

Linear regression

Cost function, Gradient descent

The objective of linear regression is to minimize the cost function

$$J = \frac{1}{m} \sum_{i=1}^m (\hat{y}^i - y^i)^2$$

$$\hat{y} = b + Wx = W_0 + W_1x$$

- J : cost function
- \hat{y} : predicted target value
- y : actual target value
- x : input value
- b : bias
- W : weight

Linear regression

(Batch) Gradient Descent

$$W_j := W_j - \alpha \frac{\partial J}{\partial W}$$

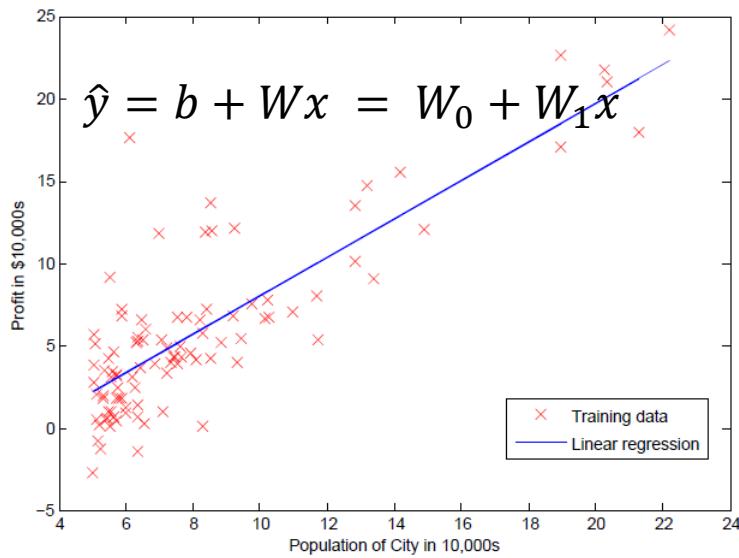
- α : learning rate
- $\frac{\partial J}{\partial W}$: gradient

$$\frac{\partial J}{\partial W_0} = \frac{2}{m} \sum_{i=1}^m (\hat{y}^i - y^i)$$

$$\frac{\partial J}{\partial W_1} = \frac{2}{m} \sum_{i=1}^m (\hat{y}^i - y^i) x$$

Linear regression

Linear regression fit



- Weights
 W_0, W_1, \dots
- Performance measure
 - Cost

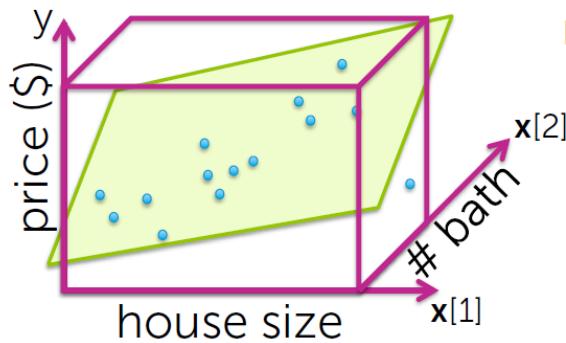
$$J = \frac{1}{m} \sum_{i=1}^m (\hat{y}^i - y^i)^2$$

- Root Mean Square Error (RMSE)

$$RMSE = \sqrt{\frac{1}{m} \sum_{i=1}^m (\hat{y}^i - y^i)^2}$$

Linear regression

With multiple variables



$$\hat{y} = W_0 + W_1x_1 + W_2x_2 + \dots$$

$$\frac{\partial J}{\partial W_0} = \frac{2}{m} \sum_{i=1}^m (\hat{y}^i - y^i)$$

$$\frac{\partial J}{\partial W_1} = \frac{2}{m} \sum_{i=1}^m (\hat{y}^i - y^i) x_1$$

$$\frac{\partial J}{\partial W_2} = \frac{2}{m} \sum_{i=1}^m (\hat{y}^i - y^i) x_2$$

<Machine Learning, Emily Fox & Carlos Guestrin>

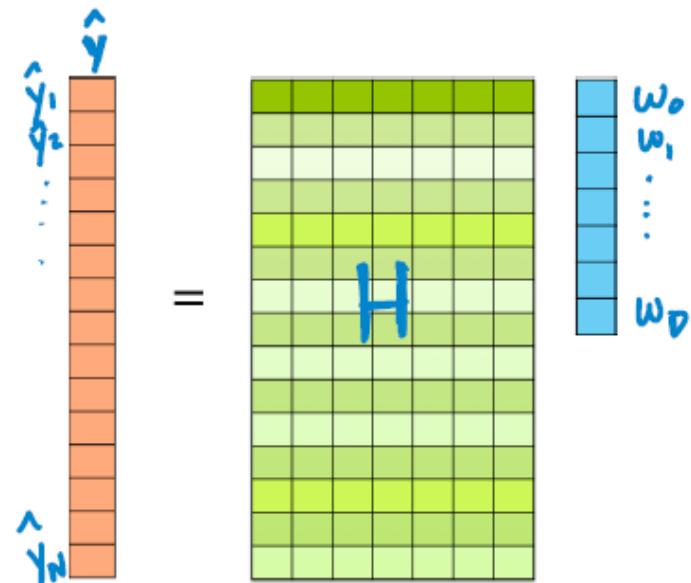
Linear regression

Matrix calculation

$$J = \frac{1}{m} \sum_{i=1}^m (\hat{y}^i - y^i)^2$$

$$J = \frac{1}{m} (\hat{Y} - Y)^T (\hat{Y} - Y)$$

$$J = \frac{1}{m} (W \cdot X - Y)^T (W \cdot X - Y)$$



<Machine Learning, Emily Fox & Carlos Guestrin>

Linear regression

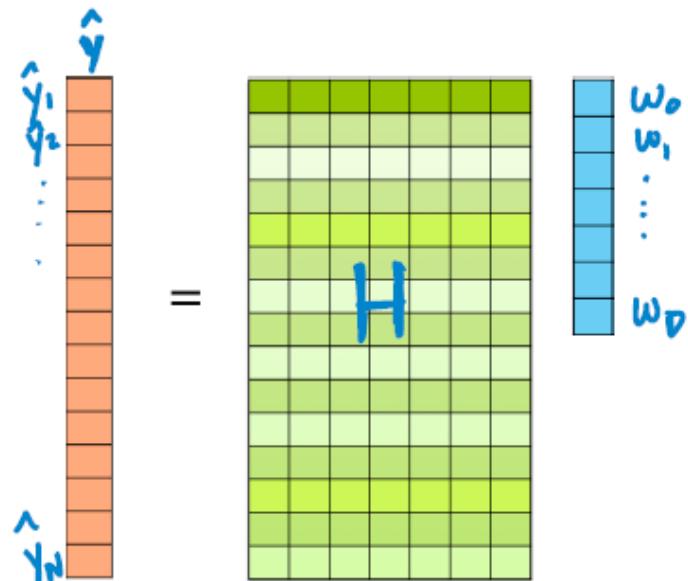
Matrix calculation for gradient

$$\frac{\partial J}{\partial W_j} = \frac{2}{m} \sum_{i=1}^m (\hat{y}^i - y^i) x_j$$

$$\frac{\partial J}{\partial W} = \frac{2}{m} (\hat{Y} - Y)^T \cdot X$$

$$\frac{\partial J}{\partial W} = \frac{2}{m} (X \cdot W - Y)^T \cdot X$$

$$\frac{\partial J}{\partial W} = \frac{2}{m} (X^T \cdot X \cdot W - X^T \cdot Y)$$



<Machine Learning, Emily Fox & Carlos Guestrin>

Linear regression

Normal equation

To find the value of θ that minimizes the cost function, there is a *closed-form solution*

—in other words, a mathematical equation that gives the result directly. This is called the *Normal Equation*.

$$W = (X^T \cdot X)^{-1} \cdot X^T \cdot Y$$

Normal equation

$$W = (X^T \cdot X)^{-1} \cdot X^T \cdot Y$$

Gradient descent

- Need to choose learning rate α
- Need to feature scaling
- Need to many iterations
- Works well even large number of features

Normal equation

- No Need to choose α
 - No Need to feature scaling
 - No need to iterations
 - Needs to compute inverse matrix
- $$(X^T \cdot X)^{-1}$$

Feature Scaling

Feature normalization

- Min-Max scaling: $(X - \text{min.value}) / (\text{max.value} - \text{min.value})$
 - shifted and rescaled. 0~1
- Standardization: $(X - \text{mean.value}) / (\text{std.value})$
 - zero mean and unit standard deviation
 - does not bound a specific range.
 - less affected by outliers

Feature normalization for test data

$$(X - \text{mean.value}) / (\text{std.value})$$

In order to normalize validation / test data set, you need to use mean value and standard deviation value for training data set.

Avoiding Overfitting through Regularization

- **Sum of squares**

Ridge : L2 norm

$$W_0^2 + W_1^2 + W_2^2 + \dots = \sum_{j=1}^n W_j^2 = \|W\|_2^2$$

Cost function $+\lambda\|W\|_2^2$

- **Sum of absolute value**

Lasso: L1 norm

$$|w_0| + |w_1| + |w_2| + \dots = \sum_{j=1}^n |w_j| = \|W\|_1$$

Cost function $+\lambda\|W\|_1$

Regularized linear regression

Ridge regression (L2 penalty)

Cost function

$$J = \frac{1}{m} \sum_{i=1}^m (\hat{y}^i - y^i)^2 + \boxed{\frac{\lambda}{m} \sum_{j=1}^n W_j^2}$$

Regularized linear regression

Ridge regression (L2 penalty)

Gradient

$$\frac{\partial J}{\partial W_0} = \frac{2}{m} \sum_{i=1}^m (\hat{y}^i - y^i) x_0^i \quad \text{when } j=0$$

$$\frac{\partial J}{\partial W_j} = \frac{2}{m} \sum_{i=1}^m (\hat{y}^i - y^i) x_j^i + \boxed{\frac{2\lambda}{m} W_j} \quad \text{when } j>1$$

Regularized linear regression

Ridge regression (L2 penalty)

Gradient descent

Repeat {

$$W_0 = W_0 - \alpha \frac{2}{m} \sum_{i=1}^m (\hat{y}^i - y^i) x_0^i \quad \text{when } j=0$$

$$W_j = W_j - \alpha \left[\frac{2}{m} \sum_{i=1}^m (\hat{y}^i - y^i) x_j^i + \frac{2\lambda}{m} W_j \right] \quad \text{when } j > 1$$

}

Regularized linear regression

Ridge regression (L2 penalty)

Matrix form

Cost function

$$J = \frac{1}{m} [(W \cdot X - Y)^T (W \cdot X - Y) + \lambda W^T \cdot W]$$

Gradient Descent

$$\frac{\partial J}{\partial W} = \frac{2}{m} [(X \cdot W - Y)^T \cdot X + \lambda W]$$

I: identity matrix

$$\frac{\partial J}{\partial W} = \frac{2}{m} [(X \cdot W - Y)^T \cdot X + \lambda I \cdot W]$$

Regularized linear regression

Ridge regression (L2 penalty)

Normalization equation for ridge regression

$$W = (X^T \cdot X + \lambda I)^{-1} \cdot X^T \cdot Y$$

when $\lambda=0$, $W = (X^T \cdot X)^{-1} \cdot X^T \cdot Y$

When $\lambda=\infty$, $W=0$

Regularized linear regression

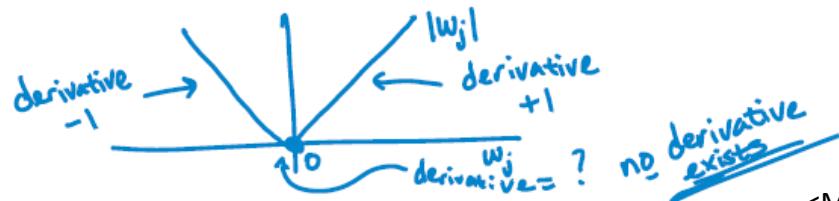
Lasso regression (L1 penalty)

Matrix form

Cost function

$$J = \frac{1}{m} [(W \cdot X - Y)^T (W \cdot X - Y) + \lambda |W|]$$

Gradient Descent



Coordinate descent method

<Machine Learning, Emily Fox & Carlos Guestrin>

Bias / Variance Trade off

Large λ : high bias, low variance

Small λ : low bias, high variance

Bias

This part of the generalization error is due to wrong assumptions, such as assuming that the data is linear when it is actually quadratic. A high-bias model is most likely to underfit the training data.

Variance

This part is due to the model's excessive sensitivity to small variations in the training data. A model with many degrees of freedom (such as a high-degree polynomial model) is likely to have high variance, and thus to overfit the training data.

<Hands-On ML, Aurelien Geron>

Logistic Regression

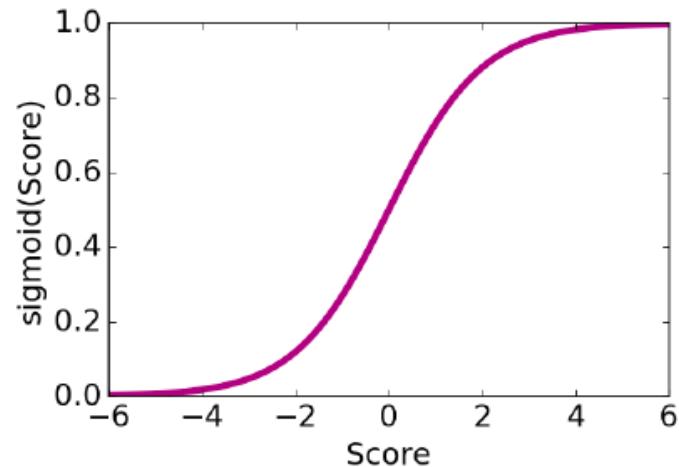
Given x , the probability if $y = 1$

$$\hat{y} = P(y = 1|x)$$

$$P = \sigma(W^T x + b)$$

Logistic function (sigmoid)

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$



<Machine Learning, Emily Fox & Carlos Guestrin>

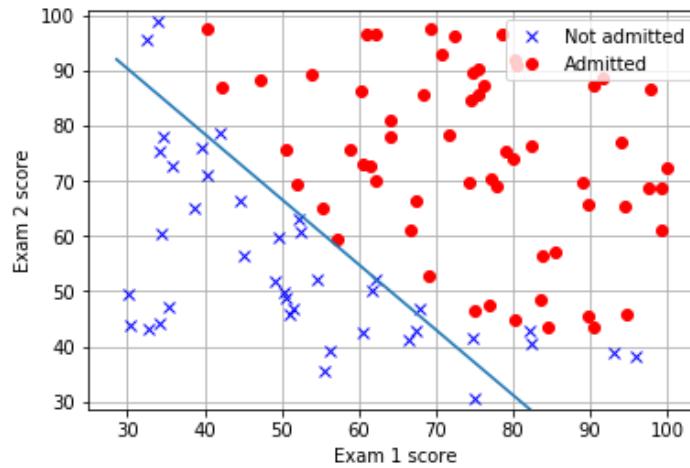
Logistic regression

Model Prediction

$$P(y = 1|x)$$

$$= \frac{1}{1+\exp(-W^T x - b)}$$

$$\hat{y} = \begin{cases} 1 & \text{if } P \geq 0.5 \\ 0 & \text{if } P < 0.5 \end{cases}$$



Logistic Regression

Cost function

$$J = -\frac{1}{m} \sum_{i=1}^m [y^i \log(\hat{y}^i) + (1 - y^i) \log(1 - \hat{y}^i)]$$

If $y = 1$ $J = -\frac{1}{m} \sum_{i=1}^m [y^i \log(\hat{y}^i)]$

If $y = 0$ $J = -\frac{1}{m} \sum_{i=1}^m [\log(1 - \hat{y}^i)]$

<Machine Learning, Emily Fox & Carlos Guestrin>

Logistic Regression

Regularization

Cost function

$$J = -\frac{1}{m} \sum_{i=1}^m [y^i \log(\hat{y}^i) + (1 - y^i) \log(1 - \hat{y}^i)] + \frac{\lambda}{m} \sum_{j=1}^n W_j^2$$

gradient

$$\frac{\partial J}{\partial W_0} = -\frac{1}{m} \sum_{i=1}^m [(\hat{y}^i - y^i)x_j^i] \quad \text{for } j=0$$

$$\frac{\partial J}{\partial W_j} = -\frac{1}{m} \sum_{i=1}^m [(\hat{y}^i - y^i)x_j^i] + \frac{2\lambda}{m} W_j \quad \text{for } j \geq 1$$

Logistic Regression

Gradient (or derivative of cost function)

$$\frac{\partial J}{\partial W} = -\frac{1}{m} \sum_{i=1}^m [(\hat{y}^i - y^i)x_j^i]$$

$$\frac{\partial J}{\partial W} = -\frac{1}{m} \sum_{i=1}^m [(P - y^i)x_j^i]$$

$$\frac{\partial J}{\partial W} = -\frac{1}{m} \sum_{i=1}^m [(\sigma(W^T x + b) - y^i)x_j^i]$$

Logistic Regression

Regularization

Cost function

$$J = -\frac{1}{m} \sum_{i=1}^m [y^i \log(\hat{y}^i) + (1 - y^i) \log(1 - \hat{y}^i)] + \frac{2\lambda}{m} \sum_{j=1}^n W_j^2$$

gradient

$$\frac{\partial J}{\partial W_0} = -\frac{1}{m} \sum_{i=1}^m [(\hat{y}^i - y^i)x_j^i] \quad \text{for } j=0$$

$$\frac{\partial J}{\partial W_j} = -\frac{1}{m} \sum_{i=1}^m [(\hat{y}^i - y^i)x_j^i] + \frac{\lambda}{m} W_j \quad \text{for } j \geq 1$$

Logistic Regression

$$\hat{y} = \sigma(W^T x + b)$$

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

$$J = -\frac{1}{m} \sum_{i=1}^m [y^i \log(\hat{y}^i) + (1 - y^i) \log(1 - \hat{y}^i)]$$

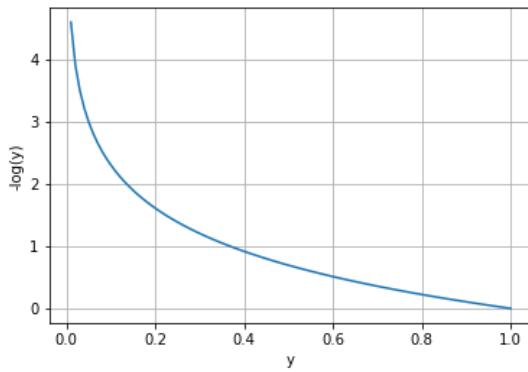
<Machine Learning, Emily Fox & Carlos Guestrin>

Logistic Regression

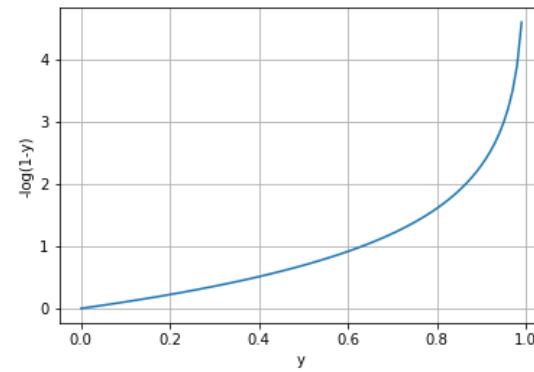
Loss function Given x , want $\hat{y}^i \approx y^i$

$$L = -[y^i \log(\hat{y}^i) + (1 - y^i) \log(1 - \hat{y}^i)]$$

$$\text{If } y = 1, \quad L = -\log(\hat{y}^i)$$



$$\text{If } y = 0, \quad L = -\log(1 - \hat{y}^i)$$



Confusion matrix

		Predicted label	
		(+1)	(-1)
True label	(+1)	1503 tp	203 fn
	(-1)	905 fp	1845 tn

Accuracy : $(tp + tn) / (tp + fp + fn + tn) * 100$

Precision : $tp / \text{predict positive} = tp / (tp + fp) * 100$

Recall (true positive rate) : $tp / \text{actual positive} = tp / (tp + fn) * 100$

True negative rate : $tn / \text{actual negative} = tn / (fp + tn) * 100$

Precision

Fraction of positive predictions that are correct.

$$\text{precision} = \frac{\# \text{ true positives}}{\# \text{ true positives} + \# \text{ false positives}}$$

Recall

Fraction of positive data points correctly classified

$$\text{Recall} = \frac{\# \text{ true positives}}{\# \text{ true positives} + \# \text{ false negatives}}$$

Mini-batch Gradient Descent

- Batch Gradient Descent

Computes the gradients based on full training set

Ex) Offline learning

- Stochastic Gradient Descent

Computes just one instance

Ex) Online learning

- Mini-batch Gradient Descent

Computes the gradients on small random sets of instances

<Hands-On ML, Aurelien Geron>

Stochastic gradient descent

**Batch
Gradient
descent**

$$\frac{\partial J}{\partial W} = -\frac{1}{m} \sum_{i=1}^m [(\hat{y}^i - y^i)x_j^i] \quad (\text{Sum over data points})$$

**Stochastic
Gradient
descent**

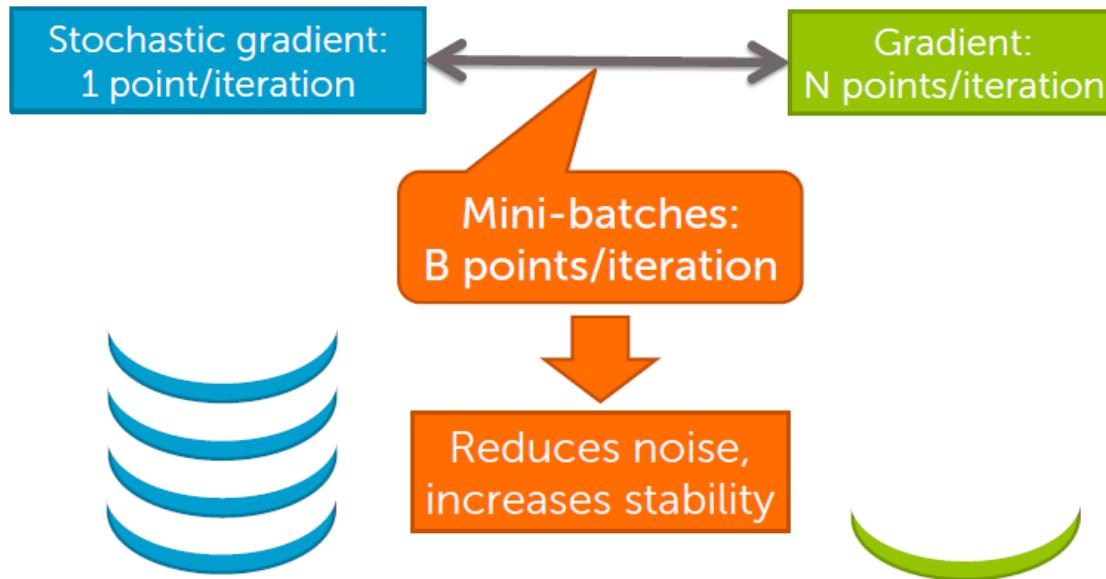
$$\frac{\partial J}{\partial W} = -(\hat{y}^i - y^i)x_j^i \quad (\text{Each time, pick different data point})$$

Batch vs. Stochastic gradient

Algorithm	Time per iteration	Total time to convergence for large data		
		In theory	In practice	Sensitivity to parameters
Batch Gradient	Slow for large data	Slower	Often slower	Moderate
Stochastic gradient	Always fast	Faster	Often faster	Very high

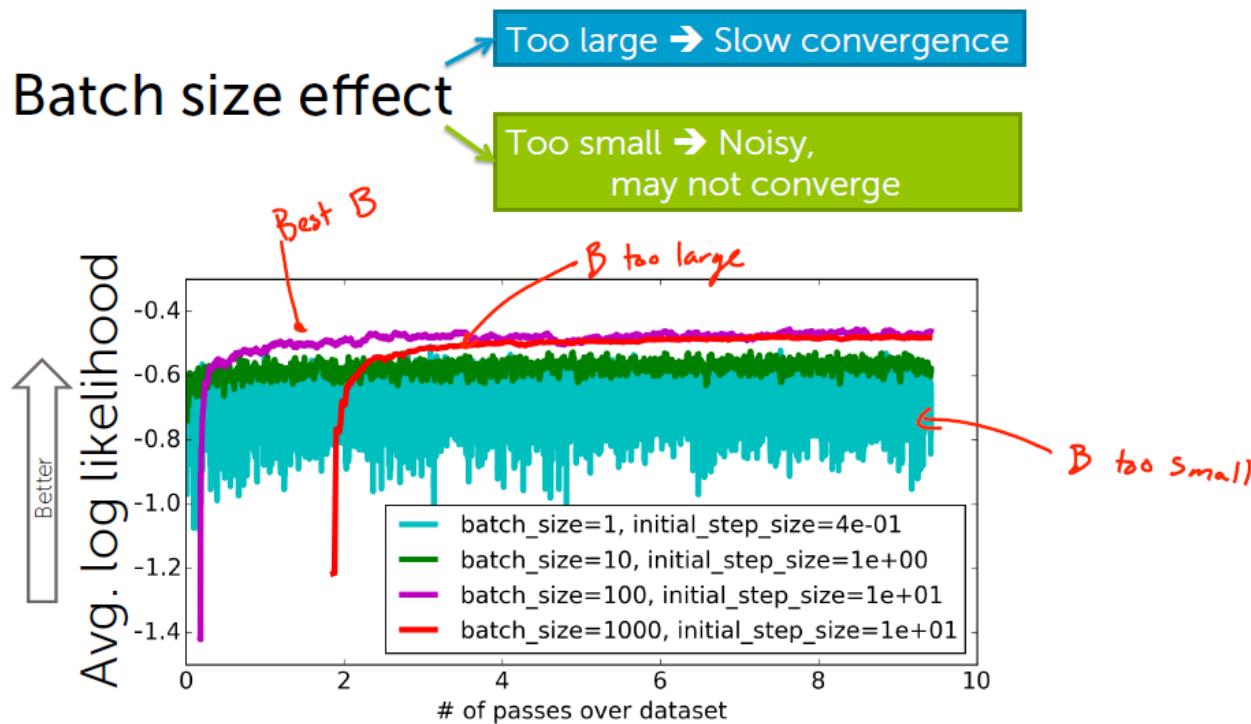
<Machine Learning, Emily Fox & Carlos Guestrin>

Batch / stochastic: two extremes



<Machine Learning, Emily Fox & Carlos Guestrin>

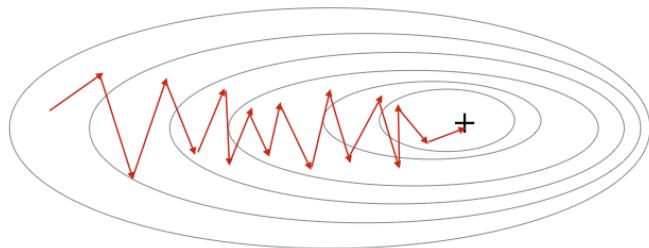
Batch size effect in mini-batch



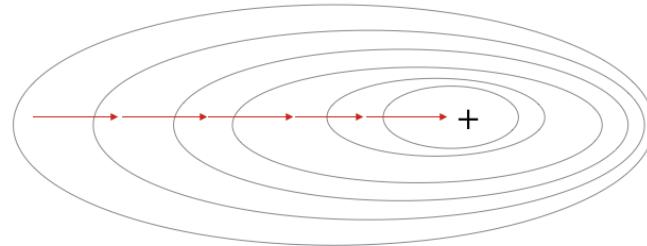
<Machine Learning, Emily Fox & Carlos Guestrin>

Converge

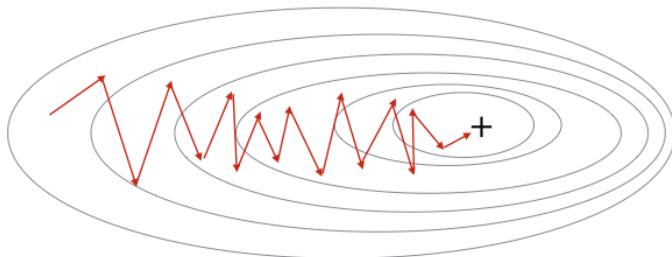
Stochastic Gradient Descent



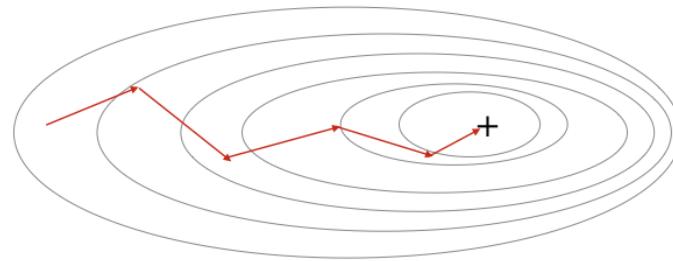
Gradient Descent



Stochastic Gradient Descent



Mini-Batch Gradient Descent



Softmax function

Softmax score for class k

$$s_k(x) = W^k \cdot x$$

Softmax function

$$p_k = \frac{\exp(s_k(x))}{\sum_{j=1}^K \exp(s_j(x))}$$

- K is the number of classes.
- $\mathbf{s}(\mathbf{x})$ is a vector containing the scores of each class for the instance \mathbf{x} .
- p_k is the estimated probability that the instance \mathbf{x} belongs to class k given the scores of each class for that instance.

Softmax classifier

- Multiclass classifier
- Provides normalized class probabilities

$$p_k = \frac{\exp(s_k(x))}{\sum_{j=1}^K \exp(s_j(x))}$$

$$\hat{y} = \operatorname{argmax}_k p_k$$

It predicts the class with the highest estimated probability.

