

Elena Georgieva - 4/23/2021

Homework - Categorization and Model Comparison Part A (60/100 points)

by *Todd Gureckis* and *Brenden Lake*

Computational Cognitive Modeling

NYU class webpage: <https://brendenlake.github.io/CCM-site/> (<https://brendenlake.github.io/CCM-site/>)

email to course instructors: instructors-ccm-spring2021@nyuccl.org (<mailto:instructors-ccm-spring2021@nyuccl.org>)

This homework is due before midnight on Apr 19, 2021.

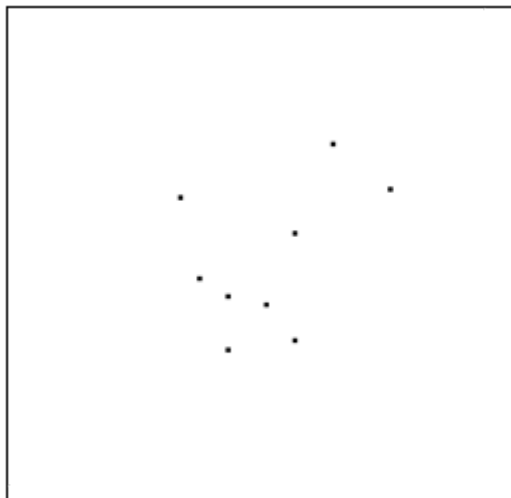
```
In [1]: import string as str
import os
import numpy as np
import seaborn as sns
import pandas as pd
import math
from random import random, randint, shuffle, uniform
from scipy.optimize import fmin, brute
```

Background and Theory

In this homework we explore the cognitive mechanisms that support unsupervised pattern categorization in humans. In addition, we use this as an example of testing and comparing between models.

A simple (classic) unsupervised categorization experiment

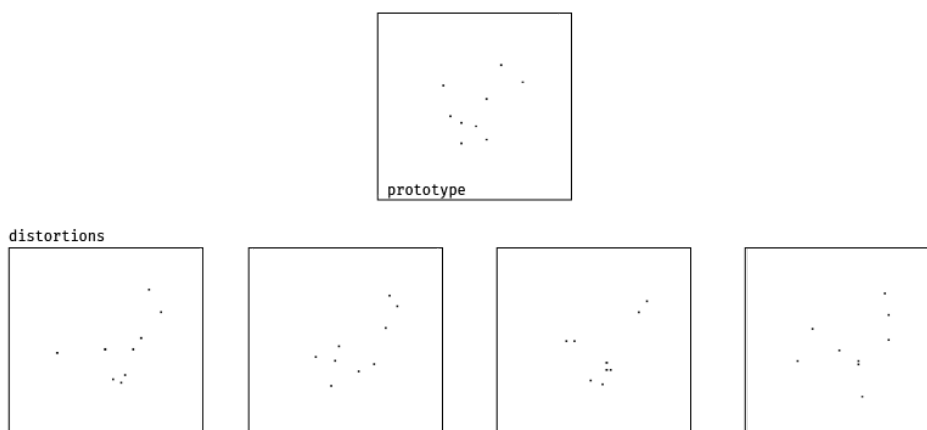
In Posner and Keele (1968) report a now classic categorization experiment with humans. In the task participants viewed visual stimuli that are clouds of points (known as dot patterns) similar to a scatter plot of data on a graph. An examples of the stimuli is shown here:



Experiment Design

The experiment was divided into a training and test phase. During the training phase, for each subject a single random dot pattern was generated and considered to be the underlying "prototype" structure. A prototype is like a common template or reference pattern. The key is that participants never get to see the "prototype" pattern directly during training. Instead they see what are known as "distortions" of the prototype. A distortion of a pattern is made by adding random spatial noise to each point in a pattern to kind of "wiggle" the points away from their original position.

For example, here is a random prototype (top) and a bunch of random distortions of the prototype made by adding or subtracting small random values from the $\langle x, y \rangle$ value of each point in the pattern.



Posner and Keele created distortions that added more or less random noise. For example, "high" distortions add a lot of randomness to the underlying template pattern whereas "low" distortions add only a little bit of noise.

Training Phase

In the training phase of the experiment subjects view 10 training examples one at a time which are "high" distortions of a randomly generated prototype. The instructions are that subjects should look at these patterns, and that they come from a single category similar to if you viewed a series of pictures of dogs they would all come from the category `dog`. Subjects were try to figure out the pattern that related the different images to one another. Try it for yourself by looking at each of the "distortions" patterns above one by one and trying to detected the common structure.

Test Phase

During the test phase, participants view a series of dot patterns one at a time and have to judge: **Does the given pattern come from the same general category or family you studied earlier or is it a new pattern that is different?** This is an unsupervised categorization task because the subject has to abstract what the common structure is from the given patterns and then use that information to make classification decisions about new patterns.

Unknown to participants the set of test items varied in a specific way with respect the training patterns. In particular, there were five particular types of patterns presented during test.

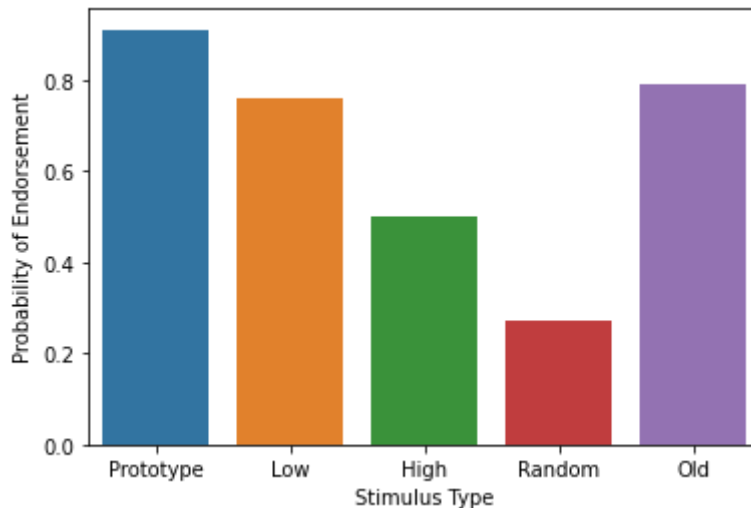
- The first type were "old" patterns which were identical to those presented during the training phase.
- The second type were "random" patterns which were from a complete new randomly generated prototype (thus had nothing to do with the items presented during training).
- The third type were new "high distortions" of the underlying prototype that was used to create the study set. These are thus similar to the "old" items but do not match exactly.
- The fourth type were "low distortions" of the underlying prototype that was used to crete the study set. These are more similar to the prototype pattern than the "high" distortions are.
- Finally the actual prototype used to generate the items during training was presented. This pattern is interesting because the prototype pattern was never seen exactly during training. However, people saw many high distortions of this item during training and given the instructions to detect what the common structure of the training patterns is, they may have learned something about this latent or hidden pattern.

Typical Results:

This graph show example results that are typical for an experiment like this:

```
In [2]: df = pd.DataFrame({"Stimulus Type":["Prototype", 'Low', 'High', 'Random', 'Old']  
sns.barplot(x="Stimulus Type",y="Probability of Endorsement",data=df)
```

```
Out[2]: <matplotlib.axes._subplots.AxesSubplot at 0x10f5b2550>
```



The height of the bars indicates the probability of endorsing a pattern as a member of the category during the test phase (high values mean that at test a subject is more likely to agree "yes this pattern fits with the one I learned"). Notice that the "old" items (the exact patterns studied during learning) are endorsed at a relatively high rate. In contrast the "random" patterns (those coming from a completely different underlying pattern) are endorsed at a very low rate. The low and high distortions are endorsed at intermediate rates (with the low distortions preferred). Interestingly, the prototype pattern is endorsed most strongly even though it was never presented during the study period! It is like during learning people figured out the underlying pattern that generated the stimuli!

