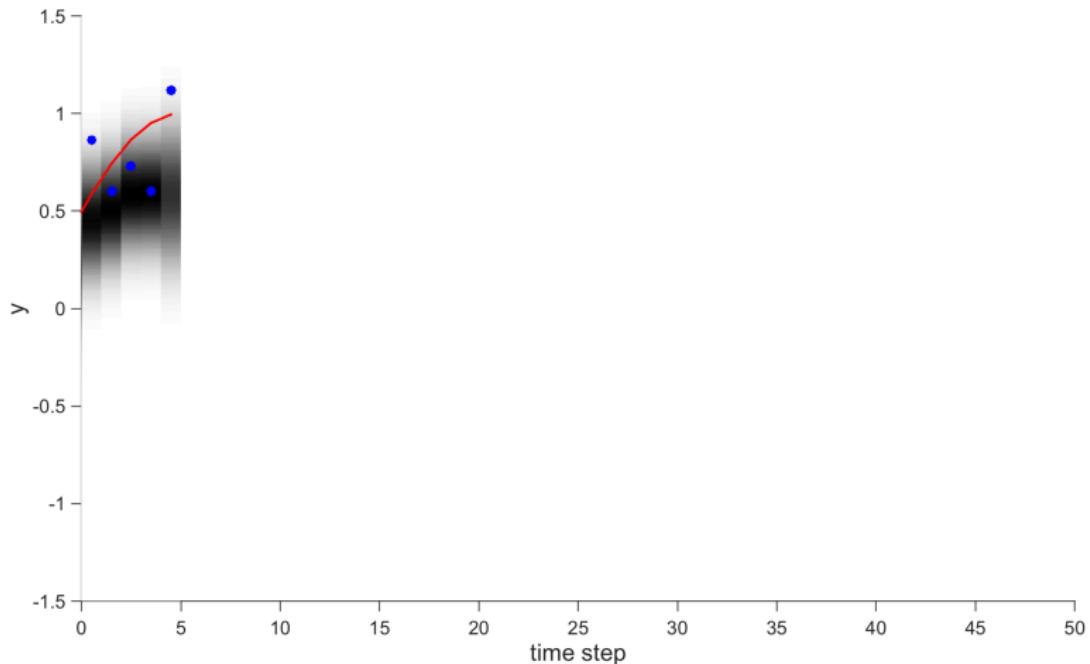
observed noisy data y_t , ground truth sinusoid



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

Kalman Filter Demo

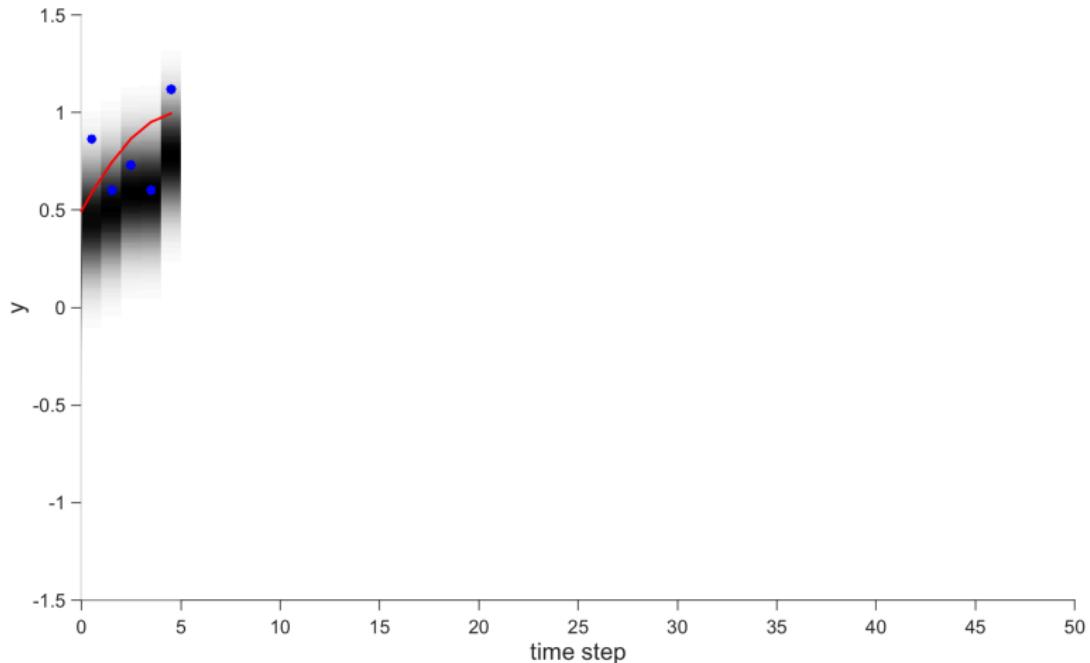
observed noisy data y_t , ground truth sinusoid



observe next data point y_5

Kalman Filter Demo

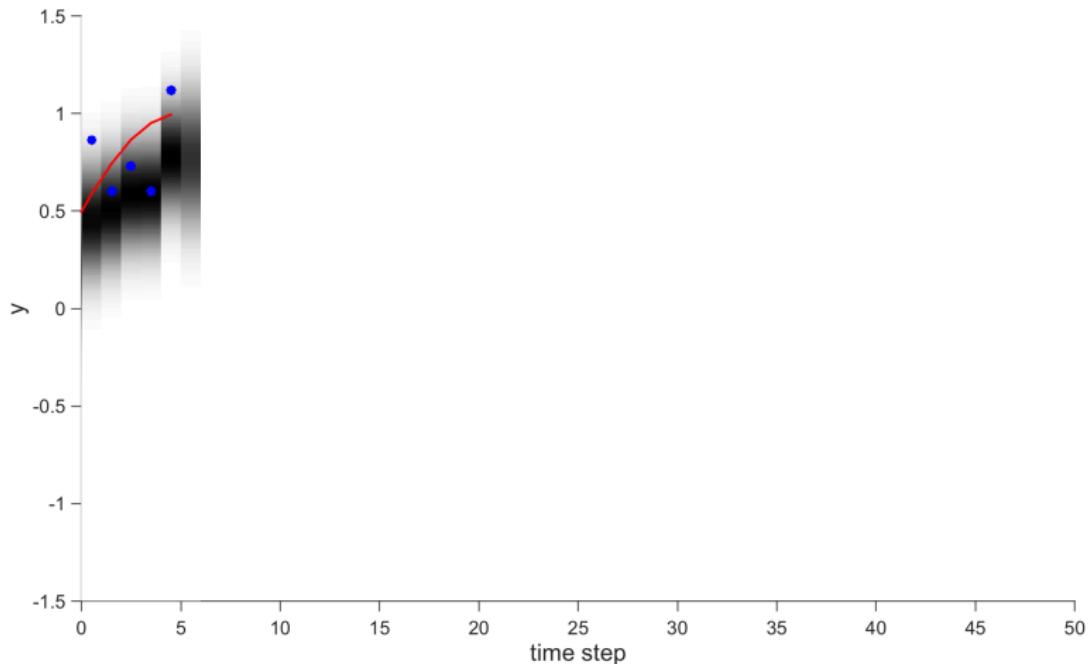
observed noisy data y_t , ground truth sinusoid



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

Kalman Filter Demo

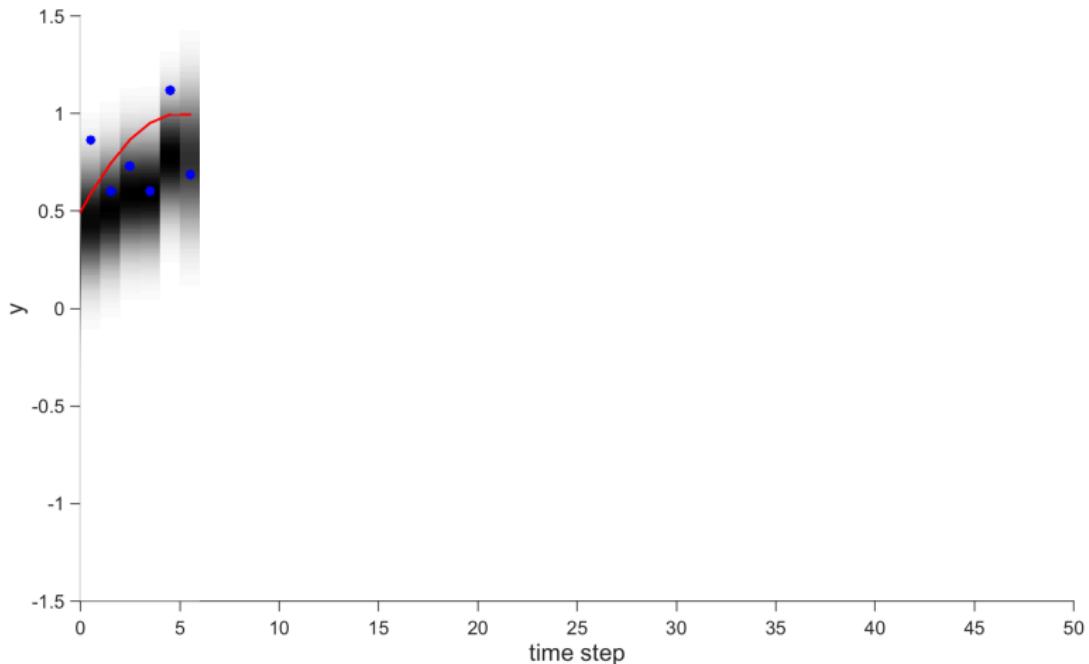
observed noisy data y_t , ground truth sinusoid



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

Kalman Filter Demo

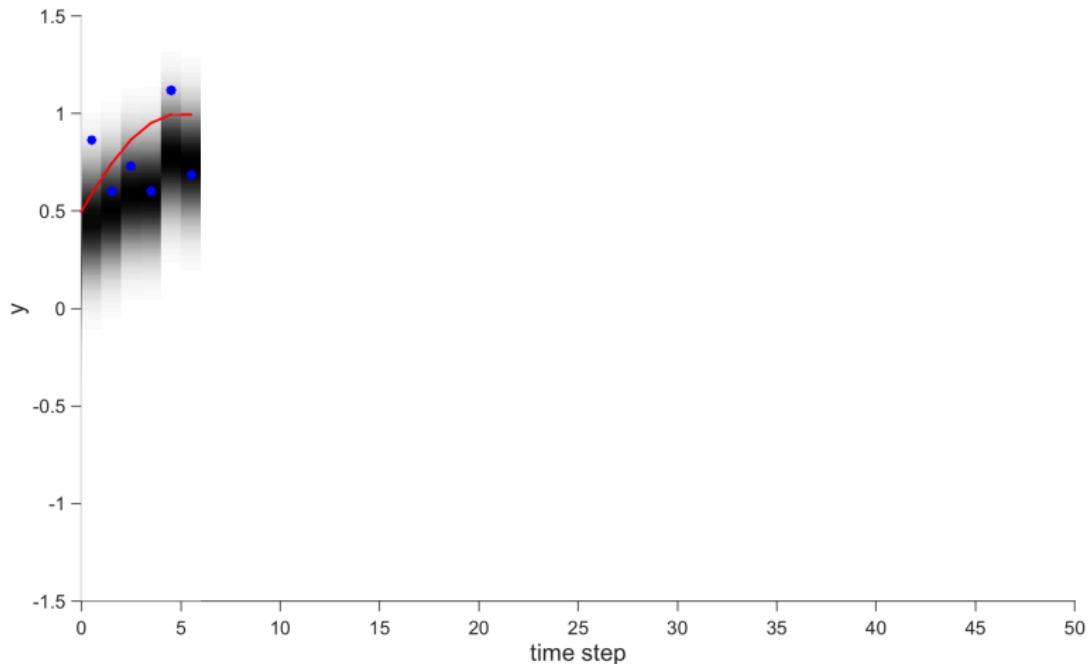
observed noisy data y_t , ground truth sinusoid



observe next data point y_6

Kalman Filter Demo

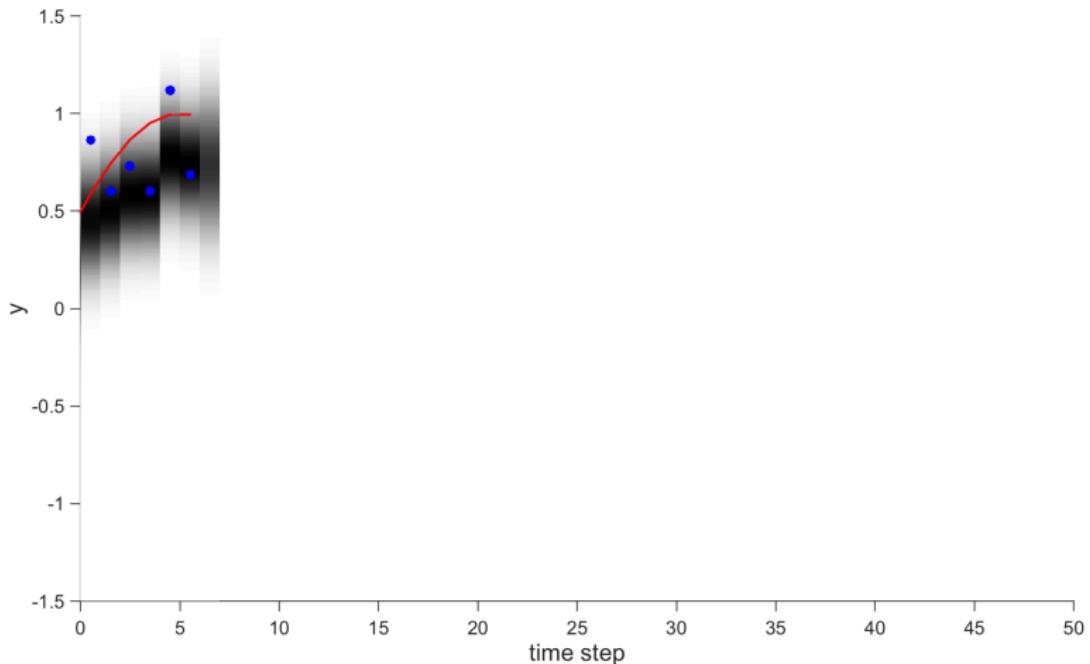
observed noisy data y_t , ground truth sinusoid



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

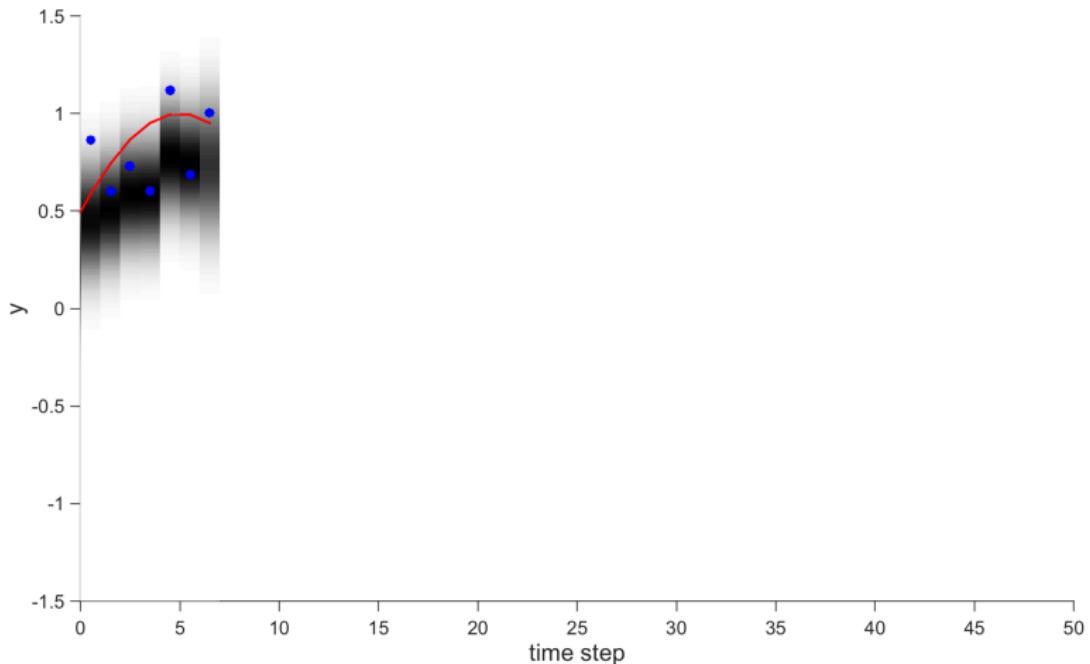
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



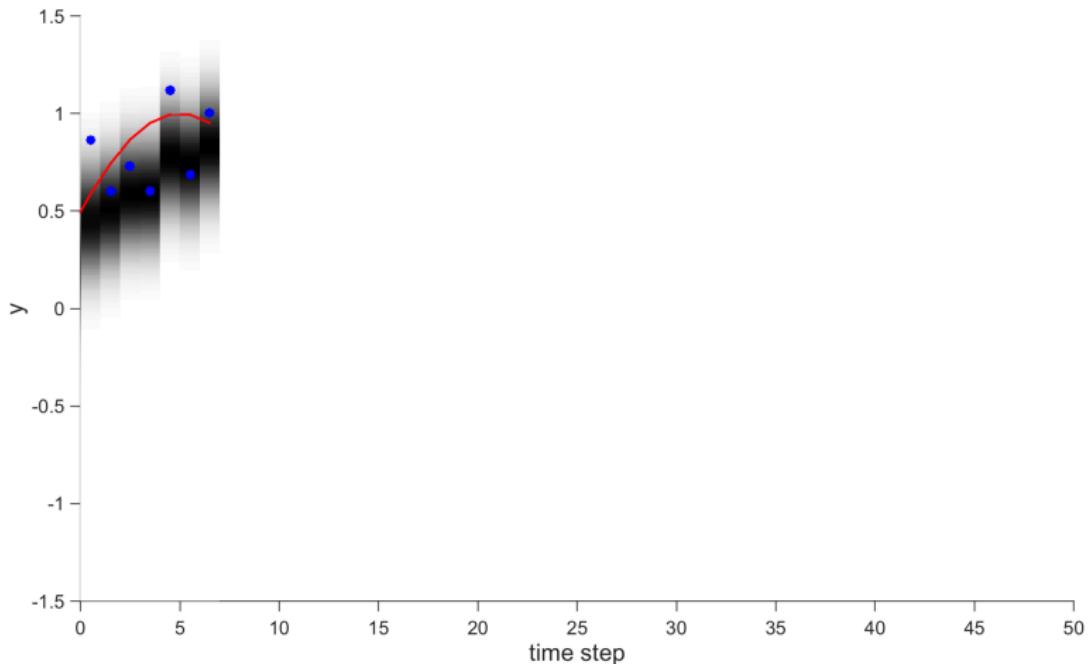
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



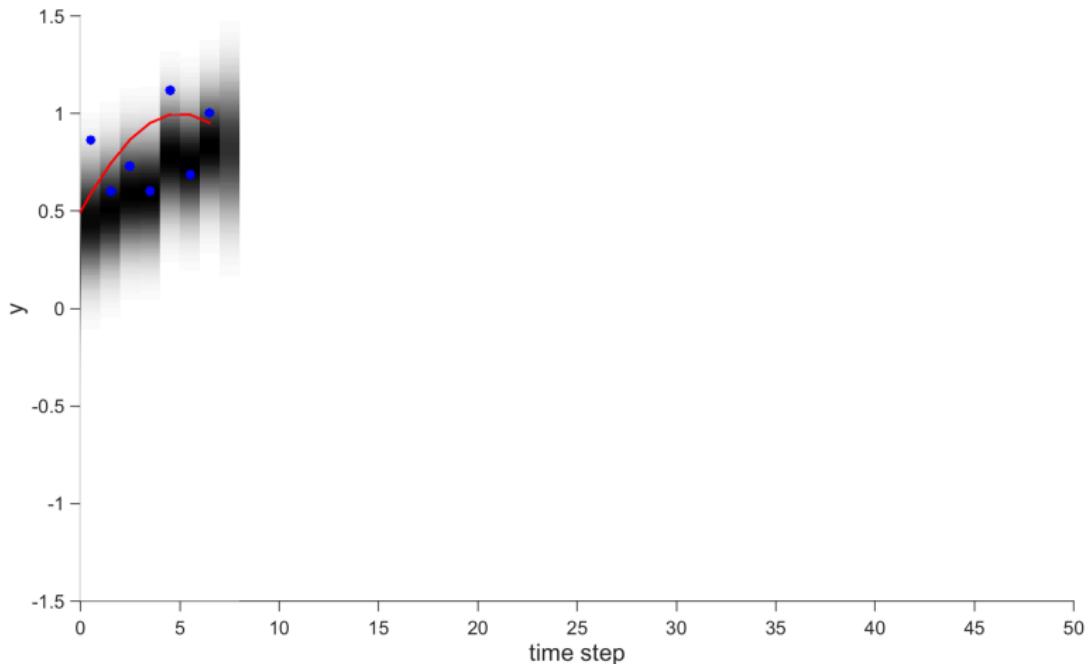
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



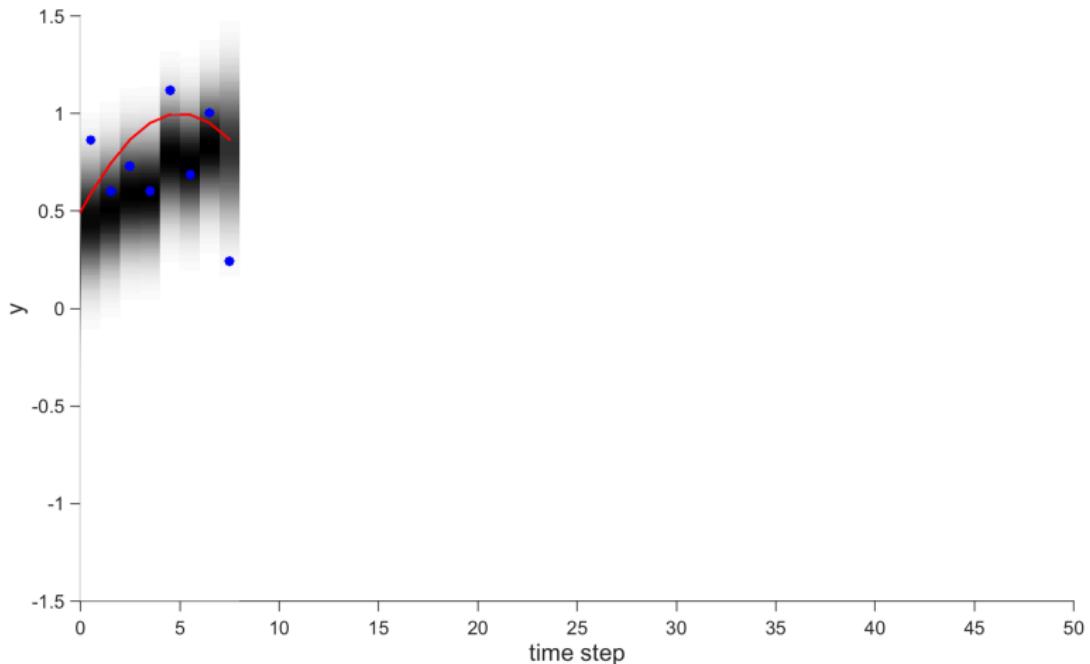
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



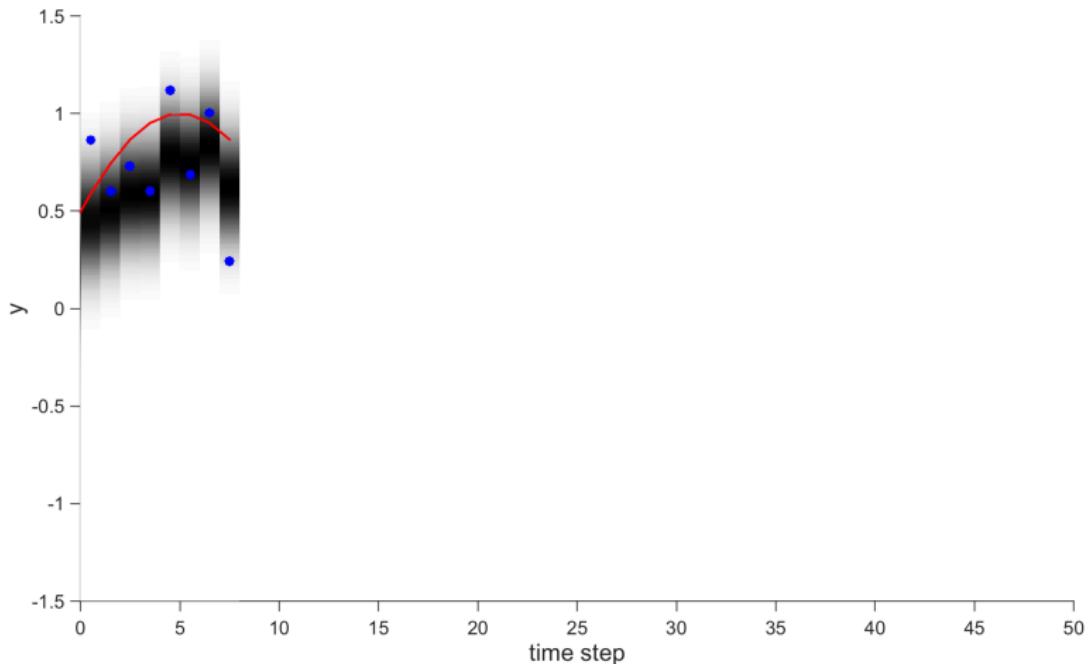
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



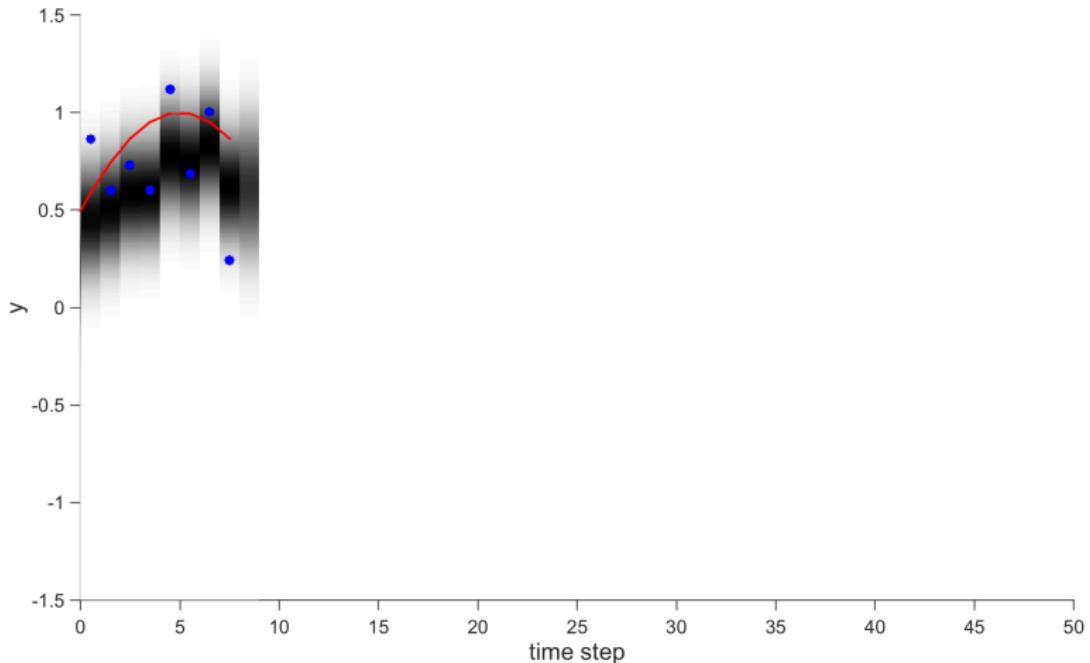
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



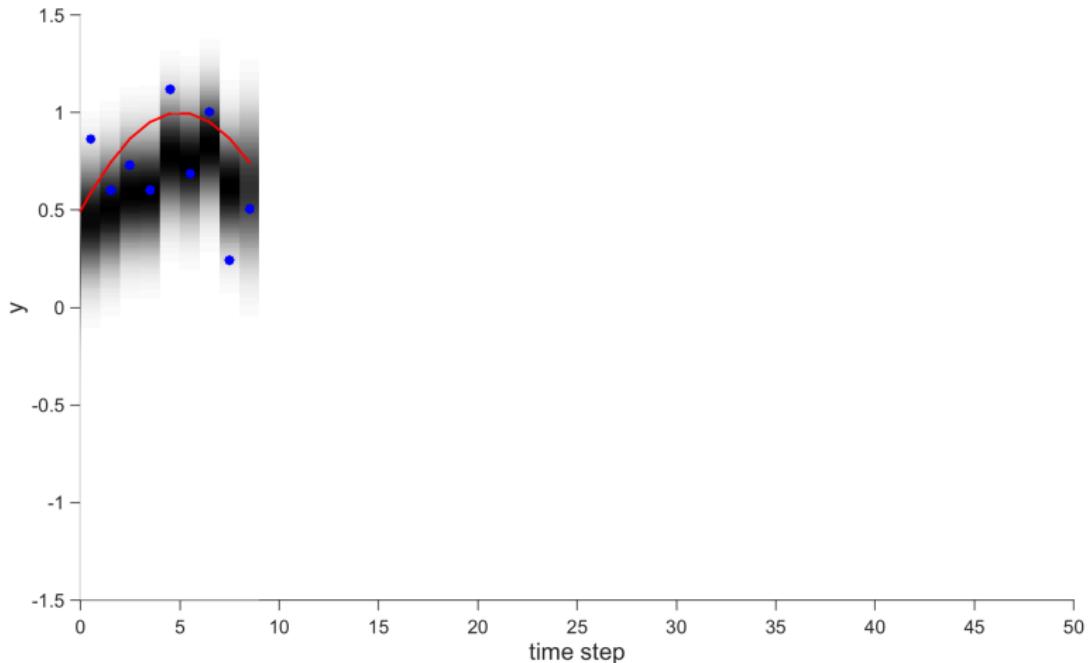
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



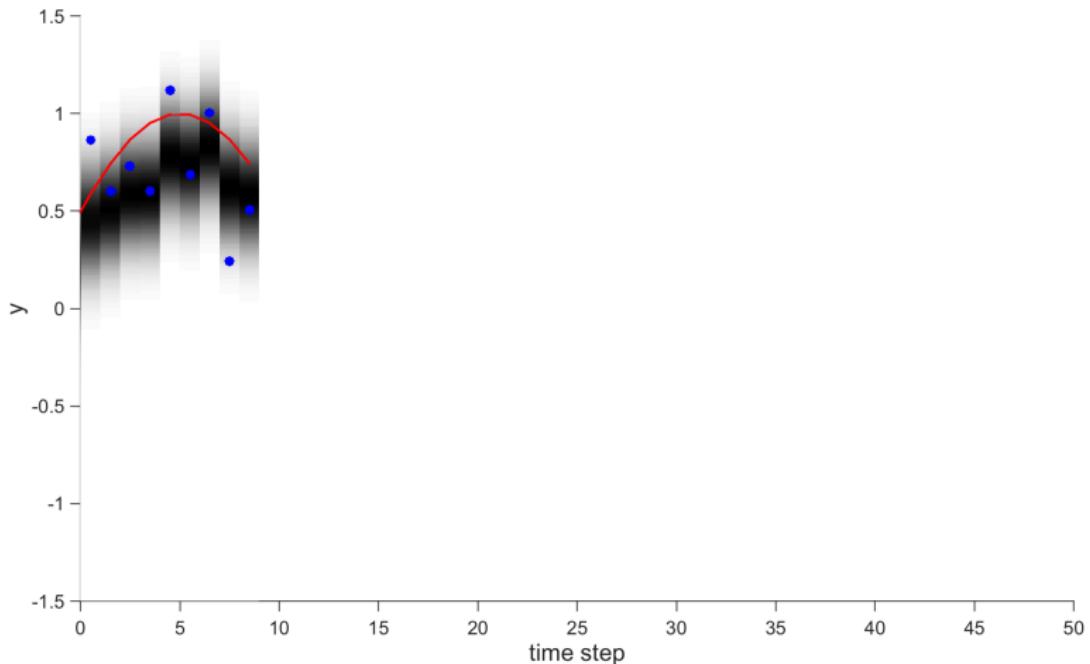
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



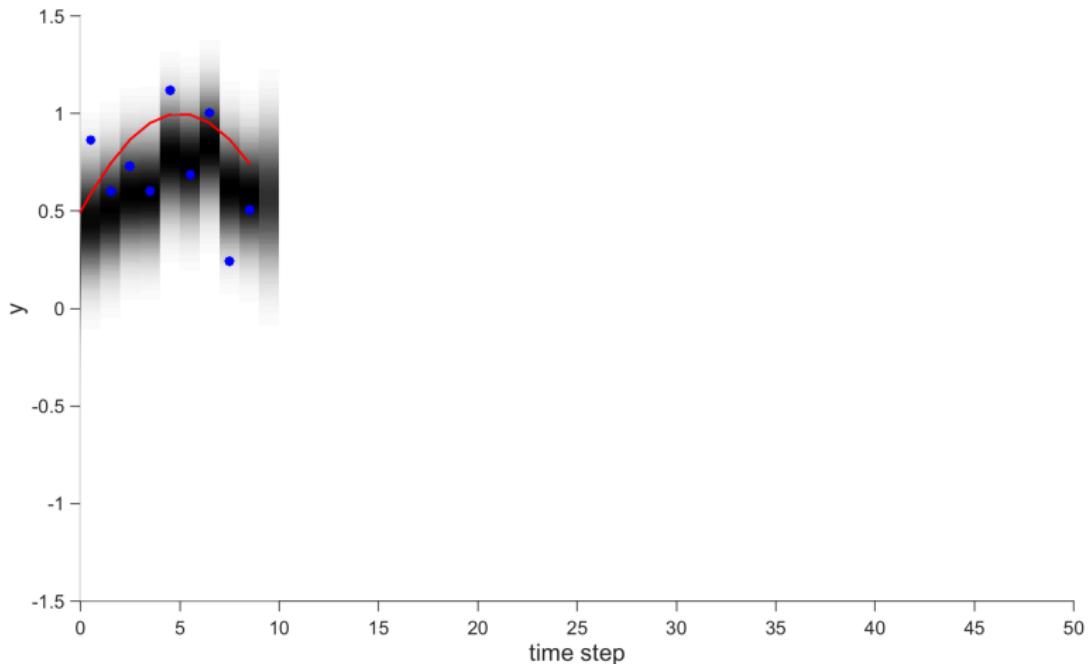
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



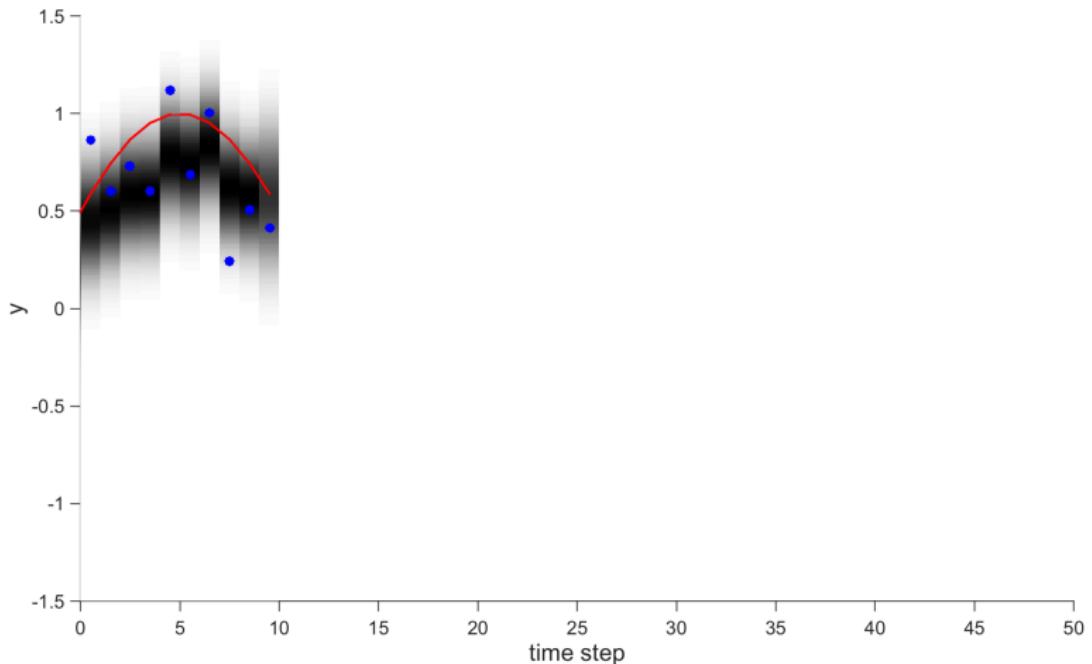
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



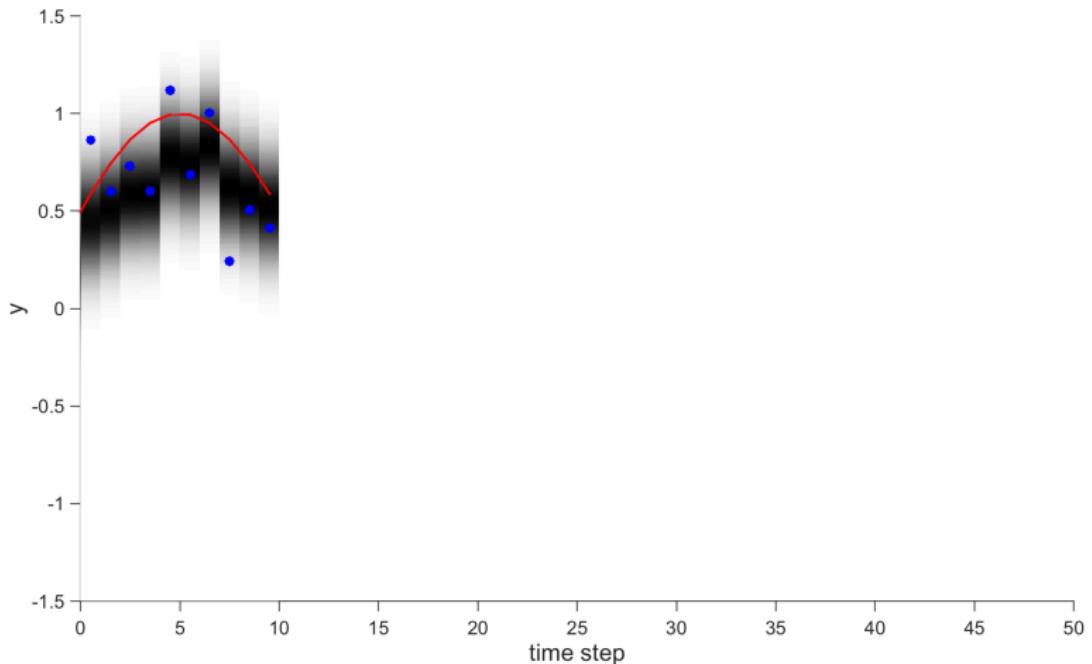
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



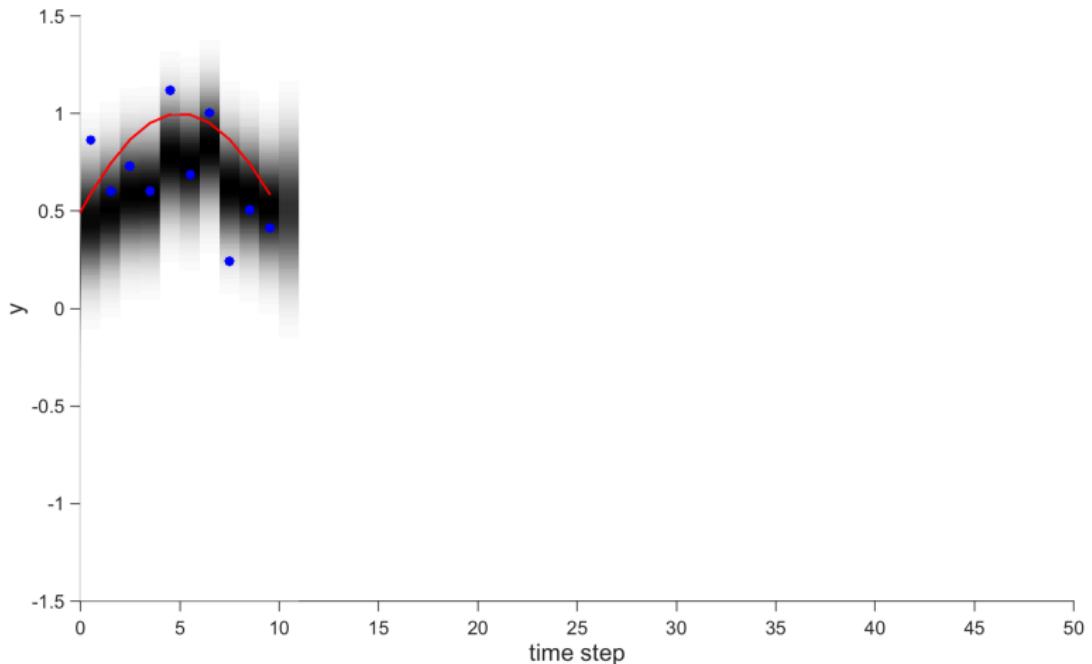
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



