

Bias and Variance

Train set error:	1%	15%	13%	0-5%
Dev set error:	11%	16%	30%	1%
	high variance	high bias	high bias	low bias
			high variance	low variance

Optimal (Bayes) error: 15%

Variance \rightarrow overfitting

bias \rightarrow underfitting

Basic recipe for machine learning

High bias $\xrightarrow{\lambda}$ Bigger network

$\downarrow N$

High variance \rightarrow More data
Regularisation

Bias variance trade off

\uparrow \downarrow
 \downarrow \uparrow

Bigger network is

\rightarrow Deep learning to be low variance.

Regularization \rightarrow Overfitting's engeller?

$\lambda \Rightarrow$ (du kagile durs over) smooth high bias

$$\min_{w, b} J(w, b) \quad w \in \mathbb{R}^{n \times x}, b \in \mathbb{R}$$

$\lambda =$ regularization parameter

backward is in

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

$$L_2 \text{ regularization} \quad \|w\|_2^2 = \sum_{j=1}^{n \times x} \hat{w}_j^2 = w^T w$$

$$L_1 \text{ regularization} \quad \frac{\lambda}{2m} \sum_{j=1}^{n \times x} |w_j| = \frac{\lambda}{2m} \|w\|_1$$

\rightarrow could add more weights
which regularize the weights

W

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

$$\|w^{[l]}\|_F^2 = \sum_{i=1}^{n^{[l-1]}} \sum_{j=1}^{n^{[l]}} (w_{ij}^{[l]})^2$$

"Frobenius norm"

$$W = (w^{[1]} \dots w^{[L-1]})$$

$$dw^{[l]} = (\text{from backward prop}) + \frac{\lambda}{m} w^{[l]}$$

$$w^{[l]} := w^{[l]} - \mathcal{L} dw^{[l]}$$

$$\frac{dJ}{dw^{[l]}} = dw^{[l]}$$

"weight decay"

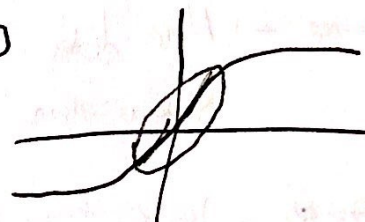
$$w^{[l]} := w^{[l]} - \mathcal{L} \left[(\text{from backward}) + \frac{\lambda}{m} w^{[l]} \right]$$

$$\left(1 - \frac{2\lambda}{m}\right) w^{[l]} = w^{[l]} - \frac{2\lambda}{m} w^{[l]} - \mathcal{L} (\text{from backward prop})$$

How does regularization prevent overfitting?

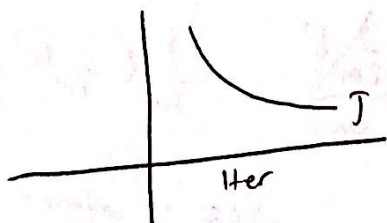
$$\lambda \uparrow \quad w^{[l]} \downarrow \quad \downarrow \quad \underline{z}^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$$

Every layer is linear



2 hidden layers with linear & yolostr.

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



Dropout regularization

Implementing dropout ("Inverted dropout")

Illustrate with layer $l=3$

keep-prob = 0.8

16 Feb
0.8 around 1
0.2 around 0 den

0.2 change eliminated

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

$$u_3 = \text{np.multiply}(u_3, d_3) \quad \# \text{ } a_3 \text{ } \leftarrow d_3$$

↓ elementwise

$$a_3 / = \text{keep-prob}$$

↳ a_3 in shape i layer
%0.8 den
0.2 den olan bir
matrisi aynı d_3 ile
çarpılır.

