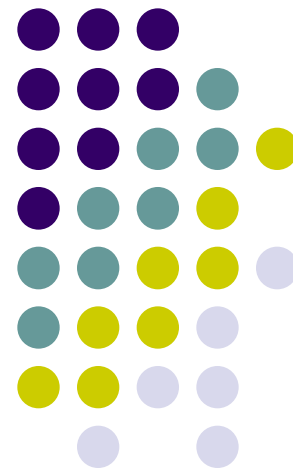# Probabilistic Graphical Models
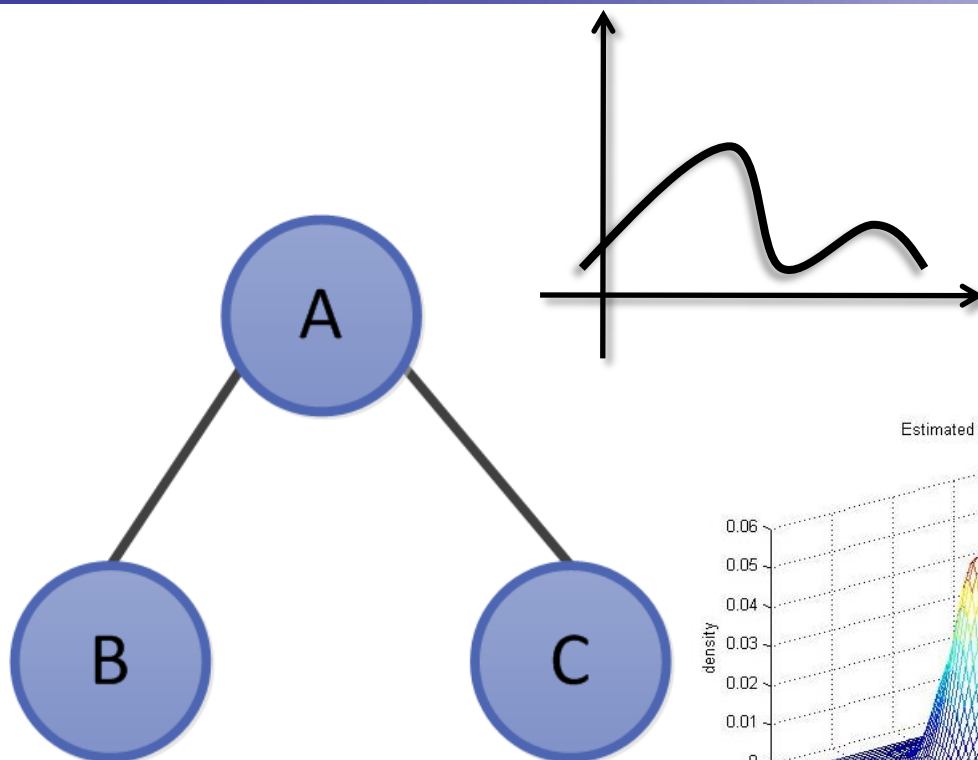
## Kernel Graphical Models

**Ankur Parikh**

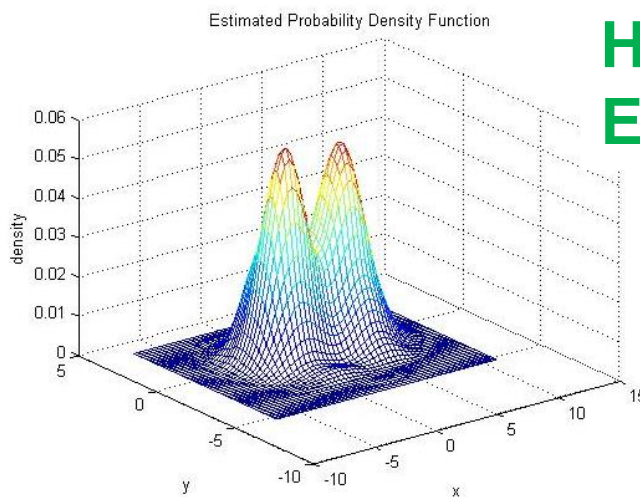**Lecture 20, April 1, 2013**

$$\mathbb{E}[\phi(X)]$$

$\mu_X$

# Nonparametric Graphical Models



**How do we make a conditional probability table out of this?**
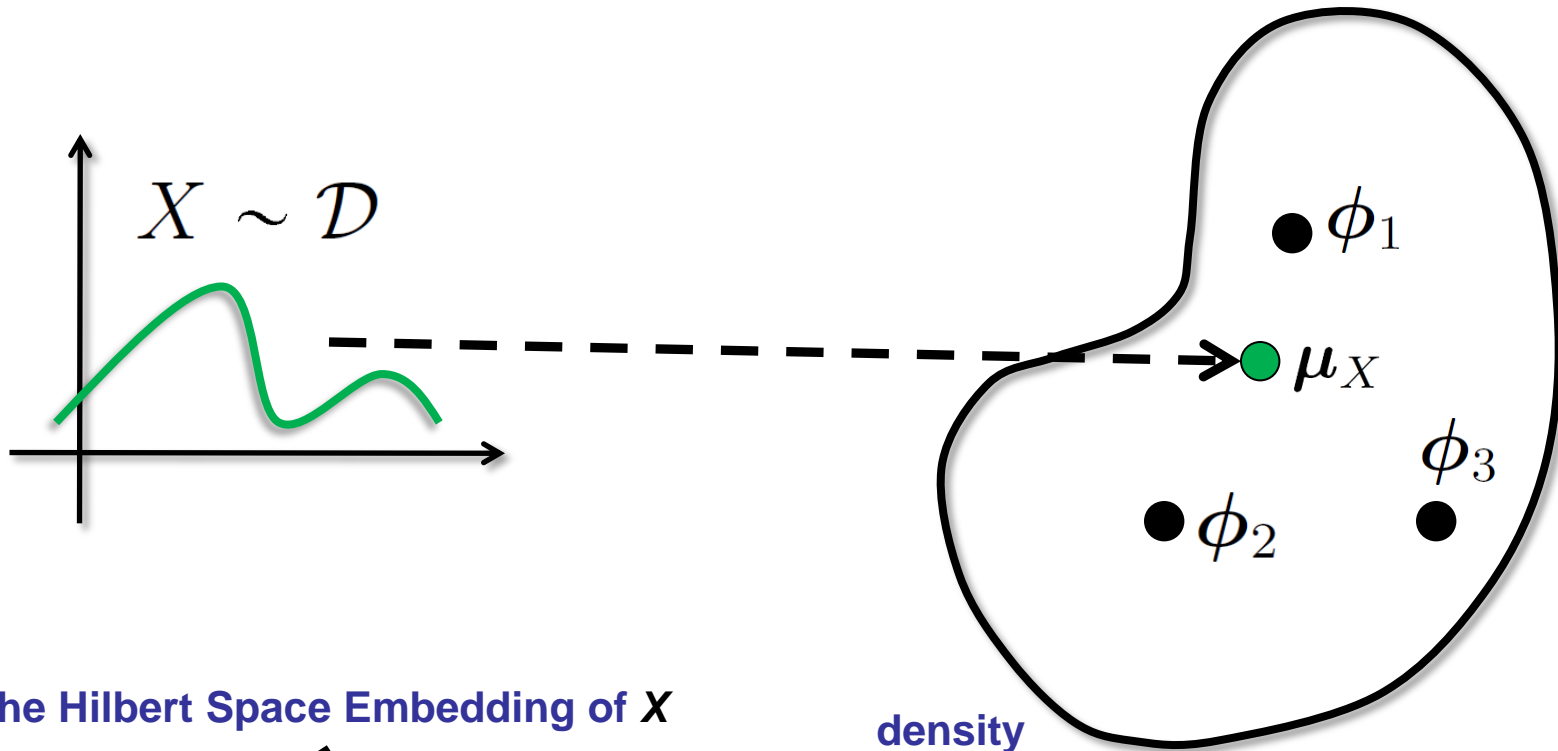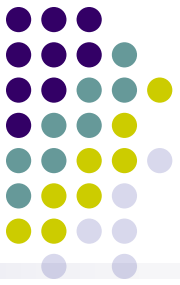
**Hilbert Space Embeddings!!!!!**

Estimated Probability Density Function

- **How to learn parameters?**
- **How to perform inference?**

# Review: Embedding Distribution of One Variable[Smola et al. 2007]



$X \sim \mathcal{D}$

$\phi_1$

$\mu_X$

$\phi_3$

$\phi_2$

**The Hilbert Space Embedding of X**

**density**

$$\boldsymbol{\mu}_X(\cdot) = \mathbb{E}_{X \sim \mathcal{D}}[\phi_X] = \int \boldsymbol{p}_{\mathcal{D}}(X)\phi_X(\cdot)\,dX$$
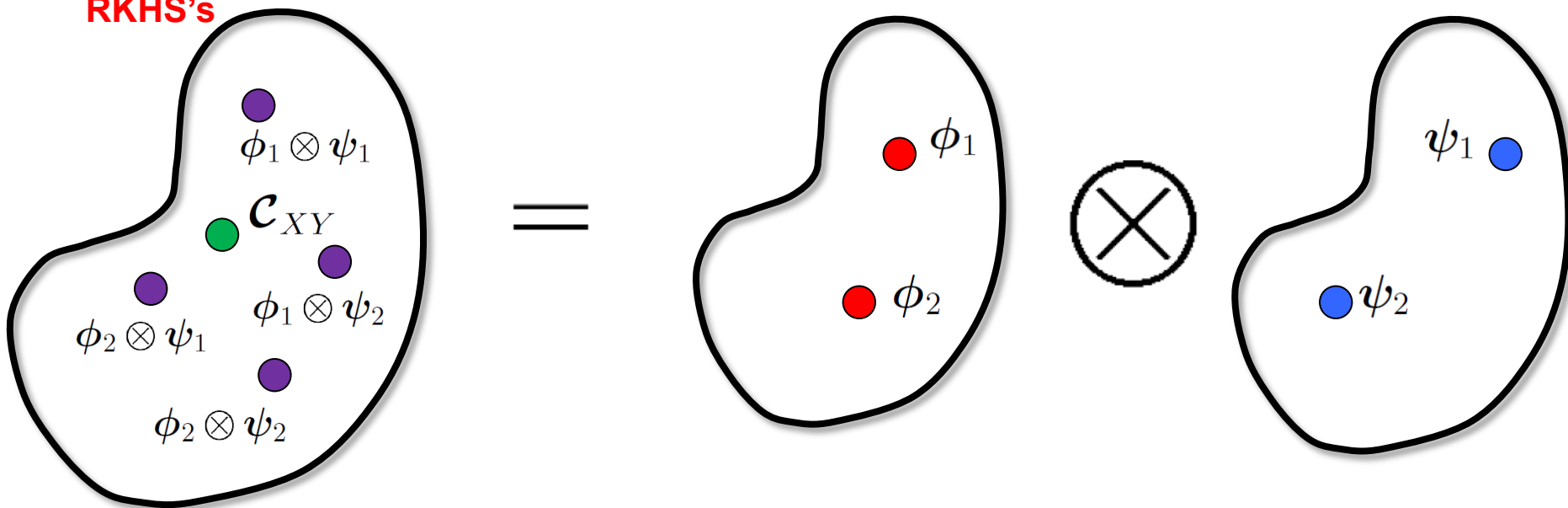
# Review: Cross Covariance Operator [Smola et al. 2007]

