

1. data set: training dat, dev, test:

data set: helps for finding high performance neural network.
 deep learning is a highly iterative process, try ideas-code-experiment to find better hyperparameter values
 -->good data setting will faster this process.

Note: one domain from one application area often do not transfer to other application areas, and hard to find good hyperparameters for one application in the start.-->highly iterative process to find hyperparameters.

1.1 data set splitted into: 3 part (data size m is small)

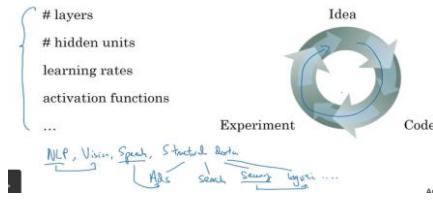
- >1. Training set: 60%, train algorithm;
- >2. dev set: 20%, evaluate algorithm performance with different models/hyperparameters
- >3. Test set: 20%, run final algorithm, get final algorithm performance

data splitting fraction:

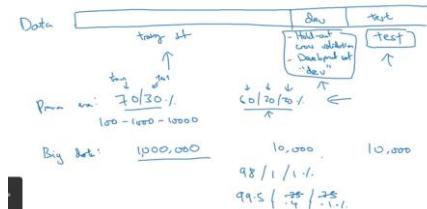
when have very large data set, dev, test set fraction could be smaller.
 as dev set just need to be big enough to evaluate different algorithm performance choices and quickly decide which one is doing better.
 (test set: just need to be big enough to have high confidence for final algorithm performance)

bigger data size, dev, test set fraction smaller. small data size more reasonable to use 60%, 20%, 20% for splitting.

Applied ML is a highly iterative process



Train/dev/test sets



1.2: train , test set mismatch distribution

>1. Make sure dev and test set come from the same distribution.
 Dev set used to evaluate different algorithm choice and decide hyperparameter, so algorithm performance showed in dev, have to keep aligned in test set./real cared data

>2. Training set may mismatch dev/test set distribution:
 to get more data for training set , may use all sorts of data. have deviation from dev/test set.

but as long as dev and test set is same distribution, the progress in machine learning algorithm will be faster

(note: dev set: is the representative of real cared object, + singular number evaluation-->decide algorithm target)

>3. Not having a test set might be ok (only dev set)
 test set is to give unbiased estimate performance of selected final algorithm . if do not need unbiased algorithm performance, might be ok do not have test set.
 (note: or dev set is large enough->algorithm will not overfit dev set: could cover all real cared data feature)

note: since algorithm have fit dev set, performance get in dev set is actually biased.

train , dev, test set also helps for judging bias, variance problem, helping for improving algorithm performance.

2. Bias and variance

Train set error:

Dev set error:

-->
Variance: algorithm's ability to generalize from training set to dev/test set. Reflected by gap between J_dev and J_train:

(note: algorithm generalizability is decided by:
 >>1. algorithm architecture: esp. when overfit training set-->regularization (influence J_dev mainly)
 >>2. training set: big enough to cover dev/test set feature: -->bigger training set (influence J_dev mainly))

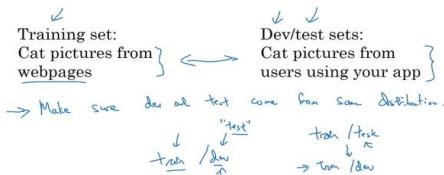
Bias: algorithm's performance on training set. reflected by gap between J_train and **Bayes error/ J_human (human error).**

may have both bias and variance error at same time: happens when overfit part of training set and underfit another part of training set. e.g: decision boundary most part is straight line -underfitting most training set, yet in the middle overfit some training examples.

judging bias or variance based on training error and dev error, assumption/conditions is: training set and dev set are from same distribution.

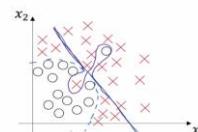
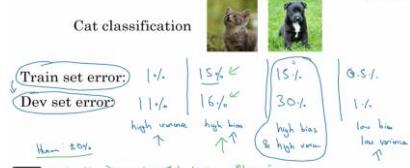
(note: if training set come from different distribution from dev set : add another 'training-dev' set to judge variance problem)

Mismatched train/test distribution



Not having a test set might be okay. (Only dev set.)

Bias and Variance



Note: 贝叶斯误差

贝叶斯误差 (bayes error rate) 是指在现有特征集上，任意可以基于特征输入进行随机输出的分类器所能达到最小误差。也可以叫做最小误差。

1. 贝叶斯误差的定义有两个关键点：

1) 给定特征集后的最小误差：即可以认为我们的训练集无限大且已经按真实分布穷举了所有可能的特征组合后，任何分类器所能达到的误差下限。产生贝叶斯误差的本质原因是特征集不足以推理出准确预测值，否则贝叶斯误差为0。

2) 概率特性：“最小误差”是限定在概率随机性条件下的（如果是开了实例级别的“天眼”的最小误差，显然无论任何情形下最小误差都是0，这个概念将没有任何意义。贝叶斯误差只能开概率级别的“天眼”。即：2.1）如果场景是给定特征输入，输出值是唯一的，则贝叶斯误差为0。例如特征是人的净高和鞋底厚度，回归目标是人的头顶距离地面高度。那贝叶斯误差为0。在训练集准确的情况下，线性回归也可以达到相同误差。2.2）反之，如果输出值是有一定随机性的，贝叶斯误差是此时选概率最大输出作为输出值所能达到的误差。比如，特征为头发长短，目标值为男、女分类label。假设在长发下，真实世界中女性概率90%，10%男性。反之女性10%，男性90%。那么贝叶斯概率误差即为认定长发为女，短发为男的误差。显然其贝叶斯误差为10%。

比较重要的是理解输入分布是真实的，但输出只能是一个值，所以会有误差。如果输出也可以是分布，那就没有误差了。

2. 贝叶斯误差的作用：

直观上可以这么理解，贝叶斯误差是在给定特征的情况下，假设数据无限（且准确），依靠统计所能得到的最小误差。它是通过增加数据集/优化数据集分布/提升模型学习能力/防止过拟合等措施后所能达到的误差下限。**如果当前算法已经能达到接近贝叶斯误差的误差，则在不动特征的前提下（比如DNN输入特征为原始文本/像素时，输入特征就动不了）。**我们已经没有继续优化的~~意义~~了。-- 然而，事实上，贝叶斯误差无法求得（因为求得的前提是你知道真实分布，你知道真实分布还做什么机器学习呀）。目前就我知道的，业界并不会真正使用该指标，该指标更多的还是出现在学术探讨中。

3. Basic recipe for machine learning

>1. first , High bias: algorithm performance on training set----->bigger network
(more hidden units/ layer number) or training longer/ new NN architecture

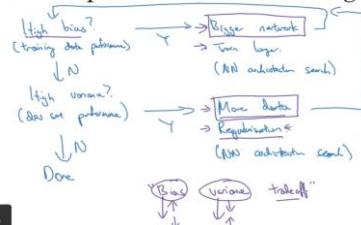
Algorithm should at least fit or overfit training set.

>2. High variance: algorithm performance on dev set ----->more data regularization, / new NN architecture

Above two steps' tool, do not hurting the other one when decreasing bias/variance, so no need to consider trade off bias and variance.

(note: does not mean do not influence. e.g: bigger network: will reduce both bias and variance.)

Basic recipe for machine learning



4. Regularization for overfitting

>add weight w in cost function: L2 regularization more popular (logistic regression) / Frobenius norm (sum of square of element in matrix)

>add weight w in cost function: L1 regularization--->
weight will be sparse, W vector will have a lot of 0 in it.?

>3. Could also add b in regularization, but normally do not do it, because:
weight W is high dimensional parameter vector (logistic regression) / matrix (NN),
especially with a high variance problem, could not fit all the parameters well, whereas b is
just a single value/ vector, so almost all the parameters are in w rather than b.
In practice, adding b does not make much difference, just one parameter over a very large
number of parameters.

>4. Regularization parameter: lambda: in python will use 'lambda' as lambda is a reserved keyword in python.

>> L2 regularization also called weight decay: as in gradient descent , L2 also involved--->
 $W := (1-\lambda/m) W - dW$
origin W matrix is multiplying a fracton <1

I1 相比于 I2 为什么容易获得稀疏解?

Neural network

$$\text{J}(\mathbf{w}^{(0)}, \mathbf{b}^{(0)}, \dots, \mathbf{w}^{(L)}, \mathbf{b}^{(L)}) = \frac{1}{m} \sum_{i=1}^m f(\mathbf{y}_i^{(i)}, \hat{\mathbf{y}}_i) + \frac{\lambda}{2m} \sum_{l=1}^L \|\mathbf{w}^{(l)}\|_F^2$$

$$\|\mathbf{w}^{(0)}\|_F^2 = \sum_{i=1}^m \sum_{j=1}^{n^{(0)}} (w_{ij}^{(0)})^2 \quad \mathbf{w}: (n^{(0)}, n^{(1)})$$

Frobenius norm

$$\|\cdot\|_F^2 \quad \|\cdot\|_F^2$$

$$\frac{\partial \text{J}}{\partial w^{(0)}} = \Delta w^{(0)}$$

$$\Delta w^{(0)} = (\text{chain backprop}) + \frac{\lambda}{m} \mathbf{w}^{(0)}$$

$$\rightarrow \mathbf{w}^{(0)} := \mathbf{w}^{(0)} - \Delta \mathbf{w}^{(0)}$$

"Weight Decay"

$$\underline{\mathbf{w}^{(0)}} = \mathbf{w}^{(0)} - \frac{1}{m} \left[(\text{chain backprop}) + \frac{\lambda}{m} \mathbf{w}^{(0)} \right]$$

$$= \mathbf{w}^{(0)} - \frac{1}{m} \left[\frac{\partial \text{J}}{\partial \mathbf{w}^{(0)}} \right] - \frac{1}{m} \left[(\text{chain backprop}) \right]$$

note:

$J_0 := f(w)$: without regularization

L1 = L0 + L1: with L1 regularization

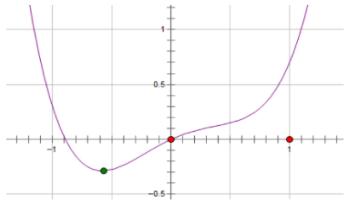
$L_2 = L_0 + L_2$: with L_2 regularization

minimize loss: try to make $dL_i/dw = 0$:

L2: $dw_2 = dw_0 + \lambda/m * w = dw_0$ (when $w == 0$)
 L1: $dw_1 = dw_0 + C * 1/1$ (when $w == 0$, could set C to make $dw_1 == 0$)

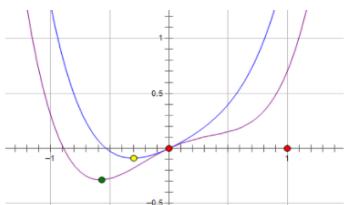
L1: $dw_1 = dw_0 + C * -1/1$ (when $w == 0$, could set C to make
 (只要系数C大于1且 $dw_0(w=0) < 0$; $w=0$ 就会变成一个极小值点。)

假设费用函数 L 与某个参数 x 的关系如图所示：



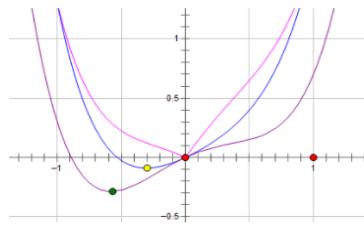
则最优的 x 在绿点处， x 非零。

现在施加 L2 regularization，新的费用函数 ($L + Cx^2$) 如图中蓝线所示：



最优的 x 在黄点处， x 的绝对值减小了，但依然非零。

而如果施加 L1 regularization，则新的费用函数 ($L + C|x|$) 如图中粉线所示：



最优的 x 就变成了 0。这里利用的就是绝对值函数的尖峰。

两种 regularization 不能把最优的 x 变成 0，取决于原先的费用函数在 0 点处的导数。

如果本来导数不为 0，那么施加 L2 regularization 后导数依然不为 0，最优的 x 也不会变成 0。

而施加 L1 regularization 时，只要 regularization 项的系数 C 大于原先费用函数在 0 点处的绝对值， $x = 0$ 就会变成一个极小值点。

上面只分析了一个参数 x ，事实上 L1 regularization 会使得许多参数的最优值变成 0，这样模型就稀疏了。

Yihua

2018-12-11

你为什么不正则化项但经过原点的loss函数来讲解会让很多人不明白。应该用：不如L1正则化项loss函数不过原点，加了L2通过调整C的值能使loss函数最小值在原点，而加了L2无济于事调整C也使函数的最小值到了原点，这样的例子来讲才能让人明白。

由 4 回复 学院 幸报

王斐 Maigo (作者) 回复 Yihua
嗯.....其实我把损失函数向上平移一下就好了.....

2018-12-11

<https://zhuanlan.zhihu.com/p/35356992>

Help understanding

L1和L2正则化：

我们所说的正则化，就是在原来的 loss function 的基础上，加上了一些正则化项或者称为模型复杂度惩罚项。现在我们还是以最熟悉的线性回归为例子。

优化目标：

$$\min \frac{1}{N} * \sum_{i=1}^N (y_i - \omega^T x_i)^2 \text{ 式子 (1)}$$

加上 L1 正则项 (lasso 回归) :

$$\min \frac{1}{N} * \sum_{i=1}^N (y_i - \omega^T x_i)^2 + C|\omega|_1 \text{ 式子 (2)}$$

加上 L2 正则项 (岭回归) :

$$\min \frac{1}{N} * \sum_{i=1}^N (y_i - \omega^T x_i)^2 + C|\omega|_2^2 \text{ 式子 (3)}$$

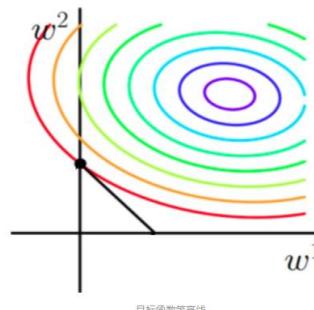
结构风险最小化角度：

结构风险最小化：在经验风险最小化的基础上（也就是训练误差最小化），尽可能采用简单的模型，以此提高泛化预测精度。

那现在我们就看看加了 L1 正则化和 L2 正则化之后，目标函数求解的时候，最终解有什么变化。

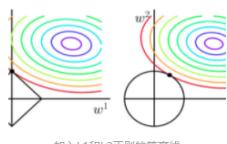
图像解释（假设 X 为一个二维样本，那么要求解参数 ω 也是二维）：

- 原函数曲线等高线（同颜色曲线上，每一组 ω_1, ω_2 带入值都相同）



目标函数等高线

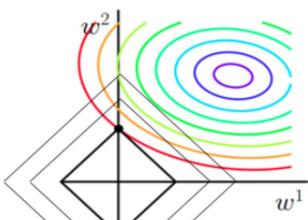
- L1 和 L2 加入后的函数图像：



加入 L1 和 L2 正则的等高线

从上边两幅图中我们可以看出：

- 如果不加 L1 和 L2 正则化的时候，对于线性回归这种目标函数凸函数的话，我们最终的结果就是最里边的紫色的小圆圈等高线上的点。
- 当加入 L1 正则化的时候，我们先画出 $|\omega_1| + |\omega_2| = F$ 的图像，也就是一个菱形，代表这些曲线上的点算出来的 1 范数 $|\omega_1| + |\omega_2|$ 都为 F 。那我们现在的目标是不仅是原曲线算得值要小（越来越接近中心的紫色圆圈），还要使得这个菱形越小越好（ F 越小越好）。那么和原来一样的话，过中心紫色圆圈的那个菱形明显很大，因此我们要取到一个恰好值。那么如何求值呢？



带 L1 正则化的目标函数求解

1. 以同一条原曲线目标等高线来说，现在以最外圈的红色等高线为例，我们看到，对于红色曲线上的每个点都可以做一个菱形，根据上图可知，当这个菱形与某条等高线相切（只有一个交点）的时候，这个菱形最小，上图相切对比较大的两个菱形对应的 1 范数更大。

用公式说这个时候能使得在相同的 $\frac{1}{N} * \sum_{i=1}^N (y_i - \omega^T x_i)^2$ 下，

由于相切的时候的 $C|\omega|_1$ 小，即 $|\omega_1| + |\omega_2|$ 小，所以：

能够使得 $\frac{1}{N} * \sum_{i=1}^N (y_i - \omega^T x_i)^2 + C|\omega|_1$ 更小。

2. 有了 1. 的说明，我们可以看出，最终加入 L1 范数得到的解，一定是某个菱形和某条原函数等高线的切点。现在有个比较重要的结论来了，我们

经过观察可以看到，几乎对于很多原函数等高曲线，和某个菱形相交的时候及其容易相交在坐标轴（比如上图），也就是说最终的结果，解的某些维度及其容易是 0，比如上图最终解是 $\omega = (0, x)$ ，这也就是我们所说的 L1 更容易得到稀疏解（解向量中 0 比较多）的原因。

5. How regularization work.

Cost function $J = J_0 + \text{regularization}$

>1. lambda is very big --> many hidden units weight is close 0, eliminating these hidden units influence, close to logistic regression (for each layer) --> simpler NN (might do not change depth) --> prone to underfitting
find a reasonable lambda --> to just fit

>Intuition 2: e.g. activation function is tanh: regularization make to use linear part of activation function - fast gradient descent.

if lambda is large --> w is smaller --> $Z = w \cdot x + b$ is smaller, reach the region of z, corresponds to activation linear part $g(z)$ is almost linear,

--> every layer activation function will be like linear function, just like linear regression, --> whole NN is like linear network, could not fit complicated non-linear decision boundary. --> prevent overfitting

Note:

>1. cost function changed when intro. regularization --> when debugging/plotting cost vs iteration, need to use new cost J , which should be single decreasing, if use old J without regularization, J curve may not show decrease, not good for debugging;

>2. When check algorithm performance on dev/test, still use original cost without regularization:

regularization is only added in cost function to reduce hidden units influence and train algorithm in training set, while real cost error still be the old formula without regularization.

(note:

Variance: --> trained algorithm ability to generalize from training set to dev set.

High variance: due to: 1. training set feature could not be representative of dev/test set, and/or algorithm learned fake features of training set; 2. algorithm learned fake features of training set due to: 'overfit' (may caused by too complicated algorithm or too small training set)

--> solution:

1. first Get more training data (more similar to dev/test set distribution - need do error analysis for mismatch ceiling performance): to cover all features of dev set, then if still have variance

2. intro. regularization: try to reduce weight of fake features algorithm has learned. - since no idea which feature is fake, reduce all features weight.

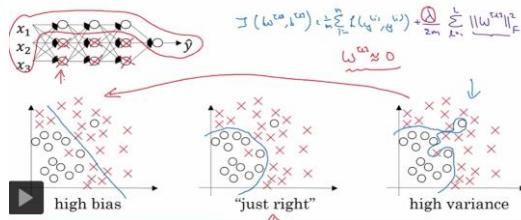
reducing weight --> simpler architecture: due to:

1. simpler computation in each layer: less units weight - may less hidden units involved - less fake features

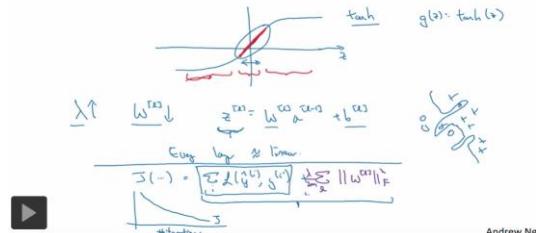
2. activation function in each layer: prone to like linear function - simpler fitting, learn less complicated function - less fake features

)

How does regularization prevent overfitting?



How does regularization prevent overfitting?



6. Drop out

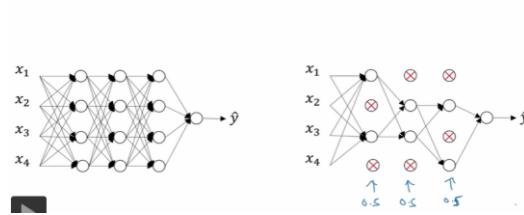
6.1 Idea: go through each layer unit and set some probability of eliminating a node in neural network --> get simpler NN

Remove nodes in layers in both forward and backward propagation (one iteration), and train NN on examples.

6.2 Function: have regularization effect as much smaller networks are being trained.

Note: nodes removed in layers is randomly --> different iteration will eliminate different nodes in layers. (even same iteration, for different batch samples, eliminated nodes also could be different acc. to possibility matrix: same dimension as a[i])
and for each example, would train it using one of these neural reduced networks.
(different iteration, use different architecture actually)

Dropout regularization



6.3 Implementing method:

>1. inverted dropout.

e.g.: adding drop out on layer 3:

dropout vector: for each example and for each hidden unit, (value is true or false), have same dimension with layer output $a[1]$ (n_l, m)

>>1. have same dimension as $A[3]$: $d3 = np.random.rand(A[3].shape[0], A[3].shape[1]) < \text{keep-prob}$

>>2. keep-prob is a real value, is the probability that a given hidden unit will be kept.

if keep-prob=0.8 --> for each example and for each hidden unit, there's 0.8 chance that the corresponding $d3$ will be one.

>>3. Activation function/layer output: $A[3] = np.multiply(A[3], d3)$ or $A[3] *= d3$ (element wise)

when use python, dropout vector is a boolean array: value is true and false rather than 1 and 0, but multiply operation works and will interpret the true and false values as one and zero.

>>4. $A[3]/=\text{keep-prob}$ -- called inverted dropout, as nomatter what value keep-prob is, $A[3]$ expected mean value remains same.

in order to not reduce the expected value of $Z[4]$? (drop out is meant to make algorithm simpler, should output less mean, why inverted to expected same mean??)

If keep-prob=0.8 --> $A[3]$ only keep 80% units, others is 0

-->1. $Z[4] = w[4] \cdot A[3] + b[4]$, to keep $E(Z[4])$ do not change, $A[3]/0.8$, to correct / just a bump that back up by the roughly 20% that need, $E[A[3]]$ also do not change much. activation function expected value do not change. (activation expectation and variance influence gradient exploring/vanishing)

-->: Inverted dropout $A[3]/=\text{keep-prob}$, makes test time easier as have less of a scaling problem.

--> even do not implement dropout at test time, the expected value of these activations $A[1]$ do not change.

do not need to add extra scaling parameter at test time, which is different from that at training time. (note: at time test, do not implement dropout, just set keep-prob=1)

Implementing dropout ("Inverted dropout")

Illustrate with layer l=3. keep-prob = 0.8
 $\rightarrow \underline{d3} = np.random.rand(a[3].shape[0], a[3].shape[1]) < \text{keep-prob}$
 $\underline{a3} = np.multiply(a[3], d3)$ # $a3 *= d3$.
 $\rightarrow \underline{a3} /= \text{keep-prob}$ ←
 50 units → 10 units short off
 $\underline{z^{(3)}} = w^{(3)} \cdot \underline{a3} + b^{(3)}$ ↓
 ↓ retain by 20% Test
 $\underline{z^{(3)}} = 0.8$

Making predictions at test time

$\underline{a^{(3)}} = X$
 No drop out
 $\underline{z^{(3)}} = w^{(3)} \cdot a^{(3)} + b^{(3)}$
 $\underline{a^{(2)}} = g^{(2)}(\underline{z^{(3)}})$
 $\underline{z^{(2)}} = w^{(2)} \cdot \underline{a^{(2)}} + b^{(2)}$
 $\underline{a^{(1)}} = \dots$
 $\underline{y} = \dots$ ↓
 ↓ keep-prob

Note: drop out

1. for different example, zero out different hidden units; dropout vector dimension is: $[n_l, m]$;

2. for each iteration, also randomly zero out different hidden units: randomly set dropout vector in every iteration: for each hidden units in each layer and each example, same dimension as $Z[1] / A[1]$

>6.5 At test time, no drop out-set keep-prob=1:

as do not want output is random, otherwise just add noise to prediction.

