

# Latent Variable Models (contd) and Deep Neural Networks

CS771: Introduction to Machine Learning

# MLE for Latent Variable Models

Optimization is difficult in general because the objective has a complex form, complex derivatives, so FOO etc won't apply easily

- Original problem:  $\Theta_{MLE} = \operatorname{argmax}_{\Theta} \log p(\mathbf{X}|\Theta) = \operatorname{argmax}_{\Theta} \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta)$
- The EM approach for solving MLE problem in LVMs

MAP for  $\Theta$  can also be done by adding a log prior term for  $\Theta$

$$\Theta_{MLE} = \operatorname{argmax}_{\Theta} \mathbb{E}_{p(\mathbf{Z}|\Theta, \mathbf{X})} [\log p(\mathbf{X}, \mathbf{Z}|\Theta)]$$

- Can be easily shown (exercise) that original objective (log-lik) can be written as

$$\log p(\mathbf{X}|\Theta) = \mathcal{L}(q, \Theta) + KL(q||p_{\mathbf{Z}})$$

Holds for any distribution  $q(\mathbf{Z})$  over  $\mathbf{Z}$

$$\mathcal{L}(q, \Theta) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \right\}$$

KL is always non-negative

$$KL(q||p_{\mathbf{Z}}) = - \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \Theta)}{q(\mathbf{Z})} \right\}$$



$$\log p(\mathbf{X}|\Theta) \geq \mathcal{L}(q, \Theta)$$

Thus  $\mathcal{L}(q, \Theta)$  is a **lower bound** on the log likelihood of the LVM

EM maximizes this lower bound to get (an approximate) MLE



# Maximizing the Lower Bound

- As we saw,  $\mathcal{L}(q, \Theta)$  depends on  $q$  and  $\Theta$ . Consider ALT-OPT.
- Let's maximize  $\mathcal{L}(q, \Theta)$  w.r.t.  $q$  with  $\Theta$  fixed at  $\Theta^{\text{old}}$

The posterior distribution of  $\mathbf{Z}$  given older parameters  $\Theta^{\text{old}}$  (will need this posterior to get the expectation of CLL)

Since  $\log p(\mathbf{X}|\Theta) = \mathcal{L}(q, \Theta) + KL(q||p_z)$  is constant when  $\Theta$  is held fixed at  $\Theta^{\text{old}}$

$$\hat{q} = \operatorname{argmax}_q \mathcal{L}(q, \Theta^{\text{old}}) = \operatorname{argmin}_q KL(q||p_z) = p_z = p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}})$$

- Now let's maximize  $\mathcal{L}(q, \Theta)$  w.r.t.  $\Theta$  with  $q$  fixed at  $\hat{q} = p_z = p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}})$

$$\begin{aligned} \Theta^{\text{new}} &= \operatorname{argmax}_{\Theta} \mathcal{L}(\hat{q}, \Theta) = \operatorname{argmax}_{\Theta} \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}})} \right\} \\ &= \operatorname{argmax}_{\Theta} \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) \\ &= \operatorname{argmax}_{\Theta} \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}})} [\log p(\mathbf{X}, \mathbf{Z}|\Theta)] \\ &= \operatorname{argmax}_{\Theta} Q(\Theta, \Theta^{\text{old}}) \end{aligned}$$

Maximization of expected CLL w.r.t. the posterior distribution of  $\mathbf{Z}$  given older parameters  $\Theta^{\text{old}}$



# EM: An Illustration

- As we saw, EM maximizes the lower bound  $\mathcal{L}(q, \Theta)$  to get  $\Theta^{(MLE)}$  in two steps
- Step 1 sets  $\hat{q} = p(\mathbf{Z}|\Theta, \mathbf{X})$  using  $\Theta^{\text{old}}$ . KL becomes zero and  $\mathcal{L}(\hat{q}, \Theta^{\text{old}}) = \log p(\mathbf{X}|\Theta^{\text{old}})$
- Step 2 maximizes  $\mathcal{L}(\hat{q}, \Theta)$  w.r.t.  $\Theta$  which gives the new  $\Theta$ .

