

Introduction

- Gradient Descent
- when to use deep learning
- MMSE
- Decision Theory
 - MAP
 - ML
- Entropy
- Regulatization
- Train/Test Split

Simple DNN

- Forward Propagation
- Back Propagation
 - idea
 - summary
 - regularizer
- Multi-Layer Perceptron
 - A Big Picture
 - Activation Function
 - Cost Function
 - Optimizer
- Regularization
- Dropout
- Normalization
- Dimensionality Reduction
 - PCA
 - LDA
- Hyper parameter
- Universal Approximation Theorem

CNN

- Introduction
 - Kernel(Filter)
 - Conv layer
 - Pooling layer
 - De-Conv(Transpose convolution) layer
 - Transposed 2D convolution with no padding, stride of 2 and kernel of 3
 - Group Convolution
 - Seperable Convolutions
 - Full Connect layer
- # of Parameters computation
- Visualization
- Different Nets
 - LeNet(1998)

- AlexNet(2012)
- VGG
- Resnet
- GoogleNet (Inception)
- MobileNet
- Comparation
- CNN in NLP
- Coding
- Keras

RNN

- Introduction
- Vanish Gradient Problem
- LSTM
- GRU
- Back-Propagation through Time (BPTT)
- Word Embedding
- BERT

TensorFlow and Keras

- sequential model

Pytorch

EE599 Deep Learning

Introduction

Gradient Descent

Steepest Gradient Descent

$$\widehat{\mathbf{W}}_{n+1} = \widehat{\mathbf{W}}_n - \eta \nabla_{\mathbf{W}} E(\widehat{\mathbf{W}}_n)$$

Stochastic Gradient Descent

SGD: use averages over data (typically subset of data points) to approximate ensemble averaging and therefore approximate the true gradient with a noisy (stochastic) approximation

$$\widehat{\mathbf{W}}_{n+1} = \widehat{\mathbf{W}}_n - \eta \hat{\nabla}_{\mathbf{W}} E(\widehat{\mathbf{W}}_n)$$

single point stochastic gradient descent

Single Point Stochastic Gradient Descent:

$$\begin{aligned}
 -\frac{1}{2}\nabla_{\mathbf{w}}E &= \mathbf{r}_{xy} - \mathbf{R}_x \mathbf{w} \\
 &= \mathbb{E} \{y(u)\mathbf{x}(u) - \mathbf{x}(u)\mathbf{x}^t(u)\mathbf{w}\} \\
 &\approx y_n \mathbf{x}_n - \mathbf{x}_n \mathbf{x}_n^t \mathbf{w} \\
 &= (y_n - \mathbf{x}_n^t \hat{\mathbf{w}}_n) \mathbf{x}_n
 \end{aligned}$$

$$\hat{\mathbf{w}}_{n+1} = \hat{\mathbf{w}}_n + \eta(y_n - \hat{\mathbf{w}}_n^t \mathbf{x}_n) \mathbf{x}_n$$

when to use deep learning

Don't

- You have a good statistical model for the inference problem and you can derive and implement the ideal inference function $\mathbf{f}(\cdot)$
- A simpler form of supervised ML performs well
 - eg., Linear regression, logistical regression, Naive Bayes Classifier, etc

Do

- Tough to model accurately
- You can use a classical approach, but there are many “clamps” in your algorithm (conditional statements... hacking)
- Modeling is good, but implementing $\mathbf{f}(\cdot)$ exactly is prohibitively complex

MMSE

Minimum Mean-Square Error Estimation (MMSE)

estimate $\mathbf{y}(u)$ from $\mathbf{x}(u)=\mathbf{x}$ “cross covariance matrix”
 $K_{yx} = \mathbb{E} \{(y(u) - m_y)(x(u) - m_x)^t\} = [K_{xy}]^t$

Affine MMSE:

$$\begin{aligned}
 \min_{f(x)=Fx+b} \mathbb{E} \{ \|y(u) - f(x(u))\|^2 \} &\quad F_{AMMSE} = K_{yx} K_x^{-1} && b_{AMMSE} = m_y - F_{AMMSE} m_x \\
 &\quad \hat{y} = K_{yx} K_x^{-1} (x - m_x) + m_y && AMMSE = \text{tr} (K_y - K_{yx} K_x^{-1} K_{xy})
 \end{aligned}$$

Linear MMSE:

$$\begin{aligned}
 \min_{f(x)=Fx} \mathbb{E} \{ \|y(u) - f(x(u))\|^2 \} &\quad F_{LMMSE} = R_{xy} R_x^{-1} && LMMSE = \text{tr} (R_y - R_{yx} R_x^{-1} R_{xy}) \\
 \min_F \mathbb{E} \{ \|y(u) - Fx(u)\|^2 \} &\quad \hat{y} = R_{yx} R_x^{-1} x
 \end{aligned}$$

Linear MMSE (typ. means are zero):

$$\hat{y} = \mathbf{w}^t \mathbf{x}$$

$$\mathbf{w} = \mathbf{F}^t = \mathbf{R}_x^{-1} \mathbf{r}_{xy}$$

$$\mathbf{r}_{xy} = \mathbf{R}_{xy} = \mathbb{E} \{ \mathbf{x}(u) y(u) \} = [\mathbb{E} \{ y(u) \mathbf{x}^t(u) \}]^t$$

$$\text{LMMSE} = \sigma_y^2 - \mathbf{r}_{xy}^t \mathbf{R}_x^{-1} \mathbf{r}_{xy}$$

Decision Theory

- **Bayesian Decision Theory**

- Bayes decision rule
- MAP rule - minimum error probability rule
- Maximum Likelihood
- Likelihood, Negative-Log-Likelihood, Likelihood ratios
- Neyman-Pearson test and the ROC
- Detection and False Alarm trade off

Bayes risk for decision rule r

$$\text{Risk}(r) = \int_{\mathcal{Z}} p(\mathbf{z}) \left[\sum_j r(A_j|\mathbf{z}) C(A_j|\mathbf{z}) \right] d\mathbf{z}$$

Cost for taking action m given observation \mathbf{z}

$$C(A_j|\mathbf{z}) = \sum_i C(\mathcal{H}_i, A_j) p(\mathcal{H}_i|\mathbf{z})$$

$$C_{i,j} = C(\mathcal{H}_i, A_j) = \text{Cost of deciding } \mathcal{H}_j \text{ when } \mathcal{H}_i \text{ is true}$$

Bayes decision rule (minimizes Bayes risk)

$$r_{\text{bayes}}(A_m|\mathbf{z}) = \begin{cases} 1 & m = \arg \min_j C(A_j|\mathbf{z}) \\ 0 & \text{else} \end{cases}$$

APP factoring

$$p(\mathcal{H}_m|\mathbf{z}) = \frac{p(\mathbf{z}|\mathcal{H}_m)\pi_m}{p(\mathbf{z})}$$

a posteriori probability (APP)

$\equiv p(\mathbf{z}|\mathcal{H}_m)\pi_m$

equivalent for making decisions ($p(\mathbf{z})$ not dependent on hypothesis)

MAP

MAP Estimation (parameter estimation)

Theta is a parameter set with a known/modeled a priori distribution:
like ML, but with a prior distribution on Theta

$$\hat{\Theta}_{MAP} = \hat{\Theta}_{MAP}(y) = \arg \max_{\theta} p_{y(u)|\Theta(u)}(y|\theta) p_{\Theta(u)}(\theta)$$

May be based on a conditional or joint distribution:

$$\hat{\Theta}_{MAP} = \hat{\Theta}_{MAP}(y, \mathbf{x}) = \arg \max_{\theta} p_{y(u)|\mathbf{x}(u), \Theta(u)}(y|\mathbf{x}, \theta) p_{\Theta(u)|\mathbf{x}(u)}(\theta|\mathbf{x})$$

$$\hat{\Theta}_{MAP} = \hat{\Theta}_{MAP}(y, \mathbf{x}) = \arg \max_{\theta} p_{y(u), \mathbf{x}(u), \Theta(u)}(y, \mathbf{x}, \theta) p_{\Theta(u)}(\theta)$$

ML

Maximum Likelihood Parameter Estimation

Theta is a (non-random) parameter set

$$\hat{\Theta}_{ML} = \hat{\Theta}_{ML}(y) = \arg \max_{\Theta} p_{y(u)}(y; \Theta)$$

Alternate notation: $p_{y(u)}(y|\Theta)$

recall: likelihood of Theta (given y)

The ML estimate is the parameters values that best explain the observation y under the statistical model for y(u)

May be based on a conditional or joint distribution:

$$\hat{\Theta}_{ML} = \hat{\Theta}_{ML}(y, \mathbf{x}) = \arg \max_{\Theta} p_{y(u)|\mathbf{x}(u)}(y|\mathbf{x}; \Theta)$$

$$\hat{\Theta}_{ML} = \hat{\Theta}_{ML}(y, \mathbf{x}) = \arg \max_{\Theta} p_{\mathbf{x}(u), y(u)}(\mathbf{x}, y; \Theta)$$

Properties of ML Estimators

- Asymptotically Gaussian:
 - For large amounts of data, the ML estimate is Gaussian with mean equal to the true parameter (models matched)
- Consistent:
 - The limit in probability of the ML estimate is the true parameter (model matched)
 - The ML estimate minimizes the KL Divergence between the model distribution and the empirical data distribution. KL divergence measures the difference between two distribution (Info.Theory).
 - Minimizing KL divergence in this case also corresponds to minimizing the cross entropy

Entropy

Entropy:

$$\begin{aligned} H(X(u)) &= \mathbb{E} \left\{ \log_2 \left(\frac{1}{p_{X(u)}(X(u))} \right) \right\} \\ &= \sum_k p_{X(u)}(k) \log_2 \left(\frac{1}{p_{X(u)}(k)} \right) \\ &= \sum_k p_k \log_2 \left(\frac{1}{p_k} \right) \end{aligned}$$

Intuition:

events with low probability have large information
— e.g., “it will snow in Phoenix tomorrow”

the entropy is the average information learned when the value of $X(u)$ is revealed.

weather report in Phoenix has low entropy (almost always the same), whereas in Sioux City, SD it has high entropy (highly variant weather)

Examples:

fair die:

$$H(X(u)) = \log_2(1/6) = 2.58 \text{ bits/roll}$$

loaded die:

$$\begin{aligned} H(X(u)) &= -0.4 \log_2(0.4) - 0.1 \log_2(0.1) - 0.01 \log_2(0.01) \\ &\quad - 0.09 \log_2(0.09) - 0.25 \log_2(0.25) - 0.15 \log_2(0.15) \\ &= 2.15 \text{ bits/roll} \end{aligned}$$

Cross-Entropy

Cross-Entropy

$$\text{CE}(p, \tilde{p}) = \mathbb{E}_p \left\{ \log \left(\frac{1}{\tilde{p}_x(X(u))} \right) \right\}$$

Binary Cross Entropy

Define the binary cross entropy function:

$$\begin{aligned} \text{BCE}(p, \tilde{p}) &\triangleq \mathbb{E}_p \left\{ \log \left(\frac{1}{\tilde{p}(y)} \right) \right\} \\ &= p \log \frac{1}{\tilde{p}} + (1 - p) \log \frac{1}{(1 - \tilde{p})} \\ &= -[p \log(\tilde{p}) + (1 - p) \log(1 - \tilde{p})] \end{aligned}$$

Can show that minimizing BCE minimizes a measure of difference between p and p-tilde

Regulation

enforce penalty on weights to bias toward a prior distribution. --> to reduce over-fitting (by smaller weight)

Regularization Interpretation...

$$\max_{\theta} p_{y(u)|x(u), \Theta(u)}(y|x, \theta) p_{\Theta(u)|x(u)}(\theta|x) \iff \min_w \|y - Xw\|^2 + \lambda \|w\|^2$$

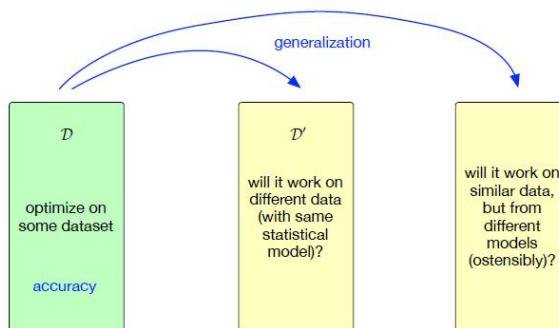
The a-priori Gaussian distribution on the weights leads to “L2 regularization”

penalizes large w — even if large w cause smaller squared error

this can be viewed a method to combat **over-fitting**

lambda is called the regularization coefficient in this context

Larger lambda ==> penalize larger weights more aggressively (at expense of SE)



in the end, we do not really care about the performance on the dataset we have (it is labeled, after all)

we care about performance on similar data that has no labels

Accuracy/Generalization trade-off (aka bias-variance trade):

generally, optimizing accuracy to the extreme will cause reduced generalization capability

Train/Test Split

Training Set -- for define trainable parameters
Validation Set -- for define hyper-parameters
Test Set -- for verify model performance
Mini-Batch -- for do one SGD update per mini-batch
epoch -- one training run through all data set
iteration -- number of mini-batches per epoch

Simple DNN

Input an example from a dataset.

The network will take that example and apply some complex computations to it using randomly initialised variables (called weights and biases).

A predicted result will be produced.

Comparing that result to the expected value will give us an error.

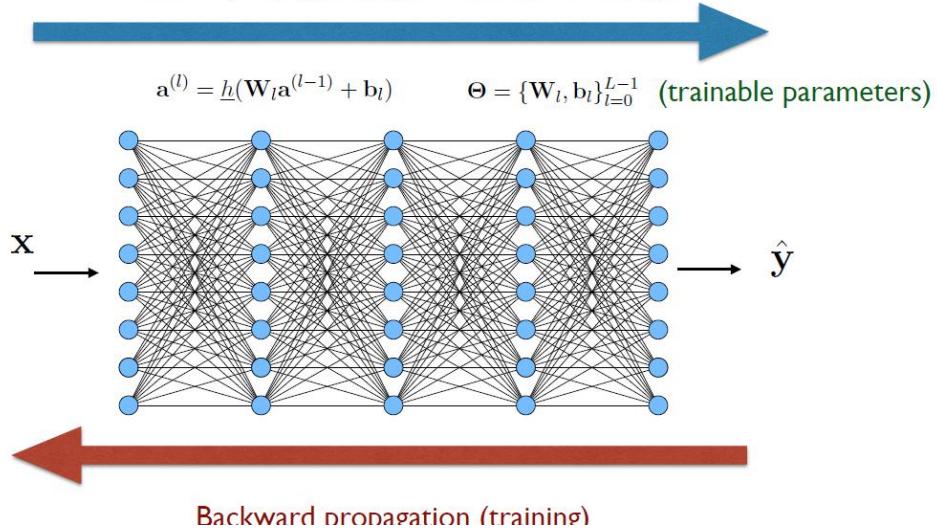
Propagating the error back through the same path will adjust the variables.

Steps 1–5 are repeated until we are confident to say that our variables are well-defined.

A prediction is made by applying these variables to a new unseen input.

MLPs

Forward propagation (inference and training)



Forward Propagation

Back Propagation

idea

do SGD on all trainable parameters:

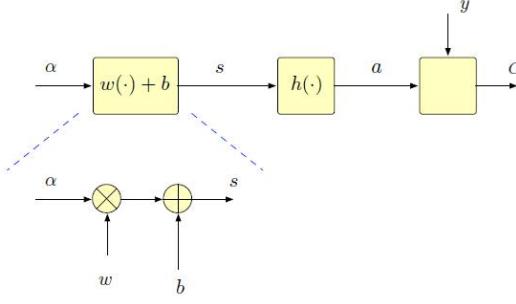
$$\text{weight updates: } w_{i,j}^{(l)} \leftarrow w_{i,j}^{(l)} - \eta \frac{\partial C}{\partial w_{i,j}^{(l)}}$$

$$\text{bias updates: } b_i^{(l)} \leftarrow b_i^{(l)} - \eta \frac{\partial C}{\partial b_i^{(l)}}$$

all layers, all indices: $\forall l \in \{1, 2, \dots, L\}, i, j$

Backprop is an algorithm for computing these, starting at the neural network output and propagating backward to the input layer using the **chain rule**

To be specified, one step of forward neural net can be extracted as follow:



want to compute: $\frac{\partial C}{\partial w}$ and $\frac{\partial C}{\partial b}$

Note that: $\frac{\partial s}{\partial w} = \alpha$ $\frac{\partial s}{\partial b} = 1$

$$\frac{\partial C}{\partial w} = \frac{\partial C}{\partial s} \frac{\partial s}{\partial w} = \alpha \frac{\partial C}{\partial s}$$

$$\frac{\partial C}{\partial b} = \frac{\partial C}{\partial s} \frac{\partial s}{\partial b} = \frac{\partial C}{\partial s}$$

so we can find $\frac{\partial C}{\partial s}$ and convert to the desired partials easily

$$\begin{aligned}\frac{\partial C}{\partial s} &= \frac{\partial C}{\partial a} \frac{\partial a}{\partial s} \\ &= \frac{\partial C}{\partial a} \frac{\partial h(s)}{\partial s} \\ &= \dot{C}(a) \dot{h}(s)\end{aligned}$$

shorthand

$$\begin{aligned}\dot{C}(a) &= \frac{\partial C}{\partial a} \\ \dot{h}(s) &= \frac{dh(s)}{ds}\end{aligned}$$

$$\begin{aligned}\frac{\partial C}{\partial w} &= \dot{C}(a) \dot{h}(s) \alpha \\ \frac{\partial C}{\partial b} &= \dot{C}(a) \dot{h}(s)\end{aligned}$$

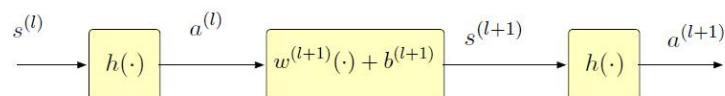
we know how to compute: $\frac{\partial C}{\partial s}$

$$\frac{\partial C}{\partial w} = \frac{\partial C}{\partial s} \frac{\partial s}{\partial w} = \alpha \frac{\partial C}{\partial s}$$

and how to convert to:

$$\frac{\partial C}{\partial b} = \frac{\partial C}{\partial s} \frac{\partial s}{\partial b} = \frac{\partial C}{\partial s}$$

now consider:



<p>Shorthand:</p> $\frac{\partial C}{\partial s^{(l)}} \leftarrow \frac{\partial C}{\partial s^{(l+1)}}$	$\delta^{(l)} \triangleq \frac{\partial C}{\partial s^{(l)}}$ <div style="border: 1px solid black; padding: 10px; margin-top: 10px;"> $\delta^{(l)} \leftarrow \delta^{(l+1)}$ </div>
$\frac{\partial C}{\partial w^{(l)}} = \frac{\partial C}{\partial s^{(l)}} \frac{\partial s^{(l)}}{\partial w^{(l)}} = \frac{\partial C}{\partial s^{(l)}} a^{(l-1)}$	$\frac{\partial C}{\partial w^{(l)}} = \delta^{(l)} a^{(l-1)}$
$\frac{\partial C}{\partial b^{(l)}} = \frac{\partial C}{\partial s^{(l)}} \frac{\partial s^{(l)}}{\partial b^{(l)}} = \frac{\partial C}{\partial s^{(l)}}$	$\frac{\partial C}{\partial b^{(l)}} = \delta^{(l)}$
<div style="border: 1px solid black; padding: 10px;"> $w^{(l)} \leftarrow w^{(l)} - \eta \delta^{(l)} a^{(l-1)}$ </div> <div style="border: 1px solid black; padding: 10px; margin-top: 10px;"> $b^{(l)} \leftarrow b^{(l)} - \eta \delta^{(l)}$ </div>	

© Keith M. Chugg, 2019

11

$\delta^{(l)} = \frac{\partial C}{\partial s^{(l)}} = \frac{\partial C}{\partial s^{(l+1)}} \frac{\partial s^{(l+1)}}{\partial s^{(l)}}$	$\delta^{(l)} \leftarrow \delta^{(l+1)}$
$= \delta^{(l+1)} \frac{\partial s^{(l+1)}}{\partial a^{(l)}} \frac{\partial a^{(l)}}{\partial s^{(l)}}$	<p>Shorthand:</p> $\dot{a}^{(l)} \triangleq \dot{h}(s^{(l)})$
$= \delta^{(l+1)} w^{(l+1)} \dot{h}(s^{(l)})$	
$= \dot{h}(s^{(l)}) w^{(l+1)} \delta^{(l+1)}$	
<div style="border: 1px solid black; padding: 10px;"> $w^{(l)} \leftarrow w^{(l)} - \eta \delta^{(l)} a^{(l-1)}$ </div> <div style="border: 1px solid black; padding: 10px; margin-top: 10px;"> $b^{(l)} \leftarrow b^{(l)} - \eta \delta^{(l)}$ </div>	

Note: you update the weights and biases only after all the deltas are computed

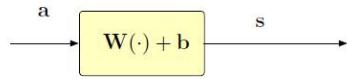
$$w^{(l)} \leftarrow w^{(l)} - \eta \delta^{(l)} a^{(l-1)}$$

$$b^{(l)} \leftarrow b^{(l)} - \eta \delta^{(l)}$$

work for matrix (vector case)

now we just need to handle the vector case...

