# Neural Networks for Machine Learning

## Lecture 7a
## Modeling sequences: A brief overview

Geoffrey Hinton
Nitish Srivastava,
Kevin Swersky
Tijmen Tieleman
Abdel-rahman Mohamed
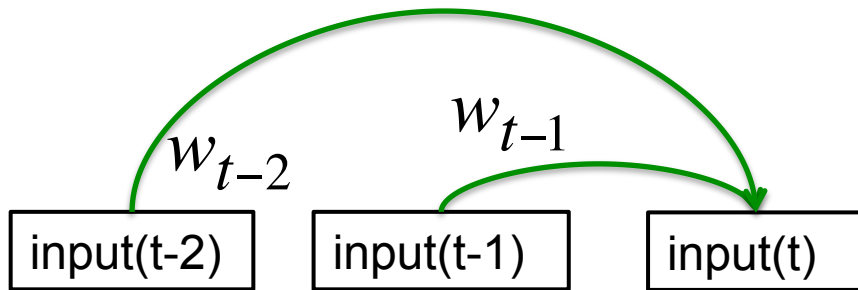
# Getting targets when modeling sequences

- When applying machine learning to sequences, we often want to turn an input sequence into an output sequence that lives in a different domain.
  - *E. g.* turn a sequence of sound pressures into a sequence of word identities.
- When there is no separate target sequence, we can get a teaching signal by trying to predict the next term in the input sequence.
  - The target output sequence is the input sequence with an advance of 1 step.
  - This seems much more natural than trying to predict one pixel in an image from the other pixels, or one patch of an image from the rest of the image.
  - For temporal sequences there is a natural order for the predictions.
- Predicting the next term in a sequence blurs the distinction between supervised and unsupervised learning.
  - It uses methods designed for supervised learning, but it doesn't require a separate teaching signal.

# Memoryless models for sequences
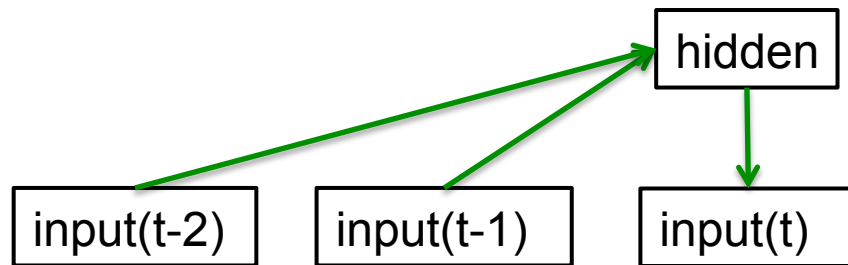
- Autoregressive models
  Predict the next term in a sequence from a fixed number of previous terms using "delay taps".

  *Linear autoregressive models: take weighted average of those previous terms*

$$w_{t-2} \qquad w_{t-1}$$

| input(t-2) | input(t-1) | input(t) |

- Feed-forward neural nets
  These generalize autoregressive models by using one or more layers of non-linear hidden units.
  *e.g.* Bengio's first language model.

