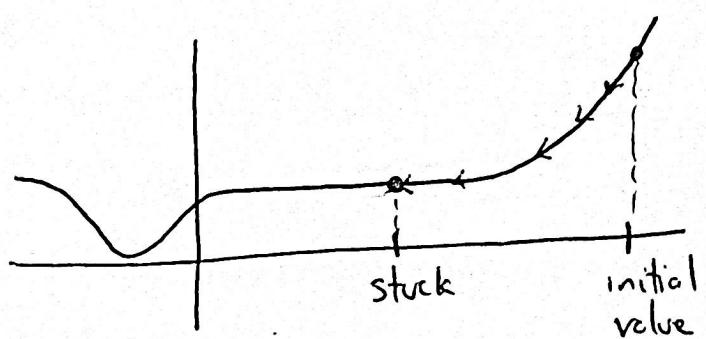


Vanishing Gradient Problem

In practice, the function $J_{W,B}(X,Y)$ for a neural network N is non convex, so there may be flat regions one may "get stuck" in during (stochastic) gradient descent



Factors that affect "flatness" of cost function

① Choice of loss function

SE	lot of flat regions
CE	less
BCE	least

② Choice of activations

the ones we discussed
so far are good
choices

Note: Using several sigmoids
in hidden layers can
cause flatness.

Suggestion: Use ReLU in hidden layers

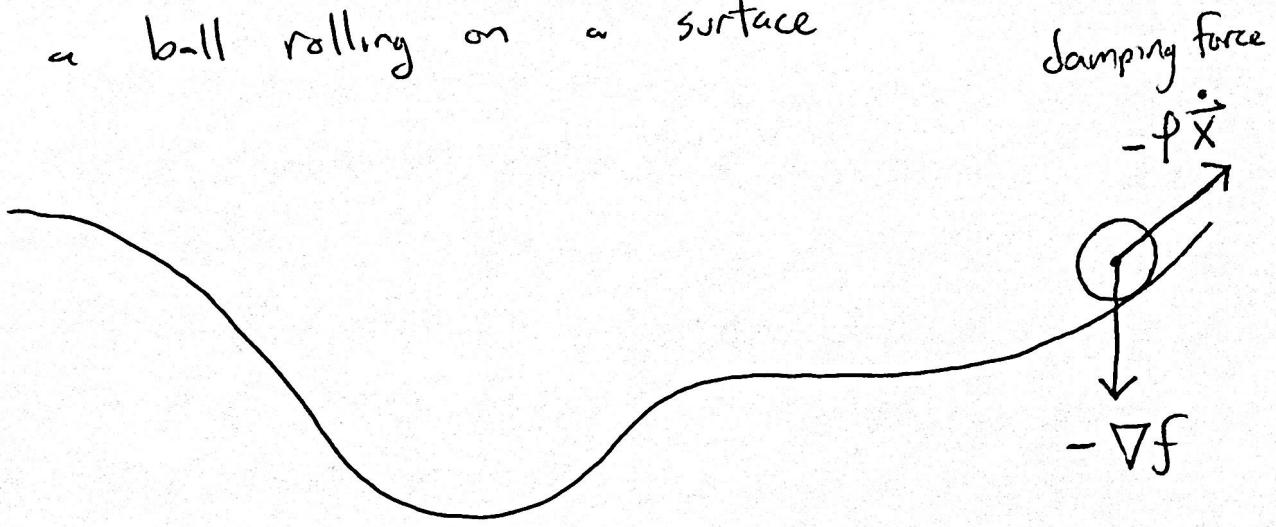
③ Momentum

④ Weight Initialization

Momentum Method

Motivation (from physics)

Consider a ball rolling on a surface



f = potential function $f(\vec{x}) = mgx_3$

$$\sum \vec{F} = m\vec{a} \implies \ddot{\vec{x}} = -f \dot{\vec{x}} - \nabla f(\vec{x})$$

