

Applications of Exponential Families

Thanks to Yasemin Altun, Thomas Hofmann, Vishy Vishwanathan

Alexander J. Smola

Statistical Machine Learning Program
Canberra, ACT 0200 Australia
Alex.Smola@nicta.com.au

ICONIP 2006, Hong Kong, October 3

- 1 Conditional Models**
 - Log-partition Function, Densities, and Expectations
 - Inner Products and Kernels
 - Examples of Kernels
- 2 Gaussian Process Classification
 - Feature map
 - Examples
- 3 Gaussian Process Regression
 - Homoscedastic Model
 - Heteroscedastic Model
- 4 Conditional Random Fields
 - Model Structure
 - Kernel Expansion
 - Connections to Hidden Markov Models

Conditional Models

Conditional Density

$$p(x|\theta) = \exp(\langle \phi(x), \theta \rangle - g(\theta))$$

$$p(y|x, \theta) = \exp(\langle \phi(x, y), \theta \rangle - g(\theta|x))$$

Log-partition function

$$g(\theta|x) = \log \int_y \exp(\langle \phi(x, y), \theta \rangle) dy$$

Sufficient Criterion

$p(x, y|\theta)$ is a member of the exponential family itself.

Key Idea

Avoid computing $\phi(x, y)$ directly, only evaluate inner products

$$k((x, y), (x', y')) := \langle \phi(x, y), \phi(x', y') \rangle$$

Conditional Models

Conditional Density

$$p(x|\theta) = \exp(\langle \phi(x), \theta \rangle - g(\theta))$$

$$p(y|x, \theta) = \exp(\langle \phi(x, y), \theta \rangle - g(\theta|x))$$

Log-partition function

$$g(\theta|x) = \log \int_y \exp(\langle \phi(x, y), \theta \rangle) dy$$

Sufficient Criterion

$p(x, y|\theta)$ is a member of the exponential family itself.

Key Idea

Avoid computing $\phi(x, y)$ directly, only evaluate inner products

$$k((x, y), (x', y')) := \langle \phi(x, y), \phi(x', y') \rangle$$

Conditional Models

Conditional Density

$$p(x|\theta) = \exp(\langle \phi(x), \theta \rangle - g(\theta))$$

$$p(y|x, \theta) = \exp(\langle \phi(x, y), \theta \rangle - g(\theta|x))$$

Log-partition function

$$g(\theta|x) = \log \int_y \exp(\langle \phi(x, y), \theta \rangle) dy$$

Sufficient Criterion

$p(x, y|\theta)$ is a member of the exponential family itself.

Key Idea

Avoid computing $\phi(x, y)$ directly, only evaluate inner products

$$k((x, y), (x', y')) := \langle \phi(x, y), \phi(x', y') \rangle$$

Conditional Models

Conditional Density

$$p(x|\theta) = \exp(\langle \phi(x), \theta \rangle - g(\theta))$$

$$p(y|x, \theta) = \exp(\langle \phi(x, y), \theta \rangle - g(\theta|x))$$

Log-partition function

$$g(\theta|x) = \log \int_y \exp(\langle \phi(x, y), \theta \rangle) dy$$

Sufficient Criterion

$p(x, y|\theta)$ is a member of the exponential family itself.

Key Idea

Avoid computing $\phi(x, y)$ directly, only evaluate inner products

$$k((x, y), (x', y')) := \langle \phi(x, y), \phi(x', y') \rangle$$

Conditional Distributions

Maximum a Posteriori Estimation

$$-\log p(\theta|X) = \sum_{i=1}^m -\langle \phi(x_i), \theta \rangle + mg(\theta) + \frac{1}{2\sigma^2} \|\theta\|^2 + c$$

$$-\log p(\theta|X, Y) = \sum_{i=1}^m -\langle \phi(x_i, y_i), \theta \rangle + g(\theta|x_i) + \frac{1}{2\sigma^2} \|\theta\|^2 + c$$

Solving the Problem

- The problem is strictly convex in θ .
- Direct solution impossible if we cannot compute $\phi(x, y)$.
- Solve convex problem in expansion coefficients.
- Expand θ in a linear combination of $\phi(x_i, y)$.

Conditional Distributions

Maximum a Posteriori Estimation

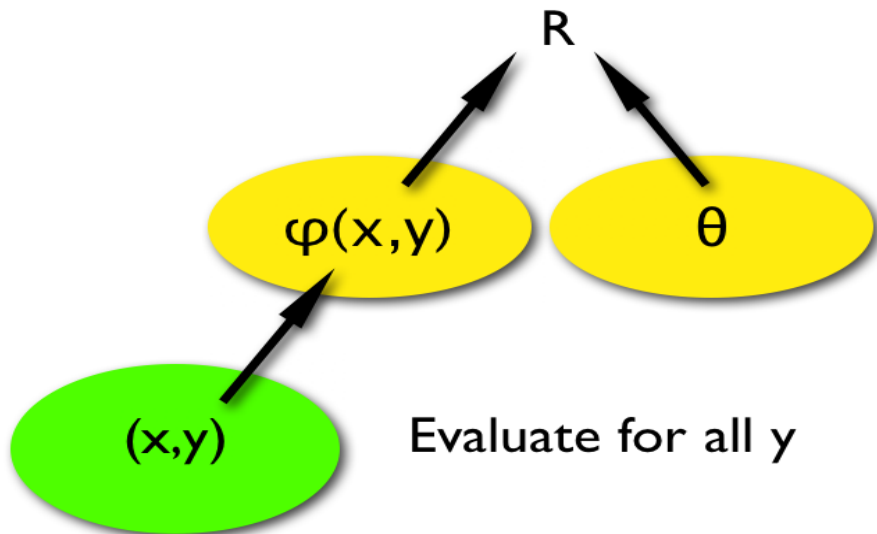
$$-\log p(\theta|X) = \sum_{i=1}^m -\langle \phi(x_i), \theta \rangle + mg(\theta) + \frac{1}{2\sigma^2} \|\theta\|^2 + c$$

$$-\log p(\theta|X, Y) = \sum_{i=1}^m -\langle \phi(x_i, y_i), \theta \rangle + g(\theta|x_i) + \frac{1}{2\sigma^2} \|\theta\|^2 + c$$

Solving the Problem

- The problem is strictly convex in θ .
- Direct solution impossible if we cannot compute $\phi(x, y)$.
- Solve convex problem in expansion coefficients.
- Expand θ in a linear combination of $\phi(x_i, y)$.

Joint Feature Map



Representer Theorem

Objective Function

$$-\log p(\theta|X, Y) = \sum_{i=1}^m -\langle \phi(x_i, y_i), \theta \rangle + g(\theta|x_i) + \frac{1}{2\sigma^2} \|\theta\|^2 + c$$

Decomposition

- Decompose θ into $\theta = \theta_{\parallel} + \theta_{\perp}$ where

$$\theta_{\parallel} \in \text{span}\{\phi(x_i, y) \text{ where } 1 \leq i \leq m \text{ and } y \in \mathcal{Y}\}$$

- Both $g(\theta|x_i)$ and $\langle \phi(x_i, y_i), \theta \rangle$ are independent of θ_{\perp} .

Theorem

$-\log p(\theta|X, Y)$ is minimized for $\theta_{\perp} = 0$, hence $\theta = \theta_{\parallel}$.

Corollary

If $|\mathcal{Y}| < \infty$ we have a parametric optimization problem.

Representer Theorem

Objective Function

$$-\log p(\theta|X, Y) = \sum_{i=1}^m -\langle \phi(x_i, y_i), \theta \rangle + g(\theta|x_i) + \frac{1}{2\sigma^2} \|\theta\|^2 + c$$

Decomposition

- Decompose θ into $\theta = \theta_{\parallel} + \theta_{\perp}$ where

$$\theta_{\parallel} \in \text{span}\{\phi(x_i, y) \text{ where } 1 \leq i \leq m \text{ and } y \in \mathcal{Y}\}$$

- Both $g(\theta|x_i)$ and $\langle \phi(x_i, y_i), \theta \rangle$ are independent of θ_{\perp} .

Theorem

$-\log p(\theta|X, Y)$ is minimized for $\theta_{\perp} = 0$, hence $\theta = \theta_{\parallel}$.

Corollary

If $|\mathcal{Y}| < \infty$ we have a *parametric optimization problem*.

Representer Theorem

Objective Function

$$-\log p(\theta|X, Y) = \sum_{i=1}^m -\langle \phi(x_i, y_i), \theta \rangle + g(\theta|x_i) + \frac{1}{2\sigma^2} \|\theta\|^2 + c$$

Decomposition

- Decompose θ into $\theta = \theta_{\parallel} + \theta_{\perp}$ where

$$\theta_{\parallel} \in \text{span}\{\phi(x_i, y) \text{ where } 1 \leq i \leq m \text{ and } y \in \mathcal{Y}\}$$

- Both $g(\theta|x_i)$ and $\langle \phi(x_i, y_i), \theta \rangle$ are independent of θ_{\perp} .

Theorem

$-\log p(\theta|X, Y)$ is minimized for $\theta_{\perp} = 0$, hence $\theta = \theta_{\parallel}$.

Corollary

If $|\mathcal{Y}| < \infty$ we have a *parametric optimization problem*.

Representer Theorem

Objective Function

$$-\log p(\theta|X, Y) = \sum_{i=1}^m -\langle \phi(x_i, y_i), \theta \rangle + g(\theta|x_i) + \frac{1}{2\sigma^2} \|\theta\|^2 + c$$

Decomposition

- Decompose θ into $\theta = \theta_{\parallel} + \theta_{\perp}$ where

$$\theta_{\parallel} \in \text{span}\{\phi(x_i, y) \text{ where } 1 \leq i \leq m \text{ and } y \in \mathcal{Y}\}$$

- Both $g(\theta|x_i)$ and $\langle \phi(x_i, y_i), \theta \rangle$ are independent of θ_{\perp} .

Theorem

$-\log p(\theta|X, Y)$ is minimized for $\theta_{\perp} = 0$, hence $\theta = \theta_{\parallel}$.

Corollary

If $|\mathcal{Y}| < \infty$ we have a *parametric optimization problem*.

