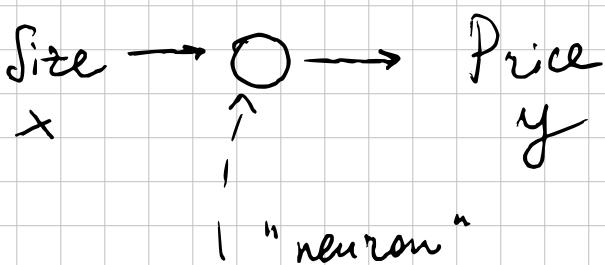


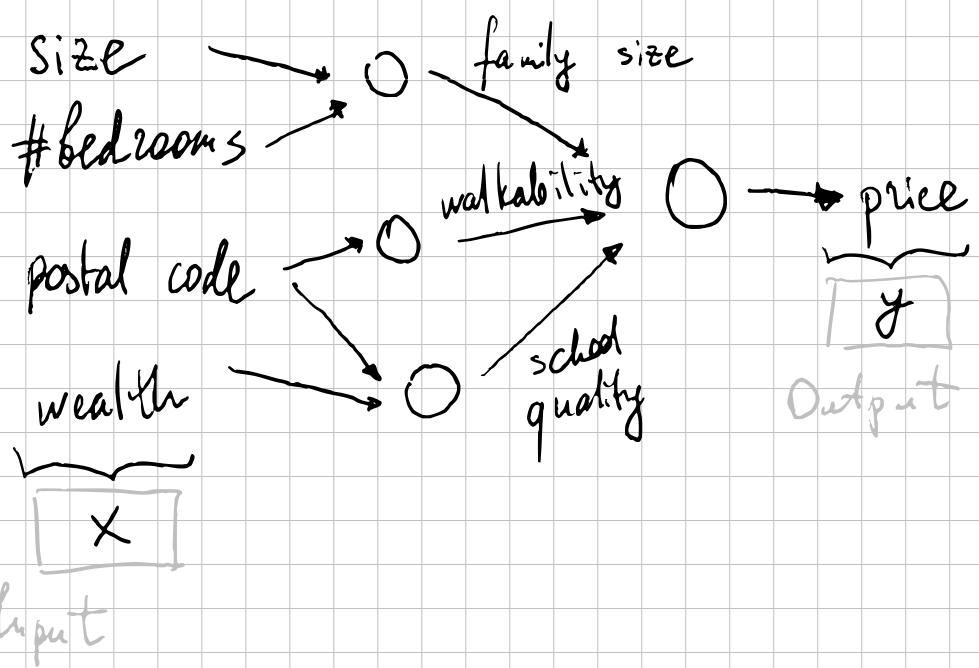
Housing Price Prediction



function that fits the price may be represented as simplest neural network



Another example is



Supervised Learning

Input (x)	Output (y)	Application
Home features	Price	Real estate
Ad, user info	Click on ad? (0/1)	Online advertising
Image	Object	Photo tagging
Audio	Text transcript	Speech recognition
English	Chinese	Machine translation
Image, Radar info	Position of other cars	Autonomous driving

For ad we may use

- Standard NN

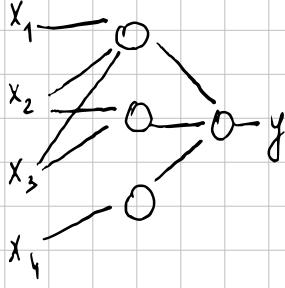
For images

- Convolutional NN (CNN)

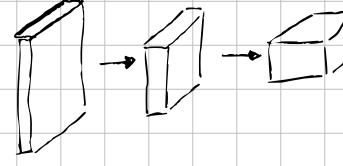
For sequential data

- Recurrent NN (RNN)

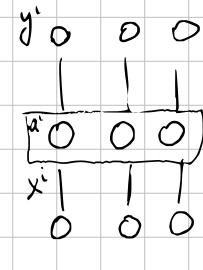
Neural Network Examples



Standard NN



Convolutional NN



Recurrent NN

Supervised Learning

Structured Data

(databases of data)

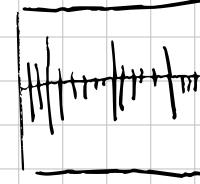
every feature has
very well-defined
meaning

Size #bedrooms ... Price

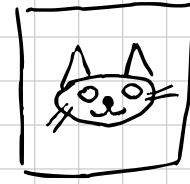
100 2 1000

150 1 1200

Unstructured Data



Audio



Image

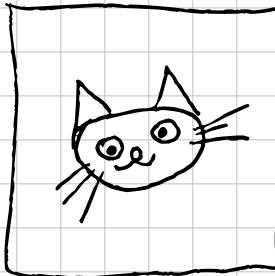
Four scores 2
seven years ago...

Text

Drivers of Deep Learning progress

- Data
- Computation
- Algorithms

Binary classification



$\rightarrow 1(\text{cat}) \text{ vs } 0(\text{non-cat})$

y

In binary classification our goal is to learn a classifier that can input an image, represented by a feature vector X , and predict whether the corresponding label y is 0 or 1.

Notation

- (X, y) where n_x
 $X \in \mathbb{R}^{n_x}$ feature vector
 $y \in \{0, 1\}$ label

- m training examples: $(X^{(1)}, y^{(1)})$,
 $(X^{(2)}, y^{(2)}) \dots (X^{(m)}, y^{(m)})$

- m_{train} training samples
- m_{test} test examples

$$X = \begin{bmatrix} | & | & | \\ X^{(1)} & X^{(2)} & \dots & X^{(m)} \\ | & | & | & | \end{bmatrix} \quad \begin{array}{c} \uparrow \\ n_x \\ \downarrow \end{array}$$

$\xleftarrow{\hspace{1cm}} m \xrightarrow{\hspace{1cm}}$

$$X \in \mathbb{R}^{n_x \cdot m}$$

$$y = [y^{(1)}, y^{(2)} \dots y^{(m)}]$$

$$y \in \mathbb{R}^{1 \cdot m}$$

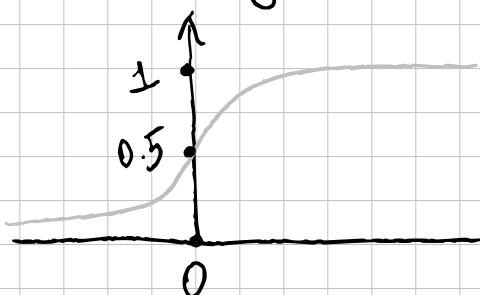
Logistic regression

Given X , we want to predict
 \hat{y} which is an estimate of y ,
i.e. $\hat{y} = P(y=1 | X)$

$$x \in \mathbb{R}^{n_x}$$

Parameters: $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$

Output: $\hat{y} = G(w^T x + b)$



$$G(z) = \frac{1}{1+e^{-z}}$$

- If z large

$$G(z) \approx \frac{1}{1+0} \approx 1$$

- If z small or

- large neg. number

$$G(z) \approx \frac{1}{1+e^{-z}} = \frac{1}{1+\text{Big}} \approx 0$$

Logistic Regression Cost Function

We need to find

$$\hat{y} = g(w^T x + b),$$

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

Given $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$,
 want $\hat{y}^{(i)} \approx y^{(i)}$

Loss (error) function:

$$L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$$

\downarrow Example

Function we need to define that
 shows how good predicted
 output (\hat{y}) compared to the
 ground truth (y)

In logistic regression usually
 uses:

$$L(\hat{y}, y) = - (y \log \hat{y} + (1-y) \log(1-\hat{y}))$$

If $\bullet y=1 \quad L(\hat{y}, y) = - \log \hat{y}$
 $\bullet y=0 \quad L(\hat{y}, y) = - \log(1-\hat{y})$

Cost function: $J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$

$$= - \frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)})]$$

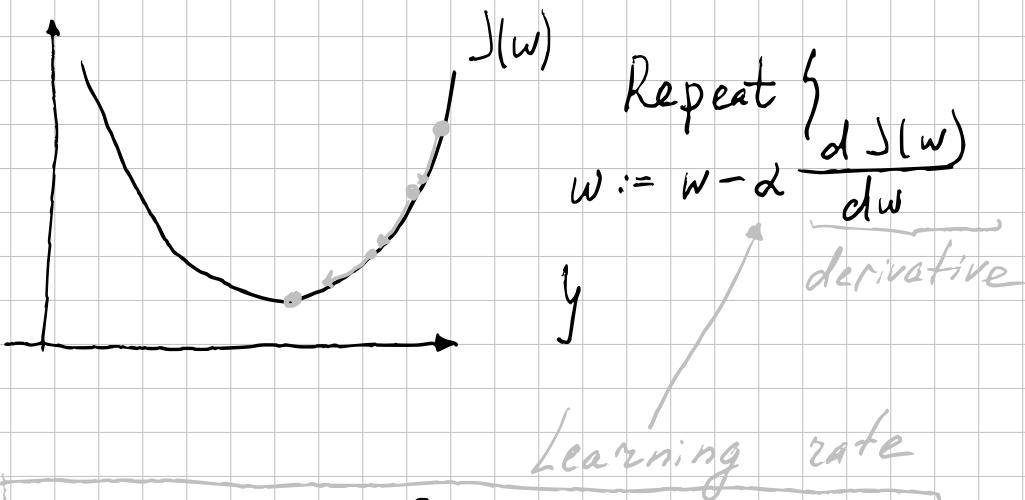
Gradient Descent

Want to find w, b that minimize $J(w, b)$

$J(w, b)$ is a multi-dimensional surface

Gradient descent starts at some initial point and takes a step towards the steepest direction

For 1-dimensional case

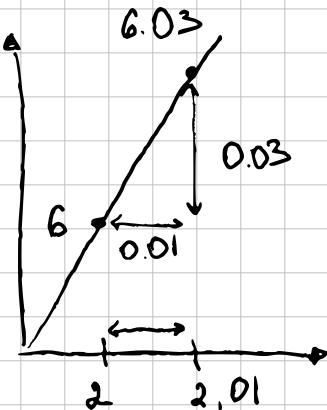


$$w := w - \lambda \frac{d J(w, b)}{d w}$$

$J(w, b)$

$$b := b - \lambda \frac{d J(w, b)}{d b}$$

Derivatives



$$f(a) = 3a$$

$$a = 2 \quad f(a) = 6$$

$$a = 2 + 0.01 \quad f(a) = 6.03$$

Slope (derivative) of $f(a)$
at $a = 2$ is 3

Slope is defined as $\frac{\text{height}}{\text{width}}$

$$a = 5 \quad f(a) = 15$$

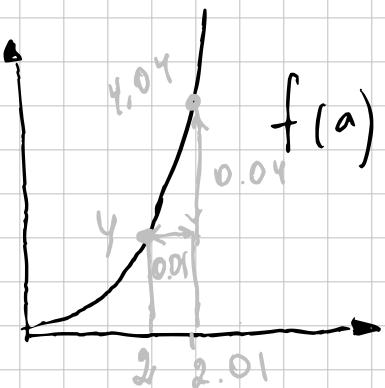
nudge $a = 5 + 0.01 \quad f(a) = 15.03$

Slope at $a = 5$ is 3

$$\frac{df(a)}{da} = 3 = \frac{d}{da} f(a)$$

Derivative is operating with
infinitely small nudges

More derivatives



$$f(a) = a^2$$

$$a = 2 \quad f(a) = 4$$

$$a = 2 + 0.01 \quad f(a) \approx$$

$$\approx 4.0401 \approx 4.04$$

$$\frac{d}{da} f(a) = 4 \text{ when } a = 2$$

$$a = 5 \quad f(a) = 25$$

$$a = 5 + 0.001 \quad f(a) \approx 25.010$$

$$\frac{d}{da} f(a) = 10 \text{ when } a = 5$$

$$\frac{d}{da} a^2 = 2a$$

Additional examples

$$f(a) = a^2 \quad \frac{d}{da} f(a) = 2a$$

$$f(a) = a^3 \quad \frac{d}{da} f(a) = 3a^2$$

$$f(a) = \log_e(a) = \ln(a) \quad \frac{d}{da} f(a) = \frac{1}{x}$$

Computation graph

Let's say we want to compute $J(a, b, c) = 3(a + bc)$

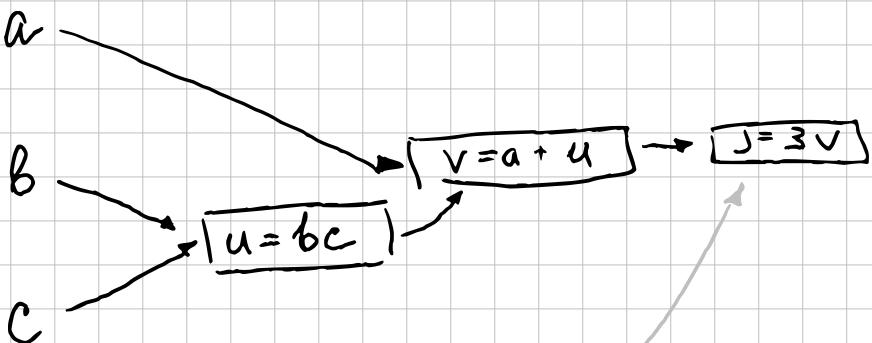
We can compute this function in a few steps:

$$1. \quad u = bc$$

$$2. \quad v = a + u$$

$$3. \quad J = 3v$$

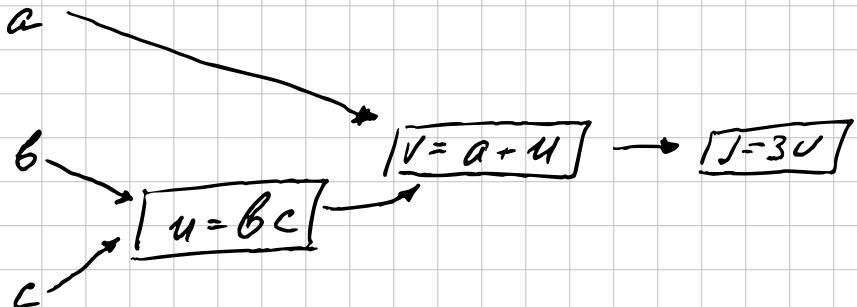
We can take this 3 steps and draw them in a computational graph



In case of logistic regression J is a cost function

Computing derivatives with a computation graph

The computation graph



Let's compute $\frac{dJ}{dv}$?

$$J = 3v$$

$$\frac{d}{dv} J = 3$$

What is $\frac{dJ}{da}$?

$$\frac{dJ}{da} = \frac{dJ}{dv} \cdot \frac{dv}{da} = 3 \cdot 1 = 3$$

Chain rule

Final output variable is a variable that we want to optimise

$$\frac{dJ}{du} = \frac{dJ}{dv} \frac{dv}{du} \frac{du}{db} = 3 \cdot 1 = 3$$

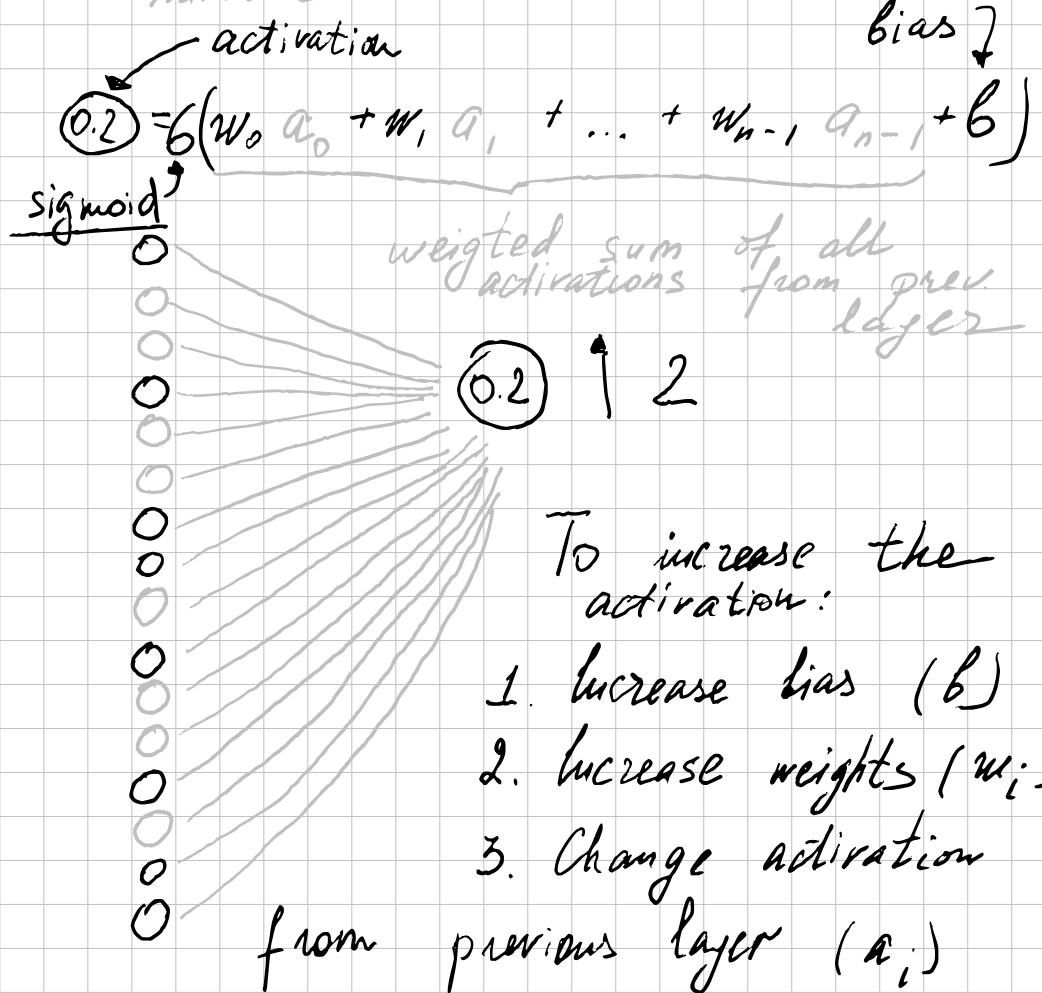
$$\frac{dJ}{db} = \frac{dJ}{dv} \frac{dv}{du} \frac{du}{db} = 3c$$

$$\frac{dJ}{dc} = \frac{dJ}{dv} \frac{dv}{du} \frac{du}{dc} = 3b$$

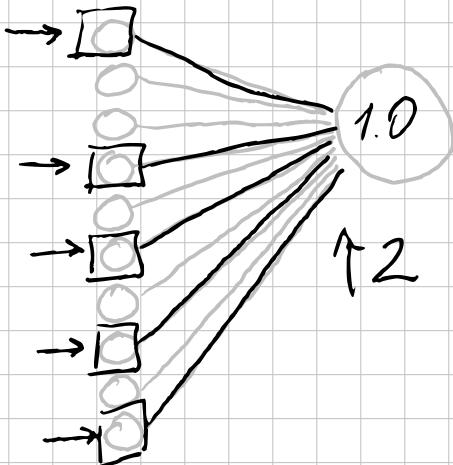
Back Propagation



We are trying to recognise hand-written numbers



Increasing neurons weight for neurons with bigger activation value will bigger affect the resulting activation



neurons in Π
have bigger a_i

- Increase w_i in proportion to a_i
- Change a_i in proportion to w_i

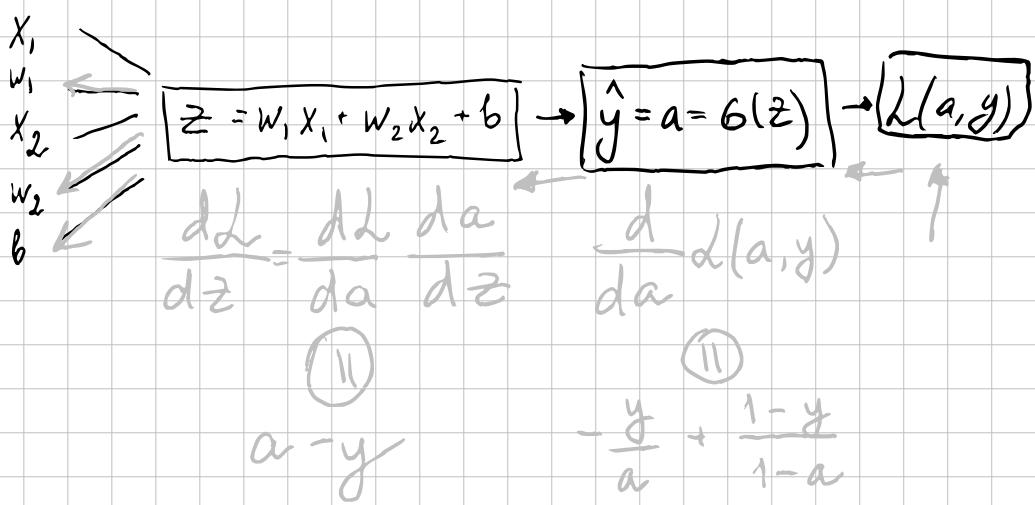
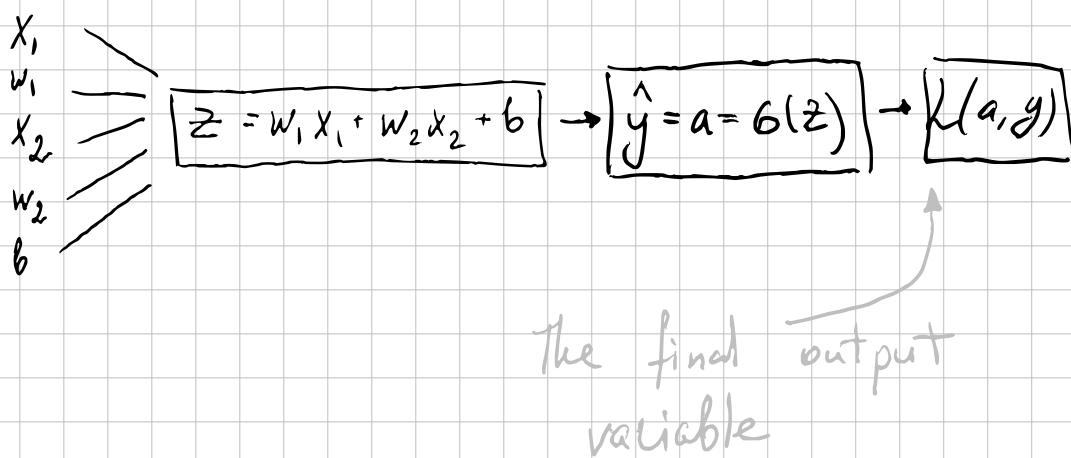
Logistic Regression Gradient

Desceat

$$z = w^T x + b$$

$$\hat{y} = a = g(z)$$

$$L(a, y) = -(y \log(a) + (1-y) \log(1-a))$$



$$\frac{dL}{dw_1} = \frac{dd}{dz} x_1$$

$$\frac{dh}{db} = \frac{dd}{dz}$$

$$\frac{dd}{dw_2} = \frac{dd}{dz} x_2$$

$$w_i = w_{i_1} - \alpha dw_{i_1}$$

$$\omega_2 = \omega_2 - 2 d\omega_2$$

$$b = b - \alpha \nabla_b$$

Gradient descent for logistic regression

One step of Gradient Descent:

$$J=0; dw_1=0; dw_2=0; db=0$$

For $i=1$ to m :

$$z^{(i)} = w^T x^{(i)} + b$$

$$a^{(i)} = \sigma(z^{(i)})$$

$$J += -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log(1-a^{(i)})]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1 += x_1^{(i)} dz^{(i)} \quad \uparrow n=2$$

$$dw_2 += x_2^{(i)} dz^{(i)}$$

$$db += dz^{(i)}$$

$$J / m$$

$$dw_1 / m; dw_2 / m; db / m.$$

$$w_1 := w_1 - \alpha dw_1$$

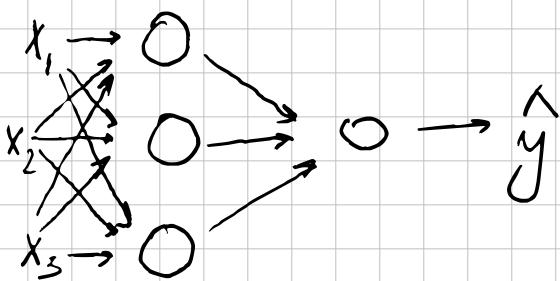
$$w_2 := w_2 - \alpha dw_2$$

$$b := b - \alpha db$$

Problem with the approach:

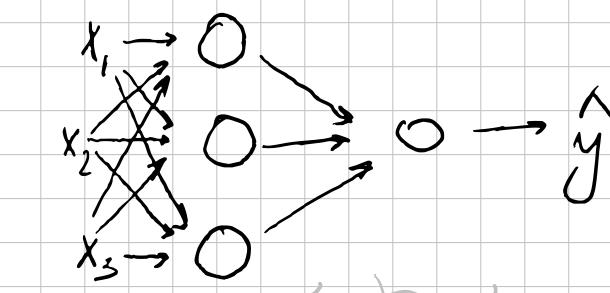
- multiple for loops, i.e. complexity is $O(n \times m)$

Neural Network



Notation

$$w^{[1]} \quad b^{[1]} \quad \left\{ \begin{array}{l} [1] \\ [1] \end{array} \right. - \text{first layer}$$



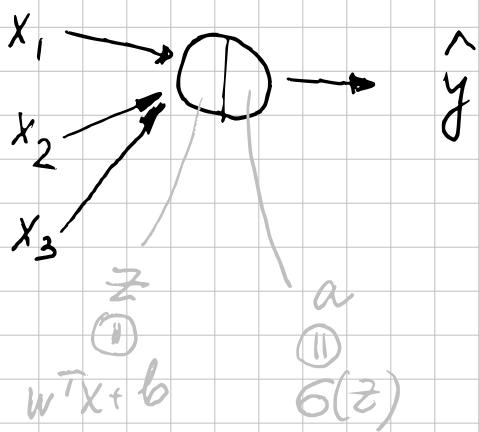
Input layer hidden layer Output layer

2 layers
neural network

(we don't
count input +
layer)

Alternative notation for
input layer is $a^{[0]}$

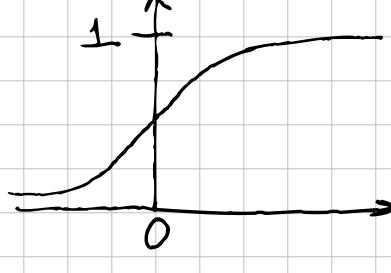
Neuron



$a_i^{[l]} \leftarrow$ layer
 $a_i \leftarrow$ node in layer

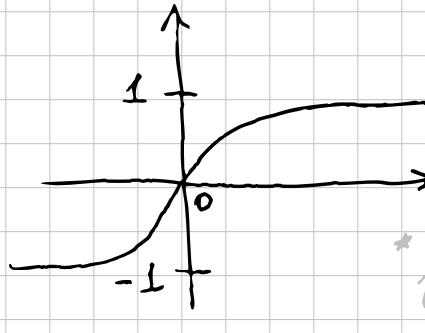
Activation functions

Examples:



Sigmoid

$$a(z) = \frac{1}{1+e^{-z}}$$



Hyperbolic tangent

$$a(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

* Usually works better than Sigmoid as shifts mean to 0.

One possible exception for it is the last layer if output is binary classification ∈ {0, 1}

Problem of Sigmoid and tanh is gradient vanishing: values that are really big or really small usually are close to each other after transformation.

ReLU

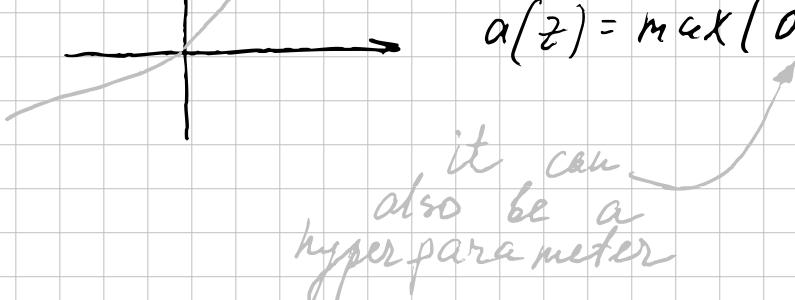
$$a(z) = \max(0, z)$$



Sigmoid is very natural for the output layer for binary classif. {0, 1}. ReLU is increasingly default choice for other units.

Leaky ReLU

$$a(z) = \max(0, 0.01z, z)$$



it can also be a hyperparameter

Why do you need Non-linear Activation Functions?

If we replace non-linear transformations with linear functions, we would be able to substitute all layers with the only layer.

No matter how many layers NN has, if it only has linear activations, it will only output linear function.

