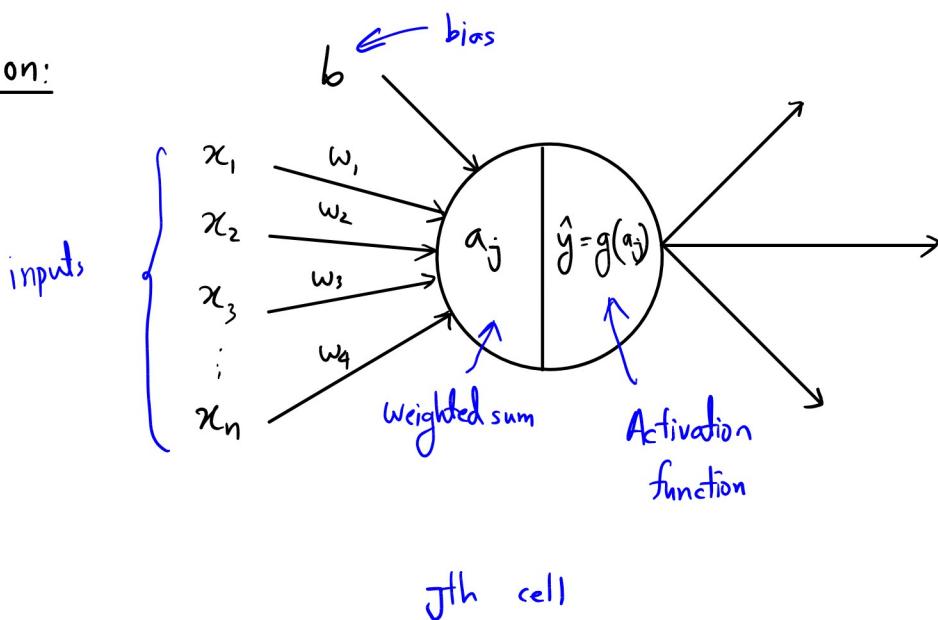


Neural Network

Book Ch 7, 7.1

Neuron:



jth cell

Weighted sum, $a_j = x_1 w_1 + x_2 w_2 + \dots + x_n w_n + b$

Let, $w_0 = b$ and $x_0 = 1$,

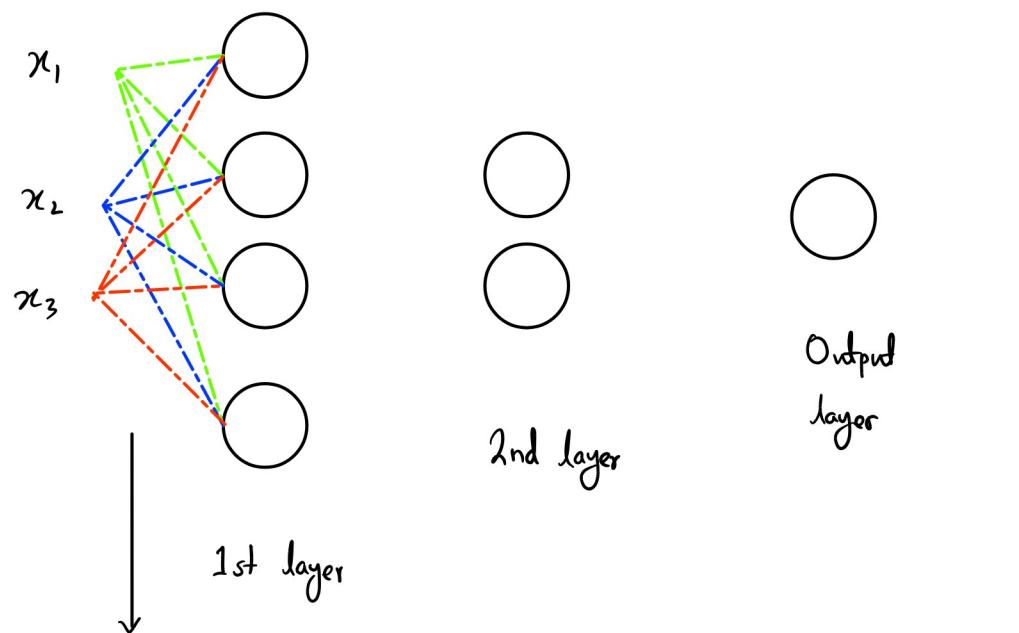
[Replacing the bias, book page 144]

$$a_j = \underbrace{x_0 w_0}_{\text{bias}} + x_1 w_1 + x_2 w_2 + \dots + x_n w_n$$

