

# Bayesian Non-linear Regression

Václav Šmíd

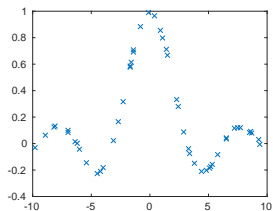
May 5, 2020

# Limit of linear models – basis functions

Fit by a linear combination of basis functions:

$$y_i = a\phi_1(x_i) + b\phi_2(x_i) + e_i$$

What should be the basis function?



# Limit of linear models – basis functions

Fit by a linear combination of basis functions:

$$y_i = a\phi_1(x_i) + b\phi_2(x_i) + e_i$$

What should be the basis function?

Adaptive basis functions:

$$y_i = a\phi_1(\psi_1, x_i) + b\phi_2(\psi_2, x_i) + e_i$$

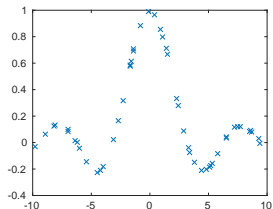
where  $\phi_j$  are non-linear functions with parameters  $\psi_k$ .

Estimating new set of parameters

$$\theta = [a, b, \psi_1, \psi_2]$$

$$\frac{d(y_i - a\phi_1(\psi_1, x_i) + b\phi_2(\psi_2, x_i))^2}{d\theta} = 0.$$

$$\hat{\theta} = ?$$



# Limit of linear models – basis functions

Fit by a linear combination of basis functions:

$$y_i = a\phi_1(x_i) + b\phi_2(x_i) + e_i$$

What should be the basis function?

Adaptive basis functions:

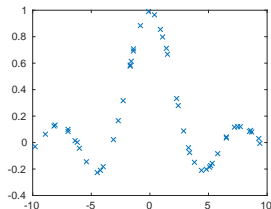
$$y_i = a\phi_1(\psi_1, x_i) + b\phi_2(\psi_2, x_i) + e_i$$

where  $\phi_j$  are non-linear functions with parameters  $\psi_k$ .

Estimating new set of parameters

$$\theta = [a, b, \psi_1, \psi_2]$$

$$\frac{d(y_i - a\phi_1(\psi_1, x_i) + b\phi_2(\psi_2, x_i))^2}{d\theta} = 0.$$
$$\hat{\theta} = ?$$



What is the shape of  $\phi_j$ ?

$$\phi_j(x_i, W, b) = \tanh(Wx_i + b),$$

$$\phi_j(x_i, W, b) = \max(0, Wx_i + b),$$

⋮

# Limit of linear models – basis functions

Fit by a linear combination of basis functions:

$$y_i = a\phi_1(x_i) + b\phi_2(x_i) + e_i$$

What should be the basis function?

Adaptive basis functions:

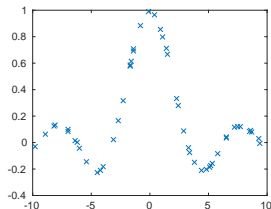
$$y_i = a\phi_1(\psi_1, x_i) + b\phi_2(\psi_2, x_i) + e_i$$

where  $\phi_j$  are non-linear functions with parameters  $\psi_k$ .

Estimating new set of parameters

$$\theta = [a, b, \psi_1, \psi_2]$$

$$\frac{d(y_i - a\phi_1(\psi_1, x_i) + b\phi_2(\psi_2, x_i))^2}{d\theta} = 0.$$
$$\hat{\theta} = ?$$



What is the shape of  $\phi_j$ ?

$$\phi_j(x_i, W, b) = \tanh(Wx_i + b),$$

$$\phi_j(x_i, W, b) = \max(0, Wx_i + b),$$

⋮

Nested functions

$$\tanh(W_1 \tanh(W_2 x_i + b_2) + b_1),$$

# Neural networks

Feed forward NN:

$$z_1 = \sigma_1(W_1 x + b_1),$$

$$z_2 = \sigma_2(W_2 z_1 + b_2), \dots$$

$$y = \sigma_2(w_m z_m + b_m) + e$$

with vector-valued

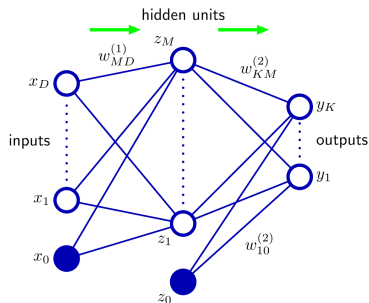
- **activation** functions  $\sigma_j()$ ,
- **weights**  $w_j$
- **biases**  $b_j$ .