Alternative: run prediction on test time many times with different hidden units randomly dropped out and average cross them.--->computationally inefficient and will give roughly same result of without dropout.

7 Why drop out work

Intuition 1: on every iteration, randomly zero out hidden units, works on smaller NN;

regularization effect:

Intuition 2: can not rely on any one feature, so have to spread out weights: (weights smaller??)

for one node: $Z[l]_i : W[l]_i \cdot A[l-1] + b[l]_i$
Input of layer: $A[l-1]$: any one of input feature/ anyone of its own input, could go away at random, so reluctant to put all of its bets/ too much weight on any one input (feature), as it might go away-->spread weight to all input feature/unite.

Spreading weight tend to have an effect of shrinking the L2 of weights--->sort of regularization effect:

e.g: without drop out $W = [w1; w2]$, with drop out $W' = 1/2 [w1+w2, w1+w2]$
---> $\|W\|^2 - \|W'\|^2 = 1/2 * (w1-w2)^2 \geq 0$
---> Weight L2 with dropout, becomes smaller.

Note: L2 penalty on different weights/ are different, depending on the size of the activations being multiplied that way.

(for activation function tanh, sigmoid, smaller weight--->smaller Z---> prone to be the region where activation function is linear--->simpler whole neural network).

Summarize:

it's possible to show dropout has a similar effect to L2 regularization (**penalizing weight**), only the L2 regularization applied to different weights can be different and even more adaptive to the scale of different inputs.

(note: drop out 'even' weight for activations/hidden units in each layer, while L2 regularization do not even weight--balance cost and weight norm)

Note:

1. Keep-prob: hyperparameter, can be vary by layer, set low value for layers with many parameter/units, and 1 for layer do not want dropout.
2. Could apply dropout to input also: eliminate input feature. normally set 1 for input layer.

3. Keep-prob: is hyperparameter, need to search using dev set. (how to carry out, separately for each layer, or make a keep -prob matrix for all layers and try on dev?)

4. Drop out use more often in computer vision, where input image as too many axles/features (applied also on input?). and no enough data, so always overfitting.

>5. Dropout have regularization effect, and used to prevent overfitting.

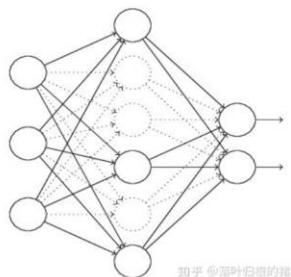
disadvange: debugging tool (plotting J vs iteration)

>1. Cost function is no longer well-defined: every iteration, randomly zero out hidden nodes. while well-defined cost function singal decrease every iteration---> loose debugging tool (plotting J vs iteration)

-->turn off dropout, run code make sure well-defined cost funtion is monotonically decreasing. -->turn on drop out-hope works well.

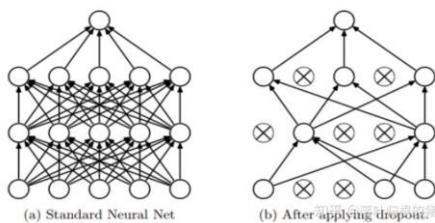
2. 什么是dropout? —— Dropout 会随机删除一些神经元, 以在不同批量上训练不同的神经网络架构。

在训练时, 每次随机 (如50%概率) 忽略隐层的某些节点; 这样, 我们相当于随机从 2^H 个模型中采样选择模型; 同时, 由于每个网络只见过一个训练数据 (每次都是随机的新网络), 所以类似bagging 的做法, 这就是我为什么将它分类到「结合多种模型」中:



此外, 而不同模型之间权值共享 (共同使用这 H 个神经元的连接权值), 相当于一种权值正则方法, 实际效果比 L2 regularization 更好。 ↗

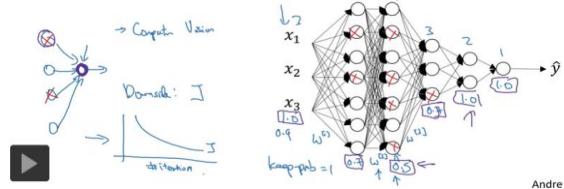
我们在前向传播的时候, 让某个神经元的激活值以一定的概率 p 停止工作, 这样可以使模型泛化性更强, 因为它不太依赖某些局部的特征。



Why does drop-out work?

云课堂

Intuition: Can't rely on any one feature, so have to spread out weights. → Shrink weights.



note:

dropout for regularization: reducing fake feature influence learned by algorithm from training set.

Reducing each feature weight of each units in each layer, by spreading weights to all units--could not rely on specify feature due to random drop out.-->

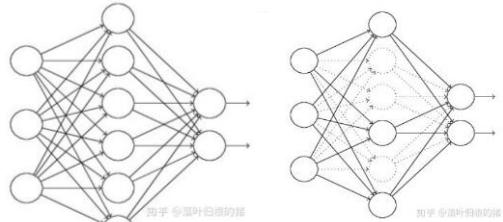
1. reducing weight norm $\|W\|_2$ -regularization effect

2. more robust to hidden units

(note: in this way, seems to force each of the units in each layer learn similar feature-redundant??

over fitting due to fake feature in hidden units-->drop out make each hidden units learn similar feature, reduce generating fake features?)

3. dropout如何工作的?



输入是 x 输出是 y , 正常的流程是: 我们首先把 x 通过网络前向传播, 然后把误差反向传播以决定如何更新参数让网络进行学习。使用Dropout之后, 过程变成如下:

(1) 首先随机 (临时) 删掉网络中一半的隐藏神经元, 输入输出神经元保持不变 (图中虚线为部分临时被删除的神经元)

(2) 然后把输入 x 通过修改后的网络前向传播, 然后把得到的损失结果通过修改的网络反向传播。一小批训练样本执行完这个过程后, 在没有被删除的神经元上按照随机梯度下降法更新对应的参数 (w, b)。

(3) 然后继续重复这一过程:

a. 恢复被删掉的神经元 (此时被删除的神经元保持原样, 而没有被删除的神经元已经有所更新)

b. 从隐藏层神经元中随机选择一个一半大小的子集临时删掉 (备份被删除神经元的参数)。

c. 对一小批训练样本, 先前向传播然后反向传播损失并根据随机梯度下降法更新参数 (w, b) (没有被删除的那一部分参数得到更新, 删掉的神经元参数保持被删除前的结果)。

不断重复这一过程。

4. 为什么dropout可以解决过拟合:

(1) 取平均的作用: 先回到标准的模型即没有dropout, 我们用相同的训练数据去训练5个不同的神经网络, 一般会得到5个不同的结果, 此时我们可以采用“5个结果取均值”或者“多数取胜的投票策略”去决定最终结果。例如3个网络判断结果集为数字9, 那么很有可能真正的结果就是数字9, 其它两个网络给出了错误结果。这种“综合起来取平均”的策略通常可以有效防止过拟合问题。因为不同的网络可能产生不同的过拟合, 取平均则有可能让一些“相反的”拟合互相抵消, dropout掉不同的隐藏神经元就类似在训练不同的网络, 随机删掉一半隐藏神经元导致网络结构已经不同, 整个dropout过程就相当于对很多个不同的神经网络取平均。而不同的网络产生不同的过拟合, 一些互为“反向”的拟合相互抵消就可以对整体上减少过拟合。

(2) 减少神经元之间复杂的共通关系: 因为dropout程序导致两个神经元不一定每次都在一个dropout网络中出现。这样权重的更新不再依赖于看固定节点的共同作用, 阻止了某些特征仅仅在其它特征存在时才有效的状况。迫使网络去学习更加鲁棒的特征, 这些特征在其它的神经元的随机子集中也存在。换句话说假如我们的神经网络是在做出某种预测, 它不应该对一些肯定的像素片段太过敏感, 即使丢失特定的像素, 它也应该可以从众多其它线索中学到一些共同的特征。从这个角度看dropout就有点像L1, L2正则, 减少权重使得网络对丢失特征神经元连接的鲁棒性提高。

(3) Dropout类似于性别在生物进化中的角色: 物种为了生存往往倾向于适应这种环境, 环境突变则会导致物种难以做出及时反应, 性别的出现可以繁衍出适应新环境的变种, 有效的阻止过拟合, 即避免环境改变时物种可能面临的灭绝。

8. Other regularization methods

8.1. More data: by data augmentation (getting more data will be expensive) by: rotating, cropping randomly, distortion
-->in this way, training set is a bit redundant, get less info. than brand new examples.

8.2. Early stopping

>1. definition: stopping the neural network training earlier.

>2. Plotting:

J_{train} vs iteration: signal decrease;

J_{dev} vs iteration: decrease then increase (overfit).

early stopping: stop iteration when dev and train error seems ok.

>3.Idea:

Weight W start from very small close to 0: randomly initialize to small value, and increase to bigger value as iteration times increase.

early stopping iteration --> weight is not so big, mid-size W, similar to L2 regularization by picking a neural network with smaller norm for parameters 2, -->overfitting less.

optimizing J and prevent overfitting is orthogonalized.

>4 disadvantage:

Machine learning optimizing Standard process: orthogonalization (正交化)
>>>1. Optimize cost function J : using gradient descent /momentum...: J_{train} go to flatten after enough iterations:
in this step , only need to focus on w and b, to make J small. no consider other things (hyperparameter, or architecture..).

>>>2. after finish J optimization, need to reduce overfitting: intro. regularization/more data, etc.

this is completely separated task to reduce variance,
optimizing J and prevent overfitting is orthogonalized: only work on one task at one time.

In early stopping, could not work on above two task independently: minimize J_{train} -bias, minimize J_{dev} -variance.

idea is using early stop fix two problem bias and variance at same time.---->make

things to try more complicated to think about.

alternative is to use L2 , then could train neural network as long as possible. whiel disadvantage is have to try a lot of values of regularization parameter lambda, -make searching over many values of lambda more computationally expensive.

>5.Advantage: running the gradient descent process just once, get to try out values of small , mid-size and large W.

Without needing to try a lot of values of L2 regularization hyperparameter lambda.

9. Normalizing training set-speed training

9.1 Normalize method:

normalizing training examples m: 输入的每个feature, 在所有样本上平均, 归一化, 以使所有样本分布在(0,1)范围内。

1st normalize mean and 2nd normalize variance for each feature

--> each feature have 0 mean, 1 variance over all examples:

Note: normalize dev/test by using same u and variance training set used for normalization.-->all data go through same transformation.

9.2: normalization training set -->motivation:

when feature are all on similar scales, cost function will be more round and easier to optimize.

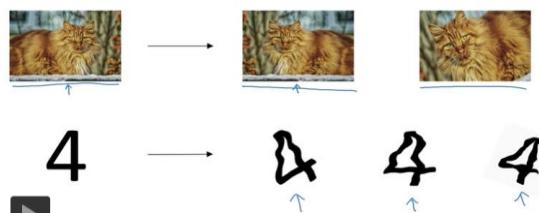
make cost function more symmetric, -->less sensitive about initialization point/start point for gradient descent and could use bigger learning rate --> less step to find optima, algorithm run faster.: Derivative of J over each parameter w_i is similar-->gradient descent step for each parameter is similar-->less sensitive to start point and total cost move forward faster

if cost function is not symmetric, while doing gradient descent, start point/initialization weight is important (derivative on different start point may differ much, and derivative direction not always point to optima, gradient descent step also differ much) and the tour to final optima is more tortuous.

-->when features scales from a large range, necessary to do normalization. or if features already in similar range (around 0-1), then do not need normalization -yet better normalize also.

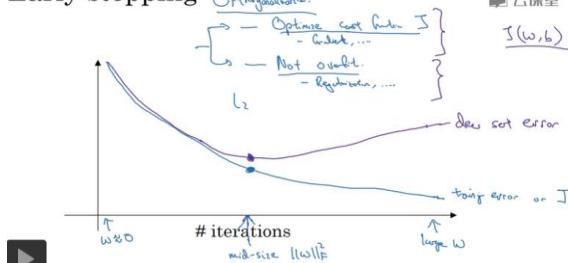
Data augmentation

云课堂

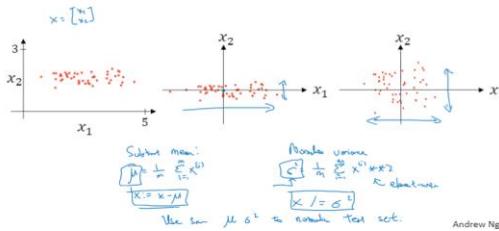


Early stopping

云课堂



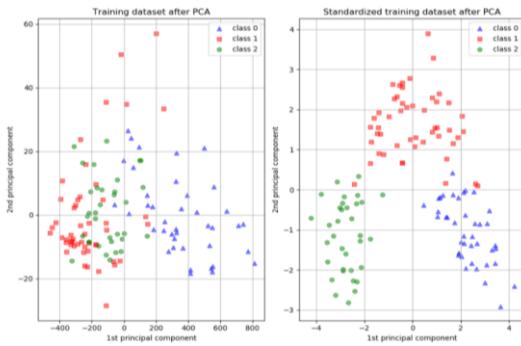
Normalizing training sets



normalization

什么时候需要feature scaling?

- 涉及或包含距离计算的算法，比如K-means、KNN、PCA、SVM等，一般需要feature scaling，因为
- zero-mean一般可以增加样本间余弦距离或者内积结果的差异，区分力更强，假设数据集中分布在第一象限远点的右上角，将其平移到原点处，可以想象样本间余弦距离的差异被放大了。在模版匹配中，zero-mean可以明显提高响应结果的区分度。
- 就欧式距离而言，增大某个特征的尺度，相当于增加了其在距离计算中的权重。如果有明确的先验知识表明某个特征很重要，那么适当增加其权重可能对向好。但如果没用这样的先验，或者目的就是想知道哪些特征更重要，那么就需要先feature scaling，对各特征等权视之。
- 增大尺度的同时增大了该特征维度上的方差，PCA算法倾向于关注方差较大的特征所在的坐标轴方向，其他特征可能会被忽视。因此，在PCA前做Standardization效果可能更好，如下图所示，图片来自sklearn learn-Importance of Feature Scaling。



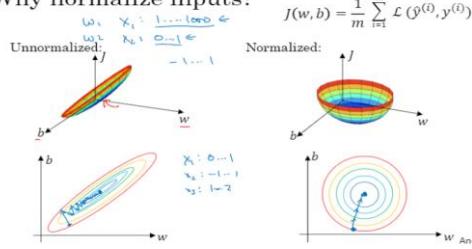
多维情况下可以分解成多个上图，每个维度上分别下降，参数 W 为向量，但学习率只有一个，即所有参数维度共用同一个学习率（暂不考虑为每个维度都分配单独学习率的算法）。“收敛”意味着在每个参数维度上都取得极小值，每个参数维度上的偏导数都为0，但是每个参数维度上的下降速度是不同的，为了每个维度上都能收敛，学习率应取所有维度在当前位置合适步长中最小的那个。下面讨论feature scaling对gradient descent的作用。

- zero center与参数初始化相配合，缩短初始参数位置与local minimum间的距离，加快收敛。模型的最终参数未知的，所以一般随机初始化，比如从0值的均匀分布或高斯分布中采样得到，对线性模型而言，其分界面初始位置大致在原点附近，bias经常初始化为0，则分界面直接通过原点。同时，为了收敛，学习率不会很大，而每个数据集的特征分布是不一样的，如果其分布集中且距离原点较远，比如位于第一象限远点的右上角，分界面可能需要花费很多步骤才能“爬到”数据集所在的位置。所以，无论什么数据集，先平移到原点，再配合参数初始化，可以保证分界面一定会穿过数据集。此外，outliers常分布在数据集的外围，与分界面从外部向内部运动相比，从中心区域开始运动可能受outliers的影响更小。
- 对于采用均方误差损失LMS的线性模型，损失函数恰为二阶，如下图所示

$$E(W) = \frac{1}{2P} \sum_{p=1}^P \left| d^p - \sum_i w_i x_i^p \right|^2$$

不同方向上的下降速度变化不同（阶数不同，曲率不同），恰由输入的协方差矩阵决定。通过 scaling改变了损失函数的形状，减小不同方向上的曲率差异。将每个维度上的下降分解来看，给定一个下降步长，如果不小心，有的维度下降的多，有的下降的少，有的还可能在上升，损失函数的整体表现可能是上升也可能下降，就会不稳定。scaling后不同方向上的曲率相对更接近，更容易选择到合适的学习率，使下降过程相对更稳定。

Why normalize inputs?



request%255Fid%2522%253A%2522162873938816780271537906%2522%252C%2522scm%2522%253A%252220140713.130102334.pc%255Fall.%2522%2527D&request_id=162873938816780271537906&bir_id=0&utm_medium=distribute.pc_search_result.none-task-blog-

- 损失函数中含有正则项时，一般需要feature scaling：对于线性模型 $y = w \cdot x + b$ 而言， x 的任何线性变换（平移、放缩），都可以被 w 和 b 吸收掉，理论上，不会影响模型的拟合能力。但是，如果损失函数中含有正则项，如 $\lambda |w|^2$ ， λ 为超参数，其对 w 的每一个参数施加同样的惩罚，但对于某一根特征 x_i 而言，其 scale 越大，系数 w_i 越小，其在正则项中的惩罚就会变小，相对于对 w ，惩罚变小，即损失函数会相对忽视那些 scale 增大的特征，这并不合理，所以需要 feature scaling，使损失函数平等看待每一根特征。

- 梯度下降法，需要feature scaling。梯度下降的参数更新公式如下，

$$W(t+1) = W(t) - \eta \frac{dE(W)}{dW}$$

$E(W)$ 为损失函数，收敛速度取决于：参数的初始位置到 local minima 的距离，以及学习率 η 的大小。一维情况下，在 local minima 附近，不同学习率对梯度下降的影响如下图所示，

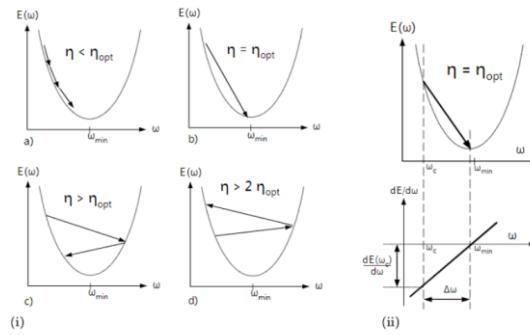


Fig. 6. Gradient descent for different learning rates.

多维情况下可以分解成多个上图，每个维度上分别下降，参数 W 为向量，但学习率只有1个，即所有参数维度共用同一个学习率（暂不考虑为每个维度都分配单独学习率的算法）。**收敛**意味着在每个参数维度上都取得极小值，每个参数维度上的偏导数都为0，但是每个参数维度上的下降速度是不同的，为了每个维度上都能收敛，学习率应取所有维度在当前位置合适步长中最小的那个。**下面讨论feature scaling对gradient descent的作用。**

• 另有从Hessian矩阵特征值以及condition number角度的理解，详见Lecun paper-Efficient BackProp中的Convergence of Gradient Descent一节，有清晰的数学描述，同时还介绍了白化的作用——解除特征间的线性相关性，使每个维度上的梯度下降可独立看待。

• 文章开篇的椭圆形和圆形等高线图，仅在采用均方误差的线性模型上适用，其他损失函数或更复杂的模型，如深度神经网络，损失函数的error surface可能很复杂，并不能简单地用椭圆和圆来刻画，所以用它来解释feature scaling对所有损失函数的梯度下降的作用，似乎过于简化，见Hinton video-3.2 The error surface for a linear neuron。

• 对于损失函数不是均方误差的情况，只要权重 w 与输入特征 x 间是相关关系，损失函数对 w 的偏导必然会有因子 x_i ， w 的梯度下降速度就会受到特征 x 尺度的影响。理论上为每个参数都设置上自适应的学习率，可以吸收掉 x 尺度的影响，但在实践中出于计算量的考虑，往往还是所有参数共用一个学习率，此时 x 尺度不同可能会导致不同方向上的下降速度悬殊较大，学习率不容易选择，下降过程也可能不稳定。通过 scaling 可对不同方向上的下降速度有所控制，使下降过程相对更稳定。

• 对于传统的神经网络，对输入做feature scaling也很重要，因为采用sigmoid等有饱和区的激活函数，如果输入分布范围很广，参数初始化时没有适配好，很容易直接陷入饱和区，导致梯度消失。所以，需要对输入做Standardization或映射到 $[0, 1], [-1, 1]$ ，配合精心设计的参数初始化方法，对域进行控制。但自从有了Batch Normalization，每次线性变换改变特征分布后，都会重新进行Normalization，似乎可以不太需要对网络的输入进行feature scaling了？但习惯上还是会做feature scaling。

什么时候不需要Feature Scaling?

- 与距离计算无关的概率模型，不需要feature scaling，比如Naive Bayes；
- 与距离计算无关的基于树的模型，不需要feature scaling，比如决策树、随机森林等，树中节点的选择只关注当前特征在哪里切分对分类更好，即只在意特征内部的相对大小，而与特征间的相对大小无关。

10: Gradients vanishing/exploding

Definition: in deep NN, while training, derivative/slop sometimes get either very big or very small.

E.g: activation function for all layers: $g(z) = z$ linear \rightarrow

$y = W[L] \cdot W[L-1] \cdots W[1] \cdot x$

make $W[i]$ all = [1.5, 0, 0, 1.5] $\rightarrow y = W[L] \cdot x$

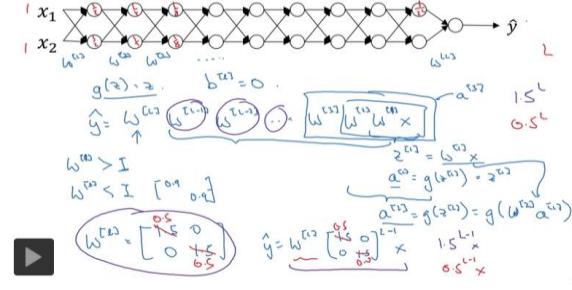
in a very deep network, output y value will explode;

if $w[i] = [0.5, 0, 0, 0.5] \rightarrow$ in a very deep network, output y will vanish

in a very deep network, if weight is little bigger than 1, output y value/activation value will explode; if weight value is a bit less than 1, output value will vanish.

same argument also applied to derivatives, exploding/vanishing as function of the layer numbers. \rightarrow difficult for learning, as gradient descent step too big or too small.

Vanishing/exploding gradients



11. Weight initialization for deep networks

Weight initialization: partial solution for gradient descent exploding/vanishing: does not solve it entirely, but helps a lot;

By better or more careful choice of the random initialization for neural network.

\rightarrow make derivative of activation function in each hidden layer, close to 1

e.g:

forward: $A_1 = g_1(W_1 \cdot x)$, $A_2 = g_2(W_2 \cdot A_1) \dots A[L] = g_L(W_L \cdot A[L-1])$

backward: $dA[L-1] = W[L] \cdot dA[L]$, $dA[L-2] = W[L-1] \cdot dA[L-1] \cdot g'_L \cdot (L-1) = W[L-1] \cdot W[L] \cdot dA[L]$, \dots , $dA[1] = g'_1 \cdot (L-1) \dots$

1. output/activation of each layer will not explode, when $g[i]$ not explode /vanish (hidden layer activation function: sigmoid, tanh, $\rightarrow -1 <= g[i] <= 1$, ok for not exploding tanh not exploding only when $z = wx + b: \sim 1$)

2. derivative will not explode or vanish only when derivative of activation function $g'[i]$ close to 1.

(note: so more often choose hidden layer activation function Relu)

11.1 : Singal neuron example

$z = w_1 \cdot x_1 + \dots + w_m \cdot x_n$,

in case Z is too large or too small \rightarrow larger is n , smaller is w_i

could set weight variance:

$V(w) = 1/\text{sqrt}(n)$ (n is the input feature of this neuron) for tanh activation function

$V(w) = 2/\text{sqrt}(n)$ for Relu activation function

$W[i] = \text{np.random.randn}(\text{shape}) * \text{np.sqrt}(1 \text{ or } 2/n-1)$

\rightarrow if input feature have 0 mean and 1 variance, then Z also take on a similar scale with above initialized w .

