# ML vs DL

# Outline

- An overview of DL components
  - Historical remarks: early days of neural networks
  - Modern building blocks: units, layers, activations functions, loss functions, etc.
  - Reverse-mode automatic differentiation (aka backpropagation)

- Similarities and differences between GMs and NNs
  - Graphical models vs. computational graphs
  - Sigmoid Belief Networks as graphical models
  - Deep Belief Networks and Boltzmann Machines

- Combining DL methods and GMs
  - Using outputs of NNs as inputs to GMs
  - GMs with potential functions represented by NNs
  - NNs with structured outputs

- Bayesian Learning of NNs
  - Bayesian learning of NN parameters
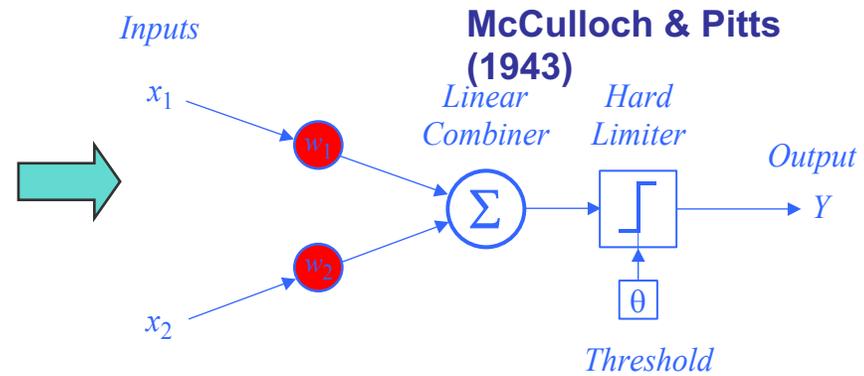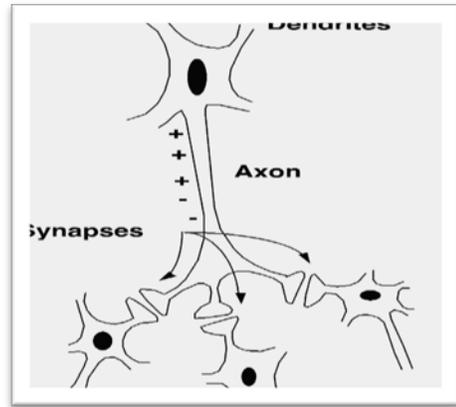  - Deep kernel learning

# Outline

- An overview of DL components
  - Historical remarks: early days of neural networks
  - Modern building blocks: units, layers, activations functions, loss functions, etc.
  - Reverse-mode automatic differentiation (aka backpropagation)

- Similarities and differences between GMs and NNs
  - Graphical models vs. computational graphs
  - Sigmoid Belief Networks as graphical models
  - Deep Belief Networks and Boltzmann Machines

- Combining DL methods and GMs
  - Using outputs of NNs as inputs to GMs
  - GMs with potential functions represented by NNs
  - NNs with structured outputs

- Bayesian Learning of NNs
  - Bayesian learning of NN parameters
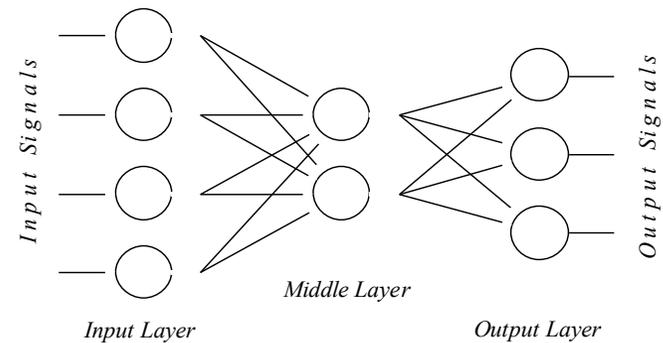  - Deep kernel learning

# Perceptron and Neural Nets

❑ From biological neuron to artificial neuron (perceptron)



Inputs

McCulloch & Pitts (1943)

Linear Combiner    Hard Limiter

$x_1$

$w_1$

$\Sigma$

Output

$Y$

$w_2$

$x_2$

$\theta$

Threshold

❑ From biological neuron network to artificial neuron networks



Synapse

Synapse    Dendrites

Axon    Axon

Soma

Soma

Dendrites

Synapse

Input Signals    Output Signals

Middle Layer

Input Layer    Output Layer

# The perceptron learning algorithm



❑ Recall the nice property of sigmoid function

❑ Consider regression problem f: X➔Y, for scalar $\frac{d\sigma}{dt} = \sigma(1 - \sigma)$

❑ We used to maximize the conditional data likelihood $y = f(x) + \epsilon$

$$\vec{w} = \arg\max_{\vec{w}} \ln \prod_i P(y_i | x_i; \vec{w})$$

❑ Here …

$$\vec{w} = \arg\min_{\vec{w}} \sum_i \frac{1}{2}(y_i - \hat{f}(x_i; \vec{w}))^2$$

# The perceptron learning algorithm

$$\frac{\partial E_D[\vec{w}])}{\partial w_j} = \frac{\partial}{\partial w_i} \frac{1}{2} \sum_l (t_d - o_d)^2$$

$$= \frac{1}{2} \sum_d 2(t_d - o_d) \frac{\partial}{\partial w_i}(t_d - o_d)$$

$$= \sum_d (t_d - o_d)\left(-\frac{\partial o_d}{\partial w_i}\right)$$

$$= -\sum_d (t_d - o_d) \frac{\partial o_d}{\partial net_i} \frac{\partial net_d}{\partial w_i}$$

$$= -\sum_d (t_d - o_d) o_d (1 - o_d) x_d^i$$

$x_d$ = input

$t_d$ = target output

$o_d$ = observed output

$w_i$ = weight i

Batch mode:

Do until converge:

   1. compute gradient $\nabla E_D[w]$

2. $\vec{w} = \vec{w} - \eta \nabla E_D[\vec{w}]$

Incremental mode:

Do until converge:

▪ For each training example $d$ in $D$

   1. compute gradient $\nabla E_d[w]$
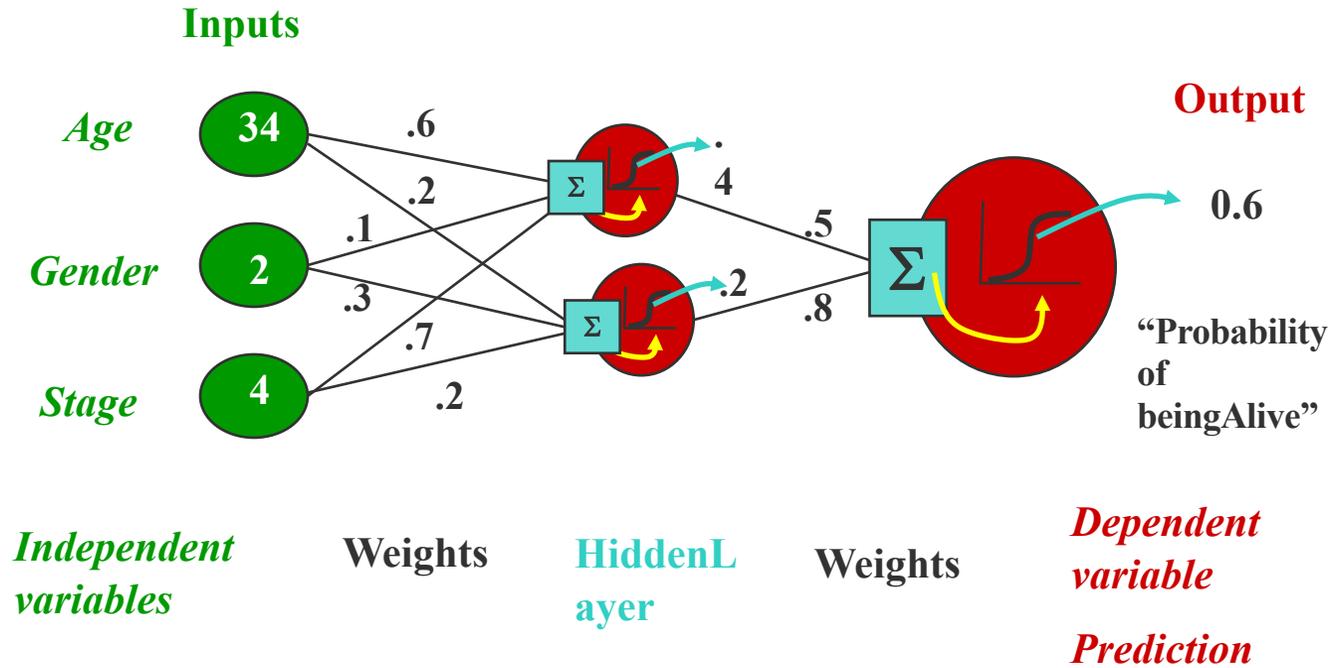
2. $\vec{w} = \vec{w} - \eta \nabla E_d[\vec{w}]$

   where

$$\nabla E_d[\vec{w}] = -(t_d - o_d)o_d(1 - o_d)\vec{x}_d$$

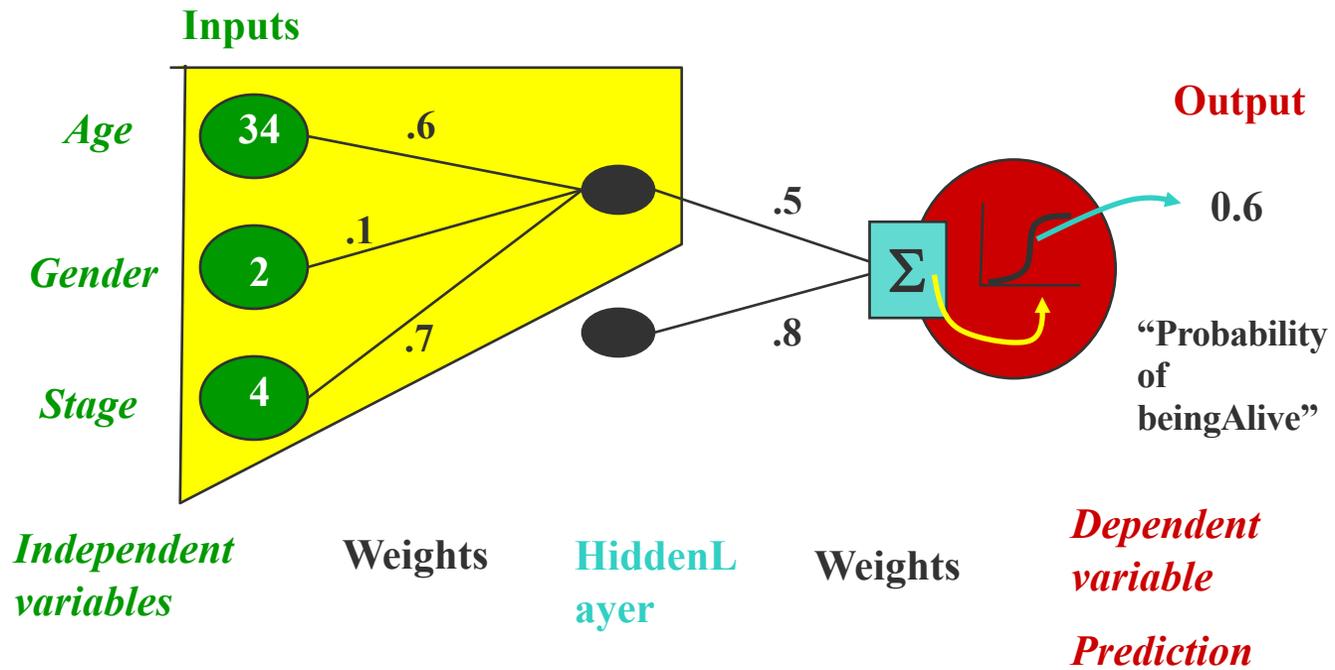# Neural Network Model



Inputs

*Age* 34 .6
.2
*Gender* 2 .1
.3
*Stage* 4 .7
.2

Σ .
4
Σ .2

.5
.8

Σ

Output

0.6

"Probability of beingAlive"