Alternating between them until convergence to some local optima

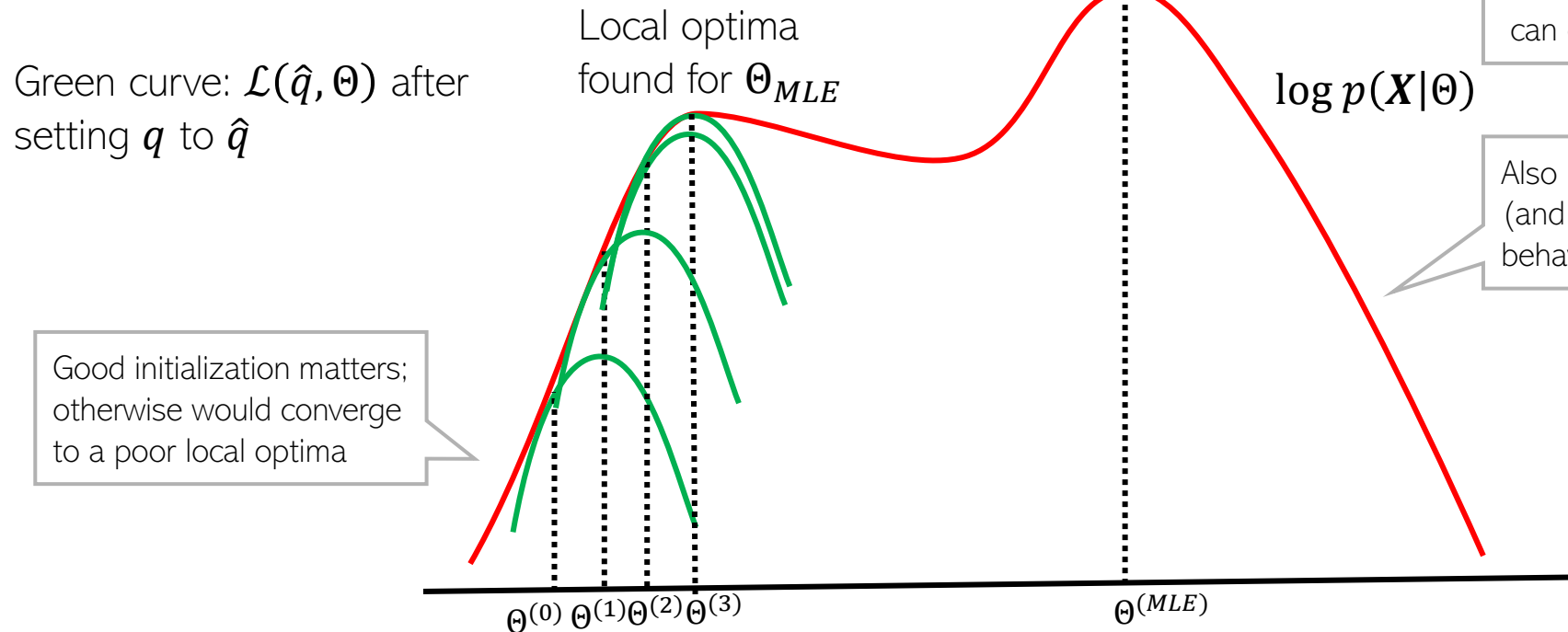
Note that  $\Theta$  only changes in Step 2 so the objective  $\log p(\mathbf{X}|\Theta)$  can only change in Step 2



Also kind of similar to Newton's method (and has second order like convergence behavior like Newton's in some cases)

Unlike Newton's method, we don't construct and optimize a quadratic approximation, but a lower bound

Even though original MLE problem  $\text{argmax}_{\Theta} \log p(\mathbf{X}|\Theta)$  could be solved using gradient methods, EM often works faster and has cleaner updates



# The EM Algorithm in its general form..

- Maximization of  $\mathcal{L}(q, \Theta)$  w.r.t.  $q$  and  $\Theta$  gives the EM algorithm (Dempster, Laird, Rubin, 1977)

## The EM Algorithm

- 1 Initialize  $\Theta$  as  $\Theta^{(0)}$ , set  $t = 1$
- 2 Step 1: Compute **posterior** of latent variables given current parameters  $\Theta^{(t-1)}$

$$p(\mathbf{z}_n^{(t)} | \mathbf{x}_n, \Theta^{(t-1)}) = \frac{p(\mathbf{z}_n^{(t)} | \Theta^{(t-1)}) p(\mathbf{x}_n | \mathbf{z}_n^{(t)}, \Theta^{(t-1)})}{p(\mathbf{x}_n | \Theta^{(t-1)})} \propto \text{prior} \times \text{likelihood}$$

- 3 Step 2: Now maximize the **expected complete data log-likelihood** w.r.t.  $\Theta$

$$\Theta^{(t)} = \arg \max_{\Theta} \mathcal{Q}(\Theta, \Theta^{(t-1)}) = \arg \max_{\Theta} \sum_{n=1}^N \mathbb{E}_{p(\mathbf{z}_n^{(t)} | \mathbf{x}_n, \Theta^{(t-1)})} [\log p(\mathbf{x}_n, \mathbf{z}_n^{(t)} | \Theta)]$$

- 4 If not yet converged, set  $t = t + 1$  and go to step 2.

- Note: If we can take the MAP estimate  $\hat{\mathbf{z}}_n$  of  $\mathbf{z}_n$  (not full posterior) in Step 2 and maximize the CLL in Step 3 using that, i.e., do  $\arg \max_{\Theta} \sum_{n=1}^N \left[ \log p(\mathbf{x}_n, \hat{\mathbf{z}}_n^{(t)} | \Theta) \right]$  this will be ALT-OPT



# The Expected CLL

- Expected CLL in EM is given by (assume observations are i.i.d.)

$$\begin{aligned} Q(\Theta, \Theta^{old}) &= \sum_{n=1}^N \mathbb{E}_{p(\mathbf{z}_n|\mathbf{x}_n, \Theta^{old})} [\log p(\mathbf{x}_n, \mathbf{z}_n|\Theta)] \\ &= \sum_{n=1}^N \mathbb{E}_{p(\mathbf{z}_n|\mathbf{x}_n, \Theta^{old})} [\log p(\mathbf{x}_n|\mathbf{z}_n, \Theta) + \log p(\mathbf{z}_n|\Theta)] \end{aligned}$$

Was indeed the case of GMM:  $p(\mathbf{z}_n|\Theta)$  was multinoulli,  $p(\mathbf{x}_n|\mathbf{z}_n, \Theta)$  was Gaussian

- If  $p(\mathbf{z}_n|\Theta)$  and  $p(\mathbf{x}_n|\mathbf{z}_n, \Theta)$  are exponential family distributions, then  $Q(\Theta, \Theta^{old})$  has a very simple form
- In resulting expressions, replace terms containing  $\mathbf{z}_n$ 's by their respective expectations, e.g.,
  - $\mathbf{z}_n$  replaced by  $\mathbb{E}_{p(\mathbf{z}_n|\mathbf{x}_n, \hat{\Theta})}[\mathbf{z}_n]$
  - $\mathbf{z}_n \mathbf{z}_n^\top$  replaced by  $\mathbb{E}_{p(\mathbf{z}_n|\mathbf{x}_n, \hat{\Theta})}[\mathbf{z}_n \mathbf{z}_n^\top]$
- However, in some LVMs, these expectations are intractable to compute and need to be approximated (beyond the scope of CS771)



