

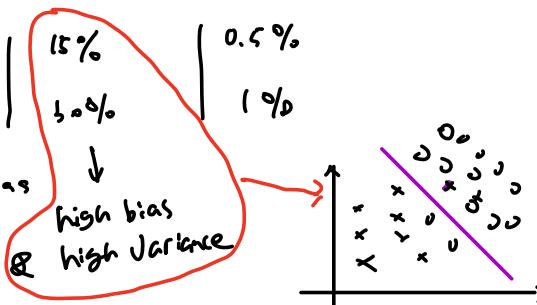
Week 1. Practical Neural Network

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

- train/dev/test sets for big data
- make sure dev and test come from the same distribution

• bias and variance

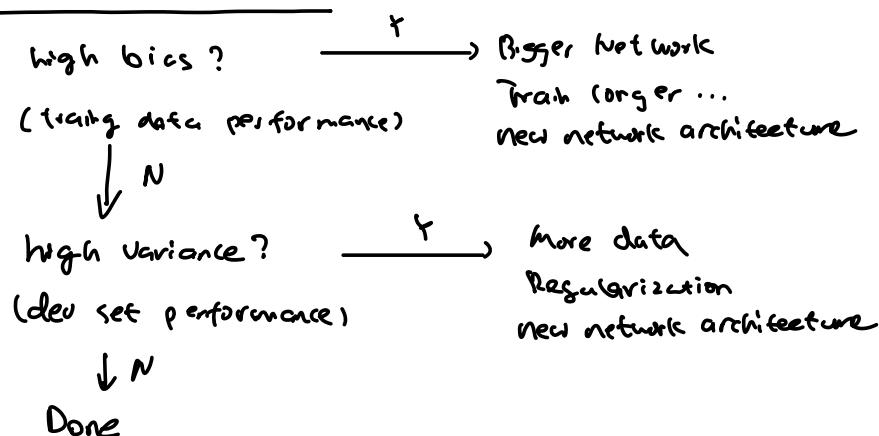
Train set error: 1% | 15%
 Dev set error: 11% | 16% | 15% | 0.5%
 ↓ ↓ ↓ ↓
 high variance high bias high bias high bias & high variance



Optimal (Bayes) error

In high dimension, less variance

• basic recipe for MC



In deep learning
 Bias Variance
 You can change one and
 doesn't hurt the other.
 (no trade-off)

Regularization Frobenius Norm: $\|A\|_F^2 = \sum_{ij} a_{ij}^2 = \text{trace}(A^T A)$

for logistics Regression

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2$$

$$\text{L}_2 \text{ regularization: } \|w\|_2^2 = \sum_{i=1}^n w_i^2 = w^T w$$

L₁ regularization: $\|w\|_1 \rightarrow w$ will be sparse

For Neural Networks

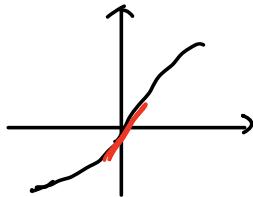
$$J(w, b, \dots, w^{(L)}, b^{(L)}) = \frac{1}{m} \sum_{i=1}^m l(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2$$

frobenius Norm.

add up every element in matrix.

Q: Why L₂-regularization can prevent overfitting?

A:



tanh.

$$\downarrow z^{(l+1)} = \underbrace{w^{(l)} a^{(l-1)} + b^{(l)}}_{\text{become linear}} \downarrow \lambda \top w^{(l)} \downarrow$$

Tips of implement:



J has new definition when you add the penalty.
So please plot the figure using the new J.

Dropout Regularization *

- Inverted dropout

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

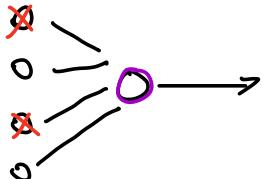
$$g_3 = np.multiply(a_3, d_3)$$

$$a_3 / = \text{keep prob} \leftarrow \text{Scaling}$$

Q: Why does "dropout" work?

