Introduction to Nonparametric Bayesian Modeling and Gaussian Process Regression

Piyush Rai

Dept. of CSE, IIT Kanpur

(Mini-course: lecture 3)

Nov 07, 2015

Piyush Rai (IIT Kanpur)

Nonparametric Bayesian Modeling and Gaussian Process Regression

Recap

Piyush Rai (IIT Kanpur)

Nonparametric Bayesian Modeling and Gaussian Process Regression

イロト イロト イヨト イヨト 二日

Optimization vs Inference

All ML problems require estimating parameters given data. Primarily two views:

1. Learning as Optimization

- Parameter θ is a fixed unknown
- Seeks a point estimate (single best answer) for $\boldsymbol{\theta}$

 $\hat{\theta} = \arg\min_{\theta} \text{Loss}(\mathcal{D}; \theta)$ subject to constraints on θ

• Probabilistic methods such as MLE and MAP also fall in this category

2. Learning as (Bayesian) Inference

- Parameter θ is a random variable with a prior distribution $P(\theta)$
- Seeks a posterior distribution over the parameters

$$\frac{P(\theta \mid \mathcal{D})}{P(\mathcal{D})} = \frac{P(\mathcal{D} \mid \theta)P(\theta)}{P(\mathcal{D})}$$

Piyush Rai (IIT Kanpur)

(本部) (本語) (本語) (二語

- $\bullet\,$ Prior distribution specifies our prior belief/knowledge about parameters θ
- Bayesian inference updates the prior and gives the posterior



伺 ト イヨト イヨト

- Prior distribution specifies our prior belief/knowledge about parameters $\boldsymbol{\theta}$
- Bayesian inference updates the prior and gives the posterior



周 ト イ ヨ ト イ ヨ ト

- $\bullet\,$ Prior distribution specifies our prior belief/knowledge about parameters θ
- Bayesian inference updates the prior and gives the posterior



マロト イヨト イヨト

- $\bullet\,$ Prior distribution specifies our prior belief/knowledge about parameters θ
- Bayesian inference updates the prior and gives the posterior



- Prior distribution specifies our prior belief/knowledge about parameters θ
- Bayesian inference updates the prior and gives the posterior



• Posterior $P(\theta|D)$ quantifies uncertainty in the parameters



э

(本間) (本語) (本語)

• Posterior $P(\theta|D)$ quantifies uncertainty in the parameters



• More robust predictions by averaging over the posterior $P(\theta|D)$

$$P(d_{test}|\hat{\theta})$$
 vs $P(d_{test}|\mathcal{D}) = \int P(d_{test}|\theta)P(\theta|\mathcal{D})d\theta$

(人間) シスヨン スヨン

• Posterior $P(\theta|D)$ quantifies uncertainty in the parameters



• More robust predictions by averaging over the posterior $P(\theta|D)$

$$P(d_{test}|\hat{\theta})$$
 vs $P(d_{test}|\mathcal{D}) = \int P(d_{test}|\theta)P(\theta|\mathcal{D})d\theta$

• Allows inferring hyperparameters of the model and doing model comparison

・ 同 ト ・ ヨ ト ・ ヨ ト

• Posterior $P(\theta|D)$ quantifies uncertainty in the parameters



• More robust predictions by averaging over the posterior $P(\theta|D)$

$$P(d_{test}|\hat{\theta})$$
 vs $P(d_{test}|\mathcal{D}) = \int P(d_{test}|\theta)P(\theta|\mathcal{D})d\theta$

- Allows inferring hyperparameters of the model and doing model comparison
- Offers a natural way for informed data acquisition (active learning)
 - Can use the predictive posterior of unseen data points to guide data selection

くぼう くほう くほう

• Posterior $P(\theta|D)$ quantifies uncertainty in the parameters



• More robust predictions by averaging over the posterior $P(\theta|D)$

$$P(d_{test}|\hat{\theta})$$
 vs $P(d_{test}|\mathcal{D}) = \int P(d_{test}|\theta)P(\theta|\mathcal{D})d\theta$

- Allows inferring hyperparameters of the model and doing model comparison
- Offers a natural way for informed data acquisition (active learning)
 - Can use the predictive posterior of unseen data points to guide data selection
- Can do nonparametric Bayesian modeling

Piyush Rai (IIT Kanpur)

< 同 > < 三 > < 三 >

• How big/complex my model should be? How many parameters suffice?



