# MAP 534
# Introduction to machine learning
## Bayesian machine learning

Alain Durmus

Motivations

Bayesian statistics

- Decide what the input-output pairs are.
- Decide how to encode inputs and outputs. This defines the input space X, and the output space Y and the dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{N}$.
- Choose a class of hypotheses/representations $\mathcal{F} = \{f_w : X \rightarrow Y \ : \ w \in \mathbb{R}^d\}$.
- Choose a loss function $\ell$.
- Define the error function

$$E_\ell(w) = N^{-1} \sum_{i=1}^{N} \ell(y_i, f_w(x_i)) \ . \tag{1}$$

- Choose an algorithm to solve

$$\text{minimize } E_\ell \ . \tag{2}$$

- How to do so: vanish the gradient or gradient descent (see later)...

- Decide what the input-output pairs are.
- Decide how to encode inputs and outputs. This defines the input space X, and the output space Y and $\mathcal{D}$.
- Choose a class of hypotheses/representations $\mathcal{F} = \{f_w : X \to Y \ : \ w \in \mathbb{R}^d\}$.
- Choose a loss function $\ell$.
- Define the error function

$$E_\ell(w) = N^{-1} \sum_{i=1}^n \ell(y_i, f_w(x_i)) \ . \tag{3}$$

- Choose an algorithm to solve

$$\hat{w} \in \operatorname{argmin} E_\ell \ . \tag{4}$$

- Prediction:

$$\hat{y}_{\mathsf{pred}} = f_{\hat{w}}(x_{\mathsf{new}}) \ . \tag{5}$$

- This framework is called the deterministic discriminative setting.

- Consider $N$ observations $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ such that

$$x_i \in X = [0, 1] \text{ and } y_i \in Y = \mathbb{R} . \tag{6}$$

- We consider here that

$$\mathcal{P}_d = \left\{ f_w(x) = \sum_{i=1}^d w_i x^{(i)} \right\} . \tag{7}$$

- This corresponds to the choice of basis function $\phi_j(x) = x^j$.
- Justification: polynomials can approximate continuous any function on $[0, 1]$.
- LSE:

$$\hat{w} \in \underset{w}{\operatorname{argmin}} \, E(w) , \quad E(w) = \frac{1}{2N} \sum_{j=1}^N \left\{ y_j - \sum_{i=1}^d w_i x_j^i \right\}^2 . \tag{8}$$

- Two main questions still remain:
  - Can we weight the possible choices for $\mathcal{F}/d$?
  - Can we quantify the uncertainty of the prediction?
- The two questions are related and addressed with the use of Bayesian statistics:
  - For these two problems, we give some weights/probabilities on models/coefficients based on a priori knowledge.
  - Regarding the second point, Bayesian inference "sees" the parameter $w$ as random!

- No modeling view discussed previously: no probability!
- Here we consider a statistical model on the observations $\{(x_i, y_i)\}_{i=1}^N$.
- This model as in your statistics course is specified by a likelihood, i.e., a family of parametrized probability density functions (p.d.f.)

$$\{(x, y) \mapsto L_w(x, y) \ : \ w \in \Theta \subset \mathbb{R}^d\} \ .$$

- Examples:
  - Regression:
  $$Y_i = f_w(X_i) + \epsilon_i \ , \quad \epsilon_i \overset{\text{iid}}{\sim} N(0, 1) \ . \tag{9}$$
  Likelihood:
  $$L_w(x, y) = ? \tag{10}$$

- No modeling view discussed previously: no probability!
- Here we consider a statistical model on the observations $\{(x_i, y_i\}_{i=1}^{N}$.
- This model as in your statistics course is specified by a likelihood, i.e., a family of parametrized probability density functions (p.d.f.)

$$\{(x, y) \mapsto \mathrm{L}_w(x, y) \ : \ w \in \Theta \subset \mathbb{R}^d\} .$$

- Examples:
  - Classification ($Y = \{0, 1\}$):

$$Y_i = \mathbb{1}\{f_w(X_i) + \epsilon_i \geq 0\} \ , \quad \epsilon_i \overset{\text{iid}}{\sim} \mathsf{N}(0, 1) . \tag{11}$$

    Likelihood:

$$\mathrm{L}_w(x, y) = ? \tag{12}$$

- No modeling view discussed previously: no probability!
- Here we consider a statistical model on the observations $\{(x_i, y_i)\}_{i=1}^{N}$.
- Here, we fix a distribution $p$ for $x$ but does not matter, so the likelihood has the form:

$$\{(x, y) \mapsto L_w(x, y) = p_w(y|x)p(x) : w \in \Theta \subset \mathbb{R}^d\} .$$

- MLE:
$$\hat{w} \in \underset{w}{\operatorname{argmax}}\{\log p_w(y|x)\} \tag{13}$$

- Prediction:
$$y_{\text{pred}} = \underset{y}{\operatorname{argmax}} \, p_{\hat{w}}(y|x_{\text{new}}) . \tag{14}$$

- In our examples, what would $y_{\text{pred}}$ be?
- We only care about the conditional $y|x$!
- This framework is referred to as the probabilistic discriminative framework.