hidden

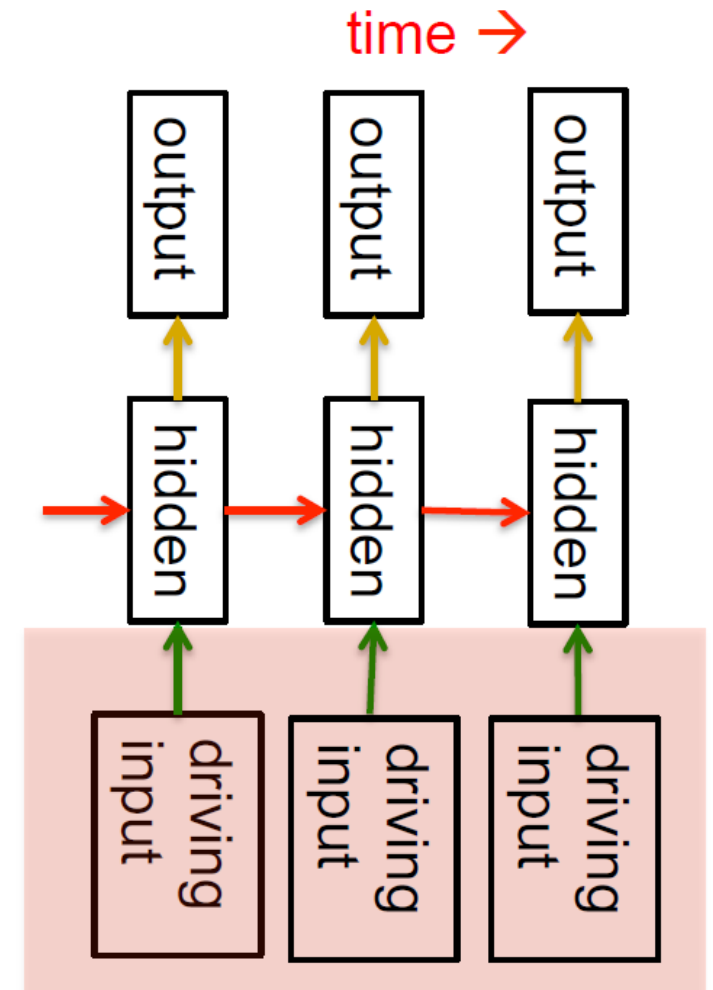| input(t-2) | input(t-1) | input(t) |

# Beyond memoryless models

- If we give our generative model some <u>hidden state</u>, and if we give this hidden state its own internal dynamics, we get a much more interesting kind of model.
  - It can <u>store information in its hidden state for a long time</u>.
  - If the dynamics is noisy and the way it generates outputs from its hidden state is noisy, we can never know its exact hidden state.
  - The best we can do is to <u>infer a probability distribution over the space of hidden state vectors</u>.
- This inference is only tractable for two types of hidden state model.
  - The next three slides are mainly intended for people who already know about these two types of hidden state model. They show how RNNs differ.
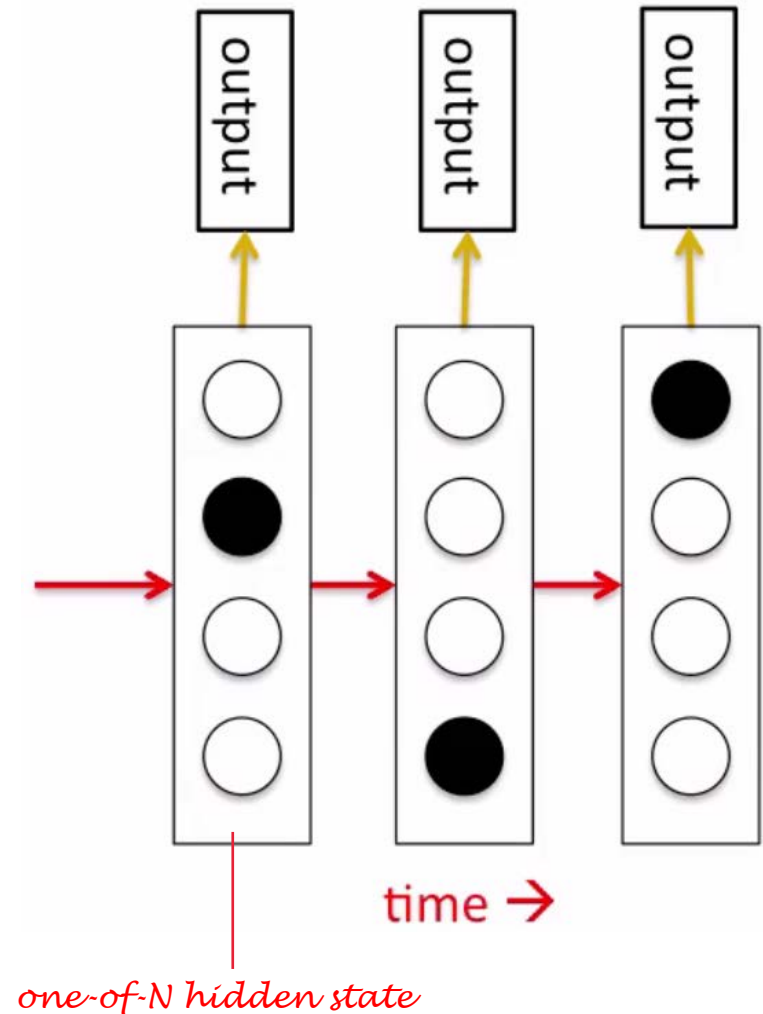  - Do not worry if you cannot follow the details.

# Linear Dynamic Systems (Engineers love them!)

- These are generative models. They have a <u>real-valued hidden state</u> that cannot be observed directly.
    - The <u>hidden state</u> has linear dynamics with <u>Gaussian noise</u> and produces the observations using a <u>linear model with Gaussian noise</u>.
    - There may also be driving inputs.

- To predict the next output (so that we can shoot down the missile) we need to infer the hidden state.
    - A <u>linearly transformed Gaussian is a Gaussian</u>. So the distribution over the hidden state given the data so far is Gaussian. It can be computed using "<u>Kalman filtering</u>".

time →

| output | output | output |

| hidden | hidden | hidden |

| driving input | driving input | driving input |

# Hidden Markov Models (Computer scientists love them!)

- Hidden Markov Models have a <u>discrete one-of-N hidden state</u>. Transitions between states are <u>stochastic</u> and controlled by a <u>transition matrix</u>. The outputs produced by a state are stochastic.

  - We cannot be sure which state produced a given output. So the state is "hidden".
  - It is easy to represent a probability distribution across N states with N numbers.

- To predict the next output we need to <u>infer the probability distribution over hidden states</u>.

  - HMMs have <u>efficient</u> algorithms for inference and learning.
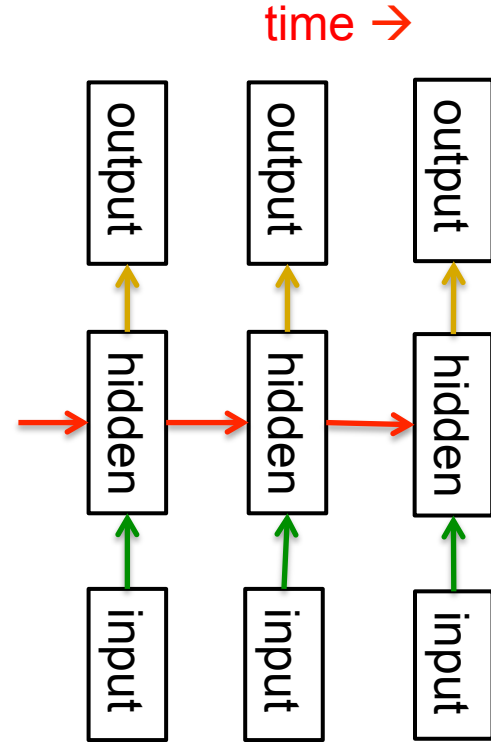


one-of-N hidden state

# A fundamental limitation of HMMs

- Consider what happens when a hidden Markov model generates data.
  - At each time step it must select one of its hidden states. So with N hidden states it can only remember log(N) bits about what it generated so far.

- Consider the information that the first half of an utterance contains about the second half:
  - The syntax needs to fit (e.g. number and tense agreement).
  - The semantics needs to fit. The intention needs to fit.
  - The accent, rate, volume, and vocal tract characteristics must all fit.

- All these aspects combined could be 100 bits of information that the first half of an utterance needs to convey to the second half. $2^{100}$ is big!

# Recurrent neural networks

- RNNs are very powerful, because they combine two properties:
  - Distributed hidden state that allows them to store a lot of information about the past efficiently.
  - Non-linear dynamics that allows them to update their hidden state in complicated ways.
- With enough neurons and time, RNNs can compute anything that can be computed by your computer.

time →

# Do generative models need to be stochastic?

- <u>Linear</u> dynamical systems and hidden Markov models are <u>stochastic</u> models.
  - But the <u>posterior</u> probability distribution over their hidden states given the observed data so far is a <u>deterministic</u> function of the data.

