

LeNet: visual classification network that recognized digits for zipcodes

Representation Learning: learn from input image/output labels

- early layer weights task independent

→ customize models trained for diff tasks w less data

Multi-task Learning - sharing model weights across tasks improves performance on both

Deep nets have no obvious perf ceiling

$\hat{P}(Y|X) :=$ model approximation

Generative: computes full joint $P(X,Y)$

- can generate new data pairs (x_i, y_i)

Discriminative: compute only model target values conditioned on data $P(Y|X)$

generative includes additional assumptions

Generative

Discriminative

- strong assumptions about $P(X,Y)$

- insights into phys data generating

- faster training

- better performance w sparse data

- biased if assumptions are violated, poor asymptotic accuracy

- high bias error

Prediction

- simplify by making strong assumption

- y takes a single value given x

Loss Function measures difference between target prediction and target data value

Linear Regression: $L_2(\hat{y}, y) = (\hat{y} - y)^2$

$L = \sum_{i=1}^n (\hat{y}_i - y_i)^2$

$\hat{y}_i = ax_i + b$

Differentiate loss to find optimum values of a, b

- want to minimize expected loss on new data: $E((\hat{y} - y)^2) \approx$ risk

- actually minimize average loss across a finite number of data points

→ empirical risk

Cannot do well w empirical risk if

- biased sample

- not enough data

Multivariate: $y = Ax, x \in \mathbb{R}^m, y \in \mathbb{R}^k$

Gradient: a vector of partial derivatives

$\nabla_x L(x) = 0$, all partials are zero, loss not changing → local optimum

Logistic Regression - binary classification

$f(x) = \frac{1}{1 + \exp(-w^T x)}$

$\hat{y} = f(x)$ is probability $x \in$ target class

Cross Entropy Loss: negative log probability that every label is correct to 0.1

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

Compares target distr y_i w model distr \hat{y}_i

$L = - \sum_{i=1}^n y_i \log \hat{y}_i$

Softmax: $f_j(x) = \frac{\exp(CS_j)}{\sum_{i=1}^K \exp(CS_i)}$

$f_j(x) = \frac{\exp(CS_j)}{\sum_{i=1}^K \exp(CS_i)}$

$f_j(x) = \frac{\exp(CS_j)}{\sum_{i=1}^K \exp(CS_i)}$

$f_j(x) = \frac{\exp(CS_j)}{\sum_{i=1}^K \exp(CS_i)}$

$f_j(x) = \frac{\exp(CS_j)}{\sum_{i=1}^K \exp(CS_i)}$

$f_j(x) = \frac{\exp(CS_j)}{\sum_{i=1}^K \exp(CS_i)}$

Bias - difference between prediction and true y

$$\text{Bias}(\hat{f}(x)) = E[\hat{f}(x) - f(x)] = \hat{f}(x) - f(x)$$

Variance - variance of predictions

$$\text{Total squared error} = E[(\hat{f}(x) - f(x))^2]$$

$$= E[(\hat{f}(x) - f(x))^2] = \text{Bias}^2 + \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

$$= \text{Variance}(\hat{f}(x)) + \text{Bias}^2$$

k fold cross validation

- average hyperparam on validation set

partition into k sets, use different set for testing

To reach loss minimum:

- follow negative gradient

$-\nabla_w L(w)$

$w_{t+1} = w_t - \alpha \nabla_w L(w)$

GD: calculating gradient requires full pass thru dataset - expensive

SGD: mini batches of size m

N updates on full pass

gradient of function orthogonal to contour

Newton's method: compute vector straight to center

Quadratic (const)

Convergence

$x_{t+1} = x_t - \frac{f'(x_t)}{f''(x_t)}$

Update:

$x_{t+1} = x_t - H_t^{-1} \nabla f(x_t)$

- taking inverse is expensive

- quickly goes to near net gradient zero - often saddle point

SGD w momentum

$p_{t+1} = \mu p_t + (1 - \mu) \nabla f(x_t)$

$w_{t+1} = w_t + p_{t+1}$

p_t : momentum

μ : momentum constant

α : learning rate

prevents oscillation, may overshoot

Nesterov: 1) step in gradient dir

2) correct it accordingly

$p_{t+1} = \mu p_t + (1 - \mu) \nabla f(x_t)$

$w_{t+1} = w_t + p_{t+1}$

RMS Prop scales gradients by inverse of moving average

$s_t = \beta s_{t-1} + (1 - \beta) g_t^2$

s_t : moving average of squared gradients

Time t

$B_t: [0, 1]$ moving average decay factor

$w_{t+1} = w_t - \frac{g_t}{\sqrt{s_t}}$

ADAGRAD

\hat{C}_t : cumulative sum of squared gradients

$C_t = \sum_{i=1}^t g_i^2$

tends to grow linearly over time, so effective learning rate $\propto 1/\sqrt{t}$