$$\mathcal{C}_{XY} = \mathbb{E}_{XY}[\phi_X \otimes \psi_Y]$$

**Embed Joint Distribution of *X* and *Y* in the Tensor Product of two RKHS's**

$\mathcal{F}$

$\mathcal{G}$

$\phi_1 \otimes \psi_1$

$\mathcal{C}_{XY}$

$\phi_2 \otimes \psi_1$

$\phi_1 \otimes \psi_2$

$\phi_2 \otimes \psi_2$

$=$

$\phi_1$

$\phi_2$

$\otimes$

$\psi_1$

$\psi_2$
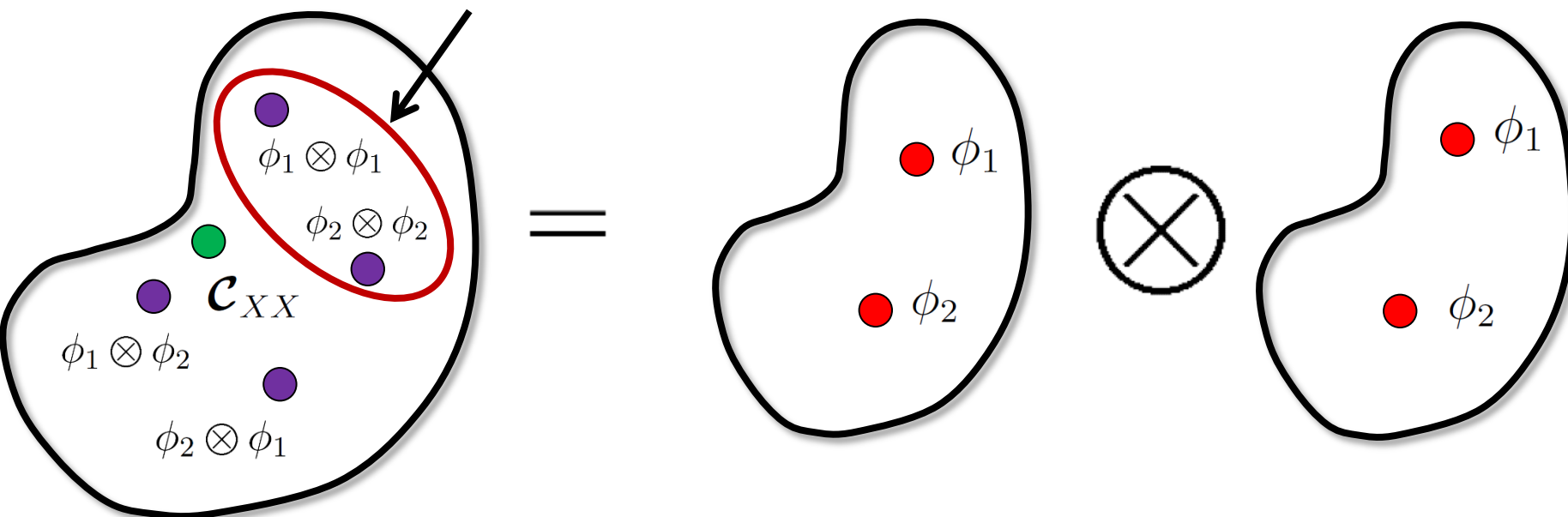
**Embedding of** $\overrightarrow{\mathbb{P}}[X, Y]$

4

# Review: Auto Covariance Operator [Smola et al. 2007]

$$\mathcal{C}_{XX} = \mathbb{E}_X[\phi_X \otimes \phi_X]$$

**Only take expectation over these**

$\phi_1 \otimes \phi_1$

$\phi_2 \otimes \phi_2$

$\mathcal{C}_{XX}$

$\phi_1 \otimes \phi_2$

$\phi_2 \otimes \phi_1$

$\mathcal{F}$

$\phi_1$

$\phi_2$

$\bigotimes$

$\mathcal{F}$

$\phi_1$

$\phi_2$

$=$

**Embedding of** $\text{Diag}(\mathbb{P}[X])$

# Review: Conditional Embedding Operator [Song et al. 2009]

- Conditional Embedding Operator:

$$\mathcal{C}_{X|Y} = \mathcal{C}_{XY}\mathcal{C}_{YY}^{-1}$$

- Has Following Property:

$$\mathbb{E}_{X|y}[\phi_X|y] = \mathcal{C}_{X|Y}\phi_y$$

- Analogous to "Slicing" a Conditional Probability Table in the Discrete Case:

$$\overrightarrow{\mathbb{P}}[X|Y = 1] = \overrightarrow{\mathbb{P}}[X|Y]\boldsymbol{\delta}_1$$
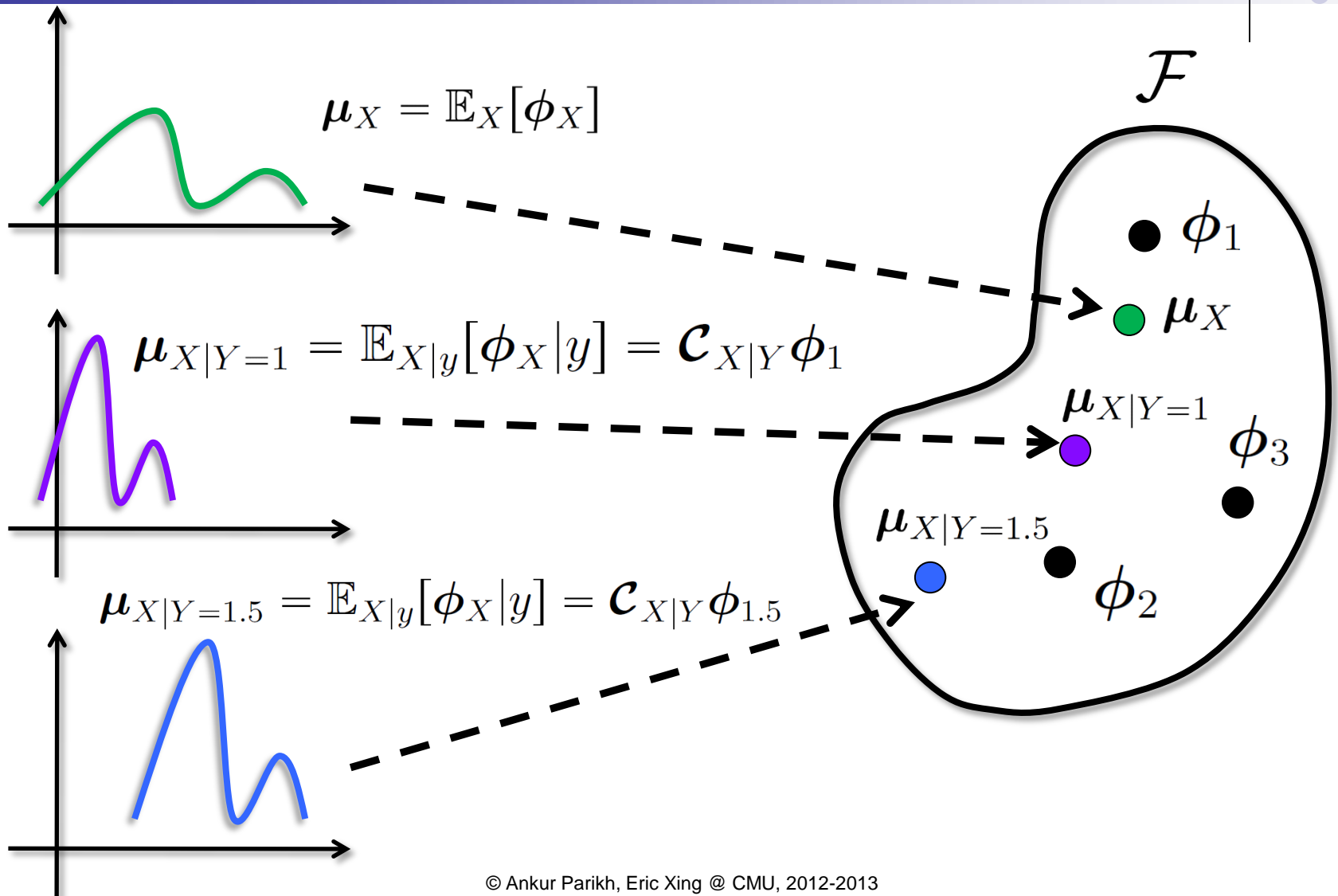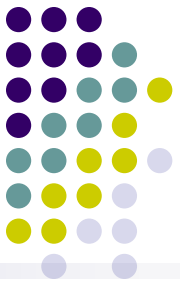
# Slicing the Conditional Probability Matrix

$$\mathcal{P}[X]$$

$$\mathcal{P}[X|Y=1] = \mathcal{P}[X|Y]\boldsymbol{\delta}_1$$

$$\mathcal{P}[X|Y=2] = \mathcal{P}[X|Y]\boldsymbol{\delta}_2$$

# "Slicing" the Conditional Embedding Operator

$$\boldsymbol{\mu}_X = \mathbb{E}_X[\boldsymbol{\phi}_X]$$

$$\boldsymbol{\mu}_{X|Y=1} = \mathbb{E}_{X|y}[\boldsymbol{\phi}_X|y] = \mathcal{C}_{X|Y}\boldsymbol{\phi}_1$$

$$\boldsymbol{\mu}_{X|Y=1.5} = \mathbb{E}_{X|y}[\boldsymbol{\phi}_X|y] = \mathcal{C}_{X|Y}\boldsymbol{\phi}_{1.5}$$

$\mathcal{F}$

$\boldsymbol{\phi}_1$

$\boldsymbol{\mu}_X$

$\boldsymbol{\mu}_{X|Y=1}$

$\boldsymbol{\phi}_3$

$\boldsymbol{\mu}_{X|Y=1.5}$

$\boldsymbol{\phi}_2$

# Why we Like Hilbert Space Embeddings

**We can marginalize and use chain rule in Hilbert Space too!!!**

**Sum Rule:**

$$\mathbb{P}[X] = \int_Y \mathbb{P}[X, Y] = \int_Y \mathbb{P}[X|Y]\mathbb{P}[Y]$$

**Chain Rule:**

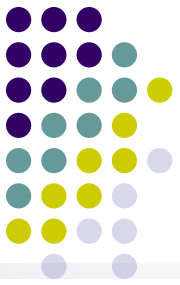$$\mathbb{P}[X, Y] = \mathbb{P}[X|Y]\mathbb{P}[Y] = \mathbb{P}[Y|X]\mathbb{P}[Y]$$

**Sum Rule in RKHS:**

$$\boldsymbol{\mu}_X = \boldsymbol{\mathcal{C}}_{X|Y}\boldsymbol{\mu}_Y$$

**Chain Rule in RKHS:**

$$\boldsymbol{\mathcal{C}}_{YX} = \boldsymbol{\mathcal{C}}_{Y|X}\boldsymbol{\mathcal{C}}_{XX} = \boldsymbol{\mathcal{C}}_{X|Y}\boldsymbol{\mathcal{C}}_{YY}$$

**We will prove these now**

# Sum Rules

- The sum rule can be expressed in two ways:

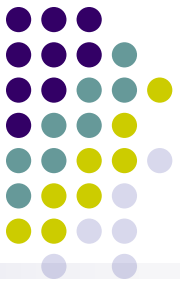- First way:

$$\mathbb{P}[X] = \sum_Y \mathbb{P}[X, Y]$$

**Does not work in RKHS, since there is no "sum" operation for an operator**

- Second way:

$$\mathbb{P}[X] = \sum_Y \mathbb{P}[X|Y]\mathbb{P}[Y]$$

**Works in RKHS!!!**

- What is special about the second way? Intuitively, it can be expressed elegantly as matrix multiplication ☺