To solve: Express as a system of 1st order ODEs

$$\dot{\vec{v}} = -f \vec{v} - \nabla f(\vec{x})$$

$$\dot{\vec{x}} = \vec{v}$$

Write this as a finite difference system

$$\left(\begin{array}{l} \text{Pick } \Delta t \text{ small} \\ \text{enough so} \\ f \Delta t < 1 \end{array} \right) \quad \begin{aligned} \vec{v}_{n+1} - \vec{v}_n &= -f \vec{v}_n \Delta t - \nabla f(\vec{x}_n) \Delta t \\ \vec{x}_{n+1} - \vec{x}_n &= \vec{v}_{n+1} \Delta t \end{aligned}$$

$$\text{Let } \varepsilon = \Delta t \quad \alpha = 1 - \rho \varepsilon \quad 0 \leq \alpha < 1$$

$$\vec{V}_{n+1} = \alpha \vec{V}_n - \varepsilon \nabla f(\vec{x}_n)$$

$$\vec{x}_{n+1} = \vec{x}_n + \varepsilon \vec{V}_{n+1}$$

Let $\tilde{\vec{V}}\varepsilon = \tilde{\vec{V}}$ (new scale for velocity)

$$\frac{\tilde{\vec{V}}_{n+1}}{\varepsilon} = \alpha \frac{\tilde{\vec{V}}_n}{\varepsilon} - \nabla f(\vec{x}_n)$$

$$\vec{x}_{n+1} = \vec{x}_n + \tilde{\vec{V}}_{n+1}$$

Define $\alpha = \varepsilon^2$

$$\tilde{\vec{V}}_{n+1} = \alpha \tilde{\vec{V}}_n - \nabla f(\vec{x}_n)$$

$$\vec{x}_{n+1} = \vec{x}_n + \tilde{\vec{V}}_{n+1}$$

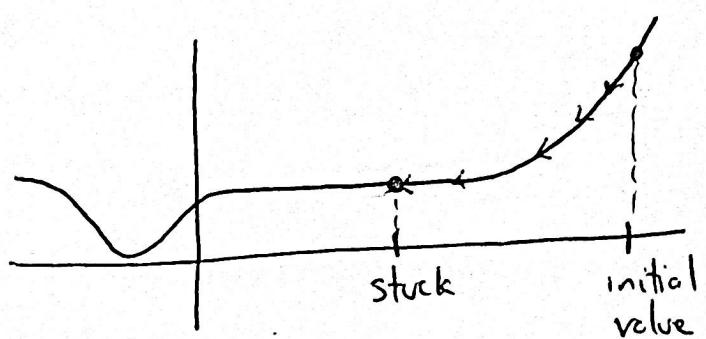
If $\alpha = 0$
this is just
gradient
descent

Weight Initialization

Neural net performance depends heavily on $(W^{(l)})$
weight initialization especially if D is
large. Depends less on bias initialization $(b^{(l)})$
but this depends on architecture.

Vanishing Gradient Problem

In practice, the function $J_{W,B}(X, Y)$ for a neural network N is non convex, so there may be flat regions one may "get stuck" in during (stochastic) gradient descent



Factors that affect "flatness" of cost function

① Choice of loss function

SE	lot of flat regions
CE	less
BCE	least

② Choice of activations

the ones we discussed
so far are good
choices

Note: Using several sigmoids
in hidden layers can
cause flatness.

