

Bayesian Multiclass Classification - Complete Hand Calculations

1. Model Setup

Consider a multiclass classification problem with K classes and N training examples:

- Training data: $D = \{(x_n, y_n)\}_{n=1}^N$ where $x_n \in \mathbb{R}^D$ and $y_n \in \{0, 1, \dots, K-1\}$
- Feature expansion: $\phi(x_n) \in \mathbb{R}^D$ (often just x_n with intercept)
- Weight matrix: $W \in \mathbb{R}^{K \times D}$ where row k contains weights w_k for class k

2. Prior Distribution

We place independent Gaussian priors on each weight:

Prior on individual weight:

$$p(w_{ij}) = N(w_{ij} \mid 0, 1/\alpha) = \sqrt{\alpha/(2\pi)} \exp(-\alpha w_{ij}^2/2)$$

Prior on weight vector w_k :

$$p(w_k) = N(w_k \mid 0, (1/\alpha)I) = (\alpha/2\pi)^{D/2} \exp(-\alpha \|w_k\|^2/2)$$

Prior on full weight matrix W :

$$p(W) = \prod_{k=1}^K \prod_{j=1}^D p(w_{kj}) = (\alpha/2\pi)^{KD/2} \exp(-\alpha/2 \sum_{k,j} w_{kj}^2)$$

Log prior:

$$\log p(W) = KD/2 \log(\alpha/2\pi) - \alpha/2 \|W\|_F^2$$

where $\|W\|_F^2 = \sum_{k,j} w_{kj}^2$ is the Frobenius norm.

3. Likelihood Function

The likelihood uses a categorical distribution with softmax link:

Softmax function:

$$\text{softmax}(f)_k = \exp(f_k) / \sum_j \exp(f_j)$$

For linear model $f_k(x_n) = w_k^T \phi(x_n)$:

$$p(y_n = k \mid x_n, W) = \text{softmax}(W\phi(x_n))_k = \exp(w_k^T \phi(x_n)) / \sum_j \exp(w_j^T \phi(x_n))$$

Likelihood for all data:

$$p(y \mid X, W) = \prod_n p(y_n \mid x_n, W) = \prod_n \prod_k [p(y_n = k \mid x_n, W)]^{I(y_n = k)}$$

Log likelihood using one-hot encoding $Y_{nk} = I(y_n = k)$:

$$\begin{aligned} \log p(y \mid X, W) &= \sum_n \sum_k Y_{nk} \log p(y_n = k \mid x_n, W) \\ &= \sum_n \sum_k Y_{nk} [w_k^T \phi(x_n) - \log \sum_j \exp(w_j^T \phi(x_n))] \end{aligned}$$

4. Posterior Distribution

Log posterior (up to constant):

$$\begin{aligned} \log p(W \mid y, X) &= \log p(y \mid X, W) + \log p(W) + \text{const} \\ &= \sum_n \sum_k Y_{nk} \log \text{softmax}(W\phi(x_n))_k - \alpha/2 \|W\|_F^2 + \text{const} \end{aligned}$$

5. MAP Estimation

To find W_{MAP} , we maximize $\log p(W \mid y, X)$ or minimize the negative log posterior.

Gradient with respect to w_k :

$$\nabla_{\{w_k\}} \log p(W \mid y, X) = \sum_n [Y_{nk} - p(y_n = k \mid x_n, W)] \phi(x_n) - \alpha w_k$$

Gradient in matrix form:

$$\nabla_W \log p(W \mid y, X) = \Phi^T (Y - P) - \alpha W$$

where:

- Φ is $N \times D$ design matrix with rows $\phi(x_n)^T$
- Y is $N \times K$ one-hot encoded label matrix
- P is $N \times K$ matrix with $P_{nk} = p(y_n = k \mid x_n, W)$

6. Laplace Approximation

The Hessian is computed as:

Block structure of Hessian: For weights w_k and w_l , the Hessian block is:

$$H_{\{kl\}} = -\partial^2 \log p(W | y, X) / \partial w_k \partial w_l^T$$

If $k = l$:

$$H_{\{kk\}} = \sum_n p_{nk}(1 - p_{nk}) \phi(x_n) \phi(x_n)^T + \alpha I$$

If $k \neq l$:

$$H_{\{kl\}} = -\sum_n p_{nk} p_{nl} \phi(x_n) \phi(x_n)^T$$

Laplace approximation:

$$p(W | y, X) \approx N(W | W_{MAP}, \Sigma)$$

where Σ^{-1} is the Hessian at W_{MAP} .