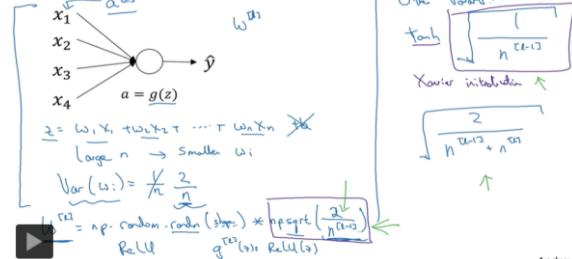
(note: $D(z) = D(w_1) + \dots + D(w_n) + D(x_1) + \dots + D(x_n)$; in conditions W_i and X_i are independent) does not solve but definitely helps reduce the vanishing, exploding gradients problem, as set each layer weight not too bigger than 1, and not too much less than 1.

Note: weight variance could be another hyperparameter to tune-but have low priority than other hyperparameters.

$W[i] = \text{np.random.randn}(\text{shape}) * \text{np.sqrt}(v(w))$

$E(X), D(X)$ definition:

Single neuron example



$$E(X) = 0, D(X) = 1 = E((X-E(X))^2) = E(X^2) - [E(X)]^2$$

$$E(Z) = E(w \cdot x), \text{ if } w, x \text{ unrelated}, \rightarrow E(Z) = E(w) \cdot E(x) = 0$$

$$D(Z) = D(w \cdot x), \text{ if } w, x \text{ unrelated} \rightarrow D(Z) = E(W^2)$$

$$D(W) = E((W-E(W))^2)$$

$$D(X+Y) = D(X) + D(Y) \text{ if } X, Y \text{ are independent.}$$

<https://www.zhihu.com/search?type=content&q=%E6%A6%82%E7%8E%87%E7%BB%9F%E8%AE%A1%20%E6%9C%9F%E6%9C%9B%20%E6%96%B9%E5%B7%AE>

一、数学期望

(加权平均值)

1.1 离散型

$$X = \sum_{k=1}^n p_k \cdot x_k$$

1.2 连续型随机变量

$$E(X) = \int_{-\infty}^{\infty} x f(x) dx$$

二、数学期望的性质

- (1) $E(c) = c$, c 是常数;
- (2) $E(X+c) = E(X)+c$, c 是常数;
- (3) $E(cX) = c \cdot E(X)$, c 是常数;
- (4) $E(X+Y) = E(X)+E(Y)$,
- (5) $\text{若 } X, Y \text{ 相互独立, 则 } E(XY) = E(X)E(Y)$.

三、方差

意义: 随机变量偏离期望的程度

定义: 偏差的平方的期望。

计算方法: $E[(X-E(X))^2]$

1. 离散型:
 $D(X) = (x_1-E(X))^2 \cdot p_1 + (x_2-E(X))^2 \cdot p_2 + \dots = \sum_{k=1}^n (x_k-E(X))^2 \cdot p_k$

2. 连续型:
 $D(X) = \int_{-\infty}^{\infty} (x-E(X))^2 \cdot f(x) dx$

注意:

$$\therefore D(X) = E(X^2) - [E(X)]^2$$

性质:

- (1) $D(c) = 0$, c 是常数;
- (2) $D(X+c) = D(X)$, c 是常数;
- (3) $D(cX) = c^2 \cdot D(X)$, c 是常数;
- (4) $D(X \pm Y) = D(X) + D(Y) \pm 2E[(X-E(X))(Y-E(Y))]$
- (5) $\text{若 } X, Y \text{ 相互独立, 则 } D(X \pm Y) = D(X) + D(Y)$.

四、常见分布的期望与方差

分布	参数	分布律或概率密度	数学期望	方差
0-1 分布	p	$P(x=k) = p^k(1-p)^{1-k}$ $\{k=0,1\}$	p 知乎 @fzshi如烟海	$p(1-p)$

二项分布	n, p	$P(x=k) = C_n^k p^k (1-p)^{n-k}$	np 知乎 @fzshi如烟海	$np(1-p)$
------	--------	----------------------------------	----------------------	-----------

泊松分布	λ	$P(x=k) = \frac{\lambda^k e^{-\lambda}}{k!}$	λ 知乎 @fzshi如烟海	λ
------	-----------	--	---------------------------	-----------

均匀分布	a, b	$f(x) = \frac{1}{b-a}, (a < x < b)$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$
正态分布	μ, σ	$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$	μ	σ^2
指数分布	λ	$f(x) = \begin{cases} \lambda e^{-\lambda x}, & x > 0 \\ 0, & \text{其他} \end{cases}$	知乎 @fzshi如烟海	$\frac{1}{\lambda^2}$

深度学习中常见的梯度爆炸

梯度消失和梯度爆炸

反向传播的详细推导可以参见：

小明说：反向传播算法详解
zhuanlan.zhihu.com



从反向传播的推导公式我们知道

$$dW^{[l]} = \frac{1}{m} dZ^{[l]} \cdot A^{[l-1]T} \quad (1)$$

$$db^{[l]} = \frac{1}{m} np.sum(dZ^{[l]}, axis=1, keepdim=True) \quad (2)$$

$$dZ^{[l]} = W^{[l+1]T} dZ^{[l+1]} * g^{[l]}(Z^{[l]}) \quad (3)$$

从上面式子（2）我们得到了反向传播的核心递推公式。假设网络共 L 层，我们计算第 l 层权重的梯度 $dW^{[l]}$ 需要计算 $W^{[l+1]T} \cdot W^{[l+2]T} \cdots W^{[L]T}$ 和 $g^{[l]}(Z^{[l]}) \cdot g^{[l+1]}(Z^{[l+1]}) \cdots g^{[L-1]}(Z^{[L-1]})$ 。

可以看到梯度的大小与两个因素相关：权重矩阵和激活函数。

举个极端例子，假如权重矩阵都是一维的，且权重矩阵的值都是 t ，那么：

$$W^{[l+1]T} \cdot W^{[l+2]T} \cdots W^{[L]T} = t^{L-l}$$

那么对于一些比较深的网络，对于一些靠近输入层的层， $L-l$ 较大。若 $|t| < 1$ ，则 $|t^{L-l}|$ 的值接近于 0，导致 $dW^{[l]}$ 和 $db^{[l]}$ 接近于 0，发生梯度消失。若 $|t| > 1$ ，则 $|t^{L-l}|$ 的值很大，导致 $dW^{[l]}$ 和 $db^{[l]}$ 很大，发生梯度爆炸。

另外从激活函数方面来看，若激活函数的导数绝对值都小于 1，那么导致梯度消失；若激活函数的导数绝对值都大于 1，那么导致梯度爆炸。这也是为啥 ReLU 激活函数表现好的原因，因为梯度刚好是 1。

下面我们来讨论一下深度学习中常用的参数初始化方法，主要参考 Keras 的官方文档。

首先我们讨论将权值矩阵初始化为常数可以吗？

Constant

keras.initializers.Constant(value=0)

将权量初始值设为一个常数的初始化器。

每层所有神经元相等，效果等价于一个神经元，这无疑极大限制了网络的能力。

特别地，全 0 初始化，根据式（3），反向传播时所有的梯度为 0，权重不会更新，一直保持为 0。

很明显，在深度学习中将权值矩阵初始化为常数是不合适的。

通过上面的讨论，全 0、常数、过大、过小的权重初始化都是不好的，那什么样的初始化是好的呢？

什么样的初始化是好的？

• 因为对参数 w 的大小和正负缺乏先验知识， w 应该为随机数，且期望 $E(w) = 0$ ；

• 在限定期望 $E(w) = 0$ 后，为了防止梯度消失和梯度爆炸，参见公式（3），权重不易过大或过小，所以要对权重的方差 $Var(w)$ 有所控制；

• 参见公式（1）， $dW^{[l]}$ 还与 $A^{[l-1]T}$ 相关，所以我们希望不同激活层输出的方差相同，即 $Var(a^{[l]}) = Var(a^{[l-1]})$ ，这也意味不同激活层输入的方差相同，即 $Var(z^{[l]}) = Var(z^{[l-1]})$ ；

• 权重初始化时，权重的数值范围（方差）应考虑到前向和后向两个过程。数值太大，前向时可能陷入饱和区，反向时可能梯度爆炸；数值太小，反向时可能梯度消失。

glorot (Xavier)

论文地址：[Understanding the difficulty of training deep feedforward neural networks](#)

核心思想：正向传播时，激活值的方差保持不变；反向传播时，关于状态值的梯度的方差保持不变。

glorot (Xavier)

论文地址：[Understanding the difficulty of training deep feedforward neural networks](#)

核心思想：正向传播时，激活值的方差保持不变；反向传播时，关于状态值的梯度的方差保持不变。

假设现在有一批输入数据 $X = \{x^{(1)}, x^{(2)}, \dots\}$, $x^{(1)}, x^{(2)}, \dots$ 是独立同分布的，且 $E(x) = 0$ 。

假设同一层（例如第 l 层有 M 个输入节点， N 个输出节点）的权重

$$w_{11}^{[l]}, w_{12}^{[l]}, \dots, w_{ij}^{[l]}, \dots, w_{NM}^{[l]}$$

因为每层的权重 $w_{11}^{[l]}, w_{12}^{[l]}, \dots, w_{ij}^{[l]}, \dots, w_{NM}^{[l]}$ 和输入 $a_1^{[l-1]}, \dots, a_m^{[l-1]}, \dots, a_N^{[l-1]}$ 相互独立，所以两者之积 $w_{ij}^{[l]} a_m^{[l-1]}$ 构成的随机变量 $z_1^{[l]}, \dots, z_n^{[l]}, \dots, z_N^{[l]}$ 亦相互独立，且同分布。

因为 $E(x) = 0$ ， $E(w_{ij}^{[l]}) = 0$ ，所以 $E(a_m^{[l]}) = 0$ 。

可以推得： $Var(w_{ij}^{[l]} a_m^{[l-1]}) = Var(w_{ij}^{[l]}) Var(a_m^{[l-1]})$

因为初始化时期望 $E(w) = 0$ ， w 在 0 附近，而且 $E(x) = 0$ ，所以 $g'(z_n^{[l]}) \approx 1$ ， $tanh$ 在初始化时激活函数的活动范围在 0 附近类似于线性函数。

$$Var(dW^{[l]}) = Var(x) Var(dw^{[l]}) \prod_{i=1}^{l-1} n^{[i-1]} Var(W^{[i]}) \prod_{i=l-1}^L n^{[i]} Var(W^{[i]})$$

he初始化

论文地址：[Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification](#)

核心思想：正向传播时，状态值的方差保持不变；反向传播时，关于激活值的梯度的方差保持不变。

Xavier 假定激活函数在 0 附近为线性函数，这对于 ReLU 是不成立的。**he 初始化**是针对 ReLU 激活函数的初始化方法。

因为 $W^{[l]}$ 是 0 均值，所以在前向传播过程中有：

$$Var(z^{[l]}) = n^{[l-1]} Var(W^{[l]} a^{[l-1]}) = n^{[l-1]} Var(W^{[l]}) E((a^{[l-1]})^2)$$

其中 $(a^{[l-1]})^2$ 表示 $a^{[l-1]}$ 的平方。因为激活函数是 ReLU 函数，所以 $a^{[l-1]}$ 不是 0 均值的。所以 $E((a^{[l-1]})^2) \neq Var(a^{[l-1]})$ 。

如果我们令 $W^{[l-1]}$ 为以 0 为中心的对称分布且 $b^{[l-1]} = 0$ ，则 $z^{[l-1]}$ 也将是以 0 为中心的对称分布，且 $E(z^{[l-1]}) = 0$ 。

那么当激活函数为 ReLU 时有： $E((a^{[l-1]})^2) = \frac{1}{2} E((z^{[l-1]})^2) = \frac{1}{2} Var(z^{[l-1]})$

那么有 $Var(z^{[l]}) = n^{[l-1]} Var(W^{[l]}) E((a^{[l-1]})^2) = \frac{1}{2} n^{[l-1]} Var(W^{[l]}) Var(z^{[l-1]})$

所以 $Var(z^{[l]}) = Var(z^{[1]}) \prod_{i=1}^l \frac{1}{2} n^{[i-1]} Var(W^{[i]})$

为了使前向传播不同激活层输出的方差相同，反向传播中不同层输入的梯度相同

即有对于任意的 i, j , 有:

$$Var(a^{[l]}) = Var(a^{[l]})$$

$$Var(dz^{[l]}) = Var(dz^{[l]})$$

解得对于同一个 $W^{[l]}$: $\forall i, n^{[l-1]} Var(W^{[l]}) = 1, n^{[l]} Var(W^{[l]}) = 1$

综合以上两式得 $\forall i, Var(w^{[l]}) = \frac{2}{n^{[l-1]} + n^{[l]}}$.

如果我们使用均匀分布来初始化，假设均匀分布的范围为 $[-a, a]$ ，则均匀分布的方差为:

$$Var(uniform) = \frac{[a - (-a)]^2}{12} = \frac{a^2}{3}.$$

$$\Leftrightarrow \frac{a^2}{3} = \frac{2}{n^{[l-1]} + n^{[l]}} \text{ 脱得: } a = \sqrt{\frac{6}{n^{[l-1]} + n^{[l]}}}$$

所以最终 Xavier 初始化均匀分布形式为 $(-\sqrt{\frac{6}{n^{[l-1]} + n^{[l]}}}, \sqrt{\frac{6}{n^{[l-1]} + n^{[l]}}})$.

若是使用正态分布，正态分布方差为:

$$Var(normal) = \sigma^2$$

$$\Leftrightarrow \sigma^2 = \frac{2}{n^{[l-1]} + n^{[l]}} \text{ 即可.}$$

所以 Xavier 的正态分布形式为 $(\mu, \sigma^2) = (0, \frac{2}{n^{[l-1]} + n^{[l]}})$.

令前向传播过程中各层状态值的方差保持不变，有:

$$\forall l, \frac{1}{2} n^{[l-1]} Var(W^{[l]}) = 1$$

前向过程对应的高斯分布为 $N(0, \sqrt{\frac{2}{n^{[l-1}}})}$ ，均匀分布为 $U(-\sqrt{\frac{6}{n^{[l-1}}}, \sqrt{\frac{6}{n^{[l-1}}})}$

在反向传播过程中，有:

$dz^{[l]} = f'(z^{[l]}) da^{[l]}$, 因为 f' 是 ReLU 激活函数，所以 $f'(z^{[l]})$ 为 0 或 1 的概率是相同的。

所以 $E(dz^{[l]}) = \frac{1}{2} E(da^{[l]}) = 0$. (由期望的定义可以推得)

所以 $E((dz^{[l]})^2) = Var(dz^{[l]}) = \frac{1}{2} Var(da^{[l]})$ (由方差的定义可以推得)

于是: $Var(da^{[l-1]}) = n^{[l]} Var(W^{[l]}) Var(dz^{[l]}) = \frac{1}{2} n^{[l]} Var(W^{[l]}) Var(da^{[l]})$

所以 $Var(da^{[l-1]}) = Var(da^{[l]}) \prod_{i=l}^L \frac{1}{2} [n^{[i]} Var(W^{[i]})]$

令前向传播过程中关于激活值的梯度的方差保持不变，有:

$$\forall l, \frac{1}{2} n^{[l]} Var(W^{[l]}) = 1$$

反向过程对应的高斯分布为 $N(0, \sqrt{\frac{2}{n^{[l]}}})$ ，均匀分布为 $U(-\sqrt{\frac{6}{n^{[l]}}}, \sqrt{\frac{6}{n^{[l]}}})$.

原论文中指出不管是按照前向过程还是反向传播过程得到的初始化方法对结果影响不大。两个过程得到的结果只是一个是使用输入层单元个数，一个是使用输出层单元个数。

12. Numerical approximation of gradients_two side difference

for checking if gradient descent is working properly. Use two side difference.

Gradient descent estimation:

>1. Two side estimation: error $\sim \varepsilon^2$ (while ε close to 0)
while 1 side estimation: error $\sim \varepsilon$

13. Gradient check:

to debug if implementing gradient descent correctly.

Process:

>1. Take all parameters $w[1], b[1], \dots, w[L], b[L]$ and reshape into a big vector θ .
--> $J(w, b) = J(\theta) = J(\theta_1, \theta_2, \theta_3, \dots)$

>2. Take $d w[1], d b[1], \dots, d w[L], d b[L]$ and reshape into a big vector $d\theta$.

>3. for every element in vector θ

>>> 3.1 Compute estimated $d\theta_i$ by two side difference:
 $d\theta_i \text{ app} = \{J(\theta_1, \theta_2, \theta_3, \dots, \theta_i + \varepsilon, \dots) - J(\theta_1, \theta_2, \theta_3, \dots, \theta_i - \varepsilon, \dots)\} / 2\varepsilon$

>>> 3.2 Compute derivative $d\theta_i$ by formula:
 $d\theta_i = dJ/d\theta_i$

--> after go through all element in vector θ , get two vectors: $d\theta_{\text{app}}, d\theta_i$

>4. check similarity of two vectors: $d\theta_{\text{app}}, d\theta_i$ by vector distance:

$$\|d\theta_{\text{app}} - d\theta_i\|_2 / (\|d\theta_{\text{app}}\|_2 + \|d\theta_i\|_2)$$

Euclidean distance: $\|d\theta_{\text{app}} - d\theta_i\|_2$

sum of element in distance vector and then take square root.

Euclidean length: $\|d\theta\|_2$. used here just incase any of these two vectors is too small or large, turn result into a ration

could set: $\varepsilon = 10^{-7}$.

result: 10^{-7} , may ok;

result: 10^{-5} : concerned, double check components of vector, make sure none of the components are too large.

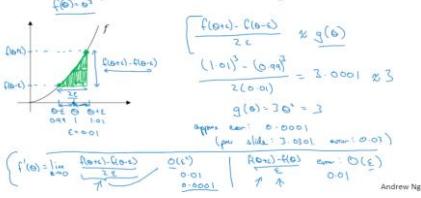
result 10^{-3} : worry there is a bug.--->

look at each component of these two vectors, find the component of $d\theta_{\text{app}}$ that is very different from $d\theta_i$, then to track down derivative computation might be incorrect.

Note:

Implement forward, backward, and do grad check, and do debug if need.

Checking your derivative computation



Gradient check for a neural network

Take $[W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}]$ and reshape into a big vector θ .
 ↓↓↓
 $J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = J(\theta)$

Take $[dW^{[1]}, db^{[1]}, \dots, dW^{[L]}, db^{[L]}]$ and reshape into a big vector $d\theta$.
 ↓↓↓

Is $d\theta$ the gradient of $J(\theta)$?

Gradient checking (Grad check)

$$\text{for each } i: \frac{J(\theta_1, \theta_2, \dots, \theta_i + \varepsilon, \dots) - J(\theta_1, \theta_2, \dots, \theta_i - \varepsilon, \dots)}{2\varepsilon}$$

$$\approx \frac{\partial \theta_i}{\partial \theta_i} = \frac{1}{2\varepsilon}$$

$$\text{Check: } \frac{\|d\theta_{\text{app}} - d\theta_i\|_2}{\|d\theta_{\text{app}}\|_2 + \|d\theta_i\|_2} \approx \frac{\|d\theta_{\text{app}} - d\theta_i\|_2}{2\varepsilon}$$

$$\varepsilon = 10^{-3} \quad \approx \frac{10^{-3} - \text{exact!}}{10^{-3}} \leftarrow$$

$$\approx 10^{-3} - \text{wrong!} \leftarrow$$

14. Grad check tips:

>1. First, don't use grad check in training, only to debug.

computing $d\theta_i$ _approx , for all the values of i, this is a very slow computation.

So to implement gradient descent, you'd use backprop to compute $d\theta$ and just use backprop to compute the derivative. And it's only when you're debugging that you would compute this to make sure it's close to $d\theta$. But once you've done that, then you would turn off the grad check, and don't run this during every iteration of gradient descent, because that's just much too slow.

>2. Second, if an algorithm fails grad check, look at the components, look at the individual components, and try to identify the bug.

Remember, different components of θ_i correspond to different components of $b[i]$ and $w[i]$.

if $d\theta_i$ _approx is very far from $d\theta$, look at the different values of i to see which are the values of $d\theta_i$ _approx that are really very different than the values of $d\theta$. e.g: if find that the values of $d\theta_i$ _approx or $d\theta_i$ they're very far off, all correspond to $db[i]$ for some layer , but the components for $dw[i]$ are quite close, then maybe you find that the bug is in how you're computing db , the derivative with respect to parameters b.

This doesn't always let you identify the bug right away, but sometimes it helps you give you some guesses about where to track down the bug.

>3. Remember regularization:

when doing grad check, remember your regularization term if you're using regularization.

(note: debugging is used in training ; while regularization also only used in training to penalize weights.)

$J(\theta) = \text{original cost} + \text{regularization (w)}$

$d\theta$: gradient of J over θ -consider regularization term in derivative computation

grad check purpose is to check cost J derivative computation during gradient descent-minimizing cost and updating parameter.-->regularization need to be involved in grad check if cost J have regularization, to check if optimizing parameter derivative computation ok or not.

>4. grad check doesn't work with dropout

as in every iteration, dropout is randomly eliminating different subsets of the hidden units. There isn't an easy to compute cost function J that dropout is doing gradient descent on.

It turns out that dropout can be viewed as optimizing some cost function J, but it's cost function J defined by summing over all exponentially large subsets of nodes they could eliminate in any iteration. So the cost function J is very difficult to compute, and you're just sampling the cost function every time you eliminate different random subsets in those we use dropout. So it's difficult to use grad check to double check your computation with dropouts.

So implement grad check without dropout, can set keep_prob and dropout to be equal to 1.0. And then turn on dropout and hope that my implementation of dropout was correct.

There are some other things could do: like fix the pattern of nodes dropped and verify that grad check for that pattern of algorithm is correct. but in practice do not usually do that. So my recommendation is turn off dropout, use grad check to double check that your algorithm is at least correct without dropout, and then turn on dropout.

>5. Run at beginning (at random initialization); perhaps run grad check again after some training /iterations.

This is a subtlety. It is not impossible, rarely happens, but it's not impossible that your implementation of gradient descent is correct when w and b are close to 0, so at random initialization. But it gets more inaccurate when w and b become large.?

So one thing you could do, I don't do this very often, but one thing you could do is:

>>>1. run grad check at random initialization
>>>2. train the network for a while so that w and b have some time to wander away from 0, from your small random initial values.

>>>3. run grad check again after you've trained for some number of iterations.

Gradient checking implementation notes

网易云课堂

- Don't use in training – only to debug
- If algorithm fails grad check, look at components to try to identify bug.
- Remember regularization.
- Doesn't work with dropout.

Run at random initialization; perhaps again after some training.



Andrew Ng

This week summary:

Learned about how to set up your train, dev, and test sets, how to analyze bias and variance and what things to do if you have high bias versus high variance versus maybe high bias and high variance.

You also saw how to apply different forms of regularization, like L2 regularization and dropout on your neural network.

So some tricks for speeding up the training of your neural network. And then finally, gradient checking.

1. Mini-batch gradient descent

applying machine learning is a highly empirical process, is a highly iterative process. In which you just had to train a lot of models to find one that works really well.

So, it really helps to really train models quickly. One thing that makes it more difficult is that Deep Learning tends to work best in the regime of big data. We are able to train neural networks on a huge data set and training on a large data set is just slow. So, what you find is that having fast optimization algorithms, having good optimization algorithms can really speed up the efficiency of you and your team.

Vectorization allows to efficiently compute on m examples. But when m is too large, it's too slow.
Only do one gradient step after computing all examples.

