

Image Classification Pipe line

모두의연구소
박은수 Research Director

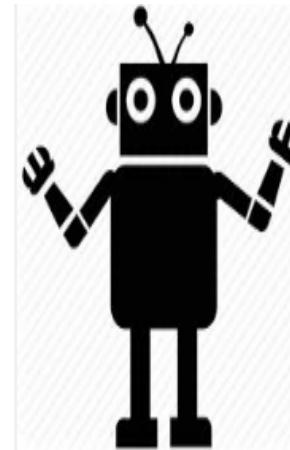
딥러닝의 30가지 적용사례

<https://brunch.co.kr/@itschloe1/23>

물체인식에서의 딥러닝



인간 vs. 기계



물체인식에서의 딥러닝



$$43214 \times 245124 = ?$$

물체인식에서의 딥러닝

A screenshot of a Google search results page. The search bar at the top contains the query "43214 x 245124". Below the search bar, there are several search filters: 전체 (selected), 지도, 이미지, 뉴스, 동영상, 더보기 ▾, and 검색 도구. The main search results area shows the text "검색결과 약 125개 (0.65초)". Below this, a snippet of a search result is shown, containing the equation "43 214 x 245 124 =" followed by the result "10592788536".



물체인식에서의 딥러닝



고양이랑 강아지 구분하기

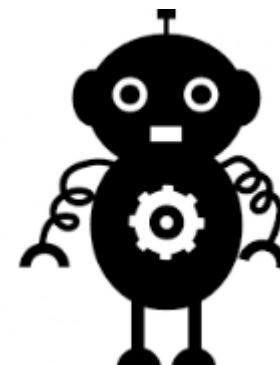


물체인식에서의 딥러닝

왼쪽 고양이
오른쪽 강아지 ~ !!



아...



물체인식에서의 딥러닝

컴퓨터 비전에서의 물체인식



{고양이, 강아지, 버스, 자동차, ..., 집}

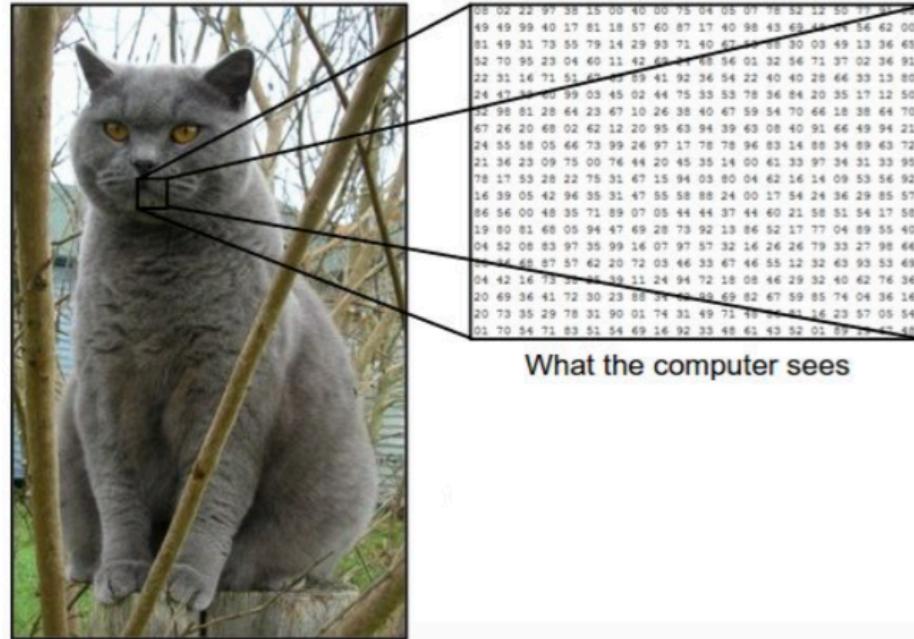
→ 고양이

물체인식에서의 딥러닝

컴퓨터 비전에서의 물체인식

영상은 $[0, 255]$ 범위의 값
을 갖는 RGB 3채널로 구
성되어 있다

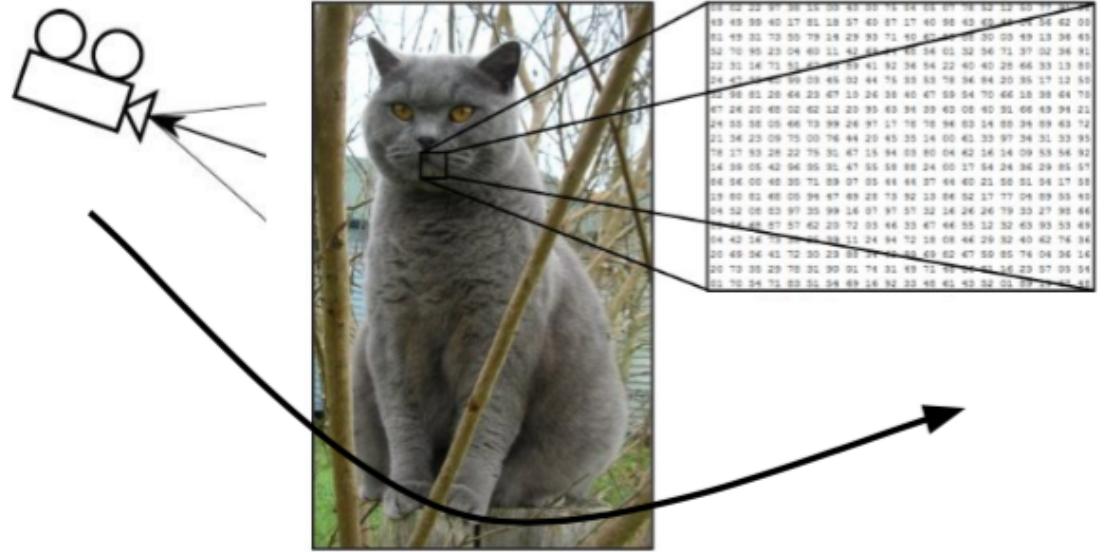
옆의 고양이는
300x100x3 의 배열
(3은 RGB (Red,
Green, Blue) 채널)



물체인식에서의 딥러닝

컴퓨터 비전에서의 물체인식의 어려움

카메라 위치에 따른
변화



물체인식에서의 딥러닝

컴퓨터 비전에서의 물체
인식의 어려움

조명에 따른 변화



물체인식에서의 딥러닝

컴퓨터 비전에서의 물체인식의 어려움

객체 형태의 변화



* From : Lecture slide 2 of "CS231n : Convolutional Neural Networks for Visual Recognition"

물체인식에서의 딥러닝

컴퓨터 비전에서의 물체인식의 어려움

가려짐



* From : Lecture slide 2 of "CS231n : Convolutional Neural Networks for Visual Recognition"

물체인식에서의 딥러닝

컴퓨터 비전에서의 물체인식의 어려움
배경과의 유사성



물체인식에서의 딥러닝

컴퓨터 비전에서의 물체인식의 어려움

같으면서 다른 형태의 물체



물체인식에서의 딥러닝



컴퓨터 비전에서의 물체인식의 어려움

```
def predict(image):
    # *****
    return class_label
```

고양이와 개를 구분하는 알고리즘을
만들어 보세오

물체인식에서의 딥러닝



컴퓨터 비전에서의 물체인식의 어려움

죄송합니다

어려운걸 알고있습니다

Data-driven 방법

1. 이미지와 그에 해당하는 레이블(label)을 모음
2. 머신러닝 방법론으로 분류기를 학습
3. 학습된 분류기를 테스트 이미지에 평가

```
def train(train_images, train_labels):  
    # build a model for images -> labels...  
    return model  
  
def predict(model, test_images):  
    # predict test_labels using the model...  
    return test_labels
```

Example training set



첫번째 분류기 : Nearest Neighbor Classifier



간단한 분류기를 만들어 봅시다 :
Nearest Neighbor Classifier

```
def train(train_images, train_labels):
    # build a model for images -> labels...
    return model

def predict(model, test_images):
    # predict test_labels using the model...
    return test_labels
```

- 모든 학습용 이미지와 그들의 레이블(label)을 저장
- 학습용 이미지 중 가장 유사한 이미지의 레이블 반환

첫번째 분류기 : Nearest Neighbor Classifier

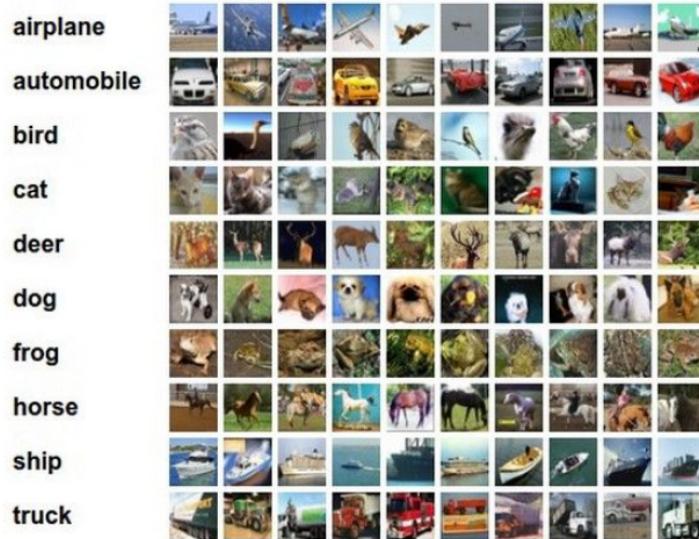
Example dataset: **CIFAR-10**

10 labels

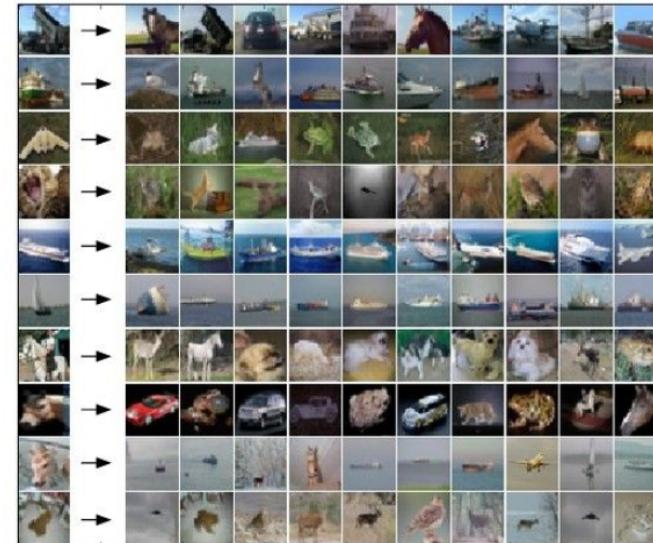
50,000 training images

10,000 test images.

32x32 크기임



For every test image (first column),
examples of nearest neighbors in rows



가장 유사한 그림은 어떻게 찾는가?

Distance 를 측정

L1 distance: $d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$

test image				training image				pixel-wise absolute value differences			
56	32	10	18	10	20	24	17	46	12	14	1
90	23	128	133	8	10	89	100	82	13	39	33
24	26	178	200	12	16	178	170	12	10	0	30
2	0	255	220	4	32	233	112	2	32	22	108

- = add → 456

```
import numpy as np

class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where each row is an example. Y is 1-dimension of size N """
        # the nearest neighbor classifier simply remembers all the training data
        self.Xtr = X
        self.ytr = y

    def predict(self, X):
        """ X is N x D where each row is an example we wish to predict label for """
        num_test = X.shape[0]
        # lets make sure that the output type matches the input type
        Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

        # loop over all test rows
        for i in xrange(num_test):
            # find the nearest training image to the i'th test image
            # using the L1 distance (sum of absolute value differences)
            distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
            min_index = np.argmin(distances) # get the index with smallest distance
            Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

        return Ypred
```

Nearest Neighbor classifier

```
import numpy as np

class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where each row is an example. Y is 1-dimension of size N """
        # the nearest neighbor classifier simply remembers all the training data
        self.Xtr = X
        self.ytr = y

    def predict(self, X):
        """ X is N x D where each row is an example we wish to predict label for """
        num_test = X.shape[0]
        # lets make sure that the output type matches the input type
        Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

        # loop over all test rows
        for i in xrange(num_test):
            # find the nearest training image to the i'th test image
            # using the L1 distance (sum of absolute value differences)
            distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
            min_index = np.argmin(distances) # get the index with smallest distance
            Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

        return Ypred
```

Nearest Neighbor classifier

remember the training data

```
import numpy as np

class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where each row is an example. Y is 1-dimension of size N """
        # the nearest neighbor classifier simply remembers all the training data
        self.Xtr = X
        self.ytr = y

    def predict(self, X):
        """ X is N x D where each row is an example we wish to predict label for """
        num_test = X.shape[0]
        # lets make sure that the output type matches the input type
        Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

        # loop over all test rows
        for i in xrange(num_test):
            # find the nearest training image to the i'th test image
            # using the L1 distance (sum of absolute value differences)
            distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
            min_index = np.argmin(distances) # get the index with smallest distance
            Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

        return Ypred
```

Nearest Neighbor classifier

for every test image:

- find nearest train image with L1 distance
- predict the label of nearest training image

```

import numpy as np

class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where each row is an example. Y is 1-dimension of size N """
        # the nearest neighbor classifier simply remembers all the training data
        self.Xtr = X
        self.ytr = y

    def predict(self, X):
        """ X is N x D where each row is an example we wish to predict label for """
        num_test = X.shape[0]
        # lets make sure that the output type matches the input type
        Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

        # loop over all test rows
        for i in xrange(num_test):
            # find the nearest training image to the i'th test image
            # using the L1 distance (sum of absolute value differences)
            distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
            min_index = np.argmin(distances) # get the index with smallest distance
            Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

    return Ypred

```

질문 : 학습용 이미지의 수와
우리가 만든 Nearest
Neighbor Classifier의 처리
속도와의 관계는?

```

import numpy as np

class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where each row is an example. Y is 1-dimension of size N """
        # the nearest neighbor classifier simply remembers all the training data
        self.Xtr = X
        self.ytr = y

    def predict(self, X):
        """ X is N x D where each row is an example we wish to predict label for """
        num_test = X.shape[0]
        # lets make sure that the output type matches the input type
        Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

        # loop over all test rows
        for i in xrange(num_test):
            # find the nearest training image to the i'th test image
            # using the L1 distance (sum of absolute value differences)
            distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
            min_index = np.argmin(distances) # get the index with smallest distance
            Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

        return Ypred

```

질문 : 학습용 이미지의 수와
우리가 만든 Nearest
Neighbor Classifier의 처리
속도와의 관계는?

