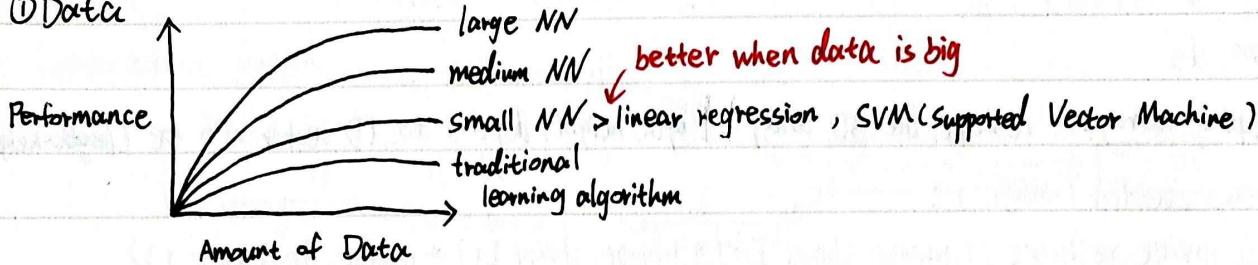


Deep Learning.ai Summary

1- Neural Network and Deep Learning

1.1 Deep Learning is taking off because:

① Data



② Computation

GPUs, powerful CPUs, Distributed Computing, ASICs

③ Algorithm

Creative algorithms has appeared that changed the way NN works

Ex: RELU is much better than SIGMOID

1.2 Logistic Regression

① A most popular binary classification technique.

② We have an input X (image) and we want to know the probability that the image belongs to class 1 (i.e. a cat). For a given X vector, the output will be:

$$y = w^T x + b \quad * W \text{ is } n \times \text{dimensional vector, } b \text{ is a real number}$$

To let the result y make sense as a probability between 0 and 1, Using sigmoid function

$$\hat{y} = \sigma(w^T x + b) \quad \sigma(z) = \frac{1}{1 + e^{-z}} \quad \begin{array}{ll} \text{if } z \rightarrow \infty, e^z \rightarrow \infty \Rightarrow \sigma(z) = 1 \\ \text{if } z \rightarrow -\infty, e^z \rightarrow 0 \Rightarrow \sigma(z) = 0 \end{array}$$

③ Logistic Regression Cost Function

To train the parameters w and b , we need a cost function

$$J(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2 \quad * \text{Problem: non-convex, resulting in multiple local optima}$$

a convex function: $J(\hat{y}, y) = -[y \log \hat{y} + (1-y) \log(1-\hat{y})]$ * This is a loss function that

measures a single training example, to measure the entire training set, we need a cost function:

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

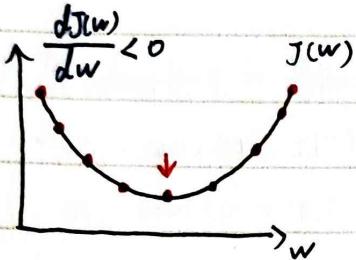
1.3. Gradient Descent

To train / learn w and b in such a way that the cost function is minimized

Steps: ① Initialize w and b (usually "0")

② Take a step in the steepest downhill direction

③ Repeat step ② until global optimum is achieved



The updated equations for the parameters of logistic regression are:

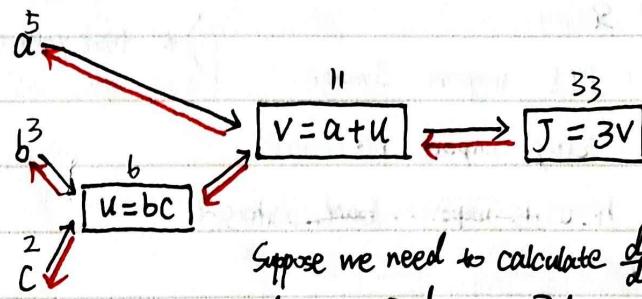
$$w := w - \alpha \frac{dJ(w, b)}{dw} \quad b := b - \alpha \frac{dJ(w, b)}{db}$$

* α is the learning rate that controls how big a step we should take after each iteration.

1.4. Computation Graph

$$J(a, b, c) = 3(a + \underbrace{bc}_v) \quad \underbrace{u}_{\substack{u \\ v}} \quad J = \beta v$$

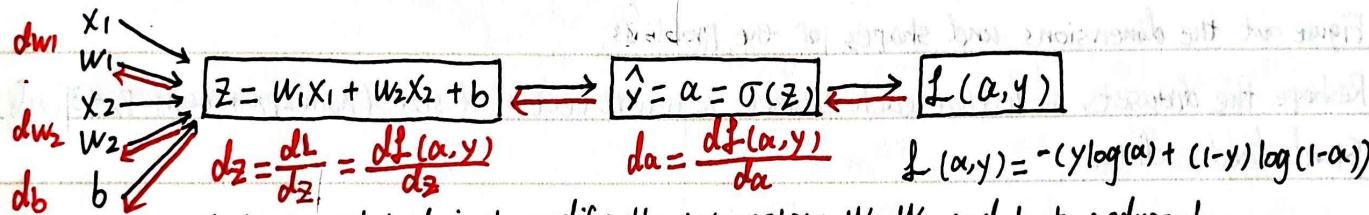
$$u = bc \\ v = a + u \\ J = \beta v$$



Suppose we need to calculate $\frac{dJ}{da}$: (chain rule)

$$1.5. \text{ Logistic Regression Gradient Descent (a single training example)} \quad \textcircled{1} \frac{dJ}{dv} = 3 \quad \textcircled{2} \frac{dv}{da} = 1 \quad \textcircled{3} \frac{dJ}{da} = \frac{dJ}{dv} \times \frac{dv}{da} = 3$$

For two features (x_1, x_2) , the computation graph:



what we want to do is to modify the parameters w_1, w_2 and b to reduce loss

$$\text{Steps to calculate } dw_1: \textcircled{1} da = \frac{df}{dx} = -\frac{y}{a} + \frac{1-y}{1-a}$$

$$\text{(Similarly, we can } \textcircled{2} dz = \frac{df}{da} \times \frac{da}{dz} = [-\frac{y}{a} + \frac{1-y}{1-a}] \times [a(1-a)] = a-y \text{)}$$

$$\text{calculate } dw_2 \text{ and } b \quad \textcircled{3} dw_1 = [\frac{df}{da} \times \frac{da}{dz}] \times \frac{dz}{dw_1} = (a-y) \times \frac{dz}{dw_1}$$

$$w_1 := w_1 - \alpha dw_1 \quad w_2 := w_2 - \alpha dw_2 \quad b := b - \alpha db$$

1.6 Gradient Descent on 'm' Examples

Using logistic regression in Python would be slow because you have to run the loop 'm' times to get the output

To get rid of 'for' loops in codes, we use vectorization

$$\textcircled{1} \quad Z = \underbrace{\text{np.dot}(W.T, X) + b}_{W^T X} \quad * X \text{ contains the features for all the training examples while } W \text{ is the Coefficient matrix for these examples}$$

$$\textcircled{2} \quad A = \text{np. exp}(-Z) \quad * A \text{ is the sigmoid of } Z$$

$$\textcircled{3} \quad dz = A - Y \quad * \text{Calculate the loss and then use backpropagation to minimize the loss}$$

$$\textcircled{4} \quad dw = \text{np. dot}(X, dz.T) / m \quad * \text{Finally, calculate the derivative of the parameters and update them.}$$

$$db = dz.sum() / m$$

$$w = w - \alpha dw$$

$$b = b - \alpha db$$

1.7. Broadcasting in Python

• obj.sum(axis=0) sums the columns

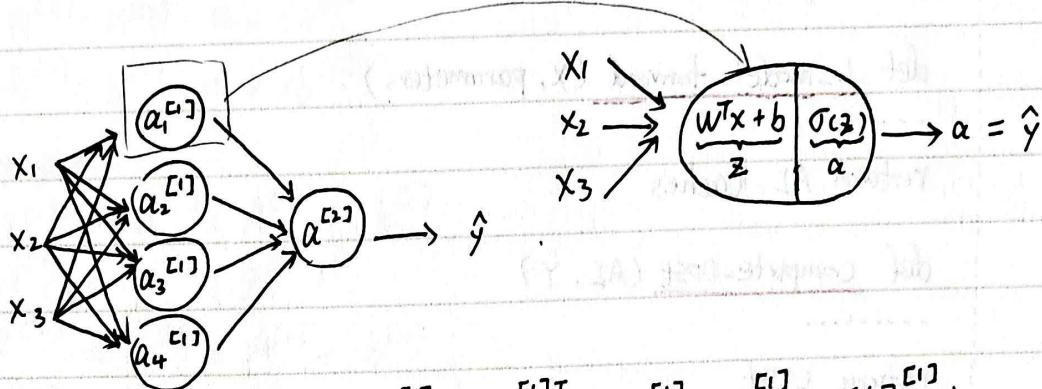
obj.sum(axis=1) sums the rows

General principle: $\begin{matrix} (m, n) \\ \text{matrix} \end{matrix} \xrightarrow{\quad} \begin{matrix} (1, n) \Rightarrow (m, n) \\ (m, 1) \Rightarrow (m, n) \end{matrix}$

Example: $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix} = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$

1.8 Neural Network

2 layers



The equations for the first hidden layer

$$\left\{ \begin{array}{l} z_1^{[1]} = W_1^{[1]T} x + b_1^{[1]}, \quad a_1^{[1]} = \sigma(z_1^{[1]}) \\ z_2^{[1]} = W_2^{[1]T} x + b_2^{[1]}, \quad a_2^{[1]} = \sigma(z_2^{[1]}) \\ z_3^{[1]} = W_3^{[1]T} x + b_3^{[1]}, \quad a_3^{[1]} = \sigma(z_3^{[1]}) \\ z_4^{[1]} = W_4^{[1]T} x + b_4^{[1]}, \quad a_4^{[1]} = \sigma(z_4^{[1]}) \end{array} \right.$$

$$\left. \begin{array}{l} z^{[1]} = W^{[1]} x + b^{[1]}, \\ A^{[1]} = \sigma(z^{[1]}) \\ z^{[2]} = W^{[2]} x + b^{[2]}, \\ A^{[2]} = \sigma(z^{[2]}) \end{array} \right\}$$

The vectorized form for calculating the output

* this will significantly reduce the computation time

1.9. Activation Function

① Sigmoid:

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

Used in the output layer
for binary classification

Pros

If z is either very large or very small, the gradient of the derivative

② tanh

$$a = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

Better than Sigmoid
: effectively center your data

or the slope of this function becomes very small

③ ReLu

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

Update parameters faster

as slope is 1 when $z > 0$

slope = 0 when $x < 0$

✓ (most time choose softmax)

* Softmax function : multiclass logistic regression

Sigmoid function : two-class logistic regression

2.1

Gradient Descent for NN :

Steps : Repeat :

Compute predictions ($y^{(i)}$, $i=1, \dots, m$)

Get derivatives : $dW^{[1]}, db^{[1]}, dW^{[2]}, db^{[2]}$

Update : $W^{[1]} = W^{[1]} - \alpha * dW^{[1]}$

$$b^{[1]} = b^{[1]} - \alpha * db^{[1]}$$

$$W^{[2]} = W^{[2]} - \alpha * dW^{[2]}$$

$$b^{[2]} = b^{[2]} - \alpha * db^{[2]}$$

Forward propagation :

$$Z^{[1]} = W^{[1]} * A^{[0]} + b^{[1]}$$

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

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

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

Backpropagation :

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

$$dW^{[2]} = (dZ^{[2]} * A^{[1]}.T) / m$$

$$db^{[2]} = \text{Sum}(dZ^{[2]}) / m$$

$$dZ^{[1]} = (W^{[2]}.T * dZ^{[2]}) * g'^{[1]}(Z^{[1]})$$

$$dW^{[1]} = (dZ^{[1]} * A^{[0]}.T) / m$$

$$db^{[1]} = \text{Sum}(dZ^{[1]}) / m$$

2.2

Random Initialization

$$W^{[1]} = np.random.randn((2, 2)) * 0.01$$

* Should NOT initialize to 0

$$b^{[1]} = np.zeros((2, 1))$$

2.3. Forward Propagation in a Deep Neural Network

$$Z^{[l]} = W^{[l]} A^{[l-1]} + B^{[l]}$$

* have been vectorized for

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

'm' training examples

2.4. The generalized form of dimensions of W, b and their derivatives :

$$W^{[l]} = (n^{[l]}, n^{[l-1]})$$

* m is the number of

$$b^{[l]} = (n^{[l]}, 1)$$

training examples

$$dW^{[l]} = (n^{[l]}, n^{[l-1]})$$

$$db^{[l]} = (n^{[l]}, 1)$$

Dimension of $Z^{[l]}, A^{[l]}, dZ^{[l]}, dA^{[l]} = (n^{[l]}, m)$

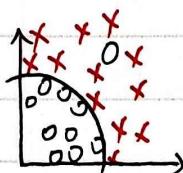
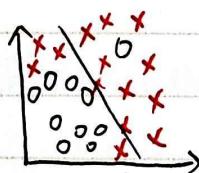
2.5 Vectorized equations for backward propagation in deep NN

$$\begin{aligned} dz^{[l]} &= dA^{[l]} * g'(z^{[l]}) \\ dw^{[l]} &= \frac{1}{m} * (dz^{[l]} * A^{[l-1]}.T) \\ db^{[l]} &= \frac{1}{m} * np.sum(dz^{[l]}, axis=1, keepdims=True) \\ dA^{[l-1]} &= w^{[l]}.T * dz^{[l]} \end{aligned}$$

(course 2)

3. Improving Neural Network - Hyperparameter Tuning, Regularization, and More

3.1.



training set error (%)	1	15	15	0.5
dev set error (%)	11	16	30	1
result	High Variance overfit	High Bias Underfit	High Variance and High Bias Underfit	Low Bias and Low Variance Reasonable

Q: Does the model have high bias?