For Gaussian noise, maximum log-likelihood is

$$\hat{\theta} = \arg \min \mathcal{L}(x, y, \theta), \quad \mathcal{L} = \mathbf{e}^T \mathbf{e} = \sum_{i=1}^n (y_i - \sigma_1(w_1 \sigma_2(\dots) + b_1))^2.$$

MSE (mean square error) loss function with unknowns

$$\theta = [w_1, b_1, w_2, b_2, \dots].$$



# Neural networks

Feed forward NN:

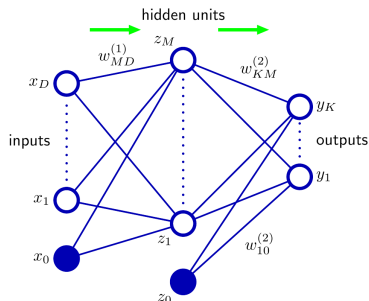
$$z_1 = \sigma_1(W_1 x + b_1),$$

$$z_2 = \sigma_2(W_2 z_1 + b_2), \dots$$

$$y = \sigma_2(w_m z_m + b_m) + e$$

with vector-valued

- **activation** functions  $\sigma_j()$ ,
- **weights**  $w_j$
- **biases**  $b_j$ .



For Gaussian noise, maximum log-likelihood is

$$\hat{\theta} = \arg \min \mathcal{L}(x, y, \theta), \quad \mathcal{L} = \mathbf{e}^T \mathbf{e} = \sum_{i=1}^n (y_i - \sigma_1(w_1 \sigma_2(\dots) + b_1))^2.$$

MSE (mean square error) loss function with unknowns

$\theta = [w_1, b_1, w_2, b_2, \dots]$ . Gradient descent method

$$\hat{\theta}^{(\tau+1)} = \hat{\theta}^{(\tau)} - \eta \nabla \mathcal{L}(\hat{\theta}^{(\tau)}),$$

where  $\eta$  is the (small) learning rate.

# Example

Trivial NN with one hidden layer:

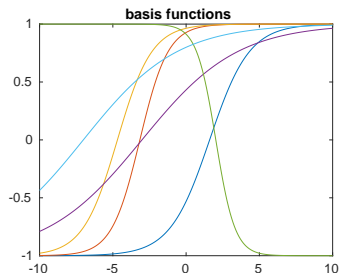
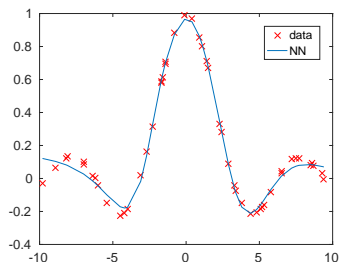
$$y_i = \sum_{i=1}^6 w_{2,i} \tanh(w_{1,j}x_i + b_{1,j}) + b_2,$$

tanh activation function on hidden layer and linear activation function on output.

Training by GD:

1. random initialization,
2. 50000 steps,
3. rate  $\eta = 0.001$ ,

Main issue: reliability, slow convergence,...



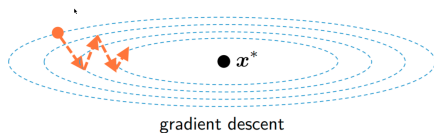


# Faster gradient descent

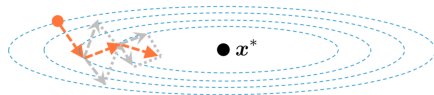
In general, gradient descent requires  $O(1/\epsilon)$  steps

$$\frac{2L(\mathcal{L}\{\hat{\theta}^{(0)}\} - \mathcal{L}\{\theta^*\})}{\epsilon} \leq \tau_{max}$$

where  $L$  is the Lipschitz constant of  $\mathcal{L}$ , for **convex** function.



gradient descent



heavy-ball method

**Heavy-ball (momentum):** accumulate velocity

$$\hat{\theta}^{(\tau+1)} = \hat{\theta}^{(\tau)} - \eta \nabla L(\hat{\theta}^{(\tau)}) + \beta(\hat{\theta}^{(\tau)} - \hat{\theta}^{(\tau-1)})$$

has theoretical asymptotic number of steps  $O(1/\sqrt{\epsilon})$ .

- ▶ Nesterov: theoretically the fastest first-order method. Tuning:  $\eta, \beta$  (via  $L$ ?)

## Second-order: Newton method

Optimize:

$$\hat{\theta} = \arg \min_{\theta} \mathcal{L}(\theta)$$

using Taylor expansion

$$\mathcal{L}(\theta^{(\tau)} + \mathbf{h}) \approx \mathcal{L}(\theta^{(\tau)}) + \nabla \mathcal{L}(\theta^{(\tau)}) \mathbf{h} + \frac{1}{2} \mathbf{h}^T H_{\mathcal{L}}(\theta_k) \mathbf{h}$$

where  $H_{\mathcal{L}}(\theta) = \nabla^2 \mathcal{L}(\theta)$ .

