

# Sequence Modelling

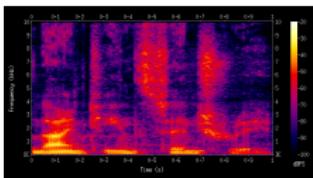
Rich Turner

# Sequence data

\*



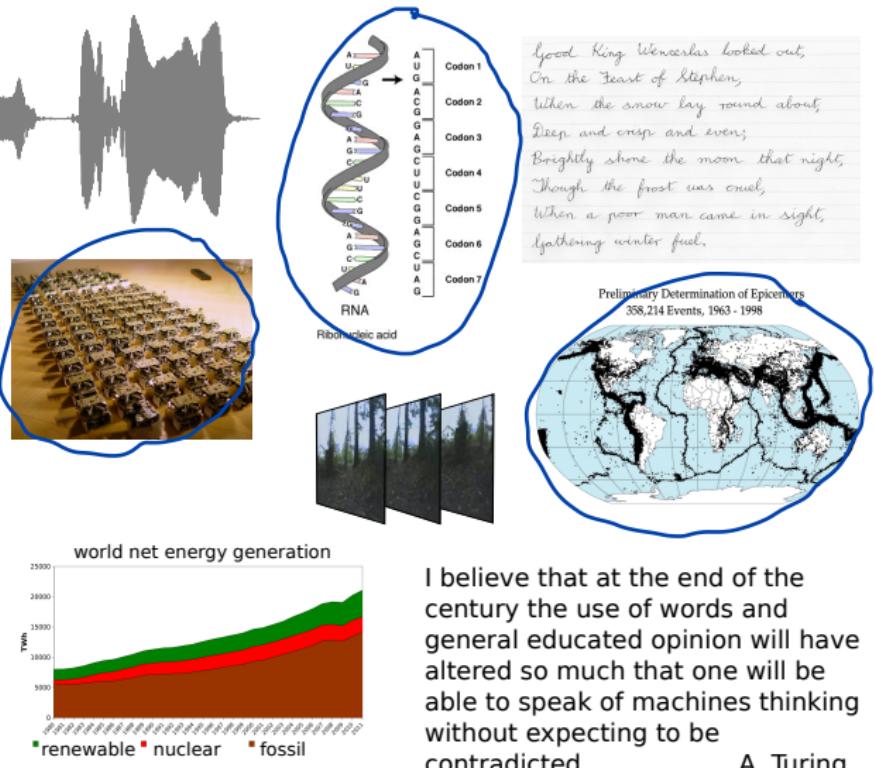
\*



\*



Some images taken from wikipedia



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

A. Turing

## Goals of sequence modelling

Predict future items in sequence

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

Remove noise from a sequence

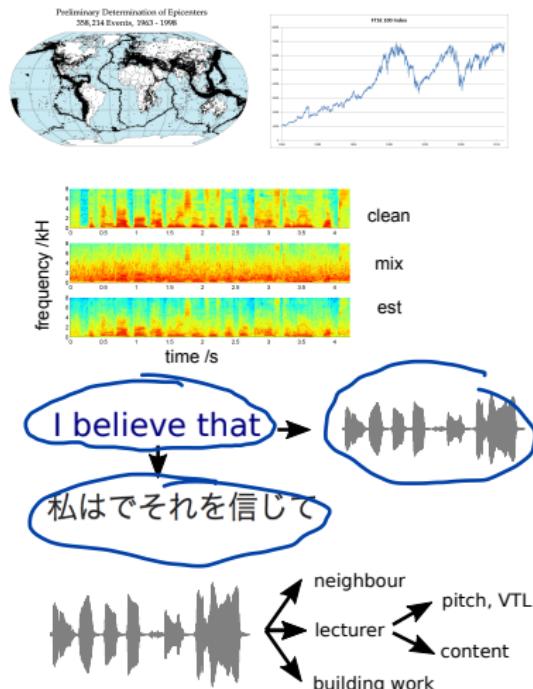
$$p(\underbrace{y'_1, \dots, y'_t} | \underbrace{y_1, \dots, y_t})$$

Predict one sequence from another

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

Discover underlying latent variables

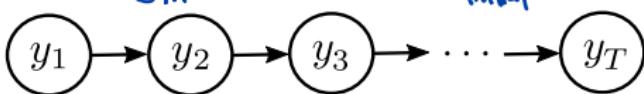
$$p(x_1, \dots, x_t | y_1, \dots, y_t)$$



## Markov models

### First order Markov

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



$$p(y_2=k|y_1=l) = p(y_3=k|y_2=l)$$

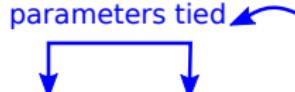
$$p(y_{1:T}) = N(y_{1:T} | \mu_{1:T}, \Sigma_{1:T, 1:T})$$

## Markov models

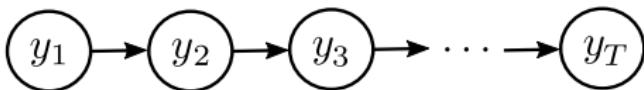
First order Markov

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

parameters tied

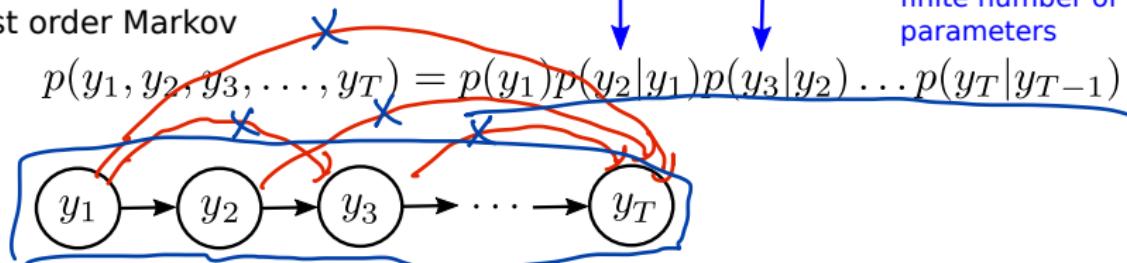


$\infty$  number of variables  
finite number of parameters



## Markov models

First order Markov



Markov model = conditional independence relationship + product rule

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

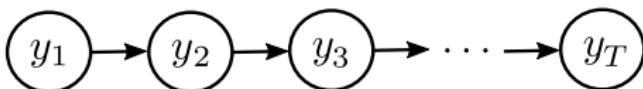
$$p(y_{1:T}) = \prod_{t=1}^T p(y_t | y_{1:t-1})$$
$$p(y_{1:T}) = p(y_1) p(y_2 | y_1) p(y_3 | y_1, y_2) \dots p(y_T | y_{1:T-1})$$
$$p(y_3 | y_2)$$
$$p(y_T | y_{T-1})$$

## Markov models

First order Markov

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

parameters tied  
∞ number of variables  
finite number of parameters



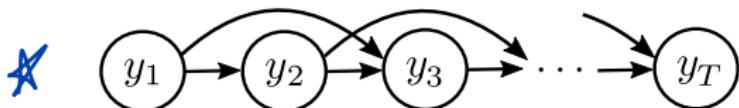
Markov model = conditional independence relationship + product rule

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

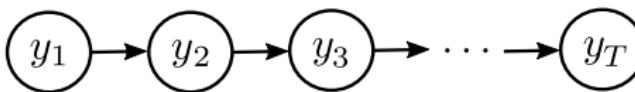
$$p(y_1, y_2, y_3, \dots, y_T) = p(y_1)p(y_2|y_1)p(y_3|y_2, y_1)\dots p(y_T|\underline{y_{T-1}, y_{T-2}})$$



## Markov models for discrete data: n-gram models

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

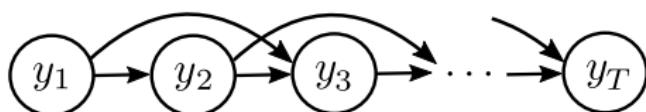
$$p(y_t = k | y_{t-1} = l) = T_{k,l}$$

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

100 x 100 x 100

## Some questions about n-gram models

First order Markov (bi-gram)

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

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

$$\begin{aligned} p(\underbrace{y_2 = k}_{\text{marginal}} | \theta) &= \sum_{l=1}^K p(y_2 = k, y_1 = l | \theta) \\ &= \sum_{l=1}^K p(y_2 = k | y_1 = l, \theta) p(y_1 = l | \theta) \\ &= \sum_{l=1}^K T_{kl} \quad \pi_l^0 \\ &= [T \pi^0]_k \end{aligned}$$

Sum rule  
product rule

## 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?

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

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

$$\lim_{t \rightarrow \infty} p(y_t = k | \theta) = \sum_l T_{k,l} p(y_{t-1} = l | \theta)$$

||

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

$\underline{\underline{M}} \underline{\underline{e}} \underline{\underline{\mu}} = \lambda \underline{\underline{\mu}}$

## 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

First order Markov (bi-gram)

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

	A	B	C ← end	in
A	2/5	3/5	0	5
B	2/5	1/5	2/5	5
C	0	1/3	2/3	3

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?

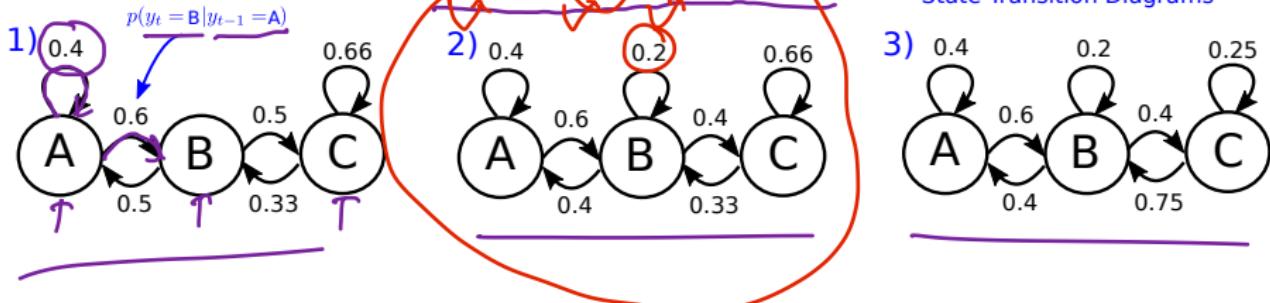
$$p(y_t = k) = \sum_{l=1}^K p(y_t = k | y_{t-1} = l) p(y_{t-1} = l)$$

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

eigenvectors of transition matrix with eigenvalue = 1

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

~~ABAAABBABCCCB~~



## 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?

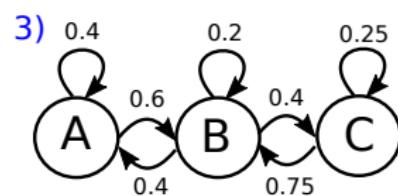
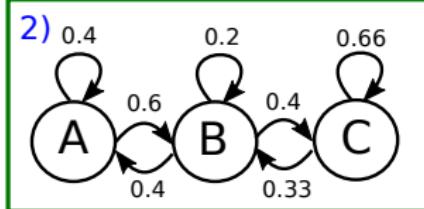
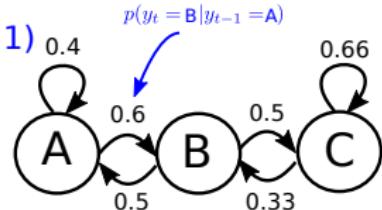
$$p(y_t = k) = \sum_{l=1}^K p(y_t = k | y_{t-1} = l) p(y_{t-1} = l)$$

eigenvectors of transition matrix with eigenvalue = 1

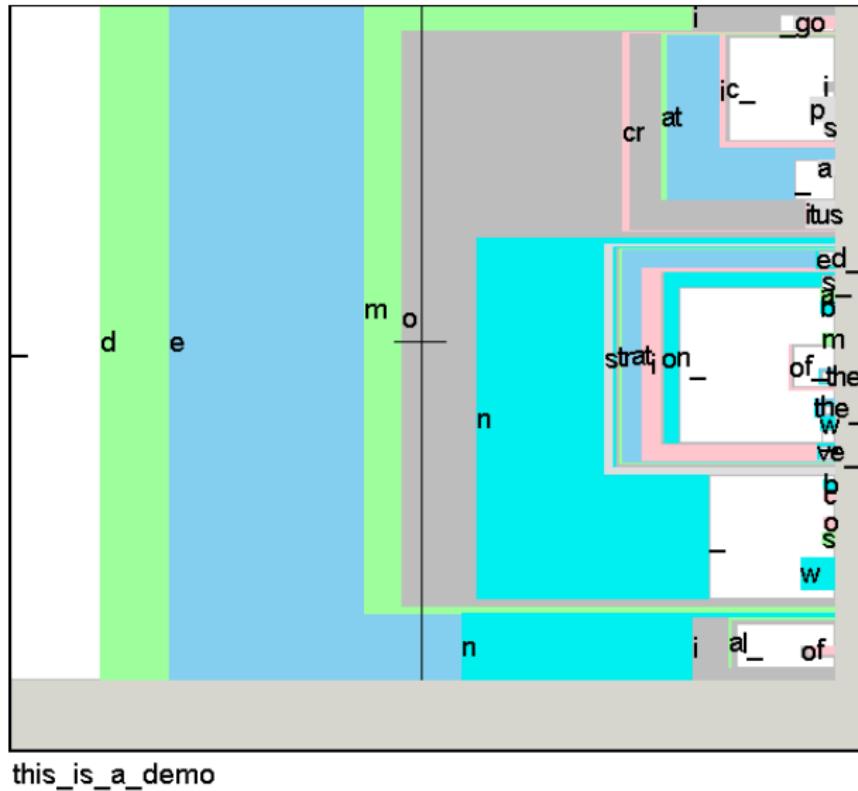
$$\pi_k^\infty = \sum_{l=1}^K T_{k,l} \pi_l^\infty$$

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

ABAAABBABCCCB



## Example application of n-grams: text modelling for dasher



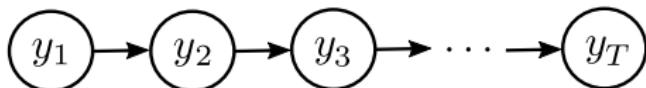
<http://www.inference.phy.cam.ac.uk/dasher/>

<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

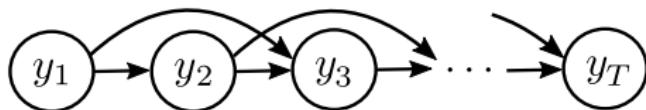
$$p(y_t = k | y_{t-1} = l) = T_{k,l}$$

transition probabilities  
(stochastic matrix)

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

### Second order Markov (tri-gram)

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



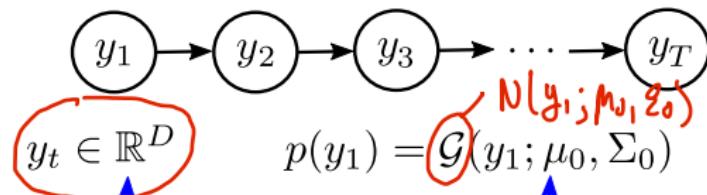
$$p(y_t = k | y_{t-1} = l, y_{t-2} = m) = T_{k,l,m}$$

n-grams require large  
multidimensional arrays

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

First order Markov (AR(1))

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



$$\underline{y}_t = \underline{W} \underline{x}_t + \underline{\varepsilon}_t$$

$$\underline{y}_t = \underline{\Lambda} \underline{y}_{t-1} + \underline{\varepsilon}_t$$

$\uparrow$        $\downarrow$        $\varepsilon_t \sim N(0, \Sigma)$

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

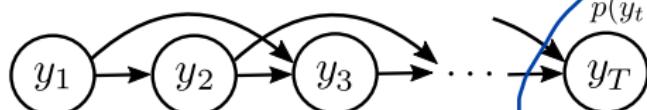
transition density

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

$$\begin{pmatrix} \Lambda & \Sigma \\ M_0 & \Sigma_0 \end{pmatrix}$$

Second order Markov (AR(2))

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



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

joint distribution over all variables  
is always multivariate Gaussian

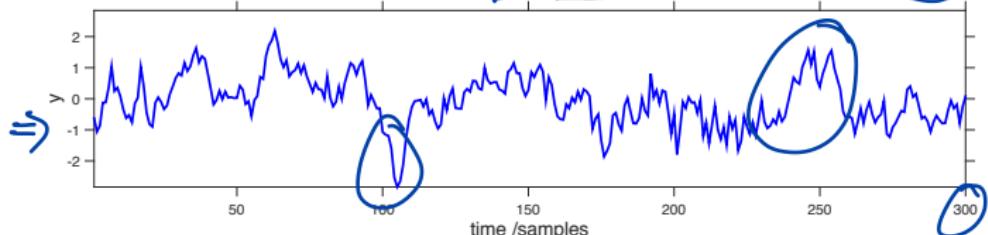
$$p(\underline{y}_{1:T} | \theta) = \mathcal{G}(\underline{y}_{1:T}; \underline{\mu}_{1:T}, \underline{\Sigma}_{1:T, 1:T})$$

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

First order Markov (AR(1))

$$y_t = 0.9 \times y_{t-1} + \xi_t \sqrt{0.01} \quad \xi_t \sim \mathcal{N}(0, 1)$$

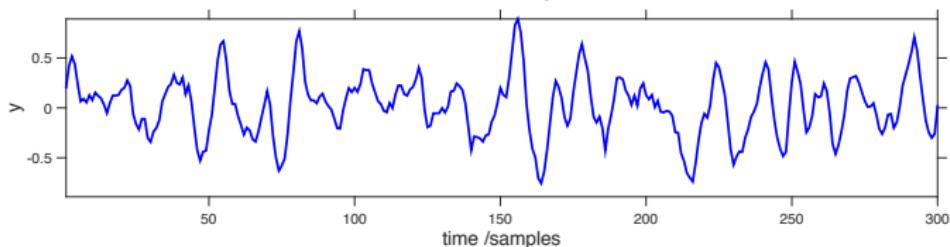
$$y_t \in \mathbb{R}^1 \quad p(y_t | y_{t-1}) = \mathcal{G}(y_t; \underline{\lambda} y_{t-1}, \underline{\sigma^2}) \quad \lambda = 0.9 \quad \sigma^2 = 0.01$$



Second order Markov (AR(2))

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

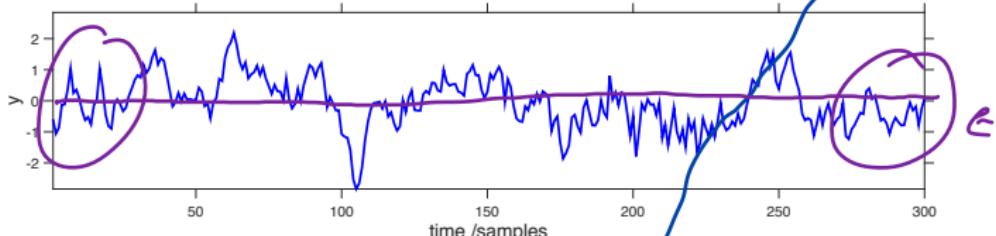
$$[\lambda_1, \lambda_2] = [1.57, -0.78] \quad \sigma^2 = 0.01$$



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

$$\text{First order Markov (AR(1)) } y_t = \lambda y_{t-1} + \varepsilon_t \quad \varepsilon_t \sim \mathcal{G}(0, 1)$$

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



What is the stationary distribution of this process?

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

$$= \mathbb{E}(y_t) = \mathbb{E}(\lambda y_{t-1} + \varepsilon_t) \quad \begin{matrix} \downarrow \\ p(y_{t-1}, \varepsilon_t) \end{matrix}$$

$$= \mu_t = \lambda \mu_{t-1} + \sigma \mathbb{E}(\varepsilon_t) \xrightarrow{\sigma=0} \mu_t = \lambda \mu_{t-1} \quad \begin{matrix} \downarrow \\ 0.9 \end{matrix}$$

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

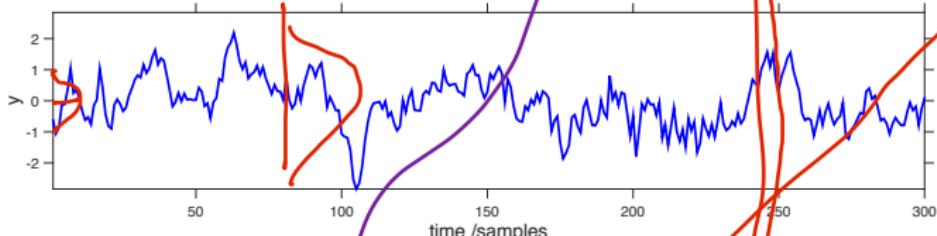
# Markov models for continuous data: Auto-Regressive (AR) Gaussian models

First order Markov (AR(1))

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

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

$$\lambda = 0.9 \quad \sigma^2 = 0.01$$



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

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

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

$$\mathbb{E}(y_{t-1} \varepsilon_t) = \mathbb{E}(y_{t-1}) \mathbb{E}(\varepsilon_t) \quad (\text{because } \varepsilon_t \sim \mathcal{G}(0, 1))$$

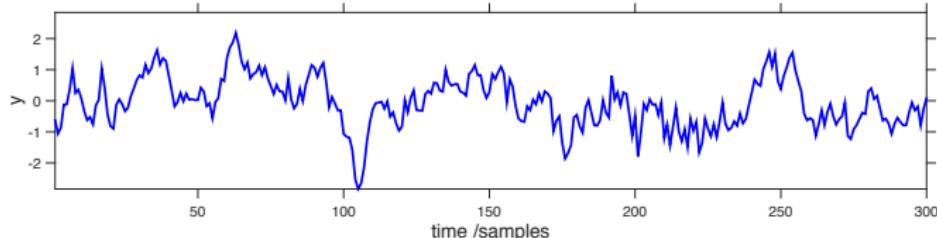
$$\mathbb{E}(y_{t-1}) = \mu_\infty \quad (\text{because it's stationary})$$

$$\mathbb{E}(y_t^2) = \lambda^2 \sigma_\infty^2 + \sigma^2 \quad \Rightarrow \quad \sigma_\infty^2 = \frac{\sigma^2}{1 - \lambda^2} = \frac{0.01}{1 - 0.9^2} = 0.111$$

## 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?  $p(y_\infty) = ?$

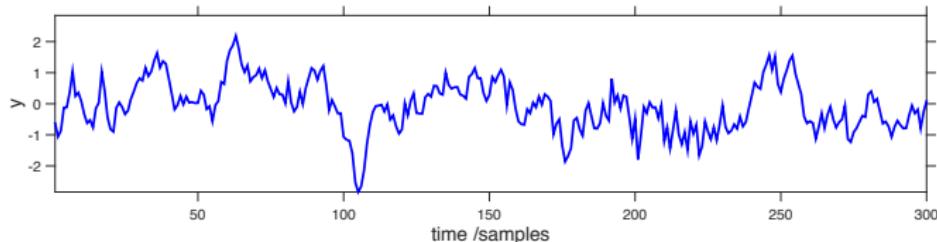
Everything is linear Gaussian  $\Rightarrow$  must be Gaussian  $p(y_\infty) = \mathcal{G}(y_\infty; \mu_\infty, \sigma_\infty^2)$

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

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

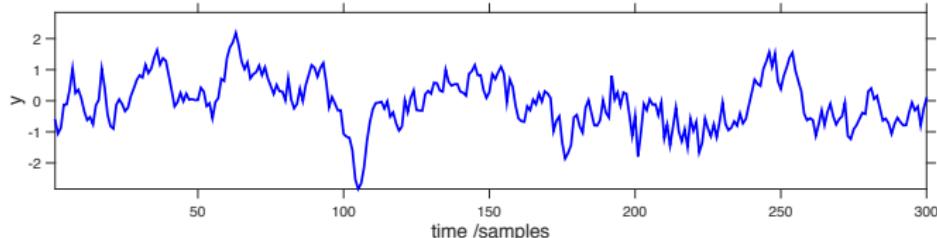
$$y_t = \lambda y_{t-1} + \sigma \epsilon_t \quad \epsilon_t \sim \mathcal{G}(0, 1)$$

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

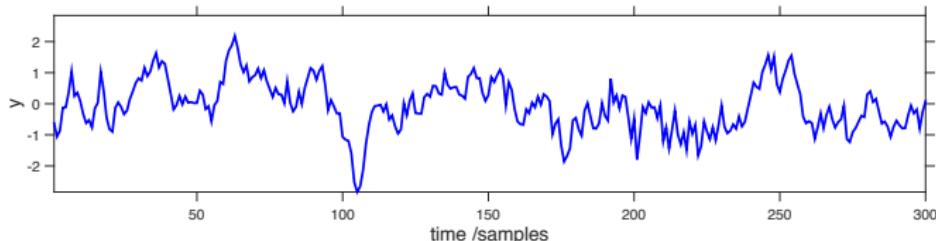
$$y_t = \lambda y_{t-1} + \sigma \epsilon_t \quad \epsilon_t \sim \mathcal{G}(0, 1)$$

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

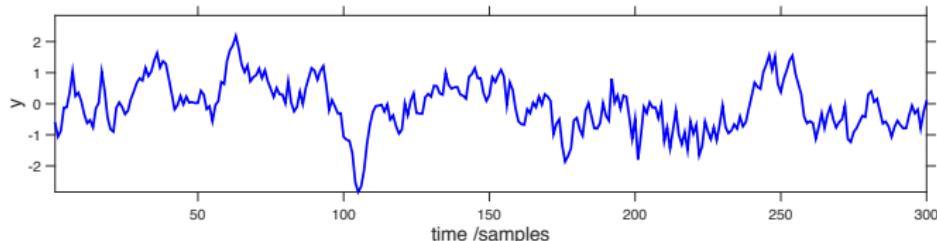
$$y_t = \lambda y_{t-1} + \sigma \epsilon_t \quad \epsilon_t \sim \mathcal{G}(0, 1)$$

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

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

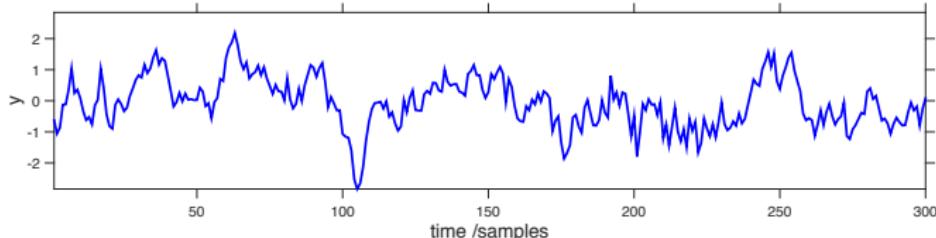
$$\text{Mean: } \langle y_t \rangle = \lambda \langle y_{t-1} \rangle + \sigma \langle \epsilon_t \rangle = 0 \quad \mu_\infty = 0$$

$$\text{Variance: } \langle y_t^2 \rangle$$

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

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

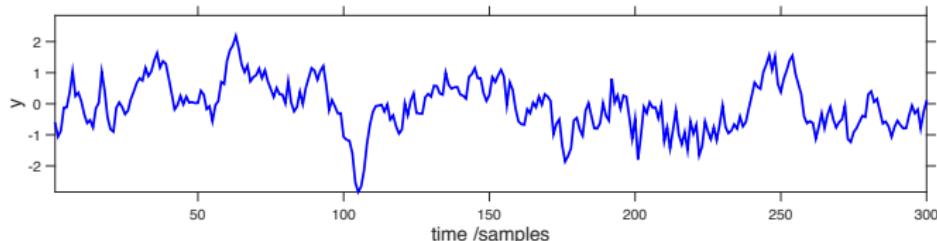
Mean:  $\langle y_t \rangle = \lambda \langle y_{t-1} \rangle + \sigma \langle \epsilon_t \rangle = 0 \quad \mu_\infty = 0$

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

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

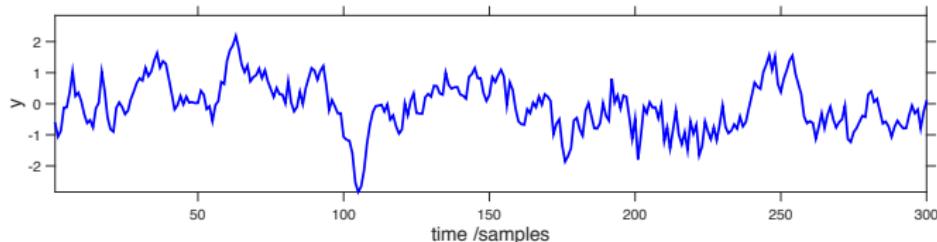
Mean:  $\langle y_t \rangle = \lambda \langle y_{t-1} \rangle + \sigma \langle \epsilon_t \rangle = 0 \quad \mu_\infty = 0$

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

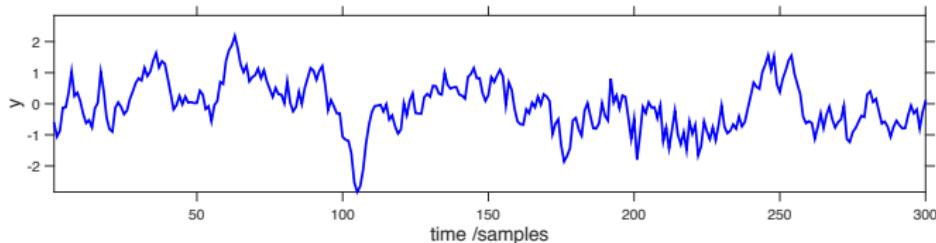
Variance:  $\langle y_t^2 \rangle = \langle (\lambda y_{t-1} + \sigma \epsilon_t)^2 \rangle = \lambda^2 \langle y_{t-1}^2 \rangle + 2\lambda \sigma \langle y_{t-1} \epsilon_t \rangle + \sigma^2 \langle \epsilon_t^2 \rangle$

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

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

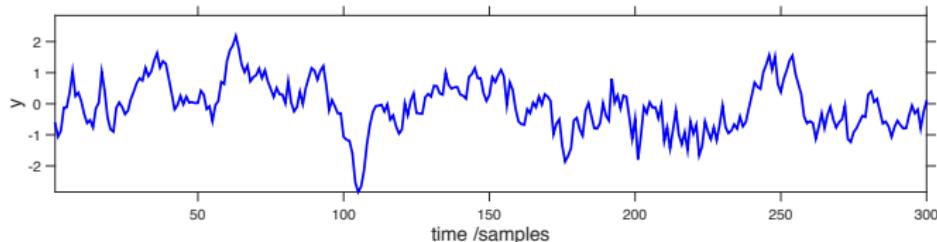
Variance:  $\langle y_t^2 \rangle = \langle (\lambda y_{t-1} + \sigma \epsilon_t)^2 \rangle = \lambda^2 \langle y_{t-1}^2 \rangle + 2\lambda \sigma \langle y_{t-1} \epsilon_t \rangle + \sigma^2 \langle \epsilon_t^2 \rangle$

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

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

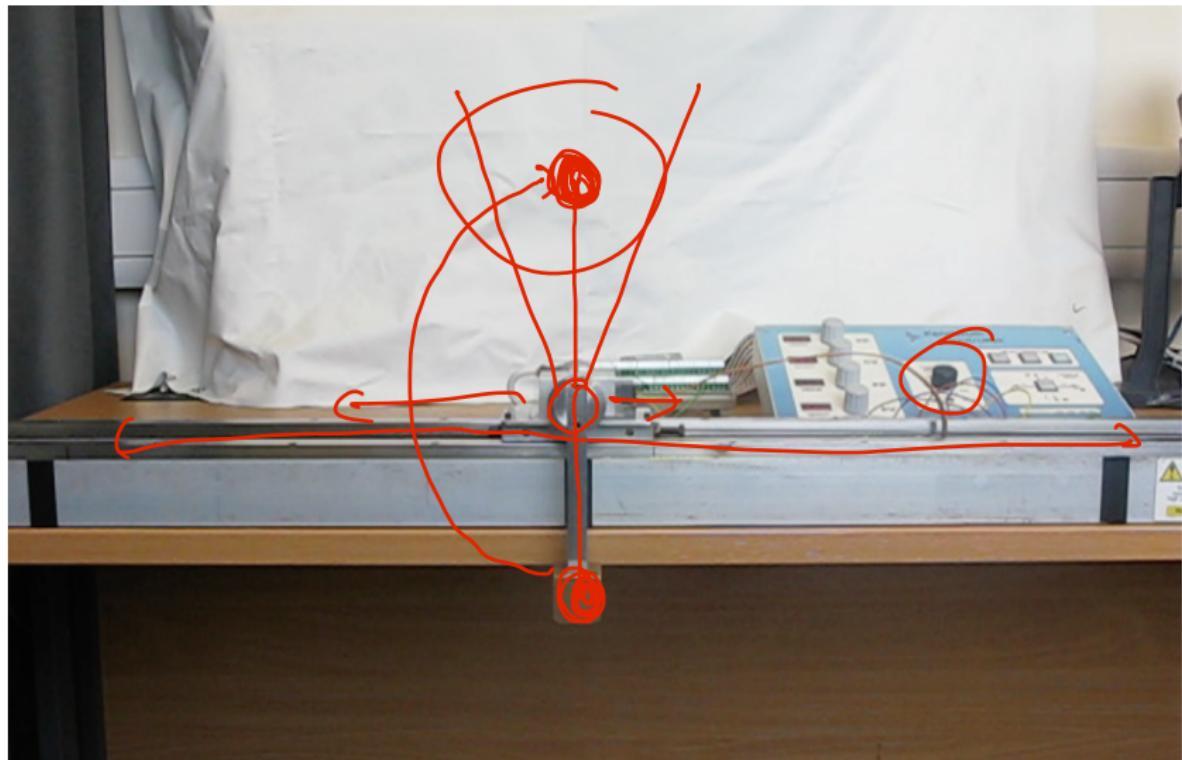
$$y_t = \lambda y_{t-1} + \sigma \epsilon_t \quad \epsilon_t \sim \mathcal{G}(0, 1)$$

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



## 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

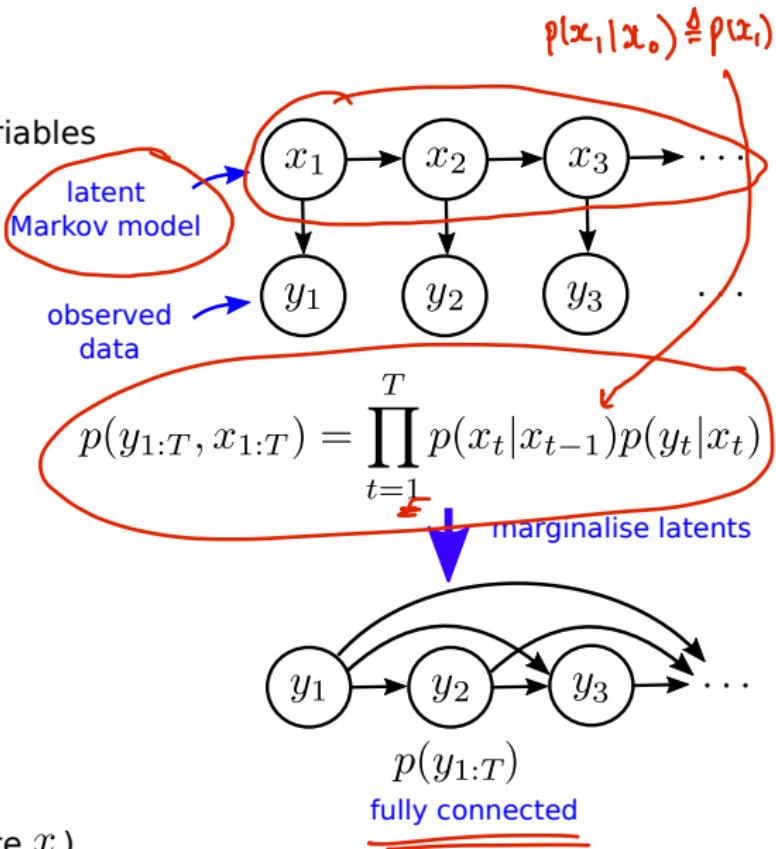
$x$  topics

$y$  words

Two prevalent Examples:

Hidden Markov Models (discrete  $x$ )

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



## Hidden Markov models: discrete hidden state

Discrete Hidden State

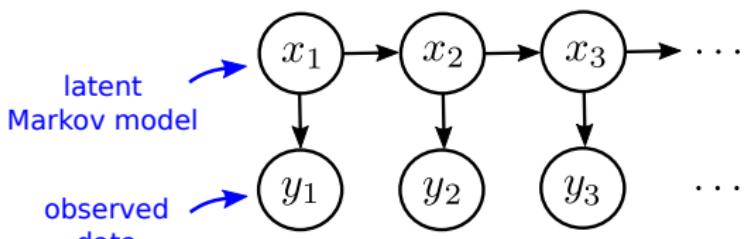
$$x_t \in \{1, \dots, K\}$$

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

E.g. in examples below

$$\underline{\underline{K=2}}$$

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

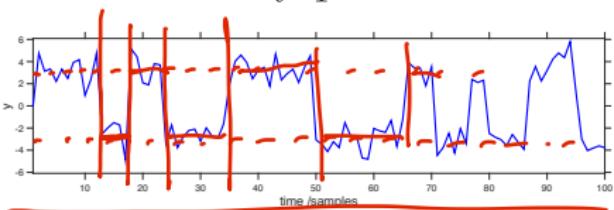


$$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t | x_{t-1}) p(y_t | x_t)$$

Continuous Observed State

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

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



## Hidden Markov models: discrete hidden state

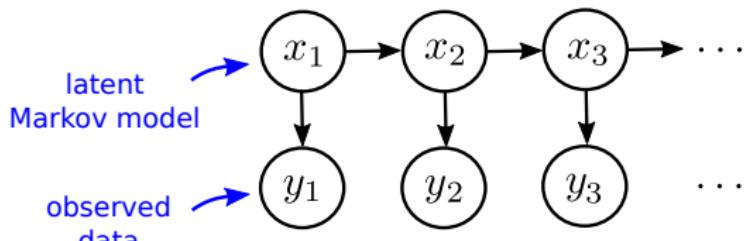
Discrete Hidden State

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

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

E.g. in examples below     $K = 2$

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

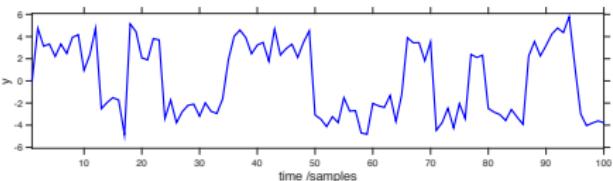


$$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t | x_{t-1}) p(y_t | x_t)$$

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

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

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

ABBBBAAABAAACCCCCCBBBBBCCCCCCCCCCCC  
AAABBBBAABAAABBCCCCCCCCCCCCCCCCCBBA  
AACCCCCCBABCCCCCAABBAABABCCCCC

# Sequence Modelling Recap...

## Markov Models

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

⇒ Discrete  $y \Rightarrow$  bigram models

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

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

Continuous  $y \Rightarrow$  Autoregressive models

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

$$\parallel \quad p(y_t|y_{t-1}) = G(y_t; \mu_{t-1}, \Sigma_{t-1})$$

## Hidden Markov Models

$$\underset{\text{observed}}{p(y_{1:T}, x_{1:T})} = \prod_{t=1}^T p(x_t|x_{t-1}) \underset{\text{hidden / latent}}{p(y_t|x_t)} \quad \text{where } p(x_1|x_0) \stackrel{\Delta}{=} p(x_1)$$

discrete hidden state HMM:

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

$$p(y_t|x_t=k) \stackrel{\text{e.g.}}{=} G(y_t; \mu_k, \Sigma_k)$$

## Hidden Markov models: discrete hidden state

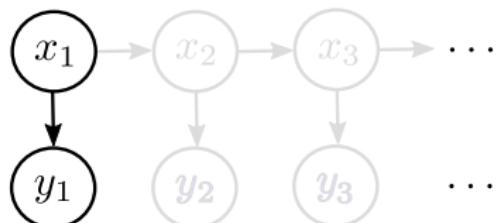
Discrete Hidden State, Continuous Observed State

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

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

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

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



Consider  $T = 1$

Q1: What type of distribution is  $\underline{\underline{p(y_1)}}$ ?

## Hidden Markov models: discrete hidden state

Discrete Hidden State, Continuous Observed State

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

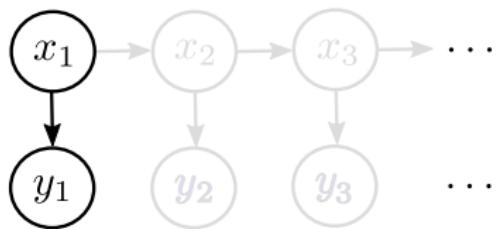
$$p(x_1 = k) = \pi_k^0$$

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

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

Q1: What type of distribution is  $p(y_1)$ ?  $p(y_1, x_1) = p(y_1 | x_1)p(x_1)$

$$p(y_1) = \sum_k p(y_1 | x_1 = k) p(x_1 = k) \leftarrow \text{sum rule + product rule}$$



Consider  $T = 1$

## Hidden Markov models: discrete hidden state

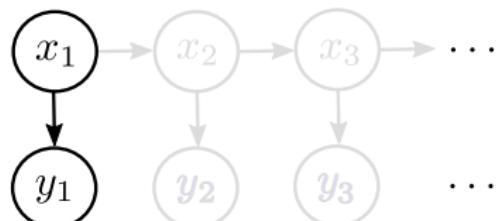
Discrete Hidden State, Continuous Observed State

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

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

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

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



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

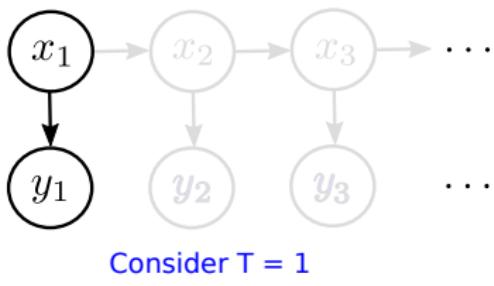
Discrete Hidden State, Continuous Observed State

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

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

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

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



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?

## Hidden Markov models: discrete hidden state

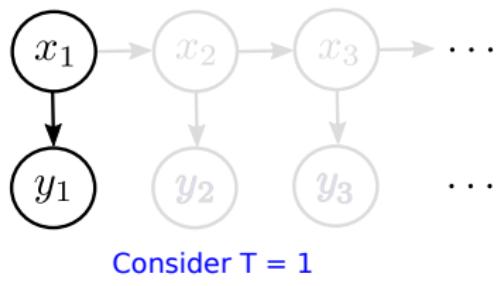
Discrete Hidden State, Continuous Observed State

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

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

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

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



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

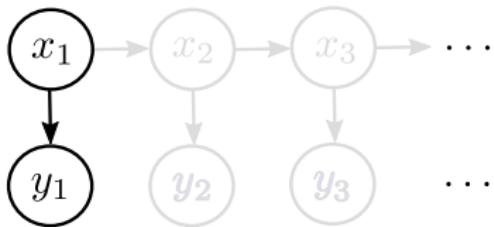
Discrete Hidden State, Continuous Observed State

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

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

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

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



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

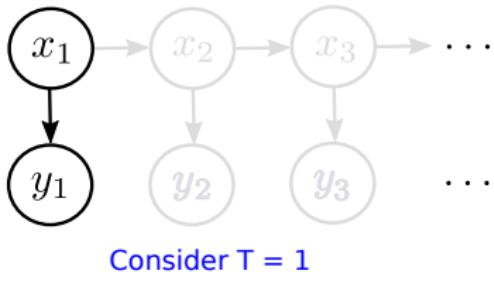
Discrete Hidden State, Continuous Observed State

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

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

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

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



Consider  $T = 1$

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

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

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

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

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

## Hidden Markov models: discrete hidden state

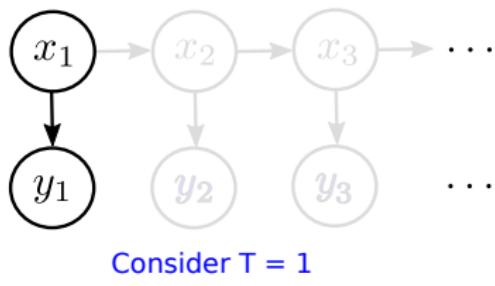
Discrete Hidden State, Continuous Observed State

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

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

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

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



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

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

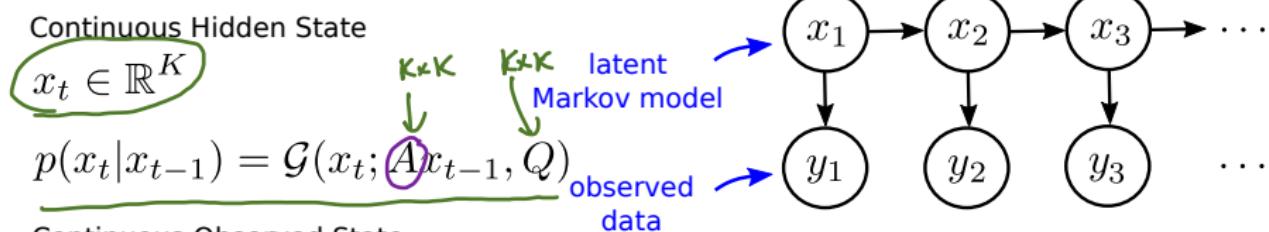
Q2: What distribution does  $p(y_t)$  converge to after a long time?

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)



Continuous Observed State  
 $y_t \in \mathbb{R}^D$

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

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

E.g. simple example

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

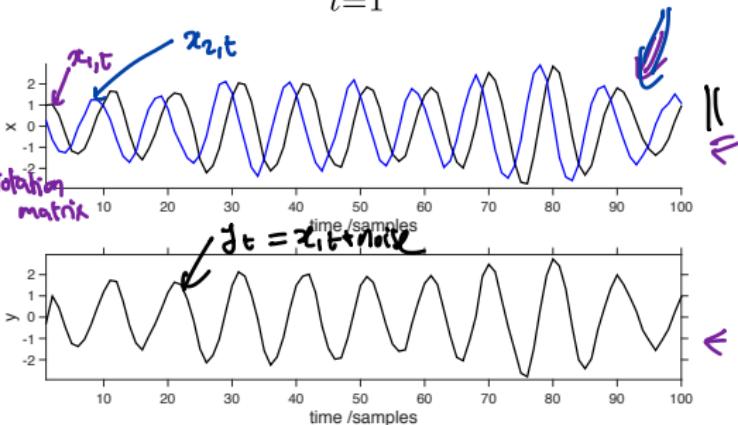
$\lambda = 0.99$        $\theta = 2\pi/10$

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

$C = [1, 0]$        $R = 0.01$

dynamics model

obs. model

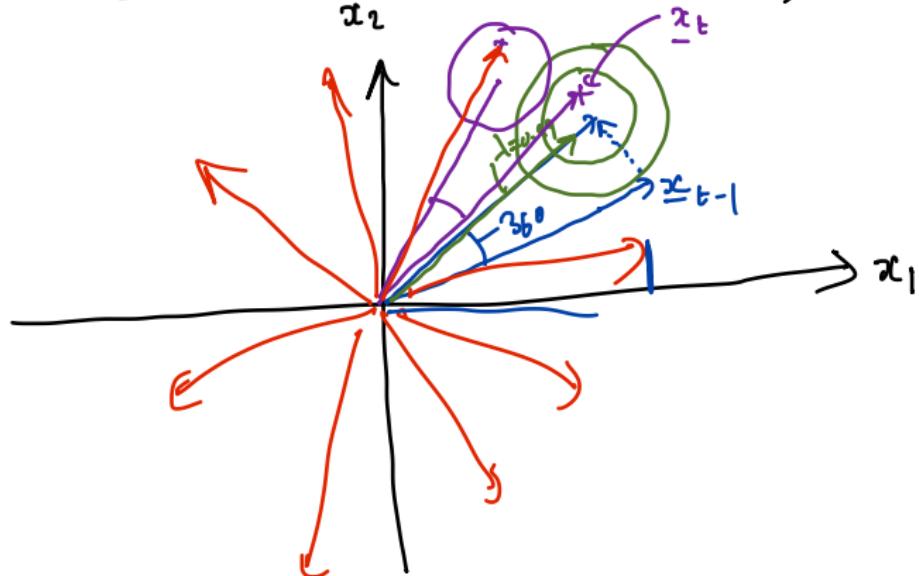


$$\underline{x}_t = \begin{bmatrix} x_{1t} \\ x_{2t} \end{bmatrix} = \lambda \begin{bmatrix} \cos\theta & \sin\theta \\ -\sin\theta & \cos\theta \end{bmatrix} \begin{bmatrix} x_{1,t-1} \\ x_{2,t-1} \end{bmatrix} + (1 - \lambda^2)^{1/2} \underline{\varepsilon}_t$$

$\lambda = 0.99$        $\theta = \frac{\pi}{10}$

$\underline{\varepsilon}_t \sim G(0, I)$

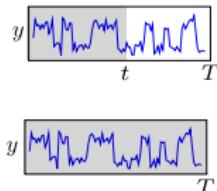
$$y_t = \underline{x}_{1t} + (0.01)^{1/2} n_t \quad n_t \sim G(0, 1)$$



# Varieties of Inference

## Distributional estimates

future data available?



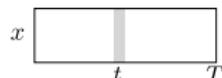
online

filter

smoother

batch

infer single state or sequence?



marginal

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



joint

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

learning

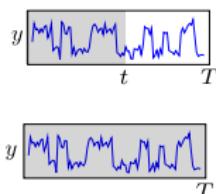
$$p(x_t | y_{1:T})$$

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

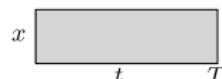
# Varieties of Inference

## 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

$$x_t^* = \arg \max_{x_t} p(x_t | y_{1:T})$$

most probable state @ t

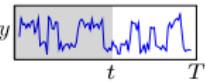
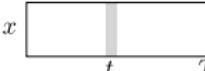
$$x'_{1:T} = \arg \max_{x_{1:T}} p(x_{1:T} | y_{1:T})$$

most probable sequence

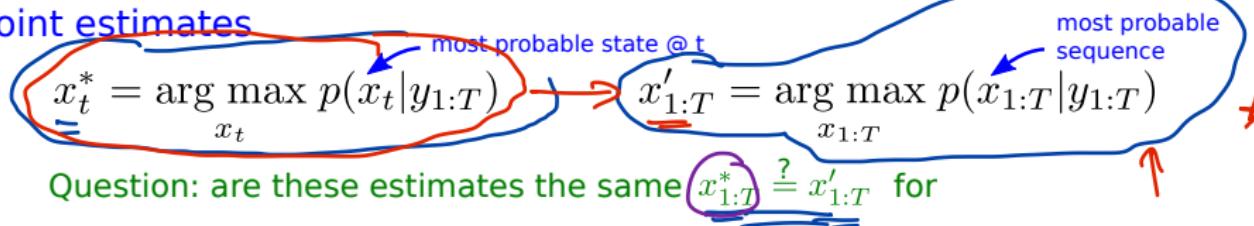
# Varieties of Inference

## Distributional estimates

future data available?

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

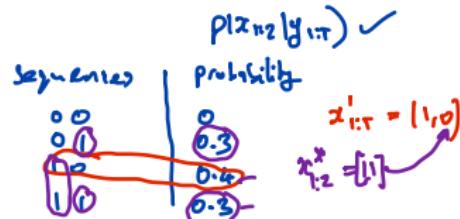
## Point estimates



1. Linear Gaussian State Space Models?

2. Discrete Hidden State HMMs?

$$T=2, K=2$$



$$p(x_{1:T} | y_{1:T}) = N(x_{1:T}; \mu_{1:T}, \Sigma_{1:T, 1:T})$$

$$x_{1:T}^t = \mu_{1:T}$$

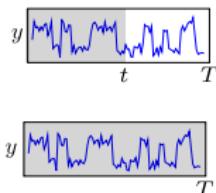
$$p(x_t | y_{1:T}) = \int dx_{1:T \neq t} N(x_{1:T}; \mu_{1:T}, \Sigma_{1:T, 1:T}) = N(x_t; \underline{\mu_t}, \underline{\Sigma_t})$$

$$x_t^* = \underline{\mu_t} \leftarrow x_{1:T}^t = \underline{\mu_{1:T}}$$

## 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

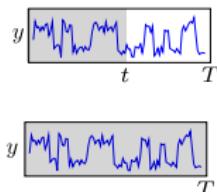
Question: are these estimates the same  $x_{1:T}^* \stackrel{?}{=} x'_{1:T}$  for

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?



