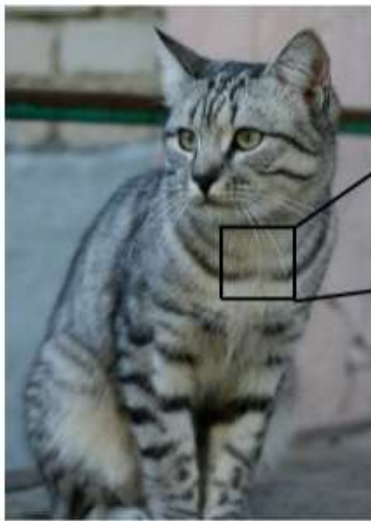


# Lecture 2 Image Classification pipeline

## Classification's Challenges

### The Problem: Semantic Gap

### The Problem: Semantic Gap



This image by NIKSA is licensed under CC-BY 2.0

1	100	112	100	111	104	99	106	99	96	103	112	110	104	97	93	87
2	91	88	102	106	104	79	98	103	95	105	123	116	110	105	94	95
3	76	85	89	105	120	105	87	96	95	99	115	112	106	103	99	85
4	98	81	81	93	120	131	127	100	95	98	107	98	98	93	101	94
5	106	91	81	64	69	85	88	85	101	107	109	98	79	84	96	95
6	114	100	85	95	95	69	64	94	64	87	112	120	98	74	84	91
7	123	127	147	103	65	81	88	65	52	64	74	84	102	93	85	82
8	128	127	144	140	109	95	86	78	62	65	63	68	73	86	101	
9	125	133	148	137	119	121	117	94	65	79	89	65	54	64	72	88
10	127	125	131	147	133	127	126	131	111	86	89	75	61	64	72	84
11	115	114	109	123	158	148	131	118	113	109	100	92	74	65	72	78
12	88	93	90	97	100	147	131	118	113	114	113	100	106	95	77	80
13	63	77	86	81	77	79	102	123	117	115	117	125	125	130	115	87
14	62	65	82	89	78	71	88	101	124	126	119	101	107	114	131	119
15	63	65	75	88	89	71	82	91	120	130	125	105	81	99	110	118
16	87	65	71	87	100	95	69	45	76	130	125	107	92	94	105	112
17	110	97	82	86	117	123	116	66	41	51	95	93	89	95	102	107
18	104	146	112	88	82	120	124	104	76	48	45	66	88	101	102	100
19	117	170	117	120	93	86	114	132	112	87	65	55	78	82	99	94
20	110	120	134	161	130	100	109	118	121	114	114	87	60	53	89	86
21	120	112	96	117	150	144	120	115	104	107	102	93	87	81	72	79
22	123	107	95	88	83	112	153	149	122	109	104	75	88	107	112	89
23	122	121	100	88	82	86	94	117	145	148	153	102	58	78	93	107
24	122	164	148	103	71	56	78	83	93	103	119	110	102	61	60	84

What the computer sees

An image is just a big grid of numbers between [0, 255]:

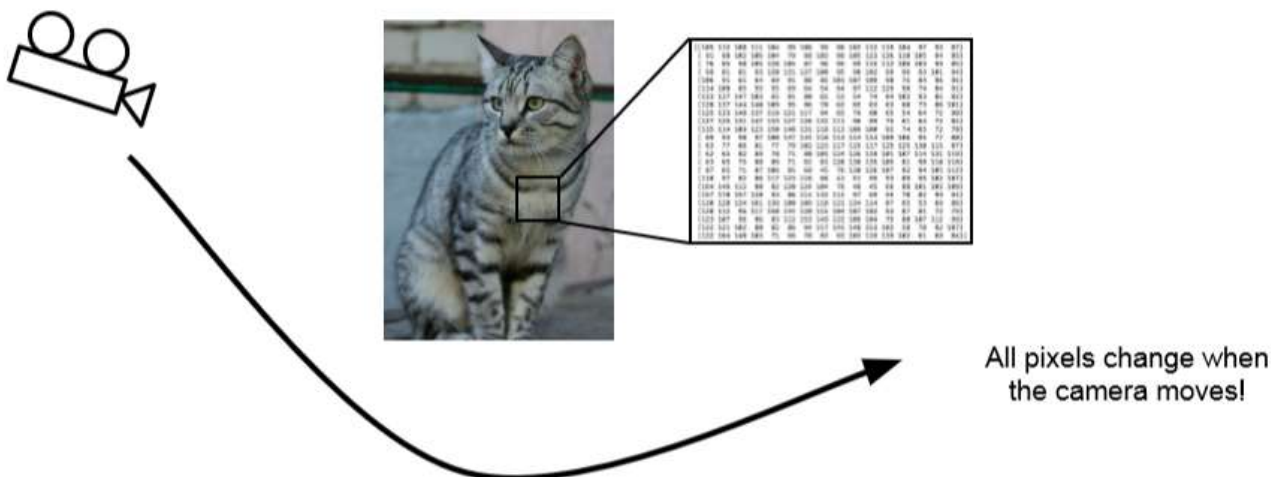
e.g. 800 x 600 x 3  
(3 channels RGB)

