

CHƯƠNG 3

TỐI ƯU HOÁ MÔ HÌNH MẠNG NEURAL (P1)

Khoa Khoa học và Kỹ thuật thông tin
Bộ môn Khoa học dữ liệu

NỘI DUNG

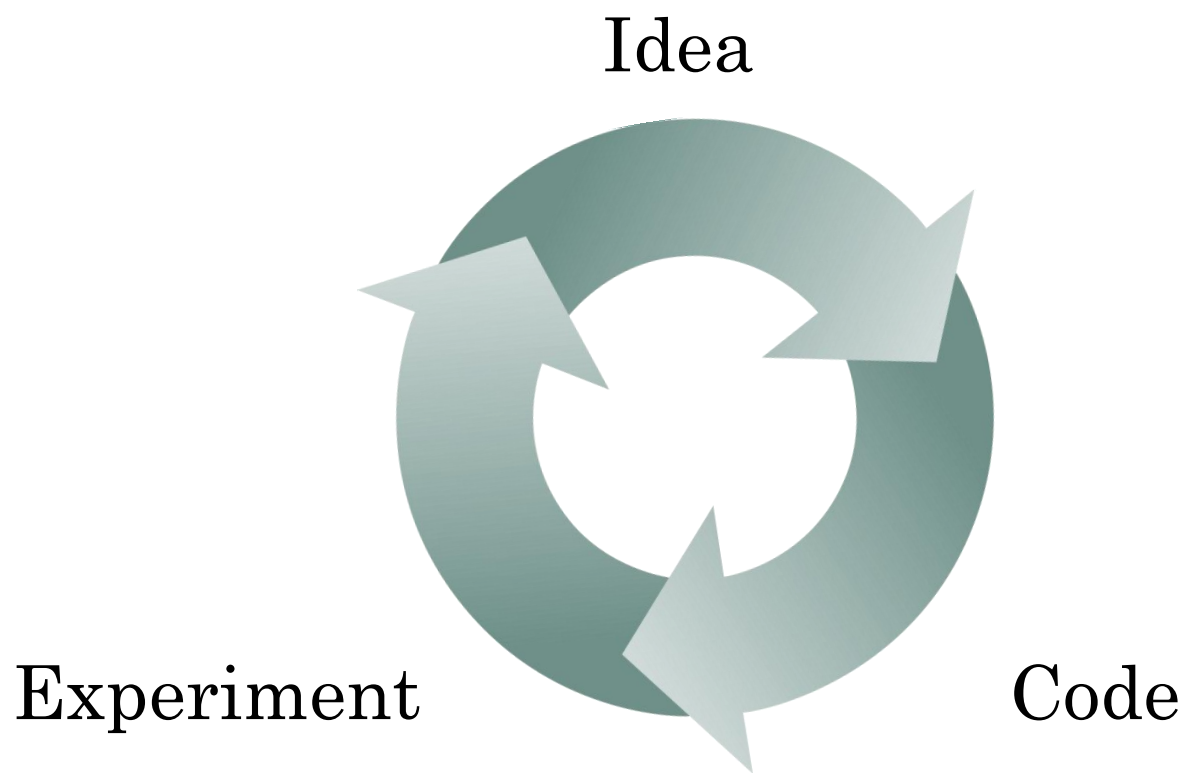
1. Phân chia dữ liệu huấn luyện.
2. Deep learning error.
3. Chuẩn hoá mạng neural (Regulazation)
4. Khởi tạo trọng số.

PHÂN CHIA DỮ LIỆU HUẤN LUYỆN

Các siêu tham số khác trong mô hình

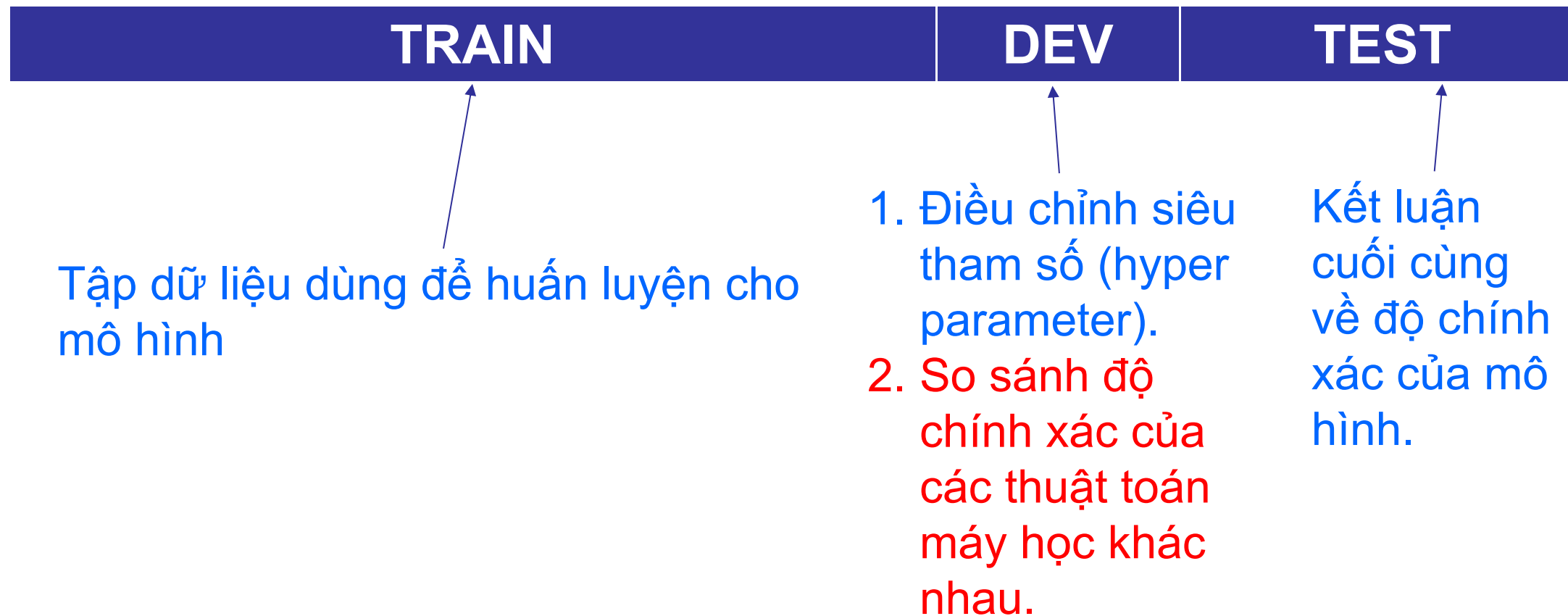
- Momentum.
- Epochs.
- Mini batch size.
- Regularization.
-

Điều chỉnh siêu tham số



- Quá trình áp dụng deep learning là 1 quá trình thực nghiệm.
 - + Điều chỉnh các siêu tham số (ví dụ như learning rate, iteration, etc) sao cho hàm chi phí đạt giá trị nhỏ nhất.
- Ngoài ra, việc phân chia dữ liệu thành các tập train/dev/test ảnh hưởng rất lớn đến hiệu quả của mô hình

Phân chia dữ liệu



Cách chia dữ liệu

	TRAIN	DEV	TEST
	%	%	%
Truyền thống (~10K)	60	20	20
Big data (1M)	98	1	1
Hơn 1M dữ liệu	99.5	0.25	0.25

Trường hợp dữ liệu train và dev/test không khớp

Train

- Hình chụp các con mèo từ các trang web.
- Bộ văn bản tiếng Việt chính quy (từ Wikipedia)

Dev/test

- Hình các con mèo tự chụp bằng camera.
- Bộ văn bản tiếng Việt trên mạng xã hội.