- Here we consider a statistical model on the observations $\{(x_i, y_i)\}_{i=1}^{N}$.
- We can also take a generative model approach:

$$\{(x, y) \mapsto L_w(x, y) = p_w(x|y)p(y) \; : \; w \in \Theta \subset \mathbb{R}^d\} \;.$$

- Still the MLE:

$$\hat{w} \in \operatorname*{argmax}_{w}\{\log p_w(y|x)\} \tag{15}$$

- Prediction:

$$y_{\mathsf{pred}} = \operatorname*{argmax}_{y} p_{\hat{w}}(y|x_{\mathsf{new}}) \;. \tag{16}$$

- What is $p_w(y|x)$? $p(y)$?
- Answer: Bayes theorem/formula and $p(y)$ is a prior to choose (details further...)
- This framework is referred to as the probabilistic generative framework.
- Pros: access to $p_{\hat{w}}$ which allows detection of outliers.
- Cons: computational demanding...
- Example in the next course!

- The statistical model associated with the LSE is:

$$Y_i \stackrel{\text{iid}}{\sim} \sum_{i=1}^{d} w_j \phi_j(X_i) + \sigma^2 Z_i , \quad i \in \{1, \ldots, N\} , \tag{17}$$

  where
  - $Z_i \stackrel{\text{iid}}{\sim} N(0, 1)$;
  - $X_i$ are i.i.d random variables with unknown distribution;
  - $w$ is the parameter to infer.

- Then, if we are just interested in inferring $w$, the log-likelihood is

$$\ell(w) = \frac{1}{2\sigma} \sum_{i=1}^{N} \left\{ y_i - \sum_{i=1}^{d} w_j \phi_j(x_i) \right\}^2 , \tag{18}$$

  where the observations are $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{N}$.

- Therefore, maximizing the log-likelihood leads to the same solution as minimizing the error function.

- MLE: Only point estimate! No uncertainty quantification!
- Confident intervals on the coefficients of $w$ defined by bootstrap:

$$\mathcal{D}_{\text{rand},i} \subset \mathcal{D}_{\text{train}} \text{ uniformly random }, \quad \text{and consider } \hat{w}(d, \mathcal{D}_{\text{rand},i}) . \quad (19)$$

- Consider the intervals which contains ..% of the solutions.
- Analysis of variance (ANOVA, using an F-test): test the null hypothesis that a model $\mathcal{M}_1$ is sufficient to explain the data against the alternative hypothesis that a more complex model $\mathcal{M}_2$.
- Do not generalize well and have pathologies[1]:
  - counter-intuitive behavior on confident intervals for some simple models;
  - p-values tend to overstate the evidence against the null no matter how large the sample size;
  - p-values is sensible to slight changes in the statistical models;
  - many frequentist methods regarding uncertainties/model selection does not follow likelihood principle: are based on hypothetical future observations.
- All these problems can be addressed by Bayesian inference!

---

[1] [Mur13, Section 6.6]

- Consider $N$ observations $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{N}$ such that

$$x_i \in X = [0, 1] \text{ and } y_i \in Y = \mathbb{R} . \tag{20}$$

- We consider here that

$$\mathcal{P}_d = \left\{ f_w(x) = \sum_{i=1}^{d} w_i x^i \right\} . \tag{21}$$

- We would like to quantify uncertainties with respect to
  - the parameters $w$;
  - our prediction;
  - the hypothesis class complexity $d$.

- This generalizes to any parametrized hypothesis class

$$\mathcal{F} = \{ f_w : X \to Y \ : \ w \in \Theta \} . \tag{22}$$

- Do we think that all hypothesis/models are equally probable... before we see any data?
- Here an hypothesis is a fixed function $f_w$ for some fixed parameter $w$.
- What does the probability of a model/hypothesis even mean?
- Do we need to choose a single "best" model $\mathcal{F}$ or can we consider several $\mathcal{F}_1, \mathcal{F}_2$ for our predictions?
- We need a framework to answer such questions.

- Bayes rule tells us how to do inference about hypothesis (the uncertain quantities) from data (measured quantities).

- Learning and prediction can be seen as a form of inference

$$p(\text{hypothesis}|\text{data}) = \frac{p(\text{data}|\text{hypothesis})p(\text{hypothesis})}{p(\text{data})} \ . \qquad (23)$$

- $p(\text{data}|\text{hypothesis})$ is the likelihood associated with the family of hypothesis we first consider.

- $p(\text{hypothesis}|\text{data})$ is called the posterior distribution of the hypothesis.

- However, in contrast to frequentist statistics, we choose a prior on our hypothesis!

- Bayes inference recipe:
    - Consider a staistical model for $\mathcal{D}$ parametrized by $w \in \Theta$:
    $$p(\mathcal{D}|w) = \mathrm{L}_w(x, y) \, . \tag{24}$$
    - We treat the likelihood as the conditional distribution of the data given the parameter!
    - Choose a prior for $w$, $p(w)$.
    - Consider the posterior:
    $$p(w|\mathcal{D}) \propto p(\mathcal{D}|w)p(w) = \mathrm{L}_w(x, y)p(w) \, . \tag{25}$$
    - All the conclusions are then drawn from the posterior.

- Consider the observations $\{y_i\}_{i=1}^N$ be i.i.d from

$$Y_i = \text{Ber}(q) , \; q \in [0, 1] \text{ is the parameter to infer} . \tag{26}$$

- The likelihood is

$$p_q(y) = ? \tag{27}$$

- We chose as a prior $\text{Beta}(\alpha, \beta)$, for $\alpha, \beta > 0$:

$$p(q|y) \propto ? \tag{28}$$

- The posterior distribution for $q$ is ?...

- Notation here if $p(q|y)$ is a conditional density:

$$p(q|y) \propto h(q, y) , \; \text{if } p(q|y) = h(q, y) \Big/ \int h(q, y) \mathrm{d}q . \tag{29}$$

- Bayes inference recipe:
  - Consider a staistical model for $\mathcal{D}$ parametrized by $w \in \Theta$:
  $$p(\mathcal{D}|w) = \mathrm{L}_w(x, y) . \tag{30}$$
  - We treat the likelihood as the conditional distribution of the data given the parameter!
  - Choose a prior for $w$, $p(w)$.
  - Consider the posterior:
  $$p(w|\mathcal{D}) \propto p(\mathcal{D}|w)p(w) = \mathrm{L}_w(x, y)p(w) . \tag{31}$$
  - All the conclusions are then drawn from the posterior.
  - The posterior is known up to a multiplicative constant:
  $$\mathrm{Z}(\mathcal{D}) = \int p(\mathcal{D}|w)p(w)\mathrm{d}w = \int \mathrm{L}_w(x, y)p(w)\mathrm{d}w . \tag{32}$$

  This constant is also known as the marginal likelihood.
  - In many models, this constant can not be computed and the posterior does not belong to "common" distribution.

- A robot, in order to behave intelligently, should be able to represent beliefs about propositions in the world:
  - charging station is at location (x,y,z)
  - that cat is hostile...
- Using probabilistic models, we want to represent the strengths of these beliefs, and be able to manipulate these beliefs based on a priori.
- The prior distribution models this prior knowledge.
- Data are then used to update our knowledge and give the posterior.
- Probabilistic learning can also be used for calibrated models and prediction uncertainty - getting systems that know what they do not know.

- Choosing a prior and following the Bayesian paradigm, we do not believe all models are equally probable to explain the data.

- We may believe that a simpler model is more probable than a complex one based on Occam's razzor (Aristotle, Ockham, Newton, Russel...)
  *We consider it a good principle to explain the phenomena by the simplest hypothesis possible.*

  - Ptolemy (c. AD 90 - c. 168) -

- Bayesian allows us to consider/combine a collection of hypothesis/models:

  - We do not know what particular function generated the data.
  - More than one of our models can perfectly fit the data.
  - We believe more than one of our models could have generated the data.
  - We want to reason in terms of a set of possible explanations, not just one.

- The first Bayesian estimator, the maximum a posterior estimator (MAP):

$$\hat{w}_{\mathsf{MAP}} \in \operatorname{argmax} p(w|\mathcal{D}) \ . \tag{33}$$

- The MAP is not fully Bayesian (not an admissible estimator)...
- The usual Bayesian estimator is the posterior mean:

$$\hat{w}_{\mathsf{post}} = \int w p(w|\mathcal{D}) \mathrm{d}w \ . \tag{34}$$

- To quantify the uncertainties over $w$ we consider $1 - \alpha$-credible region for $\alpha \in (0, 1)$.
- $C_\alpha$ is set to be a $1 - \alpha$-credible region if

$$\int \mathbb{1} \left\{ w \in C_\alpha \right\} p(w|\mathcal{D}) \mathrm{d}w \geq 1 - \alpha \ . \tag{35}$$

- Consider the observations $\{y_i\}_{i=1}^N$ be i.i.d from

$$Y_i = \text{Ber}(q) \ , \ q \in [0, 1] \text{ is the parameter to infer .} \tag{36}$$

- We chose as a prior $\text{Beta}(\alpha, \beta)$, for $\alpha, \beta > 0$:

$$p(q|y) \propto q^{\alpha - 1 + \sum_{i=1}^N y_i} (1 - q)^{\beta - 1 + N - \sum_{i=1}^N y_i} \ . \tag{37}$$