*Independent variables*   **Weights**   HiddenLayer   **Weights**   *Dependent variable*

*Prediction*

# "Combined logistic models"



Inputs

Age — 34 — .6

Gender — 2 — .1

Stage — 4 — .7

Output

0.6

"Probability of beingAlive"

.5

.8

Σ

Independent variables

Weights

HiddenLayer

Weights

Dependent variable

Prediction

# "Combined logistic models"



**Inputs**

*Age* 34

*Gender* 2

*Stage* 4

.2

.3

.2

.5

.8

Σ

**Output**

0.6

"Probability of beingAlive"

*Independent variables*

Weights

HiddenLayer

Weights

*Dependent variable*

*Prediction*

# "Combined logistic models"

# Not really, no target for hidden units...



Independent variables: Age 34, Gender 2, Stage 4

Weights: .6, .2, .1, .3, .7, .2

Hidden Layer

Weights: .5, .8, .4, .2

Dependent variable: "Probability of being Alive" 0.6

Prediction

# Backpropagation: Reverse-mode differentiation

❑ Artificial neural networks are nothing more than complex functional compositions that can be represented by **computation graphs**:



$$x \quad \text{(1)} \quad \begin{array}{c} \text{(2)} \\ \text{(3)} \end{array} \quad \text{(4)} \quad \text{(5)} \quad f(x) \qquad \frac{\partial f_n}{\partial x} =$$

Input variables

Intermediate computations

Outputs

# Backpropagation: Reverse-mode differentiation

- Artificial neural networks are nothing more than complex functional compositions that can be represented by **computation graphs**:
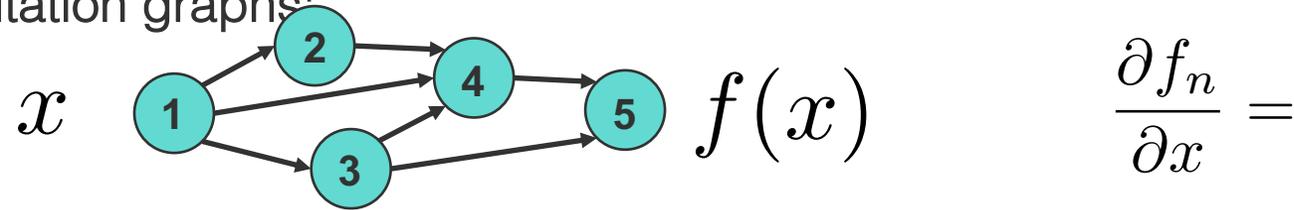
$$x \quad \text{1} \to \text{2} \to \text{4} \to \text{5} \quad f(x) \qquad \frac{\partial f_n}{\partial x} =$$

- By applying the chain rule and using reverse accumulation, we get

$$\frac{\partial f_n}{\partial x} = \sum_{i_1 \in \pi(n)} \frac{\partial f_n}{\partial f_{i_1}} \frac{\partial f_{i_1}}{\partial x} = \sum_{i_1 \in \pi(n)} \frac{\partial f_n}{\partial f_{i_1}} \sum_{i_2 \in \pi(i_1)} \frac{\partial f_{i_1}}{\partial f_{i_2}} \frac{\partial f_{i_1}}{\partial x} = \dots$$
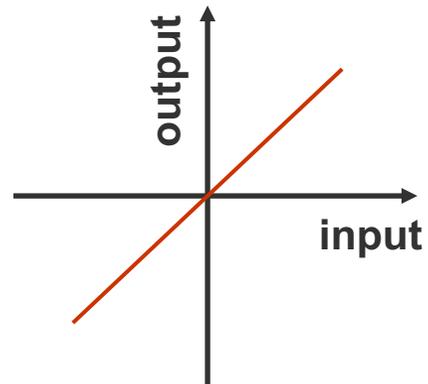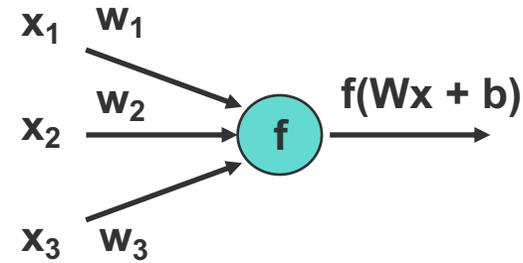
- The algorithm is commonly known as backpropagation
- What if some of the functions are stochastic?
- Then use stochastic backpropagation!
  (to be covered in the next part)
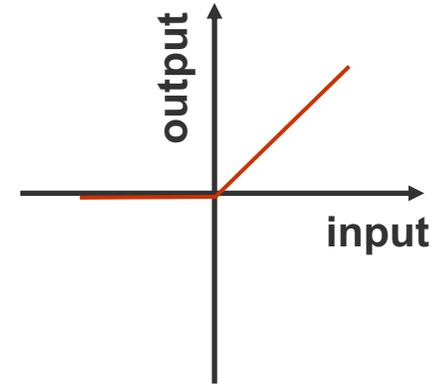- Modern packages can do this *automatically* (more later)

# Modern building blocks of deep networks

- ❑ Activation functions
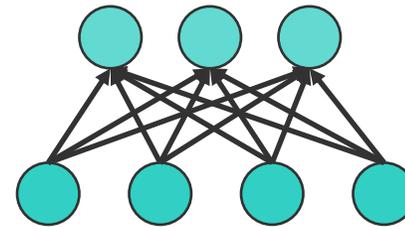  - ❑ Linear and ReLU
  - ❑ Sigmoid and tanh
  - ❑ Etc.

$x_1$ $w_1$

$x_2$ $w_2$    f    $f(Wx + b)$

$x_3$ $w_3$

**Linear**

**Rectified linear (ReLU)**

# Modern building blocks of deep networks

- ❑ Activation functions
  - ❑ Linear and ReLU
  - ❑ Sigmoid and tanh
  - ❑ Etc.
- ❑ Layers
  - ❑ Fully connected
  - ❑ Convolutional & pooling
  - ❑ Recurrent
  - ❑ ResNets
  - ❑ Etc.

fully connected

convolutional

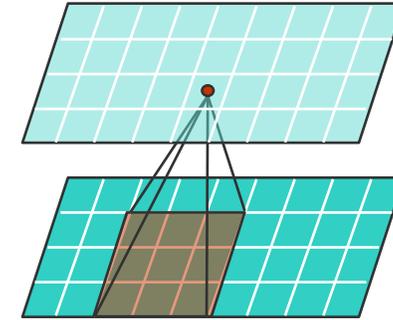recurrent

blocks with residual connections

# Modern building blocks of deep networks
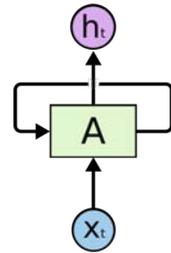
- Activation functions
  - Linear and ReLU
  - Sigmoid and tanh
  - Etc.
- Layers
  - Fully connected
  - Convolutional & pooling
  - Recurrent
  - ResNets
  - Etc.
- Loss functions
  - Cross-entropy loss
  - Mean squared error
  - Etc.

**Putting things together:**



(a part of GoogleNet)

# Modern building blocks of deep networks

- Activation functions
  - Linear and ReLU
  - Sigmoid and tanh
  - Etc.
- Layers
  - Fully connected
  - Convolutional & pooling
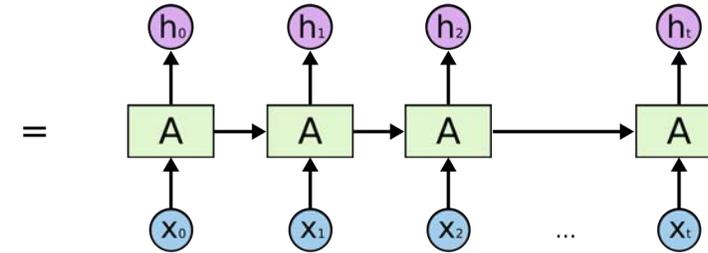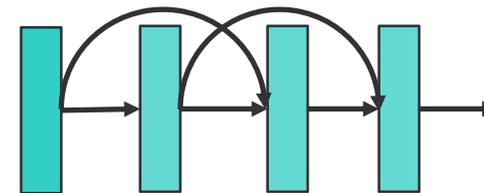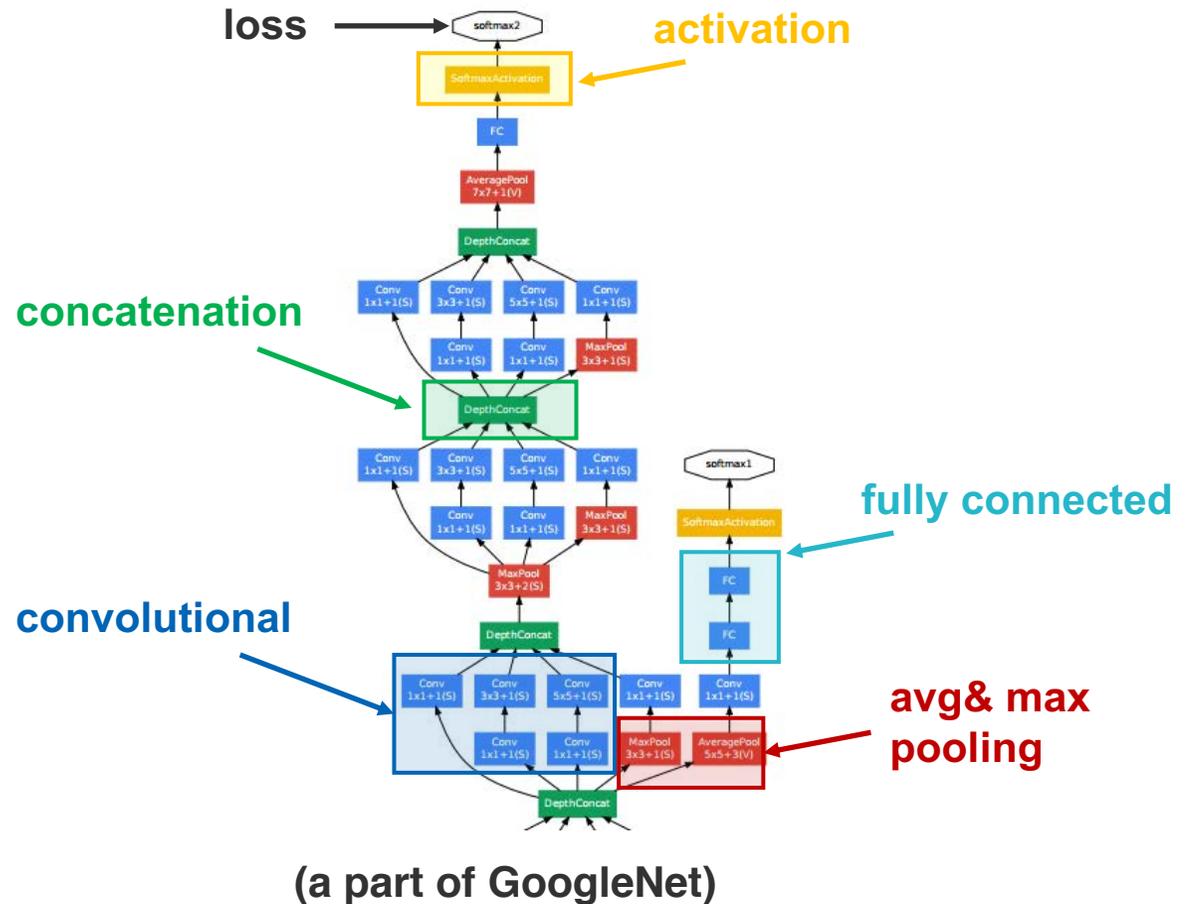  - Recurrent
  - ResNets
  - Etc.
- Loss functions
  - Cross-entropy loss
  - Mean squared error
  - Etc.