Dữ liệu từ tập train và tập dev/test phải từ cùng một nguồn, cùng một tính chất với nhau (có cách gọi khác là same distribution). Nếu không thì việc đánh giá mô hình mạng neural sẽ không thể chính xác được.

Tập test có cần thiết không

- Tập dev thường dùng để đánh giá và so sánh độ chính xác giữa các mô hình với nhau.
- Tập test dùng để kiểm tra lại việc đánh giá trên tập dev xem kết quả có chính xác hay không (đôi khi còn gọi là unbiased estimate).
 - + Nếu như kết quả giữa dev và test không chênh lệch nhiều thì có thể chấp nhận độ chính xác của mô hình, ngược lại thì sẽ xảy ra hiện tượng overfitting.
- Trên thực tế, có khi sẽ không nhất thiết phải dùng tập test set mà chỉ cần dev set là đủ → kiểm tra độ chính xác của mô hình.
 - + Nếu chỉ có train set và dev set mà không có test set thì tập dev set sẽ được gọi luôn thành test set.

DEEP LEARNING ERROR

Underfitting và Overfitting

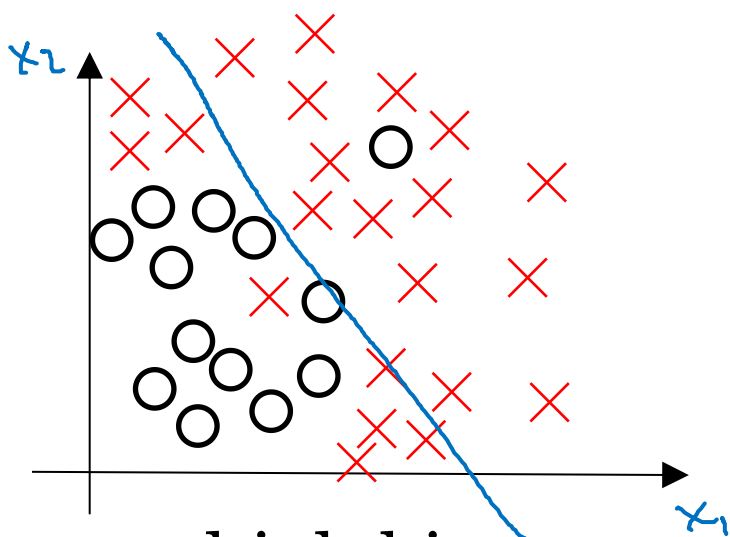
Underfitting

- Mô hình xây dựng được chưa tổng quát được trên tổng thể bộ dữ liệu huấn luyện → độ chính xác khi dự đoán của mô hình thấp.
- **High bias.**
- Độ chính xác trên tập train thấp.

Overfitting

- Mô hình xây dựng thể hiện quá chi tiết trên bộ dữ liệu huấn luyện → các dữ liệu nhiễu cũng được đưa vào → khi gặp các dạng nhiễu khác, mô hình sẽ không còn chính xác → trên bộ dev (hoặc test) độ chính xác khi dự đoán sẽ thấp.
- **High variance.**
- Độ chính xác trên tập train cao, nhưng tập dev/test thấp.

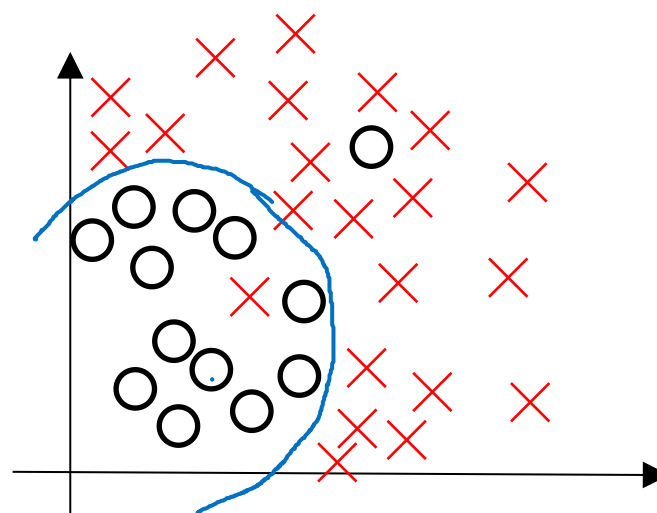
Bias and Variance – Deep learning error



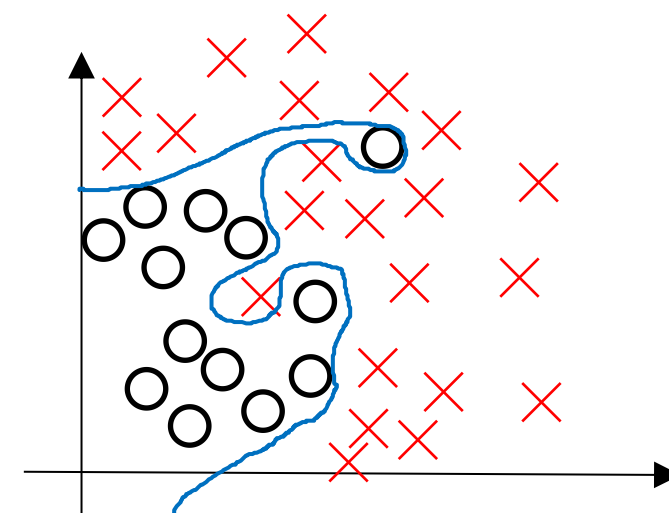
high bias

underfitting

Underfitting: Đường thẳng dự đoán không khớp với lại dữ liệu dự đoán → không tổng quát được trên dữ liệu.



→ “just right”



high variance

overfitting

Overfitting: Đường thẳng dự đoán quá khớp với lại dữ liệu dự đoán (kể cả nhiễu).

Ví dụ

Cat classification



Train set error:

1%

15%

15%

1%

Dev set error:

11%

16%

30%

0.5%

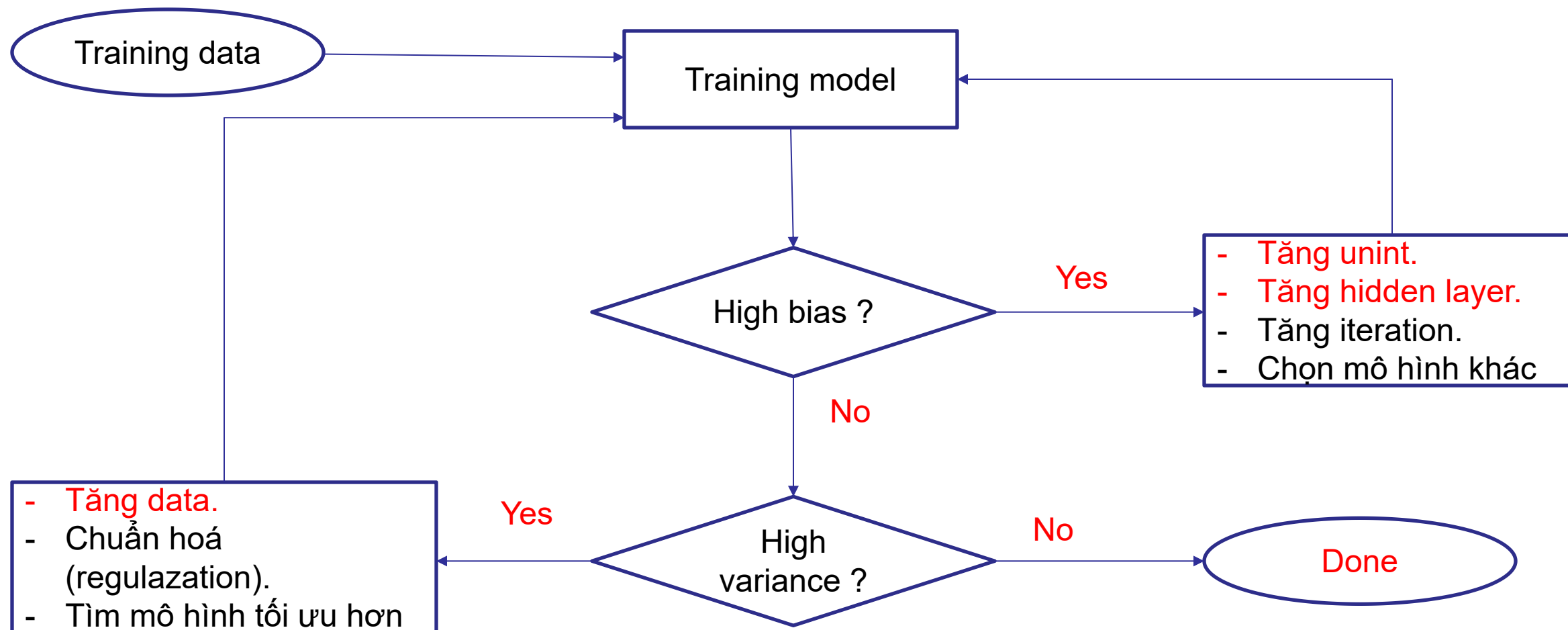
high variance

high bias

high bias
+ high variance

low bias
+ low variance

Quy trình tối ưu cho mô hình



Bias and Variance trade off

- Đối với máy học truyền thống, Bias và Variance tỉ lệ nghịch với nhau:
 - + Tăng bias sẽ làm giảm variance.
 - + Giảm bias sẽ làm tăng variance.
- Tuy nhiên, trong kỹ nguyên Big data và deep learning:
 - + Tăng kích thước mạng neural (số lớp ẩn, số unit) sẽ làm *giảm bias mà không ảnh hưởng nhiều tới variance*.
 - + Tăng dữ liệu sẽ làm giảm variance mà *không ảnh hưởng nhiều tới bias*.

CHUẨN HOÁ MẠNG NEURAL – GIẢM DEEP LEARNING ERROR: OVERFITTING

Regularzation

- Để hạn chế việc **overfitting** (high variance) trên mô hình, ta có các kỹ thuật sau:
 - + Tăng cường data.
 - + Chuẩn hoá mô hình mạng neural (regularzation).
- **Regularzation** là một kỹ thuật điều chỉnh tham số cho mạng neural nhằm giải quyết quá trình overfitting.

Nhắc lại về Logistic regression

— Hàm chi phí:

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

Trong đó:

$$W \in \mathbb{R}^{n \times x}, b \in \mathbb{R}.$$

Chuẩn hoá L2 cho Logistic regression

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

L2 Regularization
Ridge Regression.

$$\|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w$$

b là một số đơn lẻ (single number) → xét về kích thước, b nhỏ hơn rất nhiều so với w – vốn chứa rất nhiều tham số. Do đó, việc chuẩn hoá b không ảnh hưởng đáng kể so với chuẩn hóa w → có thể bỏ đi.

Chuẩn hoá L1 cho Logistic regression

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \boxed{\frac{\lambda}{2m} * \|w\|_1}$$

L1 Regularization $\|w\|_1 = \sum_{j=1}^{n_x} |w_j|$

Lasso Regression

Trong trường hợp này, **w** là một ma trận thưa (**sparse matrix**) – là một ma trận mà hầu hết các phần tử bằng 0. Ngược lại với ma trận thưa là ma trận dày đặc (**dense matrix**), nơi mà hầu hết các phần tử khác 0.

MỘT SỐ LƯU Ý

- L1 regularization sẽ khiến cho tham số của mô hình bị thưa, do đó trong thực tế, người ta thường sử dụng chuẩn hoá L2 regularization.
- λ được gọi là tham số chuẩn hoá (regularization parameter).
 - + Cài đặt thông số cho λ sử dụng tập dev (development test) hoặc là sử dụng cross validation.
 - + λ cũng là một loại siêu tham số (hyper parameter) trong mô hình.

Regularization trong mạng neural

– Hàm chi phí của mô hình mạng neural L lớp:

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

– L2 Regularization:

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

Frobenius norm



Frobenius norm

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

Chiều của w : $(n^{[l]}, n^{[l-1]})$

REGULARIZATION TRONG MẠNG NEURAL

Gradient Descent:

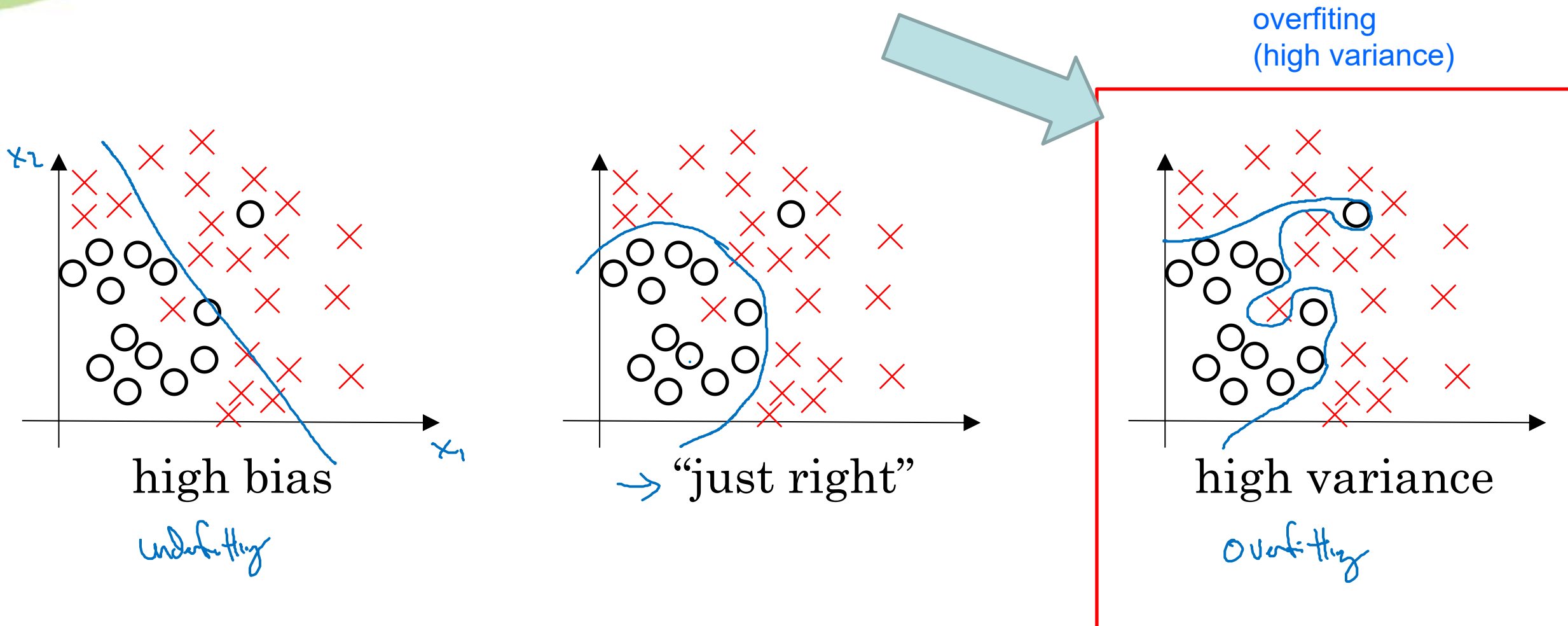
$$W^{[l]} = W^{[l]} - \alpha * dW^{[l]} \quad \text{Với: } dW^{[l]} = \langle back_propagation \rangle + \frac{\lambda}{m} w^{[l]}$$

$$= W^{[l]} - \boxed{\alpha * \frac{\lambda}{m} w^{[l]}} - \alpha * \langle back_propagation \rangle$$

