

_ff76ae8a7ee526eb18cc9b3acd0cbc3f_C2W1L03



deeplearning.ai

Regularizing your neural network

Regularization

prevent over-fitting
 reduce variation

Logistic regression

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

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

λ = regularization parameter
 lambda lambda

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

L_2 regularization $\|w\|_2^2 = \sum_{j=1}^n w_j^2 = w^T w \leftarrow$



w will be sparse

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

if not for compress model won't be very helpful.

more or less only a normalizing const

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

\rightarrow why doing this? w high dim
 b one param. won't make much difference. Don't bother to include it.

Andrew Ng

Neural network

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

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

"Frobenius norm"

$$\| \cdot \|_2 \quad \| \cdot \|_F$$

$$dw^{(l)} = (\text{from backprop}) + \frac{\lambda}{n} w^{(l)}$$

$$\rightarrow w^{(l)} := w^{(l)} - \alpha \frac{\partial J}{\partial w^{(l)}}$$

$$\frac{\partial J}{\partial w^{(l)}} = dw^{(l)}$$

back propagation

L2 regularization = weight decay

"Weight decay"

$$\begin{aligned} w^{(l)} &:= w^{(l)} - \alpha \left[(\text{from backprop}) + \frac{\lambda}{n} w^{(l)} \right] \\ &= w^{(l)} - \frac{\alpha \lambda}{n} w^{(l)} - \alpha (\text{from backprop}) \\ &= \underbrace{\left(1 - \frac{\alpha \lambda}{n}\right)}_{< 1} w^{(l)} - \alpha (\text{from backprop}) \end{aligned}$$

↳ make it smaller.

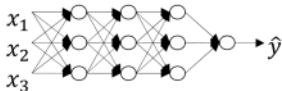
Andrew Ng

Neural network

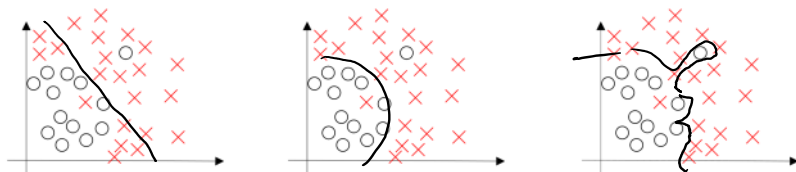
$$J(w^{(0)}, b^{(0)}, \dots, w^{(L)}, b^{(L)}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2$$

Andrew Ng

How does regularization prevent overfitting?



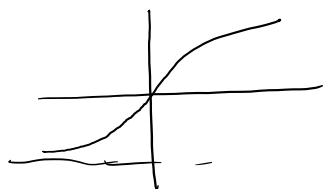
$$J(w^{(0)}, b^{(0)}, \dots, w^{(L)}, b^{(L)}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2$$



zero out lots of input of the hidden unit \Rightarrow
reduce impact of lots of hidden unit. $\Rightarrow w^{(l)} = 0$.
not an accurate description though.

Andrew Ng

How does regularization prevent overfitting?



$$\lambda \uparrow \Rightarrow w^{(2)} \downarrow \quad z^{(2)} = w^{(2)} a^{(1)} - b^{(2)}$$

Way layer is roughly linear
if regularized

⇒ simpler function

Andrew Ng

Implementation tip: add term to the loss function to
plot cost function — gradient descent.
⇒ should ↓.

★ with regularization remember to use the new def
of $J = \sum L_i + \lambda \|w\|^2$ remember to
include this part
in plotting.

dropout

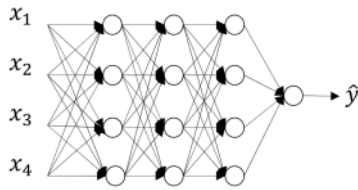


deeplearning.ai

Regularizing your
neural network

Dropout
regularization

Dropout regularization



↑ 0.5 ↑ 0.5 ↑ 0.5

Andrew Ng

Implementing dropout ("Inverted dropout")

Illustrate with layer $l=3$. keep-prob = $\frac{0.8}{0.2}$

activation $\rightarrow d3 = \text{np.random.rand}(a3.\text{shape}[0], a3.\text{shape}[1]) < \text{keep-prob}$ If $1 < 0.8 \rightarrow 0.8 \Rightarrow$ prob 0.2 drop.

$\rightarrow a3 = \text{np.multiply}(a3, d3)$ # $a3 \times d3$. true/false or 1/0.

★ $\rightarrow a3 = \frac{a3}{\text{keep-prob}}$ scale it up. Divide by keep-prob.

50 units. \rightarrow 10 units shut off example

$z^{(4)} = w^{(4)} \cdot a^{(3)} + b^{(4)}$ reduced by 20%. $1 = 0.8$ Test

so it won't change the expected value of $z^{(4)}$

Andrew Ng

layer $l=3$ -
keep-prob. = 0.8.
 np.random.rand
 \Rightarrow prob 0.2 drop.

At test time, the inverted dropout is easier \rightarrow has less scaling problem

Making predictions at test time

$a^{(0)} = X$

No drop out.