We wish that  $\theta^{(\tau+1)} = \theta^{(\tau)} + \mathbf{h}$  is an optimum, i.e.  $\nabla_{\mathbf{h}} \mathcal{L}(\theta_k + \mathbf{h}) \equiv 0$  :

$$\nabla \mathcal{L}(\theta^{(\tau)}) + H_{\mathcal{L}}(\theta^{(\tau)}) \mathbf{h} = 0 \quad \Leftrightarrow \quad \mathbf{h} = - \left( H_{\mathcal{L}}(\theta^{(\tau)}) \right)^{-1} \nabla \mathcal{L}(\theta^{(\tau)})$$

yielding

$$\theta^{(\tau+1)} = \theta^{(\tau)} - H_{\mathcal{L}}(\theta^{(\tau)})^{-1} \nabla \mathcal{L}(\theta^{(\tau)}).$$

with theoretical asymptotic number of steps  $O(\log(\log \epsilon))$ . (Expensive steps!)

Approximation of the Hessian: **LBFGS**.

# Stochastic Gradient Descent

Original loss function

$$\mathcal{L}(y, x, \theta) = \sum_{i=1}^n (y_i - \sigma_1(w_1 \sigma_2(\dots) + b_1))^2.$$

is replaced by:

$$\tilde{\mathcal{L}}(y, x, \theta) = \sum_{i \in \mathcal{I}} (y_i - \sigma_1(w_1 \sigma_2(\dots) + b_1))^2.$$

where  $\mathcal{I} \subset \{1, \dots, n\}$ ,  $|\mathcal{I}| \ll n$ . For random samples of indices  $j = 1, \dots, m$ ,

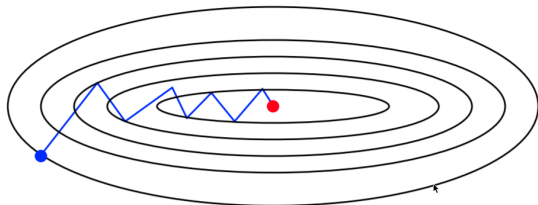
$$\nabla_{\theta} \mathcal{L}(y, x, \theta) = \mathbb{E}(\nabla \tilde{\mathcal{L}}(y, x, \theta))$$

yielding

$$\hat{\theta}^{(\tau+1)} = \hat{\theta}^{(\tau)} - \eta \nabla \tilde{\mathcal{L}}(\hat{\theta}^{(\tau)}),$$

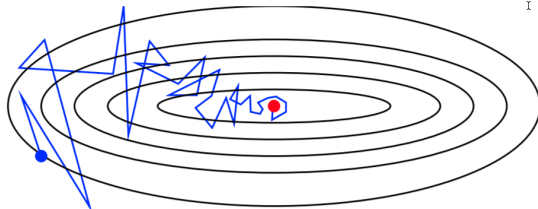
# Stochastic Gradient Descent

**Deterministic gradient:**



**Stochastic gradient:** will converge only if  $\eta_\tau \rightarrow 0$ .

For constant  $\eta_\tau$  it “walks” around optima.



# Adaptive Learning Rate SGD

**AdaGrad** (Duchi, 2011) method uses estimate of the Hessian

$$H_{\mathcal{L}}(\hat{\theta}) \approx \text{diag}(\sqrt{\mathbf{r}_{\tau+1}}),$$
$$\mathbf{r}_{\tau+1} = \mathbf{r}_{\tau} + \left[ \nabla \tilde{\mathcal{L}}(\hat{\theta}^{(\tau)}) \right]^2.$$

accumulates all values from the beginning (infinite window) .

# Adaptive Learning Rate SGD

**AdaGrad** (Duchi, 2011) method uses estimate of the Hessian

$$H_{\mathcal{L}}(\hat{\theta}) \approx \text{diag}(\sqrt{\mathbf{r}_{\tau+1}}),$$
$$\mathbf{r}_{\tau+1} = \mathbf{r}_{\tau} + \left[ \nabla \tilde{\mathcal{L}}(\hat{\theta}^{(\tau)}) \right]^2.$$

accumulates all values from the beginning (infinite window) .

**RMSPprop** (Hinton, 2012) methods adds forgetting

$$H_{\mathcal{L}}(\hat{\theta}) \approx \text{diag}(\sqrt{\mathbf{r}_{\tau+1}}),$$
$$\mathbf{r}_{\tau+1} = \rho \mathbf{r}_{\tau} + (1 - \rho) \left[ \nabla \tilde{\mathcal{L}}(\hat{\theta}^{(\tau)}) \right]^2.$$

# Adaptive Learning Rate SGD

**AdaGrad** (Duchi, 2011) method uses estimate of the Hessian

$$H_{\mathcal{L}}(\hat{\theta}) \approx \text{diag}(\sqrt{\mathbf{r}_{\tau+1}}),$$
$$\mathbf{r}_{\tau+1} = \mathbf{r}_{\tau} + \left[ \nabla \tilde{\mathcal{L}}(\hat{\theta}^{(\tau)}) \right]^2.$$

accumulates all values from the beginning (infinite window) .