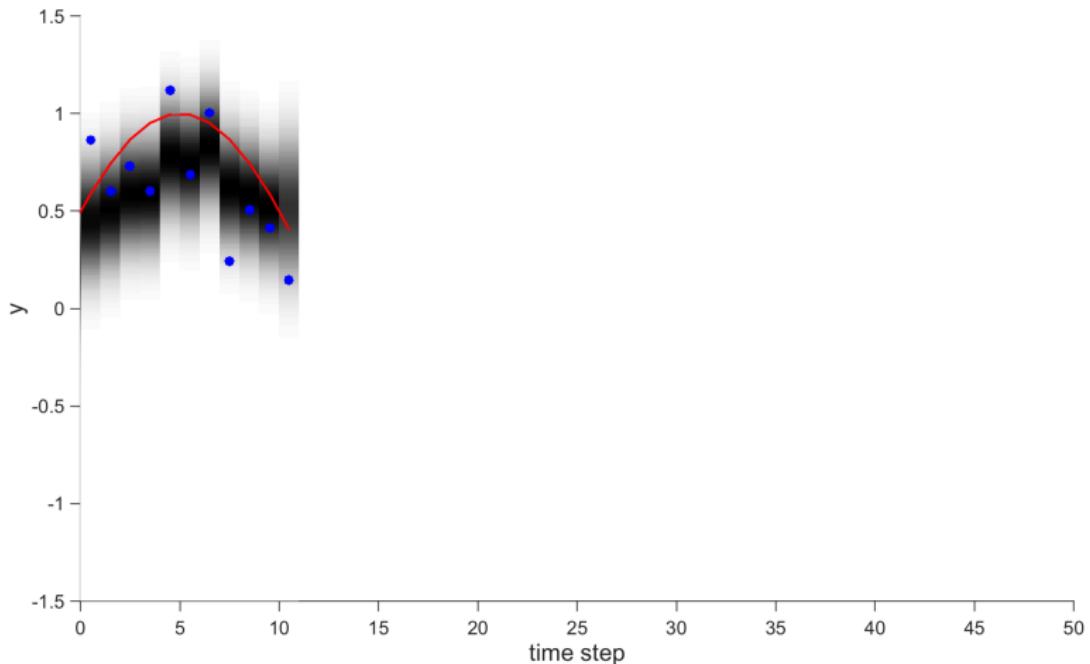
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



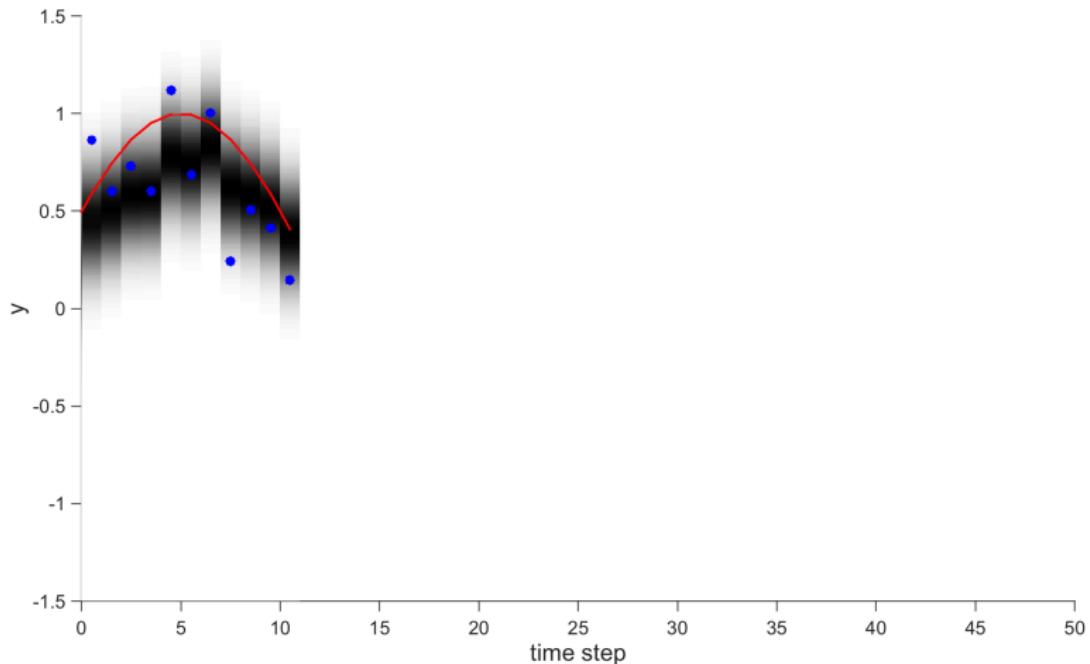
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



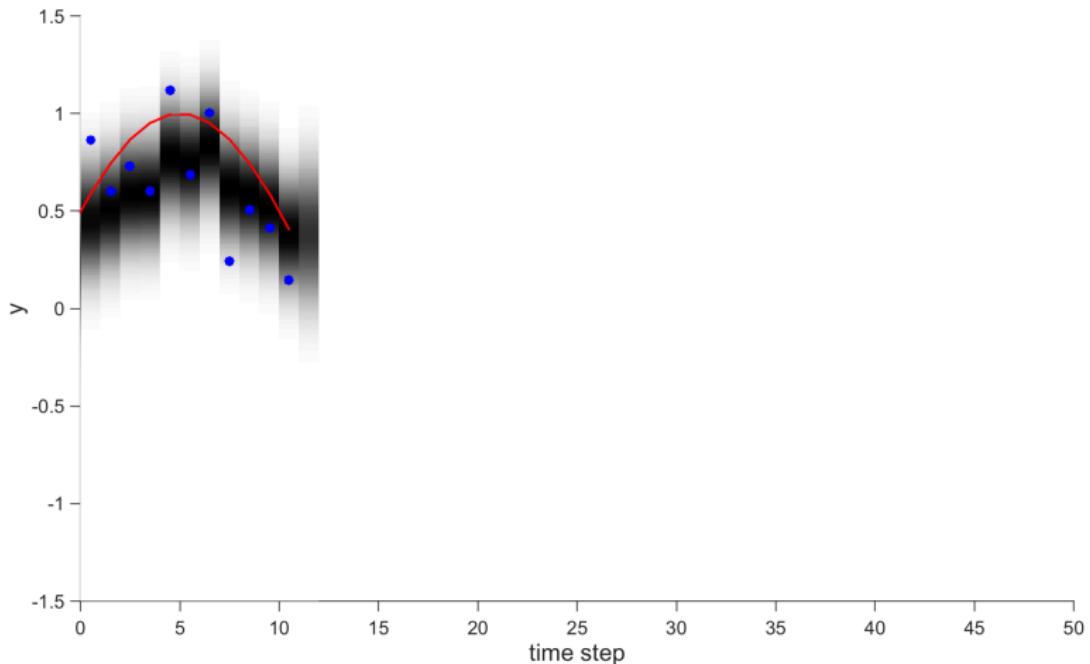
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



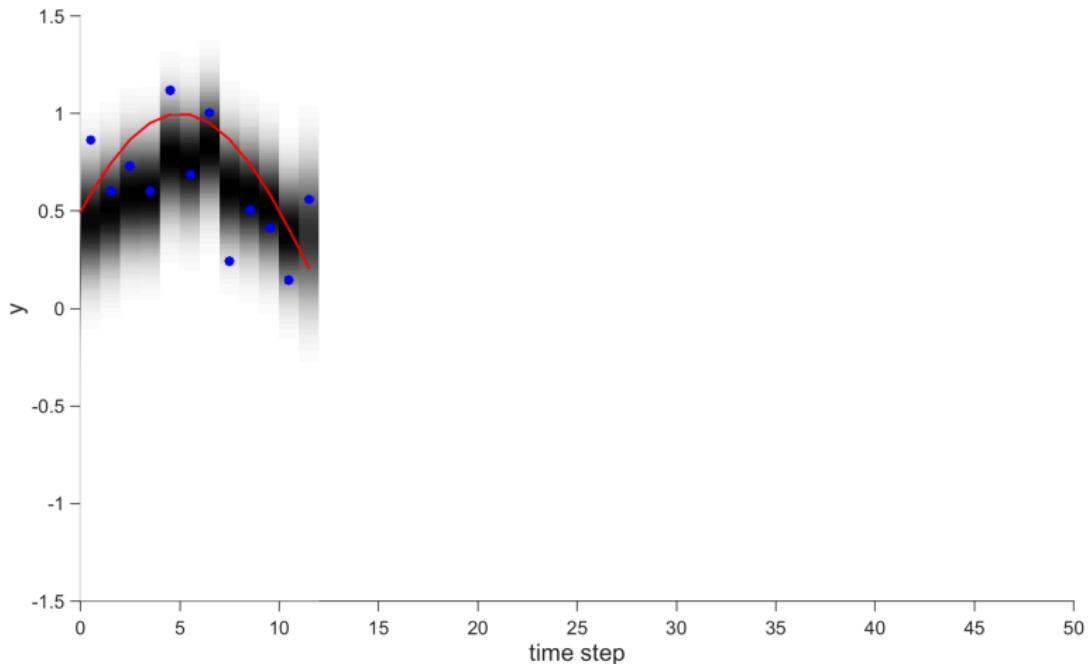
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



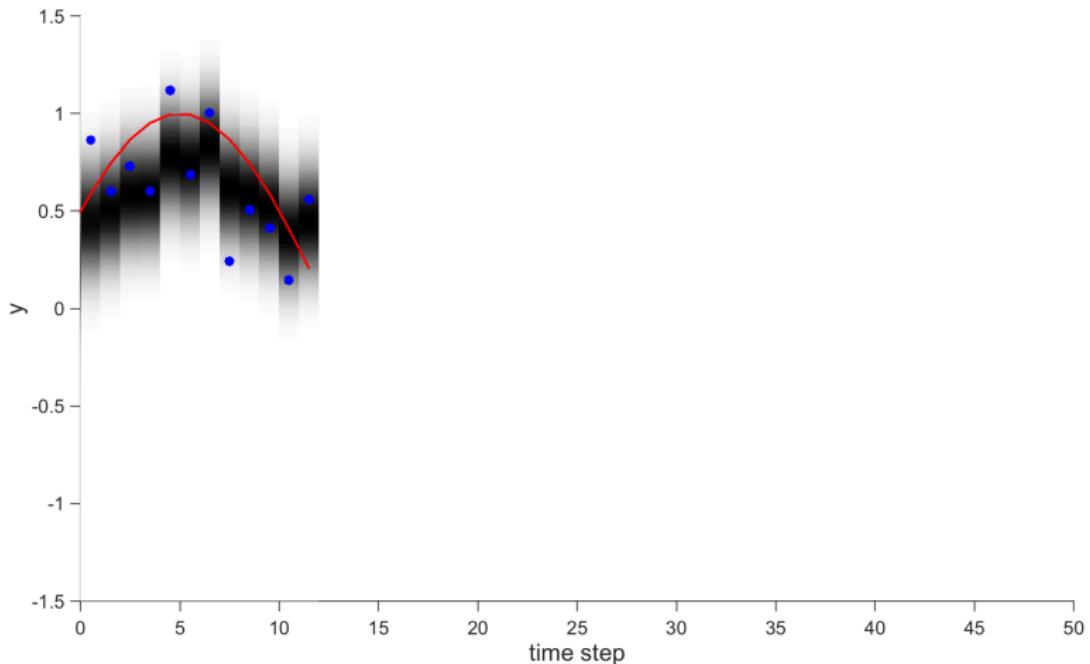
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



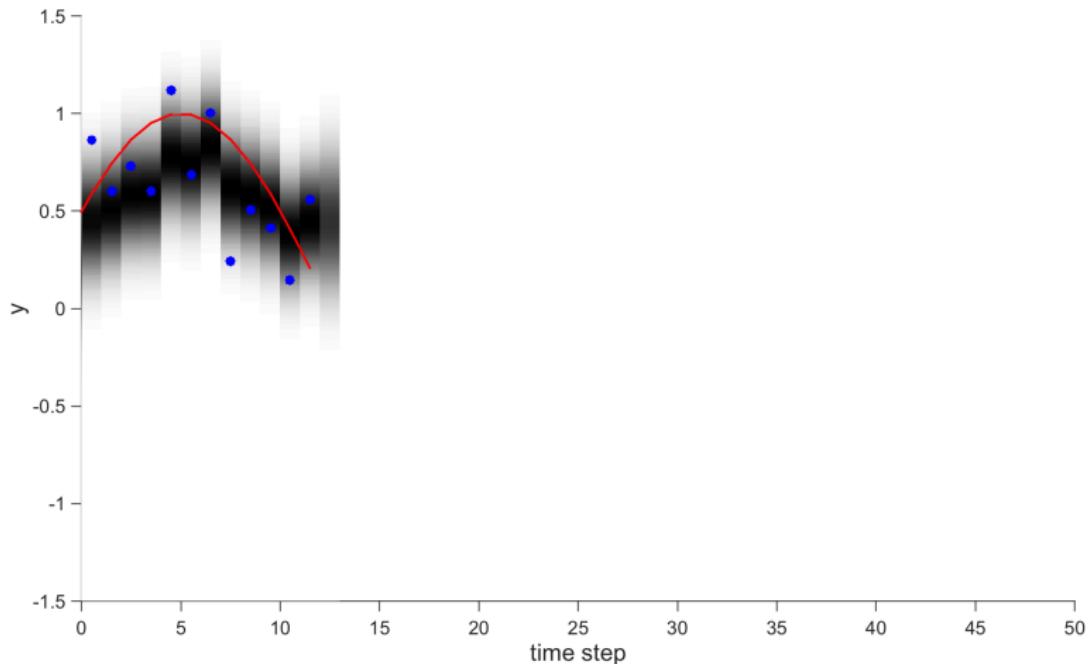
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



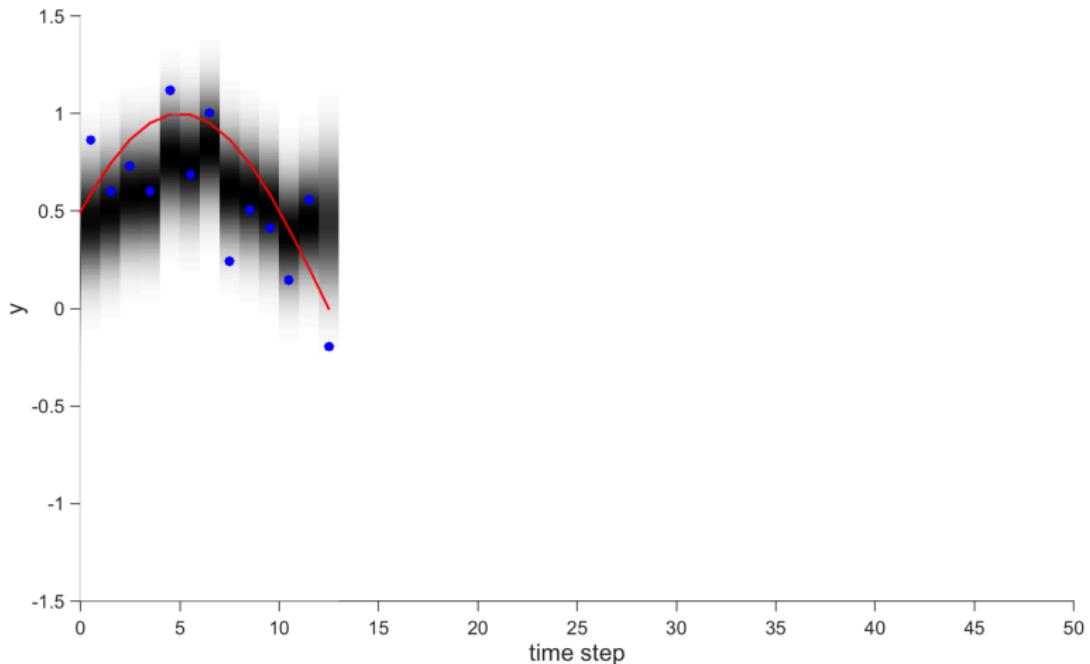
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



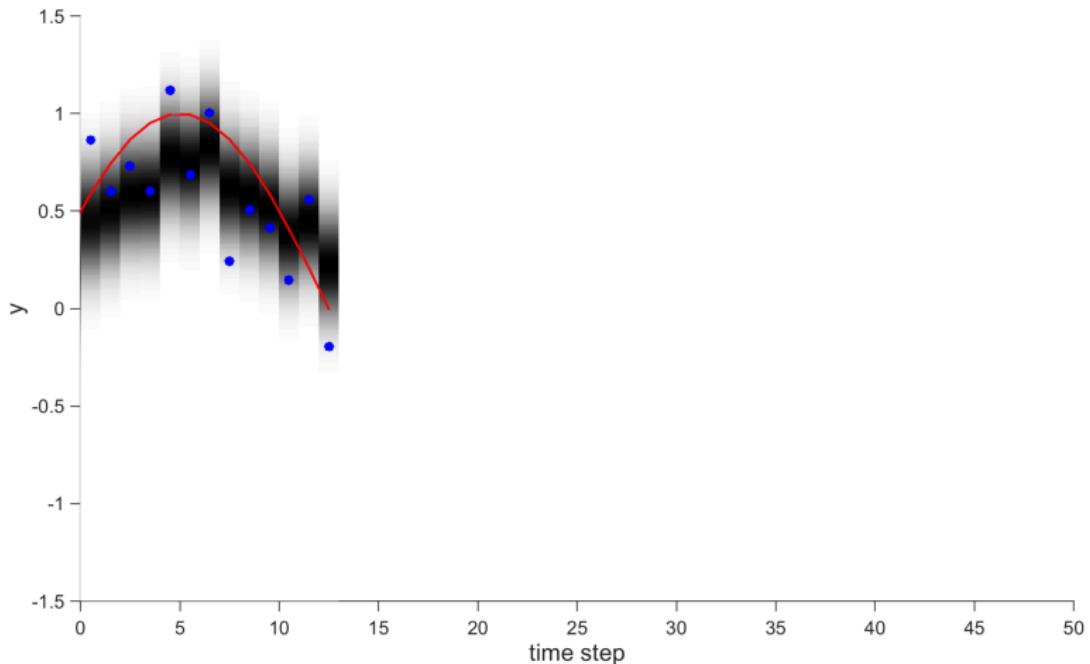
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



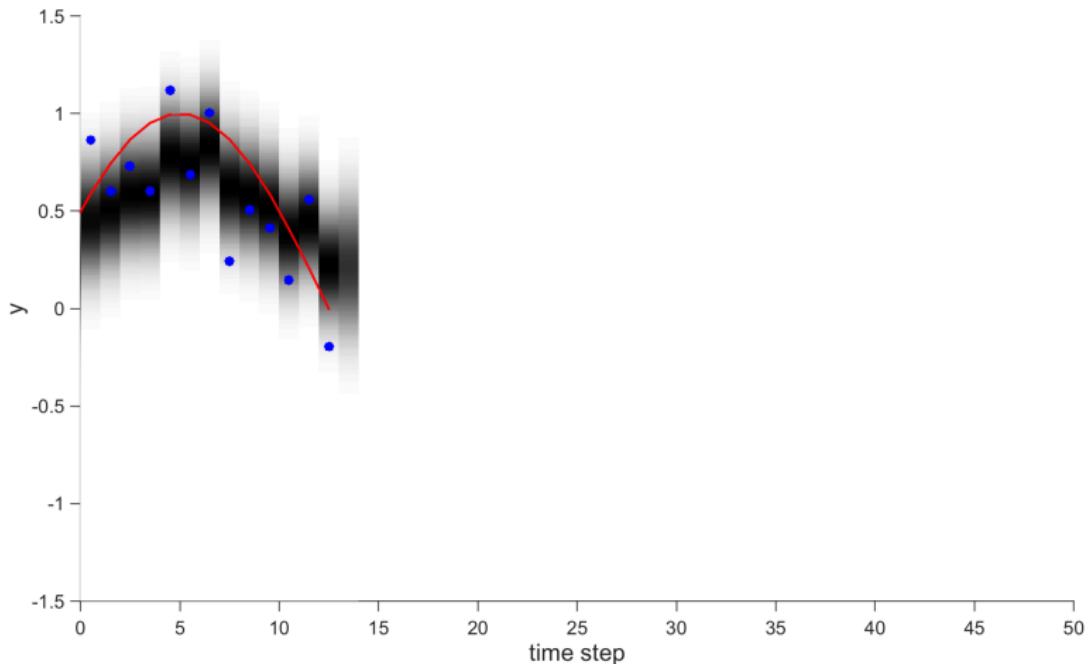
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



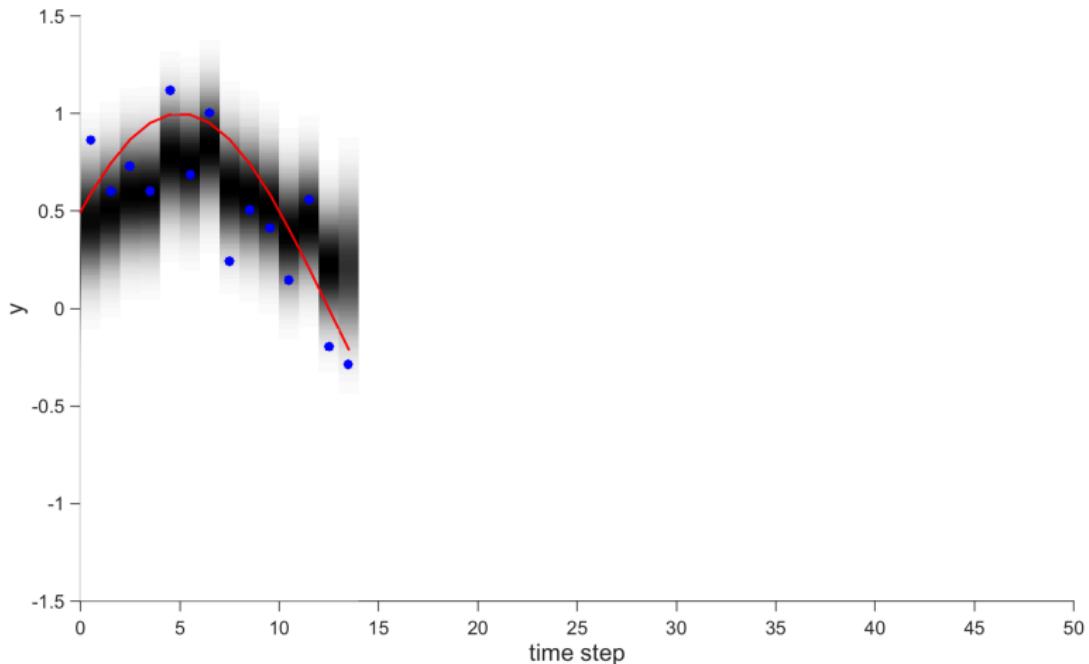
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



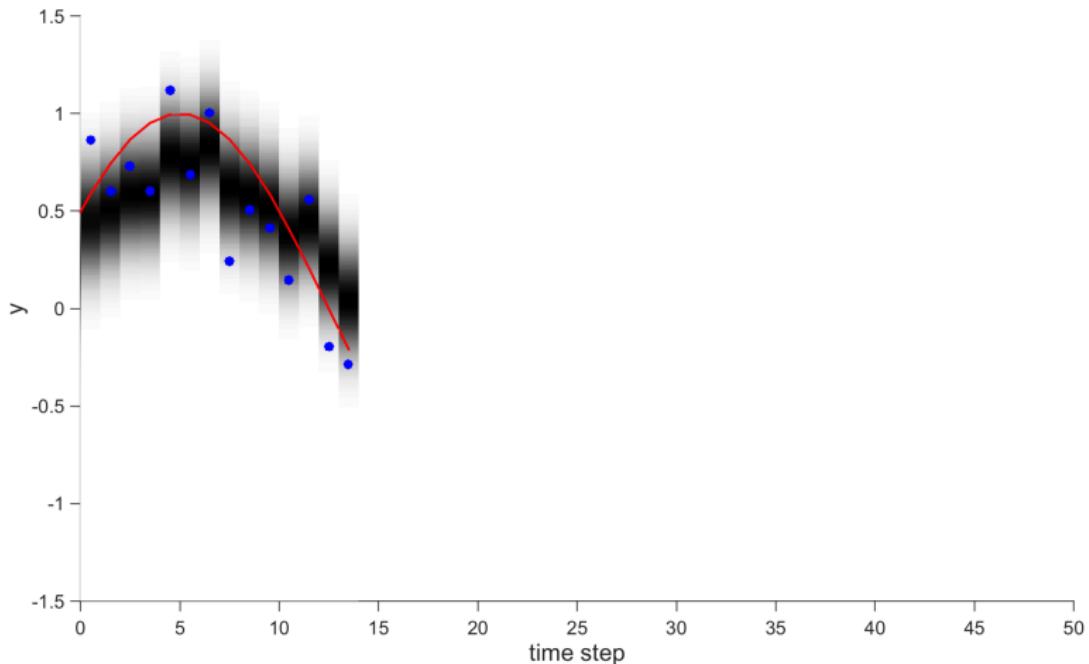
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



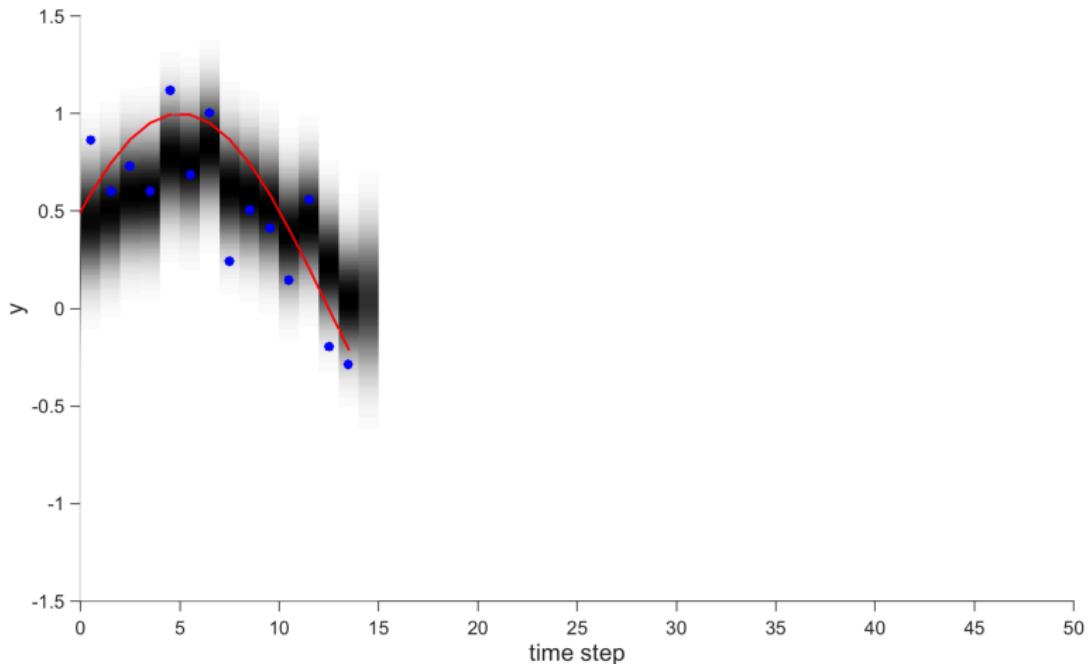
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



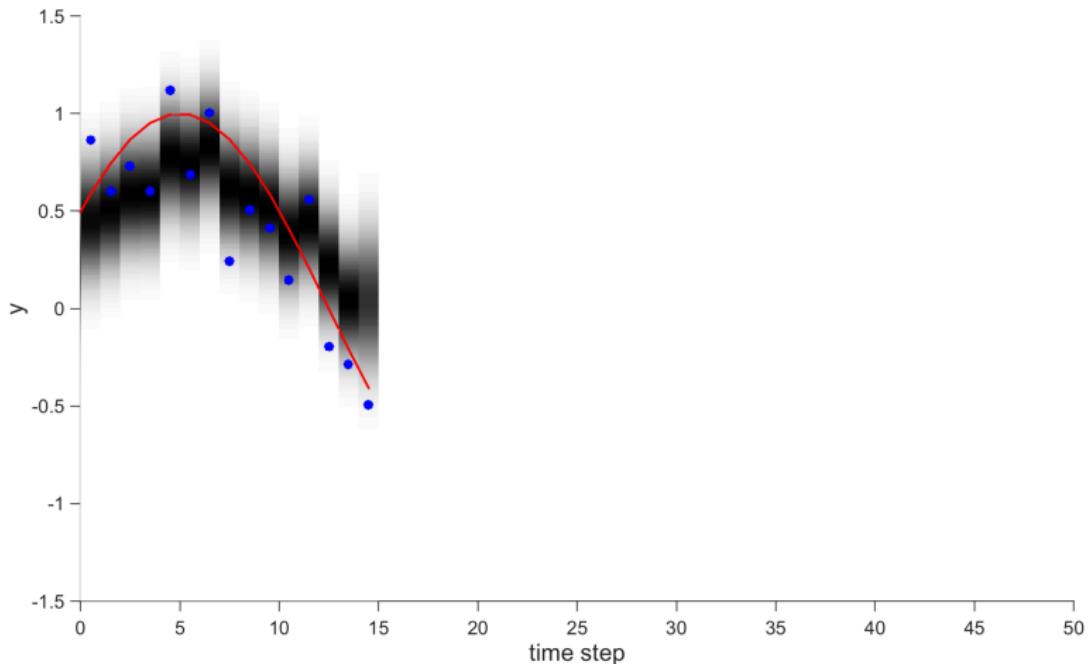
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



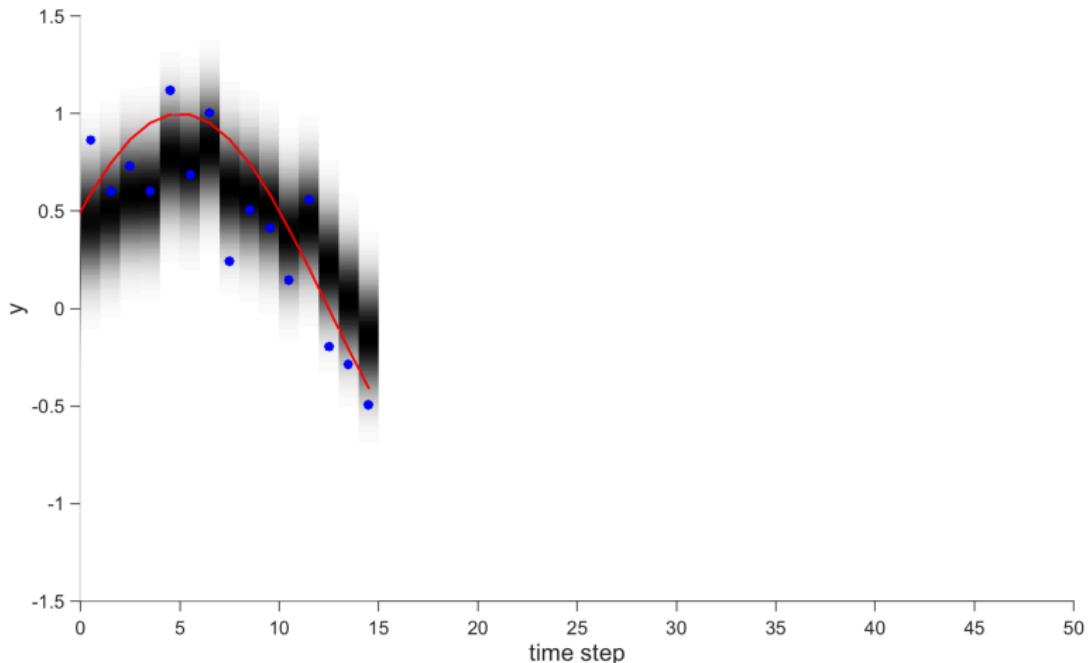
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



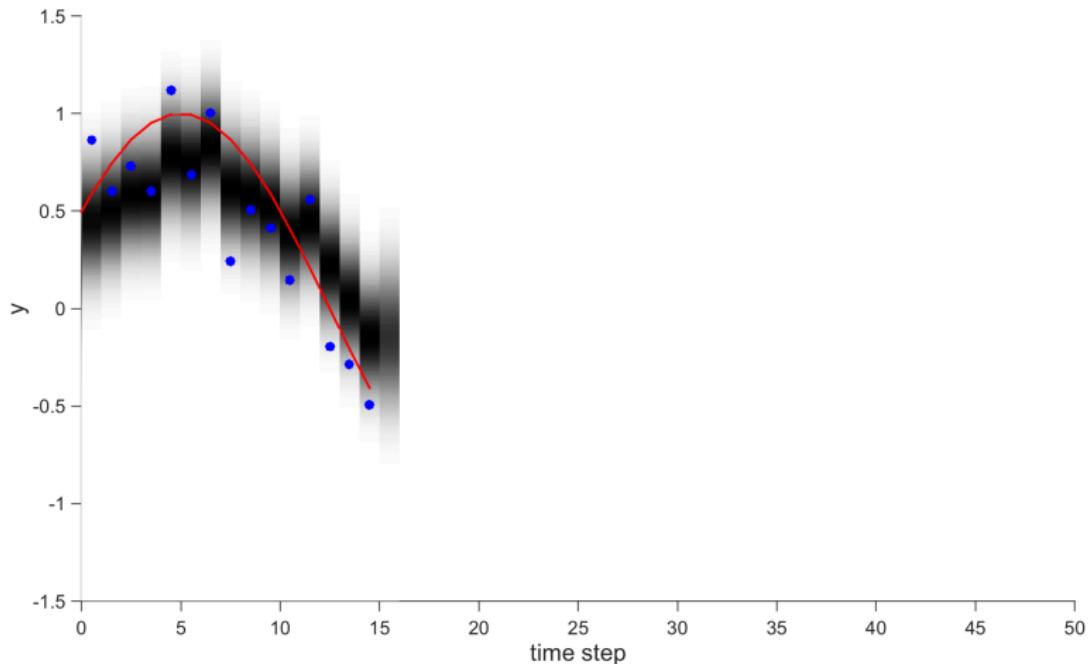
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



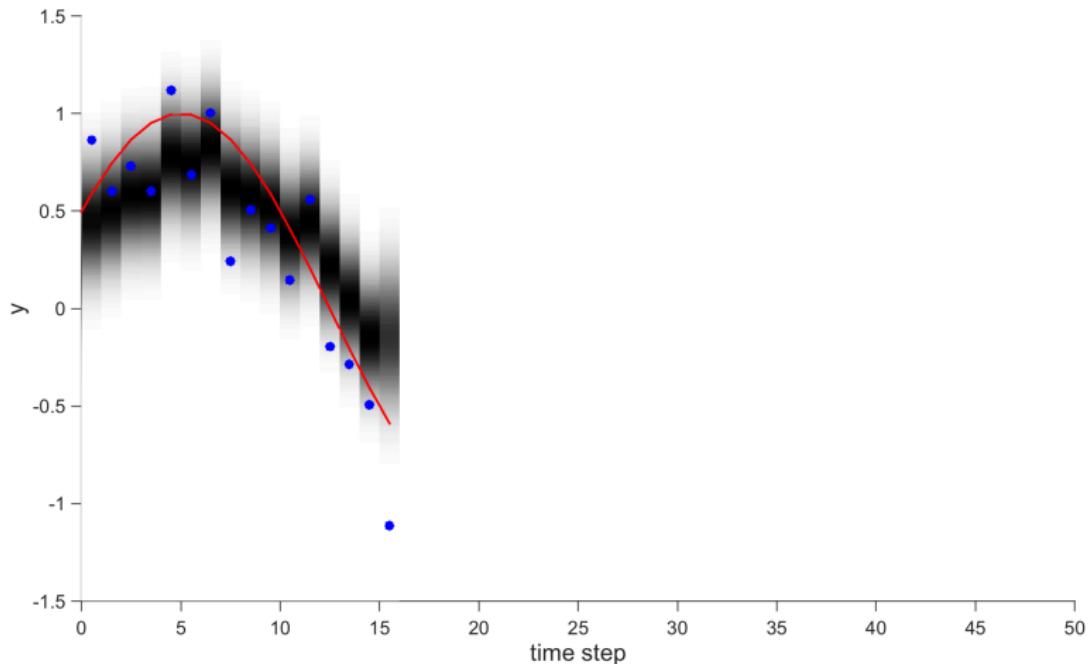
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



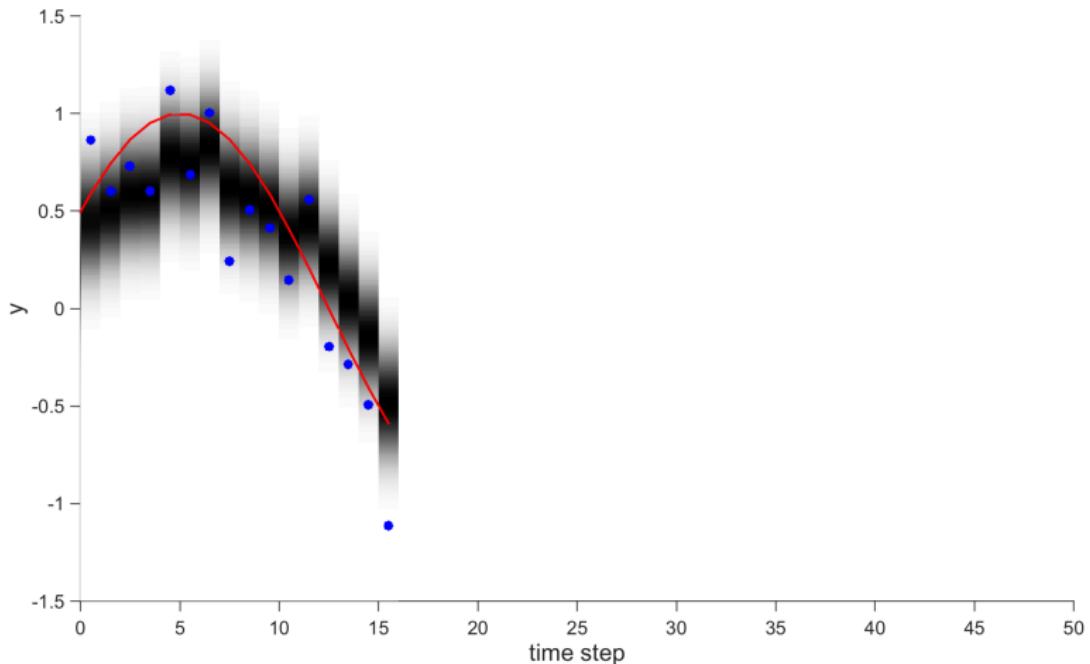
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



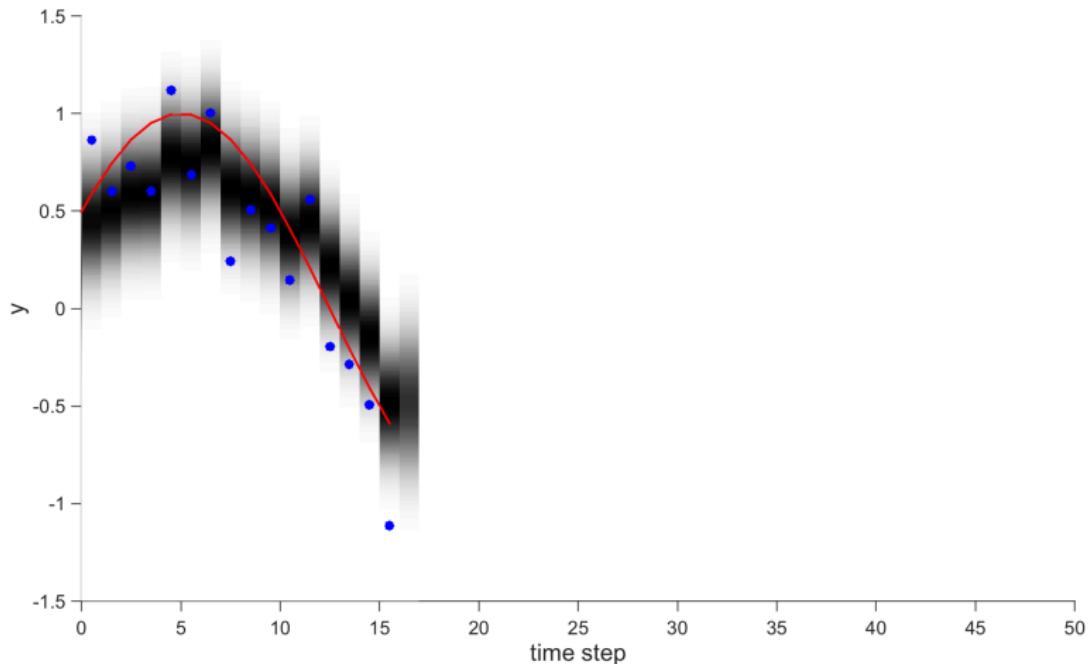
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