$z^{(1)} = w^{(1)} a^{(0)} + b^{(1)}$
 $a^{(1)} = g^{(1)}(z^{(1)})$
 $z^{(2)} = w^{(2)} a^{(1)} + b^{(2)}$
 $a^{(2)} = \dots$
 \downarrow
 \hat{y}

$1 = \text{keep-prob}$

Andrew Ng

test time NO DROP OUT

\rightarrow the fact that we do this in training makes sure that there's no scaling problem even if we don't implement dropout at test time

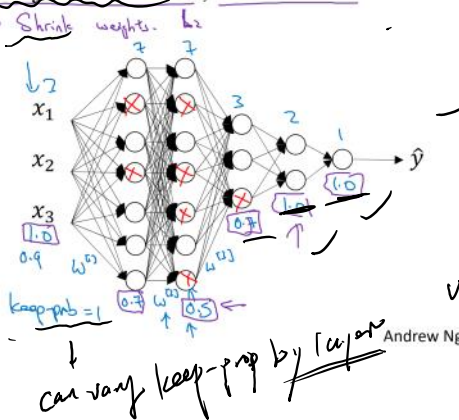
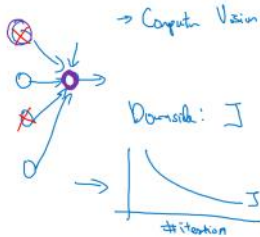


Regularizing your neural network

Understanding dropout

Why does drop-out work?

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



any feature could go away at random.
reluctant to put weight on any one particular input

for each node toss a coin.

with some probability. remove it

"Diminished network"

P_1, P_2, P_3 param. \Rightarrow keep-prob. $\frac{1}{2}$.

more worry about overfitting \Rightarrow lower loss prop.

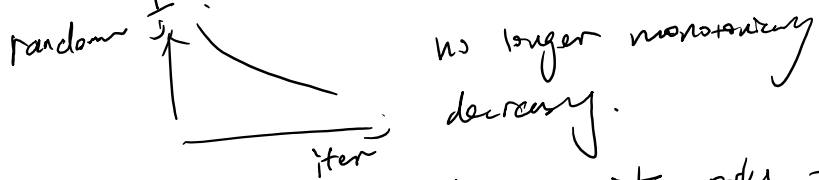
similar to playing with λ .

Occasionally, dropout at input feature but this is not common practice

* Many dropout has to do with CV cuz so many input.


Unless the application has over-fitting problem \Rightarrow Don't use dropout

Another problem of Dropout = Loss function J not well defined.



(What to do?) Do it without dropout, it works then

⇒ ~~that~~ add dropout.

 other regularization



deeplearning.ai

Regularizing your
neural network

Other regularization methods

Data augmentation



flipping images

transform
cheap way to
get more data.

4



4

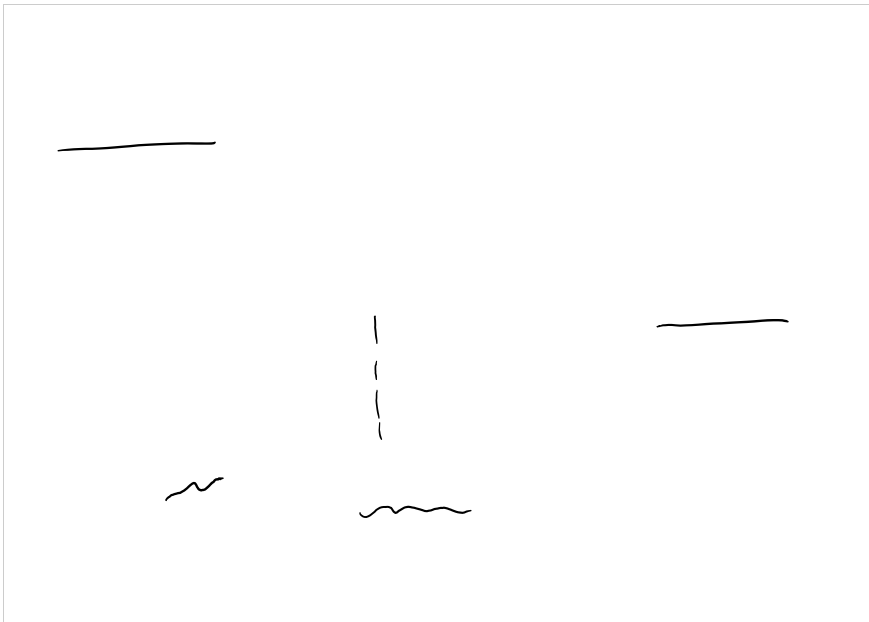
4

4



Andrew Ng

random rotation & distortion.



Also plot dev set error
find: Dev set error $\downarrow \Rightarrow$ 'T'
stop half way

Downside of early stopping [- Optimize $\Rightarrow J(w, b)$
- avoid overfitting.] a different tool for this tool

work on
two tasks
independently

Orthogonalize

\rightarrow 应用 \rightarrow 防止过拟合。

7.68.

\rightarrow alternative — L2 regularization

will talk
about it later.

\rightarrow Downside: computationally expensive
searching for λ .

7. With the inverted dropout technique, at test time:

☒ You do not apply dropout (do not randomly eliminate units), but keep the $1/\text{keep_prob}$ factor in the calculations used in training.

This should not be selected

☐ You do not apply dropout (do not randomly eliminate units) and do not keep the $1/\text{keep_prob}$ factor in the calculations used in training

☐ You apply dropout (randomly eliminating units) and do not keep the $1/\text{keep_prob}$ factor in the calculations used in training

☐ You apply dropout (randomly eliminating units) but keep the $1/\text{keep_prob}$ factor in the calculations used in training.

4. You are working on an automated check-out kiosk for a supermarket, and are building a classifier for apples, bananas and oranges. Suppose your classifier obtains a training set error of 0.5%, and a dev set error of 7%. Which of the following are promising things to try to improve your classifier? (Check all that apply.)

☒ Increase the regularization parameter λ

Correct

☐ Decrease the regularization parameter λ

Un-selected is correct

☐ Get more training data

This should be selected

☐ Use a bigger neural network

Un-selected is correct