- The posterior distribution for $q$ is $\text{Beta}(\alpha + \sum_{i=1}^N y_i, \beta + N - \sum_{i=1}^N y_i)$.
- MAP:

$$\hat{w}_{\text{MAP}} = \frac{\alpha - 1 + \sum_i y_i}{\alpha + \beta - 2 + N} \ . \tag{38}$$

- Posterior mean:

$$\hat{w}_{\text{post}} = \frac{\alpha + \sum_i y_i}{\alpha + \beta + N} \ . \tag{39}$$

- We consider the statistical model associated with the LSE is:

$$Y_i \overset{\text{iid}}{\sim} \sum_{i=1}^{d} w_j \phi_j(X_i) + \sigma^2 Z_i , \quad i \in \{1, \dots, N\} , \tag{40}$$

where
- $Z_i \overset{\text{iid}}{\sim} N(0, 1)$;
- $X_i$ are i.i.d random variables with unknown distribution;
- $w$ is the parameter to infer.

- What does it mean to choose a prior on the hypothesis here?

- A hypothesis $f_w$ is a choice of a model structure $\mathcal{F}$ (first block) and a parameter value (second block) $w$.

- Consider the linear regression example:

$$f_w(x) = \sum_{i=1}^{d} w_j \phi_j(x) \,, \tag{41}$$

- The number $d$ and the choices of basis functions $\{\phi_j\}$ constitute the model structure;

- The coefficient $w$, the parameter value.

- Setting a prior $p(w)$ determines what functions this model can generate.

- For the moment $\mathcal{F}$ is fixed but we can also set a prior on the model structure (see after)!

- What is the posterior in this case? ✎

- The likelihood setting $\beta = \sigma^{-2}$ the precision:

$$L(\mathcal{D}|w) = N_n(\text{vector}(f_w(x_i)), \beta^{-1}I_n) \tag{42}$$

$$= (\beta/2\pi)^{N/2} \exp\left(-\frac{\beta}{2}\sum_{i=1}^{N}(y_i - \phi(x_i)^{\text{T}}w)^2\right) \tag{43}$$

$$= (\beta/2\pi)^{N/2} \exp\left(-\frac{\beta}{2}\|y - \Phi_x w\|^2\right), \tag{44}$$

with $\Phi_x =$??.

- If we choose $p(w) = N_d(m_0, S_0)$, we get

$$p(w|\mathcal{D}) = N(m_N, S_n), \tag{45}$$

**Proof of** (49).

$\square$

- The likelihood setting $\beta = \sigma^{-2}$ the precision:

$$L(\mathcal{D}|w) = N_n(\text{vector}(f_w(x_i)), \beta^{-1}I_n) \tag{46}$$

$$= (\beta/2\pi)^{N/2} \exp\left(-\frac{\beta}{2}\sum_{i=1}^{N}(y_i - \phi(x_i)^T w)^2\right) \tag{47}$$

$$= (\beta/2\pi)^{N/2} \exp\left(-\frac{\beta}{2}\|y - \Phi_x w\|^2\right), \tag{48}$$

with $\Phi_x =$??.

- If we choose $p(w) = N_d(m_0, S_0)$, we get

$$p(w|\mathcal{D}) = N(m_N, S_n),$$
$$m_N = S_N(S_0^{-1}m_0 + \beta\Phi_x^T y), \quad S_N = (S_0^{-1} + \beta\Phi_x^T\Phi_x)^{-1}. \tag{49}$$

**Proof of** (49).

$\square$

- Consider the linear regression example:

$$f_w(x) = \sum_{i=1}^{d} w_j \phi_j(X_i) . \tag{50}$$

- Based on the posterior $p(w|\mathcal{D})$, $\mathcal{D} = \{(y_i, x_i)\}_{i=1}^{N}$, how to make our predictions?

- First frequentist-like option:

$$y_{\text{pred}} = f_{\hat{w}}(x_{\text{new}}) , \tag{51}$$

where $\hat{w}$ is either the MAP or the posterior mean.

- Not really Bayesian...

- Indeed, Bayesian inference is also guided by the aim to give an "optimal" prediction.

- To define what we mean by an "optimal" prediction, we rely on decision theory.

- Given a dataset $\mathcal{D} = \{(y_i, x_i)\}_{i=1}^{N}$ with a probabilistic model

$$\{(\tilde{x}, \tilde{y}) \mapsto L_w(\tilde{x}, \tilde{y}) \ : \ w \in \Theta\} \ ,$$

  we would like to find the best estimator for the prediction $y_{\text{pred}}$ based on $x_{\text{new}}$.