Let's

$$a^{[1]} = z$$

$$\cdot a^{[1]} = z^{[1]} = W^{[1]} X + b^{[1]}$$

$$\cdot a^{[2]} = z^{[2]} = W^{[2]} \cancel{X} + b^{[2]} \quad \textcircled{=} \quad \textcircled{=}$$

$$\textcircled{=} W^{[2]} (W^{[1]} X + b^{[1]}) + b^{[2]} \quad \textcircled{=} \quad \textcircled{=}$$

$$\underbrace{W^{[2]} W^{[1]} X}_{W'} + \underbrace{W^{[2]} b^{[1]} + b^{[2]}}_{b'} \quad \textcircled{=} \quad \textcircled{=}$$

$$\underline{\underline{W' X + b'}} \quad \textcircled{=} \quad \textcircled{=}$$

Derivatives of activation functions

• Sigmoid

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

$$\frac{d}{dz} g(z) = \frac{-1}{(1+e^{-z})^2} \cdot e^{-z} \cdot (-1) \quad \textcircled{=}$$

$$\textcircled{=} \frac{e^{-z}}{(1+e^{-z})^2} = \frac{e^{-z}}{1+e^{-z}} \cdot \frac{1}{1+e^{-z}} = \frac{e^{-z} + 1 - 1}{1+e^{-z}} \quad \textcircled{+}$$

$$\textcircled{\bullet} \frac{1}{1+e^{-z}} = \left(1 - \frac{1}{1+e^{-z}}\right) \cdot \frac{1}{1+e^{-z}} = (1 - g(z))g'(z)$$

• Hyperbolic tangent

$$g(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\frac{d}{dz} g(z) = \left(\frac{e^z - e^{-z}}{e^z + e^{-z}} \right)' \quad \textcircled{=}$$

$$\textcircled{=} \frac{(e^z - e^{-z})'(e^z + e^{-z}) - (e^z - e^{-z})(e^z + e^{-z})'}{(e^z + e^{-z})^2}$$

$$\textcircled{=} \frac{(e^z + e^{-z})^2 - (e^z - e^{-z})^2}{(e^z + e^{-z})^2} = 1 - \left(\frac{e^z - e^{-z}}{e^z + e^{-z}} \right)^2$$

$$\textcircled{=} 1 - (\tanh(z))^2$$

• ReLU

$$g(z) = \max(0, z)$$

$$\frac{d}{dz} g(z) = \begin{cases} 0, & z \leq 0 \\ 1, & z > 0 \end{cases}$$

• Leaky ReLU

$$g(z) = \max(0.01z, z)$$

$$g'(z) = \begin{cases} 0.01, & z \leq 0 \\ 1, & z > 0 \end{cases}$$

Gradient Descent for neural networks

NN with a single hidden layer

Parameters

$$w^{(1)} \quad (n^{(1)}, n^{(0)})$$

$$b^{(1)} \quad (n^{(1)}, 1)$$

$$w^{(2)} \quad (n^{(2)}, n^{(1)})$$

$$b^{(2)} \quad (n^{(2)}, 1)$$

Cost function

$$J(w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i)$$

Gradient descent :

repeat {

forward prop.

backward prop.

compute predicts ($\hat{y}^{(i)}$, $i = 1 \dots m$)

$$\frac{dw^{(1)}}{dw^{(1)}} = \frac{dJ}{dw^{(1)}},$$

$$\frac{db^{(1)}}{db^{(1)}} = \frac{dJ}{db^{(1)}}$$

:

$$w^{(1)} = w^{(1)} - \alpha \frac{dw^{(1)}}{dw^{(1)}}$$

:

y

logistic regression

$x \rightarrow$

$$w \rightarrow \boxed{z = w^T x + b} \rightarrow \boxed{a = g(z)} \rightarrow \boxed{d(a, y)}$$

b

NN

$$\begin{matrix} x_1 & 0 \\ x_2 & \rightarrow 0 \rightarrow 0 \rightarrow \\ x_3 & 0 \end{matrix}$$

$$x \rightarrow \boxed{z^{(1)} = w^{(1)} x + b^{(1)}} \rightarrow \boxed{a = g(z)} \rightarrow \boxed{z^{(2)} = w^{(2)} x + b^{(2)}}$$

$w^{(2)}$

$b^{(2)}$

$\ominus \boxed{a = g(z^{(2)})} \rightarrow d(a^{(2)}, y)$

$$\frac{dd}{dz^{(2)}} = a^{(2)} - y$$

$$\frac{dL}{dw^{(2)}} = dz^{(2)} a^{(1)T} \quad \frac{dd}{db^{(2)}} = dz^{(2)}$$

$$dz^{(1)} = w^{(2)T} dz^{(2)} * g^{(1)'}(z^{(1)})'$$

Random initialization

Let's initialize everything with 0. In that case all derivatives are identical. All units in the same layer compute the same function. (symmetry breaking prob.)

The solution to this is initialize every parameter at random.

Deep L-layer Neural Network

Number of hidden layers can be considered as another hyperparameter

L (number of layers in the network)

$n^{[l]}$ (number of nodes / units in layer l)

$a^{[l]}$ (activations in layer l)

$$a^{[l]} = g^{[l]}(z^{[l]})$$

$w^{[l]}$ (weights for $z^{[l]}$)

Forward propagation in a deep network

$$x : z^{(1)} = w^{(1)} x + b^{(1)}$$

$$a^{(1)} = g^{(1)}(z^{(1)})$$

$$z^{(2)} = w^{(2)} a^{(1)} + b^{(2)}$$

$$a^{(2)} = g^{(2)}(z^{(2)})$$

... until you get to the output layer

$$z^{(n)} = w^{(n)} a^{(n-1)} + b^{(n)}$$

$$a^{(n)} = g^{(n)}(z^{(n)})$$

General rule is

$$z^{(l)} = w^{(l)} a^{(l-1)} + b^{(l)}$$

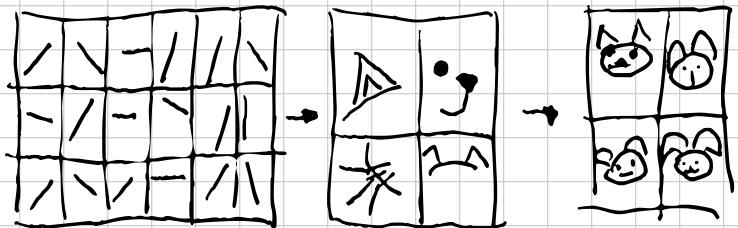
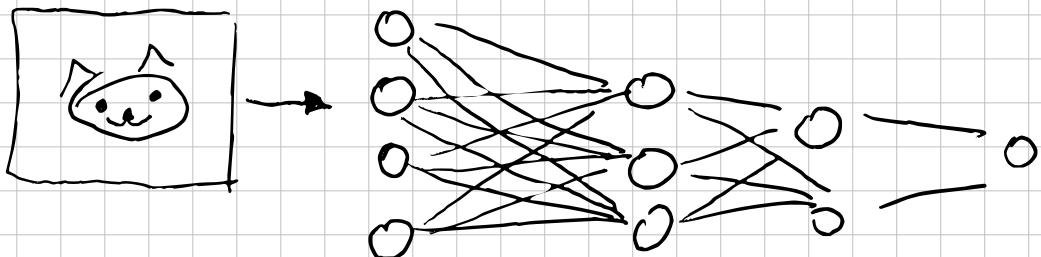
$$a^{(l)} = g^{(l)}(z^{(l)})$$

$$w^{(l)} : \left(h^{(l)}, h^{(l-1)} \right)$$

$$a^{(l)} (X \text{ is } A^{(0)}) : \left(h^{(l)}, 1 \right)$$

$$b^{(l)} : \left(h^{(l)}, 1 \right)$$

Why Deep Representations?



Allows to compose more
simple features into more
complicated ones

Building blocks of Deep Neural networks

layer l : $W^{(l)}$ $b^{(l)}$

forward: input $a^{(l-1)}$ output $a^{(l)}$

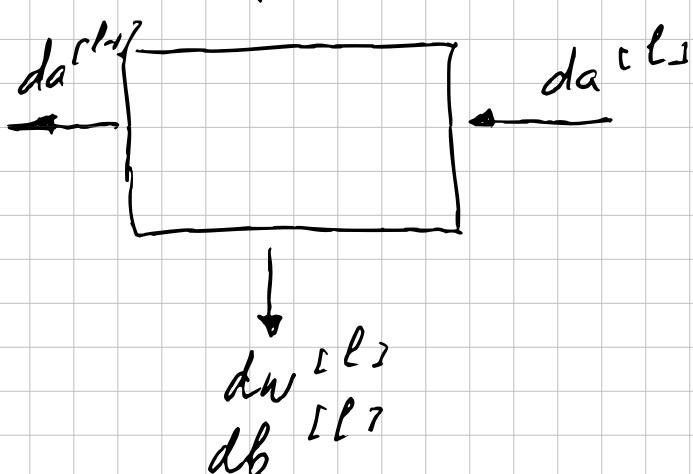
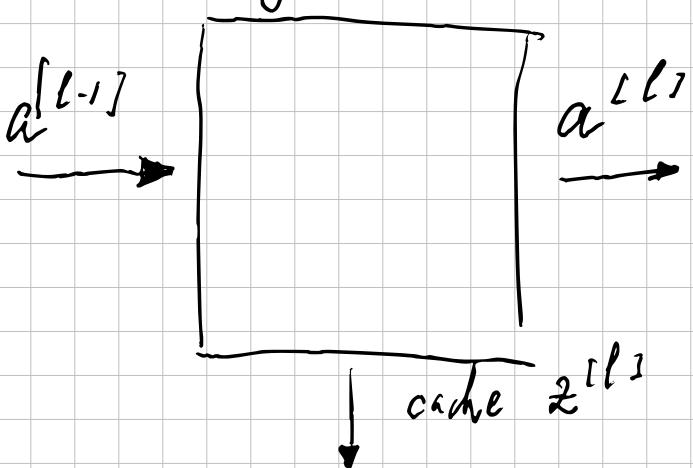
$$z^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)}$$

cache $z^{(l)}$

$$a^{(l)} = g^{(l)}(z^{(l)})$$

backward: input $da^{(l)}$,
output $da^{(l-1)}$

layer l



Forward propagation for layer ℓ

Input $a^{[\ell-1]}$

Output $a^{[\ell]}$, cache $(z^{[\ell]})$

$$z^{[\ell]} = W^{[\ell]} a^{[\ell-1]} + b^{[\ell]}$$

$$a^{[\ell]} = g^{[\ell]}(z^{[\ell]})$$

Backward propagation for layer ℓ

Input $da^{[\ell]}$

Output $da^{[\ell-1]}, dW^{[\ell]}, db^{[\ell]}$

$$dz^{[\ell]} = da^{[\ell]} * g^{[\ell]}'(z^{[\ell]})$$

prime,
i.e. derivative

$$dW^{[\ell]} = dz^{[\ell]} \cdot a^{[\ell-1]T}$$

$$db^{[\ell]} = dz^{[\ell]}$$

$$da^{[\ell-1]} = W^{[\ell]T} \cdot dz^{[\ell]}$$

Parameters vs. hyperparameters

Parameters: $W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}, \dots$

Hyperparameters: learning rate
d, # iterations, # hidden layers, # hidden units, choice of activation function

Forward propagation

$$z_{ji}^{[l]} = \sum_k w_{jk}^{[l]} a_{ki}^{[l-1]} + b_j^{[l]}$$

$$a_{ji}^{[l]} = g_j^{[l]}(z_{1i}^{[l]}, \dots, z_{ji}^{[l]}, \dots, z_{ni}^{[l]})$$

$$a_{k-2,i}^{[l-1]} \xrightarrow{w_{jk-2}^{[l]}} a_{ji}^{[l]}$$

$$a_{k-1,i}^{[l-1]} \xrightarrow{w_{jk}^{[l]}} a_{j+1,i}^{[l]}$$

$$a_{ki}^{[l-1]} \xrightarrow{w_{jk}^{[l]}} a_{j+2,i}^{[l]}$$

Vectorized version (for one example i)

$$\begin{bmatrix} z_{1i}^{[l]} \\ z_{ji}^{[l]} \\ z_{ni}^{[l]} \end{bmatrix} = \begin{bmatrix} w_{11}^{[l]} & \dots & w_{1k}^{[l]} & \dots & w_{1n}^{[l]} \\ w_{j1}^{[l]} & \dots & w_{jk}^{[l]} & \dots & w_{jn}^{[l]} \\ w_{n1}^{[l]} & \dots & w_{nk}^{[l]} & \dots & w_{nn}^{[l]} \end{bmatrix} \begin{bmatrix} a_{1i}^{[l-1]} \\ \vdots \\ a_{ki}^{[l-1]} \\ \vdots \\ a_{ni}^{[l-1]} \end{bmatrix} + \begin{bmatrix} b_1^{[l]} \\ \vdots \\ b_j^{[l]} \\ \vdots \\ b_n^{[l]} \end{bmatrix}$$

$$\begin{bmatrix} a_{1i}^{[l]} \\ \vdots \\ a_{ji}^{[l]} \\ \vdots \\ a_{ni}^{[l]} \end{bmatrix} = \begin{bmatrix} g_1^{[l]}(z_{1i}^{[l]}, \dots, z_{ji}^{[l]}, \dots, z_{ni}^{[l]}) \\ \vdots \\ g_n^{[l]}(\dots) \end{bmatrix}$$

Or

$$z_{:,i}^{[l]} = W^{[l]} a_{:,i}^{[l-1]} + b^{[l]}$$

$$a_{:,i}^{[l]} = g^{[l]}(z_{:,i}^{[l]})$$

Vectorizing over all examples

$$z^{[l]} = [z_{:,1}^{[l]} \dots z_{:,i}^{[l]} \dots z_{:,n}^{[l]}] \quad (=)$$

$$= W^{[l]} \begin{bmatrix} a_{:,1}^{[l-1]} & \dots & a_{:,i}^{[l-1]} & \dots & a_{:,n}^{[l-1]} \end{bmatrix} +$$

$$+ [b^{[l]} \dots b^{[l]} \dots b^{[l]}] \quad (=)$$

$$= W^{[l]} A^{[l-1]} + \text{broadcast}(b^{[l]})$$

Cost function

$$J = f(\hat{y}, y) = f(A^{[l]}, y)$$

Backward propagation

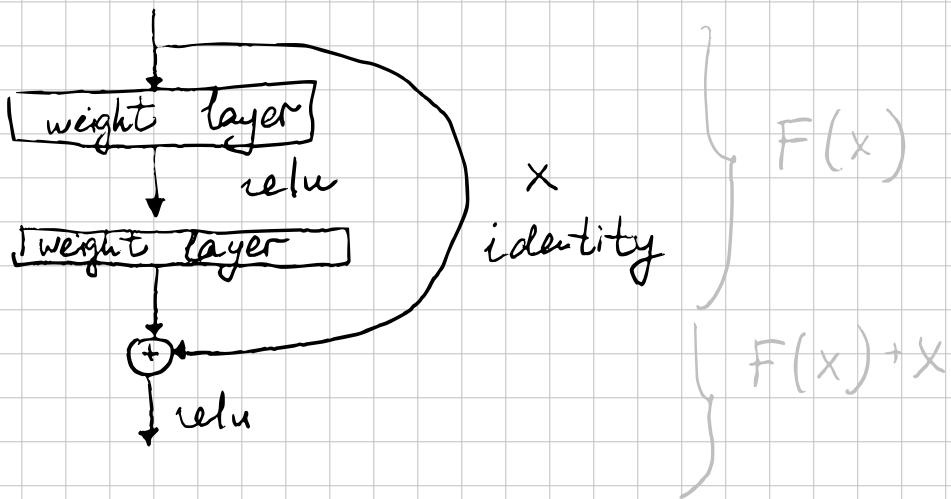
$$\frac{\partial J}{\partial w_{jk}^{[l]}} = \sum_i \frac{\partial J}{\partial z_{ji}^{[l]}} \frac{\partial z_{ji}^{[l]}}{\partial w_{jk}^{[l]}} = \sum_i \frac{\partial J}{\partial z_{ji}^{[l]}} \cdot a_{ki}^{[l-1]}$$

$$\frac{\partial J}{\partial b_j^{[l]}} = \sum_i \frac{\partial J}{\partial z_{ji}^{[l]}} \frac{\partial z_{ji}^{[l]}}{\partial b_j^{[l]}} = \sum_i \frac{\partial J}{\partial z_{ji}^{[l]}}$$

$$\frac{\partial J}{\partial a_{pi}^{[l]}} = \sum_i \frac{\partial J}{\partial z_{ji}^{[l]}} \frac{\partial z_{ji}^{[l]}}{\partial a_{pi}^{[l]}} = \sum_i \frac{\partial J}{\partial z_{ji}^{[l]}}$$

$$\frac{\partial J}{\partial a_{ki}^{[l-1]}} = \sum_j \frac{\partial J}{\partial z_{ji}^{[l]}} \frac{\partial z_{ji}^{[l]}}{\partial a_{ki}^{[l-1]}} = \sum_j \frac{\partial J}{\partial z_{ji}^{[l]}} \cdot w_{jk}^{[l]}$$

Deep Residual Learning



$H(x)$ is the true function we want to learn

Lets pretend we want to learn $F(x) = H(x) - x$ instead

The original function is then

$$\underline{F(x) + x}$$

Depthwise Separable Convolution

A standard convolution is an operator of size

$$D_k * D_k * C_{in}$$

where D_k - convolutional kernel size, C_{in} - input number of channels

Computational complexity of the convolutional layer is $D_k * D_k * C_{in} \star$

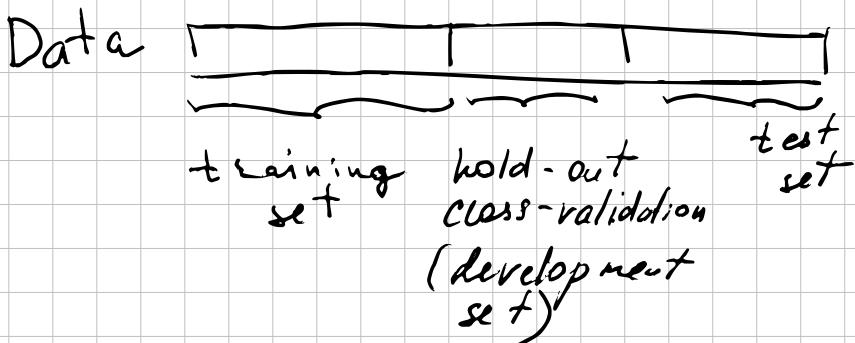
$$\star D_f * D_f * C_{out}$$

Practical aspects of deep learning

You may want to decide on how many of each you want:

- # layers
- # hidden units
- learning rates
- activation functions

...



Usually people do: 60/20/20 (%)
split

Dev. set helps to evaluate multiple algorithm choices

Mismatch train/test distribution

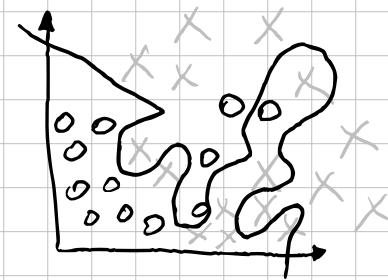
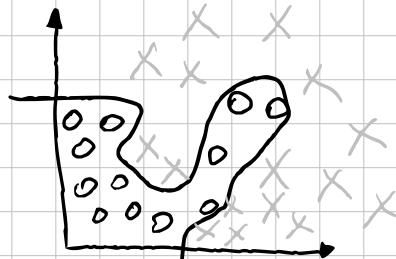
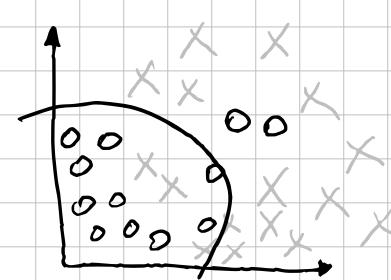
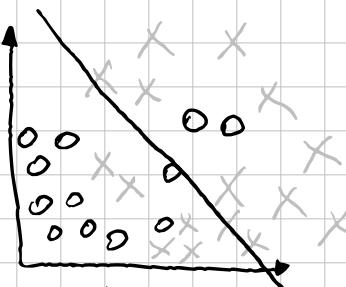
Training set: cat pictures from webpages

Dev/test sets: cat pictures from users using your app

Rule of thumb: make sure dev and test come from the same distribution.

Not having a test set might be okay. (Only dev. set)

Bias / Variance



Cat classification



$y=1$



$y=0$

The key metrics for high-dimens. data are:

train set error: 1%

dev set error: 11%

may signal/lse of
overfitting as
good on train and
bad on dev
(high variance)

15%

16%

not doing
well on both

smts? so underfitting
(high bias)

15%

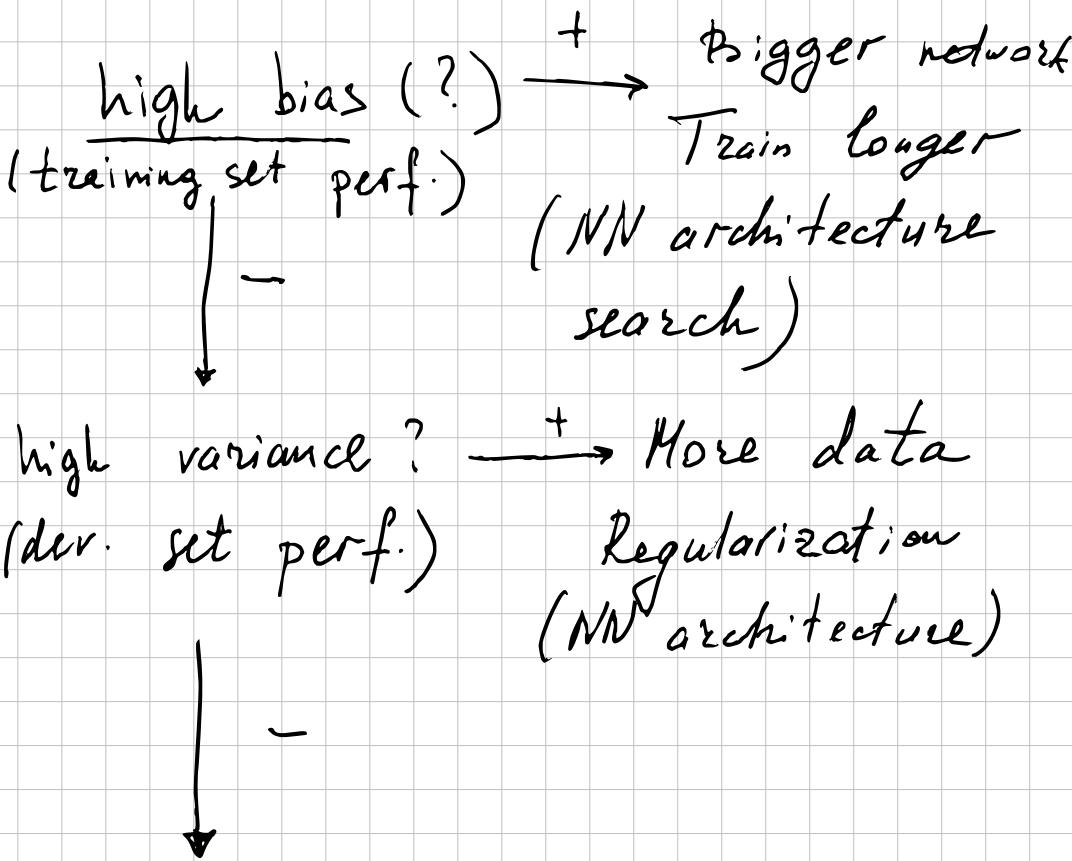
30%

high bias

&

high variance

Basic recipe for Machine Learning



"bias variance tradeoff"

Regularization

$$\min_{w,b} J(w,b), \quad w \in \mathbb{R}^{n_x}, b \in \mathbb{R}$$

For logistic regression:

$$J(w,b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \underbrace{\frac{\lambda}{2m} \|w\|^2}_{\text{regularization factor}}$$

$$\|w\|^2 = \sum_{j=1}^{n_x} w_j^2 = w^\top w$$

$$\frac{\lambda}{2m} \sum_{i=1}^{n_x} |w_i| = \frac{\lambda}{2m} \|w\|_1,$$

L2 regular.

L1 regular.

λ -regularization parameter

For neural network:

$$J(w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|^2$$

$$\|w^{(l)}\|^2 = \sum_{i=1}^{n_l} \sum_{j=1}^{n_{l-1}} (w_{ij}^{(l)})^2$$

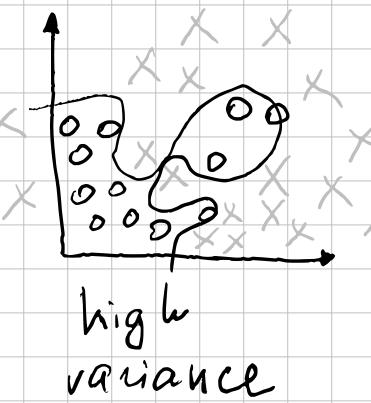
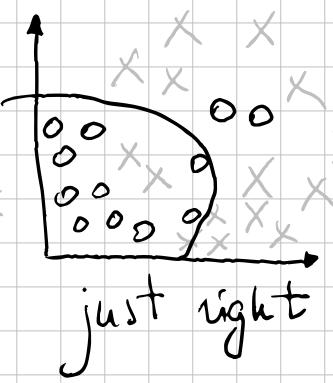
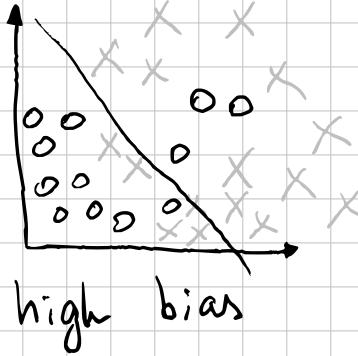
Frobenius norm

Gradient descent:

$$dw^{(l)} = -(\text{from backprop}) + \frac{1}{m} W^{(l)}$$

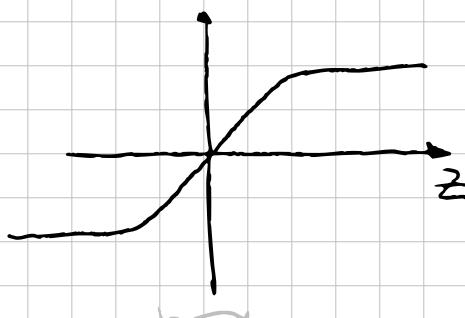
$$w^{(l)} := w^{(l)} - \alpha dw^{(l)}$$

Why regularization reduces overfitting?



$$J(w^{[l]}, b^{[l]}) = \frac{1}{m} \sum_{i=1}^m l(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$$

Another intuition



$$g(z) = \tanh(z)$$

z taking values from this interval

If λ is large ($\lambda \uparrow$)

my parameters are small ($w^{[l]} \downarrow$)
because they are penalized
by the cost function

$$\underbrace{z^{[l]}}_{\rightarrow} = w^{[l]} a^{[l-1]} + b^{[l]}$$

Every layer will be roughly \approx
linear

If w is small $z^{[l]}$

takes values from small limited range

Dropout regularization

You eliminate some nodes (randomly) from the network and train this reduced network.

- Inverted dropout
- Illustrate with $L=3$

$$\text{keep_prob} = 0.8$$

$$d_3 = \text{np.random.rand}(a_3.\text{shape}[0], a_3.\text{shape}[1]) < \text{keep_prob}$$

$$a_3 = \text{np.multiply}(a_3, d_3)$$

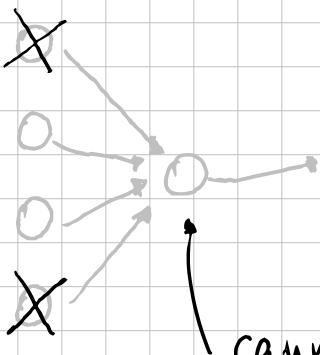

zeroing-out
elements

Scaling a_3 distribution back
 $a_3 / \text{keep_prob}$

At test time we don't add
any drop outs.