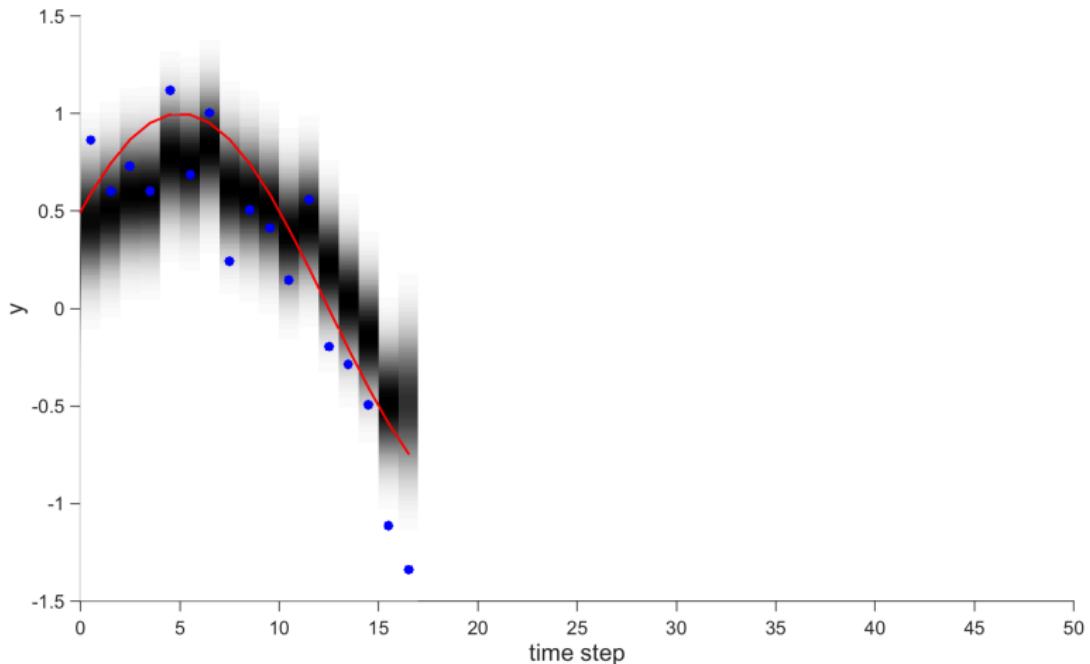
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



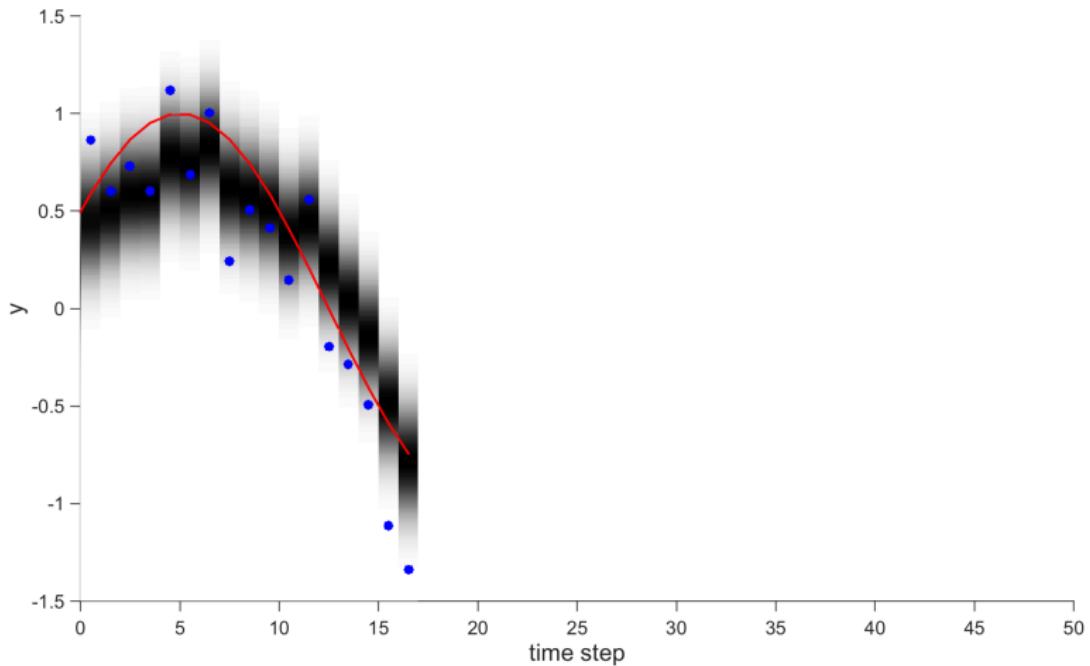
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



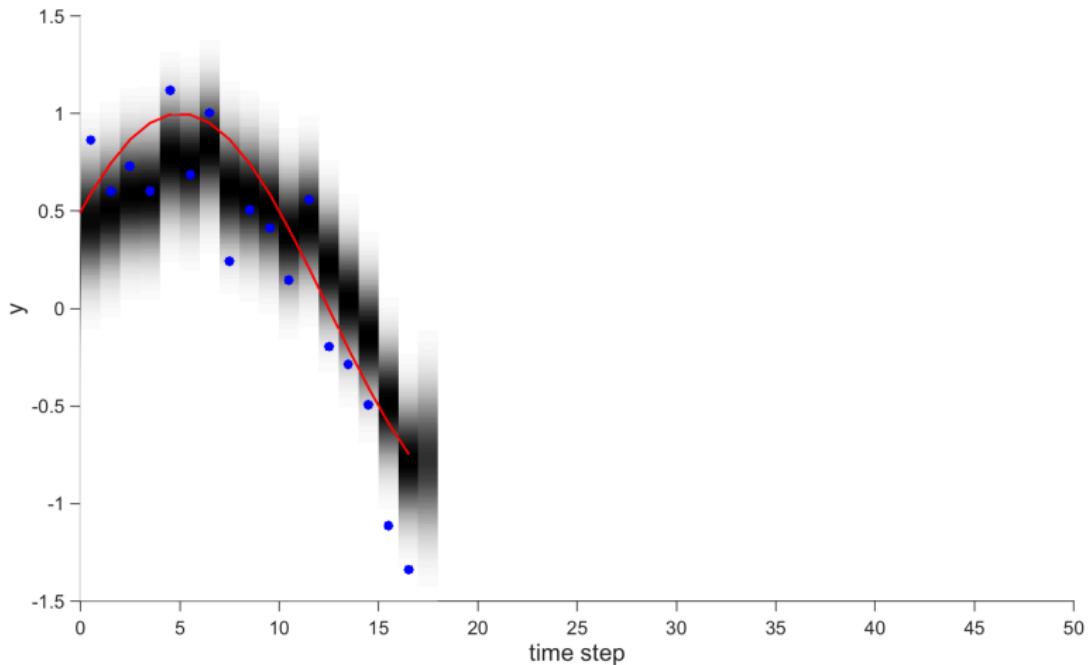
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



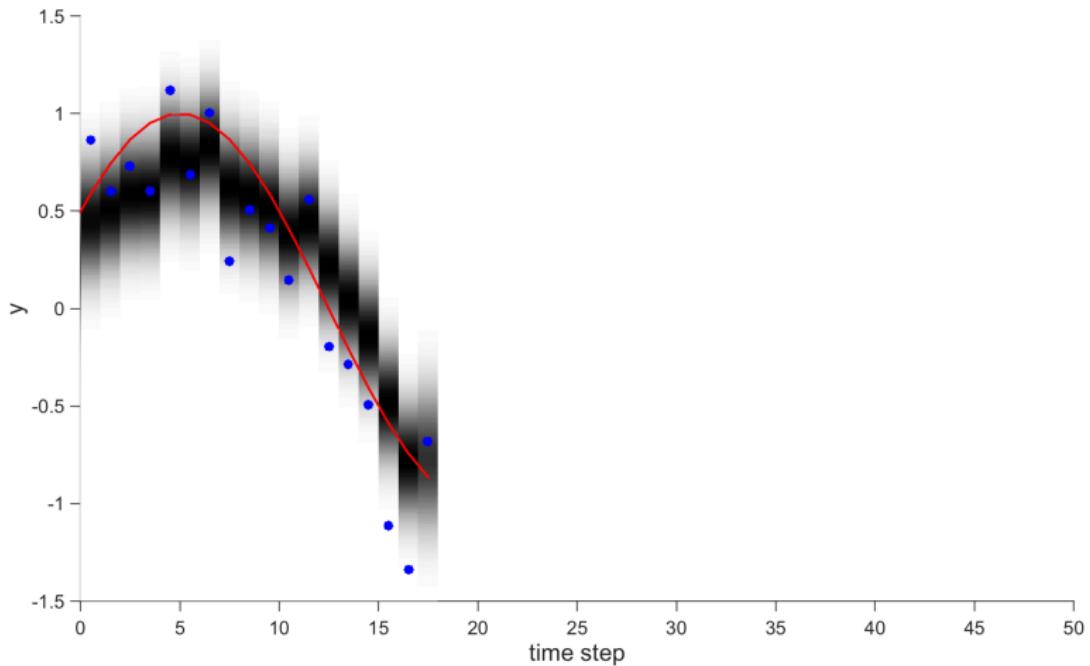
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



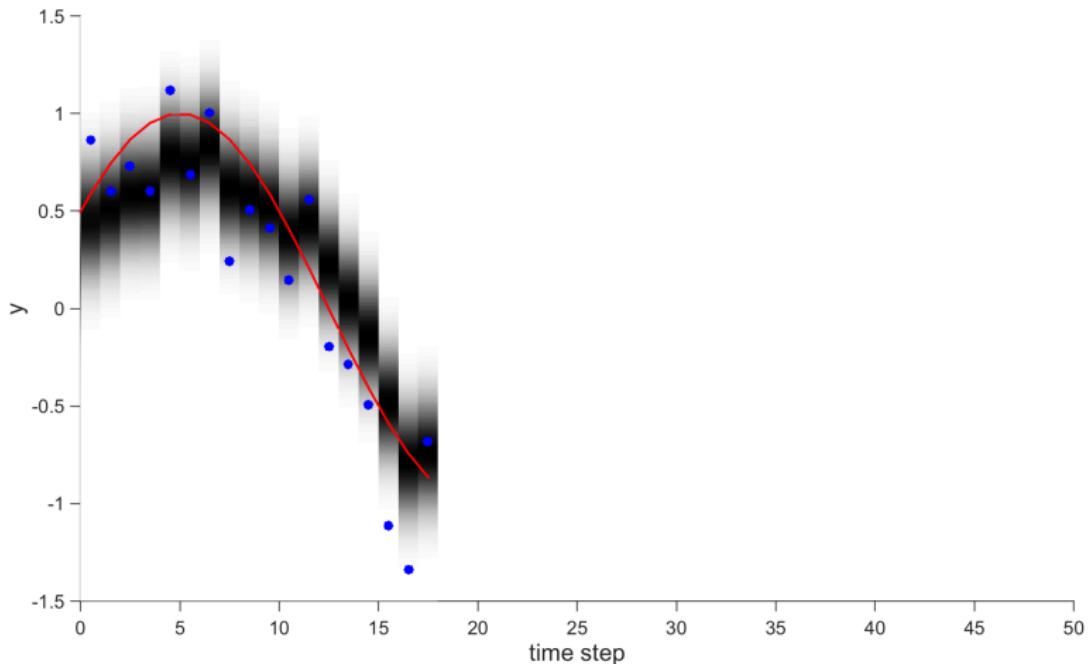
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



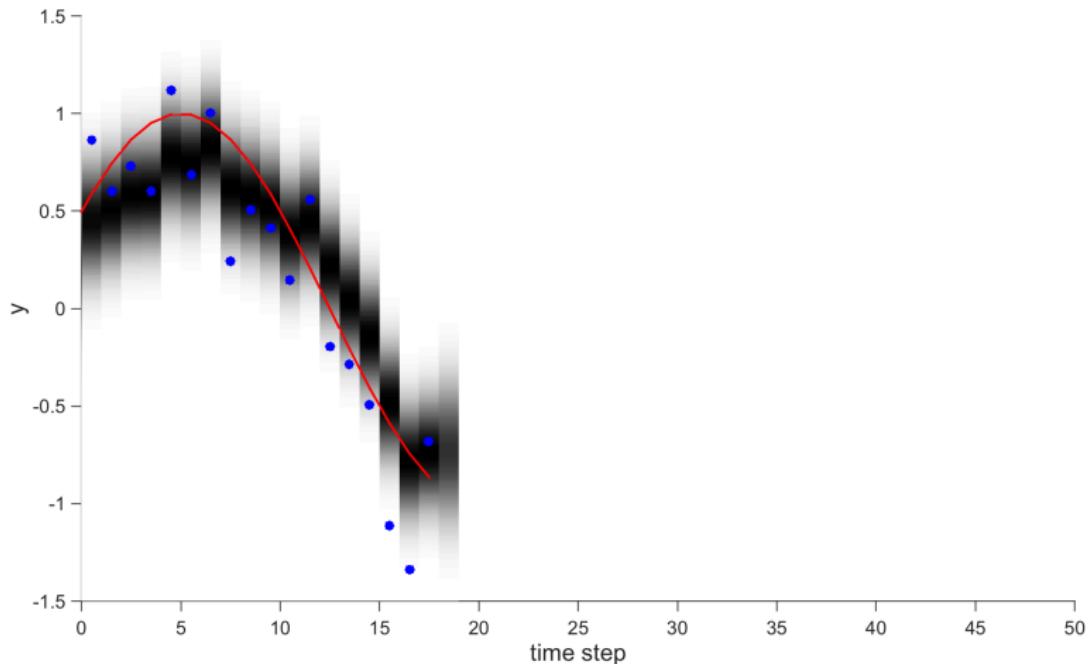
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



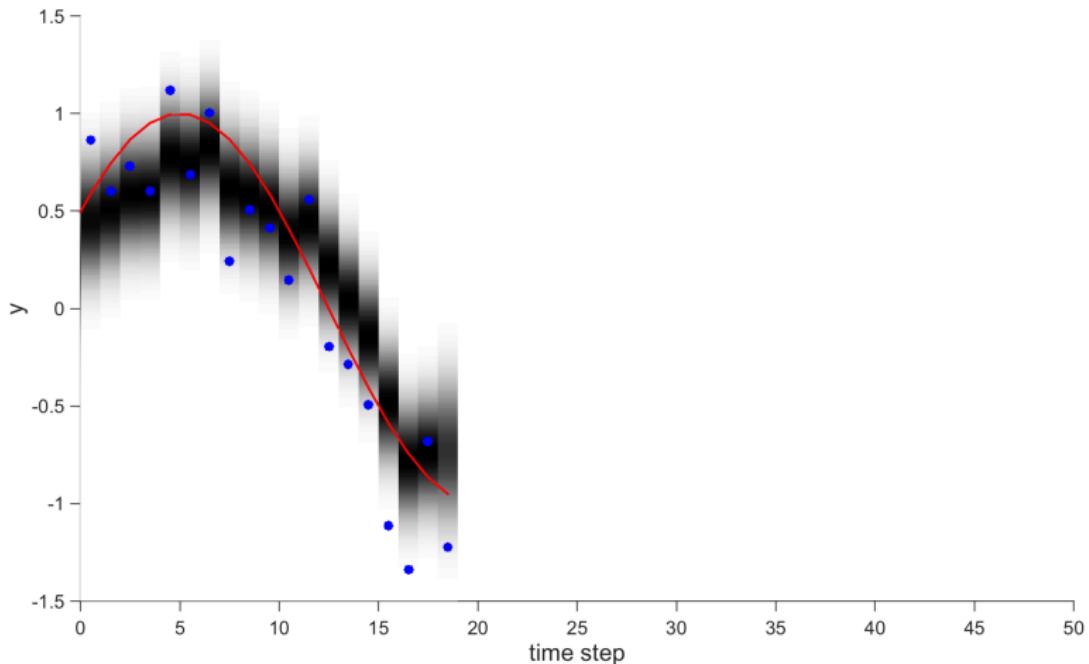
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



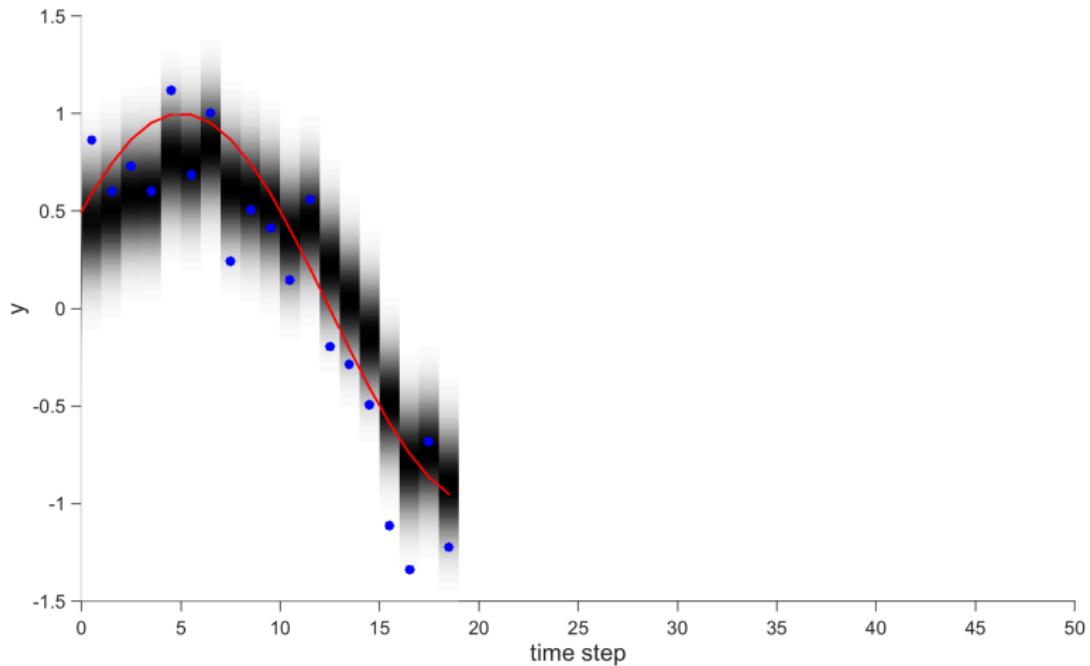
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



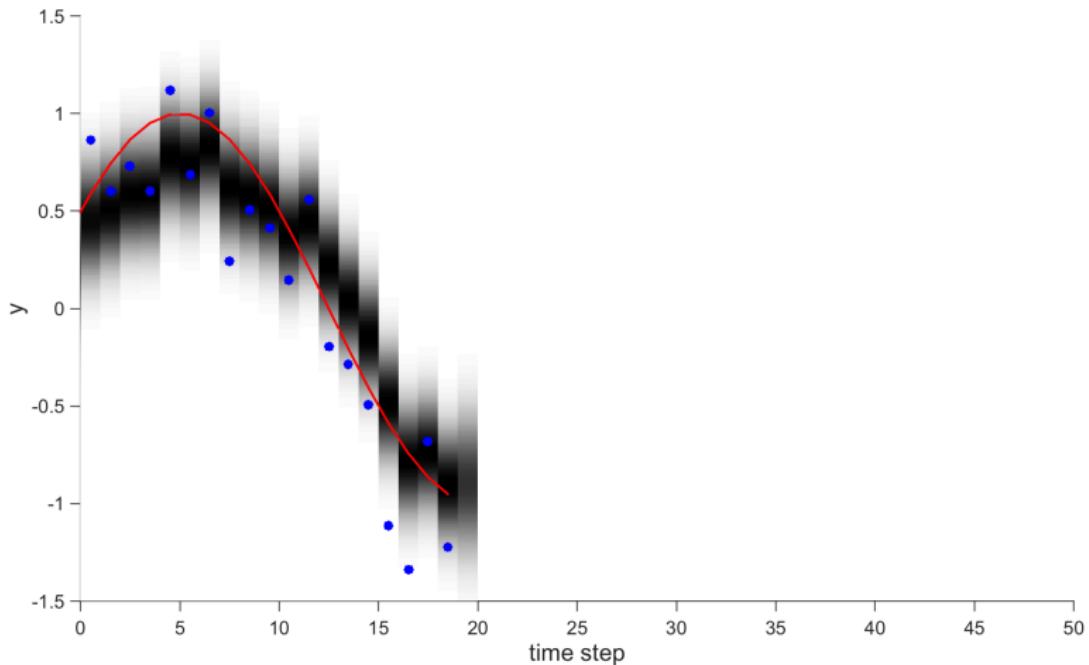
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



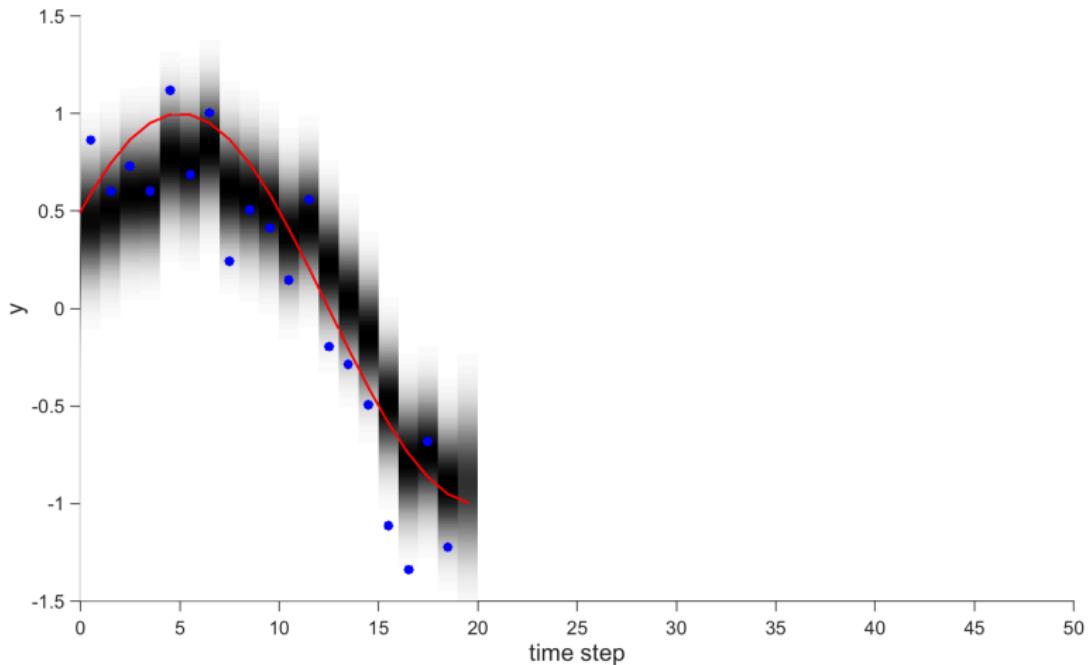
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



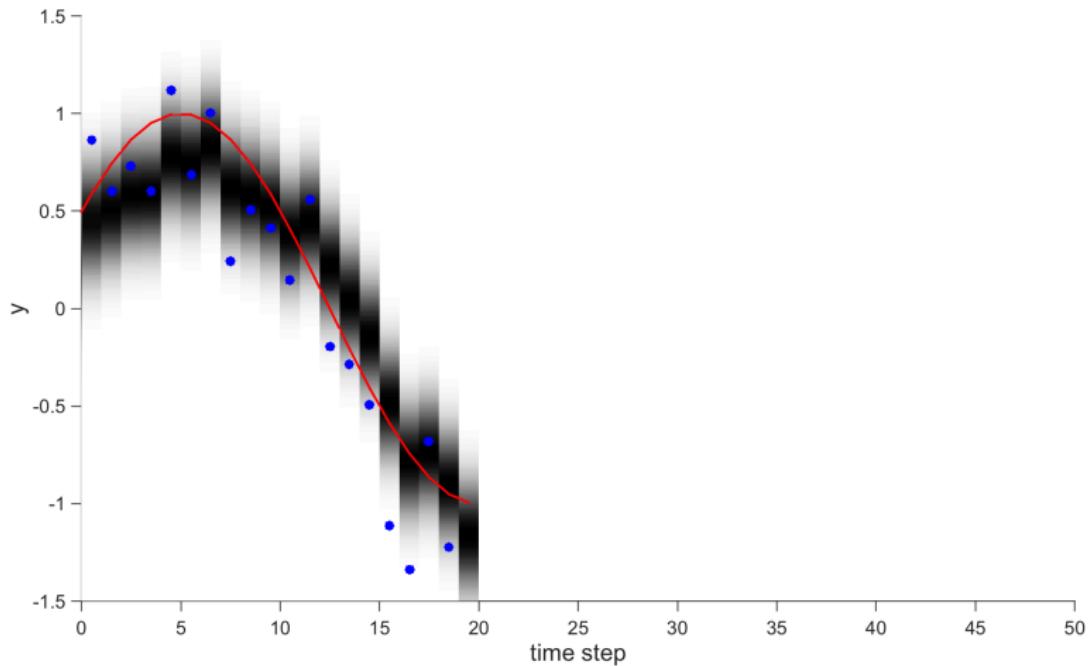
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



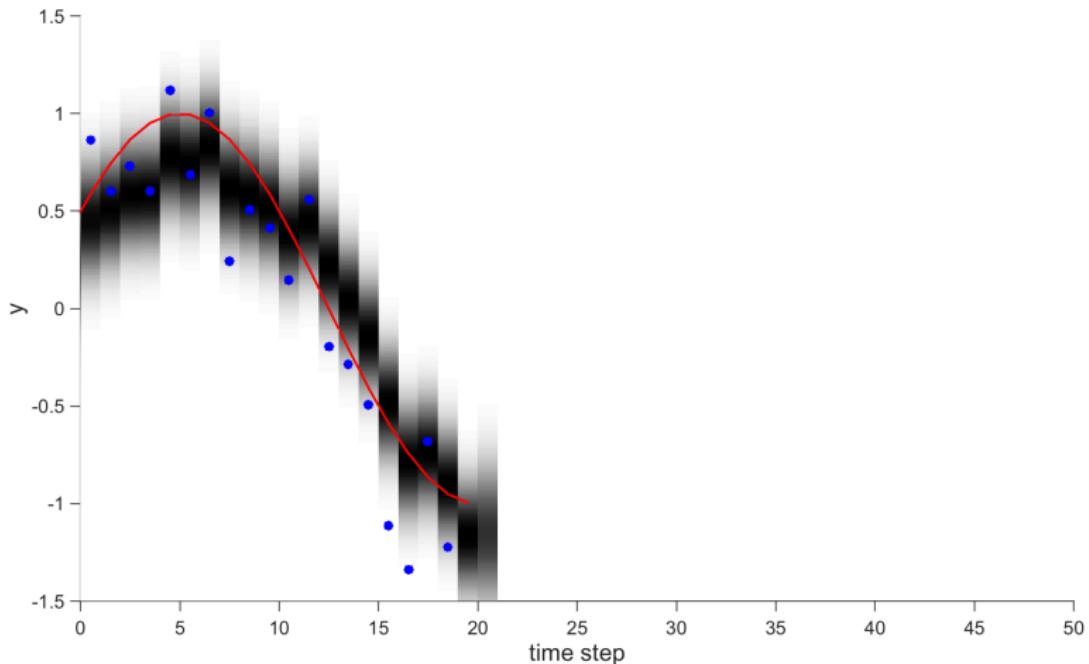
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



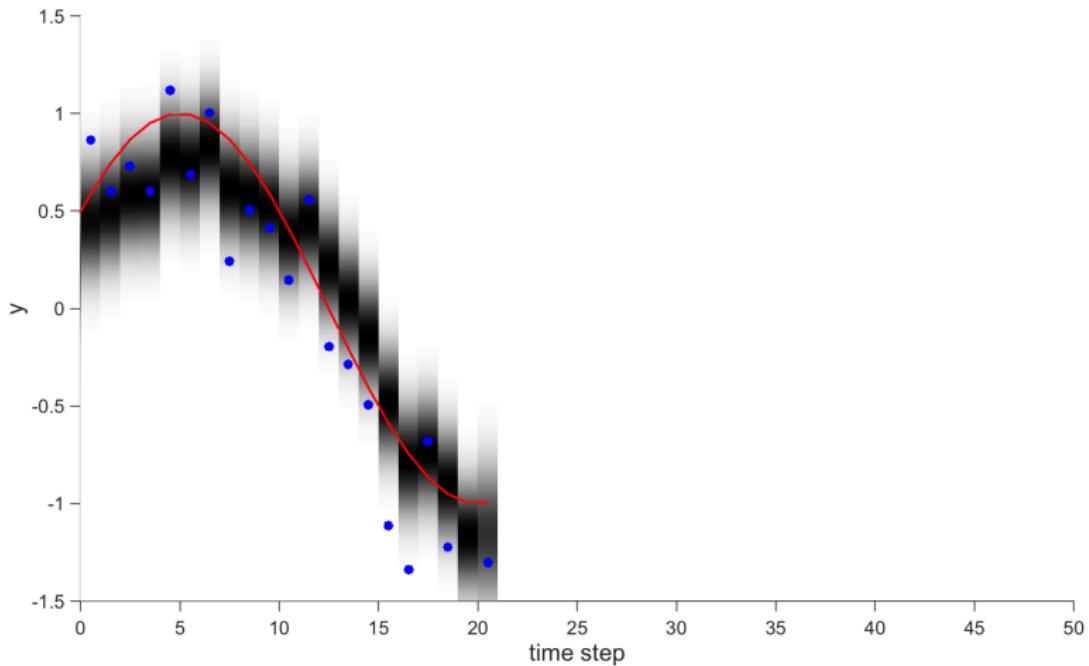
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



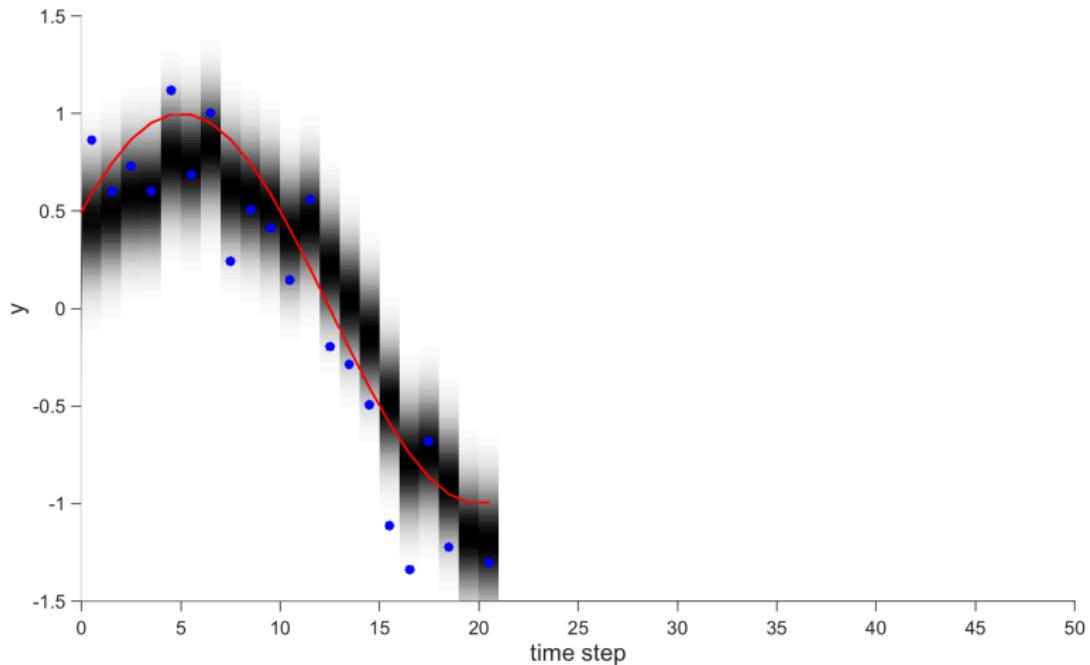
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



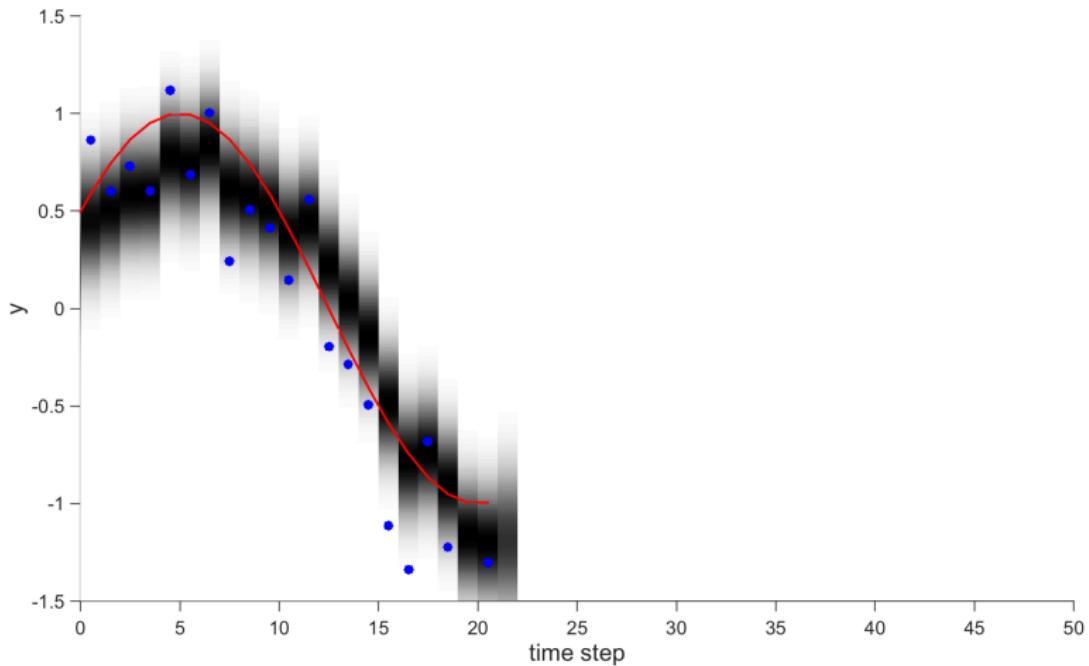
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



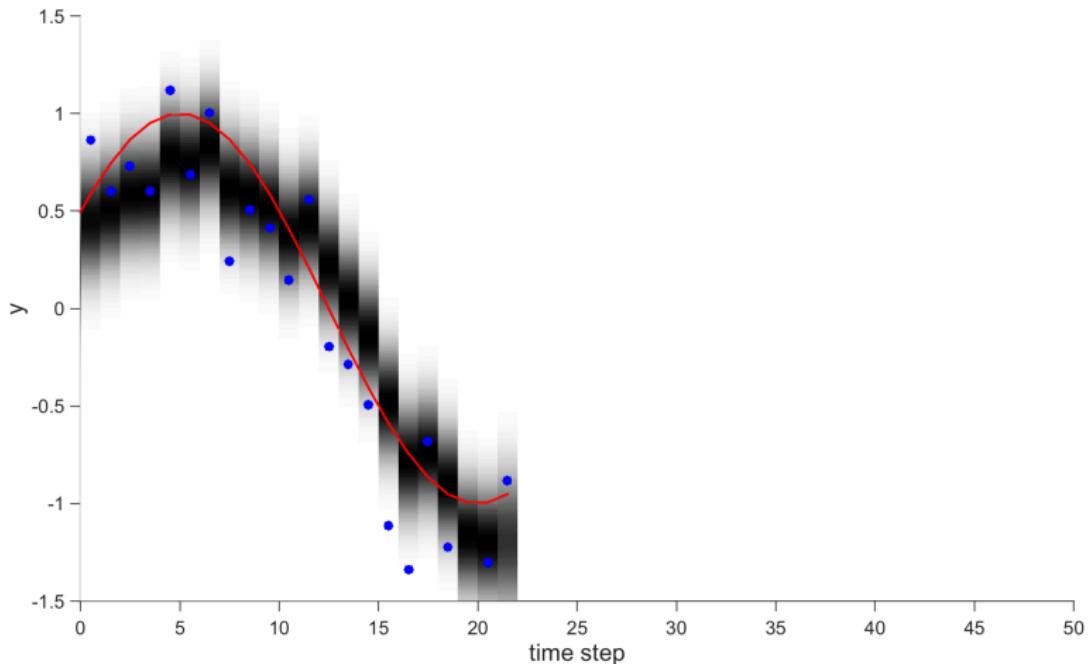
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



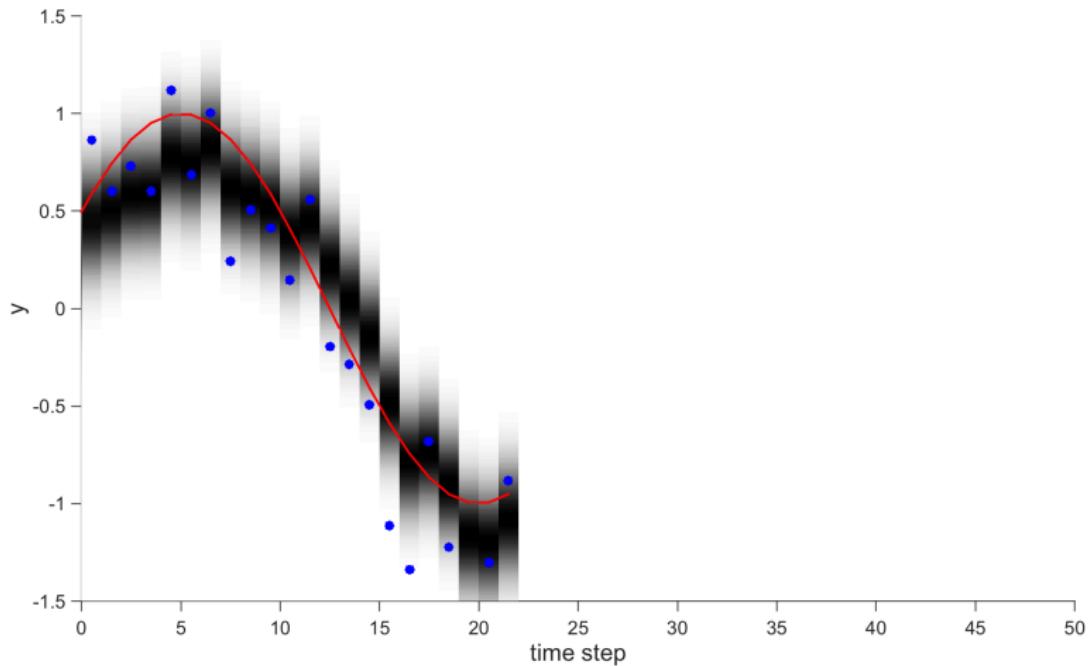
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



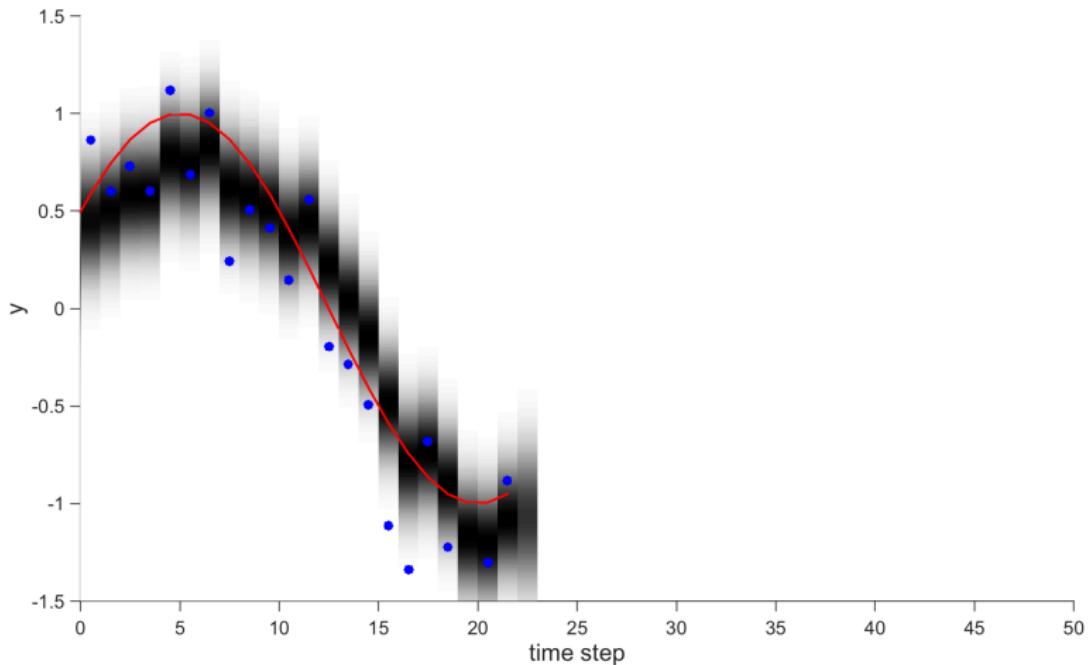
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