Suggestion: Use ReLU in hidden layers

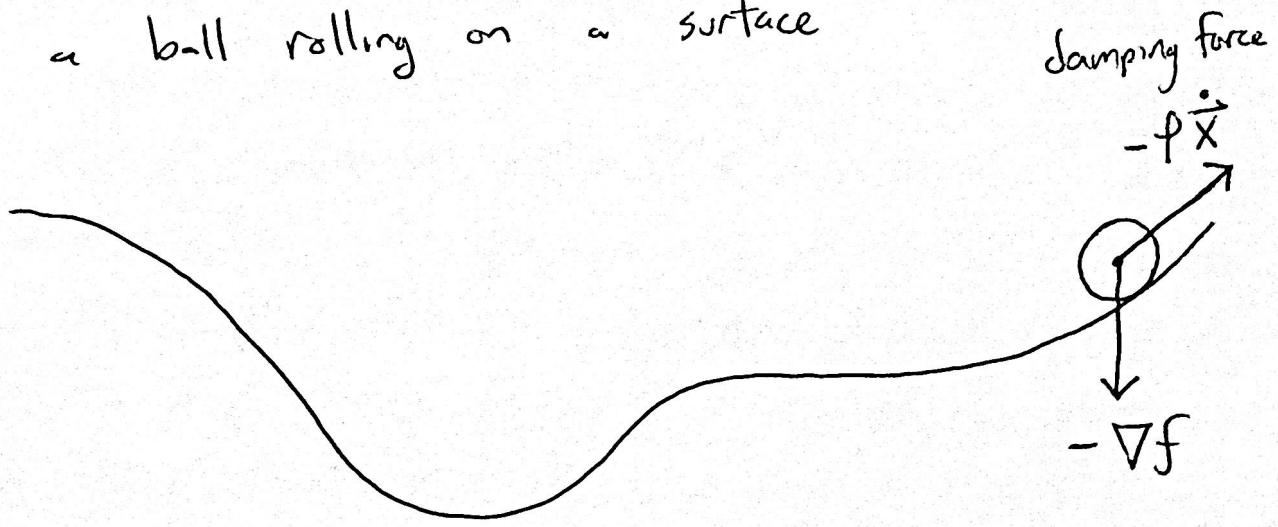
③ Momentum

④ Weight Initialization

Momentum Method

Motivation (from physics)

Consider a ball rolling on a surface



$$f = \text{potential function} \quad f(\vec{x}) = mgx_3$$

$$\sum \vec{F} = m\vec{a} \implies \ddot{\vec{x}} = -f \dot{\vec{x}} - \nabla f(\vec{x})$$

To solve: Express as a system of 1st order ODEs

$$\begin{aligned}\dot{\vec{v}} &= -f \vec{v} - \nabla f(\vec{x}) \\ \dot{\vec{x}} &= \vec{v}\end{aligned}$$

Write this as a finite difference system

$$\left(\begin{array}{l} \text{Pick } \Delta t \text{ small} \\ \text{enough so} \\ f \Delta t < 1 \end{array} \right) \quad \begin{aligned}\vec{v}_{n+1} - \vec{v}_n &= -f \vec{v}_n \Delta t - \nabla f(\vec{x}_n) \Delta t \\ \vec{x}_{n+1} - \vec{x}_n &= \vec{v}_{n+1} \Delta t\end{aligned}$$

$$\text{Let } \varepsilon = \Delta t \quad \alpha = 1 - \rho \varepsilon \quad 0 \leq \alpha < 1$$

$$\vec{V}_{n+1} = \alpha \vec{V}_n - \varepsilon \nabla f(\vec{x}_n)$$

$$\vec{x}_{n+1} = \vec{x}_n + \varepsilon \vec{V}_{n+1}$$

Let $\tilde{\vec{V}}\varepsilon = \tilde{\vec{V}}$ (new scale for velocity)

$$\frac{\tilde{\vec{V}}_{n+1}}{\varepsilon} = \alpha \frac{\tilde{\vec{V}}_n}{\varepsilon} - \nabla f(\vec{x}_n)$$

$$\vec{x}_{n+1} = \vec{x}_n + \tilde{\vec{V}}_{n+1}$$

Define $\alpha = \varepsilon^2$

$$\tilde{\vec{V}}_{n+1} = \alpha \tilde{\vec{V}}_n - \nabla f(\vec{x}_n)$$

$$\vec{x}_{n+1} = \vec{x}_n + \tilde{\vec{V}}_{n+1}$$

If $\alpha = 0$
this is just
gradient
descent

Weight Initialization

Neural net performance depends heavily on $(W^{(l)})$
weight initialization especially if D is
large. Depends less on bias initialization $(b^{(l)})$
but this depends on architecture.

- Claim:
- ① If the initialized values of $W^{(l)}$ are too small, the variance of $s^{(l)}$ decreases at each layer
 - ② If initialized weights are too large, variance $s^{(l)}$ amplifies at each layer

A Solution (More of a suggestion)

Xavier Initialization: Based on the assumption that $\text{Var}(s^{(l)})$ should remain constant at each layer.

For $p(n_0, \dots, n_D)$ entries of $W^{(l)}$ should be sampled from $N(x; 0, \frac{2}{n_{l-1} + n_l})$.

Many other "suggestions" exist

- Read §8.4 of the deep learning book
- Read §6.3 of Deep Learning Architectures

Bias can usually be initialized to $\vec{0}$.

