

# The GPML Toolbox version 3.4

Carl Edward Rasmussen & Hannes Nickisch

June 10, 2014

## Abstract

The GPML toolbox is an Octave 3.2.x and Matlab 7.x implementation of inference and prediction in Gaussian process (GP) models. It implements algorithms discussed in Rasmussen & Williams: *Gaussian Processes for Machine Learning*, the MIT press, 2006 and Nickisch & Rasmussen: *Approximations for Binary Gaussian Process Classification*, JMLR, 2008.

The strength of the function lies in its flexibility, simplicity and extensibility. The function is flexible as firstly it allows specification of the properties of the GP through definition of mean function and covariance functions. Secondly, it allows specification of different inference procedures, such as e.g. exact inference and Expectation Propagation (EP). Thirdly it allows specification of likelihood functions e.g. Gaussian or Laplace (for regression) and e.g. cumulative Logistic (for classification). Simplicity is achieved through a single function and compact code. Extensibility is ensured by modular design allowing for easy addition of extension for the already fairly extensive libraries for inference methods, mean functions, covariance functions and likelihood functions.

This document is a technical manual for a developer containing many details. If you are not yet familiar with the GPML toolbox, the *user documentation* and examples therein are a better way to get started.

# Contents

<b>1</b>	<b>Gaussian Process Training and Prediction</b>	<b>3</b>
<b>2</b>	<b>The gp Function</b>	<b>4</b>
<b>3</b>	<b>Inference Methods</b>	<b>9</b>
3.1	Exact Inference . . . . .	10
3.2	Laplace's Approximation . . . . .	11
3.3	Expectation Propagation . . . . .	11
3.4	Kullback Leibler Divergence Minimisation . . . . .	12
3.5	Variational Bayes . . . . .	13
3.6	FITC Approximations . . . . .	14
<b>4</b>	<b>Likelihood Functions</b>	<b>15</b>
4.1	Prediction . . . . .	15
4.2	Interface . . . . .	16
4.3	Implemented Likelihood Functions . . . . .	18
4.4	Usage of Implemented Likelihood Functions . . . . .	19
4.5	Compatibility Between Likelihoods and Inference Methods . . . . .	20
4.6	Gaussian Likelihood . . . . .	20
4.6.1	Exact Inference . . . . .	22
4.6.2	Laplace's Approximation . . . . .	22
4.6.3	Expectation Propagation . . . . .	22
4.6.4	Variational Bayes . . . . .	22
4.7	Warped Gaussian Likelihood . . . . .	23
4.8	Gumbel Likelihood . . . . .	24
4.9	Laplace Likelihood . . . . .	24
4.10	Student's t Likelihood . . . . .	27
4.11	Cumulative Logistic Likelihood . . . . .	27
4.12	GLM Likelihoods: Poisson, Weibull, Gamma, Exponential, Inverse Gaussian and Beta	28
4.12.1	Inverse Link Functions . . . . .	28
4.12.2	Poisson Likelihood . . . . .	29
4.12.3	Weibull Likelihood . . . . .	30
4.12.4	Gamma Likelihood . . . . .	30
4.12.5	Exponential Likelihood . . . . .	31
4.12.6	Inverse Gaussian Likelihood . . . . .	31
4.12.7	Beta Likelihood . . . . .	31
<b>5</b>	<b>Mean Functions</b>	<b>32</b>
5.1	Interface . . . . .	32
5.2	Implemented Mean Functions . . . . .	33
5.3	Usage of Implemented Mean Functions . . . . .	33
<b>6</b>	<b>Covariance Functions</b>	<b>35</b>
6.1	Interface . . . . .	35
6.2	Implemented Covariance Functions . . . . .	38
6.3	Usage of Implemented Covariance Functions . . . . .	38

# 1 Gaussian Process Training and Prediction

The `gpm1` toolbox contains a single user function `gp` described in section 2. In addition there are a number of supporting structures and functions which the user needs to know about, as well as an internal convention for representing the posterior distribution, which may not be of direct interest to the casual user.

**Inference Methods** An inference method is a function which computes the (approximate) posterior, the (approximate) negative log marginal likelihood and its partial derivatives w.r.t.. the hyperparameters, given a model specification (i.e., GP mean and covariance functions and a likelihood function) and a data set. Inference methods are discussed in section 3. New inference methods require a function providing the desired inference functionality and possibly extra functionality in the likelihood functions applicable.

**Hyperparameters** The hyperparameters is a struct controlling the properties of the model, i.e.. the GP mean and covariance function and the likelihood function. The hyperparameters is a struct with the three fields `mean`, `cov` and `lik`, each of which is a vector. The number of elements in each field must agree with number of hyperparameters in the specification of the three functions they control (below). If a field has no entries it can either be empty or non-existent.

**Likelihood Functions** The likelihood function specifies the form of the likelihood of the GP model and computes terms needed for prediction and inference. For inference, the required properties of the likelihood depend on the inference method, including properties necessary for hyperparameter learning, section 4.

**Mean Functions** The mean function is a cell array specifying the GP mean. It computes the mean and its derivatives w.r.t.. the part of the hyperparameters pertaining to the mean. The cell array allows flexible specification and composition of mean functions, discussed in section 5. The default is the zero function.

**Covariance Functions** The covariance function is a cell array specifying the GP covariance function. It computes the covariance and its derivatives w.r.t.. the part of the hyperparameters pertaining to the covariance function. The cell array allows flexible specification and composition of covariance functions, discussed in section 6.

Inference methods, see section 3, compute (among other things) an approximation to the posterior distribution of the latent variables  $f_i$  associated with the training cases,  $i = 1, \dots, n$ . This approximate posterior is assumed to be Gaussian, and is communicated via a struct `post` with the fields `post.alpha`, `post.sW` and `post.L`. Often, starting from the Gaussian prior  $p(f) = \mathcal{N}(f|m, K)$  the approximate posterior admits the form

$$q(f|\mathcal{D}) = \mathcal{N}(f|\mu = m + K\alpha, V = (K^{-1} + W)^{-1}), \text{ where } W \text{ diagonal with } W_{ii} = s_i^2. \quad (1)$$

In such cases, the entire posterior can be computed from the two vectors `post.alpha` and `post.sW`; the inference method may optionally also return  $L = \text{chol}(\text{diag}(s)\text{diag}(s)K + I)$ .

If on the other hand the posterior doesn't admit the above form, then `post.L` returns the matrix  $L = -(K + W^{-1})^{-1}$  (and `post.sW` is unused). In addition, a sparse representation of the posterior may be used, in which case the non-zero elements of the `post.alpha` vector indicate the active entries.

## 2 The gp Function

The gp function is typically the only function the user would directly call.

```
4a <gp.m 4a>≡
1 function [varargout] = gp(hyp, inf, mean, cov, lik, x, y, xs, ys)
2 <gp function help 4b>
3 <initializations 5b>
4 <inference 6b>
5 if nargin==7                                     % if no test cases are provided
6     varargout = {nlZ, dnlZ, post};                % report -log marg lik, derivatives and post
7 else
8     <compute test predictions 7a>
9 end
```

It offers facilities for training the hyperparameters of a GP model as well as predictions at unseen inputs as detailed in the following help.

```
4b <gp function help 4b>≡ (4a)
1 % Gaussian Process inference and prediction. The gp function provides a
2 % flexible framework for Bayesian inference and prediction with Gaussian
3 % processes for scalar targets, i.e. both regression and binary
4 % classification. The prior is Gaussian process, defined through specification
5 % of its mean and covariance function. The likelihood function is also
6 % specified. Both the prior and the likelihood may have hyperparameters
7 % associated with them.
8 %
9 % Two modes are possible: training or prediction: if no test cases are
10 % supplied, then the negative log marginal likelihood and its partial
11 % derivatives w.r.t. the hyperparameters is computed; this mode is used to fit
12 % the hyperparameters. If test cases are given, then the test set predictive
13 % probabilities are returned. Usage:
14 %
15 %     training: [nlZ dnlZ                ] = gp(hyp, inf, mean, cov, lik, x, y);
16 % prediction: [ymu ys2 fmu fs2         ] = gp(hyp, inf, mean, cov, lik, x, y, xs);
17 %             or: [ymu ys2 fmu fs2 lp] = gp(hyp, inf, mean, cov, lik, x, y, xs, ys);
18 %
19 % where:
20 %
21 %     hyp      column vector of hyperparameters
22 %     inf      function specifying the inference method
23 %     cov      prior covariance function (see below)
24 %     mean     prior mean function
25 %     lik      likelihood function
26 %     x        n by D matrix of training inputs
27 %     y        column vector of length n of training targets
28 %     xs       ns by D matrix of test inputs
29 %     ys       column vector of length nn of test targets
30 %
31 %     nlZ      returned value of the negative log marginal likelihood
32 %     dnlZ     column vector of partial derivatives of the negative
33 %             log marginal likelihood w.r.t. each hyperparameter
34 %     ymu      column vector (of length ns) of predictive output means
35 %     ys2      column vector (of length ns) of predictive output variances
36 %     fmu      column vector (of length ns) of predictive latent means
37 %     fs2      column vector (of length ns) of predictive latent variances
38 %     lp       column vector (of length ns) of log predictive probabilities
39 %
```

```

40 % post      struct representation of the (approximate) posterior
41 %           3rd output in training mode or 6th output in prediction mode
42 %           can be reused in prediction mode gp(..., cov, lik, x, post, xs,...)
43 %
44 % See also covFunctions.m, infMethods.m, likFunctions.m, meanFunctions.m.
45 %
46 <gpml copyright 5a>

```

```

5a <gpml copyright 5a>≡ (4b 9 10 16 19 21a 32 33 35 39)
1 % Copyright (c) by Carl Edward Rasmussen and Hannes Nickisch, 2014-06-10.
2 % File automatically generated using noweb.

```

Depending on the number of input parameters, gp knows whether it is operated in training or in prediction mode. The highlevel structure of the code is as follows. After some initialisations, we perform inference and decide whether test set predictions are needed or only the result of the inference is demanded.

```

5b <initializations 5b>≡ (4a)
1 <minimalist usage 5c>
2 <process input arguments 5d>
3 <check hyperparameters 6a>

```

If the number of input arguments is incorrect, we echo a minimalist usage and return.

```

5c <minimalist usage 5c>≡ (5b)
1 if nargin<7 || nargin>9
2   disp('Usage: [n1Z dn1Z      ] = gp(hyp, inf, mean, cov, lik, x, y);')
3   disp('      or: [ymu ys2 fmu fs2  ] = gp(hyp, inf, mean, cov, lik, x, y, xs);')
4   disp('      or: [ymu ys2 fmu fs2 lp] = gp(hyp, inf, mean, cov, lik, x, y, xs, ys);')
5   return
6 end

```

Set some useful default values for empty arguments, and convert inf and lik to function handles and mean and cov to cell arrays if necessary. Initialize variables.

```

5d <process input arguments 5d>≡ (5b)
1 if isempty(mean), mean = {@meanZero}; end % set default mean
2 if ischar(mean) || isa(mean, 'function_handle'), mean = {mean}; end % make cell
3 if isempty(cov), error('Covariance function cannot be empty'); end % no default
4 if ischar(cov) || isa(cov, 'function_handle'), cov = {cov}; end % make cell
5 cov1 = cov{1}; if isa(cov1, 'function_handle'), cov1 = func2str(cov1); end
6 if strcmp(cov1,'covFITC') && isfield(hyp,'xu'), cov{3} = hyp.xu; end %use hyp.xu
7 if isempty(inf) % set default inference method
8   if strcmp(cov1,'covFITC'), inf = @infFITC; else inf = @infExact; end
9 else
10  if iscell(inf), inf = inf{1}; end % cell input is allowed
11  if ischar(inf), inf = str2func(inf); end % convert into function handle
12 end
13 if strcmp(cov1,'covFITC') % only infFITC* are possible
14  if isempty(strfind(func2str(inf),'infFITC')==1)
15    error('Only infFITC* are possible inference algorithms')
16  end
17 end % only one possible class of inference algorithms
18 if isempty(lik), lik = {@likGauss}; end % set default lik
19 if ischar(lik) || isa(lik, 'function_handle'), lik = {lik}; end % make cell
20 if iscell(lik), likstr = lik{1}; else likstr = lik; end
21 if ~ischar(likstr), likstr = func2str(likstr); end
22
23 D = size(x,2);

```

Check that the sizes of the hyperparameters supplied in `hyp` match the sizes expected. The three parts `hyp.mean`, `hyp.cov` and `hyp.lik` are checked separately, and define empty entries if they don't exist.

```
6a <check hyperparameters 6a>≡ (5b)
1 if ~isfield(hyp,'mean'), hyp.mean = []; end % check the hyp specification
2 if eval(feval(mean{:})) ~= numel(hyp.mean)
3     error('Number of mean function hyperparameters disagree with mean function')
4 end
5 if ~isfield(hyp,'cov'), hyp.cov = []; end
6 if eval(feval(cov{:})) ~= numel(hyp.cov)
7     error('Number of cov function hyperparameters disagree with cov function')
8 end
9 if ~isfield(hyp,'lik'), hyp.lik = []; end
10 if eval(feval(lik{:})) ~= numel(hyp.lik)
11     error('Number of lik function hyperparameters disagree with lik function')
12 end
```

Inference is performed by calling the desired inference method `inf`. In training mode, we accept a failure of the inference method (and issue a warning), since during hyperparameter learning, hyperparameters causing a numerical failure may be attempted, but the minimize function may gracefully recover from this. During prediction, failure of the inference method is an error.

```
6b <inference 6b>≡ (4a)
1 try % call the inference method
2     % issue a warning if a classification likelihood is used in conjunction with
3     % labels different from +1 and -1
4     if strcmp(likstr,'likErf') || strcmp(likstr,'likLogistic')
5         if ~isstruct(y)
6             uy = unique(y);
7             if any( uy~=+1 & uy~=-1 )
8                 warning('You try classification with labels different from {+1,-1}')
9             end
10        end
11    end
12    if nargin>7 % compute marginal likelihood and its derivatives only if needed
13        if isstruct(y)
14            post = y; % reuse a previously computed posterior approximation
15        else
16            post = inf(hyp, mean, cov, lik, x, y);
17        end
18    else
19        if nargin<=1
20            [post nlZ] = inf(hyp, mean, cov, lik, x, y); dnlZ = {};
21        else
22            [post nlZ dnlZ] = inf(hyp, mean, cov, lik, x, y);
23        end
24    end
25 catch
26    msgstr = lasterr;
27    if nargin > 7, error('Inference method failed [%s]', msgstr); else
28        warning('Inference method failed [%s] .. attempting to continue',msgstr)
29        dnlZ = struct('cov',0*hyp.cov, 'mean',0*hyp.mean, 'lik',0*hyp.lik);
30        varargout = {NaN, dnlZ}; return % continue with a warning
31    end
32 end
```

We copy the already computed negative log marginal likelihood to the first output argument, and if desired report its partial derivatives w.r.t. the hyperparameters if running in inference mode.

Predictions are computed in a loop over small batches to avoid memory problems for very large test sets.

7a  $\langle \text{compute test predictions } 7a \rangle \equiv$  (4a)

```

1 alpha = post.alpha; L = post.L; sW = post.sW;
2 if issparse(alpha) % handle things for sparse representations
3   nz = alpha ~= 0; % determine nonzero indices
4   if issparse(L), L = full(L(nz,nz)); end % convert L and sW if necessary
5   if issparse(sW), sW = full(sW(nz)); end
6 else nz = true(size(alpha,1),1); end % non-sparse representation
7 if numel(L)==0 % in case L is not provided, we compute it
8   K = feval(cov{:}, hyp.cov, x(nz,:));
9   L = chol(eye(sum(nz))+sW*sW'.*K);
10 end
11 Lchol = all(all(tril(L,-1)==0)&diag(L)')>0)&isreal(diag(L)); % L contains chol?
12 ns = size(xs,1); % number of data points
13 nperbatch = 1000; % number of data points per mini batch
14 nact = 0; % number of already processed test data points
15 ymu = zeros(ns,1); ys2 = ymu; fmu = ymu; fs2 = ymu; lp = ymu; % allocate mem
16 while nact<ns % process minibatches of test cases to save memory
17   id = (nact+1):min(nact+nperbatch,ns); % data points to process
18    $\langle \text{make predictions } 7b \rangle$ 
19   nact = id(end); % set counter to index of last processed data point
20 end
21 if nargin<9
22   varargout = {ymu, ys2, fmu, fs2, [], post}; % assign output arguments
23 else
24   varargout = {ymu, ys2, fmu, fs2, lp, post};
25 end

```

In every iteration of the above loop, we compute the predictions for all test points of the batch.

7b  $\langle \text{make predictions } 7b \rangle \equiv$  (7a)

```

1 kss = feval(cov{:}, hyp.cov, xs(id,:), 'diag'); % self-variance
2 if strcmp(cov1,'covFITC') % cross-covariances
3   Ks = feval(cov{:}, hyp.cov, x, xs(id,:)); Ks = Ks(nz,:); % res indep. of x
4 else
5   Ks = feval(cov{:}, hyp.cov, x(nz,:), xs(id,:)); % avoid computation
6 end
7 ms = feval(mean{:}, hyp.mean, xs(id,:));
8 N = size(alpha,2); % number of alphas (usually 1; more in case of sampling)
9 Fmu = repmat(ms,1,N) + Ks'*full(alpha(nz,:)); % conditional mean fs|f
10 fmu(id) = sum(Fmu,2)/N; % predictive means
11 if Lchol % L contains chol decomp => use Cholesky parameters (alpha,sW,L)
12   V = L'\(repmat(sW,1,length(id)).*Ks);
13   fs2(id) = kss - sum(V.*V,1)'; % predictive variances
14 else % L is not triangular => use alternative parametrisation
15   fs2(id) = kss + sum(Ks.*(L*Ks),1)'; % predictive variances
16 end
17 fs2(id) = max(fs2(id),0); % remove numerical noise i.e. negative variances
18 Fs2 = repmat(fs2(id),1,N); % we have multiple values in case of sampling
19 if nargin<9
20   [Lp, Ymu, Ys2] = feval(lik{:},hyp.lik,[],Fmu(:),Fs2(:));
21 else
22   Ys = repmat(ys(id),1,N);
23   [Lp, Ymu, Ys2] = feval(lik{:},hyp.lik,Ys(:),Fmu(:),Fs2(:));
24 end
25 lp(id) = sum(reshape(Lp, [],N),2)/N; % log probability; sample averaging
26 ymu(id) = sum(reshape(Ymu, [],N),2)/N; % predictive mean ys|y and ..

