

Introduction to Probabilistic Machine Learning

Summer School on Methods for Statistical Evaluation of AI

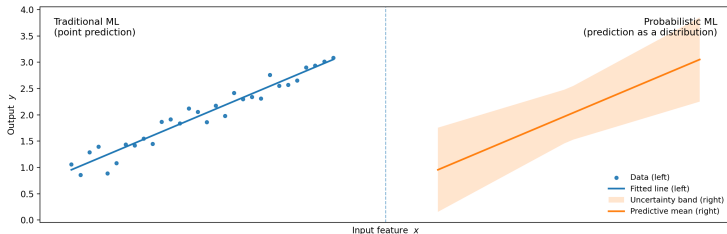
Andrés Masegosa

[Material made by Thomas D. Nielsen and Helge Langseth]

Introduction

From Predictions to Uncertainty-Aware Predictions

- Traditional ML: learns a function $f(x) \approx y$
 - Outputs: single prediction
 - Ignores uncertainty
- Probabilistic ML:
 - Models **distributions**, not just points
 - Explicitly represents **uncertainty**
 - Principles: *Models + Inference*



Prediction as a distribution (uncertainty included).

**In the real world, predictions are not enough. We need to know
*how confident they are.***

- **Healthcare:** Misdiagnosis risk - uncertainty flags when to call a human expert.
- **Autonomous driving:** Preventing Fatal Errors - uncertainty flags when driving under conditions not included in the train data.
- **General ML systems:** Uncertainty enables robustness, outlier detection, and better decision-making.

Probabilistic ML = Predictions + Confidence → Safer, more trustworthy AI.

Bayesian Machine Learning = Probabilistic model + Bayesian inference

- **Likelihood-part:** A probabilistic model typically defined by $p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta})$.
- **Prior:** $p(\boldsymbol{\theta})$ reflects our *a priori* belief about the parameters $\boldsymbol{\theta}$.

Bayesian Machine Learning = Probabilistic model + Bayesian inference

- **Likelihood-part:** A probabilistic model typically defined by $p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta})$.
- **Prior:** $p(\boldsymbol{\theta})$ reflects our *a priori* belief about the parameters $\boldsymbol{\theta}$.

Now we can calculate the posterior over $\boldsymbol{\theta}$ given training data \mathcal{D} ,

$$p(\boldsymbol{\theta} \mid \mathcal{D}) = \frac{p(\boldsymbol{\theta}) p(\mathcal{D} \mid \boldsymbol{\theta})}{p(\mathcal{D})},$$

... and, e.g., the predictive distribution of a new observation \mathbf{x}' :

$$p(\mathbf{y}' \mid \mathbf{x}' \mathcal{D}) = \int_{\boldsymbol{\theta}} p(\mathbf{y}' \mid \mathbf{x}', \boldsymbol{\theta}) p(\boldsymbol{\theta} \mid \mathcal{D}) d\boldsymbol{\theta}.$$

Bayesian Machine Learning = Probabilistic model + Bayesian inference

- **Likelihood-part:** A probabilistic model typically defined by $p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta})$.
- **Prior:** $p(\boldsymbol{\theta})$ reflects our *a priori* belief about the parameters $\boldsymbol{\theta}$.

Now we can calculate the posterior over $\boldsymbol{\theta}$ given training data \mathcal{D} ,

$$p(\boldsymbol{\theta} | \mathcal{D}) = \frac{p(\boldsymbol{\theta}) p(\mathcal{D} | \boldsymbol{\theta})}{p(\mathcal{D})},$$

... and, e.g., the predictive distribution of a new observation \mathbf{x}' :

$$p(\mathbf{y}' | \mathbf{x}' \mathcal{D}) = \int_{\boldsymbol{\theta}} p(\mathbf{y}' | \mathbf{x}', \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathcal{D}) d\boldsymbol{\theta}.$$

Being Bayesian means maintaining a distribution over $\boldsymbol{\theta}$.

Using a point-estimate for $\boldsymbol{\theta}$ is not **Bayesian** ML.

A Bayesian linear regression with univariate explanatory variables:

Likelihood – $p(\mathcal{D} | \theta)$: $p(y_i | x_i, \mathbf{w}, \sigma_y^2) = \mathcal{N}(w_0 + w_1 \cdot x_i, \sigma_y^2)$

Note! The observation noise, σ_y^2 , is known, so the parameter-set is simply $\theta = \{\mathbf{w}\}$.

Prior – $p(\theta)$: $p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \sigma_w^2)$

Bayesian Linear regression – Full model:

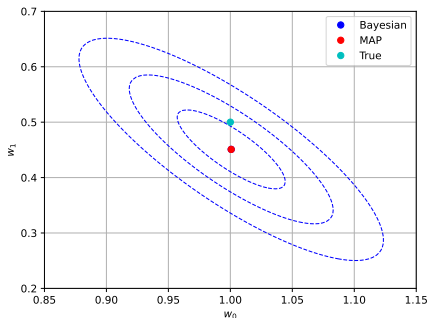
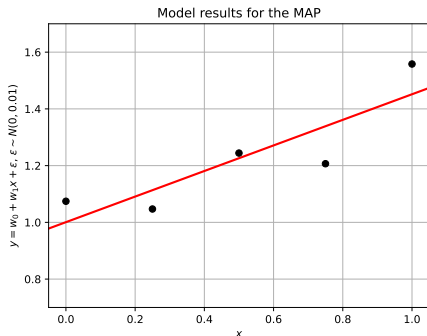
$$p(\mathcal{D}, \theta) = p(\{y_i\}_{i=1}^n, \mathbf{w} | \{\mathbf{x}_i\}_{i=1}^n, \sigma_y^2, \sigma_w^2) = \underbrace{p(\mathbf{w} | \sigma_w^2)}_{p(\theta)} \underbrace{\prod_{i=1}^n p(y_i | \mathbf{w}, \mathbf{x}_i, \sigma_y^2)}_{p(\mathcal{D} | \theta)}$$

Example: Linear regression – MAP vs (fully) Bayesian

Bayes linear regression with some fake data:

- We have generated $N = 5$ examples from $y_i = 1.0 + 0.5 \cdot x_i + \epsilon_i$, $\epsilon_i \sim \mathcal{N}(0, 0.1^2)$.
- Weights unknown a priori, so here we use the vague priors $w_j \sim \mathcal{N}(0, 10^2)$.

Results for the MAP and the fully Bayesian model:



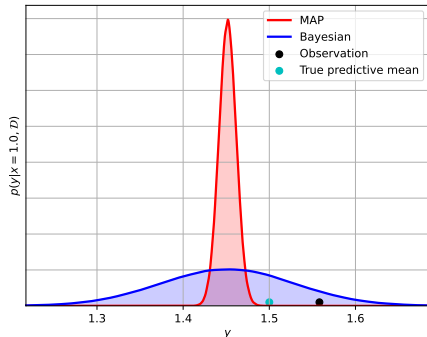
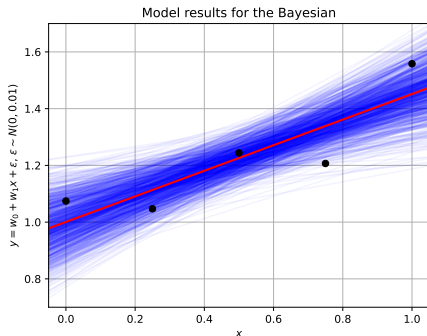
- **MAP:** Reasonable point estimate; No model uncertainty;
- **Bayes:** Model uncertainty around same MAP estimate;

Example: Linear regression – MAP vs (fully) Bayesian

Bayes linear regression with some fake data:

- We have generated $N = 5$ examples from $y_i = 1.0 + 0.5 \cdot x_i + \epsilon_i$, $\epsilon_i \sim \mathcal{N}(0, 0.1^2)$.
- Weights unknown a priori, so here we use the vague priors $w_j \sim \mathcal{N}(0, 10^2)$.

Results for the MAP and the fully Bayesian model:



- **MAP:** Reasonable point estimate; No model uncertainty; Predictive uncertainty degenerated to observation noise: poor fit wrt. true value and observation.
- **Bayes:** Model uncertainty around same MAP estimate; Captures model uncertainty well; Predictive distribution reasonable.

Bayesian inference is in principle easy using Bayes' rule:

$$p(\boldsymbol{\theta} | \mathcal{D}) = \frac{p(\boldsymbol{\theta}) p(\mathcal{D} | \boldsymbol{\theta})}{p(\mathcal{D})} = \frac{p(\mathcal{D}, \boldsymbol{\theta})}{\int_{\boldsymbol{\theta}} p(\boldsymbol{\theta}) p(\mathcal{D} | \boldsymbol{\theta}) d\boldsymbol{\theta}}$$

Note! This can only be solved analytically for **some simple models** (e.g., linear regression), but typically not for the really interesting models.

We need to approximate $p(\boldsymbol{\theta} | \mathcal{D})$

What we want:

- Computationally efficient;
- Well-founded approach;
- Easy integration with other frameworks.

What we don't want:

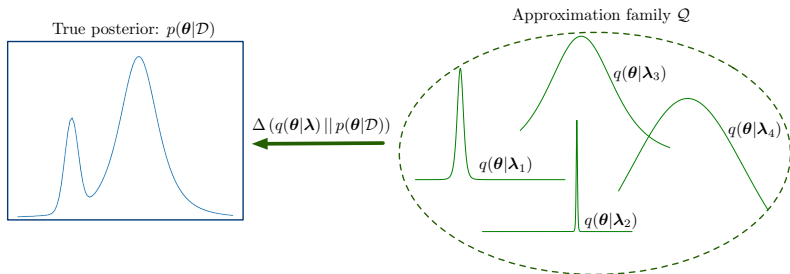
- Non scalable solutions;
- Widely applicable.

FUNDAMENTAL assumption:

It will always be *computationally efficient* to evaluate $p(\mathcal{D}, \boldsymbol{\theta})$ at any given point $\{\mathcal{D}, \boldsymbol{\theta}\}$, e.g., using the simple factorization $p(\mathcal{D}, \boldsymbol{\theta}) = p(\boldsymbol{\theta}) \cdot p(\mathcal{D} | \boldsymbol{\theta}) = p(\boldsymbol{\theta}) \prod_i p(\mathbf{x}_i | \boldsymbol{\theta})$.

Variational Bayes: Approximate inference by optimization

Variational Inference: Approximate the true posterior distribution $p(\boldsymbol{\theta} | \mathcal{D})$ with a **variational distribution** from a tractable family of distributions \mathcal{Q} . The family is indexed by the parameters $\boldsymbol{\lambda}$.



- **General goal:** Somehow approximate $p(\boldsymbol{\theta} \mid \mathcal{D})$ with a $q(\boldsymbol{\theta} \mid \mathcal{D})$.
 - **Note!** We use $q(\boldsymbol{\theta})$ as a short-hand for $q(\boldsymbol{\theta} \mid \mathcal{D})$.

Formalization of approximate inference through optimization:

Given a family of tractable distributions \mathcal{Q} and a distance measure between distributions Δ , choose

$$\hat{q}(\boldsymbol{\theta}) = \arg \min_{q \in \mathcal{Q}} \Delta(q(\boldsymbol{\theta}) \parallel p(\boldsymbol{\theta} \mid \mathcal{D})).$$

Decisions to be made:

- 1 How to define $\Delta(\cdot \parallel \cdot)$ so that we end up with a high-quality solution?
 - How to work with $\Delta(q(\boldsymbol{\theta}) \parallel p(\boldsymbol{\theta} \mid \mathcal{D}))$ when we don't know what $p(\boldsymbol{\theta} \mid \mathcal{D})$ is?
- 2 How to define a family of distributions \mathcal{Q} that is both flexible enough to generate good approximations and restrictive enough to support efficient calculations?

Standard choice when working with probability distributions

The **Kullback-Leibler divergence** is the standard distance measure:

$$\text{KL}(f||g) = \int_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) \log \left(\frac{f(\boldsymbol{\theta})}{g(\boldsymbol{\theta})} \right) d\boldsymbol{\theta} = \mathbb{E}_{\boldsymbol{\theta} \sim f} \left[\log \left(\frac{f(\boldsymbol{\theta})}{g(\boldsymbol{\theta})} \right) \right].$$

Notice that while $\text{KL}(f||g)$ obeys the positivity criterion, it satisfies neither symmetry nor the triangle inequality. It is thus **not a proper distance measure**.

Two alternative KL definitions: $\text{KL}(q||p)$ or $\text{KL}(p||q)$?

Information-projection

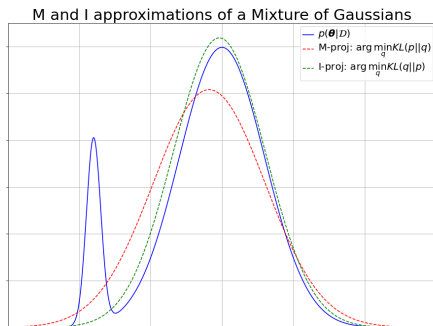
- Minimizes $\text{KL}(q||p) = -\mathbb{E}_{\theta \sim q}[\log p(\theta | \mathcal{D})] - \mathcal{H}_q$.
- Preference given to q that has:
 - 1 High q -probability allocated to p -probable regions.
 - 2 Small q in any region where p is small.
“ $p(\theta | \mathcal{D}) \approx 0 \implies q(\theta) \approx 0$ ”.
 - 3 High entropy (\sim variance)

Moment-projection

- Minimizes $\text{KL}(p||q) = -\mathbb{E}_{\theta \sim p}[\log q(\theta)] - \mathcal{H}_p$.
- Preference given to q that has:
 - 1 High q -probability allocated to p -probable regions.
 - 2 $q(\theta) > 0$ in any region where p is non-negligible.
“ $p(\theta | \mathcal{D}) > 0 \implies q(\theta) > 0$ ”
 - 3 No explicit focus of entropy

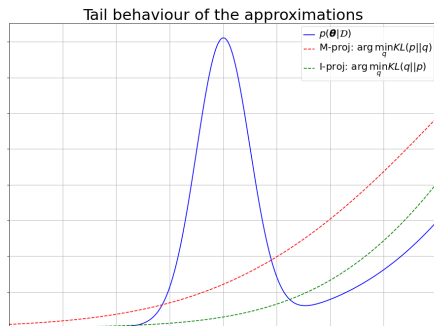
Cheat-sheet:

- **KL-divergence:** $\text{KL}(f||g) = \mathbb{E}_f \left[\log \left(\frac{f(\theta)}{g(\theta)} \right) \right] = -\mathbb{E}_f [\log (g(\theta))] - \mathcal{H}_f$.
- **Entropy:** $\mathcal{H}_f = -\int_{\theta} f(\theta) \log (f(\theta)) d\theta = -\mathbb{E}_f [\log (f(\theta))]$.
- **Intuition:** Cheat a bit (measure-zero, limit-zero-rates, etc.) and think “If $g(\theta_0) \approx 0$, then $-\mathbb{E}_{\theta \sim f}[\log g(\theta)]$ becomes ‘huge’ unless $f(\theta_0) \approx 0$ ” because $\lim_{x \rightarrow 0^+} \log(x)$ diverges, while $\lim_{x \rightarrow 0^+} x \cdot \log(x) = 0$.



Example: Approximating a Mix-of-Gaussians by a single Gaussian

- **Moment projection** – optimizing $KL(p||q)$ – has slightly larger variance.
- Similar mean values, but **Information projection** – optimizing $KL(q||p)$ – focuses mainly on the most prominent mode.