- works well on sets w wide range of gradient magnitudes

- less effective w strong feature dependencies

ADAM: momentum + RMS Prop

compute moving averages of gradient + squared gradient

$p_t = \beta_1 p_{t-1} + (1 - \beta_1) g_t$

$s_t = \beta_2 s_{t-1} + (1 - \beta_2) g_t^2$

$w_{t+1} = w_t - \frac{p_t}{\sqrt{s_t}}$

Bias bc moments initialized to 0

Network gets more complex, more local minima near global minima

Evaluate backprop

Run output back, bc output Jacobian is a row vector

also common sub expr

Multidimensional arrays incorrectly called tensors

CS 182 MT1

2 Types of Jacobians:

- loss wrt input vector: data path

- loss wrt model params: model path

Convolution: effect of 2 signal on other

$(h * f)(x, y) = \sum_{i=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} h(i, j) f(x-i, y-j)$

Correlation: similarity

$(h \circ f)(x, y) = \sum_{i=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} h(i, j) f(x+i, y+j)$

Num Params: $(n \cdot m \cdot R + 1) K$

$n \times m$ filter

R input feature map

K output feature map - num filters

FCC: n input, m output $(n+1) \times m$

Convolutional Filters:

Conv \times Cin \times Fin \times Fw

output channels

input channels

Filter height

Filter width

Input \rightarrow Low level Feature \rightarrow high level Feature \rightarrow final output classifier

Conv Output Size: $(N-F)/\text{stride} + 1$

Image: $N \times N$, Filter: $F \times F$

K : # filters, S : stride

Conv Layer:

Input of W_1, H_1, D_1

Produces W_2, H_2, D_2

$W_2 = W_1 - F + 2P/S + 1$

$H_2 = (H_1 - F + 2P)/S + 1$

$D_2 = K$

F, F, D , weights per filter

Pooling down samples (no params)

LeNet-5 for handwritten digit classification

5x5 conv filters at stride 1

[conv-pool] x 2 - conv - fc

1998

AlexNet: similar architecture, bigger/deeper

1st layer: 48 11x11 filters w stride 4

output: [55x55x96]

2nd: [227x227x3]

~35K params

[conv-pool-norm] x 2 - [conv] x 3

- pool - [fc] x 3

- first use of ReLUs, GPUs, dropout

~60 million params

Transfer Learning

1) Train a Image Net

2) - small dataset

- retrain only classifier, i.e. last (conv) softmax layer

- medium dataset

- perform fine tuning

- used old weights as initialization, train full network or only some of higher layers

- may be retrain lower layers

- cover param in VGG
- batch norm after every conv layer
- no dropout
- trend toward no pool/FC layers
- smaller filters, deeper architectures

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

- squashes to [0, 1]
- can kill grads
- Key in LSTMs
- good for logical fn, learning non linear control
- bad for image net (ReLU)
- not zero centered

tanh: numbers to $[-1, 1]$

- zero centered
- kills gradients when saturated
- used in LSTMs for bounded, signal values
- not as good for binary fn

ReLU $f(x) = \max(0, x)$

- does not saturate in + region
- converges faster than sigmoid/tanh
- not suitable for logical fn
- not for control in RNN
- not zero centered
- leaky ReLU: $f(x) = \max(0, x), x$
- will not die

Parametric ReLU: $f(x) = \max(0, x)$

Exponential Linear Units

$f(x) = \begin{cases} x & x \geq 0 \\ \alpha \exp(x) - 1 & x < 0 \end{cases}$

- doesn't die
- closer to zero near outputs

Maxout Neuron

- non linearity
- linear regime, doesn't saturate, doesn't die
- doubles # parameters/neuron

Sigmoids good for smooth fn (robot control), and logical fn (and/or)

If initial weights are too small/large, activations will vanish/explode during fwd pass

Small: variance ↓: non-linearly

big: variance ↑: activations saturated, gradients → 0

Xavier Initialization keeps variance the same

- pick weights $N(0, \frac{1}{n})$, $n = \# \text{ input neurons}$

Batch Norm

$\hat{x}^{(b)} = \frac{x^{(b)} - \bar{x}}{\sqrt{\text{Var}[x^{(b)}]}}$

$\bar{x} = \frac{1}{n} \sum x^{(b)}$

- improves gradient flow thru net
- allows higher learning rates
- reduces stoch dependence on initialization
- reduces need for dropout
- limits magnitude of gradient
- randomly set some neurons to zero in forward pass
- forces net to have redundant representation

