

Advanced Probabilistic Machine Learning SSY316

Bayesian Linear Regression

Alexandre Graell i Amat

`alexandre.graell@chalmers.se`

<https://sites.google.com/site/agraellamat>

October 31 and November 7, 2023



CHALMERS

Supervised machine learning

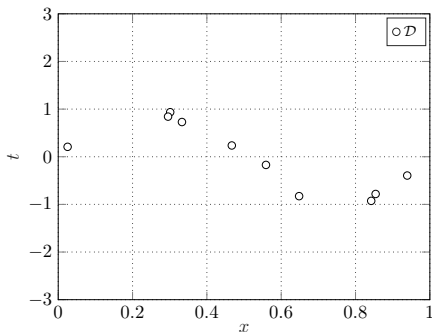
Goal: Given a training set $\mathcal{D} = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N)\}$, identify an algorithm to predict the label t for a new (yet unseen) point \mathbf{x}

- Learning a model from labeled data
 - Predicting output of new data based on the learned model
-
- \mathbf{x}_i : free variables (features, covariates, domain points, explanatory variables)
 - t_i : target variables, dependent on \mathbf{x}_i (dependent variables, labels, responses)

Supervised machine learning

Goal: Given a training set $\mathcal{D} = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N)\}$, identify an algorithm to predict the label t for a new (yet unseen) point \mathbf{x} .

Impossible if no information on the mechanism relating \mathbf{x} and t !



We may assume that \mathbf{x} and t are related via a function $t = \tilde{f}(\mathbf{x})$

Goal: Find best possible approximation of $\tilde{f}(\mathbf{x})$, $f(\mathbf{x})$ (prediction model).
Predict the outcome t for \mathbf{x} as $\hat{t} = f(\mathbf{x})$.

Supervised machine learning

Goal: Identify a predictive algorithm that minimizes the error in the prediction of a new label t for an unobserved x (generalization loss).

Two approaches:

- Frequentist approach
- Bayesian approach

Rephrase supervised machine learning as a Bayesian inference problem.

Inference

Learning: Treat \mathbf{x} and t as random variables

Predict t given the observation of \mathbf{x} under assumption joint distribution $p(\mathbf{x}, t)$ is known.

- Loss function $\ell(t, \hat{t})$: cost (loss or risk) incurred when the correct value is t while the estimate is \hat{t}
- Quadratic loss function:

$$\ell(t, \hat{t}) = (t - \hat{t})^2$$

- Optimal prediction $\hat{t}^*(\mathbf{x})$: $\hat{t}(\mathbf{x})$ that minimizes generalization loss (generalization error)

$$L_p(\hat{t}) = \mathbb{E}_{\mathbf{x}t \sim p_{\mathbf{x}, t}}[\ell(t, \hat{t}(\mathbf{x}))]$$

Inference

- **Solution:**

$$\begin{aligned}\hat{t}^*(\mathbf{x}) &= \arg \min_{\hat{t}} L_p(\hat{t}) \\ &= \arg \min_{\hat{t}} \mathbb{E}_{\mathbf{x}, \mathbf{t} \sim p_{\mathbf{x}, \mathbf{t}}} [\ell(\mathbf{t}, \hat{t}(\mathbf{x}))] \\ &= \arg \min_{\hat{t}} \mathbb{E}_{\mathbf{x} \sim p_{\mathbf{x}}} \left[\mathbb{E}_{\mathbf{t} \sim p_{\mathbf{t}|\mathbf{x}}} [\ell(\mathbf{t}, \hat{t}(\mathbf{x}))] \right] \\ &= \arg \min_{\hat{t}} \mathbb{E}_{\mathbf{t} \sim p_{\mathbf{t}|\mathbf{x}}} [\ell(\mathbf{t}, \hat{t}(\mathbf{x}))]\end{aligned}$$

Optimal prediction a function of $p(\mathbf{t}|\mathbf{x})$ and the loss function.

- For the loss function $\ell(t, \hat{t}) = (t - \hat{t})^2$,

$$\hat{t}^*(\mathbf{x}) = \mathbb{E}_{\mathbf{t}|\mathbf{x}}[\mathbf{t}|\mathbf{x}]$$

Supervised machine learning

Goal (revisited): Obtaining predictor $\hat{t}(\mathbf{x})$ that performs close to optimal predictor $\hat{t}^*(\mathbf{x})$ based only on training set \mathcal{D} (without knowledge of $p(\mathbf{x}, t)$).

