

Machine Learning

MIRI Master

Lluís A. Belanche

`belanche@cs.upc.edu`



Soft Computing Research Group
Dept. de Ciències de la Computació (Computer Science)
Universitat Politècnica de Catalunya

Spring Semester 2016-2017

LLECTURE 6: Discriminative classifiers - Logistic regression and beyond

Discriminative classifiers

Outline

1. **Maximum Likelihood (ML) framework (I)**
 - Gentle motivation
 - Example 1: ML method for the Gaussian
 - Example 2: ML method for coin flipping
 - Example 3: ML method for a biased die
2. **Generalized Linear Models (GLM)**
 - Generative vs. discriminative classifiers
 - Basic GLM setting
 - The Logistic Regression model
 - The Poisson regression model
3. **Maximum Likelihood (ML) framework (II)**
 - Maximum Likelihood theory
 - Bias and Variance of a ML estimate
 - Some properties of ML
 - Examples

Maximum likelihood

Gentle motivation

We have a r.v. X (e.g., the height of a randomly chosen Dutch)

The population has some distribution, which is assumed to have a special form. A common choice for a continuous distribution is the Gaussian (or normal) density*:

$$f(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{1}{2} \frac{(x - \mu)^2}{\sigma^2} \right),$$

and written $X \sim N(x; \mu, \sigma^2)$, where μ , the mean, and σ^2 , the variance, are the **parameters** of the distribution

(*) The probability that $X \in (x - \Delta x, x + \Delta x)$ approaches $2f(x)\Delta x$ as $\Delta x \rightarrow 0$.

Maximum likelihood

Gentle motivation

- Suppose we take an i.i.d sample $D = \{x_1, \dots, x_N\}$ of the r.v. X
- From the sample, we wish to **estimate** μ (it could be σ^2)
- It is not clear *a priori* what is the best way to do this:
 1. the average of D ?
 2. the median of D ?
 3. the average of the minimum and the maximum in D ?

Maximum likelihood

Example 1: the Gaussian

The **likelihood** of observing a particular x_n is $f(x_n; \mu, \sigma^2)$

The **likelihood** of seeing all the sample D is $\prod_{n=1}^N f(x_n; \mu, \sigma^2)$

Viewing this as *a function of the parameters*, we define

$$\mathcal{L}(\mu, \sigma^2) = P(D | \mu, \sigma^2) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2} \frac{(x_n - \mu)^2}{\sigma^2}\right)$$

“how likely it is that the population has parameters μ and σ^2 *given* the observed data sample D ”

The **maximum likelihood** estimates for the parameters are the values $\hat{\mu}$ and $\hat{\sigma}^2$ that maximize $\mathcal{L}(\mu, \sigma^2; D)$

Maximum likelihood

Example 1: the Gaussian

It is sometimes convenient (and equivalent) to maximize the “log-likelihood” :

$$l(\mu, \sigma^2) := \ln \mathcal{L}(\mu, \sigma^2) = \sum_{n=1}^N \ln f(x_n; \mu, \sigma^2)$$

In the Gaussian case, we get:

$$l(\mu, \sigma^2) = \sum_{n=1}^N \left[\ln \frac{1}{\sqrt{2\pi}} - \ln \sigma - \frac{1}{2} \left(\frac{x_n - \mu}{\sigma} \right)^2 \right]$$

Maximum likelihood

Example 1: the Gaussian

1. We compute $\frac{\partial l}{\partial \mu} = \frac{1}{\sigma^2} \sum_{n=1}^N (x_n - \mu)$,
2. We make $\frac{\partial l}{\partial \mu} = 0$, obtaining $\hat{\mu} = \frac{1}{N} \sum_{n=1}^N x_n$, the average of the sample
3. Also $\frac{\partial^2 l}{\partial \mu \partial \mu} = -\frac{N}{\sigma^2} < 0$, and therefore we have found a maximum

The estimate for the variance is $\hat{\sigma}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \hat{\mu})^2$

Maximum likelihood

Example 2: coin flipping

- Suppose we flip a coin that turns up heads with probability p . Find the ML estimate for p
- We take a sample $D = \{x_1, \dots, x_N\}$ out of N flips and get n_1 heads and $N - n_1$ tails. The number of heads follows a binomial $B(N, p)$
- The likelihood is $\mathcal{L}(p) = \binom{N}{n_1} p^{n_1} (1 - p)^{N - n_1}$
- The log-likelihood is $l(p) = \ln \binom{N}{n_1} + n_1 \ln p + (N - n_1) \ln(1 - p)$

$$\frac{\partial l}{\partial p} = \frac{n_1}{p} - \frac{N - n_1}{1 - p} = 0, \quad \text{therefore} \quad \hat{p} = \frac{n_1}{N}$$