```

```
27 ys2(id) = sum(reshape(Ys2,[],N),2)/N; % .. variance
```



### 3 Inference Methods

Inference methods are responsible for computing the (approximate) posterior `post`, the (approximate) negative log marginal likelihood `nlZ` and its partial derivatives `dnlZ` w.r.t. the hyperparameters `hyp`. The arguments to the function are hyperparameters `hyp`, mean function `mean`, covariance function `cov`, likelihood function `lik` and training data `x` and `y`. Several inference methods are implemented and described this section.

```

9  <infMethods.m 9>≡
1  % Inference methods: Compute the (approximate) posterior for a Gaussian process.
2  % Methods currently implemented include:
3  %
4  %   infExact          Exact inference (only possible with Gaussian likelihood)
5  %   infLaplace        Laplace's Approximation
6  %   infEP             Expectation Propagation
7  %   infVB             Variational Bayes Approximation
8  %   infKL             Kullback-Leibler optimal Approximation
9  %
10 %   infFITC           Large scale regression with approximate covariance matrix
11 %   infFITC_Laplace   Large scale inference with approximate covariance matrix
12 %   infFITC_EP        Large scale inference with approximate covariance matrix
13 %
14 %   infMCMC           Markov Chain Monte Carlo and Annealed Importance Sampling
15 %                   We offer two samplers.
16 %                   - hmc: Hybrid Monte Carlo
17 %                   - ess: Elliptical Slice Sampling
18 %                   No derivatives w.r.t. to hyperparameters are provided.
19 %
20 %   infLOO            Leave-One-Out predictive probability and Least-Squares Approxim.
21 %
22 % The interface to the approximation methods is the following:
23 %
24 %   function [post nlZ dnlZ] = inf..(hyp, cov, lik, x, y)
25 %
26 % where:
27 %
28 %   hyp           is a struct of hyperparameters
29 %   cov           is the name of the covariance function (see covFunctions.m)
30 %   lik           is the name of the likelihood function (see likFunctions.m)
31 %   x             is a n by D matrix of training inputs
32 %   y             is a (column) vector (of size n) of targets
33 %
34 %   nlZ           is the returned value of the negative log marginal likelihood
35 %   dnlZ          is a (column) vector of partial derivatives of the negative
36 %                 log marginal likelihood w.r.t. each hyperparameter
37 %   post          struct representation of the (approximate) posterior containing
38 %       alpha     is a (sparse or full column vector) containing  $\text{inv}(K) * (\mu - m)$ ,
39 %                 where  $K$  is the prior covariance matrix,  $m$  the prior mean,
40 %                 and  $\mu$  the approx posterior mean
41 %       sW        is a (sparse or full column) vector containing diagonal of  $\text{sqrt}(W)$ 
42 %                 the approximate posterior covariance matrix is  $\text{inv}(\text{inv}(K) + W)$ 
43 %       L         is a (sparse or full) matrix,  $L = \text{chol}(sW * K * sW + \text{eye}(n))$ 
44 %
45 % Usually, the approximate posterior to be returned admits the form
46 %  $N(\mu = m + K * \alpha, V = \text{inv}(\text{inv}(K) + W))$ , where  $\alpha$  is a vector and  $W$  is diagonal;
47 % if not, then  $L$  contains instead  $-\text{inv}(K + \text{inv}(W))$ , and  $sW$  is unused.
48 %

```

```

49 % For more information on the individual approximation methods and their
50 % implementations, see the separate inf??.m files. See also gp.m
51 %
52 (gpml copyright 5a)

```

Not all inference methods are compatible with all likelihood functions, e.g.. exact inference is only possible with Gaussian likelihood. In order to perform inference, each method needs various properties of the likelihood functions, section 4.

### 3.1 Exact Inference

For Gaussian likelihoods, GP inference reduces to computing mean and covariance of a multivariate Gaussian which can be done exactly by simple matrix algebra. The program `inf/infExact.m` does exactly this. If it is called with a likelihood function other than the Gaussian, it issues an error. The Gaussian posterior  $q(f|\mathcal{D}) = \mathcal{N}(f|\mu, V)$  is exact.

```

10 (inf/infExact.m 10)≡
1  function [post nlZ dnlZ] = infExact(hyp, mean, cov, lik, x, y)
2
3  % Exact inference for a GP with Gaussian likelihood. Compute a parametrization
4  % of the posterior, the negative log marginal likelihood and its derivatives
5  % w.r.t. the hyperparameters. See also "help infMethods".
6  %
7  (gpml copyright 5a)
8  %
9  % See also INFMETHODS.M.
10
11 if iscell(lik), likstr = lik{1}; else likstr = lik; end
12 if ~ischar(likstr), likstr = func2str(likstr); end
13 if ~strcmp(likstr,'likGauss') % NOTE: no explicit call to likGauss
14     error('Exact inference only possible with Gaussian likelihood');
15 end
16
17 [n, D] = size(x);
18 K = feval(cov{:}, hyp.cov, x); % evaluate covariance matrix
19 m = feval(mean{:}, hyp.mean, x); % evaluate mean vector
20
21 sn2 = exp(2*hyp.lik); % noise variance of likGauss
22 if sn2<1e-6 % very tiny sn2 can lead to numerical trouble
23     L = chol(K+sn2*eye(n)); s1 = 1; % Cholesky factor of covariance with noise
24     pL = -solve_chol(L,eye(n)); % L = -inv(K+inv(sW^2))
25 else
26     L = chol(K/sn2+eye(n)); s1 = sn2; % Cholesky factor of B
27     pL = L; % L = chol(eye(n)+sW*sW'.*K)
28 end
29 alpha = solve_chol(L,y-m)/s1;
30
31 post.alpha = alpha; % return the posterior parameters
32 post.sW = ones(n,1)/sqrt(sn2); % sqrt of noise precision vector
33 post.L = pL;
34
35 if nargin>1 % do we want the marginal likelihood?
36     nlZ = (y-m)'*alpha/2 + sum(log(diag(L))) + n*log(2*pi*s1)/2; % -log marg lik
37     if nargin>2 % do we want derivatives?
38         dnlZ = hyp; % allocate space for derivatives
39         Q = solve_chol(L,eye(n))/s1 - alpha*alpha'; % precompute for convenience
40         for i = 1:numel(hyp.cov)

```

```

41     dn1Z.cov(i) = sum(sum(Q.*feval(cov{:}), hyp.cov, x, [], i))/2;
42 end
43 dn1Z.lik = sn2*trace(Q);
44 for i = 1:numel(hyp.mean),
45     dn1Z.mean(i) = -feval(mean{:}), hyp.mean, x, i)*alpha;
46 end
47 end
48 end

```

### 3.2 Laplace's Approximation

For differentiable likelihoods, Laplace's approximation, approximates the posterior by a Gaussian centered at its mode and matching its curvature `inf/infLaplace.m`.

More concretely, the mean of the posterior  $q(\mathbf{f}|\mathcal{D}) = \mathcal{N}(\mathbf{f}|\boldsymbol{\mu}, \mathbf{V})$  is – defining  $\ell_i(\mathbf{f}_i) = \ln p(\mathbf{y}_i|\mathbf{f}_i)$  and  $\ell(\mathbf{f}) = \sum_{i=1}^n \ell_i(\mathbf{f}_i)$  – given by

$$\boldsymbol{\mu} = \arg \min_{\mathbf{f}} \phi(\mathbf{f}), \text{ where } \phi(\mathbf{f}) = \frac{1}{2}(\mathbf{f} - \mathbf{m})^\top \mathbf{K}^{-1}(\mathbf{f} - \mathbf{m}) - \ell(\mathbf{f}) \stackrel{c}{=} -\ln[p(\mathbf{f})p(\mathbf{y}|\mathbf{f})], \quad (2)$$

which we abbreviate by  $\boldsymbol{\mu} \leftarrow \mathcal{L}(\ell)$ . The curvature  $\frac{\partial^2 \phi}{\partial \mathbf{f}^2} = \mathbf{K}^{-1} + \mathbf{W}$  with  $W_{ii} = -\frac{\partial^2}{\partial f_i^2} \ln p(\mathbf{y}_i|\mathbf{f}_i)$  serves as precision for the Gaussian posterior approximation  $\mathbf{V} = (\mathbf{K}^{-1} + \mathbf{W})^{-1}$  and the marginal likelihood  $Z = \int p(\mathbf{f})p(\mathbf{y}|\mathbf{f})d\mathbf{f}$  is approximated by  $Z \approx Z_{LA} = \int \tilde{\phi}(\mathbf{f})d\mathbf{f}$  where we use the 2nd order Taylor expansion at the mode  $\boldsymbol{\mu}$  given by  $\tilde{\phi}(\mathbf{f}) = \phi(\boldsymbol{\mu}) + \frac{1}{2}(\mathbf{f} - \boldsymbol{\mu})^\top \mathbf{V}^{-1}(\mathbf{f} - \boldsymbol{\mu}) \approx \phi(\mathbf{f})$ .

Laplace's approximation needs derivatives up to third order for the mode fitting procedure (Newton method)

$$d_k = \frac{\partial^k}{\partial \mathbf{f}^k} \log p(\mathbf{y}|\mathbf{f}), \quad k = 0, 1, 2, 3$$

and

$$d_k = \frac{\partial}{\partial \rho_i} \frac{\partial^k}{\partial \mathbf{f}^k} \log p(\mathbf{y}|\mathbf{f}), \quad k = 0, 1, 2$$

evaluated at the latent location  $\mathbf{f}$  and observed value  $\mathbf{y}$ . The likelihood calls (see section 4)

- `[d0, d1, d2, d3] = lik(hyp, y, f, [], 'infLaplace')`

and

- `[d0, d1, d2] = lik(hyp, y, f, [], 'infLaplace', i)`

return exactly these values.

### 3.3 Expectation Propagation

The basic idea of Expectation Propagation (EP) as implemented in `inf/infEP.m` is to replace the non-Gaussian likelihood terms  $p(\mathbf{y}_i|\mathbf{f}_i)$  by Gaussian functions  $t(\mathbf{f}_i; \nu_i, \tau_i) = \exp(\nu_i \mathbf{f}_i - \frac{1}{2} \tau_i \mathbf{f}_i^2)$  and to adjust the natural parameters  $\nu_i, \tau_i$  such that the following identity holds:

$$\frac{1}{Z_{t,i}} \int \mathbf{f}^k q_{-i}(\mathbf{f}) \cdot t(\mathbf{f}; \nu_i, \tau_i) d\mathbf{f} = \frac{1}{Z_{p,i}} \int \mathbf{f}^k q_{-i}(\mathbf{f}) \cdot p(\mathbf{y}_i|\mathbf{f}) d\mathbf{f}, \quad k = 1, 2$$

with the so-called cavity distributions  $q_{-i}(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mathbf{m}, \mathbf{K}) \prod_{j \neq i} t(\mathbf{f}_j; \nu_j, \tau_j) \propto \mathcal{N}(\mathbf{f}|\boldsymbol{\mu}, \mathbf{V})/t(\mathbf{f}_i; \nu_i, \tau_i)$  equal to the posterior divided by the  $i$ th Gaussian approximation function and the two normalisers

$Z_{t,i} = \int q_{-i}(f) \cdot t(f_i; \nu_i, \tau_i) df$  and  $Z_{p,i} = \int q_{-i}(f) \cdot p(y_i|f_i) df$ . The moment matching corresponds to minimising the following local KL-divergence

$$\nu_i, \tau_i = \arg \min_{\nu, \tau} \text{KL}[q_{-i}(f)p(y_i|f_i)/Z_{p,i} \| q_{-i}(f)t(f_i; \nu, \tau)/Z_{t,i}].$$

In order to apply the moment matching steps in a numerically safe way, EP requires the derivatives of the expectations w.r.t. the Gaussian mean parameter  $\mu$

$$d_k = \frac{\partial^k}{\partial \mu^k} \log \int p(y|f) \mathcal{N}(f|\mu, \sigma^2) df, \quad k = 0, 1, 2$$

and the  $i$ th likelihood hyperparameter  $\rho_i$

$$d = \frac{\partial}{\partial \rho_i} \log \int p(y|f) \mathcal{N}(f|\mu, \sigma^2) df$$

which can be obtained by the likelihood calls (see section 4)

- `[d0, d1, d2] = lik(hyp, y, mu, s2, 'infEP')`

and

- `d = lik(hyp, y, mu, s2, 'infEP', i).`

### 3.4 Kullback Leibler Divergence Minimisation

Another well known approach to approximate inference implemented `inf/infKL.m` in attempts to directly find the closest Gaussian  $q(f|\mathcal{D}) = \mathcal{N}(f|\mu, V)$  to the exact posterior  $p(f|\mathcal{D})$  w.r.t. to some proximity measure or equivalently to maximise a lower bound  $Z(\mu, V)$  to the marginal likelihood  $Z$  as described in Nickisch & Rasmussen *Approximations for Binary Gaussian Process Classification*, JMLR, 2008. In particular, one minimises  $\text{KL}(\mathcal{N}(f|\mu, V) \| p(f|\mathcal{D}))$  which amounts to minimising  $-\ln Z(\mu, V)$  as defined by:

$$\begin{aligned} -\ln Z &= -\ln \int p(f)p(y|f)df = -\ln \int q(f|\mathcal{D}) \frac{p(f)}{q(f|\mathcal{D})} p(y|f)df \\ &\stackrel{\text{Jensen}}{\leq} \int q(f|\mathcal{D}) \ln \frac{q(f|\mathcal{D})}{p(f)} df - \int q(f|\mathcal{D}) \ln p(y|f) df =: -\ln Z(\mu, V) \\ &= \text{KL}(\mathcal{N}(f|\mu, V) \| \mathcal{N}(f|\mathbf{m}, \mathbf{K})) - \sum_{i=1}^n \int \mathcal{N}(f_i|\mu_i, \nu_{ii}) \ln p(y_i|f_i) df_i, \nu_{ii} = [V]_{ii} \\ &= \frac{1}{2} \left( \text{tr}(\mathbf{V}\mathbf{K}^{-1}) - \ln |\mathbf{V}\mathbf{K}^{-1}| \right) + \frac{1}{2} (\mu - \mathbf{m})^\top \mathbf{K}^{-1} (\mu - \mathbf{m}) - \sum_{i=1}^n \ell^{\text{KL}}(\mu_i, \nu_{ii}) \end{aligned}$$

where  $\ell_v^{\text{KL}}(\mu_i) = \int \mathcal{N}(f_i|\mu_i, \nu_{ii}) \ell_i(f_i) df_i$  is the convolution of the log likelihood  $\ell_i$  with the Gaussian  $\mathcal{N}$  and  $\mathbf{v} = \text{dg}(\mathbf{V})$ . Equivalently, one can view  $\ell^{\text{KL}}$  as a smoothed version of  $\ell$  with univariate smoothing kernel  $\mathcal{N}$ .

From Challis & Barber *Concave Gaussian Variational Approximations for Inference in Large Scale Bayesian Linear Models*, AISTATS, 2011 we know that the mapping  $(\mu, L) \mapsto -\ln Z(\mu, L^\top L)$  is jointly convex whenever the likelihoods  $f_i \mapsto \mathbb{P}(y_i|f_i)$  are log concave. In particular, this implies that every  $(\mu_i, \sqrt{\nu_{ii}}) \mapsto \ell^{\text{KL}}(\mu_i, \nu_{ii})$  is jointly convex.

We use an optimisation algorithm similar to EP (section 3.3) where we minimise the local KL-divergence the other way round  $\mu_i, \sqrt{v_{ii}} = \arg \min_{\mu, \sigma} \text{KL}[\mathcal{N}(f|\mu_i, v_{ii}) \| q_{-i}(f)p(y_i|f_i)/Z_{p,i}]$ . This view was brought forward by Tom Minka *Convex Divergence measures and message passing*, MSR-TR, 2005. The KL minimisation constitutes a jointly convex 2d optimisation problem solved by `klmin` using a scaled Newton approach which is included as a sub function in `inf/infKL.m`. The smoothed likelihood  $\ell^{\text{KL}}(\mu_i, v_{ii})$  is implemented as a meta likelihood in `likKL`; it uses Gaussian-Hermite quadrature to compute the required integrals. Note that – as opposed to EP – Gaussian-Hermite quadrature is appropriate since we integrate against the  $\ln \mathbb{P}(y_i|f_i)$  (which can be well approximated by a polynomial) instead of  $\mathbb{P}(y_i|f_i)$  itself. The algorithm is – again unlike EP – provably convergent for log-concave likelihoods (e.g. `likGauss`, `likLaplace`, `likSech2`, `likLogistic`, `likPoisson`) since it can be regarded as coordinate descent with guaranteed decrease in the objective in every step. Due to the complex update computations, `infKL` can be quite slow although it has the same  $\mathcal{O}(n^3)$  asymptotic complexity as EP and Laplace.

