

MARS: Insights from Carts, Insights for Nets

Trent McConaghy

Berlin Machine Learning Group
Feb 2, 2015

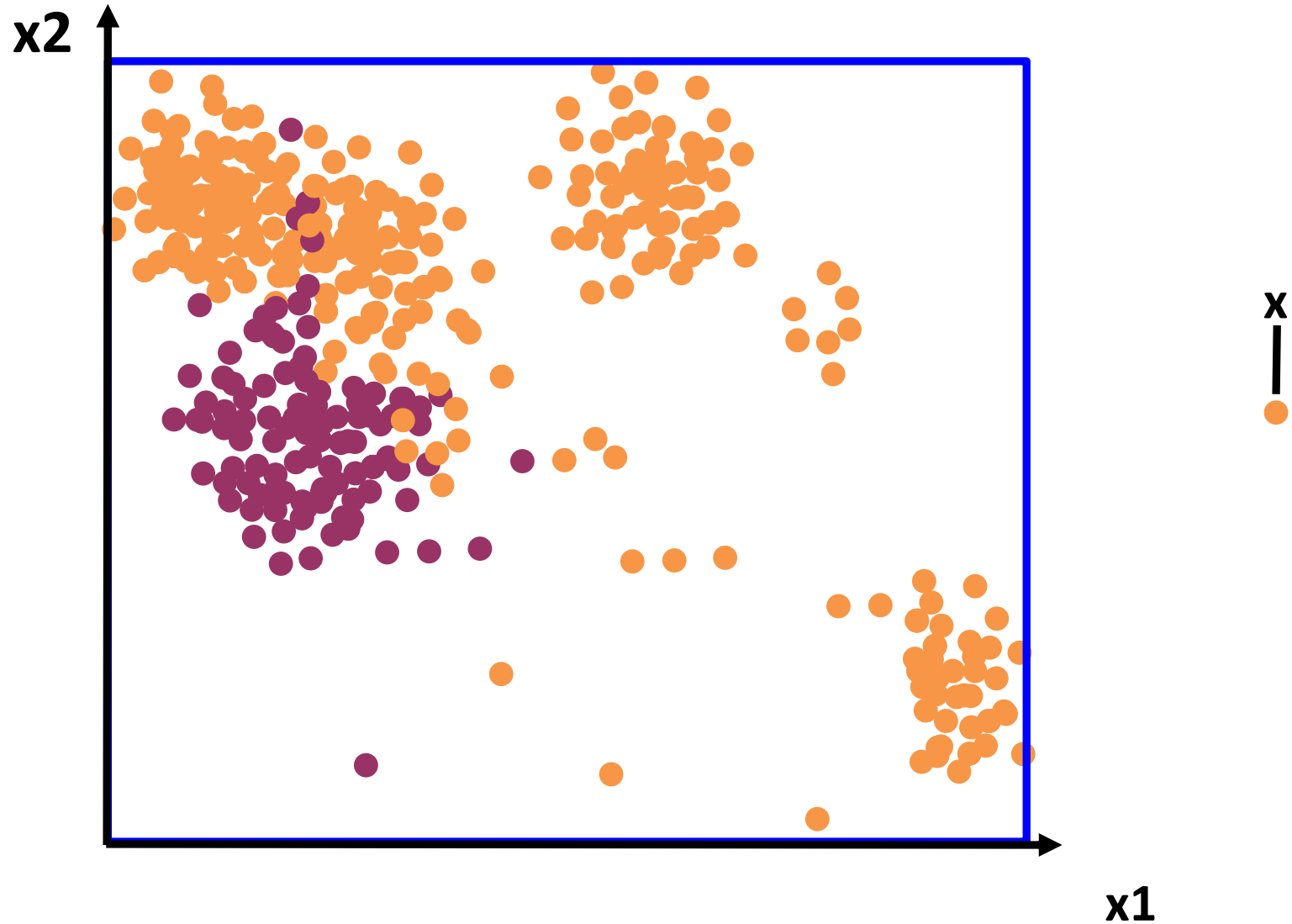
Outline

- Carts
- Mars
- Nets

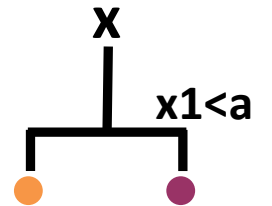
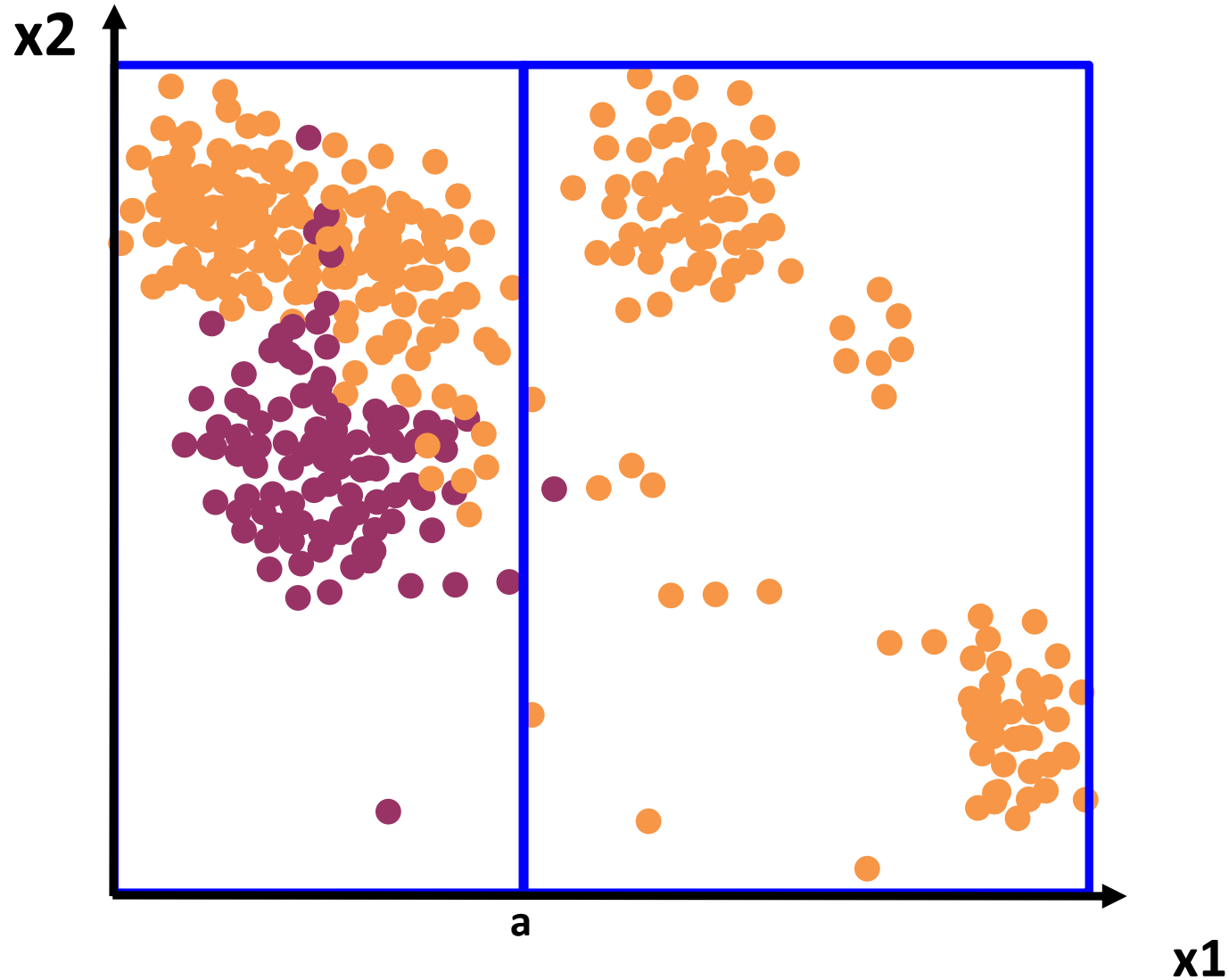
Outline

