

# Introduction to Neural Networks

# Objectives

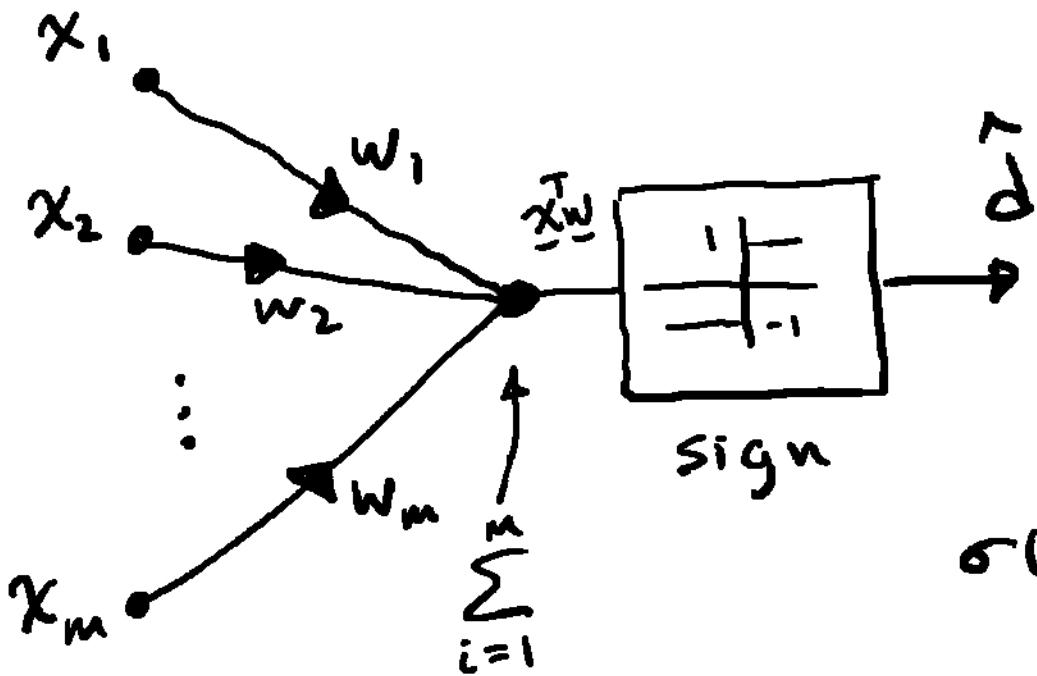
- relate neural networks to linear classifiers
- define structure of multilayer neural network
- overview procedure for training neural networks

# The "neuron" generalizes a linear classifier 2

Feature:  $\underline{x}^T = [x_1 \ x_2 \ \dots \ x_n]$  Weights:  $\underline{w}^T = [w_1 \ w_2 \ \dots \ w_m]$