Example: Approximating a Mix-of-Gaussians by a single Gaussian

- **Moment projection** – optimizing $\text{KL}(p||q)$ – has slightly larger variance.
- Similar mean values, but **Information projection** – optimizing $\text{KL}(q||p)$ – focuses mainly on the most prominent mode.
- **M-projection** is **zero-avoiding**, while **I-projection** is **zero-forcing**.

VB uses information projections:

Variational Bayes relies on **information projections**, i.e., approximates $p(\boldsymbol{\theta} \mid \mathcal{D})$ by

$$\hat{q}(\boldsymbol{\theta}) = \arg \min_{q \in \mathcal{Q}} \text{KL}(q(\boldsymbol{\theta}) \parallel p(\boldsymbol{\theta} \mid \mathcal{D}))$$

- **Positives:**

- Clever interpretation when used for Bayesian machine learning.
 - We will end up with an objective that lower-bounds the marginal log likelihood, $\log p(\mathcal{D})$.
- Very efficient when combined with cleverly chosen \mathcal{Q} .

- **Negatives:**

- May result in *zero-forcing* behaviour.
 - Typical choice of \mathcal{Q} can make this issue even more prominent.

Notice how we can rearrange the KL divergence as follows:

$$\text{KL} (q(\boldsymbol{\theta})||p(\boldsymbol{\theta} | \mathcal{D})) \quad = \quad \mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta} | \mathcal{D})} \right]$$

Notice how we can rearrange the KL divergence as follows:

$$\text{KL} (q(\boldsymbol{\theta})||p(\boldsymbol{\theta} | \mathcal{D})) = \mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta} | \mathcal{D})} \right] = \mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{q(\boldsymbol{\theta}) \cdot p(\mathcal{D})}{p(\boldsymbol{\theta} | \mathcal{D}) \cdot p(\mathcal{D})} \right]$$

Notice how we can rearrange the KL divergence as follows:

$$\begin{aligned}\text{KL}(q(\boldsymbol{\theta})||p(\boldsymbol{\theta}|\mathcal{D})) &= \mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}|\mathcal{D})} \right] = \mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{q(\boldsymbol{\theta}) \cdot p(\mathcal{D})}{p(\boldsymbol{\theta}|\mathcal{D}) \cdot p(\mathcal{D})} \right] \\ &= \log p(\mathcal{D}) - \mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}, \mathcal{D})} \right]\end{aligned}$$

ELBO: Evidence Lower-Bound

Notice how we can rearrange the KL divergence as follows:

$$\begin{aligned}\text{KL}(q(\boldsymbol{\theta})||p(\boldsymbol{\theta}|\mathcal{D})) &= \mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}|\mathcal{D})} \right] = \mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{q(\boldsymbol{\theta}) \cdot p(\mathcal{D})}{p(\boldsymbol{\theta}|\mathcal{D}) \cdot p(\mathcal{D})} \right] \\ &= \log p(\mathcal{D}) - \mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}, \mathcal{D})} \right] = \log p(\mathcal{D}) - \mathcal{L}(q)\end{aligned}$$

Evidence Lower Bound (ELBO): $\mathcal{L}(q) = -\mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}, \mathcal{D})} \right] = \mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{p(\boldsymbol{\theta}, \mathcal{D})}{q(\boldsymbol{\theta})} \right].$

Notice how we can rearrange the KL divergence as follows:

$$\begin{aligned}\text{KL}(q(\boldsymbol{\theta})||p(\boldsymbol{\theta}|\mathcal{D})) &= \mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}|\mathcal{D})} \right] = \mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{q(\boldsymbol{\theta}) \cdot p(\mathcal{D})}{p(\boldsymbol{\theta}|\mathcal{D}) \cdot p(\mathcal{D})} \right] \\ &= \log p(\mathcal{D}) - \mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}, \mathcal{D})} \right] = \log p(\mathcal{D}) - \mathcal{L}(q)\end{aligned}$$

Evidence Lower Bound (ELBO): $\mathcal{L}(q) = -\mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}, \mathcal{D})} \right] = \mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{p(\boldsymbol{\theta}, \mathcal{D})}{q(\boldsymbol{\theta})} \right].$

