



**Lecture 11: Radial Basis Functions,
Gaussian Processes, Relevance Vector Machines,
Feature Selection, Dimensionality Reduction, LVQ
& Information-theoretic Co-clustering**

Prof. Krishna R. Pattipati
Dept. of Electrical and Computer Engineering
University of Connecticut
Contact: krishna@engr.uconn.edu (860) 486-2890

Fall 2018
November 26 and December 3, 2018



Lecture Outline

- Radial Basis Functions
- Gaussian Processes
- Relevant Vector Machines
- Feature Selection & Dimensionality Reduction
- Review of k-means and GMM-based Clustering
- Learning Vector Quantization
- Information-Theoretic Co-clustering
- Summary



What are Radial Basis Functions ?

- **Radial Basis Functions (RBFs)**

MLP

Activation function is determined by the inner product

$$y = \underline{w}^T \underline{x} \quad \text{and} \quad \hat{z} = g(y) = g(\underline{w}^T \underline{x})$$

RBF (Recall Kernel Idea of SVM and PNN)

Activation of hidden units is determined by the distance between the input vector and a prototype vector. Activation is highest at the center and radiates outward with intensity decreasing toward zero as the distance from the center

increases, $\phi(\|\underline{x} - \underline{x}^n\|)$



Characterizing RBF Networks

- Direct Connections between inputs and outputs also help a great deal
- *Hyperspheres* versus *Hyperplanes*
- Hidden units in RBF are *local* versus *global* in MLP
- Simple architecture
- Faster to train

Like Gaussian mixtures, except that **weights can be positive or negative**

$$y_k(\underline{x}) = \sum_{j=1}^M w_{kj} \phi_j(\underline{x}) + w_{k0} = \sum_{j=0}^M w_{kj} \phi_j(\underline{x})$$

or $\underline{y} = W \underline{\phi}(\underline{x})$ W is a C by $(M + 1)$ matrix

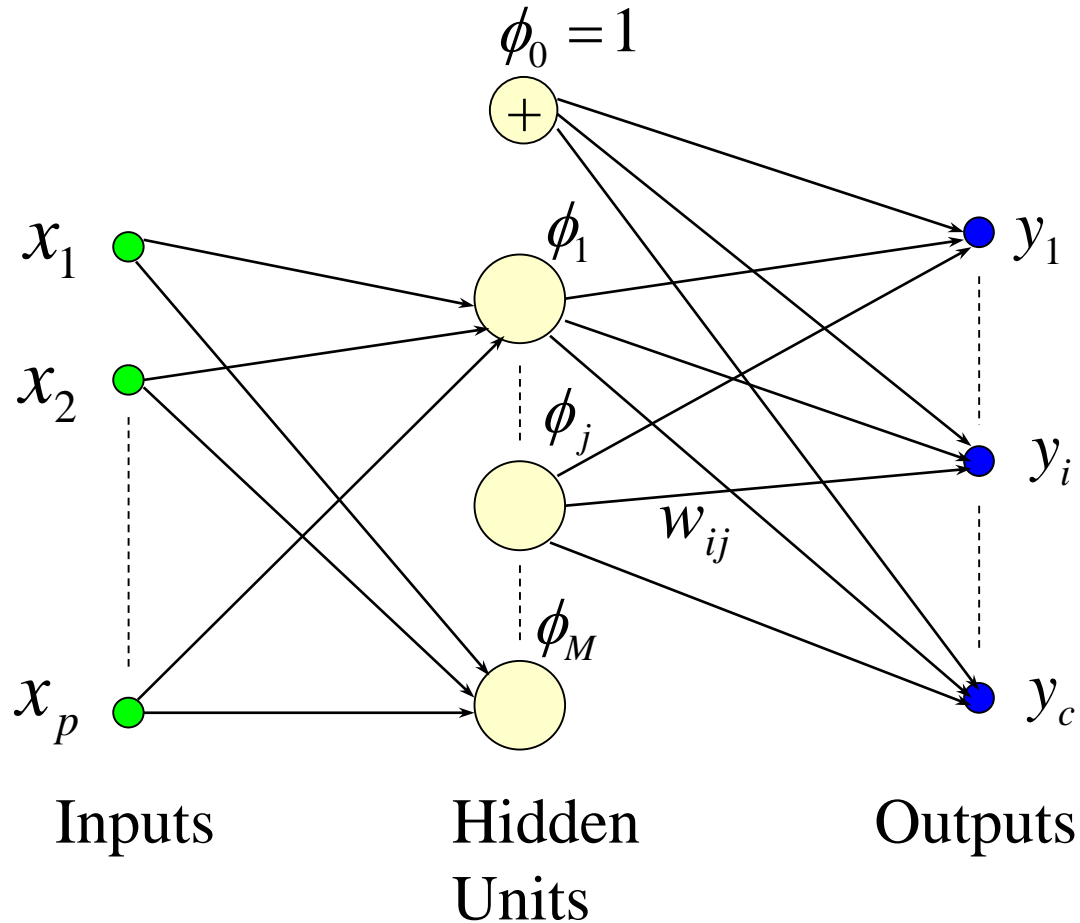
Typically, the basis functions $\phi_j(\underline{x})$, $j \geq 1$ are given by

$$\phi_j(\underline{x}) = \exp\left\{-\frac{\|\underline{x} - \underline{\mu}_j\|^2}{2\sigma_j^2}\right\}; \quad \varphi_j(\underline{x}) = \exp\left\{-\frac{1}{2}(\underline{x} - \underline{\mu}_j)^T \Sigma_j^{-1}(\underline{x} - \underline{\mu}_j)\right\}$$



Structure of an RBF Network

- In fact, RBF network is a three layer nonlinear regression system





Training RBF Weights - 1

- If know $\{\phi_j(\underline{x})\}$ and training data $\{\underline{x}^n, \underline{z}^n\}_{n=1}^N$, then the weights that minimize the squared error

$$E = \frac{1}{2} \sum_{n=1}^N \left\| \underline{y}(\underline{x}^n, W) - \underline{z}^n \right\|^2 + \underbrace{\frac{\lambda}{2} \text{tr}(W^T W)}_{\text{Regularization term}}$$

can be obtained as follows

$$\begin{aligned} E &= \frac{1}{2} \sum_{n=1}^N \left\| W \underline{\phi}(\underline{x}^n) - \underline{z}^n \right\|^2 + \frac{\lambda}{2} \text{tr}(W^T W) \\ E &= \frac{1}{2} \sum_{n=1}^N \text{tr} \left[\left(W \underline{\phi}(\underline{x}^n) - \underline{z}^n \right) \left(W \underline{\phi}(\underline{x}^n) - \underline{z}^n \right)^T \right] + \frac{\lambda}{2} \text{tr}(W^T W) \\ &= \frac{1}{2} \sum_{n=1}^N \text{tr} \left[W \underline{\phi}(\underline{x}^n) \underline{\phi}^T(\underline{x}^n) W^T - \underline{z}^n \underline{\phi}^T(\underline{x}^n) W^T - W \underline{\phi}(\underline{x}^n) \underline{z}^{nT} + \underline{z}^n \underline{z}^{nT} \right] \\ &\quad + \frac{\lambda}{2} \text{tr}(W^T W) \end{aligned}$$



Training RBF Weights - 2

$$\nabla_w E = W \sum_{n=1}^N \underline{\phi}(x^n) \underline{\phi}^T(x^n) - \sum_{n=1}^N \underline{z}^n \underline{\phi}^T(x^n) + \lambda W$$

or
$$\left[\sum_{n=1}^N \underline{\phi}(x^n) \underline{\phi}^T(x^n) \right] W^T + \lambda W^T = \sum_{n=1}^N \underline{\phi}(x^n) \underline{z}^{nT}$$

Let
$$\Phi =_N \begin{bmatrix} \underline{\phi}^T(x^1) \\ \underline{\phi}^T(x^2) \\ \vdots \\ \underline{\phi}^T(x^n) \end{bmatrix}; \quad Z =_N \begin{bmatrix} \underline{z}^{1T} \\ \underline{z}^{2T} \\ \vdots \\ \underline{z}^{NT} \end{bmatrix}$$

$$(\Phi^T \Phi + \lambda I_{M+1}) W^T = \Phi^T Z$$

$$\underline{y}(x) = W \underline{\phi}(x) = Z^T \underbrace{(\Phi \Phi^T + \lambda I_N)^{-1} \Phi \underline{\phi}(x)}_{\text{all inner products}} = D^T \underbrace{\Phi \underline{\phi}(x)}_{\text{kernel vector } \underline{g}(x)}$$

M+1 by M+1 M+1 by C M+1 by N N by C

$$W^T = (\Phi^T \Phi + \lambda I_{M+1})^{-1} \Phi^T Z = \Phi^T \underbrace{(\Phi \Phi^T + \lambda I_N)^{-1}}_{\text{Kernel Matrix } G + \lambda I_N} Z = \Phi^T D$$

Solution via: G-S, SVD, RLS, LMS, ... *D NxN matrix*



Optimal Weights of RBF - 1

Thus, we can obtain each column of D independently

$$y_k(\underline{x}) = \sum_{n=1}^N d_{n,k} g(\underline{x}, \underline{x}^n) = \underline{d}_k^T \underline{g}(\underline{x}); \underline{g}(\underline{x}) = \begin{bmatrix} g(\underline{x}, \underline{x}^1) \\ g(\underline{x}, \underline{x}^2) \\ \vdots \\ g(\underline{x}, \underline{x}^N) \end{bmatrix} = \begin{bmatrix} \underline{\phi}^T(\underline{x}^1) \underline{\phi}(\underline{x}) \\ \underline{\phi}^T(\underline{x}^2) \underline{\phi}(\underline{x}) \\ \vdots \\ \underline{\phi}^T(\underline{x}^N) \underline{\phi}(\underline{x}) \end{bmatrix}$$

$$\begin{bmatrix} g(\underline{x}^1, \underline{x}^1) + \lambda & g(\underline{x}^1, \underline{x}^2) & \cdots & g(\underline{x}^1, \underline{x}^N) \\ g(\underline{x}^2, \underline{x}^1) & g(\underline{x}^2, \underline{x}^2) + \lambda & \cdots & g(\underline{x}^2, \underline{x}^N) \\ \vdots & \vdots & \ddots & \vdots \\ g(\underline{x}^N, \underline{x}^1) & g(\underline{x}^N, \underline{x}^2) & \cdots & g(\underline{x}^N, \underline{x}^N) + \lambda \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_N \end{bmatrix} = \begin{bmatrix} z^1 \\ z^2 \\ \vdots \\ z^N \end{bmatrix}$$

$\underline{d}_k = (G + \lambda I_N)^{-1} \underline{z}_k = (\Phi \Phi^T + \lambda I_N)^{-1} \underline{z}_k; k = 1, 2, \dots, C$ A column of D Corresponding column of Z

$$y_k(\underline{x}) = \underline{d}_k^T \underline{g}(\underline{x}) = \underline{z}_k^T (G + \lambda I_N)^{-1} \underline{g}(\underline{x}) = \underline{g}^T(\underline{x}) (G + \lambda I)^{-1} \underline{z}_k$$

Kernel Ridge Regression



Selecting Basis Functions

- Key questions:
 - Why Gaussian basis functions?
 - How to choose M ?
 - How to choose the basis function parameters?
- **Why Gaussian basis functions?**

The motivation comes from regularization theory. Consider a single output problem for simplicity and consider the error function

$$J = \frac{1}{2} \sum_{n=1}^N (y(\underline{x}^n) - z^n)^2 + \frac{\lambda}{2} \int_{\underline{x}} |Py|^2 d\underline{x} \quad \dots\dots(1)$$

where P is some linear operator (e.g., differential operator) and λ is a regularization parameter. λ controls the smoothness of the fit.