- Recurrent neural networks are <u>deterministic</u>.
  - So think of the hidden state of an RNN as the <u>equivalent</u> of the deterministic probability distribution over hidden states in a linear dynamical system or hidden Markov model.

# Recurrent neural networks

- What kinds of behaviour can RNNs exhibit?
  - They can oscillate. Good for motor control?
  - They can settle to point attractors. Good for retrieving memories?
  - They can behave chaotically. Bad for information processing?
  - RNNs could potentially learn to implement lots of small programs that each capture a nugget of knowledge and run in parallel, interacting to produce very complicated effects.
- But the computational power of RNNs makes them very hard to train.
  - For many years we could not exploit the computational power of RNNs despite some heroic efforts (e.g. Tony Robinson's speech recognizer).

# Neural Networks for Machine Learning

## Lecture 7b
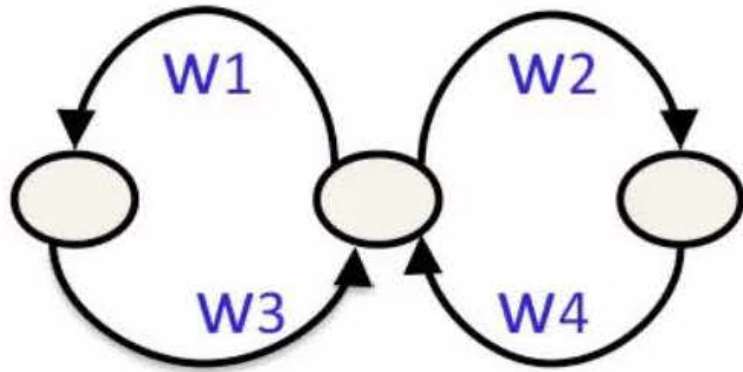## Training RNNs with backpropagation

Geoffrey Hinton
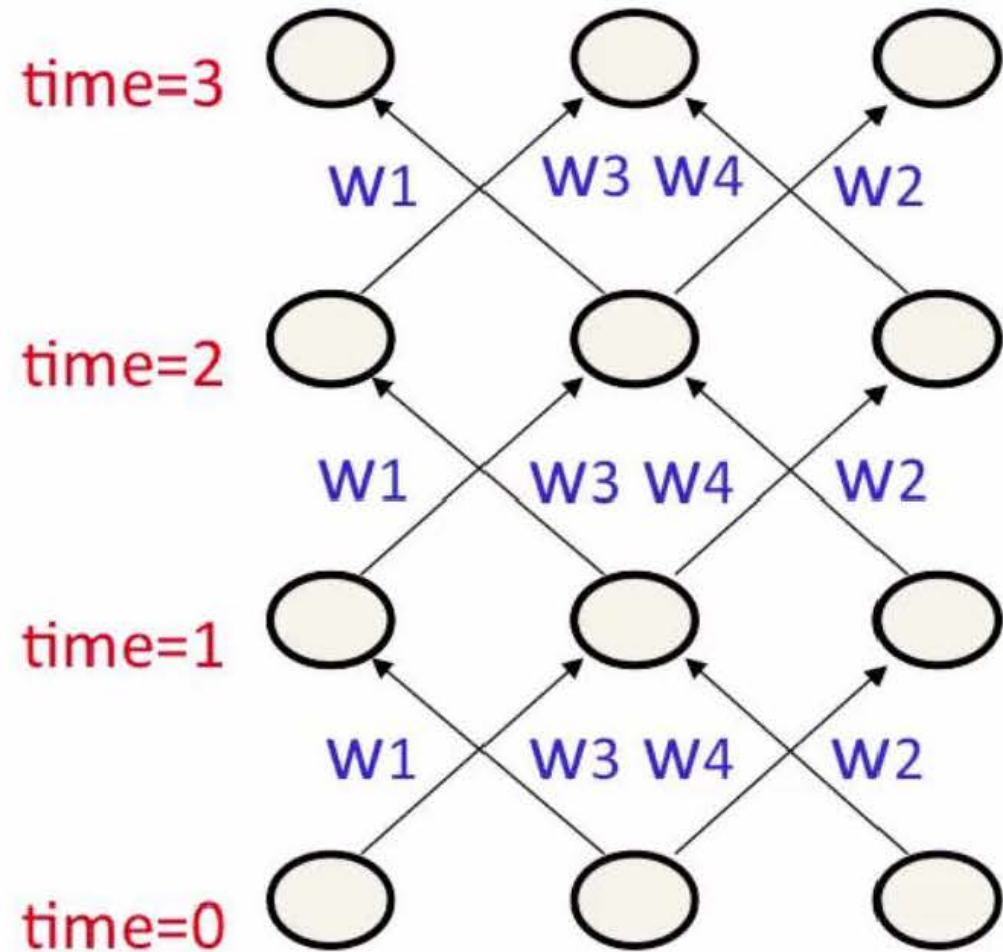Nitish Srivastava,
Kevin Swersky
Tijmen Tieleman
Abdel-rahman Mohamed

# The equivalence between feedforward nets and recurrent nets



Assume that there is a time delay of 1 in using each connection.

The recurrent net is just a layered net that keeps reusing the same weights.

# Reminder: Backpropagation with weight constraints

- It is easy to modify the backprop algorithm to incorporate linear constraints between the weights.

- We compute the gradients as usual, and then modify the gradients so that they satisfy the constraints.