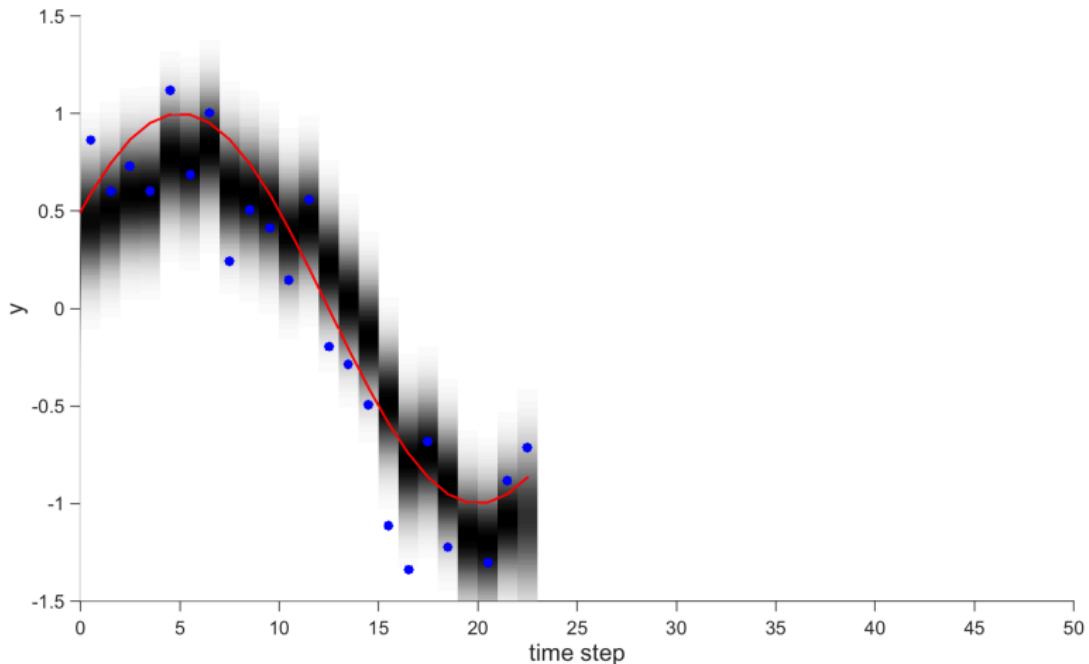
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



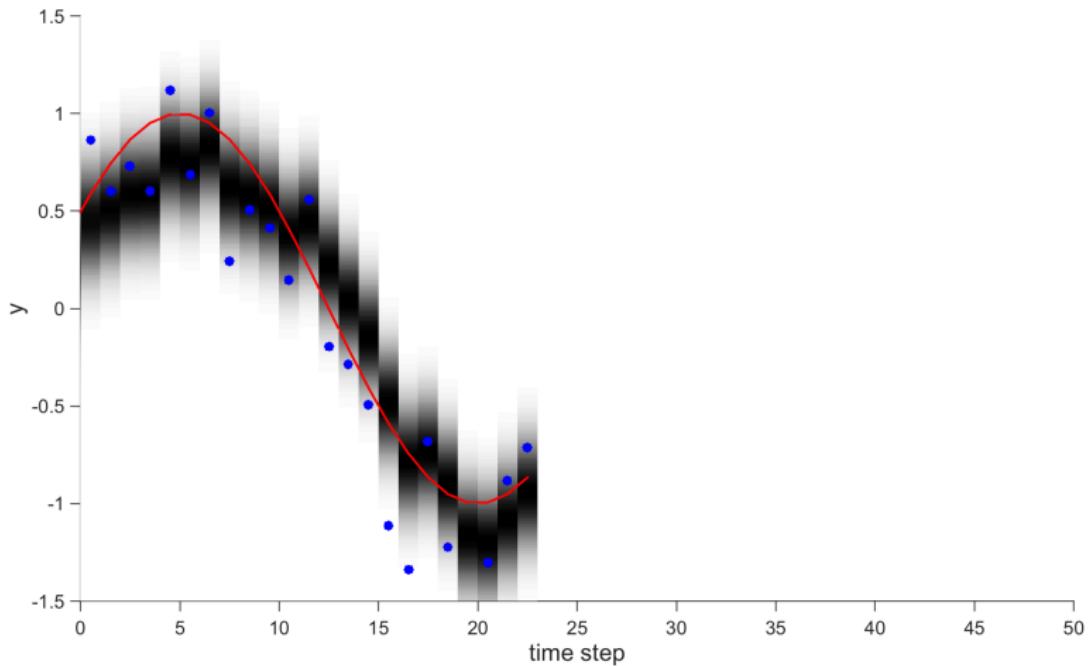
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



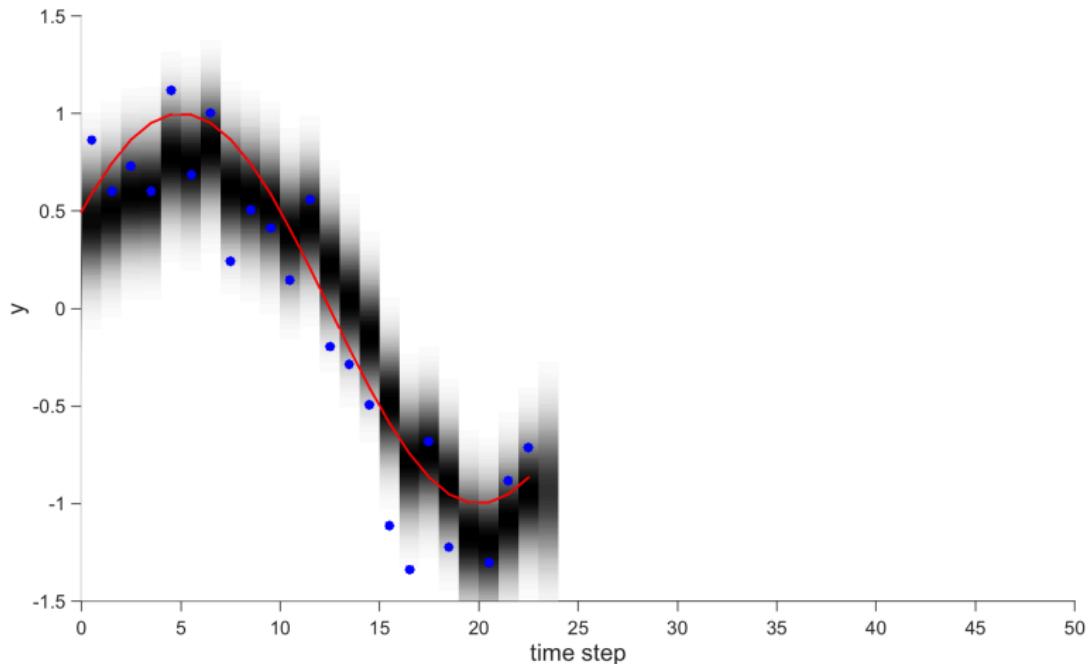
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



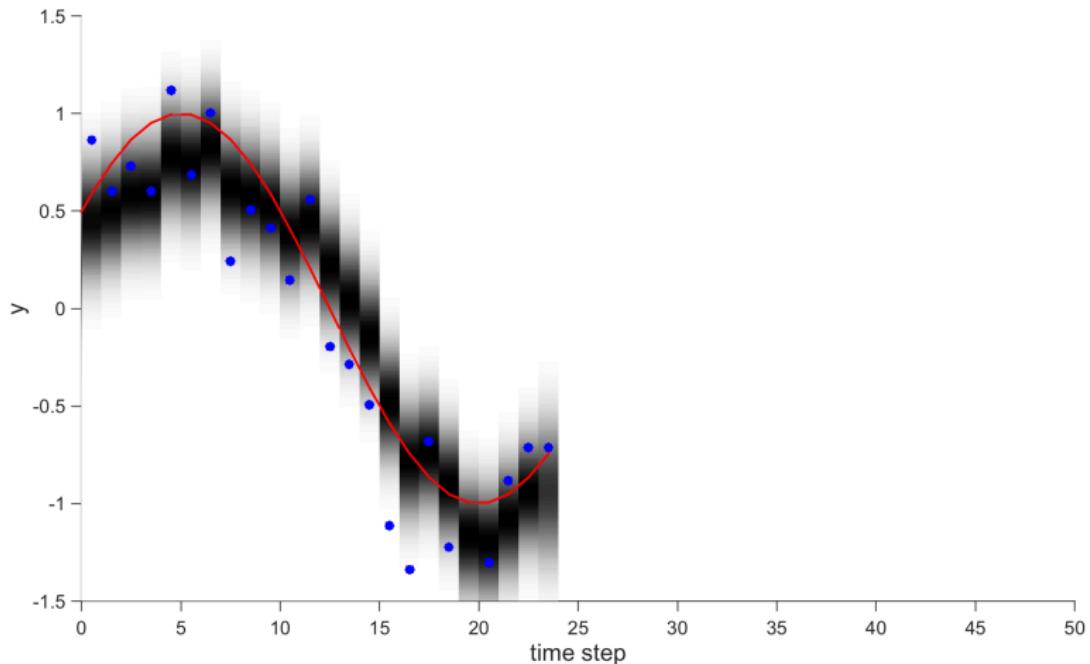
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



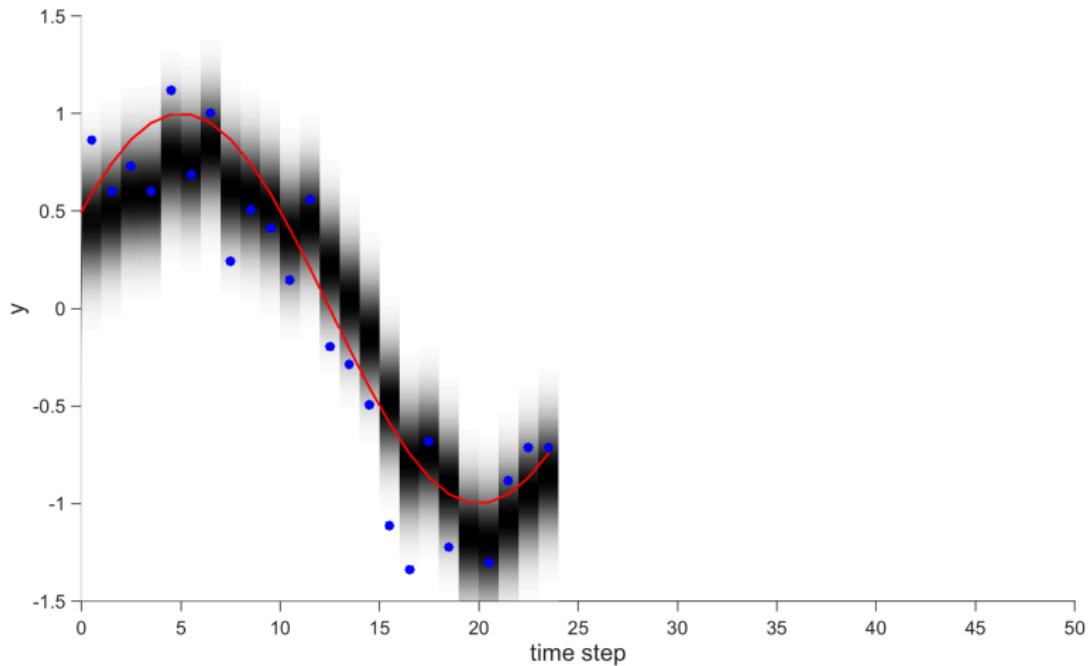
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



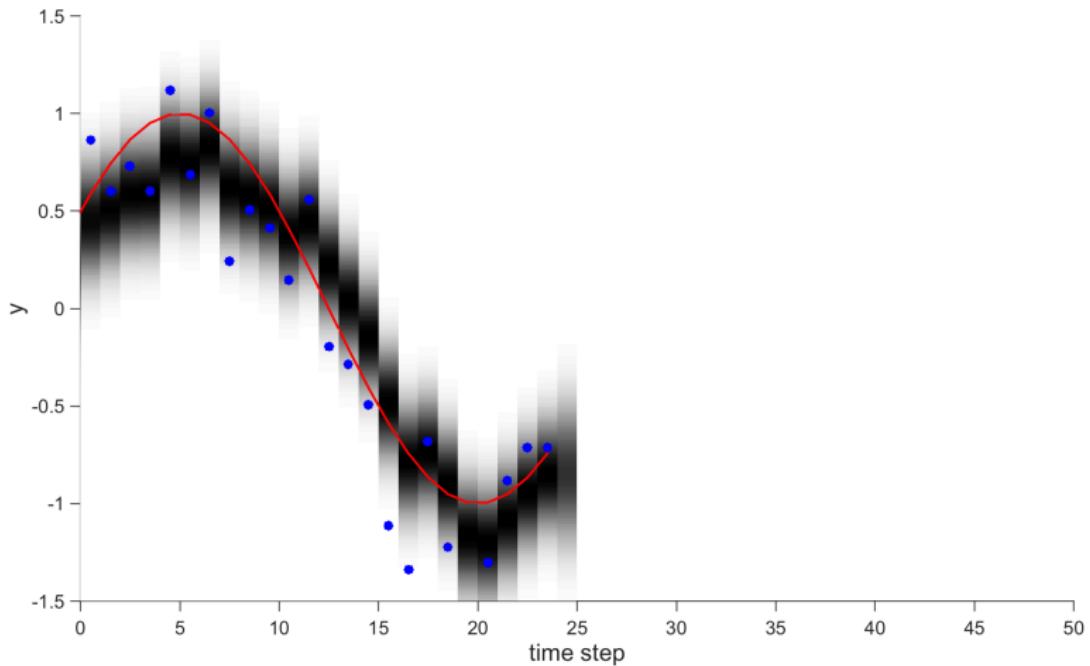
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



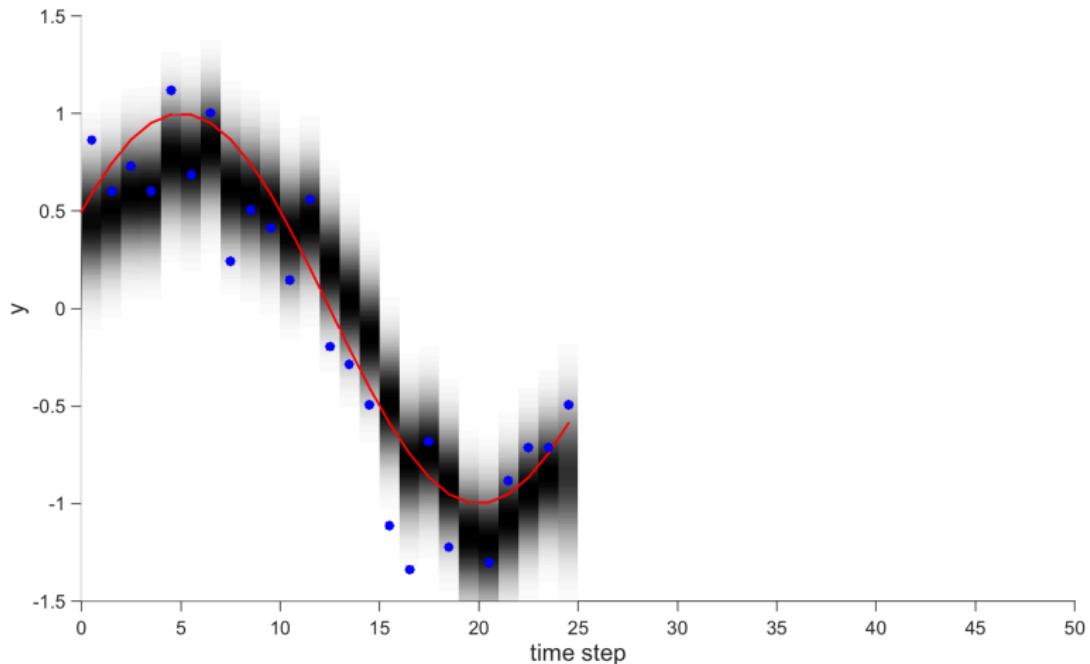
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



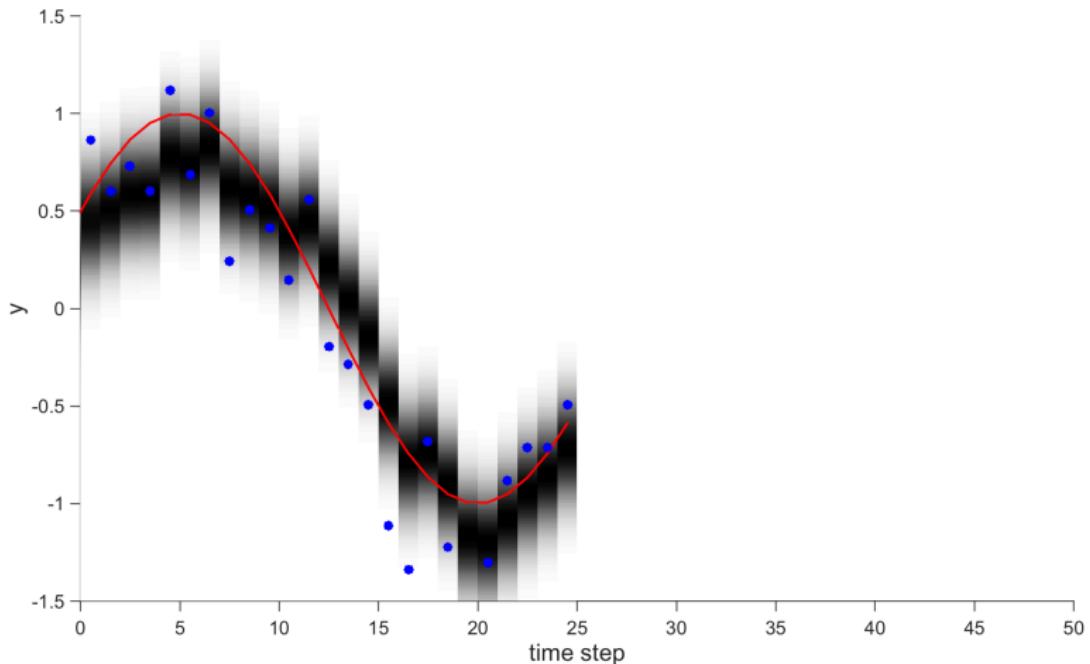
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



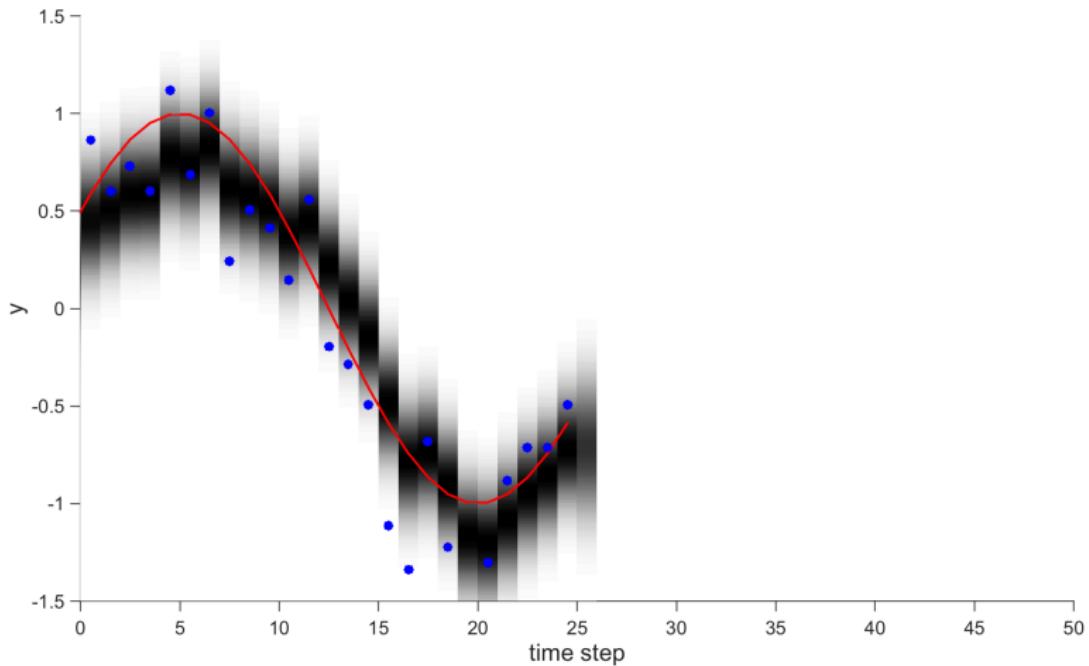
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



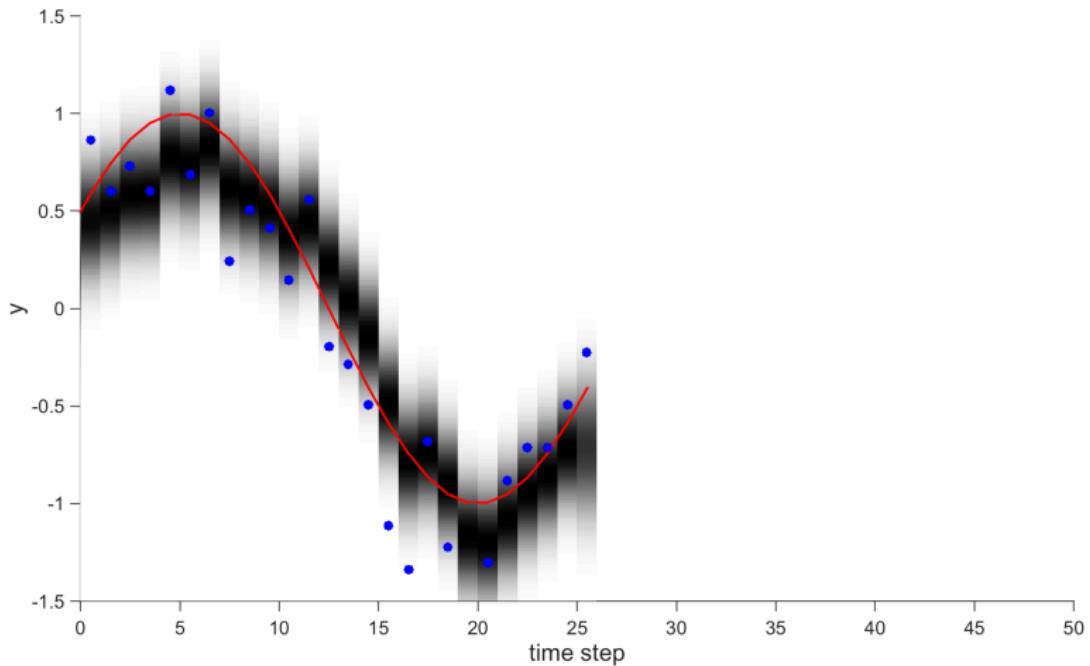
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



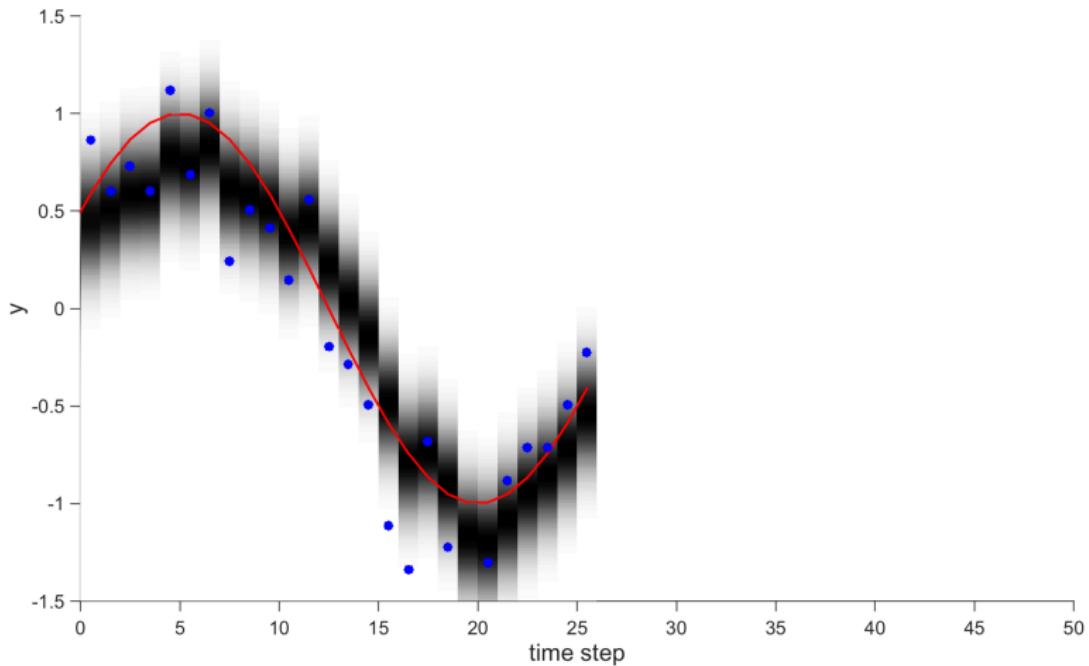
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



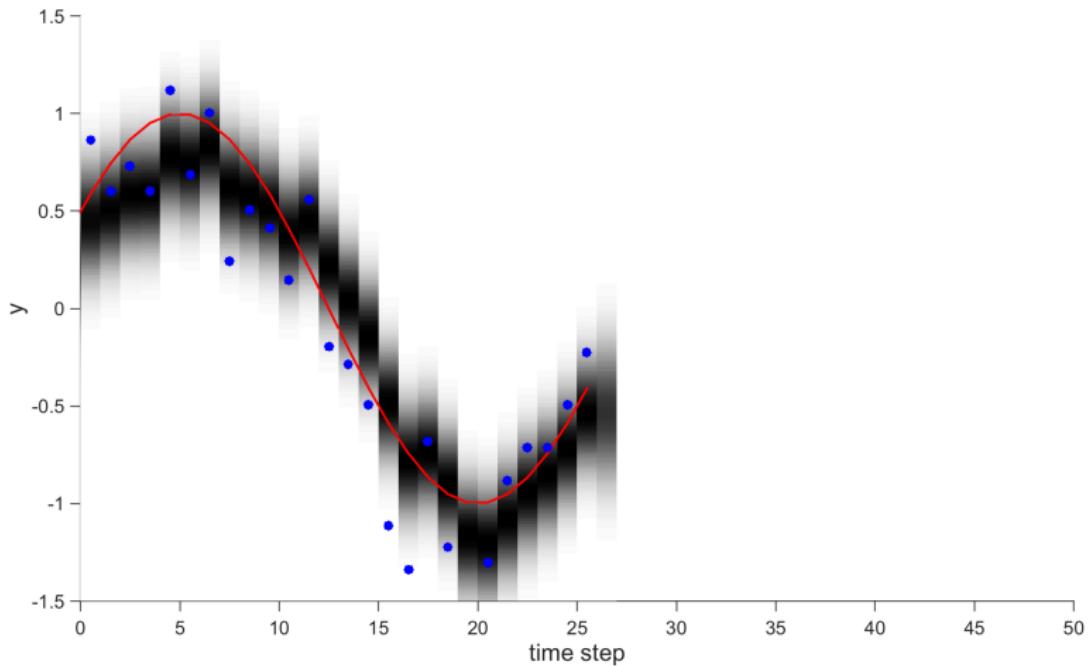
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



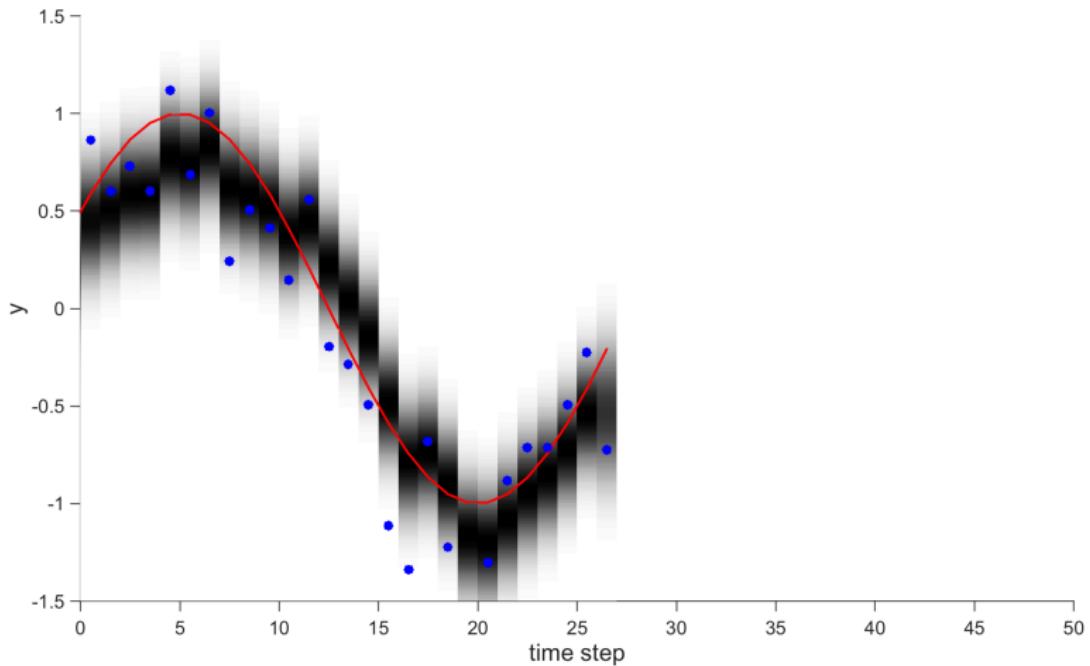
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



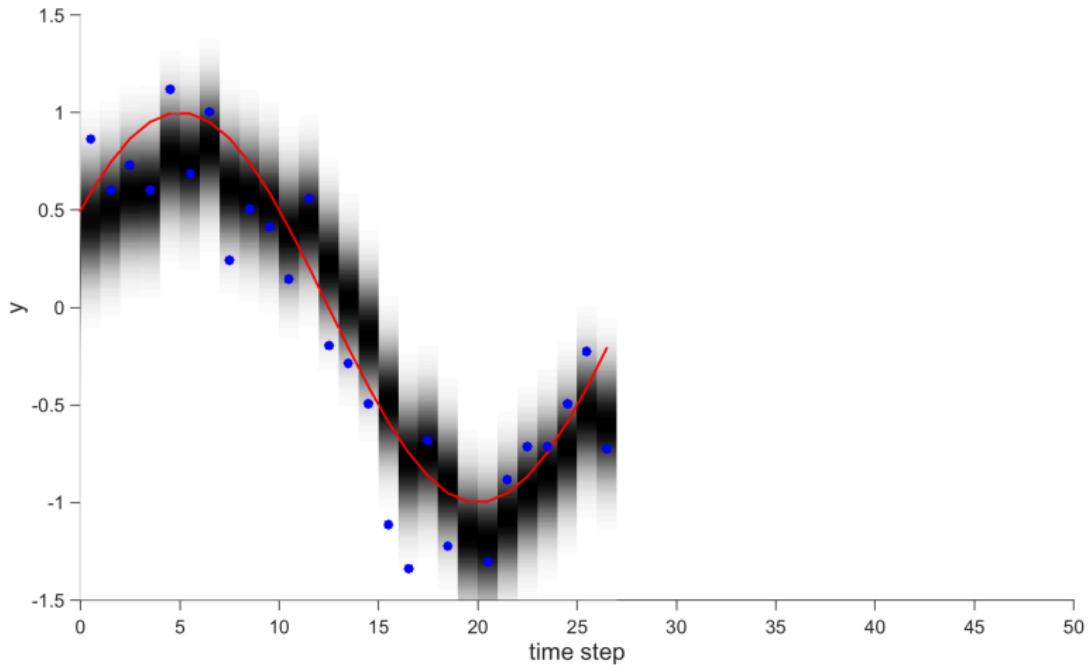
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



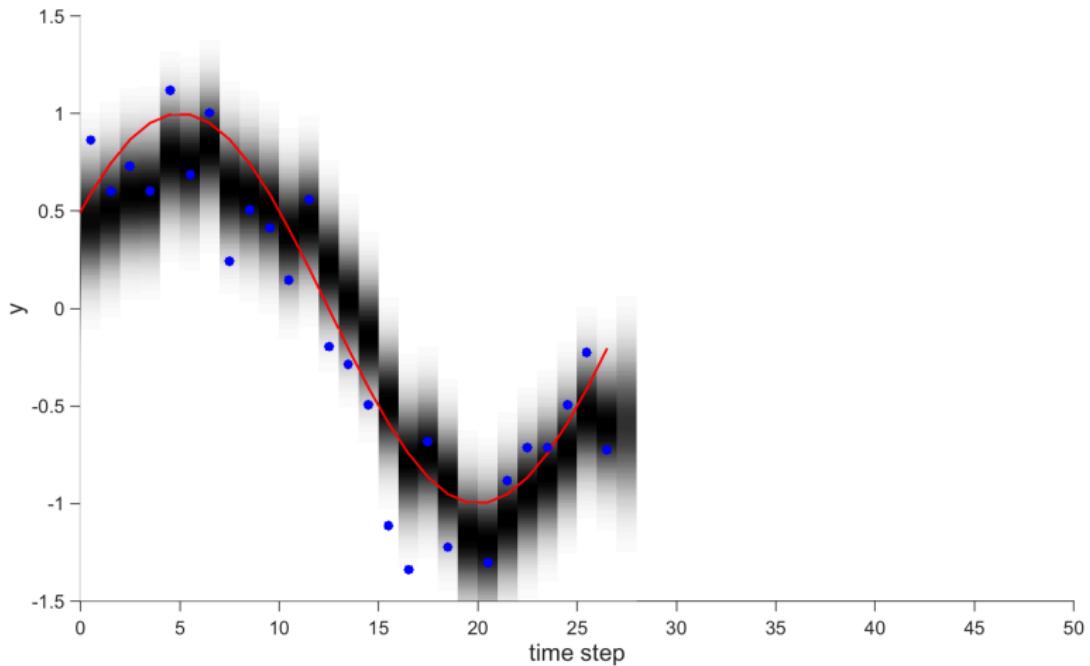
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



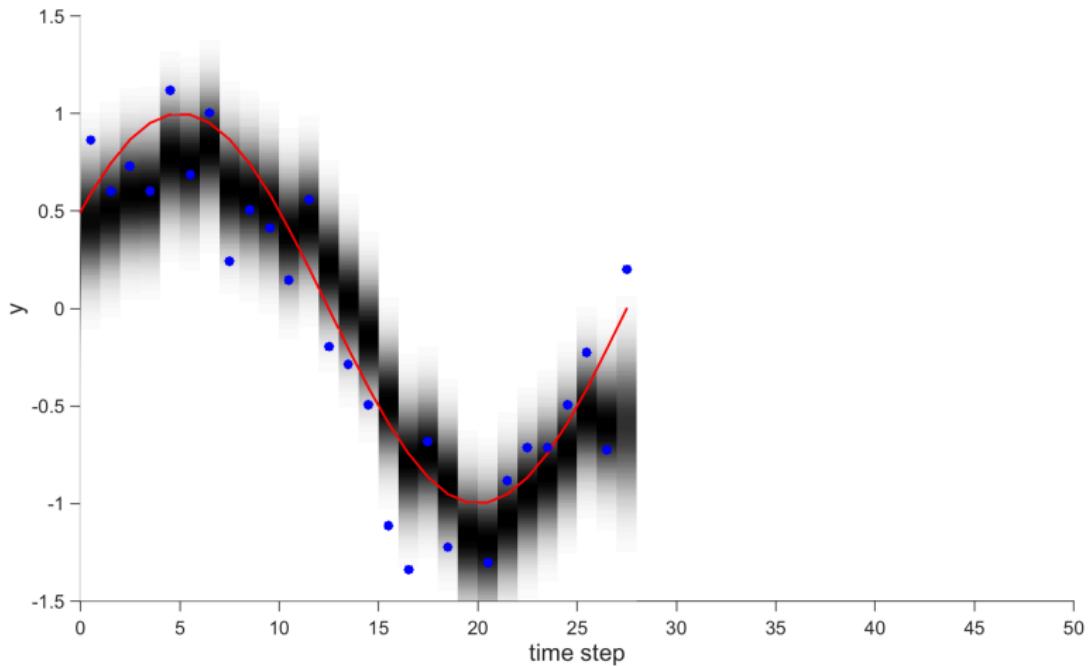
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



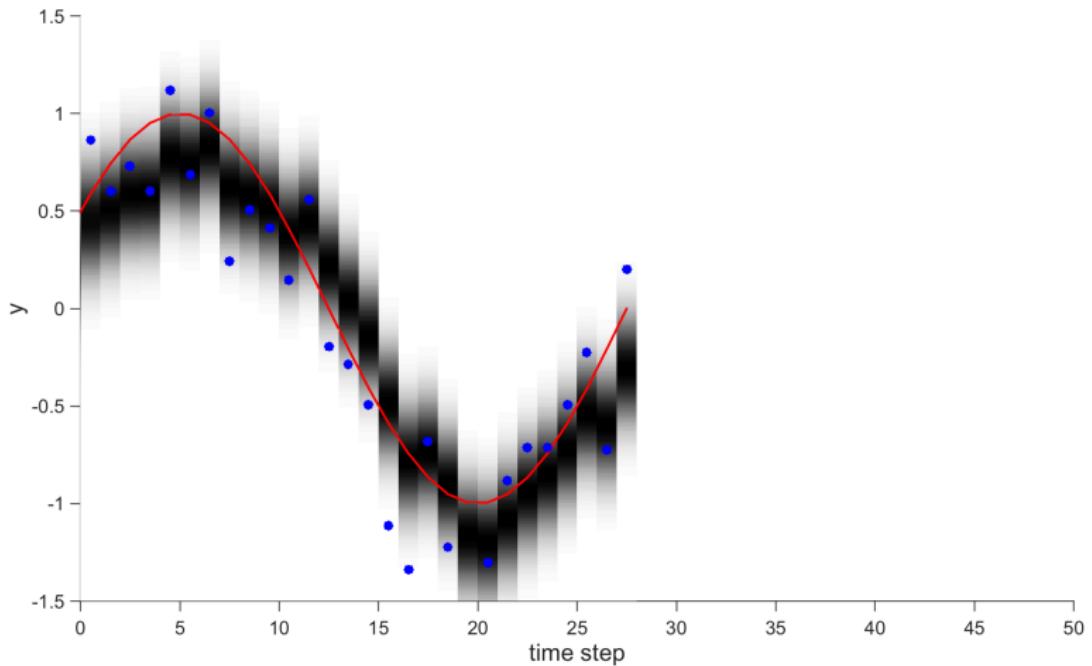
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



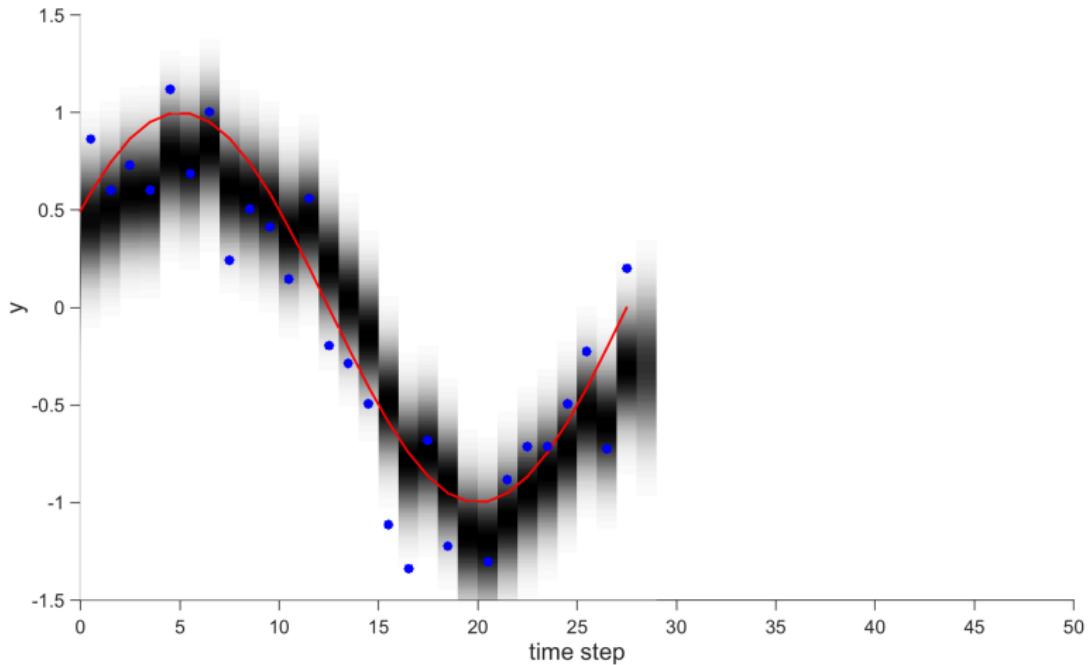
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