Understanding dropout

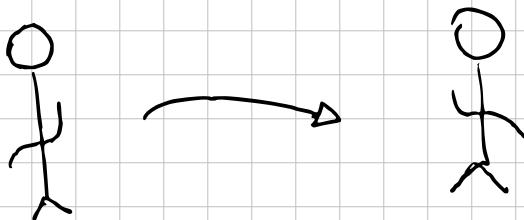
Intuition: Can't rely on any one feature, so have to spread out weights.



cannot rely on any prev.
nodes feature as the
feature may go away

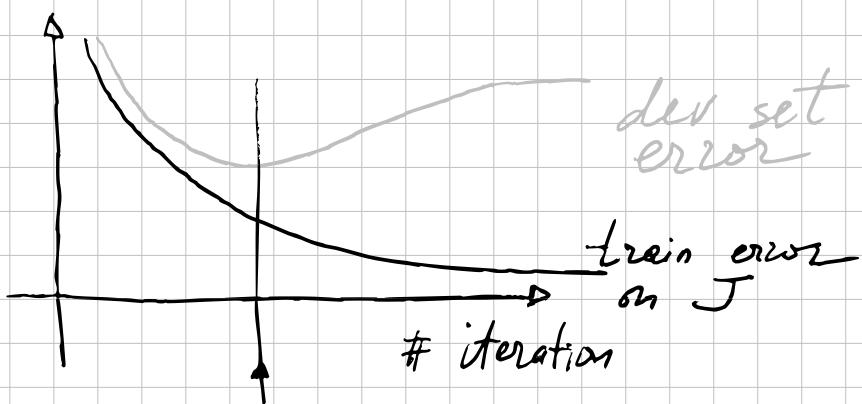
Other regularisation methods

Data augmentation



- Image flip
- Randomly crop the image
- Random distortions / transformations

Early stopping



stop here
as it was
the best parameter

Normalizing training set

$$X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

steps:

1. Subtract mean

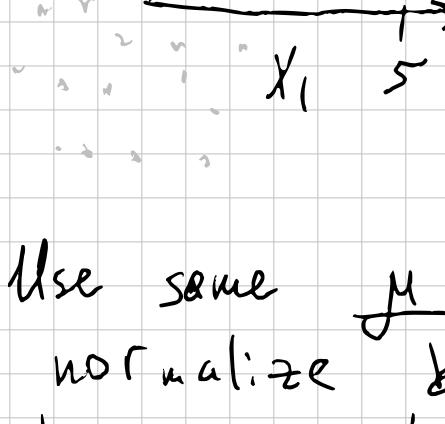
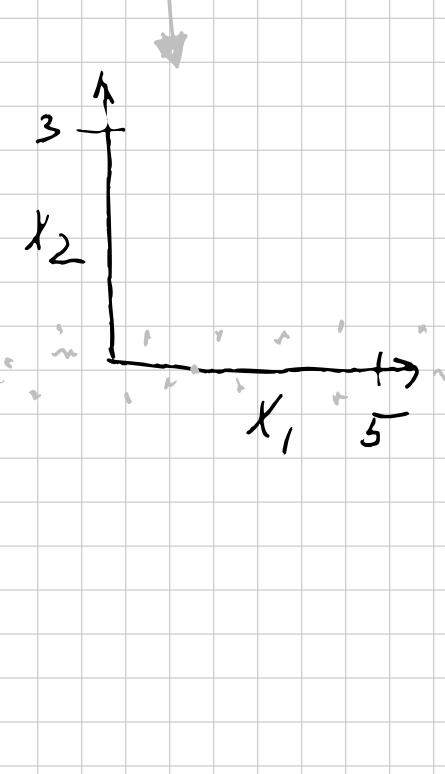
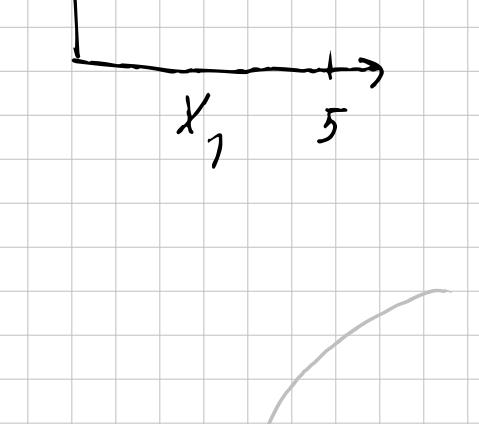
$$M = \frac{1}{m} \sum^m x^{(i)}$$

$$x = x - \mu$$

Normalize

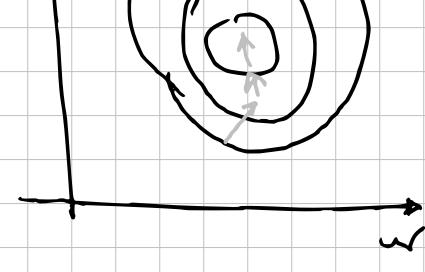
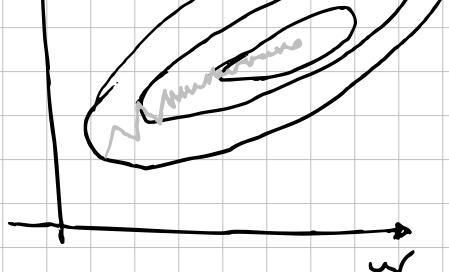
$$b = m \sum_{i=1}^n \binom{x}{i}$$

$$X = X/G$$



<u>training sets</u>	Non-normalized loss function	Normalized loss function
b		

A graph on a grid showing a function $f(x)$. The curve starts at a local minimum at $x = b$, rises to a local maximum at $x = a$, and then decreases. A vertical arrow points upwards from the label b to the minimum point on the curve.

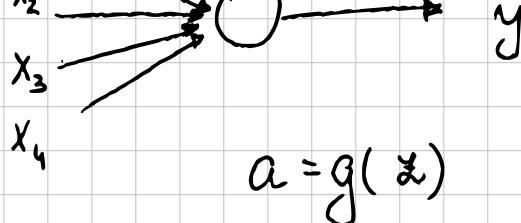


Vanishing / Exploding Gradients

- $w^{[l]} > 1$ the activations can explode in a really deep network
- $w^{[l]} < 1$ the activations will decrease exponentially in deep networks

The same happens with derivatives.

Weights Initialization for Deep Networks



$$a = g(z)$$

$$z = w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_n x_n + b$$

$b=0$ in this example

Large $n \rightarrow$ Smaller w_i

$$\text{Var}(w) = \frac{1}{n}$$

$$w^{[l]} = \text{np.random.rand(shape)} \odot \text{np.sqrt}\left(\frac{1}{n^{[l-1]}}\right)$$

$$\text{ReLU} \quad g^{[l]}(z) = \text{ReLU}(z)$$

Other variance:

$$\text{ReLU} = \frac{2}{n^{[l-1]}}$$

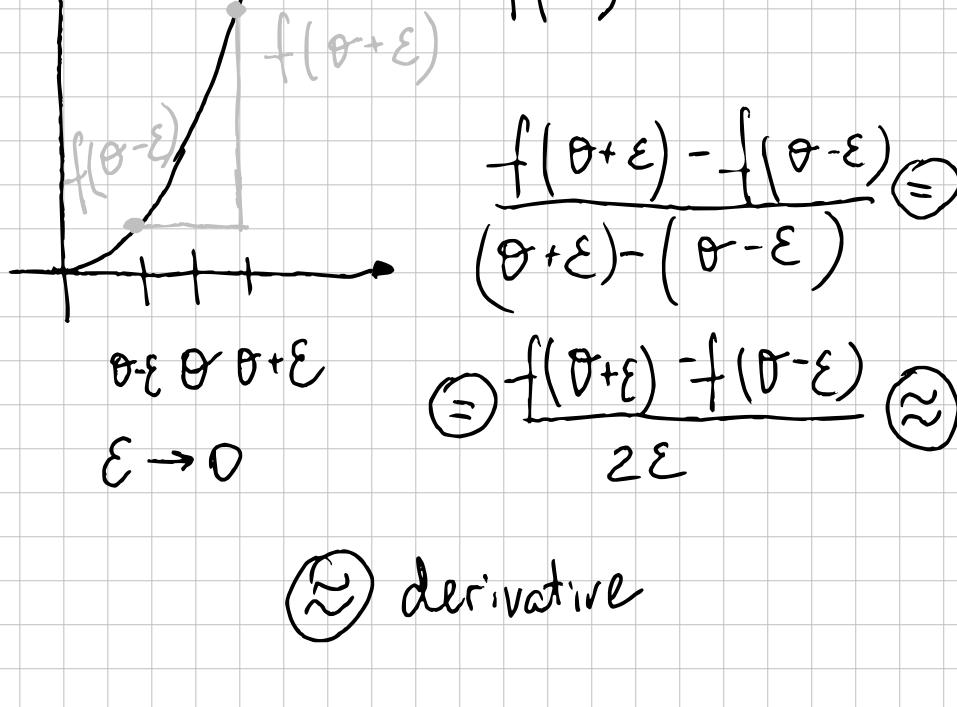
$$\tanh = \sqrt{\frac{1}{n^{[l-1]}}} \quad \begin{cases} \text{Xavier} \\ \text{initialisation} \end{cases}$$

$$\sqrt{\frac{2}{n^{[l-1]} + n^{[l]}}} - \text{another popular}$$

Numeric Approximation

of Gradients

Checking your derivative computation



Gradient check for a neural network

Take $w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}$ and reshape into a big vector

θ .

Take $d\theta^{[1]}, d\theta^{[2]}, \dots, d\theta^{[L]}$ and reshape into a big vector $d\theta$.

Q: Is $d\theta$ a gradient of $J(\theta)$?

for each i :

$$d\theta_{\text{approx}}[i] = \frac{J(\theta_1, \theta_2, \dots, \overset{i}{\theta_i + \epsilon}, \dots) - J(\theta_1, \theta_2, \dots, \overset{i}{\theta_i - \epsilon}, \dots)}{2\epsilon} \approx d\theta[i]$$

$$d\theta_{\text{approx}} \approx d\theta \quad \leftarrow \text{The actual check}$$

$$\frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2} < \text{some small epsilon, like } 10^{-7}$$

error :	10^{-7}	- great
	10^{-5}	so-so
	10^{-3}	- worry

Gradient Checking

tricks & hints

- Don't use in training - only to debug;
- If algorithm fails grad check, look at components to try to identify bugs;

For example db may be far off, but dW may be just fine

- Remember regularization;
- Does not work with dropout (implement check without dropout)
- Run at random initialisation; perhaps again after some training.

Mini-batch Gradient

Descent

Batch vs. mini-batch gradient descent

Vectorization allows you to efficiently compute on m examples.

$$X = [x^{(1)} \ x^{(2)} \dots x^{(m)}]$$

$$y = [y^{(1)} \ y^{(2)} \dots y^{(m)}]$$

X dimension is (n_x, m)

y dimension is $(1, m)$

What if $m = 5 \cdot 10^7$?

mini-batches are small (tiny, baby) split ups of the training set

$x^{(1)}$ - is the first mini-batch split

For example, if mini-batch size is 10^3

$$X = [x^{(1)} \dots x^{(1000)} / x^{(1001)} \dots x^{(2000)} / \dots x^{(m)}]$$

y set is split up in the same manner.

Notation recap.

$x^{(i)}$ - i-th example

$z^{(l)}$ - l-th layer

$x^{(t)}$ - t-th mini-batch

for $t = 1, \dots, 5000$:

vectorized implementation 1. forward prop on $x^{(t)}$

$$z^{(L)} = w^{(L)} x^{(t)} + b^{(L)}$$

$$A^{(L)} = g^{(L)}(z^{(L)})$$

2. compute cost function

$$J = \frac{1}{1000} \sum_{i=1}^{1000} L(\hat{y}^{(i)}, y^{(i)})$$

3. backprop. to (as if $m = 1000$)

compute gradients

4. Update

$$w^{(l)} := w^{(l)} - \frac{\partial J}{\partial w^{(l)}}$$

$$b^{(l)} := b^{(l)} - \frac{\partial J}{\partial b^{(l)}}$$

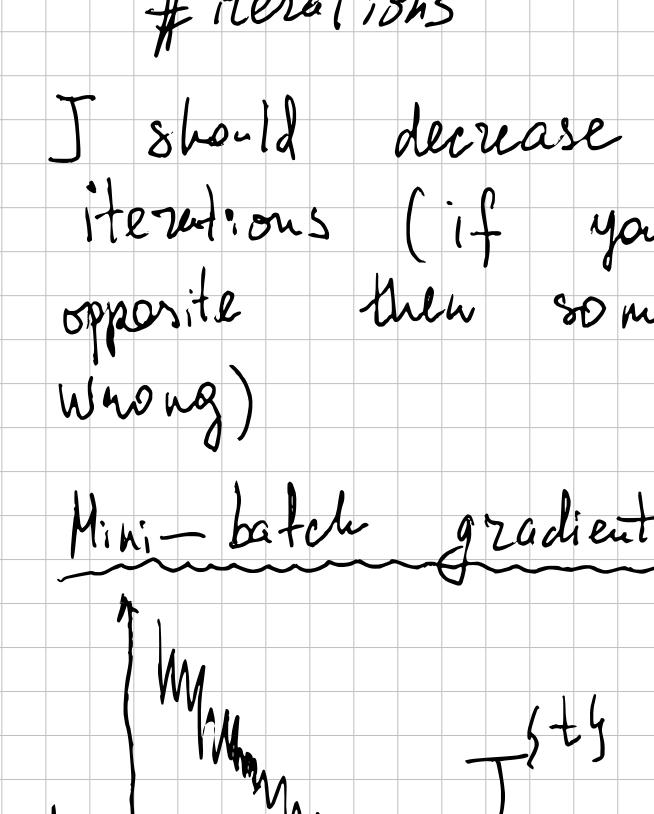
It is called "1 epoch", i.e.

one pass through the training set.

Understanding Mini-batch Gradient Descent

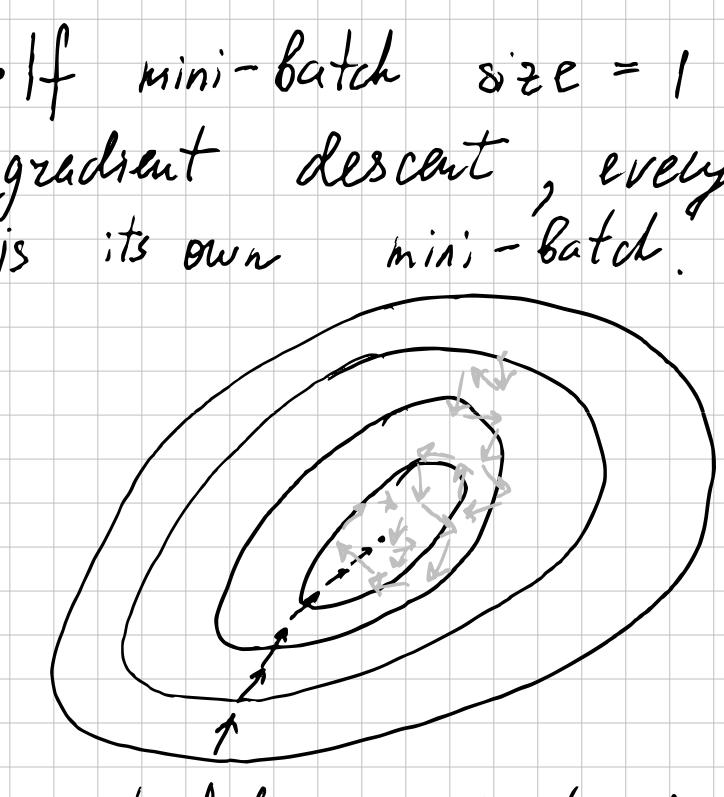
Batch gradient descent

You go through the entire training set



J should decrease over iterations (if you see opposite then something is wrong)

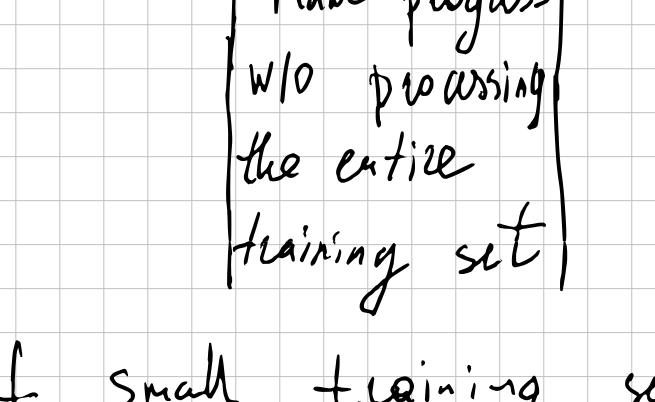
Mini-batch gradient descent



It's okay if it does not go down every iteration.

Choosing your mini-batch size

- If mini-batch size = m : batch gradient descent , in this extreme case (x^{s14}, y^{s14}) = (X, Y)
- If mini-batch size = 1 : Stochastic gradient descent , every example is its own mini-batch.



Batch gradient descent

Stochastic gradient descent

In practice mini-batch gradient descent is in-between

- | Batch grad. desc | Mini-batch grad. desc | Stochastic grad. desc. |
|----------------------|--|-----------------------------------|
| • Too long per iter. | • Fast learning w/o processing the entire training set | • lose speedup from vectorization |

more common values

- If small training set : use batch gradient descent ($m \leq 2 \cdot 10^3$)

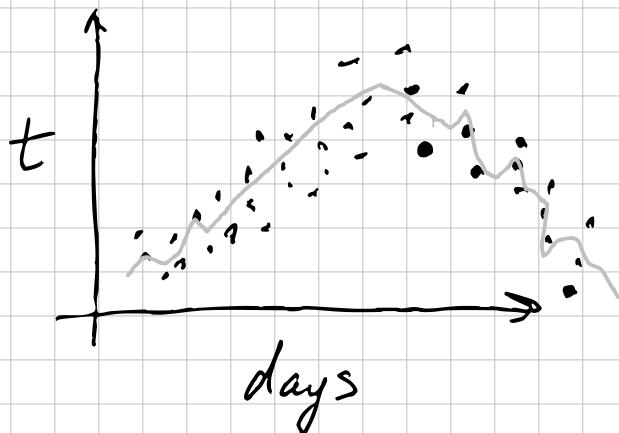
- Otherwise typical mini-batch sizes: $64, 128, 256, 512, 1024$

more common values

- ! Make sure mini-batch fits in CPU / GPU memory

Exponentially weighted moving averages

Temperature in London



$$\theta_1 = 4^\circ\text{C}$$

$$\theta_2 = 9^\circ\text{C}$$

:

$$\theta_{180} = 15^\circ\text{C}$$

$$V_0 = 0$$

$$V_1 = 0.9V_0 + 0.1\theta_1$$

$$V_2 = 0.9V_1 + 0.1\theta_2$$

:

$$V_t = 0.9V_{t-1} + 0.1\theta_t$$

Exponentially
weighted
average
(~~the~~ or the plot)

$$V_t = \beta V_{t-1} + (1-\beta)\theta_t$$

V_t is an approximate average over $\left(\frac{1}{1-\beta}\right)$ days' temperature

For example,

$\beta = 0.9$ is ≈ 10 days' t°

$\beta = 0.98$ is ≈ 50 days' t°

$\beta = 0.5$ is ≈ 2 days' t°

Understanding exponentially weighted average

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

$$\begin{cases} V_{100} = 0.9 V_{99} + 0.1 \theta_{100} \\ V_{99} = 0.9 V_{98} + 0.1 \theta_{99} \\ V_{98} = 0.9 V_{97} + 0.1 \theta_{98} \\ \dots \end{cases}$$

$$V_{100} = 0.1 \theta_{100} + \underline{0.9 V_{99}} \quad \textcircled{1}$$

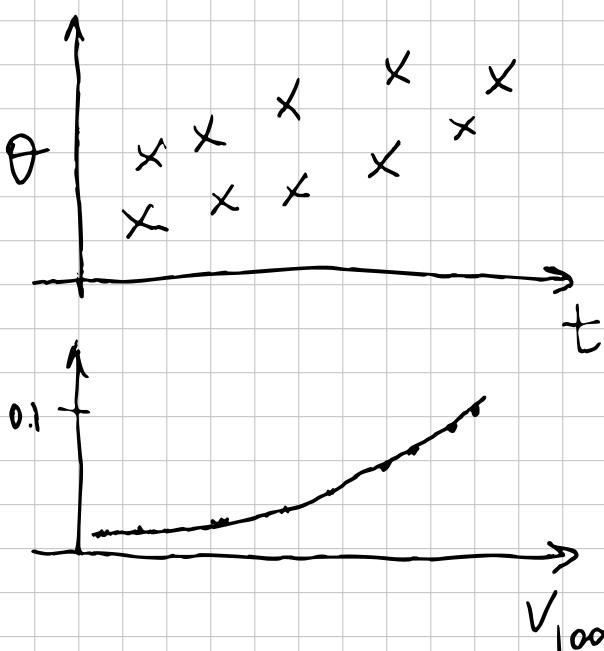
$$\textcircled{2} \quad 0.1 \theta_{100} + 0.9 (0.1 \theta_{99} + \underline{0.9 V_{98}}) \quad \textcircled{3}$$

$$\textcircled{4} \quad 0.1 \theta_{100} + 0.9 (0.1 \theta_{99} + 0.9 (0.1 \theta_{98} + \underline{0.9 V_{97}}))$$

and so on.

$$V_{100} = 0.1 \theta_{100} + 0.09 \theta_{99} + 0.081 \theta_{98} \quad \textcircled{5}$$

$$\textcircled{6} \quad 0.1 \cdot (0.9)^2 \theta_{97} + \dots$$



$$(1-\varepsilon)^{1/\varepsilon} \approx \frac{1}{e}$$

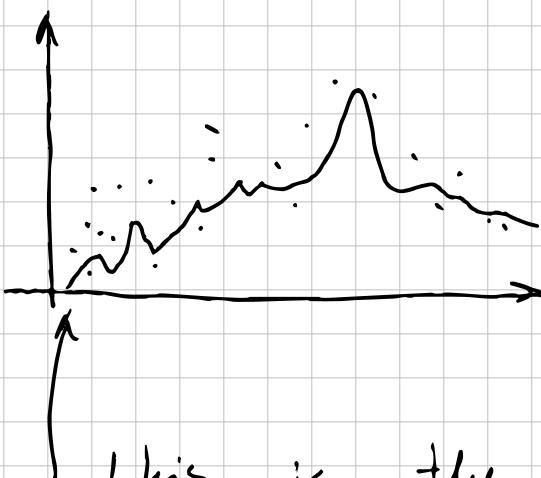
$$0.9^{\infty} \approx \frac{1}{e}$$

$$0.98^{50} \approx \frac{1}{e}$$

Bias correction

Using $V_t = \beta V_{t-1} + (1-\beta) \theta_t$

we get curve that is
low



this is the problem

Original algorithm

$$V_0 = 0$$

$$V_t = \underbrace{0.9 V_{t-1}}_0 + 0.1 \theta_t$$

...

and so on.

Instead of taking V_t
take $\frac{V_t}{1-\beta^t}$ — bias correction

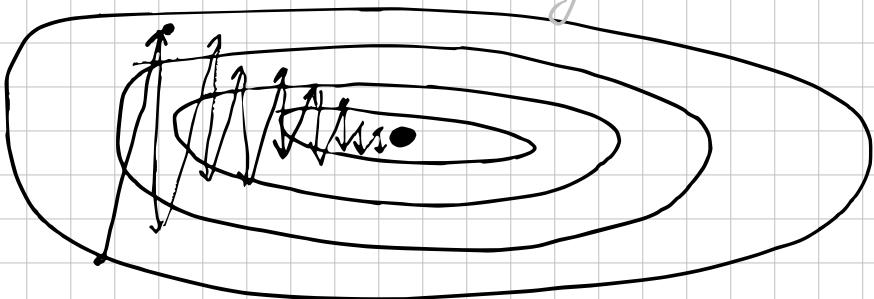
$$V_t = \beta \frac{V_{t-1}}{1-\beta^t} + (1-\beta) \theta_t$$

$\frac{V_t}{1-\beta^t}$ becomes closer to $\underline{V_t}$

when t is big

Gradient descent with momentum

✓ regular
gradient descent



You want:

↑ slower learning ← faster learning

Momentum:

On iteration t :

Compute dW, dt on current batch / mini-batch

$$V_{dw} = \beta V_{dw} + (1 - \beta) dW$$

$$V_{db} = \beta V_{db} + (1 - \beta) db$$

$$W := W - \alpha V_{dw}$$

$$b := b - \alpha V_{db}$$

Hyperparameters: α, β

In practice people don't bother with bias correction.

RMS prop
(root mean square)

On iteration t :

Compute dW, db on current
batch (mini-batch) element-wise

$$S_{dw} = \beta S_{dw} + (1 - \beta) dW^2$$

$$S_{db} = \beta S_{db} + (1 - \beta) db^2$$

$$w := w - \alpha \frac{dW}{\sqrt{S_{dw}}}$$

$$b := b - \alpha \frac{db}{\sqrt{S_{db}}}$$

Adam optimization

algorithm

$$V_{dw} = 0, S_{dw} = 0$$

$$V_{db} = 0, S_{db} = 0$$

On iteration t :

Compute dW, db using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dW \quad \left. \begin{array}{l} \text{momentum} \\ \text{grad.} \end{array} \right\}$$

$$V_{db} = \beta_1 V_{db} + (1 - \beta_1) db \quad \left. \begin{array}{l} \text{grad.} \\ \text{desc.} \end{array} \right\}$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dW^2 \quad \left. \begin{array}{l} \text{RHS prop} \end{array} \right\}$$

$$S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2 \quad \left. \begin{array}{l} \text{RHS prop} \end{array} \right\}$$

$$V_{dw}^c = \frac{V_{dw}}{(1 - \beta_1)^t}$$

$$V_{db}^c = \frac{V_{db}}{(1 - \beta_1)^t}$$

$$S_{dw}^c = \frac{S_{dw}}{(1 - \beta_2)^t}$$

$$S_{db}^c = \frac{S_{db}}{(1 - \beta_2)^t}$$

$$w = w - \alpha \frac{V_{dw}^c}{\sqrt{S_{dw}^c + \epsilon}} \quad \left. \begin{array}{l} \text{bias} \\ \text{correction} \end{array} \right\} \quad \text{to avoid devision by } 0$$

$$b = b - \alpha \frac{V_{db}^c}{\sqrt{S_{db}^c + \epsilon}}$$

Hyperparameters:

α : needs to be tuned

β_1 : 0.9 (first moment)

β_2 : 0.999 (second moment)

ϵ : 10^{-8}

Adam stands for adaptive

moment estimation

Learning rate decay

1 epoch = 1 pass through the data

$$\alpha = \frac{1}{1 + \text{decay_rate} \times \text{epoch_number}} \cdot \alpha_0$$

another hyper parameter

Other learning rate decay methods

$$1. \alpha = 0.95^{\text{epoch-number}} \cdot \alpha_0$$

- exponential decay

$$2. \alpha = \frac{k}{\sqrt{\text{epoch-number}}} \cdot \alpha_0 \quad \text{or}$$

$$\frac{k}{\sqrt{t}} \alpha_0$$



4. Manual decay

The problem of local optima

Most likely local points are saddle points

Problem of plateaus



plateaus are parts of functions where the gradient is close to 0 for a long time

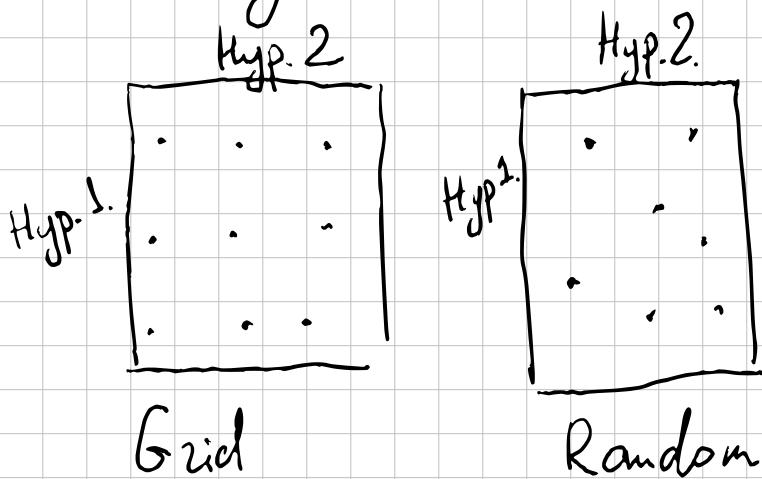
Notes:

- Unlikely to get stuck in a bad local optima
- Plateaus can make learning slow

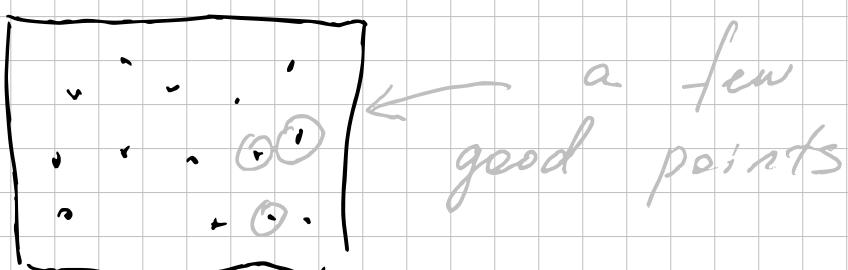
Hyperparameter Tuning

Hyperparameters

- learning rate (α)
 - # hidden-unit
 - mini-batch size
 - momentum (β)
 - learning-rate decay
- Try random values: don't use a grid



Coarse to fine

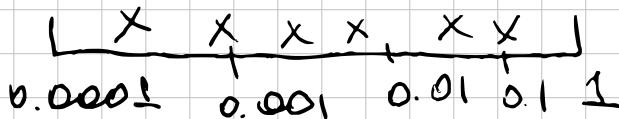


Zoom-in "good" region
and perform thorough
search there

Using an Appropriate Scale to pick hyperparameters

1. Picking hyperparameters at random

$$\lambda = 0.0001$$



$$r = -4 \times \text{np.random.rand}()$$

$$\lambda = 10^r$$

Sampling on a log scale

2. Hyperparameters for exponentially weighted averages

$$\beta = 0.9 \dots 0.999$$

$$\left\{ \begin{array}{c} \\ \downarrow \\ \end{array} \right\}$$

avg. 10 d-ts over d-ts
over

We want to explore

$$\underline{1-\beta} = 0.1 \dots 0.001$$



uniformly random sampling

$$r \in [-3, -1]$$

$$1-\beta = 10^r$$

$$\beta = 1 - 10^r$$

Tips & tricks

2 ways to tune
hyperparameters

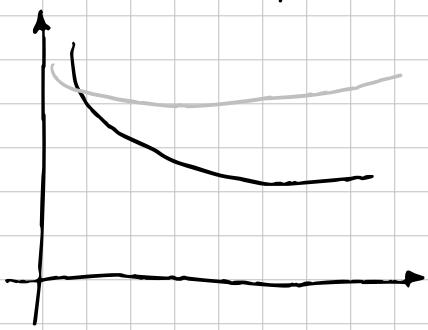
Babysitting one
model



"Panda"
approach

* "One kid
(trial)"

Train many
models in parallel



Model 1

Model 2

"Caviar"
approach

* "A lot of
kids (trials)"

** requires a lot
of computational
power

Normalizing activations

in a network

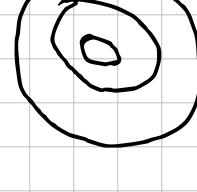
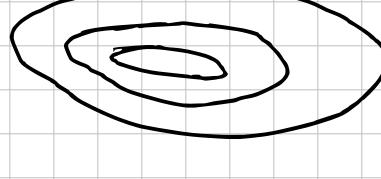
regular normalization

$$\mu = \frac{1}{m} \sum_i x^{(i)}$$

$$x = x - \mu$$

$$G^2 = \frac{1}{m} \sum_i x^{(i)}^2$$

$$x = x / G$$



Can we normalize $a^{(i)}$ so
to train $w^{(i)}, b^{(i)}$ faster

Normalize $\tilde{z}^{(i)}$
(in practice this is done often)

Given some intermediate values

in NN $\underbrace{z^{(1)}, \dots, z^{(m)}}_{z^{(l)(i)}}$

$$\mu = \frac{1}{m} \sum_i z^{(i)}$$

$$G^2 = \frac{1}{m} \sum_i (z_i - \mu)^2$$

$$\tilde{z}_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{G^2 + \epsilon}}$$

$$\tilde{z}^{(i)} = \gamma \tilde{z}_{\text{norm}}^{(i)} + \beta$$

γ, β - learnable parameters
of model

$$\text{If } \gamma = \sqrt{G^2 + \epsilon}$$

$$\beta = \mu$$

$$\text{then } \tilde{z}^{(i)} = z^{(i)}$$

use $\tilde{z}^{(i)}$ instead of $z^{(i)}$

Fitting Batch Norm to a network



Node computes 2 functions:

z, a :

$$X \xrightarrow{\text{batch norm}} z \xrightarrow{\beta, \gamma} \tilde{z} \xrightarrow{} a$$

Normalised value

Parameters

$$w^{(i)}, b^{(i)}, \beta^{(i)}, \gamma^{(i)}$$

Working with mini-batches

$$X^{(l)} \xrightarrow{w^{(l)}, b^{(l)}} z^{(l)} \xrightarrow{\beta^{(l)}, \gamma^{(l)}} \tilde{z}^{(l)} \xrightarrow{} a^{(l)} \rightarrow \dots$$

$$z^{(l)} = w^{(l)} a^{(l-1)} + b^{(l)}$$



Normalization will disregard any constant added to $z^{(l)}$ when will subtract

mean, therefore if using Batch Norm parametered $b^{(l)}$ can be dropped

$$z^{(l)} = w^{(l)} a^{(l-1)}$$

z_{norm}

$$\tilde{z}^{(l)} = \gamma^{(l)} z_{\text{norm}}^{(l)} + \beta^{(l)}$$

$$\beta^{(l)} \quad \gamma^{(l)}$$

dims $(n^{(l)}, 1)$ $(n^{(l)}, 1)$

Implementing gradient descent

for $t=1 \dots \# \text{mini-batches}$:

- Compute forward prop on X in each hidden layer,

$$\text{use BN to replace } z^{(l)} \text{ with } \tilde{z}^{(l)}$$

- Use backprop to compute $dW^{(l)}, d\beta^{(l)}, d\gamma^{(l)}$

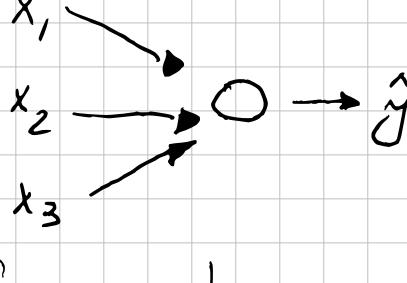
- Update params

$$w^{(l)} := w^{(l)} - \alpha dW^{(l)}$$

$$\beta^{(l)} := \beta^{(l)} - \alpha d\beta^{(l)}$$

...
 $\gamma^{(l)}$

Why does Batch Norm work?