Generalized Linear Models (GLM)

- Very general and classical technique for fitting **linear models**
- Genuine representatives of **discriminative** methods
- Work for many **target types**:
 - binary (two-class) and nominal (multi-class)
 - proportions and counts
 - ordinal (ordered classes)
 - continuous
- Admit **general predictors** (categorical ones are binarized)

Generalized Linear Models

- A GLM is a linear predictor of a convenient function of the **conditional expectation** of the target variable:

$$h(\mathbb{E}[T_n|\mathbf{X}_n]) = \boldsymbol{\beta}^\top \mathbf{X}_n + \beta_0$$

- This convenient function h is typically smooth and invertible and called the **link function**
- We **optimize** the $\boldsymbol{\beta}$ and β_0 parameters directly (without distributional assumptions on the \mathbf{x})
- The T_n are taken as independent and drawn from a distribution of the **exponential family** (Poisson, Gaussian, Bernoulli, Gamma, ...)

Generalized Linear Models

The modeller chooses a suitable distribution for T_n given X_n :

1. Gaussian distribution (h is identity): **linear regression**
 2. Bernoulli distribution (h is logit): **logistic regression**
 3. Poisson distribution (h is \ln): **Poisson regression**
- This generality comes at a cost: in general we need an iterative procedure for the optimization of the β and β_0 parameters
 - A popular procedure is to set it up as a ML problem and use a preferred numerical optimization method (e.g. Newton-Raphson)

Logistic regression

Introduction

We are in the case of binary classification ($K = 2$). We **model** the posterior probability for class ω_1 as:

$$P(\omega_1|\mathbf{x}) = g(\boldsymbol{\beta}^\top \mathbf{x} + \beta_0)$$

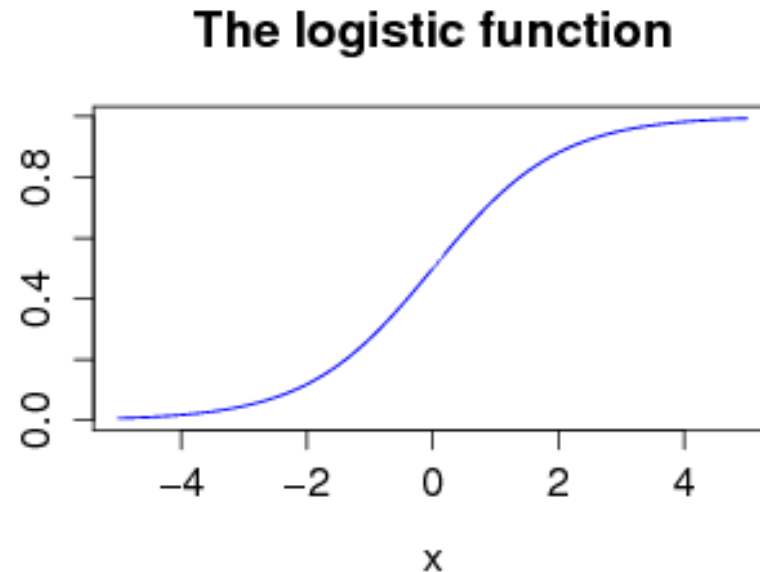
where $g(z) = \frac{\exp(z)}{1 + \exp(z)} = \frac{1}{1 + \exp(-z)}$ is the **logistic function**

obviously $P(\omega_2|\mathbf{x}) = 1 - P(\omega_1|\mathbf{x}) = 1 - g(\boldsymbol{\beta}^\top \mathbf{x} + \beta_0)$

Logistic regression

The logistic function

The logistic function is a C^∞ function $g : \mathbb{R} \longrightarrow (0, 1)$



This function is a bijection (one-to-one), with inverse $g^{-1}(z) = \ln \left(\frac{z}{1-z} \right)$, for $z \in (0, 1)$ (the **logit function**)

Logistic regression

Logistic Regression as a GLM

Each $T_n \sim \text{Ber}(p_n)$, where $p_n = g(\beta^\top \mathbf{X}_n + \beta_0)$,

$$P(T|\mathbf{X}, \beta) = \begin{cases} p_n & \text{if } T = 1 \ (\mathbf{X} \in \omega_1) \\ 1 - p_n & \text{if } T = 0 \ (\mathbf{X} \in \omega_2) \end{cases}$$

$$= p_n^T (1 - p_n)^{(1-T)}$$

Note $g(\beta^\top \mathbf{X}_n + \beta_0) = \mathbb{E}[T_n] = p_n$ (GLM setting)

