

Day 6: Neural networks, backpropagation

Introduction to Machine Learning Summer School
June 18, 2018 - June 29, 2018, Chicago

Instructor: Suriya Gunasekar, TTI Chicago

25 June 2018



THE UNIVERSITY OF
CHICAGO

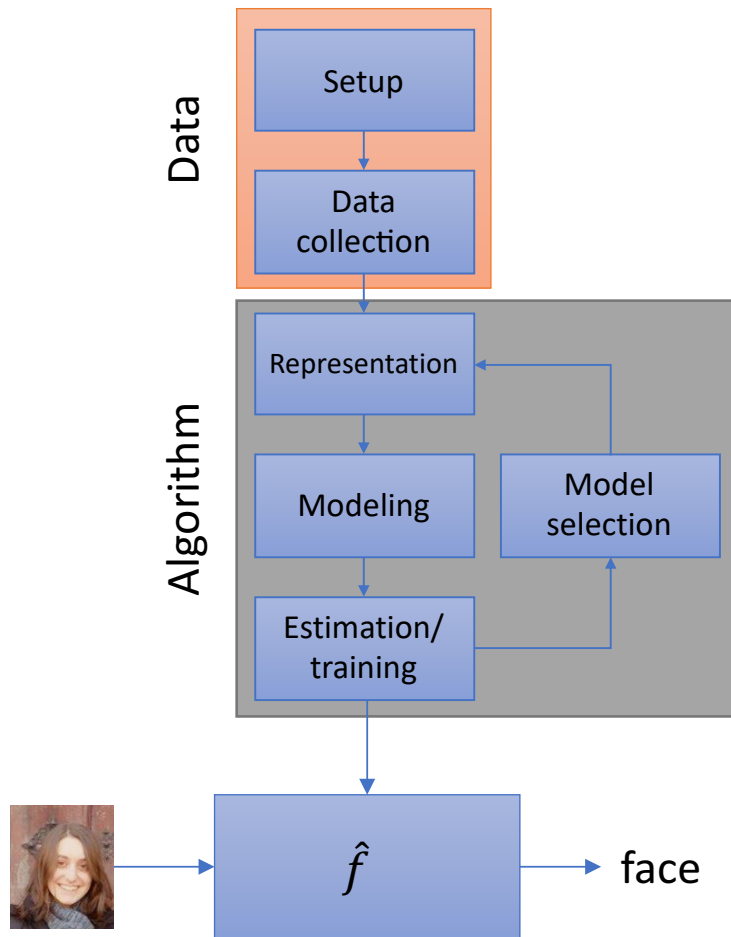


Schedule

- 9:00am-10:25am – Lecture 6.a: Review of week 1, introduction to neural networks
- 10:30am-11:30am – Invited Talk - Greg Durrett (also the TTIC colloquium talk)
- 11:30am-12:30pm – Lunch
- 12:30pm-2:00pm – Lecture 6.b: Backpropagation
- 2:00pm-5:00pm – Programming

Review of week 1

Supervised learning – key questions



- **Data:** what kind of data can we get? how much data can we get?
- **Model:** what is the correct model for my data? – want to minimize the effort put into this question!
- **Training:** what resources - computation/memory - does the algorithm need to estimate the model \hat{f} ?
- **Testing:** how well will \hat{f} perform when deployed? what is the computational/memory requirement during deployment?

Linear regression

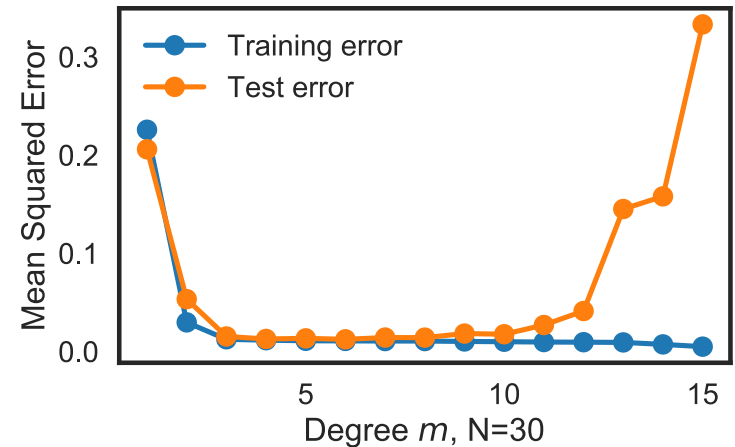
- Input $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d$, output $y \in \mathbb{R}$, want to learn $f: \mathcal{X} \rightarrow \mathbb{R}$
- Training data $S = \{(\mathbf{x}^{(i)}, y^{(i)}): i = 1, 2, \dots, N\}$
- Parameterize candidate $f: \mathcal{X} \rightarrow \mathbb{R}$ by linear functions,
 $\mathcal{H} = \{\mathbf{x} \rightarrow \mathbf{w} \cdot \mathbf{x}: \mathbf{w} \in \mathbb{R}^d\}$
- Estimate \mathbf{w} by minimizing loss on training data

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} J_S^{LS}(\mathbf{w}) := \sum_{i=1}^N (\mathbf{w} \cdot \mathbf{x}^{(i)} - y^{(i)})^2$$

- $J_S^{LS}(\mathbf{w})$ is convex in $\mathbf{w} \rightarrow$ minimize $J_S^{LS}(\mathbf{w})$ by setting gradient to 0
 - $\nabla_{\mathbf{w}} J_S^{LS}(\mathbf{w}) = \sum_{i=1}^N (\mathbf{w} \cdot \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$
 - Closed form solution $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X} \mathbf{y}$
- Can get non-linear functions by mapping $\mathbf{x} \rightarrow \phi(\mathbf{x})$ and doing linear regression on $\phi(\mathbf{x})$

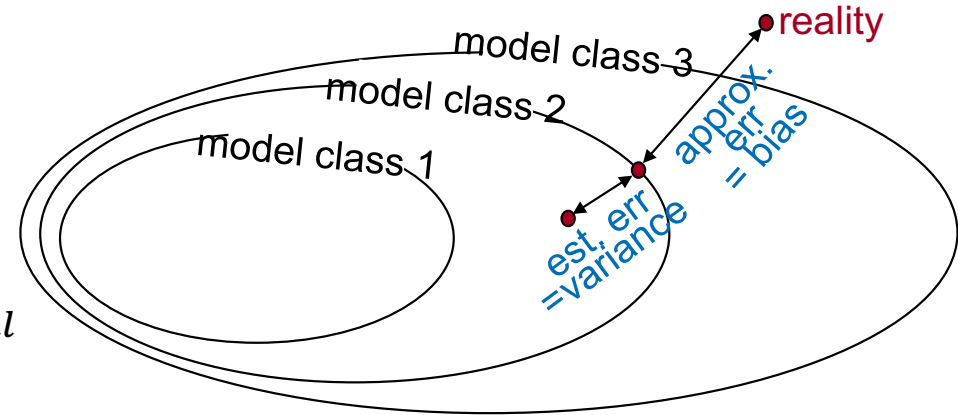
Overfitting

- For same amount of data, more complex models (e.g., higher degree polynomials) overfit more
- or need more data to fit more complex models
- complexity \approx number of parameters



Model selection

- m model classes $\{\mathcal{H}_1, \mathcal{H}_2, \dots, \mathcal{H}_m\}$
- $S = S_{train} \cup S_{val} \cup S_{test}$
- Train on S_{train} to pick best $\hat{f}_r \in \mathcal{H}_r$
- Pick \hat{f}^* based on validation loss on S_{val}
- Evaluate test loss $L_{S_{test}}(\hat{f}^*)$



Regularization

- Complexity of model class can also be controlled by norm of parameters – smaller range of values allowed
- Regularization for linear regression

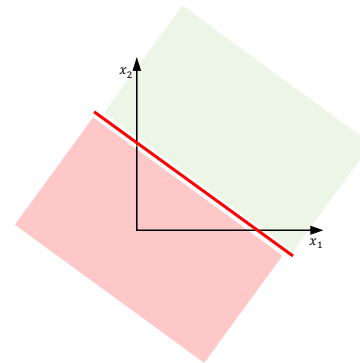
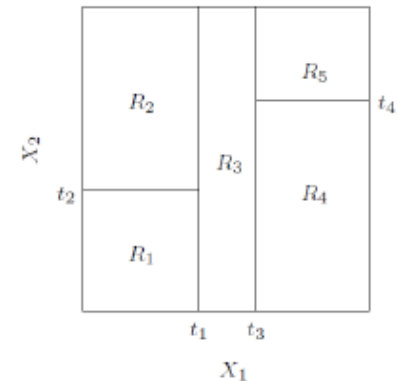
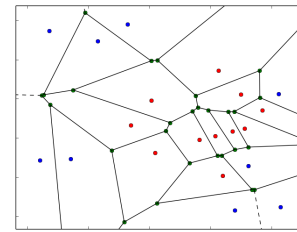
$$\operatorname{argmin}_{\mathbf{w}} J_S^{LS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$$

$$\operatorname{argmin}_{\mathbf{w}} J_S^{LS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$$

- Again do model selection to pick λ – using S_{val} or cross-validation

Classification

- Output $y \in \mathcal{Y}$ takes discrete set of values, e.g., $\mathcal{Y} = \{0,1\}$ or $\mathcal{Y} = \{-1,1\}$ or $\mathcal{Y} = \{spam, nospam\}$
 - Unlike regression, label-values do not have meaning
- Classifiers divide the space of input \mathcal{X} (often \mathbb{R}^d) to “regions” where each region is assigned a label
- Non-parametric models
 - k-nearest neighbors – regions defined based on nearest neighbors
 - decision trees – structured rectangular regions
- Linear models – classifier regions are halfspaces



Classification – logistic regression

Logistic loss

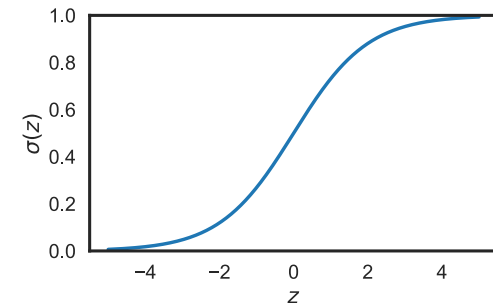
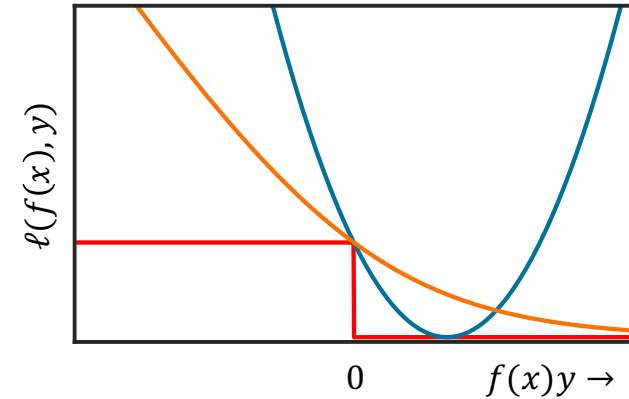
$$\ell(f(\mathbf{x}), y) = \log(1 + \exp(-f(\mathbf{x})y))$$

