

Machine Learning

Bayesian Decision Theory, Maximum Likelihood,
and Regularized Empirical Risk Minimization

Prof. Matthias Hein

Machine Learning Group
Department of Mathematics and Computer Science
Saarland University, Saarbrücken, Germany

Lecture 5, 06.11.2013

Regression: $y_i = f(x_i) + \varepsilon_i$,

Bayesian: Model for distribution $p(y|X = x, f) \rightarrow$ model for distribution of ε .

Likelihood: $p(y|X = x, f)$ models: how *likely* is the output y given the point x and the function value $f(x)$?

Maximum likelihood estimation:

Use the function f which maximizes the likelihood.

$$f^*(x) := \arg \max_{f \in \mathcal{F}} \mathbb{E}_{Y|X=x} \left[p(Y|X = x, f) \right]$$

Correspondence to loss function:

$$L(y, f(x)) = -\log p(y|X = x, f(x)) + c,$$

⇒ Maximizing the likelihood $p(y|X = x, f(x))$ is **equivalent** to minimizing the loss $L(y, f(x))$.

Example:

$$p(y|X = x, f(x)) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - f(x))^2}{2\sigma^2}\right).$$

and so the corresponding loss function $L(y, f(x))$ is given as

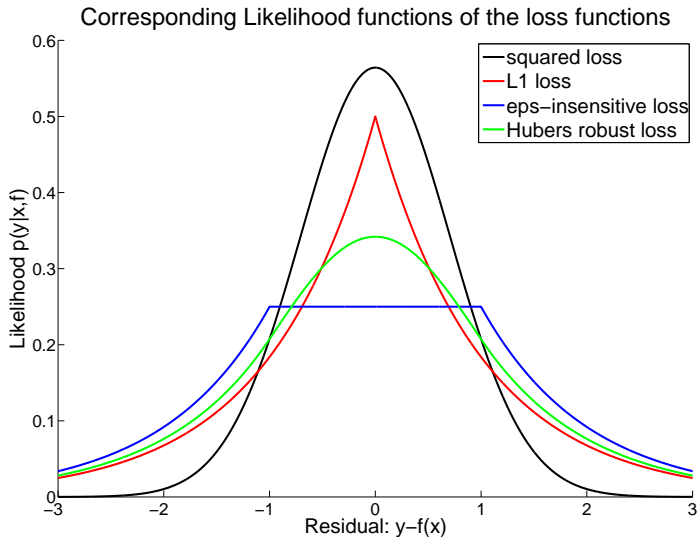
$$L(y, f(x)) = -\log p(y|X = x, f(x)) = \log(\sqrt{2\pi\sigma^2}) + \frac{(y - f(x))^2}{2\sigma^2}.$$

Bayesian interpretation of loss functions III

With: $\varepsilon = y - f(x)$,

	$L(\varepsilon)$	$p(y x, f(x))$
squared loss	$ \varepsilon ^2$	$\frac{1}{\sqrt{\pi}} e^{-(y-f(x))^2}$,
L_1 - loss	$ \varepsilon $	$\frac{1}{2} e^{- y-f(x) }$,
σ -insensitive	$(\varepsilon - \sigma) \mathbb{1}_{ \varepsilon > \sigma}$	$\frac{1}{2+2\sigma} e^{-(y-f(x) - \varepsilon) \mathbb{1}_{ y-f(x) > \sigma}}$,
Huber's robust loss	$\begin{cases} \frac{1}{2\sigma} \varepsilon ^2 & \text{if } \varepsilon \leq \sigma \\ \varepsilon - \frac{\sigma}{2} & \text{if } \varepsilon > \sigma \end{cases}$	$\begin{cases} e^{-\frac{(y-f(x))^2}{2\sigma}} & \text{if } y-f(x) \leq \sigma \\ e^{-(y-f(x) - \frac{\sigma}{2})} & \text{if } y-f(x) > \sigma \end{cases}$

Bayesian interpretation of loss functions IV



- The optimal classifier for classification is the **Bayes classifier**. Extensions to **cost-sensitive learning** and the **multi-class** setting possible.
- Two schemes for solving multi-class problems: **one-versus-all** and **one-versus-one**.
- Discussion of the optimal function in regression (loss-dependent).
- **Bayesian** interpretation of loss functions \Rightarrow **Maximizing the likelihood** is equivalent to **minimizing the corresponding loss**.

