



Neural Networks

in 1 lecture!

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

January 31st, 2005

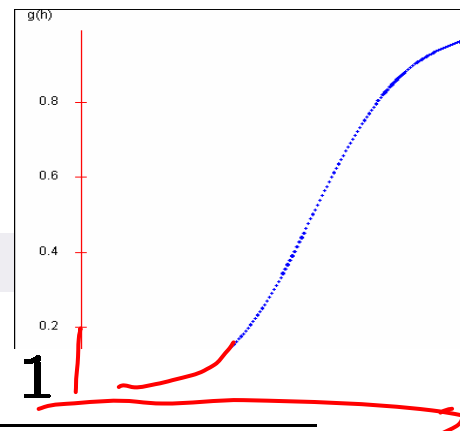
Announcements



- First homework due Wednesday
 - **Beginning** of class
 - If using late days (24 hours per late day), timestamp and hand to Sharon Cavlovich, Wean Hall 5315

w_0 is the weight for constant: $x_0^j = 1 \forall j$

Logistic regression



- $P(Y|X)$ represented by:

$$P(Y = 1 | x, W) = \frac{1}{1 + e^{-(w_0 + \sum_i w_i x_i)}} = g(w_0 + \sum_i w_i x_i)$$

Sigmoid

- Learning rule – MLE:

$$\frac{\partial \ell(W)}{\partial w_i} \underset{\text{data}}{=} \sum_j x_i^j [y^j - P(Y^j = 1 | x^j, W)] = \sum_j x_i^j [y^j - g(w_0 + \sum_i w_i x_i^j)]$$

$$\underline{w_i} \leftarrow w_i + \eta \sum_j x_i^j \delta^j$$

how well j is classified

learning rate

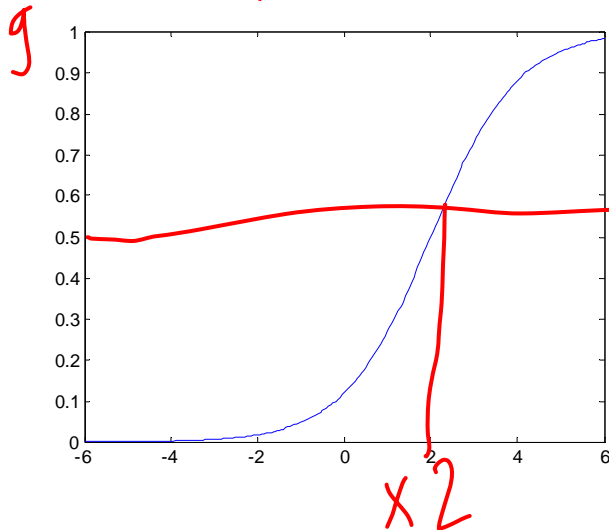
$$\delta^j = y^j - g(w_0 + \sum_i w_i x_i^j)$$

Sigmoid

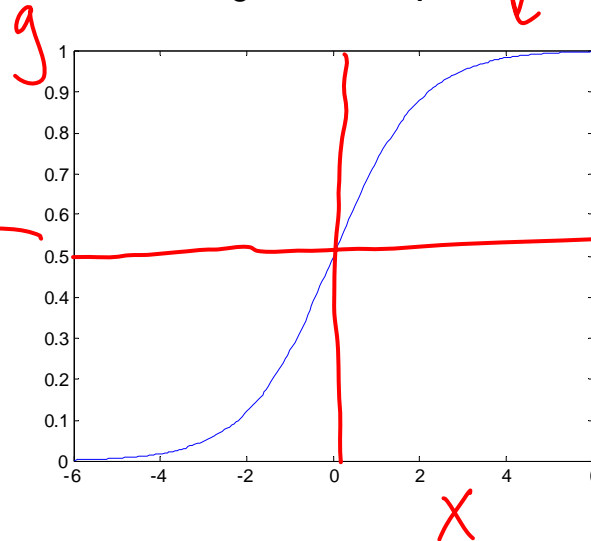
$$g(w_0 + \sum_i w_i x_i) = \frac{1}{1 + e^{-(w_0 + \sum_i w_i x_i)}}$$