- Carts
- Mars
- Nets

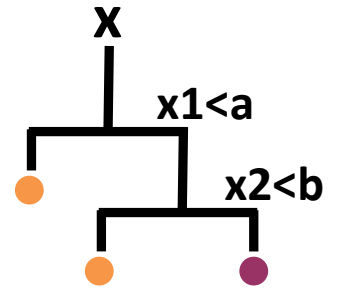
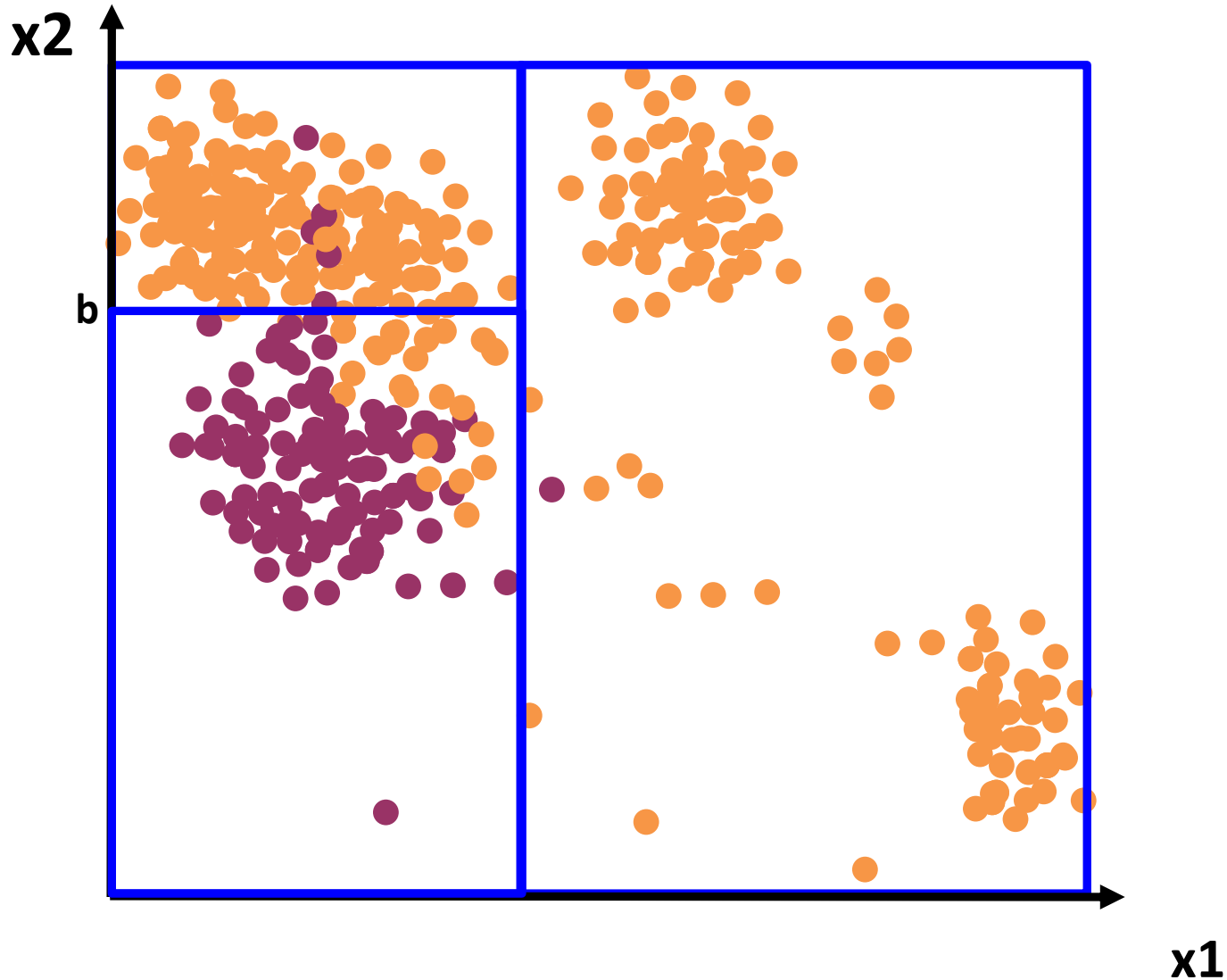
CART Classifier



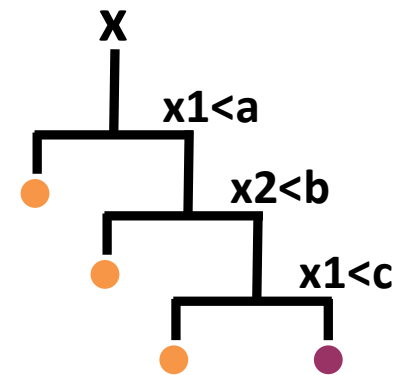
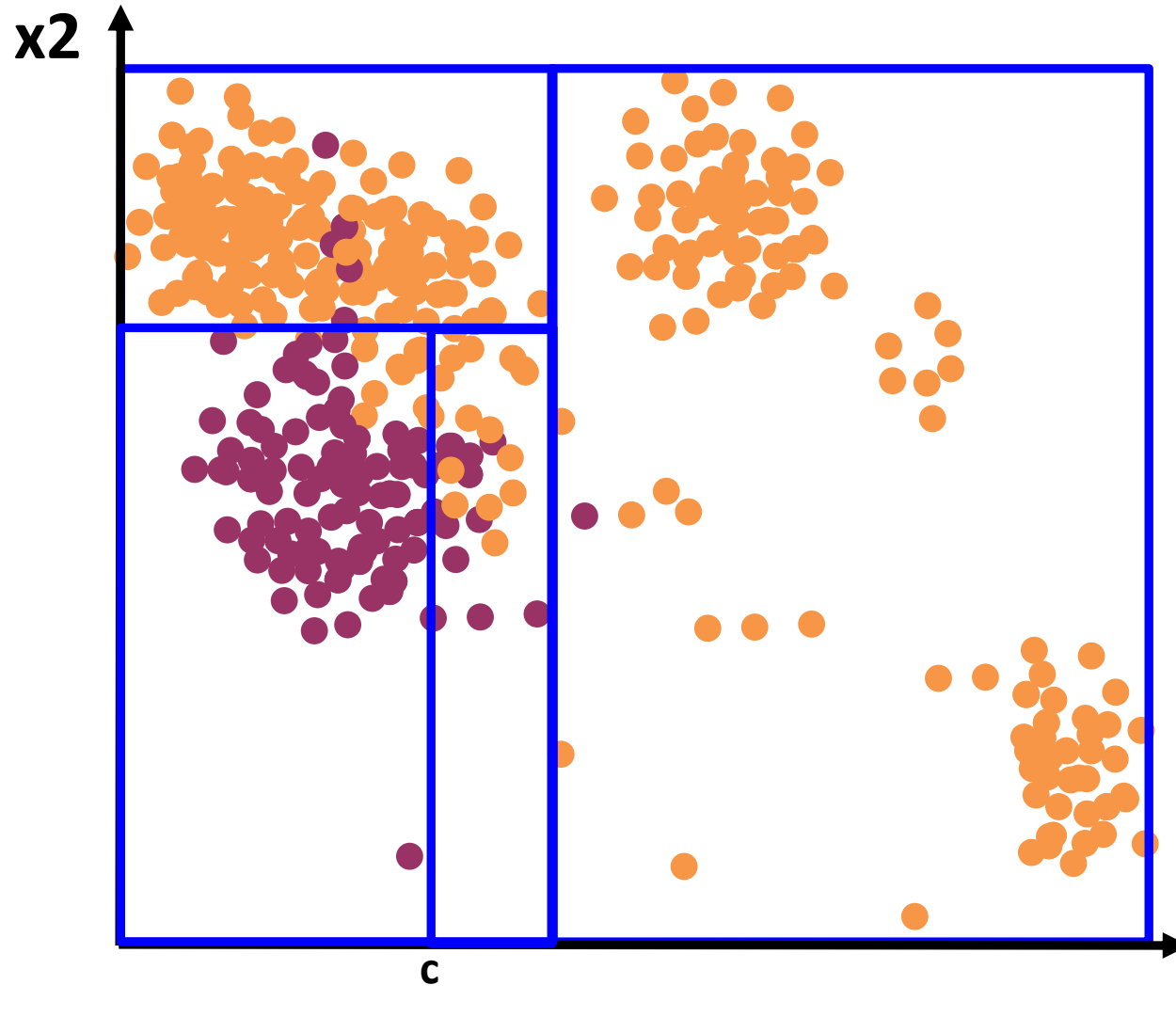
CART Classifier