**RMSProp** (Hinton, 2012) methods adds forgetting

$$H_{\mathcal{L}}(\hat{\theta}) \approx \text{diag}(\sqrt{\mathbf{r}_{\tau+1}}),$$
$$\mathbf{r}_{\tau+1} = \rho \mathbf{r}_{\tau} + (1 - \rho) \left[ \nabla \tilde{\mathcal{L}}(\hat{\theta}^{(\tau)}) \right]^2.$$

**ADAM** (Kingma&Ba, 2014) combines adaptive rate with adaptive momentum

# Adaptive Learning Rate SGD

**AdaGrad** (Duchi, 2011) method uses estimate of the Hessian

$$H_{\mathcal{L}}(\hat{\theta}) \approx \text{diag}(\sqrt{\mathbf{r}_{\tau+1}}),$$
$$\mathbf{r}_{\tau+1} = \mathbf{r}_{\tau} + \left[ \nabla \tilde{\mathcal{L}}(\hat{\theta}^{(\tau)}) \right]^2.$$

accumulates all values from the beginning (infinite window) .

**RMSProp** (Hinton, 2012) methods adds forgetting

$$H_{\mathcal{L}}(\hat{\theta}) \approx \text{diag}(\sqrt{\mathbf{r}_{\tau+1}}),$$
$$\mathbf{r}_{\tau+1} = \rho \mathbf{r}_{\tau} + (1 - \rho) \left[ \nabla \tilde{\mathcal{L}}(\hat{\theta}^{(\tau)}) \right]^2.$$

**ADAM** (Kingma&Ba, 2014) combines adaptive rate with adaptive momentum

**Bayesian filtering** (Aichison, 2018) explains ADAM, as an extended state in Kalman filter.



# Adaptive Learning Rate SGD

**AdaGrad** (Duchi, 2011) method uses estimate of the Hessian

$$H_{\mathcal{L}}(\hat{\theta}) \approx \text{diag}(\sqrt{\mathbf{r}_{\tau+1}}),$$
$$\mathbf{r}_{\tau+1} = \mathbf{r}_{\tau} + \left[ \nabla \tilde{\mathcal{L}}(\hat{\theta}^{(\tau)}) \right]^2.$$

accumulates all values from the beginning (infinite window) .

**RMSProp** (Hinton, 2012) methods adds forgetting

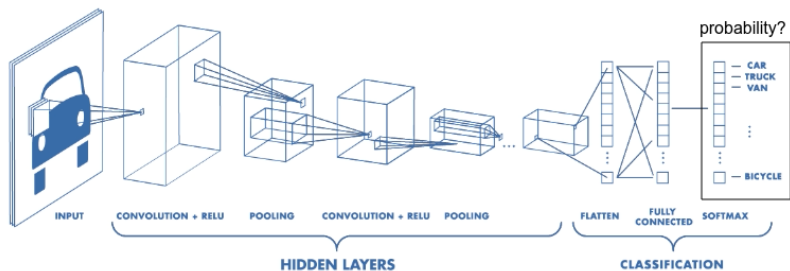
$$H_{\mathcal{L}}(\hat{\theta}) \approx \text{diag}(\sqrt{\mathbf{r}_{\tau+1}}),$$
$$\mathbf{r}_{\tau+1} = \rho \mathbf{r}_{\tau} + (1 - \rho) \left[ \nabla \tilde{\mathcal{L}}(\hat{\theta}^{(\tau)}) \right]^2.$$

**ADAM** (Kingma&Ba, 2014) combines adaptive rate with adaptive momentum

**Bayesian filtering** (Aichison, 2018) explains ADAM, as an extended state in Kalman filter.

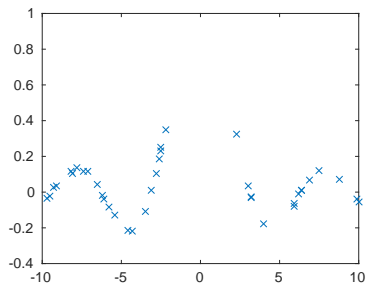
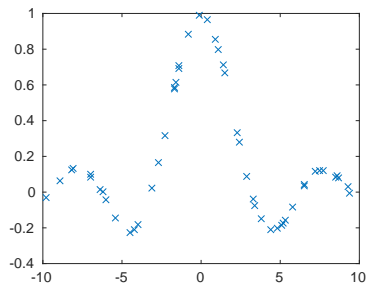
**Controversy:** adaptation can help but can also harm convergence