(let's just repeat what we did)



$$\mathbf{s} = \mathbf{W}\mathbf{a} + \mathbf{b}$$

$$\begin{bmatrix} s_0 \\ s_1 \\ s_2 \end{bmatrix} = \begin{bmatrix} w_{00} & w_{01} & w_{02} \\ w_{10} & w_{11} & w_{12} \\ w_{20} & w_{21} & w_{22} \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix} + \begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix} \quad s_n = \left(\sum_m w_{nm} a_m \right) + b_n$$

Suppose we know: $\nabla_{\mathbf{s}} C = \boldsymbol{\delta}$

How to convert to: $\frac{\partial C}{\partial w_i} \quad \frac{\partial C}{\partial b_i}$

$$\frac{\partial C(s_0(w_0, w_1), s_1(w_0, w_1))}{\partial w_0} = \frac{\partial C(s_0(w_0, w_1), s_1(w_0, w_1))}{\partial s_0} \frac{\partial s_0}{\partial w_0} + \frac{\partial C(s_0(w_0, w_1), s_1(w_0, w_1))}{\partial s_1} \frac{\partial s_1}{\partial w_0}$$

$$\frac{\partial C}{\partial w_0} = \frac{\partial C}{\partial s_0} \frac{\partial s_0}{\partial w_0} + \frac{\partial C}{\partial s_1} \frac{\partial s_1}{\partial w_0}$$

$$\delta_i^{(l)} = \sum_m \delta_m^{(l+1)} \frac{\partial s_m^{(l+1)}}{\partial s_i^{(l)}}$$

$$= \sum_m \delta_m^{(l+1)} w_{mi} \dot{h}(s_i^{(l)})$$

$$\dot{a}_i^{(l)} = \dot{h}(s_i^{(l)})$$

shorthand:

$$= \left(\sum_m \delta_m^{(l+1)} w_{mi} \right) \dot{h}(s_i^{(l)})$$

$$\dot{\mathbf{a}}^{(l)} = \dot{h}(\mathbf{s}^{(l)})$$

$$= \left(\sum_m \delta_m^{(l+1)} w_{mi} \right) \dot{a}_i^{(l)}$$

$$\mathbf{v} \odot \mathbf{w} = \begin{bmatrix} v_0 w_0 \\ v_1 w_1 \\ \vdots \\ v_{D-1} w_{D-1} \end{bmatrix}$$

$$\delta^{(l)} = \dot{\mathbf{a}}^{(l)} \odot \left[\left(\mathbf{W}^{(l+1)} \right)^t \boldsymbol{\delta}^{(l+1)} \right]$$

Hadamard product: `*.` in matlab or `*` for np.arrays

summary

$$\mathbf{a}^{(l)} = \underline{h} \left(\mathbf{W}^{(l)} \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)} \right)$$

$$\dot{\mathbf{a}}^{(l)} = \dot{h} \left(\mathbf{W}^{(l)} \mathbf{a}_{i-1} + \mathbf{b}_i \right)$$

$$\boldsymbol{\delta}^{(L)} = \dot{\mathbf{a}}^{(L)} \odot \dot{\underline{C}} \left(\mathbf{y}, \mathbf{a}^{(l)} \right)$$

$$\boldsymbol{\delta}^{(l)} = \dot{\mathbf{a}}^{(l)} \odot \left[\left(\mathbf{W}^{(l+1)} \right)^t \boldsymbol{\delta}^{(l+1)} \right]$$

$$\mathbf{W}^{(l)} \leftarrow \mathbf{W}^{(l)} - \eta \boldsymbol{\delta}^{(l)} \left[\mathbf{a}^{(l-1)} \right]^t$$

$$\mathbf{b}^{(l)} \leftarrow \mathbf{b}^{(l)} - \eta \boldsymbol{\delta}^{(l)}$$

specialized to vectorized output activation and additive cost

$$\mathbf{a}_l = \text{act}(\mathbf{W}_l \mathbf{a}_{l-1} + \mathbf{b}_l) \quad (\text{activations})$$

$$\dot{\mathbf{a}}_l = \dot{\text{act}}(\mathbf{W}_l \mathbf{a}_{l-1} + \mathbf{b}_l) \quad (\text{derivative activations})$$

$$\boldsymbol{\delta}_L = \dot{\mathbf{a}}_L \odot \dot{\text{cost}}(\mathbf{y}, \mathbf{a}_L) \quad (\text{delta initialization})$$

$$\boldsymbol{\delta}_l = \dot{\mathbf{a}}_l \odot [\mathbf{W}_{l+1}^t \boldsymbol{\delta}_{l+1}] \quad (\text{delta recursion})$$

$$\mathbf{W}_l \leftarrow \mathbf{W}_l - \eta \boldsymbol{\delta}_l \mathbf{a}_{l-1}^t \quad (\text{weight SGD update})$$

$$\mathbf{b}_l \leftarrow \mathbf{b}_l - \eta \boldsymbol{\delta}_l \quad (\text{bias SGD update})$$

specialized to vectorized output activation and additive cost

using batch

so it's like compute several data points with same \mathbf{w} and $\boldsymbol{\delta}$ and use average as the update delta

For each data sample: $(\mathbf{x}[n], \mathbf{y}[n])$

Do FF and BP to compute activations, a-dots, and deltas	$\mathbf{a}_l[n] = \text{act}(\mathbf{W}_l \mathbf{a}_{l-1}[n] + \mathbf{b}_l)$ (activations)
	$\dot{\mathbf{a}}_l[n] = \dot{\text{act}}(\mathbf{W}_l \mathbf{a}_{l-1}[n] + \mathbf{b}_l)$ (derivative activations)
	$\delta_L[n] = \dot{\mathbf{a}}_L[n] \odot \text{cost}(\mathbf{y}[n], \mathbf{a}_L[n])$ (delta initialization)
	$\boldsymbol{\delta}_l[n] = \dot{\mathbf{a}}_l[n] \odot [\mathbf{W}_{l+1}^t \boldsymbol{\delta}_{l+1}[n]]$ (delta recursion)

Do one SGD update after finishing the batch	$\mathbf{W}_l \leftarrow \mathbf{W}_l - \eta \frac{1}{B} \sum_{n=0}^{B-1} \boldsymbol{\delta}_l[n] \mathbf{a}_{l-1}^t[n]$
	$\mathbf{b}_l \leftarrow \mathbf{b}_l - \eta \frac{1}{B} \sum_{n=0}^{B-1} \boldsymbol{\delta}_l[n]$

regularizer

For typical weight-penalty regularizers (e.g., L1 and L2), these are direct functions of the weights

Example: $C_{\text{reg}} = C + \lambda \sum_l \|\mathbf{W}^{(l)}\|^2 = C + \lambda \sum_l \sum_{i,j} [\mathbf{W}_{i,j}^{(l)}]^2$

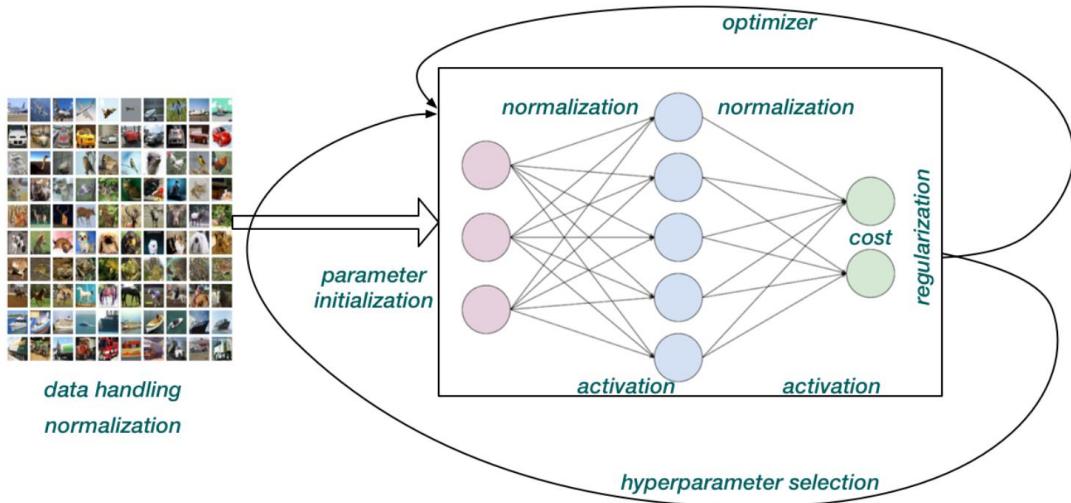
$$\frac{\partial C_{\text{reg}}}{\partial w_{i,j}^{(l)}} = \frac{\partial C}{\partial w_{i,j}^{(l)}} + 2\lambda w_{i,j}^{(l)}$$

Minor change to gradient update:

$$\mathbf{W}^{(l)} \leftarrow \mathbf{W}^{(l)} - \eta \left(\boldsymbol{\delta}^{(l)} [\mathbf{a}^{(l-1)}]^t + 2\lambda \mathbf{W}^{(l)} \right)$$

Multi-Layer Perceptron

A Big Picture



Activation Function

middle layers Activation Function

sigmoid & tanh

Squashing activations

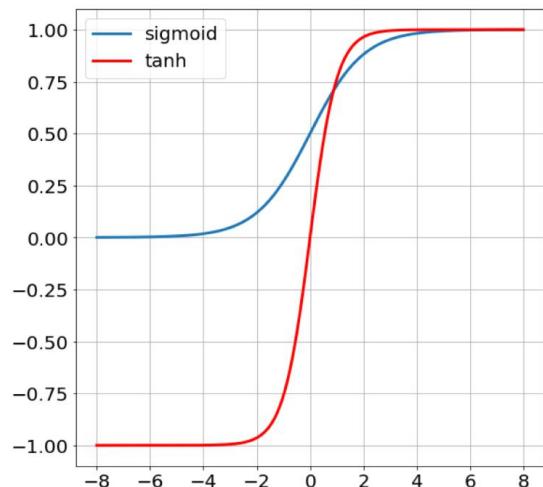
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

Sigmoid

$$\begin{aligned} \tanh(x) &= \frac{e^x - e^{-x}}{e^x + e^{-x}} \\ &= 2\sigma(2x) - 1 \end{aligned}$$

Hyperbolic Tangent
(just rescaled sigmoid)

problem: vanish



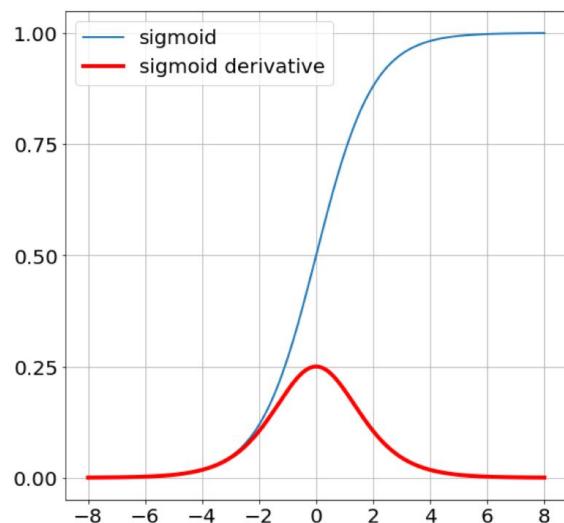
Vanishing gradients

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

Recall from BP (take example L=3)

$$\delta^{(1)} = \mathbf{H}'^{(1)} \mathbf{W}^{(2)\top} \mathbf{H}'^{(2)} \mathbf{W}^{(3)\top} \mathbf{H}'^{(3)} \nabla_{\mathbf{a}^{(3)}} C$$

All these numbers are <= 0.25, so gradients become very small!!!



more important: **ReLU family -- will not vanish**

ReLU family of activations

	$x \geq 0$	$x < 0$
Rectified Linear Unit (ReLU)	x	0
Exponential Linear Unit (ELU)	x	$\alpha(e^x - 1)$
Leaky ReLU	x	αx

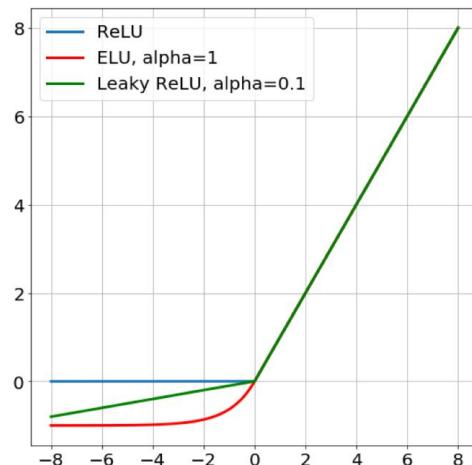
Biologically inspired - *neurons firing vs not firing*

Solves vanishing gradient problem

Non-differentiable at 0, replace with anything in [0,1]

ReLU can die if $x < 0$

Leaky ReLU solves this, but inconsistent results
ELU saturates for $x < 0$, so less resistant to noise



Clevert, Djork-Arné; Unterthiner, Thomas; Hochreiter, Sepp (2015-11-23). "Fast and

Maxout networks - Generalization of ReLU

Normally:

$$\underset{N^{(l)} \times 1}{\mathbf{s}^{(l)}} = \underset{N^{(l)} \times N^{(l-1)}}{\mathbf{W}^{(l)}} \underset{N^{(l-1)} \times 1}{\mathbf{a}^{(l-1)}} + \underset{N^{(l)} \times 1}{\mathbf{b}^{(l)}} \quad \underset{N^{(l)} \times 1}{\mathbf{a}^{(l)}} = h \left(\underset{N^{(l)} \times 1}{\mathbf{s}^{(l)}} \right)$$

For maxout:

$$\underset{k \times N^{(l)} \times 1}{\mathbf{s}^{(l)}} = \underset{k \times N^{(l)} \times N^{(l-1)}}{\mathbf{W}^{(l)}} \underset{N^{(l-1)} \times 1}{\mathbf{a}^{(l-1)}} + \underset{k \times N^{(l)} \times 1}{\mathbf{b}^{(l)}}$$

$$\underset{N^{(l)} \times 1}{\mathbf{a}^{(l)}} = \max_{\text{rows}} \underset{k \times N^{(l)} \times 1}{\mathbf{s}^{(l)}}$$

Learns the activation function itself

Better approximation power

Takes more computation

Goodfellow, Ian J.; Warde-Farley, David; Mirza, Mehdi; Courville, Aaron; Bengio, Yoshua

Example of maxout

$$k = 2, \ N^{(l)} = 5$$

$$\begin{aligned} s^{(l)}_{2 \times 5} &= \begin{bmatrix} 2.2 & -1.4 & 9.7 & -1.3 & 0 \\ -4 & -2.1 & 3.8 & -3 & 0.2 \end{bmatrix} \\ a^{(l)} &= [2.2 \quad -1.4 \quad 9.7 \quad -1.3 \quad 0.2]^T \end{aligned}$$

ReLU is a special case of maxout with k=2 and 1 set of (\mathbf{W}, b) = all 0s

Output Activation Function

Output layer activation - Softmax

$$\mathbf{a} = \begin{bmatrix} \frac{e^{s_1}}{\sum_{i=1}^N e^{s_i}} \\ \vdots \\ \frac{e^{s_N}}{\sum_{i=1}^N e^{s_i}} \end{bmatrix}$$

Extending logistic regression to multiple classes
Network output is a probability distribution!

Compare to ideal output probability distribution

$$\mathbf{y}^{(L)} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

Cost Function

Cross-entropy Cost

One Hot label vs. Softmax

Quadratic Cost (MSE)

$$C = - \sum_{i=1}^{N^{(L)}} y_i^{(L)} \ln a_i^{(L)}$$

↓ ↓
Ground truth labels Network outputs

Minimizing cross-entropy is the same as minimizing KL-divergence between the probability distributions \mathbf{y} and \mathbf{a}

For binary labels, this reduces to:

$$C = - \left[y^{(L)} \ln (a^{(L)}) + (1 - y^{(L)}) \ln (1 - a^{(L)}) \right]$$

Softmax and cross-entropy with one-hot labels

Recall:

$$\delta^{(L)} = \mathbf{H}'^{(L)T} \nabla_{\mathbf{a}^{(L)}} C$$

$$\delta_r^{(L)} = a_r^{(L)} - 1$$

$$\delta_{\neq r}^{(L)} = a_{\neq r}^{(L)}$$

Combining: $\boldsymbol{\delta}^{(L)} = \mathbf{a}^{(L)} - \mathbf{y}^{(L)}$

This makes beautiful sense as the error vector!