Smoothing Functions

- Some examples of P

- $$\int_{\underline{x}} |Py|^2 d\underline{x} = \int_{\underline{x}} \sum_{i=1}^p \left(\frac{\partial^2 y}{\partial x_i^2} \right)^2 d\underline{x} \Rightarrow \text{reduce curvature}$$

- $$\int_{\underline{x}} |Py|^2 d\underline{x} = \int_{\underline{x}} \sum_{i=1}^p \sum_{m=0}^{\infty} a_m \left(\frac{\partial^m y}{\partial x_i^m} \right)^2 d\underline{x}$$

It can be shown that if $a_m = \frac{\sigma^{2m}}{m! 2^m}$, then $y(\underline{x}^n)$ is a sum of Gaussian RBFs.

- $$\int_{\underline{x}} |Py|^2 d\underline{x} = \int_{\underline{x}} \sum_{i=1}^p \left(\frac{\partial y}{\partial x_i} \right)^2 d\underline{x} \Rightarrow \text{to reduce sensitivity w.r.t. } x_i$$



Generalized Regression Network

- As in PNN, we can select $M = N$. In practice, take $M < N$ basis functions (e.g., via K-means, GMM,... clustering)
- Have different σ_i^2 for each unit.

$$\Rightarrow y(\underline{x}) = \sum_{n=1}^M w_n e^{-\|\underline{x} - \underline{\mu}_n\|^2 / 2\sigma_n^2} \Rightarrow \text{Generalized RBF}$$

- ***Relationships to GRNN and Regression***

- Suppose have training data $\{\underline{x}^n, \underline{z}^n\}$
- Joint density estimate via Gaussian kernels (recall PNN)

$$\hat{p}(\underline{x}, \underline{z}) = \frac{1}{N} \sum_{n=1}^N \frac{1}{(2\pi\sigma_x^2)^{p/2} (2\pi\sigma_z^2)^{c/2}} \exp \left\{ \frac{-\|\underline{x} - \underline{x}^n\|^2}{2\sigma_x^2} - \frac{\|\underline{z} - \underline{z}^n\|^2}{2\sigma_z^2} \right\}$$



Nadaraya – Watson Estimator

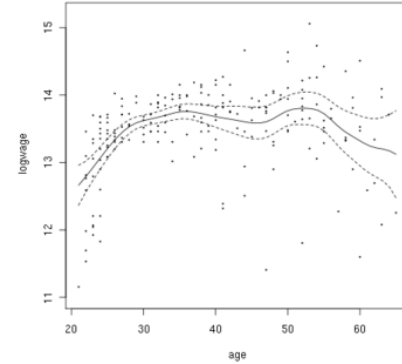
- **MMSE estimate = Conditional mean**

$$\underline{y}(\underline{x}) = E \{ \underline{z} / \underline{x} \} = \int \underline{z} p(\underline{z} / \underline{x}) d \underline{z}$$

$$= \frac{\int \underline{z} p(\underline{x}, \underline{z}) d \underline{z}}{\int p(\underline{x}, \underline{z}) d \underline{z}}$$

$$= \frac{\sum_{n=1}^N \underline{z}^n \exp \left\{ -\|\underline{x} - \underline{x}^n\|^2 / 2\sigma_x^2 \right\}}{\sum_{n=1}^N \exp \left\{ -\|\underline{x} - \underline{x}^n\|^2 / 2\sigma_x^2 \right\}} = \sum_{n=1}^N k(\underline{x}, \underline{x}^n) \underline{z}^n$$

$$k(\underline{x}, \underline{x}^n) = \frac{\exp \left\{ -\|\underline{x} - \underline{x}^n\|^2 / 2\sigma_x^2 \right\}}{\sum_{i=1}^N \exp \left\{ -\|\underline{x} - \underline{x}^i\|^2 / 2\sigma_x^2 \right\}}$$



Estimated Regression Function.

Nadaraya-Watson Estimator



RBF includes GMM

□ Gaussian Mixture Model

$$\hat{p}(\underline{x}, \underline{z}) = \sum_{j=1}^M P_j \frac{1}{(2\pi\sigma_x^2)^{p/2} \cdot (2\pi\sigma_z^2)^{c/2}} \exp \left\{ \frac{-\|\underline{x} - \underline{\mu}_j\|^2}{2\sigma_x^2} - \frac{\|\underline{z} - \underline{v}_j\|^2}{2\sigma_z^2} \right\}$$

$$y(\underline{x}) = \frac{\sum_{j=1}^M P_j \underline{v}_j \exp \left\{ \frac{-\|\underline{x} - \underline{\mu}_j\|^2}{2\sigma_x^2} \right\}}{\sum_{j=1}^M P_j \exp \left\{ \frac{-\|\underline{x} - \underline{\mu}_j\|^2}{2\sigma_x^2} \right\}}$$



RBF Networks for Classification

- RBF Networks for classification

Let
$$p(\underline{x} | z = k) = \sum_{j=1}^M p(\underline{x} | j)P(j | z = k)$$

$$p(z = k | \underline{x}) = \frac{p(\underline{x} | z = k)P(z = k)}{p(\underline{x})}$$

$$= \frac{\sum_{j=1}^M P(j | z = k) p(\underline{x} | j) \frac{P_j}{P_j} P(z = k)}{\sum_{l=1}^M p(\underline{x} | l) P_l}$$

$$= \sum_{j=1}^M \left[\frac{p(\underline{x} | j) P_j}{\sum_{l=1}^M p(\underline{x} | l) P_l} \right] \left[\frac{P(j | z = k) P(z = k)}{P_j} \right] = \sum_{j=1}^M \phi_j(\underline{x}) w_{kj}$$

Normalized RBF	$\phi_j(\underline{x}) = P(j \underline{x})$	$w_{kj} = P(z = k j)$
	RBF center given \underline{x}	Output class given RBF center



Probabilistic Interpretation of Normalized RBF

- Note that
 - $\phi_j(\underline{x}) = P(j | \underline{x}) \Rightarrow$ posterior probability of mixture component j given \underline{x}
 - $w_{kj} = P(z = k | j) \Rightarrow$ posterior probability of class k given mixture component j

Alternately,

$$\begin{aligned} P(z = k | \underline{x}) &= \sum_{j=1}^M P(z = k, j | \underline{x}) = \sum_{j=1}^M P(z = k | j, \underline{x}) P(j | \underline{x}) \\ &= \sum_{j=1}^M P(z = k | j) P(j | \underline{x}) = \sum_{j=1}^M w_{kj} \phi_j(\underline{x}) \end{aligned}$$



Optimizing RBF Parameters - 1

- **How to choose the basis function parameters ?**
 - Supervised training . . . Batch mode

$$y_k(\underline{x}) = \sum_{j=0}^M w_{kj} \phi_j(\underline{x})$$

where

$$\phi_j(\underline{x}) = \sum_{j=1}^M \exp \left\{ -\frac{\|\underline{x} - \underline{\mu}_j\|^2}{2\sigma_j^2} \right\} \quad j \geq 1$$

$$\phi_0(\underline{x}) = 1$$

Criterion:
$$\min_{\{W, \{\underline{\mu}_j\}, \{\sigma_j\}\}} J = \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^C \underbrace{\left[z_k^n - \sum_{j=0}^M w_{kj} \phi_j(\underline{x}^n) \right]^2}_{e_k^n}$$



Selecting Basis Function Parameters - 2

$$\frac{\partial J}{\partial w_{kj}} = -\sum_{n=1}^N e_k^n \varphi_j(\underline{x}^n) \quad \text{or, } \nabla_w J = -\underline{e}^n \underline{\varphi}^{nT} \quad \dots(1)$$

where

$$\underline{e}^n = \begin{bmatrix} e_1^n \\ e_2^n \\ \vdots \\ e_C^n \end{bmatrix}; \quad \underline{\varphi}^n = \begin{bmatrix} \varphi_0 \\ \varphi_1 \\ \vdots \\ \varphi_M \end{bmatrix}$$

$$\nabla_{\underline{\mu}_j} J = -\sum_{n=1}^N \sum_{k=1}^C e_k^n w_{kj} \varphi_j(\underline{x}^n) \frac{(\underline{x}^n - \underline{\mu}_j)}{\sigma_j^2} \quad \dots(2)$$

$$\nabla_{\sigma_j} J = -\sum_{n=1}^N \sum_{k=1}^C e_k^n w_{kj} \varphi_j(\underline{x}^n) \frac{\|\underline{x}^n - \underline{\mu}_j\|^2}{\sigma_j^3} \quad \dots(3)$$



Optimization Algorithms

- A variety of implementations are possible:
 - i) NLP techniques using (1), (2) and (3)
 - ii) EM or coordinate descent-type algorithms:
 - a) For a given $\{\underline{\mu}_j\}, \{\sigma_j\}$, find W via least squares
 - b) For a given W , find $\{\underline{\mu}_j\}, \{\sigma_j\}$ via NLP techniques or by setting gradients to zero for $\{\underline{\mu}_j\}$ and by grid search over $\{\sigma_j\}$
 - iii) Same as i) or ii) with M and initial $\{\underline{\mu}_j\}, \{\sigma_j\}$ via unsupervised techniques... recall K-means, GMM,...
 - iv) Extended Kalman Filter (treat parameters as states with identity transition matrix and process noise, and observations $(\underline{x}, \underline{z})$ as training data).



