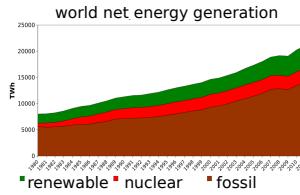
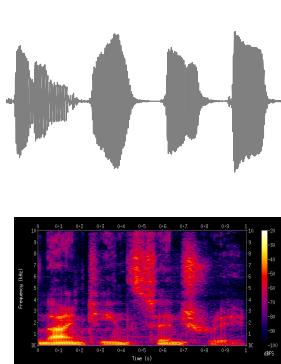


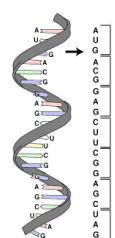
# Sequence Modelling

Rich Turner and José Miguel Hernández-Lobato

## Sequence data



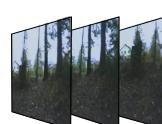
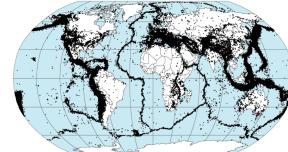
Some images taken from wikipedia



Ribonucleic acid

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

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

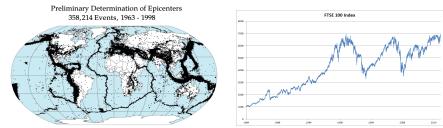


I believe that at the end of the century the use of words and general educated opinion will have altered so much that one will be able to speak of machines thinking without expecting to be contradicted.  
A. Turing

## Goals of sequence modelling

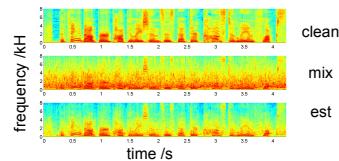
Predict future items in sequence

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



Remove noise from a sequence

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



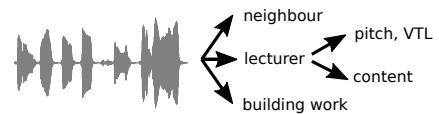
Predict one sequence from another

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



Discover underlying latent variables

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

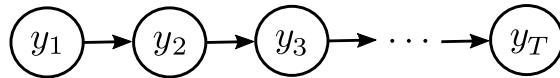


## Markov models

First order Markov

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

parameters tied  $\infty$  number of variables  
finite number of parameters

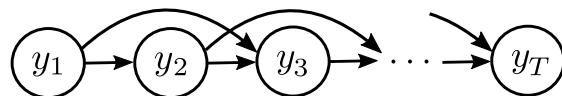


Markov model = conditional independence relationship + product rule

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

Second order Markov

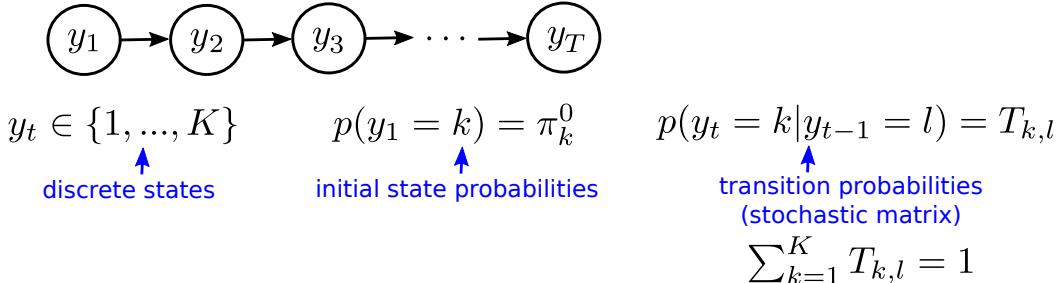
$$p(y_1, y_2, y_3, \dots, y_T) = p(y_1)p(y_2|y_1)p(y_3|y_2, y_1) \dots p(y_T|y_{T-1}, y_{T-2})$$



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

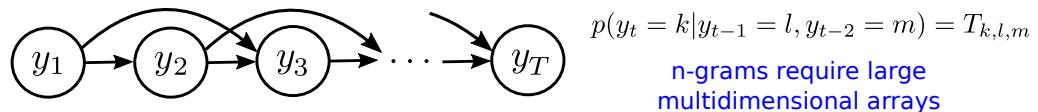
First order Markov (bi-gram)