High training error results in high bias.

Solutions: Try bigger networks, try different NN architectures, train models for a longer period of time

Q: Does the model have high variance?

High dev set error results in high variance

Solutions: Get more data, use regularization, try different NN architectures

3.2

Regularization

The cost function for a NN: $J(w^{[1]}, b^{[1]}, \dots w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^n \frac{1}{2} (\hat{y}^{(i)}, y^{(i)})$

Adding a regularization term: $J(w^{[1]}, b^{[1]}, \dots w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^n \frac{1}{2} (\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{i=1}^L \|w^{[i]}\|^2$

* L₂ regularization

$$* \|W^{[i]}\|^2 = \sum_{j=1}^{n_{[i-1]}} \sum_{k=1}^{n_{[i]}} (W_{jk}^{[i]})^2$$

Regularized form of update equations:

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

$$w^{[l]} = w^{[l]} - \alpha [dw^{[l]} + \frac{\lambda}{m} W^{[l]}]$$

* That is why L₂ regularization is also known as weight decay

Other regularization methods:

Dropout Regularization only use during training time, both during forward and backward propagation

Data Augmentation

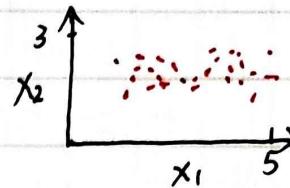
Early Stopping

3.3. Setting up your optimization problem

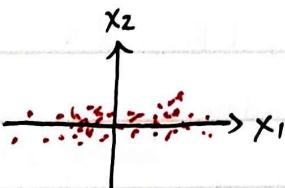
① Normalizing inputs

What we want to do is to :

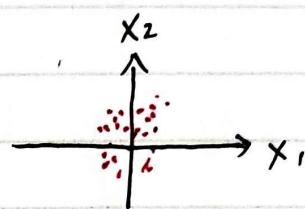
Ex: Changing the 2 input features with scatter plot looks like :



① Subtract the mean
from the input features: $\begin{cases} \mu = \frac{1}{m} \sum_{i=1}^m x^{(i)} \\ x := x - \mu \end{cases}$



② normalize the variance: $\begin{cases} \sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2 \\ x := \frac{x - \mu}{\sigma} \end{cases}$



Normalizing the inputs makes the cost function symmetric. This makes it easier for the gradient descent algorithm to find the global minima more quickly

3.4. How to initialize weights in deep NN to avoid vanishing / exploding gradients ?

Initialize weights randomly.

Reason: Break symmetry, making sure different hidden units learn different patterns

3.5. Optimization Algorithm

1). Mini-Batch Gradient Descent (make faster)

For a large training set, typical mini-batch sizes are : 64, 128, 256, 512

2). Exponentially weighted averages

$$V_t = [\beta * V_{t-1} + (1-\beta) * \theta_t] / (1 - \beta t) \quad \text{a bias correction term}$$

* This takes a lot less memory as we are overwriting the previous values

3). Gradient Descent with Momentum (reduce oscillations of Mini-Batch gradient descent)

Momentum takes into account the past gradients to smooth out the update.

The stored 'direction' of the previous gradients is the exponentially weighted average of the gradient on the previous steps

* Choose of β : 0.9 but the optimal β should be tuned by several values

4). Adam

① Most fast

② Low memory requirements

③ Work well even with little tuning of hyperparameters (except α)

Convolutional Neural Network

1. Edge Detection Example 灰度图, 每一格代表灰度值

$$3 \times 1 + 1 \times 1 + 2 \times 1 + 0 \times 0 + 5 \times 0 + 7 \times 0 + 1 \times (-1) + 8 \times (-1) + 2 \times (-1) = -5$$

3	1	0	0	1	-1
1	1	5	0	8	-1
2	1	1	0	2	-1
0	1	1	3	1	1
4	2	1	6	2	8
2	4	5	2	3	9

6×6

filter

$$\begin{matrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{matrix} = \begin{matrix} -5 & -4 & 0 & 8 \\ -10 & -2 & 2 & 3 \\ 0 & -2 & -4 & -7 \\ -3 & -2 & -3 & -16 \end{matrix}$$

3×3 4×4

Convolution, Not multiplication

python: cv2.filter2D

tensorflow: tf.nn.conv2d

2. Vertical and Horizontal Edge Detection

bright \leftarrow

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

→ dark

vertical

bright

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

→ dark

Horizontal

1	0	-1
2	0	-2
1	0	-1

Sobel filter

Adv: put a little more weight to the central row, which makes it more robust

3	0	-3
10	0	-10
3	0	-3

Scharr filter

W_1	W_2	W_3
W_4	W_5	W_6
W_7	W_8	W_9

Just learn them using back propagation

3. Padding : to solve two problems :

① Shrink output

② Throw away a lot of information from the edges of the image.

Padding:

from 6×6 to 8×8

Input: $n \times n$
 Filter: $f \times f$
 Output: $(n-f+1) \times (n-f+1)$

Padding = p

Input: $n \times n$
 Filter: $f \times f$
 Output: $(n+2p-f+1) \times (n+2p-f+1)$

Valid convolutions: no padding

Same convolutions: the output size is the same as the input size

$$\Rightarrow n + 2p - f + 1 = n \Rightarrow p = \frac{f-1}{2}$$

* f is usually odd

4. Strided Convolutions

Suppose we choose a stride of 2. While convoluting through the image, we will take two steps - both in the horizontal and vertical directions separately, the dimensions for stride s will be:

• Input: $n \times n$ Stride: s
 • Padding: p Filter size: $f \times f$ Output: $[(n+2p-f)/s+1] \times [(n+2p-f)/s+1]$

Stride helps to reduce the size of the image

5. Convolutions over volume

① Instead of 2-D image, we have a 3-D image of shape $6 \times 6 \times 3$

• Input: $6 \times 6 \times 3$
 • Output: $3 \times 3 \times 3$ channel * The number of channels in Input and output should be the same

② Instead of a single filter, we can use multiple filters as well.

Ex: use two $3 \times 3 \times 3$ filters. Then the output should be: $4 \times 4 \times 2$

Input: $6 \times 6 \times 3$

Generalized dimensions:

Input: $n \times n \times n_c$ * n_c : num of channels

Filter: $f \times f \times n_c'$ * n_c' : num of filters

Padding: p

Stride: s

Output: $[(n+2p-f)/s+1] \times [(n+2p-f)/s+1] \times n_c'$

6. One layer of a convolutional network

No matter how big the image is, the parameters only depend on the filter size.

Summary of notation:

If layer l is a convolution layer:

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

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

$p^{[l]}$ = padding

Output: $n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

$s^{[l]}$ = stride

$n_{H/W}^{[l]} = \frac{n_{H/W}^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}}$

$n_C^{[l]}$ = num of filters

$+ 1$

when using mini-batch gradient descent

Each filter is : $f^{[l]} \times f^{[l]} \times n_c^{[l-1]}$

Activations: $\alpha^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$

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

Bias: $n_c^{[l]} = (1, 1, 1, n_c^{[l]})$

\downarrow num of training examples

$A^{[l]} \rightarrow M \times n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$

Pooling Layers

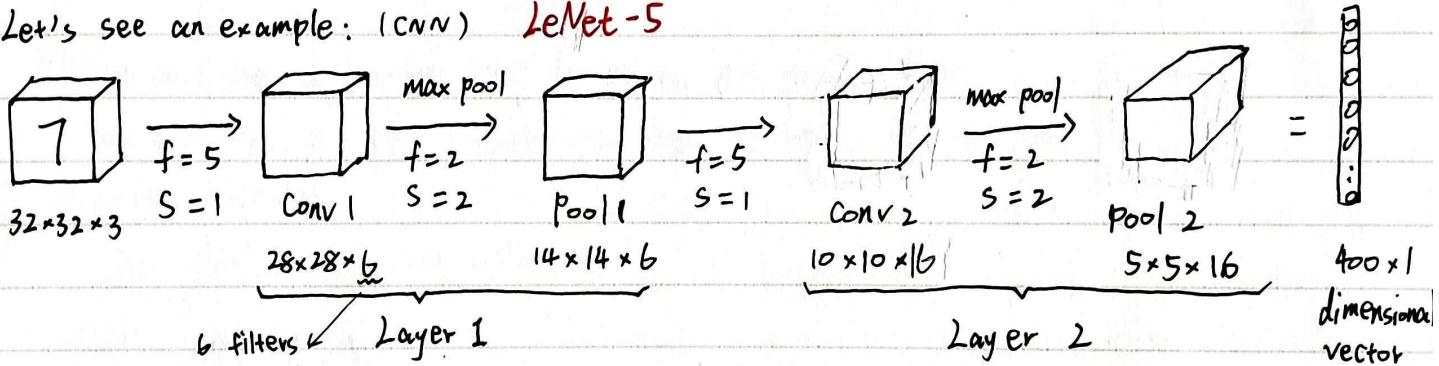
To reduce the size of the inputs and hence speed up the computation

- ① Max pooling
- ② Average pooling

Why Convolutions?

1. Parameter sharing → speed up the training of the model
2. Sparsity of connections

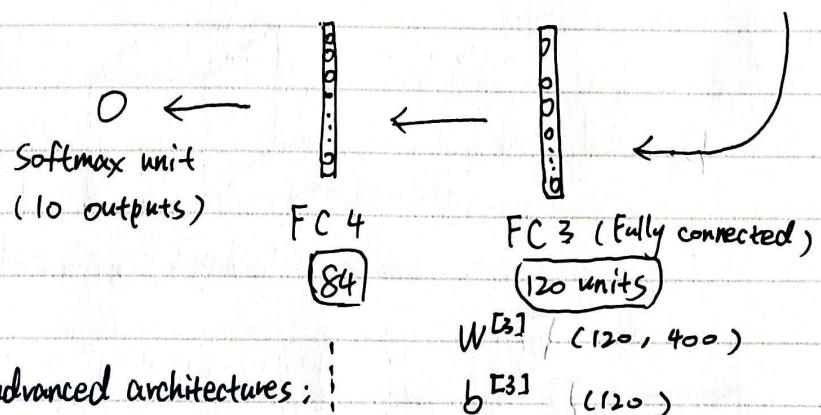
Let's see an example: (CNN) LeNet-5



* Choose of hyperparameters:

Do some research on literature

* Deeper NN, n_w & $n_h \downarrow$, while $n_c \uparrow$



| Three classic architectures: | More advanced architectures: |

LeNet-5

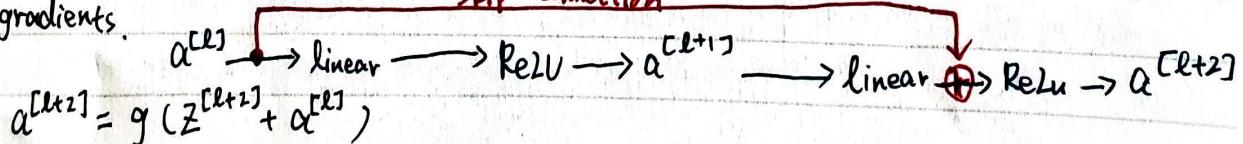
much larger network

AlexNet

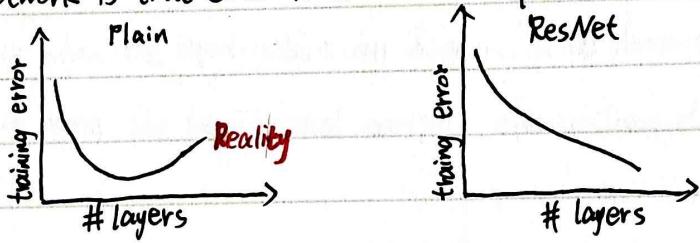
VGG-16 simplify hyperparameters

ResNet

ResNet: Using residual blocks. Skip connections where we take activations from one layer and feed it to another layer that is even more deeper in the network. \Rightarrow Solving the problem of vanishing and exploding gradients.



