

```

from google.colab import files
uploaded = files.upload()

import pymc as pm
import matplotlib.pyplot as plt
import arviz as az
import pandas as pd
from scipy import special, stats
import numpy as np
import seaborn as sns

```

4. (a) You have been asked by the government of Bangladesh to determine whether the use of contraceptives by women in Bangladesh varies by district. The data come from surveys conducted in Bangladesh [2] (a) Develop three models for these data: pooled, unpoled, and hierarchical for all districts to predict usage of contraceptives. Use only district and age_centered as predictor variables. For each model briefly explain your choice of priors. (b) For each model evaluate the sampling performance and discuss your findings. (c) Plot the posterior distributions for the parameters for each model and discuss what they show regarding the question posed by WHO. (d) Plot each of the predictions with age_centered on the x-axis and the expected proportion of women using contraception on the y-axis with overlaid plots for each of the districts, as appropriate. Briefly explain these results as you will report them to the government of Bangladesh.

```
bang = pd.read_csv('bangladesh.csv', header=0)
```

```
bang
```

```

# HalfStudentT distribution for the error term provides robustness
sigma = pm.HalfStudentT('sigma', nu=2, sigma=8)

# Likelihood function assuming Bernoulli distribution for the binary outcome
y = pm.Bernoulli("y", p=pm.math.sigmoid(theta), observed=use_contraception, dims="obs_id")

# Sampling from the model
with unpooled_model:
    unpooled_trace = pm.sample(1000, tune=1000, target_accept=0.95)

```

In the un-pooled model, I assigned individual intercepts to each district with a broad normal prior, reflecting my neutral expectation and allowing for variability. The age effect also has a neutral prior. To handle outliers, I used a HalfStudentT distribution for the error term. This model captures both the unique effects of each district and the influence of age on contraceptive use.

```
#heirarchical model
```

```
age_centered = bang['age_centered'].values
use_Contraception = bang['use.contraception'].values
districts = pd.Categorical(bang['district']).codes
```

```
# Define coordinates
coords = {
    "district": np.unique(districts),
    "obs_id": np.arange(len(use_Contraception))
}
```

```
with pm.Model(coords=coords) as hierarchical_model:
    hierarchical_trace = pm.sample(1000, tune=1000, target_accept=0.95)
```

```

# Hyperpriors for the group means
a = pm.Normal("a", mu=0.0, sigma=5.0)
b = pm.Normal("b", mu=0.0, sigma=1.0)

# Hyperpriors for the group standard deviations

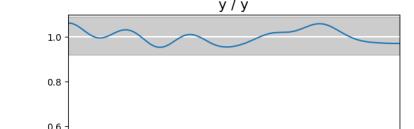
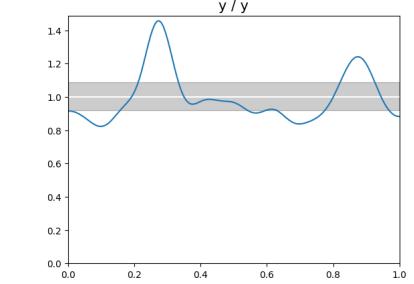
```

```

100.00% [4000/4000 00:01<00:00]
100.00% [2000/2000 00:00<00:00]
100.00% [2000/2000 00:01<00:00]

```

```
<Axes: title={'center': 'y / y'}
```



woman	district	use.contraception	living.children	age_centered	urban	
0	1	1	0	4	18.4400	1
1	2	1	0	1	-5.5599	1
2	3	1	0	3	1.4400	1
3	4	1	0	4	8.4400	1
4	5	1	0	1	-13.5599	1
...
1929	1930	61	0	4	14.4400	0
1930	1931	61	0	3	-4.5599	0
1931	1932	61	0	4	14.4400	0
1932	1933	61	0	1	-13.5599	0

(a) Develop three models for these data: pooled, unpoled, and hierarchical for all districts to predict usage of contraceptives. Use only district and age_centered as predictor variables. For each model briefly explain your choice of priors.

```
# Pooled model
```

```

use = pd.Categorical(bang['use.contraception'])
age = pd.Categorical(bang['age_centered'])
age = bang.loc[:, 'age_centered'].values
district = pd.Categorical(bang['district'])
districts = pd.Categorical(bang['district']).codes # Convert district to categorical codes
coords = {'district': np.arange(len(np.unique(districts))), # Number of unique districts
          'obs_id': np.arange(len(use_contraception)) # Number of observations
         }
coords = {'district': district.unique(), "age": age.codes, "obs_id": np.arange(age.size)}
coords = {'district': district.unique(), 'obs_id': np.arange(age.size)}

with pm.Model(coords=coords) as pooled:
    # un-centered
    age_centered = bang['age_centered'].values
    use_contraception = bang['use.contraception'].values
    districts = pd.Categorical(bang['district']).codes # Convert district to categorical codes
    coords = {
        "district": np.arange(len(np.unique(districts))), # Number of unique districts
        "obs_id": np.arange(len(use_contraception)) # Number of observations
    }

    with pm.Model(coords=coords) as unpoled_model:
        # Mutable data container for district
        district_idx = pm.Data("district_idx", districts, dims="obs_id")

        # Normal priors for the district effect
        beta_district = pm.Normal("beta_district", mu=0.0, sigma=10.0, dims="district")

        # Normal prior for the effect of age
        beta_age = pm.Normal("beta_age", mu=0.0, sigma=10.0)

        # Model equation
        theta = beta_district[district_idx] + beta_age * age_centered

    mu = beta0 + pm.math.dot(age, beta1)

    # HalfStudentT distribution for the error term provides robustness
    sigma_a = pm.HalfStudentT('sigma_a', nu=2, sigma=8)
    sigma_b = pm.HalfStudentT('sigma_b', nu=2, sigma=8)

    # Varying intercepts and slopes for each district
    a_district = pm.Normal("a_district", mu=a, sigma=sigma_a, dims="district")
    b_district = pm.Normal("b_district", mu=b, sigma=sigma_b, dims="district")

    # Model equation
    theta = a_district[district_idx] + b_district[district_idx] * age_centered

    # Model error
    sigma_y = pm.HalfStudentT("sigma_y", nu=2, sigma=8)

    # Likelihood
    y = pm.Bernoulli("y", p=pm.math.sigmoid(theta), observed=use_contraception, dims="obs_id")

    # Sampling
    with hierarchical_model:
        hierarchical_trace = pm.sample(1000, tune=1000, target_accept=0.95)

```

In the hierarchical model, I used varying intercepts and slopes for each district, guided by hyperpriors. This captures both district-specific effects and the overarching trend of age on contraceptive use. The HalfStudentT distributions provide robustness against outliers.

(b) For each model evaluate the sampling performance and discuss your findings.

```

#heirarchical model
with hierarchical_model:
    hierarchical_trace = pm.sample(1000, tune=1000, target_accept=0.95)

# Pooled
pooled_pp = pm.sample_posterior_predictive(trace_pooled, model = pooled)
az.plot_bpv(pooled_pp)

# Unpooled
unpooled_pp = pm.sample_posterior_predictive(unpooled_trace, model = unpoled_model)
az.plot_bpv(unpooled_pp)

```

As expected, without the use of district as a predictor variable, the unpoled model is not nearly as accurate as the others. We see the best performance within the pooled model with p-values that lie entirely in an acceptable range. Additionally, the heirarchical model seems to have a degree of accuracy that sits between the two. While the heirarchical may hypothetically have the best performance in some circumstances, the accuracy of this model might suggest suboptimal choice of priors or tuning for that model.

(c) Plot the posterior distributions for the parameters for each model and discuss what they show regarding the question posed by WHO.

```
az.plot_posterior(trace_pooled)
```

```

#μ = beta[age_idx]
#θ = pm.Deterministic("θ", pm.math.sigmoid(μ))

y = pm.Bernoulli("y", p=θ, observed=use.contraception, dims="obs_id")
#y = pm.Bernoulli("y", p=θ, observed=use, dims="obs_id")

trace_pooled = pm.sample(1000, cores=4, random_seed=1234, tune=1000)

# In the pooled model I developed for contraceptive use based on age, I chose relatively uninformative priors to let the data primarily inform the results. I set normal distributions with a mean of 0 for both the intercept and slope, indicating no initial bias towards any particular age effect. The standard deviation of 10 gives a wide possible range, showcasing my uncertainty about the true parameter values. This approach ensures that my model's conclusions rely heavily on the data and not on prior assumptions.

# un-pooled model
age_centered = bang['age_centered'].values
use_contraception = bang['use.contraception'].values
districts = pd.Categorical(bang['district']).codes # Convert district to categorical codes
coords = {
    "district": np.arange(len(np.unique(districts))), # Number of unique districts
    "obs_id": np.arange(len(use_contraception)) # Number of observations
}

with pm.Model(coords=coords) as unpoled_model:
    # Mutable data container for district
    district_idx = pm.Data("district_idx", districts, dims="obs_id")

    # Normal priors for the district effect
    beta_district = pm.Normal("beta_district", mu=0.0, sigma=10.0, dims="district")

    # Normal prior for the effect of age
    beta_age = pm.Normal("beta_age", mu=0.0, sigma=10.0)

    # Model equation
    theta = beta_district[district_idx] + beta_age * age_centered

    # HalfStudentT distribution for the error term provides robustness
    sigma_a = pm.HalfStudentT('sigma_a', nu=2, sigma=8)
    sigma_b = pm.HalfStudentT('sigma_b', nu=2, sigma=8)

    # Varying intercepts and slopes for each district
    a_district = pm.Normal("a_district", mu=a, sigma=sigma_a, dims="district")
    b_district = pm.Normal("b_district", mu=b, sigma=sigma_b, dims="district")

    # Model equation
    theta = a_district[district_idx] + b_district[district_idx] * age_centered

    # Model error
    sigma_y = pm.HalfStudentT("sigma_y", nu=2, sigma=8)

    # Likelihood
    y = pm.Bernoulli("y", p=pm.math.sigmoid(theta), observed=use_contraception, dims="obs_id")

    # Sampling
    with hierarchical_model:
        hierarchical_trace = pm.sample(1000, tune=1000, target_accept=0.95)

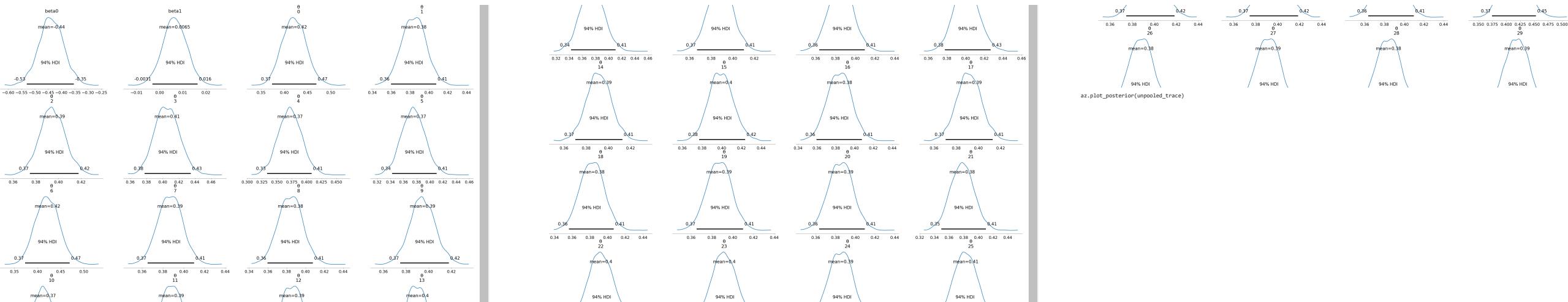
```