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

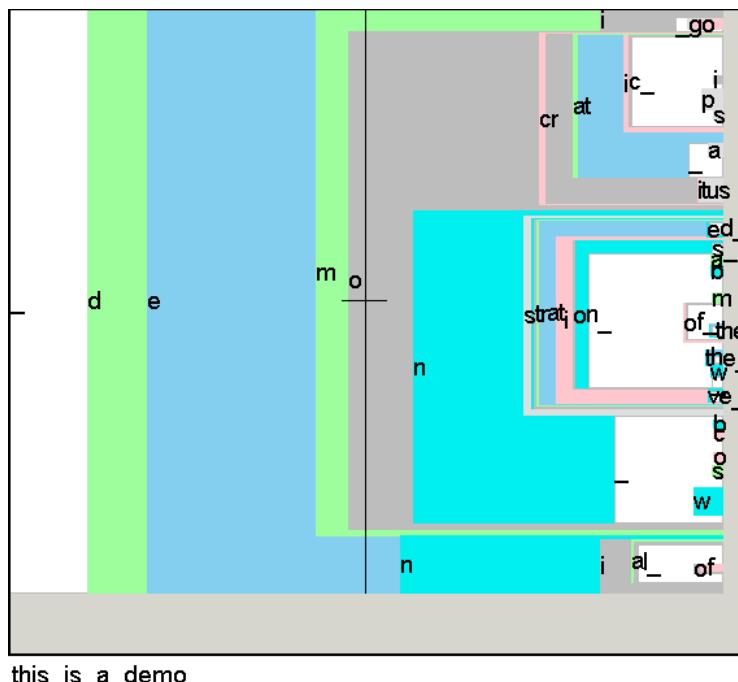


Second order Markov (tri-gram)

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



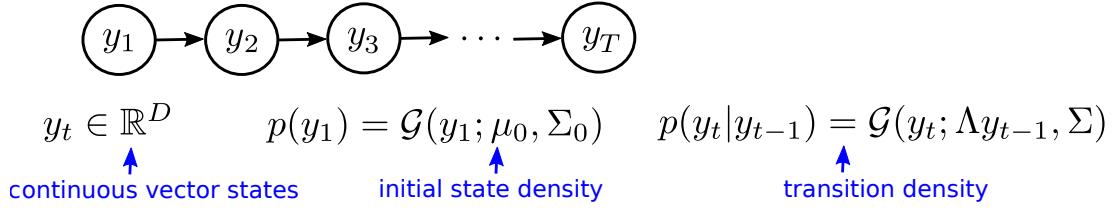
Example application of n-grams: text modelling for dasher



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

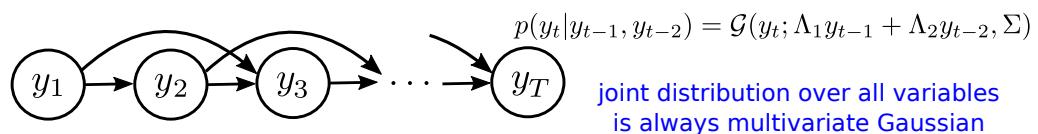
First order Markov (AR(1))

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



Second order Markov (AR(2))

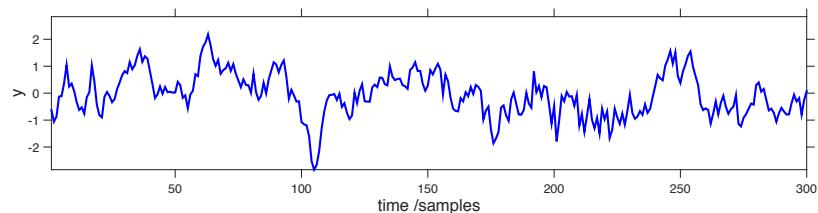
$$p(y_1, y_2, y_3, \dots, y_T) = p(y_1)p(y_2|y_1)p(y_3|y_2, y_1)\dots p(y_T|y_{T-1}, y_{T-2})$$



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

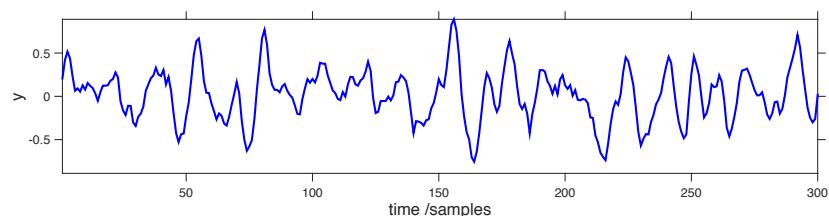
First order Markov (AR(1))