If you training set is

But dev set



You may generalize
badly

This change is called
"covariate shift"

* Consider a deeper network

Batch norm reduces the
distribution: no matter how
hidden layers parameters are
changing the mean and
variance stays the same.

Later layers of the NN to
have more "stable ground"

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch
 - This adds some noise to the values $z^{(l)}$ within that mini-batch. So, similar to dropout, it adds some noise to each hidden layer's activations.
 - This has a slight regularization effect.

Batch Norm at Test time

$$\bar{\mu} = \frac{1}{m} \sum_i z^{(i)}$$

$$G^2 = \frac{1}{m} \sum_i (z^{(i)} - \bar{\mu})^2$$

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \bar{\mu}}{\sqrt{G^2 + \epsilon}}$$

$$\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

$\bar{\mu}, G^2$: estimate using exponentially weighted average across mini-batches)

$$\begin{array}{ccc} X^{111}, & X^{121}, & X^{131} \\ | & | & | \\ \mu^{111\{l\}}, & \mu^{121\{l\}}, & \mu^{131\{l\}} \xrightarrow{} \mu^{L\{l\}} \\ \theta_1 & \theta_2 & \theta_3 \end{array}$$

$$G^2 \xrightarrow{111\{l\}} G^2 \xrightarrow{121\{l\}} G^2 \xrightarrow{131\{l\}} G^2$$

*

$$z_{\text{norm}} = \frac{z - \bar{\mu}}{\sqrt{G^2 + \epsilon}}$$

$$\tilde{z} = \gamma z_{\text{norm}} + \beta$$

Softmax regression

Helps to detect multiple classes

if dim is $(4, 1)$

Softmax layer

In the last layer ℓ

$$z^{[\ell]} = w^{[\ell]} a^{[\ell-1]} + b^{[\ell]}$$

• Activation function:

$$t = e^{z^{[\ell]}} \quad (\text{dim is } (4, 1))$$

$$a^{[\ell]} = \frac{e^{z^{[\ell]}}}{\sum_i t_i} = \frac{t_i}{\sum_i t_i}$$

Understanding softmax

"hard max" gets max and considers it as $\underline{1}$ and everything else $\underline{0}$.

Softmax regression generalizes logistic regression to C classes.

- * If $C = 2$, softmax reduces to logistic regression.

Loss function

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

\hat{y} is probability

$$L(\hat{y}, y) = - \sum_i^4 y_i \log \hat{y}_i$$

to make it small is to make \hat{y}_2 bigger.

$$J(w^{(1)}, b^{(1)}, \dots) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

Backprop:

$$\delta z^{(l)} = \hat{y} - y$$

Orthogonalisation

Orthogonalisation refers to a process of making hyperparameters independent from each other (i.e. change of one does not affect the ^{other} chain of assumptions) in ML

- Fit training set well on cost function



- Fit dev set well on cost function



- Fit test set well on cost function



- Performs well in real world

Single number valuation metric

Ex.

Classifier	Precision	Recall
A	95%	90%
B	88%	85%

Is B better than A ?

It is not clear what of them is better

Possible solution: use F1-score as a combined metric for precision and recall

Classifier F1

A 92.4%

B 91.0%

* F1 is a harmonic mean
$$\left(\frac{2}{\frac{1}{P} + \frac{1}{R}} \right)$$

Satisficing and Optimising metric

Ex.

Classifier	Accuracy	Running time
A	99%	80 ms
B	92%	95 ms
C	95%	1500 ms

One thing to do is to
combine accuracy and
running time in one metric

Another option: maximize accuracy as a subject to running time ($\leq 100\text{ ms}$)

The diagram consists of two curved arrows originating from the words 'accuracy' and 'running time' in the main text. The arrow from 'accuracy' points downwards to the word 'optimising'. The arrow from 'running time' points downwards to the word 'satisficing'.

Pick one as optimising
and other metrics as
satisficing

Train | Dev | Test distributions

dev - development / hold out
cross-validation set

dev and test sets should
come from the same
distribution

You spend a lot of time
optimising dev set (development
time), but when it comes
to test it may not
generalise good (if dev
and test are from
different dist.)

Size of train/dev/test

98%	1%	1%
train	dev	test
for	<u>10^6 examples</u>	

- * Reasonable to use smaller sets for larger data sets.

Set your test set to be big enough to give high confidence in the overall performance of your system.

It is fine to do not have test set. (it is not recommended though)

Convolutional Neural Networks

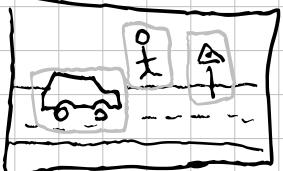
Computer Vision Problems

- Image Classification



64×64

- Object detection



- Neural Style Transfer

Input in computer vision, usually, is really big.

Convolution operation

Edge detection

$$\begin{array}{c}
 \text{convolution} \\
 \boxed{\begin{array}{ccccccccc} 3 & 0 & 1 & 2 & 7 & 4 \\ 1 & 5 & 8 & 9 & 3 & 1 \\ 2 & 7 & 2 & 5 & 1 & 3 \\ 0 & 1 & 3 & 1 & 7 & 8 \\ 4 & 2 & 1 & 6 & 2 & 8 \\ 2 & 4 & 5 & 2 & 3 & 9 \end{array}} * \boxed{\begin{array}{ccc} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{array}} = \boxed{\begin{array}{ccc} & & \\ & & \\ & & \end{array}}
 \end{array} \\
 \text{filter/kernel} \\
 6 \times 6 \times 1 \text{ (grayscale)} \\
 \text{---}$$

$$\begin{array}{c}
 \text{---} \\
 \text{---} \\
 \begin{array}{c}
 \begin{array}{cccc} -5 & -4 & 0 & 8 \\ -10 & -2 & 2 & 3 \\ 0 & -2 & -4 & -7 \\ -3 & -2 & -3 & -16 \end{array} \\
 4 \times 4
 \end{array}
 \end{array} \\
 \begin{array}{cccc} 3.1 & 0.0 & 1.0 & -1 \\ 1.1 & 5.0 & 8.0 & -1 \\ 2.1 & 7.0 & 2.0 & -1 \end{array}$$

python: conv-forward

tensorflow: tf.nn.conv2d

keras: Conv2D

Vertical edge detection examples

$$\begin{array}{c}
 \begin{array}{cccccc} 10 & 10 & 10 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 \end{array} * \begin{array}{ccc} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{array} = \begin{array}{cccccc} 0 & 30 & 30 & 0 & 0 & 0 \\ 0 & 30 & 30 & 0 & 0 & 0 \\ 0 & 30 & 30 & 0 & 0 & 0 \\ 0 & 30 & 30 & 0 & 0 & 0 \end{array} \\
 \begin{array}{c} \boxed{\square} \\ \boxed{\square} \end{array} \quad \begin{array}{c} \boxed{\square} \\ \boxed{\square} \\ \boxed{\square} \end{array} \quad \begin{array}{c} \boxed{\square} \\ \boxed{\square} \\ \boxed{\square} \\ \boxed{\square} \end{array}
 \end{array}$$

$$\begin{array}{c}
 \begin{array}{cccccc} 0 & 0 & 0 & 10 & 10 & 10 \\ 0 & 0 & 0 & 10 & 10 & 10 \\ 0 & 0 & 0 & 10 & 10 & 10 \\ 0 & 0 & 0 & 10 & 10 & 10 \\ 0 & 0 & 0 & 10 & 10 & 10 \\ 0 & 0 & 0 & 10 & 10 & 10 \end{array} * \begin{array}{ccc} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{array} = \begin{array}{cccccc} 0 & -30 & -30 & 0 & 0 & 0 \\ 0 & -30 & -30 & 0 & 0 & 0 \\ 0 & -30 & -30 & 0 & 0 & 0 \\ 0 & -30 & -30 & 0 & 0 & 0 \end{array} \\
 \begin{array}{c} \boxed{\square} \\ \boxed{\square} \end{array} \quad \begin{array}{c} \boxed{\square} \\ \boxed{\square} \\ \boxed{\square} \end{array} \quad \begin{array}{c} \boxed{\square} \\ \boxed{\square} \\ \boxed{\square} \\ \boxed{\square} \\ \boxed{\square} \end{array}
 \end{array}$$

In this filter there is a difference between light-to-dark and dark-to-light transitions.

If you don't care you may use absolute values.

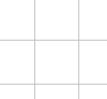
Filters

$$\begin{array}{ccc}
 \begin{array}{ccc} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{array} & & \begin{array}{ccc} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{array}
 \end{array}$$

Vertical

Horizontal

$$\begin{array}{c}
 \begin{array}{cccccc} 10 & 10 & 10 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 \\ 0 & 0 & 0 & 10 & 10 & 10 \\ 0 & 0 & 0 & 10 & 10 & 10 \\ 0 & 0 & 0 & 10 & 10 & 10 \end{array} * \begin{array}{ccc} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{array} = \begin{array}{cccccc} 0 & 0 & 0 & 0 & 0 & 0 \\ 30 & 10 & -10 & -30 & 0 & 0 \\ 30 & 10 & -10 & -30 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{array}
 \end{array}$$



Generally learning a filter is more robust approach.

Sobel filter

$$\begin{array}{c}
 w_1 \quad w_2 \quad w_3 \\
 w_4 \quad w_5 \quad w_6 \\
 w_7 \quad w_8 \quad w_9
 \end{array}$$

We can treat filter values as parameters

and learn them at

backpropagation step

3 0 -3

10 0 -10

3 0 -3

Sobel filter

$$\begin{array}{c}
 w_1 \quad w_2 \quad w_3 \\
 w_4 \quad w_5 \quad w_6 \\
 w_7 \quad w_8 \quad w_9
 \end{array}$$

We can treat filter values as parameters

and learn them at

backpropagation step

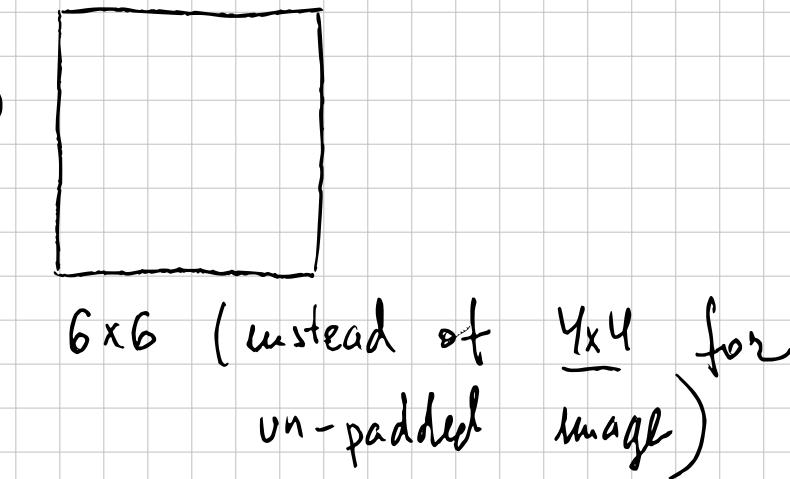
Padding

$$n \times n \text{ (image)} * f \times f \text{ (filter)} = (n-f+1) \times (n-f+1)$$

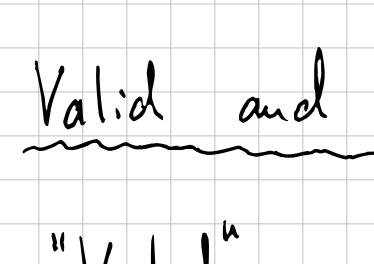
Problems:

- Shrinking output
- Throwing away a lot of information

To solve it we can pad image



$$6 \times 6 \rightarrow 4 \times 4$$



6x6 (instead of 4x4 for un-padded image)

by convention:

- padded values are 0
- $p = \text{padding} = \frac{f-1}{2}$ (in this case)

Valid and Same Convolutions

"Valid" means no padding

$$6 \times 6 * 3 \times 3 \rightarrow 4 \times 4$$

"Same" means to pad so that output size is the same as the input size.

$$(n + 2p - f + 1) \times (n + 2p - f + 1)$$

$$n + 2p - f + 1 = n$$

$$2p = f - 1$$

$$p = \frac{f-1}{2}$$

$$\boxed{p = \frac{f-1}{2}}$$

By convention f is usually odd.

There are 2 reasons:

- If f is even you need an asymmetric padding

- Even filters don't have centered position

Strided convolution

Stride = 2

$$\begin{matrix} 2 & 3 & 7 & 4 & 6 & 2 & 9 \\ 6 & 6 & 9 & 8 & 7 & 4 & 3 \\ 3 & 4 & 8 & 3 & 8 & 9 & 7 \\ 7 & 8 & 3 & 6 & 6 & 3 & 4 \\ 4 & 2 & 1 & 8 & 3 & 4 & 6 \\ 3 & 2 & 4 & 1 & 9 & 8 & 3 \\ 0 & 1 & 3 & 9 & 2 & 1 & 4 \end{matrix}$$

$$* \begin{matrix} 3 & 4 & 4 \\ -1 & 0 & 2 \\ -1 & 0 & 3 \end{matrix}$$

3×3

7×7

Step 1

91

=

$$\begin{matrix} 2 & 3 & 7 & 4 & 6 & 2 & 9 \\ 6 & 6 & 9 & 8 & 7 & 4 & 3 \\ 3 & 4 & 8 & 3 & 8 & 9 & 7 \\ 7 & 8 & 3 & 6 & 6 & 3 & 4 \\ 4 & 2 & 1 & 8 & 3 & 4 & 6 \\ 3 & 2 & 4 & 1 & 9 & 8 & 3 \\ 0 & 1 & 3 & 9 & 2 & 1 & 4 \end{matrix}$$

$$* \begin{matrix} 3 & 4 & 4 \\ -1 & 0 & 2 \\ -1 & 0 & 3 \end{matrix}$$

3×3

7×7

Step 2

91 100

=

$$\begin{matrix} 2 & 3 & 7 & 4 & 6 & 2 & 9 \\ 6 & 6 & 9 & 8 & 7 & 4 & 3 \\ 3 & 4 & 8 & 3 & 8 & 9 & 7 \\ 7 & 8 & 3 & 6 & 6 & 3 & 4 \\ 4 & 2 & 1 & 8 & 3 & 4 & 6 \\ 3 & 2 & 4 & 1 & 9 & 8 & 3 \\ 0 & 1 & 3 & 9 & 2 & 1 & 4 \end{matrix}$$

$$* \begin{matrix} 3 & 4 & 4 \\ -1 & 0 & 2 \\ -1 & 0 & 3 \end{matrix}$$

3×3

7×7

Step 3

91 100 83

=

$$\begin{matrix} 2 & 3 & 7 & 4 & 6 & 2 & 9 \\ 6 & 6 & 9 & 8 & 7 & 4 & 3 \\ 3 & 4 & 8 & 3 & 8 & 9 & 7 \\ 7 & 8 & 3 & 6 & 6 & 3 & 4 \\ 4 & 2 & 1 & 8 & 3 & 4 & 6 \\ 3 & 2 & 4 & 1 & 9 & 8 & 3 \\ 0 & 1 & 3 & 9 & 2 & 1 & 4 \end{matrix}$$

$$* \begin{matrix} 3 & 4 & 4 \\ -1 & 0 & 2 \\ -1 & 0 & 3 \end{matrix}$$

3×3

7×7

Step 4

91 100 83

= 69

=

and so on

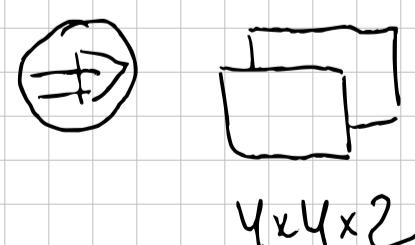
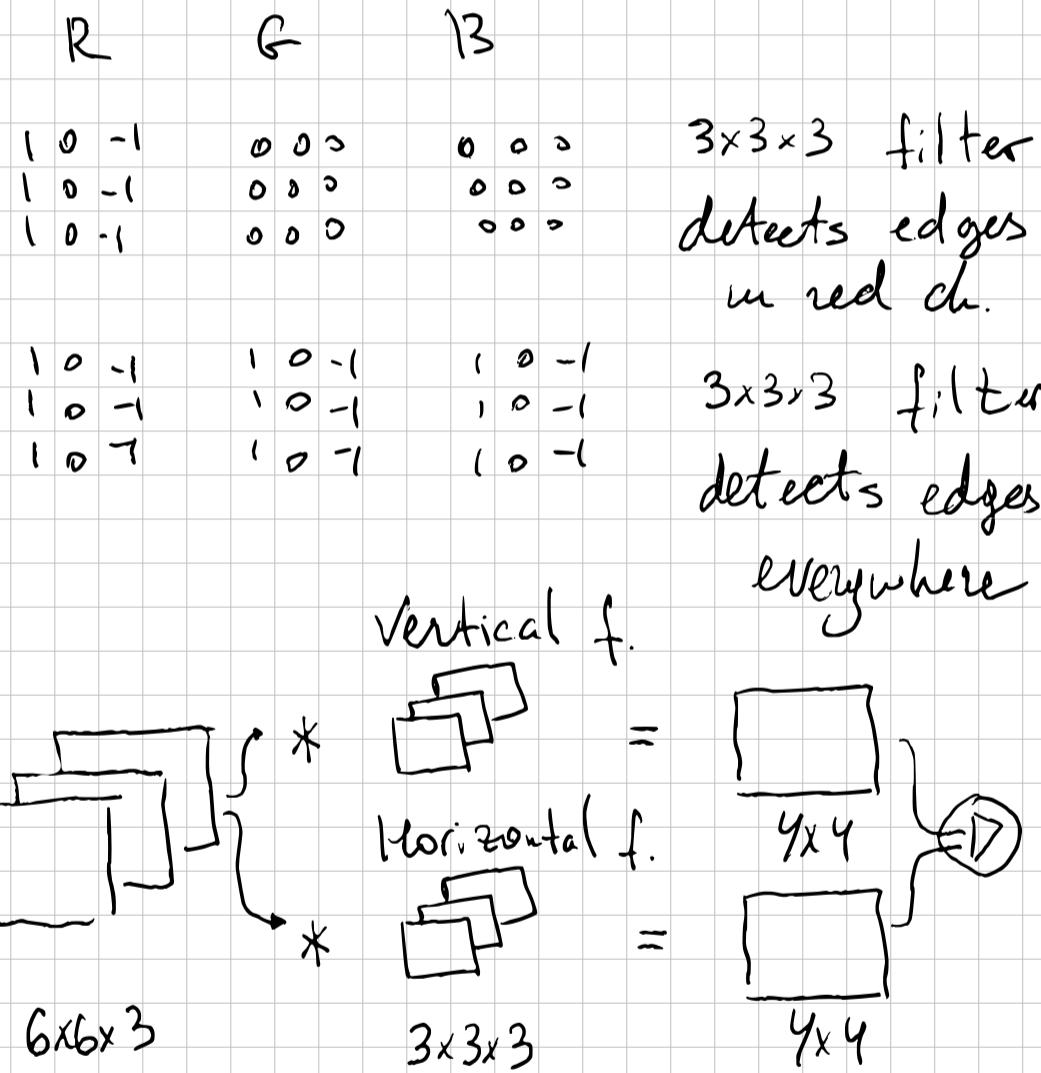
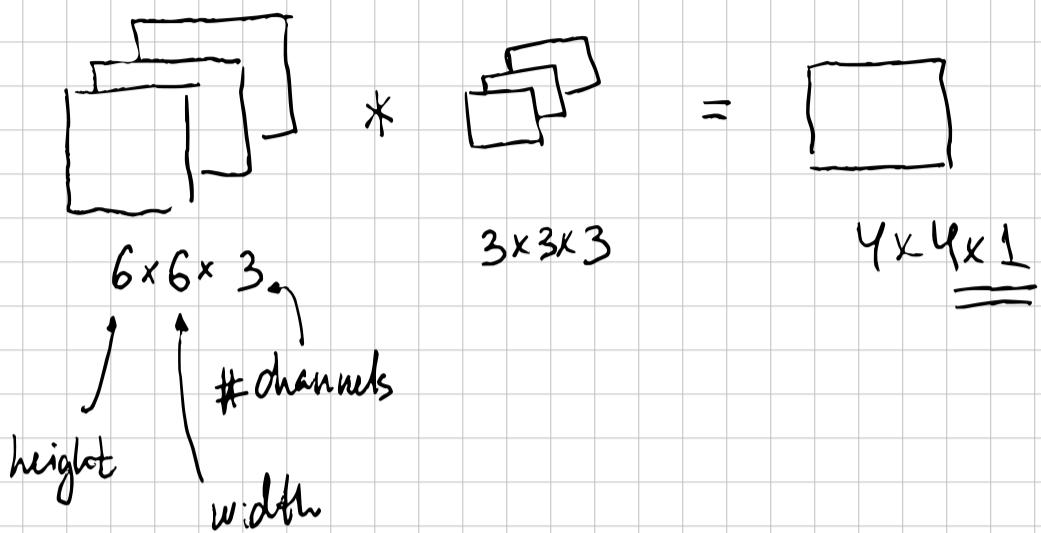
$n \times n * f \times f$
padding p stride s

$$\left\lfloor \frac{n + 2p - f}{s} \right\rfloor + 1 \quad (\times)$$

$$\left\lfloor \frac{n + 2p - f}{s} \right\rfloor + 1 \quad (\times)$$

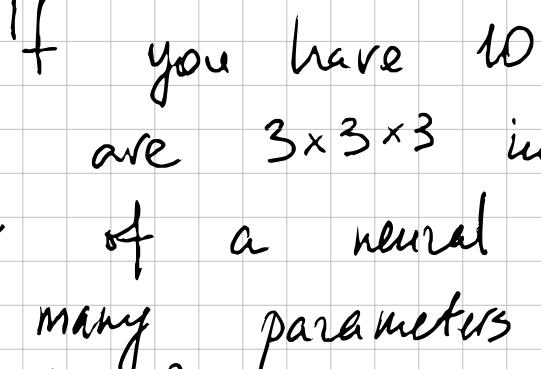
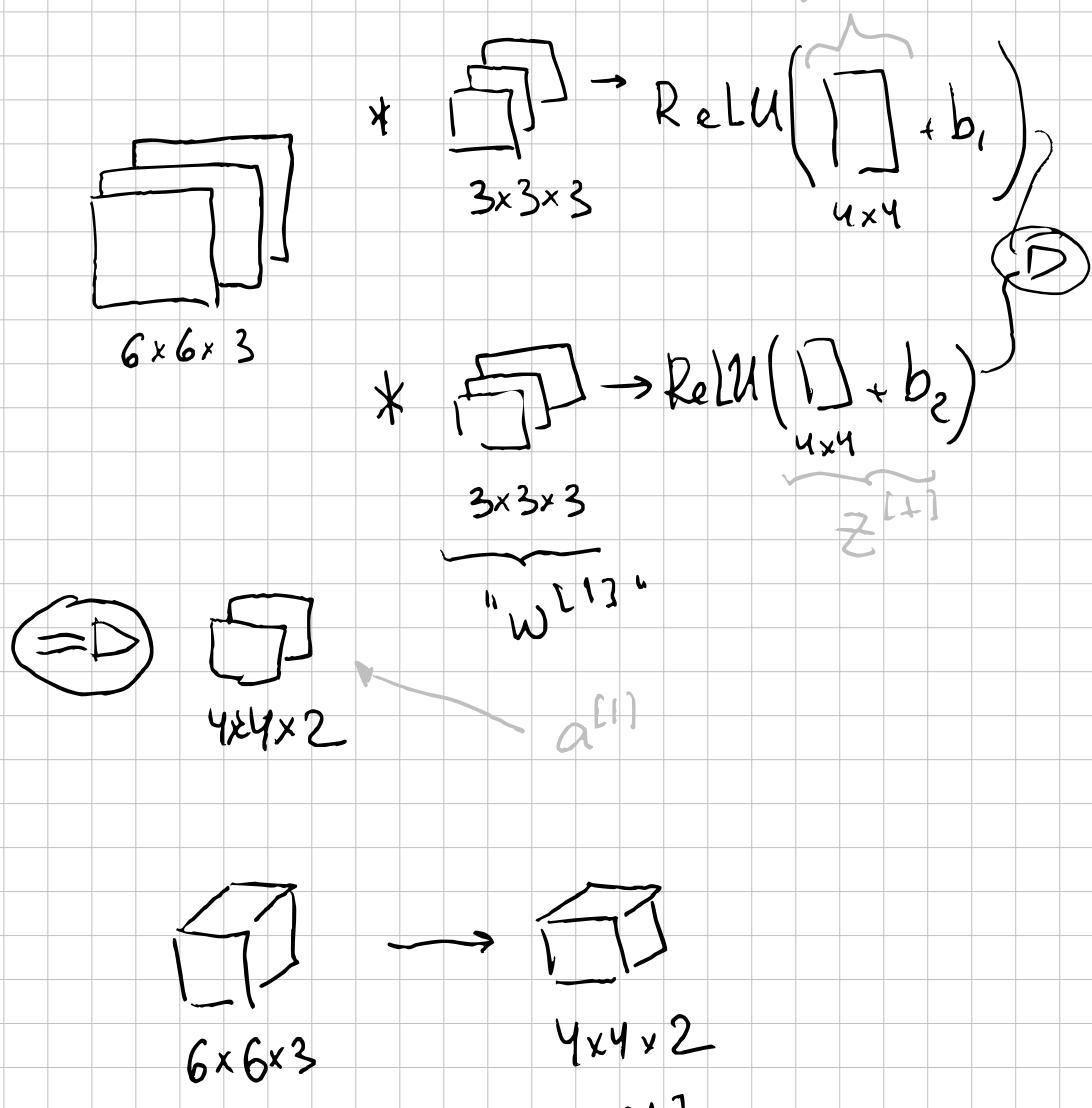
Convolution over volumes

Convolutions on RGB images



Sometimes people say depth
instead of channels.

