# PAC-Bayes Generalization Bound:
# A Complete Self-Contained Proof

Generated by Alethfeld v5.1

January 5, 2026

**Abstract**

We derive a PAC-Bayes bound on the expected generalization error for hypotheses that achieve zero training error. The bound relates the expected population loss to the KL divergence between posterior and prior. This proof is **fully self-contained**: we prove the core Gibbs variational inequality from the non-negativity of KL divergence.

## 1 Setup

**Definition 1** (Weight Space and Distributions). Let $\mathcal{W}$ be the weight space of a learning algorithm. We define:

- $P$: prior distribution on $\mathcal{W}$

- $\mathcal{S} = \{w \in \mathcal{W} : \text{train\_error}(w) = 0\}$: the set of solutions

- $Q$: posterior distribution, defined as $P$ restricted to $\mathcal{S}$

- $h(w)$: generalization error (population loss) of weight vector $w$

- $\beta = m - 1$ where $m$ is the number of training points

## 2 Core Lemma: Gibbs Variational Inequality

**Lemma 2** (Change of Measure / Donsker-Varadhan). *Let $P$ and $Q$ be distributions on the same space with $Q \ll P$. Let $h(w)$ be any measurable function such that $\mathbb{E}_{w \sim P}[e^{\beta h(w)}] < \infty$ for some $\beta > 0$. Then:*

$$\mathbb{E}_{w \sim Q}[h(w)] \leq \frac{\text{KL}(Q \| P) + \ln \mathbb{E}_{w \sim P}[e^{\beta h(w)}]}{\beta}$$

*Proof.* We prove this from the non-negativity of KL divergence.

⟨2⟩**1.** Define the **tilted distribution**:

$$\tilde{P}(w) = \frac{e^{\beta h(w)} P(w)}{Z}, \quad \text{where } Z = \mathbb{E}_{w \sim P}[e^{\beta h(w)}]$$

[definition]

⟨2⟩**2.** $\tilde{P}$ is a valid probability distribution since:

$$\int \tilde{P}(w) \, dw = \frac{1}{Z} \int e^{\beta h(w)} P(w) \, dw = \frac{Z}{Z} = 1$$

[algebraic; from ⟨2⟩1]

⟨2⟩**3.** By non-negativity of KL divergence:

$$\mathrm{KL}(Q\|\tilde{P}) \geq 0$$

[known result; from ⟨2⟩2]

⟨2⟩**4.** Expanding the KL divergence:

$$\mathrm{KL}(Q\|\tilde{P}) = \mathbb{E}_{w \sim Q}\left[\ln \frac{Q(w)}{\tilde{P}(w)}\right] = \mathbb{E}_{w \sim Q}\left[\ln Q(w) - \ln \tilde{P}(w)\right]$$

[definition; from ⟨2⟩3]

⟨2⟩**5.** From the definition of $\tilde{P}$:

$$\ln \tilde{P}(w) = \beta h(w) + \ln P(w) - \ln Z$$

[substitution; from ⟨2⟩1, ⟨2⟩4]

⟨2⟩**6.** Therefore:

$$\mathrm{KL}(Q\|\tilde{P}) = \mathbb{E}_{w \sim Q}\left[\ln Q(w) - \beta h(w) - \ln P(w) + \ln Z\right]$$
$$= \mathrm{KL}(Q\|P) - \beta \mathbb{E}_{w \sim Q}[h(w)] + \ln Z$$

[algebraic; from ⟨2⟩4, ⟨2⟩5]

⟨2⟩**7.** From $\mathrm{KL}(Q\|\tilde{P}) \geq 0$:

$$\mathrm{KL}(Q\|P) - \beta \mathbb{E}_{w \sim Q}[h(w)] + \ln Z \geq 0$$

[modus ponens; from ⟨2⟩3, ⟨2⟩6]

⟨2⟩**8.** Rearranging:
$$\beta \mathbb{E}_{w \sim Q}[h(w)] \leq \mathrm{KL}(Q\|P) + \ln Z$$

Dividing by $\beta > 0$:

$$\mathbb{E}_{w \sim Q}[h(w)] \leq \frac{\mathrm{KL}(Q\|P) + \ln \mathbb{E}_{w \sim P}[e^{\beta h(w)}]}{\beta} \qquad \square$$

[algebraic; from ⟨2⟩7]

$$\square$$

# 3   Main Result

**Theorem 3** (PAC-Bayes Bound for Zero Training Error)**.** *Under the setup above:*

$$\mathbb{E}_{w \sim Q}[h(w)] \leq \frac{\mathrm{KL}(Q\|P) + \ln \mathbb{E}_{w \sim P}[e^{\beta h(w)}]}{\beta}$$

*where for posterior $Q$ restricted to solutions:*

$$\mathrm{KL}(Q\|P) = \ln \frac{1}{\text{evidence}}$$

*Proof.* We apply Lemma 2 with specific choices for the posterior.

⟨1⟩**1. Assumption.** Let $P$ be a prior distribution on the weight space $\mathcal{W}$. [assumption]

⟨1⟩**2. Assumption.** Let $Q$ be the posterior distribution, defined as $P$ restricted to the set of solutions $\mathcal{S} = \{w : \text{train\_error}(w) = 0\}$. [assumption]

⟨1⟩**3. Assumption.** Let $h(w)$ denote the generalization error (population loss) of weight vector $w$. [assumption]

⟨1⟩**4. Assumption.** Let $\beta = m - 1$ where $m$ is the number of training points. [assumption]

⟨1⟩**5. Gibbs Variational Inequality** (Lemma 2): For any prior $P$, posterior $Q$, and $\beta > 0$:

$$\mathbb{E}_{w \sim Q}[h(w)] \leq \frac{\text{KL}(Q \| P) + \ln \mathbb{E}_{w \sim P}[e^{\beta h(w)}]}{\beta}$$

[✓ VERIFIED; proved above]

⟨1⟩**6.** Since $Q$ is $P$ restricted to solutions $\mathcal{S}$, we have:

$$Q(w) = \frac{P(w)}{P(\mathcal{S})} \quad \text{for } w \in \mathcal{S}, \qquad Q(w) = 0 \text{ otherwise}$$

[definition; from 3, 3]

⟨1⟩**7.** Computing the KL divergence:

$$\begin{aligned}
\text{KL}(Q \| P) &= \mathbb{E}_{w \sim Q}\left[\ln \frac{Q(w)}{P(w)}\right] \\
&= \mathbb{E}_{w \sim Q}\left[\ln \frac{P(w)/P(\mathcal{S})}{P(w)}\right] \\
&= \mathbb{E}_{w \sim Q}\left[\ln \frac{1}{P(\mathcal{S})}\right] \\
&= \ln \frac{1}{P(\mathcal{S})}
\end{aligned}$$

[algebraic; from 3]

⟨1⟩**8.** $P(\mathcal{S})$ is the prior probability of all solutions, i.e., the **evidence**. Thus:

$$\text{KL}(Q \| P) = \ln \frac{1}{\text{evidence}}$$

[definition; from 3]

⟨1⟩**9.** For $w \in \text{supp}(Q) = \mathcal{S}$, the training error is 0, so $h(w)$ measures the population (generalization) error directly. [modus ponens; from 3, 3]

⟨1⟩**10.** Applying Lemma 2 with our setup:

$$\mathbb{E}_{w \sim Q}[h(w)] \leq \frac{\ln(1/\text{evidence}) + \ln \mathbb{E}_{w \sim P}[e^{(m-1)h(w)}]}{m - 1}$$

[substitution; from 3, 3, 3]

⟨1⟩**11. QED.** The expected generalization error over solutions is bounded by the sum of log-inverse-evidence and the log-moment-generating-function of the loss under the prior, normalized by $\beta = m - 1$. [from 3, 3]

□

# 4 Interpretation

The bound decomposes into two terms:

1. **Complexity term**: $\ln(1/\text{evidence}) = -\ln P(\mathcal{S})$

   This measures how "surprising" it is that a random weight vector from the prior achieves zero training error. Smaller solution sets (more complex hypotheses) yield larger complexity penalties.

2. **Prior MGF term**: $\ln \mathbb{E}_{w \sim P}[e^{\beta h(w)}]$

   This is the log-moment-generating-function of the generalization error under the prior. It captures how well the prior concentrates on low-error hypotheses.

# 5 Proof Metadata

| | |
|---|---|
| **Graph ID:** | `graph-18d478-4253a9` |
| **Version:** | 21 |
| **Proof Mode:** | strict-mathematics |
| **Total Steps:** | 19 |
| **Assumptions:** | 5 |
| **Claims:** | 13 |
| **Verified:** | 1 (Gibbs lemma) |
| **Taint Status:** | **Clean** (fully self-contained) |

# 6 References

- Donsker, M. D. & Varadhan, S. R. S. (1975). Asymptotic evaluation of certain Markov process expectations for large time. *Comm. Pure Appl. Math.*

- McAllester, D. A. (1999). PAC-Bayesian model averaging. *Proceedings of the 12th Annual Conference on Computational Learning Theory.*

- Catoni, O. (2007). *PAC-Bayesian Supervised Classification.* Institute of Mathematical Statistics Lecture Notes.

*Generated by Alethfeld Proof Orchestrator v5.1*