# Another LVM: Probabilistic PCA (PPCA)

- Assume  $\mathbf{x}_n \in \mathbb{R}^D$  as a linear mapping of a latent var  $\mathbf{z}_n \in \mathbb{R}^K$  + Gaussian noise

A “reverse” generative way of thinking about PCA (low-dim  $\mathbf{z}_n$  generating high-dim  $\mathbf{x}_n$ )

$D \times 1$  offset/bias

$D \times K$  matrix

Drawn from a zero-mean  $D$ -dim Gaussian  $\mathcal{N}(\mathbf{0}, \sigma^2 I_D)$

This linear mapping can be replaced by more powerful nonlinear mapping of the form  $\mathbf{x}_n = \mathbf{f}(\mathbf{z}_n) + \boldsymbol{\epsilon}_n$  where  $\mathbf{f}$  can be modeled using a deep neural net (e.g., models like variational autoencoders)

$$\mathbf{x}_n = \boldsymbol{\mu} + \mathbf{W}\mathbf{z}_n + \boldsymbol{\epsilon}_n$$

- Equivalent to saying  $p(\mathbf{x}_n | \mathbf{z}_n, \boldsymbol{\mu}, \mathbf{W}, \sigma^2) = \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu} + \mathbf{W}\mathbf{z}_n, \sigma^2 I_D)$
- Assume a zero-mean  $K$ -dim Gaussian prior on  $\mathbf{z}_n$ , so  $p(\mathbf{z}_n) = \mathcal{N}(\mathbf{z}_n | \mathbf{0}, I_K)$
- We would like to do MLE for  $\Theta = (\boldsymbol{\mu}, \mathbf{W}, \sigma^2)$
- ILL for this model  $p(\mathbf{x}_n | \boldsymbol{\mu}, \mathbf{W}, \sigma^2)$  is also a Gaussian (thanks to Gaussian properties)

$$p(\mathbf{x}_n | \boldsymbol{\mu}, \mathbf{W}, \sigma^2) = \int p(\mathbf{x}_n | \mathbf{z}_n, \boldsymbol{\mu}, \mathbf{W}, \sigma^2) p(\mathbf{z}_n) d\mathbf{z}_n = \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2 I_D)$$

- Maximizing ILL w.r.t.  $\Theta = (\boldsymbol{\mu}, \mathbf{W}, \sigma^2)$  is possible but requires solving eig decomp. problem
- We can use ALT-OPT/EM to estimate  $\Theta = (\boldsymbol{\mu}, \mathbf{W}, \sigma^2)$  more efficiently without eig decomp.

PRML 12.2.1

# Learning PPCA using EM

- Instead of maximizing the ILL  $p(\mathbf{x}_n | \boldsymbol{\mu}, \mathbf{W}, \sigma^2) = N(\mathbf{x}_n | \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_D)$ , let's use ALT-OPT/EM
- EM will instead maximize expected CLL, with CLL (assume  $\boldsymbol{\mu} = \mathbf{0}$ ) given by

$$\log p(X, Z | W, \sigma^2) = \log \prod_{n=1}^N p(\mathbf{x}_n, \mathbf{z}_n | \mathbf{W}, \sigma^2) = \log \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n, \mathbf{W}, \sigma^2) p(\mathbf{z}_n)$$

- Using  $p(\mathbf{x}_n | \mathbf{z}_n, \mathbf{W}, \sigma^2) = \mathcal{N}(\mathbf{x}_n | \mathbf{W}\mathbf{z}_n, \sigma^2 \mathbf{I}_D)$  and  $p(\mathbf{z}_n) = \mathcal{N}(\mathbf{z}_n | \mathbf{0}, \mathbf{I}_K)$

$$\text{CLL} = - \sum_{n=1}^N \left\{ \frac{D}{2} \log \sigma^2 + \frac{1}{2\sigma^2} \|\mathbf{x}_n\|^2 - \frac{1}{\sigma^2} \mathbf{z}_n^\top \mathbf{W}^\top \mathbf{x}_n + \frac{1}{2\sigma^2} \text{trace}(\mathbf{z}_n \mathbf{z}_n^\top \mathbf{W}^\top \mathbf{W}) + \frac{1}{2} \text{trace}(\mathbf{z}_n \mathbf{z}_n^\top) \right\}$$

- Expected CLL will require  $\mathbb{E}[\mathbf{z}_n]$  and  $\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top]$  w.r.t. conditional posterior of  $\mathbf{z}_n$

Using the fact that  $p(\mathbf{x}_n | \mathbf{z}_n)$  and  $p(\mathbf{z}_n)$  are Gaussians and the CP is just the reverse conditional  $p(\mathbf{z}_n | \mathbf{x}_n)$  and must also be Gaussian

$$\begin{aligned} p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}) &= \mathcal{N}(\mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n, \sigma^2 \mathbf{M}^{-1}) \quad \text{where } \mathbf{M} = \mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}_K \\ \mathbb{E}[\mathbf{z}_n] &= \mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n \\ \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] &= \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^\top + \text{cov}(\mathbf{z}_n) = \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^\top + \sigma^2 \mathbf{M}^{-1} \end{aligned}$$



# Learning PPCA using EM

- The EM algo for PPCA alternates between two steps
  - Compute CP of  $\mathbf{z}_n$  given parameters  $\Theta = (\mathbf{W}, \sigma^2)$  and required expectatuions

$$p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}) = \mathcal{N}(\mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n, \sigma^2 \mathbf{M}^{-1}) \quad \text{where } \mathbf{M} = \mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}_K$$

$$\mathbb{E}[\mathbf{z}_n] = \mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n$$

$$\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] = \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^\top + \text{cov}(\mathbf{z}_n) = \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^\top + \sigma^2 \mathbf{M}^{-1}$$