선형적임 ☹

이건 좋지않다.

- 실제 어플리케이션에선 테스트 타임이 중요하다.
- CNN은 반대임. 트레이닝이 오래 걸리고, 테스트는 빠름

Aside: Approximate Nearest Neighbor

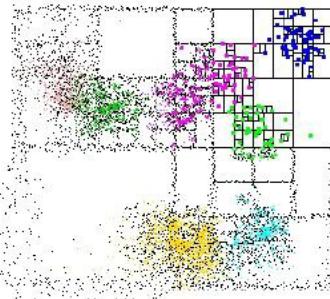
find approximate nearest neighbors quickly

ANN: A Library for Approximate Nearest Neighbor Searching

David M. Mount and Sunil Arya

Version 1.1.2

Release Date: Jan 27, 2010



What is ANN?

ANN is a library written in C++, which supports data structures and algorithms for both exact and approximate nearest neighbor searching in arbitrarily high dimensions.

In the nearest neighbor problem a set of data points in d-dimensional space is given. These points are preprocessed into a data structure, so that given any query point q , the nearest or generally k nearest points of P to q can be reported efficiently. The distance between two points can be defined in many ways. ANN assumes that distances are measured using any class of distance functions called Minkowski metrics. These include the well known Euclidean distance, Manhattan distance, and max distance.

Based on our own experience, ANN performs quite efficiently for point sets ranging in size from thousands to hundreds of thousands, and in dimensions as high as 20. (For applications in significantly higher dimensions, the results are rather spotty, but you might try it anyway.)

The library implements a number of different data structures, based on kd-trees and box-decomposition trees, and employs a couple of different search strategies.

The library also comes with test programs for measuring the quality of performance of ANN on any particular data sets, as well as programs for visualizing the structure of the geometric data structures.

FLANN - Fast Library for Approximate Nearest Neighbors

- Home
- News
- Publications
- Download
- Changelog
- Repository

What is FLANN?

FLANN is a library for performing fast approximate nearest neighbor searches in high dimensional spaces. It contains a collection of algorithms we found to work best for nearest neighbor search and a system for automatically choosing the best algorithm and optimum parameters depending on the dataset.

FLANN is written in C++ and contains bindings for the following languages: C, MATLAB and Python.

News

- (14 December 2012) Version 1.8.0 is out bringing incremental addition/removal of points to/from indexes
- (20 December 2011) Version 1.7.0 is out bringing two new index types and several other improvements.
- You can find binary installers for FLANN on the [Point Cloud Library](#) project page. Thanks to the [PCL developers!](#)
- Mac OS X users can install flann through MacPorts (thanks to Mark Moll for maintaining the Portfile)
- New release introducing an easier way to use custom distances, kd-tree implementation optimized for low dimensionality search and experimental MPI support
- New release introducing new C++ templated API, thread-safe search, save/load of indexes and more.
- The FLANN license was changed from LGPL to BSD.

How fast is it?

In our experiments we have found FLANN to be about one order of magnitude faster on many datasets (in query time), than previously available approximate nearest neighbor search software.

Publications

More information and experimental results can be found in the following papers:

- Marius Muja and David G. Lowe, "Scalable Nearest Neighbor Algorithms for High Dimensional Data", Pattern Analysis and Machine Intelligence (PAMI), Vol. 36, 2014. [[PDF](#)] [[BibTeX](#)]
- Marius Muja and David G. Lowe, "Fast Matching of Binary Features", Conference on Computer and Robot Vision (CRV) 2012. [[PDF](#)] [[BibTeX](#)]
- Marius Muja and David G. Lowe, "Fast Approximate Nearest Neighbors with Automatic Algorithm Configuration", in International Conference on Computer Vision Theory and Applications (VISAPP'09), 2009. [[PDF](#)] [[BibTeX](#)]

어떤 Distance 함수를 사용할까?



Distance 함수는 여러개입니다. 여기서
어떤 걸 사용하는게 좋을까요?

이렇게 변형가능한 우리의 설정 같은 것들을 하이버 파라미터라고 합니다

일반적으로 아래의 2개를 사용합니다

L1 (Manhattan) distance

$$d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$

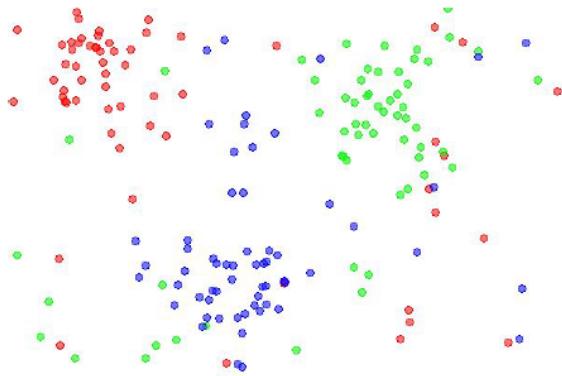
L2 (Euclidean) distance

$$d_2(I_1, I_2) = \sqrt{\sum_p (I_1^p - I_2^p)^2}$$

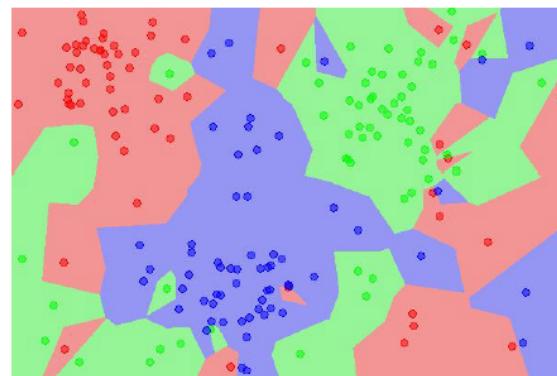
Nearest Neighbor은 알겠습니다. 제일 비슷한 것 찾아서 그 비슷한 것의 레이블을 반환하는 것이군요

k-Nearest Neighbor은 무엇인가?

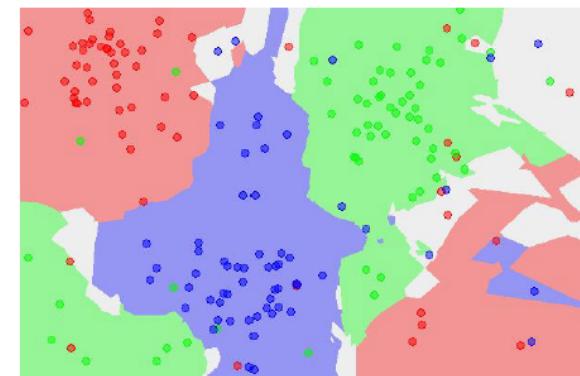
the data



NN classifier



5-NN classifier



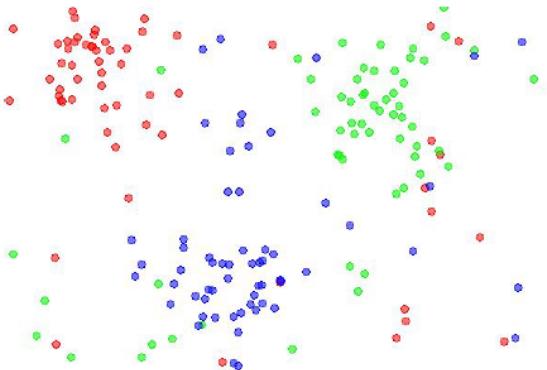
http://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm

Nearest Neighbor은 알겠습니다. 제일 비슷한 것 찾아서 그 비슷한 것의 레이블을 반환하는 것이군요

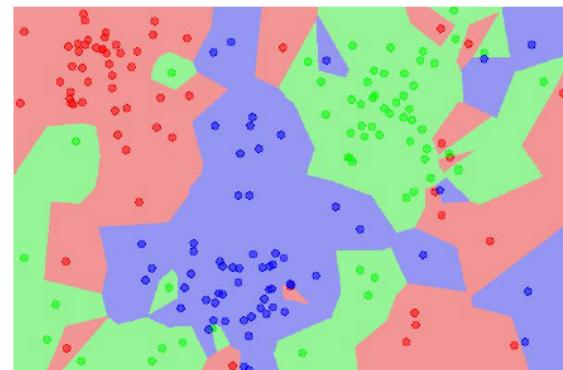
k-Nearest Neighbor은 무엇인가?

k개의 Nearest Neighbor를 찾고, k개 중 많이 나온 레이블을 선택 (voting)

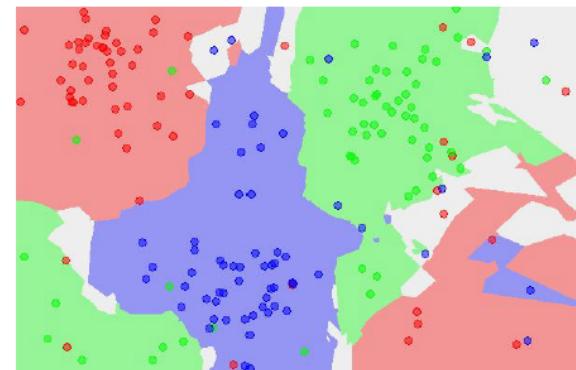
the data



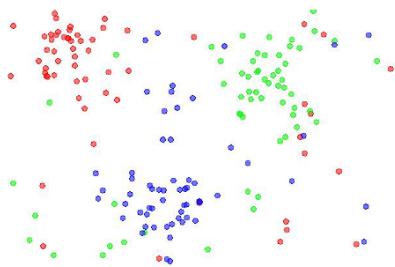
NN classifier



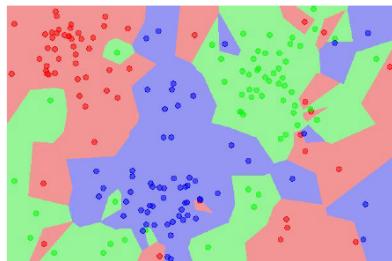
5-NN classifier



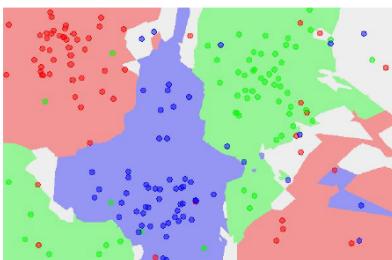
the data



NN classifier



5-NN classifier



이제 NN (Nearest Neighbor) 분류기
를 만들었다고 치자

질문 : NN 분류기(Nearest Neighbor Classifier)를 학습할때 사용한 이미지를 테스트 할때 넣어보면 정확도는 ?

질문 : k-NN 분류기(Nearest Neighbor Classifier)에 학습할때 사용한 이미지를 테스트 할때 넣어보면 정확도는 ?

고민이 생깁니다 ...



- Distance 함수로는 뭘 사용하는게 좋을까?

- k 값을 뭘로 결정하는게 좋을까?

이러한 파라미터들을 하이퍼 파라미터 (hyper parameter)라고 합니다

고민이 생깁니다 ...