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

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

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

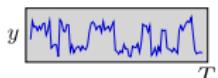
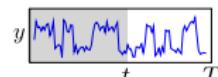
1. Linear Gaussian State Space Models?  $x_{1:T}^* = x'_{1:T}$  (Gaussian)
2. Discrete Hidden State HMMs?  $x_{1:T}^* \neq x'_{1:T}$

# Varieties of Inference

## Distributional estimates

future data available?

look at this next



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

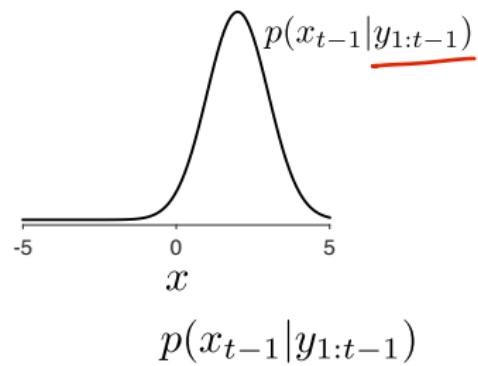
$$x_t^* = \arg \max_{x_t} p(x_t|y_{1:T}) \quad \text{most probable state @ t}$$
$$x'_{1:T} = \arg \max_{x_{1:T}} p(x_{1:T}|y_{1:T}) \quad \text{most probable sequence}$$

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

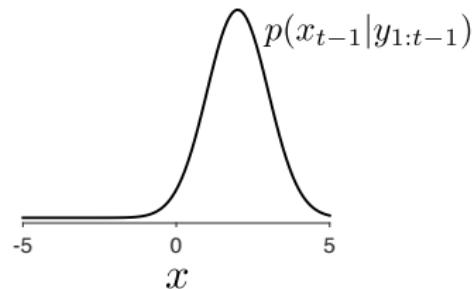
1. Linear Gaussian State Space Models?  $x_{1:T}^* = x'_{1:T}$  (Gaussian)

2. Discrete Hidden State HMMs?  $x_{1:T}^* \neq x'_{1:T}$

## Inference: Kalman Filter



## Inference: Kalman Filter

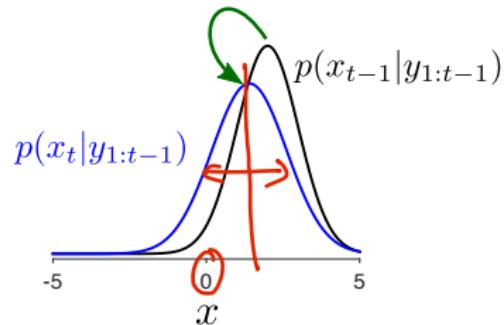


diffuse via dynamics

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

sum for discrete hidden state

## Inference: Kalman Filter

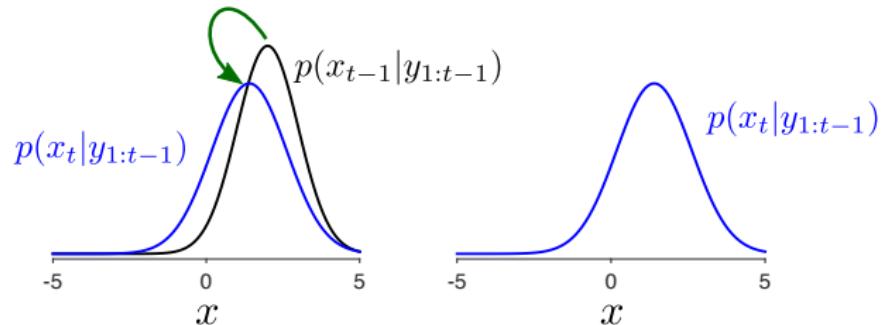


diffuse via dynamics

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

sum for discrete hidden state

## Inference: Kalman Filter



diffuse via dynamics

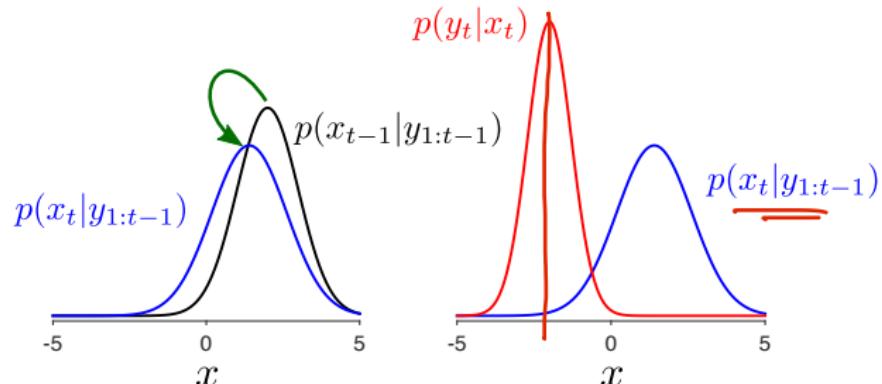
combine with likelihood

sum for discrete hidden state

$$p(x_t | y_{1:t-1}) = \int p(x_t | x_{t-1}) p(x_{t-1} | y_{1:t-1}) dx_{t-1}$$
$$\underline{p(x_t | y_{1:t})} \propto \underline{\text{prior}} \underline{\text{likelihood}}$$

Bayes' Rule

## Inference: Kalman Filter



diffuse via dynamics  
combine with likelihood

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

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

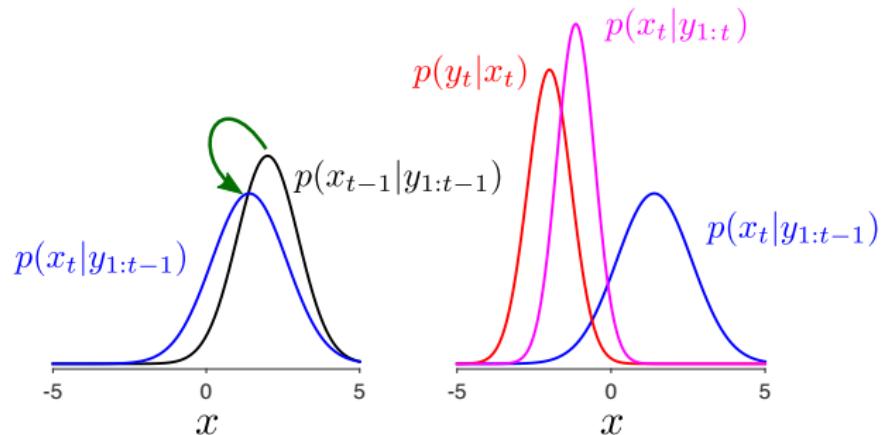
$p(x_t|y_{1:t}) \propto p(x_t|y_{1:t-1})p(y_t|x_t)$

sum for discrete hidden state

Bayes' Rule

prior      likelihood

## Inference: Kalman Filter



diffuse via dynamics  
combine with likelihood

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

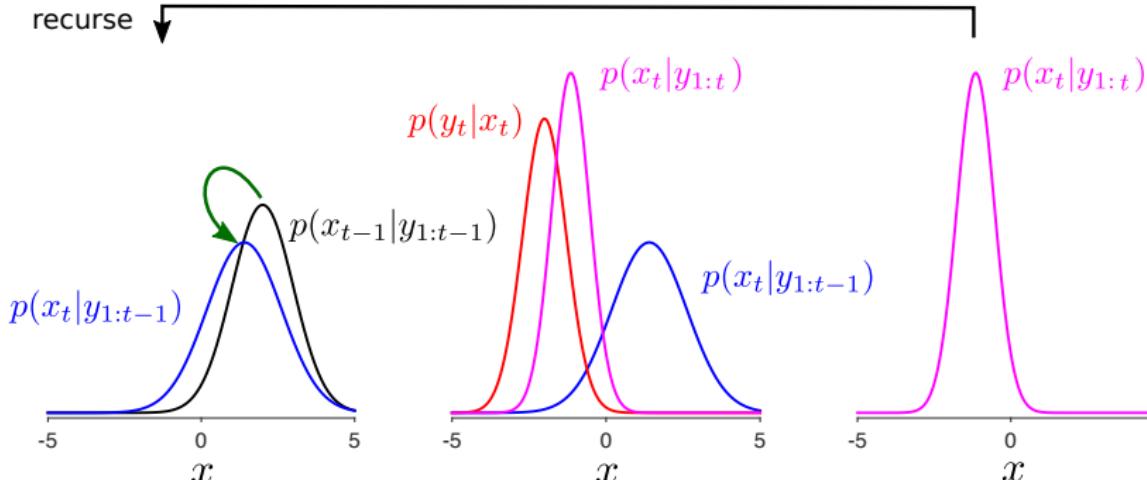
sum for discrete hidden state

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

prior      likelihood

Bayes' Rule

## Inference: Kalman Filter



diffuse via dynamics

combine with likelihood

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

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

sum for discrete hidden state

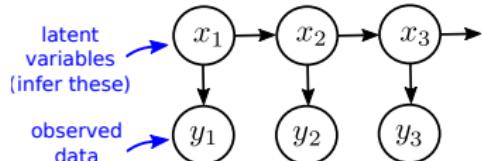
$p(x_t|y_{1:t}) \propto p(x_t|y_{1:t-1})p(y_t|x_t)$

prior      likelihood

Bayes' Rule

## Inference: Derivation of General Filtering Equations

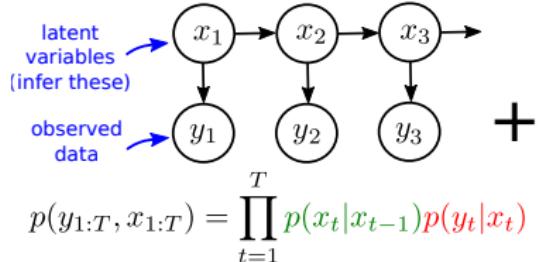
Model



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

# Inference: Derivation of General Filtering Equations

Model



Rules of probability

product rule

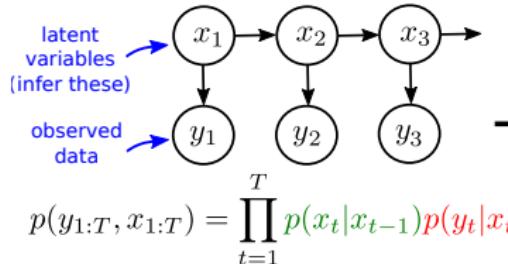
$$p(A|B, C) = \frac{1}{p(B|C)} p(B|A, C)p(A|C)$$

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

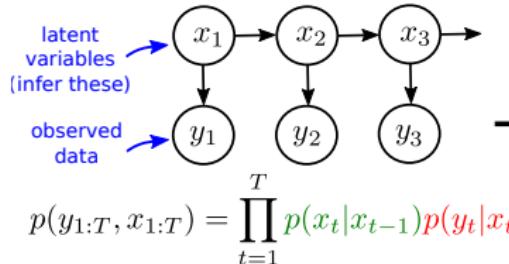
$$p(A|C) = \sum_B p(A, B|C)$$

Inference

= ?

# Inference: Derivation of General Filtering Equations

Model



Rules of probability

product rule

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

sum rule

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

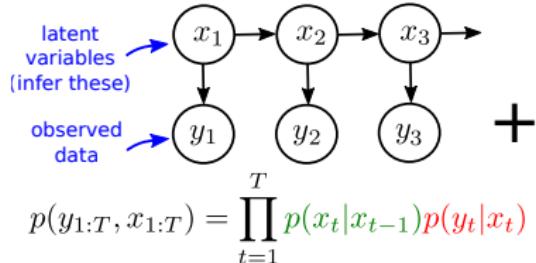
Inference

= ?

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

# 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

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

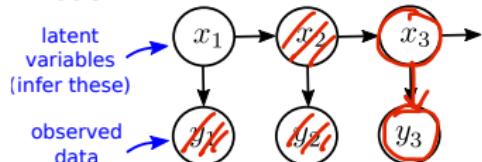
Inference

= ?

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

# Inference: Derivation of General Filtering Equations

Model



Rules of probability

product rule

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

sum rule

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

Inference

= ?

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

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

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

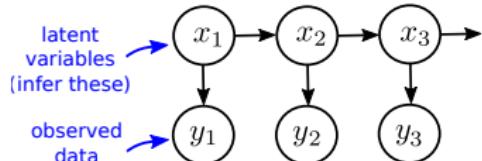
$$= p(y_t|x_t)$$

product rule

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

# 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

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

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

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

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

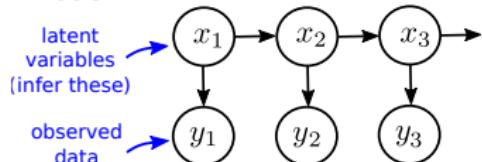
product rule  
 $A = x_t$   $B = y_t$   $C = y_{1:t-1}$

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

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

# 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

= ?

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

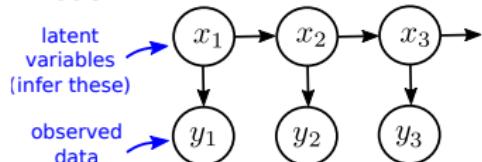
$$= \frac{1}{p(y_t|y_{1:t-1})} p(y_t|x_t, y_{1:t-1})p(x_t|y_{1:t-1}) \quad \begin{matrix} \text{product rule} \\ A = x_t \ B = y_t \ C = y_{1:t-1} \end{matrix}$$

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

$$\propto p(y_t|x_t)p(x_t|y_{1:t-1}) \quad \begin{matrix} \text{constant of proportionality} \\ p(y_t|y_{1:t-1}) \ (\text{see learning}) \end{matrix}$$

# 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

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

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

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

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

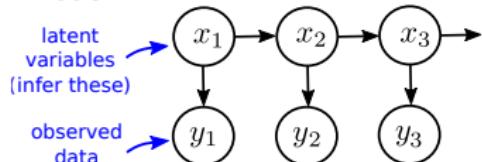
$$\propto p(y_t|x_t)p(x_t|y_{1:t-1})$$

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

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

# 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

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

Inference

= ?

sum rule

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

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

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

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

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

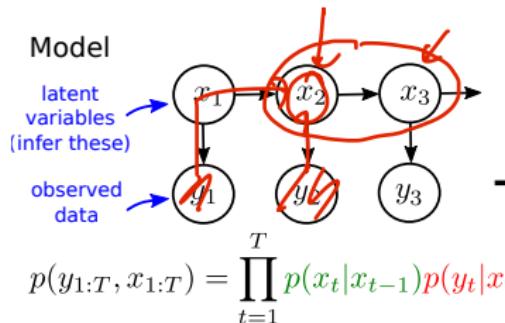
constant of proportionality  $p(y_t|y_{1:t-1})$  (see learning)

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

sum rule

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

# Inference: Derivation of General Filtering Equations



**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)

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

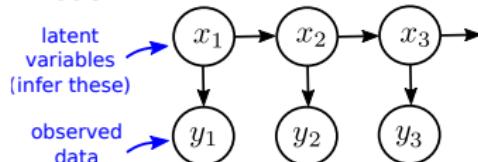
$$\quad \begin{matrix} \text{sum rule} \\ A = x_t \ B = x_{t-1} \ C = y_{1:t-1} \end{matrix}$$

$$= \int p(x_t|x_{t-1}, y_{1:t-1})p(x_{t-1}|y_{1:t-1})dx_{t-1} \quad \begin{matrix} \text{product rule} \\ \underbrace{\hspace{10em}}_{\parallel} \end{matrix}$$

$$p(x_t|x_{t-1})$$

# 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)

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

$$\quad \begin{matrix} \text{sum rule} \\ A = x_t \ B = x_{t-1} \ C = y_{1:t-1} \end{matrix}$$

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

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

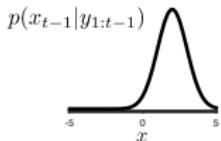
conditional independence from model

## Inference: Kalman Filter

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

diffuse via  
dynamics

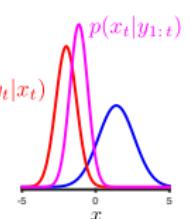
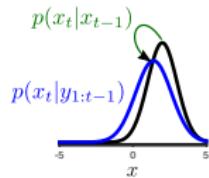
$$p(x_t|y_{1:t-1}) = \int p(x_t|x_{t-1})p(x_{t-1}|y_{1:t-1})dx_{t-1}$$



combine  
with  
likelihood

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

prior              likelihood



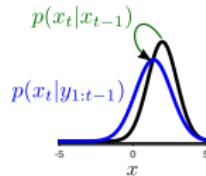
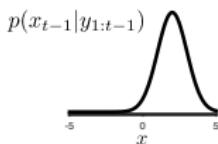
## Inference: Kalman Filter

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

most recent data used in prediction  
variable being predicted

diffuse via dynamics

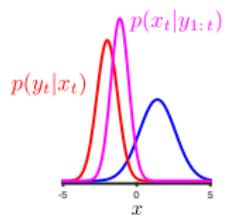
$$p(x_t|y_{1:t-1}) = \int p(x_t|x_{t-1})p(x_{t-1}|y_{1:t-1})dx_{t-1}$$



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

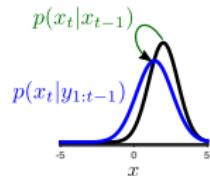
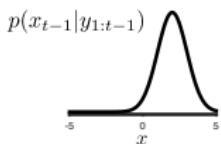
$$p(x_{t-1}|y_{1:t-1}) = \mathcal{G}(x_{t-1}; \mu_{t-1}^{t-1}, V_{t-1}^{t-1})$$

most recent data used in prediction  
variable being predicted

diffuse via dynamics

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

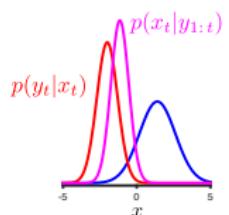
$$p(x_t|y_{1:t-1}) = \mathcal{G}(x_t; \mu_t^{t-1}, V_t^{t-1})$$



combine with likelihood

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

prior                      likelihood



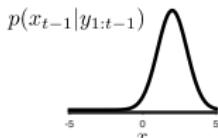
## Inference: Kalman Filter

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

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

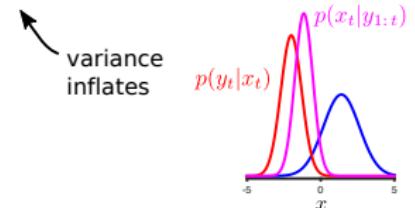
$$p(x_t|y_{1:t-1}) = \mathcal{G}(x_t; \mu_t^{t-1}, V_t^{t-1}) \quad \mu_t^{t-1} = A\mu_{t-1}^{t-1}$$

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

combine  
with  
likelihood

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

prior              likelihood



## 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

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

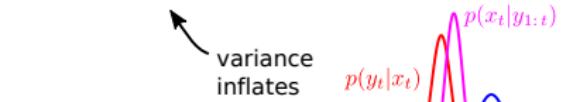
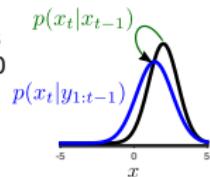
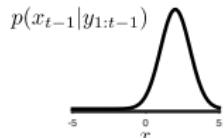
$$V_t^{t-1} = AV_{t-1}^{t-1}A^\top + Q$$

combine  
with  
likelihood

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

prior                      likelihood

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



# 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

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

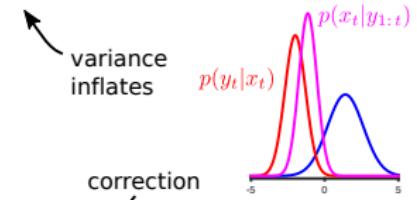
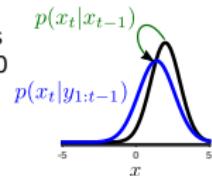
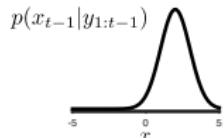
prior              likelihood

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

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

$$V_t^t = V_t^{t-1} - K_t C V_t^{t-1}$$

Kalman gain →  $K_t = V_t^{t-1}C^\top(CV_t^{t-1}C^\top + R)^{-1}$

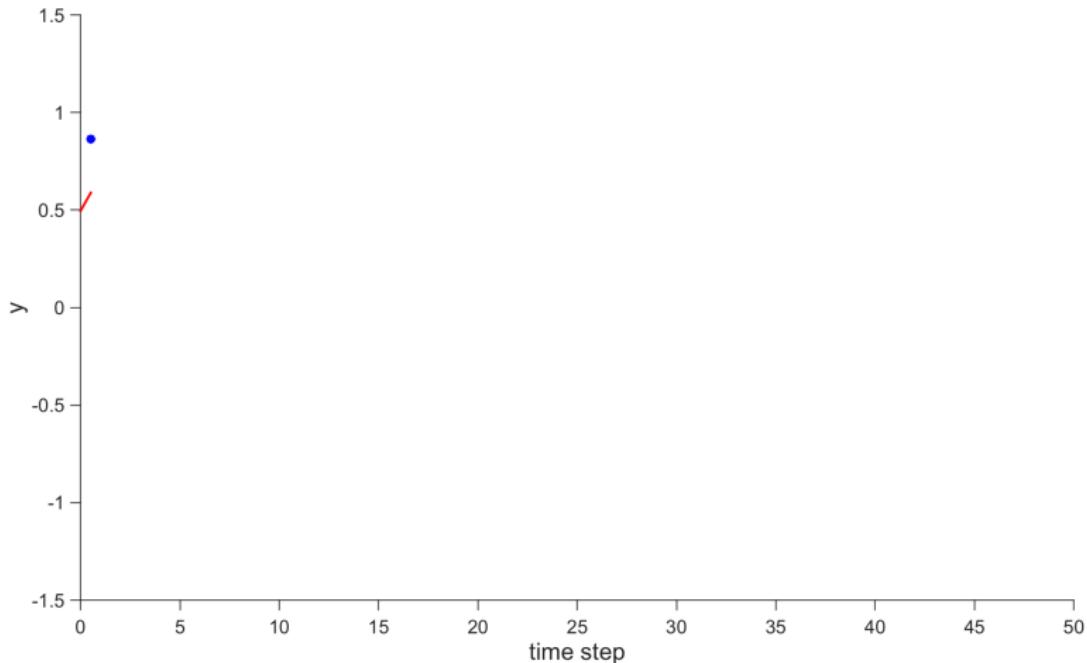


## Kalman Filter Demo

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

## Kalman Filter Demo

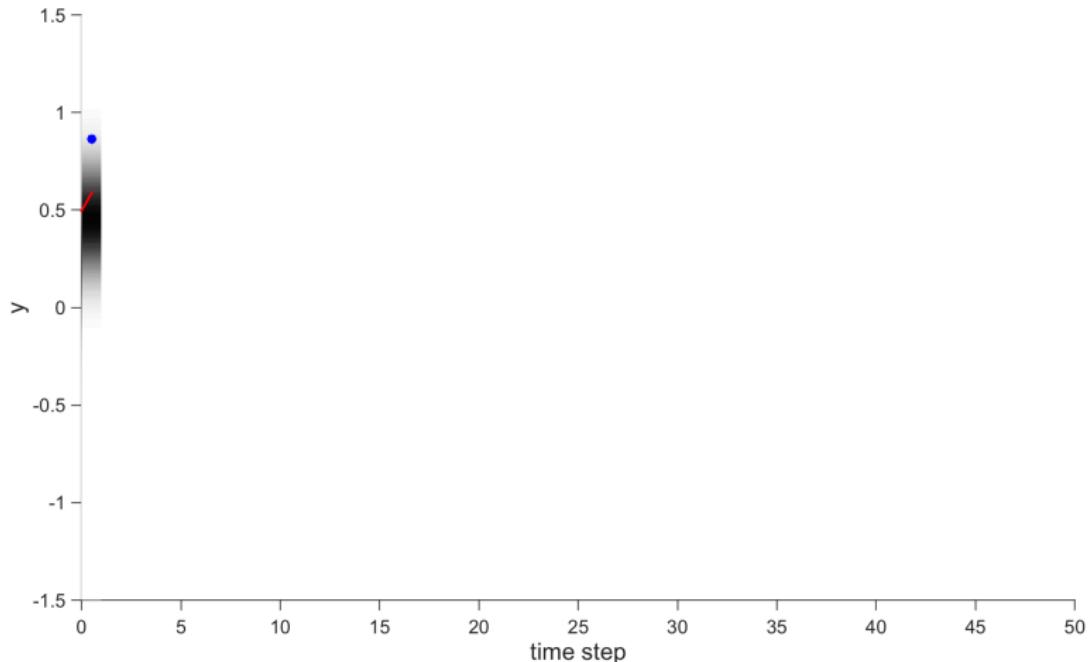
observed noisy data  $y_t$ , ground truth sinusoid



observe first data point  $y_1$

## Kalman Filter Demo

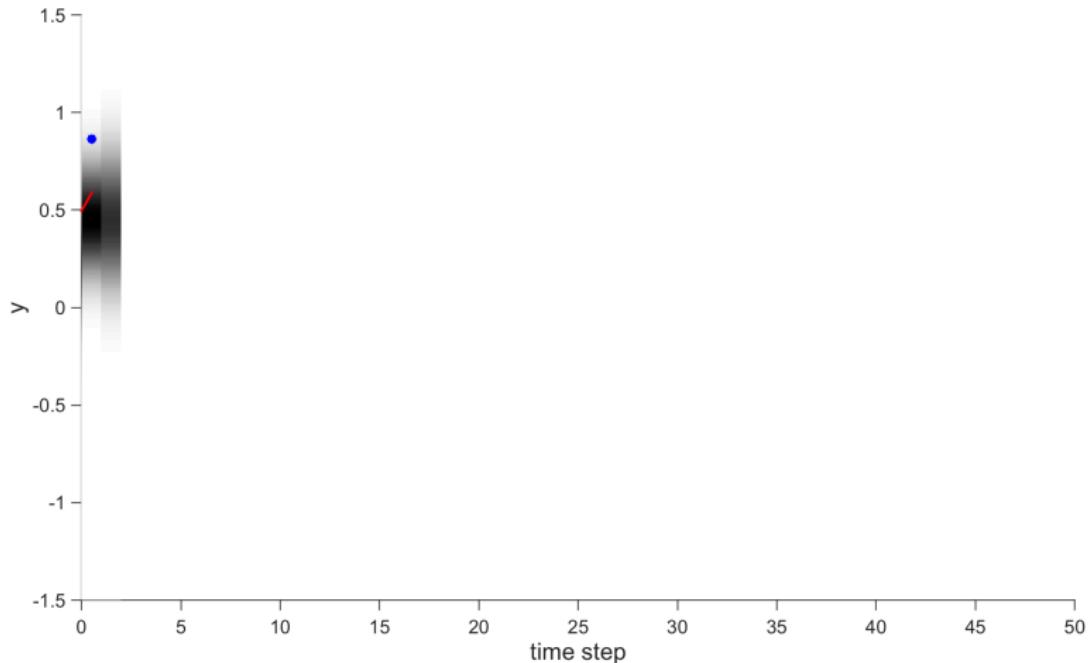
observed noisy data  $y_t$ , ground truth sinusoid



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

## Kalman Filter Demo

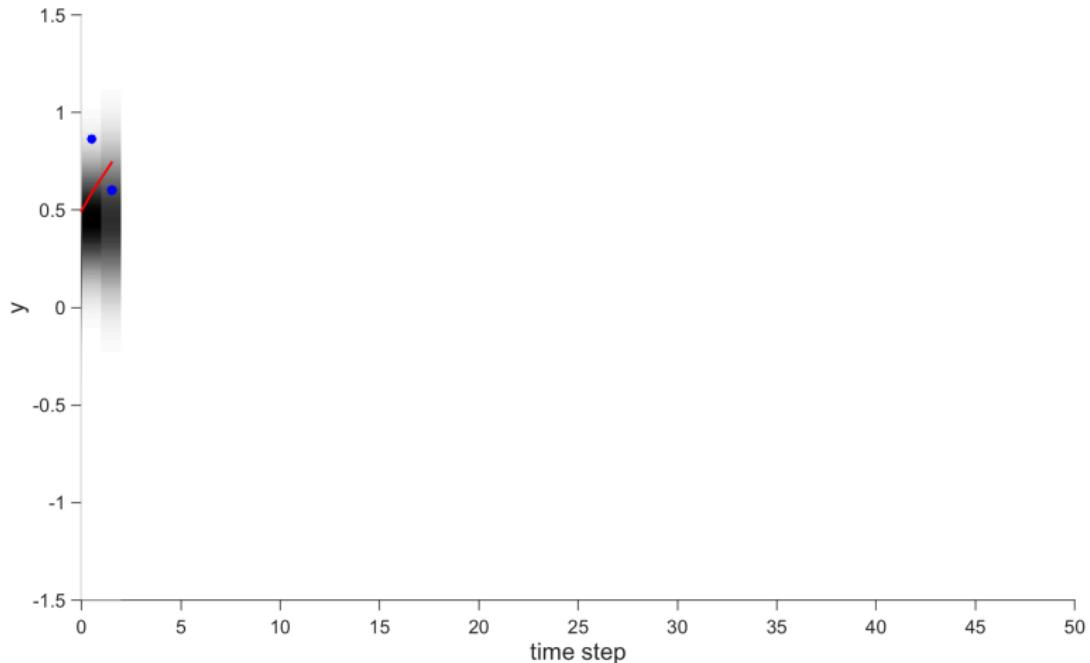
observed noisy data  $y_t$ , ground truth sinusoid



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

## Kalman Filter Demo

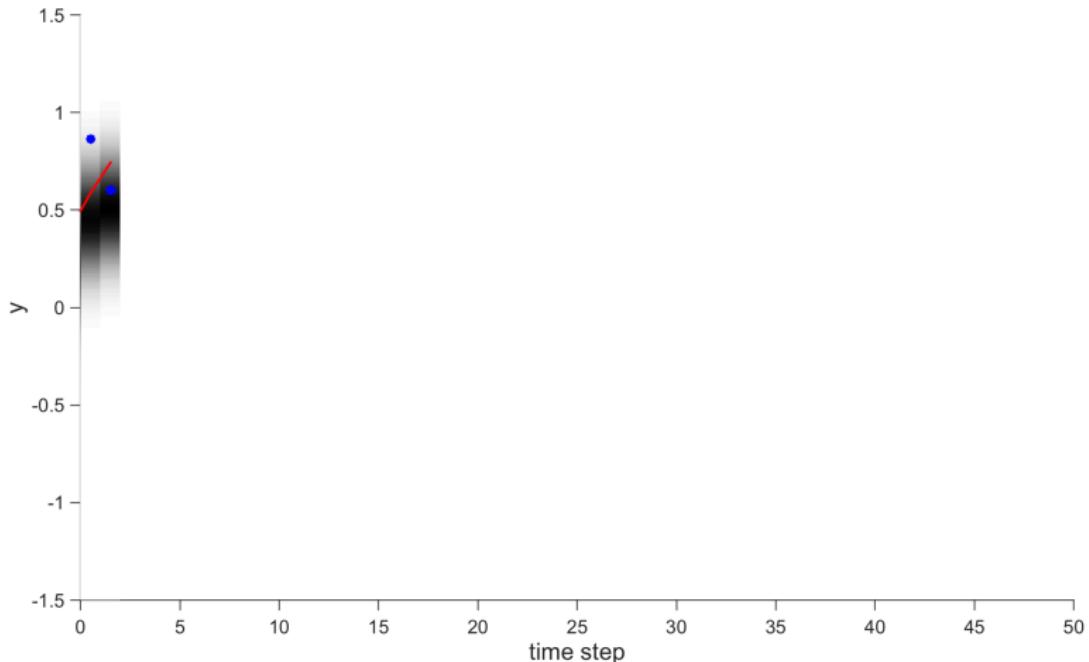
observed noisy data  $y_t$ , ground truth sinusoid



observe next data point  $y_2$

## Kalman Filter Demo

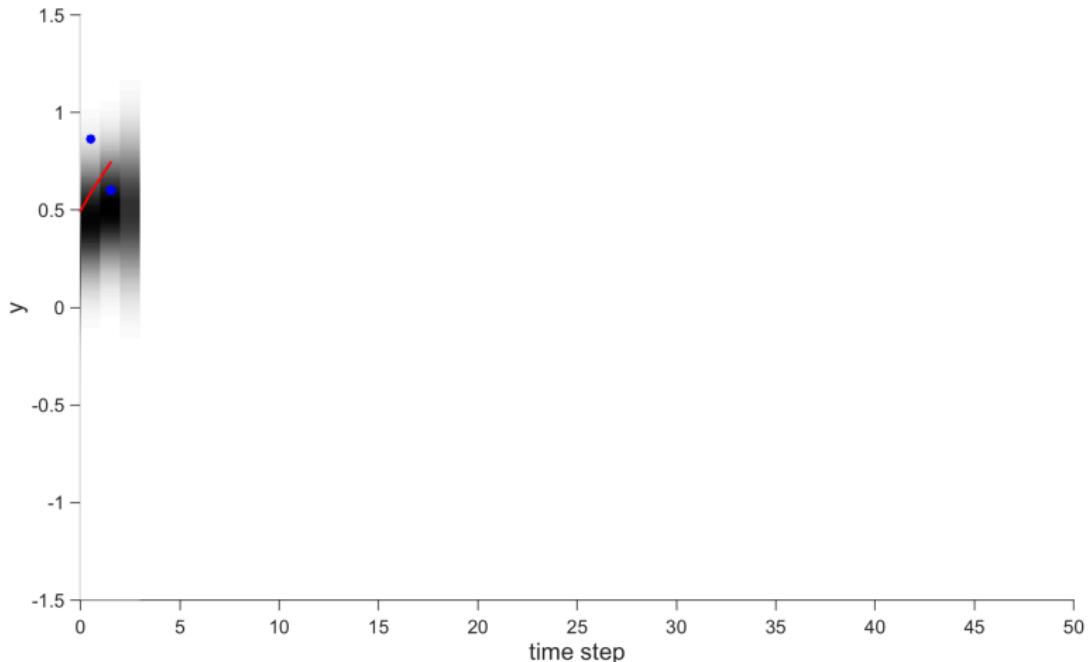
observed noisy data  $y_t$ , ground truth sinusoid



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

## Kalman Filter Demo

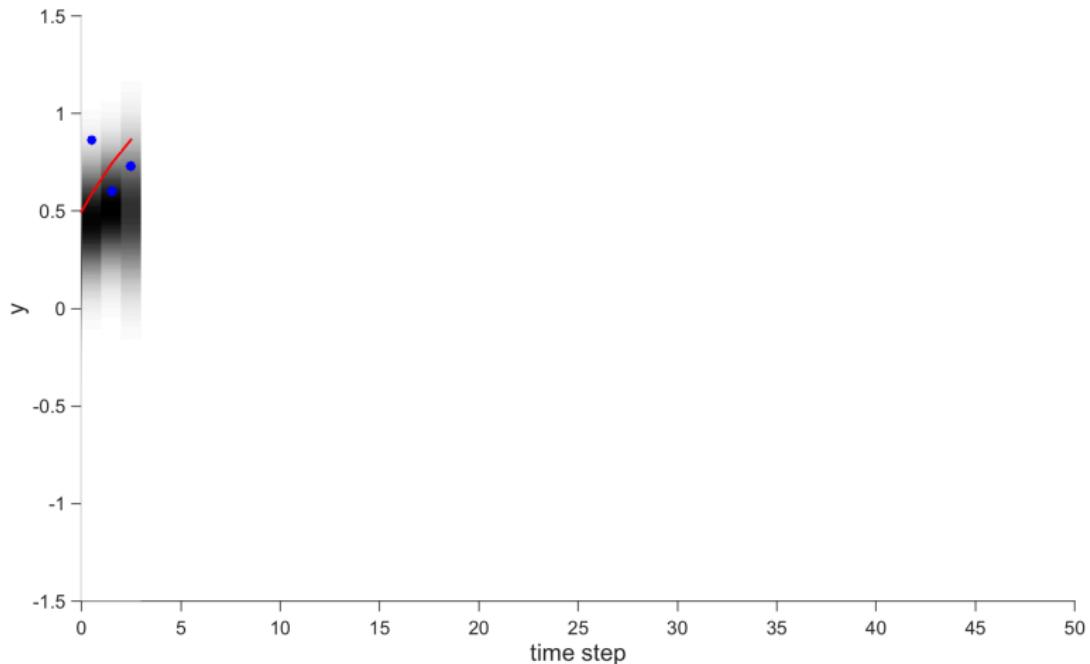
observed noisy data  $y_t$ , ground truth sinusoid



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

## Kalman Filter Demo

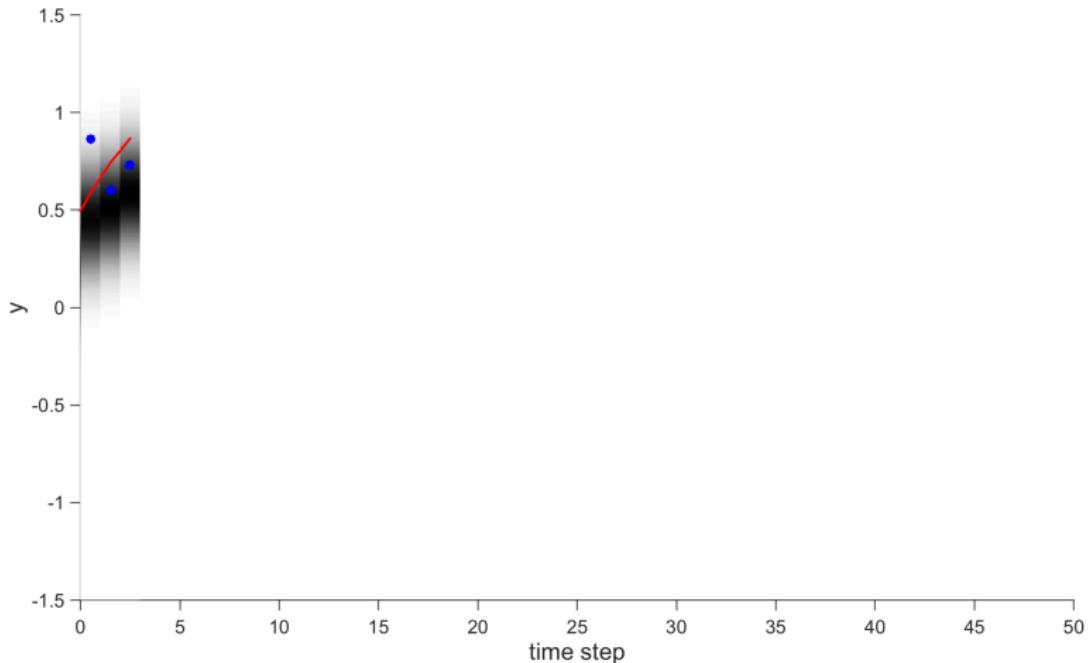
observed noisy data  $y_t$ , ground truth sinusoid



observe next data point  $y_3$

## Kalman Filter Demo

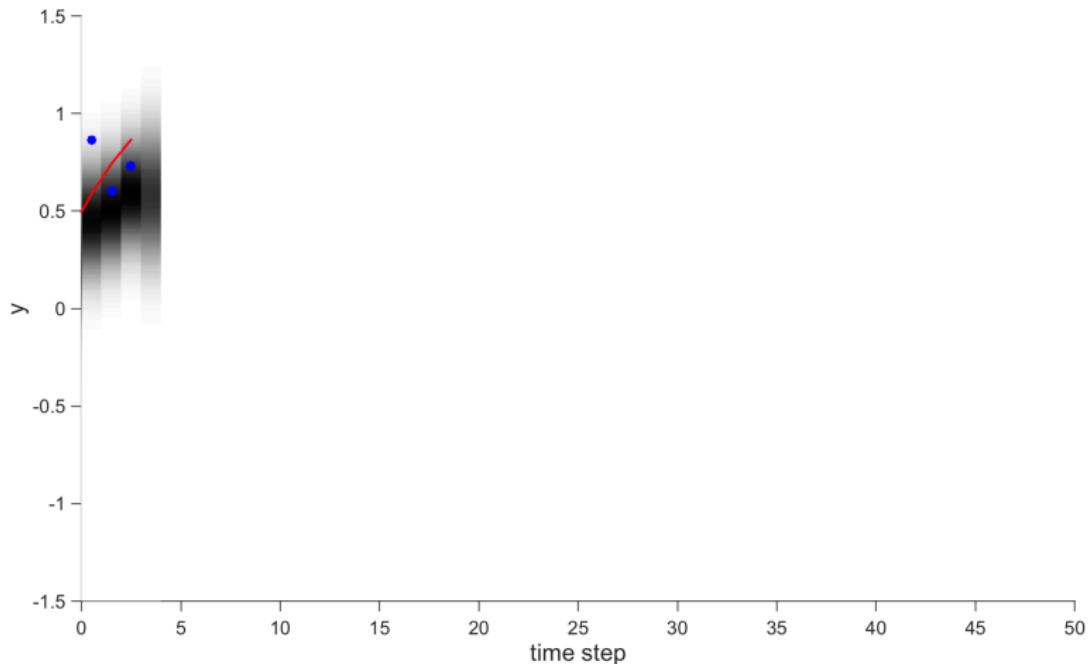
observed noisy data  $y_t$ , ground truth sinusoid



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

## Kalman Filter Demo

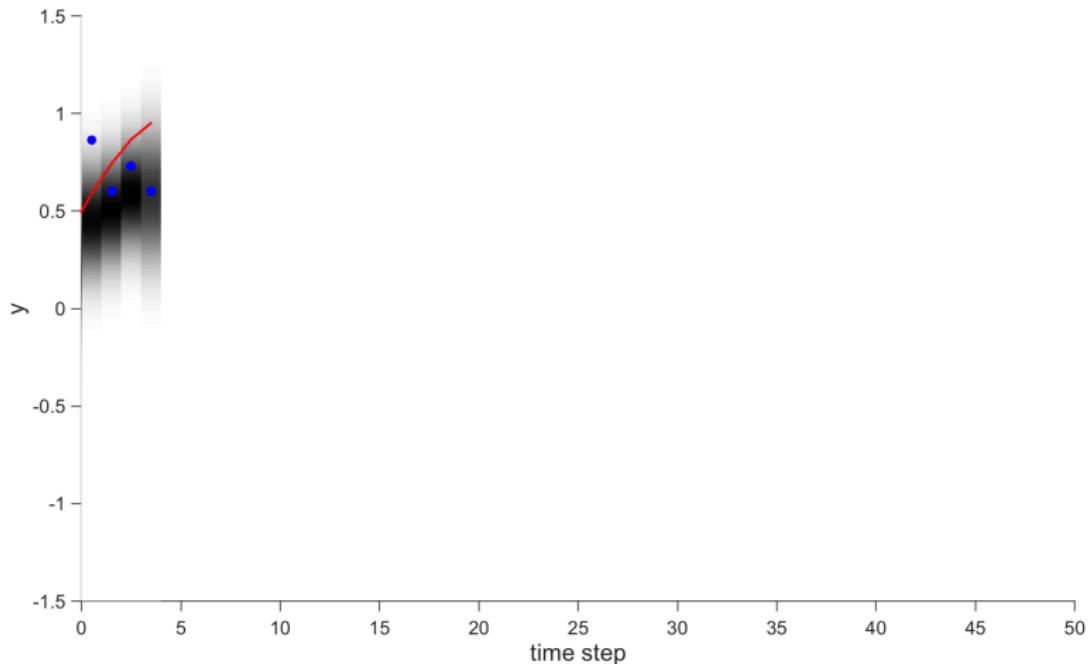
observed noisy data  $y_t$ , ground truth sinusoid



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

## Kalman Filter Demo

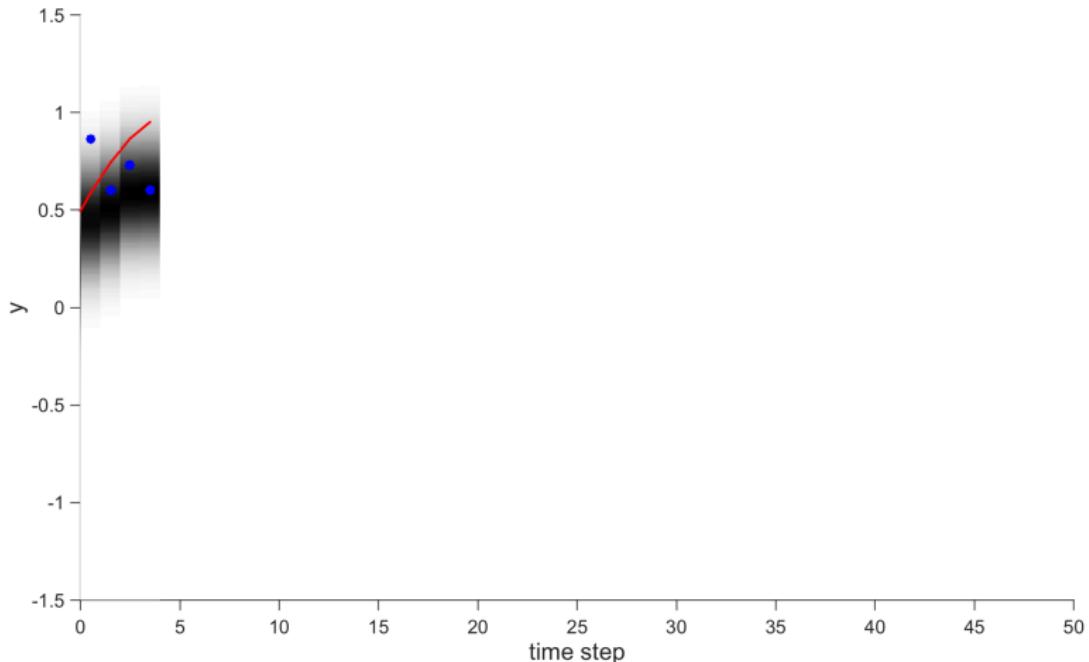
observed noisy data  $y_t$ , ground truth sinusoid



observe next data point  $y_4$

## Kalman Filter Demo

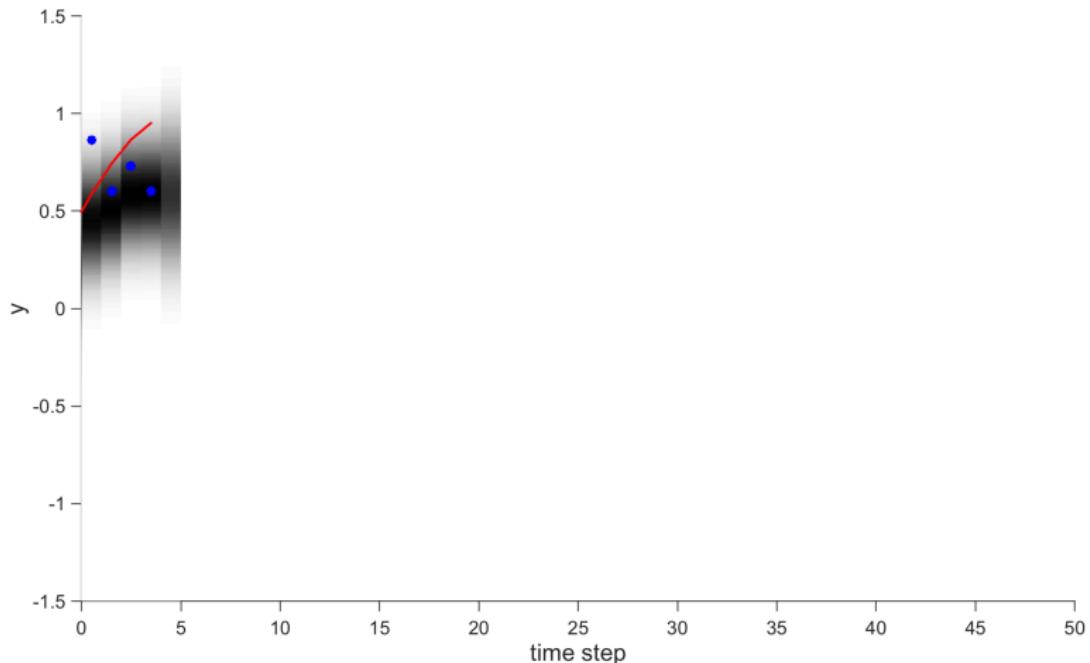
observed noisy data  $y_t$ , ground truth sinusoid



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

## Kalman Filter Demo

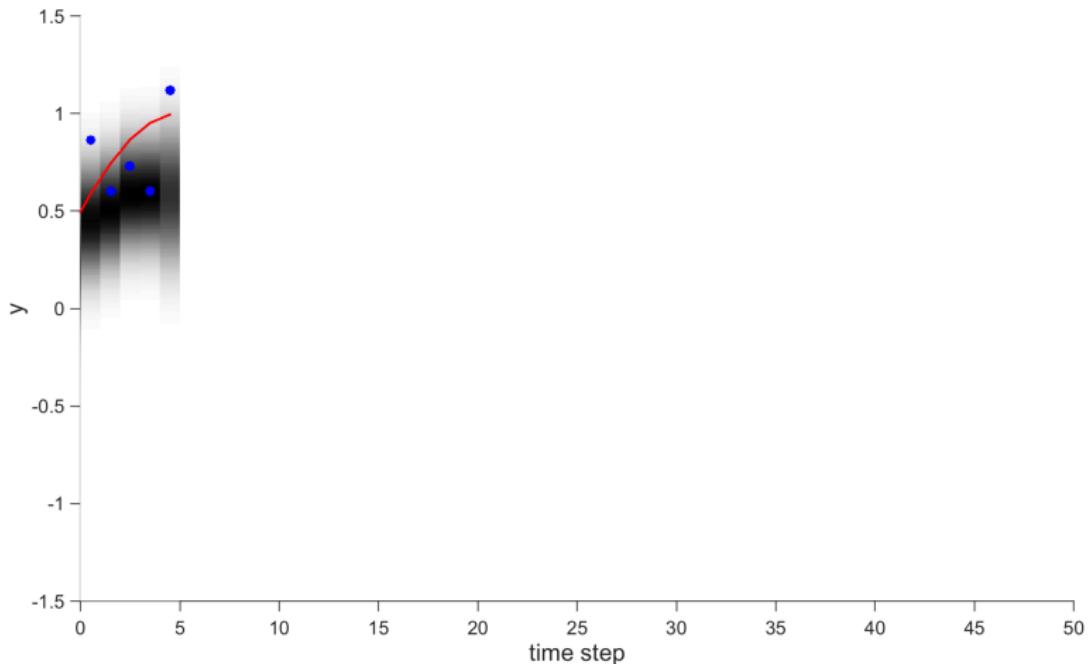
observed noisy data  $y_t$ , ground truth sinusoid



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

## Kalman Filter Demo

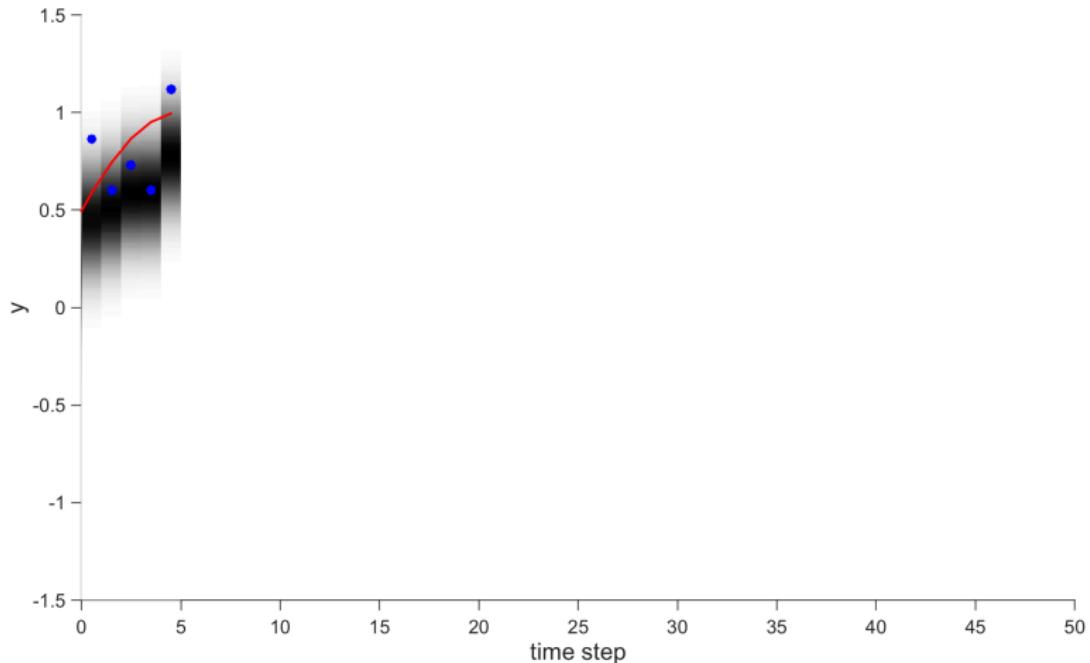
observed noisy data  $y_t$ , ground truth sinusoid



observe next data point  $y_5$

## Kalman Filter Demo

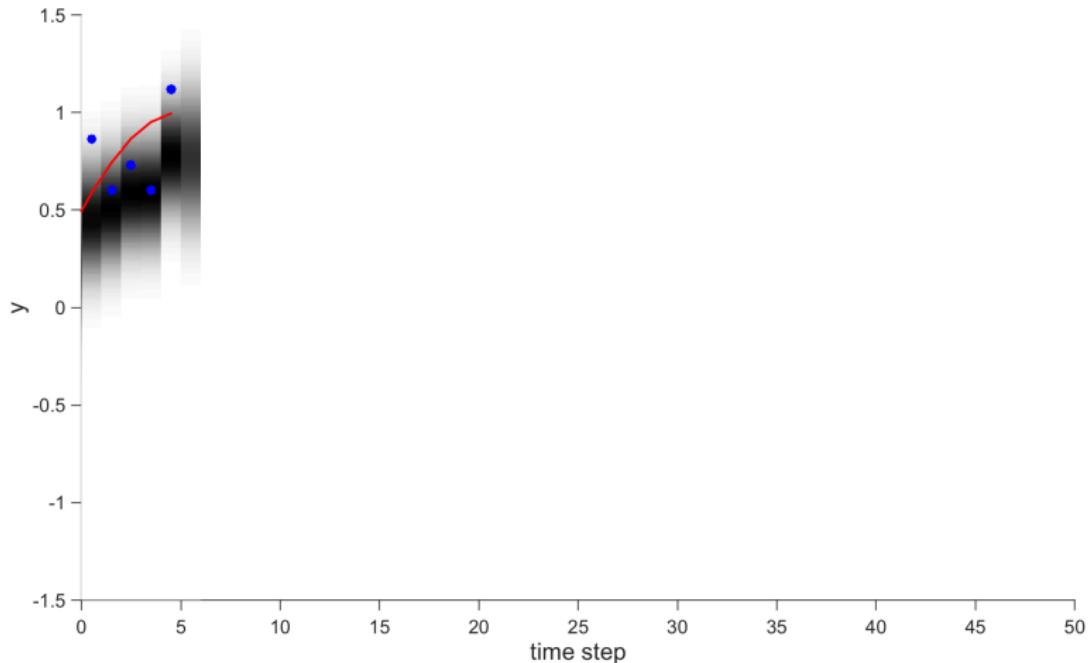
observed noisy data  $y_t$ , ground truth sinusoid