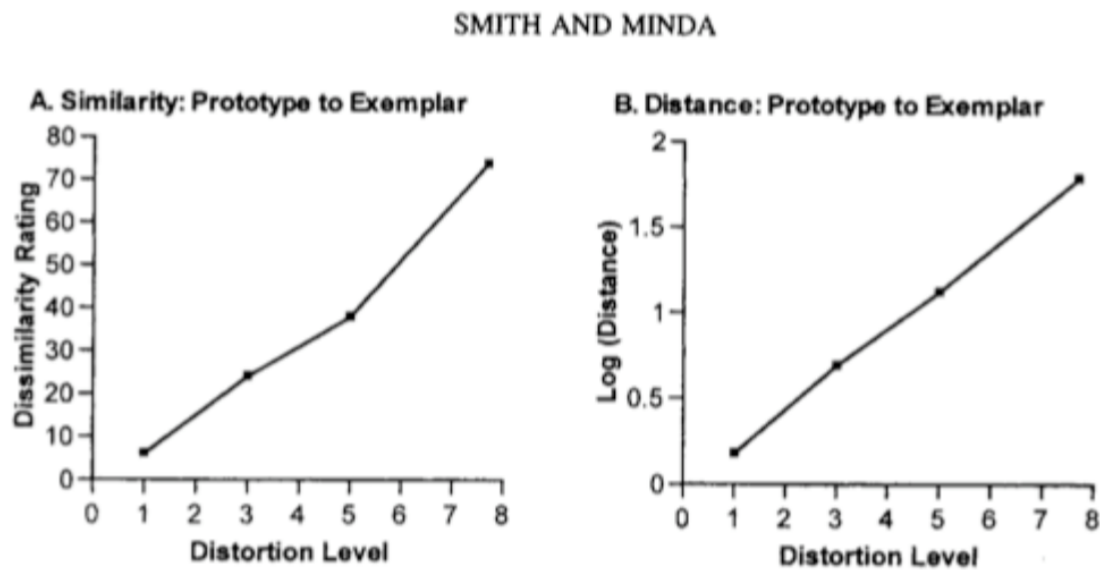
Categorization: Memory for examples or abstractions?

The dot pattern experiments are interesting and have been replicated in various ways perhaps hundreds of times. How do people make these judgments? What information do people store during the study phase that would predict their performance in the test phase? What rule do they use for

combining information from memory in order to make these kind of judgements? We are going to build up a model of categorization in a few simple steps.

Step 1: How are dot patterns represented in the mind?

Our first question concerns how people represent the dot patterns and the similarity between different dot patterns. A variety of work has suggested that the psychological similarity between pairs of dot patterns follows roughly a logarithmic transform of the average euclidean distance between the pairs of points (plus one). This was established by having people view pairs of dot patterns and rate how similar they seem. For example here is a plot from a paper by Smith and Minda showing a strong contgruence between dissimilarity ratings between pairs of stimuli and $\log(\text{distance})$



In light of this lets define the similarity between two dot patterns i and j as s_{ij} and let it equal the following equation:

$$s_{ij} = \log\left(1 + \frac{\sqrt{\sum_d [(i_{d_x} - j_{d_x})^2 + (i_{d_y} - j_{d_y})^2]}}{9}\right)$$

where i_{d_x} is the x position of the d -th dot in pattern i and j_{d_x} is the x position for pattern j (likewise for i_{d_y}). Because it can sometimes be ambiguous which dot aligns which which one in a pattern we choose the dot which are closest in the two patterns to compute this score.

Step 2: What is stored in memory?

The next consideration is what people actually store in memory during the training phase of the experiment. There are of course many alternatives. People could store an "average" of the points

seen so far, or they could store each individual pattern that they have seen, or they could store nothing and try to figure it out at test, or they could store some verbal description of what the shapes "look like", or the shape of the outer edge of the dot-cloud (the "convex null"), etc...

There are, however, two leading theories which have attracted considerable debate in the cognitive science literature: the prototype and exemplar theory.

Exemplar models

Exemplar models are a general class of psychological models related to nearest neighbor algorithms. The most important feature of these models is the idea that people have what appears to be a nearly infinite memory for the past and as a result you can store all past experiences or examples in memory. This seems crazy as we are forgetting things all the time but actually psychology is unclear about if we actually forget things or if we simply lose the ability to retrieve a memory (i.e., more like losing the pointer to the memory rather than decay).

As mentioned in lecture, nearest neighbor classifiers use a similarity function (similar to the ones described above) to retrieve from memory the nearest labeled example and to predict the category membership based on the label for this item. This nearest neighbor algorithm can be relaxed slightly to consider k -nearest neighbors. According to this algorithm you find the k neighbor examples (with $k > 1$) to the current pattern and response based on what the majority of these examples say.

Now we can go a bit further and say that you compute the similarity to all past examples but *weight* their vote according to their similarity. So instead of picking the label of the closest or k -closest examples we compare the current pattern using a global match to all examples in the memory and weight their predictions based on similarity. Pretty neat!

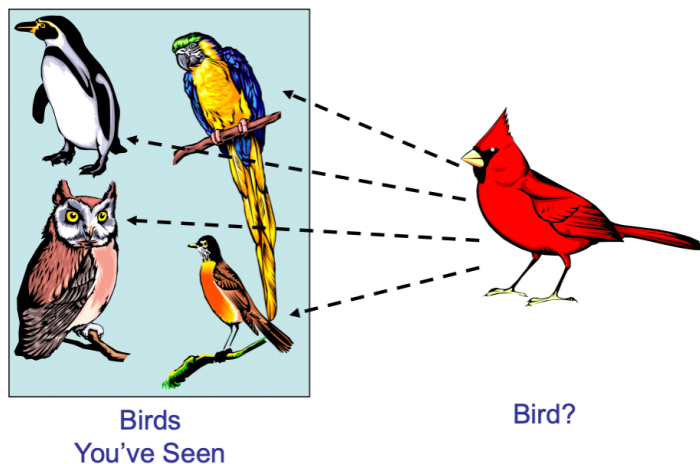
Ok, but how does this help us in the case of **unsupervised** categorization such as in the dot pattern case? Here what we will assume is that we compute this similarity of the to-be-categorized item (the test item) to all the examples stored in memory and compare it to some criterion value. If the sum of the similarity to all the examples falls below this criterion then we assume the pattern is new and doesn't match what we learned. If it is above the criterion we judge the item is a good example of the category.

In the exemplar model we will consider the probability of endorsing an item is going to be determined by the following equation:

$$P(A|i) = \frac{\sum_j e^{-c \cdot s_{ij}}}{\sum_j e^{-c \cdot s_{ij}} + k}$$

where $P(A|i)$ is the probability of endorsing pattern i as a member of the category seen during study. s_{ij} is the similarity between pattern i and pattern j which is an example stored in memory during the study phase. k is the criterion against which the summed similarity is being compared. If

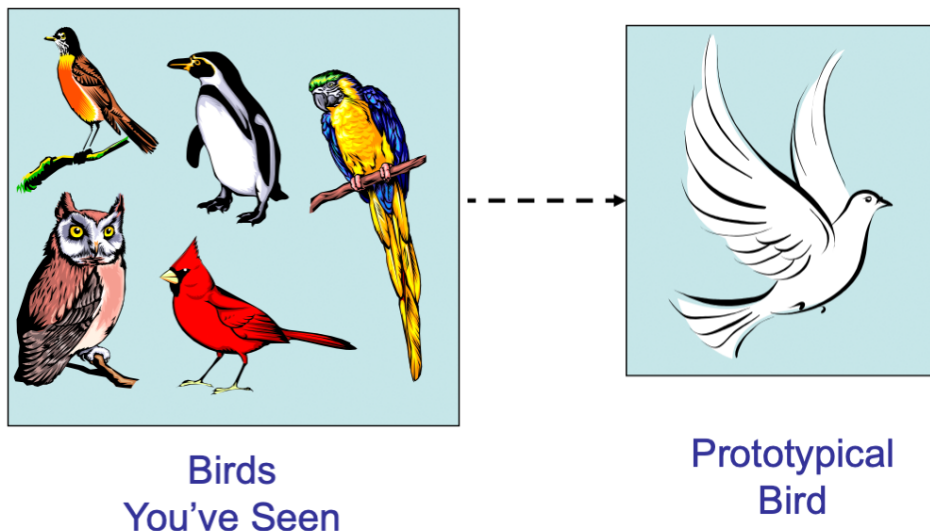
k is zero then you endorse the item as a member of the category all the time irrespective of the similarity and if k gets large you become more and more less likely to endorse the item (i.e., more likely to say no).