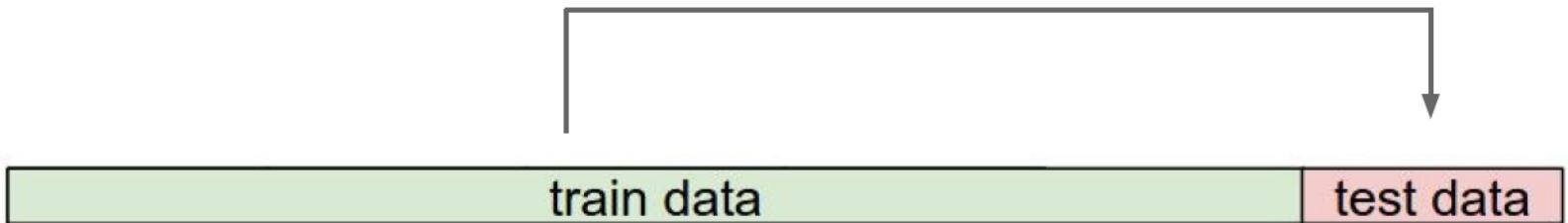
- Distance 함수로는 뭘 사용하는게 좋을까?
- k 값을 뭘로 결정하는게 좋을까?

이러한 파라미터들을 하이퍼 파라미터 (hyper parameter)라고 합니다

실제 풀고자 하는 문제에 따라 다르다

그러니 가능한 경우의 수들을 실험해 보고 결과
를 살펴봐야 한다

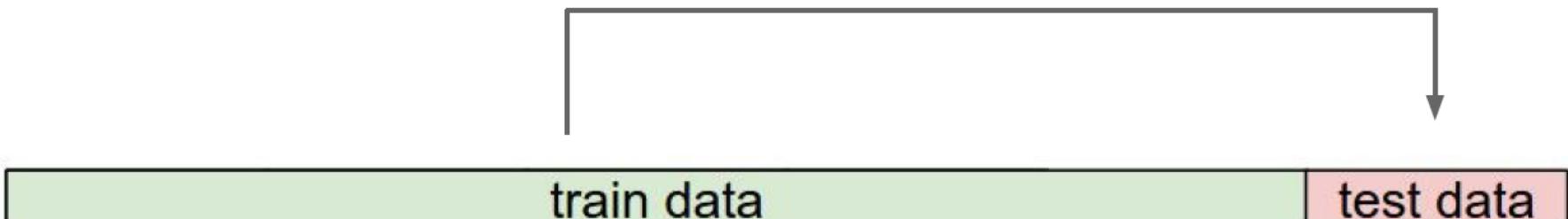
하이퍼 파라미터를 선택 한 후 테스트 데
이터를 통해 확인

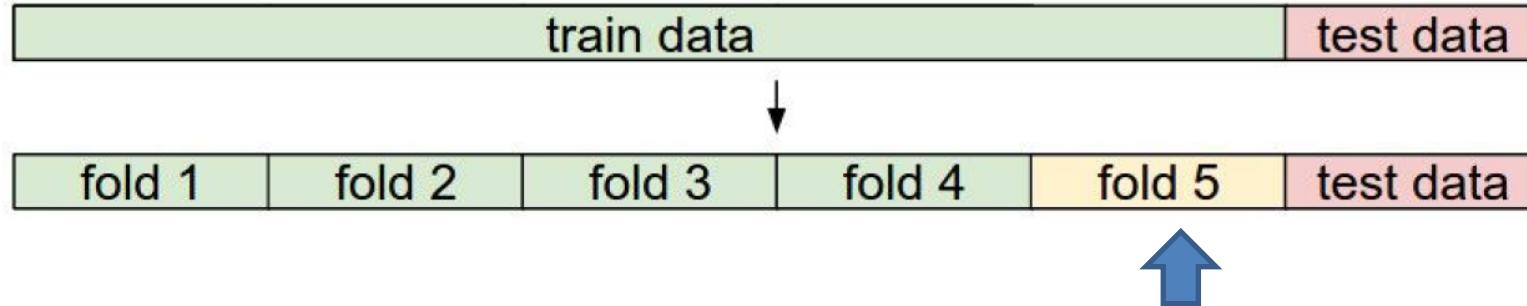


하이퍼 파라미터를 선택 한 후 테스트 데 이터를 통해 확인

안 좋은 생각이다. 테스트 데이터는 실제 환경에서도 잘 작동됨
을 보여야 할 최후의 보루

마지막을 위해 아껴라~!!



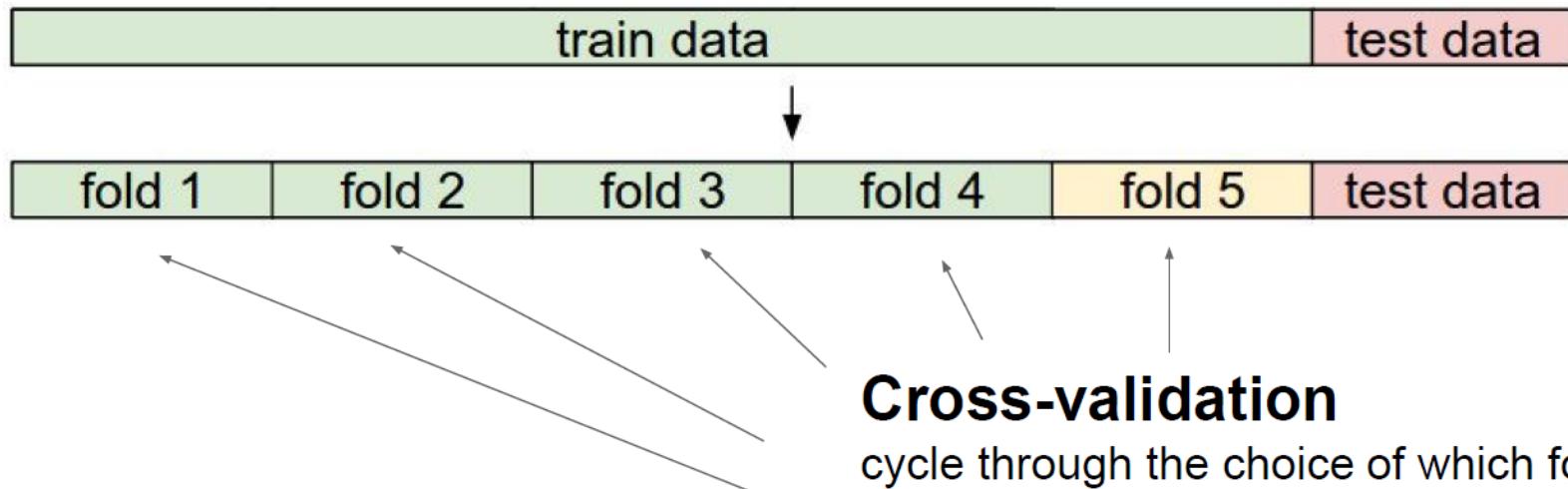


트레이닝(training)에 사용 할 데이터의 일정
부분을 떼어 내어 하이퍼파라미터를 튜닝하자

Validation Data

즉, fold1, fold2, fold3, fold4를 트레이닝에 사용하
고 테스트를 fold5에 수행해본다

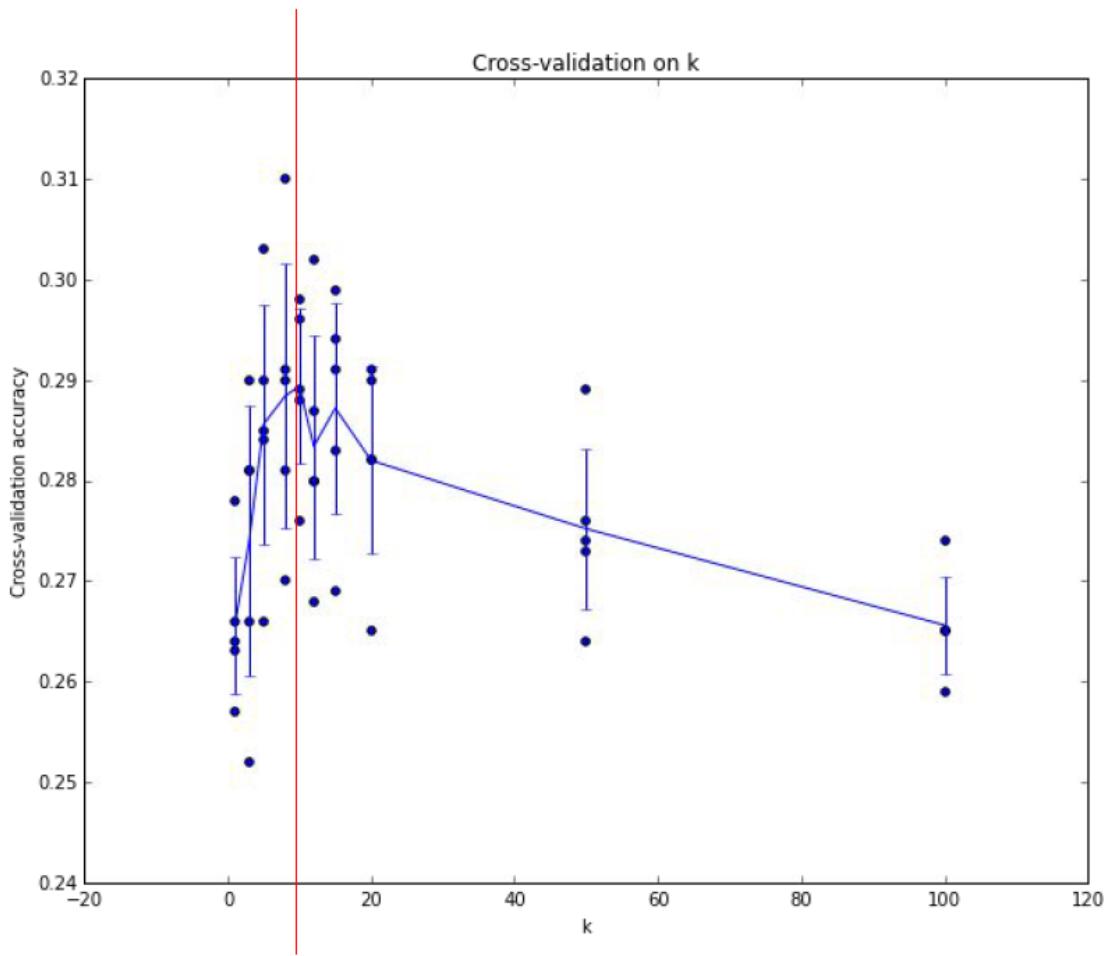
Cross-Validation



Cross-validation

cycle through the choice of which fold
is the validation fold, average results.

이제 테스트에 루프를 돌면서 트레이닝에 사용할 데이터와 테스
트에 사용할 데이터 셋을 바꿔가면서 실험해 봅니다



Example of
5-fold cross-validation
for the value of **k**.

Each point: single
outcome.

The line goes
through the mean, bars
indicated standard
deviation

(Seems that $k \approx 7$ works best
for this data)

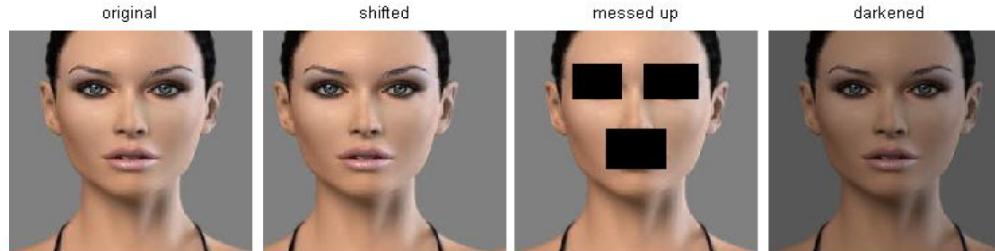
- 참고하면 좋은 동영상
- Lecture 2 | Image Classification
 - <https://www.youtube.com/watch?v=OoUX-nOEjG0&list=PL3FW7Lu3i5JvHM8ljYj-zLfQRF3EO8sYv>
- Andrew Ng: Artificial Intelligence is the New Electricity
 - <https://www.youtube.com/watch?v=21EiKfQYZXc&t=3105s>

Linear Classification

모두의연구소
박은수 Research Director

k-Nearest Neighbor 로 영상인식은 하지 않습니다

- 이걸 적용할 경우 테스트 결과가 매우 안좋습니다
- 단순한 픽셀 단위의 Distance metrics로는 좋은 결과를 기대하기 어렵습니다



(all 3 images have same L2 distance to the one on the left)

그럼 뭘로 하나요 ?

픽셀값들을 조합하여 단순히 각 레이블에 대한 점수를
부여하면 어떻까요? Score function?