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

## Kalman Filter Demo

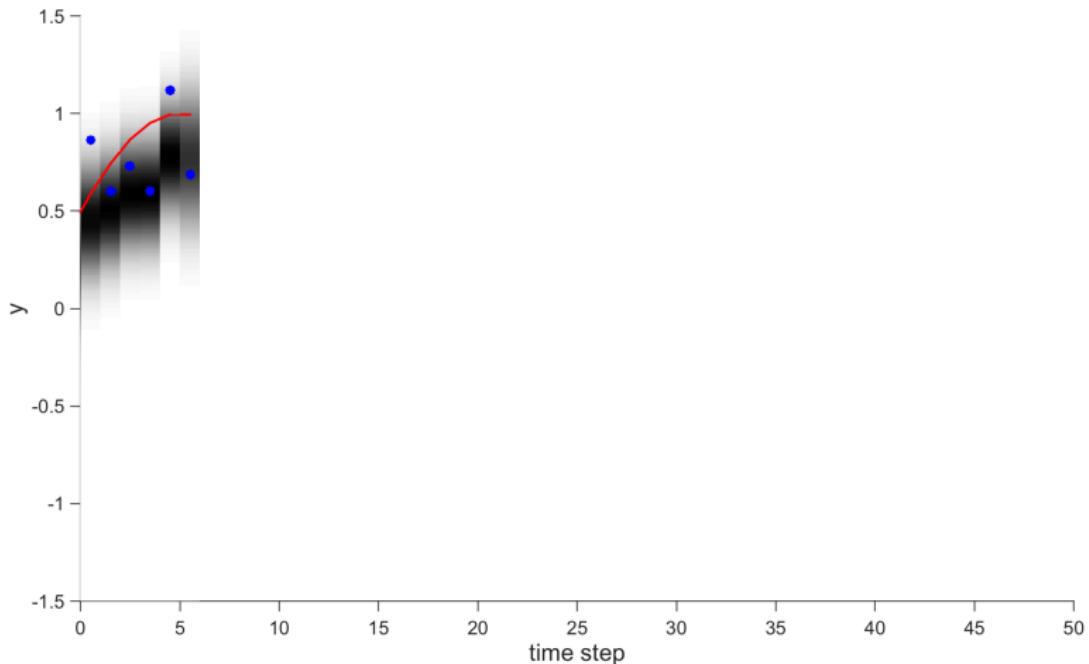
observed noisy data  $y_t$ , ground truth sinusoid



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

## Kalman Filter Demo

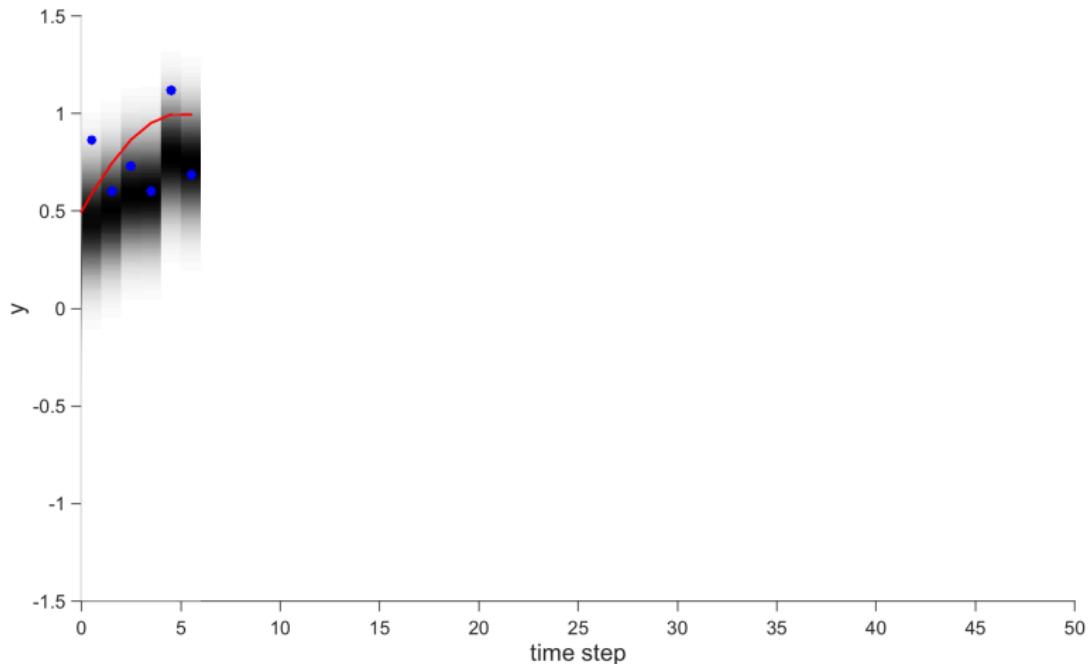
observed noisy data  $y_t$ , ground truth sinusoid



observe next data point  $y_6$

## Kalman Filter Demo

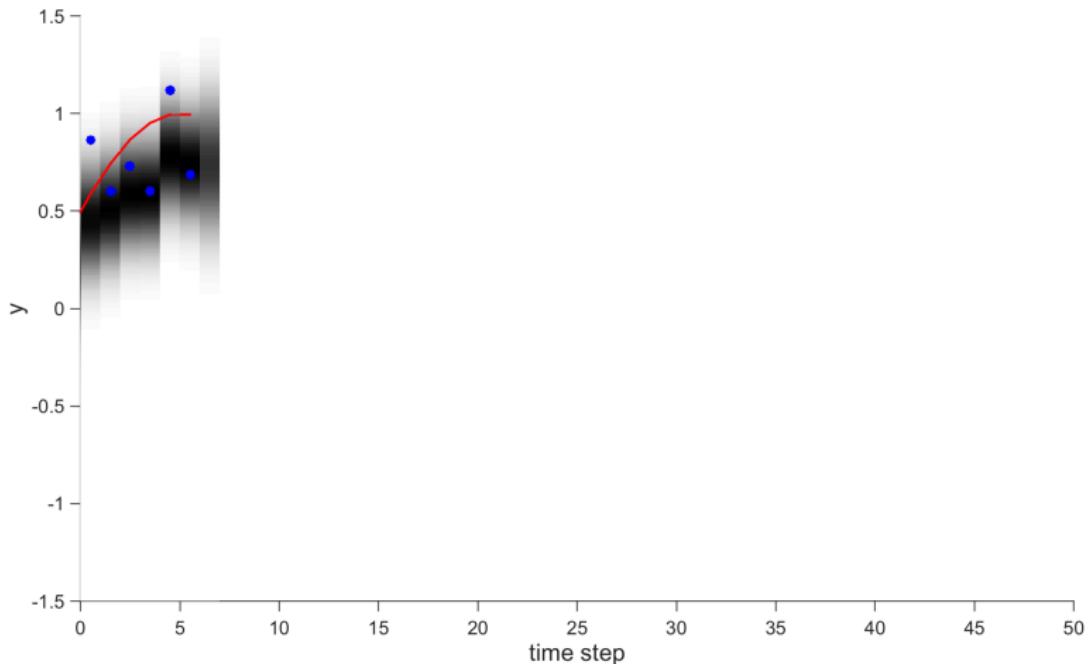
observed noisy data  $y_t$ , ground truth sinusoid



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

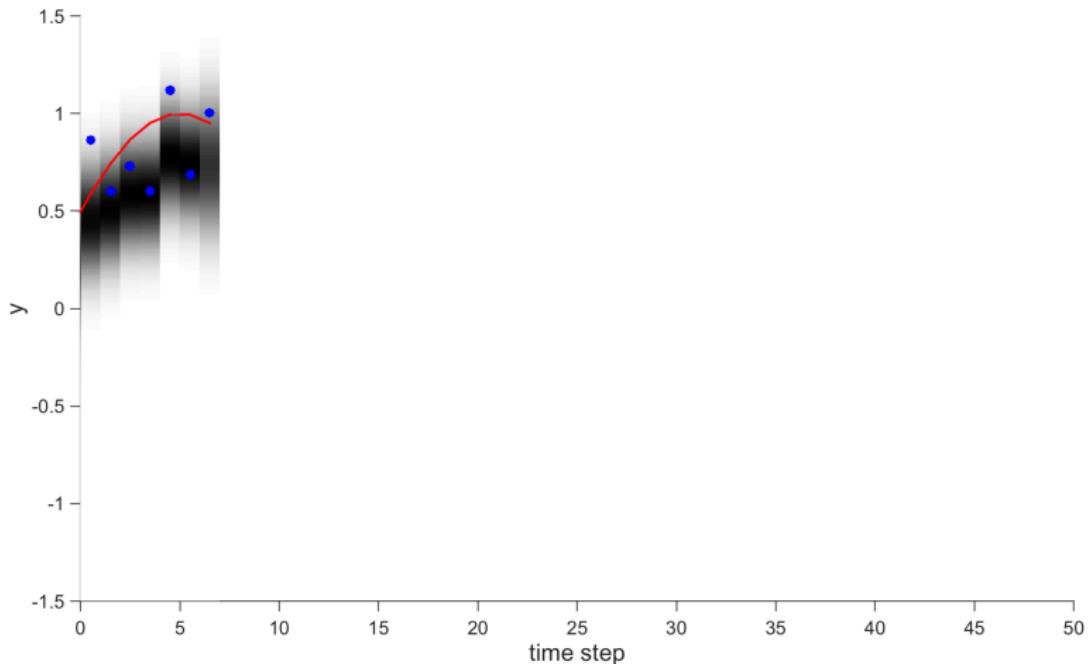
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



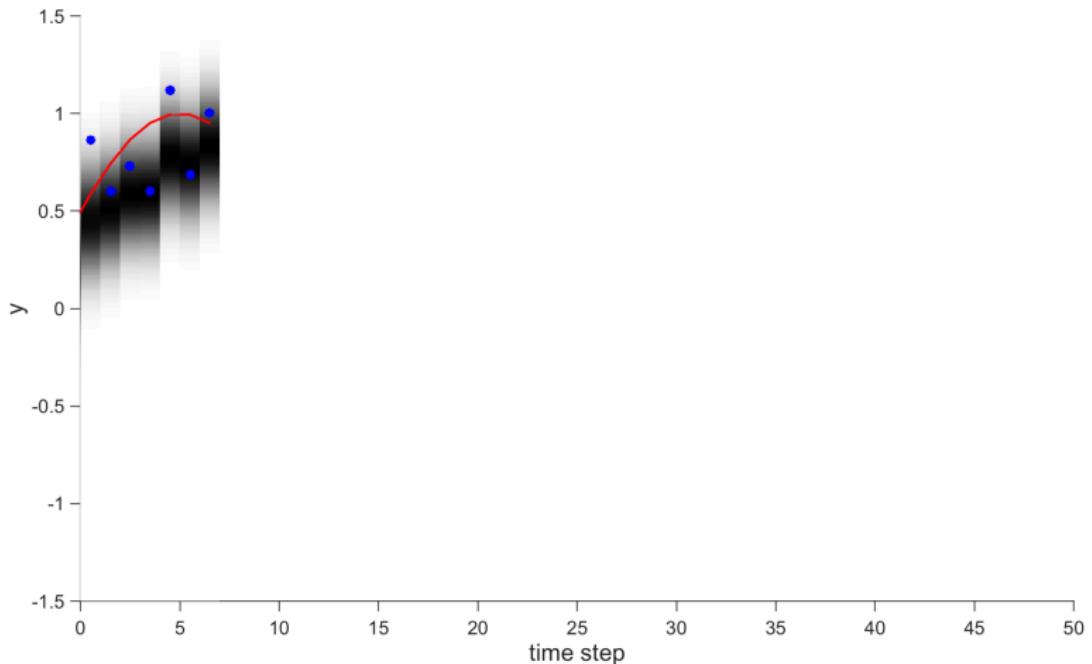
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



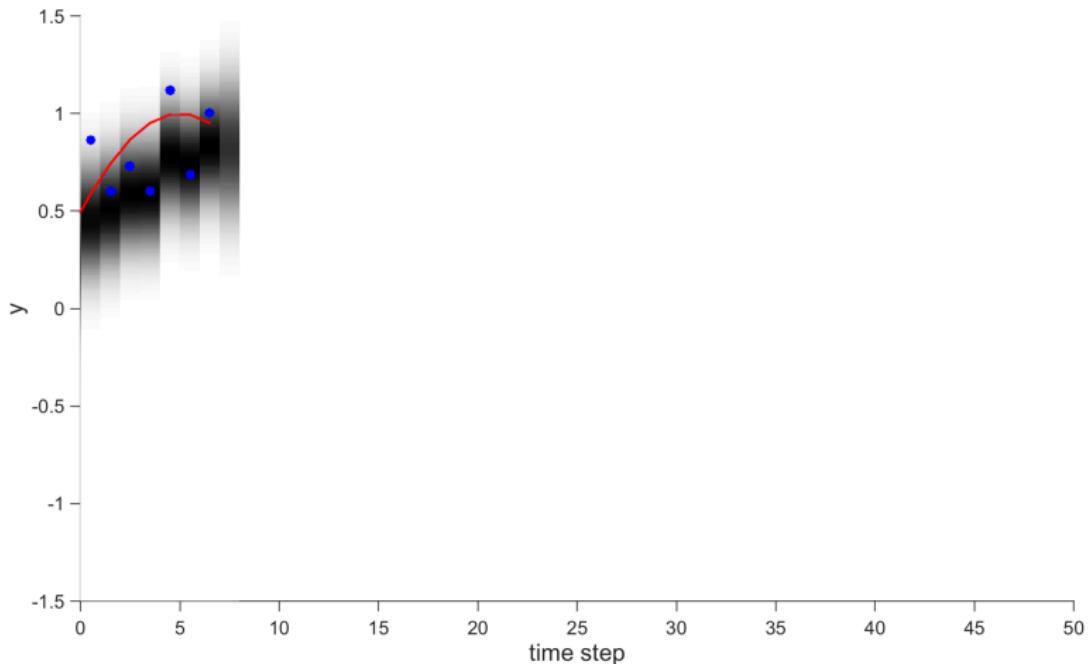
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



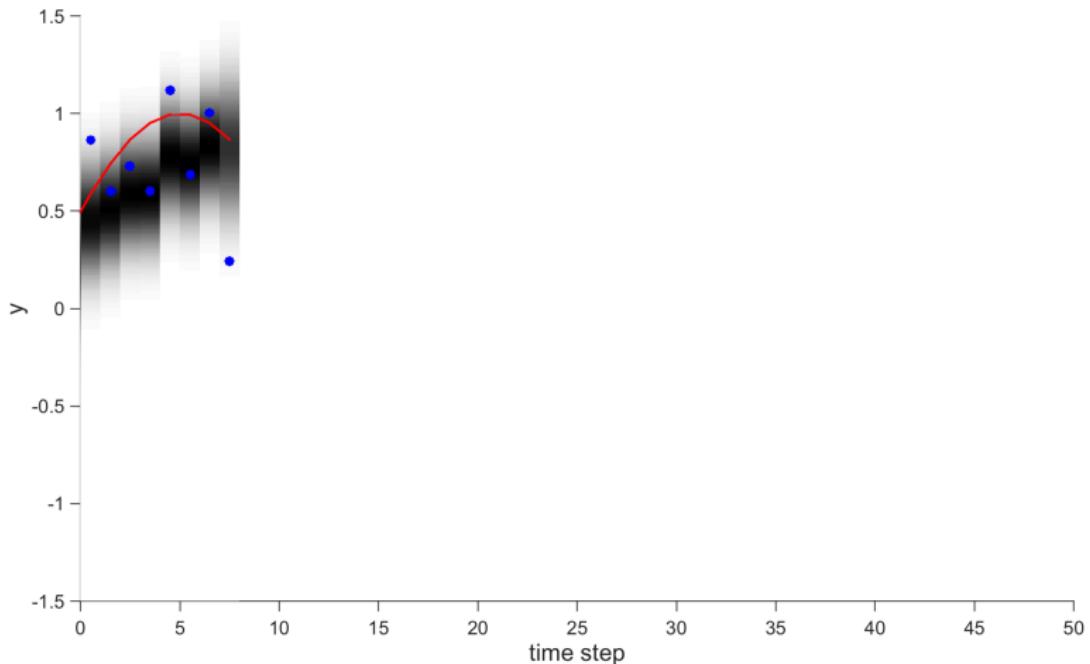
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



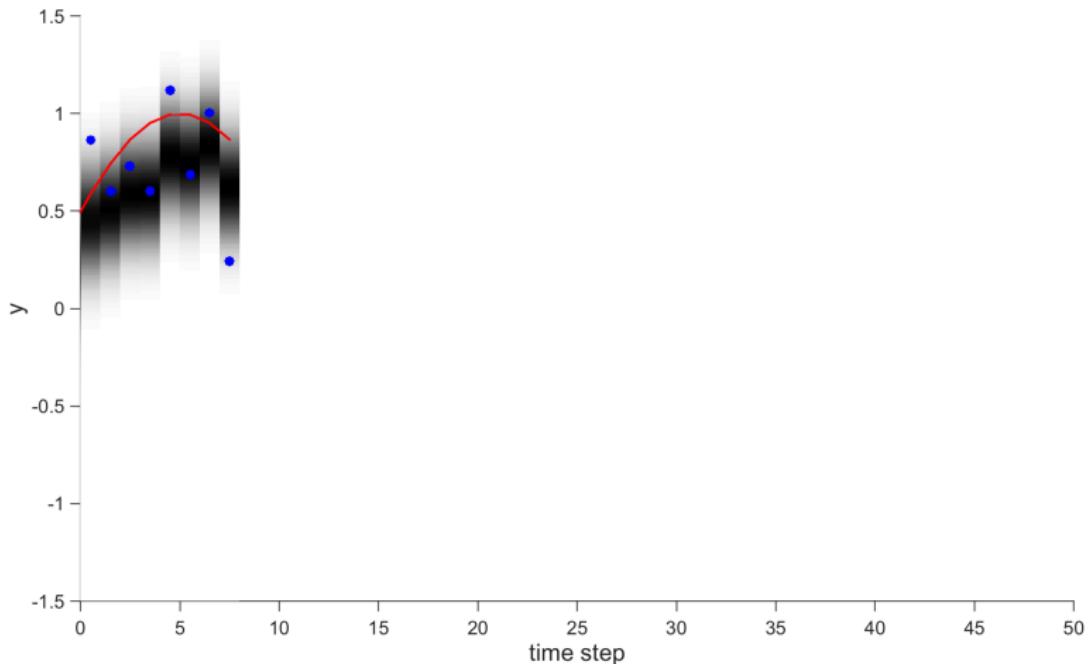
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



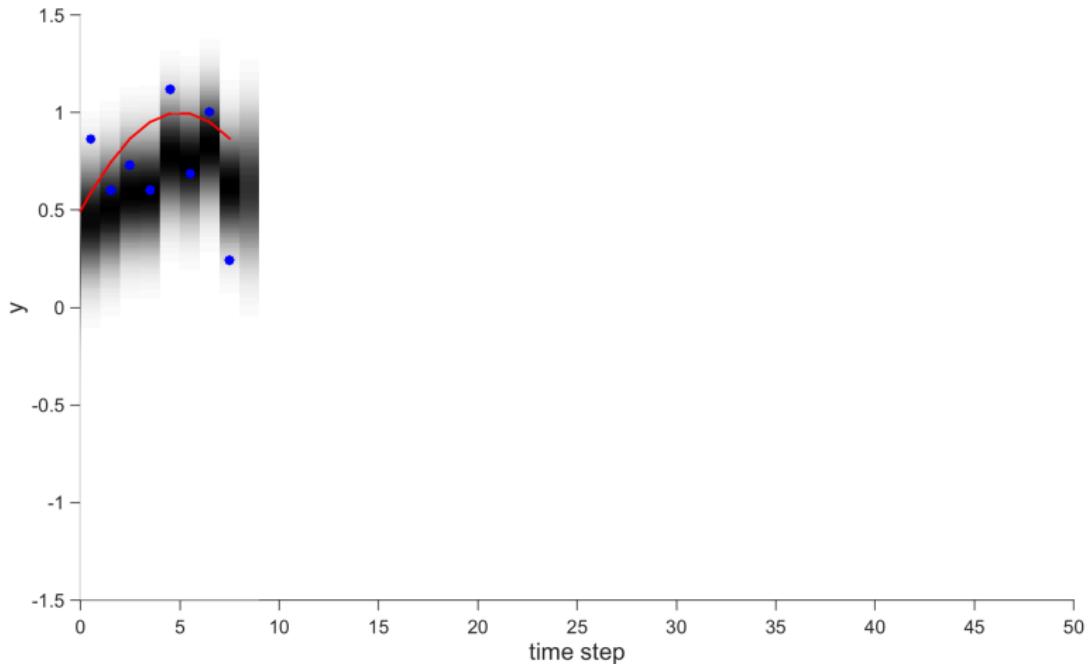
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



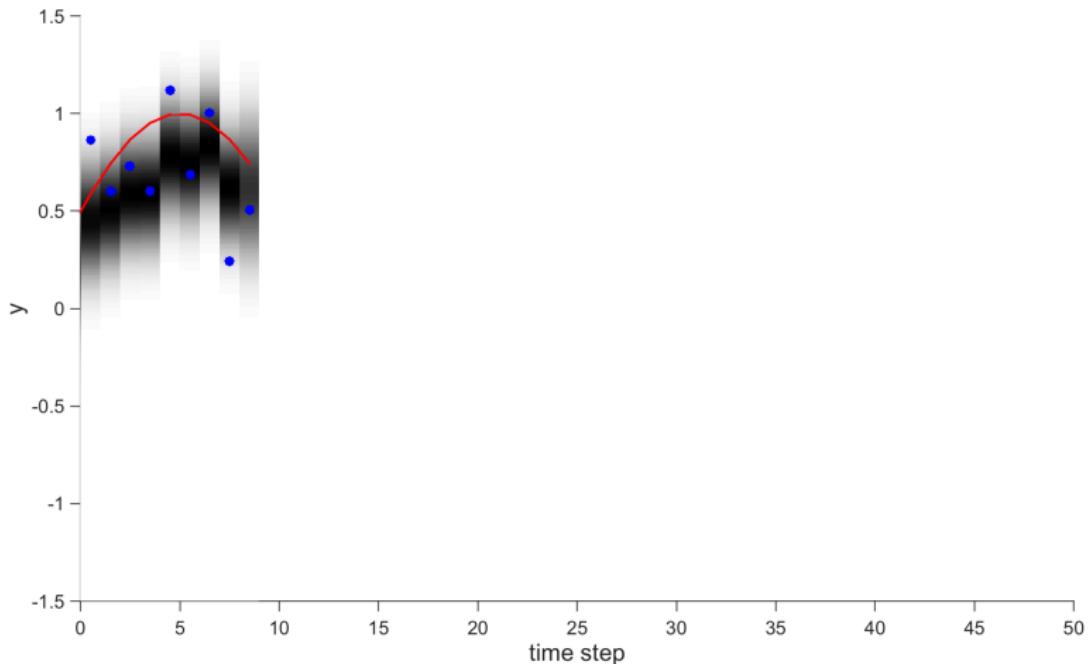
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



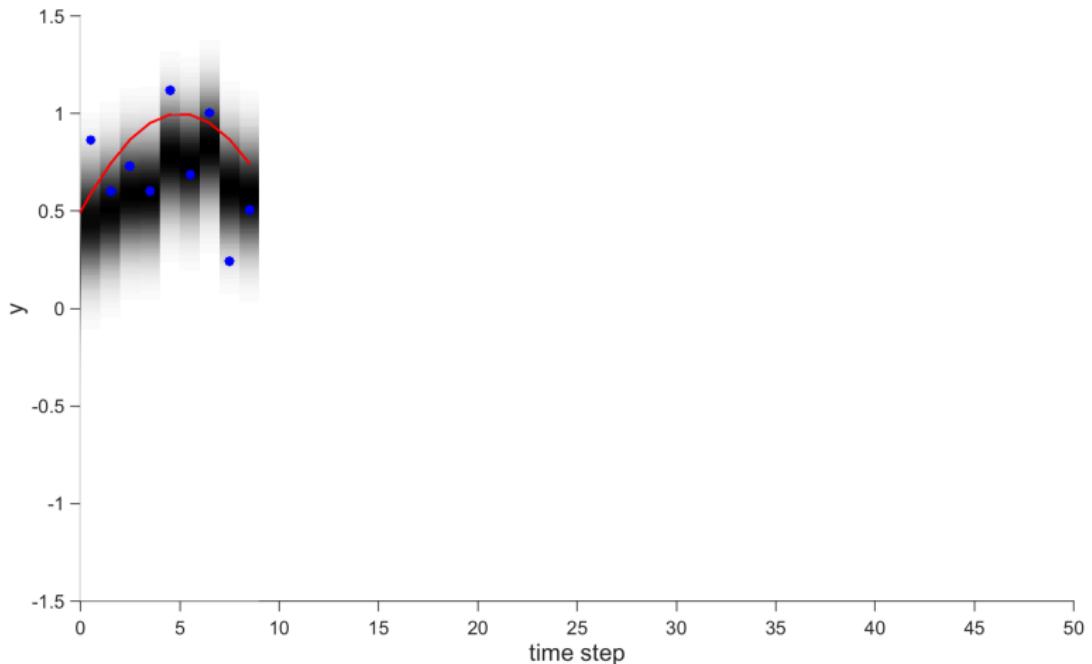
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



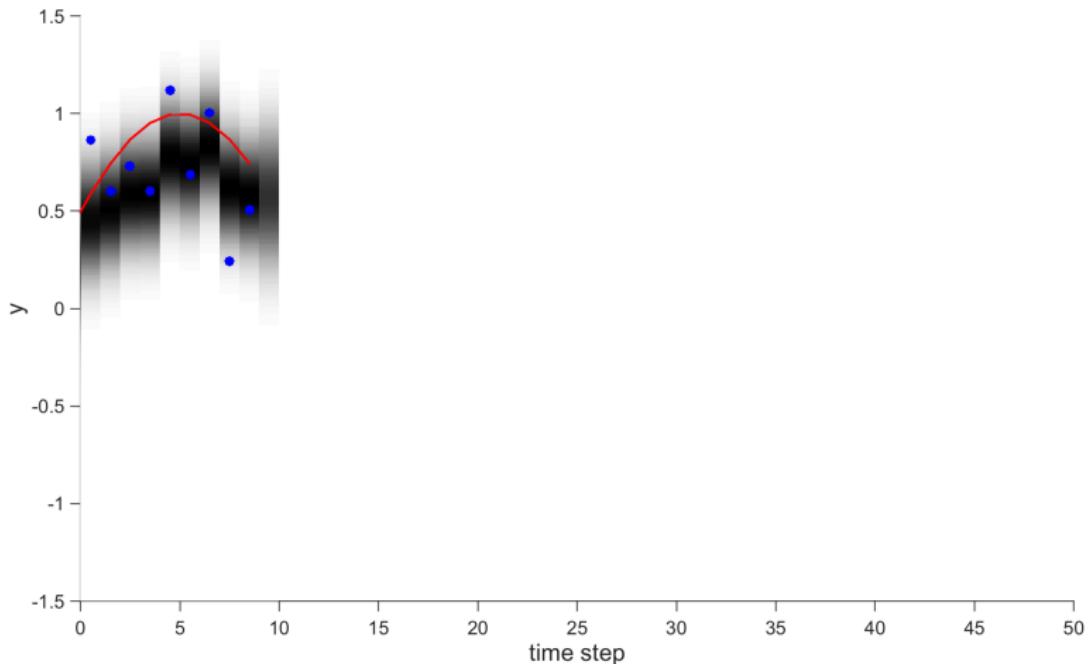
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



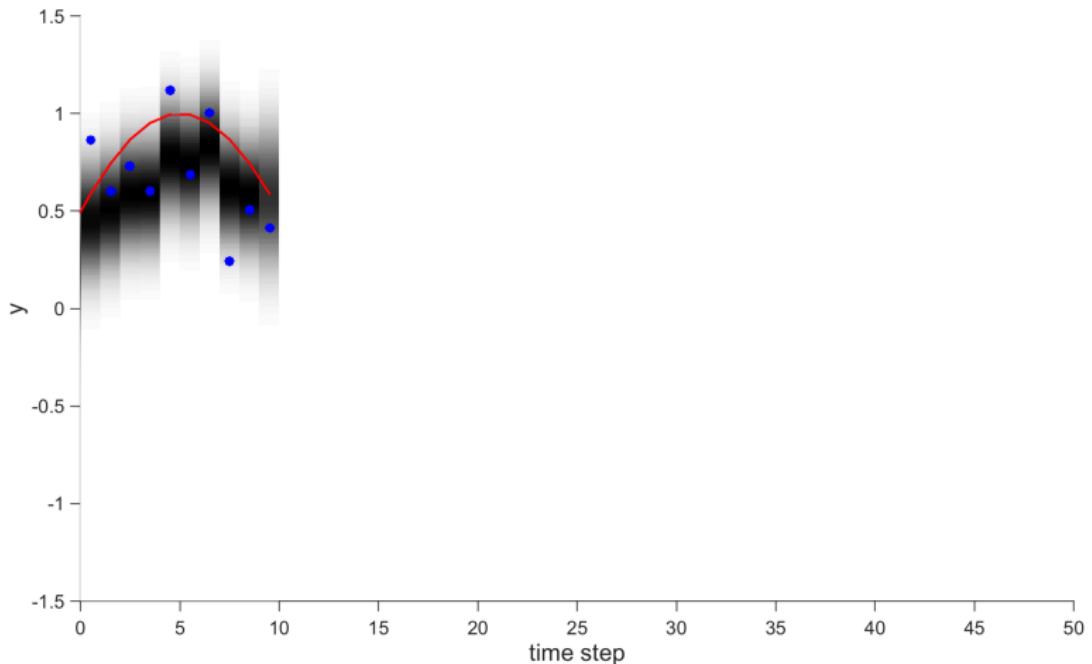
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



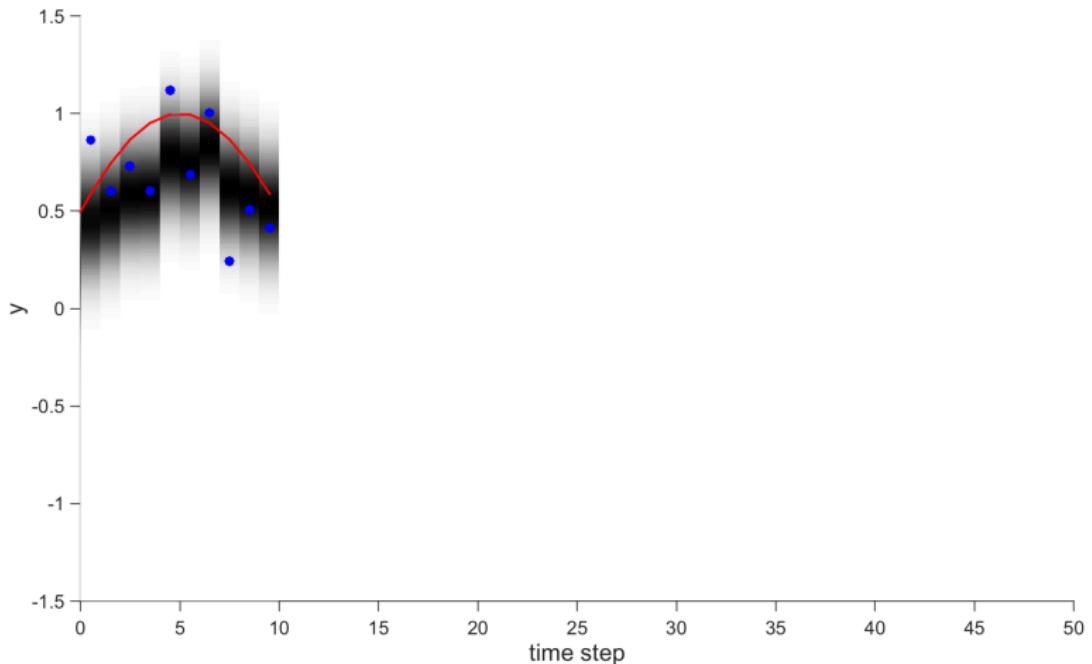
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



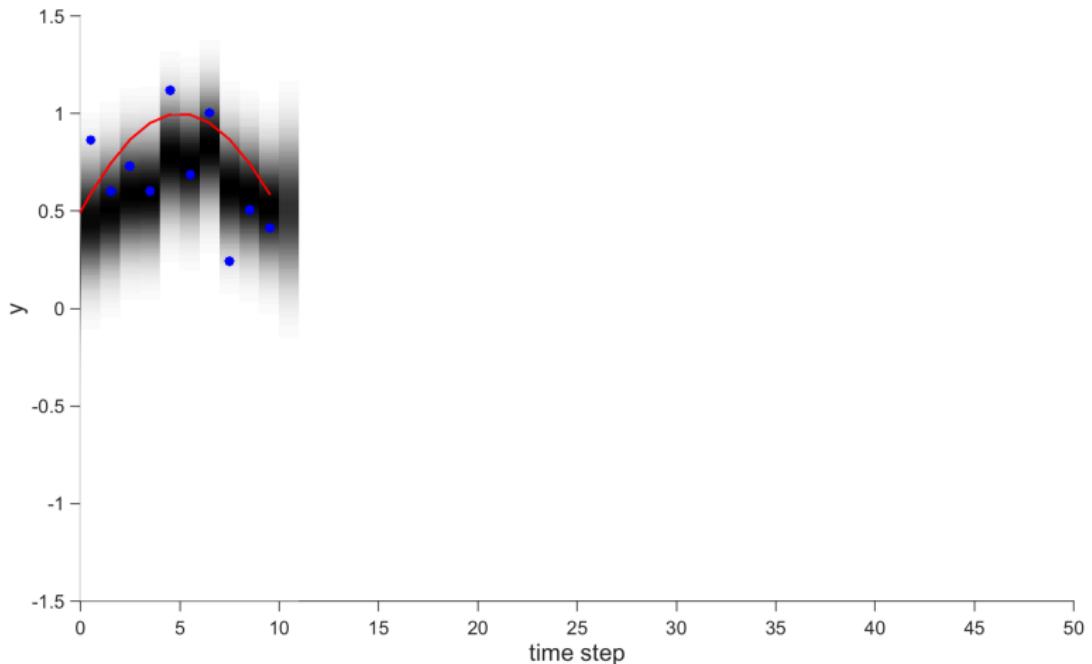
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



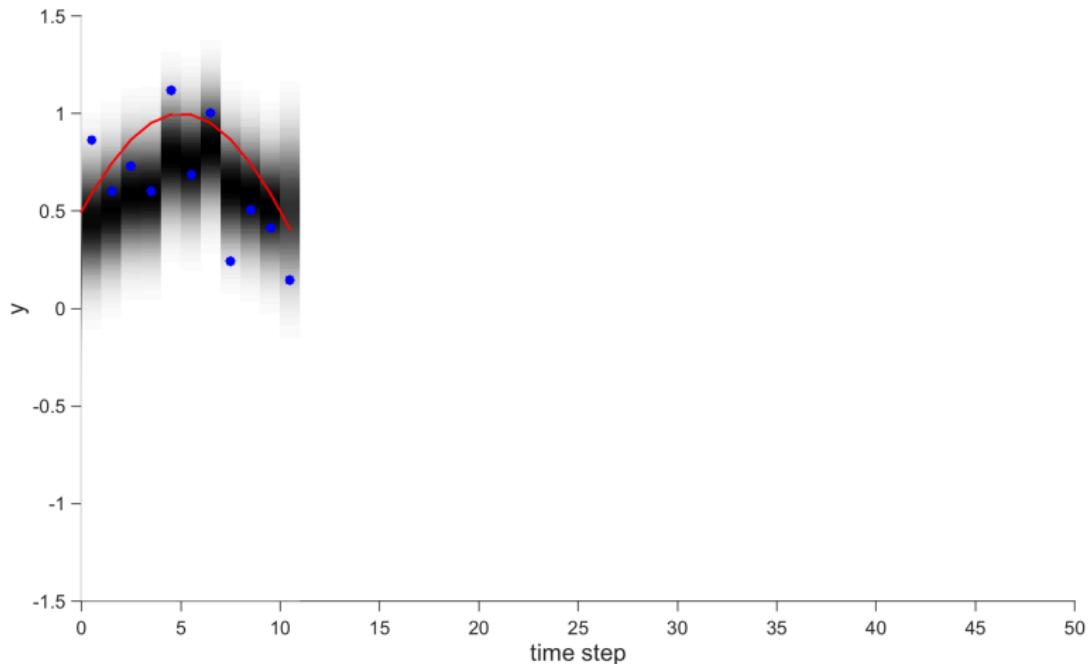
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



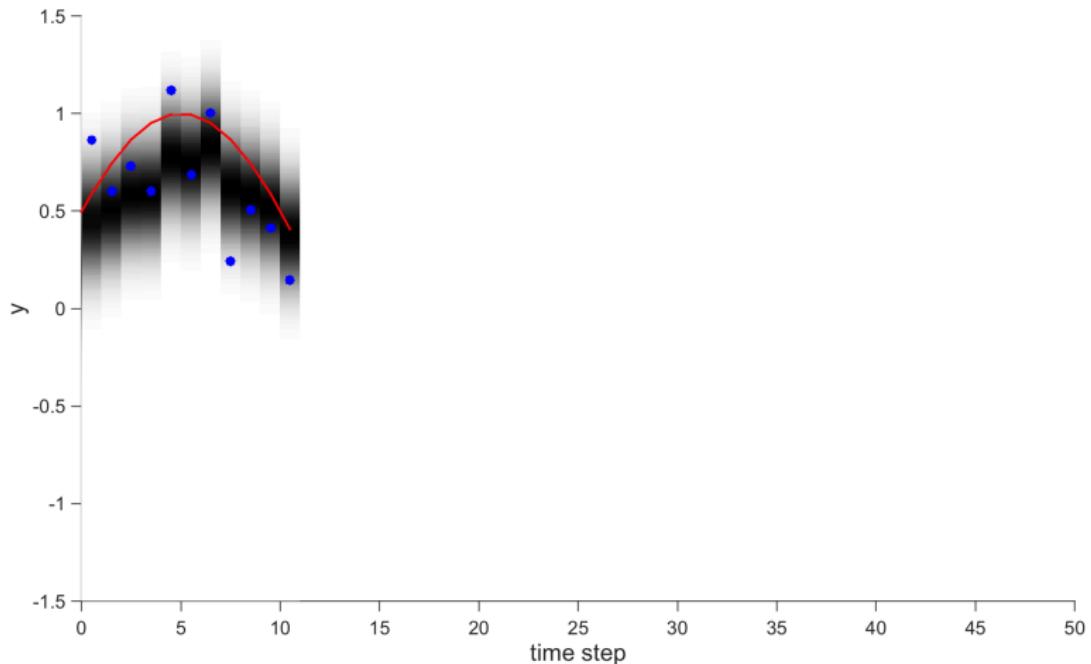
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



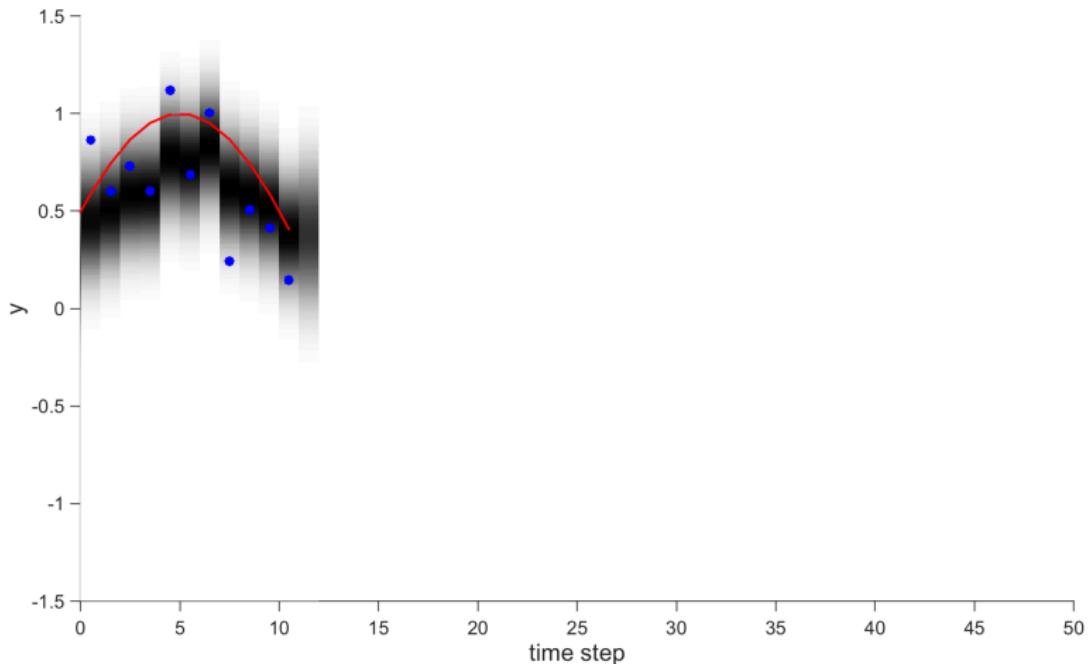
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



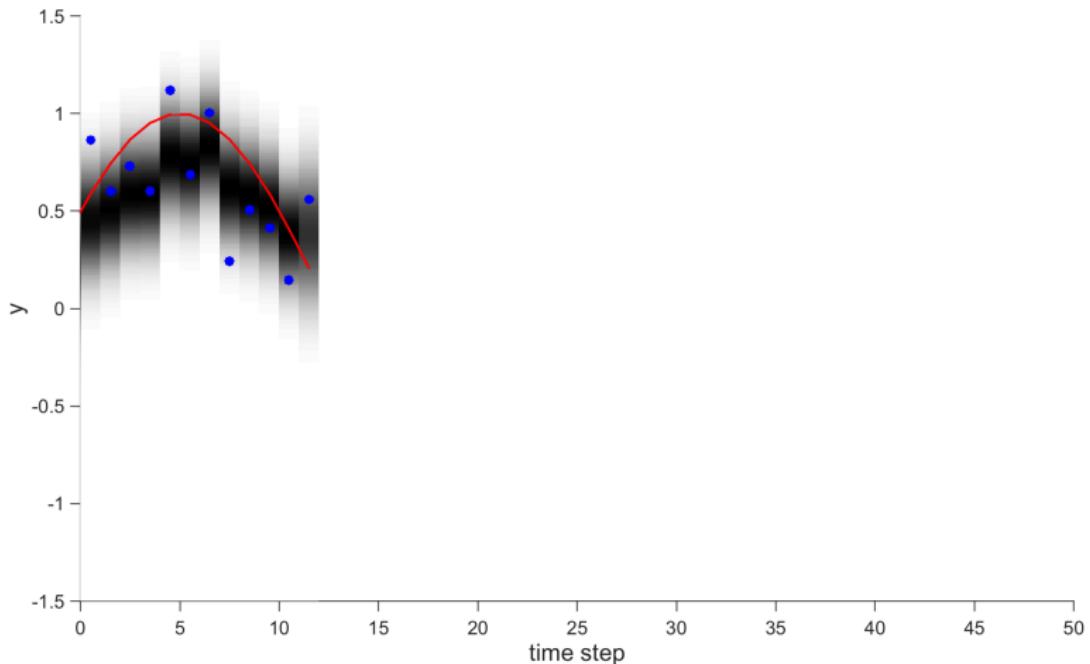
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



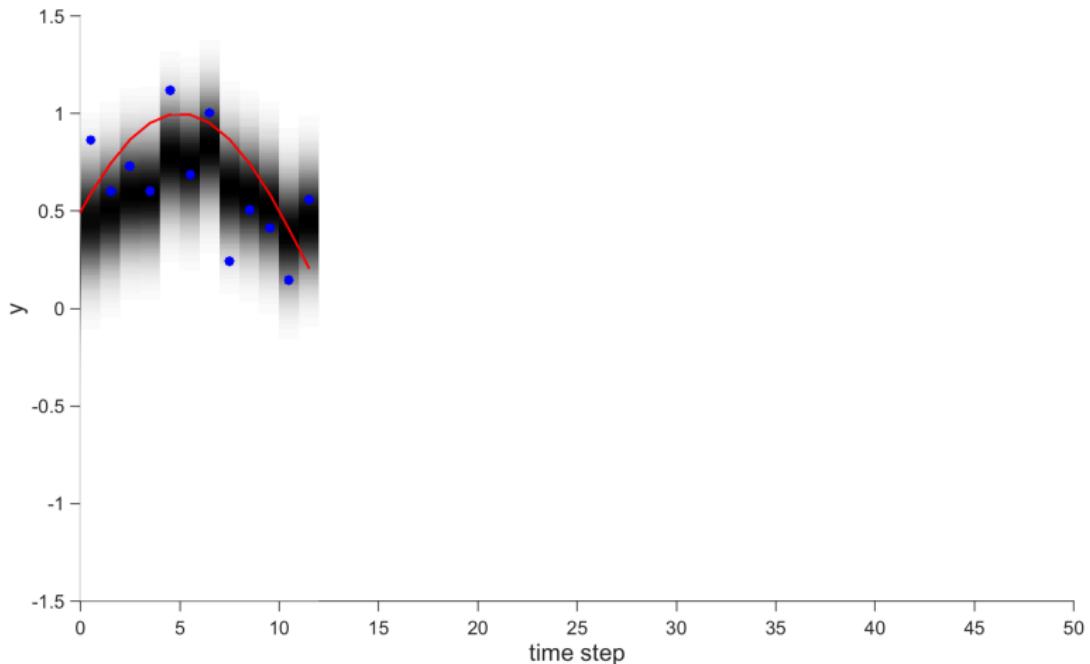
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



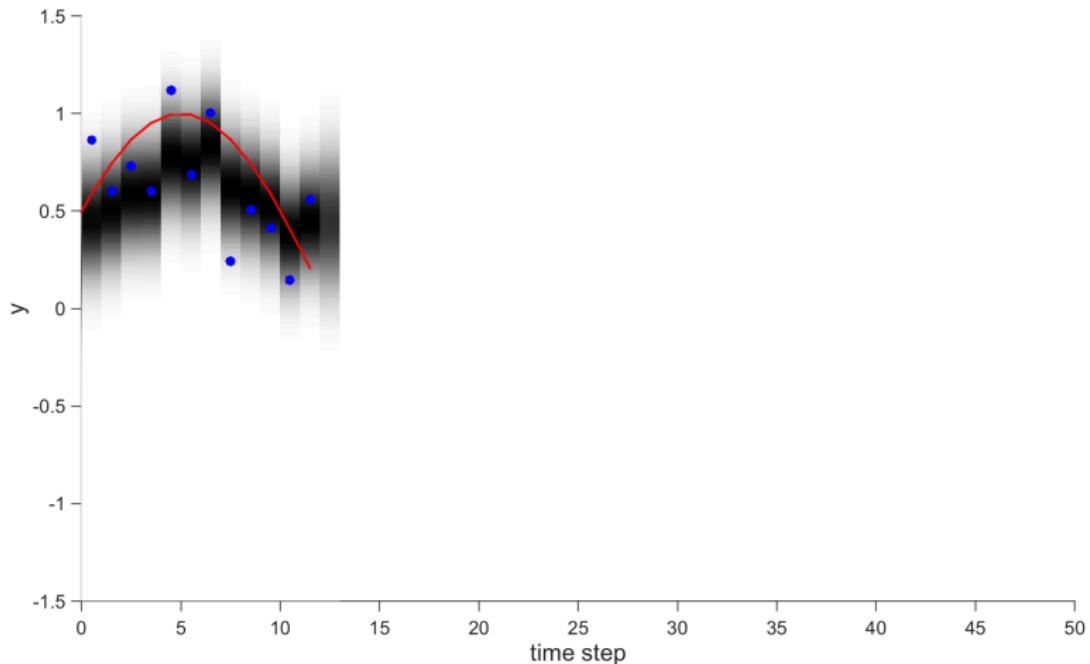
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



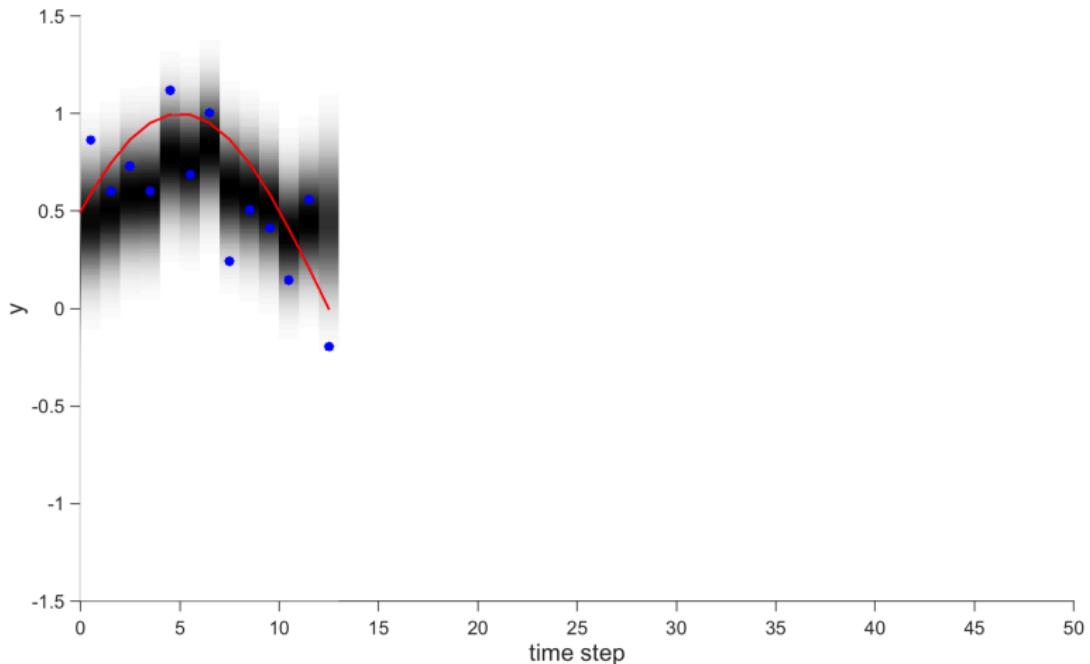
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



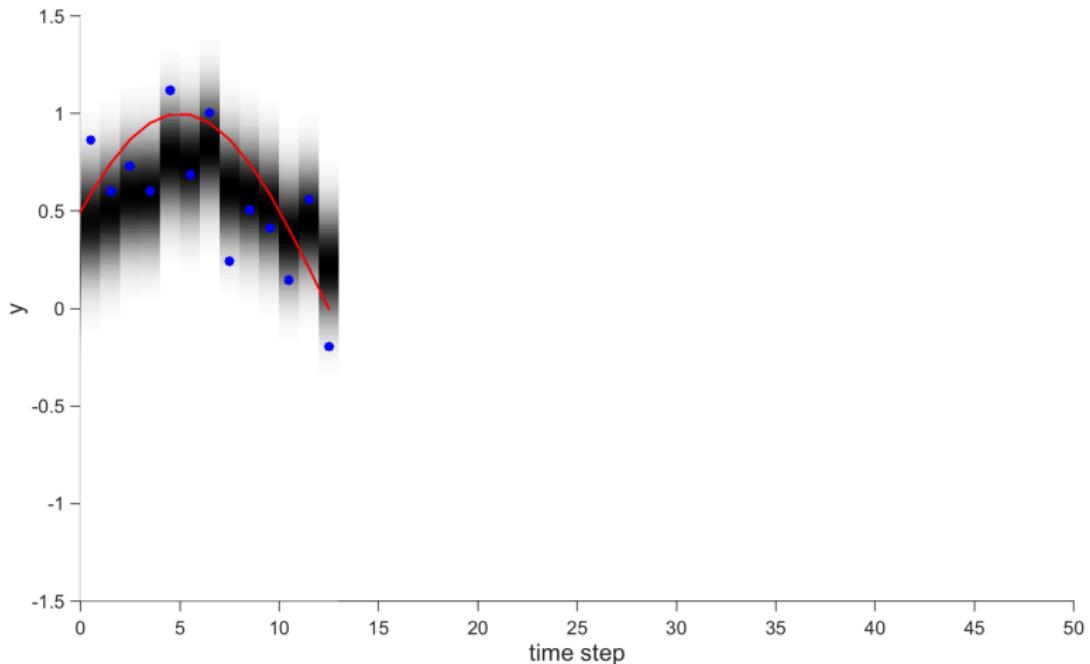
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