# Deep Learning

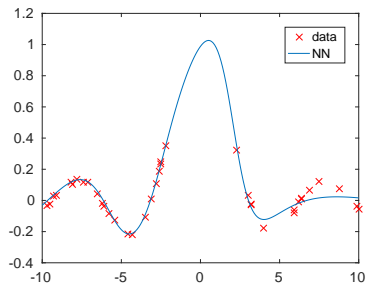
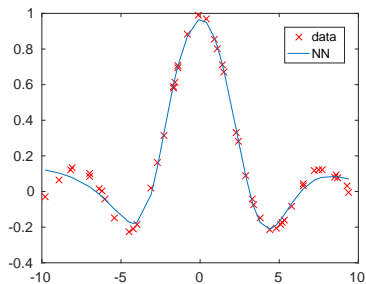
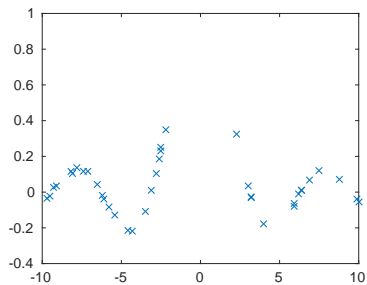
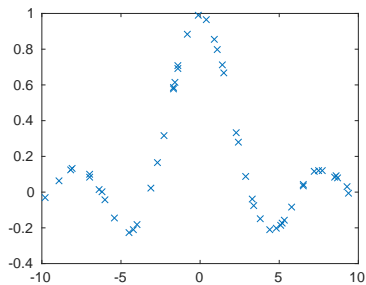


- ▶ Large networks with many layers
- ▶ Special layers that allow to compute gradients
- ▶ Training by a first-order methods
- ▶ Excellent at supervised tasks (regression)