Example of a layer



Q: If you have 10 filters that are $3 \times 3 \times 3$ in one layer of a neural network, how many parameters does that layer have?

$$\begin{array}{c} \begin{matrix} \square \\ \square \end{matrix} \\ 3 \times 3 \times 3 \\ + \\ \text{bias} \\ \hline 28 \end{array} \quad \times 10 \text{ times}$$

A: 280 parameters

If layer l is a convolutional layer:

$$f^{[l]} = \text{filter size}$$

$$p^{[l]} = \text{padding}$$

$$s^{[l]} = \text{stride}$$

$$\text{Input: } w_H^{[l-1]} \times w_W^{[l-1]} \times w_C^{[l-1]}$$

$$\text{Output: } w_H^{[l]} \times w_W^{[l]} \times w_C^{[l]}$$

$$w_H^{[l]} = \left\lfloor \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$$

$$w_W^{[l]} = \left\lfloor \frac{n_W^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$$

$$w_C^{[l]} = \text{number of filters}$$

Each filter is: $f^{[l]} \times f^{[l]} \times w_C^{[l-1]}$

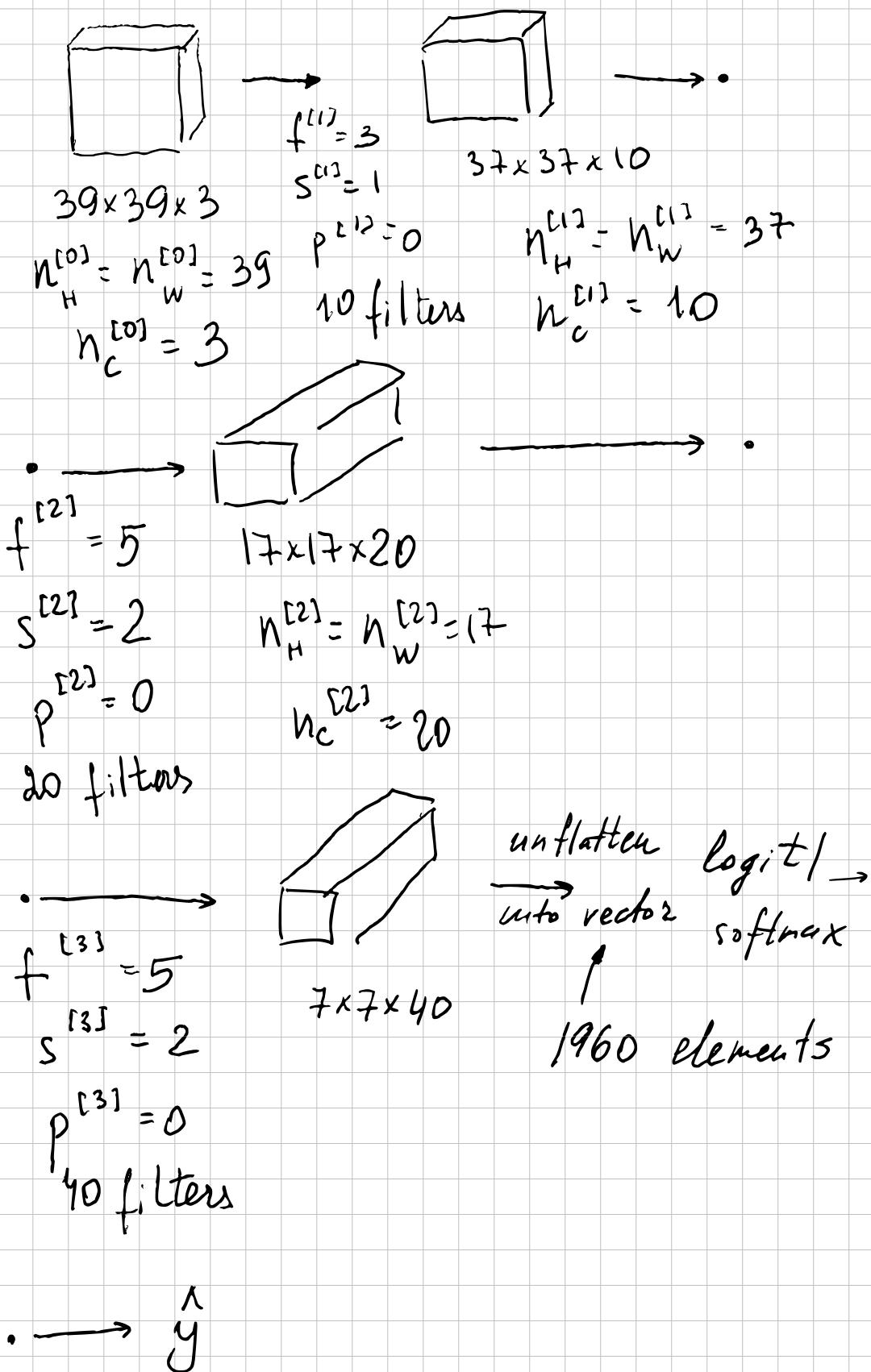
Activations: $a^{[l]} \rightarrow w_H^{[l]} \times w_W^{[l]} \times w_C^{[l]}$

Weights: $f^{[l]} \times f^{[l]} \times w_C^{[l-1]} \times w_C^{[l]}$

number of filters in layer l

Bias: $w_C^{[l]} - (1, 1, 1, w_C^{[l]})$

Simple Convolutional Network Examples



Types of layers in CNN:

- Convolution (CONV)
- Pooling (POOL)
- Fully connected (FC)

Pooling

Max pooling

$$\begin{array}{cccc} 1 & 3 & 2 & 1 \\ 2 & 9 & 1 & 1 \\ \hline 1 & 3 & 2 & 3 \\ 5 & 6 & 1 & 2 \end{array} \longrightarrow \begin{array}{cc} 9 & 2 \\ 6 & 3 \end{array}$$

4×4

Hyperparameters:

$$f = 2$$

$$s = 2$$

Take maximum in the given region.

$$\begin{array}{ccccc} 1 & 3 & 2 & 1 & 3 \\ 2 & 9 & 1 & 1 & 5 \\ 1 & 3 & 2 & 3 & 2 \\ 8 & 3 & 5 & 1 & 0 \\ \hline 5 & 6 & 1 & 2 & 9 \end{array} \xrightarrow{\begin{array}{|c|c|} \hline f & = 3 \\ \hline s & = 1 \\ \hline \end{array}} \begin{array}{cc} 9 & 9 \\ 9 & 5 \\ 8 & 6 \\ 9 & \end{array} \quad 3 \times 3$$

If 3d then pooling preserves number of channels, i.e. pooling computation is done separately for every channel.

Average pooling

$$\begin{array}{cccc} 1 & 3 & 2 & 1 \\ 2 & 9 & 1 & 1 \\ \hline 1 & 4 & 2 & 3 \\ 5 & 6 & 1 & 2 \end{array} \longrightarrow \begin{array}{cc} 5.75 & 1.25 \\ 4 & 2 \end{array}$$

Hyperparameters:

- f : filter size
- s : stride
- Max or average pooling
- p : padding (usually rarely used)

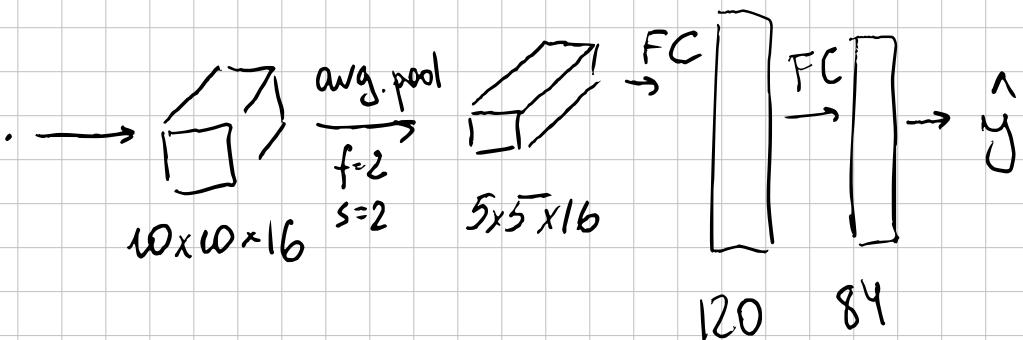
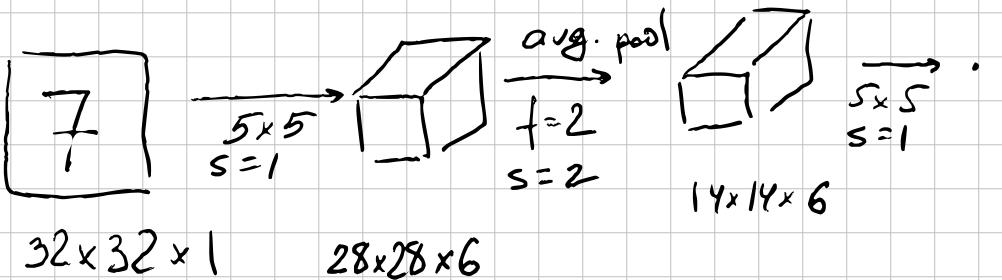
Input

$$n_H \times n_W \times n_C \rightarrow \left\lceil \frac{n_H - f + 1}{s} \right\rceil \times \left\lceil \frac{n_W - f + 1}{s} \right\rceil + n_C$$

(assuming there is no padding)

Classic CNNs

LeNet - 5

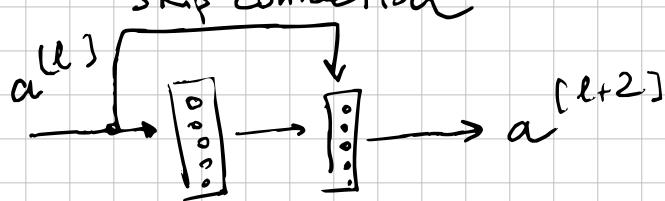


60k parameters

AlexNet

VGG - 16

Residual network



All $a \geq 0$ and use ReLU

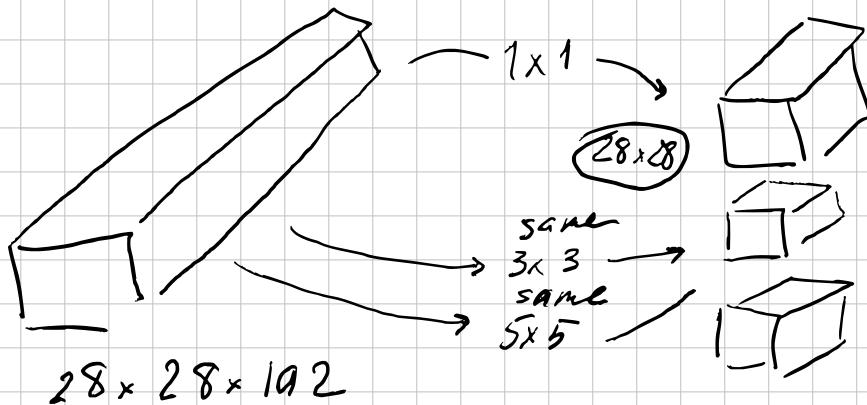
$$a^{[l+2]} = g\left(z^{[l+2]} + \underbrace{a^{[l]}}_{=}\right) = \\ = g\left(\underbrace{w^{[l+2]} a^{[l+1]}}_{=} + b^{[l+2]} + a^{[l]}\right)$$

If $w^{[l+2]} = 0$, $b^{[l+2]} = 0$
 then $\underline{a^{[l+2]} = g(a^{[l]}) = a^{[l]}}$

i.e. residual networks

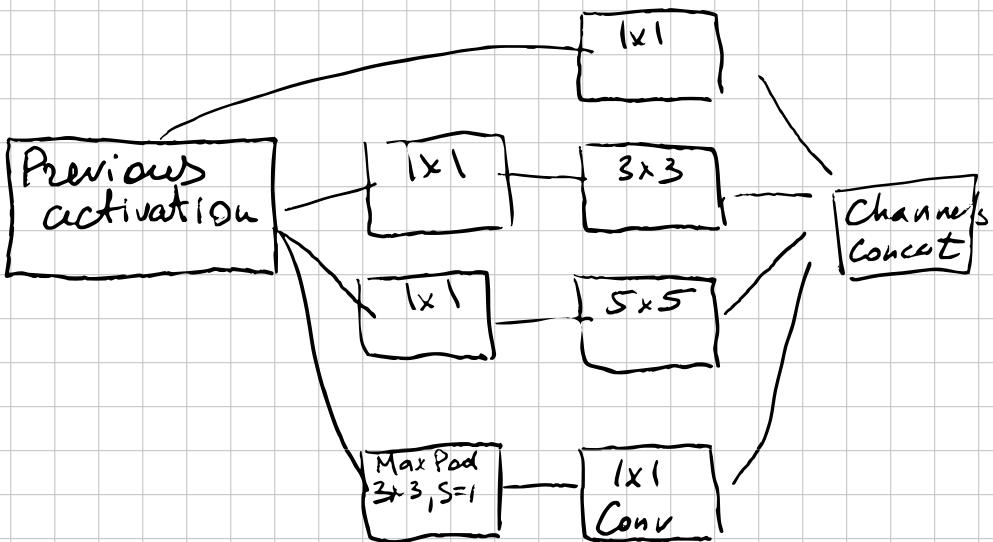
can learn to ignore
 the residual block and it
 does not hurt the network.

Inception network



Stack all
of them
together

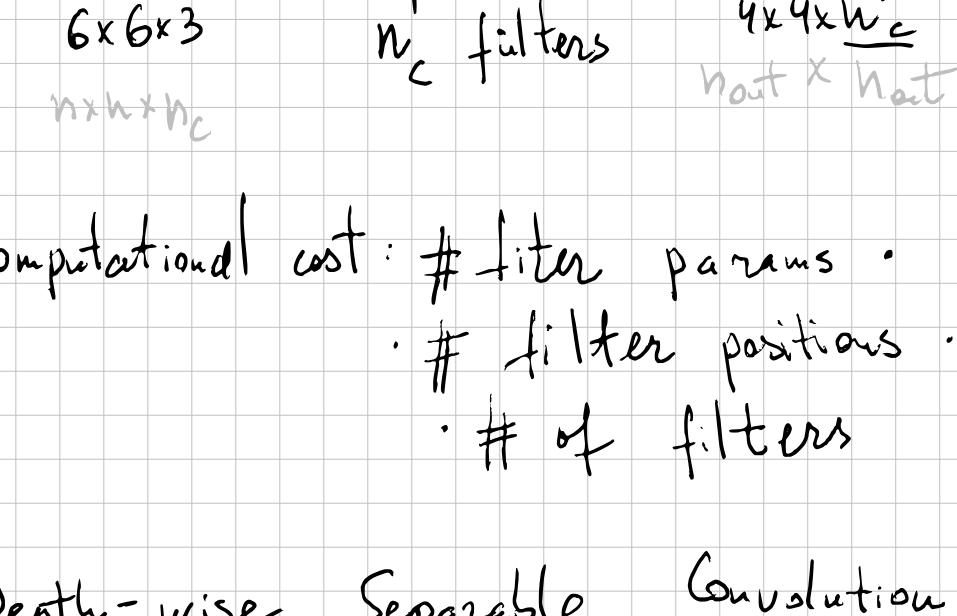
Inception module



MobileNet

- Low computational cost at deployment
- Useful for mobile and embedded vision applications
- Key idea: Normal vs. depthwise-separable convolutions

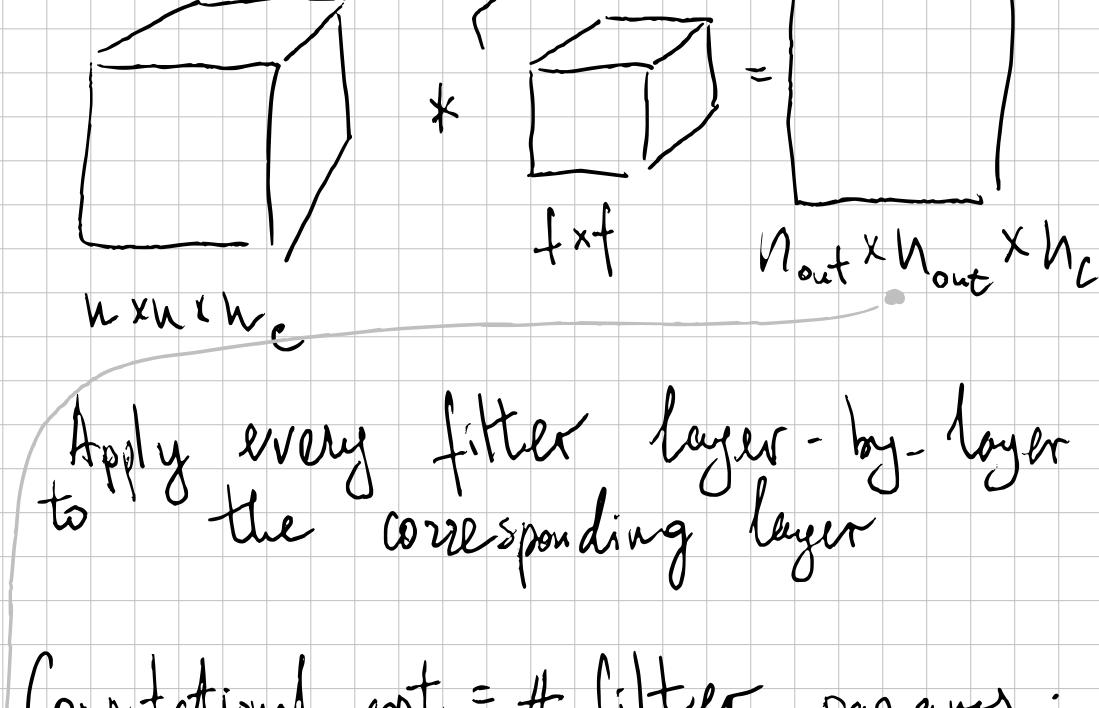
Normal convolution



Computational cost: # filter params ·

- # filter positions ·
- # of filters

Depth-wise Separable Convolution



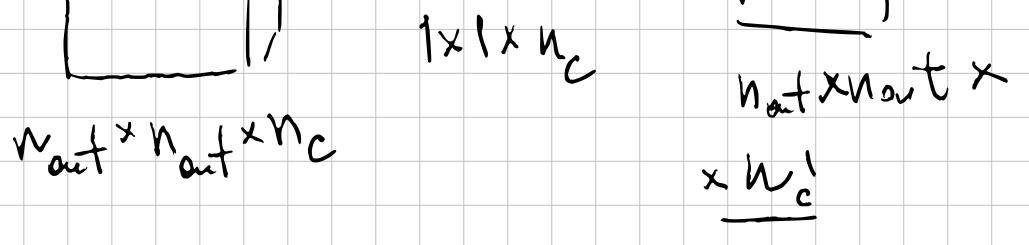
Apply every filter layer-by-layer to the corresponding layer

Computational cost = # filter params ·

filter positions ·

of filters

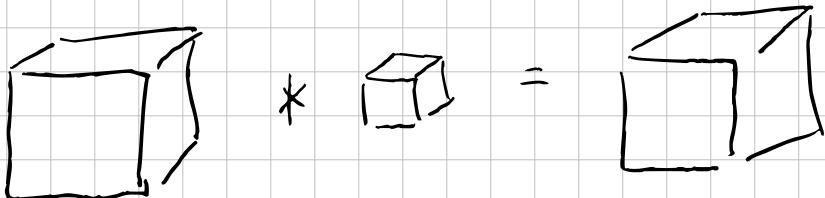
Pointwise conv.



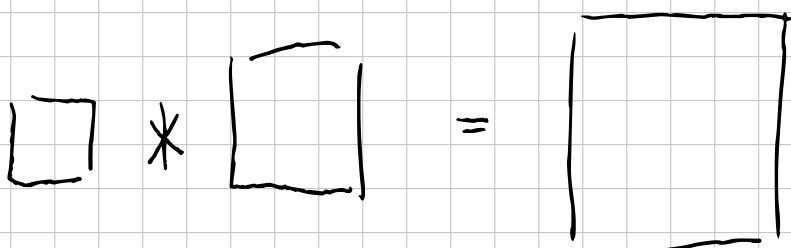
$\times n_c!$

Transpose convolution

- Normal convolution



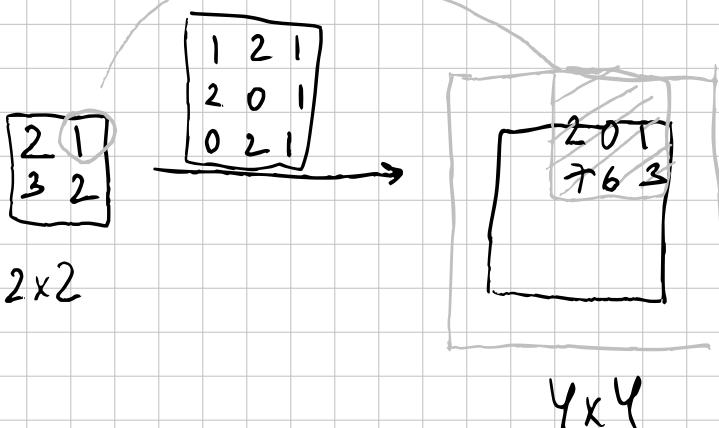
- Transpose convolution



filter $f \times f = 3 \times 3$

padding = 1

stride = 2



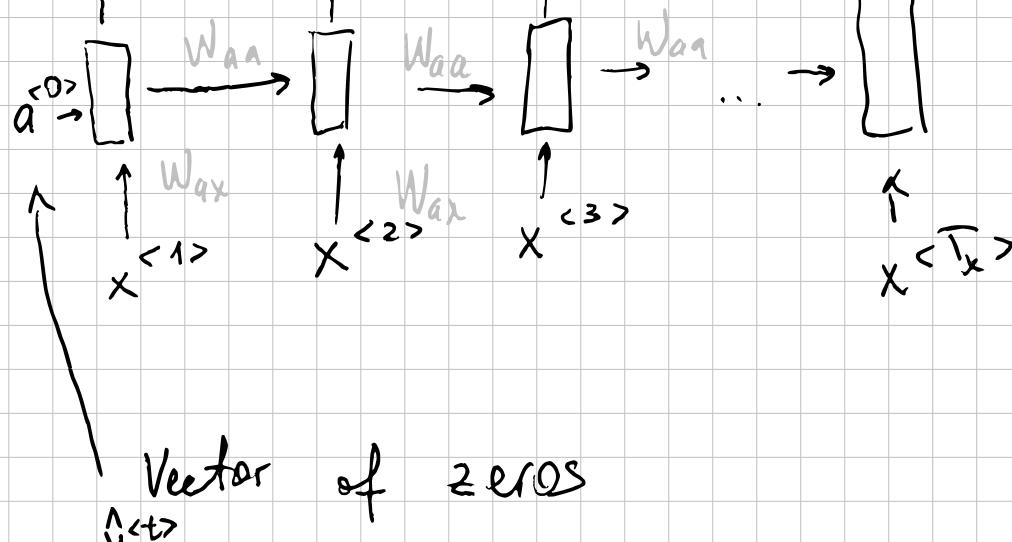
Sequence models

$x^{<t>}$ index of a word in the sequence

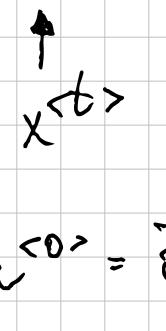
T_y = length of the output sequence

T_x = input sequence length

$T_x^{(i)}$ = input sequence length of a i -th example



Vector of zeros



$$a^{<0>} = \vec{0}$$

$$a^{<1>} = g(W_{aa} a^{<0>} + W_{ax} x^{<1>} + b_a)$$

$$\hat{y}^{<1>} = g(W_{ya} a^{<1>} + b_y)$$

The activation functions usually are : ReLU, tanh

$$a^{<t>} = g(W_{aa} a^{<t-1>} + W_{ax} x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g(W_{ya} a^{<t>} + b_y)$$

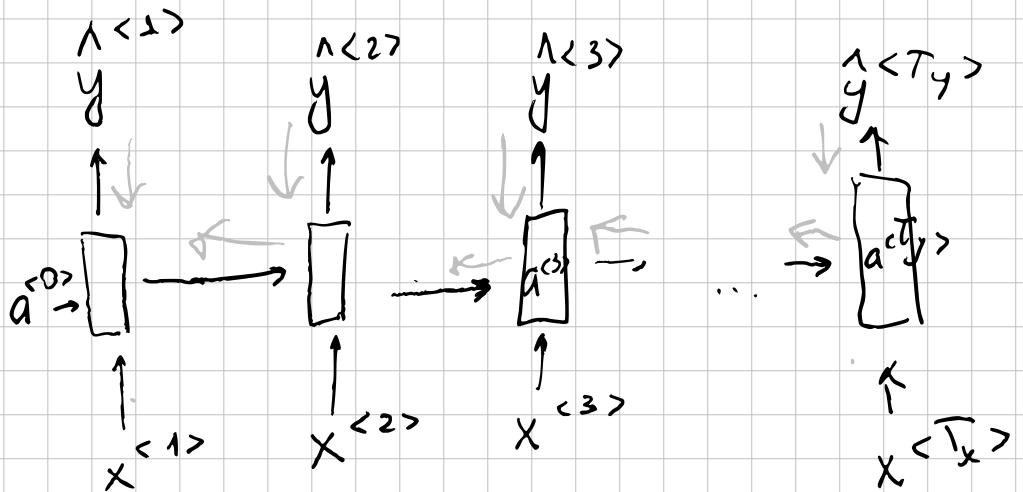
* alternative way to write :

$$a^{<t>} = g(W_a [a^{<t-1>}], x^{<t>}) + b_a$$

$$[W_{aa}; W_{ax}] = W_a$$

$$[a^{<t-1>}, x^{<t>}] = \left[\frac{a^{<t-1>}}{x^{<t>}} \right]$$

Backprop. in RNN



$$\begin{aligned} \mathcal{L}^{<t>} (\hat{y}^{<t>} | y^{<t>}) &= -y^{<t>} \log \hat{y}^{<t>} - \\ &- (1 - \hat{y}^{<t>}) \log (1 - \hat{y}^{<t>}) \end{aligned}$$

$$\mathcal{L}(\hat{y}, y) = \sum_{t=1}^{T_y} \mathcal{L}^{<t>} (\hat{y}^{<t>} | y^{<t>})$$

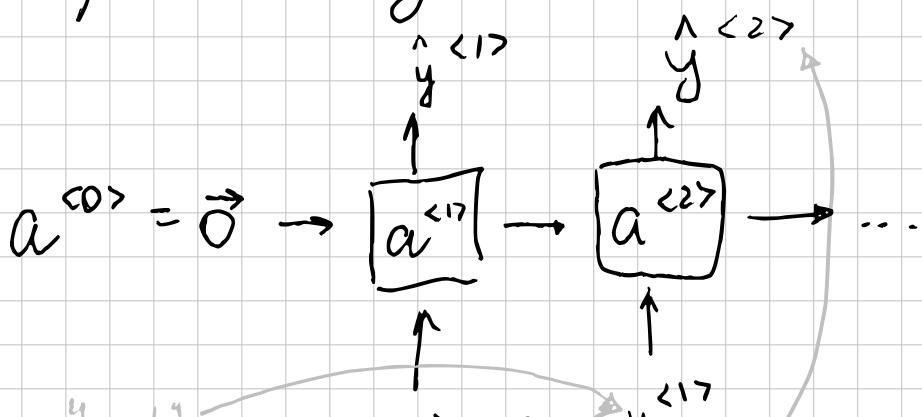
Language model and sequence generation

Training set: large corpus of English text.

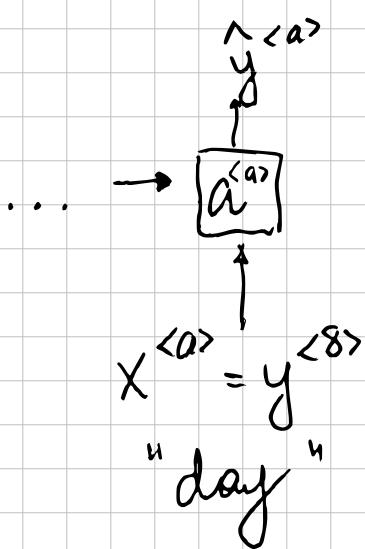
$$x^{(t)} = y^{(t-1)}$$

Ex.