The sum is with respect to an exponential sum which has some deeper relation to research on categorization that we do not have time to discuss. However, it is basically the idea that very close matches ($s_{ij} = 0$) are especially strong and things that are less similar count less. You can think of it as the particular weighted nearest neighbor algorithm we think the mind uses. c is a free parameter that controls that weighting function and is often fitted to data.

Prototype models

The prototype model is different than the exemplar model because it assumes that instead of storing each of the training patterns in memory exactly, instead people store a single summary representation. For example, people might store a mentally computed "average" pattern. When you think about how you would perform the task you might think that you kind of compare the training patterns to one another and then compute some summary.

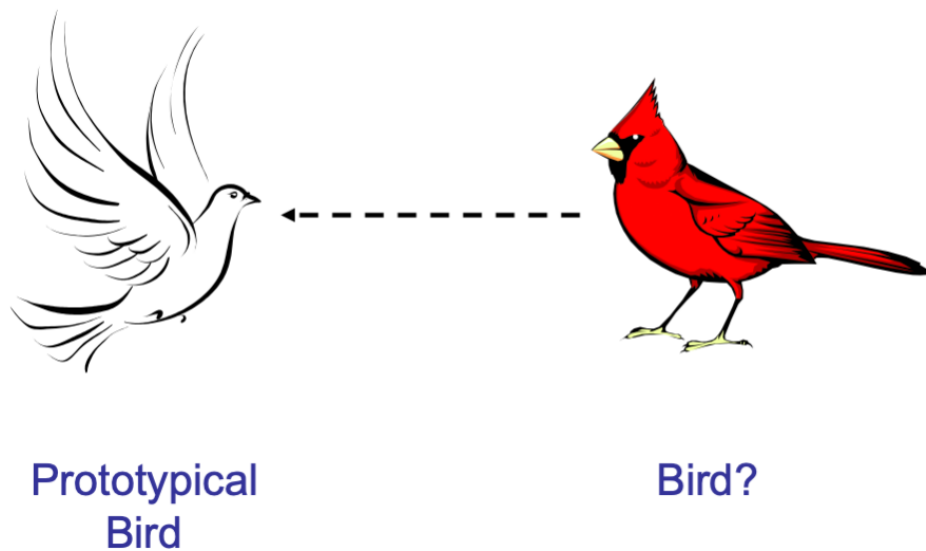


In the case of the dot pattern stimuli one way to do this is to store a special trace in memory called the prototype which is the average of all the patterns seen so far (averaging the $\langle x, y \rangle$ position of each point to find an average dot location).

According to the prototype model the probability of endorsing a test item pattern as a member of the category that was studied during training is:

$$P(A|i) = \frac{e^{-c \cdot s_{ip}}}{e^{-c \cdot s_{ip}} + k}$$

Note that nearly everything about this equation is the same except there is no longer a sum! Instead we simply compute the similarity between the test pattern and this special "prototype" pattern (p) which has been averaged from the training examples.



Parameters

k and c are "free parameters" in both the exemplar and prototype model which are assumed to modulate or alter the core psychological processes. These parameters might vary between subjects and as a function of condition. Thus, in order to assess the ability of the model to account for the data we often "fit" these parameters to our data.

Model Comparison

With these ideas in mind, in this homework you are going to compare the exemplar and prototype model to account for some data from an actual dot pattern categorization task collected with human subjects. The goal is that by doing the homework you would develop some useful code that would let you more or less plug in a model that you might come across in your research, fit it to data, and verify that the fits are good, etc...

Reading in some data

The `data/` folder that comes with this homework contains data from 14 human subjects who participated in a dot pattern classification task. The data describing each subject is in a text file (`.dat`) indexed by subject number (e.g., `1.dat` , `2.dat` , etc...).

The organization of these files is as follows:

The first 44 lines of the file contain a description of the stimulus that the subject saw on a given trial including the x, y coordinate of each dot. The first columns of these 44 lines is the number of the pattern (1–44). The second column is the type of pattern using the following codes:

- 1 = prototype
- 2 = 10 "high distortions" of the prototype that were used as study patterns during learning
- 3 = 10 new "high distortions" of the prototype presented during test
- 4 = 4 "low distortions" of the prototype that were presented during test
- 5 = 20 random items that come from different prototypes that were presented at test

The next 18 values of each row are the coordinates of the dots (with the x, y coordinates in sequence). So `[x1, y1, x2, y2, x3, y3, ...]` .

The following 40 lines of the file show the sequence of items presented during the study phase. This is not all that important for our purposes, but basically the last column is which pattern was displayed (indexed from the patterns just described. Each of 10 "high distortions" were presented four times each during study in a random order.

Finally the remaining lines of the file report the results of the test phase. The first column is the subject number, the second columns is the condition number, the next is the trial number in the experiment, the other columns worth mentioning are the last column (the pattern number from the beginning of the file), the second to last column (the type of stimulus it is according to the codes described above), and the reaction time in milliseconds.

In addition, participants in this experiment were assigned to one of two conditions: a recognition condition and a categorization condition. These conditions differed only in the instructions given to participants at the start of the test phase. In the recognition condition participants were told they would view a series of patterns and they should respond "yes" only if the patterns was **exactly** one they say in the previous study phase. In the categorization condition, participants were asked to respond "yes" only if the pattern belonged to the same general category or pattern that they observed in the training phase.

The following graph computes the probability of endorsement in the data set as a function of stimulus type and condition (Cat or Rec instructions).

```

In [3]: #####
tcurve
#####
getcurve(filename):
    prototypes = []
    low = []
    old = []
    high = []
    random = []
    mydata = readfile(filename)
    cond = mydata[-1][1]
    for line in mydata:
        if line[4] == 2 and len(line) == 9:
            if line[7] == 1:
                prototypes.append(line[5])
            if line[7] == 2:
                old.append(line[5])
            if line[7] == 3:
                high.append(line[5])
            if line[7] == 4:
                low.append(line[5])
            if line[7] == 5:
                random.append(line[5])

    #print([len(prototypes), len(low), len(high), len(random), len(old)])
    # print(prototypes)
    # print(low)
    # print(high)
    # print(random)
    # print(old)
    return [np.average(prototypes), np.average(low), np.average(high), np.average(random), np.average(old)]

readfile(filename):
    results = []
    fp = open(filename, 'r')
    for line in fp.readlines():
        myline = list(map(int, line.split(' ')[:-1]))
        results.append(myline[:])
    fp.close()
    return results

get_all_filenames(directoryname):
    files = filter(lambda x: x[-4:] == '.dat' and x[0] != '.', os.listdir(os.path.join(".", directoryname)))
    fn = map(lambda x: os.path.join(".", directoryname, x), files)
    # process each file and drop last 5 trials
    return list(fn)

create_df(subjnum, cond, pattern):
    nobs = len(pattern)
    df = pd.DataFrame({"Subject": [subjnum]*nobs, "Condition": [cond]*nobs, "Stimulus": [pattern]*nobs,
                       'Prototype', 'Low', 'High', 'Random', 'Old'], "Probability": [0.5]*nobs)
    return df