Using It

Expansion

$$\theta = \sum_{i=1}^m \sum_{y \in \mathcal{Y}} \alpha_{iy} \phi(\mathbf{x}_i, y)$$

Inner Product

$$\langle \phi(x, y), \theta \rangle = \sum_{i=1}^m \sum_{y \in \mathcal{Y}} \alpha_{iy} k((x, y), (x_i, y))$$

Norm

$$\|\theta\|^2 = \sum_{i,j=1}^m \sum_{y,y' \in \mathcal{Y}} \alpha_{iy} \alpha_{jy'} k((x_i, y), (x_j, y'))$$

Log-partition function

$$g(\theta|x) = \log \sum_{y \in \mathcal{Y}} \exp(\langle \phi(x, y), \theta \rangle)$$

Using It

Expansion

$$\theta = \sum_{i=1}^m \sum_{y \in \mathcal{Y}} \alpha_{iy} \phi(\mathbf{x}_i, y)$$

Inner Product

$$\langle \phi(\mathbf{x}, y), \theta \rangle = \sum_{i=1}^m \sum_{y' \in \mathcal{Y}} \alpha_{iy'} k((\mathbf{x}, y), (\mathbf{x}_i, y'))$$

Norm

$$\|\theta\|^2 = \sum_{i,j=1}^m \sum_{y,y' \in \mathcal{Y}} \alpha_{iy} \alpha_{jy'} k((\mathbf{x}_i, y), (\mathbf{x}_j, y'))$$

Log-partition function

$$g(\theta | \mathbf{x}) = \log \sum_{y \in \mathcal{Y}} \exp(\langle \phi(\mathbf{x}, y), \theta \rangle)$$

Using It

Expansion

$$\theta = \sum_{i=1}^m \sum_{y \in \mathcal{Y}} \alpha_{iy} \phi(\mathbf{x}_i, y)$$

Inner Product

$$\langle \phi(\mathbf{x}, y), \theta \rangle = \sum_{i=1}^m \sum_{y' \in \mathcal{Y}} \alpha_{iy'} k((\mathbf{x}, y), (\mathbf{x}_i, y'))$$

Norm

$$\|\theta\|^2 = \sum_{i,j=1}^m \sum_{y,y' \in \mathcal{Y}} \alpha_{iy} \alpha_{jy'} k((\mathbf{x}_i, y), (\mathbf{x}_j, y'))$$

Log-partition function

$$g(\theta | \mathbf{x}) = \log \sum_{y \in \mathcal{Y}} \exp(\langle \phi(\mathbf{x}, y), \theta \rangle)$$

Using It

Expansion

$$\theta = \sum_{i=1}^m \sum_{y \in \mathcal{Y}} \alpha_{iy} \phi(\mathbf{x}_i, y)$$

Inner Product

$$\langle \phi(\mathbf{x}, y), \theta \rangle = \sum_{i=1}^m \sum_{y' \in \mathcal{Y}} \alpha_{iy'} k((\mathbf{x}, y), (\mathbf{x}_i, y'))$$

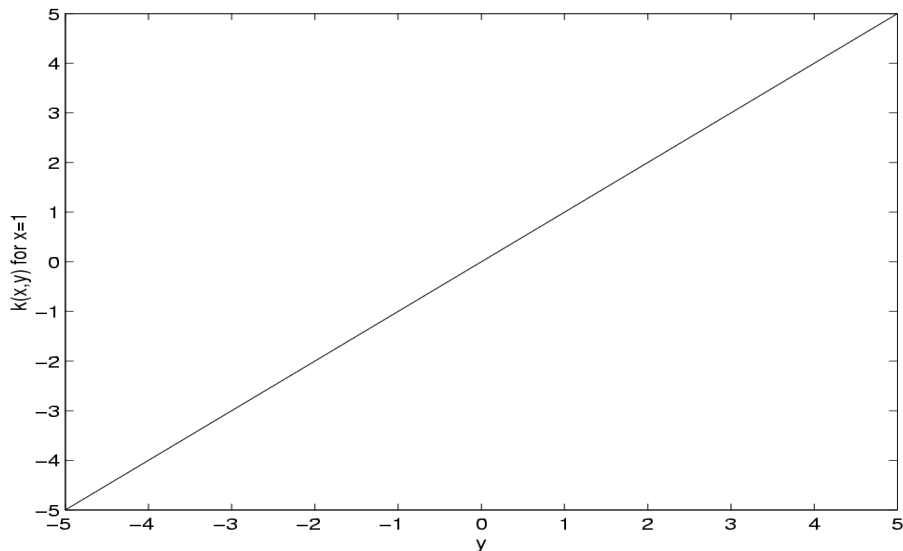
Norm

$$\|\theta\|^2 = \sum_{i,j=1}^m \sum_{y,y' \in \mathcal{Y}} \alpha_{iy} \alpha_{jy'} k((\mathbf{x}_i, y), (\mathbf{x}_j, y'))$$

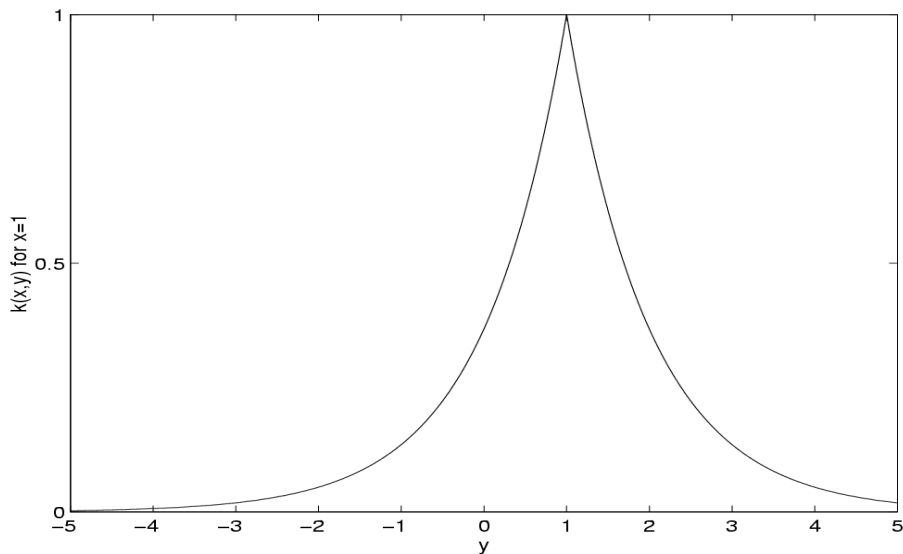
Log-partition function

$$g(\theta | \mathbf{x}) = \log \sum_{y \in \mathcal{Y}} \exp(\langle \phi(\mathbf{x}, y), \theta \rangle)$$