Mini-batch one epoch:
one epoch: a single pass through training set:

1.0: Minibatch process:

Run one epoch:

>1. for each mini-batch examples, run forward and backward, do one iteration, updating parameter;

>2. go through all mini-batches, -->algorithm do multiple times(= training data size m / mini-batch size) iterations, updated parameter multi times;

Run multi epochs:

>3. repeat a certain number epoch.

1.1: Cost function:

J_{train} :
for batch gradient descent, J_{train} decrease signally;

While for mini-batch, J_{train} decrease with noise:
Noise due to : $J_{\text{train}} = \text{cost function with parameters updated in last iteration using last batch examples, computed loss on new batch examples}$ - due to different batch examples used, J_{train} do not guarantee decrease very iteration.

there are harder mini-batch (J_{train} big), and easy mini-batch (J_{train} small)

(note: J_{train} decrease signally on conditions: learning rate is small enough.

input feature are normalized-->algorithm/cost function treat each input feature equally-->w similar, Δw_i similar also

-->gradient step : learning rate* Δw_i (same learning rate for all feature weight derivatives) also similar and easy find one learn rate fit all feature weight derivatives/step-->cost function decrease.

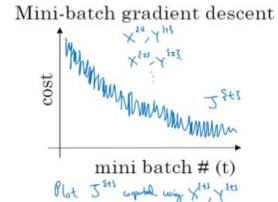
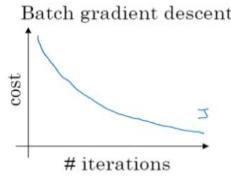
in theory, if have separate reasonable learn rate for each (input) feature weight, could make sure cost function decrease. -->using one learn rate for all features, in case overshooting for large feature derivatives-normalize all features to make their weight derivatives similar.

mini-batch: cost function noise:

all training set feature normalize first-->feature weight derivative: average on mini-batch-->gradient step, update feature weight, update cost: cost function on this mini-batch examples decrease.

yet with these weight, cost on another unfitted mini-batch: cost may bigger than before updating on previous mini-batch; weight derivative, average on previous mini-batch example, may overshooting on cost of next mini-batch.)

Training with mini batch gradient descent



1.2 Mini-batch size:

>1. size = m, batch gradient descent:

gradient descent: take relatively low noise, large steps, keep marching to global optima-blue curve.

>2. size = 1, stochastic gradient descent: every example is its own mini-batch.

Take gradient descent on each one example.

gradient descent: every iteration, take gradient descent on one single example, so most of the time hit towards the global minimum, but sometimes hit in the wrong direction.
(note: gradient descent direction based on one example, could not fit all the rest examples)

-->stochastic gradient descent can be extremely noisy, on average take good direction, but sometimes head in the wrong direction as well.

-->stochastic gradient descent won't ever converge, always just kind of oscillate and wander around the region of the minimum, but won't ever head to the minimum and stay there.

Choosing your mini-batch size

If small tiny set: use batch gradient descent.
(m < 200)

Typical minibatch sizes:

$\rightarrow 64, 128, 256, 512$

$\frac{1024}{2^6}$

Make sure minibatch fits in CPU/GPU memory.



>3. size = 1-m: mini-batch:

if batch size = m, processing on a huge training set on every iteration, takes too long per iteration.

if m is small, batch gradient descent works fine.

if size=1, stochastic gradient descent, make process just after processing 1 example, and noisiness can be reduced by just using a smaller learning rate (but noise will not be eliminated as long as batch size < m?). but huge disadvantage of stochastic gradient descent: loose almost all your speed up from vectorization-processing a single training example at a time, is very inefficient.

mini-batch is between, size not too big/small -->gives the fastest learning:

do get a lot of vectorization, faster than processing the examples one at a time;
could also make process without needing to wait till process enter training set.

mini-batch gradient descent:

>>>1. not guaranteed to always head toward the minimum, but tends to head more consistently in direction of the minimum than stochastic gradient descent.

>>>2. Does not always exactly converge or oscillate in a very small region - can always reduce the learning rate slowly.

Summarize:

>1. if small training set, just use batch gradient descent; m < 2000

>2. big training set, take mini-batch size: 64, 128, ... 512 (2^n)

consider the way computer memory is laid out and accessed, sometimes could run faster if mini-batch size is a power of 2

>3. make sure all mini-batch fit in CPU/GPU memory (?), otherwise algorithm performance suddenly falls off a cliff, much worse.

>4. mini-batch is actually hyperparameter, need to figure out which one is most sufficient of reducing the cost function J: try a different values and then pick one makes gradient descent optimization algorithm as efficient as possible.

(note: mini-batch only better than stochastic when well vectorized, how to choose min-batch size in coding:
try different size and see running time, performance vs iteration on dev set?)

1. Mini-batch gradient descent

>1. Batch gradient descent: Use all m examples in each iteration;

>2. Stochastic gradient descent: Use a single example in each iteration;

>3. Mini-batch gradient descent: Use b examples in each iteration

Parameter: b - "mini batch size", type choice: 2-100 (common 10)

1.1 Mini-batch idea: Somewhat in-between Batch gradient descent and Stochastic gradient descent. Rather than using one example at a time or m examples at a time we will use b examples at a time.

Mini-batch gradient descent

Say $b = 10, m = 1000$,

Repeat {^K}

for $i = 1, 11, 21, 31, \dots, 991$ {

$\theta_j := \theta_j - \frac{1}{10} \sum_{k=1}^{10} (h_\theta(x^{(k)}) - y^{(k)}) x_j^{(k)}$

(for every $j = 0, \dots, n$)

}

}

$m = 300, 000, 000$

$b = 10$

→ 10 examples

→ 1 example

Vectorization

1.2 Mini-batch algorithm

e.g: $b=10$, $m=1000$.

>>1. inner loop: one epoch (go through all training example one time)
 >>> 1. one iteration: gradient descent on $b=10$ examples
 >>> 2. go continually $100 (m/b)$ iterations, each iterations with next, different 10 examples
 >>-> go through all examples 1 time, $100 (m/b)$ iterations with b example for each iteration.
 each iteration use different b examples.

>>2. outer loop: repeat inner loop: --> go through again all the examples.

Mini-batch faster than batch gradient descent:

1.3 Comparison: Mini-batch vs Batch, stochastic,

>>1. Mini-batch vs Batch

Mini-batch faster:

(Each iteration) start making progress in modifying the parameters after looking at just b examples rather than needing to wait till have scan through every single training example.

>>2. Mini-batch vs Stochastic gradient descent

Mini-batch outperform stochastic : only if have a good vectorized implementation

In each iteration, sum over b examples (for the derivation of J over parameter θ) in a more vectorized way will allow to partially parallelize computation over the b examples.
 using appropriate vectorization to compute the derivatives, can sometimes partially use the good numerical algebra libraries and parallelize gradient computations over the b examples.
 Whereas with Stochastic gradient descent, just looking at one example at a time is less efficient to parallelize over.

>>3. Disadvantage of mini-batch:

Have one more parameter to fiddle with: mini-batch size b , which may therefore take time.

But if have a good vectorized implementation this can sometimes run even faster than Stochastic gradient descent.

Summary:

Mini-batch gradient descent algorithm in some sense does something that's somewhat in between what Stochastic gradient descent and batch gradient descent does.

And if choose a reasonable value of b (usually use 10... anywhere from say 2 to 100 would be reasonably common), and use a good vectorized implementation, sometimes it can be faster than both Stochastic gradient descent and faster than Batch gradient descent.

2. Exponentially weighted average:

Algorithm faster than mini-batch:

2.1 e.g: temp. trend:

>>1. plotting temp of everyday, vs days, --> a little big noisy
 >>2. if want to compute the trends the local average / moving average / exponentially weighted average of temp,
 $v_0 = 0$
 $v_1 = 0.9 * v_0 + 0.1 * t_1$
 $v_2 = 0.9 * v_1 + 0.1 * t_2$
 ...
 --> $v_t = \beta * v_{t-1} + (1-\beta) * T_t$: an approximately average over $1/(1-\beta)$ days' temp.
 plotting is more smoother as now is averaging over more days of temp.

$\beta=0.5$, about 2days average temp.- yellow line

$\beta=0.9$, about 10 days average temp.- red line

$\beta=0.98$, about 50 days average temp.- green line

--> β is bigger (weighting for previous data)

>>>1. plotting curve shift to further right, as now averaging over a much larger window of temp.
 >>>2. by averaging over large window, this exponentially weighted average formula, adapts more slowly when the temp. changes.

giving lot of weight to the previous days' temp, much smaller weight to the value of right now.

find β works best: balancing sensitive to right now value and cover previous data.

Mini-batch gradient descent

→ Batch gradient descent: Use all m examples in each iteration

→ Stochastic gradient descent: Use 1 example in each iteration

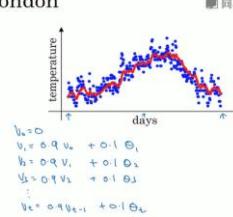
Mini-batch gradient descent: Use b examples in each iteration

$b = \text{mini-batch size}$.
 $b = \frac{m}{10}$ examples
 $\theta_j := \theta_j - \alpha \frac{1}{10} \sum_{i=1}^{10} (h_\theta(x^{(i)}) - y^{(i)}) \times_j^{(i)}$

$i := i + 1$

Temperature in London

$$\begin{aligned}\theta_1 &= 40^\circ\text{F } 45^\circ\text{C} \\ \theta_2 &= 49^\circ\text{F } 9^\circ\text{C} \\ \theta_3 &= 45^\circ\text{F } 13^\circ\text{C} \\ &\vdots \\ \theta_{100} &= 60^\circ\text{F } 15^\circ\text{C} \\ \theta_{101} &= 56^\circ\text{F } 11^\circ\text{C} \\ &\vdots\end{aligned}$$



Exponentially weighted averages

$$v_t = \beta v_{t-1} + (1-\beta) t_t$$

$\beta = 0.9$: ≈ 10 days' weight
 $\beta = 0.99$: ≈ 50 days' weight
 $\beta = 0.999$: ≈ 2 days' weight

$v_t \approx$ approximately
 older now
 $\rightarrow \frac{1}{1-\beta}$ days' temperature.

$$\frac{1}{1-\beta} = \infty$$

4.2: how many previous data to average: $1/(1-\beta)$

$$1/e = 0.35 = \beta^{1/(1-\beta)} = (1-\beta)^{1/(1-\beta)}$$

$$\beta=0.9, \rightarrow 1/(1-\beta)=0.9^{10}$$

-->take about 10 days, to decay from now day value weight $(1-\beta)$ to around $1/3$.

Note: decay factor of previous i day: β^i !

--> only take last i days , whose decay factor is β^i , take i make $\beta^i = 1/e \approx 1/3$, --> get i $\approx 1/(1-\beta)$

when $i > 1/(1-\beta)$, decay $< 1/e$, 30% , then do not consider that days data contribute to sum.

as if computing an exponentially weighted average that focus on just the last 10 $(1/(1-\beta))$ days. as for more previous days, their weight decays to more than $1/3$ of the weight of the current day.

How to set β , balancing weight put on current day data.

Note: sum of each previous i-th day's weight: $(1-\beta) * \beta^i$
 $\text{sum} = (1-\beta) * (1+\beta+\beta^2+\dots+\beta^{i-1}) = (1-\beta) * (1-\beta^i)/(1-\beta) = 1-\beta^{i-1} \rightarrow \approx 1$, when i is large enough.

for t days, $V(t) = (1-\beta) * T_t + (1-\beta) * \beta^1 * T_{t-1} + \dots + (1-\beta) * \beta^{t-1} * T_{t-1} + \dots + (1-\beta) * \beta^{t-1} * T_{t-1} + \beta^t * T_t$
 V(0)
 all days T_{t-1} weight sum = $(1-\beta) * (1+\beta+\beta^2+\dots+\beta^{t-1}) = 1-\beta^t$
 V_0 weight = β^t
 total weight sum = 1

4. Understanding exponential weighted average

$$\begin{aligned}V_t &= \beta V_{t-1} + (1-\beta) t_t \\ &= (1-\beta) * V_{t-1} + (1-\beta) * \beta^1 * t_t + \dots + (1-\beta) * \beta^{t-1} * t_t + \dots + (1-\beta) * \beta^{t-1} * t_t + \beta^t * t_t\end{aligned}$$

4.1: Plotting:

>>1. T_t vs time/day

vertical axis: temp of each day: T_0, \dots, T_t ,
 horizontal axis: days: $1, \dots, t, \dots$

>>2. Exponentially decaying function: $((1-\beta) * \beta^t, (1-\beta) * \beta^{t-1}, \dots, (1-\beta) * \beta^1)$

vertical axis: range: $[(1-\beta) * \beta^t, (1-\beta) * \beta^1]$.

-weight for each day data, sum together to get current day t average data.

Note: decay factor of previous i day: β^i)

horizontal axis: days: $1, \dots, t, \dots$

>>3. take the element wise product between above two functions and sum it up:
 take daily temp and multiply with corresponding exponentially decaying function/weight, and then summing it up.

Exponentially weighted averages

$$v_t = \beta v_{t-1} + (1-\beta) t_t$$

$$\begin{aligned}V_{100} &= 0.9 v_{99} + 0.1 \theta_{100} \\ v_{99} &= 0.9 v_{98} + 0.1 \theta_{99} \\ v_{98} &= 0.9 v_{97} + 0.1 \theta_{98} \\ &\vdots \\ v_{10} &= 0.9 v_9 + 0.1 \theta_{10} \\ v_9 &= 0.9 v_8 + 0.1 \theta_9 \\ &\vdots \\ v_0 &= 0.9 v_0 + 0.1 \theta_0\end{aligned}$$

decay?

$$\frac{1}{1-\beta} = \frac{1}{0.9} = \frac{1}{0.1} = 10$$

$\beta = 0.9$
 $\beta = 0.99$
 $\beta = 0.999$

Andrew Ng

Implementing exponentially weighted averages

$$\begin{aligned}v_0 &= 0 \\ v_1 &= \beta v_0 + (1-\beta) \theta_1 \\ v_2 &= \beta v_1 + (1-\beta) \theta_2 \\ v_3 &= \beta v_2 + (1-\beta) \theta_3 \\ &\vdots \\ &\rightarrow v_0 = 0 \\ K_{\beta} &= \frac{\beta}{1-\beta} \\ v_{\beta} &:= \beta v + (1-\beta) \theta_{\beta} \\ v_{\beta} &:= \beta v + (1-\beta) \theta_{\beta} \\ &\vdots \\ &\rightarrow v_0 = 0 \\ K_{\beta} &= \frac{\beta}{1-\beta} \\ v_{\beta} &:= \beta v + (1-\beta) \theta_{\beta} \leftarrow \frac{\beta}{1-\beta} = \frac{1}{0.9} = 10\end{aligned}$$

4.3: Implementing exponentially weighted averages:

advantage: takes very little memory, only one row number in compute memory:
 $V_{\theta} = 0$
repeat:
get next θ_t
 $V_{\theta} := \beta * V_{\theta} + (1-\beta) * \theta_t$

Summarize:

Exponentially weighted average is not the most accurate way to compute an average, while most accurate way is to store all past 10/50 days data and average them-huge storage and computation cost. while using exponentially weighted average, only need to store one value V_{θ} , storage and computation efficiently, and only need one row coding.

5. Bias correction in exponentially weighted average:

$$V_t = \beta * V_{t-1} + (1-\beta) * T_t \\ = (1-\beta) * V_{t-1} + (1-\beta) * \beta^1 * T_{t-1} + \dots + (1-\beta) * \beta^{i-1} * T_{t-i} + \dots + (1-\beta) * \beta^{t-1} * T_{t-1} + \beta^t * V_0$$

Bias in the initial phase : caused by $V_0=0$ --> (should be green line, while is purple line with bias in initial phase)

$$V_1 = (1-\beta) * T_1 \\ V_2 = \beta * V_1 + (1-\beta) * T_2 = \beta(1-\beta) * T_1 + (1-\beta) * T_2 \\ \text{no term of } \beta^i * T_i \text{, which influence } V_i \text{ when } i \text{ is small } (\beta^i \text{ is big}).$$

Bias correction:

$$V_t := 1 / (1-\beta^t)$$

>1. adding bias correction: initial phase change from purple line to green line.

>2. as t becomes bigger, $\beta^t * V_0$ nearly zero, no bias -purple line no bias in late phase
>3. as t becomes bigger, bias correction factor also close to 1, no bias correction function-green line aligned with purple line.

Note: bias correction use $1 / (1-\beta^t)$ -->sum of each day's weight =1

all days T_{-1} weight sum= $(1-\beta)(1+\beta+\beta^2+\dots+\beta^{t-1})=1-\beta^t$, using correction -->sum =1

$$\text{for } t \text{ days, } V(t) = (1-\beta) * T_{-t} + (1-\beta) * \beta^1 * T_{-(t-1)} + \dots + (1-\beta) * \beta^{i-1} * T_{-(t-i)} + \dots + (1-\beta) * \beta^{t-1} * T_{-1} + \beta^t * V_0$$

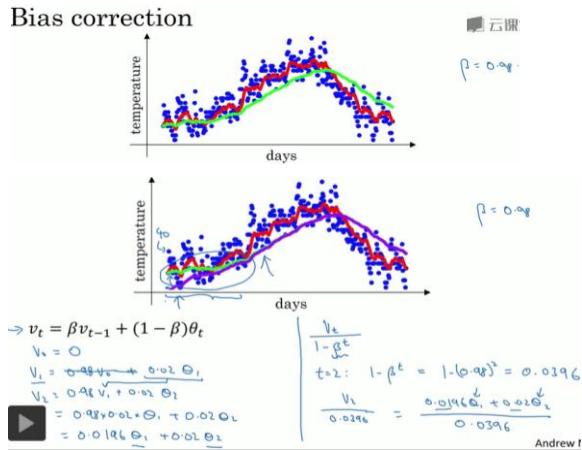
$$\text{all days } T_{-1} \text{ weight sum}= (1-\beta)(1+\beta+\beta^2+\dots+\beta^{t-1})=1-\beta^t$$

$$V_0 \text{ weight}= \beta^t$$

$$\text{total weight sum} = 1$$

Summarize:

In machine learning, for most implementations of the exponentially weighted average, people don't often bother to implement bias corrections , because most people would rather just wait that initial period and have a slightly more biased assessment , and then go from there. But if concerned about the bias during this initial phase, while your exponentially weighted moving average is still warming up, then bias correction can help you get a better estimate early on.



6. Gradient descent with momentum

compute an exponentially weighted average of gradients.

6.1 Gradient descent problem:

During gradient descent, may have up and down oscillation,-->

>1. slows down gradient descent, and prevents from using a much larger learning rate--which may end up overshooting and diverging:

>2. Oscillation could reduce by small learning rate

>3. another view: want learning a bit slower in oscillation - e.g vertical direction (less oscillation), but faster learning in no oscillation -e.g horizontal direction: move fast from left to minimum.-right ex.g

6.2: Momentum in gradient descent:

>1. Principle: using a moving average of dw, db

>2. Function-Intuition: smooth out the step of gradient descent : oscillation reduced almost to 0 as averaged over positive and negatives (in the vertical direction)-average will close to 0; while in the horizontal direction, all the derivatives are pointing to the right direction, so average in this direction will still be big

(note: oscillation: due to weight derivative change direction <--learning rate/gradient step too big for this weight/direction;

-->by using moving average of weight derivative: will even out the direction of weight derivative: the derivative value pe se may no much change, but direction: +/- more stable: do not change so frequently-->show as oscillation)

-->gradient oscillate much less and moving quickly in horizontal direction.

(optimizing in horizontal direction-this parameter updating no change, optimizing in oscillato directioin-this oscillating parameter derivative-faster)

6.3. Formula:

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

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

$$W := -a * V_{dw}$$

$$b := -a * V_{db}$$

oscillation root cause: gradient step 'a*dw' too big

1. value $|dw|$ is big-->updated paramter w' and previous parameter w will be on two sides of optima parameter w_0 --> $dw < 0$: negative and positive direction; -->derivative direction changed momentum: even out derivative direction dw / db (no change on learning rate, which used for all parameters);

RTMS: reduce dw value :

(note 1: momentum: using exponentially weighted average of gradients over iterations.-parameter derivative direction even out: positive and negative (-dw:是向量, 有方向,momentum 主要改变方向, 随带会改变其大小)

for the oscillation direction-specific parameter derivation: gradient step $|a*dw|$ too big, update parameter's derivative direction positive and negative changed from -dw to -dw'; Exponential weighted derivatives even direction-positive and negative close to 0): positive and negative direction in past and current iteration will sum together, even out direction, very close to 0--also make $|dw|$ close to 0: to keep dw always one direction positive or negative while $|dw|$ value is small. (now learning on this derivative-parameter will be quite slow, how to deal with this?-->not slow,even out derivative direction and value, -->with step dw/a , march towards to optima; now the gradient step dw/a: is a better step, previous too big to cause overshooting, even dw while no change learning rate, to actually make a better'gradient step' now for this oscillation weight/direction)

for no oscillation direction: parameter derivation average nearly no change.)

(note2: exponentially weighted average of gradients: not use on mini-batch examples/batch examples within one iteration: as need to consider all examples influence on parameter derivatives of one step, not few-1/ β batch /iterations examples) --> could use on mini-batch

(J_train decrease with noise due to different mini batch examples.-different big enough not able to use exponentially weighted average of gradient in different batch??)

note 3: oscillation in some direction-specific parameter derivative: during iterations:

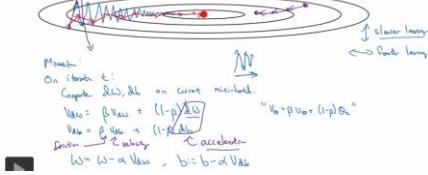
all oscillations actually could be solved by reducing learning rate to small enough--perfect solution is separate, unique learning rate for each parameter, in every iteration.

solution: separate learning rate for each parameter-->could be replaced by using same learning rate with normalized input feature.

solution: unique learning rate for each iteration-->could be replaced by learning rate decay during iterations.

-->oscillation should be improved by all layer input featuring scaling & learning rate decay during iterations ?-->then no need of momentum)

Gradient descent example



Implementation details

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

On iteration t:

Compute dW, dB on the current mini-batch

$$\rightarrow v_{dw} = \beta * v_{dw} + (1-\beta) * dW$$

$$\rightarrow v_{db} = \beta * v_{db} + (1-\beta) * db$$

$$W = W - \alpha v_{dw}, b = b - \alpha v_{db}$$

perparameters: α , β
 $\beta = 0.9$
alpha or low ≈ 10 gradient

6.4: Implementation details:

Initialize
 $V_{dw} = 0; V_{db} = 0$

repeat iterations: for one iteration:

>1. on iteration t, compute parameter derivatives dw, db on mini-batch examples;

>2. update parameter derivatives:
 $V_{dw} = \beta * V_{dw} + (1-\beta) * dw$
 $V_{db} = \beta * V_{db} + (1-\beta) * db$

>3. Update parameter
 $W := -a * V_{dw}$
 $b := -a * V_{db}$

Hyperparameter: a, β

β : common take 0.9 - pretty robust value, could try different values.

Many updated formula in paper:

$V_{dw} = \beta * V_{dw} + * dw$
--> V_{dw} scaled by $1/(1-\beta)$,
while updating parameter, learning rate needs to change by a corresponding value $1/(1-\beta)$

Note: recommend official formula $V_{dw} = \beta * V_{dw} + (1-\beta) * dw$.
otherwise need to tune not only β , learning rate have to tune correspondingly.

6.3 Momentum intuition:

could take gradient descent as: ball rolling down from bowl-cost function.
 dw, db : ball acceleration
 V_{dw}, V_{db} : velocity-plays a role of fraction, prevents ball from speeding up without limit.

So rather than previous steps taking every single step independently, now the ball can roll downhill and gain momentum by accelerating down.

7. RMSProp: root mean square prop

>1. Idea: exponentially weighted average of derivative squares. - damping out oscillation by slowing down parameter update, parameters that caused these oscillation

>2. e.g: right pid: gradient update in vertical direction and moving in horizontal direction.
Take parameter b as vertical direction, w as horizontal;
on one iteration t: `compute_dw_db_on_current_mini-batch`.