Gaussian Processes (GP) for Regression

Regression:

Let $y(\underline{x}) = \underline{w}^T \underline{\phi}(\underline{x}) = \underline{\phi}^T(\underline{x})\underline{w}$; $p(\underline{w}) = N(\underline{0}, \sigma_w^2 I) \Rightarrow p(y(\underline{x}) | \underline{x}) = N(0, \sigma_w^2 \underline{\phi}^T(\underline{x})\underline{\phi}(\underline{x}))$...note inner product

$$\text{Let } \underline{y}^N = [y(\underline{x}^1), y(\underline{x}^2), \dots, y(\underline{x}^N)]^T; \Phi_N = \begin{bmatrix} \underline{\phi}^T(\underline{x}^1) \\ \underline{\phi}^T(\underline{x}^2) \\ \vdots \\ \underline{\phi}^T(\underline{x}^N) \end{bmatrix}; \underline{X}^N = [\underline{x}^1, \underline{x}^2, \dots, \underline{x}^N]$$

$$\underline{y}^N = \Phi_N \underline{w} \Rightarrow p(\underline{y}^N | \underline{X}^N) = N(\underline{0}, \sigma_w^2 \Phi_N \Phi_N^T) = N(\underline{0}, \Sigma_{y^N})$$

$$\sigma_w^2 \Phi \Phi^T = K \Rightarrow k_{ij} = k(\underline{x}^i, \underline{x}^j) = \sigma_w^2 \underline{\phi}^T(\underline{x}^i) \underline{\phi}(\underline{x}^j) = \text{Cov}(y(\underline{x}^i), y(\underline{x}^j)) \Rightarrow \text{it is a Kernel}$$

$$\text{Typical } k(\underline{x}_i, \underline{x}_j) = \theta_0 \exp[-(\underline{x}^i - \underline{x}^j)^T \text{Diag}(\sigma_l^{-2})(\underline{x}^i - \underline{x}^j)] + \theta_1 + \theta_2 \underline{x}^i \underline{x}^j$$

Measured Targets: $z^n = y(\underline{x}^n) + v^n$; $p(v^n) = N(0, \sigma_v^2)$; $\underline{z}^N = [z^1, z^2, \dots, z^N]^T$

$$\underline{z}^N = \underline{y}^N + \underline{v}^N \Rightarrow p(\underline{z}^N | \underline{y}^N) = N(\underline{y}^N, \sigma_v^2 I_N)$$

$$\Rightarrow p(\underline{z}^N | \underline{X}^N) = \int_{\underline{y}^N} p(\underline{z}^N | \underline{y}^N) p(\underline{y}^N | \underline{X}^N) d\underline{y}^N = N(\underline{0}, \Sigma_y^N + \sigma_v^2 I_N) = N(\underline{0}, \Sigma_z^N)$$

$$\underline{z}^N = \underline{y}^N + \underline{v}^N = \Phi^N \underline{w} + \underline{v}^N$$

$$\text{Note: } p(\underline{z}^{N+1} | \underline{X}^{N+1}) = N(\Sigma_y^{N+1} + \sigma_v^2 I_{N+1}) = N(\underline{0}, \Sigma_z^{N+1}) = \begin{bmatrix} \overbrace{\Sigma_y^N + \sigma_v^2 I_N}^{\Sigma_z^N} & \overbrace{\sigma_w^2 \Phi_N \underline{\phi}(\underline{x}^{N+1})}^k \\ \underbrace{\sigma_w^2 \underline{\phi}^T(\underline{x}^{N+1}) \Phi_N^T}_{k^T} & \sigma_v^2 + \sigma_w^2 \underline{\phi}^T(\underline{x}^{N+1}) \underline{\phi}(\underline{x}^{N+1}) \end{bmatrix}$$



GP-based Prediction

Prediction \Rightarrow Conditional mean and variance

$$p(\underline{z}^{N+1} | \underline{z}^N, \underline{X}^{N+1}) = \frac{p(\underline{z}^{N+1}, \underline{X}^{N+1})}{p(\underline{z}^N, \underline{X}^{N+1})} = \frac{p(\underline{z}^{N+1} | \underline{X}^{N+1})}{\int_{\underline{z}^{N+1}} p(\underline{z}^{N+1} | \underline{X}^{N+1}) d\underline{z}^{N+1}}$$

$$p(\underline{z}^{N+1} | \underline{z}^N, \underline{X}^{N+1}) = N(\underbrace{\underline{k}^T (\Sigma_z^N)^{-1} \underline{z}^N}_{\underline{\alpha}}, \underbrace{\sigma_v^2 + \sigma_w^2 \underbrace{\underline{\phi}^T(\underline{x}^{N+1}) \Phi_N^T (\Sigma_z^N)^{-1} \underline{\phi}(\underline{x}^{N+1}) - \underline{k}^T (\Sigma_z^N)^{-1} \underline{k}}_{\text{Var}(\underline{z}^{N+1} | \underline{z}^N, \underline{X}^{N+1})})$$

Note: $E\{\underline{z}^{N+1} | \underline{z}^N, \underline{X}^{N+1}\} = \underbrace{\sigma_w^2 \underline{\phi}^T(\underline{x}^{N+1}) \Phi_N^T (\Sigma_z^N)^{-1} \underline{z}^N}_{\underline{k}^T(\underline{x}^N, \underline{x}^{N+1})} = \underline{k}^T(\underline{x}^N, \underline{x}^{N+1}) \underline{\alpha}$

where $\underline{\alpha} = (\Sigma_z^N)^{-1} \underline{z}^N$

$$\underline{k}^T(\underline{x}^N, \underline{x}^{N+1}) = \sigma_w^2 [\underbrace{\underline{\phi}^T(\underline{x}^{N+1}) \underline{\phi}(\underline{x}^1), \underline{\phi}^T(\underline{x}^{N+1}) \underline{\phi}(\underline{x}^2), \dots, \underline{\phi}^T(\underline{x}^{N+1}) \underline{\phi}(\underline{x}^N)}_{\text{inner products}}]$$

$$\underline{\phi}^T(\underline{x}^{N+1}) \underline{\phi}(\underline{x}^j) = \sigma_f^2 \exp\left(-\frac{1}{2} (\underline{x}^{N+1} - \underline{x}^j)^T \text{diag}(l_i^{-2}) (\underline{x}^{N+1} - \underline{x}^j)\right) + \sigma_y^2 I_{p+1}$$

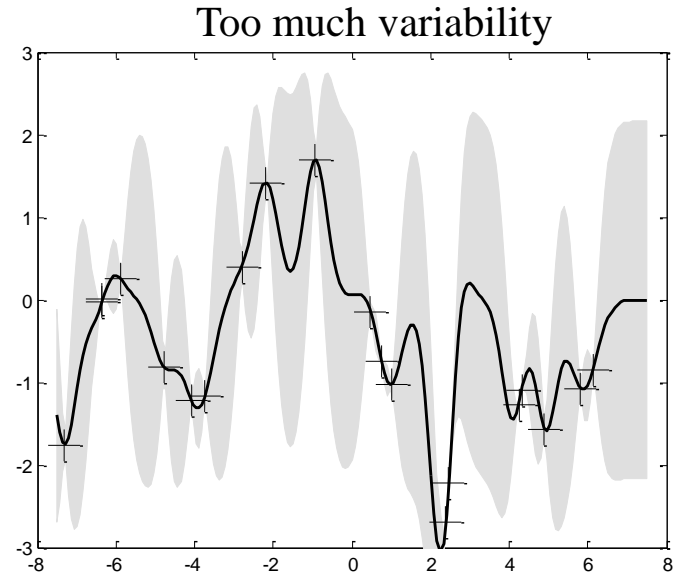
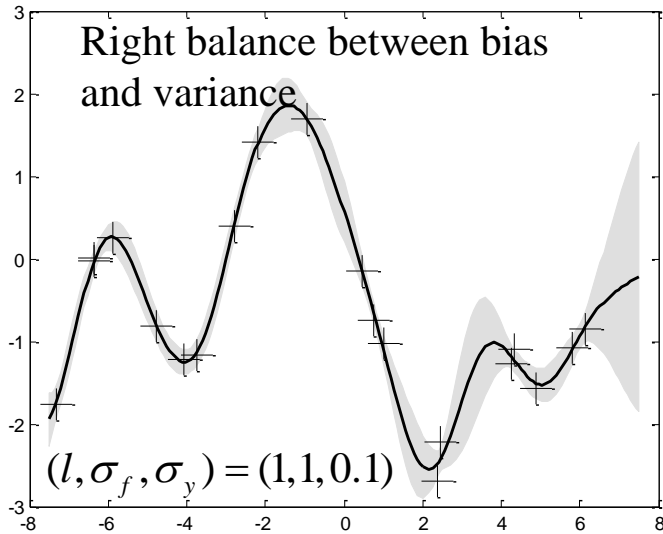
l_i = horizontal scale over which the function changes along coordinate i

σ_f^2 = controls the vertical scale of the function

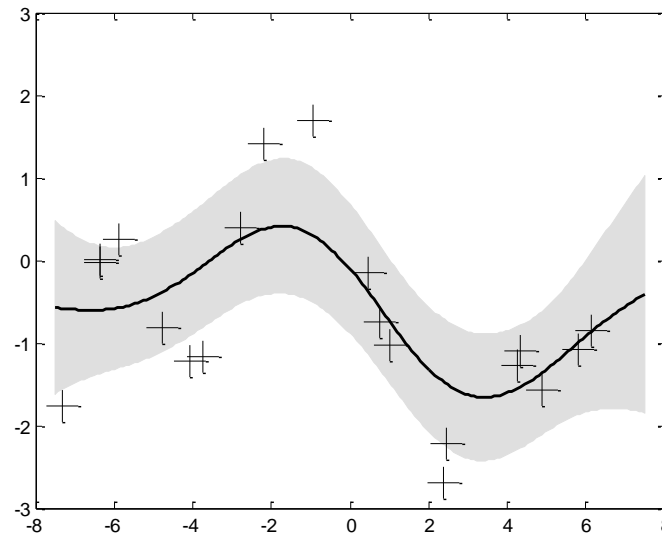
σ_y^2 = noise variance



GP-based Prediction



Too much bias





GP-based Learning

Learning Kernel Parameters $\{l_i, \sigma_f^2, \sigma_y^2\}$ via ML:

$$\ln p(\underline{z}^N, \underline{X}^N | \underline{\theta}) = -\frac{N}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_z^N| - \frac{1}{2} \underline{z}^{N^T} (\Sigma_z^N)^{-1} \underline{z}^N$$

$$\begin{aligned} \frac{\partial}{\partial \theta_j} \ln p(\underline{z}^N, \underline{X}^N | \underline{\theta}) &= -\frac{1}{2} \text{tr}[(\Sigma_z^N)^{-1} \frac{\partial \Sigma_z^N}{\partial \theta_j}] + \frac{1}{2} \underline{z}^{N^T} (\Sigma_z^N)^{-1} \frac{\partial \Sigma_z^N}{\partial \theta_j} (\Sigma_z^N)^{-1} \underline{z}^N \\ &= \frac{1}{2} \text{tr}[(\underline{\alpha} \underline{\alpha}^T - (\Sigma_z^N)^{-1}) \frac{\partial \Sigma_z^N}{\partial \theta_j}] \end{aligned}$$

where $\underline{\alpha} = (\Sigma_z^N)^{-1} \underline{z}^N$

Rank ordering features:

Large scale $l_i \Rightarrow$ feature i is not relevant



GP Implementation via Kalman Filter

$$\underline{w}^{n+1} = \underline{w}^n; \underline{w}^0 \sim N(\underline{0}, \sigma_w^2 I)$$

$$z^n = \underline{\phi}^T(\underline{x}^n) \underline{w}^n + v^n; v^n \sim N(0, \sigma_v^2)$$

Measurement Update :

$$\begin{aligned} \hat{\underline{w}}^{N/N} &= \Sigma^{N|N} \left(\Sigma^{N-1|N-1} \right)^{-1} \hat{\underline{w}}^{N-1/N-1} + \frac{\Sigma^{N|N} \underline{\phi}(\underline{x}^N)}{\sigma_v^2} z^N \\ &= \left(I + \Sigma^{N-1|N-1} \frac{\underline{\phi}(\underline{x}^N) \underline{\phi}^T(\underline{x}^N)}{\sigma_v^2} \right)^{-1} \hat{\underline{w}}^{N-1/N-1} + \frac{\Sigma^{N|N} \underline{\phi}(\underline{x}^N)}{\sigma_v^2} z^N \\ &= \hat{\underline{w}}^{N-1/N-1} + K_n \left(z^N - \underline{\phi}^T(\underline{x}^N) \hat{\underline{w}}^{N-1/N-1} \right) \end{aligned}$$

where $K_n = \frac{\Sigma^{N-1|N-1} \underline{\phi}(\underline{x}^N)}{\sigma_v^2 + \underline{\phi}^T(\underline{x}^N) \Sigma^{N-1|N-1} \underline{\phi}(\underline{x}^N)} = \frac{\Sigma^{N|N} \underline{\phi}(\underline{x}^N)}{\sigma_v^2} = \mathbf{Kalman\ Gain}$

So, $E\{z^{N+1} | \underline{z}^N, \underline{X}^{N+1}\} = \underline{\phi}^T(\underline{x}^{N+1}) \hat{\underline{w}}^{N/N}$

$Var\{z^{N+1} | \underline{z}^N, \underline{X}^{N+1}\} = \underline{\phi}^T(\underline{x}^{N+1}) \Sigma^{N|N} \underline{\phi}(\underline{x}^{N+1}) + \sigma_v^2 = \mathbf{Innovation\ variance}$

Note: This is recursive, but lose Kernel advantage (requires $\phi(\underline{x})$)



GP for Classification

Classification: Remember continuous output goes through a logistic

$$p(z | y) = \sigma(y)^z [1 - \sigma(y)]^{1-z}; \sigma(y) = [1 + \exp(-y)]^{-1} = \hat{z}; y(\underline{x}) = \underline{w}^T \underline{\phi}(\underline{x})$$

$$p(\underline{y}^N | \underline{X}^N) = N(\underline{0}, \sigma_w^2 \Phi \Phi^T + \varepsilon I_N) = N(\underline{0}, \Sigma_y^N); \varepsilon \text{ for numerical stability}$$

$$p(z^{N+1} = 1 | \underline{z}^N, \underline{X}^{N+1}) = \int_{y(\underline{x}^{N+1})} p(z^{N+1} = 1 | y(\underline{x}^{N+1})) p(y(\underline{x}^{N+1}) | \underline{z}^N, \underline{X}^{N+1}) dy(\underline{x}^{N+1})$$

$$= \int_{y(\underline{x}^{N+1})} \sigma(y(\underline{x}^{N+1})) p(y(\underline{x}^{N+1}) | \underline{z}^N, \underline{X}^{N+1}) dy(\underline{x}^{N+1})$$

⇒ **integration of logistic**

$$\begin{aligned} \underline{w}^{n+1} &= \underline{w}^n; \underline{w}^0 \sim N(\underline{0}, \sigma_w^2 I) \\ z^n &= \text{sgn} \left[\underbrace{\underline{\phi}^T(\underline{x}^n) \underline{w}^n}_{y^n} + v^n \right]; v^n \sim N(0, \varepsilon) \end{aligned}$$

Laplace Approximation: Assume $p(y(\underline{x}^{N+1}) | \underline{z}^N, \underline{X}^{N+1})$ is Gaussian

$$p(y(\underline{x}^{N+1}) | \underline{z}^N, \underline{X}^{N+1}) = \int_{\underline{y}^N} p(y(\underline{x}^{N+1}) | \underline{y}^N, \underline{x}^{N+1}) p(\underline{y}^N | \underline{z}^N, \underline{X}^N) d\underline{y}^N$$

$$p(y(\underline{x}^{N+1}) | \underline{y}^N) = N(\underline{k}^T (\Sigma_y^N)^{-1} \underline{y}^N, \sigma_w^2 \underline{\phi}^T(\underline{x}^{N+1}) \underline{\phi}(\underline{x}^{N+1}) + \varepsilon - \underline{k}^T (\Sigma_y^N)^{-1} \underline{k})$$

$$\underline{k}^T = [k(\underline{x}^1, \underline{x}^{N+1}), k(\underline{x}^2, \underline{x}^{N+1}), \dots, k(\underline{x}^N, \underline{x}^{N+1})]$$

$$p(\underline{y}^N | \underline{z}^N, \underline{X}^N) \propto p(\underline{z}^N | \underline{y}^N) p(\underline{y}^N | \underline{X}^N) = \prod_{n=1}^N [1 - \sigma(y(\underline{x}^n))] \left[\frac{\sigma(y(\underline{x}^n))}{1 - \sigma(y(\underline{x}^n))} \right]^{z^n} N(\underline{0}, \Sigma_y^N)$$



GP and Probit Approximation

Use the mode of $\ln p(\underline{y}^N | \underline{z}^N, \underline{X}^N)$ as mean and $(-\text{Hessian})^{-1}$ as covariance :

$$\ln p(\underline{y}^N | \underline{z}^N, \underline{X}^N) \propto -\sum_{n=1}^N \ln(1 + e^{y(x^n)}) + (\underline{z}^N)^T \underline{y}^N - \frac{1}{2} \ln |\Sigma_y^N| - \frac{1}{2} \underline{y}^{N^T} (\Sigma_y^N)^{-1} \underline{y}^N$$

$$\nabla_{\underline{y}^N} \ln p(\underline{y}^N | \underline{z}^N, \underline{X}^N) = -\underline{\sigma}(\underline{y}^N) + \underline{z}^N - (\Sigma_y^N)^{-1} \underline{y}^N$$

where $\underline{\sigma}(\underline{y}^N) = [\sigma(y(x^1)), \sigma(y(x^2)), \dots, \sigma(y(x^N))]^T$

$$\nabla_{\underline{y}^N}^2 \ln p(\underline{y}^N | \underline{z}^N, \underline{X}^{N+1}) = -\text{Diag}[\sigma(y(x^n))(1 - \sigma(y(x^n)))] - (\Sigma_y^N)^{-1} = -D_N - (\Sigma_y^N)^{-1}$$

$$p(\underline{y}^N | \underline{z}^N, \underline{X}^{N+1}) \approx N(\underline{y}^N; \hat{\underline{y}}^N, [D_N + (\Sigma_y^N)^{-1}]^{-1}); \hat{\underline{y}}^N = \Sigma_y^N [\underline{z}^N - \underline{\sigma}(\hat{\underline{y}}^N)]$$

so, $p(y(x^{N+1}) | \underline{z}^N, \underline{X}^{N+1}) \approx N(y(x^{N+1}); \hat{y}(x^{N+1}), \Sigma_{\hat{y}})$

$$\hat{y}(x^{N+1}) = \underline{k}^T (\Sigma_y^N)^{-1} \Sigma_y^N (\underline{z}^N - \underline{\sigma}(\hat{\underline{y}}^N)) = \underline{k}^T (\underline{z}^N - \underline{\sigma}(\hat{\underline{y}}^N))$$

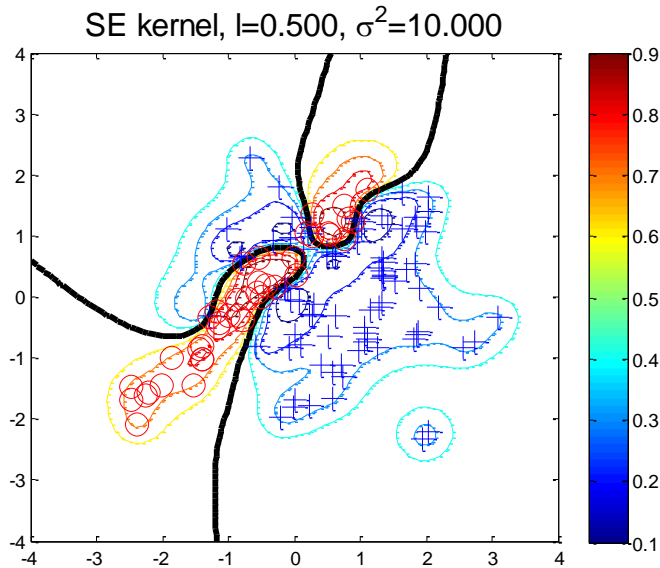
$$\Sigma_{\hat{y}} = \sigma_w^2 \underline{\phi}^T(x^{N+1}) \underline{\phi}(x^{N+1}) + \varepsilon - \underline{k}^T [D_N + (\Sigma_y^N)^{-1}]^{-1} \underline{k}$$

$$\text{Finally, } p(z^{N+1} = 1 | \underline{z}^N, \underline{X}^{N+1}) = \int_a \sigma(a) N(a | \hat{y}(x^{N+1}), \Sigma_{\hat{y}}) da = \Phi\left(\frac{\hat{y}(x^{N+1})}{\sqrt{1 + \pi \Sigma_{\hat{y}}/8}}\right)$$