Parametric Approach



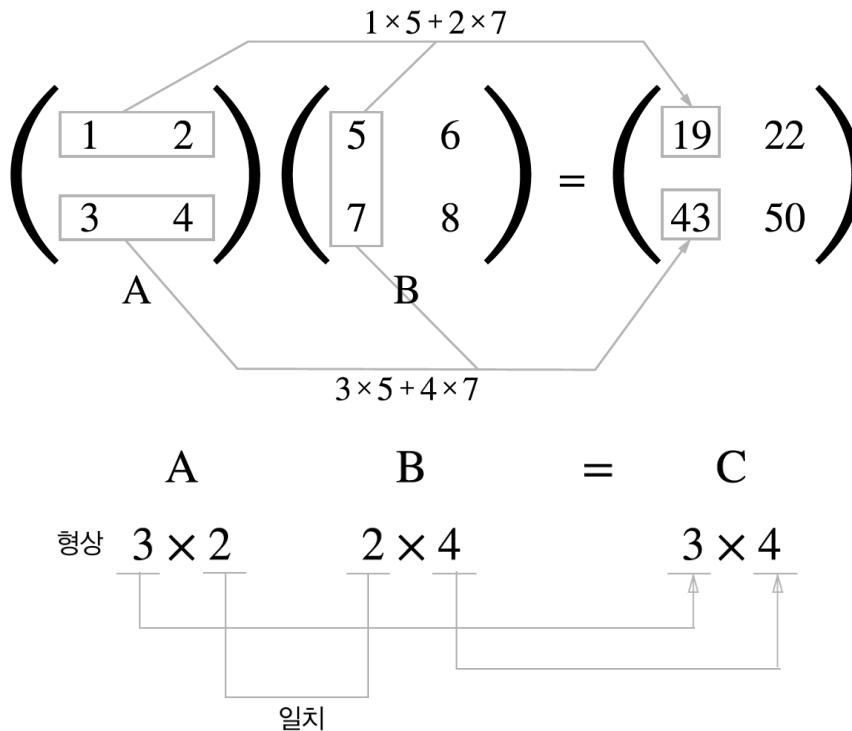
- 분류기의 구성
 - Score function
 - Loss function
 - Optimization



들어가기 전에

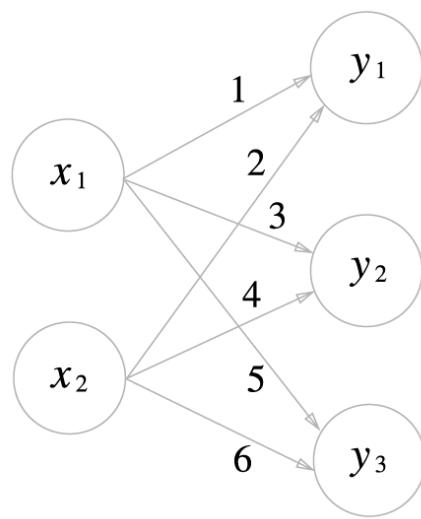
다차원 배열의 계산

- **다차원 배**



다차원 배열의 계산

- **다차원 배열**



$$\begin{array}{c}
 \left(\begin{array}{ccc} 1 & 3 & 5 \\ 2 & 4 & 6 \end{array} \right) \\
 X \quad W = Y \\
 \hline
 2 \quad 2 \times 3 \quad 3 \\
 \text{일치}
 \end{array}$$

Matrix 연산 리뷰 : Matrix-vector

Example

$$\begin{bmatrix} 1 & 3 \\ 4 & 0 \\ 2 & 1 \end{bmatrix}_{3 \times 2} \times \begin{bmatrix} 1 \\ 5 \end{bmatrix}_{2 \times 1} = \begin{bmatrix} 16 \\ 4 \\ 7 \end{bmatrix}_{3 \times 1} \text{ matrix}$$

$$1 \times 1 + 3 \times 5 = 16$$

$$4 \times 1 + 0 \times 5 = 4$$

$$2 \times 1 + 1 \times 5 = 7$$

Details:

$$\begin{array}{c} A \\ \times \\ \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \\ = \\ \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \end{array}$$

A : $m \times n$ matrix (m rows, n columns)
 x : $n \times 1$ matrix (n-dimensional vector)
 y : $m \times 1$ matrix (m-dimensional vector)

To get y_i , multiply A 's i^{th} row with elements of vector x , and add them up.

Matrix 연산 리뷰 : Matrix-vector

Example

$$\begin{bmatrix} 1 & 2 & 1 & 5 \\ 0 & 3 & 0 & 4 \\ -1 & -2 & 0 & 0 \end{bmatrix}_{3 \times 4} \begin{bmatrix} 1 \\ 3 \\ 2 \\ 1 \end{bmatrix}_{4 \times 1} = \begin{bmatrix} 14 \\ 13 \\ -7 \end{bmatrix}_{3 \times 1}$$

$$\left. \begin{aligned} 1 \times 1 + 2 \times 3 + 1 \times 2 + 5 \times 1 &= 14 \\ 0 \times 1 + 3 \times 3 + 0 \times 2 + 4 \times 1 &= 13 \\ -1 \times 1 + (-2) \times 3 + 0 \times 2 + 0 \times 1 &= -7 \end{aligned} \right]$$

House sizes:
 → 2104
 → 1416
 → 1534
 → 852

Matrix $\begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix}_{4 \times 2}$

Vector $\begin{bmatrix} -40 \\ 0.25 \end{bmatrix}_{2 \times 1}$

$h_{\theta}(x) = -40 + 0.25x$

$h_{\theta}(x)$

$h_{\theta}(2104)$

$h_{\theta}(1416)$

$h_{\theta}(1534)$

$h_{\theta}(852)$

$h_{\theta}(x) = \begin{bmatrix} -40 \times 1 + 0.25 \times 2104 \\ -40 \times 1 + 0.25 \times 1416 \\ -40 \times 1 + 0.25 \times 1534 \\ -40 \times 1 + 0.25 \times 852 \end{bmatrix}_{4 \times 1}$

$\text{prediction} = \text{Data Matrix} \otimes \text{Parameters}$

for $i = 1:1000$, prediction(i) = ...

Matrix 연산 리뷰 : Matrix-matrix

Example

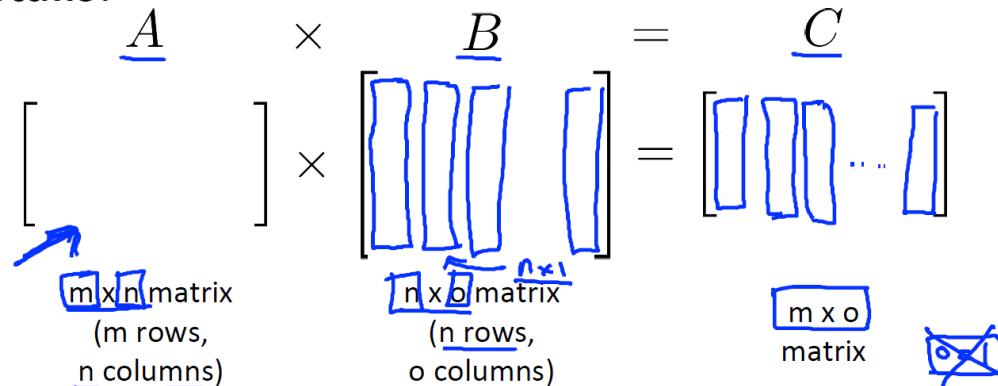
$$\begin{bmatrix} 1 & 3 & 2 \\ 4 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 3 \\ 0 & 1 \\ 5 & 2 \end{bmatrix} = \begin{bmatrix} 11 & 10 \\ 9 & 14 \end{bmatrix}$$

(2x3) (3x2)

$$\begin{bmatrix} 1 & 3 & 2 \\ 4 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 0 \\ 5 \end{bmatrix} = \begin{bmatrix} 11 \\ 9 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 3 & 2 \\ 4 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} 3 \\ 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 10 \\ 14 \end{bmatrix}$$

Details:



The i^{th} column of the matrix C is obtained by multiplying A with the i^{th} column of B . (for $i = 1, 2, \dots, o$)

Matrix 연산 리뷰 : Matrix-matrix

House sizes:

$$\left\{ \begin{array}{r} 2104 \\ 1416 \\ 1534 \\ 852 \end{array} \right.$$

Have 3 competing hypotheses:

$$1. h_{\theta}(x) = -40 + 0.25x$$

$$2. h_{\theta}(x) = 200 + 0.1x$$

$$3. h_{\theta}(x) = -150 + 0.4x$$

Matrix

$$\begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix}$$

Matrix

$$\begin{bmatrix} -40 \\ 0.25 \\ 200 \\ 0.1 \\ -150 \\ 0.4 \end{bmatrix}$$

$$\begin{bmatrix} 486 \\ 314 \\ 344 \\ 353 \\ 173 \\ 285 \\ 692 \\ 342 \\ 464 \\ 191 \end{bmatrix}$$

Prediction
of 1st
h_θ

Predictions
of 2nd
h_θ

선형대수 리뷰

<http://cs229.stanford.edu/section/cs229-linalg.pdf>

좀 더 깊은 리뷰를 보시려
면 참고하세요

2.1 Vector-Vector Products

Given two vectors $x, y \in \mathbb{R}^n$, the quantity $x^T y$, sometimes called the *inner product* or *dot product* of the vectors, is a real number given by

$$x^T y \in \mathbb{R} = \begin{bmatrix} x_1 & x_2 & \cdots & x_n \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \sum_{i=1}^n x_i y_i.$$

Observe that inner products are really just special case of matrix multiplication. Note that it is always the case that $x^T y = y^T x$.

Given vectors $x \in \mathbb{R}^m, y \in \mathbb{R}^n$ (not necessarily of the same size), $xy^T \in \mathbb{R}^{m \times n}$ is called the *outer product* of the vectors. It is a matrix whose entries are given by $(xy^T)_{ij} = x_i y_j$, i.e.,

$$xy^T \in \mathbb{R}^{m \times n} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} \begin{bmatrix} y_1 & y_2 & \cdots & y_n \end{bmatrix} = \begin{bmatrix} x_1 y_1 & x_1 y_2 & \cdots & x_1 y_n \\ x_2 y_1 & x_2 y_2 & \cdots & x_2 y_n \\ \vdots & \vdots & \ddots & \vdots \\ x_m y_1 & x_m y_2 & \cdots & x_m y_n \end{bmatrix}.$$

As an example of how the outer product can be useful, let $\mathbf{1} \in \mathbb{R}^n$ denote an n -dimensional vector whose entries are all equal to 1. Furthermore, consider the matrix $A \in \mathbb{R}^{m \times n}$ whose columns are all equal to some vector $x \in \mathbb{R}^m$. Using outer products, we can represent A compactly as,

$$A = \begin{bmatrix} | & | & | & | \\ x & x & \cdots & x \\ | & | & \cdots & | \end{bmatrix} = \begin{bmatrix} x_1 & x_1 & \cdots & x_1 \\ x_2 & x_2 & \cdots & x_2 \\ \vdots & \vdots & \ddots & \vdots \\ x_m & x_m & \cdots & x_m \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} \begin{bmatrix} 1 & 1 & \cdots & 1 \end{bmatrix} = x\mathbf{1}^T.$$

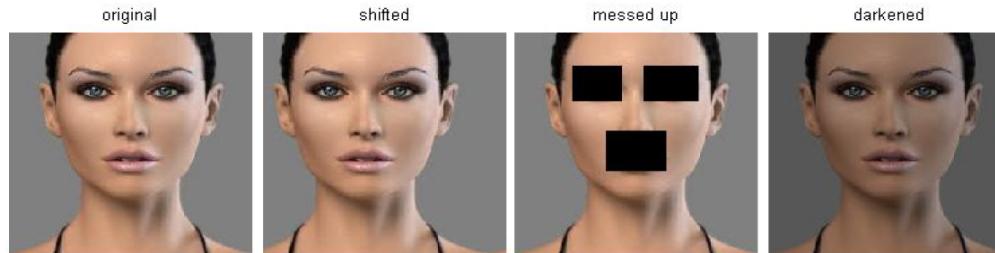
2.2 Matrix-Vector Products

Given a matrix $A \in \mathbb{R}^{m \times n}$ and a vector $x \in \mathbb{R}^n$, their product is a vector $y = Ax \in \mathbb{R}^m$. There are a couple ways of looking at matrix-vector multiplication, and we will look at each of them in turn.

If we write A by rows, then we can express Ax as,

$$y = Ax = \begin{bmatrix} \cdots & a_1^T & \cdots \\ \cdots & a_2^T & \cdots \\ \vdots & \vdots & \ddots & \vdots \\ \cdots & a_m^T & \cdots \end{bmatrix} x = \begin{bmatrix} a_1^T x \\ a_2^T x \\ \vdots \\ a_m^T x \end{bmatrix}.$$

다시 돌아와서 ...