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

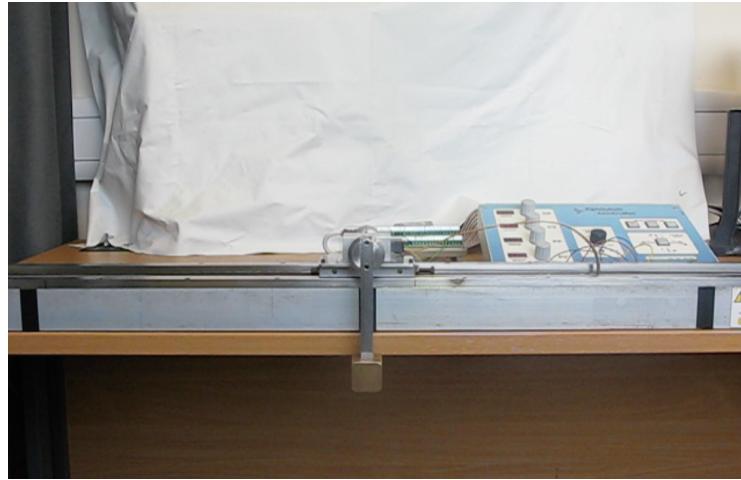


Second order Markov (AR(2))

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



## 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

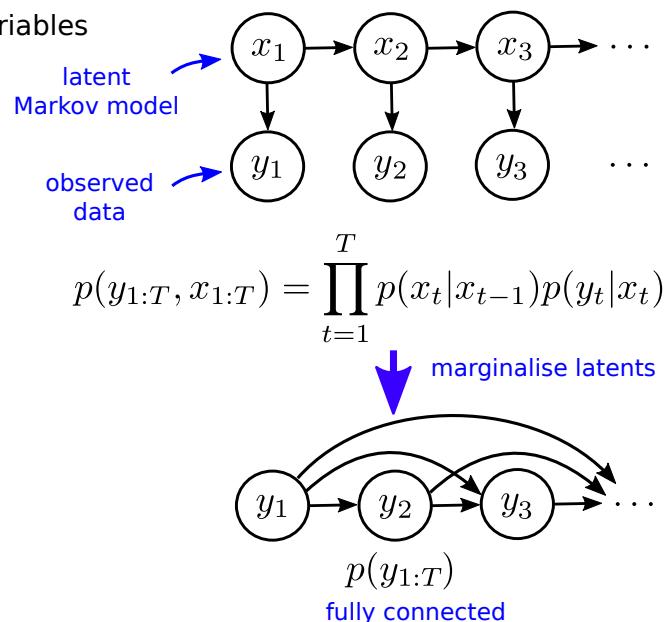
$x$  topics

$y$  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

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

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

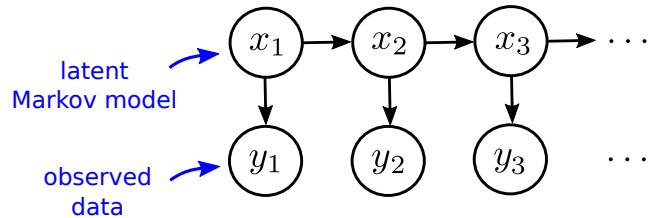
E.g. in examples below  $K = 2$

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

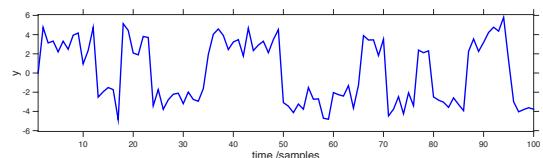
Continuous Observed State

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

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



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



Discrete Observed State

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

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

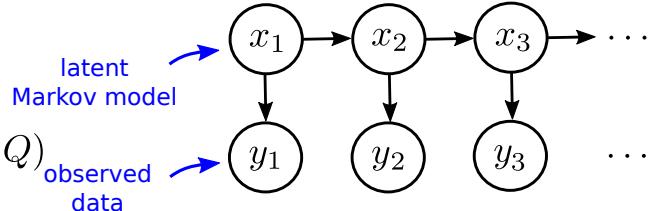
AAABBBAAABAAACCCCCBBBBBCCCCCCCCCC  
AAABBBBAABAAABBCCCCCCCCCCCCCCCCCBBA  
AACCCTCCBABCCCCCAABBAABABCCCCC

