

Setting Up Your ML Application

Train / Dev / Test Sets

Data	training set	dev	test
		hold-out	
		cross-validation	
		- development set	

Mismatched train/test distribution

make sure dev and test sets come from the same distribution.

devset: what we make decisions based on (which algorithm to choose)

a problem \Rightarrow training set

high quality cat pictures from web

test set

cat pictures from users using your app
(low resolution images)

if you only have a train / test set, your train acts like a dev because you use that set to tune hyperparameters and algorithmic choices

Bias and Variance

bias \Rightarrow high bias: many mistakes on training data (underfitting)

variance \Rightarrow high variance: very sensitive to small fluctuations in training set (overfitting)

Train set error: 1%
Dev set error: 11%

} Overfitting: not generalizing well (high variance)

Train set error: 15%
Dev set error: 16%

} Underfitting: performs poorly (high bias)
but generalizes well

Train set error: 15%
Dev set error: 30%

} high bias & high variance (you can get BOTH)

BUT, all of this analysis is based on the assumption that

HUMAN ERROR $\approx 0.0\%$

OR MORE
GENERALLY

Optimal (Bayes) ERROR: 0.0%

Basic Recipe for ML

High Bias? (diagnose by: looking at training data performance)

- Bigger network (almost always helps) (doesn't hurt variance)
- Train Longer
- Try different NN architecture

High Variance? (diagnose by: dev set performance)

- Get more data (doesn't hurt bias)
- Regularization
- Try finding better NN Architecture

Bias / Variance Tradeoff: In the modern DNN era, getting bigger network almost always reduces bias w/o hurting variance and getting more data almost always reduces variance w/o hurting bias.

Regularizing your NN

Regularization

Logistic Regression (Developing the Argument)

$$\min_{w,b} J(w,b)$$

$w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$

λ : regularization parameter

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

\nwarrow norm

L2 Regularization: $\|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w$
(Euclidean Norm)

L1 Regularization: $\frac{1}{m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1$

Sparse model
feature selection

L1 regularization produces a sparse model, a model that has most of its parameters equal to zero.

L1 makes feature selection by deciding which features are essential for prediction and which ones are not.

λ : Regularization Parameter

Another hyperparameter \Rightarrow we can decide its value with validation set.

Neural Network

$$J(w^{[0]}, b^{[0]}, \dots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|^2$$

$$\|w^{[l]}\|^2 = \sum_{i=1}^{n^{[l]}} \|w_i^{[l]}\|^2 = \sum_{r=1}^{n^{[l]}} \sum_{j=1}^{n^{[l-1]}} (W_{i,j}^{[l]})^2$$

recall

$$W^{[l]} \in \mathbb{R}^{n^{[l]} \times n^{[l-1]}}$$

"Frobenius Norm"

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

sum of squared
of elements of
a matrix

ith row
row vector

for the width $L-1$
layer

we add gradient of regularization term
to gradients of weights

Change In Backprop

$$\frac{\partial}{\partial W} \left(\frac{1}{2} \frac{\lambda}{m} W^2 \right) = \frac{\lambda}{m} W$$

$$\delta W^{[l]} = \frac{1}{m} \delta W^{[l]} A^{[l-1]T} + \frac{\lambda}{m} W^{[l]}$$

How Does This Affect Backprop?

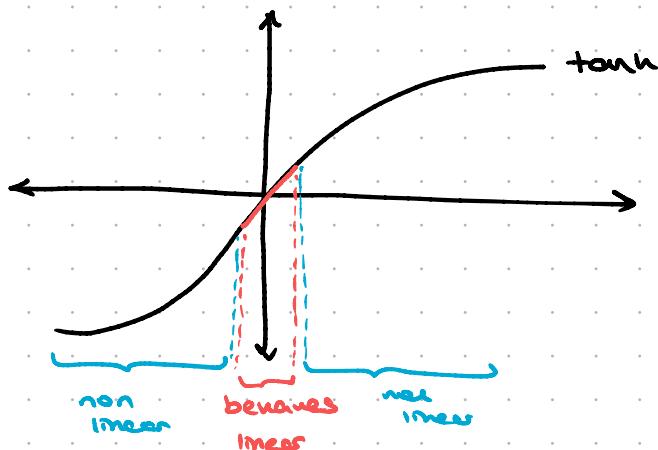
$$\begin{aligned} dW^{[c]} &= (\text{from backprop}) + \frac{\lambda}{m} W^{[c]} \\ W^{[c]} &:= W^{[c]} - \alpha dW^{[c]} \\ &:= W^{[c]} - \frac{\alpha \lambda}{m} W^{[c]} - \alpha (\text{from backprop}) \end{aligned}$$

this is why L2 regularization has been called "weight decay"

Why Does Regularization Help with Overfitting?

regularization incentivizes the model to set W close to zero, i.e. reducing the impact of some hidden layers.