Optimizer

Hessian

Hessian: Matrix of double derivatives

Given $f(\mathbf{w}) : \mathbb{R}^n \rightarrow \mathbb{R}$

$$\mathbf{H} = \begin{bmatrix} \frac{\partial^2 f}{\partial w_1^2} & \cdots & \frac{\partial^2 f}{\partial w_1 \partial w_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial w_m \partial w_1} & \cdots & \frac{\partial^2 f}{\partial w_m^2} \end{bmatrix}$$

For continuous functions:
 $\frac{\partial^2 f}{\partial w_i \partial w_j} = \frac{\partial^2 f}{\partial w_j \partial w_i}$

Hessian is symmetric

Condition

Conditioning in terms of Hessian

Conditioning is measured as the condition number of the Hessian of the cost function

$$\kappa = \frac{\sigma_{\max}}{\sigma_{\min}} = \frac{|d_{\max}|}{|d_{\min}|}$$

I'm using d for eigenvalues

Examples:

$$C(w_1, w_2) = (w_1 - 2)^2 + (w_2 - 3)^2 \quad \kappa = 1 \Rightarrow \text{Well-conditioned}$$

$$C(w_1, w_2) = (w_1 - 2)^2 + 9(w_2 - 3)^2 \quad \kappa = 9 \Rightarrow \text{Ill-conditioned}$$

will condition lead to fast gradient

Momentum

Normal update:

$$p \leftarrow p - \eta \nabla_p C$$

Momentum update:

$$v \leftarrow \alpha v - \eta \nabla_p C$$

$$p \leftarrow p + v$$

Equivalent to smoothing the update by a low pass filter

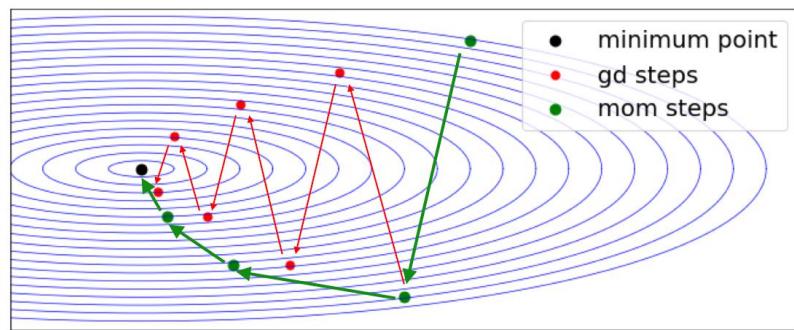
$$\alpha \in [0, 1]$$

How much of past history should be applied, typically ~ 0.9

α (Momentum factor) normally **increase** with the step goes

η (Learning Rate) **decrease** with step goes

Gradient descent with momentum
converges quickly even when ill-conditioned



Nesterov Momentum

Nesterov Momentum

Normal update:

$$p \leftarrow p - \eta \nabla_p C$$

Momentum update:

$$v \leftarrow \alpha v - \eta \nabla_p C$$

$$p \leftarrow p + v$$

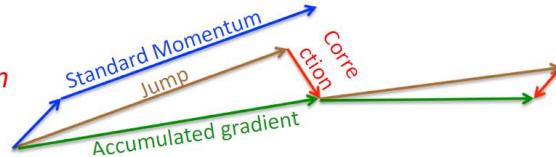
Nesterov Momentum update:

$$\tilde{p} \leftarrow p + \alpha v$$

$$v \leftarrow \alpha v - \eta \nabla_{\tilde{p}} C(\tilde{p})$$

$$p \leftarrow p + v$$

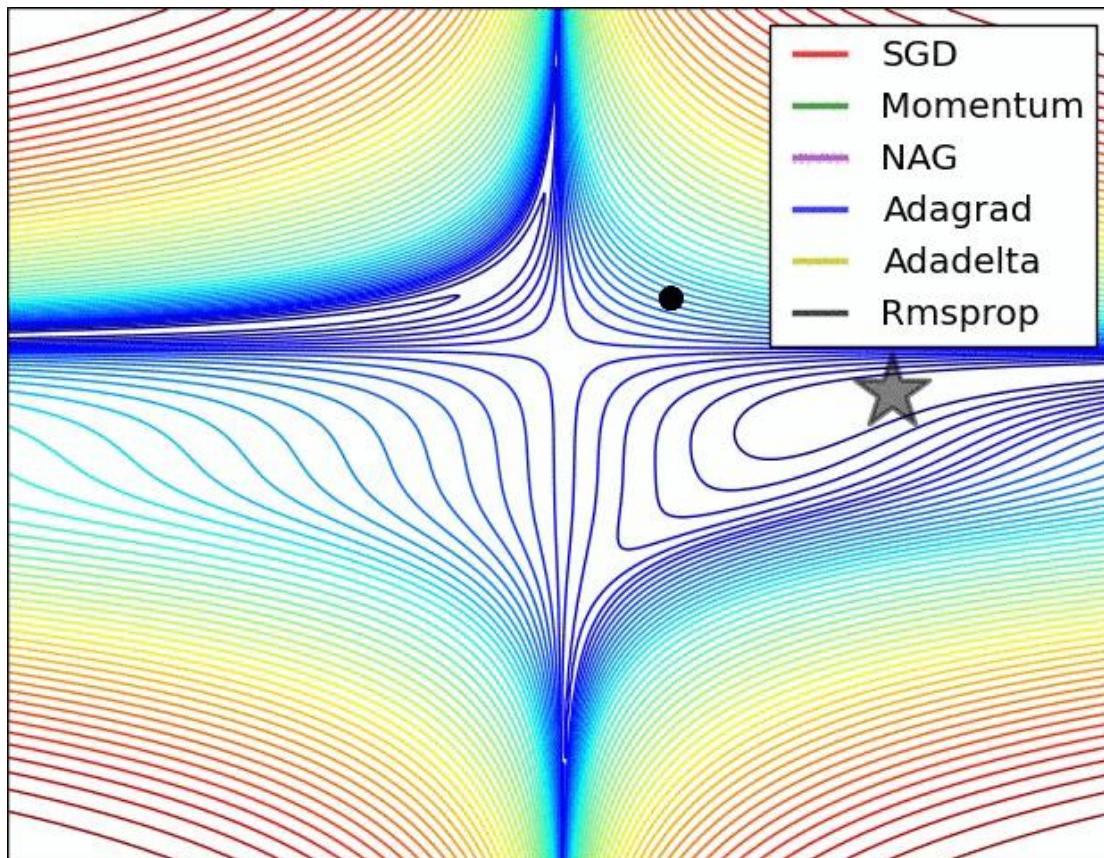
*Correction factor applied to momentum
May give faster convergence*

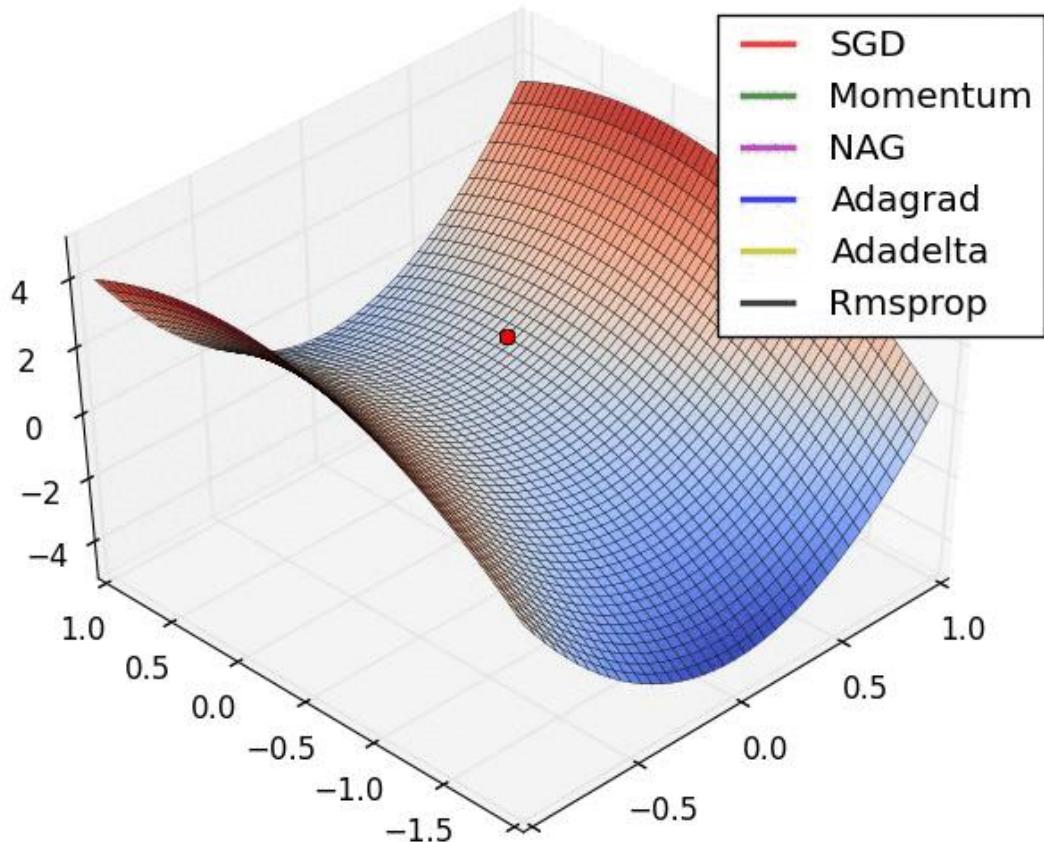


Different Kinds of Optimizers

Adaptive Optimizers

- ✓ Learning rate should be reduced with time. We don't want to overshoot the minimum point and oscillate.
- ✓ Learning rate should be small for sensitive parameters, and vice-versa.
- ✓ Momentum factor should be increased with time to get smoother updates.





RMSprop

Scale each gradient by an exponentially weighted moving average of its past history

$$v \leftarrow \rho v + (1 - \rho) (\nabla_p C)^2$$

$$p \leftarrow p - \frac{\eta}{\sqrt{v} + \epsilon} \nabla_p C$$

Default $\rho = 0.9$

ϵ is just a small number to prevent division by 0. Can be machine epsilon

Adam (more complex compared to RMSprop)

Adam

RMSprop with momentum and bias correction

At time step (t+1):

$$\begin{aligned}
 v_1 &\leftarrow \rho_1 v_1 + (1 - \rho_1) \nabla_p C && \text{Momentum} \\
 v_2 &\leftarrow \rho_2 v_2 + (1 - \rho_2) (\nabla_p C)^2 && \text{Exponentially weighted} \\
 &&& \text{moving average of past history} \\
 \tilde{v}_1 = \frac{v_1}{1 - \rho_1^t} && \left. \begin{array}{l} \text{Bias corrections to make} \\ \mathbb{E}\{\tilde{v}_1\} = \mathbb{E}\{\nabla_p C\} \end{array} \right\} \\
 \tilde{v}_2 = \frac{v_2}{1 - \rho_2^t} && \left. \begin{array}{l} \mathbb{E}\{\tilde{v}_2\} = \mathbb{E}\{(\nabla_p C)^2\} \end{array} \right\} \\
 p &\leftarrow p - \frac{\eta}{\sqrt{\tilde{v}_2} + \epsilon} \tilde{v}_1 && \text{Defaults:} \\
 &&& \eta=0.001, \rho_1 = 0.9, \rho_2 = 0.999
 \end{aligned}$$

SGD

(may include momentum and learning rate)

more basic, need more experience

AdaMax

The v_t factor in the Adam update rule scales the gradient inversely proportionally to the ℓ_2 norm of the past gradients (via the v_{t-1} term) and current gradient $|g_t|^2$:

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) |g_t|^2$$

We can generalize this update to the ℓ_p norm. Note that Kingma and Ba also parameterize β_2 as β_2^p :

$$v_t = \beta_2^p v_{t-1} + (1 - \beta_2^p) |g_t|^p$$

Norms for large p values generally become numerically unstable, which is why ℓ_1 and ℓ_2 norms are most common in practice. However, ℓ_∞ also generally exhibits stable behavior. For this reason, the authors propose AdaMax (Kingma and Ba, 2015) and show that v_t with ℓ_∞ converges to the following more stable value. To avoid confusion with Adam, we use u_t to denote the infinity norm-constrained v_t :

$$\begin{aligned}
 u_t &= \beta_2^\infty v_{t-1} + (1 - \beta_2^\infty) |g_t|^\infty \\
 &= \max(\beta_2 \cdot v_{t-1}, |g_t|)
 \end{aligned}$$

We can now plug this into the Adam update equation by replacing $\sqrt{\hat{v}_t + \epsilon}$ with u_t to obtain the AdaMax update rule:

$$\theta_{t+1} = \theta_t - \frac{\eta}{u_t} \hat{m}_t$$

Note that as u_t relies on the max operation, it is not as suggestible to bias towards zero as m_t and v_t in Adam, which is why we do not need to compute a bias correction for u_t . Good default values are again $\eta = 0.002$, $\beta_1 = 0.9$, and $\beta_2 = 0.999$.

Adadelta (fastest)

Adadelta [13] is an extension of Adagrad that seeks to reduce its aggressive, monotonically decreasing learning rate. Instead of accumulating all past squared gradients, Adadelta restricts the window of accumulated past gradients to some fixed size w .

Instead of inefficiently storing w previous squared gradients, the sum of gradients is recursively defined as a decaying average of all past squared gradients. The running average $E[g^2]_t$ at time step t then depends (as a fraction γ similarly to the Momentum term) only on the previous average and the current gradient:

$$E[g^2]_t = \gamma E[g^2]_{t-1} + (1 - \gamma) g_t^2.$$

We set γ to a similar value as the momentum term, around 0.9. For clarity, we now rewrite our vanilla SGD update in terms of the parameter update vector $\Delta\theta_t$:

$$\begin{aligned}\Delta\theta_t &= -\eta \cdot g_{t,i} \\ \theta_{t+1} &= \theta_t + \Delta\theta_t\end{aligned}$$

The parameter update vector of Adagrad that we derived previously thus takes the form:

$$\Delta\theta_t = -\frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t.$$

We now simply replace the diagonal matrix G_t with the decaying average over past squared gradients $E[g^2]_t$:

$$\Delta\theta_t = -\frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t.$$

As the denominator is just the root mean squared (RMS) error criterion of the gradient, we can replace it with the criterion short-hand:

$$\Delta\theta_t = -\frac{\eta}{RMS[g]_t} g_t.$$

The authors note that the units in this update (as well as in SGD, Momentum, or Adagrad) do not match, i.e. the update should have the same hypothetical units as the parameter. To realize this, they first define another exponentially decaying average, this time not of squared gradients but of squared parameter updates:

$$E[\Delta\theta^2]_t = \gamma E[\Delta\theta^2]_{t-1} + (1 - \gamma) \Delta\theta_t^2.$$

The root mean squared error of parameter updates is thus:

$$RMS[\Delta\theta]_t = \sqrt{E[\Delta\theta^2]_t + \epsilon}.$$

Since $RMS[\Delta\theta]_t$ is unknown, we approximate it with the RMS of parameter updates until the previous time step. Replacing the learning rate η in the previous update rule with $RMS[\Delta\theta]_{t-1}$ finally yields the Adadelta update rule:

$$\begin{aligned}\Delta\theta_t &= -\frac{RMS[\Delta\theta]_{t-1}}{RMS[g]_t} g_t \\ \theta_{t+1} &= \theta_t + \Delta\theta_t\end{aligned}$$

With Adadelta, we do not even need to set a default learning rate, as it has been eliminated from the update rule.

which optimizer is best?

RMSprop is an extension of Adagrad that deals with its radically diminishing learning rates. It is **identical to Adadelta**, except that Adadelta uses the RMS of parameter updates in the numerator update rule. Adam, finally, adds bias-correction and momentum to RMSprop. Insofar, **RMSprop**, **Adadelta**, and **Adam** are very similar algorithms that do well in similar circumstances.

Regularization

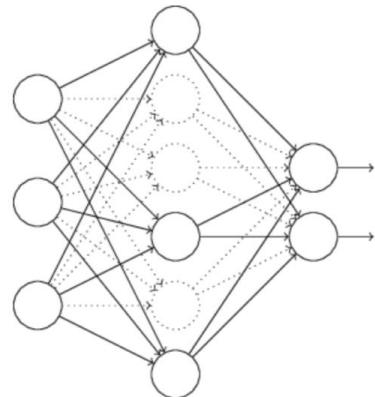
bias-variance trade-off (in EE660) (for a given MSE)

Dropout

For every minibatch during training, delete a random selection of input and hidden nodes

Keep probability p_k : Fraction of units to keep (i.e. not drop)

But why dropout??



Esemble method turn to drop out

Return to dropout

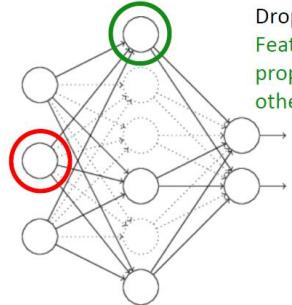
Dropout is an ensemble method using the same shared weights
=> Computationally cheaper

Improves robustness, since individual nodes have to work without depending on other deleted nodes => Regularization

How much dropout?

Original paper recommendations:

Input $p_k = 0.8$



Hidden $p_k = 0.5$

Drop hidden:
Features still propagate through other hidden nodes

Drop output: Cannot classify!

0.5 gives maximum regularization effect:
P. Baldi, P. J. Sadowski, "Understanding Dropout", Proc. NIPS 2013

But you should try other values as well!

Output $p_k = 1$

Conv layer $p_k = 0.7$

Drop input: Feature gone forever!
Some people use input $p_k = 1$

N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov,
"Dropout: A simple way to prevent neural networks from overfitting,"
Journal of Machine Learning Research, vol. 15, pp. 1929–1958, 2014.

- All nodes and weights are present in inference. Multiply weights by p_k during inference to compensate
 - This assumes linearity, but still works!
- Dropout is like multiplicative 0-1 noise for each node
- Dropout is best for big networks
 - Increase capacity, increase generalization power, but prevent overfitting

Normalization

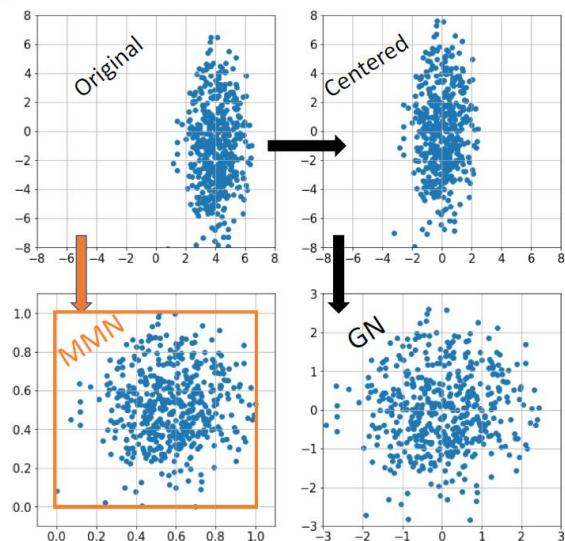
Essential preprocessing

$$\mathbf{X}_{:,i} = \frac{\mathbf{X}_{:,i} - \mu_i}{\sigma_i}$$

Gaussian normalization:
Each feature is a unit Gaussian

$$\mathbf{X}_{:,i} = \frac{\mathbf{X}_{:,i} - m_i}{M_i - m_i}$$

Minmax normalization:
Each feature is in [0,1]



Dimensionality Reduction

SVD: singular value decomposition

PCA

Principal Component Analysis (PCA)

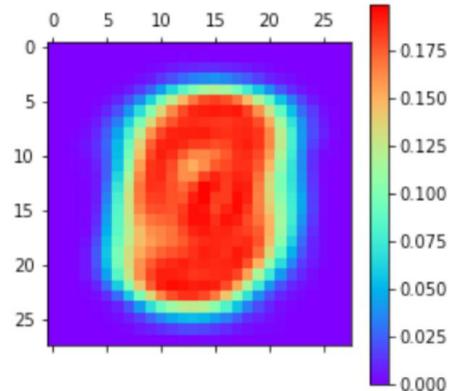
*Keep features with maximum variance
=> features which change the most*

Assume sorted

$$\mathbf{U}_{f \times f} \quad \Sigma_{f \times f} \quad \mathbf{V}_{f \times f} = \text{SVD} \{ \text{Cov}(\mathbf{X}) \}$$

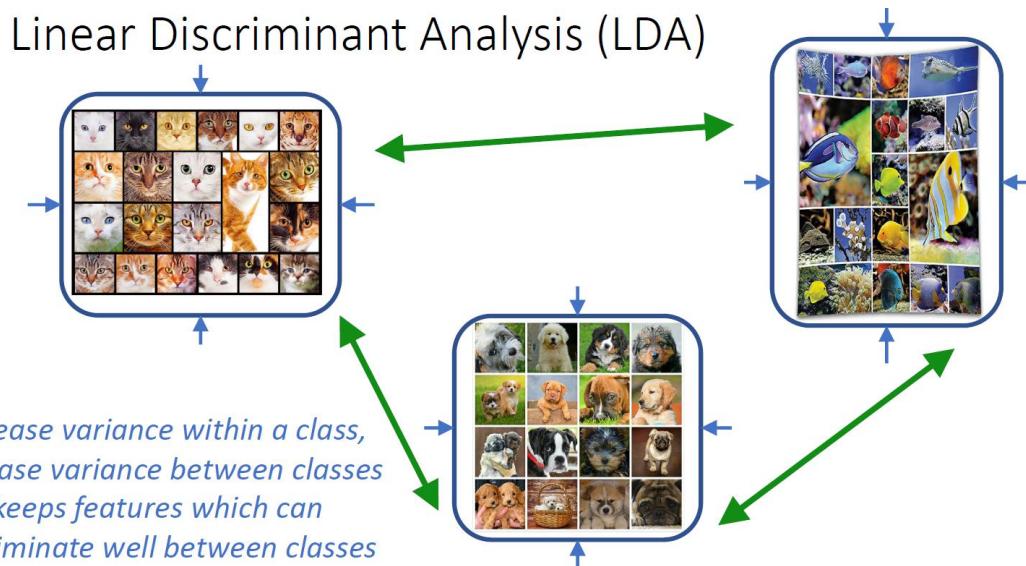
$$\mathbf{X}_{\text{PCA}} = \mathbf{X}_{n \times f} \quad \mathbf{U}_{f \times f'}^{[:, :f']}$$