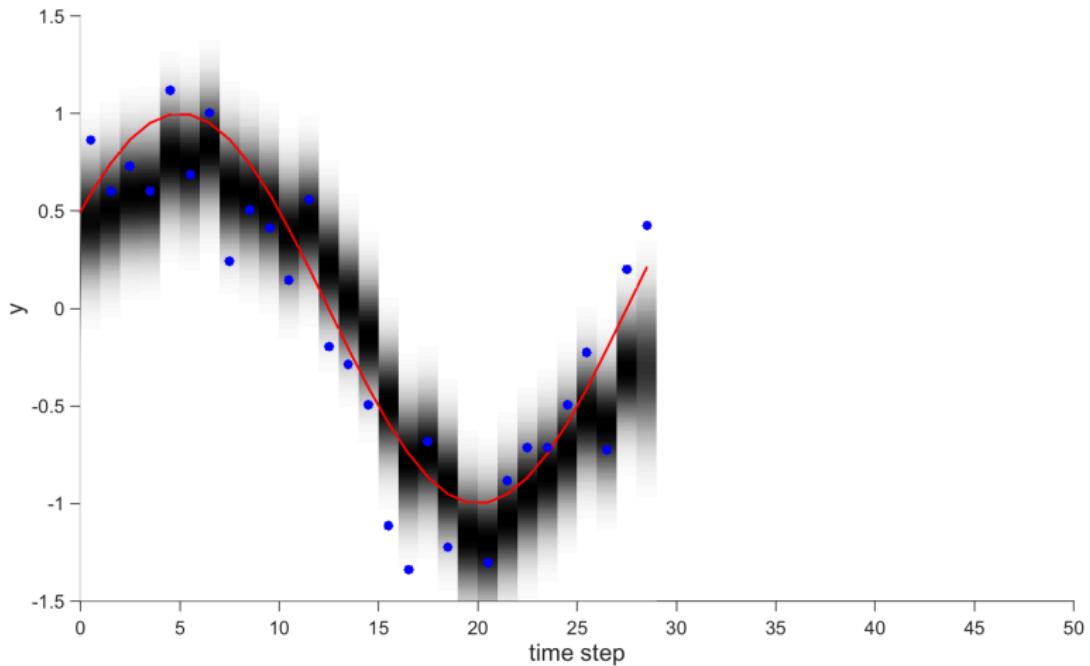
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



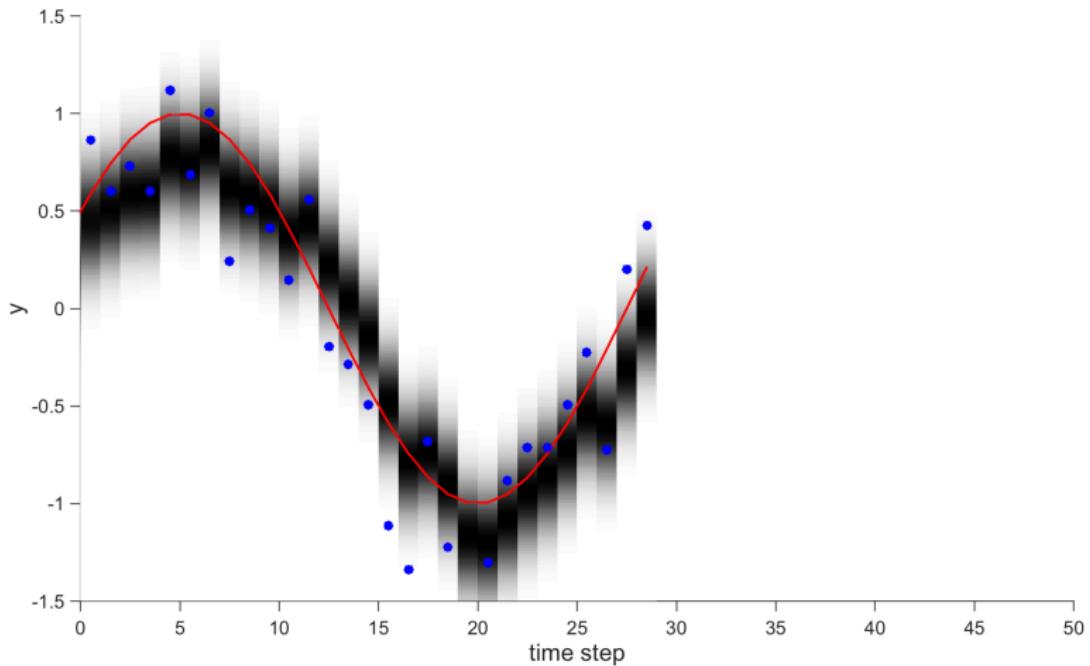
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



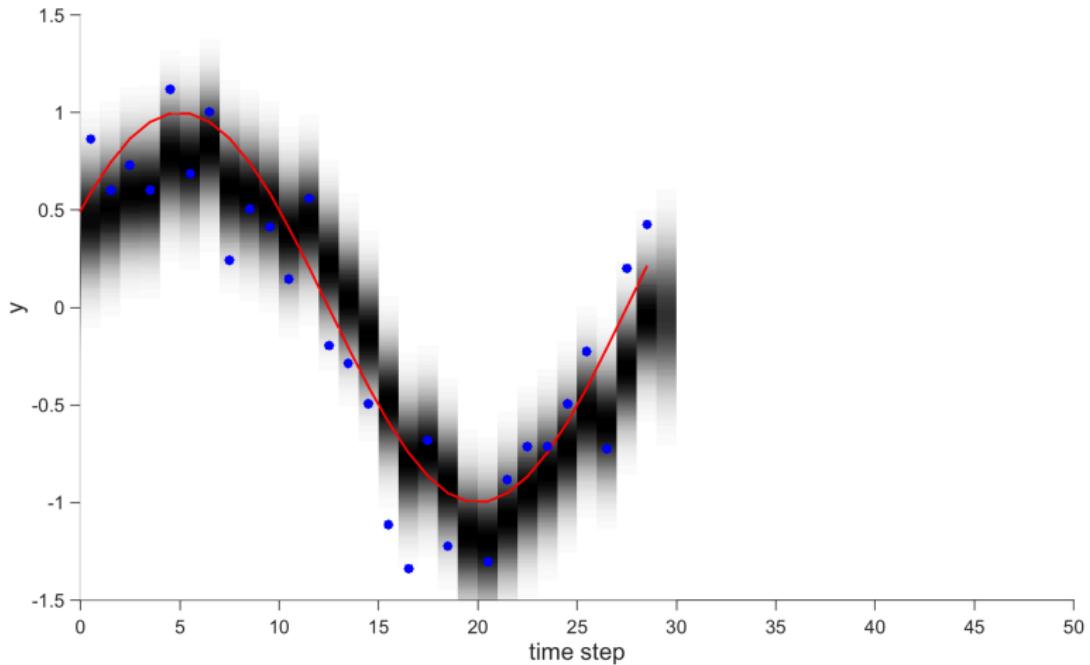
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



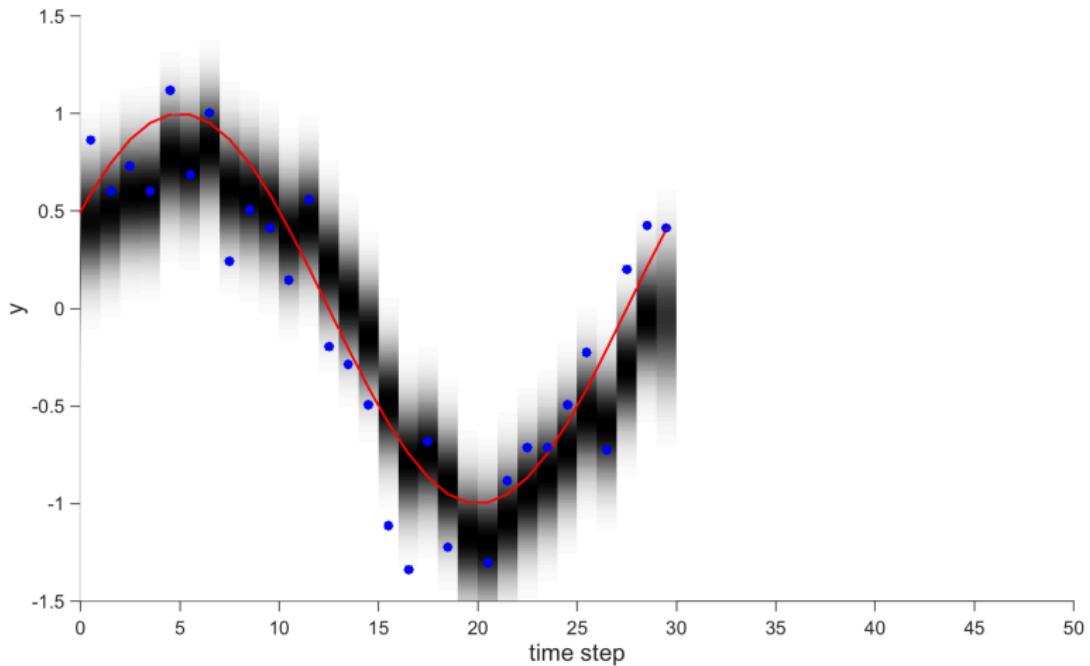
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



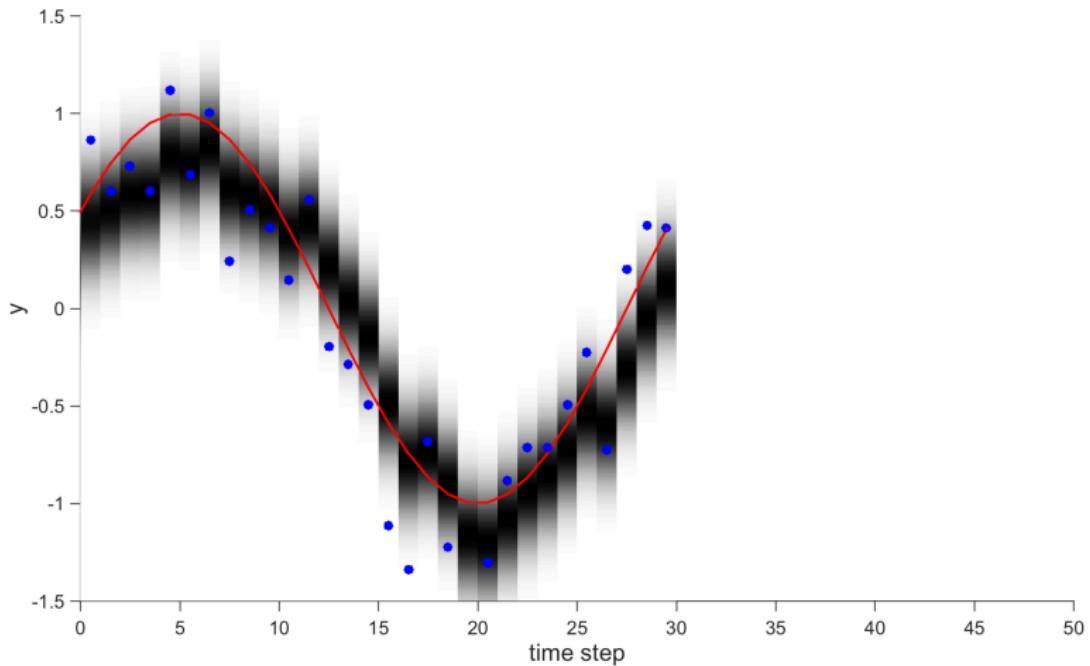
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



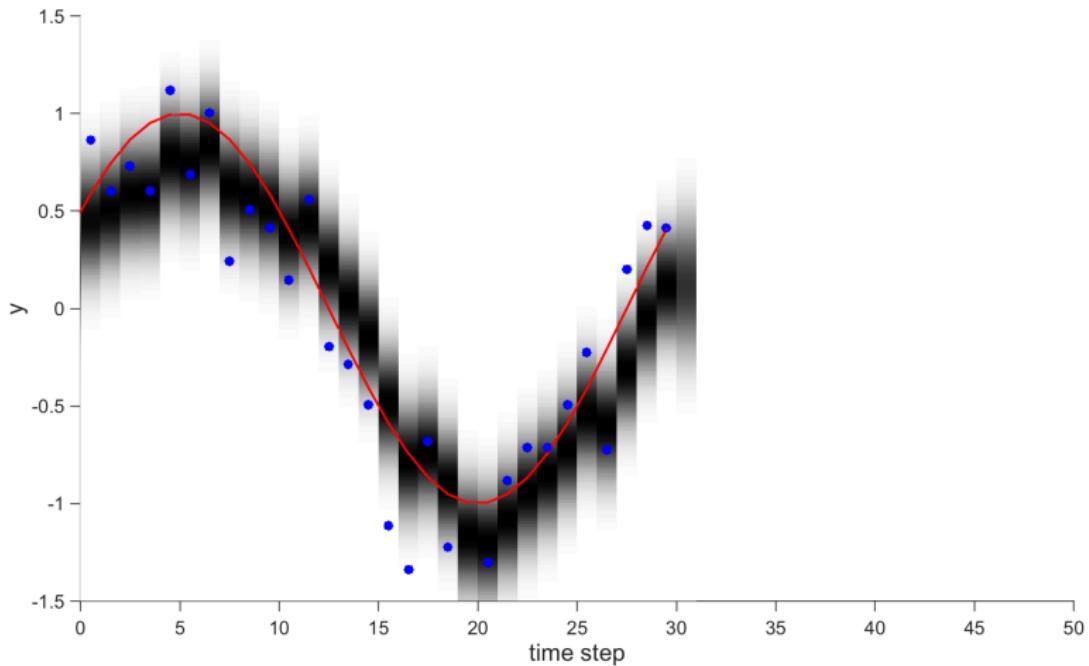
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



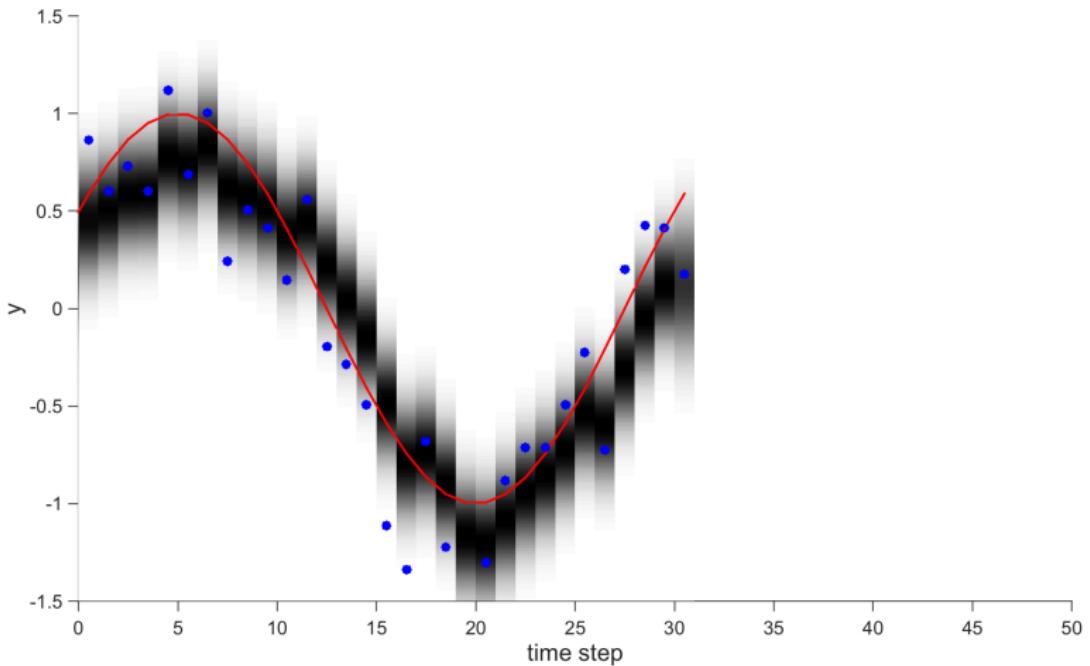
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



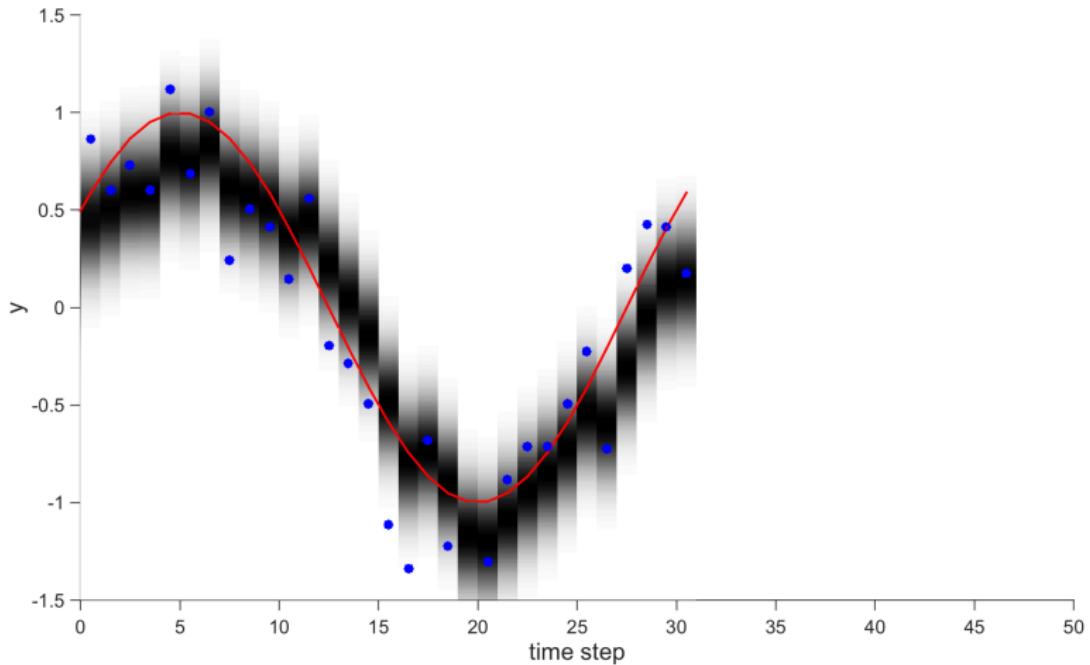
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



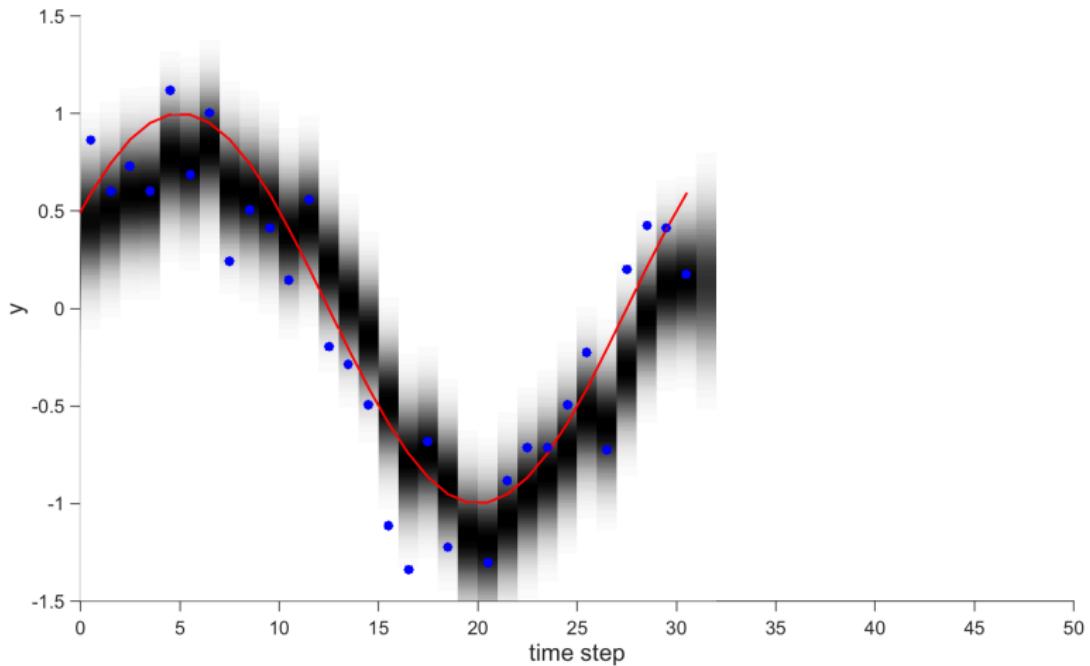
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



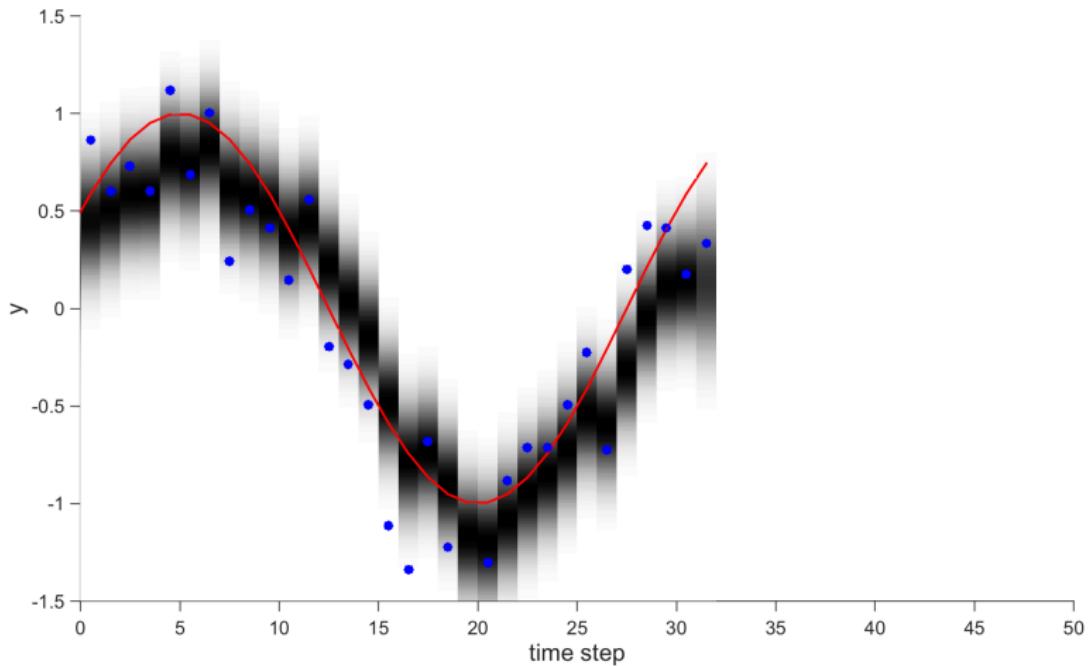
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



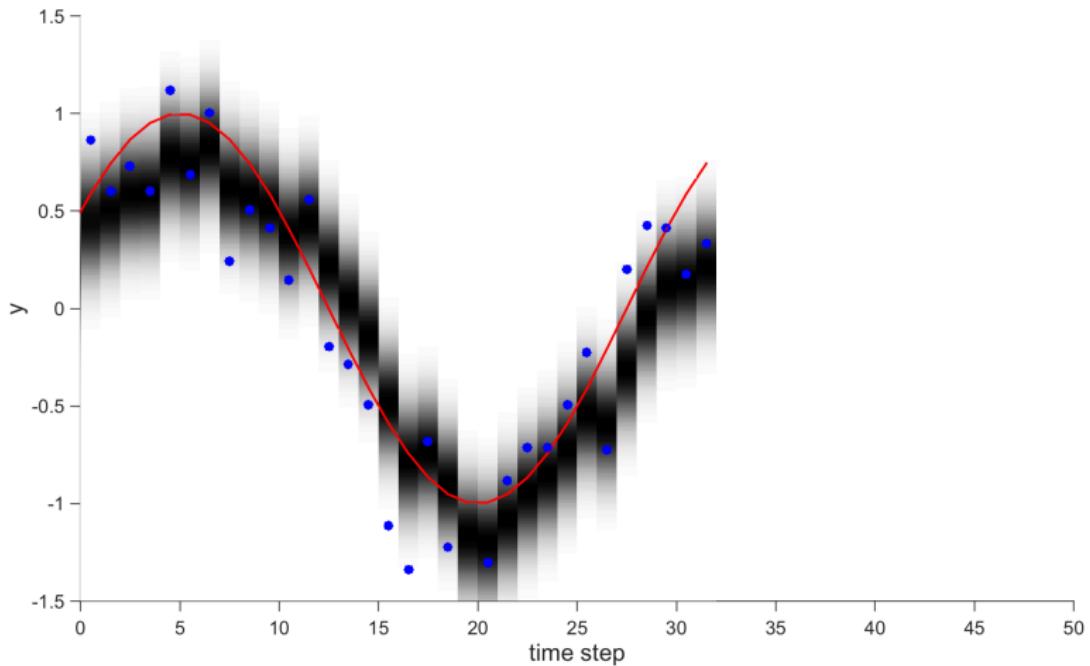
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



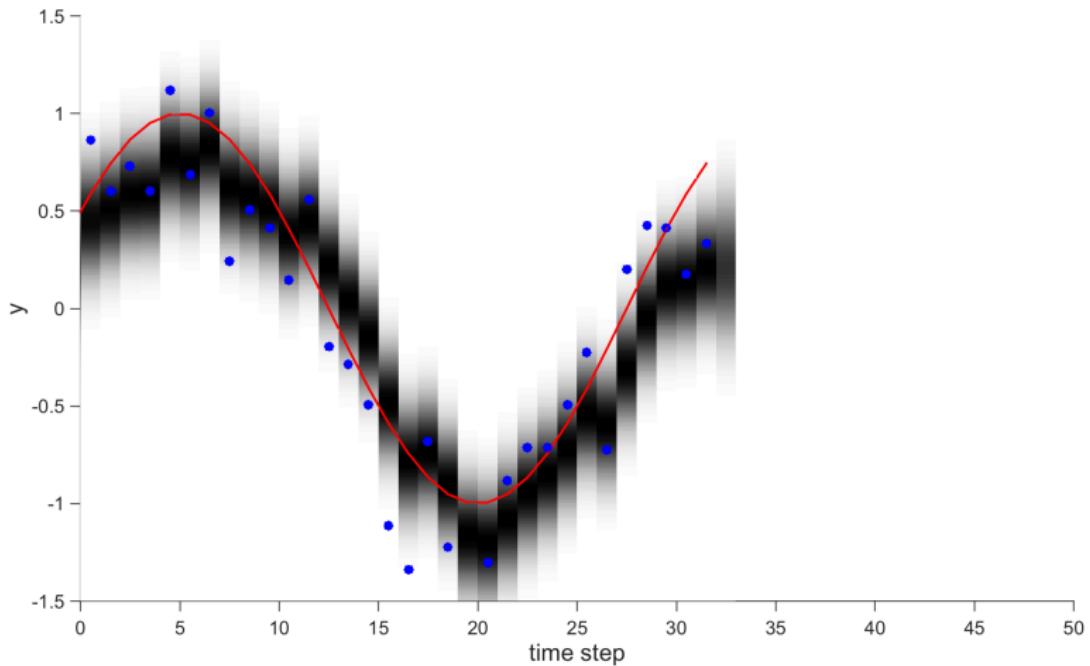
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



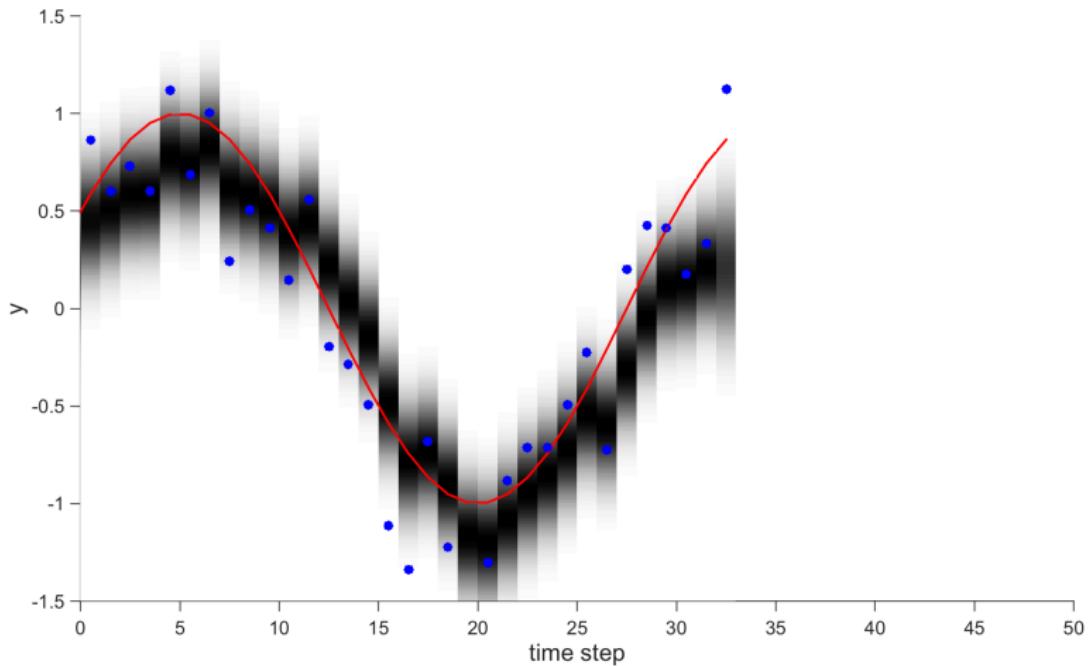
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



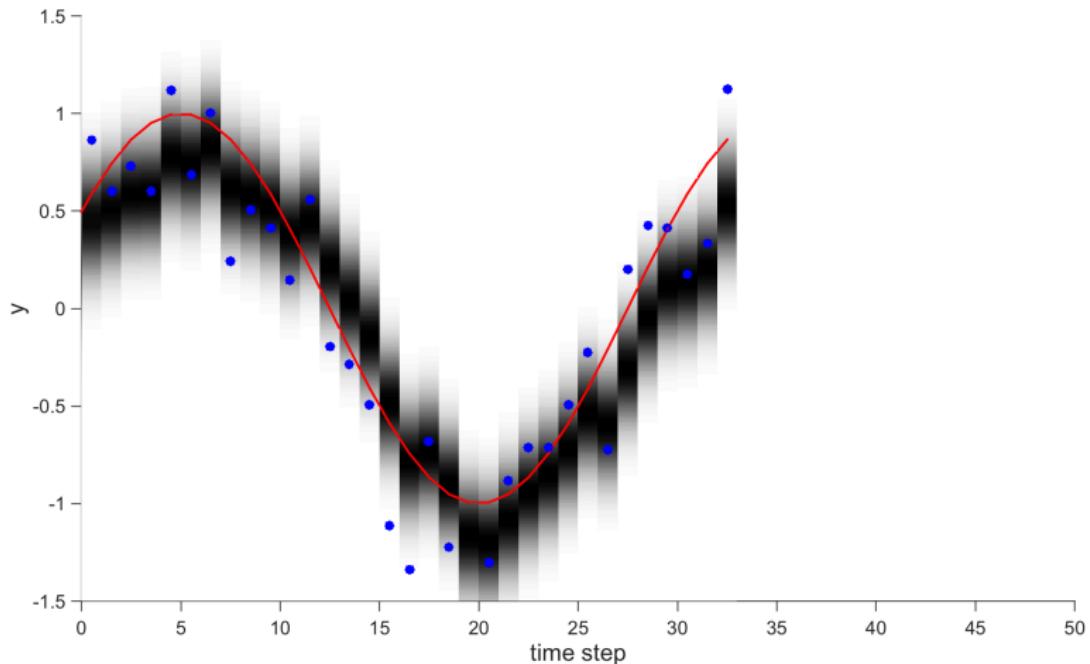
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



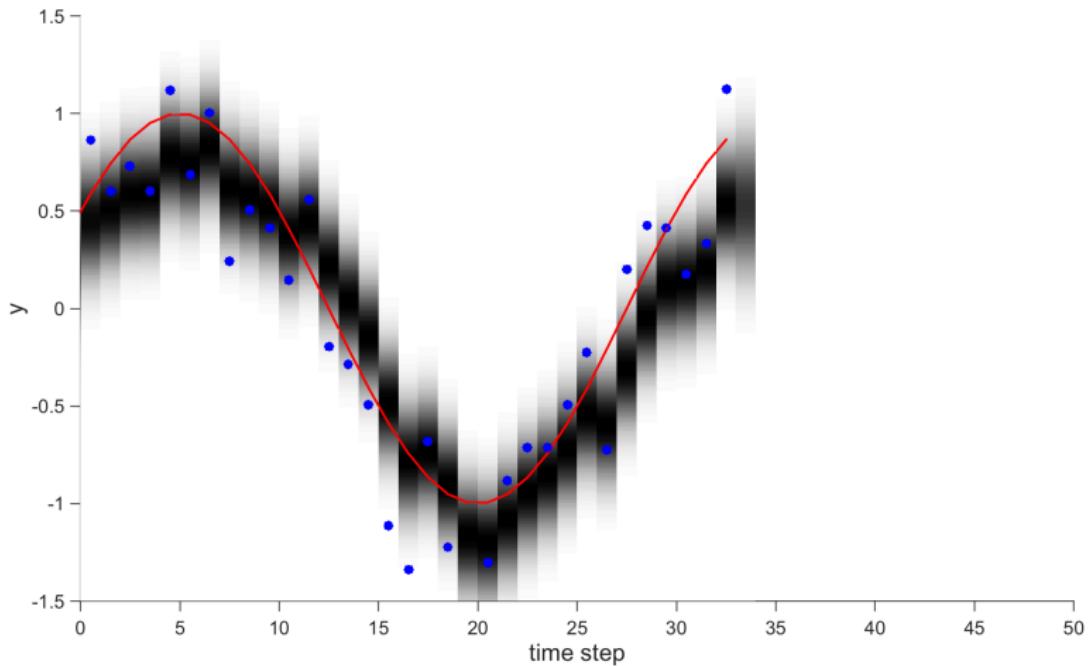
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



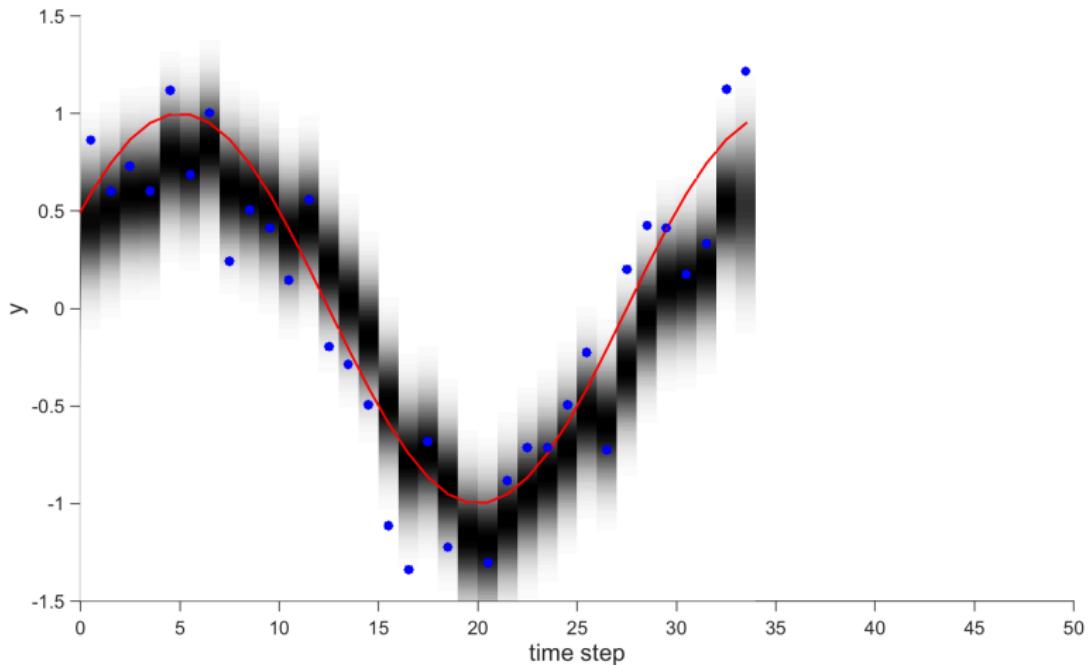
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



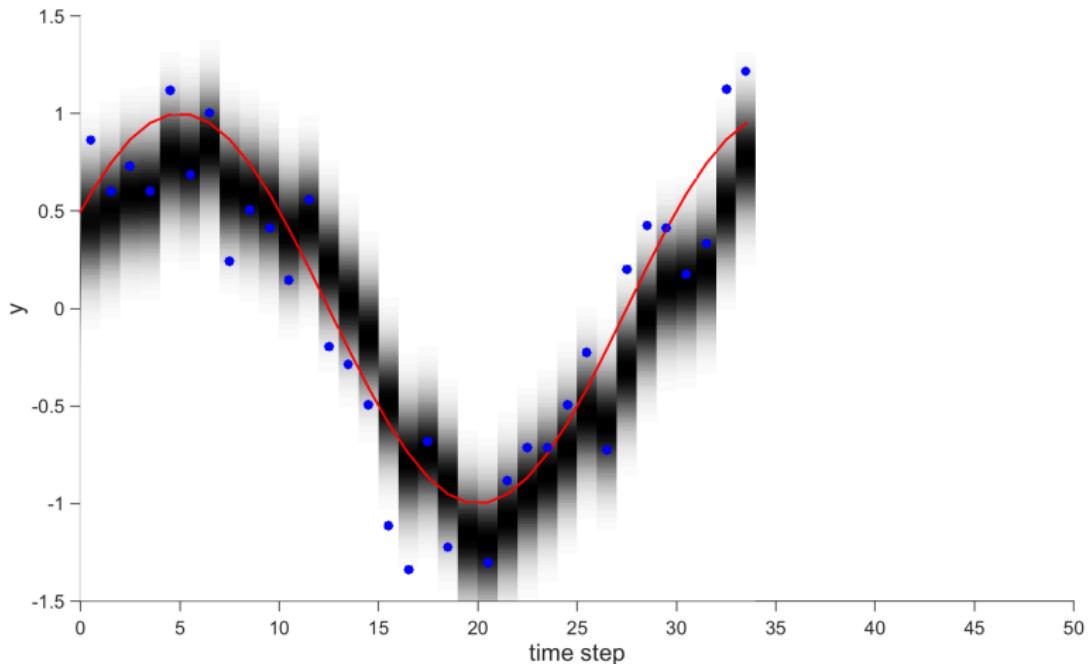
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



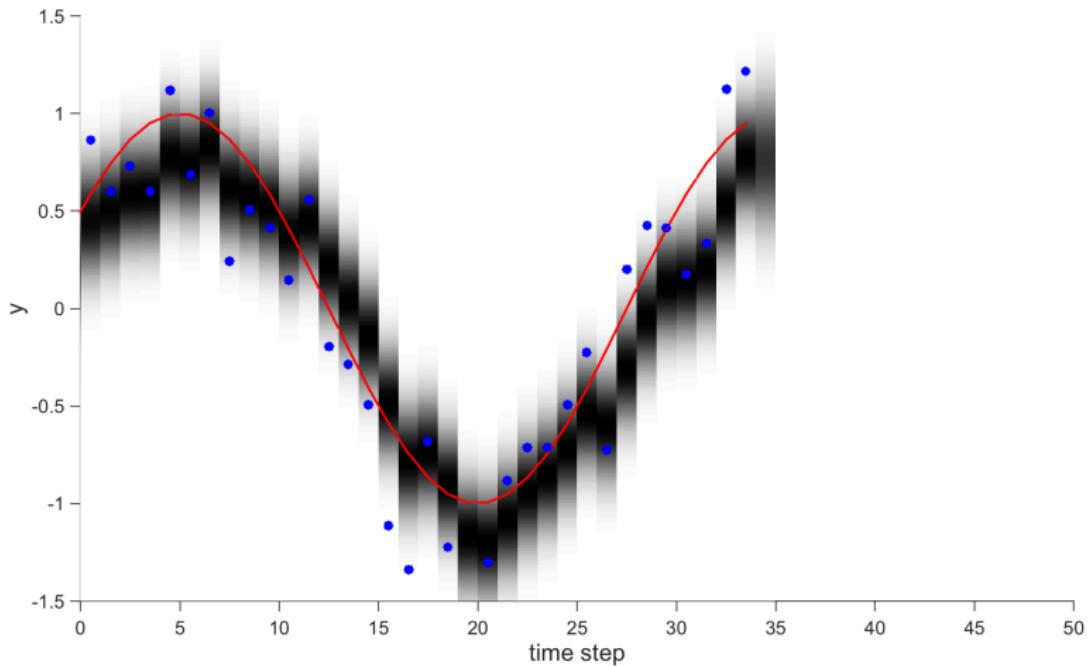
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