- $\mathcal{X} = \mathbb{R}^d$, $\mathcal{Y} = \{-1, 1\}$, $S = \{(\mathbf{x}^{(i)}, y^{(i)}) : i = 1, 2, \dots, N\}$
- Linear model $f(\mathbf{x}) = f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x}$
- Output classifier $\hat{y}(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x})$
- Empirical risk minimization

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\text{argmin}} \sum_i \log(1 + \exp(-\mathbf{w} \cdot \mathbf{x}^{(i)} y^{(i)}))$$

- Probabilistic formulation: $\Pr(y = 1 | \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w} \cdot \mathbf{x})}$
- Multi-class generalization: $\mathcal{Y} = \{1, 2, \dots, m\}$
- $\Pr(y | \mathbf{x}) = \frac{\exp(-\mathbf{w}_y \cdot \mathbf{x})}{\sum_{y'} \exp(-\mathbf{w}_{y'} \cdot \mathbf{x})}$

- Can again get non-linear decision boundaries by mapping $\mathbf{x} \rightarrow \phi(\mathbf{x})$



Classification – maximum margin classifier

Separable data

- Original formulation

$$\hat{\mathbf{w}} = \operatorname{argmax}_{\mathbf{w} \in \mathbb{R}^d} \min_i \frac{y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)}}{\|\mathbf{w}\|}$$

- Fixing $\|\mathbf{w}\| = 1$

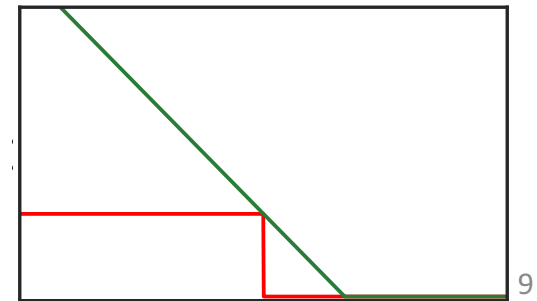
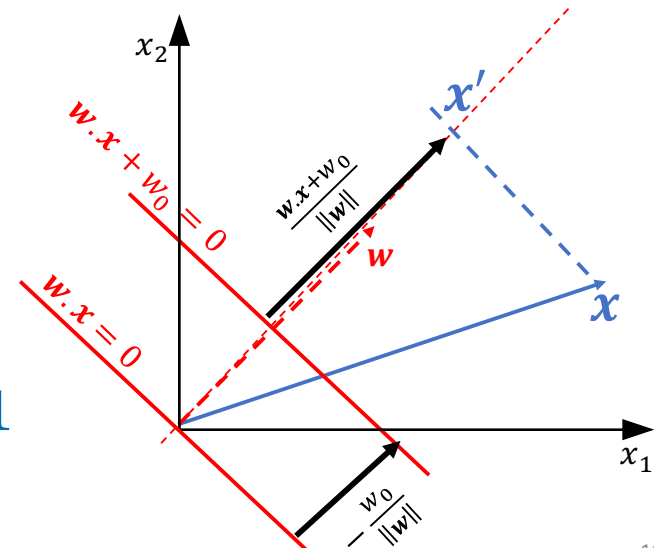
$$\hat{\mathbf{w}} = \operatorname{argmax}_{\mathbf{w}} \min_i y^{(i)} (\mathbf{w} \cdot \mathbf{x}^{(i)}) \quad \text{s.t.} \quad \|\mathbf{w}\| = 1$$

- Fixing $\min_i y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} = 1$

$$\tilde{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \|\mathbf{w}\|^2 \quad \text{s.t.} \quad \forall i, y^{(i)} (\mathbf{w} \cdot \mathbf{x}^{(i)}) \geq 1$$

Slack variables for non-separable data

$$\begin{aligned} \hat{\mathbf{w}} &= \operatorname{argmin}_{\mathbf{w}, \{\xi_i \geq 0\}} \|\mathbf{w}\|^2 + \lambda \sum_i \xi_i \quad \text{s.t.} \quad \forall i, y^{(i)} (\mathbf{w} \cdot \mathbf{x}^{(i)}) \geq 1 - \xi_i \\ &= \operatorname{argmin}_{\mathbf{w}, \{\xi_i \geq 0\}} \|\mathbf{w}\|^2 + \lambda \sum_i \max(0, 1 - y^{(i)} (\mathbf{w} \cdot \mathbf{x}^{(i)})) \end{aligned}$$



Kernel trick

- Using representer theorem $\mathbf{w} = \sum_{i=1}^N \beta_i \mathbf{x}^{(i)}$

$$\begin{aligned} \min_{\mathbf{w}} \|\mathbf{w}\|^2 + \lambda \sum_i \max(0, 1 - y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)}) \\ \equiv \min_{\boldsymbol{\beta} \in \mathbb{R}^N} \boldsymbol{\beta}^\top \mathbf{G} \boldsymbol{\beta} + \lambda \sum_i \max(0, 1 - y^{(i)} (\mathbf{G} \boldsymbol{\beta})_i) \end{aligned}$$

$\mathbf{G} \in \mathbb{R}^{N \times N}$ with $G_{ij} = \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}$ is called the gram matrix

- Optimization depends on $\mathbf{x}^{(i)}$ only through $G_{ij} = \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}$
- For prediction $\hat{\mathbf{w}} \cdot \mathbf{x} = \sum_i \beta_i \mathbf{x}^{(i)} \cdot \mathbf{x}$, we again only need $\mathbf{x}^{(i)} \cdot \mathbf{x}$
- Function $K(\mathbf{x}, \mathbf{x}') = \mathbf{x} \cdot \mathbf{x}'$ is called the Kernel
- When learning non-linear classifiers using feature transformations $\mathbf{x} \rightarrow \phi(\mathbf{x})$ and $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w} \cdot \phi(\mathbf{x})$
 - Classifier fully specified in terms of $K_{\phi}(\mathbf{x}, \mathbf{x}') = K(\phi(\mathbf{x}), \phi(\mathbf{x}'))$
 - $\phi(\mathbf{x})$ itself can be very very high dimensional (maybe even infinite dimensional)

Optimization

- ERM+regularization optimization problem

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} J_S^\lambda(\mathbf{w}) := \sum_{i=1}^N \ell(\mathbf{w} \cdot \phi(\mathbf{x}^{(i)}), y^{(i)}) + \lambda \|\mathbf{w}\|$$

- If $J_S^\lambda(\mathbf{w})$ is convex in \mathbf{w} , then $\hat{\mathbf{w}}$ is optimum if and only if gradient at $\hat{\mathbf{w}}$ is 0, i.e., $\nabla J_S^\lambda(\hat{\mathbf{w}}) = 0$
- Gradient descent: start with initialization \mathbf{w}^0 and iteratively update
 - $\mathbf{w}^{t+1} = \mathbf{w}^t - \eta^t \nabla J_S^\lambda(\mathbf{w}^t)$
 - where $\nabla J_S^\lambda(\mathbf{w}^t) = \sum_i \nabla \ell(\mathbf{w}^t \cdot \phi(\mathbf{x}^{(i)}), y^{(i)}) + \lambda \nabla \|\mathbf{w}^t\|$
- Stochastic gradient descent
 - use gradients from only one example
 - $\mathbf{w}^{t+1} = \mathbf{w}^t - \eta^t \hat{\nabla}^{(i)} J_S^\lambda(\mathbf{w}^t)$
 - where $\hat{\nabla}^{(i)} J_S^\lambda(\mathbf{w}^t) = \nabla \ell(\mathbf{w}^t \cdot \phi(\mathbf{x}^{(i)}), y^{(i)}) + \lambda \nabla \|\mathbf{w}^t\|$ for a random sample $(\mathbf{x}^{(i)}, y^{(i)})$

Other classification models

- **Optimal unrestricted predictor**
 - Regression + squared loss $\rightarrow f^{**}(\mathbf{x}) = \mathbf{E}[y|\mathbf{x}]$
 - Classification + 0-1 loss $\rightarrow \hat{y}^{**}(\mathbf{x}) = \operatorname{argmax}_c \Pr(y = c|\mathbf{x})$
- **Discriminative models:** directly model $\Pr(y|\mathbf{x})$, e.g., logistic regression
- **Generative models:** model full joint distribution $\Pr(y, \mathbf{x}) = \Pr(\mathbf{x}|y) \Pr(y)$
- **Why generative models?**
 - One conditional might be simpler to model with prior knowledge, e.g., compare specifying $\Pr(\text{image}|\text{digit} = 1)$ vs $\Pr(\text{digit} = 1|\text{image})$
 - Naturally handles missing data
- **Two examples of generative models**
 - Naïve Bayes classifier
 - Hidden Markov model

Other classifiers

- **Naïve Bayes classifier:** with d features $x = [x_1, x_2, \dots, x_d]$ where each x_1, x_2, \dots, x_d can take one of K values $\rightarrow C K^d$ parameters
 - **NB assumption:** features are independent given class $y \rightarrow C K d$ params.

$$\Pr(x_1, x_2, \dots, x_d | y) = \Pr(x_1 | y) \Pr(x_2 | y) \dots \Pr(x_d | y) = \prod_{k=1}^d \Pr(x_k | y)$$

- Training amounts to averaging samples across classes
- **Hidden Markov model:** variable length input/observations $\{x_1, x_2, \dots, x_m\}$ (e.g., words) and variable length output/state $\{y_1, y_2, \dots, y_m\}$ (e.g., tags)
 - **HMM assumption:** a) current state conditioned on immediate previous state is conditionally independent of all other variables, and (b) current observation conditioned on current state is conditionally independent of all other variables.

$$\Pr(x_1, x_2, \dots, x_m, y_1, y_2, \dots, y_m) = \Pr(y_1) \Pr(x_1 | y_1) \prod_{k=2}^m \Pr(y_k | y_{k-1}) \Pr(y_k | x_k)$$