# Sum Rule (Matrix Form)

- Sum Rule

$$\mathbb{P}[X] = \sum_Y \mathbb{P}[X|Y]\mathbb{P}[Y]$$

- Equivalent view using Matrix Algebra

$$\mathcal{P}[X] = \mathcal{P}[X|Y] \times \mathcal{P}[Y]$$

$$\begin{pmatrix} \mathbb{P}[X=0] \\ \mathbb{P}[X=1] \end{pmatrix} = \begin{pmatrix} \mathbb{P}[X=0|Y=0] & \mathbb{P}[X=0|Y=1] \\ \mathbb{P}[X=1|Y=0] & \mathbb{P}[X=1|Y=1] \end{pmatrix} \times \begin{pmatrix} \mathbb{P}[Y=0] \\ \mathbb{P}[Y=1] \end{pmatrix}$$
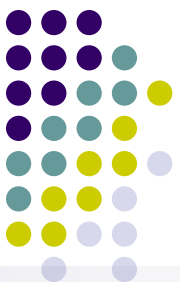
# Important Notation for this Lecture

- We will use the calligraphic P to denote that the probability is being treated as a matrix/vector/tensor

- Probabilities

$$\mathbb{P}[X, Y] = \mathbb{P}[X|Y]\mathbb{P}[Y]$$

- Probability Vectors/Matrices/Tensors

$$\mathcal{P}[X] = \mathcal{P}[X|Y]\mathcal{P}[Y]$$

# Chain Rule (Matrix Form)

- Chain Rule

$$\mathbb{P}[X, Y] = \mathbb{P}[X|Y]\mathbb{P}[Y] = \mathbb{P}[Y|X]\mathbb{P}[Y]$$

**Means on diagonal**

- Equivalent view using Matrix Algebra

$$\boldsymbol{\mathcal{P}}[X, Y] = \boldsymbol{\mathcal{P}}[X|Y] \times \boldsymbol{\mathcal{P}}[\oslash Y]$$

$$\begin{pmatrix} \mathbb{P}[X=0, Y=0] & \mathbb{P}[X=0, Y=1] \\ \mathbb{P}[X=1, Y=0] & \mathbb{P}[X=1, Y=1] \end{pmatrix} =$$

$$\begin{pmatrix} \mathbb{P}[X=0|Y=0] & \mathbb{P}[X=0|Y=1] \\ \mathbb{P}[X=1|Y=0] & \mathbb{P}[X=1|Y=1] \end{pmatrix} \times \begin{pmatrix} \mathbb{P}[Y=0] & 0 \\ 0 & \mathbb{P}[Y=1] \end{pmatrix}$$

- Note how diagonal is used to keep **Y** from being marginalized out.

# Example

- What about?

$$\boldsymbol{\mathcal{P}}[B|A]\boldsymbol{\mathcal{P}}[\oslash A]\boldsymbol{\mathcal{P}}[C|A]^{\top}$$

$$\boldsymbol{\mathcal{P}}[B,C]$$

- **Only if B and C are conditionally independent given A!!!**
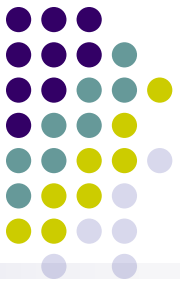
# Different Proof of Matrix Sum Rule with Expectations

- Let's now derive the matrix sum rule differently.

- Let $\boldsymbol{\delta_i}$ denote an indicator vector, that is 1 in the $i^{th}$ position.

$$\boldsymbol{\delta}_0 = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \quad \boldsymbol{\delta}_1 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

$$\boldsymbol{\mathcal{P}}[X] = \mathbb{E}_X[\boldsymbol{\delta}_X] = \mathbb{P}[X = 0]\boldsymbol{\delta}_0 + \mathbb{P}[X = 1]\boldsymbol{\delta}_1$$

$$\begin{aligned} \boldsymbol{\mathcal{P}}[X|Y = y] &= \mathbb{E}_{X|Y=y}[\boldsymbol{\delta}_X] \\ &= \mathbb{P}[X = 0|Y = y]\boldsymbol{\delta}_0 + \mathbb{P}[X = 1|Y = y]\boldsymbol{\delta}_1 \end{aligned}$$

# Random Variables?

$$\mathcal{P}[X] = \mathbb{E}_X[\delta_X]$$

**Remember this is a probability vector. It is not a random variable.**

**This is a random vector**

$$\boldsymbol{\mu}_X = \mathbb{E}_X[\boldsymbol{\phi}_X]$$

**Similarly, this is a function in an RKHS. It is not a random variable.**

**This is a random function**

# Expectation Proof of Matrix Sum Rule Cont.

$$\mathcal{P}[X|Y]\mathcal{P}[Y] = \mathcal{P}[X|Y]\mathbb{E}_Y[\boldsymbol{\delta}_Y]$$

$$= \mathbb{E}_Y[\mathcal{P}[X|Y]\boldsymbol{\delta}_Y]$$

$$= \mathbb{E}_Y[\mathbb{E}_{X|Y}[\boldsymbol{\delta}_X]]$$

$$= \mathbb{E}_{XY}[\boldsymbol{\delta}_X]$$

$$= \mathcal{P}[X]$$

**This is a conditional probability matrix, so it is not a random variable (despite the misleading notation), and thus the Expectation can be pulled out**

**This is a random variable**

# Proof of RKHS Sum Rule

- Now apply the same technique to the RKHS Case.

$$\mathcal{C}_{X|Y}\boldsymbol{\mu}_Y$$

$$= \mathcal{C}_{X|Y}\mathbb{E}_Y[\boldsymbol{\psi}_Y]$$

$$= \mathbb{E}_Y[\mathcal{C}_{X|Y}\boldsymbol{\psi}_Y] \quad \textbf{Move expectation outside}$$

$$= \mathbb{E}_Y[\mathbb{E}_{X|Y}[\boldsymbol{\phi}_X|Y]] \quad \textbf{Property of conditional embedding}$$

$$= \mathbb{E}_{XY}[\boldsymbol{\phi}_X] \quad \textbf{Property of Expectation}$$

$$= \boldsymbol{\mu}_X \quad \textbf{Definition of Mean Map}$$

# Kernel Graphical Models [Song et al. 2010, Song et al. 2011]

- The idea is to replace the CPTs with RKHS operators/functions.

- Let's do this for a simple example first.



- **We would like to compute** $\mathbb{P}[A = a, D = d]$

# Consider the Discrete Case

A $\qquad$ B $\qquad$ C $\qquad$ D

$$\mathcal{P}(A) \qquad \mathcal{P}(B|A) \qquad \mathcal{P}(C|B) \qquad \mathcal{P}(D|C)$$

# Inference as Matrix Multiplication

A          B          C          D

$$\boldsymbol{\mathcal{P}}(D) = \boldsymbol{\mathcal{P}}(A)\boldsymbol{\mathcal{P}}(B|A)^{\top}\boldsymbol{\mathcal{P}}(C|B)^{\top}\boldsymbol{\mathcal{P}}(D|C)^{\top}$$

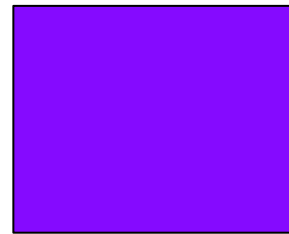**Oops....we accidentally integrated out A**

# Put A on Diagonal Instead



$$\mathcal{P}(\oslash A) \quad \mathcal{P}(B|A) \quad \mathcal{P}(C|B) \quad \mathcal{P}(D|C)$$

# Now it works



$$\boldsymbol{\mathcal{P}}(A, D) = \boldsymbol{\mathcal{P}}(\oslash A)\boldsymbol{\mathcal{P}}(B|A)^\top \boldsymbol{\mathcal{P}}(C|B)^\top \boldsymbol{\mathcal{P}}(D|C)^\top$$
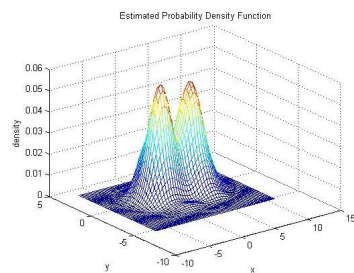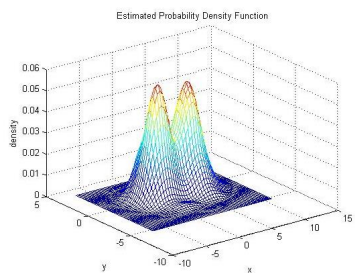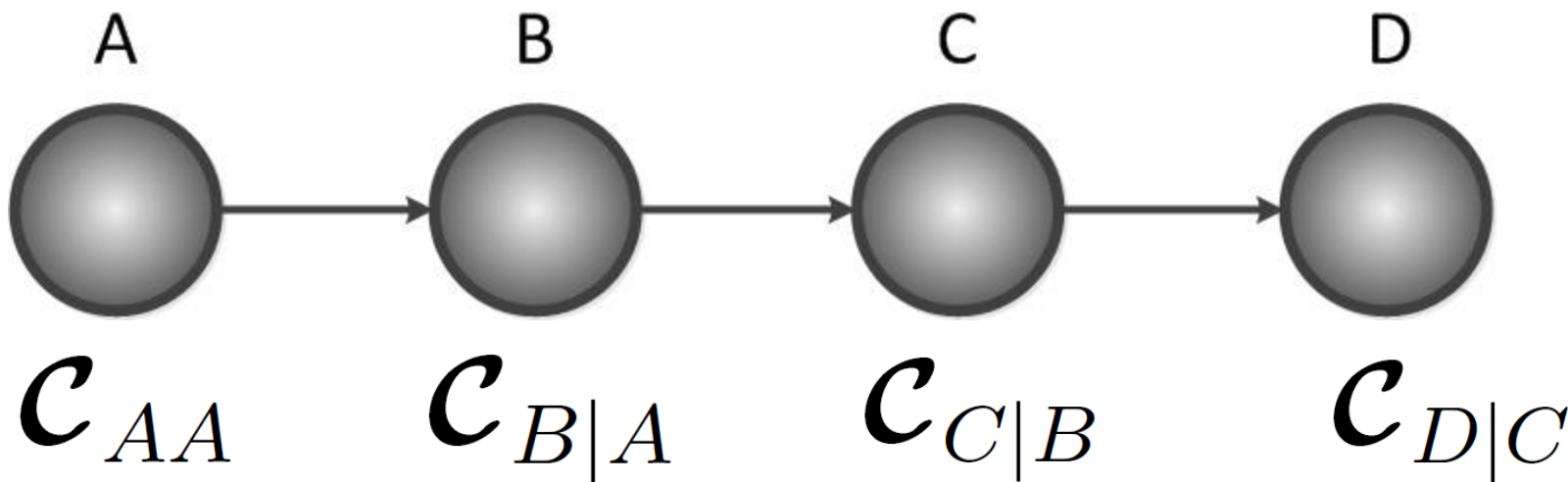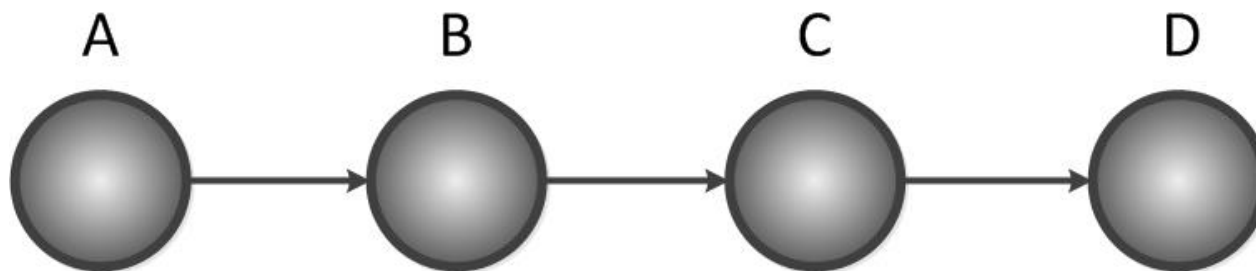
# Introducing evidence