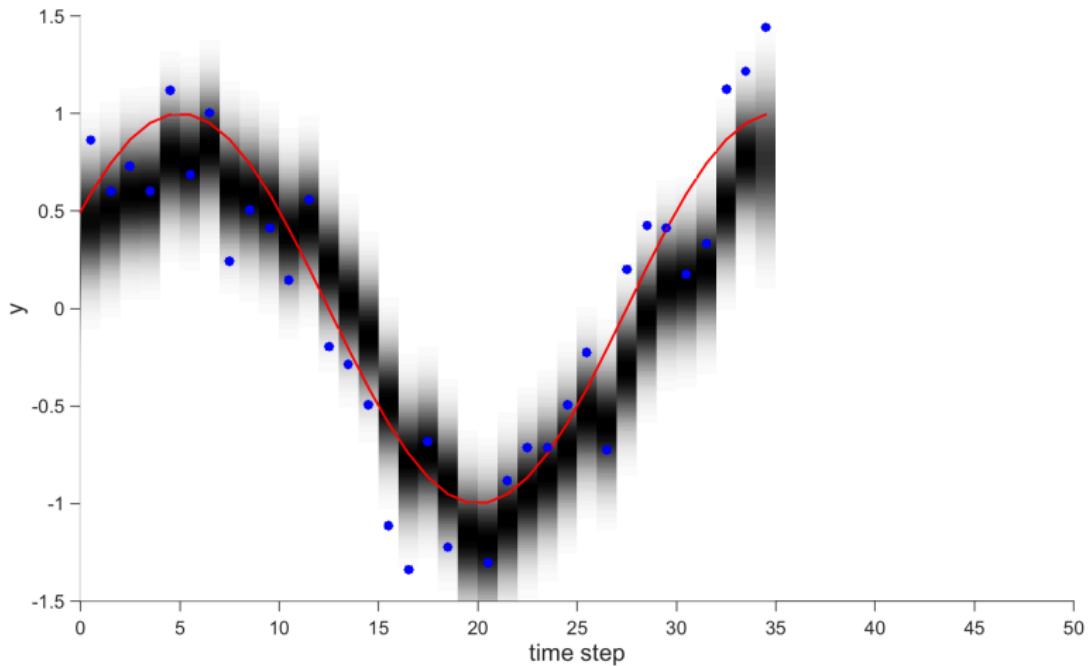
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



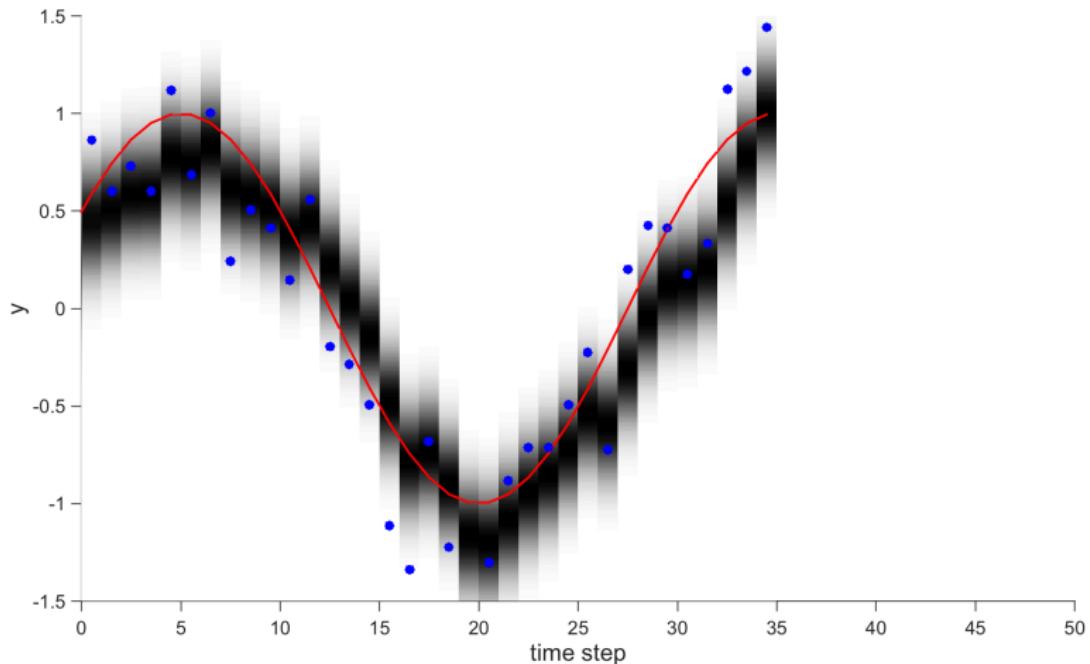
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



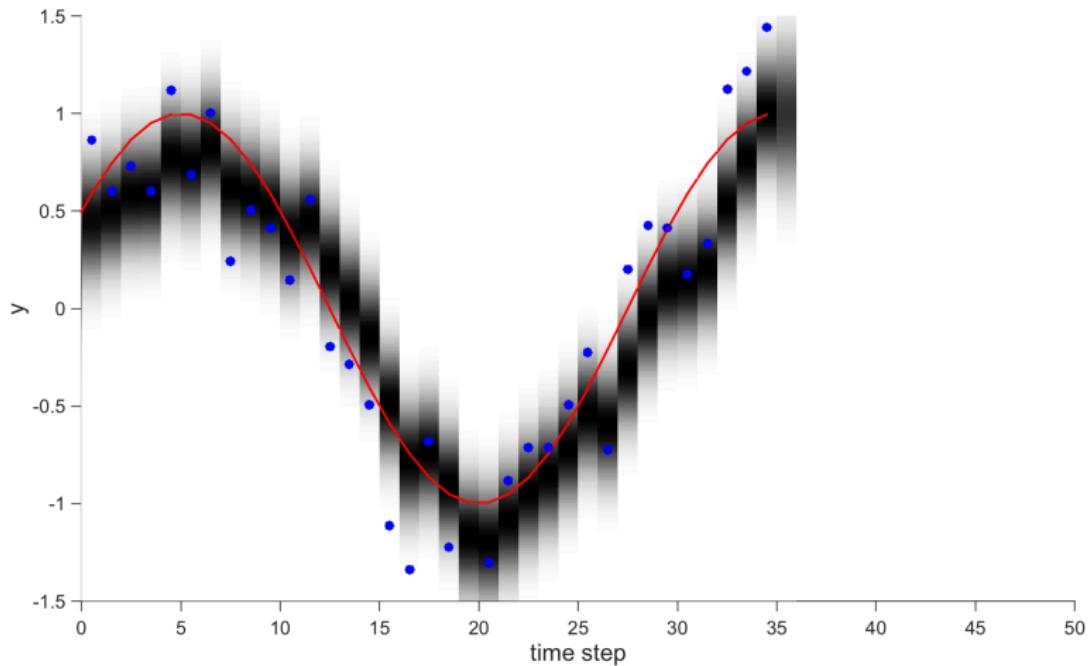
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



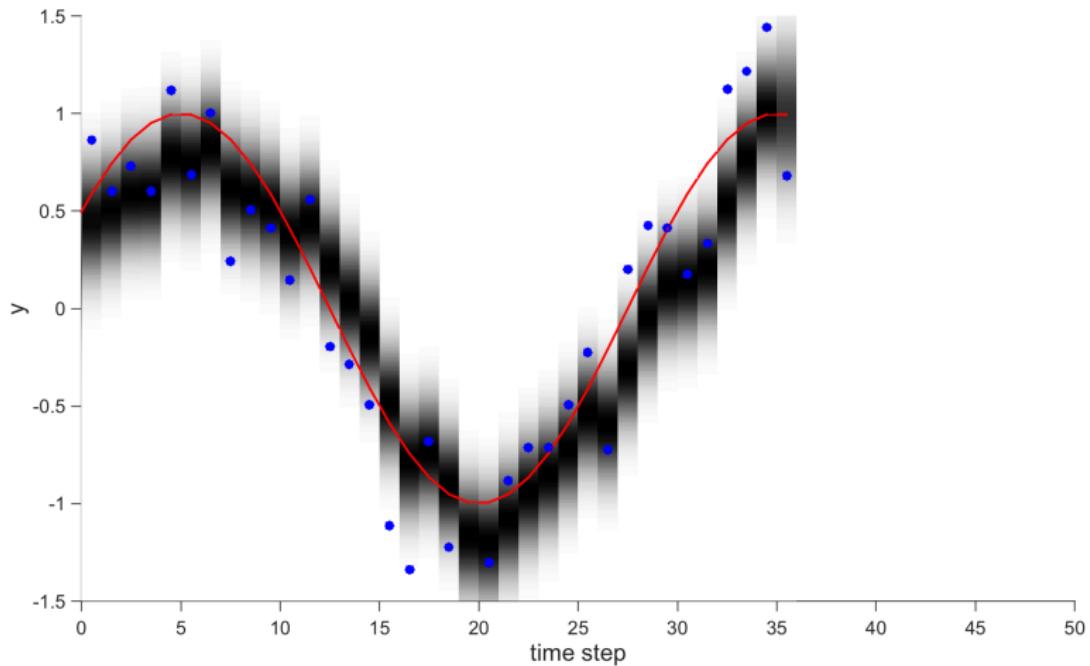
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



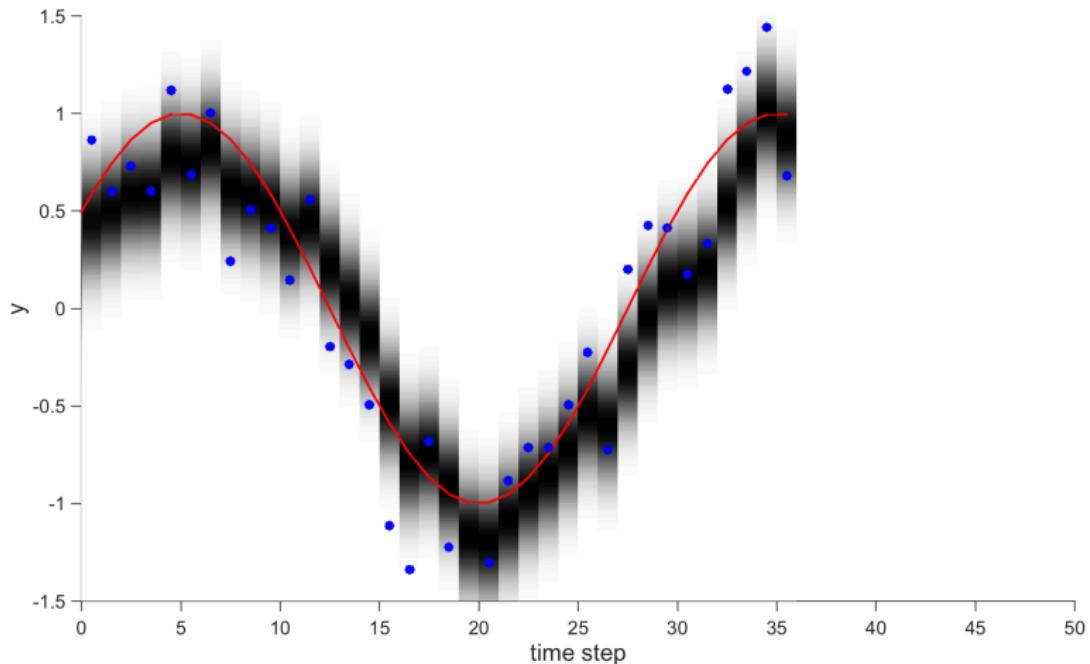
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



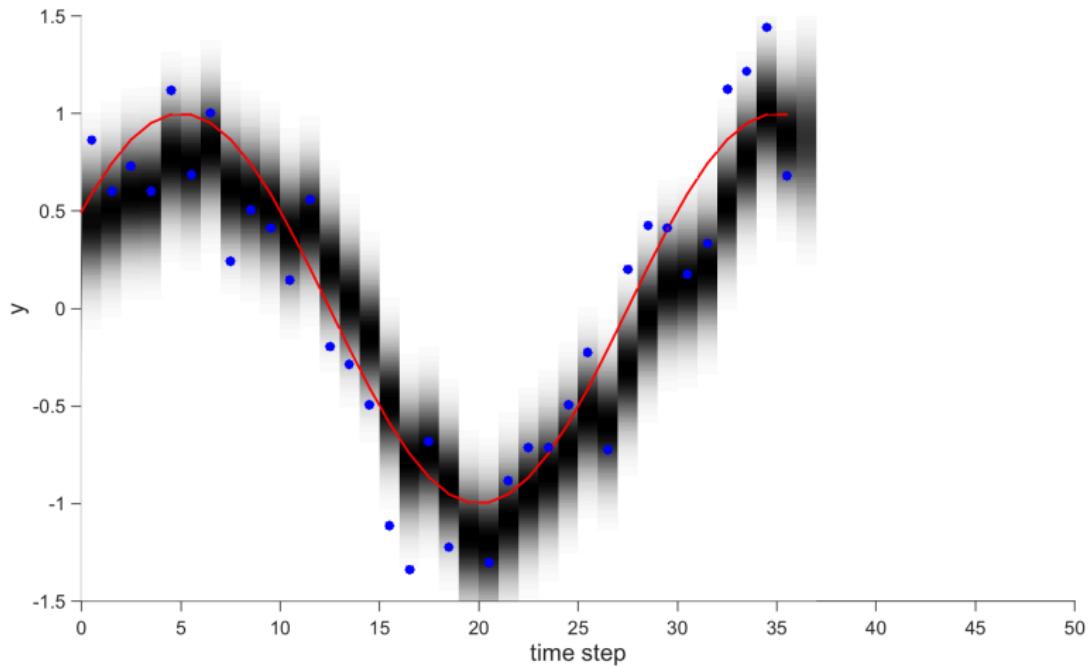
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



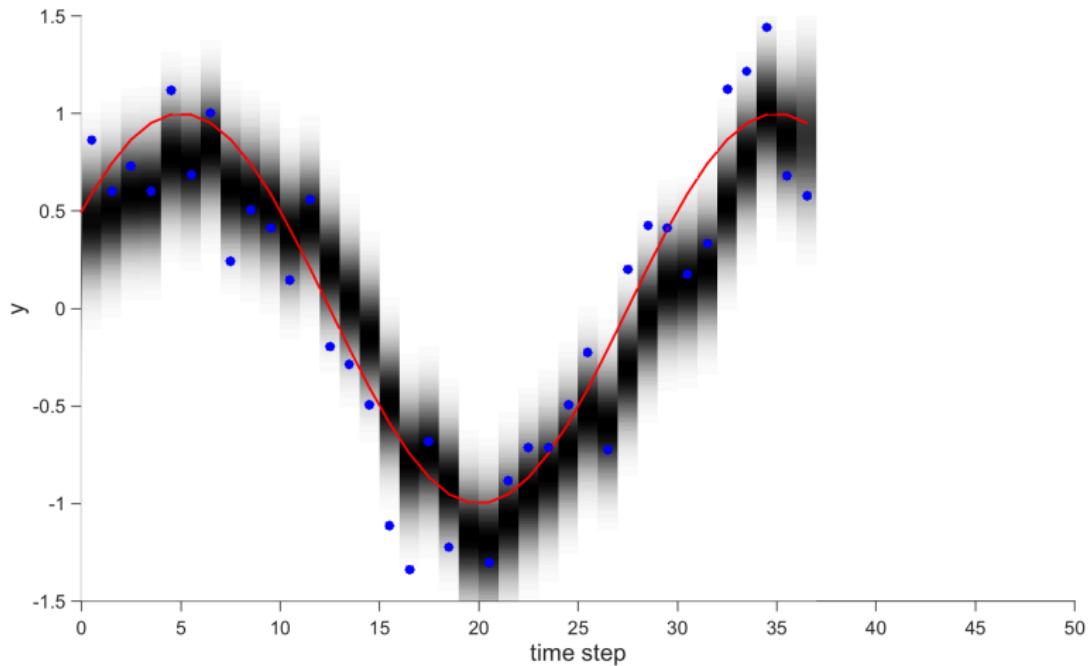
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



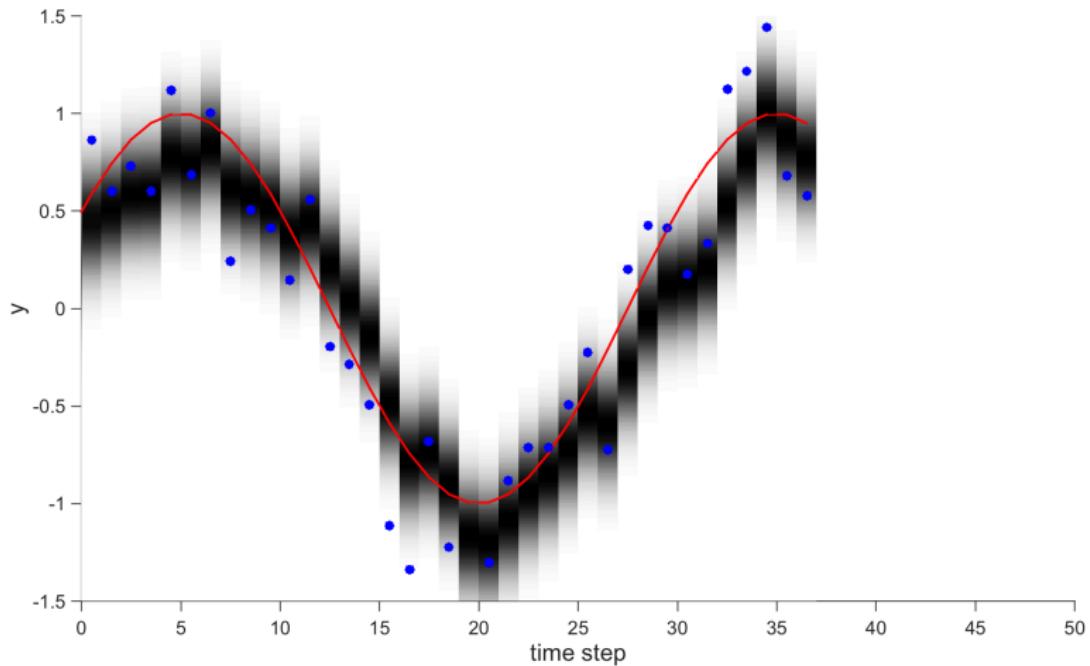
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



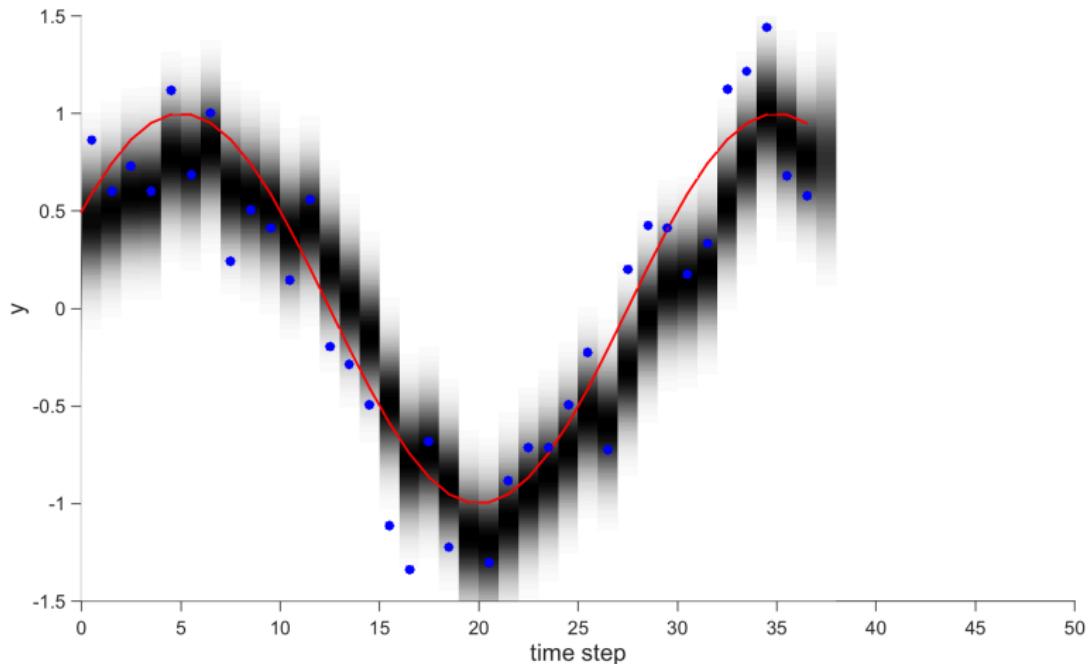
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



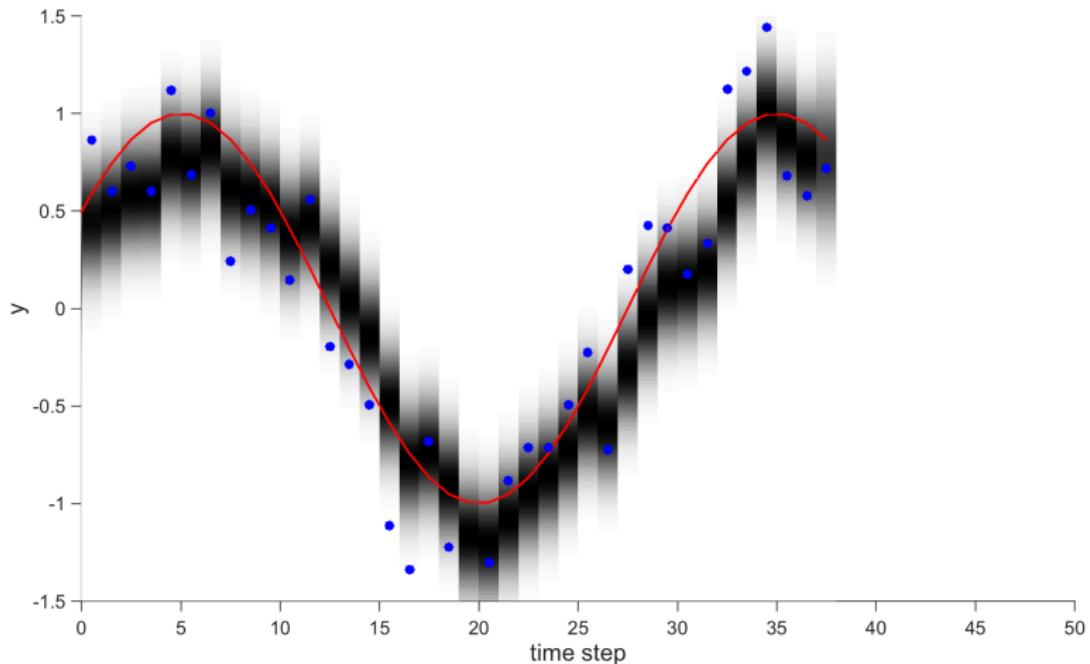
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



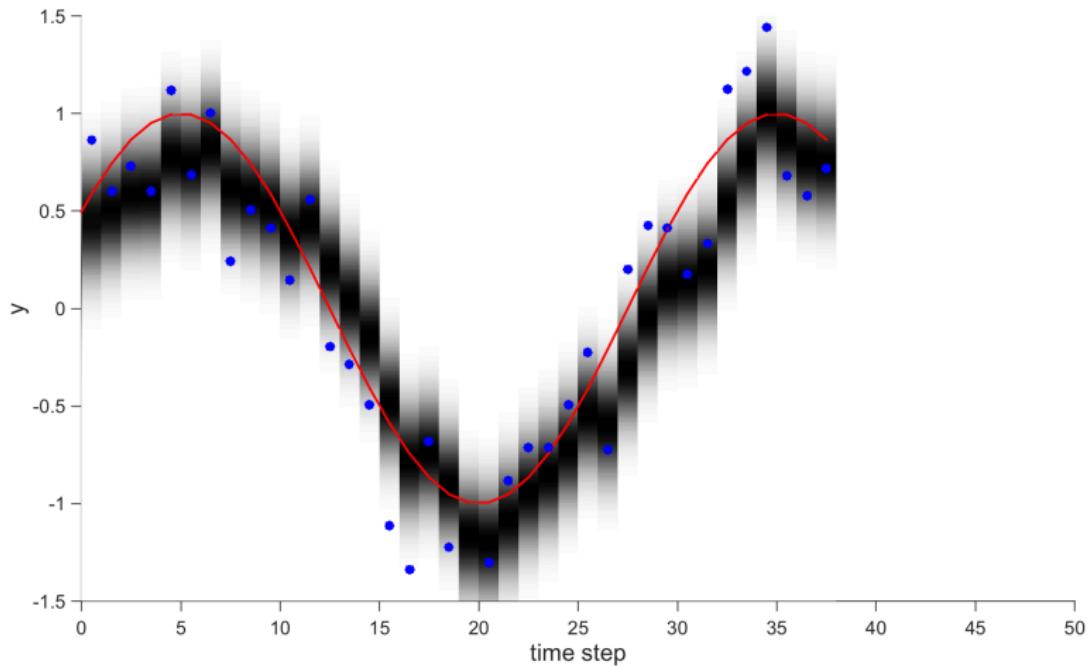
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



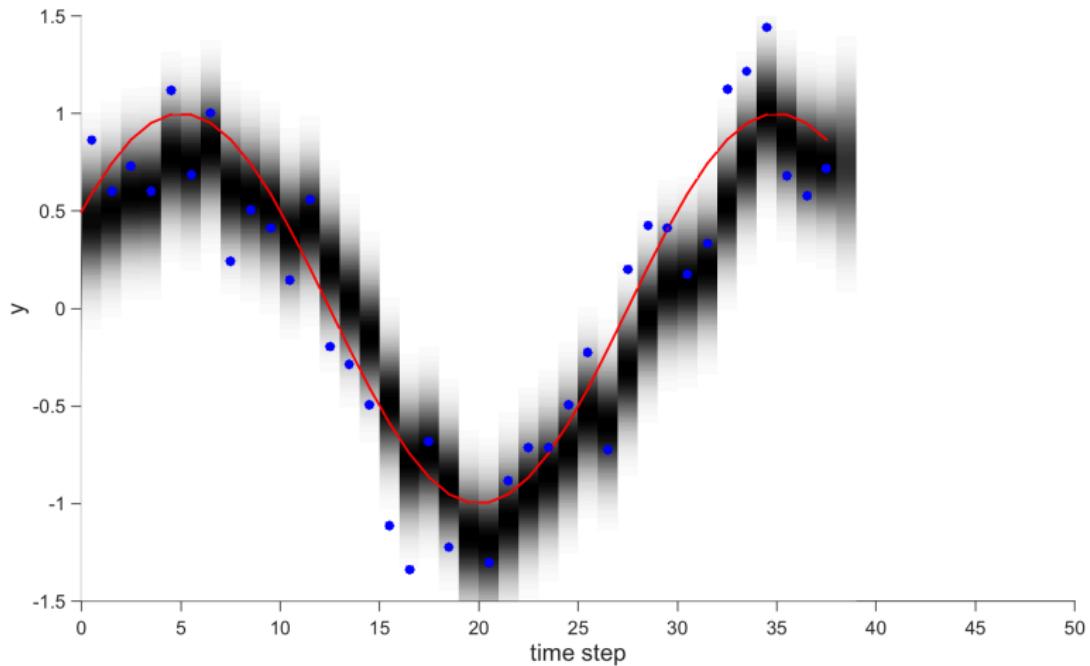
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



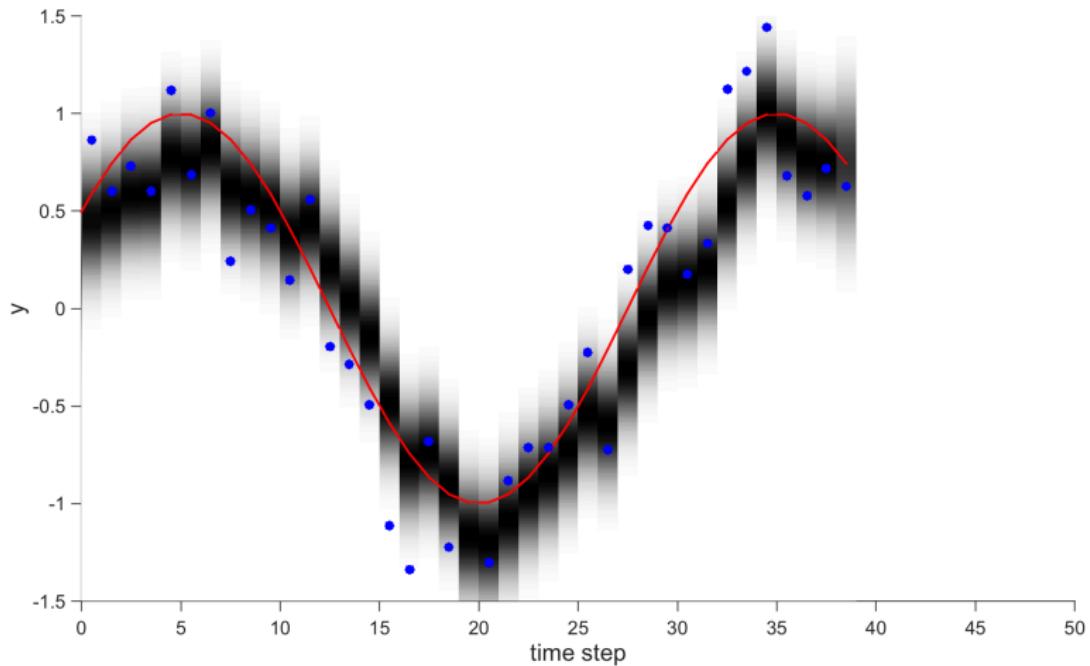
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



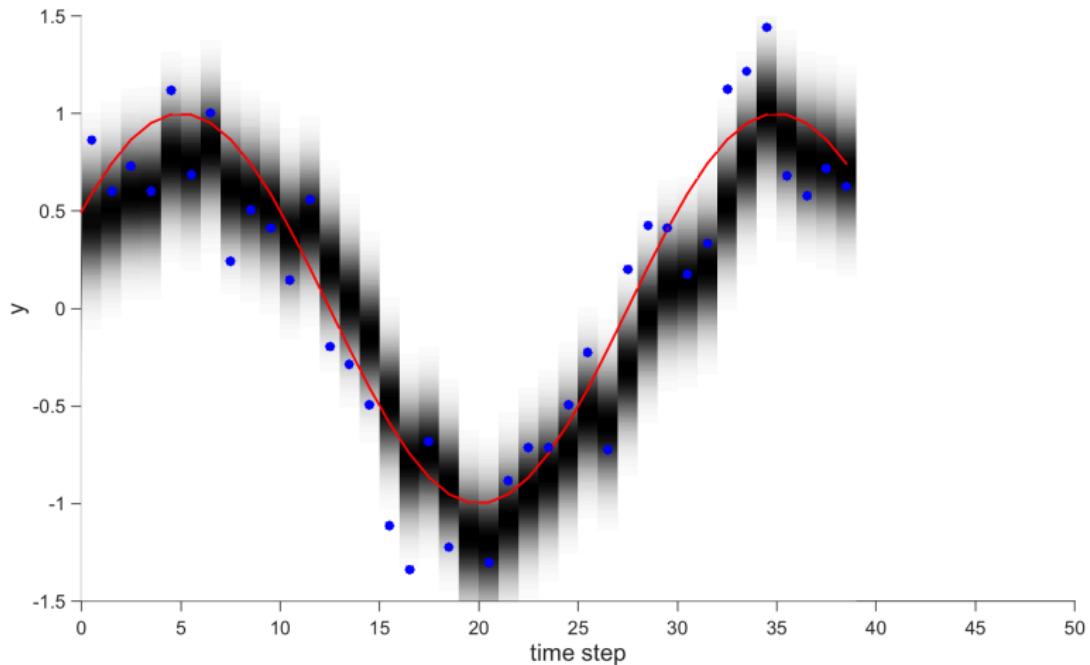
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



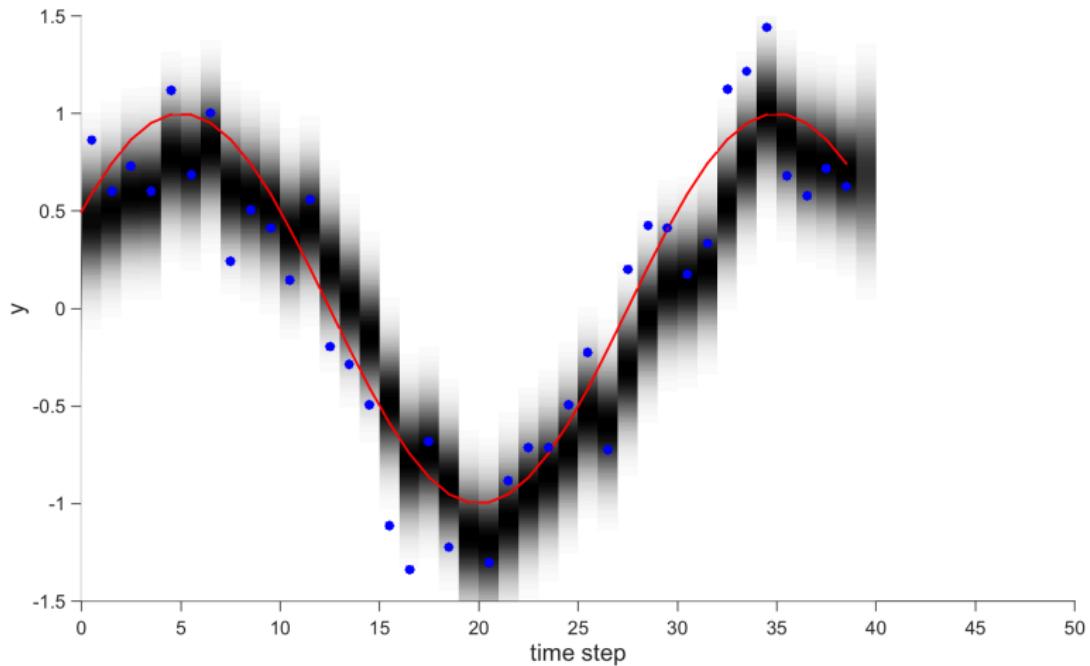
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



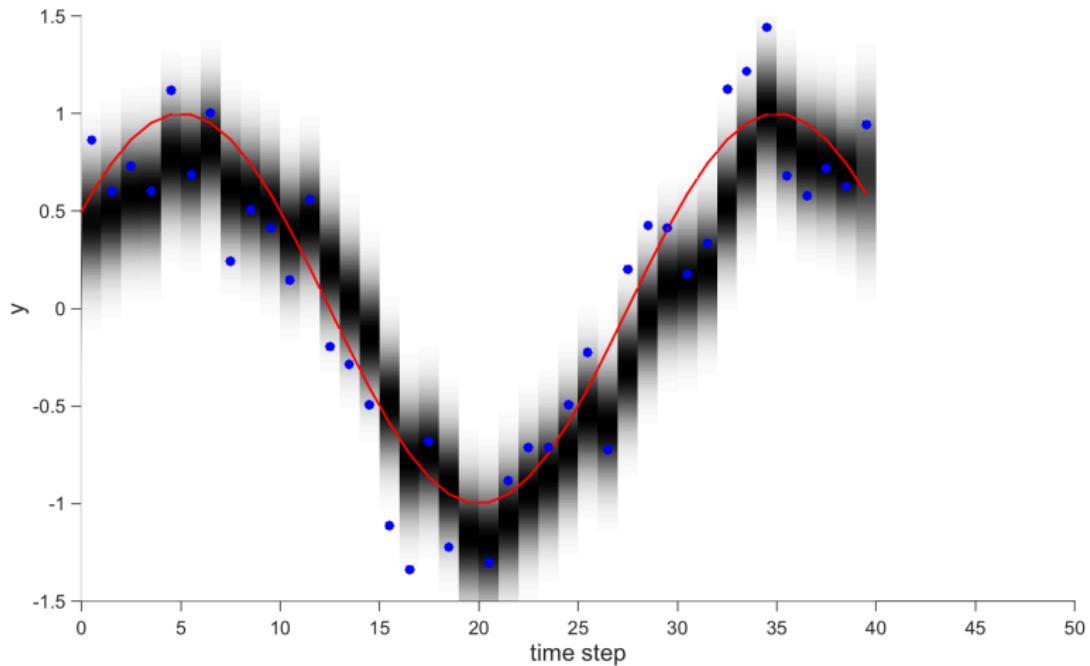
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



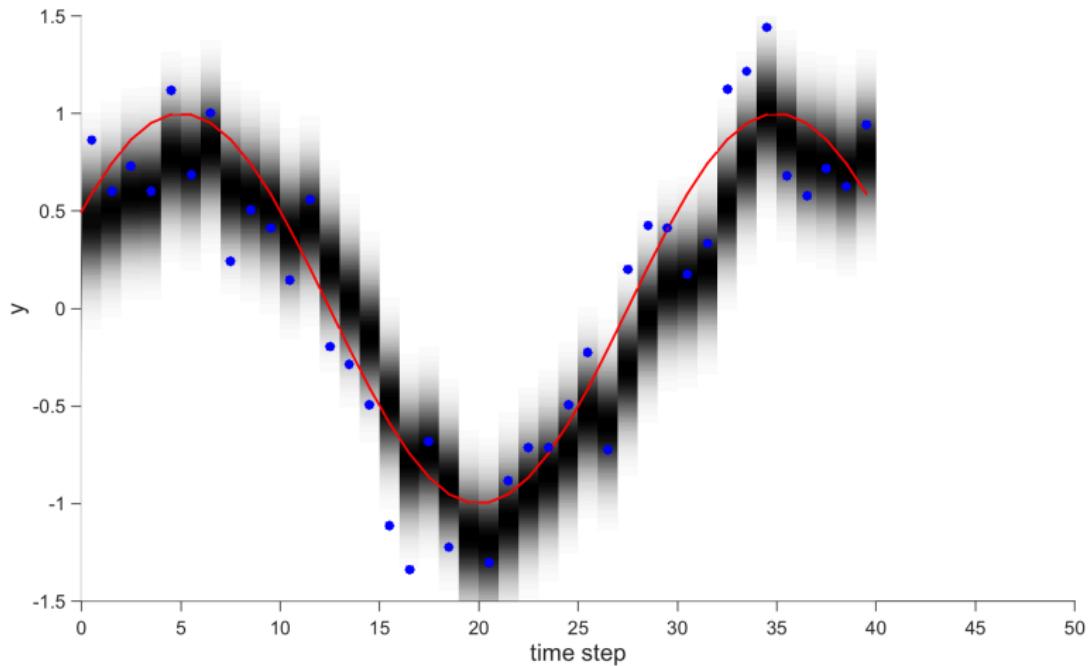
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



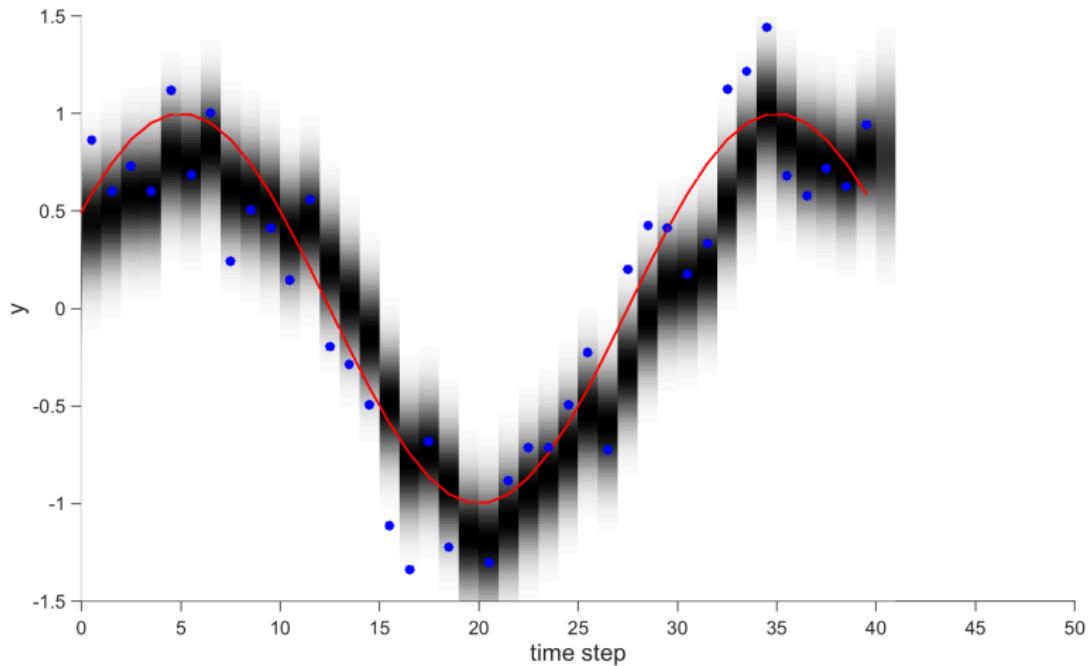
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