$S_{dw} = \beta * S_{dw} + (1-\beta) * dw^2$ (element-wise squaring operation)
 $S_{db} = \beta * S_{db} + (1-\beta) * db^2$

$W := -a * dw / \sqrt{S_{dw}}$
 $b := -a * db / \sqrt{S_{db}}$

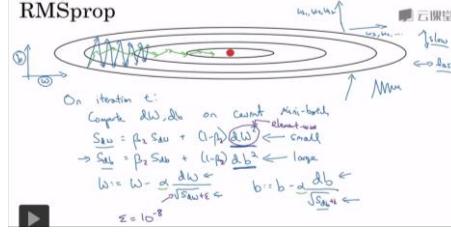
-->if oscillate in parameter w/b direction is big, --> S_{dw}/S_{db} is large, slow down parameter w/b update;
reduce oscillation.
(note: oscillation: due to gradient step a* dw too big, -->
>>1. oscillate: derive direction change
>>2. oscillation magnitude: big if value|a* dw| big, otherwise small oscillation
RMSP: reduce big oscillation only: if was small oscillation, may make it a little bigger... but no care small oscillation.)

-->if original gradient descent has no oscillate in parameter w/b, dw/db is small--> S_{dw}/S_{db} is smaller,
faster move in this direction than direction have oscillation.
(faster than direction have oscillation: with RMSP, gradient step a* dw/db is bigger in direction |dw/db|
small-->faster move in non-oscillation direction than oscillation direction
possible be slower than without RMSP in non-oscillation direction-if |dw|>1)

Note:
>1. here take b, w as vertical and horizontal direction, in practice w, b is high dimensional, so the oscillation may happen on parameter w_i, b_i.
but in dimension getting these oscillations, will end up computing a larger sum: $S_{dw,i}/S_{db,i}$, so end up damping out in which there are these oscillations.

>2. for ensuring numerical stability, in case s_dw/s_db too small, -->divide sqrt(s_dw + ε)
could set ε = 10^{-8}

RMSprop



(note: root cause for oscillation: parameter updating with gradient descent : $w = w - a * dw$ gradient step : $a * dw$

1. oscillation happens when: $w' = w - a * dw$ and w are on the two side of optima parameter w_0 -over shooting <---gradient step $a * dw$ is too big.
2. there are parameter gradient descent step ok: $a * dw_i$: ok, no overshooting for parameter w_i , while learning rate a is same for all parameters->
for oscillation direction: parameter derivative dw_i is big;
for no oscillation direction: parameter derivative dw_i is small;

for solving oscillation, reduce value $a * dw$ while keep a no change.

>1. oscillation direction-parameter: fracter $1/dw_i$: <---gradient step reduce; help improve oscillation
-->2. no oscillation direction-parameter: fracter $1/dw_i$: may >1-->gradient step increase; ->while turn to be oscillation?? this will be even out in next iteration.

with RMSProp, -->
>1.gradient descent end up oscillate less in parameter that cause oscillation (e.g vertical) and keep going in no oscillation (e.g horizontal) direction.
>2. could use larger learning rate, fast learning.
(If faster learning on no-oscillation direction?
as increase learning rate will make original oscillation direction more worse
-->dw kind of 'normalized' $|dw| \approx 1$
need to find learning rate a, make gradient decent step $a * dw$ fit well for loss function with all parameters direction)

两者的区别如下:

1. 前者利用了历史梯度的指数移动加权平均, 没有显式的考虑梯度的范数大小, 后者是显式地对历史梯度各个分量的二范数进行指数移动加权平均, 显式地在梯度各个分量的数值大小上做平衡, 在梯度值较大的分量上减小更新步伐, 在梯度值较小的分量上增大更新步伐, 没有显式的考虑梯度方向。
2. 从RMSprop公式上看, 它可以实现学习率自适应调整, 变化较大的梯度分量上的学习率会自动减小, 变化较小的梯度分量上的学习率会自动增大。

8. Adam optimization algorithm

Adam: Adaptive Moment Estimation

RMSprop and Momentum, work well across a wide range of deep learning architectures.

8.1 Principle: take Momentum and RMSprop together.

8.2 Formula:

>1. initialization: $V_{dw} = 0, V_{db} = 0, S_{dw} = 0, S_{db} = 0$

>2. for each iteration on mini-batch

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

$S_{dw} = \beta_2 * S_{dw} + (1-\beta_2) * dw^2$ (element-wise squaring operation)
 $S_{db} = \beta_2 * S_{db} + (1-\beta_2) * db^2$

Adam optimization algorithm



$V_{dw} = 0, S_{dw} = 0, V_{db} = 0, S_{db} = 0$

On iterate t:
Compute $\delta W, \delta b$ using current mini-batch.
 $V_{dw} = \beta_1 V_{dw} + (1-\beta_1) \delta W$, $V_{db} = \beta_1 V_{db} + (1-\beta_1) \delta b$ ← "moment" β_1
 $S_{dw} = \beta_2 S_{dw} + (1-\beta_2) \delta W^2$, $S_{db} = \beta_2 S_{db} + (1-\beta_2) \delta b^2$ ← "RMSprop" β_2
 $V_{dw,corr} = V_{dw} / (1-\beta_1^{t+1})$, $V_{db,corr} = V_{db} / (1-\beta_1^{t+1})$
 $S_{dw,corr} = S_{dw} / (1-\beta_2^{t+1})$, $S_{db,corr} = S_{db} / (1-\beta_2^{t+1})$
 $W := W - \alpha \frac{V_{dw,corr}}{\sqrt{S_{dw,corr}} + \epsilon}$
 $b := b - \alpha \frac{V_{db,corr}}{\sqrt{S_{db,corr}} + \epsilon}$

>3. Bias correction: on both momentum and RMSprop

$V_{dw_correct} = V_{dw} / (1-\beta_1^{t+1})$; $V_{db_correct} = V_{db} / (1-\beta_1^{t+1})$

$S_{dw_correct} = S_{dw} / (1-\beta_2^{t+1})$; $S_{db_correct} = S_{db} / (1-\beta_2^{t+1})$

>4. Parameter update

$W := -a * V_{dw} / \sqrt{S_{dw_correct} + \epsilon}$

$b := -a * V_{db} / \sqrt{S_{db_correct} + \epsilon}$

8.3. Hyperparameter in Adam:

>1. a: learning rate , need to be tune.

>2. β_1 , 0.9 : moving averaging; default
(note: consider previous 10 iterations dw direction)

>2. β_1 , 0.999 : RMSprop; default

(note: consider previous 1000 iterations $|dw|^{1/2} \rightarrow$ gradient step $a * dw$ much more stable-cons value, less sensitive to current dw status, why? gradient step should be smaller when reaching optima)

>4. ϵ : 10^{-8} ; default. Does not matter much.

Hyperparameters choice:

$\rightarrow \alpha$: needs to be tune
 $\rightarrow \beta_1$: 0.9 $\rightarrow (\underline{dw})$
 $\rightarrow \beta_2$: 0.999 $\rightarrow (\underline{dw^2})$
 $\rightarrow \epsilon$: 10^{-8}

Adam: Adaptive moment estimation

Adam: commonly used learning algorithm, proven to be very effective for many different neural networks.

>2. for each iteration on mini-batch

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

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

$S_{dw} = \beta_2 * S_{dw} + (1 - \beta_2) * dw^2$ (element-wise squaring operation)

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

Note:
momentum: reduce gradient descent oscillation by exponentially weighted average over derivative, of the past 1/(1-β1) iterations. --> reduce dw/db oscillation, by averaging the positive and negative of dw_i/db_i that caused these oscillations (may change derivative direction, and slight 'scalling' effect by multiplying weight).

RMSprop: reduce gradient descent oscillation by exponential weighted average over square derivative, of the past 1/(1-β2) iterations.

--> reduce dw/db oscillation, by scaling the positive and negative of dw_i/db_i that caused these oscillations (but do not change derivative direction-positive or negative direction), using factor 1/sqrt (s_dw_i / db_i); make derivation caused oscillation smaller by scaling based on averaged past iteration derivation length.

Adam: for oscillations caused by certain parameters: average these positive and negative parameter derivatives and then scaling by factor 1/sqrt (s_dw_i / db_i)

(note: for big oscillations:

--> RMSP: --> by scaling make it oscillate less

--> Momentum: --> by averaging derivative direction, make this small oscillations to nearly no oscillation

--> RMSP: by scaling: make no-oscillation direction, by small derivative bigger--> gradient step bigger in non-oscillation direction

--> finally make all weights: gradient decent march toward non-oscillate and stop sooner.)

e.g.: $\beta_1=\beta_2=0.5$, dw has oscillation, weight for i1, i2, i3 value: $(1-\beta_1)\beta_1^{1/2}, (1-\beta_1)\beta_1, (1-\beta_1)$

--> $V_{dw,i3} = 7/8$: derivation change from negative to positive, direction change

--> $S_{dw,i3} \sim 25$, RMSprop only--> $dw_{i3} = dw_{i3}/5 = -0.8$: derivation direction no change, yet size becomes much less

--> $V_{dw} + s_{dw}$ together --> $v_{dw,i3} / \sqrt{s_{dw,i3}} = 7/8 / 5 \sim 0.2$: derivative direction change, and size scaling to less.

db (iterative 1)=3, db_i2=5, db_i3=4; db no oscillation

--> $V_{db,i3} = 19/8$: direction no change

--> $S_{db,i3} = 67/8$, RMSprop only--> $dw_{i3} = dw_{i3}/\sqrt{s_{db,i3}} \sim 4/3$: derivation direction no change, yet size becomes much less

--> $V_{dw} + s_{dw}$ together --> $v_{dw,i3} / \sqrt{s_{dw,i3}} \sim 19/24$: derivative direction change, and size scaling to less.

9. Learning rate decay

9.1 Gradient descent on mini-batch in each iteration:

gradient decrease with noise in each mini-batch, while learning rate is fixed value, algorithm end up wandering around and never converge.
 (note: this noise will help algorithm go out of local optima/-saddle point.?)

Reduce learning rate in the late phase, will make algorithm end up in a tighter region around this minimum.

Intuition: in the initial phase, could afford to take much bigger steps-big learning rate, but as learning approaches converges, then having a slower learning rate allows to take smaller steps.

9.2 Learning rate decap implementation:

1 epoch: 1 pass through training set.

Learning rate: $a = 1 / (1 + \text{decay rate} * \text{epoch number}) * a_0$

hyperameter:
 >1. decay rate, need to tune
 >2, a_0: need to tune

9.3 other decay method:

exponentially decay:
 $a = 0.95^{\text{epoch}} * a_0$;

$a = K(\text{constant}) / \sqrt{\text{epoch-number}} * a_0$; or $a = k(\text{constant}) / \sqrt{\text{mini-batch number t}} * a_0$

or discrete learning rate: some steps have some learning rate.

Manually decay:

change learning rate manually after algorithm run some time: tuning learning rate by hand, hour by hour, or day by day.
 works only if training a small number of models.

Recommend:

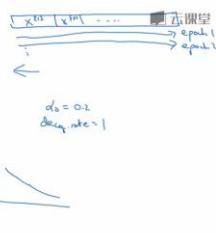
a_0 have high priority to tune, have huge impact. Learning rate decay does help, but a little bit lower down in tuning list

Learning rate decay

1 epoch = 1 pass through data.

$$a_t = \frac{1}{1 + \text{decay_rate} * \text{epoch_num}} * a_0$$

Epoch	LR
1	0.1
2	0.67
3	0.5
4	0.4
...	...



Other learning rate decay methods

hyperameter

$$\left\{ \begin{array}{l} a = 0.95^{\text{epoch}} * a_0 \quad \text{- exponentially decay} \\ a = \frac{k}{\text{epoch}} * a_0 \quad \text{or } \frac{k}{\sqrt{t}} * a_0 \\ a = \text{---} \\ a = \text{---} \end{array} \right.$$

discrete staircase

Manual Decay

10: Local optima in deep learning

10.1: Intuition : local optima

>1. Intuition: for low dimension parameter

Previous thinking for local optima: for low dimension parameter:
 cost function: have many local optima, where gradient easy to get stuck to find global optima.

>2. This intuition incorrect for deep learning:

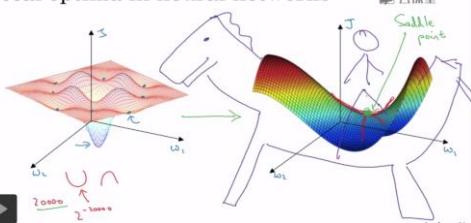
Most points of 0 gradients, are not local optima, but saddle points instead.

for a function of high dimension, when gradient = 0, then in this parameter direction it can either be convex light function, or concave light function.
 so in this high dimension space, for 0 gradient point to be a local optima, all parameter directions need to be convex. the chance of this happening is very small. more likely some curve bend up in one direction, and bend down in other direction, rather than all bend upwards.

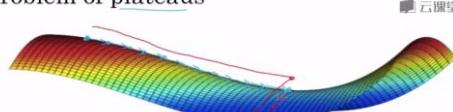
--> in very high dimensional space, actually more likely to run into a saddle point.

A lot of intuitions about low-dimensional spaces- like cost function local optima, really do not transfer to the very high-dimensional spaces

Local optima in neural networks



Problem of plateaus



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

10.2 Plateaus problem for deep learning

Local optimum is not a problem for deep learning., instead will meet plateaus problem.

Plateaus: the region where the derivative is close to zero for a long time.

E.g:

in the plateau region , surface is platten, gradient descent in all direciton close to 0,
take a very long time to slowly find way to maybe the saddle point on the plateau -blue curve.
and then because of a random perturbation of left/right, then finally find way off the plateau- red line.
(mini-batch noise on cost function helps for going out of plateau.)

Plateau problem for deep learning: take a very long stop/step to the saddle point and then get off this plateau.

Summary:

>1. Unlikely to get stuck in a bad local optima:
on conditions: as long as training a reasonably large neural network-->(high dimension space)
A lot of parameters, and cost function is defined over a relatively high dimensional space.

>2. Plateaus can make learning slow.
this is where Momentum / RMSprop / Adam can really help learning algorithm as well!!!
and this is cases where more sophisticated optimization algorithms like Adam can actually speed up the rate at which you could move down the plateau and then get off the plateau.

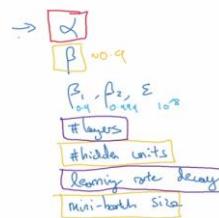
2-3_Hyperparameter tuning

1. Hyperparameters

1.1. Hyperparameters

- α : learning rate: most import hyperparameter to tune
- β , $=0.9$ good default, momentum
- mini-batch size: tune to make sure optimization algorithm running efficiently
- hidden units: w, b
- layers
- learning rate decay,
- $\beta_1=0.9$, $\beta_2=0.999$, $\epsilon=10^{-8}$: RTM, Adam: never tune in Adam
- red: most important to tune
- orange: second import to tune
- Purple: third important

Hyperparameters



1.2 Try random values:

>1: Grid:

sample points in grid (fix parameter 1, sampling parameter 2)-sample a row number, and then systematically explore these values (fix sampled parameter 2, explore parameter 1)-explore vertically. then try all these points see which one works best.

works ok when hyperparameter number was relatively small.

>2: Random choose for deep learning

in deep learning, use random choosing instead. and try these random points on algorithm.

>3 Reason for random choosing for deep learning

>>1. as could not know in advance which hyperparameter is the most important.

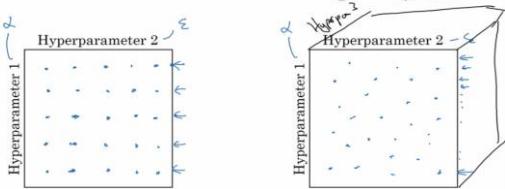
E.g., right pic.: hp1 is learning rate, hp2 is ϵ in RTM. ϵ is not so important, 5 points: one hp1 with 5 different hp2, result no big difference as rely more on hp1, same in 5 points-->25 points in grid actually could only test 5 hp1 values.

by using random points: right pic.: all points have unique hp1, hp2-->none of two points forced to have same hp2/hp1, both hp1, hp2 tested 25 values instead of 5. have more chance to find a value that works really well.

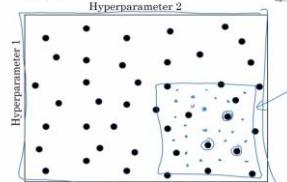
when have more parameters, sample parameter in higher dimension space-->

random choosing ensures that you are more richly exploring set of possible values for the most important hyperparameters

Try random values: Don't use a grid



Coarse to fine



>4: Coarse to fine

except above random choosing, second common practice in Sampling hyperparameters is to use a coarse to fine sampling scheme.

>>>1. In above random choosing sampling, run algorithm and find a points/ a few other points around it tended to work really well.

>>>2. Zoom in to a smaller region of these hyperparameters points, and then sample more density within this space.

maybe again randomly, but to then focus more resources on searching within this parameter region-blue square, if suspecting best setting /hyperparameters maybe in this region.

after doing a coarse sample of this entire square, then result tells to focus on a smaller square, then sample more densely into smaller square.

Coarse to fine search is frequently used.

2. Using an appropriate scale to pick hyperparameters

While searching hyperparameters, sampling at random over the range of hyperparameters, could allow to search over the hyperparameter space more efficiently.

But sampling random does not mean sampling 'uniformly' at random over the range of valid values. instead it's important to pick the appropriate scale. on which to explore the hyperparameters.

2.1 : Picking hyperparameters at random

>1. : Sampling uniformly at random over range considered

e.g:

hidden units in layers: n_l : could picking some number values at random within number line: 50-100

hidden layer number L: sampling uniformly at random along 2-4 , might be reasonable/ or could use grid search for value 2, 3, 4.

Picking hyperparameters at random

$$\rightarrow n_l = 50, \dots, 100$$

$\xrightarrow{[50, 100]}$

$$\rightarrow \# \text{layers} \quad L = 2 - 4$$

$\xrightarrow{[2, 3, 4]}$

Appropriate scale for hyperparameters

$$\begin{aligned} & d = 0.0001 \dots 1 \\ & \xrightarrow{[0.0001, 1]} \\ & \log_{10}(r) = -4 * np.random.rand() \quad \leftarrow r \in [-4, 0] \\ & \frac{10^{\log_{10}(r)}}{10^{-4}} = 10^r \quad \leftarrow 10^{-4} \dots 10^4 \\ & \frac{10^r}{10^{-4}} = 10^{r+4} \quad \leftarrow 10^0 \dots 10^4 \end{aligned}$$

Andrew Ng

2.2: Appropriate scale for hyperparameters

>Log scale:

e.g learning rate: considering range: 0.0001 - 1

if sampling uniformly at random, 90% percent sample would be between 0.1-1, -->90% of the resources to search between 0.1-1, only 10% to search 0.0001-0.1
-->more reasonable to search on a log scale, and sampling uniformly at random on this new scale. now have more resource to deviated to 0.0001 - 0.1.

code:

```
r=-4* np.random.rand() --->[-4, 0];  
a=10^r; --->[10^-4, 10^0]
```

summarize:

if try to sample on range $10^a - 10^b$, on log scale, then sample r uniformly at random, between a and b.

if sample on log scale, take the low value , take logs to figure out what is a, and take log of high value to find b, then sampling r uniformly between a and b.

2.3 Hyperparameters for exponentially weighted averages

β : considered range: 0.9 - 0.999

Remember β means: -->using $1/(1-\beta)$ previous data:

β : considered average days range is : 10 - 1000

Using linear scale on range 0.9-0.999, does not make sense;

better explore: $1-\beta$: 0.1 - 0.0001. use log scale: --> [-3, -1]

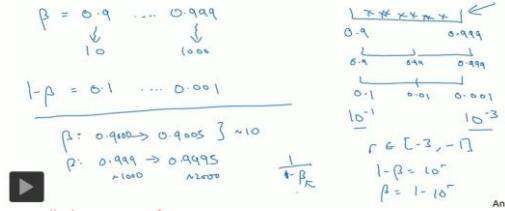
Reason do not use linear scale directly is:

when sampling close to 1, the sensitivity of the result changes, even with very small changes to β .

e.g when β change from 0.9 (10 days)to 0.9005 (10.5 days), there is hardly any change in result, but when β change from 0.999 (1000days)to 0.9995 (2000 days average), this will have a huge impact on what algorithm is doing. ($1/(1-\beta)$ is very sensitive to β when it is close to 0)

need to sample more densely in the region of when β close to 1, so can more efficient in terms of how distribute the samples to explore the space of possible outcomes more efficiently.

Hyperparameters for exponentially weighted averages



Summary:

In case you don't end up making the right scaling decision on some hyperparameter choice, don't worry to much about it.

Even if you sample on the uniform scale, where some of the scale would have been superior, you might still get okay results, especially if you use a coarse to fine search: in later iterations, will focus more on the most useful range of hyperparameter values to sample.

3: Hyperparameters tuning in practice

Item purpose: Final tips and tricks for how to organize hyperparameter search process.

3.1 Re-test hyperparameters occasionally;

>1. cross-fertilization: NLP, Vision, Speech, Ads, logistics.....

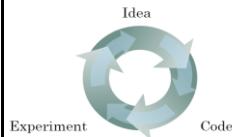
Deep learning today is applied to many different application areas and that intuitions about hyperparameter settings from one application area may or may not transfer to a different one. There is a lot of cross-fertilization among different applications' domains, so for example, I've seen ideas developed in the computer vision community, such as Convnets or ResNets, which we'll talk about in a later course, successfully applied to speech. I've seen ideas that were first developed in speech successfully applied in NLP, and so on.

So one nice development in deep learning is that people from different application domains do read increasingly research papers from other application domains to look for inspiration for cross-fertilization.

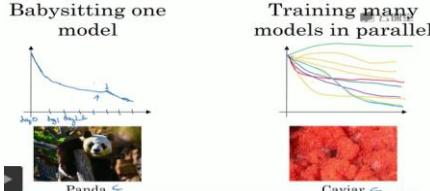
>Intuitions do get stale. Re-evaluate occasionally.

In terms of your settings for the hyperparameters, though, I've seen that intuitions do get stale. So even if you work on just one problem, say logistics, you might have found a good setting for the hyperparameters and kept on developing your algorithm, or maybe seen your data gradually change over the course of several months, or maybe just upgraded servers in your data center. And because of those changes, the best setting of your hyperparameters can get stale. So I recommend maybe just retesting or reevaluating your hyperparameters at least once every several months to make sure that you're still happy with the values you have.

Re-test hyperparameters occasionally



- NLP, Vision, Speech, Ads, logistics,
- Intuitions do get stale. Re-evaluate occasionally.



3.2 Two method in searching for hyperparameters,

>1. babysit one model.

>>Conditions: have maybe a huge data set but not a lot of computational resources- CPUs or GPU, train only one model or a very small number of models at a time.

>>Idea: gradually babysit that model even as it's training : watching performance and patiently nudging parameter up or down

Every day you kind of look at algorithm performance and try nudging up and down your parameters: kind of babysitting the model one day at a time even as it's training over a course of many days or over the course of several different weeks..

E.G: on first initialize your parameter as random and then start training. And gradually watch learning curve, maybe the cost function or your dataset error or something else. gradually decrease over the first day. Then at the end of day one, you might say, gee, looks it's learning quite well, I'm going to try increasing the learning rate a little bit and see how it does. And then maybe it does better. And then that's your Day 2 performance. And after two days you say, okay, it's still doing quite well. Maybe I'll fill the momentum term a bit or decrease the learning variable a bit now, and then you're now into Day 3. And every day you kind of look at it and try nudging up and down your parameters. And maybe on one day you found your learning rate was too big. So you might go back to the previous day's model, and so on. But you're kind of babysitting the model one day at a time even as it's training over a course of many days or over the course of several different weeks.

also called panda approach:

When pandas have children, they have very few children, usually one child at a time, and then they really put a lot of effort into making sure that the baby panda survives. So that's really babysitting. One model or one baby panda.

