

Deep learning - learn the feature & model from data
old ML: create a program that learn from examples
Design the features by hand.

Supervised Learning = labeled examples of correct behavior i.e. ground truth input/output

Reinforcement Learning = learning system receives a reward signal, tries to learn to maximize the reward system

Unsupervised learning = no labeled examples, looking for interesting patterns
NN can learn useful representation in supervised, unsupervised, or reinforcement learning.

e.g. generate image, GPT2 (next word) robot

Linear Model $X = (\vec{x}_1, \dots, \vec{x}_N)$ $\vec{t} = (\vec{t}_1, \dots, \vec{t}_N)$

linear regression $\vec{y}^{(i)} = \vec{W}^T \vec{x}^{(i)} + b = (\vec{w}_1, \dots, \vec{w}_D) (\vec{x}_1^{(i)}, \dots, \vec{x}_D^{(i)}) + b$

loss: $L(\vec{y}^{(i)}, \vec{t}^{(i)}) = \frac{1}{2} (\vec{y}^{(i)} - \vec{t}^{(i)})^2$

cost: $\Sigma(\vec{W}) = \frac{1}{N} \Sigma L(\vec{y}^{(i)}, \vec{t}^{(i)})$

$\vec{y} = \vec{X}\vec{W} + b\vec{I} = (\vec{x}_1^{(1)}, \dots, \vec{x}_D^{(1)}) (\vec{w}_1, \dots, \vec{w}_D) + (b) = (\vec{y}^{(1)})$

$\Sigma(\vec{W}) = \frac{1}{2N} \|\vec{X}\vec{W} + b\vec{I} - \vec{t}\|^2 = \frac{1}{2N} (\dots)^T (\dots)$

$= \frac{1}{2N} \|\vec{X}\vec{W} - \vec{t}\|^2 \quad \frac{\partial \Sigma}{\partial \vec{W}} = \frac{1}{N} \vec{X}^T (\vec{X}\vec{W} - \vec{t})$

Gradient Descent

① Direct Solution $\nabla_{\vec{W}} \Sigma(\vec{W}) = 0 \Rightarrow \vec{W} = (\vec{X}\vec{X})^{-1} \vec{X}^T \vec{t}$

② Gradient Descent Update Rule: $\vec{W} \leftarrow \vec{W} - \alpha \nabla_{\vec{W}} \Sigma(\vec{W})$

Direction = negative of gradient's sign

Size = proportional to gradient's magnitude

③ GD update rule for Linear Regression

$\vec{W} \leftarrow \vec{W} - \frac{\alpha}{N} \vec{X}^T (\vec{X}\vec{W} - \vec{t})$

$\vec{W} \leftarrow \vec{W} - \frac{\alpha}{N} \sum_{i=1}^N \vec{x}^{(i)} (\vec{W}^T \vec{x}^{(i)} - t^{(i)})$

④ More efficient: Each GD cost $O(ND)$, rather than $O(D^3)$ for matrix inversion. Much cheaper if D is large (i.e. High dimensional data)

Feature Maps = learn non-linear relation using LR

⑤ Polynomial feature mapping: $y = W_0 + W_1 x_1 + \dots + W_m x_m^m$

$\psi(x) = [1 \ x_1 \ \dots \ x_m^m]^T$ $y = W^T \psi(x)$ is linear in $\psi(x)$

⑥ Generalization: to unseen data; the degree of polynomial (m) is a hyperparameter. Tune it using validation set.

Train \rightarrow Validation \rightarrow test tune hyperparameters

Logistic Regression - binary linear

Linear model: $Z = \vec{W}^T \vec{x}$ logistic activation: $y = \sigma(Z) = \frac{1}{1+e^Z}$

Cross Entropy: $LCE(y, t) = -t \log y - (1-t) \log(1-y)$

Classification accuracy = discontinuous \rightarrow SGD

LCE: it is a proper scoring rule \rightarrow opt y is true probability

Softmax Regression weighted matrix k class
linear model: $Z = W\vec{x}$ $\vec{x} \times k$ (classes)
softmax activation: $\vec{y} = \text{softmax}(\vec{z})$ $y_k = \frac{e^{z_k}}{\sum_{m=1}^k e^{z_m}}$
CE loss: $LCE(\vec{y}, \vec{t}) = -\vec{t}^T \log(\vec{y})$