Softmax in matrix

- for $x \in \mathbb{R}^{1 \times n}$, $\text{softmax}(x) = \text{softmax}([x_1 \quad x_2 \quad \dots \quad x_n]) = \left[\frac{e^{x_1}}{\sum_j e^{x_j}} \quad \frac{e^{x_2}}{\sum_j e^{x_j}} \quad \dots \quad \frac{e^{x_n}}{\sum_j e^{x_j}} \right]$
- for a matrix $x \in \mathbb{R}^{m \times n}$, x_{ij} maps to the element in the i^{th} row and j^{th} column of x , thus we have:

$$\text{softmax}(x) = \text{softmax} \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots & x_{1n} \\ x_{21} & x_{22} & x_{23} & \dots & x_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & x_{m3} & \dots & x_{mn} \end{bmatrix} = \begin{bmatrix} \frac{e^{x_{11}}}{\sum_j e^{x_{1j}}} & \frac{e^{x_{12}}}{\sum_j e^{x_{1j}}} & \frac{e^{x_{13}}}{\sum_j e^{x_{1j}}} & \dots & \frac{e^{x_{1n}}}{\sum_j e^{x_{1j}}} \\ \frac{e^{x_{21}}}{\sum_j e^{x_{2j}}} & \frac{e^{x_{22}}}{\sum_j e^{x_{2j}}} & \frac{e^{x_{23}}}{\sum_j e^{x_{2j}}} & \dots & \frac{e^{x_{2n}}}{\sum_j e^{x_{2j}}} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{e^{x_{m1}}}{\sum_j e^{x_{mj}}} & \frac{e^{x_{m2}}}{\sum_j e^{x_{mj}}} & \frac{e^{x_{m3}}}{\sum_j e^{x_{mj}}} & \dots & \frac{e^{x_{mn}}}{\sum_j e^{x_{mj}}} \end{bmatrix} = \begin{pmatrix} \text{softmax(first row of } x) \\ \text{softmax(second row of } x) \\ \dots \\ \text{softmax(last row of } x) \end{pmatrix}$$

<Deep Learning, Andrew Ng>

Backpropagation

Forward pass

$$f(u) = 5u$$

$$x(u) = 1 - u$$

$$u(t) = t^2$$

$$f(t) = 5(1 - t^2) = 5 - 5t^2$$

$$\frac{df}{dt} = -10t$$

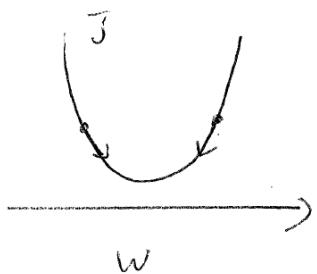
Backward pass

$$\frac{df}{dt} = \frac{df}{dx} \cdot \frac{dx}{du} \cdot \frac{du}{dt}$$

$$\frac{df}{dt} = (5) (-1) (2t)$$

$$= -10t$$

Linear regression & chain rule



$$W := W - \alpha \frac{\partial J}{\partial w}$$

$$J = \frac{1}{m} \sum (\hat{y} - y)^2$$

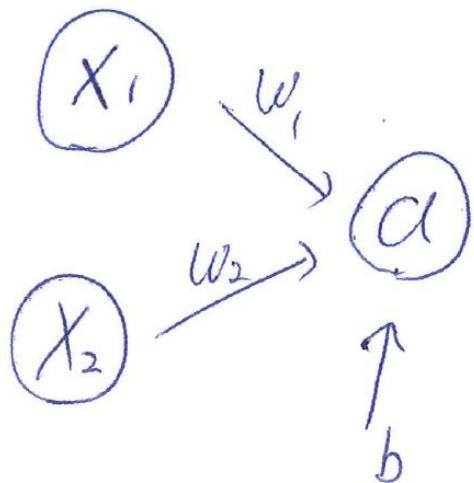
$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w}$$

$$\frac{\partial J}{\partial \hat{y}} = \frac{2}{m} \sum (\hat{y} - y)$$

$$\frac{\partial \hat{y}}{\partial w} = x \quad \text{since } \hat{y} = w \cdot x$$

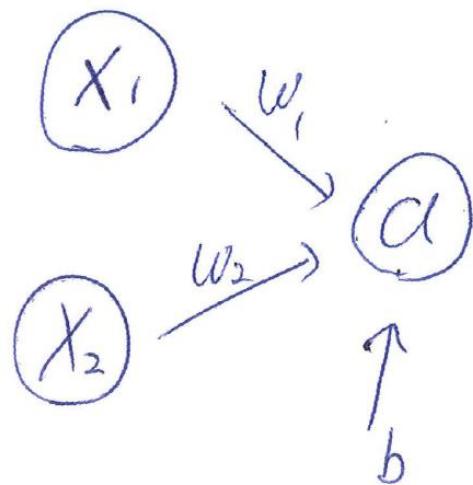
$$\therefore \frac{\partial J}{\partial w} = \frac{2}{m} \sum (\hat{y} - y) \cdot x$$

Logistic regression



$$\begin{aligned} & X_1 \quad w_1 \\ & X_2 \quad w_2 \\ & b \end{aligned} \rightarrow z = w_1 x_1 + w_2 x_2 + b \rightarrow \hat{y} = q = \sigma(z) \rightarrow L(\hat{y}, y)$$

Logistic regression



$$z = w^T x + b$$

$$\hat{y} = a = \sigma(z) = \frac{1}{1 + \exp(-z)}$$

$$L(\hat{y}, y) = -[y \log(a) + (1-y) \log(1-a)]$$

Two neurons (Logistic regression)

$$a^{[v]} \xrightarrow{O} a^{[u]}$$

$$a^{[v]} = \sigma(w \cdot a^{[u]} + b)$$

$$L = (a^{[v]} - y)^2$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial a^{[u]}} \frac{\partial a^{[u]}}{\partial w}$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial a^{[u]}} \frac{\partial a^{[u]}}{\partial b}$$

$$z^{[v]} = w \cdot a^{[u]} + b$$

$$a^{[v]} = \sigma(z^{[v]})$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial a^{[u]}} \frac{\partial a^{[u]}}{\partial z^{[u]}} \frac{\partial z^{[u]}}{\partial w}$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial a^{[u]}} \frac{\partial a^{[u]}}{\partial z^{[u]}} \frac{\partial z^{[u]}}{\partial b}$$

Derivative of da/dz

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

$$\begin{aligned}\frac{\partial a}{\partial z} &= \frac{d}{dz} [(1+e^{-z})^{-1}] \\&= e^{-z} (1+e^{-z})^{-2} \\&= \frac{1}{1+e^{-z}} \cdot \frac{e^{-z}}{1+e^{-z}} \\&= \frac{1}{1+e^{-z}} \cdot \left[\frac{1+e^{-z}-1}{1+e^{-z}} \right] \\&= \frac{1}{1+e^{-z}} \cdot \left[1 - \frac{1}{1+e^{-z}} \right] \\&= a [1-a]\end{aligned}$$

$$z = w^T x + b$$

$$\hat{y} = a = \sigma(z) = \frac{1}{1+\exp(-z)}$$

$$L(\hat{y}, y) = -[y \log(a) + (1-y) \log(1-a)]$$

Derivative of dL/dz

$$\frac{\partial L}{\partial a} = "da" = -\frac{y}{a} + \frac{1-y}{1-a}$$

$$"da" = \frac{\partial L(a, y)}{\partial a}$$

$$\frac{\partial L}{\partial z} = "dz" = \frac{\partial L(a, y)}{\partial z}$$

$$dz = \frac{\partial L}{\partial a} \frac{\partial a}{\partial z}$$

$$= \left[-\frac{y}{a} + \frac{1-y}{1-a} \right] [a(1-a)]$$

$$= a - y$$

$$z = w^T x + b$$

$$\hat{y} = a = \sigma(z) = \frac{1}{1 + \exp(-z)}$$

$$L(\hat{y}, y) = -[y \log(a) + (1-y) \log(1-a)]$$

Derivative of dL/dW

$$\frac{\partial L}{\partial w_1} = "dL/dw_1"$$

$$= \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$
$$= dz \cdot \frac{\partial z}{\partial w_1}$$

$$= (a-y) \cdot x_1$$

$$dL/dw_2 = \frac{\partial L}{\partial w_2}$$

$$= (a-y) \cdot x_2$$

$$dL/db = \frac{\partial L}{\partial b}$$

$$= (a-y) \cdot 1$$

$$= a - y$$

$$z = w^T x + b$$

$$\hat{y} = a = \sigma(z) = \frac{1}{1 + \exp(-z)}$$

$$L(\hat{y}, y) = -[y \log(a) + (1-y) \log(1-a)]$$

Gradient Descent on m observations

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

$$a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)$$

$$\frac{\partial}{\partial w_1} J(w, b) = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial}{\partial w_1} L(a^{(i)}, y)}_{d w_1^{(i)}}$$

Gradient Descent on m observations

Loop to calculate derivatives

For $i = 1$ to m :

$$z^{(i)} = w^\top x^{(i)} + b$$

$$\alpha^{(i)} = \sigma(z^{(i)})$$

$$J += -[y^{(i)} \log \alpha^{(i)} + (1-y^{(i)}) \log (1-\alpha^{(i)})]$$

$$d\hat{z}^{(i)} = \alpha^{(i)} - y^{(i)}$$

$$dw_1 += d\hat{z}^{(i)} \cdot x_i^{(i)}$$

$$dw_2 += d\hat{z}^{(i)} \cdot 1$$

$$db += d\hat{z}^{(i)}$$

After the loop

$$J = J/m$$

$$dw_1 = dw_1/m$$

$$dw_2 = dw_2/m$$

$$db = db/m$$

Gradient Descent

After the loop

$$J = J/m$$

$$dw_1 := dw_1/m$$

$$dw_2 := dw_2/m$$

$$db := db/m$$

for iteration in range (1000):

$$w_1 := w_1 - \alpha dw_1$$

$$w_2 := w_2 - \alpha dw_2$$

$$b := b - \alpha db$$

Gradient Descent using matrix form

For iter in range(1000):

$$z = np.dot(X, w) + b \quad : (m, 1)$$

$$A = \sigma(z) \quad : (m, 1)$$

$$d\hat{z} = A - Y \quad : (m, 1)$$

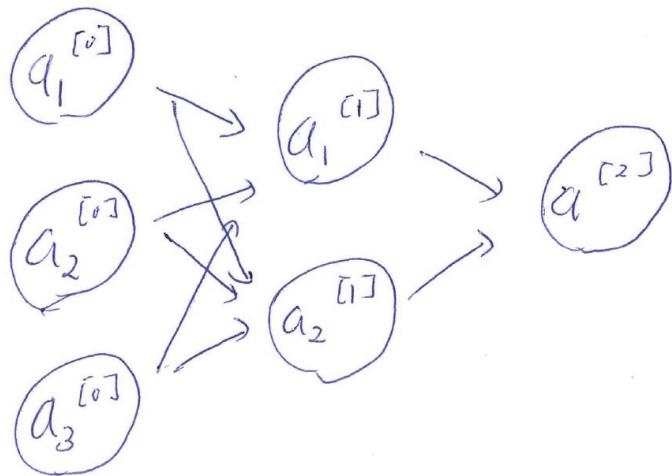
$$dw = \frac{1}{m} np.dot(X.T, d\hat{z}) \quad : (n, 1)$$

$$db = \frac{1}{m} np.sum(d\hat{z})$$

$$w := w - \alpha dw$$

$$b := b - \alpha db$$

Neural Network: 1 hidden layer

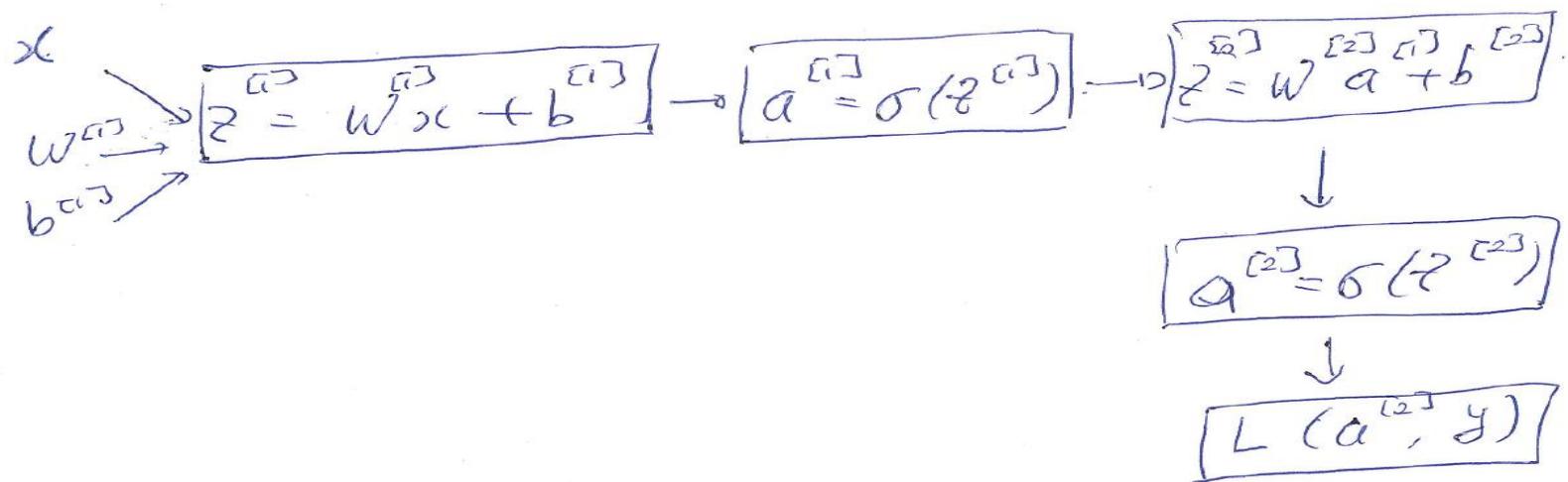


Parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$

Cost function: $J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]})$

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

Neural Network: 1 hidden layer



Derivative of $dL/dW^{[2]}$

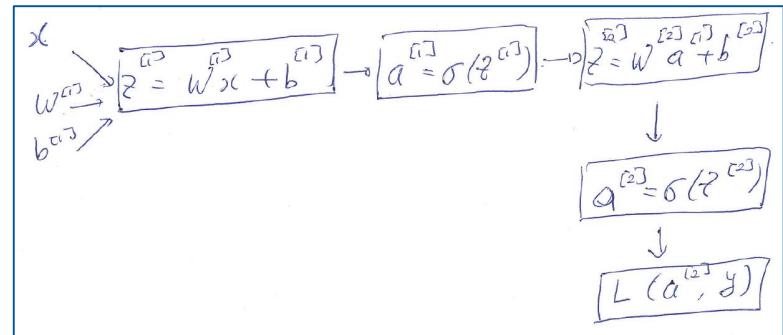
$$L(\hat{y}, y) = -[y \log(a^{[2]}) + (1-y) \log(1-a^{[2]})]$$

$$da^{[2]} = -\frac{y}{a^{[2]}} + \frac{1-y}{1-a^{[2]}}$$

$$dz^{[2]} = a^{[2]} - y$$

$$dW^{[2]} = (a^{[2]} - y) a^{[1]} = dz^{[2]} \cdot a^{[1]}$$

$$db^{[2]} = dz^{[2]} = a^{[2]} - y$$



Derivative of $dL/dZ^{[1]}$

$$da^{[1]} = \frac{\partial L}{\partial a^{[1]}} = \underbrace{\frac{\partial L}{\partial a^{[2]}}}_{\cdot} \frac{\partial a^{[2]}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial a^{[1]}}$$

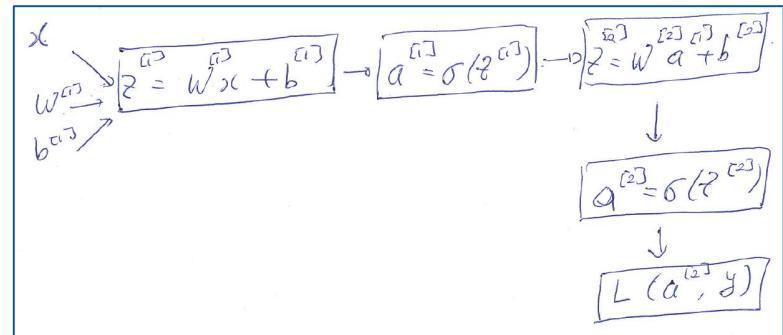
$$= dz^{[2]} \cdot w^{[2]}$$

$$= (a^{[2]} - y) \cdot w^{[2]}$$

$$dz^{[1]} = \frac{\partial L}{\partial z^{[1]}} = \underbrace{\frac{\partial L}{\partial a^{[2]}}}_{\cdot} \frac{\partial a^{[2]}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial a^{[1]}} \frac{\partial a^{[1]}}{\partial z^{[1]}}$$

$$= da^{[1]} \cdot a^{[1]}(1-a^{[1]})$$

$$= (a^{[2]} - y) \cdot w^{[2]} \cdot a^{[1]}(1-a^{[1]})$$



Derivative of $dL/dW^{[1]}$

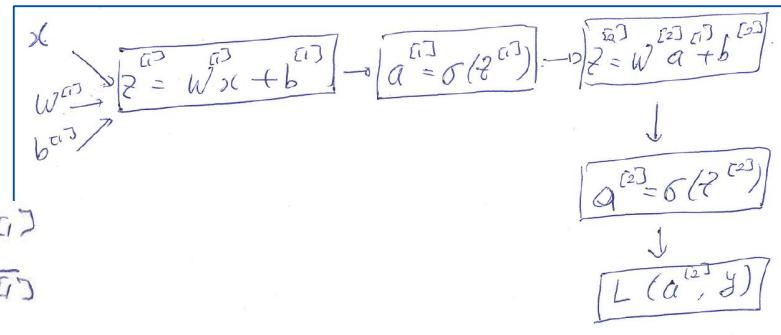
$$dW^{[1]} = \frac{\partial L}{\partial W^{[1]}} = \underbrace{\frac{\partial L}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial a^{[1]}} \frac{\partial a^{[1]}}{\partial z^{[1]}} \frac{\partial z^{[1]}}{\partial W^{[1]}}}_{\checkmark}$$

$$= dz^{[1]} \cdot \frac{\partial z^{[1]}}{\partial W^{[1]}}$$

$$= (a^{[1]} - y) \cdot w^{[2]} \cdot a^{[1]} \cdot (1 - a^{[1]}) \cdot x$$

$$= dz^{[1]} \cdot x$$

$$= dz^{[1]} \cdot a^{[0]}$$

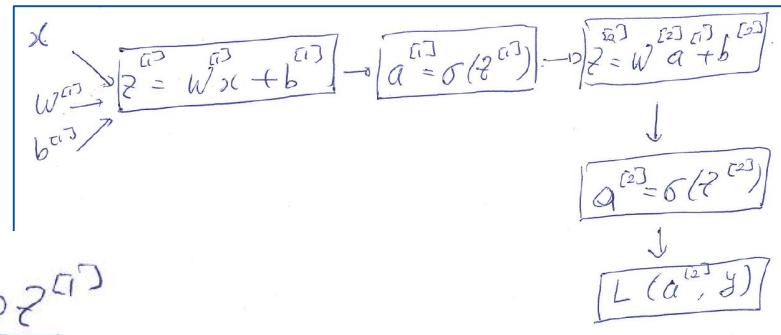


Derivative of $dL/db^{[1]}$

$$db^{[1]} = \frac{\partial L}{\partial b^{[1]}} = \underbrace{\frac{\partial L}{\partial a^{[2]}}}_{\text{Chain rule}} \frac{\partial a^{[2]}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial a^{[1]}} \frac{\partial a^{[1]}}{\partial z^{[1]}} \frac{\partial z^{[1]}}{\partial b^{[1]}}$$

$$= dz^{[1]} \cdot \frac{\partial z^{[1]}}{\partial b^{[1]}}$$

$$= dz^{[1]}$$



Forward propagation

$$z^{[1]} = w^{[1]} \cdot x + b^{[1]}$$

$$= w^{[1]} \cdot A^{[0]} + b^{[1]}$$

$$A^{[1]} = g^{[1]}(z^{[1]})$$

$$= \sigma(z^{[1]})$$

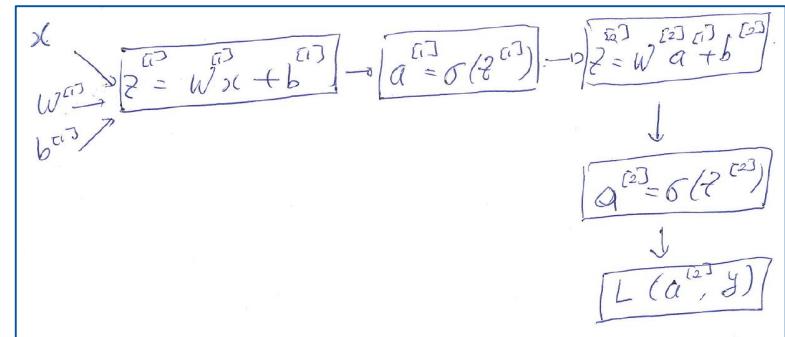
$$z^{[2]} = w^{[2]} \cdot A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(z^{[2]})$$

$$= \sigma(z^{[2]})$$

$$= \hat{Y}$$

© Inderen Jeong



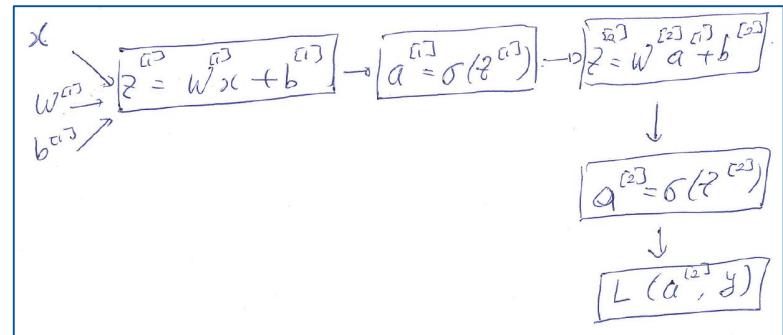
Backward propagation

$$dA^{[2]} = -\frac{Y}{A^{[2]}} + \frac{1-Y}{1-A^{[2]}} \rightarrow \text{skip}$$

$$dz^{[2]} = A^{[2]} - Y$$

$$dW^{[2]} = \frac{1}{m} dz^{[2]} A^{[1]}$$

$$db^{[2]} = \frac{1}{m} \text{np.sum}(dz^{[2]}, \text{axis}=1, \text{keepdims=True})$$



Backward propagation

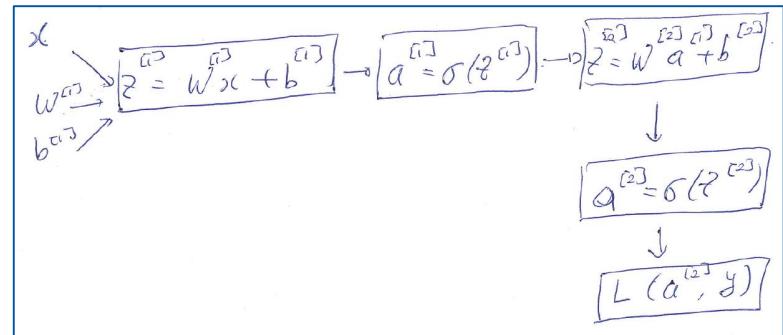
$$dA^{[1]} = dZ^{[2]} \cdot W^{[2]}$$

$$dZ^{[1]} = dA^{[1]} \cdot A^{[1]}(1-A^{[1]})$$

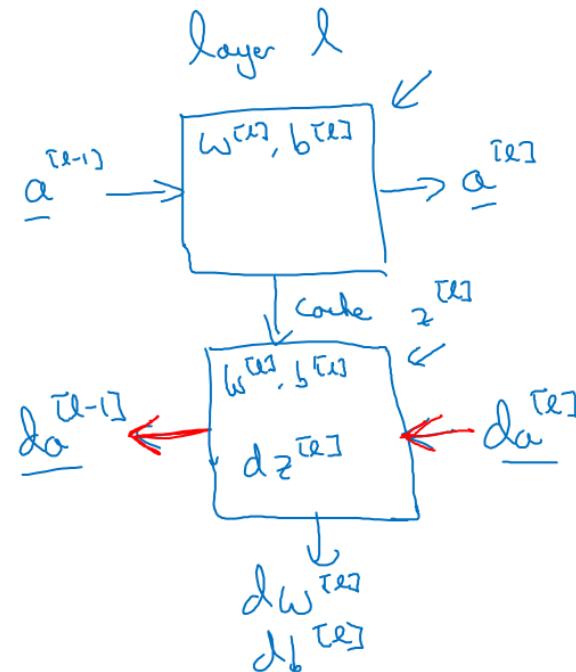
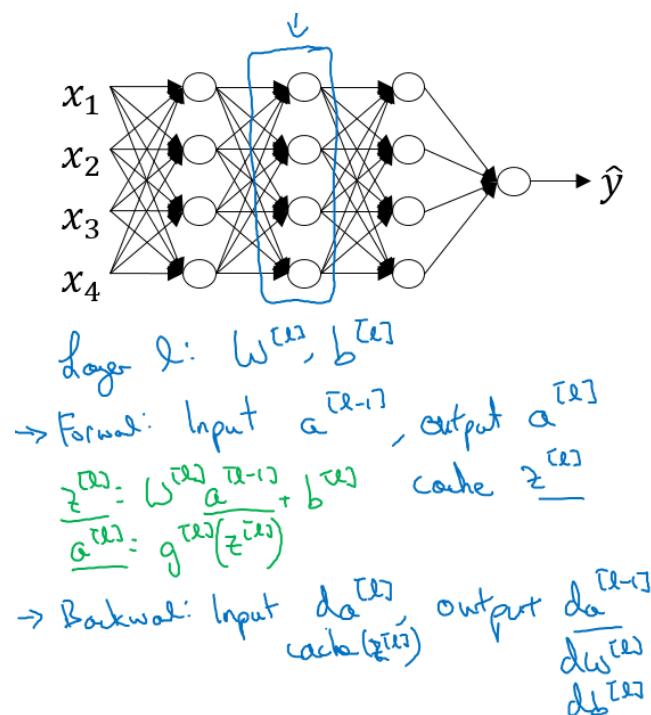
$$dW^{[1]} = \frac{1}{m} dZ^{[1]} \cdot A^{[0]}$$

$$= \frac{1}{m} dZ^{[1]} \cdot X$$

$$db^{[1]} = \frac{1}{m} \text{np.sum}(dZ^{[1]}, axis=1, keepdims=True)$$



Forward and backward functions



<deep learning, Andrew Ng>

Forward and backward propagation

$$Z^{[1]} = W^{[1]}X + b^{[1]}$$

$$A^{[1]} = g^{[1]}(Z^{[1]})$$

$$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(Z^{[2]})$$

⋮

$$A^{[L]} = g^{[L]}(Z^{[L]}) = \hat{Y}$$

$$dZ^{[L]} = A^{[L]} - Y$$

$$dW^{[L]} = \frac{1}{m} dZ^{[L]} A^{[L]T}$$

$$db^{[L]} = \frac{1}{m} np.\text{sum}(dZ^{[L]}, axis=1, keepdims=True)$$

$$dZ^{[L-1]} = dW^{[L]T} dZ^{[L]} g'^{[L]}(Z^{[L-1]})$$

$$\vdots$$
$$dZ^{[1]} = dW^{[L]T} dZ^{[2]} g'^{[1]}(Z^{[1]})$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} A^{[1]T}$$

$$db^{[1]} = \frac{1}{m} np.\text{sum}(dZ^{[1]}, axis=1, keepdims=True)$$

<deep learning, Andrew Ng>

Initialize weights

Neural network with hidden layers

Weights should be initialized randomly rather than to all zeros.

If weights are initialized to be zero, then, each neuron in the first hidden layer will perform the same computation. So even after multiple iterations of gradient descent, each neuron in the layer will be computing the same thing as other neurons.

<deep learning, Andrew Ng>

Initialize weights

Pitfall: all zero initialization

Assume that approximately half of the weights will be positive and half of them will be negative.

Set all the initial weights to zero.

If every neuron in the network computes the same output, then they will also all compute the same gradients during backpropagation and undergo the exact same parameter updates.

In other words, there is no source of asymmetry between neurons if their weights are initialized to be the same.

<cs231n, Stanford university>

Initialize weights

Small random numbers

Want the weights to be very close to zero, but not identically zero.

Assume that the neurons are all random and unique in the beginning, so they will compute distinct updates and integrate themselves as diverse parts of the full network.

Initialize the weights of the neurons to small numbers and refer to doing so as symmetry breaking.

$W = 0.01 * np.random.randn(D, H)$

`randn` samples from a zero mean, unit standard deviation Gaussian.

<cs231n, Stanford university>

Initialize weights

Logistic regression

Logistic regression does not have a hidden layer.

If you initialize the weights to zeros, the first example x fed in the logistic regression will output zero but the derivatives of the logistic regression depend on the input x (because there is no hidden layer) which is not zero.

So at the second iteration, the weights values follow x 's distribution and are different from each other if x is not a constant vector.

<deep learning, Andrew Ng>

Xavier initialization

When using sigmoid activation function

Multiplies weights by

$$\frac{\sqrt{2}}{\sqrt{n_{inputs} + n_{outputs}}}$$

<Hands-On ML, Aurelien Geron>

He initialization

When using ReLU activation function (and its variants, including the ELU activation)

Multiplies weights by

$$\sqrt{2} \frac{\sqrt{2}}{\sqrt{n_{inputs} + n_{outputs}}}$$

<Hands-On ML, Aurelien Geron>

Initialization for Hyperbolic tangent

When using tanh activation function

Multiplies weights by

$$4 \frac{\sqrt{2}}{\sqrt{n_{inputs} + n_{outputs}}}$$

Weights initialization

Activation function	Uniform distribution $[-r, r]$	Normal distribution
Logistic	$r = \sqrt{\frac{6}{n_{\text{inputs}} + n_{\text{outputs}}}}$	$\sigma = \sqrt{\frac{2}{n_{\text{inputs}} + n_{\text{outputs}}}}$
Hyperbolic tangent	$r = 4\sqrt{\frac{6}{n_{\text{inputs}} + n_{\text{outputs}}}}$	$\sigma = 4\sqrt{\frac{2}{n_{\text{inputs}} + n_{\text{outputs}}}}$
ReLU (and its variants)	$r = \sqrt{2}\sqrt{\frac{6}{n_{\text{inputs}} + n_{\text{outputs}}}}$	$\sigma = \sqrt{2}\sqrt{\frac{2}{n_{\text{inputs}} + n_{\text{outputs}}}}$

<Hands-On ML, Aurelien Geron>

L2 Regularization

Most common form of regularization

Penalize the squared magnitude of all parameters directly.

Heavily penalize peaky weight vectors and prefer diffuse weight vectors.

Network to use all of its inputs a little rather than some of its inputs a lot.

During gradient descent parameter update, every weight is decayed linearly towards zero.

<cs231n, Stanford university>

Logistic regression with L2 Regularization

Cost function

$$J = -\frac{1}{m} \sum_{i=1}^m [y^i \log(\hat{y}^i) + (1 - y^i) \log(1 - \hat{y}^i)] + \boxed{\frac{\lambda}{m} \sum_{j=1}^n W_j^2}$$

Gradient

$$\frac{\partial J}{\partial w_0} = -\frac{1}{m} \sum_{i=1}^m [(\hat{y}^i - y^i)x_j^i] \quad \text{for } j=0$$

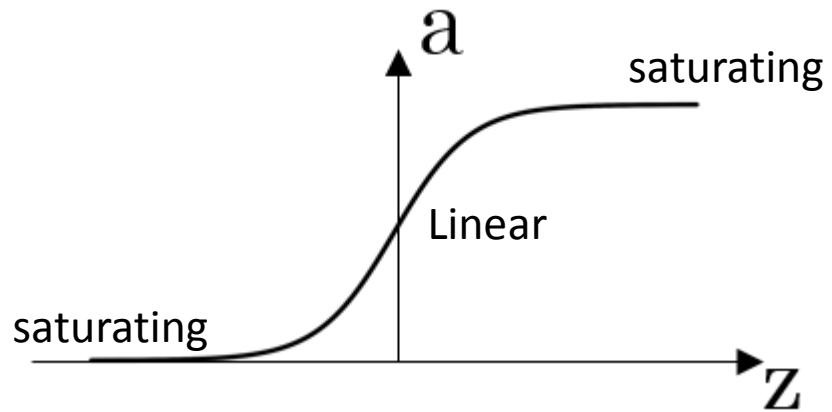
$$\frac{\partial J}{\partial w_j} = -\frac{1}{m} \sum_{i=1}^m [(\hat{y}^i - y^i)x_j^i] + \boxed{\frac{2\lambda}{m} W_j} \quad \text{for } j \geq 1$$

Neural network with L2 regularization

$$J = -\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \log(a^{[L](i)}) + (1 - y^{(i)}) \log(1 - a^{[L](i)}) \right)$$

$$J_{regularized} = \underbrace{-\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \log(a^{[L](i)}) + (1 - y^{(i)}) \log(1 - a^{[L](i)}) \right)}_{\text{cross-entropy cost}} + \boxed{\underbrace{\frac{1}{m} \frac{\lambda}{2} \sum_l \sum_k \sum_j W_{k,j}^{[l]2}}_{\text{L2 regularization cost}}}$$

Sigmoid

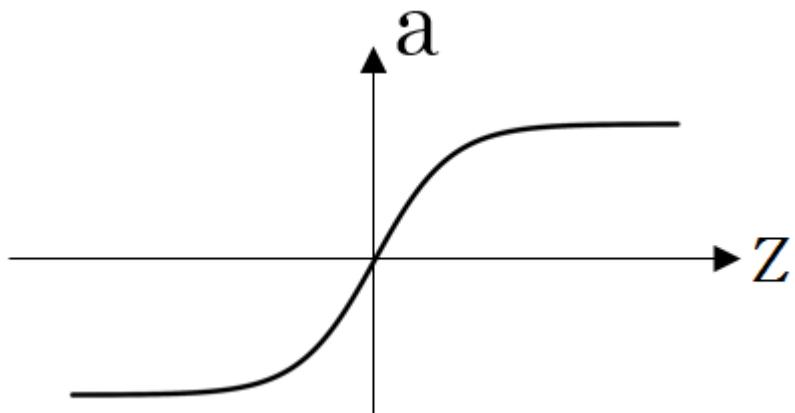


$$a = \sigma(z) = \frac{1}{1 + \exp(-z)}$$

- It can be saturated(0 or 1) when z is large.
- Gradients will be almost zero.
- The output is not zero-centered.

<Deep learning, Andrew Ng>

tanh



$$a = \tanh(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}$$

- The output is zero-centered.
- It is simply a scaled sigmoid function.