2) training large ensemble of models that share params

no dropout during test time

must scale activations so that for each neuron output at test time = expected output at training time

Inverted Dropout

- divide dropout mask by p at training time

Ensemble Learning

Bagging (Bootstrap Aggregation)

- train base models on bootstrap samples
- take majority vote for classification

average output for regression

- models trained independently
- reduces variance in prediction

Boosting

- ordinal learners
- each tries to reduce error on examples misclassified by earlier learners
- models are dependent, trained sequentially
- ADABOOST: weigh hard samples more

Gradient Boost: use residual to train later models

- reduces bias, possibly variance

Bagging often used in deep learning

models, but rarely used

→ parallelizes, models don't have much bias

The Ensemble: independent models

- prediction avg: avg prediction probs, or vote
- always works
- param avg: avg params, almost never works

Model Snapshot: train one model, take snapshots of params

- param avg often works
- snapshots close in parameter space

Gradient noise seems to help

- use validation set for hyperparameter tuning with each training block
- coarse → fine
- want ratio of weight updates / weight magnitudes to be ~ 0.001

Classification

Classification + Localization: single object

Object Detection: multiple

Instance Segmentation

Semantic Segmentation - label every pixel

1. **Classify every pixel**
 - extract patch, CNN on multiple pixels
 - repeat same computation multiple times
2. **CNN**
 - a bunch of layers to predict all pixels at once
 - convolutions can be expensive
3. **FC CNN**
 - downsampling/upsampling inside net

Downsampling: pooling, strided conv

Upsampling: nearest neighbor

Nearest Neighbor

1 2 → 1 1 2 2

3 4 → 1 1 4 4

3 4 → 2 2 4 4

Red of Naïv

1 2 → 1 0 2 0

3 4 → 3 0 4 0

0 0 0 0

Max Unpooling

- remember which element was max

- corresponding pairs of up/down sampling layers

Transpose Convolution

4. **UNET**: FCNN + Residuals

- residual connection made by copying input & concat w upsampled layer

Classification + Localization

bounding box

Image → Conv pool → conv feature → softmax (out)

Can also share FC layers

Per-Class Regression

- 1 bounding box for each class
- choose bounding box by predicting class (a bit)

Class Activation

- 1 bounding box total

Object Detection

- use metric called mean average precision = mAP
- mAP is # from [0, 100], high is good

Detection w Reg version

- depending on image, need variable sized outputs

Detection w Classification

- need to test many priors and scales
- only look at promising regions of image

Classification + Region Proposals

R-CNN

- 1) Input Image
- 2) Extract region proposals
- 3) Compute CNN features
- 4) Classify regions

Problems!

- 1) Finding Region proposals
- 2) Classifying each part of image is time/space consuming

Fast R-CNN - 250x speedup

- share computation of conv layers between proposals
- put whole image thru convnet before extracting regions
- region proposal after conv

Features, pool

Fast R-CNN - 250x speedup

- after CNN, include a Region Proposal NW
- network for external region proposals

RPN: slides small window on feature map

- classify object or not

- use N anchor boxes at each location

- regression gives offset from anchor boxes
- classification gives prob that each proposed anchor shows object

State of the Art: Single Shot Detection

- base boxes centered at each grid cell within each
- regress from each of B base boxes to final box w 5 numbers $(cx, cy, dx, dy, \text{conf}, \text{class})$

Output: $7 \times 7 \times (5 \cdot B + C)$

B base boxes, C classes

RethaNet

- 1) Eval pass of ResNet/Conv net
- 2) each level of downsampling, single shot detection

Recurrent NNs

- introduce cycles and a notion of time
- designed to process sequences of data
- produce sequences of outputs
- can unroll RNNs to support backprop (BPTT)

Layers are often stacked vertically

deep RNNs

- each layer has same parameters

Diagram showing layers x_0, x_1, x_2 with weights w_{01}, w_{12} and biases b_1, b_2 .

time

Lot of flexibility!

- vanilla NN
- many: image captioning
- many: sentiment classification
- many: machine translation
- many: video classification per frame

At each step: $h_t = f(h_{t-1}, x_t)$

same fn, params used at each step

State: hidden vector h

Diagram showing hidden states h_0, h_1, h_2 and inputs x_0, x_1, x_2 with weights w_{hx} and w_{hh} .