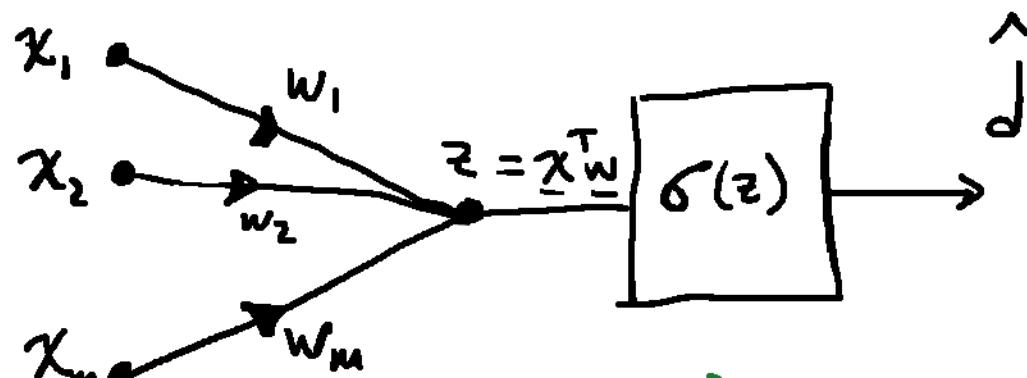
Linear classifier

$$\hat{d} = \text{sign}(\underline{x}^T \underline{w})$$



Neuron

$$\hat{d} = \sigma(\underline{x}^T \underline{w})$$



Common  $\sigma(z)$

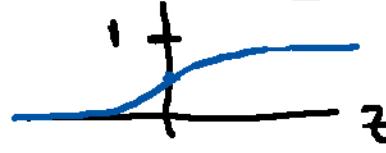
ReLU

$$\sigma(z) = \max\{0, z\}$$

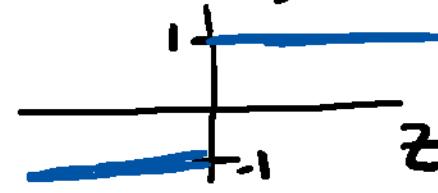


Logistic

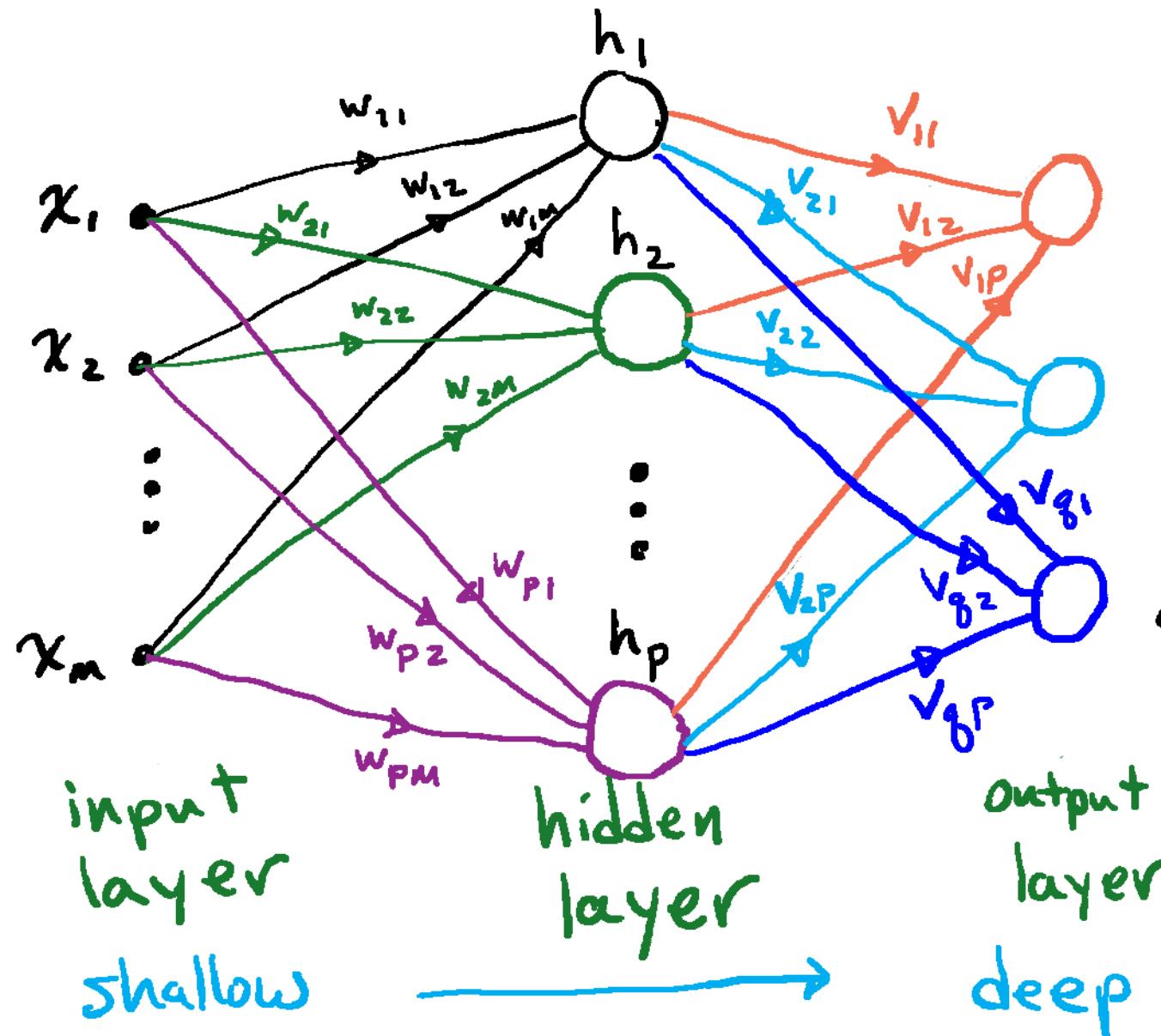
$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



$$\sigma(z) = \text{sign}(z)$$



# A neural network is a network of neurons 3



multiple outputs →  
solve multiple problems

$$\hat{d}_1$$

$$\hat{d}_2$$

$$\hat{d}_Q$$

$$h_q = \sigma \left( \sum_{j=1}^m w_{qj} x_j \right)$$

$$\hat{d}_m = \sigma \left( \sum_{n=1}^P v_{mn} h_n \right)$$

$$= \sigma \left( \sum_{n=1}^P v_{m,n} \sigma \left( \sum_{j=1}^m w_{nj} x_j \right) \right)$$

"deep" learning →  
many hidden layers

Two issues must be addressed to use NN's 4

1) Network structure: number of layers, number of hidden nodes in each layer

Open question. Universal approximation theorem (1991):

Three layer network can approximate any function arbitrarily well given enough hidden nodes and the right weights

2) Choosing the weights

Stochastic gradient descent and backpropagation

Nonconvex  $\rightarrow$  local minima

Backpropagation updates each layer in sequence<sup>5</sup>

- work back from deep (output) to shallow (input)

$$\min_{w_{qj}, v_{lj}, \dots} \sum_{i=1}^N \sum_{q=1}^Q \frac{1}{2} (\hat{d}_{i,q} - d_{i,q})^2 \quad N \text{ training samples, } Q \text{ outputs}$$

- 1) Initial guess on  $w_{qj}, v_{lj}$  (etc)
- 2) Randomly choose training sample  $i$
- 3) Calculate  $h_{i,p}, \hat{d}_{i,q}$
- 4) Gradient descent update  $v_{lj}$ , then  $w_{qj}$  (deep to shallow)