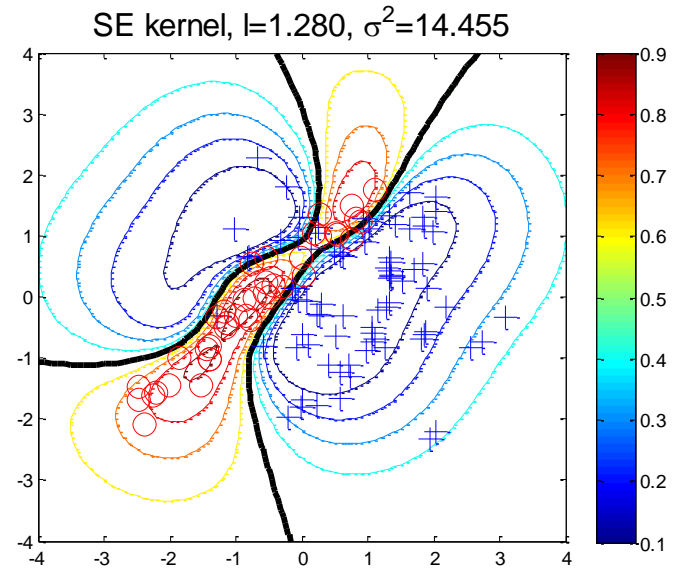
Approximate $\sigma(a) \approx \Phi\left(\sqrt{\frac{\pi}{8}} a\right)$ by requiring logistic and probit to have same slope at $a = 0$.



GP for Classification



$$(l, \sigma_f, \sigma_y) = (0.3, 1, 3.16)$$



$$(l, \sigma_f, \sigma_y) = (1.280, 1, 3.8)$$

SE: Squared Exponential or Exponentiated Quadratic



Relevance Vector Machine

Classification : Remember continuous output goes through a logistic

$$p(z | y) = \sigma(y)^z [1 - \sigma(y)]^{1-z}; \sigma(y) = [1 + \exp(-y)]^{-1} = \hat{z}; y(\underline{x}) = \underline{w}^T \underline{\phi}(\underline{x})$$

$$\ln p(\underline{w} | \underline{z}^N, \underline{X}^N) = \ln p(\underline{z}^N | \underline{w}, \underline{X}^N) + \ln p(\underline{w}) + \text{const.}$$

$$p(\underline{w}) = N(\underline{0}, \Sigma_w); \Sigma_w = \text{Diag}(\sigma_w^2) \dots \text{need to be learned}$$

$$\begin{aligned} \nabla_{\underline{w}} \left[z^n \ln \sigma(y^n) + (1 - z^n) \ln(1 - \sigma(y^n)) \right] \\ = \left(z^n [1 - \sigma(y^n)] - (1 - z^n) \sigma(y^n) \right) \underline{\phi}(\underline{x}^n) \\ = \left(z^n - \sigma(y^n) \right) \underline{\phi}(\underline{x}^n) \end{aligned}$$

$$p(\underline{z}^N | \underline{w}, \underline{X}^N) = \prod_{n=1}^N \sigma(y^n)^{z^n} [1 - \sigma(y^n)]^{1-z^n} = \prod_{n=1}^N [\hat{z}^n]^{z^n} [1 - \hat{z}^n]^{1-z^n}$$

$$\ln p(\underline{w} | \underline{z}^N, \underline{X}^N) = \sum_{n=1}^N \underbrace{z^n \ln \sigma(y^n) + (1 - z^n) \ln(1 - \sigma(y^n))}_{\text{negative cross entropy}} - \frac{1}{2} \underline{w}^T \Sigma_w^{-1} \underline{w} + \text{const.}$$

$$\nabla_{\underline{w}} \ln p(\underline{w} | \underline{z}^N, \underline{X}^N) = \Phi^T (\underline{z}^N - \hat{\underline{z}}^N) - \Sigma_w^{-1} \underline{w} \Rightarrow \underline{w}^* = \Sigma_w \Phi^T (\underline{z}^N - \hat{\underline{z}}^{*N})$$

$$\nabla_{\underline{w}}^2 \ln p(\underline{w} | \underline{z}^N, \underline{X}^N) = -(\Phi^T D_N \Phi + \Sigma_w^{-1}) \Rightarrow \Sigma_w^* = (\Phi^T D_N^* \Phi + \Sigma_w^{-1})^{-1}; D_N = \text{diag}[\hat{z}^n (1 - \hat{z}^n)]$$



Learning RVM Classifiers

$$\nabla_{\underline{w}} \ln p(\underline{w} | \underline{z}^N, \underline{X}^N) = \Phi^T (\underline{z}^N - \hat{\underline{z}}^N) - \Sigma_w^{-1} \underline{w} \Rightarrow \underline{w}^* = \Sigma_w \Phi^T (\underline{z}^N - \hat{\underline{z}}^N)$$

$$\nabla_{\underline{w}}^2 \ln p(\underline{w} | \underline{z}^N, \underline{X}^N) = -(\Phi^T D_N \Phi + \Sigma_w^{-1}) \Rightarrow \Sigma_w^* = (\Phi^T D_N^* \Phi + \Sigma_w^{-1})^{-1}; D_N = \text{diag}[\hat{z}^n(1 - \hat{z}^n)]$$

EM or IRLS – like Algorithm between steps 1 and 2:

1. For given \underline{w} and Σ_w , obtain \underline{w}^* and Σ_w^*

2. $p(\underline{z}^N | \underline{X}^N) = \int_{\underline{w}} p(\underline{z}^N | \underline{w}, \underline{X}^N) p(\underline{w} | \underline{X}^N) d\underline{w} \approx (2\pi)^{M/2} |\Sigma_w^*|^{1/2} p(\underline{z}^N | \underline{w}^*, \underline{X}^N) p(\underline{w}^*) \dots \text{Laplace}$

$$\frac{\partial \ln p(\underline{z}^N | \underline{X}^N)}{\partial \sigma_{w_i}} = \frac{\partial}{\partial \sigma_{w_i}} \left[\frac{1}{2} \ln |\Sigma_w^*| - \frac{w_i^{*2}}{2\sigma_{w_i}^2} - \ln \sigma_{w_i} \right] = 0$$

$$\Rightarrow \frac{1}{\sigma_{w_i}^3} [(\Phi^T D_N^* \Phi + \Sigma_w^{-1})^{-1}]_{ii} + \frac{w_i^{*2}}{\sigma_{w_i}^3} - \frac{1}{\sigma_{w_i}} = 0 \Rightarrow \sigma_{w_i}^2 = w_i^{*2} + [\Sigma_w^*]_{ii}$$

Logistic regression with L_2 , L_1 criteria, RVM and SVM have similar performance



RVM for Regression

Estimate : $y(\underline{x}) = \underline{w}^T \underline{\phi}(\underline{x})$

Likelihood : $p(z | \underline{x}, \underline{w}, \beta) = N(z; y(\underline{x}), \beta^{-1})$

Prior : $p(\underline{w} | \Sigma_w) = N(\underline{0}, \Sigma_w); \Sigma_w = \text{Diag}(\sigma_{w_i}^2)$

$$\underline{w}^{n+1} = \underline{w}^n = \underline{w}; \underline{w} \sim N(\underline{0}, \text{Diag}(\sigma_{w_i}^2))$$

$$\underline{z}^n = \underline{\phi}^T(\underline{x}^n) \underline{w}^n + v^n; v^n \sim N(0, \beta^{-1})$$

Like estimation with unknown noise variances

$$p(\underline{w} | \underline{z}^N, \underline{X}^N, \beta, \Sigma_w) \approx N(\underline{w}, \underline{w}_{MAP}, H^{-1}); \underline{w}_{MAP} = \beta H^{-1} \Phi^T \underline{z}^N; H = (\Sigma_w)^{-1} + \beta \Phi^T \Phi; \Phi = \begin{bmatrix} \underline{\phi}^T(\underline{x}^1) \\ \underline{\phi}^T(\underline{x}^2) \\ \vdots \\ \underline{\phi}^T(\underline{x}^N) \end{bmatrix}; N \text{ by } M \text{ vector}$$

$$p(\underline{z}^N | \underline{X}^N, \beta, \Sigma_w) = \int_{\underline{w}} p(\underline{z}^N | \underline{X}^N, \beta, \underline{w}) p(\underline{w} | \Sigma_w) d\underline{w} = N(\underline{z}^N; \underline{0}, \Sigma_z^N); \Sigma_z^N = \beta^{-1} I_N + \Phi \Sigma_w \Phi^T$$

$$NLL = -\ln p(\underline{z}^N | \underline{X}^N, \beta, \Sigma_w) = \frac{1}{2} \left\{ N \ln 2\pi + \ln |\Sigma_z^N| + (\underline{z}^N)^T (\Sigma_z^N)^{-1} \underline{z}^N \right\}$$

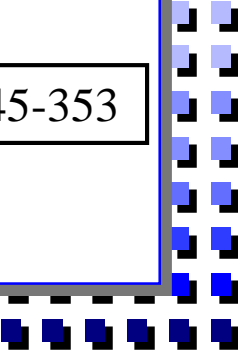
$$= \frac{1}{2} \left\{ N \ln 2\pi - N \ln \beta + \beta (\underline{z}^N - \Phi \underline{w}_{MAP})^T (\underline{z}^N - \Phi \underline{w}_{MAP}) + \ln |H| + \underline{w}_{MAP}^T \Sigma_w^{-1} \underline{w}_{MAP} \right\}$$

$$\underline{z}^N = \Phi \underline{w} + v^n; v^n \sim N(0, \beta^{-1})$$

$$(\sigma_{w_i}^2)^{new} = \frac{w_{MAP,i}^2 \sigma_{w_i}^2}{\sigma_{w_i}^2 - (H^{-1})_{ii}} \Rightarrow \text{feature that does not reduce posterior variance will be removed.}$$

