

Terminology

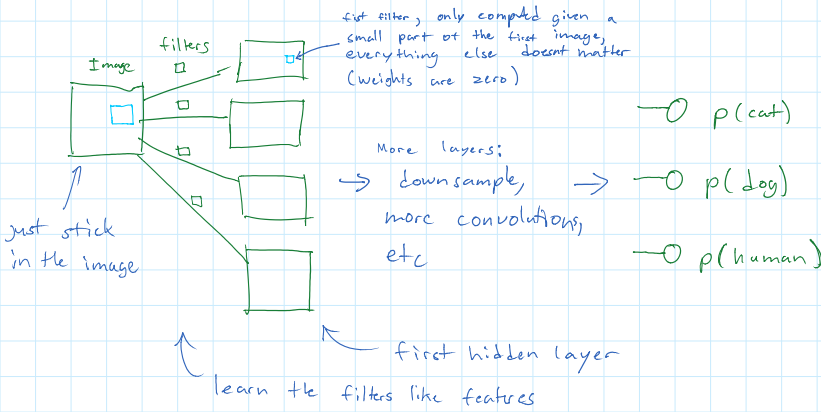
Regression	modelling continuous function
Classification	modelling a discrete categorical function
Supervised Learning	learning a function from $(x_i, y_i)$ pairs
Unsupervised Learning	learning a structure from a data set $\{X\}$
Overfitting	model learns peculiarities of training set
Underfitting	model insufficiently complex
Sigmoid	function that looks like an 'S'
Logistic Function	$f(x) = 1 / (1 + e^{-x})$
Early Stopping	stop training when validation error starts to increase
Momentum	$V_{i+1} = \mu V_i - \epsilon \nabla L(\theta_i)$
Weight Decay	L2 regularization on weights

Announcements

- No class Monday (Thanksgiving)
- Assignment 4 & Marked Assignment 2 due Wednesday (October 15)

Convolutional Neural Networks (CNNs)

- used for image/vision tasks



Refresher On Convolutions

Image (1000x1000)      Filter (2x2)

1	2	1	0	3	...
4	5	-2	0	2	
1	3	4	5	3	
⋮					⋮

Filter (2x2):

1	0
1	0

Step 1:  $1+4 = 5$   
 Step 2:  $2+5 = 7$   
 Step 3:  $1-2 = -1$

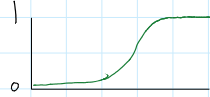
"responds" to situations similar to the filter (dot product is large)

• exactly a NN but with a lot of weights set to 0, i.e. only looking at locality with filters

Activation Functions

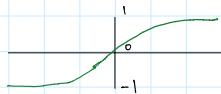
Logistic Function:

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



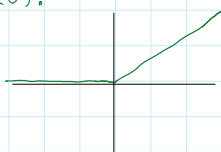
Hyperbolic Tangent:

$$\sigma(x) = \tanh(x)$$



Rectified Linear Units (RELU):

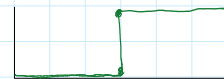
$$\sigma(x) = \max\{0, x\}$$



• this one is better

} These are like Sigmoid functions

These are related to/emulate: threshold function:



aka heavyside step function (non differentiable)

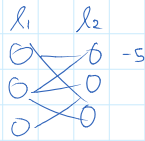
• initialization is important

say  $\text{Rand}(-100, 100)$

: we would rather near "scale" of sigmoid  $(-1, 1)$   
 ("saturated" if you're far away)

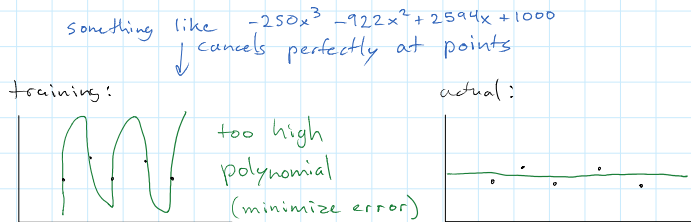
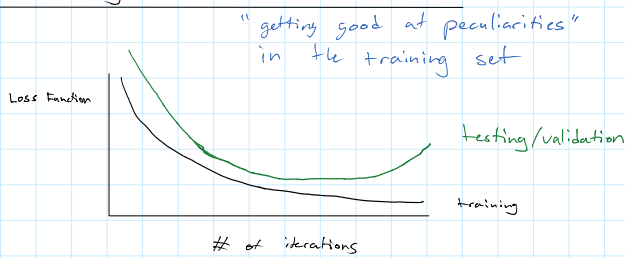
relu : not "scaled" with size

vanishing gradient: gradients less than 1, going back layers, learning will be slow because steps will be small



derivative can flow through another path

## Overfitting and Regularization



to try and fix this:

Regularization:

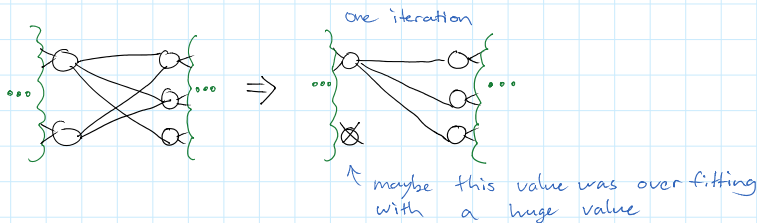
$$L = \sum_i \|y_i - \hat{y}_i\|_2^2 + \lambda \sum_{\text{all weights}} |w_j|^2$$

⇒ "Weight Decay"

- make weights small
- helps prevent overfitting
- could also use L1

## Dropout (a regularization method) (2012)

- During training, at every iteration, "Drop out" (set to zero) each unit with probability  $P$ .
- During prediction, multiply weights by  $(1-P)$



• kind of like taking the geometric mean of a bunch of different neural networks

## Optimization

- Use Stochastic Gradient Descent

$$W_{t+1} = W_t - \alpha_t \nabla L(W_t) \leftarrow \text{gradient of loss}$$

↑ learning rate

kinda like velocity

- Momentum

$$W_{t+1} = W_t + V_{t+1}$$

$$W_{t+1} = W_t - \alpha_t \nabla L(W_t) + \epsilon V_t$$

← retain some momentum  
from the previous step

"Second order method"

Example with a ball