**Putting things together:**



**(a part of GoogleNet)**

- **Arbitrary combinations of the basic building blocks**
- **Multiple loss functions – multi-target prediction, transfer learning, and more**
- **Given enough data, deeper architectures just keep improving**
- **Representation learning: the networks learn increasingly more abstract representations of the data that are "disentangled," i.e., amenable to linear separation.**

# Feature learning

❑ Successful learning of intermediate representations
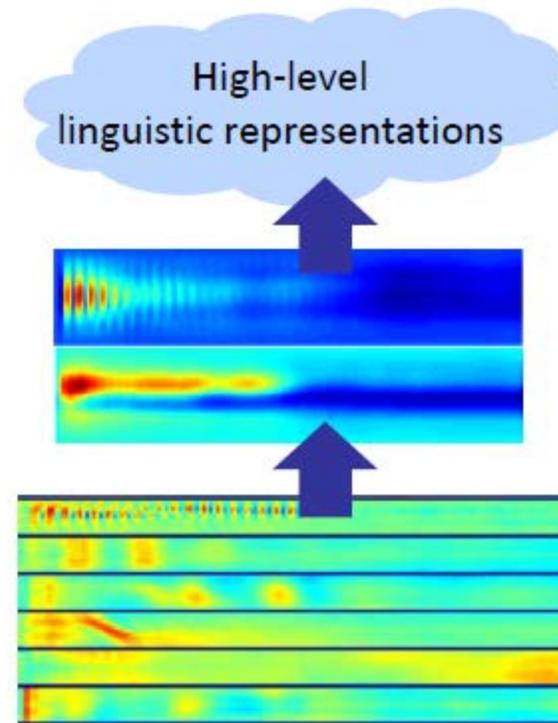[Lee et al ICML 2009, Lee et al NIPS 2009]
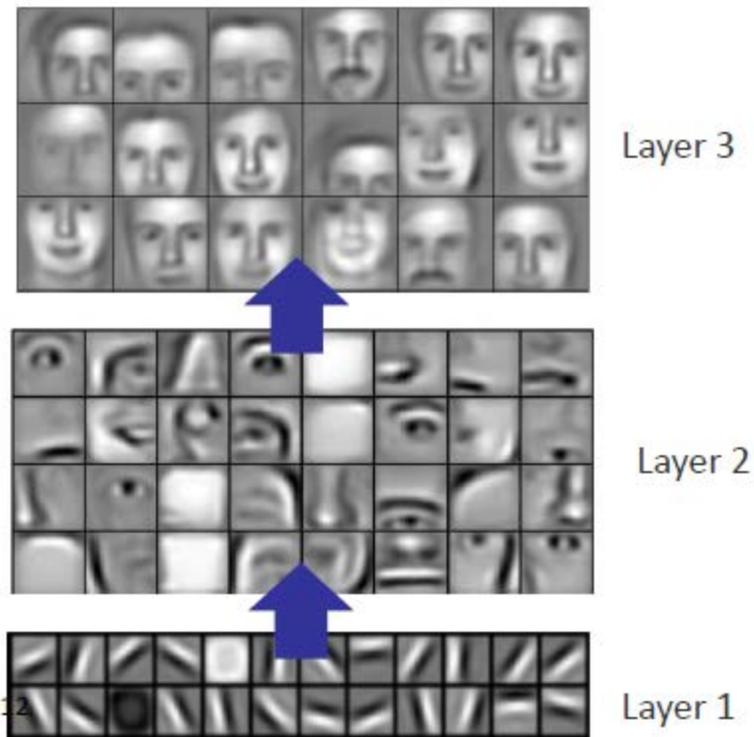


- **Arbitrary combinations of the basic building blocks**
- **Multiple loss functions – multi-target prediction, transfer learning, and more**
- **Given enough data, deeper architectures just keep improving**
- **Representation learning: the networks learn increasingly more abstract representations of the data that are "disentangled," i.e., amenable to linear separation.**

# Outline

- An overview of the DL components
    - Historical remarks: early days of neural networks
    - Modern building blocks: units, layers, activations functions, loss functions, etc.
    - Reverse-mode automatic differentiation (aka backpropagation)

- Similarities and differences between GMs and NNs
    - Graphical models vs. computational graphs
    - Sigmoid Belief Networks as graphical models
    - Deep Belief Networks and Boltzmann Machines

- Combining DL methods and GMs
    - Using outputs of NNs as inputs to GMs
    - GMs with potential functions represented by NNs
    - NNs with structured outputs

- Bayesian Learning of NNs
    - Bayesian learning of NN parameters
    - Deep kernel learning

# Graphical models vs. Deep nets

## Graphical models

- <u>Representation</u> for encoding meaningful knowledge and the associated uncertainty in a graphical form



## Deep neural networks

- <u>Learn representations</u> that facilitate computation and performance on the end-metric (intermediate representations are not guaranteed to be meaningful)
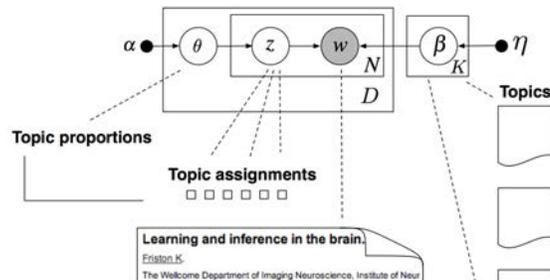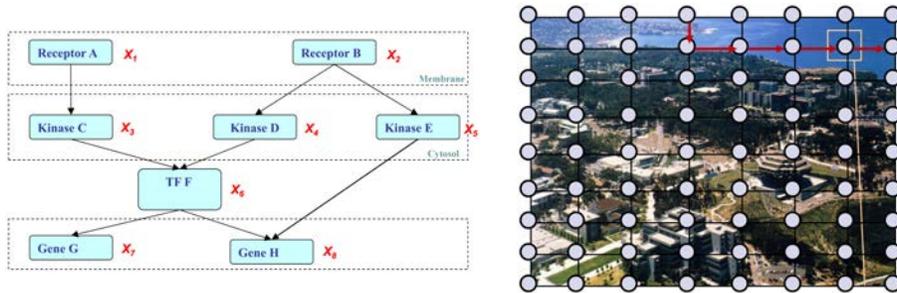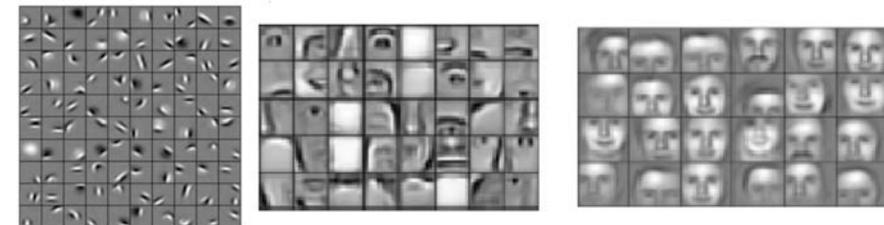
# Graphical models vs. Deep nets

## Graphical models

- Representation for encoding meaningful knowledge and the associated uncertainty in a graphical form

- Learning and inference are based on a rich toolbox of well-studied (structure-dependent) techniques (e.g., EM, message passing, VI, MCMC, etc.)

- Graphs represent models

## Deep neural networks

- Learn representations that facilitate computation and performance on the end-metric (intermediate representations are not guaranteed to be meaningful)

- Learning is predominantly based on the gradient descent method (aka backpropagation); Inference is often trivial and done via a "forward pass"

- Graphs represent computation

# Graphical models vs. Deep nets

## Graphical models
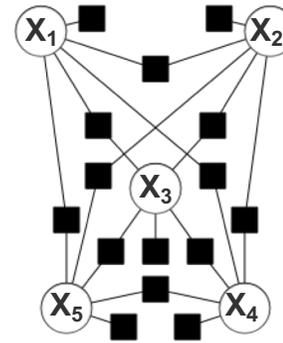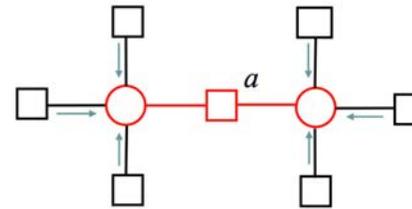
### Utility of the graph

- A vehicle for synthesizing a global loss function from local structure
  - potential function, feature function, etc.
- A vehicle for designing sound and efficient inference algorithms
  - Sum-product, mean-field, etc.
- A vehicle to inspire approximation and penalization
  - Structured MF, Tree-approximation, etc.
- A vehicle for monitoring theoretical and empirical behavior and accuracy of inference

### Utility of the loss function

- A major measure of quality of the learning algorithm and the model

$$\log P(X) = \sum_i \log \phi(x_i) + \sum_{i,j} \log \psi(x_i, x_j)$$

$$m_{i\to a}(x_i) = \prod_{c\in N(i)\setminus a} m_{c\to i}(x_i)$$

$$b_a(X_a) \propto f_a(X_a) \prod_{i\in N(a)} m_{i\to a}(x_i)$$

$$m_{a\to i}(x_i) = \sum_{X_a\setminus x_i} f_a(X_a) \prod_{j\in N(a)\setminus i} m_{j\to a}(x_j)$$

$$Q(H) \sim P(H|V)$$

$$\Omega := \left\{ \theta \in \mathbb{R}^d \,\middle|\, A(\theta) < +\infty \right\}$$

$$F_0:$$

$$T:$$

$$\Omega(F_0) := \left\{ \theta \in \Omega \,\middle|\, \theta_{(s,t)} = 0 \,\forall\, (s,t) \in E \right\}. \quad \Omega(T) := \left\{ \theta \in \Omega \,\middle|\, \theta_{(s,t)} = 0 \,\forall\, (s,t) \notin E(T) \right\}$$

$$\theta = \operatorname{argmax}_\theta P_\theta(V)$$

# Graphical models vs. Deep nets



Images from Distill.pub

## Deep neural networks

Utility of the network

- A vehicle to conceptually synthesize complex decision hypothesis
  - stage-wise projection and aggregation
- A vehicle for organizing computational operations
  - stage-wise update of latent states
- A vehicle for designing processing steps and computing modules
  - Layer-wise parallelization
- No obvious utility in evaluating DL inference algorithms

Utility of the Loss Function

- Global loss? Well it is complex and non-convex...

# Graphical models vs. Deep nets

## Graphical models

Utility of the graph

- A vehicle for synthesizing a global loss function from local structure
  - potential function, feature function, etc.
- A vehicle for designing sound and efficient inference algorithms
  - Sum-product, mean-field, etc.
- A vehicle to inspire approximation and penalization
  - Structured MF, Tree-approximation, etc.
- A vehicle for monitoring theoretical and empirical behavior and accuracy of inference

Utility of the loss function

- A major measure of quality of the learning algorithm and the model

## Deep neural networks

Utility of the network

- A vehicle to conceptually synthesize complex decision hypothesis
  - stage-wise projection and aggregation
- A vehicle for organizing computational operations
  - stage-wise update of latent states
- A vehicle for designing processing steps and computing modules
  - Layer-wise parallelization
- No obvious utility in evaluating DL inference algorithms

Utility of the Loss Function

- Global loss? Well it is complex and non-convex...