- Closeness measured as

$$L_p(\hat{t}) - L_p(\hat{t}^*)$$

Supervised machine learning: Frequentist approach

- Training data points $(\mathbf{x}_i, \mathbf{t}_i) \in \mathcal{D}$ are i.i.d RVs drawn from a true (unknown) distribution $p(\mathbf{x}, t)$,

$$(\mathbf{x}_i, \mathbf{t}_i) \sim_{\text{i.i.d}} p(\mathbf{x}, t) \quad i = 1, \dots, N$$

- Since $p(\mathbf{x}, t)$ unknown, cannot find optimal prediction via

$$\hat{t}^*(\mathbf{x}) = \arg \min_{\hat{t}} \mathbb{E}_{\mathbf{t}|\mathbf{x}}[\ell(\mathbf{t}, \hat{t})|\mathbf{x}]$$

Supervised machine learning: Frequentist approach

Solutions:

1. **Separate learning and inference:** Learn an approximation of $p(t|\mathbf{x})$ based on \mathcal{D} ($p_{\mathcal{D}}(t|\mathbf{x})$) and use it in

$$\hat{t}^*(\mathbf{x}) = \arg \min_{\hat{t}} \mathbb{E}_{t|\mathbf{x}}[\ell(t, \hat{t})|\mathbf{x}]$$

to obtain

$$\hat{t}_{\mathcal{D}}(\mathbf{x}) = \arg \min_{\hat{t}} \mathbb{E}_{t \sim p_{\mathcal{D}}(t|\mathbf{x})}[\ell(t, \hat{t})|\mathbf{x}]$$

2. **Direct inference via empirical risk minimization (ERM):** Learn directly an approximation of optimal decision rule ($\hat{t}_{\mathcal{D}}(\cdot)$) by minimizing an empirical estimate of generalization loss,

$$\hat{t}_{\mathcal{D}}^*(\mathbf{x}) = \arg \min_{\hat{t}} L_{\mathcal{D}}(\hat{t}) = \arg \min_{\hat{t}} \frac{1}{N} \sum_{i=1}^N \ell(t_i, \hat{t}(\mathbf{x}_i))$$

Separate learning and inference

How do we learn an approximation $p_{\mathcal{D}}(t|\mathbf{x})$ of $p(t|\mathbf{x})$ based on \mathcal{D} ?

Idea:

- Select a family of parametric probabilistic models (hypothesis class)
- Learn model parameters to fit \mathcal{D}

Linear regression

- We are interested in values of a function $t(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}$,
 $\mathbf{x} = (x_1, \dots, x_d)^\top$
- We have some observations of this mapping, $\mathcal{D} = \{(\mathbf{x}_i, t_i)\}_{i=1}^N$

Goal: predict t for new input $\mathbf{x} \rightarrow$ Learn an accurate prediction function $\hat{t}(\mathbf{x})$ (regression function) from \mathcal{D} .

Linear regression:

Assumes relationship between \mathbf{x} and t is linear:

$$\begin{aligned} t(\mathbf{x}) &= w_0 + w_1 x_1 + \dots + w_d x_d \\ &= \mu(\mathbf{x}, \mathbf{w}) \end{aligned}$$

Linear regression

Linear regression (2):

Linear combinations of fixed nonlinear functions of the input variables,

$$\mu(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^M w_j \phi_j(\mathbf{x})$$

- $\phi_j(\mathbf{x})$: basis functions
- w_0 : bias, allows for any fixed offset
- Convenient to define dummy basis function $\phi_0(\mathbf{x}) = 1$, so that

$$\mu(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^M w_j \phi_j(\mathbf{x}) = \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x})$$

with $\mathbf{w} = (w_0, \dots, w_M)^\top$ and $\boldsymbol{\phi}(\mathbf{x}) = (\phi_0(\mathbf{x}), \dots, \phi_M(\mathbf{x}))^\top$

Linear regression IBM

Examples of basis functions

- Polynomial linear regression:

$$\mu(x, \mathbf{w}) = \sum_{j=0}^M w_j x^j = \mathbf{w}^\top \phi(x)$$

with $\mathbf{w} = (w_0, \dots, w_M)^\top$ and $\phi(x) = (1, x, x^2, \dots, x^M)^\top$

- Gaussian basis functions:

$$\phi_j(x) = \exp\left(-\frac{(x - \alpha_j)^2}{2s^2}\right)$$

- Sigmoidal basis functions:

$$\phi_j(x) = \sigma\left(\frac{x - \alpha_j}{s}\right)$$

with

$$\sigma(a) = \frac{1}{1 + e^{-a}}$$

Linear regression

- Can express our **uncertainty** over value of $\mu(\mathbf{x}, \mathbf{w})$ as

$$t(\mathbf{x}, \mathbf{w}) = \mathbf{w}^\top \phi(\mathbf{x}) + \varepsilon = \mu(\mathbf{x}, \mathbf{w}) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \beta^{-1})$$

- Given \mathbf{x} , t modeled by **probabilistic model**

$$t|\mathbf{x} = \mathbf{x} \sim \mathcal{N}(\mu(\mathbf{x}, \mathbf{w}), \beta^{-1})$$

equivalently,

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = p(t|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(t|\mu(\mathbf{x}, \mathbf{w}), \beta^{-1}), \quad \boldsymbol{\theta} = (\mathbf{w}, \beta)$$

Discriminative vs. generative probabilistic models

Discriminative probabilistic models:

- Posterior (predictive) distribution $p(t|\mathbf{x})$ assumed to belong to a hypothesis class $p(t|\mathbf{x}, \boldsymbol{\theta})$ defined by $\boldsymbol{\theta}$
- $\boldsymbol{\theta}$ learned from \mathcal{D}
- Once model is learned, compute

$$\hat{t}_{\mathcal{D}}(\mathbf{x}) = \arg \min_{\hat{t}} \mathbb{E}_{t \sim p_{\mathcal{D}}(t|\mathbf{x})} [\ell(t, \hat{t}) | \mathbf{x}]$$