7. Predictive Distribution

*For latent functions at x^**

$$f_k(x^*) | y, X \sim N(\mu_{k^*}, \sigma^2_{k^*})$$

where:

$$\mu_{k^*} = w_{k^*}^T \phi(x^*)$$

$$\sigma^2_{k^*} = \phi(x^*)^T \Sigma_k \phi(x^*)$$

and Σ_k is the k -th diagonal block of Σ .

For class probabilities (via Monte Carlo):

$$p(y^* = k | y, x^*) \approx 1/S \sum_{s=1}^S \text{softmax}(W^{(s)} \phi(x^*))_k$$

where $W^{(s)} \sim N(W_{MAP}, \Sigma)$.

8. Example: K=3 classes, D=2 features

Weight matrix:

```
W = [[w_11, w_12],... # weights for class 0
      [w_21, w_22],... # weights for class 1
      [w_31, w_32]]... # weights for class 2
```

Flattened representation:

```
w_flat = [w_11, w_12, w_21, w_22, w_31, w_32]^T
```

For a single data point $x_n = [1, x_{n1}]^T$ (with intercept):

Linear outputs:

```
f_0(x_n) = w_11 + w_12 x_n1
```

```
f_1(x_n) = w_21 + w_22 x_n1
```

```
f_2(x_n) = w_31 + w_32 x_n1
```

Softmax probabilities:

```
p_0 = exp(f_0) / (exp(f_0) + exp(f_1) + exp(f_2))
```

```
p_1 = exp(f_1) / (exp(f_0) + exp(f_1) + exp(f_2))
```

```
p_2 = exp(f_2) / (exp(f_0) + exp(f_1) + exp(f_2))
```

If $y_n = 1$ (one-hot: $[0, 1, 0]$):

```
log likelihood contribution = 0·log(p_0) + 1·log(p_1) + 0·log(p_2) = log(p_1)
```

Gradient for w_1 :

```
 $\nabla_{\{w_1\}} \log p(y_n \mid x_n, W) = (1 - p_1) x_n$ 
```

Gradient for w_0 and w_2 :

```
 $\nabla_{\{w_0\}} \log p(y_n \mid x_n, W) = -p_0 x_n$ 
```

```
 $\nabla_{\{w_2\}} \log p(y_n \mid x_n, W) = -p_2 x_n$ 
```

9. Numerical Computation Steps

1. **Initialize:** $W^{(0)} = 0$

2. **Iterate until convergence:**

- Compute $P = \text{softmax}(XW^T)$
- Gradient: $g = X^T(Y - P) - \alpha W$
- Hessian: H (block structure as above)
- Update: $W^{(t+1)} = W^{(t)} - H^{-1}g$

3. **At convergence:** $W_{\text{MAP}} = W^{(\text{final})}$

4. **Posterior covariance:** $\Sigma = H^{-1}$

5. **Predictions:** Monte Carlo sampling

10. Decision Theory

Expected utility for predicting class k:

$$EU(k) = \sum_j p(y^* = j \mid y, x^*) U(j, k)$$

Optimal decision:

$$k^* = \text{argmax}_k EU(k)$$

For 0/1 utility ($U(j,k) = I(j=k)$):

$$k^* = \text{argmax}_k p(y^* = k \mid y, x^*)$$

11. Complete Numerical Example

Let's work through a small example with $K=3$ classes, $N=3$ data points, and $D=2$ features.

Data

```
X = [[1, 0], ... # x_0 (intercept=1, feature=0)
      [1, 1], ... # x_1 (intercept=1, feature=1)
      ...,
      [1, -1]] ... # x_2 (intercept=1, feature=-1)

y = [0, 1, 2] ... # Class labels

Y = [[1, 0, 0], # One-hot encoding
      [0, 1, 0],
      [0, 0, 1]]
```

Initial weights ($W_0 = 0$)

```
W = [[0, 0], ... # w_0 for class 0
      [0, 0], ... # w_1 for class 1
      [0, 0]]    # w_2 for class 2
```

Step 1: Compute linear outputs $f = XW^T$

```
f = [[0, 0, 0], ... # f(x_0) for all classes
      [0, 0, 0], ... # f(x_1) for all classes
      [0, 0, 0]] ... # f(x_2) for all classes
```

Step 2: Apply softmax

For each row, softmax gives:

```
p_00 = exp(0)/(exp(0)+exp(0)+exp(0)) = 1/3
p_01 = exp(0)/(exp(0)+exp(0)+exp(0)) = 1/3
p_02 = exp(0)/(exp(0)+exp(0)+exp(0)) = 1/3
```

```
P = [[1/3, 1/3, 1/3],
      [1/3, 1/3, 1/3],
      [1/3, 1/3, 1/3]]
```