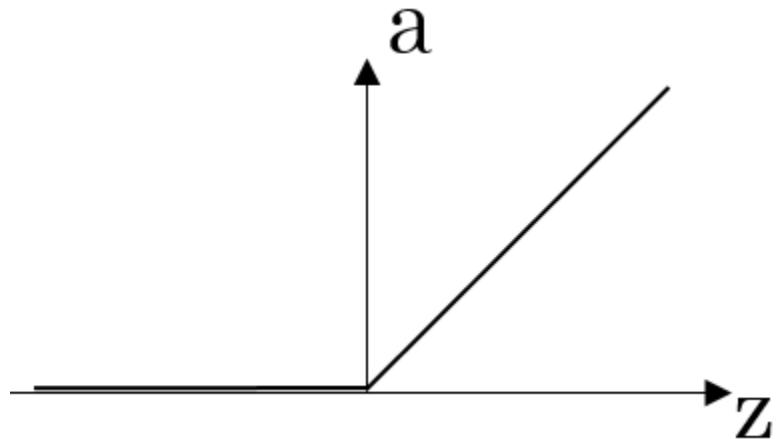
$$\tanh(z) = 2\sigma(2z) - 1$$

<Deep learning, Andrew Ng>

ReLU

Rectified Linear Unit

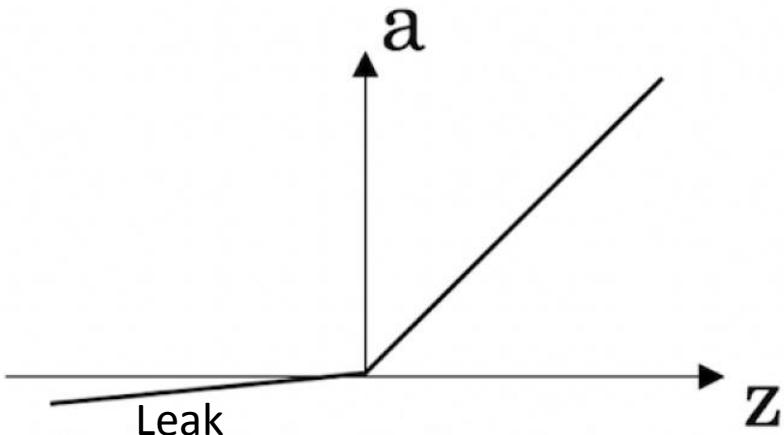
$$a = \text{ReLU}(z) = \max(0, z)$$



<Deep learning, Andrew Ng>

- It is zero when z is ≤ 0
- It is linear with slope 1 when z is > 0 .
- It can learn fast
- It is most commonly used in CNN.
- Some network can be dead since neuron will not activate across entire dataset.

Leaky ReLU



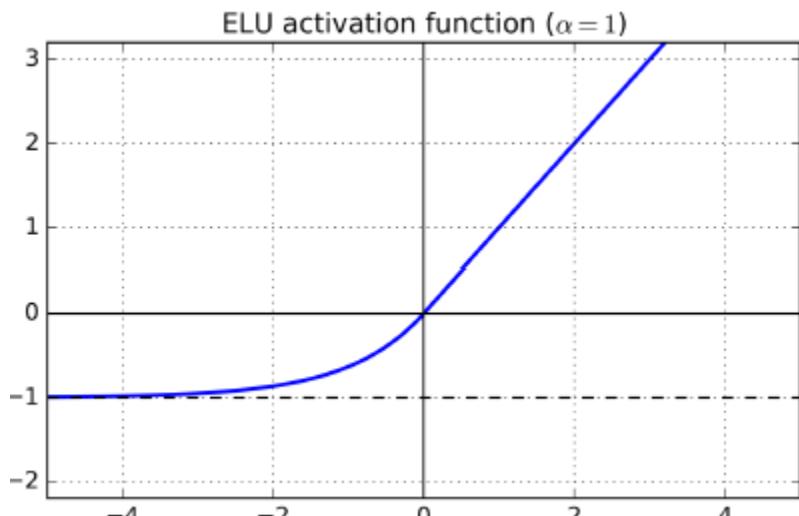
$$a = \text{Leaky ReLU}(z) = \max(0.01z, z)$$

- It is linear with small negative slope when z is ≤ 0
- It is linear with slope 1 when z is > 0 .

<Deep learning, Andrew Ng>

ELU

Exponential Linear Unit



<Hands-On ML, Aurelien Geron>

$$a = ELU(z) = \begin{cases} \alpha(\exp(z) - 1) & \text{if } z < 0 \\ z & \text{if } z \geq 0 \end{cases}$$

- Negative values when $z < 0 \rightarrow$ average output closer to 0
- Nonzero gradient for $z < 0 \rightarrow$ avoid dying units issue
- It is smooth everywhere including $z=0$

Derivatives of activation functions

Sigmoid function

$$g(z) = a = \frac{1}{1 + \exp(-z)}$$

$$\frac{d}{dz} g(z) = \frac{1}{1 + \exp(-z)} \left(1 - \frac{1}{1 + \exp(-z)} \right)$$

$$g'(z) = g(z)(1 - g(z))$$

$$g'(z) = a(1 - a)$$

Derivatives of activation functions

Tanh function

$$g(z) = a = \tanh(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}$$

$$\frac{d}{dz} g(z) = 1 - (\tanh(z))^2$$

$$g'(z) = 1 - g(z)^2$$

$$g'(z) = 1 - a^2$$

Derivatives of activation functions

ReLU and Leaky ReLU

ReLU

$$a = g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

Leaky ReLU

$$a = g(z) = \max(0.01z, z)$$

$$g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

Derivatives of activation functions

ELU

$$a = ELU(z) = \begin{cases} \alpha(\exp(z) - 1) & \text{if } z < 0 \\ z & \text{if } z \geq 0 \end{cases}$$

$$g'(z) = \begin{cases} \alpha \exp(z) & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

L2 regularization

Cost function

$$J = -\frac{1}{m} \sum_{i=1}^m [y^i \log(a^{[L](i)}) + (1 - y^i) \log(1 - a^{[L](i)})] + \frac{\lambda}{2m} \sum_{l=1}^L \|W^{[l]}\|^2$$

Gradient

$$dW^{[l]} = (\text{from backpropagation}) + \frac{\lambda}{m} W^{[l]}$$

$$W^{[l]} = W^{[l]} - \alpha dW^{[l]}$$

$$W^{[l]} = W^{[l]} - \alpha[(\text{from backprop}) + \frac{\lambda}{m} W^{[l]}]$$

$$W^{[l]} = W^{[l]} - \frac{\alpha\lambda}{m} W^{[l]} - \alpha(\text{from backprop})$$

Per-layer regularization:

It is not common to regularize different layers to different amount.

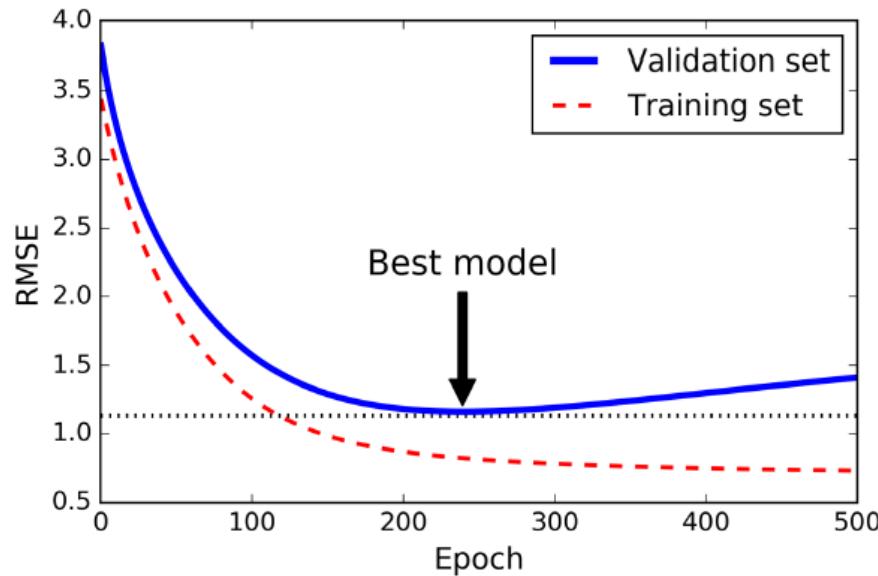
Elastic net regularization

Combine L1 regularization with L2 regularization

$$\lambda_1|W| + \lambda_2 W^2$$

Early Stopping

Stop training as the validation error reaches a minimum.



<Hands-On ML, Aurelien Geron>

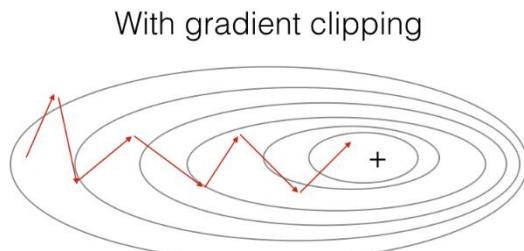
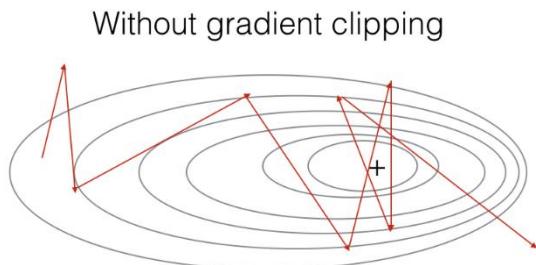
Max norm constraints

Enforce an absolute upper bound on the magnitude of the weight vector for every neuron.

$$\|W\|_2 < c$$

Gradient is clipping to prevent exploding gradients.

Every element of the gradient vector is clipped to lie between some range [-N,N].

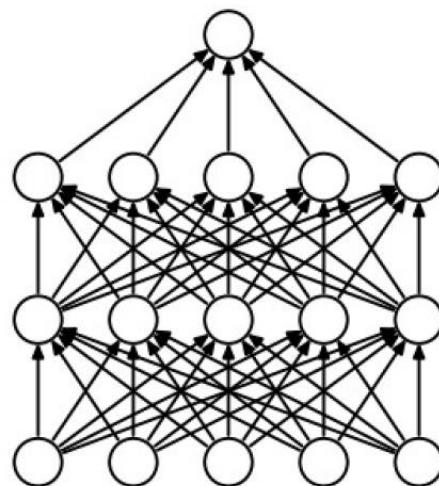


<Deep learning, Andrew Ng>

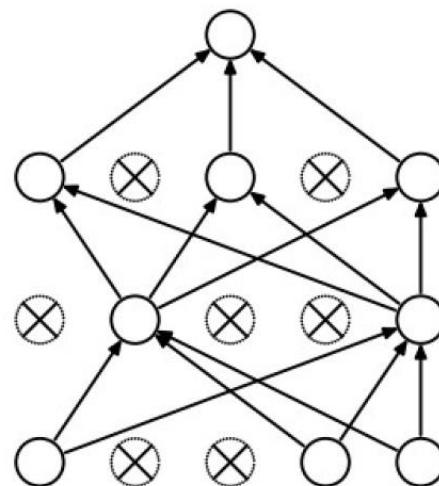
Drop out

$$px + (1 - p)0$$

Keep a neuron active with some probability p and set to zero others



(a) Standard Neural Net



(b) After applying dropout.

<cs231n, Stanford university>

Forward propagation with Drop out

1. Initialize matrix D1
2. Convert entries of D1 to 0 or 1(using keep prob as threshold)
3. Shut down some neurons of A1
4. Scale the value of neurons that haven't been shut down

```
keep_prob = 0.5  
D1 = np.random.rand(A1.shape[0],A1.shape[1])  
D1 = np.where(D1<keep_prob,1,0)  
A1 = np.multiply(D1,A1)  
A1 = A1/keep_prob
```

Backward propagation with dropout

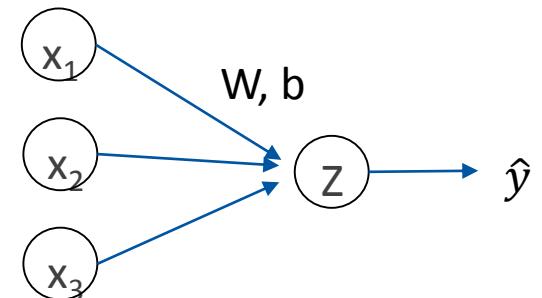
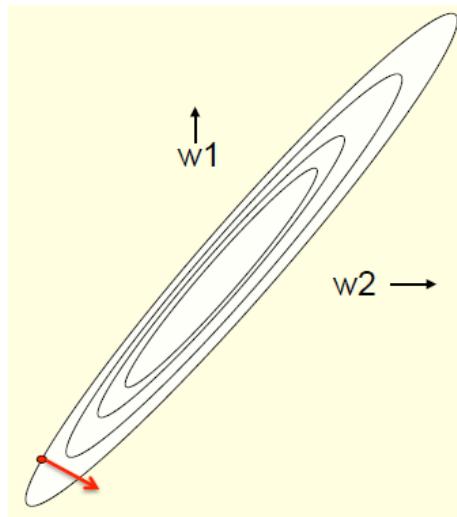
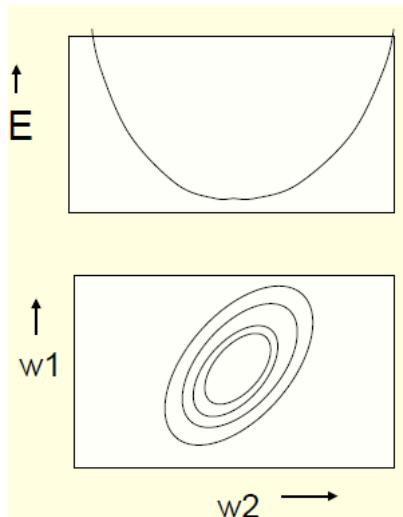
1. Apply mask D_1 to shut down the same neurons as during the forward propagation
2. Scale the value of neurons that haven't been shut down

```
dA1 = np.multiply(D1,dA1)
dA1 = dA1/keep_prob
dZ1 = np.multiply(dA1, np.int64(A1 > 0))
```

Normalizing inputs to speed up learning

In linear and logistic regression, we normalize input X to speed up learning.

$$X = (X - \mu)/\sigma$$

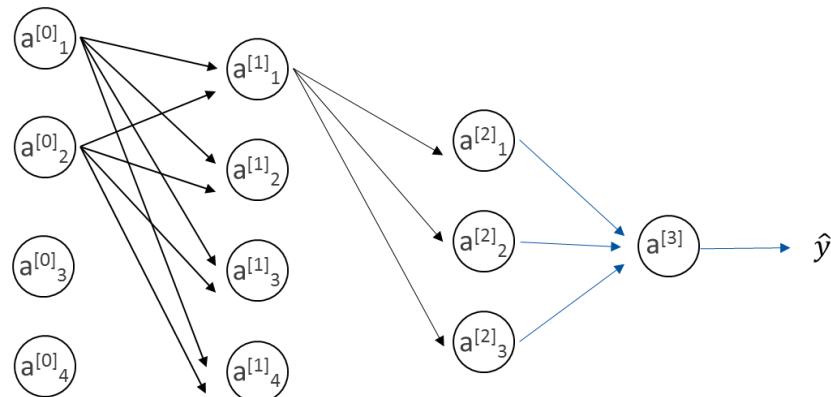


<Neural Networks, Geoffrey Hinton>

Normalizing activations in a network

In Neural Network, we may normalize $a^{[-1]}$ to train $W^{[l]}, b^{[l]}$ faster.

In actual, we normalize $Z^{[-1]}$ instead of $a^{[-1]}$.



Batch Normalization

Address the vanishing / exploding gradients problems

Address the problem that the distribution of each layer's inputs changes during training.

Adding an operation in the model just before the activation function of each layer, zero-centering and normalizing the inputs, then scaling and shifting the result using two new parameters per layer (one for scaling, the other for shifting)

Implementing Batch Norm

Given some intermediate values in NN,

$$\mu = \frac{1}{m} \sum_i Z^{(i)}$$

$$\sigma^2 = \frac{1}{m} \sum_i (Z^{(i)} - \mu)^2$$

$$Z_{norm}^{(i)} = \frac{Z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\tilde{Z}^{(i)} = \gamma Z_{norm}^{(i)} + \beta$$

If $\gamma = \sqrt{\sigma^2 + \epsilon}$

$$\beta = \mu$$

then

$$\tilde{Z}^{(i)} = Z_{norm}^{(i)}$$

ϵ is a smoothing term to avoid division by zero, a tiny number

γ is scaling parameter for a layer

β is shifting parameter for a layer

γ and β are learnable parameters for model

Adding Batch Norm to a network

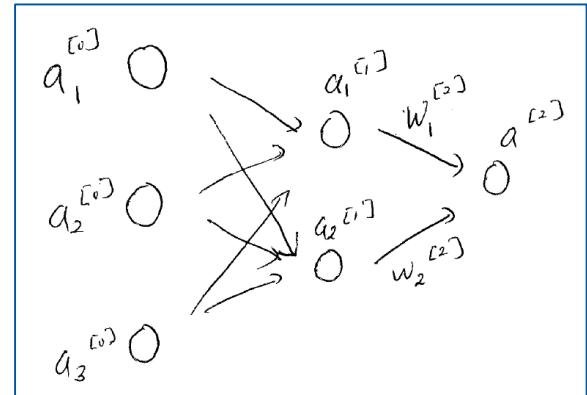
$$d^{[0]} \xrightarrow{w^{[0]}, b^{[0]}} z^{[1]} \xrightarrow{\beta^{[1]}, \gamma^{[1]}} \hat{z}^{[1]} \xrightarrow{\text{Batch Norm}} a^{[1]}$$

$$a^{[1]} \xrightarrow{w^{[1]}, b^{[1]}} z^{[2]} \xrightarrow{\beta^{[2]}, \gamma^{[2]}} \hat{z}^{[2]} \xrightarrow{\text{Batch Norm}} a^{[2]}$$

Parameters:

$$\left\{ w^{[0]}, b^{[0]}, w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]} \right. \\ \left. \beta^{[0]}, \gamma^{[0]}, \beta^{[1]}, \gamma^{[1]}, \dots, \beta^{[L]}, \gamma^{[L]} \right\}$$

$$\beta^{[l]} := \beta^{[l]} - \alpha d\beta^{[l]}$$



$$a_1^{[0]} = \sigma(z_1^{[0]})$$

$$z_1^{[0]} = w_{11}^{[0]} \cdot a_1^{[0]} + w_{12}^{[0]} \cdot a_2^{[0]} + w_{13}^{[0]} \cdot a_3^{[0]} + b_1^{[0]}$$

$$a_2^{[0]} = \sigma(z_2^{[0]})$$

$$z_2^{[0]} = w_{21}^{[0]} \cdot a_1^{[0]} + w_{22}^{[0]} \cdot a_2^{[0]} + w_{23}^{[0]} \cdot a_3^{[0]} + b_2^{[0]}$$

$$a^{[1]} = \sigma(z^{[1]})$$

$$z^{[1]} = w_1^{[1]} \cdot a_1^{[0]} + w_2^{[1]} \cdot a_2^{[0]} + b^{[1]}$$

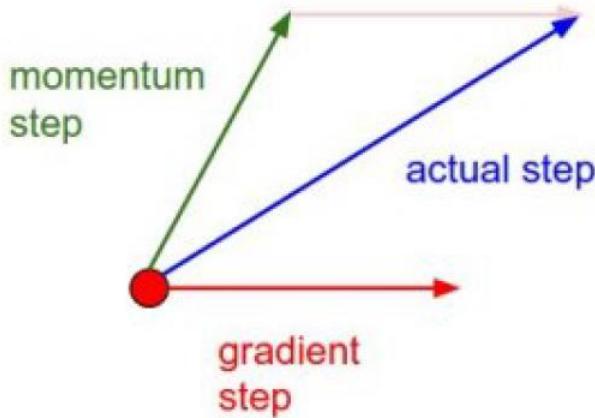
Pros & Cons for Batch Normalization

- Vanishing gradients problems are strongly reduced
- Less sensitive to the weight initialization
- Be able to use much larger learning rates
- Play as a regularization
- Add complexity to the model
- Require Extra computation

Momentum update

Gradient descent update

$$x := x - \text{learning rate} * dx$$



Momentum update

$$v = \beta * v - \text{learning rate} * dx$$

$$x := x + v$$

v : velocity

β : momentum

Adagrad

Adaptive learning rate method

$$v := v + dx^{**2}$$

$$x := x - \text{learning_rate} * dx / (\sqrt{v} + \text{eps})$$

v : sum of squared gradients

Normalized the parameter update step, element-wise.

Weights with high gradients will have effective learning rate reduces,

Weights with low gradients will have effective learning rate increases.

Gradient descent with momentum

$$V_{dW} = 0, V_{db} = 0$$

For every iteration t:

Compute V_{dW} , V_{db} using dW and db

$$V_{dW} := \beta_1 \cdot V_{dW} + (1 - \beta_1) \cdot dW$$

$$V_{db} := \beta_1 \cdot V_{db} + (1 - \beta_1) \cdot db$$

$$W := W - \alpha \cdot V_{dW}$$

$$b := b - \alpha \cdot V_{db}$$

RMS prop

Root Mean square prop

$$S_{dW} \neq 0, S_{db} \neq 0, \epsilon = 10^{-8}$$

For every iteration t:

Compute S_{dW} , S_{db} using dW and db

$$S_{dW} := \beta_2 \cdot S_{dW} + (1 - \beta_2) \cdot dW^2$$

$$S_{db} := \beta_2 \cdot S_{db} + (1 - \beta_2) \cdot db^2$$

$$W := W - \alpha \cdot \frac{dW}{\sqrt{S_{dW} + \epsilon}}$$

$$b := b - \alpha \cdot \frac{db}{\sqrt{S_{db} + \epsilon}}$$

Adam Optimization

Adaption moment estimation

$$V_{dW} = 0, V_{db} = 0, S_{dW} = 0, S_{db} = 0, \epsilon = 10^{-8}$$

For every iteration t:

Compute V_{dW} , V_{db} , S_{dW} , S_{db} using dW and db

$$V_{dW} := \beta_1 \cdot V_{dW} + (1 - \beta_1) \cdot dW$$

$$V_{db} := \beta_1 \cdot V_{db} + (1 - \beta_1) \cdot db$$

$$S_{dW} := \beta_2 \cdot S_{dW} + (1 - \beta_2) \cdot dW^2$$

$$S_{db} := \beta_2 \cdot S_{db} + (1 - \beta_2) \cdot db^2$$

$$W := W - \alpha \cdot V_{dW}$$

$$b := b - \alpha \cdot V_{db}$$

Adam Optimization

Bias correction

$$V_{dW}^{cor} = \frac{V_{dW}}{1 - \beta_1^t}$$

$$V_{db}^{cor} = \frac{V_{db}}{1 - \beta_1^t}$$

$$S_{dW}^{cor} = \frac{S_{dW}}{1 - \beta_2^t}$$

$$S_{db}^{cor} = \frac{S_{db}}{1 - \beta_2^t}$$

$$W := W - \alpha \cdot \frac{V_{dW}^{cor}}{\sqrt{S_{dW}^{cor} + \varepsilon}}$$

$$b := b - \alpha \cdot \frac{V_{db}^{cor}}{\sqrt{S_{db}^{cor} + \varepsilon}}$$

Adam optimization

Hyper parameters choices

α : need to be tune

β_1 : 0.9 (momentum)

β_2 : 0.999 (RMS prop)

ϵ : 10^{-8}

Learning rate decay

Slowing decreasing α

If learning rate is large, parameters decay too aggressively and the parameters are unable to reach the local minimum.

1 epoch : 1 pass through data set, iteration number

- $1/t$ decay
- Exponentially decay
- Step decay: discrete staircase
- Manual decay

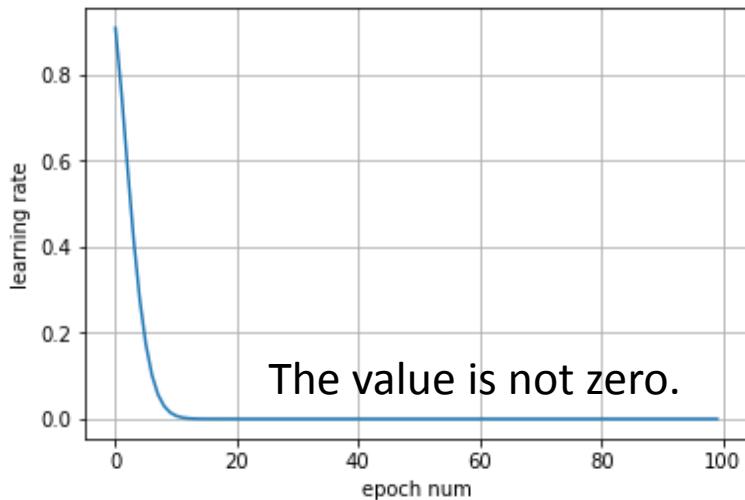
1/t decay

$$\alpha = \frac{\alpha_0}{1 + decay-rate \times epoch-num}$$

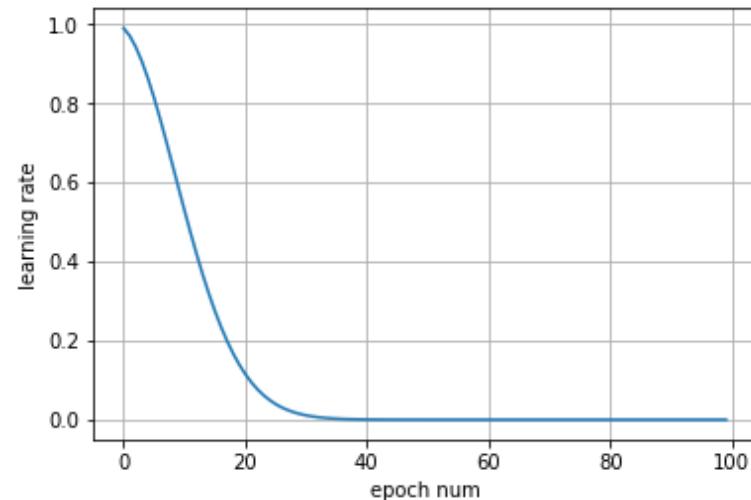
$$\alpha = \frac{\alpha_0}{1 + kt}$$

$$\alpha_0 = 1$$

decay rate = 0.1



decay rate = 0.01

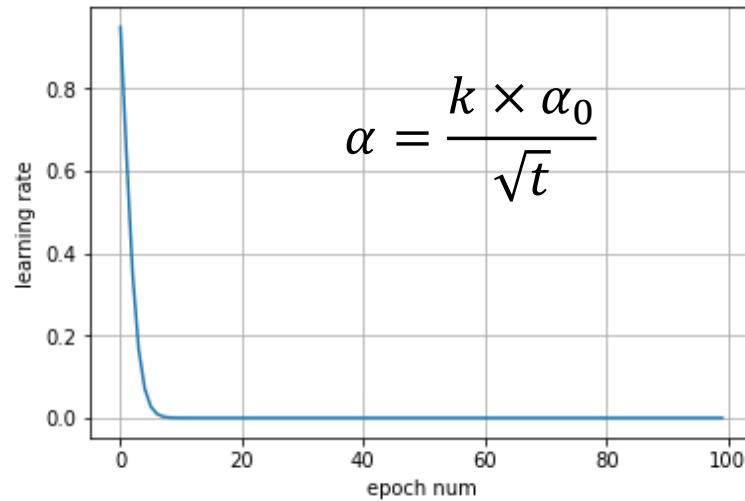


1/t decay

$$\alpha = \frac{k \times \alpha_0}{\sqrt{epoch - num}}$$

$$\alpha_0 = 1$$

$$k = 0.95$$

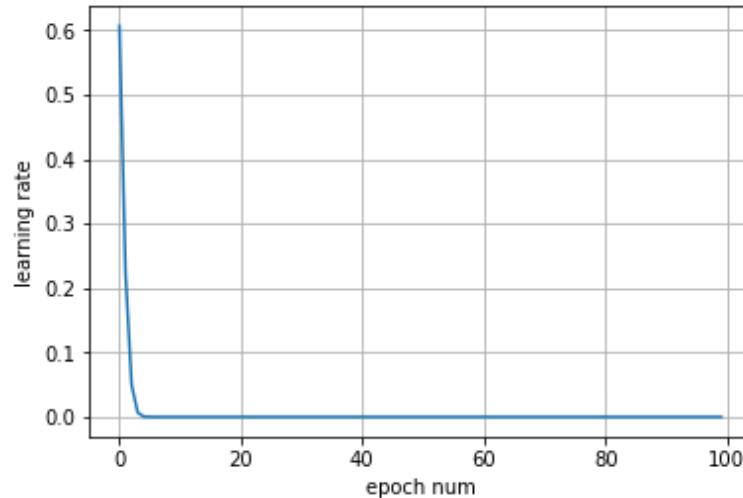


Exponential decay

$$\alpha = \alpha_0 \cdot \exp(-kt)$$

$$\alpha_0 = 1$$

$$k = 0.5$$

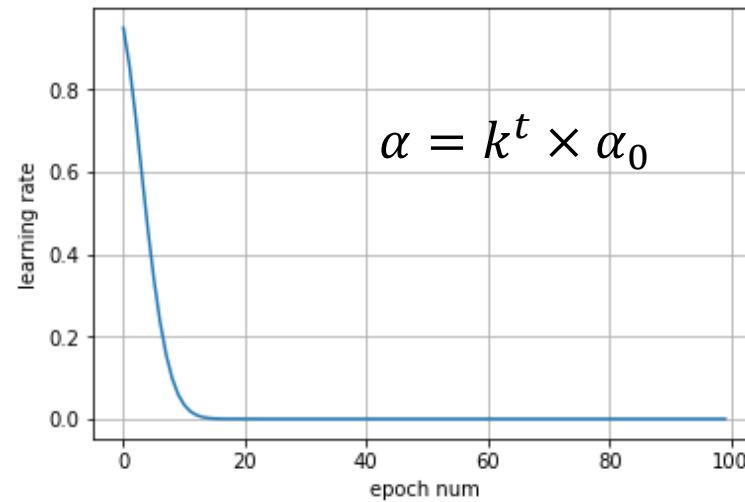


Exponentially decay

$$\alpha = k^{epoch_num} \times \alpha_0$$

$$\alpha_0 = 1$$

$$k = 0.95$$



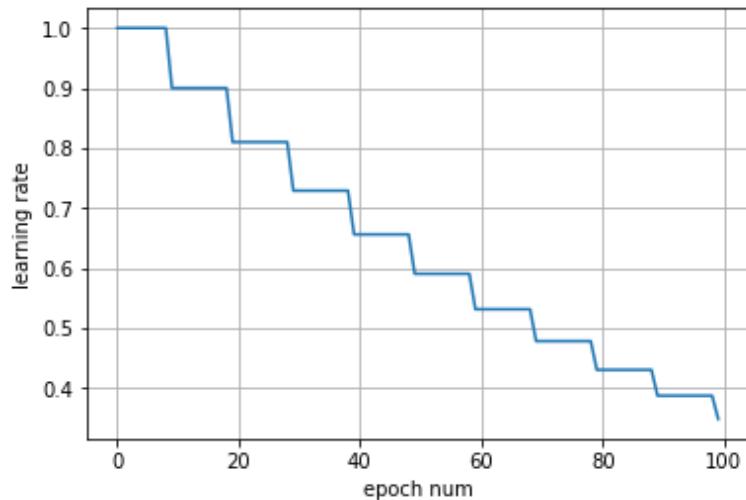
Discrete staircase

$\alpha = k \times \alpha$ with every few epochs

$$\alpha_0 = 1$$

$$k=0.9$$

every 10 epochs



Learning rate decay

Annealing learning rate

Hyper parameter k: [0.01, 0.1, 0.5, 0.9, 0.95, 0.99]

$$\alpha = \frac{\alpha_0}{1 + kt}$$

$$\alpha = \alpha_0 \cdot \exp(-kt)$$

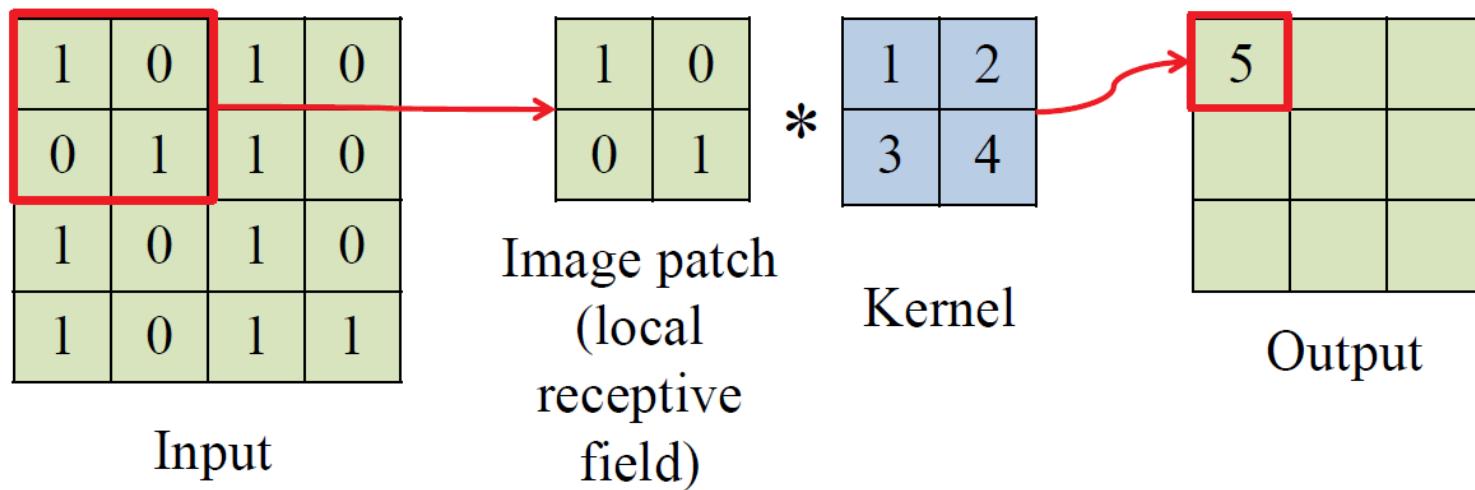
$$\alpha = \frac{k \times \alpha_0}{\sqrt{t}}$$

$$\alpha = k^t \times \alpha_0$$

$\alpha = k \times \alpha$ with every t

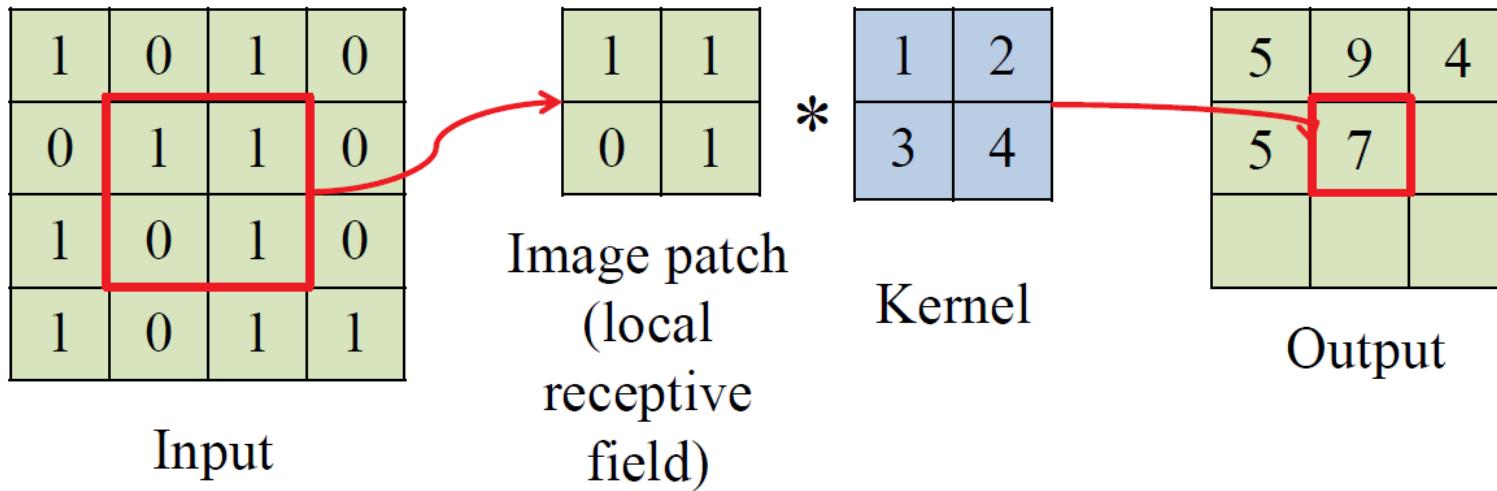
Convolutions

Convolution is a dot product of a **kernel** (or filter) and a patch of an image (**local receptive field**) of the same size