$h_t = \tanh(W_{hh} h_{t-1} + W_{hx} x_t)$

$y_t = W_{yh} h_t$

Sequence Generation

- top k sequences generated so far are remembered for next t

Sequence of RNN

Each word of sentence comes from different part of image

W_{hh} is multiplied by gradient at each time step

largest λ of $W_{hh} > 1$, gradients grow exponentially

$\lambda_{max} < 1$, gradients shrink exponentially

LSTMs

- non-linear, linearly transformed hidden states as well as memory cell c_t that is not transformed

LSTM encapsulates RNN

RNN: $h_t^* = \tanh(W^* h_{t-1}^*)$

LSTM: $\begin{pmatrix} i \\ f \\ o \end{pmatrix} = \begin{pmatrix} \text{sig} \\ \text{sig} \\ \text{tanh} \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ c_{t-1} \end{pmatrix}$

LSTM: $\begin{pmatrix} i \\ f \\ o \end{pmatrix} = \begin{pmatrix} \text{sig} \\ \text{sig} \\ \text{tanh} \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ c_{t-1} \end{pmatrix}$

$c_t = f \cdot c_{t-1} + i \cdot g$ $g = \text{output}$

$h_t = o \cdot \tanh(c_t)$ h output

- remember exactly what came in.

$c_t = f \cdot c_{t-1} + i \cdot g$ $g = \text{output}$

$h_t = o \cdot \tanh(c_t)$ h output

"gate" - how much info to let thru

c_t is filtered version of c_{t-1}

h_t is output: $\tanh(c_t) \times \dots$

LSTM

- 1) decide what to forget
- 2) decide what new things to remember
- 3) decide what to output

- cell c_t , gradient grows linearly w time

- path not well behaved

tSNE visualization (OCN²)

"Stochastic Neighbor Embedding"

- locally, pairwise distances are preserved
- similar things end up in similar place
- doesn't cluster data

DBSCAN - density based clustering

What do convnets learn?

- visualize patches that activate neurons
- visualize weights
- visualize representation space
- occlusion experiments
- deconv approaches
- optimization over image approaching

Occlusion - in part of image to see what classification most heavily depends on

We can generate an image that maximizes some class score

1) feed zeros

2) set gradient of some vector to be $\begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$ to not prop

3) image update

4) feed image thru net (in class we're interested in)

5) go to 2.

3 regularizers

- 1) penalize high freq 2) clip pixels w small norm
- 3) clip pixels w small contribution

Saliency Map: image that shows each pixel's unique quality

Deconv - maps features to pixels

Excitatory Input: positive influence on gradient

Inhibitory: negative influence

Linear Nature is primary cause of NN vulnerability to adversarial perturbation

Deep Dream modifies image to boost all activated features

Guided Backprop gives good results

requires image regularization - pressure to look like normal image

backprop: $R_i \cdot \frac{\partial f}{\partial x_i} \cdot \frac{\partial f}{\partial x_{i+1}}$

semantics - concerned w meaning of texts

Propositional / Formal semantics - block of text converted into formula in logical language

Vector repr - texts embedded into high-dimensional space

Prop: dog barks man → bites (dog, man)
 Vec: dog barks man → (0.2, -0.3, 1.5...)

Prop: allows for logical inferences

Vector: word embeddings, bag of words
 By of Words: sentence is sum of vectors of words in that sentence → no idea of order

Embedding: depend on word similarity → contexts words that share similar contexts have similar embeddings

One-hot encoding: vec size = vocab size

Vectors of contexts: words of nearby words

- vec size = vocab size
- better captures semantics - similar words have large inner products

Dimension Reduction

Latent Semantic Analysis:

1. treat words as equally similar/dissimilar
2. ignore word order

encodes set of docs in matrix

T_{ij} is count of word j in document i .

most entries are 0.

N docs T allows for efficient sparse matrix methods

Word features

Complete $T \times U \times V$

$T \approx N \begin{bmatrix} U & V \end{bmatrix} K$

K features

K embeddings

EF U, V, S is a diagonal matrix of singular values, if T is a document & U of T .

value is embedding of document in latent space

$C = UTV^T$ is decoding of document from its embedding $T \approx UVS$

SVD gives best possible reconstructions of documents C' from their embeddings.

$C' = U'V'U'$

LSA: auto encoding method

LSA uses squared L2 loss for deep net

Word2vec - 2 layer NN: negative sampling

uses words a few positions away from each center word - "skip grams"

considers all words as center words, and all context words

can predict center word from context word, vice versa - "skip-gram"

minimizes softmax loss $-\log p(c|c')$

for each output word given an input word

$p(c|c') = \frac{\exp(c_j^T v_i)}{\sum_k \exp(c_k^T v_i)}$

j - output word

i - input word

loss is largest for most similar word (j to i)

→ don't want similar words consecutive by deep net actually uses cross entropy - log of softmax

input x : one-hot encoded input word

y : index of output word

captures more info about relations and properties than LSA

LSA encodes context words in N -dim hidden layer

Criticisms:

- cross entropy loss puts much emphasis on small (unlikely) combinations of word/context
- expensive to normalize softmax over all words
- uses heuristic down weighting of frequent words

Co-occurrence Matrix

C_{ij} counts # windows containing both words i, j

$GLOVE$ - NN used for dimensionality reduction

$C_{ij} = \# \text{times word } j \text{ occurs in context of word } i$

minimizes: $J(C) = \sum_{i,j} f(C_{ij})(u_i^T v_j + b_i + b_j + \gamma C_{ij})^2$

u_i = word embedding, v_j is context word embedding

b_i, b_j are bias vectors, $f(\cdot)$ satisfies $f(0) \geq 0, f(x)$ non-decreasing, $f(x)$ saturates

LSM

Lexical Semantics: focuses on meaning of individual words

Computational Semantics: meaning depends on words, and on how they're combined

Skip Thought Embeddings: use sequence-to-sequence RNNs to next + previous sentences

Embedding is ~~output of boundary layer~~

output state vector of boundary layer

Prev Sentence

Cur Sentence

Embedding is output of this unit

Encoding doesn't require backprop

allows us to encode sentences + paragraphs

Siamese Networks for Semantic Relationships

trained on pairs of sentences a, b w similarity labels y

Sick Semantic entailment tasks:

- semantic sentences for relations: ENTAILMENT, CONTRADICTION, NEUTRAL

Hard Attention: attend to single input location

no GO - neural RL

Soft Attention:

- compute weighted combo over some inputs using attention w
- can use backprop to train w

Supervised Learning:

- independent input samples, each sample x receives a label y
- pair (x, y) assigned a loss that is differentiable

RL

- learner visits sequence of states in an epoch
- at time t , learner performs action a_t and receives reward r_t from environment
- agent maximizes sum of rewards over 1 epoch

Attention for Recognition: RNN based model

8a: input cat image

8b: glimpses are natural images

- part of image
- glimpse trace on some digit images

Soft Attention for Translation

RNN for Captioning:

- take layer w/1 before classification layer in CNN (features), pass into RNN, start generating text
- parts of description provided by certain parts of image.
- take feature w/ before spatial (CFC) layer
- pass spatially ordered feature into RNN
- generate weights for features h, d
- each time, generate attention/data vector
- attention made a feature, go into next iteration

Soft Attention: generate continuous set of weights, d_{in} = feature dim

1-based: generate set of probabilities that you attend to specific feature vector

Attention weights have same spatial dimensions as last CNN spatial layer

Salience: product of upstream gradient + feature layer activations

→ "importance"

Attention maps are explanations of not behavior - identify influential parts of input stream

DRAW: attend to arbitrary regions of input, then start narrowing down to classify

To generate: looks at partial studies it's drawn itself

3 params: x, y , width of box

Models improve accuracy, reduce computation

- learn to predict salience

Translation

Sequence to Sequence RNNs

- keep n-best list of partial sentences, along with their partial softmax scores

BLEUScore - compare machine translations against human generated translations, allowing for variation

Unigram Precision: count unique n-grams in reference sentence

unigrams in test sentence

Modified Unigram Precision: clip count by max occurred in any reference sentence

Attention heads find start of sentence, Subject/object relations, prepos

N-gram precision: count n-grams in reference sentence

n-grams in test sentence

Unigrams capture adequacy

N-grams capture fluency

BLEUScore: take weighted geometric mean of n-gram precisions up to some length C_n

$BLEU = BP \cdot \exp(\sum w_n \log P_n)$

$BP = \begin{cases} 1 & \text{if } C_1 \leq C_1^r \\ C_1 / C_1^r & \text{else} \end{cases}$

C_1^r = max

Sequence 2 Sequence

Encoder (input)

Decoder (output)

reversed LSTMs perform the best

→ first part of sentence is most important, so we use longer term dependencies from output to input sentence

Beam Search Width: how many phrases we keep track of

Criticism: bottleneck at last unit of encoder

Soft Attention for Translation - Bahdanau 2015

- for each output word, focus attention on a subset of all input words
- alignment score: how well do input words match; match output words to position i
- make weights: set max over alignment scores
- context vector (input to decoder)

Decoder

Yes Y_t

No Y_t

RNN

compares latent states & encoder hidden

h_t, h_{t-1}, h_{t-2}

x_t, x_{t-1}, x_{t-2}

x_1, x_2, x_3

encoder (bidirectional RNN)

Criticism: attention fn is complex, but seems to be best map on word similarity

data path is complicated, attention path is simpler recurrent net between output states

Long, Pham, Manning 2015 translation

- stacked LSTM w arbitrary depths
- global attention model sits above encoder/decoder, not recurrent
- simpler than Bahdanau's
- LOCAL attention model
- computes best aligned position first
- compute context vector centered at that position
- multiply & add must be recurrent

Paraphrase

- Sequence models generate linear structure, but there can easily encode things by "relating phrases"

SQuAD Parser

- 1) First tried training a basic seq2seq model on human annotated training time books. Poor results. (co-sim)
- 2) training on parse trees by Bank parser
- 3) added attention model (90.5 F1)
- 4) created synthetic dataset of high core ideas parse trees. SQuAD

- depth = 3, layer dim = 256
- dropout between 12, 243
- no part of speech tags!
- input reversing

Problems w Recurrent NW:

- sequential training + inference, hard to parallelize
- long range dependencies have to be remembered across many single time steps
- difficulty to learn hierarchical structure

Self Attention

info flows from within same row w

- constant path length, supports hierarchical info flow, trivial to parallelize
- weights inputs into propagation
- can entirely replace word-based recurrence

Transformer

- stacked with N=6
- position (or only): words are consumed then attention, location is lost

Attention

Encoder-decoder: queries from previous decoder layer, memory keys + values come from output of encoder

Self Attention Layer: all K, V, Q from output of previous layer in encoder

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

Multi-Headed Attention allows for parallel transformation - attend to multiple points of input seq w/ heads

Positional Encoding: Spatial encoding

- encoding vector has same dimension as model
- components are all sinusoidal fns of position
- can learn by Long & Beyer

Large Scale Summarization

- like translation, but no encoder

Memory Networks support Cross Attention

- reading + comprehension
- dialog
- learning from dialog
- reading w Attention over Memory

Frameworks:

Input Feature Map: convert input data to internal feature representation

Generalization: update memories given new input

1) produce new output (in feature representation space) given the memories

2) (response) convert output O into a response seen by outside world

standard

QA Memory NW

Input data: $m, x, c, c' = C(x, c)$

Query Embedding: $u = Bq + \text{embedding}$

Attention: $P_i = \text{softmax}(u^T m_i)$

Output: $O = \sum P_i c_i$

Prediction: $\hat{O} = \text{softmax}(W \text{cot}(u))$

Model trained on 2e on series of assertions/questions; learns A, B, C, W

Position Encoding: allow an embedded words to carry info about their location in the input

computed & used

Memory Networks - using KVStore from context

- sentence level encoding: fine text
- input broken into sentences, encoded in BoW as key, value = standard mean NN
- window level: encode window of W words in BoW as key, use center word as value
- KB Triple: typically have form "subject relation object" - key is subject-relation pair, value is object
- original natural languages datasets

standardized using Simple Questions dataset

SQuAD 1.0 - Stanford QA Dataset

- 100k+ questions posed by crowd workers on a set of Wikipedia articles
- longer answers

SQuAD 2.0 - 1.0 + active neural comp

- 1.0 could cheat by providing answer of appropriate type
- 2.0 includes unanswerable q's, preventing educated guesses

QANet

- conv input + transformer
- transformer encoder
- conv layers + linear
- in NW
- uses LayerNorm to reduce covariate shift by scaling and biasing activations

Context Query Attention Block:

- calculate similarity between query, context & embeddings
- row wise softmax $\Rightarrow A$

Query Context Attention Block:

- column wise softmax $\Rightarrow B$

Model Encoder, blocks take context, 2 attention matrices

unigram outputs sum over word positions
for start + end of answer span
non recurrent design, training time doesn't grow w input length
backtranslation provides more depth in training data
multiple attention heads
bilinear combination of context & query

Goal-Directed Dialog Systems

- engage the user
- help user w goal-directed tasks
- learn from sample dialog how to respond to user queries

Task Summary

- 1) Issue API calls, w regulated fields
- 2) Updating API calls users update requests
- 3) displaying options retrieved from KB
- 4) Providing extra info
- 5) Conducting full dialogs

Datasets:

Restaurant Reservations:

- contains two KBs of 4200 Cnbs, 600 restaurants each
- Dialog State Tracking Challenge
- Restaurant booking data from real users
- every dialog turn labeled with a state to be predicted

Online Concierge Service

- Data from real online concierge service
- users make requests through chat interface handled by human operators who make API calls

BERT: language model trained on large dataset of natural language fine-tuned for particular tasks

Generative Pre-Training (GPT): language model

- transformer-based generator only
- single left-to-right transformer
- uses left-to-right transformer
- dis cards multi-head attention (middle) step
- initially trained on large text corpus to minimize language modeling loss
- $$L(u) = \sum_i \log P(u_i | u_{1:i-1}, u_{i+1:T}; \theta)$$

