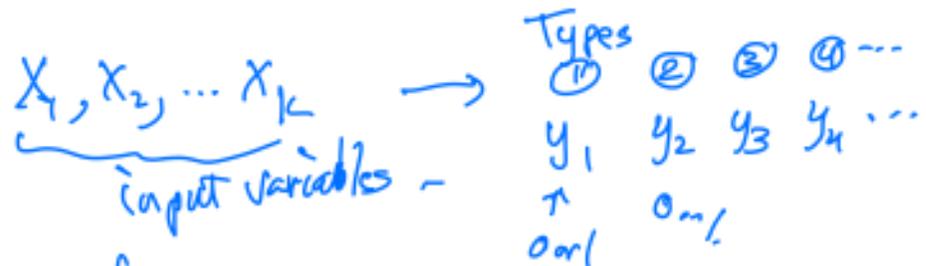


Another way to classify data using a smoothing function:



Predicted y -values: $\hat{y}_1, \hat{y}_2, \hat{y}_3, \hat{y}_4, \dots$
 want to be between 0 & 1.

Parameters θ for model for each type
 $\theta^{(j)}$ for type y_j .

Do regression: answer.

$$\hat{y}_j = \frac{e^{\theta^{(j)T} x}}{\sum_j e^{\theta^{(j)T} x}}$$

soft max function.

See that $\sum_j \hat{y}_j = 1$.

Think that \hat{y}_j = probability that the data point is of type (j) (i.e. $y_j = 1$).

Your real prediction:

$$\max_j \frac{e^{\theta^{(j)T} x}}{\sum_k e^{\theta^{(k)T} x}} = \text{prob. that your answer is correct.}$$

See examples of SVM classifiers using polynomial regression inputs:

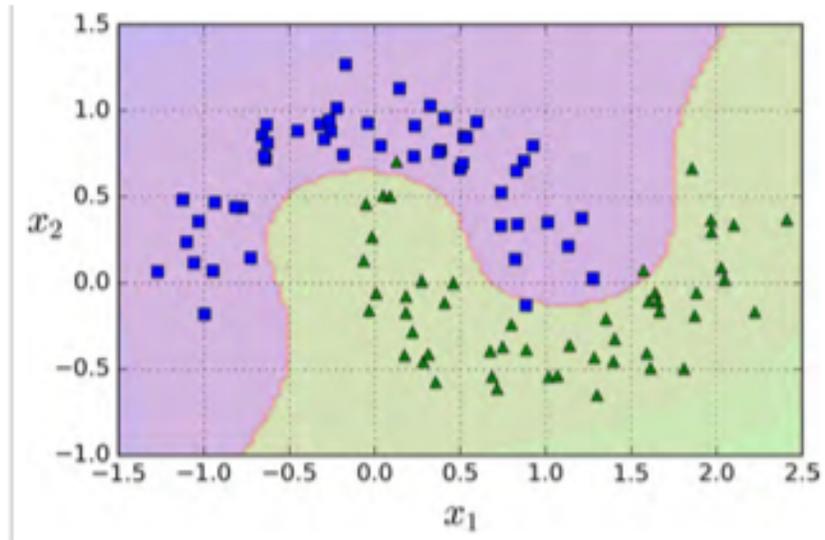


Figure 5-6. Linear SVM classifier using polynomial features

Another Classification Algorithm:
Decision Tree.

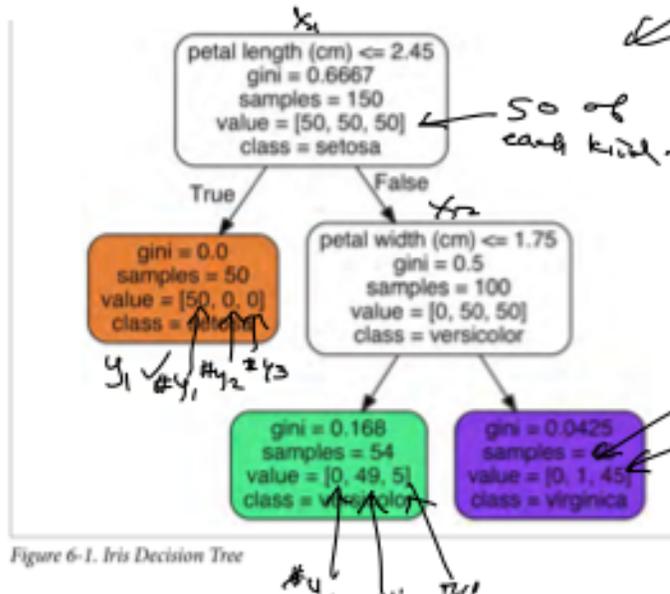


Figure 6-1. Iris Decision Tree

example from the book trying to classify the type of iris (setosa, versicolor, virginica) from petal length, petal width, sepal length, ...