VB focuses on ELBO:

$$\log p(\mathcal{D}) = \mathcal{L}(q) + \text{KL}(q(\boldsymbol{\theta})||p(\boldsymbol{\theta}|\mathcal{D}))$$

Since $\log p(\mathcal{D})$ is constant wrt. the distribution q it follows:

- We can minimize $\text{KL}(q(\boldsymbol{\theta})||p(\boldsymbol{\theta}|\mathcal{D}))$ by maximizing $\mathcal{L}(q)$
- This is **computationally simpler** because it uses $p(\boldsymbol{\theta}, \mathcal{D})$ and not $p(\boldsymbol{\theta}|\mathcal{D})$.
- $\mathcal{L}(q)$ is a **lower bound** of $\log p(\mathcal{D})$ because $\text{KL}(q(\boldsymbol{\theta})||p(\boldsymbol{\theta}|\mathcal{D})) \geq 0$.

\rightsquigarrow Look for $\hat{q}(\boldsymbol{\theta}) = \arg \max_{q \in \mathcal{Q}} \mathcal{L}(q)$.

Notice how we can rearrange the KL divergence as follows:

$$\begin{aligned}\text{KL}(q(\boldsymbol{\theta})||p(\boldsymbol{\theta}|\mathcal{D})) &= \mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}|\mathcal{D})} \right] = \mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{q(\boldsymbol{\theta}) \cdot p(\mathcal{D})}{p(\boldsymbol{\theta}|\mathcal{D}) \cdot p(\mathcal{D})} \right] \\ &= \log p(\mathcal{D}) - \mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}, \mathcal{D})} \right] = \log p(\mathcal{D}) - \mathcal{L}(q)\end{aligned}$$

Evidence Lower Bound (ELBO): $\mathcal{L}(q) = -\mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}, \mathcal{D})} \right] = \mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{p(\boldsymbol{\theta}, \mathcal{D})}{q(\boldsymbol{\theta})} \right].$

Summary:

- We started out looking for $\arg \min_{q \in \mathcal{Q}} \text{KL}(q(\boldsymbol{\theta})||p(\boldsymbol{\theta}|\mathcal{D}))$.
- Didn't know how to calculate $\text{KL}(q(\boldsymbol{\theta})||p(\boldsymbol{\theta}|\mathcal{D}))$ because $p(\boldsymbol{\theta}|\mathcal{D})$ is unknown.
- Still, we can find the optimal approximation by maximizing $\mathcal{L}(q)$:

$$\arg \max_{q \in \mathcal{Q}} \mathcal{L}(q) = \arg \min_{q \in \mathcal{Q}} \text{KL}(q(\boldsymbol{\theta})||p(\boldsymbol{\theta}|\mathcal{D})).$$

- It all makes sense: We aim to maximize $\mathcal{L}(q)$, which is a lower-bound of $\log p(\mathcal{D})$.

Variational Bayes w/ Mean Field

What we have ...

We now have the first building-block of the approximation:

$$\Delta(q || p) = \text{KL} (q(\boldsymbol{\theta}) || p(\boldsymbol{\theta} | \mathcal{D})) ,$$

and avoided the issue with $p(\boldsymbol{\theta} | \mathcal{D})$ by focusing on $\mathcal{L}(q)$.

We still need the set \mathcal{Q} :

Very often you will see the **mean field assumption**, which states that \mathcal{Q} consists of distributions that **factorize** according to the equation

$$q(\boldsymbol{\theta} | \boldsymbol{\lambda}) = \prod_i q_i(\theta_i | \lambda_i).$$

This may seem like a very restricted set, but it often works well anyway ...