Chain rule is key for deriving gradients

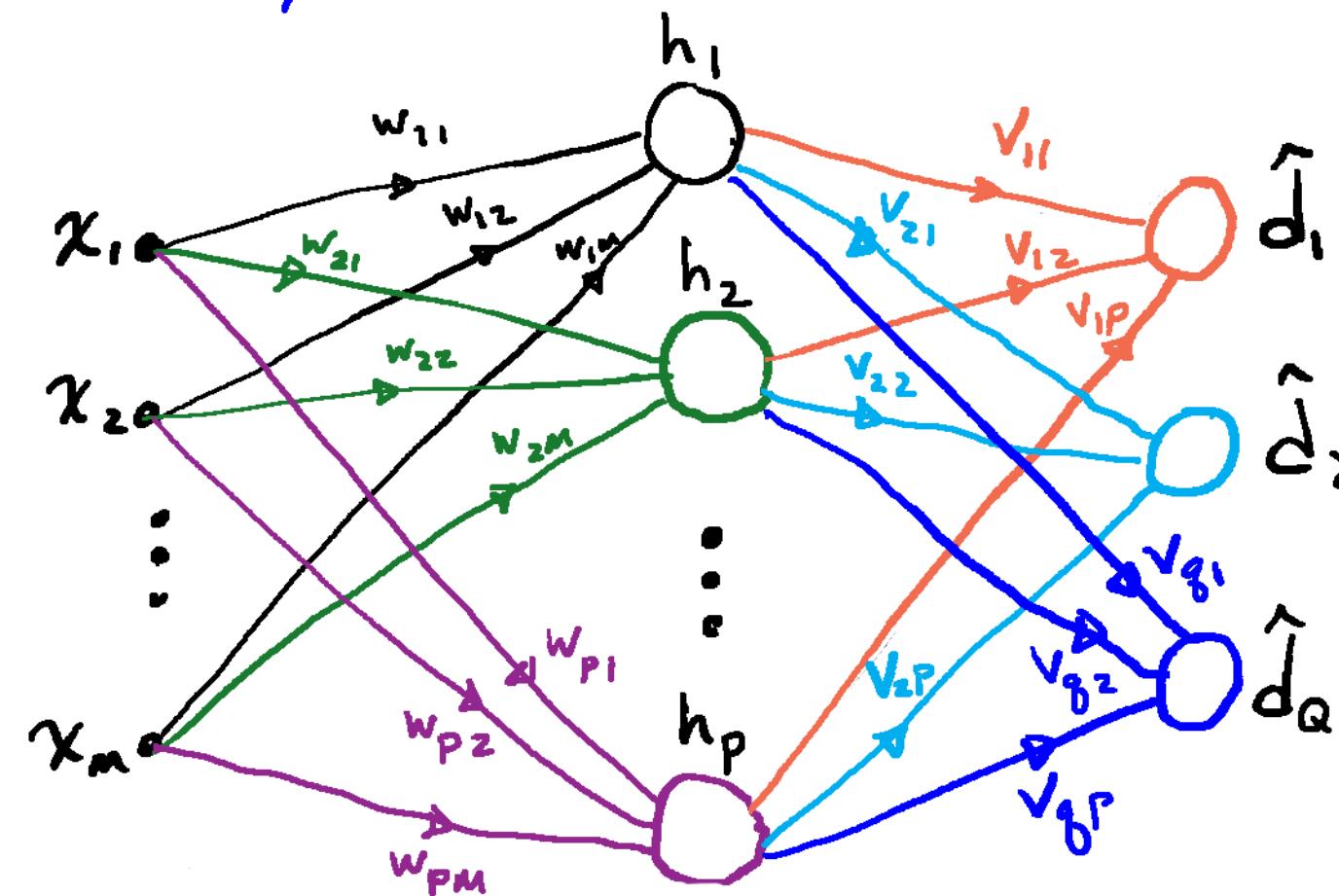
Copyright 2019  
Barry Van Veen

# The Backpropagation Algorithm for Training Neural Networks

# Objectives

- Method for learning weights in a neural network
- Apply stochastic gradient descent (SGD)
- Find gradients of squared error loss with respect to weights

Backpropagation uses SGD to train the weights 2



$$\hat{d}_q = \sigma \left( \sum_{n=1}^p v_{q,n} h_n \right)$$

$$h_n = \sigma \left( \sum_{k=1}^m w_{n,k} x_k \right) \quad (= \frac{1}{1+e^{-z}} \text{ here})$$

$N$  training samples

$$(x_1^i, x_2^i \dots x_m^i, d^i, i=1, 2, \dots N)$$

- 1) initialize  $w_{n,k}$ ,  $v_{q,n}$
- 2) for  $t=1, 2, 3, \dots$

Pick  $i_t \in \{1, 2, \dots, N\}$  randomly  
Compute  $h_j^{i_t}, \hat{d}_k^{i_t}$

$$v_{k,t}^{t+1} = v_{k,t}^t - \alpha_t \frac{\partial f^{i_t}}{\partial v_{k,t}}$$

$$w_{m,j}^{t+1} = w_{m,j}^t - \alpha_t \frac{\partial f^{i_t}}{\partial w_{m,j}}$$

Converged?

Gradients are obtained using the chain rule

Squared error loss:  $f(\underline{v}, \underline{w}) = \sum_{i=1}^N \frac{1}{2} \sum_{g=1}^G (\hat{d}_g^i - d_g^i)^2$

$$\frac{\partial f^i}{\partial v_{k,e}} : \left( \frac{\partial f^i}{\partial \hat{d}_k^i} \right) \frac{\partial \hat{d}_k^i}{\partial v_{k,e}}$$



$$\frac{\partial f^i}{\partial v_{k,e}} = (\hat{d}_k^i - d_k^i) \cdot \sigma' \left( \sum_{n=1}^P h_n^i v_{k,n} \right) \cdot \frac{\partial}{\partial v_{k,e}} \left( \sum_{n=1}^P h_n^i v_{k,n} \right)$$

$$\frac{\partial f^i}{\partial v_{k,e}} = (\hat{d}_k^i - d_k^i) \hat{d}_k^i (1 - \hat{d}_k^i) h_e^i$$

$$= \delta_k^i h_e^i$$

$$\begin{aligned} \sigma'(z) &= \frac{1}{(1+e^{-z})^2} e^{-z} = \frac{1}{1+e^{-z}} \left( 1 - \frac{1}{1+e^{-z}} \right) \\ &= \sigma(z) (1 - \sigma(z)) \end{aligned}$$

$$\begin{aligned} \text{But } \sigma \left( \sum_{n=1}^P h_n^i v_{k,n} \right) (1 - \sigma \left( \sum_{n=1}^P h_n^i v_{k,n} \right)) \\ &= \hat{d}_k^i (1 - \hat{d}_k^i) \end{aligned}$$

layer input  
"error" at output

$$v_{k,e}^{t+1} = v_{k,e}^t - \alpha_t \delta_k^i h_e^i$$

Use a similar approach for  $w_{m,j}$

$$\frac{\partial f^i}{\partial w_{m,j}} = \sum_{g=1}^Q \frac{\partial f^i}{\partial \hat{d}_g^i} \cdot \frac{\partial \hat{d}_g^i}{\partial h_m^i} \frac{\partial h_m^i}{\partial w_{m,j}}$$

$\frac{\partial f^i}{\partial \hat{d}_g^i} = (\hat{d}_g^i - d_g^i)$

loss / output      output / hidden      hidden / prev. layer weights

$$\frac{\partial \hat{d}_g^i}{\partial h_m^i} = \frac{\partial}{\partial h_m^i} \sigma \left( \sum_{n=1}^P h_n^i v_{g,n}^i \right) = \sigma' \left( \sum_{n=1}^P h_n^i v_{g,n}^i \right) \cdot v_{g,m}^i = \hat{d}_g^i (1 - \hat{d}_g^i) v_{g,m}^i$$

$$\frac{\partial h_m^i}{\partial w_{m,j}} = \frac{\partial}{\partial w_{m,j}} \sigma \left( \sum_{k=1}^N w_{m,k} x_k^i \right) = \sigma' \left( \sum_{k=1}^N w_{m,k} x_k^i \right) x_j^i = h_m^i (1 - h_m^i) x_j^i$$

$$\frac{\partial f^i}{\partial w_{m,j}} = \sum_{g=1}^Q (\hat{d}_g^i - d_g^i) \hat{d}_g^i (1 - \hat{d}_g^i) v_{g,m}^i h_m^i (1 - h_m^i) x_j^i = \underline{\gamma_m^i x_j^i}$$

$$w_{m,j}^{t+1} = w_{m,j}^t - \alpha^t \gamma_m^i x_j^i$$

layer input      backpropagated "error"

# Backpropagation involves simple computations

5

Initialize  $w_{m,j}^0, v_{k,e}^0$

Iterate for  $t=0,1,2,\dots$

SGD

choose  $i_t \in \{1, 2, \dots, N\}$   
at random

Forward  
Net

compute  $h_m^{i_t}, \hat{d}_g^{i_t}$  from  
 $x_k^{i_t}, w_{m,j}^t, v_{k,e}^t$

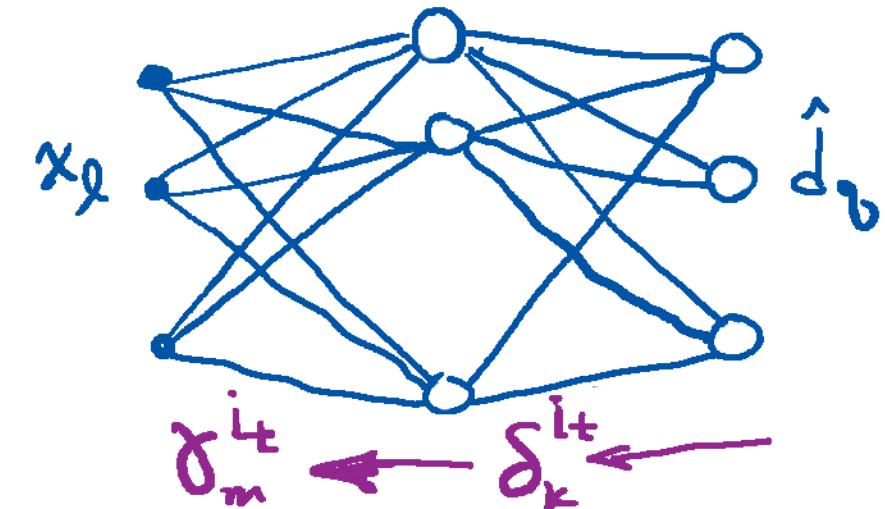
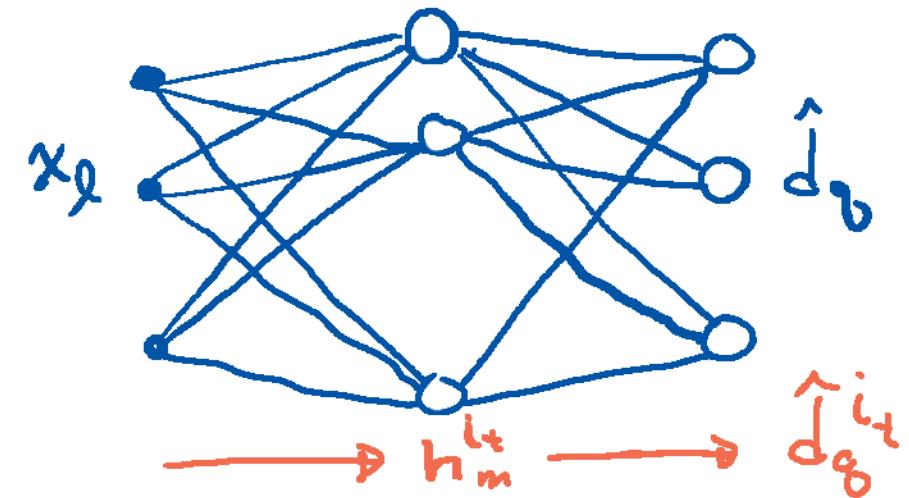
Backward  
Updates

$$\delta_k^{i_t} = (\hat{d}_k^{i_t} - d_k^{i_t}) \hat{d}_k^{i_t} (1 - \hat{d}_k^{i_t})$$

$$v_{k,e}^{t+1} = v_{k,e}^t - \alpha_t \delta_k^{i_t} h_g^{i_t}$$

$$\delta_m^{i_t} = \sum_{g=1}^G \delta_g^{i_t} v_{g,m}^{t+1} h_m^{i_t} (1 - h_m^{i_t})$$

$$w_{m,j}^{t+1} = w_{m,j}^t - \alpha_t \delta_m^{i_t} x_j^{i_t}$$



Variations follow the same pattern<sup>6</sup>

- other activation functions
- deeper networks

Cost is a non convex function of weights

- converge to local minima
- empirical "tricks" e.g. batch SGD, normalization to accelerate convergence