\therefore weighted sum of a single unit, $a_j = \sum_k x_k w_k$

$$\begin{aligned}\therefore \text{Output of } &= g \left([x_0 \ x_1 \ x_2 \ \dots \ x_n] \begin{bmatrix} w_{1,j} \\ w_{2,j} \\ \vdots \\ w_{n,j} \end{bmatrix} \right) \\ &= g(x \quad w_{*,j})\end{aligned}$$

Now, considering the following neural network



$$W^{(1)} = \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} & w_{1,4} \\ w_{2,1} & w_{2,2} & w_{2,3} & w_{2,4} \\ w_{3,1} & w_{3,2} & w_{3,3} & w_{3,4} \end{bmatrix} \quad 3 \times 4 \quad \text{Weight matrix}$$

$$\begin{aligned} \text{So, } a^{(1)} &= x \times W^{(1)} \\ &= [1 \times 3] \times [3 \times 4] \quad \leftarrow \text{matrix shape} \\ &= [1 \times 4] \quad \leftarrow \dots \end{aligned}$$

After using activation function on each element,

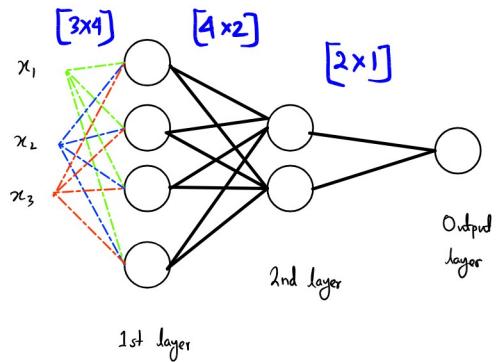
$$\begin{aligned} g(a^{(1)}) &= g([1 \times 4]) \\ &= [1 \times 4] \end{aligned}$$



This $[1 \times 4]$ contains the output of 4 neurons in the 1st layer.

These will be the input of 2nd layer.

Now we can easily find out the required weight matrix shape.



∴ Output of this NN,

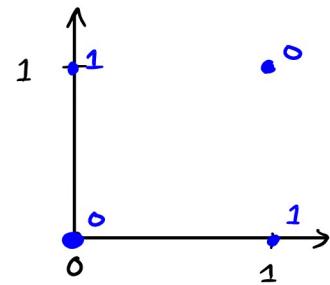
$$\begin{aligned}
 O &= g(g(g(x \times w^{(1)}) \times w^{(2)}) \times w^{(3)}) \\
 &= g(g(g([1 \times 3] [3 \times 4]) \times [4 \times 2]) \times [2 \times 1]) \\
 &= g(g([1 \times 4] \times [4 \times 2]) \times [2 \times 1]) \\
 &= g([1 \times 2] \times [2 \times 1]) \\
 &= g([1 \times 1])
 \end{aligned}$$

$$\text{Output} = [1 \times 1]$$

Determining the size of matrix are very important.

* Using the NN to solve the XOR problem

Ch 7.2



* We can not find a boundary line in this case using LR or a single neuron.

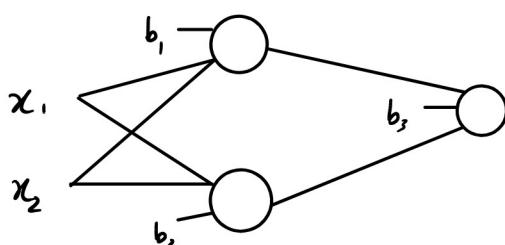
$$X = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \quad y = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

$$y = Xw + b$$

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} ? \\ ? \end{bmatrix} + ?$$

* No such w exists.