A    B    C    D

- Introduce evidence with delta vectors

$$\mathcal{P}(A = a, D = d) = \delta_a^\top \mathcal{P}(A, D) \delta_d$$

# Now with Kernels



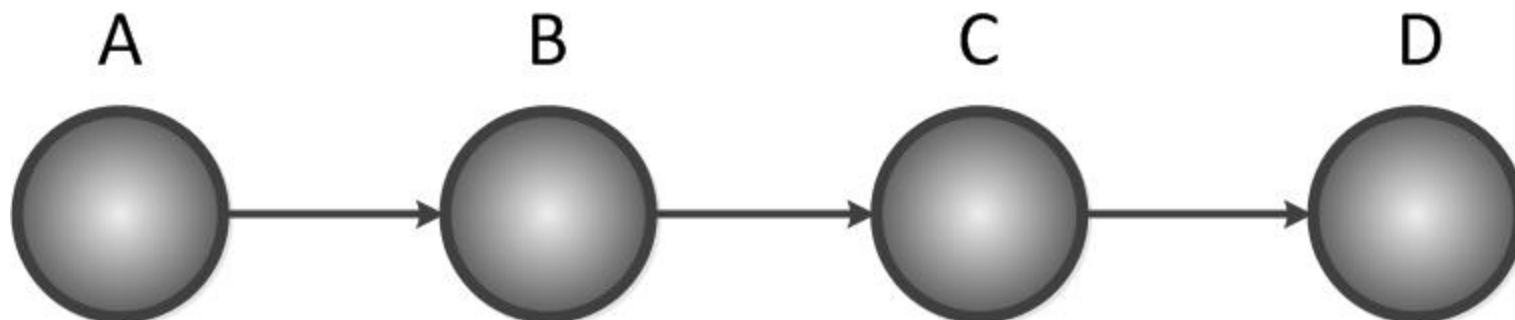$$\mathcal{C}_{AA} \qquad \mathcal{C}_{B|A} \qquad \mathcal{C}_{C|B} \qquad \mathcal{C}_{D|C}$$

# Sum-Product with Kernels



$$\mathcal{C}_{AB} = \mathcal{C}_{AA}\mathcal{C}_{B|A}^{\top}$$

$$\mathcal{C}_{AD} = \mathcal{C}_{AA}\mathcal{C}_{B|A}^{\top}\mathcal{C}_{B|C}^{\top}\mathcal{C}_{C|D}^{\top}$$

# Sum-Product with Kernels



$$\text{some number} = \phi_a^\top \boldsymbol{\mathcal{C}}_{A,D} \phi_d$$

# What does it mean to evaluate the mean map at a point?

- Consider just evaluating one random variable *X* at a particular evidence value using the Gaussian RBF Kernel:

$$
\begin{aligned}
\langle \boldsymbol{\mu}_X, \boldsymbol{\phi}_{\bar{x}} \rangle &= \mathbb{E}_X \left[ \langle \boldsymbol{\phi}_X, \boldsymbol{\phi}_{\bar{x}} \rangle \right] \\
&= \mathbb{E}_X \left[ \boldsymbol{K}(X, \bar{x}) \right] \\
&= \mathbb{E}_X \left[ \exp \left( \frac{-\|X - \bar{x}\|_2^2}{\sigma^2} \right) \right]
\end{aligned}
$$

- What does this looks like?

# Kernel Density Estimation!
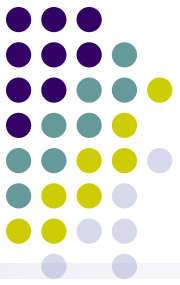
- Consider Kernel Density Estimate at point $\bar{x}$:

$$\mathbb{P}_{kde}[X = \bar{x}] \propto \mathbb{E}\left[\exp\left(\frac{-\|X - \bar{x}\|_2^2}{\sigma^2}\right)\right]$$

- And its empirical estimate:

$$\widehat{\mathbb{P}}_{kde}[X = \bar{x}] \propto \frac{1}{N}\sum_{n=1}^{N}\exp\left(-\frac{\|X - \bar{x}\|_2^2}{\sigma^2}\right)$$

- So evaluating the mean map at a point is like an unnormalized kernel density estimate. To find the "MAP" assignment, we can evaluate on a grid of points, and then pick the one with the highest value.
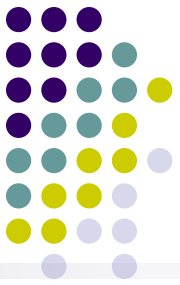
# Multiple Variables

- Kernel Density Estimation with Gaussian RBF Kernel in Multiple Variables is:

$$\mathbb{P}_{kde}[\boldsymbol{X}_{1:\mathcal{O}} = \bar{x}_{1:\mathcal{O}}] \propto \mathbb{E}\left[\prod_{o=1}^{\mathcal{O}} \exp\left(-\frac{\|X_O - \bar{x}_o\|_2^2}{\sigma^2}\right)\right]$$

- Like evaluating a "Huge" Covariance Operator using Gaussian RBF Kernel (without normalization):

$$\langle \boldsymbol{\mathcal{C}}_{X_1,\ldots,X_{\mathcal{O}}}, \boldsymbol{\phi}_{\bar{x}_1} \otimes \boldsymbol{\phi}_{\bar{x}_2} \ldots \otimes \boldsymbol{\phi}_{\bar{x}_{\mathcal{O}}} \rangle$$

# What is the problem with this?

- The empirical estimate is very inaccurate because of curse of dimensionality

$$\widehat{\mathbb{P}}_{kde}[\boldsymbol{X}_{1:\mathcal{O}} = \bar{x}_{1:\mathcal{O}}] \propto \frac{1}{N} \sum_{n=1}^{N} \prod_{o=1}^{\mathcal{O}} \exp\left(-\frac{\|X_O^{(n)} - \bar{x}_o\|_2^2}{\sigma^2}\right)$$

- Empirically computing the "huge" covariance operator will have the same problem.

- But then what is the point of Hilbert Space Embeddings?

# We can factorize the "Huge" Covariance Operator

- Hilbert Space Embeddings allow us to factorize the huge covariance operator using the graphical model structure that kernel density estimation does not do.

$$\langle \mathcal{C}_{X_1,\ldots,X_{\mathcal{O}}}, \phi_{\bar{x}_1} \otimes \phi_{\bar{x}_2} \ldots \otimes \phi_{\bar{x}_{\mathcal{O}}} \rangle$$

**Factorizes into smaller covariance/conditional embedding operators using the graphical model that are more efficient to estimate.**

$$\mathcal{C}_{AA} \qquad \mathcal{C}_{B|A} \qquad \mathcal{C}_{C|B} \qquad \mathcal{C}_{D|C}$$
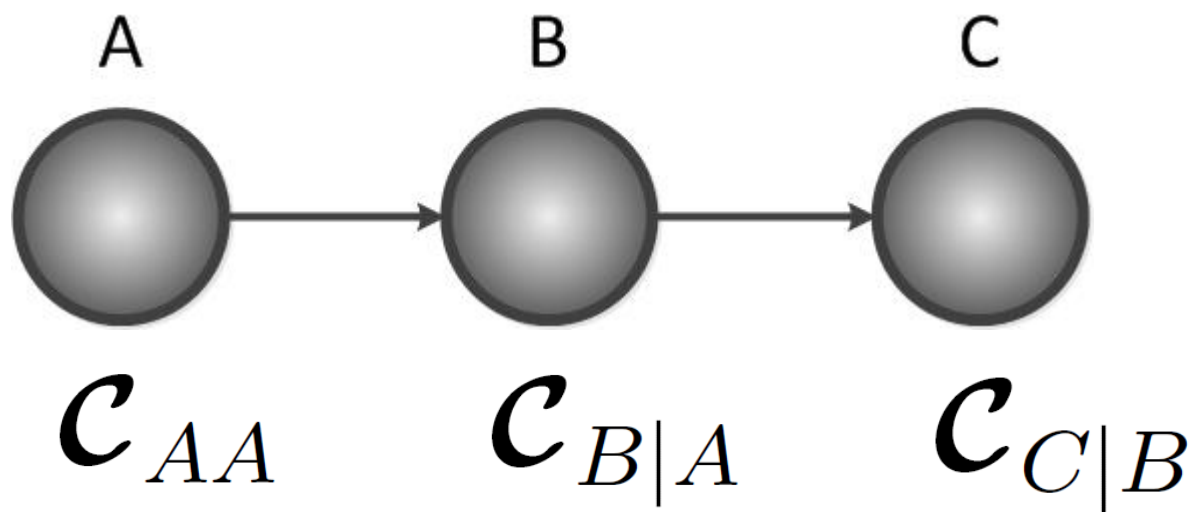
# Kernel Graphical Models: The Overall Picture

Naïve way to represent joint distribution of discrete variables is to store and manipulate a "huge" probability table.

Discrete Graphical Models allow us to factorize the "huge" joint distribution table into smaller factors.

Naïve way to represent joint distribution for many continuous variables is to use multivariate kernel density estimation.

Kernel Graphical Models allow us to factorize joint distributions of continuous variables into smaller factors.

# Consider an Even Simpler Graphical Model



$$A \qquad B \qquad C$$

$$\mathcal{C}_{AA} \qquad \mathcal{C}_{B|A} \qquad \mathcal{C}_{C|B}$$

**We are going to show how to estimate these operators from data.**

# The Kernel Matrix

$$K_{XX} = \begin{bmatrix} \langle \phi_{x_1}, \phi_{x_1} \rangle & \cdots & \langle \phi_{x_1}, \phi_{x_N} \rangle \\ & & \\ & & \\ \langle \phi_{x_N}, \phi_{x_1} \rangle & \cdots & \langle \phi_{x_N}, \phi_{x_N} \rangle \end{bmatrix}$$

$N$

$N$

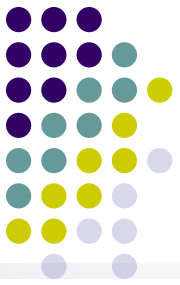# Empirical Estimate Auto Covariance

$$\mathcal{C}_{XX} = \mathbb{E}_X[\phi_X \otimes \phi_X]$$

$$\hat{\mathcal{C}}_{XX} = \frac{1}{N} \sum_{n=1}^{N} \phi_{x_n} \otimes \phi_{x_n}$$

$$\hat{\mathcal{C}}_{XX} = \frac{1}{N} \Phi_X \Phi_X^{\top}$$

**Defined on next slide**

# Conceptually,

$$\mathbf{\Phi}_X := \begin{bmatrix} \phi_{x_1} \phi_{x_2} \cdots \phi_{x_n} \end{bmatrix}$$

$N$

$\infty$

# Conceptually,

$$\sum_{n=1}^{N} \boldsymbol{v}_i \boldsymbol{\phi}_{x_n} = \begin{bmatrix} \boldsymbol{\phi}_{x_1} \boldsymbol{\phi}_{x_2} & \cdots & \boldsymbol{\phi}_{x_n} \end{bmatrix} \begin{pmatrix} \boldsymbol{v}_1 \\ \boldsymbol{v}_2 \\ \ldots \\ \boldsymbol{v}_N \end{pmatrix}$$

$\boldsymbol{\Phi}_X$

# Conceptually,

$$\begin{pmatrix} \phi_{x_1}^\top f \\ \phi_{x_2}^\top f \\ \dots \\ \phi_{x_n}^\top f \end{pmatrix} = \Phi_X^\top \begin{bmatrix} \longleftarrow \phi_{x_1} \longrightarrow \\ \longleftarrow \phi_{x_2} \longrightarrow \\ \dots \\ \longleftarrow \phi_{x_n} \longrightarrow \end{bmatrix} f$$