change weight constant

$w_0=2, w_1=1$

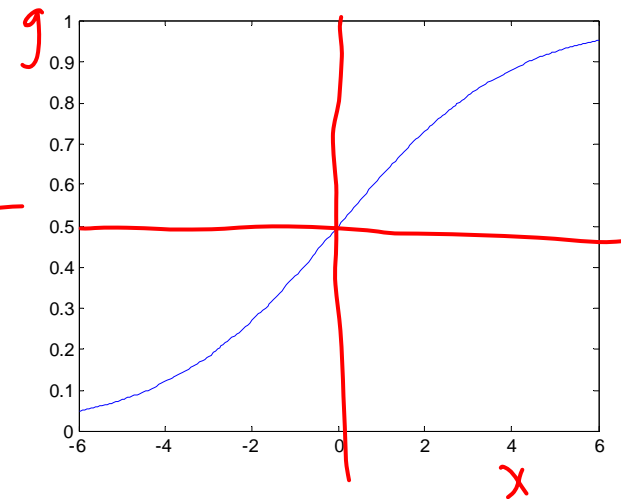


$w_0=0, w_1=1$

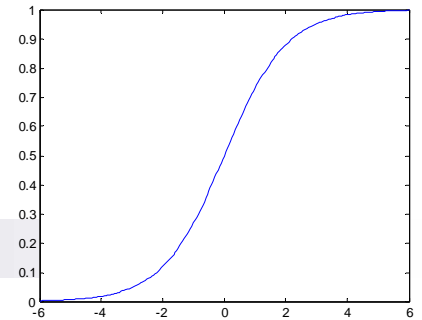


change w_1

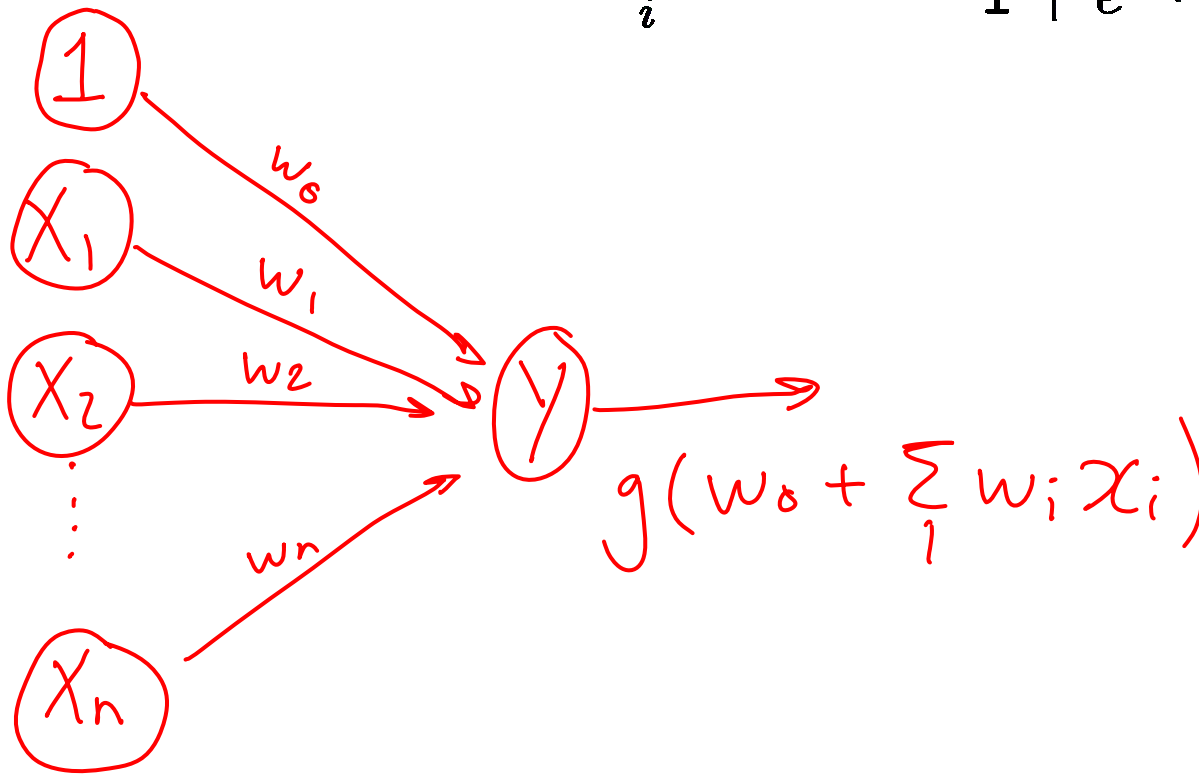
$w_0=0, w_1=0.5$



Perceptron as a graph



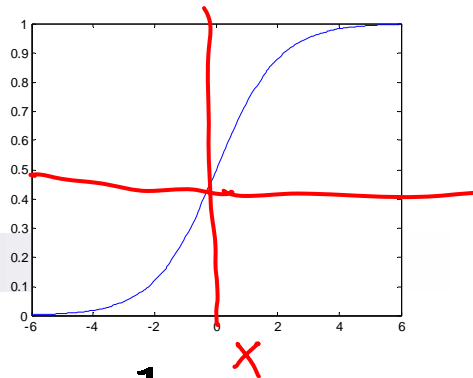
$$g(w_0 + \sum_i w_i x_i) = \frac{1}{1 + e^{-(w_0 + \sum_i w_i x_i)}}$$



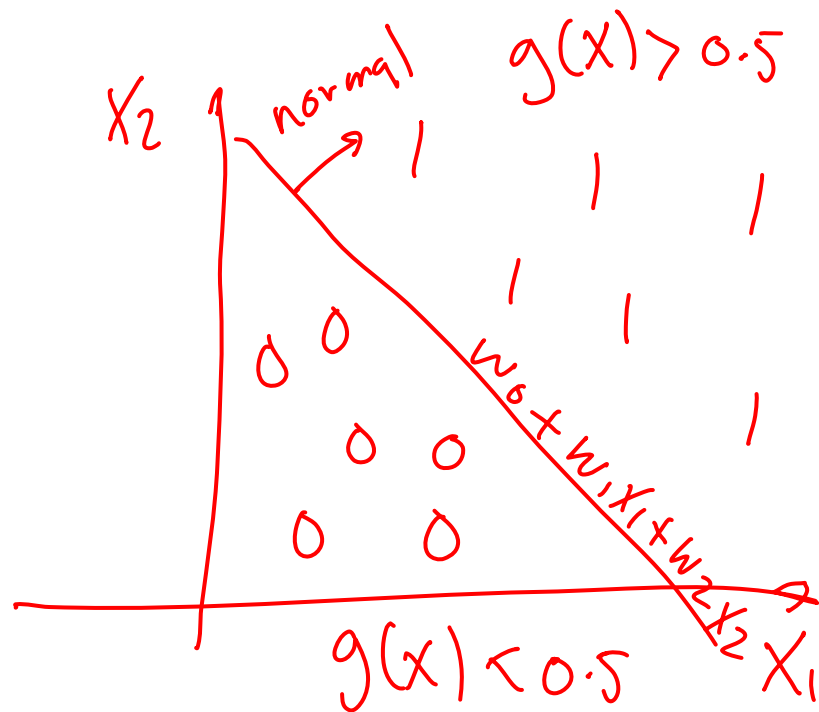
Linear perceptron classification region

$$g(x) \geq 0.5 \quad \text{when } x > 0$$

$$g(x) < 0.5 \quad \text{when } x < 0$$



$$g(w_0 + \sum_i w_i x_i) = \frac{1}{1 + e^{-(w_0 + \sum_i w_i x_i)}}$$



$$g(w_0 + \sum_i w_i x_i)$$

$g(\text{line})$

Optimizing the perceptron

- Trained to minimize sum-squared error

$$\ell(W) = \frac{1}{2} \sum_j [y^j - g(w_0 + \sum_i w_i x_i^j)]^2$$

$$\frac{\partial \ell(W)}{\partial w_i} = \sum_j -[y^j - g(w_0 + \sum_i w_i x_i^j)] *$$