Example: Linear regression

1. Learn $\boldsymbol{\theta} = (\mathbf{w}, \beta)$ based on \mathcal{D}
2. Optimal prediction (under quadratic loss function):

$$\hat{t}_{\mathcal{D}}(\mathbf{x}) = \mathbb{E}_{t \sim p(t|\mathbf{x}, \boldsymbol{\theta}_{\mathcal{D}}^*)} [t | \mathbf{x}]$$

i.e.,

$$\hat{t}_{\mathcal{D}}(\mathbf{x}) = \mu(\mathbf{x}, \mathbf{w}_{\mathcal{D}}^*)$$

Discriminative vs. generative probabilistic models

Generative probabilistic models:

- Models $p(\mathbf{x}, t)$ as being part of a parametric family $p(\mathbf{x}, t|\boldsymbol{\theta})$
- Models also $p(\mathbf{x})$
- **Generative**: capacity to generate a realization of \mathbf{x} by using marginal $p(\mathbf{x}|\boldsymbol{\theta})$
- Once model is learned, obtain $p(t|\mathbf{x}, \boldsymbol{\theta})$ applying Bayes' and compute

$$\hat{t}_{\mathcal{D}}(\mathbf{x}) = \arg \min_{\hat{t}} \mathbb{E}_{t \sim p_{\mathcal{D}}(t|\mathbf{x})} [\ell(t, \hat{t})|\mathbf{x}]$$

Observations:

- Make **stronger assumptions** \rightarrow may lead to more significant **bias**
- Ability to deal with **missing data** or **latent variables** (**semi-supervised learning**)

Maximum likelihood learning

- **Hypothesis class** $p(t|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(\mu(\mathbf{x}, \mathbf{w}), \beta^{-1})$: e.g., definition of the polynomial degree M (defines *capacity* of the class)
- **Specific model** $p(t|\mathbf{x}, \boldsymbol{\theta})$: selection of $\boldsymbol{\theta} = (\mathbf{w}, \beta)$ (learned from \mathcal{D})

How do we learn the model parameters $\boldsymbol{\theta}$ (for a given hypothesis class)?

Idea: Maximum likelihood (ML) learning

Select $\boldsymbol{\theta}$ such that \mathcal{D} has *maximum probability* of being observed

- Need to write distribution observed labels $t_{\mathcal{D}} \in \mathcal{D}$ given $x_{\mathcal{D}}$,

$$\begin{aligned} p(t_{\mathcal{D}} | x_{\mathcal{D}}, \mathbf{w}, \beta) &= \prod_{i=1}^N p(t_i | \mathbf{x}_i, \mathbf{w}, \beta) \\ &= \prod_{i=1}^N \mathcal{N}(t_i | \mu(\mathbf{x}_i, \mathbf{w}), \beta^{-1}) \end{aligned}$$

Maximum likelihood learning

Log-likelihood function:

$$\begin{aligned}\ln p(t_{\mathcal{D}}|x_{\mathcal{D}}, \mathbf{w}, \beta) &= \sum_{i=1}^N \ln p(t_i|x_i, \mathbf{w}, \beta) \\ &= -\frac{\beta}{2} \sum_{i=1}^N (t_i - \mu(\mathbf{x}_i, \mathbf{w}))^2 + \frac{N}{2} \ln \frac{\beta}{2\pi}\end{aligned}$$

ML criterion (cross-entropy loss or log-loss):

$$\begin{aligned}(\mathbf{w}_{\text{ML}}, \beta_{\text{ML}}) &= \arg \max_{\mathbf{w}, \beta} p(t_{\mathcal{D}}|x_{\mathcal{D}}, \mathbf{w}, \beta) \\ &= \arg \max_{\mathbf{w}, \beta} \frac{1}{N} \ln p(t_{\mathcal{D}}|x_{\mathcal{D}}, \mathbf{w}, \beta) \\ &= \arg \min_{\mathbf{w}, \beta} -\frac{1}{N} \ln p(t_{\mathcal{D}}|x_{\mathcal{D}}, \mathbf{w}, \beta) \\ &= \arg \min_{\mathbf{w}, \beta} \frac{\beta}{2N} \sum_{i=1}^N (t_i - \mu(\mathbf{x}_i, \mathbf{w}))^2 - \frac{1}{2} \ln \frac{\beta}{2\pi}\end{aligned}$$

Maximum likelihood learning

If only interested in learning **posterior mean**:

$$\begin{aligned} \mathbf{w}_{\text{ML}} &= \arg \min_{\mathbf{w}} \frac{\beta}{2N} \sum_{i=1}^N (t_i - \mu(\mathbf{x}_i, \mathbf{w}))^2 \\ &= \arg \min_{\mathbf{w}} \frac{1}{N} \sum_{i=1}^N (t_i - \mu(\mathbf{x}_i, \mathbf{w}))^2 \\ &= \arg \min_{\mathbf{w}} L_{\mathcal{D}}(\mathbf{w}) \end{aligned}$$

$L_{\mathcal{D}}(\mathbf{w})$: training loss

Criterion coincides with **direct inference via empirical risk minimization** if we parametrize the predictor as $\hat{t}(\mathbf{x}) = \mu(\mathbf{x}, \mathbf{w})$!

Minimizing $L_{\mathcal{D}}(\mathbf{w})$ can be solved in **closed form**: $\mathbf{w}_{\text{ML}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}_{\mathcal{D}}$, with $\Phi = (\phi(\mathbf{x}_1) \ \phi(\mathbf{x}_2) \cdots \phi(\mathbf{x}_N))^T$

Maximum likelihood learning

To predict value of t for new input \mathbf{x} , we use \mathbf{w}_{ML} :

$$\hat{t} = \mathbf{w}_{\text{ML}}^{\top} \phi(\mathbf{x}) = \left((\Phi^{\top} \Phi)^{-1} \Phi^{\top} \mathbf{t}_{\mathcal{D}} \right)^{\top} \phi(\mathbf{x})$$

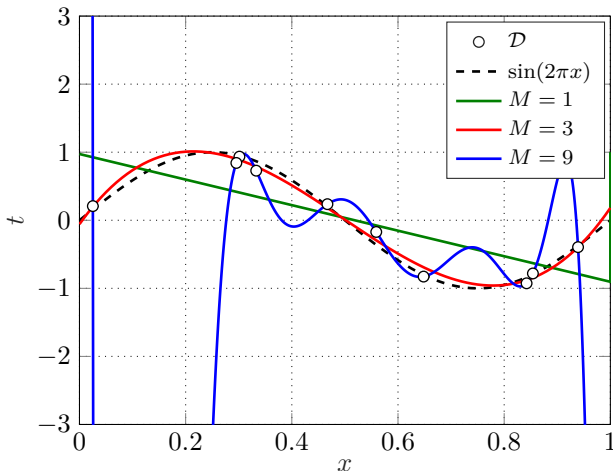
Once \mathbf{w}_{ML} is determined, we can determine β_{ML} as:

$$\beta_{\text{ML}} = \arg \min_{\beta} \frac{\beta}{2N} \sum_{i=1}^N (t_i - \mu(\mathbf{x}_i, \mathbf{w}_{\text{ML}}))^2 - \frac{1}{2} \ln \frac{\beta}{2\pi}$$

i.e.,

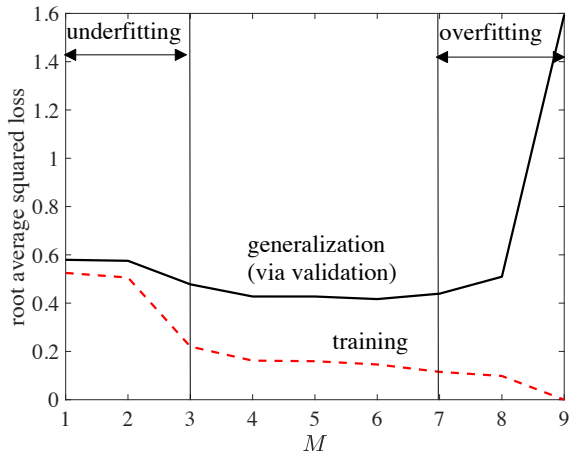
$$\frac{1}{\beta_{\text{ML}}} = \frac{1}{N} \sum_{i=1}^N (t_i - \mu(\mathbf{x}_i, \mathbf{w}_{\text{ML}}))^2 = L_{\mathcal{D}}(\mathbf{w}_{\text{ML}})$$

Maximum likelihood learning



- $p(x, t) = p(x)p(t|x)$, $x \sim \mathcal{U}(0, 1)$, $t|x \sim \mathcal{N}(\sin(2\pi x), 0.1)$
- Optimal predictor under quadratic loss: $\hat{t}^* = \sin(2\pi x)$
- ML predictor: $\hat{t}_{\text{ML}}(x) = \mu(x, w_{\text{ML}})$

Training and generalization loss



- **Training loss:** $L_{\mathcal{D}}(\mathbf{w}_{\text{ML}}) = \frac{1}{N} \sum_{i=1}^N (t_i - \mu(x_i, \mathbf{w}_{\text{ML}}))^2$
- **Generalization loss:** $L_p(\mathbf{w}_{\text{ML}}) = \mathbb{E}_{\mathbf{x}, \mathbf{t}}[\ell(\mathbf{t}, \hat{t}_{\text{ML}}(\mathbf{x}))] = \mathbb{E}_{\mathbf{x}, \mathbf{t}}[\ell(\mathbf{t}, \mathbf{w}_{\text{ML}}^{\text{T}} \phi(\mathbf{x}))]$

Training and generalization loss

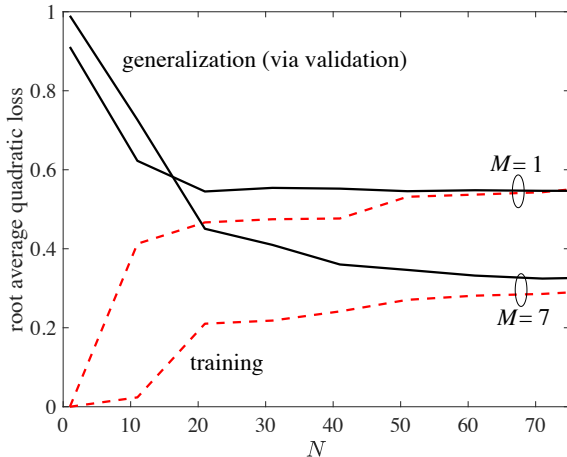
- As the number of data points increases, overfitting is **avoided**
- When \mathcal{D} is **big** compared to the **number of parameters** in θ ,

$$L_{\mathcal{D}}(\mathbf{w}) \simeq L_p(\mathbf{w})$$

- For large N ,

$$\mathbf{w}_{\text{ML}} \longrightarrow \mathbf{w}^* = \arg \min_{\mathbf{w}} L_p(\mathbf{w})$$

Training and generalization loss



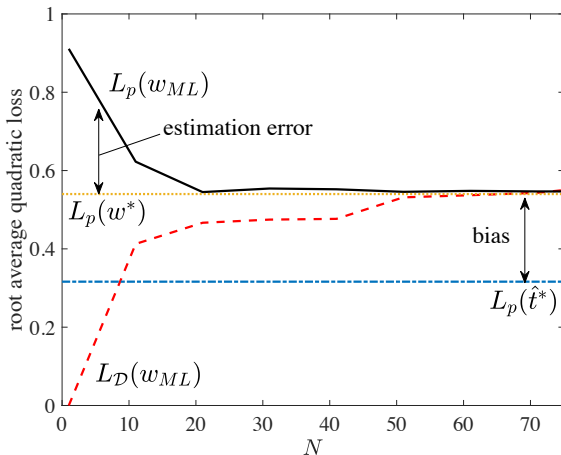
- Squared root of $L_{\mathcal{D}}(\mathbf{w}_{\text{ML}})$ and $L_p(\mathbf{w}_{\text{ML}})$
- Overfitting decreases with increasing N

Bias and estimation error

$$L_p(\mathbf{w}_{\text{ML}}) = L_p(\hat{t}^*) + \underbrace{(L_p(\mathbf{w}^*) - L_p(\hat{t}^*))}_{\text{bias}} + \underbrace{(L_p(\mathbf{w}_{\text{ML}}) - L_p(\mathbf{w}^*))}_{\text{estimation error}}$$

- $L_p(\hat{t}^*)$: generalization loss of optimal predictor (**minimum achievable loss**)
- **Bias** (approximation error): caused by choice of **hypothesis class**
- **Estimation error** (generalization gap): Caused by the fact that N is **not large enough**

Bias and estimation error



- Estimation error decreases with N and vanishes for large N

Learning and validation

We would like to choose model (hyperparameters and parameters) such that generalization error $L_p(\hat{t}) = \mathbb{E}_{\mathbf{x}, \mathbf{t}}[\ell(\mathbf{t}, \hat{t}(\mathbf{x}))]$ is minimized

... but it depends on $p(\mathbf{x}, t)$ (unknown)!

Learning and validation



Validation: Divide available data into three sets

- **Training set** \mathcal{D} : To fit models
- **Validation set** \mathcal{V} : To choose hypothesis class via evaluation of approximation of the generalization error

$$L_p \approx \frac{1}{N_v} \sum_{i=1}^{N_v} \ell(t_i, \mu(\mathbf{x}_i, \mathbf{w}))$$

- For selected hypothesis class retrain θ based on $\mathcal{D} \cup \mathcal{V}$
- **Test set** \mathcal{T} : To produce estimate generalization error obtained with final model

Typically: 50% – 25% – 25%

Cross-validation

Pitfall of validation:

Part of available data not used for training!

Validation suitable when having plenty of data.

Alternative: k -fold cross-validation

1. Randomly partition data points into k partitions (folds)
2. For each partition $\kappa \in \{1, \dots, k\}$, train model over all other $k - 1$ partitions
3. Compute generalization error on the κ -th partition
4. Generalization error estimated as average over all partitions
5. Choose hypothesis class that minimizes estimate of generalization error in step 4

Supervised machine learning: The frequentist approach

Goal: Given a training set $\mathcal{D} = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N)\}$, identify an algorithm to predict the label t for a new (yet unseen) point \mathbf{x} .

Linear regression model:

$$t(\mathbf{x}, \mathbf{w}) = \mathbf{w}^\top \phi(\mathbf{x}) + \varepsilon = \mu(\mathbf{x}, \mathbf{w}) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \beta^{-1})$$

Assumptions:

1. \mathbf{x} : observed random variable
2. t : observed random variable
3. ε : unknown random variable
4. β : (un)known deterministic variable
5. \mathbf{w} : unknown deterministic

Supervised machine learning: The Bayesian approach

Goal: Given a training set $\mathcal{D} = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N)\}$, identify an algorithm to predict the label t for a new (yet unseen) point \mathbf{x} .

Linear regression model:

$$t(\mathbf{x}, \mathbf{w}) = \mathbf{w}^\top \phi(\mathbf{x}) + \varepsilon = \mu(\mathbf{x}, \mathbf{w}) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \beta^{-1})$$
$$\mathbf{w} \sim p(\mathbf{w})$$

Assumptions:

1. \mathbf{x} : observed random variable
2. t : observed random variable
3. ε : unknown random variable
4. β : (un)known deterministic variable
5. \mathbf{w} : unknown random variable

Provides the uncertainty of the prediction

Maximum a posteriori learning

- **Observation:** ML learning introduces a tension between **bias** (larger M) and **estimation error** (smaller M)
- **MAP:** Enables a **finer control** of bias and estimation error

Key idea: Leverage **prior information** available on the behavior of parameters in the absence, or presence, of overfitting.

- **Observation:** A large value of $\|w\|$ a manifestation of **overfitting** \rightarrow introduce **prior** on w that gives **lower probability to larger values**,

$$w \sim \mathcal{N}(0, \alpha^{-1} I)$$

$$\text{i.e., } p(w) = \mathcal{N}(w|0, \alpha^{-1} I) = \left(\frac{\alpha}{2\pi}\right)^{(M+1)/2} \exp\left(-\frac{\alpha}{2} w^T w\right)$$

How do we choose the prior distribution?

Common approach: choose a conjugate prior

- Likelihood function

$$p(t_{\mathcal{D}}|x_{\mathcal{D}}, \mathbf{w}, \beta) = \prod_{i=1}^N \mathcal{N}(t_i | \mu(\mathbf{x}_i, \mathbf{w}), \beta^{-1})$$

Exponential of a quadratic function of \mathbf{w} \longrightarrow Conjugate given by a Gaussian distribution of the form

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$$

Maximum a posteriori learning

MAP criterion:

$$\begin{aligned}(\mathbf{w}_{\text{MAP}}, \beta) &= \arg \max_{\mathbf{w}, \beta} p(t_{\mathcal{D}}, \mathbf{w} | x_{\mathcal{D}}, \beta) \\ &= \arg \max_{\mathbf{w}, \beta} p(\mathbf{w}) \prod_{i=1}^N p(t_i | \mathbf{x}_i, \mathbf{w}, \beta)\end{aligned}$$

Equivalently,

$$(\mathbf{w}_{\text{MAP}}, \beta) = \arg \min_{\mathbf{w}, \beta} - \sum_{i=1}^N \ln p(t_i | \mathbf{x}_i, \mathbf{w}, \beta) - \ln p(\mathbf{w})$$

If β a known constant,

$$\mathbf{w}_{\text{MAP}} = \arg \min_{\mathbf{w}} - \sum_{i=1}^N \ln p(t_i | \mathbf{x}_i, \mathbf{w}, \beta) - \ln p(\mathbf{w})$$

MAP: Equivalent to maximizing posterior distribution of \mathbf{w} given the available data, $p(\mathbf{w} | \mathcal{D})$.

Maximum a posteriori learning

Assuming $\mathbf{w} \sim \mathcal{N}(0, \alpha^{-1} \mathbf{I})$, i.e., $p(\mathbf{w}) = \left(\frac{\alpha}{2\pi}\right)^{(M+1)/2} \exp\left(-\frac{\alpha}{2} \mathbf{w}^\top \mathbf{w}\right)$:

$$\begin{aligned}\mathbf{w}_{\text{MAP}} &= \arg \min_{\mathbf{w}} - \sum_{i=1}^N \ln p(t_i | \mathbf{x}_i, \mathbf{w}, \beta) - \ln p(\mathbf{w}) \\&= \arg \min_{\mathbf{w}} - \left(-\frac{\beta}{2} \sum_{i=1}^N (\mu(\mathbf{x}_i, \mathbf{w}) - t_i)^2 + \frac{N}{2} \ln \frac{\beta}{2\pi} \right) - \ln p(\mathbf{w}) \\&= \arg \min_{\mathbf{w}} \frac{\beta}{2} \sum_{i=1}^N (\mu(\mathbf{x}_i, \mathbf{w}) - t_i)^2 - \ln p(\mathbf{w}) \\&= \arg \min_{\mathbf{w}} \frac{1}{N} \sum_{i=1}^N (\mu(\mathbf{x}_i, \mathbf{w}) - t_i)^2 - \frac{2}{N\beta} \ln p(\mathbf{w}) \\&= \arg \min_{\mathbf{w}} \frac{1}{N} \sum_{i=1}^N (\mu(\mathbf{x}_i, \mathbf{w}) - t_i)^2 + \frac{\alpha}{N\beta} \mathbf{w}^\top \mathbf{w} \\&= \arg \min_{\mathbf{w}} L_{\mathcal{D}}(\mathbf{w}) + \frac{\lambda}{N} \|\mathbf{w}\|^2 \longrightarrow \mathbf{w}_{\text{MAP}} = (\lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{t}_{\mathcal{D}}\end{aligned}$$

Maximum a posteriori learning

MAP: Equivalent to maximizing posterior distribution of \mathbf{w} given the available data, $p(\mathbf{w}|\mathcal{D})$.

Consider $p(\mathbf{w}|t_{\mathcal{D}}, x_{\mathcal{D}}, \beta)$:

- Due to **conjugate Gaussian prior** distribution, **posterior also Gaussian**,

$$p(\mathbf{w}|t_{\mathcal{D}}, x_{\mathcal{D}}, \beta) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N)$$

with

$$\mathbf{m}_N = \mathbf{S}_N (\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{t})$$

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \Phi^T \Phi$$

- We assume $\mathbf{w} \sim \mathcal{N}(0, \alpha^{-1} \mathbf{I})$,

$$\mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{t}_{\mathcal{D}}$$

$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi$$

Maximum a posteriori learning

Applying Bayes',

$$\begin{aligned} p(\mathbf{w}|t_{\mathcal{D}}, x_{\mathcal{D}}, \beta) &= \frac{p(t_{\mathcal{D}}|x_{\mathcal{D}}, \mathbf{w}, \beta)p(\mathbf{w})}{p(t_{\mathcal{D}})} \\ &\propto p(\mathbf{w})p(t_{\mathcal{D}}|x_{\mathcal{D}}, \mathbf{w}, \beta) \\ &= p(\mathbf{w}) \prod_{i=1}^N p(t_i|x_i, \mathbf{w}, \beta) \end{aligned}$$

- $p(\mathbf{w})$: knowledge about \mathbf{w} before observing any data
- $p(t_{\mathcal{D}}|\mathbf{w})$: how likely the observed data is for a particular parameter value
- $p(\mathbf{w}|t_{\mathcal{D}})$: knowledge about \mathbf{w} from the observed data and the model

Maximum a posteriori learning

MAP criterion for \mathbf{w} :

$$\begin{aligned}\mathbf{w}_{\text{MAP}} &= \arg \max_{\mathbf{w}} p(\mathbf{w} | t_{\mathcal{D}}, x_{\mathcal{D}}, \beta) \\ &= \arg \max_{\mathbf{w}} p(\mathbf{w}) \prod_{i=1}^N p(t_i | \mathbf{x}_i, \mathbf{w}, \beta) \\ &= \arg \min_{\mathbf{w}} - \sum_{i=1}^N \ln p(t_i | \mathbf{x}_i, \mathbf{w}, \beta) - \ln p(\mathbf{w}) \\ &= \arg \min_{\mathbf{w}} L_{\mathcal{D}}(\mathbf{w}) + \frac{\lambda}{N} \|\mathbf{w}\|^2\end{aligned}$$

with $\lambda = \alpha/\beta$

Maximum a posteriori learning

$$\mathbf{w}_{\text{MAP}} = \arg \min_{\mathbf{w}} L_{\mathcal{D}}(\mathbf{w}) + \frac{\lambda}{N} \|\mathbf{w}\|^2$$

Observations:

- As $N \rightarrow \infty$ MAP estimate tends to ML estimate
- MAP criterion (ridge regression), modifies ML criterion by adding the quadratic (or Tikhonov) regularization function

$$R(\mathbf{w}) = \|\mathbf{w}\|^2$$

Solution:

$$\mathbf{w}_{\text{MAP}} = \mathbf{m}_N = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}_{\mathcal{D}}$$

Regularization: Decreases model capacity \rightarrow Can control overfitting.

Example: Bayesian linear regression (Bishop Figure 3.7)

Problem: fitting straight line to noisy measurements generated from

$$f(x, \mathbf{a}) = a_0 + a_1 x, \quad a_0 = -0.3, \quad a_1 = 0.5$$

by adding Gaussian noise $\mathcal{N}(0, 0.04)$.

Model:

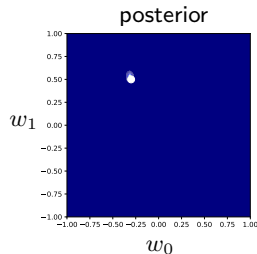
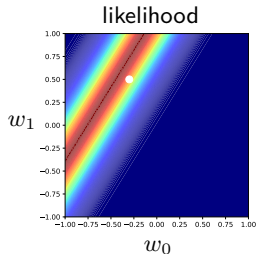
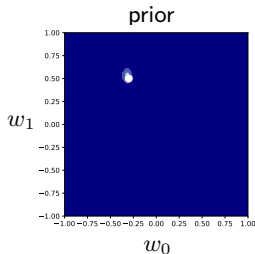
$$t(x, \mathbf{w}) = \underbrace{w_0 + w_1 x}_{\mathbf{w}^\top \mathbf{x}} + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, 0.04)$$

with prior

$$p(\mathbf{w}) = \mathcal{N} \left(\mathbf{w} \mid \begin{pmatrix} 0 & 0 \end{pmatrix}^\top, \alpha^{-1} \mathbf{I}_2 \right), \quad \alpha = 2$$

Example: Bayesian linear regression

Plot after **two hundred** measurements



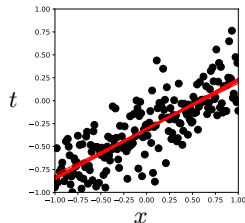
$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_{199}, \mathbf{S}_{199})$$

$$p(t_{200} | \mathbf{w}) = \mathcal{N}(t_{200} | w_0 + w_1 x_{200}, \beta^{-1})$$

$$p(\mathbf{w} | y_1^{200}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_{200}, \mathbf{S}_{200})$$

t

x



Bayesian linear regression

Frequentist approach: Aims at identifying a **specific value** for θ of a probabilistic model to derive a predictor

$$\hat{t}^*(x) = \arg \min_{\hat{t}} \mathbb{E}_{\mathbf{t}|\mathbf{x}}[\ell(\mathbf{t}, \hat{t})|\mathbf{x}]$$