A Way to do SGD on neural net N (for $P(n_0, \dots, n_D)$)

Pick $\alpha > 0$, $\lambda > 0$, $\mu > 0$, max-iters, batch-size
step size reg. parameter momentum

For $0 \leq l \leq D-1$ initialize $W^{(l)}$ by sampling $N(x; 0, \frac{2}{n_{l+1} + n_l})$

$$b^{(0)} = \overline{0}$$

initialize $V_w^{(l)} = np.zeros(W^{(l)}.shape)$

$V_b^{(l)} = np.zeros(b^{(l)}.shape)$

epochs = 0

while epochs \leq max-iters:

randomly pick $B \subseteq \{1, \dots, m\}$ with $|B| = \text{batch-size}$

for $i \in B$:

compute $[(x_i^{(0)}, s_i^{(0)}), (x_i^{(1)}, s_i^{(1)}), \dots, (x_i^{(D-1)}, s_i^{(D-1)}), (x_i^{(0)})]$

compute $\delta_i^{(D-1)}, \dots, \delta_i^{(0)}$

compute $\frac{\partial J}{\partial W^{(l)}}, \frac{\partial J}{\partial b^{(l)}}$ (batch approximations)

for $0 \leq l \leq D-1$:

$$V_w^{(l)} = \mu V_w^{(l)} - \alpha \frac{\partial J}{\partial W^{(l)}}$$

$$V_b^{(l)} = \mu V_b^{(l)} - \alpha \frac{\partial J}{\partial b^{(l)}}$$

$$W^{(l)} = W^{(l)} + V_w^{(l)}$$

$$b^{(l)} = b^{(l)} + V_b^{(l)}$$

append cost to list of costs

return W, B, costs

Convolutional Neural Networks

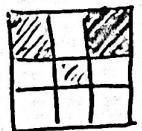
Idea: NN architecture should be chosen to complement the type & structure of input data if possible

E.g. ① \vec{x} = statistics on houses

(no real structure so use fully connected architecture)

② \vec{x} = flattened vector of pixel info on an image

E.g. image is 3×3



$\rightarrow (2, 0, 0, 0, 1, 0, 3, 0, 0)$

then x_4 is "nearby" x_1 and x_7

Q: How can we incorporate this extra structure?

A: "Convolution" and "Pooling"

Def: A discrete 1-dim'l signal is a sequence

$$y = (y_i)_{i \in \mathbb{Z}} \quad \text{of real numbers}$$

and is L¹-finite if $\sum_k |y_k| < \infty$

compact support if $\exists N \text{ s.t.}$

$$|k| > N \implies y_k = 0$$

A 1-dim'l kernel (or filter) is a compact support signal of "weights"

$$w = (w_i)_{i \in \mathbb{Z}}$$

Note: Usually want w to be a probability distribution so $\sum_i w_i = 1$

The convolution of y and w is

$$z = y * w$$

$$\text{where } z_j = \sum_{k \in \mathbb{Z}} y_{j+k} w_k$$

(often called the
"cross correlation" in
signal processing)

(well defined because
 w has compact support)

E.g. (Moving average signal)

$$w_i = \begin{cases} \frac{1}{2} & i=0,1 \\ 0 & \text{else} \end{cases}$$

$$(y * w)_i = \sum_{k \in \mathbb{Z}} y_{i+k} w_k = \frac{y_i + y_{i+1}}{2}$$

$$\text{Similarly } w_i = \begin{cases} \frac{1}{M} & i=0, \dots, M-1 \\ 0 & \text{else} \end{cases}$$

$$(y * w)_i = \frac{y_i + y_{i+1} + \dots + y_{i+M-1}}{M}$$

$$\text{E.g. } y = \begin{pmatrix} 1, 2, 3, 4, 5, 6 \\ 0 \end{pmatrix} \quad w = \begin{pmatrix} \frac{1}{2}, \frac{1}{2}, \\ 0 \end{pmatrix}$$

$$y * w = \begin{pmatrix} \frac{1}{2}, \frac{3}{2}, \frac{5}{2}, \frac{7}{2}, \frac{9}{2}, \frac{11}{2}, 3 \\ -1 \quad 0 \quad 1 \quad 5 \end{pmatrix}$$