$\frac{\partial y_k}{\partial z_k} = y_k(1-y_k)$

$\frac{\partial y_k}{\partial z_j} = -y_k y_j$ if $j \neq k$

Multi-layer Perceptrons single neurons are linear classifier and have limited expressive power

$\frac{\partial y_k}{\partial z_j} = -y_k y_j$ if $j \neq k$

XOR not linear separable $\Psi(x) = [x_1 \ x_2 \ x_1 x_2]^T$

Prove by contradiction: Suppose there is a feasible solution. The average of all transaction of A is the vector $(0.25, 0.25, \dots, 0.25)$. Therefore, we must classify this vector as A. The average of all translations of B is the vector above. Therefore, we must classify this vector as B. Contradiction!

Convex: $x_1, x_2 \in S \Rightarrow \lambda x_1 + (1-\lambda)x_2 \in S \quad 0 \leq \lambda \leq 1$

$\lambda_1 x_1 + \dots + \lambda_N x_N \in S, \lambda_i > 0, \lambda_1 + \lambda_2 + \dots + \lambda_N = 1$

Multilayer Perceptron $f((1-\lambda)x_0 + \lambda x_1) \leq (1-\lambda)f(x_0) + \lambda f(x_1)$

Directed Acyclic Graph; Feed forward NN

Fully connected layer $y = \phi(Wx+b)$

(N input unit, M output unit $\Rightarrow M \times N$ weight matrix)

Activation

ReLU $\sigma(z) = \max\{0, z\}$ nondifferentiable at 0

Soft ReLU (softplus) $\sigma(z) = \log(1+e^z)$

• does not produce zero gradient for negative inputs

• differentiable everywhere computational expensive

Learning XOR 别忘了每一个单元的值是对应的 Activation!

$x_1 \ x_2 \ x_1 x_2 \ h_1 \ x_2 - x_1 \ h_2 \ h_1 + h_2 \ y$

$\begin{matrix} 1 & 0 & 1 & 1 & -1 & 0 & 1 \\ 0 & 1 & -1 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 0 & 0 & 0 & 0 &$

Distributed Representation

Language modelling → Predicting next word

Build a generative model:

- the observation model, $P(\alpha|s)$
- The prior, $P(s)$
- Infer $P(s|\alpha)$ using Bayes rules posterior prob

Learn a distribution $P_\theta(s)$ of sentences

Given some observed sentences $s^{(1)}, \dots, s^{(N)}$

Learn θ using maximum likelihood that maximizes the probability of the observed s

Prob. of sentences: $P(s^{(1)}, \dots, s^{(N)}) = \prod_{i=1}^N P(s^{(i)})$

Use log to prevent prob becoming vanishingly small:

$$\log \prod_{i=1}^N P(s^{(i)}) = \sum_{i=1}^N \log P(s^{(i)})$$

$$P(s) = P(W_1, \dots, W_T) = P(W_1)P(W_2|W_1)\dots P(W_T|W_1, \dots, W_{T-1})$$

Markov Assumption: The probability of the next word depends on a fixed # of prev words.

Memoryless model: $P(W_t|W_1, \dots, W_{t-1}) = P(W_t|W_{t-1}, W_{t-2}, W_{t-1})$

Auto Regressive Model: supervised prediction prob.

Auto: use its prev output/values to predict the next value. The model is self referential

Regressive: regression. predict the current value based on past values.

3-Gram Language Model

$$P(W_3=\text{Cat}|W_1=\text{the}, W_2=\text{fat}) = \frac{\text{count}(\text{the cat})}{\text{count}(\text{the fat})}$$

Limitation: (row) many N-grams are rarely seen in training data

Data Sparsity = $N \uparrow$, # of possible N-gram exp \uparrow

Sol: shorter context - smooth prob by adding imaginary count

Combine predictions from N-gram models with diff value of N . b^1 must also be local min

N-gram is a localist representation: Each N-gram

is treated as a distinct & independent entity

The info is not shared btw N-grams as the dog sit

the cat sit

Neural Language Model

Distributed representation: Words are represented as vectors in a continuous space.

word → attribute → next word

Loss function: $-\sum_{i=1}^V t_i \log p_i$ ← prediction

t_i = one hot encoded target (1 for correct, 0 for wrong)

p_i = the predicted probability

Bengio's Neural Language Model:

Input layer → Embedding layer → Hidden layer → Output layer.