  - So if the weights started off satisfying the constraints, they will continue to satisfy them.

$$To \ \ constrain: \quad w_1 = w_2$$

$$we \ \ need: \quad \Delta w_1 = \Delta w_2$$

$$compute: \quad \frac{\partial E}{\partial w_1} \quad and \quad \frac{\partial E}{\partial w_2}$$

$$use \quad \boxed{\frac{\partial E}{\partial w_1} + \frac{\partial E}{\partial w_2}} \quad for \ w_1 \ and \ w_2$$

# Backpropagation through time
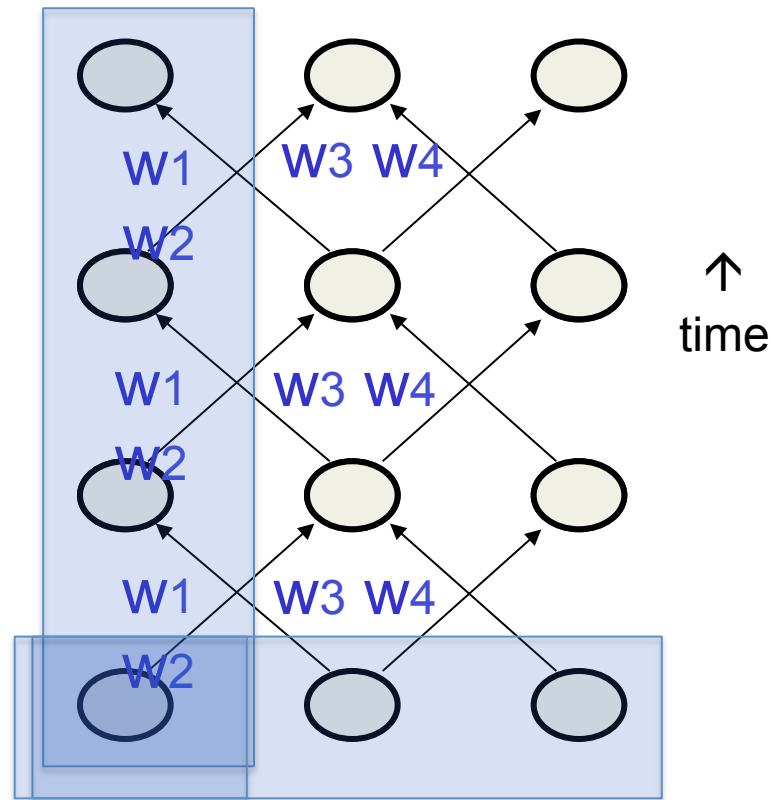
- We can think of the recurrent net as a <u>layered</u>, <u>feed-forward</u> net with <u>shared weights</u> and then train the <u>feed-forward net with weight constraints</u>.
- We can also think of this training algorithm in the time domain:
  - The forward pass builds up a stack of the activities of all the units at each time step.
  - The backward pass peels activities off the stack to compute the error derivatives at each time step.
  - After the backward pass we add together the derivatives at all the different times for each weight.

# An irritating extra issue

- We need to specify the initial activity state of all the hidden and output units.
- We could just fix these initial states to have some default value like 0.5.
- But it is better to treat the initial states as learned parameters.
- We learn them in the same way as we learn the weights.
  - Start off with an initial random guess for the initial states.
  - At the end of each training sequence, backpropagate through time all the way to the initial states  to get the gradient of the error function with respect to each initial state.
  - Adjust the initial states by following the negative gradient.

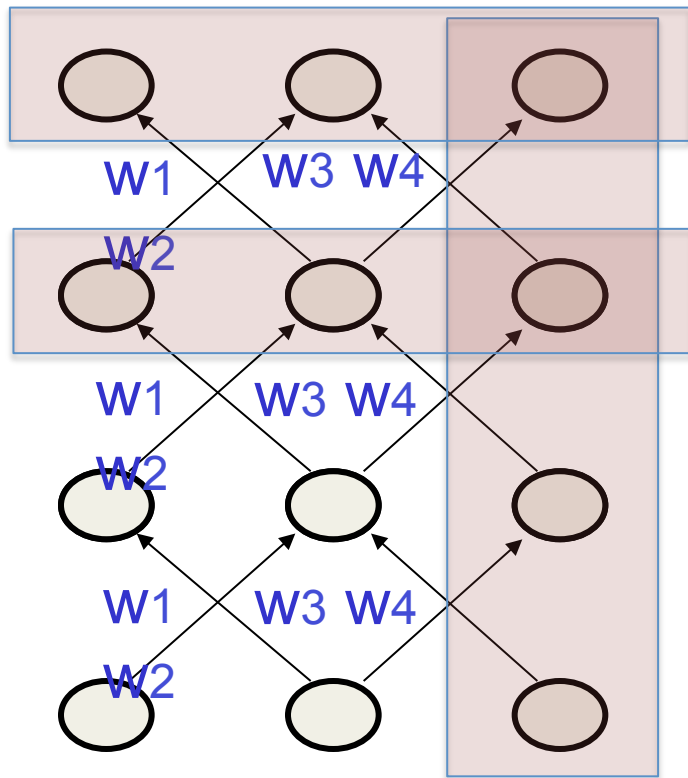# Providing input to recurrent networks

- We can specify inputs in several ways:
    - Specify the initial states of all the units.
    - Specify the initial states of a subset of the units.
    - Specify the states of the same subset of the units at every time step.
        - This is the natural way to model most sequential data.

W1   W3  W4

W2

W1   W3  W4

W2

W1   W3  W4

W2

↑
time

# Teaching signals for recurrent networks

- We can specify targets in several ways:
  - Specify desired final activities of all the units
  - Specify desired activities of all units for the last few steps
    - Good for learning attractors
    - It is easy to add in extra error derivatives as we backpropagate.
  - Specify the desired activity of a subset of the units.
    - The other units are input or hidden units.