$f' < f$



LDA

need labels of all data, change data based on it's label



Linear Discriminant Analysis (LDA)

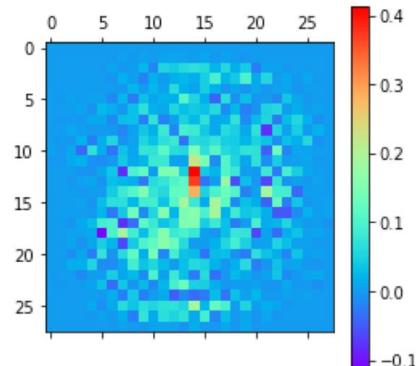
Intra (within) class scatter
Minimize

$$S_{\text{intra}} = \sum_{c=1}^C (\tilde{\mathbf{X}}_c - \boldsymbol{\mu}_c)^T (\tilde{\mathbf{X}}_c - \boldsymbol{\mu}_c)$$

Inter (between) class scatter
Maximize

$$S_{\text{inter}} = \sum_{c=1}^C N_c (\boldsymbol{\mu}_c - \boldsymbol{\mu}) (\boldsymbol{\mu}_c - \boldsymbol{\mu})^T$$

$$\begin{aligned} \mathbf{U}\Sigma\mathbf{V} &= \text{SVD} \left\{ S_{\text{intra}}^{-1} S_{\text{inter}} \right\} \\ \mathbf{X}_{\text{LDA}} &= \mathbf{X}_{n \times f} \mathbf{U}[:, :f']_{f \times f'} \end{aligned}$$



C: Number of classes
N_c: Number of samples in class c
 $\boldsymbol{\mu}_c$: Feature mean of class c

Hyper parameter

things Neural Network don't learning by themself

Continuous Examples:

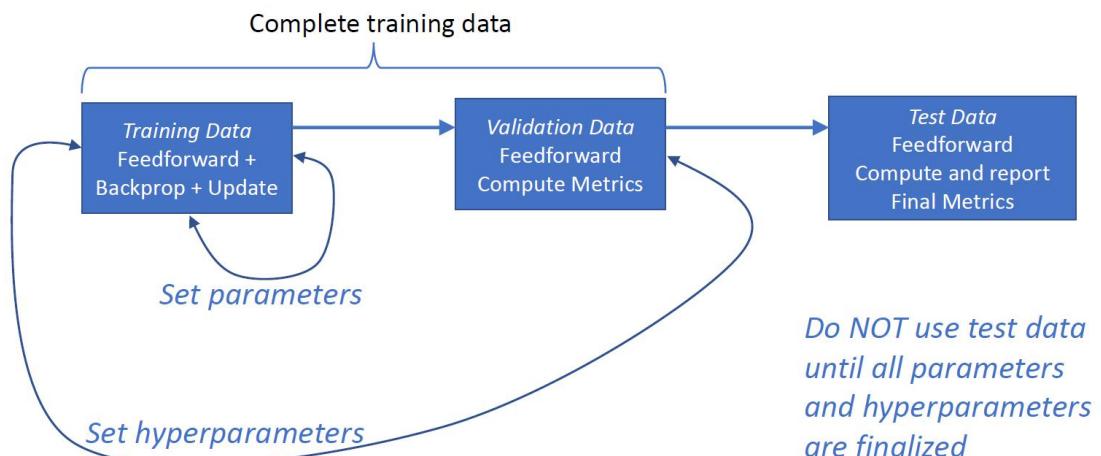
- Learning rate η
- Momentum α
- Regularization λ



Discrete Examples:

- What optimizer: *SGD, Adam, RMSprop, ...?*
- What regularization: *L2, L1, both, ...?*
- What initialization: *Glorot, He, normal, uniform, ...?*
- Batch size M: *32, 64, 128, ...?*
- Number of MLP hidden layers: *1, 2, 3, ...?*
- How many neurons in each layer?
- What kinds of data augmentation?

Using data properly

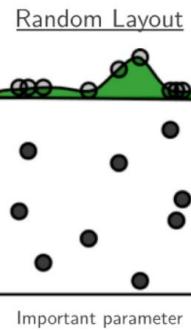
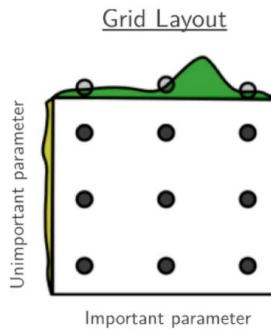


use cross validation to set hyper-parameters

Hyper-parameter search

Grid search for (η, α) :

$(0.1, 0.5), (0.1, 0.9), (0.1, 0.99),$
 $(0.01, 0.5), (0.01, 0.9), (0.01, 0.99),$
 $(0.001, 0.5), (0.001, 0.9), (0.001, 0.99)$



Random search for (η, α) :

$(0.07, 0.68), (0.002, 0.94), (0.008, 0.99),$
 $(0.08, 0.88), (0.005, 0.62), (0.09, 0.75),$
 $(0.14, 0.81), (0.006, 0.76), (0.01, 0.8)$

Random search samples more values of the important hyperparameter

Pic courtesy: J. Bergstra, Y. Bengio, "Random Search for Hyperparameter Optimization"

Learning Rate Schedules

$$\eta_{t+1} = \frac{\eta_0}{1 + kt}$$

Initial η →
Epochs →

Fractional decay

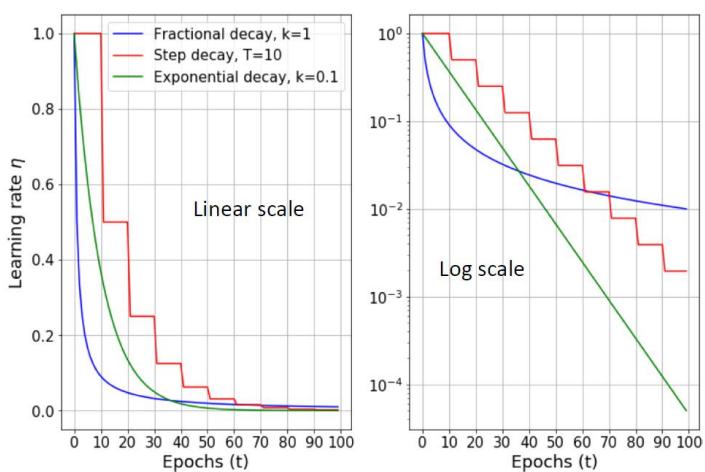
Used in Keras

$$\eta_{T+1} = \eta_0 \left(\frac{1}{2}\right)^{\lfloor t/T \rfloor}$$

Step decay

$$\eta_{t+1} = \eta_0 e^{(-kt)}$$

Exponential decay



Other Learning rate strategies

Potential project topic

Warmup: Increase η for a few epochs at the beginning. Works well for large batch sizes.

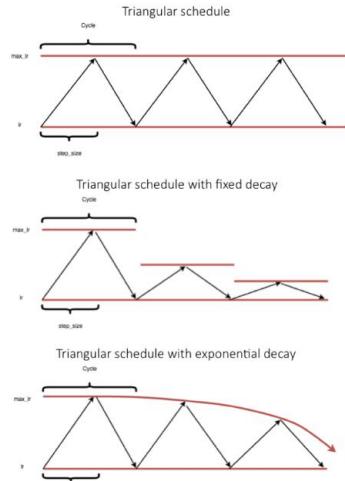
Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, Kaiming He. "Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour". arXiv:1706.02677

Cosine Scheduling

I. Loshchilov, F. Hutter, "SGDR: Stochastic Gradient Descent with Warm Restarts", Proc. ICLR 2017.

Triangular Scheduling

L. N. Smith, "Cyclical Learning Rates for Training Neural Networks", arXiv:1506.01186
Pic courtesy <https://www.jeremyjordan.me/nn-learning-rate/>



Universal Approximation Theorem

basic idea: Neural Network can simulate all function

<http://neuralnetworksanddeeplearning.com/chap4.html>

Let $\varphi(\cdot)$ be a nonconstant, bounded, and monotone-increasing continuous function. Let I_{m_0} denote the m_0 -dimensional unit hypercube $[0, 1]^{m_0}$. The space of continuous functions on I_{m_0} is denoted by $C(I_{m_0})$. Then, given any function $f \in C(I_{m_0})$ and $\varepsilon > 0$, there exist an integer m_1 and sets of real constants α_i , b_i , and w_{ij} , where $i = 1, \dots, m_1$ and $j = 1, \dots, m_0$ such that we may define

$$F(x_1, \dots, x_{m_0}) = \sum_{i=1}^{m_1} \alpha_i \varphi \left(\sum_{j=1}^{m_0} w_{ij} x_j + b_i \right) \quad (4.88)$$

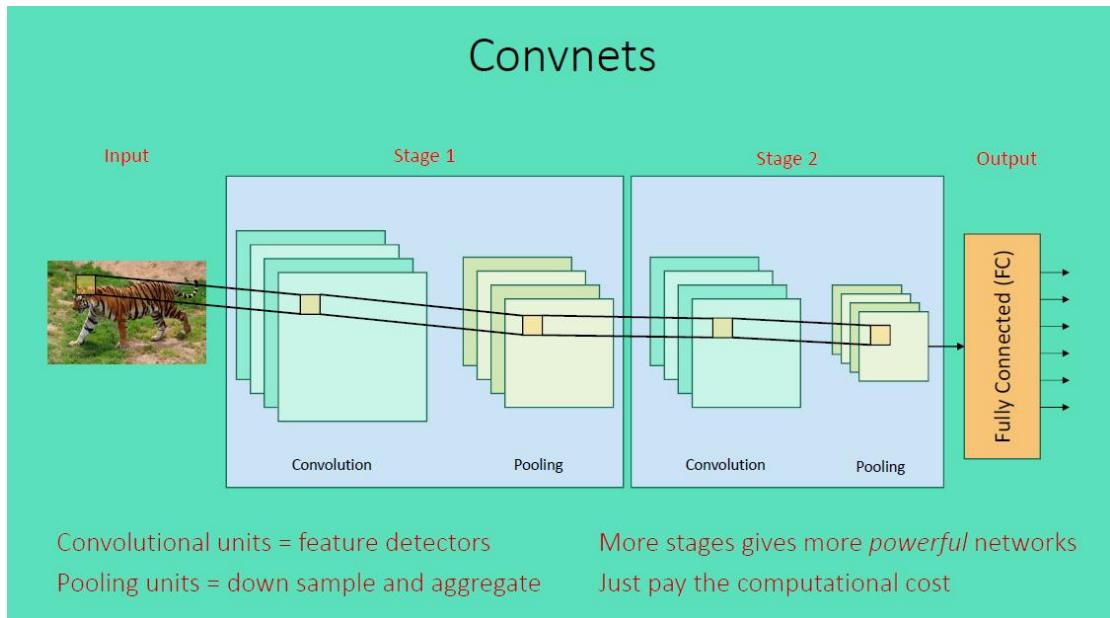
as an approximate realization of the function $f(\cdot)$; that is,

$$|F(x_1, \dots, x_{m_0}) - f(x_1, \dots, x_{m_0})| < \varepsilon$$

for all x_1, x_2, \dots, x_{m_0} that lie in the input space.

A single hidden layer MLP with squashing activation in the hidden layer and linear output layer can approximate any “engineering function”

CNN



Introduction

good introduction: <http://cs231n.github.io/convolutional-networks/#pool>

In summary:

A ConvNet architecture is in the simplest case a list of Layers that transform the image volume into an output volume (e.g. holding the class scores)

There are a few distinct types of Layers (e.g. CONV/FC/RELU/POOL are by far the most popular)

Each Layer accepts an input 3D volume and transforms it to an output 3D volume through a differentiable function

Each Layer may or may not have parameters (e.g. CONV/FC do, RELU/POOL don't)

Each Layer may or may not have additional hyperparameters (e.g. CONV/FC/POOL do, RELU doesn't)

Kernel(Filter)

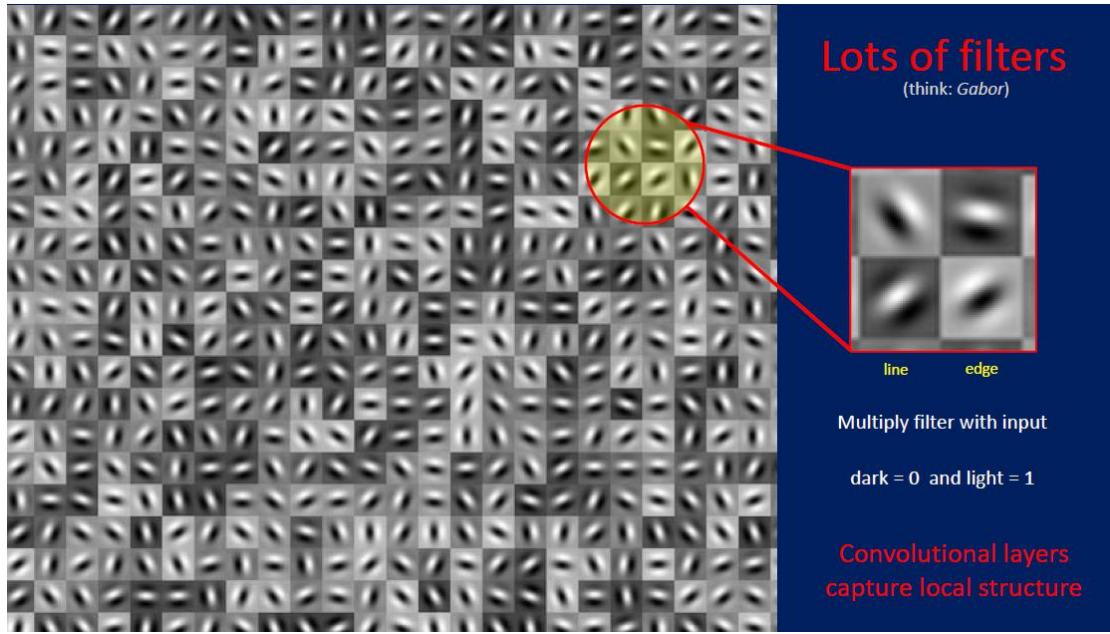
Slide

Padding

size:

small (3×3) deep

large (5×5) fast reduction, shallow, more parameters



Conv layer

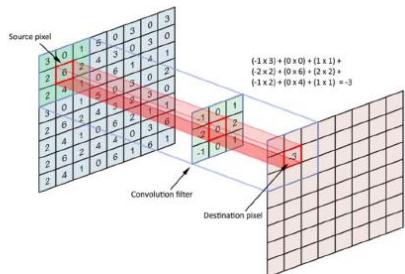
image: 3-D in volume

basic: each neural is only connect to some of the neural in next layer

Basic Function of Conv-layer:

also like $\Sigma(wx) + b$ (if b exist?)

Discrete “convolution” is a local dot product.



$$\begin{aligned} S(i, j) &= (I * K)(i, j) \\ &= \sum_m \sum_n I(i+m, j+n) K(m, n) \end{aligned}$$

K has small dimension $\Rightarrow \therefore$ fast compute

use dot production to compute kernel* cut part

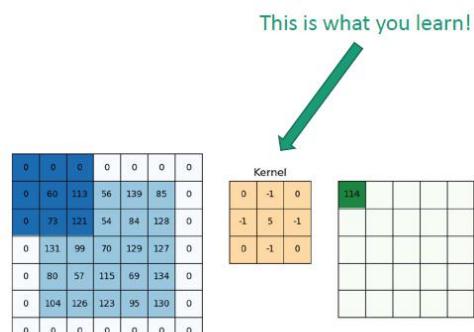
Intuitively, the network will learn filters that activate when they see some type of visual feature such as an edge of some orientation or a blotch of some color on the first layer, or eventually entire honeycomb or wheel-like patterns on higher layers of the network.

Hyper-Parameters

depth

number of filters we want to use, they will look into the same region

stride



commonly--> 1

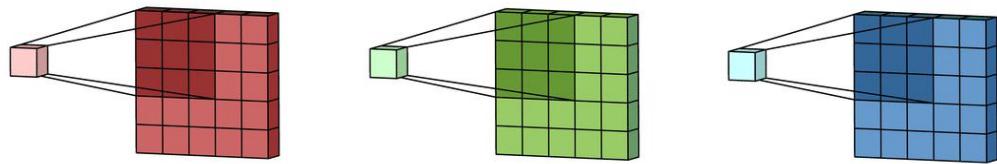
zero-padding

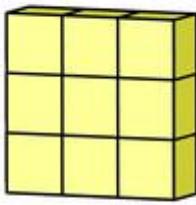
we could:

1. control the output volume size
2. remain the edge information

For three channel

three filters to combine a kernel, three output add up together, and add bias





dilation

it's possible to have filters that have spaces between each cell, called dilation.

in one dimension a filter w of size 3 would compute over input x the

following: $w[0]*x[0] + w[1]*x[1] + w[2]*x[2]$. This is dilation of 0. For dilation 1

the filter would instead compute $w[0]*x[0] + w[1]*x[2] + w[2]*x[4]$

so the computed block is not continuous

compute output size(O)

by input size(**W**), filter size(**F**) and stride(**S**) and amount of zero padding(**P**)

$$O = (W - F + 2P) / S + 1$$

input W^*W output O^*O

Parameter Sharing(weight sharing)

one feature is useful to compute at some spatial position (x,y) , then it should also be useful to compute at a different position $(x2,y2)$. In other words, denoting a single 2-dimensional slice of depth as a **depth slice**

we are going to constrain the neurons in each depth slice to use the same weights and bias
for example: for a $11*11$ kernel in RGB with 96 filters

there are $96 * 11 * 11 * 3 + 96 = 34944$ parameters in total

all $55 * 55$ neurons in each depth slice will now be using the same parameters.

questions? **each filter has three channel or just one?**

matrix. If we were to use a kernel K of size 3 on the reshaped 4×4 input to get a 2×2 output, the equivalent transformation matrix would be:

$$\begin{bmatrix} k_{1,1} & k_{1,2} & k_{1,3} & 0 & k_{2,1} & k_{2,2} & k_{2,3} & 0 & k_{3,1} & k_{3,2} & k_{3,3} & 0 & 0 & 0 & 0 & 0 \\ 0 & k_{1,1} & k_{1,2} & k_{1,3} & 0 & k_{2,1} & k_{2,2} & k_{2,3} & 0 & k_{3,1} & k_{3,2} & k_{3,3} & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & k_{1,1} & k_{1,2} & k_{1,3} & 0 & k_{2,1} & k_{2,2} & k_{2,3} & 0 & k_{3,1} & k_{3,2} & k_{3,3} & 0 \\ 0 & 0 & 0 & 0 & 0 & k_{1,1} & k_{1,2} & k_{1,3} & 0 & k_{2,1} & k_{2,2} & k_{2,3} & 0 & k_{3,1} & k_{3,2} & k_{3,3} \end{bmatrix}$$

a new method: separable convolution

Xception

convert 3*3 kernel to 3*1 and 1*3 (used in Enmedded system)

Pooling layer

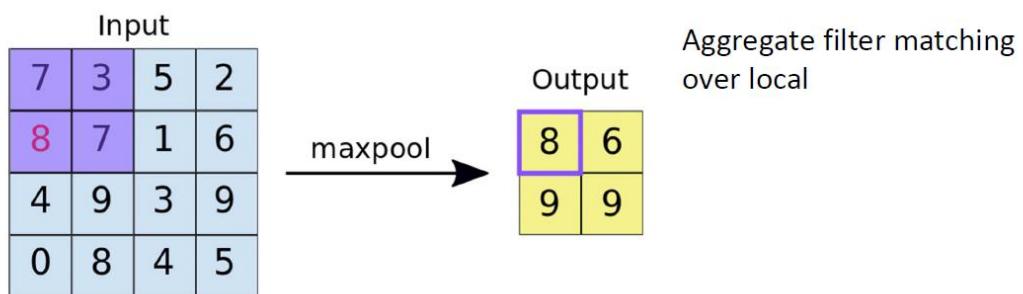
Its function is to progressively reduce the spatial size of the representation to **reduce the amount of parameters and computation in the network**, and hence to also control overfitting.