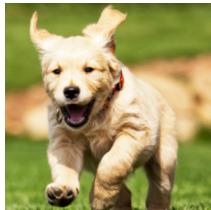


(all 3 images have same L2 distance to the one on the left)

픽셀값들을 조합하여 단순히 각 레이블에 대한 점수를
부여하면 어떻까요? **Score function?**

Parametric approach

- Score function



$$f(\text{dog image}, W)$$



- 강아지 점수
- 고양이 점수



$$f(\text{cat image}, W)$$



- 강아지 점수
- 고양이 점수

:

:

:

점수가
높은것으
로 분류

Parametric approach

- Score function : Simple Linear Classifier

$$f(\mathbf{x}, \mathbf{W}) = \mathbf{W}\mathbf{x} \quad 3,072 \times 1$$

32x32x3 영상 =
3,072 픽셀



1	155
2	200
...	...
3071	110
3072	78

\mathbf{x}

$$f(\mathbf{x}, \mathbf{W})$$

- 강아지 점수
- 고양이 점수

Parametric approach

- Score function : Simple Linear Classifier

$$2 \times 1 \text{ } f(\mathbf{x}, \mathbf{W}) = \mathbf{W} \mathbf{x} \quad 3,072 \times 1$$

32x32x3 영상 =
3,072 픽셀



1	155
2	200
...	...
3071	110
3072	78

\mathbf{x}

$$f(\mathbf{x}, \mathbf{W})$$

- 강아지 점수
- 고양이 점수

0.3
0.7

Parametric approach

- Score function : Simple Linear Classifier

$$2 \times 1 \quad f(\mathbf{x}, \mathbf{W}) = \mathbf{W}\mathbf{x} \quad 3,072 \times 1$$

2x3,072

32x32x3 영상 =
3,072 픽셀



	155
1	200
2	...
...	...
3071	110
3072	78

X

$$f(\mathbf{x}, \mathbf{W})$$

- 강아지 점수
- 고양이 점수

0.2	0.1	...	0.3	0.7
0.7	0.8	...	0.9	0.4

0.3
0.7

Parametric approach

- Score function : Simple Linear Classifier

$$f(x, W) = Wx + b$$

3,072x1
 2x1 2x3,072

32x32x3 영상 =
3,072 픽셀



1	155
2	200
...	...
3071	110
3072	78

x

0.2	0.1	...	0.3	0.7
0.7	0.8	...	0.9	0.4

- 강아지 점수
- 고양이 점수

0.3
0.7

Parametric approach

- Score function : Simple Linear Classifier

$$f(x, W) = Wx + b$$

$$32 \times 32 \times 3 \text{ 영상} = \\ 3,072 \text{ 픽셀}$$



2x3,072

강아지 파라미터

0.2	0.1	...	0.3	0.7
0.7	0.8	...	0.9	0.4

고양이 파라미트

W

X	155
	200
	...
	110
	78

0.1
0.2

b

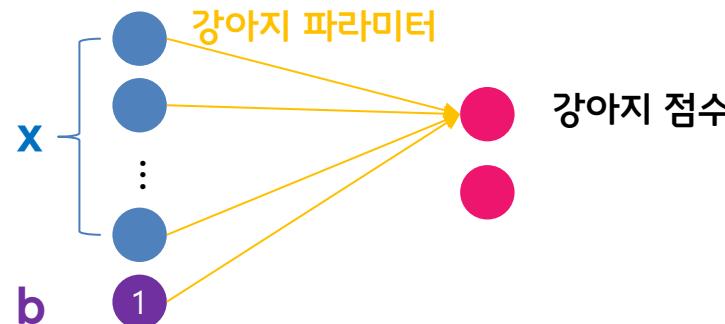
0.3

강아지 점수
고양이 점수

Parametric approach

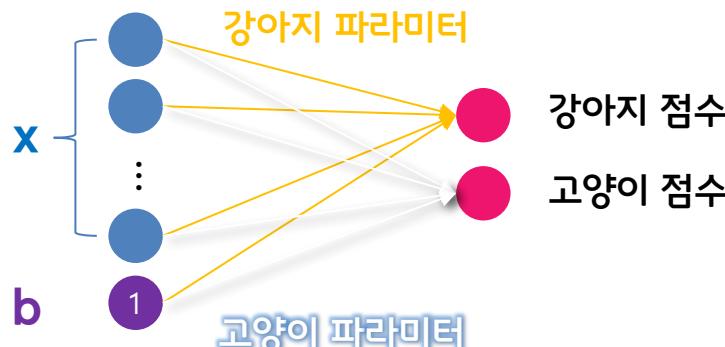
- Score function : Simple Linear Classifier

$$\begin{array}{c}
 \text{강아지 파라미터} \\
 \text{고양이 파라미터} \\
 \hline
 \begin{matrix} 0.2 & 0.1 & \dots & 0.3 & 0.7 \\ 0.7 & 0.8 & \dots & 0.9 & 0.4 \end{matrix}
 \end{array}
 \times \mathbf{W} \times \mathbf{x} + \mathbf{b} = \begin{array}{c} \text{강아지 점수} \\ \text{고양이 점수} \\ \hline \begin{matrix} 0.3 \\ 0.7 \end{matrix} \end{array}$$



Parametric approach

- Score function : Simple Linear Classifier



선형분류기의 뉴럴 네트워크 표현

Parametric approach

- Score function : Simple Linear Classifier

The diagram shows the calculation of a dog's score (강아지 점수) based on its parameters (고양이 파라미터). The formula is:

$$\text{강아지 점수} = \text{고양이 파라미터} \times \mathbf{W} + \mathbf{b}$$

The parameters (고양이 파라미터) are represented by a green matrix:

0.2	0.1	...	0.3	0.7
0.7	0.8	...	0.9	0.4

The weight matrix (\mathbf{W}) is:

155
200
...
110
78

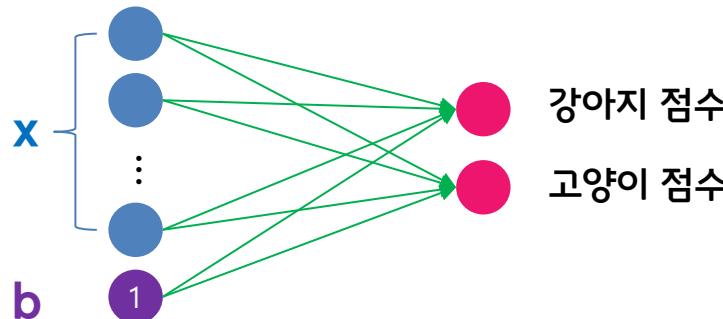
The bias vector (\mathbf{b}) is:

0.1
0.2

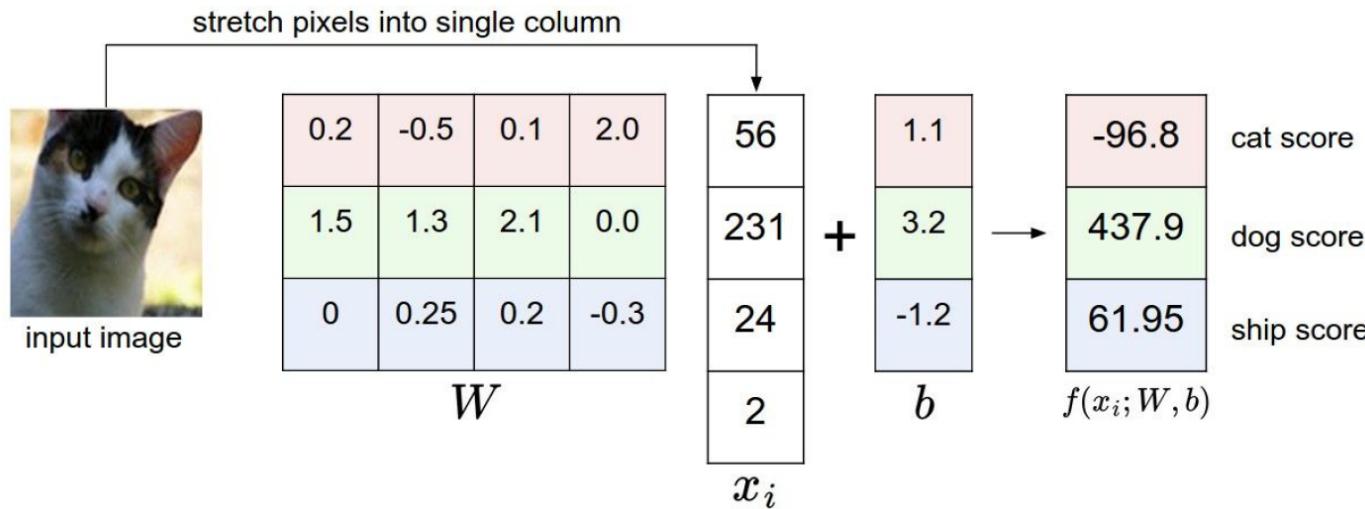
The resulting scores (강아지 점수) are:

0.3
0.7

선형분류기의 뉴럴 네트워크 표현



Example with an image with 4 pixels, and 3 classes (cat/dog/ship)



$$f(x, W) = \begin{matrix} \boxed{Wx} \\ \begin{matrix} 10 \times 1 & 10 \times 3072 \end{matrix} \end{matrix} \quad \begin{matrix} 3072 \times 1 \\ (+b) \quad 10 \times 1 \end{matrix}$$

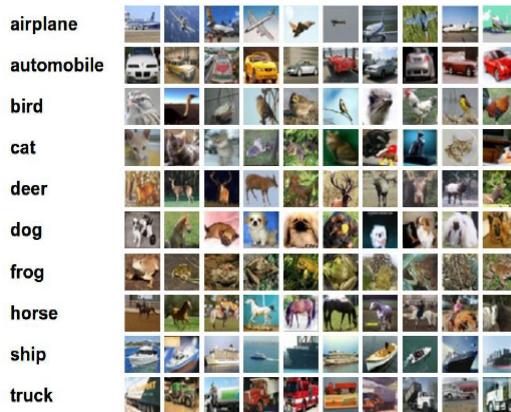
[32x32x3]
array of numbers 0...1

parameters, or “weights”

10 numbers,
indicating class
scores

W를 열어보면 ...

Interpreting a Linear Classifier



$$f(x_i, W, b) = Wx_i + b$$

Example trained weights
of a linear classifier
trained on CIFAR-10:



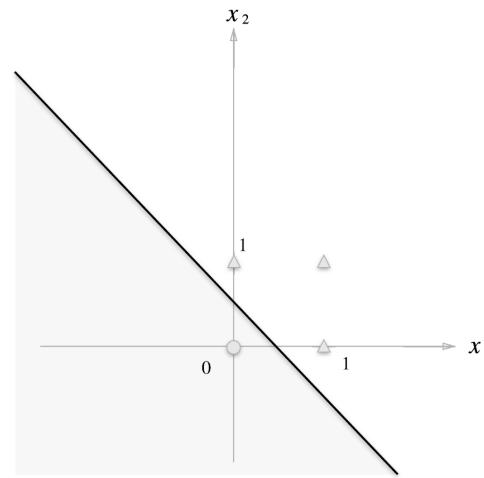
W를 어떤 템플릿으로 볼 수도 있겠군요...

Decision Boundary의 해석

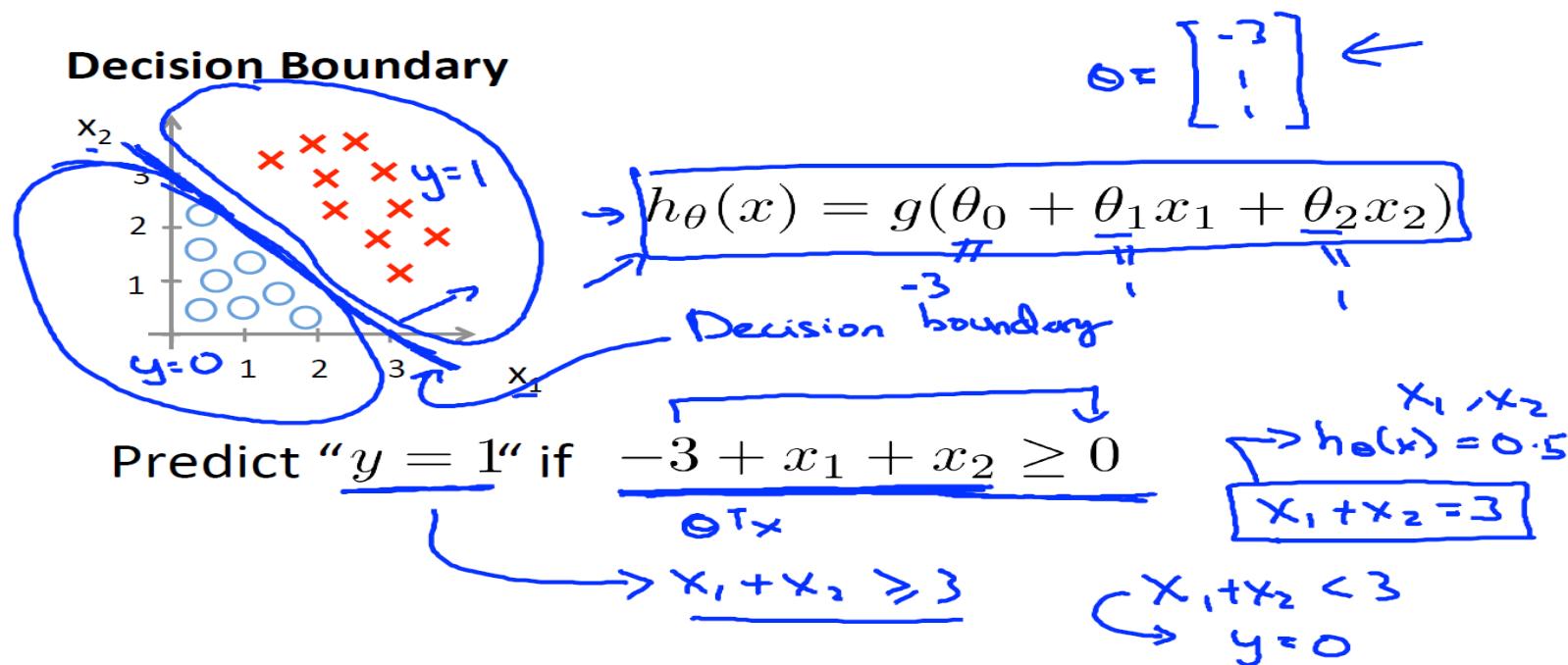
OR 게이트 :

x_1	x_2	y
0	0	0
1	0	1
0	1	1
1	1	1

$$y = \begin{cases} 0 & (-0.5 + x_1 + x_2 \leq 0) \\ 1 & (-0.5 + x_1 + x_2 > 0) \end{cases}$$

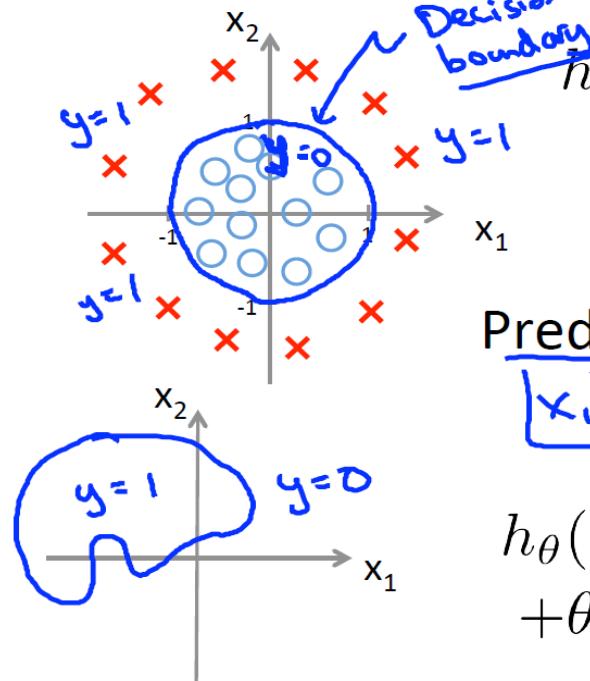


Decision Boundary의 해석



Decision Boundary의 해석

Non-linear decision boundaries



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

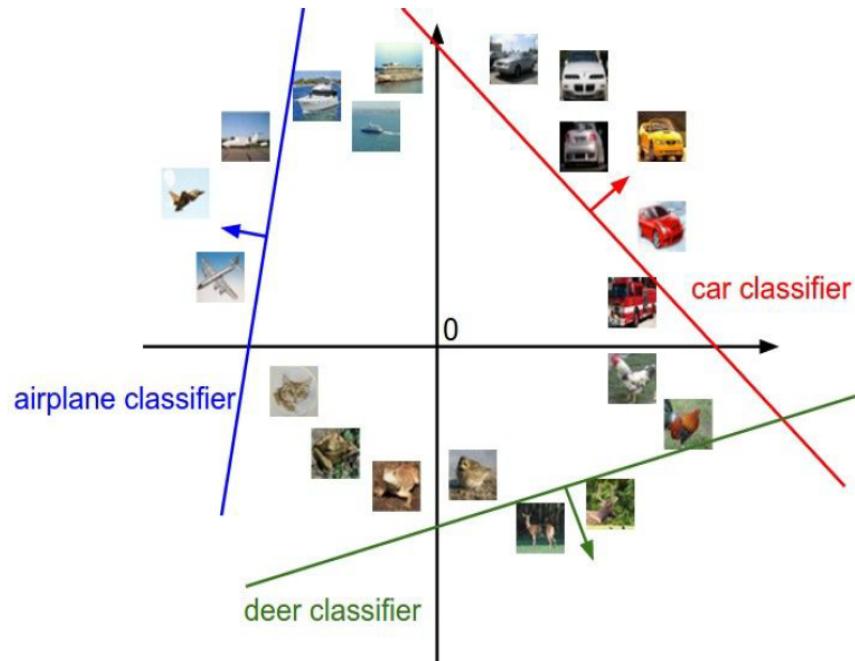
$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Predict " $y = 1$ " if $\frac{-1 + x_1^2 + x_2^2 \geq 0}{x_1^2 + x_2^2 \geq 1}$

$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_1^2 x_2 + \theta_5 x_1^2 x_2^2 + \theta_6 x_1^3 x_2 + \dots)$$

같은 방식으로 해석될 수 있습니다

Interpreting a Linear Classifier



$$f(x_i, W, b) = Wx_i + b$$



[32x32x3]
array of numbers 0...1
(3072 numbers total)

지금까지 : Score function (Linear Classifier) 을 알아봤습니다



$$f(x_i, W, b) = Wx_i + b$$

Example class
scores for 3
images, with a
random W:

airplane	-3.45	-0.51	3.42
automobile	-8.87	6.04	4.64
bird	0.09	5.31	2.65
cat	2.9	-4.22	5.1
deer	4.48	-4.19	2.64
dog	8.02	3.58	5.55
frog	3.78	4.49	-4.34
horse	1.06	-4.37	-1.5
ship	-0.36	-2.09	-4.79
truck	-0.72	-2.93	6.14

Coming up:

- Loss function
- Optimization

$$f(x, W) = Wx$$

- Loss Function : 현재 내가 구성한 W 가 좋은지 안좋은지 측정할 수 있게 해줍니다 (좋으면 0)
- Optimization : random W 로 시작해서 Loss를 최소로 갖게 W 를 변경해 줍니다

기타 :

**영상인식과 관련하여 우리는
어디까지 왔는가 ?**

Course Instructors



Fei-Fei Li



Andrej Karpathy

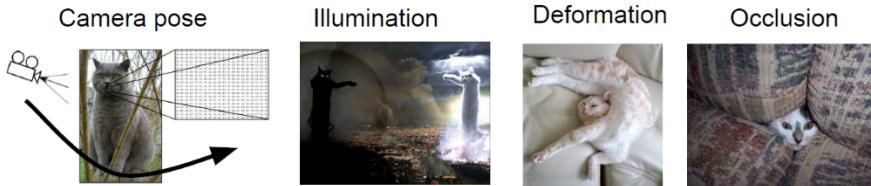


<http://karpathy.github.io/2012/10/22/state-of-computer-vision/>

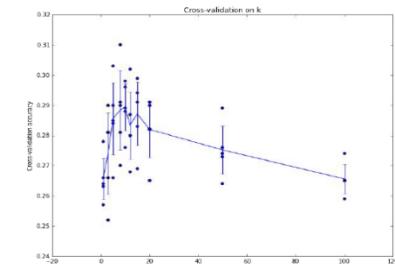
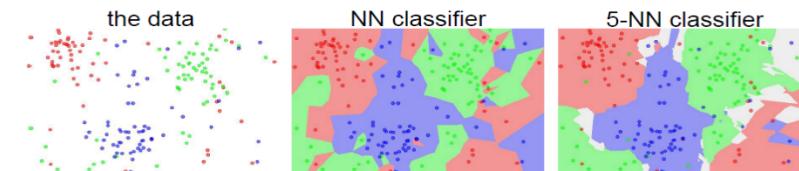
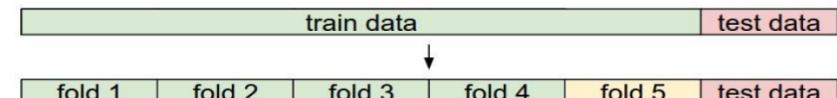
Loss Function

모두의연구소
박은수 Research Director

돌아보기 ...



Background clutter Intraclass variation



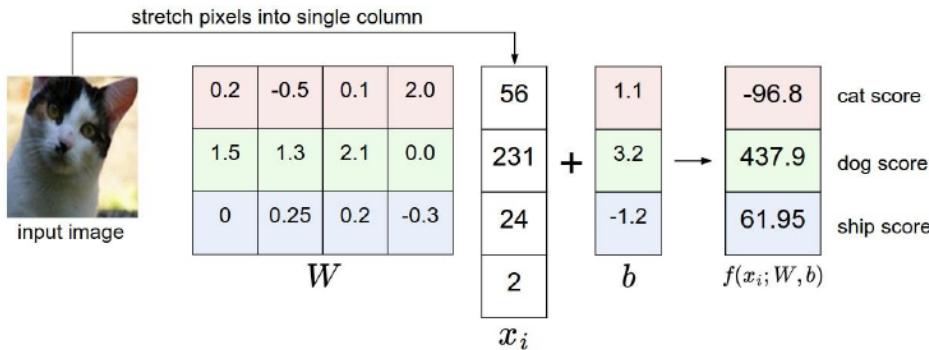
돌아보기 ...



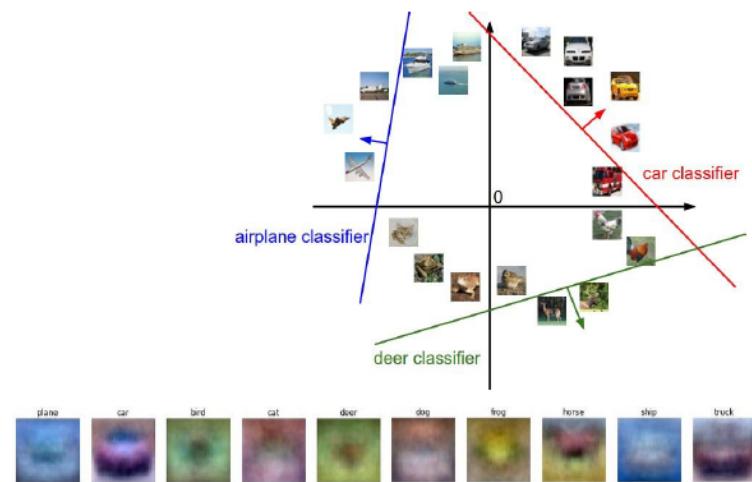
[32x32x3]

array of numbers 0...1
(3072 numbers total)

image parameters
 $f(\mathbf{x}, \mathbf{W})$



10 numbers, indicating class scores

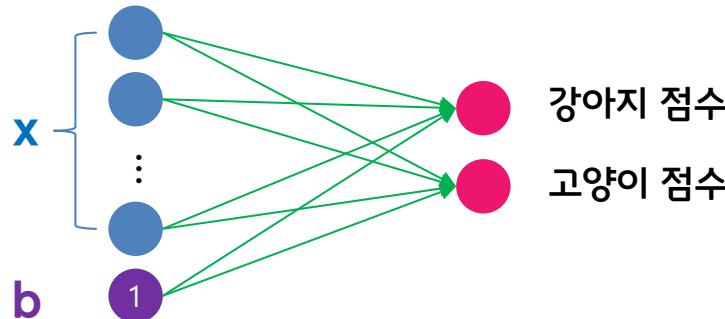


돌아보기 ...

- Score function

$$\begin{array}{c} \text{강아지 파라미터} \\ \text{고양이 파라미터} \end{array} \begin{array}{c} 0.2 \\ 0.1 \\ \dots \\ 0.3 \\ 0.7 \end{array} \begin{array}{c} 0.7 \\ 0.8 \\ \dots \\ 0.9 \\ 0.4 \end{array} \times \begin{array}{c} 155 \\ 200 \\ \dots \\ 110 \\ 78 \end{array} + \begin{array}{c} 0.1 \\ 0.2 \end{array} = \begin{array}{c} 0.3 \\ 0.7 \end{array}$$

\mathbf{W} \mathbf{x} \mathbf{b}



선형분류기의
뉴럴 네트워크 표현

- 분류기의 구성
 - Score function
 - **Loss function**
 - Optimization

Loss Function

- 측정된 Score가 얼마나 안좋은지를 나타내는 하는 함수
- Optimization은 이 측정된 Loss function의 값을 줄이는 방법

이제 해야 할 것은

Score



airplane	-3.45	-0.51	3.42
automobile	-8.87	6.04	4.64
bird	0.09	5.31	2.65
cat	2.9	-4.22	5.1
deer	4.48	-4.19	2.64
dog	8.02	3.58	5.55
frog	3.78	4.49	-4.34
horse	1.06	-4.37	-1.5
ship	-0.36	-2.09	-4.79
truck	-0.72	-2.93	6.14

- 계산된 Score가 얼마나
후진지 알려줄 함수를
만들어야 합니다
→ Loss function

- Loss function을 최소
화 할 수 있는 방법을
만들어야 합니다
(Optimization)

첫번째로 살펴 볼 Loss function은
SVM Loss입니다

$$L_i = \sum_{j \neq y_i} \max(0, f(x_i, W)_j - f(x_i, W)_{y_i} + \Delta)$$

Suppose: 3 training examples, 3 classes.