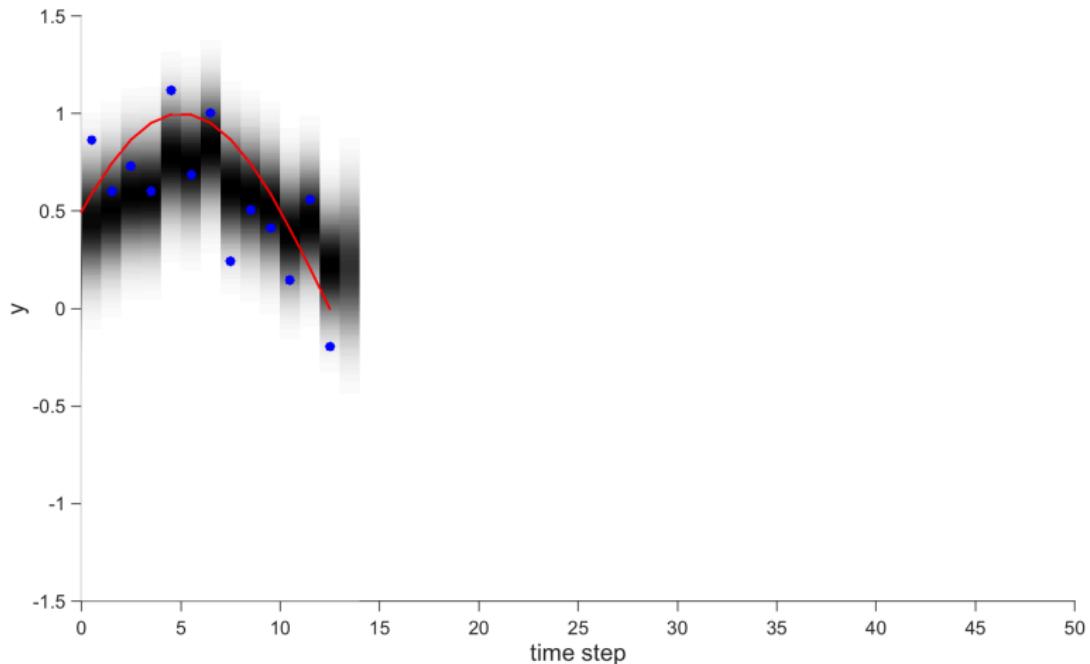
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



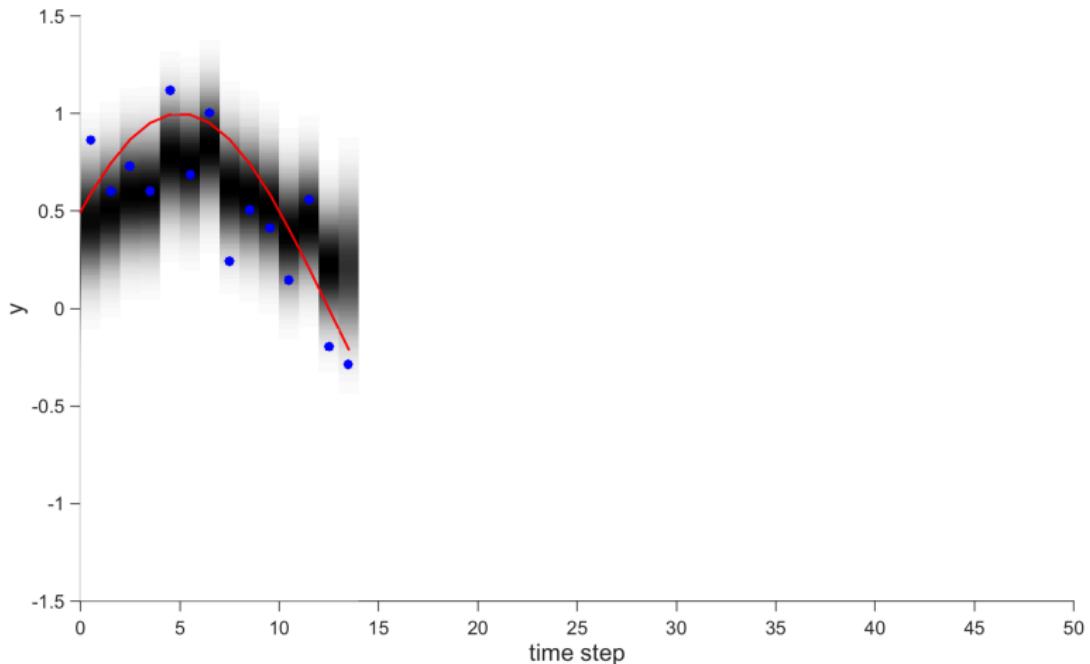
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



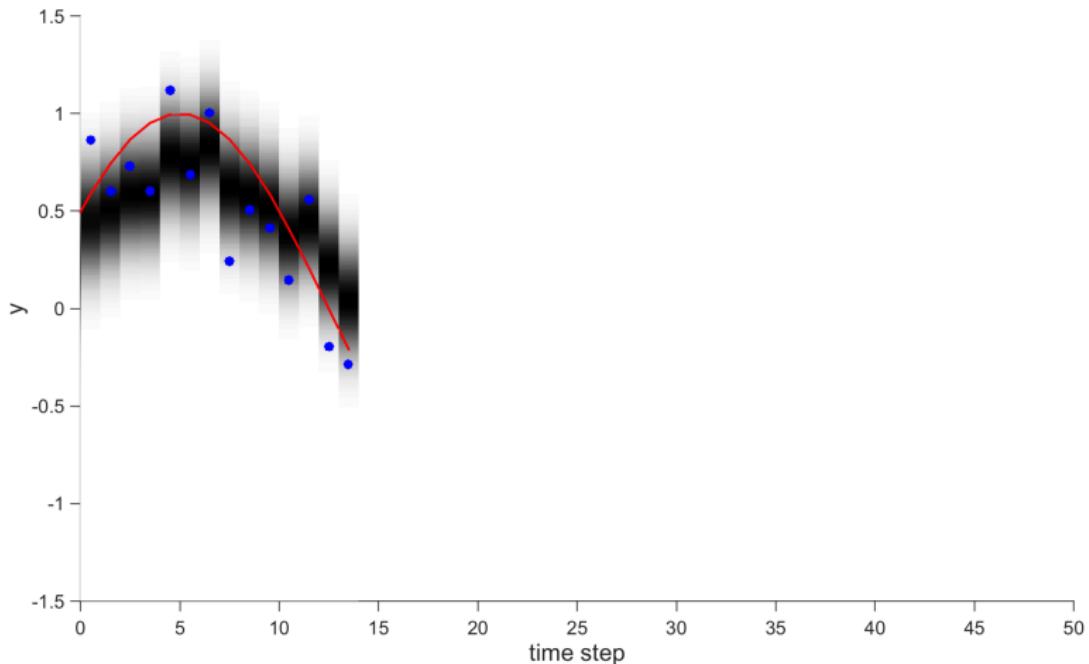
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



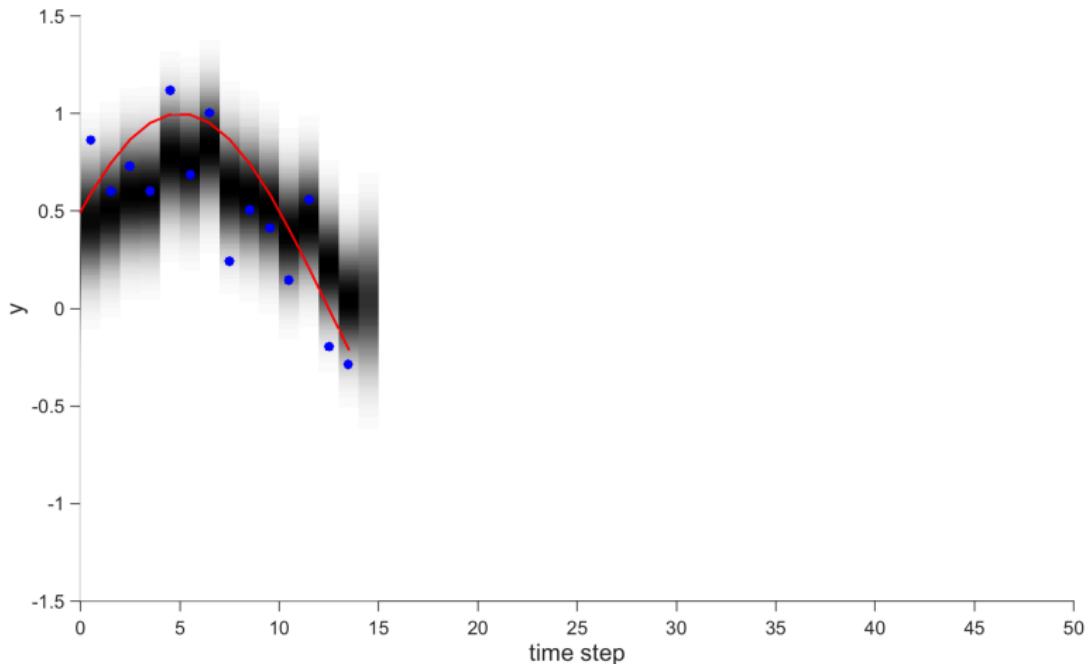
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



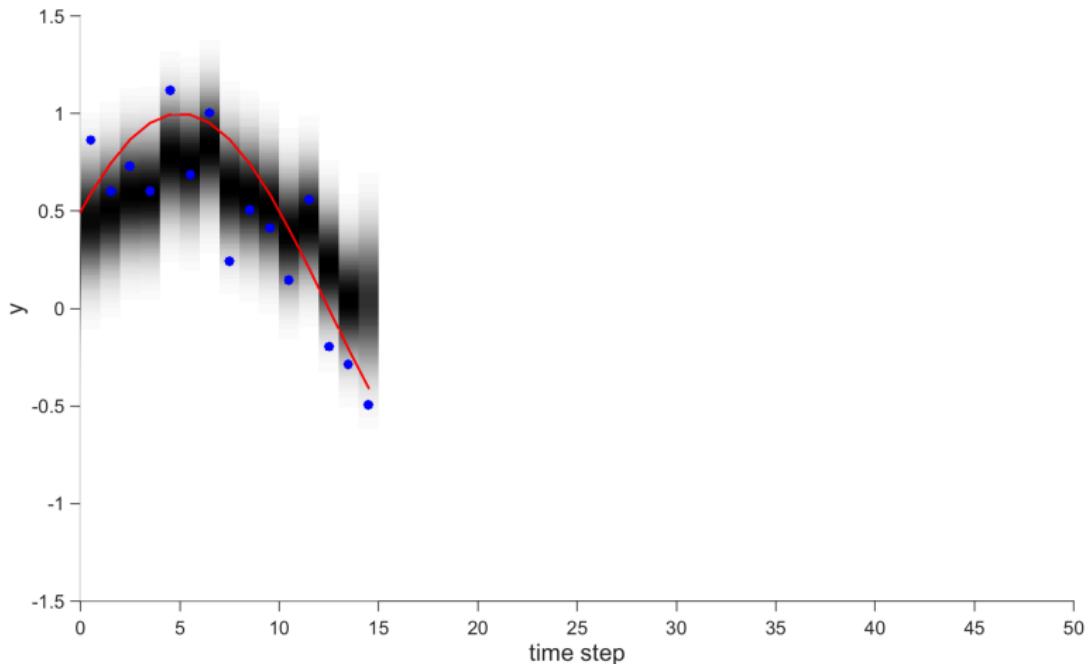
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



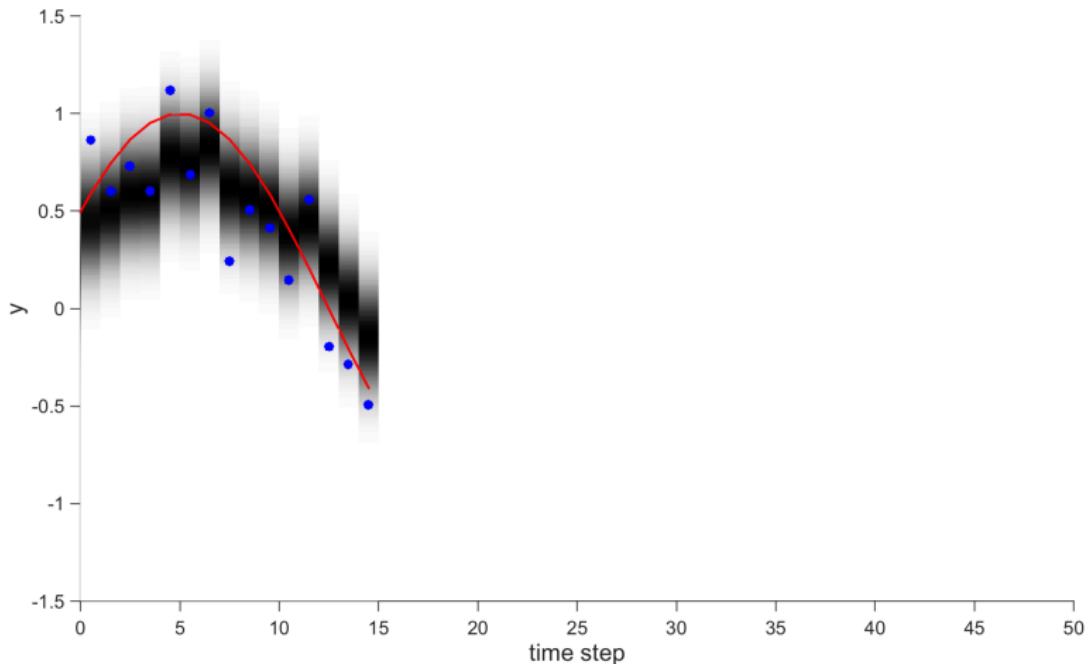
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



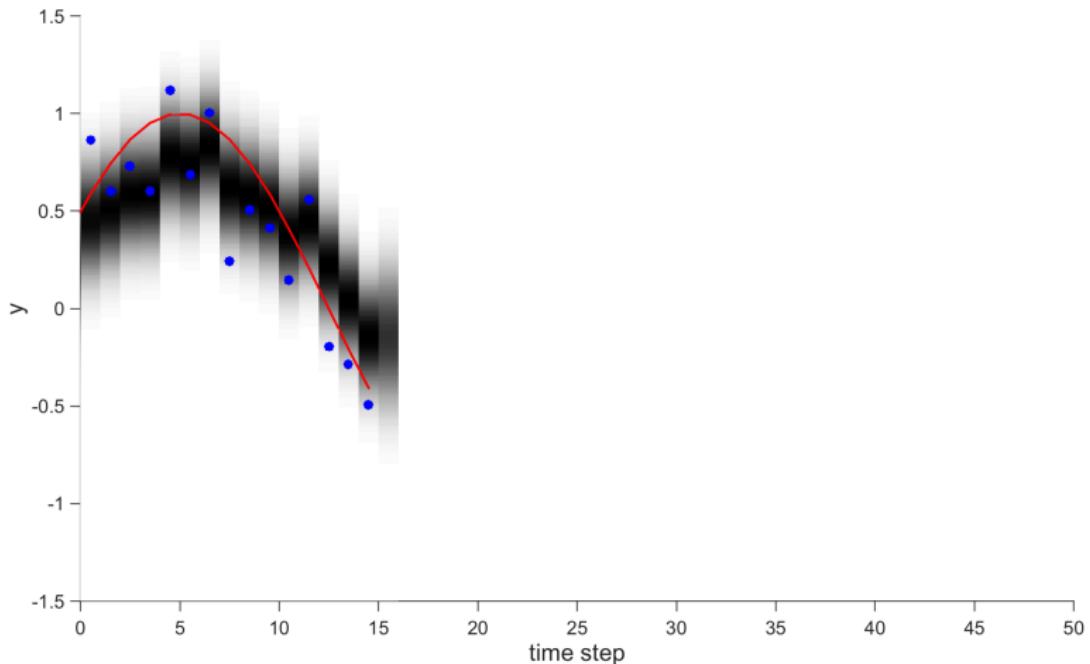
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



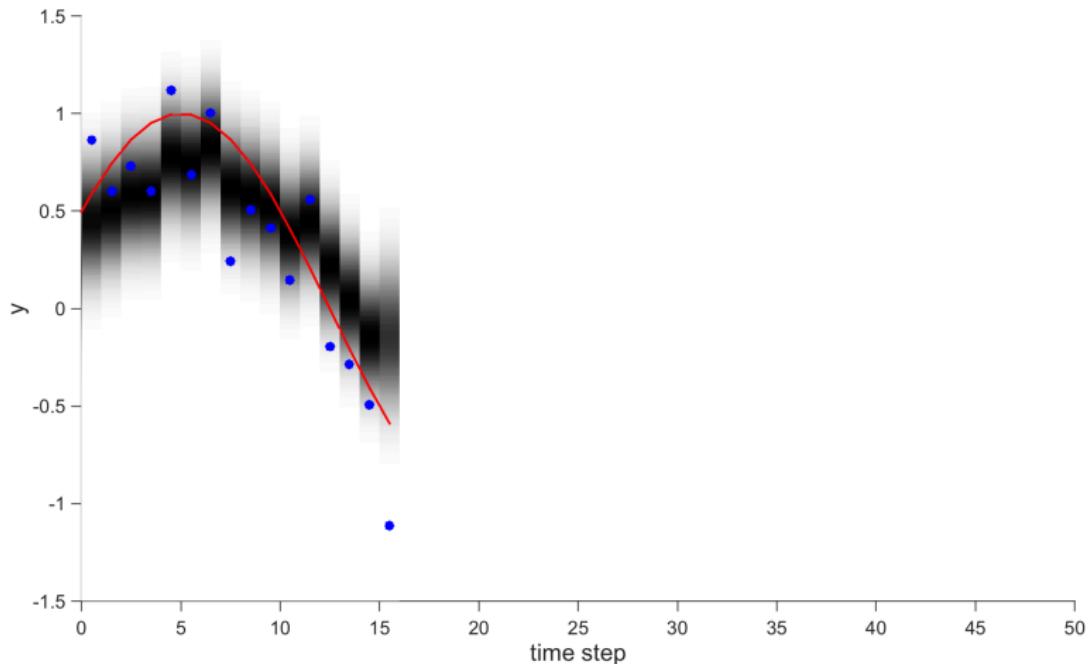
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



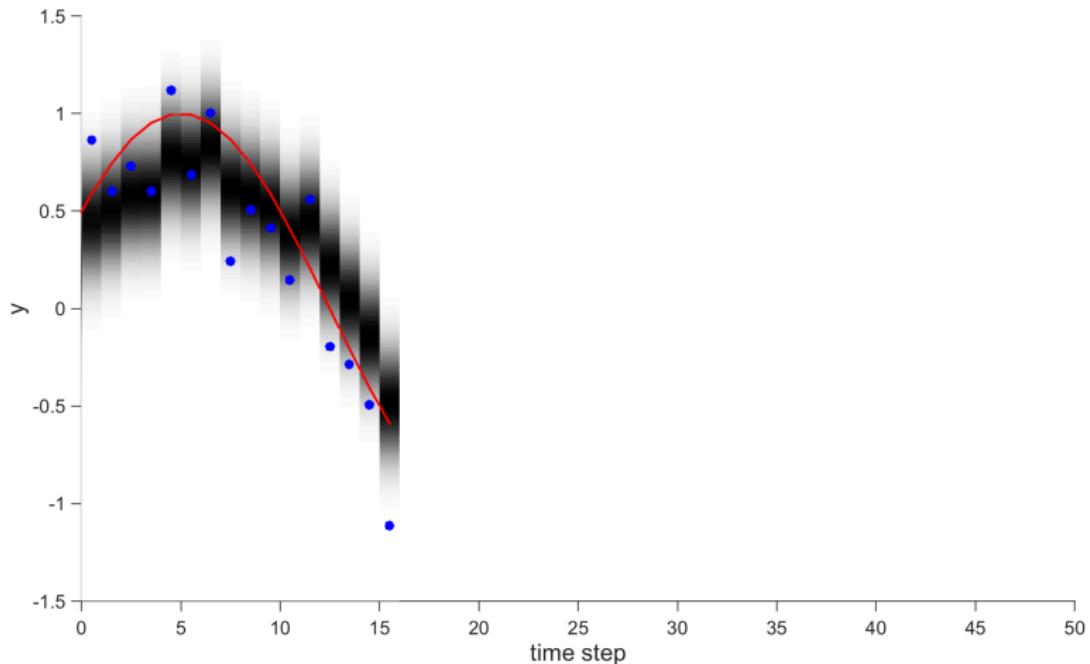
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



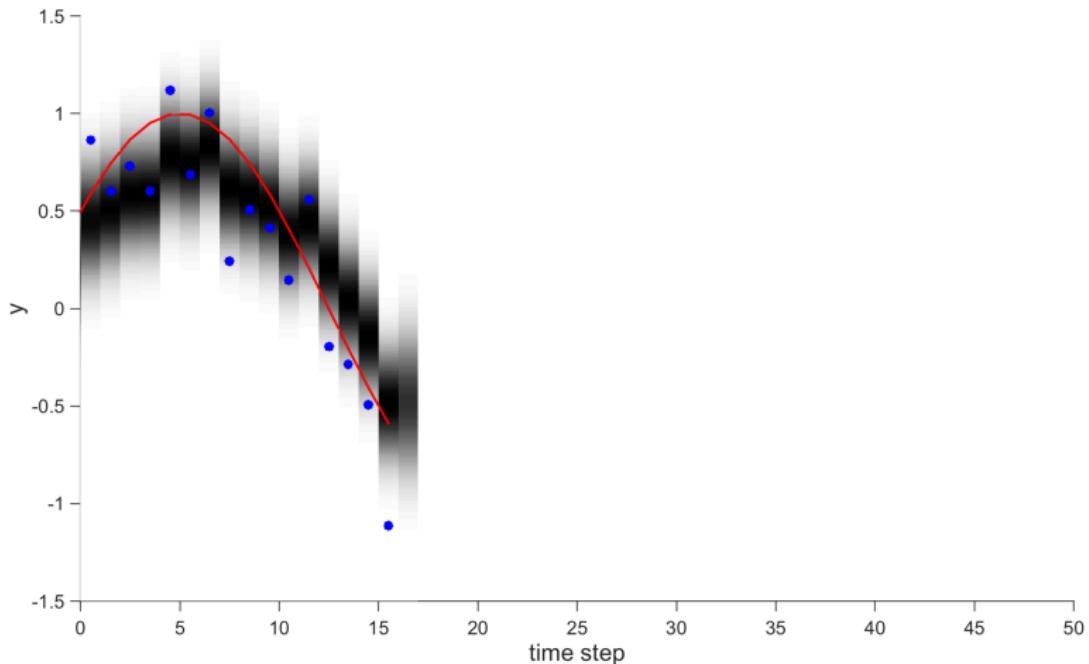
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



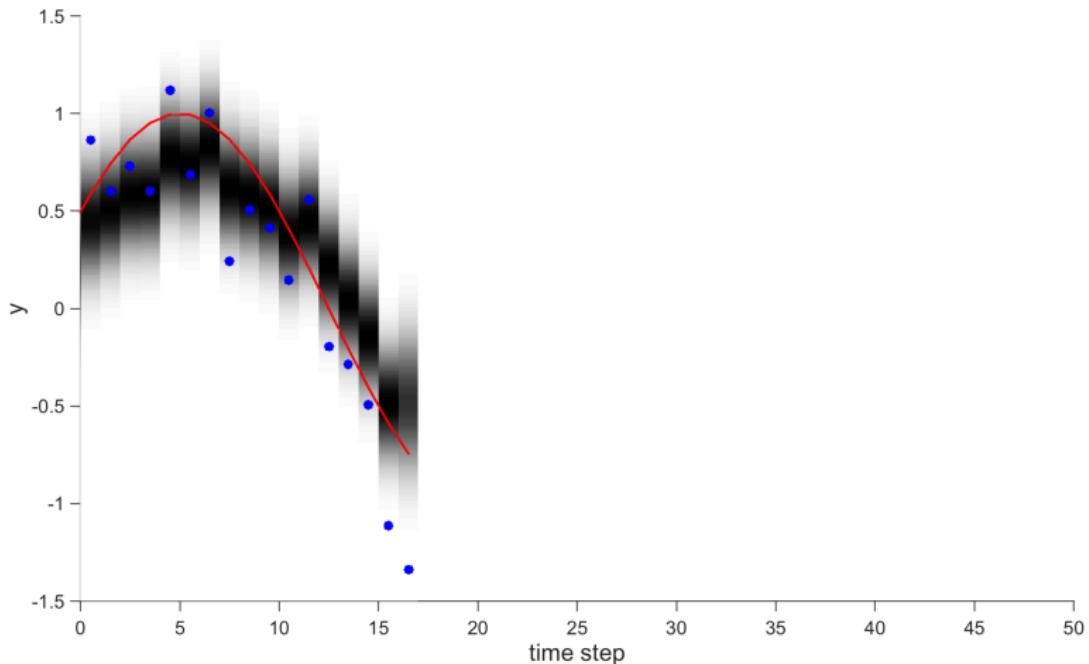
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



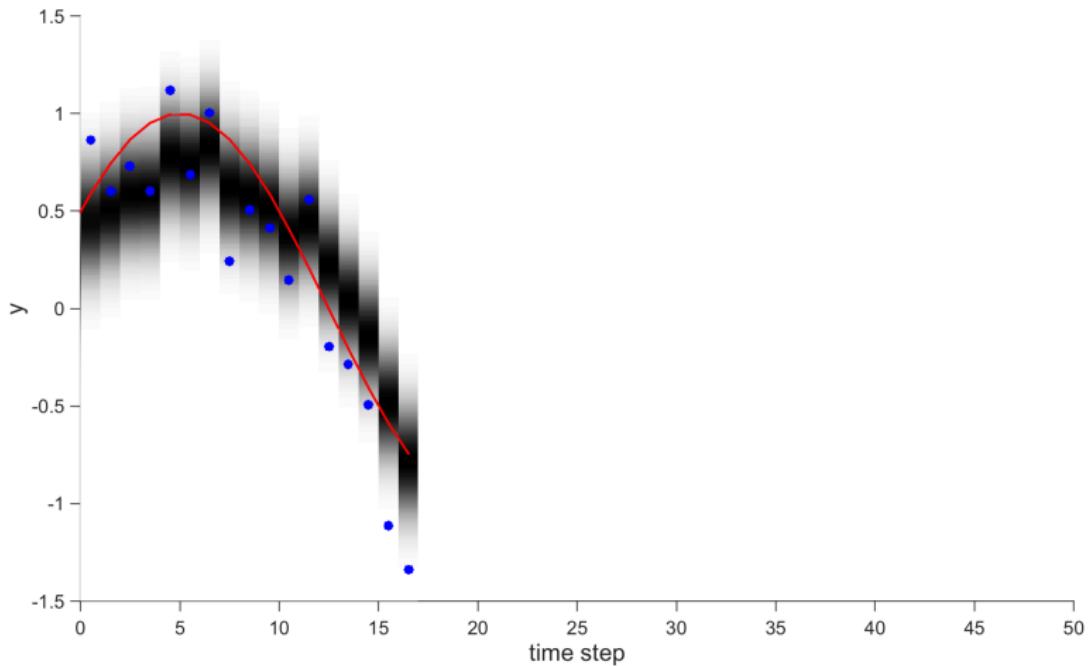
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



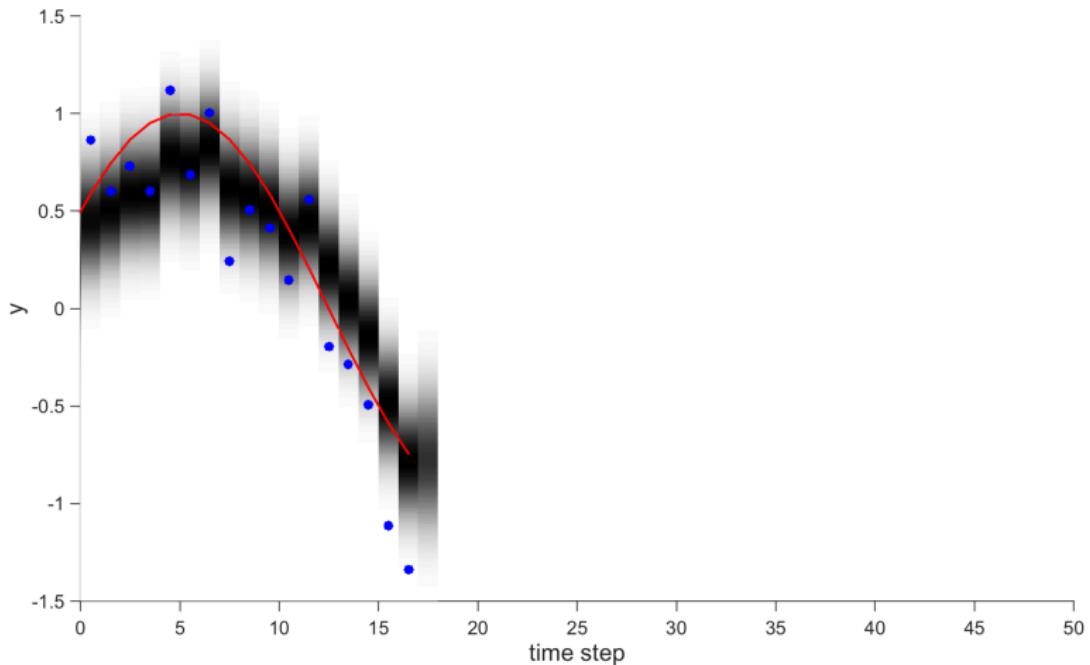
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



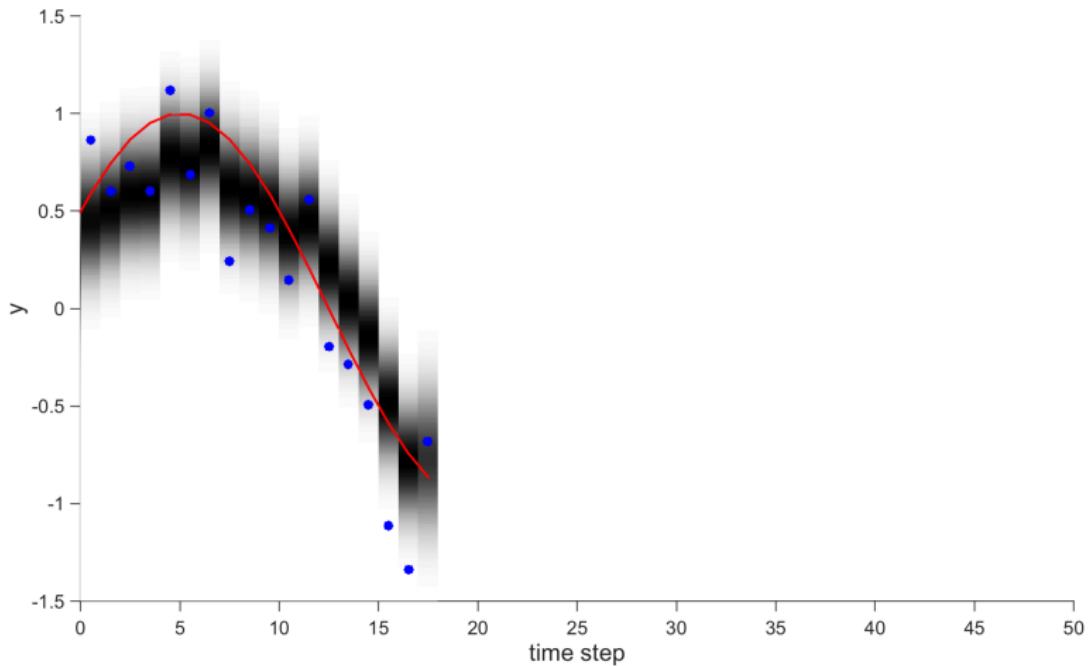
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



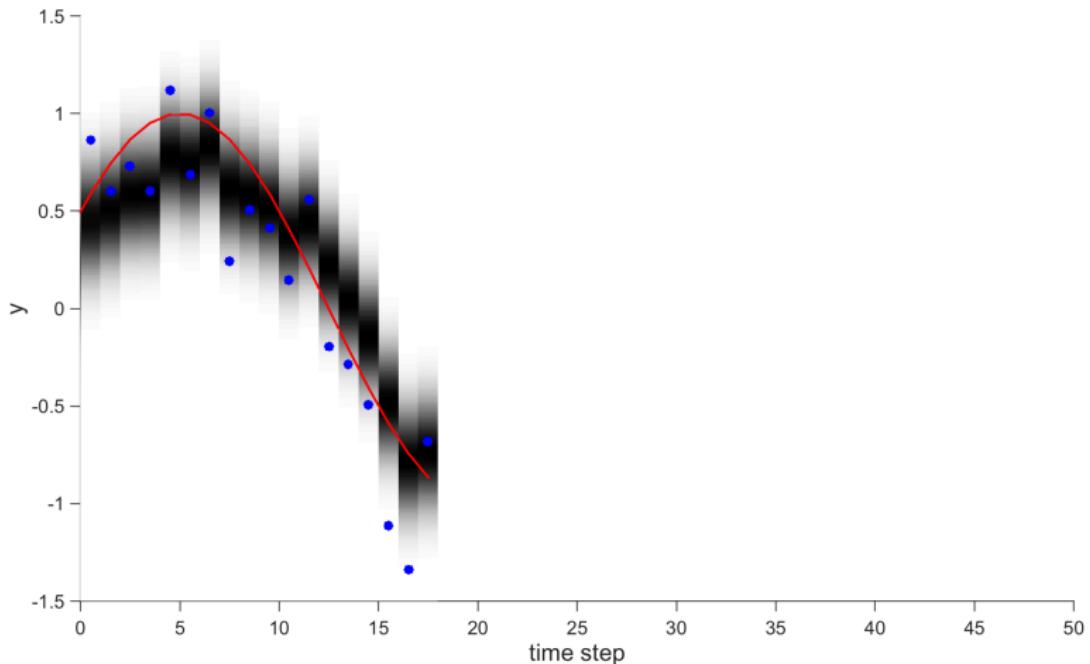
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



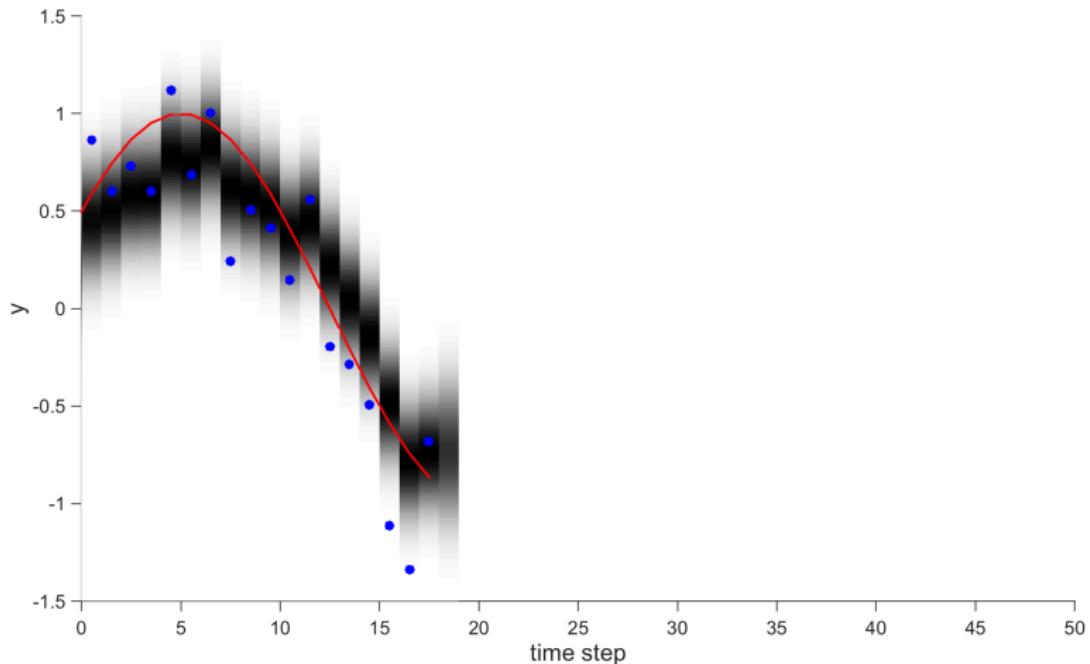
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



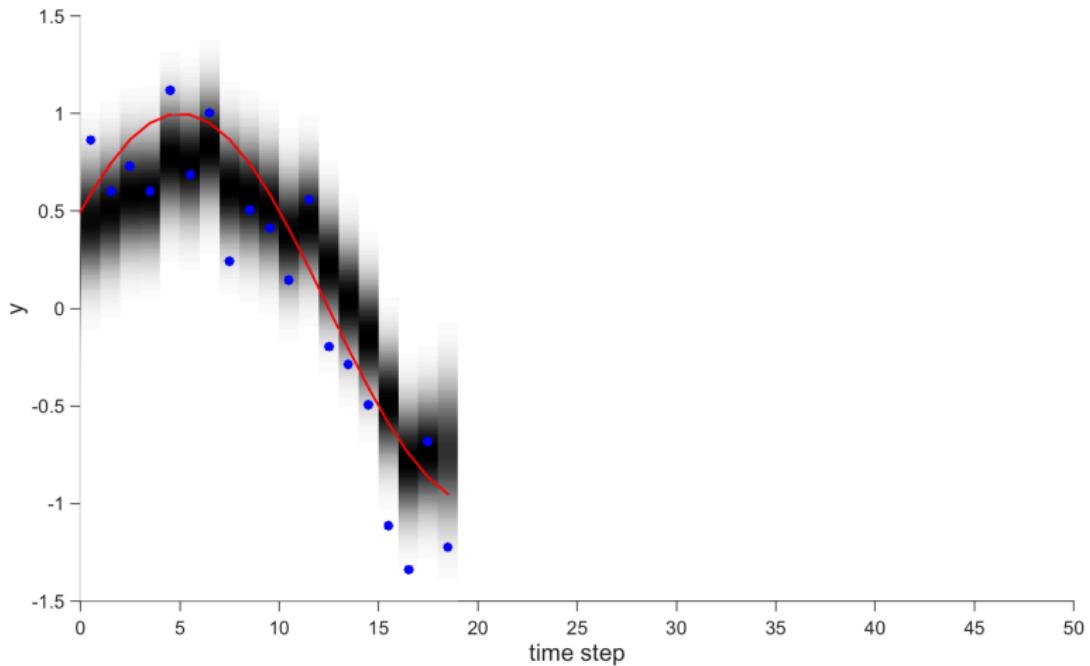
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



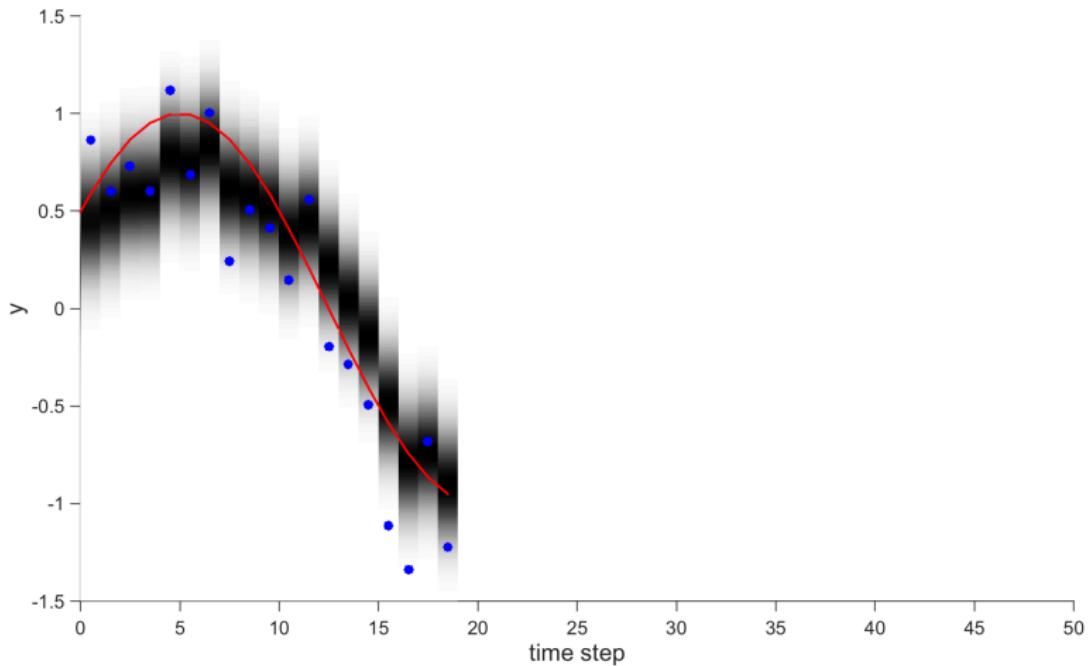
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



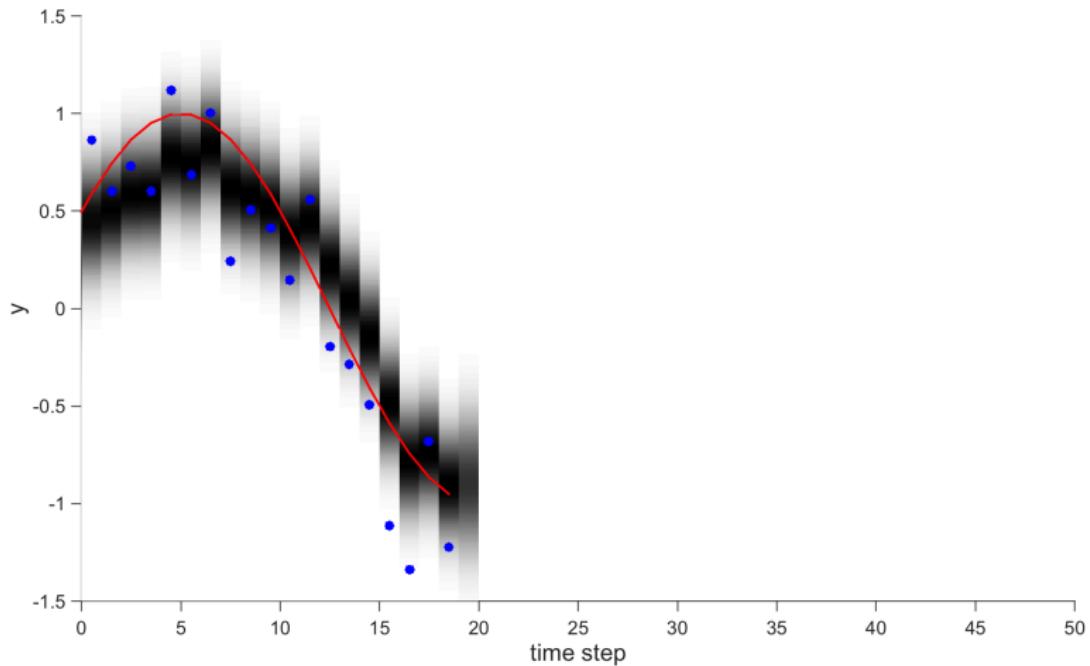
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



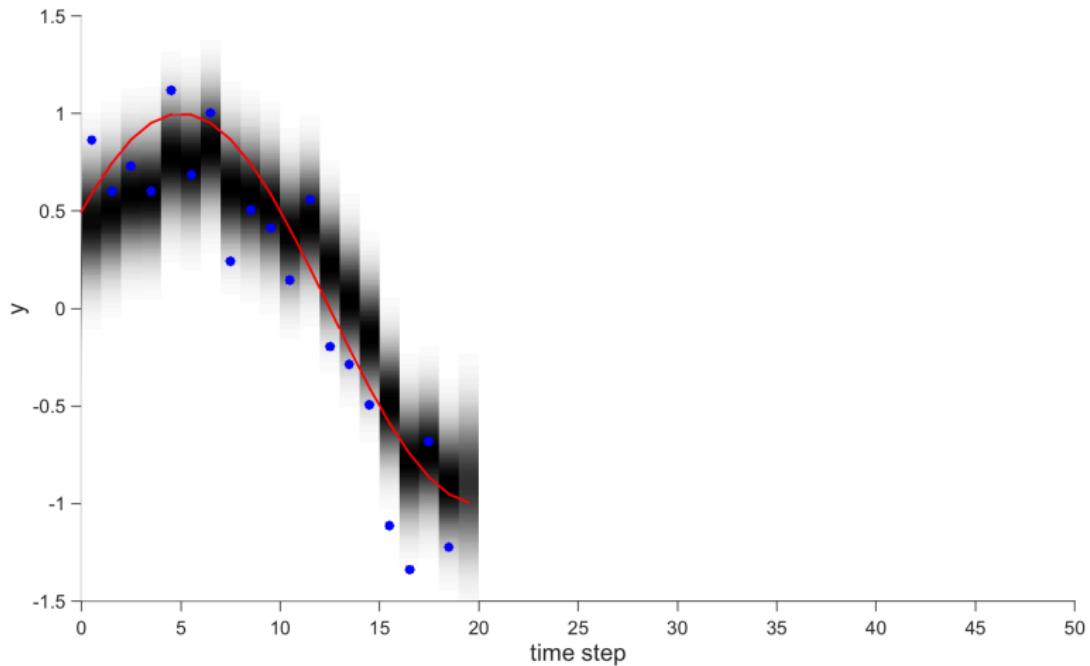
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



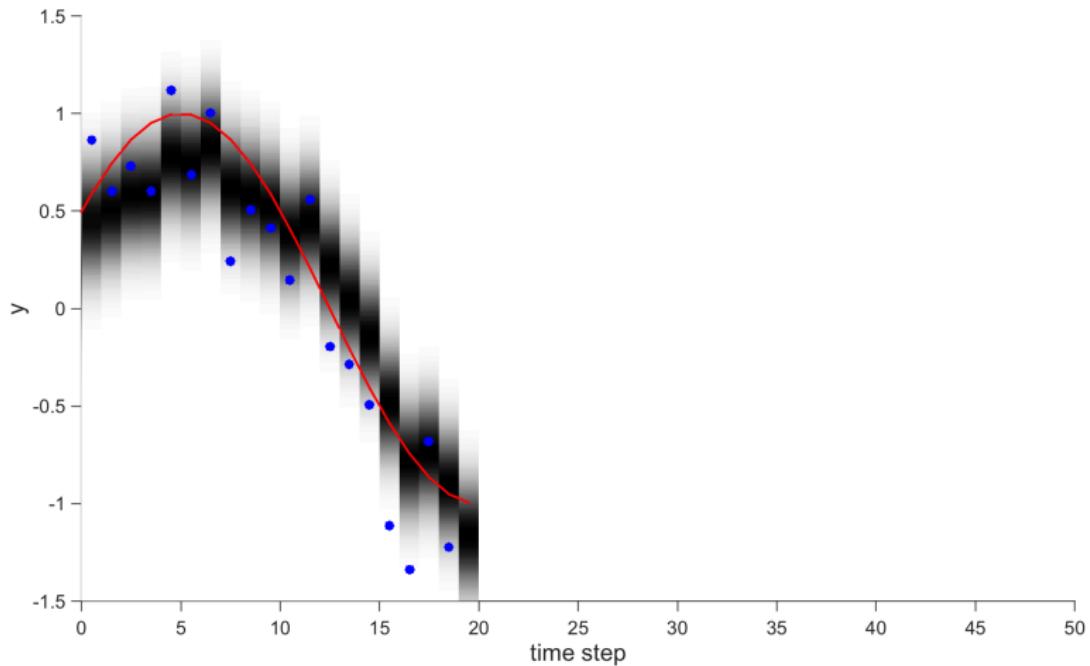
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



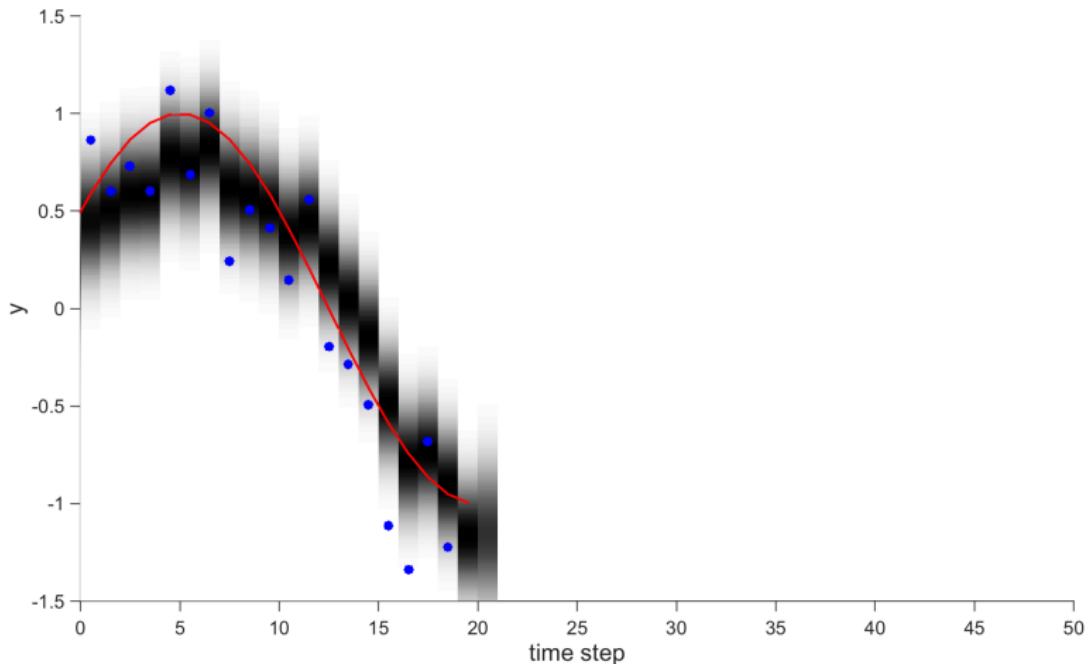
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



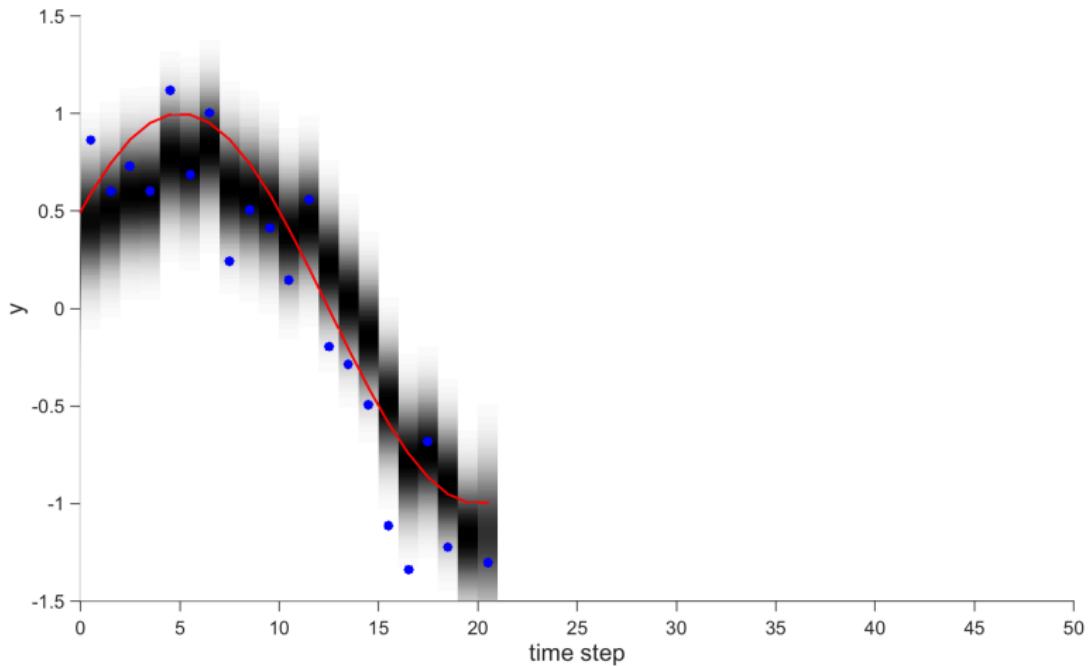
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



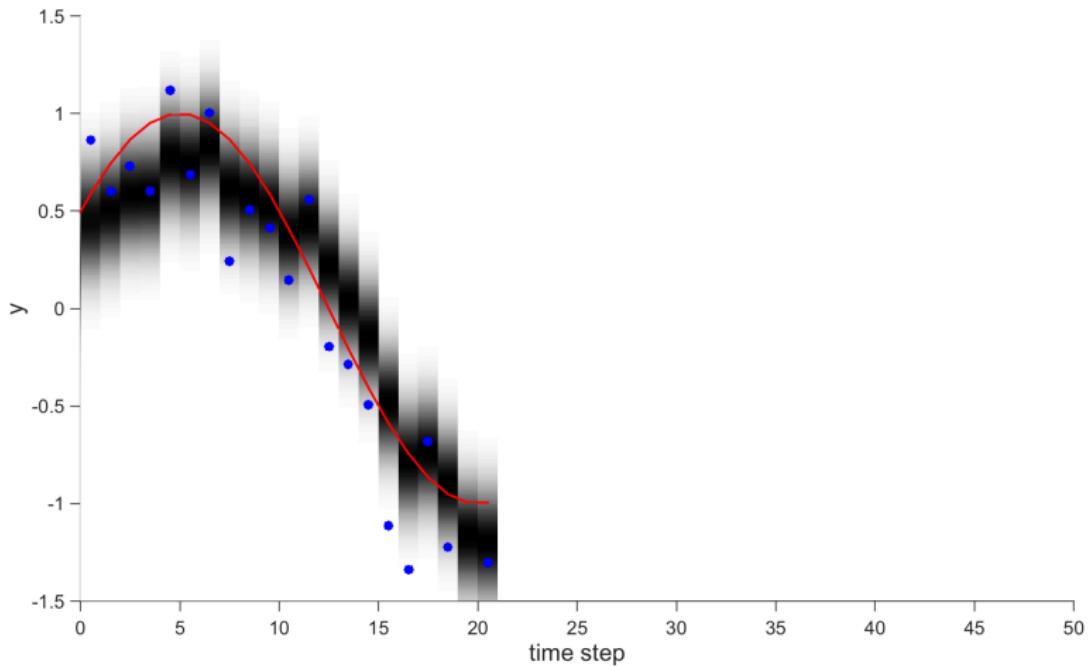
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