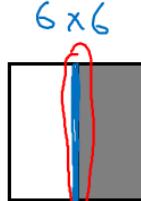
<Intro to deep learning, national research university>

Convolutions



Vertical edge detection

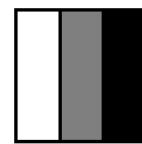
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0



$$\begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array}$$

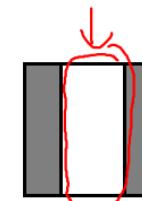
3×3

*



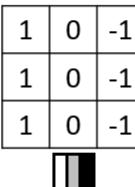
0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0

$\uparrow 4 \times 4$



<Deep Neural Network, Andrew Ng>

Vertical edge detection

$$\begin{array}{ccccccc} 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 & 0 \end{array}$$


*

$$\begin{array}{cccc} 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \end{array}$$


$$\begin{array}{ccccccc} 0 & 0 & 0 & 10 & 10 & 10 & 0 \\ 0 & 0 & 0 & 10 & 10 & 10 & 0 \\ 0 & 0 & 0 & 10 & 10 & 10 & 0 \\ 0 & 0 & 0 & 10 & 10 & 10 & 0 \\ 0 & 0 & 0 & 10 & 10 & 10 & 0 \\ 0 & 0 & 0 & 10 & 10 & 10 & 0 \end{array}$$


*

$$\begin{array}{cccc} 0 & -30 & -30 & 0 \\ 0 & -30 & -30 & 0 \\ 0 & -30 & -30 & 0 \\ 0 & -30 & -30 & 0 \end{array}$$


Horizontal edge detection

10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10

6×6



*

1	1	1
0	0	0
-1	-1	-1

=

0	0	0	0
30	10	-10	-30
30	10	-10	-30
0	0	0	0

<Deep Neural Network, Andrew Ng>

Edge detection



Original
image

Kernel

$$\ast \begin{array}{|c|c|c|} \hline -1 & -1 & -1 \\ \hline -1 & 8 & -1 \\ \hline -1 & -1 & -1 \\ \hline \end{array} = \begin{array}{c} \text{Edge} \\ \text{detection} \end{array}$$
A high-contrast, black-and-white image of the same dog's head, where the edges are highlighted in white, indicating the boundaries between different regions of the face.

Sums up to 0 (black color)
when the patch is a solid fill

<Intro to deep learning, national research university>

Sharpening

Original image 

Kernel

$$\begin{matrix} * & \begin{matrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{matrix} & = & \begin{matrix} \text{Edge detection image} \end{matrix} \end{matrix}$$

Edge detection

$$\begin{matrix} * & \begin{matrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{matrix} & = & \begin{matrix} \text{Sharpening image} \end{matrix} \end{matrix}$$

Sharpening

Doesn't change an image for solid fills
Adds a little intensity on the edges

<Intro to deep learning, national research university>

Blurring



Original
image

$$\begin{array}{c} \text{Kernel} \\ \hline \begin{matrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{matrix} \\ * = \end{array} \quad \begin{array}{c} \text{Edge} \\ \text{detection} \\ \hline \text{Resulting edge-detected image} \end{array}$$
$$\begin{array}{c} \text{Kernel} \\ \hline \begin{matrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{matrix} \\ * = \end{array} \quad \begin{array}{c} \text{Sharpening} \\ \hline \text{Resulting sharpened image} \end{array}$$
$$\begin{array}{c} \text{Kernel} \\ \hline \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix} \\ * \frac{1}{9} = \end{array} \quad \begin{array}{c} \text{Blurring} \\ \hline \text{Resulting blurred image} \end{array}$$

<Intro to deep learning, national research university>

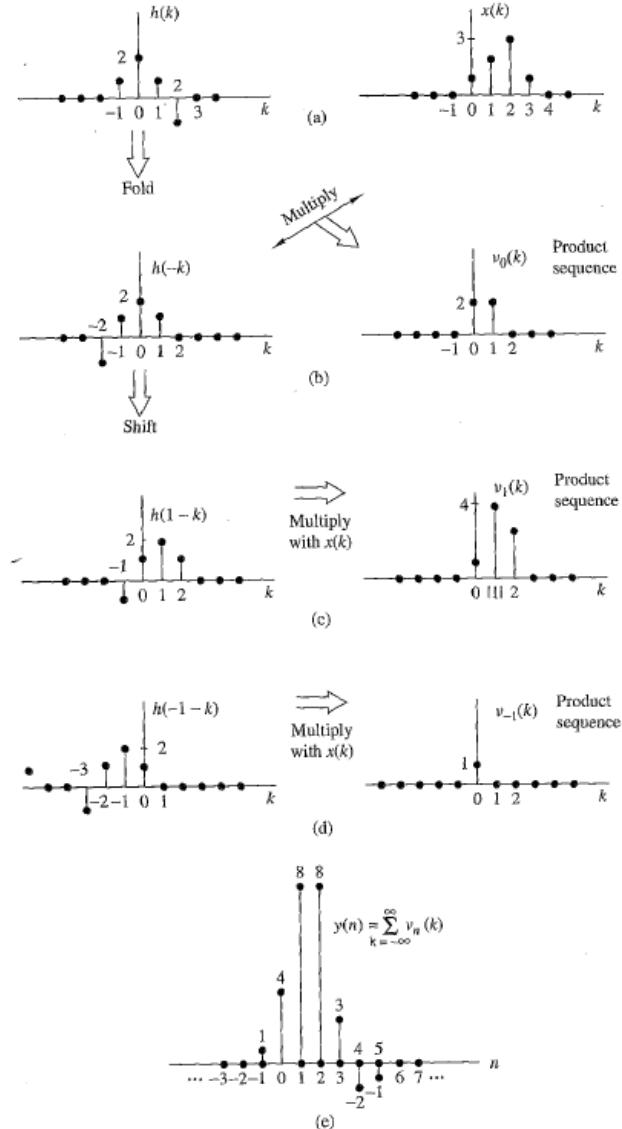
Convolution

$$y(n) = \sum_{k=-\infty}^{\infty} x(k)h(n-k)$$

<Digital signal processing, J.G., Proakis et al.>

Convolution is a mathematical operation that slides one function over another and measures the integral of their pointwise multiplication. It is heavily used in signal processing.

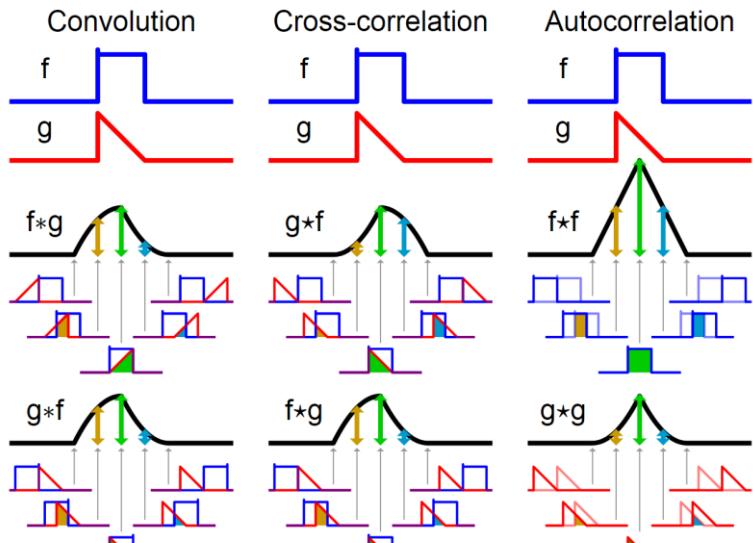
<Hands-On ML, A. Geron>



Convolution

Convolutional layers actually use cross-correlations.

<Hands-On ML, A. Geron>



<Wikipedia>

Convolution as correlation

$$\begin{array}{|c|c|c|c|} \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 1 & 0 \\ \hline 0 & 0 & 0 & 1 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 1 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 0 & 0 & 0 \\ \hline 0 & 1 & 0 \\ \hline 0 & 0 & 2 \\ \hline \end{array}$$

Input Kernel Output

$$\begin{array}{|c|c|c|c|} \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 1 \\ \hline 0 & 0 & 1 & 0 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 1 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 0 & 0 & 0 \\ \hline 0 & 0 & 1 \\ \hline 0 & 1 & 0 \\ \hline \end{array}$$

Input Kernel Output

Convolution as translation

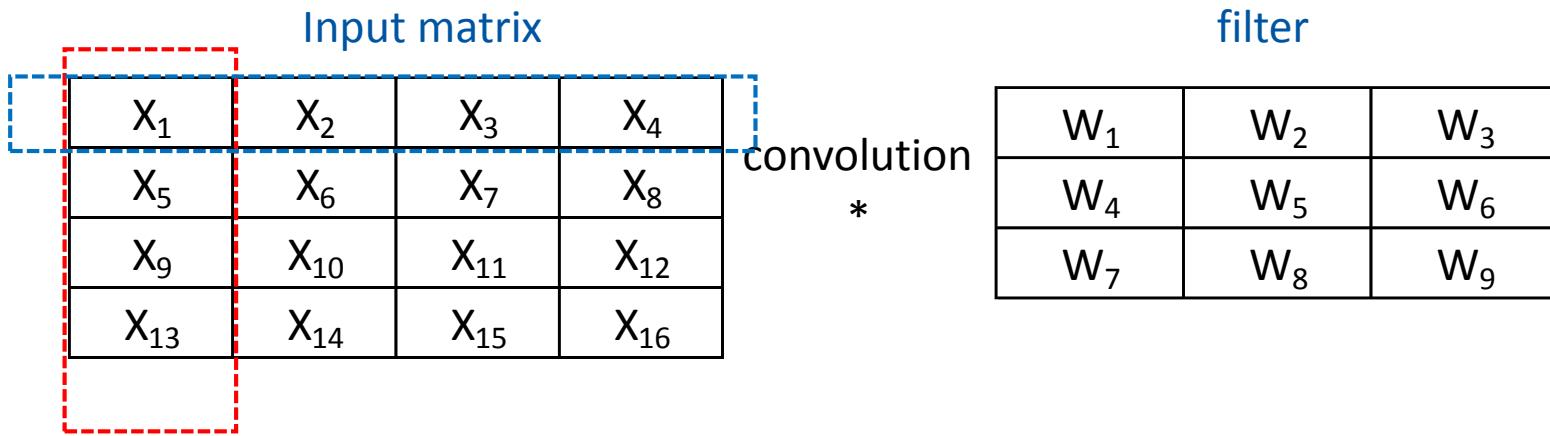
$$\begin{array}{|c|c|c|c|} \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 1 & 0 \\ \hline 0 & 0 & 0 & 1 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 1 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 0 & 0 & 0 \\ \hline 0 & 1 & 0 \\ \hline 0 & 0 & 2 \\ \hline \end{array}$$

Input Kernel Output

$$\begin{array}{|c|c|c|c|} \hline 1 & 0 & 0 & 0 \\ \hline 0 & 1 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 1 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 2 & 0 & 0 \\ \hline 0 & 1 & 0 \\ \hline 0 & 0 & 0 \\ \hline \end{array}$$

Input Kernel Output

Edge pixels?



What about information from edge pixels?

Padding

0	0	0	0	0	0
0	X_1	X_2	X_3	X_4	0
0	X_5	X_6	X_7	X_8	0
0	X_9	X_{10}	X_{11}	X_{12}	0
0	X_{13}	X_{14}	X_{15}	X_{16}	0
0	0	0	0	0	0

*

W_1	W_2	W_3
W_4	W_5	W_6
W_7	W_8	W_9

=

Z_1	Z_2	Z_3	Z_4
Z_5	Z_6	Z_7	Z_8
Z_9	Z_{10}	Z_{11}	Z_{12}
Z_{13}	Z_{14}	Z_{15}	Z_{16}

Padding

To detect the information of edge or border

To prevent shrinking(compression) of size of next layer

- Valid convolution: without padding

$$(n \times n) * (f \times f) \rightarrow (n-f+1) \times (n-f+1)$$

$$(6 \times 6) * (3 \times 3) \rightarrow (4 \times 4)$$

- Same convolution: with zero padding, output size is the same as the input size

$$(n \times n) * (f \times f) \rightarrow (n+2p-f+1) \times (n+2p-f+1)$$

$$(6 \times 6) * (3 \times 3) \rightarrow (6 \times 6)$$

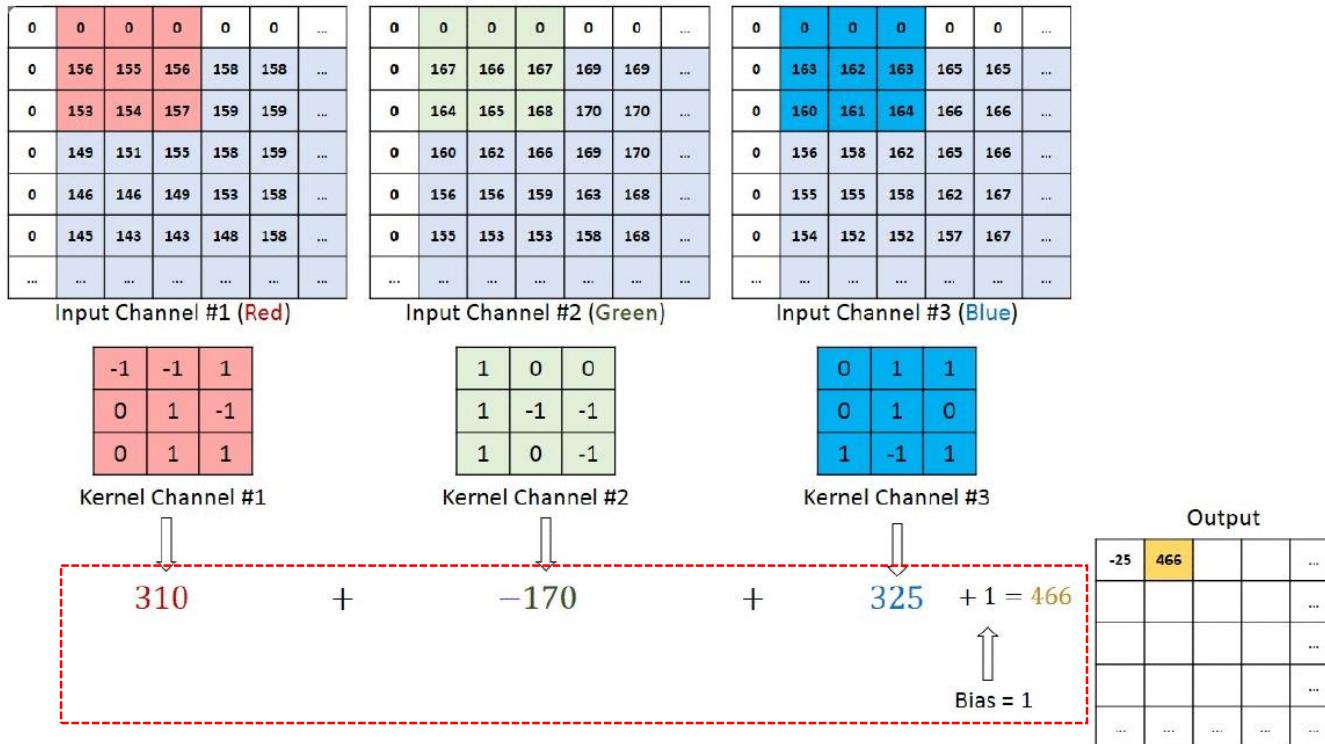
Summary of convolutions

- $n \times n$ image
- $f \times f$ filter
- padding p
- stride s

Output size:

$$\left[\frac{n + 2p - f}{s} + 1 \right] \times \left[\frac{n + 2p - f}{s} + 1 \right]$$

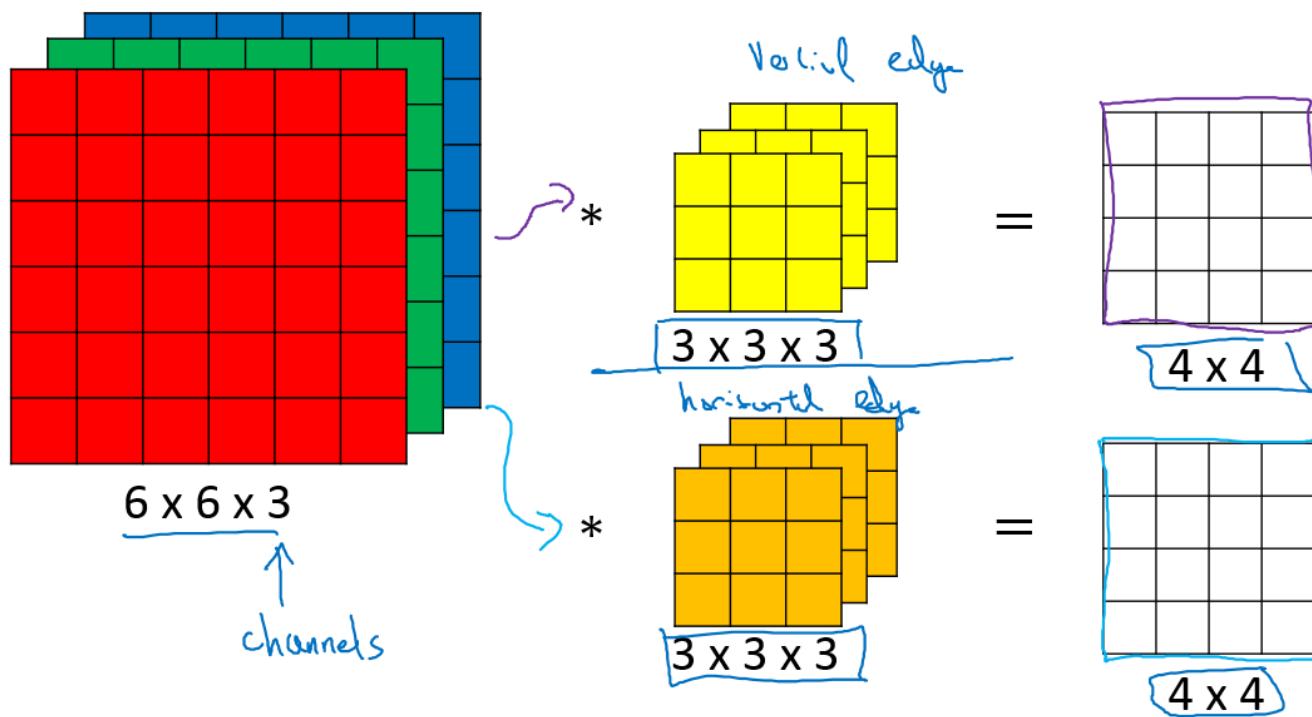
Convolutions on RGB images



<understanding convolutional layers in convolutional neural networks>

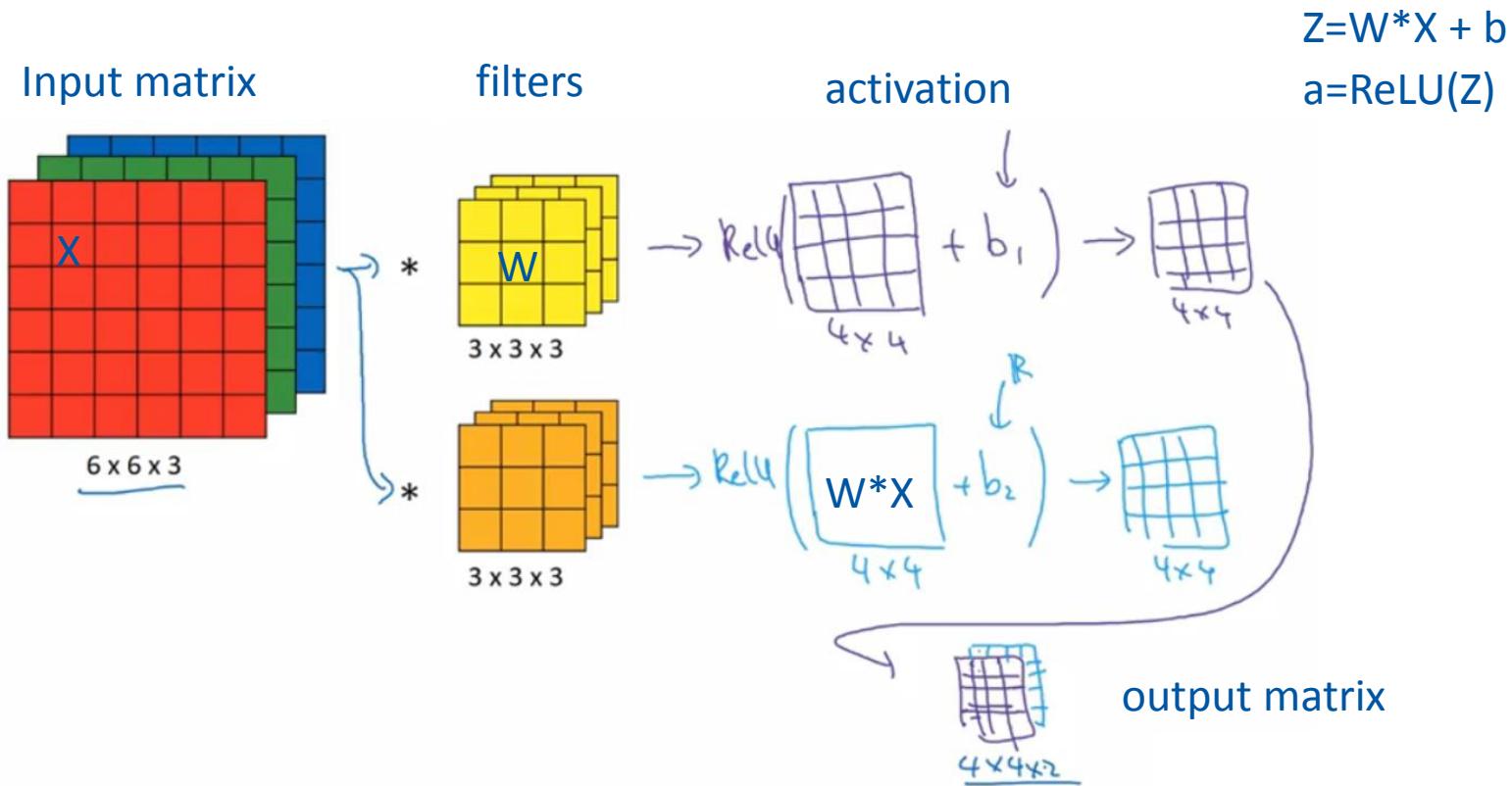
Machinelearningguru.com/computer_vision/basics/convolution/convolution_layer.html

Multiple filters (channels)



<deep learning, Andrew Ng>

Multiple filters (channels)



<deep learning, Andrew Ng>

Number of parameters in one layer

If there are 10 filters that are $3 \times 3 \times 3$ in one layer of a neural network, how many parameters does that layer have?

For 1 filter:

$3 \times 3 \times 3 = 27$ for W

1 for bias

For 10 filters:

Then, $28 * 10 = 280$

Summary of size in matrix

If layer l is a convolution layer:

$f^{[l]}$ = filter size

$p^{[l]}$ = padding

$s^{[l]}$ = stride

Weights: $f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]}$

bias: $1 \times 1 \times 1 \times n_c^{[l]}$

Input size: $n_H^{[l-1]} \times n_W^{[l-1]} \times n_c^{[l-1]}$

output size: $n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$

$$n_H^{[l]} = \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1$$

$$n_W^{[l]} = \frac{n_W^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1$$

Pooling layers

No parameters to learn

- To subsample (shrink) the input image in order to reduce the computational load, the memory usage, the number of parameters (reduce overfitting)
- Max Pooling : more commonly use
- Average Pooling

Pooling

2	3	1	9
4	7	3	5
8	2	2	2
1	3	4	5

Hyperparameters

$$f=2$$

$$s=2$$

Max:

7	9
8	5

Average:

4	4.5
3.25	3.25

<Deep Learning, Andrew Ng>

Size of Pooling

- Hyperparameters
 - f : filter size
 - s : stride
 - p : padding (usually $p=0$)
- Max or average pooling
- No parameters to learn

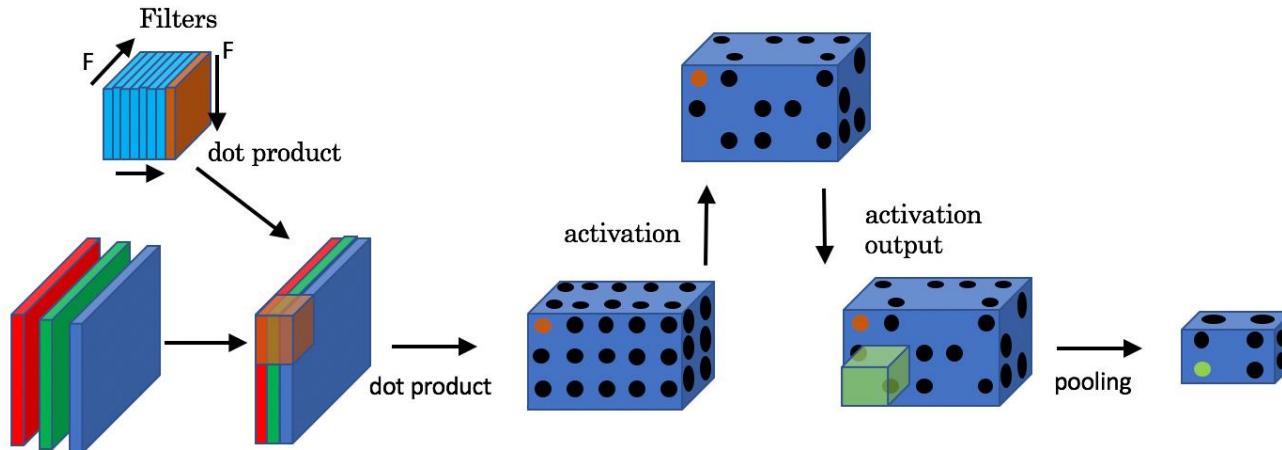
$$\text{Input size: } n_H^{[l-1]} \times n_W^{[l-1]} \times n_c^{[l-1]}$$

$$\text{output size: } n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$$

$$n_H^{[l]} = \frac{n_H^{[l-1]} - f^{[l]}}{s^{[l]}} + 1$$

$$n_W^{[l]} = \frac{n_W^{[l-1]} - f^{[l]}}{s^{[l]}} + 1$$

Convolution and Pooling



Input

Convolution

Take each filter
And convolve around the
image. Store the result in the
conv output.

Conv Output

Each point in the matrix
represents a dot product
from the convolution layer.

Pooling

We pool after the
activation to reduce the
number of parameters

Pooling Output

Each colored point
represents a different class.
In this case C = 4.

<Deep Learning, Andrew Ng>

Summary of Pooling

- Hyperparameters
 - f : filter size
 - s : stride
 - p : padding (usually $p=0$)
- Max or average pooling
- No parameters to learn

$$\text{Input size: } n_H^{[l-1]} \times n_W^{[l-1]} \times n_c^{[l-1]}$$

$$\text{output size: } n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$$

$$n_H^{[l]} = \frac{n_H^{[l-1]} - f^{[l]}}{s^{[l]}} + 1$$

$$n_W^{[l]} = \frac{n_W^{[l-1]} - f^{[l]}}{s^{[l]}} + 1$$

Forward propagation step

convolution activation pooling convolution activation pooling
X → Z1 → A1 → P1 → Z2 → A2 → P2

W1 W2
b1 b2

flatten linear activation linear activation
→ A3 → Z4 → A4 → Z5 → A5
W4 W5
b4 b5

Backward propagation step

The diagram illustrates the backward pass of a neural network, showing the flow of gradients (dZ) from the output layer back to the input layer. The layers are represented by arrows pointing left, indicating the direction of gradient flow.