每个图片都是由很多像素点组成的。就像主讲人说的，图片对于计算机来说就是a **gigantic grid of numbers**

所以Semantic Gap的意义就是我们给图片标的语义标签与计算机实际看到的像素值之间有巨大的差距。

### Viewpoint variation

## Challenges: Viewpoint variation



同一物体的不同角度得到的像素值序列是不一样的，但是是表示同样的物体，所以我们的算法就必须对视角这样的问题有较强的鲁棒性robust

## Illumination



不同的光照条件也会影响，所以我们的算法也需要对这样的情况就要较强的鲁棒性。

## Deformation(变形)



Our algorithms should be robust to these different kinds of transforms.

## Occlusion(遮挡)





This image is CC0.1.0 public domain



This image is CC0.1.0 public domain



This image by jorison is licensed under CC-BY 2.0

## Background Clutter



This image is CC0.1.0 public domain



This image is CC0.1.0 public domain

猫身上的纹理或者猫本身的颜色与背景非常相似，那么我们的算法也需要对这种情况作出正确的判断。

## Intraclass variation



This image is CC0.1.0 public domain

类内的差距，比如说同样都是猫，但是不同的猫具有不一样的颜色或者纹理什么的。

## 处理过程

---

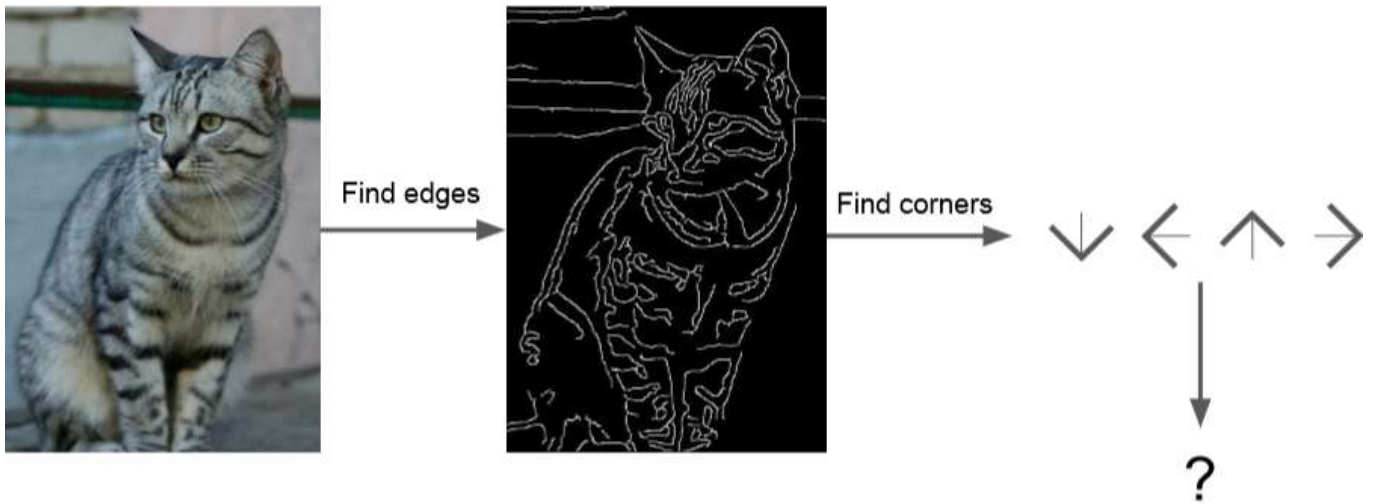
### An image classifier

图像分类的时候可以做一个这样的API接口用于接收图像并返回图像的标签。

同样，并没有现成的函数或者算法可以做到直接去classification

### Attempts have been made

Attempts have been made



一些图像预处理的方法，比如就是想图中做的，先用函数算出边缘，然后得到某些线或者点的交汇方式并作出标签，但是这样的方法并不是一个特别好的方法。

### 曼哈顿距离和欧几里得距离

#### L1 distance

**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	add → 456
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	

这里的曼哈顿距离主要在意的是测试图片的每个像素与训练图片的对应位置像素的差距。这是一种简单的最临近分类器。

```
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
```