- Parameters estimated using MLE dynamic programming

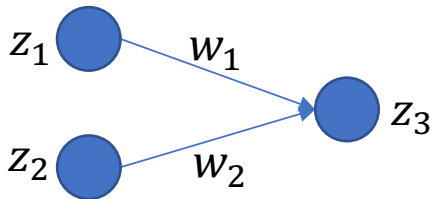
Today

Introduction to neural networks

Backpropagation

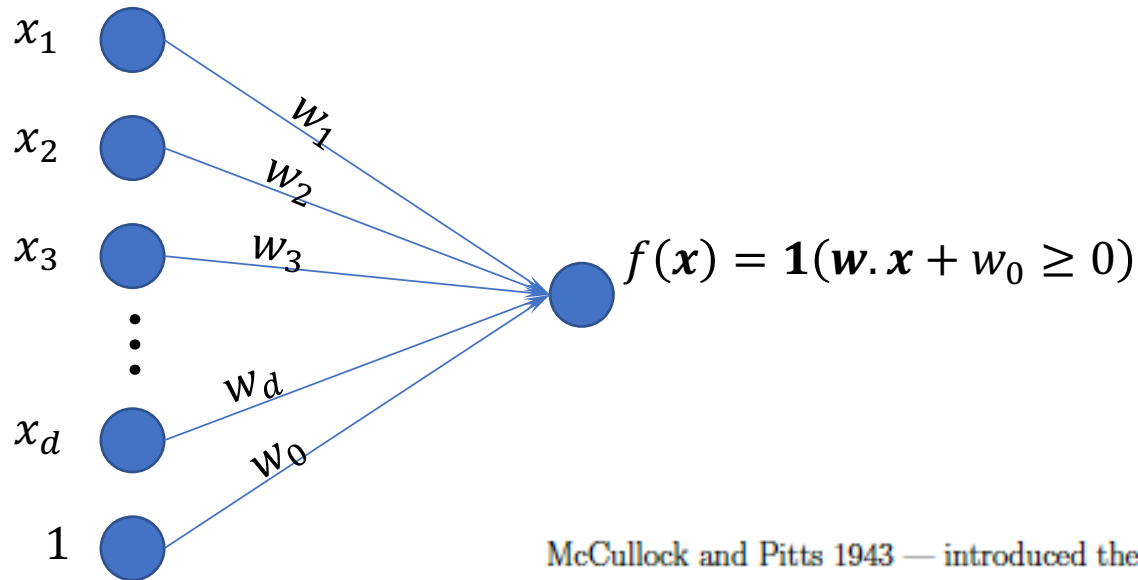
Graph notation

- General variables
 - can be input variables like x_1, x_2, \dots, x_d
 - prediction \hat{y}
 - or any intermediate computation (we will see examples soon)



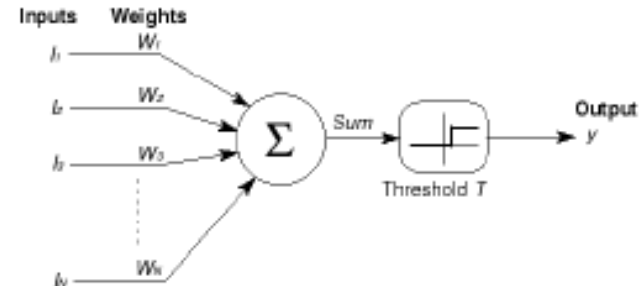
denotes computation $z_3 = \sigma(w_1 z_1 + w_2 z_2)$ for some “activation” function σ (specified apriori)

Linear classifier



McCullock and Pitts 1943 — introduced the linear threshold “neuron”.

- Biological analogy:
single neuron – stimuli
reinforce synaptic
connections

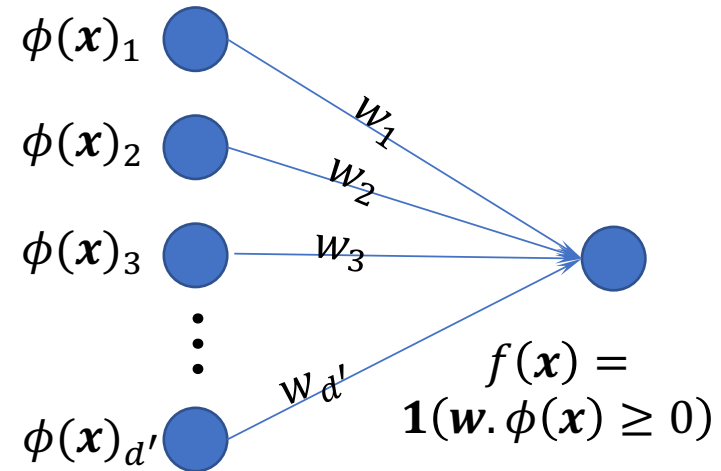


Shallow learning

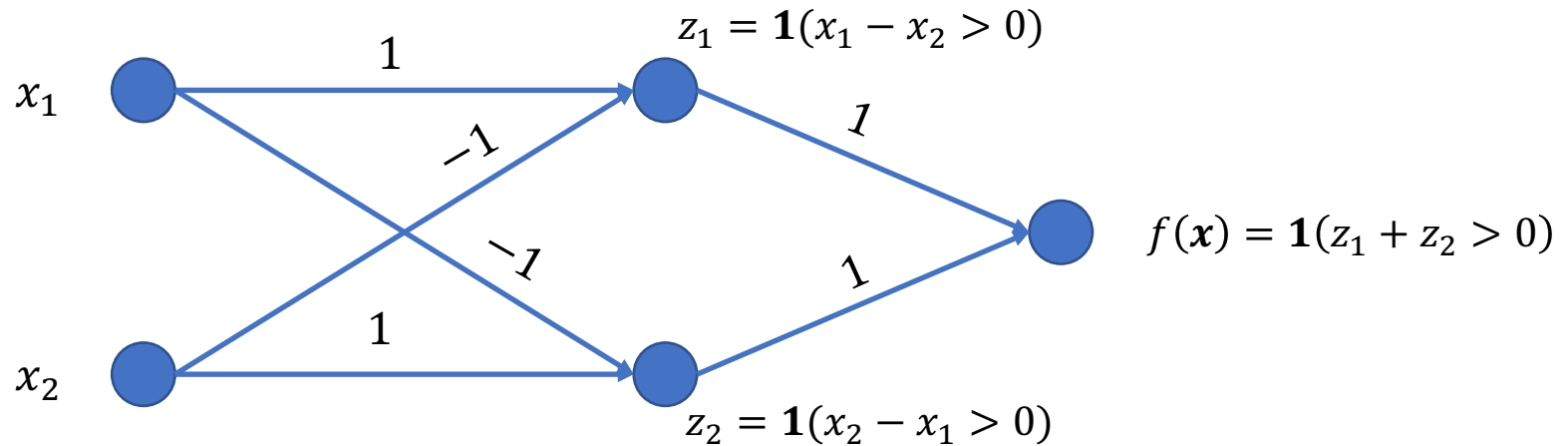
- We already saw how to use linear models to get non-linear decision boundaries
- Feature transform: map $\mathbf{x} \in \mathbb{R}^d$ to $\phi(\mathbf{x}) \in \mathbb{R}^{d'}$ and use

$$f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w} \cdot \phi(\mathbf{x})$$

- Shallow learning: hand-crafted and non-hierarchical ϕ
 - Polynomial regression with squared or logistic loss, $\phi(x)_p = x^p$
 - Kernel SVM: $K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x}) \cdot \phi(\mathbf{x}')$

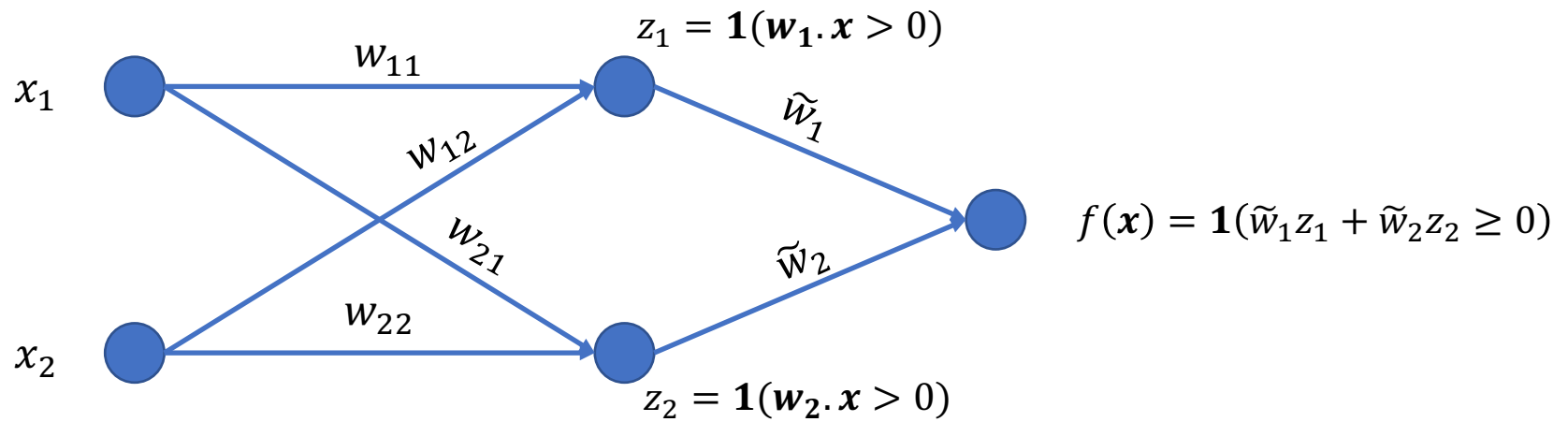


Combining Linear Units

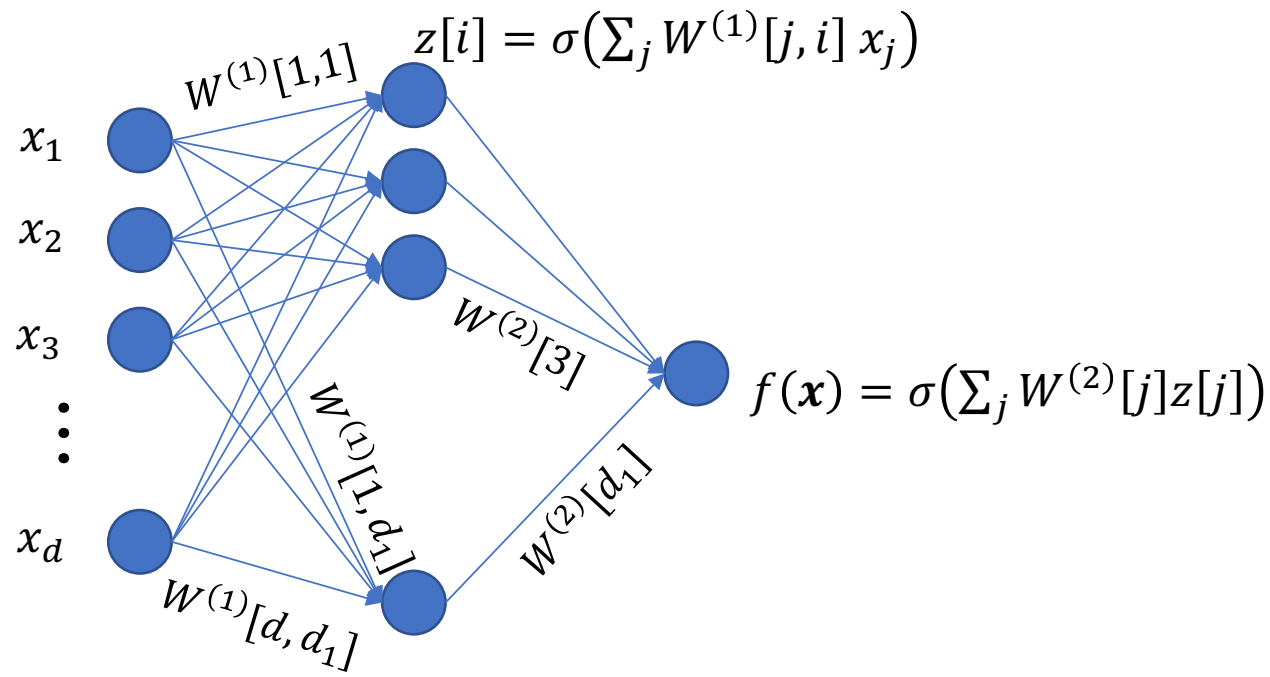


- The network represents the function
 $f(\mathbf{x}) = (x_1 \text{ and not}(x_2)) \text{ or } (x_2 \text{ and not}(x_1))$
- **Not a linear function of \mathbf{x}**