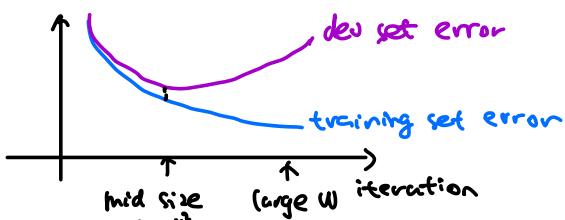
A:

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



Other regularization methods

- Data augmentation : Rotation / flipping / random cropping
- Early stopping



$\|w\|_F$

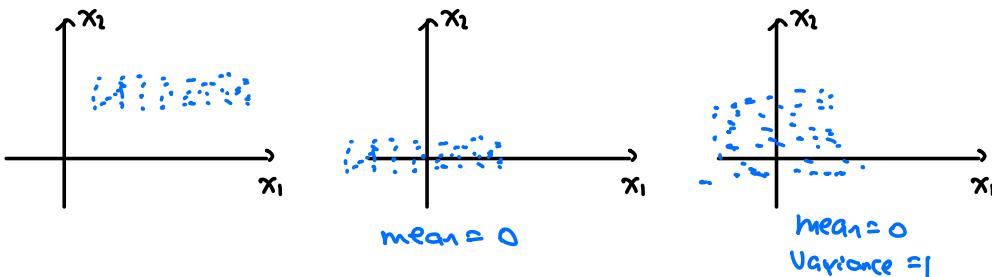
- Process of ML

- Optimize cost function J
 - Gradient descent, ...
 - Not overfitting
 - Regularization, ...
- Downside of early stopping:
Combined two tasks, not orthogonalization

Optimization

Normalizing training sets

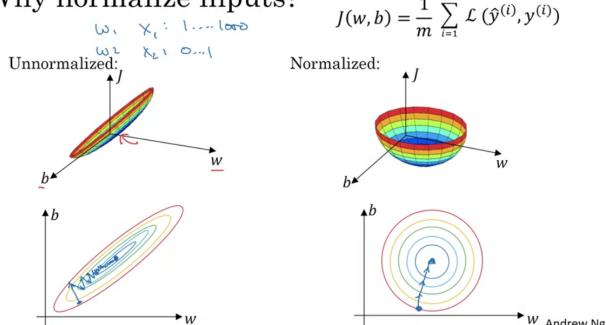
$$x = \frac{x - \mu}{\sigma}$$



- use same μ and σ to normalize training data and testing data!

Q: Why normalize input?

A: Why normalize inputs?

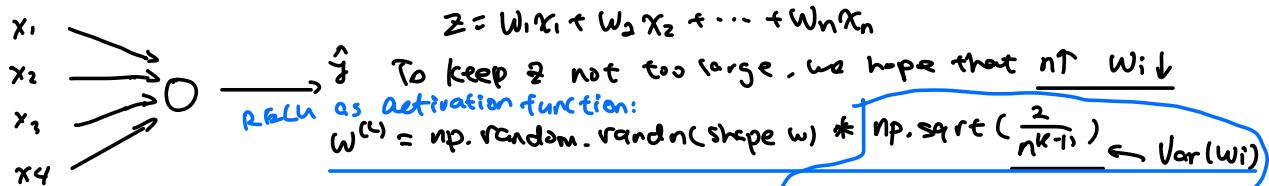


It makes the cost function faster to optimize

Vanishing / exploding gradients

- $w^T > I \rightarrow \text{explode}$
 - $w^T < I \rightarrow \text{vanish}$
- layer T , $w^T \uparrow$
- partial solution: carefully choose initial weight

Weight Initialization for deep network



tanh:
 np. sqrt ($\frac{1}{n^{e-1}}$)
 or np. square $\frac{2}{n^{e-1} \cdot f_n}$
 As a hyperparameter
 (but not priority one)

Numerical Approximation of Gradients

Aim: Prerequisite for Gradient checking (*)

$$\frac{f(0+\epsilon) - f(0-\epsilon)}{2\epsilon} \approx \underline{g(0)}$$

$$\frac{(1.01)^3 - (0.99)^3}{2(0.01)} = 3.0001 \approx 3$$

$$g(0) = 30^2 = 3$$

approx error: 0.0001
 (prev slide: 3.0301, error: 0.03)

$$f'(0) = \lim_{\epsilon \rightarrow 0} \frac{f(0+\epsilon) - f(0-\epsilon)}{2\epsilon} \quad \left| \quad \frac{f(0+\epsilon) - f(0)}{\epsilon} \right.$$

↓

error: $O(\epsilon^2)$

↓ error: $O(\epsilon)$

Conclusion: two-sided difference formula is much more accurate.

Gradient Cheesecake

Aim: Debug, see if the back prop algorithm has been performed correctly.

STEP : ① Take $w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}$ and reshape into a big vector θ
 Concatenate to a giant vector θ

$$\overline{J}(w^{(1)}, b^{(1)}, \dots, w^{(n)}, b^{(n)}) = J(\omega)$$

② Take $dW^{(1)}, dB^{(1)}, \dots, dW^{(n)}, dB^{(n)}$ and reshape into a big vector Θ

Gradient checking (Grad check)

$$\text{for each } i: \quad \frac{\partial J(\theta_0, \theta_1, \dots, \theta_i + \varepsilon, \dots) - J(\theta_0, \theta_1, \dots, \theta_i - \varepsilon, \dots)}{2\varepsilon} \approx \frac{\partial J}{\partial \theta_i}[i]$$

$$\text{Check } \rightarrow \frac{\|d_{\text{target}} - d\|_2}{\|d_{\text{target}}\|_2 + \|db\|_2} = \frac{10^{-7}}{10^{-5}} = 10^{-2} = 0.01 \text{ -- great!}$$

Gradient checking implementation Notes

- Don't use in training — only to debug
- If algorithm fails grad check, look at components to try to identify bug.
different i
- Remember regularization:
- Doesn't work with dropout. → turn off dropout, use grad check, then turn on it.
- Run at random initialization; perhaps again after some training.

Some takeaways

① Initialization

- Different initializations lead to very different results
- Random initialization is used to break symmetry and make sure different hidden units can learn different things
- Resist initializing to values that are too large
- "He initialization" works well for networks with ReLU activations

$$w_L = (\text{Layer_dims}[L], \text{layer_dims}[L-1]) \times \sqrt{\frac{2}{\text{dimension of the previous layer}}}$$

② Regularization

L_2 - regularization

- The cost computation: A regularization term is added to the cost
- The backprop function: There are extra terms in the gradients with respect to weight matrices.
- Weights end up similar ('weight decay')

Dropout

- Only use dropout during training
- Apply dropout both during forward & backward prop.
- During training time, divide each dropout layer by keep_prob to keep the same expected value for the activations.
- Regularization will drive your weights to lower values.

③ Gradient Checking

- Gradient Checking doesn't work with drop out. Work before dropout, make sure the gradient has been calculated correctly! Then turn on the gradient descent.

Optimization Algorithms

Mini-batch gradient descent

- Vectorization allows you to efficiently compute on m examples.

Assume $m = 5,000,000$, then the mini-batch is

$$X_{m,n} = \underbrace{[x^{(1)}, x^{(2)}, \dots, x^{(n)}]}_{x^{(1:n)} \text{ (n x 1000)}} \quad \underbrace{x^{(100)}, \dots, x^{(2000)}}_{x^{(100:2000)}} \quad \dots \quad \underbrace{x^{(5000)}}_{x^{(5000:n)}} \\ Y_{1:m} = \underbrace{[y^{(1)}, y^{(2)}, \dots, y^{(n)}]}_{y^{(1:n)}} \quad \underbrace{y^{(100)}, \dots, y^{(2000)}}_{y^{(100:2000)}} \quad \dots \quad \underbrace{y^{(m)}}_{y^{(5000:n)}}$$

- Algorithm

for $t = 1, 2, \dots, 5000$

Forward prop on $X^{(t)}$

$$z^{(l)} = W^{(l)} X^{(t)} + b^{(l)}$$

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

:

$$z^{(L)} = W^{(L)} X^{(t)} + b^{(L)}$$

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

Vectorized implement

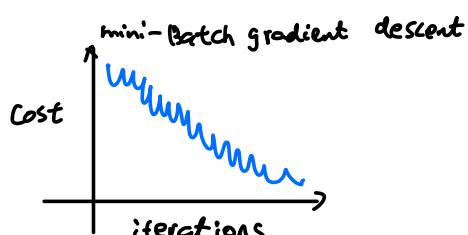
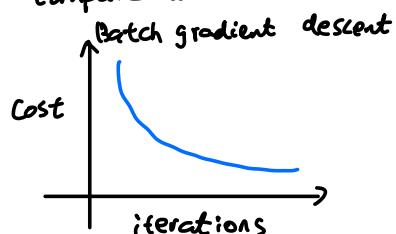
$$\text{Compute cost } J^{(t)} = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|W^{(l)}\|_F^2$$

Backprop to compute gradient $J^{(t)}$

$$w^{(l)} = w^{(l)} - \lambda d_w^{(l)}$$

$$b^{(l)} = b^{(l)} - \lambda d_b^{(l)}$$

- Compare with Batch



- How to choose mini-batch size?

If mini-batch size = $m \Rightarrow$ Batch gradient descent
 $= 1 \Rightarrow$ Stochastic gradient descent



Advantage: ① Vectorization
 ② make process without using entire training set

Guide line:

If small training set: Use batch gradient descent
 Otherwise, mini-batch size: 64, 128, 256, 512 ... (2^n)

Exponentially weighted average

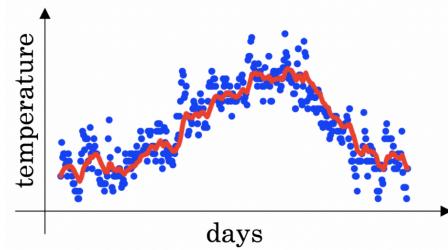
formula:

$$V_0 = 0$$

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

V_t is approximately $\frac{1}{1-\beta}$ days

- Increase β will shift the red line slightly to the right
- Decrease β will create more oscillation within the red line



Algorithm:

$$V_0 = 0$$

Repeat {

 get next θ_t

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

Bias correction Exponentially weighted Averages

Initial phase

$$\frac{V_t}{1-\beta^t}$$

Gradient descent with momentum

Algorithm:

On Iteration t :

Compute d_w, d_b on current mini-batch

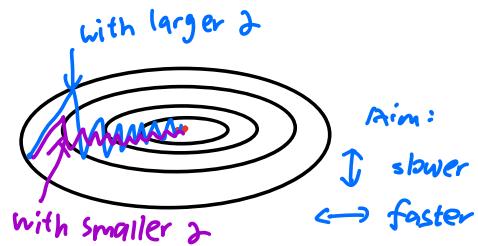
$$V_{dw} = \beta V_{dw} + (1-\beta) d_w$$

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

$$w = w - \alpha V_{dw}$$

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

Hyperparameters: α, β (0.9 is robust)



RMS prop

Algorithm:

On Iteration t :

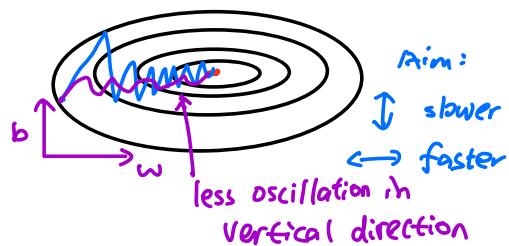
Compute d_w, d_b on current mini-batch

$$S_{dw} = \beta_1 S_{dw} + (1-\beta_1) d_w^2 \quad \text{element-wise}$$

$$S_{db} = \beta_2 S_{db} + (1-\beta_2) d_b^2 \quad \text{large}$$

$$w = w - \frac{d_w}{\sqrt{S_{dw} + \epsilon}}$$

$$b = b - \frac{d_b}{\sqrt{S_{db} + \epsilon}} \quad \text{small}$$



Then, we can choose a larger learning rate α to speedup the convergence

Adam optimization algorithm (combined momentum and RMS prop)

$$V_{dw} = 0, V_{dw}^2 = 0, V_{db} = 0, V_{db}^2 = 0$$

On Iteration t :

Compute d_w, d_b using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1-\beta_1) d_w, \quad V_{db} = \beta_2 V_{db} + (1-\beta_2) d_b$$

$$S_{dw} = \beta_1 S_{dw} + (1-\beta_1) d_w^2, \quad S_{db} = \beta_2 S_{db} + (1-\beta_2) d_b^2$$

$$V_{dw}^{\text{correct}} = \frac{V_{dw}}{1-\beta_1 t}, \quad V_{db}^{\text{correct}} = \frac{V_{db}}{1-\beta_2 t} \quad (t \text{ counts the number of steps taken of Adam})$$

$$S_{dw}^{\text{correct}} = \frac{S_{dw}}{1-\beta_1 t}, \quad S_{db}^{\text{correct}} = \frac{S_{db}}{1-\beta_2 t}$$

$$w = w - \alpha \frac{V_{dw}^{\text{correct}}}{\sqrt{S_{dw}^{\text{correct}}} + \epsilon} \quad b = b - \alpha \frac{V_{db}^{\text{correct}}}{\sqrt{S_{db}^{\text{correct}}} + \epsilon}$$

Hyperparameter choices:

α : needs to be tuned

$\beta_1: 0.9$ ($d\omega$)

$\beta_2: 0.999$ ($d\omega^2$)

$\eta: 10^{-8}$

Learning rate decay

1 epoch = 1 pass through data

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

e.g.

Epoch	α
1	0.1
2	0.067
3	0.05
4	0.04
:	:

Fixed Interval Scheduling

$$\alpha = \frac{1}{1 + \text{decay_rate} \times \left\lfloor \frac{\text{epochNum}}{\text{time Interval}} \right\rfloor} \alpha_0$$

$$\alpha_0 = 0.2$$

$$\text{decay_rate} = 1$$

hyperparameters

The problem of Local Optima

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

Hints on practicals

• Mini-batch

- Shuffling and Partitioning are two steps required to build mini-batches
(know how to handle the last mini-batch if batch-size != 0)
- Powers of two are often chosen to be the mini-batch size, e.g. 16, 32, 64, 128...

• Momentum

- Momentum takes past gradients into account to smooth out the steps of gradient descent, It can be applied with batch gradient descent, mini-batch descent, stochastic gradient descent

- You have to tune a momentum hyperparameter β and a learning rate α

- Adam
 - Relatively low memory requirements
 - Usually works well even with little tuning of hyperparameters (except β)

Week 3 Hyperparameter Tuning - Batch Normalization and Programming Frameworks

. Hyperparameter tuning

• Tuning Process

α ~~★★★~~

β ~~★★~~

$$\beta_1, \beta_2, \epsilon = 0.9, 0.999, 10^{-8}$$

layers ~~★~~

hidden units ~~★★~~

learning rate decay ~~★~~

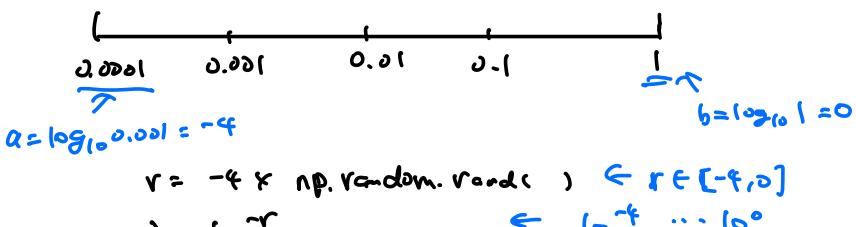
mini-batch size ~~★★~~

Δ Try random values: Don't use a grid!

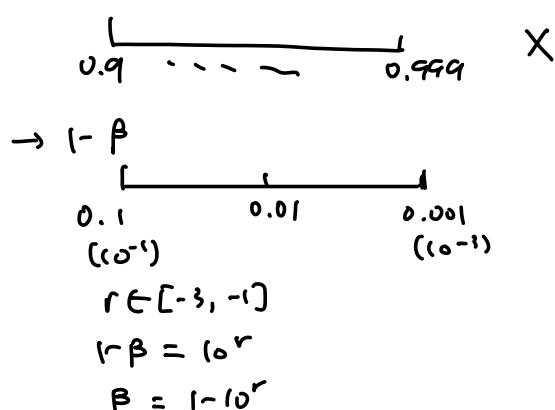
Δ Coarse to fine

• Using an appropriate scale to pick hyperparameters

log-scale()



try to find β



• Pandas model & Cuvier model

• Batch Normalization

$$x = \frac{x - \mu}{\sigma}$$

• Normalizing Activations in a network

Implement Batch Norm

Give some intermediate values in Neural Network $z^{(1)}, \dots, z^{(n)}$ $\Rightarrow z^{(l+1)}$

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

$$\sigma^2 = \frac{1}{n} \sum_i (z^{(i)} - \mu)^2$$

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

$$\tilde{z}^{(i)} = \gamma \hat{z}^{(i)}_{\text{norm}} + \beta \quad \text{learnable parameters of model}$$

If

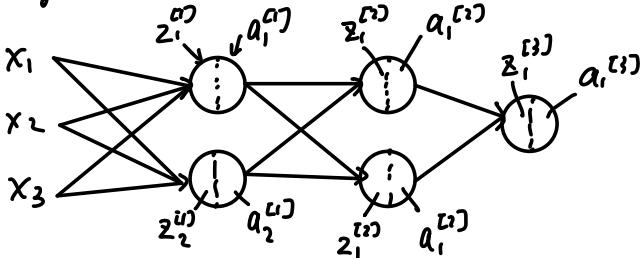
$$\gamma = \sqrt{\sigma^2 + \epsilon}, \beta = \mu$$

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

use $\tilde{z}^{(l+1)}$ instead of $z^{(l+1)}$

• Fitting Batch Norm into a Neural Network

Adding Batch norm to a network



$$x \xrightarrow{w^{(1)}, b^{(1)}} z^{(1)} \xrightarrow{\beta^{(1)}, \gamma^{(1)}} \tilde{z}^{(1)} \xrightarrow{a^{(1)} = \psi(\tilde{z}^{(1)})} z^{(2)} \xrightarrow{w^{(2)}, b^{(2)}, \beta^{(2)}, \gamma^{(2)}} \tilde{z}^{(2)} \xrightarrow{a^{(2)} = \psi(\tilde{z}^{(2)})} z^{(3)} \xrightarrow{w^{(3)}, b^{(3)}, \beta^{(3)}, \gamma^{(3)}} \tilde{z}^{(3)} \xrightarrow{a^{(3)} = \psi(\tilde{z}^{(3)})} \dots$$

$$\text{Parameter: } w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(L)}, b^{(L)}, \beta^{(1)}, \gamma^{(1)}, \dots, \beta^{(L)}, \gamma^{(L)}$$

$$* \quad z^{(l)} = w^{(l)} a^{(l-1)} + b^{(l)} \quad \leftarrow$$

$$z_{\text{normal}}^{(l)} = w^{(l)} a^{(l-1)}$$

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

Algorithm

for $t = 1, \dots, \text{num MiniBatch}$:
 Compute forward prop on $x^{(t)}$:

In each hidden layer, use BN to replace $\bar{z}^{(t)}$ with $\hat{z}^{(t)}$

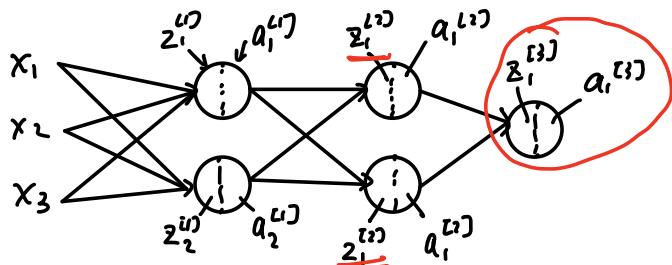
Use backprop to compute $dW^{(t)}, db^{(t)}, d\beta^{(t)}, df^{(t)}$

Update parameters $\begin{cases} W = W - \gamma dW^{(t)} \\ b = b - \gamma db^{(t)} \\ \beta = \beta - \gamma d\beta^{(t)} \\ f = f - \gamma df^{(t)} \end{cases}$ ← Can use momentum, RMSprop.
 Adam to replace GD

Why does Batch Norm work?

Covariate shift $X \rightarrow Y$

If X distribution change, then we should re-train the network



$z^{(2)}$ depends on $\bar{z}^{(1)}$ and $\bar{z}^{(2)}$ would keep changing. However, after apply batch normalization, though the change of $\bar{z}^{(2)}$ still happens, the means and variance of $z^{(2)}$ remain unchanged, its just affects by hyperparameter f and β



reduce the dependence among layers



speed up the learning in the whole network

Batch normalization as regularization

(weak, just a side effect)

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch $\bar{z}^{(t)} \rightarrow \hat{z}^{(t)}$
- This add some noise to the value $\bar{z}^{(t)}$ within that mini-batch. So similar to dropout, it adds some noise to each hidden layer's activations.
- This has a slight regularization effect.

- Use larger mini-batch size, the effect of regularization will be reduced.

Multiclass Classification

- Softmax Regression

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

Activation function:

$$t = e^{(z^{(L)})}$$

$$a^{(L)} = \frac{e^{(z^{(L)})}}{\sum_{i=1}^C t_i}, \quad a_i^{(L)} = \frac{t_i}{\sum_{i=1}^C t_i}$$

e.g.

$$\begin{aligned} z^{(L)} &= \begin{bmatrix} 5 \\ 2 \\ 1 \\ 3 \end{bmatrix} \\ t &= \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix} \\ \sum_{i=1}^4 t_i &= 176.3 \\ \overbrace{a^{(L)}}^{4 \times 1} &= \frac{t}{176.3} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix} \end{aligned}$$

In summary, Softmax: $a^{(L)} = g^{(L)}(z^{(L)})$

- Training a softmax classifier

Softmax regression generalizes logistic regression to C classes

Loss function

$$y^{(i)} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \quad a^{(i)} = y^{(i)} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$$L(\hat{y}, y) = - \sum_{j=1}^C y_j \log \hat{y}_j$$

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

Backprop:

$$\begin{aligned} d_z^{(L)} &= \hat{y} - y \\ \hookrightarrow \frac{\partial J}{\partial z^{(L)}} \end{aligned}$$

Deep Learning Framework

Tensorflow

Code example

```
import numpy as np
import tensorflow as tf

w = tf.Variable(0, dtype=tf.float32)
x = np.array([1.0, -10.0, 25.0], dtype=np.float32)
optimizer = tf.keras.optimizers.Adam(0.1)

def training(x, w, optimizer):
    def cost_fn():
        return x[0] * w ** 2 + x[1] * w + x[2]
    for i in range(1000):
        optimizer.minimize(cost_fn, [w])
    return w
w = training(x, w, optimizer)
```



`tf.data.Dataset.from_tensor_slices()`

`x_train.element_spec`

`next(iter(x_train))`

`transform → use map`

`x_train.map(normalize)`

`x = tf.constant(np.random.rand(3,1), name='x')`