$$\frac{\partial}{\partial w_i} g(w_0 + \sum_i w_i x_i^j)$$

→
derivative
of
Sigmoids

Derivative of sigmoid

$$\frac{\partial \ell(W)}{\partial w_i} = - \sum_j [y^j - g(w_0 + \sum_i w_i x_i^j)] x_i^j g'(w_0 + \sum_i w_i x_i^j)$$

$$g(x) = \frac{1}{1 + e^{-x}}$$

$$; \quad \frac{\partial g(x)}{\partial x} = \frac{e^{-x}}{(1 + e^{-x})^2}$$

$$\frac{\partial g(x)}{\partial x} = g(x)(1 - g(x))$$

$$= \frac{1 - 1 + e^{-x}}{(1 + e^{-x})^2}$$

$$= \boxed{\frac{1}{(1 + e^{-x})}} \left(\frac{-1}{1 + e^{-x}} + \frac{1 + e^{-x}}{1 + e^{-x}} \right)$$

The perceptron learning rule

$$w_i \leftarrow w_i + \eta \sum_j x_i^j \delta^j$$

learning rate

training example feature

difference

$$\delta^j = [y^j - g(w_0 + \sum_i w_i x_i^j)] g^j (1 - g^j)$$

$$g^j = g(w_0 + \sum_i w_i x_i^j)$$

■ Compare to MLE:

$$w_i \leftarrow w_i + \eta \sum_j x_i^j \delta^j$$

$$\delta^j = [y^j - g(w_0 + \sum_i w_i x_i^j)]$$

Perceptron, linear classification, Boolean functions

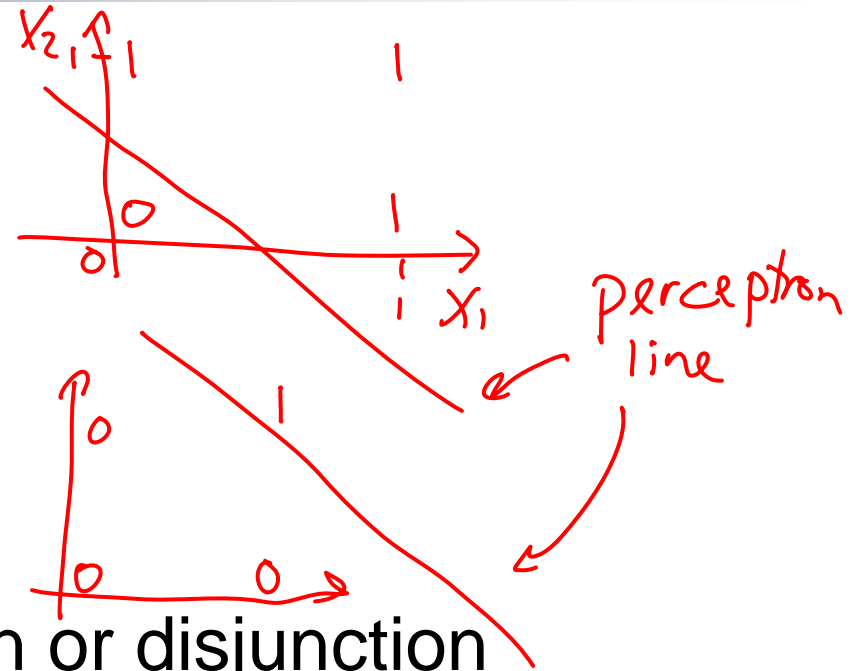
- Can learn $x_1 \vee x_2$

- Can learn $x_1 \wedge x_2$

- Can learn any conjunction or disjunction

$$x_1 \wedge \neg x_2 \wedge x_3 \wedge x_5 \wedge \neg x_6$$

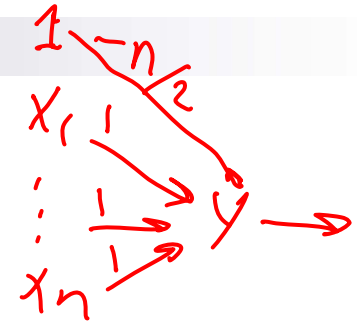
$$x_1 \vee \neg x_2 \vee x_3 \vee x_5 \vee \neg x_6$$



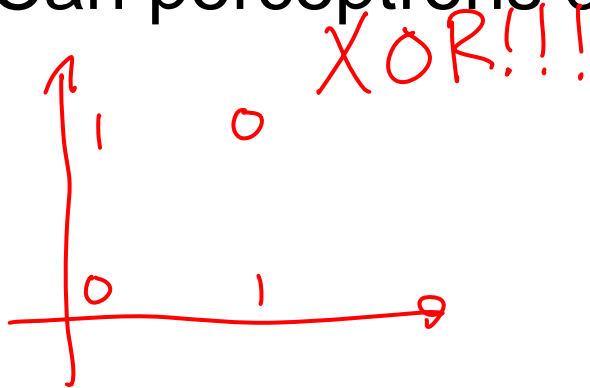
Perceptron, linear classification, Boolean functions

- Can learn majority

$$\sum_{i=1}^n x_i > \frac{n}{2}$$



- Can perceptrons do everything?