The benefit of training a residual network is that even if we train deeper networks, the training errors does not increase



To calculate activation using

$$\text{a residual block: } a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$$

$$a^{[l+2]} = g(W^{[l+2]} * a^{[l+1]} + b^{[l+2]} + a^{[l]})$$

Now, say $W^{[l+2]} = 0$ and the bias $b^{[l+2]}$ is also 0, then:

$$a^{[l+2]} = g(a^{[l]})$$

Networks in Networks and 1×1 convolutions

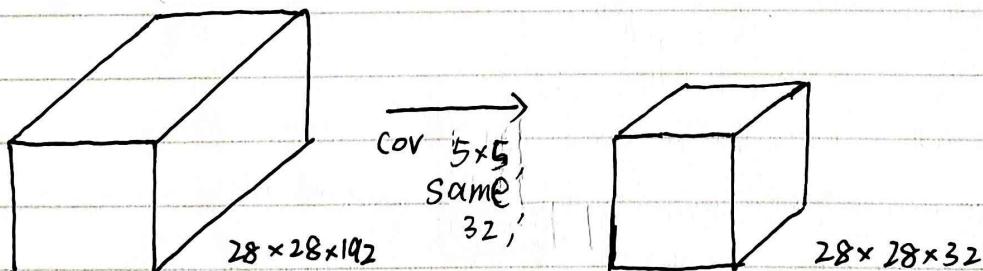
- The basic idea of using 1×1 Convolution is to reduce the number of channels from the image
We use a pooling layer to shrink the height and width of the image

* You can use a pooling layer to reduce n_h , n_w , and n_c .

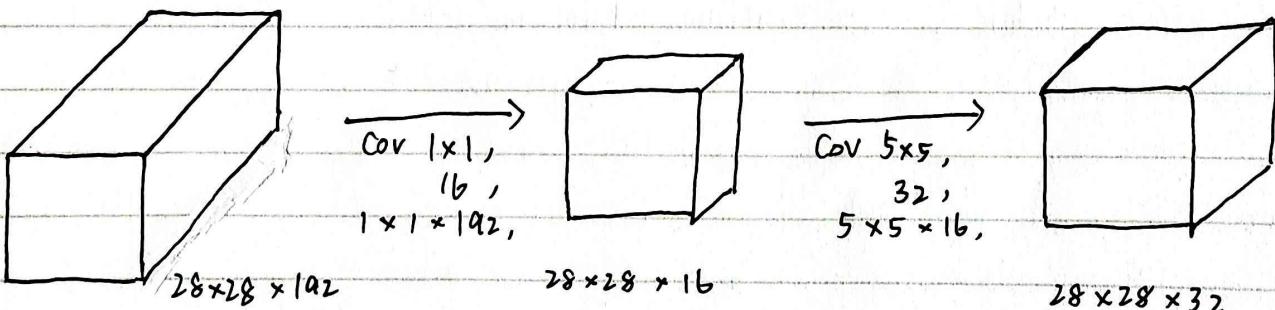
You can use a 1×1 convolutional layer to reduce n_c , but not n_h , n_w

Inception Network

Key idea: significant reduction



$$\text{Number of multipliers} = 28 \times 28 \times 32 \times 5 \times 5 \times 192 = 120 \text{ million}$$



$$\text{Number of multiplies for first convolution: } 28 \times 28 \times 16 \times 1 \times 1 \times 192 = 2.4 \text{ million}$$

$$\text{for second convolution: } 28 \times 28 \times 32 \times 5 \times 5 \times 16 = 10 \text{ million}$$

$$\text{Total: } 2.4 \text{ million} + 10 \text{ million} = 12.4 \text{ million}$$

GANs

1. Generative Algorithm

e.g. Assuming this email is spam, how likely are these features?

⇒ Predict features given a certain label

⇒ Model the distribution of individual classes

(MNIST)

2. One neural network, called the generator, generates new data instances, while the discriminator, evaluates them for authenticity.

The goal of the generator is to generate passable hand-written digits : to lie without being caught. The goal of the discriminator is to identify images coming from the generator as fake

3. For MNIST, the discriminator network is a standard convolutional network that can categorize the images fed to it, a binomial classifier labeling images as real or fake. The generator is an inverse convolutional network.

A standard convolutional classifier takes an image and downsamples it to produce a probability, the generator takes a vector of random noise and upsamples it to an image.

4. Autoencoders encode input data as vectors. They create a hidden, or compressed, representation of the raw data. They are useful in dimensionality reduction.

5. Steps : ① Given a label, they predict the associated features (Naive Bayes)
② Given a hidden representation, they predict the associated features (VAE, GAN)
③ Given some features, they predict the rest (inpainting, imputation)

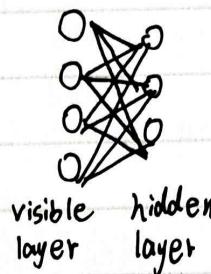
6. Tips : 1). When you train the discriminator, hold the generator values constant
When you train the generator, hold the discriminator values constant
2). The two NNs must have a similar "skill level"

RBM_s (Restricted Boltzmann Machines)

Geoffrey Hinton

1. Definition: ~ is an algorithm useful for dimensionality reduction, classification, regression, collaborative filtering, feature learning and topic modeling

2. Structure: RBMs are shallow, two-layer NNs



Entropy: a common way to measure impurity

Entropy $H(x)$ of a random variable X :

$$H(x) = - \sum_{i=1}^n P(X=i) \log_2 P(X=i)$$

* n : # of possible values for X

- $H(x)$ is the expected number of bits needed to encode a randomly drawn value of X
- Most efficient code assigns $-\log_2 P(X=i)$ bits to encode the message $X=i$

Information Gain:

- 1). We want to determine which attribute in a given set of training feature vectors is most useful for discriminating between the classes to be learned
- 2). Information gain tells us how important a given attribute of the feature vectors is.
- 3) We will use it to decide the ordering of attributes in the nodes of a decision tree.

$$\text{Gain}(S, A) = I_S(A, Y) = H_S(Y) - H_S(Y|A) \quad \text{ii)}$$

$$\text{the same as: } I(X, Y) = H(X) - H(X|Y) = H(Y) - H(Y|X) \quad \text{i)}$$

- i). Information gain is the mutual information between input attribute A and target variable Y
- ii). Information gain is the expected reduction in entropy of target variable Y for data sample S , due to sorting on variable A .

Using Information Gain to Construct a Decision Tree:

Choose the attribute A with highest information gain for the full training set at the root of the tree

Summary : Info Gain (X, A) = $H(X) - H(X|A)$

$$= H(X) - \sum_{v \in \text{values}(A)} \underbrace{\frac{|\{x \in X \mid x_A = v\}|}{|X|}}_{\text{fraction of instances with value } v \text{ in attribute } A} \times \underbrace{H(\{x \in X \mid x_A = v\})}_{\text{entropy of those instances}}$$

Entropy :

$$H(X) = - \sum_{c \in \text{classes}} \underbrace{\frac{|\{x \in X \mid \text{class}(x) = c\}|}{|X|}}_{\text{fraction of instances of class } C} \log_2 \frac{|\{x \in X \mid \text{class}(x) = c\}|}{|X|}$$

Attributes with Many Values

Use Gain Ratio() instead

$$\text{Gain Ratio}(X, A) = \frac{\text{Info Gain}(X, A)}{\text{Split Information}(X, A)}$$

$$\text{Split Information}(X, A) = - \sum_{v \in \text{values}(A)} \frac{|X_v|}{|X|} \log_2 \frac{|X_v|}{|X|}$$

* Where X_v is a subset of X for which A has value v .

Evaluation

1. Classification Metrics : accuracy = $\frac{\# \text{ correct predictions}}{\# \text{ test instances}}$

$$\text{error} = 1 - \text{accuracy} = \frac{\# \text{ incorrect predictions}}{\# \text{ test instances}}$$

2. Confusion Matrix

(Predicted Class)

$$\text{accuracy} = \frac{TP + TN}{P + N}$$

		Yes	No
(Actual Class)	Yes	TP	FN
	No	FP	TN

$$\text{Precision} = \frac{TP}{TP + FP}$$

Probability that a randomly selected result is relevant

$$\text{recall} = \frac{TP}{TP + FN}$$

Probability that a randomly selected relevant doc is retrieved

Linear Regression

- Cost Function $J(\theta) = \frac{1}{2n} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2$

Fit by solving $\min_{\theta} J(\theta)$

Gradient Descent :

1). Initialize θ

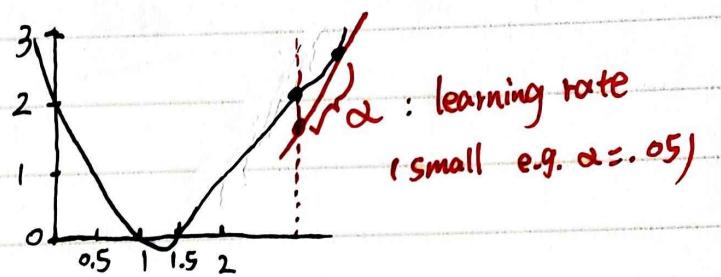
2). Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

Simultaneous update for $j = 0 \dots d$

$$\text{For linear regression : } \frac{\partial}{\partial \theta_j} J(\theta) = \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2$$

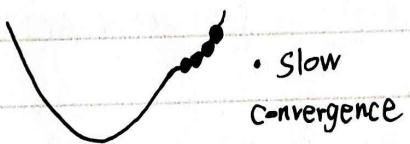
$$= \frac{1}{n} \sum_{i=1}^n \left(\sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right) x_j^{(i)}$$



Hence, repeat until convergence : $\theta_j \leftarrow \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$ simultaneous update for $j = 0 \dots d$

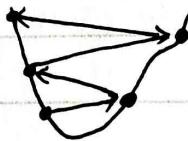
3). Choosing α

α too small



- Slow convergence

α too large



← Increasing value for $J(\theta)$

- May overshoot the minimum
- May fail to converge, even diverge

* To see if gradient descent is working, print out $J(\theta)$ each iteration

The value should decrease each iteration, if it doesn't, adjust α

- Basic Linear Model : $h_\theta(x) = \sum_{j=0}^d \theta_j x_j$

- Generalized Linear Model : $h_\theta(x) = \sum_{j=0}^d \theta_j \phi_j(x)$

basis function allows use of linear regression techniques to fit non-linear datasets

① Polynomial basis functions : $\phi_j(x) = x^j$

② Gaussian basis functions : $\phi_j(x) = \exp\left\{-\frac{(x-\mu_j)^2}{2s^2}\right\}$

③ Sigmoidal : $\phi_j(x) = \sigma\left(\frac{x-\mu_j}{s}\right)$

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

Vectorization :

1. Benefits:
 - ① More compact equations
 - ② Faster code (using optimized matrix libraries)

2. Consider our model for n instances: $h(x^{(i)}) = \sum_{j=0}^d \theta_j x_j^{(i)}$

Let: $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}$ $X = \begin{bmatrix} 1 & x_1^{(1)} & \cdots & x_d^{(1)} \\ 1 & x_1^{(2)} & \cdots & x_d^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \cdots & x_d^{(n)} \end{bmatrix}$

dimensions: $\mathbb{R}^{(d+1) \times 1}$

3. For the linear regression cost function:

$$J(\theta) = \frac{1}{2n} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2 = \frac{1}{2n} \sum_{i=1}^n (\theta^T x^{(i)} - y^{(i)})^2 = \frac{1}{2n} \underbrace{(X\theta - y)^T}_{\mathbb{R}^{1 \times n}} \underbrace{(X\theta - y)}_{\mathbb{R}^{n \times 1}}$$

4. Closed Form Solution

- Instead of using GD, solve for optimal θ analytically
- The solution is when $\frac{\partial}{\partial \theta} J(\theta) = 0$

Derivation:

$$\begin{aligned} J(\theta) &= \cancel{\frac{1}{2n}} (X\theta - y)^T (X\theta - y) \propto ((X\theta)^T - y^T)(X\theta - y) \\ &\stackrel{\text{Constant can be deleted}}{=} (X\theta)^T (X\theta) - \cancel{(X\theta)^T y} - \cancel{y^T (X\theta)} + y^T y = \theta^T X^T X \theta - 2y^T \theta + y^T y \end{aligned}$$

Take derivative and set equal to 0, then solve for θ

$$\begin{aligned} \frac{\partial}{\partial \theta} (\theta^T X^T X \theta - 2y^T \theta + \cancel{y^T y}) &= 0 \quad * \frac{\partial (X^T X)}{\partial \theta} = 2X \\ X^T X \theta - X^T y &= 0 \\ \theta &= (X^T X)^{-1} X^T y \end{aligned}$$

Gradient Descent

- Requires multiple iterations
- Need to choose α
- Works well when n is large
- Can support incremental learning

Closed Form Solution

- Non-iterative
- No need for α
- Slow if n is large
- 只适用于 linear-regression

Why we choose GD?