# Can we trust the result?

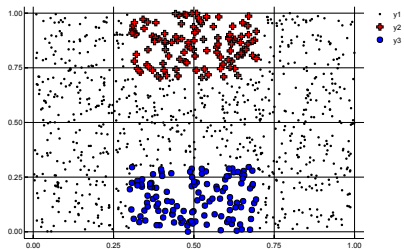


# Can we trust the result?

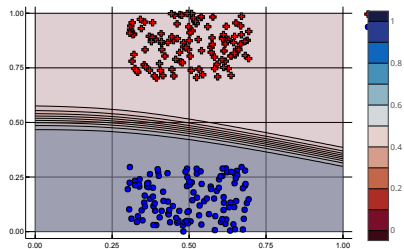


# False confidence in classification

Labeled Data

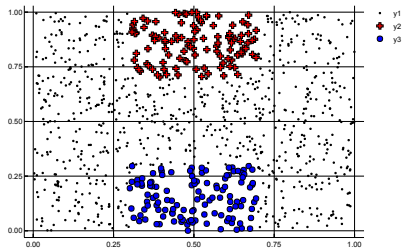


Prediction with Deep-NN

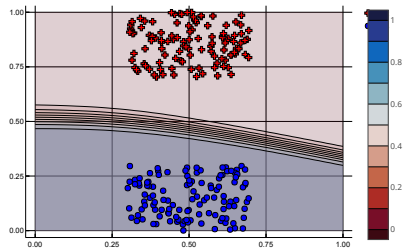


# False confidence in classification

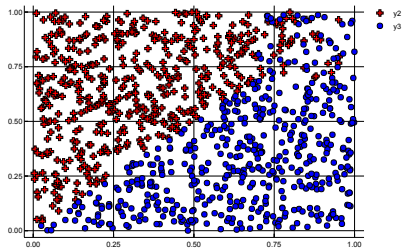
Labeled Data



Prediction with Deep-NN

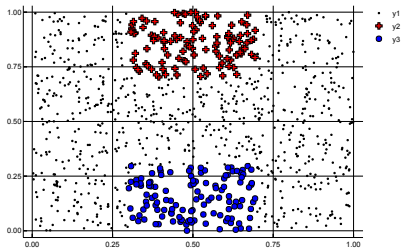


True labels

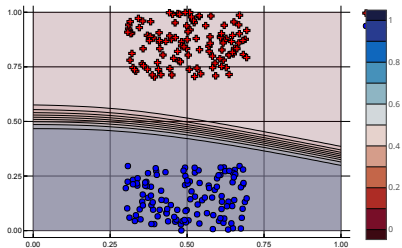


# False confidence in classification

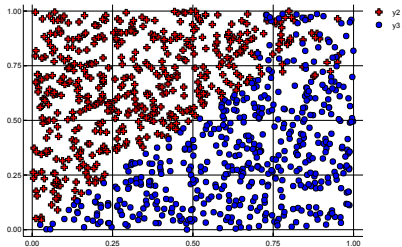
Labeled Data



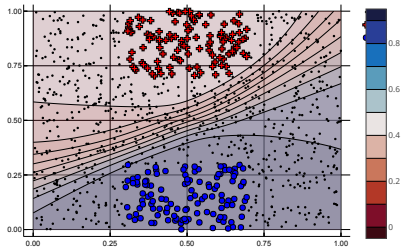
Prediction with Deep-NN



True labels



$\int p(y|\theta, x)p(\theta)d\theta$  from 500 HMC



# Posterior distribution of the weights

Replace parameter estimate  $\hat{\theta}$  by  $p(\theta)$ . Harder than with linear models:



# Posterior distribution of the weights

Replace parameter estimate  $\hat{\theta}$  by  $p(\theta)$ . Harder than with linear models:

**Laplace** build covariance around  $\hat{\theta}$ , (MacKay, 1992):

$$p(\theta) \approx \mathcal{N}(\hat{\theta}, \Sigma), \quad \Sigma = (-\nabla\nabla \log p(\hat{x}))^{-1},$$

# Posterior distribution of the weights

Replace parameter estimate  $\hat{\theta}$  by  $p(\theta)$ . Harder than with linear models:

**Laplace** build covariance around  $\hat{\theta}$ , (MacKay, 1992):

$$p(\theta) \approx \mathcal{N}(\hat{\theta}, \Sigma), \quad \Sigma = (-\nabla\nabla \log p(\hat{x}))^{-1},$$

**HMC** (Metropolis Hastings) uses gradient for solving Hamiltonian ODE

- ▶ golden standard for posterior estimates
- ▶ computationally expensive

# Posterior distribution of the weights

Replace parameter estimate  $\hat{\theta}$  by  $p(\theta)$ . Harder than with linear models:

**Laplace** build covariance around  $\hat{\theta}$ , (MacKay, 1992):

$$p(\theta) \approx \mathcal{N}(\hat{\theta}, \Sigma), \quad \Sigma = (-\nabla\nabla \log p(\hat{x}))^{-1},$$

**HMC** (Metropolis Hastings) uses gradient for solving Hamiltonian ODE

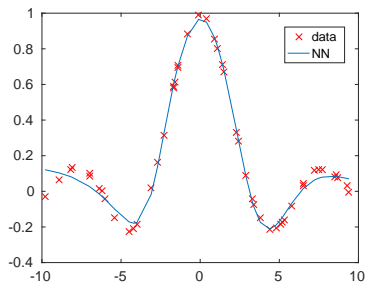
- ▶ golden standard for posterior estimates
- ▶ computationally expensive

**Langevin Dynamics** (Welling, Teh, 2011):

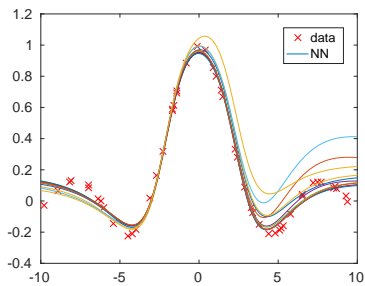
$$\hat{\theta}^{(\tau+1)} = \hat{\theta}^{(\tau)} - \eta \nabla \tilde{L}(\hat{\theta}^{(\tau)}) + e_t, \quad e_t \sim \mathcal{N}(0, \epsilon),$$

where  $\epsilon$  and  $\eta$  needs to be carefully balanced. (Asymptotic proof of acceptance rate=1).