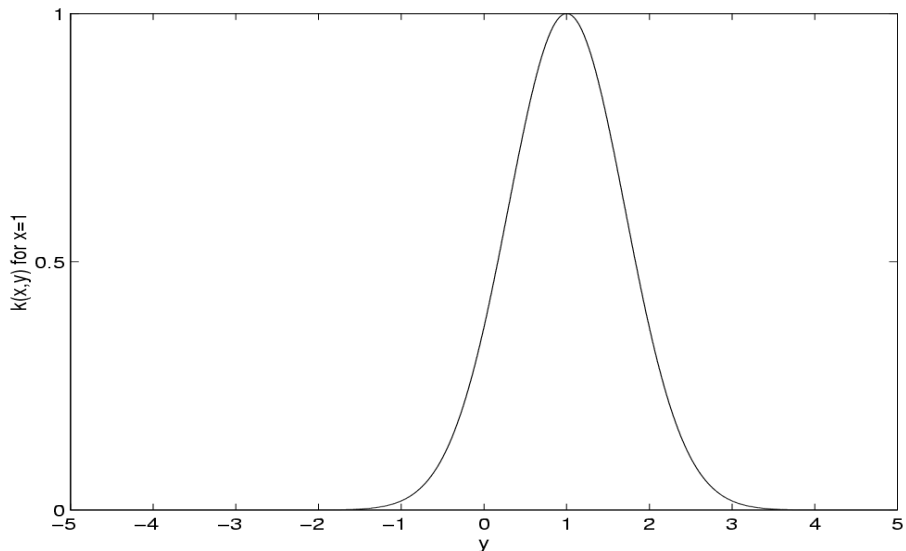
Linear Kernel



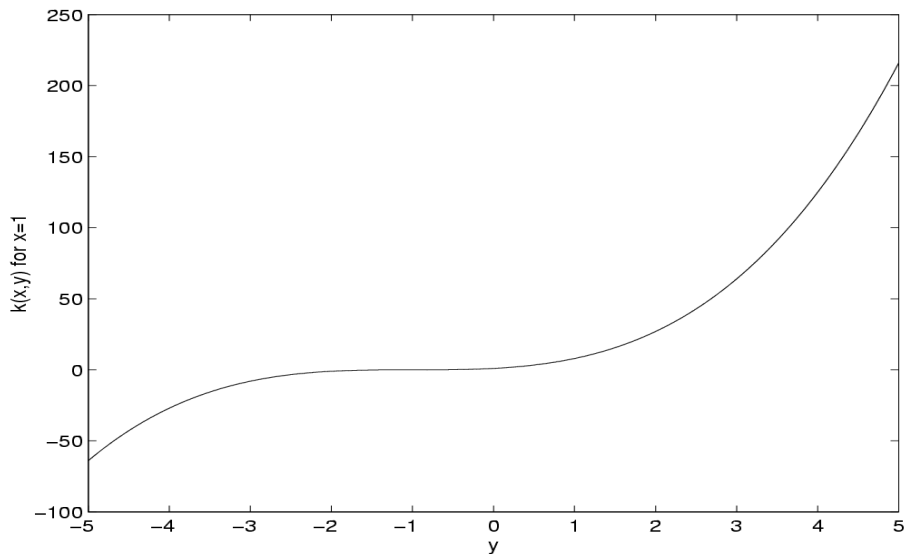
Laplace Kernel Covariance



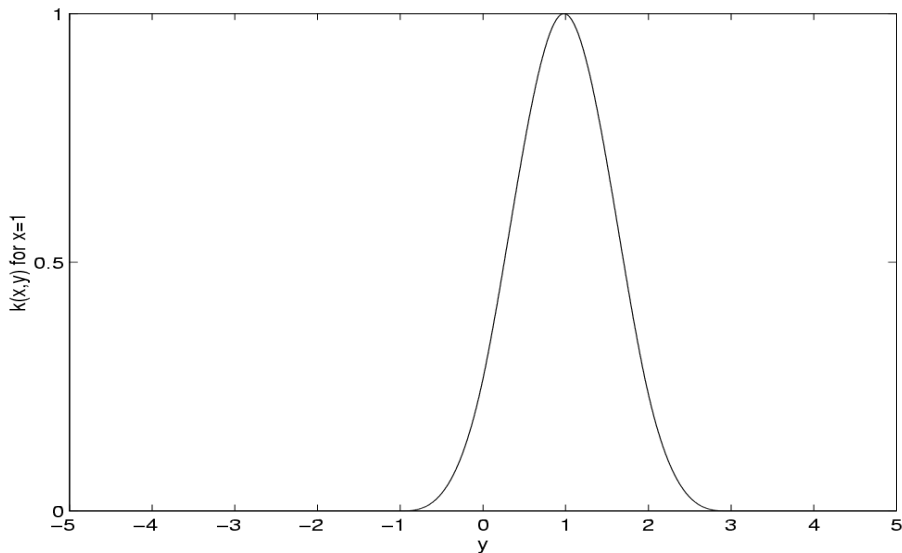
Gaussian Kernel



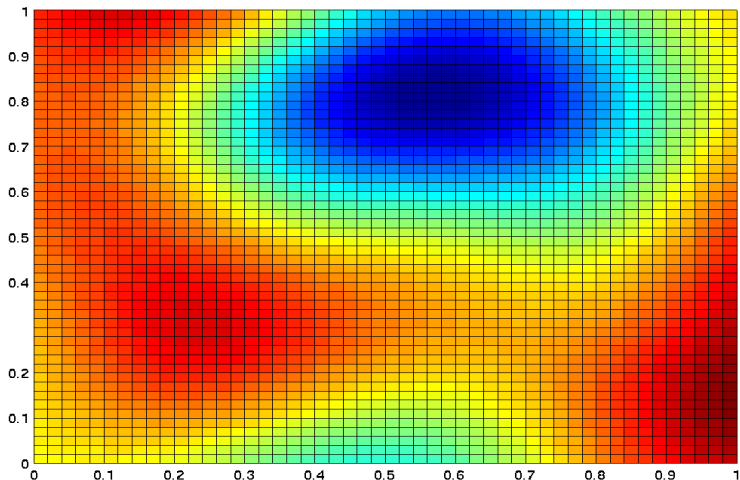
Polynomial (Order 3)