Max pooling | average pooling | L2-norm pooling

Extents--> # of points pooling together

Stride--> # of step to move in next pooling

Pooling layers downsample



Downsampling does more than just reduce size
It concentrates the image for deeper feature detection

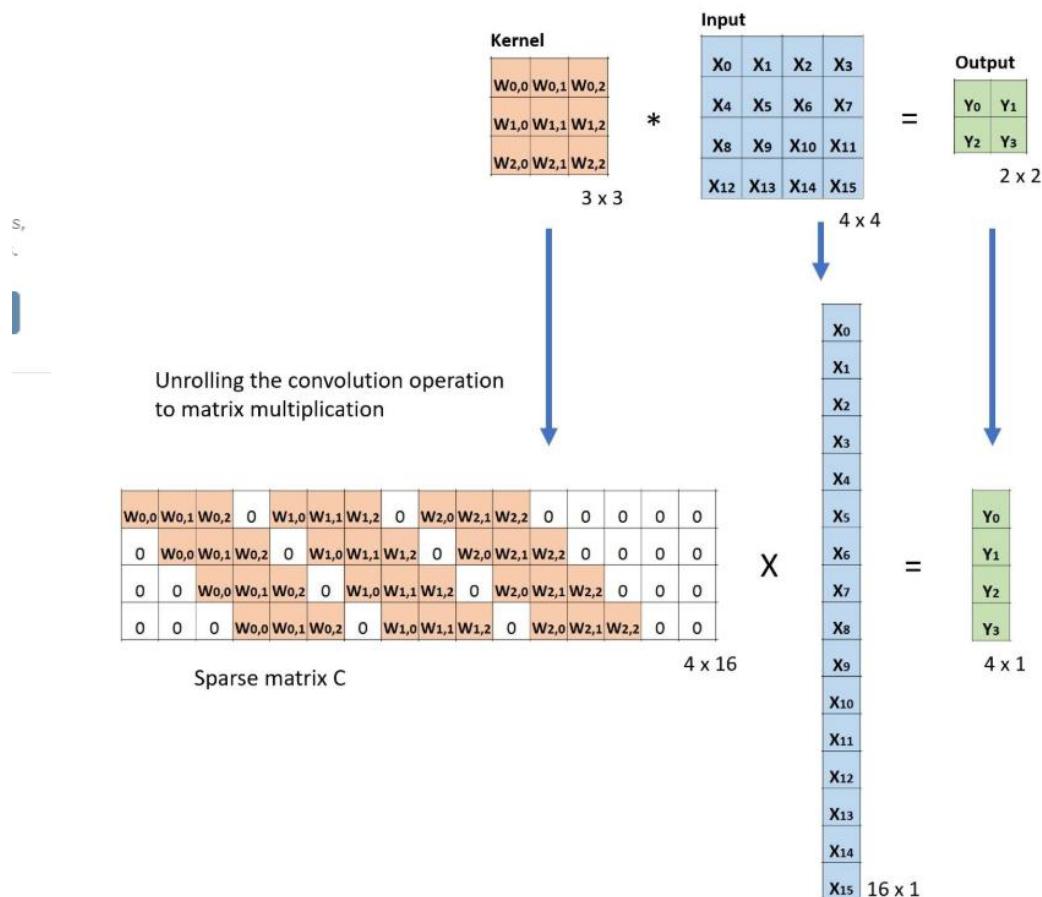
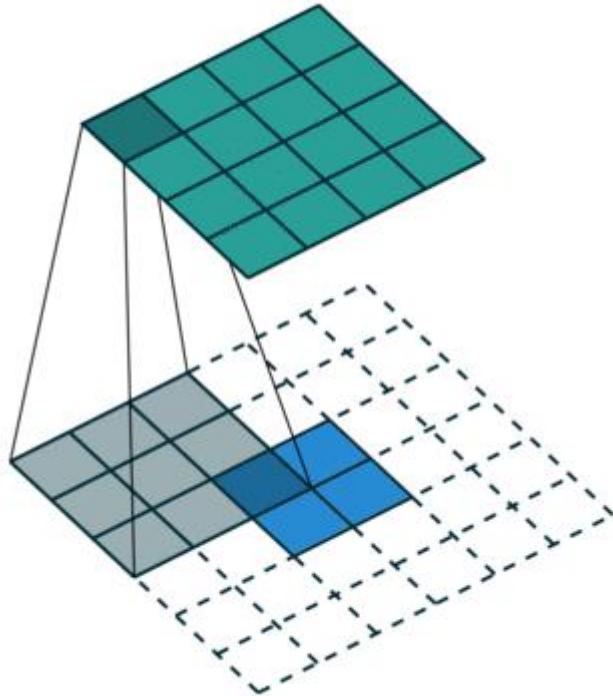
De-Conv(Transpose convolution) layer

convolution: a pixels block into a point --> down sampling

de-convolution: a point to a block --> up sampling, from a low resolution to a higher one

use one point in the previous layer to generate 3*3 blocks with the help of filter

For an example in the image below, we apply transposed convolution with a 3 x 3 kernel over a 2 x 2 input padded with a 2 x 2 border of zeros using unit strides. The up-sampled output is with size 4 x 4.



Matrix multiplication for convolution: from a Large input image (4×4) to a Small output image (2×2).

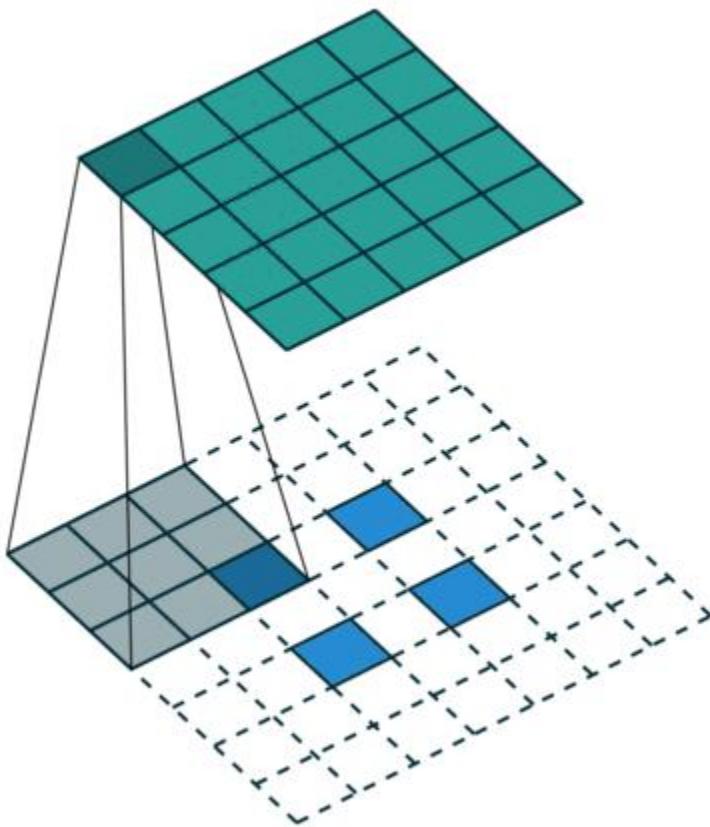
$$\begin{array}{c}
 \begin{array}{|c|c|c|c|} \hline
 W_{0,0} & 0 & 0 & 0 \\ \hline
 W_{0,1} & W_{0,0} & 0 & 0 \\ \hline
 W_{0,2} & W_{0,1} & W_{0,0} & 0 \\ \hline
 0 & W_{0,2} & W_{0,1} & W_{0,0} \\ \hline
 W_{1,0} & 0 & W_{0,2} & W_{0,1} \\ \hline
 W_{1,1} & W_{1,0} & 0 & W_{0,2} \\ \hline
 W_{1,2} & W_{1,1} & W_{1,0} & 0 \\ \hline
 0 & W_{1,2} & W_{1,1} & W_{1,0} \\ \hline
 W_{2,0} & 0 & W_{1,2} & W_{1,1} \\ \hline
 W_{2,1} & W_{2,0} & 0 & W_{1,2} \\ \hline
 W_{2,2} & W_{2,1} & W_{2,0} & 0 \\ \hline
 0 & W_{2,2} & W_{2,1} & W_{2,0} \\ \hline
 0 & 0 & W_{2,2} & W_{2,1} \\ \hline
 0 & 0 & 0 & W_{2,2} \\ \hline
 0 & 0 & 0 & 0 \\ \hline
 \end{array} \\
 \times \quad \quad \quad = \quad \quad \quad \rightarrow \quad \quad \quad 4 \times 1 \quad \quad \quad 16 \times 4
 \end{array}$$

$\begin{array}{|c|} \hline Y_0 \\ \hline Y_1 \\ \hline Y_2 \\ \hline Y_3 \\ \hline \end{array}$ $\begin{array}{|c|} \hline X_0 \\ \hline X_1 \\ \hline X_2 \\ \hline X_3 \\ \hline X_4 \\ \hline X_5 \\ \hline X_6 \\ \hline X_7 \\ \hline X_8 \\ \hline X_9 \\ \hline X_{10} \\ \hline X_{11} \\ \hline X_{12} \\ \hline X_{13} \\ \hline X_{14} \\ \hline X_{15} \\ \hline \end{array}$ $\begin{array}{|c|c|c|c|} \hline X_0 & X_1 & X_2 & X_3 \\ \hline X_4 & X_5 & X_6 & X_7 \\ \hline X_8 & X_9 & X_{10} & X_{11} \\ \hline X_{12} & X_{13} & X_{14} & X_{15} \\ \hline \end{array}$ $4 \times 1 \quad \quad \quad 4 \times 4$

Sparse matrix C^T

Matrix multiplication for convolution: from a Small input image (2×2) to a Large output image (4×4).

map e.g. 4-dimensional space to 25-dimensional space



Transposed 2D convolution with no padding, stride of 2 and kernel of 3

Group Convolution

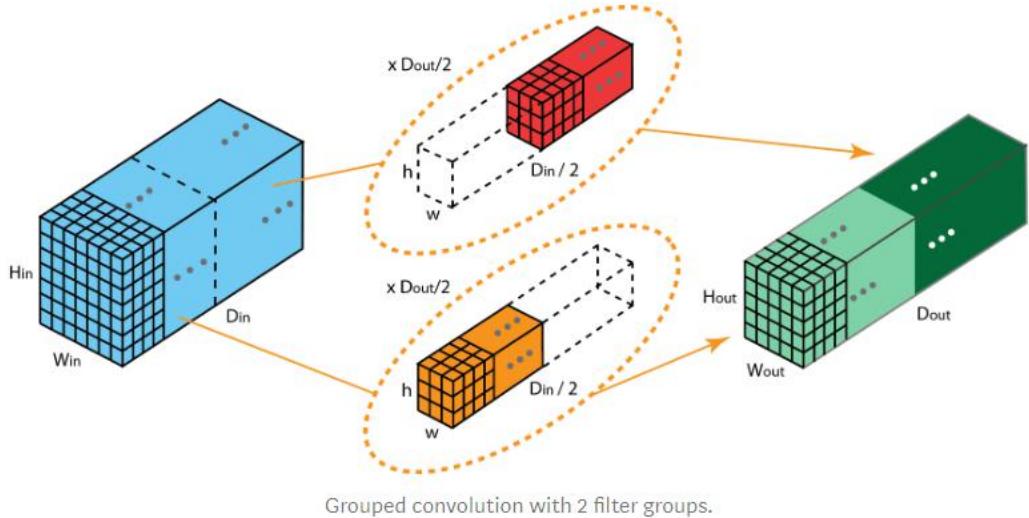
In each filter group, the depth of each filter is only half of that in the nominal 2D convolutions. They are of depth $D_{in} / 2$. Each filter group contains $D_{out} / 2$ filters. The first filter group (red) convolves with the first half of the input layer ($[:, :, 0:D_{in}/2]$), while the second filter group (blue) convolves with the second half of the input layer ($[:, :, D_{in}/2:D_{in}]$). As a result, each filter group creates $D_{out}/2$ channels. Overall, two groups create $2 \times D_{out}/2 = D_{out}$ channels. We then stack these channels in the output layer with D_{out} channels.

Advantage:

efficient training.

model is more efficient

Grouped convolution may provide a better model than a nominal 2D convolution.



this is a bit similar like separable convolutions

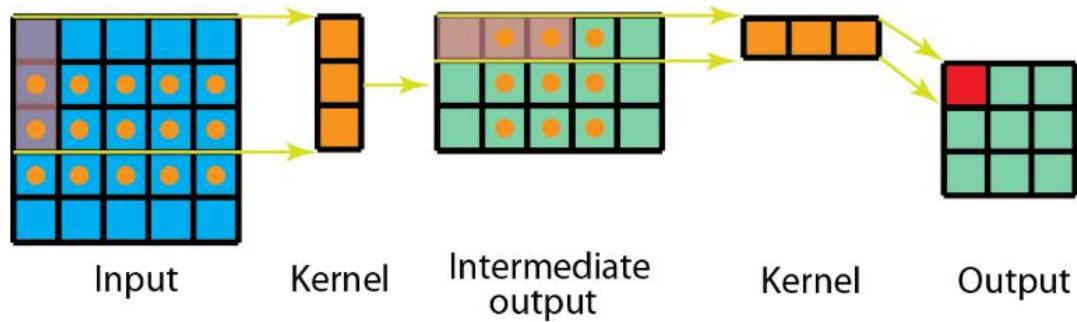
If the number of filter groups is the same as the input layer channel, each filter is of depth $D_{in} / D_{in} = 1$. This is the same filter depth as in depthwise convolution.

Seperable Convolutions

Efficiency!

in MobileNet: <https://arxiv.org/pdf/1704.04861.pdf>

change the 3×3 kernel to 3×1 and 1×3 two kernel This would require 6 instead of 9 parameters while doing the same operations.



Spatially separable convolution with 1 channel.

Although spatially separable convolutions save cost, it is rarely used in deep learning. One of the main reason is that not all kernels can be divided into two, smaller kernels. If we replace all traditional convolutions by the spatially separable convolution, we limit ourselves for searching all possible kernels during training. The training results may be sub-optimal.

Depthwise Separable Convolutions

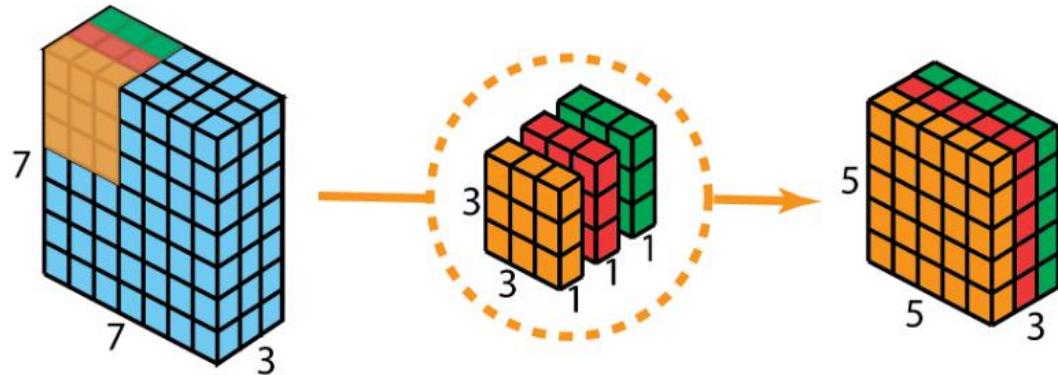
much commonly used in deep learning: **depthwise convolution** and **1×1 convolution**

First, we apply **depthwise convolution** to the input layer. Instead of using a single filter of size $3 \times 3 \times 3$ in 2D convolution, we used **3 kernels, separately**. Each filter has size $3 \times 3 \times 1$. Each kernel convolves with 1 channel of the input layer (1 channel only, not all channels!). We then

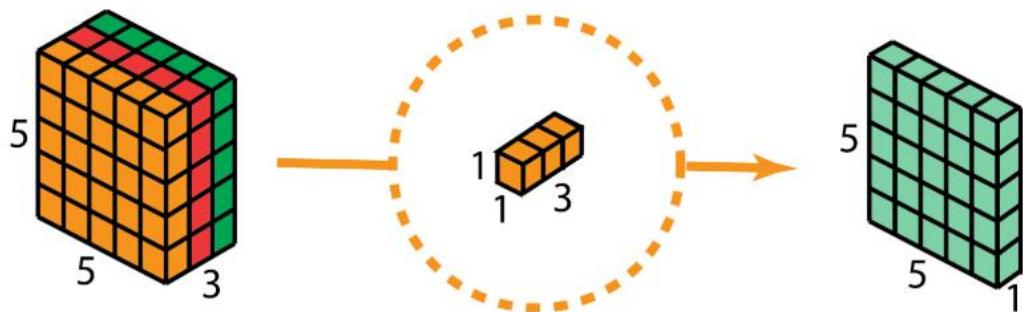
stack these maps together to create a $5 \times 5 \times 3$ image. After this, we have the output with size $5 \times 5 \times 3$.

As the second step of depthwise separable convolution, to extend the depth, we apply the 1×1 convolution with kernel size $1 \times 1 \times 3$. Convolving the $5 \times 5 \times 3$ input image with each $1 \times 1 \times 3$ kernel provides a map of size $5 \times 5 \times 1$.

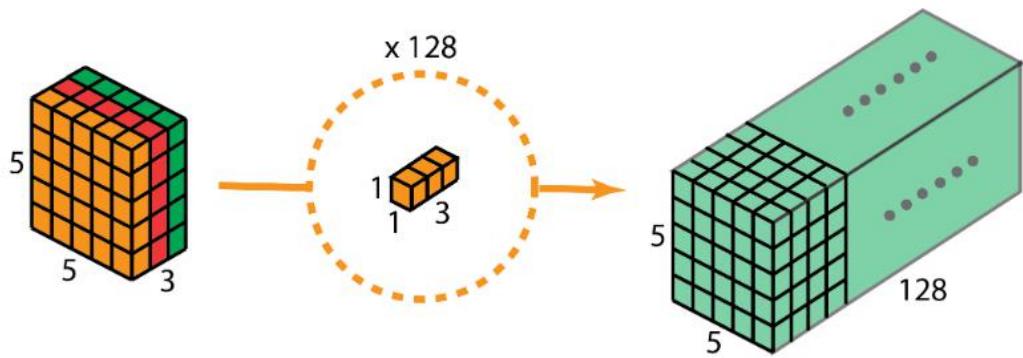
Thus, after applying 128 1×1 convolutions, we can have a layer with size $5 \times 5 \times 128$.



Depthwise separable convolution – first step: Instead of using a single filter of size $3 \times 3 \times 3$ in 2D convolution, we used 3 kernels, separately. Each filter has size $3 \times 3 \times 1$. Each kernel convolves with 1 channel of the input layer (1 channel only, not all channels!). Each of such convolution provides a map of size $5 \times 5 \times 1$. We then stack these maps together to create a $5 \times 5 \times 3$ image. After this, we have the output with size $5 \times 5 \times 3$.



Thus, after applying 128 1x1 convolutions, we can have a layer with size 5 x 5 x 128.