# Neural Networks for Machine Learning

# Lecture 7c
# A toy example of training an RNN
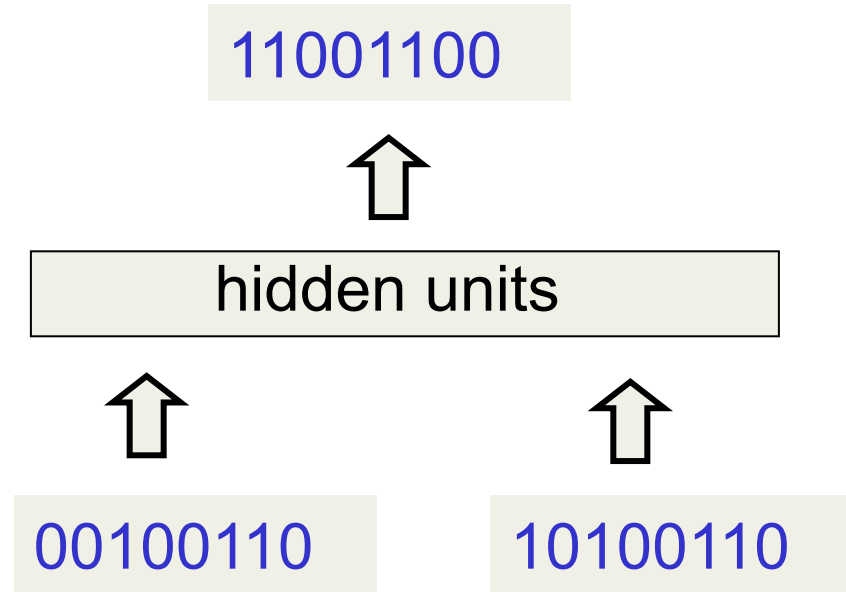
Geoffrey Hinton

Nitish Srivastava,
Kevin Swersky
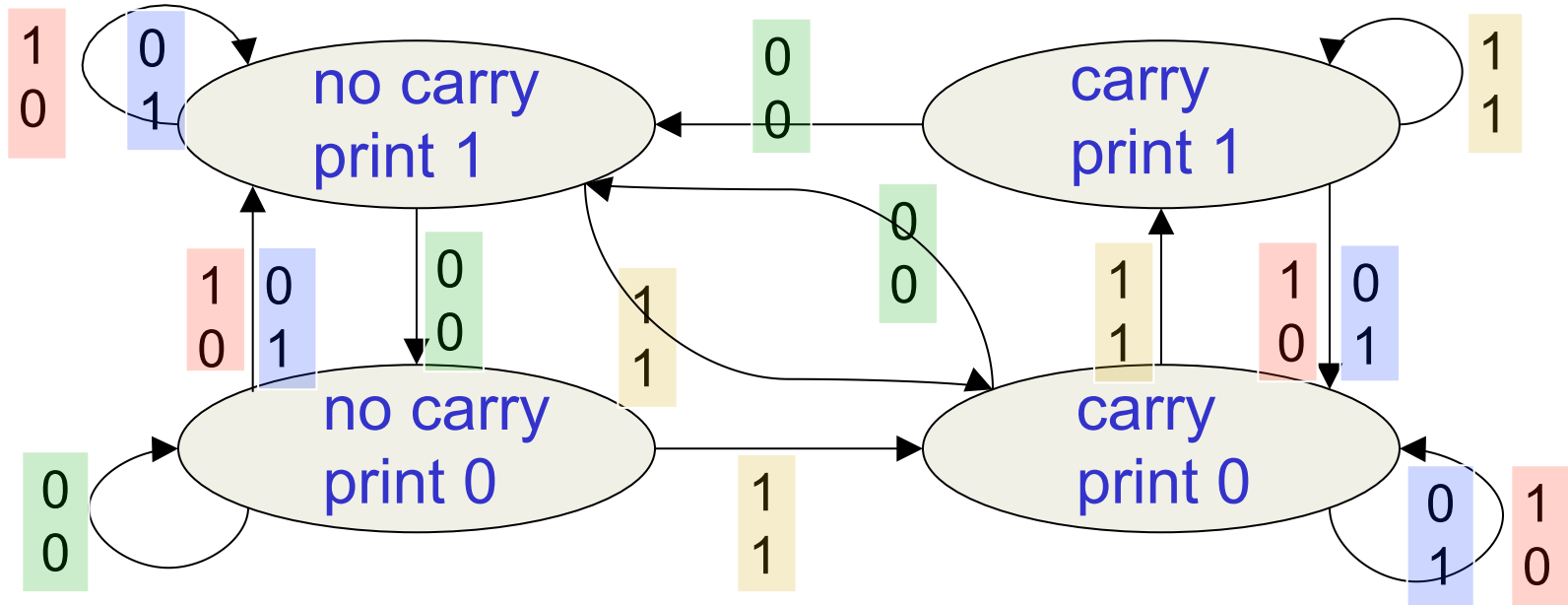Tijmen Tieleman
Abdel-rahman Mohamed

# A good toy problem for a recurrent network

- We can train a <u>feedforward net</u> to do <u>binary addition</u>, but there are obvious regularities that it cannot capture efficiently.
  - We must decide in advance the <u>maximum number of digits</u> in each number.
  - The processing applied to the beginning of a long number does not generalize to the end of the long number because     it uses different weights.
- As a result, feedforward nets <u>do not generalize well</u> on the binary addition task.

11001100

⇧

| hidden units |

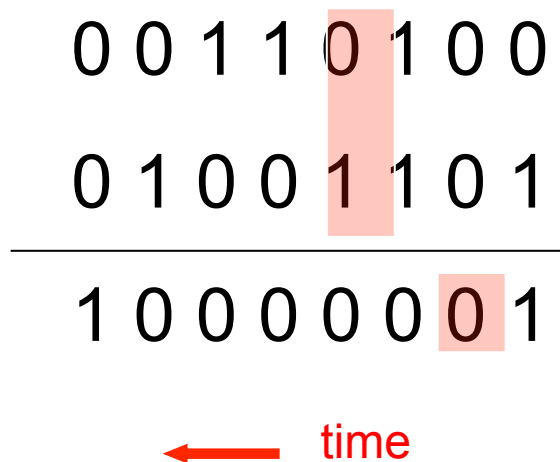⇧          ⇧

00100110          10100110

# The algorithm for binary addition



This is a finite state automaton. It decides what transition to make by looking at the next column. It prints after making the transition. It moves from right to left over the two input numbers.
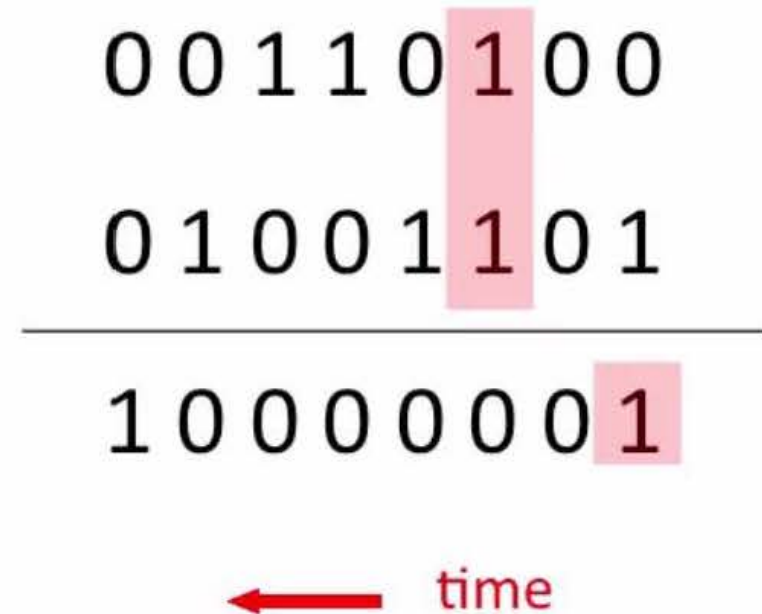
# A recurrent net for binary addition

