

Understanding Model Calibration: Expected Calibration Error and Temperature Scaling

Subhasish Bandyopadhyay

November, 2024

Introduction

Deep neural networks have achieved remarkable success in classification tasks, but their confidence scores often don't align with their actual accuracy. A model might be 95% confident in a prediction that's only correct 70% of the time. This misalignment between confidence and accuracy is a calibration problem, and it has serious implications for real-world applications where decision-making depends on reliable uncertainty estimates. In this article, we'll explore Expected Calibration Error (ECE) as a metric for measuring calibration and temperature scaling as a simple yet effective post-processing technique to improve it.

Mathematical Foundation of Model Calibration

Formal Definition

Consider a classification model that outputs a probability distribution over K classes. For an input x , the model produces:

- A predicted class: $\hat{y} = \arg \max_k p(y = k|x)$
- A confidence score: $\hat{p} = \max_k p(y = k|x)$

A model is perfectly calibrated if:

$$P(Y = \hat{y} \mid \hat{P} = p) = p, \quad \forall p \in [0, 1]$$

This means that among all predictions where the model reports confidence p , the actual accuracy should be exactly p . For example, if we collect all predictions where the model is 70% confident, exactly 70% of them should be correct.

The Calibration Gap

The calibration gap at confidence level p is defined as:

$$\text{Gap}(p) = \left| P(Y = \hat{y} \mid \hat{P} = p) - p \right|$$

This measures how far off the model's confidence is from its actual accuracy at that confidence level.

Expected Calibration Error (ECE) - Deep Dive

Mathematical Formulation

Since we work with finite samples, we can't evaluate calibration at every possible confidence value. Instead, ECE discretizes the probability space into M bins and computes a weighted average of calibration gaps:

$$\text{ECE} = \sum_{m=1}^M \frac{n_m}{n} \times |\text{acc}(B_m) - \text{conf}(B_m)|$$

Where:

- B_m is the m -th bin with boundaries $[(m-1)/M, m/M]$
- $n_m = |\{i : \hat{p}_i \in B_m\}|$ (number of samples in bin m)
- n = total number of samples
- $\text{acc}(B_m) = \frac{1}{n_m} \sum_{i \in B_m} \mathbb{I}(\hat{y}_i = y_i)$ (accuracy in bin m)
- $\text{conf}(B_m) = \frac{1}{n_m} \sum_{i \in B_m} \hat{p}_i$ (average confidence in bin m)
- \mathbb{I} is the indicator function

Detailed ECE Computation Example

```
1 # Bin 1: [0.0, 0.33] - Low confidence
2 # Contains 20 samples
3 # Average confidence: 0.25
4 # Actual accuracy: 0.30
5 # Contribution: (20/100) * |0.30 - 0.25| = 0.01
6 # Bin 2: [0.33, 0.67] - Medium confidence
7 # Contains 30 samples
8 # Average confidence: 0.50
9 # Actual accuracy: 0.45
10 # Contribution: (30/100) * |0.45 - 0.50| = 0.015
11 # Bin 3: [0.67, 1.0] - High confidence
12 # Contains 50 samples
13 # Average confidence: 0.85
14 # Actual accuracy: 0.70
15 # Contribution: (50/100) * |0.70 - 0.85| = 0.075
16 # ECE = 0.01 + 0.015 + 0.075 = 0.10
```

Implementation Details

```
1 def compute_ece(probs, labels, n_bins=10):
2     """ Compute Expected Calibration Error
3     Args:
4         probs: np.array of shape (n_samples, n_classes) - probability
5               distributions
6         labels: np.array of shape (n_samples,) - true labels
7         n_bins: number of bins for discretization
8     Returns:
9         ece: scalar Expected Calibration Error
10    """
11    # Clip probabilities to avoid numerical issues
12    probs = np.clip(probs, 1e-5, 1-1e-5)
13    # Extract confidence (maximum probability) and predictions
14    confidences = np.max(probs, axis=1) # p_i for each sample
15    predictions = np.argmax(probs, axis=1) # _i for each sample
16    accuracies = predictions == labels # binary array of correct/incorrect
17    # Create bins
18    bin_boundaries = np.linspace(0, 1, n_bins + 1)
19    bin_lowers = bin_boundaries[:-1]
20    bin_uppers = bin_boundaries[1:]
21    ece = 0.0
22    for bin_lower, bin_upper in zip(bin_lowers, bin_uppers):
23        # Find samples in this bin
24        in_bin = (confidences >= bin_lower) & (confidences < bin_upper)
25        prop_in_bin = np.mean(in_bin) # n_m / n
26        if prop_in_bin > 0:
27            # Compute accuracy and average confidence in bin
28            accuracy_in_bin = np.mean(accuracies[in_bin]) # acc(B_m)
29            avg_confidence_in_bin = np.mean(confidences[in_bin]) # conf(B_m)
30            # Add weighted absolute difference to ECE
31            ece += prop_in_bin * np.abs(avg_confidence_in_bin - accuracy_in_bin)
32
33    return ece
```