- ① Can apply GD to any model, can be done automatically by software packages (Theano, TensorFlow)
- ② Much faster than solving the linear system

Improving Learning :

① Feature Scaling : Ensure that feature have similar scales \Rightarrow GD converges much faster

② Feature Standardization : Rescale features to have zero mean and unit variance

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{s_j}$$

μ_j : the mean of feature j $\mu_j = \frac{1}{n} \sum_{i=1}^n x_j^{(i)}$

s_j : the SD of feature j

③ Regularization : A method for automatically controlling the complexity of the learned hypothesis

Idea: Penalize for large values of θ :

- Can incorporate into the cost function

- Works well when we have a lot of features, each that contributes a bit to predicting the label

Linear regression objective function

$$J(\theta) = \underbrace{\frac{1}{2n} \sum_{i=1}^n (\hat{y}(x^{(i)}) - y^{(i)})^2}_{\text{model fit to data}} + \underbrace{\frac{\lambda}{2} \sum_{j=1}^d \theta_j^2}_{\text{regularization}}$$

1). λ is the regularization parameter ($\lambda > 0$)

2). No regularization on θ_0

Understanding: $\sum_{j=1}^d \theta_j^2 = \|\theta_{1:d}\|_2^2$ \Rightarrow the magnitude of the feature coefficient vector

$$\sum_{j=1}^d (\theta_j - \bar{\theta})^2 = \|\theta_{1:d} - \bar{\theta}\|_2^2 \quad \Rightarrow \text{L}_2 \text{ regularization pulls coefficient toward } 0$$

Gradient update: $\frac{\partial}{\partial \theta_j} J(\theta) \quad \theta_j \leftarrow \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n (\hat{y}(x^{(i)}) - y^{(i)}) x_j^{(i)} - \alpha \lambda \theta_j$
regularization

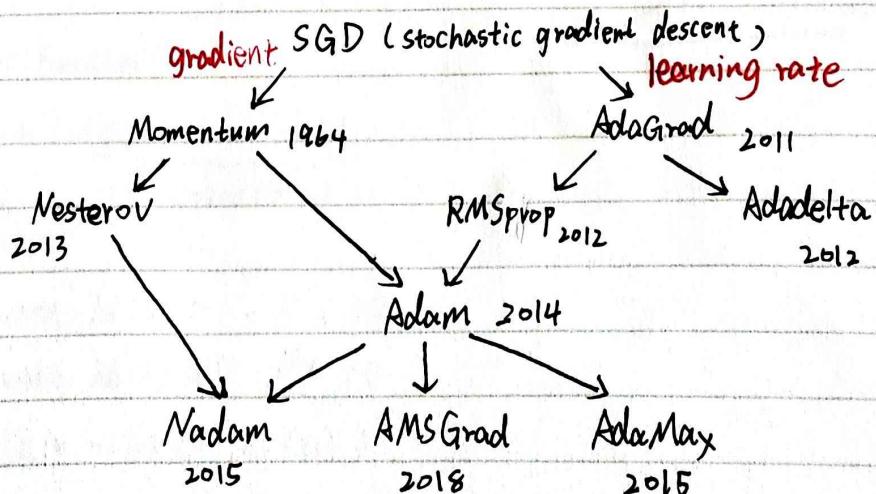
Gradient Descent Optimization

3 main ways: 1). modifying the learning rate α

2). modifying the gradient component $\frac{\partial L}{\partial w}$

3). both

Evolutionary Map:



Support Vector Machines & Kernels

Doing really well with linear decision surfaces

1. Why might predictions be wrong?

- i) True non-determinism
- ii) Partial observability
- iii) Noise in the observation \times (Measurement error, instrument limitations)
- iv) Representational bias
- v) Algorithmic bias

2. Strengths of SVMs:

- i) Good generalization
- ii) Works well with few training instances
- iii) Find globally best model
- iv) Efficient algorithms
- v) Amenable to the kernel trick

3. Why Maximize Margin?

i). Increasing margin reduces capacity

ii). If the following holds:

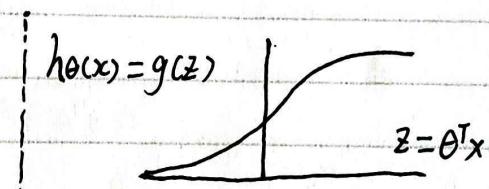
- H is sufficiently constrained in size
- and/or the size of the training data set n is large, then low training error is likely to be evidence of low generalization error
then low training error is likely to be evidence of low generalization error

4. Alternative view of Logistic Regression

$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

The Logistic Regression Objective Function

$$J(\theta) = - \sum_{i=1}^n \left[\underbrace{y_i \log h_\theta(x_i)}_{\text{cost}_1(\theta^T x_i)} + \underbrace{(1-y_i) \log(1-h_\theta(x_i))}_{\text{cost}_0(\theta^T x_i)} \right]$$



Intuition: If $y=1$, we want $h_\theta(x) \approx 1$. $\theta^T x \geq 0$

If $y=0$, we want $h_\theta(x) \approx 0$, $\theta^T x \leq 0$

Cost: $-y_i \log h_\theta(x_i) - (1-y_i) \log(1-h_\theta(x_i))$

$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

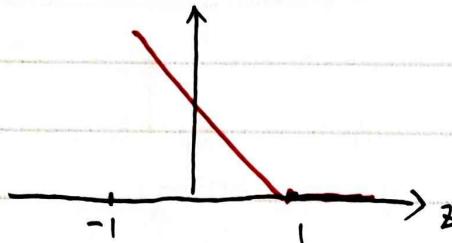
Logistic Regression to SVMs:

$$\text{Logistic Regression : } \min_{\theta} - \sum_{i=1}^n [y_i \log h_{\theta}(x_i) + (1-y_i) \log(1-h_{\theta}(x_i))] + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2$$

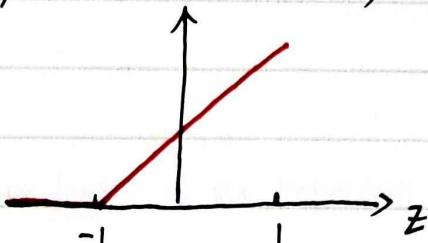
$$\text{SVMs : } \min_{\theta} C \sum_{i=1}^n [y_i \text{cost}_1(\theta^T x_i) + (1-y_i) \text{cost}_0(\theta^T x_i)] + \frac{1}{2} \sum_{j=1}^d \theta_j^2$$

* you can think C as similar to $\frac{1}{\lambda}$

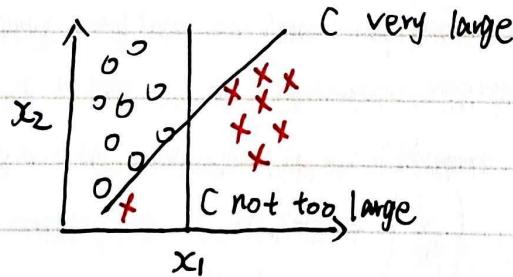
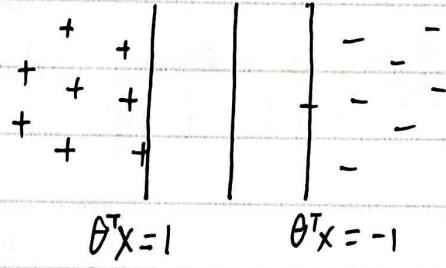
If $y=1$ (want $\theta^T x \geq 1$):



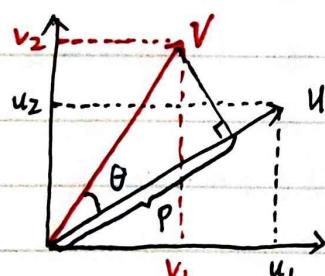
If $y=0$ (want $\theta^T x \leq -1$)



$$\text{margin} = \frac{2}{\|\theta\|_2}$$



Vector Inner Product



$$u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$\|u\|_2 = \text{length}(u) \in \mathbb{R} = \sqrt{u_1^2 + u_2^2}$$

$$u^T v = v^T u = u_1 v_1 + u_2 v_2 = \|u\|_2 \|v\|_2 \cos \theta = P \|u\|_2$$

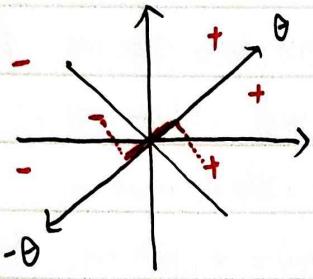
$$\text{where } P = \|v\|_2 \cos \theta$$

Understanding the Hyperplane

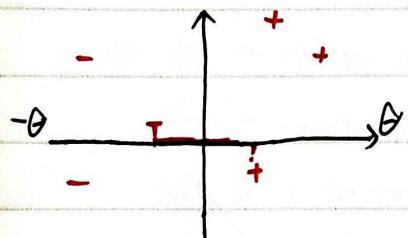
$$\min_{\theta} \frac{1}{2} \sum_{j=1}^d \theta_j^2 \quad \text{s.t. } \theta^T x_i \geq 1 \text{ if } y_i=1 \\ \theta^T x_i \leq -1 \text{ if } y_i=-1$$

Assume $\theta_0=0$ so that the hyperplane is centered at the origin, and that $d=2$

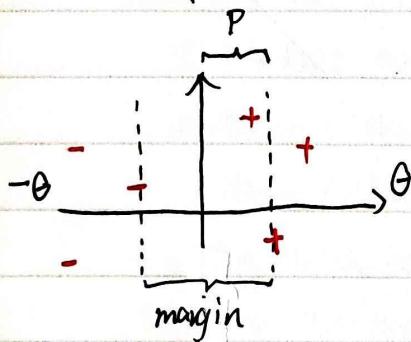
$$\theta^T x = \|\theta\|_2 \underbrace{\|x\|_2 \cos \theta}_P = P \|\theta\|_2$$



Since p is small, therefore $\|\theta\|_2$ must be large
to have $P\|\theta\|_2 \geq 1$ (or ≤ -1)



since p is larger, $\|\theta\|_2$ can be smaller in
order to have $P\|\theta\|_2 \geq 1$ (or ≤ -1)



$P\|\theta\|_2 = \pm 1$. P is the length of the projection of the SVs onto
Therefore, $P = \frac{1}{\|\theta\|_2}$

$$\text{margin} = 2P = \frac{2}{\|\theta\|_2}$$

Feature mappings : an easy way to use linear regression to learn nonlinear dependencies.

$$\text{eq. } y = w_3x^3 + w_2x^2 + w_1x + w_0 \Rightarrow \text{Polynomial regression}$$

* Algorithmically, polynomial regression is no different from linear regression

we define a feature mapping ϕ . $\phi(x) = \begin{pmatrix} x \\ x^2 \\ x^3 \end{pmatrix}$, compute the predictions as

$$y = w^T \phi(x) \text{ instead of } w^T x$$

Limitations: 1). The features need to be known in advance

2). In high dimensions, the feature representations can get very large

Solution: Using Neural Network to learn nonlinear predictors directly from the raw inputs.

Generalization error

* Hyperparameter : are values that we can't include in the training procedure itself, but which we need to set using other means. In practice, we tune hyperparameters by partitioning the dataset into three different subsets:

1). training set. to train the model

2). validation set. to eliminate the generalization error of each hyperparameter setting

3). test set. - - - - - of the final model

PCA : principal component analysis

- unsupervised learning algorithm
- A linear model, with a Closed-form Solution.

Def.: Choosing a subspace to maximize the projected variance, or minimize the reconstruction error

1). Learning a subspace

Q: How to choose a good subspace S ?