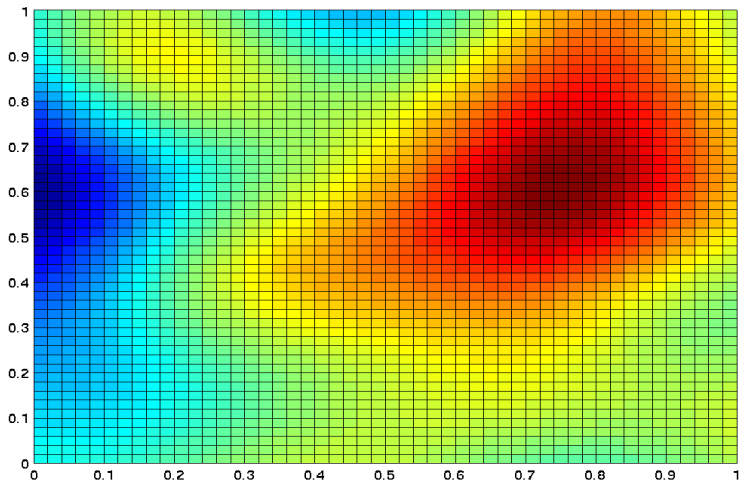
B_3 -Spline Kernel



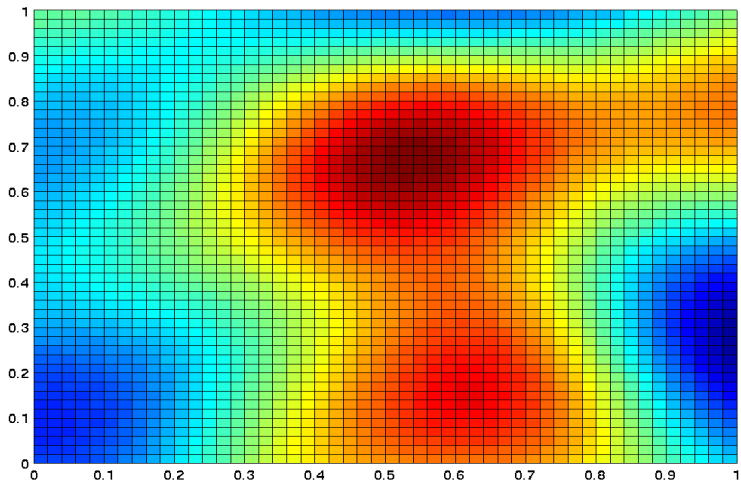
Sample from Gaussian RBF



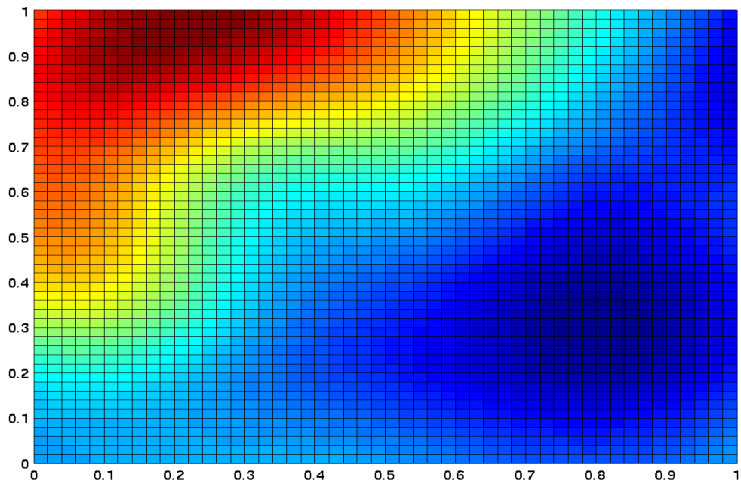
Sample from Gaussian RBF



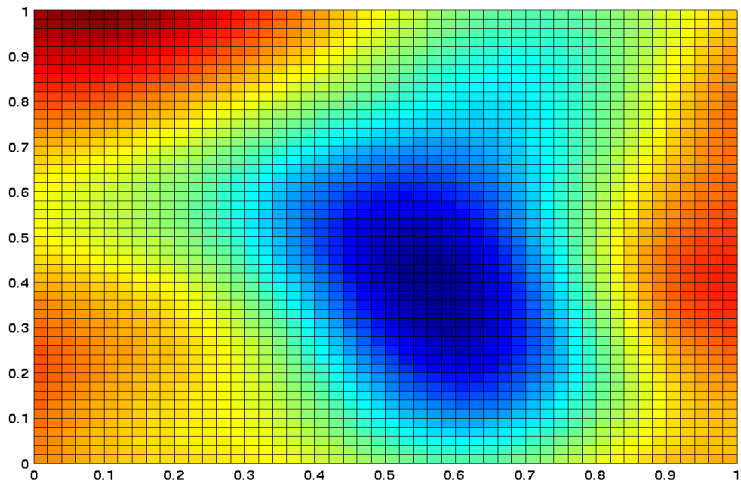
Sample from Gaussian RBF



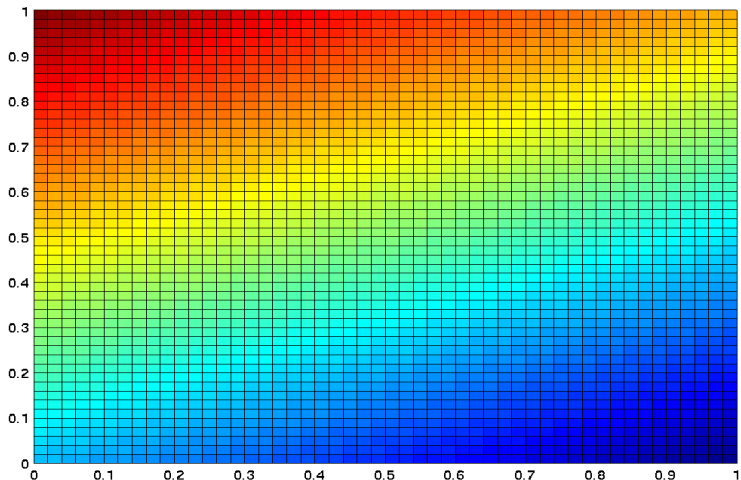
Sample from Gaussian RBF



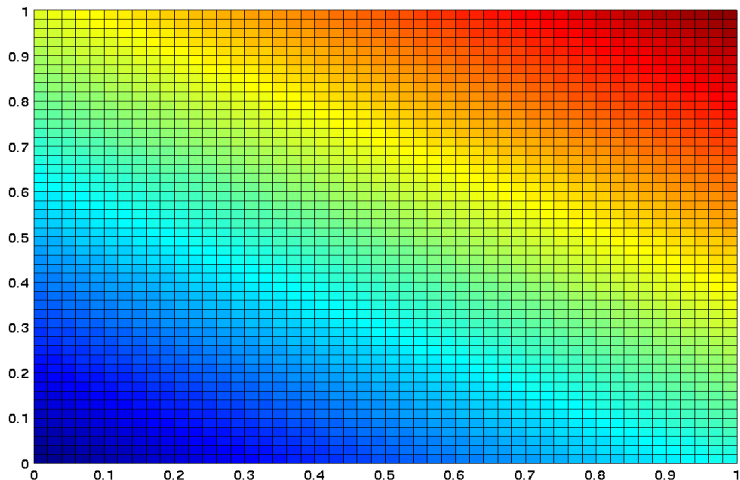
Sample from Gaussian RBF



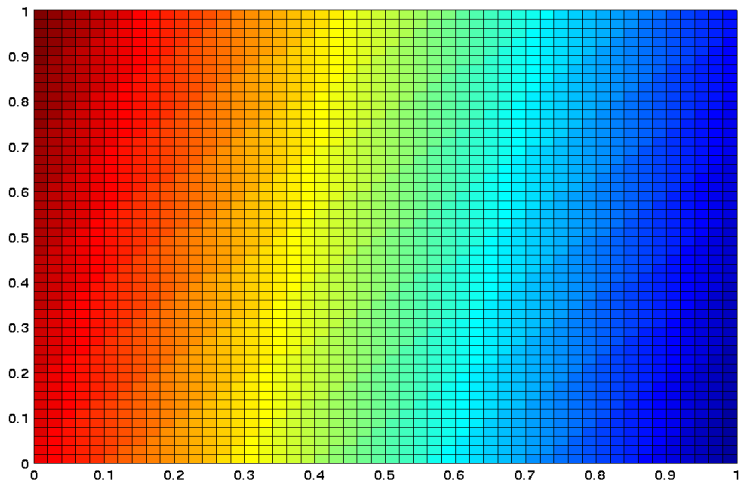
Sample from linear kernel



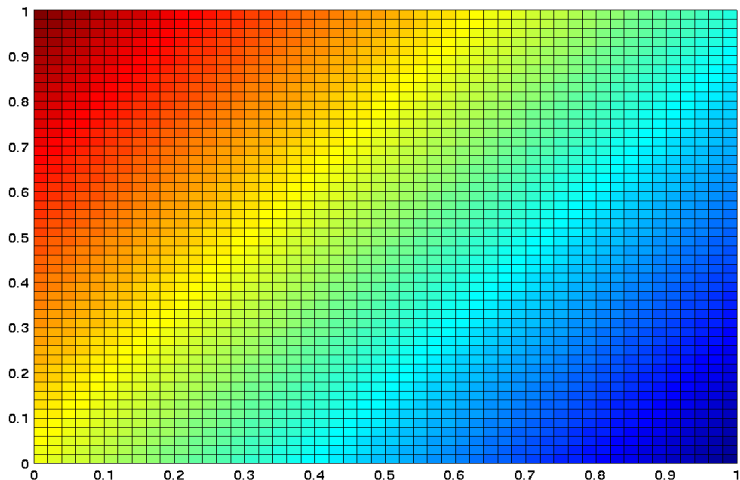
Sample from linear kernel



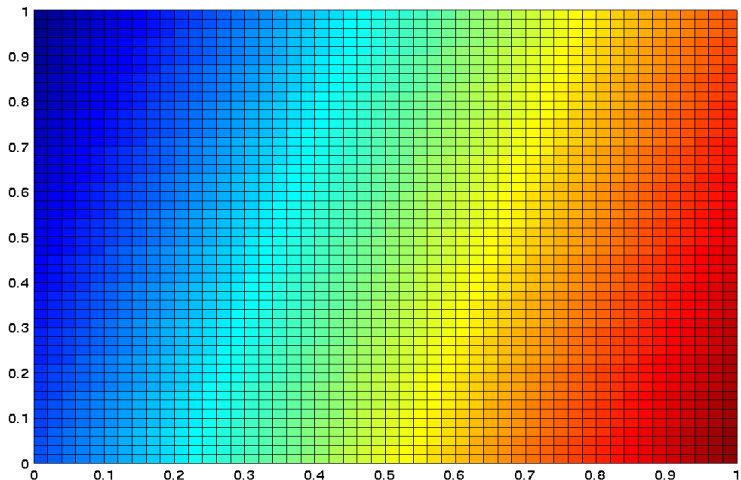
Sample from linear kernel



Sample from linear kernel



Sample from linear kernel



Mini Summary

Choose a suitable sufficient statistic $\phi(x, y)$

- Conditionally multinomial distribution leads to Gaussian Process multiclass classifier.
- Conditionally Gaussian leads to Gaussian Process regression. **Note:** we estimate mean and variance.
- Conditionally Poisson distributions yield spatial Poisson regression.

Solve the optimization problem

This is typically convex.

The bottom line

Instead of choosing $k(x, x')$ choose $k((x, y), (x', y'))$.

Mini Summary

Choose a suitable sufficient statistic $\phi(x, y)$

- Conditionally multinomial distribution leads to Gaussian Process multiclass classifier.
- Conditionally Gaussian leads to Gaussian Process regression. **Note:** we estimate mean and variance.
- Conditionally Poisson distributions yield spatial Poisson regression.

Solve the optimization problem

This is typically convex.

The bottom line

Instead of choosing $k(x, x')$ choose $k((x, y), (x', y'))$.

Mini Summary

Choose a suitable sufficient statistic $\phi(x, y)$

- Conditionally multinomial distribution leads to Gaussian Process multiclass classifier.
- Conditionally Gaussian leads to Gaussian Process regression. **Note:** we estimate mean and variance.
- Conditionally Poisson distributions yield spatial Poisson regression.

Solve the optimization problem

This is typically convex.

The bottom line

Instead of choosing $k(x, x')$ choose $k((x, y), (x', y'))$.

Mini Summary

Choose a suitable sufficient statistic $\phi(x, y)$

- Conditionally multinomial distribution leads to Gaussian Process multiclass classifier.
- Conditionally Gaussian leads to Gaussian Process regression. **Note:** we estimate mean and variance.
- Conditionally Poisson distributions yield spatial Poisson regression.

Solve the optimization problem

This is typically convex.

The bottom line

Instead of choosing $k(x, x')$ choose $k((x, y), (x', y'))$.

Mini Summary

Choose a suitable sufficient statistic $\phi(x, y)$

- Conditionally multinomial distribution leads to Gaussian Process multiclass classifier.
- Conditionally Gaussian leads to Gaussian Process regression. **Note:** we estimate mean and variance.
- Conditionally Poisson distributions yield spatial Poisson regression.

Solve the optimization problem

This is typically convex.

The bottom line

Instead of choosing $k(x, x')$ choose $k((x, y), (x', y'))$.

Mini Summary

Choose a suitable sufficient statistic $\phi(x, y)$

- Conditionally multinomial distribution leads to Gaussian Process multiclass classifier.
- Conditionally Gaussian leads to Gaussian Process regression. **Note:** we estimate mean and variance.
- Conditionally Poisson distributions yield spatial Poisson regression.

Solve the optimization problem

This is typically convex.

The bottom line

Instead of choosing $k(x, x')$ choose $k((x, y), (x', y'))$.

- 1 Conditional Models
 - Log-partition Function, Densities, and Expectations
 - Inner Products and Kernels
 - Examples of Kernels
- 2 **Gaussian Process Classification**
 - Feature map
 - Examples
- 3 Gaussian Process Regression
 - Homoscedastic Model
 - Heteroscedastic Model
- 4 Conditional Random Fields
 - Model Structure
 - Kernel Expansion
 - Connections to Hidden Markov Models

Gaussian Process Classification

Sufficient Statistic

We pick $\phi(x, y) = \phi(x) \otimes \mathbf{e}_y$, that is

$$k((x, y), (x', y')) = k(x, x')\delta_{yy'}, \text{ where } y, y' \in \{1, \dots, n\}$$

Kernel Expansion

By the representer theorem we get that

$$\theta = \sum_{i=1}^m \sum_y \alpha_{iy} \phi(x_i, y)$$

Optimization Problem

Not too messy and convex.

Gaussian Process Classification

Sufficient Statistic

We pick $\phi(x, y) = \phi(x) \otimes \mathbf{e}_y$, that is

$$k((x, y), (x', y')) = k(x, x')\delta_{yy'}, \text{ where } y, y' \in \{1, \dots, n\}$$

Kernel Expansion

By the representer theorem we get that

$$\theta = \sum_{i=1}^m \sum_y \alpha_{iy} \phi(x_i, y)$$

Optimization Problem

Not too messy and convex.

Gaussian Process Classification

Sufficient Statistic

We pick $\phi(x, y) = \phi(x) \otimes \mathbf{e}_y$, that is

$$k((x, y), (x', y')) = k(x, x')\delta_{yy'}, \text{ where } y, y' \in \{1, \dots, n\}$$

Kernel Expansion

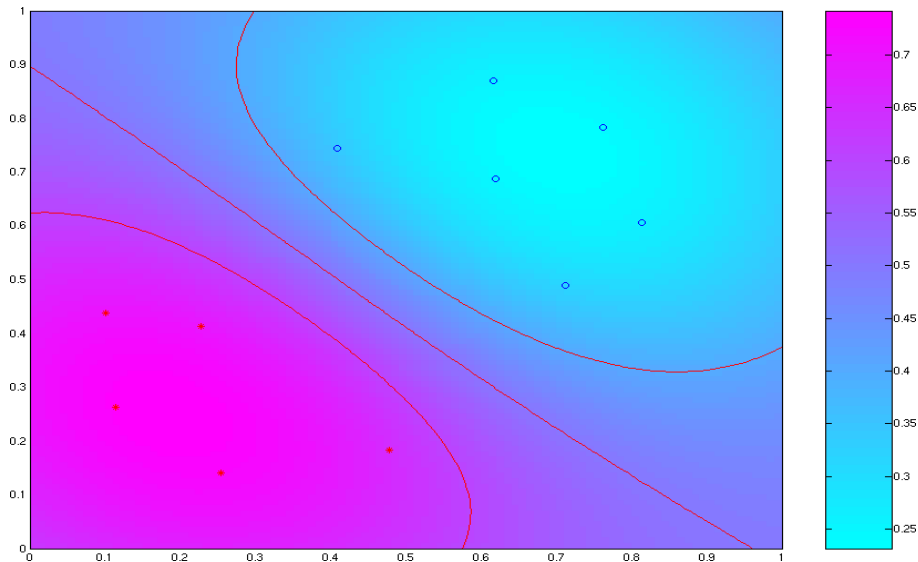
By the representer theorem we get that

$$\theta = \sum_{i=1}^m \sum_y \alpha_{iy} \phi(x_i, y)$$

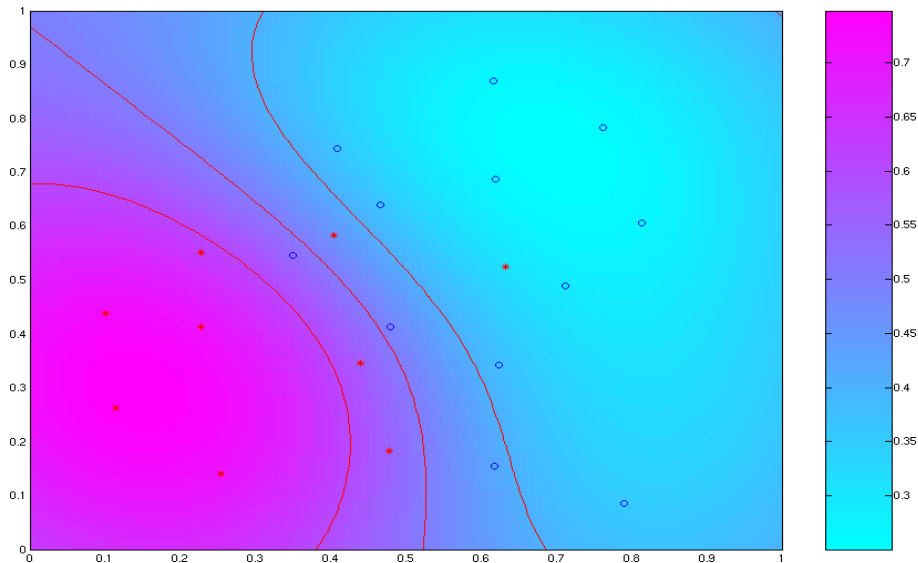
Optimization Problem

Not too messy and convex.