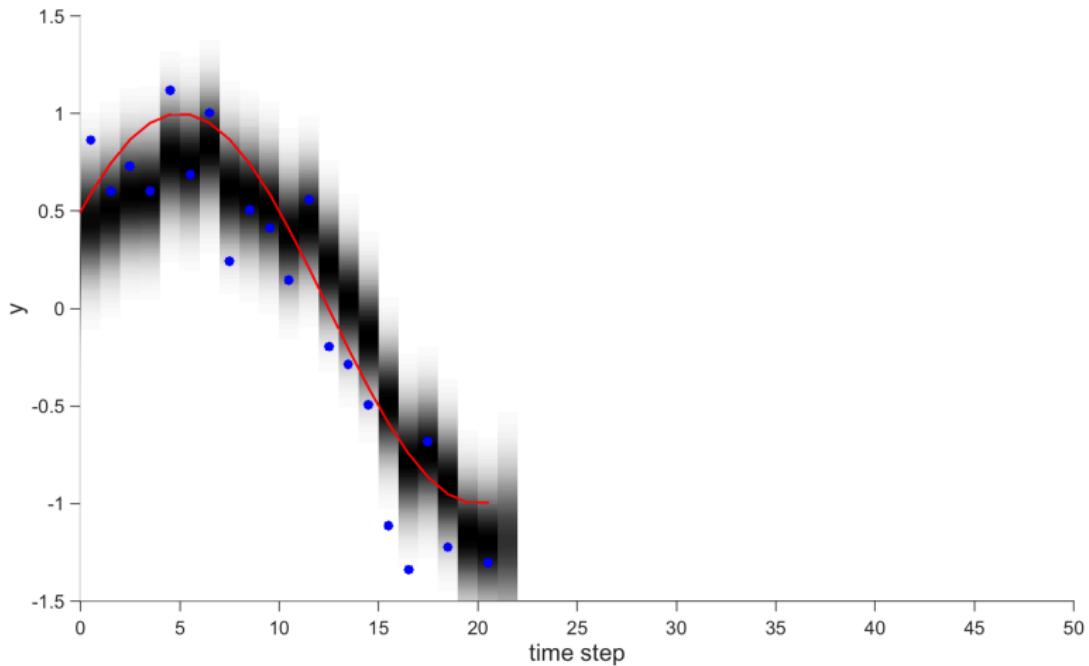
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



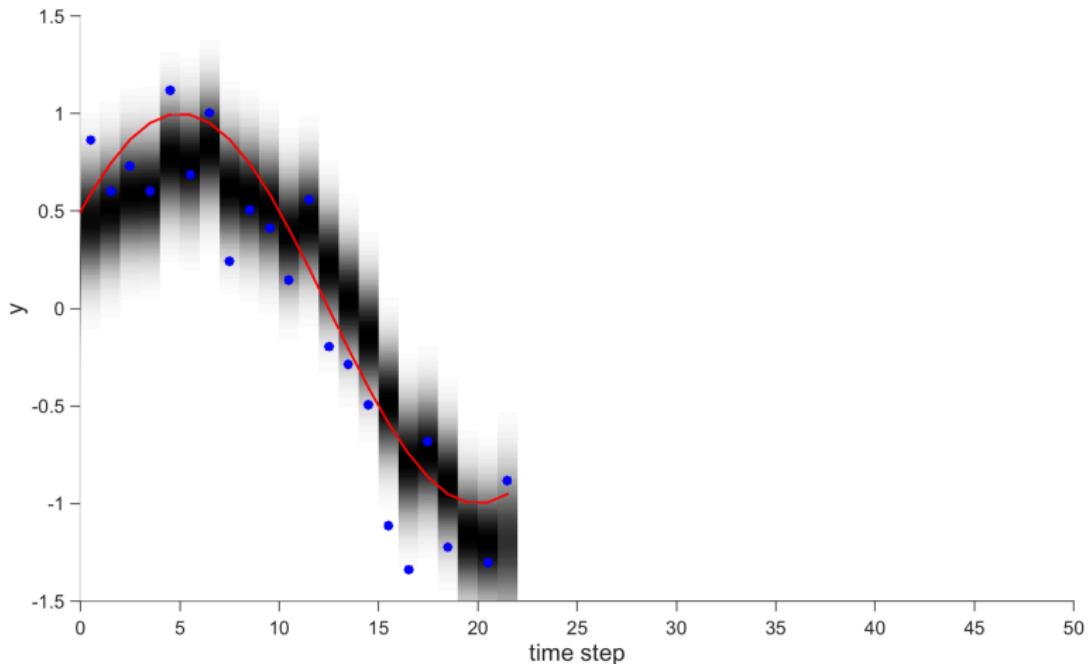
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



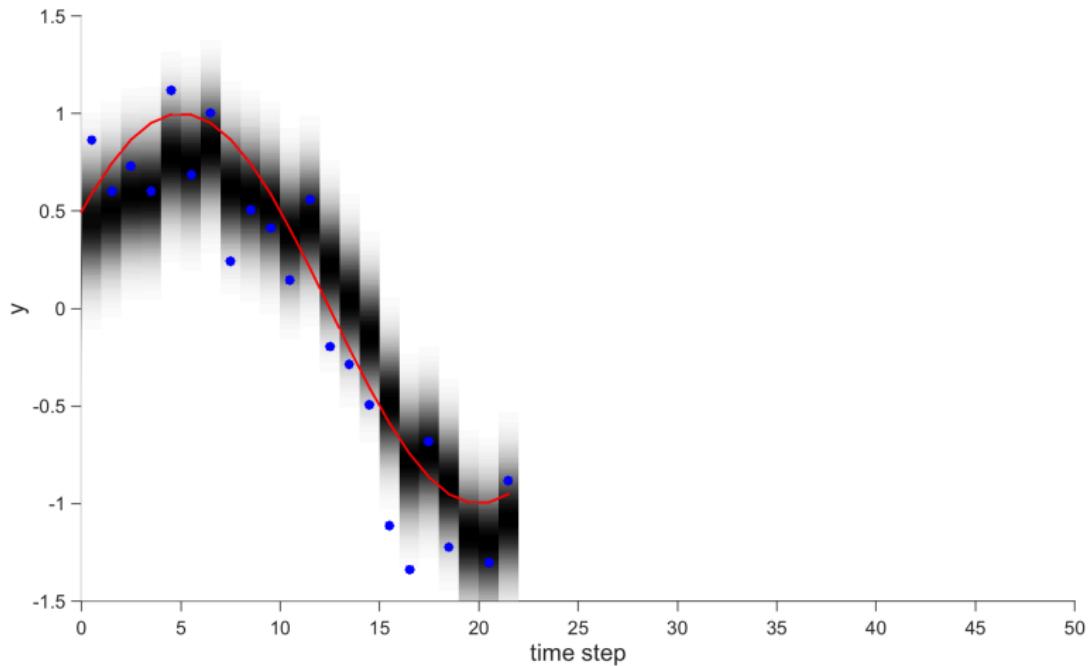
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



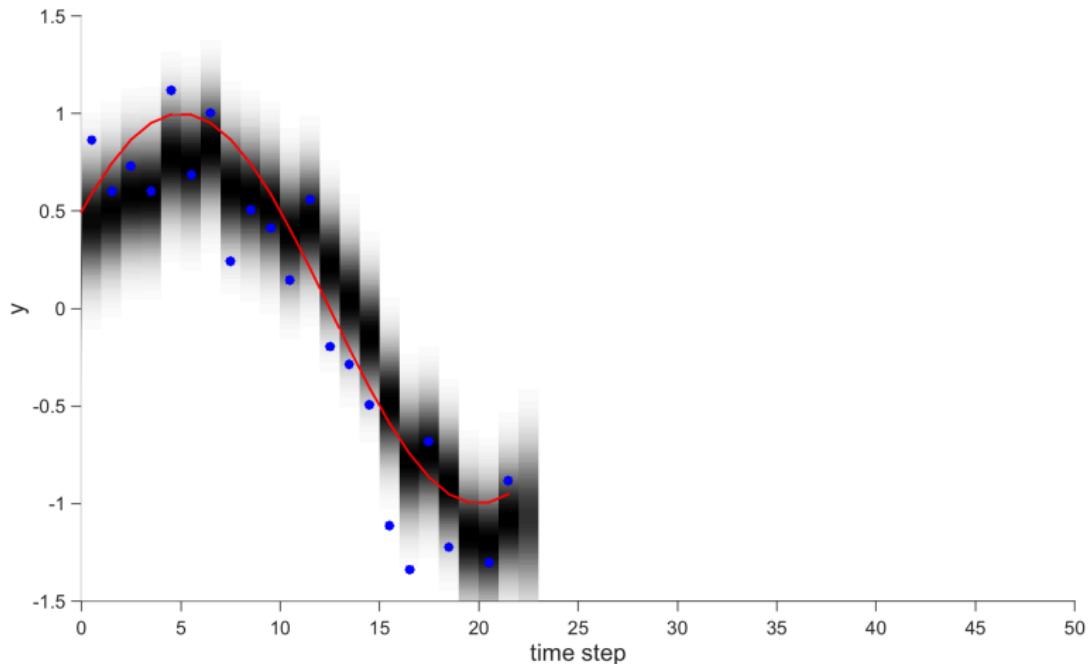
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



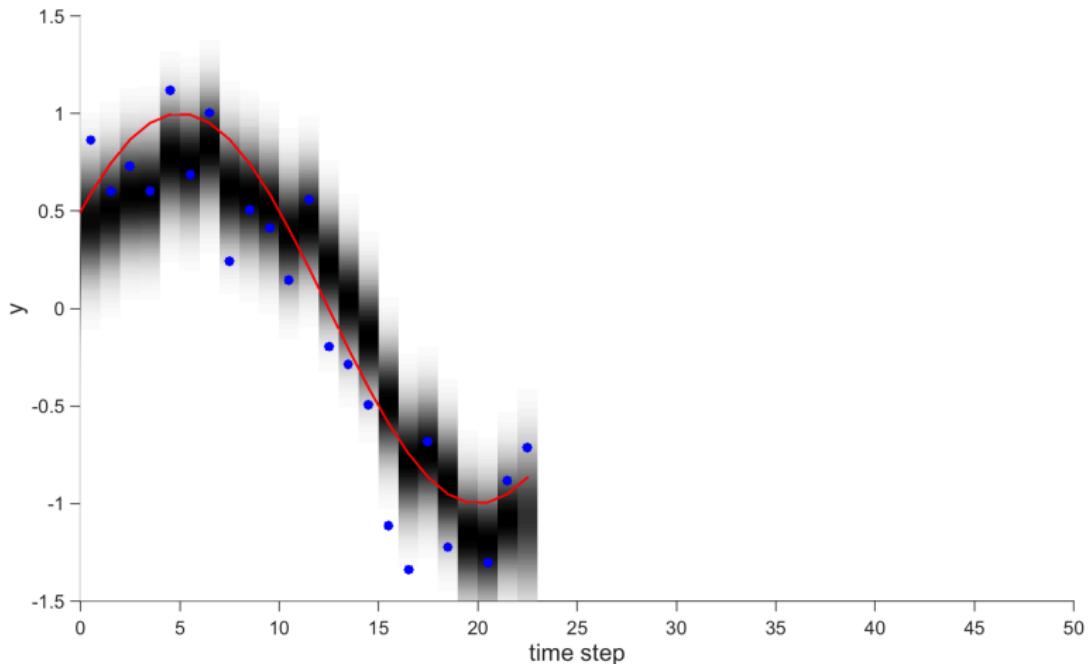
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



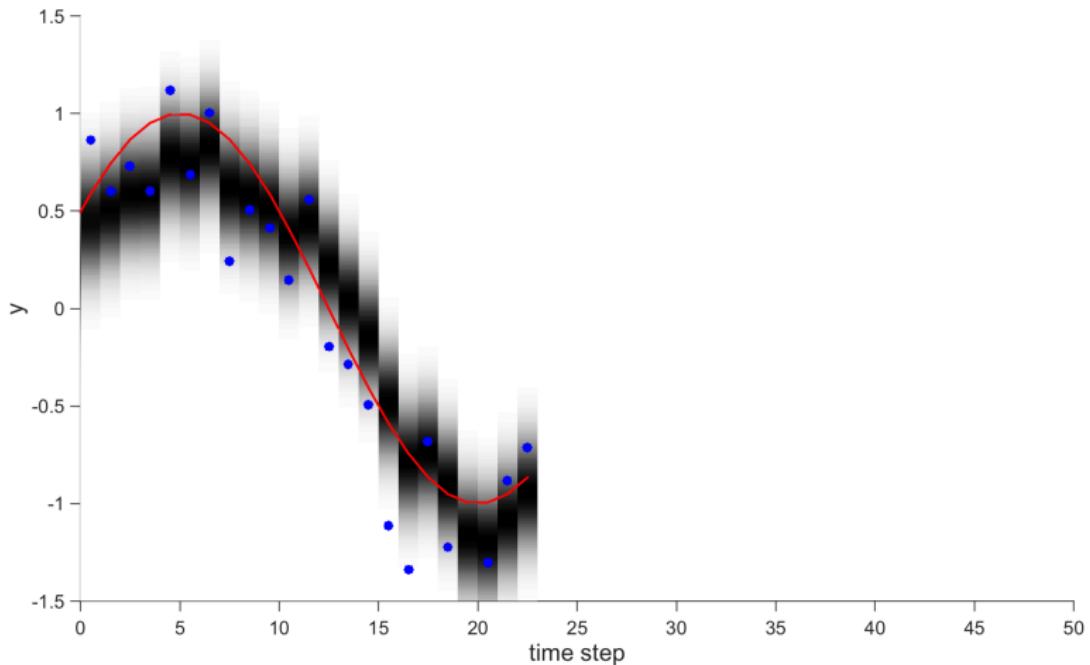
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



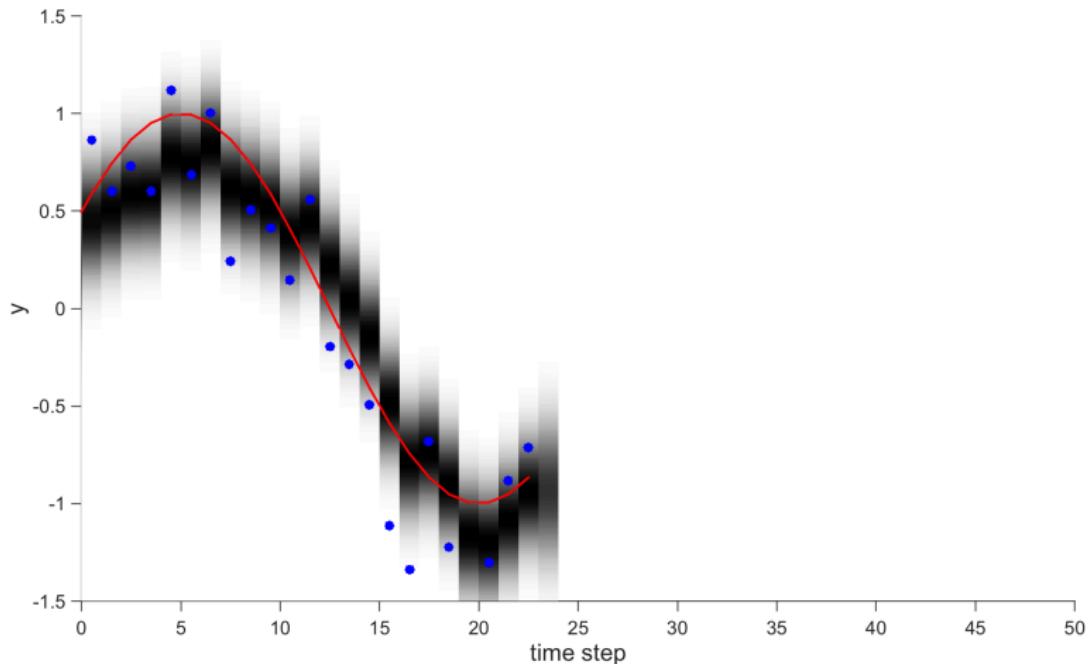
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



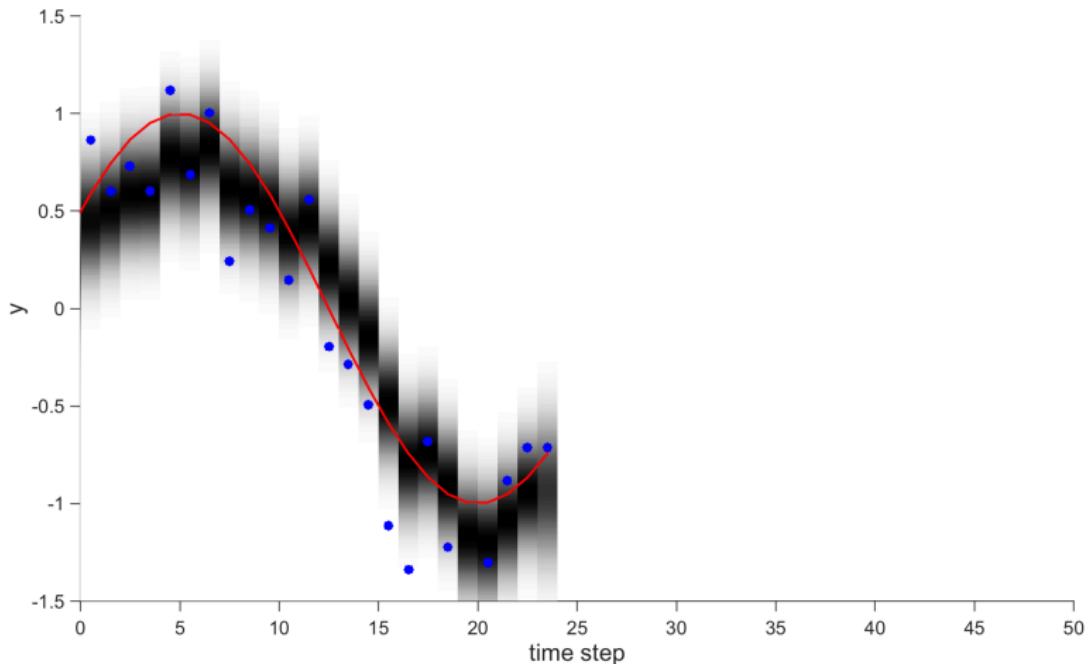
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



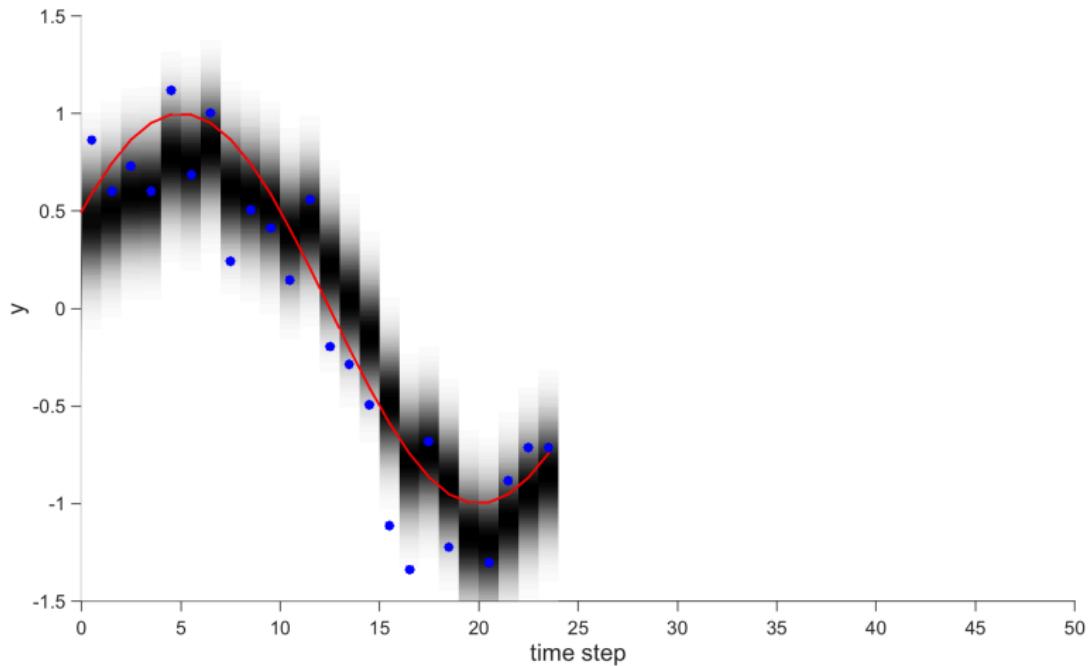
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



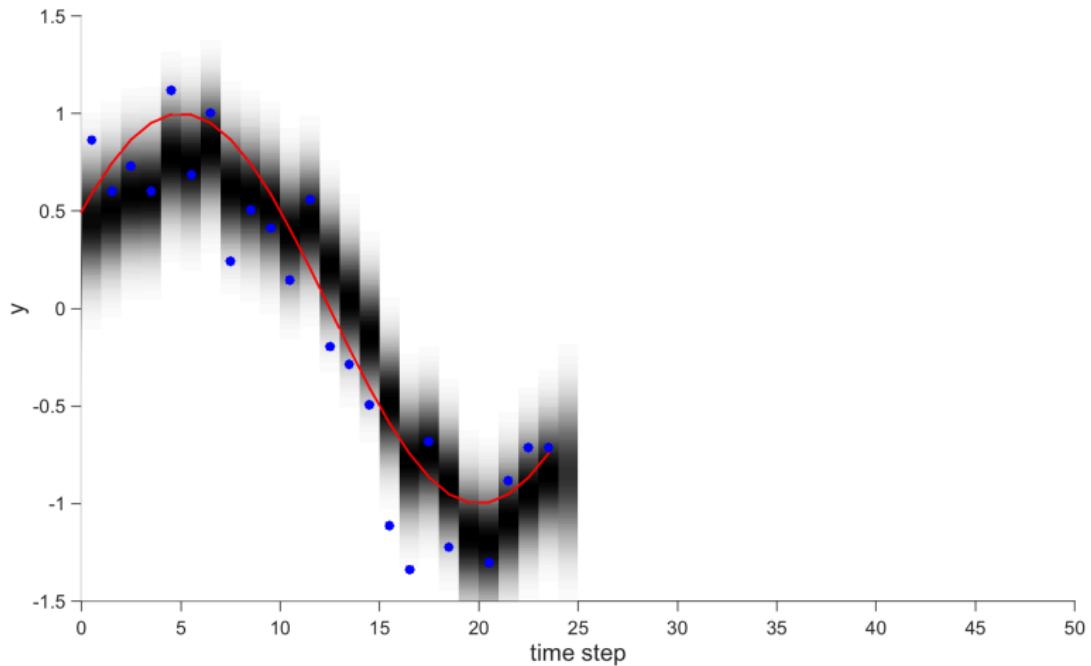
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



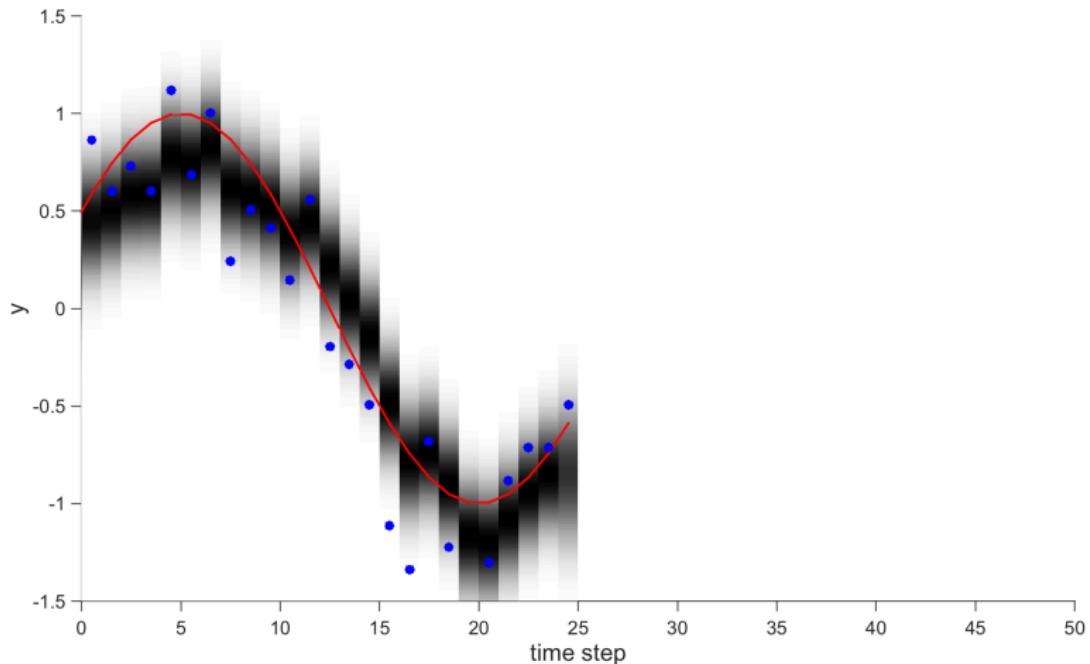
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



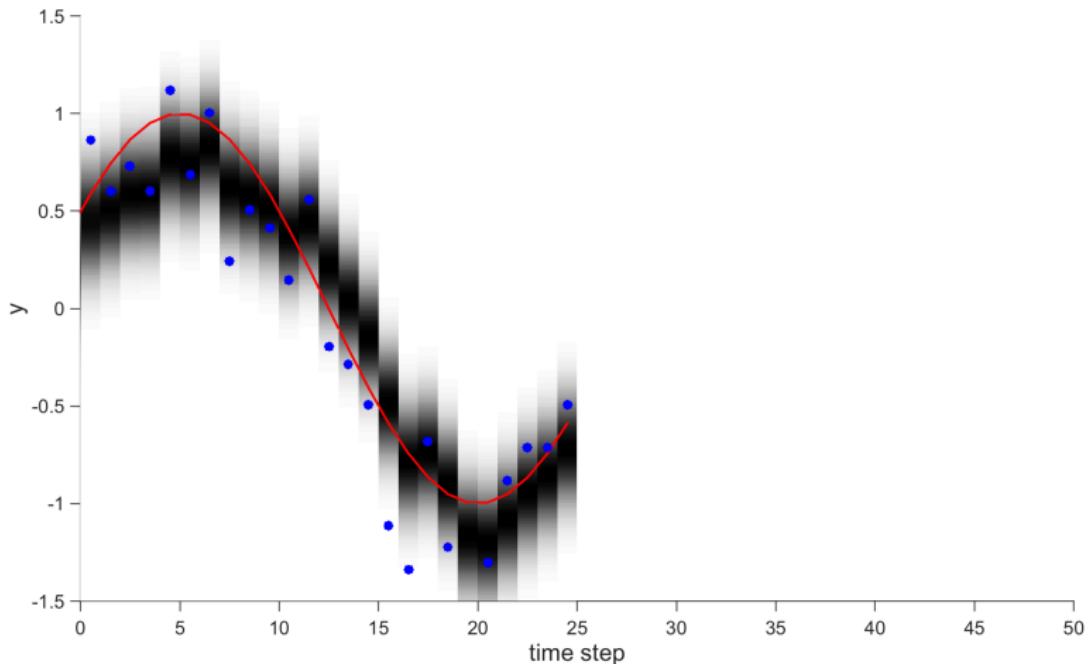
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



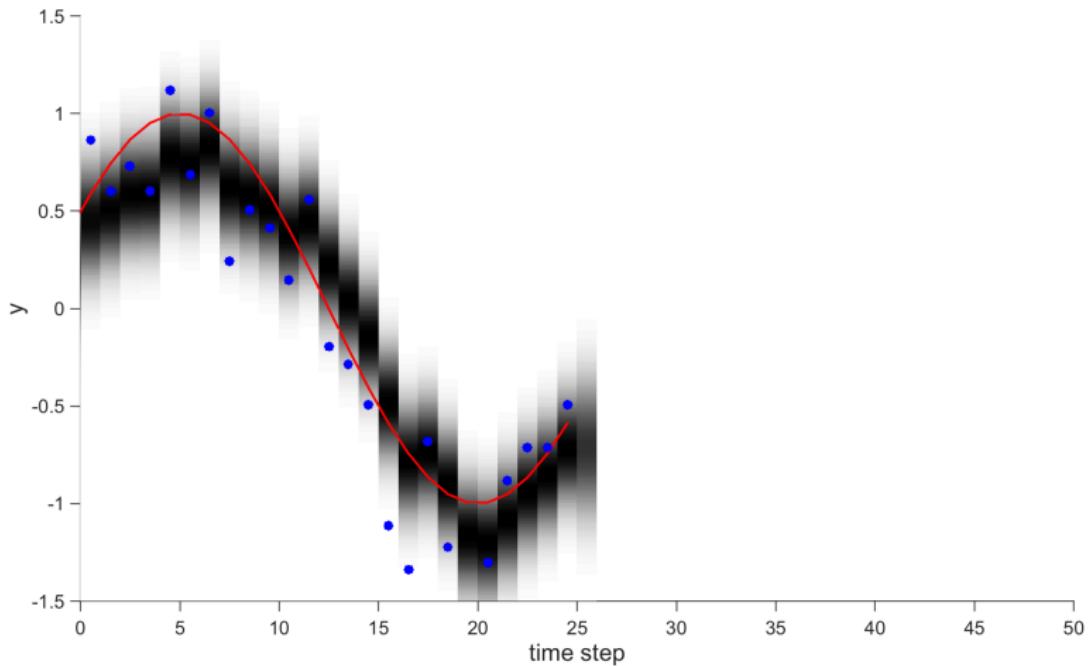
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



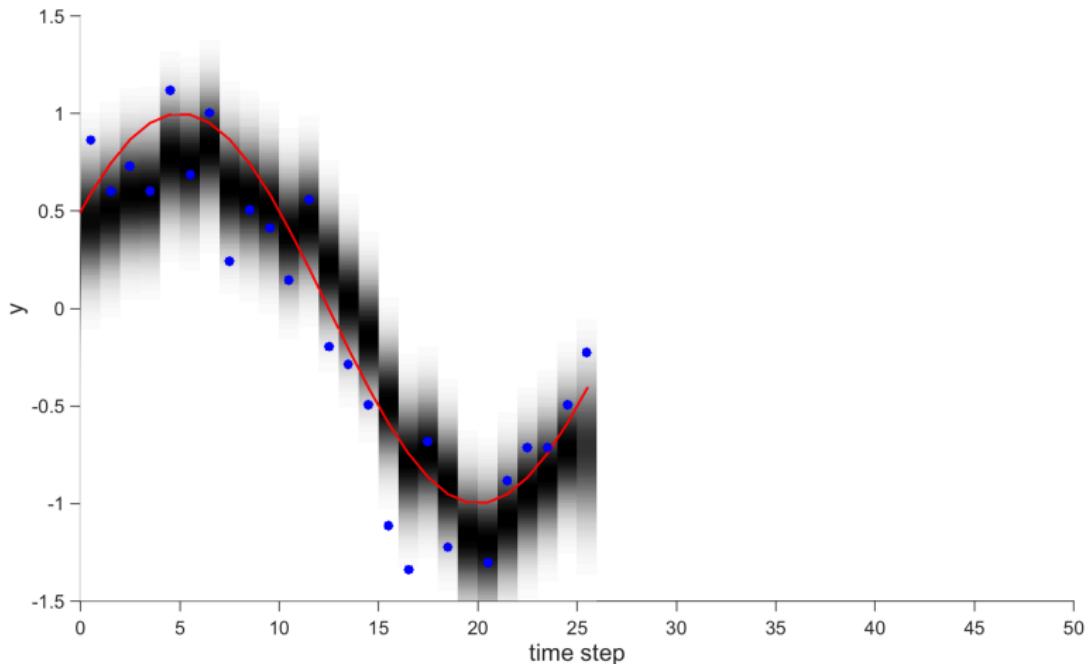
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



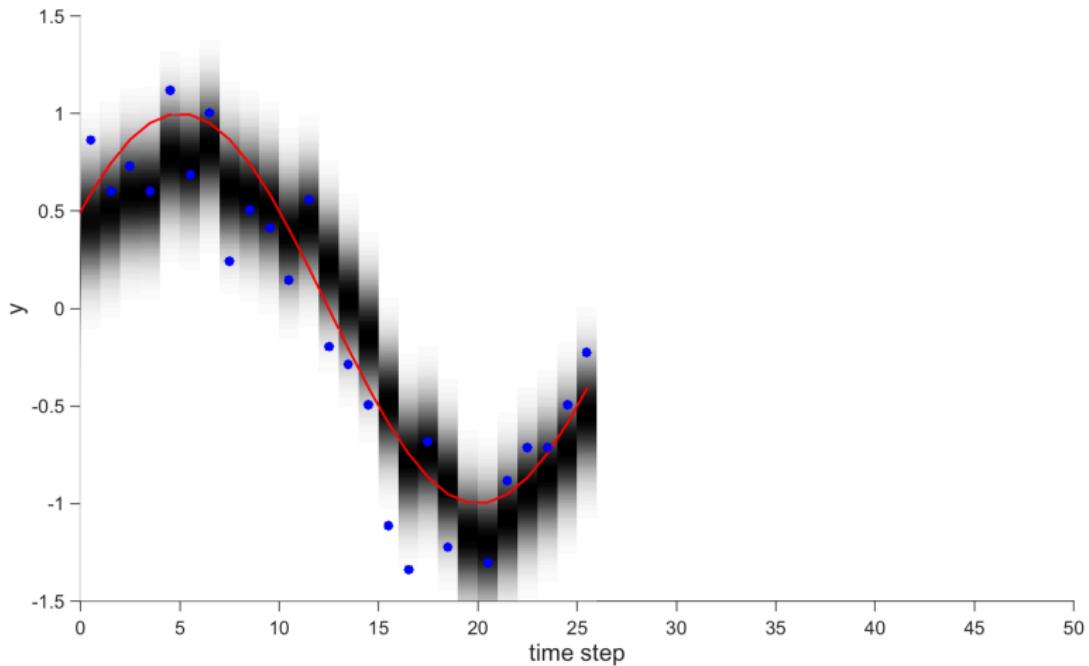
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



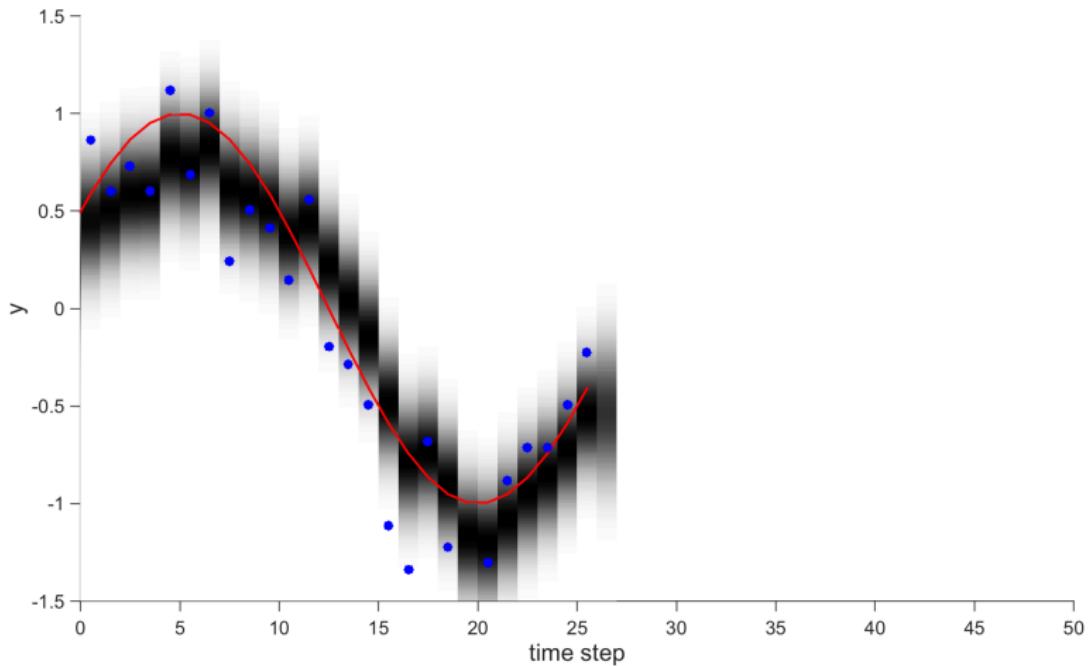
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



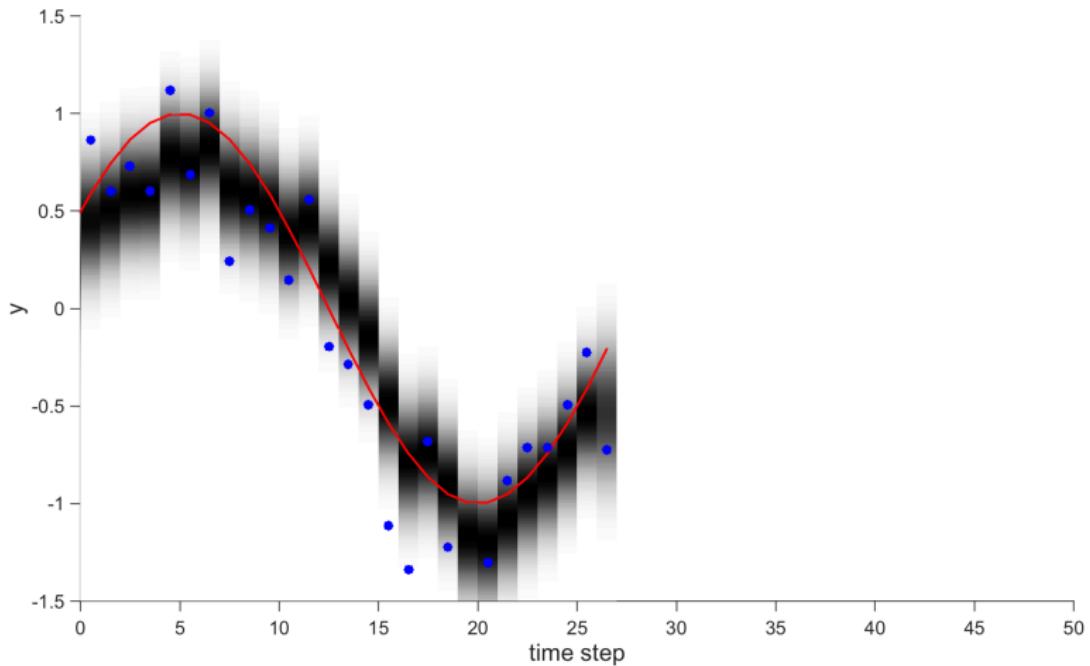
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



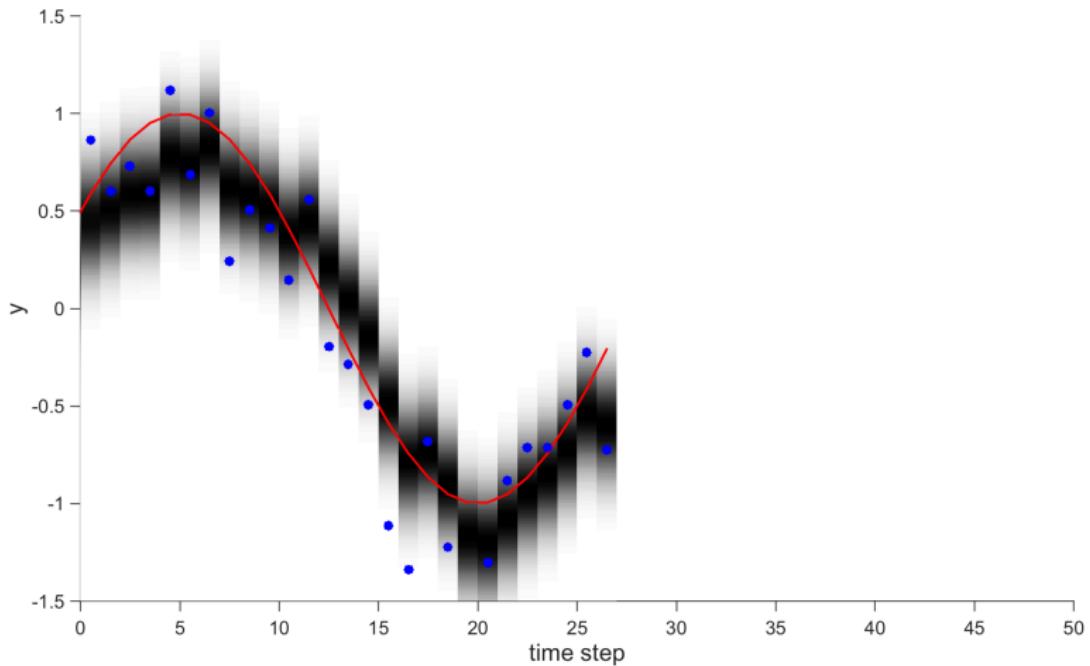
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



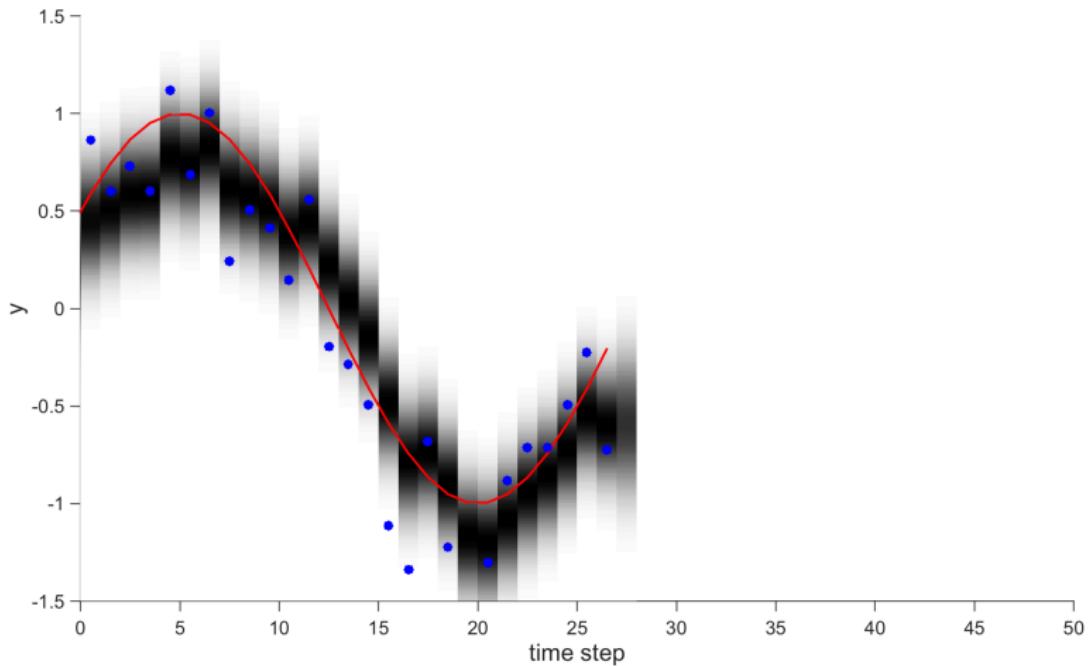
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



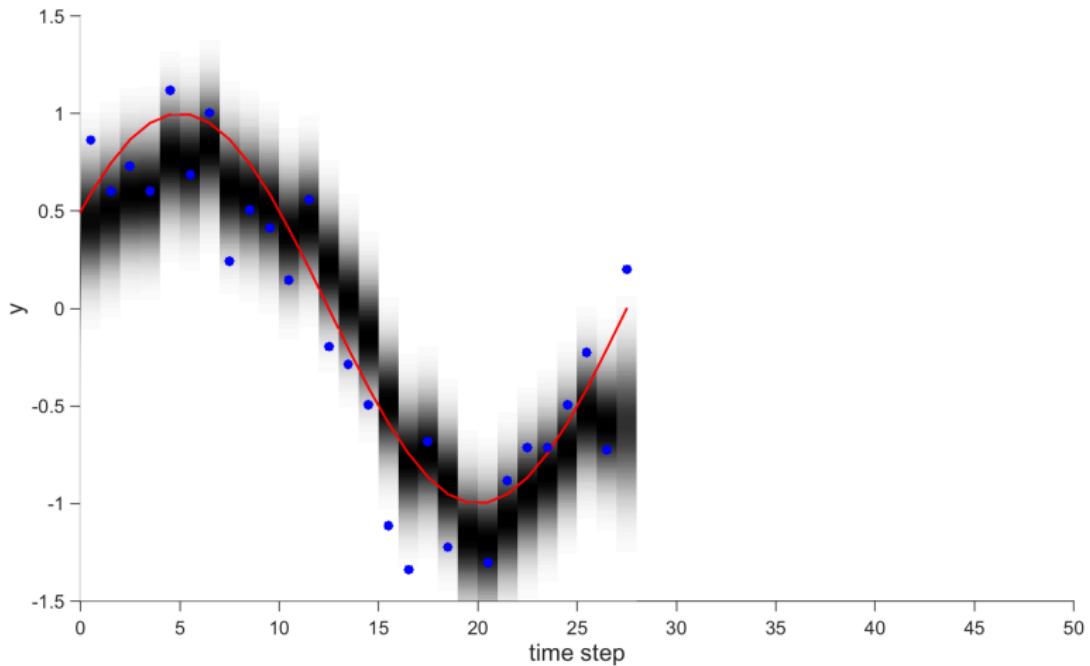
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



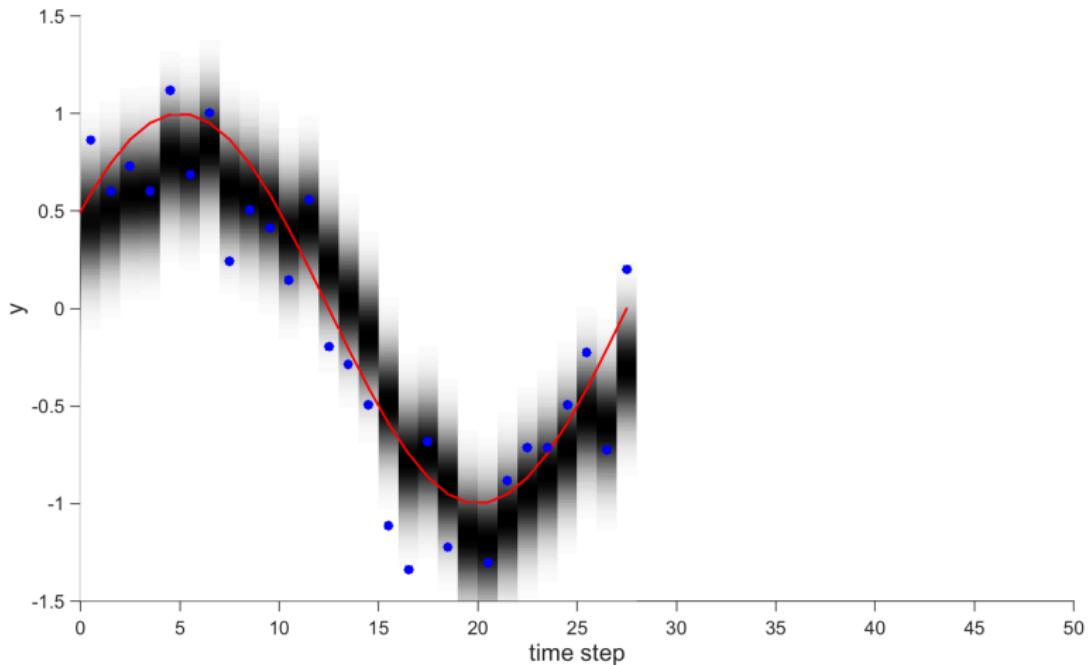
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



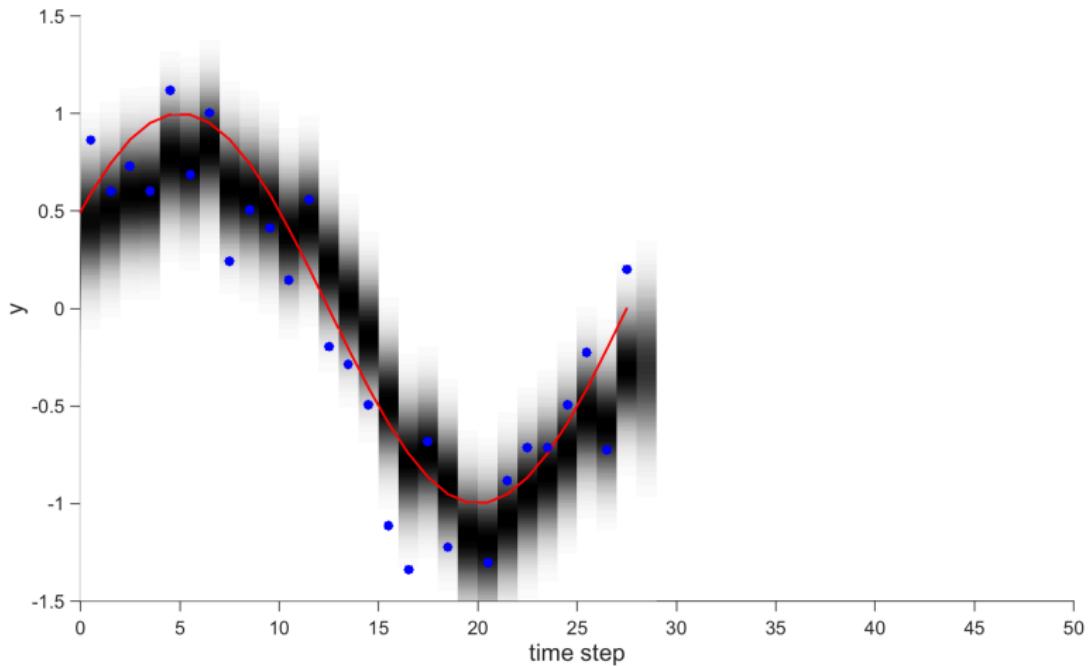
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



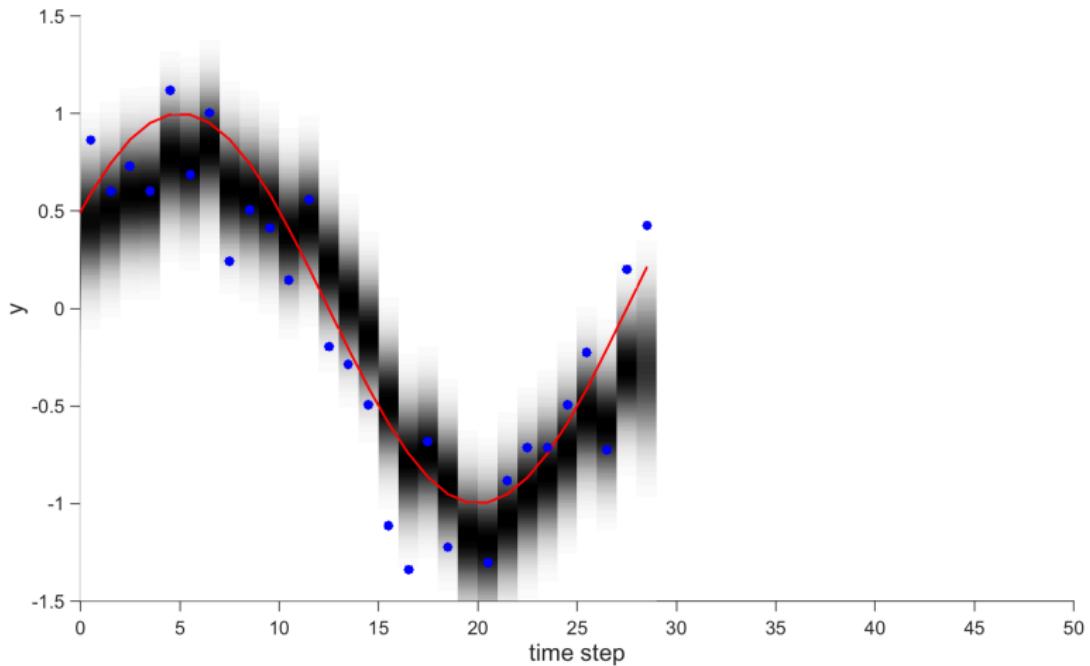
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