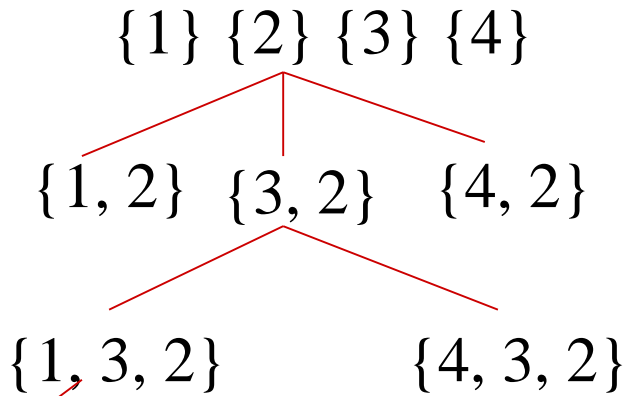
$$(\beta^{new})^{-1} = \frac{(\underline{z}^N - \Phi \underline{w}_{MAP})^T (\underline{z}^N - \Phi \underline{w}_{MAP}) \sigma_{w_i}^2}{(N - M) \sigma_{w_i}^2 + \sum_{n=1}^N (H^{-1})_{ii}}$$

See Bishop's book, pp. 168-170 and 345-353



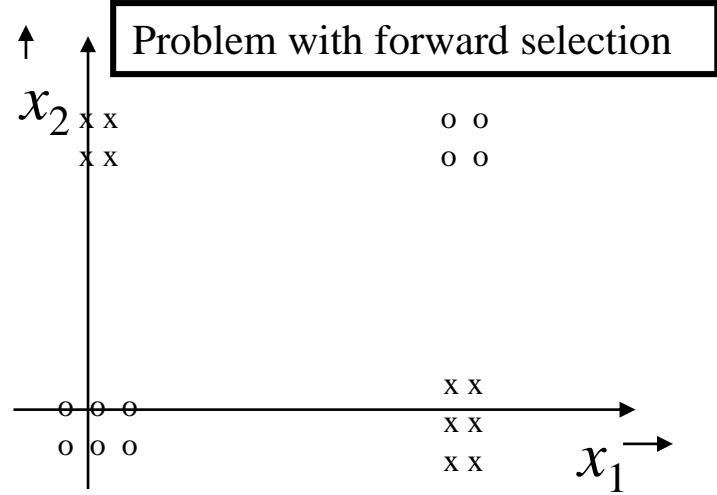


Forward and Backward Selection



{4, 1, 3, 2}

$$\frac{p(p+1)}{2} \text{ evaluations}$$



Individually, x_1 and x_2 do not provide discrimination but together they do.

Another strategy is sequential backward elimination. In the strategy, one feature is deleted from the set, chosen from among the candidates as the one which gives the *smallest* decrement ΔS_k^2 . This overcomes the problem of sequential forward selection methods, but at a higher computational cost.



Criterion for Binary Features

Clearly, one can devise hybrid strategies. At a stage k , add l features via forward selection and delete r features using the backward algorithm.

ΔS_k^2 when features are binary: $\Delta S_k^2 = \frac{\|\text{cov}(x_k, \underline{z})\|^2}{\text{var}(x_k)}$

$$\text{prob}\{x_k | z = i\} = \begin{cases} 1 - \mu_{ik} & \text{for } x_k = 0 \\ \mu_{ik} & \text{for } x_k = 1 \end{cases} \quad \text{prob}\{x_k\} = \begin{cases} 1 - \mu_k & \text{for } x_k = 0 \\ \mu_k & \text{for } x_k = 1 \end{cases};$$

$$\text{cov}(x_k, \underline{z}) = \begin{bmatrix} \pi_1 E(x_k | 1) \\ \pi_2 E(x_k | 2) \\ \cdot \\ \cdot \\ \pi_C E(x_k | K) \end{bmatrix} - E(x_k) \begin{bmatrix} \pi_1 \\ \pi_2 \\ \cdot \\ \cdot \\ \pi_C \end{bmatrix} = \begin{bmatrix} \pi_1(\mu_{1k} - \mu_k) \\ \pi_2(\mu_{2k} - \mu_k) \\ \cdot \\ \cdot \\ \pi_C(\mu_{Ck} - \mu_k) \end{bmatrix} \Rightarrow \Delta S_k^2 = \frac{\sum_{l=1}^C \pi_l^2 (\mu_{lk} - \mu_k)^2}{\mu_k(1 - \mu_k)}$$

$$\mu_k = P(x_k = 1) = \sum_{j=1}^C P(x_k = 1 | z = j)P(z = j) = \sum_{j=1}^C \pi_j \mu_{jk} \leftarrow$$



Orthogonal Matching Pursuit (OMP)

OMP: Greedy Forward Selection

$$D = \{ \{ \underline{x}^1, z^1 \}, \{ \underline{x}^2, z^2 \}, \dots, \{ \underline{x}^N, z^N \} \}$$

$$X = \begin{bmatrix} (\underline{x}^1)^T \\ (\underline{x}^2)^T \\ \vdots \\ (\underline{x}^N)^T \end{bmatrix} = \begin{bmatrix} \tilde{\underline{x}}_1 & \tilde{\underline{x}}_2 & \dots & \tilde{\underline{x}}_p \end{bmatrix}; N \times p; \underline{z} = \begin{bmatrix} z^1 \\ z^2 \\ \vdots \\ z^N \end{bmatrix}; N \times 1$$

Start with set of features, $S_o = \emptyset$. At iteration k , select a feature j such that

$$j_k = \arg \min_j \min_{\underline{w}} \| X_{S_{k-1} \cup \{j\}} \underline{w} - \underline{z} \|^2 = \arg \max_j \frac{|\tilde{\underline{x}}_j^T \underline{r}_{k-1}|^2}{\tilde{\underline{x}}_j^T \tilde{\underline{x}}_j}; \text{residual } \underline{r}_{k-1} = \underline{z} - X_{S_{k-1}} \underline{w}_{k-1}$$

$$\underline{w}_k = \arg \min_{\underline{w}} \| X_{S_k} \underline{w} - \underline{z} \|^2$$

- Implement using SVD
- L_1 regression (LASSO, compressed sensing) is an alternate method for feature selection
- Adaboost is a form of forward selection procedure

$$L_M(\underline{\alpha}) = \sum_{n=1}^N \exp \left\{ -z^n \sum_{m=1}^M \alpha_m g_m(\underline{x}^n) \right\}; g_m(\underline{x}^n) = \text{discriminant} (MLP, SVM, LDA, QDA, Logistic, \dots)$$



Alternative Feature Selection Methods

- Alternate criterion: Mutual Information

- Select feature (or feature subsets) with the maximum mutual information

$$k = \arg \max_{x_i} (H(\underline{z}) - H(\underline{z} | x_i))$$

- Applicable to both classification and regression

- Optimization: Select features via

- Local search
- Tabu search
- Genetic Algorithms
- JMI and CMMI

- References

- Guyon, I., and A. Elisseeff, "An Introduction to Variable and Feature Selection," *JMLR*, Vol. 3, 2003, pp. 1157-1182.
- Brown, J., A. Pocock, M-J. Zhao and M. Lujan, "Conditional Likelihood Maximization: A Unifying Framework for Information Theoretic Feature Selection," *JMLR*, Vol. 13, 2012, pp. 27-66.

T_k = Current Test set used on the path leading to node k

Joint Mutual Information (JMI) Criterion: (works well)

$$\begin{aligned} j_k &= \arg \max_j \sum_{l \in T_k} I(\{t_l, t_j\}; z) = \arg \max_j \sum_{l \in T_k} [I(t_l; z) + I(t_j; z | t_l)] \\ &= \arg \max_j \sum_{l \in T_k} I(t_j; z | t_l) \end{aligned}$$

Conditional Mutual Information Maximization (CMMI) Criterion:

$$j_k = \arg \max_j \min_{l \in T_k} [I(t_j; z | t_l)]$$



Dimensionality Reduction Methods

- Principal Component Analysis (PCA) ~ SVD
Consider the input data matrix (mean removed)

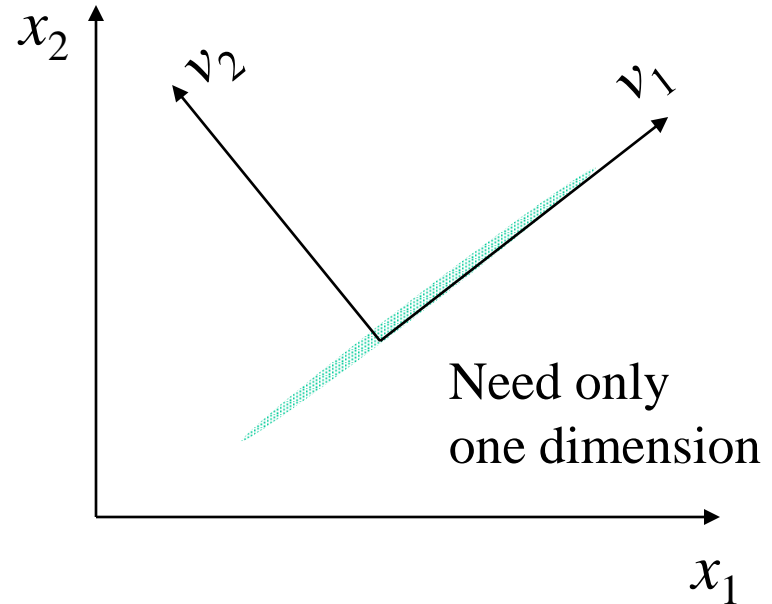
$$X = \begin{bmatrix} (\underline{x}_1 - \bar{x})^T \\ (\underline{x}_2 - \bar{x})^T \\ \vdots \\ (\underline{x}_N - \bar{x})^T \end{bmatrix} \quad \text{an } N \times P \text{ matrix}$$

$$X = U\Sigma V^T \quad \Sigma \ni \sigma_1 > \sigma_2 > \dots > \sigma_P$$

$$X^T X = \sum_{i=1}^N (\underline{x}_i - \bar{x})(\underline{x}_i - \bar{x})^T = V\Sigma^2 V^T = V\Lambda V^T$$

If $\sigma_1 > \sigma_2 > \dots > \sigma_r \gg \sigma_{r+1}$, the best r feature vectors are

$z_i = \underline{v}_i^T \underline{x}; i = 1, 2, \dots, r \sim$ first r principal components define the new features



Model Selection in PCA

- Applications of PCA: Pre-processing; Visualization; Data compression
- “Scree” plot: Sum of discarded eigenvalues of $X^T X$. If pick $L < r$ components

$$J(L) = \sum_{k=L+1}^p \lambda_k$$

“plot of $J(L)$ looks like the side of a mountain and “scree” refers to the debris fallen at the base”