- Layer 1:** convolution → activation → pooling
- Layer 2:** convolution → activation → pooling
- Layer 3:** flatten → linear → activation
- Layer 4:** linear → activation

Red labels indicate specific gradient components:

- Layer 1:** dW1, db1
- Layer 2:** dW2, db2
- Layer 3:** dW4, db4
- Layer 4:** dW5, db5

Pooling backward

average

```
dA := dP / (distribute_value) * np.ones of the shape
```

<Deep Learning, Andrew Ng>

$dP = 1$ (scalar)

shape : (n_H, n_W) of the output matrix

distribute_value = $n_H * n_W$

```
In [15]: dP = 1
```

```
In [24]: n_H, n_W = (2,2)
```

```
In [25]: distribute_value = n_H*n_W
```

```
In [26]: dA = dP*np.ones((n_H,n_W))/float(distribute_value)
```

```
In [27]: dA
```

```
Out[27]: array([[ 0.25,  0.25],  
                 [ 0.25,  0.25]])
```

Pooling backward

max :

$$dA := \text{mask} * dP$$

<Deep Learning, Andrew Ng>

```
In [28]: A = np.arange(4).reshape((2, 2))
A
Out[28]: array([[0, 1],
 [2, 3]])
```

Pooling forward

```
In [31]: P = np.max(A)
P
Out[31]: 3
```

```
In [29]: mask = (A==np.max(A))
mask
Out[29]: array([[False, False],
 [False, True]], dtype=bool)
```

Pooling backward

$dP=1$

```
In [33]: dA = mask*dP
dA
Out[33]: array([[0, 0],
 [0, 1]])
```

Convolution layer backward

dA

<Deep Learning, Andrew Ng>

$$dA+ = \sum_{h=0}^{n_H} \sum_{w=0}^{n_W} W_c \times dZ_{hw}$$

```
da [vert_start:vert_end, horiz_start:horiz_end, :] += W[:, :, :, c] * dZ[i, h, w, c]
```

dW

$$dW_c+ = \sum_{h=0}^{n_H} \sum_{w=0}^{n_W} a_{slice} \times dZ_{hw}$$

```
dW[:, :, :, c] += a_slice * dZ[i, h, w, c]
```

Convolution layer backward

db

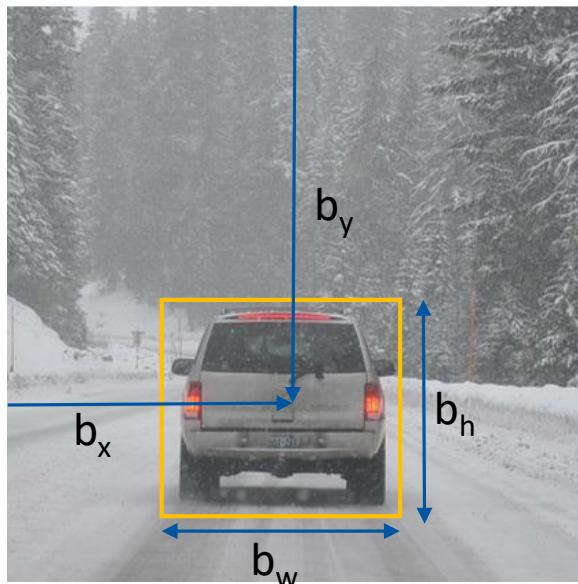
<Deep Learning, Andrew Ng>

$$db = \sum_h \sum_w dZ_{hw}$$

```
db[:, :, :, c] += dZ[i, h, w, c]
```

Classification with localization

(0,0)



(1,1)

Bounding box

$$b_x = 0.5$$

$$b_y = 0.7$$

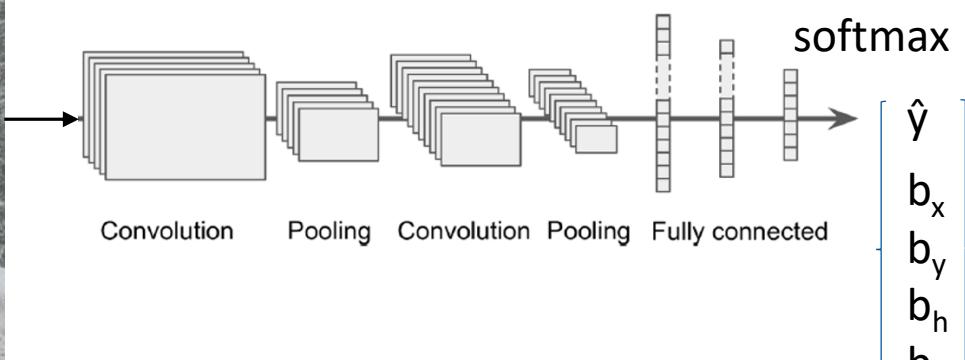
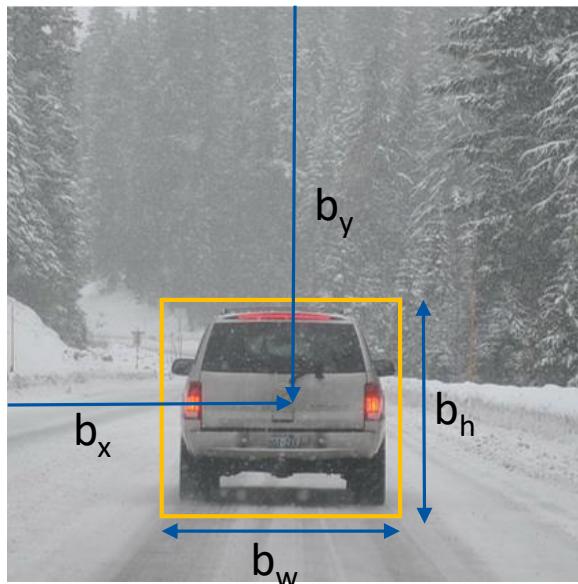
$$b_h = 0.3$$

$$b_w = 0.4$$

<deep learning, Andrew Ng>

Classification with localization

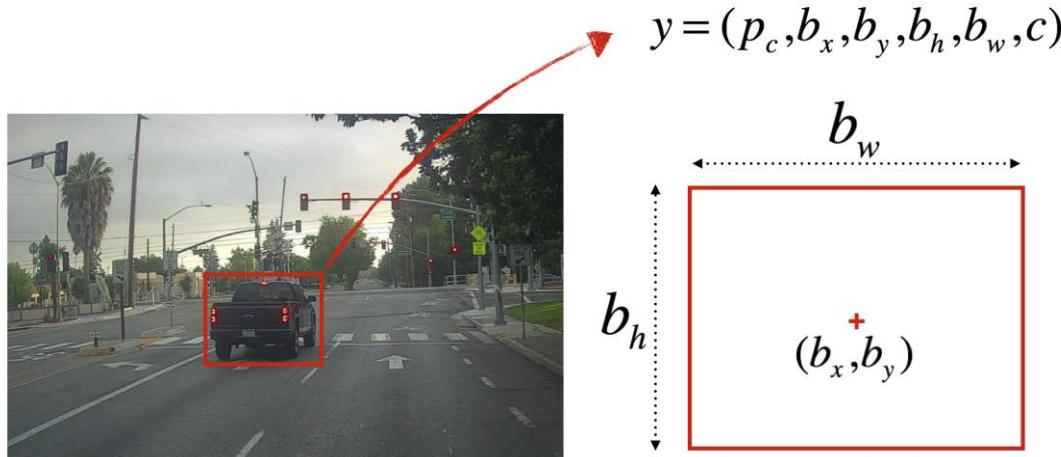
(0,0)



<deep learning, Andrew Ng>

(1,1)

Example of bounding box



$p_c = 1$: confidence of an object being present in the bounding box

$c = 3$: class of the object being detected (here 3 for “car”)

<deep learning, Andrew Ng>

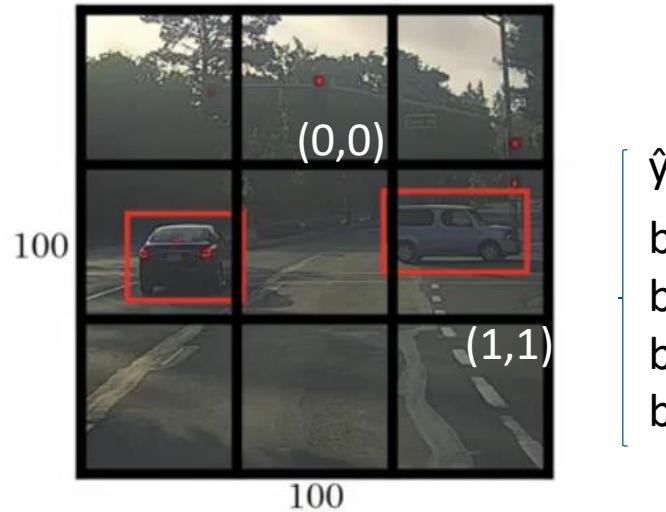
Sliding windows detection



Instead of one size,
Many different size of windows
(from small one to larger one) are
needed to apply.

<Deep Learning, Andrew Ng>

Specify bounding boxes



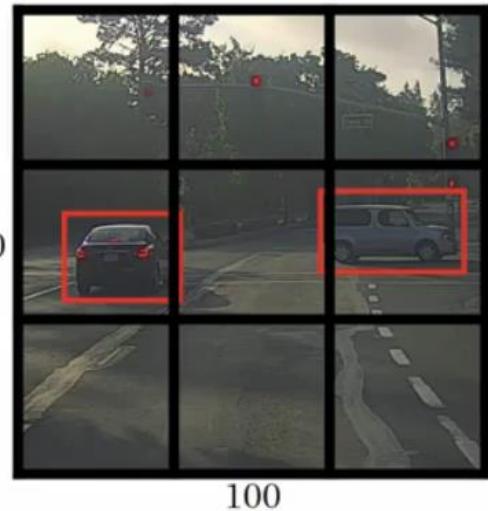
$$\begin{bmatrix} \hat{y} \\ b_x \\ b_y \\ b_h \\ b_w \end{bmatrix}$$

$$\begin{aligned} 0 < b_x, b_y < 1 \\ 0 < b_h, b_w \end{aligned}$$

<Deep Learning, Andrew Ng>

Redmon et al., 2015, You Only Look Once: Unified real-time object detection

YOLO (you only look once)



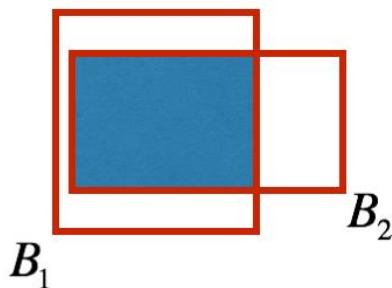
For each grid cell:

$$\mathbf{y} = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

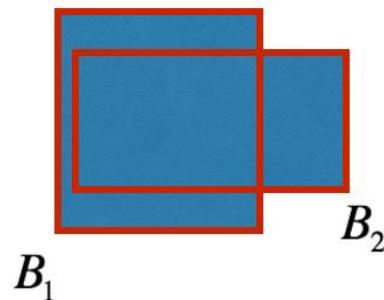
Redmon et al., 2015, You Only Look Once: Unified real-time object detection

Intersection over Union (IoU)

Intersection



Union



Intersection over Union

$$IoU = \frac{B_1 \cap B_2}{B_1 \cup B_2} = \frac{\text{Area of intersection}}{\text{Area of union}} = P_c$$

The diagram illustrates the formula for IoU. It shows two overlapping bounding boxes, B_1 and B_2 , outlined in red and filled with blue. The intersection of the two boxes is highlighted with a dashed red border and a solid blue fill. The union of the two boxes is shown as a larger rectangle, also outlined in red and filled with blue. The formula $IoU = \frac{B_1 \cap B_2}{B_1 \cup B_2}$ is written above the boxes, with arrows pointing from the intersection and union regions to their respective terms in the formula.

“Correct” if $IoU \geq 0.5$

More generally, IoU is a measure of the overlap between two bounding boxes.

<Deep Learning, Andrew Ng>

Non-max suppression

Each output prediction is:

$$\begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \end{bmatrix}$$

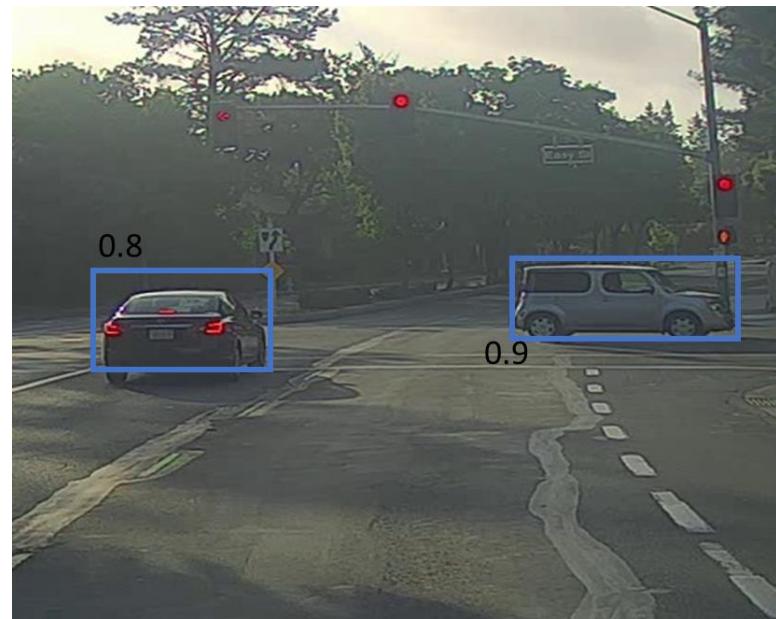
Discard all boxes with $p_c \leq 0.6$
Among remaining boxes,
Pick the box with the largest p_c .
Output that as a prediction.

Non-max suppression

Before non-max suppression



After non-max suppression



<Deep Learning, Andrew Ng>

Face recognition

Input image :



- Has a database of K persons
- Get an input image
- Output identity (name or id) if the image is any of the K persons (or “not recognized”)
- OneShotLearning

<Deep Learning, Andrew Ng>

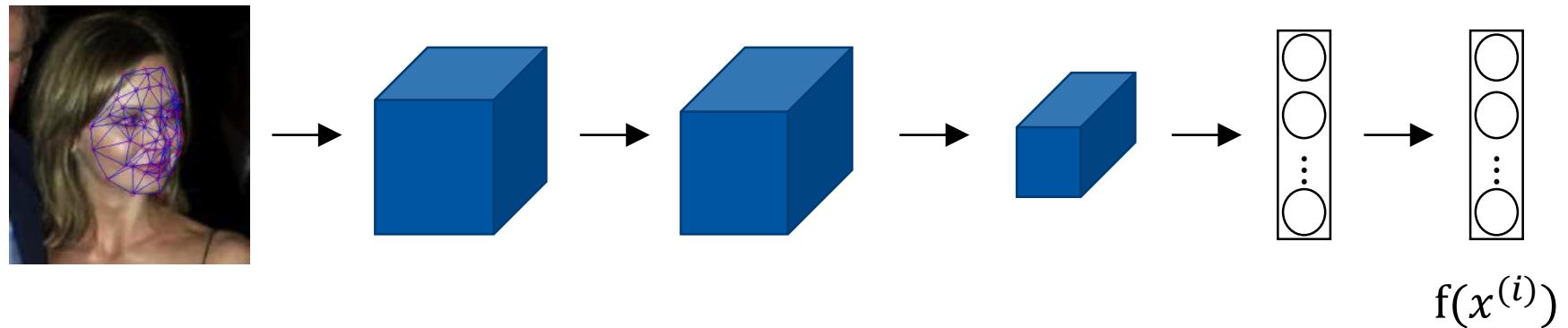
Distance between images



- $d(\text{img1}, \text{img2})$ = degree of difference between images
- Image similarity
 - $d(\text{img1}, \text{img2}) \leq \tau \rightarrow \text{same}$
 - $d(\text{img1}, \text{img2}) > \tau \rightarrow \text{different}$

<Deep Learning, Andrew Ng>

Deep face feature vector



If $x^{(i)}, x^{(j)}$ are same person, $d(f(x^{(i)}) - f(x^{(j)}))$ is small.

If $x^{(i)}, x^{(j)}$ are different persons, $d(f(x^{(i)}) - f(x^{(j)}))$ is large.

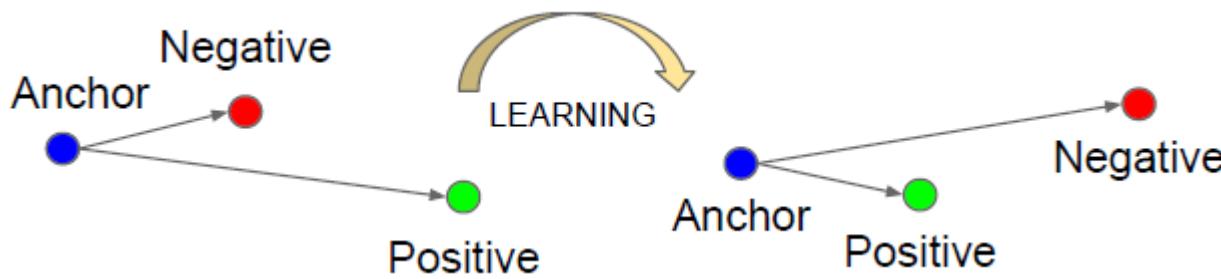
Distance between two feature vectors

- χ^2 distance =
$$\frac{(f(x^{(i)}) - f(x^{(j)}))^2}{f(x^{(i)}) + f(x^{(j)})}$$
- Siamese network: $d(x^{(i)} - x^{(j)}) = \|f(x^{(i)}) - f(x^{(j)})\|_2^2$

Same as L2 norm

Taigman et. al., 2014. DeepFace closing the gap to human level performance

Triplet Loss



$$\|f(x_i^a) - f(x_i^p)\|_2^2 + \alpha < \|f(x_i^a) - f(x_i^n)\|_2^2$$

Schroff et al., 2015, FaceNet: A unified embedding for face recognition and clustering

Triplet Loss

Loss is minimized:

$$L = \sum_i^N \left[\|f(x_i^a) - f(x_i^p)\|_2^2 - \|f(x_i^a) - f(x_i^n)\|_2^2 + \alpha \right]$$

Schroff et al., 2015, FaceNet: A unified embedding for face recognition and clustering

Triplet Selection

For fast convergence,

Given x_i^a ,

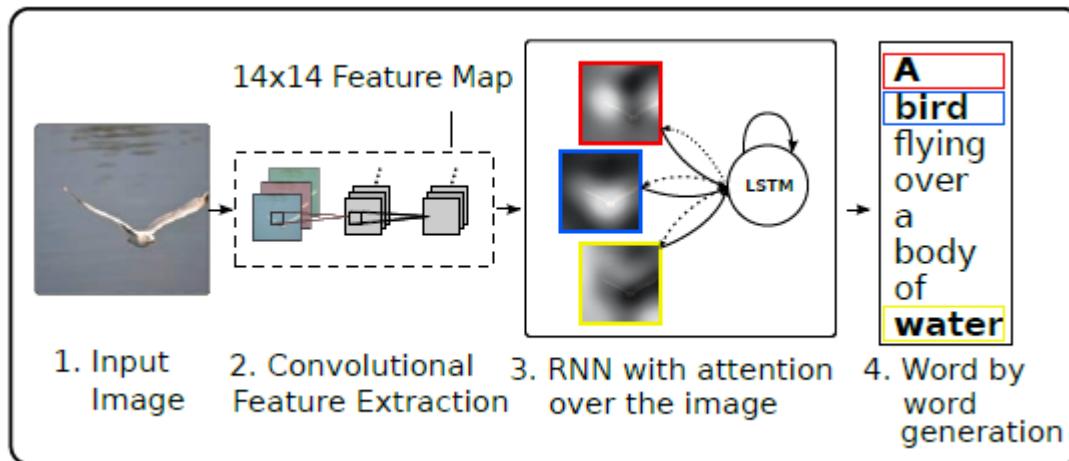
Select hard positive x_i^p , to make $\|f(x_i^a) - f(x_i^p)\|_2^2$ maximum.

Select hard negative x_i^n , to make $\|f(x_i^a) - f(x_i^n)\|_2^2$ minimum.

$$\|f(x_i^a) - f(x_i^p)\|_2^2 + \alpha < \|f(x_i^a) - f(x_i^n)\|_2^2$$

Neural Image Caption Generation with Visual Attention

Visualized attentional maps



Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, Kelvin Xu et al., 2016

How to separate words from a sentence?

Tokenization

Tokenization is a process that splits an input sequence into tokens.

We can split token by space, punctuation, a set of rule.

Tokenization

In Düsseldorf I took my hat off. But I can't put it back on.

WhitespaceTokenizer



WordPunctTokenizer



TreebankWordTokenizer



<http://text-processing.com/demo/tokenize/>

Token normalization

Same token for different forms of words

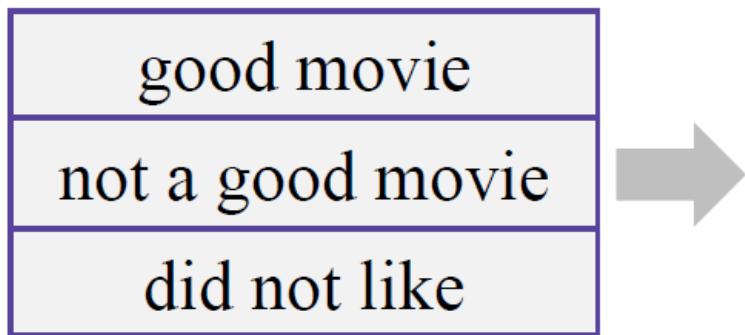
- Examples
 - wolf, wolves → wolf
 - talk, talks → talk
- Stemming
 - removes and replaces suffixes to get to the root form of a word, which is called as stem.
- Lemmatization
 - returns the base or dictionary form of a word, which is known as lemma.

Coursera: Natural Language Processing, National Research University Higher School of Economics

Transforming tokens into features

Bag of words (BOW)

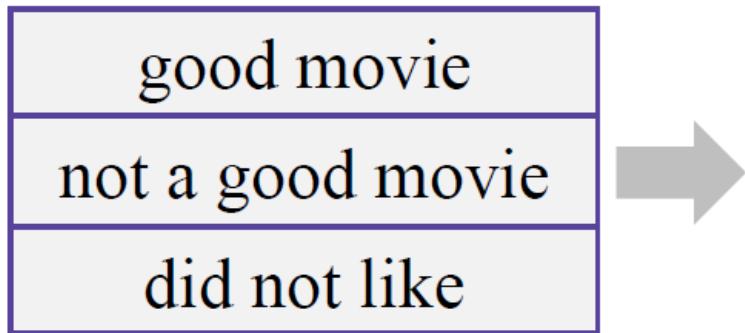
For each token, we have a feature column, which is called text vectorization.



Coursera: Natural Language Processing, National Research University Higher School of Economics

Preserve some ordering

N-grams: Token pairs, triplets, etc.



Coursera: Natural Language Processing, National Research University Higher School of Economics

Remove some n-grams

- High frequency n-grams
 - Articles, prepositions, etc. (example: and, a, the)
 - They are called stop-words. They do not help to discriminate texts.
- Low frequency n-grams
 - Typos, rare words

Word Embedding

- Convert texts into numbers
- Map a word to a vector using a dictionary
- Applications
 - Sentiment analysis of reviews (amazon, movie review)
 - Document or news classification or clustering (google)

<https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/>

Frequency based Embedding

- Count Vector
 - Frequency : Number of times a word has appeared in the document
 - Presence : Has the word appeared in the document?
- TF-IDF Vector
- Co-occurrence Vector

<https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/>

Word Matrix

	Document 1	Document 2	Document 3	Document 4	Document 5	Document 6	Document 7	Document 8
Term(s) 1	10	0	1	0	0	0	0	2
Term(s) 2	0	2	0	0	0	18	0	2
Term(s) 3	0	0	0	0	0	0	0	2
Term(s) 4	6	0	0	4	6	0	0	0
Term(s) 5	0	0	0	0	0	0	0	2
Term(s) 6	0	0	1	0	0	1	0	0
Term(s) 7	0	1	8	0	0	0	0	0
Term(s) 8	0	0	0	0	0	3	0	0

↑
Document Vector

Word Vector
(Passage Vector)

How to make Term (word) features

1. All words in a dictionary
2. Unique words in corpus(all documents)

<https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/>

TF-IDF

- N-gram with smaller frequency can be more discriminating because it can capture a specific issue in the text
- Term frequency (TF)
 - Frequency for term (or n-gram) t in document d

term frequency	$f_{t,d} / \sum_{t' \in d} f_{t',d}$
----------------	--------------------------------------

- Inverse document frequency (IDF)
 - $N = D$: total number of documents in corpus
 - $\{d \in D : t \in d\}$: Number of documents where the term t appears

$$\text{idf}(t, D) = \log \frac{N}{|\{d \in D : t \in d\}|}$$

Coursera: Natural Language Processing, National Research University Higher School of Economics

TF-IDF

$$Tfidf(t,d,D) = tf(t,d) \times idf(t,D)$$

A high weight in TF-IDF means a high term frequency (in a given document) and a low document frequency of the term in a all collection of documents

good movie
not a good movie
did not like



good movie	movie	did not	...
0.17	0.17	0	...
0.17	0.17	0	...
0	0	0.47	...

Term frequency & Inverse document frequency

1-gram

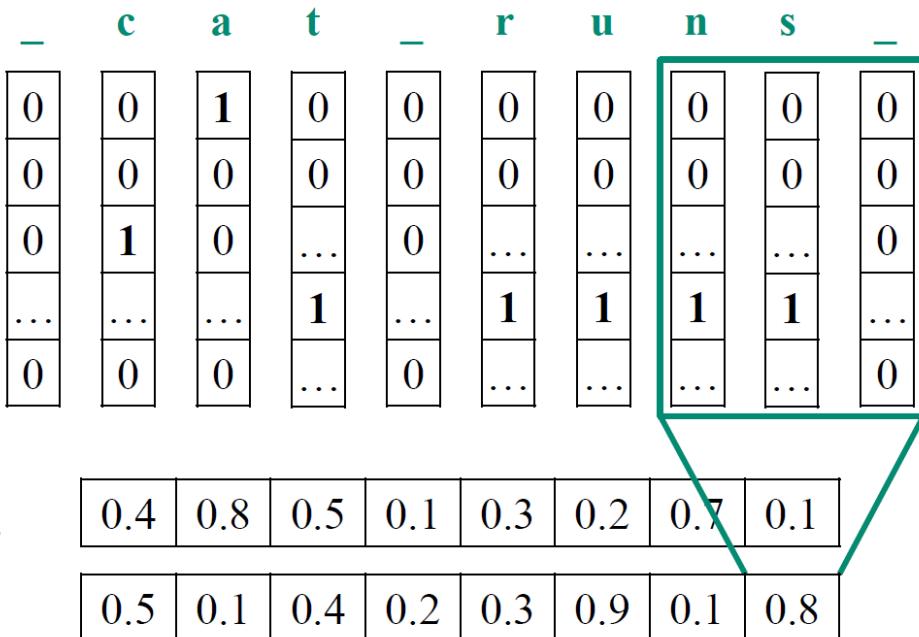
							term frequency				
text	good	movie	not	a	did	not	like	I	it	one	
good movie											
not a good movie											
did not like											
I like it											
good one											
							inverse document frequency				
text	good	movie	not	a	did	not	like	I	it	one	
good movie											
not a good movie											
did not like											
I like it											
good one											

Text as a sequence of characters

One-hot encoding characters, length ~ 70

_	c	a	t	_	r	u	n	s	_
0	0	1	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	1	0	...	0	0
...	1	...	1	1	1	1	...
0	0	0	...	0	0

1D convolutions on characters



Coursera: Natural Language Processing, National Research University Higher School of Economics

Max pooling

Provides position invariance for character n-gram

Filter #1

0.4	0.8	0.5	0.1	0.3	0.2	0.7	0.1
-----	-----	-----	-----	-----	-----	-----	-----

Filter #2

0.5	0.1	0.4	0.2	0.3	0.9	0.1	0.8
-----	-----	-----	-----	-----	-----	-----	-----

Filter #3

0.4	0.7	0.3	0.7	0.5	0.5	0.9	0.4
-----	-----	-----	-----	-----	-----	-----	-----

0.8	0.5	0.3	0.7
-----	-----	-----	-----

--	--	--	--

--	--	--	--

Repeat 1D convolutions + pooling

Pooling
output

0.8	0.5	0.3	0.7
0.5	0.4	0.9	0.8
0.7	0.7	0.5	0.9

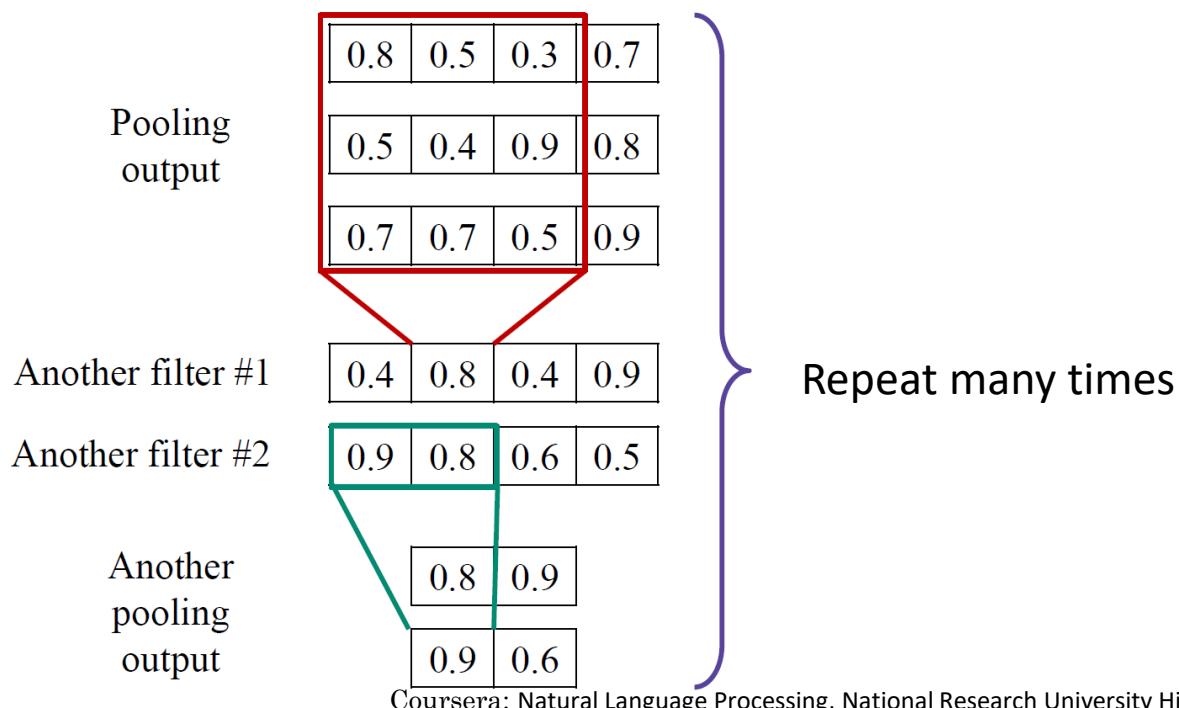
Another filter #1

0.4	0.8	0.4	0.9
-----	-----	-----	-----

Another filter #2

0.9	0.8	0.6	0.5
-----	-----	-----	-----

Repeat 1D convolutions + pooling



Example of 1D convolution architecture for characters

- Characters of text : 1014
- Apply 1D convolution + max pooling 6 times
- Kernels widths : 7,7,3,3,3,3
- # of Filters at each step : 1024
- Output: 1024 x 34 matrix of features
- Apply multiclass classification for sentiment analysis

Problem of one-hot representation

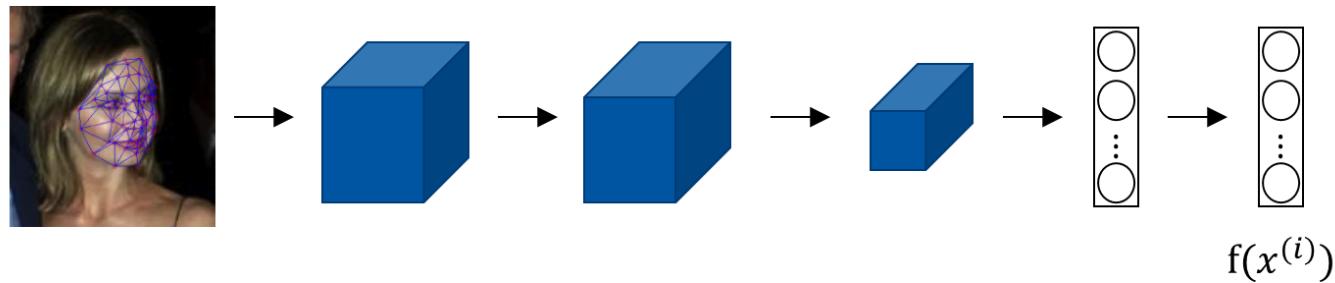
- It treats each word as a thing unto itself, and it doesn't allow an algorithm to easily generalize the cross words.
- Example
 - Sally Johnson is an orange farmer.
 - Robert Lin is an apple farmer.
 - Robert Lin is a durian cultivator.
- product of any two word vector is zero.