Mỗi lần cập nhật trọng số, thì giá trị trọng số W sẽ **nhỏ đi một lượng** $\alpha * \frac{\lambda}{m} w^{[l]}$
→ **weigh decay**

TẠI SAO REGULARIZATION GIẢM ĐƯỢC OVERFITTING

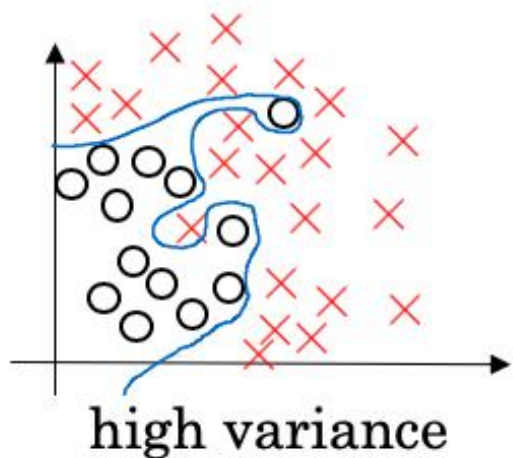
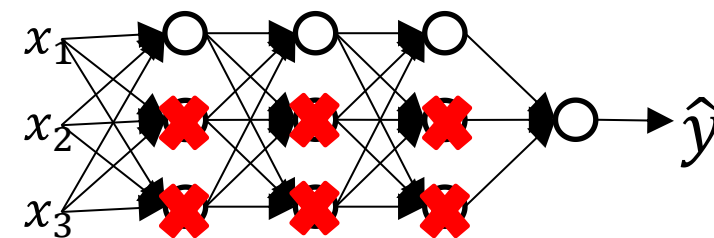
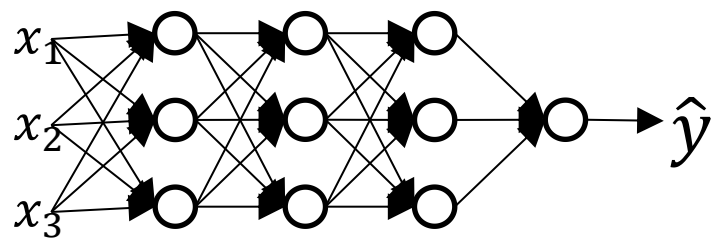
Nhắc lại về các dạng deep learning error



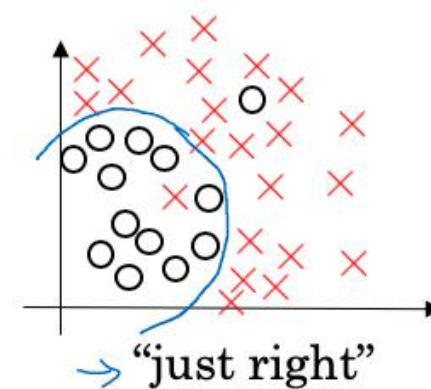
Giải thích 1: Trên hàm chi phí

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

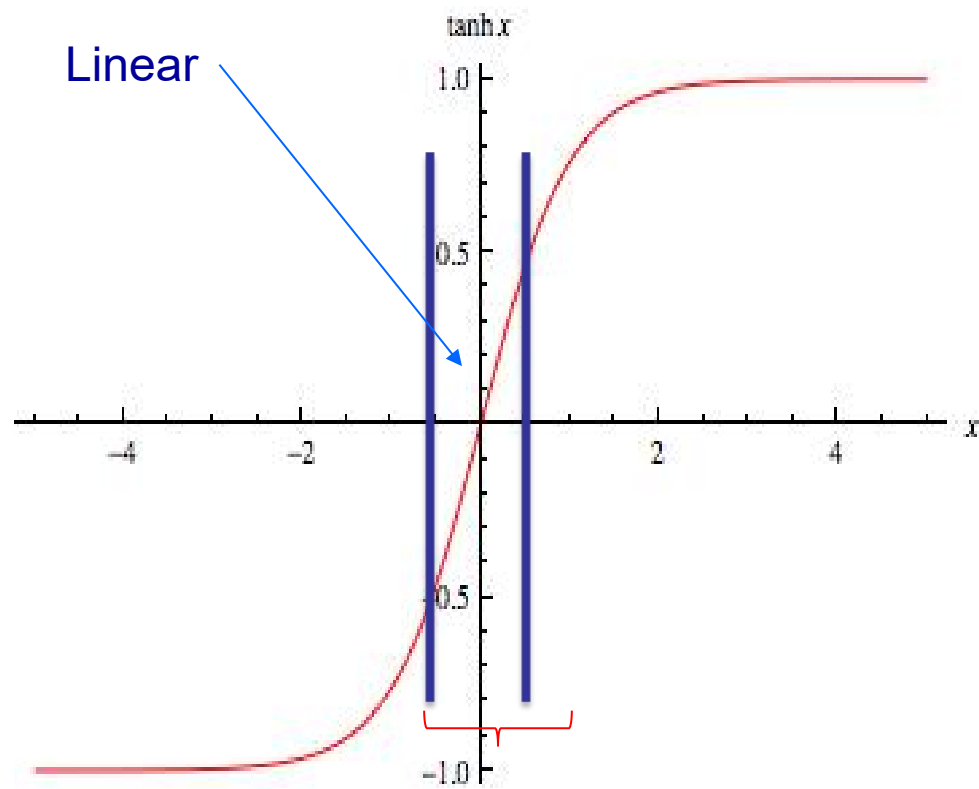
$$J(W^{[l]}, b^{[l]}) = \frac{1}{m} * L(\hat{y}, y) + \frac{\lambda}{2m} \sum_{l=0}^L ||w^{[l]}||_F^2$$



Khi có chuẩn hoá (regularization), thì trong quá trình tính Gradient descent, một số trọng số $W^{[l]}$ (ứng với các unit) sẽ giảm về gần bằng 0 ($W^{[l]} \approx 0$)
 → một số neural (unit) bị vô hiệu hoá.
 Khi đó, đường dự đoán (mô hình) sẽ “bớt khớp” với dữ liệu hơn



Giải thích 2: Trên hàm truyền



Z càng nhỏ \rightarrow hàm tanh có xu hướng quay về linear

Mỗi lần cập nhật trọng số, thì giá trị trọng số W sẽ nhỏ đi một lượng $\alpha * \frac{\lambda}{m} w^{[l]}$
 $\rightarrow \lambda$ càng lớn thì $w^{[l]}$ càng nhỏ.