But using a NN,



$$g(x) = \text{ReLU}(x) = \max(0, x)$$

$$\Rightarrow \max(0, xw^{(1)} + b^{(1)}) w^{(2)} b^{(2)} = y$$

$$= \max \left(0, \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} ? & ? \\ ? & ? \end{bmatrix} + \begin{bmatrix} ? & ? \end{bmatrix} \right) \begin{bmatrix} ? \\ ? \end{bmatrix} + ? = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

If we use the following weights & Bias,

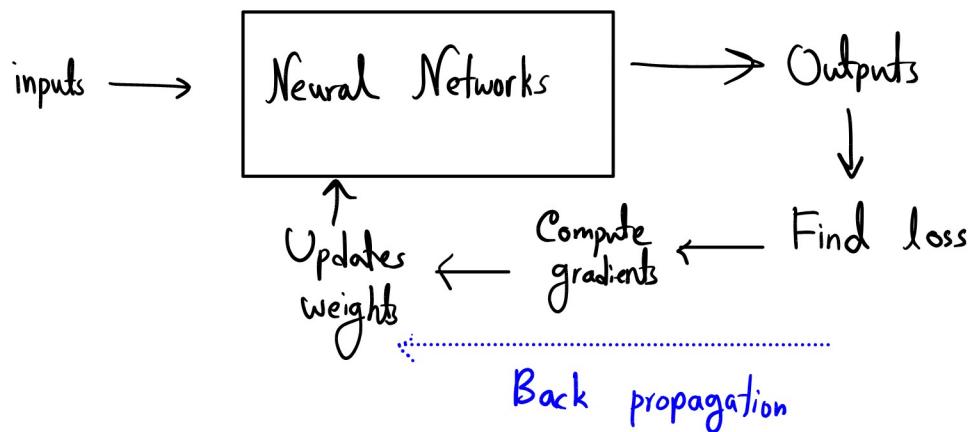
$$= \max \left(0, \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} + \begin{bmatrix} 0 & -1 \end{bmatrix} \right) \begin{bmatrix} 1 \\ -2 \end{bmatrix} + 0$$

$$= \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$

$$= \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = y$$

Note: There are multiple solution for this XOR problem. You can use other weights & other activation functions too.

Forward propagation



Component of a Neural Networks:

Activation functions:

[Ch 7-3.1]

In neural network, we denote activation function by $g(x)$ which can be any function.

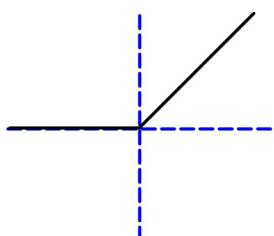
$g(x)$ can be linear function. But a non-linear function is required. Why?

Because, with linear activation functions, the neural network becomes just a linear transformation, no matter how many layer it has. So, it can not do complex mapping like the XOR problem.

Also, in a linear activation function, we can not use back propagation as the derivative of a linear function is constant.

Common Activation functions:

$$\text{ReLU: } h = \max(0, z)$$



Active only when input is positive

Gradient:

1 when $z = +ve$

0 when $z = -ve$

But, ReLU is non differentiable when $z = 0$.

in this case, left derivative = 0

right " " = 0

A few non differentiable point is not a problem as

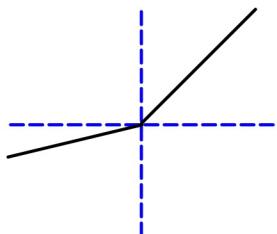
- ① Rarely we get a point with gradient 0 anyway
- ② Software simply return either left or right derivative.

ReLU has a dead neuron problem. As ReLU outputs 0 for all negative inputs, which can make some neurons stop activating entirely.

When this happens, their gradient becomes 0, and their weights never updates. Thus these neurons are dead.

Leaky ReLU:

$$h = \max(0.1z, z)$$



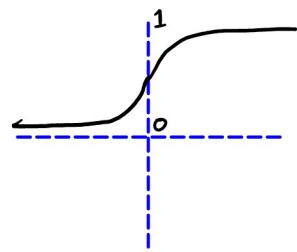
- Strong activation when z is positive
- weak negative activation, $\therefore z$ is negative

Gradient:

- 1 when positive
- 0.1 when negative

Sigmoid:

$$f(z) = \frac{1}{1 + e^{-z}}$$

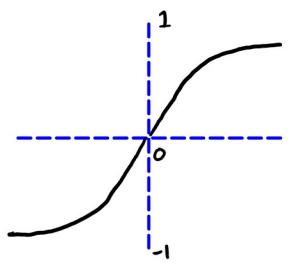


→ Maps z into probability
→ 0/1 switch

Gradient:

Saturates across most of its domains.

Tanh:



$$f(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

→ -1/+1 switch

Gradient:

Saturates across most of its domains.

Optimizes slightly better as around $z=0$, the gradient is closer to 1.

NN weights aka parameters are learnt through out the training phase.

Hyperparameters:

→ The parameters that we can control to facilitate the learning are called hyperparameter.