Need to choose a vector μ and a $D \times K$ matrix U with orthonormal columns

Set μ to the mean of the data, $\mu = \frac{1}{N} \sum_{i=1}^N x^{(i)}$

Two criteria:

① Minimize the reconstruction error $\min \frac{1}{N} \sum_{i=1}^N \|x^{(i)} - \tilde{x}^{(i)}\|^2$

② Maximize the variance of the code vectors

$$\max_j \text{Var}(z_j) = \frac{1}{N} \sum_j \sum_i (z_j^{(i)} - \bar{z}_j)^2 = \frac{1}{N} \sum_j \|z^{(i)} - \bar{z}\|^2 = \frac{1}{N} \sum_j \|z^{(i)}\|^2$$

* \bar{z} denotes the mean

These two criteria are equivalent

by unitarity: $\|\tilde{x}^{(i)} - \mu\| = \|Uz^{(i)}\| = \|z^{(i)}\|$

by the Pythagorean Theorem: $\underbrace{\frac{1}{N} \sum_{i=1}^N \|\tilde{x}^{(i)} - \mu\|^2}_{\text{Projected Variance}} + \underbrace{\frac{1}{N} \sum_{i=1}^N \|x^{(i)} - \tilde{x}^{(i)}\|^2}_{\text{Reconstruction Error}} = \underbrace{\frac{1}{N} \sum_{i=1}^N \|x^{(i)} - \mu\|^2}_{\text{Constant}}$

Stochastic Gradient Descent

mini-batch : size S : a hyperparameter that needs to be set
a reasonable value might be $S = 100$

In stochastic training, the learning rate also influences the fluctuations due to the stochasticity of the gradients.

Typical Strategy:

- 1). Use a large learning rate early in training so you can get close to the optimum
- 2). Gradually decay the learning rate to reduce the fluctuations
逐漸地

* By reducing the learning rate, you reduce the fluctuations, which can appear to make the loss drop suddenly

Generalization

Decomposition of generalization error, or expected loss into bias, variance, and Bayes error.

$$E[(y - t)^2] = \underbrace{(y^* - E[y])^2}_{\text{bias}} + \underbrace{\text{Var}(y)}_{\text{Variance}} + \underbrace{\text{Var}(t)}_{\text{Bayes Error}}$$

An overly simple model might have high bias and low variance

An overly complex model might have low bias and high variance

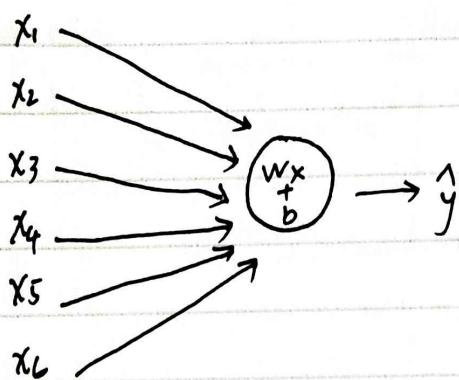
Q: How to achieve low bias and low variance?

- Solu:
- ① Data augmentation
 - ② Reduce the number of parameters
 - ③ Weight decay
 - ④ Early Stopping
 - ⑤ Ensembles
 - ⑥ Stochastic regularization (e.g. dropout)
 - ⑦ Reduce the # of layers or the # of parameters per layer
 - ⑧ Adding a linear bottleneck layer

Q: Why are neural networks vulnerable to adversarial examples?

Ex: Linear Regression:

The forward propagation: $\hat{y} = Wx + b$



Network has been trained and converged to:

$$W = (1, 3, -1, 2, 2, 3), b = 0$$

$$\star x = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}, w = \begin{bmatrix} \cdots \end{bmatrix}$$

$$\text{i.e. } x = \begin{bmatrix} 1 \\ -1 \\ 2 \\ 0 \\ 3 \\ -2 \end{bmatrix} \quad \hat{y} = Wx + b = 1 - 3 - 2 + 0 + 6 - 6 = [-4]$$

(Q): How to change $x \rightarrow x^*$ such that \hat{y} changes radically but $x^* \approx x$? a prob of adversarial example

The intuition comes from the derivative

$$\frac{\partial \hat{y}}{\partial x} : \text{the impact on } \hat{y} \text{ of small changes of } x \quad \frac{\partial \hat{y}}{\partial x} = W^T$$

$$x^* = x + \epsilon w^T \quad \begin{array}{l} \text{小变动} \\ \downarrow \text{value of perturbation} \end{array} \quad \hat{y}^* = Wx^* + b = Wx + \epsilon w^T \\ (= 0) \quad = Wx + \epsilon w^T \quad \text{always } > 0$$

Ex:

$$x^* = x + \epsilon w^T = \begin{bmatrix} 1 + 0.2 \\ -1 + 0.6 \\ 2 - 0.2 \\ 0 + 0.4 \\ 3 + 0.4 \\ -2 + 0.6 \end{bmatrix} = \begin{bmatrix} 1.2 \\ -0.4 \\ 1.8 \\ 0.4 \\ 3.4 \\ -1.4 \end{bmatrix}$$

* A very slight change in x^* will

$$\hat{y}^* = Wx^* = 1.2 - 1.2 - 1.8 + 0.8 + 6.8 - 4.2 = [1.6] \quad \text{has pushed } \hat{y} \text{ from } -4 \text{ to } 1.6$$

Insights: ①: If w is large $\Rightarrow x^* \neq x$

$\hookrightarrow \text{sign}(w)$ instead of w^*

②: As x grows in dimensions, the impact of $+\epsilon \text{ sign}(w)$ on \hat{y} increases

Images have high dimensions: $64 \times 64 \times 3$

Conclusion: The fast way to generate adversarial example: (fast gradient sign method)

$$x^* = x + \epsilon \text{ sign}(\nabla_x \hat{y}(W, x, y))$$

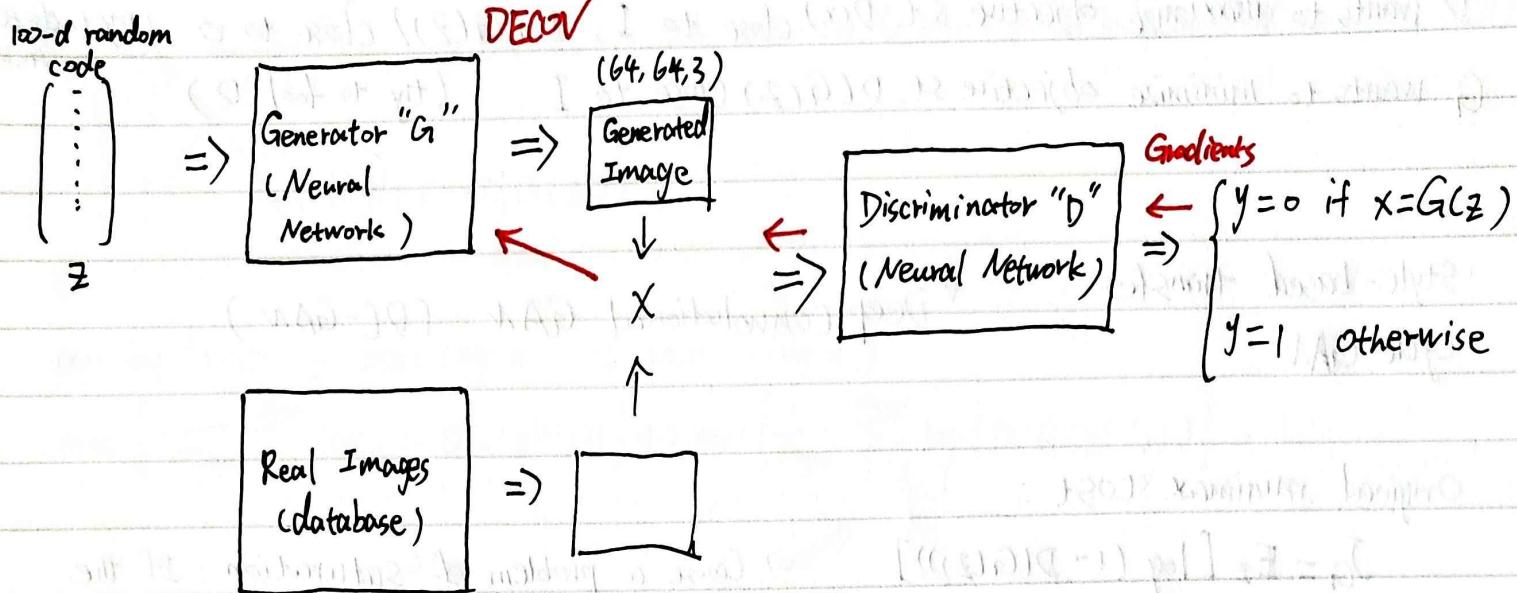
We are linearizing the cost function in the proximity of the parameters

We are pushing the pixel images in one direction that is going to impact highly the output

GANs

Goal: Collect a lot of data, use it to train a model to generate similar data from scratch

Intuition: # of parameters of the model << # of data



- The way we will train the "G" is that we will backpropagate the gradient to the "D" to train the "D", using a binary cross entropy. If the image is real, then the generator couldn't get the gradient, because x doesn't depend on z or on the features and parameters of the generator.
- Run Adam simultaneously on two minibatches (true data / generated data)

Training procedure, we want to minimize : Labels : $\begin{cases} Y_{\text{real}} \text{ is always 1} \\ Y_{\text{gen}} \text{ is always 0} \end{cases}$

i): The cost of the "D"

$$J(D) = - \underbrace{\frac{1}{m_{\text{real}}} \sum_{i=1}^{m_{\text{real}}} Y_{\text{real}}^{(i)} \cdot \log(D(x^{(i)}))}_{\text{Cross-entropy 1}} - \underbrace{\frac{1}{m_{\text{gen}}} \sum_{i=1}^{m_{\text{gen}}} (1 - Y_{\text{gen}}^{(i)}) \cdot \log(1 - D(G(z^{(i)})))}_{\text{Cross-entropy 2}}$$

"D" should correctly label real data as 1
--- generated data as 0

We both want the D to correctly identify real data and also correctly identify fake data.

ii): The cost of the "G"

$$J(G) = -J(D) = \frac{1}{m_{\text{gen}}} \sum_{i=1}^{m_{\text{gen}}} \log(1 - D(G(z^{(i)})))$$

G should try to fool D:
by minimizing the opposite of what D is trying to minimize

We want D to go up to be very good and G to go up at the same time until the equilibrium: a certain point where D will always output one-half like, random probabilities.

Train GAN jointly via minimax game:

$$\min_{\theta_g} \max_{\theta_d} [E_{x \sim P_{\text{data}}} \log D_{\theta_d}(x) + E_{z \sim P(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z)))]$$

D wants to maximize objective s.t. $D(x)$ close to 1, $D(G(z))$ close to 0 (try to detect fake)
G wants to minimize objective st. $D(G(z))$ close to 1 (try to fool D)

Style-based transfer
cycle GAN

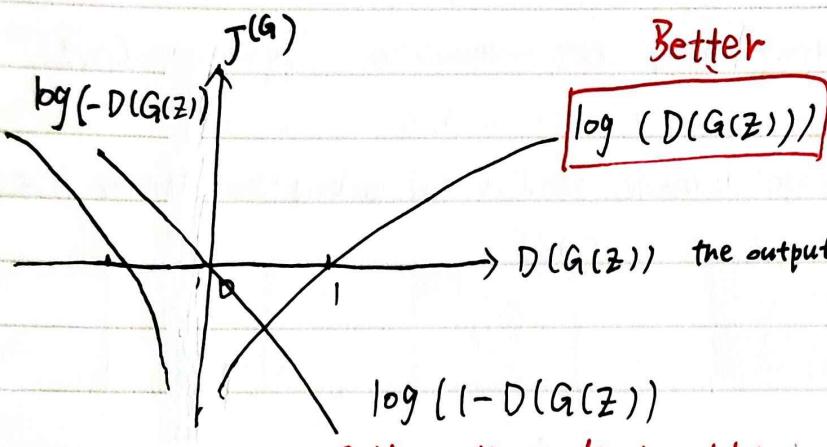
Deep convolutional GAN (DC-GAN)

Original minimax cost:

$$J_G = E_z [\log (1 - D(G(z)))] \Rightarrow \text{Cause a problem of saturation; If the generated sample is really bad, the discriminator's prediction is close to 0, and the generator's cost is flat.}$$

Modified generator cost:

$$J_G = E_z [-\log D(G(z))]$$



$$\min \log(1-x) = \max \log x = \min(-\log x)$$

$$\min \left[\frac{1}{m_{\text{gen}}} \sum_{i=1}^{m_{\text{gen}}} \log(1-D(G(z^{(i)}))) \right] \Leftrightarrow \max \left[\frac{1}{m_{\text{gen}}} \sum_{i=1}^{m_{\text{gen}}} \log(D(G(z^{(i)}))) \right]$$

$$\min \left[-\frac{1}{m_{\text{gen}}} \sum_{i=1}^{m_{\text{gen}}} \log(D(G(z^{(i)}))) \right]$$

New training procedure, we want to minimize:

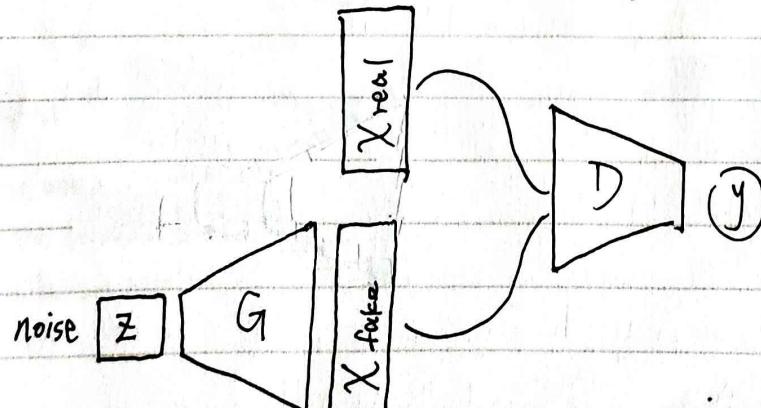
$$J(D) = \underbrace{-\frac{1}{m_{\text{real}}} \sum_{i=1}^{m_{\text{real}}} y_{\text{real}}^{(i)} \cdot \log(D(x^{(i)}))}_{\text{Cross-entropy 1}} - \underbrace{\frac{1}{m_{\text{gen}}} \sum_{i=1}^{m_{\text{gen}}} (1-y_{\text{gen}}^{(i)}) \cdot \log(1-D(G(z^{(i)})))}_{\text{Cross-entropy 2}}$$

D should correctly label real data as 1 D ... generated data as 0

$$J(G) = -\frac{1}{m_{\text{gen}}} \sum_{i=1}^{m_{\text{gen}}} \log(D(G(z^{(i)}))) \Rightarrow G \text{ should try to fool } D \text{ by minimizing this}$$

GANs' training tips:

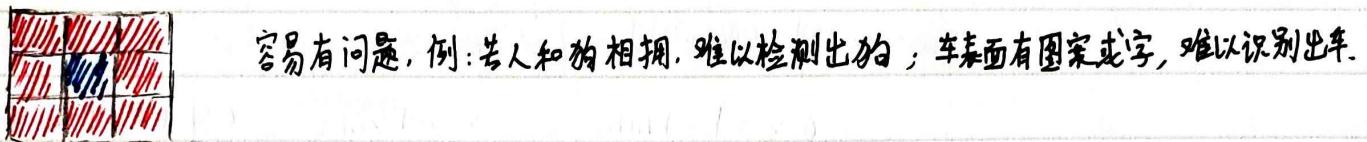
- 1). Modification of the cost function
 - 2). Keep D up-to-date with respect to G
(k update for D / 1 update for G)
 - 3). Use virtual Batchnorm
 - 4). One-sided label smoothing
- for num_iterations =
for k iterations =
update D
update G



The generator turns noise into an imitation of the data to try to trick the discriminator

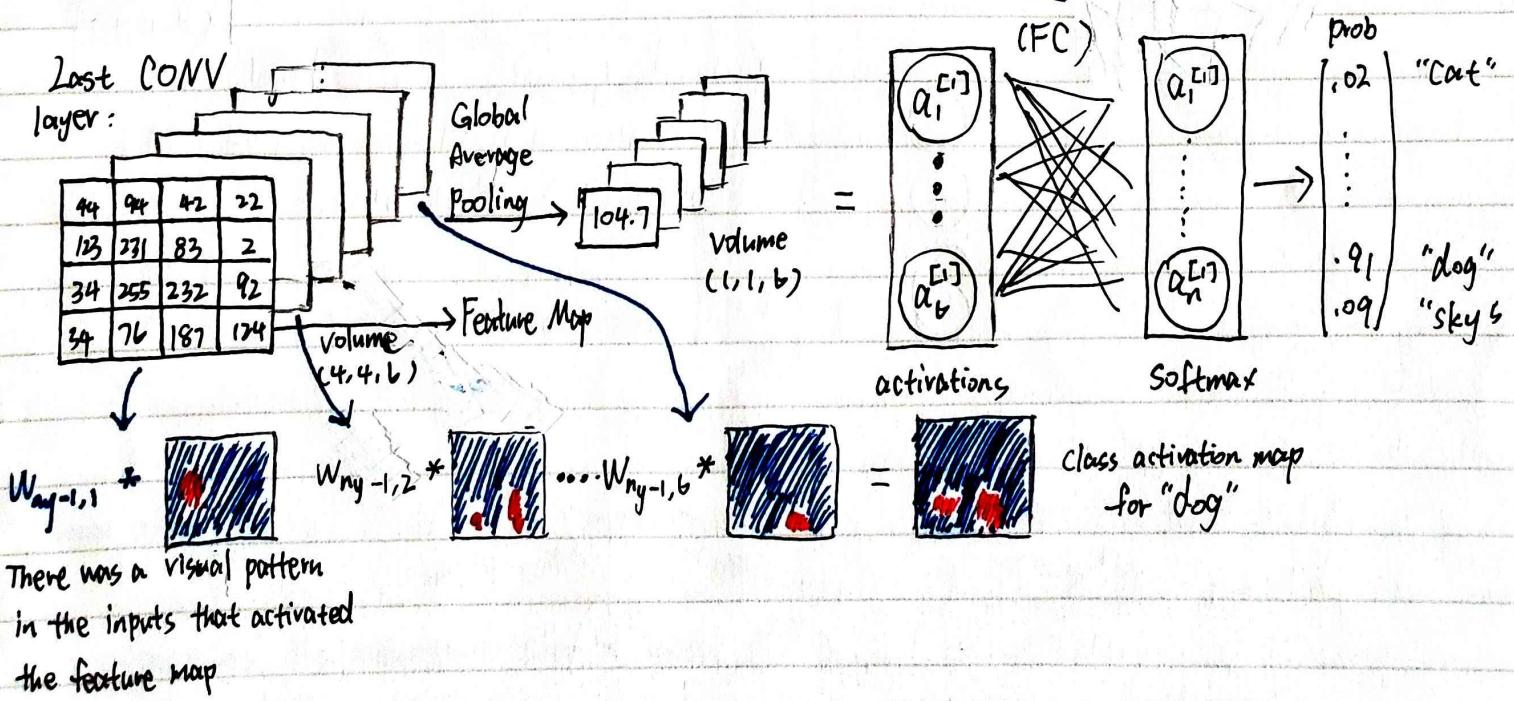
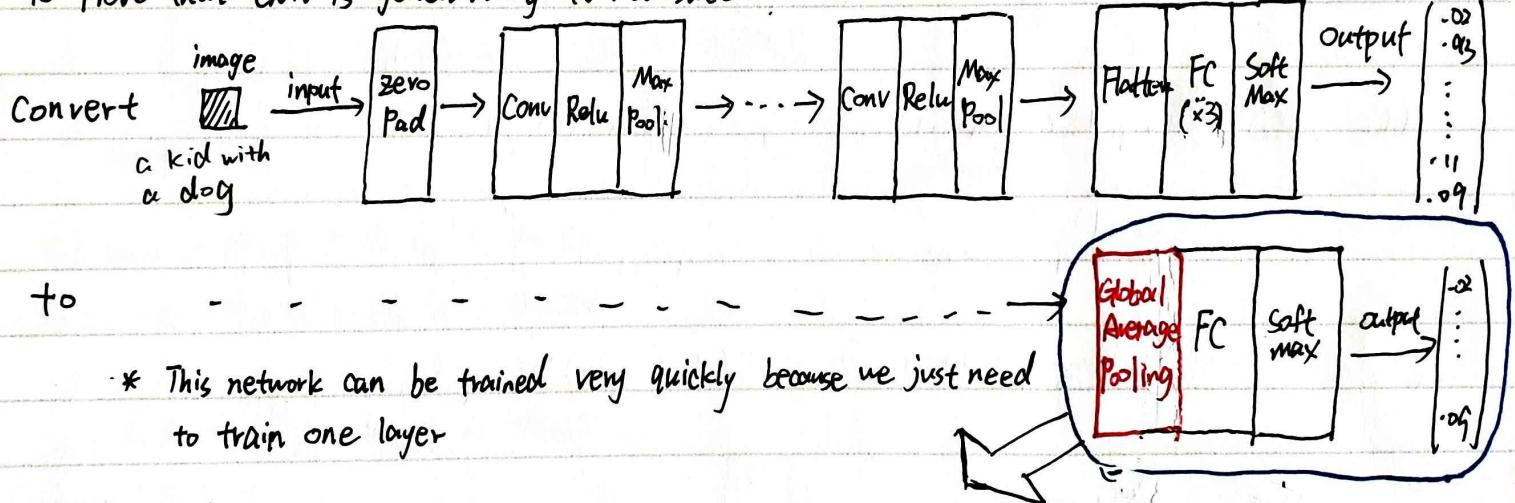
The D tries to identify real data from fakes created by the generator

1. Saliency maps: a common technique to quickly visualize what the network is looking at.
2. A method more precise but a little slower: occlusion sensitivity



3. Class activation map shows that CNNs have a very good localization ability even if they were trained only on image level labels.

To Prove that CNN is generalizing to localization:



For non-image like speech detection, attention model.

Visualize NNs from the inside:

A. With gradient ascent (class model visualization)

B. With dataset search

C. The deconvolution and its applications \star

Deconvolution: will up-sample the height and width of an image

See Conv first: (increase)

(Goal:) Frame the convolution as a simple matrix vector mathematical operation

① 1D Conv:

$$x = \begin{bmatrix} 0 \\ 0 \\ x_1 \\ x_2 \\ \vdots \\ x_8 \\ 0 \\ 0 \end{bmatrix} \xrightarrow{\text{Input}} \begin{array}{|c|} \hline \text{The input is a one-dimensional vector which has two paddings on both sides.} \\ \hline \end{array} \xrightarrow{\text{1 filter, whose filter size is 4}} \begin{array}{|c|} \hline \text{Stride = 2 } (W_1, W_2, W_3, W_4) \\ \hline \end{array} \xrightarrow{\text{* output}} \begin{bmatrix} y_1 \\ \vdots \\ y_5 \end{bmatrix} = \frac{n_x - f + 2P}{s} + 1 \\ = \frac{12 - 8 + 4}{2} + 1$$

$$y_1 = w_1 \cdot 0 + w_2 \cdot 0 + w_3 \cdot x_1 + w_4 \cdot x_2$$

$$y_2 = w_1 \cdot x_1 + w_2 \cdot x_2 + w_3 \cdot x_3 + w_4 \cdot x_4$$

\vdots

$$y_5 = w_1 \cdot x_7 + w_2 \cdot x_8 + w_3 \cdot 0 + w_4 \cdot 0$$

Conv: $y = wx$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \end{bmatrix} = \begin{bmatrix} w_1 & w_2 & w_3 & w_4 & 0 & \cdots & 0 \\ 0 & 0 & w_1 & w_2 & w_3 & w_4 & 0 \cdots 0 \\ & & \ddots & & & & \\ & & & \ddots & & & \\ 0 & \cdots & 0 & w_1 & w_2 & w_3 & w_4 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ x_1 \\ x_2 \\ \vdots \\ x_8 \\ 0 \\ 0 \end{bmatrix}$$

For deconvolution: $x = \underbrace{w^{-1}y}_{(12, 5)}$

Some assumptions: ① W is an orthogonal matrix ex: filter = $(1, 0, 0, -1)$

$$\therefore w^{-1} = w^T \quad w^T w = I \Rightarrow x = \underline{w^T y}$$

* often times, the assumption doesn't hold.

If you need the weights, you need the assumption. otherwise you don't need.

(B) 1D deconv

$$\begin{bmatrix} 0 \\ 0 \\ x_1 \\ x_2 \\ \vdots \\ x_8 \\ 0 \\ 0 \end{bmatrix} \text{Crop } \left\{ \begin{bmatrix} w_1 & 0 & 0 \\ w_2 & 0 & 0 \\ w_3 & w_1 & 0 \\ w_4 & w_2 & 0 \\ 0 & w_3 & w_1 \\ w_4 & w_2 & 0 \\ 0 & w_3 & w_1 \\ w_4 & 0 & 0 \end{bmatrix} \right\} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \end{bmatrix}$$

W^T

They are totally the same equations

Explain : ① In order to be able to efficiently

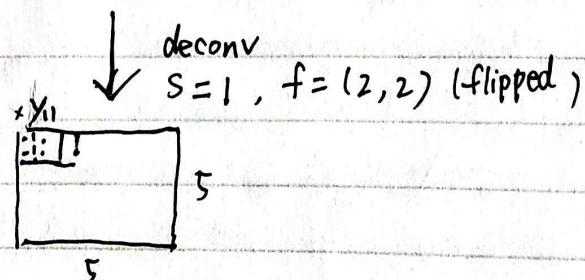
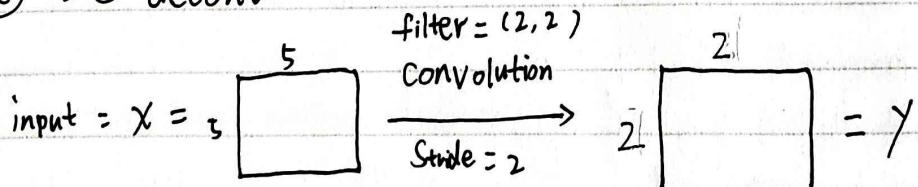
compute the deconvolution the same way as we've learned to compute the convolution, we wanted to have the weights scattered from left to right with a stride moving from left to right. So we used a sub-pixel version of Y by inserting zeros in the middle. And we divided the stride by two. So instead of having a stride of two as we had in our convolution, we have a stride of one in our deconvolution.