- By estimator, here, we mean a function $\mathcal{D} \mapsto \hat{y}^{\mathcal{D}}$ which outputs a function:

$$y_{\text{pred}} = \hat{y}^{\mathcal{D}}(x_{\text{new}}) \ . \tag{52}$$

- How to compare estimator?
- We need a loss function $\ell : Y \times Y \to \mathbb{R}_+$ and a prior on $w$, $w \mapsto p(w)$.
- We define then the conditional risk (given $\mathcal{D}$ and $w$) as

$$\mathrm{cR}(\hat{y}^{\mathcal{D}}, w) = \mathbb{E}_{(Y_{\textbf{new}}, X_{\textbf{new}}) \sim L_w}[\ell(Y_{new}, \hat{y}^{\mathcal{D}}(X_{\text{new}}))] \tag{53}$$

$$= \int \ell(y_{\text{new}}, \hat{y}^{\mathcal{D}}(x_{\text{new}}))L_w(x_{\text{new}}, y_{\text{new}})\mathrm{d}(x_{\text{new}}, y_{\text{new}}) \ . \tag{54}$$

- Given a dataset $\mathcal{D} = \{(y_i, x_i)\}_{i=1}^N$ with a probabilistic model

$$\{(\tilde{x}, \tilde{y}) \mapsto L_w(\tilde{x}, \tilde{y}) \ : \ w \in \Theta\} \,,$$

  we would like to find the best estimator for the prediction $y_{\text{pred}}$ based on $x_{\text{new}}$.

- By estimator, here, we mean a function $\mathcal{D} \mapsto \hat{y}^{\mathcal{D}}$ which outputs a function:

$$y_{\text{pred}} = \hat{y}^{\mathcal{D}}(x_{\text{new}}) \,. \tag{55}$$

- We define then the conditional risk (given $\mathcal{D}$ and $w$) as

$$\mathrm{cR}(\hat{y}^{\mathcal{D}}, w) = \mathbb{E}_{(Y_{\text{new}}, X_{\text{new}}) \sim L_w}[\ell(Y_{\text{new}}, \hat{y}^{\mathcal{D}}(X_{\text{new}}))] \,. \tag{56}$$

- An ideal estimator is the one which minimizes the integrated/Bayesian risk:

$$\mathrm{IR} = \mathbb{E}_{\mathcal{D}, w}[\mathrm{R}(\hat{y}^{\mathcal{D}}, w)] = \int \mathrm{R}(\hat{y}^{\mathcal{D}}, w) L_w(\mathcal{D}) p(w) \mathrm{d}\mathcal{D} \mathrm{d}w \,. \tag{57}$$

- Here $L_w(\mathcal{D})$ is the complete likelihood $L_w(x, y) = \prod_{i=1}^N L_w(x_i, y_i)$.

- Given a dataset $\mathcal{D} = \{(y_i, x_i)\}_{i=1}^N$ with a probabilistic model $\{(\tilde{x}, \tilde{y}) \mapsto \mathrm{L}_w(\tilde{x}, \tilde{y}) : w \in \Theta\}$, we would like to find the best estimator for the prediction $y_{\mathrm{pred}}$ based on $x_{\mathrm{new}}$.

- By estimator, here, we mean a function $\mathcal{D} \mapsto \hat{y}^{\mathcal{D}}$ which outputs a function:

$$y_{\mathrm{pred}} = \hat{y}^{\mathcal{D}}(x_{\mathrm{new}}) . \tag{58}$$

- In the case $\ell(y_1, y_2) = (y_1 - y_2)^2 / 2$, we can show that the best estimator is

$$\hat{y}^{\mathcal{D}}_{\star, \mathrm{L}^2} = \int \tilde{y}_{\mathrm{new}} \mathrm{L}_w(\tilde{y}_{\mathrm{new}} | x_{\mathrm{new}}) p(w | \mathcal{D}) \mathrm{d}(\tilde{y}_{\mathrm{new}}, w) , \tag{59}$$

where $p(w|\mathcal{D})$ is the posterior distribution associated with prior $p$.

- It is called the Bayes estimator.

- Given a dataset $\mathcal{D} = \{(y_i, x_i)\}_{i=1}^{N}$ with a probabilistic model $\{(\tilde{x}, \tilde{y}) \mapsto L_w(\tilde{x}, \tilde{y}) : w \in \Theta\}$, we would like to find the best estimator for the prediction $y_{\mathsf{pred}}$ based on $x_{\mathsf{new}}$.

- By estimator, here, we mean a function $\mathcal{D} \mapsto \hat{y}^{\mathcal{D}}$ which outputs a function:

$$y_{\mathsf{pred}} = \hat{y}^{\mathcal{D}}(x_{\mathsf{new}}) . \tag{60}$$

- In the case $\ell(y_1, y_2) = (y_1 - y_2)^2/2$, we can show that the best estimator is

$$\hat{y}_{\star,\mathsf{L}^2}^{\mathcal{D}} = \int \tilde{y}_{\mathsf{new}} L_w(\tilde{y}_{\mathsf{new}}|x_{\mathsf{new}}) p(w|\mathcal{D}) \mathrm{d}(\tilde{y}_{\mathsf{new}}, w) , \tag{61}$$

  where $p(w|\mathcal{D})$ is the posterior distribution associated with prior $p$.

- This gives rise to the posterior predictive distribution:

$$p_{\mathsf{post}}(\tilde{y}_{\mathsf{new}}|\mathcal{D}) = \int L_w(\tilde{y}_{\mathsf{new}}|x_{\mathsf{new}}) p(w|\mathcal{D}) \mathrm{d}(w) . \tag{62}$$

- With this notation:

$$\hat{y}_{\star,\mathsf{L}^2}^{\mathcal{D}} = \int \tilde{y}_{\mathsf{new}} p_{\mathsf{post}}(\tilde{y}_{\mathsf{new}}|\mathcal{D}, x_{\mathsf{new}}) \mathrm{d}\tilde{y}_{\mathsf{new}} . \tag{63}$$

- This distribution give a point estimate for our prediciton but also completely characterizes the uncertainties about our predictions!

- Given a dataset $\mathcal{D} = \{(y_i, x_i)\}_{i=1}^N$ with a probabilistic model $\{(\tilde{x}, \tilde{y}) \mapsto L_w(\tilde{x}, \tilde{y}) : w \in \Theta\}$, we would like to find the best estimator for the prediction $y_{\mathrm{pred}}$ based on $x_{\mathrm{new}}$.

- By estimator, here, we mean a function $\mathcal{D} \mapsto \hat{y}^{\mathcal{D}}$ which outputs a function:

$$y_{\mathrm{pred}} = \hat{y}^{\mathcal{D}}(x_{\mathrm{new}}) . \tag{64}$$

- In the case $\ell(y_1, y_2) = (y_1 - y_2)^2/2$, we can show that the best estimator is

$$\hat{y}^{\mathcal{D}}_{\star,L^2}(x_{\mathrm{new}}) = \int \tilde{y}_{\mathrm{new}} L_w(\tilde{y}_{\mathrm{new}}|x_{\mathrm{new}}) p(w|\mathcal{D}) \mathrm{d}(\tilde{y}_{\mathrm{new}}, w) , \tag{65}$$

where $p(w|\mathcal{D})$ is the posterior distribution associated with prior $p$.

**Proof of** (65).

$\Box$

- The likelihood setting $\beta = \sigma^{-2}$ the precision:

$$L(\mathcal{D}|w) = N_n(\text{vector}(f_w(x_i)), \beta^{-1}I_n) \qquad (66)$$

$$= (\beta/2\pi)^{N/2} \exp\left(-\frac{\beta}{2N}\sum_{i=1}^{N}(y_i - f_w(x_i))^2\right) . \qquad (67)$$

- If we choose $p(w) = N_d(m_0, S_0)$, we get

$$p(w|\mathcal{D}) = N(m_N, S_n) , \quad m_N = S_N(S_0^{-1}m_0 + \beta\Phi_x^{\mathrm{T}}y) , \quad S_n = (S_0^{-1} + \beta\Phi_x^{\mathrm{T}}\Phi_x)^{-1} . \qquad (68)$$

- Since $p(y_{\text{new}}|w, x_{\text{new}}) = N(f_w(x_{\text{new}}), \beta^{-1})$, we get that the predictive posterior is

$$p(y_{\text{new}}|x_{\text{new}}) = N(\phi(x_{\text{new}})^{\mathrm{T}}m_N, \beta^{-1} + \phi(x_{\text{new}})S_N\phi(x_{\text{new}})) , \qquad (69)$$

- the Bayes estimator:

$$\hat{y}_{\star,\mathrm{L}^2} = \phi(x_{\text{new}})^{\mathrm{T}}m_N . \qquad (70)$$

- Proof in practical sessions.

- What if we are unsure which model is right? So far we assumed we were able to start by making a definite choice of model.

- We can compare models based on marginal likelihoods (also known as model evidence) for each model - this is the probability the model assigns to the observed data.

- This is the normalizing constant in Bayes rule which we ignored previously.

- Let us say that we have two models $\mathcal{F}_1, \mathcal{F}_2$.

- Question: given some data, can we say if one of them is most probable?

- Examples:

$$\mathcal{F}_1 = \left\{ f_w = \sum_{i=1}^{d_1} w_j \phi_j(X_i) \ : \ w \in \mathbb{R}^{d_1} \right\} \ , \quad \mathcal{F}_2 = \left\{ f_w = \sum_{i=1}^{d_2} w_j \phi_j(X_i) \ : \ w \in \mathbb{R}^{d_2} \right\}$$

$$(71)$$

$d_1$ or $d_2$ should be privileged?

- Solution: Bayesian paradigm.

- We treat the prior $p(w|\mathcal{F}_1), p(w|\mathcal{F}_2)$ used for $\mathcal{F}_1, \mathcal{F}_2$ as likelihood/conditional probability.

- We set some prior on the models $\mathcal{F}_i$, $i = 1, 2$.

- Let us say that we have two models $\mathcal{F}_1, \mathcal{F}_2$.

- Question: given some data, can we say if one of them is most probable?

- Solution: Bayesian paradigm.

- We treat the likelihoods $p(\mathcal{D}|w, \mathcal{F}_1), p(\mathcal{D}|w, \mathcal{F}_1)$ and priors $p(w|\mathcal{F}_1), p(w|\mathcal{F}_2)$ used for $\mathcal{F}_1, \mathcal{F}_2$ as likelihood/conditional probability.

- We set some prior on the models $\mathcal{F}_i$, $i = 1, 2$.

- In most cases, the uniform prior is chosen...

- The posterior distribution for $(w, \mathcal{F}_i)$ is by Bayes theorem:

$$p(w, \mathcal{F}_i|\mathcal{D}) = p(\mathcal{F}_i)p(w|\mathcal{F}_i)p(\mathcal{D}|w, \mathcal{F}_i)/p(\mathcal{D}) , \qquad (72)$$

$$p(\mathcal{D}) = \sum_i \int_w p(\mathcal{F}_i)p(w|\mathcal{F}_i)p(\mathcal{D}|w, \mathcal{F}_i)\mathrm{d}w . \qquad (73)$$

- The posterior distribution for $\mathcal{F}_i$ is then:

$$p(\mathcal{F}_i|\mathcal{D}) = \frac{p(\mathcal{F}_i)}{p(\mathcal{D})} \int p(w|\mathcal{F}_i)p(\mathcal{D}|w, \mathcal{F}_i)\mathrm{d}w . \qquad (74)$$

- Let us say that we have two models $\mathcal{F}_1, \mathcal{F}_2$.

- Question: given some data, can we say if one of them is most probable?

- The posterior distribution for $(w, \mathcal{F}_i)$ is by Bayes theorem:

$$p(w, \mathcal{F}_i|\mathcal{D}) = p(\mathcal{F}_i)p(w|\mathcal{F}_i)p(\mathcal{D}|w, \mathcal{F}_i)/p(\mathcal{D}) , \qquad (75)$$

$$p(\mathcal{D}) = \sum_i \int_w p(\mathcal{F}_i)p(w|\mathcal{F}_i)p(\mathcal{D}|w, \mathcal{F}_i)\mathrm{d}w . \qquad (76)$$

- The posterior distribution for $\mathcal{F}_i$ is then:

$$p(\mathcal{F}_i|\mathcal{D}) = \frac{p(\mathcal{F}_i)Z_i(\mathcal{D})}{p(\mathcal{D})} , \quad Z_i(\mathcal{D}) = \int p(w|\mathcal{F}_i)p(\mathcal{D}|w, \mathcal{F}_i)\mathrm{d}w . \qquad (77)$$

- If $p(\mathcal{F}_i) = 1/2$, the posterior distribution for $\mathcal{F}_i$ simplifies:

$$p(\mathcal{F}_i|\mathcal{D}) = \frac{Z_i(\mathcal{D})}{Z_1(\mathcal{D}) + Z_2(\mathcal{D})} . \qquad (78)$$

- This easily generalizes to finite number of models.

- From looking at the equation of posterior distribution, the marginal likelihood is given by

$$Z(\mathcal{D}, \mathcal{M}) = \int p(\mathcal{D}|w, \mathcal{M}) p(w|\mathcal{M}) \mathrm{d}w = \int \mathrm{L}_w^{\mathcal{M}}(x, y) p(w|\mathcal{M}) \mathrm{d}w \ . \quad (79)$$

- Second level inference : model comparison

$$p(\mathcal{M}|\mathcal{D}) \propto Z(\mathcal{D}, \mathcal{M}) p(\mathcal{M}) \ . \quad (80)$$

- Represents some belief/probability on our models given $\mathcal{D}$.

- Model selection:

$$\mathcal{M}^\star = \underset{\mathcal{M}}{\mathrm{argmax}} \, p(\mathcal{M}|\mathcal{D}) \ . \quad (81)$$

- The likelihood setting $\beta = \sigma^{-2}$ the precision:

$$L(\mathcal{D}|w) = N_n(\text{vector}(f_w(x_i)), \beta^{-1}I_n) \tag{82}$$

$$= (\beta/2\pi)^{N/2} \exp\left(-\frac{\beta}{2N}\sum_{i=1}^{N}(y_i - f_w(x_i))^2\right). \tag{83}$$