- Fraction of variance explained

$$F(L) = \frac{\sum_{k=1}^L \lambda_k}{\sum_{k=1}^p \lambda_k}$$

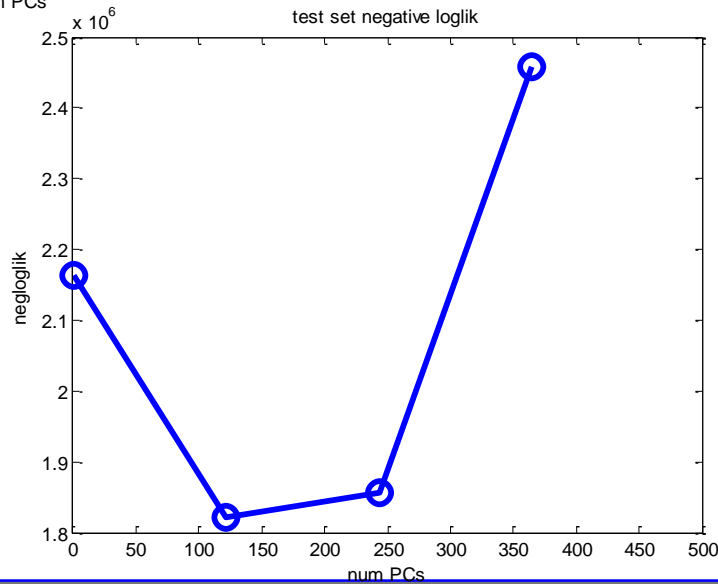
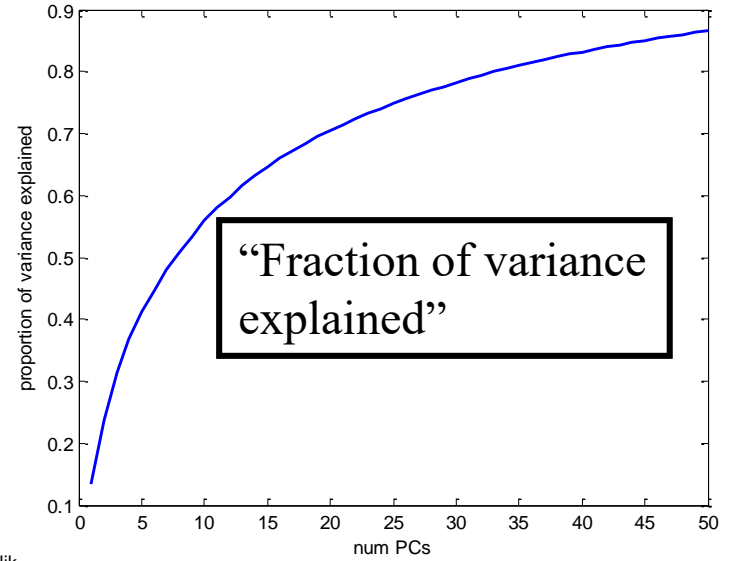
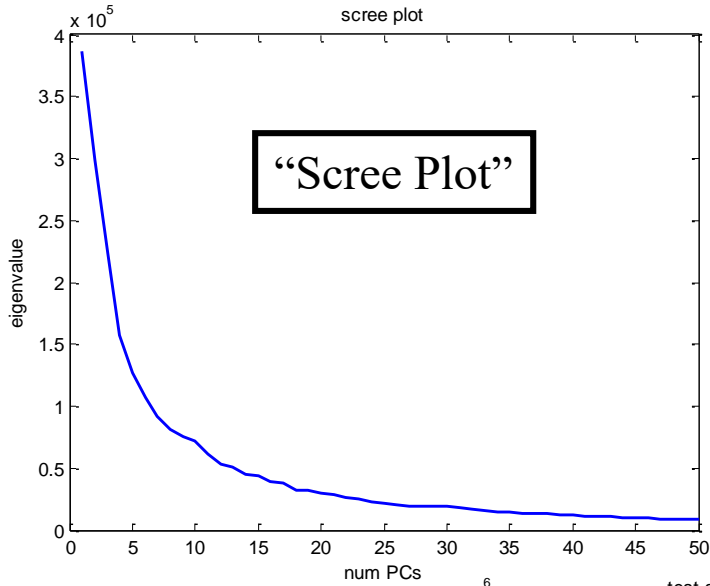
“Ideas are applicable to K-means as well”

- Profile likelihood

$$NLL(L) = \left(\sum_{k=1}^L \left(\ln \sigma(L) + \frac{1}{2} \left(\frac{\lambda_k - \mu_1(L)}{\sigma(L)} \right)^2 \right) + \sum_{k=L+1}^p \left(\ln \sigma(L) + \frac{1}{2} \left(\frac{\lambda_k - \mu_2(L)}{\sigma(L)} \right)^2 \right) \right)$$
$$\mu_1(L) = \left(\frac{\sum_{k=1}^L \lambda_k}{L} \right); \mu_2(L) = \left(\frac{\sum_{k=L+1}^p \lambda_k}{p-L} \right); \sigma^2(L) = \frac{\sum_{k=1}^L (\lambda_k - \mu_1(L))^2 + \sum_{k=L+1}^p (\lambda_k - \mu_2(L))^2}{p}$$

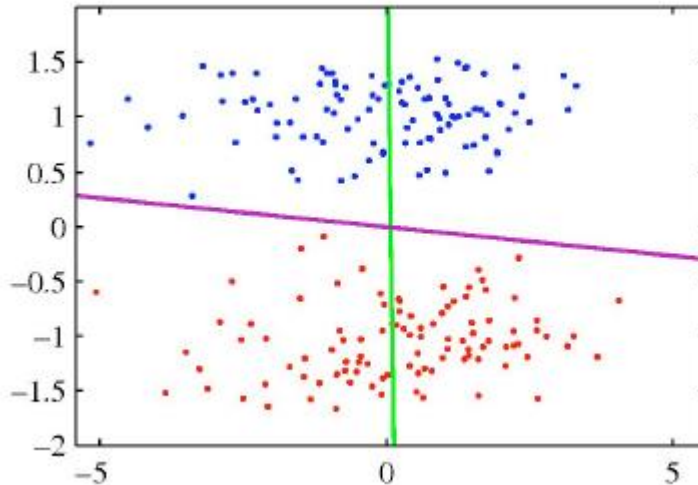


Illustration of Model Selection Metrics



PCA versus Fisher's Discriminant

- PCA versus Fisher's Linear Discriminant Analysis
 - Unsupervised (PCA) versus supervised (LDA)



LDA: Green
PCA: Magenta

- PCA chooses direction of maximum variance leading to strong class overlap (unsupervised)
- LDA takes into account the class labels (supervised), leading to a projection into the green curve

Probabilistic PCA - 1

Latent (Hidden) Variables : $\underline{z} \sim N(\underline{0}, I_z)$

Measurements : $p(\underline{x} | \underline{z}) = N(W\underline{z} + \underline{\mu}, \sigma^2 I_x)$; W a $p \times r$ matrix

Example : $\underline{x} = W\underline{z} + \underline{\mu} + \underline{v}$; $\underline{v} = N(\underline{0}, \sigma^2 I_x)$

$$\Rightarrow p(\underline{x}) = \int N(W\underline{z} + \underline{\mu}, \sigma^2 I_x) N(\underline{0}, I_z) d\underline{z} = N(\underline{\mu}, WW^T + \sigma^2 I_x) = N(\underline{\mu}, \Sigma_{xx})$$

$$\Rightarrow p(\underline{z} | \underline{x}) = N[E(\underline{z} | \underline{x}), \Sigma_z]$$

where $E(\underline{z} | \underline{x}) = \Sigma_{zx} \Sigma_{xx}^{-1} (\underline{x} - \underline{\mu}) + \underline{\mu} = W^T \Sigma_{xx}^{-1} (\underline{x} - \underline{\mu}) + \underline{\mu} = \frac{M^{-1} W^T}{\sigma^2} (\underline{x} - \underline{\mu})$

$$\Sigma_z = I_z - \Sigma_{zx} \Sigma_{xx}^{-1} \Sigma_{zx}^T = I_z - W^T \Sigma_{xx}^{-1} W = (I_z + \frac{W^T W}{\sigma^2})^{-1} = M^{-1}$$

$W^T \Sigma_{xx}^{-1} = \frac{M^{-1} W^T}{\sigma^2}$... Recall Kalman gain can be written in terms of

predicted or updated covariance

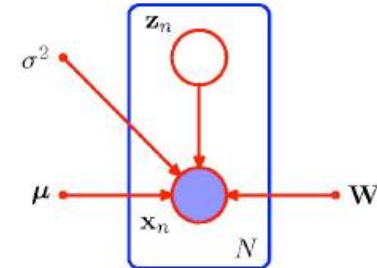
ML PCA : $\ln(p(\underline{X}^N | \underline{\mu}, W, \sigma^2)) = -\frac{N}{2} \{ p \ln(2\pi) + \ln |\Sigma_{xx}| + \text{Tr}(\Sigma_{xx}^{-1} S) \}$

where $S = \sum_{i=1}^N (\underline{x}^i - \underline{\mu})(\underline{x}^i - \underline{\mu})^T = U \Lambda_s U^T$

Recall Kernel Trick

Solution : $\underline{\mu} = \bar{\underline{x}}$; $W = U_r (\Lambda_r - \sigma^2 I)^{1/2} R$; $\sigma^2 = \frac{1}{p-r} \sum_{i=r+1}^p \lambda_i$; R orthogonal

Kernel PCA : Re place \underline{x}_i by $\phi(\underline{x}_i)$



$$\begin{aligned} & I_z - W^T \Sigma_{xx}^{-1} W \\ &= I_z - W^T (\sigma^2 I_x + WW^T)^{-1} W \\ &= (I_z + \frac{W^T W}{\sigma^2})^{-1} = M^{-1} \end{aligned}$$

$$\begin{aligned} & \begin{bmatrix} \underline{x} \\ \underline{z} \end{bmatrix} \sim N \left(\begin{bmatrix} \underline{\mu} \\ \underline{0} \end{bmatrix}; \begin{bmatrix} WW^T + \sigma^2 I_x & W \\ W^T & I_z \end{bmatrix} \right) \\ & K = W^T (\sigma^2 I_x + WW^T)^{-1} \\ &= \frac{W^T}{\sigma^2} \left(I_x + \frac{WW^T}{\sigma^2} \right)^{-1} \\ &= \frac{1}{\sigma^2} \left(W^T - \frac{W^T W}{\sigma^2} (I_z + \frac{W^T W}{\sigma^2})^{-1} W^T \right) \\ &= \frac{(I_z + \frac{W^T W}{\sigma^2})^{-1} W^T}{\sigma^2} = \frac{M^{-1} W^T}{\sigma^2} \end{aligned}$$