Setup:

- We have observed \mathcal{D} , and can calculate the full joint $p(\boldsymbol{\theta}, \mathcal{D}) = p(\boldsymbol{\theta}) \cdot p(\mathcal{D} | \boldsymbol{\theta})$.
 - We use the ELBO as our objective, and assume $q(\boldsymbol{\theta})$ factorizes.
 - We posit a *variational family* of distributions $q_i(\cdot | \boldsymbol{\lambda}_i)$, i.e., we choose the distributional form, while wanting to optimize the parameterization $\boldsymbol{\lambda}_i$.
-
- We then aim to solve the following continuous maximization problem:

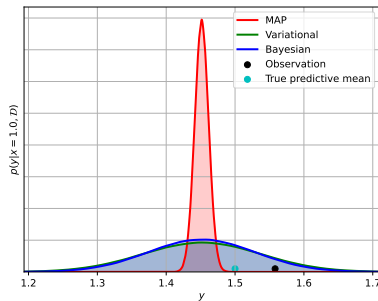
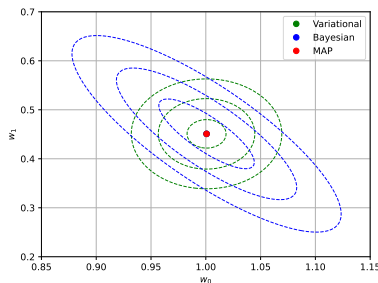
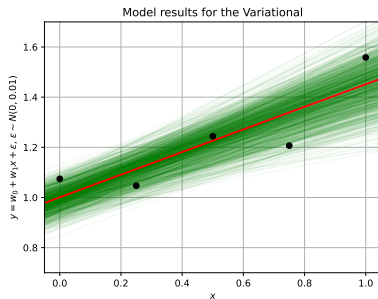
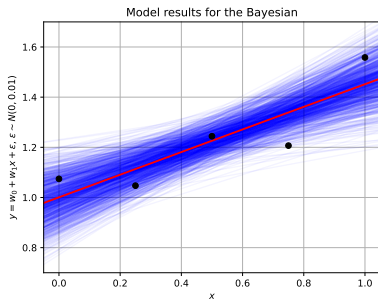
$$\arg \max_{\boldsymbol{\lambda}} \mathcal{L}(\boldsymbol{\lambda}) = \arg \min_{q \in \mathcal{Q}} \text{KL} (q(\boldsymbol{\theta}) || p(\boldsymbol{\theta} | \mathcal{D})) .$$

(Stochastic) Gradient ascent algorithm for maximizing a function $\mathcal{L}(\boldsymbol{\lambda})$:

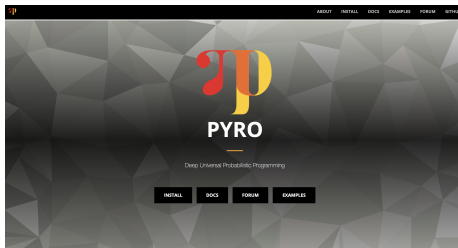
- 1 Initialize $\boldsymbol{\lambda}^{(0)}$ randomly.
- 2 For $t = 1, \dots$:

$$\boldsymbol{\lambda}^{(t)} \leftarrow \boldsymbol{\lambda}^{(t-1)} + \rho_t \cdot \mathcal{L}' \left(\boldsymbol{\lambda}^{(t-1)} \right)$$

Bayes linear regression with likelihood $y_i \mid \{w_0, w_1, x_i, \sigma_y^2\} = \mathcal{N}(w_0 + w_1 x_i, \sigma_y^2)$.



Probabilistic programming: Pyro



Pyro's main features (www.pyro.ai) :

- Initially developed by UBER (the car riding company).
- Community of contributors and a dedicated team at Broad Institute (US).
- Rely on Pytorch (Deep Learning Framework).
- Enable GPU accelaration and distributed learning.

Pyro

Pyro (pyro.ai) is a Python library for probabilistic machine learning integrated with PyTorch.

- Modeling:**
 - Directed graphical models
 - Neural networks (via `nn.Module`)
 - ...
- Inference:**
 - Variational inference – including BBVI, SVI
 - Monte Carlo – including Importance sampling and Hamiltonian Monte Carlo
 - ...
- Criticism:**
 - Point-based evaluations
 - Posterior predictive checks
 - ...

<https://github.com/PGM-Lab/2025-MSE-AI>