イロト イポト イヨト イヨト

• How big/complex my model should be? How many parameters suffice?



• Model-selection or cross-validation, can often be expensive and impractical

• How big/complex my model should be? How many parameters suffice?



- Model-selection or cross-validation, can often be expensive and impractical
- Nonparametric Bayesian Models: Allow unbounded number of parameters

• How big/complex my model should be? How many parameters suffice?



- Model-selection or cross-validation, can often be expensive and impractical
- Nonparametric Bayesian Models: Allow unbounded number of parameters
 - The model can grow/shrink adaptively as we observe more and more data

イロト イポト イヨト イヨト

• How big/complex my model should be? How many parameters suffice?



- Model-selection or cross-validation, can often be expensive and impractical
- Nonparametric Bayesian Models: Allow unbounded number of parameters
 - The model can grow/shrink adaptively as we observe more and more data
 - We "let the data speak" how complex the model needs to be

Piyush Rai (IIT Kanpur)

・ 同 ト ・ ヨ ト ・ ヨ ト ・

What's a Nonparametric Bayesian Model?

- An NPBayes model is NOT a model with no parameters!
- It has potentially infinite many (unbounded number of) parameters
- It has the ability to "create" new parameters if data requires so..



- A E N A E N

What's a Nonparametric Bayesian Model?

- An NPBayes model is NOT a model with no parameters!
- It has potentially infinite many (unbounded number of) parameters
- It has the ability to "create" new parameters if data requires so..



 Some non-Bayesian models are also nonparametric. For example: nearest neighbor regression/classification, kernel SVMs, kernel density estimation

向下 イヨト イヨト

What's a Nonparametric Bayesian Model?

- An NPBayes model is NOT a model with no parameters!
- It has potentially infinite many (unbounded number of) parameters
- It has the ability to "create" new parameters if data requires so..



- Some non-Bayesian models are also nonparametric. For example: nearest neighbor regression/classification, kernel SVMs, kernel density estimation
- NPBayes models offer the benefits of both Bayesian modeling and nonparametric modeling

Piyush Rai (IIT Kanpur)

周下 イヨト イヨト

Examples of NPBayes Models

Some modeling problems and NPBayes models of choice¹:

Distributions on functions	Gaussian process
Distributions on distributions	Dirichlet process
	Polya Tree
Clustering	Chinese restaurant process
	Pitman-Yor process
Hierarchical clustering	Dirichlet diffusion tree
	Kingman's coalescent
Sparse binary matrices	Indian buffet processes
Survival analysis	Beta processes
Distributions on measures	Completely random measures

¹Table courtesy: Zoubin Ghahramani

Piyush Rai (IIT Kanpur)

. . .

Nonparametric Bayesian Modeling and Gaussian Process Regression

. . .

- 4 回 ト - 4 回 ト

- A Gaussian Process (GP) is a distribution over functions $f: f \sim GP(\mu, \mathbf{\Sigma})$
- .. such that f's value at a finite set of points x_1, \ldots, x_N is jointly Gaussian $\{f(x_1), f(x_2), \ldots, f(x_N)\} \sim \mathcal{N}(\mu, \mathbf{K})$

- A Gaussian Process (GP) is a distribution over functions $f: f \sim GP(\mu, \mathbf{\Sigma})$
- .. such that f's value at a finite set of points x_1, \ldots, x_N is jointly Gaussian $\{f(x_1), f(x_2), \ldots, f(x_N)\} \sim \mathcal{N}(\mu, \mathsf{K})$
- If $\mu = \mathbf{0}$, a GP is fully specified by its covariance (kernel) matrix K

(人間) (人) (人) (人) (日) (日) (日)

- A Gaussian Process (GP) is a distribution over functions $f: f \sim GP(\mu, \mathbf{\Sigma})$
- .. such that f's value at a finite set of points x_1, \ldots, x_N is jointly Gaussian $\{f(x_1), f(x_2), \ldots, f(x_N)\} \sim \mathcal{N}(\mu, \mathsf{K})$
- If $\mu = \mathbf{0}$, a GP is fully specified by its covariance (kernel) matrix K
- Covariance matrix defined by a kernel function $k(x_n, x_m)$. Some examples:

•
$$k(\mathbf{x}_n, \mathbf{x}_m) = \exp\left(-\frac{||\mathbf{x}_n - \mathbf{x}_m||^2}{2\sigma^2}\right)$$
: Gaussian kernel
• $k(\mathbf{x}_n, \mathbf{x}_m) = v_0 \exp\left\{-\left(\frac{|\mathbf{x}_n - \mathbf{x}_m|}{r}\right)^{\alpha}\right\} + v_1 + v_2\delta_{nm}$

(人間) (人) (人) (人) (日) (日) (日)

- A Gaussian Process (GP) is a distribution over functions $f: f \sim GP(\mu, \mathbf{\Sigma})$
- .. such that f's value at a finite set of points x_1, \ldots, x_N is jointly Gaussian $\{f(x_1), f(x_2), \ldots, f(x_N)\} \sim \mathcal{N}(\mu, \mathbf{K})$
- If $\mu = \mathbf{0}$, a GP is fully specified by its covariance (kernel) matrix K
- Covariance matrix defined by a kernel function $k(x_n, x_m)$. Some examples:

•
$$k(\mathbf{x}_n, \mathbf{x}_m) = \exp\left(-\frac{||\mathbf{x}_n - \mathbf{x}_m||^2}{2\sigma^2}\right)$$
: Gaussian kernel
• $k(\mathbf{x}_n, \mathbf{x}_m) = v_0 \exp\left\{-\left(\frac{|\mathbf{x}_n - \mathbf{x}_m|}{r}\right)^{\alpha}\right\} + v_1 + v_2\delta_{nm}$

• GP based modeling also allows learning the kernel hyperparameters from data

Piyush Rai (IIT Kanpur)

イロン イロン イヨン イヨン 三日

Left: some functions drawn from a GP prior $\mathcal{N}(\mathbf{0},\mathbf{K})$

Right: posterior over these functions after observing 5 examples $\{x_n, y_n\}$



通 とう きょう うちょう

• Training data: $\{x_n, y_n\}_{n=1}^N$. Response is a noisy function of the input

$$y_n = f(\boldsymbol{x}_n) + \epsilon_n$$

• Assume a zero-mean Gaussian error

$$p(\epsilon | \sigma^2) = \mathcal{N}(\epsilon | \mathbf{0}, \sigma^2)$$

• Leads to a Gaussian likelihood model for the responses

$$p(y_n|f(\boldsymbol{x}_n)) = \mathcal{N}(y_n|f(\boldsymbol{x}_n), \sigma^2)$$

・ 同 ト ・ ヨ ト ・ ヨ ト …

• Training data: $\{x_n, y_n\}_{n=1}^N$. Response is a noisy function of the input

$$y_n = f(\boldsymbol{x}_n) + \epsilon_n$$

• Assume a zero-mean Gaussian error

$$p(\epsilon | \sigma^2) = \mathcal{N}(\epsilon | 0, \sigma^2)$$

• Leads to a Gaussian likelihood model for the responses

$$p(y_n|f(\boldsymbol{x}_n)) = \mathcal{N}(y_n|f(\boldsymbol{x}_n), \sigma^2)$$

• Denote $\mathbf{y} = [y_1, \dots, y_N]^\top \in \mathbb{R}^N$, $\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]^\top \in \mathbb{R}^N$ and write $p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2 \mathbf{I}_N)$

• Training data: $\{x_n, y_n\}_{n=1}^N$. Response is a noisy function of the input

$$y_n = f(\mathbf{x}_n) + \epsilon_n$$

• Assume a zero-mean Gaussian error

$$p(\epsilon | \sigma^2) = \mathcal{N}(\epsilon | 0, \sigma^2)$$

• Leads to a Gaussian likelihood model for the responses

$$p(y_n|f(\boldsymbol{x}_n)) = \mathcal{N}(y_n|f(\boldsymbol{x}_n), \sigma^2)$$

- Denote $\mathbf{y} = [y_1, \dots, y_N]^\top \in \mathbb{R}^N$, $\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]^\top \in \mathbb{R}^N$ and write $p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2 \mathbf{I}_N)$
- In GP regression, we assume *f* drawn from a GP

$$p(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mathbf{0},\mathbf{K})$$

Piyush Rai (IIT Kanpur)

御下 人居下 人居下 二日

• The likelihood model

$$p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2 \mathbf{I}_N)$$

• The prior distribution

$$p(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mathbf{0},\mathbf{K})$$

• The marginal distribution over the responses y

$$p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f})d\mathbf{f} = \mathcal{N}(\mathbf{y}|\mathbf{0},\sigma^2\mathbf{I}_N+\mathbf{K})$$