② I flipped my weights. Instead of w_1, w_2, w_3, w_4 , now I have w_4, w_3, w_2, w_1 .

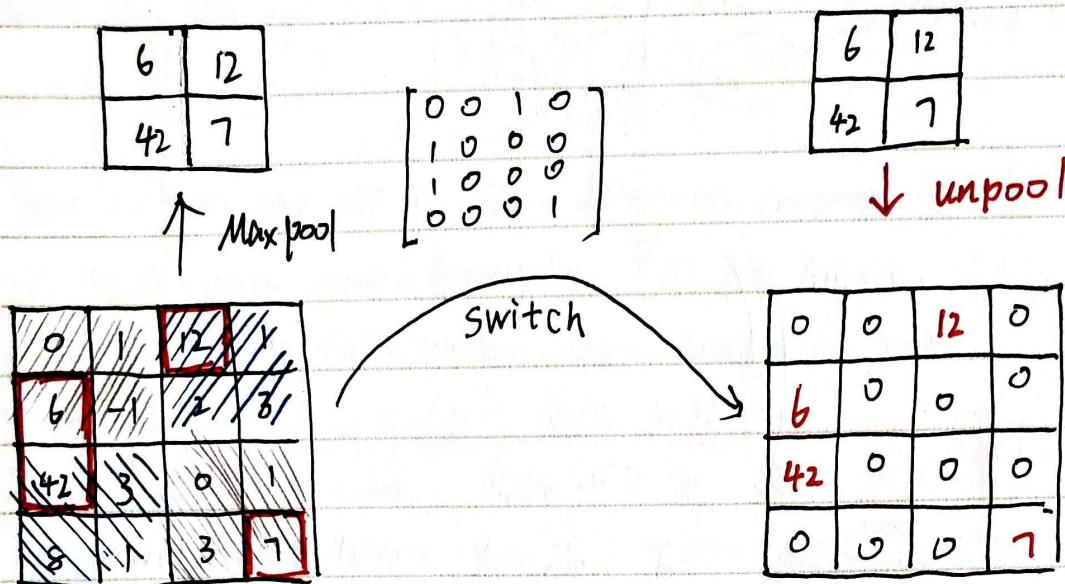
Summary :

To deconvolve, just flip all the weights, insert zeros, sub-pixel, and divide the stride

(C) 2D deconv



Unpool :



Q Learning

In Reinforcement Learning we want to obtain a function $Q(s, a)$ that predicts best action a in state s in order to maximize a cumulative reward.

This function can be estimated using Q-Learning, which iteratively updates $Q(s, a)$ using the Bellman Equation

$$Q(s, a) = \underbrace{r}_{\text{immediate reward}} + \gamma \max_{a'} \underbrace{Q(s', a')}_\text{future reward}$$

Initialize Q-table

- Q-table : matrix of entries representing "how good is it to take action a in state s "

Choose action from Q

- Policy : function telling us what's the best strategy to adopt

Perform action

- Bellman equation satisfied by the optimal Q-table

Measure Reward

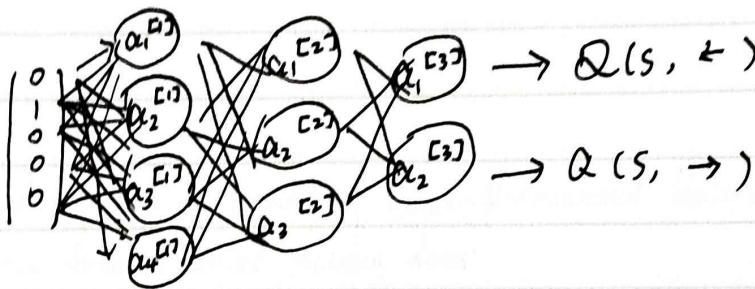
Deep Q-Learning (DQN) :

Update Q

Find a Q-function to replace the Q-table
→ Neural Network

DQN

Replace Q-table by : $s = \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}$



How to train the NN? Some differences compared with the classic supervised learning :

- ① The Q-scores change dynamically
- ② No labels.

- The DQN is a regression problem, not a classification prob.

① Loss function : $L = (y - Q(s, \leftarrow))^2$

② Target value. Assume $Q(s, \leftarrow) > Q(s, \rightarrow)$ Case 1

Then define the target as $y = r_{\leftarrow} + \gamma \max(Q(s^{\text{next}}, a'))$

Immediate reward for taking action \leftarrow in state s

↓ Discounted maximum future reward when you are in state s^{next}

* Considering $Q(s^{\text{next}}, a')$ fixed, until we get close to there, and our gradient is small, then we will update it and we'll fix it.

- We have two networks in parallel, one that is fixed and another is not.

Case 2

$$Q(s, \leftarrow) < Q(s, \rightarrow)$$

$$y = r_{\rightarrow} + \gamma \max(Q(s^{\text{next}}, a'))$$

③ Backpropagation

Compute $\frac{\partial^2}{\partial w^2}$ and update W using stochastic gradient descent

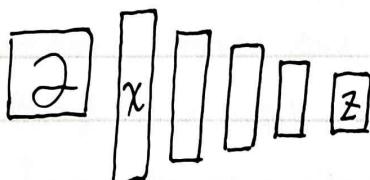
Summary of DQN implementation:

- Initialize your Q-network parameters
- Loop over episodes :
 - Start from initial state s .
 - Loop over time-steps :
 - Forward propagate s in the Q-network
 - Execute action a (that has the maximum $Q(s, a)$ output of Q-network)
 - Observe reward r and next state s' (immediate reward)
 - Compute targets y by forward propagating state s' in the Q-network, then compute loss.
 - Update parameters with gradient descent.

Autoencoders :

1. Background :

- ① Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data.
- ② "Encoder" learns mapping from the data, x , to a low-dimensional latent space, z .

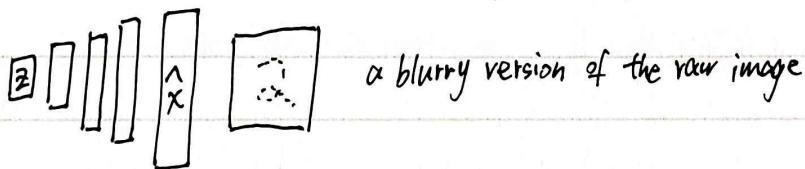


* Low-dimensional z allows us to find the most rich features in the image

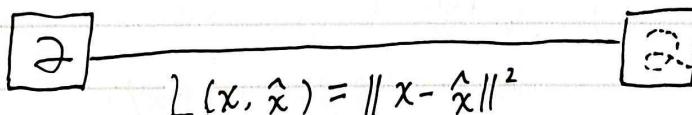
- ③ How can we learn this latent space?

Train the model to use these features to reconstruct the original data.

- ④ "Decoder" learns mapping back from latent z , to a reconstructed observation \hat{x} .



* A blurry version of the raw image

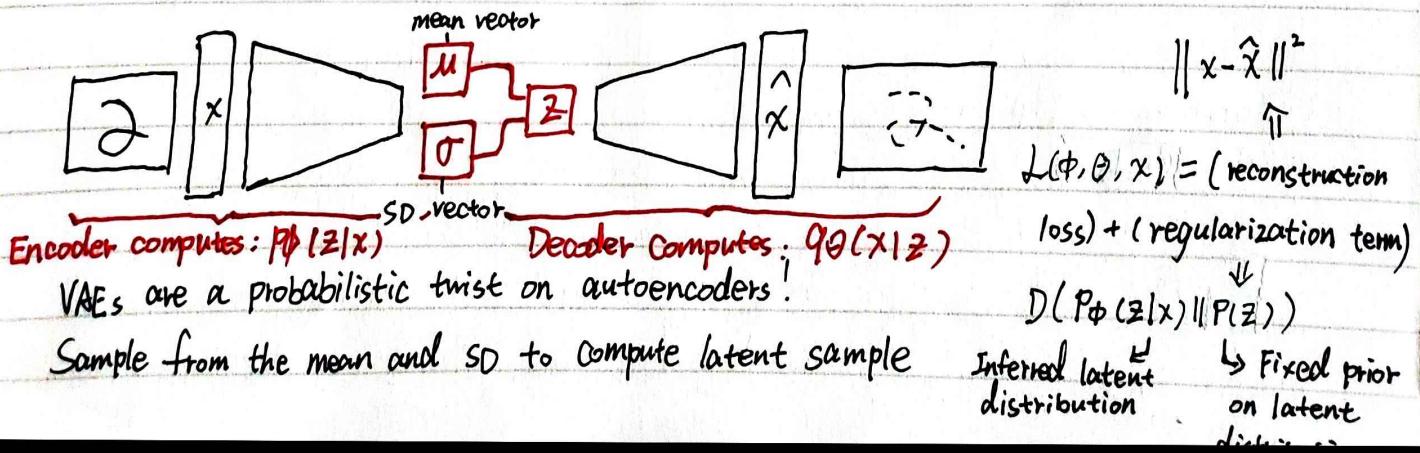


* Loss function doesn't use any labels

2. Summary :

- 1). Bottleneck hidden layer forces network to learn a compressed latent representation
- 2). Reconstruction loss forces the latent representation to capture (or encode) as much "information" about the data as possible
- 3). Autoencoding = Automatically encoding data

Variational Autoencoders (VAEs)



$$D(P_\phi(z|x) \| P(z)) = -\frac{1}{2} \sum_{j=0}^{k-1} (\sigma_j + \mu_j^2 - 1 - \log \sigma_i)$$

Problem: We cannot backpropagate gradients through sampling layers!

Solu: Reparametrizing the sampling layer

Key idea: $z \sim N(\mu, \sigma^2)$

\therefore We can't just simply sample from this distribution in a backward pass because we can't compute that backwards pass and we can't compute a gradient through it. So Let's instead consider a different alternative :

Consider the sample latent vector as a sum of :

i). a fixed μ vector

ii). a fixed σ vector, scaled by random constants drawn from the prior distribution $\Rightarrow z = \mu + \sigma \cdot \epsilon$ * where $\epsilon \sim N(0, 1)$

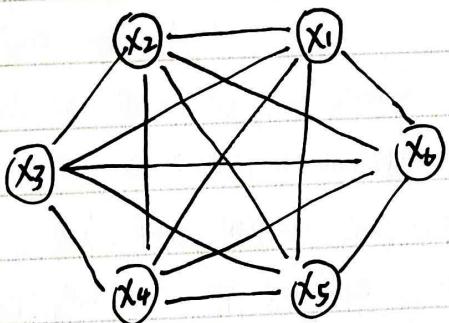
LSTM : Long - Term Short Term Memory

- The idea is that a network's activations are its short-term memory and its weights are its long-term memory
 - This is a kind of special RNN
 - 有三种门来保持及控制细胞状态.:
 - 决定从细胞状态中舍弃哪些信息: 遗忘门 forget gate
 - 决定在细胞状态中保存哪些信息: 输入门 input gate
 - 输出门 output gate
- } sigmoid

DBN : Deep Belief Network

通过采用逐层训练的方式,为整个网络赋予了较好的初始权值,解决了深层次网络的优化问题. 其中起重要作用的是 RBM (Restricted Boltzmann Machine).