Featured representation: word embedding

	Man	Woman	King	Queen	Apple	Orange
Gender	-1	1	-0.95	0.97	0	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0
Age	0.03	0.02	0.7	0.69	0.03	-0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97
Size						
Cost						
alive						
verb						

Word embedding vector

Deep face feature vector



If $x^{(i)}, x^{(j)}$ are same person, $d(f(x^{(i)}) - f(x^{(j)}))$ is small.

If $x^{(i)}, x^{(j)}$ are different persons, $d(f(x^{(i)}) - f(x^{(j)}))$ is large.

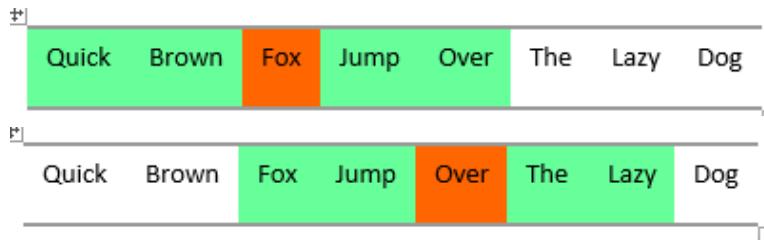
Taigman et. al., 2014. DeepFace closing the gap to human level performance

Transfer learning and word embeddings

- Learn word embeddings from large text corpus. (1-100B words)
(Or download pre-trained embedding online.)
- Transfer embedding to new task with smaller training set. (say, 100k words)
- Continue to fine-tune the word embeddings with new data.

Co-occurrence Matrix

- Hypothesis: Similar words tend to occur together and will have similar context.
- Example
 - Apple is a fruit. Mango is a fruit.
- Co-occurrence
 - For a given corpus, the co-occurrence of a pair of words is the number of times they have appeared together in a context window
- Context window
 - Context window is specified by a number and the direction



<https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/>

Co-occurrence Matrix

Corpus = He is not lazy. He is intelligent. He is smart.

	He	is	not	lazy	intelligent	smart
He	0	4	2	1	2	1
is	4	0	1	2	2	1
not	2	1	0	1	0	0
lazy	1	2	1	0	0	0
intelligent	2	2	0	0	0	0
smart	1	1	0	0	0	0

Context windows = 2

He	is	not	lazy	He	is	intelligent	He	is	smart
He	is	not	lazy	He	is	intelligent	He	is	smart
He	is	not	lazy	He	is	intelligent	He	is	smart
He	is	not	lazy	He	is	intelligent	He	is	smart

<https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/>

Co-occurrence Matrix

Matrix size

- $V \times V$
 - Not practical
- $V \times N$
 - N is a subset of V and can be obtained by removing irrelevant words like stop words
- $V \times k$
 - k is k principal components out of V using PCA

PCA to decompose Co-occurrence matrix

$X = U \cdot S \cdot V^T$ (Singular value decomposition)

U and S represent word vector

V presents word context

U is principal component.

$$\begin{pmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & \\ \vdots & \ddots & & \\ x_{m1} & & & x_{mn} \end{pmatrix}_{m \times n} \approx \begin{pmatrix} u_{11} & \dots & u_{1r} \\ \vdots & \ddots & \\ u_{m1} & & u_{mr} \end{pmatrix}_{m \times r} \begin{pmatrix} s_{11} & 0 & \dots \\ 0 & \ddots & \\ \vdots & & s_{rr} \end{pmatrix}_{r \times r} \begin{pmatrix} v_{11} & \dots & v_{1n} \\ \vdots & \ddots & \\ v_{r1} & & v_{rn} \end{pmatrix}_{r \times n}$$

$m \times m \qquad m \times n \qquad n \times n$

<https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/>

Frequency based Embedding

- Count Vector
 - Frequency : Number of times a word has appeared in the document
 - Presence : Count if the word appeared in the document
- TF-IDF Vector
- Co-occurrence Vector

<https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/>

Analogy

	Man	Woman	King	Queen	Apple	Orange
Gender	-1	1	-0.95	0.97	0	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0
Age	0.03	0.02	0.7	0.69	0.03	-0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97

Man → Woman vs. King → ? (Queen)

$$e_{\text{Man}} - e_{\text{Woman}} \approx e_{\text{King}} - e_{\text{Queen}}$$

Linguistic regularities in continuous space word representations, Mikolov et. al., 2013,

Coursera: Deep learning Specialization, Andrew Ng

Analogy using word vectors

Man → Woman vs. King → ? (Queen)

$$e_{\text{Man}} - e_{\text{Woman}} \approx e_{\text{King}} - e_{\text{Queen}}$$

$$e_{\text{Man}} - e_{\text{Woman}} = \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad e_{\text{King}} - e_{\text{Queen}} = \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$e_{\text{Man}} - e_{\text{Woman}} \approx e_{\text{King}} - e_w$$

Man:Woman as Boy:Girl
Ottawa:Canada as Nairobi:Kenya
Big:Bigger as Tall:Taller
Yen:Japan as Ruble:Russia

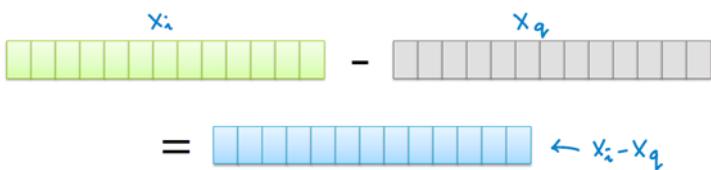
Find word w : $\arg \max (e_w, e_{\text{King}} - e_{\text{Man}} + e_{\text{Woman}})$

Distance between two word vectors

(non-scaled) Euclidean distance

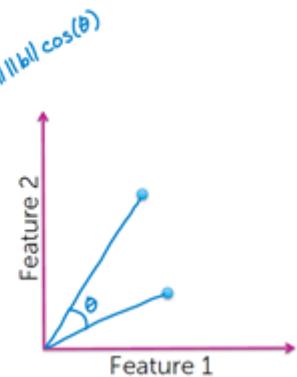
Defined in terms of inner product

$$\text{distance}(\mathbf{x}_i, \mathbf{x}_q) = \sqrt{(\mathbf{x}_i - \mathbf{x}_q)^T (\mathbf{x}_i - \mathbf{x}_q)}$$
$$= \sqrt{(\mathbf{x}_i[1] - \mathbf{x}_q[1])^2 + \dots + (\mathbf{x}_i[d] - \mathbf{x}_q[d])^2}$$



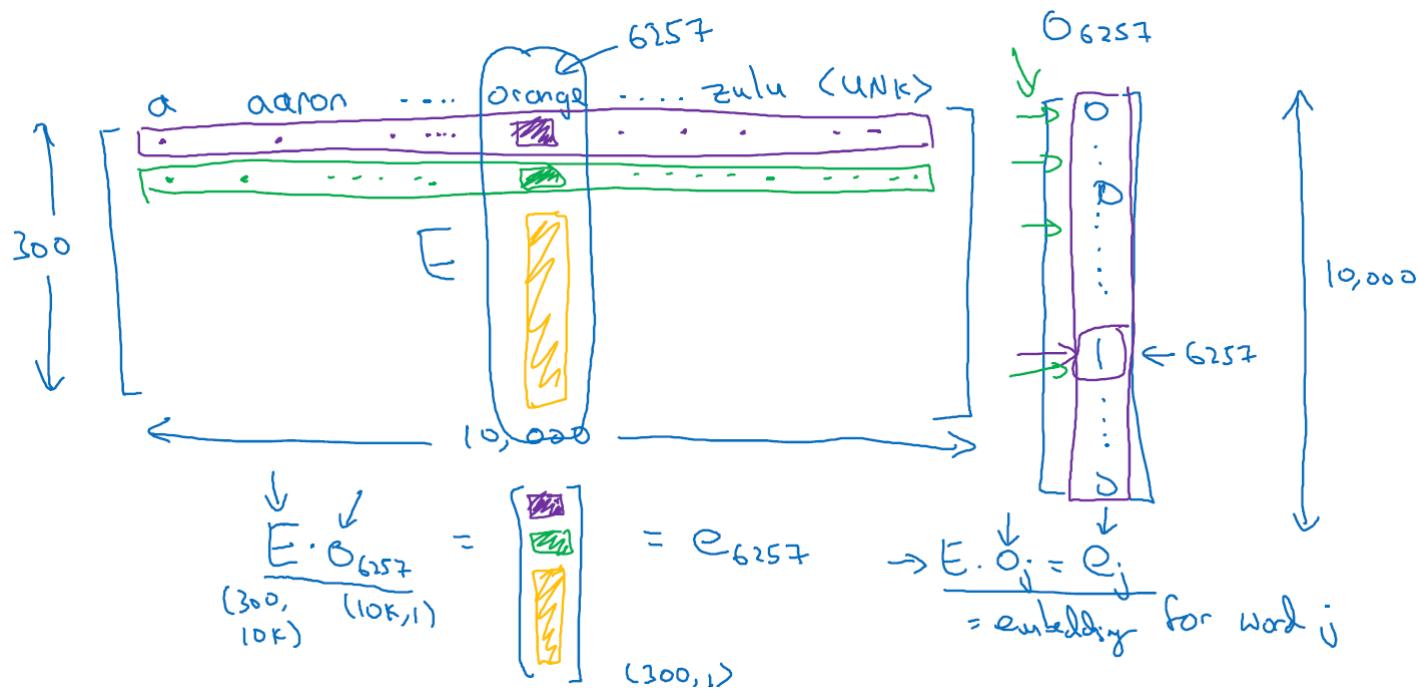
Cosine similarity – normalize

$$\text{Similarity} = \frac{\sum_{j=1}^d \mathbf{x}_i[j] \mathbf{x}_q[j]}{\sqrt{\sum_{j=1}^d (\mathbf{x}_i[j])^2} \sqrt{\sum_{j=1}^d (\mathbf{x}_q[j])^2}}$$
$$\xrightarrow{\text{a}^T b = \|a\| \|b\| \cos(\theta)}$$
$$\frac{\mathbf{x}_i^T \mathbf{x}_q}{\|\mathbf{x}_i\| \|\mathbf{x}_q\|} = \cos(\theta)$$
$$= \left(\frac{\mathbf{x}_i}{\|\mathbf{x}_i\|} \right)^T \left(\frac{\mathbf{x}_q}{\|\mathbf{x}_q\|} \right) \xrightarrow{\text{first normalize}}$$



Coursera: Machine Learning, Emily Fox & Carlos Guestrin

Pull out word vector from Embedding matrix



Coursera: Deep learning Specialization, Andrew Ng

Language model using word vector

I want a glass of orange $\underline{\quad}$.

4343 9665 1 3852 6163 6257
 $e_{4343} = E o_{4343}$

I $o_{4343} \rightarrow E \rightarrow e_{4343}$

want $o_{9665} \rightarrow E \rightarrow e_{9665}$

a $o_1 \rightarrow E \rightarrow e_1$

glass $o_{3852} \rightarrow E \rightarrow e_{3852}$

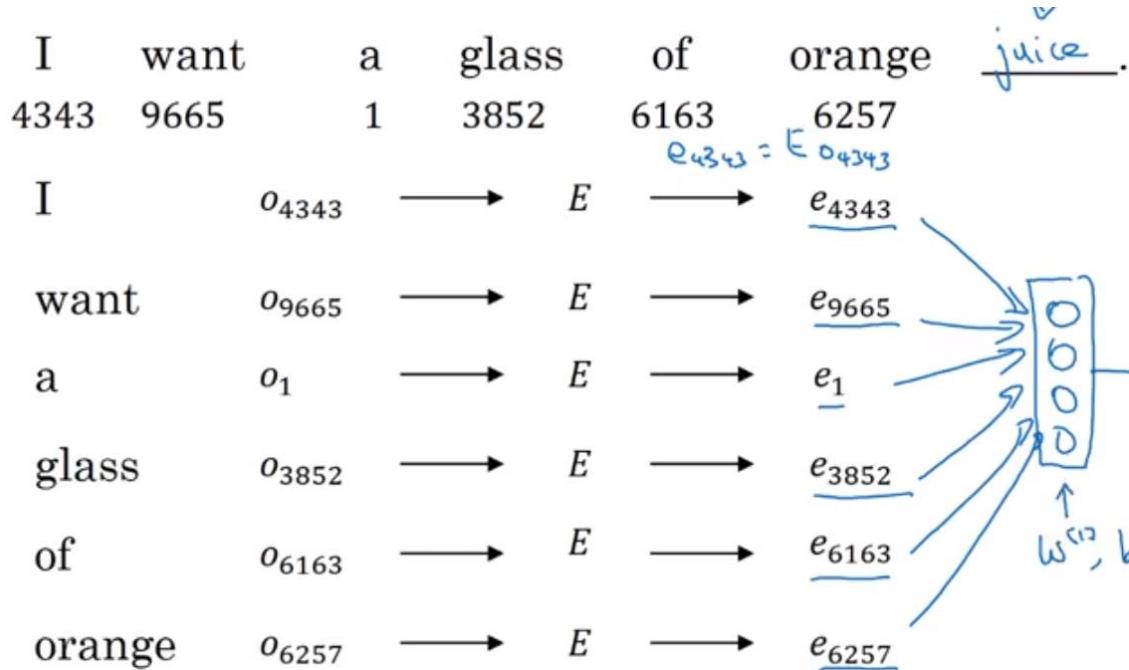
of $o_{6163} \rightarrow E \rightarrow e_{6163}$

orange $o_{6257} \rightarrow E \rightarrow e_{6257}$

Language model using word vector

A neural probabilistic language model , Bengio et. al., 2003

$$\hat{P}(w_1^T) = \prod_{t=1}^T \hat{P}(w_t | w_1^{t-1})$$

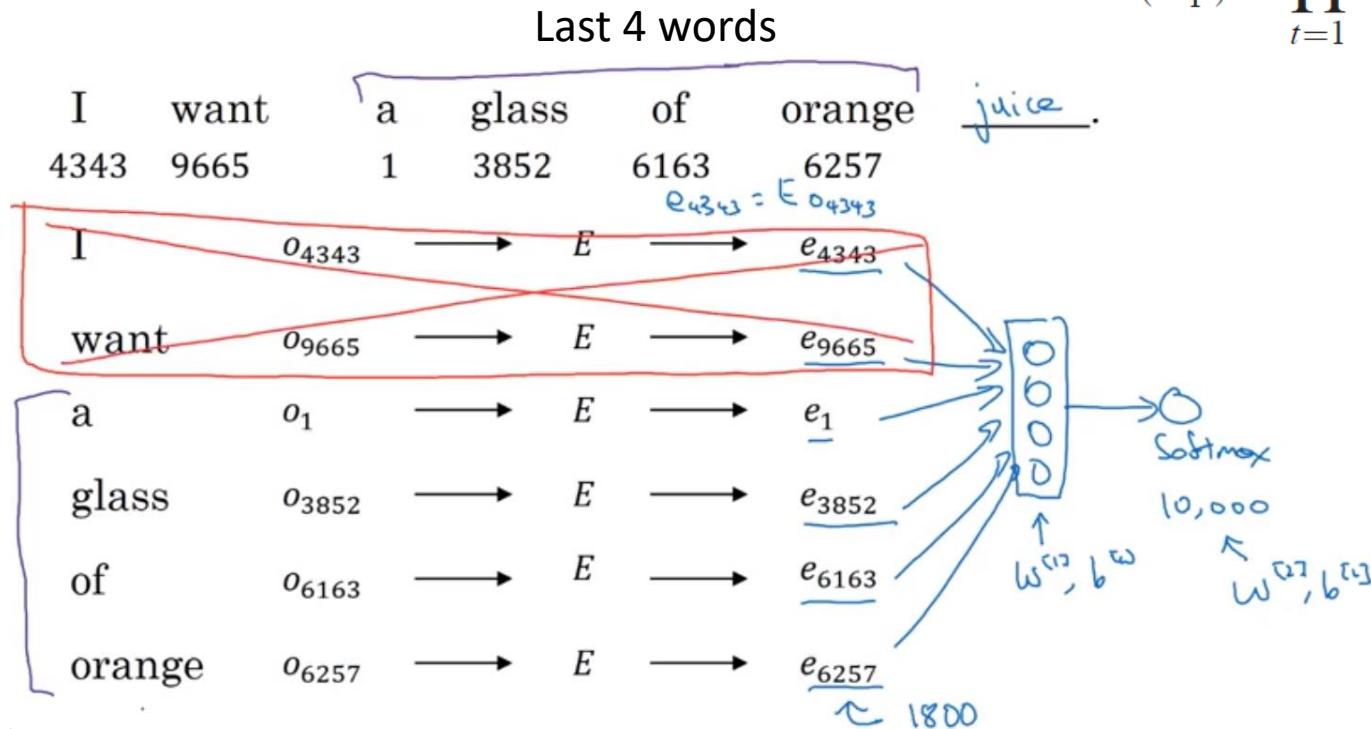


Coursera: Deep learning Specialization, Andrew Ng

Language model using word vector

A neural probabilistic language model , Bengio et. al., 2003

$$\hat{P}(w_1^T) = \prod_{t=1}^T \hat{P}(w_t | w_1^{t-1})$$



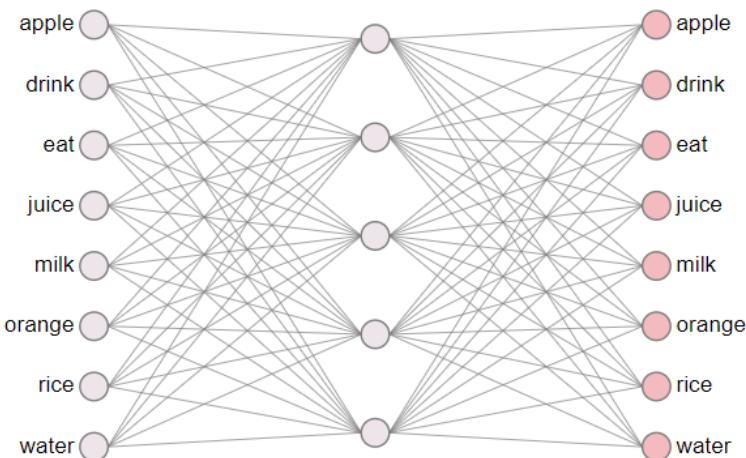
Coursera: Deep learning Specialization, Andrew Ng

Word2vec : Prediction based on Embedding

- CBOW (Continuous Bag of words)
 - Predicts the current word based on the context.
- Skip-gram
 - Predicts surrounding words given the current word.

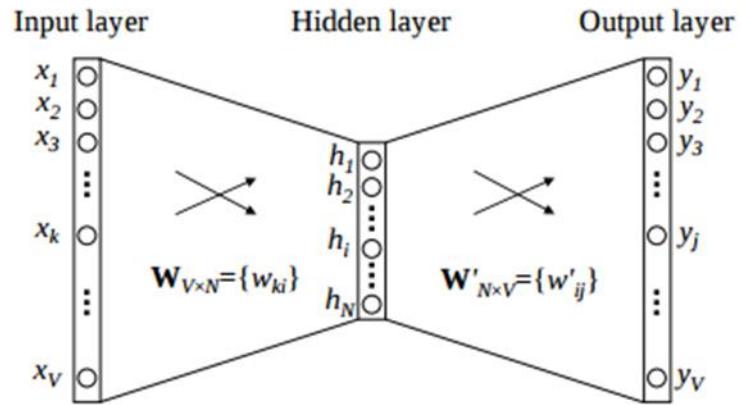
CBOW(Continuous Bag of words)

Predicts the current word based on the context



<https://ronxin.github.io/wevi/>

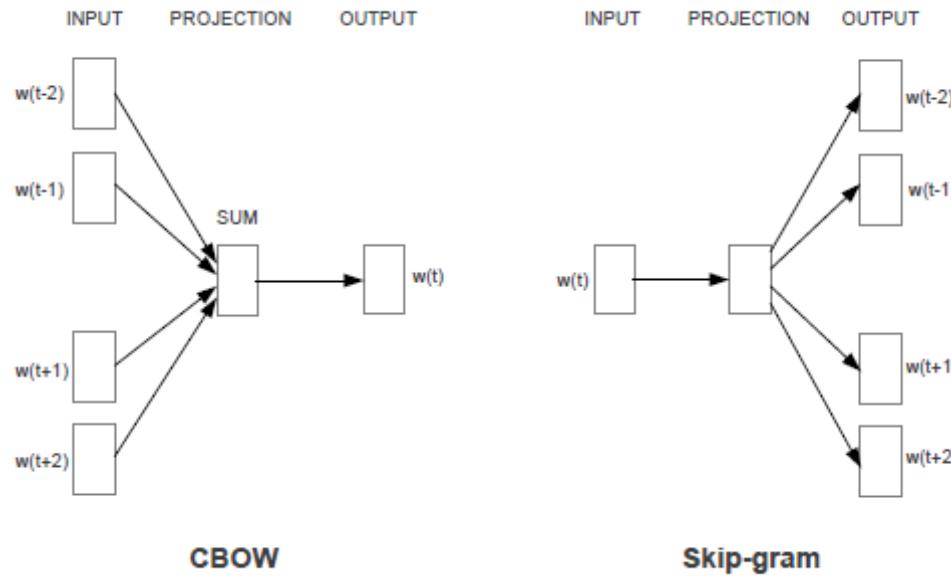
word vector :
the weight between the hidden layer and the output layer



Linear activation (no activation function between any layers)

<https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/>

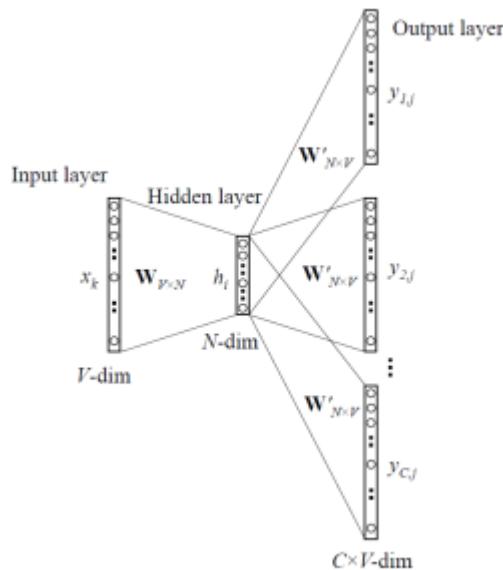
Word2vec : Prediction based on Embedding



Efficient Estimation of Word Representations in Vector Space, Tomas Mikolov et al., 2013

Skip-gram

Predicts surrounding words given the current word



<https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/>

Context/target pair

I want a glass of orange juice to go along with my cereal.
target

Context:

Last 4 words : a glass of orange

4 words on left & right : a glass of orange, to go along with

Last 1 word : orange

Nearby 1 word → skip-gram

Skip-gram

I want a glass of orange juice to go along with my cereal.

Content

Orange

Orange

Orange

Target

Juice

Glass

to

Next 1 word

Left/right 2 windows

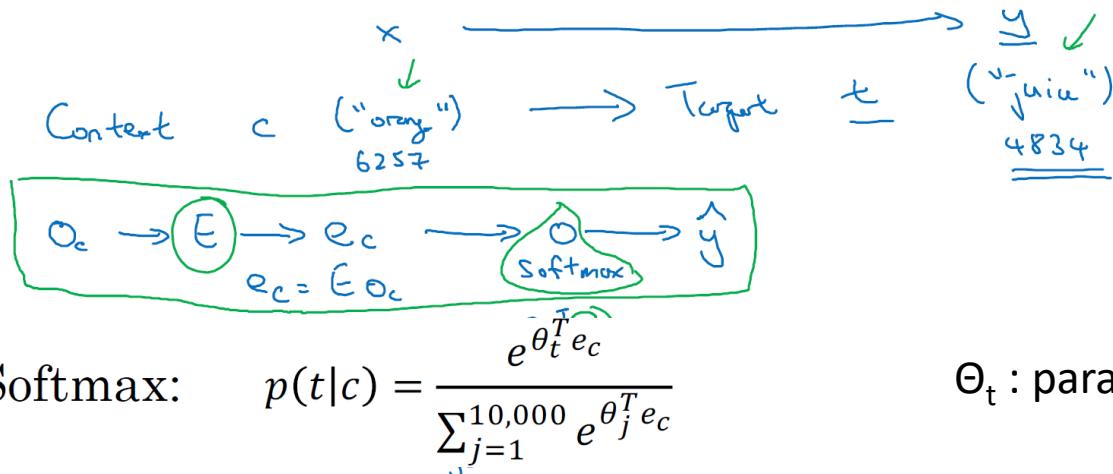
Skip-gram

How to sample content word?

- Randomly select → there is high chance to pick up most of stop-words (a, the, and, etc)
- Select from less-frequent words

Softmax-classification

Vocab size = 10,000k



Θ_t : parameter associated with target t

Cost: $\mathcal{L}(\hat{y}, y) = - \sum_{i=1}^{10,000} y_i \log \hat{y}_i$

Coursera: Deep learning Specialization, Andrew Ng

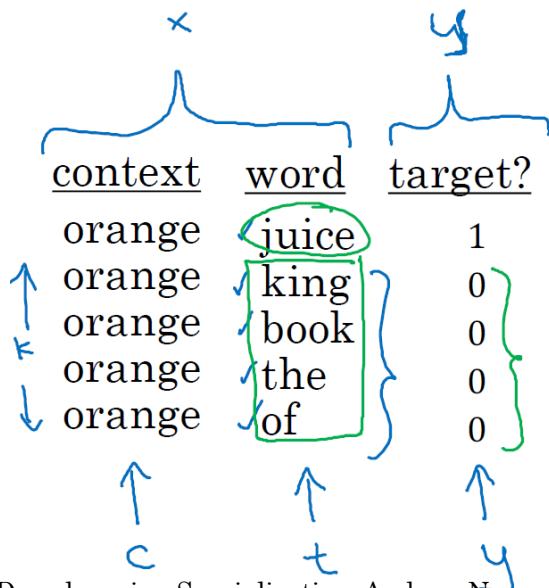
Negative sampling

I want a glass of orange juice to go along with my cereal.

Binary classification

$$P(y=1 | c, t) = \sigma(\theta_t^T e_c)$$

Sigmoid function



Randomly pick up
from dictionary

or

$$P(w_i) = \frac{f(w_i)^{2/4}}{\sum_{j=1}^{10,000} f(w_j)^{2/4}}$$

Coursera: Deep learning Specialization, Andrew Ng

Word2vec

vocabulary_size=7 and embedding_size=3

<i>anarchism</i>	0.5	0.1	-0.1
<i>originated</i>	-0.5	0.3	0.9
<i>as</i>	0.3	-0.5	-0.3
<i>a</i>	0.7	0.2	-0.3
<i>term</i>	0.8	0.1	-0.1
<i>of</i>	0.4	-0.6	-0.1
<i>abuse</i>	0.7	0.1	-0.4

	Man	Woman	King	Queen	Apple	Orange
Gender	-1	1	-0.95	0.97	0	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0
Age	0.03	0.02	0.7	0.69	0.03	-0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97

<http://adventuresinmachinelearning.com/word2vec-tutorial-tensorflow/>

GloVe: Global Vectors for Word Representation

GloVe: Global Vectors for Word Representation, Jeffrey Pennington et al., 2014

I want a glass of orange juice to go along with my cereal.

$X_{ij} = \# \text{ times } j \text{ appears in content of } i$

content target target content

Shallow Window-based methods (making predictions within local context windows)

GloVe: Model

GloVe: Global Vectors for Word Representation, Jeffrey Pennington et al., 2014

Loss function
$$J = \sum_{i,j=1}^V f(X_{ij}) (w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij})^2$$

minimize
$$\sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij}) (\Theta_i^T e_j + b_i + b_j' - \log X_{ij})^2$$

1. $f(0) = 0$. If f is viewed as a continuous function, it should vanish as $x \rightarrow 0$ fast enough that $\lim_{x \rightarrow 0} f(x) \log_2 x$ is finite.
2. $f(x)$ should be non-decreasing so that rare co-occurrences are not overweighted.
3. $f(x)$ should be relatively small for large values of x , so that frequent co-occurrences are not overweighted.

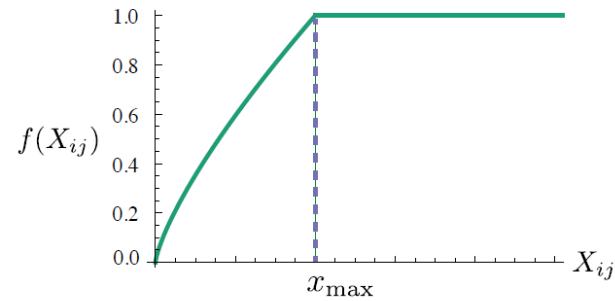


Figure 1: Weighting function f with $\alpha = 3/4$.

GloVe: Global Vectors for Word Representation

GloVe: Global Vectors for Word Representation, Jeffrey Pennington et al., 2014

- Unsupervised word representation
- Combined
 - Count-based method
 - Prediction-based method
- Outperforms
 - Word analogies
 - Word similarity
 - Named entity recognition

Sentiment classification problem

x

The dessert is excellent.

y



Service was quite slow.



Good for a quick meal, but nothing special.

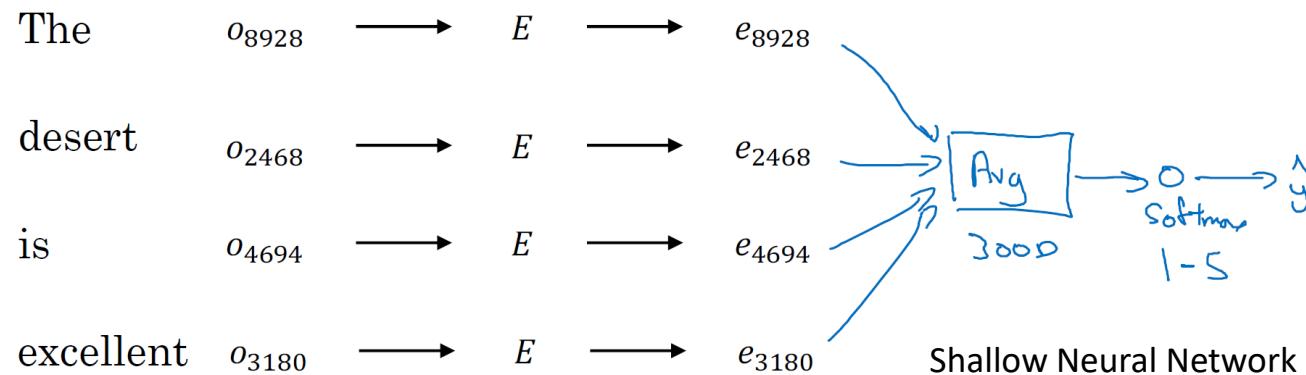


Completely lacking in good taste, good service, and good ambience.