$$Z^{[l]} = W^{[l]} * a^{[l-1]} + b^{[l]}$$

$$g(Z^{[l]}) = \tanh(Z^{[l]})$$

$\rightarrow \lambda$ càng lớn thì $w^{[l]}$ càng nhỏ

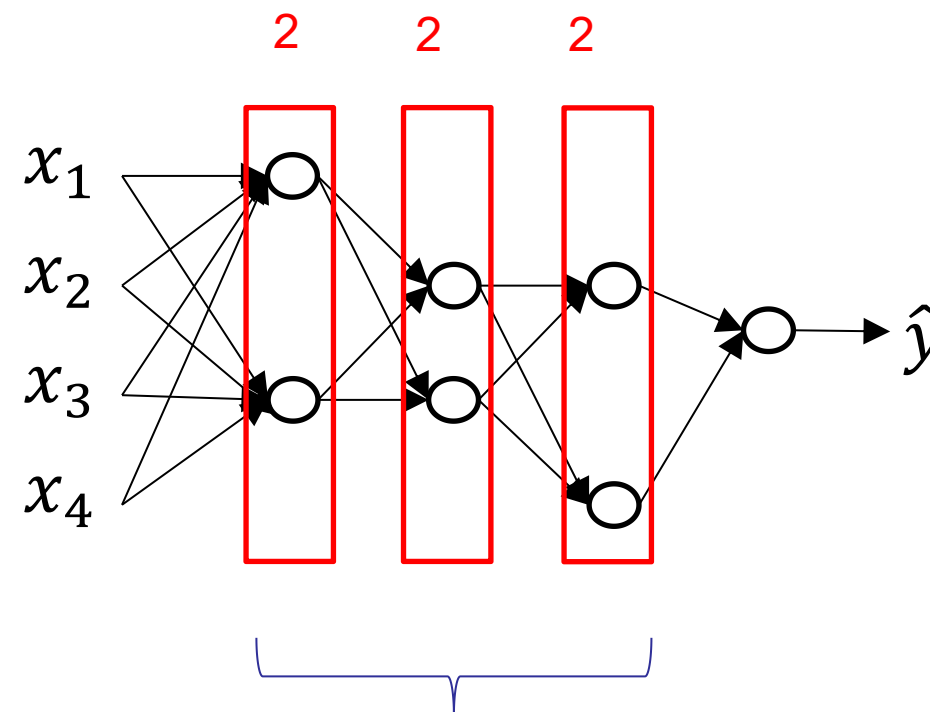
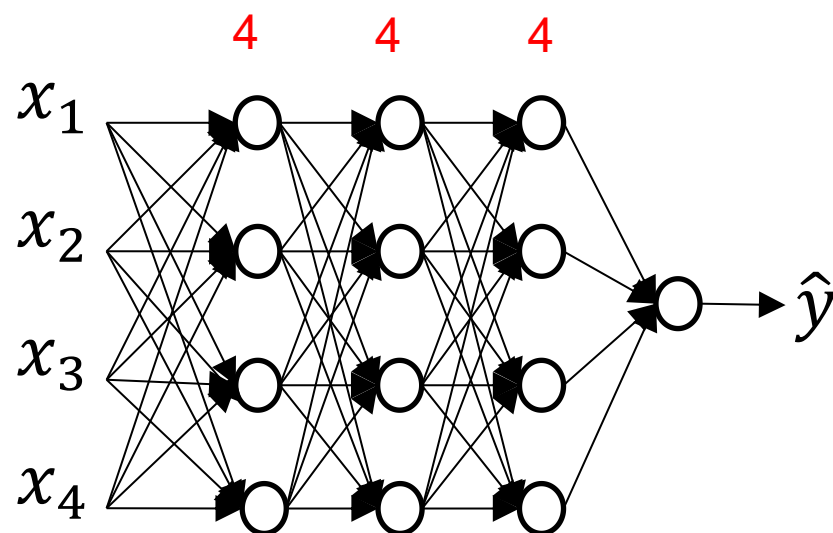
$\rightarrow Z^{[l]}$ càng nhỏ $\rightarrow g(Z^{[l]})$ càng nhỏ.

Dropout regularization

- Bên cạnh L2 regularization, dropout là một kỹ thuật regularization khác nhằm hạn chế overfitting.
- Dropout là kỹ thuật bỏ đi một số unit trong mỗi lớp nhằm đơn giản hoá mạng neural:
 - + Bỏ đi một số units => giảm $W^{[l]}$.
- Kỹ thuật này do Google phát triển và đã được cấp bằng sáng chế.

<https://jmlr.org/papers/v15/srivastava14a.html>

Minh họa



Bỏ đi khoảng 50% ($p=0.5$) cho
mỗi layer

Ví dụ

Giả sử, đang xét lớp thứ 3 trong mạng neural ($l=3$), xác suất giữ là 0.8 ($\text{keep_prob}=0.8$) \rightarrow khoảng 20% unit sẽ bị tắt.

```
d3 = np.random.randn(a3.shape[0], a3.shape[1]) < keep_prob
a3 = np.multiply(a3, d3)
a3 = a3 / keep_prob
```

Giả sử, a_3 có 50 unit, khi thực hiện bỏ khoảng 20% unit xong, a_3 còn lại khoảng 40 unit ($20\% \cdot 50 = 10$).

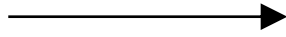
Như vậy, ở lớp tiếp theo:

$z^{[4]} = W^{[4]} \cdot a^{[3]} + b^{[4]}$, lớp $z^{[4]}$ sẽ chỉ còn khoảng 80% unit, do ta nhân với $a^{[3]}$ ($a^{[3]}$ đã dropout 20%)

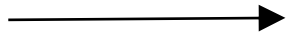
Một số kỹ thuật khác nhằm làm giảm overfitting

- Tăng dữ liệu (data augmentation).
 - + Chi phí đắt.
- Early stopping.

Ví dụ về Tăng dữ liệu ảnh



4



Ví dụ về Tăng dữ liệu TEXT

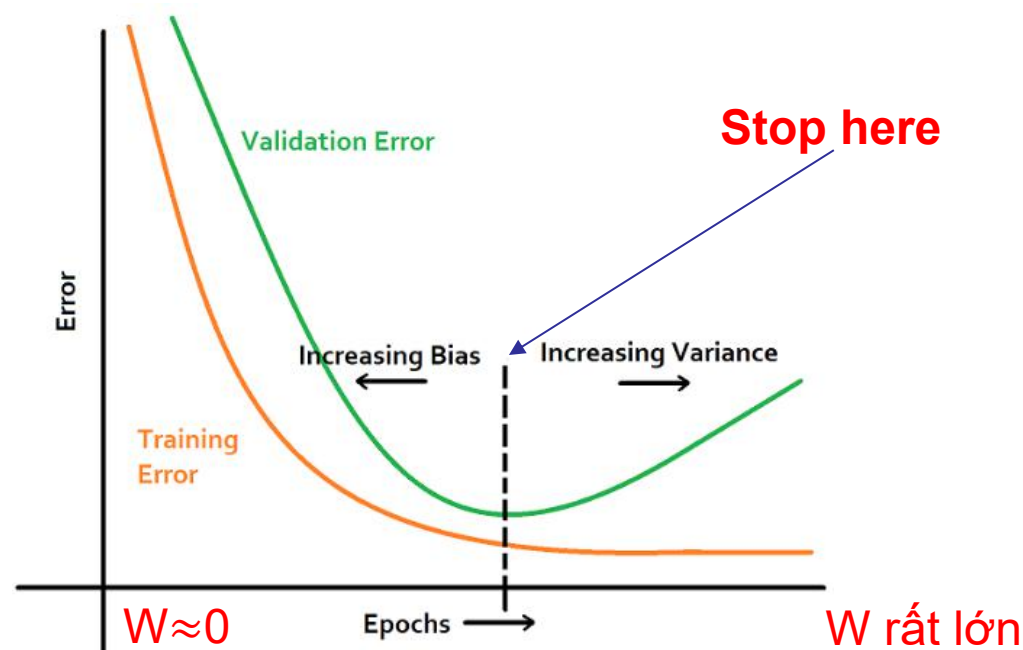
Operation	Sentence
None	A sad, superior human comedy played out on the back roads of life.
SR	A <i>lamentable</i> , superior human comedy played out on the <i>backward</i> road of life.
RI	A sad, superior human comedy played out on <i>funniness</i> the back roads of life.
RS	A sad, superior human comedy played out on <i>roads</i> back <i>the</i> of life.
RD	A sad, superior human out on the roads of life.