BM:



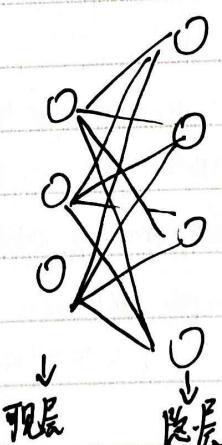
$x_1 \ x_2 \ x_3$: 可见层

$x_4 \ x_5 \ x_6$: 隐层

全连接无向图

训练时间特别长

RBM :



显元用于接受输入

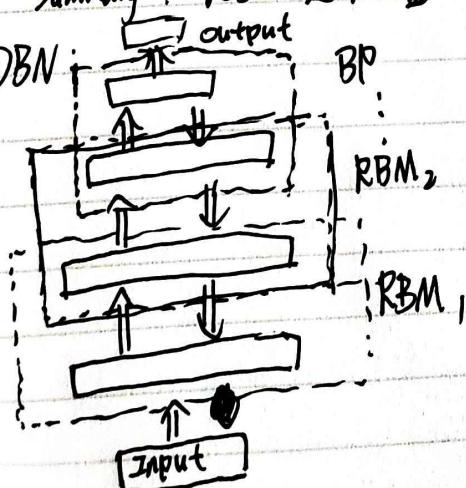
隐元用于提取特征

RBM 可以通过学习将数据表示成概率模型, 还可以生成新数据

层与层全连接, 层内无连接的二分图

\Rightarrow 计算量大大减少

DBN:



Summary: RBM 是非监督学习的利器, 可以用于: 降维, 学习提取特征, 自编码器, 以及 DBN

1). 最底部 RBM 以原始输入数据进行训练

2). 将底部 RBM 抽取的特征作为顶部 RBM 的输入继续训练.

3). 重复这个过程训练尽可能多的 RBM 层

AlexNet 模型的特点：

- 1). 激活函数使用了 ReLU
- 2). 防止过拟合，使用了 Dropout 和 Data augmentation
- 3). 多 GPU 实现，LRN 归一化层的使用
- 其他：4). Overlapping Pooling 重叠池化

LRN: Local Response Normalization 局部归一化

Four modern approaches to generative modeling:

- ① GANs
- ② Reversible architectures
- ③ Autoregressive models (Neural language models, RNN language models)
- ④ Variational autoencoders

CNN vs. RNN

- 1). The RNN has a memory, so it can use information from all past time steps.

The CNN has a limited context. That is - the advantage of RNNs over CNNs is that their memory lets them learn arbitrary long-distance dependencies.

- 2). But training the RNN is very expensive since it requires a for loop over time steps.

The CNN only requires a series of convolutions

一些 ideas:

1). CNN 中，是 set stride 效果好还是 pooling 效果好？

i). stride 要小于卷积核的 size。当 size of the kernel = 3，stride = 2 时，输出矩阵 size 降低为输入矩阵 size 的一半，可以利用这一特性代替池化层。

在 ResNet 中，每隔几层就 set the stride = 2 来压缩 size，降维的同时减小计算量

2). 对于卷积的思考：

i) Dilated convolution 空洞卷积，让卷积核遍到更大的范围

Deformable convolution 可变形卷积核

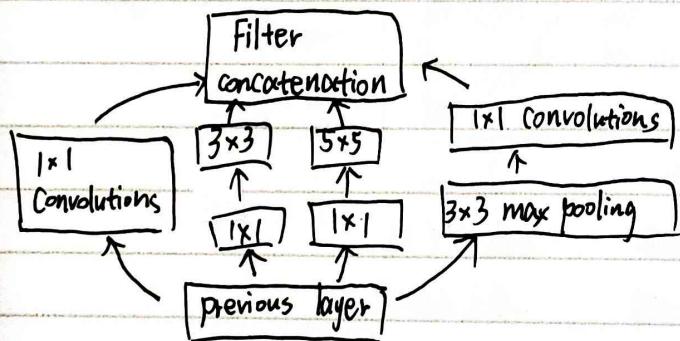
3). 网络的变迁

① AlexNet：创新点：① ReLU；② Dropout 动态卷积过程？

② ZFNet：通过反卷积进行卷积网络可视化方法

③ GoogLeNet：利用 Inception 机制，在增大网络规模（深度+宽度）的同时防止过拟合及大量增加的计算量。

Inception 机制将多个不同尺度的卷积核，池化进行整合，形成一个 Inception 模块，对图像进行多尺度处理。



1x1 卷积，不改变图像的高度和宽度，只改变通道数
所以用它来对数据降维，

为降低网络参数，用了 2 个方法：

- 1). 去除最后的全连接层，用全局平均池化代替
- 2). 精心设计的 inception 模块

④ VGG 创新点：采用小尺寸的卷积核来模拟更大尺寸的卷积核

e.g. 2 层连续的 3×3 卷积层可以达到 1 层 5×5 卷积层的感受野，但是所需参数量会更少。

⑤ ResNet：用跨层连接 (shortcut Connections) 构造残差项 (Residual Representations)

(有争议) ? 残差网络不是一个单一的超深网络，而是多个网络指代级的隐式集成，大多由一些浅层的网络组成，因此不能解决梯度消失问题

⇒ 进阶 ResNeXt

全连接层的作用：

Adding a fully-connected layer is (usually) a cheap way of learning non-linear combinations of the high-level features as represented by the output of the convolutional layer.

④ DenseNet 一种具有密集连接的卷积神经网络

优点：网络更窄，参数更少，减轻过拟合。对小数据集的学习有效果。

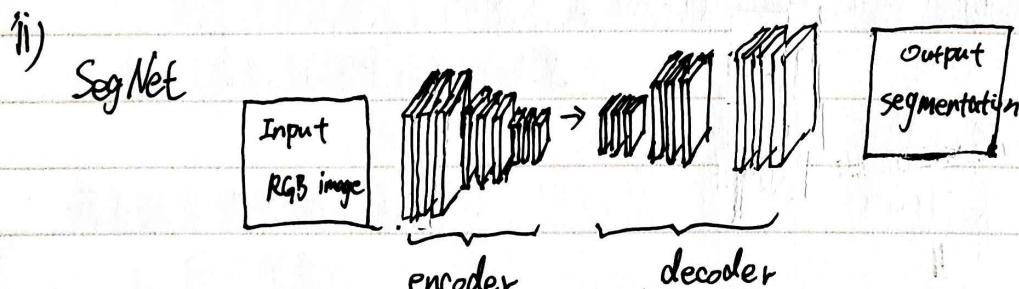
一些优化算法：

- 1). Mini-batch Gradient Descent
- 2). AdaGrad (adaptive gradient) 自适应梯度
- 3). Adam (Adaptive moment estimation)
- 4). RMSProp
- 5). NAG

图像语义分割和图像识别：全卷积网络：能接受任意尺寸的输入图像，并产生相同尺寸的输出图像。
输入图像和输出图像的像素一一对应，这种网络支持端到端、像素到像素的训练。

i) 思路：在网络的最后，通过反卷积实现上采样

DeepLab：
①用上采样的滤波器进行卷积（atrous 卷积）以实现密集的、对像素级的预测。
②用 atrous 空间金字塔下采样技术，以实现对物体的多尺度分割
③用概率图模型，实现更精准的目标边界定位。

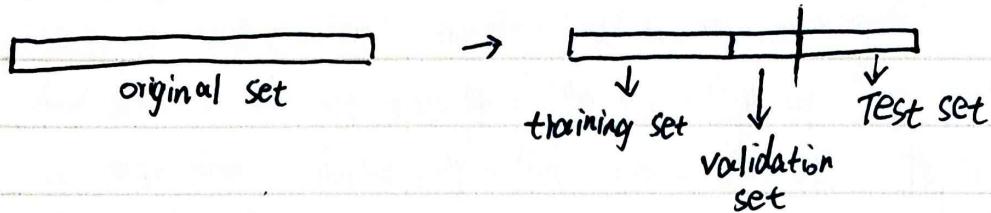


产生语义信息
的特征图像 将编码器网络输出的低分辨率
特征图像映射回输入图像的尺寸

交叉验证法：用于评估模型的预测性能，在一定程度上减小过拟合

1. 留出法 holdout cross validation

when dataset $> 100,000$



2. K 折交叉验证 K-fold cross validation 当数据集 $< 100,000$ observations

(K一般取10)

① 将原始数据分为K份

② 每一次选其中1份作为测试集，剩余 $K-1$ 份作为训练集

③ 重复k次，这样每个子集都能有一次机会作为测试集，其余机会为训练集。

计算K组测试结果的平均值作为模型精度的估计

3. 留一法 Leave one out Cross Validation

把K换成m

4. Bootstrap

通过自取样法，即在含有m个样本的数据集中进行m次有放回地随机抽样，组成新数据集作为训练集

超参数有哪些调优方法

1. 网络搜索

2. 随机搜索

3. 贝叶斯优化算法：

1. When using MSE, increasing capacity lowers bias but raises variance
2. Bayesian methods typically generalize much better when limited train data is available, but they typically suffer from high computational cost when the number of training examples is large
3. Pooling is essential for handling inputs of varying size.
4. Zero padding allows us to control the kernel width and output size independently to stop shrinkage
5. To evaluate convolutional architectures, use randomized weights and only train the last layer.
6. The most effective sequence models in practice are
gated RNNs {
 LSTMs : learn long-term dependencies more easily than simple RNNs
 Gated recurrent units.
Adding a bias of 1 to the forget gate makes the LSTMs as strong as gated RNNs (GRN) variants
7. Good baseline optimization algorithms:
 - ① SGD with momentum and a decaying learning rate
 - ② ADAM
8. If there are < 10 million samples in the training set,
 - ① almost always use early stopping
 - ② Dropout is a good choice and batch normalization is a possible alternative

The LR is the most important hyperparameter because it controls the effective model capacity in a more complicated way than other hyperparameters

As long as your training error is low, you can always ↓ generalization error by collecting more training data.

Standardizing pixel ranges is the only strict preprocessing required for Computer vision

Contrast normalization : a safe computer vision preprocessing step
- Global contrast normalization (GCN)

$$\text{Weight decay} = 1e-1 \Rightarrow \text{only in NN}$$

$$L(x, w) = \text{MSE}(M(x, w), y) + \underbrace{wd \cdot \sum w^2}_{\text{L2 regularization}}$$

$$\frac{dwd \cdot w^2}{w} = \cancel{\cancel{2}} \underbrace{wd}_{\downarrow} w$$

(delete) $1e-5$

L2 regularization

~~Adam~~; RMSprop

$$\text{Adam} = \text{Momentum} + \text{RMSprop}$$

Random Forest

- A universal ML tech.
- It can predict sth. that can be of any kind
 - { Category \Rightarrow classification \Rightarrow classifier
 - { Continuous variable \Rightarrow regression \Rightarrow regressor
- It can predict with columns of any kind
 - (both structured and unstructured data)
- Does not generally overfit badly
- Do not need a separate validation set in general
- Has few statistical assumptions
- Requires very few pieces of feature engineering.

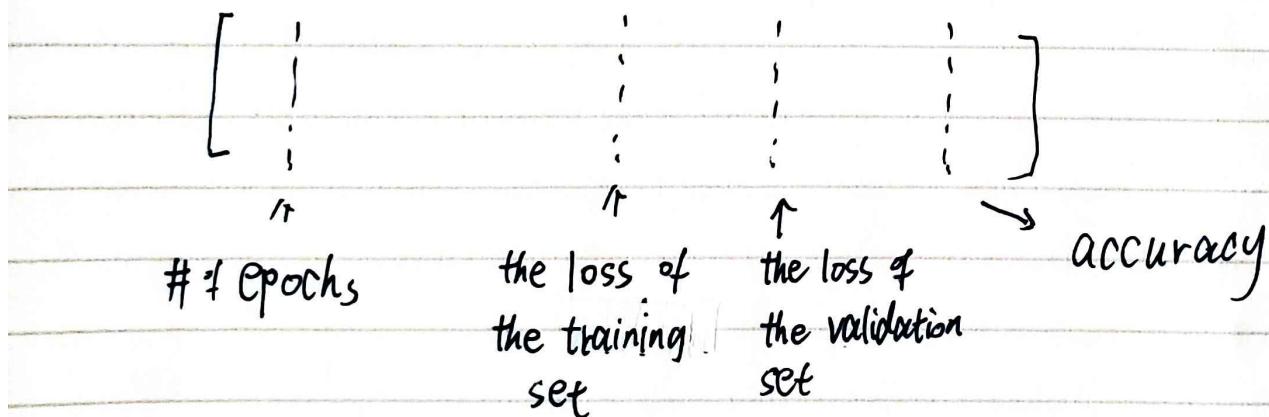
Deep Learning :

1. Infinitely flexible function (DL has multiple hidden layers)
 2. All-purpose parameter fitting (Gradient Descent)
 3. Fast and scalable (GPU)

Paper.

semantic style transfer and turning two-bit doodles into fine artworks by Alex J. Champandard

learn.fit(
 ↓ ↓
 LR # of epochs



用 NN 对 cat data 和 cont data 建模：

对类别变量 (categorical) : embeddings

对连续变量 (continuous) : NN

MAE (L_1 损失) 对异常点有更好的鲁棒性.

MSE (L_2 损失)