• The likelihood model

$$p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2 \mathbf{I}_N)$$

• The prior distribution

$$p(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mathbf{0},\mathbf{K})$$

• The marginal distribution over the responses y

$$p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f})d\mathbf{f} = \mathcal{N}(\mathbf{y}|\mathbf{0},\sigma^2\mathbf{I}_N + \mathbf{K})$$

Marginal and Conditional of Gaussians

If, $p(x) = \mathcal{N}(x|\mu, \Lambda^{-1})$, and $p(y|x) = \mathcal{N}(y|Ax + b, L^{-1})$, then the marginal, p(y), and the conditional p(x|y) distributions are also Gaussians and are given by,

$$\begin{split} p(y) &= \mathcal{N}(y|A\mu + b, L^{-1} + A\Lambda^{-1}A^{\top}), \\ p(x|y) &= \mathcal{N}(x|\Sigma(A^{\top}L(y-b) + \Lambda\mu), \Sigma). \end{split}$$

Piyush Rai (IIT Kanpur)

Making Predictions

• Recall, the marginal distribution over the responses $\mathbf{y} = [y_1, \dots, y_N]$ $p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \sigma^2 \mathbf{I}_N + \mathbf{K}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{C}_N)$

• Adding the response y_{*} of a new test point x_{*}

$$p([\mathbf{y}, y_*]) = \mathcal{N}([\mathbf{y}, y_*]|\mathbf{0}, \mathbf{C}_{N+1})$$

where the $(N+1) \times (N+1)$ matrix C_{N+1} is given by

$$\mathbf{C}_{N+1} = \begin{bmatrix} \mathbf{C} & \mathbf{k}_* \\ \mathbf{k}_*^\top & \mathbf{c} \end{bmatrix}$$

and
$$\mathbf{k}_{*} = [k(\mathbf{x}_{*}, \mathbf{x}_{1}), \dots, k(\mathbf{x}_{*}, \mathbf{x}_{N})], \ \mathbf{c} = k(\mathbf{x}_{*}, \mathbf{x}_{*}) + \sigma^{2}$$

N+1
$$\mathbf{C}_{\mathbf{N+1}} = \mathbf{C}_{\mathbf{N}} \mathbf{k}_{\mathbf{k}}$$

Making Predictions on Test Data

• Recall $p([y, y_*]) = \mathcal{N}([y, y_*]|\mathbf{0}, \mathbf{C}_{N+1})$. The predictive distribution will be

$$p(y_*|\mathbf{y}) = \frac{p([\mathbf{y}, y_*])}{p(\mathbf{y})}$$

$$p(y_*|\mathbf{y}) = \mathcal{N}(y_*|m(\mathbf{x}_*), \sigma^2(\mathbf{x}_*))$$

$$m(\mathbf{x}_*) = \mathbf{k}_*^\top \mathbf{C}_N^{-1} \mathbf{y}$$

$$\sigma^2(\mathbf{x}_*) = c - \mathbf{k}_*^\top \mathbf{C}_N^{-1} \mathbf{k}_*$$

イロン イロン イヨン イヨン 三日

Making Predictions on Test Data

• Recall $p([y, y_*]) = \mathcal{N}([y, y_*]|0, C_{N+1})$. The predictive distribution will be

$$p(y_*|\mathbf{y}) = \frac{p([\mathbf{y}, y_*])}{p(\mathbf{y})}$$

$$p(y_*|\mathbf{y}) = \mathcal{N}(y_*|m(\mathbf{x}_*), \sigma^2(\mathbf{x}_*))$$

$$m(\mathbf{x}_*) = \mathbf{k}_*^\top \mathbf{C}_N^{-1} \mathbf{y}$$

$$\sigma^2(\mathbf{x}_*) = c - \mathbf{k}_*^\top \mathbf{C}_N^{-1} \mathbf{k}_*$$

Partitioned Gaussians

If x_a and x_b are Gaussian variables respectively with means μ_a and μ_b , then the conditional distribution, $p(x_a|x_b)$ is also Gaussian with mean and variance respectively, $\mu_{a|b}$ and $\Sigma_{a|b}$, given by,

$$\mu_{a|b} = \mu_a + \Sigma_{ab} \Sigma_{bb}^{-1} (x_b - \mu_b),$$

$$\Sigma_{a|b} = \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba}.$$

Piyush Rai (IIT Kanpur)