>2 Train many models in parallel.

train different model with different hyperparameters setting at the same time, and generate different learning curve.

so this way you can try a lot of different hyperparameter settings and then just maybe quickly at the end pick the one that works best.

also called caviar (鱼子酱) strategy:

like fish that lay over 100 million eggs in one mating season. But the way fish reproduce is they lay a lot of eggs and don't pay too much attention to any one of them, just see that hopefully one of them, or maybe a bunch of them, will do well.

Summary:

>1. So the way to choose between these two approaches is really a function of how much computational resources you have.

If have enough computers to train a lot of models in parallel, then by all means take the caviar approach and try a lot of different hyperparameters and see what works.

Otherwise, use babysitting, but could baby sitting on more than one models:

In some application domains, e.g. online advertising settings and computer vision applications, there's just so much data and the models you want to train are so big that it's difficult to train a lot of models at the same time. It's really application dependent of course, but those communities use the panda approach a little bit more: kind of baby sitting on a single model along and nudging the parameters up and down and trying to make this one model work.

Although, even the panda approach, having trained one model and then seen it work or not work, maybe in the second week or the third week, maybe I should initialize a different model and then baby that one along just like even pandas, I guess, can have multiple children in their lifetime, even if they have only one, or a very small number of children, at any one time.

2.3-4 Normalizing activations in a network_Batch Normalization

Item purpose:

Make neural network much more robust to the choice of hyperparameters. It doesn't work for all neural networks, but when it does, it can make the hyperparameter search much easier and also make training go much faster.

In the rise of deep learning, one of the most important ideas has been an algorithm called batch normalization. Batch normalization makes hyperparameter search problem much easier, makes neural network much more robust. The choice of hyperparameters is a much bigger range of hyperparameters that work well, and will also enable you to much more easily train even very deep networks.

2. Implementing Batch Norm:

1. Normalizing: logistic and deep model

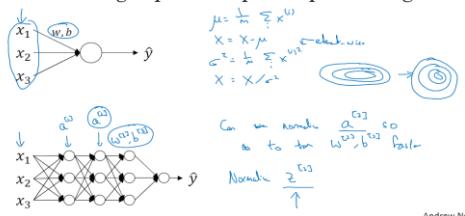
In logistic regression, normalizing input feature could make cost function from elongated to more round over parameter w, b, easier for algorithm like gradient descent to optimize. But in deep model, not only have input features x, but have activation input a[i].

--> for any hidden layer, normalize $a[i-1]/Z[i-1]$ to make training of w[i], b[i] more efficient? $z[i-1]_i \text{ norm} = z[i-1]_i - u) / \sqrt{\sigma^2 + \epsilon}$: mean 0 and variance 1; dimension: $(n_i \times 1)$

Note: have choices to normalize a[i] or Z[i], better take z[i] as default. (why?)

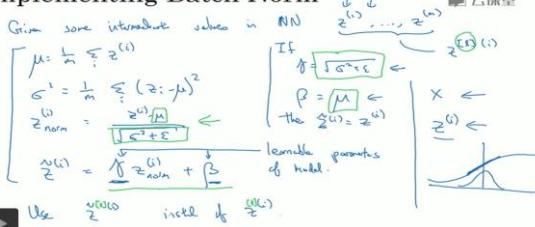
对每层进行Batch Normalization的时候一般都是先进行normalization再传递active function的，好处就是可以让输入更多地落在active function (比如sigmoid和tanh)的微分比较大的地方也就是零附近。

Normalizing inputs to speed up learning



Andrew Ng

Implementing Batch Norm



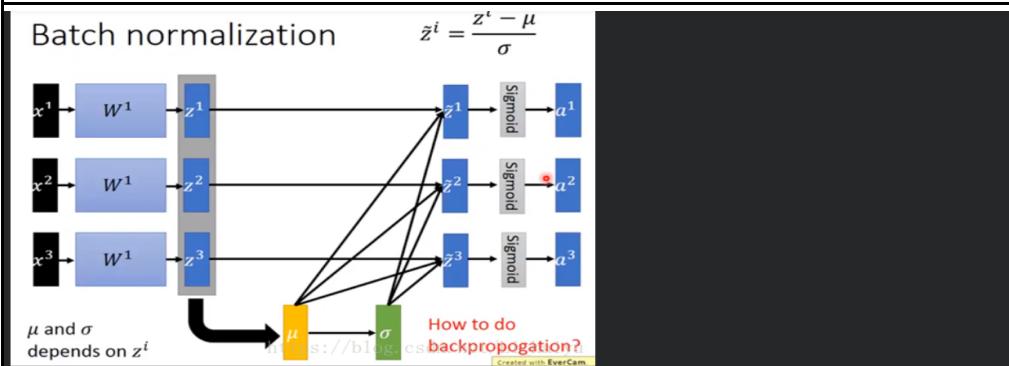
Intuition:

>1. Normalizing the input features x can help learning in a neural network, and what batch norm does is it applies that normalization process not just to the input layer, but to the values even deep in some hidden layer in the neural network. So it will apply this type of normalization to normalize the mean and variance of some of your hidden units' values, z.

>2. But one difference between the training input and these hidden unit values is you might not want your hidden unit values be forced to have mean 0 and variance 1.

e.g: while using sigmoid activation function, might want z to have a larger variance or have a mean that's different than 0, so have a wider range not only under sigmoid function linear area, in order to better take advantage of the nonlinearity of the sigmoid function rather than have all your values be in just this linear regime.

So, parameters y and beta, can now make sure that z[i] values have the range of values that you want, normalizes mean and variance of these hidden unit values, z[i] to have some fixed mean and variance. (separate y, beta for each units in layer i?-vector of output z size?)



在每层的active function之前都进行一次batch normalization那么做了batch normalization之后的方向传播又是怎么进行的呢？

在方向后传播的时候，一路反向传播回来的时候是会通过o然后通过μ然后去update Z的，如果在BP过程中没有考虑μ和σ的话是不对的，因为μ和σ是依赖于Z的，改动了Z就相当于改动了μ和σ，所以在training的时候μ和σ是会被考虑进去的。

现在我们把Z normalize成z，然后在active function输入之前还可以增加β和γ参数，这两个参数在training的时候也会被学习到。

2.3-5 Batch norm into neural work

1. Add batch norm to neural work:

>1. Norm add between Z & a(Z).

Each hidden units can take it as computing two things: z and a.

In batch norm:

$X \rightarrow w[1], b[1] \rightarrow Z[1] \rightarrow y[1], \beta[1] \rightarrow -Z[1] \rightarrow A[1]$, $w[2], b[2] \rightarrow Z[2] \rightarrow y[2], \beta[2] \rightarrow -Z[2] \dots$

normalization: between $Z[l]$ and $A[l]$, and use \tilde{z} not the z_{norm}

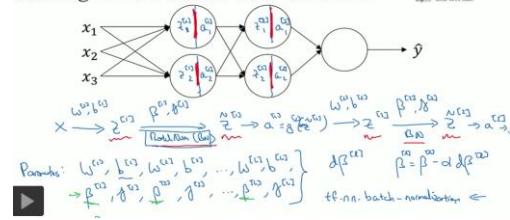
Parameters: $w[l], b[l], y[l], \beta[l]$.

$b[l]$ could be removed now, and \tilde{z} mean now is controlled by $\beta[l]$

>2. gradient descent with batch norm: could still use momentum, RMSProp, Adam;

>3. Batch norm application in coding: could use one code in python: tf.nn.BatchNorm

Adding Batch Norm to a network



2. Working with mini batch

In practice, Batch Norm is usually applied with mini-batches of your training set.

>1. Batch Norm in mini batch:

>>1. first mini-batch:
 $X[1] \rightarrow w[1], b[1] \rightarrow Z[1] \rightarrow y[1], \beta[1] \rightarrow -Z[1] \rightarrow A[1]$, $w[2], b[2] \rightarrow Z[2] \rightarrow y[2], \beta[2] \rightarrow -Z[2] \dots$

Note: computer mean and variance of the $Z[l]$ on just this mini batch and then re-scale $Z[l]$ by y, β to get $\tilde{Z}[l]$, all this is done on the first mini-batch.

finish one step gradient descent on the mini-batch. then go to second mini-batch, do something similar.

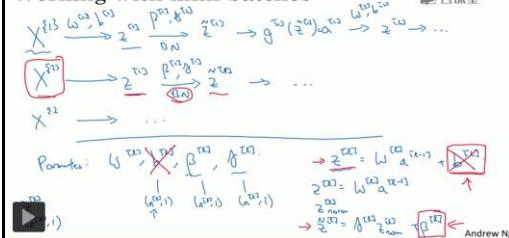
$X[2] \rightarrow w[1], b[1] \rightarrow Z[1] \rightarrow y[1], \beta[1] \rightarrow -Z[1] \rightarrow A[1]$, $w[2], b[2] \rightarrow Z[2] \rightarrow y[2], \beta[2] \rightarrow -Z[2] \dots$

Note: normalizing $Z[l]$ using just the data in second mini-batch: computing mean μ and variance σ^2 on second mini-batch data, re-scale $Z[l]$ by y, β to get $\tilde{Z}[l]$.

>2. $b[l]$ removed as parameter or set to 0: as any constant added will get cancelled out by the mean normalization step.
 $\beta[l]$ instead used to decide mean of $\tilde{Z}[l]$

>3. $Y[l]$ & $\beta[l]$: dimension: $((n[l], 1) / (n[l], m))$. Each hidden units in layer l (in all dimension) have 1 Y value for setting variance: scale the mean and variance of each of the hidden values (total $n[l]$ values)

Working with mini-batches



Implementing gradient descent

for $t = 1 \dots \text{num Mini-Batches}$

Compute forward pass on $X^{(t)}$.
In each hidden layer, use BN to map $Z^{(t)}$ with $\tilde{Z}^{(t)}$.
Use backprop to compute $dW^{(t)}, dB^{(t)}, dY^{(t)}$.
Update params: $W^{(t+1)} := W^{(t)} - \alpha dW^{(t)}$
 $\beta^{(t+1)} := \beta^{(t)} - \alpha d\beta^{(t)}$
 $y^{(t+1)} := \dots$

Works w/ momentum, RMSProp, Adam.

3. Implementing batch norm in gradient descent with mini-batch

>1. one epoch: go through all iterations in one epoch (iteration number = batch size / mini-size)

>2. Inner loop: for each iteration/mini-batch

>>1. forward pro:

in each hidden layer, use BN to replace $Z[l]$ with $\tilde{Z}[l]$:

Make sure in that mini-batch, the values and version of the normalized mean and variance is $\tilde{Z}[l]$:

>>2. back pro: compute $dW[l], dY[l], d\beta[l]$

>>3. update parameters: $w[l] := w[l] - \alpha dW[l]; Y[l] := a * dY[l], \beta[l] := a * d\beta[l]$

in this step could use momentum, RMSProp, Adam for fast learning move

Note: Instead of taking this gradient descent update mini-batch, could use the updates given by these other algorithms as we discussed in the previous week's videos (Momentum, RMSProp, ADAM). Some of these other optimization algorithms as well can be used to update the parameters β and y that Batch Norm added to algorithm.

2.3-6 Why does Batch Norm work

Normalizing all input features X , to take on a similar range of values can speed up learning. So, one intuition behind why batch norm works is, this is doing a similar thing, but further values in your hidden units and not just for your input there.

1. Learning on shifting input distribution

Weight, later or deeper that network. more robust to changes of weight in earlier layers

>1. covariate shift: Data distribution changing influence

e.g: shallow NN like logistic regression, for cat classification:

Have trained algorithm with training set of black cats, --> now apply classifier to colored cat classifier may not do well.

algorithm trained on one data set, may not do well on data set generalized from original data set - decision boundary may not apply well to the new generalized data set.

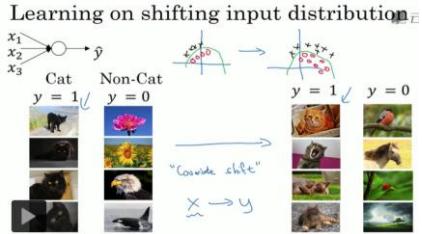
Covariate shift: data distribution changing.

for the learned X to Y mapping, if the distribution of X changes, then might need to retrain learning algorithm.

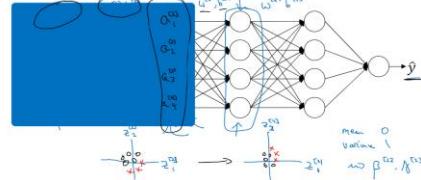
(note: normally require training set, dev set and test set same distribution, training set cover all the dev/test set features --> algorithm generalize better on dev/test set.)

(dev set and test set have to be same distribution; training set may come from different distribution - as need large data to train, data may from other distribution.)

And this is true even if the ground truth function, mapping from X to Y , remains unchanged, which it is in this example, because the ground truth function is: is this picture a cat or not. And the need to retrain your function becomes even more acute or it becomes even worse if the ground truth function shifts as well.



Why this is a problem with neural networks?



2. covariate shift in neural network

2.1 e.g:

layer 3: $w[3], b[3]$, input $a[2]$
 >1. Layer 3 purpose: take input $a[2]$ and find a way to map $a[2]$ to output $a[3] = y^*$.
 --> $w[3], b[3], w[4], b[4]$ learned and updated to do this mapping well

>2. while whole network also adapting param in earlier layers $w[1], \dots, b[2]$.
 when these parameters change, $a[2]$ as input to layer 3 also change-->

layer 3 have covariate shift problem. (input **distribution** changed)

- covariate shift: changes in previous layer (hyperparameter, training data) will influence current layer hyperparameter : eg. Layer L, $W[l], b[l]$ have trained based on input $a[l-1]$. (who trained on batch $X[t]$ training data), to estimate output y . if $a[l-1]$ changed, $W[l], b[l]$ normally have to change to estimate same y (if output do not change).

cost = sum cost in each layer c_l over all layers

--> cost error in hidden layers could not pass to next/cost function?-backprop;

covariate shift is forward prop, how to influence on backprop? -->

parameter update-->output of each layer change-->cost error shifting,

2.2 Batch norm helps covariate shifting

Batch norm: Make parameter e.g weight in the later or deeper network layer, more robust to changes of weight in earlier layers.

>1. batch norm: reduce the amount that the distribution of the hidden units values shifting around.

when parameters in the neural network earlier layer changes --> the $Z[l]$ will change.

Batch norm will make sure hidden unit values ($Z[l]$)'s mean and variance remain same, -->holding hidden unit values in certain region under all changes: limit the amount to which updating the parameters in earlier layer can affect the distribution of hidden unit values. which is the next layer sees and has to learn on.

So batch norm reduces the problem of the input values changing, it really causes these values to become more stable, so that the later layers of the neural network has more firm ground to stand on.

>2. And even as the earlier layers keep learning, the amounts that this forces the later layers to adapt to as early as layer changes is reduced or, it weakens the coupling between what the early layers parameters has to do and what the later layers parameters have to do. -->allows each layer of the network to learn by itself, a little bit more independently of other layers, and this has the effect of speeding up of learning in the whole network.

>3. Especially from the perspective of one of the later layers of the neural network, the earlier layers don't get to shift around as much, because they're constrained to have the same mean and variance. And so this makes the job of learning on the later layers easier.

2.3 Batch norm as regularization

>1. Each mini-batch X_t , has the values Z_t , has the values Z_l , scaled by the mean and variance computed on just that one mini-batch.

--> as the mean and variance computed on just that mini-batch as opposed to computed on the entire data set, that mean and variance has a little bit of noise in it, the scaling process of $Z[l]$ is a little bit noisy as well.

Batch Norm as regularization

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch.
- This adds some noise to the values $Z[l]$ within that minibatch. So similar to dropout, it adds some noise to each hidden layer's activations.
- This has a slight regularization effect.

>2. This adds some noise to the values $Z[l]$ within that minibatch.

So similar to dropout, it adds some noise to each hidden layer's activations.

dropout noise: takes a hidden unit (activation output) and multiplies it by zero with some probability.

So your dropout has multiple of noise because it's multiplied by zero or one, whereas batch norm has multiples of noise because of scaling by the standard deviation, as well as additive noise because it's subtracting the mean. while the estimates of the mean and the standard deviation are noisy.

So similar to dropout, batch norm therefore has a slight regularization effect : by adding noise to the hidden units, it's forcing the downstream hidden units not to rely too much on any one hidden unit.

note: regularization effect: reduce weight

>3. This has a slight regularization effect:

Because the noise added is quite small, and therefore it is a very slight regularization effect.

Could choose to use batch norm together with dropout if want the more powerful regularization effect of dropout.

Note:

>1. using a bigger mini-batch size, will reducing this noise and therefore also reducing this regularization effect.

So that's one strange property of dropout which is that by using a bigger mini-batch size, you reduce the regularization effect.

>2. Do not use batch norm as a regularizer.

That's really not the intent of batch norm, but sometimes it has this extra intended or unintended effect on your learning algorithm. But don't turn to batch norm as a regularization. Use it as a way to normalize your hidden units activations and therefore speed up learning.

>3. Batch norm handles data one mini-batch at a time. It computes mean and variances on mini-batches.

So at test time, you try and make predictions, try and evaluate the neural network, you might not have a mini-batch of examples, you might be processing one single example at the time. So, at test time you need to do something slightly differently to make sure your predictions make sense.

2.3-7 Batch norm at test

1. mean & variance in minibatch-one layer: $u=1/m \sum z^{(i)}$: m mini batch size;
 u and 6^2 in layer l, for scaling $Z^{(l)}$, are computed on the entire mini batch examples.

2. mean & variance in test-each layer:
when have only one example in test dev, normalization by using it's own u and 6^2 does not make sense($u=z$, $D(z)=0$, \rightarrow no value). \rightarrow separate estimate of u and 6^2 :
Using expansion weight average-across mini batch (save storage): averager the values in layer l over previous values used/computed (in different iterations/mini-batch):
e.g:
 $X^{(1)} \rightarrow 6^2 / u^{(1)} - 1$: mean computed and used in scaling in layer l, with first mini-batch example.
 $X^{(2)} \rightarrow 6^2 / u^{(1)} - 1$: mean computed and used in scaling in layer l, with second mini-batch example.
 \dots
 $X^{(t)} \rightarrow 6^2 / u^{(1)} - t$: mean computed and used in scaling in layer l, with t-th mini-batch example.

keep running average of u and 6^2 that computed and used in each layer as train the neural network across different mini-batches.

and use in test example: $Z_{norm} = z - u / \sqrt{6^2 + \epsilon}$

Batch Norm at test time

2.3-8 Softmax layer introduction

1. Softmax:

Softmax: generalization of logistic regression, used when have multiple classification, :make predictions where you're trying to recognize one of multiple classes, rather than just recognize two classes.

Output: probability of each of classes: $p(c_i | x)$.

C: classes number;

Output vector dimension: (C,1)

Output sum of all output is =1.

softmax layer used in neural network as output layer, in order to generate these outputs.

2. Softmax layer

outlayer: $Z^{[L]} = W^{[L]} \cdot A^{[L-1]} + b^{[L]}$

activation function:

$t_j = e^{Z^{[L]}_j}$, same dimension as $Z^{[L]}$

$A^{[L-1]}_j = t_j / \sum t_j$ over all j classes/units in output layer
 $= e^{Z^{[L]}_j} / \sum e^{Z^{[L]}_j}$ over all elements.

$A^{[L]}_j$: Same dimension as $Z^{[L]}$; normalize t.

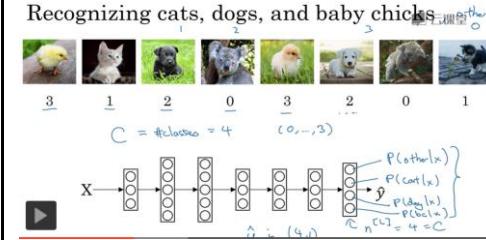
activation function: $A^{[L]} = g(Z^{[L]})$

take a vector $Z^{[L]}$ and output a vector $A^{[L]}$.

(note: for m examples: $Z^{[L]}$: C x m matrix:

$\rightarrow A^{[L]} = e^{Z^{[L]}} / (\text{np. sum}(e^{Z^{[L]}}), \text{dim}=1)$

)



3. softmax example: no hidden layer

$$Z^{[1]} = W^{[1]}X + b$$

$$A^{[1]} = g(W^{[1]})$$

softmax layer with C output classes can represent this type of decision boundaries-several linear decision boundaries, allows it to separate out the data into C classes.

Coloring in the right pic./input based on which one of the c outputs have the highest probability.
can kind of see it like a generalization of logistic regression, with sort of linear decision boundaries

Note:

>1. without any hidden layer, decision boundary between any two classes will be linear (e.g in right)
>2 with hidden layers, can learn more complex non-linear decision boundaries to separate out multiple different classes.