A Toy Example



Noisy Data



Mini Summary

Feature Map

- Conditionally multinomial models $y|x$
- Feature map is $e_y \otimes \phi(x)$
- Kernel $k((x, y), (x', y')) = \delta_{y,y'} k(x, x')$
- Could use different interaction between labels.

Optimization Problem

- Convex problem
- Solve in dual space by Newton's method
- Could solve in primal space if $\phi(x, y)$ can be computed efficiently.

Caveat

- True posterior is only approximated by **mode** of posterior.
- Would need sampling methods for exact inference.

- 1 Conditional Models
 - Log-partition Function, Densities, and Expectations
 - Inner Products and Kernels
 - Examples of Kernels
- 2 Gaussian Process Classification
 - Feature map
 - Examples
- 3 Gaussian Process Regression**
 - Homoscedastic Model
 - Heteroscedastic Model
- 4 Conditional Random Fields
 - Model Structure
 - Kernel Expansion
 - Connections to Hidden Markov Models

Recall: Maximum a Posteriori Estimation

$$-\log p(\theta|X, Y) = \sum_{i=1}^m -\langle \phi(x_i, y_i), \theta \rangle + g(\theta|x_i) + \frac{1}{2\sigma^2} \|\theta\|^2 + c$$

Domain

- Continuous domain of observations $\mathcal{Y} = \mathbb{R}$
- We want to have a conditionally normal distribution $y|x$.
- Log-partition function $g(\theta|x)$ easy to compute in closed form as normal distribution.

Recall: Maximum a Posteriori Estimation

$$-\log p(\theta|X, Y) = \sum_{i=1}^m -\langle \phi(x_i, y_i), \theta \rangle + g(\theta|x_i) + \frac{1}{2\sigma^2} \|\theta\|^2 + c$$

Domain

- Continuous domain of observations $\mathcal{Y} = \mathbb{R}$
- We want to have a conditionally normal distribution $y|x$.
- Log-partition function $g(\theta|x)$ easy to compute **in closed form** as normal distribution.

Recall: Maximum a Posteriori Estimation

$$-\log p(\theta|X, Y) = \sum_{i=1}^m -\langle \phi(x_i, y_i), \theta \rangle + g(\theta|x_i) + \frac{1}{2\sigma^2} \|\theta\|^2 + c$$

Domain

- Continuous domain of observations $y = \mathbb{R}$
- We want to have a conditionally normal distribution $y|x$.
- Log-partition function $g(\theta|x)$ easy to compute **in closed form** as normal distribution.

Recall: Maximum a Posteriori Estimation

$$-\log p(\theta | X, Y) = \sum_{i=1}^m -\langle \phi(x_i, y_i), \theta \rangle + g(\theta | x_i) + \frac{1}{2\sigma^2} \|\theta\|^2 + c$$

Domain

- Continuous domain of observations $y = \mathbb{R}$
- We want to have a conditionally normal distribution $y|x$.
- Log-partition function $g(\theta|x)$ easy to compute **in closed form** as normal distribution.

Recall: Maximum a Posteriori Estimation

$$-\log p(\theta|X, Y) = \sum_{i=1}^m -\langle \phi(x_i, y_i), \theta \rangle + g(\theta|x_i) + \frac{1}{2\sigma^2} \|\theta\|^2 + c$$

Domain

- Continuous domain of observations $\mathcal{Y} = \mathbb{R}$
- We want to have a conditionally normal distribution $y|x$.
- Log-partition function $g(\theta|x)$ easy to compute **in closed form** as normal distribution.

Standard Model

Key Idea

- For fixed parameter θ we want to have a normal distribution with fixed variance.
- Exponential family model of $y|x$ with $c(x)y - \frac{1}{2\sigma^2}y^2$ in exponent.

Sufficient Statistic

Pick $\phi(x, y) = (y\phi(x), y^2)$, that is

$$k((x, y), (x', y')) = k(x, x')yy' + y^2y'^2 \text{ where } y, y' \in \mathbb{R}$$

Traditionally the variance is fixed.

Inference Problem

After straightforward algebra we get standard GP regression model.

Standard Model

Key Idea

- For fixed parameter θ we want to have a normal distribution with fixed variance.
- Exponential family model of $y|x$ with $c(x)y - \frac{1}{2\sigma^2}y^2$ in exponent.

Sufficient Statistic

Pick $\phi(x, y) = (y\phi(x), y^2)$, that is

$$k((x, y), (x', y')) = k(x, x')yy' + y^2y'^2 \text{ where } y, y' \in \mathbb{R}$$

Traditionally the variance is fixed.

Inference Problem

After straightforward algebra we get standard GP regression model.

Standard Model

Key Idea

- For fixed parameter θ we want to have a normal distribution with fixed variance.
- Exponential family model of $y|x$ with $c(x)y - \frac{1}{2\sigma^2}y^2$ in exponent.

Sufficient Statistic

Pick $\phi(x, y) = (y\phi(x), y^2)$, that is

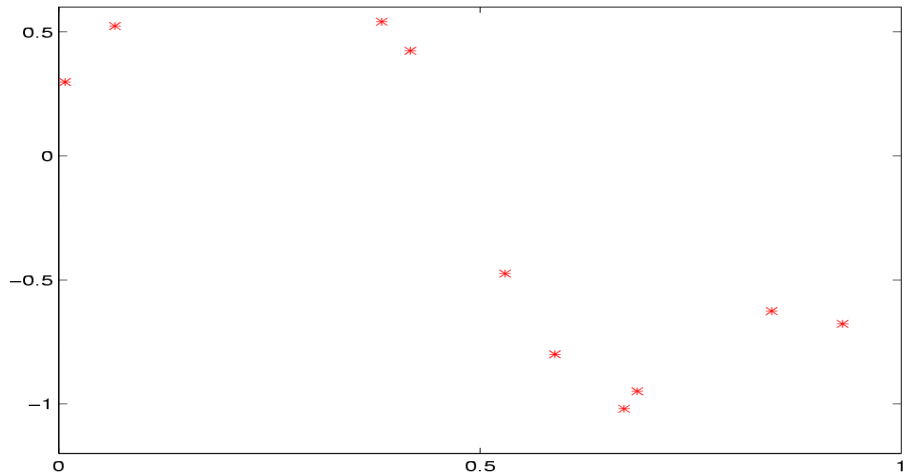
$$k((x, y), (x', y')) = k(x, x')yy' + y^2y'^2 \text{ where } y, y' \in \mathbb{R}$$

Traditionally the variance is fixed.

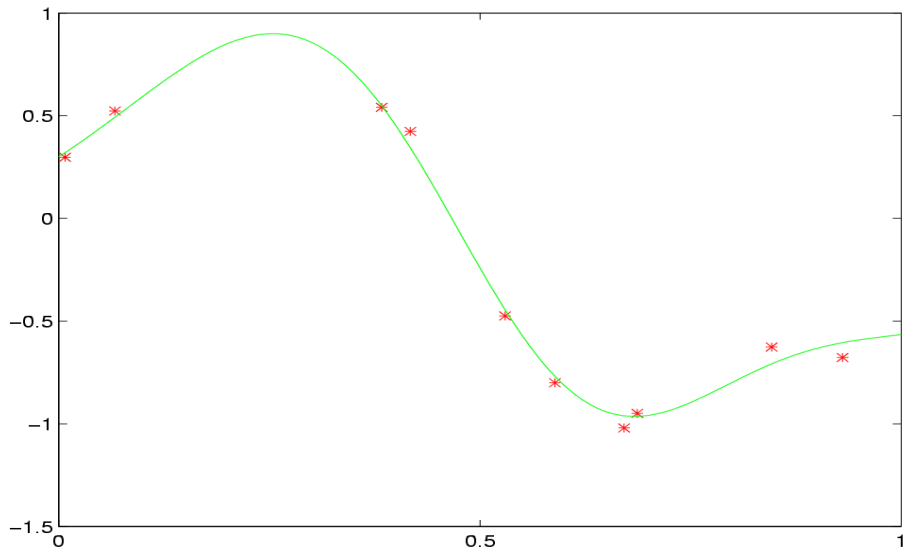
Inference Problem

After straightforward algebra we get standard GP regression model.