### 3.5 Variational Bayes

One can drive the bounding even further by means of local quadratic lower bounds to the log likelihood  $\ell(f) = \ln p(y|f)$ . Suppose that we use a super-Gaussian likelihood  $p(y|f)$  i.e. likelihoods that can be lower bounded by Gaussians of any width  $w$  (e.g. `likLaplace`, `likT`, `likLogistic`, `likSech2`). Formally, that means that there are  $b, z \in \mathbb{R}$  such that

$$\rho(f) = \ln p(y|f - z) - bf$$

is symmetric and  $\sqrt{f} \mapsto \rho(f)$  is a convex function for all  $f \geq 0$ . As a result, we obtain the following exact representation of the likelihood

$$\ell(f) = \ln p(y|f) = \max_{w>0} \left( (b + wz)f - \frac{wf^2}{2} - \frac{1}{2}h(\gamma) \right),$$

which can be derived by convex duality and assuming the likelihoods to be super-Gaussian. Details can be found in papers by Palmer et al. *Variational EM Algorithms for Non-Gaussian Latent Variable Models*, NIPS, 2006 and Nickisch & Seeger *Convex Variational Bayesian Inference for Large Scale Generalized Linear Models*, ICML, 2009.

The bottom line is that we can treat the variational bounding as a sequence of Laplace approximations with the “variational Bayes” log likelihood

$$\ell^{\text{VB}}(f_i) = \ell(g_i) + b_i(f_i - g_i), \mathbf{g} = \text{sgn}(\mathbf{f} - \mathbf{z}) \odot \sqrt{(\mathbf{f} - \mathbf{z})^2 + \mathbf{v}} + \mathbf{z}$$

instead of the usual likelihood  $\ell(f_i) = \ln p(y_i|f_i)$  i.e. we solve  $\mu \leftarrow \mathcal{L}(\ell_v^{\text{VB}})$  instead of  $\mu \leftarrow \mathcal{L}(\ell)$ . See section 3.2. In the code of `inf/infVB.m`, the likelihood is implemented in the function `likVB`.

At the end, the optimal value of  $\mathbf{W}$  can be obtained analytically via  $w_i = |b_i - \ell'(g_i)|/|g_i - z_i|$ .

For the minimisation in `inf/infVB.m`, we use a provably convergent double loop algorithm, where in the inner loop a nonlinear least squares problem (convex for log-concave likelihoods) is solved using `inf/infLaplace.m` such that  $\mu \leftarrow \mathcal{L}(\ell_v^{\text{VB}})$  and in the outer loop, we compute  $\mathbf{v} \leftarrow \text{dg}((\mathbf{K}^{-1} + \mathbf{W})^{-1})$ . The only requirement to the likelihood function is that it returns the values  $\mathbf{z}$  and  $\mathbf{b}$  required by the bound which are delivered by the call (see section 4)

- `[b,z] = lik(hyp, y, [], ga, 'infVB')`

The negative marginal likelihood upper bound  $-\ln Z_{\text{VB}}$  is obtained by integrating the prior times the exact representation of the likelihood

$$p(y|f) = \max_{\gamma>0} q(y|f, \gamma), \quad q(y|f, \gamma) = \mathcal{N}(f|\mathbf{v}, \gamma) \exp \left( -\frac{h(\gamma)}{2} - \frac{v^2}{2\gamma} \right) \sqrt{2\pi\gamma}, \quad \gamma = \frac{1}{w}, \quad \mathbf{v} = \mathbf{b}\gamma + \mathbf{z}$$

w.r.t. the latent variables  $\mathbf{f}$  yielding

$$\begin{aligned} -\ln Z_{VB} &= -\ln \int \mathcal{N}(\mathbf{f}|\mathbf{m}, \mathbf{K}) \prod_{i=1}^n q_i(y_i|\mathbf{f}_i, \gamma_i) d\mathbf{f} \\ &= -\ln \mathcal{N}(\mathbf{m}|\mathbf{v}, \mathbf{K} + \mathbf{\Gamma}) + \frac{1}{2} \left( h(\boldsymbol{\gamma}) - \mathbf{w}^\top \mathbf{v}^2 - \mathbf{1}^\top \ln 2\pi\boldsymbol{\gamma} \right). \end{aligned}$$

### 3.6 FITC Approximations

One of the main problems with GP models is the high computational load for inference computations. In a setting with  $n$  training points  $\mathbf{x}$ , exact inference with Gaussian likelihood requires  $O(n^3)$  effort; approximations like Laplace of EP consist of a sequence of  $O(n^3)$  operations.

There is a line of research with the goal to alleviate this burden by using approximate covariance functions  $\tilde{\mathbf{k}}$  instead of  $\mathbf{k}$ . A review is given by Candela and Rasmussen *A Unifying View of Sparse Approximate Gaussian Process Regression*, JMLR, 2005. One basic idea in those approximations is to work with a set of  $m$  inducing inputs  $\mathbf{u}$  with a reduced computational load of  $O(nm^2)$ . In the following, we will provide a rough idea of the FITC approximation used in the toolbox. Let  $\mathbf{K}$  denote the  $n \times n$  covariance matrix between the training points  $\mathbf{x}$ ,  $\mathbf{K}_u$  the  $m \times n$  covariance matrix between the  $n$  training points and the  $m$  inducing points, and  $\mathbf{K}_{uu}$  the  $m \times m$  covariance matrix between the  $m$  inducing points. The FITC approximation to the covariance is given by

$$\mathbf{K} \approx \tilde{\mathbf{K}} = \mathbf{Q} + \mathbf{G}, \mathbf{G} = \text{diag}(\mathbf{g}), \mathbf{g} = \text{diag}(\mathbf{K} - \mathbf{Q}), \mathbf{Q} = \mathbf{K}_u^\top \mathbf{Q}_{uu}^{-1} \mathbf{K}_u, \mathbf{Q}_{uu} = \mathbf{K}_{uu} + \sigma_{n_u}^2 \mathbf{I},$$

where  $\sigma_{n_u}$  is the noise from the inducing inputs. Note that  $\tilde{\mathbf{K}}$  and  $\mathbf{K}$  have the same diagonal elements  $\text{diag}(\tilde{\mathbf{K}}) = \text{diag}(\mathbf{K})$ ; all off-diagonal elements are the same as for  $\mathbf{Q}$ . The toolbox offers FITC versions for regression with Gaussian likelihood `inf/infFITC.m`, as well as for Laplace's approximation `inf/infFITCLaplace.m` and expectation propagation `inf/infFITCEP.m`.

## 4 Likelihood Functions

A likelihood function  $p_\rho(y|f)$  (with hyperparameters  $\rho$ ) is a conditional density  $\int p_\rho(y|f)dy = 1$  defined for scalar latent function values  $f$  and outputs  $y$ . In the GPML toolbox, we use iid. likelihoods  $p_\rho(y|f) = \prod_{i=1}^n p_\rho(y_i|f_i)$ . The approximate inference engine does not explicitly distinguish between classification and regression likelihoods: it is fully generic in the likelihood allowing to use a single code in the inference step.

Likelihood functionality is needed both during inference and while predicting.

### 4.1 Prediction

A prediction at  $\mathbf{x}_*$  conditioned on the data  $\mathcal{D} = (X, y)$  (as implemented in `gp.m`) consists of the predictive mean  $\mu_{y_*}$  and variance  $\sigma_{y_*}^2$  which are computed from the the latent marginal moments  $\mu_{f_*}, \sigma_{f_*}^2$  i.e. the Gaussian marginal approximation  $\mathcal{N}(f_*|\mu_{f_*}, \sigma_{f_*}^2)$  via

$$p(y_*|\mathcal{D}, \mathbf{x}_*) = \int p(y_*|f_*)p(f_*|\mathcal{D}, \mathbf{x}_*)df_* \approx \int p(y_*|f_*)\mathcal{N}(f_*|\mu_{f_*}, \sigma_{f_*}^2)df_*. \quad (3)$$

The moments are given by  $\mu_{y_*} = \int y_*p(y_*|\mathcal{D}, \mathbf{x}_*)dy_*$  and  $\sigma_{y_*}^2 = \int (y_* - \mu_{y_*})^2p(y_*|\mathcal{D}, \mathbf{x}_*)dy_*$ . The likelihood call

$$\bullet \text{ [lp, ymu, ys2] = lik(hyp, [], fmu, fs2)}$$

does exactly this. Evaluation of the logarithm of  $p_{y_*} = p(y_*|\mathcal{D}, \mathbf{x}_*)$  for values  $y_*$  can be done via

$$\bullet \text{ [lp, ymu, ys2] = lik(hyp, y, fmu, fs2)}$$

where `lp` contains the number  $\ln p_{y_*}$ .

Using the moments of the likelihood  $\mu(f_*) = \int y_*p(y_*|f_*)dy_*$  and  $\sigma^2(f_*) = \int (y_* - \mu(f_*))^2p(y_*|f_*)dy_*$  we obtain for the predictive moments the following (exact) expressions

$$\begin{aligned} \mu_{y_*} &= \int \mu(f_*)p(f_*|\mathcal{D}, \mathbf{x}_*)df_*, \text{ and} \\ \sigma_{y_*}^2 &= \int \left[ \sigma^2(f_*) + (\mu(f_*) - \mu_{y_*})^2 \right] p(f_*|\mathcal{D}, \mathbf{x}_*)df_*. \end{aligned}$$

1. The binary case is simple since  $y_* \in \{-1, +1\}$  and  $1 = p_{y_*} + p_{-y_*}$ . Using  $\pi_* = p_{+1}$ , we find

$$\begin{aligned} p_{y_*} &= \begin{cases} \pi_* & y_* = +1 \\ 1 - \pi_* & y_* = -1 \end{cases} \\ \mu_{y_*} &= \sum_{y_*=\pm 1} y_*p(y_*|\mathcal{D}, \mathbf{x}_*) = 2 \cdot \pi_* - 1 \in [-1, 1], \text{ and} \\ \sigma_{y_*}^2 &= \sum_{y_*=\pm 1} (y_* - \mu_{y_*})^2p(y_*|\mathcal{D}, \mathbf{x}_*) = 4 \cdot \pi_*(1 - \pi_*) \in [0, 1]. \end{aligned}$$

2. The continuous case for homoscedastic likelihoods depending on  $r_* = y_* - f_*$  only and having noise variance  $\sigma^2(f_*) = \sigma_n^2$  is also simple since the identity  $p(y_*|f_*) = p(y_* - f_*|0)$  allows to substitute  $y_* \leftarrow y_* + f_*$  yielding  $\mu(f_*) = f_* + \int y_*p(y_*|0)dy_*$  and assuming  $\int y_*p(y_*|0)dy_* = 0$  we arrive at

$$\begin{aligned} \mu_{y_*} &= \mu_{f_*}, \text{ and} \\ \sigma_{y_*}^2 &= \sigma_{f_*}^2 + \sigma_n^2. \end{aligned}$$

3. The generalised linear model (GLM) case is also feasible. Evaluation of the predictive distribution is done by quadrature

$$p_{y_*} = \int p(y_*|f_*)p(f_*|\mathcal{D}, \mathbf{x}_*)df_* \approx \int p(y_*|f_*)\mathcal{N}(f_*|\mu_{f_*}, \sigma_{f_*}^2)df_*.$$

For GLMs the mean is given by  $\mu(f_*) = g(f_*)$  and the variance is usually given by a simple function of the mean  $\sigma^2(f_*) = v(g(f_*))$ , hence we use Gaussian-Hermite quadrature with  $\mathcal{N}(f_*|\mu_{f_*}, \sigma_{f_*}^2) \approx p(f_*|\mathcal{D}, \mathbf{x}_*)$  to compute

$$\begin{aligned}\mu_{y_*} &= \int g(f_*)p(f_*|\mathcal{D}, \mathbf{x}_*)df_*, \text{ and} \\ \sigma_{y_*}^2 &= \int \left[ v(g(f_*)) + (g(f_*) - \mu_{y_*})^2 \right] p(f_*|\mathcal{D}, \mathbf{x}_*)df_* \neq v(\mu_{y_*}).\end{aligned}$$

4. Finally the warped Gaussian likelihood predictive distribution with strictly monotonically increasing warping function  $g$  is given by the expression

$$p(y_*|\mathcal{D}, \mathbf{x}_*) = g'(y_*)\mathcal{N}\left(g(y_*)|\mu_{f_*}, \sigma_n^2 + \sigma_{f_*}^2\right)$$

so that the predictive moments can be computed by Gaussian-Hermite quadrature.

In the following, we will detail how and which likelihood functions are implemented in the GPML toolbox. Further, we will mention dependencies between likelihoods and inference methods and provide some analytical expressions in addition to some likelihood implementations.

## 4.2 Interface

The likelihoods are in fact the most challenging object in our implementation. Different inference algorithms require different aspects of the likelihood to be computed, therefore the interface is rather involved as detailed below.

```
16 <likFunctions.m 16>≡
1 % likelihood functions are provided to be used by the gp.m function:
2 %
3 %   likErf           (Error function, classification, probit regression)
4 %   likLogistic      (Logistic,           classification, logit regression)
5 %   likUni           (Uniform likelihood, classification)
6 %
7 %   likGauss         (Gaussian, regression)
8 %   likGaussWarp     (Warped Gaussian, regression)
9 %   likGumbel        (Gumbel likelihood for extremal values)
10 %  likLaplace        (Laplacian or double exponential, regression)
11 %  likSech2          (Sech-square, regression)
12 %  likT              (Student's t, regression)
13 %
14 %  likPoisson        (Poisson regression, count data)
15 %  likGamma          (Nonnegative regression, positive data)
16 %  likExp            (Nonnegative regression, positive data)
17 %  likInvGauss       (Nonnegative regression, positive data)
18 %  likBeta           (Beta regression, interval data)
19 %
20 %  likMix            (Mixture of individual covariance functions)
21 %
22 % The likelihood functions have three possible modes, the mode being selected
```



```

23 % as follows (where "lik" stands for any likelihood function in "lik/lik*.m"):
24 %
25 % 1) With one or no input arguments:           [REPORT NUMBER OF HYPERPARAMETERS]
26 %
27 %     s = lik OR s = lik(hyp)
28 %
29 % The likelihood function returns a string telling how many hyperparameters it
30 % expects, using the convention that "D" is the dimension of the input space.
31 % For example, calling "likLogistic" returns the string '0'.
32 %
33 %
34 % 2) With three or four input arguments:           [PREDICTION MODE]
35 %
36 %     lp = lik(hyp, y, mu) OR [lp, ymu, ys2] = lik(hyp, y, mu, s2)
37 %
38 % This allows to evaluate the predictive distribution. Let  $p(y_*|f_*)$  be the
39 % likelihood of a test point and  $N(f_*|\mu, s2)$  an approximation to the posterior
40 % marginal  $p(f_*|x_*, x, y)$  as returned by an inference method. The predictive
41 % distribution  $p(y_*|x_*, x, y)$  is approximated by.
42 %      $q(y_*) = \int N(f_*|\mu, s2) p(y_*|f_*) df_*$ 
43 %
44 %     lp = log( q(y) ) for a particular value of y, if s2 is [] or 0, this
45 %                               corresponds to log( p(y|\mu) )
46 %     ymu and ys2               the mean and variance of the predictive marginal q(y)
47 %                               note that these two numbers do not depend on a particular
48 %                               value of y
49 % All vectors have the same size.
50 %
51 %
52 % 3) With five or six input arguments, the fifth being a string [INFERENCE MODE]
53 %
54 % [varargout] = lik(hyp, y, mu, s2, inf) OR
55 % [varargout] = lik(hyp, y, mu, s2, inf, i)
56 %
57 % There are three cases for inf, namely a) infLaplace, b) infEP and c) infVB.
58 % The last input i, refers to derivatives w.r.t. the ith hyperparameter.
59 %
60 % a1) [lp, dlp, d2lp, d3lp] = lik(hyp, y, f, [], 'infLaplace')
61 % lp, dlp, d2lp and d3lp correspond to derivatives of the log likelihood
62 % log(p(y|f)) w.r.t. to the latent location f.
63 %     lp = log( p(y|f) )
64 %     dlp = d log( p(y|f) ) / df
65 %     d2lp = d^2 log( p(y|f) ) / df^2
66 %     d3lp = d^3 log( p(y|f) ) / df^3
67 %
68 % a2) [lp_dhyp, dlp_dhyp, d2lp_dhyp] = lik(hyp, y, f, [], 'infLaplace', i)
69 % returns derivatives w.r.t. to the ith hyperparameter
70 %     lp_dhyp = d log( p(y|f) ) / ( dhyp_i)
71 %     dlp_dhyp = d^2 log( p(y|f) ) / (df dhyp_i)
72 %     d2lp_dhyp = d^3 log( p(y|f) ) / (df^2 dhyp_i)
73 %
74 %
75 % b1) [lZ, d1Z, d21Z] = lik(hyp, y, mu, s2, 'infEP')
76 % let  $Z = \int p(y|f) N(f|\mu, s2) df$  then
77 %     lZ = log(Z)
78 %     d1Z = d log(Z) / dmu
79 %     d21Z = d^2 log(Z) / dmu^2
80 %

```

```

81 % b2) [dlZhyp] = lik(hyp, y, mu, s2, 'infEP', i)
82 % returns derivatives w.r.t. to the ith hyperparameter
83 % dlZhyp = d log(Z) / dhyp_i
84 %
85 %
86 % c1) [b,z] = lik(hyp, y, [], ga, 'infVB')
87 % ga is the variance of a Gaussian lower bound to the likelihood p(y|f).
88 % p(y|f) \ge exp( b*(f+z) - (f+z).^2/(2*ga) - h(ga)/2 ) \propto N(f|b*ga-z,ga)
89 % The function returns the linear part b and z.
90 %
91 % Cumulative likelihoods are designed for binary classification. Therefore, they
92 % only look at the sign of the targets y; zero values are treated as +1.
93 %
94 % Some examples for valid likelihood functions:
95 %     lik = @likLogistic;
96 %     lik = {'likMix','likUni','likErf'}
97 %     lik = {@likPoisson,'logistic'};
98 %
99 % See the help for the individual likelihood for the computations specific to
100 % each likelihood function.
101 %
102 % (gpml copyright 5a)