Simple sentiment classification model

The dessert is excellent
8928 2468 4694 3180

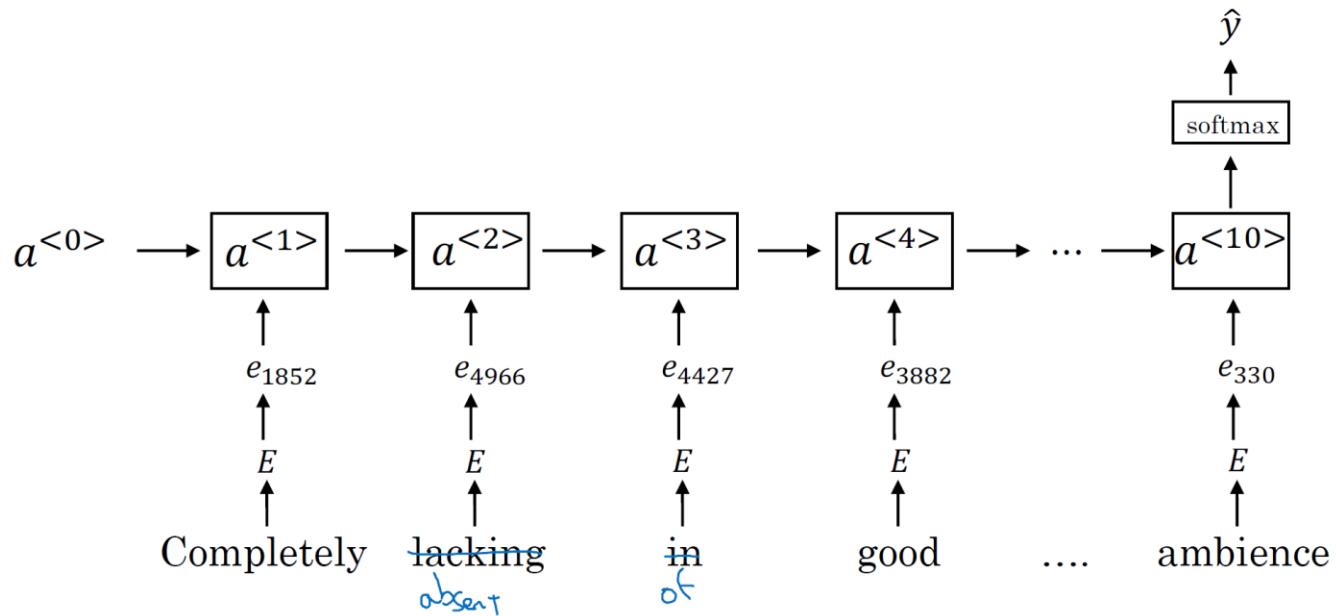


Problem in Simple sentiment classification model

Completely lacking in **good** taste, **good** service,
and **good** ambience.



RNN for sentiment classification



Character-level language model

Vocabulary = [a, b, c, ..., z, [], ., ', ;, ..., A, ..., Z]

- Advantage:
 - Do not worry about unknown token.
- Disadvantage:
 - Much longer sequence
 - Is not as good as word level language models at capturing long range dependencies between how
 - the earlier parts of the sentence also affect the later part of the sentence.
 - more computationally expensive to train

Example of Name entity recognition

Name entity recognition can be used to find people's names, companies names, times, locations, countries names, currency names, and so on.

x: (Harry Potter) and (Hermione Granger) invented a new spell.

y: | | o | | o o o o

Name entity recognition

x: Harry Potter) and Hermione Granger invented a new spell.

Position of index $x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad \dots \quad x^{<t>} \quad \dots \quad x^{<n>}$

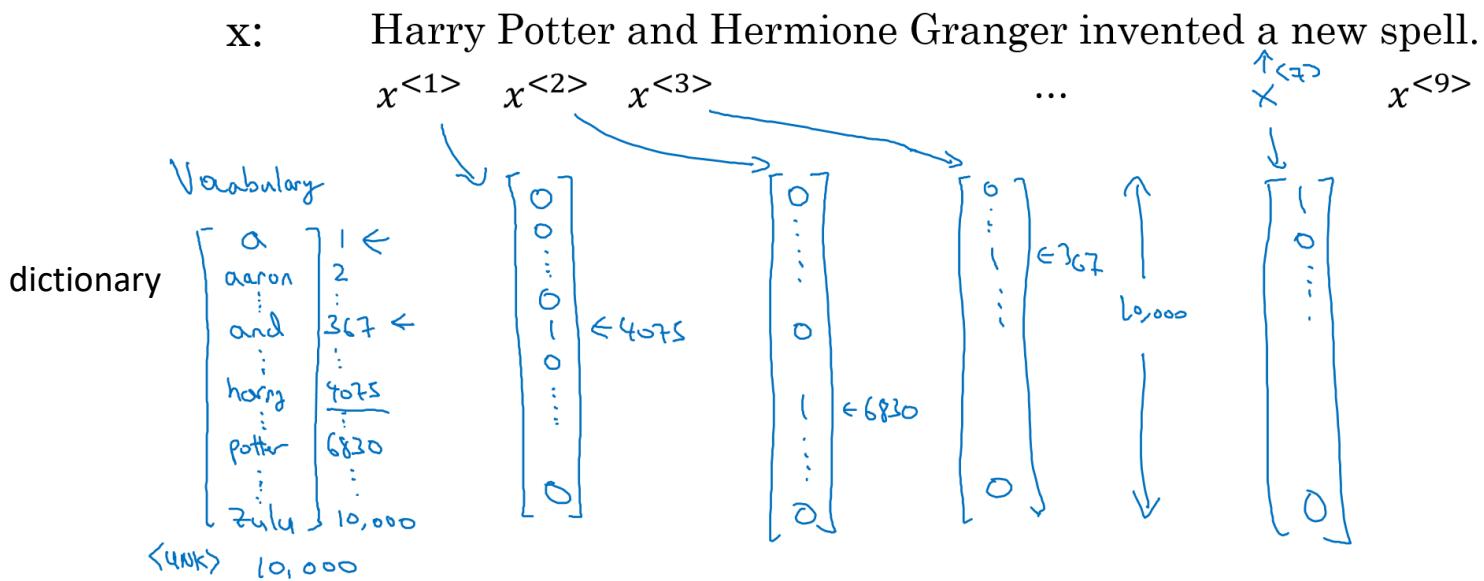
$$\text{length } T_x = 9$$

y: | | 0 | | 0 0 0 0
 $y^{<1>} \quad y^{<2>} \quad y^{<3>} \quad \dots \quad y^{<9>}$

$$\text{length } T_y = 9$$

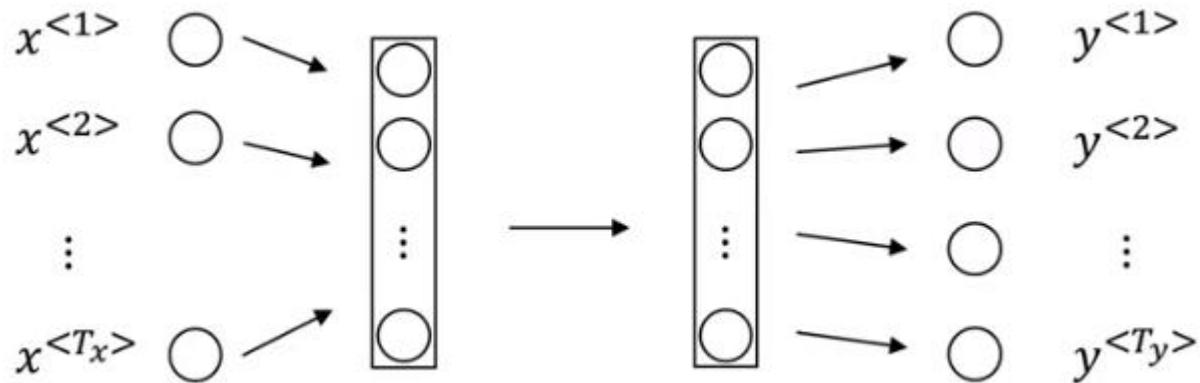
One-hot vector to represent words

Map: $x^{<t>} \rightarrow y^{<t>}$



Coursera: Deep learning Specialization, Andrew Ng

Why not a Neural Network?



Problems:

- Arbitrary sequence length : Inputs, outputs can be different lengths in different examples.
- Doesn't share features learned across different positions of text.
- Large number of parameters

Coursera: Deep learning Specialization, Andrew Ng

Language model

Cats sleep average 15 hours per a day. <EOS>

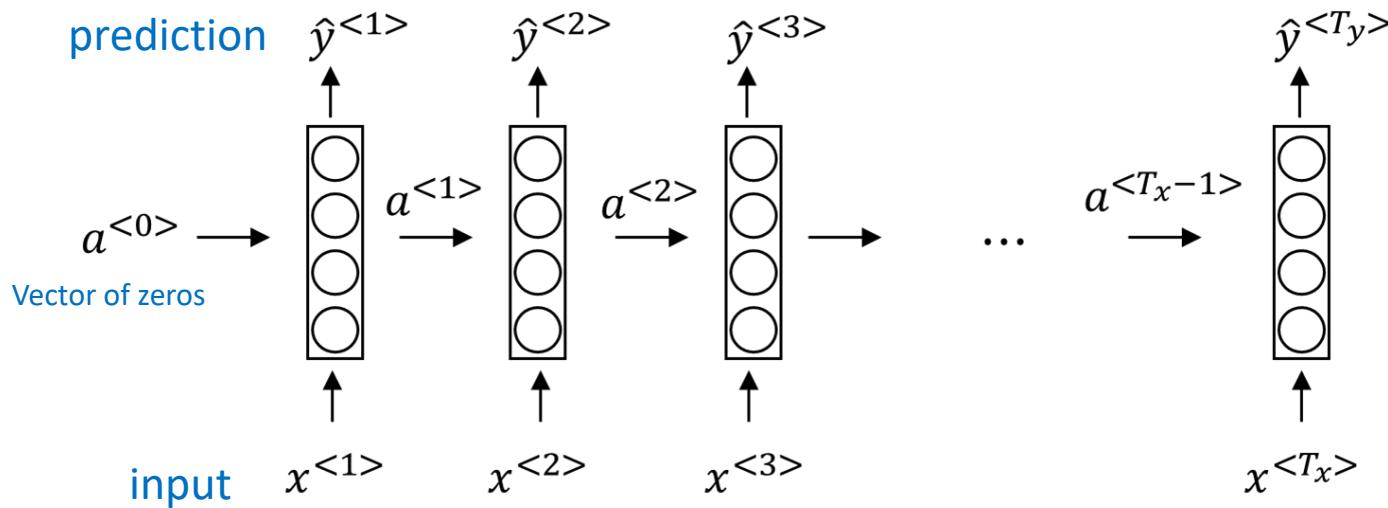
$$\begin{aligned} P(\text{text}) &= P(x_0, \dots, x_n) = \\ &= P(x_0)P(x_1|x_0)P(x_2|x_0, x_1)\dots P(x_n|\dots) \end{aligned}$$

$$\hat{P}(w_1^T) = \prod_{t=1}^T \hat{P}(w_t | w_1^{t-1})$$

A neural probabilistic language model , Bengio et. al., 2003

Coursera: Natural Language Processing, National Research University Higher School of Economics

Recurrent Neural Networks



$x^{<t>}$: current input (word embedding)

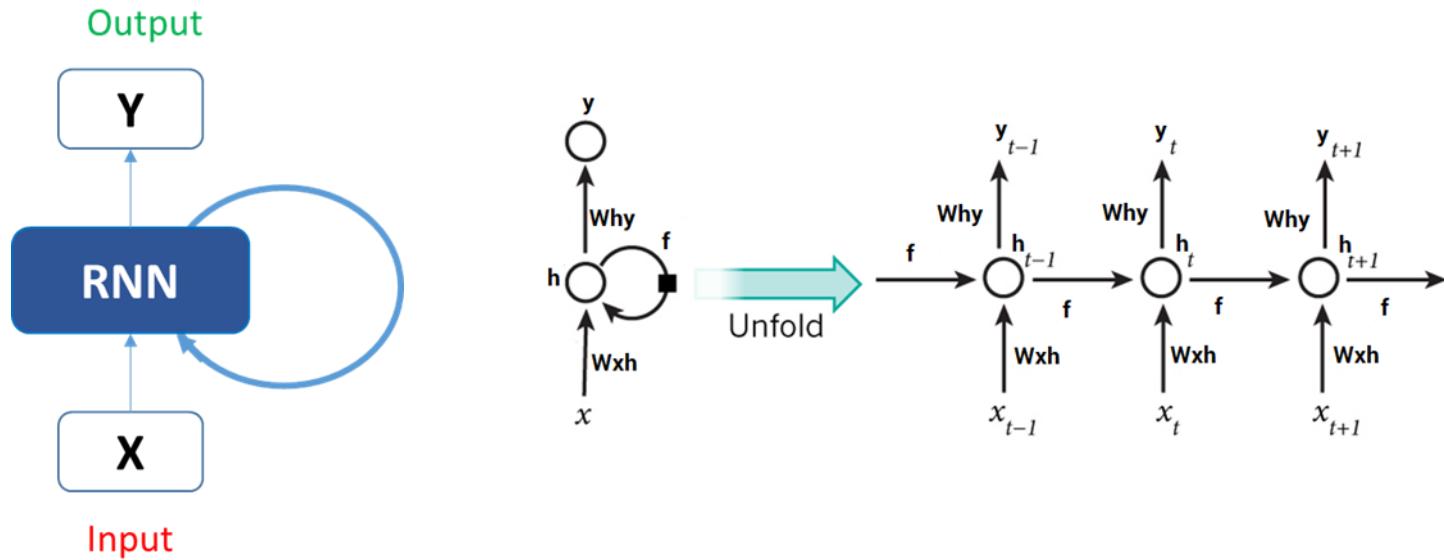
$a^{<t-1>}$: activation value in previous hidden state

$a^{<t>}$: activation value in current hidden state

$\hat{y}^{<t>}$: prediction (probability distribution for the next word)

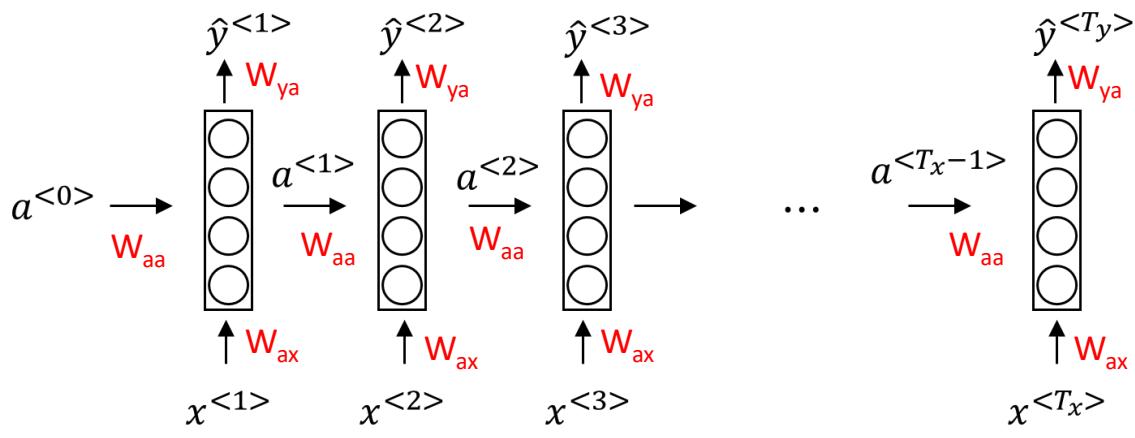
Fixed number of inputs at each time step
→ Solve the arbitrary sequence length

Recurrent Neural Networks



Fundamentals of Deep Learning – Introduction to Recurrent Neural Networks
DISHASHREE GUPTA, 2017

Forward Propagation



$$a^{<t>} = g(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g(W_{ya}a^{<t>} + b_y)$$

Coursera: Deep learning Specialization, Andrew Ng

$x^{<t>}$: current input
 $a^{<t-1>}$: activation in previous hidden state
 $a^{<t>}$: activation in current hidden state
 $y^{<t>}$: prediction

g_1, g_2 : activation function
 g_1 : tanh, ReLU
 g_2 : sigmoid, softmax

$$a^{<1>} = g_1(W_{aa}a^{<0>} + W_{ax}x^{<1>} + b_a)$$

$$\hat{y}^{<1>} = g_2(W_{ya}a^{<1>} + b_y)$$

All the parameters are shared across the different time steps → solve large number of parameters

Simplified RNN notation

$$a^{} = g(W_{aa}a^{} + W_{ax}x^{} + b_a)$$

$$\hat{y}^{} = g(W_{ya}a^{} + b_y)$$

$x^{}$: current input
 $a^{}$: activation in previous hidden state
 $a^{}$: activation in current hidden state
 $y^{}$: prediction

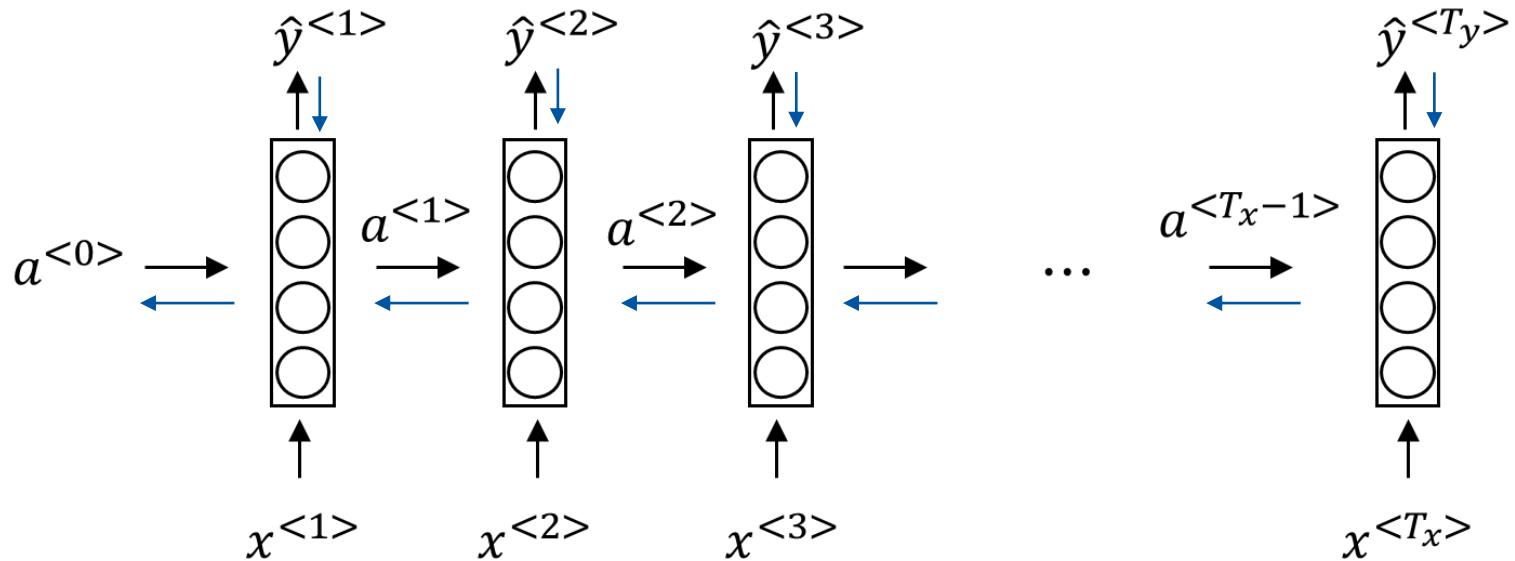
$$a^{} = g(W_a[a^{}, x^{}] + b_a)$$

$$\hat{y}^{} = g(W_ya^{} + b_y)$$

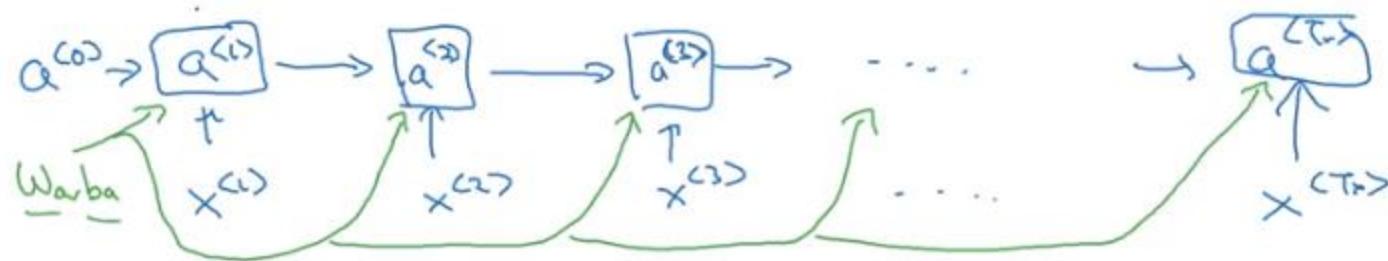
$$[W_{aa}; W_{ax}] = W_a$$

$$[W_{aa}; W_{ax}] \begin{bmatrix} a^{} \\ x^{} \end{bmatrix} = W_{aa}a^{} + W_{ax}x^{}$$

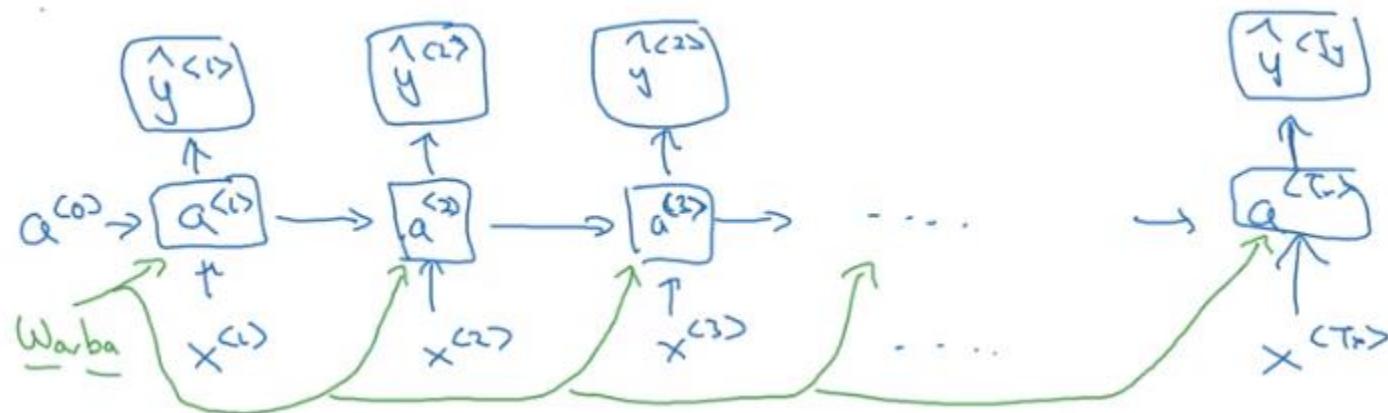
Forward propagation and backpropagation



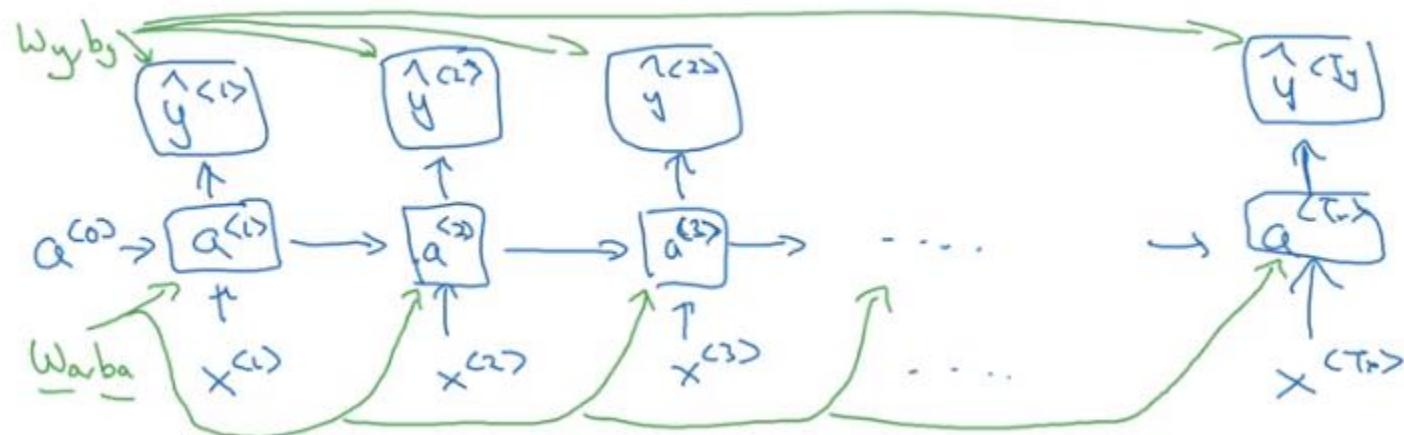
Backpropagation through time



Backpropagation through time

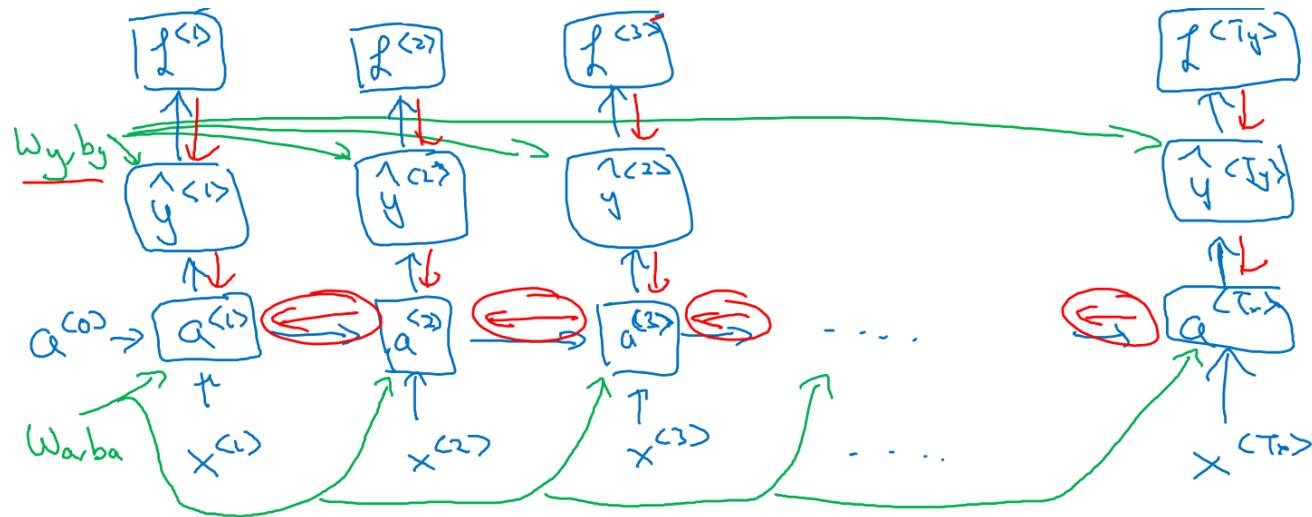


Backpropagation through time



Coursera: Deep learning Specialization, Andrew Ng

Backpropagation through time(BPTT)

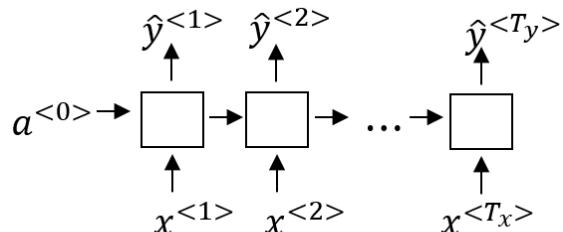


$$\text{Cross-entropy loss} \quad L^{(t)}(\hat{y}^{(t)}, y^{(t)}) = -y^{(t)} \log(\hat{y}^{(t)}) - (1 - y^{(t)}) \log(1 - \hat{y}^{(t)})$$

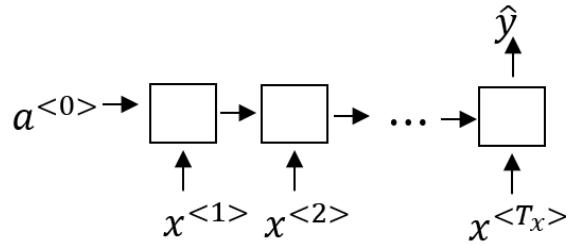
$$\text{Total loss} \quad L(\hat{y}, y) = \sum_{t=1}^{T_y} L^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

Coursera: Deep learning Specialization, Andrew Ng

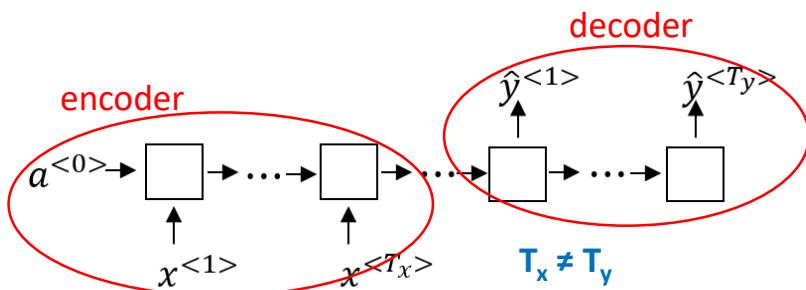
RNN types



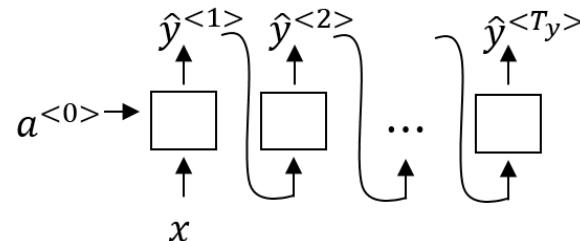
Many to many
Name entity recognition



Many to one
Sentiment classification

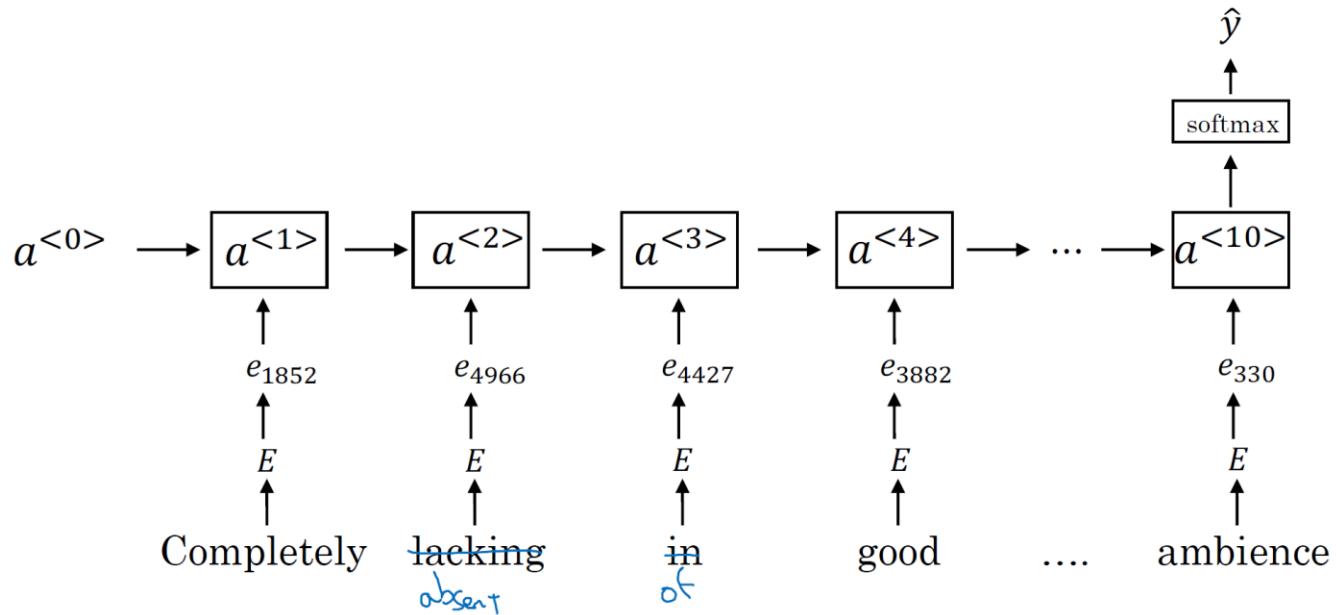


Many to many
Machine translation



One to many
Music generation

RNN for sentiment classification



Long range of dependence

The cat, which already ate ..., was full.