Making Predictions on Test Data

• Recall $p([y, y_*]) = \mathcal{N}([y, y_*]|\mathbf{0}, \mathbf{C}_{N+1})$. The predictive distribution will be

$$p(y_*|\mathbf{y}) = \frac{p([\mathbf{y}, y_*])}{p(\mathbf{y})}$$

$$p(y_*|\mathbf{y}) = \mathcal{N}(y_*|m(\mathbf{x}_*), \sigma^2(\mathbf{x}_*))$$

$$m(\mathbf{x}_*) = \mathbf{k}_*^{\top} \mathbf{C}_N^{-1} \mathbf{y}$$

$$\sigma^2(\mathbf{x}_*) = c - \mathbf{k}_*^{\top} \mathbf{C}_N^{-1} \mathbf{k}_*$$

Partitioned Gaussians

If x_a and x_b are Gaussian variables respectively with means μ_a and μ_b , then the conditional distribution, $p(x_a|x_b)$ is also Gaussian with mean and variance respectively, $\mu_{a|b}$ and $\Sigma_{a|b}$, given by,

$$\mu_{a|b} = \mu_a + \Sigma_{ab} \Sigma_{bb}^{-1} (x_b - \mu_b),$$

$$\Sigma_{a|b} = \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba}.$$

• Note that for GP regression, exact inference is possible at test time!

Piyush Rai (IIT Kanpur)

Interpreting GP predictions..

• Let's look at the predictions made by GP regression

$$p(y_*|\mathbf{y}) = \mathcal{N}(y_*|m(\mathbf{x}_*), \sigma^2(\mathbf{x}_*))$$

$$m(\mathbf{x}_*) = \mathbf{k}_*^{\top} \mathbf{C}_N^{-1} \mathbf{y}$$

$$\sigma^2(\mathbf{x}_*) = c - \mathbf{k}_*^{\top} \mathbf{C}_N^{-1} \mathbf{k}_*$$

- Two interpretations for the mean prediction $m(x_*)$
 - An SVM like interpretation

$$m(\boldsymbol{x}_*) = \boldsymbol{k}_*^{\top} \boldsymbol{\mathsf{C}}_N^{-1} \boldsymbol{y} = \boldsymbol{k}_*^{\top} \boldsymbol{\alpha} = \sum_{n=1}^N k(\boldsymbol{x}_*, \boldsymbol{x}_n) \alpha_n$$

. .

・ 同 ト ・ ヨ ト ・ ヨ ト

where lpha is akin to the weights of support vectors

• A nearest neighbors interpretation

$$m(\boldsymbol{x}_*) = {\boldsymbol{k}_*}^{\top} {\boldsymbol{\mathsf{C}}_N}^{-1} \boldsymbol{y} = {\boldsymbol{w}}^{\top} \boldsymbol{y} = \sum_{n=1}^N w_n y_n$$

where \boldsymbol{w} is akin to the weights of the neighbors

Piyush Rai (IIT Kanpur)

• Recall, the marginal distribution over the responses $\mathbf{y} = [y_1, \dots, y_N]$

$$p(\boldsymbol{y}|\sigma^2, \theta) = \mathcal{N}(\boldsymbol{y}|\boldsymbol{0}, \sigma^2 \boldsymbol{I}_N + \boldsymbol{K}_{\theta})$$

• Can maximize the (log) marginal likelihood w.r.t. σ^2 and the kernel hyperparameterss θ and get point estimates of the hyperparameters

$$\log p(\mathbf{y}|\sigma^2, \theta) = -\frac{1}{2} \log |\sigma^2 \mathbf{I}_N + \mathbf{K}_\theta| - \frac{1}{2} \mathbf{y}^\top (\sigma^2 \mathbf{I}_N + \mathbf{K}_\theta)^{-1} \mathbf{y} + \text{const}$$

(4 回) (4 回) (4 回)

• Recall, the marginal distribution over the responses $\mathbf{y} = [y_1, \dots, y_N]$

$$p(\boldsymbol{y}|\sigma^2, \theta) = \mathcal{N}(\boldsymbol{y}|\boldsymbol{0}, \sigma^2 \boldsymbol{I}_N + \boldsymbol{K}_{\theta})$$

• Can maximize the (log) marginal likelihood w.r.t. σ^2 and the kernel hyperparameterss θ and get point estimates of the hyperparameters

$$\log p(\mathbf{y}|\sigma^2, \theta) = -\frac{1}{2} \log |\sigma^2 \mathbf{I}_N + \mathbf{K}_\theta| - \frac{1}{2} \mathbf{y}^\top (\sigma^2 \mathbf{I}_N + \mathbf{K}_\theta)^{-1} \mathbf{y} + \text{const}$$

• Note: Can also put hyperpriors on the hyperparameters and infer the hyperparameters in a fully Bayesian manner

- 本間 と 本臣 と 本臣 と 二臣

Gaussian Process Classification

Binary classification problem: Given a data set $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, with binary class labels $y_i \in \{-1, +1\}$, infer class label probabilities at new points.



