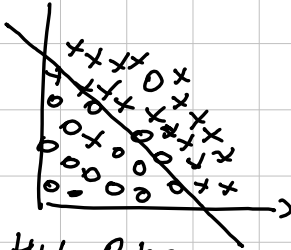


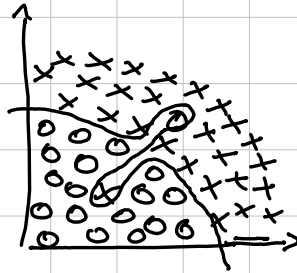
Deep Learning: A Course on Courses

- Bias / Variance / Over fit / Under fit



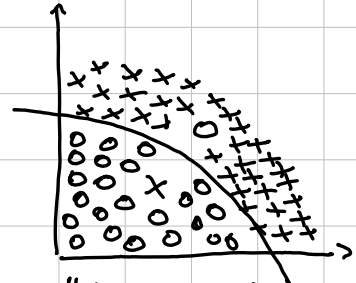
High Bias
= Underfit

- High error^x on training set
- Similar error^x on dev set



High variance
= Overfit

- Low error^x on training set
- High error^x on dev set



"Just right"

- Low-ish error^x on training set
- Low-ish error^x on dev set

^x: Relative error to optimal error called "Bayes error"^u

What to do?

- High bias:
- more powerful / bigger network / model
 - train longer
 - bigger network should be able to fit training set well (as long as somebody can do this well.)

High variance? (good train set \rightarrow bad dev set perf.)

- Get more data?
- Regularization
- Better NN architecture

Knowing if we have high bias or high variance problem tells you what to do!

- E.g. high bias: more training data will not help (as much)
 \rightarrow more powerful network or train longer/better.

Regularization

\rightarrow Helps fit overfitting / high variance

Ex: L_2 regularization for logistic regression
(add $\frac{\lambda}{2m} \|w\|_2^2$ to cost function)

For neural net: Frobenius norm of matrix

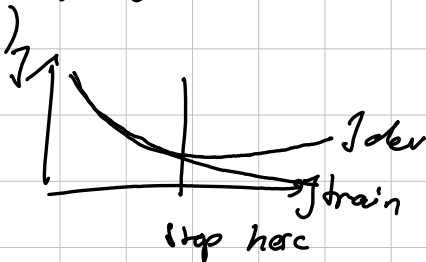
"weight decay" \rightarrow makes weight update smaller in backprop.

Dropout

- Scale layer weights by inverse of dropout probability \rightarrow inverse dropout.

Other techniques

- Augment dataset (e.g. by mirroring & distorting images, cropping)
- Early stopping



- Mixes hyperparam search and preventing overfitting
- Andrew Ng likes to separate the two.

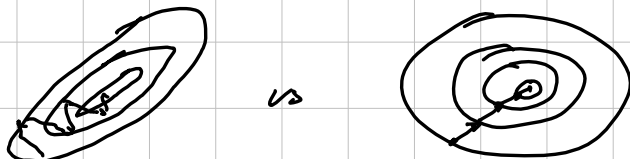
Why regularization works:

Moves weights so that network is closer to a linear function (ex. tanh activation, which is linear around 0)

(I understand)

Normalizing Inputs

Inputs on different scales but learning speed



Need lower learning rate here

⇒ normalize all inputs e.g. $\mu=0$, $\sigma=1$.

Vanishing & Exploding Gradients

Deep networks: gradients a product of many weights
if every/most steps are > 1 → exploding
 < 1 → vanishing

How to "solve" it: initialize weights "correctly":

$\text{Var}(w^i) = \frac{2}{n}$ → randomly initialize weights with this variance.



o Get the idea but don't really understand the details.

Numerical Approximation & Gradient Checking

$$\text{Linear approx} \quad \frac{f(\theta + \epsilon) - f(\theta)}{\epsilon}$$

$$\text{or} \quad \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon}$$

→ Check derivative computation against linear approx as debugging check.