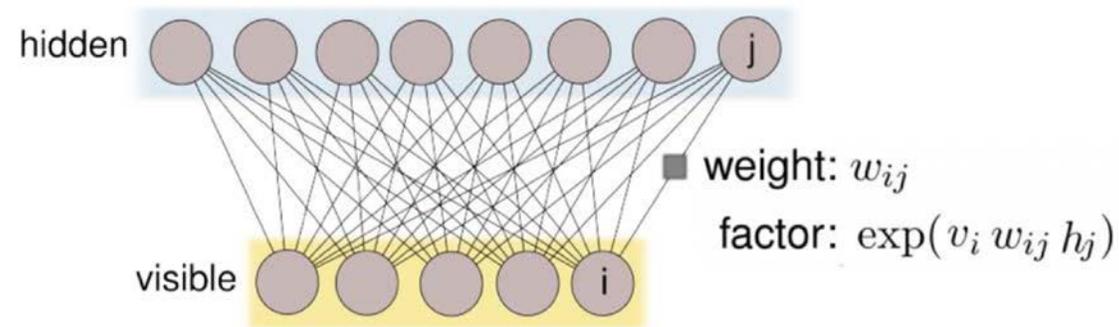| | DL | ≦ ? | ML (e.g., GM) |
|---|---|---|---|
| Empirical goal: | e.g., classification, feature learning | | e.g., latent variable inference, transfer learning |
| Structure: | Graphical | | Graphical |
| Objective: | Something aggregated from local functions | | Something aggregated from local functions |
| | | | |
| Vocabulary: | Neuron, activation function, … | | Variable, potential function, … |
| Algorithm: | A single, unchallenged, inference algorithm – Backpropagation (BP) | | A major focus of open research, many algorithms, and more to come |
| Evaluation: | On a black-box score – end performance | | On almost every intermediate quantity |
| Implementation: | Many tricks | | More or less standardized |
| Experiments: | Massive, real data (GT unknown) | | Modest, often simulated data (GT known) |

# Graphical Models vs. Deep Nets

❑ So far:
   ❑ Graphical models are representations of probability distributions
   ❑ Neural networks are function approximators (with no probabilistic meaning)

❑ Some of the neural nets are in fact proper graphical models (i.e., units/neurons represent random variables):
   ❑ Boltzmann machines (Hinton & Sejnowsky, 1983)
   ❑ Restricted Boltzmann machines (Smolensky, 1986)
   ❑ Learning and Inference in sigmoid belief networks (Neal, 1992)
   ❑ Fast learning in deep belief networks (Hinton, Osindero, Teh, 2006)
   ❑ Deep Boltzmann machines (Salakhutdinov and Hinton, 2009)

❑ Let's go through these models one-by-one

# I: Restricted Boltzmann Machines

- ❑ RBM is a Markov random field represented with a bi-partite graph
- ❑ All nodes in one layer/part of the graph are connected to all in the other; no inter-layer connections



hidden

weight: $w_{ij}$

factor: $\exp(v_i \, w_{ij} \, h_j)$

visible

- ❑ Joint distribution:

$$P(v,h) = \frac{1}{Z} \exp \left\{ \sum_{i,j} w_{ij} v_i h_i + \sum_i b_i v_i + \sum_j c_j h_j \right\}$$

# I: Restricted Boltzmann Machines

❏ Log-likelihood of a single data point (unobservables marginalized out):

$$\log L(v) = \log \sum_h \exp \left\{ \sum_{i,j} w_{ij} v_i h_i + \sum_i b_i v_i + \sum_j c_j h_j - \log(Z) \right\}$$

❏ Gradient of the log-likelihood w.r.t. the model parameters:

$$\frac{\partial}{\partial w_{ij}} \log L(v) = \sum_h P(h|v) \frac{\partial}{\partial w_{ij}} P(v,h) - \sum_{v,h} P(v,h) \frac{\partial}{\partial w_{ij}} P(v,h)$$

❏ where we have averaging over the posterior and over the joint.

# I: Restricted Boltzmann Machines

❑ Gradient of the log-likelihood w.r.t. the parameters (alternative form):

$$\frac{\partial}{\partial w_{ij}} \log L(v) = \mathbb{E}_{P(h|v)}\left[\frac{\partial}{\partial w_{ij}} P(v,h)\right] - \mathbb{E}_{P(v,h)}\left[\frac{\partial}{\partial w_{ij}} P(v,h)\right]$$

❑ Both expectations can be approximated via sampling
❑ Sampling from the posterior is exact (RBM factorizes over $h$ given $v$)
❑ Sampling from the joint is done via MCMC (e.g., Gibbs sampling)
❑ In the neural networks literature:
  ❑ computing the first term is called **the clamped / wake / positive phase**
    (the network is "awake" since it conditions on the visible variables)
  ❑ Computing the second term is called **the unclamped / sleep / free / negative phase**
    (the network is "asleep" since it samples the visible variables from the joint;
    metaphorically, it is "dreaming" the visible inputs)

# I: Restricted Boltzmann Machines

- Gradient of the log-likelihood w.r.t. the parameters (alternative form):
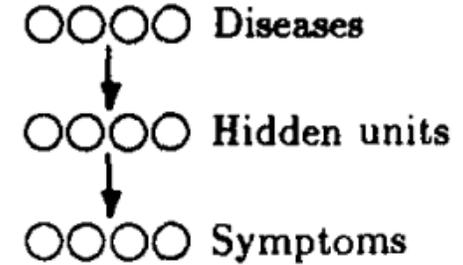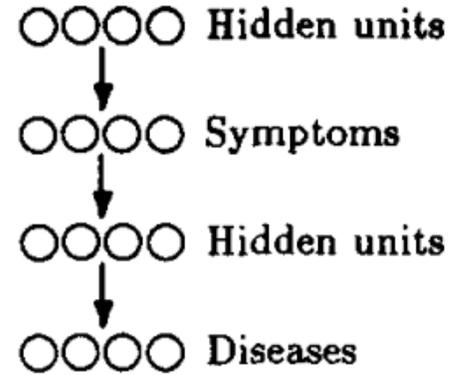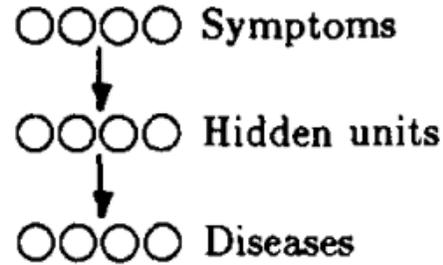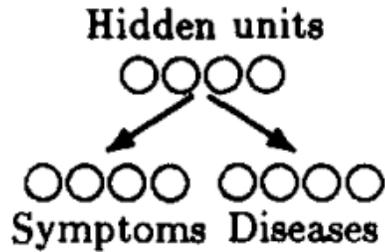
$$\frac{\partial}{\partial w_{ij}} \log L(v) = \mathbb{E}_{P(h|v)}\left[\frac{\partial}{\partial w_{ij}} P(v, h)\right] - \mathbb{E}_{P(v,h)}\left[\frac{\partial}{\partial w_{ij}} P(v, h)\right]$$

- Learning is done by optimizing the log-likelihood of the model for a given data via stochastic gradient descent (SGD)

- Estimation of the second term (the negative phase) heavily relies on the mixing properties of the Markov chain

- This often causes slow convergence and requires extra computation

# II: Sigmoid Belief Networks



from Neal, 1992

- Sigmoid belief nets are simply Bayesian networks over binary variables with conditional probabilities represented by sigmoid functions:

$$P(x_i | \pi(x_i)) = \sigma \left( x_i \sum_{x_j \in \pi(x_i)} w_{ij} x_j \right)$$

- Bayesian networks exhibit a phenomenon called "explain away effect"



**If A correlates with C, then the chance of B correlating with C decreases. ⇒ A and B become correlated given C.**
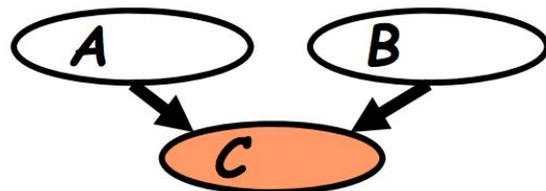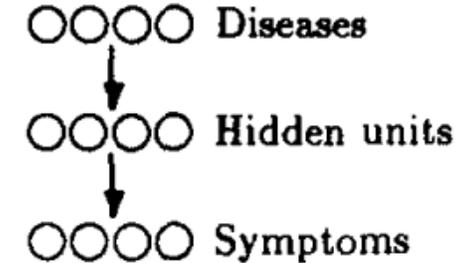
# II: Sigmoid Belief Networks



from Neal, 1992

- Sigmoid belief nets are simply Bayesian networks over binary variables with conditional probabilities represented by sigmoid functions:

$$P(x_i|\pi(x_i)) = \sigma\left(x_i \sum_{x_j \in \pi(x_i)} w_{ij}x_j\right)$$

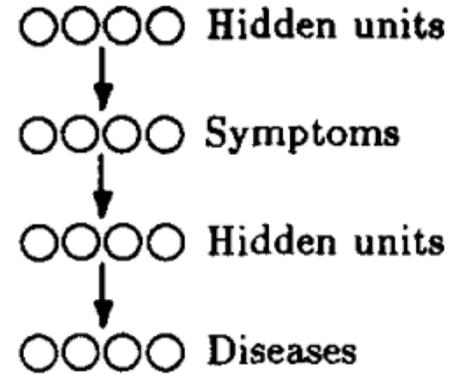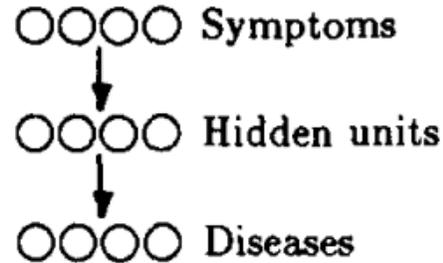- Bayesian networks exhibit a phenomenon called "explain away effect"
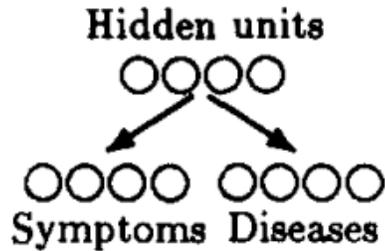


**Note:**
**Due to the "explain away effect," when we condition on the visible layer in belief networks, hidden variables all become dependent.**

# Sigmoid Belief Networks: Learning and Inference

- Neal proposed Monte Carlo methods for learning and inference (Neal, 1992):

  Approximated with Gibbs sampling

- Conditional distributions:

$$P(S_i = x \mid S_j = s_j : j \neq i)$$
$$\propto \ \sigma\left(x^* \sum_{j<i} s_j w_{ij}\right) \prod_{j>i} \sigma\left(s_j^*\left(x w_{ji} + \sum_{k<j,k\neq i} s_k w_{jk}\right)\right)$$

- **No negative phase as in RBM!**

- **Convergence is very slow, especially for large belief nets, due to the intricate "explain-away" effects…**

Equations from Neal, 1992

$$\frac{\partial L}{\partial w_{ij}} = \sum_{\tilde{v}\in\mathcal{T}} \frac{1}{P(\tilde{V}=\tilde{v})} \frac{\partial P(\tilde{V}=\tilde{v})}{\partial w_{ij}}$$

log derivative

$$= \sum_{\tilde{v}\in\mathcal{T}} \frac{1}{P(\tilde{V}=\tilde{v})} \sum_{\tilde{h}} \frac{\partial P(\tilde{S}=\langle\tilde{h},\tilde{v}\rangle)}{\partial w_{ij}}$$

prob. of the visibles via marginalization

$$= \sum_{\tilde{v}\in\mathcal{T}} \sum_{\tilde{h}} P(\tilde{S}=\langle\tilde{h},\tilde{v}\rangle \mid \tilde{V}=\tilde{v})$$
$$\cdot \frac{1}{P(\tilde{S}=\langle\tilde{h},\tilde{v}\rangle)} \frac{\partial P(\tilde{S}=\langle\tilde{h},\tilde{v}\rangle)}{\partial w_{ij}}$$

Bayes rule + rearrange sums

$$= \sum_{\tilde{v}\in\mathcal{T}} \sum_{\tilde{s}} P(\tilde{S}=\tilde{s} \mid \tilde{V}=\tilde{v}) \frac{1}{P(\tilde{S}=\tilde{s})} \frac{\partial P(\tilde{S}=\tilde{s})}{\partial w_{ij}}$$

$$= \sum_{\tilde{v}\in\mathcal{T}} \sum_{\tilde{s}} P(\tilde{S}=\tilde{s} \mid \tilde{V}=\tilde{v})$$
$$\cdot \frac{1}{\sigma\left(s_i^* \sum_{k<i} s_k w_{ik}\right)} \frac{\partial \sigma\left(s_i^* \sum_{k<i} s_k w_{ik}\right)}{\partial w_{ij}}$$

Plug-in the actual sigmoid form of the conditional prob.

$$= \sum_{\tilde{v}\in\mathcal{T}} \sum_{\tilde{s}} P(\tilde{S}=\tilde{s} \mid \tilde{V}=\tilde{v}) \, s_i^* s_j \sigma\left(-s_i^* \sum_{k<i} s_k w_{ik}\right).$$

# RBMs are infinite belief networks

- Recall the expression for the gradient of the log likelihood for RBM:

$$\frac{\partial}{\partial w_{ij}} \log L(v) = \mathbb{E}_{P(h|v)} \left[ \frac{\partial}{\partial w_{ij}} P(v, h) \right] - \mathbb{E}_{P(v,h)} \left[ \frac{\partial}{\partial w_{ij}} P(v, h) \right]$$

- To make a gradient update of the model parameters, we need compute the expectations via sampling.
  - We can sample exactly from the posterior in the first term
  - We run block Gibbs sampling to approximately sample from the joint distribution

sampling steps

# RBMs are infinite belief networks

- Gibbs sampling: alternate between sampling hidden and visible variables



- Conditional distributions $P(v|h)$ and $P(h|v)$ are represented by sigmoids
- Thus, we can think of Gibbs sampling from the joint distribution represented by an RBM as a top-down propagation in an infinitely deep sigmoid belief network!