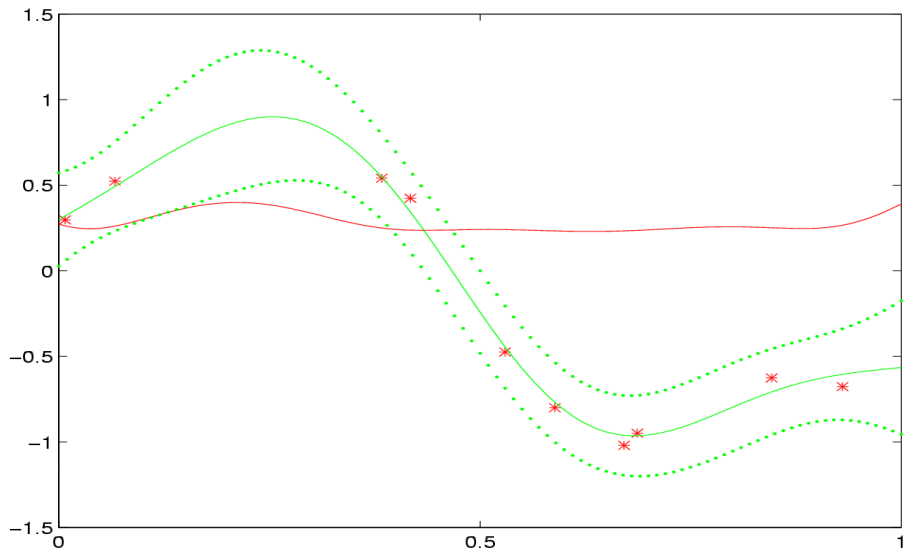
Training Data



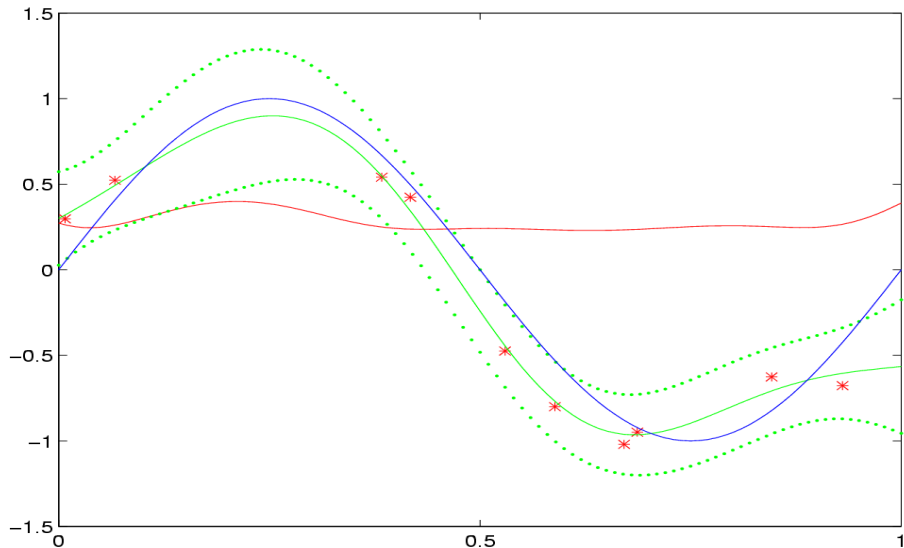
Mean $\vec{k}^\top(x)(K + \sigma^2\mathbf{1})^{-1}y$



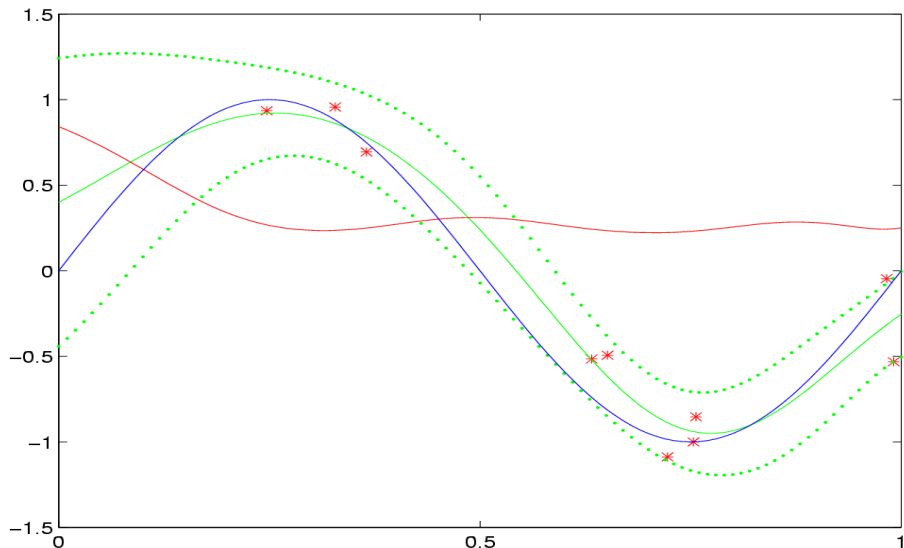
Variance $k(x, x) + \sigma^2 - \vec{k}^\top(x)(K + \sigma^2\mathbf{1})^{-1}\vec{k}(x)$



Putting everything together ...



Another Example



Heteroscedastic Regression

Key Idea

Make both linear and quadratic term in $y|x$ dependent on x .

Sufficient Statistic

Pick $\phi(x, y) = (y\phi_1(x), y^2\phi_2(x))$, that is

$k((x, y), (x', y')) = k_1(x, x')yy' + k_2(x, x')y^2y'^2$ where $y, y' \in \mathbb{R}$

We estimate mean and variance **simultaneously**.

Kernel Expansion

By the representer theorem (and more algebra) we get

$$\theta = \left(\sum_{i=1}^m \alpha_{i1} \phi_1(x_i), \sum_{i=1}^m \alpha_{i2} \phi_2(x_i) \right)$$

Heteroscedastic Regression

Key Idea

Make both linear and quadratic term in $y|x$ dependent on x .

Sufficient Statistic

Pick $\phi(x, y) = (y\phi_1(x), y^2\phi_2(x))$, that is

$k((x, y), (x', y')) = k_1(x, x')yy' + k_2(x, x')y^2y'^2$ where $y, y' \in \mathbb{R}$

We estimate mean and variance **simultaneously**.

Kernel Expansion

By the representer theorem (and more algebra) we get

$$\theta = \left(\sum_{i=1}^m \alpha_{i1} \phi_1(x_i), \sum_{i=1}^m \alpha_{i2} \phi_2(x_i) \right)$$

Heteroscedastic Regression

Key Idea

Make both linear and quadratic term in $y|x$ dependent on x .

Sufficient Statistic

Pick $\phi(x, y) = (y\phi_1(x), y^2\phi_2(x))$, that is

$k((x, y), (x', y')) = k_1(x, x')yy' + k_2(x, x')y^2y'^2$ where $y, y' \in \mathbb{R}$

We estimate mean and variance **simultaneously**.

Kernel Expansion

By the representer theorem (and more algebra) we get

$$\theta = \left(\sum_{i=1}^m \alpha_{i1} \phi_1(x_i), \sum_{i=1}^m \alpha_{i2} \phi_2(x_i) \right)$$

Heteroscedastic Regression

Optimization Problem

$$\begin{aligned} & \sum_{i=1}^m \left[-\frac{1}{4} \left[\sum_{j=1}^m \alpha_{1j} k_1(x_i, x_j) \right]^\top \left[\sum_{j=1}^m \alpha_{2j} k_2(x_i, x_j) \right]^{-1} \left[\sum_{j=1}^m \alpha_{1j} k_1(x_i, x_j) \right] \right. \\ & - \frac{1}{2} \log \det -2 \left[\sum_{j=1}^m \alpha_{2j} k_2(x_i, x_j) \right] - \sum_{j=1}^m [y_i^\top \alpha_{1j} k_1(x_i, x_j) + (y_j^\top \alpha_{2j} y_j) k_2(x_i, x_j)] \\ & \left. + \frac{1}{2\sigma^2} \sum_{i,j} \alpha_{1i}^\top \alpha_{1j} k_1(x_i, x_j) + \text{tr} [\alpha_{2i} \alpha_{2j}^\top] k_2(x_i, x_j) \right] \\ & \text{subject to } 0 \succ \sum_{i=1}^m \alpha_{2i} k(x_i, x_j) \end{aligned}$$

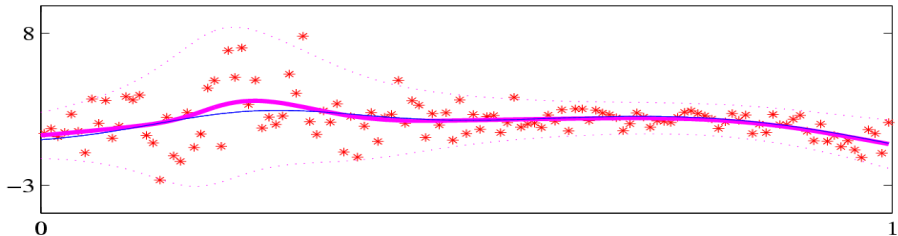
Heteroscedastic Regression

Optimization Problem

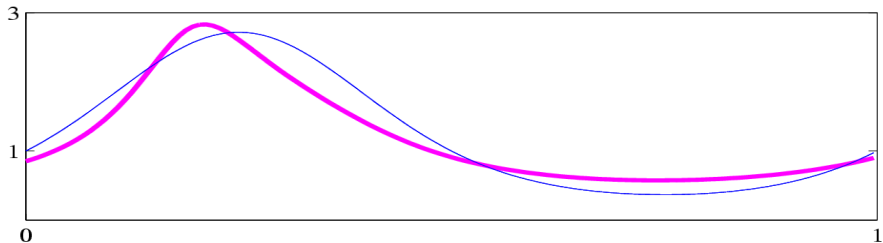
- The problem is convex
- The log-determinant from the normalization of the Gaussian distribution acts as a **barrier function**.
- We get a semidefinite program.
- Because of the barrier function we can solve it by Newton's method.

Heteroscedastic Regression

regression estimation and training data

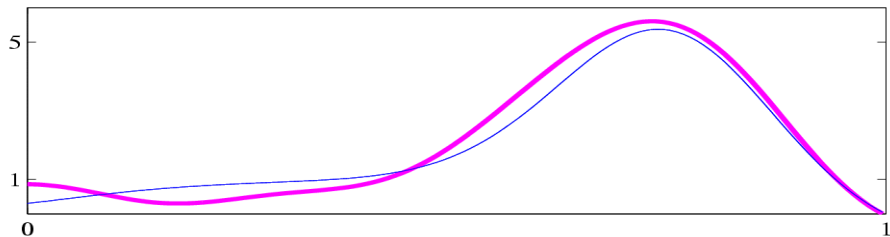


variance estimation

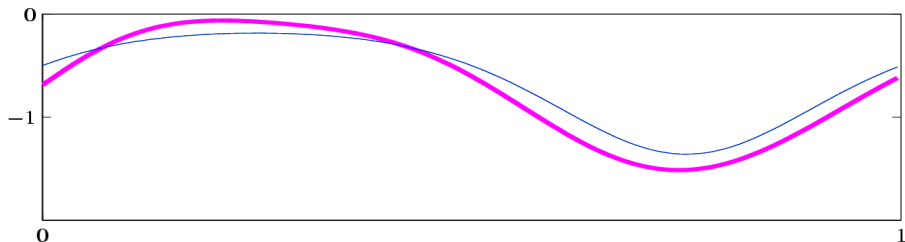


Natural Parameters

θ_1 estimation



θ_2 estimation



Sufficient Statistics

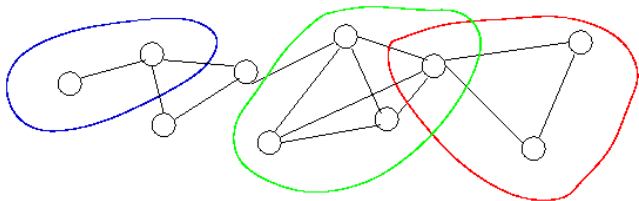
- Conditionally normal model explained if $\phi(y, x)$ has linear and quadratic terms.
- For homoscedastic model we only need to estimate the linear term. Quadratic term is fixed.
- Second order kernel in y , arbitrary kernel in x .