- **ML:** Chooses θ that maximizes probability of training data,
- **MAP:** Includes also prior information about parameter vector

Bayesian approach: Assumes θ **jointly distributed** with data \rightarrow Does not commit to a single value of θ but considers explanations provided by all possible values of θ , each weighted according to a data-dependent **belief**

Bayesian linear regression

Bayesian linear regression model:

$$t(\mathbf{x}, \mathbf{w}) = \mathbf{w}^\top \phi(\mathbf{x}) + \varepsilon = \mu(\mathbf{x}, \mathbf{w}) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \beta^{-1})$$
$$\mathbf{w} \sim p(\mathbf{w})$$

Probabilistic model:

$$p(t_{\mathcal{D}} | x_{\mathcal{D}}, \mathbf{w}, \beta) = \prod_{i=1}^N \mathcal{N}(t_i | \mu(\mathbf{x}_i, \mathbf{w}), \beta^{-1}) \quad \text{likelihood}$$
$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0) \quad \text{prior}$$

End goal: Making predictions of t for new values of \mathbf{x} \rightarrow Directly evaluate the posterior distribution $p(t | t_{\mathcal{D}}, x_{\mathcal{D}}, \mathbf{x}, \beta)$

A fully **Bayesian solution** returns the entire posterior $p(t | \mathcal{D}, \mathbf{x}, \beta)$ \rightarrow provides more information about unobserved label t .

Predictive distribution

We obtain

$$p(t|\mathcal{D}, \mathbf{x}, \beta) = \int \underbrace{p(\mathbf{w}|\mathcal{D}, \beta)}_{\text{posterior dist. of } \mathbf{w}} p(t|\mathbf{x}, \mathbf{w}, \beta) d\mathbf{w}$$

Bayesian approach: $p(t|\mathbf{x}, \mathbf{w}, \beta)$ associated with each value of \mathbf{w} weighted by the posterior belief

$$p(\mathbf{w}|\mathcal{D}, \beta) = \frac{p(\mathbf{w})p(t_{\mathcal{D}}|\mathbf{x}_{\mathcal{D}}, \mathbf{w}, \beta)}{p(t_{\mathcal{D}}|\mathbf{x}_{\mathcal{D}}, \beta)}$$

Predictive distribution

Computing $p(\mathbf{w}|\mathcal{D}, \beta)$ and $p(t|\mathcal{D}, \mathbf{x}, \beta)$ **difficult!**

For our Bayesian linear regression problem with $\mathbf{w} \sim \mathcal{N}(0, \alpha^{-1}\mathbf{I})$:

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|\mu(\mathbf{x}, \mathbf{w}), \beta^{-1})$$

$$p(\mathbf{w}|\mathcal{D}, \beta) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N)$$

with

$$\mathbf{m}_N = \beta \mathbf{S}_N \mathbf{\Phi}^\top \mathbf{t}, \quad \mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \mathbf{\Phi}^\top \mathbf{\Phi}$$

We obtain:

$$p(t|\mathbf{x}, \mathcal{D}, \beta) = \mathcal{N}(t|\mu(\mathbf{x}, \mathbf{w}_{\text{MAP}}), \sigma_N^2(\mathbf{x})),$$

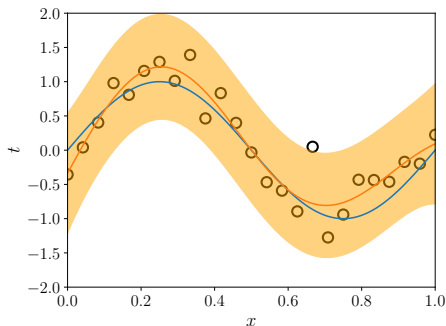
with

$$\mu(\mathbf{x}, \mathbf{w}_{\text{MAP}}) = \mathbf{w}_{\text{MAP}}^\top \phi(\mathbf{x}) = \mathbf{m}_N^\top \phi(\mathbf{x})$$

$$\sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \frac{1}{\beta} \left(\phi(\mathbf{x})^\top (\lambda \mathbf{I} + \mathbf{\Phi}^\top \mathbf{\Phi})^{-1} \phi(\mathbf{x}) \right)$$

Predictive distribution (illustration of $p(t|x, \mathcal{D}, \beta)$)

- $t|x \sim \mathcal{N}(\sin(2\pi x), \beta^{-1})$
- **Model:** $\mu(x, \mathbf{w}) = \mathbf{w}^\top \phi(x)$, with Gaussian basis functions $\phi_j(x)$
- $\alpha = 10^{-3}$, $\beta = 2$



blue: true model, $\sin(2\pi x)$

orange: mean of predictive distribution, $\mu(x, \mathbf{w}_{\text{MAP}})$

shaded orange: $\mu(x, \mathbf{w}_{\text{MAP}}) \pm \sigma_N(x)$

Reading

“Pattern recognition and machine learning,”

Chapter 1 (1.2.4–1.2.6), Chapter 3 (Intro, 3.1, 3.2, 3.3 (until 3.3.2))