

NeuralNetwo

Storysofar

- Neural networks are universal approximators
 - Can model any odd thing
 - Provided they have the right architecture
- We must train them to approximate any function
 - Specify the architecture
- Networks are trained to minimize total “error” on training set
 - $\sum_{i=1}^n \|f(x_i) - y_i\|^2$

Recap Gradient Descent Algorithm

- **I had to minimize any function**

Fairing Neural Nets by Gradient Descent



Fairing Neural Nets by Gradient Descent



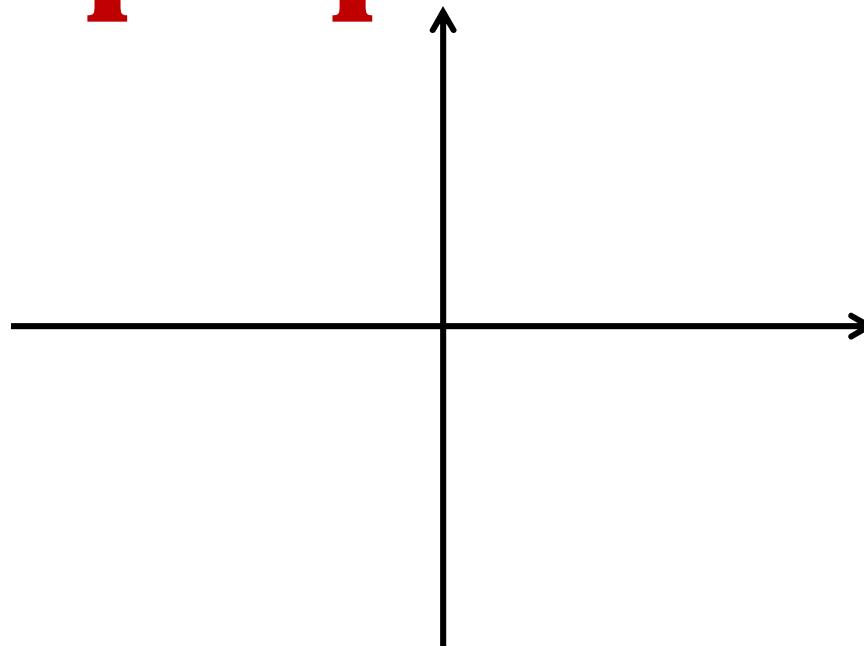
Computing

Computing

Recap Backpropagation for training

- Initialize weights and biases

Bad pop fails to separate where perception ends



- **Brady Raghavan, Slawny '89**
- **Simple problem, 3 training instances, single neuron**
- **Perception training rule trivially finds a perfect solution**

Backprop vs Perception

- Backpropagation using logistic function and

Badprop's Perception

- **Now add a fourth point**
- t

Backprop

- Consider backprop
- Contribution of a weight point to derivative of L

Badprop

- **Thefourthpointat**

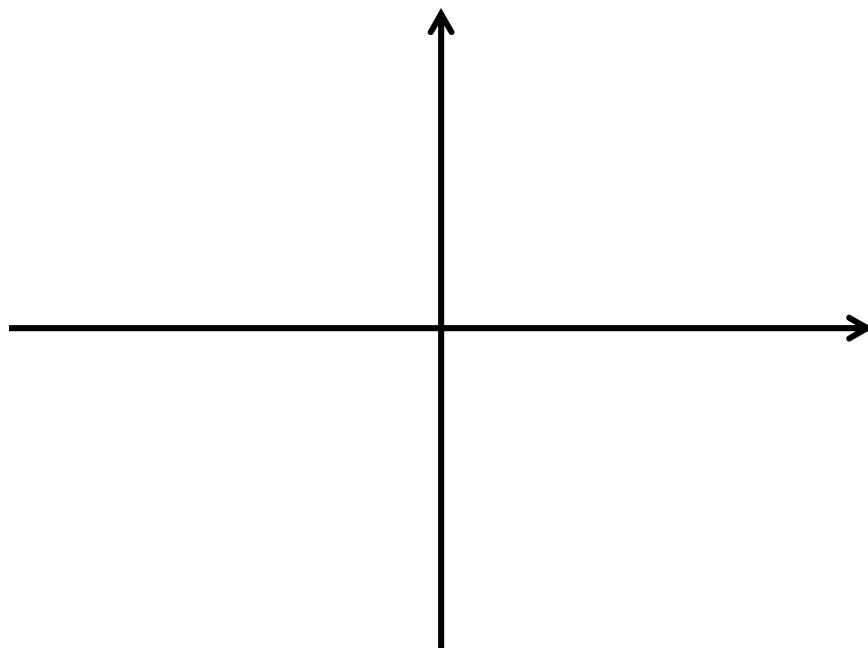
Badprop

- **Solution found by badprop**
- **Does not separate $Q \cup B \cup Q$ even though hs has $hepoints$ a line**