3. softmax example: no hidden layer

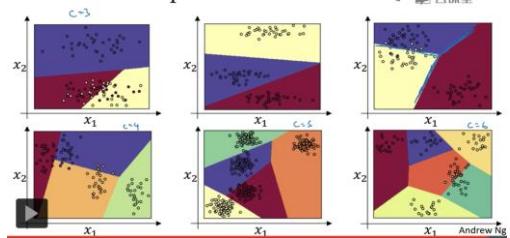
```
[x1, x2]T -> w: [neuro1, neuro2, neuro3] -> y [y1,y2,y3].
Z=W[w1, w12; w21, w22]* [x1, x2]T + b (default 0)=[y1, y2,y3]
y1=exp (w1*x1+w12*x2)/(y1+y2+y3)
y2=exp (w21*x1+w22*x2)/(y1+y2+y3)
y3=exp (w31*x1+w32*x2)/(y1+y2+y3)
y1+y2+y3=1
```

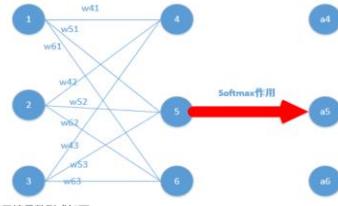
right pic: train softmax with three output labels on the data, colors in plot show output threshold in softmax classifier, and coloring in the input based on which one of the three output units have highest probability -->like generalization of logistic regression, with sort of linear boundaries (more complicated non-linear boundaries could be learned with hidden layers-activation function-->Y=g(Z)=g(w11*x1+w12*x2), but with more than two classes.

softmax derivative

Softmax layer

Softmax examples





交叉熵函数形式如下：

$$Loss = - \sum_i y_i \ln a_i$$

其中 y_i 代表我们的真实值， a_i 代表我们softmax求出的值。 i 代表的是输出结点的标号！在上面例子， i 就可以取值为4,5,6三个结点（当然我这里只是为了简单，真实应用中可能有很多结点）

现在看起来是不是感觉复杂了，居然还有累和，然后还要求导，每一个 a_i 都是softmax之后的形式！

但是实际上不是这样的，我们往往在真实中，如果只预测一个结果，那么在目标中只有一个结点的值为1，比如我认为在该状态下，我想要输出的是第四个动作（第四个结点），那么训练数据的输出就是 $a_4 = 1, a_5=0, a_6=0$ ，哎呀，这太好了。除了一个为1，其它都是0，那么所谓的求和符合，就是一个例子，我可以去掉啦！

为了形式化说明，我这里认为训练数据的真实输出为第 j 个为1，其它均为0！

那么 $Loss$ 就变成了 $Loss = -y_j \ln a_j$ 累和已经去掉了，太好了。现在我们要开始求导数了！

我们在整理一下上面公式，为了更加明白的看出相关变量的关系：

其中 $y_j = 1$ ，那么形式变为 $Loss = -\ln a_j$

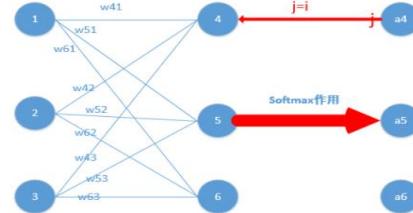
那么形式越来越简单了，求导分析如下：

参数的形式在该例子中，总共分为 $w_{41}, w_{42}, w_{43}, w_{51}, w_{52}, w_{53}, w_{61}, w_{62}, w_{63}$ 这些，那么比如我要求出 w_{41}, w_{42}, w_{43} 的偏导，就需要将 $Loss$ 函数求偏导传到结点4，然后再利用链式法则继续求导即可。举个例子此时求 w_{41} 的偏导为：

$$\text{if } j = i: \\ \frac{\partial a_j}{\partial z_i} = \frac{\partial}{\partial z_i} \left(\frac{e^{z_j}}{\sum_k e^{z_k}} \right) \\ = \frac{(e^{z_j})' \cdot \sum_k e^{z_k} - e^{z_j} \cdot e^{z_j}}{(\sum_k e^{z_k})^2} \\ = \frac{e^{z_j}}{\sum_k e^{z_k}} \cdot \frac{e^{z_j}}{\sum_k e^{z_k}} = a_j(1 - a_j)$$

$j=4$ 对应例子里就是如下图所示：

比如我选定了 j 为4，那么就是说我现在求导传到4结点这！



那么由上面求导结果再乘以交叉熵损失函数求导

$Loss = -\ln a_j$ ，它的导数为 $-\frac{1}{a_j}$ ，与上面 $a_j(1 - a_j)$ 相乘为 $a_j - 1$ （形式非常简单，这说明我只要正向求一次得出结果，然后反向传递度的时候，只需要将它结果减1即可，后面还会举例子！）那么我们可以得到 $Loss$ 对于4结点的偏导就求出了（这里假定4是我们的预计输出）

第二种情况为：

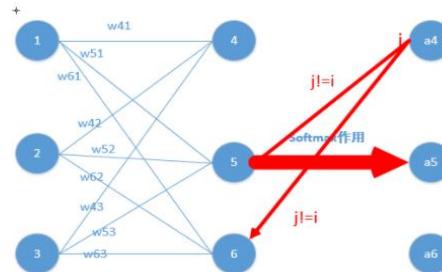
$$\begin{aligned} \frac{\partial Loss}{\partial w_{41}} &= \frac{\partial Loss}{\partial a_4} \cdot \frac{\partial a_4}{\partial z_4} \cdot \frac{\partial z_4}{\partial w_{41}} \\ &= \boxed{-\frac{1}{a_4}} \cdot \frac{\partial a_4}{\partial z_4} \cdot \boxed{\frac{\partial z_4}{\partial w_{41}}} \\ &\quad \text{则关键为求出 } \frac{\partial a_4}{\partial z_4} \end{aligned}$$

$w_{51} \dots w_{63}$ 等参数的偏导同理可以求出，那么我们的关键就在于 $Loss$ 函数对于结点4,5,6的偏导怎么求，如下：

这里分为两种情况：

$$\text{if } j \neq i: \\ \frac{\partial a_j}{\partial z_i} = \frac{\partial}{\partial z_i} \left(\frac{e^{z_j}}{\sum_k e^{z_k}} \right) \\ = \frac{0 \cdot \sum_k e^{z_k} - e^{z_j} \cdot e^{z_i}}{(\sum_k e^{z_k})^2} \\ = -\frac{e^{z_j}}{\sum_k e^{z_k}} \cdot \frac{e^{z_i}}{\sum_k e^{z_k}} = -a_j a_i$$

这里对应我的例子图如下，我这时对的是 j 不等于 i ，往前传：



那么由上面求导结果再乘以交叉熵损失函数求导

$Loss = -\ln a_j$ ，它的导数为 $-\frac{1}{a_j}$ ，与上面 $-a_j a_i$ 相乘为 a_i （形式非常简单，这说明我只要正向求一次得出结果，然后反向传递度的时候，只需要将它结果保存即可，后续例子会讲到）这里就求出了除4之外的其它所有结点的偏导，然后利用链式法则继续传递过去即可！我们的问题也就解决了！

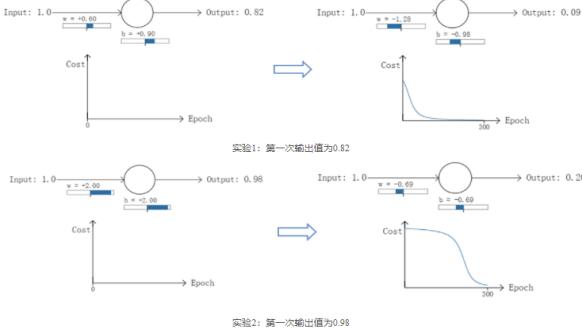
loss function:

交叉熵代价函数（Cross-entropy cost function）是用来衡量人工神经网络（ANN）的预测值与实际值的一种方式。与二次代价函数相比，它能更有效地促进ANN的训练。在介绍交叉熵代价函数之前，本文先简要介绍二次代价函数，以及其存在的不足。

1. 二次代价函数的不足

ANN的设计目的之一是为使机器可以像人一样学习知识。人在学习分析新事物时，当发现自己犯的错误越大时，改正的力度就越大。比如打篮球：当运动员发现自己投篮的位置离正确位置越远，那么他调整的投篮角度就应该越大，篮球就更容易投进篮筐。同样，我们希望：ANN在训练时，如果预测值与实际值的误差越大，那么在反向传播过程中，各参数调整的幅度就越大，从而训练更快收敛。然而，如果使用二次代价函数训练ANN，看到的实际效果是，如果误差越大，参数调整的幅度可能更小，训练更缓慢。

以一个神经元的一次分类训练为例，进行两次实验（ANN常用的激活函数为sigmoid函数，该实验也采用该函数）：输入一个相同的样本数据x=1.0（该样本对应的实际分类y=0）；两次实验各自随机初始化参数，从而在各自的第一次前向传播后得到不同的输出值，形成不同的代价（误差）：



在实验1中，随机初始化参数，使得第一次输出值0.82（该样本对应的实际值为0）；经过300次迭代训练后，输出值由0.82降到0.09，逼近实际值。而在实验2中，第一次输出值0.98，同样经过300次迭代训练，输出值却降到了0.20。

从两次实验的代价曲线中可以看出：实验1的代价随着训练次数增加而快速降低，但实验2的代价在一开始下降得非常缓慢；直观上看，初始的误差越大，收敛得越慢。

其实，误差大导致训练缓慢的原因在于使用了二次代价函数。二次代价函数的公式如下：

$$C = \frac{1}{2n} \sum_x \|y(x) - a^T(x)\|^2$$

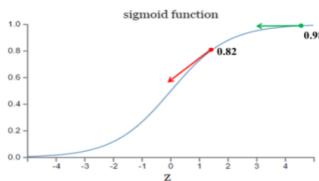
其中， C 表示代价， x 表示样本， y 表示实际值， a 表示输出值， n 表示样本的总数。为简单起见，同样一个样本为例进行说明。此时二次代价函数为：

$$C = \frac{(y - a)^2}{2}$$

目前训练ANN最有效的算法是反向传播法。简而言之，训练ANN就是通过反向传播代价，以减少代价为导向，调整参数。参数主要有：神经元之间的连接权重 w ，以及每个神经元本身的偏置 b 。调整的方式是采用梯度下降算法（Gradient descent）。沿梯度方向调整参数大小， w 和 b 的梯度推导如下：

$$\begin{aligned}\frac{\partial C}{\partial a} &= (a - y)\sigma'(z)x \\ \frac{\partial C}{\partial b} &= (a - y)\sigma'(z)\end{aligned}$$

其中， z 表示神经元的输入， σ 表示激活函数。从以上公式可以看出， w 和 b 的梯度跟激活函数的梯度成正比，激活函数的梯度越大， w 和 b 的大小调整得越快，训练收敛得就越快。而神经网络常用的激活函数为sigmoid函数，该函数的曲线如下所示：



如图所示，实验2的初始输出值（0.98）对应的梯度明显小于实验1的输出值（0.82），因此实验2的参数梯度下降得比实验1慢。这就是初始的代价（误差）越大，导致训练越慢的原因。与我们的期望不符，即：不能像人一样，错误越大，改正的幅度越大，从而学习得越快。

可能有人会说，那就选择一个梯度不变化或变化不明显的激活函数不就解决问题了吗？图标图破，那样虽然简单粗暴地解决了这个问题，但可能会引起其他更多更麻烦的问题。而且，类似sigmoid这样的函数（比如tanh函数）有很多优点，非常适合用来做激活函数，具体请自行google之。

2. 交叉熵代价函数

换个思路，我们不计算梯度函数，而是换掉二次代价函数，改用交叉熵代价函数：

$$C = -\frac{1}{n} \sum_x [y \ln a + (1 - y) \ln(1 - a)]$$

其中， x 表示样本， n 表示样本的总数。那么，重新计算参数 w 的梯度：

$$\begin{aligned}\frac{\partial C}{\partial w_j} &= -\frac{1}{n} \sum_x \left(\frac{y}{\sigma(z)} - \frac{(1-y)}{1-\sigma(z)} \right) \frac{\partial \sigma}{\partial w_j} \\ &= -\frac{1}{n} \sum_x \left(\frac{y}{\sigma(z)} - \frac{(1-y)}{1-\sigma(z)} \right) \sigma'(z) x_j \\ &= \frac{1}{n} \sum_x \frac{\sigma'(z) x_j}{\sigma(z)(1-\sigma(z))} (\sigma(z) - y) \\ &= \frac{1}{n} \sum_x x_j (\sigma(z) - y)\end{aligned}$$

其中（具体证明见附录）：

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

因此， w 的梯度公式中原来的 $\sigma'(z)$ 被消掉了；另外，该梯度公式中的 $\sigma(z) - y$ 表示输出值与实际值之间的误差。所以，当误差越大，梯度就越大，参数 w 调整得越快，训练速度也就越快。同理可得， b 的梯度为：

$$\frac{\partial C}{\partial b} = \frac{1}{n} \sum_x (\sigma(z) - y)$$

实际情况证明，交叉熵代价函数带来的训练效果往往比二次代价函数要好。

3. 交叉熵代价函数是如何产生的？

以偏置 b 的梯度计算为例，推导出交叉熵代价函数：

$$\begin{aligned}\frac{\partial C}{\partial b} &= \frac{\partial C}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial b} \\ &= \frac{\partial C}{\partial a} \cdot \sigma'(z) \cdot \frac{\partial (wx+b)}{\partial b} \\ &= \frac{\partial C}{\partial a} \cdot \sigma'(z) \\ &= \frac{\partial C}{\partial a} \cdot a(1-a)\end{aligned}$$

在第1小节中，由二次代价函数推导出来的 b 的梯度公式为：

$$\frac{\partial C}{\partial b} = (a - y)\sigma'(z)$$

为了消掉该公式中的 $\sigma'(z)$ ，我们想找到一个代价函数使得：

$$\frac{\partial C}{\partial b} = (a - y)$$

即：

$$\frac{\partial C}{\partial a} \cdot a(1-a) = (a - y)$$

对两侧求积分，可得：

$$C = -[y \ln a + (1 - y) \ln(1 - a)] + \text{constant}$$

而这就是前面介绍的交叉熵代价函数。

softmax的log似然代价函数（公式求导）

2. softmax配合log似然代价函数训练ANN

在上一篇博文交叉熵代价函数中讲到，二次代价函数在训练ANN时可能会导致训练速度变慢的问题。那就是，初始的输出值离真实值越远，训练速度越慢。这个问题可以通过采用交叉熵代价函数来解决。其实，这个问题也可以采用另一种方法解决，那就是采用softmax激活函数，并采用log似然代价函数（log-likelihood cost function）来解决。

log似然代价函数的公式为：

$$C = -\sum_k y_k \log a_k$$

其中， a_k 表示第 k 个神经元的输出值。 y_k 表示第 k 个神经元对应的真值，取值为0或1。

我们来简单理解一下这个代价函数的含义。在ANN中输入一个样本，那么只有一个神经元对应了该样本的正确类别；若这个神经元输出的概率值越高，则按照以上的代价函数公式，其产生的代价就越小；反之，则产生的代价就越高。

为了检验softmax和这个代价函数也可以解决上述所说的训练速度变慢问题，接下来的重点就是推导ANN的权重 w 和偏置 b 的梯度公式。以偏置 b 为例：

$$\begin{aligned}\frac{\partial C}{\partial b_j} &= \frac{\partial C}{\partial z_j} \cdot \frac{\partial z_j}{\partial b_j} \\ &= \frac{\partial C}{\partial z_j} \cdot \frac{\partial (w_j a_k + b_j)}{\partial b_j} \\ &= \frac{\partial C}{\partial z_j} \left(-\sum_k y_k \log a_k \right) \\ &= -\sum_k y_k \frac{1}{a_k} \frac{\partial a_k}{\partial z_j} \\ &\quad (\text{use softmax derivative})\end{aligned}$$

同理可得：

$$\frac{\partial C}{\partial w_j} = a_k^{j-1} (a_k^j - y_j)$$

从上述梯度公式可知，softmax函数配合log似然代价函数可以很好地训练ANN，不存在学习速度变慢的问题。

2.3-9 softmax understanding

1. **hardmax:** set max. probability =1, others 0.

while softmax is contrast to hardmax, more gentle mapping from z to these probabilities.

note: two classification previous talked , activation fucntion is Relu/sigmoid, last layer $z[i]$ is a number (dimension 1), -->out put is also a number -one dimension.

2. Train softmax regression

>1. Loss function singal input: $J=\sum y[i] * \log(y^-[i])$ over all classes.

singal example loss: sum of each class loss= label class * algorithm predicted probility for that class.

note: labeled $y[i]=1$, or 0 (1 means labeled $y[i]$, corresponding to $x[i]$).-->loss function $J=-y[i] * \log(y^-[i])$

min J-->max log ($y^-[i]$)-->max $y^-[i]$ -corresponding
>2, minimize J, means looks at whatever the ground truth calls in training data, and try to make the corresponding probability of that class as high as possible.--a form of **maximum likelihood estimation**.

>3. Batch cost function: batch input: $J=1/m * \sum (J[i])$ over all m training examples)

minimize batch cost function with gradient descent.

Note: labeled y and prediction matrix dimension is (C, m)

Understanding softmax

$$z^{[l]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix}$$

$$C = 4 \quad g^{[l]}(z^{[l]})$$

$$\phi^{[l]} = g^{[l]}(z^{[l]}) = \begin{bmatrix} e^5/(e^5 + e^2 + e^{-1} + e^3) \\ e^2/(e^5 + e^2 + e^{-1} + e^3) \\ e^{-1}/(e^5 + e^2 + e^{-1} + e^3) \\ e^3/(e^5 + e^2 + e^{-1} + e^3) \end{bmatrix} = \begin{bmatrix} 0.892 \\ 0.042 \\ 0.012 \\ 0.114 \end{bmatrix}$$

Softmax regression generalizes logistic regression to C classes.

Loss function

$$L(z^{[l]}, y) = -\sum_{i=1}^C y_i \log \hat{y}_i$$

$$J(\omega^{[l]}, b^{[l]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{[i]}, y^{[i]})$$

$$-y_i \log \hat{y}_i = -\log \hat{y}_{y_i}$$

$$\hat{Y} = \begin{bmatrix} \hat{y}^{[1]} \\ \vdots \\ \hat{y}^{[m]} \end{bmatrix} = \begin{bmatrix} * \\ * \\ \vdots \\ * \end{bmatrix}$$

3. Gradient descent with softmax

forward:

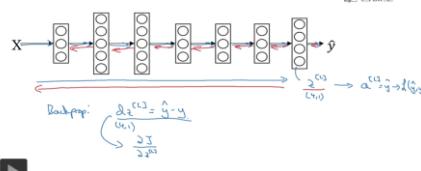
$Z[l] \rightarrow A[l] = y^A \rightarrow L(y^A, y)$

backward:
see above

Note: when use framework, usually just need to focus on getting the forward prop right, and so long as specify a program framework, the forward prop pass, the program framework will figure out how to do back prop.

(note: when $C=2$, softmax reduce to sigmoid :
softmax: $a_1 = e^z / (e^z + e^z) = 1 / (1 + e^{z_2 - z_1})$;
sigmoid: $a = 1 / (1 + e^{-z})$)

Gradient descent with softmax



2.3_10 Framework_deep learning

Have learned to implement deep learning algorithms more or less from scratch using Python and NumPY, to understand what these deep learning algorithms are really doing.

But when implement more complex models, such as convolutional neural networks or recurring neural networks, or as you start to implement very large models that is increasingly not practical, at least for most people, is not practical to implement everything yourself from scratch. Fortunately, there are now many good deep learning software frameworks that can help you implement these models.

1. Deep learning frameworks

Below are some of the leading deep learning software frameworks. Each of these frameworks has a dedicated user and developer community and I think each of these frameworks is a credible choice for some subset of applications.

Note: these frameworks are often evolving and getting better month to month,

Criteria choosing frame work

>1. easy of programming: developing neural network and itrating on it , and deploying it for production

>2. running speed:

especially training on large data sets, some frameworks will let you run and train your neural network more efficiently than others.

>3. truly open (open source with good governance): how much you trust a framework will remain open source for a long time, rather than just being under the control of a single company.

for a framework to be truly open, it needs not only to be open source but it needs good governance as well. Unfortunately, in the software industry some companies have a history of open sourcing software but maintaining single corporation control of the software. And then over some number of years, as people start to use the software, some companies have a history of gradually closing off what was open source, or perhaps moving functionality into their own proprietary cloud services. So one thing I pay a bit of attention to is how much you trust that the framework will remain open source for a long time rather than just being under the control of a single company, which for whatever reason may choose to close it off in the future even if the software is currently released under open source.

>4. But at least in the short term depending on you're preferences of language, whether prefer Python or Java or C++ or something else, and depending on what application you're working on, whether this can be division or natural language processing or online advertising or something else, Multiple of these frameworks could be a good choice.

Deep learning frameworks

- Caffe/Caffe2
- CNTK
- DL4J
- Keras
- Lasagne
- mxnet
- PaddlePaddle
- TensorFlow
- Theano
- Torch

- Choosing deep learning frameworks
- Ease of programming (development and deployment)
 - Running speed
 - - Truly open (open source with good governance)

Summarize:

on programming frameworks by providing a higher level of abstraction than just a numerical linear algebra library,any of these program frameworks can make it more efficient as develop machine learning applications.

Motivating problem

2.3.11 TensorFlow Practice

Motivation problem:
e.g. $J(w) = w^2 - 10w + 5$

a very similar structure of program can be used to train neural network where can have some complicated cost function $J(w, b)$ ---> then could use tensorflow automatically to find parameters that minimize this cost function.

1. Tensorflow coding:

>Define:
>>1. variable (w).
>>2. cost function: could use tensorflow add and multiply as well as other functions , tensorflow knows how to compute derivatives of these functions, which already built in to these functions.
-->this is why can only have to implement basically forward prop.

Note: could also use simple way of these functions in cost fuction formula, which are overloaded already in tensorflow.

2. Get training data into Tensorflow program

e.g:
define x (no y involved here): name placeholder, float 32, (3,1) vector;

note: placehoder in TensorFlow is a variable whose value assign later.
convenient way to get trainging data into the cost function (using feed_dict).
when use mini-batch, then on different iterations feed in different subsets of training sets.

```
-->x = tf.placeholder(tf.float32, [3,1])
cost = x[0][0]*w**2 + x[1][0]*w + x[2][0]

coefficients = np.array ([1], [-10], [25])

train = tf.train.GradientDescentOptimizer(0.01).minimize (cost)
# could change 'GradientDescentOptimizer' to use other method like Adma;
# 0.01 here is learning rate;
#purpose is minimize (cost)

for i in range (1000) # 1000 iterations.
session.run (train, feed_dict={x: coefficients}) # coefficients value feed to x, which control cost function.
```

Cost function computation in Tensorflow:
compute forward acc. to computing graph,
and all necessary backward functions have been built-in.

Thus using the built-in functions to compute the forward function, it can automatically do the backward functions as well to implement back prop through even very complicated functions and compute derivatives.--->make computation more efficient.

Code example

```
import numpy as np
import tensorflow as tf

coefficients = np.array([1, -10, 25])

w = tf.Variable(0, dtype=tf.float32)
x = tf.placeholder(tf.float32, [3,1])
cost = x[0][0]*w**2 + x[1][0]*w + x[2][0] # (w-5)**2
train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
init = tf.global_variables_initializer()
session = tf.Session()
with tf.Session() as session:
    session.run(init)
    session.run(train, feed_dict={x:coefficients})
    print(session.run(w))
```

Andrew Ng

3. Below three lines of code are quite idiomatic in TensorFlow,

>1.session = tf.Session ---> alternative : with tf.Session () as session: # with command in python is a little bit better at cleaning up in case there's error or exception while executing this inner loop.

>2. session.Run (int)

>3. print(session.run(w))

Summarize:

Tensorflow is powerful because all need to do is specify how to compute the cost fuction, then it can take derivatives, apply a gdient optimizer or an Adam optimizer or some other optimize rwith just pretty much one or two lines of codes.

Summary

1. Input normalization-->cost function is more like round bowl - in all paramter direction, fast gradient descent, faster iteration .

2. Batch norm: make each layer output shifting in limitation, each layer more independend to train-->faster algorithm parameter tranning.

3. Weight initialization: variance 1/n, or 2/n---->reduce gradient descent vanishing/exploring: by making variance of activation in forwad prop and derivative of state in backprop same (Tanh/sigmoid) or variance of state z in forward prop and derivative of activation in backprop same (Relu activatioin)

Loss function:

Principle: update parameter in big step when error /cost function is big.

Better: dw = $(y^{\wedge} - y) \cdot c$ (c: nothing to do with y^{\wedge}, y)

-->

1. choose cost function:

>>1. square error:

>>> ok for linear function: dw = $(y^{\wedge} - y) \cdot X.T$
>>> non-ok for logistic function: dw = $(y^{\wedge} - y) \cdot [y^{\wedge}(1-y^{\wedge})] \cdot X.T$ -->Can not garantte dw bigger when has worse prediction ($y^{\wedge} - y$) bigger

>>2. cross-error:

>>> for logistic function: dw = $(y^{\wedge} - y) \cdot X.T$

>>3. softmax error:

>>> loss = log y^{\wedge}_i : corresponding to max probabiltiy: maximum likelihood estimation

2. Choose activation function for output: linear , sigmoid, softmax

Normalizati

[BatchNormalization](#)、[LayerNormalization](#)、[InstanceNorm](#)、[GroupNorm](#)

1、综述

1.1 论文链接

1、Batch Normalization

<https://arxiv.org/pdf/1502.03167.pdf>

2、Layer Normalizaiton

<https://arxiv.org/pdf/1607.06450v1.pdf>

3、Instance Normalization

<https://arxiv.org/pdf/1607.08022.pdf>

https://github.com/DmitryUlyanov/texture_nets

1.2 介绍

归一化层，目前主要有这几个方法，Batch Normalization（2015年）、Layer Normalization（2016年）、Instance Normalization（2017年）、Group Normalization（2018年）、Switchable Normalization（2018年）；

将输入的图像shape记为[N, C, H, W]，这几个方法主要的区别就是在，

- batchNorm是在batch上，对NHW做归一化，对小batchsize效果不好；
- layerNorm在通道方向上，对CHW做归一化，主要对RNN作用明显；
- instanceNorm在图像像素上，对HW做归一化，用在风格化迁移；
- GroupNorm将channel分组，然后再做归一化；
- SwitchableNorm是将BN、LN、IN结合，赋予权重，让网络自己去学习归一化层应该使用什么方法。

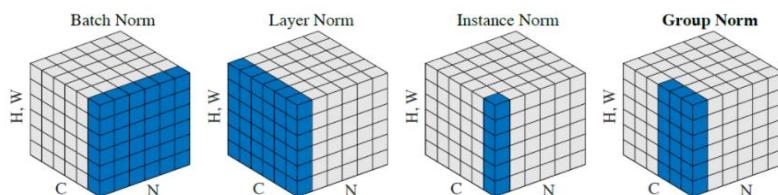


Figure 2. **Normalization methods**. Each subplot shows a feature map tensor, with N as the batch axis, C as the channel axis, and (H, W) as the spatial axes. The pixels in blue are normalized by the same mean and variance, computed by aggregating the values of these pixels.