- The network has two input units and one output unit.
- It is given <u>two input digits</u> at <u>each time step</u>.
- The desired output at each time step is the output for the column that was provided as input two time steps ago.
  - It takes <u>one time step</u> to <u>update the hidden units</u> based on the two input digits.
  - It takes <u>another time step</u> for the hidden units to cause the output.

```
0 0 1 1 0 1 0 0

0 1 0 0 1 1 0 1
_____
1 0 0 0 0 0 0 1
```
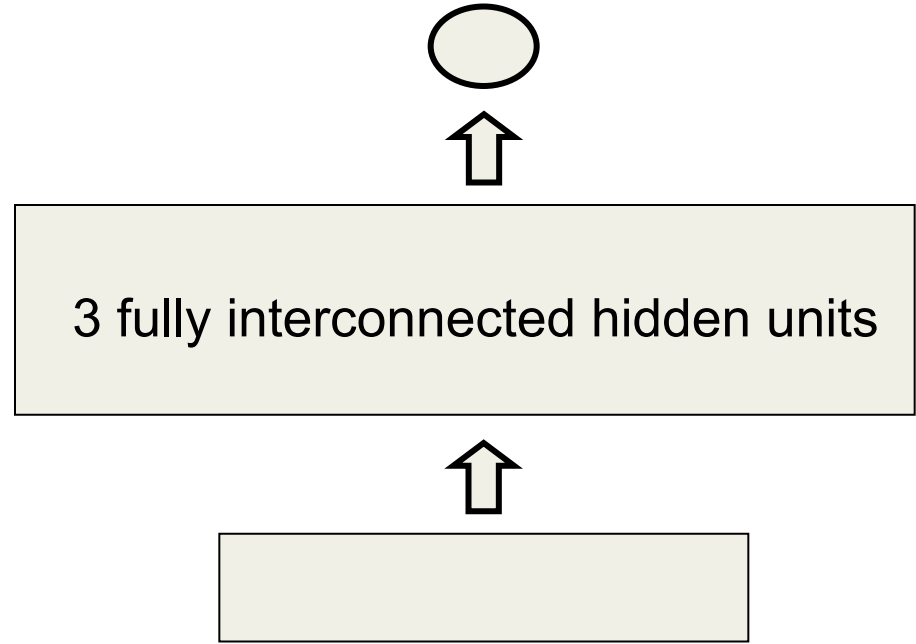
← time

# A recurrent net for binary addition

- The network has two input units and one output unit.

- It is given two input digits at each time step.

- The desired output at each time step is the output for the column that was provided as input <u>two time steps ago</u>.

  - It takes one time step to <u>update the hidden units based on the two input digits</u>.

  - It takes another time step <u>for the hidden units to cause the output</u>.

$$0\ 0\ 1\ 1\ 0\ 1\ 0\ 0$$

$$0\ 1\ 0\ 0\ 1\ 1\ 0\ 1$$

$$\rule{6cm}{0.4pt}$$

$$1\ 0\ 0\ 0\ 0\ 0\ 0\ 1$$

$\longleftarrow$ time

# The connectivity of the network

- The 3 hidden units are fully interconnected in both directions.
  - This allows a hidden activity pattern at one time step to vote for the hidden activity pattern at the next time step.
- The input units have feedforward connections that allow then to vote for the next hidden activity pattern.

3 fully interconnected hidden units

# What the network learns

- It learns <u>four distinct patterns of activity for the 3 hidden units</u>. These patterns correspond to the nodes in the finite state automaton.

  - Do not confuse units in a neural network with nodes in a finite state automaton. Nodes are like activity vectors.

  - The <u>automaton</u> is restricted to be in exactly one <u>state</u> at each time. The <u>hidden units</u> are restricted to have exactly one <u>vector of activity</u> at each time.

- A recurrent network can emulate a finite state automaton, but it is exponentially more powerful. With N hidden neurons it has 2^N possible binary activity vectors (but only N^2 weights)

  - This is important when the input stream has two separate things going on at once.

  - A <u>finite state automaton</u> needs to <u>square</u> its number of states.

  - An <u>RNN</u> needs to <u>double</u> its number of units.