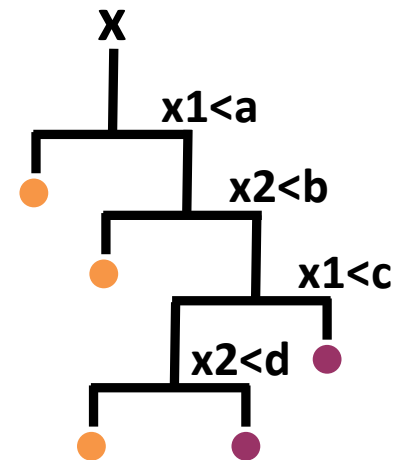
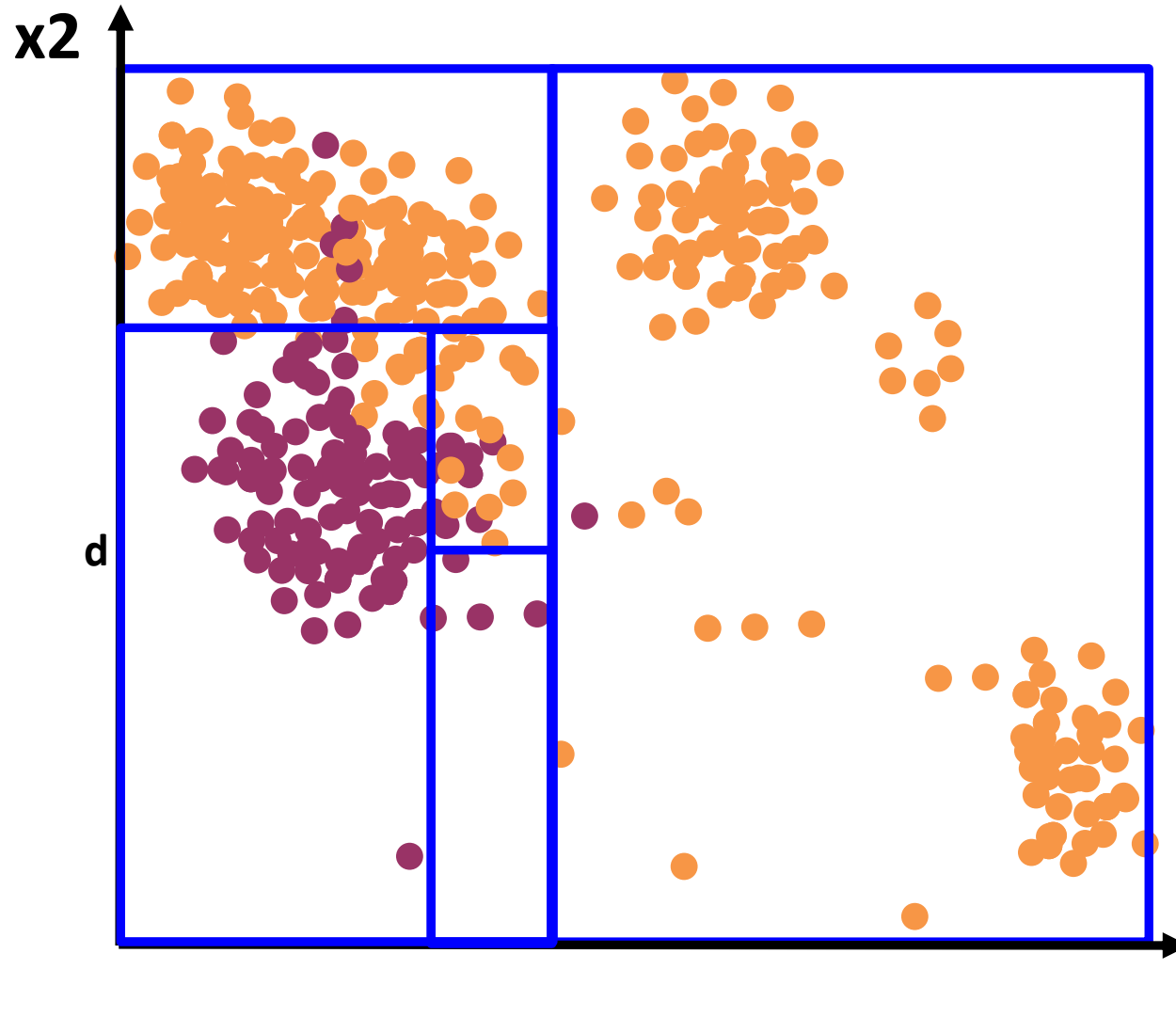
CART Classifier



CART Classifier



CART Classifier

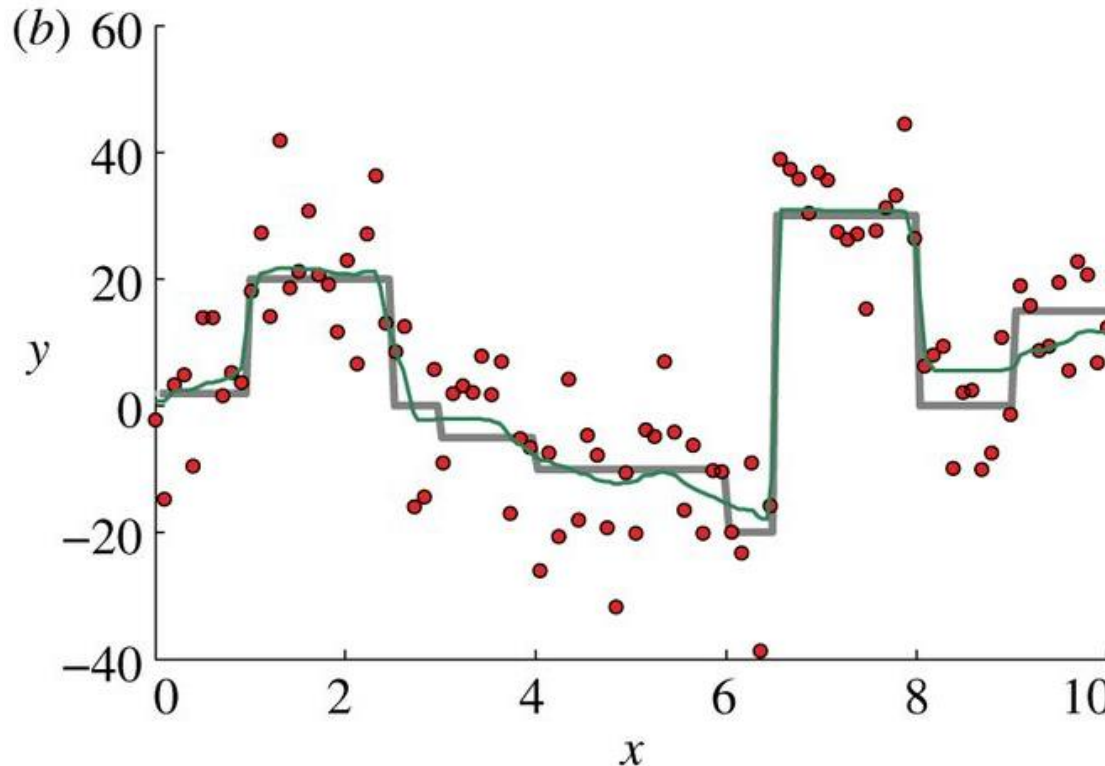


On CARTs

- Adaptively zooms in on the most important regions. Basis functions zero out the rest.
- Computationally cheap