# Rigorously,

$\mathbf{\Phi}_X$ **is an operator that maps vectors in** $\mathbb{R}^N$ **to functions in** $\mathcal{F}$

**such that:**
$$\sum_{n=1}^{N} v_i \phi_{x_n} = \mathbf{\Phi}_X v$$

**Its adjoint (transpose)** $\mathbf{\Phi}_X^\top$ **can then be derived to be:**

$$\begin{pmatrix} \langle \phi_{x_1}, f \rangle \\ \langle \phi_{x_2}^\top, f \rangle \\ \dots \\ \langle \phi_{x_n}^\top, f \rangle \end{pmatrix} = \mathbf{\Phi}_X^\top f$$

# Empirical Estimate Cross Covariance

$$\mathcal{C}_{YX} = \mathbb{E}[\phi_Y \otimes \phi_X]$$

$$\hat{\mathcal{C}}_{YX} = \frac{1}{N} \sum_{n=1}^{N} \phi_{y_n} \otimes \phi_{x_n}$$

$$\hat{\mathcal{C}}_{YX} = \frac{1}{N} \mathbf{\Phi}_Y \mathbf{\Phi}_X^\top$$

# Getting the Kernel Matrix

- It can then be shown that,

$$\boldsymbol{\Phi}_X^\top \boldsymbol{\Phi}_X = \boldsymbol{K}_{XX} \qquad \boldsymbol{K}_{XX}(i,j) := \langle \boldsymbol{\phi}_{x_i}, \boldsymbol{\phi}_{x_j} \rangle$$

- This is finite and easy to compute!! ☺

- However, note that the estimates of the covariance operators are **not** finite since:

$$\hat{\mathcal{C}}_{XX} = \frac{1}{N} \boldsymbol{\Phi}_X \boldsymbol{\Phi}_X^\top$$

# Intuition 1: Why the Kernel Trick works

$$\hat{c}_{XX} = \frac{1}{N} \sum_{n=1}^{N} \phi_{x_n} \otimes \phi_{x_n}$$

$$\hat{c}_{XX} = \frac{1}{N} \Phi_X \Phi_X^\top$$

**This operator is infinite dimensional but it has at most rank N**

$$\Phi_X^\top \Phi_X = \boldsymbol{K}_{XX}$$

**The kernel matrix is N by N, and thus the kernel trick is exploiting the low rank structure**

# Empirical Estimate of Conditional Embedding Operator

$$\hat{\mathcal{C}}_{Y|X} = \hat{\mathcal{C}}_{YX} \hat{\mathcal{C}}_{XX}^{-1} \quad \textbf{?}$$

**Sort of……**

**We need to regularize so that this is invertible**

$$\hat{c}_{YX} \qquad \hat{c}_{XX} \qquad \textbf{regularizer}$$

$$\hat{\mathcal{C}}_{Y|X} = \mathbf{\Phi}_Y \mathbf{\Phi}_X^\top (\frac{1}{N} \mathbf{\Phi}_X \mathbf{\Phi}_X^\top + \lambda I)^{-1}$$

# Return of Matrix Inversion Lemma

- Believe it or not, the matrix inversion lemma works for linear, bounded operators too.

$$(E - FH^{-1}G)^{-1} = E^{-1} + E^{-1}F(H - GE^{-1}F)^{-1}GE^{-1}$$

- Using it we get,

$$\hat{\mathcal{C}}_{Y|X} = \Phi_Y(\Phi_X^\top \Phi_X + \lambda N \boldsymbol{I})^{-1}\Phi_X^\top$$

$$\hat{\mathcal{C}}_{Y|X} = \Phi_Y(\boldsymbol{K}_{XX} + \lambda N \boldsymbol{I})^{-1}\Phi_X^\top$$

# But Our estimates are still Infinite….



$$\hat{\mathcal{C}}_{AA} = \frac{1}{N}\boldsymbol{\Phi}_A\boldsymbol{\Phi}_A^{\top}$$

$$\hat{\mathcal{C}}_{B|A} = \boldsymbol{\Phi}_B(\boldsymbol{K}_{AA} + \lambda N\boldsymbol{I})^{-1}\boldsymbol{\Phi}_A^{\top}$$

$$\hat{\mathcal{C}}_{C|B} = \boldsymbol{\Phi}_C(\boldsymbol{K}_{BB} + \lambda N\boldsymbol{I})^{-1}\boldsymbol{\Phi}_B^{\top}$$

**Lets do inference and see what happens.**

# Running Inference



$$\hat{\mathcal{C}}_{AC} = \hat{\mathcal{C}}_{AA} \hat{\mathcal{C}}_{B|A}^{\top} \hat{\mathcal{C}}_{C|B}^{\top}$$

$$\hat{\mathcal{C}}_{AC} = \frac{1}{N} \boldsymbol{\Phi}_A \boxed{\boldsymbol{\Phi}_A^{\top} \boldsymbol{\Phi}_A} (\boldsymbol{K}_{AA} + \lambda N \boldsymbol{I})^{-1} \boxed{\boldsymbol{\Phi}_B^{\top} \boldsymbol{\Phi}_B} (\boldsymbol{K}_{BB} + \lambda N \boldsymbol{I})^{-1} \boldsymbol{\Phi}_C^{\top}$$

$$\boldsymbol{K}_{AA} \qquad \boldsymbol{K}_{BB}$$

# Incorporating the Evidence

$$\phi_a^\top \hat{\boldsymbol{\mathcal{C}}}_{AC} \phi_c =$$

$$\frac{1}{N} \boxed{\phi_a^\top \boldsymbol{\Phi}_A} \boldsymbol{K}_{AA} (\boldsymbol{K}_{AA} + \lambda N \boldsymbol{I})^{-1} \times$$

$$\boldsymbol{K}_{BB} (\boldsymbol{K}_{BB} + \lambda N \boldsymbol{I})^{-1} \boxed{\boldsymbol{\Phi}_C^\top \phi_c}$$

$$\boldsymbol{K}_{AA}(1:N, a)$$

$$\boldsymbol{K}_{CC}(1:N, c)$$

# Reparameterize the Model

**Finite!!!!!**

**A**

$$\hat{\mathcal{C}}_{AA}$$

$$\Rightarrow \quad \hat{\mathcal{D}}_A = \frac{1}{N} \boldsymbol{K}_{AA}$$

**B**

$$\hat{\mathcal{C}}_{B|A}$$

$$\Rightarrow \quad \hat{\mathcal{D}}_B = (\boldsymbol{K}_{AA} + \lambda N \boldsymbol{I})^{-1} \boldsymbol{K}_{BB}$$

**C**

$$\hat{\mathcal{C}}_{C|B}$$

$$\Rightarrow \quad \hat{\mathcal{D}}_C = (\boldsymbol{K}_{BB} + \lambda N \boldsymbol{I})^{-1}$$

**Evidence:**

$$\phi_a$$
$$\phi_c$$

$$\Rightarrow \quad \boldsymbol{K}_{AA}(1:N, a)$$
$$\boldsymbol{K}_{CC}(1:N, c)$$

# Intuition 2: Why the Kernel Trick Works

$$K_{XX} = \begin{bmatrix} \langle \phi_{x_1}, \phi_{x_1} \rangle & \cdots & \langle \phi_{x_n}, \phi_{x_n} \rangle \\ & & \\ & & \\ \langle \phi_{x_n}, \phi_{x_1} \rangle & \cdots & \langle \phi_{x_n}, \phi_{x_n} \rangle \end{bmatrix}$$

$N \times N$

# Intuition 2: Why the Kernel Trick Works

**Evaluating a feature function at the N data points!!!**

$$K_{XX} = \begin{bmatrix} \phi_{x_1}(x_1) & \cdots & \phi_{x_1}(x_N) \\ & & \\ & & \\ & & \\ \phi_{x_N}(x_1) & \cdots & \phi_{x_N}(x_N) \end{bmatrix}$$
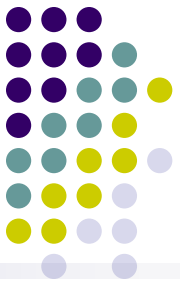
# Intuition 2: Why the Kernel Trick Works

- Generally people interpret the kernel matrix to be a similarity matrix.

- However, we can also view each row of the kernel matrix as evaluating a function at the N data points.

- Although the function may be continuous and not easily represented analytically, we only really care about what its value is on the N data points.

- Thus, when we only have a finite amount of data, the computation should be inherently finite.
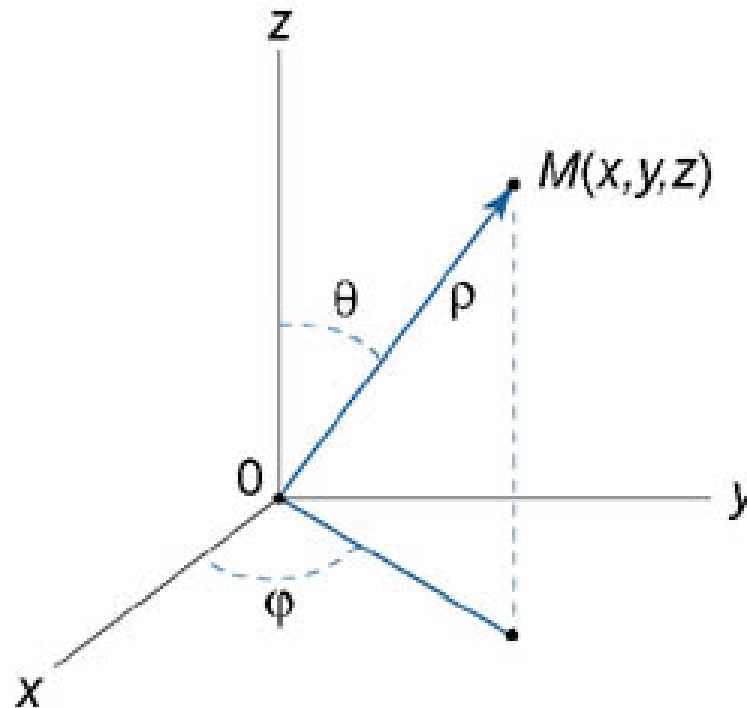
# Protein Sidechains

Goal is to predict the 3D configuration of each sidechain

http://t3.gstatic.com/images?q=tbn:ANd9GcS_nfJy1o9yrDt3
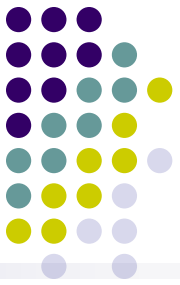7YlpK7i5s0f7QFqhPrG7-1CLm2AfWNt5wCE50plKNZd0