上图展示了最临近分类器的计算操作，基于的是曼哈顿距离的理论，训练过程主要是记录训练图片每个像素的值，而测试过程是将测试图片的每个像素与训练的model做比较。这样的结果就是，训练的时间非常短，而测试的时间非常长，这样其实在应用中是一种很难推广的方法。

## 两者的理论定义



欧几里得距离，欧氏距离，也就是我们熟知的距离，可扩展至m维

2维:  $d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$

3维:  $d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2 + (z_1 - z_2)^2}$

m维:  $d = \sqrt{\sum (x_i - x_{i2})^2}$

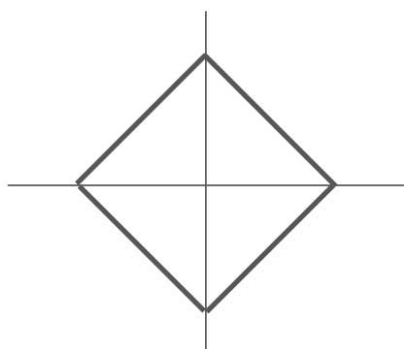
曼哈顿距离，出租车距离，在一个横竖分布的网格上，两点之间的距离即为曼哈顿距离

也就是说只能横着走或者竖着走

$d = |x_1 - x_2| + |y_1 - y_2|$

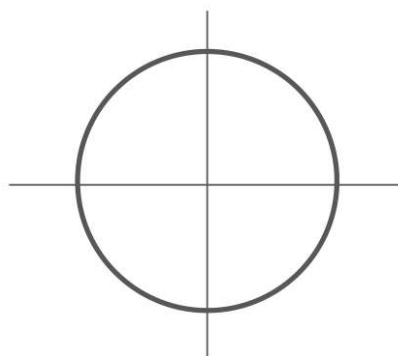
## L1 (Manhattan) distance

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



## L2 (Euclidean) distance

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



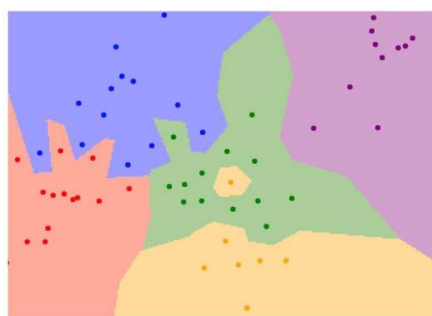
这里主要的问题在于，欧几里得距离对于任何的特征都没有偏向，所以如果我们在做分类的时候并没有去有特定的考虑的话，可以用欧几里得距离。

但是曼哈顿距离就完全不同了，他其实相当于有所偏向的欧几里得距离，也就是说某些特征比另外一些特征更加对距离有限制，所以，如果在分类的过程中对某些特征比较看重的话，是需要用曼哈顿距离的。

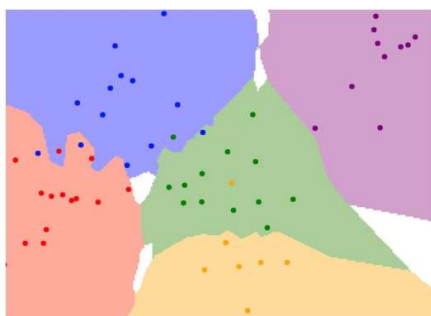
## K最近分类

### K-Nearest Neighbors

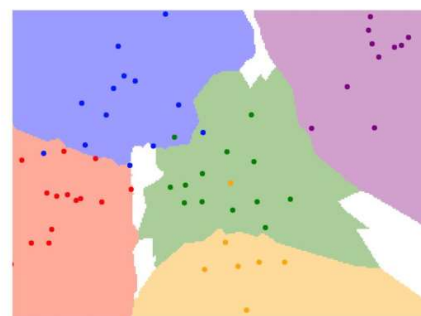
Instead of copying label from nearest neighbor,  
take **majority vote** from K closest points



K = 1



K = 3



K = 5

这里用的方法是**majority vote**最大值投票的方法。

## 超参数——Hyperparameters

### Hyperparameters

What is the best value of **k** to use?

What is the best **distance** to use?