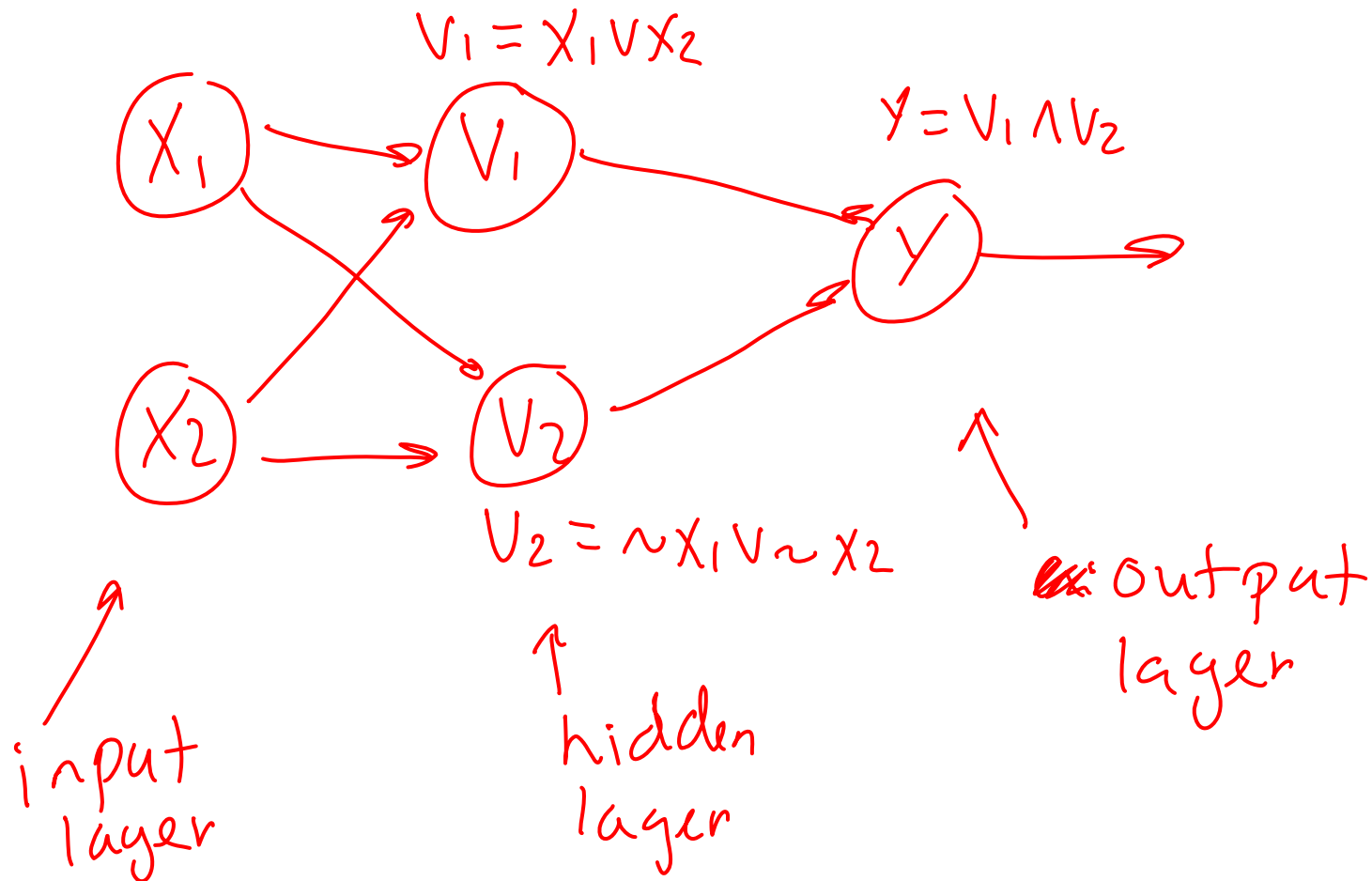


$$(x_1 \vee x_2) \wedge (\sim x_1 \vee \sim x_2)$$

no line can separate

Going beyond linear classification

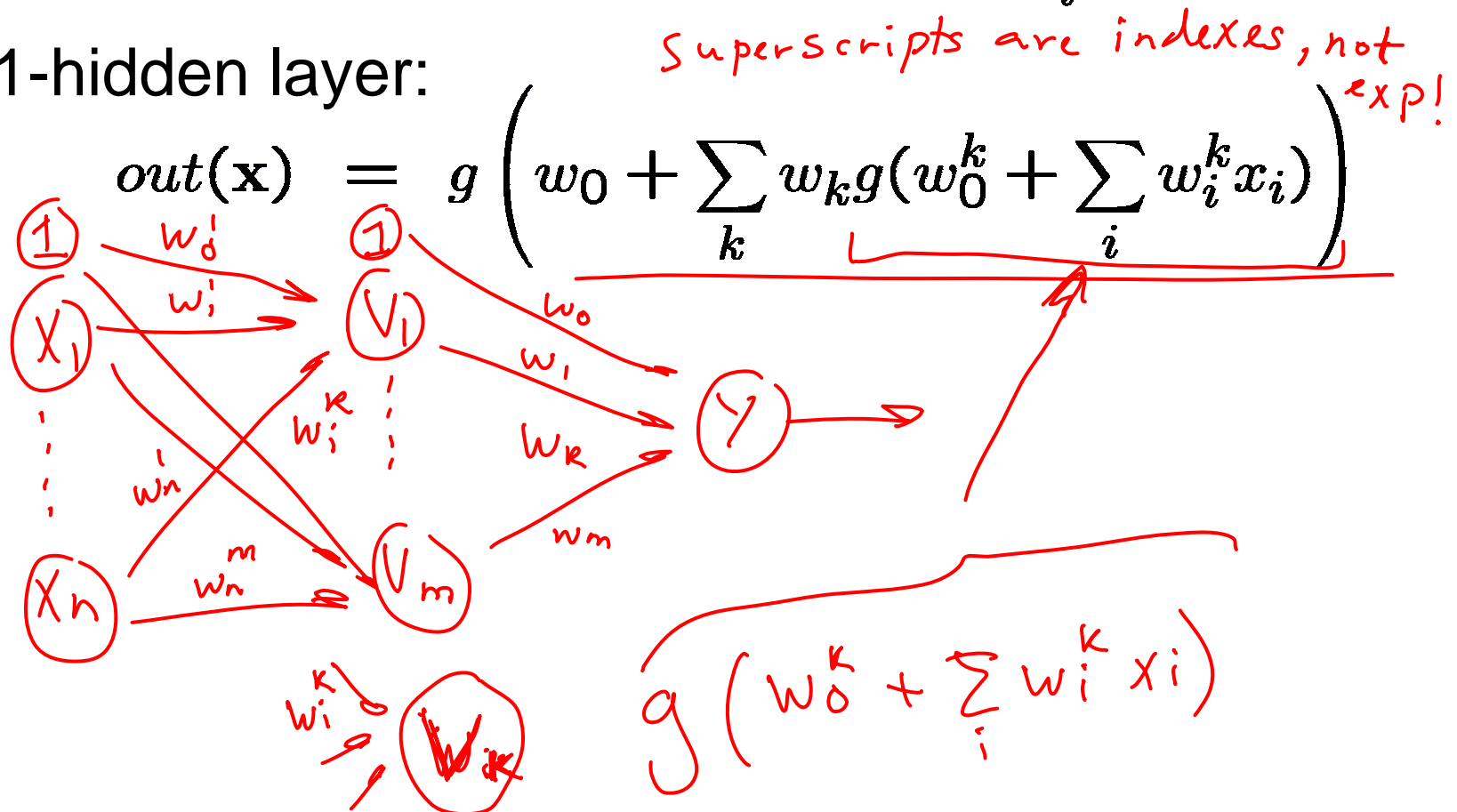
- Solving the XOR problem $(X_1 \vee X_2) \wedge (\sim X_1 \vee \sim X_2)$



Hidden layer

- Perceptron: $out(\mathbf{x}) = g(w_0 + \sum_i w_i x_i)$

- 1-hidden layer:



Gradient descent for 1-hidden layer

Dropped w_0 for simpler notation
perceptron \rightarrow

$$w_i \leftarrow w_i + \eta x_i \delta$$

$$\delta = [y - g(w_0 + \sum_i w_i x_i)] g'(1 - g)$$

$$out(\mathbf{x}) = g\left(\sum_k w_k g\left(\sum_i w_i^k x_i\right)\right)$$

$$\frac{\partial \ell(W)}{\partial w} = -[y - out(\mathbf{x})] \frac{\partial out(\mathbf{x})}{\partial w} \quad (*)$$

output layer:

$$\frac{\partial out(\mathbf{x})}{\partial w_k} = g\left(\sum_i w_i^k x_i\right) \cdot g'\left(\sum_k w_k g\left(\sum_i w_i^k x_i\right)\right)$$

$$= g\left(\sum_i w_i^k x_i\right) g\left[\sum_k w_k g\left(\sum_i w_i^k x_i\right)\right] [1 - g\left(\sum_k w_k g\left(\sum_i w_i^k x_i\right)\right)]$$