3、Layer Normalization

batch normalization存在以下缺点：

- 对batchsize的大小比较敏感，由于每次计算均值和方差是在一个batch上，所以如果batchsize太小，则计算的均值、方差不足以代表整个数据分布；
- BN实际使用时需要计算并且保存某一层神经网络batch的均值和方差等统计信息，对于对一个固定深度的前向神经网络（DNN，CNN）使用BN，很方便；但对于RNN来说，sequence的长度是不一致的，换句话说RNN的深度不是固定的，不同的time-step需要保存不同的statistics特征，可能存在一个特殊sequence比其他sequence长很多，这样training时，计算很麻烦。（参考于<https://blog.csdn.net/lqfarmer/article/details/71439314>）

与BN不同，LN是针对深度网络的某一层的所有神经元的输入按以下公式进行normalize操作。

$$\mu^l = \frac{1}{H} \sum_{i=1}^H a_i^l \quad \sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^H (a_i^l - \mu^l)^2}$$

<https://blog.csdn.net/liuxiao214>

LN用于RNN效果比较明显，但是在CNN上，不如BN。

```
1 def ln(x, s):
2     _eps = 1e-5
3     output = (x - x.mean(1)[:,None]) / tensor.sqrt((x.var(1)[:,None] + _eps))
4     output = s[None, :] * output + b[None,:]
5     return output
```

用在四维图像上，

```
1 def LayerNorm(x, gamma, beta):
2
3     # x_shape:[B, C, H, W]
4     results = 0.
5     eps = 1e-5
6
7     x_mean = np.mean(x, axis=(1, 2, 3), keepdims=True)
8     x_var = np.var(x, axis=(1, 2, 3), keepdims=True)
9     x_normalized = (x - x_mean) / np.sqrt(x_var + eps)
10    results = gamma * x_normalized + beta
11    return results
```

4. Instance Normalization

BN注重对每个batch进行归一化，保证数据分布一致，因为判别模型中结果取决于数据整体分布。

但是图像风格化中，生成结果主要依赖于某个图像实例，所以对整个batch归一化不适合图像风格化中，因而对HW做归一化。可以加速模型收敛，并且保持每个图像实例之间的独立。

公式：

$$y_{tijk} = \frac{x_{tijk} - \mu_{ti}}{\sqrt{\sigma_{ti}^2 + \epsilon}}, \quad \mu_{ti} = \frac{1}{HW} \sum_{l=1}^W \sum_{m=1}^H x_{tilm}, \quad \sigma_{ti}^2 = \frac{1}{HW} \sum_{l=1}^W \sum_{m=1}^H (x_{tilm} - \mu_{ti})^2.$$

代码：

```
1 def Instancenorm(x, gamma, beta):
2
3     # x_shape:[B, C, H, W]
4     results = 0.
5     eps = 1e-5
6
7     x_mean = np.mean(x, axis=(2, 3), keepdims=True)
8     x_var = np.var(x, axis=(2, 3), keepdims=True)
9     x_normalized = (x - x_mean) / np.sqrt(x_var + eps)
10    results = gamma * x_normalized + beta
11
12    return results
```

6、Switchable Normalization

本篇论文作者认为，

- 第一，归一化虽然提高模型泛化能力，然而归一化层的操作是人工设计的。在实际应用中，解决不同的问题原则上需要设计不同的归一化操作，并没有一个通用的归一化方法能够解决所有应用问题；
- 第二，一个深度神经网络往往包含几十个归一化层，通常这些归一化层都使用同样的归一化操作，因为手工为每一个归一化层设计操作需要进行大量的实验。

因此作者提出自适配归一化方法——Switchable Normalization (SN) 来解决上述问题。与强化学习不同，SN 使用可微分学习，为一个深度网络中的每一个归一化层确定合适的归一化操作。

公式：

SN has an intuitive expression

$$\hat{h}_{ncij} = \gamma \frac{h_{ncij} - \sum_{k \in \Omega} w_k \mu_k}{\sqrt{\sum_{k \in \Omega} w'_k \sigma_k^2 + \epsilon}} + \beta, \quad (3)$$

$$w_k = \frac{e^{\lambda_k}}{\sum_{z \in \mathcal{Z}} e^{\lambda_z}}, \quad k \in \{\text{in}, \text{ln}, \text{bn}\},$$

$$\mu_{\text{in}} = \frac{1}{HW} \sum_{i,j}^{H,W} h_{ncij}, \sigma_{\text{in}}^2 = \frac{1}{HW} \sum_{i,j}^{H,W} (h_{ncij} - \mu_{\text{in}})^2,$$

$$\begin{aligned}\mu_{\text{ln}} &= \frac{1}{C} \sum_{c=1}^C \mu_{\text{in}}, \quad \sigma_{\text{ln}}^2 = \frac{1}{C} \sum_{c=1}^C (\sigma_{\text{in}}^2 + \mu_{\text{in}}^2) - \mu_{\text{ln}}^2, \\ \mu_{\text{bn}} &= \frac{1}{N} \sum_{n=1}^N \mu_{\text{in}}, \quad \sigma_{\text{bn}}^2 = \frac{1}{N} \sum_{n=1}^N (\sigma_{\text{in}}^2 + \mu_{\text{in}}^2) - \mu_{\text{bn}}^2, \quad (5)\end{aligned}$$

代码：

```
1 def SwitchableNorm(x, gamma, beta, w_mean, w_var):  
2     # x_shape:[B, C, H, W]  
3     results = 0.  
4     eps = 1e-5  
5
```

```
6     mean_in = np.mean(x, axis=(2, 3), keepdims=True)
7     var_in = np.var(x, axis=(2, 3), keepdims=True)
8
9     mean_ln = np.mean(x, axis=(1, 2, 3), keepdims=True)
10    var_ln = np.var(x, axis=(1, 2, 3), keepdims=True)
11
12    mean_bn = np.mean(x, axis=(0, 2, 3), keepdims=True)
13    var_bn = np.var(x, axis=(0, 2, 3), keepdims=True)
14
15    mean = w_mean[0] * mean_in + w_mean[1] * mean_ln + w_mean[2] * mean_bn
16    var = w_var[0] * var_in + w_var[1] * var_ln + w_var[2] * var_bn
17
18    x_normalized = (x - mean) / np.sqrt(var + eps)
19    results = gamma * x_normalized + beta
20
21    return results
```


ion method

GroupNorm、SwitchableNorm总结

4、Group Normalization

<https://arxiv.org/pdf/1803.08494.pdf>

5、Switchable Normalization

<https://arxiv.org/pdf/1806.10779.pdf>

<https://github.com/switchablenorms/Switchable-Normalization>

2、Batch Normalization

首先，在进行训练之前，一般要对数据做归一化，使其分布一致，但是在深度神经网络训练过程中，通常以送入网络的每一个batch训练，这样每个batch具有不同的分布；此外，为了解决internal covariate shift问题，这个问题定义是随着batch normalization这篇论文提出的，在训练过程中，数据分布会发生变化，对下一层网络的学习带来困难。

所以batch normalization就是强行将数据拉回到均值为0，方差为1的正太分布上，这样不仅数据分布一致，而且避免发生梯度消失。

此外，internal covariate shift和covariate shift是两回事，前者是网络内部，后者是针对输入数据，比如我们在训练数据前做归一化等预处理操作。

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots m\}$;

Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

Algorithm 1: Batch Normalizing Transform, applied to activation x over a mini-batch.

<https://blog.csdn.net/linxi/>

算法过程：

- 沿着通道计算每个batch的均值u
- 沿着通道计算每个batch的方差 σ^2
- 对x做归一化， $x' = (x - u) / \sqrt{\sigma^2 + \epsilon}$
- 加入缩放和平移变量 γ 和 β ，归一化后的值， $y = \gamma x' + \beta$

加入缩放平移变量的原因是：保证每一次数据经过归一化后还保留原有学习来的特征，同时又能完成归一化操作，加速训练。这两个参数是用来学习的参数。

```

1 import numpy as np
2
3 def Batchnorm(x, gamma, beta, bn_param):
4
5     # x_shape:[B, C, H, W]
6     running_mean = bn_param['running_mean']
7     running_var = bn_param['running_var']
8     results = 0.
9     eps = 1e-5
10
11    x_mean = np.mean(x, axis=(0, 2, 3), keepdims=True)
12    x_var = np.var(x, axis=(0, 2, 3), keepdims=True)
13    x_normalized = (x - x_mean) / np.sqrt(x_var + eps)
14    results = gamma * x_normalized + beta
15
16    # 因为在测试时是单个图片测试，这里保留训练时的均值和方差，用在后面测试时用
17    running_mean = momentum * running_mean + (1 - momentum) * x_mean
18    running_var = momentum * running_var + (1 - momentum) * x_var
19
20    bn_param['running_mean'] = running_mean
21    bn_param['running_var'] = running_var
22
23    return results, bn_param

```

5、Group Normalization

主要是针对Batch Normalization对小batchsize效果差，GN将channel方向分group，然后每个group内做归一化，算(C//G)*H*W的均值，这样与batchsize无关，不受其约束。

公式：

Group Norm. Formally, a Group Norm layer computes μ and σ in a set \mathcal{S}_i defined as:

$$\mathcal{S}_i = \{k \mid k_N = i_N, \lfloor \frac{k_C}{C/G} \rfloor = \lfloor \frac{i_C}{C/G} \rfloor\}. \quad (7)$$

Here G is the number of groups, which is a pre-defined hyper-parameter ($G = 32$ by default). C/G is the number of channels per group. $\lfloor \cdot \rfloor$ is the floor operation, and “ $\lfloor \frac{k_C}{C/G} \rfloor = \lfloor \frac{i_C}{C/G} \rfloor$ ” means that the indexes i and k are in the same group of channels, assuming each group of channels are stored in a sequential order along the C axis. GN computes μ and σ along the (H, W) axes and along a group of $\frac{C}{G}$ channels. The computation of GN is illustrated in

Figure 2 (rightmost), which is a simple case of 2 groups ($G = 2$) each having 3 channels. <https://blog.csdn.net/liuxiao21>

代码：

```
1 def GroupNorm(x, gamma, beta, G=16):
2
3     # x_shape:[B, C, H, W]
4     results = 0.
5     eps = 1e-5
6     x = np.reshape(x, (x.shape[0], G, x.shape[1]/16, x.shape[2], x.shape[3]))
7
8     x_mean = np.mean(x, axis=(2, 3, 4), keepdims=True)
9     x_var = np.var(x, axis=(2, 3, 4), keepdims=True)
10    x_normalized = (x - x_mean) / np.sqrt(x_var + eps)
11    results = gamma * x_normalized + beta
12
13    return results
```

提出的原因：在小 batch 的样本中，Batch Normalization 估计的值不准。一般用在很大的模型中，这时 batch size 就很小。

思路：数据不够，通道来凑。每个样本的特征分为几组，每组特征分别计算均值和方差。可以看作是 Layer Normalization 的基础上添加了特征分组。

注意事项：

- 不再有 running_mean 和 running_var
- \$\gamma\$ 和 \$\beta\$ 为逐通道的

结果比较：

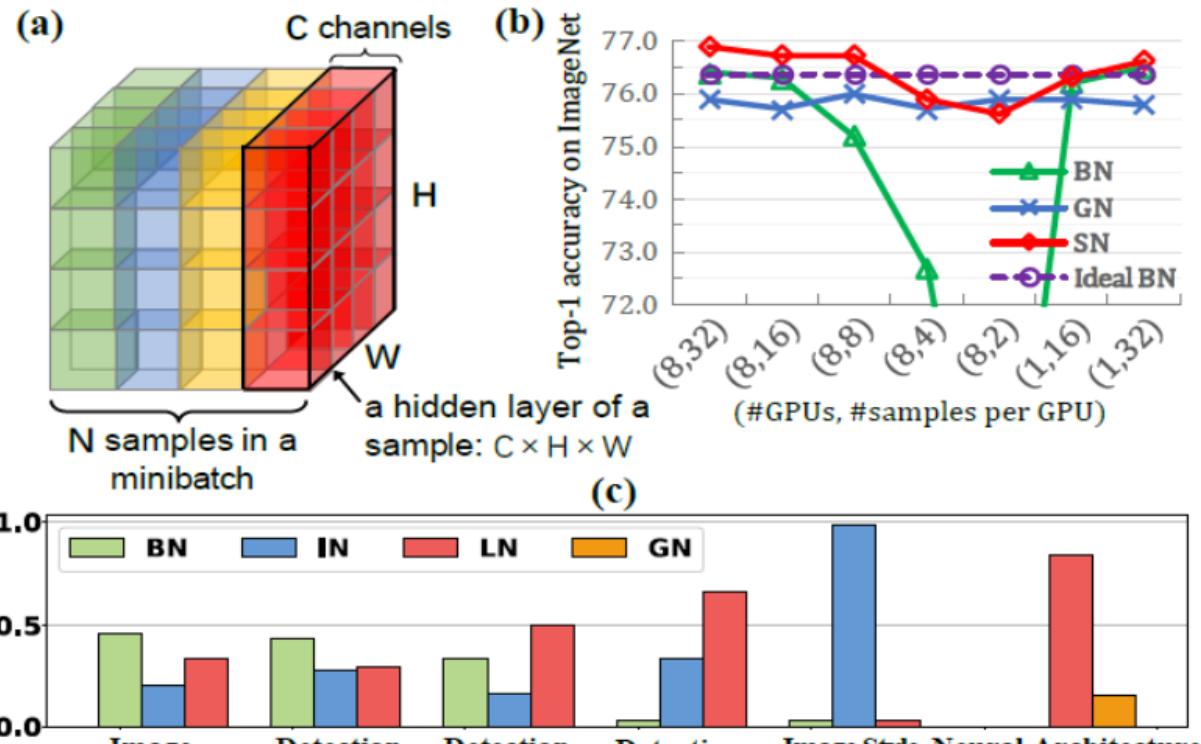


Image Classification	Detection Backbone	Detection Box Head	Detection Mask Head	Image Style Transfer	Neural Architecture Search (LSTM)
					https://blog.csdn.net/Pytxiao214

Figure 1: **(a)** In a CNN, the feature maps of a hidden convolutional layer is of size $N \times C \times H \times W$, which represents N samples in a minibatch, C channels of a layer, height and width of a channel respectively. We have $N = 4$ in (a) as an illustrative example. Different normalization methods compute statistics (means and variances) along different axes of N, C, H, W . **(b)** The top-1 accuracies of ResNet50 trained with SN on ImageNet are compared with BN and GN in different batch sizes. For instance, all methods are compared to an ideal case, ‘ideal BN’, whose accuracies are 76.4%, which is the result of (8,32) of BN. This ideal case isn’t obtainable in practice. In fact, when the minibatch size decreases, BN’s accuracies drop significantly, while SN and GN both maintain reasonably good performance. SN surpasses or is comparable to both BN and GN in all settings. **(c)** shows that SN adapts to various network architectures and tasks, by learning an importance weight to select each normalization method (*i.e.* a

ratio between 0 and 1; all the weights of each task sum to 1).
ao214

【深度学习】参数初始化

一、参数初始化原则

参数初始化作为模型训练的起点，决定了模型训练的初始位置。选择的好坏很大程度影响收敛的速度与最终的训练结果。一般来说，参数初始化遵循以下几点要求：

- 不建议同时初始化为一个值，容易出现“对称失效”
 - 最好保证均值为0，正负交错，参数大致数量相等
 - 初始化参数不要太大或太小，前者会导致梯度发散难以训练，后者会导致特征对后面的影响逐渐变小
-

二、常见的初始化方法

1、全零初始化与随机初始化

如果神经元的权重被初始化为0，在第一次更新的时候，除了输出之外，所有的中间层的节点的值都为零。一般神经网络拥有对称的结构，那么在进行第一次误差反向传播时，更新后的网络参数将会相同，在下一次更新时，相同的网络参数学习提取不到有用的特征，因此深度学习模型都不会使用0初始化所有参数。

随机初始化顾名思义，一般值大了容易饱和，小了对后面的影响太小，激活函数的效果也不好。

3、Kaiming初始化

Kaiming初始化，也称之为he初始化。是残差网络的作者何凯明在这篇论文中提出了ReLU网络的初始化方法。

(1) 正态化的kaiming初始化——he_normal

它从以 0 为中心，标准差为 $\text{stddev} = \sqrt{2 / \text{fan_in}}$ 的截断正态分布中抽取样本，其中 fan_in 是权值张量中的输入单位的数量

(2) 标准化的kaiming初始化——he_uniform

它从 $[-\text{limit}, \text{limit}]$ 中的均匀分布中抽取样本，其中 limit 是 $\sqrt{6 / \text{fan_in}}$ ，其中 fan_in 是权值张量中的输入单位的数量

4、lecun初始化

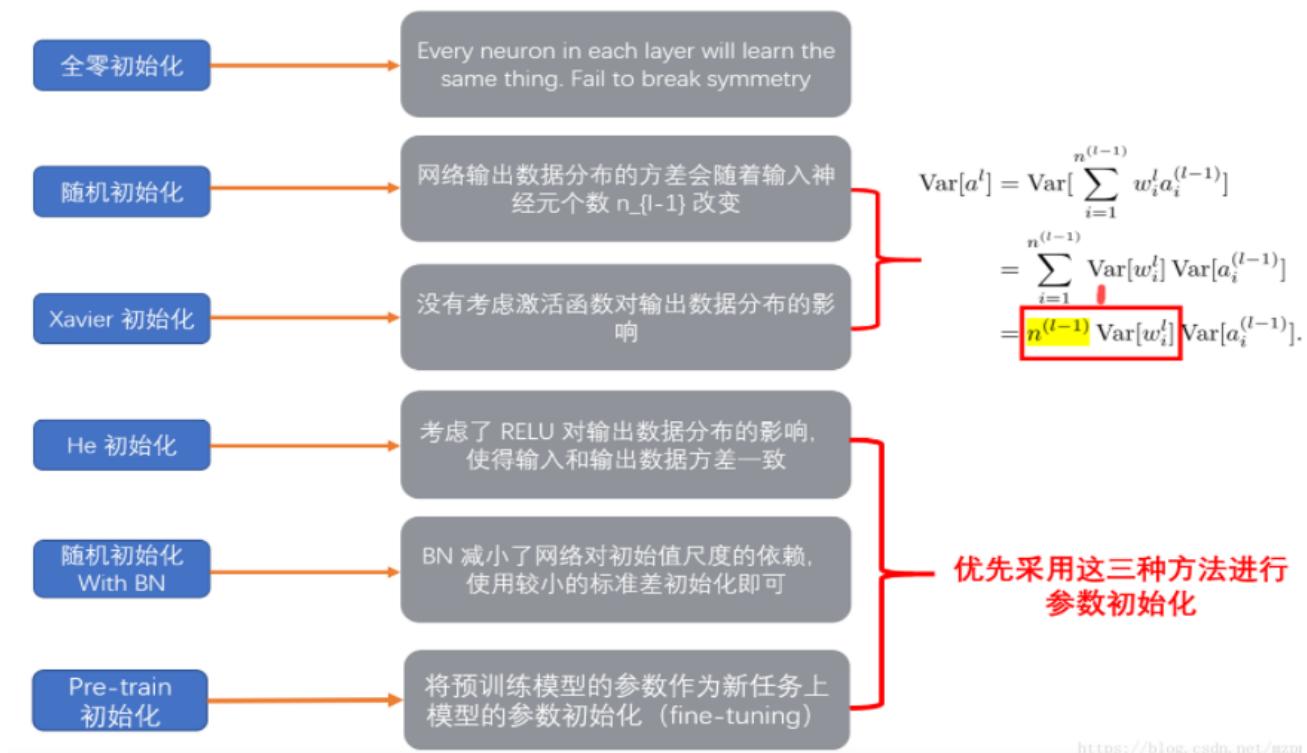
(1) 正态化的lecun初始化——lecun_normal

它从以 0 为中心，标准差为 $\text{stddev} = \sqrt{1 / \text{fan_in}}$ 的截断正态分布中抽取样本，其中 fan_in 是权值张量中的输入单位的数量。

(2) 标准化的lecun初始化——lecun_uniform

它从 $[-\text{limit}, \text{limit}]$ 中的均匀分布中抽取样本，其中 limit 是 $\sqrt{3 / \text{fan_in}}$ ，fan_in 是权值张量中的输入单位的数量。

四、初始化选择建议



<https://blog.csdn.net/mzpi>

Batch Normalization

2、Glorot初始化方法

(1) 正态化的Glorot初始化——glorot_normal

Glorot 正态分布初始化器，也称为 Xavier 正态分布初始化器。它从以 0 为中心，标准差为 $\text{stddev} = \sqrt{2 / (\text{fan_in} + \text{fan_out})}$ 的截断正态分布中抽取样本。其中 fan_in 是权值张量中的输入单位的数量， fan_out 是权值张量中的输出单位的数量。

(2) 标准化的Glorot初始化——glorot_uniform

Glorot 均匀分布初始化器，也称为 Xavier 均匀分布初始化器。

它从 $[-\text{limit}, \text{limit}]$ 中的均匀分布中抽取样本，其中 $\text{limit} = \sqrt{6 / (\text{fan_in} + \text{fan_out})}$ ， fan_in 是权值张量中的输入单位的数量， fan_out 是权值张量中的输出单位的数量。

(3) Glorot初始化器的缺点

首先有一个共识必须先提出：神经网络如果保持每层的信息流动是同一方差，那么会更加有利于优化。不然大家也不会去争先恐后地研究各种normalization方法。不过，Xavier Glorot认为还不够，应该增强这个条件，好的初始化应该使得各层的激活值和梯度的方差在传播过程中保持一致，这个被称为Glorot条件。如果反向传播每层梯度保持近似的方差，则信息能反馈到各层。而前向传播激活值方差近似相等，有利于平稳地学习。当然为了做到这一点，对激活函数也必须作出一些约定。

(1) 激活函数是线性的，至少在0点附近，而且导数为1。

(2) 激活值关于0对称。

这两个都不适用于sigmoid函数和ReLU函数，而适合tanh函数。

三、Batch Normalization

随着网络层数的增加，分布逐渐发生偏移，之所以收敛慢，是因为整体分布往非线性函数取值区间的上下限靠近。这会导致反向传播时梯度消失。BN就是通过规范化的手段，把每层神经网络任意神经元这个输入值的分布强行拉回到均值0方差1的标准正态分布，使得激活输入值落入非线性函数中比较敏感的区域。可以让梯度变大，学习收敛速度快，能大大加快收敛速度。

这张图已经能很好地反映初始化的一些思路了，由于现在batch norm比较流行所以可能对初始化要求不是那么高，一般用relu的话使用he初始化会好一些，或者直接用norm batch。

使用 RELU (without BN) 激活函数时，最好选用 He 初始化方法，将参数初始化为服从高斯分布或者均匀分布的较小随机数。

使用 BN 时，减少了网络对参数初始值尺度的依赖，此时使用较小的标准差(eg: 0.01)进行初始化即可，初始化方法没有那么重要了。

另外关于ReLU和BN的顺序，原paper建议将BN层放置在ReLU前，因为ReLU激活函数的输出非负，不能近似为高斯分布。实验表明，放在前后的差异似乎不大，甚至放在ReLU后还好一些。使用梯度下降法对参数进行优化学习时，非常关键的一个问题是如何合理地初始化参数值，为此提出了Xavier初始化和He初始化等方法，而使用BN操作后，对于参数初始化的要求就没那么高了。而且也可以使用大的学习率，另外BN也能够有效地缓解梯度消失的问题。