```

```

get_human_results():
    allres = map(getcurve, get_all_filenames('data'))
    cat = []
    rec = []
    for patt in allres:
        if patt[-1] == 0:
            cat.append(create_df(patt[-2], 'cat', patt[:-2]))
        else:
            rec.append(create_df(patt[-2], 'rec', patt[:-2]))
    cat, rec = pd.concat(cat), pd.concat(rec)
    #print("cat", cat, "rec", rec )
    return pd.concat([cat, rec])

```

```

sns.barplot(x="Stimulus Type", y="Probability of Endorsement",
hue="Condition", data=get_human_results())

```

Problem 1 (20 points)

Using your own words explain the data pattern you see in the above figure. What is different between the conditions and stimulus type? Why do you suspect these patterns exist? Your answer will need to consider the nature of the experiments, what is manipulated, and even your intuitive psychological theory about what might be going on. Your response should take 3-4 sentences and appear in a cell below. Is any feature of this data surprising to you?

In the categorical condition, participants are looking to see if the pattern is in the same general category or pattern as the training phase, as opposed to an identical pattern as in the recognition condition. Therefore, it makes sense participants in the 'categorical' condition consistently gave a higher endorsement on average across all five stimulus types. In the 'categorical' condition, the 'low' stimulus condition had the highest rating. This is likely because the small amount of variance introduced in the low condition was not enough to make the object seem like it was outside the category. Meanwhile, in the 'recognition' condition, that amount of noise made the object not identical and the participant did not endorse it. Instead, the 'old' stimulus condition had the highest probability of endorsement, because indeed some of the patterns were repeats of old patterns.

Predictions for the exemplar model

The following cells set up the exemplar model using the equations described above.

```
In [4]: #####  
# unitdist:  
# computes the euclidean distance between  
# two dots  
#####  
def unitdist(x, y):  
    x1 = np.array(x)  
    y1 = np.array(y)  
    return math.sqrt(sum(pow(x-y, 2.0)))  
  
#####  
# computeresponse  
# computes the "activation" of each  
# trace in memory  
#####  
def computeresponse(target, memory, c, k):  
    res = []  
    for mem in memory:  
        res.append(  
            math.log(1.0+np.average(list(map(lambda x, y: unitdist(x, y), t  
resp = [math.exp(-c*x) for x in res]  
pofr = sum(resp)/(sum(resp)+k)  
return pofr
```

```

In [5]: #####
# exemplar model
# stores all 10 study items in memory
# and computes the probability of endorsement
# for each item type
#####

def exemplar_model(filename, c, k):
    data = readfile(filename)
    cond = data[-1][1]
    memory = []
    for line in data:
        if len(line) == 20 and line[1] == 2:
            memory.append(np.resize(line[2:], (9, 2)))
    # print(memory)

    # prototype items
    proto = []
    for line in data:
        if len(line) == 20 and line[1] == 1:
            item = np.resize(line[2:], (9, 2))
            pofr = computeresponse(item, memory, c, k)
            proto.append(pofr)
    # print(np.average(proto))

    # old items
    old = []
    for line in data:
        if len(line) == 20 and line[1] == 2:
            item = np.resize(line[2:], (9, 2))
            pofr = computeresponse(item, memory, c, k)
            old.append(pofr)
    # print "p of r", old
    # print(np.average(old))

    # new high items
    newhigh = []
    for line in data:
        if len(line) == 20 and line[1] == 3:
            item = np.resize(line[2:], (9, 2))
            pofr = computeresponse(item, memory, c, k)
            newhigh.append(pofr)
    # print(np.average(newhigh))

    # new low items
    newlow = []
    for line in data:
        if len(line) == 20 and line[1] == 4:
            item = np.resize(line[2:], (9, 2))
            pofr = computeresponse(item, memory, c, k)
            newlow.append(pofr)
    # print(np.average(newlow))

    # random items
    random = []
    for line in data:

```

```

if len(line) == 20 and line[1] == 5:
    item = np.resize(line[2:], (9, 2))
    pofr = computeresponse(item, memory, c, k)
    random.append(pofr)
# print(np.average(random))

return [np.average(proto), np.average(newlow), np.average(newhigh), np.

```

```

In [6]: def get_exemplar_results(c_cat, k_cat, c_rec, k_rec):
    allres = {fn: readfile(fn) for fn in get_all_filenames('data')}
    cat = []
    rec = []
    for filename in allres.keys():
        if allres[filename][-1][1] == 0:
            res = exemplar_model(filename, c_cat, k_cat)
            cat.append(create_df(filename, 'cat', res[:-2]))
        else:
            res = exemplar_model(filename, c_rec, k_rec)
            rec.append(create_df(filename, 'rec', res[:-2]))
    cat, rec = pd.concat(cat), pd.concat(rec)
    return pd.concat([cat, rec])

```

First let's replot the human results:

```

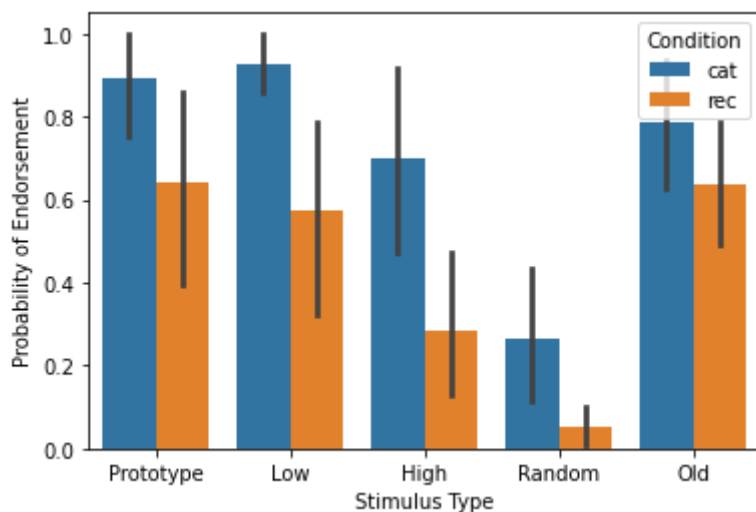
In [7]: sns.barplot(x="Stimulus Type", y="Probability of Endorsement", hue="Condition")

```

```

Out[7]: <matplotlib.axes._subplots.AxesSubplot at 0x122d24d00>

```



Problem 2 (10 points)

By hand adjust the setting of the model parameters to roughly fit the human data pattern shown above. How close can you get? What parameters did you find (report them) and your assessment of how well they fit. Was it a good fit or are there systematic problems with the fit? In addition, what are the parameter values and do they make sense in light of the

equations described above? When the parameters are the same for recognition and categorization instructions why do the bars look a little different?