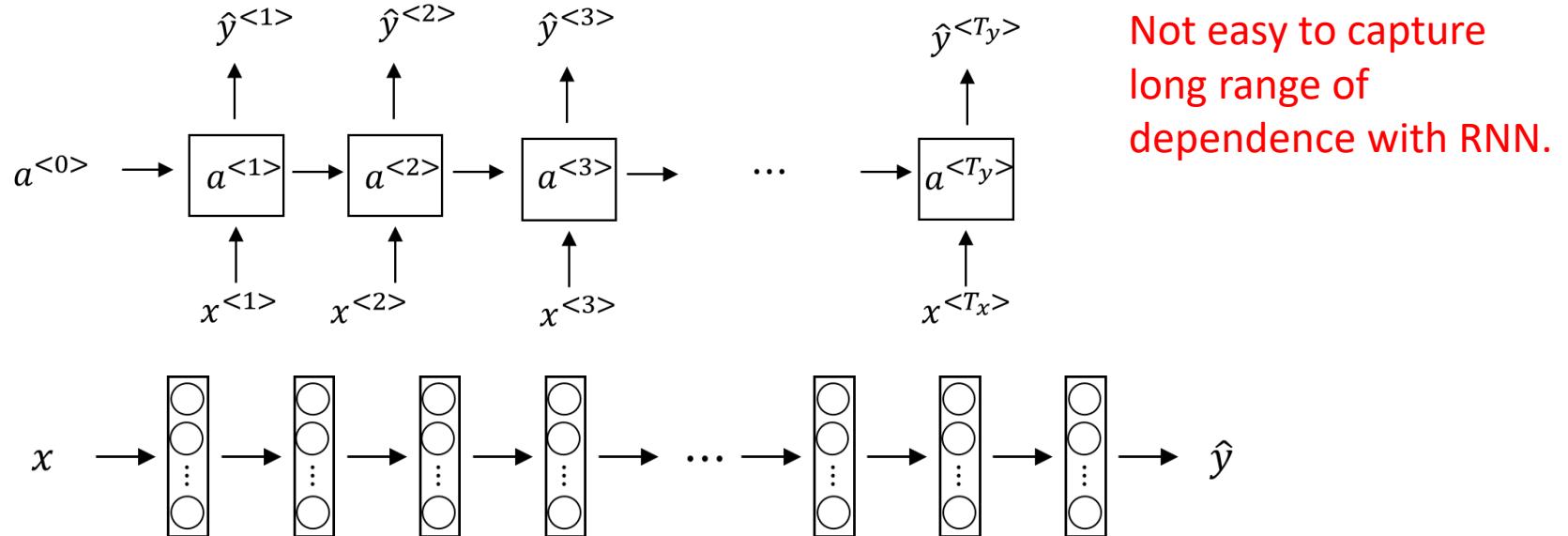
The cats, which already ate ..., were full.

long range dependencies is how the earlier parts of the sentence affects the later part of the sentence.

Vanishing gradients with RNNs

The cat, which already ate ..., was full.

The cats, which already ate ..., were full.



Gradients with RNNs

$$\left| \frac{\partial h_i}{\partial h_{i-1}} \right| < 1$$

- Gradients is vanishing after 3-4 time steps.
 - Not easy to capture long range of dependence between the earlier parts of the sentence and the later part of the sentence.

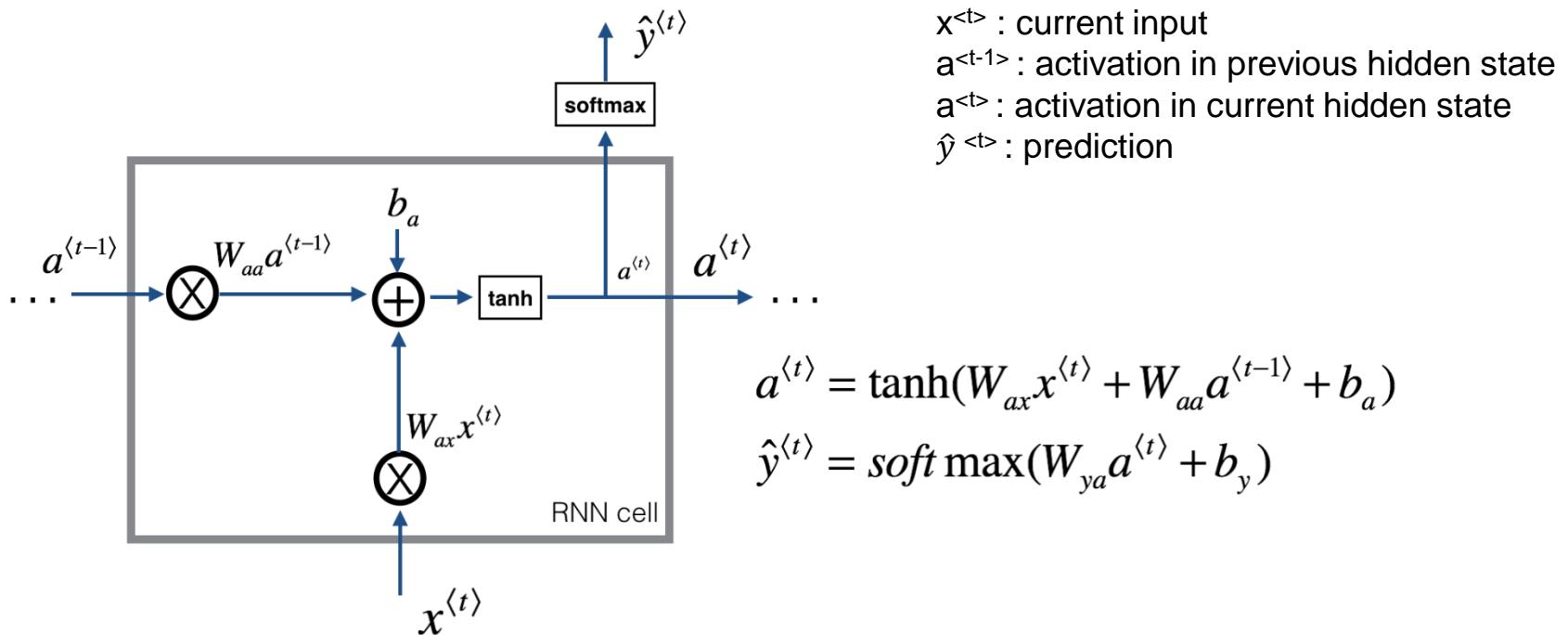
$$\left| \frac{\partial h_i}{\partial h_{i-1}} \right| > 1$$

- Exploding gradients:
 - Learning process becomes unstable
 - Gradient becomes NaN
 - Gradient clipping

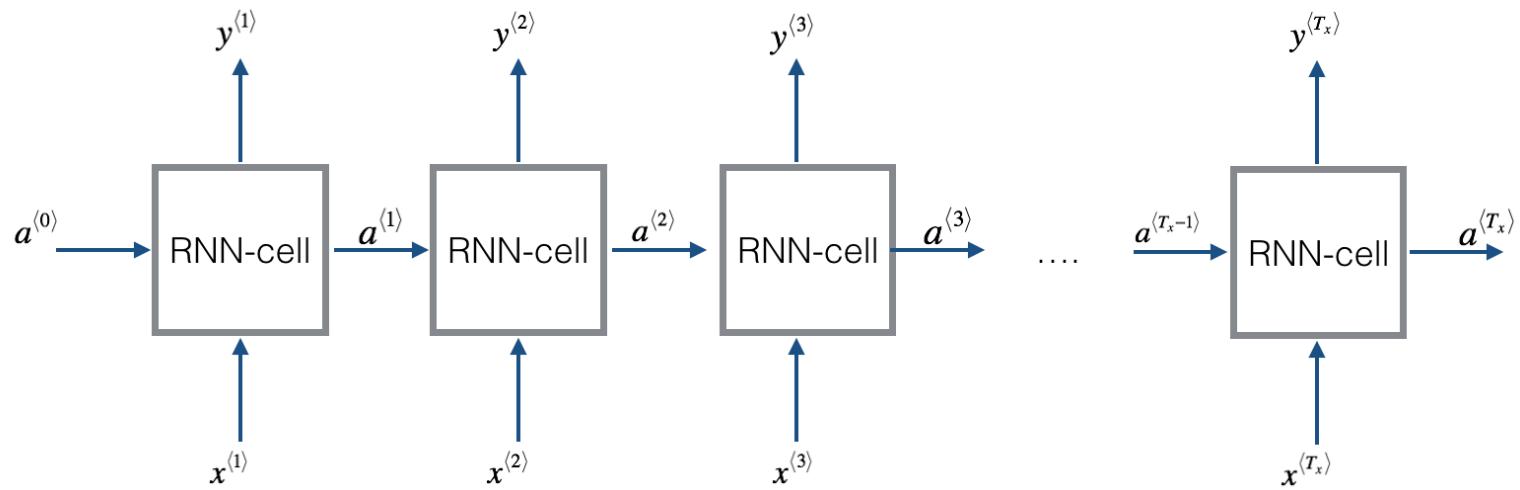
Coursera: Natural Language Processing, National Research University Higher School of Economics

Coursera: Deep learning Specialization, Andrew Ng

RNN unit



RNN forward pass



Coursera: Deep learning Specialization, Andrew Ng

Gated Recurrent Unit (simplified GRU)

The cat, which already ate ..., was full.

$$c^{<\leftrightarrow} = 1$$

$$c^{<\leftrightarrow} = 1$$

$$\Gamma_u = 1$$

$$\Gamma_u = 0$$

$$\Gamma_u = 0$$

$$\Gamma_u = 1$$

c : memory cell

$c^{<\leftrightarrow} = a^{<\leftrightarrow}$

$\tilde{c}^{<\leftrightarrow}$: candidate of activation

Γ_u : Update gate

$$\tilde{c}^{<t>} = \tanh(W_c[c^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

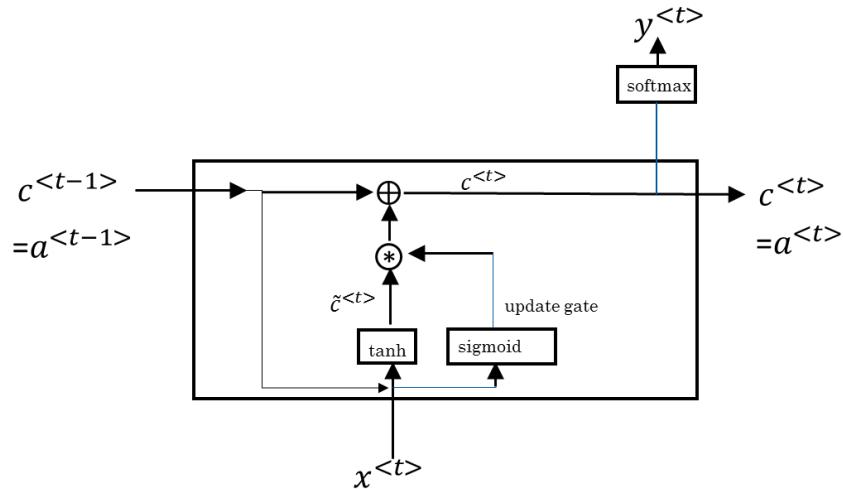
*:Element-wise multiplication

Coursera: Deep learning Specialization, Andrew Ng

On the properties of neural machine translation: Encoder-decoder approaches, Cho et al., 2014.

Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling, Chung et al., 2014.

Gated Recurrent Unit (GRU)



$$\tilde{c}^{<t>} = \tanh(W_c[c^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

$$\Gamma_u \approx 0, c^{<t>} = c^{<t-1>}$$

$c^{<t>}$ maintains many previous time-steps and helps vanishing gradient problem and allows long range dependencies.

Coursera: Deep learning Specialization, Andrew Ng

On the properties of neural machine translation: Encoder-decoder approaches, Cho et al., 2014.

Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling, Chung et al., 2014.

Full GRU

$$\tilde{h} \quad \tilde{c}^{<t>} = \tanh(W_c[\Gamma_u * c^{<t-1>}, x^{<t>}] + b_c) \quad \Gamma_r: \text{relevance gate}$$

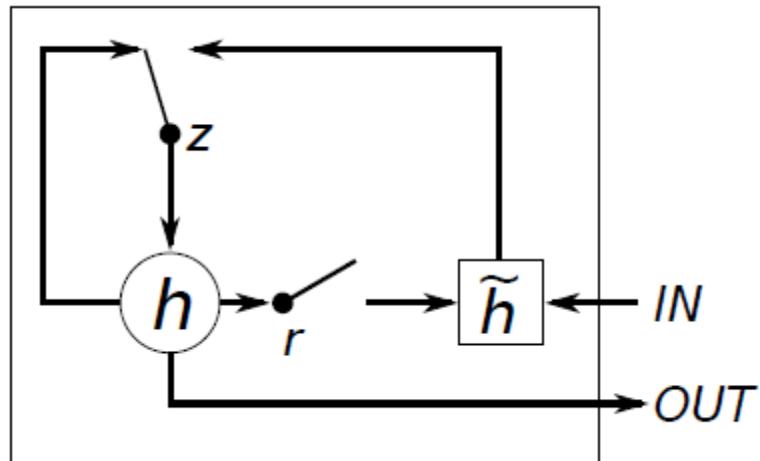
$$u \quad \Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

$$r \quad \Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r)$$

$$h \quad c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

Relevance gate is how relevant $c^{<t-1>}$ is to compute $\tilde{c}^{<t>}.$

Gated Recurrent Unit (GRU)



r : reset gate

z : update gates

h : activation in hidden state

\hat{h} : candidate activation in hidden state.

On the properties of neural machine translation: Encoder-decoder approaches, Cho et al., 2014.

Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling, Chung et al., 2014.

GRU vs. LSTM(long short term memory)

GRU

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

$$a^{<t>} = c^{<t>}$$

LSTM

$$\tilde{c}^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c)$$

update $\Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)$

forget $\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$

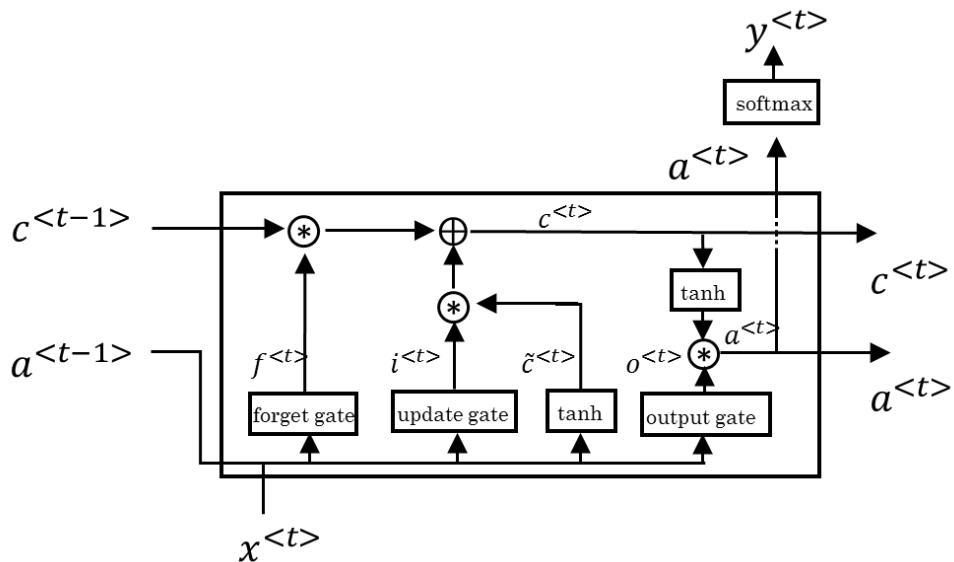
output $\Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o)$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

$$a^{<t>} = \Gamma_o * \tanh(c^{<t>})$$

Coursera: Deep learning Specialization, Andrew Ng
Long short-term memory, Hochreiter & Schmidhuber 1997.

LSTM (long short term memory) cell



$$\tilde{c}^{t-1} = \tanh(W_c[a^{t-1}, x^{t-1}] + b_c)$$

$$\Gamma_u = \sigma(W_u[a^{t-1}, x^{t-1}] + b_u)$$

$$\Gamma_f = \sigma(W_f[a^{t-1}, x^{t-1}] + b_f)$$

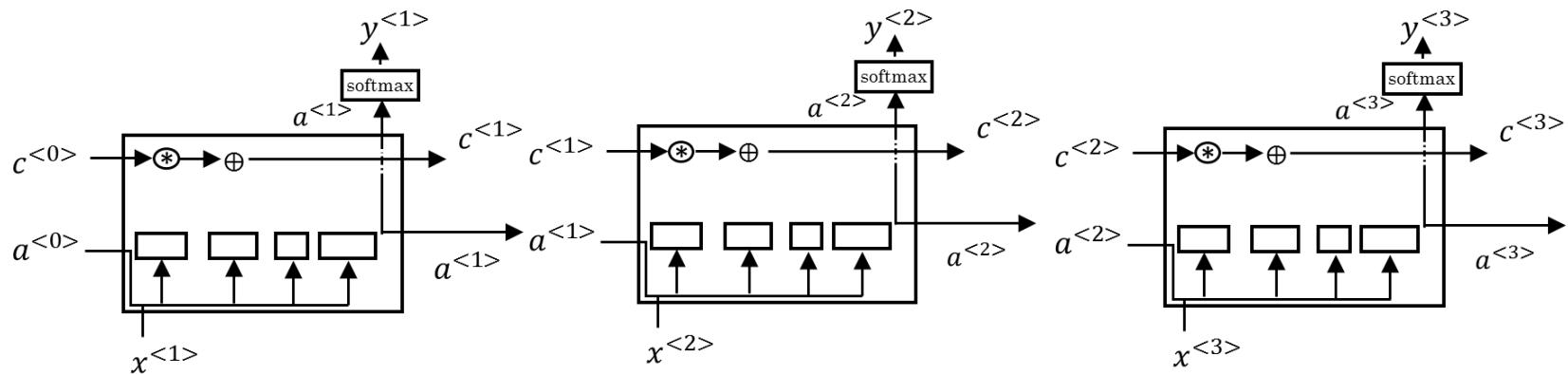
$$\Gamma_o = \sigma(W_o[a^{t-1}, x^{t-1}] + b_o)$$

$$c^{t-1} = \Gamma_u * \tilde{c}^{t-1} + \Gamma_f * c^{t-1}$$

$$a^{t-1} = \Gamma_o * \tanh(c^{t-1})$$

Coursera: Deep learning Specialization, Andrew Ng

LSTM forward



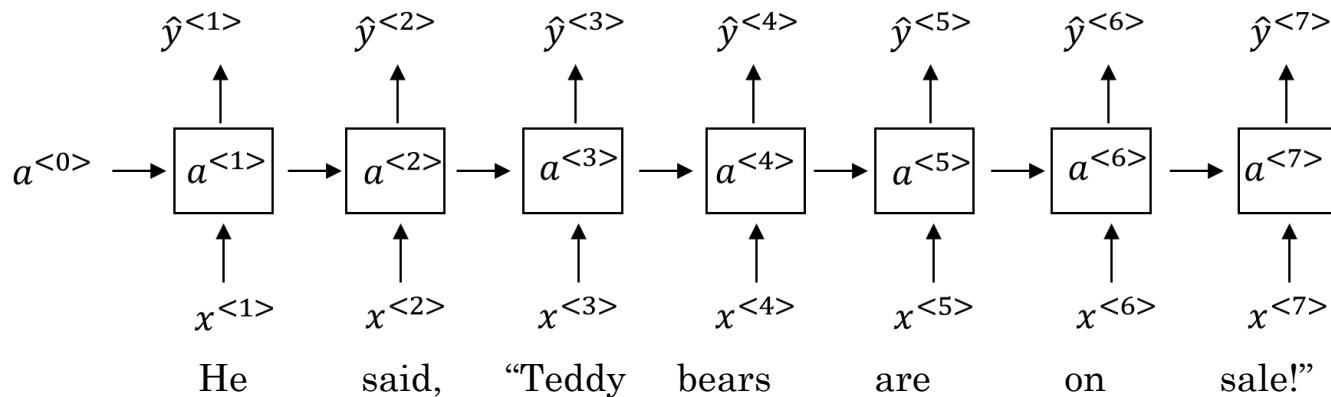
$$\Gamma_f \approx 0, c^{<3>} = c^{<0>}$$

$c^{<t>}$ maintains many previous time-steps and helps vanishing gradient problem and allows long range dependencies.

Unidirectional RNN

He said, “Teddy bears are on sale!”

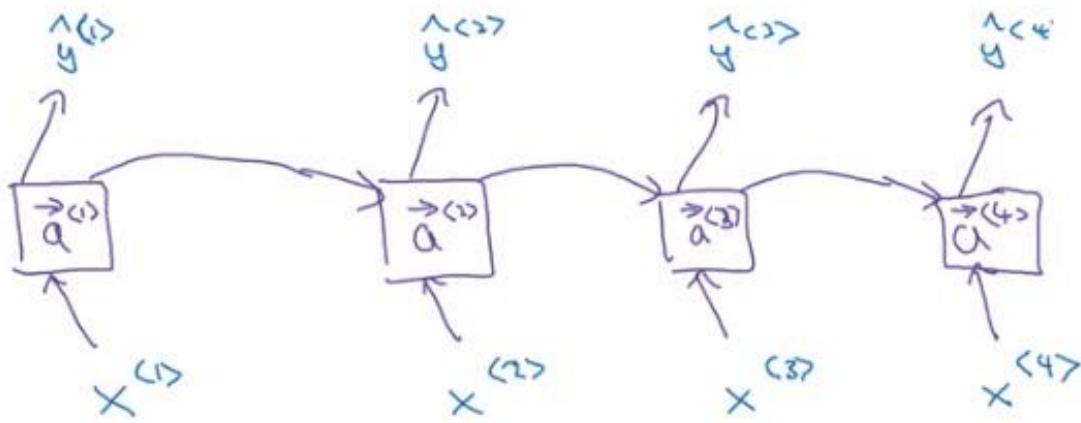
He said, “Teddy Roosevelt was a great President!”



Coursera: Deep learning Specialization, Andrew Ng

Bidirectional RNN

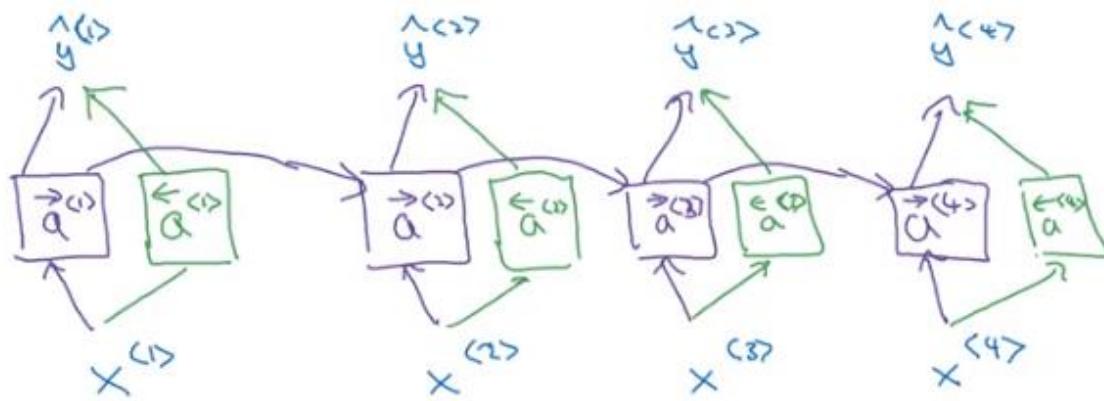
Getting information from the future



Coursera: Deep learning Specialization, Andrew Ng

Bidirectional RNN

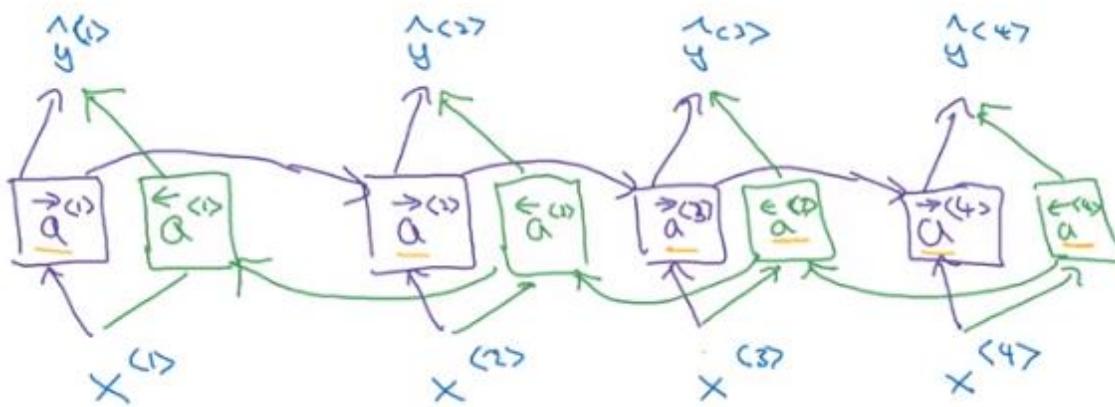
Getting information from the future



Coursera: Deep learning Specialization, Andrew Ng

Bidirectional RNN

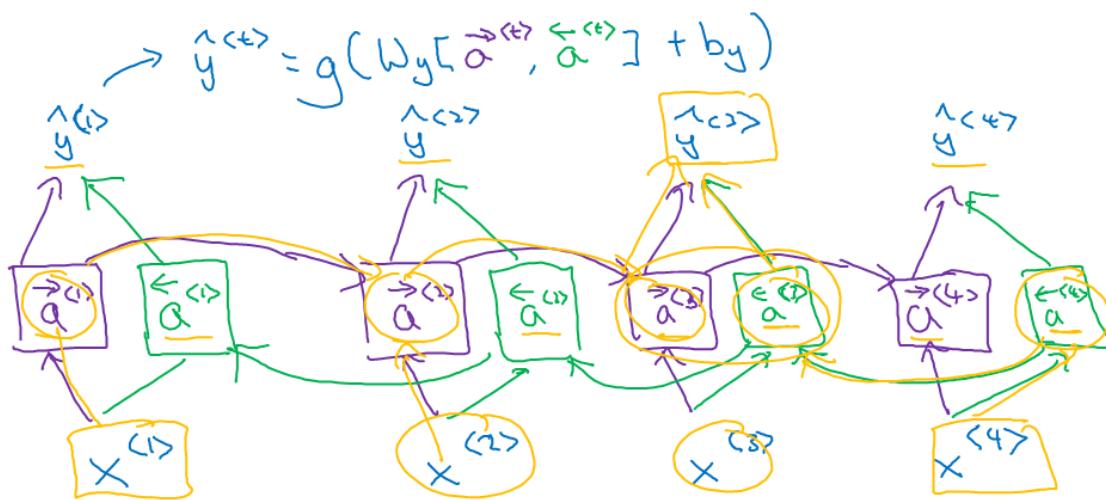
Getting information from the future



Coursera: Deep learning Specialization, Andrew Ng

Bidirectional RNN

Getting information from the future

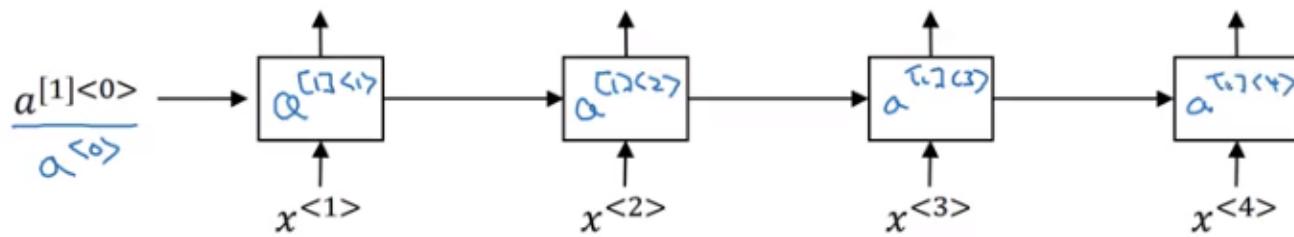


BRNN with GRU/LSTM

Entire sentence is needed
to build BRNN.

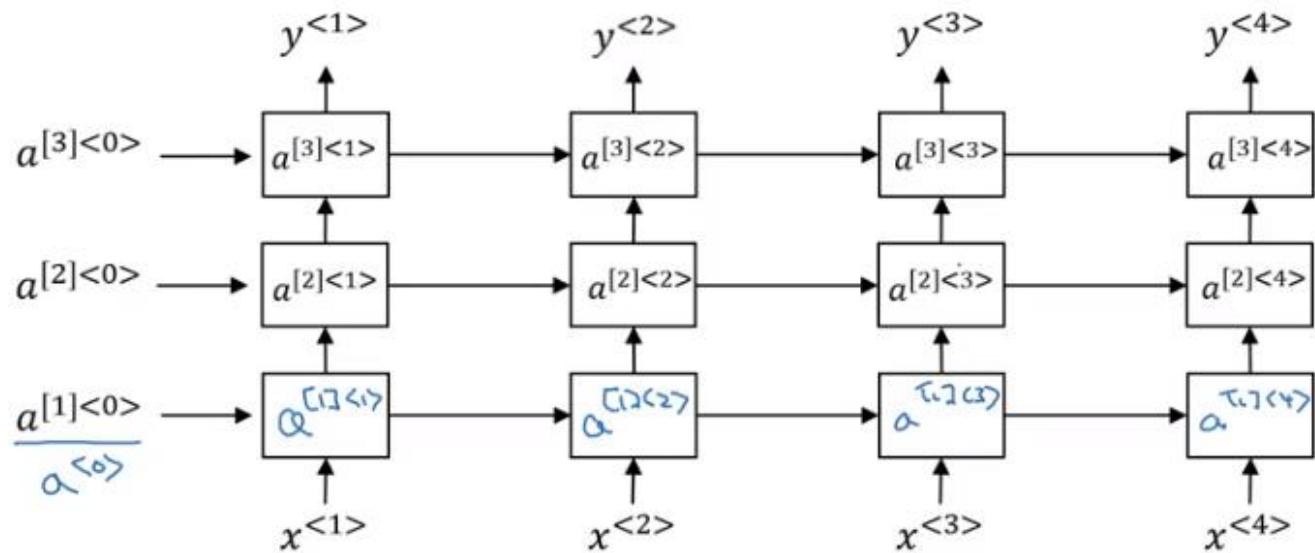
Coursera: Deep learning Specialization, Andrew Ng

Deep RNNs



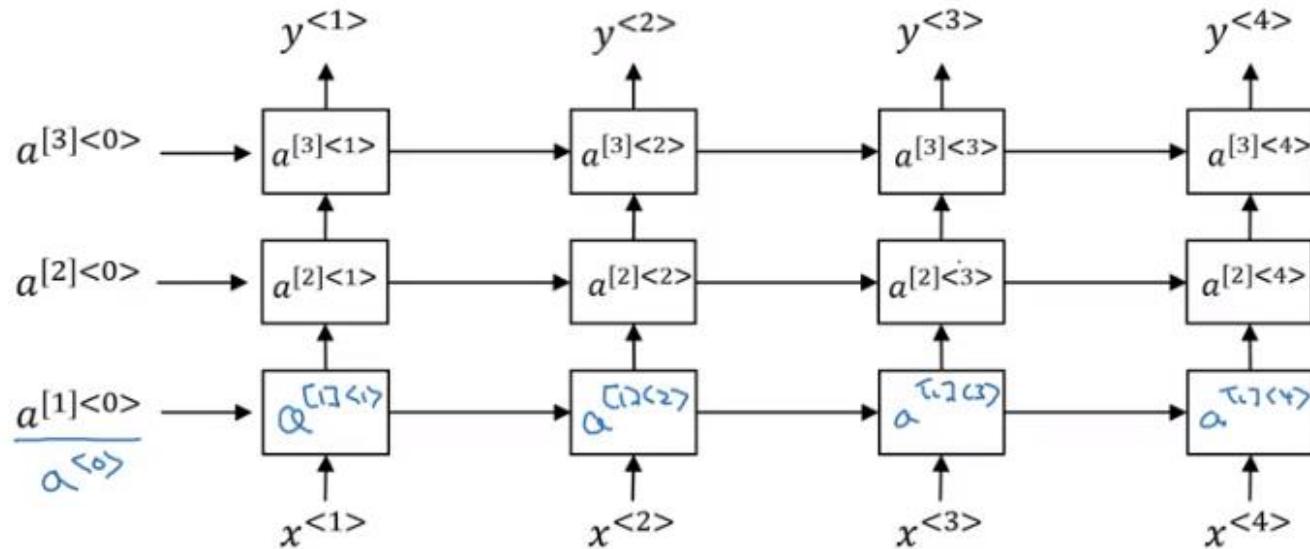
Coursera: Deep learning Specialization, Andrew Ng

Deep RNNs



Coursera: Deep learning Specialization, Andrew Ng

Deep RNNs



$$a^{[2]<3>} = g(W_a^{[2]} [a^{[2]<2>}, a^{[1]<3>}] + b_a^{[2]})$$

Coursera: Deep learning Specialization, Andrew Ng