# Neural Networks for Machine Learning

## Lecture 7d
## Why it is difficult to train an RNN

Geoffrey Hinton

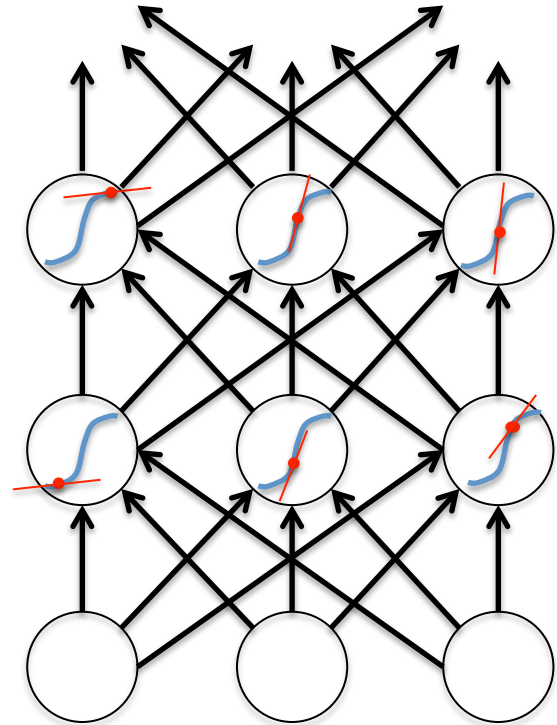Nitish Srivastava,

Kevin Swersky

Tijmen Tieleman

Abdel-rahman Mohamed

# The backward pass is linear
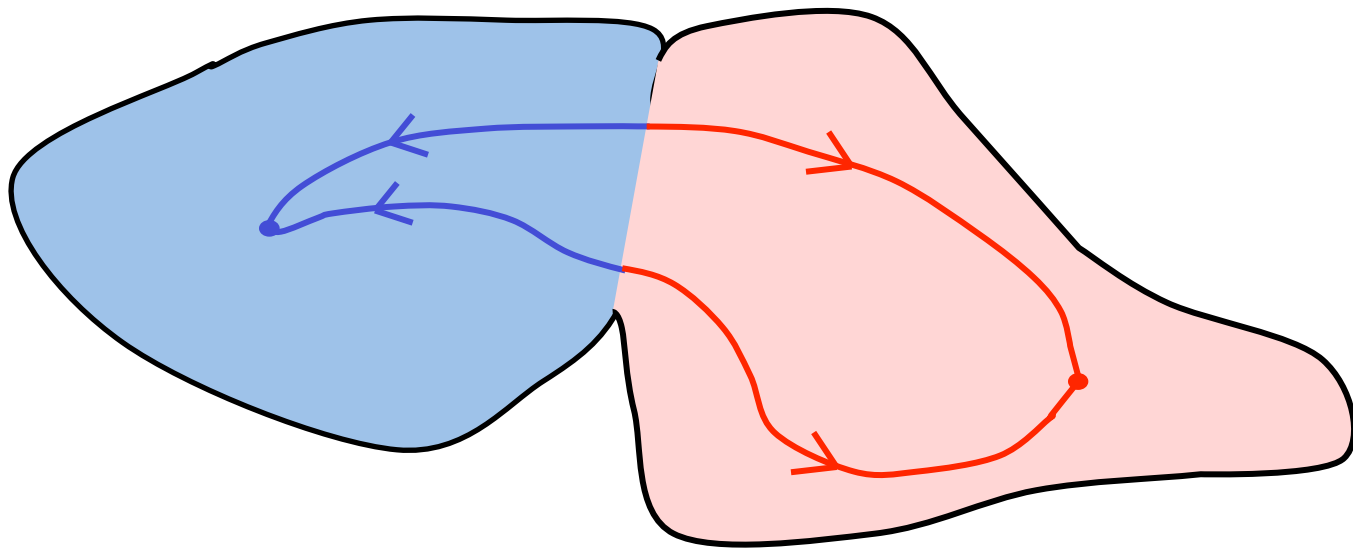
- There is a <u>big difference</u> between the forward and backward passes.

- In the <u>forward pass</u> we use <u>squashing functions</u> (like the logistic) to prevent the activity vectors from exploding.

- The <u>backward pass</u>, is completely <u>linear</u>. If you double the error derivatives at the final layer, all the error derivatives will double.

  – The forward pass determines the slope of the linear function used for backpropagating through each neuron.

# The problem of exploding or vanishing gradients

- What happens to the magnitude of the gradients as we backpropagate through many layers?
    - If the weights are small, the gradients shrink exponentially.
    - If the weights are big the gradients grow exponentially.
- Typical feed-forward neural nets can cope with these exponential effects because they only have a few hidden layers.

- In an RNN trained on long sequences (*e.g.* 100 time steps) the gradients can easily explode or vanish.
    - We can avoid this by initializing the weights very carefully.
- Even with good initial weights, its very hard to detect that the current target output depends on an input from many time-steps ago.
    - So RNNs have difficulty dealing with long-range dependencies.

# Why the back-propagated gradient blows up



- If we start a trajectory within an attractor, small changes in where we start make no difference to where we end up.

- But if we start almost exactly on the boundary, tiny changes can make a huge difference.

# Four effective ways to learn an RNN

- **Long Short Term Memory**
Make the RNN out of little modules that are <u>designed to remember values for a long time</u>.

- **Hessian Free Optimization:** Deal with the vanishing gradients problem by using a fancy optimizer that can detect directions with a tiny gradient but even smaller curvature.

  – The HF optimizer ( Martens & Sutskever, 2011) is good at this.

- **Echo State Networks:** <u>Initialize</u> the input→hidden and hidden→hidden and output→hidden connections <u>very carefully</u> so that the hidden state has a huge reservoir of weakly coupled oscillators which can be selectively driven by the input.

  – <u>ESNs</u> only need to learn the <u>hidden→output</u> connections.

- **Good initialization with momentum** <u>Initialize like in Echo State Networks</u>, but then learn all of the connections using momentum.