Problem 2

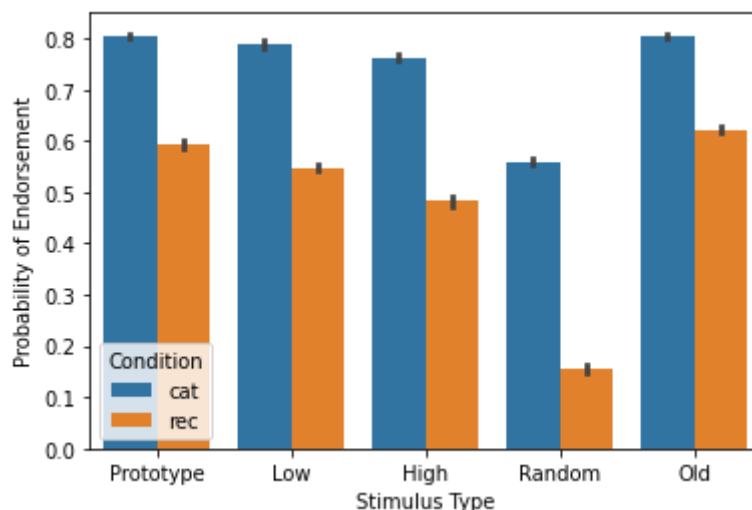
I hand-adjusted the parameters as follows: c_{cat} , k_{cat} , c_{rec} , k_{rec} = 0.8, 0.8, 1.3, 1.3

I can get moderately close to the human data pattern show above, but not precisely. I am not able to get the categorical ratings close to 0.9 while maintaining the variance between the different stimulus types. The values of the recognition condition were more of a fit. When the parameters are the same for recognition and categorization instructions, the bar look a little different because we are plotting the model predictions for the exemplar model and it is not real human data.

The following cell lets you plot the model predictions for the exemplar model fitted to the stimuli that participants in this experiment actually viewed. There is a k and a c parameter for both categorization and recognition.

```
In [8]: c_cat, k_cat, c_rec, k_rec = 0.8, 0.8, 1.3, 1.3
sns.barplot(x="Stimulus Type", y="Probability of Endorsement",
            hue="Condition", data=get_exemplar_results(c_cat, k_cat, c_rec,
```

```
Out[8]: <matplotlib.axes._subplots.AxesSubplot at 0x122e7f850>
```



Predictions for the prototype model

The following cells set up the prototype model using the equations described above.

```

In [9]: #####
        prototype model
        res an avarate of the study items in memory
        computes the probability of endorsement
        each item type
        #####
        rototypemodel(filename, c, k):
        ata = readfile(filename)
        ond = data[-1][1]
        average all the old items in memory
        emory = []
        or line in data:
            if len(line) == 20 and line[1] == 2:
                memory.append(line[2:])
        emory = [
            np.resize(list(map(np.average, np.transpose(np.array(memory)))), (9, 2)),

        prototype items
        roto = []
        or line in data:
            if len(line) == 20 and line[1] == 1:
                item = np.resize(line[2:], (9, 2))
                pofr = computeresponse(item, memory, c, k)
                proto.append(pofr)
        print(np.average(proto))

        old items
        ld = []
        or line in data:
            if len(line) == 20 and line[1] == 2:
                item = np.resize(line[2:], (9, 2))
                pofr = computeresponse(item, memory, c, k)
                old.append(pofr)
        print "p of r", old
        print(np.average(old))

        new high items
        ewhigh = []
        or line in data:
            if len(line) == 20 and line[1] == 3:
                item = np.resize(line[2:], (9, 2))
                pofr = computeresponse(item, memory, c, k)
                newhigh.append(pofr)
        print(np.average(newhigh))

        new low items
        ewlow = []
        or line in data:
            if len(line) == 20 and line[1] == 4:
                item = np.resize(line[2:], (9, 2))
                pofr = computeresponse(item, memory, c, k)
                newlow.append(pofr)
        print(np.average(newlow))

        random items
        andom = []
        or line in data:

```



```

if len(line) == 20 and line[1] == 5:
    item = np.resize(line[2:], (9, 2))
    pofr = computeresponse(item, memory, c, k)
    random.append(pofr)
print(np.average(random))

return [np.average(proto), np.average(newlow), np.average(newhigh), np.average(newhigh)]

```

```

In [10]: def get_prototype_results(c_cat, k_cat, c_rec, k_rec):
allres = {fn: readfile(fn) for fn in get_all_filenames('data')}
cat = []
rec = []
for filename in allres.keys():
    if allres[filename][-1][1] == 0:
        res = prototypemodel(filename, c_cat, k_cat)
        cat.append(create_df(filename, 'cat', res[:-2]))
    else:
        res = prototypemodel(filename, c_rec, k_rec)
        rec.append(create_df(filename, 'rec', res[:-2]))
cat, rec = pd.concat(cat), pd.concat(rec)
return pd.concat([cat, rec])

```

Again lets replot the human results for easy reference.

```

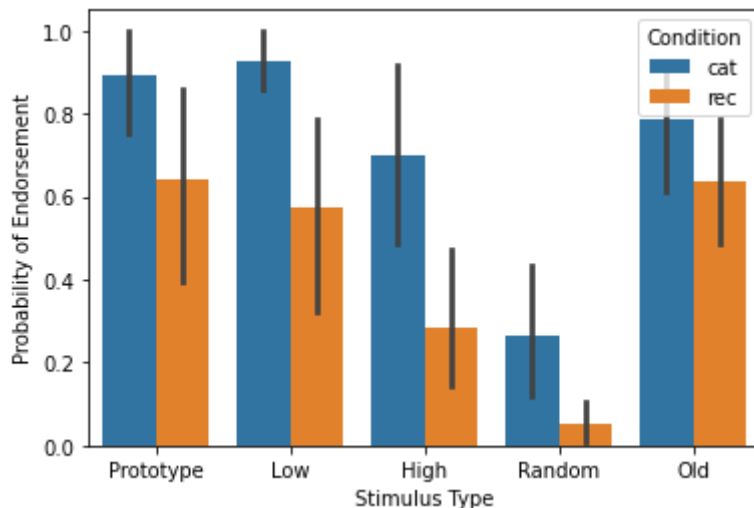
In [11]: sns.barplot(x="Stimulus Type", y="Probability of Endorsement", hue="Condition")

```

```

Out[11]: <matplotlib.axes._subplots.AxesSubplot at 0x122f32c40>

```



Problem 3 (10 points)

By hand adjust the setting of the model parameters in the next cell to roughly fit the human data pattern shown above. How close can you get? What parameters did you find (report them) and your assessment of how well they fit. Was it a good fit or are there systematic

problems with the fit? In addition, what are the parameter values and do they make sense in light of the equations described above?

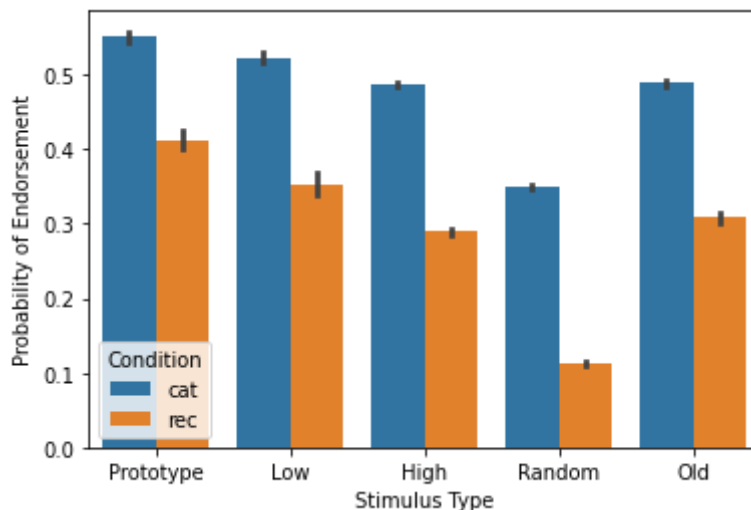
Problem 3

I hand-adjusted the setting of the model parameters and ended up with these values: c_{cat} , k_{cat} , c_{rec} , $k_{rec} = 0.4, 0.6, .8, .8$