Note: This approach does not assume/ensure that  $\mathbf{W}$  is orthonormal

- Maximize the expected CLL  $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \mathbf{W}, \sigma^2)]$  w.r.t.  $\mathbf{W}$  and  $\sigma^2$

$$\mathbb{E}[\text{CLL}] = - \sum_{n=1}^N \left\{ \frac{D}{2} \log \sigma^2 + \frac{1}{2\sigma^2} \|\mathbf{x}_n\|^2 - \frac{1}{\sigma^2} \mathbb{E}[\mathbf{z}_n^\top] \mathbf{W}^\top \mathbf{x}_n + \frac{1}{2\sigma^2} \text{trace}(\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] \mathbf{W}^\top \mathbf{W}) + \frac{1}{2} \text{trace}(\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top]) \right\}$$

Unlike standard PCA (non-probabilistic), no eigendecomposition needed to estimate  $\mathbf{W}$

$$\mathbf{W}_{\text{new}} = \left[ \sum_{n=1}^N \mathbf{x}_n \mathbb{E}[\mathbf{z}_n]^\top \right] \left[ \sum_{n=1}^N \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] \right]^{-1}$$

$$\sigma_{\text{new}}^2 = \frac{1}{ND} \sum_{n=1}^N \left\{ \|\mathbf{x}_n\|^2 - 2 \mathbb{E}[\mathbf{z}_n]^\top \mathbf{W}_{\text{new}}^\top \mathbf{x}_n + \text{tr} \left( \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] \mathbf{W}_{\text{new}}^\top \mathbf{W}_{\text{new}} \right) \right\}$$

Note: setting  $\sigma^2 = 0$  makes it equivalent to standard PCA without orthonormality constraint, but EM is **more efficient** since no eigendecomposition is needed

- Will get ALT-OPT if we use mode of the CP as  $\hat{\mathbf{z}}_n$  in the CLL

# Generative Models can generate synthetic data!

10

- Once parameters  $\Theta = (\mu, W, \sigma^2)$  are learned, we can even generate new data, e.g.,
  - Generate a random  $\mathbf{z}_n$  from  $\mathcal{N}(\mathbf{0}, I_K)$
  - Generate  $\mathbf{x}_n$  condition on  $\mathbf{z}_n$  from  $\mathcal{N}(\mu + W\mathbf{z}_n, \sigma^2 I_D)$

In addition to, of course, reducing the data dimensionality



(a) Training data



(b) Random samples

Generated using a more sophisticated generative model, not PPCA (but similar in formulation)

Methods such as variational autoencoders, GAN, diffusion models, etc are based on similar ideas



# EM: Some Comments

- Good initialization is important
- The E and M steps may not always be possible to perform exactly. Some reasons
  - CP of latent variables  $p(\mathbf{Z}|\mathbf{X}, \Theta)$  may not be easy to find and may require approx.
  - Even if  $p(\mathbf{Z}|\mathbf{X}, \Theta)$  is easy, expected CLL, i.e.,  $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)]$  may still not be tractable

$$\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)] = \int \log p(\mathbf{X}, \mathbf{Z}|\Theta) p(\mathbf{Z}|\mathbf{X}, \Theta) d\mathbf{Z}$$

Monte-Carlo EM

..and may need to be approximated, e.g., using [Monte-Carlo expectation](#)

Gradient methods may still be needed for this step

- Maximization of the expected CLL may not be possible in closed form
- EM works even if the M step is only solved approximately ([Generalized EM](#))
- Other advanced probabilistic inference algorithms are based on ideas similar to EM
  - E.g., Variational Bayesian inference a.k.a. [Variational Inference \(VI\)](#)
- EM is also related to non-convex optimization algorithms [Majorization-Maximization \(MM\)](#)
  - MM maximizes a difficult-to-optimize objective function by iteratively constructing surrogate functions that are easier to maximize (in EM, the surrogate function was the CLL)

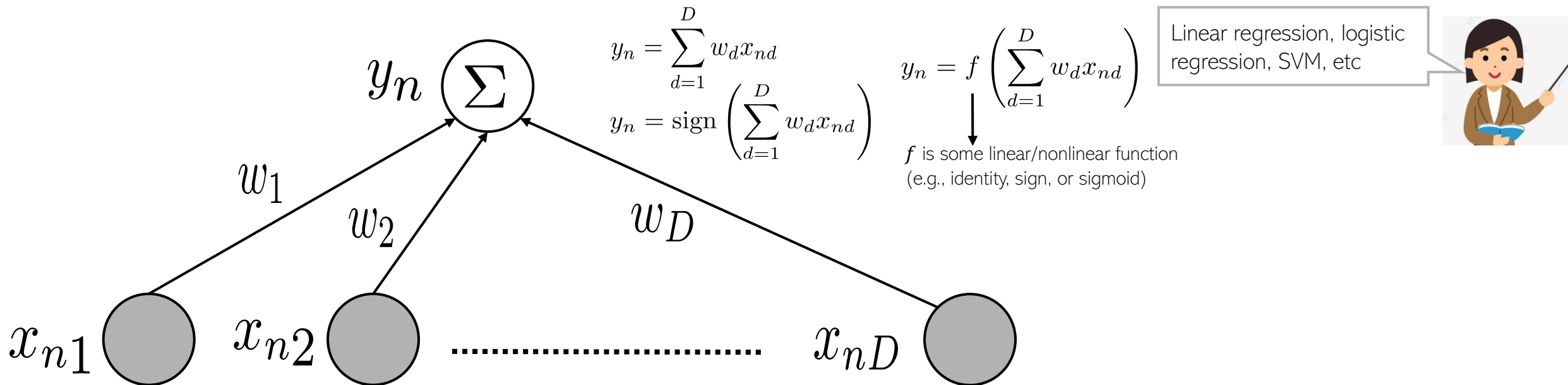


# Deep Neural Nets



# Limitation of Linear Models

- Linear models: Output produced by taking a linear combination of input features