# Neural Networks for Machine Learning

# Lecture 7e
# Long term short term memory

Geoffrey Hinton

Nitish Srivastava,

Kevin Swersky

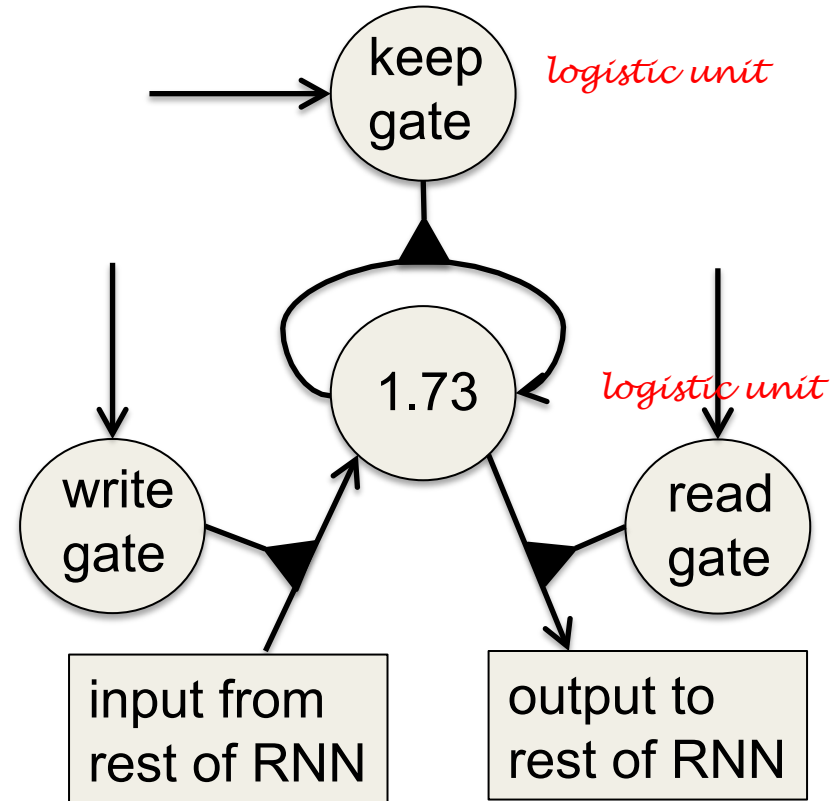Tijmen Tieleman

Abdel-rahman Mohamed

# Long Short Term Memory (LSTM)

- Hochreiter & Schmidhuber (1997) solved the problem of getting an RNN to remember things for a long time (like hundreds of time steps).

- They designed a memory cell using logistic and linear units with multiplicative interactions.

- Information gets into the cell whenever its "write" gate is on.

- The information stays in the cell so long as its "keep" gate is on.

- Information can be read from the cell by turning on its "read" gate.

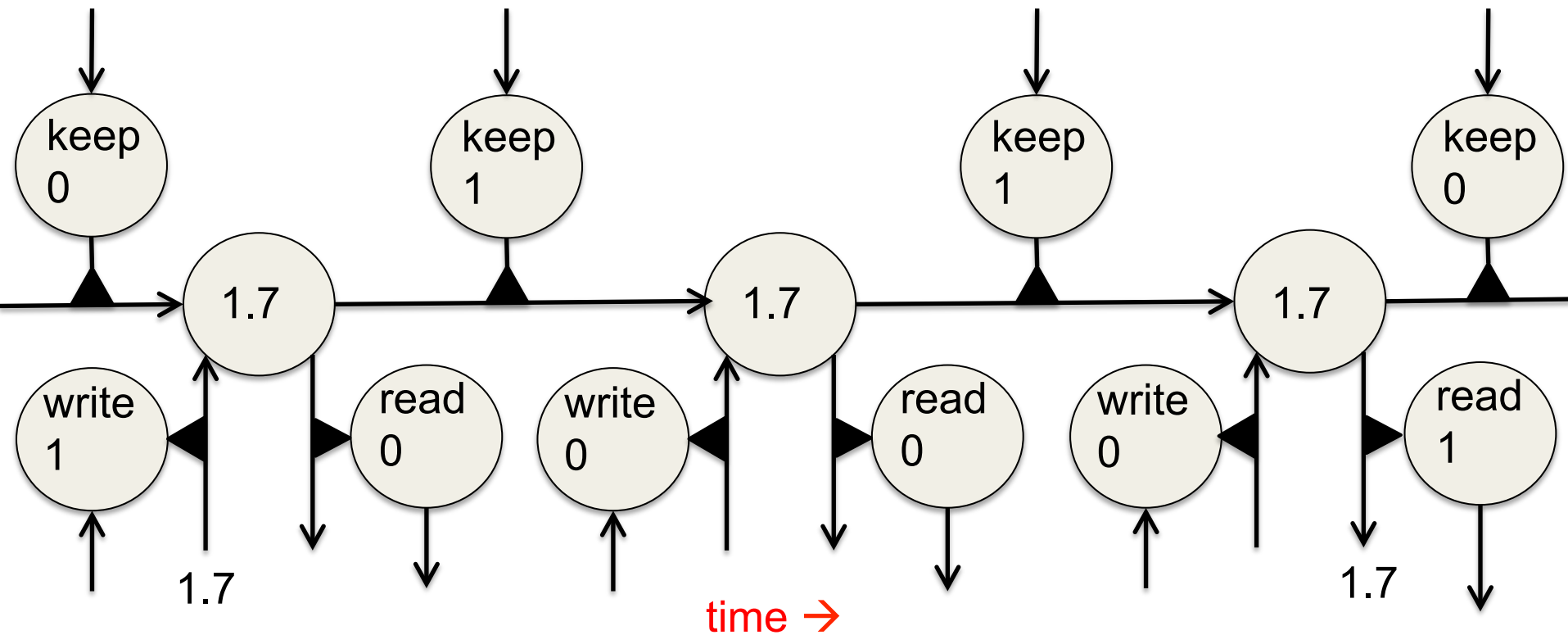# Implementing a memory cell in a neural network

To preserve information for a long time in the activities of an RNN, we use a circuit that implements an analog memory cell.

- A <u>linear unit</u> that has a self-link with a weight of 1 will maintain its state.
- Information is <u>stored</u> in the cell by activating its <u>write</u> gate.
- Information is <u>retrieved</u> by activating the <u>read</u> gate.
- We can <u>backpropagate</u> through this circuit because <u>logistics</u> are have <u>nice derivatives</u>.

Backpropagation through a memory cell

# Reading cursive handwriting

- This is a natural task for an RNN.

- The input is a sequence of (x,y,p) coordinates of the tip of the pen, where p indicates whether the pen is up or down.

- The output is a sequence of characters.

p: whether the pen is on the paper or not.

- Graves & Schmidhuber (2009) showed that RNNs with LSTM are currently the best systems for reading cursive writing.

  – They used a sequence of small images as input rather than pen coordinates.

# A demonstration of online handwriting recognition by an RNN with Long Short Term Memory (from Alex Graves)

- The movie that follows shows several different things:
- Row 1:  This shows when the characters are recognized.
  - It never revises its output so difficult decisions are more delayed.
- Row 2:  This shows the states of a subset of the memory cells.
  - Notice how they get reset when it recognizes a character.
- Row 3:  This shows the writing. The net sees the x and y coordinates.
  - Optical input actually works a bit better than pen coordinates.
- Row 4:  This shows the gradient backpropagated all the way to the x and y inputs from the currently most active character.
  - This lets you see which bits of the data are influencing the decision.