#y2
#y3

" #4, #73
+
wrong

Gini = shows how good the classification in a given box is. ← percentage of j th type in the box.
 $= 1 - \sum_{j=1}^K P_j$ ← want to be close to 0.

$$\begin{aligned} \text{1st box gini} &= 1 - \left[\left(\frac{50}{150}\right)^2 + \left(\frac{50}{150}\right)^2 + \left(\frac{50}{150}\right)^2 \right] \\ &= \frac{2}{3}. \end{aligned}$$

$$\begin{aligned} \text{orange box gini} &= 1 - \left[\left(\frac{50}{50}\right)^2 + \left(\frac{0}{58}\right)^2 + \left(\frac{0}{58}\right)^2 \right] \\ &= 0. \end{aligned}$$

$$\begin{aligned} \text{green box gini} &= 1 - \left[\left(\frac{0}{54}\right)^2 + \left(\frac{49}{54}\right)^2 + \left(\frac{5}{54}\right)^2 \right] \\ &\approx 0.168 \end{aligned}$$

Question: How is the variable and breakpoint (eg $x_2 \leq 2.783$) chosen?

Answer: Consider all possible input columns, all possible breakpoints → choose the input column (variable) and break pt so that it gives a minimum

Score for one of the resulting boxes.

Then: keep going through the decision tree -

How does this look like with the classification of the data?

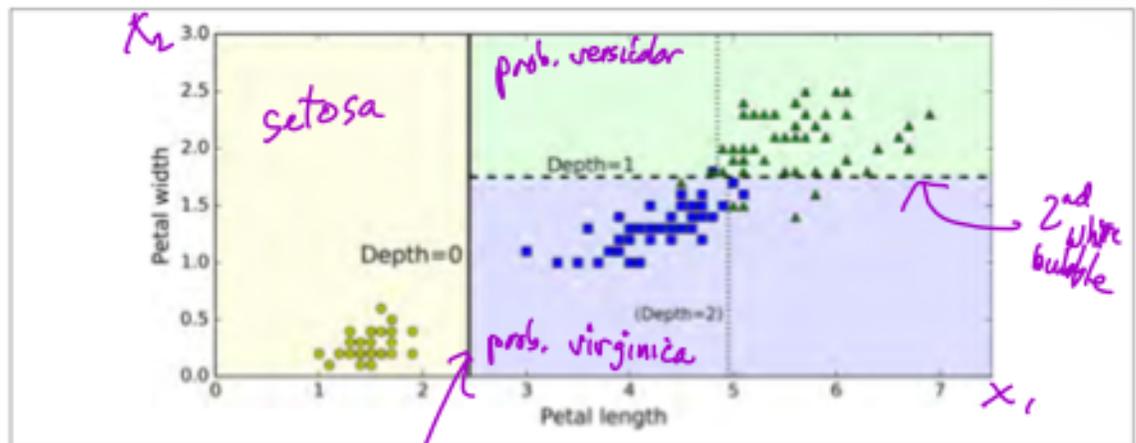


Figure 6-2. Decision Tree decision boundaries

Ensemble Learning -

several different algorithms at the same time.

- Use all the different algorithms to make a prediction of the classification of a ^{new} data point (x_1, x_2, \dots, x_n)

(Then \rightarrow the most often chosen classification is the one you choose.)

- Could also choose the classification is by weighting votes with output probabilities.
-

When we use the same training data but chop into small subsets (randomly, with replacement) → use one algorithm on this.

→ get ensemble result (eg by voting)