DỮ LIỆU CHÍNH

DỮ LIỆU TĂNG CƯỜNG

Early stopping

- Là một kỹ thuật dừng vòng lặp của gradient descent sớm để tránh trường hợp lặp quá mức dẫn đến overfitting.



KHỞI TẠO TRỌNG SỐ

Chuẩn hoá dữ liệu

— Chuẩn hoá dữ liệu huấn luyện dựa vào

+ Trung bình (mean): $\mu = \frac{1}{m} \sum_{i=1}^m X^{(i)}$

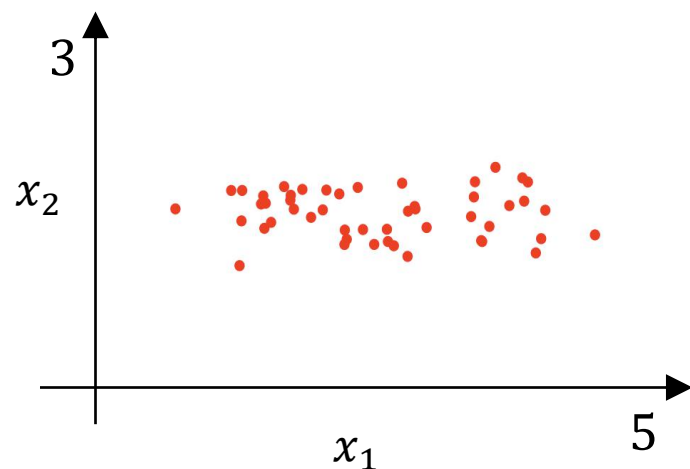
+ Phương sai (variance): $\sigma = \frac{1}{m} \sum_{i=1}^m X^{(i)} \otimes \otimes 2$

— Lưu ý: Khi chuẩn hoá tập train, thì tập dev/test phải được thực hiện tương tự (cùng μ và σ).

Chuẩn hoá dữ liệu (tt)