Can add regularization

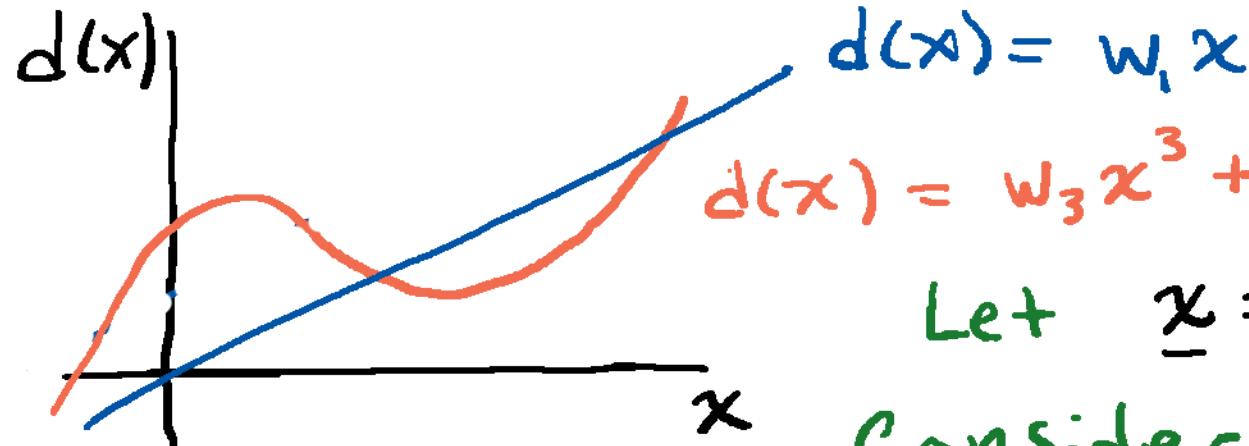
Copyright 2019  
Barry Van Veen

# Kernel Regression

# Objectives

- Why use higher-dimensional feature spaces
- Reformulate regression in terms of kernels
- Popular kernels
- Cautions and considerations

# Higher dimensional feature spaces extend regression



$$d(x) = w_3 x^3 + w_2 x^2 + w_1 x + w_0$$

Let  $\underline{x} = [x, x_2 \dots x_n]^T \in \mathbb{R}^n$

Consider  $d(\underline{x}) = \underline{\phi}^T(\underline{x}) \underline{w}$ ,  $\phi(\underline{x}) \in \mathbb{R}^P$   
 $P > n$

Example:  $\underline{x} = [x, x_2]^T$ ,  $\underline{\phi}^T(\underline{x}) = [x^2, x^2, \sqrt{2}x, x_2, x, x_1, 1]$

Finding  $\underline{w}$ : "training" data  $\underline{x}^i, d^i, i=1, 2, \dots, N$

$$\min_{\underline{w}} \sum_{i=1}^N (d^i - \underline{\phi}^T(\underline{x}^i) \underline{w})^2 + \lambda \|\underline{w}\|_2^2 \quad (\text{Ridge})$$

$$\underline{d} = [d^1 \ d^2 \ \dots \ d^N]^T$$

$$\underline{\Phi} = [\underline{\phi}(\underline{x}^1) \ \underline{\phi}(\underline{x}^2) \ \dots \ \underline{\phi}(\underline{x}^N)]^T \quad (N \times P) \Rightarrow \underline{w} = (\underline{\Phi}^T \underline{\Phi} + \lambda \underline{I})^{-1} \underline{\Phi}^T \underline{d}$$

$$\min_{\underline{w}} \|\underline{d} - \underline{\Phi} \underline{w}\|_2^2 + \lambda \|\underline{w}\|_2^2$$

Regression is a weighted sum of "Kernels" 3

$$d(\underline{x}) = \underline{\phi}^T(\underline{x}) \underline{w} = \underline{\phi}^T(\underline{x}) (\underbrace{\underline{\Phi}^T \underline{\Phi} + \lambda \underline{I}}_{P \times P})^{-1} \underline{\Phi}^T \underline{d}$$

Matrix identity:  $(\underline{\Phi}^T \underline{\Phi} + \lambda \underline{I})^{-1} \underline{\Phi}^T = \underline{\Phi}^T (\underline{\Phi} \underline{\Phi}^T + \lambda \underline{I})^{-1}$  (activity)

Thus  $d(\underline{x}) = \underline{\phi}^T(\underline{x}) \underbrace{\underline{\Phi}^T (\underline{\Phi} \underline{\Phi}^T + \lambda \underline{I})^{-1} \underline{d}}_{N \times N}$

Note:  $\begin{bmatrix} \underline{\Phi} & \underline{\Phi}^T \end{bmatrix}_{i,j} = \underline{\phi}^T(\underline{x}^i) \underline{\phi}(\underline{x}^j)$  } Define "Kernel"  
 $K(\underline{u}, \underline{v}) = \underline{\phi}^T(\underline{u}) \underline{\phi}(\underline{v})$

$$\begin{bmatrix} \underline{\phi}^T(\underline{x}) & \underline{\Phi}^T \end{bmatrix}_j = \underline{\phi}^T(\underline{x}) \underline{\phi}(\underline{x}^j)$$

Let  $\underline{\alpha} = [\alpha_1, \dots, \alpha_N]^T$   
 $= (\underline{\Phi} \underline{\Phi}^T + \lambda \underline{I})^{-1} \underline{d}$

$$d(\underline{x}) = \sum_{i=1}^N \alpha_i \underline{\phi}^T(\underline{x}) \underline{\phi}(\underline{x}^i) = \sum_{i=1}^N \alpha_i K(\underline{x}, \underline{x}^i)$$

Kernel methods find  $d(\underline{x})$  without computing  $\phi(\underline{x})^T \underline{d}$

$$d(\underline{x}) = \sum_{i=1}^N \alpha_i K(\underline{x}, \underline{x}^i) \quad \underline{\alpha} = (\underline{\Phi} \underline{\Phi}^T + \lambda \underline{\mathbb{I}})^{-1} \underline{d} \Rightarrow \underline{\mathbb{K}} = \underline{\Phi} \underline{\Phi}^T$$

$$[\underline{\mathbb{K}}]_{ij} = \underline{\phi}^T(\underline{x}^i) \underline{\phi}(\underline{x}^j) = K(\underline{x}^i, \underline{x}^j)$$

$\underline{\mathbb{K}}$  can be computed efficiently!

Ex: Monomials of degree 3  $\underline{\phi}(\underline{x}) \rightarrow \underbrace{x_1^3, x_1^{3-1}x_2, \dots}_{\dots}, \underbrace{x_3^3 x_6^2 x_8^3 \dots}_{\dots}$

$$K(\underline{u}, \underline{v}) = \underbrace{\underline{\phi}^T(\underline{u}) \underline{\phi}(\underline{v})}_{O(P)} = \underbrace{(\underline{u}^T \underline{v})^3}_{O(M)} \text{ (activity)} \quad P = \frac{(g+M-1)!}{g! (M-1)!} \text{ terms}$$

Suppose  $M=10, g=5 \rightarrow P \sim 2000$  computing  $O(P)$  vs  $O(M)$

$M=100, g=5 \rightarrow P \sim 10^8$  memory  $O(NP)$  vs  $O(N^2)$

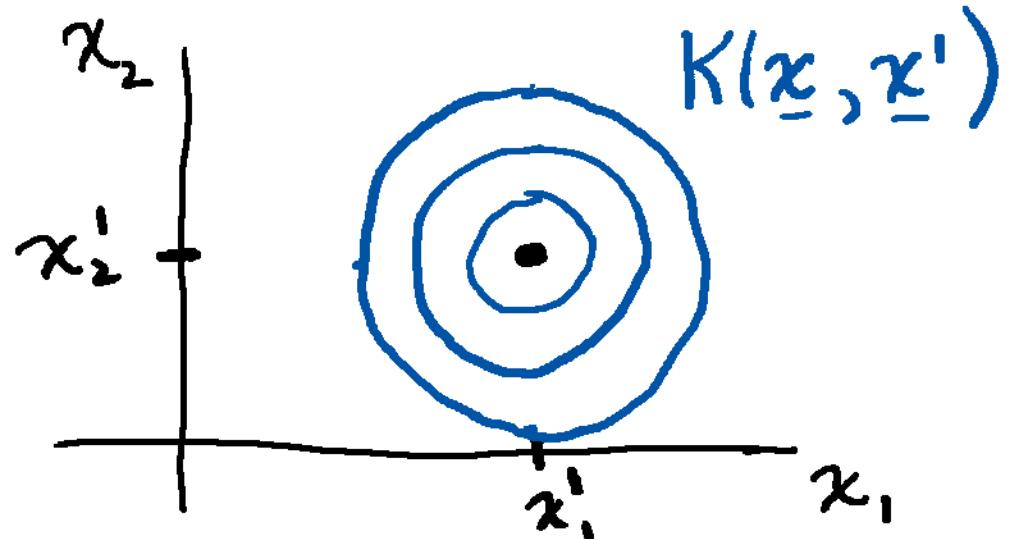
Popular kernels depend on similarity of  $\underline{u}, \underline{v}$  5

$$\underline{u}^T \underline{v} = \|\underline{u}\|_2 \|\underline{v}\|_2 \cos \theta$$

Monomials of degree  $q$ :  $K(\underline{u}, \underline{v}) = (\underline{u}^T \underline{v})^q$

Polynomials up to degree  $q$ :  $K(\underline{u}, \underline{v}) = (\underline{u}^T \underline{v} + 1)^q$

Gaussian/radial Kernel:  $K(\underline{u}, \underline{v}) = \exp\left\{-\frac{\|\underline{u} - \underline{v}\|_2^2}{2\sigma^2}\right\}$



- No explicit  $\phi(\underline{x})$
- All polynomial orders
- smoothness controlled by  $\sigma$

# Kernel regression considerations

6

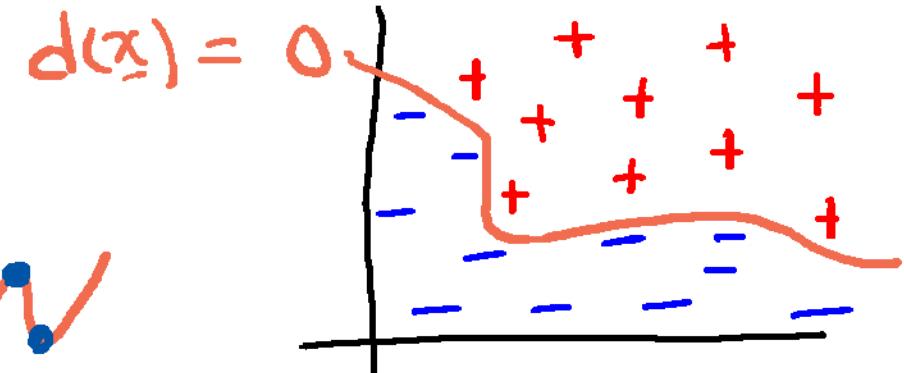
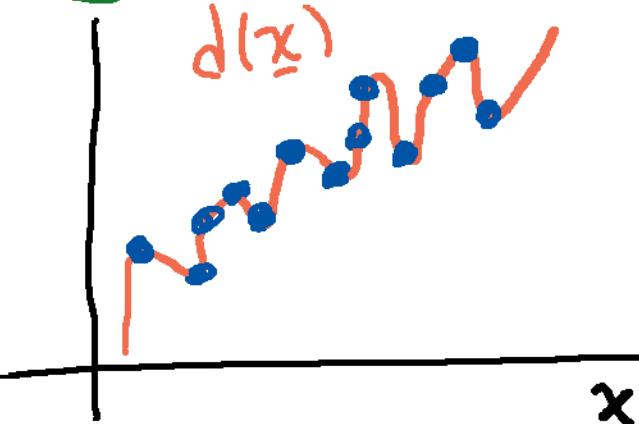
$$d(\underline{x}) = \underline{\phi}^T(\underline{x}) \underline{w} \quad \text{vs} \quad d(\underline{x}) = \sum_{i=1}^N \alpha_i K(\underline{x}, \underline{x}^i)$$

- Store and compute  $\underline{\alpha}$  ( $N \times 1$ ) vs  $\underline{w}$  ( $P \times 1$ )

- Binary classification  $\text{sign}\{d(\underline{x})\}$

- Avoid "overfitting" with  
high-D feature  
spaces

(cross-validation)



Copyright 2019  
Barry Van Veen

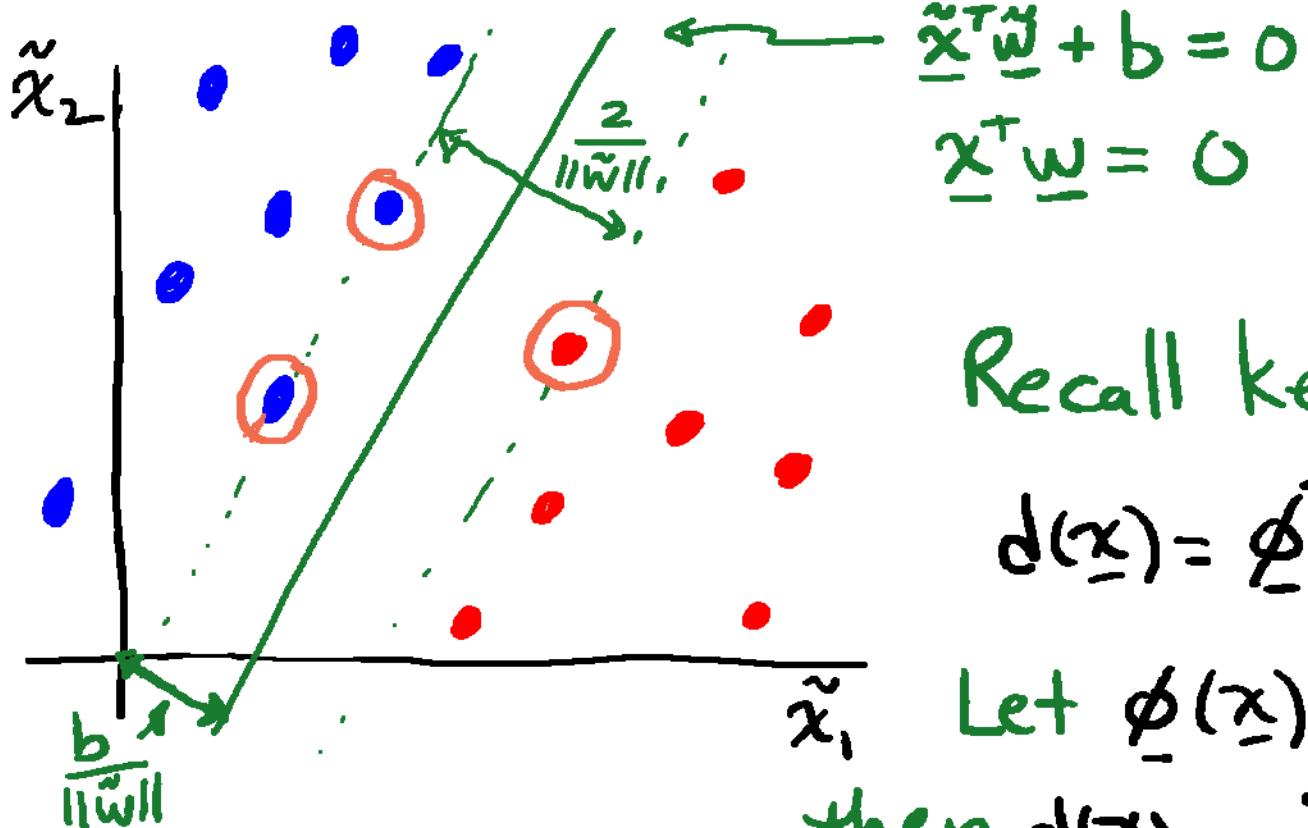
# **Kernel Based Support Vector Machines**

# Objectives

1

- reformulate linear max margin classifier in terms of support vectors
- derive kernel version of hinge loss with ridge regression
- summarize features of support vector machines

## Support vectors define max-margin classifier 2



$\underline{w} = [\tilde{w}^T \ b]^T$  depends only  
on the support vectors

Recall kernel regression

$$d(\underline{x}) = \underline{\phi}^T(\underline{x}) \underline{w} = \sum_{j=1}^N \alpha_j K(\underline{x}, \underline{x}^j)$$

Let  $\underline{\phi}(\underline{x}) = \underline{x} \Rightarrow K(\underline{x}, \underline{x}^j) = \underline{x}^T \underline{x}^j$   
then  $d(\underline{x}) = \sum_{j=1}^N \alpha_j \underline{x}^T \underline{x}^j = \underline{x}^T \underbrace{\sum_{j=1}^N \alpha_j \underline{x}^j}_{\underline{w}}$

So  $\underline{w} = \sum_{j=1}^N \alpha_j \underline{x}^j$

All  $\alpha_j = 0$  except Support vectors!

Use kernels for nonlinear decision boundaries<sup>3</sup>

High-dimensional feature space:  $\underline{x} \rightarrow \underline{\phi}(\underline{x})$

e.g.,  $\underline{\phi}(\underline{x}) = [x_1^2 \ x_2^2 \ \dots \ x_2 x_4 \ \dots \ x_{n-1} \ x_n \ 1]$

$$\hat{d}(\underline{x}) = \text{sign}(\underline{\phi}^T(\underline{x}) \underline{w})$$

Hinge loss with ridge regression

$$\min_{\underline{w}} \sum_{i=1}^N (1 - d^i \underline{\phi}^T(\underline{x}^i) \underline{w})_+ + \lambda \|\underline{w}\|_2^2$$



Claim:  $\underline{w} = \sum_{j=1}^N \underline{\phi}(\underline{x}^j) \alpha_j$  (proof in notes)

Kernel "trick" replaces  $\underline{\phi}^T(\underline{x}^i) \underline{\phi}(\underline{x}^i)$  with  $K(\underline{x}^i, \underline{x}^i)$  4

Restate:  $\underline{w} = \sum_{j=1}^N \alpha_j \underline{\phi}(\underline{x}^j)$

Hinge loss with ridge regression

$$\min_{\underline{\alpha}} \sum_{i=1}^N \left( 1 - d^i \underline{\phi}^T(\underline{x}^i) \sum_{j=1}^N \alpha_j \underline{\phi}(\underline{x}^j) \right)_+ + \lambda \underbrace{\sum_{i=1}^N \alpha_i \underline{\phi}^T(\underline{x}^i)}_{\underline{w}^T} \underbrace{\sum_{j=1}^N \alpha_j \underline{\phi}(\underline{x}^j)}_{\underline{w}}$$

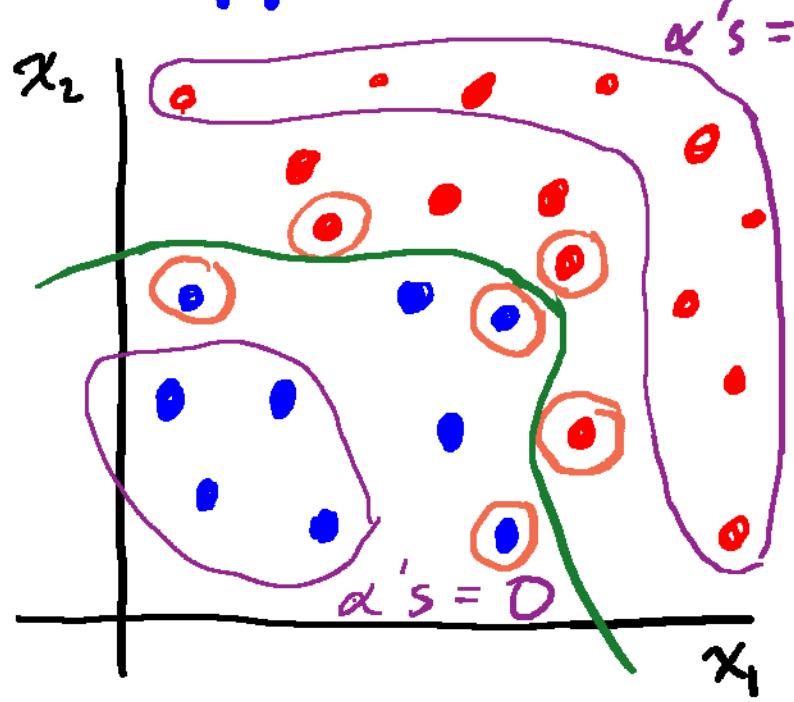
$$\min_{\underline{\alpha}} \sum_{i=1}^N \left( 1 - d^i \sum_{j=1}^N \alpha_j \underline{\phi}^T(\underline{x}^i) \underline{\phi}(\underline{x}^j) \right)_+ + \lambda \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j \underbrace{\underline{\phi}^T(\underline{x}^i) \underline{\phi}(\underline{x}^j)}_{K(\underline{x}^i, \underline{x}^j)}$$

Kernel "trick"

$$\text{SVM} \quad \min_{\underline{\alpha}} \sum_{i=1}^N \left( 1 - d^i \sum_{j=1}^N \alpha_j K(\underline{x}^i, \underline{x}^j) \right)_+ + \lambda \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j K(\underline{x}^i, \underline{x}^j)$$

# Support vector machines have sparse $\underline{\alpha}$

5

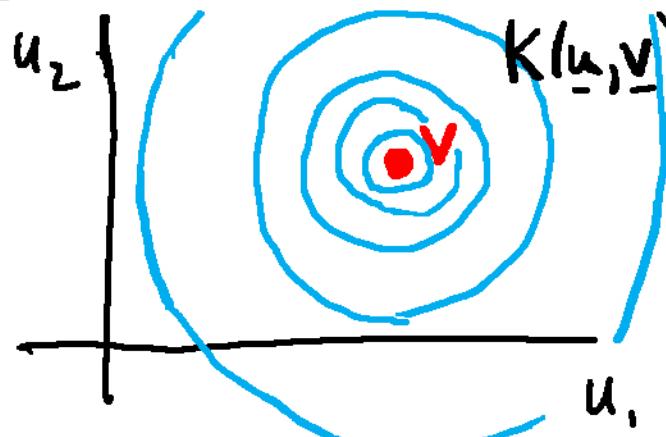


decision boundary

$$d(\underline{x}) = 0 = \underline{\phi}^T(\underline{x}) \underline{w} = \sum_{j=1}^N \alpha_j K(\underline{x}, \underline{x}_j)$$

Boundary (hinge loss) depends only on the support vectors

Ex: Gaussian kernels



$$K(\underline{u}, \underline{v}) = \exp \left\{ -\frac{\|\underline{u}-\underline{v}\|_2^2}{2\sigma^2} \right\}$$

Solve for  $\underline{\alpha}$  using gradient descent

Copyright 2019  
Barry Van Veen