Previous word Each word → vector NN Prob distn of next word.

Localist: a piece of info is stored in one place

Distributed: info is shared between related entities

more complex and difficult to interpret. capture more nuanced info

more robust to damage or noise. large dataset

Global vector (Glove) Embeddings

Learning word embeddings from global word-word co-occurrences. V is the vocabulary size. X is the co-occurrence matrix ($V \times V$ matrix).

$$X = R R^T \quad X_{ij} = Y_i^T Y_j$$

$$\text{cost function } J(R) = \sum_{ij} f(x_{ij})(Y_i^T Y_j + b_i + b_j - \log x_{ij})^2$$

$$\text{weights fix}_{ij} = \begin{cases} \left(\frac{x_{ij}}{x_{\max}}\right)^2 & \text{if } x_{ij} < x_{\max} \\ 1 & \text{if } x_{ij} \geq x_{\max} \end{cases}$$

Why log count

Emphasize the diminishing importance of very high counts. reduce the range of counts, learning more stable

emphasize the relative diff in counts

PROS: - capture global statistics

- capture semantic relationship between words

- can be trained relatively quickly on large dataset

- Pretrained Glove Embeddings available

CONS: - static embedding - words representation

does not change based on its context.

The global co M requires significant memory for storage & computation during training

struggles with words not seen during training

convex: any local is global

strong convex: unique global

non convex: mult local

local not global weight space symmetry

NN are not convex: weight space symmetry

Suppose we are at a local min θ

swap two hidden units with its weights

and biases and get θ'

θ' must also be local min

combine predictions from N-gram models with diff value of N

θ' is random start & local min is fine.

Saddle point

The gradient (1st deriv) = 0

The Hessian matrix has bot +, - eigenvalues

eg two h_1, h_2 have same incoming & outgoing weights

aft GRD, still same ⇒ forever; ② all init w-b=0

all w be identical ⇒ no learning

⇒ Initial the weights to small random values.

plateaux

saturated units

output value close to extreme

→ vanishing gradients

.. close to zero

⇒ zero gradient

choose the scale of random init

of the weights so that activation are in the middle of their dynamic range

Xavier: initialize the weights s.t. variance of weights activations is same across every layer

→ ReLU: initial the biases to a small + value

Badly conditioned curvature

Ravine w_2 Gradient along w_1 is Large GD takes large steps we should take small steps

⇒ avoid ravine: normalize your inputs to zero mean and unit variance

② Batch normalization explicitly center each activation. speed up training. so that activation have zero mean and unit variance.

⇒ Momentum

$$M_{k+1} \leftarrow \mu M_k + \alpha \frac{\partial J}{\partial \theta}$$

$$\theta \leftarrow \theta - M_{k+1}$$

How quickly the contribution of previous gradients exp decay (0.9, 0.99)

Second order information

RMSprop (Root Mean Square Propagation)

An exp ma of the squared gradient

$$S_j \leftarrow (1-\gamma)S_j + \gamma \left[\frac{\partial J}{\partial \theta_j} \right]^2$$

$$\theta_j \leftarrow \theta_j - \frac{\alpha}{\sqrt{S_j + \epsilon}} \frac{\partial J}{\partial \theta_j}$$

Δθ.

