# Introduction to Probabilistic Machine Learning Summer School on Methods for Statistical Evaluation of Al

Andrés Masegosa

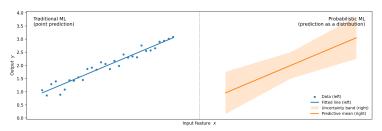
[Material made by Thomas D. Nielsen and Helge Langseth]

MseAl - 2025

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).

## Why Uncertainty Matters in Practice

# 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 Al.

# Bayesian Machine Learning

## Bayesian Machine Learning = Probabilistic model + Bayesian inference

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

# Bayesian Machine Learning

#### Bayesian Machine Learning = Probabilistic model + Bayesian inference

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

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

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

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

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

MseAl - 2025 Introduction

# Bayesian Machine Learning

## Bayesian Machine Learning = Probabilistic model + Bayesian inference

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

Now we can calculate the posterior over  $\theta$  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 x':

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

Being Bayesian means maintaining a distribution over  $\theta$ .

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

MseAl - 2025 Introduction

## Example: Linear regression

A Bayesian linear regression with univariate explanatory variables:

$$\textbf{Likelihood} - p(\mathcal{D} \mid \boldsymbol{\theta}) \textbf{:} \quad p(y_i \mid x_i, \mathbf{w}, \sigma_y^2) = \mathcal{N} \left( w_0 + w_1 \cdot x_i, \sigma_y^2 \right)$$

**Note!** The observation noise,  $\sigma_y^2$ , 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:

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

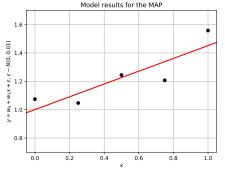
MseAl - 2025 Introduction

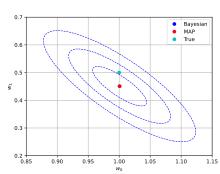
# Example: Linear regression – MAP vs (fully) Bayesian

## Bayes linear regression with some fake data:

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

## Results for the MAP and the fully Bayesian model:





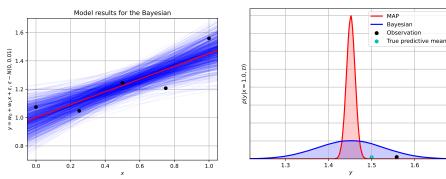
- MAP: Reasonable point estimate; No model uncertainty;
- Bayes: Model uncertainty around same MAP estimate;

# Example: Linear regression – MAP vs (fully) Bayesian

## Bayes linear regression with some fake data:

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

## 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 – Summary

Bayesian inference is in principle easy using Bayes' rule:

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

**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 $p(\theta \mid \mathcal{D})$

#### What we want:

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

#### What we don't want:

- Non scalable solutions;
- Widely applicable.

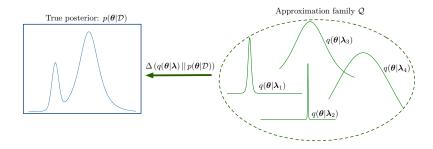
#### **FUNDAMENTAL** assumption:

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  $p(\mathcal{D}, \boldsymbol{\theta}) = p(\boldsymbol{\theta}) \cdot p(\mathcal{D} \mid \boldsymbol{\theta}) = p(\boldsymbol{\theta}) \prod_i p(\mathbf{x}_i \mid \boldsymbol{\theta})$ .



## Approximate inference through optimization – Main idea

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



# Approximate inference through optimization

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

## Formalization of approximate inference through optimization:

Given a family of tractable distributions  $\mathcal Q$  and a distance measure between distributions  $\Delta,$  choose

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

#### Decisions to be made:

- lacktriangle How to define  $\Delta(\cdot||\cdot)$  so that we end up with a high-quality solution?
  - ullet How to work with  $\Deltaig(q(m{ heta})\,||\,p(m{ heta}\,|\,\mathcal{D})ig)$  when we don't know what  $p(m{ heta}\,|\,\mathcal{D})$  is?
- ${\cal Q}$  How to define a family of distributions  ${\cal Q}$  that is both flexible enough to generate good approximations and restrictive enough to support efficient calculations?

#### Distance measure

#### Standard choice when working with probability distributions

The Kullback-Leibler divergence is the standard distance measure:

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

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

## Two alternative KL definitions: KL(q||p) or KL(p||q)?

#### Information-projection

- Minimizes  $\mathrm{KL}\left(q||p\right) = -\mathbb{E}_{\boldsymbol{\theta} \sim q}[\log p(\boldsymbol{\theta}\,|\,\mathcal{D})] \mathcal{H}_q.$
- Preference given to q that has:
  - High q-probability allocated to p-probable regions.
  - Small q in any region where p is small.

"
$$p(\boldsymbol{\theta} \mid \mathcal{D}) \approx 0 \implies q(\boldsymbol{\theta}) \approx 0$$
".

**1** High entropy ( $\sim$  variance)

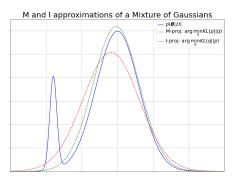
## Moment-projection

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

#### Cheat-sheet:

- KL-divergence:  $\mathrm{KL}\left(f||g\right) = \mathbb{E}_f\left[\log\left(\frac{f(\pmb{\theta})}{g(\pmb{\theta})}\right)\right] = -\mathbb{E}_f\left[\log\left(g(\pmb{\theta})\right)\right] \mathcal{H}_f.$
- Entropy:  $\mathcal{H}_f = -\int_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) \log (f(\boldsymbol{\theta})) d\boldsymbol{\theta} = -\mathbb{E}_f [\log (f(\boldsymbol{\theta}))].$
- 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  $\lim_{x \to 0^+} \log(x)$  diverges, while  $\lim_{x \to 0^+} x \cdot \log(x) = 0$ .

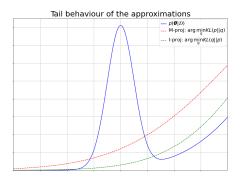
## Moment and Information projection – main difference



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

- Similar mean values, but Information projection optimizing  $\mathrm{KL}\,(q||p)$  focuses mainly on the most prominent mode.

## Moment and Information projection – main difference



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

- Similar mean values, but Information projection optimizing  $\mathrm{KL}\,(q||p)$  focuses mainly on the most prominent mode.
- M-projection is zero-avoiding, while I-projection is zero-forcing.

## Variational Bayes setup

#### VB uses information projections:

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

$$\hat{q}(\boldsymbol{\theta}) = \arg\min_{q \in \mathcal{Q}} \mathrm{KL}\left(q(\boldsymbol{\theta})||p(\boldsymbol{\theta} \mid \mathcal{D})\right)$$

#### 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 Q.

#### Negatives:

- May result in zero-forcing behaviour.
  - Typical choice of Q can make this issue even more prominent.

Notice how we can rearrange the KL divergence as follows:

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

Notice how we can rearrange the KL divergence as follows:

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

Notice how we can rearrange the KL divergence as follows:

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

$$= \log p(\mathcal{D}) - -\mathbb{E}_{\boldsymbol{\theta}\sim q}\left[\log \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta},\mathcal{D})}\right]$$

Notice how we can rearrange the KL divergence as follows:

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

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

Notice how we can rearrange the KL divergence as follows:

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

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

#### VB focuses on ELBO:

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

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

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

$$\rightsquigarrow$$
 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:

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

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

#### **Summary:**

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

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

• 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

## The mean field assumption

#### 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  $p(\theta \mid \mathcal{D})$  by focusing on  $\mathcal{L}(q)$ .

#### We still need the set 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 . . .

# Wrapping it all up: The VB algorithm under MF

# Setup:

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

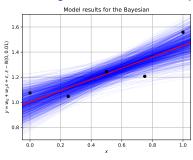
$$\arg\max_{\boldsymbol{\lambda}}\mathcal{L}(\boldsymbol{\lambda}) = \arg\min_{q \in \mathcal{Q}} \mathrm{KL}\left(q(\boldsymbol{\theta})||p(\boldsymbol{\theta} \mid \mathcal{D})\right).$$

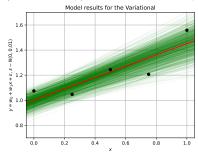
## (Stochastic) Gradient ascent algorithm for maximizing a function L ( $\lambda$ ):

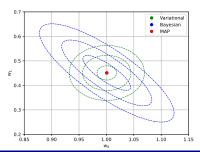
- Initialize  $\lambda^{(0)}$  randomly.
- ② For t = 1, ...:

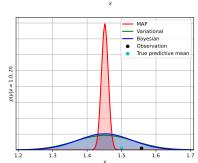
$$\boldsymbol{\lambda}^{(t)} \leftarrow \boldsymbol{\lambda}^{(t-1)} + \rho_t \cdot \mathcal{L}\left(\boldsymbol{\lambda}^{(t-1)}\right)$$

**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