# Protein Sidechains

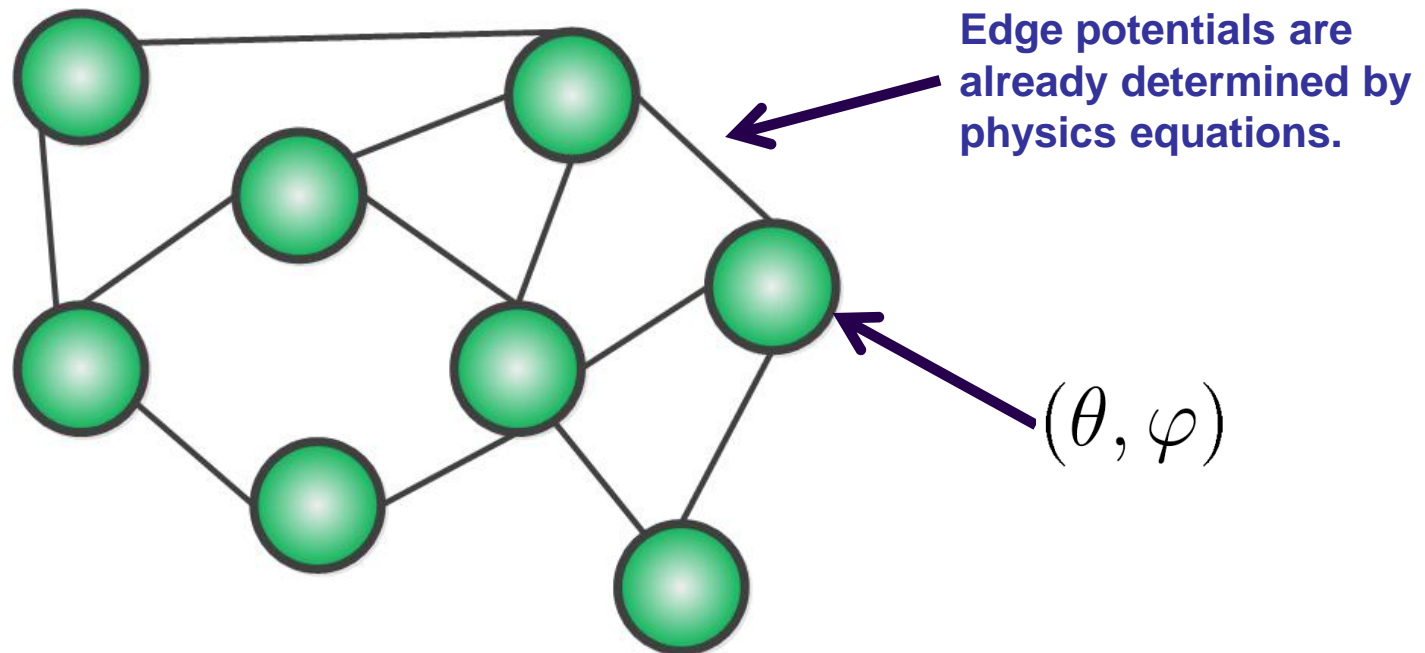- 3D configuration of the sidechain is determined by two angles (spherical coordinates).



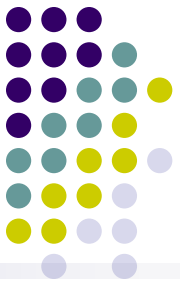**http://www.math24.net/images/triple-int23.jpg**

# The Graphical Model

- Construct a Markov Random Field.

- Each side-chain angle pair is a node. There is an edge between side-chains that are nearby in the protein.



**Edge potentials are already determined by physics equations.**
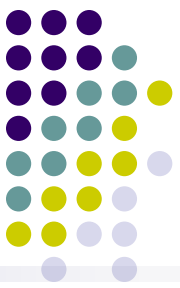
$$(\theta, \varphi)$$

# The Graphical Model

- Goal is to find the MAP assignment of all the sidechain angle pairs.

- Note that this is not Gaussian. But it is easy to define a kernel between angle pairs:
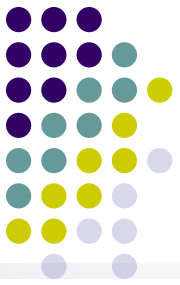
$$\boldsymbol{K}(\boldsymbol{p}_i, \boldsymbol{p}_j) = \exp\left(\boldsymbol{p}_i^\top \boldsymbol{p}_j\right)$$
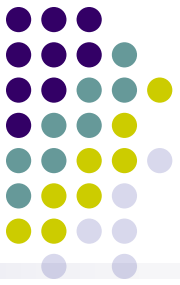
- Can then run Kernel Belief Propagation ☺

# References

- Smola, A. J., Gretton, A., Song, L., and Schölkopf, B., A Hilbert Space Embedding for Distributions, Algorithmic Learning Theory, E. Takimoto (Eds.), Lecture Notes on Computer Science, Springer, 2007.

- L. Song. Learning via Hilbert space embedding of distributions. PhD Thesis 2008.

- Song, L., Huang, J., Smola, A., and Fukumizu, K., Hilbert space embeddings of conditional distributions, International Conference on Machine Learning, 2009.

- Song, L., Gretton, A., and Guestrin, C., Nonparametric Tree Graphical Models via Kernel Embeddings, Artificial Intelligence and Statistics (AISTATS), 2010.

- Song, L., Gretton, A., Bickson, D., Low, Y., and Guestrin, C., Kernel Belief Propagation, International Conference on Artifical Intelligence and Statistics (AISTATS), 2011.
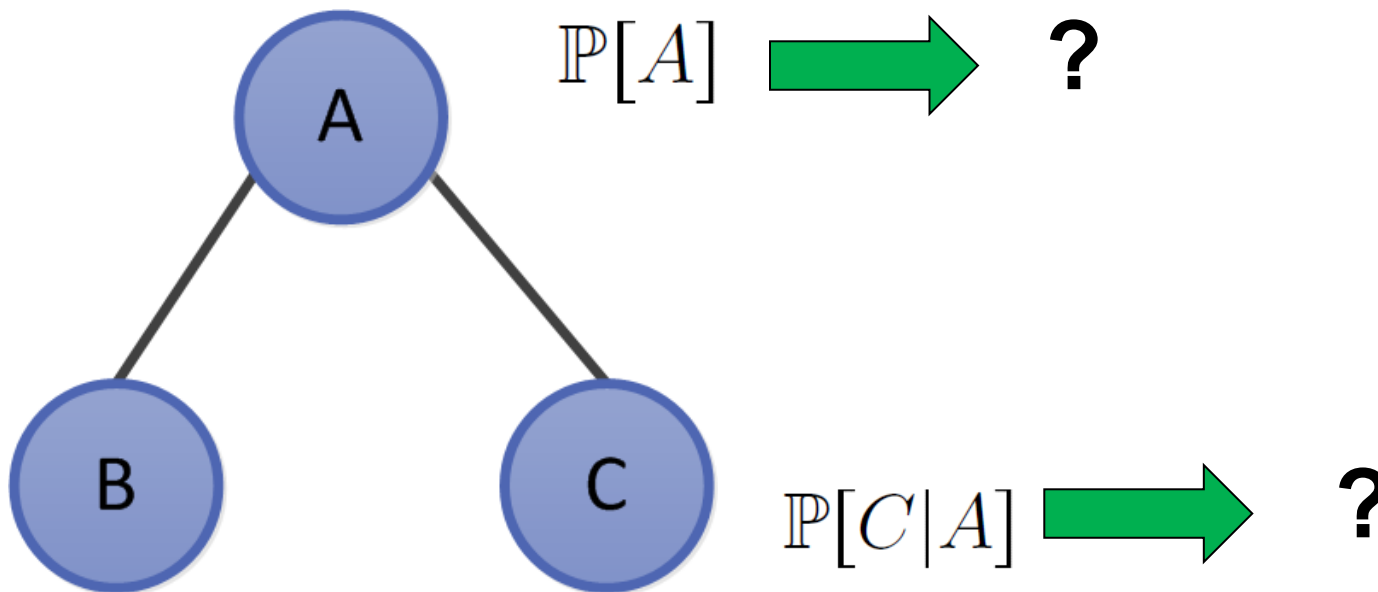
# Supplemental: Kernel Belief Propagation on Trees

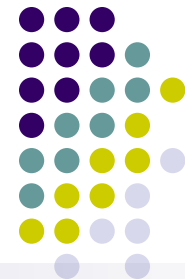# Kernel Tree Graphical Models [Song et al. 2010]

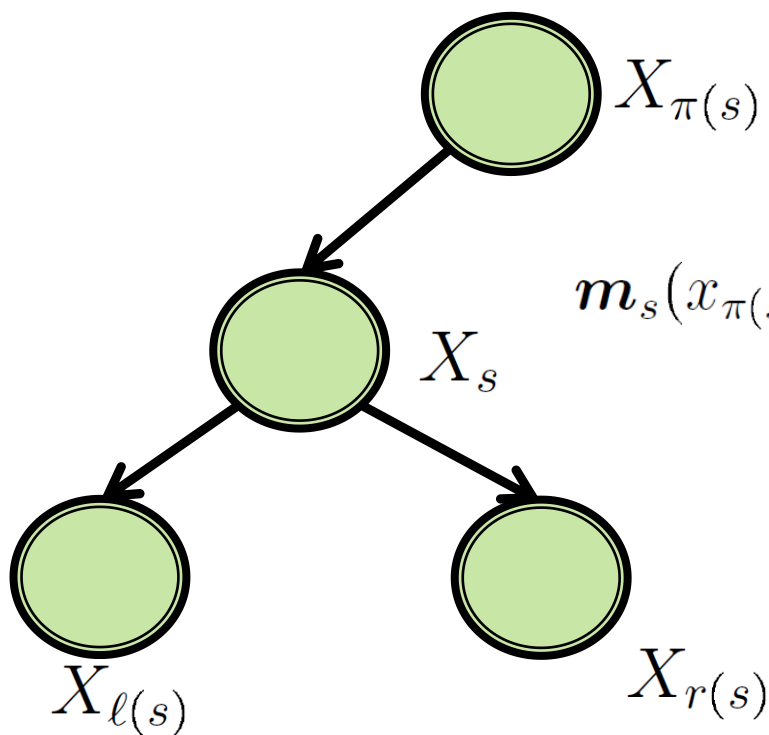- The goal is to somehow replace the CPTs with RKHS operators/functions.



$\mathbb{P}[A]$ ⟹ **?**

$\mathbb{P}[C|A]$ ⟹ **?**

- **But we need to do this in a certain way so that we can still do inference.**

# Message Passing/Belief Propagation

- We need to "matricize" message passing to apply the RKHS trick (but matrices are not enough, we need tensors ☺ )
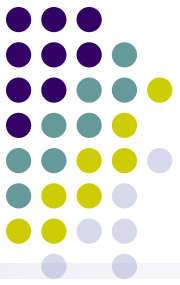
$$\boldsymbol{m}_s(x_{\pi(s)}) = \sum_{X_s} \mathbb{P}[X_s | x_{\pi(s)}] \boldsymbol{m}_{\ell(s)}(X_s) \boldsymbol{m}_{r(s)}(X_s)$$

# Outline

- Show how to represent discrete graphical models using higher order tensors

- Derive *Tensor Message Passing*

- Show how *Tensor Message Passing* can also be derived using Expectations

- Derive **Kernel Message Passing** *[Song et al. 2010]* using the intuition from *Tensor Message Passing* / Expectations

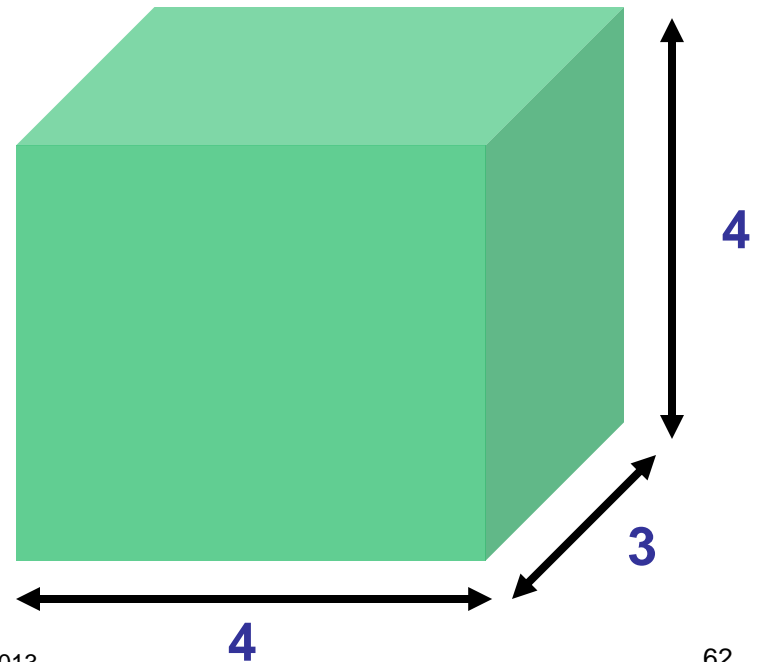- (For simplicity, we will assume a binary tree – all internal nodes have 2 children).