There are many ways to relate function values $f_i = f(\mathbf{x}_i)$ to class probabilities:

$$p(y_i|f_i) = \begin{cases} \frac{1}{1 + \exp(-y_i f_i)} & \text{sigmoid (logistic)} \\ \Phi(y_i f_i) & \text{cumulative normal (probit} \\ H(y_i f_i) & \text{threshold} \\ \epsilon + (1 - 2\epsilon)H(y_i f_i) & \text{robust threshold} \end{cases}$$

Non-Gaussian likelihood, so we need to use approximate inference methods (Laplace, EP, MCMC).

• Non-binary labels (multiclass, counts, etc.) can also be easily handled

Piyush Rai (IIT Kanpur)

伺下 イヨト イヨト

• The objective function of a soft-margin SVM looks like

$$\frac{1}{2}||\bm{w}||^2 + C\sum_{n=1}^N (1-y_n f_n)_+$$

where $f_n = \mathbf{w}^\top \mathbf{x}_n$ and y_n is the true label for \mathbf{x}_n

• The objective function of a soft-margin SVM looks like

$$\frac{1}{2}||\bm{w}||^2 + C\sum_{n=1}^N (1-y_n f_n)_+$$

where $f_n = \mathbf{w}^\top \mathbf{x}_n$ and y_n is the true label for \mathbf{x}_n

• Kernel SVM: $f_n = \sum_{m=1}^N \alpha_m k(\mathbf{x}_n, \mathbf{x}_m)$. Denote $\mathbf{f} = [f_1, \dots, f_N]^\top$

• The objective function of a soft-margin SVM looks like

$$\frac{1}{2}||\bm{w}||^2 + C\sum_{n=1}^N (1-y_n f_n)_+$$

where $f_n = \boldsymbol{w}^\top \boldsymbol{x}_n$ and y_n is the true label for \boldsymbol{x}_n

• Kernel SVM: $f_n = \sum_{m=1}^N \alpha_m k(\mathbf{x}_n, \mathbf{x}_m)$. Denote $\mathbf{f} = [f_1, \dots, f_N]^\top$

• We can write $\frac{||\boldsymbol{w}||^2}{2} = \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha} = \mathbf{f}^\top \mathbf{K}^{-1} \mathbf{f}$, and kernel SVM objective becomes

$$\frac{1}{2}\mathbf{f}^{\top}\mathbf{K}^{-1}\mathbf{f} + C\sum_{n=1}^{N}(1-y_nf_n)_+$$

• The objective function of a soft-margin SVM looks like

$$\frac{1}{2}||\bm{w}||^2 + C\sum_{n=1}^N (1-y_n f_n)_+$$

where $f_n = \mathbf{w}^\top \mathbf{x}_n$ and y_n is the true label for \mathbf{x}_n

• Kernel SVM: $f_n = \sum_{m=1}^N \alpha_m k(\mathbf{x}_n, \mathbf{x}_m)$. Denote $\mathbf{f} = [f_1, \dots, f_N]^\top$

• We can write $\frac{||\boldsymbol{w}||^2}{2} = \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha} = \mathbf{f}^\top \mathbf{K}^{-1} \mathbf{f}$, and kernel SVM objective becomes

$$\frac{1}{2}\mathbf{f}^{\top}\mathbf{K}^{-1}\mathbf{f}+C\sum_{n=1}^{N}(1-y_{n}f_{n})_{+}$$

Negative log of the likelihood p(f|X) of a GP can be written as

$$\frac{1}{2}\mathbf{f}^{\top}\mathbf{K}^{-1}\mathbf{f} - \sum_{n=1}^{N}\log p(y_n|f_n) + \text{const}$$

Piyush Rai (IIT Kanpur)