Adam (Adaptive Moment Estimation)

$$M_{k+1} \leftarrow B_1 M_k + (1-B_1) \frac{\partial J}{\partial \theta}$$

$$S_{k+1} \leftarrow B_2 S_k + (1-B_2) \left[\frac{\partial J}{\partial \theta} \right]^2$$

$$\hat{M}_{k+1} = \frac{M_{k+1}}{1-B_1^{k+1}} \quad \hat{S}_{k+1} = \frac{S_{k+1}}{1-B_2^{k+1}}$$

$$\theta_{k+1} \leftarrow \theta_k - \frac{\alpha}{\sqrt{\hat{S}_{k+1} + \epsilon}} \hat{M}_{k+1}$$

convolutional NN and Image Classification

Convolutional layer small region of the image

Pros: fewer parameters the weights are shared

effective computation, better generalization

preserve spatial structure of input

Translation Invariance: shared weights allows the layer to detect features regardless of their position

convolution operation

$$S(i,j) = (I * K)(i,j) = \sum_m \sum_n I(m,n) K(i+m, j+n)$$

flipped kernel

Implementation: not flipped kernel

sharpen blur horizontal edge

convolutional Network structure

Input layer → convolution → NonLinear Activation

→ Pooling → Next convolution

convolution is a linear operation

pooling = reduce the dimensions of the feature maps

replace the output with a summary statistic

Benefits of pooling layer • reduce D ...

{ reduce H of parameters, reduce computational load faster to train. Prevent overfitting}

• support translation invariance:

• Pooling layers expands receptive field. helps higher-layer cover larger regions and enable deeper layers to interact with more of the input.

Image classification

Conv Net architecture

Size of a Conv Net

for back prop # of units: store activations in memory

of weights: for training & testing

of connections: complexity of backprop (3 add-multiply per connection / forward 2 backward)

For a fully connected layer with M input, N output units

it has MN connections, MN weights

Fully connected layer Convolution layer

output units IWH IWH

weights $W^H IJ$ $K^2 IJ$

connections $W^H IJ$ $W^H K^2 IJ$

Alex Net • ReLU activation (faster training & convergence)

• weight decay (L^2 regularization) prevent overfitting

• Data augmentation (image translate) size of data

• SGD with momentum (with previous update)

• Dropout randomly set a fraction of input units to prevent overfitting

GoogleNet ↓ # of weights

• No fully connected layer smaller convolutions

• Break down convolutions into multiple

Inception module

use parallel convolutions with various filter sizes

capture features at multiple scales with

1x1 convolutions reduce dimensions of feature maps

allows for a deeper network without significantly more computational complexity.

$$(A+B)^T = A^T + B^T \quad (AB)^{-1} = B^{-1} A^{-1} \quad \frac{\partial x}{\partial z} = \frac{\partial x}{\partial y} \cdot \frac{\partial y}{\partial z}$$

$$(rA)^T = rA^T \quad I^{-1} = I \quad \frac{\partial x}{\partial z} = ab^T$$

$$(AB)^T = B^T A^T \quad f(w) = W^T A w + b^T w + y \quad \frac{\partial x}{\partial z} = (B+B^T)X$$

$$I^T = I \quad \nabla f = \frac{1}{2}(A^T A)w + b \quad \frac{\partial x}{\partial z} = ab^T$$

$$AA^T = A^T A = I \quad = aw+b \quad (\text{symm}) \quad a(X-AS)^T w(X-AS)$$

$$(IA)^T = \frac{1}{2}A^{-1} \quad f = WAw \quad = -2A^T w(X-AS)$$

$$I^T = I \quad \nabla f = 2A \quad X = 2w(X-AS)$$
</div

Optimization and Generalization

Stochastic Gradient Descent

batch training: $\frac{\partial \mathcal{E}}{\partial \mathbf{W}} = \frac{1}{N} \sum_{i=1}^N \frac{\partial L(y^{(i)}, \mathbf{t})}{\partial \mathbf{W}}$ summing over all training examples

SGD: Update model parameter based on the gradient calculated from a single, randomly chosen training ex.

BT: moves directly downhill on avg.

SGD: Take steps in a noisy direction, but moves downhill on avg.

Fast update (SGD), more mem/stuck local / guaranteed conv (Bi) dead/saturated unit histogram

Mini-Batch Gradient Descent: compute gradients from small, randomly sampled subsets of the training data

Faster convergence • Smoother conv • handle large dataset
variance
Vectorized operations

mini batch size S hyperparameter & saturate at large size GPU favor larger batch size \Rightarrow afford more parallelism larger batch size require more mem. Diminishing returns In convergence improvement VS. computation cost.