- A basic unit of the form  $y = f(\mathbf{w}^T \mathbf{x})$  is known as the “Perceptron” (not to be confused with the Perceptron “algorithm”, which learns a linear classification model)

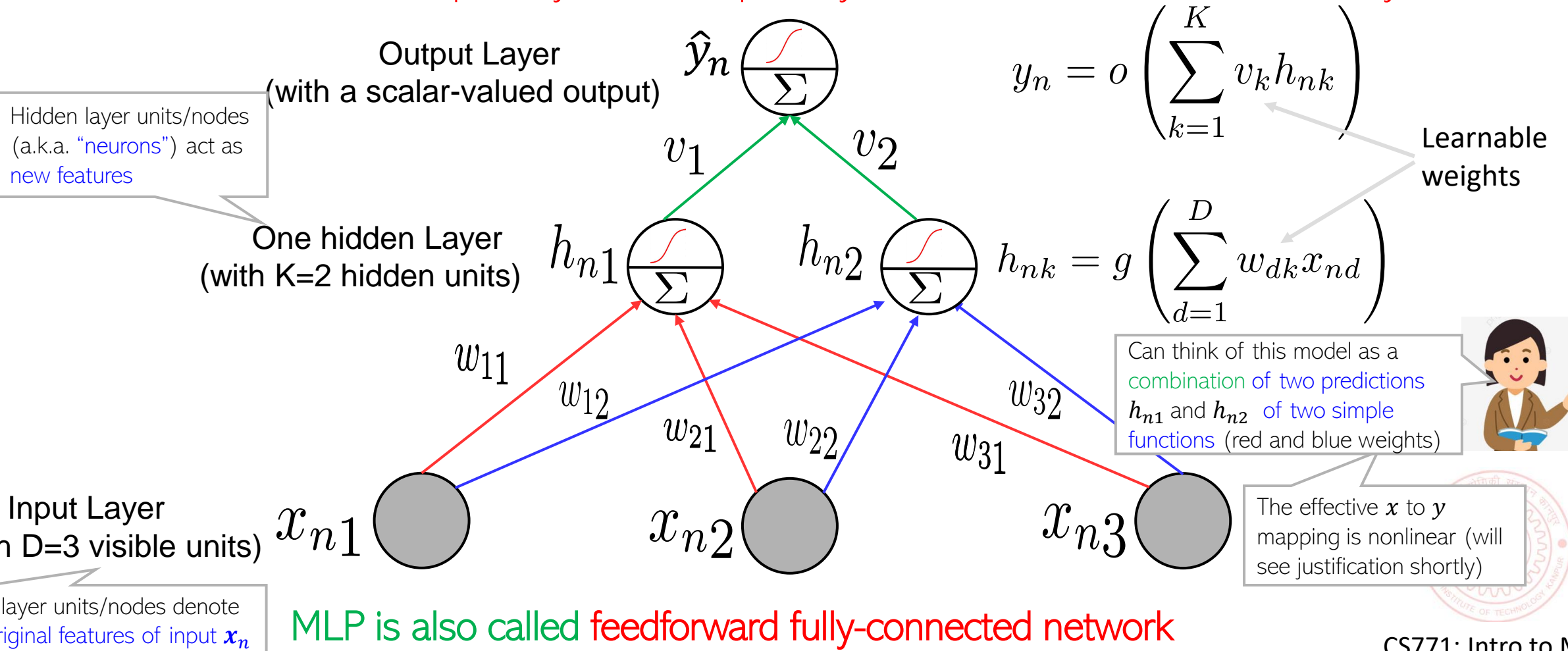
- This can't however learn nonlinear functions or nonlinear decision boundaries

Although can kernelize to make them nonlinear



# Neural Networks: Multi-layer Perceptron (MLP)

- An MLP is a network containing several Perceptron units across many layers
- An MLP consists of an **input layer**, an **output layer**, and **one or more hidden layers**



# Illustration: Neural Net with Single Hidden Layer

- Compute  $K$  pre-activations for each input  $\mathbf{x}_n$

A linear model with learnable weight vec  $\mathbf{w}_k$

$$z_{nk} = \mathbf{w}_k^\top \mathbf{x}_n = \sum_{d=1}^D w_{dk} x_{nd} \quad (k = 1, 2, \dots, K)$$