#heirarchical model
unpooled_pp = pm.sample_posterior_predictive(hierarchical_trace, model = hierarchical_model)
az.plot_bpv(unpooled_pp)

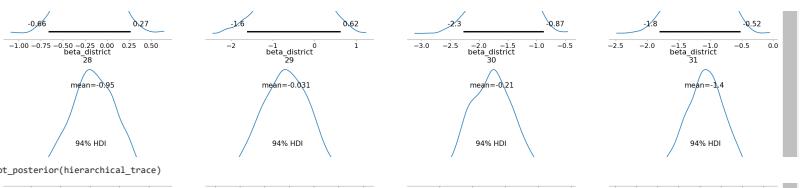
```

/usr/local/lib/python3.10/dist-packages/arviz/plots/plot_utils.py:271: UserWarning: rcParams['plot.max_subplots'] (40) is smaller than the number of vari
array[[<Axes: title={'center': 'βeta0'}>,
       <Axes: title={'center': 'βeta1'}>,
       <Axes: title={'center': 'βeta2'}>,
       <Axes: title={'center': 'βeta3'}>,
       <Axes: title={'center': 'βeta4'}>,
       <Axes: title={'center': 'βeta5'}>,
       <Axes: title={'center': 'βeta6'}>,
       <Axes: title={'center': 'βeta7'}>,
       <Axes: title={'center': 'βeta8'}>,
       <Axes: title={'center': 'βeta9'}>,
       <Axes: title={'center': 'βeta10'}>,
       <Axes: title={'center': 'βeta11'}>,
       <Axes: title={'center': 'βeta12'}>,
       <Axes: title={'center': 'βeta13'}>,
       <Axes: title={'center': 'βeta14'}>,
       <Axes: title={'center': 'βeta15'}>,
       <Axes: title={'center': 'βeta16'}>,
       <Axes: title={'center': 'βeta17'}>,
       <Axes: title={'center': 'βeta18'}>,
       <Axes: title={'center': 'βeta19'}>,
       <Axes: title={'center': 'βeta20'}>,
       <Axes: title={'center': 'βeta21'}>,
       <Axes: title={'center': 'βeta22'}>,
       <Axes: title={'center': 'βeta23'}>,
       <Axes: title={'center': 'βeta24'}>,
       <Axes: title={'center': 'βeta25'}>,
       <Axes: title={'center': 'βeta26'}>,
       <Axes: title={'center': 'βeta27'}>,
       <Axes: title={'center': 'βeta28'}>,
       <Axes: title={'center': 'βeta29'}>,
       <Axes: title={'center': 'βeta30'}>,
       <Axes: title={'center': 'βeta31'}>,
       <Axes: title={'center': 'βeta32'}>,
       <Axes: title={'center': 'βeta33'}>,
       <Axes: title={'center': 'βeta34'}>,
       <Axes: title={'center': 'βeta35'}>,
       <Axes: title={'center': 'βeta36'}>,
       <Axes: title={'center': 'βeta37'}>], dtype=object)

```



```
/usr/local/lib/python3.10/dist-packages/arviz/plots/plot_utils.py:271: UserWarning: rcParams['plot.max_subplots'] (48) is smaller than the number of variables (49)
  warnings.warn(array([{'Axes': titles['center'], 'beta_district': 'beta_district{0}'}], dtype=object)
```



These plots display the distribution for each of the betas. They demonstrate that the distributions chosen are good fits to measure the contraception use of women in Bangladesh.

(d) Plot each of the predictions with age_centered on the x-axis and the expected proportion of women using contraception on the y-axis with overlaid plots for each of the districts, as appropriate. Briefly explain these results as you will report them to the government of Bangladesh.

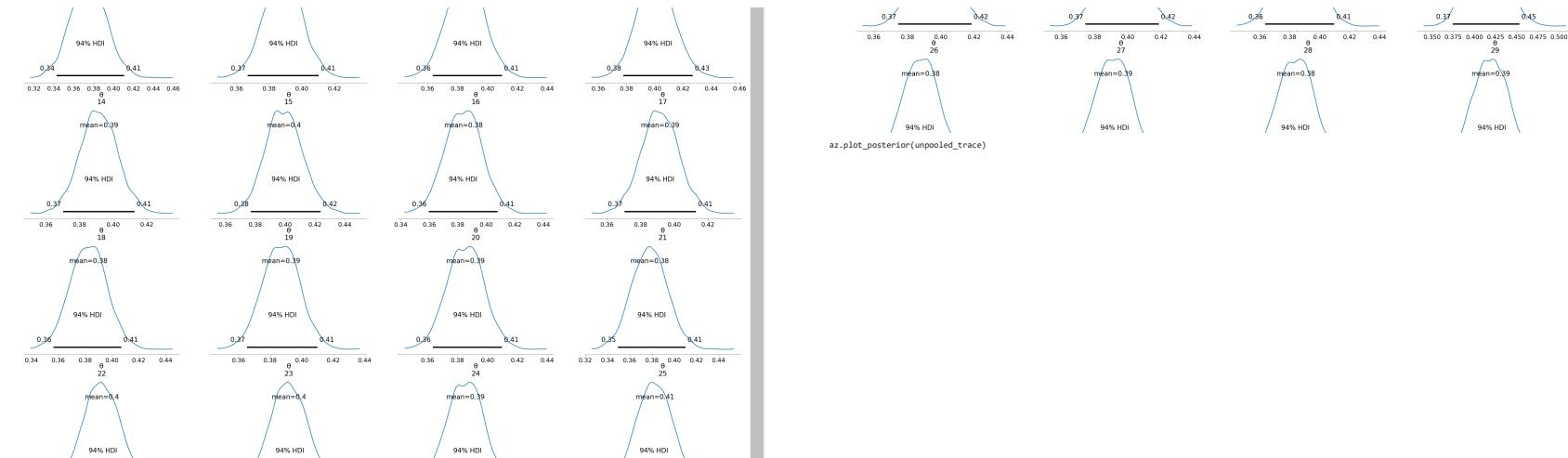
parameters from the model trace

```
beta0_pooled_samples = trace_pooled.posterior["beta0"].values
beta1_pooled_samples = trace_pooled.posterior["beta1"].values

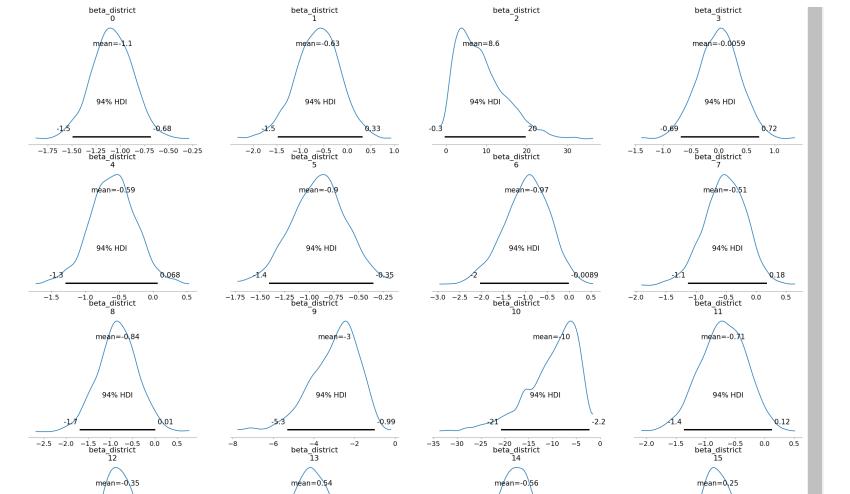
beta_district_samples = unpooled_trace.posterior["beta_district"].values
beta_age_samples = unpooled_trace.posterior["beta_age"].values

a_district_samples = hierarchical_trace.posterior["a_district"].values
b_district_samples = hierarchical_trace.posterior["b_district"].values
```

```
beta0_pooled_mean = np.mean(beta0_pooled_samples)
beta1_pooled_mean = np.mean(beta1_pooled_samples)
```



az.plot_posterior(unpooled_trace)



```
# making preds for pooled
prediction_pooled = beta0_pooled_mean + beta1_pooled_mean * age_centered

# Transforming the prediction to probability using logistic function
probability_pooled = 1 / (1 + np.exp(-prediction_pooled))
```

```
predictions = []

# Compute Predictions for Each District
for d in np.unique(districts):
    district_mean_beta = np.mean(beta_district_samples[:, d])
    age_mean_beta = np.mean(beta_age_samples)

    prediction = district_mean_beta + age_mean_beta * age_centered

    # Transform prediction to probability using logistic function
    probability = 1 / (1 + np.exp(-prediction))

    predictions.append(probability)
```

```
hierarchical_predictions = []

for d in np.unique(districts):
    a_district_mean = np.mean(a_district_samples[:, d])
    b_district_mean = np.mean(b_district_samples[:, d])

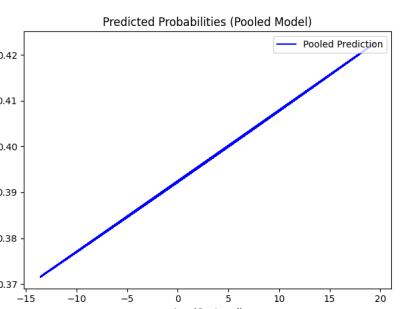
    prediction_hierarchical = a_district_mean + b_district_mean * age_centered

    probability_hierarchical = 1 / (1 + np.exp(-prediction_hierarchical))

    hierarchical_predictions.append(probability_hierarchical)
```

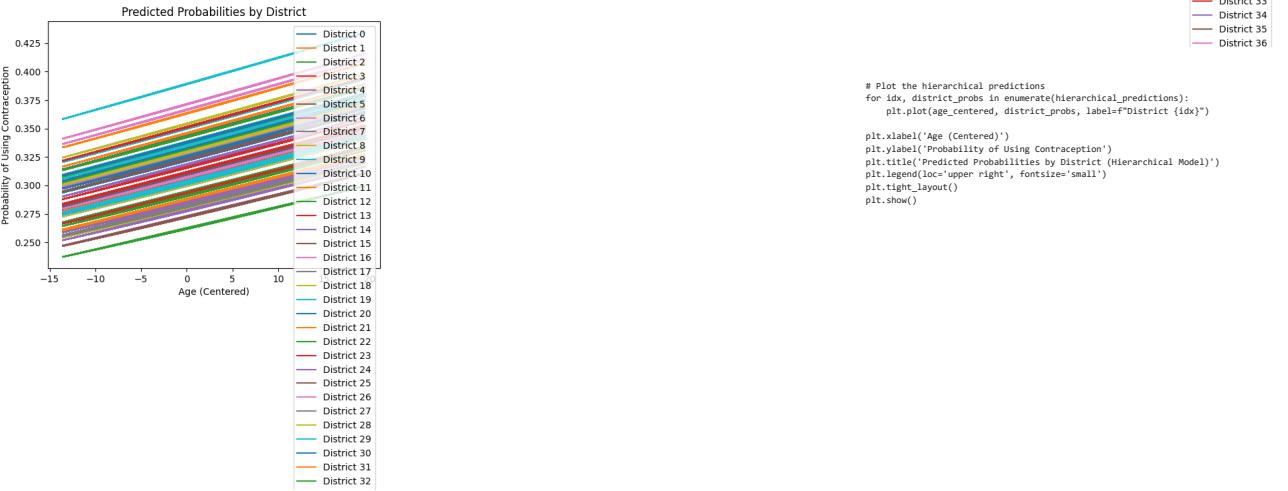
```
# plotting preds for pooled
plt.plot(age_centered, probability_pooled, label='Pooled Prediction', color='blue')

plt.xlabel('Age (Centered)')
plt.ylabel('Probability of Using Contraception')
plt.title('Predicted Probabilities (Pooled Model)')
plt.legend(loc='upper right')
plt.tight_layout()
plt.show()
```