↓
Creates a simpler + smaller network



$$\begin{aligned} \lambda \uparrow & \quad W^{[c]} \downarrow \\ z^{[c]} &= W^{[c]} \alpha^{[c-1]} + b^{[c]} \\ \text{every layer } & \approx \text{linear} \quad \longrightarrow \text{model is closer to linear regression (simpler decision boundaries)} \end{aligned}$$

Dropout Regularization

- 1) go through each layer \Rightarrow drop a unit based on a probability.
- 2) if dropped, remove in-going and outgoing connections of the unit
- 3) train the smaller network

At each step of training, dropout applies one step of backprop learning to a new version of the network. This new version is created by deactivating a randomly chosen subset of the units.

This is a low cost way to train a large ensemble of different networks.

Implementing Dropout ("Inverted Dropout")

we'll illustrate with layer $\ell = 3$

$\text{keep-prob} = 0.8$

$d3 = \text{np.random.rand}(a3.shape[0], a3.shape[1]) < \text{keep-prob}$ ↑ create a 0/1 matrix

$a3 = \text{np.multiply}(a3, d3)$

we have different units for different examples

$a3 /= \text{keep-prob}$ ↑ eliminate some hidden units

AT test time:

on expectation, 20% of elements of $a3$ will be zeroed out. \rightarrow so, to keep expected value at $z^{(u)}$ constant, divide it by keep-prob

→ don't use dropout.

↓ expected value of z .

→ because we divided by keep-prob at train time,
we ensure that the scaling

remains constant

doesn't change

→ no need for
fancy scaling

The intuition is because we are turning off units, no one unit can rely on another unit. So it has to spread the weight around, which reduces L_2 norm of weights.

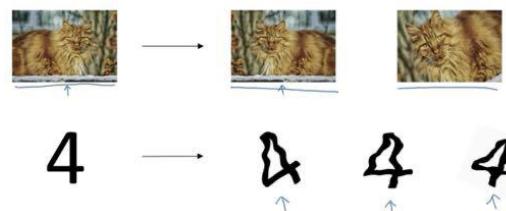
Other Regularization Methods

Augmenting training set

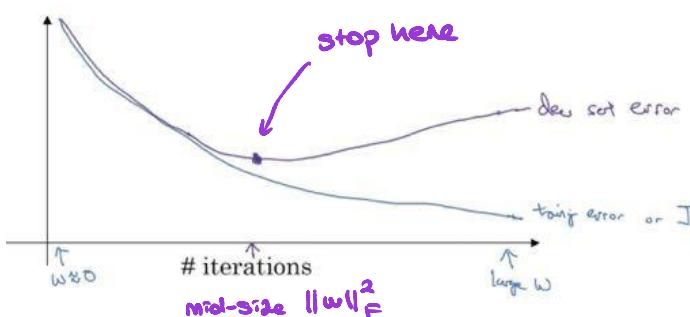
⇒ flip your images + add them to train set

⇒ random distortions of image

Note: these methods don't add a lot of value but a cheap way to get more data



Early Stopping



Orthogonalization:

- optimize cost function J
- gradient descent
- not overfit
- regularization

Early stopping works on both tasks at once. Violates orthogonalization.

Setting Up Your Optimization Problem

Normalizing Inputs

Subtract mean:

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

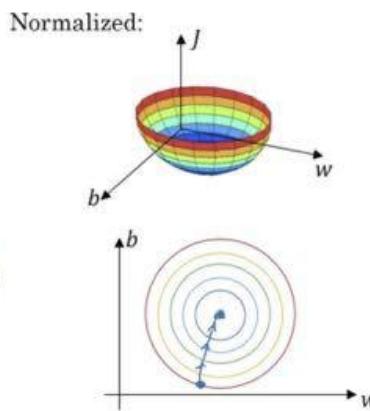
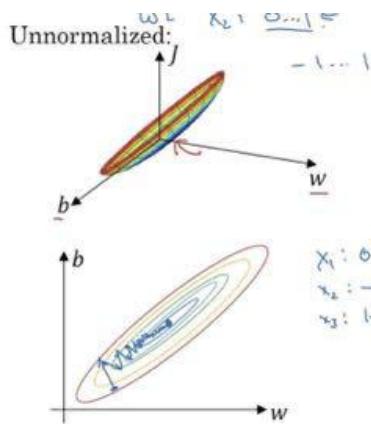
$$x := x - \mu$$

Normalize variance

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m x^{(i)} \cdot \sigma^2$$

$$x / \sigma$$

element-wise

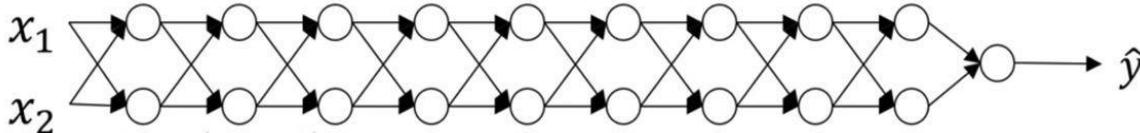


If features are on different scales, it might slow down gradient descent.