$\theta = [w^1, b^1, w^2, b^2, \dots]$, compute $J(\theta)$
 $d\theta = [dw^1, db^1, \dots]$
is $d\theta$ gradient of $J(\theta)$?

$$\forall i: d\theta_{\text{approx}}^{(i)} \approx \frac{J(\theta_1, \dots, \theta_i + \epsilon, \dots) - J(\dots, \theta_i - \epsilon, \dots)}{2\epsilon}$$

use $\epsilon \approx 10^{-7}$

$$\frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2} \approx 10^{-7} \rightarrow \text{great}$$

$$\frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2} \approx 10^{-5} \quad ? \text{ ok?}$$

$$10^{-3} \rightarrow \text{worry!}$$

Notes

- Doesn't work with dropout
- Only use for debugging, not training
- To debug, try to identify faulty component (dw^i or db^i)
- Remember regularization
- Run at start & maybe after some training

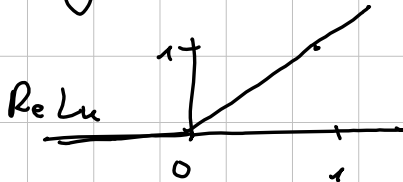
① $20.020 \overline{000} \cdot \frac{5}{1000}$ 100k obs & train ✓

② no test \rightarrow no unbiased estimated value of variance.

2023-08-30

First DeepLearning AI prof assignment.

3 layer neural net: Lin \rightarrow ReLU \rightarrow Lin \rightarrow ReLU \rightarrow Lin-Sigmoid



$$\text{ReLU}(x) = \max(x, 0).$$

Sigmoid



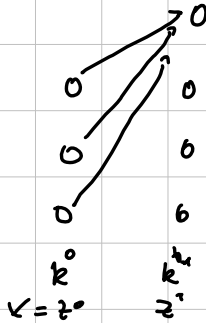
$$\text{Sigmoid}(x) = \frac{1}{1 + e^{-x}}$$

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

$$\frac{1}{1 + e^{-x}}$$

is actually the logistic function.
sigmoid is any function with this name.

Layer dims



$$z^1 = g(W_{z^0}^1 + b^1)$$

$$\dim(W^1) = \underbrace{k^1}_{\text{rows}} \times \underbrace{k^0}_{\text{cols}}$$

$$\dim(b^1) = k^1 \times 1$$

This exercise: 2 for 2 dimensions

Maybe I should do the initial course again? step through the output layers. seems like that would not hurt...

And read my own post again...

Wow, the initialization really works wonders. matrix weights initialized to $N(0, \frac{2}{k^{in}})$ direct previous layer z^1

Deep learning. AI Special - Week 1 Recap

What I learned

- Weight initialization: zero, large random, $\frac{1}{\sqrt{n}}$
→ can have significant impact on training speed.

- Bias / Variance tradeoff

High Bias: underfit.

· High error on training set

- more powerful model
- better training ← $\begin{matrix} \text{Larger} \\ \text{opt-algo} \\ \text{initialization} \end{matrix}$

High Variance: overfit

· High error on dev / eval set

- regularization ← $\begin{matrix} \text{L2} \\ \text{Dropout} \end{matrix}$
- more training data → early stopping