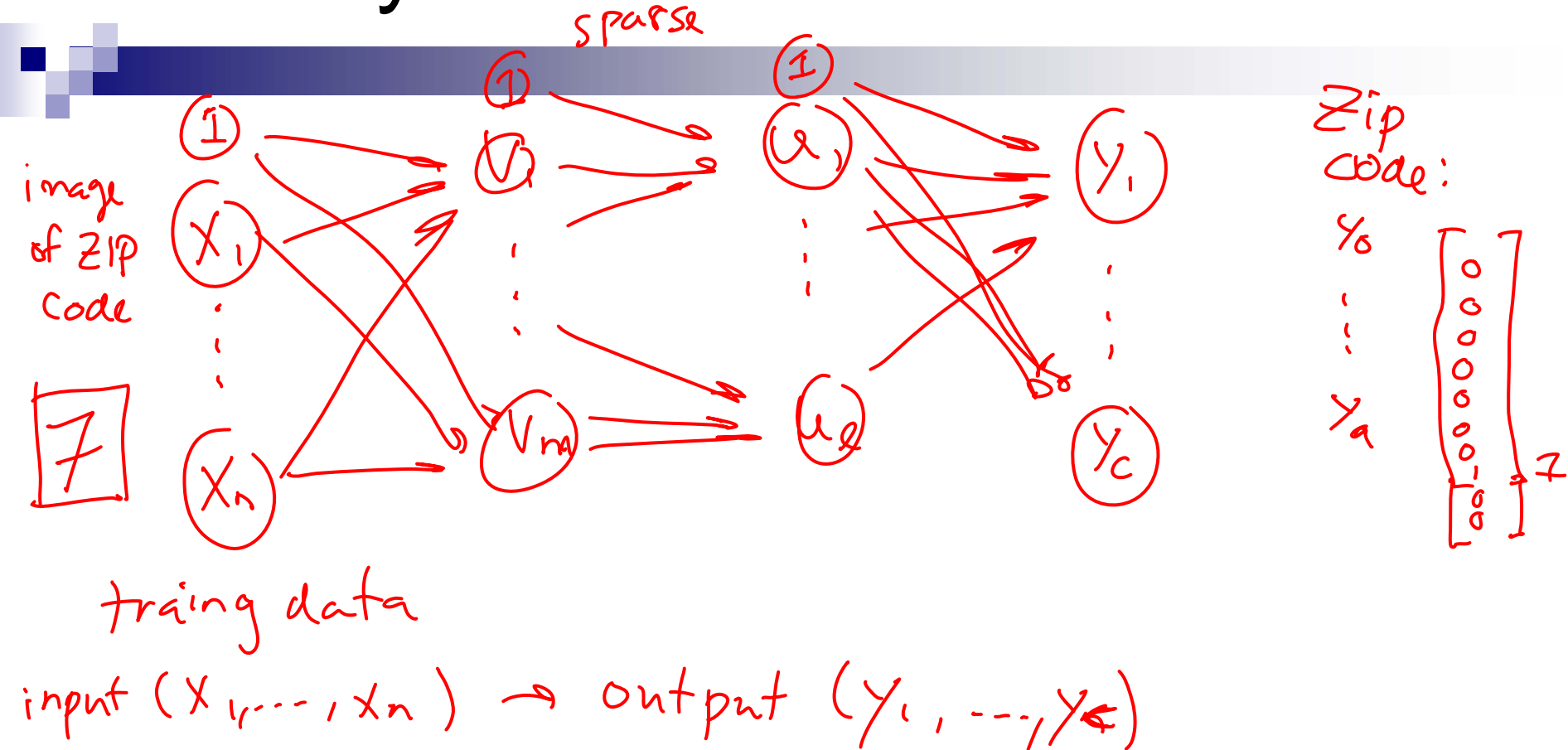
x_i in the perceptron

hidden layer:

$$\frac{\partial out(\mathbf{x})}{\partial w_i^k} = g'\left(\sum_k w_k g\left(\sum_i w_i^k x_i\right)\right) \cdot \frac{\partial \left[\sum_k w_k g\left(\sum_i w_i^k x_i\right)\right]}{\partial w_i^k}$$

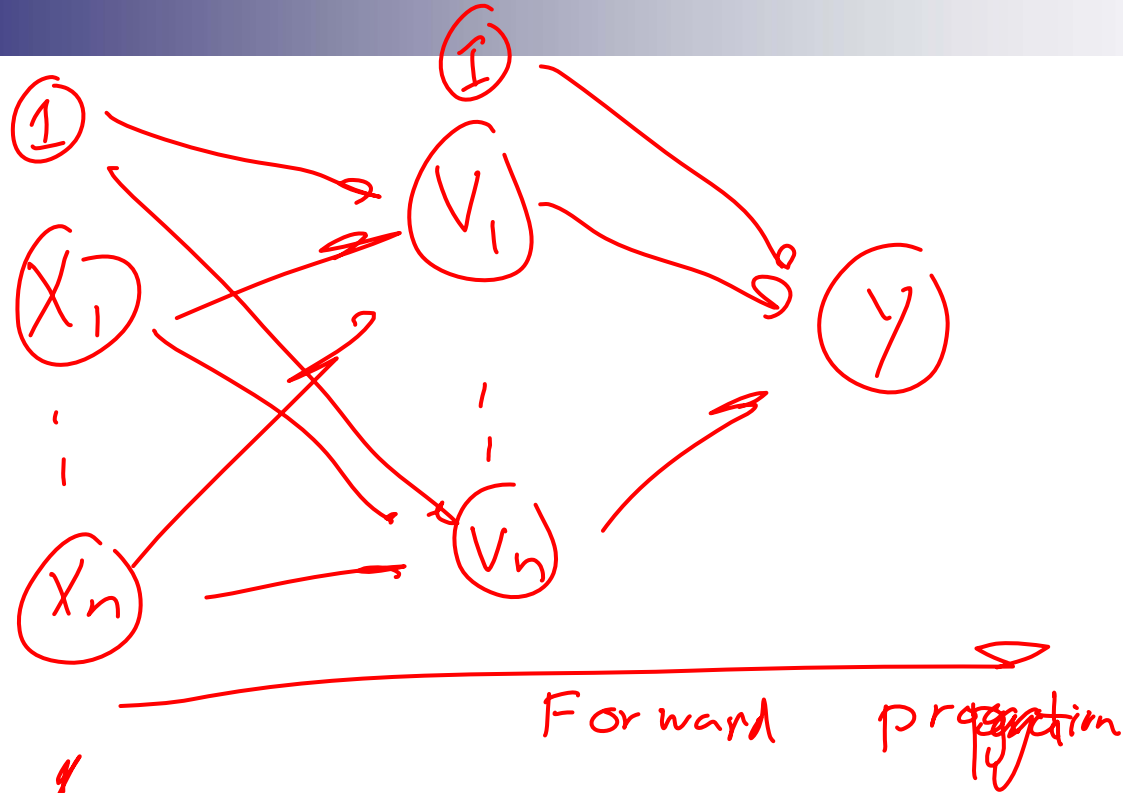
$w_k g'\left(\sum_i w_i^k x_i\right) x_i$

Multilayer neural networks



Forward propagation – prediction

given
new
example
 x



plug x \rightarrow Compute ~~each~~ hidden unit \rightarrow Compute output based on hidden

For multi layer: Recursion!