50 units \rightarrow 10 units shut off

$$z^{(4)} = w^{(4)} \cdot a^{(3)} + b^{(4)}$$

↑ reduced by 20%
 $l = 0.8$

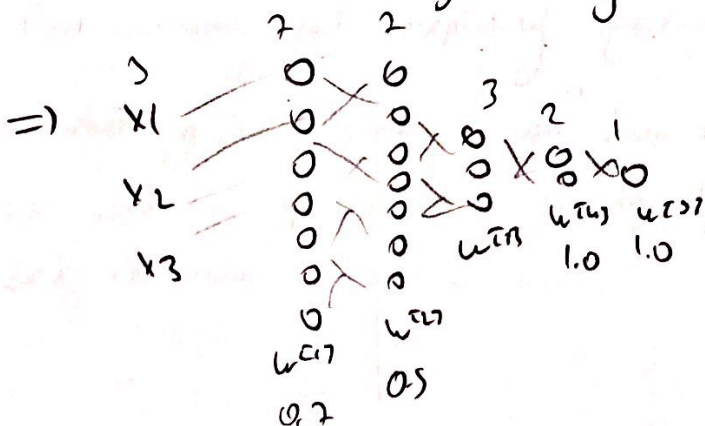
- Regularization technique
- only use dropout during training.
- Applied both forward and backward

Making predictions at test time

No drop out \rightarrow den daha random olmasın istemiyorsanız.

Why does dropout work?

Intuition: Don't rely on any one feature, so have to spread out weights



$$W^{(1)} = (3, 7)$$

$$W^{(2)} = (7, 7) \rightarrow w^2 \text{ çarpıldıkça büyük oluyor}$$

$$W^{(3)} = (7, 7)$$

$$W^{(4)} = (3, 2)$$

$$W^{(5)} = (2, 1)$$

run keep prob daha
daha olabilir 0.5

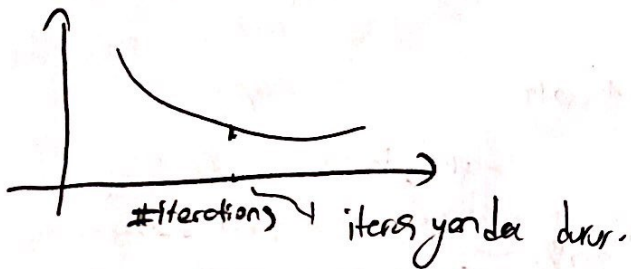
diğerleri run 0.7

Azırı öğrenme konusunda şüpheler
yeter 1.0 den daha

→ L2 and drop-out dan beska regularization teknik
ata augmentation → Bm cost regularization

- Resmilerin ters carihner
- Postgele gerintiler kuyru
- Zoomlar
- Distorsiyon (biskelme)

2) Early stopping



Normalizing training sets \rightarrow better cost function optimize

17 Feb

Subtract mean:

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

$$x := x - \mu$$

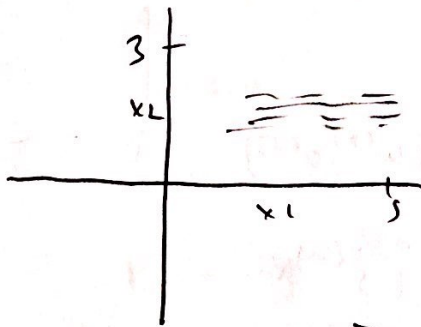
Normalize variance

$$\sigma = \frac{1}{n} \sum_{i=1}^n x^{(i)2} - 2$$

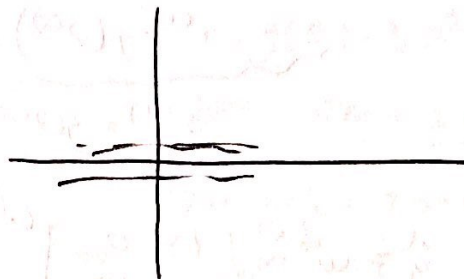
elementwise

$$x := x / \sigma$$

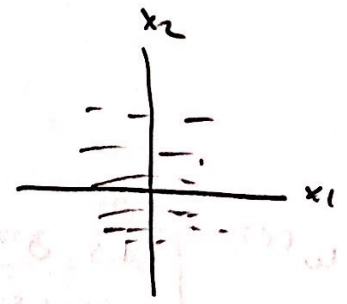
Use same μ σ^2 to normalize test set



Subtract mean



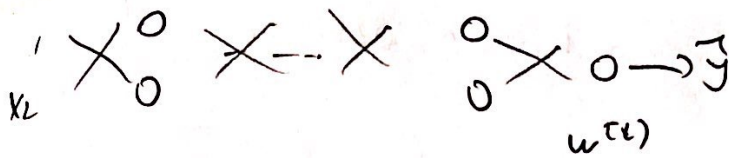
Normalize variance



If your input features came from very different scales, maybe some features are from 0 to 1, some from 1, to 1,000, then it's important to normalize your features.

Data vanishing and exploding gradients

\Rightarrow When you're training a very deep network your derivatives or your slopes can sometimes get either very, very big or small, or even exponentially small, and this makes training difficult.



For sake of simplicity,

$$g(z) = z \quad b^{(L)} = 0$$

1 layer

$$\hat{y} = w^{(L)} \cdot w^{(L-1)} \dots w^{(2)} w^{(1)} x$$

→ B. deeper exp. older
other we by training algorithm
solution

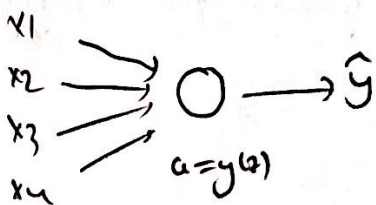
$$a^{(1)} = y(z^{(1)}) = z^{(1)}$$

$$a^{(2)} = y(z^{(2)}) = y(w^{(2)} \cdot a^{(1)})$$

$$w^{(2)} = \begin{bmatrix} 0.5 & 1.5 & 0 & 0.5 \\ 0 & 1.5 & 0 & 0.5 \end{bmatrix}$$

$$\hat{y} = w^{(2)} \begin{bmatrix} 1.5 & 0.5 \\ 0 & 1.5 \end{bmatrix}^{L-1} \cdot x \rightarrow 1.5^{L-1} x$$

Weight Initialization for Deep Network → vanishing or exploding gradients in weight learning



$$z = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

large n → smaller w_i

$$\text{Var}(w_1) = \frac{1}{n} \cdot \frac{2}{n}$$

$$w^{(L)} = \text{np.random.randn}(\text{shape}) \cdot \text{np.sqrt}\left(\frac{2}{n^{(L-1)}}\right)$$

using Relu

$$g^{(L)}(z) = \text{Relu}(z)$$

other variati,

$$\tanh \sqrt{\frac{1}{n^{(L-1)}}} \quad \text{Xavier initialization}$$

$$\sqrt{\frac{2}{n^{(L-1)} + n^{(L)}}}$$

He initialization

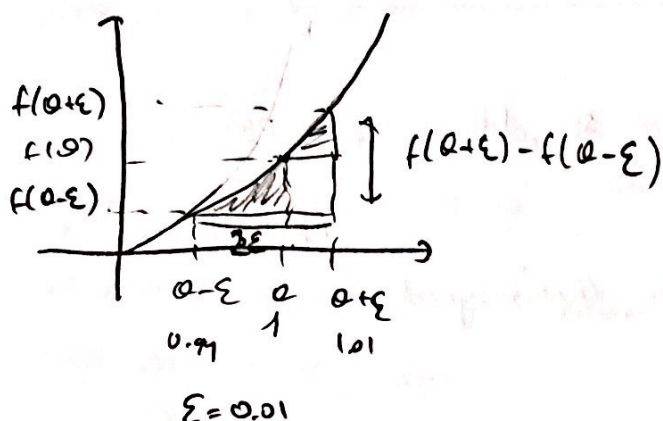
$$\text{np.random.randn}(\text{layer_dim}^{(L)}, \text{layer_dim}^{(L-1)})$$