CART Regressor

- Is a weighted sum of pairs of indicator functions
- = a piecewise constant model



CART Learning

$$B_1(\mathbf{x}) \leftarrow 1$$

For $M = 2$ to M_{\max} do: $\text{lof}^* \leftarrow \infty$ } iteratively add M_{\max} regions (basis funcs)

For $m = 1$ to $M - 1$ do:

For $v = 1$ to n do:

For $t \in \{x_{vj} | B_m(\mathbf{x}_j) > 0\}$

$$g \leftarrow \sum_{i \neq m} a_i B_i(\mathbf{x}) + a_m B_m(\mathbf{x}) H[+(x_v - t)] + a_M B_m(\mathbf{x}) H[-(x_v - t)]$$

$$\text{lof} \leftarrow \min_{a_1, \dots, a_M} \text{LOF}(g)$$

if $\text{lof} < \text{lof}^*$, then $\text{lof}^* \leftarrow \text{lof}$; $m^* \leftarrow m$; $v^* \leftarrow v$; $t^* \leftarrow t$ end if

end for

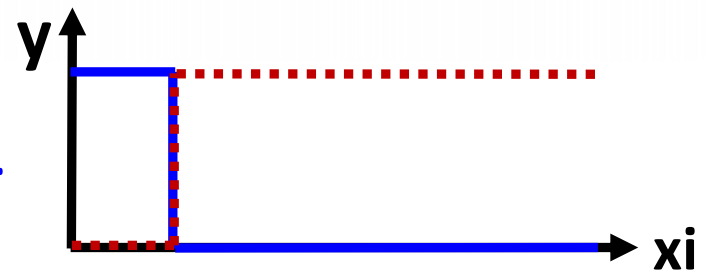
end for

end for

$B_M(\mathbf{x}) \leftarrow B_{m^*}(\mathbf{x}) H[-(x_{v^*} - t^*)]$ } add best split, in both directions

$B_{m^*}(\mathbf{x}) \leftarrow B_{m^*}(\mathbf{x}) H[+(x_{v^*} - t^*)]$

indicator function building blocks.
Discontinuous..



CARTs Live On

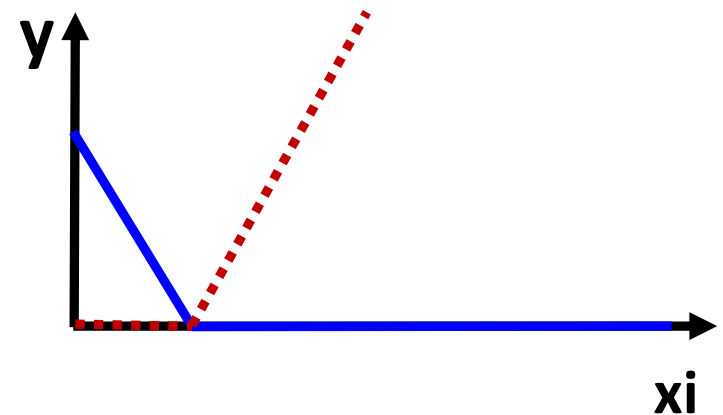
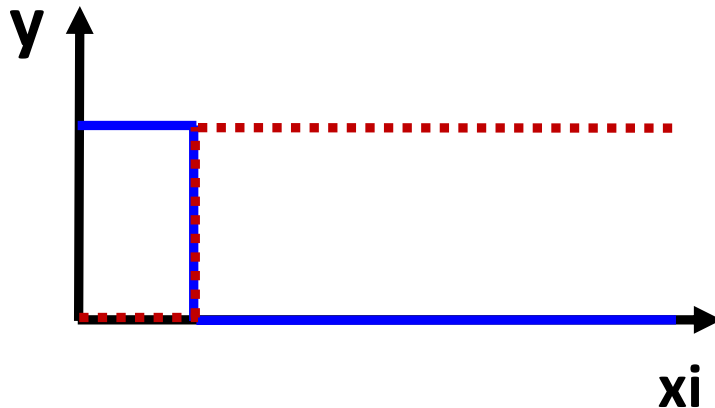
- Mostly as bagged ensembles (RFs)
- State of the art performance on many problems
- Linear in # variables, # samples

Outline

- Carts
- Mars
- Nets
- Links

MARS Learning

- Start with CART algorithm
- Indicator functions → hinge functions
- That's it!



MARS Learning

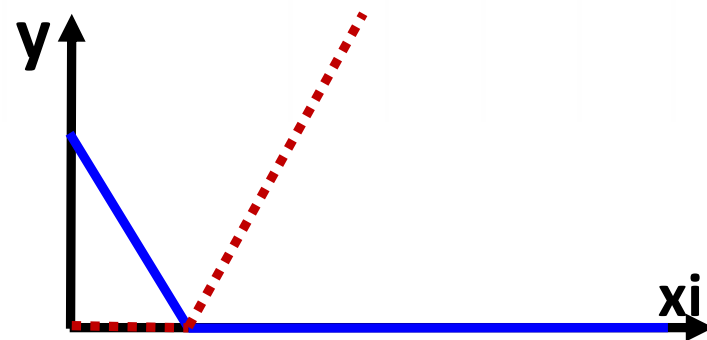
```

 $B_1(\mathbf{x}) \leftarrow 1; M \leftarrow 2$ 
Loop until  $M > M_{\max}$ :  $\text{lof}^* \leftarrow \infty$  } iteratively add  $M_{\max}$  regions (basis funcs)
  For  $m = 1$  to  $M - 1$  do:
    For  $v \notin \{v(k, m) | 1 \leq k \leq K_m\}$  } test each split: {region to split} x {split
      For  $t \in \{x_{vj} | B_m(\mathbf{x}_j) > 0\}$  } variable} x {split value}
         $g \leftarrow \sum_{i=1}^{M-1} a_i B_i(\mathbf{x}) + a_M B_m(\mathbf{x})[+(x_v - t)]_+ + a_{M+1} B_m(\mathbf{x})[-(x_v - t)]_+$ 
         $\text{lof} \leftarrow \min_{a_1, \dots, a_{M+1}} \text{LOF}(g)$ 
        if  $\text{lof} < \text{lof}^*$ , then  $\text{lof}^* \leftarrow \text{lof}; m^* \leftarrow m; v^* \leftarrow v; t^* \leftarrow t$  end if }
      end for
    end for
  end for
   $B_M(\mathbf{x}) \leftarrow B_{m^*}(\mathbf{x})[+(x_{v^*} - t^*)]_+$ 
   $B_{M+1}(\mathbf{x}) \leftarrow B_{m^*}(\mathbf{x})[-(x_{v^*} - t^*)]_+$  } add best split, in both directions
   $M \leftarrow M + 2$ 

```

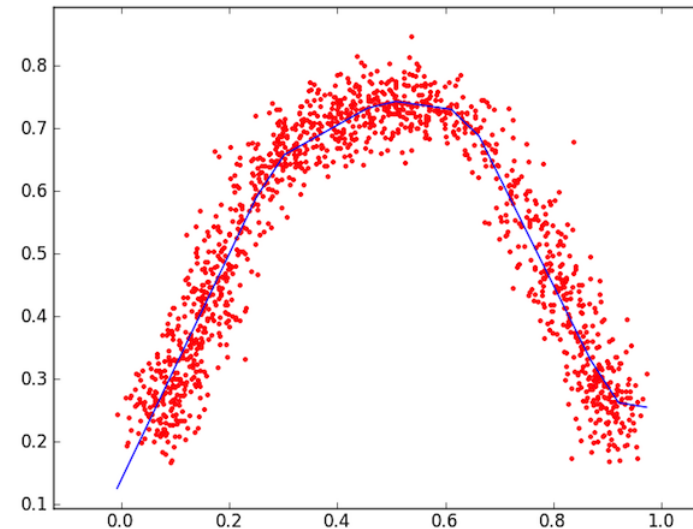
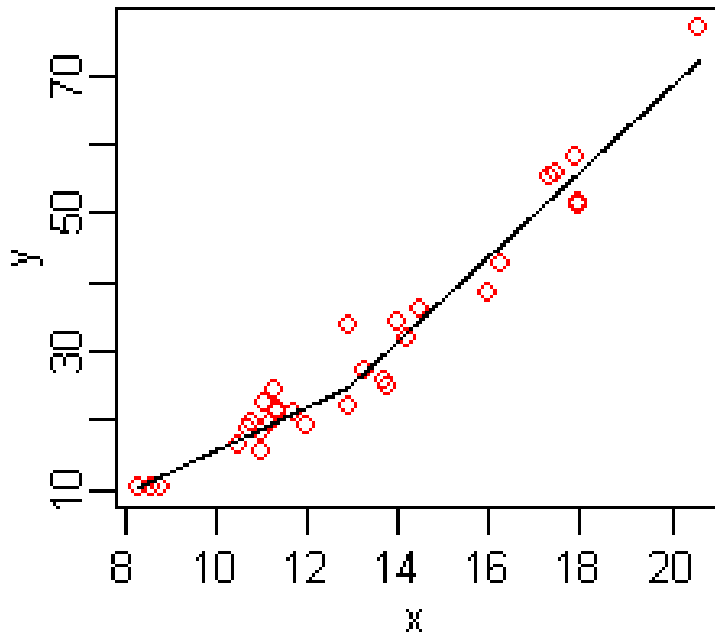
measure error,
update best split

Hinge function building blocks.



MARS Regressor

- Is a weighted sum of pairs of hinge functions
- = PWL model

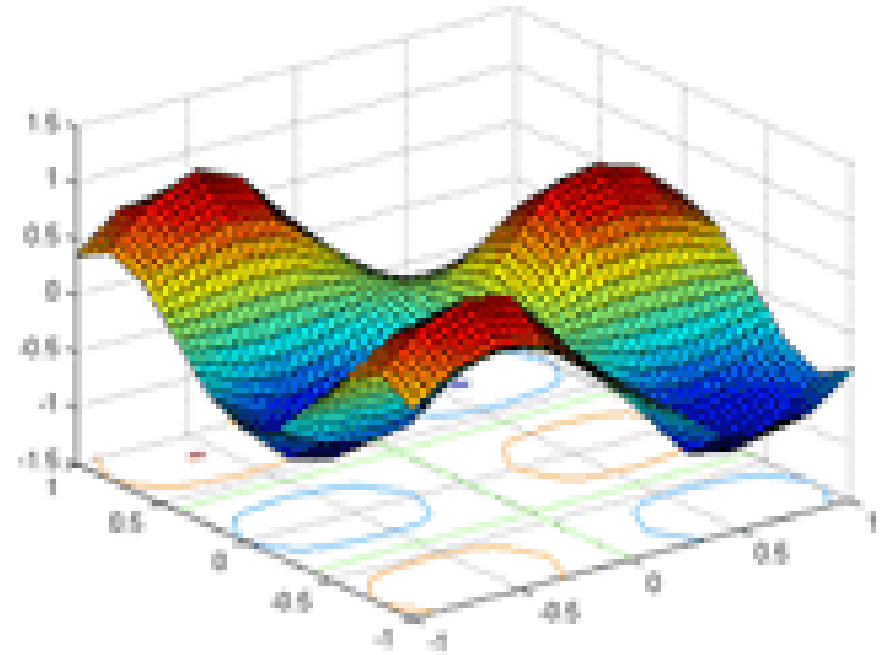
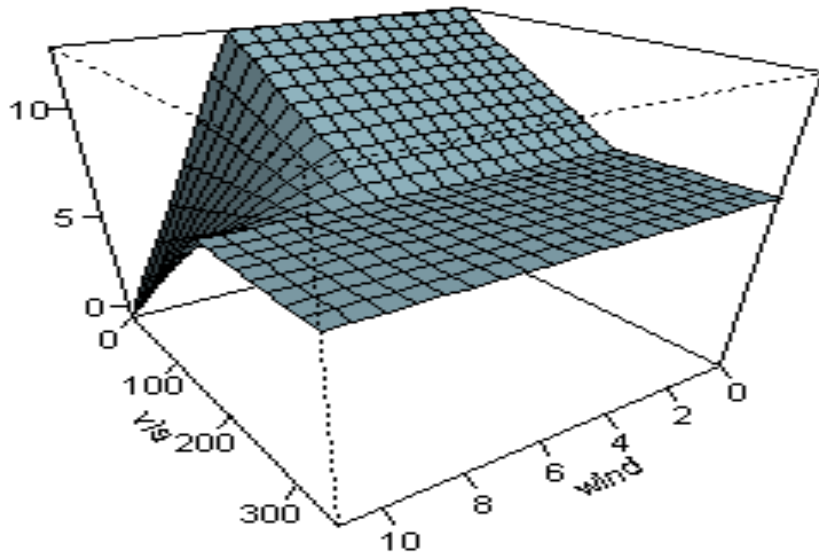


- Actually, we can have squared, cubic, .. terms
- = Piecewise polynomial = Spline

http://en.wikipedia.org/wiki/Multivariate_adaptive_regression_splines

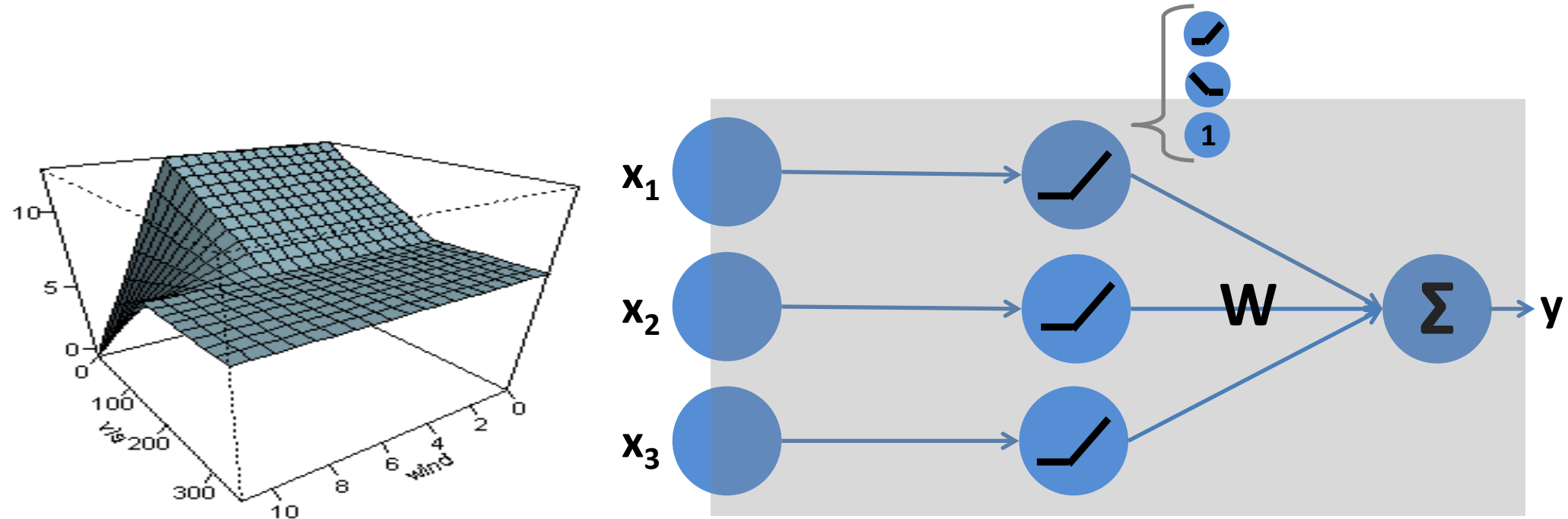
http://blog.biolab.si/wp-content/uploads/2011/12/13/earth_demo_2.png__600x470_q95_subject_location-407%2C297.png

MARS Regressor



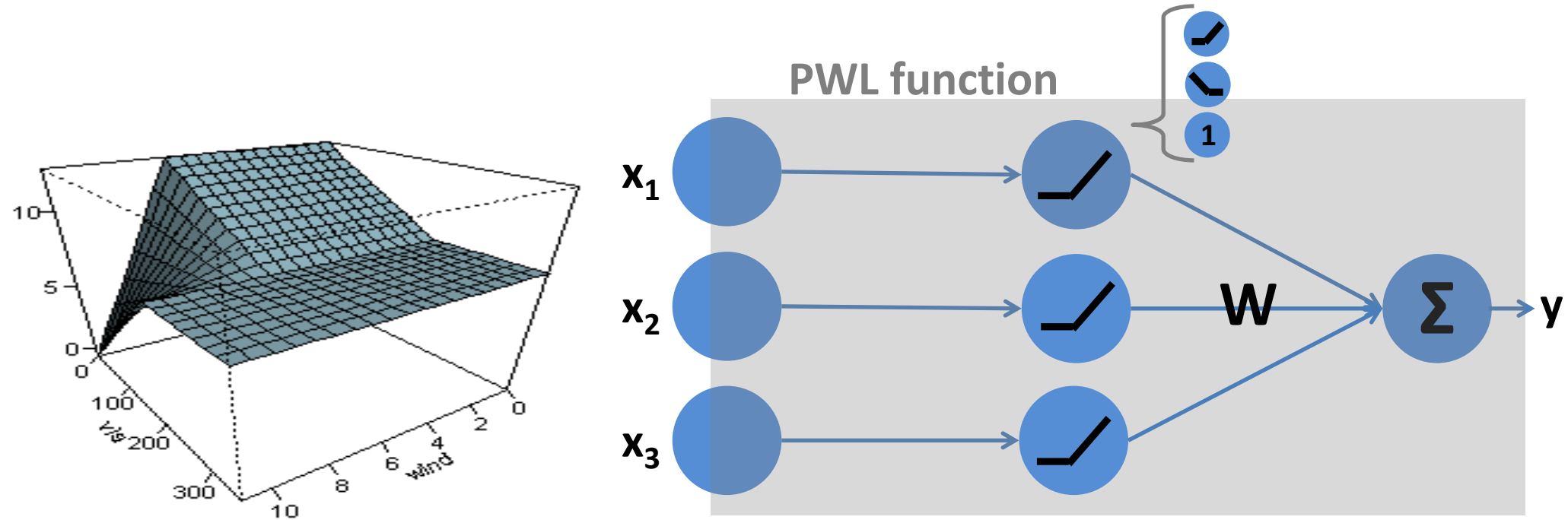
Adaptively zooms in on the most important regions.
Basis functions zero out the rest.

MARS with 1 layer



Adaptively zooms in on the most important regions.
Basis functions zero out the rest.

MARS with 1 layer



Adaptively zooms in on the most important regions.

Basis functions zero out the rest.

MARS thinking gives us intuition into PWL modeling

MARS Scalability

The diagram illustrates the scaling of the MARS complexity C^* with respect to three key parameters: the number of variables, the number of training points, and the maximum number of basis functions. The equation $C^* \sim nNM_{\max}^3(\alpha + \beta M_{\max}/L)$ is shown, with blue arrows indicating the relationship between each parameter and its corresponding term in the equation. The label 'constants' points to the terms α and β .

variables

training points

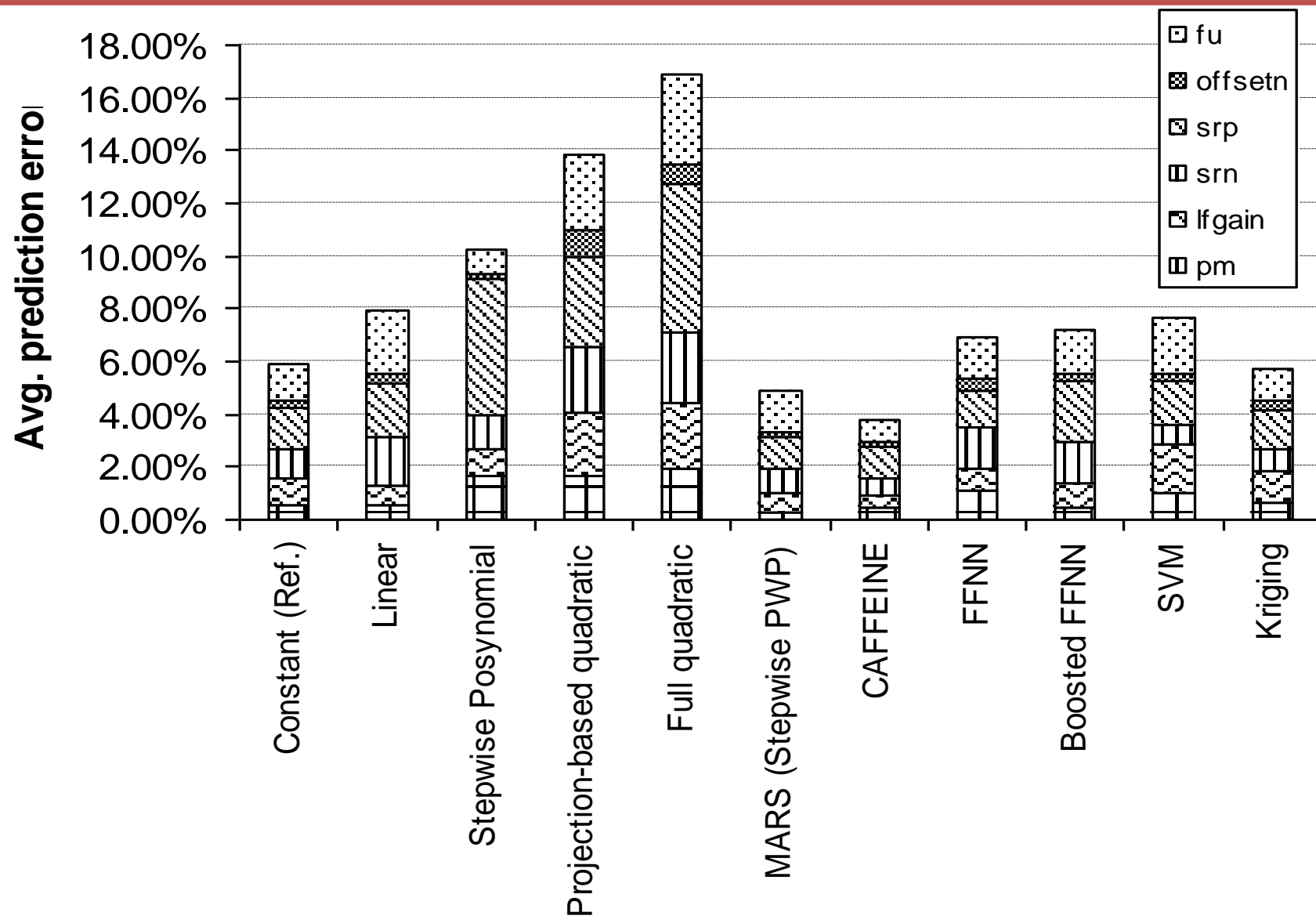
max # basis functions

$$C^* \sim nNM_{\max}^3(\alpha + \beta M_{\max}/L)$$

constants

Punchline: MARS has linear scaling in # vars and # points

MARS is Still Competitive



Aside: People Bag MARS too

Example: % error reduction by bagging

training set sizes	64		128		256		512		1024	
	<i>r</i>	<i>p</i>	<i>r</i>	<i>p</i>	<i>r</i>	<i>p</i>	<i>r</i>	<i>p</i>	<i>r</i>	<i>p</i>
kin-32fh	36.4	1	27.2	2	8.9	16	3.7	13	-3.8	22
kin-32fm	26.3	4	24.1	2	10.8	30	-8.1	28	-5.7	47
kin-32nh	40.1	1	28.7	1	23.6	1	12.3	1	6.8	1
kin-32nm	25.8	1	28.8	1	21.7	1	13.0	1	6.1	1
kin-8fh	37.3	1	16.0	4	16.8	1	13.2	2	3.8	4
kin-8fm	37.0	7	17.0	5	13.6	20	13.0	2	9.8	7
kin-8nh	20.7	5	17.6	1	11.6	1	8.3	1	9.4	4
kin-8nm	31.3	2	19.2	1	16.0	1	11.5	1	7.7	1
pumadyn-32fh	21.3	2	32.7	13	11.4	1	9.0	1	2.6	4
pumadyn-32fm	4.3	52	24.2	2	22.4	1	3.8	1	3.8	3
pumadyn-32nh	25.0	21	22.2	24	9.3	2	10.4	1	3.3	17
pumadyn-32nm	28.3	43	-23.9	44	22.5	1	16.8	1	5.3	7
pumadyn-8fh	12.9	11	3.6	59	4.5	16	5.5	1	3.6	7
pumadyn-8fm	20.2	2	15.1	2	8.4	3	4.4	2	5.9	4
pumadyn-8nh	12.1	41	19.0	2	10.5	1	10.1	1	4.0	7
pumadyn-8nm	22.7	8	30.3	1	17.0	3	13.0	3	7.5	25
averages	25.1		18.9		14.3		8.7		4.3	

Outline

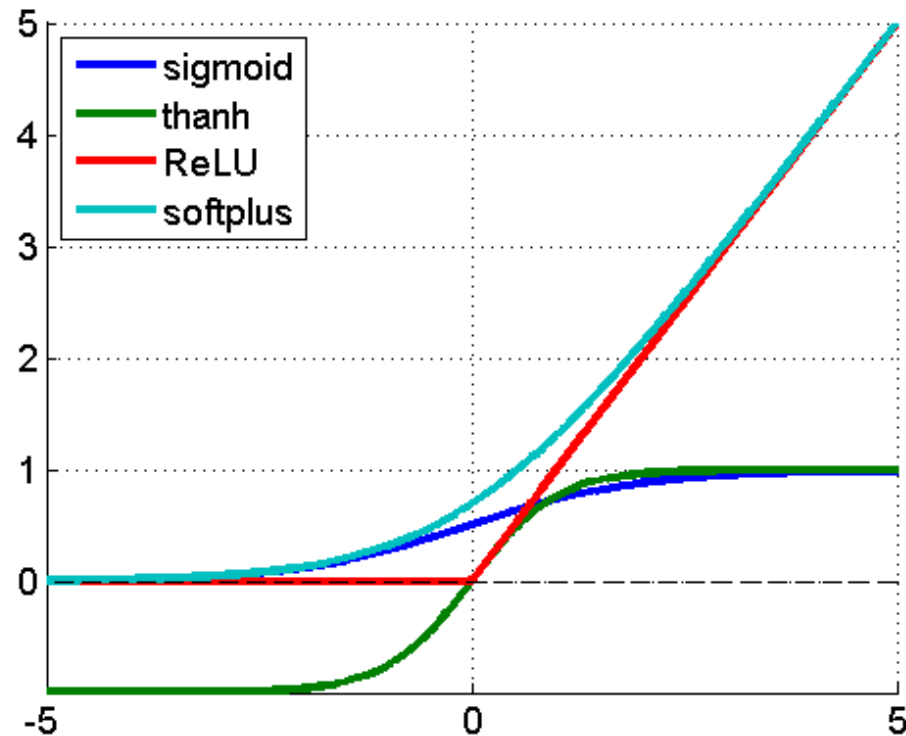
- Carts
- Mars
- Nets

Deep Nets

1. Intuition is difficult
 2. Shortage of theory
 3. Long training times
- But they work really well on many problems!
 - So people use them
 - My goal here: help 1, maybe 2, maybe 3

Deep Nets

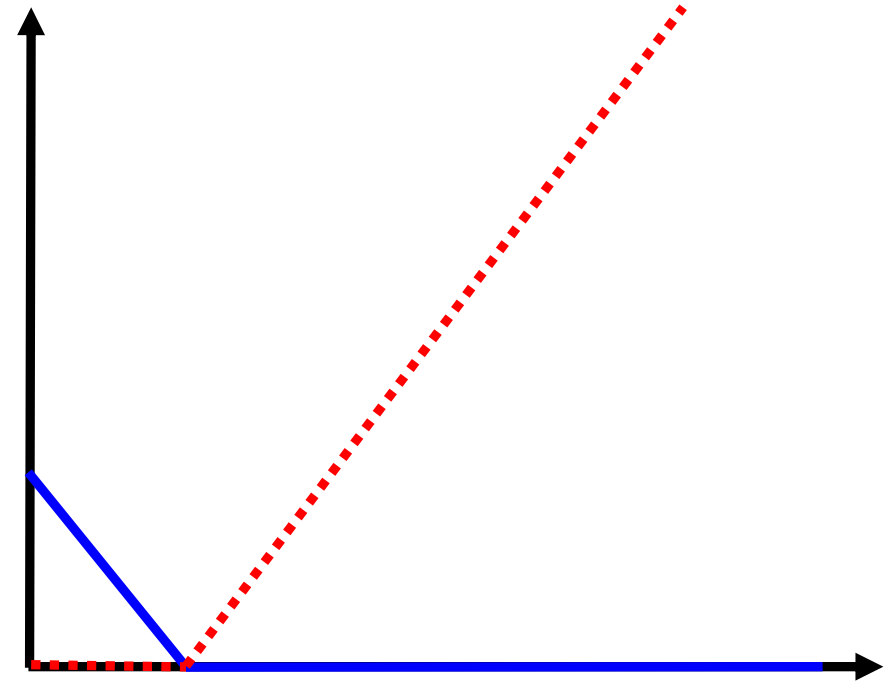
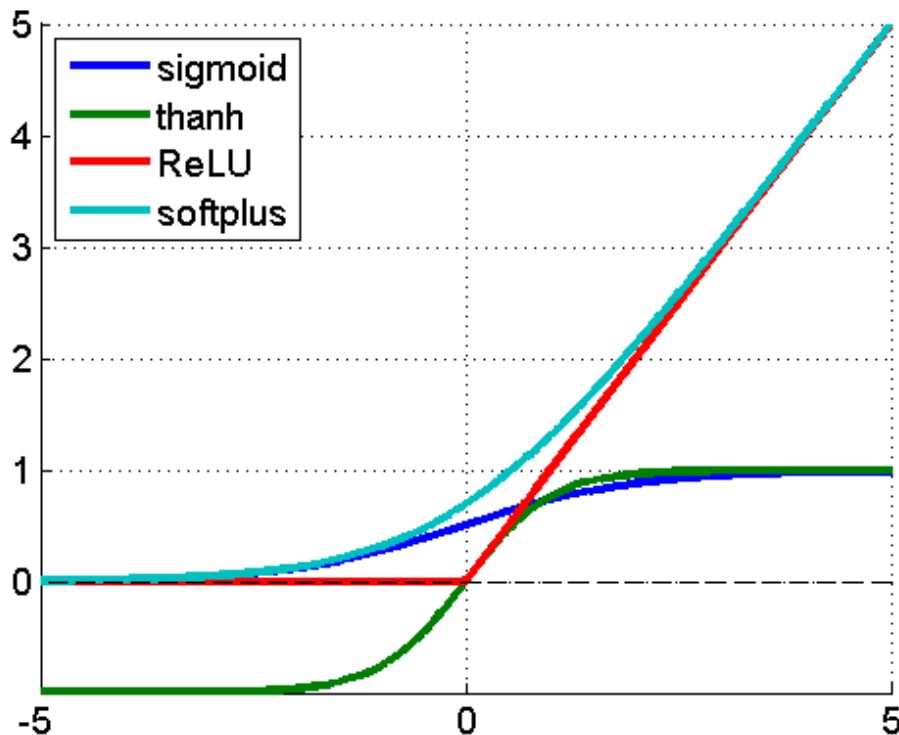
Rectified Linear Units (ReLUs)