```

### 4.3 Implemented Likelihood Functions

The following table enumerates all (currently) implemented likelihood functions that can be found at `lik/lik<NAME>.m` and their respective set of hyperparameters  $\rho$ .

lik<NAME>	regression $y_i \in \mathbb{R}$	$p_\rho(y_i f_i) =$	$\rho =$
Gauss	Gaussian	$N(y_i f_i, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_i-f_i)^2}{2\sigma^2}\right)$	$\{\ln \sigma\}$
GaussWarp	Warped Gaussian	$N(g_\theta(y_i) f_i, \sigma^2) g'_\theta(y_i)$	$\{\theta_1, \dots, \theta_{n_g}, \ln \sigma\}$
Gumbel	Gumbel	$\frac{\pi}{\sigma\sqrt{6}} \exp(-z_i - e^{-z_i}), z_i = \gamma + \frac{s \cdot \pi(y_i - f_i)}{\sigma\sqrt{6}},  s  = 1$	$\{\ln \sigma\}$
Sech2	Sech-squared	$\frac{\tau}{2 \cosh^2(\tau(y_i - f_i))}, \tau = \frac{\pi}{2\sigma\sqrt{3}}$	$\{\ln \sigma\}$
Laplace	Laplacian	$\frac{1}{2b} \exp\left(-\frac{ y_i - f_i }{b}\right), b = \frac{\sigma}{\sqrt{2}}$	$\{\ln \sigma\}$
T	Student's t	$\frac{\Gamma(\frac{\nu+1}{2})}{\Gamma(\frac{\nu}{2})} \frac{1}{\sqrt{\nu\pi}\sigma} \left(1 + \frac{(y_i - f_i)^2}{\nu\sigma^2}\right)^{-\frac{\nu+1}{2}}$	$\{\ln(\nu - 1), \ln \sigma\}$
lik<NAME>	classification $y_i \in \{\pm 1\}$	$p_\rho(y_i f_i) =$	$\rho =$
Erf	Error function	$\int_{-\infty}^{y_i f_i} N(t) dt$	$\emptyset$
Logistic	Logistic function	$\frac{1}{1 + \exp(-y_i f_i)}$	$\emptyset$
Uni	Label noise	$\frac{1}{2}$	$\emptyset$
lik<NAME>	count data $y_i \in \mathbb{N}$	$p_\rho(y_i f_i) =$	$\rho =$
Poisson	Poisson	$\mu^{y_i} \cdot \frac{e^{-\mu}}{y_i!}, \mu = e^f \text{ or } \mu = \log(1 + e^f)$	$\emptyset$
lik<NAME>	nonnegative data $y_i \in \mathbb{R}_+ \setminus \{0\}$	$p_\rho(y_i f_i) =$	$\rho =$
Weibull	Weibull, $\gamma_1 = \Gamma(1 + 1/\kappa)$	$\kappa \gamma_1 / \mu (y \gamma_1 / \mu)^{\kappa-1} \exp(-(y \gamma_1 / \mu)^\kappa)$	$\{\ln \kappa\}$
Gamma	Gamma	$\frac{\alpha^\alpha \gamma^{\alpha-1}}{\Gamma(\alpha)} \mu^{-\alpha} \exp\left(-\frac{y \alpha}{\mu}\right)$	$\{\ln \alpha\}$
Exp	Exponential	$\mu^{-1} \exp\left(-\frac{y}{\mu}\right)$	$\emptyset$
InvGauss	Inverse Gaussian	$\sqrt{\frac{\lambda}{2\pi y^3}} \exp\left(-\frac{\lambda(y-\mu)^2}{2\mu^2 y}\right)$	$\{\ln \lambda\}$
lik<NAME>	interval data $y_i \in [0, 1]$	$p_\rho(y_i f_i) =$	$\rho =$
Beta	Beta	$\frac{\Gamma(\Phi)}{\Gamma(\mu\Phi)\Gamma((1-\mu)\Phi)} y^{\mu\Phi-1} (1-y)^{(1-\mu)\Phi-1}$	$\{\ln \Phi\}$
Composite likelihood functions $[p_1(y_i f_i), p_1(y_i f_i), \dots] \mapsto p_\rho(y_i f_i)$			
Mix	Mixture	$\sum_j \alpha_j p_j(y_i f_i)$	$\{\ln \alpha_1, \ln \alpha_2, \dots\}$

## 4.4 Usage of Implemented Likelihood Functions

Some code examples taken from `doc/usageLik.m` illustrate how to use simple and composite likelihood functions to specify a GP model.

Syntactically, a likelihood function `lk` is defined by

```
lk := 'func' | @func // simple
```

```
lk := {lk} | {param, lk} | {lk, {lk, ..., lk}} // composite
```

i.e., it is either a string containing the name of a likelihood function, a pointer to a likelihood function or one of the former in combination with a cell array of likelihood functions and an additional list of parameters.

```
19 <doc/usageLik.m 19>≡
1 % demonstrate usage of likelihood functions
2 %
3 % See also likFunctions.m.
4 %
5 <gpml copyright 5a>
6 clear all, close all
7 n = 5; f = randn(n,1);          % create random latent function values
8
9 % set up simple classification likelihood functions
10 yc = sign(f);
11 lc0 = {'likErf'};      hypc0 = [];    % no hyperparameters are needed
12 lc1 = {@likLogistic}; hypc1 = [];    % also function handles are OK
13 lc2 = {'likUni'};      hypc2 = [];
14 lc3 = {'likMix',{'likUni',@likErf}}; hypc3 = log([1;2]); %mixture
15
16 % set up simple regression likelihood functions
17 yr = f + randn(n,1)/20;
18 sn = 0.1;                                % noise standard deviation
19 lr0 = {'likGauss'};    hypr0 = log(sn);
20 lr1 = {'likLaplace'}; hypr1 = log(sn);
21 lr2 = {'likSech2'};   hypr2 = log(sn);
22 nu = 4;                                % number of degrees of freedom
23 lr3 = {'likT'};       hypr3 = [log(nu-1); log(sn)];
24 lr4 = {'likMix',{lr0,lr1}}; hypr4 = [log([1;2]);hypr0;hypr1];
25
26 a = 1; % set up warped Gaussian with g(y) = y + a*sign(y).*y.^2
27 lr5 = {'likGaussWarp',['poly2']}; hypr5 = log([a;sn]);
28 lr6 = {'likGumbel','+'}; hypr6 = log(sn);
29
30 % set up Poisson regression
31 yp = fix(abs(f)) + 1;
32 lp0 = {@likPoisson,'logistic'}; hypp0 = [];
33 lp1 = {@likPoisson,'exp'};      hypp1 = [];
34
35 % set up other GLM likelihoods for positive or interval regression
36 lg1 = {@likGamma,'logistic'}; al = 2;    hyp.lik = log(al);
37 lg2 = {@likInvGauss,'exp'};   lam = 1.1; hyp.lik = log(lam);
38 lg3 = {@likBeta,'expexp'};    phi = 2.1; hyp.lik = log(phi);
39 lg4 = {@likBeta,'logit'};     phi = 4.7; hyp.lik = log(phi);
40
41 % 0) specify the likelihood function
42 lik = lc0; hyp = hypc0; y = yc;
43 % lik = lr4; hyp = hypr4; y = yr;
44 % lik = lp1; hyp = hypp1; y = yp;
```

```

45
46 % 1) query the number of parameters
47 feval(lik{:})
48
49 % 2) evaluate the likelihood function on f
50 exp(feval(lik{:},hyp,y,f))
51
52 % 3a) evaluate derivatives of the likelihood
53 [lp,dlp,d2lp,d3lp] = feval(lik{:}, hyp, y, f, [], 'infLaplace');
54
55 % 3b) compute Gaussian integrals w.r.t. likelihood
56 mu = f; s2 = rand(n,1);
57 [lZ,d1Z,d2lZ] = feval(lik{:}, hyp, y, mu, s2, 'infEP');
58
59 % 3c) obtain lower bound on likelihood
60 ga = rand(n,1);
61 [b,z] = feval(lik{:}, hyp, y, [], ga, 'infVB');

```

## 4.5 Compatibility Between Likelihoods and Inference Methods

The following table lists all possible combinations of likelihood function and inference methods.

Likelihood \ Inference	Exact	EP	Laplace	VB	KL	MCMC	LOO	Type, Output Domain	Alternative Names
	FITC	FITC-EP	FITC-Laplace						
Gaussian	✓	✓	✓	✓	✓	✓	✓	regression, $\mathbb{R}$	
Warped Gaussian		✓	✓	✓	✓	✓	✓	regression, $\mathbb{R}$	
Gumbel			✓		✓	✓	✓	regression, $\mathbb{R}$	
Sech-squared		✓	✓	✓	✓	✓	✓	regression, $\mathbb{R}$	logistic distribution
Laplacian		✓	✓	✓	✓	✓	✓	regression, $\mathbb{R}$	double exponential
Student's t			✓	✓	✓	✓	✓	regression, $\mathbb{R}$	
Mixture		✓	✓		✓	✓	✓		mixing meta likelihood
Error function		✓	✓		✓	✓	✓	classification, $\{\pm 1\}$	probit regression
Logistic function		✓	✓	✓	✓	✓	✓	classification, $\{\pm 1\}$	logit regression
Uniform		✓	✓	✓	✓	✓	✓	classification, $\{\pm 1\}$	label noise
Weibull			✓			✓	✓	positive data, $\mathbb{R}_+ \setminus \{0\}$	nonnegative regression
Gamma			✓			✓	✓	positive data, $\mathbb{R}_+ \setminus \{0\}$	nonnegative regression
Exp			✓			✓	✓	positive data, $\mathbb{R}_+ \setminus \{0\}$	nonnegative regression
Inverse Gaussian			✓			✓	✓	positive data, $\mathbb{R}_+ \setminus \{0\}$	nonnegative regression
Poisson		(✓)*	✓		✓	✓	✓	count data, $\mathbb{N}$	Poisson regression
Beta			✓			✓	✓	interval data, $[0, 1]$	beta regression

(✓)\* EP might not converge in some cases since quadrature is used.

Exact inference is only tractable for Gaussian likelihoods. Expectation propagation together with Student's t likelihood is inherently unstable due to non-log-concavity. Laplace's approximation for Laplace likelihoods is not sensible because at the mode the curvature and the gradient is undefined due to the non-differentiable peak of the Laplace distribution. Special care has been taken for the non-convex optimisation problem imposed by the combination Student's t likelihood and Laplace's approximation.

## 4.6 Gaussian Likelihood

The Gaussian likelihood is the simplest likelihood because the posterior distribution is not only Gaussian but can be computed analytically. In principle, the Gaussian likelihood would only be operated in conjunction with the exact inference method but we chose to provide compatibility with all other inference algorithms as well because it enables code testing and allows to switch between

different regression likelihoods very easily.

```

21a <lik/likGauss.m 21a>≡
1 function [varargout] = likGauss(hyp, y, mu, s2, inf, i)
2
3 % likGauss - Gaussian likelihood function for regression. The expression for the
4 % likelihood is
5 %   likGauss(t) = exp(-(t-y)^2/2*sn^2) / sqrt(2*pi*sn^2),
6 % where y is the mean and sn is the standard deviation.
7 %
8 % The hyperparameters are:
9 %
10 % hyp = [ log(sn) ]
11 %
12 % Several modes are provided, for computing likelihoods, derivatives and moments
13 % respectively, see likFunctions.m for the details. In general, care is taken
14 % to avoid numerical issues when the arguments are extreme.
15 %
16 <gpml copyright 5a>
17 %
18 % See also LIKFUNCTIONS.M.
19
20 if nargin<3, varargout = {'1'}; return; end % report number of hyperparameters
21
22 sn2 = exp(2*hyp);
23
24 if nargin<5 % prediction mode if inf is not present
25   <Prediction with Gaussian likelihood 21b>
26 else
27   switch inf
28     case 'infLaplace'
29       <Laplace's method with Gaussian likelihood 22a>
30     case 'infEP'
31       <EP inference with Gaussian likelihood 22b>
32     case 'infVB'
33       <Variational Bayes inference with Gaussian likelihood 22c>
34   end
35 end
21b <Prediction with Gaussian likelihood 21b>≡ (21a)
1 if numel(y)==0, y = zeros(size(mu)); end
2 s2zero = 1; if nargin>3, if norm(s2)>0, s2zero = 0; end, end % s2==0 ?
3 if s2zero % log probability
4   lp = -(y-mu).^2./sn2/2-log(2*pi*sn2)/2; s2 = 0;
5 else
6   lp = likGauss(hyp, y, mu, s2, 'infEP'); % prediction
7 end
8 ymu = {}; ys2 = {};
9 if nargin>1
10   ymu = mu; % first y moment
11   if nargin>2
12     ys2 = s2 + sn2; % second y moment
13   end
14 end
15 varargout = {lp,ymu,ys2};

```

The Gaussian likelihood function has a single hyperparameter  $\rho$ , the log of the noise standard deviation  $\sigma_n$ .

### 4.6.1 Exact Inference

Exact inference doesn't require any specific likelihood related code; all computations are done directly by the inference method, section 3.1.

### 4.6.2 Laplace's Approximation

```
22a  <Laplace's method with Gaussian likelihood 22a>≡ (21a)
1  if nargin<6 % no derivative mode
2  if numel(y)==0, y=0; end
3  ymmu = y-mu; dlp = {}; d2lp = {}; d3lp = {};
4  lp = -ymmu.^2/(2*sn2) - log(2*pi*sn2)/2;
5  if nargin>1
6      dlp = ymmu/sn2; % dlp, derivative of log likelihood
7      if nargin>2 % d2lp, 2nd derivative of log likelihood
8          d2lp = -ones(size(ymmu))/sn2;
9          if nargin>3 % d3lp, 3rd derivative of log likelihood
10             d3lp = zeros(size(ymmu));
11         end
12     end
13 end
14 varargout = {lp,dlp,d2lp,d3lp};
15 else % derivative mode
16     lp_dhyp = (y-mu).^2/sn2 - 1; % derivative of log likelihood w.r.t. hypers
17     dlp_dhyp = 2*(mu-y)/sn2; % first derivative,
18     d2lp_dhyp = 2*ones(size(mu))/sn2; % and also of the second mu derivative
19     varargout = {lp_dhyp,dlp_dhyp,d2lp_dhyp};
20 end
```

### 4.6.3 Expectation Propagation

```
22b  <EP inference with Gaussian likelihood 22b>≡ (21a)
1  if nargin<6 % no derivative mode
2  lZ = -(y-mu).^2/(sn2+s2)/2 - log(2*pi*(sn2+s2))/2; % log part function
3  d1Z = {}; d21Z = {};
4  if nargin>1
5      d1Z = (y-mu)./(sn2+s2); % 1st derivative w.r.t. mean
6      if nargin>2
7          d21Z = -1./(sn2+s2); % 2nd derivative w.r.t. mean
8      end
9  end
10 varargout = {lZ,d1Z,d21Z};
11 else % derivative mode
12     d1Zhyp = ((y-mu).^2/(sn2+s2)-1) ./ (1+s2./sn2); % deriv. w.r.t. hyp.lik
13     varargout = {d1Zhyp};
14 end
```

### 4.6.4 Variational Bayes

```
22c  <Variational Bayes inference with Gaussian likelihood 22c>≡ (21a)
1  % variational lower site bound
2  % t(s) = exp(-(y-s)^2/2sn2)/sqrt(2*pi*sn2)
3  % the bound has the form: (b+z/ga)*f - f.^2/(2*ga) - h(ga)/2
4  n = numel(s2); b = zeros(n,1); y = y.*ones(n,1); z = y;
5  varargout = {b,z};
```

## 4.7 Warped Gaussian Likelihood

Starting from the likelihood  $p(y|f)$  we are sometimes facing the situation where the data  $y \in \mathcal{Y} \subseteq \mathbb{R}$  is not distributed according to  $p(y|f)$  but some nonlinear transformation of the data  $g(y) = z$  so that  $z \sim p(z|f)$ . Here, the warping function  $g : \mathcal{Y} \rightarrow \mathbb{R}$  needs to be strictly monotonically increasing i.e.  $g'(y) > 0$ . Formally, we start from the fact that  $p(z|f)$  integrates to one and use the derivative  $dz = g'(y)dy$  to substitute

$$\int p(z|f)dz = 1 = \int p_g(y|f)dy, \quad p_g(y|f) = p(g(y)|f)g'(y)$$

where we have defined the log warped likelihood  $\ln p_g(y|f) = \ln p(g(y)|f) + \ln g'(y)$ . The interesting bit is that approximate inference methods such as `infExact`, `infLaplace`, `infEP`, `infVB`, `infKL` remain fully feasible only prediction and derivatives become more involved. The usual GP inference is recovered by using the identity warping function  $g : y \mapsto y$ . The construction works in principle for any likelihood but our implementation in `likGaussWarp` is limited to the Gaussian likelihood.

### Hyperparameter derivatives

Hyperparameter derivatives for `infLaplace` are obtained as follows

$$\begin{aligned} \frac{\partial}{\partial \theta} \ln \frac{\partial^k}{\partial f^k} p_g(y|f) &= \frac{\partial}{\partial \theta} \ln \frac{\partial^k}{\partial f^k} p(g(y)|f) + \frac{\partial}{\partial \theta} \frac{\partial^k}{\partial f^k} \ln g'(y), \quad k = 0, 1, 2 \\ &= -\frac{\partial^{k+1}}{\partial f^{k+1}} \ln p(g(y)|f) \frac{\partial}{\partial \theta} g(y) + \frac{\partial}{\partial \theta} \frac{\partial^k}{\partial f^k} \ln g'(y). \end{aligned}$$

Similarly for `infEP` the derivatives are given by

$$\begin{aligned} \frac{\partial}{\partial \theta} \ln \int p_g(y|f) \mathcal{N}(f|\mu, \sigma^2) df &= \frac{\partial}{\partial \theta} \ln \int p(g(y)|f) \mathcal{N}(f|\mu, \sigma^2) df + \frac{\partial}{\partial \theta} \ln g'(y) \\ &= -\frac{\partial}{\partial \mu} \ln \int p(g(y)|f) \mathcal{N}(f|\mu, \sigma^2) df \frac{\partial}{\partial \theta} g(y) + \frac{\partial}{\partial \theta} \ln g'(y). \end{aligned}$$

This trick above works for any homoscedastic likelihood where  $p(y|f) = p(y + \beta|f + \beta)$  such as `likGauss`, `likLaplace`, `likSech2` and `likT`.

### Predictive moments

As detailed in 4, the predictive distribution is – for Gaussian likelihood – given by

$$\begin{aligned} p(z_*|\mathcal{D}, \mathbf{x}_*) &= \int p(z_*|f_*) p(f_*|\mathcal{D}, \mathbf{x}_*) df_* = \int \mathcal{N}(z_*|f_*, \sigma_n^2) \mathcal{N}(f_*|\mu_{f_*}, \sigma_{f_*}^2) df_* \\ &= \mathcal{N}(z_*|\mu_{f_*}, \sigma_n^2 + \sigma_{f_*}^2), \quad \text{where } z_* = g(y_*) \\ p(y_*|\mathcal{D}, \mathbf{x}_*) &= g'(y_*) \mathcal{N}(g(y_*)|\mu_{f_*}, \sigma_n^2 + \sigma_{f_*}^2). \end{aligned}$$

Hence, the predictive moments are obtained by the 1d integrals

$$\begin{aligned} \mu_{y_*} &= \int y_* g'(y_*) \mathcal{N}(g(y_*)|\mu_{f_*}, \sigma_n^2 + \sigma_{f_*}^2) dy_* \\ &= \int g^{-1}(z_*) \mathcal{N}(z_*|\mu_{f_*}, \sigma_n^2 + \sigma_{f_*}^2) dz_*, \quad \text{and} \\ \sigma_{y_*}^2 &= \int (y_* - \mu_{y_*})^2 g'(y_*) \mathcal{N}(g(y_*)|\mu_{f_*}, \sigma_n^2 + \sigma_{f_*}^2) dy_* \\ &= \int (g^{-1}(z_*) - \mu_{y_*})^2 \mathcal{N}(z_*|\mu_{f_*}, \sigma_n^2 + \sigma_{f_*}^2) dz_*. \end{aligned}$$

## 4.8 Gumbel Likelihood

Distributions of extrema are well captured by the Gumbel distribution

$$p(y) = \frac{1}{\beta} \exp(-z - e^{-z}), \quad z = s \frac{y - \eta}{\beta}, \quad s \in \{\pm 1\}$$

with mean  $\mu = \eta + \beta\gamma$  and variance  $\sigma^2 = \pi^2\beta^2/6$  where  $\gamma = 0.57721566490153$  denotes Euler–Mascheroni’s constant. Skewness is approximately given by  $1.1395s$  where  $s$  is a sign switching between left and right skewness and kurtosis is  $12/5$ . The final expression for the Gumbel likelihood is

$$p(y|f) = \frac{\pi}{\sigma\sqrt{6}} \exp(-z - e^{-z}), \quad z = \gamma + s \frac{\pi}{\sigma\sqrt{6}}(y - f), \quad s \in \{\pm 1\}.$$

## 4.9 Laplace Likelihood

### Laplace’s Approximation

The following derivatives are needed:

$$\begin{aligned} \ln p(y|f) &= -\ln(2b) - \frac{|f - y|}{b} \\ \frac{\partial \ln p}{\partial f} &= \frac{\text{sign}(f - y)}{b} \\ \frac{\partial^2 \ln p}{(\partial f)^2} &= \frac{\partial^3 \ln p}{(\partial f)^3} = \frac{\partial^3 \ln p}{(\partial \ln \sigma_n)(\partial f)^2} = 0 \\ \frac{\partial \ln p}{\partial \ln \sigma_n} &= \frac{|f - y|}{b} - 1 \end{aligned}$$

### Expectation Propagation

Expectation propagation requires integration against a Gaussian measure for moment matching.

We need to evaluate  $\ln Z = \ln \int \mathcal{L}(y|f, \sigma_n^2) \mathcal{N}(f|\mu, \sigma^2) df$  as well as the derivatives  $\frac{\partial \ln Z}{\partial \mu}$  and  $\frac{\partial^2 \ln Z}{\partial \mu^2}$  where  $\mathcal{N}(f|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(f-\mu)^2}{2\sigma^2}\right)$ ,  $\mathcal{L}(y|f, \sigma_n^2) = \frac{1}{2b} \exp\left(-\frac{|y-f|}{b}\right)$ , and  $b = \frac{\sigma_n}{\sqrt{2}}$ . As a first step, we reduce the number of parameters by means of the substitution  $\tilde{f} = \frac{f-y}{\sigma_n}$  yielding

$$\begin{aligned} Z &= \int \mathcal{L}(y|f, \sigma_n^2) \mathcal{N}(f|\mu, \sigma^2) df \\ &= \frac{1}{\sqrt{2\pi\sigma}} \frac{\sqrt{2}}{2\sigma_n} \int \exp\left(-\frac{(f-\mu)^2}{2\sigma^2}\right) \exp\left(-\sqrt{2} \frac{|f-y|}{\sigma_n}\right) df \\ &= \frac{\sqrt{2}}{2\sigma\sqrt{2\pi}} \int \exp\left(-\frac{(\sigma_n \tilde{f} + y - \mu)^2}{2\sigma^2}\right) \exp\left(-\sqrt{2} |\tilde{f}| \right) d\tilde{f} \\ &= \frac{\sigma_n}{\sigma\sigma_n\sqrt{2\pi}} \int \exp\left(-\frac{\sigma_n^2 \left(\tilde{f} - \frac{\mu-y}{\sigma_n}\right)^2}{2\sigma^2}\right) \mathcal{L}(\tilde{f}|0, 1) d\tilde{f} \\ &= \frac{1}{\sigma_n} \int \mathcal{L}(f|0, 1) \mathcal{N}(f|\tilde{\mu}, \tilde{\sigma}^2) df \\ \ln Z &= \ln \tilde{Z} - \ln \sigma_n = \ln \int \mathcal{L}(f|0, 1) \mathcal{N}(f|\tilde{\mu}, \tilde{\sigma}^2) df - \ln \sigma_n \end{aligned}$$



with  $\tilde{\mu} = \frac{\mu - \eta}{\sigma_n}$  and  $\tilde{\sigma} = \frac{\sigma}{\sigma_n}$ . Thus, we concentrate on the simpler quantity  $\ln \tilde{Z}$ .

$$\begin{aligned}
\ln Z &= \ln \int \exp \left( -\frac{(f - \tilde{\mu})^2}{2\tilde{\sigma}^2} - \sqrt{2}|f| \right) df - \overbrace{\ln \tilde{\sigma} \sqrt{2\pi} - \ln \sqrt{2}\sigma_n}^C \\
&= \ln \left[ \int_{-\infty}^0 \exp \left( -\frac{(f - \tilde{\mu})^2}{2\tilde{\sigma}^2} + \sqrt{2}f \right) df + \int_0^{\infty} \exp \left( -\frac{(f - \tilde{\mu})^2}{2\tilde{\sigma}^2} - \sqrt{2}f \right) df \right] + C \\
&= \ln \left[ \int_{-\infty}^0 \exp \left( -\frac{f^2 - 2(\overbrace{\tilde{\mu} + \tilde{\sigma}^2 \sqrt{2}}^{m_-})f + \tilde{\mu}^2}{2\tilde{\sigma}^2} \right) df + \int_0^{\infty} \exp \left( -\frac{f^2 - 2(\overbrace{\tilde{\mu} - \tilde{\sigma}^2 \sqrt{2}}^{m_+})f + \tilde{\mu}^2}{2\tilde{\sigma}^2} \right) df \right] + C \\
&= \ln \left[ \exp \left( \frac{m_-^2}{2\tilde{\sigma}^2} \right) \int_{-\infty}^0 \exp \left( -\frac{(f - m_-)^2}{2\tilde{\sigma}^2} \right) df + \exp \left( \frac{m_+^2}{2\tilde{\sigma}^2} \right) \int_0^{\infty} \exp \left( -\frac{(f - m_+)^2}{2\tilde{\sigma}^2} \right) df \right] - \frac{\tilde{\mu}^2}{2\tilde{\sigma}^2} + C \\
&= \ln \left[ \exp \left( \frac{m_-^2}{2\tilde{\sigma}^2} \right) \int_{-\infty}^0 \mathcal{N}(f|m_-, \tilde{\sigma}^2) df + \exp \left( \frac{m_+^2}{2\tilde{\sigma}^2} \right) \left( 1 - \int_{-\infty}^0 \mathcal{N}(f|m_+, \tilde{\sigma}^2) df \right) \right] - \frac{\tilde{\mu}^2}{2\tilde{\sigma}^2} - \ln \sqrt{2}\sigma_n \\
&= \ln \left[ \exp \left( \frac{m_-^2}{2\tilde{\sigma}^2} \right) \Phi \left( \frac{m_-}{\tilde{\sigma}} \right) - \exp \left( \frac{m_+^2}{2\tilde{\sigma}^2} \right) \Phi \left( \frac{m_+}{\tilde{\sigma}} \right) + \exp \left( \frac{m_+^2}{2\tilde{\sigma}^2} \right) \right] - \frac{\tilde{\mu}^2}{2\tilde{\sigma}^2} - \ln \sqrt{2}\sigma_n
\end{aligned}$$

Here,  $\Phi(z) = \int_{-\infty}^z \mathcal{N}(f|0, 1)df$  denotes the cumulative Gaussian distribution. Finally, we have

$$\begin{aligned}
\ln Z &= \ln \left[ \exp \left( -\sqrt{2}\tilde{\mu} \right) \Phi \left( \frac{m_-}{\tilde{\sigma}} \right) + \exp \left( \sqrt{2}\tilde{\mu} \right) \Phi \left( -\frac{m_+}{\tilde{\sigma}} \right) \right] + \tilde{\sigma}^2 - \ln \sqrt{2}\sigma_n \\
&= \ln \left[ \exp \left( \underbrace{\ln \Phi(-z_+) + \sqrt{2}\tilde{\mu}}_{a_+} \right) + \exp \left( \underbrace{\ln \Phi(z_-) - \sqrt{2}\tilde{\mu}}_{a_-} \right) \right] + \tilde{\sigma}^2 - \ln \sqrt{2}\sigma_n \\
&= \ln(e^{a_+} + e^{a_-}) + \tilde{\sigma}^2 - \ln \sqrt{2}\sigma_n
\end{aligned}$$

where  $z_+ = \frac{\tilde{\mu}}{\tilde{\sigma}} + \tilde{\sigma}\sqrt{2} = \frac{\mu - \eta}{\sigma} + \frac{\sigma}{\sigma_n}\sqrt{2}$ ,  $z_- = \frac{\tilde{\mu}}{\tilde{\sigma}} - \tilde{\sigma}\sqrt{2} = \frac{\mu - \eta}{\sigma} - \frac{\sigma}{\sigma_n}\sqrt{2}$  and  $\tilde{\mu} = \frac{\mu - \eta}{\sigma_n}$ ,  $\tilde{\sigma} = \frac{\sigma}{\sigma_n}$ .

Now, using  $\frac{d}{d\theta} \ln \Phi(z) = \frac{1}{\Phi(z)} \frac{d}{d\theta} \Phi(z) = \frac{\mathcal{N}(z)}{\Phi(z)} \frac{dz}{d\theta}$  we tackle first derivative

$$\begin{aligned}
\frac{\partial \ln Z}{\partial \mu} &= \frac{e^{a_+} \frac{\partial a_+}{\partial \mu} + e^{a_-} \frac{\partial a_-}{\partial \mu}}{e^{a_+} + e^{a_-}} \\
\frac{\partial a_+}{\partial \mu} &= \frac{\partial}{\partial \mu} \ln \Phi(-z_+) + \frac{\sqrt{2}}{\sigma_n} \\
&= -\frac{\mathcal{N}(-z_+)}{\sigma \Phi(-z_+)} + \frac{\sqrt{2}}{\sigma_n} = -\frac{q_+}{\sigma} + \frac{\sqrt{2}}{\sigma_n} \\
\frac{\partial a_-}{\partial \mu} &= \frac{\partial}{\partial \mu} \ln \Phi(z_-) - \frac{\sqrt{2}}{\sigma_n} \\
&= \frac{\mathcal{N}(z_-)}{\sigma \Phi(z_-)} - \frac{\sqrt{2}}{\sigma_n} = \frac{q_-}{\sigma} - \frac{\sqrt{2}}{\sigma_n} \\
\frac{\partial a_{\pm}}{\partial \mu} &= \mp \frac{q_{\pm}}{\sigma} \pm \frac{\sqrt{2}}{\sigma_n}.
\end{aligned}$$

as well as the second derivative

$$\begin{aligned}
\frac{\partial^2 \ln Z}{\partial \mu^2} &= \frac{\frac{\partial}{\partial \mu} \left( e^{a_+} \frac{\partial a_+}{\partial \mu} \right) + \frac{\partial}{\partial \mu} \left( e^{a_-} \frac{\partial a_-}{\partial \mu} \right)}{e^{a_+} + e^{a_-}} - \left( \frac{\partial \ln Z}{\partial \mu} \right)^2 \\
\frac{\partial}{\partial \mu} \left( e^{a_{\pm}} \frac{\partial a_{\pm}}{\partial \mu} \right) &= e^{a_{\pm}} \left[ \left( \frac{\partial a_{\pm}}{\partial \mu} \right)^2 + \frac{\partial^2 a_{\pm}}{\partial \mu^2} \right] \\
\frac{\partial^2 a_+}{\partial \mu^2} &= -\frac{1}{\sigma} \frac{\frac{\partial}{\partial \mu} \mathcal{N}(-z_+) \Phi(-z_+) - \frac{\partial}{\partial \mu} \Phi(-z_+) \mathcal{N}(-z_+)}{\Phi^2(-z_+)} \\
&= -\frac{1}{\sigma} \frac{\mathcal{N}(-z_+) \Phi(-z_+) \frac{\partial -z_+/2}{\partial \mu} - \mathcal{N}^2(-z_+) \frac{\partial -z_+}{\partial \mu}}{\Phi^2(-z_+)} \\
&= \frac{\mathcal{N}(-z_+)}{\sigma^2} \cdot \frac{\Phi(-z_+) z_+ - \mathcal{N}(-z_+)}{\Phi^2(-z_+)} = -\frac{q_+^2 - q_+ z_+}{\sigma^2} \\
\frac{\partial^2 a_-}{\partial \mu^2} &= \frac{1}{\sigma} \frac{\frac{\partial}{\partial \mu} \mathcal{N}(z_-) \Phi(z_-) - \frac{\partial}{\partial \mu} \Phi(z_-) \mathcal{N}(z_-)}{\Phi^2(z_-)} \\
&= \frac{1}{\sigma} \frac{\mathcal{N}(z_-) \Phi(z_-) \frac{\partial -z_-/2}{\partial \mu} - \mathcal{N}^2(z_-) \frac{\partial z_-}{\partial \mu}}{\Phi^2(z_-)} \\
&= \frac{\mathcal{N}(z_-)}{\sigma^2} \cdot \frac{-\Phi(z_-) z_- - \mathcal{N}(z_-)}{\Phi^2(z_-)} = -\frac{q_-^2 + q_- z_-}{\sigma^2} \\
\frac{\partial^2 a_{\pm}}{\partial \mu^2} &= -\frac{q_{\pm}^2 \mp q_{\pm} z_{\pm}}{\sigma^2}
\end{aligned}$$

which can be simplified to

$$\frac{\partial^2 \ln Z}{\partial \mu^2} = \frac{e^{a_+} b_+ + e^{a_-} b_-}{e^{a_+} + e^{a_-}} - \left( \frac{\partial \ln Z}{\partial \mu} \right)^2$$

using

$$\begin{aligned}
b_{\pm} = \left( \frac{\partial a_{\pm}}{\partial \mu} \right)^2 + \frac{\partial^2 a_{\pm}}{\partial \mu^2} &= \left( \mp \frac{q_{\pm}}{\sigma} \pm \frac{\sqrt{2}}{\sigma_n} \right)^2 - \frac{q_{\pm}^2 \mp q_{\pm} z_{\pm}}{\sigma^2} \\
&= \left( \frac{q_{\pm}}{\sigma} - \frac{\sqrt{2}}{\sigma_n} \right)^2 - \frac{q_{\pm}^2}{\sigma^2} \pm \frac{q_{\pm} z_{\pm}}{\sigma^2} \\
&= \frac{2}{\sigma_n^2} - \left( \frac{\sqrt{8}}{\sigma \sigma_n} \mp \frac{z_{\pm}}{\sigma^2} \right) q_{\pm}.
\end{aligned}$$

We also need

$$\frac{\partial \ln Z}{\partial \ln \sigma_n} = \frac{e^{a_+} \frac{\partial a_+}{\partial \ln \sigma_n} + e^{a_-} \frac{\partial a_-}{\partial \ln \sigma_n}}{e^{a_+} + e^{a_-}} - \frac{2\sigma^2}{\sigma_n^2} - 1.$$

## Variational Bayes

We need  $h(\gamma)$  and its derivatives as well as  $\beta(\gamma)$ :

$$\begin{aligned}
h(\gamma) &= \frac{2}{\sigma_n^2} \gamma + \ln(2\sigma_n^2) + y^2 \gamma^{-1} \\
h'(\gamma) &= \frac{2}{\sigma_n^2} - y^2 \gamma^{-2} \\
h''(\gamma) &= 2y^2 \gamma^{-3} \\
\beta(\gamma) &= y \gamma^{-1}
\end{aligned}$$

#### 4.10 Student's t Likelihood

The likelihood has two hyperparameters (both represented in the log domain to ensure positivity): the degrees of freedom  $\nu$  and the scale  $\sigma_n$  with mean  $y$  (for  $\nu > 1$ ) and variance  $\frac{\nu}{\nu-2} \sigma_n^2$  (for  $\nu > 2$ ).

$$p(y|f) = Z \cdot \left(1 + \frac{(f-y)^2}{\nu \sigma_n^2}\right)^{-\frac{\nu+1}{2}}, \quad Z = \frac{\Gamma(\frac{\nu+1}{2})}{\Gamma(\frac{\nu}{2}) \sqrt{\nu \pi \sigma_n^2}}$$

#### Laplace's Approximation

For the mode fitting procedure, we need derivatives up to third order; the hyperparameter derivatives at the mode require some mixed derivatives. All in all, using  $r = y - f$ , we have

$$\begin{aligned}
\ln p(y|f) &= \ln \Gamma\left(\frac{\nu+1}{2}\right) - \ln \Gamma\left(\frac{\nu}{2}\right) - \frac{1}{2} \ln \nu \pi \sigma_n^2 - \frac{\nu+1}{2} \ln \left(1 + \frac{r^2}{\nu \sigma_n^2}\right) \\
\frac{\partial \ln p}{\partial f} &= (\nu+1) \frac{r}{r^2 + \nu \sigma_n^2} \\
\frac{\partial^2 \ln p}{(\partial f)^2} &= (\nu+1) \frac{r^2 - \nu \sigma_n^2}{(r^2 + \nu \sigma_n^2)^2} \\
\frac{\partial^3 \ln p}{(\partial f)^3} &= 2(\nu+1) \frac{r^3 - 3r \nu \sigma_n^2}{(r^2 + \nu \sigma_n^2)^3} \\
\frac{\partial \ln p}{\partial \ln \nu} &= \frac{\partial Z}{\partial \ln \nu} - \frac{\nu}{2} \ln \left(1 + \frac{r^2}{\nu \sigma_n^2}\right) + \frac{\nu+1}{2} \cdot \frac{r^2}{r^2 + \nu \sigma_n^2} \\
\frac{\partial Z}{\partial \ln \nu} &= \frac{\nu}{2} \frac{d \ln \Gamma(\frac{\nu+1}{2})}{d \ln \nu} - \frac{\nu}{2} \frac{d \ln \Gamma(\frac{\nu}{2})}{d \ln \nu} - \frac{1}{2} \\
\frac{\partial^3 \ln p}{(\partial \ln \nu)(\partial f)^2} &= \nu \frac{r^2(r^2 - 3(\nu+1)\sigma_n^2) + \nu \sigma_n^2}{(r^2 + \nu \sigma_n^2)^3} \\
\frac{\partial \ln p}{\partial \ln \sigma_n} &= (\nu+1) \frac{r^2}{r^2 + \nu \sigma_n^2} - 1 \\
\frac{\partial^3 \ln p}{(\partial \ln \sigma_n)(\partial f)^2} &= 2\nu \sigma_n^2 (\nu+1) \frac{\nu \sigma_n^2 - 3r^2}{(r^2 + \nu \sigma_n^2)^3}
\end{aligned}$$

#### 4.11 Cumulative Logistic Likelihood

The likelihood has one hyperparameter (represented in the log domain), namely the standard deviation  $\sigma_n$

$$p(y|f) = Z \cdot \cosh^{-2}(\tau(f-y)), \quad \tau = \frac{\pi}{2\sigma_n\sqrt{3}}, \quad Z = \frac{\pi}{4\sigma_n\sqrt{3}}$$

## Laplace's Approximation

The following derivatives are needed where  $\phi(x) \equiv \ln(\cosh(x))$

$$\begin{aligned}
\ln p(y|f) &= \ln(\pi) - \ln(4\sigma_n\sqrt{3}) - 2\phi(\tau(f-y)) \\
\frac{\partial \ln p}{\partial f} &= 2\tau\phi'(\tau(f-y)) \\
\frac{\partial^2 \ln p}{(\partial f)^2} &= -2\tau^2\phi''(\tau(f-y)) \\
\frac{\partial^3 \ln p}{(\partial f)^3} &= 2\tau^3\phi'''(\tau(f-y)) \\
\frac{\partial^3 \ln p}{(\partial \ln \sigma_n)(\partial f)^2} &= 2\tau^2(2\phi''(\tau(f-y)) + \tau(f-y)\phi'''(\tau(f-y))) \\
\frac{\partial \ln p}{\partial \ln \sigma_n} &= 2\tau(f-y)\phi'(\tau(f-y)) - 1
\end{aligned}$$

## 4.12 GLM Likelihoods: Poisson, Weibull, Gamma, Exponential, Inverse Gaussian and Beta

Data  $y$  from a space other than  $\mathbb{R}$  e.g.  $\mathbb{N}$ ,  $\mathbb{R}_+$  or  $[0, 1]$  can be modeled using generalised linear model likelihoods  $p(y|f)$  where the expected value  $\mathbb{E}[y] = \mu$  is related to the underlying Gaussian process  $f$  by means of an inverse link function  $\mu = g(f)$ . Typically, the likelihoods are from an exponential family, hence the variance  $\mathbb{V}[y] = v(\mu)$ , is a simple function of the mean  $\mu$  as well as higher order moments such as skewness  $\mathbb{S}[y] = s(\mu)$  and kurtosis  $\mathbb{K}[y] = k(\mu)$ .

Here, we directly specify the inverse link function  $\mu = g(f)$  defining the mapping from the GP  $f$  to the mean intensity  $\mu$ . For numerical reasons, we work with the log of the inverse link function  $h(f) = \ln g(f)$  and use its derivatives  $h'$ ,  $h''$  and  $h'''$  for subsequent computations. In the table below, we have summarised the GLM likelihood expressions, the moments, the range of their variables and the applicable inverse link functions.

Likelihood	$\rho =$	$v(\mu) =$	$s(\mu) =$	$k(\mu) =$	$p(y f) =$	$y \in$	$\mu \in$	Inverse Links
Poisson	$\emptyset$	$\mu$	$1/\sqrt{\mu}$	$1/\mu$	$\mu^y \exp(-\mu)/y!$	$\mathbb{N}$	$\mathbb{R}_+$	exp, logistic
Weibull	$\{\ln \kappa\}$	$\mu^2(\gamma_2/\gamma_1^2 - 1)$	$\frac{\gamma_3 - 3\gamma_1\gamma_2 + 2\gamma_1^3}{(\gamma_2 - \gamma_1^2)^{3/2}}$	$\frac{\gamma_4 - 4\gamma_1\gamma_3 + 12\gamma_1^2\gamma_2 - 3\gamma_2^2 - 6\gamma_1^4}{(\gamma_2 - \gamma_1^2)^2}$	$\kappa\gamma_1/\mu (y\gamma_1/\mu)^{\kappa-1} \exp(-(y\gamma_1/\mu)^\kappa)$	$\mathbb{R}_+ \setminus \{0\}$	$\mathbb{R}_+ \setminus \{0\}$	exp, logistic
Gamma	$\{\ln \alpha\}$	$\mu^2/\alpha$	$2/\sqrt{\alpha}$	$6/\alpha$	$\frac{\alpha^\alpha y^{\alpha-1}}{\Gamma(\alpha)} \mu^{-\alpha} \exp\left(-\frac{y\alpha}{\mu}\right)$	$\mathbb{R}_+ \setminus \{0\}$	$\mathbb{R}_+ \setminus \{0\}$	exp, logistic
Exponential	$\emptyset$	$\mu^2$	2	6	$\mu^{-1} \exp\left(-\frac{y}{\mu}\right)$	$\mathbb{R}_+ \setminus \{0\}$	$\mathbb{R}_+ \setminus \{0\}$	exp, logistic
Inv. Gauss	$\{\ln \lambda\}$	$\mu^3/\lambda$	$3\sqrt{\mu/\lambda}$	$15\mu/\lambda$	$\sqrt{\frac{\lambda}{2\pi y^3}} \exp\left(-\frac{\lambda(y-\mu)^2}{2\mu^2 y}\right)$	$\mathbb{R}_+ \setminus \{0\}$	$\mathbb{R}_+ \setminus \{0\}$	exp, logistic
Beta	$\{\ln \phi\}$	$\mu(1-\mu)/(1+\phi)$	$\frac{(2-4\mu)(1+\phi)}{\sqrt{v(\mu)(2+\phi)}}$	$6 \frac{(\phi+1)^2 - v(\mu)(5\phi+6)}{v(\mu)(\phi+2)(\phi+3)}$	$\frac{\Gamma(\phi)}{\Gamma(\mu\phi)\Gamma((1-\mu)\phi)} y^{\mu\phi-1} (1-y)^{(1-\mu)\phi-1}$	$[0, 1]$	$[0, 1]$	expexp, logit

### 4.12.1 Inverse Link Functions

Possible inverse link functions and their properties ( $\cup$  convex,  $\cap$  concave,  $\uparrow$  monotone) are summarised below:

util/glm_invlink_*	$g(f) = \mu =$	$g : \mathbb{R} \rightarrow$	$g$ is	$h(f) = \ln \mu =$	$h$ is
exp	$e^f$	$\mathbb{R}_+$	$\cup, \uparrow$	$f$	$\cup, \cap, \uparrow$
logistic	$\ln(1 + e^f)$	$\mathbb{R}_+$	$\cup, \uparrow$	$\ln(\ln(1 + e^f))$	$\cap, \uparrow$
expexp	$\exp(-e^{-f})$	$[0, 1]$	$\uparrow$	$-e^{-f}$	$\cup, \uparrow$
logit	$1/(1 + e^{-f})$	$[0, 1]$	$\uparrow$	$-\ln(1 + e^{-f})$	$\cup, \uparrow$

**Exponential inverse link:** `exp`

For  $g(f) = e^f$  things are simple since  $h(f) = f$ ,  $h'(f) = 1$  and  $h''(f) = h'''(f) = 0$ .

**Logistic inverse link:** `logistic`

For  $g(f) = \ln(1 + e^f)$  the derivatives of  $h(f)$  are given by

$$\begin{aligned} h(f) &= \ln(\ln(1 + e^f)) \\ h'(f) &= \frac{1}{\ln(1 + e^f)} s(-f), \quad s(f) = \frac{1}{1 + e^f}, \quad s'(f) = \frac{-e^f}{(1 + e^f)^2} = -s(-f)s(f) \\ h''(f) &= \frac{1}{\ln(1 + e^f)} \frac{e^{-f}}{(1 + e^{-f})^2} - \frac{1}{\ln^2(1 + e^f)} \frac{e^f}{1 + e^f} \frac{1}{1 + e^{-f}} \\ &= h'(f) [s(f) - h'(f)] \\ h'''(f) &= h''(f) [s(f) - h'(f)] + h'(f) \left[ \frac{-e^f}{(1 + e^f)^2} - h''(f) \right] \\ &= h''(f) [s(f) - 2h'(f)] - h'(f)s(f)s(-f). \end{aligned}$$

Note that  $g(f) = e^{h(f)} = \ln(1 + e^f)$  is convex and  $h(f) = \ln(\ln(1 + e^f))$  with

$$h''(f) = \frac{1}{\ln(1 + e^f)} \left( 1 - \frac{e^f}{\ln(1 + e^f)} \right) \frac{1}{1 + e^f} \frac{1}{1 + e^{-f}} \leq 0$$

is concave since  $e^f \geq \ln(1 + e^f)$  for all  $f \in \mathbb{R}$ .

**Double negative exponential inverse link:** `expexp`

For  $g(f) = \exp(-e^{-f})$  the derivatives of  $h(f)$  are given by

$$\begin{aligned} h(f) &= -e^{-f} \\ h'(f) &= -h(f) \\ h''(f) &= h(f) \\ h'''(f) &= -h(f) \end{aligned}$$

**Logit regression inverse link:** `logit`

For  $g(f) = 1/(1 + e^{-f})$  the derivatives of  $h(f)$  can be computed using the logistic inverse link function  $h_\ell(f)$  since  $h(f) = f - \exp(h_\ell(f))$

$$\begin{aligned} h(f) &= f - e^{h_\ell(f)} \\ h'(f) &= 1 - e^{h_\ell(f)} h'_\ell(f) \\ h''(f) &= -e^{h_\ell(f)} [h'_\ell(f)^2 + h''_\ell(f)] = e^{h_\ell(f)} s_\ell(-f) s_\ell^2(f) \\ h'''(f) &= -e^{h_\ell(f)} [h'_\ell(f)^3 + 3h''_\ell(f) h'_\ell(f) + h'''_\ell(f)] \end{aligned}$$

#### 4.12.2 Poisson Likelihood

Count data  $y \in \mathbb{N}^n$  can be modeled in the GP framework using the Poisson distribution  $p(y) = \mu^y e^{-\mu} / y!$  with mean/variance  $\mathbb{E}[y] = \mathbb{V}[y] = \mu$ , skewness  $\mathbb{S}[y] = 1/\sqrt{\mu}$  and kurtosis  $\mathbb{K}[y] = 1/\mu$

leading to the likelihood

$$\begin{aligned} p(y|f) &= \mu^y \exp(-\mu)/y!, \mu = g(f) \\ \Leftrightarrow \ln p(y|f) &= y \cdot \ln g(f) - g(f) - \ln \Gamma(y+1). \end{aligned}$$

For Laplace's method to work, we need the first three derivatives of the log likelihood  $\ln p(y|f)$ , where  $h(f) = \ln g(f)$

$$\begin{aligned} \ln p(y|f) &= y \cdot h(f) - \exp(h(f)) - \ln \Gamma(y+1) \\ \frac{\partial}{\partial f} \ln p(y|f) &= h'(f) [y - \exp(h(f))] \\ \frac{\partial^2}{\partial f^2} \ln p(y|f) &= h''(f) [y - \exp(h(f))] - [h'(f)]^2 \exp(h(f)) \\ \frac{\partial^3}{\partial f^3} \ln p(y|f) &= h'''(f) [y - \exp(h(f))] - 3h'(f) \cdot h''(f) \exp(h(f)) - [h'(f)]^3 \exp(h(f)) \\ &\quad h'''(f) [y - \exp(h(f))] - h'(f)[h'(f)^2 + 3h''(f)] \exp(h(f)). \end{aligned}$$

Note that if  $\ln \mu = h(f)$  is concave and  $\mu = g(f)$  is convex then the Poisson likelihood  $p(y|f)$  is log-concave in  $f$  which is the case for both `exp` and `logistic`.

#### 4.12.3 Weibull Likelihood

Nonnegative data  $y \in \mathbb{R}_+$  such as time-to-failure can be modeled in the GP framework using the Weibull distribution  $p(y) = \kappa/\lambda(y/\lambda)^{\kappa-1}e^{-(y/\lambda)^\kappa}$  with shape parameter  $\kappa > 0$ , scale parameter  $\lambda > 0$ , mean  $\mathbb{E}[y] = \lambda\gamma_1 = \mu$  where  $\gamma_j = \Gamma(1 + j/\kappa)$ , variance  $\mathbb{V}[y] = \lambda^2\gamma_2 - \mu^2 = \mu^2(\gamma_2/\gamma_1^2 - 1)$ , skewness  $\mathbb{S}[y] = (\gamma_3 - 3\gamma_1\gamma_2 + 2\gamma_1^3)/(\gamma_2 - \gamma_1^2)^{3/2}$  and kurtosis  $\mathbb{K}[y] = (\gamma_4 - 4\gamma_1\gamma_3 + 12\gamma_1^2\gamma_2 - 3\gamma_2^2 - 6\gamma_1^4)/(\gamma_2 - \gamma_1^2)^2$ . Using the substitution  $\mu = \lambda\gamma_1 \Leftrightarrow 1/\lambda = \gamma_1/\mu$ , we obtain

$$\begin{aligned} p(y|f) &= \gamma_1 \frac{\kappa}{\mu} \left( \gamma_1 \frac{y}{\mu} \right)^{\kappa-1} \exp \left( - \left( \gamma_1 \frac{y}{\mu} \right)^\kappa \right), \mu = g(f) > 0 \\ \Leftrightarrow \ln p(y|f) &= \ln \left( \gamma_1 \frac{\kappa}{\mu} \right) + (\kappa-1) \ln \left( \gamma_1 \frac{y}{\mu} \right) - \left( \gamma_1 \frac{y}{\mu} \right)^\kappa. \end{aligned}$$

Note that the Weibull likelihood  $p(y|f)$  is log-concave in  $f$  neither for the `exp` nor for the `logistic` inverse link.

#### 4.12.4 Gamma Likelihood

Nonnegative data  $y \in \mathbb{R}_+$  can be modeled in the GP framework using the Gamma distribution  $p(y) = \theta^{-\alpha}/\Gamma(\alpha)y^{\alpha-1}e^{-y/\theta}$  with shape parameter  $\alpha > 0$ , scale parameter  $\theta > 0$ , mean  $\mathbb{E}[y] = \alpha\theta = \mu$ , variance  $\mathbb{V}[y] = \alpha\theta^2 = \mu^2/\alpha$ , skewness  $\mathbb{S}[y] = 2/\sqrt{\alpha}$  and kurtosis  $\mathbb{K}[y] = 6/\alpha$ . Using the substitution  $\mu = \alpha\theta \Leftrightarrow \alpha/\mu = 1/\theta$ , we obtain

$$\begin{aligned} p(y|f) &= \frac{\alpha^\alpha y^{\alpha-1}}{\Gamma(\alpha)} \mu^{-\alpha} \exp \left( - \frac{y\alpha}{\mu} \right), \mu = g(f) > 0 \\ \Leftrightarrow \ln p(y|f) &= -\alpha \left( \ln \mu + \frac{y}{\mu} \right) - \ln Z_\alpha(y), \ln Z_\alpha(y) = \ln \Gamma(\alpha) - \alpha \ln \alpha + (1-\alpha) \ln y. \end{aligned}$$

Note that if  $\ln \mu = h(f)$  was convex and  $\mu = g(f)$  was concave then the Gamma likelihood  $p(y|f)$  would be log-concave in  $f$  which is not the case for both `exp` and `logistic`.

#### 4.12.5 Exponential Likelihood

Nonnegative data  $y \in \mathbb{R}_+$  can be modeled in the GP framework using the Exponential distribution  $p(y) = \theta^{-1} e^{-y/\theta}$  with scale parameter  $\theta > 0$ , mean  $\mathbb{E}[y] = \theta = \mu$ , variance  $\mathbb{V}[y] = \mu^2$ , skewness  $\mathbb{S}[y] = 2$  and kurtosis  $\mathbb{K}[y] = 6$ . We obtain

$$\begin{aligned} p(y|f) &= \mu^{-1} \exp\left(-\frac{y}{\mu}\right), \mu = g(f) > 0 \\ \Leftrightarrow \ln p(y|f) &= -\ln \mu - \frac{y}{\mu}. \end{aligned}$$

Note that for `exp` (but not for `logistic`) the likelihood is log-concave. The exponential distribution corresponds to the Gamma distribution with  $\alpha = 1$  and the Weibull distribution with  $\kappa = 1$ .

#### 4.12.6 Inverse Gaussian Likelihood

Nonnegative data  $y \in \mathbb{R}_+^n$  can be modeled in the GP framework using the Inverse Gaussian distribution  $p(y) = \sqrt{\lambda/(2\pi y^3)} \exp(-\lambda(y - \mu)^2/(2\mu^2 y))$  with shape parameter  $\lambda > 0$ , mean parameter  $\mu > 0$ , mean  $\mathbb{E}[y] = \mu$ , variance  $\mathbb{V}[y] = \mu^3/\lambda$ , skewness  $\mathbb{S}[y] = 3\sqrt{\mu/\lambda}$  and kurtosis  $\mathbb{K}[y] = 15\mu/\lambda$ . We obtain

$$\begin{aligned} p(y|f) &= \sqrt{\frac{\lambda}{2\pi y^3}} \exp\left(-\frac{\lambda(y - \mu)^2}{2\mu^2 y}\right), \mu = g(f) > 0 \\ \Leftrightarrow \ln p(y|f) &= -\frac{\lambda(y - \mu)^2}{2\mu^2 y} - \ln Z_\alpha(y), \ln Z_\alpha(y) = -\frac{1}{2}(\ln \lambda - \ln 2\pi y^3). \end{aligned}$$

The inverse Gaussian likelihood is in general not log-concave in  $f$  for both `exp` and `logistic`.

#### 4.12.7 Beta Likelihood

Interval data  $y \in [0, 1]^n$  can be modeled in the GP framework using the Beta distribution  $p(y) = y^{\alpha-1}(1-y)^{\beta-1}/B(\alpha, \beta)$  with shape parameters  $\alpha, \beta > 0$ , mean  $\mathbb{E}[y] = \alpha/(\alpha + \beta)$  and variance  $\mathbb{V}[y] = \alpha\beta/[(\alpha + \beta)^2(\alpha + \beta + 1)]$  and  $1/B(\alpha, \beta) = \Gamma(\alpha + \beta)/[\Gamma(\alpha)\Gamma(\beta)]$ . Reparametrising using the mean parameter  $\mu = \mathbb{E}[y] = \alpha/(\alpha + \beta)$ , the shape parameter  $\phi = \alpha + \beta$ , the variance  $\mathbb{V}[y] = \mu(1 - \mu)/(1 + \phi)$  and hence

$$\begin{aligned} p(y|f) &= \frac{\Gamma(\phi)}{\Gamma(\mu\phi)\Gamma((1-\mu)\phi)} y^{\mu\phi-1} (1-y)^{(1-\mu)\phi-1}, \mu = g(f) > 0 \\ \Leftrightarrow \ln p(y|f) &= \ln \Gamma(\phi) - \ln \Gamma(\mu\phi) - \ln \Gamma((1-\mu)\phi) + (\mu\phi - 1) \ln y + ((1-\mu)\phi - 1) \ln(1-y). \end{aligned}$$

The Beta likelihood is in general not log-concave in  $f$  for both `exp` and `logistic`.

## 5 Mean Functions

A mean function  $m_\phi : \mathcal{X} \rightarrow \mathbb{R}$  (with hyperparameters  $\phi$ ) of a GP  $f$  is a scalar function defined over the whole domain  $\mathcal{X}$  that computes the expected value  $m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$  of  $f$  for the input  $\mathbf{x}$ .

### 5.1 Interface

In the GPML toolbox, a mean function  $m : \mathcal{X} \rightarrow \mathbb{R}$  needs to implement evaluation  $\mathbf{m} = m_\phi(\mathbf{X})$  and first derivatives  $\mathbf{m}_i = \frac{\partial}{\partial \phi_i} \mathbf{m}$  with respect to the components  $i$  of the parameter  $\phi \in \Phi$  as detailed below.

```
32 <meanFunctions.m 32>≡
1 % mean functions to be use by Gaussian process functions. There are two
2 % different kinds of mean functions: simple and composite:
3 %
4 % simple mean functions:
5 %
6 %   meanZero      - zero mean function
7 %   meanOne       - one mean function
8 %   meanConst     - constant mean function
9 %   meanLinear    - linear mean function
10 %  meanPoly      - polynomial mean function
11 %
12 % composite covariance functions (see explanation at the bottom):
13 %
14 %   meanScale     - scaled version of a mean function
15 %   meanPow       - power of a mean function
16 %   meanProd      - products of mean functions
17 %   meanSum       - sums of mean functions
18 %   meanMask      - mask some dimensions of the data
19 %
20 % Naming convention: all mean functions are named "mean/mean*.m".
21 %
22 %
23 % 1) With no or only a single input argument:
24 %
25 %   s = meanNAME or s = meanNAME(hyp)
26 %
27 % The mean function returns a string s telling how many hyperparameters hyp it
28 % expects, using the convention that "D" is the dimension of the input space.
29 % For example, calling "meanLinear" returns the string 'D'.
30 %
31 % 2) With two input arguments:
32 %
33 %   m = meanNAME(hyp, x)
34 %
35 % The function computes and returns the mean vector where hyp are the
36 % hyperparameters and x is an n by D matrix of cases, where D is the dimension
37 % of the input space. The returned mean vector is of size n by 1.
38 %
39 % 3) With three input arguments:
40 %
41 %   dm = meanNAME(hyp, x, i)
42 %
43 % The function computes and returns the n by 1 vector of partial derivatives
44 % of the mean vector w.r.t. hyp(i) i.e. hyperparameter number i.
```



```

45 %
46 % See also doc/usageMean.m.
47 %
48 <gpml copyright 5a>

```

## 5.2 Implemented Mean Functions

We offer simple and composite mean functions producing new mean functions  $m(\mathbf{x})$  from existing mean functions  $\mu_j(\mathbf{x})$ . All code files are named according to the pattern `mean/mean<NAME>.m` for simple identification. This modular specification allows to define affine mean functions  $m(\mathbf{x}) = \mathbf{c} + \mathbf{a}^\top \mathbf{x}$  or polynomial mean functions  $m(\mathbf{x}) = (\mathbf{c} + \mathbf{a}^\top \mathbf{x})^2$ . All currently available mean functions are summarised in the following table.

Simple mean functions $m(\mathbf{x})$			
<NAME>	Meaning	$m(\mathbf{x}) =$	$\Phi$
Zero	mean vanishes always	0	$\emptyset$
One	mean equals 1	1	$\emptyset$
Const	mean equals a constant	$\mathbf{c}$	$\mathbf{c} \in \mathbb{R}$
Linear	mean linearly depends on $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^D$	$\mathbf{a}^\top \mathbf{x}$	$\mathbf{a} \in \mathbb{R}^D$
Poly	mean polynomially depends on $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^D$	$\sum_d \mathbf{a}_d^\top \mathbf{x}^d$	$\mathbf{a} \in \mathbb{R}^{D \times d}$
Composite mean functions $[\mu_1(\mathbf{x}), \mu_2(\mathbf{x}), \dots] \mapsto m(\mathbf{x})$			
<NAME>	Meaning	$m(\mathbf{x}) =$	$\Phi$
Scale	scale a mean	$\alpha \mu(\mathbf{x})$	$\alpha \in \mathbb{R}$
Sum	add up mean functions	$\sum_j \mu_j(\mathbf{x})$	$\emptyset$
Prod	multiply mean functions	$\prod_j \mu_j(\mathbf{x})$	$\emptyset$
Pow	raise a mean to a power	$\mu(\mathbf{x})^d$	$\emptyset$
Mask	act on components $I \subseteq [1, 2, \dots, D]$ of $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^D$ only	$\mu(\mathbf{x}_I)$	$\emptyset$

## 5.3 Usage of Implemented Mean Functions

Some code examples taken from `doc/usageMean.m` illustrate how to use simple and composite mean functions to specify a GP model.

Syntactically, a mean function `mf` is defined by

```

mn := 'func' | @func // simple

mf := {mn} | {mn, {param, mf}} | {mn, {mf, ..., mf}} // composite

```

i.e., it is either a string containing the name of a mean function, a pointer to a mean function or one of the former in combination with a cell array of mean functions and an additional list of parameters.

```

33 <doc/usageMean.m 33>≡
  1 % demonstrate usage of mean functions
  2 %
  3 % See also meanFunctions.m.
  4 %
  5 <gpml copyright 5a>
  6 clear all, close all
  7 n = 5; D = 2; x = randn(n,D); % create a random data set
  8
  9 % set up simple mean functions
 10 m0 = {'meanZero'}; hyp0 = []; % no hyperparameters are needed
 11 m1 = {'meanOne'}; hyp1 = []; % no hyperparameters are needed

```

```

12 mc = {@meanConst}; hypc = 2; % also function handles are possible
13 m1 = {@meanLinear}; hyp1 = [2;3]; % m(x) = 2*x1 + 3*x2
14 mp = {@meanPoly,2}; hypp = [1;1;2;3]; % m(x) = x1+x2+2*x1^2+3*x2^2
15
16 % set up composite mean functions
17 msc = {'meanScale',{m1}}; hypsc = [3; hyp1]; % scale by 3
18 msu = {'meanSum',{m0,mc,m1}}; hypsu = [hyp0; hypc; hyp1]; % sum
19 mpr = {@meanProd,{mc,m1}}; hyppr = [hypc; hyp1]; % product
20 mpo = {'meanPow',3,msu}; hyppo = hypsu; % third power
21 mask = [false,true]; % mask excluding all but the 2nd component
22 mma = {'meanMask',mask,m1}; hypma = hyp1(mask);
23
24 % 0) specify mean function
25 % mean = m0; hyp = hyp0;
26 % mean = msu; hyp = hypsu;
27 % mean = mpr; hyp = hyppr;
28 mean = mpo; hyp = hyppo;
29
30 % 1) query the number of parameters
31 feval(mean{:})
32
33 % 2) evaluate the function on x
34 feval(mean{:},hyp,x)
35
36 % 3) compute the derivatives w.r.t. to hyperparameter i
37 i = 2; feval(mean{:},hyp,x,i)

```

## 6 Covariance Functions

A covariance function  $k_{\psi} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  (with hyperparameters  $\psi$ ) of a GP  $f$  is a scalar function defined over the whole domain  $\mathcal{X}^2$  that computes the covariance  $k(\mathbf{x}, \mathbf{x}') = \mathbb{V}[f(\mathbf{x}), f(\mathbf{x}')] = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$  of  $f$  between the inputs  $\mathbf{x}$  and  $\mathbf{x}'$ .

### 6.1 Interface

Again, the interface is simple since only evaluation of the full covariance matrix  $\mathbf{K} = k_{\psi}(\mathbf{X})$  and its derivatives  $\mathbf{K}_i = \frac{\partial}{\partial \psi_i} \mathbf{K}$  as well as cross terms  $\mathbf{k}_* = k_{\psi}(\mathbf{X}, \mathbf{x}_*)$  and  $\mathbf{k}_{**} = k_{\psi}(\mathbf{x}_*, \mathbf{x}_*)$  for prediction are required.

```
35 <covFunctions.m 35>≡
1 % covariance functions to be use by Gaussian process functions. There are two
2 % different kinds of covariance functions: simple and composite:
3 %
4 % simple covariance functions:
5 %   covConst      - covariance for constant functions
6 %   covCos        - sine periodic covariance function (1d) with unit period
7 %   covLIN        - linear covariance function without parameters
8 %   covLINard     - linear covariance function with ARD
9 %   covLINiso     - linear covariance function
10 %  covLINone     - linear covariance function with bias
11 %  covMaternard  - Matern covariance function with nu=1/2, 3/2 or 5/2 with ARD
12 %  covMaterniso  - Matern covariance function with nu=1/2, 3/2 or 5/2
13 %  covNNone     - neural network covariance function
14 %  covNoise     - independent covariance function (i.e. white noise)
15 %  covPeriodic  - smooth periodic covariance function (1d)
16 %  covPeriodicNoDC - as above but with zero DC component and properly scaled
17 %  covPoly      - polynomial covariance function
18 %  covPPard     - piecewise polynomial covariance function (compact support)
19 %  covPPiso     - piecewise polynomial covariance function (compact support)
20 %  covRQard     - rational quadratic covariance function with ARD
21 %  covRQiso     - isotropic rational quadratic covariance function
22 %  covSEard     - squared exponential covariance function with ARD
23 %  covSEiso     - isotropic squared exponential covariance function
24 %  covSEisoU    - same as above but without latent scale
25 %  covSM        - spectral mixture covariance function
26 %  covGaborard  - Gabor covariance function with ARD
27 %  covGaborsio  - isotropic Gabor covariance function
28 %
29 % composite (meta) covariance functions (see explanation at the bottom):
30 %   covScale     - scaled version of a covariance function
31 %   covProd      - products of covariance functions
32 %   covSum       - sums of covariance functions
33 %   covADD       - additive covariance function
34 %   covMask      - mask some dimensions of the data
35 %   covPERard    - make ARD stationary covariance periodic
36 %   covPERiso    - make isotropic stationary covariance periodic
37 %
38 % special purpose (wrapper) covariance functions
39 %   covFITC      - to be used in conjunction with infFITC for large scale
40 %                  regression problems; any covariance can be wrapped by
41 %                  covFITC such that the FITC approximation is applicable
42 %
43 % Naming convention: all covariance functions are named "cov/cov*.m". A trailing
```

```

44 % "iso" means isotropic, "ard" means Automatic Relevance Determination, and
45 % "one" means that the distance measure is parameterized by a single parameter.
46 %
47 % The covariance functions are written according to a special convention where
48 % the exact behaviour depends on the number of input and output arguments
49 % passed to the function. If you want to add new covariance functions, you
50 % should follow this convention if you want them to work with the function gp.
51 % There are four different ways of calling the covariance functions:
52 %
53 % 1) With no (or one) input argument(s):
54 %
55 %     s = cov
56 %
57 % The covariance function returns a string s telling how many hyperparameters it
58 % expects, using the convention that "D" is the dimension of the input space.
59 % For example, calling "covRQard" returns the string '(D+2)'.
60 %
61 % 2) With two input arguments:
62 %
63 %     K = cov(hyp, x) equivalent to K = cov(hyp, x, [])
64 %
65 % The function computes and returns the covariance matrix where hyp are
66 % the hyperparameters and x is an n by D matrix of cases, where
67 % D is the dimension of the input space. The returned covariance matrix is of
68 % size n by n.
69 %
70 % 3) With three input arguments:
71 %
72 %     Ks = cov(hyp, x, xs)
73 %     kss = cov(hyp, xs, 'diag')
74 %
75 % The function computes test set covariances; kss is a vector of self covariances
76 % for the test cases in xs (of length ns) and Ks is an (n by ns) matrix of cross
77 % covariances between training cases x and test cases xs.
78 %
79 % 4) With four input arguments:
80 %
81 %     dKi = cov(hyp, x, [], i)
82 %     dKsi = cov(hyp, x, xs, i)
83 %     dkssi = cov(hyp, xs, 'diag', i)
84 %
85 % The function computes and returns the partial derivatives of the
86 % covariance matrices with respect to hyp(i), i.e. with
87 % respect to the hyperparameter number i.
88 %
89 % Covariance functions can be specified in two ways: either as a string
90 % containing the name of the covariance function or using a cell array. For
91 % example:
92 %
93 %     cov = 'covRQard';
94 %     cov = {'covRQard'};
95 %     cov = {@covRQard};
96 %
97 % are supported. Only the second and third form using the cell array can be used
98 % for specifying composite covariance functions, made up of several
99 % contributions. For example:
100 %
101 %     cov = {'covScale', {'covRQiso'}};

```

```

102 %      cov = {'covSum', {'covRQiso','covSEard','covNoise'}};
103 %      cov = {'covProd',{'covRQiso','covSEard','covNoise'}};
104 %      cov = {'covMask',{mask,'covSEiso'}}
105 %      q=1; cov = {'covPPiso',q};
106 %      d=3; cov = {'covPoly',d};
107 %      cov = {'covADD',{[1,2],'covSEiso'}};
108 %      cov = {@covFITC, {@covSEiso}, u};  where u are the inducing inputs
109 %
110 % specifies a covariance function which is the sum of three contributions. To
111 % find out how many hyperparameters this covariance function requires, we do:
112 %
113 %      feval(cov{:})
114 %
115 % which returns the string '3+(D+1)+1' (i.e. the 'covRQiso' contribution uses
116 % 3 parameters, the 'covSEard' uses D+1 and 'covNoise' a single parameter).
117 %
118 % See also doc/usageCov.m.
119 %
120 (gpml copyright 5a)

```

## 6.2 Implemented Covariance Functions

Similarly to the mean functions, we provide a whole algebra of covariance functions  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  with the same generic name pattern `cov/cov<NAME>.m` as before.