When using this to process signals in a NN layer,
 we usually ignore the first and last component.
 We may also ignore some parts on the inside
 i.e. may express this as

$$\begin{pmatrix} \frac{1}{2} & \frac{1}{2} & & & & & \\ & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & & & \\ & & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & & \\ & & & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \\ & & & & \frac{1}{2} & \frac{1}{2} & \\ & & & & & \frac{1}{2} & \\ & & & & & & \end{pmatrix} \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{pmatrix} = \begin{pmatrix} \frac{3}{2} \\ \frac{5}{2} \\ \frac{7}{2} \\ \frac{9}{2} \\ \frac{11}{2} \end{pmatrix} \in \mathbb{R}^5$$

maybe we only want to keep $(\frac{3}{2}, \frac{7}{2}, \frac{11}{2})$

Could express this as

$$\begin{pmatrix} \frac{1}{2} & \frac{1}{2} & & & & & \\ & \frac{1}{2} & \frac{1}{2} & & & & \\ & & \frac{1}{2} & \frac{1}{2} & & & \\ & & & \frac{1}{2} & \frac{1}{2} & & \\ & & & & \frac{1}{2} & \frac{1}{2} & \\ & & & & & \frac{1}{2} & \\ & & & & & & \end{pmatrix} \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{pmatrix} = \begin{pmatrix} \frac{3}{2} \\ \frac{7}{2} \\ \frac{11}{2} \end{pmatrix} \in \mathbb{R}^3$$

As a neural network layer

A convolutional layer (in a neural network with 1D signals) with kernel $W = (w_1, \dots, w_k)$ and stride r is a layer with $F: \mathbb{R}^{n_e} \rightarrow \mathbb{R}^{n_{e+1}}$ given by

$$F(x) = W^{(1)}x + b^{(1)}$$

where

$$W^{(1)} = \left(\begin{array}{c} w_1 \dots w_k \\ \xrightarrow[r]{} w_1 \dots w_k \\ \xrightarrow[r]{} w_1 \dots w_k \\ \ddots \\ \textcircled{O} & & w_1 \dots w_k \end{array} \right)$$

A pooling layer (with 1D inputs) is a

layer which applies a map

$$P: \mathbb{R}^N \rightarrow \mathbb{R}^K \quad K \leq N$$

defined as follows:

Let $\overline{\Pi}$ be a partition of $\{1, \dots, N\}$

i.e. $\overline{\Pi} = \{B_1, \dots, B_K\}$ s.t.

$$B_i \subseteq \{1, \dots, N\} \quad \forall i$$

$$B_i \cap B_j = \emptyset \text{ when } i \neq j$$

$$\bigcup_{i=1}^k B_i = \{1, \dots, N\}$$

define $P(x_1, \dots, x_n) = \left(\max_{i \in B_1} x_i, \dots, \max_{i \in B_k} x_i \right)$

usually $B_1 = \{1, 2, \dots, b_1\}$ $B_2 = \{b_1 + 1, b_1 + 2, \dots, b_2\}, \dots$

Note: Alternatively we could use min or avg

E.g. $P: \mathbb{R}^{10} \rightarrow \mathbb{R}^4$ $B = \{\{1, 2, 3\}, \{4, 5\}, \{6, 7\}, \{8, 9, 10\}\}$

$$P: \begin{bmatrix} 2 \\ 4 \\ 3 \\ 7 \\ 6 \\ 4 \\ 1 \\ 0 \\ 2 \\ 9 \end{bmatrix} \mapsto \begin{bmatrix} 4 \\ 7 \\ 4 \\ 9 \end{bmatrix}$$

Two Dimensional Signals

A discrete 2-dim'l signal is $y = (y_{ij})_{i,j \in \mathbb{Z}}$ and is

L'-finite if $\sum_{i,j \in \mathbb{Z}} |y_{ij}| < \infty$

finite support if $\exists N \text{ s.t. } |k| > N \text{ and } |l| > N$

$$\implies y_{k,l} = 0$$

E.g. Grayscale images can be represented as 2-dim'l signals with finite support.

A 2-dim'l kernel is a 2-dim'l signal w
of finite support

(Again we usually want this to be a probability
distribution so $\sum_{i,j} w_{ij} = 1$)

The convolution of y and w is

$$z = y * w \quad \text{where } z_{ij} = \sum_{r,k \in \mathbb{Z}} y_{i+r, j+r} w_{k,r}$$

(often z is called a feature map)