# RBMs are infinite belief networks

- RBMs are equivalent to infinitely deep belief networks



- Sampling from this is the same as sampling from the network on the right

# RBMs are infinite belief networks

- RBMs are equivalent to infinitely deep belief networks

# RBMs are infinite belief networks

- RBMs are equivalent to infinitely deep belief networks



- When we train an RBM, we are really training an infinitely deep brief net!

- It is just that the weights of all layers are tied.

- If the weights are "untied" to some extent, we get a Deep Belief Network.

# III: Deep Belief Nets

**Deep Belief Network**



$$h^3 \quad W^3$$
$$h^2 \quad W^2$$
$$h^1 \quad W^1$$
$$v$$

| Now weights are untied! |

- DBNs are hybrid graphical models (chain graphs):
  - Exact inference in DBNs is problematic due to <u>explaining away effect</u>
  - Training: greedy pre-training + ad-hoc fine-tuning; no proper joint training
  - Approximate inference is feed-forward

# Deep Belief Networks

**Deep Belief Network**



- DBNs represent a joint probability distribution

$$P(v, h^1, h^2, h^3) = P(h^2, h^3)P(h^1|h^2)P(v|h^1)$$

- Note that $P(h^2, h^3)$ is an RBM and the conditionals $P(h^1|h^2)$ and $P(v|h^1)$ are represented in the sigmoid form
- The model is trained by optimizing the log likelihood for a given data $\log P(v)$

## Challenges:

- Exact inference in DBNs is problematic due to <u>explain away effect</u>
- Training is done in two stages:
  - greedy pre-training + ad-hoc fine-tuning; no proper joint training
- Approximate inference is feed-forward (bottom-up)

# DBN: Layer-wise pre-training

- Pre-train and freeze the 1st RBM
- Stack another RBM on top and train it



- The weights weights 2+ layers remain tied
- We repeat this procedure: pre-train and untie the weights layer-by-layer…

# DBN: Layer-wise pre-training

- We repeat this procedure: pre-train and untie the weights layer-by-layer:

- The weights of 3+ layers remain tied



- and so forth

- *From the optimization perspective, this procedure loosely corresponds to an approximate block-coordinate accent on the log-likelihood*

**images from Marcus Frean, MLSS Tutorial 2010**

# DBN: Fine-tuning

- Pre-training is quite ad-hoc and is unlikely to lead to a good probabilistic model *per se*

- However, the layers of representations could perhaps be useful for some other downstream tasks!

- We can further "fine-tune" a pre-trained DBN for some other task

Setting A: Unsupervised learning (DBN → autoencoder)

1. Pre-train a stack of RBMs in a greedy layer-wise fashion

2. "Unroll" the RBMs to create an autoencoder

3. Fine-tune the parameters by optimizing the reconstruction error



**images from Hinton & Salakhutdinov, 2006**

# DBN: Fine-tuning

- Pre-training is quite ad-hoc and is unlikely to lead to a good probabilistic model *per se*

- However, the layers of representations could perhaps be useful for some other downstream tasks!

- We can further "fine-tune" a pre-trained DBN for some other task

Setting A: Unsupervised learning (DBN → autoencoder)

1. Pre-train a stack of RBMs in a greedy layer-wise fashion

2. **"Unroll" the RBMs to create an autoencoder**

3. Fine-tune the parameters by optimizing the reconstruction error



**images from Hinton & Salakhutdinov, 2006**

# DBN: Fine-tuning

- Pre-training is quite ad-hoc and is unlikely to lead to a good probabilistic model *per se*

- However, the layers of representations could perhaps be useful for some other downstream tasks!

- We can further "fine-tune" a pre-trained DBN for some other task

Setting A: Unsupervised learning (DBN → autoencoder)

1. Pre-train a stack of RBMs in a greedy layer-wise fashion

2. "Unroll" the RBMs to create an autoencoder

3. Fine-tune the parameters by optimizing the reconstruction error



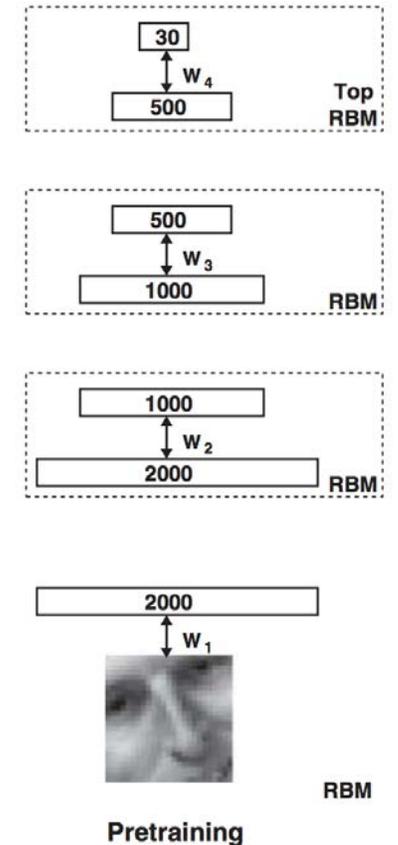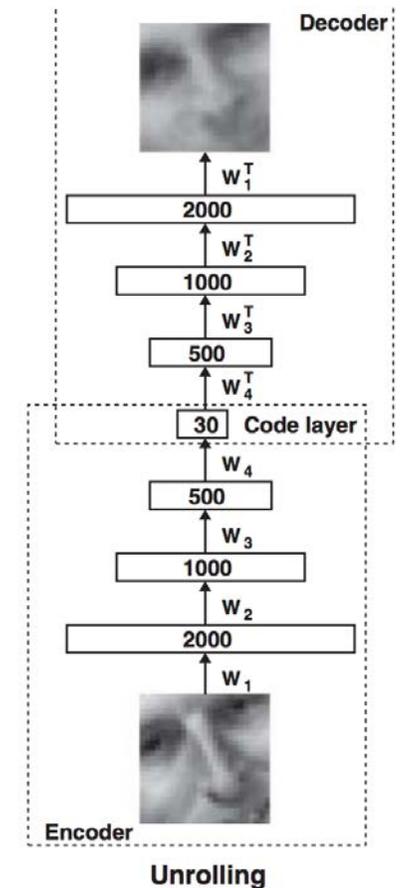**images from Hinton & Salakhutdinov, 2006**

# DBN: Fine-tuning

- Pre-training is quite ad-hoc and is unlikely to lead to a good probabilistic model *per se*

- However, the layers of representations could perhaps be useful for some other downstream tasks!

- We can further "fine-tune" a pre-trained DBN for some other task

Setting B: Supervised learning (DBN → classifier)

1. Pre-train a stack of RBMs in a greedy layer-wise fashion

2. "Unroll" the RBMs to create a feedforward classifier

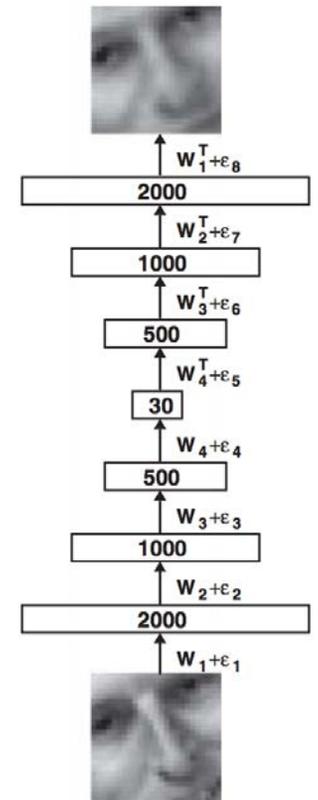3. Fine-tune the parameters by optimizing the reconstruction error

**Some intuitions about how pre-training works:**
**Erhan et al.: Why Does Unsupervised Pre-training Help Deep Learning? JMLR, 2010**

# Deep Belief Nets and Boltzmann Machines

## Deep Belief Network



- DBNs are hybrid graphical models (chain graphs):
  - Inference in DBNs is problematic due to <u>explaining away effect</u>
  - Training: greedy pre-training + ad-hoc fine-tuning; no proper joint training
  - Approximate inference is feed-forward

# Deep Belief Nets and Boltzmann Machines



Deep Belief Network

Deep Boltzmann Machine

- DBMs are fully un-directed models (Markov random fields):
  - Can be trained similarly as RBMs via MCMC (Hinton & Sejnowski, 1983)
  - Use a variational approximation of the data distribution for faster training (Salakhutdinov & Hinton, 2009)
  - Similarly, can be used to initialize other networks for downstream tasks

# Graphical models vs. Deep networks

❑ A few critical points to note about all these models:

   ❑ The primary goal of deep generative models is to represent the distribution of the observable variables. Adding layers of hidden variables allows to represent increasingly more complex distributions.

   ❑ Hidden variables are secondary (auxiliary) elements used to facilitate learning of complex dependencies between the observables.

   ❑ Training of the model is ad-hoc, but what matters is the quality of learned hidden representations.

   ❑ Representations are judged by their usefulness on a downstream task (the probabilistic meaning of the model is often discarded at the end).

❑ In contrast, classical graphical models are often concerned with the correctness of learning and inference of all variables

# An old study of belief networks from the GM standpoint

**[Xing, Russell, Jordan, UAI 2003]**

Mean-field partitions of a sigmoid belief network for subsequent GMF inference



Study focused on only inference/learning accuracy, speed, and partition

# "Optimize" how to optimize via truncation & re-opt

- ❏ Energy-based modeling of the structured output (CRF)

$$\mathbf{y}^*(\mathbf{x}; \mathbf{w}) := \arg\min_{\mathbf{y}} E(\mathbf{y}, \mathbf{x}; \mathbf{w})$$

- ❏ Unroll the optimization algorithm for a fixed number of steps (Domke, 2012)

$$\mathbf{y}^*(\mathbf{x}; \mathbf{w}) = \text{opt-alg}_{\mathbf{y}} E(\mathbf{y}, \mathbf{x}; \mathbf{w})$$



**We can backprop through the optimization steps since they are just a sequence of computations**

**Relevant recent paper:**
**Anrychowicz et al.: *Learning to learn by gradient descent by gradient descent*. 2016.**

# Dealing with structured prediction

❏ Energy-based modeling of the structured output (CRF)

$$\mathbf{y}^*(\mathbf{x}; \mathbf{w}) := \arg\min_{\mathbf{y}} E(\mathbf{y}, \mathbf{x}; \mathbf{w})$$

❏ Unroll the optimization algorithm for a fixed number of steps (Domke, 2012)

$$\mathbf{y}^*(\mathbf{x}; \mathbf{w}) = \text{opt-alg}_{\mathbf{y}} E(\mathbf{y}, \mathbf{x}; \mathbf{w})$$

❏ We can think of y* as some non-linear differentiable function of the inputs and weights → impose some loss and optimize it as any other standard computation graph using backprop!