- Apply nonlinear activation on each pre-act

Called a **hidden unit**

$$h_{nk} = g(z_{nk}) \quad (k = 1, 2, \dots, K)$$

- Apply a linear model with  $\mathbf{h}_n$  acting as features

A linear model with learnable weight vec  $\mathbf{v}$

Score of the input

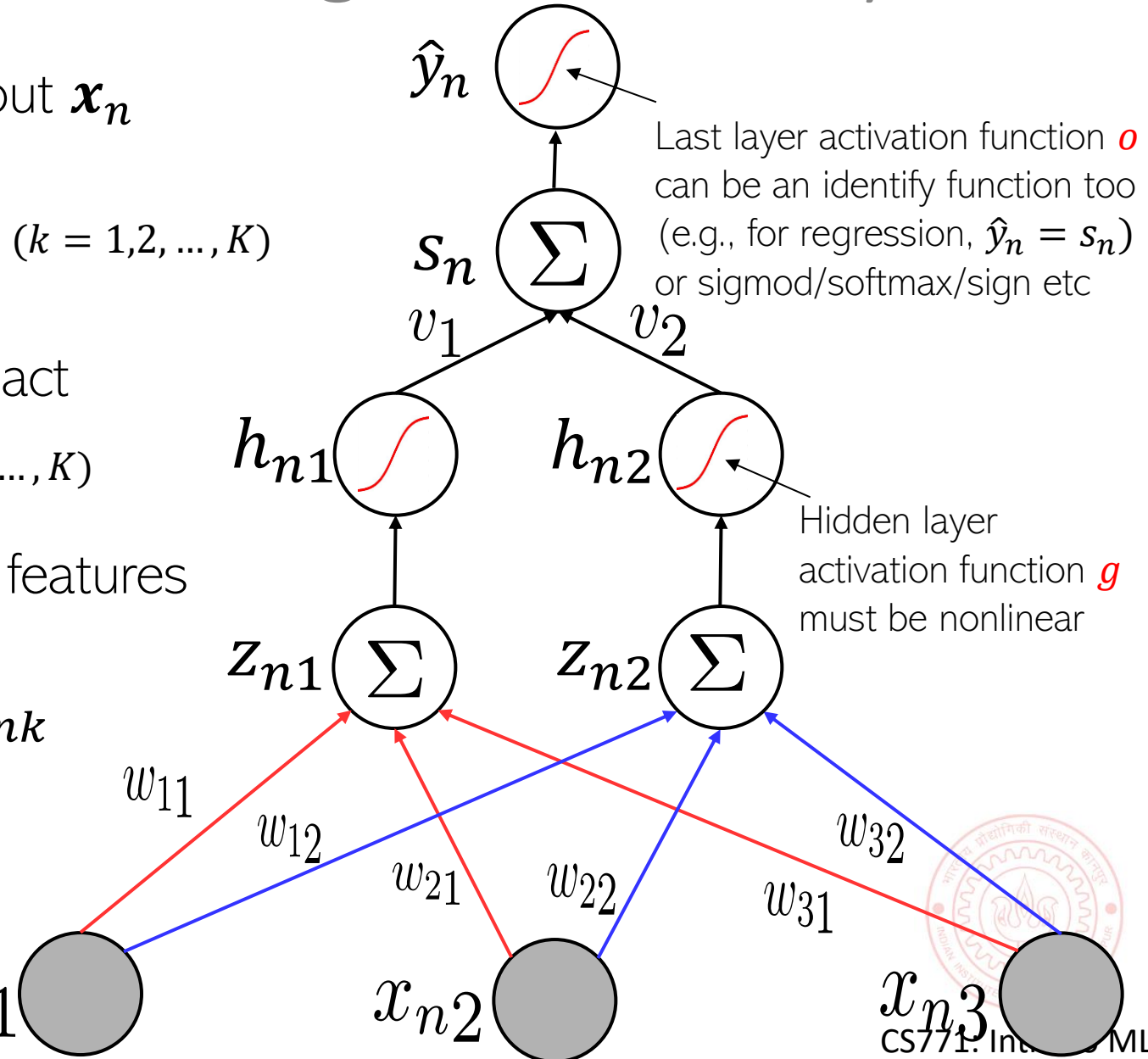
$$s_n = \mathbf{v}^\top \mathbf{h}_n = \sum_{k=1}^K v_k h_{nk}$$

- Finally, output is produced as

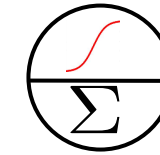
Score converted to the actual prediction

$$\hat{y}_n = o(s_n)$$

- Loss:  $\mathcal{L}(\mathbf{W}, \mathbf{v}) = \sum_{n=1}^N \ell(y_n, \hat{y}_n)$



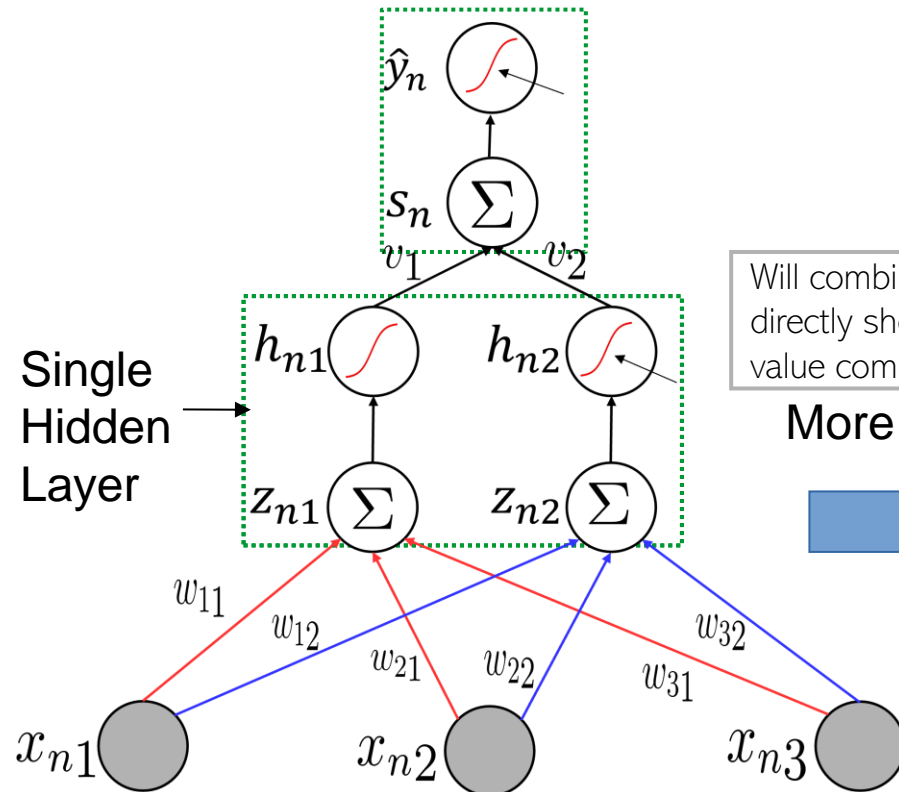
# Neural Nets: A Compact Illustration



Will denote a linear combination of inputs followed by a nonlinear operation on the result