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

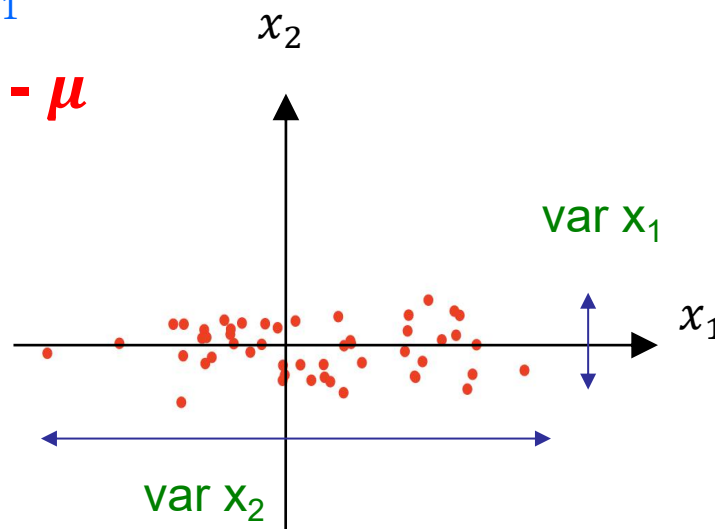
$$X = X - \mu$$



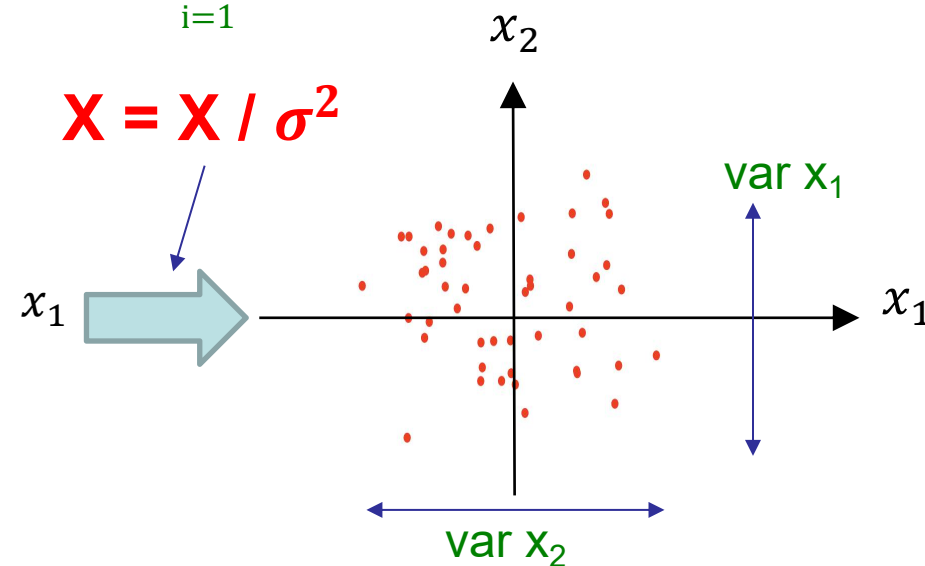
Dữ liệu ban đầu:
 $X = [x_1, x_2]$

$$\sigma = \frac{1}{m} \sum_{i=1}^m X^{(i)} \otimes \otimes 2$$

$$X = X / \sigma^2$$



$\text{var } x_1 < \text{var } x_2$



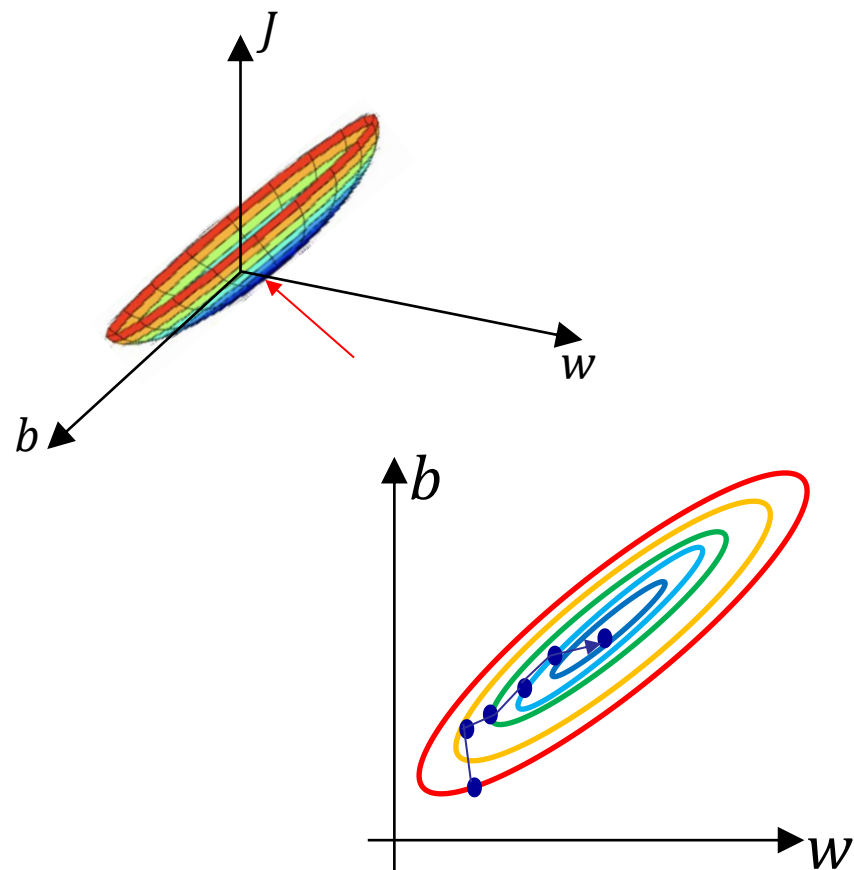
$\text{var } x_1 \approx \text{var } x_2$

Tại sao cần chuẩn hoá dữ liệu ?

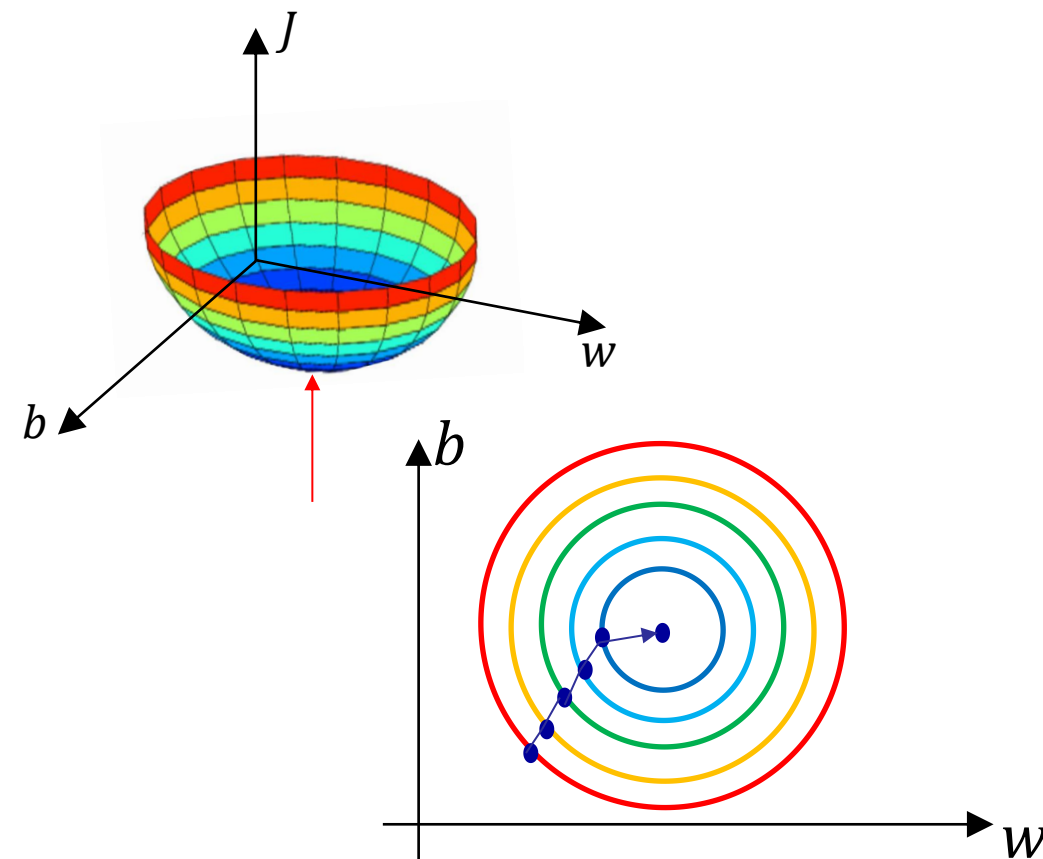
- Với dữ liệu đầu vào: $X = [x_1, x_2]$.
 - + x_1 có miền giá trị 1 đến 1000.
 - + x_2 có miền giá trị 1 đến 10.
 - giá trị w_1 và w_2 rất khác nhau.
- Nếu như không chuẩn hoá, thì giá trị của w_1 và w_2 sẽ có sự chênh lệch lớn → quá trình tính toán trên hàm chi phí J chậm.
- Như vậy, chuẩn hoá dữ liệu sẽ giúp cho quá trình gradient descent được cải thiện.

Minh họa

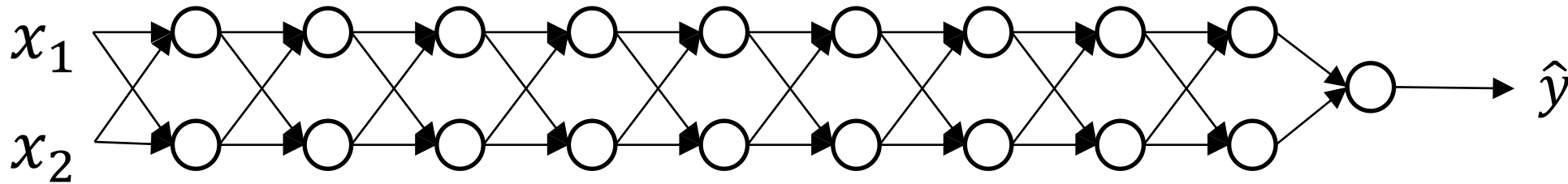
Unnormalized



Normalized



Vanishing/Exploding Gradient descent



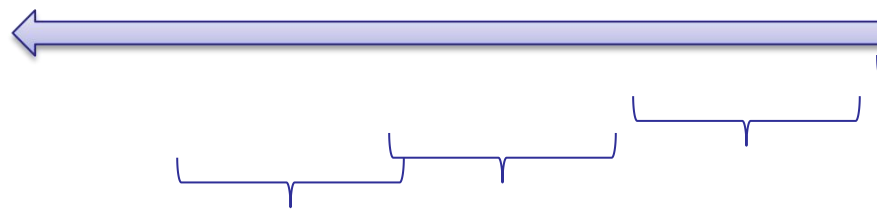
$W[1] \quad W[2] \quad W[3] \quad W[4] \quad \dots \quad W[L]$

$$\hat{y} = W[L] W[L-1] W[L-2] \dots W[3] W[2] W[1] X$$

$$Z^{[1]} = W^{[1]} X \mid a^{[1]} = g(Z^{[1]})$$

$$Z^{[2]} = W^{[2]} * a^{[1]} \mid a^{[2]} = g(Z^{[2]})$$

$$W^{[1]} * W^{[2]} * \dots * W^{[L]} = W \text{ mũ } L$$



Nhận xét

- Càng đi sâu xuống nhiều lớp, thì giá trị hàm kích hoạt (activation function) sẽ thay đổi như sau:
 - + Nếu W càng lớn thì giá trị y càng lớn (exploding).
 - + Nếu W càng nhỏ thì giá trị y càng nhỏ (vanishing).
 - Nếu khởi tạo $W^{[L]} > I$, thì giá trị hàm kích hoạt càng lớn nếu càng đi về các lớp sau.
 - Nếu khởi tạo $W^{[L]} < I$, thì giá trị hàm kích hoạt càng nhỏ dần nếu càng đi về các lớp sau.
- (I là ma trận đơn vị)*