Depthwise separable convolution — second step: apply multiple 1x1 convolutions to modify depth.

drwaback:

reduces the number of parameters in the convolution. As such, for a small model, the model capacity may be decreased significantly if the 2D convolutions are replaced by depthwise separable convolutions.

Full Connect layer

of Parameters computation

Visualization

layers that are deeper in the network visualize more training data specific features, while the earlier layers tend to visualize general patterns like edges, texture, background

Visualize different layer

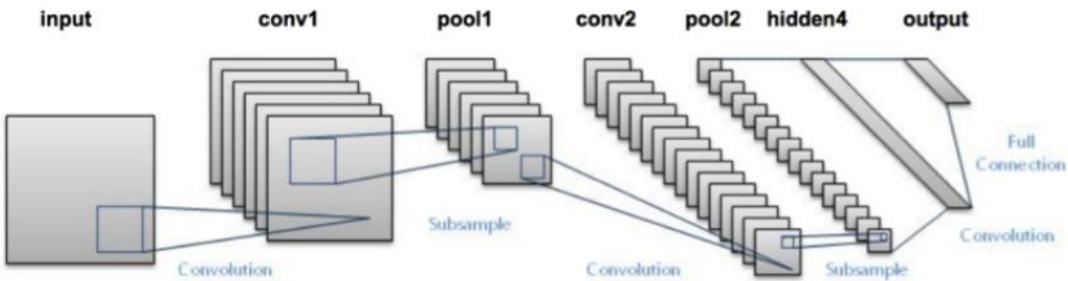
Visualize filter

Heatmap

Different Nets

LeNet(1998)

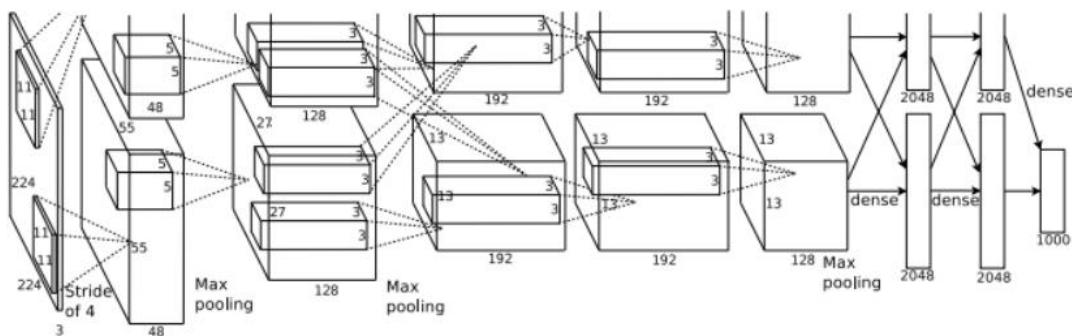
classifies digits in 32x32 pixel greyscale input images



AlexNet(2012)

deeper, with more filters per layer, and with stacked convolutional layers. It consisted 11x11, 5x5, 3x3, convolutions, max pooling, dropout, data augmentation, ReLU activations, SGD with momentum.

train on two GPUs

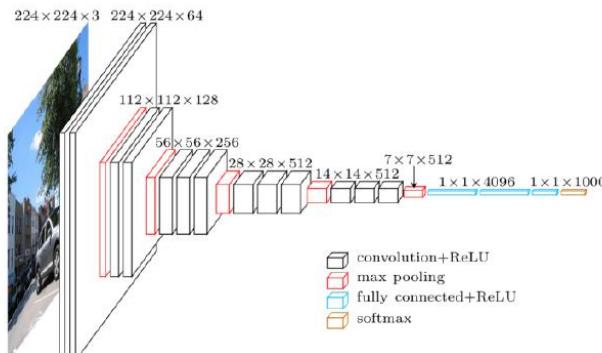


VGG

<https://arxiv.org/pdf/1409.1556.pdf>

Baseline feature extractor

16 (weighted) layers CNN 138M parameters



Original training: 3 weeks on 4 GPUs

Name	A	A-LRN	B	C	D	E
# Layers	11	11	13	16	16	19
C3D64	C3D64	C3D64	C3D64	C3D64	C3D64	C3D64
LRN	LRN	C3D64	C3D64	C3D64	C3D64	C3D64
M	M	M	M	M	M	M
C3D128	C3D128	C3D128	C3D128	C3D128	C3D128	C3D128
M	M	M	M	M	M	M
C3D256	C3D256	C3D256	C3D256	C3D256	C3D256	C3D256
C3D256	C3D256	C3D256	C3D256	C3D256	C3D256	C3D256
M	M	M	M	M	M	M
C3D512	C3D512	C3D512	C3D512	C3D512	C3D512	C3D512
C3D512	C3D512	C3D512	C3D512	C3D512	C3D512	C3D512
M	M	M	M	M	M	M
C3D512	C3D512	C3D512	C3D512	C3D512	C3D512	C3D512
C3D512	C3D512	C3D512	C3D512	C3D512	C3D512	C3D512
M	M	M	M	M	M	M
C3D512	C3D512	C3D512	C3D512	C3D512	C3D512	C3D512
C3D512	C3D512	C3D512	C3D512	C3D512	C3D512	C3D512
M	M	M	M	M	M	M
FC4096	FC4096	FC4096	FC4096	FC4096	FC4096	FC4096
FC4096	FC4096	FC4096	FC4096	FC4096	FC4096	FC4096
FC1000	FC1000	FC1000	FC1000	FC1000	FC1000	FC1000
S	S	S	S	S	S	S

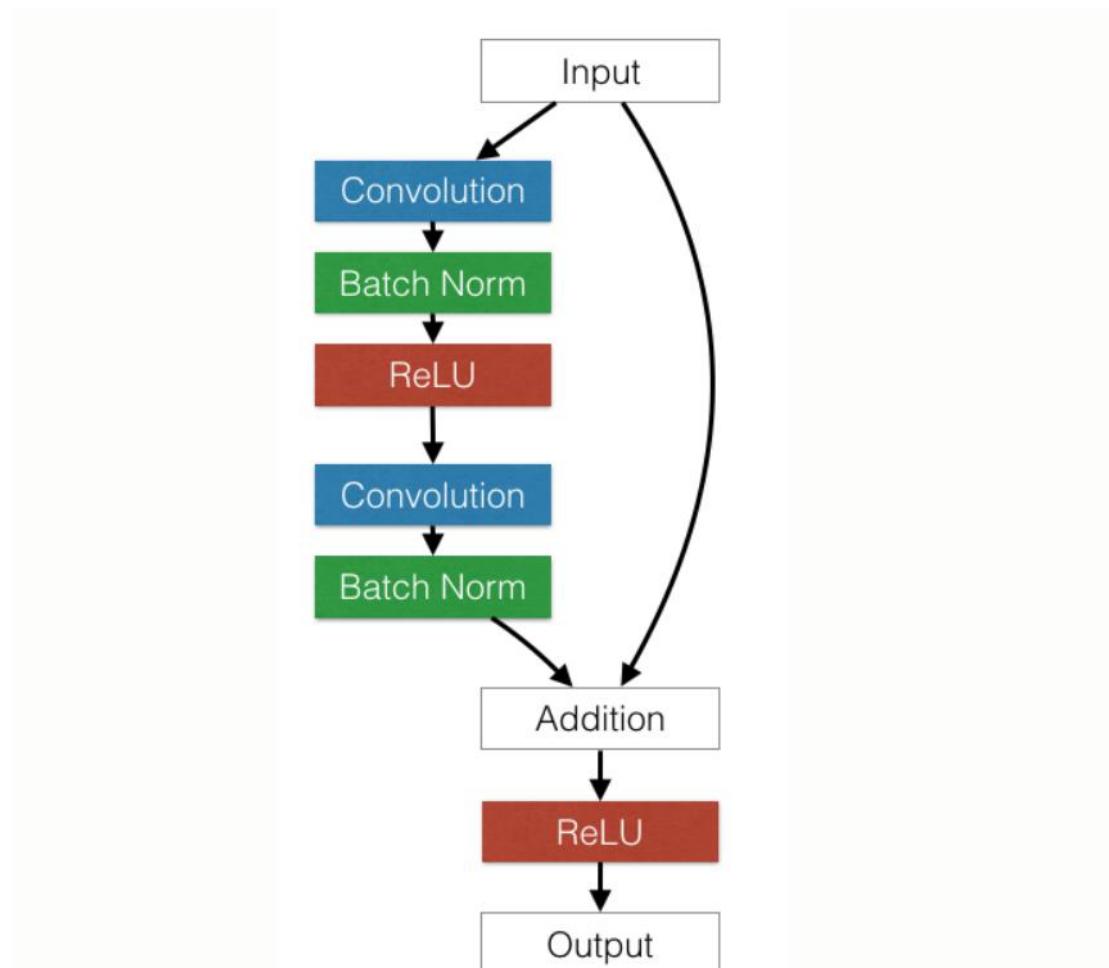
parameters ≈ 138M

Resnet

gated recurrent units

take a standard feed-forward ConvNet and add **skip connections** that bypass (or shortcut) a few convolution layers at a time. Each bypass gives rise to a residual block in which the convolution layers predict a residual that is added to the block's input tensor.

aim to avoid vanish gradient



GoogleNet (Inception)

<https://arxiv.org/pdf/1512.00567.pdf>

<https://arxiv.org/pdf/1409.4842.pdf>

no Pooling layer

Introduction

The network used a CNN inspired by LeNet but implemented a novel element which is dubbed an inception module. It used batch normalization, image distortions and RMSprop. This module is based on several very small convolutions in order to drastically reduce the number of parameters. Their architecture consisted of a 22 layer deep CNN but reduced the number of parameters from 60 million (AlexNet) to 4 million.

Speciality

1. use several small filter to stand large filter --> $n*1 + 1*n$ to replace $n*n$
2. for a single input layer, applied many different filters(some are pooling, some are $1*1$ with $1*n$ with $n*1$) and concatenated result (add together)

--> **to avoid representational bottlenecks and avoid stop locally**

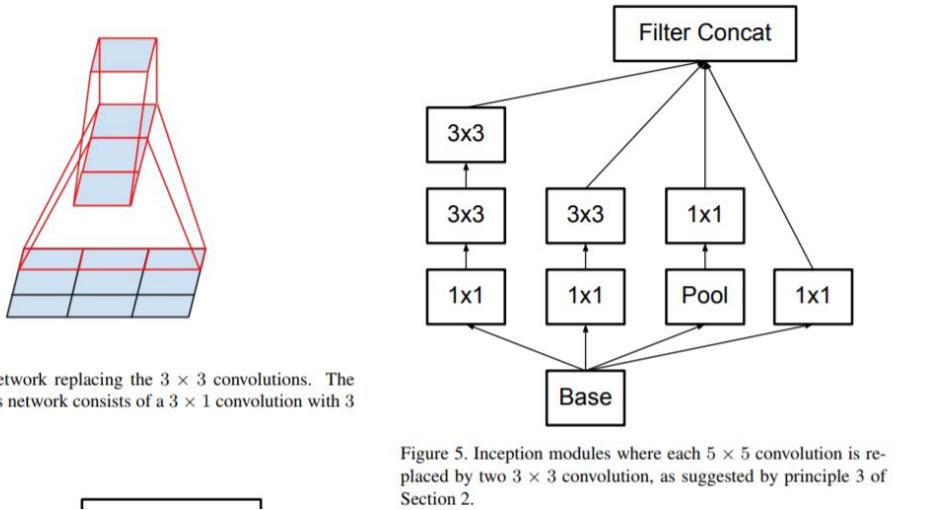
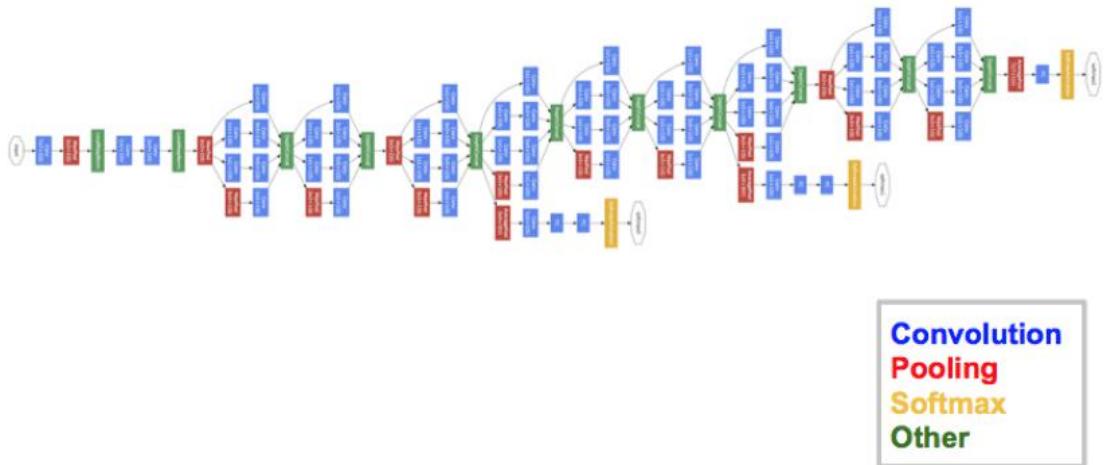


Figure 5. Inception modules where each 5×5 convolution is replaced by two 3×3 convolution, as suggested by principle 3 of Section 2.



MobileNet

<https://arxiv.org/pdf/1704.04861.pdf>

Deepwise Separable Convolution + Pointwise Convolution (1×1)

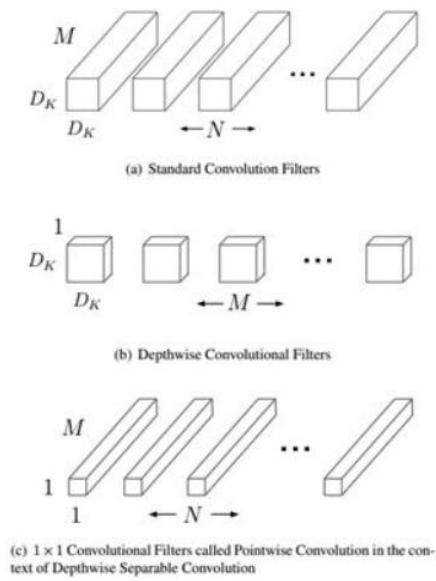


图1 Depthwise separable convolution

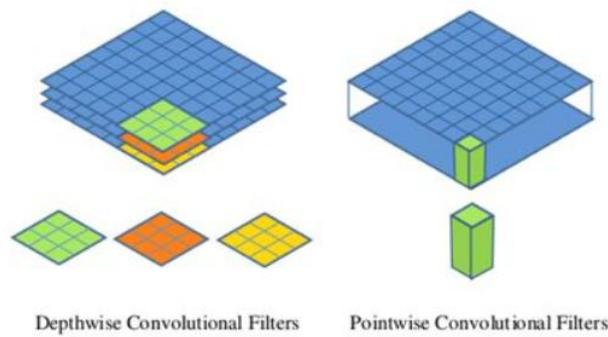


图2 Depthwise convolution和pointwise convolution

MobileNet网络结构

前面讲述了depthwise separable convolution，这是MobileNet的基本组件，但是在真正应用中会加入batchnorm，并使用ReLU激活函数，所以depthwise separable convolution的基本结构如图3所示。

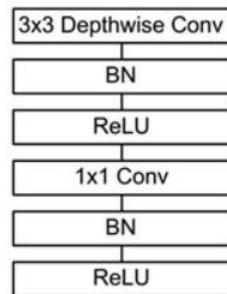


图3 加入BN和ReLU的depthwise separable convolution

MobileNet 的网络结构如表 1 所示。首先是一个 3×3 的标准卷积，然后后面就是堆积 depthwise separable convolution，并且可以看到其中的部分 depthwise convolution 会通过 strides=2 进行 down sampling。然后采用 average pooling 将 feature 变成 1×1 ，根据预测类别大小加上全连接层，最后是一个 softmax 层。如果单独计算 depthwise

convolution 和 pointwise convolution，整个网络有 28 层（这里 Avg Pool 和 Softmax 不计算在内）。我们还可以分析整个网络的参数和计算量分布，如表 2 所示。可以看到整个计算量基本集中在 1×1 卷积上，如果你熟悉卷积底层实现的话，你应该知道卷积一般通过一种 im2col 方式实现，其需要内存重组，但是当卷积核为 1×1 时，其实就不需要这种操作了，底层可以有更快的实现。对于参数也主要集中在 1×1 卷积，除此之外还有就是全连接层占了一部分参数。

MobileNet网络结构

前面讲述了depthwise separable convolution，这是MobileNet的基本组件，但是在真正应用中会加入batchnorm，并使用ReLU激活函数，所以depthwise separable convolution的基本结构如图3所示。

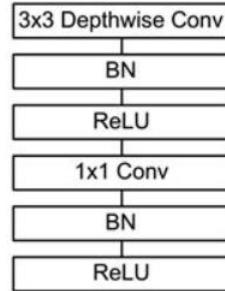


图3 加入BN和ReLU的depthwise separable convolution

Type / Stride	Filter Shape	Input Size
Conv / s2	$3 \times 3 \times 3 \times 32$	$224 \times 224 \times 3$
Conv dw / s1	$3 \times 3 \times 32$ dw	$112 \times 112 \times 32$
Conv / s1	$1 \times 1 \times 32 \times 64$	$112 \times 112 \times 32$
Conv dw / s2	$3 \times 3 \times 64$ dw	$112 \times 112 \times 64$
Conv / s1	$1 \times 1 \times 64 \times 128$	$56 \times 56 \times 64$
Conv dw / s1	$3 \times 3 \times 128$ dw	$56 \times 56 \times 128$
Conv / s1	$1 \times 1 \times 128 \times 128$	$56 \times 56 \times 128$
Conv dw / s2	$3 \times 3 \times 128$ dw	$56 \times 56 \times 128$
Conv / s1	$1 \times 1 \times 128 \times 256$	$28 \times 28 \times 128$
Conv dw / s1	$3 \times 3 \times 256$ dw	$28 \times 28 \times 256$
Conv / s1	$1 \times 1 \times 256 \times 256$	$28 \times 28 \times 256$
Conv dw / s2	$3 \times 3 \times 256$ dw	$28 \times 28 \times 256$
Conv / s1	$1 \times 1 \times 256 \times 512$	$14 \times 14 \times 256$
5× Conv dw / s1	$3 \times 3 \times 512$ dw	$14 \times 14 \times 512$
Conv / s1	$1 \times 1 \times 512 \times 512$	$14 \times 14 \times 512$
Conv dw / s2	$3 \times 3 \times 512$ dw	$14 \times 14 \times 512$
Conv / s1	$1 \times 1 \times 512 \times 1024$	$7 \times 7 \times 512$
Conv dw / s2	$3 \times 3 \times 1024$ dw	$7 \times 7 \times 1024$
Conv / s1	$1 \times 1 \times 1024 \times 1024$	$7 \times 7 \times 1024$
Avg Pool / s1	Pool 7×7	$7 \times 7 \times 1024$
FC / s1	1024×1000	$1 \times 1 \times 1024$
Softmax / s1	Classifier	$1 \times 1 \times 1000$

表1 MobileNet的网络结构

Comparation

models in the context of image classification. We measured the time for a full forward and backward pass for a mini-batch of 32 images for ResNet, VGG A, VGG D, Batch-Normalized Inception, and Inception v3 on a NVIDIA Titan X. Also listed is the top-1 single-crop validation error on the Imagenet-2012 dataset.