$$\leftarrow \text{np.sqrt}(2. / \text{layer_dim}^{(L-1)})$$

Numerical approximation of gradients

Checking your derivative computation

$$f(x) = x^3$$



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

$$\frac{(1.01)^3 - (0.99)^3}{2 \cdot (0.01)} = 3.001$$

$$g(x) = 3x^2 = 3$$

$$\text{approx error} = 0.0001$$

$$(\text{per slide} = 3.0001 \text{ error } 0.0001)$$

check your work

$$\frac{f(x+\epsilon) - f(x)}{\epsilon} = \frac{(1.01)^3 - 1}{0.01} = 3.0301$$

$$f'(x) = \lim_{\epsilon \rightarrow 0} \frac{f(x+\epsilon) - f(x-\epsilon)}{2\epsilon} = \frac{0(\epsilon^2)}{0.01} = \frac{0.0001}{0.01} = 0.01$$

Gradient checking

Take $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}$ and reshape into a big vector ϕ

concatenate = ϕ de brillon

$$J(w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}) = J(\phi)$$

Take $dw^{(1)}, db^{(1)}, dw^{(2)}, db^{(2)}$ and reshape into a big vector $d\phi$

concatenate

for checking (Grad check) $J(\theta) = J(w_1, w_2, \dots)$
 ↪ joint parameter

for each i :

$$d\theta_{\text{approx}}[i] = \frac{J(w_1, w_2, \dots, w_i + \epsilon, \dots) - J(w_1, w_2, \dots, w_i - \epsilon)}{2\epsilon}$$

$$\approx d\theta[i] = \frac{\partial J}{\partial w_i} \quad \Bigg| \quad d\theta_{\text{app}} \approx d\theta \quad \parallel \quad w_2$$

check $\frac{\|d\theta_{\text{app}} - d\theta\|_2}{\|d\theta_{\text{app}}\|_2 \|d\theta\|_2} \approx 10^{-7}$ or smaller that's great
 NP: linear approx (...)
 $\epsilon = 10^{-7}$
 $10^{-5} \rightarrow$ worry

for find lots of bugs in implementation of neural nets,

- grad check estimate better than 1%
- Select how often to debug
- $d\theta_{\text{app}}$ very far from $d\theta$.
- Remember regularization

def grad-check(x, theta, epsilon):
 $\theta^+ = \theta + \epsilon$
 $\theta^- = \theta - \epsilon$
 $J^+ = \text{forward_prob}(x, \theta^+)$
 $J^- = \text{forward_prob}(x, \theta^-)$
 $\text{grad_app} = (J^+ - J^-) / (2 * \epsilon)$ check
 $\text{grad_true} = \text{backward_prob}(x, \theta)$

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

bias estimator estimate

- grad check dropout it's better to disable
- dropout during J and J grad for debugging
- keep-prop=1.0 'a cyclical cycle control oddball'

1) high variance versus high bias versus
bigger network

- Get more training data
- Add regularization
- Get more test data

2) training set error - 0.5%
dev set error 2% Overfitting observed

- increase the regularization parameter lambda
- Get more training data

3) Weight decay

A regularization technique that results in gradient descent shrinking the weights on every iteration. (L2 regn such as) adam, adam

4) \uparrow λ regularization parameter weights close to 0

5) Useful for reducing variance

- Dropout
- L2 regularization
- Data augmentation

6) Why do we normalize the input x ?

It makes the cost function faster to optimize

Model's gradient weight's \odot gradient fails to break symmetry due

random gradient to large weights due

randomly initialize network He

Xavier glb

He $\rightarrow \text{np.random.randn}(L-1, L-1) \cdot \sqrt{2/(L-1)}$

Xavier de 2 yaline 1

\rightarrow Regularization will help you reduce overfitting

\rightarrow will drive your weights to lower values

\rightarrow L2 regularization and Dropout are two very effective reg. tech.