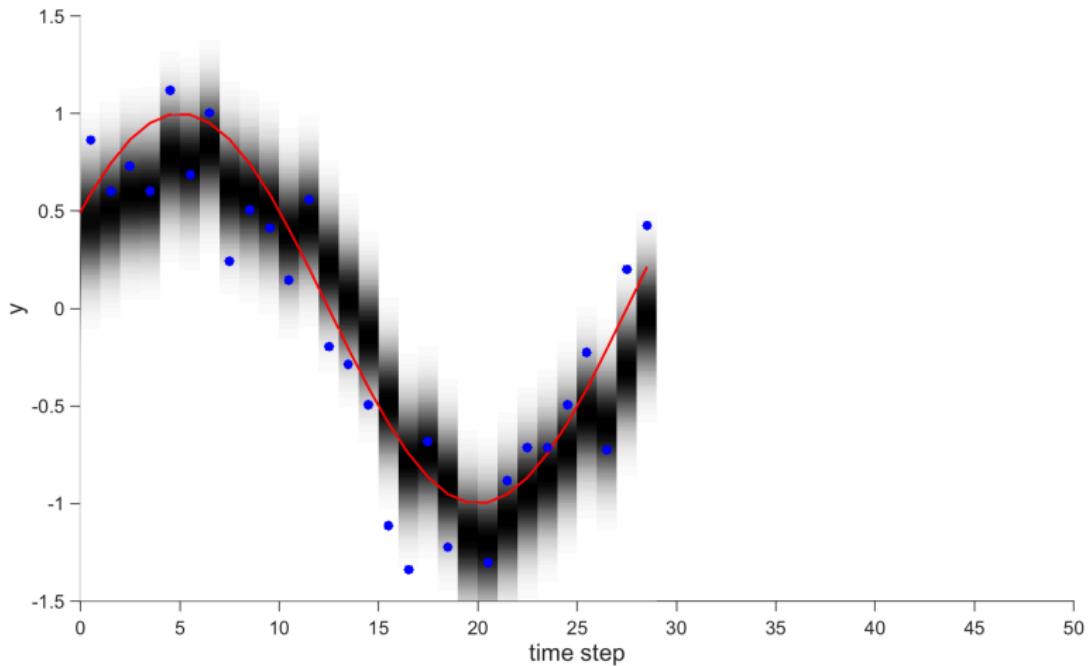
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



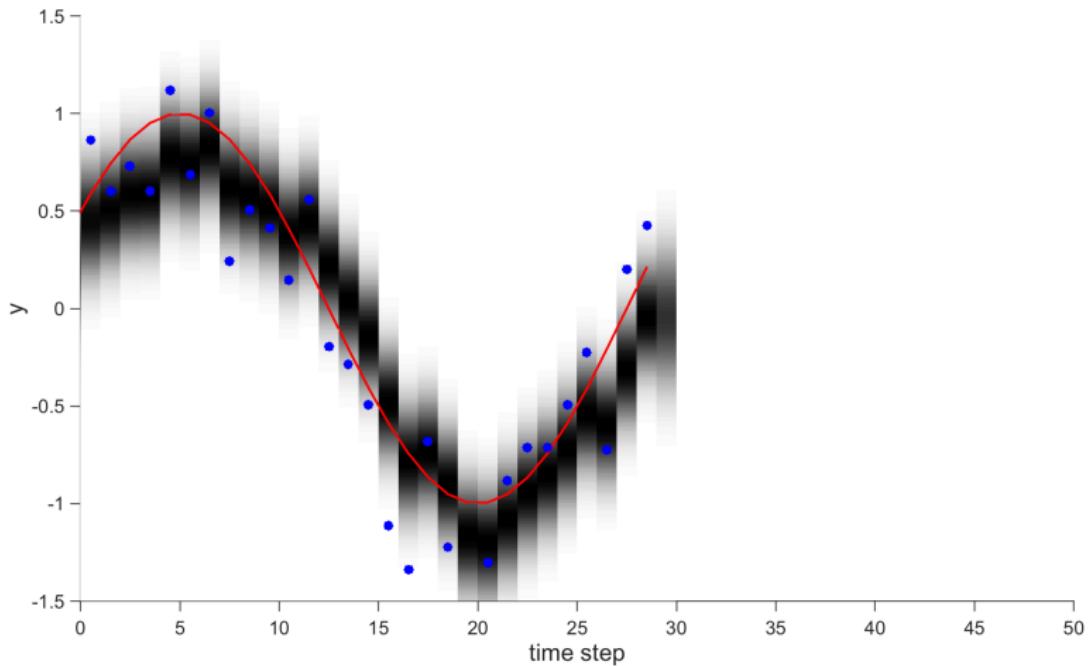
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



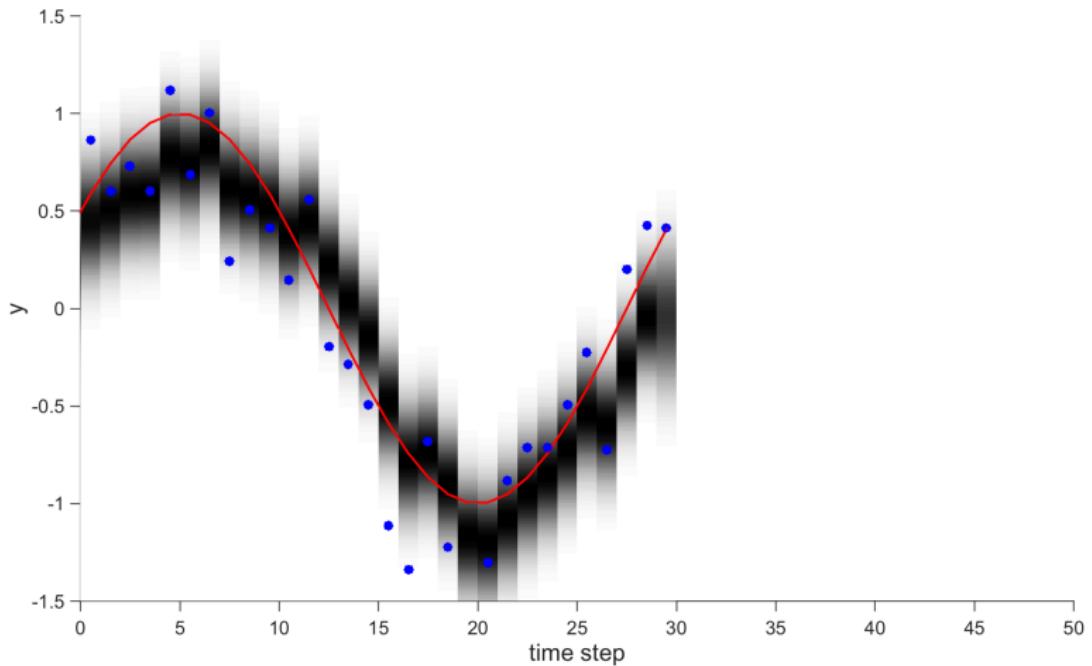
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



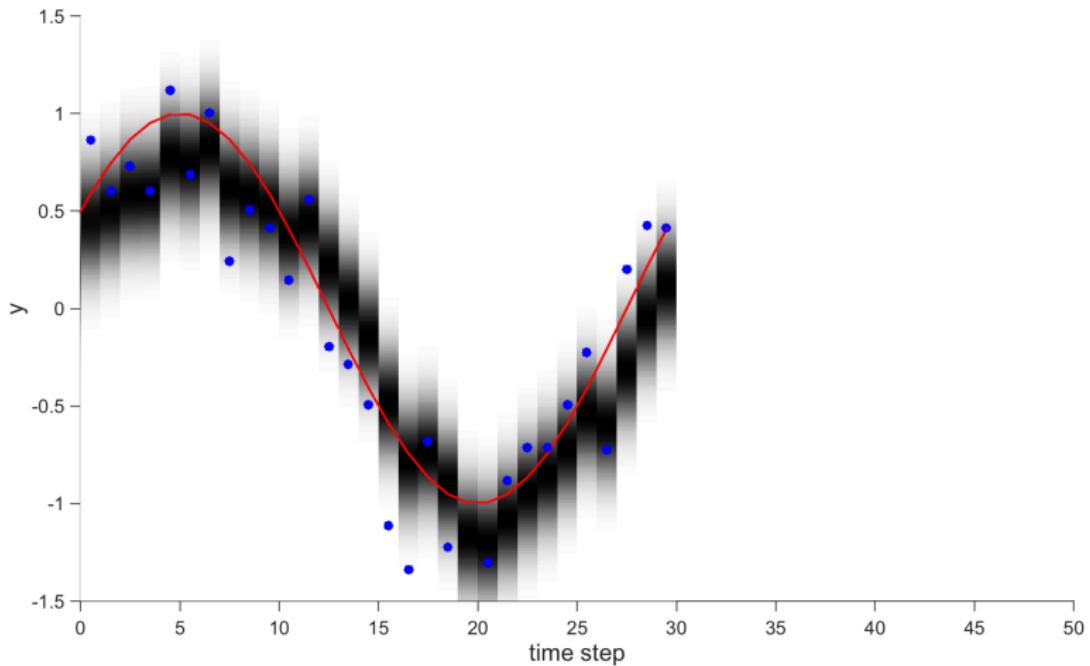
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



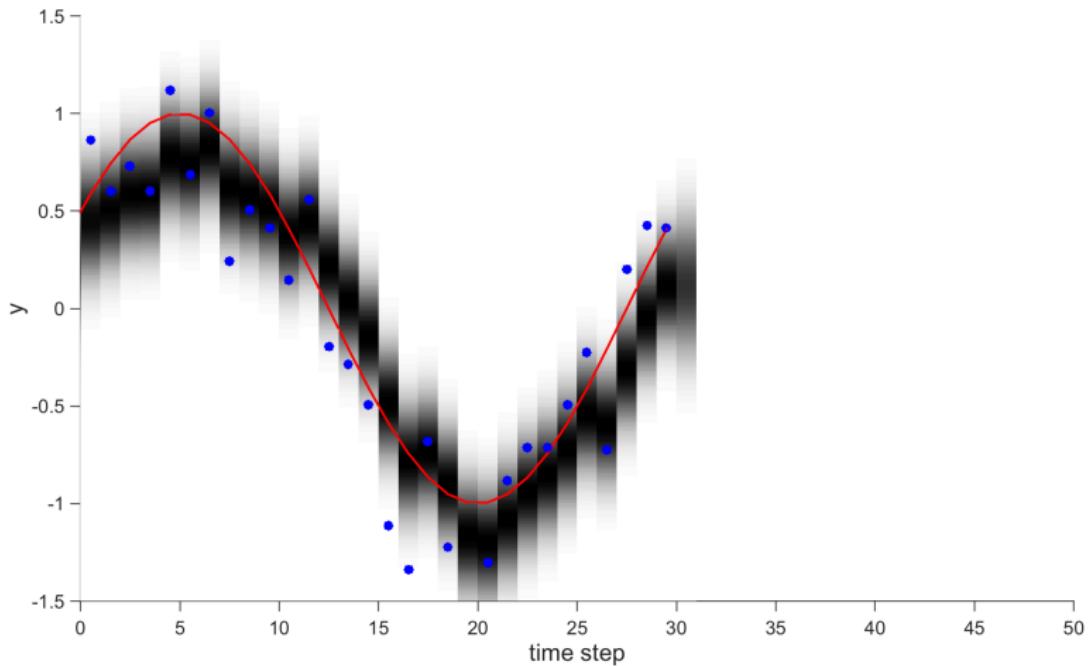
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



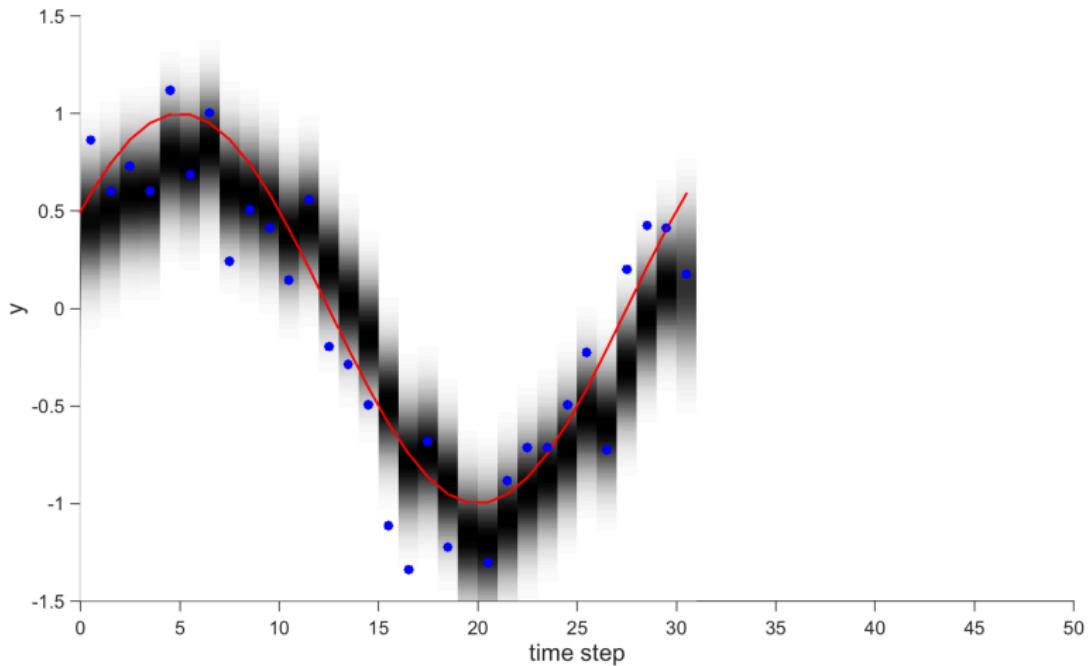
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



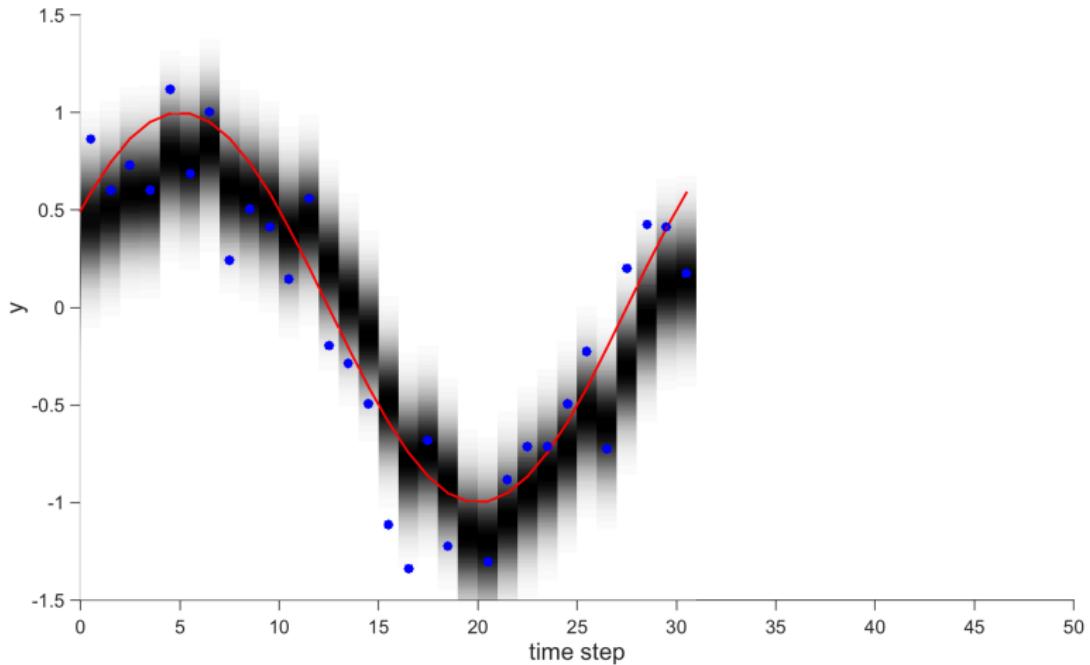
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



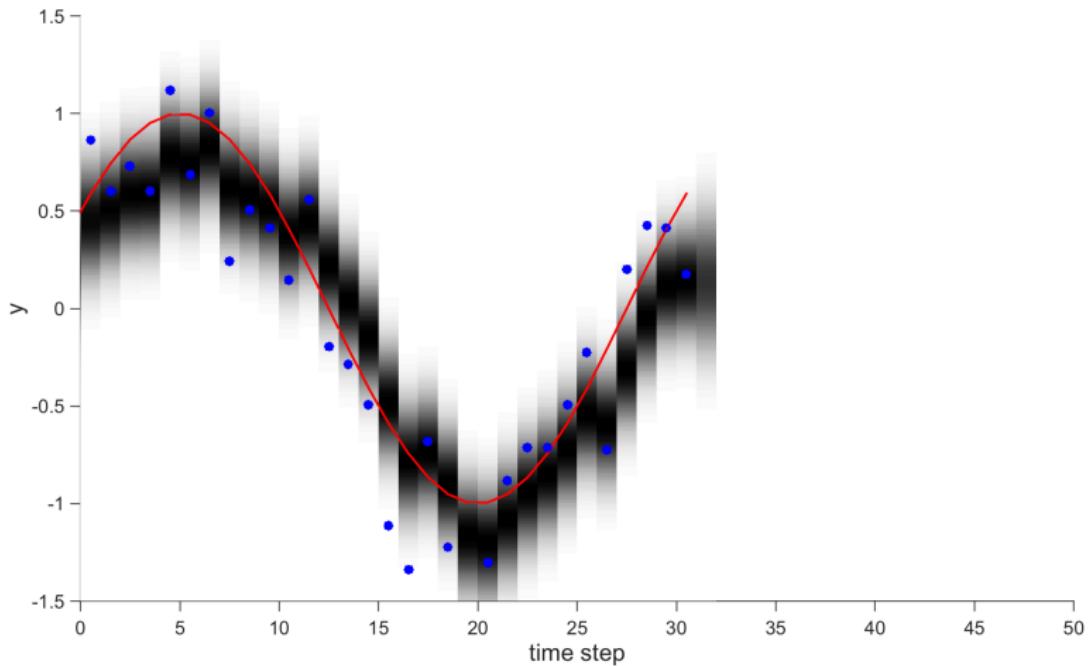
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



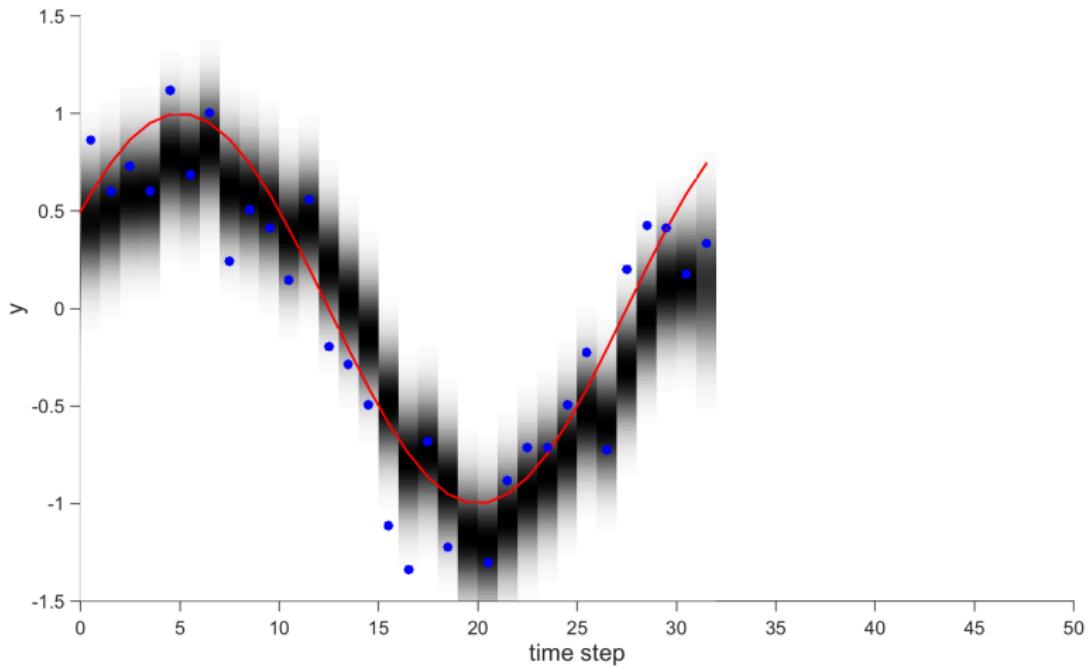
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



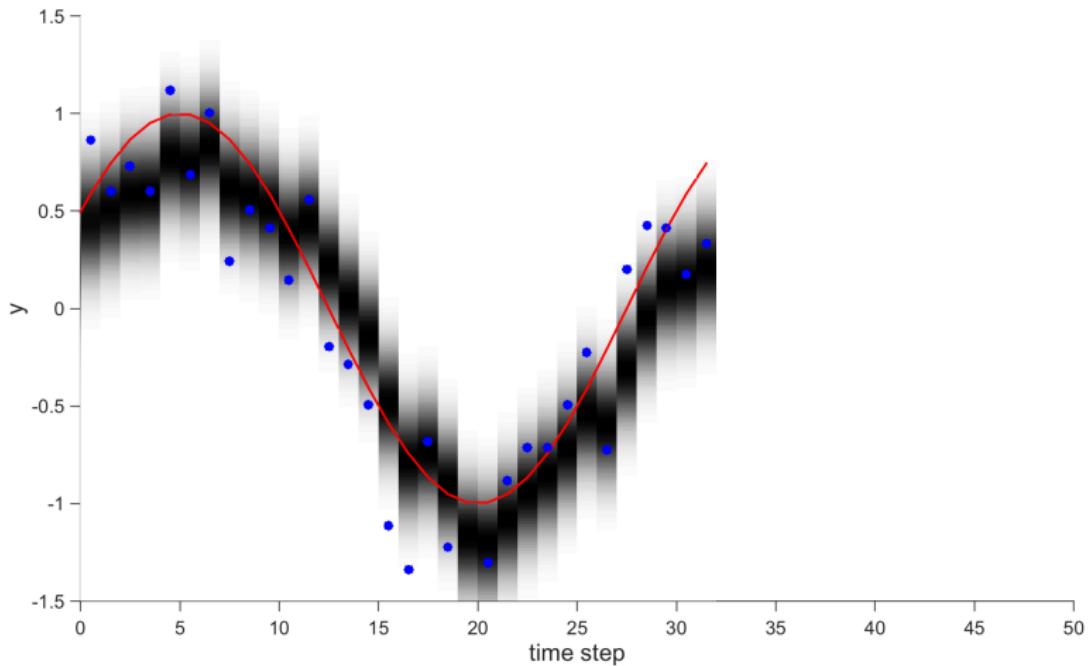
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



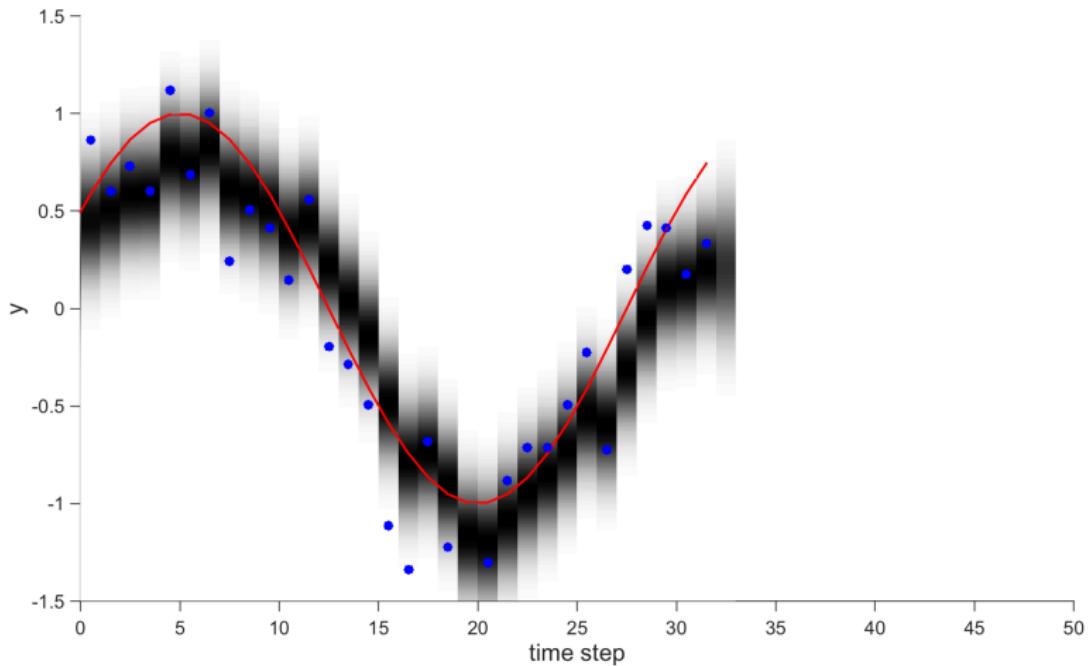
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



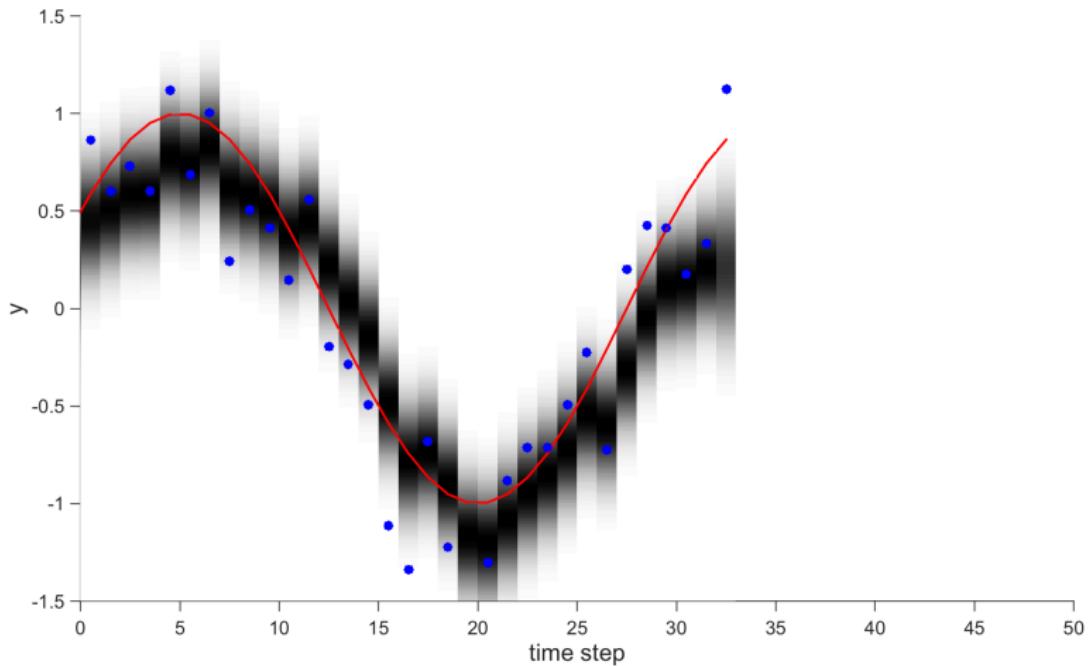
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



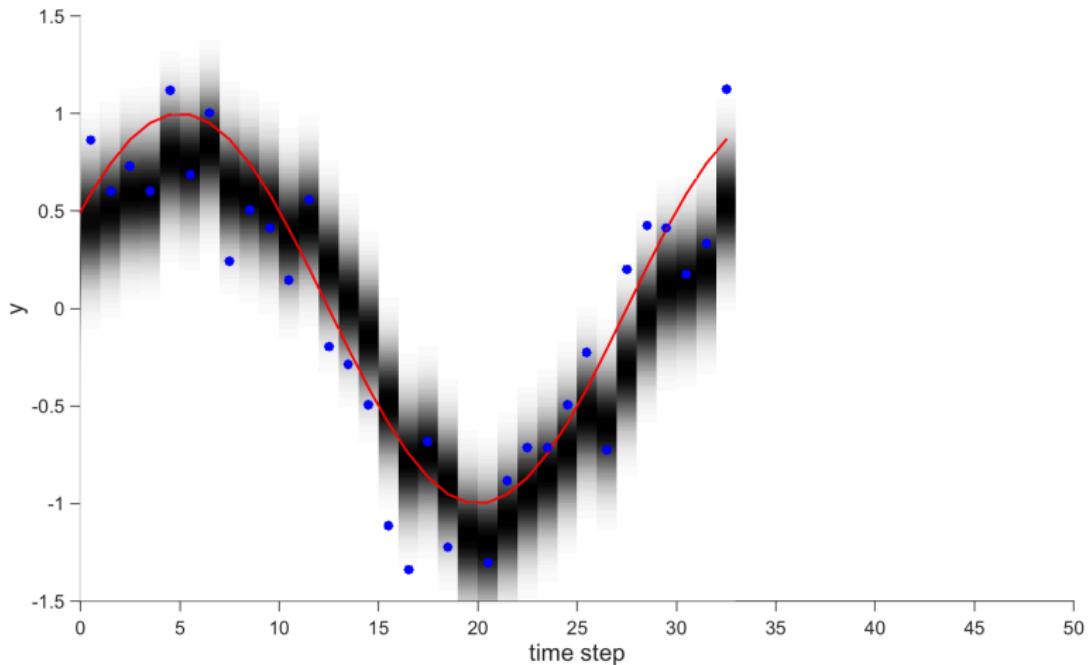
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



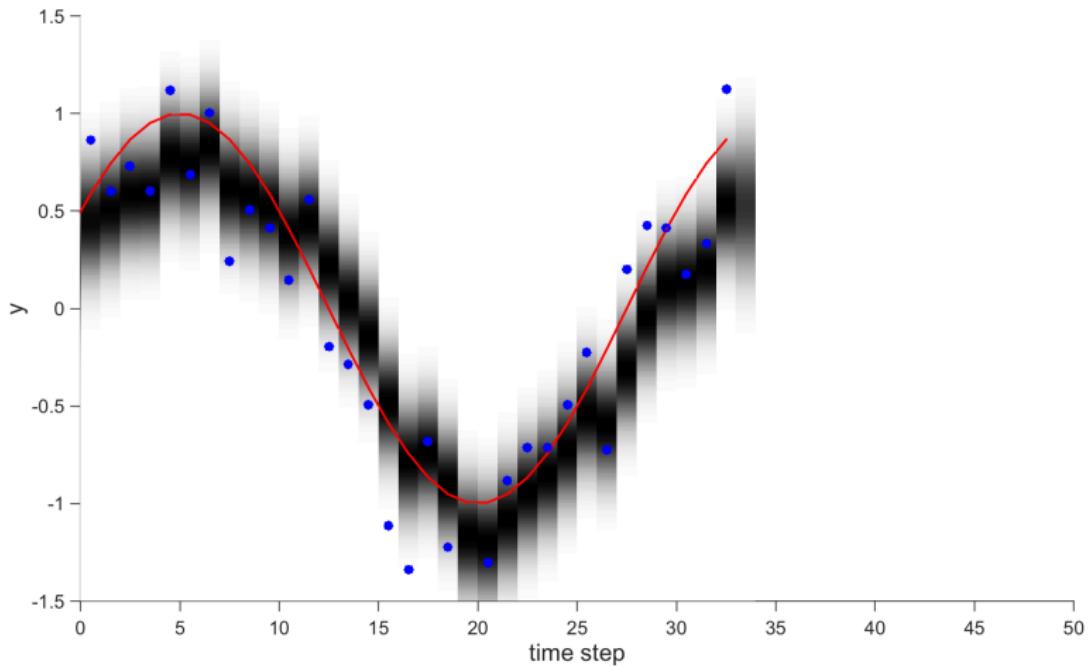
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



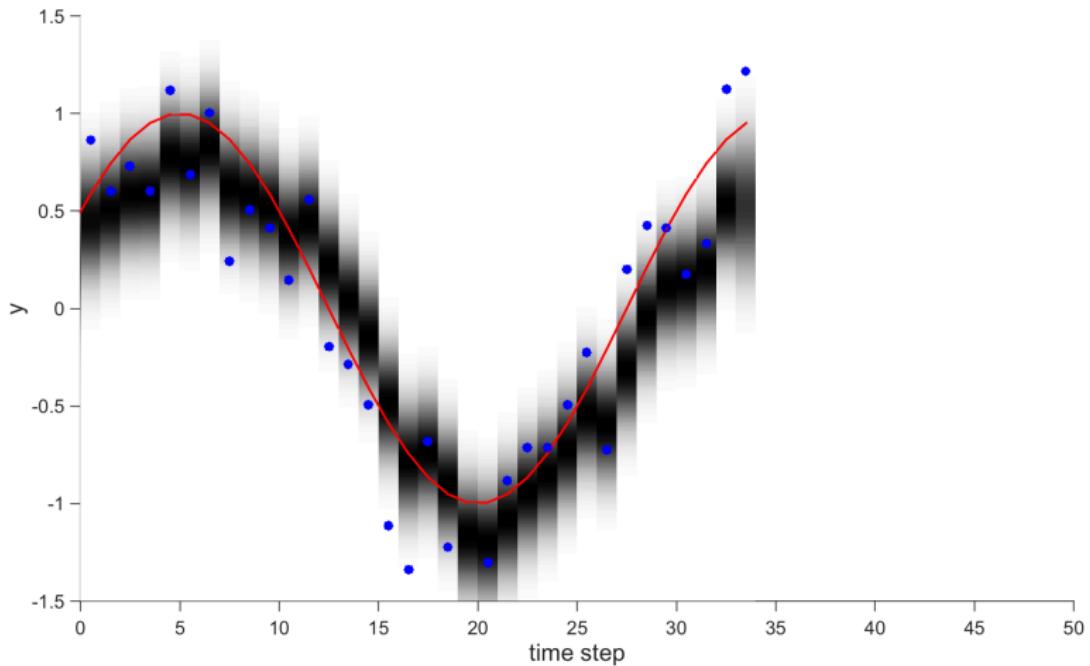
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



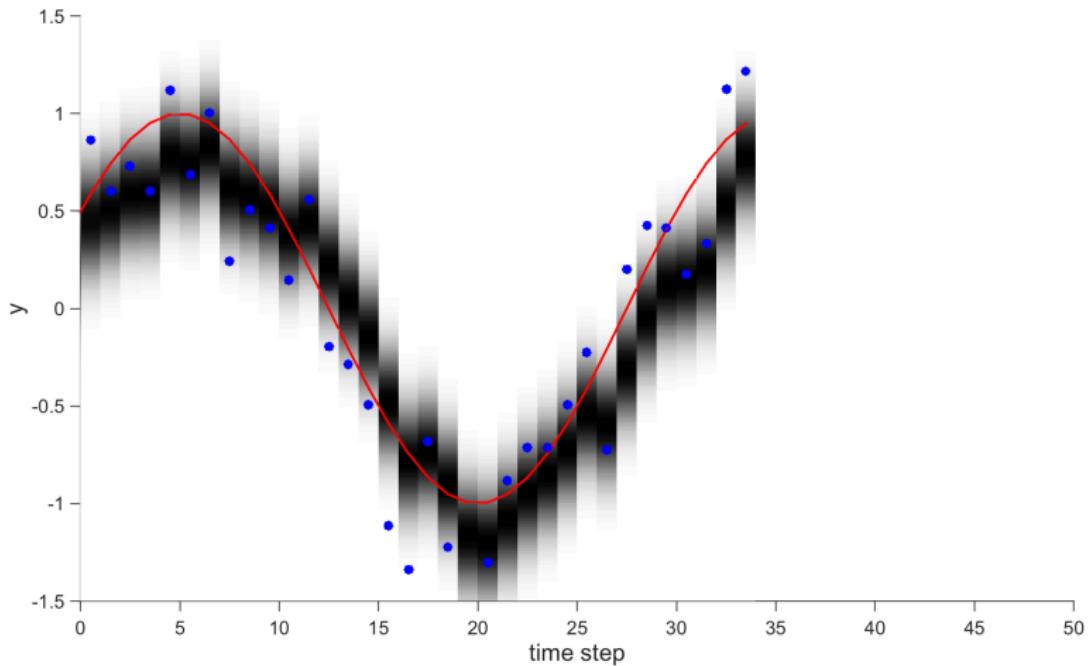
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



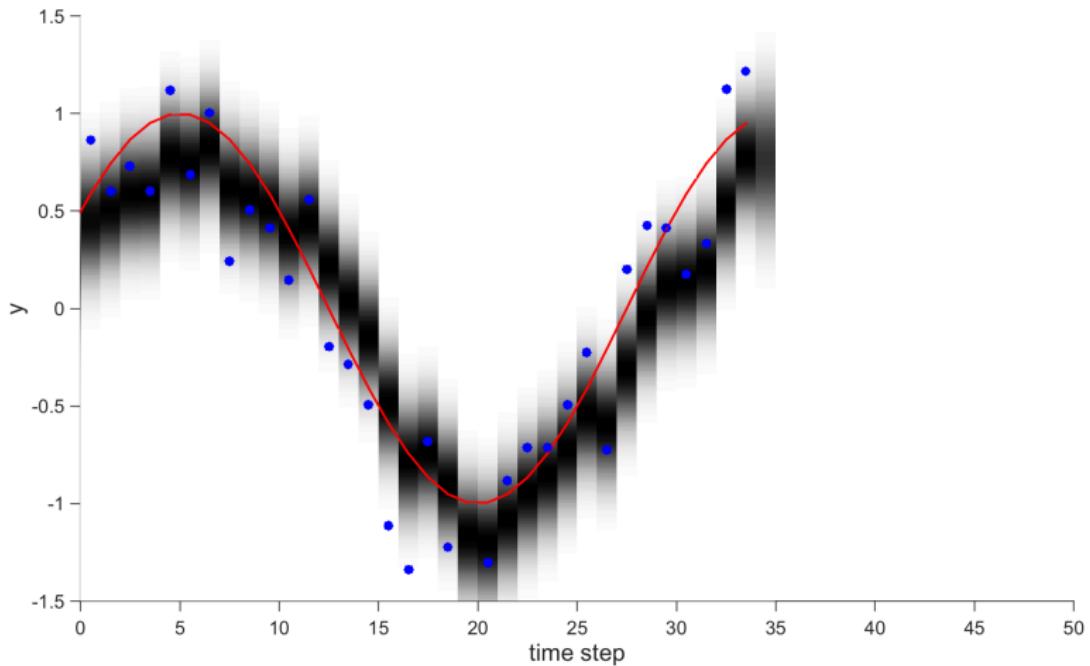
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



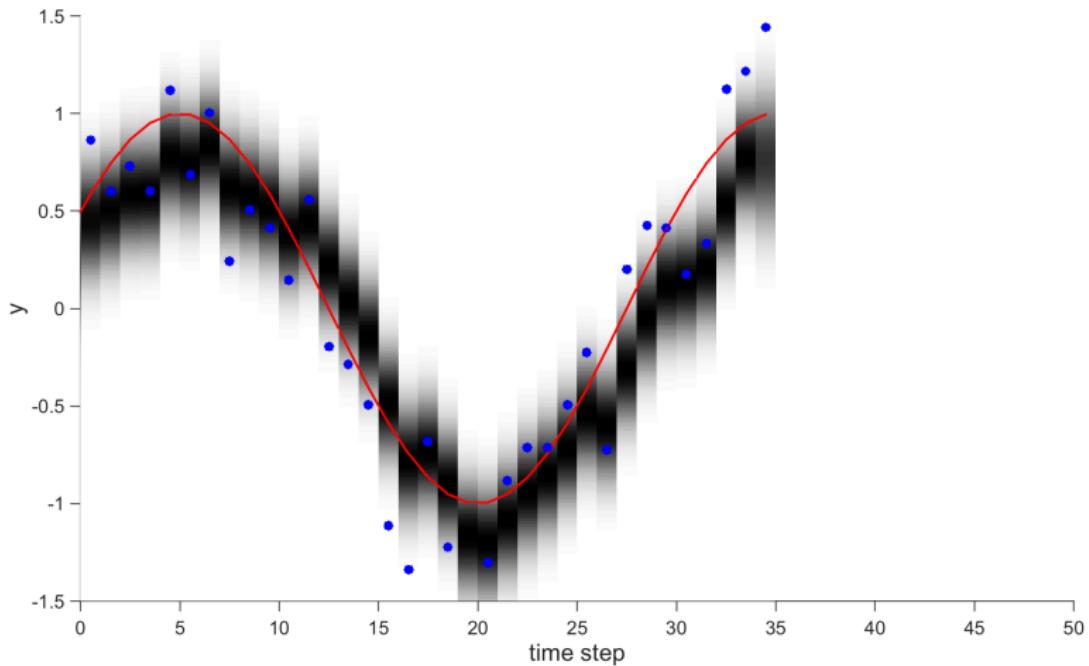
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



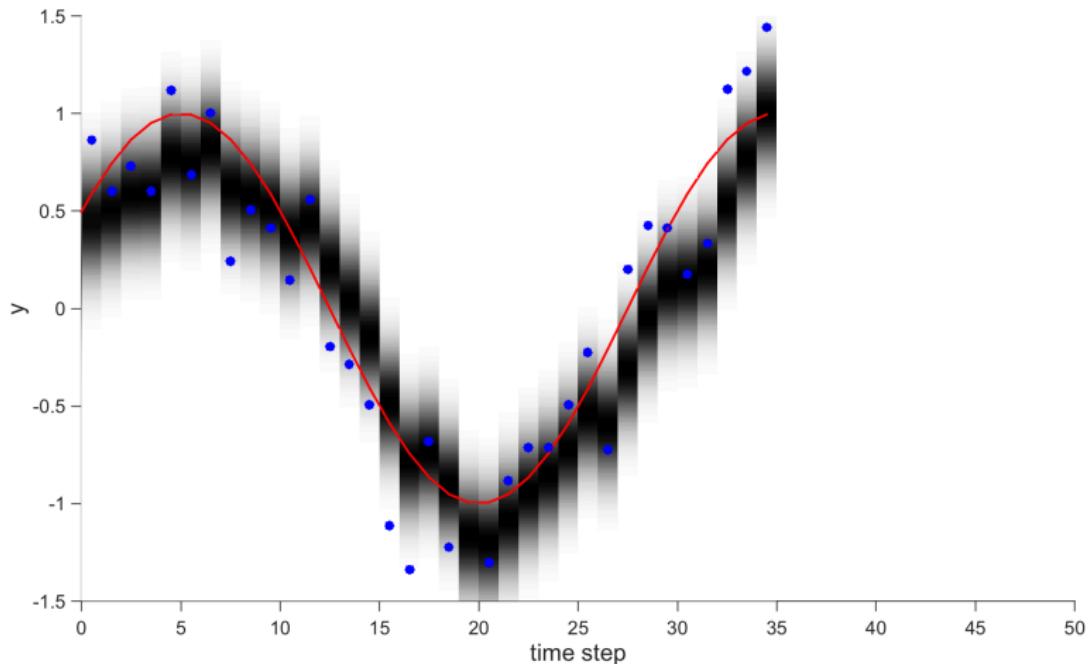
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



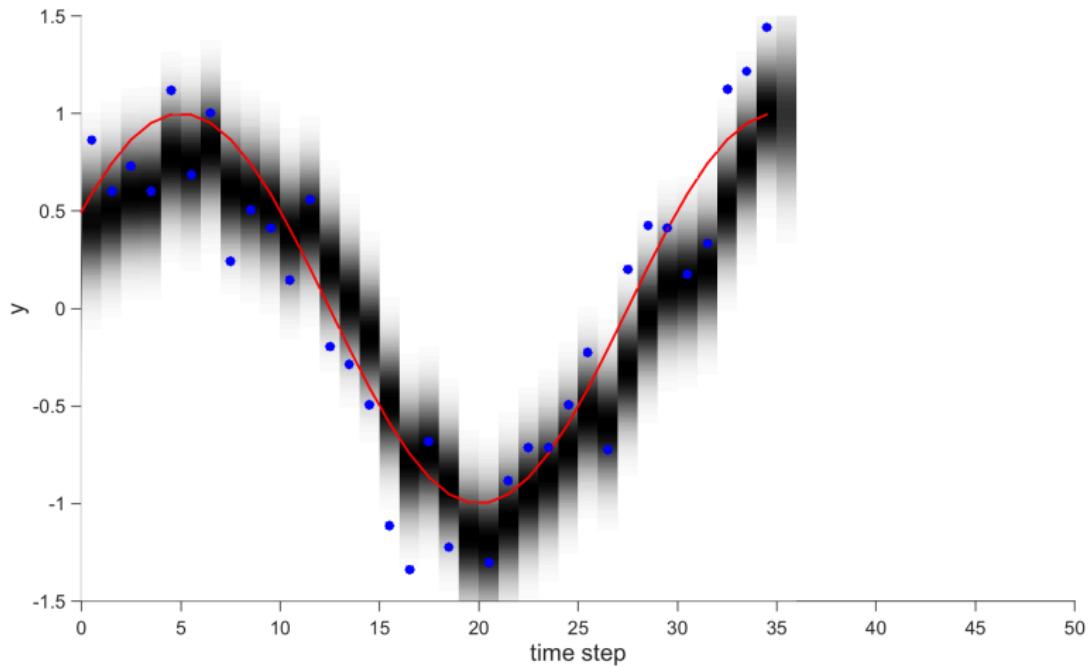
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



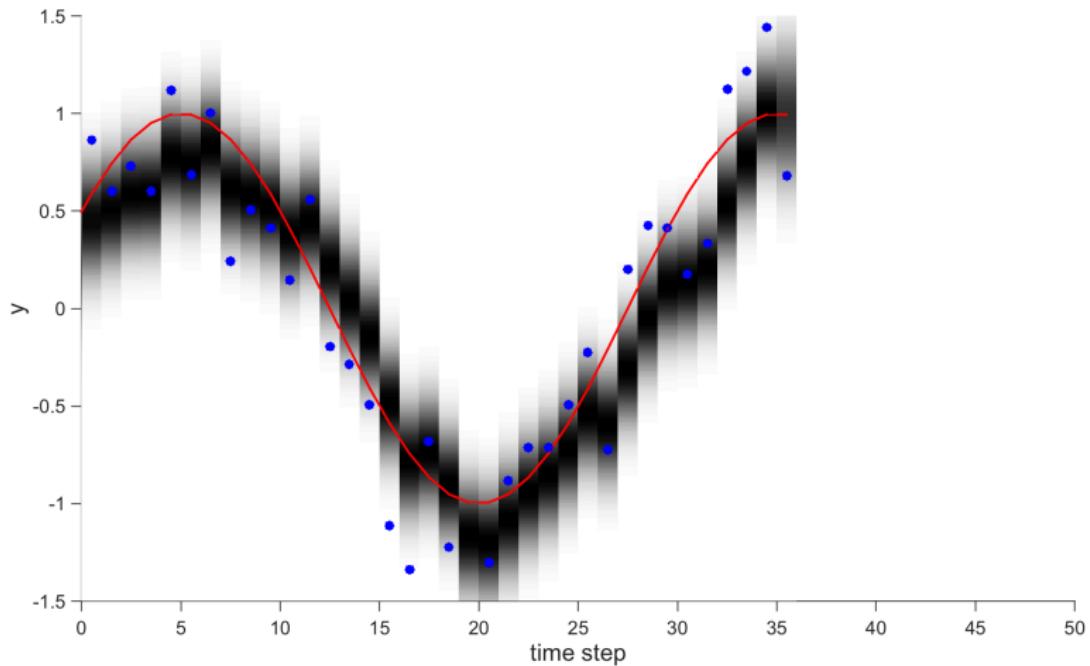
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



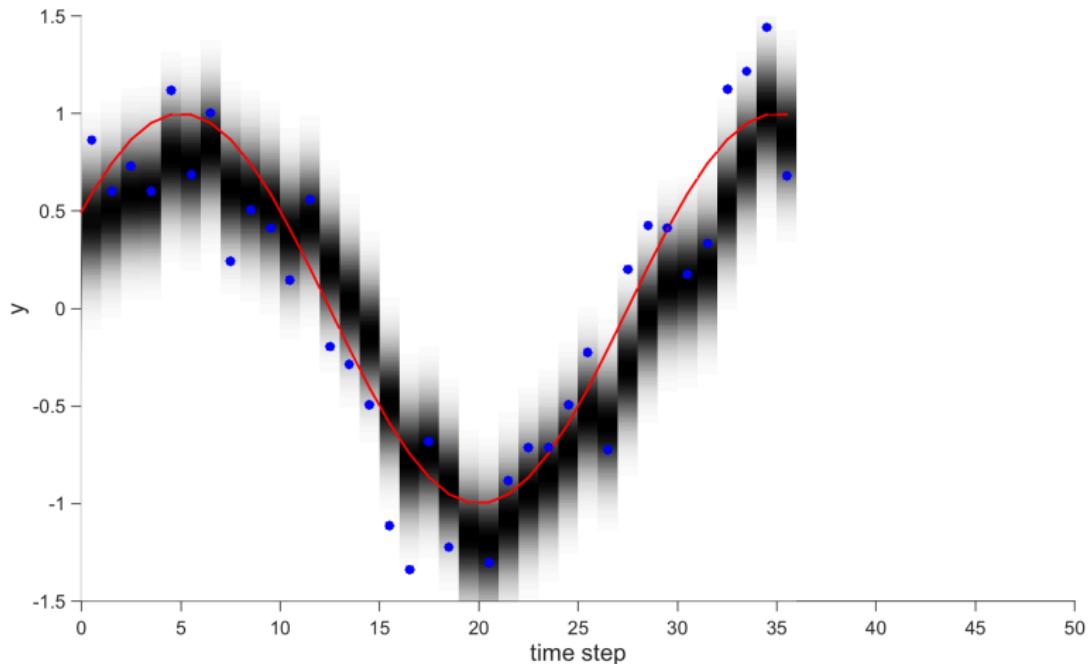
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



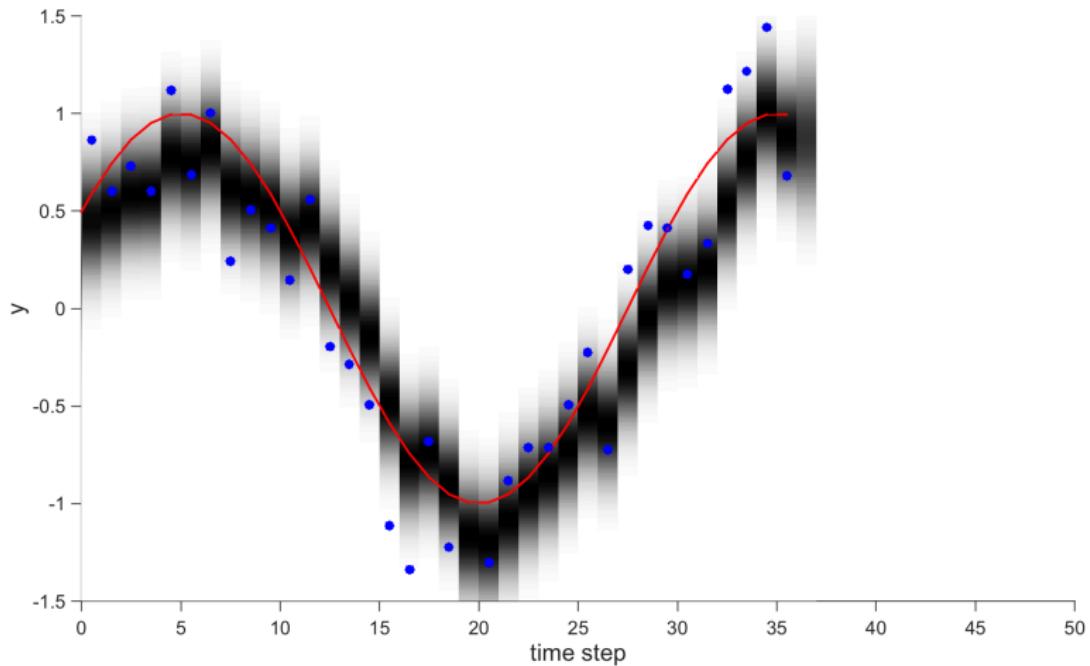
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



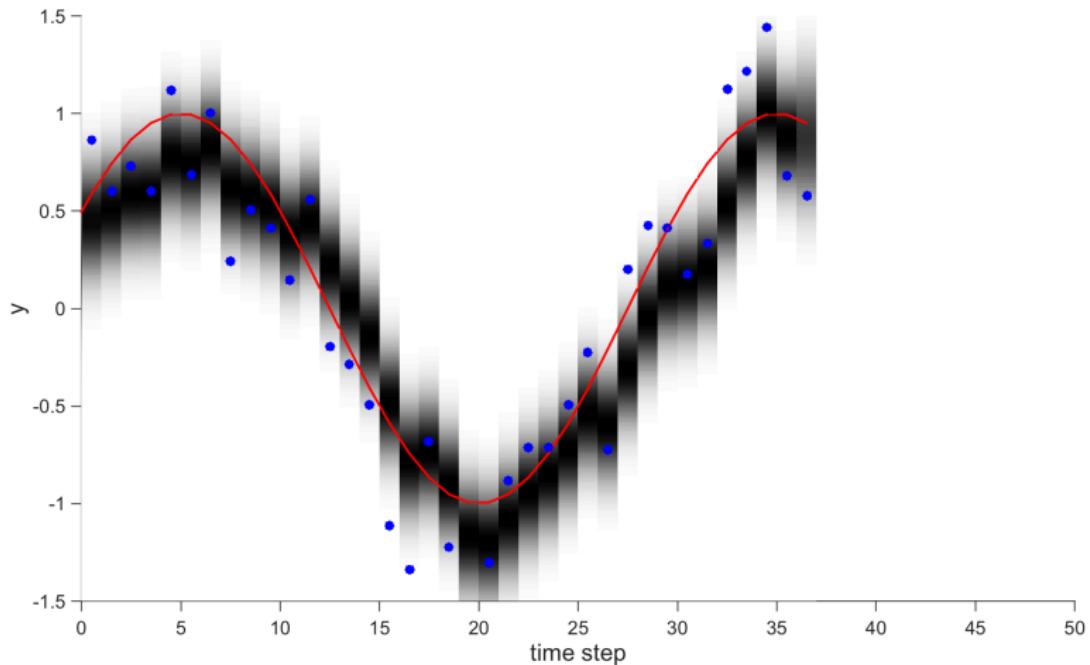
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