Bad pop fails to separate where perception succeeds

- **Brady Raghavan, Slawny '89**
- **Several**

A more complex problem



Bad pop fails to separate even when possible

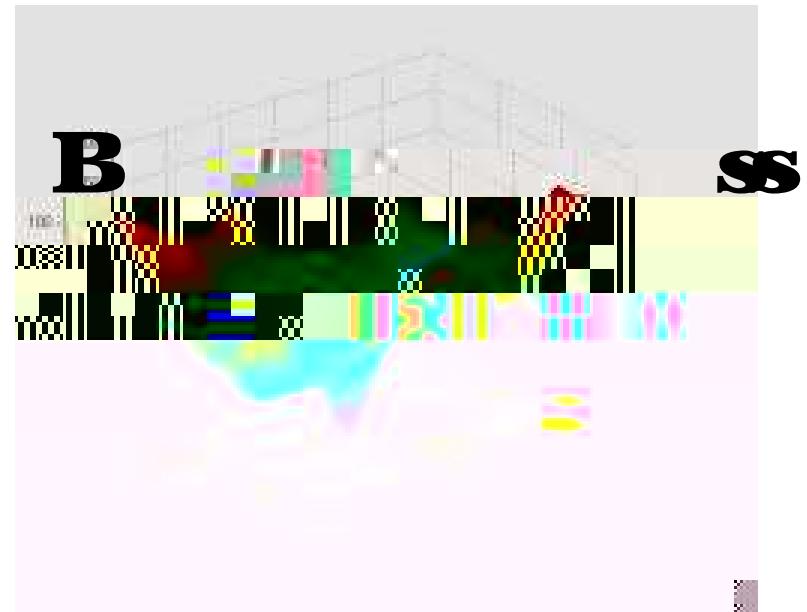
- This is not o dhto ino' c

**Bad pop fails to separate even when
possible**



The Error Surface

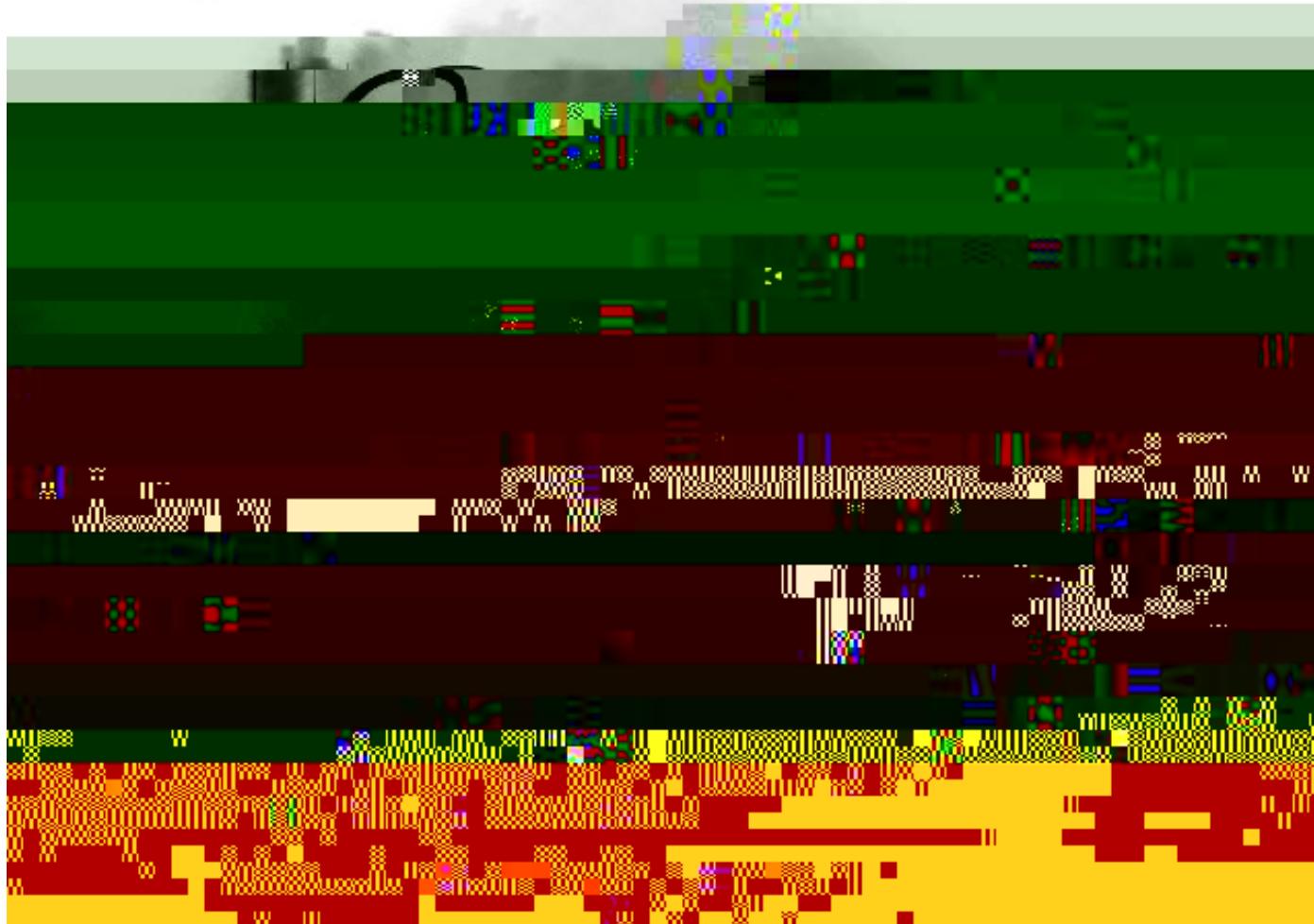
- The example (and statements)
each assumed the loss
objective had a sea **B** ss