## Hidden Markov models: continuous hidden state (LGSSMs)

Continuous Hidden State

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

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



Continuous Observed State

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

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

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

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

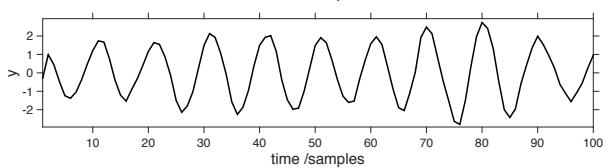
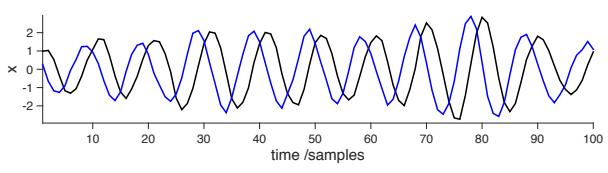
$$A = \lambda \begin{bmatrix} \cos(\theta) & \sin(\theta) \\ -\sin(\theta) & \cos(\theta) \end{bmatrix}$$

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

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

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

model  
dynamics  
model  
obs.



## Varieties of Inference

### Distributional estimates

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

Diagram illustrating the types of distributional estimates:

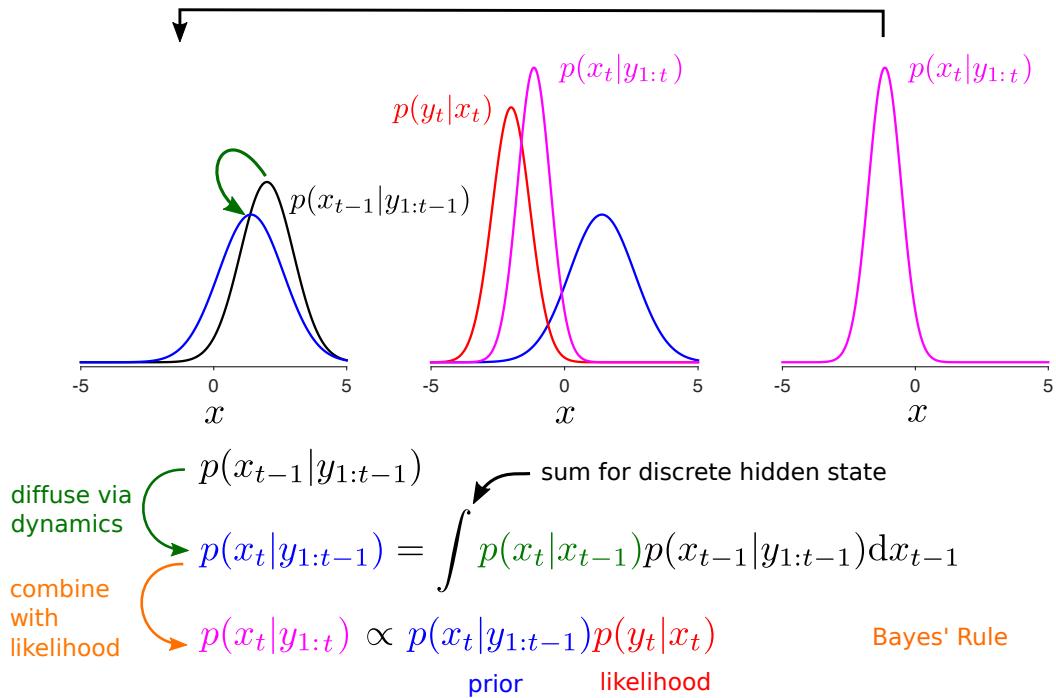
- filter:** Marginal distribution at time  $t$  given data up to  $t$ . Shows a wavy signal  $y$  from  $t$  to  $T$ .
- smoother:** Marginal distribution at time  $t$  given all data up to  $T$ . Shows a wavy signal  $y$  from  $T$  back to  $t$ .
- joint:** Joint distribution of the entire sequence  $x_{1:T}$  given data up to  $t$ .
- marginal:** Marginal distribution of a single state  $x_t$  given data up to  $t$ .

### Point estimates

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

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

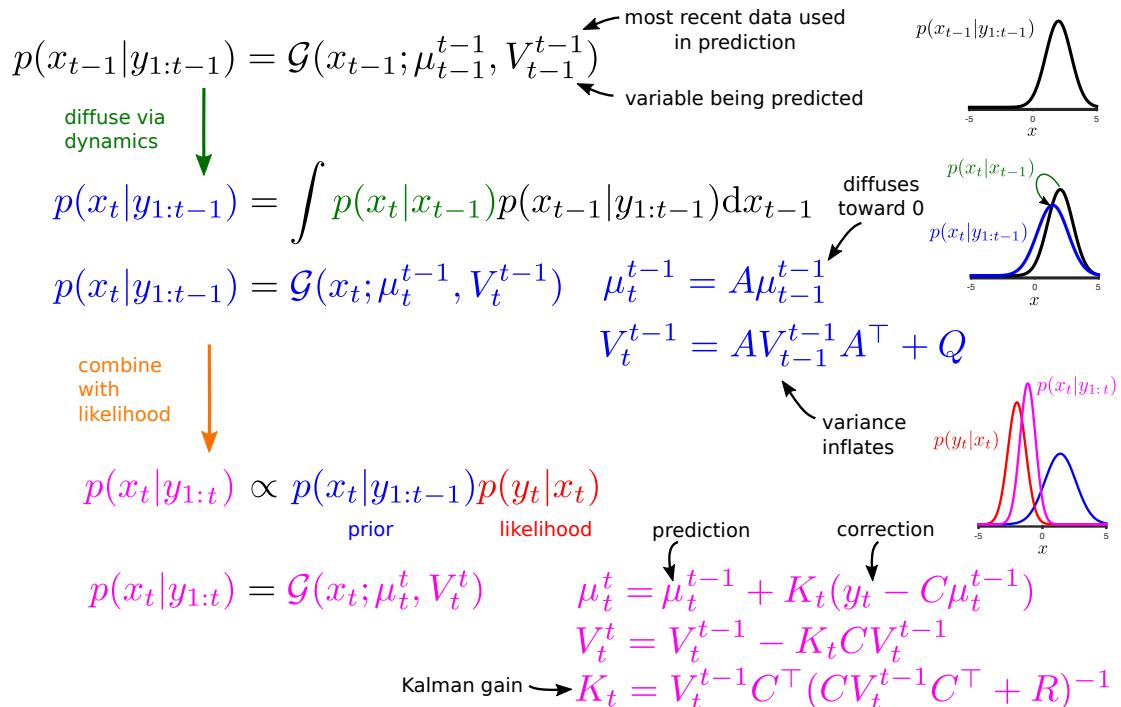
## Inference: Kalman Filter



## Inference: Derivation of General Filtering Equations

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

## Inference: Kalman Filter

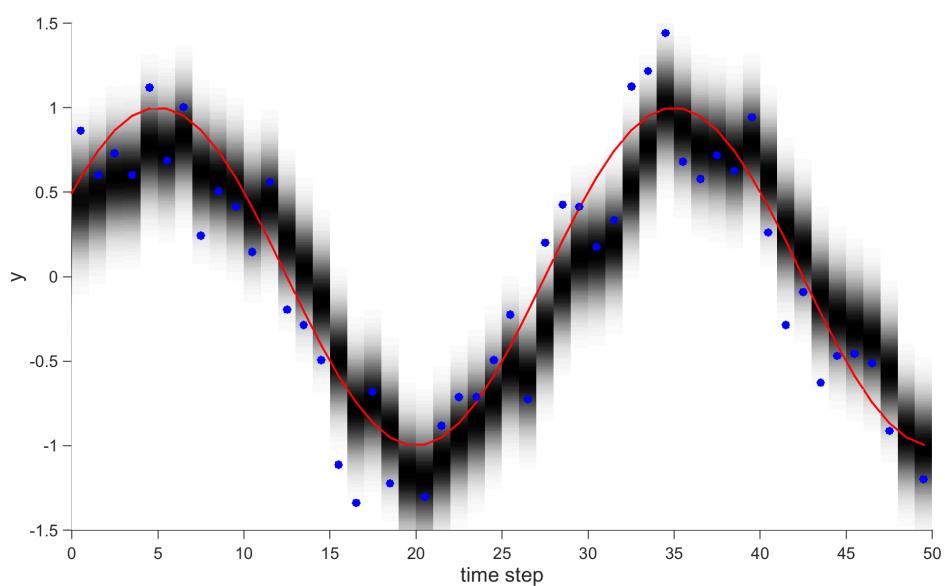


## Kalman Filter Demo

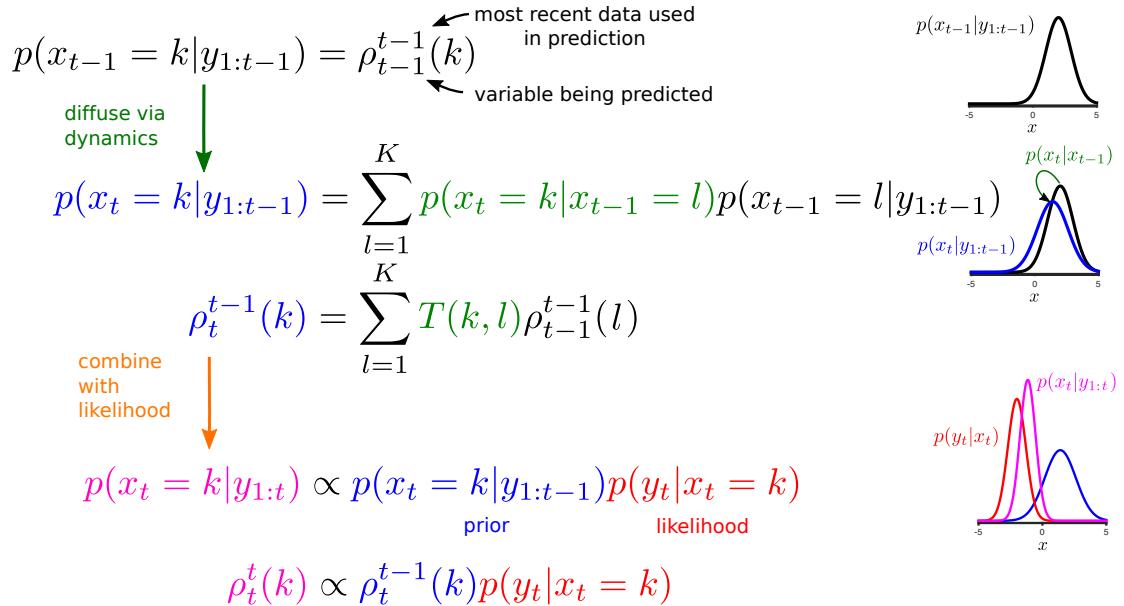
- ▶ data:  $y_t = \sin(\omega t) + \sigma_y \epsilon_t$  where  $\sigma_y^2 = 0.1$
- ▶ model:  $x_t = \lambda x_{t-1} + \sigma \eta$  and  $y_t = x_t + \sigma_y \eta'_t$  where  $\lambda = 0.99$  and  $\sigma^2 = 1 - \lambda^2$
- ▶ demo shows how the Kalman filter processes the data to form estimates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



## Inference: Forward Algorithm



When implementing, take care with numerical underflow/overflow.

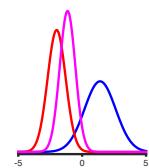
## Computing the likelihood

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

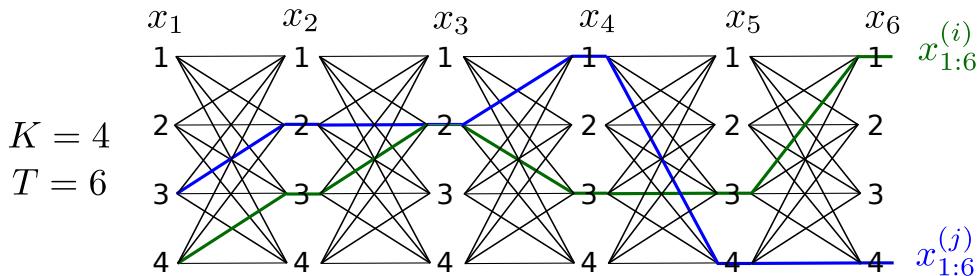
## The magic of the Forward Algorithm: Dynamic Programming

What's going on here?

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

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

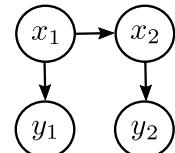
Trellis diagram represents possible sequences:



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

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

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

$$\overbrace{E(\theta; x_{1:T}, y_{1:T})}^{\text{simple form}}$$

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

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

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

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

requires posterior moments: marginals and pairwise marginals