- Model Architecture Related
 - ① Number of layers
 - ② Number of neuron per layer
 - ③ Activation functions in each layer

- Training Related:
 - ① Optimizer
 - ② Number of epoch
 - ③ Learning Rate
 - ④ Batch Size
 - ⑤ Regularization & Normalization

& many more.

Optimizer:

These are algorithms that are used to update weights in order to reduce loss.

Common optimizers:

- | | | |
|---|--|---|
| <div style="display: inline-block; transform: rotate(-90deg); transform-origin: left top; position: absolute; left: -10px; top: 100px;">First
Order
Methods</div> <div style="display: inline-block; vertical-align: middle; margin-left: 10px;"> <ul style="list-style-type: none"> • Gradient Descent • Stochastic Gradient Descent • Mini-batch Gradient Descent • Momentum </div> | <div style="display: inline-block; border-left: 1px solid black; border-right: 1px solid black; padding: 0 10px; margin: 0 10px;"></div> | <ul style="list-style-type: none"> • Use first order derivatives only • Maintain a single learning rate |
|---|--|---|

- | |
|---|
| <div style="display: inline-block; transform: rotate(-90deg); transform-origin: left top; position: absolute; left: -10px; top: 100px;">Adaptive
methods</div> <div style="display: inline-block; vertical-align: middle; margin-left: 10px;"> <ul style="list-style-type: none"> • Adagrad • RMS prop • Adadelta • Adam </div> |
|---|

- AdaGrad: → provide learning rate for each parameter.
→ improves learnability but computationally expensive.
- RMS prop: → provide learning rate for each parameter.
→ adapts based on the mean of the recent magnitudes of the gradient of each weight
- AdaDelta: → Similar as RMS prop but works with second order derivatives.

* Adam: Adaptive Moment Estimation

- Most popular and mostly used
- Straight forward to implement
- Little memory requirements
- Need little to no manual tuning.

*⇒ It combines Momentum & RMSProp for efficient & adaptive update.

Regularization:

[Ch 5.7]

→ A set of strategies used in ML to reduce Generalization Error so that model can perform well on unseen data, not just training data.

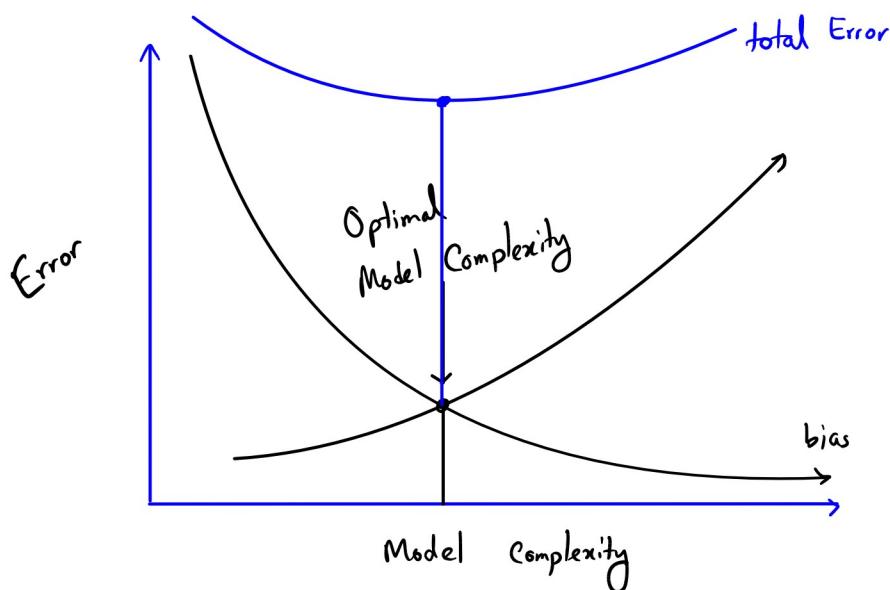
* → Bias Variance Tradeoff

Bias: • Error from overly simple assumption by simple models.

- High bias means the model makes strong assumptions and often results in Underfitting.
- Observation doesn't matter.