With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1

Suppose: 3 training examples, 3 classes.

With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1

A **loss function** tells how good our current classifier is

Given a dataset of examples

$$\{(x_i, y_i)\}_{i=1}^N$$

Where x_i is image and
 y_i is (integer) label

Loss over the dataset is a sum of loss over examples:

$$L = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i)$$

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1

Multiclass SVM loss:

Given an example (x_i, y_i) where x_i is the image and where y_i is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$
$$= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1

Multiclass SVM loss:

Given an example (x_i, y_i) where x_i is the image and where y_i is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9		

Multiclass SVM loss:

Given an example (x_i, y_i) where x_i is the image and where y_i is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$\begin{aligned}
 L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \\
 &= \max(0, 5.1 - 3.2 + 1) \\
 &\quad + \max(0, -1.7 - 3.2 + 1) \\
 &= \max(0, 2.9) + \max(0, -3.9) \\
 &= 2.9 + 0 \\
 &= 2.9
 \end{aligned}$$

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9	0	

Multiclass SVM loss:

Given an example (x_i, y_i) where x_i is the image and where y_i is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$\begin{aligned}
 L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \\
 &= \max(0, 1.3 - 4.9 + 1) \\
 &\quad + \max(0, 2.0 - 4.9 + 1) \\
 &= \max(0, -2.6) + \max(0, -1.9) \\
 &= 0 + 0 \\
 &= 0
 \end{aligned}$$

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9	0	12.9

Multiclass SVM loss:

Given an example (x_i, y_i) where x_i is the image and where y_i is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$\begin{aligned}L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \\&= \max(0, 2.2 - (-3.1) + 1) \\&\quad + \max(0, 2.5 - (-3.1) + 1) \\&= \max(0, 6.3) + \max(0, 6.6) \\&= 6.3 + 6.6 \\&= 12.9\end{aligned}$$

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9	0	12.9

Multiclass SVM loss:

Given an example (x_i, y_i) where x_i is the image and where y_i is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Loss over full dataset is average:

$$L = \frac{1}{N} \sum_{i=1}^N L_i$$

$$\begin{aligned} L &= (2.9 + 0 + 12.9)/3 \\ &= \mathbf{5.27} \end{aligned}$$

SVM (Hinge) Loss는 어떤 의미를 갖는가

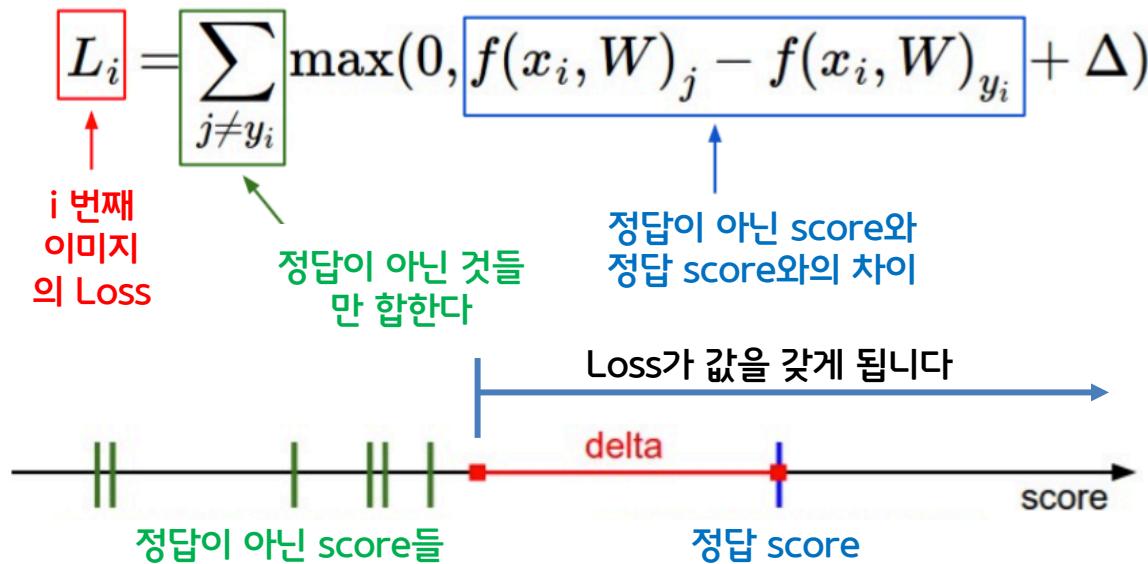
$$L_i = \sum_{j \neq y_i} \max(0, f(x_i, W)_j - f(x_i, W)_{y_i} + \Delta)$$

i 번째 이미지의 Loss
정답이 아닌 것들만 합한다
정답이 아닌 score와 정답 score와의 차이



- 최소한 Loss가 0이 아닌 값을 가지려면 정답이 아닌 녀석들이 저 붉은 delta 위치에는 오거나 정답 Score보다 크면 되겠네요

SVM (Hinge) Loss는 어떤 의미를 갖는가



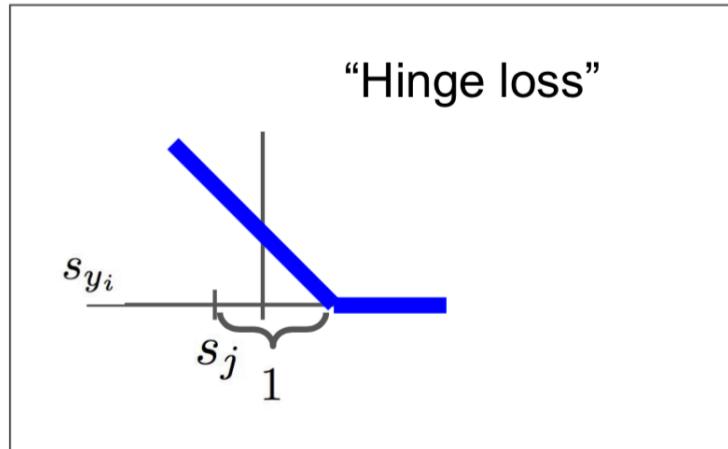
- 최소한 Loss가 0이 아닌 값을 가지려면 정답이 아닌 녀석들이 저 붉은 delta 위치에는 오거나 정답 Score보다 크면 되겠네요

Suppose: 3 training examples, 3 classes.
 With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1

Multiclass SVM loss:



$$\begin{aligned}
 L_i &= \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases} \\
 &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)
 \end{aligned}$$

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9	0	12.9

Multiclass SVM loss:

Given an example (x_i, y_i) where x_i is the image and where y_i is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Q: What happens to loss if car scores change a bit?

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9	0	12.9

Multiclass SVM loss:

Given an example (x_i, y_i) where x_i is the image and where y_i is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Q2: what is the min/max possible loss?

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9	0	12.9

Multiclass SVM loss:

Given an example (x_i, y_i) where x_i is the image and where y_i is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Q3: At initialization W is small so all $s \approx 0$. What is the loss?

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9	0	12.9

Multiclass SVM loss:

Given an example (x_i, y_i) where x_i is the image and where y_i is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Q4: What if the sum was over all classes?
(including $j = y_i$)

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9	0	12.9

Multiclass SVM loss:

Given an example (x_i, y_i) where x_i is the image and where y_i is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Q5: What if we used mean instead of sum?

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9	0	12.9

Multiclass SVM loss:

Given an example (x_i, y_i) where x_i is the image and where y_i is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Q6: What if we used

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)^2$$

Multiclass SVM Loss: Example code

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

```
def L_i_vectorized(x, y, W):
    scores = W.dot(x)
    margins = np.maximum(0, scores - scores[y] + 1)
    margins[y] = 0
    loss_i = np.sum(margins)
    return loss_i
```

$$f(x, W) = Wx \quad \text{Score 함수}$$

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1)$$

N 장의 이미지에 대한 Loss

이제 score를 만드는 법을 배웠고, 계산된 score를 통해
잘못 분류된 값들을 이용해서 Loss를 정의했다.

이제 Loss를 줄이는 W는 찾으면 된다

$$f(x, W) = Wx \quad \text{Score 함수}$$

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1)$$

N 장의 이미지에 대한 Loss

이제 score를 만드는 법을 배웠고, 계산된 score를 통해
잘못 분류된 값들을 이용해서 Loss를 정의했다.

이제 Loss를 줄이는 W 는 찾으면 된다

그러나,

$L = 0$ 이 되는 W 를 찾았다면 이 값이 W 가 고
유한 값인가? 여러개가 될 수도 있지 않나?

$L = 0$ 인 W 값은 무수하군요

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9		

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Before:

$$\begin{aligned} &= \max(0, 1.3 - 4.9 + 1) \\ &\quad + \max(0, 2.0 - 4.9 + 1) \\ &= \max(0, -2.6) + \max(0, -1.9) \\ &= 0 + 0 \\ &= 0 \end{aligned}$$

With W twice as large:

$$\begin{aligned} &= \max(0, 2.6 - 9.8 + 1) \\ &\quad + \max(0, 4.0 - 9.8 + 1) \\ &= \max(0, -6.2) + \max(0, -4.8) \\ &= 0 + 0 \\ &= 0 \end{aligned}$$

Weight Regularization

그럼 값에 제한을 둡시다

너무 멀리서 찾지마 작은데서 찾아...

너무 커지지마 ...

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \lambda R(W)$$

λ = regularization strength
(hyperparameter)

In common use:

L2 regularization

$$R(W) = \sum_k \sum_l W_{k,l}^2$$

L1 regularization

$$R(W) = \sum_k \sum_l |W_{k,l}|$$

Elastic net (L1 + L2)

$$R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$$

Max norm regularization (might see later)

Dropout (will see later)

Weight Regularization



L2 regularization: motivation

$$R(W) = \sum_k \sum_l W_{k,l}^2$$

$$x = [1, 1, 1, 1]$$

$$w_1 = [1, 0, 0, 0]$$

$$w_2 = [0.25, 0.25, 0.25, 0.25]$$

- Same Score
- But L2 penalty?

- $w_1 : 1$
- $w_2 : 0.25$

Less
Loss

$$w_1^T x = w_2^T x = 1$$

L2 regularization: small and widely distributed W
→ **Use all features with small weights**

그런데 저 1은 어디서 나온건가요?

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + \boxed{1})$$

$$L = \frac{1}{N} \sum_i \sum_{j \neq y_i} \left[\max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + \boxed{\Delta}) \right] + \lambda R(W)$$

1을 일단 델타
로 두고 버그
를 쓸어 담아
봅시다



$$L = \frac{1}{N} \sum_i \sum_{j \neq y_i} \left[\max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + \Delta) \right] + \lambda R(W)$$

Do we have to cross-validate both
 Δ and λ ?

$$L = \frac{1}{N} \sum_i \sum_{j \neq y_i} \left[\max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + \boxed{\Delta}) \right] + \lambda R(W)$$

$$\begin{aligned}
 & \max(0, \delta f + \Delta) \\
 \rightarrow & \max(0, \delta f + \Delta_2) && \text{if } \Delta \rightarrow \Delta_2 \\
 \rightarrow & \max(0, \frac{\Delta_2}{\Delta} \delta f + \Delta_2) && \text{then scale the weights by } \frac{\Delta_2}{\Delta} \\
 = & \max(0, \frac{\Delta_2}{\Delta} (\delta f + \Delta)) \\
 = & \frac{\Delta_2}{\Delta} \max(0, \delta f + \Delta)
 \end{aligned}$$

Softmax Loss



$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad \text{where} \quad s = f(x_i; W)$$

Want to maximize the log likelihood, or (for a loss function) to minimize the negative log likelihood of the correct class:

강아지 점수 2.3

$$L_i = -\log P(Y = y_i | X = x_i)$$

고양이 점수 -1.2

in summary: $L_i = -\log\left(\frac{e^{sy_i}}{\sum_j e^{s_j}}\right)$

Softmax Loss



$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

강아지 점수	2.3		
--------	-----	--	--

\exp

9.97

고양이 점수	-1.2		
--------	------	--	--

0.3

전부다 양수가 됨

Softmax Loss



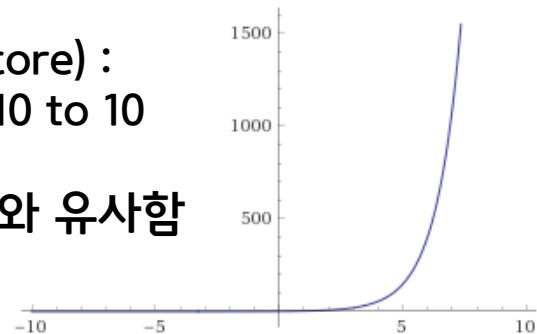
강아지 점수	2.3	\exp	9.97
고양이 점수	-1.2		0.3

전부다 양수가 됨

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

$\exp(score)$:
from -10 to 10

Max함수와 유사함



Softmax Loss



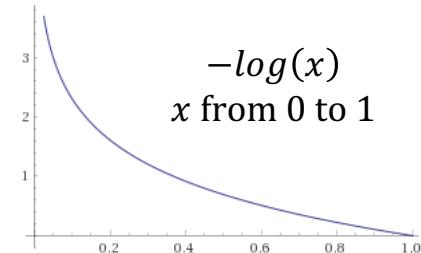
$$L_i = -\log\left(\frac{e^{s y_i}}{\sum_j e^{s j}}\right)$$

강아지 점수	2.3	exp	9.97	normalize	0.97
고양이 점수	-1.2		0.3		0.03
전부다 양수가 됨					확률

Softmax Loss



$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$



강아지 점수

2.3

\exp

9.97

normalize

0.97

고양이 점수

-1.2

0.3

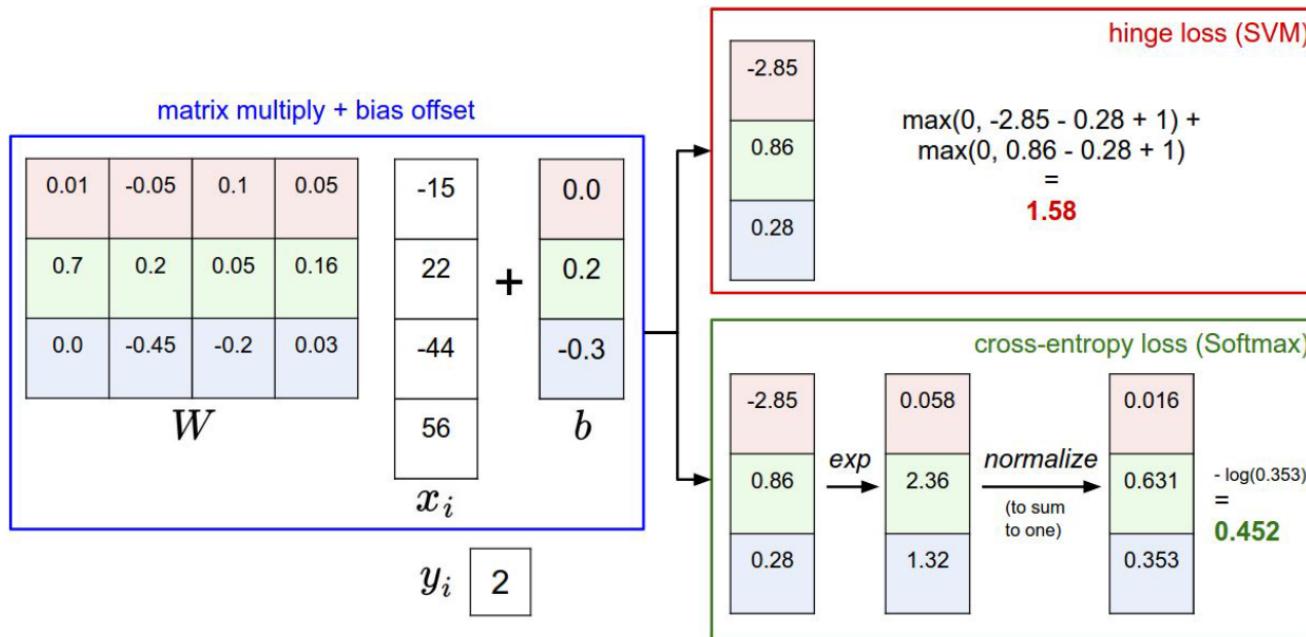
전부다 양수가 됨

0.03

확률

$$L_i = -\log(0.03) = 3.5$$

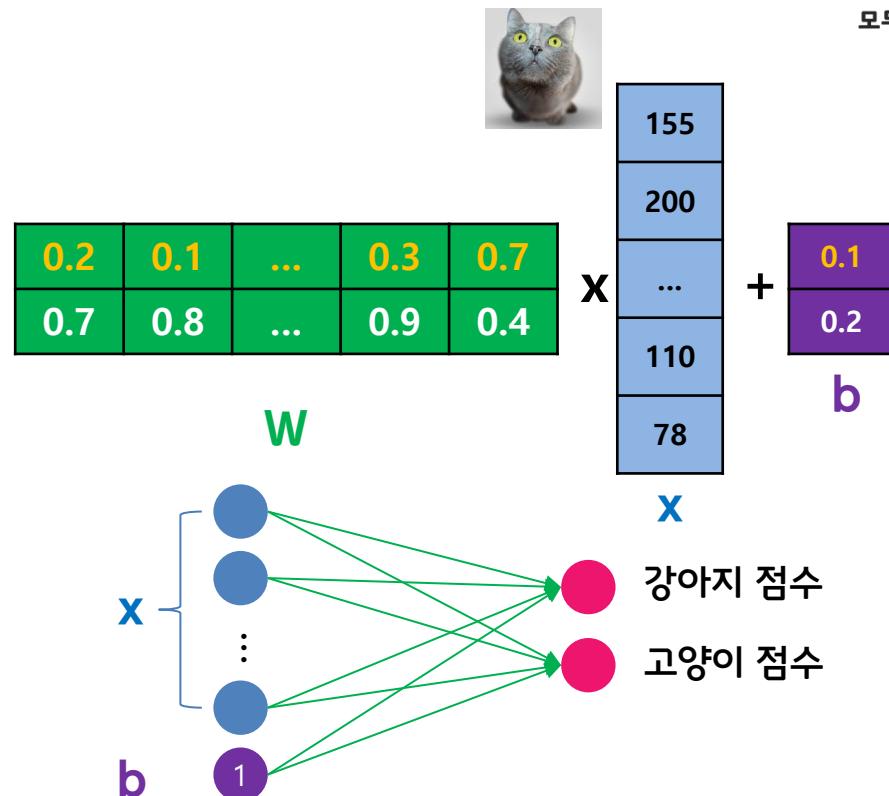
Loss Function : Softmax Loss



강아지와 고양이 분류해보기

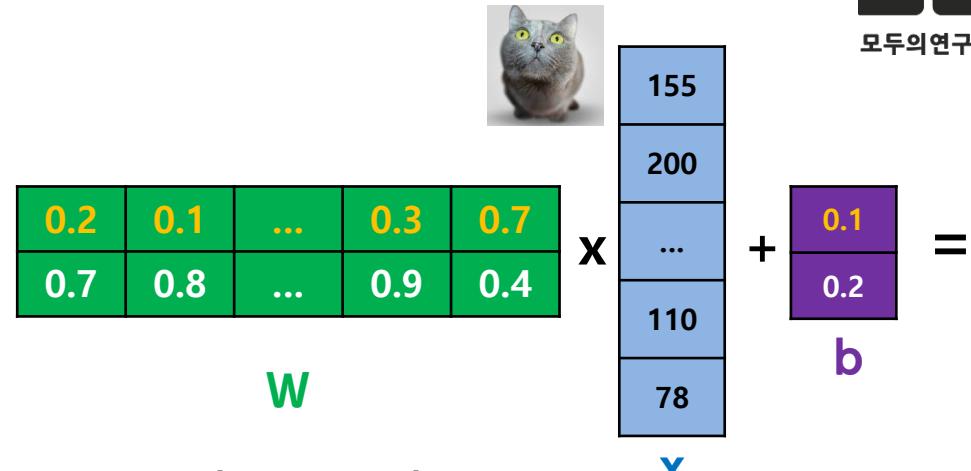
- 분류기의 구성
 - Score function
 - Loss function

고양이가 입력이면 고양이
점수가 높아야 함

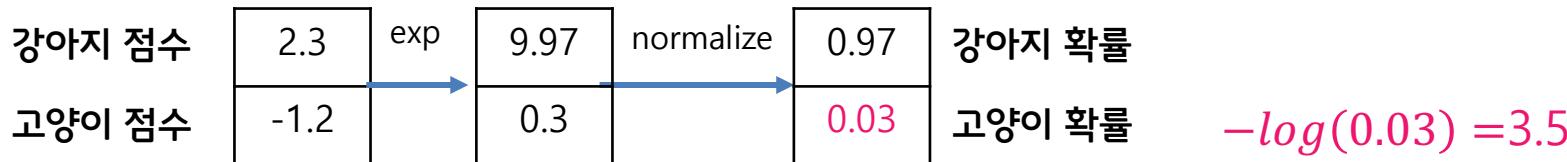


강아지와 고양이 분류해보기

- 분류기의 구성
 - Score function
 - Loss function



Cross-entropy loss (Softmax)



현재의 분류기는 3.5만큼 안 좋음. 이 loss 값을 줄이는게 목표

Loss Function : Softmax Loss



Softmax vs. SVM

- Interpreting the probabilities from the Softmax

$$[1, -2, 0] \rightarrow [e^1, e^{-2}, e^0] = [2.71, 0.14, 1] \rightarrow [0.7, 0.04, 0.26]$$

suppose the weights W were only half as large:

$$[0.5, -1, 0] \rightarrow [e^{0.5}, e^{-1}, e^0] = [1.65, 0.37, 1] \rightarrow [0.55, 0.12, 0.33]$$

Loss Function : Softmax Loss

Softmax vs. SVM

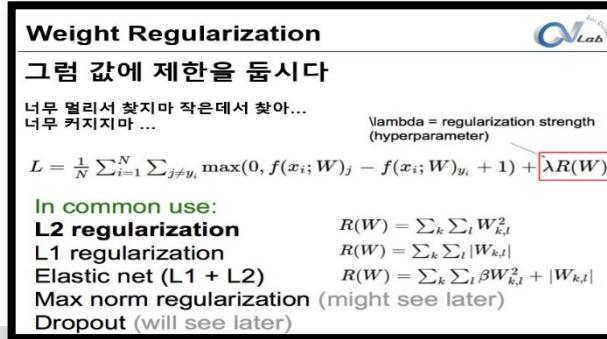
- Interpreting the probabilities from the Softmax

$$[1, -2, 0] \rightarrow [e^1, e^{-2}, e^0] = [2.71, 0.14, 1] \rightarrow [0.7, 0.04, 0.26]$$

suppose the weights W were only half as large:

$$[0.5, -1, 0] \rightarrow [e^{0.5}, e^{-1}, e^0] = [1.65, 0.37, 1] \rightarrow [0.55, 0.12, 0.33]$$

질문 : 만약 Softmax에서 Regularization 파라미터가 너무 너무 크면 어떻게 될 것 같나요?



Weight Regularization

그럼 값에 제한을 둘시다

너무 멀리서 찾지마 작은데서 찾아...
너무 커지지마 ...

λ ambda = regularization strength
(hyperparameter)

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \boxed{\lambda R(W)}$$

In common use:
L2 regularization $R(W) = \sum_k \sum_l W_{k,l}^2$
L1 regularization $R(W) = \sum_k \sum_l |W_{k,l}|$
Elastic net (L1 + L2) $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$
Max norm regularization (might see later)
Dropout (will see later)

Softmax vs. SVM

Softmax vs. SVM

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right) \quad L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

assume scores:

[10, -2, 3]

[10, 9, 9]

[10, -100, -100]

and $y_i = 0$

질문 : Score값에 작은 노이즈가 발생했습니다. 그래서 각 Score값에 작은 변화들이 생겼습니다. 이때 Softmax 와 SVM Loss는 각각 어떨까요?

요약

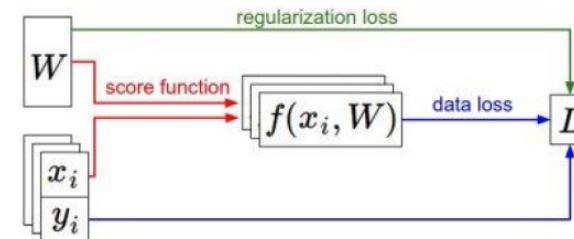
- 지금까지 W 라는 파라미터를 이용한 Parametric 접근의 선형 분류를 살펴보았습니다
- Score 함수를 선형으로 정의해서 살펴보았습니다
- Loss 함수를 살펴보았습니다 (SVM / Softmax)

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right) \quad \text{Softmax}$$

SVM

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$L = \frac{1}{N} \sum_{i=1}^N L_i + R(W) \quad \text{Full loss}$$



이제 Loss를 최소로 하는 W 를 찾는 방법만 남았습니다 :

요약



- 분류기의 구성
 - Score function
 - Loss function
 - Optimization

이제 Loss를 최소로 하는 W를 찾는 방
법만 남았습니다 :
다음 시간에



박은수 책임 연구원

E-mail : es.park@modulabs.co.kr