# Tensors

- Multidimensional arrays

- A Tensor of order **N** has **N** modes (**N** indices):
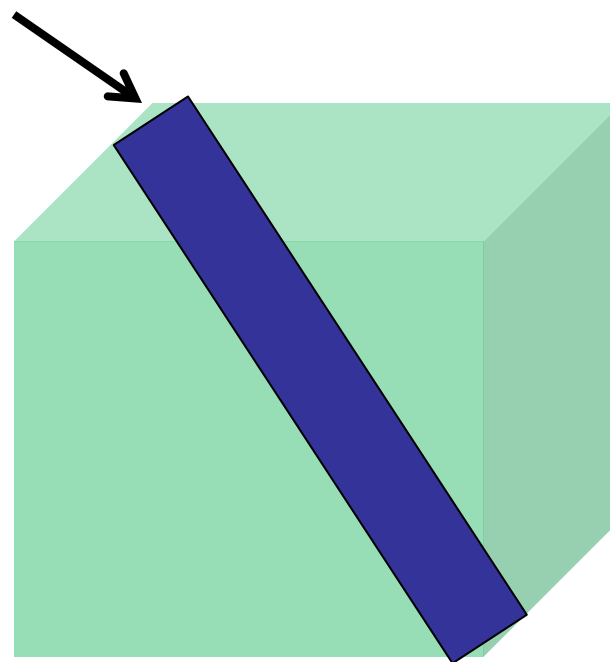
$$\mathcal{T}(i_1, ..., i_N)$$

- Each mode is associated with a dimension. In the example,

  - Dimension of mode 1 is 4
  - Dimension of mode 2 is 3
  - Dimension of mode 3 is 4

**4**

**3**

**4**

# Diagonal Tensors
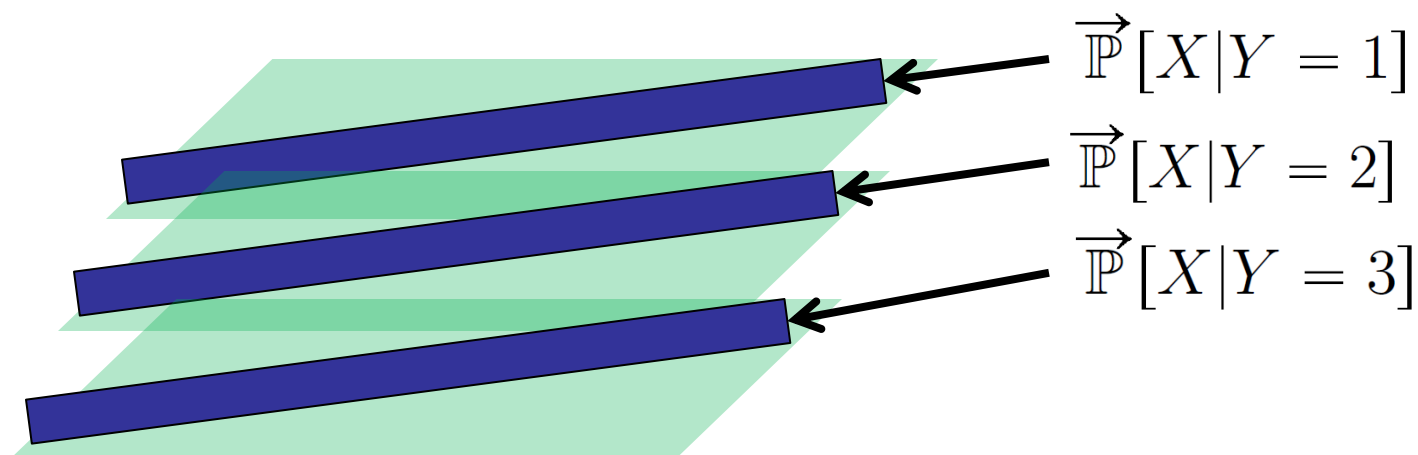
$$\overrightarrow{\mathbb{P}}\,[X]$$

$$\boldsymbol{\mathcal{T}}(i,j,k) = \begin{cases} \mathbb{P}[X=i] & \text{if } i = j = k \\ 0 & \text{otherwise} \end{cases}$$
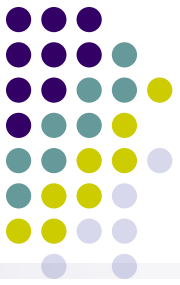
# Partially Diagonal Tensors

$$\boldsymbol{\mathcal{T}}(i, j, k) = \begin{cases} \mathbb{P}[X = i | Y = k] & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$
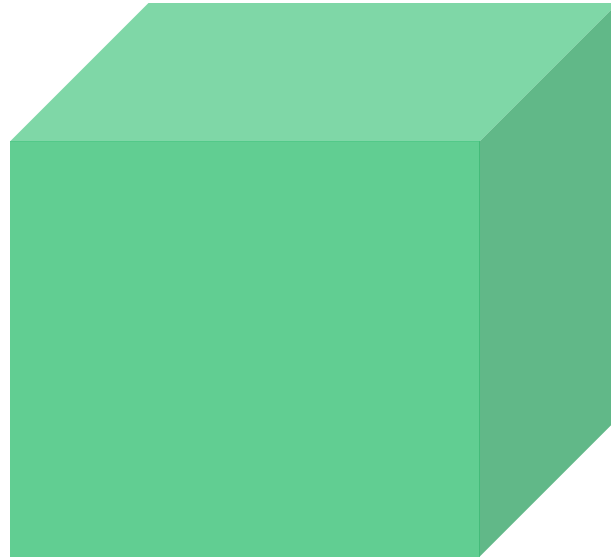


$$\overrightarrow{\mathbb{P}}[X|Y = 1]$$

$$\overrightarrow{\mathbb{P}}[X|Y = 2]$$

$$\overrightarrow{\mathbb{P}}[X|Y = 3]$$

# Tensor Vector Multiplication

- Multiplying a 3$^{rd}$ order tensor by a vector produces a matrix

$$M \qquad \mathcal{T} \qquad \boldsymbol{v}$$
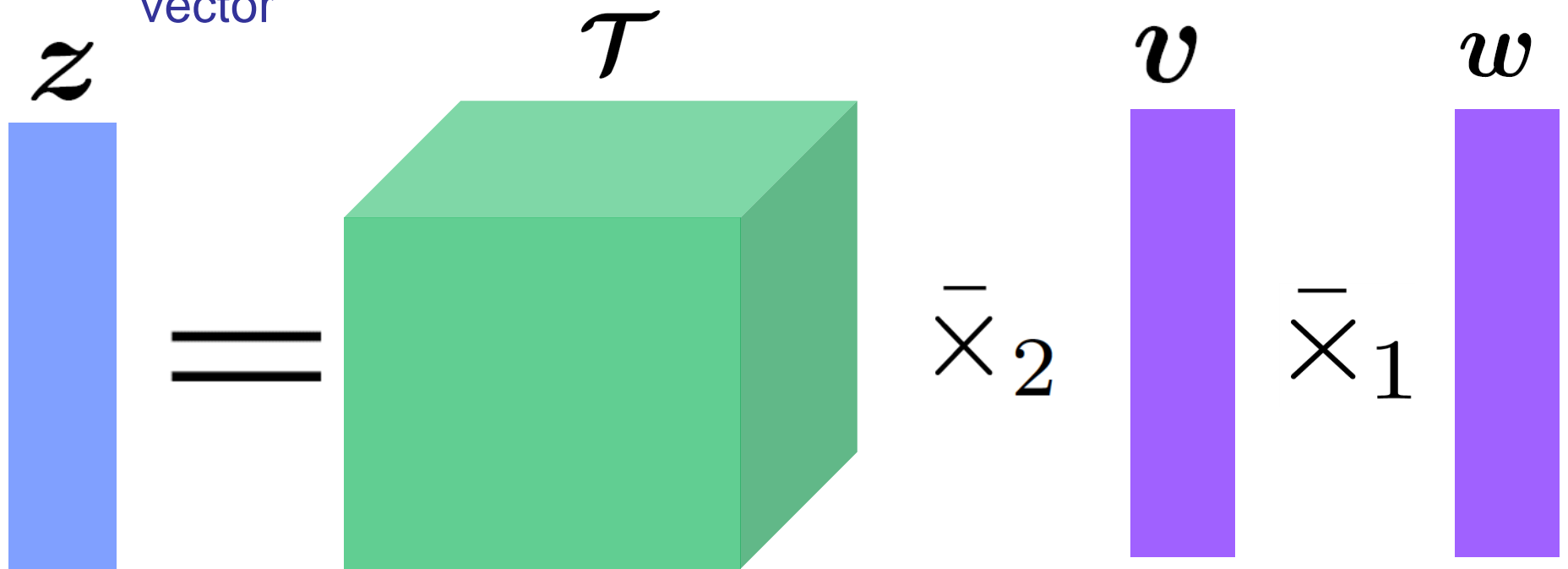
$$= \qquad \bar{\times}_1$$

$$M(j,k) = \sum_i \mathcal{T}(i,j,k)\boldsymbol{v}(i)$$

# Tensor Vector Multiplication Cont.

- Multiplying a 3$^{rd}$ order tensor by two vectors produces a vector

$$z = \mathcal{T} \quad \bar{\times}_2 \quad v \quad \bar{\times}_1 \quad w$$

$$z(k) = \sum_i \left( \sum_j \mathcal{T}(i,j,k) v(i) \right) w(j) = \sum_{i,j} \mathcal{T}(i,j,k) v(i) w(j)$$

# Conditional Probability Table At Leaf is a Matrix

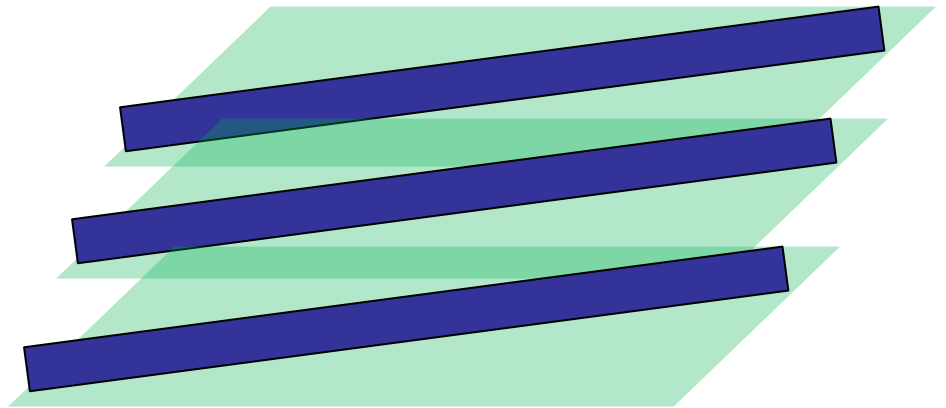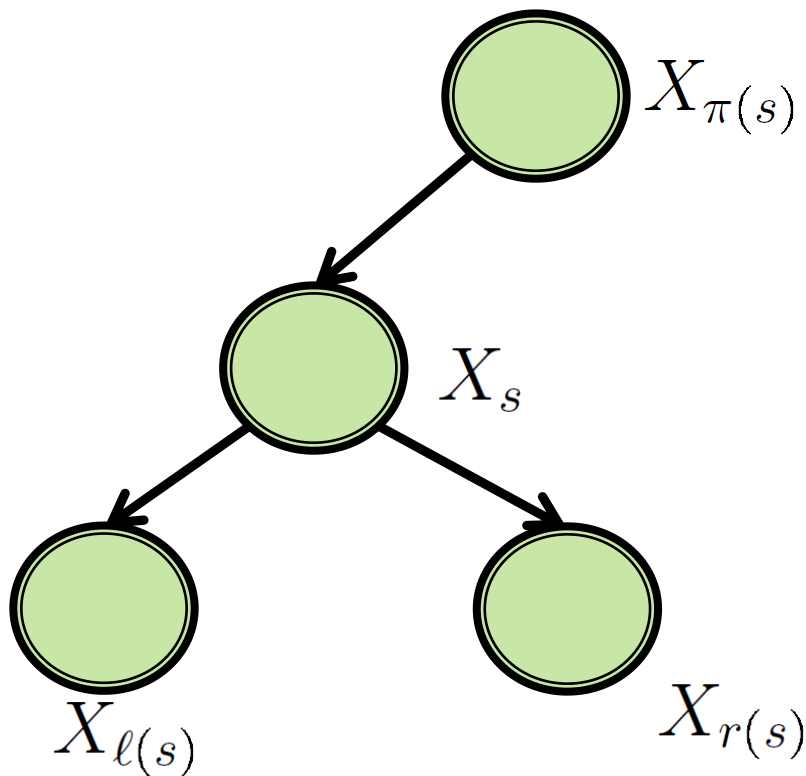$$\vec{\mathbb{P}}\left[X_s | X_{\pi(s)}\right]$$