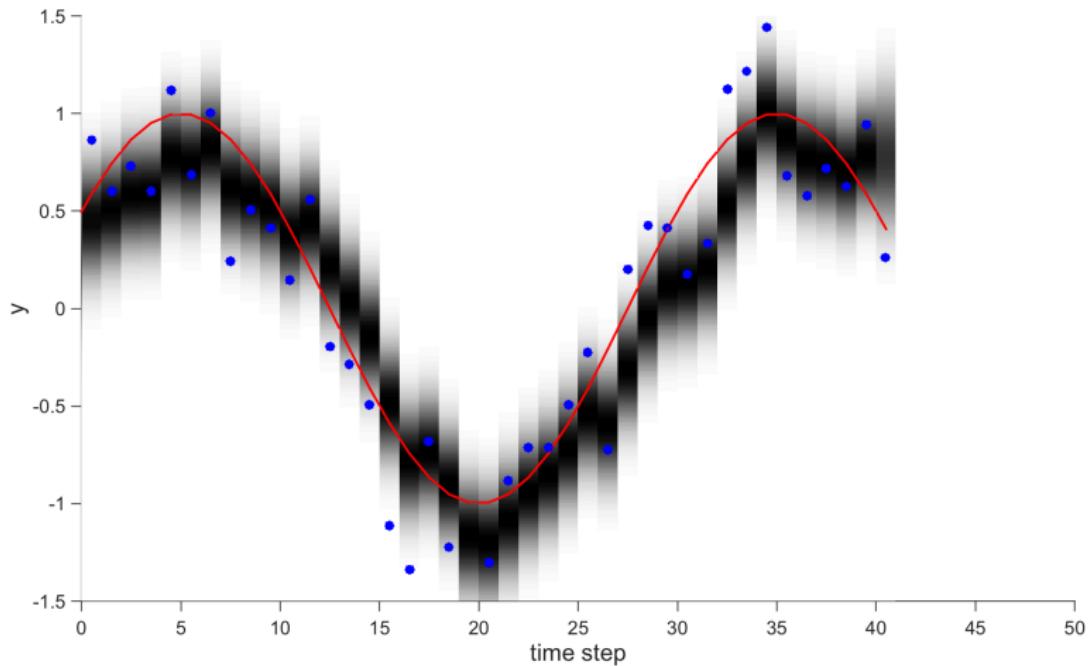
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



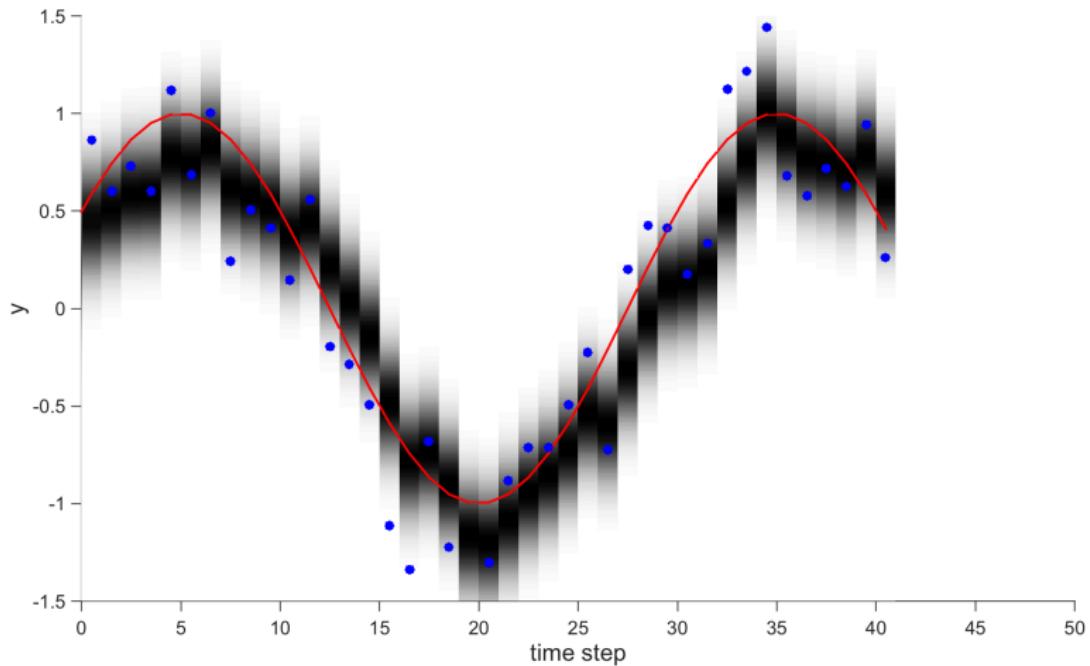
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



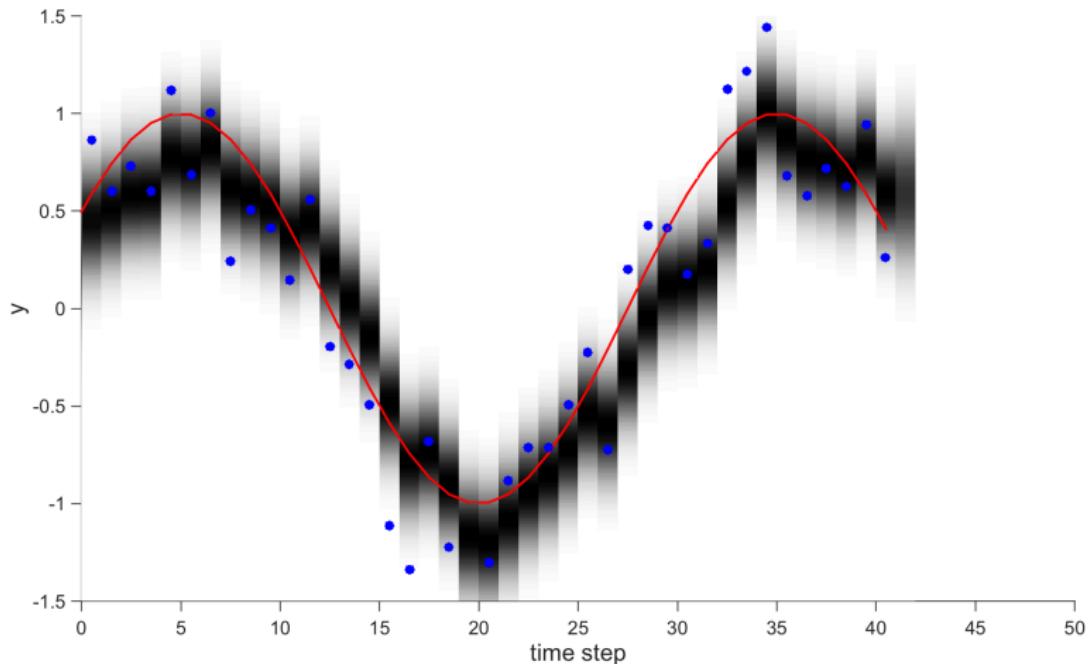
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



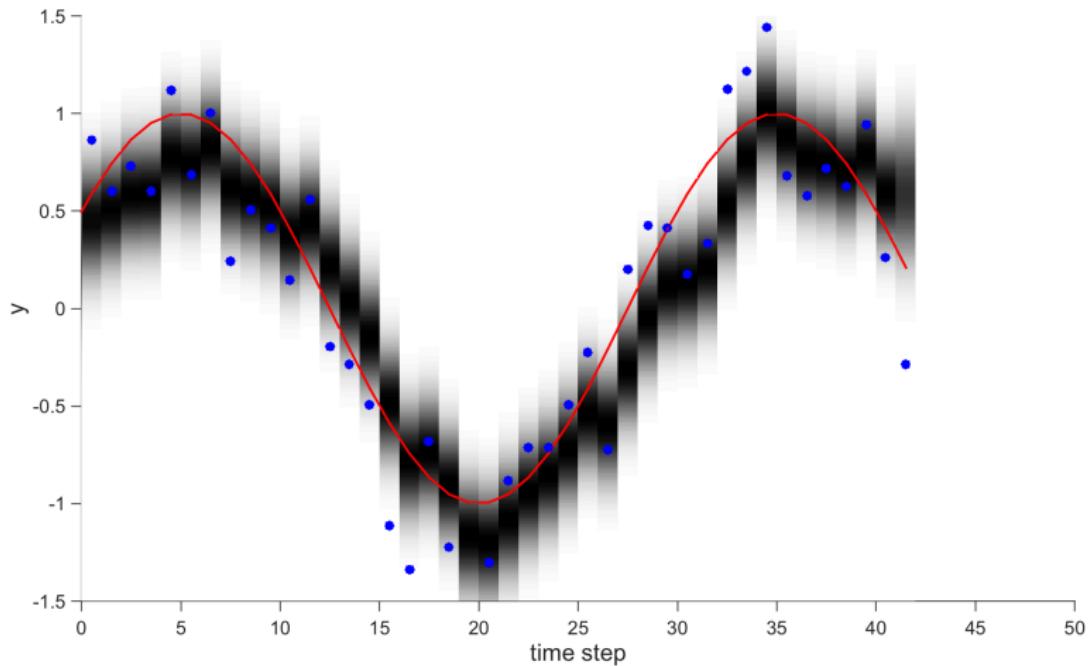
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



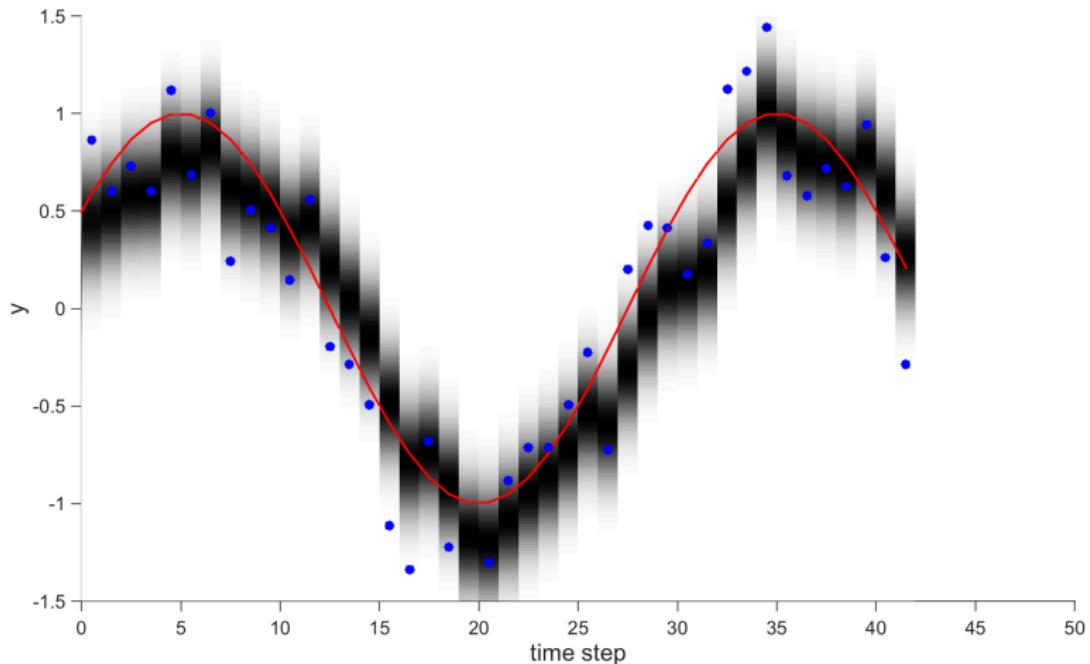
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



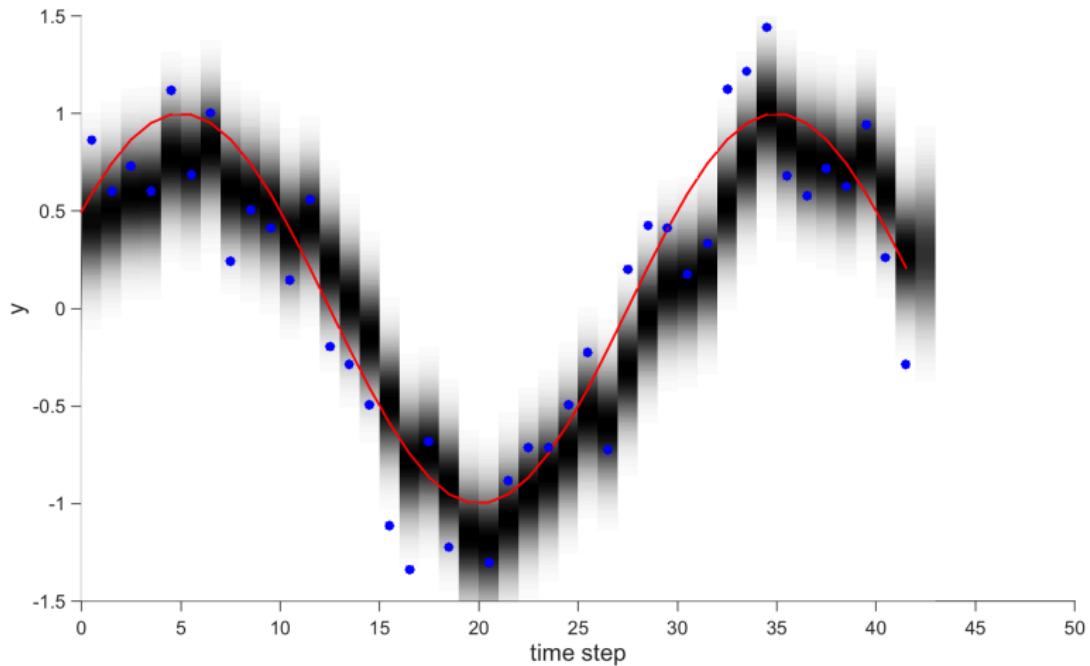
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



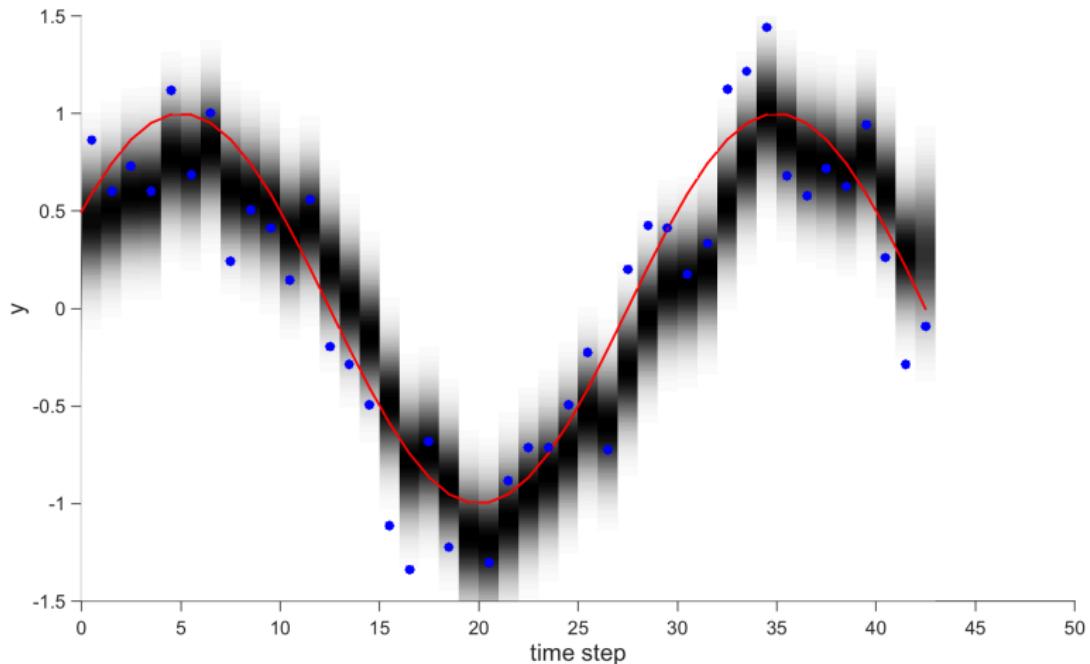
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



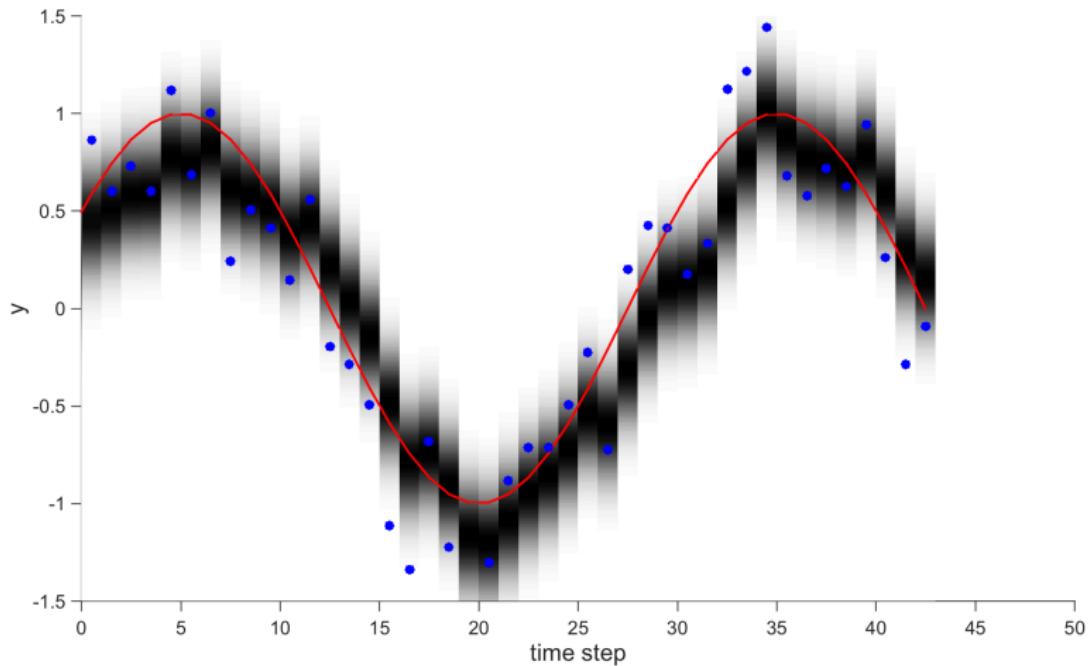
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



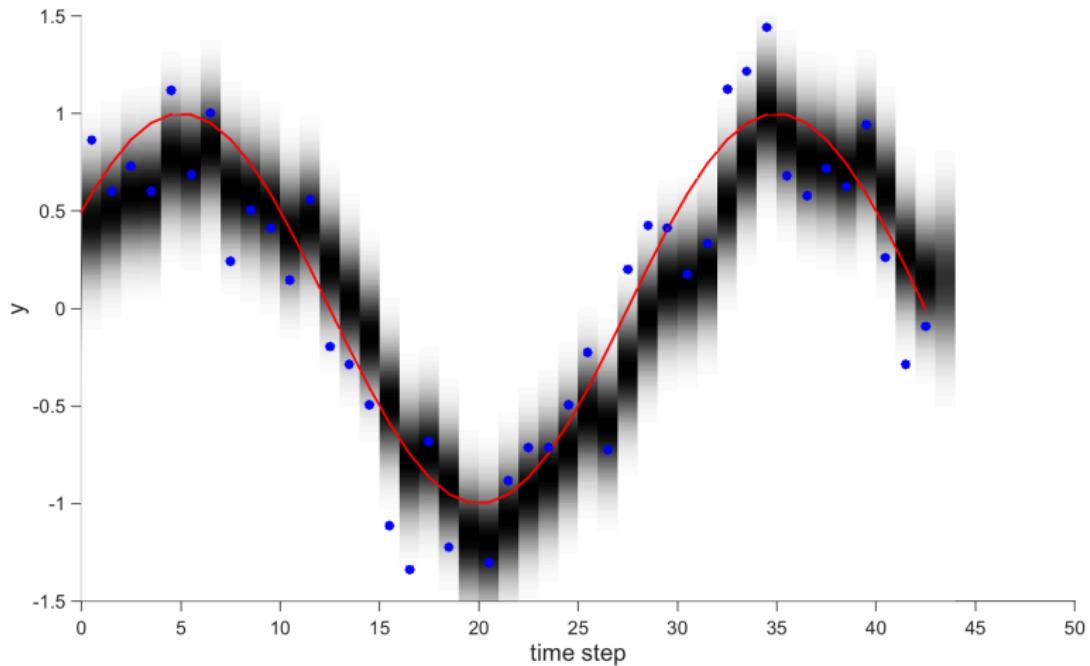
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



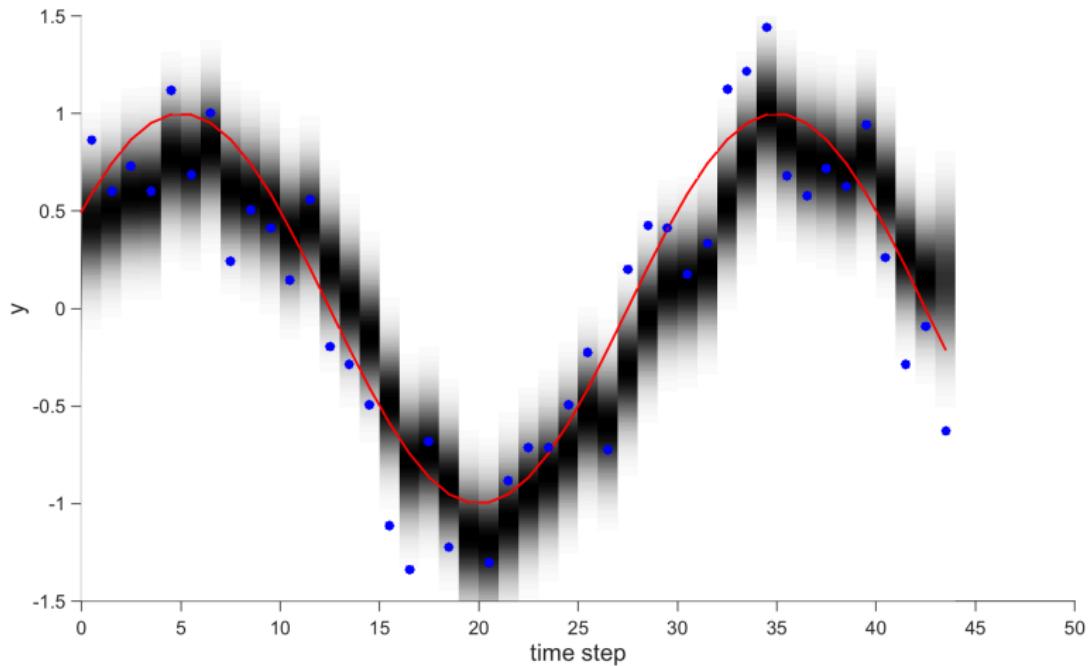
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



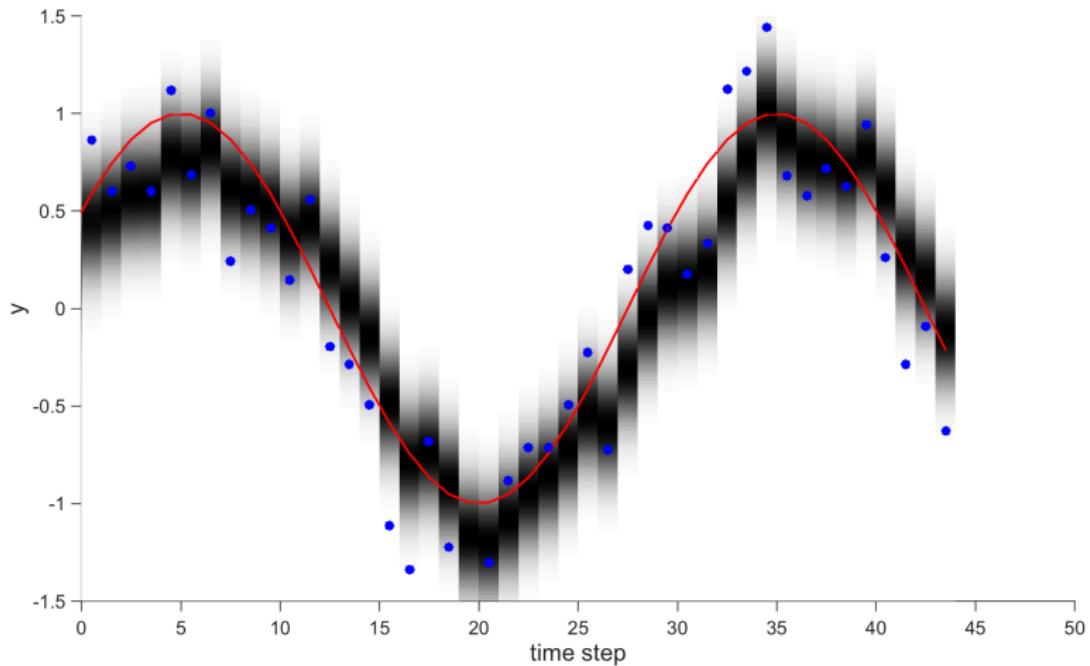
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



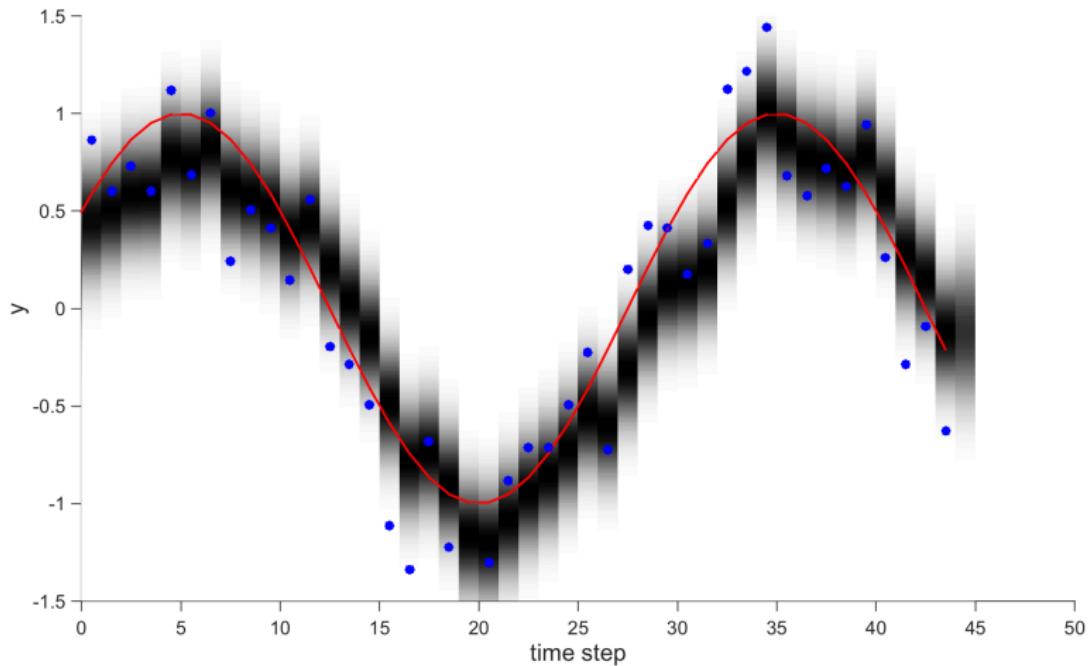
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



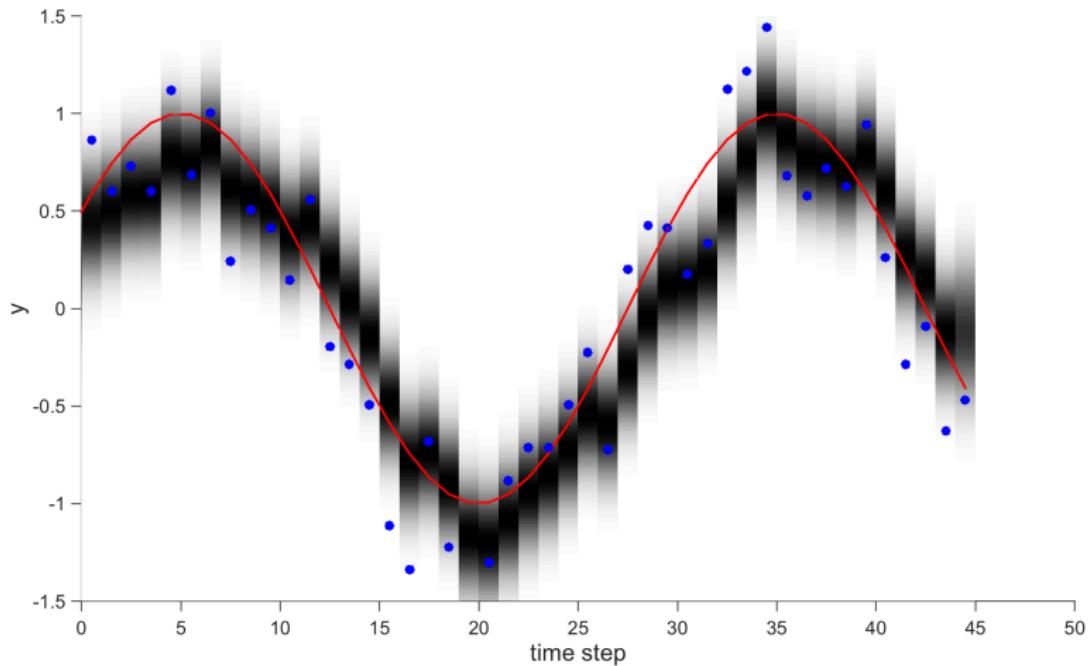
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



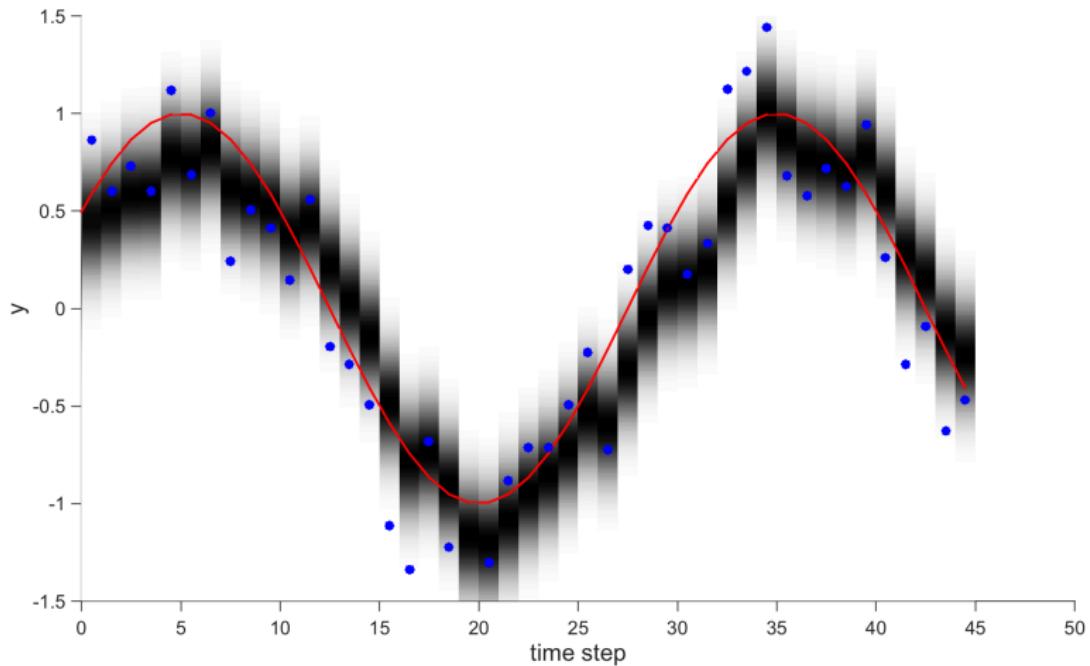
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



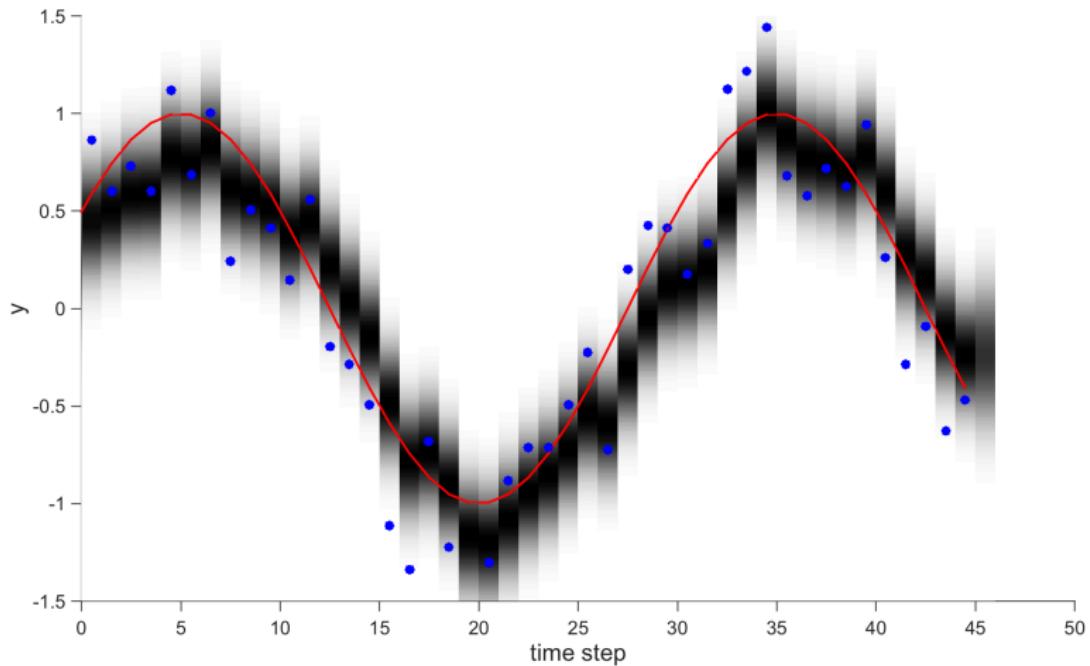
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



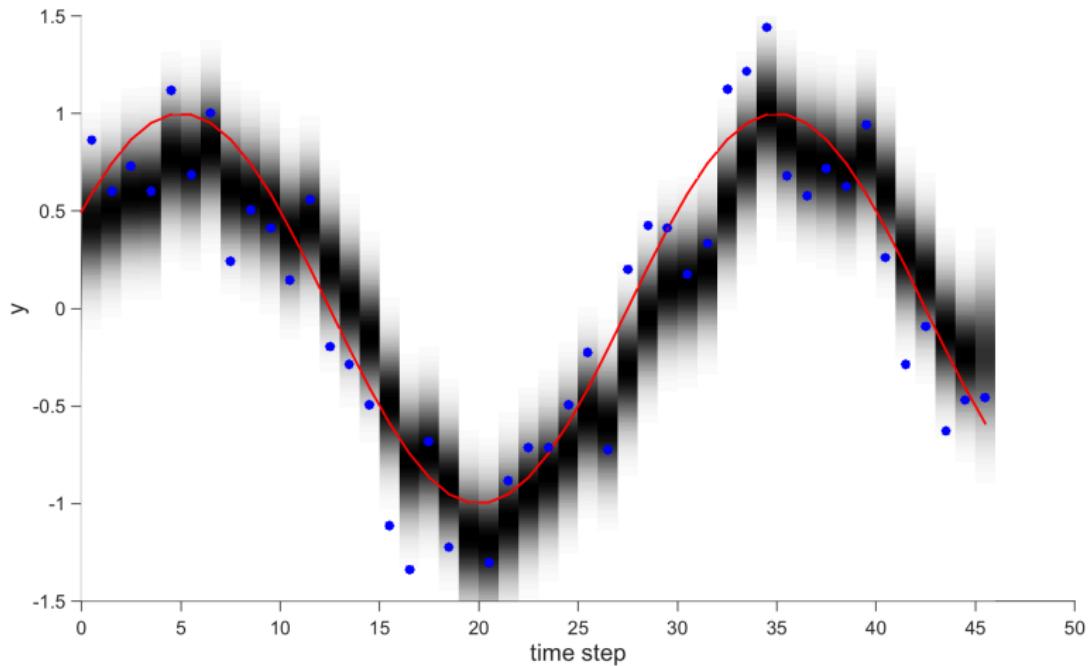
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



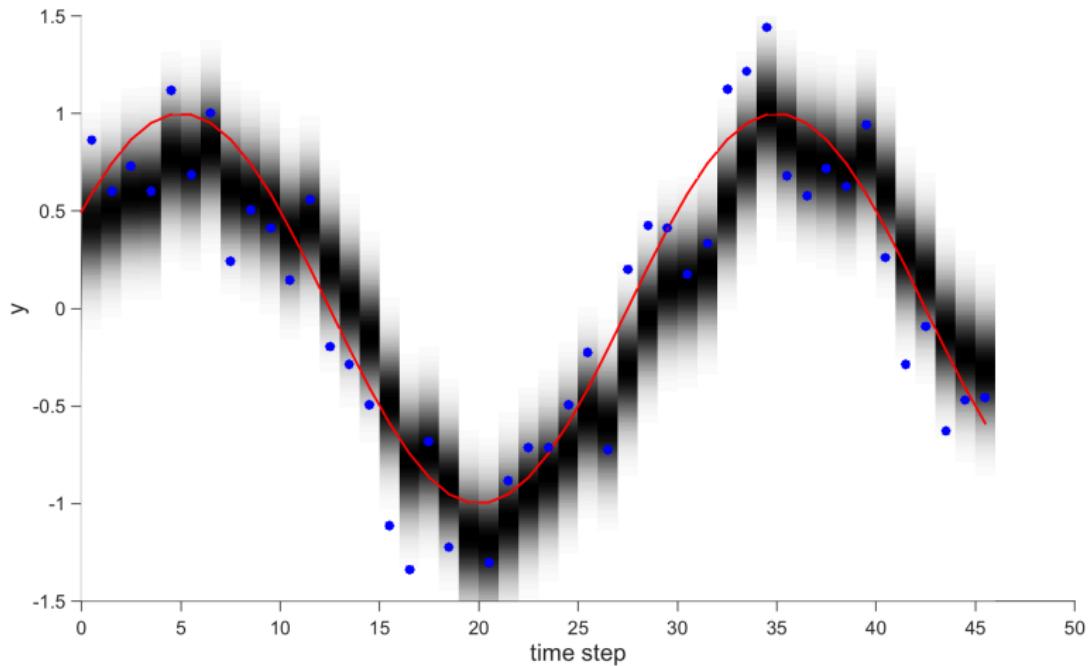
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



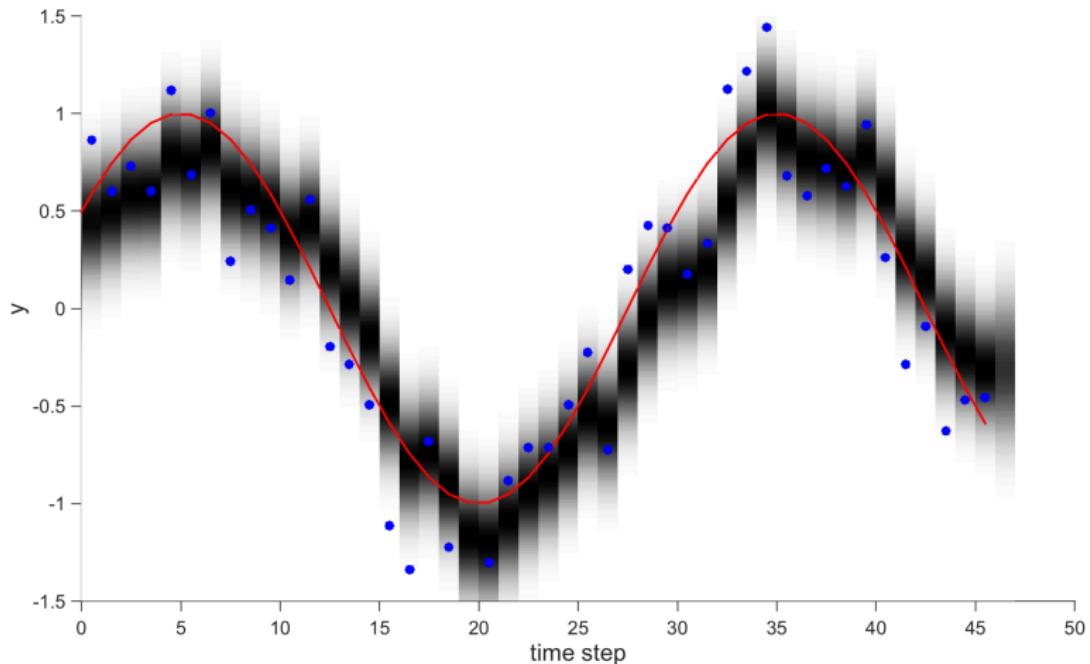
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



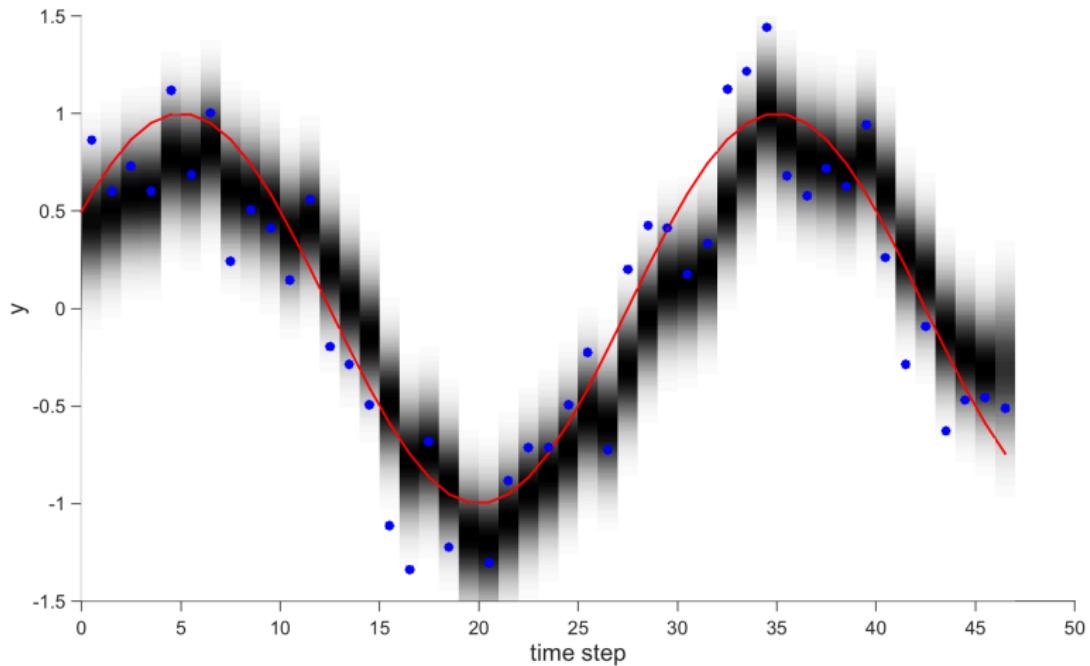
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