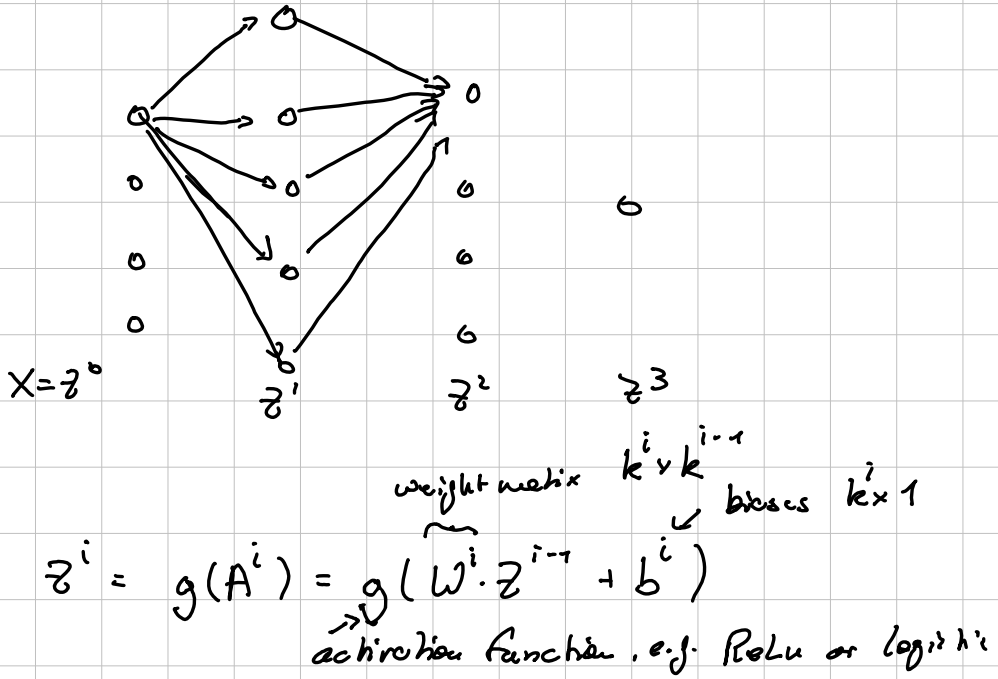
- Train / dev / test sets

└ only once, unbiased estimate.
can maybe be skipped
└ evaluate variance / overfit

Split doesn't have to be 70/30 or 60/20/20
dev and test set just need to be large enough.

- gradient checking: Compare $\frac{J(\theta+\epsilon) - J(\theta-\epsilon)}{2\epsilon}$ to $\frac{\partial J}{\partial \theta_i}$ to ensure gradient computation is correct.

Neural-net arch:



Backprop:

Just re-read my post on backprop.
It is really good! I think I could produce
quality content on this stuff. Just write
up my journey, it will show some
good shit and deepen my earnings.
How encouraging.

The details of backprop are a bit
involved, but understanding it one makes
me confident that I could again.

- Regularization: L_2 and Dropout
 - Adding L_2 and Dropout to simple neural
nets.
 - Dropout: Inverted dropout: scale
activations by $\frac{1}{\text{keep-prop}}$ to have the
same overall size.
 - Backprop: drop the same neurons
on dA_i as on A_i

2023-09-01

Last Video: Exponentially weighted averages & bias correction

$$V_t = (1-\beta) \cdot \Theta_t + \beta V_{t-1} \rightarrow \text{exp weighted avg}$$

\approx averages $\frac{1}{1-\beta}$ terms (eg. $\beta=0.02 \rightarrow 50$ terms)

Bias correction $V_t = \frac{1}{1-\beta^t} (\dots)$

corrects for initial V_0 being 0, goes to 1 for $t \rightarrow \infty$.

Gradient descent with momentum.

\rightarrow Apply exp. avg. to gradients:

$$W^e := W^e + \alpha \cdot dW^e \rightarrow V_{dW^e} = (1-\beta)V_{dW^e} + \beta dW^e$$
$$W^e \pm W^e + \alpha \cdot V_{dW^e}$$

and similar for b^r

Intuition:



Smooths out unwanted oscillations in gradient descent.

RMSProp

• Another way of having a type of momentum

• Formula for updates:

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dw^2$$

$$W = W - \alpha \frac{dw}{\sqrt{S_{dw} + \epsilon}}$$

Small smoothing, e.g. 10^{-8}

Similar for b .

Why this formula is RMSprop?