$X_{\pi(s)}$

$X_s$

# CPT At Internal Node (Non-Root) is 3$^{\text{rd}}$ Order Tensor
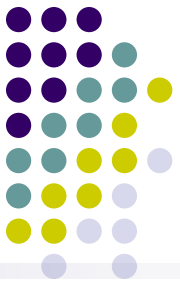
- Note that we have

$$\overrightarrow{\mathbb{P}}\left[\oslash X_s | X_{\pi(s)}\right] = \begin{cases} \mathbb{P}[X_s = i | X_{\pi(s)} = k] & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

$X_{\pi(s)}$
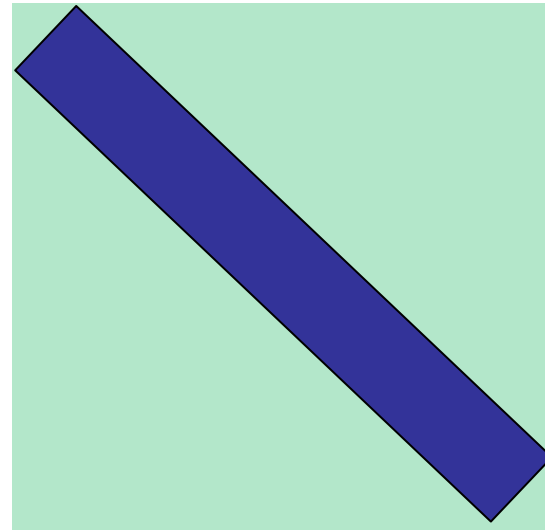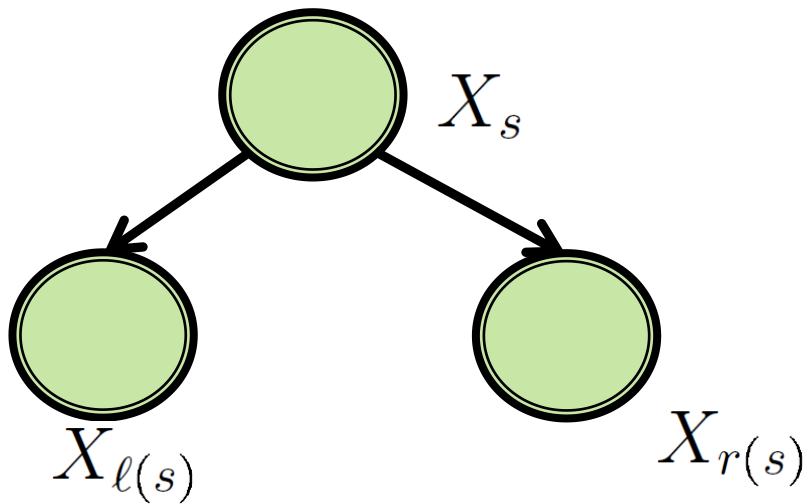
$X_s$

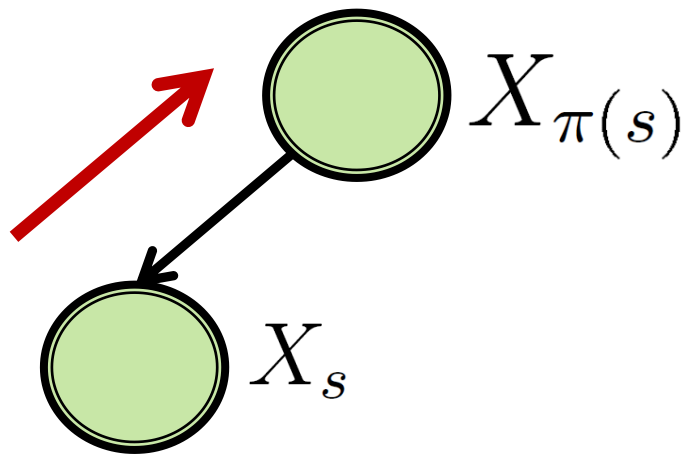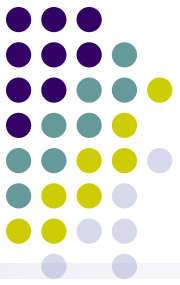$X_{\ell(s)}$

$X_{r(s)}$

# CPT At Root

- CPT at root is a matrix.

$$\overrightarrow{\mathbb{P}}\,[\oslash X_s] = \begin{cases} \mathbb{P}[X_s = i] & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

# The Outgoing Message as a Vector (at Leaf)

$$m_s = \boldsymbol{\delta}_{\bar{x}_s}^{\top} \overrightarrow{\mathbb{P}}[X_s | X_{\pi(s)}]$$

$X_{\pi(s)}$

$X_s$

**"bar" denotes evidence**

# The Outgoing Message At Internal Node

$$\boldsymbol{m}_s = \vec{\mathbb{P}}[\oslash X_s | X_{\pi(s)}] \; \bar{\times}_2 \; \boldsymbol{m}_{r(s)} \; \bar{\times}_1 \; \boldsymbol{m}_{\ell(s)}$$

$$\boldsymbol{m}_s(X_{\pi(s)} = k)$$
$$= \sum_{i,j} \mathbb{I}(i = j) \mathbb{P}[X_s = i | X_{\pi(s)} = k] \boldsymbol{m}_{x_{\ell(s)}}(X_s = i) \boldsymbol{m}_{x_{r(s)}}(X_s = j)$$

# At the Root

$$\mathbb{P}[\text{evidence}] = \boxed{\overrightarrow{\mathbb{P}}\left[\varnothing X_s\right]} \; \overline{\times}_2 \; \boxed{\boldsymbol{m}_{r(s)}} \; \overline{\times}_1 \; \boxed{\boldsymbol{m}_{\ell(s)}}$$

$$\mathbb{P}[\text{evidence}] = \sum_{i,j} \mathbb{I}(i = j)\mathbb{P}[X_s = i]\boldsymbol{m}_{x_{\ell(s)}}(X_s = i)\boldsymbol{m}_{x_{r(s)}}(X_s = j)$$

# Kernel Graphical Models [Song et al. 2010, Song et al. 2011]

- The Tensor CPTs at each node are replaced with RKHS functions/operators

**Leaf:** $\quad \overrightarrow{\mathbb{P}}\left[X_s | X_{\pi(s)}\right] \quad \Longrightarrow \quad \mathcal{C}_{s|\pi(s)}$

**Internal (non-root):** $\quad \overrightarrow{\mathbb{P}}\left[\oslash X_s | X_{\pi(s)}\right] \quad \Longrightarrow \quad \mathcal{C}_{ss|\pi(s)}$

**Root:** $\quad \overrightarrow{\mathbb{P}}\left[\oslash X_s\right] \quad \Longrightarrow \quad \mathcal{C}_{ss}$

# Conditional Embedding Operator for Internal Nodes
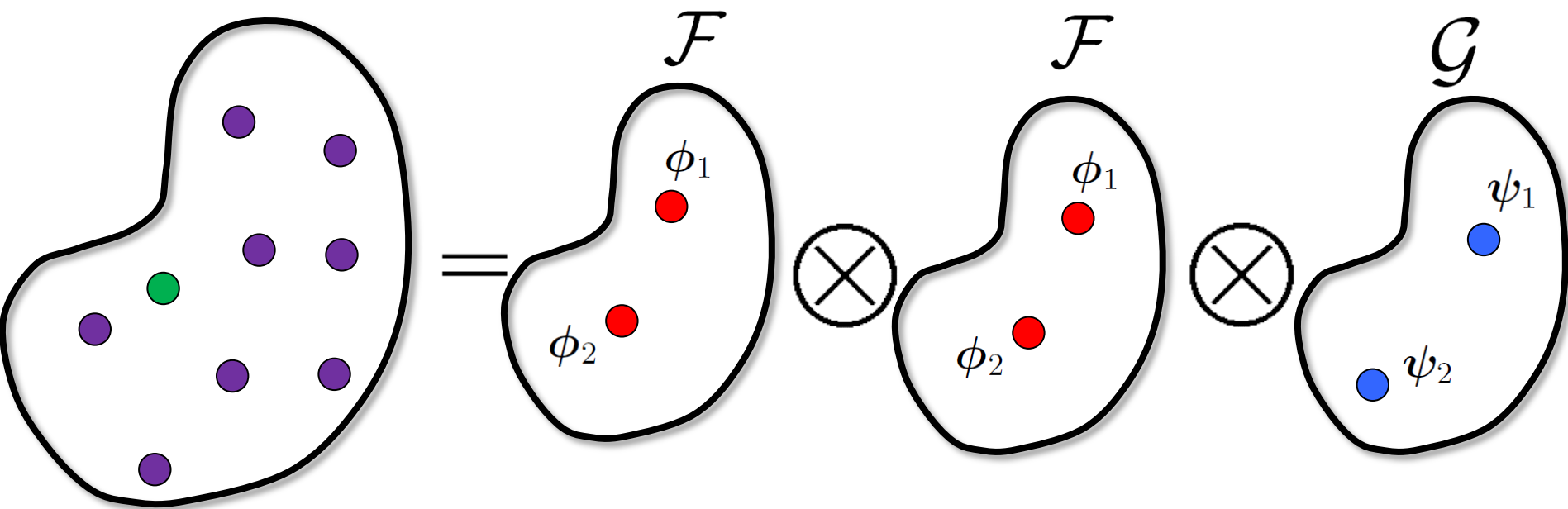
**What is** $\mathcal{C}_{ss|\pi(s)}$ **?**

$$\mathcal{C}_{XX|Y} = \mathcal{C}_{XXY}\mathcal{C}_{YY}^{-1}$$

**Embedding of** $\mathbb{P}[\oslash X_s | X_{\pi(s)}]$

# Embedding of Cross Covariance Operator in Different RKHS

$$\mathcal{C}_{XXY} = \mathbb{E}_{XY}[\boldsymbol{\phi}_X \otimes \boldsymbol{\phi}_X \otimes \boldsymbol{\psi}_Y]$$



**Embedding of** $\overrightarrow{\mathbb{P}}[\oslash X, Y]$