→ This is called bagging

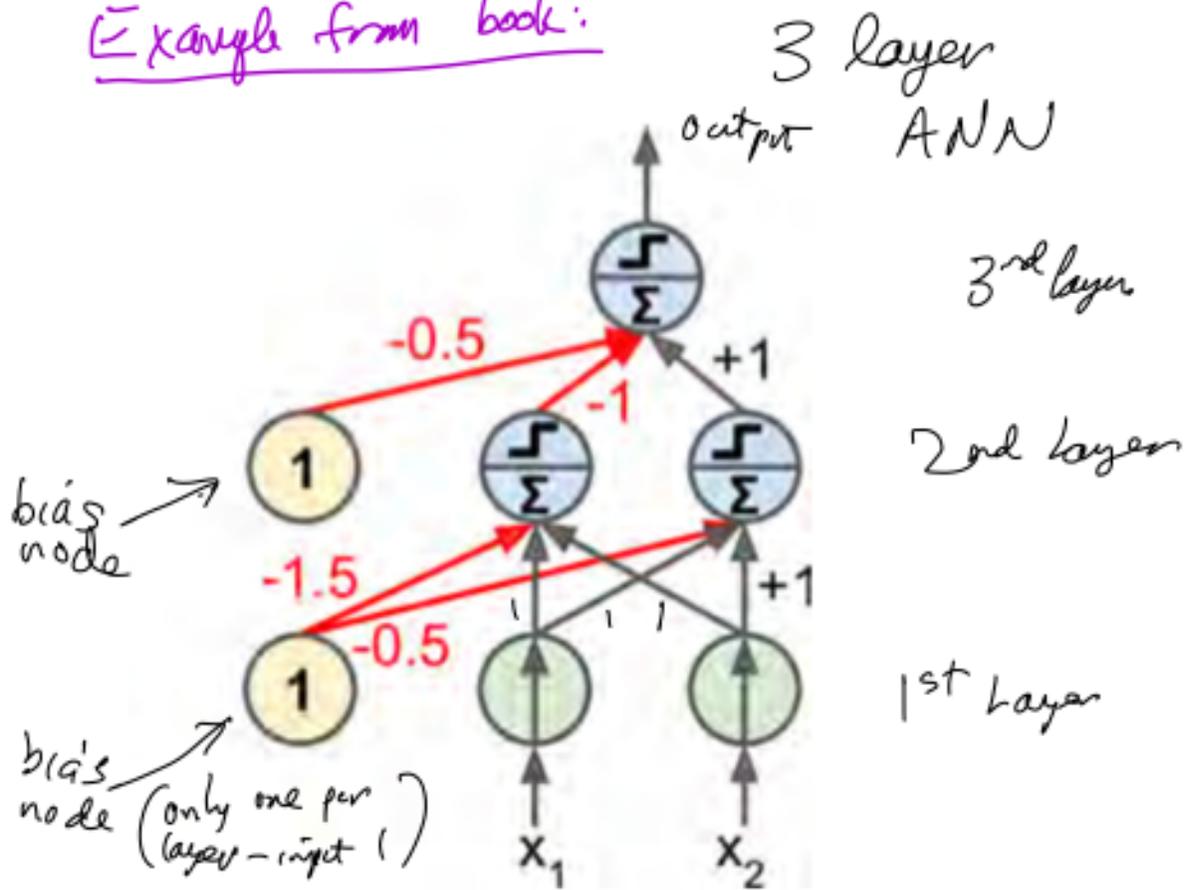
Amazing: it gives better outcome than just performing the algorithm on the whole training set.

Random Forest — uses bagging with the decision tree algorithm.

- Artificial Neural Networks (ANNs)
- Deep Learning Neural Networks —
 - same as regular neural networks but with hidden layers.

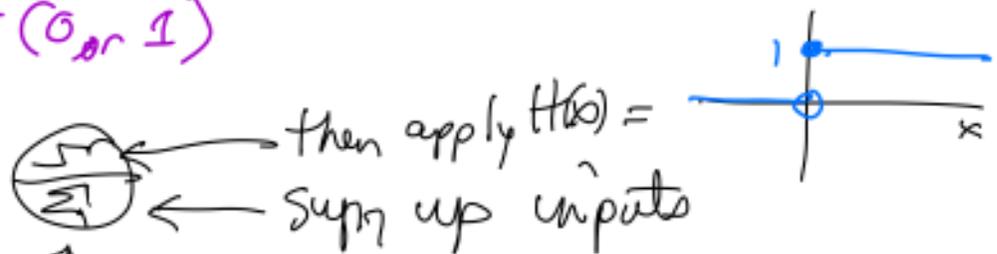
These are graph networks
 nodes = vertices
 input nodes & output nodes.

Example from book:



$x_1 = \text{input (0 or 1)}$

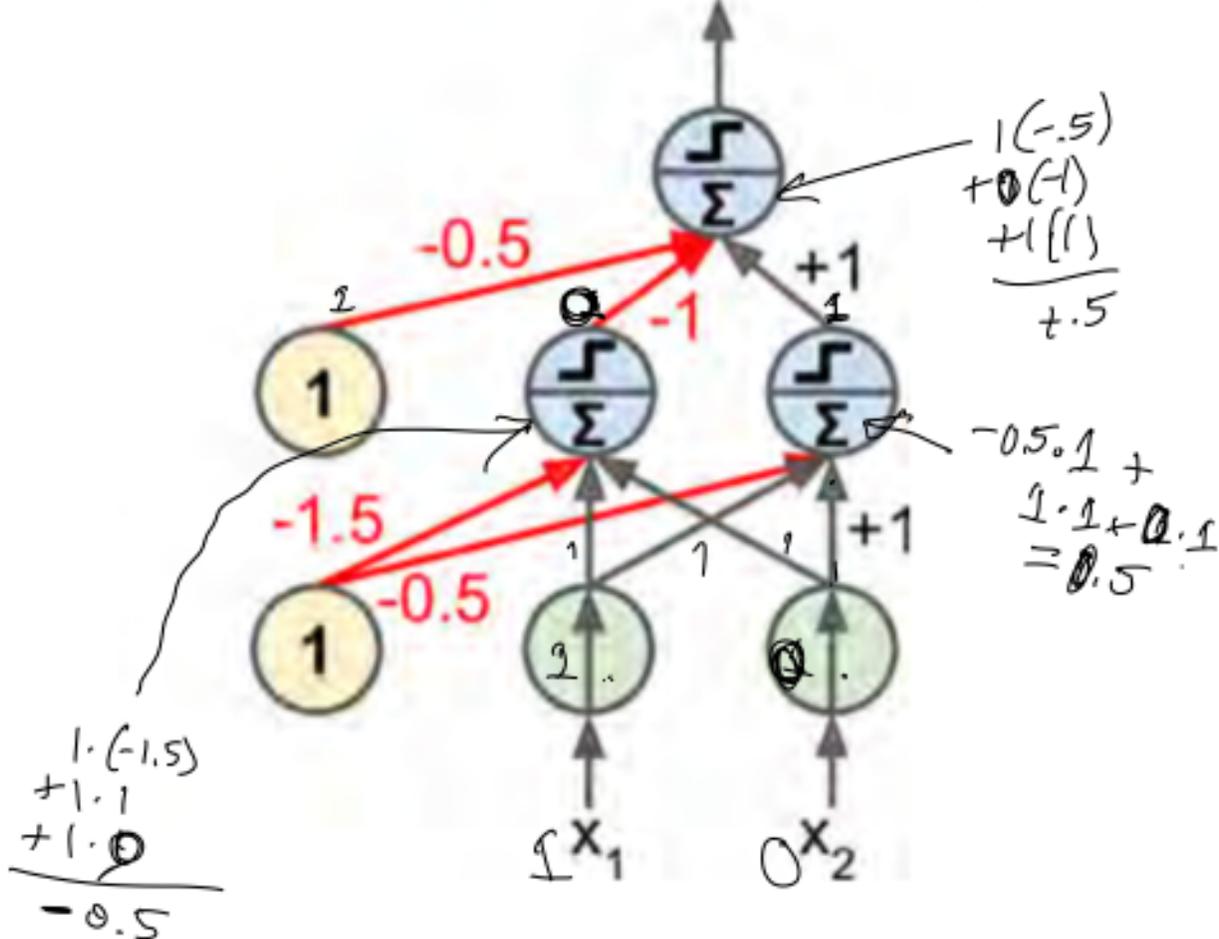
$x_2 = \text{input (0 or 1)}$



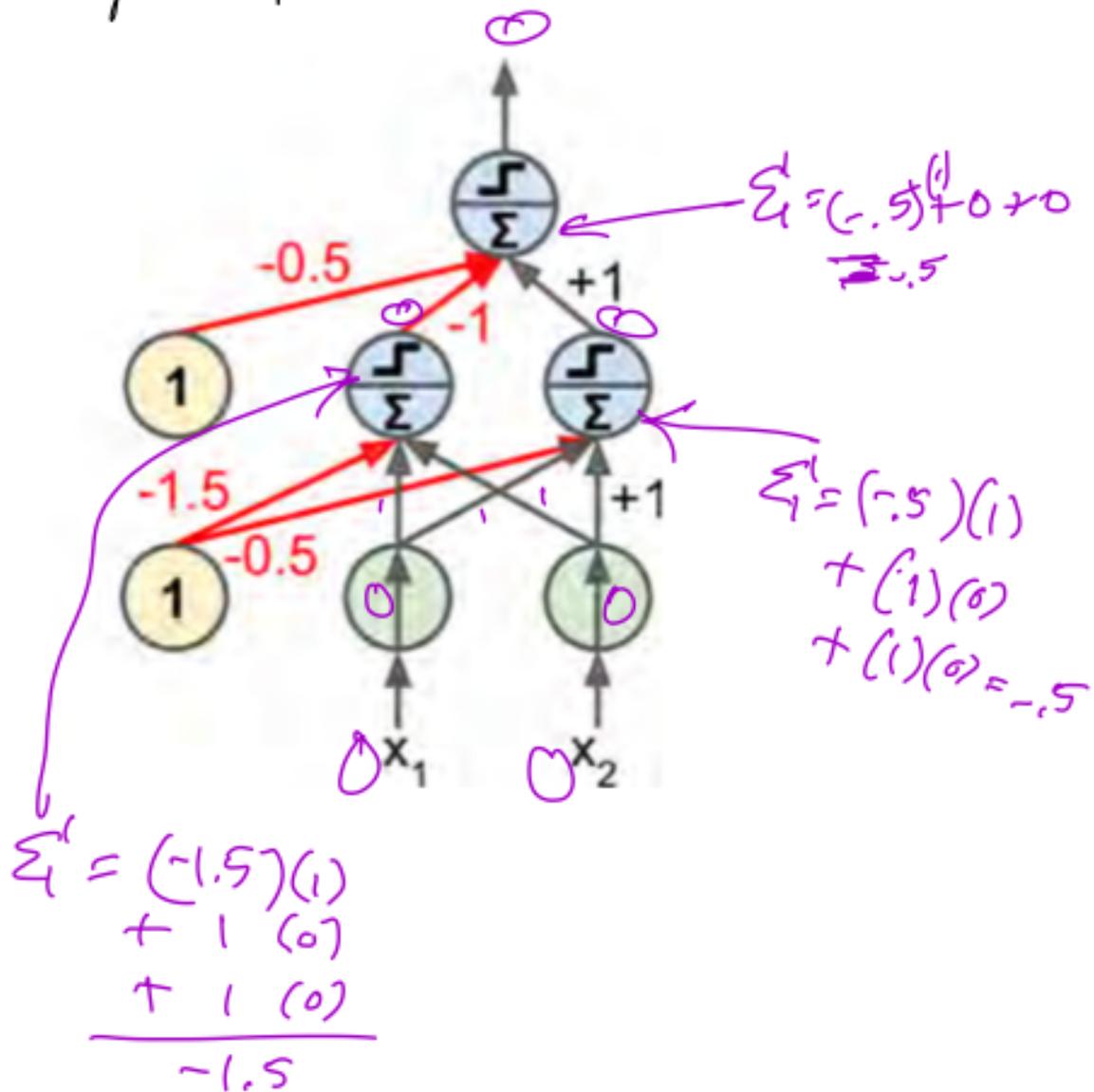
Numbers next to arrows are weights $w_{ij} =$
 weight from i^{th} node to j^{th} node. \sum_i means $\sum w_{ij} I_{ij}$

Sample input: $x_1 = 1, x_2 = 0$