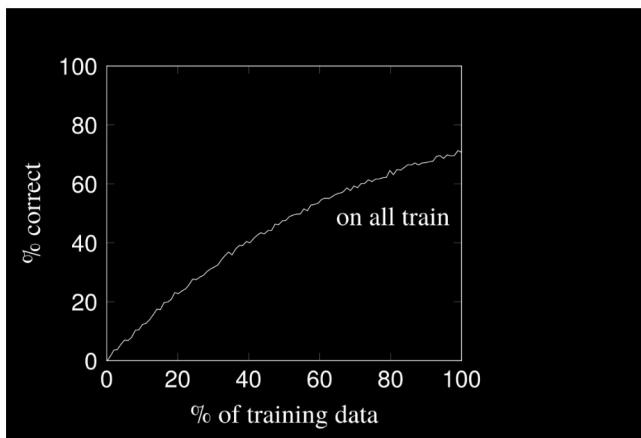
Variance: • Error from sensitivity to small fluctuations in the training data.

- High variance may result in modeling the noise in training data, causes Overfitting
- focuses too much in observation

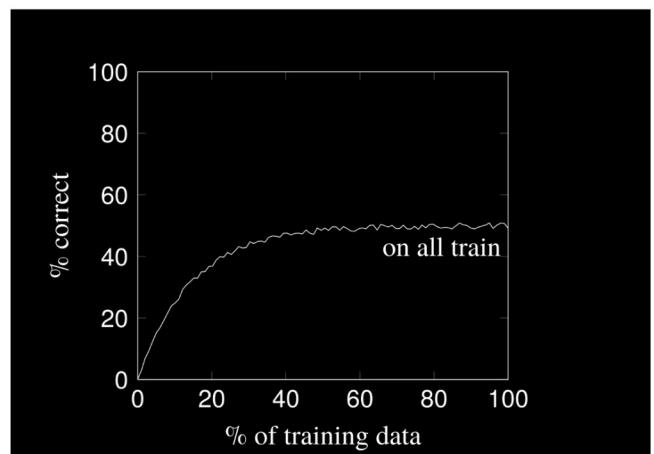


Underfitting: Low training & test accuracy.

- Insufficient training data



- Too simple model



Addressing underfitting

More complex model

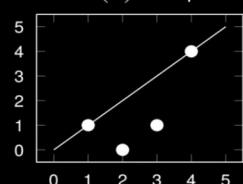
x_1	$f(x)$	$h(x)$
1	1	1
2	0	2
3	1	3
4	4	4

x_1	$f(x)$	$h(x)$
1	1	1
2	0	0
3	1	1
4	4	4

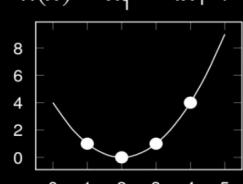
More features

x_1	x_2	$f(x)$	$h(x)$
1	1	1	1
2	4	0	0
3	9	1	1
4	16	4	4

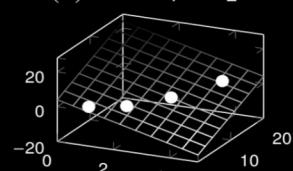
$$h(x) = x_1$$



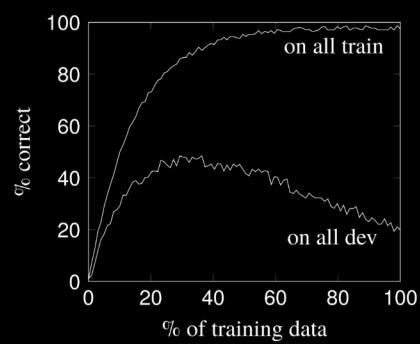
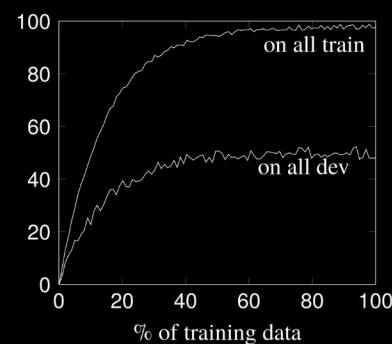
$$h(x) = x_1^2 - 4x_1 + 4$$