Cats average 15 hours of sleep a day. <EOS>



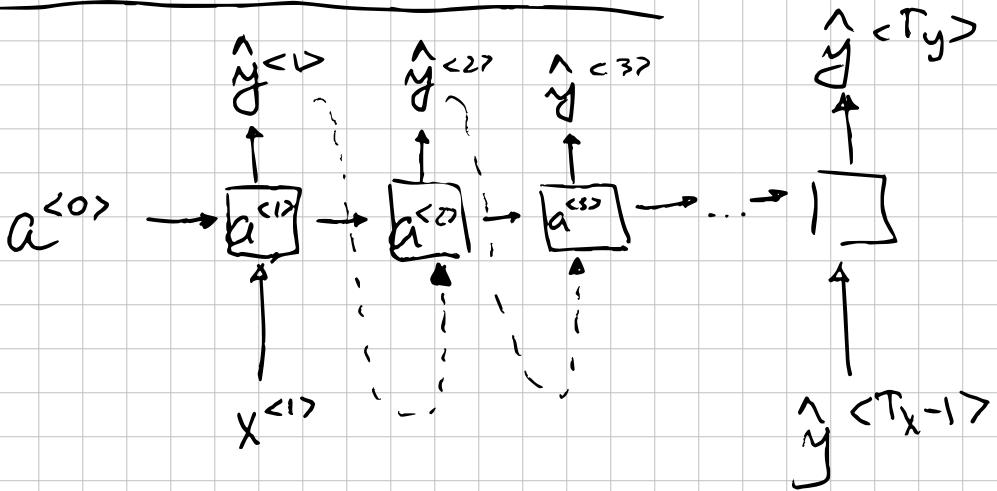
$$P(\text{"average"} | \text{"cat"})$$



$$\mathcal{L}(y^{(t)}, y^{(t)}) = - \sum_i y_i^{(t)} \log \hat{y}_i^{(t)}$$

$$\mathcal{L} = \sum_t \mathcal{L}^{(t)}(y^{(t)}, \hat{y}^{(t)})$$

Sampling a sequence from a trained RNN



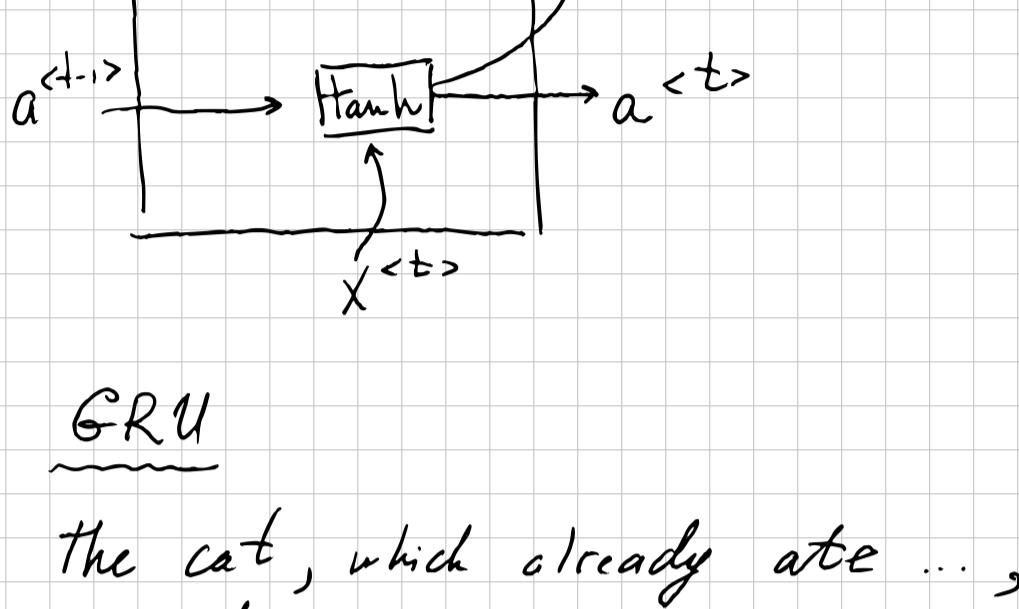
Vanishing gradients with RNNs

Gradient Clipping : "fix" the gradient if the threshold been exceeded.

Gated Recurrent Unit

RNN unit

$$a^{<t>} = g(W_a [a^{<t-1>}, x^{<t>}] + b_a)$$



GRU

the cat, which already ate ... , was full.

C = memory cell

$$\boxed{C^{<t>} = a^{<t>}}$$

$$\tilde{C}^{<t>} = \tanh (W_c [C^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma (W_u [C^{<t-1>}, x^{<t>}] + b_u)$$

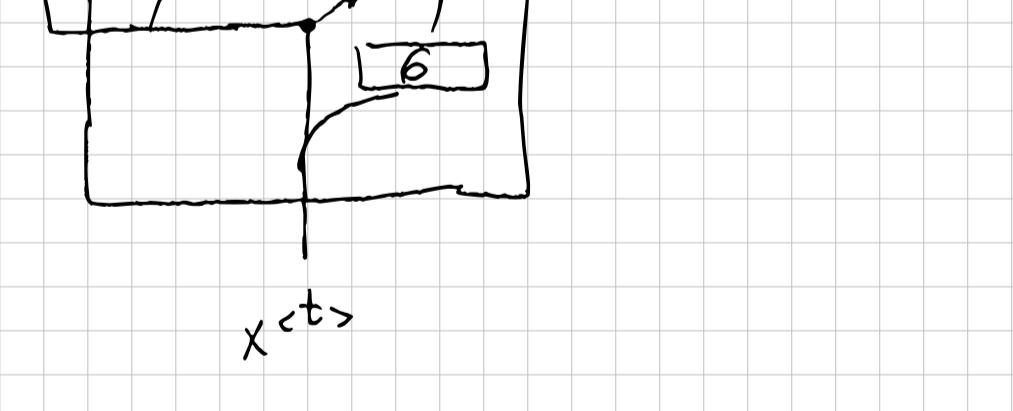
"u" stands for "update"

$$C^{<t>} = \Gamma_u * \tilde{C}^{<t>} + (1 - \Gamma_u) * C^{<t-1>}$$

allows us to

"forget" words

in the sequence we don't need.



* Helps a lot with vanishing gradient problem

Full GRU

$$\tilde{C}^{<t>} = \tanh (W_c [\Gamma_r * C^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma (W_u [C^{<t-1>}, x^{<t>}] + b_u)$$

$$C^{<t>} = \Gamma_u * \tilde{C}^{<t>} + (1 - \Gamma_u) C^{<t-1>}$$

NLP and word embeddings

Word representation

$$V = [a, \text{aaron}, \dots, 200, \langle \text{UNK} \rangle]$$

1-hot representation

$$\text{"man"} (5391) = (0 \ 0 \ \dots \ 1 \dots \ 0)$$

5391
↑

I want a glass of orange —

I want a glass of apple —

Words "orange" & "apple" can be far away from each other in 1-hot encoding which makes it really hard to generalize.

Featurized representation:

word embeddings

	Man	Woman	King	Queen	Apple	Orange
	5391	9853	4914	7157	456	6257
Gender	-1	1	-0.95	0.97	0.0	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.0
Age	0.03	0.02	0.7	0.69	0.03	-0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97
:						

t-SNE algorithm

Strategies

Man → Woman

King → ?

$$\ell_{\text{man}} - \ell_{\text{woman}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

← gender

$$\ell_{\text{king}} - \ell_{\text{queen}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

← gender

$$\ell_{\text{man}} - \ell_{\text{woman}} \approx \ell_{\text{king}} - \ell_{?}$$

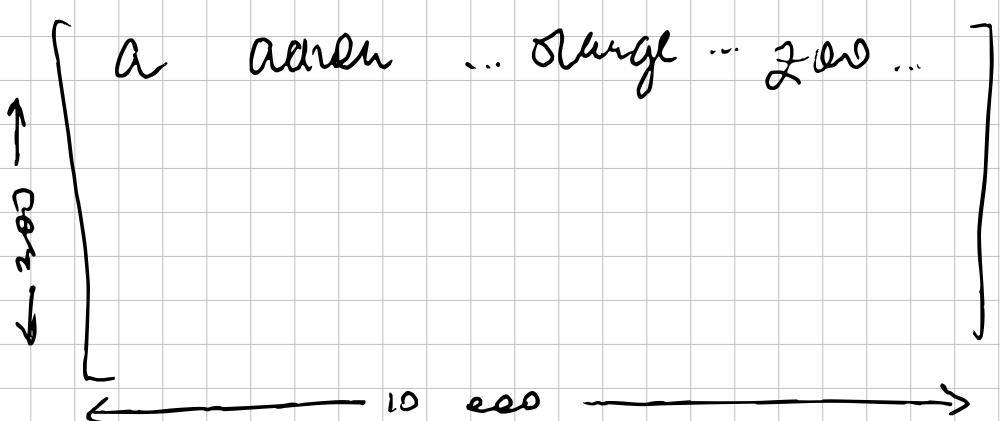
* Mikolov

Find word w :

$$\arg \max_w \text{sim}(\ell_w, \ell_{\text{king}} - \ell_{\text{man}} + \ell_{\text{woman}})$$

Usually cosine similarity is used to calculate similarity.

Embedding matrix



E - embedding matrix

Learning word embeddings

I

$O_{4343} \rightarrow \bar{E} \rightarrow e_{4343}$

want

$O_{9665} \rightarrow \bar{E} \rightarrow e_{9665}$

a

$O_1 \rightarrow E \rightarrow e_1$

glass

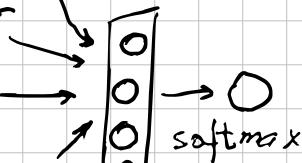
$O_{3852} \rightarrow E \rightarrow e_{3852}$

of

$O_{6163} \rightarrow \bar{E} \rightarrow e_{6163}$

orange

$O_{6257} \rightarrow \bar{E} \rightarrow e_{6257}$



Skip gram model

Context : . Last 4 words

- 4 words on left & right
- Last 1 word
- Nearby word

Word 2 Vec

I want a glass of orange juice to go along with my cereal.

Context

(some random word)

orange → juice

orange → glass

orange → my

Target

Model

Vocal size = 10,000

Context c ("orange") → Target t ("juice")

$$x \longrightarrow y$$

θ_c (one-hot vector for context)

$$\text{Soft max: } p(t|c) = \frac{e^{\theta_t^\top \theta_c}}{\sum_{j=1}^{10000} e^{\theta_j^\top \theta_c}}$$

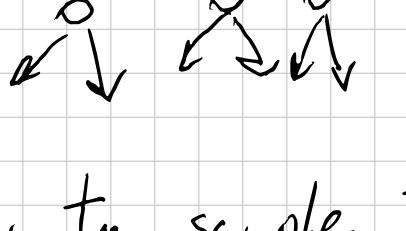
θ_t = parameter associated with output t

$$L(\hat{y}, y) = -\sum_{i=1}^{vocab} y_i \log \hat{y}_i$$

y - one-hot encoding vector

(* skip-gram model)

Can be really slow. Solution: use Hierarchical softmax



How to sample the context c ?

* Use less frequent/common words.

Negative sampling

* Softmax is really "heavy"/expensive to calculate

Given a pair of words

orange → juice

Is it a pair of context - target?

orange - juice

1

orange - king

0

orange - book

0

orange - the

0

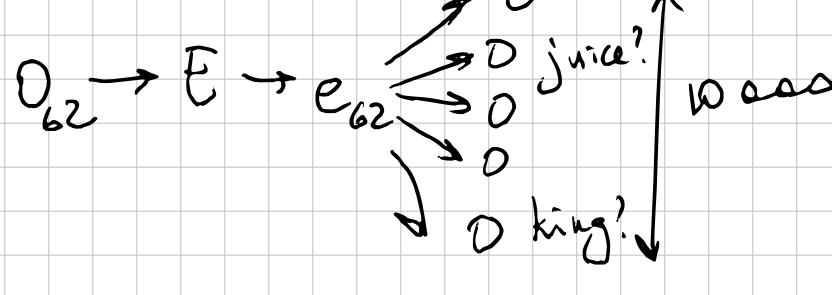
orange - of

0

$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{vocab} e^{\theta_j^T e_c}}$$

} pick second word rand. from the vocabulary

$$P(y=1 | c, t) = G(\theta_t^T e_c)$$



10,000 binary classification problem which is much faster

How to choose negative examples?

$$P(w_i) = \frac{f(w_i)}{\sum_{j=1}^{vocab} f(w_j)}^{3/4}$$

GloVe word vector

Global Vectors for word representation

c,t

x_{ij} = # times j appears in context +
of i ↑
 ↑ t
c

minimize the difference

$$\sum_{i=1}^m \sum_{j=1}^n f(x_{ij}) \underbrace{\left(\theta_i^T e_j + b_i + b_j' - \log x_{ij} \right)^2}_{\text{weight term}} \quad t \quad c$$

Sequence to sequence model

$x^{<1>} \quad x^{<2>}$

Jane visite l'Afrique en
septembre

$y^{<1>} \quad y^{<2>} \quad y^{<3>}$

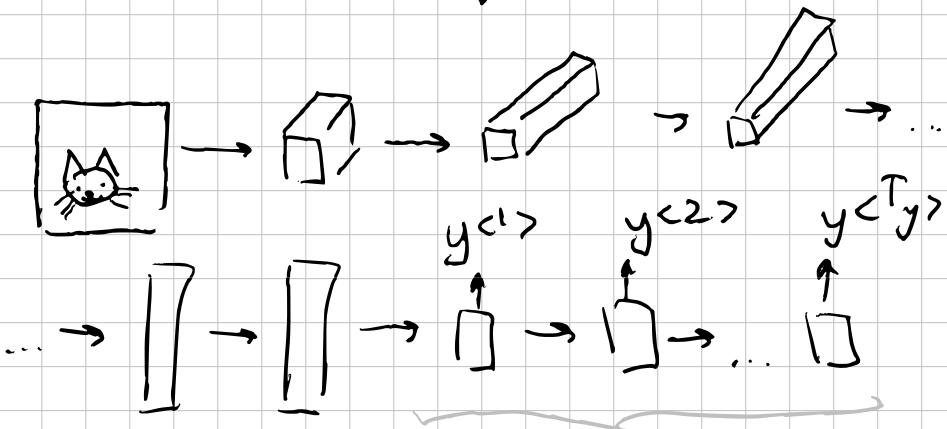
Jane is visiting Africa in
September

$a^{<0>} \rightarrow \boxed{} \rightarrow \dots \boxed{} \rightarrow \boxed{} \rightarrow \dots \rightarrow \boxed{} \rightarrow$

$\uparrow \quad \uparrow \quad \uparrow$
 $x^{<1>} \quad x^{<T_x>}$

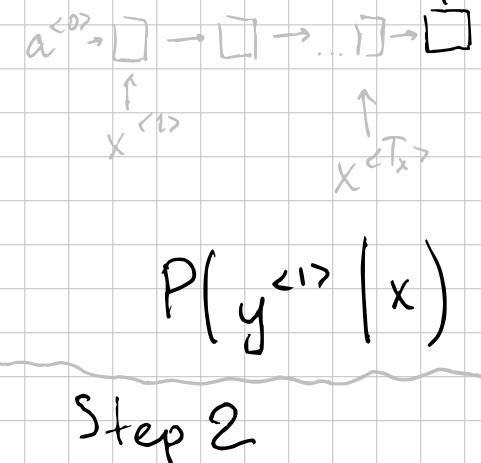
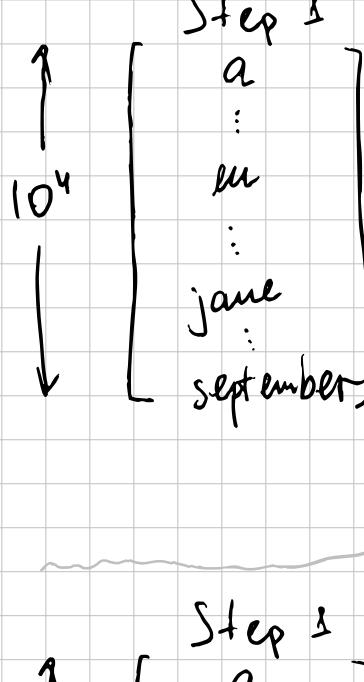
$y^{<1>} \quad y^{<T_y>}$
 $\uparrow \quad \uparrow$

Image captioning

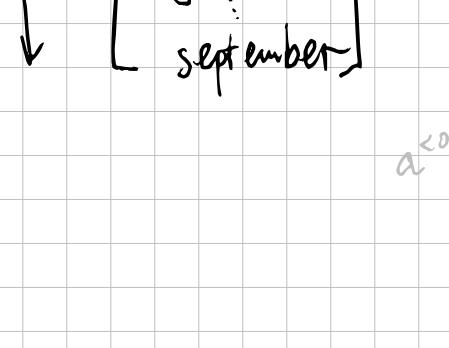


RNN

Beam Search



$$P(y^{<1>} | x)$$



Step 2

aaron
september



$\hat{y}^{<1>} \hat{y}^{<2>}$

$$P(y^{<1>}, y^{<2>} | x) = P(y^{<1>} | x) \cdot$$

$$\cdot P(y^{<2>} | x, y^{<1>})$$

Refinements to beam search

Length normalization

$$\arg \max_{y} \prod_{t=1}^{T_y} P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$$

numbers less than 1

multiplication of a plenty of small numbers result in underflow

$$\arg \max_{y} \sum_{y=1}^{T_y} \log P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$$

Usually, $\alpha = 0.7$

$$\alpha = 1$$

Tips

Beam width B ?

large B : consider more possib., better results,

slow

smaller B : worse result, faster

$$B \in [10, 100, 1000, 3000]$$

Unlike BFS/DFS, Beam Search runs faster but not guaranteed to find exact minimum.

Error Analysis in Beam search

Jane visite l'Afrique en septembre.

Human: Jane visits Africa in September

Algo: Jane visited Africa last September

RNN computes $P(y|x)$

lets compute $P(y^*|x)$ and

$P(\hat{y}|x)$

Case 1: $P(y^*|x) > P(\hat{y}|x)$

Beam search chose \hat{y} . But y^* attains higher $P(y|x)$.

Conclusion: Beam search is at fault

Case 2: $P(y^*|x) \leq P(\hat{y}|x)$

y^* is a better translation than \hat{y} . But RNN predicted $P(y^*|x) \leq P(\hat{y}|x)$

Conclusion: RNN model is at fault.

Bleu Score

How to evaluate machine translation if there are multiple equal translations?

Bleu - bilingual evaluation understudy

Le chat est sur le tapis.

Ref. 1: The cat is on the mat.

Ref. 2: There is a cat on the mat.

MT output: the the the the the.

Precision: $\frac{2}{7}$

Modified precision: $\frac{2}{7} \leq \frac{\text{Count clip}(\text{"the"})}{\text{Count clip}}$

* MT output: The cat the cat on the mat.

$\frac{\text{Count}}{\text{Count clip}}$

the cat 2 1

cat the 1 0

cat on 1 1

on the 1 1

the mat 1 1

Bigrams precision: $\frac{\sum \text{Count clip}}{\sum \text{Count}} = \frac{4}{6}$

Bleu score on unigrams

$P_1 = \frac{\sum_{\text{unigrams}} \text{Count clip}(\text{unigram})}{\sum_{\text{unigram}} \text{Count}(\text{unigram})}$

$\sum_{\text{unigram}} \text{Count}(\text{unigram})$

$P_n = \frac{\sum_{n\text{-grams}} \text{Count clip}(n\text{-gram})}{\sum_{n\text{-grams}} \text{Count}(n\text{-grams})}$

Bleu details

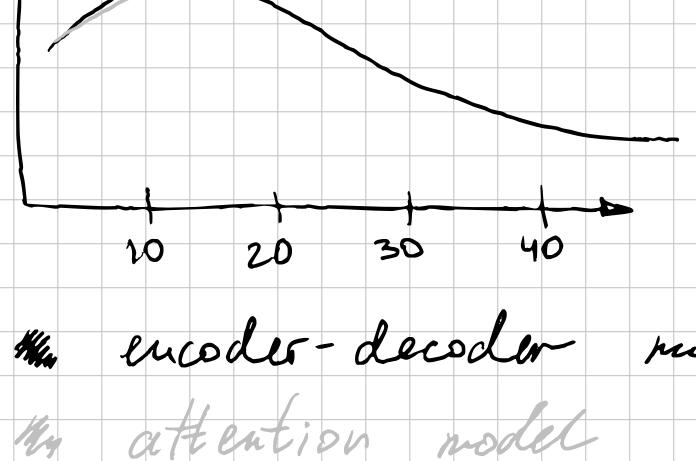
p_n - bleu score on n-grams only

Combined bleu score:

$$BP \exp \left(\frac{1}{4} \sum_{n=1}^4 p_n \right)$$

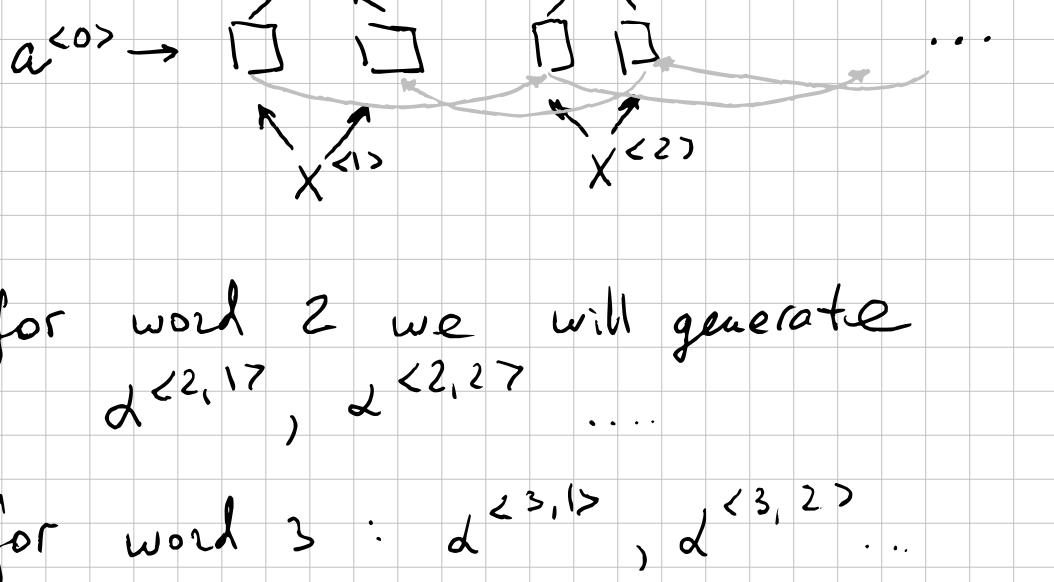
Attention model

Encoder-decoder architecture



encoder-decoder model
attention model

Attention model processes the input by parts.

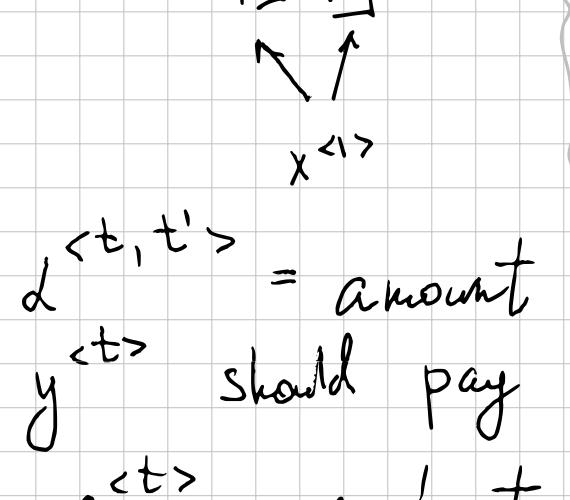


for word 2 we will generate
 $d^{<2,1>} , d^{<2,2>} \dots$

for word 3 : $d^{<3,1>} , d^{<3,2>} \dots$

$$a^{<t>} = (\vec{a}^{<t>} , \overset{\leftarrow}{a}^{<t>})$$

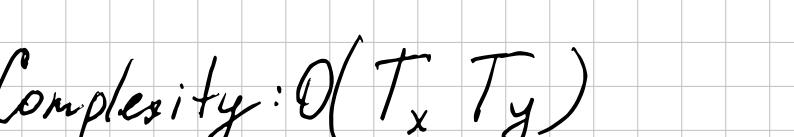
forward activations backward activations



$d^{<t,t'>} = \text{amount of "attention"}$
 $y^{<t>} \text{ should pay to } a^{<t'>}$

$c^{<t>} = \text{context vector}$

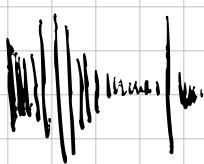
$$d^{<t,t'>} = \frac{\exp(e^{<t,t'>})}{\sum_{t'=1}^{T_x} \exp(e^{<t,t'>})}$$



Complexity: $O(T_x T_y)$

Speech recognition

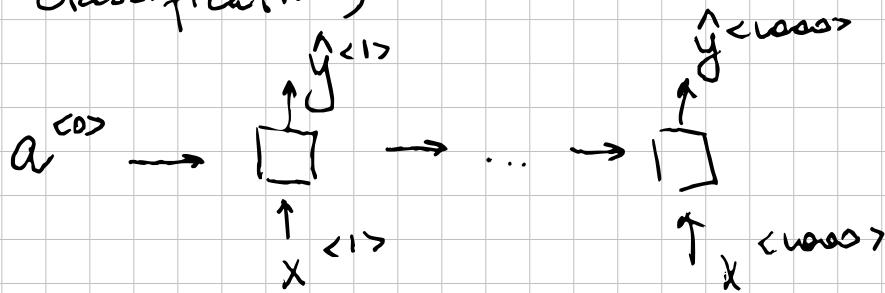
$X \rightarrow y$
 audio clip transcript



"the quick brown
fox"

CTC cost for speech recognition

(connectionist temporal
classification)



CTC allows to generate
special output

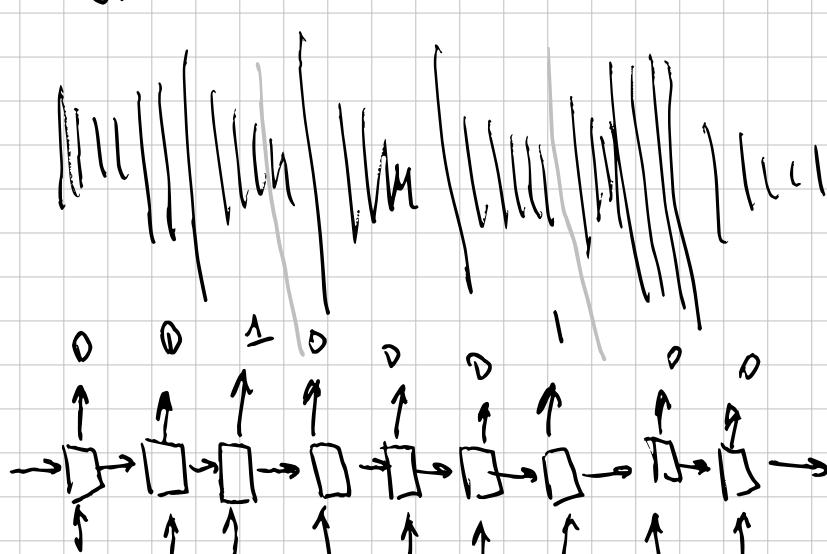
ttt_h_eee_l_g_

"blank"
"space"

collapsing repeating characters
not separated by blank.

the \rightarrow g

Trigger Word Detection



Transformers

- Attention + CNN
 - Self-attention
 - Multi-Head Attention

Self-Attention

$A(q, k, v)$ = attention-based vector representation of a word



Calculate for each word

$$A^{(1)}, A^{(2)}, \dots, A^{(n)}$$

Jane visite l'Afrique en septembre

Transformers Attention

$$A(q, k, v) = \sum_i \frac{\exp(q \cdot k^{(i)})}{\sum_j \exp(q \cdot k^{(j)})} v^{(i)}$$

$x^{(3)}$ 1. Associate each of the words with Query, key, Value.
l'Afrique $(q^{(3)}, k^{(3)}, v^{(3)})$

$$q^{(3)} = W^Q x^{(3)}$$

$$k^{(3)} = W^K x^{(3)}$$

$$v^{(3)} = W^V x^{(3)}$$

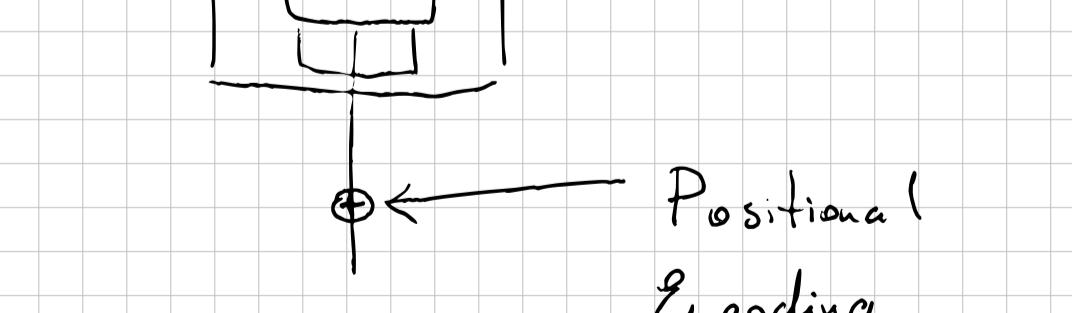
$q^{(3)}$ is a question you ask about the subject

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

↑ factor to scale dot-product so it won't get explode

Multi-Head Attention

Multi-Head(Q, K, V)



Instead of W^Q, W^K, W^V we will use W_i^Q, W_i^K, W_i^V to be able to address multiple

questions.

$$\text{Multi-Head} = \text{concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_n)W_0$$

$$\text{head}_i = \text{Attention}(W_i^Q Q, W_i^K K, W_i^V V)$$

Transformer network

1. Embeddings get fit into encoder

Positional encoding \oplus

$$\text{PE}(\text{pos}, 2t) = \sin\left(\frac{\text{pos}}{10000^{\frac{2t}{d}}}\right)$$

$$\text{PE}(\text{pos}, 2t+1) = \cos\left(\frac{\text{pos}}{10000^{\frac{2t+1}{d}}}\right)$$

Add & Norm

Similar to Batch Norm

Masked Multi-Head Attention

Helps during training: blocks out parts of the sentence

Long Short Term Memory

(LSTM)

$$\tilde{C}^{<t>} = \tanh(W_C[a^{<t-1>}, x^{<t>}] + b_C)$$

$$\Gamma_u = (W_u[a^{<t-1>}, x^{<t>}] + b_u)$$

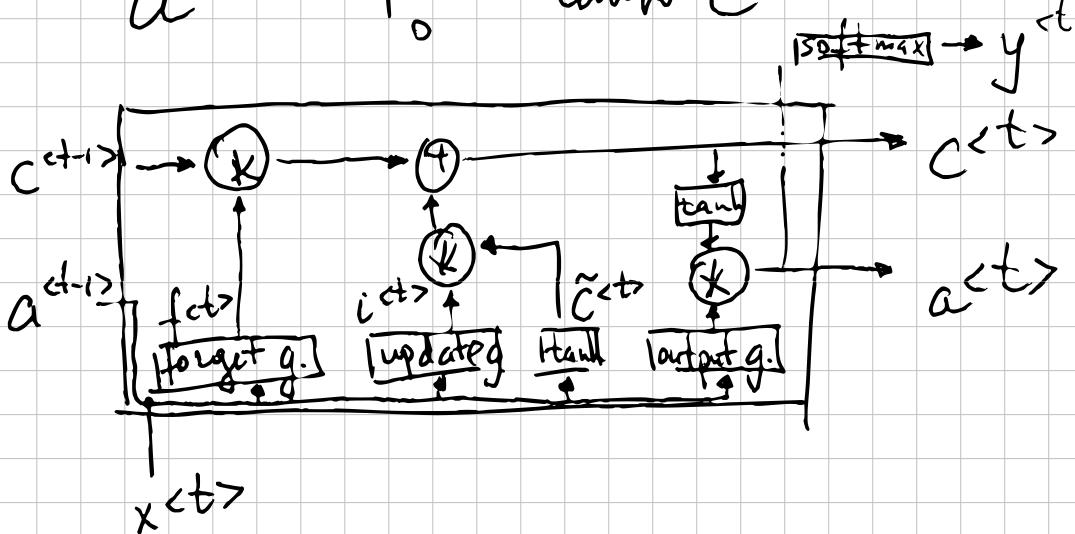
$$F_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$$

"forget" gate

$$\Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o)$$

$$C^{<t>} = \Gamma_u * \tilde{C}^{<t>} + \Gamma_f * C^{<t-1>}$$

$$a^{<t>} = \Gamma_o * \tanh C^{<t>}$$

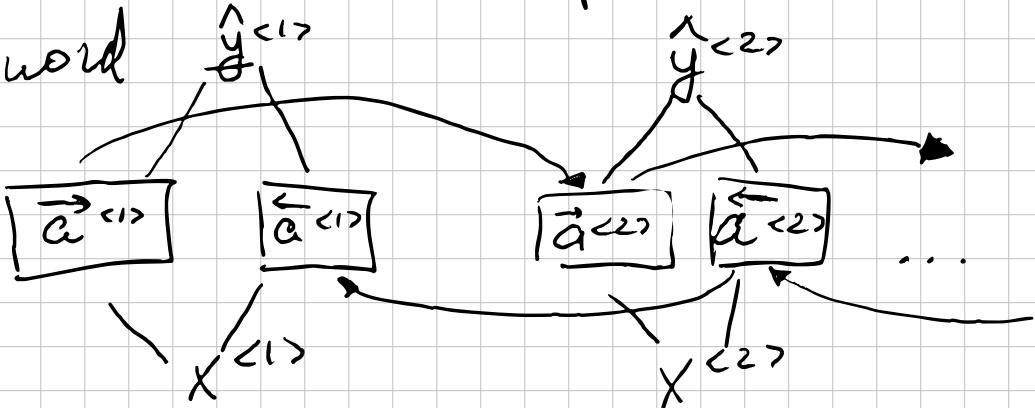


Bidirectional RNN

• He said, "Teddy Roosevelt was a great president."