❏ Similarly, message passing based inference algorithms can be truncated and converted into computational graphs (Domke, 2011; Stoyanov et al., 2011)

# Conclusion

- DL & GM: the fields are similar in the beginning (structure, energy, etc.), and then diverge to their own signature pipelines

- DL: most effort is directed to comparing different architectures and their components (models are driven by evaluating empirical performance on a downstream tasks)
  - DL models are good at learning robust hierarchical representations from the data and suitable for simple reasoning (call it "low-level cognition")

- GM: the effort is directed towards improving inference accuracy and convergence speed
  - GMs are best for provably correct inference and suitable for high-level complex reasoning tasks (call it "high-level cognition")

- Convergence of both fields is very promising!
  - Next part: a unified view of deep generative models in the GM interpretation

# Supplementary

# Outline

- An overview of DL components
    - Historical remarks: early days of neural networks
    - Modern building blocks: units, layers, activations functions, loss functions, etc.
    - Reverse-mode automatic differentiation (aka backpropagation)


- Similarities and differences between GMs and NNs
    - Graphical models vs. computational graphs
    - Sigmoid Belief Networks as graphical models
    - Deep Belief Networks and Boltzmann Machines


- Combining DL methods and GMs
    - Using outputs of NNs as inputs to GMs
    - GMs with potential functions represented by NNs
    - NNs with structured outputs


- Bayesian Learning of NNs
    - Bayesian learning of NN parameters
    - Deep kernel learning

# Combining sequential NNs and GMs



Hybrid: RNN + HMM

# Combining sequential NNs and GMs

## Hybrid: RNN + HMM

(Graves et al., 2013)

The model, inference, and learning can be **analogous** to our NN + HMM hybrid

- **Objective:** log-likelihood
- **Model:** HMM/Gaussian emissions
- **Inference:** forward-backward algorithm
- **Learning:** SGD with gradient by backpropagation

# Hybrid NNs + conditional GMs



- ❑ In a standard CRF, each of the factor cells is a parameter.
- ❑ In a hybrid model, these values are computed by a neural network.

# Hybrid NNs + conditional GMs



Hybrid: Neural Net + CRF

slide courtesy: Matt Gormley

# Hybrid NNs + conditional GMs



(Collobert & Weston, 2011)

## Hybrid: CNN + CRF

"NN + SLL"

- Model: Convolutional Neural Network (CNN) with **linear-chain CRF**
- Training objective: maximize **sentence-level likelihood** (SLL)

Figure from (Collobert & Weston, 2011)

slide courtesy: Matt Gormley

# Using GMs as Prediction Explanations

## Satellite imagery



## Meaningful attributes

| | | | |
|---|---|---|---|
| 01 | Nightlight intensity | 09 | Avg. vegetation dec. |
| 02 | Is urban | 10 | Avg. dist. to market |
| 03 | Has electricity | 11 | Avg. dist. to road |
| 04 | Has generator | 12 | Num. of rooms |
| 05 | Avg. temperature | 13 | Dist. to water src. |
| 06 | Avg. percipitation | 14 | Water usage p/ day |
| 07 | Vegetation | 15 | Is water payed |
| 08 | Avg. vegetation inc. | 16 | HH type: BQ |

Al-Shedivat, Dubey, Xing, arXiv, 2017

# Using GMs as Prediction Explanations

Satellite imagery

Area attributes

01 Nightlight intensity     09 Avg. vegetation dec.

02 Is urban     10 Avg. dist. to market

**How do we build a powerful predictive model whose predictions we can interpret in terms of semantically meaningful features?**

05 Avg. temperature     13 Dist. to water src.

06 Avg. percipitation     14 Water usage p/ day

07 Vegetation     15 Is water payed

08 Avg. vegetation inc.     16 HH type: BQ

Al-Shedivat, Dubey, Xing, arXiv, 2017

# Contextual Explanation Networks (CENs)



Al-Shedivat, Dubey, Xing, arXiv, 2017

**Context**      **Explanation**      **Prediction**

Instance 1

Instance 2

Encoder

Arbitrary neural net

| 0.9 | -0.6 | 0.3 | 0.5 | -0.2 |

| 0 | 1 | 1 | 0 | 1 |

| Unreliable water | Water src: Public tap | Walls: Unburnt bricks | Roof: Thatch, Straw | Has electricity |

Attributes

| 0 | 0 | 0 | 1 | 0 |

| -0.4 | -1.2 | -0.2 | 0.2 | -0.8 |

Not poor

Poor

Parameters: a graphical model

**Notes:**

- The final prediction is made by a linear GM.

- Each coefficient assigns a weight to a meaningful attribute.

- Allows us to judge predictions in terms of GMs produced by the context encoder.

# Contextual Explanation Networks (CENs)



Medical notes (context):
‣ regular smoker
‣ family history of diabetes
‣ [no previous heart attacks]

Diabetes    Smoking
$\theta$    Arythmia
High blood pressure    Previous heart attack

(a)

- <u>General idea</u>: Use deep neural nets to **generate parameters** for graphical models applicable in a given context (e.g., for a given patient).

- Produced GMs are used to make the final prediction ⇒ 100% fidelity and consistency.

- GMs are built on top of semantically meaningful variables (not deep embeddings!) and can be used as **explanations** for each prediction.

# CEN: Implementation Details



Workflow:

- Maintain a (sparse) dictionary of GM parameters.

- Process complex inputs (images, text, time series, etc.) using deep nets; use soft attention to either select or combine models from the dictionary.

- Use constructed GMs (e.g., CRFs) to make predictions.

- Inspect GM parameters to understand the reasoning behind predictions.

Al-Shedivat, Dubey, Xing, arXiv, 2017

# Results: imagery as context



(a)

(b)

(c) Uganda: Contextual Models

**Based on the imagery, CEN learns to select different models for urban and rural areas**

Al-Shedivat, Dubey, Xing, arXiv, 2017

# Results: classical image & text datasets

| MNIST | | CIFAR10 | |
|---|---|---|---|
| **Model** | **Err (%)** | **Model** | **Err (%)** |
| $LR_{pxl}$ | 8.00 | $LR_{pxl}$ | 60.1 |
| $LR_{hog}$ | 2.98 | $LR_{hog}$ | 48.6 |
| CNN | **0.75** | VGG | 9.4 |
| $MoE_{pxl}$ | 1.23 | $MoE_{pxl}$ | 13.0 |
| $MoE_{hog}$ | 1.10 | $MoE_{hog}$ | 11.7 |
| $CEN_{pxl}$ | **0.76** | $CEN_{pxl}$ | 9.6 |
| $CEN_{hog}$ | **0.73** | $CEN_{hog}$ | **9.2** |



**Same performance as vanilla deep networks; no compute overhead.**

**Predicting via explanation regularizes the model when there is not enough data.**

Al-Shedivat, Dubey, Xing, arXiv, 2017

# Results: classical image & text datasets

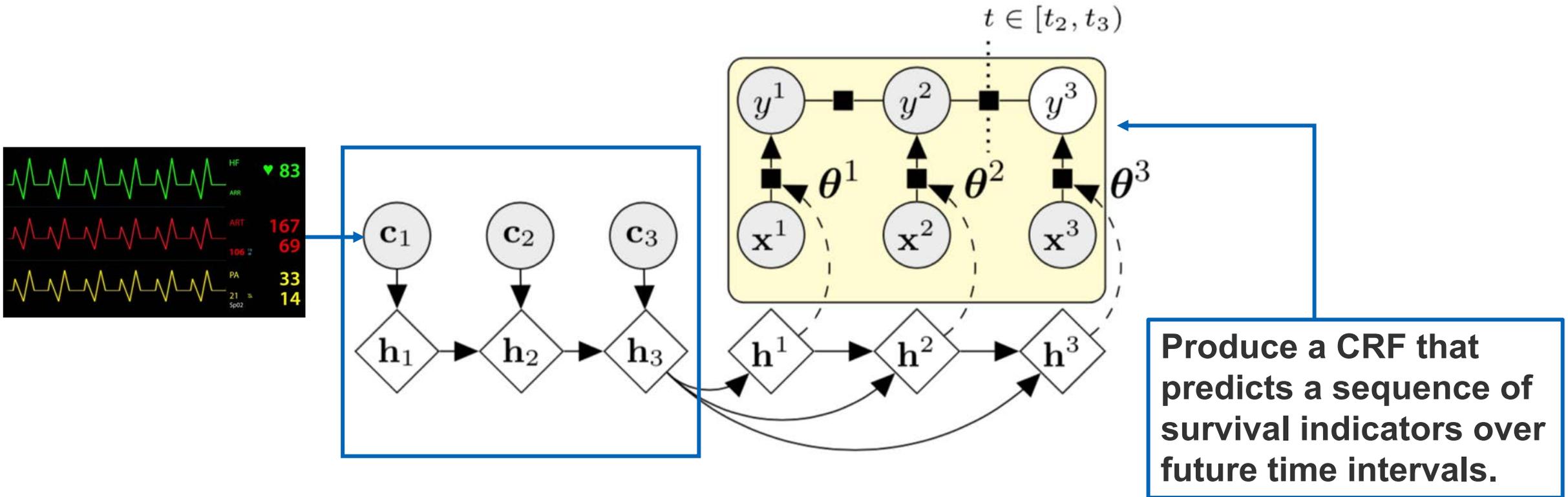| Method | Error |
|---|---|
| Paragraph Vector (Le and Mikolov, 2014) | 7.42% |
| SA-LSTM with joint training (Dai and Le, 2015) | 14.70% |
| LSTM with tuning and dropout (Dai and Le, 2015) | 13.50% |
| LSTM initialized with word2vec embeddings (Dai and Le, 2015) | 10.00% |
| SA-LSTM with linear gain (Dai and Le, 2015) | 9.17% |
| LM-TM (Dai and Le, 2015) | 7.64% |
| SA-LSTM (Dai and Le, 2015) | 7.24% |
| Virtual Adversarial (Miyato et al., 2016) | $\mathbf{5.94 \pm 0.12\%}$ |
| TopicRNN (Dieng et al., 2017) | 6.28 % |
| **CEN-bow** | $\mathbf{5.92 \pm 0.05}$ % |
| **CEN-topic** | $\mathbf{6.25 \pm 0.09}$ % |

**Semi-supervised training**

**Only supervised training (!!)**

# CEN architectures for survival analysis



Encode a sequence of observations for a patient (e.g., vitals/tests measured in ICU).

Produce a CRF that predicts a sequence of survival indicators over future time intervals.