Convergence



A quick tour of (convex) optimization



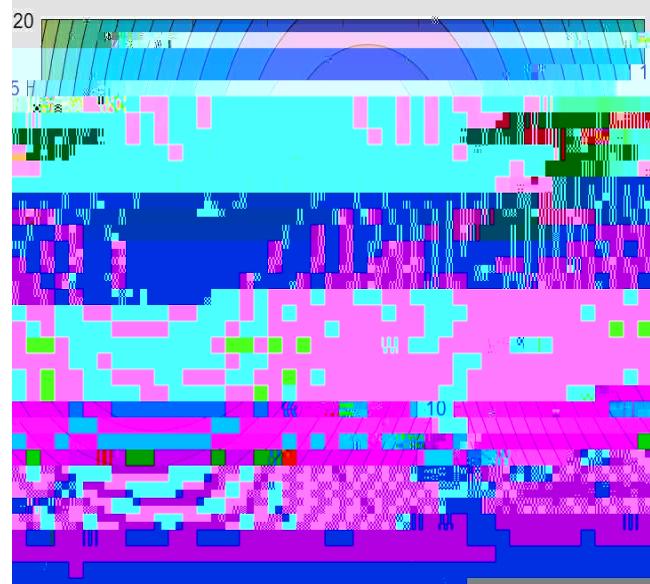
Multivariate Quadratic with Diagonal

“Deserts” are uncoupled

- **The optimum of each coordinate is not affected by the other coordinates**
 - I.e. we could optimize each coordinate independently
-

Dependence on learning rate





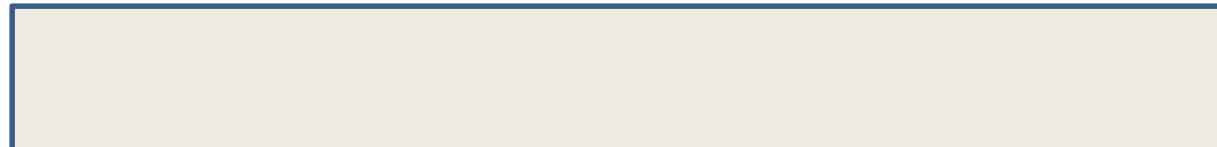
Scaling the axes

•

Returning to our problem



Returning to our problem



Minimization by Newton's method

- **Iterated localized optimization with quadratic approximations**

- **Iterated localized optimization with quadratic approximations**

Minimization by Newton's method

- Iterated localized optimization with quadratic approximations

Minimization by Newton's method

- Iterated localized optimization with quadratic approximations

Issues 2 Thelearninggate

Issues 2 The learning rate

- For complex models such as neural networks the loss function is often non-convex

Storysofar:

condadrvntthods



Story so far: learning rate

- Divergence causing learning rates may not be a bad thing

Derivative insp^{ee}e

Rprop

- **Compute the derivative in the new location**
 - If the derivative has not changed sign from the previous location, increase the step size and take a longer step
 - Δ

Rprop



Rprop

- **Compute the derivative in the new location**
 - If the derivative has changed sign
 - Return to the previous location
 - = +
 - **Shi the step**

Rprop

- **Compute the derivative in the new location**
 - If the derivative has changed sign
 - Return to the previous location
 - $= +$
 - Shrink the step
 - $=$
 - Take the smaller step forward
 - $= -$

QidRop

QuickRqp Modification1

- **It treats each dimension independently**
- **For**

QuickRqp Modification2

- It approximates the second derivative through finite differences

QuickRop

- **Updates are independent for every parameter**
- **For every layer**

QuickRop

- **Updates are independent for every parameter**
- **For every layer**

Qidpqop

A closer look at the convergence problem

-

The momentum methods

- **Maintain a single past state**

- **Inde-**

- **∞**

- aA ismo**

sum -

MomentumUpdate

- The momentum method maintains a running average of all gradients until the current step

$$\Delta$$

Training with momentum

- Initialize all weights \mathbf{W}

MomentumUpdate

- **The momentum method**

MomentumUpdate

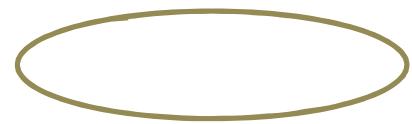
- The momentum method

$$\Delta W$$

Nesterov's Accelerated Gradient

- Change the order of operations
- At any iteration to compute the current step
 - First extend the previous step
 - Then compute the gradient step at the resultant position

Nesterov's Accelerated Gradient



- Change the order of operations
- At any iteration to compute the current step
 - First extend the previous step
 - Then compute the gradient step at the resultant position
 - Add the two to obtain the final step

FairingwithNesterov

- **Initialize weights**