Model	Top-1 err (%)	Time (ms)
VGG-A	29.6	372
VGG-D	26.8	687
ResNet-34	26.7	231
BN-Inception	25.2	192
ResNet-50	24.0	403
ResNet-101	22.4	649
Inception-v3	21.2	494

CNN in NLP

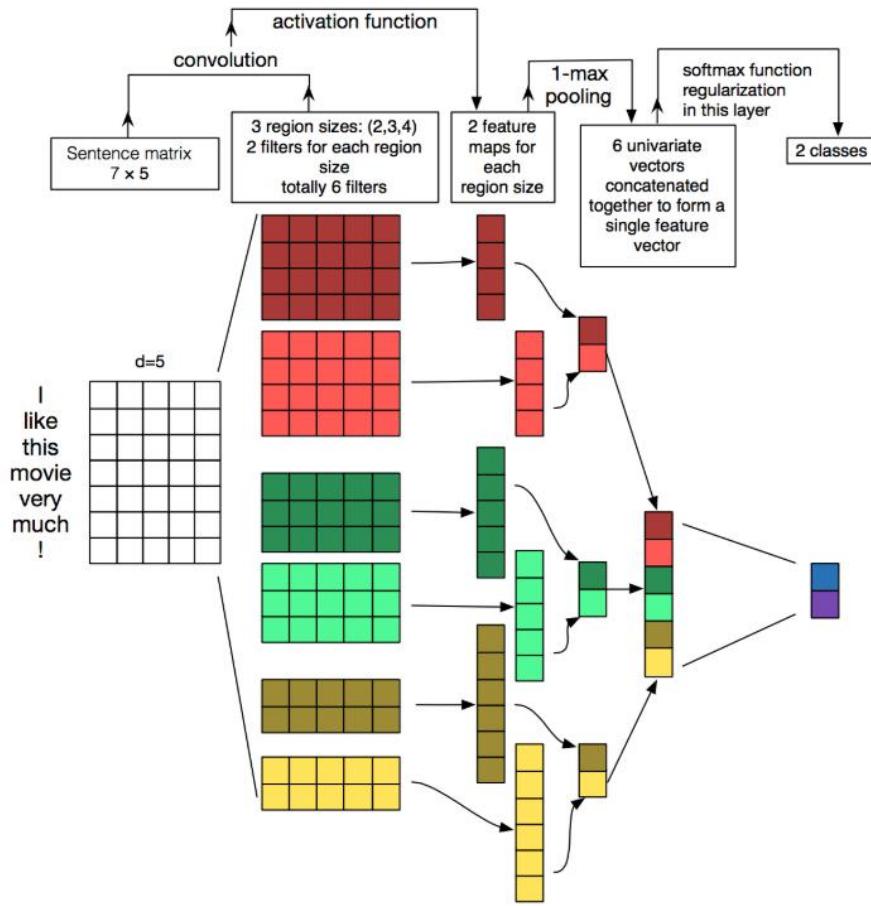


Image Reference : <http://www.wildml.com/2015/11/understanding-convolutional-neural-networks-for-nlp/>

Coding

Keras

a basic model

```
from keras.models import Sequential from keras.layers import Dense, Conv2D, Flatten# create model
model = Sequential()#add model layers
model.add(Conv2D(64, kernel_size=3, activation='relu', input_shape=(28,28,1)))
model.add(Conv2D(32, kernel_size=3, activation='relu'))
model.add(Flatten())
model.add(Dense(10, activation='softmax')) # Compiling the model takes three parameters: optimizer, loss and metrics.
model.compile(optimizer='adam', loss='categorical_crossentropy', metrics=['accuracy']) #train the model
model.fit(X_train, y_train, validation_data=(X_test, y_test), epochs=3) #predict first 4 images in the test set
model.predict(X_test[:4]) model = Sequential()
model.add(Conv2D(32, (3, 3), padding='same', input_shape=x_train.shape[1:]))
model.add(Activation('relu'))
model.add(Conv2D(32, (3, 3)))
model.add(Activation('relu'))
model.add(MaxPooling2D(pool_size=(2, 2)))
model.add(Dropout(0.25))
model.add(Conv2D(64, (3, 3), padding='same'))
model.add(Activation('relu'))
model.add(Conv2D(64, (3, 3)))
model.add(Activation('relu'))
model.add(MaxPooling2D(pool_size=(2, 2)))
model.add(Dropout(0.25))
model.add(Flatten())
```

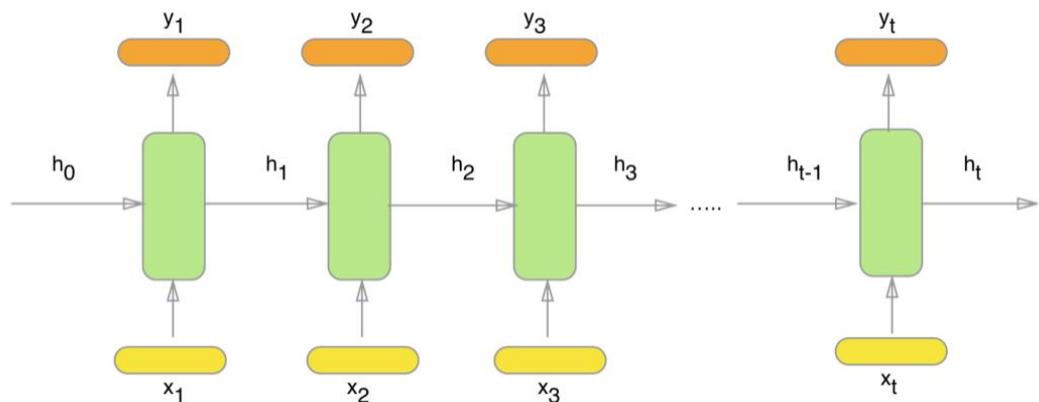
```
model.add(Dense(512)) model.add(Activation('relu')) model.add(Dropout(0.5))
model.add(Dense(num_classes)) model.add(Activation('softmax'))
```

'Flatten' layer. Flatten serves as a connection between the convolution and dense layers.

RNN

Introduction

$x_1, x_2, x_3, \dots, x_t$ represent the input words from the text, $y_1, y_2, y_3, \dots, y_t$ represent the predicted next words and $h_0, h_1, h_2, h_3, \dots, h_t$ hold the information for the previous input words.



$$1) \quad h_t = f(W^{(hh)}h_{t-1} + W^{(hx)}x_t)$$

$$2) \quad y_t = \text{softmax}(W^{(S)}h_t)$$

$$3) \quad J^{(t)}(\theta) = \sum_{i=1}^{|V|} (y_{ti}' \log y_{ti})$$

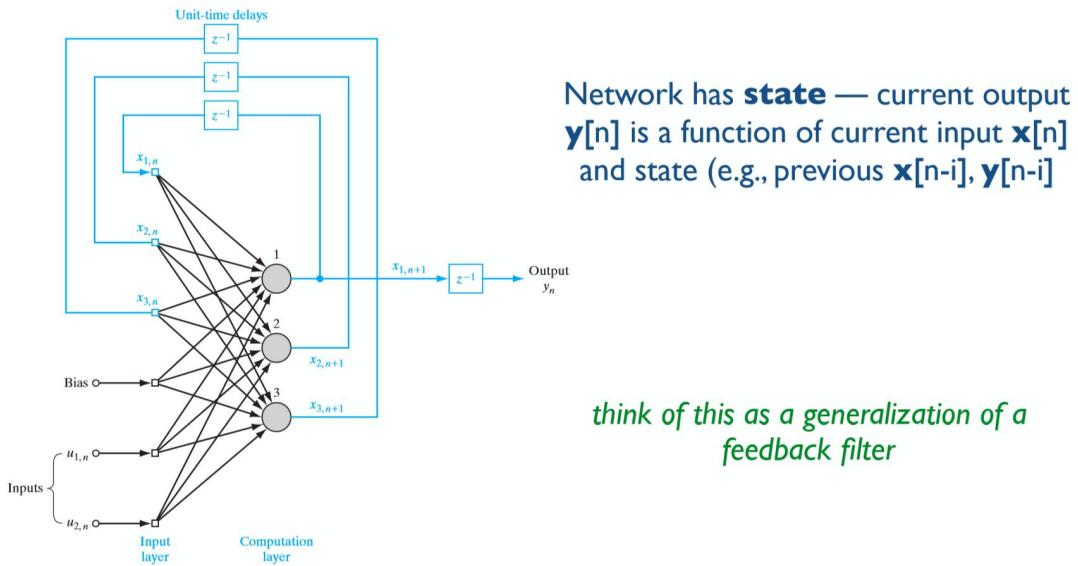
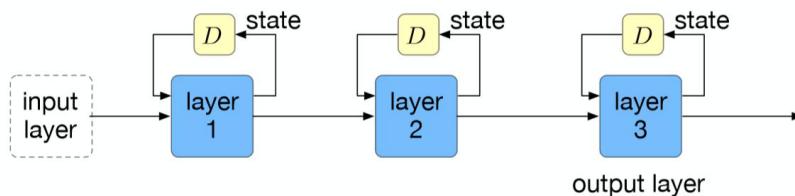


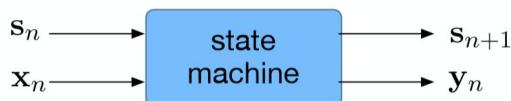
FIGURE 15.6 Fully connected recurrent network with two inputs, two hidden neurons, and one output neuron. The feedback connections are shown in red to emphasize their global role.

basic structure



STATE

state machine: network is stated:



$$s_{n+1} = \text{next_state}(s_n, x_n)$$

$$y_n = \text{output}(s_n, x_n)$$

these two functions define the state machine

s_n state at time n (before input x_n is applied)

x_n input at time n

s_{n+1} state at time $n+1$ (after input x_n is applied)

y_n output at time n

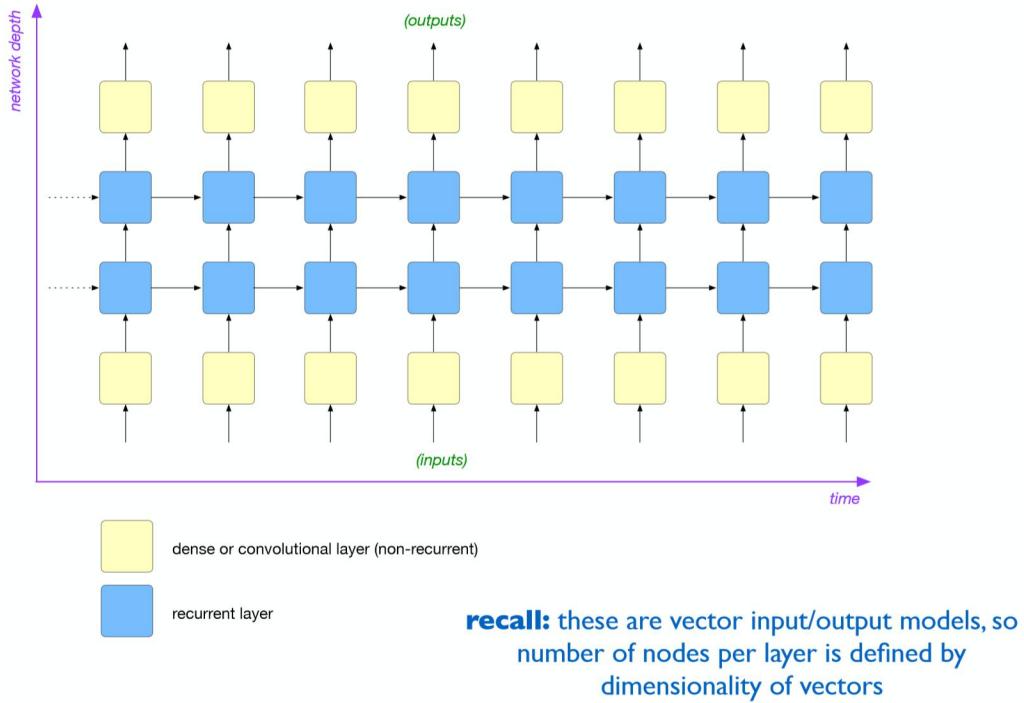
(Vanilla RNN)

Keras code

```
keras.layers.SimpleRNN(units, activation='tanh', use_bias=True,
kernel_initializer='glorot_uniform', recurrent_initializer='orthogonal', bias_initializer='zeros',
kernel_regularizer=None, recurrent_regularizer=None, bias_regularizer=None,
activity_regularizer=None, kernel_constraint=None, recurrent_constraint=None,
```

bias_constraint=None, dropout=0.0, recurrent_dropout=0.0, return_sequences=False, return_state=False, go_backwards=False, stateful=False, unroll=False)

Diagram of Neural Net



Number of nodes and parameters in each layer (blue block on upper diagram)

$$\mathbf{s}_{n+1} = h(\mathbf{Vs}_n + \mathbf{Wx}_n + \mathbf{b})$$

`keras.layers.SimpleRNN(units)`

$$\mathbf{y}_n = \mathbf{s}_{n+1}$$

$$\begin{array}{c} \mathbf{y}_n \\ N \times 1 \end{array} \quad \begin{array}{c} \mathbf{W} \\ N \times M \end{array} \quad \begin{array}{c} \mathbf{x}_n \\ M \times 1 \end{array}$$

$$\begin{array}{c} \mathbf{s}_{n+1} \\ N \times 1 \end{array} \quad \begin{array}{c} \mathbf{V} \\ N \times N \end{array} \quad \begin{array}{c} \mathbf{s}_n \\ N \times 1 \end{array}$$

N: number of “units” (nodes) in this recurrent layer

M: number of “units” (nodes) in the previous layer

Node on layer: n

Node on previous layer: m

input: m

output: n

parameters: $(m+n)*n$ (weight) + n (bias)

When training, a training window length T is selected =>

a sequence of input of length T

a sequence of label of length T

<http://www.wildml.com/2015/09/recurrent-neural-networks-tutorial-part-1-introduction-to-rnns/>

Vanish Gradient Problem

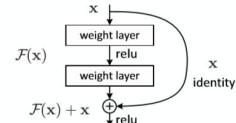
Vanishing gradient is a training-based point of view of a more fundamental issue.

The real issue is **decaying influence** or **degradation**

basically, inputs to repeated linear-nonlinear operations cannot affect outputs far into the “future”

Feedforward networks: “future” means deeper in the network

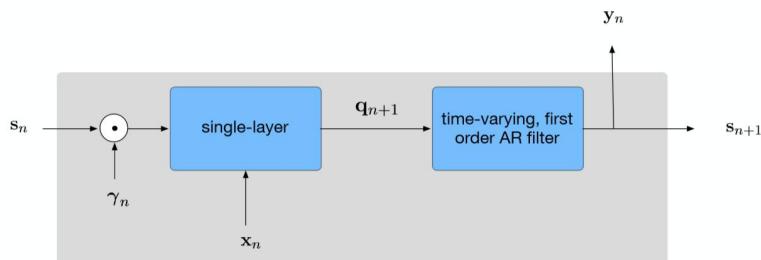
“degradation” coined by ResNet authors and Residual connections introduced to go deeper



GATE

To solve this problem we need to add GATE which is attenuating and/or filtering in the state update equation (my understanding: amplify the influence of previous state)

Generic RNN with Gates



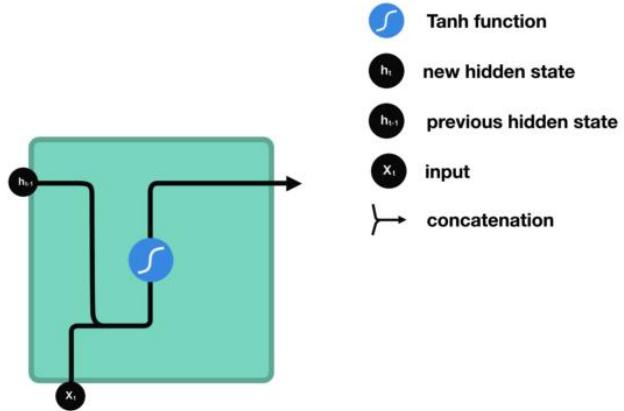
$$\mathbf{q}_{n+1} = h(\mathbf{V}[\gamma_n \odot \mathbf{s}_n] + \mathbf{W}\mathbf{x}_n + \mathbf{b}) \quad \mathbf{s}_{n+1} = \alpha_n \odot \mathbf{s}_n + \beta_n \odot \mathbf{q}_{n+1}$$

γ_n “read gate” — attenuation on current state when read for update

α_n “forget gate” — feedback coefficients for state filtering

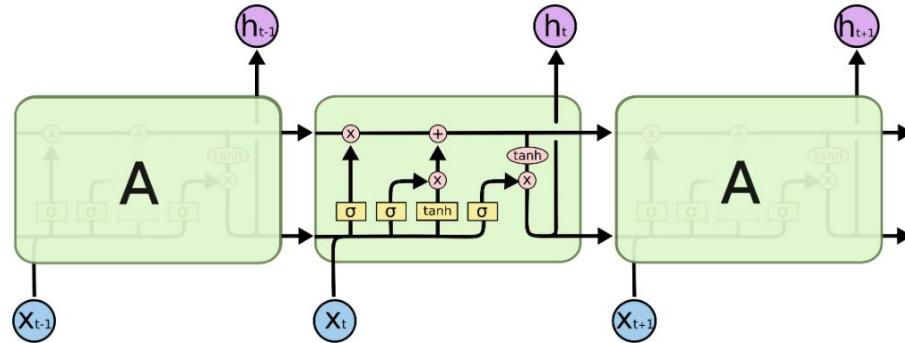
β_n “write gate” — feedforward coefficients for state filtering

All gate are trainable parameters and are learned using a single layer feedforward network (my understanding: GATE make simple $m \times n$ parameters network more complex, though if we block the process, it is still a m input and n output problem, the inner parameters (weight) are no longer $n(m+n+1)$, it because more complex, but still based on $\mathbf{V}(n \times n)$, $\mathbf{W}(m \times n)$, \mathbf{b} (n)



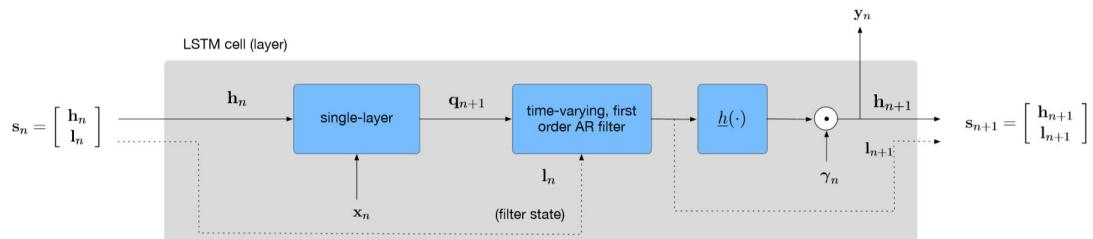
tanh make the output in between -1 → 1

LSTM



1

Long/Short-Term Memory (LSTM) Cell



$$\mathbf{q}_{n+1} = \underline{h}(\mathbf{V}\mathbf{h}_n + \mathbf{W}\mathbf{x}_n + \mathbf{b})$$

$$\mathbf{h}_{n+1} = \gamma_n \odot \underline{h}(\mathbf{l}_{n+1})$$

$$\mathbf{l}_{n+1} = \alpha_n \odot \mathbf{l}_n + \beta_n \odot \mathbf{q}_{n+1}$$

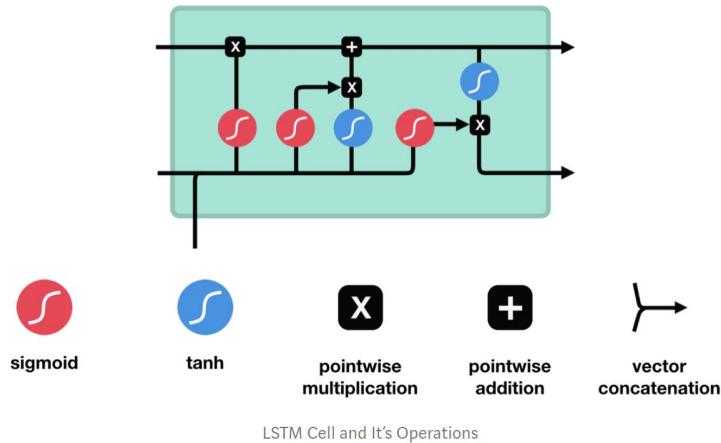
Differences from GRU:

$$\mathbf{s}_n = \begin{bmatrix} \mathbf{h}_n \\ \mathbf{l}_n \end{bmatrix} = \begin{bmatrix} \text{short-term state} \\ \text{long-term state} \end{bmatrix}$$

“forget gate” and “write gate” decoupled ~ different alpha and beta

“read gate” applied on output (gates read of next cell at n+1)

activation applied on the long-term state before applying read gate



These operations are used to allow the LSTM to keep or forget information.