# Results: survival analysis

| | SUPPORT2 | | | | | PhysioNet Challenge 2012 | | | |
|---|---|---|---|---|---|---|---|---|---|
| Model | Acc@25 | Acc@50 | Acc@75 | RAE | Model | Acc@25 | Acc@50 | Acc@75 | RAE |
| Cox | 84.1 | 73.7 | 47.6 | 0.90 | Cox | 93.0 | 69.6 | 49.1 | 0.24 |
| Aalen | 87.1 | 66.2 | 45.8 | 0.98 | Aalen | 93.3 | 78.7 | 57.1 | 0.31 |
| CRF | 84.4 | 89.3 | 79.2 | 0.59 | CRF | 93.2 | 85.1 | 65.6 | 0.14 |
| MLP-CRF | **87.7** | 89.6 | 80.1 | 0.62 | LSTM-CRF | 93.9 | 86.3 | 68.1 | **0.11** |
| MLP-CEN | 85.5 | **90.8** | **81.9** | **0.56** | LSTM-CEN | **94.8** | **87.5** | **70.1** | **0.09** |



Patient ID: 3520 (Died)    Patient ID: 1100 (Survived)

# Outline

- An overview of the DL components
  - Historical remarks: early days of neural networks
  - Modern building blocks: units, layers, activations functions, loss functions, etc.
  - Reverse-mode automatic differentiation (aka backpropagation)
- Similarities and differences between GMs and NNs
  - Graphical models vs. computational graphs
  - Sigmoid Belief Networks as graphical models
  - Deep Belief Networks and Boltzmann Machines
- Combining DL methods and GMs
  - Using outputs of NNs as inputs to GMs
  - GMs with potential functions represented by NNs
  - NNs with structured outputs
- Bayesian Learning of NNs
  - Bayesian learning of NN parameters
  - Deep kernel learning

# Bayesian learning of NNs

❑ A neural network as a probabilistic model:

  ❑ Likelihood: $p(y|x, \boldsymbol{\theta})$

    ❑ Categorical distribution for classification ⇒ cross-entropy loss

    ❑ Gaussian distribution for regression ⇒ squared loss

  ❑ Prior on parameters: $p(\boldsymbol{\theta})$

❑ Maximum a posteriori (MAP) solution:

  ❑ $\boldsymbol{\theta}^{MAP} = \text{argmax}_{\theta} \log p(y|x, \boldsymbol{\theta}) p(\boldsymbol{\theta})$

  ❑ Gaussian prior ⇒ L2 regularization

  ❑ Laplace prior ⇒ L1 regularization

❑ Bayesian learning [MacKay 1992, Neal 1996, de Freitas 2003]

  ❑ Posterior: $p(\boldsymbol{\theta}|x, y)$

  ❑ Variational inference with approximate posterior $q(\boldsymbol{\theta})$
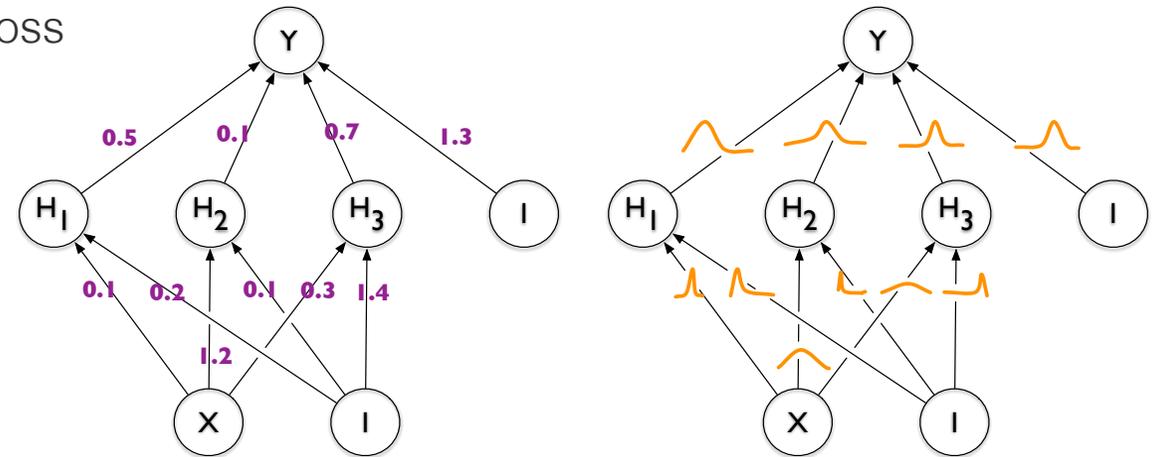


Figure courtesy: Blundell et al, 2016

# Bayesian learning of NNs

- Variational inference (in a nutshell):

$$\min_q F(D, \boldsymbol{\theta}) = \text{KL}\big(q(\boldsymbol{\theta}) \| \, p(\boldsymbol{\theta}|D)\big) - \text{E}_{q(\boldsymbol{\theta})}[\log p(D|\boldsymbol{\theta})]$$
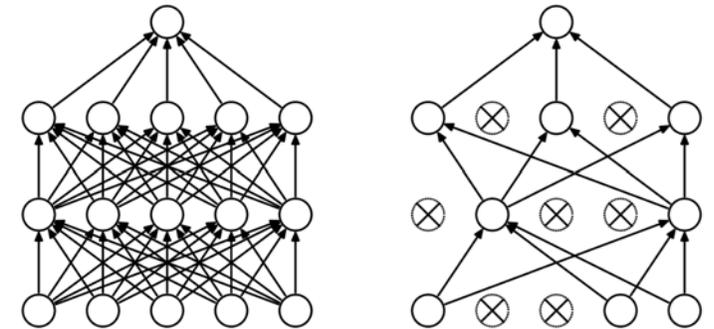
$$\min_q F(D, \boldsymbol{\theta}) = \text{KL}\big(q(\boldsymbol{\theta}) \| \, p(\boldsymbol{\theta}|D)\big) - \sum_i \log p(D|\boldsymbol{\theta}_i)$$

where $\boldsymbol{\theta}_i \sim q(\boldsymbol{\theta})$; KL term can be approximated similarly

- We can define $q(\boldsymbol{\theta})$ as a diagonal Gaussian or full-covariance Gaussian

- Alternatively, $q(\boldsymbol{\theta})$ can be defined implicitly, e.g. via dropout [Gal & Ghahramani, 2016]

$$\boldsymbol{\theta} = \boldsymbol{M} \cdot \text{diag}(\boldsymbol{z}),$$
$$\boldsymbol{z} \sim \text{Bernoulli}(p)$$



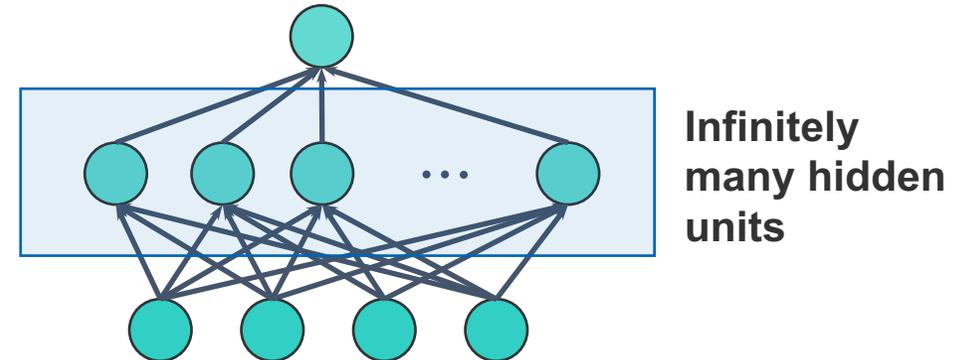- Dropping out neurons is equivalent to zeroing out columns of the parameter matrices (i.e., weights)

- $z_i = 0$ corresponds to $i$-th column of $\boldsymbol{M}$ being dropped out
  ⇒ the procedure is equivalent to dropout of unit $i$ [Hinton et al., 2012]

- Variational parameters are $\{\boldsymbol{M}, \boldsymbol{p}\}$

# "Infinitely Wide" Deep Models

- **We have seen that an "infinitely deep" network can be explained by a proper GM, How about an "infinitely wide" one?**

- Consider a neural network with a Gaussian prior on its weights an <u>infinitely</u> many hidden neurons in the intermediate layer.

- Turns out, if we have a certain Gaussian prior on the weights of such infinite network, it will be equivalent to a Gaussian process [Neal 1996].



**Infinitely many hidden units**

- Gaussian process (GP) is a distribution over functions:

$$m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})],$$
$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))],$$
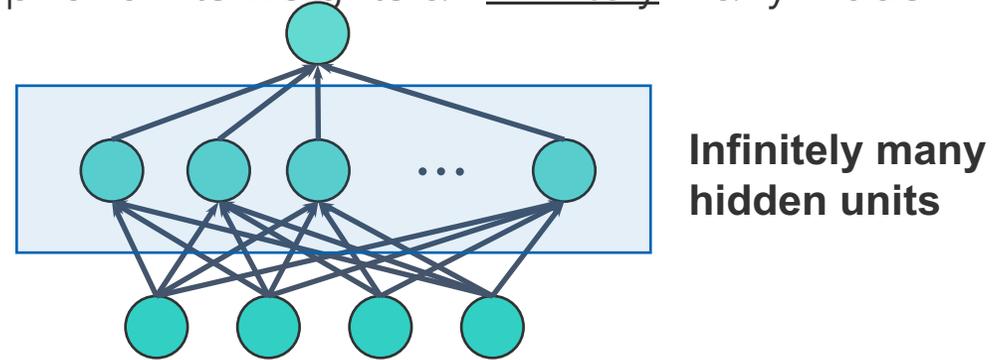$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')).$$

- When used for prediction, GPs account for correlations between the data points and can output well-calibrated predictive uncertainty estimates.

# Gaussian Process and Deep Kernel Learning

❑ Consider a neural network with a Gaussian prior on its weights an <u>infinitely</u> many hidden neurons in the intermediate layer.
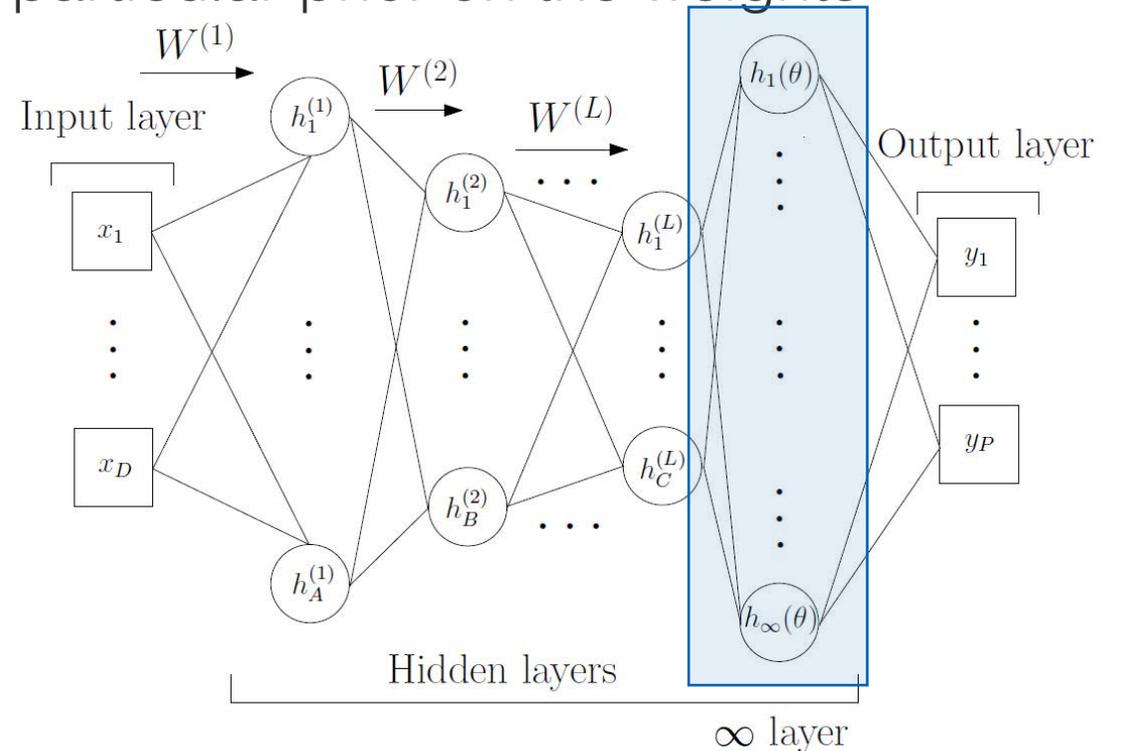


**Infinitely many hidden units**

❑ Certain classes of Gaussian priors for neural networks with infinitely many hidden units converge to Gaussian processes [Neal 1996]

❑ Deep kernel [Wilson et al., 2016]

  ❑ Combines the inductive biases of deep model architectures with the non-parametric flexibility of Gaussian processes

$$k(x_i, x_j | \phi) \rightarrow k(g(x_i, \theta), g(x_j, \theta) | \phi, \theta) \quad \text{where } K_{ij} = k(x_i, x_j)$$

  ❑ Starting from a base kernel $k(x_i, x_j | \phi)$, transform the inputs $x$ as
$$p(f|\phi) = \mathcal{N}(f|m(x), K)$$
$$p(y|f) = \mathcal{N}(y|f, \beta^{-1})$$

  ❑ Learn both kernel and neural parameters $\{\phi, \theta\}$ jointly by optimizing marginal log-likelihood (or its variational lower-bound).