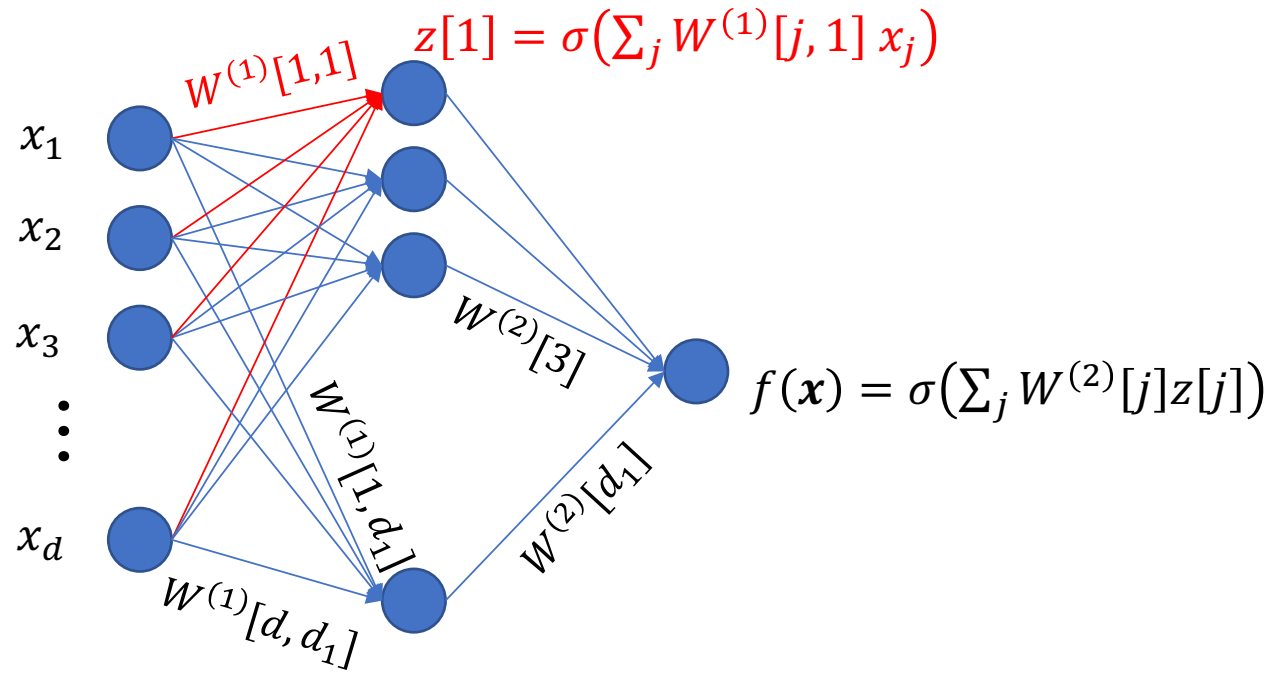
Combining Linear Units



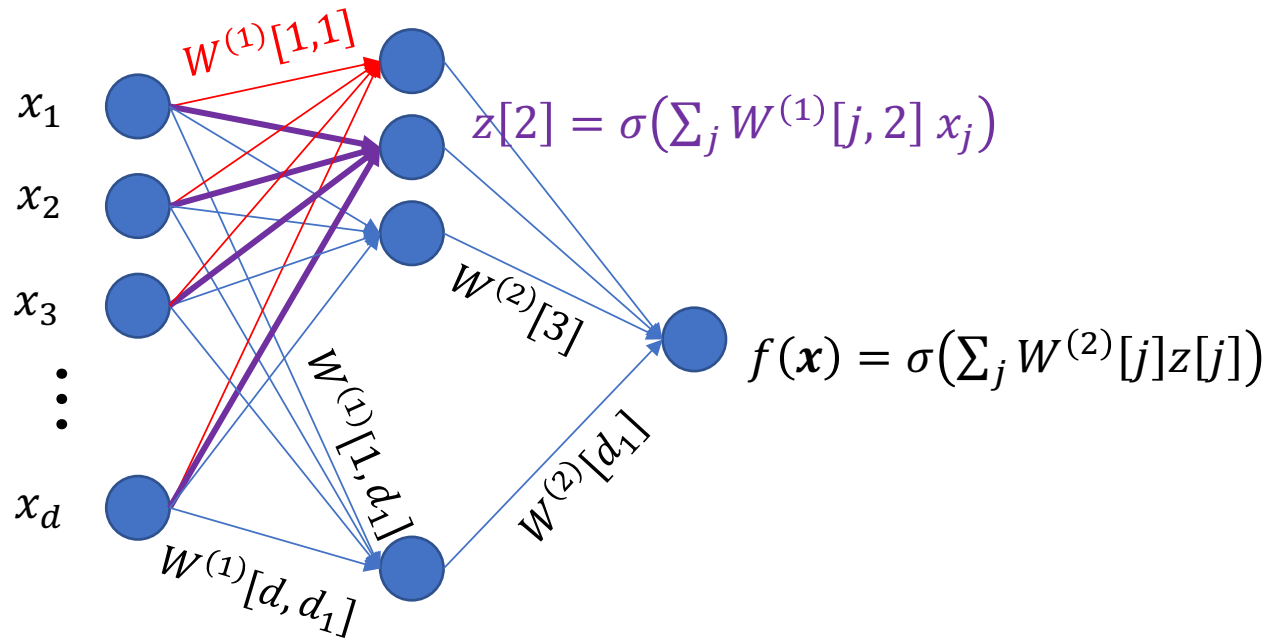
Feed-Forward Neural Networks



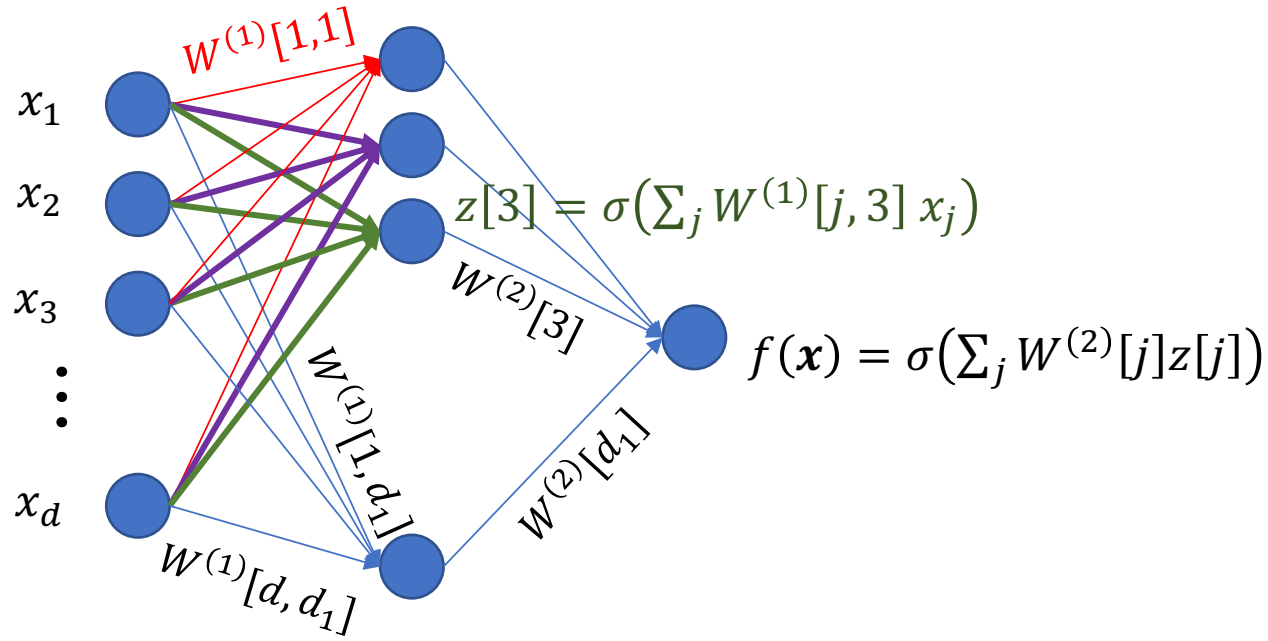
Feed-Forward Neural Networks



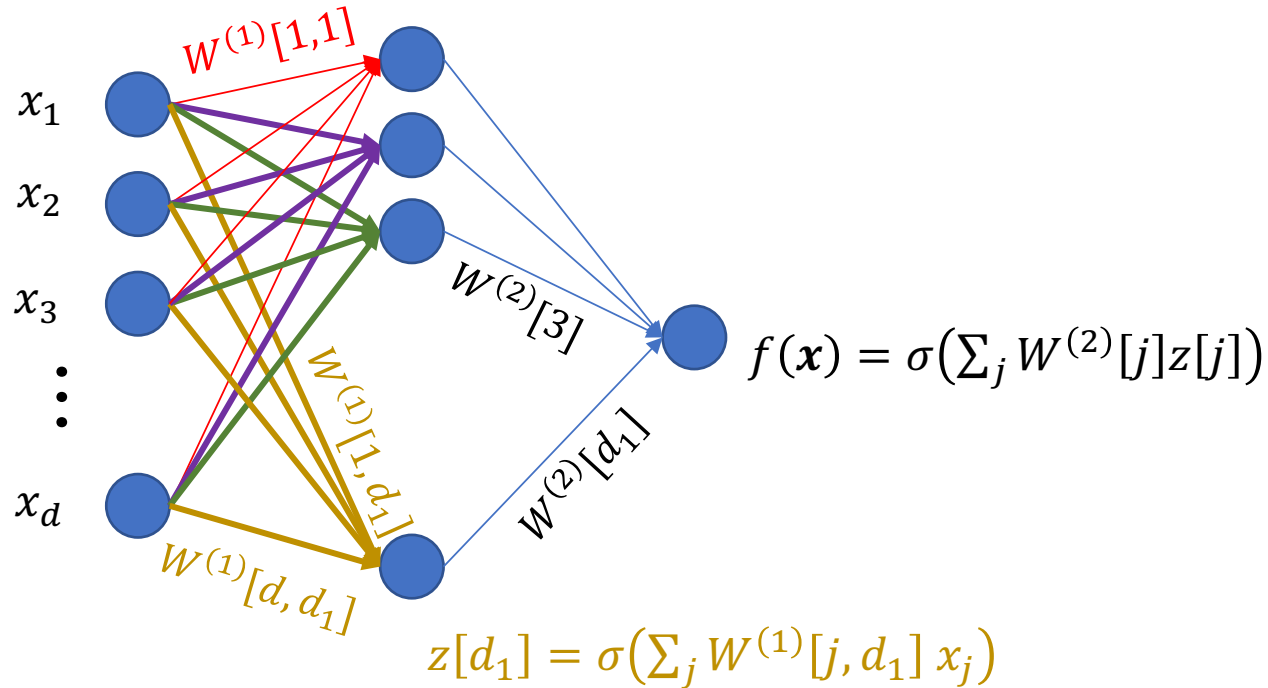
Feed-Forward Neural Networks



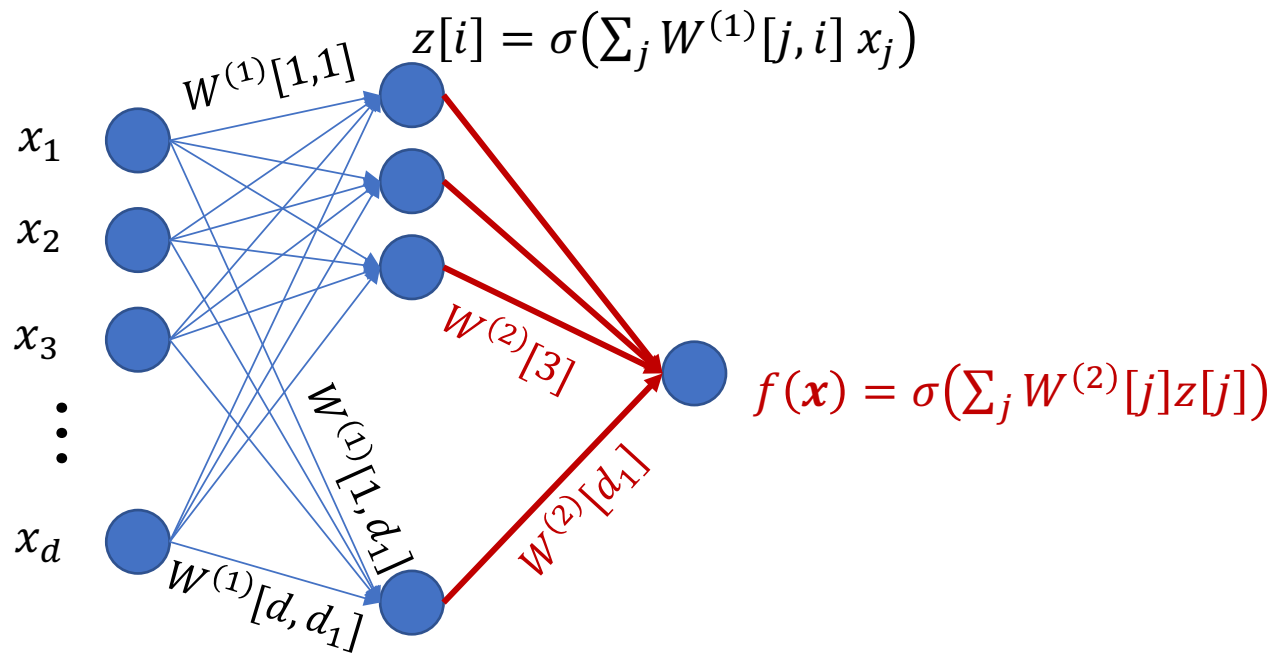
Feed-Forward Neural Networks



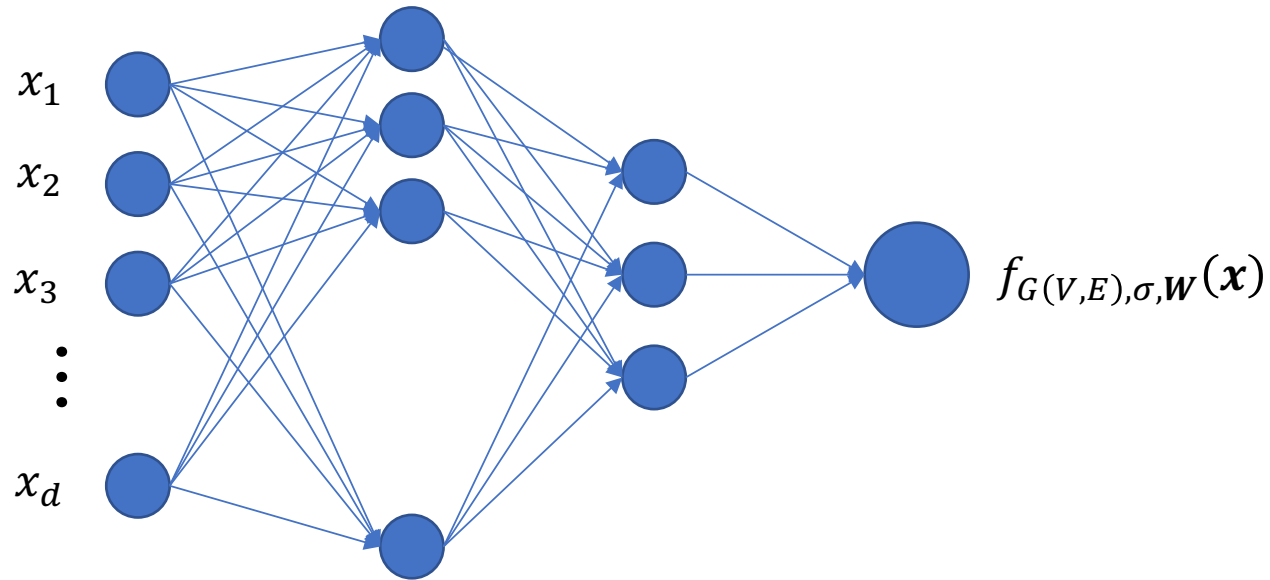
Feed-Forward Neural Networks



Feed-Forward Neural Networks



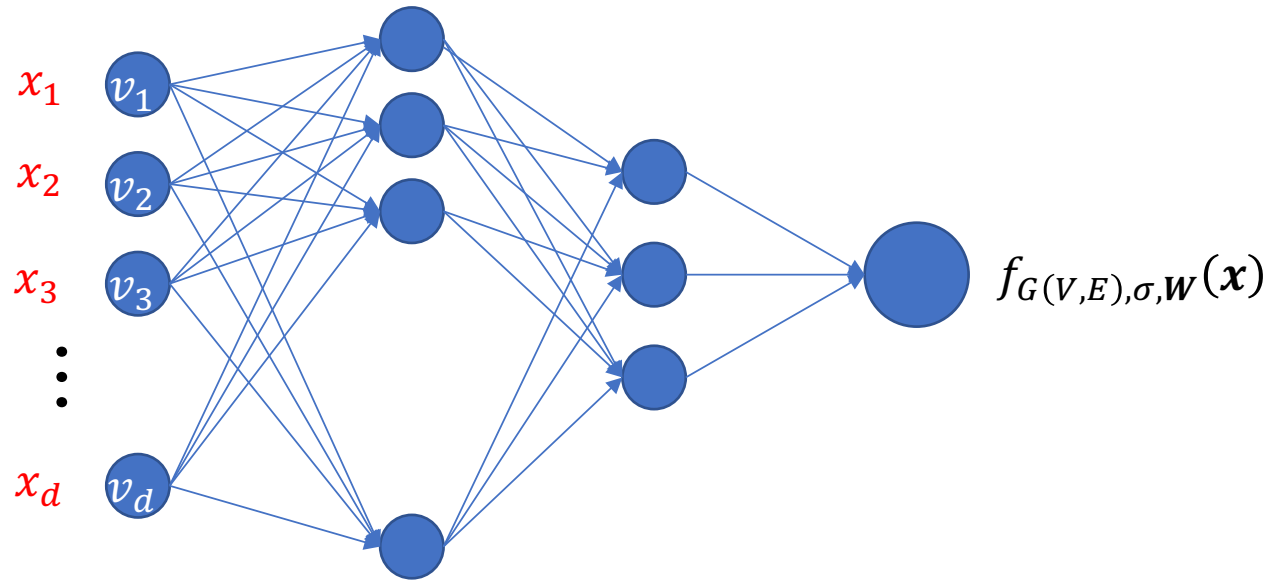
Feed-Forward Neural Networks



Architecture:

- Directed Acyclic Graph $G(V, E)$. Units (neurons) indexed by vertices in V .

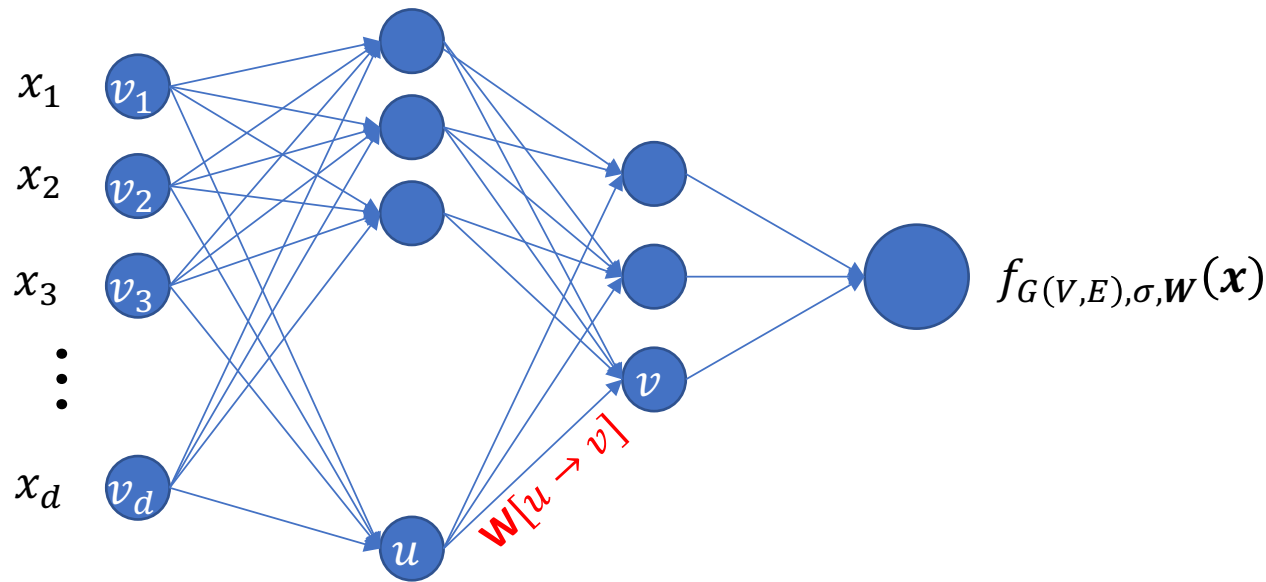
Feed-Forward Neural Networks



Architecture:

- Directed Acyclic Graph $G(V,E)$. Units (neurons) indexed by vertices in V .
 - “Input Units” $v_1 \dots v_d \in V$: no incoming edges have value $o[v_i] = x_i$

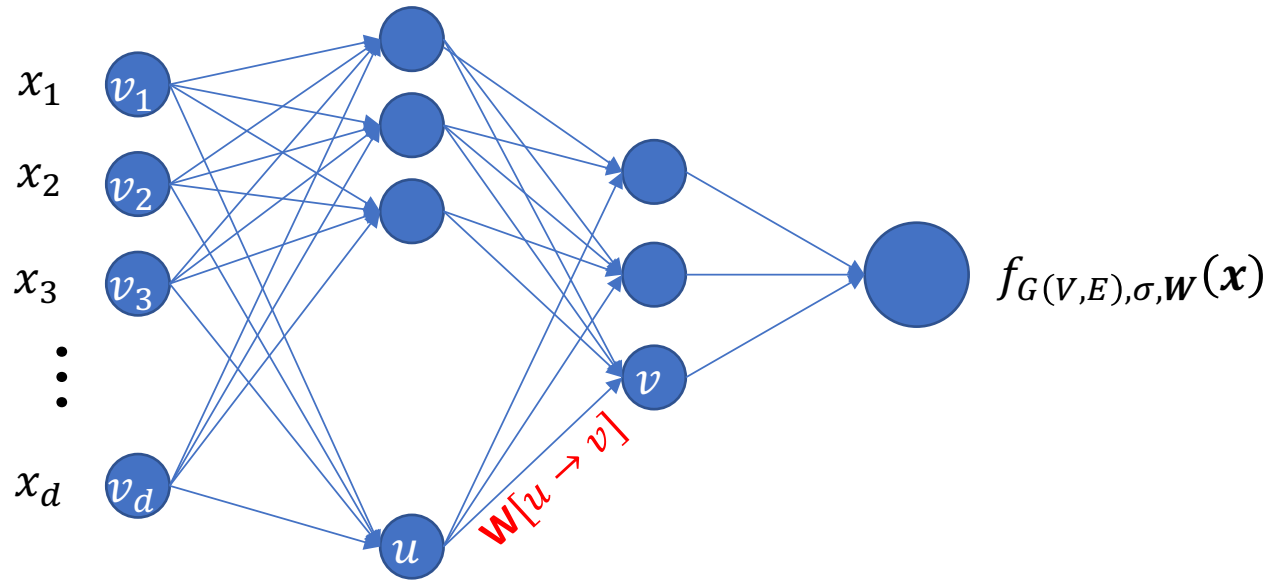
Feed-Forward Neural Networks



Architecture:

- Directed Acyclic Graph $G(V, E)$. Units (neurons) indexed by vertices in V .
 - “Input Units” $v_1 \dots v_d \in V$: no incoming edges have value $o[v_i] = x_i$
 - Each edge $u \rightarrow v$ has weight $W[u \rightarrow v]$
 - Pre-activation $a[v] = \sum_{u \rightarrow v \in E} W[u \rightarrow v] o[u]$

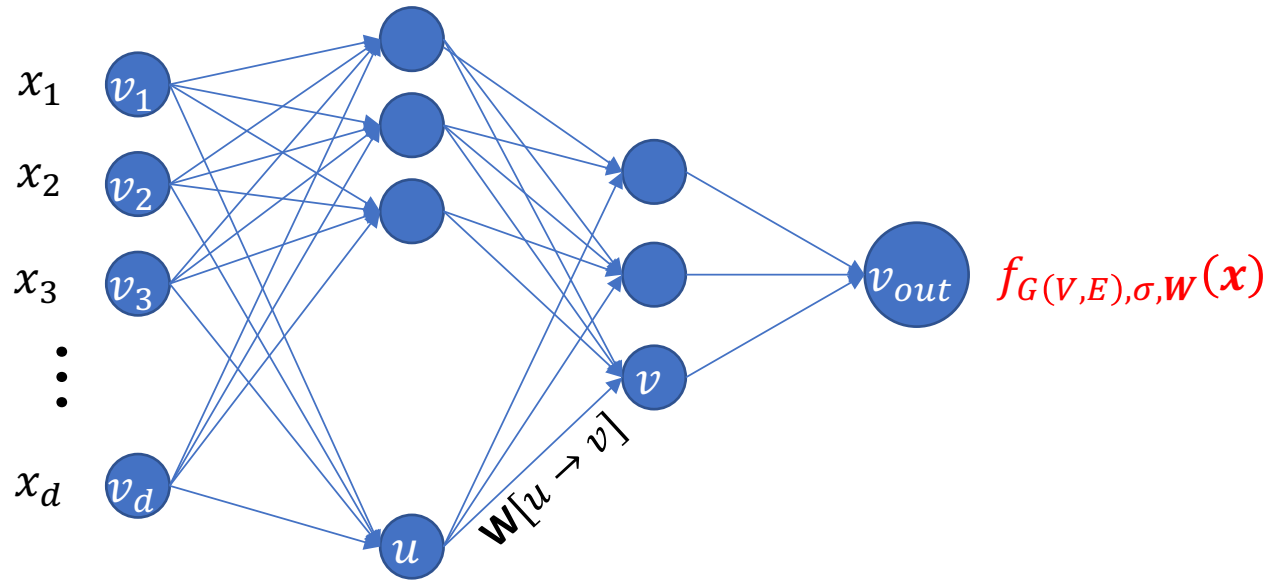
Feed-Forward Neural Networks



Architecture:

- Directed Acyclic Graph $G(V, E)$. Units (neurons) indexed by vertices in V .
 - “Input Units” $v_1 \dots v_d \in V$: no incoming edges have value $o[v_i] = x_i$
 - Each edge $u \rightarrow v$ has weight $W[u \rightarrow v]$
 - Pre-activation $a[v] = \sum_{u \rightarrow v \in E} W[u \rightarrow v] o[u]$
 - Output value $o[v] = \sigma(a[v])$

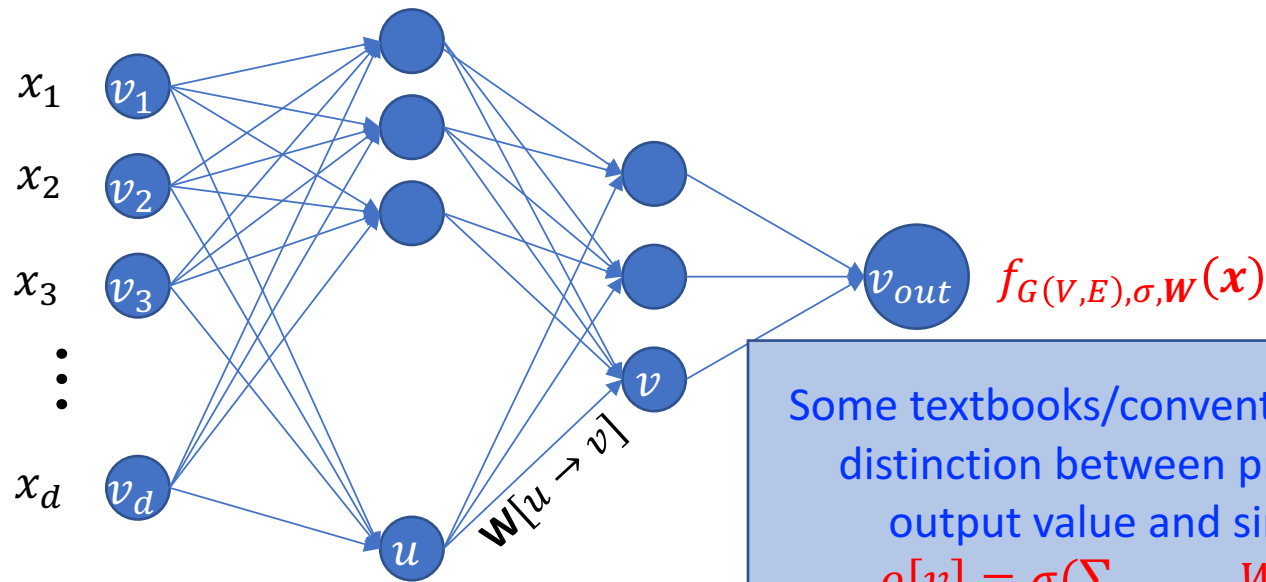
Feed-Forward Neural Networks



Architecture:

- Directed Acyclic Graph $G(V,E)$. Units (neurons) indexed by vertices in V .
 - “Input Units” $v_1 \dots v_d \in V$: no incoming edges have value $o[v_i] = x_i$
 - Each edge $u \rightarrow v$ has weight $W[u \rightarrow v]$
 - Pre-activation $a[v] = \sum_{u \rightarrow v \in E} W[u \rightarrow v] o[u]$
 - Output value $o[v] = \sigma(a[v])$
 - “Output Unit” $v_{out} \in V, f_W(x) = a[v_{out}]$

Feed-Forward Neural Networks

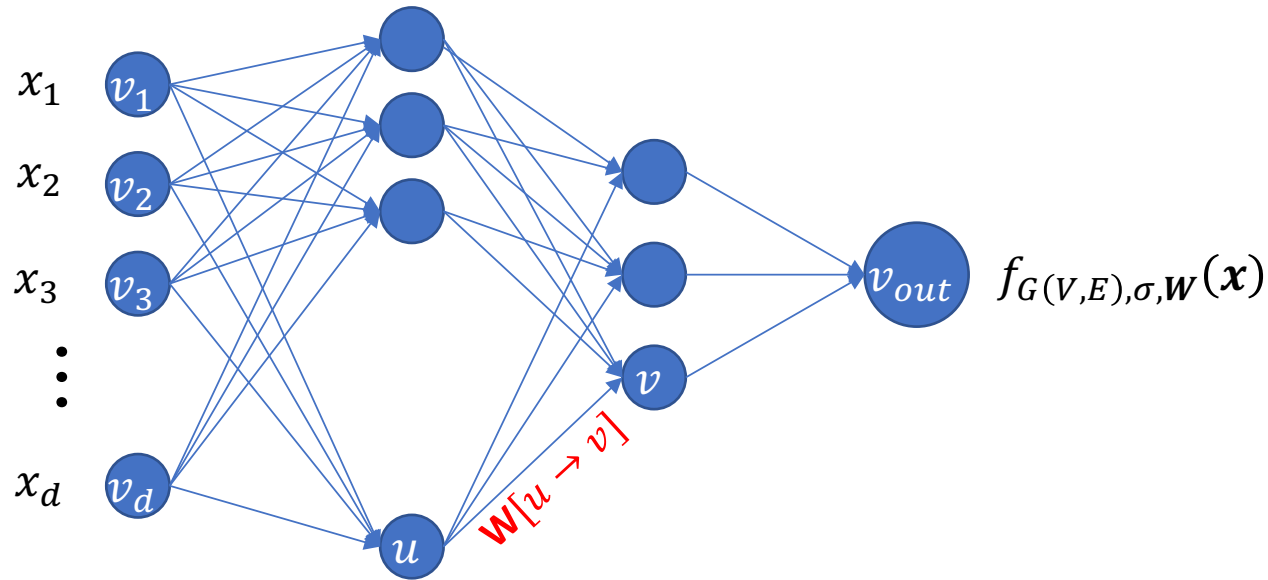


Some textbooks/convention don't make the distinction between pre-activation and output value and simply compute $o[v] = \sigma(\sum_{u \rightarrow v \in E} W[u \rightarrow v] o[u])$

Architecture:

- Directed Acyclic Graph $G(V,E)$. Units (neurons) represented by vertices in V .
 - “Input Units” $v_1 \dots v_d \in V$: no incoming edges have value $o[v_i] = x_i$
 - Each edge $u \rightarrow v$ has weight $W[u \rightarrow v]$
 - Pre-activation $a[v] = \sum_{u \rightarrow v \in E} W[u \rightarrow v] o[u]$
 - Output value $o[v] = \sigma(a[v])$
- “Output Unit” $v_{out} \in V, f_W(x) = a[v_{out}]$

Feed-Forward Neural Networks



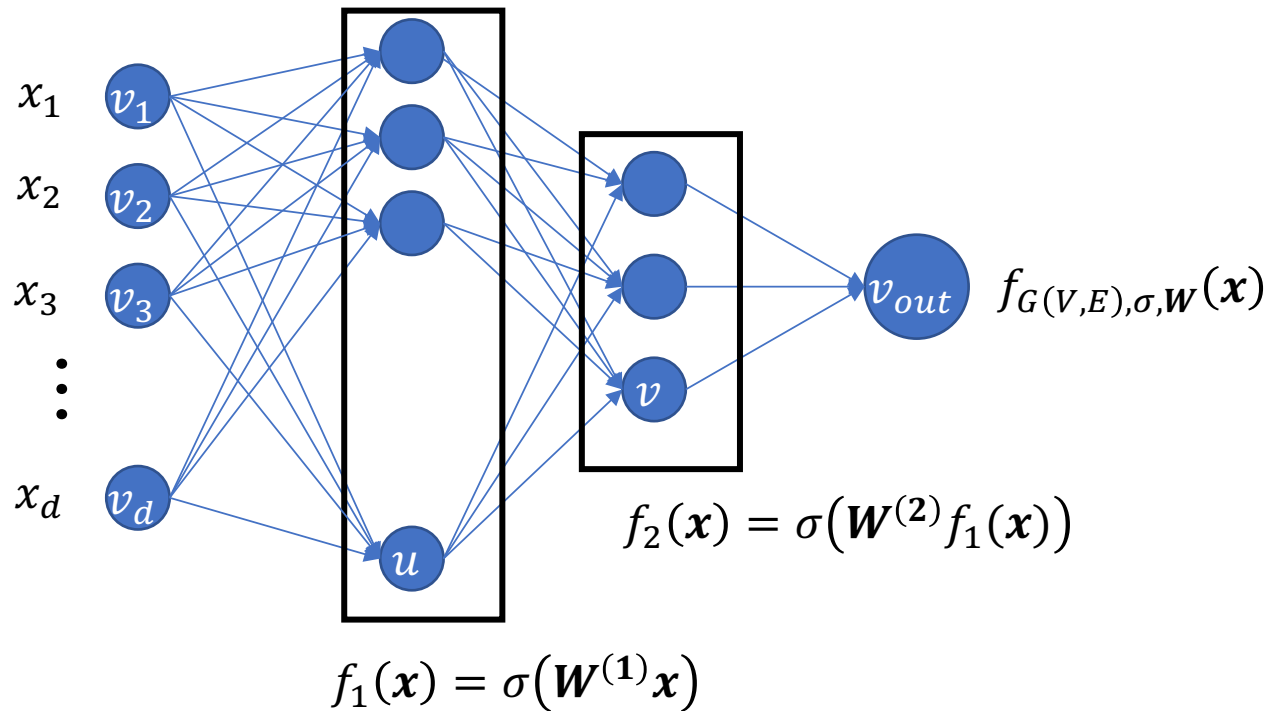
Parameters:

- Each edge $u \rightarrow v$ has weight $W[u \rightarrow v]$

Activations:

- $\sigma: \mathbb{R} \rightarrow \mathbb{R}$, for example
 - $\sigma(z) = \text{sign}(z)$ or $\sigma(z) = \frac{1}{1+\exp(-z)}$
 - $\sigma(z) = \text{ReLU}(z) = \max(0, z)$

Feed-Forward Neural Networks

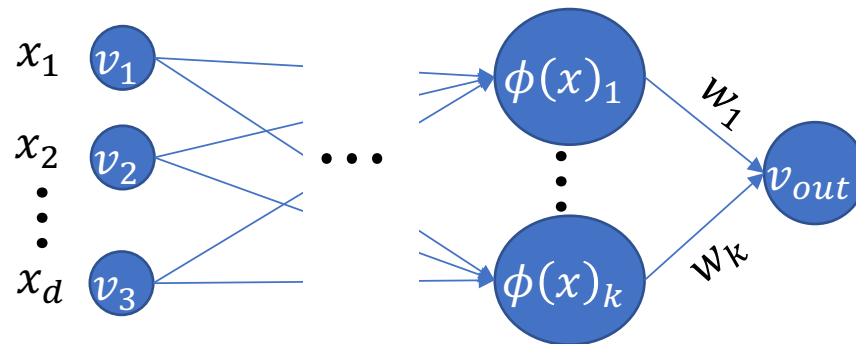


Deep learning

Generalize to hierarchy of transformations of the input, learned end-to-end jointly with the predictor.

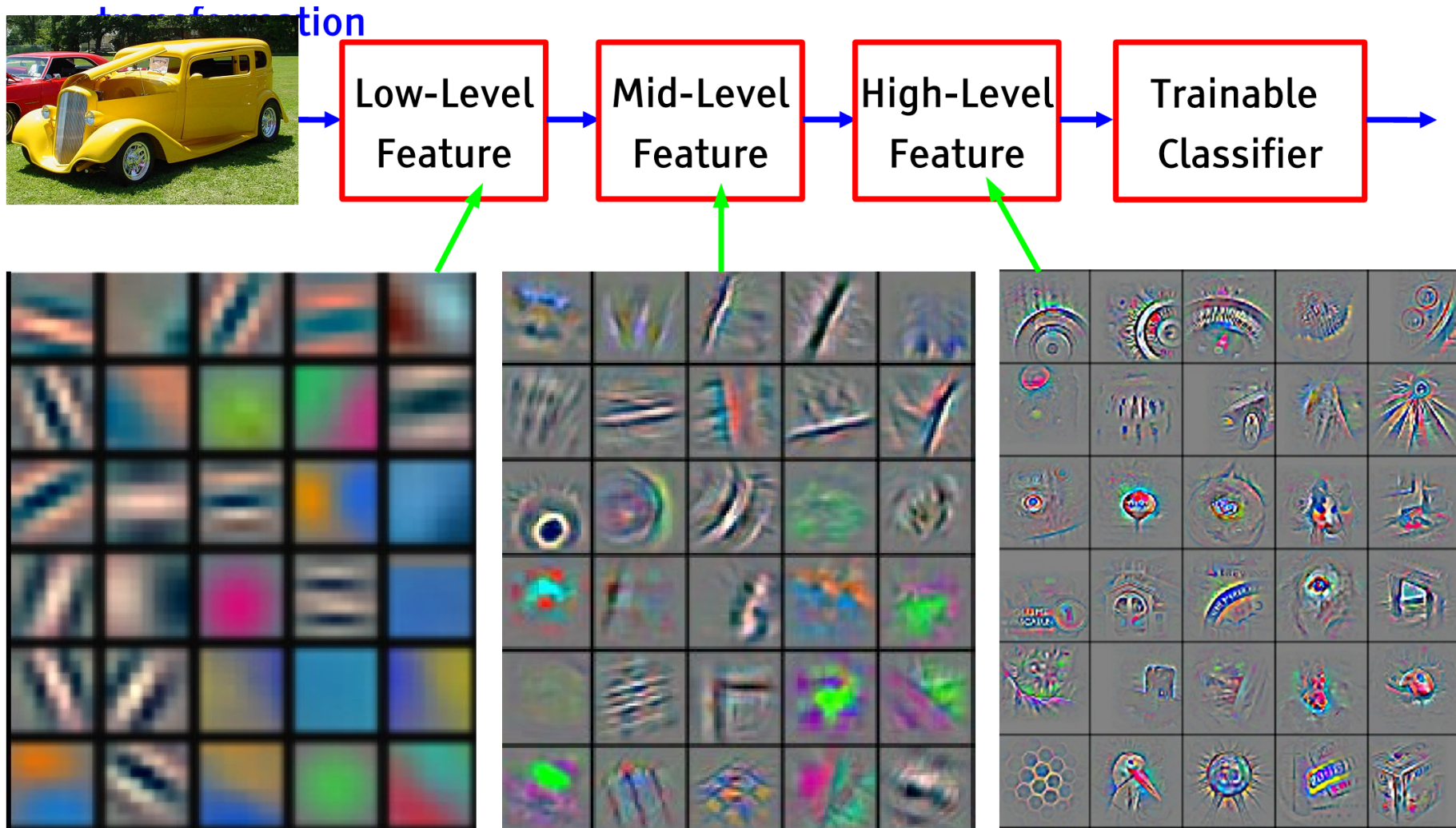
$$f_{\mathbf{W}}(\mathbf{x}) = f_L \left(f_{L-1} \left(f_{L-2} \left(\dots f_1(\mathbf{x}) \dots \right) \right) \right)$$