Problem: In order to compute the Bayes optimal learning rule we need to know the joint measure P on $\mathcal{X} \times \mathcal{Y}$,

but !

We do not know P but we have only the **training data** $(X_i, Y_i)_{i=1}^n$.

Idea: approximate the risk functional using the training data.

Assumption: Training data $(X_i, Y_i)_{i=1}^n$ is an **i.i.d.** sample of the probability measure P on $\mathcal{X} \times \mathcal{Y}$.

i.i.d. = independently and identically distributed

- $(X_i, Y_i)_{i=1}^n$ are random variables,
- **independent:** joint density factorizes

$$p((x_1, y_1); (x_2, y_2); \dots; (x_n, y_n)) = \prod_{i=1}^n p_i(x_i, y_i).$$

- **identically distributed:**

$$p_i(x, y) = p_j(x, y), \quad \forall i, j \in \{1, \dots, n\}.$$

and $p(x, y)$ is the density of the data-generating measure P on $\mathcal{X} \times \mathcal{Y}$.

Statistics: Given an i.i.d. sample $(X_i)_{i=1}^n$, use the empirical measure

$$P_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x=X_i}$$

to approximate quantities of the data generating measure.

- **empirical mean:** $\mathbb{E}_{P_n}[X] = \frac{1}{n} \sum_{i=1}^n x \mathbb{1}_{x=X_i} = \frac{1}{n} \sum_{i=1}^n X_i$,
- **empirical variance:** $\text{Var}[X] = \frac{1}{n} \sum_{i=1}^n \left(X_i - \frac{1}{n} \sum_{j=1}^n X_j \right)^2$,
- **empirical covariance:**
 $\text{Cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n X_i Y_i - \frac{1}{n} \sum_{i=1}^n X_i \frac{1}{n} \sum_{i=1}^n Y_i$,

P_n approximates P

Definition

Let $(X_i, Y_i)_{i=1}^n$ be an i.i.d. sample of P on $\mathcal{X} \times \mathcal{Y}$, which we call the **training sample**. The **empirical loss** is defined as

$$\mathbb{E}_{P_n}[L(Y, f(X))] = \frac{1}{n} \sum_{i=1}^n L(Y_i, f(X_i)).$$

Given a class of functions \mathcal{F} , **empirical risk minimization** is defined as

$$f_n = \arg \min_{f \in \mathcal{F}} \mathbb{E}_{P_n}[L(Y, f(X))] = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n L(Y_i, f(X_i)),$$

where f_n is then learning rule based on the training sample.

- if the function class is too large it is likely that we overfit the training data,
- the mapping “data” to “learning rule” can be seen as an *inverse problem*.

Definition of a **well-posed problem**:

- ▶ a solution exists,
- ▶ the solution is unique,
- ▶ the solution depends continuously on the data.

A problem which does not have one of these properties is called **ill-posed**. In particular the last two properties are most of the time not fulfilled in empirical risk minimization. In order to make problems well-posed one uses **regularization**.

Natural loss: 0-1-loss $L(y, f(x)) = \mathbb{1}_{y \neq f(x)}$.

Empirical risk minimization:

minimize the number of errors on the training set:

$$\frac{1}{n} \sum_{i=1}^n L(Y_i, f(X_i)) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{Y_i \neq f(X_i)}.$$

Problem:

- for several classes of functions, empirical risk minimization leads to NP-hard problems
 \implies use of convex margin-based loss functions.