Back-propagation – learning

Basically gradient descent

For each example x, y

For ward propagation

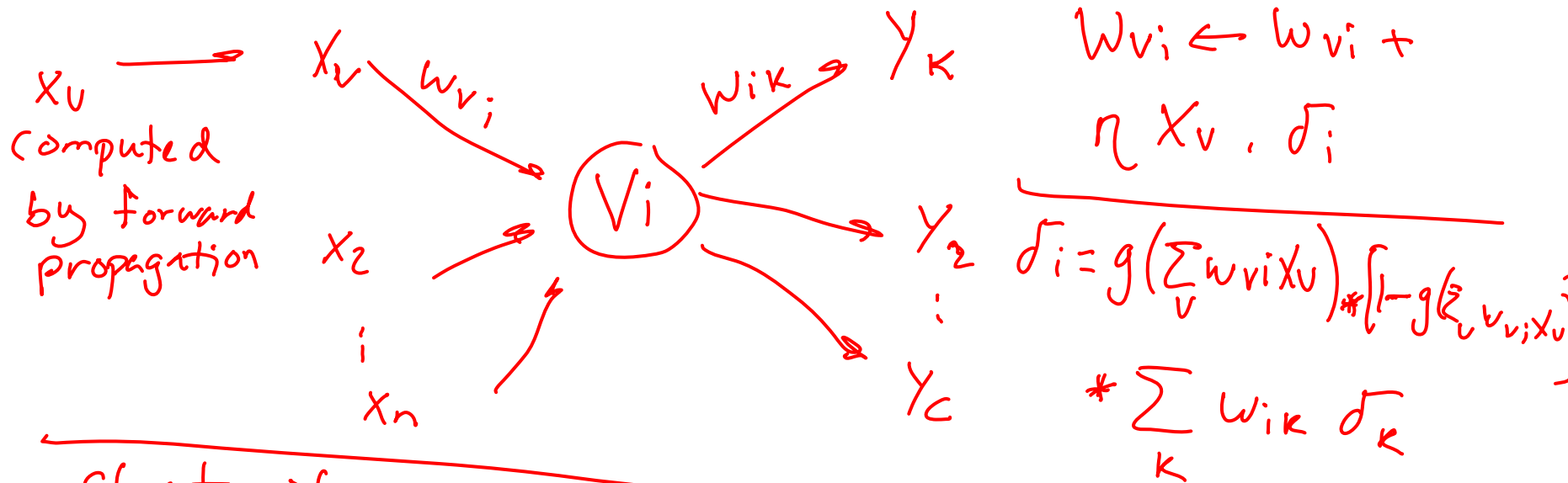
Backwards updating
weights

Back-propagation – learning

$$w_i \leftarrow w_i + \eta x_i \delta$$

$$\delta = [y - g(w_0 + \sum_i w_i x_i)] g'(1 - g)$$

Hidden unit:



Start of recursion for δ :

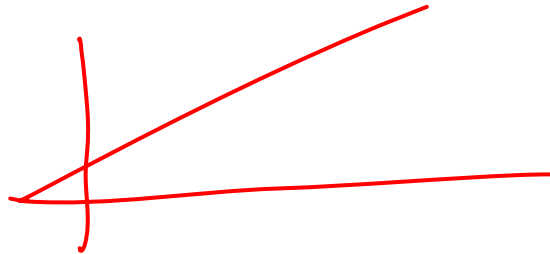
output layer : $\delta = [y - g(\text{output})]$

Many possible response functions

- Sigmoid

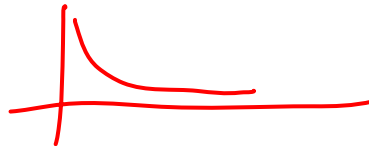


- Linear

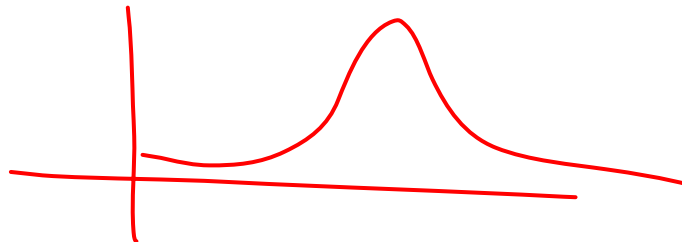


$$g = \sum_i w_i x_i$$

- Exponential



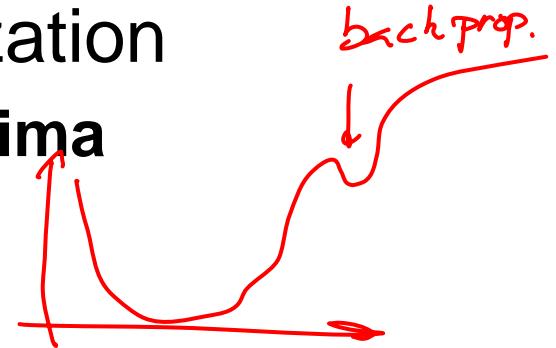
- Gaussian



- ...

Convergence of backprop

- Perceptron leads to convex optimization
 - Gradient descent reaches **global minima**
- Multilayer neural nets **not convex**
 - Gradient descent gets stuck in local minima
 - Hard to set learning rate
 - Selecting number of hidden units and layers = fuzzy process
 - NNs falling in disfavor in last few years
 - As we'll see in second half of the semester, *kernel trick* is a good alternative
 - Nonetheless, neural nets are one of the most used ML approaches

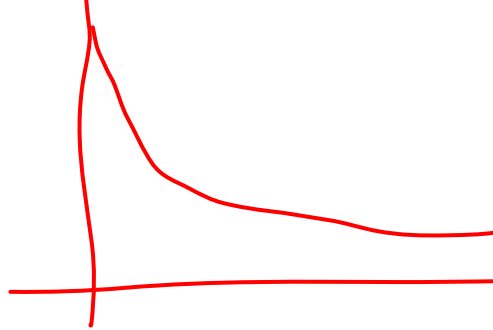


Training set error

- Neural nets represent complex functions
 - Output becomes more complex with gradient steps