Optimization

- Linear system is all we need for fixed variance
- Semidefinite program for heteroscedastic estimation
- Can be solved by Newton's method, as the log-determinant acts as barrier function.

- 1 Conditional Models
 - Log-partition Function, Densities, and Expectations
 - Inner Products and Kernels
 - Examples of Kernels
- 2 Gaussian Process Classification
 - Feature map
 - Examples
- 3 Gaussian Process Regression
 - Homoscedastic Model
 - Heteroscedastic Model
- 4 **Conditional Random Fields**
 - Model Structure
 - Kernel Expansion
 - Connections to Hidden Markov Models

Graphical Models



Corollary

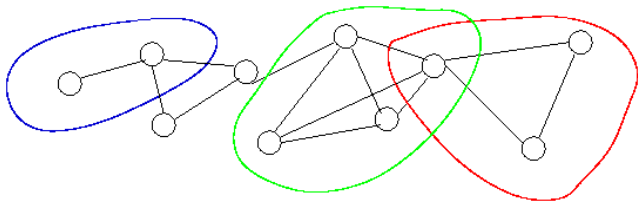
The conditional density $p(y|x)$ can be written in terms of potential functions defined on the maximal cliques in y .

Corollary

*Featuremap $\phi(x)$ decomposes via $\phi(x) = (\dots, \phi_c(x_c), \dots)$.
Consequently we can write the kernel via*

$$k(x, x') = \langle \phi(x), \phi(x') \rangle = \sum_c \langle \phi_c(x_c), \phi_c(x'_c) \rangle = \sum_c k_c(x_c, x'_c)$$

Graphical Models



Corollary

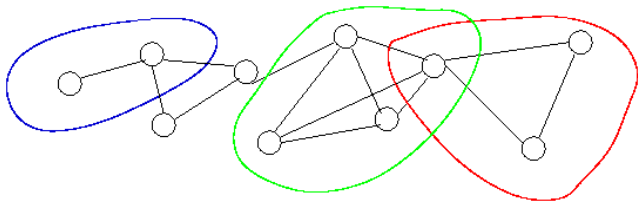
*The conditional density $p(y|x)$ can be written in terms of potential functions defined on the maximal cliques **in y** .*

Corollary

*Featuremap $\phi(x)$ decomposes via $\phi(x) = (\dots, \phi_c(x_c), \dots)$.
Consequently we can write the kernel via*

$$k(x, x') = \langle \phi(x), \phi(x') \rangle = \sum_c \langle \phi_c(x_c), \phi_c(x'_c) \rangle = \sum_c k_c(x_c, x'_c)$$

Graphical Models



Corollary

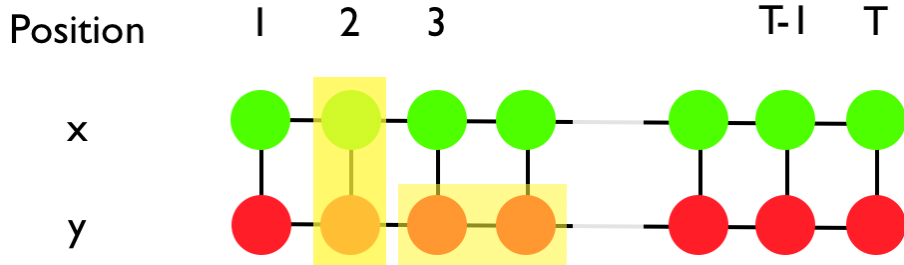
The conditional density $p(y|x)$ can be written in terms of potential functions defined on the maximal cliques *in* y .

Corollary

Featuremap $\phi(x)$ decomposes via $\phi(x) = (\dots, \phi_c(x_c), \dots)$.
Consequently we can write the kernel via

$$k(x, x') = \langle \phi(x), \phi(x') \rangle = \sum_c \langle \phi_c(x_c), \phi_c(x'_c) \rangle = \sum_c k_c(x_c, x'_c)$$

Conditional Random Fields

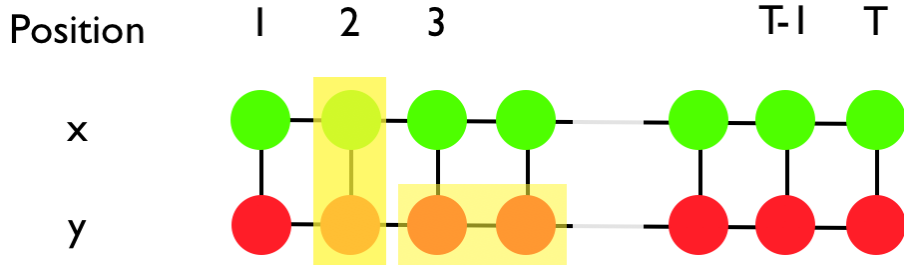


Key Points

- Cliques are (x_t, y_t) , (x_t, x_{t+1}) , and (y_t, y_{t+1})
- We can drop cliques in (x_t, x_{t+1})

$$p(y|x, \theta) = \exp\left(\sum_t \langle \phi_{xy}(x_t, y_t), \theta_{xy,t} \rangle + \langle \phi_{yy}(y_t, y_{t+1}), \theta_{yy,t} \rangle\right. \\ \left. + \langle \phi_{xx}(x_t, x_{t+1}), \theta_{xx,t} \rangle - g(\theta|x)\right)$$

Conditional Random Fields



Key Points

- Cliques are (x_t, y_t) , (x_t, x_{t+1}) , and (y_t, y_{t+1})
- We can drop cliques in (x_t, x_{t+1})

$$p(y|x, \theta) = \exp\left(\sum_t \langle \phi_{xy}(x_t, y_t), \theta_{xy,t} \rangle + \langle \phi_{yy}(y_t, y_{t+1}), \theta_{yy,t} \rangle + \langle \phi_{xx}(x_t, x_{t+1}), \theta_{xx,t} \rangle - g(\theta|x)\right)$$

Computational Issues

Key Points

- Compute $g(\theta|x)$ via dynamic programming
- Assume stationarity of the model, that is θ_c does not depend on the position of the

Dynamic Programming

$$\begin{aligned} &g(\theta|x) \\ &= \log \sum_{y_1, \dots, y_T} \prod_{t=1}^T \underbrace{\exp(\langle \phi_{xy}(x_t, y_t), \theta_{xy} \rangle + \langle \phi_{yy}(y_t, y_{t+1}), \theta_{yy} \rangle)}_{M_t(y_t, y_{t+1})} \\ &= \log \sum_{y_1} \sum_{y_2} M_1(y_1, y_2) \sum_{y_3} M_2(y_2, y_3) \dots \sum_{y_T} M_T(y_{T-1}, y_T) \end{aligned}$$

Efficient computation of $g(\theta|x)$, $p(y_t|x, \theta)$ and $p(y_t, y_{t+1}|x, \theta)$.

Computational Issues

Key Points

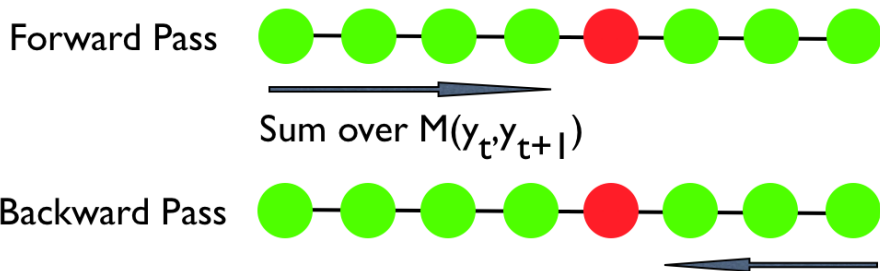
- Compute $g(\theta|x)$ via dynamic programming
- Assume stationarity of the model, that is θ_c does not depend on the position of the

Dynamic Programming

$$\begin{aligned} &g(\theta|x) \\ &= \log \sum_{y_1, \dots, y_T} \prod_{t=1}^T \underbrace{\exp(\langle \phi_{xy}(x_t, y_t), \theta_{xy} \rangle + \langle \phi_{yy}(y_t, y_{t+1}), \theta_{yy} \rangle)}_{M_t(y_t, y_{t+1})} \\ &= \log \sum_{y_1} \sum_{y_2} M_1(y_1, y_2) \sum_{y_3} M_2(y_2, y_3) \dots \sum_{y_T} M_T(y_{T-1}, y_T) \end{aligned}$$

Efficient computation of $g(\theta|x)$, $p(y_t|x, \theta)$ and $p(y_t, y_{t+1}|x, \theta)$.

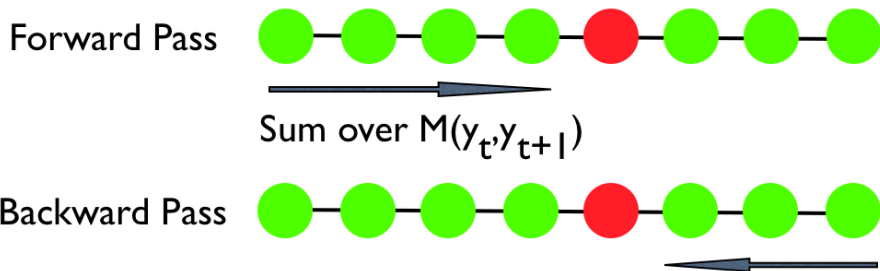
Forward Backward Algorithm



Key Idea

- Store sum over all y_1, \dots, y_{t-1} (forward pass) and over all y_{t+1}, \dots, y_T as intermediate values
- We get those values for all positions t in one sweep.
- Extend this to message passing (when we have trees).

Forward Backward Algorithm



Key Idea

- Store sum over all y_1, \dots, y_{t-1} (forward pass) and over all y_{t+1}, \dots, y_T as intermediate values
- We get those values for all positions t in one sweep.
- Extend this to message passing (when we have trees).

Minimization

Objective Function

$$-\log p(\theta|X, Y) = \sum_{i=1}^m -\langle \phi(\mathbf{x}_i, \mathbf{y}_i), \theta \rangle + g(\theta|\mathbf{x}_i) + \frac{1}{2\sigma^2} \|\theta\|^2 + c$$

$$\partial_{\theta} -\log p(\theta|X, Y) = \sum_{i=1}^m -\phi(\mathbf{x}_i, \mathbf{y}_i) + \mathbf{E} [\phi(\mathbf{x}_i, \mathbf{y}_i)|\mathbf{x}_i] + \frac{1}{\sigma^2} \theta$$

We only need $\mathbf{E} [\phi_{xy}(\mathbf{x}_{it}, \mathbf{y}_{it})|\mathbf{x}_i]$ and $\mathbf{E} [\phi_{yy}(\mathbf{y}_{it}, \mathbf{y}_{i(t+1)})|\mathbf{x}_i]$.

Kernel Trick

- Conditional expectations of $\Phi(\mathbf{x}_{it}, \mathbf{y}_{it})$ cannot be computed explicitly **but** inner products can.

$$\langle \phi_{xy}(\mathbf{x}'_t, \mathbf{y}'_t), \mathbf{E} [\phi_{xy}(\mathbf{x}_t, \mathbf{y}_t)|\mathbf{x}] = \mathbf{E} [k((\mathbf{x}'_t, \mathbf{y}'_t), (\mathbf{x}_t, \mathbf{y}_t)|\mathbf{x})$$

- Only need marginals $p(\mathbf{y}_t|\mathbf{x}, \theta)$ and $p(\mathbf{y}_t, \mathbf{y}_{t+1}|\mathbf{x}, \theta)$, which we get via dynamic programming.

Minimization

Objective Function

$$-\log p(\theta|X, Y) = \sum_{i=1}^m -\langle \phi(\mathbf{x}_i, \mathbf{y}_i), \theta \rangle + g(\theta|\mathbf{x}_i) + \frac{1}{2\sigma^2} \|\theta\|^2 + c$$

$$\partial_{\theta} -\log p(\theta|X, Y) = \sum_{i=1}^m -\phi(\mathbf{x}_i, \mathbf{y}_i) + \mathbf{E}[\phi(\mathbf{x}_i, \mathbf{y}_i)|\mathbf{x}_i] + \frac{1}{\sigma^2}\theta$$

We only need $\mathbf{E}[\phi_{xy}(\mathbf{x}_{it}, \mathbf{y}_{it})|\mathbf{x}_i]$ and $\mathbf{E}[\phi_{yy}(\mathbf{y}_{it}, \mathbf{y}_{i(t+1)})|\mathbf{x}_i]$.

Kernel Trick

- Conditional expectations of $\Phi(\mathbf{x}_{it}, \mathbf{y}_{it})$ cannot be computed explicitly **but** inner products can.

$$\langle \phi_{xy}(\mathbf{x}'_t, \mathbf{y}'_t), \mathbf{E}[\phi_{xy}(\mathbf{x}_t, \mathbf{y}_t)|\mathbf{x}] = \mathbf{E}[k((\mathbf{x}'_t, \mathbf{y}'_t), (\mathbf{x}_t, \mathbf{y}_t)|\mathbf{x})]$$

- Only need marginals $p(\mathbf{y}_t|\mathbf{x}, \theta)$ and $p(\mathbf{y}_t, \mathbf{y}_{t+1}|\mathbf{x}, \theta)$, which we get via dynamic programming.

Subspace Representer Theorem

Representer Theorem

Solutions of the MAP problem are given by

$$\theta \in \text{span}\{\phi(x_i, y) \text{ for all } y \in \mathcal{Y} \text{ and } 1 \leq i \leq n\}$$

Big Problem

$|\mathcal{Y}|$ could be huge, e.g. for sequence annotation 2^n .

Solution

- Exploit decomposition of $\phi(x, y)$ into sufficient statistics on cliques.
- Restriction of \mathcal{Y} to cliques is much smaller.

$$\theta_c \in \text{span}\{\phi_c(x_{ci}, y_c) \text{ for all } y_c \in \mathcal{Y}_c \text{ and } 1 \leq i \leq n\}$$

Rather than 2^n we now get $2^{|\mathcal{C}|}$.

Subspace Representer Theorem

Representer Theorem

Solutions of the MAP problem are given by

$$\theta \in \text{span}\{\phi(x_i, y) \text{ for all } y \in \mathcal{Y} \text{ and } 1 \leq i \leq n\}$$

Big Problem

$|\mathcal{Y}|$ could be huge, e.g. for sequence annotation 2^n .

Solution

- Exploit decomposition of $\phi(x, y)$ into sufficient statistics on cliques.
- Restriction of \mathcal{Y} to cliques is much smaller.

$$\theta_c \in \text{span}\{\phi_c(x_{ci}, y_c) \text{ for all } y_c \in \mathcal{Y}_c \text{ and } 1 \leq i \leq n\}$$

Rather than 2^n we now get $2^{|\mathcal{C}|}$.

Subspace Representer Theorem

Representer Theorem

Solutions of the MAP problem are given by

$$\theta \in \text{span}\{\phi(x_i, y) \text{ for all } y \in \mathcal{Y} \text{ and } 1 \leq i \leq n\}$$

Big Problem

$|\mathcal{Y}|$ could be huge, e.g. for sequence annotation 2^n .

Solution

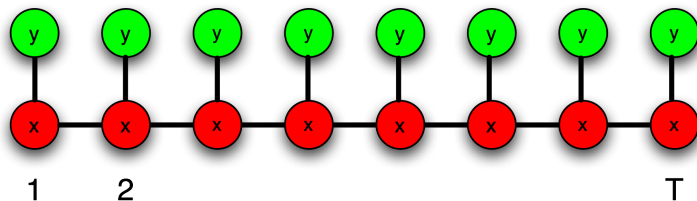
- Exploit decomposition of $\phi(x, y)$ into sufficient statistics on cliques.
- Restriction of \mathcal{Y} to cliques is much smaller.

$$\theta_c \in \text{span}\{\phi_c(x_{ci}, y_c) \text{ for all } y_c \in \mathcal{Y}_c \text{ and } 1 \leq i \leq n\}$$

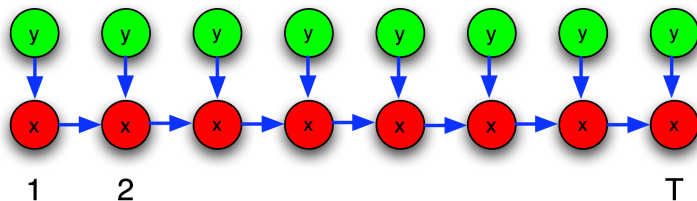
Rather than 2^n we now get $2^{|\mathcal{C}|}$.

CRFs and HMMs

Conditional Random Field: maximize $p(y|x, \theta)$

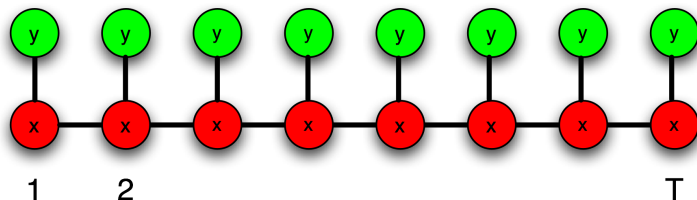


Hidden Markov Model: maximize $p(x, y|\theta)$

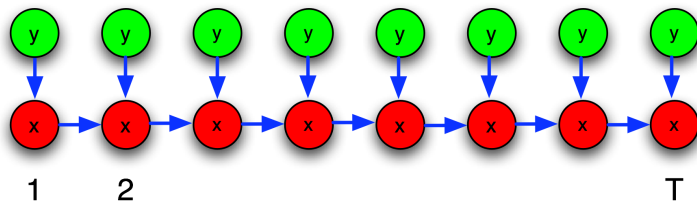


CRFs and HMMs

Conditional Random Field: maximize $p(y|x, \theta)$



Hidden Markov Model: maximize $p(x, y|\theta)$



Equivalence Theorem

Theorem

CRFs and HMMs yield *identical* probability estimates for $p(y|x, \theta)$, if the set of functions is equally expressive.

Proof.

- Write out $p_{\text{CRF}}(y|x, \theta)$ and $p_{\text{HMM}}(x, y|\theta)$, and show that they only differ in the normalization.
- This disappears when computing $p_{\text{HMM}}(y|x, \theta)$.



Consequence

Differential training for current HMM implementations.

Equivalence Theorem

Theorem

CRFs and HMMs yield *identical* probability estimates for $p(y|x, \theta)$, if the set of functions is equally expressive.

Proof.

- Write out $p_{\text{CRF}}(y|x, \theta)$ and $p_{\text{HMM}}(x, y|\theta)$, and show that they only differ in the normalization.
- This disappears when computing $p_{\text{HMM}}(y|x, \theta)$.



Consequence

Differential training for current HMM implementations.

Equivalence Theorem

Theorem

CRFs and HMMs yield *identical* probability estimates for $p(y|x, \theta)$, if the set of functions is equally expressive.

Proof.

- Write out $p_{\text{CRF}}(y|x, \theta)$ and $p_{\text{HMM}}(x, y|\theta)$, and show that they only differ in the normalization.
- This disappears when computing $p_{\text{HMM}}(y|x, \theta)$.



Consequence

Differential training for current HMM implementations.

Mini Summary

Graphical Model Structure

- Same decomposition as in **unconditional** models.
- Only need to take cliques in y into account.

Kernel Expansion

- Representer theorem is still intractable (exponential number of terms).
- Decompose along cliques (we have a representer theorem **per clique**).
- For some parts primal space optimization is more efficient (cliques in y_i alone).

Connection to Hidden Markov Models

- HMMs optimize **generative** performance.
- CRFs optimize a **discriminative** model.
- Can re-optimize HMMs for discriminative performance.

Summary

1 Conditional Models

- Log-partition Function, Densities, and Expectations
- Inner Products and Kernels
- Examples of Kernels

2 Gaussian Process Classification

- Feature map
- Examples

3 Gaussian Process Regression

- Homoscedastic Model
- Heteroscedastic Model

4 Conditional Random Fields

- Model Structure
- Kernel Expansion
- Connections to Hidden Markov Models