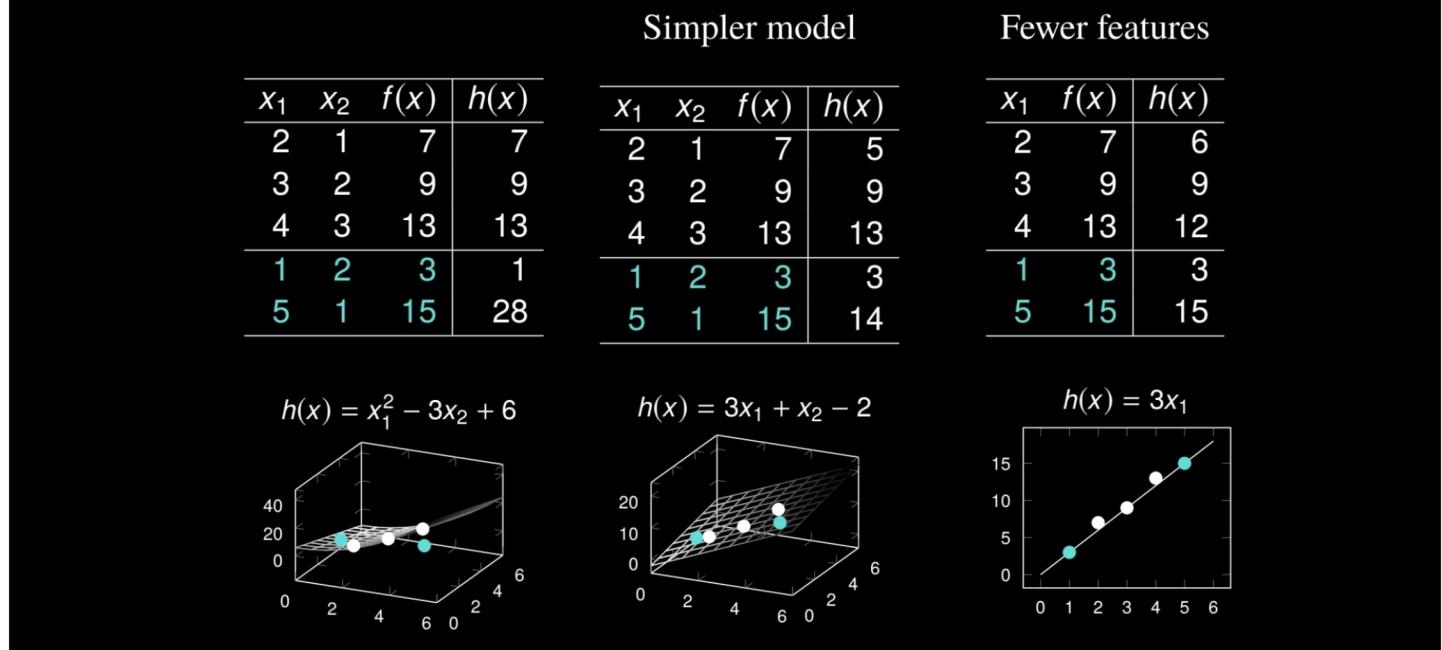
$$h(x) = -4x_1 + x_2 + 4$$



Overfitting: High training accuracy but low dev accuracy.



Addressing overfitting



What can be done to get the balance between bias & variance?

① Introduce a regularization term in the loss function

→ Adds a small bias to reduce variance → Reduces overfitting

$$\text{Loss} = \frac{\text{Original Loss}}{\text{Cross entropy or others}} + \frac{(\lambda * \text{Regularization Term})}{\text{L1 or L2}}$$

λ controls how strong the penalty is.

Common Regularization methods:

L2: Ridge Regression

$$\text{Penalty} = \lambda \sum w_i^2$$

- Adds the squared magnitude of weights to the loss
- Penalize the large weights more strongly.

- Keeps all features but shrinks their influence
- Higher $\lambda \rightarrow$ smaller weights \rightarrow smoother model
 \hookrightarrow less complex

L1: Lasso Regression

$$\text{Penalty} = \lambda \sum |w_i|$$

- Adds the absolute value of weights to the loss
- Higher $\lambda \rightarrow$ more weights becomes 0
- If a weight becomes 0, the feature doesn't affect the model.
 \hookrightarrow reduces complexity