Deep Nets

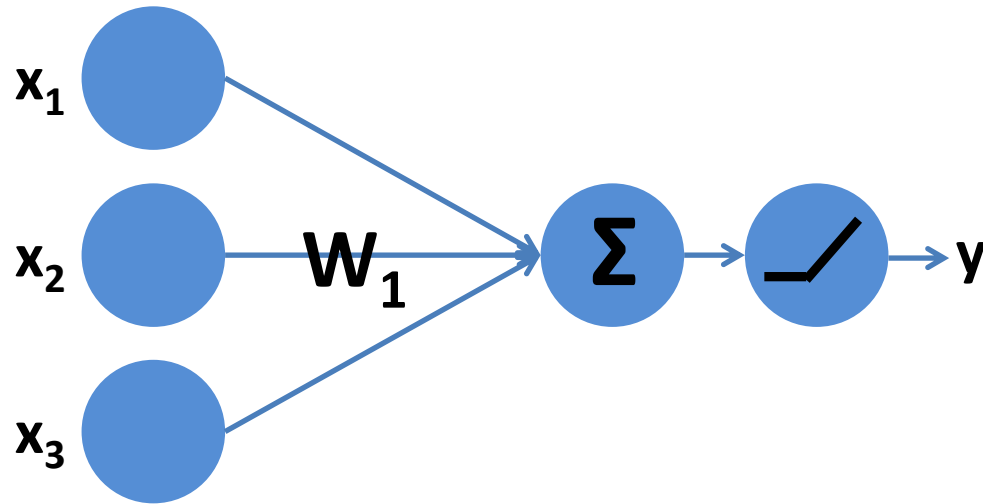
Rectified Linear Units (ReLUs)

≈ MARS Hinge Functions

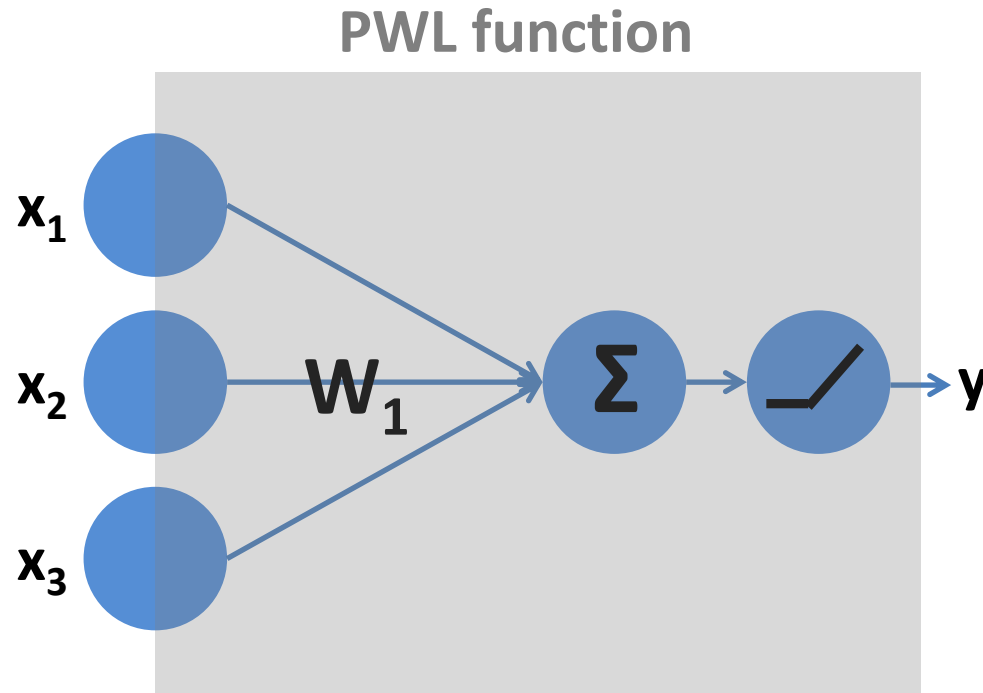


Benefit: we can translate intuition from MARS into Deep Nets. E.g. adaptively zooming in on most important regions, zero out rest

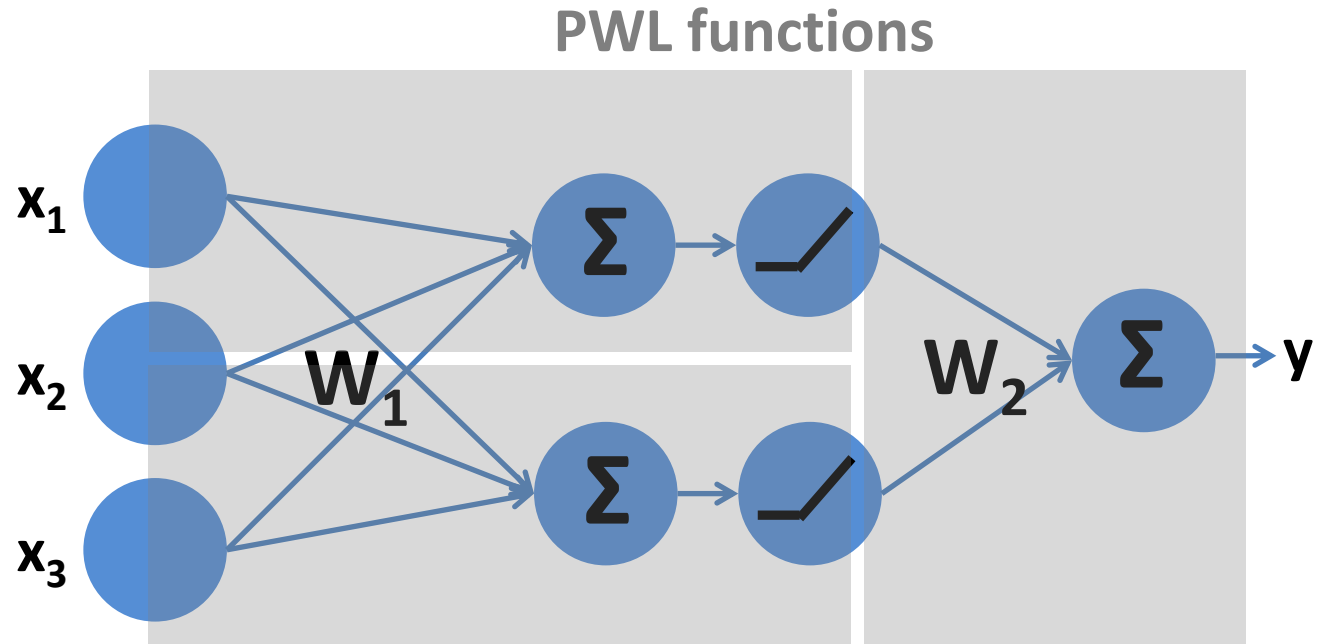
NN with 1 layer, 1 hidden node



NN with 1 layer, 1 hidden node



NN with 1 layer, 2 hidden nodes



Deep NN with 5 layers, 6 nodes per layer

PWL functions

Deep NN with 5 layers, 6 nodes per layer



Insight from MARS based PWL models can apply here.
Food for thought: what other MARS tricks can we port to Deep NNs?
Eg cascade MARS & dropout-regularize for “Deep MARS”?

Conclusion

- CARTs are old but good: fast, accurate, understood (accurate with RFs)
- MARS is old but good: fast, accurate, understood
- Deep NNs are new(ish). Very accurate, but slow & poorly understood
- Insights from CARTs led to MARS
- Insights from MARS can help Deep NNs

Appendix: MARS Software

Free [\[edit\]](#)

- Several [R](#) packages fit MARS-type models:
 - `earth` function in the [earth](#) package
 - `mars` function in the [mda](#) package
 - `polymars` function in the [polspline](#) package. Not Friedman's MARS.
- Matlab code:
 - [ARESLab: Adaptive Regression Splines toolbox for Matlab](#)
- Python
 - [Earth - Multivariate adaptive regression splines](#)
 - [py-earth](#)

Commercial [\[edit\]](#)

- [MARS](#) from Salford Systems. Based on Friedman's implementation.
- [STATISTICA Data Miner](#) from StatSoft
- [ADAPTIVEREG](#) from SAS.