- If we choose $p(w) = N_d(m_0, S_0)$, we get

$$p(w|\mathcal{D}) = N(m_N, S_n), \quad m_N = S_N(S_0^{-1}m_0 + \beta\Phi_x^Ty), \quad S_n = (S_0^{-1} + \beta\Phi_x^T\Phi_x)^{-1}. \tag{84}$$

- The marginal likelihood is in the case $p(w) = N_d(0, \alpha^{-1}I_d)$:

$$Z(\mathcal{D}) = \alpha^{d/2}\beta^{N/2}(2\pi)^{-N/2}[\det S_N]^{1/2}\exp\left(-\beta\|y\|^2/2 + \beta\left\langle S_N m_N, \Phi_x^Ty\right\rangle/2\right) \tag{85}$$

$$= \alpha^{d/2}\beta^{N/2}(2\pi)^{-N/2}[\det S_N]^{1/2}\exp\left(-\beta\|y - \Phi_x m_N\|^2/2 - \alpha\|m_N\|^2/2\right). \tag{86}$$

- Proof in small classes.

- Let us say that we have some hyperparameters $\beta$ and $\alpha$ for the likelihoods $p_w((x,y)|\beta)$ and the prior $p(w|\alpha)$ respectively.

- Question: given some data, can we make some recommandations on the choice of these hyperparameters?

- Examples: linear regression (again!) (recall $f_w(x_1) = \sum_{j=1}^{d} w_j \phi_j(x_1)$)

$$p_w((x_1, y_1)|\beta) = (2\pi\sigma^2)^{1/2} \exp(-(y_1 - f_w(x_i))^2/(2\sigma^2)) \ , \ \beta = \sigma^{-2} \ , \quad (87)$$
$$p(w|\alpha) = \alpha \|w\|^2 \ . \quad (88)$$

- Solution: Bayesian paradigm (again...).

- We set some prior on $\alpha$ and $\beta$ and treat them as parameter as it was the case for models.

- In most cases, the uniform prior is chosen uniform $p(\alpha) = 1$, $p(\beta) = 1$ (even if they do not define a well-defined distribution...).

- Let us say that we have some hyperparameters $\beta$ and $\alpha$ for the likelihoods $p_w((x, y)|\beta)$ and the prior $p(w|\alpha)$ respectively.
- Question: given some data, can we make some recommandations on the choice of these hyperparameters?
- The posterior distribution for $(w, \lambda, \beta)$ is by Bayes theorem:

$$p(w, \lambda, \beta|\mathcal{D}) = p(\alpha)p(\beta)p(w|\alpha, \beta)p(\mathcal{D}|w, \alpha, \beta)/p(\mathcal{D}) , \qquad (89)$$

$$p(\mathcal{D}) = \sum_i \int_w p(\alpha)p(\beta)p(w|\alpha, \beta)p(\mathcal{D}|w, \alpha, \beta)\mathrm{d}w \qquad (90)$$

- The posterior distribution for $\alpha, \beta$ is then:

$$p(\alpha, \beta|\mathcal{D}) = \frac{p(\alpha)p(\beta)}{p(\mathcal{D})} \int p(w|\alpha, \beta)p(\mathcal{D}|w, \alpha, \beta)\mathrm{d}w . \qquad (91)$$

- Pragamtic choice:

$$(\hat{\alpha}, \hat{\beta}) \in \operatorname{argmax} p(\alpha, \beta|\mathcal{D}) , \qquad (92)$$

this corresponds maximization of the marginal likelihood or empirical Bayes approach.
- Example: Bayesian linear regression ✎ .

- For most Bayesian inference problems, the integrals needed to do inference and prediction are not analytically tractable - hence the need for variaous approximations.

- Most of the exceptions involve conjugate priors, which combine nicely with the likelihood to give a posterior distribution of the same form.

- Basic Idea : Given likelihood function $L_w(x, y)$, choose a family of prior distributions such that integrals can be obtained tractably.

- If the prior $p(w)$ and posterior $p(w|\mathcal{D})$ belong to same family of distributions, the prior is called a conjugate prior.

- Example: if likelihood function is Gaussian, choosing Gaussian prior over mean will ensure that the posterior distribution is also Gaussian.

## Monte Carlo needs: Representing Prior and Posterior by Samples

- The complex distributions we will often use as priors, or obtain as posteriors, may not be easily represented.
- A general technique is to represent a distribution by sampling of many values drawn randomly from it. We can then
  - Visualize the distribution by viewing these sample values, or low dimensional projections of them (PCA..later).
  - Make Monte Carlo estimates for probabilities or expectations with respect to the distribution, by taking averages over these sample values.
- Obtaining a sample from the prior is easy! Obtaining a sample from the posterior is usually more difficult - nevertheless a dominant approach to Bayesian computation.

📄 Kevin P. Murphy. *Machine learning : a probabilistic perspective*. Cambridge, Mass.
[u.a.]: MIT Press, 2013. ISBN: 9780262018029 0262018020.