These are **hyperparameters**: choices about the algorithm that we set rather than learn

选择超参数的三种方式:

### Setting Hyperparameters

**Idea #1:** Choose hyperparameters that work best on the data

**BAD:**  $K = 1$  always works perfectly on training data



**Idea #2:** Split data into **train** and **test**, choose hyperparameters that work best on test data

**BAD:** No idea how algorithm will perform on new data



**Idea #3:** Split data into **train**, **val**, and **test**; choose hyperparameters on val and evaluate on test

**Better!**



显然第一和二种是不合理的:

- 对于第一种来说, 在训练阶段去设计超参数并不会确认你的超参数适应新的数据集, 更不用说你连测试都没有测试过。
- 对于第二种来说, 显然会稍微好一点, 但是仍存在着非常多的弊端, 首先, 你讲一部分的训练集转变为测试集, 虽然在一定程度上加强了超参的可靠性, 但是不得不承认这样与第一种方法其实没有什么区别, 你的超参在test上表现的再好也不能表明你一定会在新的测试集上表现的很出色。
- 对于第三种方式, 显然是最好的, 首先, 我们将数据集分为三个部分, 训练之后首先会在validation上做验证, 以保证超参的可靠性, 然后在测试集中做测试, 这样加一步会提高超参的可靠性。

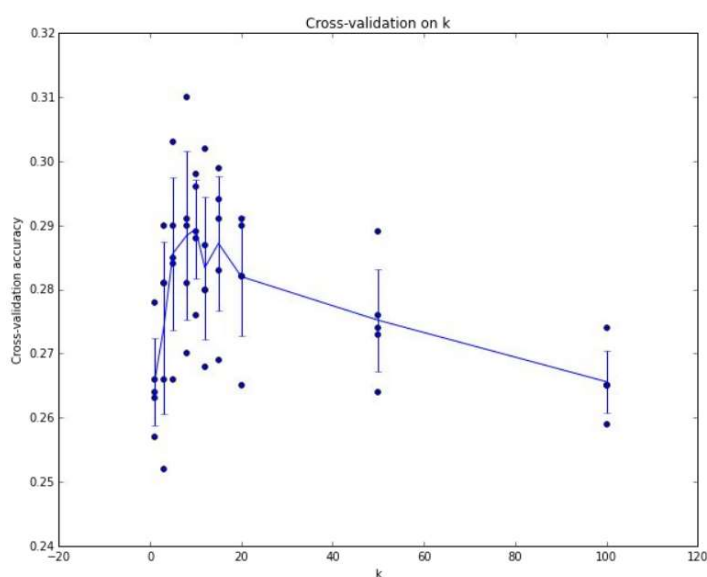
### 选择超参数的第四种方式——交叉验证

**Idea #4: Cross-Validation:** Split data into **folds**, try each fold as validation and average the results

fold 1	fold 2	fold 3	fold 4	fold 5	test
fold 1	fold 2	fold 3	fold 4	fold 5	test
fold 1	fold 2	fold 3	fold 4	fold 5	test

交叉验证也就是将同一数据集分成多个fold，每次将不同的fold作为验证集做交叉验证，这样虽然造成大量的计算量，但是相对来说是非常棒的方法。

## Setting Hyperparameters



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  $k \sim 7$  works best for this data)

这里展示的是交叉验证的结果，可以看出，就像图中说的， $k$ 大概是7的时候是最好的结果。每次我们用数据集做五次交叉验证，得到平均的准确度。

## K-Nearest Neighbor 不常使用的原因

**第一，在测试过程中花费时间太长**

**第二，像素点的距离值运算并不一定informatic——不一定有足够的信息。**



Image is  
in domain

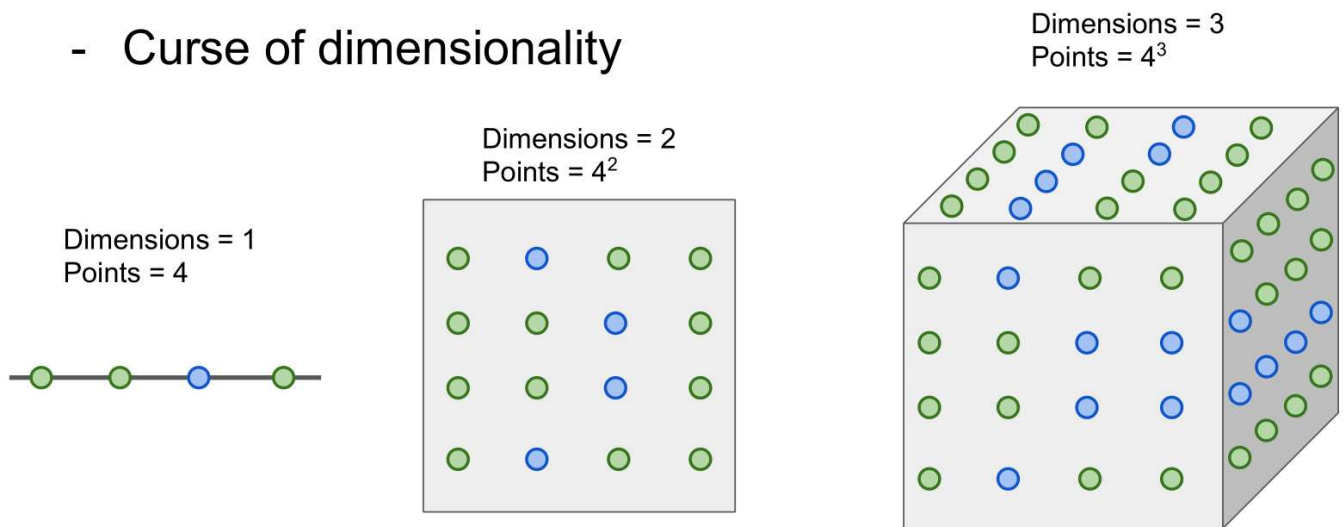
(all 3 images have same L2 distance to the one on the left)