For each task, custom linear layer is added & entire network retrained

can be applied to a variety of non-generative tasks

- entailment tasks
- question answering & commonsense reasoning

GPT-2 - depth 48 GPT trained on WebText

BERT - a recur dropped words (predicting each)

- bidirectional transformer encoder model
- trained with two types of loss
- 1) Word Prediction: 15% of input words are removed then re-predicted
- 2) Next Sentence Prediction - predict whether pairs of sentences are consecutive
- = contrastive loss: contrast true positives with near miss "negatives"

Input Representation:

- token embeddings + segment embeddings + position embeddings

Tasks:

- Sentence Pair Classification Tasks
- Single Sentence Classification Tasks
- Question/Answering Tasks
- Single Sentence Tagging Tasks

Dialog System Technology Challenge

- Sentence Selection
- Sentence generation
- Audio-Visual scene-aware dialog
- BERT performs best

CoQA: Dataset

127k+ questions w answers from 8000+ conversations

BERT, GPT, ELMO 2-pn training for downstream tasks, similar to word2vec

ELMO uses RNN: context indep. trained Lstm, Rnn, LSTM to generate features for downstream tasks (BERT)

Bi-Directional Encoder models (BERT) better than generative models (GPT) at non-gen tasks

- gen models have training efficiency & scalability advantages that may make them more accurate
- can solve entirely new kinds of tasks w text generation

Adversarial Examples

1) $n = \text{sign}(\nabla_x \ell(\theta; x, y))$

sign of gradient = perturbation
2) makes bigger change than optimum Optimization Based Attack:

- make as small as possible change
- make classifier output target class
- maximize label loss
- Perturbation must triumph noisy physical world
- need to find common δ to perturb image
- take into account spatial constraints
- apply perturbation on only part of image very weak

Adversarial perturbations are possible in physical world under different conditions and viewpoints, including distances & angles
2 different approaches to optimize objective → need to know model θ

Black-box Attacks (don't know model)

- Zero Query
- Random Perturbation
- Difference of Means
- Transferability Based

train example to fool source learning method, transfer to target method

Non-Targeted: get classifier to predict any incorrect label, easier

Targeted: get classifier to predict target label, harder

Targeted Adversarial Examples: Transferability is POOR

- create adversarial examples from white box model, transfer to black box system

Query Based Attack

- finite difference gradient estimation
- every reduced gradient estimation
- d-dimensional vector x , make 2d queries to estimate gradient
- $$FD_x(g(x), \delta) = \frac{g(x + \delta e_i) - g(x - \delta e_i)}{2\delta}$$

finite difference could be better than white box

For every reduced gradient estimation:

- random grouping - PCA
- random direction sample
- gradients along some directions (dimensions) are weak and fall off
- can find gradient in strongest dimension
- only need very few queries to find this

use GAN to help w perturbation

- want to minimize distortion
- 1) train classifier to discriminate real/artificial images
- get target label, minimize distortion
- 2) classifier classifies adversarial examples as well as possible

discriminator loss replace L2 loss to get adversarial generator to produce real images

GAN attack relatively resilient against defenses

Geometric distortion flow - more than perturbing each individual pixel

Why do Adversarial Examples Work?

- shifts attention
- inception-v3 model \approx human perception
- can train models to resist adversarial examples

Adversarial trained models unable to discriminate GAN-produced examples

can fool auto encoders: input \rightarrow latent code \rightarrow output

VAG-like generative models vulnerable to adversarial examples

Can also fool RL - perturbation on every frame (but not necessary)

Local Intrinsic Detection

- distribution of natural images (color, etc)
- adversarial examples are outliers
- natural images have context of neighbors

Generative Models

reparam Full joint $p(x, y)$

Classifier: given x, y , determine $p(x, y)$

Generator: given y , generate sample x from $p(x, y)$

Auto Encoders

- transform input x into code variable z , back to synthetic input x'
- $\dim(z) \ll \dim(x)$ = info bottleneck

Encoder: very similar to, or is neural classifier

Code: representation of input designed to support reconstruction

could be soft & class labels, often encodes within-class variability

Decoder: conditional (synthetic) input generator (generative model)

Can create decoder as inverse of encoder - activation maximization

BUT: inversion is iterative-slow

- train classifier to minimize loss
- $L(x, x')$
- $z \rightarrow x'$ not differentiable, KL

Implicit Auto Encoder

suppose we have $x = f(z, \theta)$

invert $f(z, \theta)$ to get z and minimize $L(x, x') = f(z, \theta)$

slow bc iterative inversion

slow bc high variance from individual samples

VAEs - L2 loss

build decoder as probabilistic model $P_\theta(z|x)$

invert using variational inference to get approximation to $P_\theta(z|x)$

$q_\phi(z|x)$

$z = g_\phi(x, \epsilon), \epsilon \sim N(0, 1)$

normal distr of z -values with mean μ and SD Σ

Loss F_θ :

optimize marginal likelihood of data

$$\sum_i \log P(x_i)$$

Maximize variational lower bound:

$$\mathcal{L}(\theta, \phi; x) = \mathbb{E}_{z \sim q_\phi(z|x)} [\mathcal{L}(\theta, \phi; x, z)]$$

Jensen's Inequality:

for concave $f(x)$, linear

$f(\mathbb{E} x) \geq \mathbb{E} f(x)$

at $\bar{x} = \mathbb{E} x$

$f(x) \leq \mathcal{L}(x)$

$V = \mathbb{E}_\theta - P_\theta[\mathcal{L}_\theta(\theta, \phi; x, g_\phi(z, \epsilon))]$

Update parameters θ, ϕ until conv.

Loss estimator has high variance

→ alternate estimator w KL divergence, which assumes that $q_\phi(z|x)$ allows density estimation

Simple loss: $\mathcal{L}(\theta, \phi; x) = \frac{1}{2} \sum_i \log p_\theta(x_i)$

θ, ϕ gradients are discrete - $\log q_\phi(z|x)$

Alternative:

$\mathcal{L}(\theta, \phi; x) = D_{KL}(q_\phi(z|x) || P_\theta(z))$

θ grad's approximate, ϕ is exact

BOTH: require $q_\phi(z|x)$ to provide a density and be integrable in closed form

$z = g(x, \epsilon) = \mu(x) + g(x, \epsilon)$

VAEs are efficient because problems w 3D images

say nothing about form of encoder & decoder fns

Autoregressive Models (full generative model)

Key: pixels highly correlated over various scales

generate one pixel at a time conditional on values of other pixels

Local Autoregressive

examine only a few nearby pixels at each step

PixelCNN

iteratively generate 1 pixel at a time

fast to train - no recurrence

slow at generating - pass through entire network for each pixel

only training must use masks to ignore later pixels

PixelRNN - row LSTM

low recurrence instead of 3x3 conv to allow long range dependencies

able to generate full row of pixels in 1 step

takes longer to train than PixelCNN (recur)

Row LSTM VBS reconstruct state only from tiny subset of previously computed states

Diry BiLSTM VBS all

Replace ReLUs w gating

→ PixelRNN

Image Transformer

- autoregressive design (looks back)
- use multi-headed attn network
- use large influence region of image without recurrence

computes pixel rep in layers

uses dropout for regularization

3D pos encoding, now with color coord

GANs

human compare generated image w image in dataset

→ 2 network architecture

1) generates image

2) nn discriminator to compare real & synthetic images

GANs learn transformation from known image distribution

Minimax: generator tries to increase discriminator loss, discriminator decreases

→ local Nash equilibrium found w adversarial training

Generator Network: generates realistic images

learn transformation from random distribution to target distribution

Discriminator Network

measures "goodness" of generated sample

outputs probability of image being real

Step 1: look at generated images

Step 2: look at real images

Discriminator: maximize probability real classified real, minimize generated classified real

Generator: maximize generated = real

GANs actually produce distributions similar to real distribution

minimize Jensen-Shannon divergence between $P_{\text{real}}(x)$, $x = G(z)$

Algorithm

For each training iteration:

For k steps:

Sample batch of m noise samples

Sample batch of n examples from dataset

Update D by ascending GAD

Update G by descending GAD

Maximize $\frac{1}{m} \sum_i \log(D(G(z_i)))$

be other $\log(1 - \dots)$

gives small gradients when generator samples are bad

want large gradients when generated samples are bad

In practice, must balance D/G to ensure 2 good gradients

GAD on cost fns not guaranteed to converge to actual Nash equilibrium

DCGAN

generator now from allow net

fractional - transposed conv

replace pooling w strided conv

batchnorm in D/G

remove FC layers - use fully conv

Conditional GAN

add label y as specific class

if dataset y want to sample from

ACGAN

add discriminator output class label

Text-to-Image Synthesis enabled

y class is sent as vector

→ supports interpolation on y

Problems:

Mode collapse

appearance reusing / substituting

old D/G every few iterations

multiple gradients D/G is tricky, allow generator to look ahead

Difficulty in Training

Wasserstein Distance

min cost to match distrib

can keep learning until converge

→ gradient's always there

effectively a lower bound of KL

Progressive GANs: build resolution

→ in resolution

add layers to GAN, and

upsampling layers incrementally

BigGAN

→ need a lot of power

→ sample

→ difficulty in evaluation

→ randomly pick y , class