It is harder to get close to the human data values here than in problem 2. This may be because the prototype model is different from what human systems use. Systematically, my values were all lower than those in the human data condition. The values for k and c are constants.

```
In [12]: c_cat, k_cat, c_rec, k_rec = 0.4, 0.6, .8, .8
res=get_prototype_results(c_cat, k_cat, c_rec, k_rec)
sns.barplot(x="Stimulus Type", y="Probability of Endorsement", hue="Condition")
```

```
Out[12]: <matplotlib.axes._subplots.AxesSubplot at 0x123089400>
```



```
In [ ]:
```

Fitting the models using RMSE

Next we would like to come up with a more quantitative way to assess the quality of the model fits. The first technique we will use is the "goodness of fit" measures that were discussed in lecture. One of the most common measures of goodness of fit is the Root Mean Squared Error (RMSE). This measure compares the value of each data point x to each prediction y using the following equation:

$$RMSE = \sqrt{\frac{\sum_i (x_i - y_i)^2}{N}}$$

Often the RMSE is computed between the AVERAGE prediction of the model and the AVERAGE

estimates of the behavior to all the subjects in an experiment. Using the code we developed above we can find the average endorsement curves for humans and both models like this:

```
In [13]: human_results=get_human_results()
exemplar_predictions = get_exemplar_results(c_cat, k_cat, c_rec, k_rec)
prototype_predictions = get_prototype_results(c_cat, k_cat, c_rec, k_rec)
```

```
In [14]: avghuman=human_results.groupby(['Condition', 'Stimulus Type'],as_index=False)
avghuman
```

Out[14]:

	Condition	Stimulus Type	Probability of Endorsement
0	cat	High	0.700000
1	cat	Low	0.928571
2	cat	Old	0.785714
3	cat	Prototype	0.892857
4	cat	Random	0.264286
5	rec	High	0.285714
6	rec	Low	0.571429
7	rec	Old	0.635714
8	rec	Prototype	0.642857
9	rec	Random	0.050000

```
In [15]: avgexemplar=exemplar_predictions.groupby(['Condition', 'Stimulus Type'],as_index=False)
avgexemplar
```

Out[15]:

	Condition	Stimulus Type	Probability of Endorsement
0	cat	High	0.894087
1	cat	Low	0.900813
2	cat	Old	0.901649
3	cat	Prototype	0.904897
4	cat	Random	0.841477
5	rec	High	0.772026
6	rec	Low	0.799226
7	rec	Old	0.813961
8	rec	Prototype	0.816479
9	rec	Random	0.554575

```
In [16]: avgprototype=prototype_predictions.groupby(['Condition', 'Stimulus Type'],
avgprototype
```

```
Out[16]:
```

	Condition	Stimulus Type	Probability of Endorsement
0	cat	High	0.486485
1	cat	Low	0.521748
2	cat	Old	0.488354
3	cat	Prototype	0.550461
4	cat	Random	0.348944
5	rec	High	0.290060
6	rec	Low	0.352924
7	rec	Old	0.308289
8	rec	Prototype	0.412433
9	rec	Random	0.112626

Problem 4 (20 points)

First, write a function below called `rmse` that computes the RMSE between two `numpy` vectors.

```
In [17]: def rmse(human, model):
          return np.sqrt(((human - model) ** 2).mean())
```

Write your code above. This code will then be used in the provided functions below to evaluate the fit of the prototype and exemplar models. The parameters to the model is provided as a list with `[c_cat, k_cat, c_rec, k_rec]` the implied order.

```
In [18]: def fit_exemplar_model_rmse(params, human_results):
          [c_cat, k_cat, c_rec, k_rec] = params
          predictions = get_exemplar_results(c_cat, k_cat, c_rec, k_rec)
          avgpredict=predictions.groupby(['Condition', 'Stimulus Type'],as_index=
          model_results = avgpredict['Probability of Endorsement'].values
          return rmse(human_results, model_results)

          def fit_prototype_model_rmse(params, human_results):
              [c_cat, k_cat, c_rec, k_rec] = params
              predictions = get_prototype_results(c_cat, k_cat, c_rec, k_rec)
              avgpredict=predictions.groupby(['Condition', 'Stimulus Type'],as_index=
              model_results = avgpredict['Probability of Endorsement'].values
              return rmse(human_results, model_results)
```

```
In [25]: human_results=get_human_results()
avghuman=human_results.groupby(['Condition', 'Stimulus Type'],as_index=False)
human_results = avghuman['Probability of Endorsement'].values

print(fit_exemplar_model_rmse([0.1, 2.0, 0.1, 2.0], human_results))
print(fit_prototype_model_rmse([0.1, 2.0, 0.1, 2.0], human_results))

0.3533017820219753
0.3792712020204823
```

Next adjust the parameters by hand for both the exemplar and prototype models to find values that appear to minimize the RMSE. Copy the code above for plotting the predictions of the models given your best fit parameters. Which model do you think fits better according to this fit statistic?

The values that appear to minimize Root mean squared error are: 0.2, 2.0, 1.0, 2.0 for the exemplar model, and 0.1, 0.4, 1.0, 0.4 for the prototype model. Both models fit similarly according to the fit statistic. Each gives a RMSE of 0.18. Exemplar and prototype models both somewhat strongly match human results.

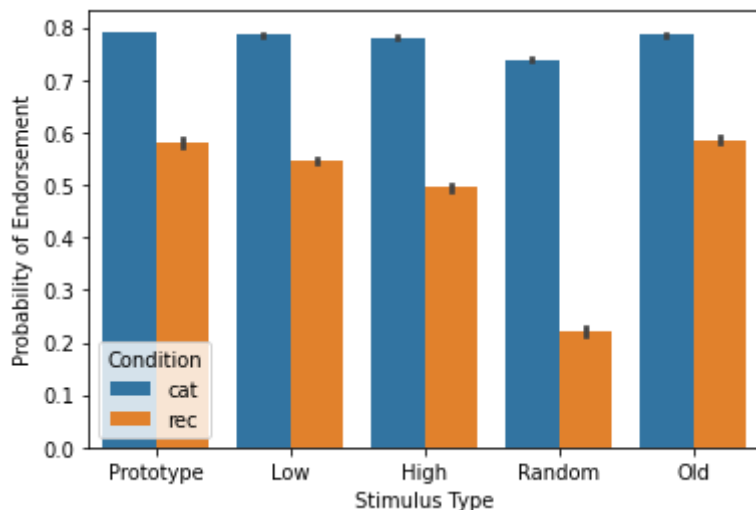
```
In [41]: human_results=get_human_results()
avghuman=human_results.groupby(['Condition', 'Stimulus Type'],as_index=False)
human_results = avghuman['Probability of Endorsement'].values

print(fit_exemplar_model_rmse([0.2, 2.0, 1.0, 2.0], human_results))
print(fit_prototype_model_rmse([0.1, 0.4, 1.0, 0.4], human_results))

0.18491067201537803
0.1848440427010307
```