像上面的四张图，是完全不相同的四张图，但是在算L2距离的时候，也就是欧几里得距离的时候，因为不考虑坐标的问题，所以其实这四张图的L2距离是完全一样的，这样可以看出，L2距离并不是合理的图像分类的算法。

第三，K近邻算法随着维度的增加，需要的数据量是平方次增加。

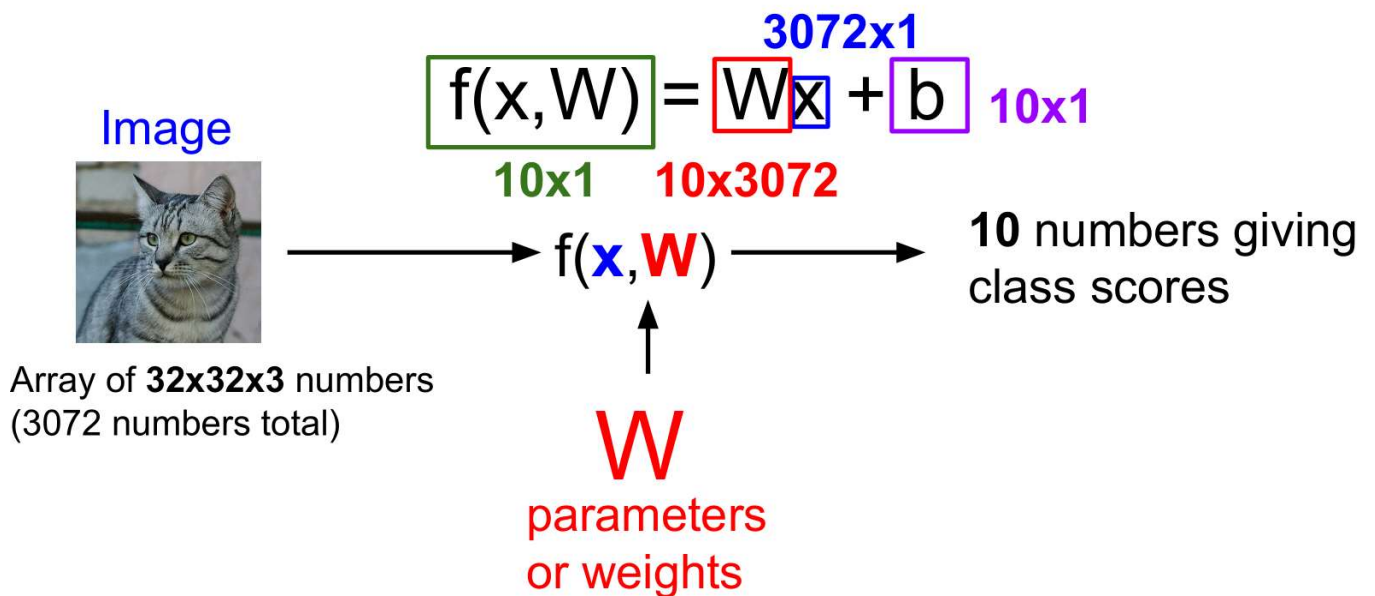
### - Curse of dimensionality



## Linear Classification——线性分类

### 新的test的方式——权重model

#### Parametric Approach: Linear Classifier



也就是说，我们每次训练只需要训练出啦权重，然后训练部分就废掉了，不用管了，在测试的过程中只需要权重参数就可以解决所有问题。

另，上图猫的三个通道指的是RGB三通道。