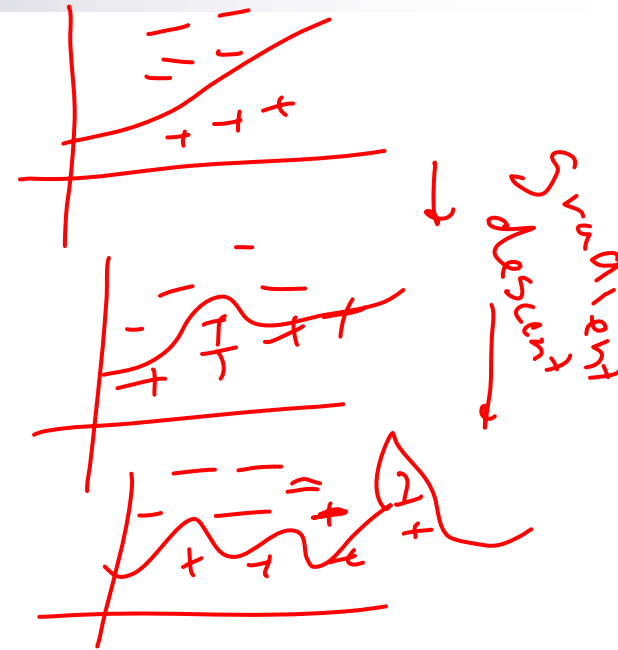
- Training set error

training error ↑

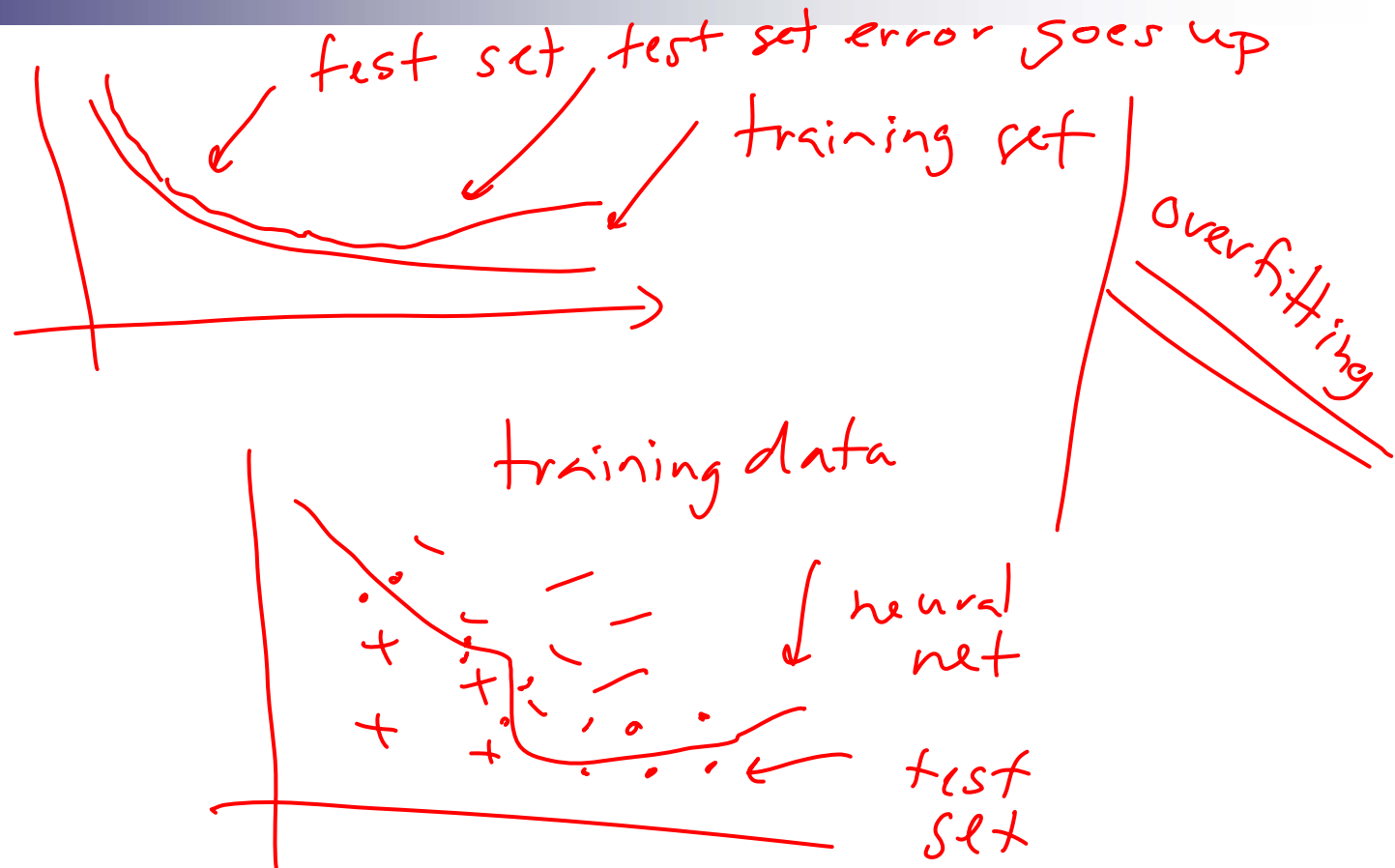


learning curve:

↓ always down



What about test set error?



Overfitting

- Output fits training data “too well”
 - Poor test set accuracy
- Overfitting the training data
 - Related to bias-variance tradeoff
 - One of central problems of ML
- Avoiding overfitting?
 - More training data
 - Regularization
 - Early stopping
 - Next few lectures



What you need to know



■ Perceptron:

- ☐ Representation
- ☐ Perceptron learning rule
- ☐ Derivation

■ Multilayer neural nets

- ☐ Representation
- ☐ Derivation of backprop
- ☐ Learning rule

■ Overfitting

- ☐ Definition
- ☐ Training set versus test set
- ☐ Learning curve
- ☐ Next few classes: strategies for dealing with overfitting