  ❑ Fast learning and inference with local kernel interpolation, structured inducing points, and Monte Carlo approximations

# Gaussian Process and Deep Kernel Learning

- ❑ By adding GP as a layer to a deep neural net, we can think of it as adding an infinite hidden layer with a particular prior on the weights

- ❑ Deep kernel learning [Wilson et al., 2016]
    - ❑ Combines the inductive biases of deep models with the non-parametric flexibility of Gaussian processes
    - ❑ GPs add powerful regularization to the network
    - ❑ Additionally, they provide predictive uncertainty estimates
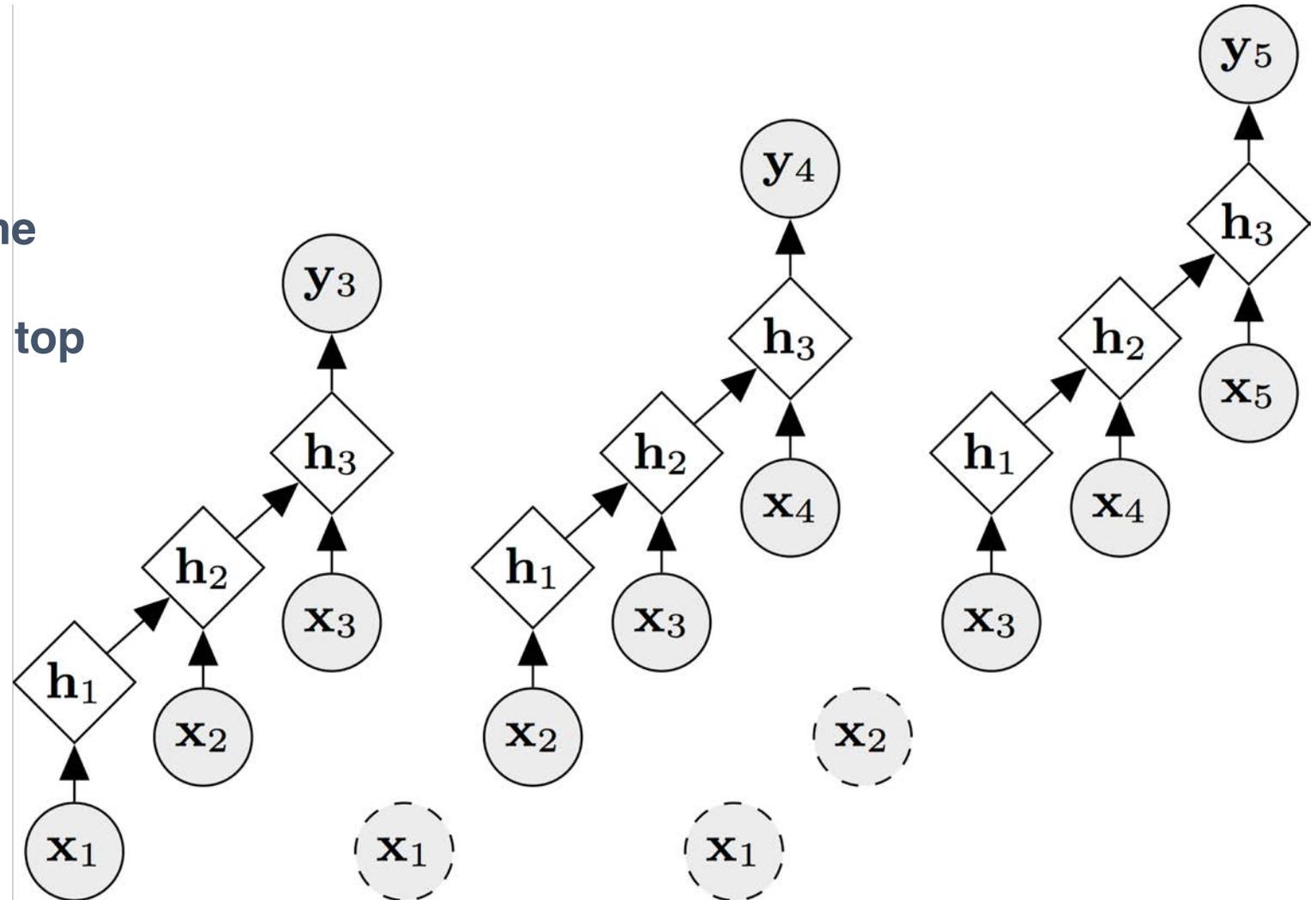
# Deep kernel learning on sequential data

What if we have data of sequential nature?

Can we still apply the same reasoning and build rich nonparametric models on top recurrent nets?
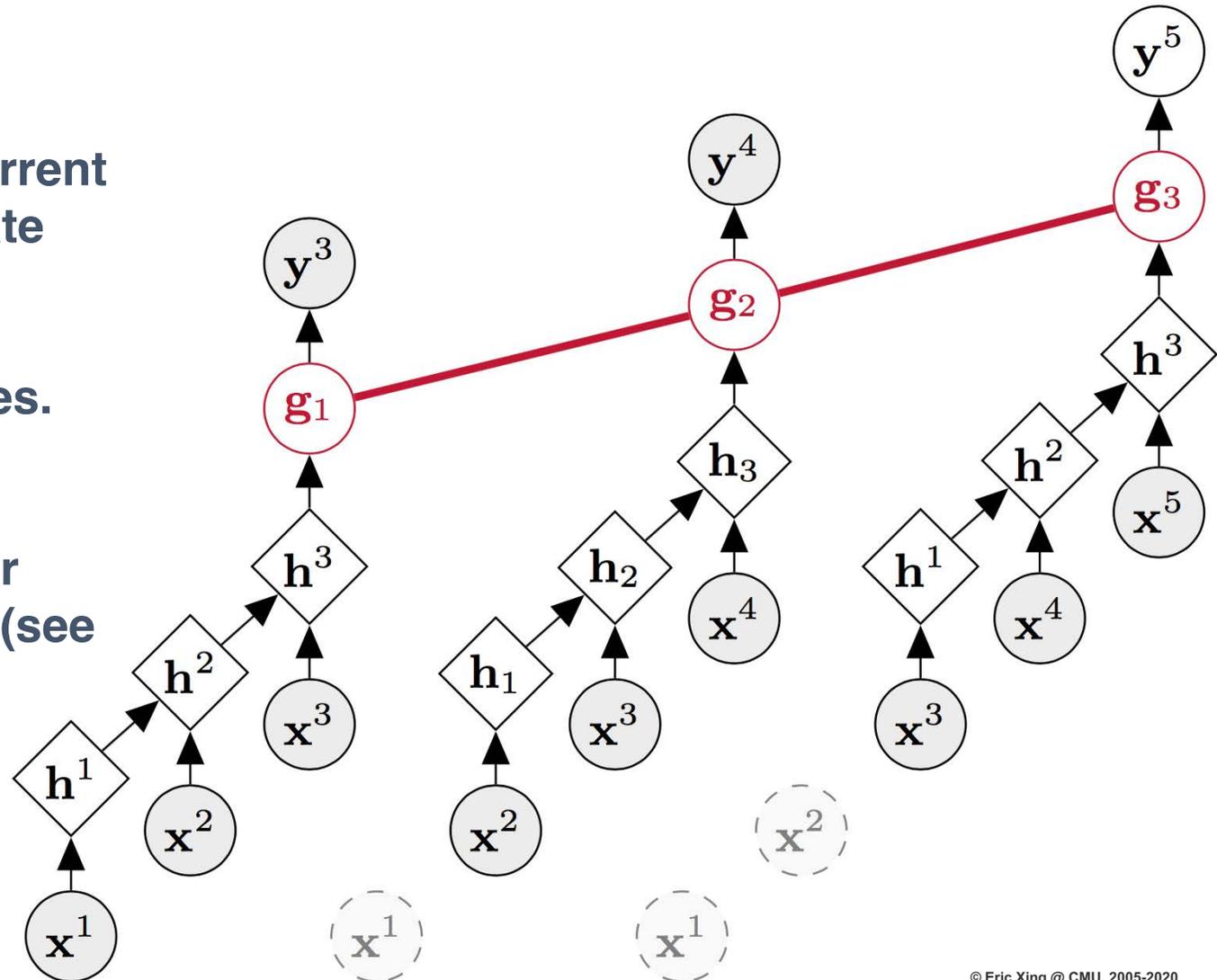
# Deep kernel learning on sequential data

The answer is **YES**!

By adding a GP layer to a recurrent network, we effectively correlate samples across time and get predictions along with well calibrated uncertainty estimates.
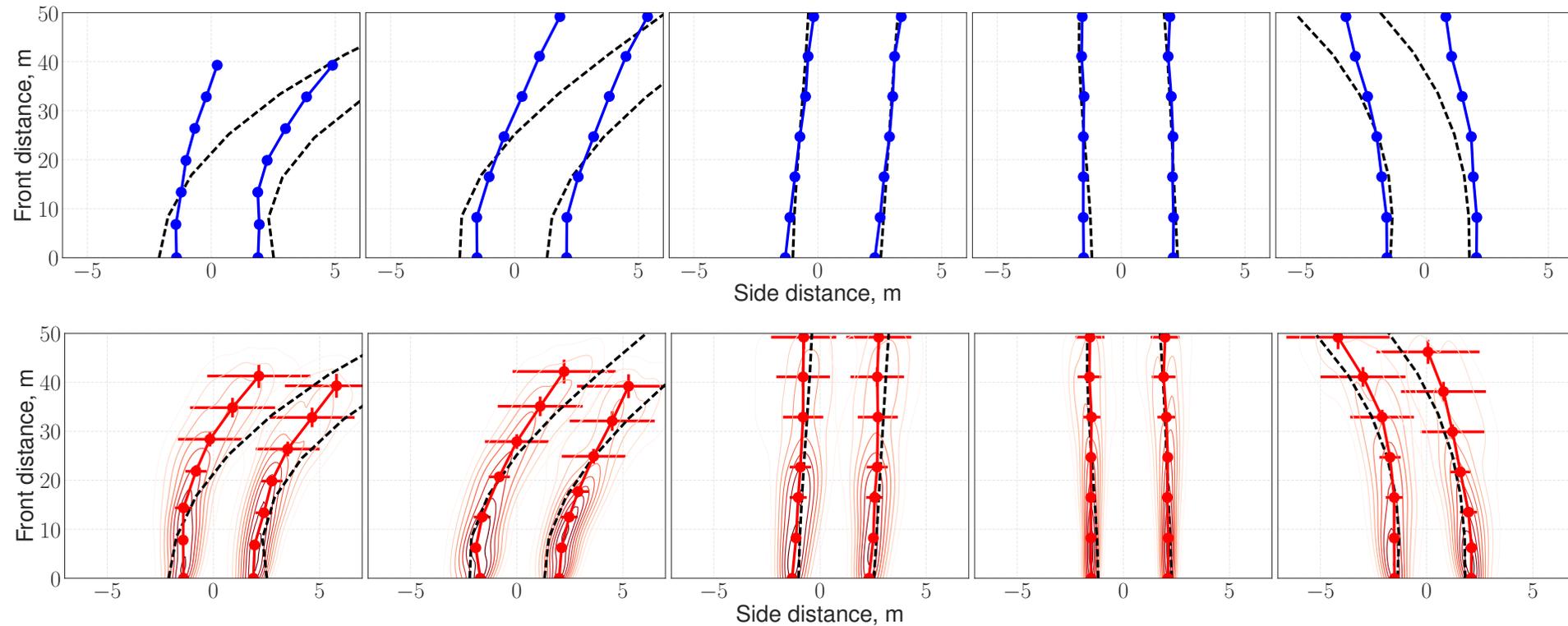
To train such model using stochastic techniques however requires some additional care (see our paper).



Al-Shedivat et al., JMLR, 2017

# Deep kernel learning on sequential data



Lane prediction: LSTM vs GP-LSTM

# Deep kernel learning on sequential data

**Lead vehicle prediction: LSTM vs GP-LSTM**