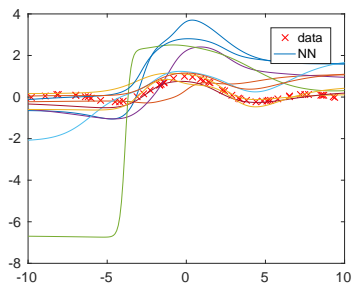
# Laplace



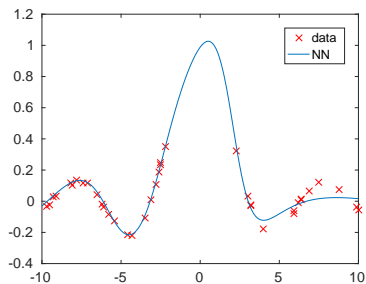
max+0.1std



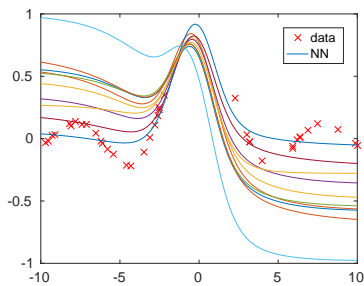
max+std



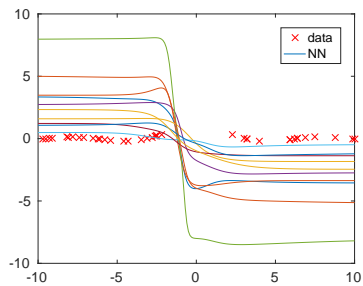
# Laplace II



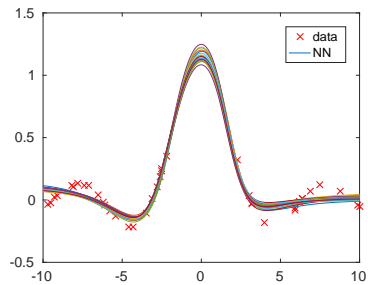
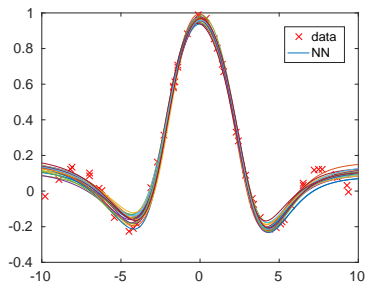
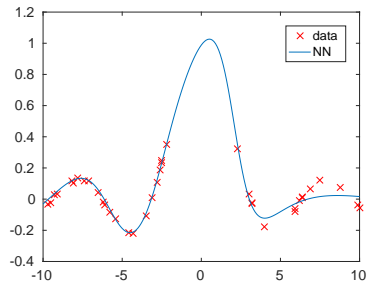
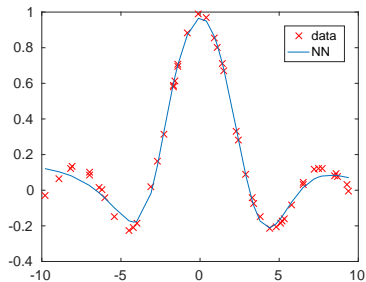
max+0.1std



max+std



# Langevin MCMC (tweaked)



# Uncertainty in the number of basis functions

- ▶ Knowing that the “true” value of the basis functions is 6 is very strong information.
- ▶ How many basis functions for an arbitrary function? Universal approximation property of NN. Infinity.
- ▶ Limit of 2 layer NN is the Gaussian process (Neal, 1994):

$$y(x) \sim GP(0, K(\theta, x, x')),$$

where  $K(x, x')$  is a covariance function of all possible  $x, x'$  pairs, and  $\theta$  is a hyperparameter.

- ▶ For a chosen vector of  $\mathbf{x}$  has Gaussian distribution with

$$y \sim \mathcal{N}(0, K(\theta, \mathbf{x}, \mathbf{x}')).$$

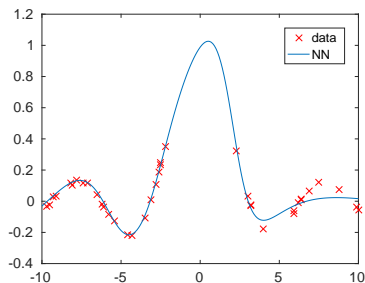
- ▶ Joint process of testing new data  $\bar{x}$  using training data  $x$  (Rasmussen, 2006)

$$\begin{bmatrix} y(x) \\ \bar{y}(\bar{x}) \end{bmatrix} \sim GP \left( 0, \begin{bmatrix} K(x, x') & K(x, \bar{x}') \\ K(\bar{x}, x') & K(\bar{x}, \bar{x}') \end{bmatrix} \right),$$

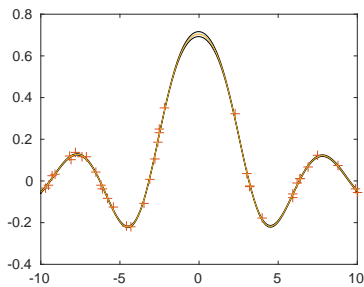
and predictive points are marginal

$$\bar{y}(\bar{x}) \sim \mathcal{N}(K(\bar{x}, x')K(x, x')^{-1}y(x), K(\bar{x}, \bar{x}') - K(\bar{x}, x')K(x, x')^{-1}K(x, \bar{x}'))$$

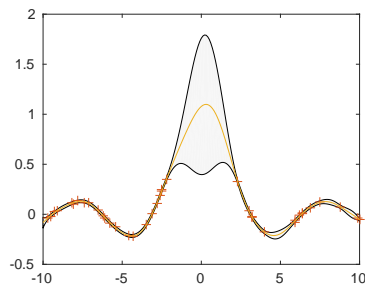
# Gaussian process and its hyperparameters



GP: empirical Bayes kernel



GP: inflated kernel





# Dropout MC

Standard Network Model:

$$z_i = \sigma_i (W_i x + b_i), \quad i = 1 : m - 1,$$
$$y = \sigma_2 (w_m z_m + b_m),$$

Dropout Network Model:

$$z_i = \sigma_i (W_i (\xi_i \circ x) + b_i),$$
$$y = \sigma_2 (w_m (\xi_m \circ z_m) + b_m)$$

where  $\xi_i$  are vectors of zeros and ones sampled from Bernoulli distribution.