→ we are identifying p_n with $P(\omega_1|\mathbf{X}_n)$

Logistic regression

Interpretation of our model

“The log of the odds is a linear function of the predictors”

Since $P(\omega_1|\mathbf{X}) = g(\boldsymbol{\beta}^\top \mathbf{X} + \beta_0)$

we have

$$\ln \left(\frac{P(\omega_1|\mathbf{X})}{P(\omega_2|\mathbf{X})} \right) = \ln \left(\frac{P(\omega_1|\mathbf{X})}{1 - P(\omega_1|\mathbf{X})} \right) = \text{logit}(P(\omega_1|\mathbf{X})) = \boldsymbol{\beta}^\top \mathbf{X} + \beta_0$$

Logistic regression

Let's go for the parameters

Suppose we have an i.i.d. sample of N *labelled* observations $S = \{(\mathbf{x}_n, t_n)\}_{n=1, \dots, N}$, where $\mathbf{x}_n \in \mathbb{R}^d, t_n \in \{0, 1\}$

1. The first thing we note is that we have $d + 1$ parameters to fit
2. In the “equivalent” generative case (LDA), we had $\frac{d(d+1)}{2} + 2d$

Let's re-write $P(\omega_1 | \mathbf{X}) = g(\boldsymbol{\beta}^\top \mathbf{X})$

with $\mathbf{X} = (1, X_1, \dots, X_d)^\top$ and $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_d)^\top$.

Logistic regression

ML for the Logistic Regression

$$\begin{aligned}l(\beta) &= \ln \mathcal{L}(\beta) \\&= \ln P(\{t_n\}_{n=1}^N | \{\mathbf{x}_n\}_{n=1}^N, \beta) \\&= \ln \prod_{n=1}^N P(t_n | \mathbf{x}_n, \beta) \\&= \sum_{n=1}^N \ln P(t_n | \mathbf{x}_n, \beta) \\&= \sum_{n=1}^N \ln \left(g(\beta^\top \mathbf{x}_n)^{t_n} (1 - g(\beta^\top \mathbf{x}_n))^{(1-t_n)} \right) \\&= \sum_{n=1}^N \ln \left((p_n)^{t_n} (1 - p_n)^{(1-t_n)} \right), \quad p_n = g(\beta^\top \mathbf{x}_n)\end{aligned}$$

Logistic regression

ML for the Logistic Regression

Now

$$\begin{aligned}(p_n)^{t_n}(1 - p_n)^{(1-t_n)} &= \left(\frac{p_n}{1 - p_n}\right)^{t_n} (1 - p_n) \\ &= \left(\exp(\boldsymbol{\beta}^\top \mathbf{x}_n)\right)^{t_n} \left(1 + \exp(\boldsymbol{\beta}^\top \mathbf{x}_n)\right)^{-1}\end{aligned}$$

Therefore

$$l(\boldsymbol{\beta}) = \sum_{n=1}^N \left[t_n \boldsymbol{\beta}^\top \mathbf{x}_n - \ln \left(1 + \exp(\boldsymbol{\beta}^\top \mathbf{x}_n) \right) \right]$$

Logistic regression

Newton-Raphson

The ML of the logistic regression does not have a closed-form solution:

$$\begin{aligned}\beta^{k+1} &= \beta^k - \left(\frac{\partial^2 l}{\partial \beta \partial \beta^\top} \right)^{-T} \left(\frac{\partial l}{\partial \beta} \right) \\ &= \beta^k + (X^\top W X)^{-1} X^\top (t - p) \\ &= (X^\top W X)^{-1} X^\top W z, \quad z = X \beta^\top + W^{-1}(t - p)\end{aligned}$$

where:

X is the matrix of the $\{\mathbf{x}_n\}$

$W = \text{diag}(p_n(1 - p_n)), \quad n = 1, \dots, N$

$\mathbf{t} = (t_1, \dots, t_N)^\top, \quad \mathbf{p} = (p_1, \dots, p_N)^\top$

since:

$$\frac{\partial l}{\partial \beta} = X^\top (t - p)$$

$$\frac{\partial^2 l}{\partial \beta \partial \beta^\top} = -X^\top W X$$

Logistic regression

Iterated Reweighted Least Squares (IRLS)