E.g. Consider $w_{ij} = \begin{cases} \frac{1}{4} & (i,j) = (0,0), (1,0), (0,1), (1,1) \\ 0 & \text{else} \end{cases}$

let $z = y * w$. Then

$$z_{ij} = \frac{y_{i,j} + y_{i+1,j} + y_{i,j+1} + y_{i+1,j+1}}{4}$$

Notation Denote $y * w = (w(y)) = ((y))$ if w is understood

Denote by $T_{a,b}$ the translation^{by (a,b)} operation so

$$(T_{a,b}(y))_{ij} = y_{i-a, j-b}$$

Proposition (Translational Equivariance of convolution)

$$C_w(T_{a,b}(y)) = T_{a,b}(C_w(y)) \quad \text{for } \forall w \text{ with compact support}$$

Proof:

$$\begin{aligned} C_w(T_{a,b}(y))_{ij} &= \sum_{k,r \in \mathbb{Z}} T_{a,b}(y)_{i+k, j+r} w_{k,r} \\ &= \sum_{k,r \in \mathbb{Z}} y_{i+k-a, j+r-b} w_{k,r} \\ &= (y * w)_{i-a, j-b} \\ &= T_{a,b}(C_w(y)) \end{aligned}$$

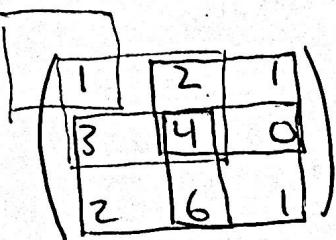
◻

Convolutional layer with 2D inputs

Let $I \in \mathbb{R}^{H \times K}$ (grayscale image represented by matrix) and consider a kernel $W \in \mathbb{R}^{h \times k}$ ($h \leq H, k \leq K$).

Define the convolution of W with I of stride $(1,1)$

$$\text{by } (I *_{(1,1)} W)_{ij} = \sum_{s,r} I_{i+s,j+r} W_{s,r} \quad \begin{matrix} 1 \leq i \leq H-h+1 \\ 1 \leq j \leq K-k+1 \end{matrix}$$

E.g.  $\left(\begin{array}{ccc|c} 1 & 2 & 1 & \\ 3 & 4 & 0 & \\ 2 & 6 & 1 & \end{array}\right) * \begin{pmatrix} -1 & 2 \\ 1 & 0 \end{pmatrix} = \begin{pmatrix} 6 & 4 \\ 7 & 2 \end{pmatrix}$.

More generally for stride (a,b) define

$$(I *_{(a,b)} W)_{ij} = \sum_{s,r} I_{\cancel{a(i-1)+1+s}, \cancel{b(j-1)+1+r}} W_{s,r}$$

Question: What are the bounds on i, j ?

Answer: if $a=b$ and $h=k$ (square kernel)

then output has size $\left(\frac{H-h}{a} + 1, \frac{K-k}{a} + 1\right)$

E.g. $a=2, h=2$

$$I = \begin{pmatrix} 1 & 2 & 3 & 2 \\ 4 & 0 & -1 & 3 \\ 0 & 1 & 2 & 1 \\ 6 & 2 & -1 & 4 \end{pmatrix} \quad W = \begin{pmatrix} -1 & 4 \\ 1 & 5 \end{pmatrix} \quad I *_{(2,2)} W = \begin{pmatrix} -1 & 4 \\ 1 & 5 \end{pmatrix} \in \mathbb{R}^{2 \times 2}$$

A feature map with 2D input $I \in \mathbb{R}^{H \times K}$
with stride (a, a) and kernel $W \in \mathbb{R}^{h \times h}$,
bias $b \in \mathbb{R}^{\frac{(H-h)}{a}+1 \times \frac{(K-h)}{a}+1}$ has output

$$G(I *_{(a,a)} W + b) \in \mathbb{R}^{\left(\frac{H-h}{a}+1\right) \times \left(\frac{K-h}{a}+1\right)}$$

~~for some activation G (almost always ReLU)~~

Note: By "unwinding" I it is possible to write

$I *_{(a,a)} W$ as a matrix + vector product

(see stackexchange 2D convolution as matrix product)

Convolution for Colored Images

A colored image I consists of 3 matrices

$I_r, I_g, I_b \in \mathbb{R}^{H \times K}$ & i, j the RGB values at

pixel (i, j) is $(I_r[i, j], I_g[i, j], I_b[i, j])$.

Represent I as a triply indexed quantity

(a 3-tensor) so $I_{i,j,c}$ $\begin{array}{l} 1 \leq i \leq H \\ 1 \leq j \leq K \\ 1 \leq c \leq 3 \end{array}$

In general, a 3-tensor can be thought of as
 a map $\{1, \dots, n_1\} \times \{1, \dots, n_2\} \times \{1, \dots, n_3\} \rightarrow \mathbb{R}$
 for some n_1, n_2, n_3 .