Works also for Gaussian distribution, can be explained by Variational Inference.

- ▶ Dropout is an approximation of GP (Gal, Ghahramani, 2016),
- ▶ Deep Neural Networks as Gaussian Processes (Lee, et. al. 2018).

# SGD is Approximate Bayesian Inference

SDG is a discretization of approximation of random walk model

$$\nabla \tilde{\mathcal{L}}(\theta) \approx \nabla \mathcal{L}(\theta) + \frac{1}{\sqrt{S}} \Delta, \quad \Delta \sim \mathcal{N}(0, C(\theta))$$

If the loss function can be approximated by quadratic function

$$\mathcal{L}(\theta) = \frac{1}{2} \theta^\top A \theta,$$

then posterior factor  $q(\theta) = \mathcal{N}(\hat{\theta}, \Sigma)$  satisfies:

$$\Sigma A + A \Sigma = \frac{\eta}{S} C(\theta).$$

Minimizing KL to  $p(\theta)$  yields (Mandt, Hoffman, Blei, 2017):

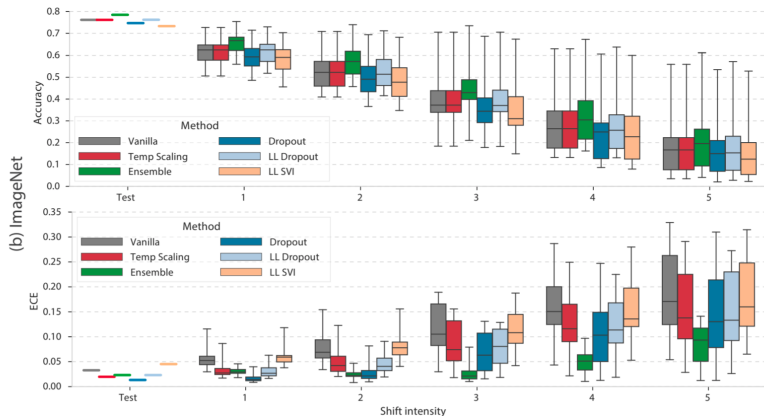
$$\eta^* = \frac{2S \dim(\theta)}{N \operatorname{tr}(C)}, \text{ or } H^* = \frac{2S}{N} C^{-1}, \text{ (matrix learning rate)}$$

Can be used to tune learning rate using

$$C_\tau = (1 - \kappa_\tau) C_{\tau-1} + \kappa_\tau \operatorname{cov}(\nabla \tilde{\mathcal{L}}).$$

# Sad story: Large scale comparison

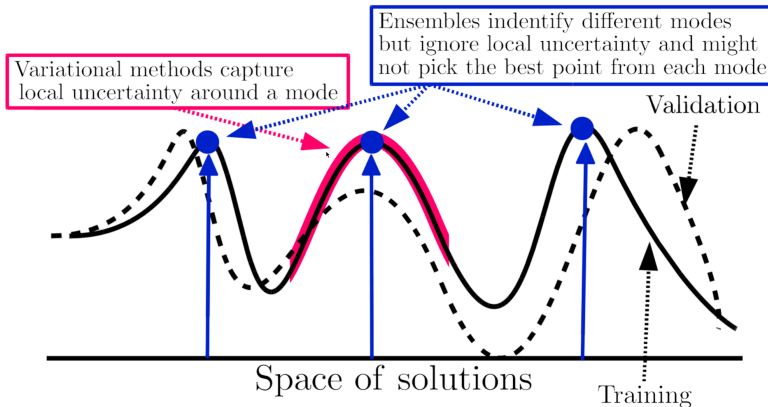
Various methods were compared in (Ovadia, et. al 2019):



The winner is **ensemble**: parallel run of NN from different starts.

# Landscape of Deep networks

Hypothesis (Fort et. al. 2019): The probability of weights in networks is multimodal:



# Assignment:

- ▶ Create a 1d regression problem with missing data
- ▶ train neural network fro minimum loss
- ▶ try one of the Bayesian approaches
  - ▶ Laplace,
  - ▶ Dropout
  - ▶ Ensemble
  - ▶ Langevin
  - ▶ HMC...