Core Concept

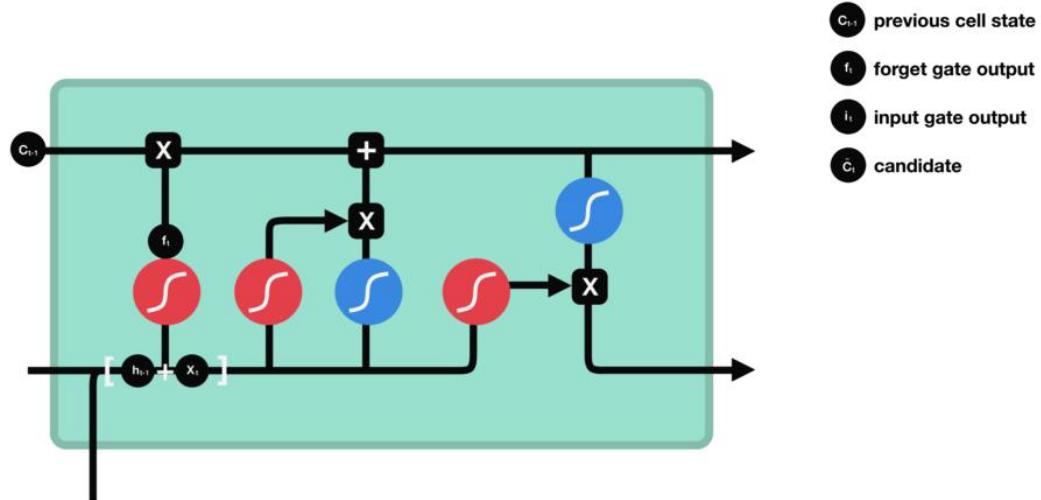
The core concept of LSTM's are the cell state, and it's various gates. The cell state act as a transport highway that transfers relative information all the way down the sequence chain. You can think of it as the "memory" of the network. The cell state, in theory, can carry relevant information throughout the processing of the sequence. So even information from the earlier time steps can make it's way to later time steps, reducing the effects of short-term memory. As the cell state goes on its journey, information get's added or removed to the cell state via gates. The gates are different neural networks that decide which information is allowed on the cell state. The gates can learn what information is relevant to keep or forget during training. there are two state value(cell and hiden in each state)

Forget Gate

This gate decides what information should be thrown away or kept. Information from the previous hidden state and information from the current input is passed through the sigmoid function. Values come out between 0 and 1. The closer to 0 means to forget, and the closer to 1 means to keep.

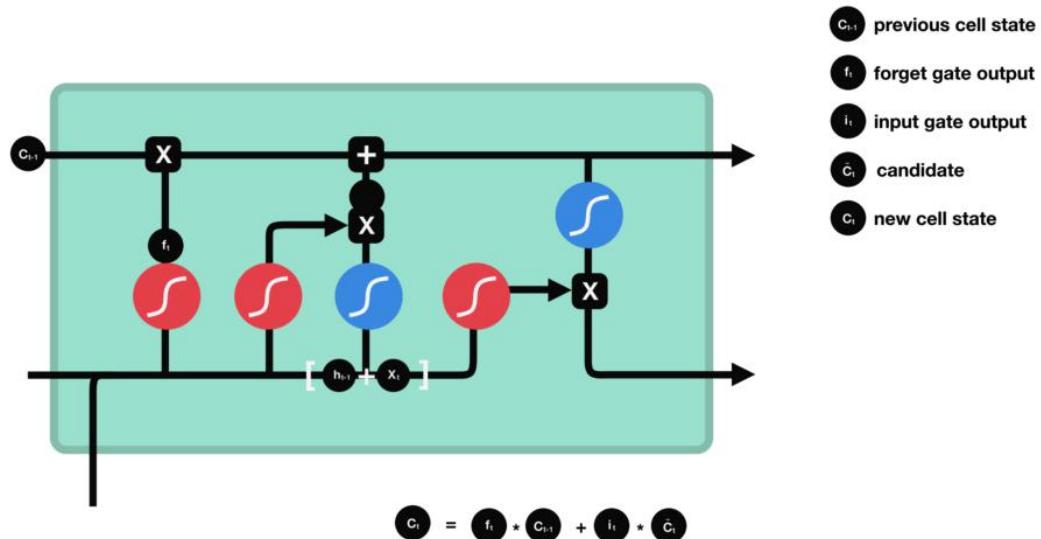
Input Gate

First, we pass the previous hidden state and current input into a sigmoid function. That decides which values will be updated by transforming the values to be between 0 and 1. 0 means not important, and 1 means important. You also pass the hidden state and current input into the tanh function to squish values between -1 and 1 to help regulate the network. Then you multiply the tanh output with the sigmoid output. The sigmoid output will decide which information is important to keep from the tanh output.



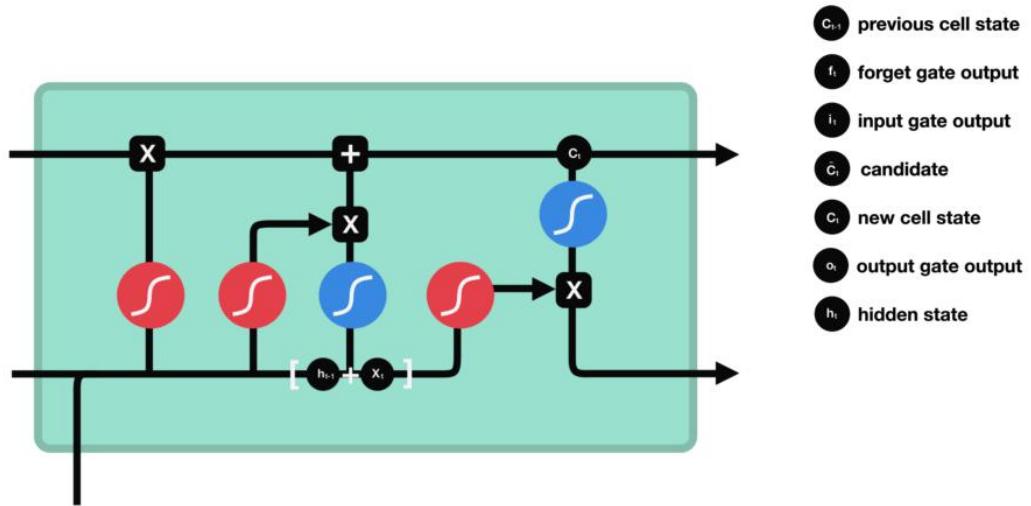
Cell Gate

First, the cell state gets pointwise multiplied by the forget vector. This has a possibility of dropping values in the cell state if it gets multiplied by values near 0. Then we take the output from the input gate and do a pointwise addition which updates the cell state to new values that the neural network finds relevant. That gives us our new cell state.

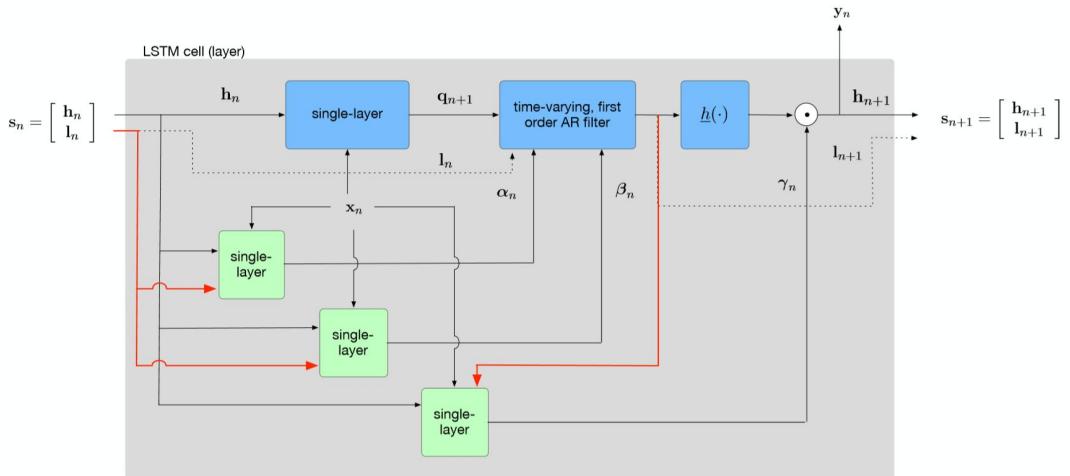


Output Gate

The output gate decides what the next hidden state should be. Remember that the hidden state contains information on previous inputs. The hidden state is also used for predictions. First, we pass the previous hidden state and the current input into a sigmoid function. Then we pass the newly modified cell state to the tanh function. We multiply the tanh output with the sigmoid output to decide what information the hidden state should carry. The output is the hidden state. The new cell state and the new hidden is then carried over to the next time step.



OTHER VERSION(Keras)



$$\alpha_n = \underline{h}_r(\mathbf{V}_\alpha \mathbf{h}_n + \mathbf{P}_\alpha \mathbf{l}_n + \mathbf{W}_\alpha \mathbf{x}_n + \mathbf{b}_\alpha)$$

$$\beta_n = \underline{h}_r(\mathbf{V}_\beta \mathbf{h}_n + \mathbf{P}_\beta \mathbf{l}_n + \mathbf{W}_\beta \mathbf{x}_n + \mathbf{b}_\beta)$$

$$\gamma_n = \underline{h}_r(\mathbf{V}_\gamma \mathbf{h}_n + \mathbf{P}_\gamma \mathbf{l}_{n+1} + \mathbf{W}_\gamma \mathbf{x}_n + \mathbf{b}_\gamma)$$

(this is an experimental layer type in keras)

or sometimes

$$\alpha_n = \underline{h}_r(\mathbf{V}_\alpha \mathbf{h}_n + \mathbf{p}_\alpha \odot \mathbf{l}_n + \mathbf{W}_\alpha \mathbf{x}_n + \mathbf{b}_\alpha)$$

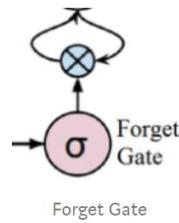
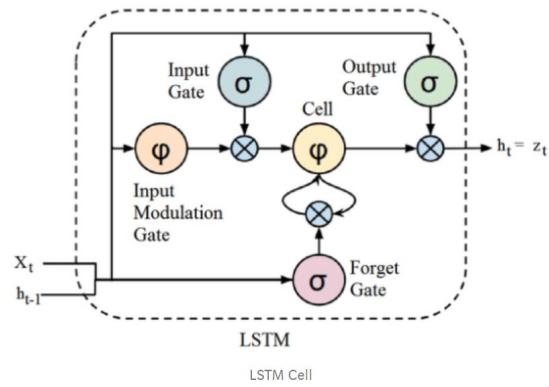
$$\beta_n = \underline{h}_r(\mathbf{V}_\beta \mathbf{h}_n + \mathbf{p}_\beta \odot \mathbf{l}_n + \mathbf{W}_\beta \mathbf{x}_n + \mathbf{b}_\beta)$$

$$\gamma_n = \underline{h}_r(\mathbf{V}_\gamma \mathbf{h}_n + \mathbf{p}_\gamma \odot \mathbf{l}_{n+1} + \mathbf{W}_\gamma \mathbf{x}_n + \mathbf{b}_\gamma)$$

the peep-hole matrices are diagonal

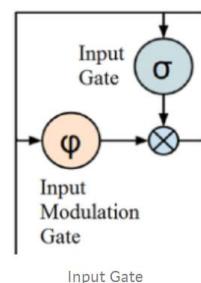
My understanding: In LSTM, each gate has a group of parameters, so four gates, four group of parameters, each group will generate an output, and they combine together to generate the final output and hidden state so the Keras version will use more parameters ($4 * n * (m+n+1)$)

other version of intorduction ()



$$f_t = \sigma(\mathbf{W}_f[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_f)$$

Forget Gate Equation

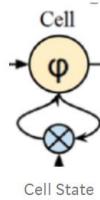


$$i_t = \sigma(\mathbf{W}_i[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_i)$$

Input Gate Equation

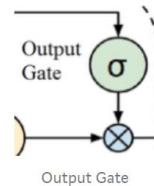
$$\tilde{c}_t = \tanh(\mathbf{W}_c[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_c)$$

Input Modulation Gate Equation



$$\mathbf{c}_t = \mathbf{f}_t \circ \mathbf{c}_{t-1} + \mathbf{i}_t \circ \tilde{\mathbf{c}}_t$$

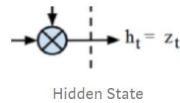
Cell State Equation



$$\mathbf{o}_t = \sigma(\mathbf{W}_o[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_o)$$

Output Gate Equation

The **focus vector** is usually called the **output gate**. Out of all the possible values from the matrix, which should be moving forward to the next hidden state?



$$\mathbf{h}_t = \mathbf{o}_t \circ \tanh(\mathbf{c}_t)$$

Hidden State Equation

LSTM code:

useful source: [Illustrated Guide to LSTM's and GRU's: A step by step explanation]

<https://towardsdatascience.com/illustrated-guide-to-lstms-and-gru-s-a-step-by-step-explanation-44e9eb85bf21>

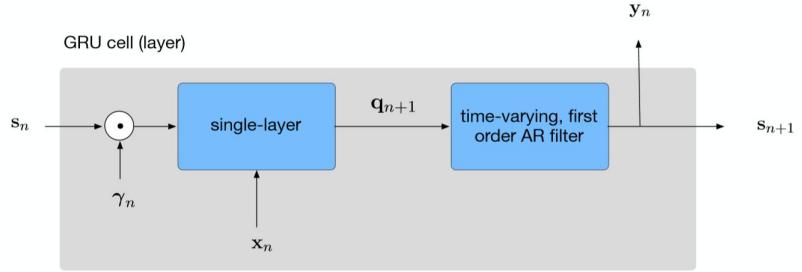
<https://towardsdatascience.com/illustrated-guide-to-recurrent-neural-networks-79e5eb8049c9>

<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

<https://medium.com/@kangeugine/long-short-term-memory-lstm-concept-cb3283934359#:~:text=This%20allows%20information%20from%20previous,with%20in%20the%20LSTM%20cell.&text=These%20gates%20determine%20which%20information,of%20%5B0%2C1%5D.>

GRU

Gated Recurrent Unit (GRU)

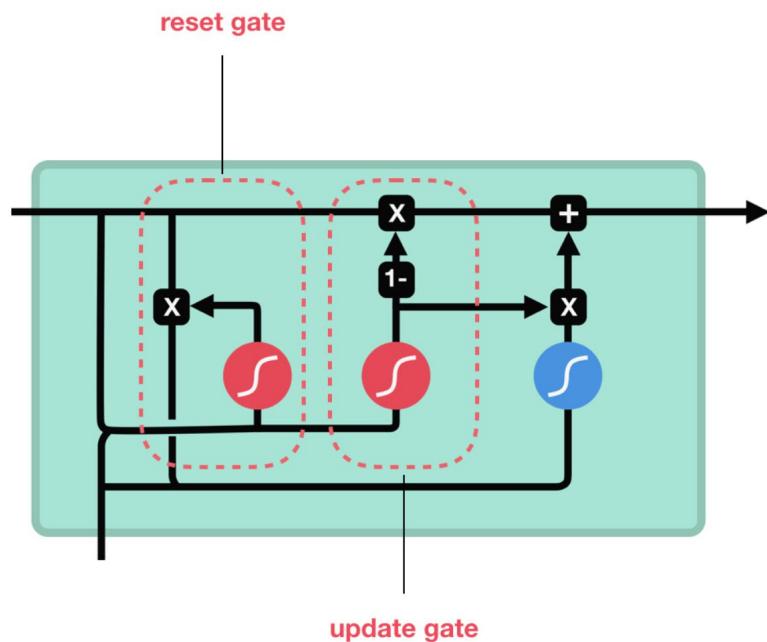


$$\mathbf{q}_{n+1} = \underline{h}(\mathbf{V}[\gamma_n \odot s_n] + \mathbf{W}\mathbf{x}_n + \mathbf{b}) \quad s_{n+1} = \alpha_n \odot s_n + (1 - \alpha_n) \odot \mathbf{q}_{n+1}$$

γ_n “read gate” — attenuation on current state when read for update

α_n “forget gate” and “write gate” coupled \sim unit DC gain filtering

$$\beta_n = 1 - \alpha_n$$



Update gate

(similar to forget+input gate)

$$\mathbf{z}_t = \underline{\sigma}(\mathbf{U}_z \mathbf{h}_{t-1} + \mathbf{W}_z \mathbf{x}_t + \mathbf{b}_z) \quad \text{“update gate”} \quad \text{“(forget gate)”}$$

Reset Gate

(similar to)

$$\mathbf{r}_t = \underline{\sigma}(\mathbf{U}_r \mathbf{h}_{t-1} + \mathbf{W}_r \mathbf{x}_t + \mathbf{b}_r) \quad \text{“regret gate”} \quad \text{ (“read gate”)}$$

$$\tilde{\mathbf{h}}_t = \underline{h}(\mathbf{U}[\mathbf{r}_t \odot \mathbf{h}_{t-1}] + \mathbf{W}\mathbf{x}_t + \mathbf{b})$$

$$\mathbf{h}_t = \mathbf{z}_n \odot \mathbf{h}_{t-1} + (\mathbf{1} - \mathbf{z}_n) \odot \tilde{\mathbf{h}}_t$$

GRU Trainable Parameters

$$\alpha_n = \underline{h}_r(\mathbf{V}_\alpha \mathbf{s}_n + \mathbf{W}_\alpha \mathbf{x}_n + \mathbf{b}_\alpha)$$

$$\mathbf{V} : \mathbf{N} \times \mathbf{N} \quad \mathbf{3}$$

$$\gamma_n = \underline{h}_r(\mathbf{V}_\gamma \mathbf{s}_n + \mathbf{W}_\gamma \mathbf{x}_n + \mathbf{b}_\gamma)$$

$$\mathbf{W} : \mathbf{N} \times \mathbf{M} \quad \mathbf{3}$$

$$\mathbf{q}_{n+1} = \underline{h}(\mathbf{V}[\gamma_n \odot \mathbf{s}_n] + \mathbf{W}\mathbf{x}_n + \mathbf{b})$$

$$\mathbf{b} : \mathbf{N} \times \mathbf{I} \quad \mathbf{3}$$

$$\mathbf{s}_{n+1} = \alpha_n \odot \mathbf{s}_n + \beta_n \odot \mathbf{q}_{n+1}$$

$$\mathbf{s}_{n+1} = \alpha_n \odot \mathbf{s}_n + (1 - \alpha_n) \odot \mathbf{q}_{n+1}$$

$$\mathbf{y}_n = \mathbf{s}_{n+1}$$

Number of Trainable Parameters: $3N(N + M + 1)$

Back-Propagation through Time (BPTT)

Word Embedding

BERT

BERT (Bidirectional Encoder Representations from Transformers), released in late 2018, is a method to **pretrain language representations** that was used to create models that NLP practitioners can then download and use for free. You can either use these models to extract high quality language features from your text data, or you can fine-tune these models on a specific task (classification, entity recognition, question answering, etc.) with your own data to produce state of the art predictions.

GRT-3

TensorFlow and Keras

sequential model

Pytorch