Hence, we should always normalize

Vanishing / Exploding Gradients

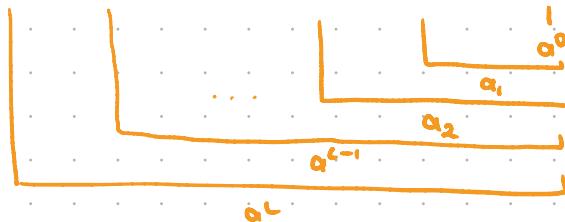
when we train very-deep networks, our gradients may be really small or really large.



Let's run a quick example.

- use linear function $g(z) = z$
- not use bias $b^{[l]} = 0$

$$\text{Then, } \hat{y} = W^{[L]} W^{[L-1]} \dots W^{[2]} W^{[1]} x$$



If we had

$$W^{[l]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$$

Then, activations would get exponentially larger and larger.

If we used

$$W^{[l]} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$$

They would get smaller and smaller.

If the weights are little bit greater than identity, in a very deep network, activations (and gradients) will grow exponentially.

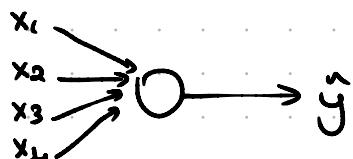
Makes training difficult

Weight Initialization Helps

Weight Initialization

why: careful choice for how weights are initialized can help with gradient vanishing/exploding problem.

Single Neuron Example



NOTE THAT

this is a
good goal

$$\text{Var}(w_i) = \frac{1}{n}$$

to prevent blowup,
we want w_i 's to be small

\downarrow
large $n \rightarrow$ smaller w_i

$$W^{[c]} = \text{np.random.randn(shape)} * \text{np.sqrt}\left(\frac{1}{n^{[c-1]}}\right)$$

if we especially use ReLU function,

$$\text{Var}(w_i) = 2/n \text{ works better}$$

number of features
in layer c

$$W^{[c]} = \text{np.random.randn(shape)} * \text{np.sqrt}\left(\frac{2}{n^{[c-1]}}\right)$$

if we use tanh function

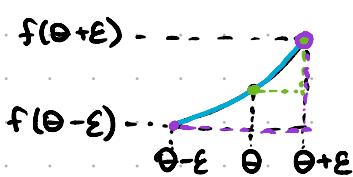
$$\text{Var}(w_i) = 1/n^{[c-1]}$$

Some other authors use

$$\text{Var}(w_i) = 2 / (n^{[c]} + n^{[c-1]})$$

Numerical Gradient Approximation

$$f(\theta) = \theta^3$$



best way to approximate is take two triangles

1 Triangle

$$\frac{f(\theta + \epsilon) - f(\theta)}{\epsilon}$$

2 Triangles

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

Approximate Err:
0.03

Error $\in O(\epsilon)$

(at $\theta = 1$)

Approx. error
0.0001

Error $\in O(\epsilon^2)$

IDEA
of writing
an article
is better

in two triangle method

$$\text{Error} \in O(\varepsilon^2)$$

in one triangle method

v.s.

$$\text{Error} \in O(\varepsilon)$$

Gradient Checking (super useful for checking backprop implementation)

1) Take $W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}$ and reshape into a big vector Θ

$$\text{then } \Rightarrow J(W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}) = J(\Theta)$$

2) Take $dW^{[1]}, db^{[1]}, \dots, dW^{[L]}, db^{[L]}$ and reshape into a big vector $d\Theta$

Remember, J is now a function of Θ .

Grad Check Implementation:

for each i :

$$d\theta_{\text{approx}}[i] = \frac{J(\theta_1, \theta_2, \dots, \theta_i + \varepsilon, \dots, \theta_n) - J(\theta_1, \theta_2, \dots, \theta_i - \varepsilon, \dots, \theta_n)}{2\varepsilon}$$

$d\theta_{\text{approx}}$ will be same dimensions as $d\theta$. We want to check if these two vectors are approximately equal

Check

$$\frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2}$$

$$\varepsilon = 10^{-7}$$

↓
What Does Values Mean?
 $\approx 10^{-7}$ great!
 10^{-5}
 10^{-3} worry...

notice that there is no square at top, so this is sum of squared of elements, then square root of the sum

HOW TO FIX

Check the vectors to see if there is a specific value of i such that

$$d\theta_{\text{approx}} \neq d\theta$$

then debug + run test again

Gradient Checking Implementation Notes

- grad check is costly. don't use in training. only use to debug
- check which value of $d\theta$ is different
- remember regularization (check that $d\theta$ includes it)
- doesn't work with dropout $\rightarrow J$ is different

↳ workaround:

implement grad check before dropout

implement dropout after checking backprop correct.