Neural Nets as Feature Learning

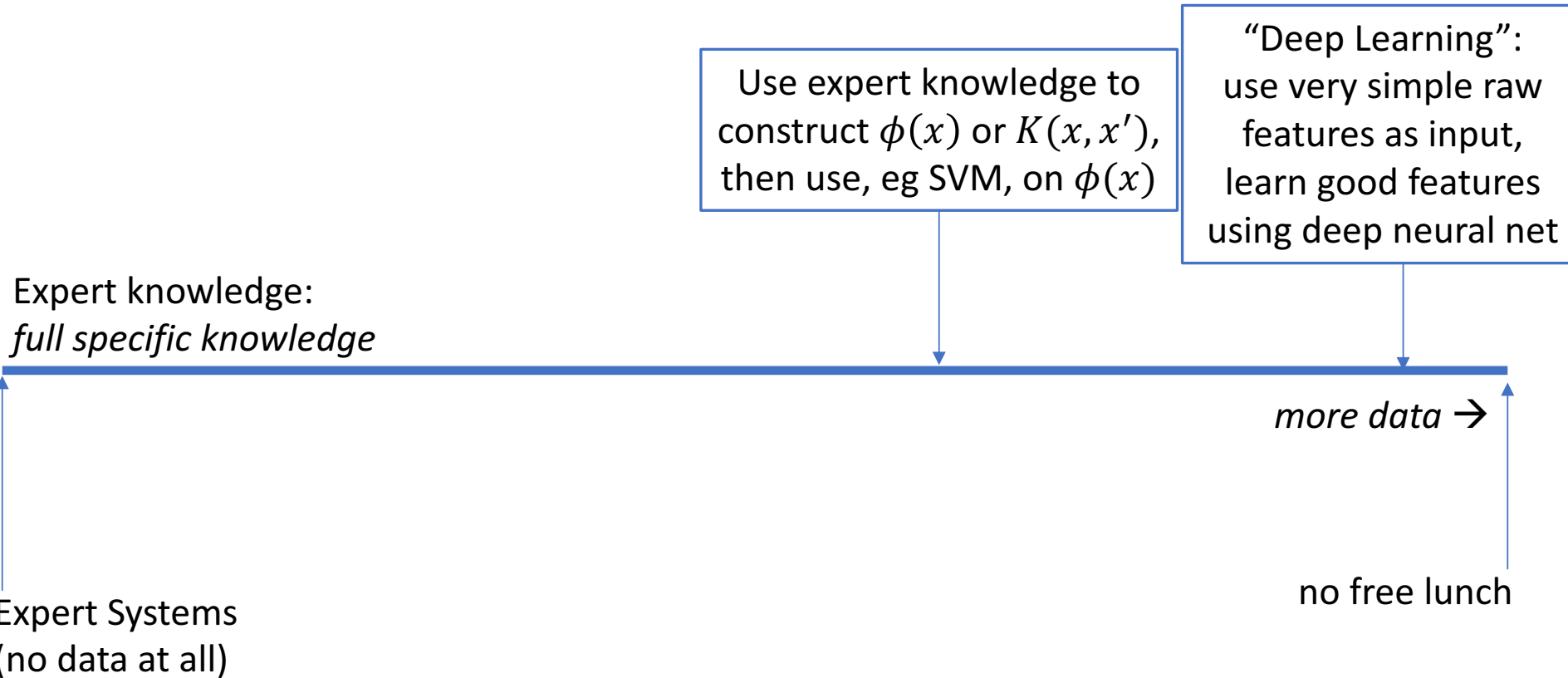


- Can think of hidden layer as “features” $\phi(x)$, then a linear predictor based on \mathbf{w} . $\phi(x)$
- “Feature Engineering” approach: design $\phi(\cdot)$ based on domain knowledge
- “Deep Learning” approach: learn features from data
- Multilayer networks with non-linear activations
 - more and more complex features

Multi-Layer Feature Learning



More knowledge or more learning



Neural networks as hypothesis class

- Hypothesis class specified by:

- Graph $G(V,E)$

- Activation function σ

} Based on architecture and fixed

- Weights \mathbf{W} , with weight $\mathbf{W}[u \rightarrow v]$ for each edge $u \rightarrow v \in E$

$$\mathcal{H} = \{ f_{G(V,E),\sigma,\mathbf{W}} \mid \mathbf{W}: E \rightarrow \mathbb{R} \}$$

- Expressive power:

$$\{ f \mid f \text{ computable in time } T \} \subseteq \mathcal{H}_{G(V,E),\text{sign}} \quad \text{with } |E| = O(T^2)$$

- Computation: empirical risk minimization

$$\hat{\mathbf{W}} = \arg \min_{\mathbf{W}} \sum_{i=1}^N \ell(f_{G(V,E),\sigma,\mathbf{W}}(\mathbf{x}^{(i)}), y^{(i)})$$

- Highly non-convex problem, even if *loss* ℓ is convex
- Hard to minimize over even tiny neural networks are hard

So how do we learn?

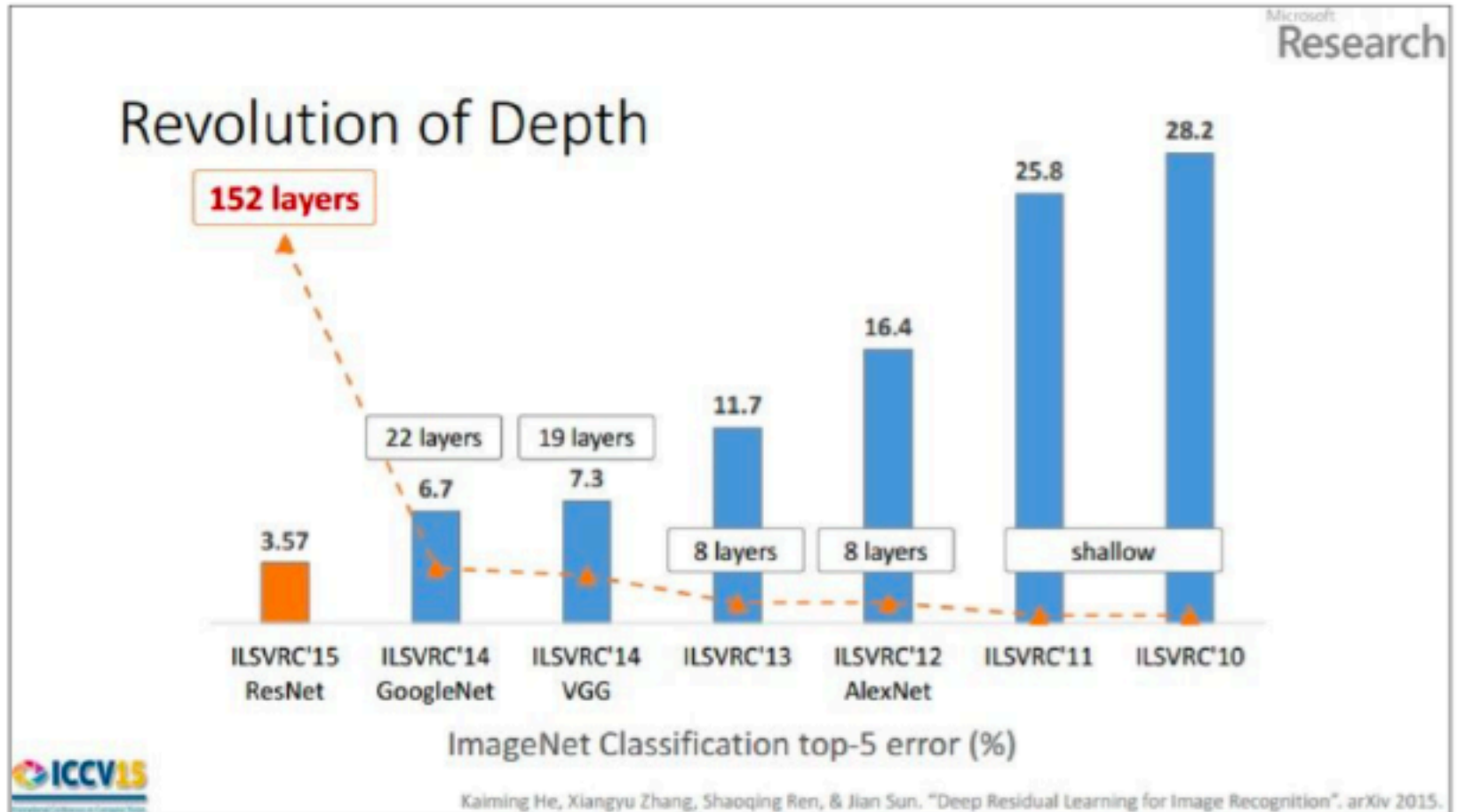
$$\widehat{\mathbf{W}} = \arg \min_{\mathbf{W}} \sum_{i=1}^N \ell(f_{G(V,E),\sigma,\mathbf{W}}(\mathbf{x}^{(i)}), y^{(i)})$$

- Stochastic gradient descent: for random $(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{S}$
$$\mathbf{W}^{(t+1)} \leftarrow \mathbf{W}^{(t)} - \eta^{(t)} \nabla \ell \left(f_{G(V,E),\sigma,\mathbf{W}^{(t)}}(\mathbf{x}^{(i)}), y^{(i)} \right)$$

(Even though its not convex)
- How do we efficiently calculate $\nabla \ell \left(f_{G(V,E),\sigma,\mathbf{W}^{(t)}}(\mathbf{x}^{(i)}), y^{(i)} \right)$?
 - Karl will tell you!
- Now a brief detour into history and resurrection of NNs

Imagenet challenge – object classification

1000 kinds of objects.



(slide from Kaiming He's recent presentation)

Object detection

PASCAL VOC Object Detection

	bicycle	bus	car	motorbike	person	20 class average
2007	36.9	23.2	34.6	27.6	21.3	17.1
2008	42.0	23.2	32.0	38.6	42.0	22.9
2009	46.8	43.8	37.2	42.0	41.5	27.9
2010	54.3	54.2	49.1	51.6	47.5	36.8
2011	58.1	57.6	54.4	58.3	51.6	40.9
2012	54.5	57.1	49.3	59.4	46.1	41.1
2013 DNN	56.3	51.4	48.7	59.8	44.4	43.2
2014 DNN						63.8
2015 ResNet	88.4	86.3	87.8	89.6	90.9	83.8
2016 ResNet						86

History of Neural Networks

- 1940s-70s:
 - Inspired by learning in the brain, and as a model for the brain (Pitts, Hebb, and others)
 - Various models, directed and undirected, different activation and learning rules
 - Perceptron Rule (Rosenblatt), Problem of XOR, Multilayer perceptron (Minsky and Papert)
 - Backpropagation (Werbos 1975)
- 1980s-early 1990s:
 - Practical Backprop (Rumelhart, Hinton et al 1986) and SGD (Bottou)
 - Relationship to distributed computing; “Connectionism”
 - Initial empirical success
- 1990s-2000s:
 - Lost favor to implicit linear methods: SVM, Boosting
- 2000-2010s:
 - revival of interest (CIFAR groups)
 - ca. 2005: layer-wise pretraining of deepish nets
 - progress in speech and vision with deep neural nets
- 2010s:
 - Computational advances allow training HUGE networks
 - ...and also a few new tricks
 - Krizhevsky et al. win ImageNet
 - Empirical success and renewed interest

Deep learning - today

State of the art performance in several tasks and are actively deployed in real systems

- Computer vision
- Speech recognition
- Machine translation
- Dialog systems
- Computer games
- Information retrieval