(faster conv; local min; GPU saturate) steps 1. optimizing batchsize $.25, .26, .27, .28$. Tune S by plotting validation error against training time for diff S . 2. Tune other para before α and S - use GPU by increasing learning rate hyper grid search

• start at a higher learning rate for quick progress
• If the error \uparrow or oscillates wildly $\rightarrow \alpha \downarrow$
• If error \downarrow consistently, but slowly $\rightarrow \alpha \uparrow$
• At the end of mini batch $\rightarrow \alpha \downarrow$ to remove fluctuation When error stop $\downarrow \rightarrow \alpha \downarrow$ loss fluctuate $\rightarrow \alpha \downarrow$ epoch

2. could cause sudden \downarrow in loss at the expense of LT bc a reduction in upward fluctuation than dramatic improve.

④ Ensembles: Avg or combine the prediction from decaying α • step decay • exponential decay $\alpha = \alpha_0 e^{-kt}$ • $1/t$ decay: $\alpha = \alpha_0 / (1+kt)$ • use a grid search on log scale (0.1, 0.01, 0.001) for diff random initialization; diff subset of training data (bagging); diff arch or hyperpara

SGD + Momentum

contour map • 同一条线相同. 越密 gradient越大. Gradient Dropout: randomly deactivate neurons during training \perp contour line 并且指向 increasing function value (GID) with probability p . At test, scale w by $1-p$

Ravine \curvearrowleft steep gradient - take large step - should \downarrow
shallow gradient - take small step - should \uparrow

Momentum steep: dampens oscillation, converge \downarrow
shallow: flat: build up velocity, accumulating small g
 $m_{k+1} \leftarrow \gamma m_k + \alpha \frac{\partial \mathcal{E}}{\partial \theta}$ $\theta \leftarrow \theta - m_{k+1}$ decay

γ : damping parameter, how quickly the prev g exp

RMS Prop (Root Mean Square Propagation)

adjust α based on mv avg of the squared gradient $s_j \leftarrow (1-\gamma) s_j + \gamma (\frac{\partial \mathcal{E}}{\partial \theta_j})^2$ $\theta_j \leftarrow \theta_j - \frac{\alpha}{s_j + \epsilon} \frac{\partial \mathcal{E}}{\partial \theta_j}$ ($\alpha \Rightarrow \alpha \downarrow$)

Adam (Adaptive Moment Estimation) smooth out overpara more para than data points landscape with infinitely many possible global ma of past gradient $m_{k+1} \leftarrow B_1 m_k + (1-B_1) \frac{\partial \mathcal{E}}{\partial \theta}$ ma of past squared gradient $s_{k+1} \leftarrow B_2 s_k + (1-B_2) [\frac{\partial \mathcal{E}}{\partial \theta}]^2$ est var of gradients correct the bias in ma:

$\hat{m}_{k+1} = \frac{m_{k+1}}{1-B_1 m_k}$ $\hat{s}_{k+1} = \frac{s_{k+1}}{1-B_2 s_k}$ estimates are initial or no bias the estimates are sys lower than their true value update parameter: $\theta_{k+1} \leftarrow \theta_k - \frac{\alpha}{\sqrt{s_{k+1} + \epsilon}} \hat{m}_{k+1}$ prevent θ

Problems

- incorrect gradient summing over all training examples
- local min slow progress
- symmetries instability oscillations
- slow linear training curve
- cost increases fluctuation in training curve
- fluctuation in training curve
- initial W randomly $\alpha \downarrow$, momentum decay $\alpha \downarrow$ iterate avg initial scale of w relu Adam, second order opt
- activation chard
- histogram ill condition

Diagnostics

- finite diff (hard)
- random restart
- visualize w
- slow, linear training curve
- cost increases
- fluctuation in training curve
- $\alpha \downarrow$, momentum decay $\alpha \downarrow$ iterate avg
- initial scale of w relu Adam, second order opt
- activation chard
- histogram ill condition

Work around

- Fix, or auto diff
- random restart
- Initial W randomly $\alpha \downarrow$, momentum decay $\alpha \downarrow$, iterate avg
- initial scale of w relu Adam, second order opt
- activation chard
- histogram ill condition

Minimum Singular values reflects the "compressibility" of the data, can be well-approx by a lower-dim subspace

Explicit regularization to discourage large weights reduce model's capacity, effectively smooth out the double descent curve and eliminate the second descent.