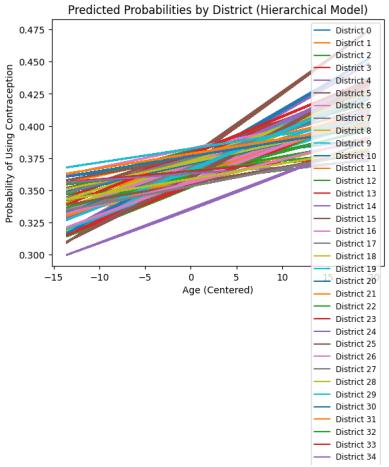


```
# Un-pooled: predictions for contraception based on age_centered
for idx, district_probs in enumerate(predictions):
    plt.plot(age_centered, district_probs, label=f'District {idx}')

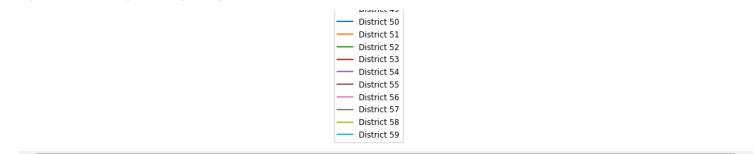

plt.xlabel('Age (Centered)')
plt.ylabel('Probability of Using Contraception')
plt.title('Predicted Probabilities by District')
plt.legend()
plt.show()
```



```
<ipython-input-21-5b58207d13f4>:9: UserWarning: Tight layout not applied. The bottom and top margins cannot be made large enough to accommodate all axes
plt.tight_layout()
```



Based on each of these above graphs, we can see how each model predicts the probability of contraception use in Bangladeshi women. In the unpooled model, we can see that the model, since it's pooled, doesn't differentiate between districts in predicting the odds of using contraception. On the other hand, the hierarchical and pooled models show the difference between districts when predicting the probability of using contraception. The hierarchical model shows different intercepts and slopes between districts, while the pooled model has different intercepts but the same slope for its logistic regression.



```
District 33
District 34
District 35
District 36
```

```
# Plot the hierarchical predictions
for idx, district_probs in enumerate(hierarchical_predictions):
    plt.plot(age_centered, district_probs, label=f'District {idx}')


plt.xlabel('Age (Centered)')
plt.ylabel('Probability of Using Contraception')
plt.title('Predicted Probabilities by District (Hierarchical Model)')
plt.legend(loc='upper right', fontsize='small')
plt.tight_layout()
plt.show()
```

```
import pymc as pm
import matplotlib.pyplot as plt
import arviz as az
import pandas as pd
from scipy import special, stats
import numpy as np
import seaborn as sns
```

```
from google.colab import files
uploaded = files.upload()
```

Choose Files [2 files]
 • algae.csv(text/csv) - 17987 bytes, last modified: 10/29/2023 - 100% done
 • algaeTest.csv(text/csv) - 20683 bytes, last modified: 10/29/2023 - 100% done
Saving algae.csv to algae.csv
Saving algaeTest.csv to algaeTest.csv

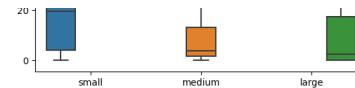
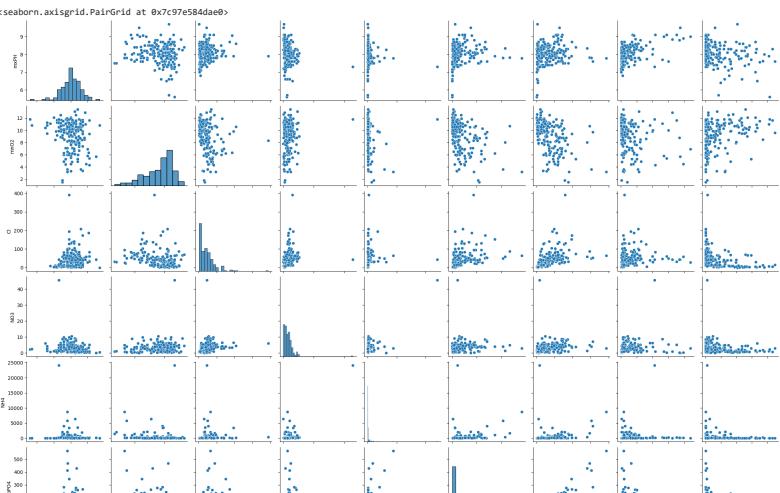
or each of these predictor variables observations, the data set also contains measurements of 7 different harmful algae, labeled a1-a7 but not named. For this work we will use the a1 measurements as the response variable. With algae.csv data and the stated goal of the EU perform the following tasks and answer the stated questions. (a) Perform exploratory data analysis (EDA) to include as minimum a review of the summary statistics for the data. Briefly explain how this EDA informs the next steps of your work.

```
#algae_train = pd.read_csv('algae.csv', header=0)
```

```
algae_train = pd.read_csv('algae.csv')

predictors = algae_train.columns[0:11]
for p in predictors:
    sum_stats = (algae_train.loc[:, [p, 'a1']].groupby(p).agg(["mean", "std", "count"]))

algae_test = pd.read_csv('algaeTest.csv', header=0)
```



EDA Comments: Having a general idea of the trends in the data, it appears there might be an instance of collinearity, so I should remove PO4 when it comes to cleaning it later. Otherwise, this seems alright. And of course I can see that certain chemical can contribute to algae growth as seen here.

b) Clean the data by identifying missing values and use the following strategy to estimate missing values: 1. For continuous (real-valued) data use the mean; 2. For discrete (integer valued) use the median; and 3. For categorical or nominal data use the mode. Then scale the predictor variables (subtract the mean and divide by the standard deviation) and code the categorical variables. Explain why you are doing this.

Reparameterize the response variable using a log transform and explain why you are doing this.

#drop unneeded vars

```
algae_train.drop(['a2', 'a3', 'a4', 'a5', 'a6', 'a7'], axis=1, inplace=True)
```

```
pd.isnull(algae_train).any()
```

```
algae_train['season']
```

```
0   winter
1   spring
2   autumn
3   spring
4   autumn
...
195  autumn
196  spring
197  autumn
198  winter
199  summer
Name: season, Length: 200, dtype: object
```

```
print(sum_stats)
```

```
//1.59998 0.0000000000000000 0.0000000000000000 1
[189 rows x 3 columns]
   a1      mean      std count
Chla
0.200  66.000000  NaN  1
0.300  47.600000 24.688999  4
0.400  31.800000  NaN  1
0.500  42.715667 24.978984  6
0.600  41.633333 12.287121  3
...
88.255  0.000000  NaN  ...
92.667  7.200000  NaN  1
93.683  12.300000  NaN  1
98.817  1.200000  NaN  1
110.456  0.000000  NaN  1
```

```
[131 rows x 3 columns]
# EGA: checking pairwise comparisons between each variable of interest
sns.pairplot(algae_train, vars= ['mxPH', 'mnO2', 'C1', 'NO3', 'NH4', 'PO4', 'Chla', 'a1' ])
```

```
#Performing analysis of categorical variables against the response
sns.boxplot(data=algae_train, x="season", y="a1", hue = 'size')
plt.figure()
sns.boxplot(data=algae_train, x="size", y="a1", hue = 'size')
plt.figure()
sns.boxplot(data=algae_train, x="speed", y="a1", hue = 'size')
```

```
# there are no discrete columns, and the categorical columns have no missing data
for col in predictors[3:]:
    algae_train[col] = algae_train[col].replace(np.nan, np.mean(algae_train[col]))
    algae_train[col] /= np.mean(algae_train[col])
    algae_train[col] /= np.std(algae_train[col])
season_vals = {'winter': 0, 'spring': 1, 'summer': 2, 'autumn': 3}
size_vals = {'small': 0, 'medium': 1, 'large': 2}
speed_vals = {'low': 0, 'medium': 1, 'high': 2}
algae_train['season'] = algae_train['season'].replace(season_vals)
algae_train['size'] = algae_train['size'].replace(size_vals)
algae_train['speed'] = algae_train['speed'].replace(speed_vals)
```

```
algae_train['a1'].loc[algae_train['a1']>0] = np.log(algae_train['a1'].loc[algae_train['a1']>0])
```

```
<ipython-input-9-ce9999891a6>:13: SettingWithCopyWarning:
  A value is trying to be set on a copy of a slice from a DataFrame
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy.
  algae_train['a1'].loc[algae_train['a1']>0] = np.log(algae_train['a1'].loc[algae_train['a1']>0])
```

```
pd.isnull(algae_train).any()
```

```
          season  size  speed  mxPH  mnO2   C1   NO3   NH4   oP04   P04   Chla   a1   a2   a3   a4   a5   a6   a7
0        0     0     1 -0.019710  0.287463 3.770155e-01  0.788574  0.039380  0.347256  0.250877 1.817928e+00  0.000000  0.0  0.0  34.2  8.3  0.0
1        1     0     1  0.568222 -0.470990 3.100179e-01 -0.532115 -0.067408  3.926568  0.327449 -6.933585e-01  0.336472  7.6  4.8  1.9  6.7  0.0  2.1
2        3     0     1  0.148270  0.961644 -7.94360e-01  0.546315 -0.079388  0.575746  0.384111  8.218554e-02 1.193922  53.6  1.9  0.0  0.0  0.0  9.7
3        1     0     1  0.097876 -1.819352 7.408577e-01 -0.261574 -0.206962 -0.137187  0.006389 -6.343127e-01  0.131402  41.0 18.9  0.0  1.4  0.0  1.4
4        3     0     1  0.081078 -0.049627 2.573017e-01  1.903289 -0.137386 -0.169012 -0.314805 -1.751483e-01  2.219203  2.9  7.5  0.0  7.5  4.1  1.0
...
195      3     2     1  0.652213 -0.302445 -5.768510e-01  0.146906 -0.214371 -0.220735 -0.450477 -5.869834e-01  2.541602  21.7  5.6  0.0  1.0  0.0  0.0
196      1     2     1  0.484232  0.624554 -6.439566e-01 -0.021982 -0.193022 -0.422968 -0.484118 -4.694576e-01  2.890372  7.0  1.7  0.0  4.8  10.3  1.0
197      3     2     1  0.316251 -0.892353 2.116488e+00 -0.081213 -0.226508  0.052431  0.018262 8.964706e-01  0.000000  15.9  2.4  1.0  0.0  0.0
198      0     2     1 -0.019710 -0.639535 1.560767e-16  0.000000  0.000000  0.000000 -8.93074e-17  0.000000  12.5  3.7  1.0  0.0  0.0  4.9
199      2     2     1  0.820193 -1.018762 8.614062e-01 -0.128704 -0.243472 -0.106032  0.020582 2.184211e-01  0.875469  10.5  9.0  7.8  0.0  0.0  5.8
200 rows x 18 columns
```

So the variables are scaled by the SD and mean to standardize the data. Effectively, doing so transforms all values into z-scores and makes it easier to identify outliers as well as each data point's position relative to the mean.

Additionally, one hot encoding enables the transformation of categorical variables into discrete variables and makes it possible to incorporate them into the pymc model in the next step.

Log transformations, as used on the response variable here, help to unskew the data and make its distribution more symmetrical, enabling stronger analysis.

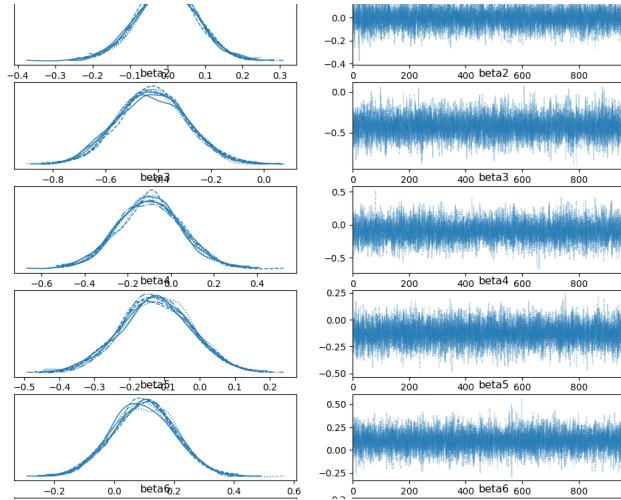
(c) Build a robust Bayesian regression model to predict the response and explain your choices for the priors.

graphics_path = '/Users/addyg/Desktop/UVA_Courses/DS6040 Bayes/'

```
with pm.Model() as robust:
    betas0 = pm.Normal('beta0', mu=0, sigma=10)
    betas1 = pm.Normal('beta1', mu=0, sigma=10)
    betas2 = pm.Normal('beta2', mu=0, sigma=10)
    betas3 = pm.Normal('beta3', mu=0, sigma=10)
    betas4 = pm.Normal('beta4', mu=0, sigma=10)
    betas5 = pm.Normal('beta5', mu=0, sigma=10)
    betas6 = pm.Normal('beta6', mu=0, sigma=10)
    betas7 = pm.Normal('beta7', mu=0, sigma=10)
    betas8 = pm.Normal('beta8', mu=0, sigma=10)
    betas9 = pm.Normal('beta9', mu=0, sigma=10)
    betas10 = pm.Normal('beta10', mu=0, sigma=10)
    betas11 = pm.Normal('beta11', mu=0, sigma=10)
```

```
o = pm.Uniform("o", 10**-5, 10)
v = pm.HalfNormal("v", 20)
```

```
ph = pm.MutableData("ph", algaes_train.iloc[:, 3])
mn = pm.MutableData("mn", algaes_train.iloc[:, 4])
cl = pm.MutableData("cl", algaes_train.iloc[:, 5])
```



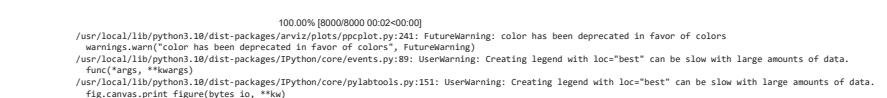
Axes: title='center': 'a1 / a1'

a1 / a1

```
with robust:
    trace_robust.extend(pm.sample_posterior_predictive(trace_robust))

fig, ax = plt.subplots()
ax.plot_ppc(trace_robust, color = 'b', observed_rug=True, ax=ax)
ax.axvline(algaes_train['a1'].mean(), ls="--", color="r", label="Observed mean")
ax.legend(fontsize=10);
```

```
az.plot_posterior(trace_robust)
```

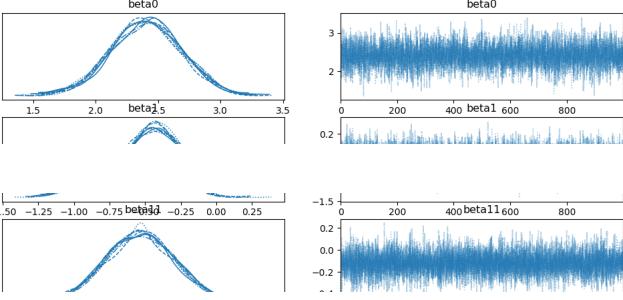


In the above plots, it's apparent that the Bayesian P-values that this robust model created distributions that run clearly to expected distribution of Bayes p-values. The posterior predictive plot demonstrates that the sampling worked well, since the observed and predicted values are similar.

(e) Plot the posterior distributions for the parameters and discuss your results. Plot the response variable vs. one of the predictors using counterfactuals and discuss this plot.

az.plot_posterior(trace_robust)

```
array([[<Axes: title='center': 'beta0'>,
       <Axes: title='center': 'beta0'>],
      [<Axes: title='center': 'beta1'>,
       <Axes: title='center': 'beta1'>],
      [<Axes: title='center': 'beta2'>,
       <Axes: title='center': 'beta2'>],
      [<Axes: title='center': 'beta3'>,
       <Axes: title='center': 'beta3'>],
      [<Axes: title='center': 'beta4'>,
       <Axes: title='center': 'beta4'>],
      [<Axes: title='center': 'beta5'>,
       <Axes: title='center': 'beta5'>],
      [<Axes: title='center': 'beta6'>,
       <Axes: title='center': 'beta6'>],
      [<Axes: title='center': 'beta7'>,
       <Axes: title='center': 'beta7'>],
      [<Axes: title='center': 'beta8'>,
       <Axes: title='center': 'beta8'>],
      [<Axes: title='center': 'beta9'>,
       <Axes: title='center': 'beta9'>],
      [<Axes: title='center': 'beta10'>,
       <Axes: title='center': 'beta10'>],
      [<Axes: title='center': 'beta11'>,
       <Axes: title='center': 'beta11'>],
      [<Axes: title='center': 'beta12'>,
       <Axes: title='center': 'beta12'>],
      [<Axes: title='center': 'beta13'>,
       <Axes: title='center': 'beta13'>],
      [<Axes: title='center': 'beta14'>,
       <Axes: title='center': 'beta14'>],
      [<Axes: title='center': 'beta15'>,
       <Axes: title='center': 'beta15'>],
      [<Axes: title='center': 'beta16'>,
       <Axes: title='center': 'beta16'>],
      [<Axes: title='center': 'beta17'>,
       <Axes: title='center': 'beta17'>],
      [<Axes: title='center': 'beta18'>,
       <Axes: title='center': 'beta18'>],
      [<Axes: title='center': 'beta19'>,
       <Axes: title='center': 'beta19'>],
      [<Axes: title='center': 'beta20'>,
       <Axes: title='center': 'beta20'>],
      [<Axes: title='center': 'beta21'>,
       <Axes: title='center': 'beta21'>],
      [<Axes: title='center': 'beta22'>,
       <Axes: title='center': 'beta22'>],
      [<Axes: title='center': 'beta23'>,
       <Axes: title='center': 'beta23'>],
      [<Axes: title='center': 'beta24'>,
       <Axes: title='center': 'beta24'>],
      [<Axes: title='center': 'beta25'>,
       <Axes: title='center': 'beta25'>]], dtype=object)
```



For this model, I chose to use normal priors for the beta distributions because the data has already been scaled in part b and a normal distribution is a reasonable expectation.

(d) Evaluate the sampling and choice of priors using Bayesian p-value and posterior predictive plots. Briefly explain your results.

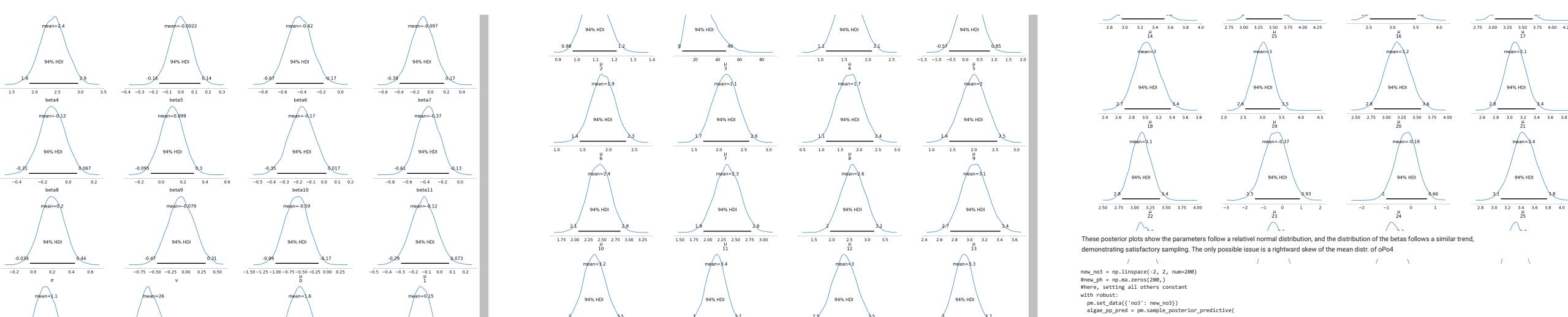
az.plot_ppv(trace_robust, figsize = (10,8))

```
array([[<Axes: title='center': 'beta0'>,
       <Axes: title='center': 'beta0'>],
      [<Axes: title='center': 'beta1'>,
       <Axes: title='center': 'beta1'>],
      [<Axes: title='center': 'beta2'>,
       <Axes: title='center': 'beta2'>],
      [<Axes: title='center': 'beta3'>,
       <Axes: title='center': 'beta3'>],
      [<Axes: title='center': 'beta4'>,
       <Axes: title='center': 'beta4'>],
      [<Axes: title='center': 'beta5'>,
       <Axes: title='center': 'beta5'>],
      [<Axes: title='center': 'beta6'>,
       <Axes: title='center': 'beta6'>],
      [<Axes: title='center': 'beta7'>,
       <Axes: title='center': 'beta7'>],
      [<Axes: title='center': 'beta8'>,
       <Axes: title='center': 'beta8'>],
      [<Axes: title='center': 'beta9'>,
       <Axes: title='center': 'beta9'>],
      [<Axes: title='center': 'beta10'>,
       <Axes: title='center': 'beta10'>],
      [<Axes: title='center': 'beta11'>,
       <Axes: title='center': 'beta11'>],
      [<Axes: title='center': 'beta12'>,
       <Axes: title='center': 'beta12'>],
      [<Axes: title='center': 'beta13'>,
       <Axes: title='center': 'beta13'>],
      [<Axes: title='center': 'beta14'>,
       <Axes: title='center': 'beta14'>],
      [<Axes: title='center': 'beta15'>,
       <Axes: title='center': 'beta15'>],
      [<Axes: title='center': 'beta16'>,
       <Axes: title='center': 'beta16'>],
      [<Axes: title='center': 'beta17'>,
       <Axes: title='center': 'beta17'>],
      [<Axes: title='center': 'beta18'>,
       <Axes: title='center': 'beta18'>],
      [<Axes: title='center': 'beta19'>,
       <Axes: title='center': 'beta19'>],
      [<Axes: title='center': 'beta20'>,
       <Axes: title='center': 'beta20'>],
      [<Axes: title='center': 'beta21'>,
       <Axes: title='center': 'beta21'>],
      [<Axes: title='center': 'beta22'>,
       <Axes: title='center': 'beta22'>],
      [<Axes: title='center': 'beta23'>,
       <Axes: title='center': 'beta23'>],
      [<Axes: title='center': 'beta24'>,
       <Axes: title='center': 'beta24'>],
      [<Axes: title='center': 'beta25'>,
       <Axes: title='center': 'beta25'>]], dtype=object)
```

beta0

beta2

beta3

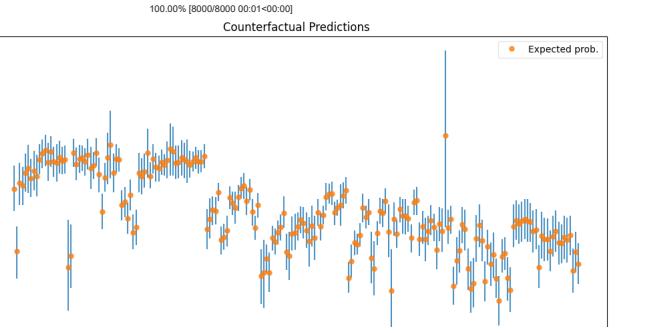


These posterior plots show the parameters follow a relative normal distribution, and the distribution of the betas follows a similar trend, demonstrating satisfactory sampling. The only possible issue is a rightward skew of the mean distr. of o|P4

```

new_no3 = np.linspace(-2, 2, num=200)
#new_ph = np.ma.zeros(200)
#Here, setting all others constant
with robust:
    pm.set_data('no3': new_no3)
    algae_pp_pred = pm.sample_posterior_predictive(

```



Though, this is not ideal, a general negative trend can be observed between the no3 Z-score and the algae z-score based on the counterfactual plot.

(f) The data set, algaetest.csv, contains a test set. Build a second model that may contain fewer predictors but should contain at least one nonlinear element. Evaluate both models using mean square error for the test set.

```

to_replace = {'high_': 'high', 'small_': 'small', 'large_': 'large', 'low_':
             'low', 'XXXXXX': '0.0000'}

```

```

algae_test = algae_test.replace(to_replace)

```

```

algae_test.columns = algae_train.columns

i = 0
to_drop = []
for i in range(len(algae_test['NO3'])):
    if algae_test.iloc[i, 6] == '1.0590011160.599611435.000001690.00000' or (algae_test.iloc[i, 6] == '1.598006249.60010' or (algae_test.iloc[i, 6] ==
new_no3 = np.linspace(-2, 2, num=200)
#new_ph = np.ma.zeros(200)
#Here, setting all others constant
with robust:
    pm.set_data('no3': new_no3)
    algae_pp_pred = pm.sample_posterior_predictive(

```

```

preds2 = list(algae_test.columns[0:11])

for col in predictors2[3:]:
    print(col)
    idx = predictors2.index(col)
    for i in range(len(algae_test)):
        algae_test.ioc[i, idx] = float(algae_test.ioc[i, idx])
    algae_test[col] = algae_test[col].replace(np.nan, np.mean(algae_test[col]))
    algae_test[col] /= np.std(algae_test[col])

algae_test['a1'].loc[algae_test['a1']>0] = np.log(algae_test['a1'].loc[algae_test['a1']>0])

algae_test['a1'].loc[algae_test['a1']>0] = np.log(algae_test['a1'].loc[algae_test['a1']>0])

```

Based on everything that I've seen in these models, their p-values, predictive performance MSE, it seems that that these would do a decent job in predicting the likelihood of algae growth, given each of the predictor variables and outcomes. In order to better deal with these algae blooms, the linear model appears to do a better job, and I would recommend the EU use such a model when monitoring water content.

	season	size	speed	mxPH	mn02	C1	NO3	NH4	o P4	P04	Chla	a1	a2	a3	a4	a5	a6	a7		
0	winter	small	medium	0.69643	-0.25018	0.500401	1.847195	0.575631	-0.314827	0.30931	-0.23548	0.182322	0.0	0.0	0.0	23.2	46.4	0.0		
1	summer	small	medium	0.091987	-0.001242	0.996222	-0.389696	-0.000947	-0.101083	0.059933	-0.428343	1.945910	23.0	6.5	1.4	21.2	0.0	2.1		
2	spring	small	high	0.484759	-0.438137	0.708331	0.050274	1.45466	-0.140185	0.040279	-0.524772	0.336472	38.2	2.4	0.0	4.8	1.0	1.2		
3	spring	small	medium	0.204208	-0.467225	-1.195141	-0.228159	-0.717077	-0.696988	-0.613783	1.360977	55.4	8.4	0.0	0.0	0.0	0.0	0.0		
4	summer	small	medium	0.507204	1.30087	-0.614769	0.849421	-0.818174	-0.731794	0.999832	0.26149	3.346389	2.4	0.0	0.0	0.0	0.0	0.0		
...		
134	summer	medium	high	0.226652	0.407871	-0.74041	-0.907453	-0.722951	-0.392811	-0.597816	-0.598948	2.895912	1.7	2.0	0.0	1.7	5.9	0.0		
135	winter	large	low	0.574536	0.689874	0.689874	-0.067968	1.366891	-0.207637	0.717794	-1.107721	2.603597	0.995310	3.9	2.1	0.0	3.9	4.6	2.3	
136	winter	large	low	0.877532	1.112878	-0.408073	0.354425	-0.425942	-0.246145	0.255449	0.498856	0.000000	4.7	0.0	0.0	2.6	2.6	0.0		
137	summer	large	low	0.204208	-0.532137	-0.286333	-0.040683	-0.675801	-0.091620	-0.202824	0.274844	0.000000	12.0	1.7	0.0	2.7	0.0	0.0		
138	autumn	large	low	0.484759	0.830875	-0.381817	0.097193	-0.504183	-0.201995	-0.094596	0.53268	0.530628	7.0	1.2	0.0	4.8	3.1	0.0		

132 rows × 18 columns

```

algae_test['mxPH'] = np.square(algae_test['mxPH'])

with pm.Model() as robust_test:
    beta0 = pm.Normal('beta0', mu=0, sigma=10)
    beta1 = pm.Normal('beta1', mu=0, sigma=10)
    beta2 = pm.Normal('beta2', mu=0, sigma=10)
    beta3 = pm.Normal('beta3', mu=0, sigma=10)
    beta4 = pm.Normal('beta4', mu=0, sigma=10)
    beta5 = pm.Normal('beta5', mu=0, sigma=10)
    beta6 = pm.Normal('beta6', mu=0, sigma=10)

```

```

o = pm.HalfNormal('o', 50)
o = pm.Uniform('o', 10**-5, 10)
v = pm.HalfNormal('v', 20)

```

```

pm = pm.MutableData("ph", algae_test.iloc[:, 3])

```

```

mn = pm.MutableData("mn", algae_test.iloc[:, 4])
c1 = pm.MutableData("c1", algae_test.iloc[:, 5])
nh4 = pm.MutableData("nh4", algae_test.iloc[:, 7])
p04 = pm.MutableData("p04", algae_test.iloc[:, 9])
chla = pm.MutableData("chla", algae_test.iloc[:, 10])

p = pm.Deterministic("p", beta0 + beta1*ph + beta2*mn +
                     + beta3*c1 + beta4*nh4 + beta5*p04 + beta6*chla)

alg_test = pm.StudentT('a1', mu=p, sigma=o, nu=v, observed=algae_test['a1'])
trace_test = pm.sample(cores = 8, return_inferencedata = False, random_seed=1234)
trace_robust_test = pm.to_inference_data(trace=trace_test, log_likelihood=True)
robust_pp_test = pm.sample_posterior_predictive(trace_robust_test, extend_inferencedata=True,
                                                random_seed=1234)
test_y_pred = robust_pp_test['posterior_predictive']['a1']

100.00% [16000/16000 00:14:00] Sampling 8 chains, 0 divergences
100.00% [8000/8000 00:02:00]

```

#Evaluating MSE

```

from sklearn.metrics import mean_squared_error

```

```

pred = y_pred_test.mean(dim=['chain', 'draw'])
mse = mean_squared_error(pred, algae_train['a1'])
mse

```

```

1.2065609866571974

```

```

pred_test = test_y_pred.mean(dim=['chain', 'draw'])
mse_test = mean_squared_error(pred_test, algae_test['a1'])
mse_test

```

```

0.2711914607782464

```

(g) Summarize your results and give your answer to the EU's goal.

```

import pymc as pm
import matplotlib.pyplot as plt
import arviz as az
import pandas as pd
from scipy import special, stats
import numpy as np
import seaborn as sns

```

```

from google.colab import files
uploaded = files.upload()

```

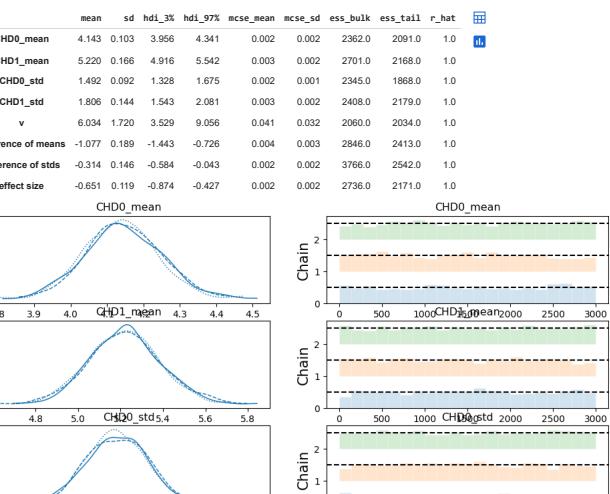
Choose File: CHDdata.csv
 • CHDdata.csv(text/csv) - 21018 bytes, last modified: 10/24/2023 - 100% done
 Saving CHDdata.csv to CHDdata.csv

```
chd = pd.read_csv('CHDdata.csv', header=0)
```

1. (10) Based on the analysis we did in class comparing systolic blood pressure among men in the South African study [1], you are now asked to extend this analysis by comparing low density lipoprotein (ldl) levels between the same two groups in the study: those with and without coronary heart disease. Perform this analysis and comment on whether you find a statistically significant difference between the groups for this variable. To do this analysis you will use the CHD data set (CHDdata.csv).

```
summary_stats = (chd.loc[:, ["chd", "ldl"]]
                 .groupby("chd")
                 .agg(["mean", "std", "count"]))
summary_stats
```

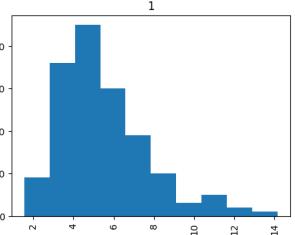
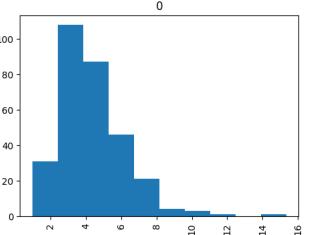
```
ldl_df = chd[["ldl", "chd"]]
ldl_df.hist("ldl", by="chd", figsize=(12, 4))
```



```

y0 = ldl_df.loc[ldl_df["chd"]==0]["ldl"]
y1 = ldl_df.loc[ldl_df["chd"]==1]["ldl"]

```



```

CHD0_std = pm.Uniform("CHD0_std", lower=0_low, upper=0_high)
v = pm.Exponential("v", 1/25)
CHD0 = pm.StudentT("No_CHD", nu_nv, mu=CHD0_mean, sigma=CHD0_std, observed=y0)
CHD1 = pm.StudentT("CHD", nu_nv, mu=CHD1_mean, sigma=CHD1_std, observed=y1)
diff_of_means = pm.Deterministic("difference of means", CHD0_mean - CHD1_mean)
diff_of_stds = pm.Deterministic("difference of stds", CHD0_std - CHD1_std)
effect_size = pm.Deterministic(
    "effect size", diff_of_means / np.sqrt((CHD0_std ** 2 + CHD1_std ** 2) / 2)
)

chd_trace=pm.sample(random_seed = random_seed, cores = cores)

```

100.00% [6000/6000 00:12<00:00 Sampling 3 chains, 0 divergences]

```

az.plot_trace(chd_trace, kind = "rank_bars")
az.summary(chd_trace)

```

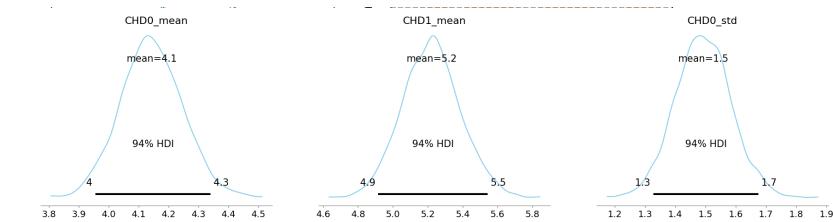
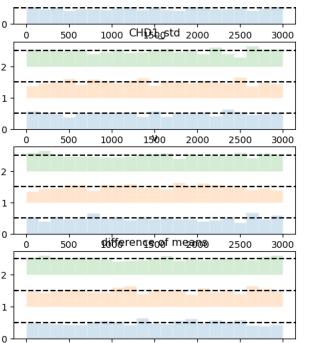
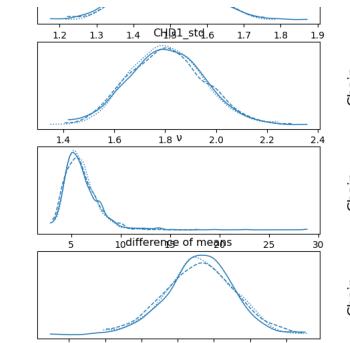
```

random_seed = 100
cores = 3

μ_prior = chd['ldl'].mean()
σ_prior = chd['ldl'].std() * 100

with pm.Model() as model:
    CHD0_mean = pm.Normal("CHD0_mean", mu=μ_prior, sigma=σ_prior)
    CHD1_mean = pm.Normal("CHD1_mean", mu=μ_prior, sigma=σ_prior)
    CHD0_std = pm.Uniform("CHD0_std", lower=0_low, upper=0_high)

```



```

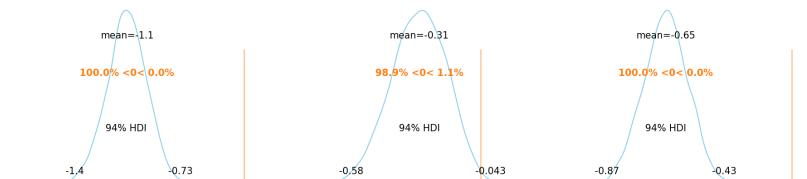
az.plot_posterior(
    chd_trace,
    var_names=["CHD0_mean", "CHD1_mean", "CHD0_std", "CHD1_std", "v"],
    color="#87ceeb",
);
az.plot_forest(chd_trace, var_names = ['CHD0_mean','CHD1_mean'])

```

```

array([<Axes: titles: '94.0% HDI'>], dtype=object)
difference of means

```



```

#chd_trace
chd_pp = pm.sample_posterior_predictive(chd_trace, model = model)
az.plot_pp(chd_pp)

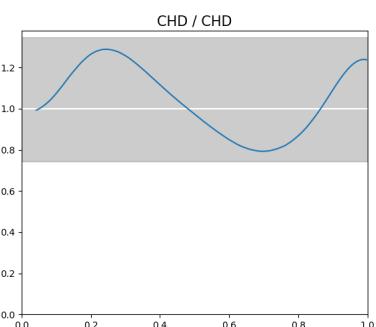
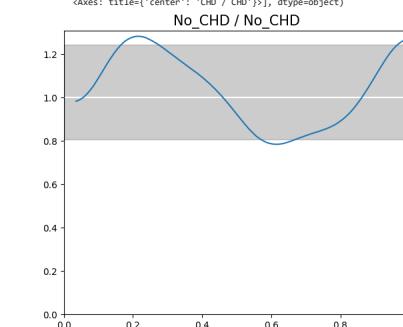
```

!!!!!!! Note: I ran very similar code as the prof and these plots were generated once, but it seems that pymc was updated or something and now I can't create the plots when I reran this and I haven't been able to resolve the issue

```

array([<Axes: title='No_CHD / No_CHD'>,
       <Axes: title='CHD / CHD'>], dtype=object)

```

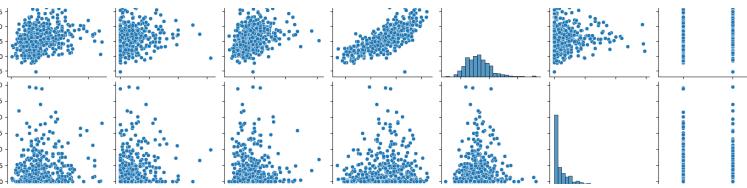
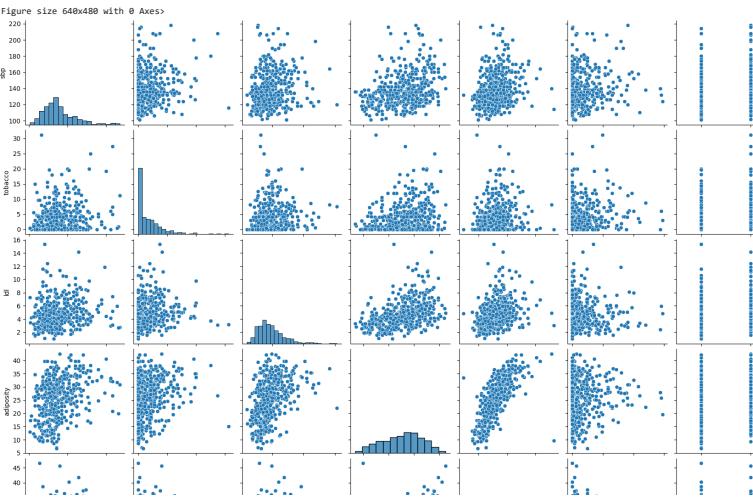


Based on the various plots above, there is a clear difference in likelihood to develop CHD between groups with lower and higher LDL content. In particular, the 94% HDI of the posterior difference of means shows that group with higher mean LDL content is more likely to have CHD. The same can be seen with the secondary 94%HDI set of graphs, and the subsequent difference of means graph. The effect size, in addition to the difference of means, seems to show the same general trends as well.

3. (30) The South African Health Ministry is now asking you to extend your work from problem 1 and perform a full analysis of the data in the coronary heart disease study [1]. Their goal is to understand factors that may be associate with coronary heart disease to inform further more focused research on those factors. For this work you will again use the CHD data set (CHDdata.csv). Perform the following steps and answer the associated questions.

(a) Perform a simple exploratory analysis of the data by at least viewing summary statistics. Comment on what the EDA implies for your analysis.

```
#EDA: checking how each continuous var is compared against each other continuous var and also the response variable CHD
# even though chd is discrete I think its a good idea to see how each variable affects it
sns.pairplot(chd_vars = ['sbp', 'tobacco', 'ldl', 'adiposity', 'obesity', 'alcohol', 'chd'])
#Performing analysis of categorical variables against the response
sns.boxplot(data=chd, x="famhist", y="chd")
plt.figure()
```



Based on the figures above, its clear that there's some correlations between high sbp levels, tobacco use, ldl, obesity and alcohol use. But at the same time there doesn't appear to be any collinearity so high that any predictor needs to be dropped.

(b) Preprocess the data by scaling the predictor variables using a z-transformation (subtract the mean divide by the standard deviation). Explain why this is appropriate. Code the categorical variables and explain your coding

```
preds_num = ['sbp', 'tobacco', 'ldl', 'adiposity', 'typea', 'obesity', 'alcohol', 'age']
for p in preds_num:
    chd[p] -= np.mean(chd[p])
    chd[p] /= np.std(chd[p])
chd = pd.get_dummies(chd, columns=['famhist'])
fam_vals = {'Absent': 0, 'Present': 1}
chd['famhist'] = chd['famhist'].replace(fam_vals)
```

chd

	sbp	tobacco	ldl	adiposity	famhist	typea	obesity	alcohol	age	chd
0	0.058564	1.823073	0.478412	-0.295503	1	-0.418470	-0.176786	3.277738	0.629336	1
1	0.277089	-0.790237	-0.159680	0.412140	0	0.193344	0.671373	-0.612745	1.383115	1
2	-0.992806	-0.774980	-0.609245	0.884332	1	-0.112563	0.735519	-0.541183	0.218184	0
3	1.546985	0.842264	0.807126	1.624141	1	-0.214532	1.412621	0.295062	1.040488	1
4	-0.211332	2.171805	-0.595977	0.305351	1	0.703189	-0.012856	1.647775	0.423760	1
...
457	3.696039	-0.705234	0.595263	0.812281	0	1.111065	0.571590	-0.696983	1.040488	0
458	2.133091	0.123004	-0.159680	0.861173	0	-0.112563	0.609602	0.068519	0.629336	1
459	-1.481228	-0.138545	-1.522877	-1.309364	0	-1.338191	-1.414575	0.392385	0.834912	0

Here, scaling the data makes the distribution appear normal and thereby we can avoid the need to use the individual means for each column of the data. Encoding the categorical variables allows for numerical analysis to be run on those variables.

462 rows × 10 columns

(c) Develop and use Bayesian GLM model with all the predictor variables. Explain your choice of priors and their parameters. Show the graphical representation of your model.

```
with pm.Model() as chd_lin_model:
    beta0 = pm.Normal('beta0', mu=0, sigma=10)
    beta1 = pm.Normal('beta1', mu=0, sigma=10)
    beta2 = pm.Normal('beta2', mu=0, sigma=10)
    beta3 = pm.Normal('beta3', mu=0, sigma=10)
    beta4 = pm.Normal('beta4', mu=0, sigma=10)
    beta5 = pm.Normal('beta5', mu=0, sigma=10)
    beta6 = pm.Normal('beta6', mu=0, sigma=10)
    beta7 = pm.Normal('beta7', mu=0, sigma=10)
    beta8 = pm.Normal('beta8', mu=0, sigma=10)
    beta9 = pm.Normal('beta9', mu=0, sigma=10)
```

```
μ = pm.Deterministic("μ", beta0 + beta1*chd.iloc[:, 0] + beta2*chd.iloc[:, 1] + beta3*chd.iloc[:, 2] + beta4*chd.iloc[:, 3]
                     + beta5*chd.iloc[:, 4] + beta6*chd.iloc[:, 5] + beta7*chd.iloc[:, 6] + beta8*chd.iloc[:, 7] + beta9*chd.iloc[:, 8])

chd_data = chd.iloc[:, 9]

chd_pre = pm.Bernoulli('chd_pre', p = pm.math.sigmoid(μ))

chd_post = pm.Bernoulli('chd_post', p = pm.math.sigmoid(μ), observed=chd_data)

trace_glm = pm.sample(1000, cores=8, random_seed=1234, return_inferencedata=True)
glm_trace = pm.to_inference_data(trace_glm, log_likelihood=True)
glm_pp = pm.sample_posterior_predictive(glm_trace,
                                         var_names=["chd_pre", "chd_post"],
                                         random_seed = 1234)

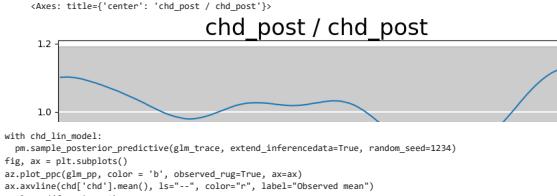
g3=pm.model_to_graphviz(chd_lin_model)
graphics_path = '/content/drive/MyDrive/'
g3.render(graphics_path+"chd_lin_model", format='png', cleanup=True)
g3
```

In this linear model, I model the correlation of 9 variables with the outcome chd_data. Each variable is associated with a coefficient (beta) that describes its relationship with the outcome. A deterministic function computes the expected value (μ) based on a linear combination of these variables and their respective coefficients. The outcome is binary, and hence a Bernoulli distribution is used for both the prior prediction (chd_pre) and the posterior likelihood (chd_post), with a logistic sigmoid function linking the linear combination to a probability. Sampling provides posterior estimates for all model parameters.

(d) Evaluate both the sampling used by your model, as well as, the prior and posterior predictive results.

```
with chd_lin_model:
    chd_prior = pm.sample_prior_predictive()
az.plot_dist(chd_prior.prior['chd_pre'])
```

1 / \ chd post / \ chd pre \ \



100.00% [8000/8000 00:02<00:00]
/usr/local/lib/python3.10/dist-packages/arviz/plots/ppcplot.py:241: FutureWarning: color has been deprecated in favor of colors
warnings.warn("color has been deprecated in favor of colors", FutureWarning)

Based on the above plots, we can see that the model has p-values that lie in the accepted range, demonstrating that the model is satisfactory. The posterior predictive sampling, with similarities between the observed and predicted parameters, supports the validity of the model.

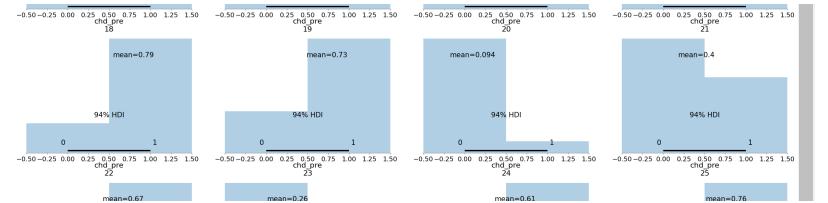
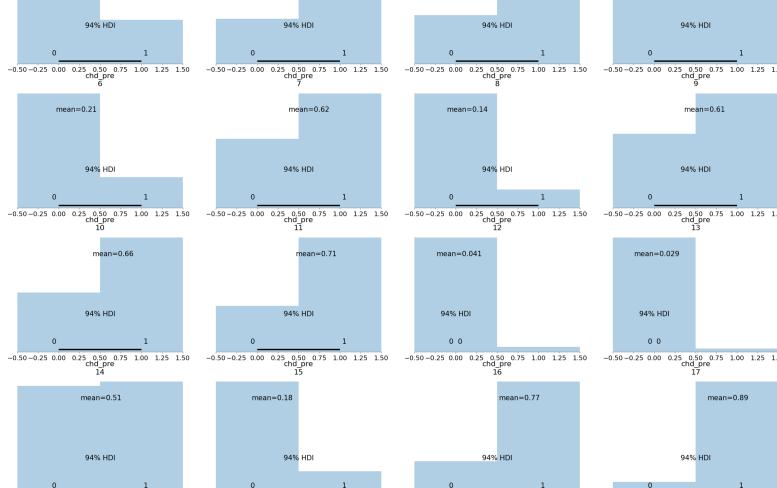
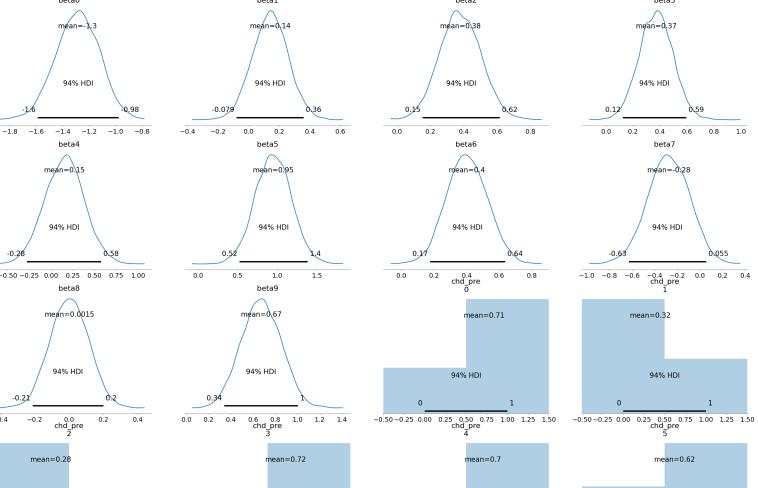
(e) Provide plots of the posterior distributions for the parameters and discuss what these results imply for the relevant predictor variables and the overall model. Show posterior analyses to include relevant odds ratio(s) and describe what this analysis means for the goal of your study.

az.plot_posterior(glm_trace)

```

/usr/local/lib/python3.10/dist-packages/arviz/plots/plot_utils.py:271: UserWarning: rcParams['plot.max_subplots'] (40) is smaller than the number of vari
array([{'Axes': titles['center': 'beta0'], 'x_min': -1.6, 'x_max': 0.98, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'beta1'], 'x_min': -0.079, 'x_max': 0.36, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'beta2'], 'x_min': 0.15, 'x_max': 0.62, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'beta3'], 'x_min': 0.12, 'x_max': 0.59, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'beta4'], 'x_min': -0.15, 'x_max': 0.15, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'beta5'], 'x_min': 0.5, 'x_max': 1.4, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'beta6'], 'x_min': -0.4, 'x_max': 0.4, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'beta7'], 'x_min': -1.0, 'x_max': 0.28, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'beta8'], 'x_min': -0.26, 'x_max': 0.58, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'beta9'], 'x_min': 0.5, 'x_max': 1.4, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'chd_pre1'], 'x_min': -0.5, 'x_max': 1.0, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'chd_pre2'], 'x_min': -0.3, 'x_max': 0.6, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'chd_pre3'], 'x_min': -0.17, 'x_max': 0.64, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'chd_pre4'], 'x_min': -1.0, 'x_max': 0.4, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'chd_pre5'], 'x_min': -0.26, 'x_max': 0.58, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'chd_pre6'], 'x_min': 0.5, 'x_max': 1.4, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'chd_pre7'], 'x_min': -0.4, 'x_max': 0.4, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'chd_pre8'], 'x_min': -1.0, 'x_max': 0.28, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'chd_pre9'], 'x_min': -0.23, 'x_max': 0.2, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'chd_pre10'], 'x_min': 0.34, 'x_max': 1.4, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'chd_pre11'], 'x_min': -0.5, 'x_max': 1.0, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'chd_pre12'], 'x_min': -0.3, 'x_max': 0.6, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'chd_pre13'], 'x_min': -0.17, 'x_max': 0.64, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'chd_pre14'], 'x_min': -1.0, 'x_max': 0.4, 'y_min': 0.0, 'y_max': 1.0}, {'Axes': titles['center': 'chd_pre15'], 'x_min': -0.26, 'x_max': 0.58, 'y_min': 0.0, 'y_max': 1.0}, {"Axes": titles["center": "chd_pre16"], "x_min": 0.5, "x_max": 1.4, "y_min": 0.0, "y_max": 1.0}, {"Axes": titles["center": "chd_pre17"], "x_min": -0.4, "x_max": 0.4, "y_min": 0.0, "y_max": 1.0}, {"Axes": titles["center": "chd_pre18"], "x_min": -0.5, "x_max": 1.0, "y_min": 0.0, "y_max": 1.0}, {"Axes": titles["center": "chd_pre19"], "x_min": -0.3, "x_max": 0.6, "y_min": 0.0, "y_max": 1.0}, {"Axes": titles["center": "chd_pre20"], "x_min": -0.5, "x_max": 1.0, "y_min": 0.0, "y_max": 1.0}, {"Axes": titles["center": "chd_pre21"], "x_min": -0.3, "x_max": 0.6, "y_min": 0.0, "y_max": 1.0}, {"Axes": titles["center": "chd_pre22"], "x_min": -0.5, "x_max": 1.0, "y_min": 0.0, "y_max": 1.0}, {"Axes": titles["center": "chd_pre23"], "x_min": -0.3, "x_max": 0.6, "y_min": 0.0, "y_max": 1.0}, {"Axes": titles["center": "chd_pre24"], "x_min": -0.5, "x_max": 1.0, "y_min": 0.0, "y_max": 1.0}, {"Axes": titles["center": "chd_pre25"], "x_min": -0.3, "x_max": 0.6, "y_min": 0.0, "y_max": 1.0}, {"Axes": titles["center": "chd_pre26"], "x_min": -0.5, "x_max": 1.0, "y_min": 0.0, "y_max": 1.0}, {"Axes": titles["center": "chd_pre27"], "x_min": -0.3, "x_max": 0.6, "y_min": 0.0, "y_max": 1.0}, {"Axes": titles["center": "chd_pre28"], "x_min": -0.5, "x_max": 1.0, "y_min": 0.0, "y_max": 1.0}, {"Axes": titles["center": "chd_pre29"], "x_min": -0.3, "x_max": 0.6, "y_min": 0.0, "y_max": 1.0}], dtype=object)

```



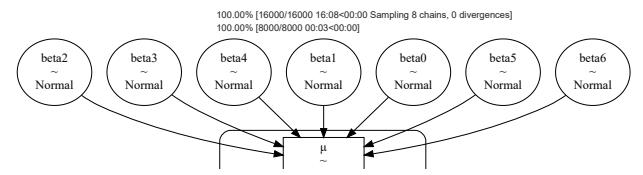
Based on the above plots, it's clear that the explanatory variables show a normal distribution, so Bernoulli distribution for the prior was a good choice. The HDI charts produced seem to indicate that the odds ratios for the parameters are acceptable. Beta are also normally distributed, showing that there shouldn't have been an issue in that area either.

In terms of the goals of the study, this helps us to see that the model being produced is accurate and should help our stakeholders achieve their goals.

```

_, ax = plt.subplots(figsize=(12, 6))
for i in range(1, 26):
    ax = plt.subplot(4, 6, i)
    chd_be = chd
    chd_be['age'] = np.square(chd_be['age'])
    OR = np.mean(chd_be['age'])
    lb, ub = np.percentile(OR, 2.5), np.percentile(OR, 97.5)
    plt.hist(np.exp(OR), bins=10, density=True)
    plt.axvline(np.exp(lb), color = 'r')
    plt.axvline(np.exp(ub), color = 'r')

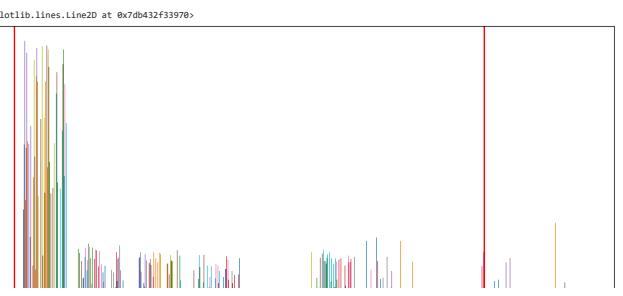
```



I kept the betas as normal, except for the use of a half normal on age, which was squared to become the nonlinear element. All other parameters remained the same for consistency.

(g) For the second model evaluate the sampling and choice of priors using Bayesian p-value and posterior predictive plots. Briefly explain your results.

```
az.plot_bpv(glm_pp2, figsize = (10,8))
```



The observed posterior values tend to be somewhat greater than the choice prior values, so it may be that the priors are underestimating the prevalence of coronary heart disease among South African men.

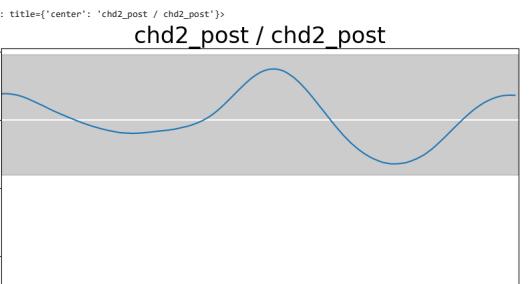
(f) Develop a second model to help in your study of this problem. This new model can use fewer predictor variables but must contain at least one nonlinear component. Explain your choice of priors and their parameters, and show the graphical representation of your new model.

```

chd2 = chd
chd2['age'] = np.square(chd2['age'])

with pm.Model() as glm2:
    beta0 = pm.Normal('beta0', mu=0, sigma=10)

```



```

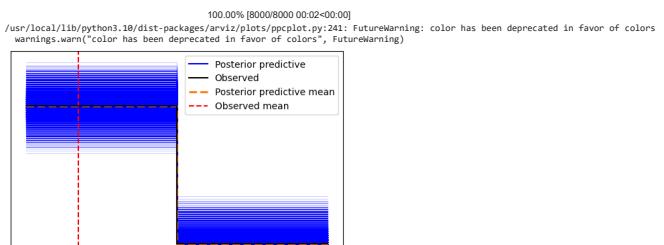
    beta1 = pm.Normal('beta1', mu=0, sigma=10)
    beta2 = pm.Normal('beta2', mu=0, sigma=10)
    beta3 = pm.Normal('beta3', mu=0, sigma=10)
    beta4 = pm.Normal('beta4', mu=0, sigma=10)
    beta5 = pm.Normal('beta5', mu=0, sigma=10)
    beta6 = pm.Normal('beta6', mu=0, sigma=10)

    p = pm.Deterministic('p', beta0 + beta1*chd2.iloc[:, 8] + beta2*chd2.iloc[:, 1] + beta3*chd2.iloc[:, 2]
                         + beta4*chd2.iloc[:, 3] + beta5*chd2.iloc[:, 6] + beta6*chd2.iloc[:, 8])

    chd2_data = chd2.iloc[:, 9]
    chd2_pre = pm.Bernoulli('chd2_pre', p=pm.math.sigmoid(p))
    chd2_post = pm.Bernoulli('chd2_post', p=pm.math.sigmoid(p), observed=chd2_data)

    trace_glm2 = pm.sample(1000, cores=8, random_seed=1234, target_accept=.99, return_inferencedata=False)
    glm_trace2 = pm.to_inference_data(trace_glm2, log_likelihood=True)
    glm_pp2 = pm.sample_posterior_predictive(glm_trace2, var_names=['chd2_pre', 'chd2_post'],
                                             random_seed = 1234)
    g3f=pm.model_to_graphviz(glm2)
    graphics_path = '/content/drive/MyDrive/'
    g3f.render(graphics_path+'GLM_Model', format='png', cleanup=True)
    g3f

```

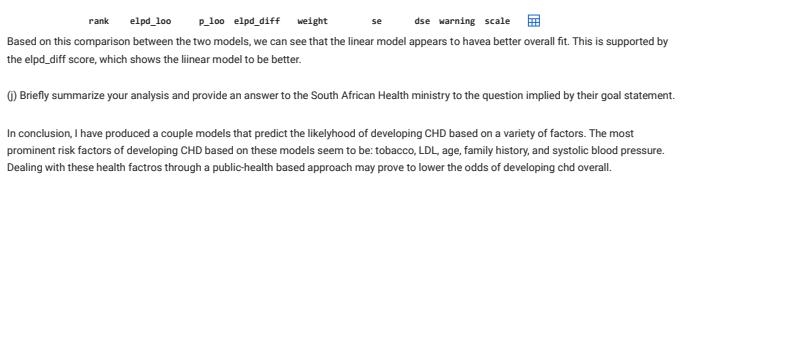
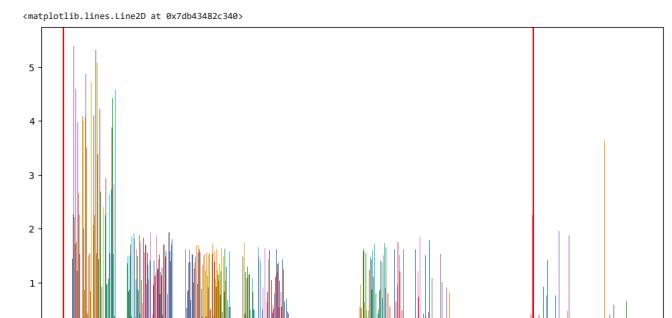
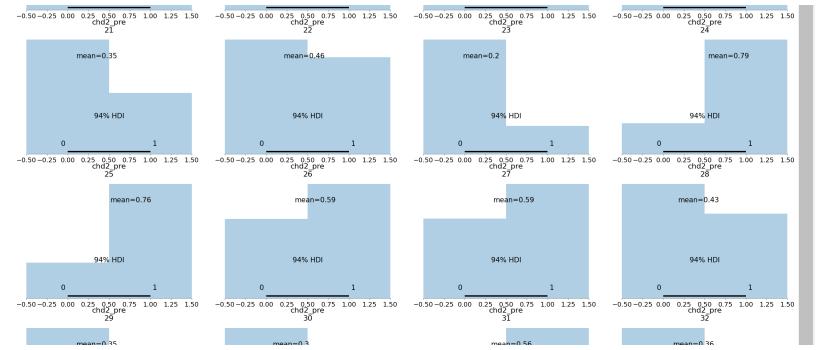
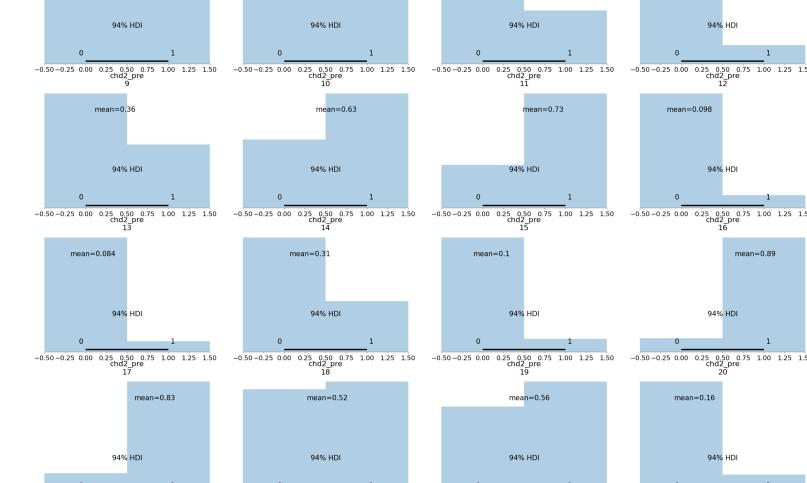
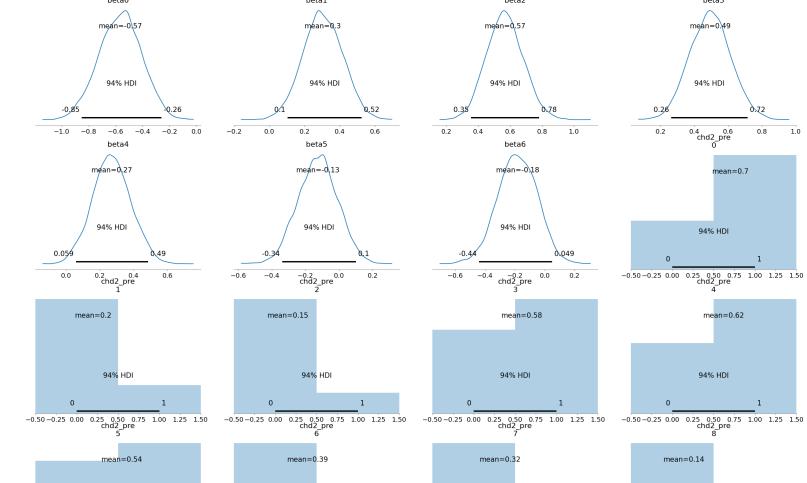


This model appears to be pretty decent. The produced Bayesian p-value is very close to 1, which appears satisfactory. Moreover, there are two clusters present within the posterior predictive plot, so it seems that the priors and parameters were well chosen.

(h) For the second model provide plots of the posterior distributions for the parameters and discuss what these results imply for the relevant predictor variables and the overall model. Show posterior analyses to include relevant odds ratio(s) and describe what this analysis means for the goal of your study.

```
az.plot_posterior(glm_trace2)
```

```
/usr/local/lib/python3.10/dist-packages/arviz/plots/plot_utils.py:271: UserWarning: rcParams['plot.max_subplots'] (40) is smaller than the number of vari
array([{'Axes': title='beta0', 'center': 'beta0'}, {'Axes': title='beta1', 'center': 'beta1'}, {'Axes': title='beta2', 'center': 'beta2'}, {'Axes': title='beta3', 'center': 'beta3'}, {'Axes': title='beta4', 'center': 'beta4'}, {'Axes': title='beta5', 'center': 'beta5'}, {'Axes': title='beta6', 'center': 'beta6'}, {"Axes": title="beta7", "center": "beta7"}, {"Axes": title="beta8", "center": "beta8"}, {"Axes": title="chd2_pre\n0", "center": "chd2_pre0"}, {"Axes": title="chd2_pre\n1", "center": "chd2_pre1"}, {"Axes": title="chd2_pre\n2", "center": "chd2_pre2"}, {"Axes": title="chd2_pre\n3", "center": "chd2_pre3"}, {"Axes": title="chd2_pre\n4", "center": "chd2_pre4"}, {"Axes": title="chd2_pre\n5", "center": "chd2_pre5"}, {"Axes": title="chd2_pre\n6", "center": "chd2_pre6"}, {"Axes": title="chd2_pre\n7", "center": "chd2_pre7"}, {"Axes": title="chd2_pre\n8", "center": "chd2_pre8"}, {"Axes": title="chd2_pre\n9", "center": "chd2_pre9"}, {"Axes": title="chd2_pre\n10", "center": "chd2_pre10"}, {"Axes": title="chd2_pre\n11", "center": "chd2_pre11"}, {"Axes": title="chd2_pre\n12", "center": "chd2_pre12"}, {"Axes": title="chd2_pre\n13", "center": "chd2_pre13"}, {"Axes": title="chd2_pre\n14", "center": "chd2_pre14"}, {"Axes": title="chd2_pre\n15", "center": "chd2_pre15"}, {"Axes": title="chd2_pre\n16", "center": "chd2_pre16"}, {"Axes": title="chd2_pre\n17", "center": "chd2_pre17"}, {"Axes": title="chd2_pre\n18", "center": "chd2_pre18"}, {"Axes": title="chd2_pre\n19", "center": "chd2_pre19"}, {"Axes": title="chd2_pre\n20", "center": "chd2_pre20"}, {"Axes": title="chd2_pre\n21", "center": "chd2_pre21"}, {"Axes": title="chd2_pre\n22", "center": "chd2_pre22"}, {"Axes": title="chd2_pre\n23", "center": "chd2_pre23"}, {"Axes": title="chd2_pre\n24", "center": "chd2_pre24"}, {"Axes": title="chd2_pre\n25", "center": "chd2_pre25"}, {"Axes": title="chd2_pre\n26", "center": "chd2_pre26"}, {"Axes": title="chd2_pre\n27", "center": "chd2_pre27"}, {"Axes": title="chd2_pre\n28", "center": "chd2_pre28"}, {"Axes": title="chd2_pre\n29", "center": "chd2_pre29"}, {"Axes": title="chd2_pre\n30", "center": "chd2_pre30"}, {"Axes": title="chd2_pre\n31", "center": "chd2_pre31"}, {"Axes": title="chd2_pre\n32", "center": "chd2_pre32"}], dtype=object)
```



```
_, ax = plt.subplots(figsize=(12, 6))
b = glm_trace2.posterior['chd2_pre']
OR = np.mean(b, axis=0)
lb = np.percentile(OR, 2.5), np.percentile(OR, 97.5)
plt.hist(np.exp(OR), bins=10, density=True)
plt.axvline(np.exp(lb), color = 'r')
plt.axvline(np.exp(ub), color = 'r')
```

```
(i) Compare the two models using PSIS - leave-one-out cross validation and WAIC.

az.compare(("Linear": glm_trace, 'Nonlinear':glm_trace2))
```