1. Initialize $\beta_0 \leftarrow \ln \left(\frac{P(\omega_1)}{1 - P(\omega_1)} \right)$ and $\beta_i \leftarrow 0, i = 1, \dots, d$ (*null model*)
2. Iterate until convergence:
 - a) Update $\mathbf{p} \leftarrow (p_1, \dots, p_N)^\top$, where $p_n = g(\boldsymbol{\beta}^\top \mathbf{x}_n)$
 - b) Update $W^{-1} \leftarrow \text{diag} \left(\frac{1}{p_n(1-p_n)} \right), n = 1, \dots, N$
 - c) Update $\mathbf{z} \leftarrow X\boldsymbol{\beta}^\top + W^{-1}(\mathbf{t} - \mathbf{p})$
 - d) Update $\boldsymbol{\beta} \leftarrow (X^\top W X)^{-1} X^\top W \mathbf{z}$
3. return $\hat{\boldsymbol{\beta}}$

Logistic regression

The Deviance and the AIC

In the context of Generalized Linear Models,

$$-2l(\hat{\beta}) = -2 \ln \mathcal{L}(\hat{\beta}) \sim \chi_{v=N-d-1}^2$$

is called the **deviance** (in ML, this is the **error**)

Null deviance: deviance of the null model (just with constant term)

Residual deviance: deviance of the proposed model

AIC: deviance with complexity penalization $-2l(\hat{\beta}) + 2d$

Actually, $2d \approx 2\|\hat{\beta}\|_0$ is a rudimentary form of **regularization**

An example of R's `glm()` in action

```
glm (formula = chd ~ age + sbp + ldl + adiposity + alcohol + tobacco +  
      obesity + famhist + typea, family = binomial, data = SAheart.learn)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-1.9287	-0.8283	-0.3762	0.8983	2.4722

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	-7.066884	1.654749	-4.271	1.95e-05	***
age	0.060443	0.015789	3.828	0.000129	***
sbp	0.005366	0.007127	0.753	0.451531	
ldl	0.195473	0.078493	2.490	0.012763	*
adiposity	-0.009770	0.035926	-0.272	0.785658	
alcohol	-0.001653	0.005979	-0.277	0.782146	
tobacco	0.090333	0.033137	2.726	0.006411	**
obesity	-0.028118	0.053949	-0.521	0.602229	
famhistPresent	0.912713	0.283802	3.216	0.001300	**
typea	0.041579	0.015356	2.708	0.006777	**

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Null deviance: 401.48 Residual deviance: 308.55

AIC: 328.55 Number of Fisher Scoring iterations: 5

Logistic regression

An example of R's `glm()` in action

Use of the AIC to simplify the model:

Coefficients:

(Intercept)	age	ldl	tobacco	famhistPresent
-7.12058	0.05943	0.17677	0.08927	0.89552
typea				
0.04012				

Null Deviance: 401.5

Residual Deviance: 310.1 AIC: 322.1

Poisson regression

Introduction

In many statistical studies, one tries to relate a count to some scientific variables:

1. Number of cardio-vascular accidents among people over 60 in a US state \sim average income in the state
2. Number of bicycles in a Danish household \sim distance to the city centre
3. Number of incoming calls to a complaints number \sim hourly interval

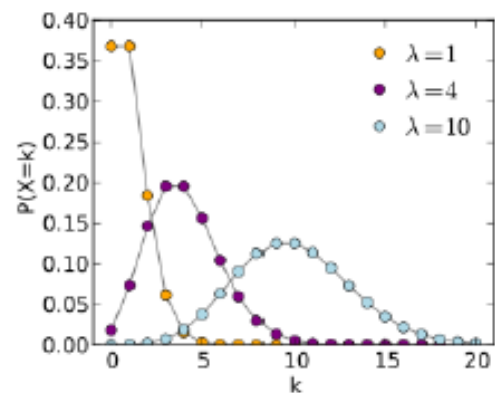
When the response is a count which does not have any natural upper bound, the logistic regression is not appropriate.

→ The **Poisson regression** is a natural alternative

Poisson regression

The Poisson distribution

This is a discrete distribution $X \sim \text{Pois}(\lambda)$ with probability mass function:



$$P(X = k) = e^{-\lambda} \frac{\lambda^k}{k!}, \quad \lambda > 0 \in \mathbb{R}, k = 0, 1, 2, \dots$$

Poisson regression

Poisson regression

- We consider independent Poisson random variables T_1, \dots, T_N with $T_n \sim \text{Pois}(\lambda_n)$. We know that $\mathbb{E}[T_n] = \lambda_n$
- We have an i.i.d. sample of N observations $\mathbf{x}_n \in \mathbb{R}^d$ (e.g. distance to the city centre)
- We have a corresponding sample t_1, \dots, t_N , where each t_n is drawn from T_n (e.g. number of bicycles)
- The idea is to model λ_n as $\exp(\boldsymbol{\beta}^\top \mathbf{x}_n + \beta_0)$ (the link function is \ln)
- The Poisson regression model is $T_n \sim \text{Pois}(\exp(\boldsymbol{\beta}^\top \mathbf{x}_n + \beta_0))$ or $\ln \lambda_n = \boldsymbol{\beta}^\top \mathbf{x}_n + \beta_0$

Poisson regression

ML for the Poisson regression

- Proceeding like the previous case for Logistic regression, we arrive at:

$$l(\beta) = \sum_{n=1}^N \left[-\exp(\beta^T \mathbf{x}_n) + t_n \beta^T \mathbf{x}_n - \ln(t_n!) \right]$$

- Again, this expression has no closed-form solution; however, we can still use NR (because minus l is convex)

Maximum likelihood

Some theory

1. Consider a random vector X_1, \dots, X_N each of which having mass (or density) $f(x_n; \theta)$ and a random sample $\{x_1, \dots, x_N\}$ thereof
2. The **likelihood function** is defined as the product:

$$\mathcal{L}(\theta) = f(x_1; \theta) f(x_2; \theta) \cdots f(x_n; \theta) = \prod_{n=1}^N f(x_n; \theta)$$

(**note:** x_i are numbers while X_i are random variables)

3. The likelihood is considered a function of θ for fixed data (whereas the mass or density are considered a function of x for fixed θ)

It can be read as the probability of observing the available sample for different values of the θ parameter

Maximum likelihood

Some theory

- The **ML estimator** (MLE) $\hat{\Theta}(X_1, \dots, X_N)$ is the *function* of the random vector that maximizes the likelihood with respect to θ
- The **ML estimate** (MLE) $\hat{\theta}_N = \hat{\Theta}(x_1, \dots, x_N)$ gives a *point* estimation for the available data
(the MLE is a function while the MLe is a number)
- How good are these estimates?
 - **Unbiasedness**: is the expected value of the estimate the number being estimated?
 - **Efficiency**: is the variance of some estimate smaller than that of another estimate?

Maximum likelihood

Bias

Let $\hat{\theta}_N = \hat{\Theta}(x_1, \dots, x_N)$. The **bias** of $\hat{\theta}_N$ is:

$$\text{Bias}(\hat{\theta}) = \mathbb{E}[\hat{\theta}_N] - \theta$$

We say that an estimate $\hat{\theta}_N$ is **unbiased** if:

$$\mathbb{E}[\hat{\theta}_N] = \theta$$