Neural Tangent Kernel - capture how network's output changes with tiny weight changes.

$\Delta Z = J \Delta W$ $\Delta W = -\alpha \nabla J(W)$ $J: \frac{\partial \mathcal{E}}{\partial W}$

$\nabla J(W) = \frac{1}{N} J^T \nabla L(Z)$ $\Delta Z = J \left[-\frac{\alpha}{N} J^T \nabla L(Z) \right] = -\frac{\alpha}{N} J^T \nabla L(Z)$

When width $\rightarrow \infty$, J, k become constant \rightarrow nn become like a linear model with a fixed kernel under certain condition, nn is guaranteed to learn the training data perfectly; work well on unseen data limitation: infinite width network, full batch gradient descent; compute NTK expensive

GD can find a global minimum of the training loss for sufficiently over parameterized nn.

Generalization

test error train error # train example # para

Technique to prevent overfitting

① Data Augmentation: augment on training set, not test set

② Reduce para: little bottleneck layer, # connection ↓

③ Weight decay ($\alpha \downarrow$ \rightarrow noise, Amply express power) the influence of individual input features, cause oscillation during training and exacerbate vanishing/exploding regularization $E_{reg} = \mathcal{E} + \lambda R = \mathcal{E} + \frac{\lambda}{2} \sum_j W_j^2$

$\tilde{W} \leftarrow \tilde{W} - \alpha \left(\frac{\partial \mathcal{E}}{\partial \tilde{W}} + \lambda \frac{\partial R}{\partial \tilde{W}} \right) = (1-\alpha\lambda) \tilde{W} - \alpha \frac{\partial \mathcal{E}}{\partial \tilde{W}}$

L² reg: $E_{reg} = \mathcal{E} + \frac{\lambda}{2} \sum_j W_j^2$ $E_{reg} = \mathcal{E} + \frac{\lambda}{2} \sum_j \|W_j\|_2^2$ encourage $N \downarrow$ but $\neq 0$ weight encourage $w=0$

④ Early stopping: w starts small and grow nn gradually becomes more non linear. mul model. consequences: • slower training • unstable gradient • unstable gradient • difficulty in hyper tuning • reduced model performance • \uparrow sensitivity to initialization

⑤ Ensembles: Avg or combine the prediction from decaying α • step decay • exponential decay $\alpha = \alpha_0 e^{-kt}$ • $1/t$ decay: $\alpha = \alpha_0 / (1+kt)$ • use a grid search on log scale (0.1, 0.01, 0.001) for diff random initialization; diff subset of training data (bagging); diff arch or hyperpara

⑥ Stochastic regularization: inject noise into network to have 0 mean and 1 variance done for each feature independently.

For a mini-batch B of size m , for each unit in a hidden layer:

Drop connect out size = (In size - K size + 2 * pad) / stride + 1

Batch normalization $\bar{z}_i = \frac{1}{m} \sum_i z_i$ mean $\mu_B = \frac{1}{m} \sum_i z_i$ normalization $\hat{z}_i = \frac{z_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$ variance $\sigma_B^2 = \frac{1}{m} \sum_i (z_i - \mu_B)^2$ scaled shift $z_i^B = \gamma \hat{z}_i + \beta$

The random in SGD updates act like a regularizer. overparameterization & Normalization

Non convexity & overparameterization

NN can generalize well even overparameterized Factors beyond capacity control play a crucial role Non convexity many local minima, saddle, plateaus

Implicit learning rate decay

Tikhonov Regularization $J_{reg}(\tilde{W}) = J(\tilde{W}) + \frac{\lambda}{2} \|\tilde{W}\|^2$