Feature maps with 3D input

Given a colored image $I = (I_{i,j,c})$, let

$W = (W_{p,s,c})$ be a 3-tensor kernel

$$\begin{array}{l} 1 \leq i \leq H \\ 1 \leq j \leq K \end{array} \quad \begin{array}{l} 1 \leq p \leq h \\ 1 \leq s \leq k \end{array}$$

$$1 \leq c \leq 3$$

Define $I^{*(W)}$ to be the 3-tensor with

$$(I^{*(W)})_{i,j,c} = \sum_s \sum_p I_{i+p, j+s, c} W_{p,s,c}$$

(could also define this for arbitrary stride)

\therefore a feature map with 3D input $\xrightarrow{*}$. takes
 I as input and outputs

$$G(I^*W + b)$$

where G is an activation & b has the
 same shape as I^*W .

For each layer we would like to use multiple feature maps.

Let (W_t, b_t) be a list of kernels and
 $1 \leq t \leq f$ biases (all of same shape)
 $f = \# \text{ features}$

Think of each kernel as detecting a certain "feature"
 e.g. horizontal lines / vertical lines

Given (W_t, b_t) $1 \leq t \leq f$ we can encode this
 as a 4-tensor $W_{p,s,c,t}$ $\begin{array}{l} 1 \leq p \leq h \\ 1 \leq s \leq h \\ 1 \leq c \leq 3 \\ 1 \leq t \leq f \end{array}$

Given a colored image I , we get features

$$z_{i,j,c,t} = G \left[\left(\sum_{s,p} I_{i+p, j+s, c} W_{p,s,c,t} \right) + b_{i,j,c,t} \right]$$

\therefore later layers need to take 4-tensors as inputs (and output a 4-tensor)

Given a 4-tensor input $Z_{i,j,c,t}^{(l-1)}$

and $W_{p,s,c,t,t'}^{(l)}$ $b_{i,j,c,t'}^{(l)}$ (weights & biases)

$t \leftrightarrow$ input features (layer $l-1$)

$t' \leftrightarrow$ output features (layer l)

Define $Z_{i,j,c,t'}^{(l)} = G \left[\left(\sum_{s,p,t} Z_{i+p,j+s,c,t}^{(l-1)} W_{p,s,c,t,t'}^{(l)} \right) + b_{i,j,c,t'}^{(l)} \right]$

This is a convolutional layer.