Adam optimization

→ Combination of momentum and RMSprop

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dw$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dw^2$$

$$W = W - \alpha \frac{1}{1 - \beta_1} \frac{V_{dw}}{\sqrt{S_{dw} + \epsilon}}$$

Learning rate decay

→ different schemes for decreasing the learning rate over time. e.g.

$$\alpha = \frac{1}{1 + \text{decay-rate} \times \text{epoch-number}}$$

Local optima & saddle points

MN optimization unlikely to get stuck in a local minimum - with large parameter spaces it is very unlikely that all parameters are stopped upwards. However much more likely to hit a saddle point - this can also slow down learning bc gradients are small. Momentum / Adadelta / RMSprop helps with this.

Hyperparameter Search

Typical Hyperparameters

η α : learning rate

$\beta_1, \beta_2, \epsilon$: Parameters for momentum / Adam optimizer

P_1 # Layers

P_1 # hidden units per layer

P_1 m : minibatch-size

How to search:

- ① Use random values, don't use a grid
→ random values give many more values per H.P.
- ② Consider using a coarse-to-fine approach



Choose linear or logarithmic samples correctly

e.g. learning rate $0.0001 - 1 \rightarrow$ logarithmic, we want \sim same # points between $[10^{-4}, 10^{-3}]$ as between $[10^{-1}, 10^0]$.

\Rightarrow Sample uniformly from the log space.
e.g. $x \in [-6, 0] \rightarrow 10^x$

For parameters close to 1, e.g. $\beta \in (0.9, 0.999]$
sample $1-\beta$ from logspace.

For same params, e.g. # layers and # neurons
sampling from linear space should work well.

Dados vs Caviar Strategy:

Dados: Retrain a single network

Caviar: Many different networks in parallel.

\Rightarrow Dados if you must, Caviar if you can.

Normalized Batch

→ Normalize linear units z^c in the hidden layers and add parameters β^c and γ^c as explicit scale & bias, which are also learned.

$$\tilde{z}^c = \frac{z^c - \mu^c}{\sqrt{\sigma^c + \epsilon}} \quad \mu^c, \sigma^c \text{ computed per sample / minibatch}$$

$$a^c = \underbrace{\beta^c}_{k \times 1} \cdot \underbrace{\tilde{z}^c}_{k \times 1} + \underbrace{\gamma^c}_{k \times 1} \quad (\text{elementwise product})$$

replace by minibatches.

→ eliminate b^c as we now have two bias terms

→ μ^c and σ^c normally computed for mini-batches

→ for test, usually a exp-weighted average of μ^c and σ^c is used that is computed across mini-batches.

→ Why it works? Explicitly controls bias and scale of hidden layers. Makes them change more slowly which benefits other layers which have to adapt to these changes.

Softmax Regression

For multiclass classification

Last-layer activation function $g(z) = \frac{e^z}{\sum_{i=1}^n e^{z_i}}$

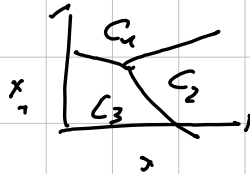
-2 normalized e^z .

-> Bit model as it is across the full layer, not element-wise

-> Generalization of logistic regression across more classes.

-> Cost function $C(y, \hat{y}) = -\sum_i y_i \log \hat{y}_i$

-> Gives linear decision boundaries for single-layer networks, e.g.



$$\frac{\partial C}{\partial z^c} = \hat{y}^c - y^c = \frac{\partial C}{\partial a^c} \cdot \frac{\partial a^c}{\partial z^c} \dots$$

ML Frameworks

- Many options to pick from
- Criteria:
 - 1) Code is easy to write + read
 - 2) Performance
 - 3) Truly open: OS + Good Governance

This course: TensorFlow.

Practical Aspects of ML Blog post

- Basically a recap of the different topics
- Links to copied lecture notes and copied notebooks
- Some thoughts on ^{my} handwritten programming assignments and how these could be better, yet doable
 - Provide step-by-step instructions, but no code
 - Provide successive links, doing a very small amount of points for accessing them.
 - Remove the training wheels over time.
- Questions on the material that go deeper
 - > but don't overdo it.