$\tilde{W} \leftarrow (1-\alpha\lambda) \tilde{W} - \alpha \frac{\partial J}{\partial \tilde{W}}$

Double descent

Forward: $y^1 \rightarrow h^1 \rightarrow y^2 \rightarrow h^2 \rightarrow y^3$ $z^t = Wh^{t-1} + ux^t$ $h^t = \phi(z^t)$

Backpropagation Through time: $x^1 \rightarrow z^1 \rightarrow h^1 \rightarrow y^1 \rightarrow r^1$ $x^2 \rightarrow z^2 \rightarrow h^2 \rightarrow y^2 \rightarrow r^2$ $x^3 \rightarrow z^3 \rightarrow h^3 \rightarrow y^3 \rightarrow r^3$ $r^t = vh^t$ $y^t = \phi(r^t)$ $L(y^1, y^2, y^3, \dots, t^1, t^2, t^3)$

Batch Norm TS Invariant to Weight Scaling

scale \tilde{W}_j by $\gamma > 0$, preactivation Z scale by γ normalized z will remain same

Effective learning rate $\tilde{W}_j - \alpha \nabla J(\tilde{W}_j^{(k)})$

$\tilde{W}_j = \tilde{W}_j / \|\tilde{W}_j\|$ $\nabla J^{(k+1)} = \frac{1}{\|\tilde{W}_j\|} \nabla J(\tilde{W}_j^{(k)})$ very small

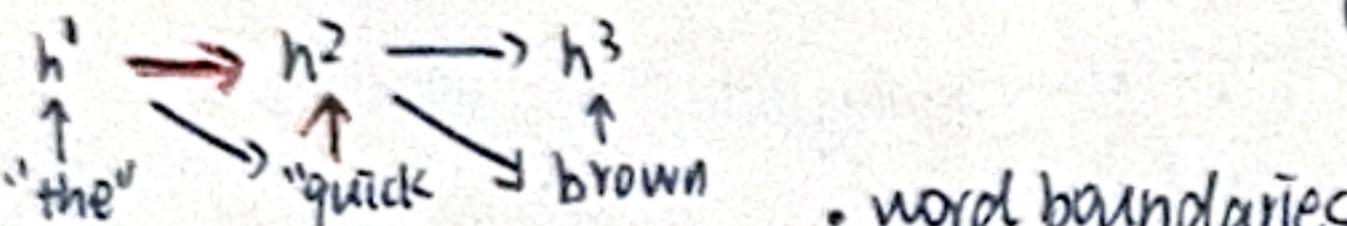
Forward: $\tilde{L} = 1 \rightarrow \tilde{y}_t = \tilde{L} L'(y_t) \rightarrow \tilde{r}_t = \tilde{y}_t \phi'(r_t)$

Back: $\tilde{L} = 1 \rightarrow \tilde{y}_t = \tilde{L} L'(y_t) \rightarrow \tilde{r}_t = \tilde{y}_t \phi'(r_t)$

$\tilde{h}^t = \tilde{r}^t \tilde{v} + \tilde{z}^{t+1} \tilde{w}$ $\tilde{z}^t = \tilde{h}^t \phi'(z^t)$

Language Modelling using RNN

Generating using model output as input next step
Training the model with teacher forcing = use target output from the training data as input



Challenge: vocabulary size challenge. ↑ sequence length

- multiplicate interactions to model character combinations - typically require deeper network limitation.
- maintain local grammar structure
- Lack global semantic coherence

Neural Machine Translation

- word by word challenge: word orders, sentence length

- need future words to resolve ambiguity
- translation require full sentence context

Sequence to Sequence Architecture

Two distinct networks with different weights

Encoder: process entire input sentence, compressing it into a code vector

Decoder: generate French translation word by word

- handle different sentence length naturally
- learns language specific word order
- capture complete sentence context

Exploding and Vanishing gradients

why RNN suffer ↑ during BP through time, gradients are multiplied by the same weight

matrix repeatedly. This repeated multiplication causes the gradients to explode or vanish.

BP = Linear operations compound multiplicatively

thin line between vanishing and exploding gradients

FP = squashes the activations

- with non-linear activation: functions
- protection against exp growth.