Step 3: Compute log likelihood

```
log L = Σ_n Σ_k Y_nk log P_nk
        = 1×log(1/3) + 0×log(1/3) + 0×log(1/3) + ... # n=0
        + 0×log(1/3) + 1×log(1/3) + 0×log(1/3) + ... # n=1
        + 0×log(1/3) + 0×log(1/3) + 1×log(1/3) + ... # n=2
        = 3 × log(1/3) = -3.296
```

Step 4: Compute gradient

$$\nabla_W \log p(W|y,X) = X^T(Y - P) - \alpha W$$

$$Y - P = \begin{bmatrix} 2/3 & -1/3 & -1/3 \\ -1/3 & 2/3 & -1/3 \\ -1/3 & -1/3 & 2/3 \end{bmatrix}$$

$$X^T(Y - P) = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & -1 \end{bmatrix} \times \begin{bmatrix} 2/3 & -1/3 & -1/3 \\ -1/3 & 2/3 & -1/3 \\ -1/3 & -1/3 & 2/3 \end{bmatrix}$$

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

With $\alpha = 1$:

$$\nabla_W = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 0 & -1 \end{bmatrix}$$

Step 5: Update weights (gradient ascent)

$$W_1 = W_0 + \eta \times \nabla_W \dots \text{ (with } \eta = 0.1 \text{)}$$

$$W_1 = \begin{bmatrix} 0 & 0 \\ 0 & 0.1 \\ 0 & -0.1 \end{bmatrix}$$

Step 6: Compute new predictions

$$f = XW_1^T = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0.1 & -0.1 \\ 0 & -0.1 & 0.1 \end{bmatrix}$$

For $x_1 = [1, 1]$:

$$p_{10} = \exp(0) / (\exp(0) + \exp(0.1) + \exp(-0.1)) = 0.315$$

$$p_{11} = \exp(0.1) / (\exp(0) + \exp(0.1) + \exp(-0.1)) = 0.347$$

$$p_{12} = \exp(-0.1) / (\exp(0) + \exp(0.1) + \exp(-0.1)) = 0.338$$

After Convergence (hypothetical)

Suppose after optimization we get:

```
W_MAP = [[0.5, -0.3],
          [0.2, 0.8],
          [-0.1, -0.5]]
```

Predictions at $x^* = [1, 0.5]$

```
f*(x*) = x* × W_MAP^T = [1, 0.5] × [[0.5, 0.2, -0.1],
                                     [-0.3, 0.8, -0.5]]
..... = [0.35, 0.6, -0.35]
```

Softmax:

```
p_0 = exp(0.35)/(exp(0.35)+exp(0.6)+exp(-0.35)) = 0.331
p_1 = exp(0.6)/(exp(0.35)+exp(0.6)+exp(-0.35)) = 0.427
p_2 = exp(-0.35)/(exp(0.35)+exp(0.6)+exp(-0.35)) = 0.242
```

Decision: $\text{argmax}_k p_k = 1$

Laplace Approximation Example

For the Hessian at W_MAP , each block H_{kl} has form:

$$H_{kk} = \sum_n p_{nk}(1-p_{nk})x_n x_n^T + \alpha I$$

Example for $k=0$:

If $p_{00} = 0.7$, $p_{10} = 0.2$, $p_{20} = 0.1$:

$$H_{00} = 0.7 \times 0.3 \times \begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \end{bmatrix} + 0.2 \times 0.8 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \end{bmatrix} + 0.1 \times 0.9 \times \begin{bmatrix} 1 \\ -1 \end{bmatrix} \begin{bmatrix} 1 & -1 \end{bmatrix} + 1 \times I$$

$$..... = \begin{bmatrix} 0.21+0.16+0.09+1, & 0+0.16-0.09 \\ 0+0.16-0.09, & 0+0.16+0.09+1 \end{bmatrix}$$

.....

$$..... = \begin{bmatrix} 1.46, & 0.07 \\ 0.07, & 1.25 \end{bmatrix}$$

Monte Carlo Prediction

Draw samples $W^{(s)} \sim N(W_MAP, \Sigma)$, then:

$$p(y^*=k|x^*,y) \approx 1/S \sum_s \text{softmax}(W^{(s)}x^*)_k$$

Example with $S=3$ samples:

$$W^{(1)}: p^{(1)} = [0.32, 0.44, 0.24]$$

$$W^{(2)}: p^{(2)} = [0.35, 0.41, 0.24]$$

$$W^{(3)}: p^{(3)} = [0.33, 0.43, 0.24]$$

$$\text{Average: } p = [0.333, 0.427, 0.240]$$

This completes the mathematical derivation with concrete numerical examples for Bayesian multiclass classification.