• Thus GPs can be interpreted as a Bayesian analogue of kernel SVMs

イロン イヨン イヨン イヨン

- Thus GPs can be interpreted as a Bayesian analogue of kernel SVMs
- Both GP and SVM need dealing with (storing/inverting) large kernel matrices
 - Various approximations proposed to address this issue (applicable to both)

- Thus GPs can be interpreted as a Bayesian analogue of kernel SVMs
- Both GP and SVM need dealing with (storing/inverting) large kernel matrices
 - Various approximations proposed to address this issue (applicable to both)
- Ability to learn the kernel hyperparameters in GP is very useful, e.g.,

- Thus GPs can be interpreted as a Bayesian analogue of kernel SVMs
- Both GP and SVM need dealing with (storing/inverting) large kernel matrices
 - Various approximations proposed to address this issue (applicable to both)
- Ability to learn the kernel hyperparameters in GP is very useful, e.g.,
 - · Learning the kernel bandwidth for Gaussian kernels

$$k(\mathbf{x}_n, \mathbf{x}_m) = \exp\left(-\frac{||\mathbf{x}_n - \mathbf{x}_m||^2}{2\sigma^2}\right)$$

- Thus GPs can be interpreted as a Bayesian analogue of kernel SVMs
- Both GP and SVM need dealing with (storing/inverting) large kernel matrices
 - Various approximations proposed to address this issue (applicable to both)
- Ability to learn the kernel hyperparameters in GP is very useful, e.g.,
 - · Learning the kernel bandwidth for Gaussian kernels

$$k(\mathbf{x}_n, \mathbf{x}_m) = \exp\left(-\frac{||\mathbf{x}_n - \mathbf{x}_m||^2}{2\sigma^2}\right)$$

• Doing feature selection (via Automatic Relevance Determination)

$$k(\mathbf{x}_n, \mathbf{x}_m) = \exp\left(-\sum_{d=1}^{D} \frac{(\mathbf{x}_{nd} - \mathbf{x}_{md})^2}{2\sigma_d^2}\right)$$

Piyush Rai (IIT Kanpur)

(人間) システン イラン

- Thus GPs can be interpreted as a Bayesian analogue of kernel SVMs
- Both GP and SVM need dealing with (storing/inverting) large kernel matrices
 - Various approximations proposed to address this issue (applicable to both)
- Ability to learn the kernel hyperparameters in GP is very useful, e.g.,
 - · Learning the kernel bandwidth for Gaussian kernels

$$k(\mathbf{x}_n, \mathbf{x}_m) = \exp\left(-\frac{||\mathbf{x}_n - \mathbf{x}_m||^2}{2\sigma^2}\right)$$

• Doing feature selection (via Automatic Relevance Determination)

$$k(\mathbf{x}_n, \mathbf{x}_m) = \exp\left(-\sum_{d=1}^{D} \frac{(\mathbf{x}_{nd} - \mathbf{x}_{md})^2}{2\sigma_d^2}\right)$$

Learning compositions of kernels for more flexible modeling

$$\mathbf{K} = \mathbf{K}_{\theta_1} + \mathbf{K}_{\theta_2} + \dots$$

Piyush Rai (IIT Kanpur)

- 4 同 ト 4 国 ト - 4 国 ト

- Nonlinear Dimensionality Reduction: Gaussian Process Latent Variable Models
- **Bayesian Optimization:** Optimizing functions that have an unknown functional form and are expensive to evaluate
- **Deep Gaussian Processes:** Data assumed to be an output of a multivariate GP, inputs to each GP are outputs of another GP, and so on..
- Many applications: Robotics and control, vision, spatial statistics, and so on..

Resources on Gaussian Processes

• Book: Gaussian Processes for Machine Learning (freely available online)



- MATLAB Packages: Useful to play with, build applications, extend existing models and inference algorithms for GPs (both regression and classification)
 - GPML: http://www.gaussianprocess.org/gpml/code/matlab/doc/
 - GPStuff: http://research.cs.aalto.fi/pml/software/gpstuff/

Piyush Rai (IIT Kanpur)

- 4 同 6 4 日 6 4 日 6

- Nonparametric Bayesian models for mixture modeling (clustering): Dirichlet Processes and Chinese Restaurant Process
- Nonparametric Bayesian models for latent factor modeling (dimensionality reduction): Beta Processes and Indian Buffet Process

< ロ > < 同 > < 回 > < 回 > < 回 > <

Thanks! Questions?

э

・ロト ・個ト ・ヨト ・ヨト