Temperature Scaling: The Mathematical Framework

Softmax and Logits

Neural networks for K -class classification typically output logits $z \in \mathbb{R}^K$, which are converted to probabilities via softmax:

$$p_i = \frac{\exp(z_i)}{\sum_{j=1}^K \exp(z_j)}$$

The maximum probability (confidence) is:

$$\hat{p} = \max_i p_i = \max_i \left[\frac{\exp(z_i)}{\sum_j \exp(z_j)} \right]$$

Temperature-Modified Softmax

Temperature scaling introduces a scalar parameter $T > 0$:

$$p_i^{(T)} = \frac{\exp(z_i/T)}{\sum_{j=1}^K \exp(z_j/T)}$$

This can be understood through several mathematical perspectives:

1. **Entropy Perspective** The entropy of the distribution changes with temperature:

$$H(p^{(T)}) = - \sum_i p_i^{(T)} \log(p_i^{(T)})$$

- As $T \rightarrow \infty$: $H(p^{(T)}) \rightarrow \log(K)$ (maximum entropy, uniform distribution)
- As $T \rightarrow 0$: $H(p^{(T)}) \rightarrow 0$ (minimum entropy, one-hot distribution)
- $T = 1$: Original entropy

2. **Gradient Perspective** The gradient of log-probabilities with respect to temperature:

$$\frac{\partial \log(p_i^{(T)})}{\partial T} = \frac{z_i - \sum_j p_j^{(T)} z_j}{T^2}$$

This shows how confidence changes: if z_i is larger than the weighted average, decreasing T increases p_i .

3. **Preservation of Accuracy** Crucially, temperature scaling preserves the argmax:

$$\arg \max_i p_i^{(T)} = \arg \max_i \exp(z_i/T) = \arg \max_i z_i$$

This is because $\exp(\cdot)$ is monotonic and $T > 0$. Therefore:

- The predicted class remains unchanged
- Top-1 accuracy is preserved
- Only confidence scores are adjusted

Optimal Temperature Selection

The optimal temperature T^* minimizes ECE on a validation set:

$$T^* = \arg \min_T \text{ECE}(p^{(T)}, y_{\text{val}})$$

This optimization problem has nice properties:

- **Convexity**: ECE is often approximately convex in T near the optimum
- **Smoothness**: The objective is differentiable with respect to T
- **Bounded**: Practical T values typically lie in $[0.5, 5.0]$

Mathematical Analysis of Temperature Effects

Consider the logit difference between the top two classes:

$$\Delta z = z_{\text{max}} - z_{\text{second}}$$

The confidence after temperature scaling is approximately:

$$\hat{p}^{(T)} \approx \sigma(\Delta z/T) = \frac{1}{1 + \exp(-\Delta z/T)}$$

Where σ is the sigmoid function. This reveals:

- Large T : $\hat{p}^{(T)} \rightarrow 0.5 + \frac{\Delta z}{4T}$ (linear approximation near 0.5)
- Small T : $\hat{p}^{(T)} \rightarrow 1$ or 0 (saturation)

The rate of confidence change:

$$\left| \frac{d\hat{p}}{dT} \right| \propto \Delta z \times \hat{p}(1 - \hat{p})/T^2$$

Advanced Calibration Metrics

Maximum Calibration Error (MCE)

While ECE measures average calibration error, MCE captures the worst-case:

$$\text{MCE} = \max_m |\text{acc}(B_m) - \text{conf}(B_m)|$$

Adaptive ECE (ACE)

Instead of equal-width bins, ACE uses equal-frequency bins:

- Each bin contains exactly n/M samples

This ensures all bins contribute equally to the metric, avoiding issues with sparse bins.

Classwise ECE

For multi-class problems, we can compute ECE per class:

$$\text{ECE}_k = \sum_m \frac{n_m^k}{n^k} \times |\text{acc}_k(B_m) - \text{conf}_k(B_m)|$$

Where calculations are restricted to samples where class k was predicted.

Theoretical Insights

Why Neural Networks Are Miscalibrated

Recent research suggests several factors:

- **Capacity and Overfitting:** Modern networks have capacity to memorize training data, leading to overconfident predictions
- **Cross-Entropy Loss:** The negative log-likelihood loss drives logits to infinity for correctly classified examples:

$$L = -\log(p_y) \rightarrow 0 \quad \text{implies} \quad z_y \rightarrow \infty$$

- **Batch Normalization:** Can amplify confidence by reducing internal covariate shift
- **Model Size:** Deeper and wider networks tend to be more miscalibrated:

$$\text{ECE} \propto f(\text{depth, width, parameters})$$

Statistical Properties

Consistency Temperature scaling is consistent: as $n \rightarrow \infty$, the estimated T^* converges to the true optimal temperature.

Sample Complexity The number of validation samples needed for reliable temperature estimation:

$$n_{\text{val}} \geq O\left(\frac{K \log(K)}{\varepsilon^2}\right)$$

Where K is the number of classes and ε is the desired ECE accuracy.

Implementation: Complete Framework

```
1 import numpy as np
2 import torch
3 import torch.nn.functional as F
4 from scipy.optimize import minimize_scalar
5
6 class TemperatureScaling:
7     def __init__(self, n_bins=15):
8         self.n_bins = n_bins
9         self.temperature = 1.0
10
11     def compute_ece(self, probs, labels, n_bins=None):
12         """ Enhanced ECE computation with additional metrics """
13         if n_bins is None:
14             n_bins = self.n_bins
15         probs = np.clip(probs, 1e-5, 1-1e-5)
16         confidences = np.max(probs, axis=1)
17         predictions = np.argmax(probs, axis=1)
18         accuracies = predictions == labels
19         # Standard ECE
20         bin_boundaries = np.linspace(0, 1, n_bins + 1)
21         bin_lower = bin_boundaries[:-1]
22         bin_upper = bin_boundaries[1:]
23         ece = 0.0
24         mce = 0.0 # Maximum Calibration Error
25         bin_stats = []
26         for bin_lower, bin_upper in zip(bin_lower, bin_upper):
27             in_bin = (confidences >= bin_lower) & (confidences < bin_upper)
28             prop_in_bin = np.mean(in_bin)
29             if prop_in_bin > 0:
30                 accuracy_in_bin = np.mean(accuracies[in_bin])
31                 avg_confidence_in_bin = np.mean(confidences[in_bin])
32                 calibration_gap = np.abs(avg_confidence_in_bin -
33                                         accuracy_in_bin)
34                 ece += prop_in_bin * calibration_gap
35                 mce = max(mce, calibration_gap)
36                 bin_stats.append({
37                     'range': (bin_lower, bin_upper),
38                     'count': np.sum(in_bin),
39                     'accuracy': accuracy_in_bin,
40                     'confidence': avg_confidence_in_bin,
41                     'gap': calibration_gap
42                 })
43         return {
44             'ece': ece,
45             'mce': mce,
46             'bin_stats': bin_stats
47         }
48
49     def find_optimal_temperature(self, val_logits, val_labels):
50         """ Find optimal temperature using Brent's method """
51         val_logits_tensor = torch.tensor(val_logits)
52         def objective(temp):
53             # Compute ECE for given temperature
54             scaled_probs = F.softmax(val_logits_tensor / temp, dim=1).numpy()
55             return self.compute_ece(scaled_probs, val_labels)['ece']
56         # Optimize using scipy's Brent method
57         result = minimize_scalar(objective, bounds=(0.1, 10.0), method='bounded')
58         self.temperature = result.x
59         return self.temperature
60
61     def calibrate(self, logits):
62         """ Apply temperature scaling to logits """
63         return F.softmax(torch.tensor(logits) / self.temperature, dim=1).numpy()
```

Experimental Validation

Typical Results

On standard datasets, temperature scaling achieves:

Model	Dataset	ECE (Before)	ECE (After)	Optimal T
ResNet-110	CIFAR-10	0.142	0.031	2.31
DenseNet-121	CIFAR-10	0.187	0.042	2.84
ResNet-50	ImageNet	0.121	0.019	1.73

Reliability Diagrams

The improvement can be visualized through reliability diagrams, where perfect calibration appears as a diagonal line:

- Perfect calibration: accuracy = confidence
- Before scaling: typically curved above diagonal (overconfident)
- After scaling: closely follows diagonal

Limitations and Extensions

Limitations of Temperature Scaling

- **Single Parameter:** Cannot correct class-specific miscalibration
- **Global Adjustment:** Applies uniform scaling across all confidence levels
- **Validation Dependency:** Requires held-out data for temperature selection

Advanced Calibration Methods

- **Vector Scaling** Uses a different temperature per class:

$$p_i^{(T)} = \frac{\exp(z_i/T_i)}{\sum_j \exp(z_j/T_j)}$$

- **Matrix Scaling** Applies a linear transformation:

$$z' = Wz + b$$

$$p = \text{softmax}(z')$$

- **Platt Scaling** Fits a sigmoid to logits:

$$p' = \sigma(az + b)$$

Conclusion

Expected Calibration Error provides a rigorous framework for measuring model calibration, while temperature scaling offers an elegant solution with strong theoretical foundations. The mathematical simplicity of temperature scaling modifying logits by a single scalar belies its effectiveness in practice. The key insights are:

- ECE quantifies the expected difference between confidence and accuracy
- Temperature scaling preserves accuracy while adjusting confidence
- The optimal temperature can be found efficiently through optimization

- The method is theoretically grounded and empirically validated

For practitioners, implementing temperature scaling should be standard practice when deploying classification models, especially in risk-sensitive applications. The minimal computational overhead and significant calibration improvements make it an essential component of reliable machine learning systems.