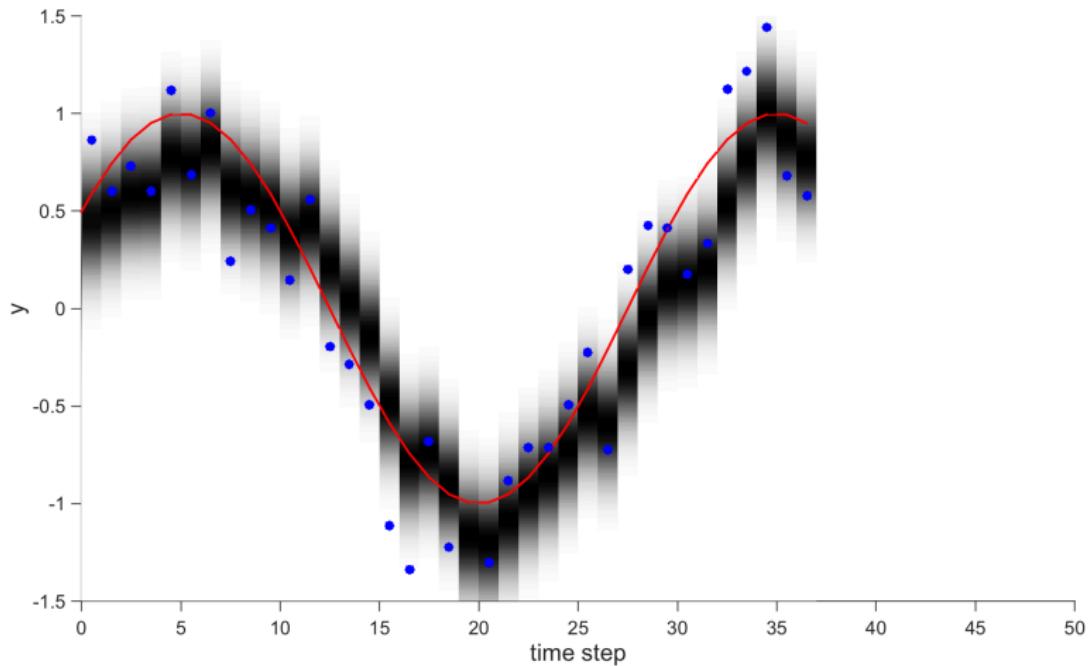
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



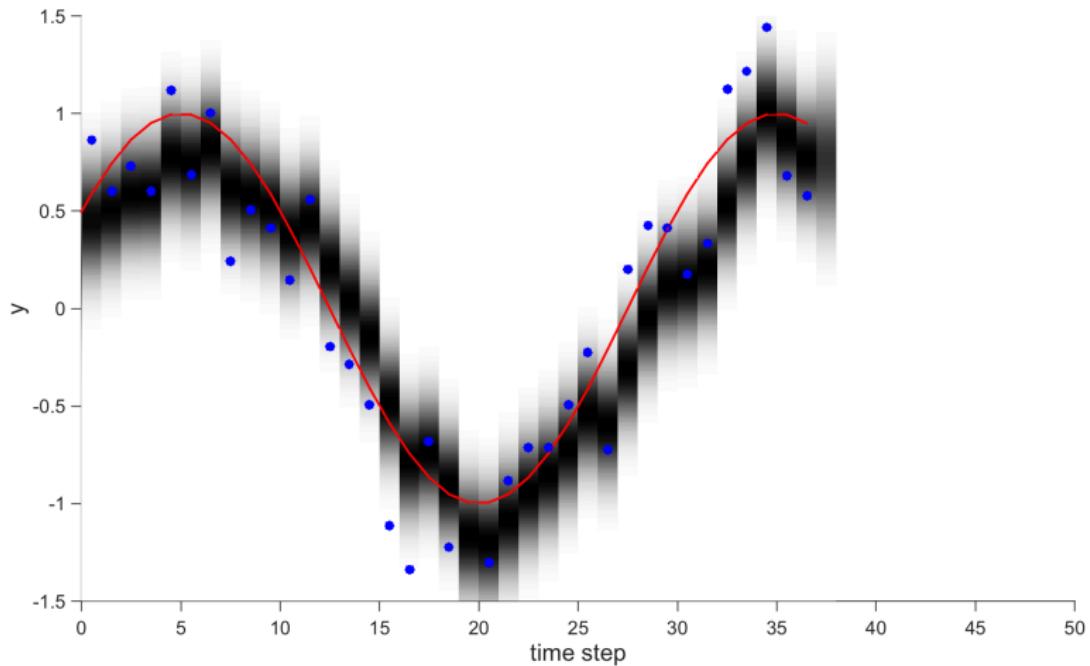
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



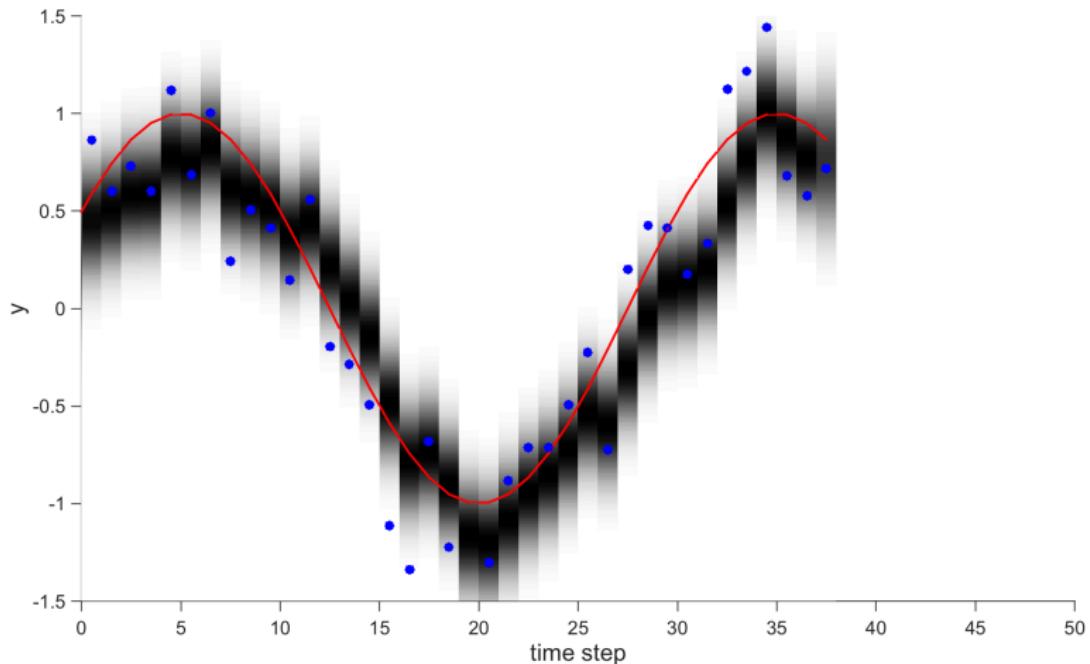
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



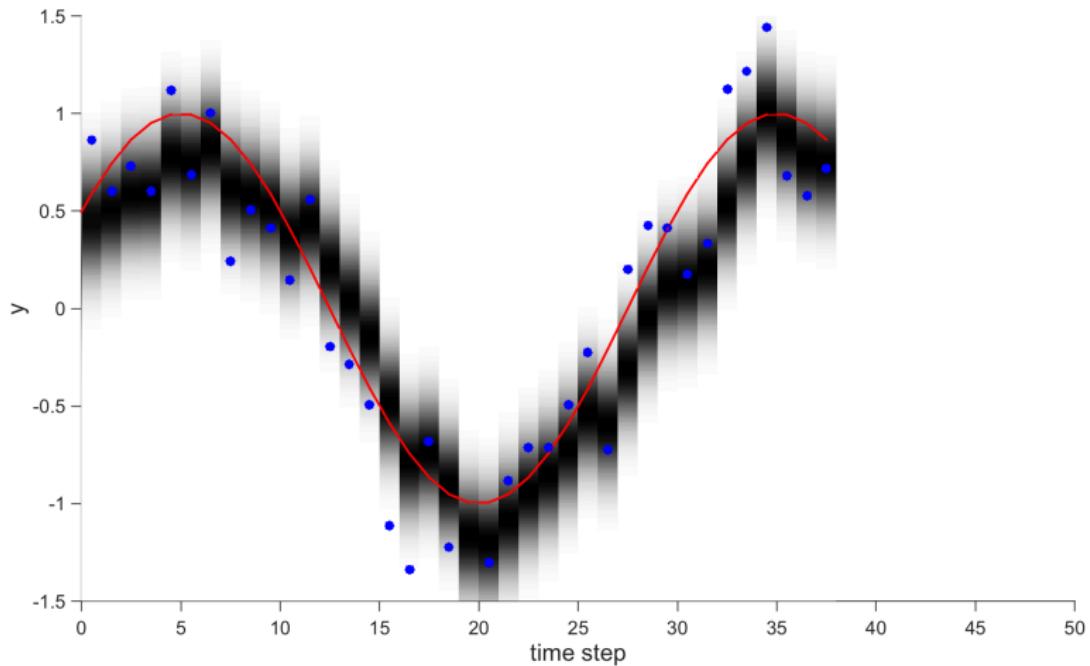
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



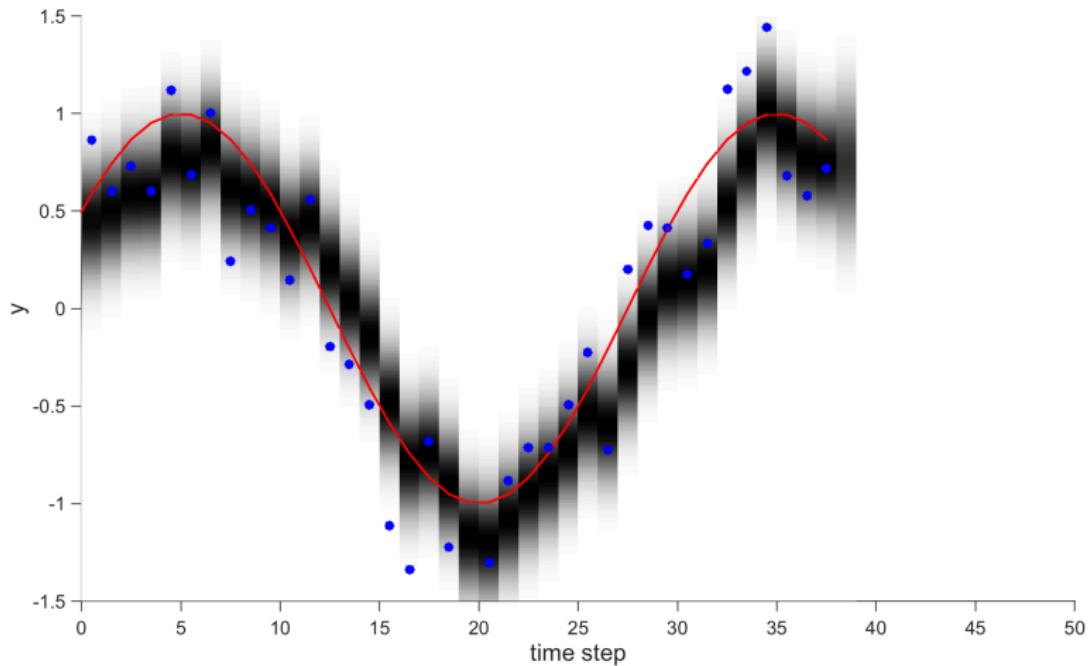
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



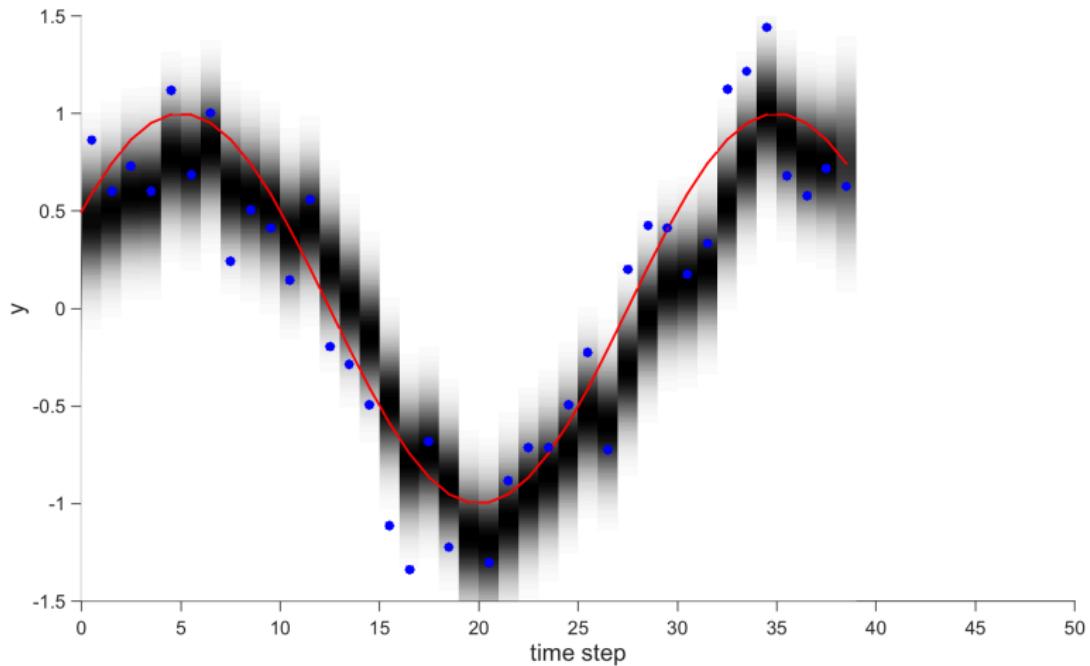
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



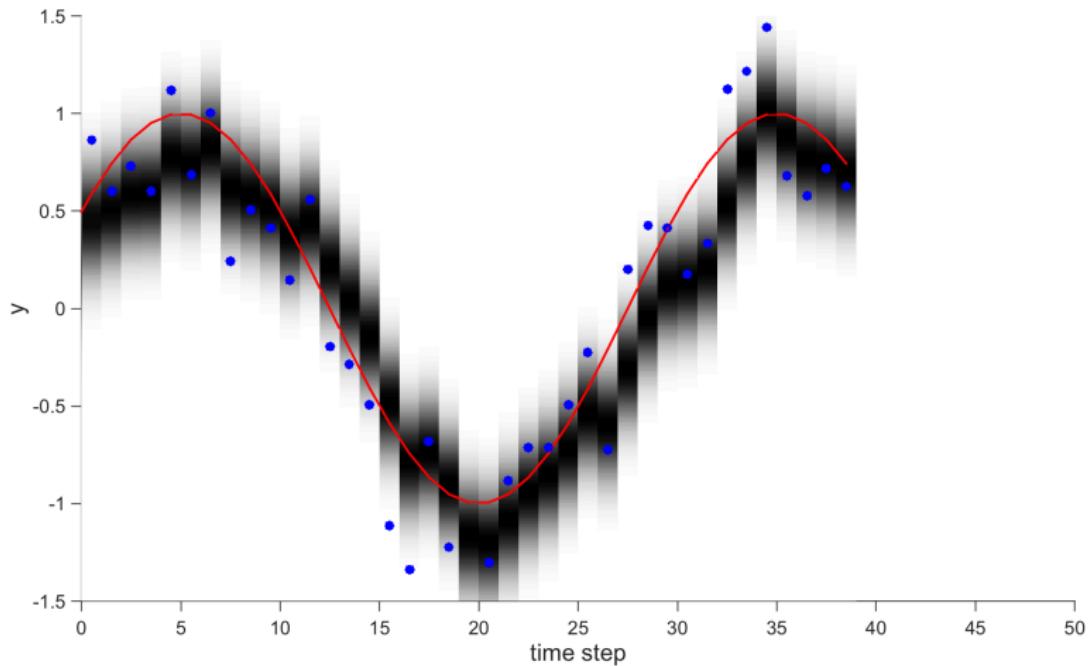
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



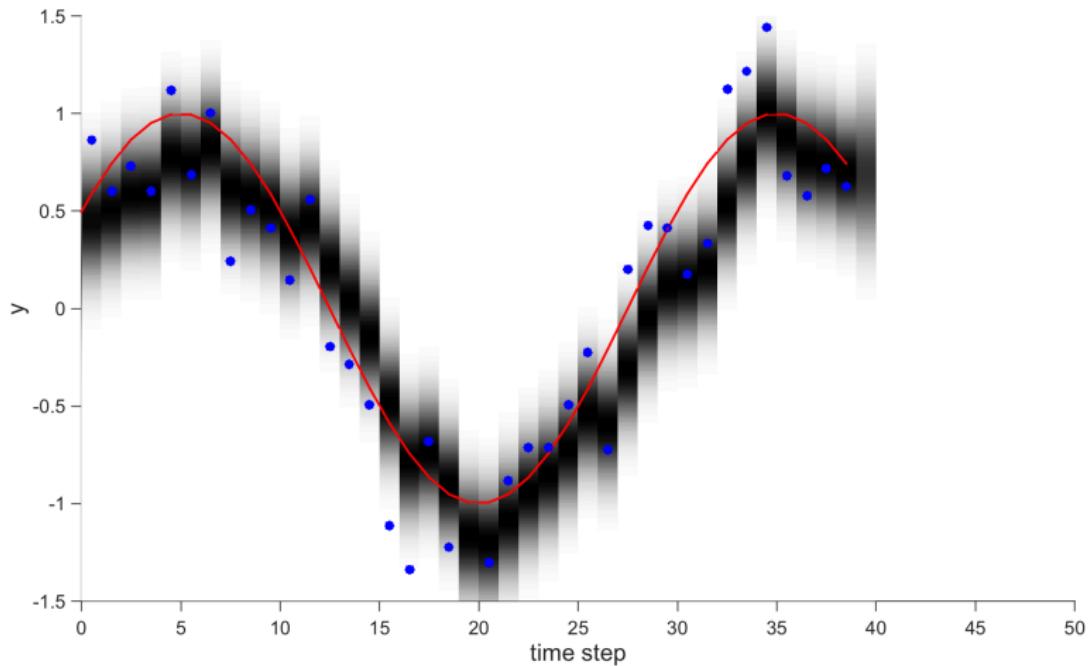
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



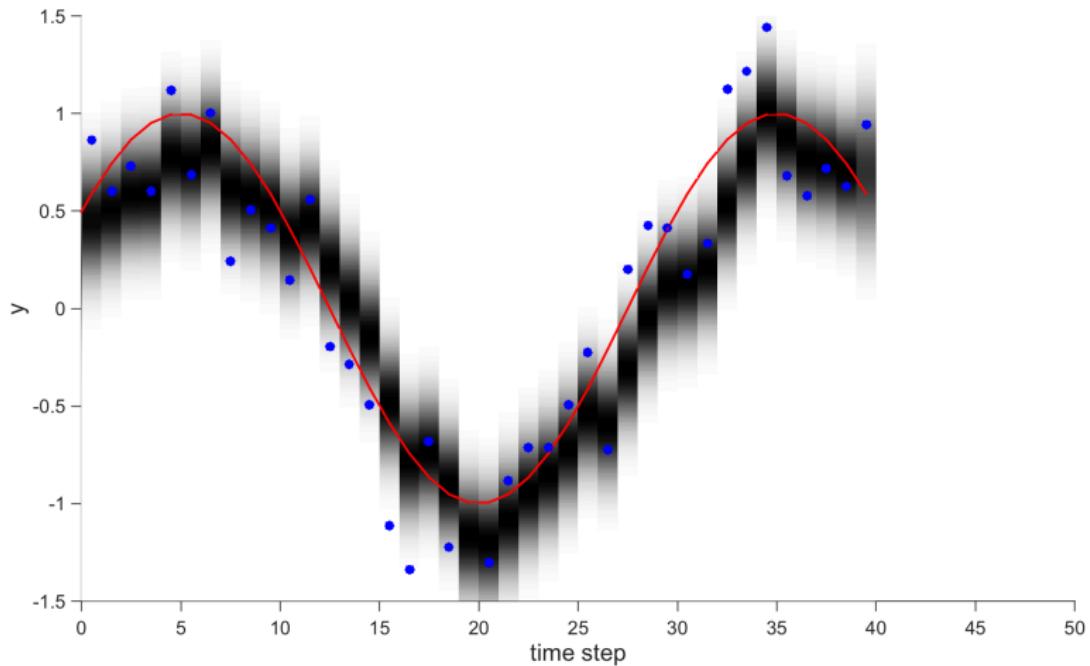
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



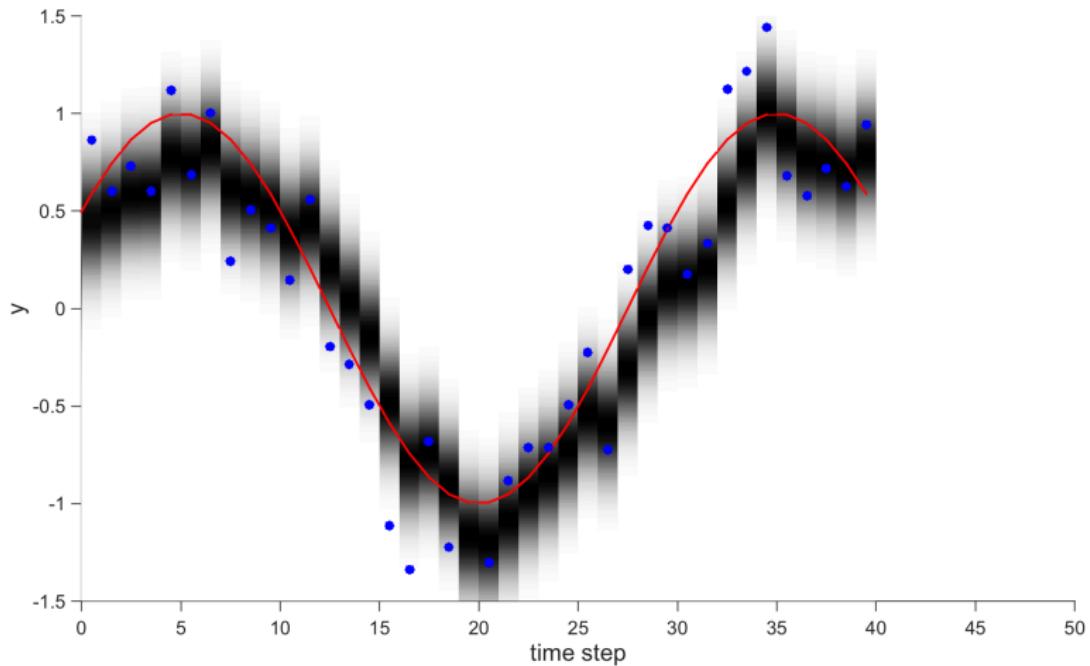
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



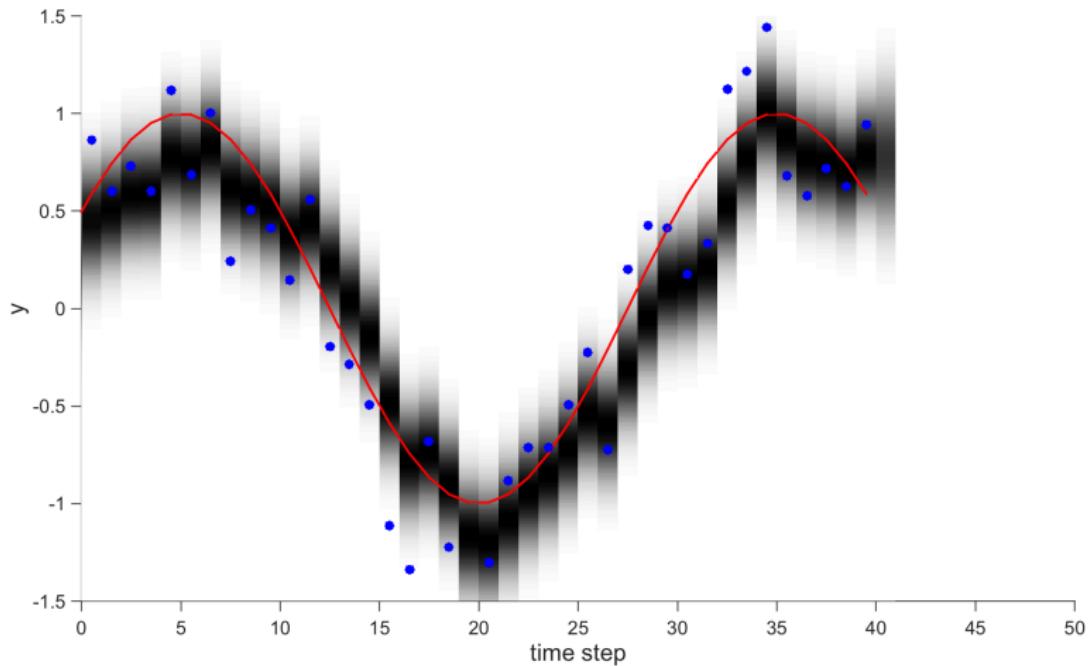
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



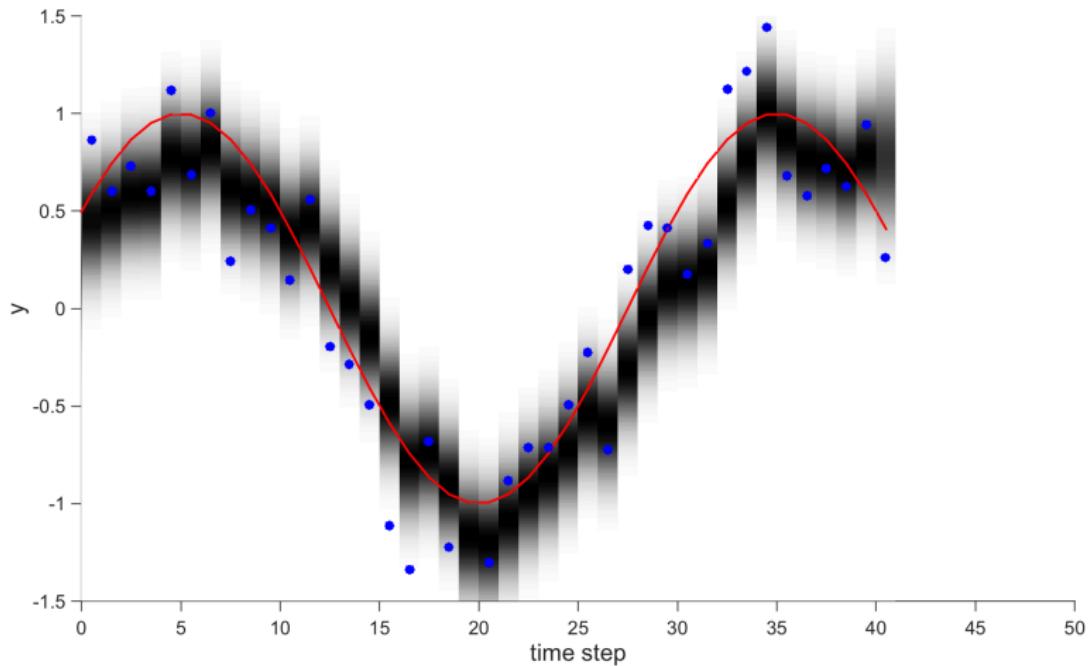
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



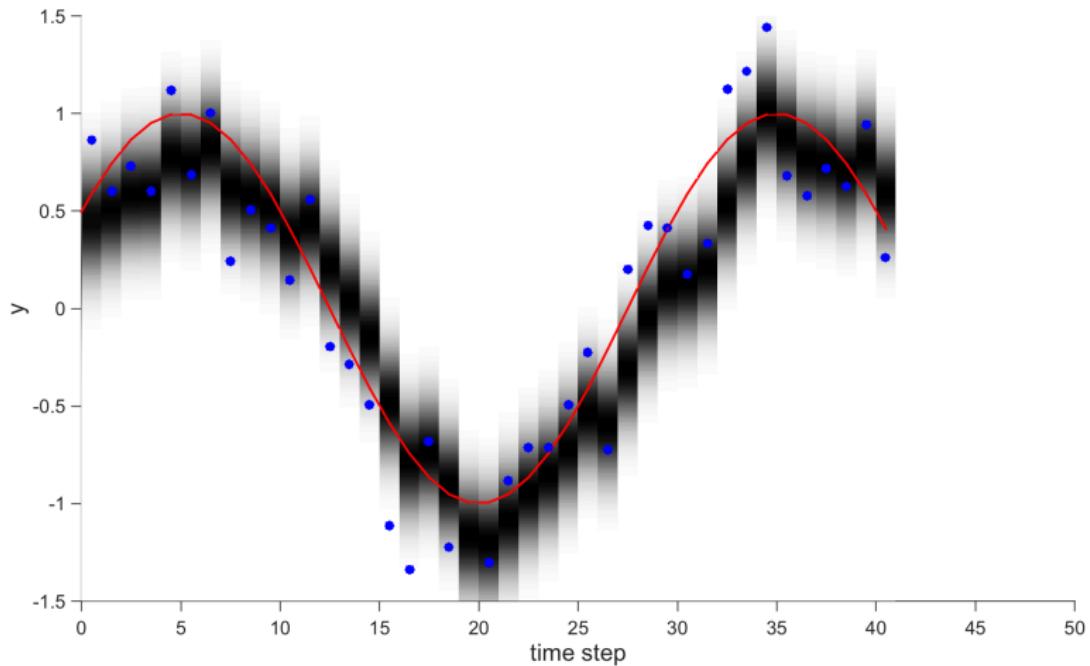
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



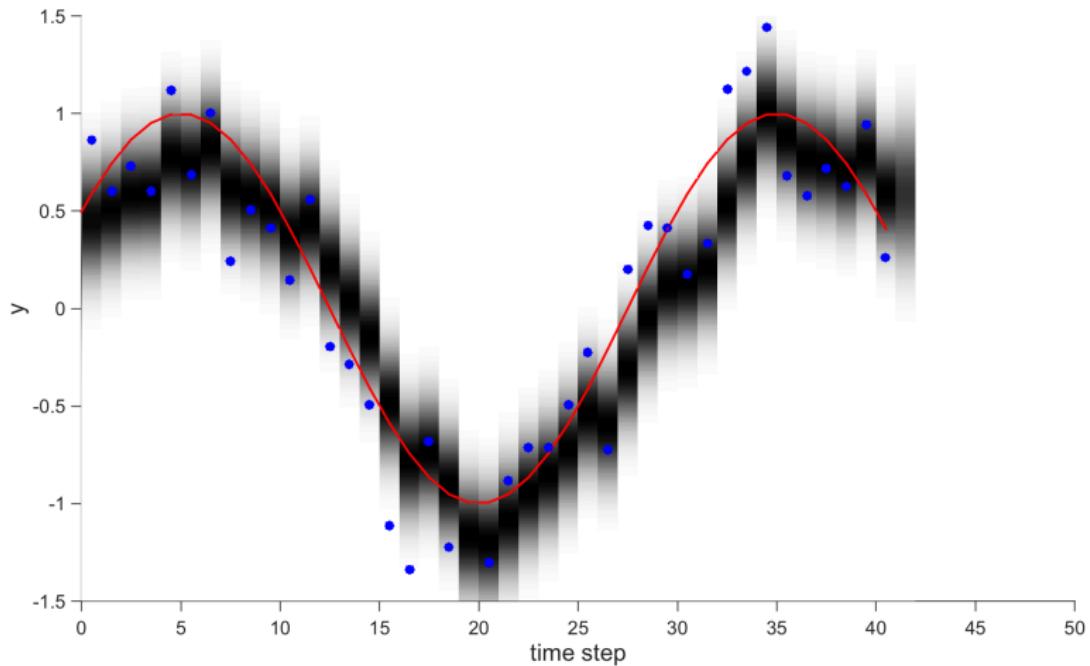
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



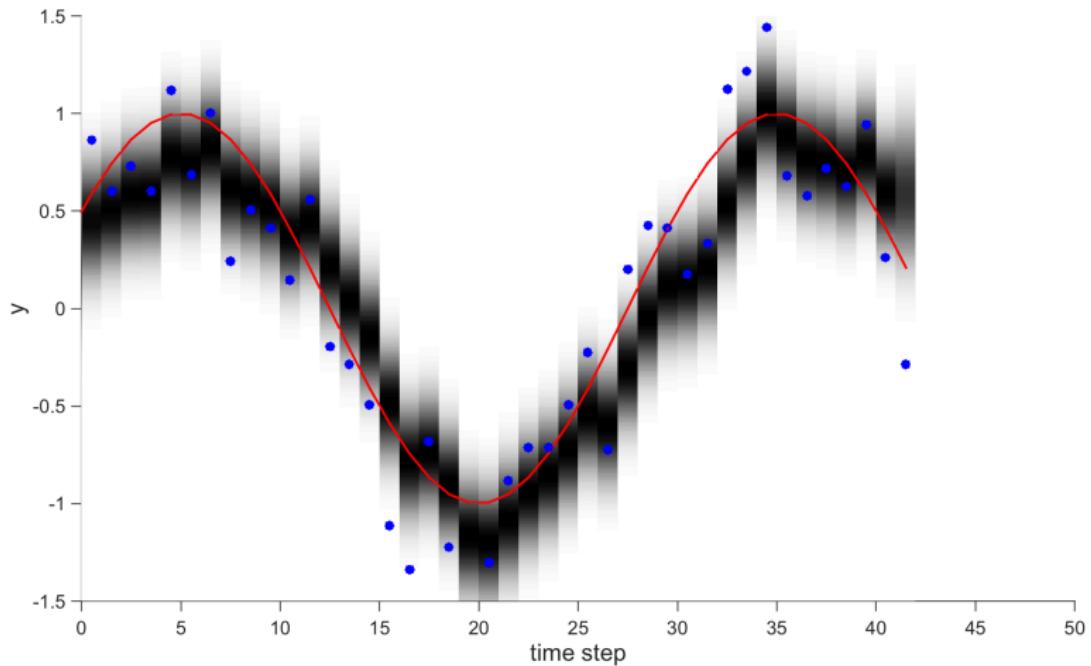
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



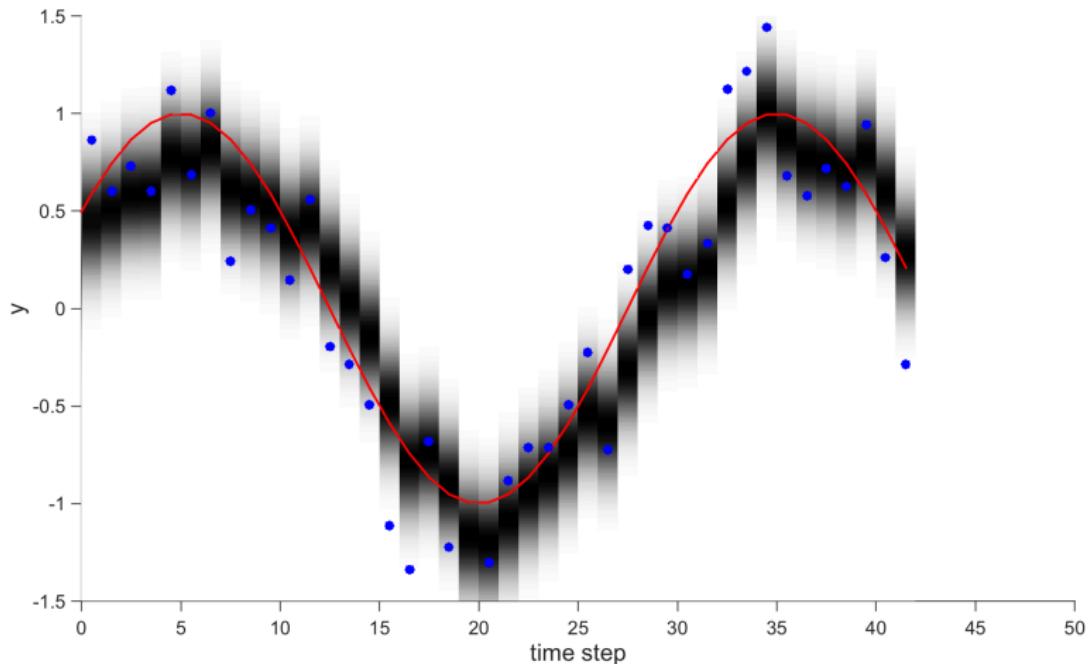
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



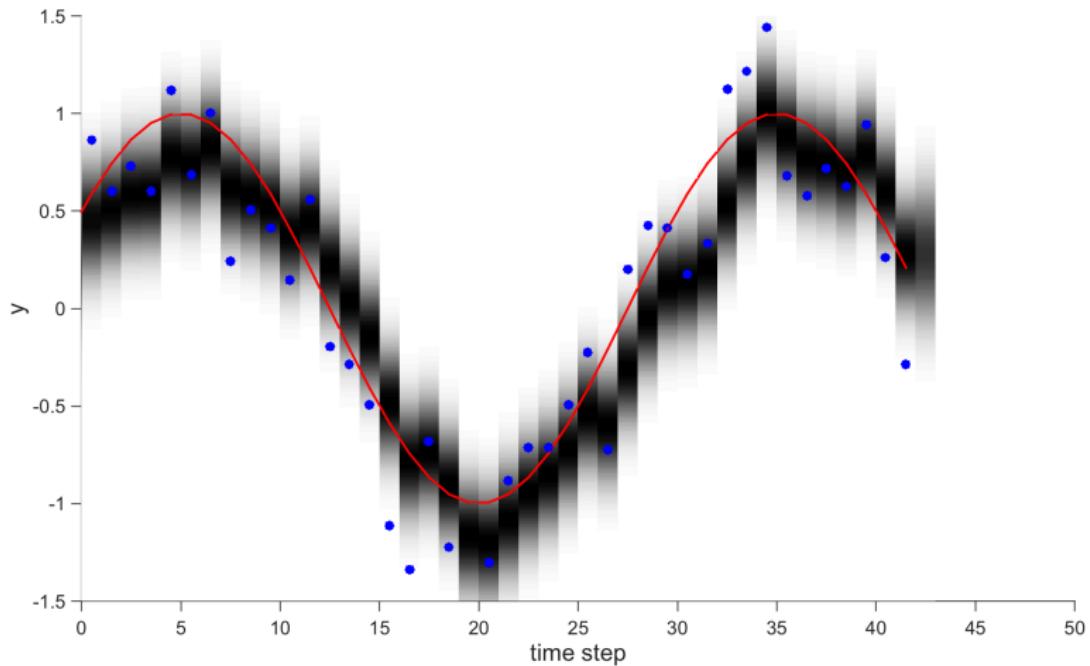
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



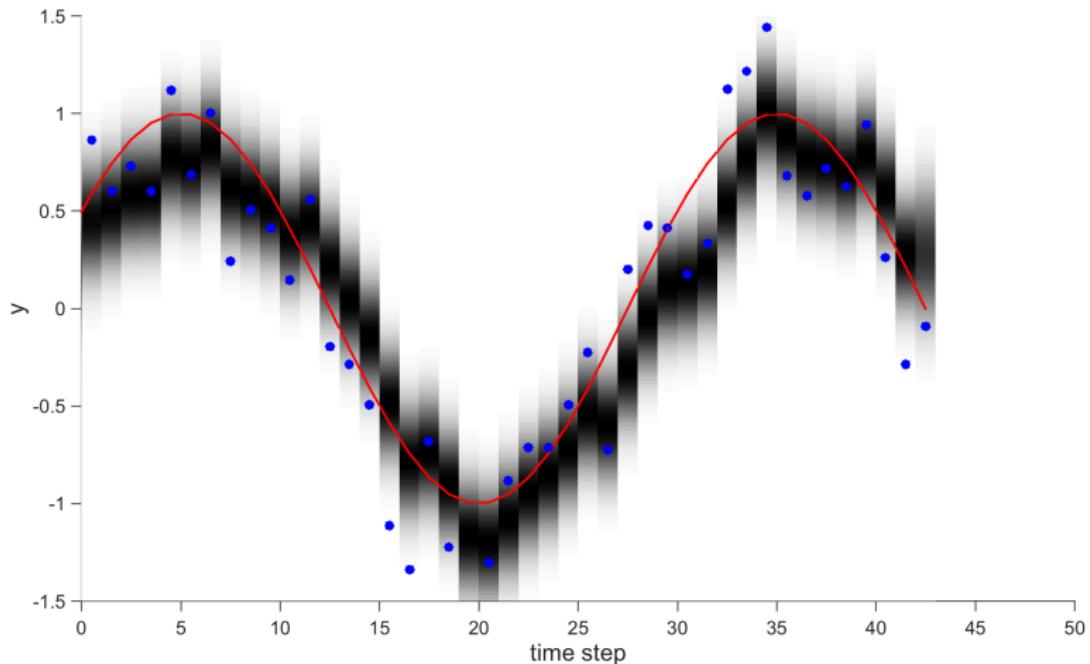
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



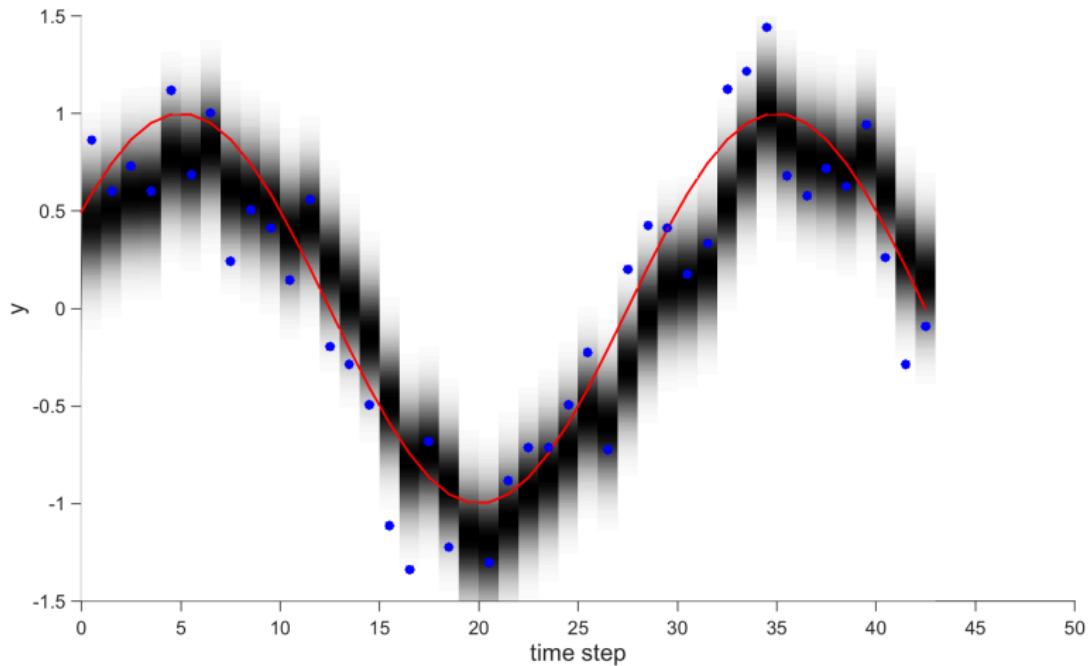
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



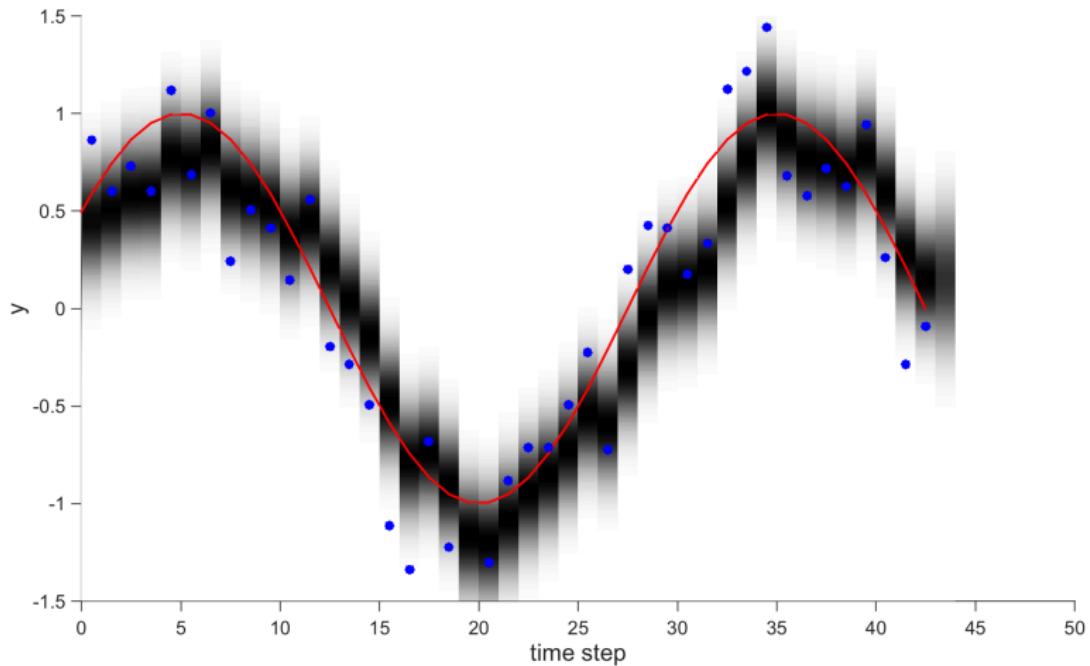
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



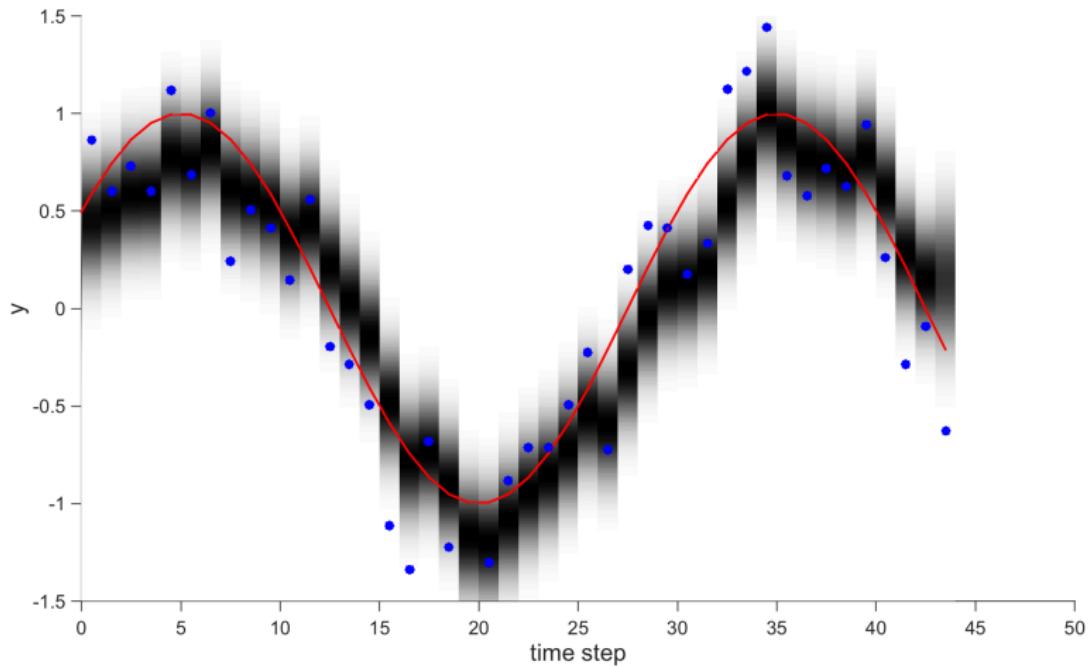
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



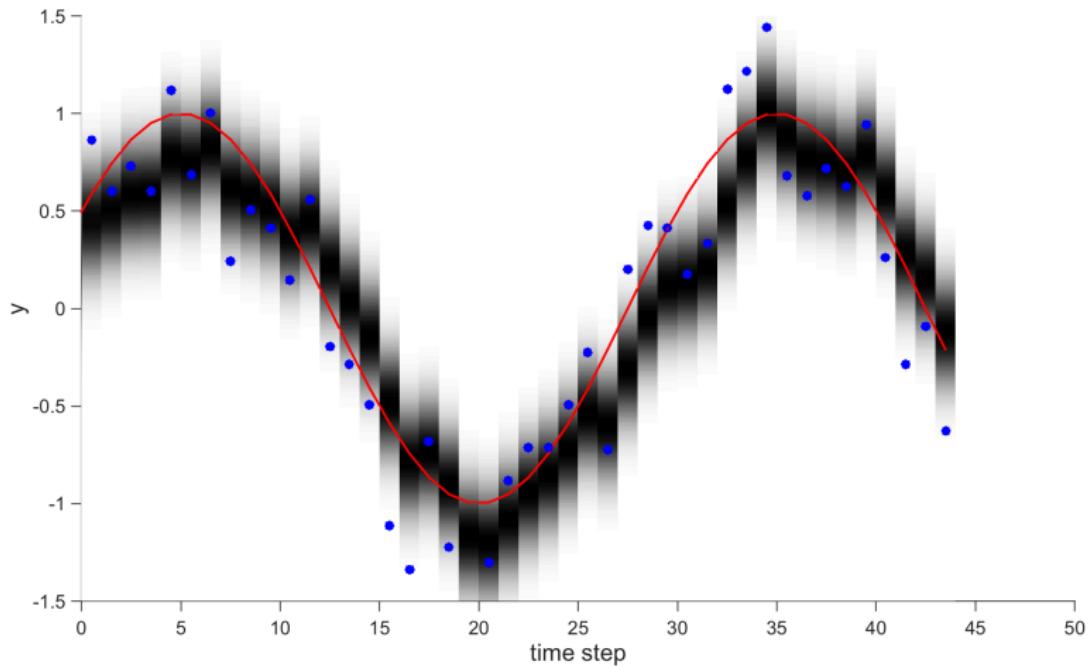
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



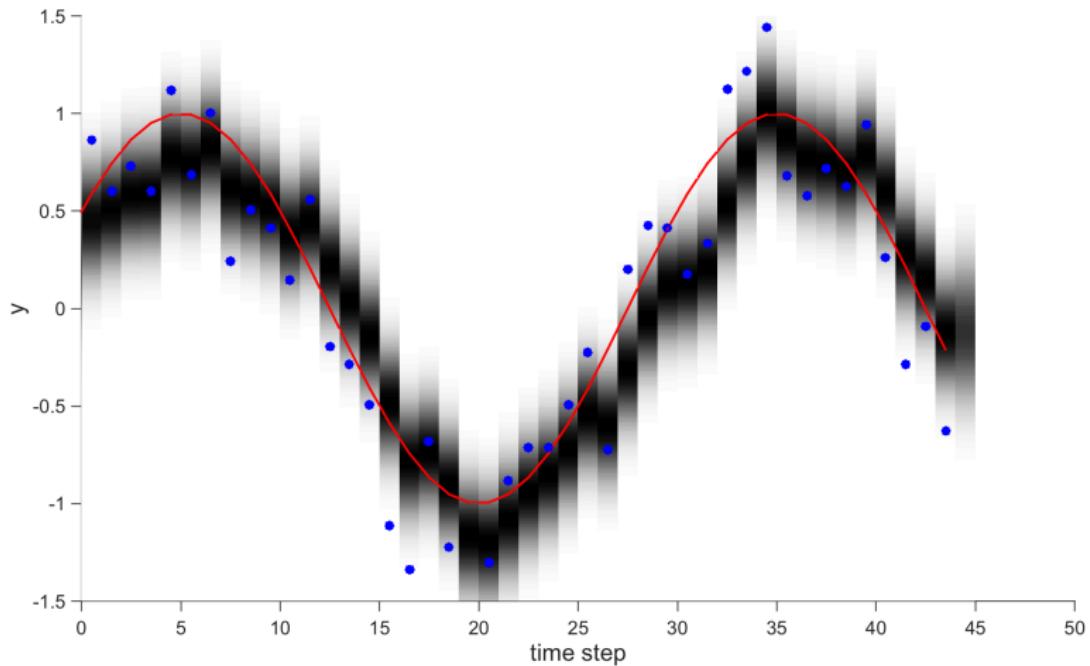
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



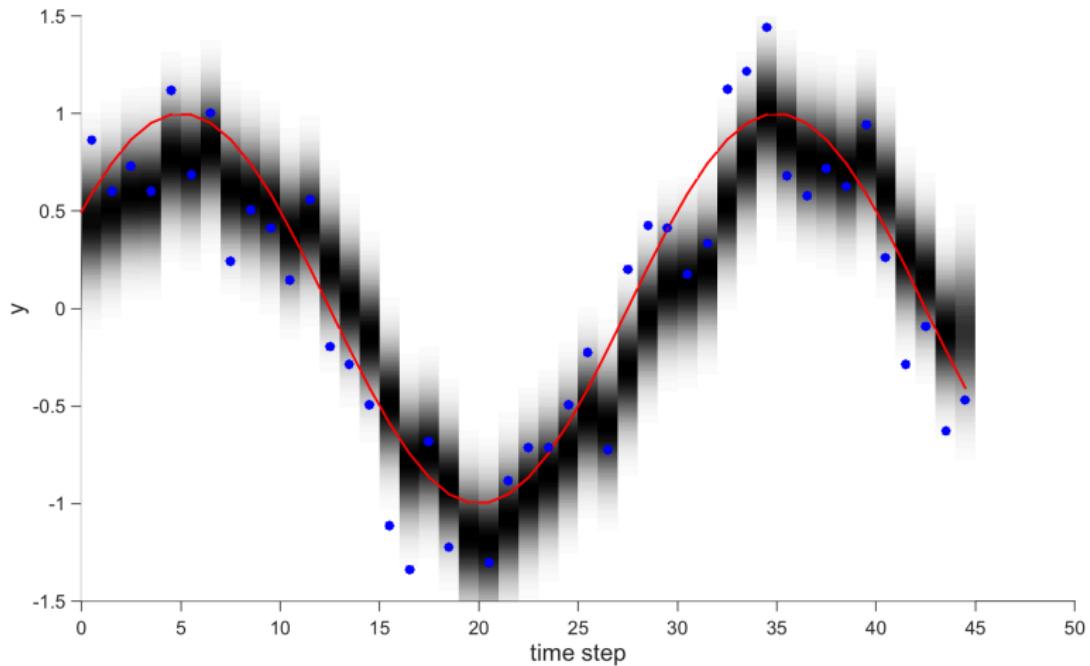
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



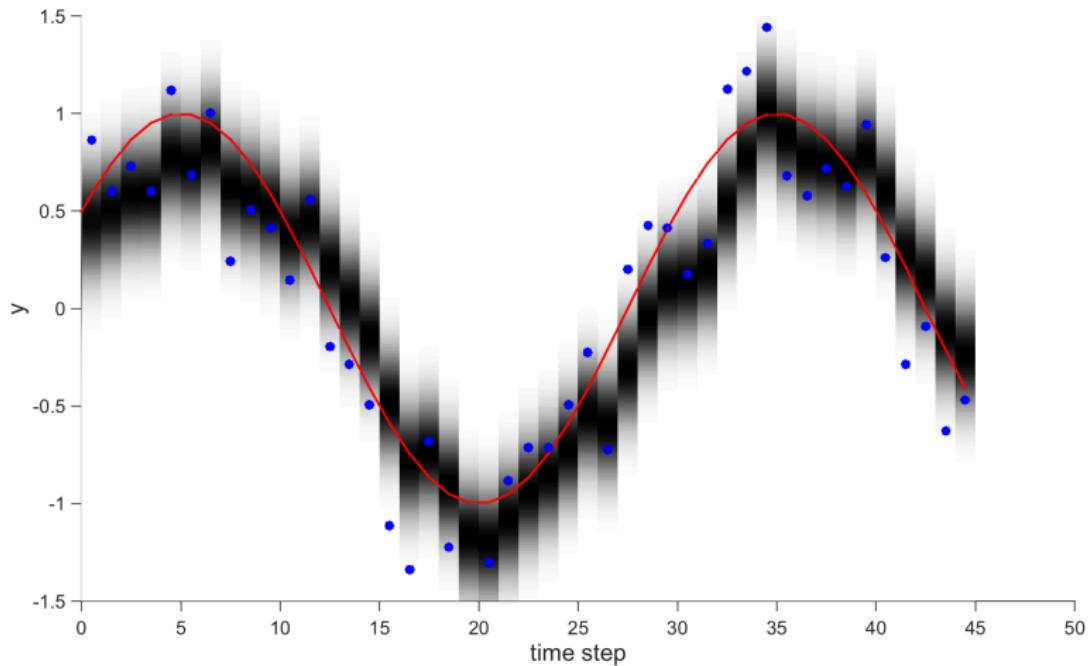
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



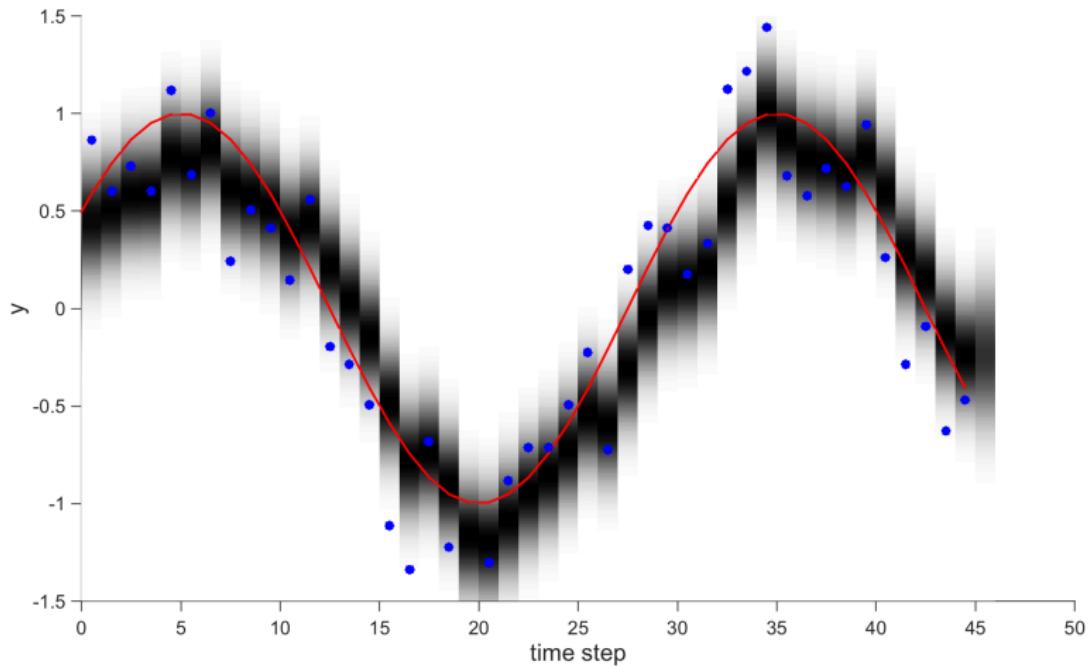
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



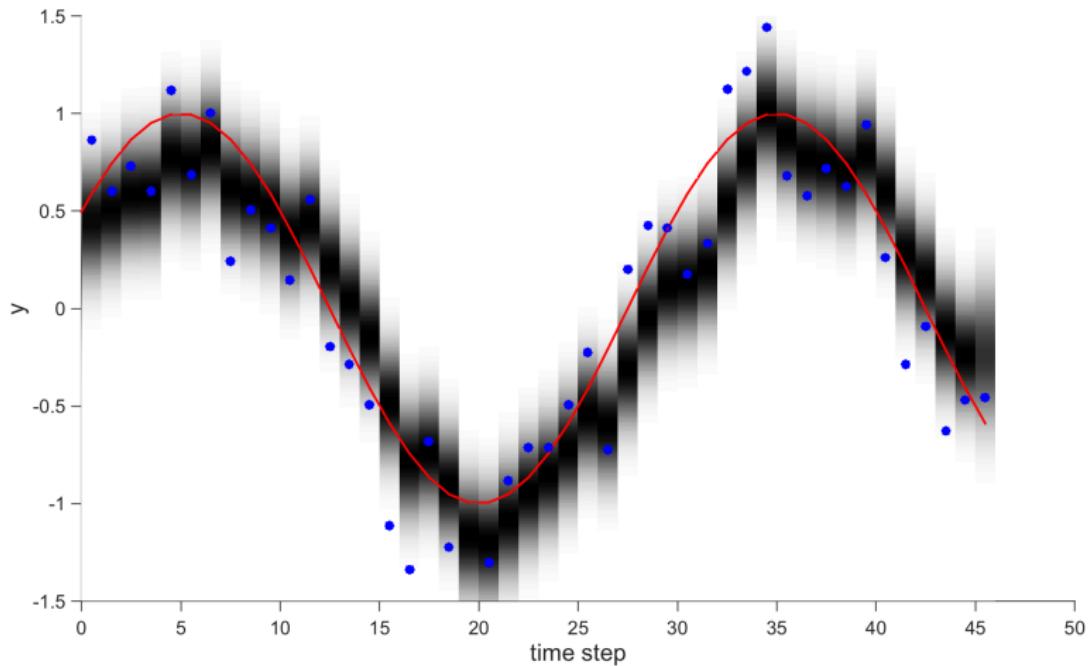
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



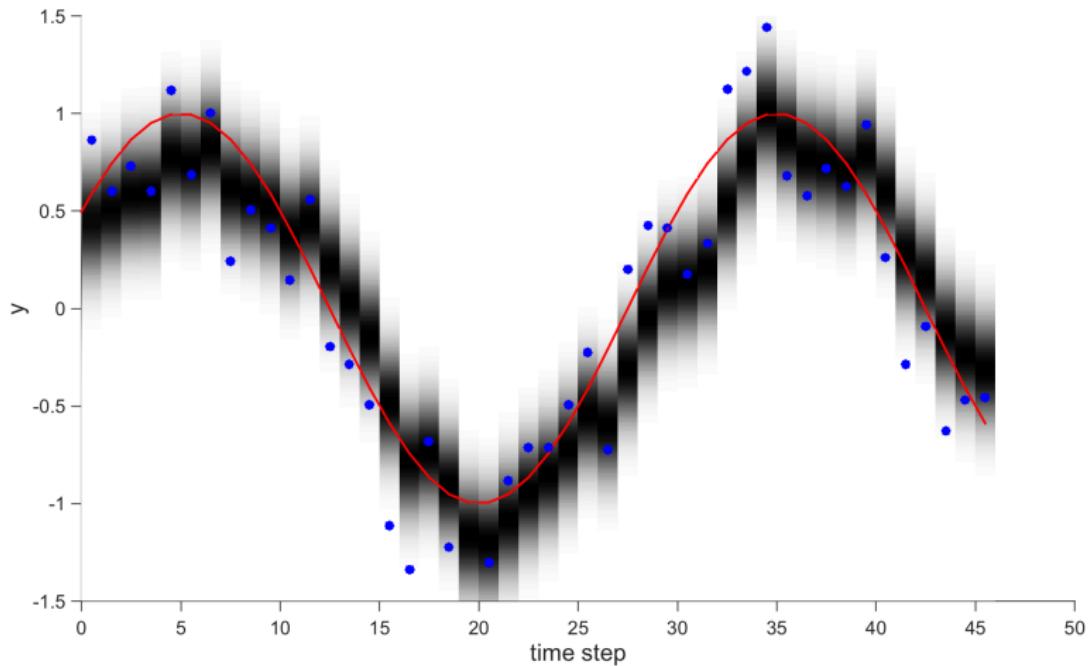
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



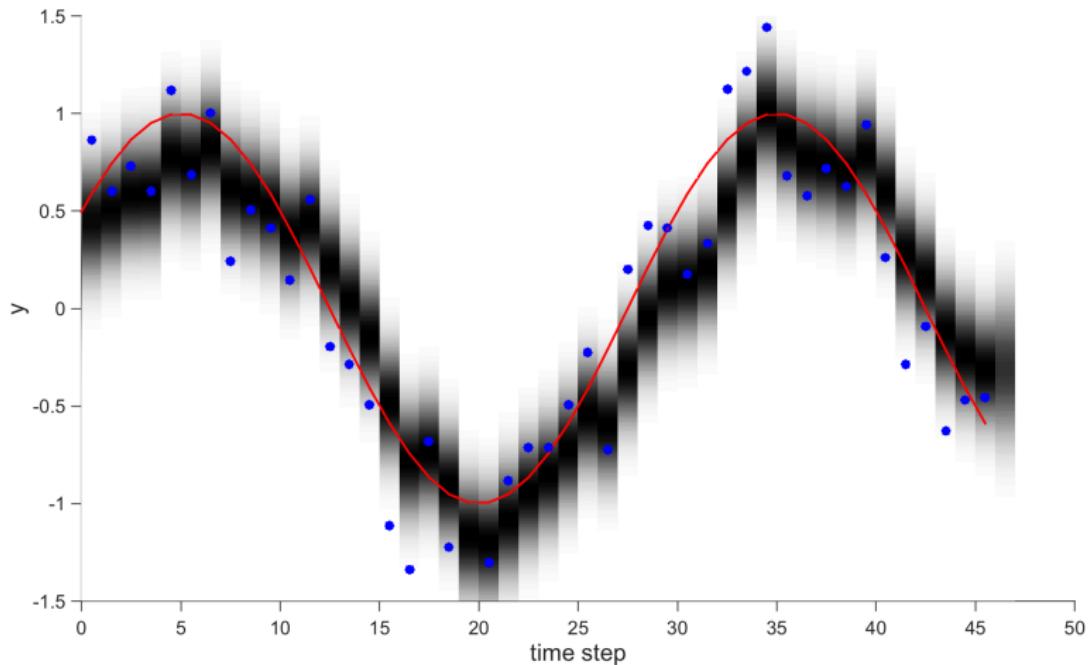
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



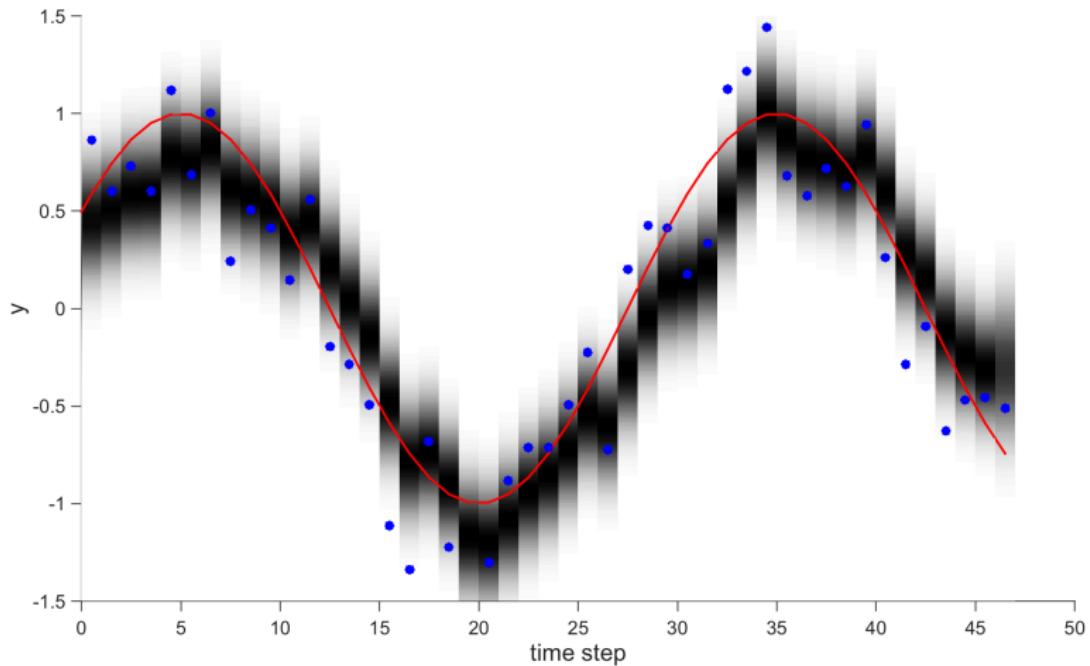
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



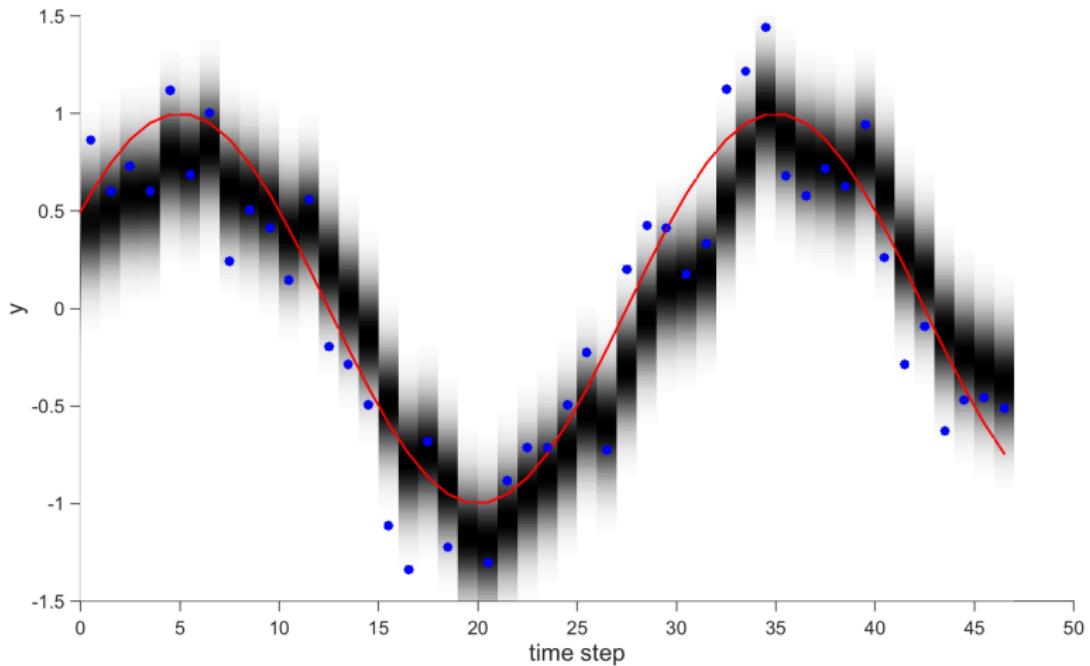
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



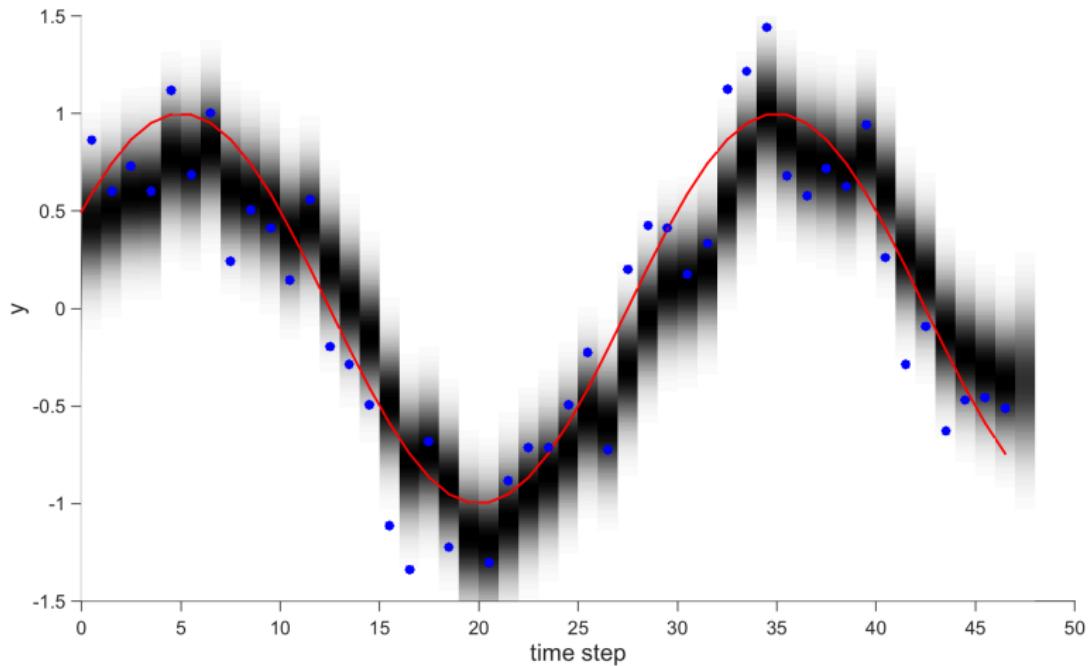
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



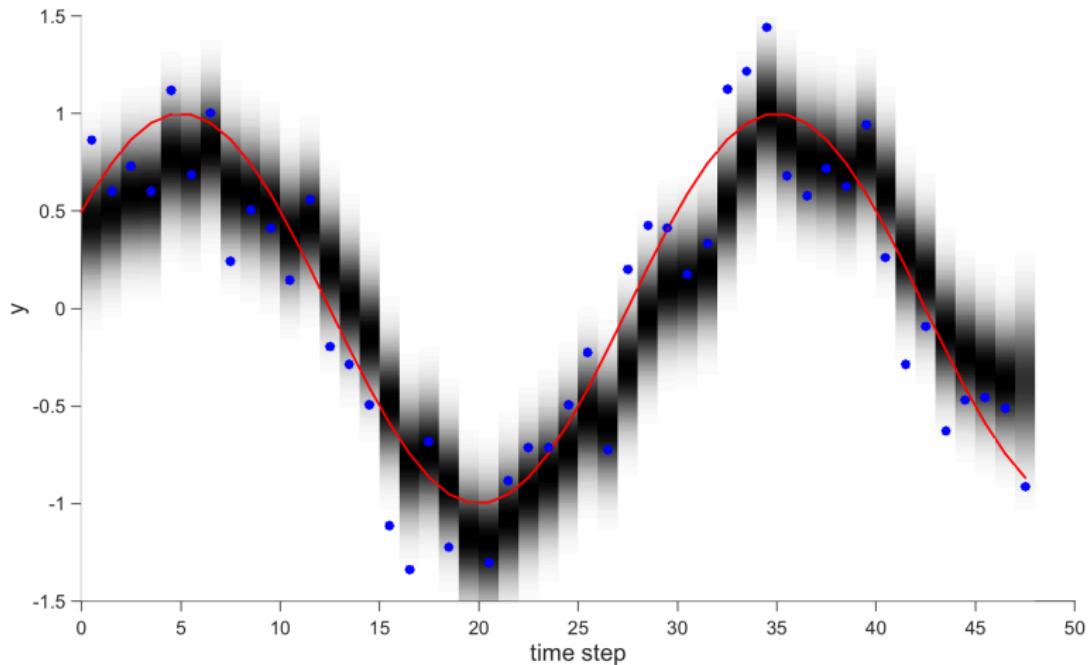
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



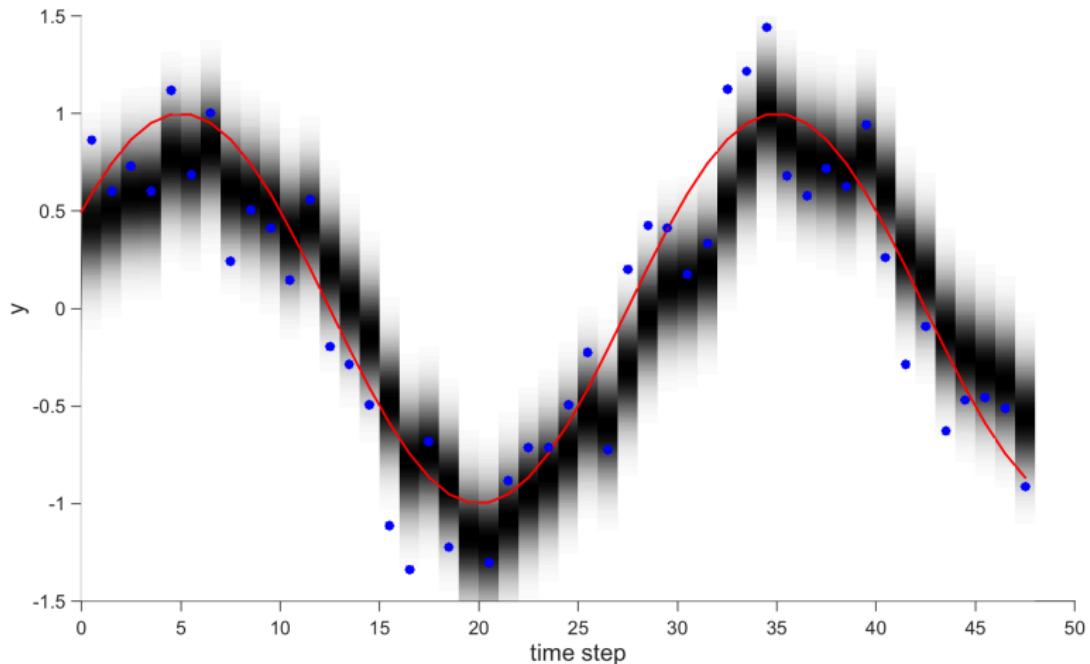
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



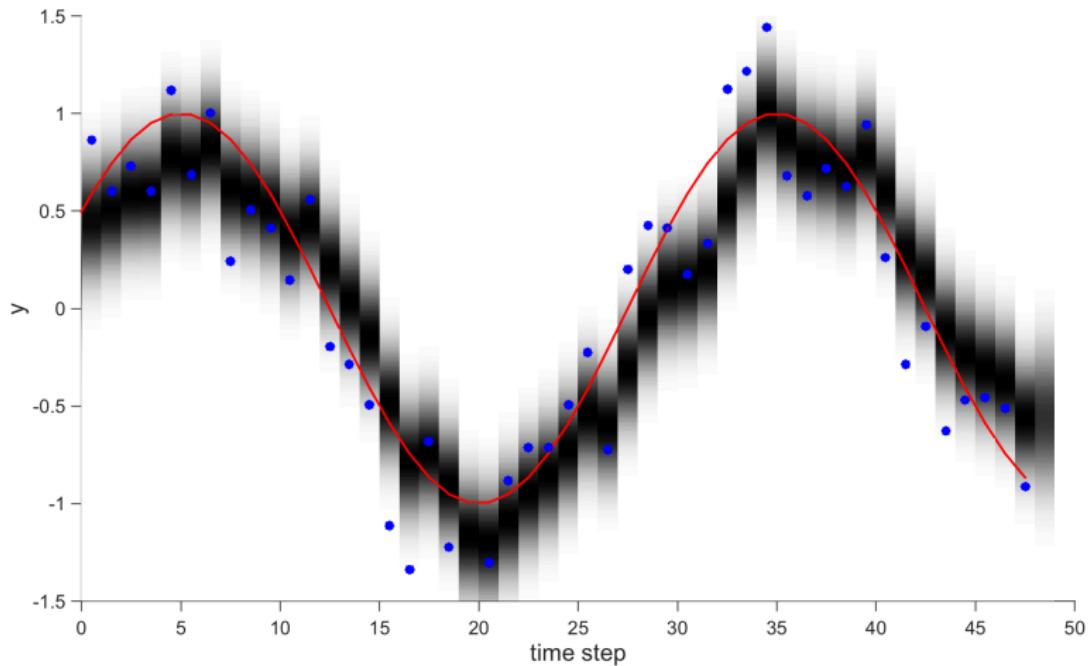
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



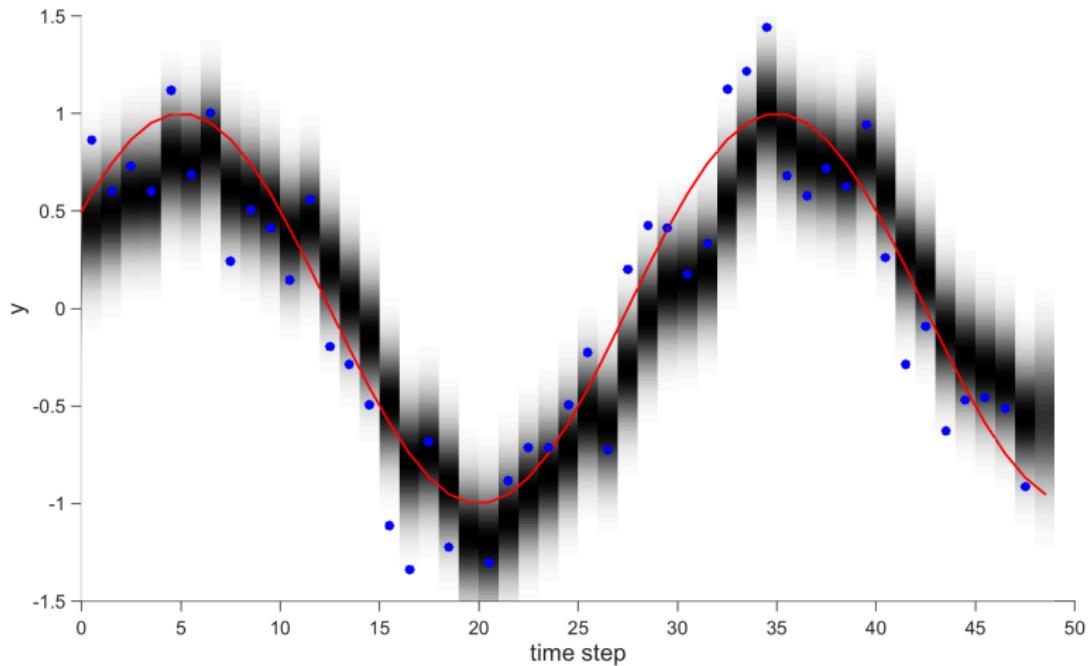
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



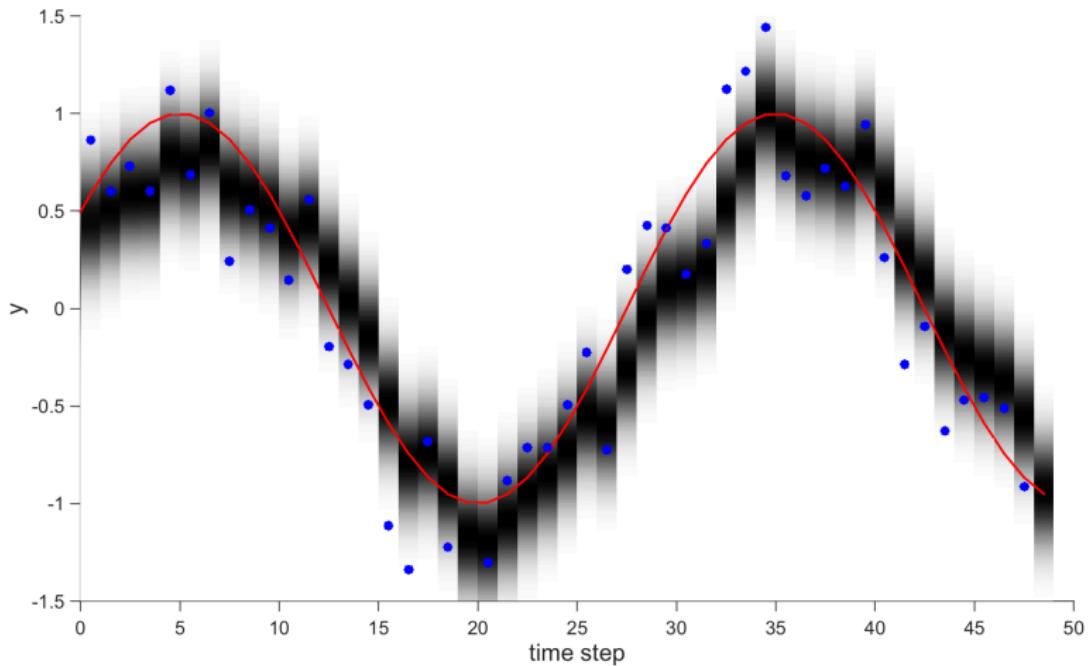
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



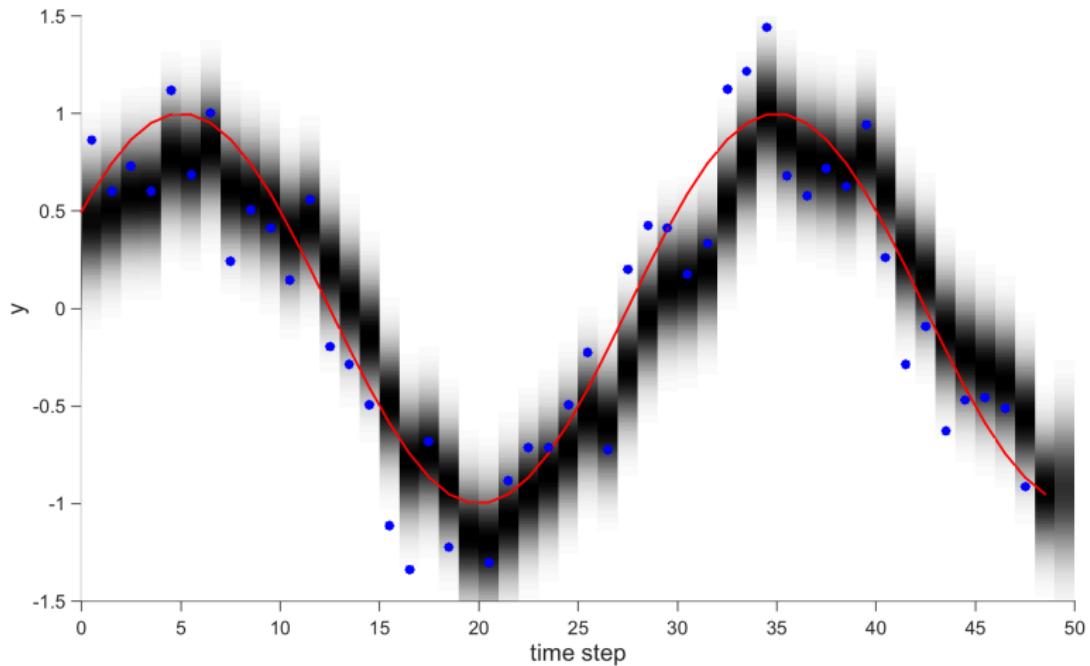
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



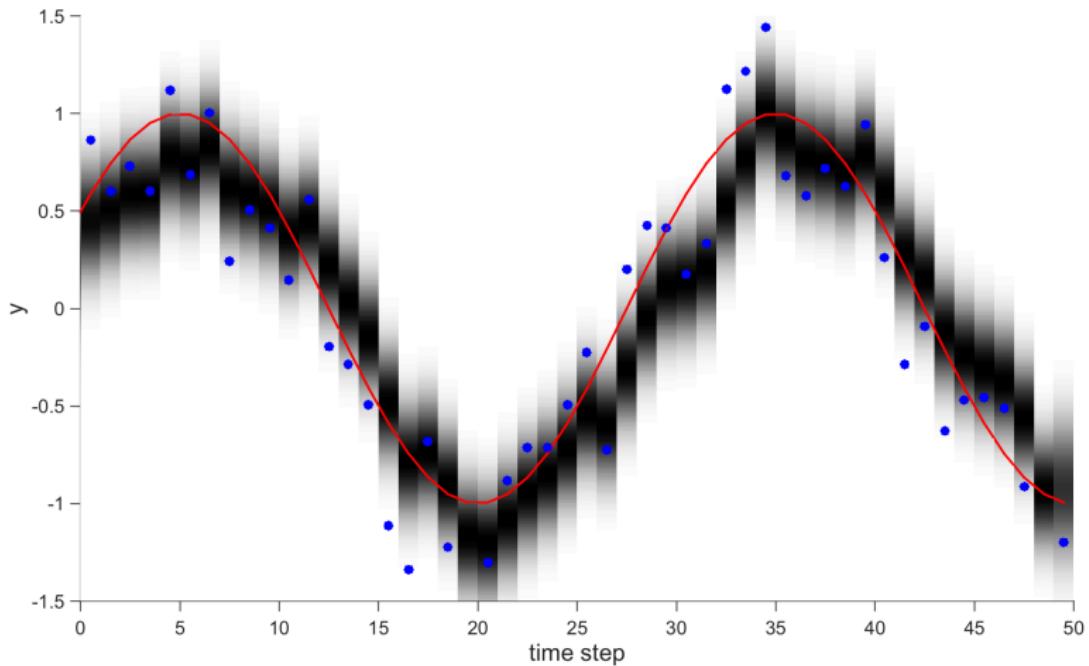
## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid



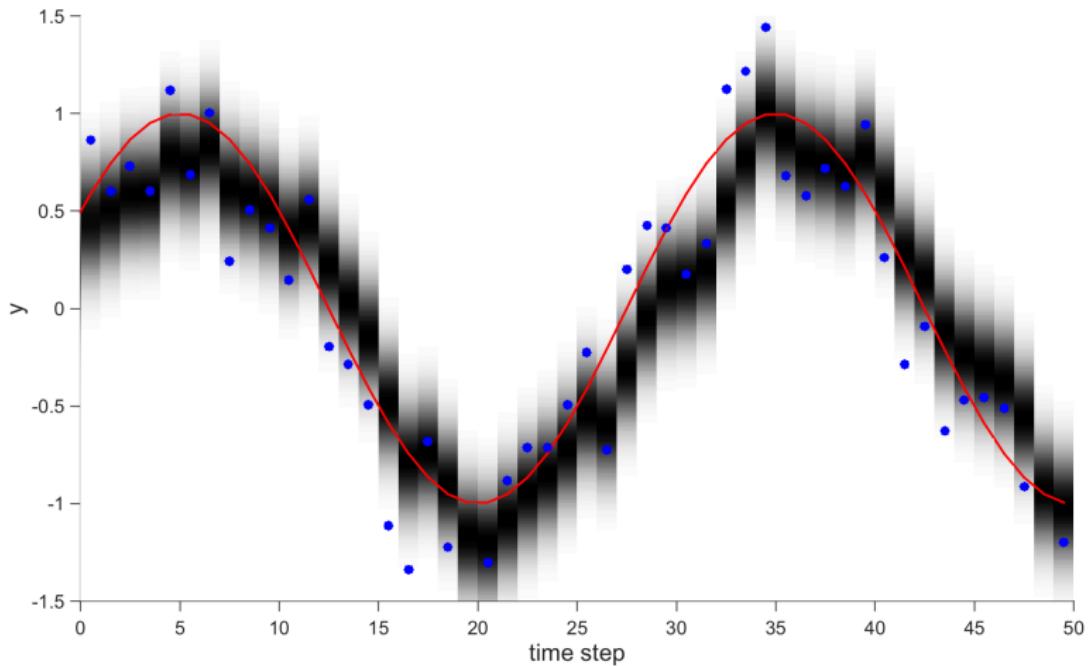
## Kalman Filter Demo

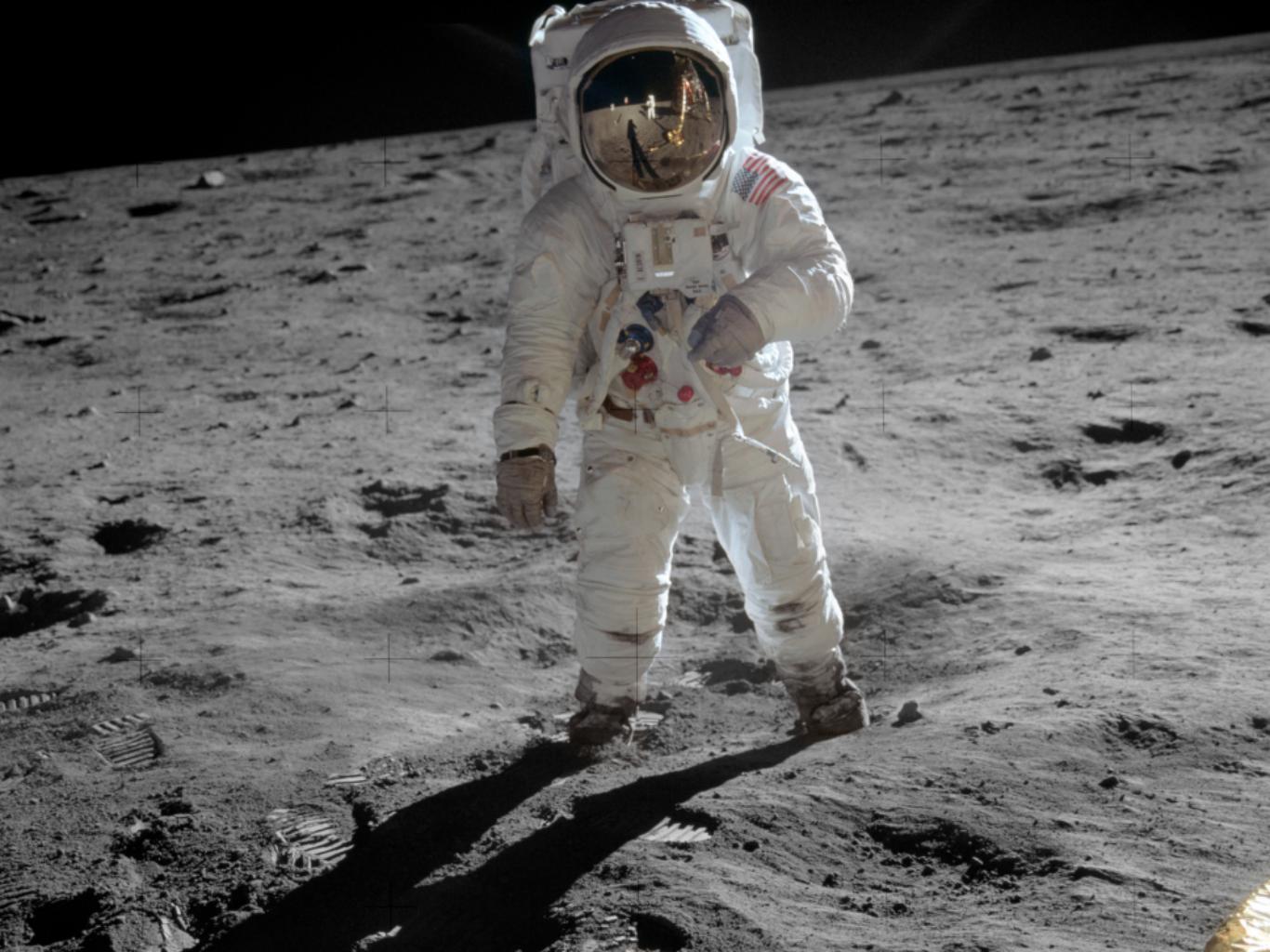
observed noisy data  $y_t$ , ground truth sinusoid



## Kalman Filter Demo

observed noisy data  $y_t$ , ground truth sinusoid





Course Survey: please complete this!



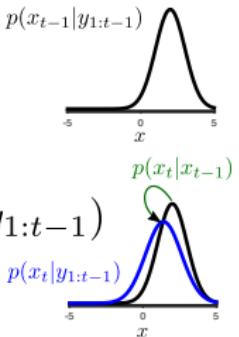
<https://tinyurl.com/3F8survey2020>

## Inference: Forward Algorithm

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

diffuse via  
dynamics

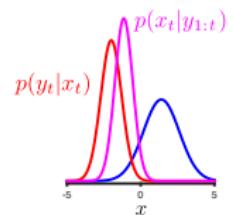
$$p(x_t = k | y_{1:t-1}) = \sum_{l=1}^K p(x_t = k | x_{t-1} = l) p(x_{t-1} = l | y_{1:t-1})$$



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

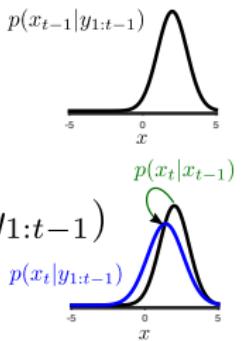
$$p(x_{t-1} = k | y_{1:t-1}) = \rho_{t-1}^{t-1}(k)$$

most recent data used  
in prediction

variable being predicted

diffuse via  
dynamics

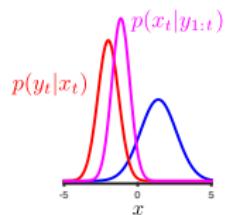
$$p(x_t = k | y_{1:t-1}) = \sum_{l=1}^K p(x_t = k | x_{t-1} = l) p(x_{t-1} = l | y_{1:t-1})$$



combine  
with  
likelihood

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

prior                      likelihood



## Inference: Forward Algorithm

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

diffuse via dynamics



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

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

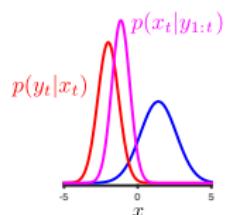
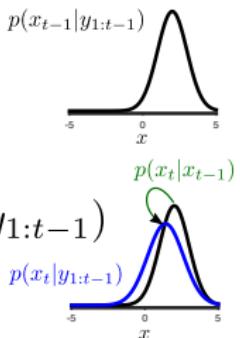
combine with likelihood



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

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

$$\downarrow$$

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

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

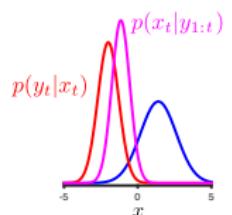
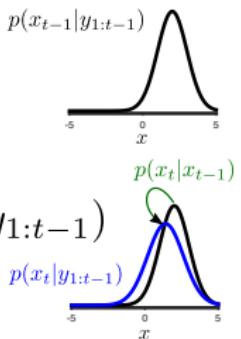
combine with likelihood

$$\downarrow$$

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

prior                      likelihood

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



## 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

$$\downarrow$$

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

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

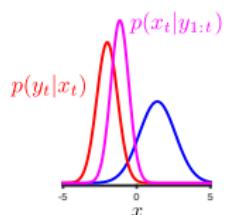
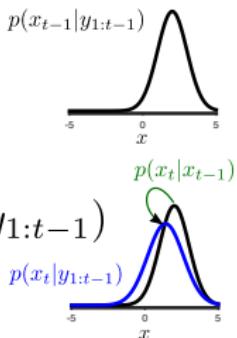
combine with likelihood

$$\downarrow$$

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

prior                      likelihood

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



When implementing, take care with numerical underflow/overflow.

## Computing the likelihood

How can we compute the likelihood efficiently?

## Computing the likelihood

How can we compute the likelihood efficiently?

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

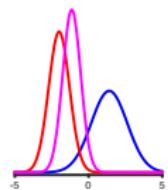
## Computing the likelihood

How can we compute the likelihood efficiently?

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

already returned by Kalman Filter/Forward algorithm

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



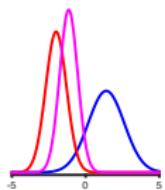
## Computing the likelihood

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

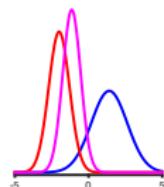
## Computing the likelihood

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

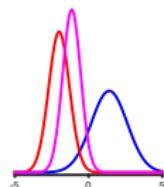
## Computing the likelihood

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:  $x_{1:T}^{(k)}$

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

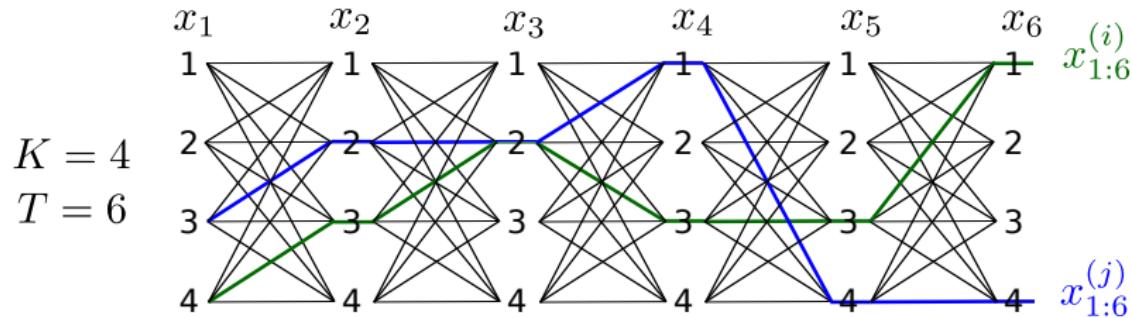
## The magic of the Forward Algorithm: Dynamic Programming

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:



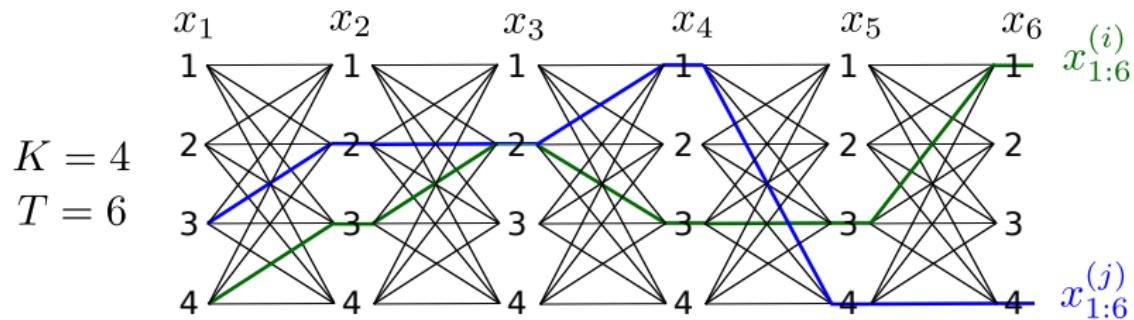
## The magic of the Forward Algorithm: Dynamic Programming

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:  $K^T$

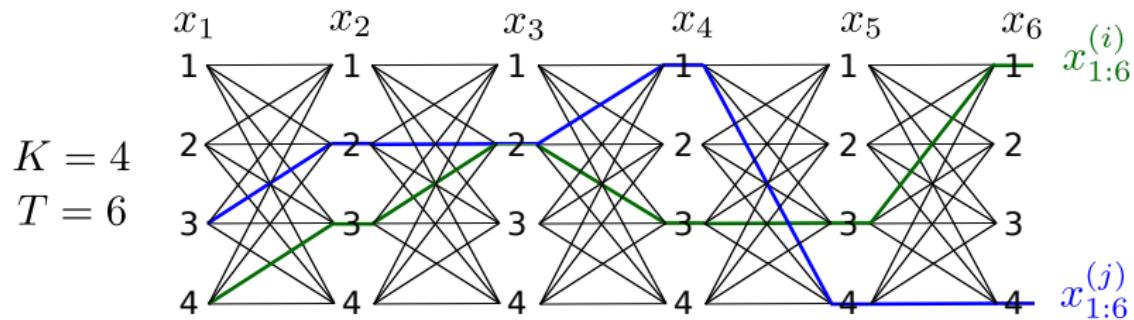
## The magic of the Forward Algorithm: Dynamic Programming

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)

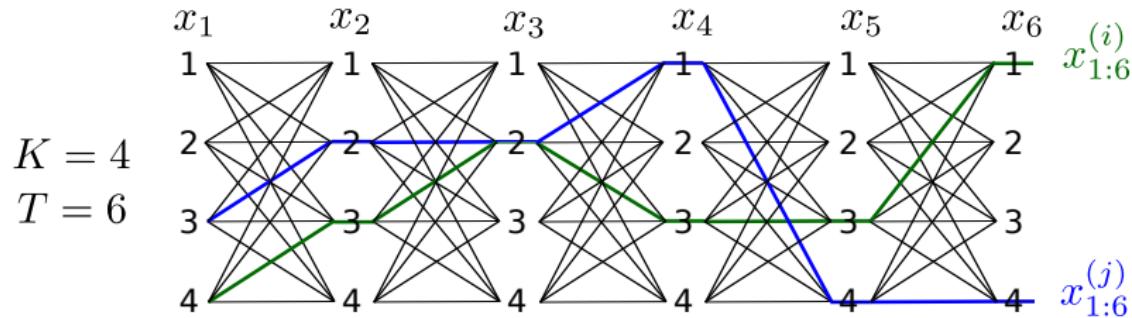
## The magic of the Forward Algorithm: Dynamic Programming

What's going on here?

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

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

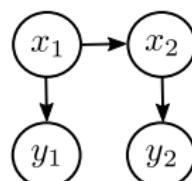
Trellis diagram represents possible sequences:



Exponential number of sequences:  $K^T$

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

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



## Maximum Likelihood Learning of HMMs: simple once inference is solved

log-likelihood:

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

## Maximum Likelihood Learning of HMMs: simple once inference is solved

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

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

## Maximum Likelihood Learning of HMMs: simple once inference is solved

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

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

show gradient depends  
on simple moments  
of posterior:

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

## Maximum Likelihood Learning of HMMs: simple once inference is solved

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

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

show gradient depends  
on simple moments  
of posterior:

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

## Maximum Likelihood Learning of HMMs: simple once inference is solved

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

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

show gradient depends  
on simple moments  
of posterior:

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

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

## 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:

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

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

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

## 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:

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

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

## 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:

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

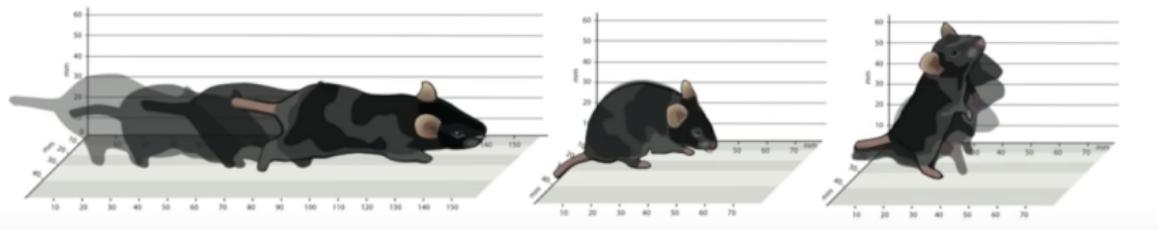
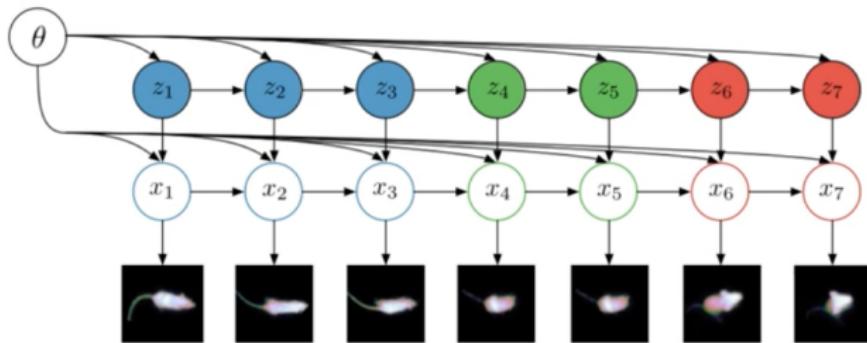
## Summary of Sequence Modelling

- ▶ **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>