• He said, "Teddy bears are on sale"

It is not enough to take a look just at the word (without looking ahead) to determine the semantic of the word

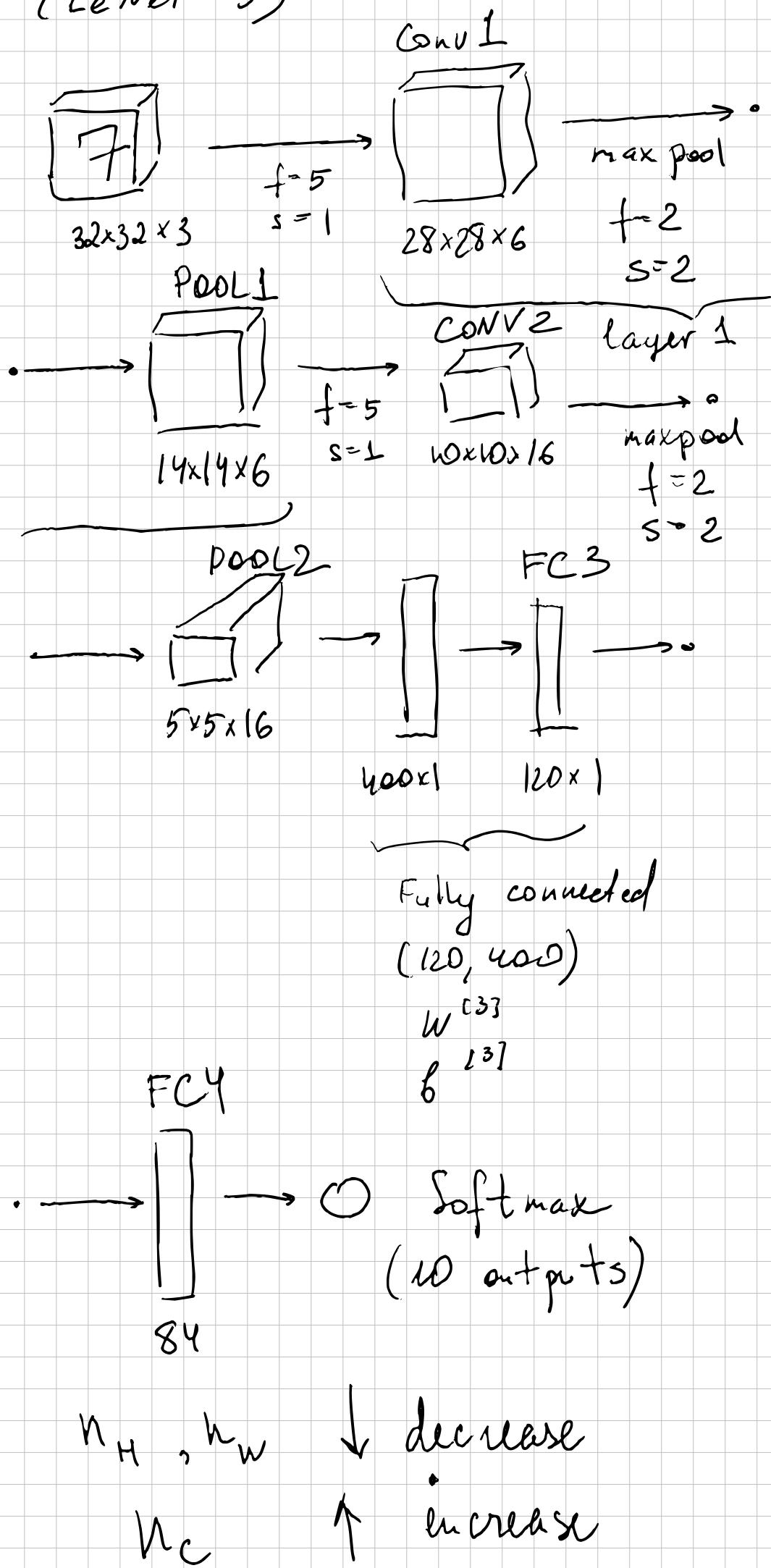


$$\hat{y}^{<t>} = g(w_y [\bar{a}^{<t>} \leftarrow \bar{a}^{<t>}] + b_y)$$

stacylic graph

Neural network example

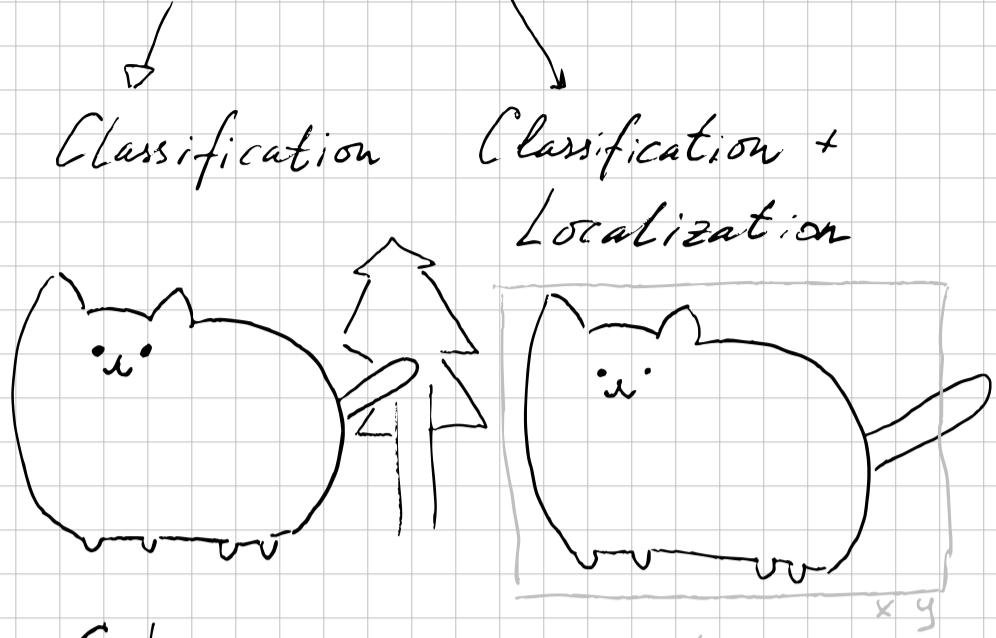
(LeNet - 5)



Why convolutions?

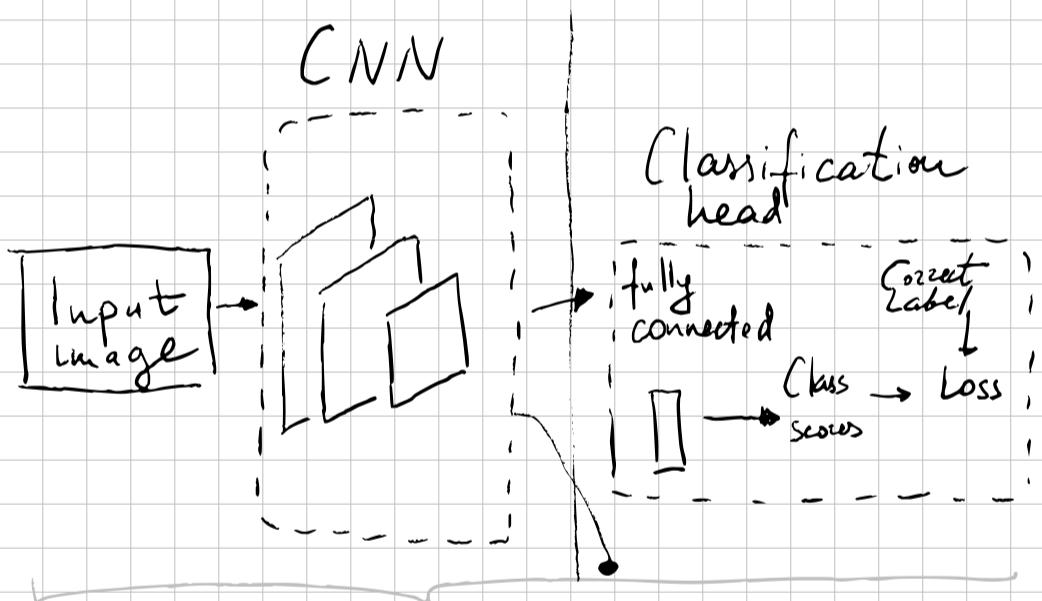
- FC have a lot of parameters, even be unfeasibly large
- Parameter sharing: a feature detector (such as vertical edge detector) that's useful in one part of the image is probably useful in another part of the image.
- Sparsity of connections: In each layer, each output value depends only on a small number of inputs.

Object detection task



What and where
the single object
exist in an image

Basic classification



Classification

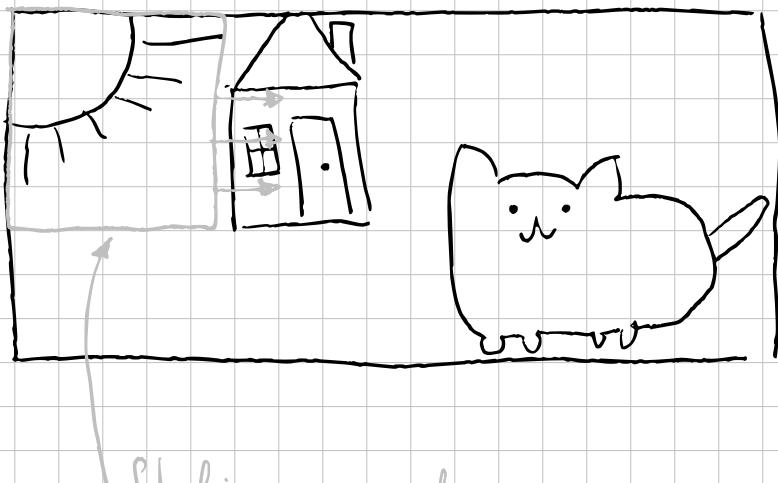
For localization we need to add extra
step



⊕ localization

Sliding window approach

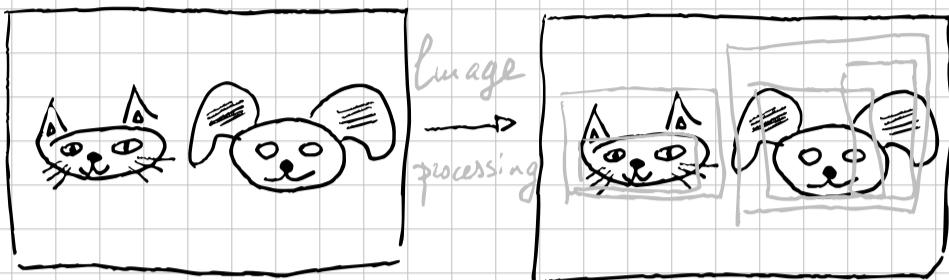
We take an area of the image of a fixed size and "slide" it over the image testing every position independently



Computationally expensive to cover all possible cases (locations, scales, aspect ratios)

Region-based approaches

- Find some image regions where object can be potentially found

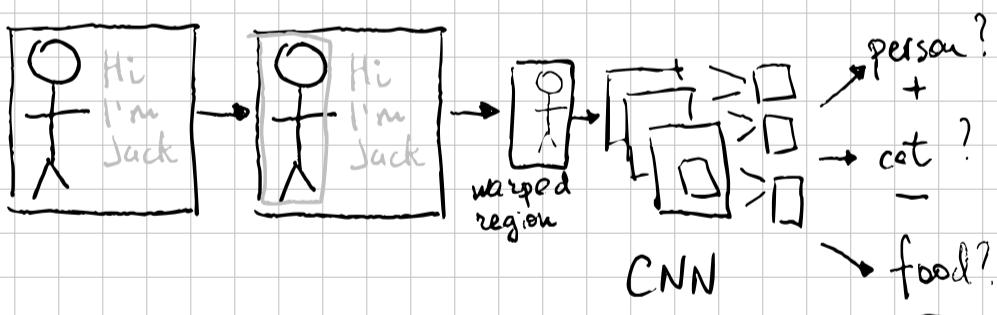


- Relatively fast

R-CNN

Pipeline

- ① Generate region proposals with selective search
- ② Resize each region proposal
- ③ Predict class probabilities with linear SVM (including background class)
- ④ Predict bounding box refinements

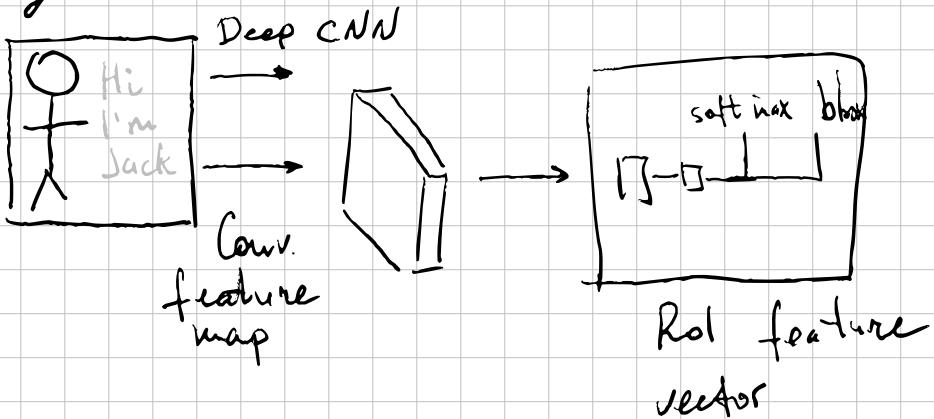


Downsides ↓ :

- Inference is quite slow ($\sim 47 \frac{\text{sec}}{\text{image}}$)
- Region proposal algorithm is fixed and not trainable
- Complicated and slow training mechanism requiring a lot of space

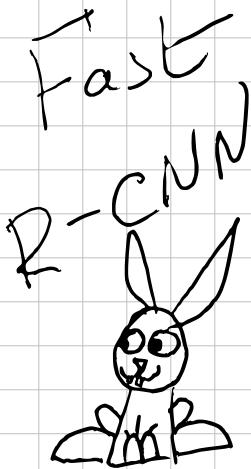
Fast R-CNN

Instead of feeding region proposals we feed the whole image to the CNN

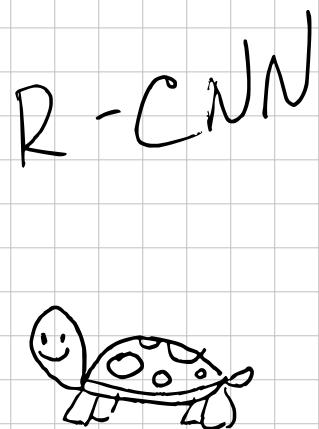


Pipeline

- ① Extract feature map through CNN
- ② Generate region proposals on feature map using selective search
- ③ ROI pooling
- ④ Predict class probabilities with softmax
- ⑤ Predict bounding box refinements



VS.



With
selective
search

2.3 $\frac{\text{sec}}{\text{image}}$

49 $\frac{\text{sec}}{\text{image}}$

w/o
selective
search

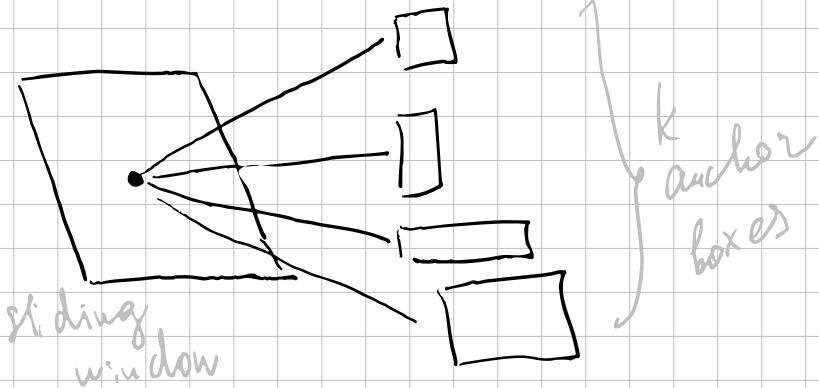
0.32 $\frac{\text{sec}}{\text{image}}$

47 $\frac{\text{sec}}{\text{image}}$

Faster R - CNN

Pipeline

- ① Extract feature map using CNN
- ② Generate region proposals on feature map using RPN
- ③ ROI pooling
- ④ Predict class probabilities with Softmax
- ⑤ Predict bounding box refinements

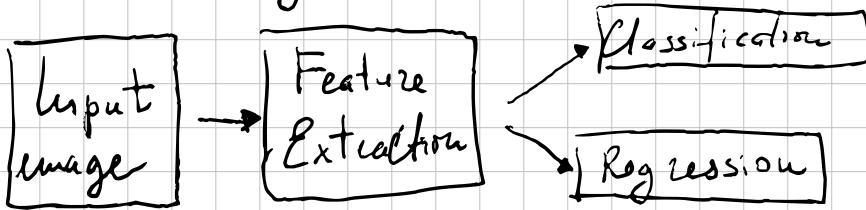


Single - Stage detectors

- 2 - Stage Detector



- 1-Stage Detector



We go directly
from image to
detection

YOLO

Each grid cell predicts : B
bounding boxes and C classification scores

Each bounding box consists of :

- Center coordinates (x, y)
- Width and height (w, h)
- Confidence score

Final prediction : $S \times S \times (5B + C)$,
where

S - grid size

B - number of base boxes

C - number of categories

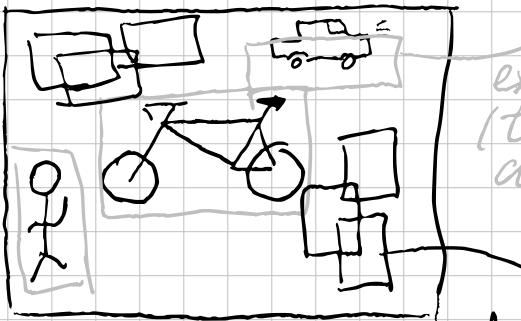
SSD (single short multibox
detector)

We apply CNN in this model
and reduce feature space.

- Each grid cell predicts: K bounding boxes
- Each bounding box consists of $C+4$ numbers:
 - 4 offsets relative to the original default box shape
 - C classification scores
- Final prediction for feature map $M \times N \times K \times (C+4)$
 - $M \times N$ - grid
 - K - Number of base boxes
 - C - number of categories (+ background)

Retina Net

Problem that led to model invention : Foreground - background imbalance



Small amount
of positive
examples
(target
classes)

Large number
of negative
examples
(background)

The architecture of Retina Net consists of four main parts:

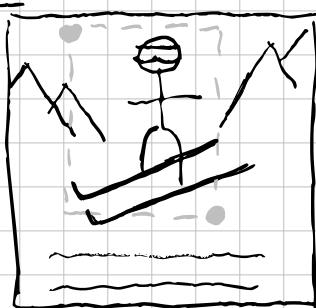
- Backbone
- FPN
- Classification Head
- Regression Head

Single Stage Detectors

anchor-free

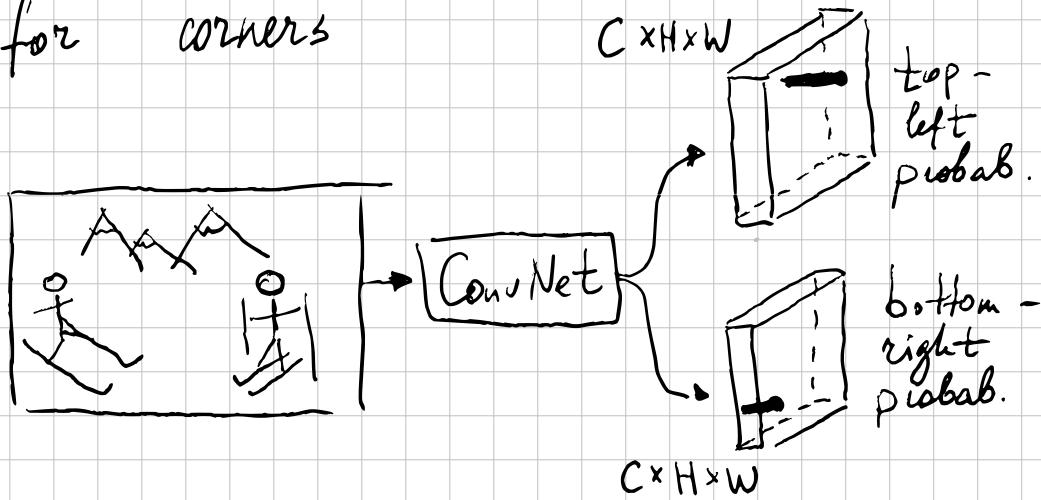
CornerNet

Represents an object as a pair of keypoints – the top-left and the bottom-right corners



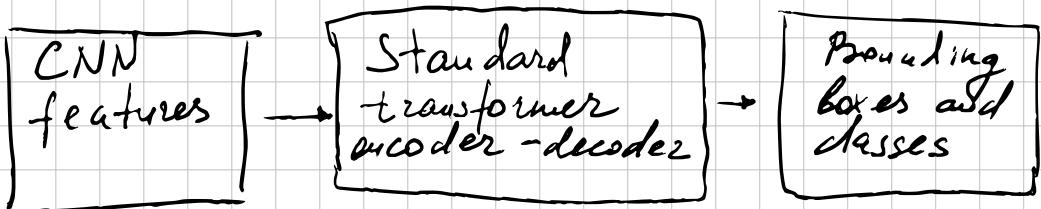
Approach

1. Extract feature map through CNN
2. Predict two sets of heatmaps for corners



3. Predict an embedding vector for each detected corner
4. Predict offsets to slightly adjust the locations of the corners

DETR



Approach :

- Extract feature map through CNN
 - Transform feature map into sequence
 - Pass sequence to encoder
 - Pass to decoder
 - Predict coordinates and class
- through FFN

Non-Maximum Suppression

NMS is an algorithm for prediction filtering. It allows selecting the best bounding box from several predictions.

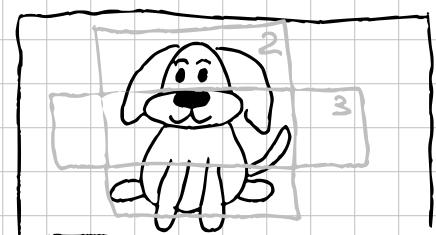
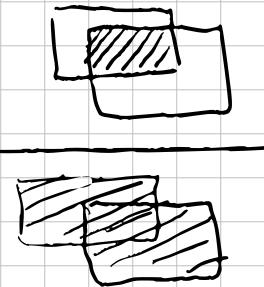


Boxes Scores

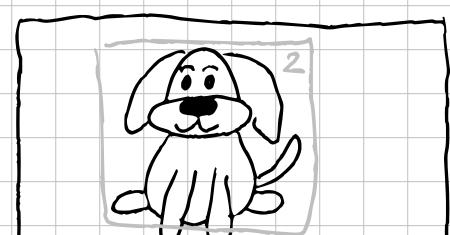
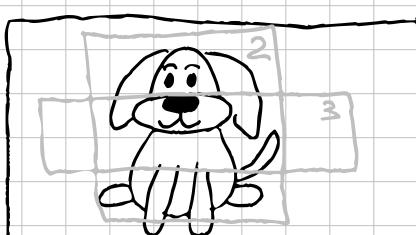
1	0.6
2	0.9
3	0.55

To find the best bounding box we consider using intersection over union (iou).

$$IOU = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$



$$IOU(1,2) > \alpha$$



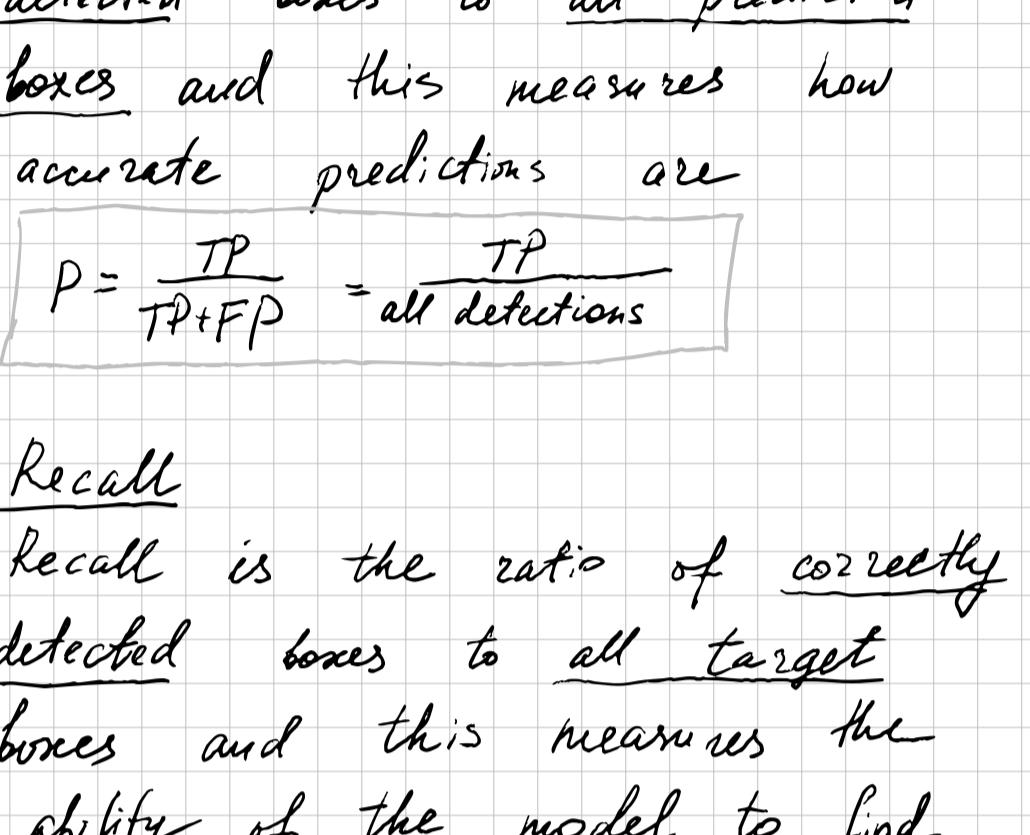
$$IOU(2,3) > \alpha$$

If IOU is bigger than pre-defined parameter α we simply remove the box

Evaluating

Let's define:

- True Positive (TP) - a valid detection ($10U \geq \alpha$)
- False Positive (FP) - an invalid detection ($10U < \alpha$)
- False Negative (FN) - Ground-truth missed by the model at $\alpha = 10U$ threshold



Ground-truth
Predicted

Precision

Precision is the ratio of correctly detected boxes to all predicted boxes and this measures how accurate predictions are

$$P = \frac{TP}{TP+FP} = \frac{TP}{\text{all detections}}$$

Recall

Recall is the ratio of correctly detected boxes to all target boxes and this measures the ability of the model to find all relevant objects

$$R = \frac{TP}{TP+FN} = \frac{TP}{\text{all ground-truths}}$$

Mean Average Precision

1. We need to run our detector on all images from the test set

2. For every image: sort all found objects based on classification score (the most confident box first)

3. For every predicted box we are going to assign a target box, based on $10U$

Example

Ground-truth boxes

Predicted boxes

Score	Matched with GT	TP/FP	Precision	Recall
0.99	+	TP	1	1/3
0.9	-	FP	1/2	1/3
0.8	+	TP	2/3	2/3
0.5	+	TP	3/4	1
0.3	-	FP	3/5	1

The AP is calculated individually for each class. There are as many AP values as there are classes. These AP values are averaged to produce Mean Average Precision (mAP) metric.

$$mAP = \frac{1}{n} \sum_{i=1}^n AP_i, \text{ for } n \text{ classes}$$

AP@0.5 means Average Precision at $10U$ threshold of 0.5.