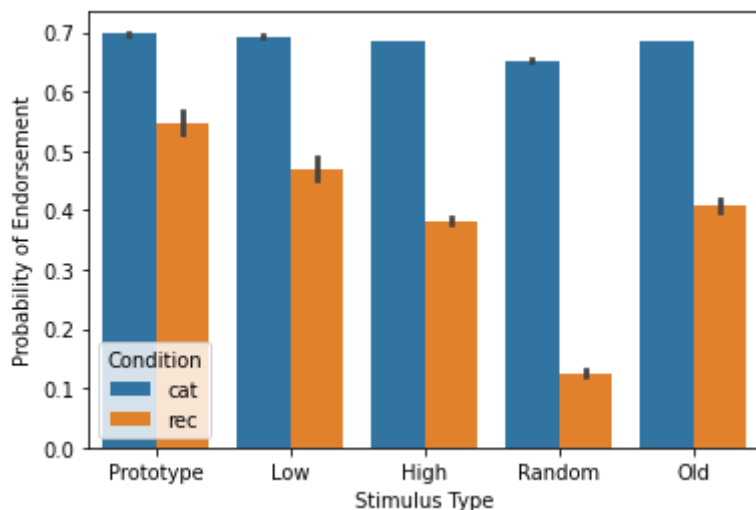
```
In [42]: # exemplar
c_cat, k_cat, c_rec, k_rec = 0.2, 2.0, 1.0, 2.0
sns.barplot(x="Stimulus Type", y="Probability of Endorsement",
            hue="Condition", data=get_exemplar_results(c_cat, k_cat, c_rec,
```

Out[42]: <matplotlib.axes._subplots.AxesSubplot at 0x12356eac0>



```
In [43]: # prototype
c_cat, k_cat, c_rec, k_rec = 0.1, 0.4, 1.0, 0.4
res=get_prototype_results(c_cat, k_cat, c_rec, k_rec)
sns.barplot(x="Stimulus Type", y="Probability of Endorsement", hue="Condi
```

Out[43]: <matplotlib.axes._subplots.AxesSubplot at 0x123151df0>



BONUS (5 points)

Read about the scipy `fmin` function. Use fmin to algorithmically search for the best parameters for each model using the RMSE evaluation function described above.

Elena Georgieva - 4/23/2021

Homework - Categorization and Model Comparison Part B (40/100 points)

by *Todd Gureckis* and *Brenden Lake*

Computational Cognitive Modeling

NYU class webpage: <https://brendenlake.github.io/CCM-site/> (<https://brendenlake.github.io/CCM-site/>)

email to course instructors: instructors-ccm-spring2021@nyuccl.org (<mailto:instructors-ccm-spring2021@nyuccl.org>)

This homework is due before midnight on Apr 19, 2021.

```
In [1]: from IPython.display import display
import string as str
import os
import numpy as np
import seaborn as sns
import pandas as pd
import math
from random import random, randint, shuffle, uniform
from scipy.optimize import fmin, brute
from scipy.special import comb # gets the combinations function
from exemplarproto import * # this grabs much of the code from Part A of t
```

Fitting the models using maximum likelihood

As mentioned in the lecture, RMSE is not always an ideal mechanism for fitting models. One reason is that it is insensitive to the number of observations that define each data point. For example, remember in our experiment that participants saw the prototype item four times at test. In contrast, there were 20 different "new" patterns. This means there are five times as many trials contributing to the "new" bar in this graph as for the prototype patterns. Since RMSE measures the raw deviation of the average model predictions from those of the model it doesn't take into account these issues. Thus, we would like to also evaluate these two models using maximum likelihood.

The key to this is going to be the provided function below which computes the likelihood of a particular set of data under a binomial probability model.:

```
In [2]: #####
# computeLogLikelihood
# N = number of observations
# S = number of "successes" (i.e., endorsements)
# p = predicted probability of successes by the model
#####
def computeLogLikelihood(N, S, p):
    p = p if p > 0.0 else 0.0+1e-10
    p = p if p < 1.0 else 1.0-1e-10
    try:
        result = math.log(comb(N, S)) + (S*math.log(p) + (N-S)*math.log(1.0-p))
    except:
        print(N, S, p) # this shouldn't happen but just in case
        result = 0
    return result

def pandas_ll(row):
    return computeLogLikelihood(row['Total'], row['N_Yes'], row['Probabilit
```

A short explanation may be in order: the models predictions take the form of probabilities of endorsement for each of the prototype, low, high, random, and old items. If you find out that people endorse the prototype on 2 out of 2 trials how likely is this outcome given that the model (for a particular set of parameters) predicts an endorsement of $p=0.8$? Three numbers are required to do this for each data point N , the number of trials/presentations within the stimulus class, S the number of successes observed ($S \leq N$), and p the predicted probability. Then you can turn the crank on the above `computeLogLikelihood()` function which returns the probability that you would get S successes in N trials if the true probability was p (make sure you understand what is happening in `computeLogLikelihood`). You can sum these log likelihoods for each stimulus class (prototype, low, high, random, old) to compute a total log(likelihood) of the data for any given model with any set of parameters. For this homework will we focus on fitting the group data rather than to individuals.

To get the data formatted into an appropriate shape for fitting likelihoods we provide a function `get_human_results_ll()` which returns a Pandas data frame containing the number of times a pattern of a particular type was endorsed and the number of times it was presented for each subject.


```
In [3]: human_res=get_human_results_ll()
human_res
```

Out[3]:

	Subject	Condition	Stimulus Type	N_Yes	Total
0	./data/8.dat	cat	Prototype	4	4
1	./data/8.dat	cat	Low	4	4
2	./data/8.dat	cat	High	10	10
3	./data/8.dat	cat	Random	15	20
4	./data/8.dat	cat	Old	20	20
...
0	./data/7.dat	rec	Prototype	3	4
1	./data/7.dat	rec	Low	4	4
2	./data/7.dat	rec	High	8	10
3	./data/7.dat	rec	Random	0	20
4	./data/7.dat	rec	Old	19	20

70 rows × 5 columns

This reorganizes the data per condition.

```
In [4]: human_data=human_res.groupby(['Condition', 'Stimulus Type']).sum()
human_data
```

Out[4]:

		N_Yes	Total
Condition	Stimulus Type		
cat	High	49	70
	Low	26	28
	Old	110	140
	Prototype	25	28
	Random	37	140
rec	High	20	70
	Low	16	28
	Old	89	140
	Prototype	18	28
	Random	7	140

Finally these function allow us to compute the negative log likelihood of the data given the model.

```
In [5]: def fit_exemplar_model_nll(params, human_results):
    [c_cat, k_cat, c_rec, k_rec] = params
    k_cat = k_cat if k_cat > 0.0 else 0.0
    k_rec = k_rec if k_rec > 0.0 else 0.0
    predictions = get_exemplar_results(c_cat, k_cat, c_rec, k_rec)
    model = predictions.groupby(
        ['Condition', 'Stimulus Type'], as_index=False).mean()
    fitted_data = pd.merge(model, human_results)
    return -1.0*fitted_data.apply(pandas_ll, axis=1).sum()

def fit_prototype_model_nll(params, human_results):
    [c_cat, k_cat, c_rec, k_rec] = params
    k_cat = k_cat if k_cat > 0.0 else 0.0
    k_rec = k_rec if k_rec > 0.0 else 0.0
    predictions = get_prototype_results(c_cat, k_cat, c_rec, k_rec)
    model = predictions.groupby(
        ['Condition', 'Stimulus Type'], as_index=False).mean()
    fitted_data = pd.merge(model, human_results)
    return -1.0*fitted_data.apply(pandas_ll, axis=1).sum()
```

Problem 5 (20 points)

The cell blocks below allow you to fit the exemplar model and the prototype model to the dataset we considered in Part A of the homework. Make sure you understand and follow the code provided above and in the provided library (exemplarproto.py). Next, try altering the parameters to minimize the negative log likelihood score. When you think you have found the best fit parameters for both the exemplar and prototype models report your final parameter values along with the plot of the resulting model predictions. In a markdown cell describe which model you believe fits better. Is this conclusion the same or different from what you considered in Part 4 of the homework? If the fit looks different, why?

Problem 5

The final parameter values I ended up with: [2.0, 0.1, 2.0, 0.2] for the exemplar model and [1.1, 0.1, 1.0, 0.2] for the prototype model. These minimize the negative log likelihood score, which were 93.7 and 56.5 for the exemplar and prototype models, respectively. In general, I believe Prototype is a better model because I was better able to reduce the negative log likelihood score. Also, looking at the plot of the human data, it seems the prototype model is more similar. This is the same outcome I got for part 4 or the first Python notebook for this assignment.