Ví dụ

Exploding

$$W^{[L]} = \begin{pmatrix} 1.5 & 1 \\ 1 & 1.5 \end{pmatrix}$$

Giả sử mạng neural 10 lớp.

$$W^{[10]} \approx \begin{pmatrix} 1.5^{10} & 1 \\ 1 & 1.5^{10} \end{pmatrix} \\ \approx \begin{pmatrix} 57.66 & 1 \\ 1 & 57.66 \end{pmatrix}$$

Vanishing

$$W^{[L]} = \begin{pmatrix} 0.5 & 1 \\ 1 & 0.5 \end{pmatrix}$$

Giả sử mạng neural 10 lớp.

$$W^{[10]} \approx \begin{pmatrix} 0.5^{10} & 1 \\ 1 & 0.5^{10} \end{pmatrix} \\ \approx \begin{pmatrix} 0.00097 & 1 \\ 1 & 0.00097 \end{pmatrix}$$

Khởi tạo trọng số cho mạng neural để tránh Vanishing/exploding GD

Khởi tạo trọng số $W^{[l]}$:

$$W^{[1]} = np.random.randn(shape(W)) \otimes \underbrace{Variance}$$

Variance cho ReLU: $np.sqrt\left(\frac{2}{n^{[l-1]}}\right)$

Làm cho trọng số W quá lớn hoặc quá nhỏ so với 1 (phương sai không quá lớn)

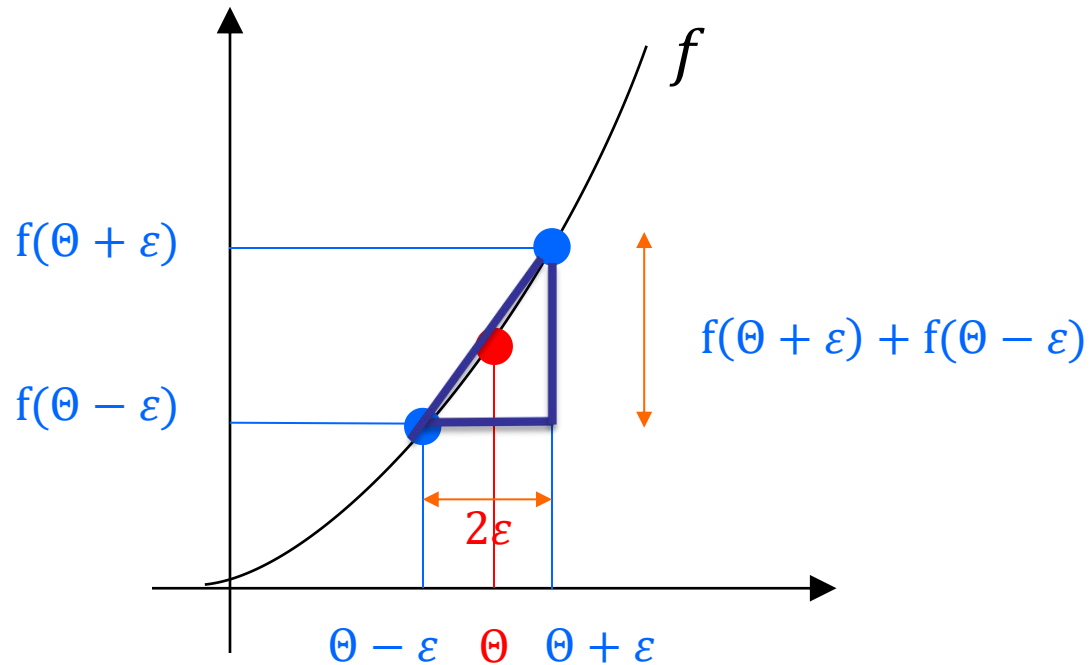
Variance cho TanH: $np.sqrt\left(\frac{1}{n^{[l-1]}}\right)$

Xavier 's initialization

Variance cho TanH (phiên bản khác): $np.sqrt\left(\frac{2}{n^{[l-1]} * n^{[l]}}\right)$

Bengio et al.

Gradient Checking – Lý thuyết



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

Gradient checking

- Lấy các trọng số $W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}$ và chuyển thành vector Θ .
+ $J(W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}) = J(\Theta)$.
- Lấy các trọng số $dW^{[1]}, db^{[1]}, \dots, dW^{[L]}, db^{[L]}$ và chuyển thành vector $d\Theta$.
+ $dJ(W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}) = dJ(\Theta)$.

Gradient Descent checking example

For i from 1 to L :

$$d_{\Theta_{\text{Approx}}}^{[i]} = \frac{J(\Theta_1, \Theta_2, \dots, \Theta_i + \varepsilon) - J(\Theta_1, \Theta_2, \dots, \Theta_i - \varepsilon)}{2\varepsilon} \approx d_{\Theta}^{[i]}$$

$$\partial = \frac{\|d_{\Theta_{\text{Approx}}}^{[i]} - d_{\Theta}^{[i]}\|_2}{\|d_{\Theta_{\text{Approx}}}^{[i]}\|_2 + \|d_{\Theta}^{[i]}\|_2}$$

If $\partial \approx 10^{-7}$ then Create gradient descent
 else if $\partial \approx 10^{-7}$ then OK! But some bug happen!
 else if $\partial \approx 10^{-3}$ then Bug everywhere

Một vài lưu ý với Gradient checking

- Chỉ dùng cho debug, không dùng để huấn luyện.
 - + Tính $d\Theta_{\text{Approx}}^{[i]}$ tốn rất nhiều tài nguyên.
- Nếu gradient descent trả về bug, thì kiểm tra lại các thành phần của layer thứ i đó.
 - + Các thành phần như: W , b , dW , db .
- Nhớ sử dụng regularization nếu như hàm chi phí có sử dụng.
- Không hoạt động với kỹ thuật Dropout.
- Thực hiện grad check ngay sau khi khởi tạo ngẫu nhiên trọng số cho W và b , và sau khi huấn luyện qua một số lần.

TỔNG KẾT

1. Các kỹ thuật chia dữ liệu để huấn luyện mạng neural.

Chia dữ liệu thành train/dev/test (8-1-1 với dữ liệu nhỏ và 9-0.5-0.5 với dữ liệu lớn).

2. Các loại lỗi trong Deep learning.

High bias (underfitting) và high variance (overfitting).

3. Các kỹ thuật giảm overfitting (high variance).

Chuẩn hoá mạng (regurlazation), tăng dữ liệu (data augmentation), early stopping.

4. Các kỹ thuật tăng cường tốc độ cho gradient descent.

Normalizing input data, Dropout, Parameter Initialization, Grad checking.

TÀI LIỆU THAM KHẢO

1. Khoá học *Neural Network and Deep learning*, deeplearning.ai.
2. Ian Goodfellow, Yoshua Bengio, Aaron Courville, *Deep learning*, MIT Press, 2016.
3. Andrew Ng., *Machine Learning Yearning*. Link:
<https://www.deeplearning.ai/machine-learning-yearning/>
4. Vũ Hữu Tiệp, *Machine Learning cơ bản*, NXB Khoa học và Kỹ thuật, 2018.