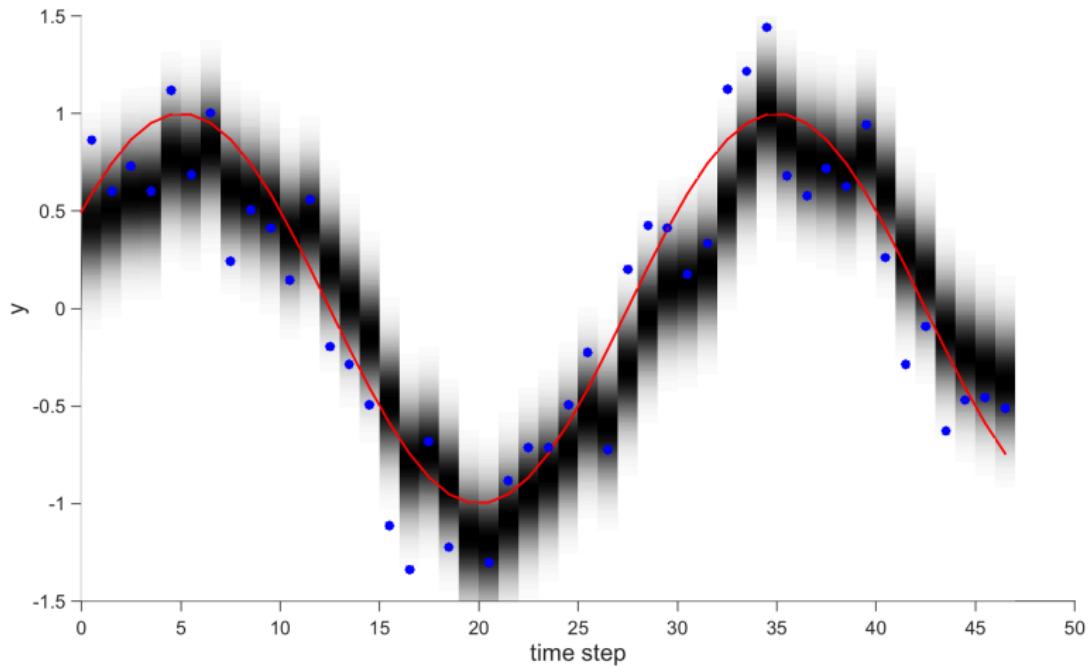
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



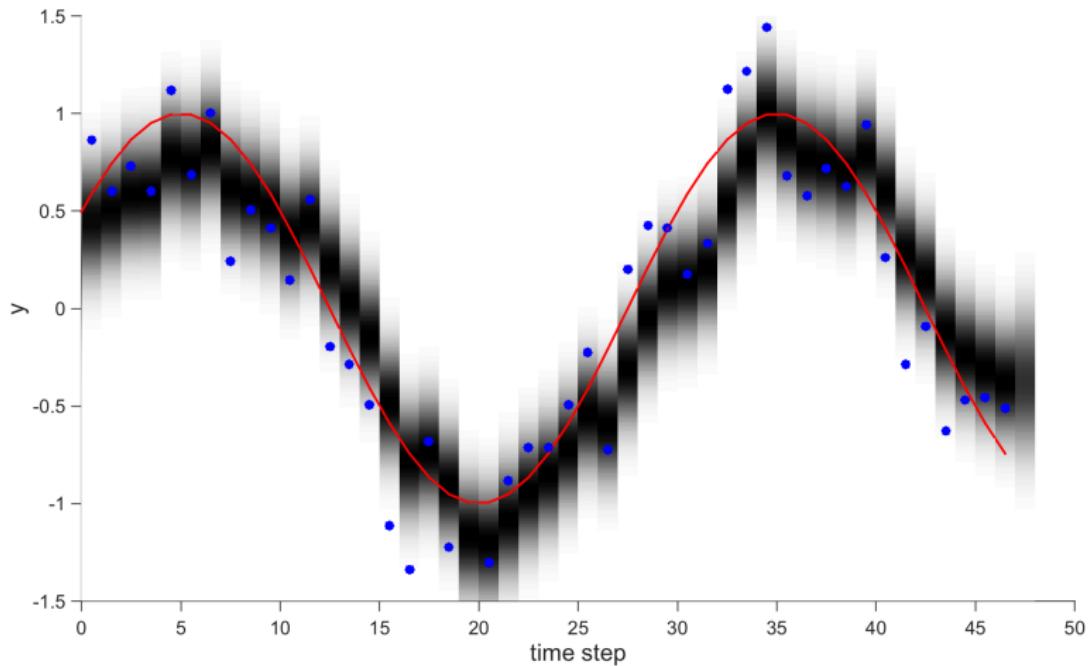
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



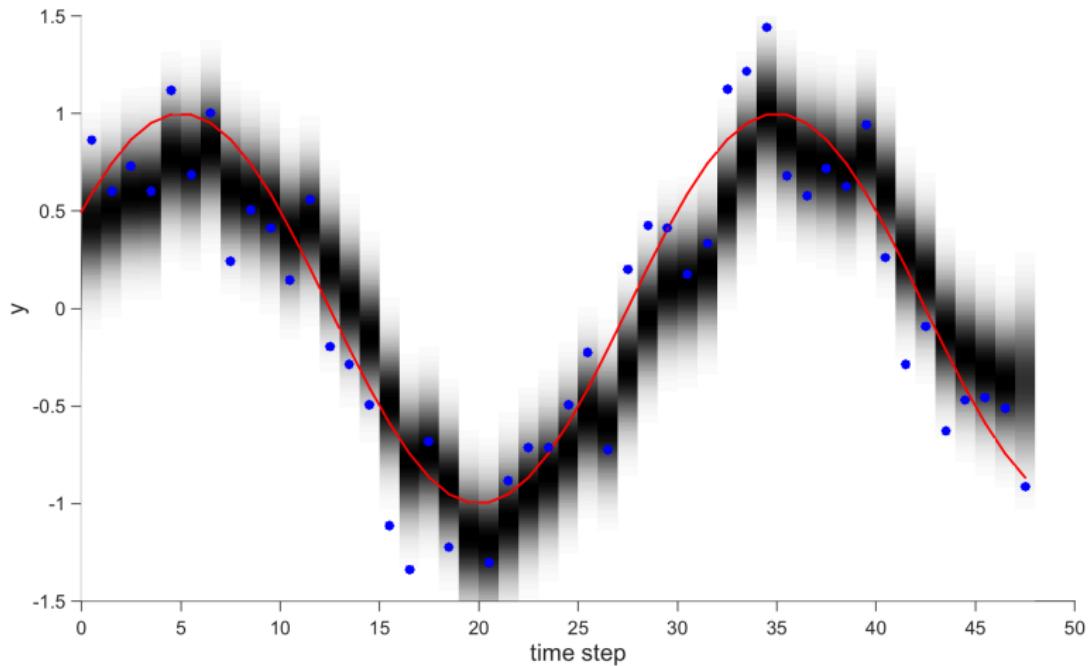
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



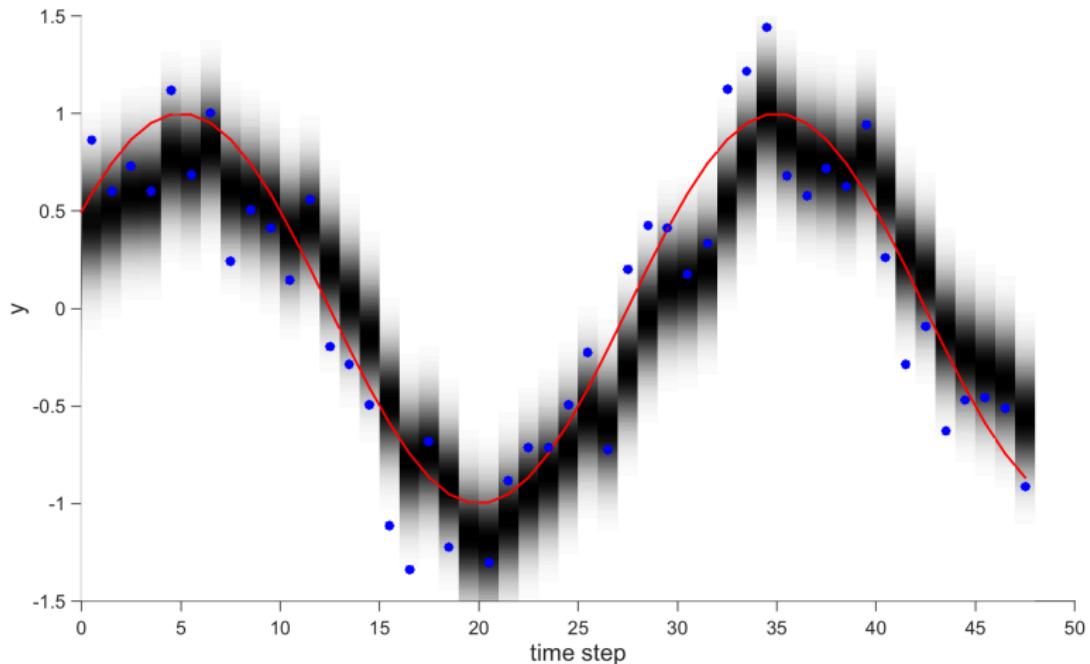
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



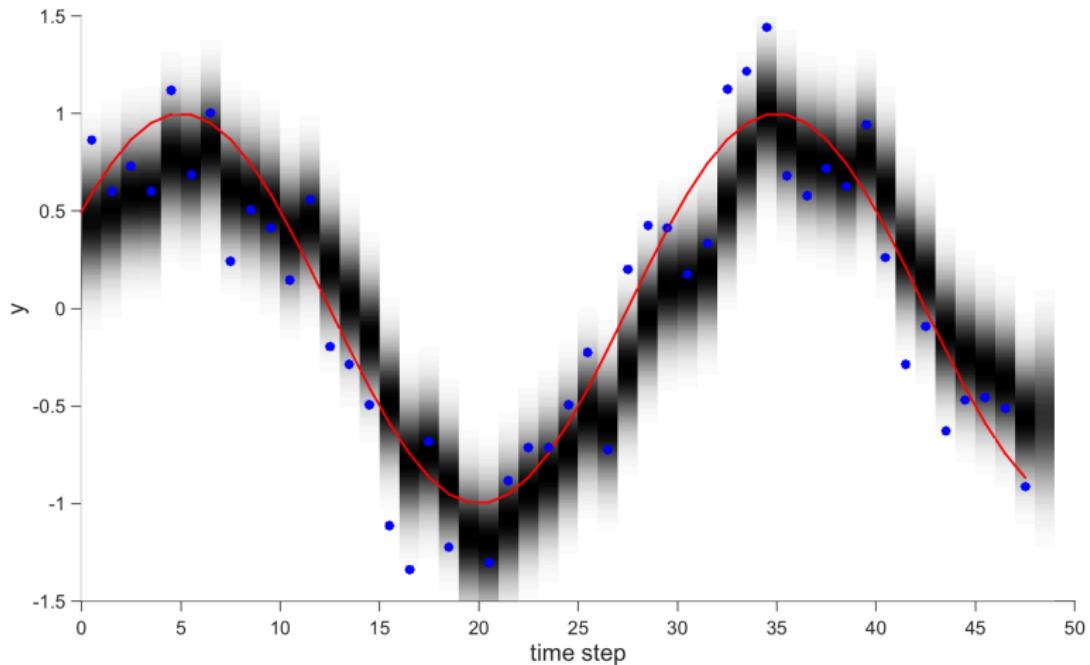
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



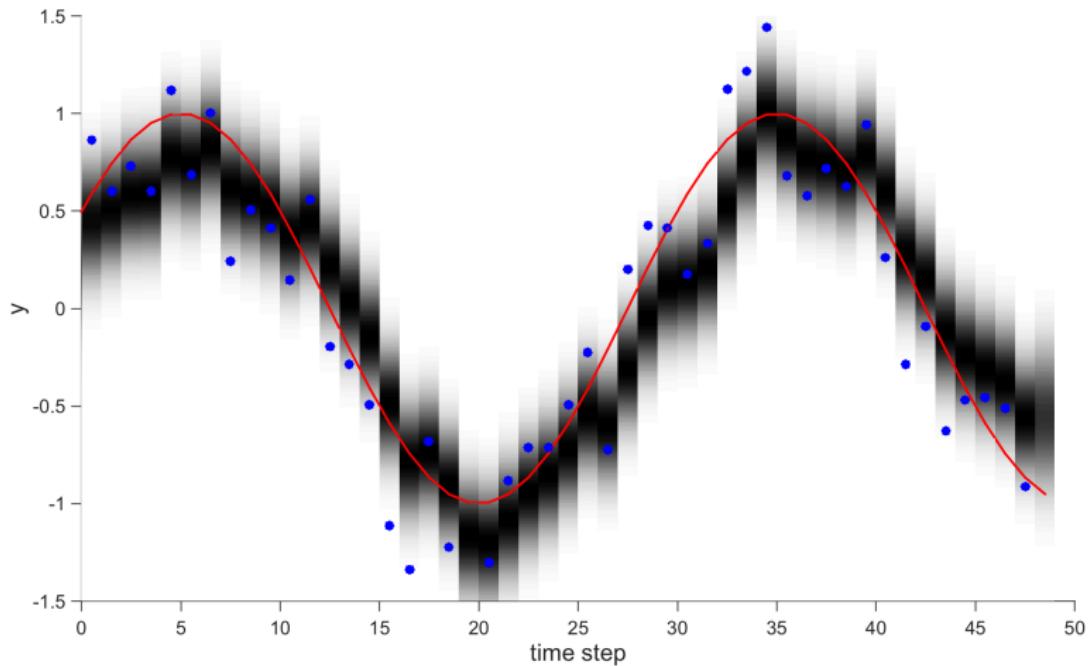
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



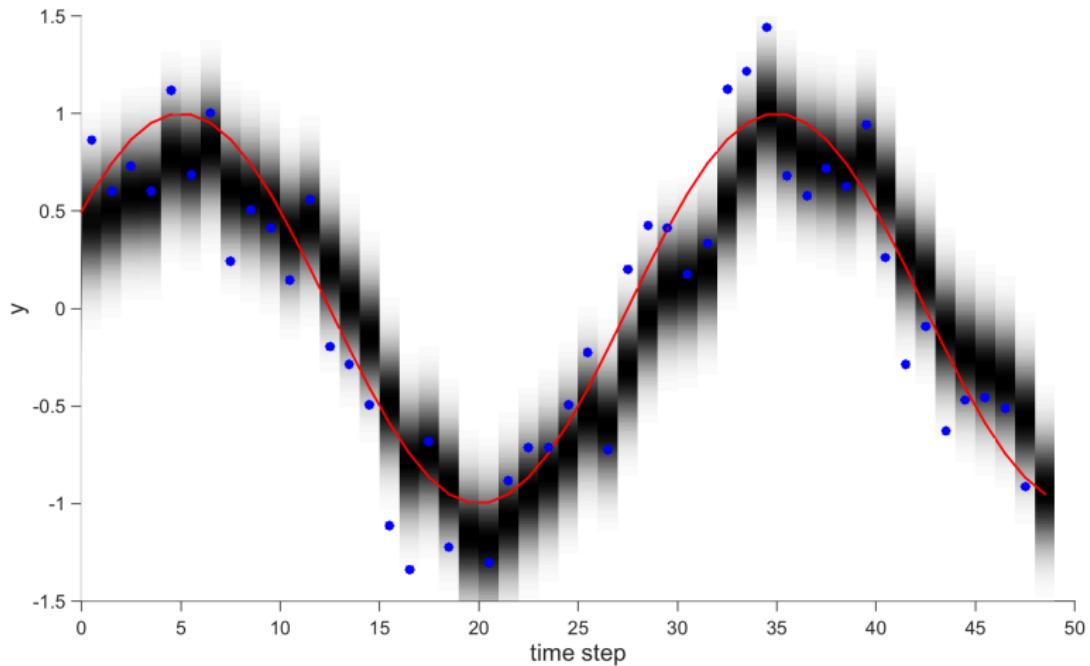
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



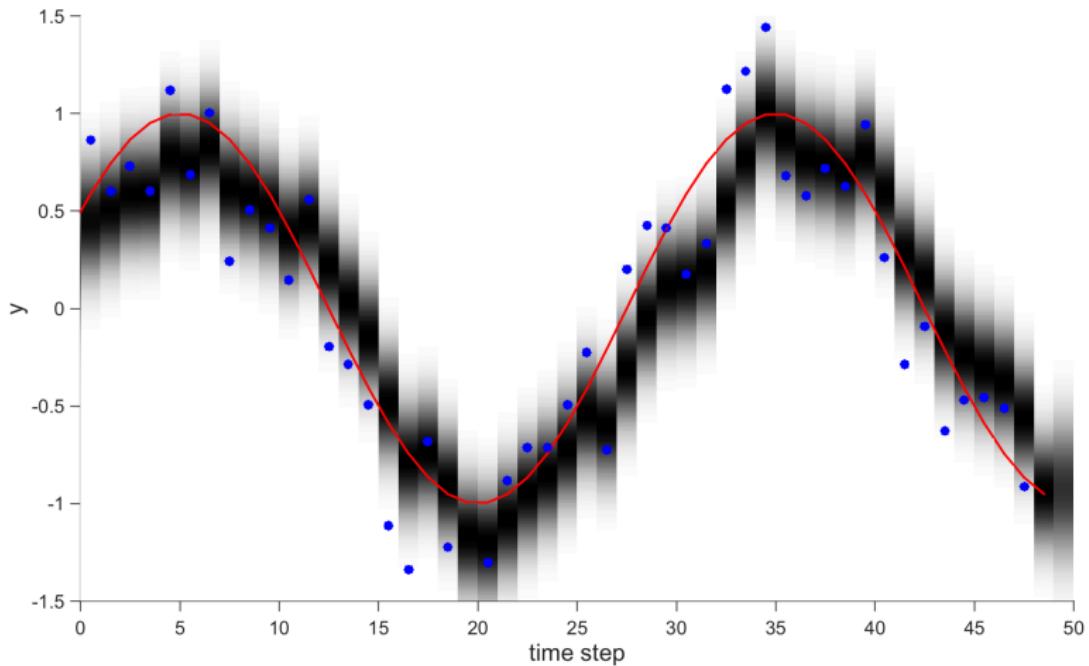
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



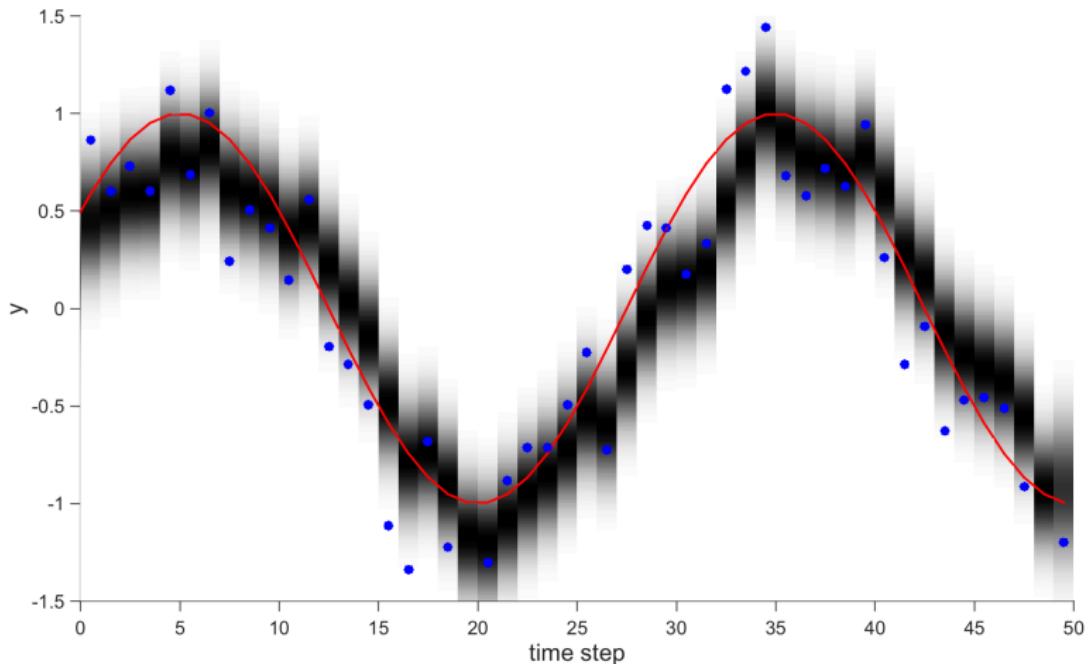
Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid



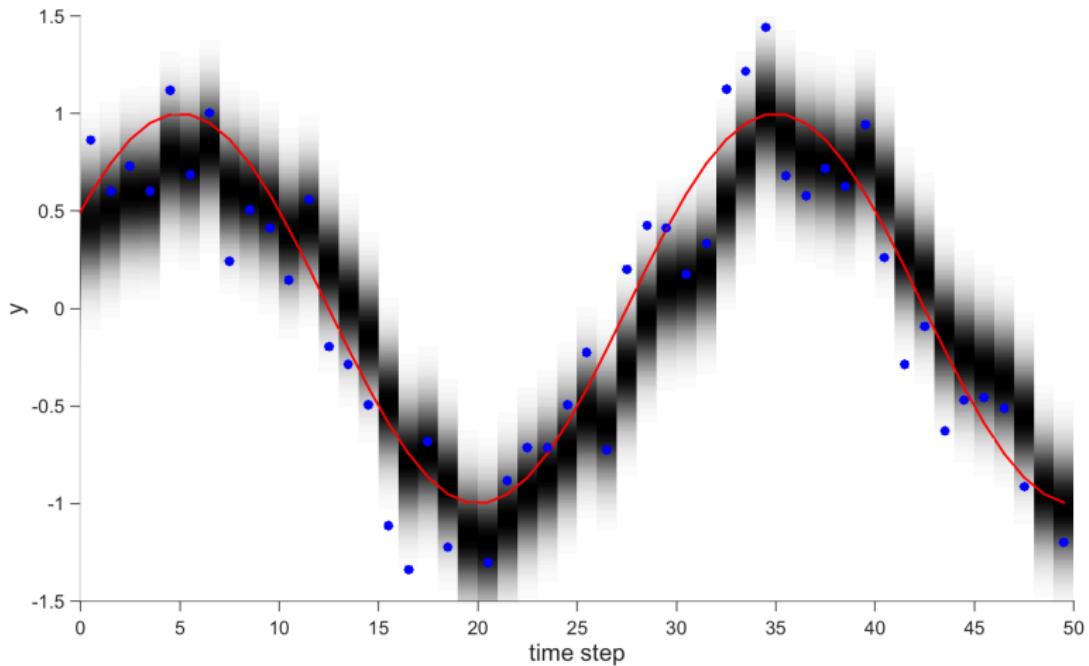
Kalman Filter Demo

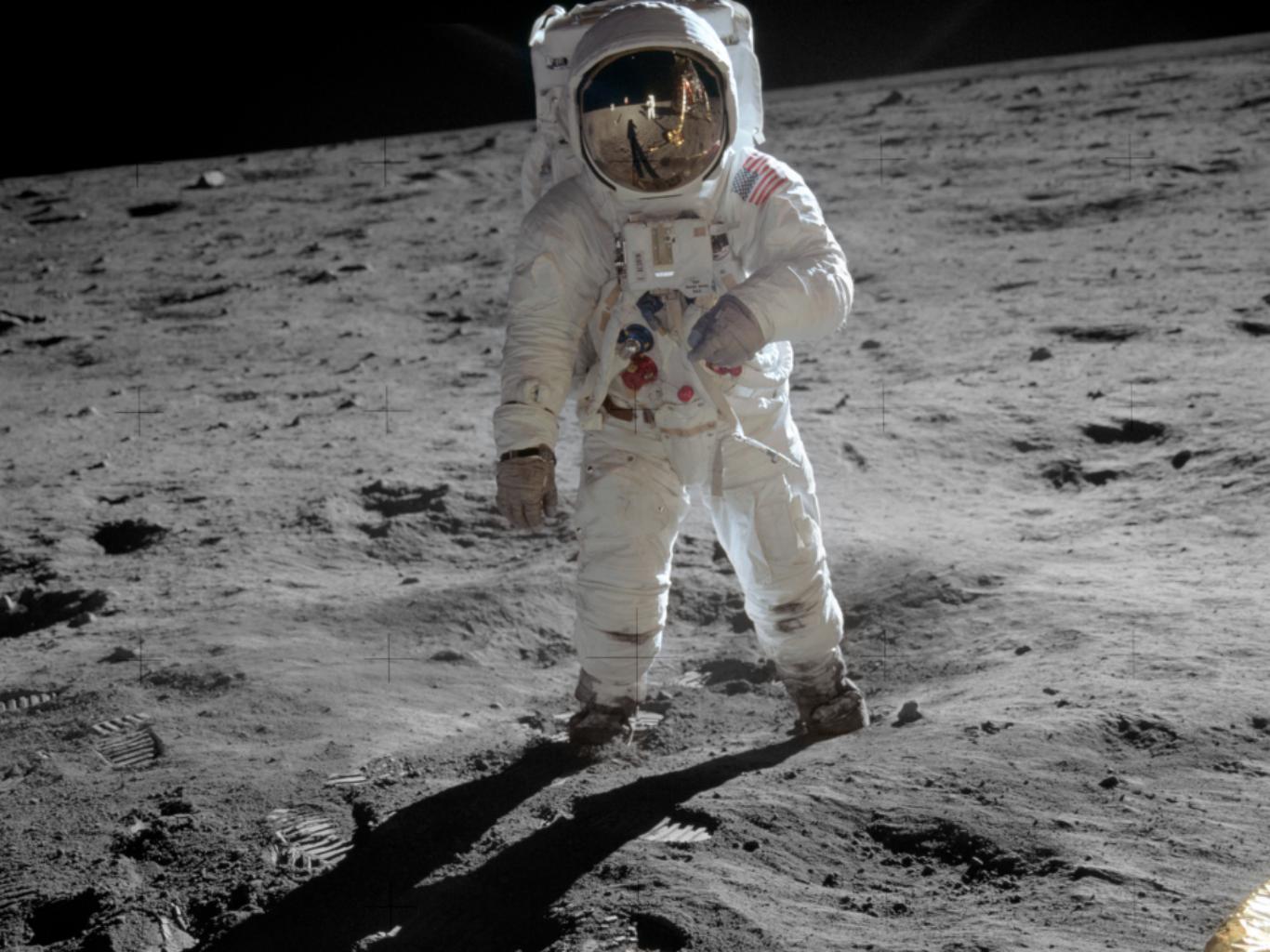
observed noisy data y_t , ground truth sinusoid



Kalman Filter Demo

observed noisy data y_t , ground truth sinusoid





Course Survey: please complete this!



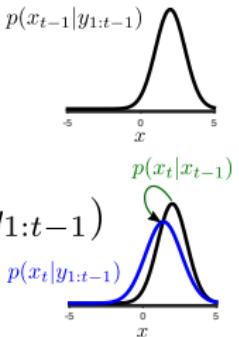
https://cambridge.eu.qualtrics.com/jfe/form/SV_cGCxoPq5McrRPv0

Inference: Forward Algorithm

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

diffuse via
dynamics

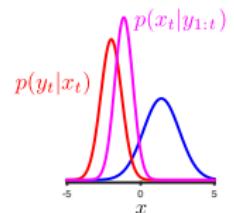
$$p(x_t = k | y_{1:t-1}) = \sum_{l=1}^K p(x_t = k | x_{t-1} = l) p(x_{t-1} = l | y_{1:t-1})$$



combine
with
likelihood

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

prior likelihood



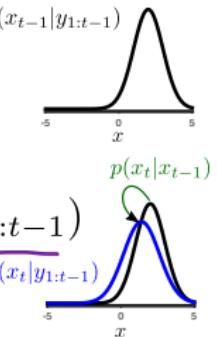
Inference: Forward Algorithm

$$p(\underline{x_{t-1}} = k | \underline{y_{1:t-1}}) = \rho_{t-1}^{t-1}(k) \quad \begin{array}{l} \text{most recent data used} \\ \text{in prediction} \end{array}$$

diffuse via dynamics

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

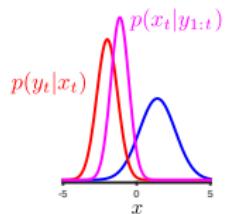
variable being predicted



combine
with
likelihood

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

prior likelihood



Inference: Forward Algorithm

$$p(x_{t-1} = k | y_{1:t-1}) = \rho_{t-1}^{t-1}(k) \quad \begin{matrix} \leftarrow \text{most recent data used} \\ \text{in prediction} \end{matrix}$$

diffuse via dynamics

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

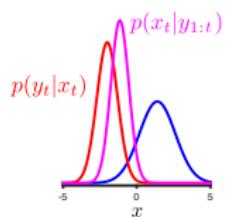
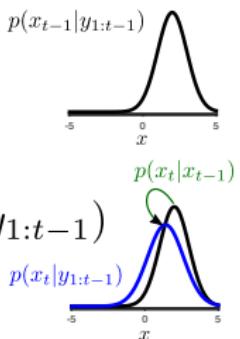
$$\rho_t^{t-1}(k) = \sum_{l=1}^K T(k, l) \rho_{t-1}^{t-1}(l)$$

combine with likelihood

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

prior

likelihood



Inference: Forward Algorithm

$$p(x_{t-1} = k | y_{1:t-1}) = \rho_{t-1}^{t-1}(k) \quad \begin{array}{l} \text{most recent data used} \\ \text{in prediction} \end{array}$$

variable being predicted

diffuse via dynamics

$$\downarrow$$

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

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

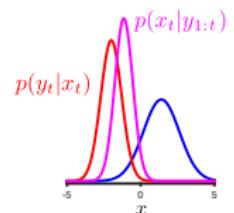
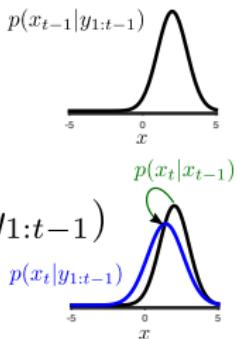
combine with likelihood

$$\downarrow$$

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

prior likelihood

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



Inference: Forward Algorithm

$$p(x_{t-1} = k | y_{1:t-1}) = \rho_{t-1}^{t-1}(k) \quad \begin{array}{l} \text{most recent data used} \\ \text{in prediction} \end{array}$$

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

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

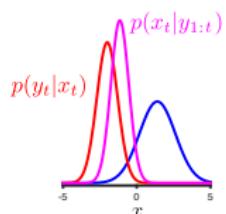
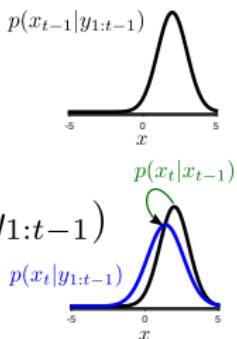
prior

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?

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

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

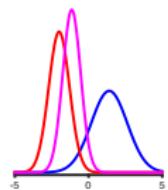
Computing the likelihood

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

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



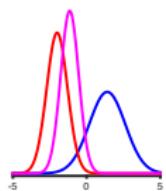
Computing the likelihood

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

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



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

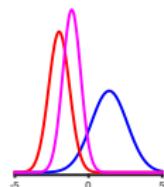
Computing the likelihood

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

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



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

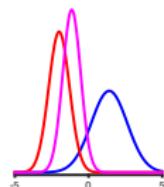
Computing the likelihood

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

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



$p(y_t | y_{1:t-1})$ 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: $x_{1:T}^{(k)}$

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

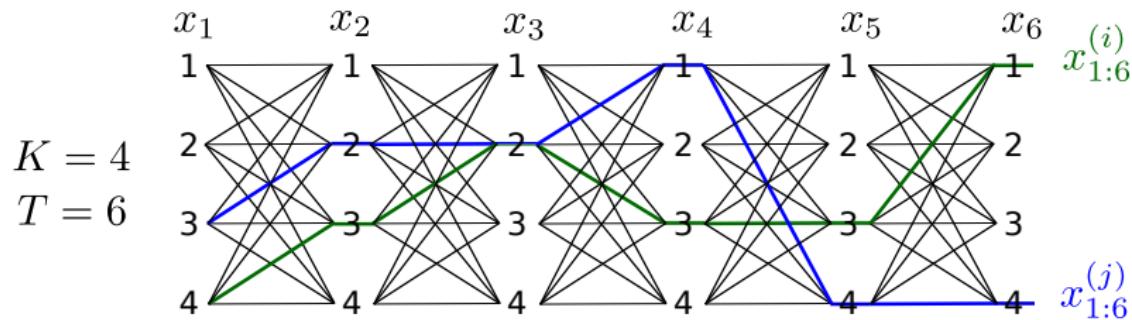
The magic of the Forward Algorithm: Dynamic Programming

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:



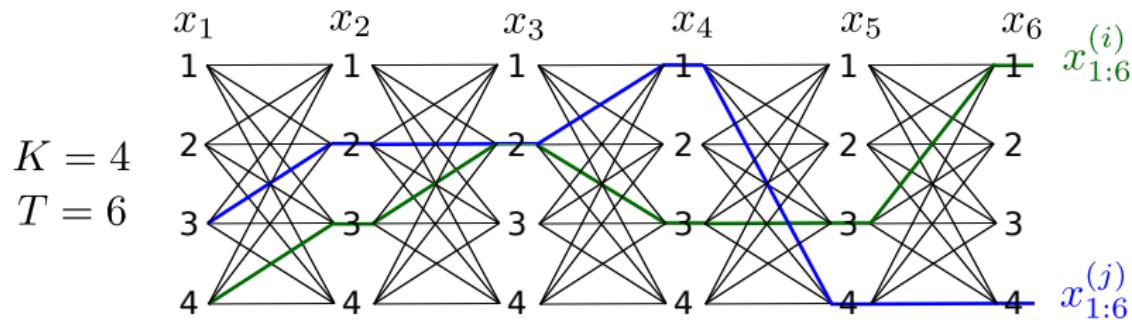
The magic of the Forward Algorithm: Dynamic Programming

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

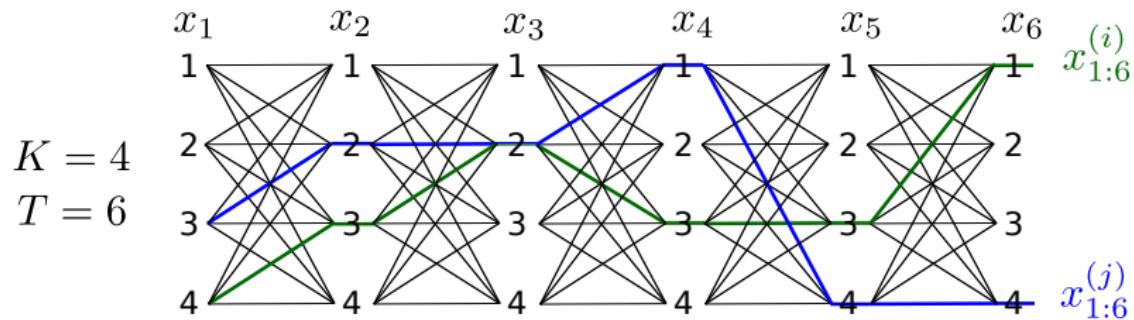
The magic of the Forward Algorithm: Dynamic Programming

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)

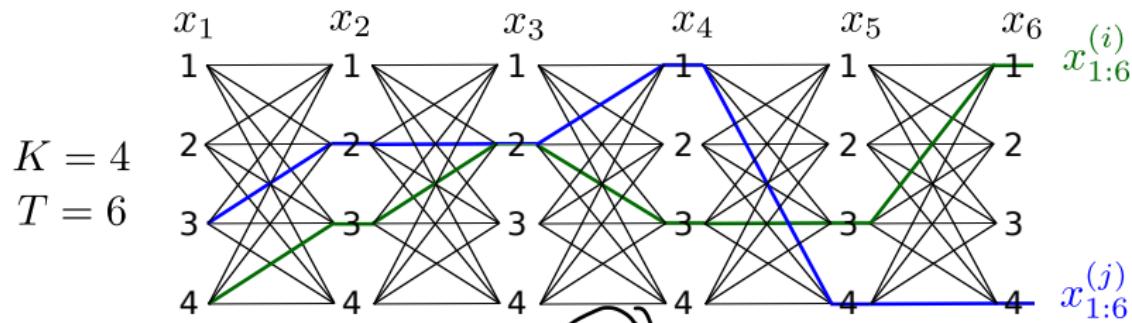
The magic of the Forward Algorithm: Dynamic Programming

What's going on here?

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

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

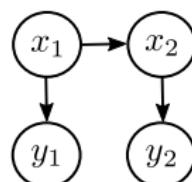
Trellis diagram represents possible sequences:



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:

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

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

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:

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

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

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

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:

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

show gradient depends
on simple moments
of posterior:

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

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

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

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

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

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

show gradient depends
on simple moments
of posterior:

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

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

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

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

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

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:

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

show gradient depends
on simple moments
of posterior:

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

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

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

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

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

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

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

↑
requires posterior moments: marginals and pairwise marginals

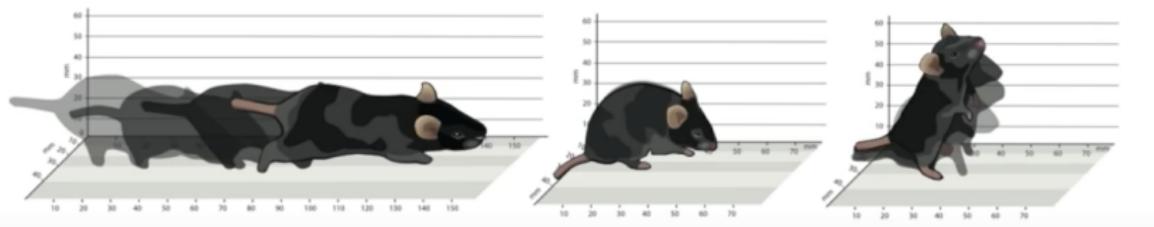
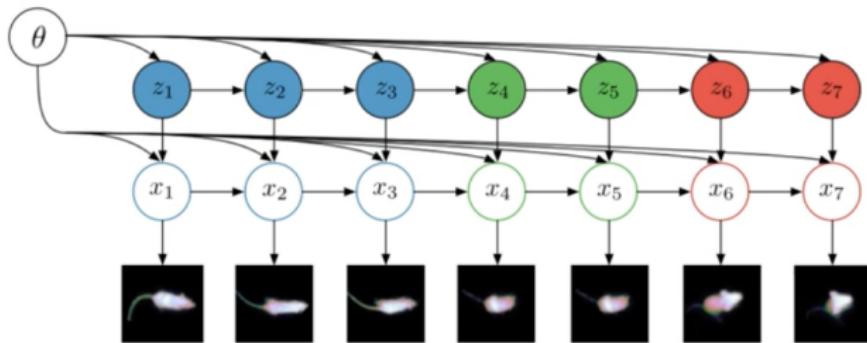
Summary of Sequence Modelling

- ▶ **Markov models:** class of probabilistic models for sequence data
 - ▶ **N-Gram models** (discrete data) and **Gaussian auto-regressive models** (continuous data)
 - ▶ simple to perform maximum likelihood fitting
 - ▶ unnatural for many tasks e.g. removing additive noise, separating signals, data containing latent variables
- ▶ **Hidden Markov Models:** more flexible class of probabilistic model
 - ▶ generalise Markov models and naturally support a wider range of tasks (removal of noise, source separation, representation learning)
 - ▶ different varieties: discrete vs. continuous latent variables (**discrete HMMs** and **linear Gaussian state space models**)
 - ▶ inference in these models requires **dynamic programming** / message passing e.g. smoothing via the **forwards-backwards** recursions or **Kalman filtering-smoothing** recursions
 - ▶ maximum-likelihood fitting requires smoothing as a subroutine

Hidden Markov Models for unsupervised high dimensional video understanding



Hidden Markov Models for unsupervised high dimensional video understanding



<https://www.youtube.com/watch?v=btr1poCYIzw>