Standard loss: squared loss $L(y, f(x)) = (y - f(x))^2$.

$$f_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i))^2$$

- $\mathcal{F} = \{f(x) = \langle w, x \rangle \mid w \in \mathbb{R}^d\}$, **linear least squares regression**.

Empirical risk minimization:

- function class \mathcal{F} too large \rightarrow overfitting,
- function class \mathcal{F} too small \rightarrow underfitting,

Idea: Use regularization together with a rather large function class \mathcal{F} .

Regularized empirical risk minimization

Definition

Let

- $(X_i, Y_i)_{i=1}^n$ be the training sample,
- \mathcal{F} a fixed function class,
- $L(y, f(x))$ the loss function,
- $\Omega : \mathcal{F} \rightarrow \mathbb{R}_+$ the **regularization functional**.

Then **regularized empirical risk minimization** is defined as

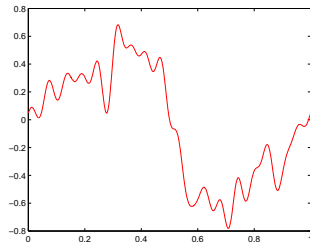
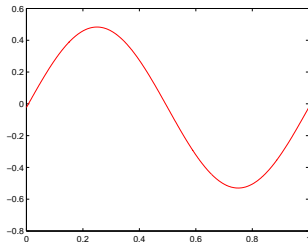
$$f_{n,\lambda} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n L(Y_i, f(X_i)) + \lambda \Omega(f),$$

where $\lambda \in \mathbb{R}_+$ is called the **regularization parameter**.

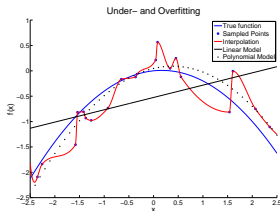
\implies This form of regularization is called **Tikhonov regularization**.

Trade-off between **fit of the data** and **complexity of the learning rule**.

Complexity of a function



Left: Relatively simple function, very smooth,
Right: Complex function, less smooth.



Equivalent formulation (Ivanov regularization):

Proposition

If the loss $L(y, f(x))$ and the regularization function $\Omega(f)$ are convex in f and the set $\{f \mid \Omega(f) < r\}$ is non-empty for every $r > 0$ and \mathcal{F} is a convex set, then regularized empirical risk minimization is equivalent to the following problem:

$$f_{n,r} = \arg \min_{f \in \mathcal{F}, \Omega(f) \leq r} \frac{1}{n} \sum_{i=1}^n L(Y_i, f(X_i)),$$

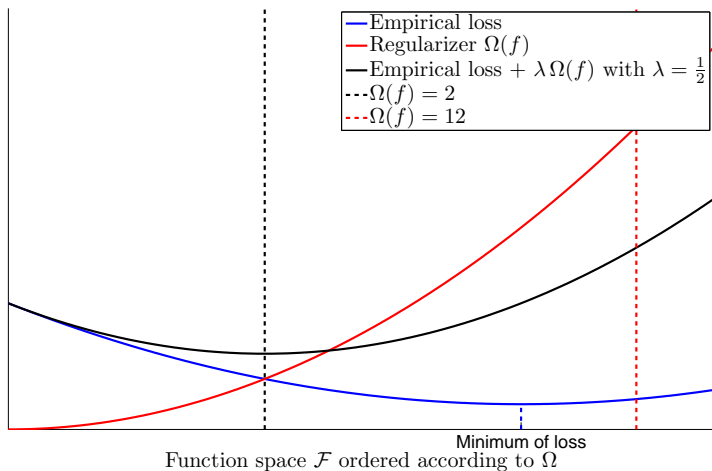
in the sense that there exists for each r a corresponding λ such that $f_{n,r} = f_{n,\lambda}$.

Proof:

- use of duality in convex optimization,
- in the lecture notes: proof for finite-dimensional function classes.

Tikhonov versus Ivanov regularization

Tikhonov versus Ivanov regularization



General Principle: prefer less complex function, as measured by Ω , if they have the same loss.

“Occam's razor”:

Pluralitas non est ponenda sine necessitas. (Plurality should not be posited without necessity.),

or similarly:

“Having two competing theories which make exactly the same predictions, the one that is simpler is the better.”

Regularized empirical risk minimization IV

Regularization parameter λ : controls trade-off between fit and complexity.

Limits: $\lambda \rightarrow 0$ and $\lambda \rightarrow \infty$.

$$\lambda \rightarrow 0 \quad \arg \min_{f \in \mathcal{F}} \Omega(f),$$
$$\arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n L(Y_i, f(X_i))$$

$$\lambda \rightarrow \infty \quad \arg \min_{\{f \in \mathcal{F} \mid \Omega(f)=0\}} \frac{1}{n} \sum_{i=1}^n L(Y_i, f(X_i))$$

Example: $f : \mathbb{R}^d \rightarrow \mathbb{R}$,

$$\Omega(f) = \int_{\mathbb{R}^d} \sum_{i=1}^d \left(\frac{\partial f}{\partial x^i} \right)^2 dx = \int_{\mathbb{R}^d} \|\nabla f\|^2 dx$$

$$\Omega(f) = 0 \quad \Longleftrightarrow \quad \exists c \in \mathbb{R}, \text{ such that } f(x) = c, \forall x \in \mathbb{R}^d.$$

Related principle: **Structural risk minimization** proposed by Vapnik.

- empirical risk minimization over nested function classes \mathcal{F}_n , such that $\mathcal{F}_1 \subset \mathcal{F}_2 \subset \dots$,
- as the size of the sample n increases one also allows more complex functions.

Example: start with the linear functions and then add polynomials of increasing order as n increases.

Relation I:

Empirical risk minimization
corresponds to
maximum likelihood estimation.

Relation II:

Regularized empirical risk minimization
corresponds to
maximum a posteriori estimation.

Maximum Likelihood Estimation

Problem: Given samples x_1, \dots, x_n identify the probability measure $p(x)$ which generated this sample.

General problem too difficult \implies parametric model of $p(x)$.

General Idea:

- parametric model of the data generating probability measure:

$$p(x | \theta) \quad (\text{the likelihood}).$$

- i.i.d. data $x_1, \dots, x_n \implies p(x_1, \dots, x_n | \theta) = \prod_{i=1}^n p(x_i | \theta)$.
- find parameter $\theta \in \Theta$ by maximizing the likelihood (resp. the log-likelihood)

$$\begin{aligned} \arg \max_{\theta \in \Theta} \prod_{i=1}^n p(x_i | \theta) &= \arg \max_{\theta \in \Theta} \log \left(\prod_{i=1}^n p(x_i | \theta) \right) \\ &= \arg \max_{\theta \in \Theta} \sum_{i=1}^n \log (p(x_i | \theta)) \end{aligned}$$

Maximum Likelihood Estimation - Example

Gaussian model:

$p(x|\mu) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$, the variance σ^2 is assumed to be known.

Maximum likelihood estimation of μ :

$$\begin{aligned}\arg \max_{\mu \in \mathbb{R}} \sum_{i=1}^n \log(p(x_i | \mu)) &= \arg \max_{\mu \in \mathbb{R}} \sum_{i=1}^n \left(-\frac{\log(2\pi\sigma^2)}{2} - \frac{(x_i - \mu)^2}{2\sigma^2} \right) \\ &= \arg \min_{\mu \in \mathbb{R}} \sum_{i=1}^n (x_i - \mu)^2\end{aligned}$$

The objective is convex in $\mu \implies$ Every local minimum is a global minimum.

The mean parameter μ^* maximizing the likelihood is:

$$\mu^* = \frac{1}{n} \sum_{i=1}^n x_i.$$

Maximum Likelihood:

- model for the conditional distribution = the **likelihood**: $p(y|x, f)$,
 f denotes the parameter of the model,
 \mathcal{F} is the set of parameters.
- i.i.d. sample of the data $D = (X_i, Y_i)_{i=1}^n$.

Definition

The **maximum likelihood** solution f_{ML} is then defined as

$$f_{ML} = \arg \max_{f \in \mathcal{F}} P(D|f) = \arg \max_{f \in \mathcal{F}} \prod_{i=1}^n P(Y_i|X_i, f),$$

Model for $p(y x, f(x))$	\longrightarrow	$L(y, f(x)) = -\log p(y x, f(x))$,
Loss function $L(y, f(x))$	$\xrightarrow[\text{correspondence}]{\text{generally no}}$	$p(y x, f(x)) = e^{-L(y, f(x))}$.

Proposition

Given an i.i.d. training sample $(X_i, Y_i)_{i=1}^n$, a class of functions \mathcal{F} and a likelihood $p(y|x, f)$, then the **maximum likelihood solution** f_{ML} agrees with the solution of **empirical risk minimization** f_n for the loss function $L(y, f(x)) = -\log p(y|x, f)$.

- output space \mathcal{Y} is discrete: likelihood is probability $P(y|x, f)$,
- output space \mathcal{Y} is continuous: likelihood is density $p(y|x, f)$.

Proof: By assumption we know $L(y, f(x)) = -\log P(y|x, f)$, then

$$\begin{aligned} f_{ML} &= \arg \max_{f \in \mathcal{F}} P(D|f) = \arg \max_{f \in \mathcal{F}} \prod_{i=1}^n P(Y_i|X_i, f) \\ &= \arg \max_{f \in \mathcal{F}} \log \left[\prod_{i=1}^n P(Y_i|X_i, f) \right] \\ &= \arg \max_{f \in \mathcal{F}} \sum_{i=1}^n \log P(Y_i|X_i, f) \\ &= \arg \min_{f \in \mathcal{F}} - \sum_{i=1}^n \log P(Y_i|X_i, f) \\ &= \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n L(Y_i, f(X_i)) = f_n, \end{aligned}$$

Maximum A Posteriori Estimation

Idea: integrate **prior belief** on the model parameter θ

Realization: θ is random, we have a **prior distribution** $p(\theta)$.

MAP Estimation:

- prior distribution $p(\theta)$
- using Bayes rule

$$p(\theta | x_1, \dots, x_n) = \frac{p(x_1, \dots, x_n | \theta)p(\theta)}{p(x_1, \dots, x_n)} = \frac{p(x_1, \dots, x_n | \theta)p(\theta)}{\int_{\Theta} p(x_1, \dots, x_n | \theta)p(\theta)d\theta}.$$

The denominator is called the **partition function**.

- find parameter θ by maximizing the a posteriori distribution $p(\theta | x_1, \dots, x_n)$

$$\begin{aligned}\arg \max_{\theta \in \Theta} \prod_{i=1}^n p(\theta | x_1, \dots, x_n) &= \arg \max_{\theta \in \Theta} \log \left(p(x_1, \dots, x_n | \theta)p(\theta) \right) \\ &= \arg \max_{\theta \in \Theta} \sum_{i=1}^n \log (p(x_i | \theta)) + \log (p(\theta)).\end{aligned}$$

Maximum A Posteriori Estimation - Example

Gaussian model:

$p(x|\mu) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$, the variance σ^2 is assumed to be known.

$p(\mu) = \frac{1}{\sqrt{2\pi\sigma_\mu^2}} e^{-\frac{(\mu-\mu_0)^2}{2\sigma_\mu^2}}$, the variance σ_μ^2 is assumed to be known.

MAP estimation of μ :

$$\begin{aligned}\arg \max_{\mu \in \mathbb{R}} p(\mu | x_1, \dots, x_n) &= \arg \max_{\mu \in \mathbb{R}} \sum_{i=1}^n \log(p(x_i | \mu)) + \log(p(\mu)) \\ &= \arg \min_{\mu \in \mathbb{R}} \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 + \frac{1}{2\sigma_\mu^2} (\mu - \mu_0)^2\end{aligned}$$

The objective is convex in μ The MAP estimate of the mean parameter μ^* is:

$$\mu^* = \frac{1}{1 + \frac{\sigma^2}{n\sigma_\mu^2}} \frac{1}{n} \sum_{i=1}^n x_i + \frac{1}{1 + \frac{n\sigma_\mu^2}{\sigma^2}} \mu_0.$$

Ingredients:

- likelihood function $P(y|x, f)$,
- set of parameters \mathcal{F} ,
- **prior** $P(f)$ over the function/parameter space \mathcal{F} .
probability measure over the class of functions
 \implies which functions are more “likely” than others,
 \implies encodes **a priori** knowledge.

Ingredients:

- likelihood function $P(y|x, f)$,
- set of parameters \mathcal{F} ,
- **prior** $P(f)$ over the function/parameter space \mathcal{F} .
probability measure over the class of functions
 \implies which functions are more “likely” than others,
 \implies encodes **a priori** knowledge.

Bayes theorem: Transform the prior and the likelihood into the **posterior probability** $P(f|D)$,

$$P(f|D) = \frac{P(D|f)P(f)}{P(D)},$$

where $P(D) = \int_{\mathcal{F}} P(D|f)P(f)df$.

Definition

The **maximum a posteriori** estimator for f is defined as

$$f_{MAP} = \arg \max_{f \in \mathcal{F}} P(f|D) = \arg \max_{f \in \mathcal{F}} \prod_{i=1}^n P(Y_i|X_i, f)P(f),$$

where we have discarded $P(D)$ since it is a constant.

Given a prior over functions $P(f)$ we define the following regularization functional $\Omega(f)$,

$$\Omega(f) = -\log P(f) \quad \implies \quad P(f) = e^{-\Omega(f)},$$

Proposition

The MAP estimator f_{MAP} agrees with the minimizer of $f_{\lambda=\frac{1}{n},n}$ of the regularized empirical risk minimization if

Loss function:	$L(y, f(x))$	$= -\log P(y x, f),$
regularization functional:	$\Omega(f)$	$= -\log P(f).$

Proof.

By assumption we know $L(y, f(x)) = -\log P(y|x, f)$ and $\Omega(f) = -\log P(f)$, then

$$\begin{aligned} f_{MAP} &= \arg \max_{f \in \mathcal{F}} P(f|D) = \arg \max_{f \in \mathcal{F}} \prod_{i=1}^n P(Y_i|X_i, f)P(f) \\ &= \arg \max_{f \in \mathcal{F}} \sum_{i=1}^n \log P(Y_i|X_i, f) + \log P(f) \\ &= \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n L(Y_i, f(X_i)) + \Omega(f) = f_{n, \lambda = \frac{1}{n}}. \end{aligned}$$

where we have used that the logarithm is a strictly increasing function. \square

Posterior predictive distribution:

$$p(x|D) = \int_{\Theta} p(x, \theta | D) d\theta = \int_{\Theta} p(x | \theta, D) p(\theta | D) d\theta$$

- $p(x|D)$ is **not** the true data-generating distribution !
- if the posterior $p(\theta | D)$ is very peaked, this is roughly the same as $p(x | \theta_{\text{MAP}})$

Learning setting:

$$p(y | x, D) = \int_{\Theta} p(y | x, \theta) p(\theta | D) d\theta$$

Bayesians consider the full distribution $p(\theta | D)$ against the point estimate in the MAP estimation.