Besides a long list of simple covariance functions, we also offer a variety of composite covariance functions as shown in the following table.

Simple covariance functions $k(\mathbf{x}, \mathbf{x}')$			
<NAME>	Meaning	$k(\mathbf{x}, \mathbf{x}') =$	$\Psi$
Noise	additive measurement noise	$\sigma_f^2 \delta(\mathbf{x} - \mathbf{x}')$	$\ln \sigma_f$
Const	covariance equals a constant	$\sigma_f^2$	$\ln \sigma_f$
LIN	linear, $\mathcal{X} \subseteq \mathbb{R}^D$	$\mathbf{x}^\top \mathbf{x}'$	$\emptyset$
LINard	linear with diagonal weighting, $\mathcal{X} \subseteq \mathbb{R}^D$	$\mathbf{x}^\top \mathbf{\Lambda}^{-2} \mathbf{x}'$	$\{\ln \lambda_1, \dots, \ln \lambda_D\}$
LINiso	linear with isotropic weighting, $\mathcal{X} \subseteq \mathbb{R}^D$	$\mathbf{x}^\top \mathbf{x}' / \ell^2$	$\ln \ell$
LINone	linear with bias, $\mathcal{X} \subseteq \mathbb{R}^D$	$(\mathbf{x}^\top \mathbf{x}' + 1) / \ell^2$	$\ln \ell$
Poly	polynomial covariance, $\mathcal{X} \subseteq \mathbb{R}^D$	$\sigma_f^2 (\mathbf{x}^\top \mathbf{x}' + c)^d$	$\{\ln c, \ln \sigma_f\}$
SEard	full squared exponential, $\mathcal{X} \subseteq \mathbb{R}^D$	$\sigma_f^2 \exp(-\frac{1}{2}(\mathbf{x} - \mathbf{x}')^\top \mathbf{\Lambda}^{-2}(\mathbf{x} - \mathbf{x}'))$	$\{\ln \lambda_1, \dots, \ln \lambda_D, \ln \sigma_f\}$
SEiso	diagonal squared exponential, $\mathcal{X} \subseteq \mathbb{R}^D$	$\sigma_f^2 \exp(-\frac{1}{2\ell^2}(\mathbf{x} - \mathbf{x}')^\top (\mathbf{x} - \mathbf{x}'))$	$\{\ln \ell, \ln \sigma_f\}$
SEisoU	squared exponential, $\mathcal{X} \subseteq \mathbb{R}^D$	$\exp(-\frac{1}{2\ell^2} \mathbf{x}^\top \mathbf{x}')$	$\ln \ell$
RQard	rational quadratic, $\mathcal{X} \subseteq \mathbb{R}^D$	$\sigma_f^2 (1 + \frac{1}{2\alpha}(\mathbf{x} - \mathbf{x}')^\top \mathbf{\Lambda}^{-2}(\mathbf{x} - \mathbf{x}'))^{-\alpha}$	$\{\ln \lambda_1, \dots, \ln \lambda_D, \ln \sigma_f, \ln \alpha\}$
RQiso	rational quadratic, $\mathcal{X} \subseteq \mathbb{R}^D$	$\sigma_f^2 (1 + \frac{1}{2\alpha\ell^2}(\mathbf{x} - \mathbf{x}')^\top (\mathbf{x} - \mathbf{x}'))^{-\alpha}$	$\{\ln \ell, \ln \sigma_f, \ln \alpha\}$
Maternard	Matérn, $\mathcal{X} \subseteq \mathbb{R}^D$ , $f_1(t) = 1$ , $f_3(t) = 1 + t$ , $f_5(t) = f_3(t) + \frac{t^2}{3}$	$\sigma_f^2 f_d(r_d) \exp(-r_d)$ , $r_d = \sqrt{d}(\mathbf{x} - \mathbf{x}')^\top \mathbf{\Lambda}^{-2}(\mathbf{x} - \mathbf{x}')$	$\{\ln \lambda_1, \dots, \ln \lambda_D, \ln \sigma_f\}$
Materniso	Matérn, $\mathcal{X} \subseteq \mathbb{R}^D$ , $f_1(t) = 1$ , $f_3(t) = 1 + t$ , $f_5(t) = f_3(t) + \frac{t^2}{3}$	$\sigma_f^2 f_d(r_d) \exp(-r_d)$ , $r_d = \sqrt{\frac{d}{\ell^2}}(\mathbf{x} - \mathbf{x}')^\top (\mathbf{x} - \mathbf{x}')$	$\{\ln \ell, \ln \sigma_f\}$
NNone	neural net, $\mathcal{X} \subseteq \mathbb{R}^D$ , $f(\mathbf{x}) = 1 + \mathbf{x}^\top \mathbf{\Lambda}^{-2} \mathbf{x}$	$\sigma_f^2 \sin^{-1} \left( \frac{\mathbf{x}^\top \mathbf{\Lambda}^{-2} \mathbf{x}'}{\sqrt{f(\mathbf{x})f(\mathbf{x}')}} \right)$	$\{\ln \ell, \ln \sigma_f\}$
Periodic	periodic, $\mathcal{X} \subseteq \mathbb{R}$	$\sigma_f^2 \exp \left( -\frac{2}{\ell^2} \sin^2[\pi \  \mathbf{x} - \mathbf{x}' \  / p] \right)$	$\{\ln \ell, \ln p, \ln \sigma_f\}$
PeriodicNoDC	periodic, $\mathcal{X} \subseteq \mathbb{R}$ , rescaled and DC component removed	$\sigma_f^2 \frac{\kappa(\mathbf{x} - \mathbf{x}') - \frac{1}{2} \int_0^{\mathbf{x} - \mathbf{x}'} \kappa(t) dt}{\kappa(0) - \frac{1}{2} \int_0^0 \kappa(t) dt}$ , $\kappa(t) = \exp(-\frac{2}{\ell^2} \sin^2[\pi t / p])$	$\{\ln \ell, \ln p, \ln \sigma_f\}$
Cos	periodic cosine, $\mathcal{X} \subseteq \mathbb{R}$	$\sigma_f^2 \cos(\pi \  \mathbf{x} - \mathbf{x}' \  / p)$	$\{\ln p, \ln \sigma_f\}$
PPard	compact support, piecewise polynomial $f_v(r)$ , $\mathcal{X} \subseteq \mathbb{R}^D$	$\sigma_f^2 \max(0, 1 - r)^{j+v} \cdot f_v(r)$ , $r^2 = (\mathbf{x} - \mathbf{x}')^\top \mathbf{\Lambda}^{-2}(\mathbf{x} - \mathbf{x}')$	$\{\ln \lambda_1, \dots, \ln \lambda_D, \ln \sigma_f\}$
PPiso	compact support, piecewise polynomial $f_v(r)$ , $\mathcal{X} \subseteq \mathbb{R}^D$	$\sigma_f^2 \max(0, 1 - r)^{j+v} \cdot f_v(r)$ , $r = \frac{\  \mathbf{x} - \mathbf{x}' \ }{\ell}$ , $j = \lfloor \frac{D}{2} \rfloor + v + 1$	$\{\ln \ell, \ln \sigma_f\}$
SM	spectral mixture, $\mathcal{X} \subseteq \mathbb{R}^D$ , $\mathbf{w} \in \mathbb{R}_+^Q$ , $\mathbf{M}, \mathbf{V} \in \mathbb{R}_+^{D \times Q}$	$\mathbf{w}^\top \left( \prod_{d=1}^D \exp(-\frac{1}{2} \mathbf{v}_d r^2) \odot \cos(\mathbf{m}_d r) \right)$ , $r = 2\pi \  \mathbf{x} - \mathbf{x}' \ $	$\{\ln \mathbf{w}, \ln \mathbf{M}, \ln \mathbf{V}\}$
	spectral mixture, $\mathcal{X} \subseteq \mathbb{R}^D$ , $\mathbf{W} \in \mathbb{R}_+^{D \times Q}$ , $\mathbf{M}, \mathbf{V} \in \mathbb{R}_+^{D \times Q}$	$\prod_{d=1}^D \mathbf{w}_d \left( \exp(-\frac{1}{2} \mathbf{v}_d r^2) \odot \cos(\mathbf{m}_d r) \right)$ , $r = 2\pi \  \mathbf{x} - \mathbf{x}' \ $	$\{\ln \mathbf{W}, \ln \mathbf{M}, \ln \mathbf{V}\}$
Gaborard	anisotropic Gabor function, $\mathcal{X} \subseteq \mathbb{R}^D$ , $\lambda, \mathbf{p} \in \mathbb{R}_+^D$	$\prod_{d=1}^D \exp(-\frac{r_d^2}{2\lambda_d^2}) \cos(2\pi r_d / p_d)$ , $r_d =  x_d - x'_d $	$\{\ln \lambda, \ln \mathbf{p}\}$
Gaboriso	isotropic Gabor function, $\mathcal{X} \subseteq \mathbb{R}^D$ , $\ell, p \in \mathbb{R}_+$	$\exp(-\frac{r^2}{2\ell^2}) \cos(2\pi r / p)$ , $r = \  \mathbf{x} - \mathbf{x}' \ $	$\{\ln \ell, \ln p\}$
Composite covariance functions $[k_1(\mathbf{x}, \mathbf{x}'), k_2(\mathbf{x}, \mathbf{x}'), \dots] \mapsto k(\mathbf{x}, \mathbf{x}')$			
<NAME>	Meaning	$k(\mathbf{x}, \mathbf{x}') =$	$\Psi$
Scale	scale a covariance	$\alpha k(\mathbf{x}, \mathbf{x}')$	$\alpha \in \mathbb{R}$
Sum	add up covariance functions	$\sum_i k_i(\mathbf{x}, \mathbf{x}')$	$\emptyset$
Prod	multiply covariance functions	$\prod_i k_i(\mathbf{x}, \mathbf{x}')$	$\emptyset$
PERard	turn ARD stationary into a periodic, $\mathcal{X} \subseteq \mathbb{R}^D$	$\kappa(\mathbf{u}(\mathbf{x}), \mathbf{u}(\mathbf{x}'))$ , $\mathbf{u}(\mathbf{x}) = [\sin \mathbf{x}_p, \cos \mathbf{x}_p]$ , $\mathbf{x}_p = 2\pi \text{diag}(\mathbf{p}^{-1}) \mathbf{x}$	$\{\ln p_1, \dots, \ln p_D\}$
PERiso	turn isotropic stationary into a periodic, $\mathcal{X} \subseteq \mathbb{R}^D$	$\kappa(\mathbf{u}(\mathbf{x}), \mathbf{u}(\mathbf{x}'))$ , $\mathbf{u}(\mathbf{x}) = [\sin \mathbf{x}_p, \cos \mathbf{x}_p]$ , $\mathbf{x}_p = 2\pi \mathbf{x} / p$	$\ln p$
Mask	act on components $I \subseteq \{1, 2, \dots, D\}$ of $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^D$ only	$\kappa(\mathbf{x}_I, \mathbf{x}'_I)$	$\emptyset$
ADD	additive, $\mathcal{X} \subseteq \mathbb{R}^D$ , index degree set $\mathcal{D} = \{1, \dots, D\}$	$\sum_{d \in \mathcal{D}} \sigma_{f_d}^2 \sum_{ I =d} \prod_{i \in I} \kappa(x_i, x'_i; \psi_i)$	$\{\psi_1, \dots, \psi_D, \ln \sigma_{f_1}, \dots, \ln \sigma_{f_{ D }}\}$

The spectral mixture covariance `covSM` was introduced by Wilson & Adams *Gaussian Process Kernels for Pattern Discovery and Extrapolation*, ICML, 2013.

The periodic covariance functions `covPERiso` and `covPERard` start from a stationary isotropic or ARD covariance function that depends on the data only through a distance  $r^2 = (\mathbf{x} - \mathbf{x}')^\top \mathbf{\Lambda}^{-2}(\mathbf{x} - \mathbf{x}')$  such as `covMatern*`, `covPP*`, `covRQ*`, `covSE*` where `*=ard|iso` and turn them into a periodic covariance function by embedding the data  $\mathbf{x} \in \mathbb{R}^D$  into a periodic high-dimensional space  $\mathbf{x}_p = \mathbf{u}(\mathbf{x}) \in \mathbb{R}^{2D}$  by a function  $\mathbf{u}(\mathbf{x}) = 2\pi \text{diag}(\mathbf{p}^{-1}) \mathbf{x}$ .

The additive covariance function `covADD` starts from a one-dimensional covariance function  $\kappa(x_i, x'_i; \psi_i)$  acting on a single component  $i \in \{1, \dots, D\}$  of  $\mathbf{x}$ . From that, we define covariance functions  $\kappa_I(\mathbf{x}_I, \mathbf{x}'_I) = \prod_{i \in I} \kappa(x_i, x'_i; \psi_i)$  acting on vector-valued inputs  $\mathbf{x}_I$ . The sums of exponential size can efficiently be computed using the Newton-Girard formulae. Samples functions drawn from a GP with additive covariance are additive functions. The number of interacting variables  $|I|$  is a measure of how complex the additive functions are.

## 6.3 Usage of Implemented Covariance Functions

Some code examples taken from `doc/usageCov.m` illustrate how to use simple and composite covariance functions to specify a GP model.

Syntactically, a covariance function `cf` is defined by

```
cv := 'func' | @func // simple
```

```
cf := {cv} | {cv, {param, cf}} | {cv, {cf, ..., cf}} // composite
```

i.e., it is either a string containing the name of a covariance function, a pointer to a covariance function or one of the former in combination with a cell array of covariance functions and an additional list of parameters.

```
39 <doc/usageCov.m 39>≡
1 % demonstrate usage of covariance functions
2 %
3 % See also covFunctions.m.
4 %
5 <gpml copyright 5a>
6 clear all, close all
7 n = 5; D = 3; x = randn(n,D); xs = randn(3,D); % create a data set
8
9 % set up simple covariance functions
10 cn = {'covNoise'}; sn = .1; hypn = log(sn); % one hyperparameter
11 cc = {@covConst}; sf = 2; hypc = log(sf); % function handles OK
12 cl = {@covLIN}; hyp1 = []; % linear is parameter-free
13 cla = {'covLINard'}; L = rand(D,1); hyp1a = log(L); % linear (ARD)
14 cli = {'covLINiso'}; l = rand(1); hyp1i = log(l); % linear iso
15 clo = {@covLINone}; ell = .9; hyplo = log(ell); % linear with bias
16 cp = {@covPoly,3}; c = 2; hypp = log([c;sf]); % third order poly
17 cga = {@covSEard}; hypga = log([L;sf]); % Gaussian with ARD
18 cgi = {'covSEiso'}; hypgi = log([ell;sf]); % isotropic Gaussian
19 cgu = {'covSEisoU'}; hypgu = log(ell); % isotropic Gauss no scale
20 cra = {'covRQard'}; al = 2; hypra = log([L;sf;al]); % ration. quad.
21 cri = {@covRQiso}; hypri = log([ell;sf;al]); % isotropic
22 cma = {@covMaternard,5}; hypma = log([ell;sf]); % Matern class d=5
23 cmi = {'covMaterniso',3}; hypmi = log([ell;sf]); % Matern class d=3
24 cnn = {'covNNone'}; hypnn = log([L;sf]); % neural network
25 cpe = {'covPeriodic'}; p = 2; hyppe = log([ell;p;sf]); % periodic
26 cpn = {'covPeriodicNoDC'}; p = 2; hyppe = log([ell;p;sf]); % w/o DC
27 cpc = {'covCos'}; p = 2; hypcpc = log([p;sf]); % cosine cov
28 cca = {'covPPard',3}; hypcc = hypm; % compact support poly degree 3
29 cci = {'covPPiso',2}; hypcc = hypm; % compact support poly degree 2
30 cgb = {@covGaboriso}; ell = 1; p = 1.2; hypgb=log([ell;p]); % Gabor
31 csm = {@covSM,4}; hypsm = log([w;m(:);v(:)]); % Spectral Mixture
32
33 % set up composite i.e. meta covariance functions
34 csc = {'covScale',{cgu}}; hypsc = [log(3); hypgu]; % scale by 9
35 csu = {'covSum',{cn,cc,cl}}; hypsu = [hypn; hypc; hyp1]; % sum
36 cpr = {@covProd,{cc,ccc}}; hyppr = [hypc; hypcc]; % product
37 mask = [0,1,0]; % binary mask excluding all but the 2nd component
38 cma = {'covMask',{mask,cgi{:}}}; hypma = hypgi;
39 % isotropic periodic rational quadratic
40 cpi = {'covPERiso',{@covRQiso}};
41 % periodic Matern with ARD
42 cpa = {'covPERard',{@covMaternard,3}};
43 % additive based on SEiso using unary and pairwise interactions
44 cad = {'covADD',{[1,2],'covSEiso'}};
45
46 % 0) specify a covariance function
47 cov = cma; hyp = hypma;
48
49 % 1) query the number of parameters
```

```

50 feval(cov{:})
51
52 % 2) evaluate the function on x
53 feval(cov{:},hyp,x)
54
55 % 3) evaluate the function on x and xs to get cross-terms
56 kss = feval(cov{:},hyp,xs,'diag')
57 Ks  = feval(cov{:},hyp,x,xs)
58
59 % 4) compute the derivatives w.r.t. to hyperparameter i
60 i = 1; feval(cov{:},hyp,x,[],i)

```