the expected value is taken over all possible samples of a given size N

Maximum likelihood

Variance

The **variance** of $\hat{\theta}_N$ is:

$$\text{Var}(\hat{\theta}_N) = \mathbb{E} [(\hat{\theta}_N - \mathbb{E}[\hat{\theta}_N])^2]$$

We say that an estimate $\hat{\theta}$ is **more efficient** than an estimate $\hat{\theta}'$ if $\text{Var}(\hat{\theta}) < \text{Var}(\hat{\theta}')$

There is a lower limit on the variance of estimates (Cramér-Rao bound)

Maximum likelihood

Consistency

We say that an estimate $\hat{\theta}_N$ is **consistent** if it tends in probability to θ as $N \rightarrow \infty$. Formally:

$$\forall \epsilon > 0, \lim_{N \rightarrow \infty} P(|\hat{\theta}_N - \theta| < \epsilon) = 1$$

Consistency implies that:

1. Bias decreases with sample size
2. Variance decreases with sample size

Maximum likelihood

Some general properties of ML

Asymptotic unbiasedness :

- Any ML estimate is at least asymptotically unbiased (hence, when the sample size is large, it is approximately unbiased)
- It is the best estimate as the sample size N becomes large (other estimates can be better for small N)

Asymptotic efficiency : The ML estimate is the more efficient among those that are consistent

Invariance : If $\hat{\Theta}$ is the MLE of θ and g is bijective (one-to-one), then the MLE of $g(\Theta)$ is $g(\hat{\Theta})$

Maximum likelihood

Example 1: the Gaussian

We had previously found that:

- The ML estimate for the mean μ of a Gaussian is

$$\hat{\mu}_N = \frac{1}{N} \sum_{n=1}^N x_n, \text{ the sample mean}$$

- The ML estimate for the variance σ^2 of a Gaussian is

$$\hat{\sigma}_N^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \hat{\mu}_N)^2, \text{ the sample variance}$$

How good are these estimates?

Do they have a bias? What is their variance?
Are they consistent? Efficient?

Maximum likelihood

Example: the Gaussian

$$\begin{aligned}\text{Bias}(\hat{\mu}_N) &= \mathbb{E}[\hat{\mu}_N] - \mu \\ &= \mathbb{E}\left[\frac{1}{N} \sum_{n=1}^N x_n\right] - \mu \\ &= \left(\frac{1}{N} \sum_{n=1}^N \mathbb{E}[x_n]\right) - \mu \\ &= \left(\frac{1}{N} \sum_{n=1}^N \mu\right) - \mu \\ &= \frac{N\mu}{N} - \mu = 0\end{aligned}$$

Maximum likelihood

Example: the Gaussian

$$\text{Var}(\hat{\mu}_N) = \mathbb{E}[(\hat{\mu}_N - \mathbb{E}[\hat{\mu}_N])^2]$$

$$\begin{aligned}\text{Var}(\hat{\mu}_N) &= \text{Var} \left[\frac{1}{N} \sum_{n=1}^N x_n \right] \\ &= \frac{1}{N^2} \sum_{n=1}^N \text{Var}[x_n] \\ &= \frac{1}{N^2} \sum_{n=1}^N \sigma^2 = \frac{N\sigma^2}{N^2} = \frac{\sigma^2}{N}\end{aligned}$$

Since $\lim_{N \rightarrow \infty} \frac{\sigma^2}{N} = 0$, it is **unbiased** and **consistent**

Maximum likelihood

Example: the Gaussian

Since

$$\mathbb{E}[\hat{\sigma}_N^2] = \mathbb{E} \left[\frac{1}{N} \sum_{n=1}^N (x_n - \hat{\mu}_N)^2 \right] = \left(\frac{N-1}{N} \right) \sigma^2$$

this estimator is biased! (it underestimates the true value by a factor of $\frac{N-1}{N}$)

Consider $\tilde{\sigma}_N^2 = \frac{N}{N-1} \hat{\sigma}_N^2 = \frac{1}{N-1} \sum_{n=1}^N (x_n - \hat{\mu}_N)^2$

$$\text{Now } \mathbb{E}[\tilde{\sigma}_N^2] = \mathbb{E} \left[\frac{N}{N-1} \hat{\sigma}_N^2 \right] = \frac{N}{N-1} \mathbb{E}[\hat{\sigma}_N^2] = \frac{N}{N-1} \left(\frac{N-1}{N} \right) \sigma^2 = \sigma^2$$

The expression for $\text{Var}[\hat{\sigma}_N^2]$ (or $\text{Var}[\tilde{\sigma}_N^2]$) is more complex but is of the form $O(\frac{1}{N})$. Therefore both $\hat{\sigma}_N^2$ and $\tilde{\sigma}_N^2$ are **consistent** estimators of σ^2 .

Maximum likelihood

Example: the Gaussian

It is important to note that unbiased and consistent estimators are not always the best estimators! Sometimes efficiency is a must!

When μ is unknown, the estimator:

$$\bar{\sigma}_N^2 = k_N^2 \cdot \tilde{\sigma}_N^2 \approx \frac{1}{N - 1,45} \sum_{n=1}^N (x_n - \hat{\mu}_N)^2$$

–with $k_N = \sqrt{\frac{N-1}{2}} \Gamma(\frac{N-1}{2}) / \Gamma(\frac{N}{2})$, being Γ the gamma function–

can be proven to be unbiased and more efficient than $\hat{\sigma}_N^2$ or $\tilde{\sigma}_N^2$.

Machine Learning

Syllabus

1. Introduction to Machine Learning
2. Theoretical issues (I): regression
3. Linear regression and beyond
4. Theoretical issues (II): classification
5. Generative classifiers
6. Discriminative classifiers

7. Clustering
8. Learning with kernels (I): The SVM
9. Learning with kernels (II): Kernel functions
10. Learning with kernels (III): Other kernel methods
11. Artificial neural networks (I): the MLP
12. Artificial neural networks (II): the RBF
13. Ensemble methods: Random Forests
14. Advanced topics and frontiers