Exemplar model

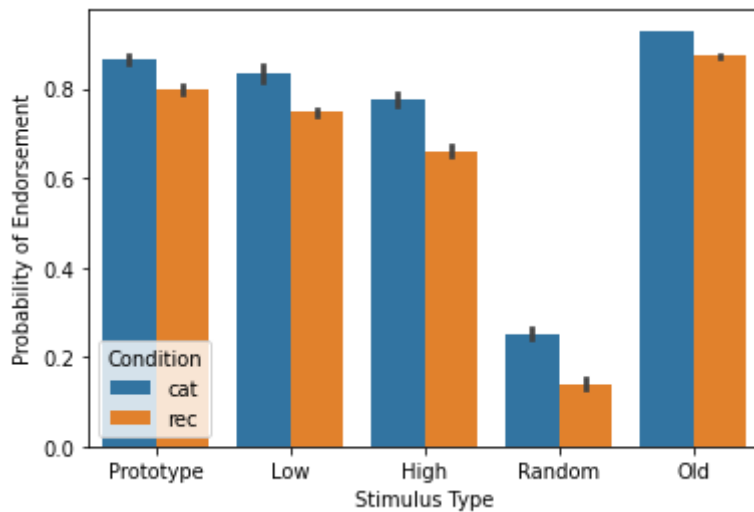
```
In [28]: human = human_res.groupby(['Condition', 'Stimulus Type'], as_index=False).s

params = [2.0, 0.1, 2.0, 0.2]
nllfit = fit_exemplar_model_nll(params, human)
print(f'The negative log score is {nllfit}')

# now plot the data
c_cat, k_cat, c_rec, k_rec = params
res = get_exemplar_results(c_cat, k_cat, c_rec, k_rec)
sns.barplot(x="Stimulus Type", y="Probability of Endorsement",
            hue="Condition", data=res)
```

The negative log score is 93.72570776237882

Out[28]: <matplotlib.axes._subplots.AxesSubplot at 0x125515a30>



Prototype Model

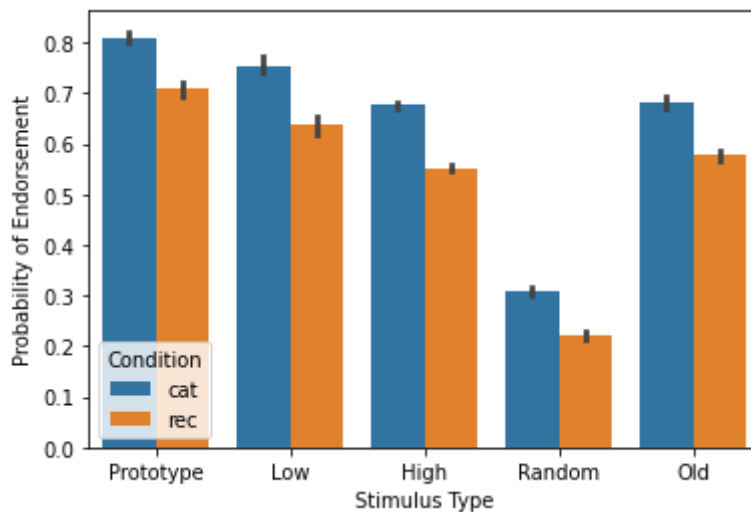
```
In [24]: human = human_res.groupby(['Condition', 'Stimulus Type'], as_index=False).s

params = [1.1, 0.1, 1.0, 0.2]
nllfit = fit_prototype_model_nll(params, human)
print(f'The negative log score is {nllfit}')

# now plot the data
c_cat, k_cat, c_rec, k_rec = params
res = get_prototype_results(c_cat, k_cat, c_rec, k_rec)
sns.barplot(x="Stimulus Type", y="Probability of Endorsement",
            hue="Condition", data=res)
```

The negative log score is 56.520395053419335

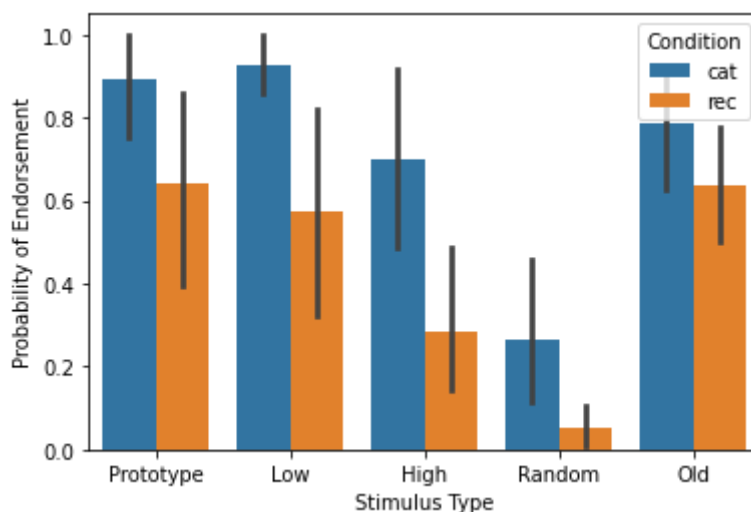
Out[24]: <matplotlib.axes._subplots.AxesSubplot at 0x124f50490>



Human data again for reference

```
In [27]: sns.barplot(x="Stimulus Type", y="Probability of Endorsement",
                    hue="Condition", data=get_human_results())
```

Out[27]: <matplotlib.axes._subplots.AxesSubplot at 0x125458d30>



Problem 6 (10 points)

A famous saying is the "All models are wrong, but some are useful" (George Box). Do you think the exemplar or prototype model provides the best account of the data? Refer to particular patterns in the data that you believe the different models do a better job with.

In general, I feel that Prototype is a better model than exemplar because I was better able to reduce the negative log likelihood score. I got a score of 56.5 for prototype, as opposed to a score of 93.7 for exemplar. Also, looking at the plot of the human data, it seems the prototype model is more similar. Prototype performs better than exemplar especially for the 'random' stimulus type. The exemplar model mistakenly endorses in the "old" stimulus type condition, because that content was shown in training. Prototype theory makes a conclusion about concept representation, the concept is represented as the 'ideal' or 'average' category example. Still, in the end, both models perform similarly well. It would be interesting to use other statistics to compare the models.

Problem 7 (5 points)

Thinking about how these models work explain why both the exemplar and prototype models have relatively high endorsement for the prototype item even though it was never presented during the training phase. In addition, explain in your own words why the models are able to explain the high endorsement rates for the old items.

The prototype model represents an object as the ideal or average category example. And the exemplar model represents an object as the best category example. Both models work somewhat well and have a high endorsement for the prototype item even if it was not presented during training! "Old" patterns are identical to those presented during the training phase, and are endorsed at a relatively high rate. It is endorsed at a high rate because the algorithm has been exposed to the information during the training phase.

Problem 8 (5 points)

Are the exemplar model and the prototype model we considered nested? Would we compare them using AIC, BIC, or the G^2 statistic (or something else)?

Both the prototype and exemplar models are fairly strong. Both models can be nested but are not necessarily so. Many categories fall into hierarchies where categories are nested inside larger more abstract categories. For example animal >> mammal >> bear >> grizzly bear. To compare these models them, it would be nice to use an appropriate model selection model. AIC, BIC, and Bayes

factor would be possible. AIC and BIC ignore the contribution of functional form to complexity. Bayes factor is sensitive to the functional form and may be a good fit. The G^2 statistic likely wouldn't be a good choice.

In []: