

# BNN for Multiclass Classification

## General Principles

Building upon the [Binary Classification BNN](#), the **BNN Multiclass Classification** model can handle dependent variables with  $K > 2$  discrete categories.

Instead of the final layer returning a single output, the final layer in a multiclass BNN returns a  $K$ -dimensional vector of scores (logits) for each observation. To transform these continuous scores into valid probabilities that sum to 1 across all  $K$  classes, we apply the **softmax** activation function. Finally, the categorized predictions are evaluated using a **Categorical** likelihood.

## Considerations

### Note

- **Output Layer Dimensions:** While binary classification network predictions can be compressed to a single output logit per observation, multiclass networks **MUST** output exactly  $K$  dimensions in their final layer, matching the number of target classes.
- **The Softmax Simplex:** Applying the softmax function across the final layer's logits guarantees that the resulting outputs form a probability simplex . This is biologically similar to independent Poisson rates strictly normalizing to fixed categorical ratios.
- **Likelihood Function:** After calculating the probabilities with softmax, we use a **Categorical** distribution as the final likelihood, matching the integer index of the observed category.
- **Improved Calibration:** Multiclass BNNs greatly reduce out-of-distribution overconfidence. Standard deep learning cross-entropy models will often assign  $>99\%$  probability to an unseen class purely due to the exponential nature of softmax. In a BNN, exploring the posterior width of the parameters yields “flat” unconfident probability profiles over  $K$  classes when the input is outside the training distribu-

tion.

## Example

Below is an example code snippet demonstrating a *Bayesian Neural Network for multiclass classification* using the Bayesian Inference (**BI**) package. This example generates a synthetic  $K = 3$  cluster dataset.

## Python

```
from BI import bi
import jax.numpy as jnp
import jax

# Setup device-----
m = bi(platform='cpu')

# Generate Synthetic Data -----
# 3 classes based on a random normal distribution split
key = jax.random.PRNGKey(42)
X = jax.random.normal(key, (300, 2))
# Rule: Q1=Class 0, Q2/Q3=Class 1, Q4=Class 2
Y = jnp.where(X[:, 0] > 0, jnp.where(X[:, 1] > 0, 0, 1), 2)

m.data_on_model = dict(X=X, Y=Y)

# Define model -----
def model(X, Y, D_H1=5, K=3):
    N, D_X = X.shape

    # First hidden layer: 2 input features -> 5 hidden units
    w1 = m.bnn.layer_linear(
        X,
        dist=m.dist.normal(0, 1, name='w1_weight', shape=(D_X, D_H1)),
        activation='tanh'
    )

    # Final output layer: 5 hidden units -> K output units
    # Note: No activation is applied automatically inside the layer function here
    w2 = m.bnn.layer_linear(
```

```

    w1,
    dist=m.dist.normal(0, 1, name='w2_weight', shape=(D_H1, K))
)

# Apply Softmax across the K dimension (axis=-1) to yield probabilities
p = jax.nn.softmax(w2, axis=-1)

# Categorical Likelihood matching indices in Y
m.dist.categorical(probs=p, obs=Y)

# Run mcmc -----
m.fit(model) # Approximate posterior distributions

# Predictions from the model -----
import matplotlib.pyplot as plt

# Create a grid to evaluate the model
n_grid = 50
x0 = jnp.linspace(X[:, 0].min() - 0.5, X[:, 0].max() + 0.5, n_grid)
x1 = jnp.linspace(X[:, 1].min() - 0.5, X[:, 1].max() + 0.5, n_grid)
xx0, xx1 = jnp.meshgrid(x0, x1)
X_grid = jnp.c_[xx0.ravel(), xx1.ravel()]

# Swap data on model temporarily to predict on the grid
m.data_on_model = dict(X=X_grid, Y=jnp.zeros(X_grid.shape[0], dtype=jnp.int32))
pred = m.sample(data = m.data_on_model)['x']
p_mean = jnp.mean(pred, axis=0)

# Plotting the posterior predictive mean (categorical blending)
fig, ax = plt.subplots(figsize=(8, 6), constrained_layout=True)
contour = ax.contourf(xx0, xx1, p_mean.reshape(n_grid, n_grid), cmap="viridis", alpha=0.6)
scatter = ax.scatter(X[:, 0], X[:, 1], c=Y, cmap="viridis", edgecolors='k')
ax.set(title="Posterior Predictive Mean", xlabel="Feature 1", ylabel="Feature 2")
fig.colorbar(contour, ax=ax)

```

/home/sosa/work/3.12venv/lib/python3.10/site-packages/tqdm/auto.py:21: TqdmWarning:

IProgress not found. Please update jupyter and ipywidgets. See <https://ipywidgets.readthedocs.org/>

BI v 0.0.46 package loaded  
jax.local\_device\_count 32

This function is still in development. Use it with caution.  
This function is still in development. Use it with caution.  
This function is still in development. Use it with caution.  
This function is still in development. Use it with caution.

```
0%|          | 0/2000 [00:00<?, ?it/s]Compiling.. : 0%|          | 0/2000 [00:00<?, ?it/s]
0%|          | 0/2000 [00:00<?, ?it/s]
Compiling.. : 0%|          | 0/2000 [00:00<?, ?it/s]
```

```
0%|          | 0/2000 [00:00<?, ?it/s]
Compiling.. : 0%|          | 0/2000 [00:00<?, ?it/s]
```

```
0%|          | 0/2000 [00:00<?, ?it/s]
Compiling.. : 0%|          | 0/2000 [00:00<?, ?it/s]
```

This function is still in development. Use it with caution.  
This function is still in development. Use it with caution.

```
Running chain 0: 0%|          | 0/2000 [00:01<?, ?it/s]
Running chain 1: 0%|          | 0/2000 [00:01<?, ?it/s]
Running chain 2: 0%|          | 0/2000 [00:01<?, ?it/s]
```

```
Running chain 3: 0%|          | 0/2000 [00:01<?, ?it/s]
Running chain 2: 5%|          | 100/2000 [00:01<00:08, 232.10it/s]Running chain 0: 5%|
Running chain 1: 5%|          | 100/2000 [00:01<00:09, 210.75it/s]
```

```
Running chain 3: 5%|          | 100/2000 [00:01<00:09, 197.42it/s]Running chain 0: 10%|
Running chain 2: 10%|         | 200/2000 [00:02<00:09, 191.34it/s]
```

```
Running chain 3: 10%|         | 200/2000 [00:02<00:09, 186.39it/s]
Running chain 1: 10%|         | 200/2000 [00:02<00:11, 151.66it/s]Running chain 0: 15%|
```

Running chain 3: 15%| | 300/2000 [00:02<00:08, 196.52it/s]

Running chain 2: 15%| | 300/2000 [00:02<00:08, 190.56it/s]

Running chain 1: 15%| | 300/2000 [00:02<00:09, 180.84it/s]

Running chain 2: 20%| | 400/2000 [00:02<00:07, 216.48it/s]Running chain 0: 20%|

Running chain 3: 20%| | 400/2000 [00:03<00:07, 214.67it/s]

Running chain 1: 20%| | 400/2000 [00:03<00:07, 203.58it/s]

Running chain 3: 25%| | 500/2000 [00:03<00:06, 237.90it/s]Running chain 0: 25%|

Running chain 2: 25%| | 500/2000 [00:03<00:07, 212.58it/s]

Running chain 1: 25%| | 500/2000 [00:03<00:06, 222.40it/s]

Running chain 3: 30%| | 600/2000 [00:03<00:05, 251.40it/s]Running chain 0: 30%|

Running chain 2: 30%| | 600/2000 [00:03<00:06, 229.74it/s]

Running chain 1: 30%| | 600/2000 [00:03<00:05, 241.48it/s]

Running chain 3: 35%| | 700/2000 [00:04<00:04, 266.37it/s]

Running chain 2: 35%| | 700/2000 [00:04<00:05, 252.78it/s]Running chain 0: 35%|

Running chain 1: 35%| | 700/2000 [00:04<00:05, 259.14it/s]

Running chain 3: 40%| | 800/2000 [00:04<00:04, 279.70it/s]

Running chain 2: 40%| | 800/2000 [00:04<00:04, 276.93it/s]Running chain 0: 40%|

Running chain 1: 40%| | 800/2000 [00:04<00:04, 270.11it/s]

Running chain 3: 45%| | 900/2000 [00:04<00:03, 302.54it/s]

Running chain 1: 45%| | 900/2000 [00:04<00:03, 312.69it/s]Running chain 0: 45%|

Running chain 2: 45%| | 900/2000 [00:04<00:03, 278.49it/s]

Running chain 3: 50%| | 1000/2000 [00:04<00:03, 295.72it/s]  
 Running chain 1: 50%| | 1000/2000 [00:05<00:03, 299.00it/s]Running chain 0: 50%|  
  
 Running chain 2: 50%| | 1000/2000 [00:05<00:03, 260.21it/s]  
  
 Running chain 3: 55%| | 1100/2000 [00:05<00:03, 283.07it/s]  
 Running chain 1: 55%| | 1100/2000 [00:05<00:03, 271.76it/s]Running chain 0: 55%|  
  
 Running chain 3: 60%| | 1200/2000 [00:05<00:02, 283.50it/s]  
  
 Running chain 2: 55%| | 1100/2000 [00:05<00:03, 231.44it/s]  
 Running chain 1: 60%| | 1200/2000 [00:05<00:03, 262.96it/s]Running chain 0: 60%|  
  
 Running chain 3: 65%| | 1300/2000 [00:06<00:02, 284.19it/s]  
  
 Running chain 2: 60%| | 1200/2000 [00:06<00:03, 226.25it/s]  
 Running chain 1: 65%| | 1300/2000 [00:06<00:02, 266.98it/s]Running chain 0: 65%|  
  
 Running chain 3: 70%| | 1400/2000 [00:06<00:02, 291.83it/s]  
  
 Running chain 2: 65%| | 1300/2000 [00:06<00:03, 223.70it/s]  
 Running chain 1: 70%| | 1400/2000 [00:06<00:02, 262.71it/s]  
  
 Running chain 3: 75%| | 1500/2000 [00:06<00:01, 292.14it/s]Running chain 0: 70%|  
  
 Running chain 3: 80%| | 1600/2000 [00:07<00:01, 286.14it/s]  
 Running chain 1: 75%| | 1500/2000 [00:07<00:01, 261.98it/s]Running chain 0: 75%|  
  
 Running chain 2: 70%| | 1400/2000 [00:07<00:02, 219.98it/s]  
  
 Running chain 3: 85%| | 1700/2000 [00:07<00:01, 290.15it/s]  
 Running chain 1: 80%| | 1600/2000 [00:07<00:01, 256.56it/s]Running chain 0: 80%|  
  
 Running chain 2: 75%| | 1500/2000 [00:07<00:02, 217.02it/s]

Running chain 3: 90% | 1800/2000 [00:07<00:00, 292.58it/s]  
Running chain 1: 85% | 1700/2000 [00:07<00:01, 252.27it/s]Running chain 0: 85%

Running chain 3: 95% | 1900/2000 [00:08<00:00, 288.06it/s]

Running chain 2: 80% | 1600/2000 [00:08<00:01, 209.16it/s]  
Running chain 1: 90% | 1800/2000 [00:08<00:00, 258.51it/s]Running chain 0: 90%

Running chain 3: 100% | 2000/2000 [00:08<00:00, 297.91it/s]Running chain 3: 100%

Running chain 2: 85% | 1700/2000 [00:08<00:01, 212.81it/s]Running chain 0: 95%  
Running chain 1: 95% | 1900/2000 [00:08<00:00, 257.39it/s]Running chain 0: 100%

Running chain 2: 90% | 1800/2000 [00:09<00:00, 220.52it/s]  
Running chain 1: 100% | 2000/2000 [00:09<00:00, 262.78it/s]Running chain 1: 100%

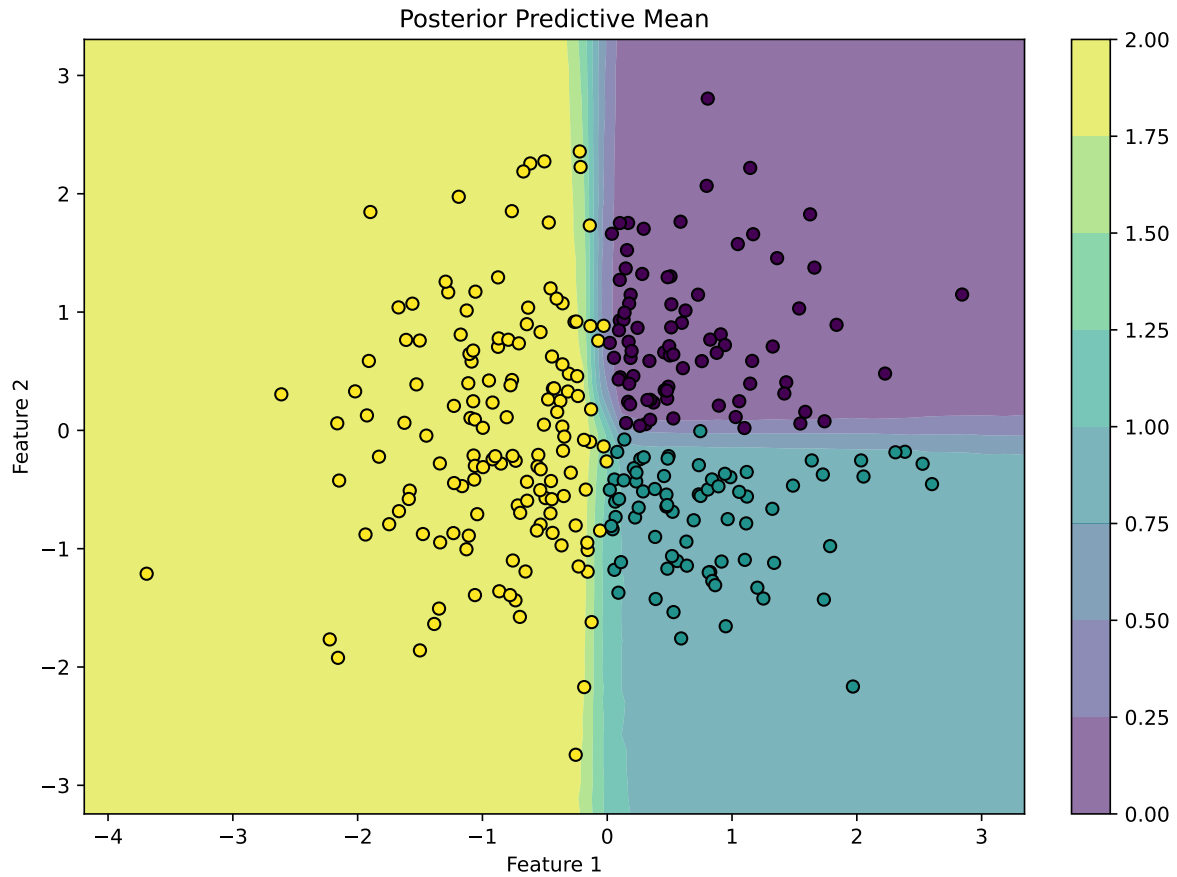
Running chain 2: 95% | 1900/2000 [00:09<00:00, 213.67it/s]

Running chain 2: 100% | 2000/2000 [00:09<00:00, 216.46it/s]Running chain 2: 100%  
/home/sosa/work/BI/BI/Main/main.py:666: UserWarning:

Sample's batch dimension size 4000 is different from the provided 1 num\_samples argument. De

This function is still in development. Use it with caution.

This function is still in development. Use it with caution.



## Julia

```
using BayesianInference
using PythonCall

# Setup device-----
m = importBI(platform="cpu")

# Generate Synthetic Data -----
np = pyimport("numpy")
jax_random = pyimport("jax.random")
jnp = pyimport("jax.numpy")

key = jax_random.PRNGKey(42)
X = jax_random.normal(key, (300, 2))
```

```

# Simple rule to partition into K=3 classes
Y = jnp.where(X[:, 0] > 0, jnp.where(X[:, 1] > 0, 0, 1), 2)

m.data_on_model["X"] = X
m.data_on_model["Y"] = Y

# Define model -----
@BI function model(X, Y)
  N, D_X = size(X)
  D_H1 = 5
  K = 3

  # First hidden layer
  w1 = m.bnn.layer_linear(
    X,
    dist=m.dist.normal(0, 1, name="w1_weight", shape=(D_X, D_H1)),
    activation="tanh"
  )

  # Final output layer
  w2 = m.bnn.layer_linear(
    w1,
    dist=m.dist.normal(0, 1, name="w2_weight", shape=(D_H1, K))
  )

  # Softmax conversion to probability simplex
  p = jax.nn.softmax(w2, axis=-1)

  # Categorical Likelihood
  m.dist.categorical(probs=p, obs=Y)
end

# Run mcmc -----
m.fit(model, num_samples=500, progress_bar=false)

```

## Mathematical Details

### Bayesian Formulation

For a multiclass classification task spanning  $N$  observations and  $K$  mutually exclusive classes, we model the probability vector  $\theta_i$  that the response  $Y_i \in \{0, 1, \dots, K - 1\}$  falls into each

respective class.

Using a single hidden layer with a hyperbolic tangent (tanh) activation function, the model is structured as:

$$\begin{aligned} Y_i &\sim \text{Categorical}(\theta_i) \\ \theta_i &= \text{Softmax}(\phi_i) \\ \phi_i &= H_i \Theta_2 \\ H_i &= \tanh(X_i \Theta_1) \\ \Theta_1 &\sim \text{Normal}(0, 1) \\ \Theta_2 &\sim \text{Normal}(0, 1) \end{aligned}$$

where:

- $Y_i$  is the observed class index for the  $i$ -th observation ( $Y_i \in \{0, 1, \dots, K - 1\}$ ).
- $\theta_i$  is the predicted probability vector for the  $i$ -th observation.
- $\phi_i$  are the  $K$ -dimensional logits.
- $X_i$  is the input row vector for the  $i$ -th observation, with features length  $D_X = 2$ .
- $H_i$  is the hidden layer representation vector for the  $i$ -th observation. It has length  $D_H = 5$ .
- $\Theta_1$  is the weight matrix of the first hidden layer ( $2 \times 5$ ).
- $\Theta_2$  is the final layer weight matrix mapping the hidden features to the logits for the  $K = 3$  classes ( $5 \times 3$ ).
- All elements within the weight matrices  $\Theta_1$  and  $\Theta_2$  are assigned independent standard Normal priors.

## Notes

### **i** Note

- For large outputs where  $K > 100$ , computing the exact softmax normalization scalar (the denominator term combining all exponentiated logits) can become computationally expensive over thousands of MCMC posterior evaluations.
- Neural networks configured with a standard Cross-Entropy loss mapping to one-hot vectors conceptually perform exactly this sequence: dot product of final weights  $\rightarrow$  Softmax  $\rightarrow$  Categorical Likelihood.

## Reference(s)

1. [PyData Berlin 2025: Introduction to Stochastic Variational Inference with NumPyro](#)