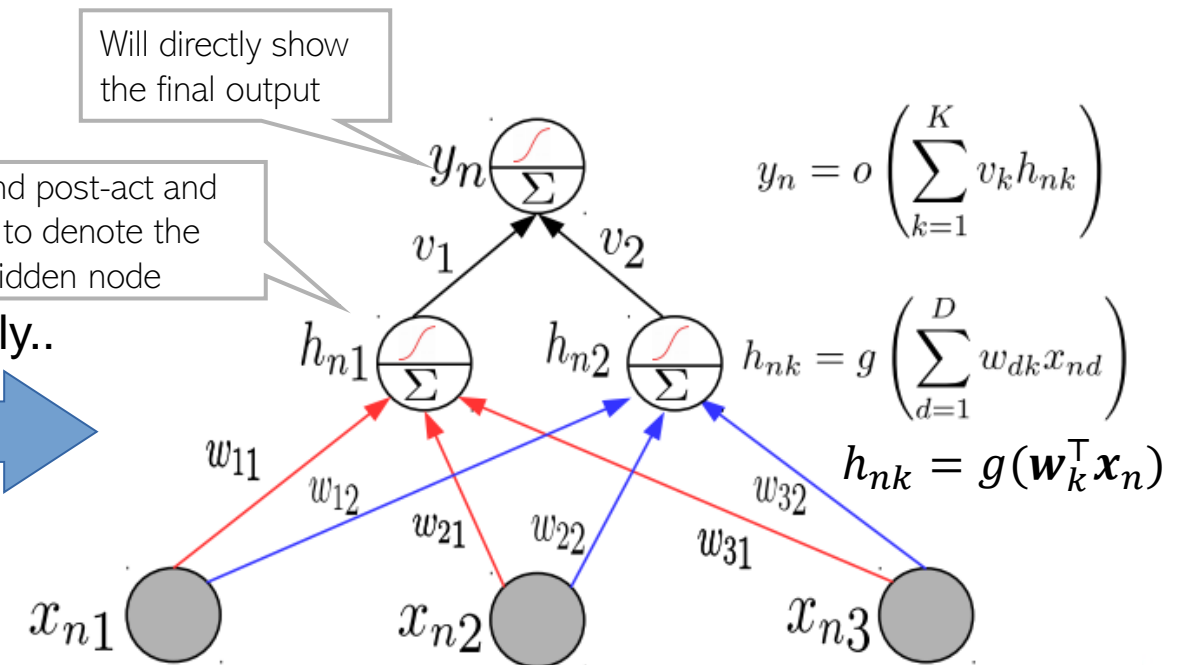
16

- Note: Hidden layer pre-act  $\mathbf{z}_{nk}$  and post-act  $\mathbf{h}_{nk}$  will be shown together for brevity



Will combine pre-act and post-act and directly show only  $\mathbf{h}_{nk}$  to denote the value computed by a hidden node

More succinctly..

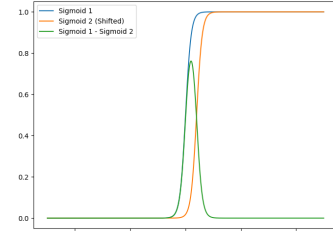


- Denoting  $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K]$ ,  $\mathbf{w}_k \in \mathbb{R}^D$ ,  $\mathbf{h}_n = g(\mathbf{W}^\top \mathbf{x}_n) \in \mathbb{R}^K$  ( $K = 2, D = 3$  above). Note:  $g$  applied elementwise on pre-activation vector  $\mathbf{z}_n = \mathbf{W}^\top \mathbf{x}_n$