AP@0.50 means AP at $10U$ threshold of 50%.

AP@0.25 means Average Precision at $10U$ threshold of 0.25.

AP@0.1 means Average Precision at $10U$ threshold of 0.1.

AP@0.05 means Average Precision at $10U$ threshold of 0.05.

AP@0.01 means Average Precision at $10U$ threshold of 0.01.

AP@0.005 means Average Precision at $10U$ threshold of 0.005.

AP@0.001 means Average Precision at $10U$ threshold of 0.001.

AP@0.0005 means Average Precision at $10U$ threshold of 0.0005.

AP@0.0001 means Average Precision at $10U$ threshold of 0.0001.

AP@0.00005 means Average Precision at $10U$ threshold of 0.00005.

AP@0.00001 means Average Precision at $10U$ threshold of 0.00001.

AP@0.000005 means Average Precision at $10U$ threshold of 0.000005.

AP@0.000001 means Average Precision at $10U$ threshold of 0.000001.

AP@0.0000005 means Average Precision at $10U$ threshold of 0.0000005.

AP@0.0000001 means Average Precision at $10U$ threshold of 0.0000001.

AP@0.00000005 means Average Precision at $10U$ threshold of 0.00000005.

AP@0.00000001 means Average Precision at $10U$ threshold of 0.00000001.

AP@0.000000005 means Average Precision at $10U$ threshold of 0.000000005.

AP@0.000000001 means Average Precision at $10U$ threshold of 0.000000001.

AP@0.0000000005 means Average Precision at $10U$ threshold of 0.0000000005.

AP@0.0000000001 means Average Precision at $10U$ threshold of 0.0000000001.

AP@0.00000000005 means Average Precision at $10U$ threshold of 0.00000000005.

AP@0.00000000001 means Average Precision at $10U$ threshold of 0.00000000001.

AP@0.000000000005 means Average Precision at $10U$ threshold of 0.000000000005.

AP@0.000000000001 means Average Precision at $10U$ threshold of 0.000000000001.

AP@0.0000000000005 means Average Precision at $10U$ threshold of 0.0000000000005.

AP@0.0000000000001 means Average Precision at $10U$ threshold of 0.0000000000001.

AP@0.00000000000005 means Average Precision at $10U$ threshold of 0.00000000000005.

AP@0.00000000000001 means Average Precision at $10U$ threshold of 0.00000000000001.

AP@0.000000000000005 means Average Precision at $10U$ threshold of 0.000000000000005.

AP@0.000000000000001 means Average Precision at $10U$ threshold of 0.000000000000001.

AP@0.0000000000000005 means Average Precision at $10U$ threshold of 0.0000000000000005.

AP@0.0000000000000001 means Average Precision at $10U$ threshold of 0.0000000000000001.

AP@0.00000000000000005 means Average Precision at $10U$ threshold of 0.00000000000000005.

AP@0.0000000000000001 means Average Precision at $10U$ threshold of 0.0000000000000001.

AP@0.00000000000000005 means Average Precision at $10U$ threshold of 0.00000000000000005.

AP@0.0000000000000001 means Average Precision at $10U$ threshold of 0.0000000000000001.

AP@0.00000000000000005 means Average Precision at $10U$ threshold of 0.00000000000000005.

AP@0.0000000000000001 means Average Precision at $10U$ threshold of 0.0000000000000001.

AP@0.00000000000000005 means Average Precision at $10U$ threshold of 0.00000000000000005.

AP@0.0000000000000001 means Average Precision at $10U$ threshold of 0.0000000000000001.

AP@0.00000000000000005 means Average Precision at $10U$ threshold of 0.00000000000000005.

AP@0.0000000000000001 means Average Precision at $10U$ threshold of 0.0000000000000001.

AP@0.00000000000000005 means Average Precision at $10U$ threshold of 0.00000000000000005.

AP@0.0000000000000001 means Average Precision at $10U$ threshold of 0.0000000000000001.

AP@0.00000000000000005 means Average Precision at $10U$ threshold of 0.00000000000000005.

AP@0.0000000000000001 means Average Precision at $10U$ threshold of 0.0000000000000001.

AP@0.00000000000000005 means Average Precision at $10U$ threshold of 0.00000000000000005.

AP@0.0000000000000001 means Average Precision at $10U$ threshold of 0.0000000000000001.

AP@0.00000000000000005 means Average Precision at $10U$ threshold of 0.00000000000000005.

AP@0.0000000000000001 means Average Precision at $10U$ threshold of 0.0000000000000001.

AP@0.00000000000000005 means Average Precision at $10U$ threshold of 0.00000000000000005.

AP@0.0000000000000001 means Average Precision at $10U$ threshold of 0.0000000000000001.

AP@0.00000000000000005 means Average Precision at $10U$ threshold of 0.00000000000000005.

AP@0.0000000000000001 means Average Precision at $10U$ threshold of 0.0000000000000001.

AP@0.00000000000000005 means Average Precision at $10U$ threshold of 0.00000000000000005.

AP@0.0000000000000001 means Average Precision at $10U$ threshold of 0.0000000000000001.

AP@0.00000000000000005 means Average Precision at $10U$ threshold of 0.00000000000000005.

AP@0.0000000000000001 means Average Precision at $10U$ threshold of 0.0000000000000001.

AP@0.00000000000000005 means Average Precision at $10U$ threshold of 0.00000000000000005.

AP@0.0000000000000001 means Average Precision at $10U$ threshold of 0.0000000000000001.

AP@0.00000000000000005 means Average Precision at $10U$ threshold of 0.00000000000000005.

AP@0.0000000000000001 means Average Precision at $10U$ threshold of 0.0000000000000001.

AP@0.00000000000000005 means Average Precision at $10U$ threshold of 0.00000000000000005.

AP@0.0000000000000001 means Average Precision at $10U$ threshold of 0.0000000000000001.

AP@0.00000000000000005 means Average Precision at $10U$ threshold of 0.00000000000000005.

AP@0.0000000000000001 means Average Precision at $10U$ threshold of 0.0000000000000001.

AP@0.00000000000000005 means Average Precision at $10U$ threshold of 0.00000000000000005.

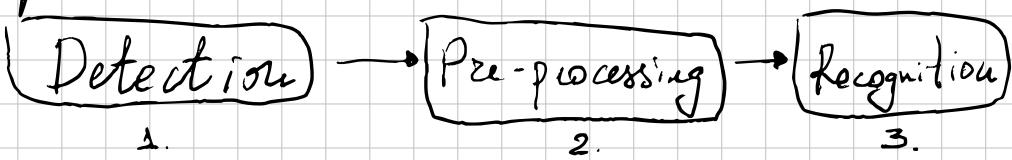
AP@0.0000000000000001 means Average Precision at $10U$ threshold of 0.0000000000000001.

AP@0.00000000000000005 means Average Precision at $10U$ threshold of 0.00000000000000005.

AP@0.0000000000000001 means Average Precision at $10U$ threshold of 0.0000000000000001.

Person Re-identification

Pipeline



The most popular face detector is by now:

- Haar Cascades



Haar-like features

- Histogram of Oriented Gradients



- HTCNN (the first successfull neural net detector)

Stages

1. Resize (using Image pyramid)

2. P-Net

3. R-Net

4. O-Net

- RetinaFace

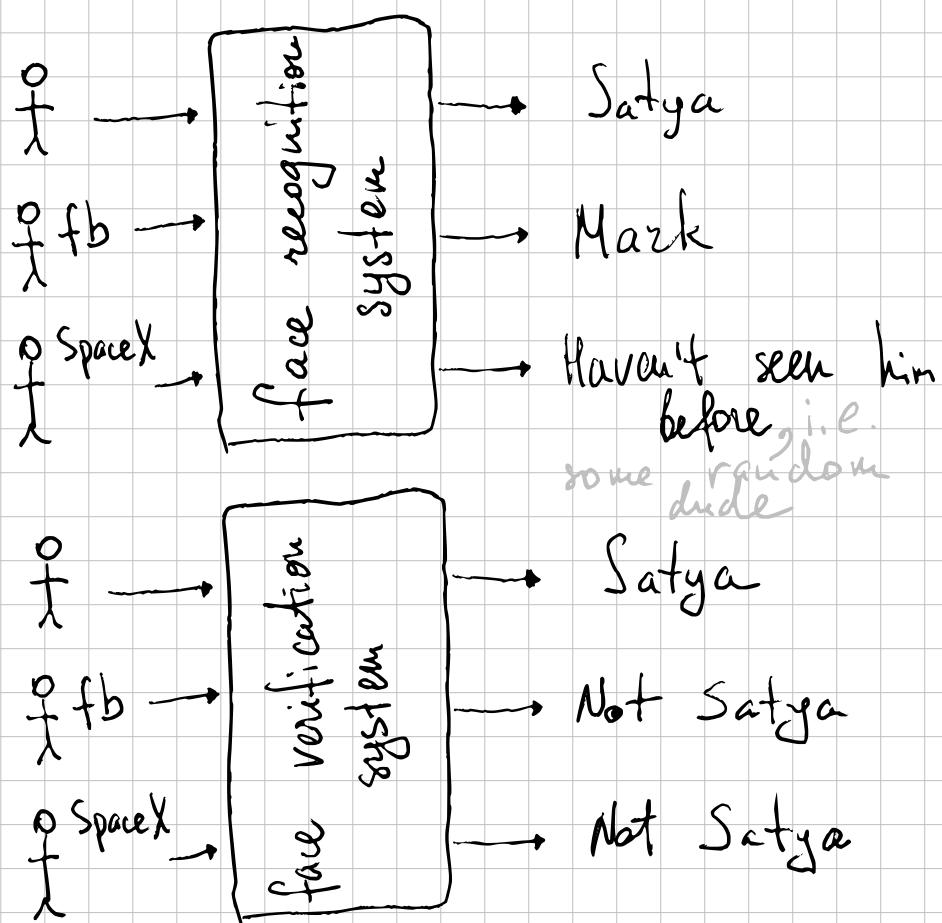
- SCRFD



Recognition vs. Verification

In general face recognition task consists of 2 subtasks:

- Identity recognition
- Identity verification (authentication)



Metrics

Verification

- Recall
- Precision

Recognition

- False Rejection Rate

$$\frac{\# \text{ false rejected}}{\# \text{ total attempts}}$$

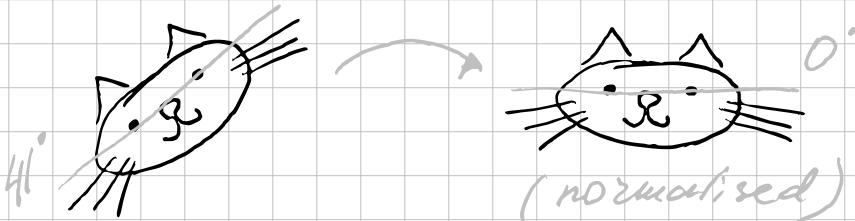
- False Acceptance Rate

$$\frac{\# \text{ false accepted}}{\# \text{ total attempts}}$$

Face Normalization

It was extremely important for faces to be uniformed.

For example, eyes should be on the same vertical level



Popular algorithms

- Eigen Faces
- Fisher Faces
- LBP

Classical Face Recognition

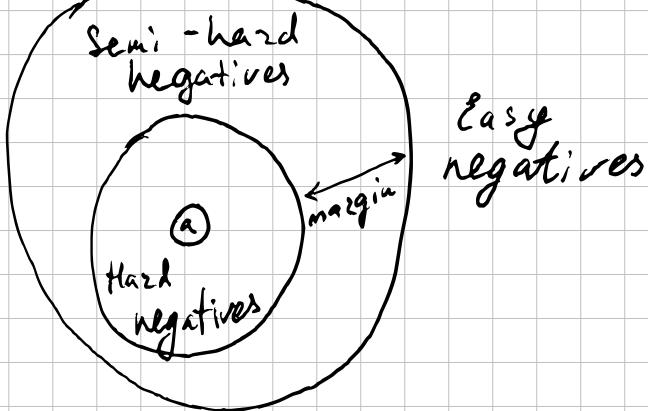
Some classical operator
applied to the image
(usually in a sliding window
manner)

Triplet Loss

For every minibatch (for every identity) :

- take equal number of identity photos
- take equal number of random photos from another identities
 - for every pair of anchor-positive let's find the hardest negative

$$\|f(x_i^a) - f(x_i^p)\|_2^2 < \|f(x_i^a) - f(x_i^n)\|_2^2$$



Object tracking

Problems

- ① Ambiguous metrics
- ② Object movement poorly predictable
- ③ Fine-tuning is costly

Tracker Types

- Online / Offline
- Single-Object (SOT) / Multi-Object (MOI)
- Neural / Classical
- Static camera / Moving camera
- Predictive / Stateless
- Single-modal / Multi-modal

Tracking features

- Coordinates of object
- Kinematics
- Physiological characteristic
- Appearance features
- Extra features

MOT

1. Starts with object detection



2. Feature Generation
3. State prediction (predicts the state of a new frame)
4. Assignment (Matching)
5. Management

Algorithms

① Sort

- Find objects
- State prediction (using Kalman filter)
- Matching (using Hungarian algorithm)
- Management

② Deep SORT

- * add cosine proximity between object last appearance and the new detection

Deep Learning for Image Segmentation

Segmentation types

1. Semantic segmentation

- All pixels have labels
- Single class not divide on objects

objects

2. Instance segmentation

- Some pixels does not have labels
- Single class divide on objects

3. Panoptic segmentation

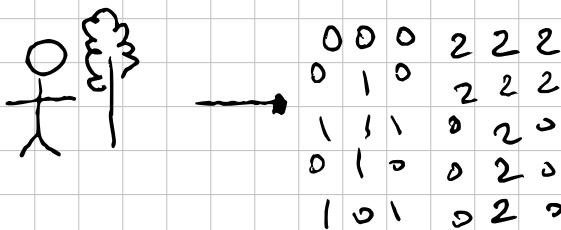
- Combines semantic and instance

Segmentation task

Send image in model and get
mask image with same
resolution

1. We need to encode labels

Encoding - transforms rgb label
mask $H \times W \times C$ using one-hot
encoding to $H \times W \times 1$



0 - ground
1 - person
2 - plant

Loss and Metrics

① Pixel accuracy

$$acc = \frac{TP + TN}{TP + TN + FP + FN}$$

This metric can provide misleading results when the class representation is small within the image, as the measure will be biased in mainly reporting how well you identify negative case (class imbalance problem)

② Dice Coefficient

$$Dice = \frac{2TP}{2TP + FP + FN} \quad \left(Dice = \frac{2|A \cap B|}{|A| + |B|} \right)$$

Dice Coefficient is 2^* the area of overlap divided by the sum of ground truth and predicted mask

③ Intersection - Over - Union

(IoU, Jaccard Index)

IoU is the area of overlap between the predicted segmentation and the ground truth divided by the area of union between the predicted segmentation

$$Jac = \frac{TP}{TP + FP + FN}$$

$$Jac = \frac{|A \cap B|}{|A \cup B|}$$

Loss

1. Distribution - based loss

- Focal loss
- Top-K
- Cross-Entropy

2. Compound loss (recombination of other 2 losses)

- ELL
- Dice CE

3. Region - based loss

- Dice
- Jaccard
- Tversky

4. Boundary - based loss

- HD loss

(aims to max.
overlap between
ground truth and
prediction)

FCN

Fully convolution network
Transforms classification
network to segmentation:

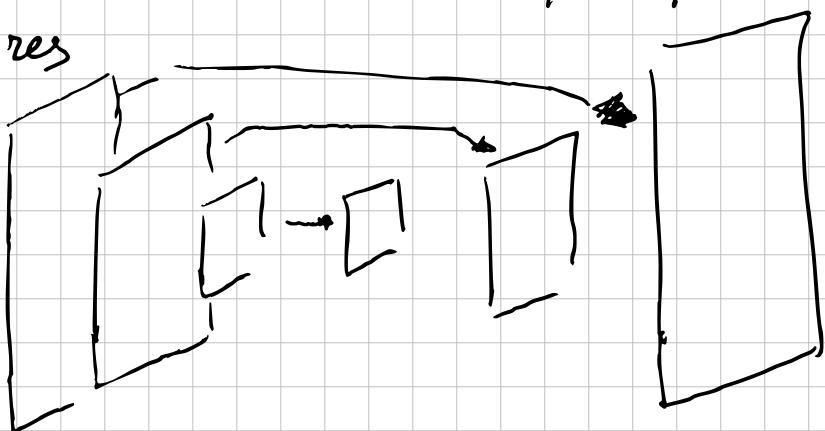
- Remove dense layer
- Add decoder part
 - Use 1×1 convolution to predict labels
 - Upsample image to default resolution

Architecture

- Fast and big upsampling gives a rough result, so we make it smoother and we mix the results with the previous encoder outputs

Segnet

- Full encoder-decoder architectures
 - Save pooling indices in encoder part and use them in decoder when upsample features

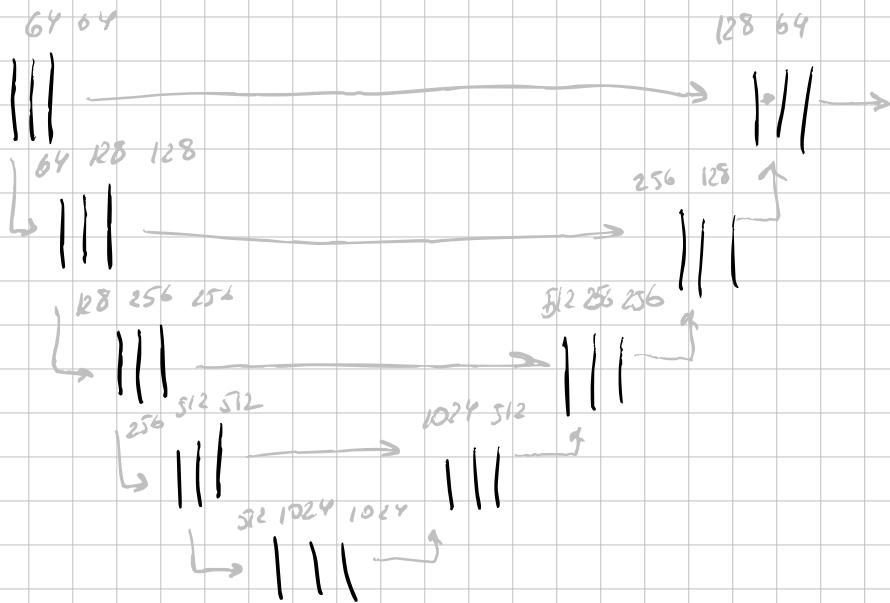


Max pooling

Saves maximum element in encoder and reuse positions from pooling layer (max unpooling)

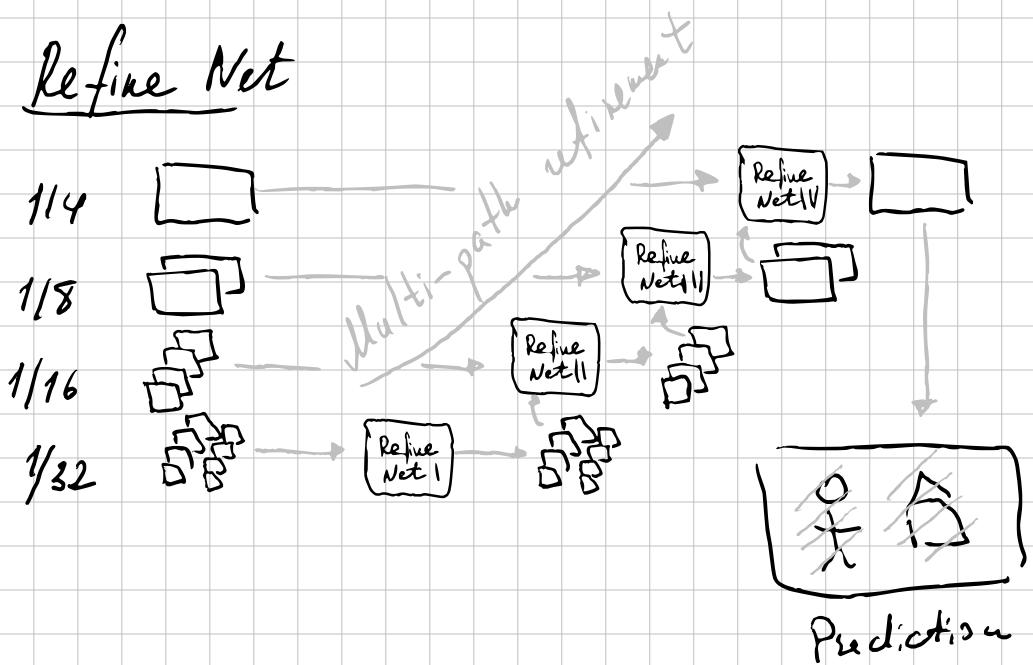
* When we use max pooling / unpooling we lost many information when unpooling and sometimes this is not enough

U-Net



- Share features from encoding to corresponding decoder layers. This sharing combine deep features from encoder which know more about all image and more locally occurring features
- Better to use pre-trained encoder from another dataset or task
 - We can train just decoder part if we have tiny dataset
 - Possible to use with multi-resolution input

Refine Net



Combined the long-range skip-connection connections from the original U-net architecture, with a lot of added skip-connections

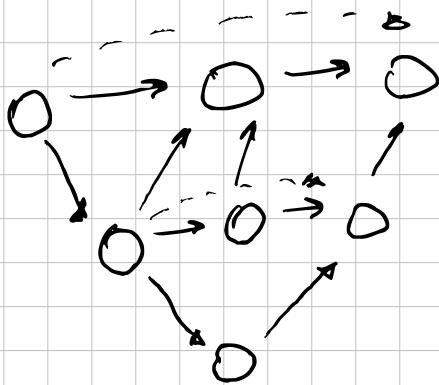
Consists of:

- Residual Conv Unit
- Multi-resolution fusion
- Chained Residual Pooling

U-Net ++

This work continues the ideas of U-net, but at the same time solves some of its problems:

1. Selected models because best results on different datasets have different depth
2. Rigid scheme of skip connections which is possible only on features of the same resolution



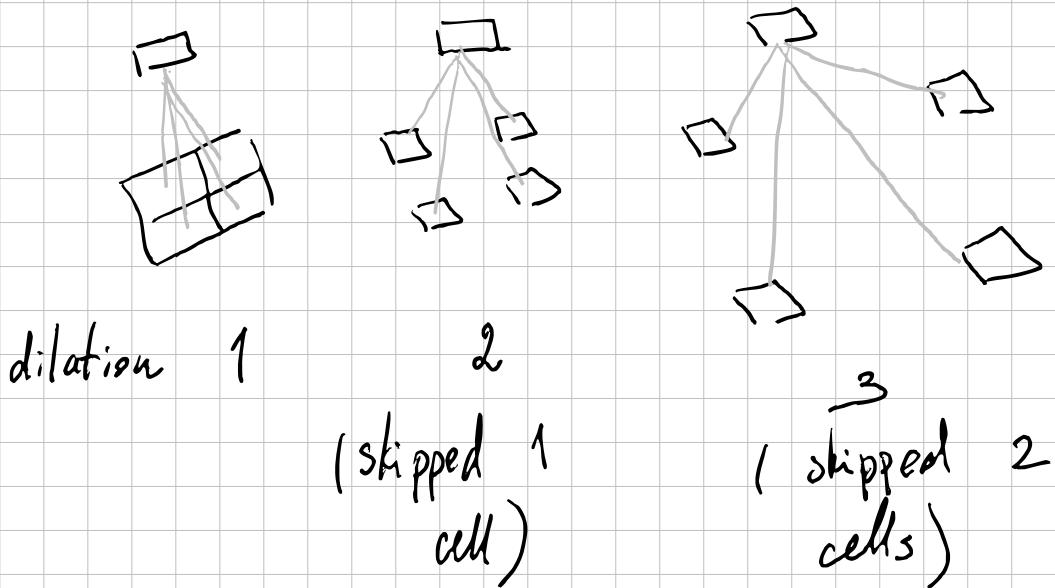
After training the whole model, we can use smaller network to improve performance

Train small u-net in u-net ++ architecture gives more IoU than isolated training

Deep Lab

Dilated convolution

Use of dilated convolution allows one to have larger receptive field (global view) with the same computation and memory costs while also preserving resolution

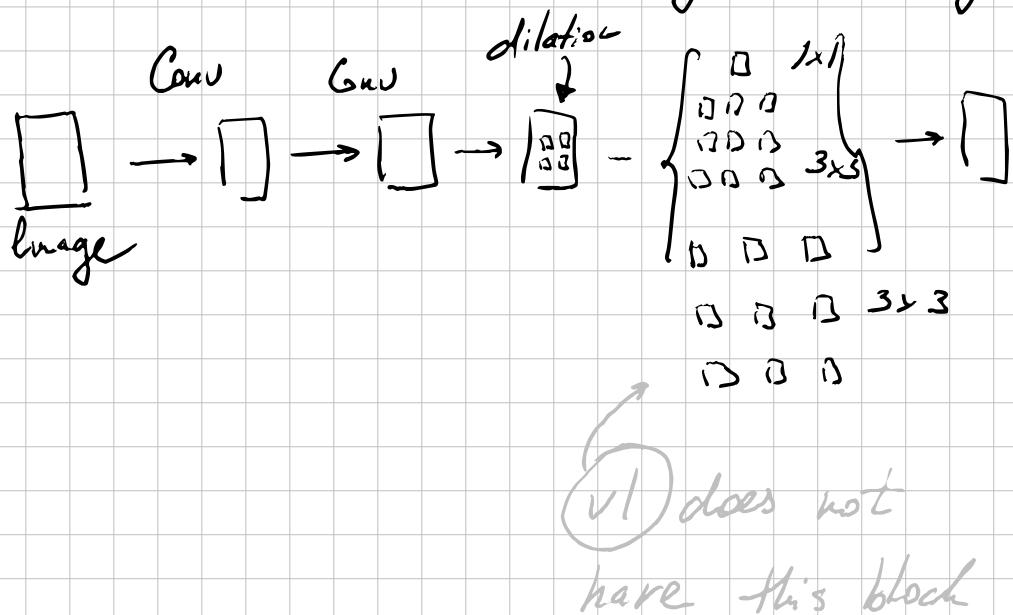


CRF

Conditional Random Field allows
to correct errors
(Post-processing step)

DeepLab v2

Atoms Spectral
Pyramidal Pooling



DeepLab v3

- Add decoder part to previous version
- Combine low-level and final encoder features
- More smoother upsample

PSPNet

Pyramid Scene Parsing Network

Use pyramid pooling module
for global scene prior construction
upon the final-layer-feature-
map of the deep neural
network

Hierarchical Multi-scale Attention

The same image with different resolutions can give different results:

- bigger scale gives better result on thin objects;
- lower scale gives better result on large objects;

So we want to combine this result using hierarchical multi-scale attention

Architecture

Hierarchical attention mechanisms learns to predict a relative weighting between adjacent scales.

* With this network we can combine different image scales

HR Net

The main body consists of several components:

- parallel multi-resolution convolutions;
- multi-resolution fusions;
- representation head.

Connects high-to-low resolution convolution streams in parallel rather than in series

Repeats multiresolution fusions to boost the high-resolution representations with the help of low-resolution repr.

Mask - R CNN

Extends Faster R-CNN

Architecture

- Generate binary mask for each class
- Change RoIPool to RoIAAlign for better score

Mask head

Collecting mask annotations, for instance, segmentation is very expensive so we want to train just using bounding boxes for all categories but use masks only for a subset of categories.

Adopt Mask-RCNN for partially-supervised training just by using better mask-head architectures without extra losses or modules.

The following factors have a big influence on generalisation of unseen classes

1. train using tight ground-truth boxes instead of combination of ground-truth and noisy proposals
2. mask-head architecture

MaskFormer

There are 2 approaches in image segmentation

per pixel classification
(for semantic segmentation)

mask classification
(for instance segmentation)

Mask classification is sufficiently general to solve both semantic- and instance-level segmentation tasks, so the authors propose a structure that transforms any per-pixel classification model to mask classification.

Architecture

- pixel-level module
- transformer module
- segmentation module