Next: Pooling layers

Note: The last formula written down last time should be

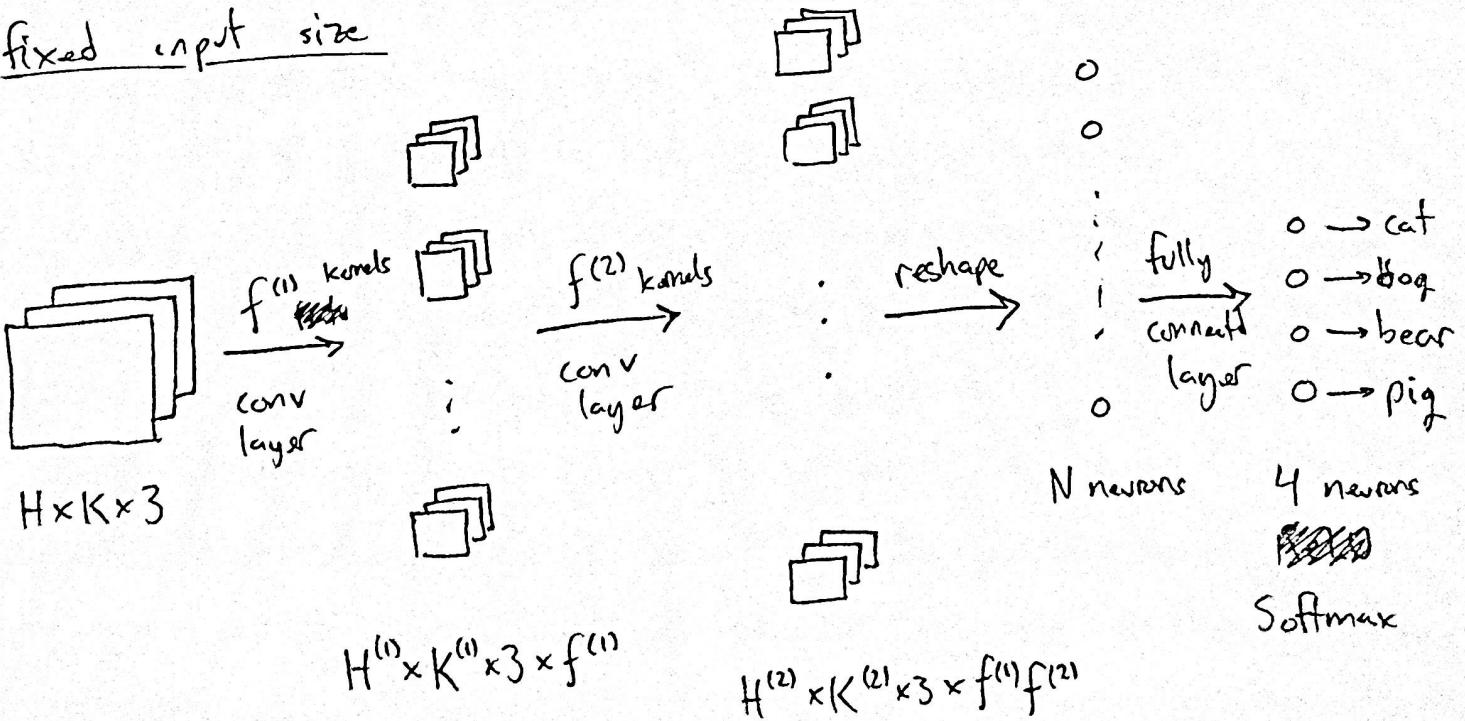
$$Z_{i,j,c,(fin, \text{fout})}^{(l)} = G \left[\left(\sum_{s,p} Z_{i+p, j+s, c, fin}^{(l-1)} W_{p,s,c,(fin, \text{fout})}^{(l)} \right) + b_{i,j,c,(fin, \text{fout})}^{(l)} \right]$$

and then "flatten" index (fin, fout) to new index f

so # features in layer ~~ℓ~~ ℓ is finfout

\therefore # features increases in a convolutional layer.

An example of a CNN to predict cat/dog/bear/pig -
fixed input size



Note: N will be very large ... this is not the most efficient approach.

In using convolution layers, we find

$$H^{(2)} \leq H^{(1)} \leq H$$

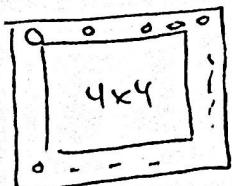
$$K^{(2)} \leq K^{(1)} \leq K$$

$$f^{(1)} f^{(2)} \geq f^{(1)} \geq 1 \quad (\# \text{ features grows})$$

Sometimes we may want to preserve the # pixels in features. Can fix this with "padding" (adding some # of 0's surrounding image)

E.g. $(4 \times 4) *_{(1,1)} (2 \times 2) = (3 \times 3)$

with 1-padding



$$*_{(1,1)} (2 \times 2) = (3 \times 3)$$

If we want to use p-padding to preserve image size with $h \times h$ filter & $H \times K$ image

$$H \times K = (H + 2p - h + 1) \times (K + 2p - h + 1)$$

so $p = \frac{h-1}{2}$ (so h must be odd)

\therefore usually people use $h = \text{odd}$ (often $h=3$)

Padding to preserve size is called "same" padding.

Pooling Layers

Local pooling Given image $I \in \mathbb{R}^{H \times K}$ / local pooling consists of partitioning I and computing some local statistic (max, min, avg)

E.g.

1	2	3	1
1	5	2	7
10	1	2	4
8	9	1	8

$\xrightarrow{(2,2)}$ max pool $\begin{pmatrix} 5 & 7 \\ 10 & 8 \end{pmatrix}$

Gives a way to reduce image size if needed

for a 4-tensor $Z_{i,j,c,f}$ pool for each fixed c,f to get a 4-tensor $Z'_{i,j,c,f}$ with smaller bounds on i,j & same bands on c,f .

Global Pooling (Local pooling where we pool over the whole image)

Given 4-tensor $Z_{i,j,c,f}$ take (max, min, avg) over all i,j to get a 2-tensor $Z'_{c,f}$

E.g.. Redo the previous example but better