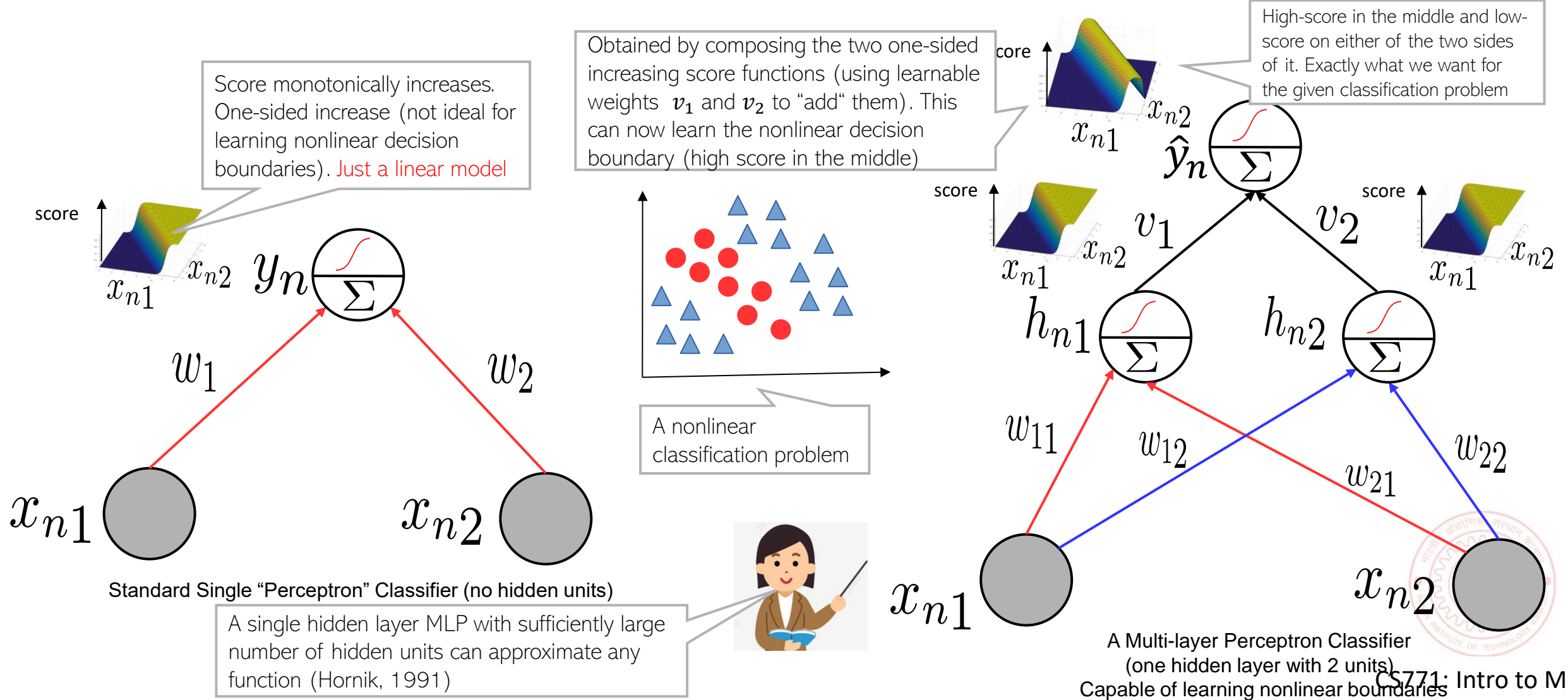




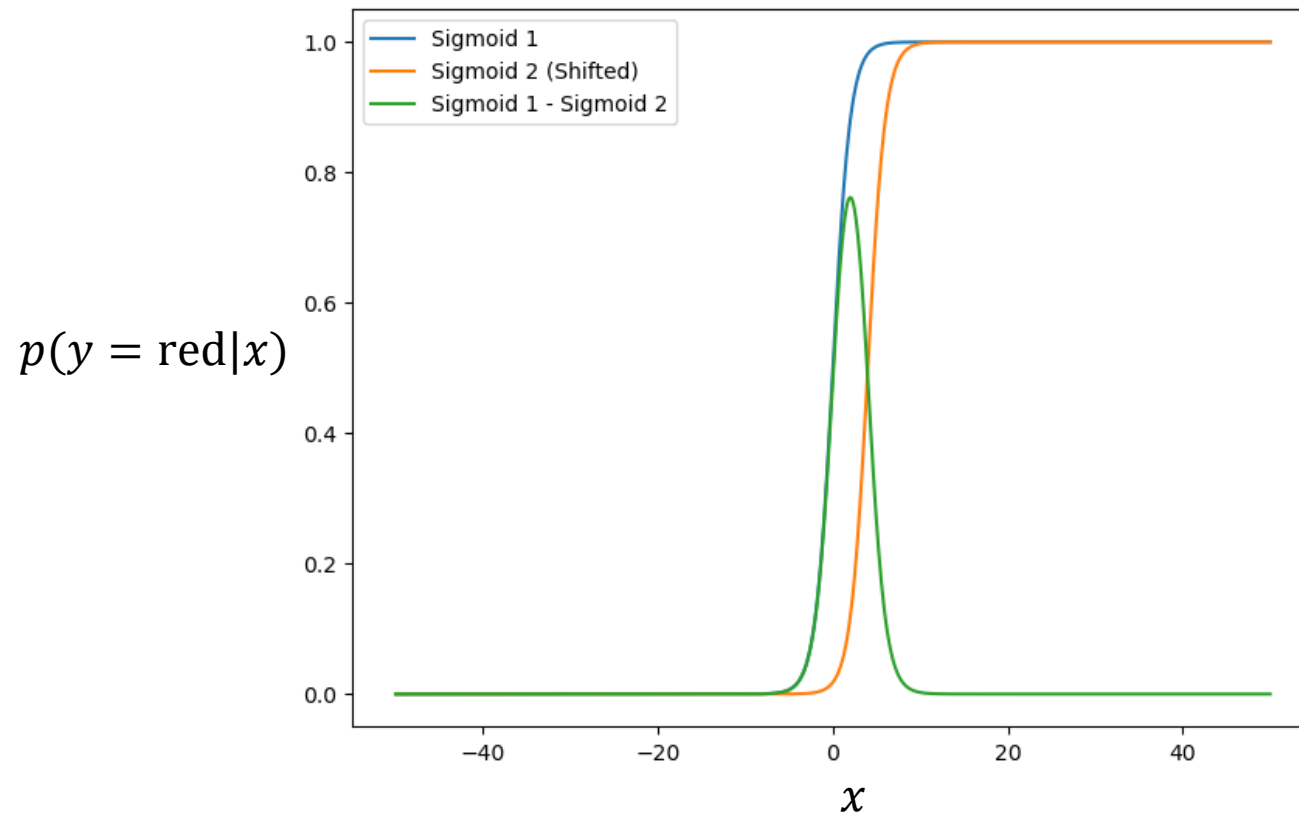
# MLP Can Learn Any Nonlinear Function



- An MLP can be seen as a composition of **multiple linear models combined nonlinearly**



# Superposition of two linear models = Nonlinear model



Two sigmoids (blue and orange) can be combined via a shift and a subtraction operation to result in a nonlinear separation boundary



Likewise, more than two sigmoids can be combined to learn even more sophisticated separation boundaries



Nonlinear separation boundary