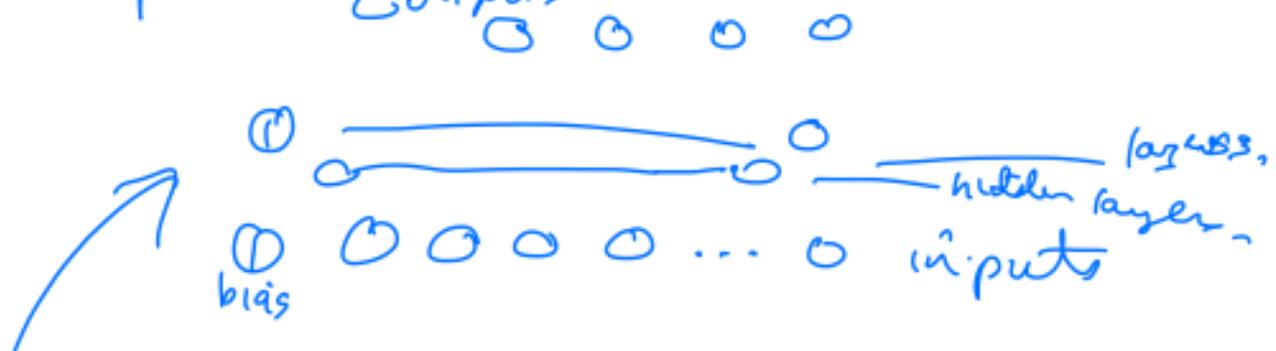
1 ← output.



Sample Input $x_1 = 0$ $x_2 = 0$



Set up Neural Net



↳ each layer has its own number of nodes —

parameters (mostly hidden)

→ weights w_{ij} .

→ Training →

reduce the weights contributing to wrong answers —

(follow backwards)

increase weights contr.

to correct answers.