Iterated Logistic Map $f(x) = rx(1-x)$

composing the same function many times can result in highly non linear behavior.

Attractors: Inside each color, g vanish

At boundary, g explode

RNN must balance stability and flexibility

tanh Activation: $y = \tanh(x)$

use mini batch during training helps stabilize gradient.

Mitigating Exploding/Vanishing Gradient

① Gradient clipping

$$\text{If } \|g\| > y: \quad \tilde{g} = \frac{y}{\|g\|} g \quad \text{rescale the } g$$

Tradeoff: bias vs stable

② Trick for sequence to sequence model

- reverse the input sequence ABC → CBA
- shorten the path between matching input-output pairs
- Progressive learning: Learn short term dependencies before long term ones.

③ Identity RNN than vanilla RNN

Won't explode g. Initialize weight matrix to identity matrix.

Use ReLU activation function

Learn much longer dependency by processing the image one pixel at a time

④ Long term Short term Memory (LSTM)

$$C_{t+1} = f_t(\text{forget gate}) C_t + i_t(\text{input gate}) \text{new input}$$

- i=0, f=1: remember
- i=1, f=1: add new input to memory
- i=0, f=0: erase reset mem to 0
- i=1, f=0: overwrite replace mem with new

$$\text{Input: } z_i = \frac{e^{z_i}}{\sum e^{z_i}}$$

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum e^{z_i}}$$

$$\text{tanh} = \frac{\sinh}{\cosh} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

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

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum e^{z_i}}$$

$$\text{tanh} = \frac{\sinh}{\cosh} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

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

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum e^{z_i}}$$

$$\text{tanh} = \frac{\sinh}{\cosh} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

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

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum e^{z_i}}$$

$$\text{tanh} = \frac{\sinh}{\cosh} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

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

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum e^{z_i}}$$

$$\text{tanh} = \frac{\sinh}{\cosh} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

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

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum e^{z_i}}$$

$$\text{tanh} = \frac{\sinh}{\cosh} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

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

Autoregressive Models

RNN: Computational cost and Parallelism

Flops # of add-multiply operations for the forward and backward passes, determined by the number of connections

② Dilated convolution

Dramatically increase CNN's receptive field by using dilated convolutions with progressively larger dilation rates in successive layers

Causal limitation: only uses past info, limit ability to consider future context. Early time steps have less context available, potentially resulting in lower quality outputs at sequence beginnings

Dilated limitation: can create checkerboard artifacts if dilation factors aren't carefully designed. May miss fine grained local patterns when dilation rate is high

Attention based RNN

Efficient due to weight sharing across timesteps

Sequential complexity: # of sequential operations required to accomplish some task

Efficient due to weight sharing across timesteps

Each hidden unit is a Gated Recurrent Unit (GRU)

update & reset gates

predict one word at a time

use output as input vector

receive a context at each time step, a weighted sum of the encoder annotation

Allows decoder to focus on different parts of the input sentence when generate each step

context $c^{(1)} = \sum a_{i,j} h^{(i)}$ ← Encoder Annotation

attention weights $a_{i,j} = \frac{\exp(e_{i,j})}{\sum_j \exp(e_{i,j})}$ relevant at

Convolutional Autoregressive Model

① Causal convolution: each prediction only depends on earlier inputs in the seq.

Implemented by masking or setting zero weights on connections to future inputs

k: # of features per layer w: kernel width $\alpha(s_h, h_j) = v_a^T \tanh(W_a s_{i-1} + U_{i,j})$

additive attention

Content based addressing: annotation function depends on annotation vector $h^{(i)}$, not position

Focus on semantic, not position

Info flow in standard Encoder Decoder: the max path length is proportional to # of time steps.

Info flow in Attention Based RNN: const path length between the encoder inputs and decoder hidden

Training Comp: $K^2 dt$ Test Comp: $K^2 dt$ Test mem: $K^2 dt$ seq op: $K^2 dt$

RNN: $K^2 dt$ AttenRNN: $K^2 dt^2$ $K^2 dt^2$ $K^2 dt^2$ $K^2 dt^2$ $K^2 dt^2$