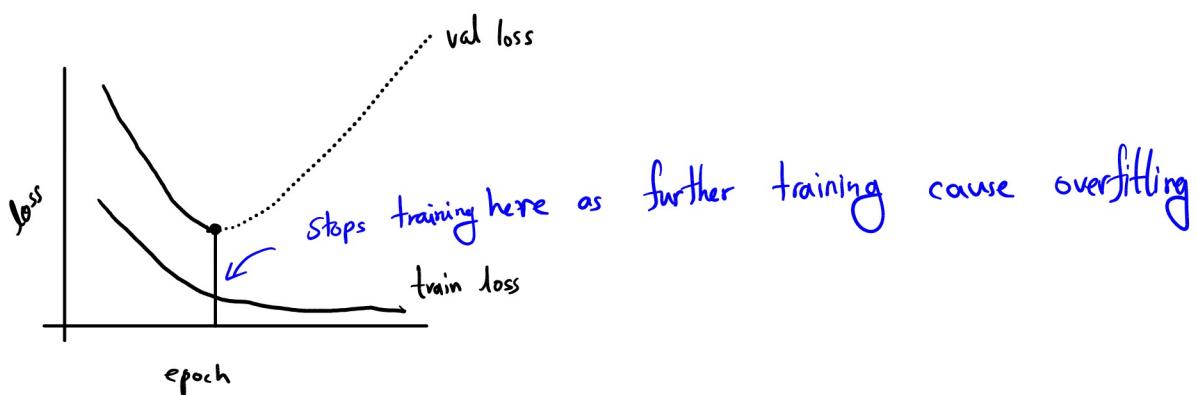
② Dropout:

Randomly drop outs (ignore) some neuron outputs during training with a probability of p .

\hookrightarrow depends on architecture
 \hookrightarrow Usually (0.2 to 0.5)

③ Early Stopping:

- Stops when validation performance stops improving.
- Returns the best current model.



④ Data Augmentation:

- Introduce new training data
- Injects noise

Other Hyperparameters:

① Learning Rate:

Lower LR = Slows down Learning but smooth converges.

Higher LR = speeds up learning but may not converge.

Decaying LR is preferred.

↳ Initially high LR, then low LR.

② Epochs:

how many time the entire training dataset has passed through the model.

③ Batch size:

a subset of training data. Weights are updated after each minibatch of training.

Smaller minibatch → too many update

Larger minibatch → too few updates, too many epoch to converge.

ML Best Practices:

Best practice:

Train	Val	Test
-------	-----	------

- Never peek the test data

Bad idea:

- ① Train on training data
- ② Evaluate on test data
- ③ Tune / Update parameters

Better idea

- ① Train on training data
- ② Tune on validation data
- ③ Evaluate on test data

Why? Indirectly using test data to guide learning. Causes overfitting.

Model has good test performance but performs poorly in real world data.

Performance Metrics

[Ch 4.7, 4.8]

$$\text{Accuracy} = \frac{\text{Correct}}{\text{Total}}$$

- A bad metric, specially for
- If a dataset contains 95% not spam and 5% spam and the model always predict not spam always, then the accuracy is 95%, which is useless

Confusion Metrics:

		Predicted	
		0	1
Actual	0	TN True Negative	FP False Positive
	1	FN False Negative	TP True Positive

Alternative of Accuracy:

$$\textcircled{1} \quad \underline{\text{Precision}} = \frac{TP}{TP + FP}$$

- Measures how many predicted positives were actually correct.
- high precision = less false positive
- Used when false alarm is costly. (Spam detection)

$$\textcircled{2} \quad \underline{\text{Recall}} = \frac{TP}{TP + FN}$$

- Measures how many actual positives were correctly classified.
- High Recall \rightarrow few false negative
- Use when missing a positive case is risky (Diseases Detection)

(3) F_1 Score:

$$F_1 = 2 * \frac{(\text{Precision} * \text{Recall})}{(\text{Precision} + \text{Recall})}$$

- Harmonic mean of precision and recall
- Balances both P and R
- Useful when dataset is imbalanced and both FP and FN matters

Other task specific Metrics:

- BLEU - Machine Translation
- WER - Automatic Speech Recognition
- Cosine Similarity - Semantic Similarity
- Euclidean distance - Spell checking.
- F-latency - Early Detection

Statistical Significance:

[Ch 4.9]

- Two models A and B has tiny score differences. (example: $A = 0.992 F_1$, $B = 0.991 F_1$)
is A actually better than B?
- We can't just trust one test set.
- So, we frame two hypothesis
 - Null hypothesis $\rightarrow H_0 : A \text{ is not truly better than } B$ (difference is luck)
 - Alternative $\rightarrow A \text{ is actually better}$
- Now, we gather multiple evaluation samples (via cross validation)
- Then we run significance test (Student's t-test) to get a p-value.
- If pvalue is small, we reject H_0 and say, 'A is significantly better than B'