Supervised PCA (Bayesian Factor Regression)

Latent (Hidden) Variables : $\underline{z} \sim N(\underline{0}, I)$

Measurements : $p(\underline{x} | \underline{z}) = N(W_x \underline{z} + \underline{\mu}_x, \sigma_x^2 I_p)$; W_x a $p \times r$ matrix

$p(t | \underline{z}) = N(\underline{w}_t^T \underline{z} + \mu_t, \sigma_t^2)$; \underline{w}_t a $r \times 1$ column vector

(\underline{x}, t) is jointly Gaussian with mean $\begin{bmatrix} \underline{\mu}_x \\ \mu_t \end{bmatrix}$

and covariance matrix $\begin{bmatrix} W_x W_x^T + \sigma_x^2 I_p & W_x \underline{w}_t \\ \underline{w}_t^T W_x^T & \underline{w}_t^T \underline{w}_t + \sigma_t^2 \end{bmatrix}$

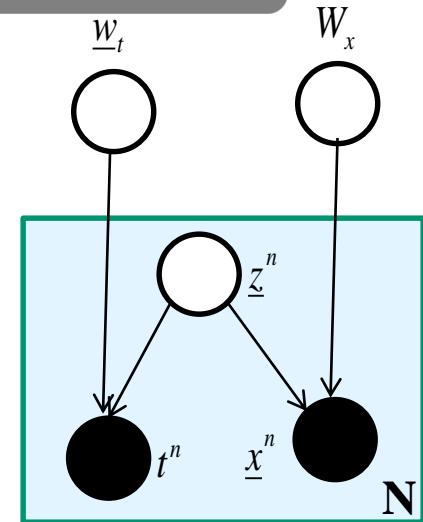
$$E(t | \underline{x}) = \mu_t + \Sigma_{tx} \Sigma_{xx}^{-1} (\underline{x} - \underline{\mu}_x) = \mu_t + \underline{w}_t^T W_x^T (W_x W_x^T + \sigma_x^2 I_p)^{-1} (\underline{x} - \underline{\mu}_x)$$

$$= \mu_t + \frac{\underline{w}_t^T W_x^T}{\sigma_x^2} \left(\frac{W_x W_x^T}{\sigma_x^2} + I_p \right)^{-1} (\underline{x} - \underline{\mu}_x) = \mu_t + \underbrace{\underline{w}_t^T \left(\frac{W_x^T W_x}{\sigma_x^2} + I_r \right)^{-1}}_{M^{-1}} \frac{W_x^T}{\sigma_x^2} (\underline{x} - \underline{\mu}_x)$$

$$\hat{\sigma}_t^2 = \sigma_t^2 + \underline{w}_t^T \underline{w}_t - \Sigma_{tx} \Sigma_{xx}^{-1} \Sigma_{tx}^T = \sigma_t^2 + \underline{w}_t^T \underline{w}_t - \underline{w}_t^T W_x^T (W_x W_x^T + \sigma_x^2 I_p)^{-1} W_x \underline{w}_t$$

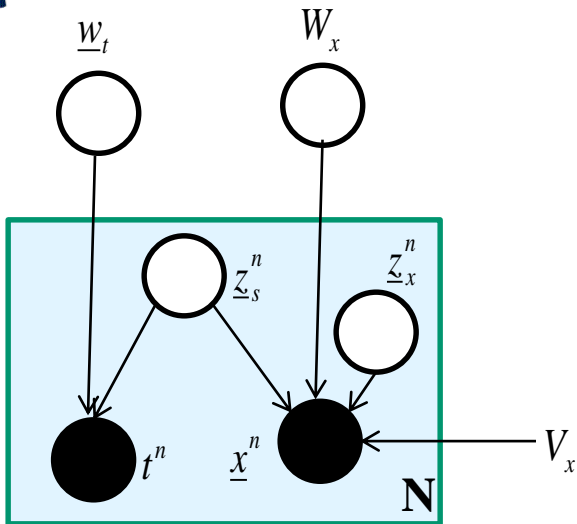
$$= \sigma_t^2 + \underline{w}_t^T \underline{w}_t - \frac{\underline{w}_t^T W_x^T}{\sigma_x^2} \left(\frac{W_x W_x^T}{\sigma_x^2} + I_p \right)^{-1} W_x \underline{w}_t = \sigma_t^2 + \underbrace{\underline{w}_t^T \left(\frac{W_x^T W_x}{\sigma_x^2} + I_r \right)^{-1}}_{M^{-1}} \underline{w}_t$$

Kernel PCA : Re place \underline{x} by $\phi(\underline{x})$



Want to estimate t based on \underline{x} and both \underline{x} and t are related to latent(state) variable \underline{z}

Partial Least Squares



PLS

Latent (Hidden) Variables : $\underline{z}_s \sim N(\underline{0}, I_s)$; $\underline{z}_x \sim N(\underline{0}, I_x)$

Measurements : $p(\underline{x} | \underline{z}) = N(W_x \underline{z}_s + V_x \underline{z}_x + \underline{\mu}_x, \sigma_x^2 I_p)$

$p(t | \underline{z}) = N(\underline{w}_t^T \underline{z}_s + \mu_t, \sigma_t^2)$

(\underline{x}, t) is jointly Gaussian with mean $\begin{bmatrix} \underline{\mu}_x \\ \mu_t \end{bmatrix}$

and covariance matrix $\begin{bmatrix} W_x W_x^T + V_x V_x^T + \sigma_x^2 I_p & W_x \underline{w}_t \\ \underline{w}_t^T W_x^T & \underline{w}_t^T \underline{w}_t + \sigma_t^2 \end{bmatrix}$

Kernel PLS : Re place \underline{x} by $\phi(\underline{x})$

$$E(t | \underline{x}) = \mu_t + \Sigma_{tx} \Sigma_{xx}^{-1} (\underline{x} - \underline{\mu}_x) = \mu_t + \underline{w}_t^T W_x^T (W_x W_x^T + V_x V_x^T + \sigma_x^2 I_p)^{-1} (\underline{x} - \underline{\mu}_x)$$

$$\hat{\sigma}_t^2 = \sigma_t^2 - \Sigma_{tx} \Sigma_{xx}^{-1} \Sigma_{tx}^T = \sigma_t^2 + \underline{w}_t^T [1 - W_x^T (W_x W_x^T + V_x V_x^T + \sigma_x^2 I_p)^{-1} W_x] \underline{w}_t$$

Good when both features and targets are noisy.



Canonical Correlation Analysis

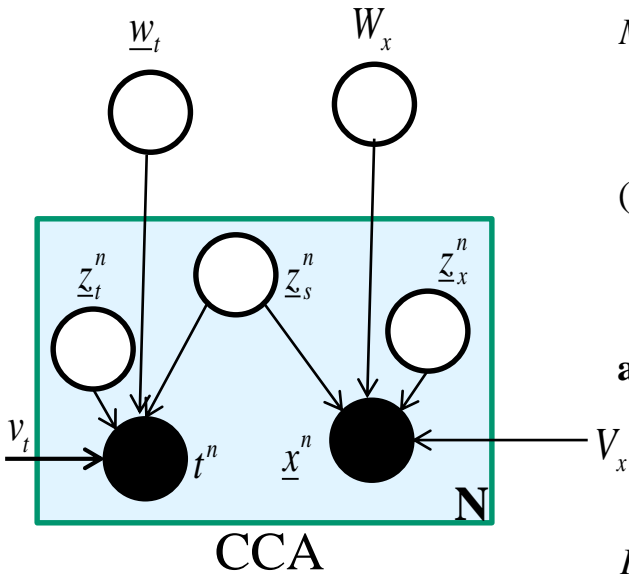
Latent (Hidden) Variables : $\underline{z}_s \sim N(\underline{0}, I_s)$; $\underline{z}_x \sim N(\underline{0}, I_x)$; $\underline{z}_t \sim N(\underline{0}, I_t)$

Measurements : $p(\underline{x} | \underline{z}) = N(W_x \underline{z}_s + V_x \underline{z}_x + \underline{\mu}_x, \sigma_x^2 I_p)$

$p(t | \underline{z}) = N(\underline{w}_t^T \underline{z}_s + \underline{v}_t^T \underline{z}_t + \mu_t, \sigma_t^2)$

(\underline{x}, t) is jointly Gaussian with mean $\begin{bmatrix} \underline{\mu}_x \\ \mu_t \end{bmatrix}$

and covariance matrix $\begin{bmatrix} \overbrace{W_x W_x^T + V_x V_x^T + \sigma_x^2 I_p}^{C_{xx}} & \overbrace{W_x \underline{w}_t}^{C_{xt}} \\ \underbrace{\underline{w}_t^T W_x^T}_{C_{tx}} & \underbrace{\underline{w}_t^T \underline{w}_t + \underline{v}_t^T \underline{v}_t + \sigma_t^2}_{C_{tt}} \end{bmatrix}$



Kernel CCA : Re place \underline{x} by $\phi(\underline{x})$

$$E(t | \underline{x}) = \mu_t + \Sigma_{tx} \Sigma_{xx}^{-1} (\underline{x} - \underline{\mu}_x) = \mu_t + \underline{w}_t^T W_x^T (W_x W_x^T + V_x V_x^T + \sigma_x^2 I_p)^{-1} (\underline{x} - \underline{\mu}_x) = \mu_t + C_{tx} C_{xx}^{-1} (\underline{x} - \underline{\mu}_x)$$

$$\hat{\sigma}_t^2 = \underline{v}_t^T \underline{v}_t + \sigma_t^2 \sigma_x^2 - \Sigma_{tx} \Sigma_{xx}^{-1} \Sigma_{tx}^T$$

$$= \sigma_t^2 + \underline{v}_t^T \underline{v}_t + \underline{w}_t^T [1 - W_x^T (W_x W_x^T + V_x V_x^T + \sigma_x^2 I_p)^{-1} W_x] \underline{w}_t = C_{tt} - C_{tx} C_{xx}^{-1} C_{xt}$$

Canonical correlation analysis (CCA) is a way of measuring the linear relationship between two sets of variables. It finds two bases, one for each variable, that are optimal with respect to correlations and, at the same time, it finds the corresponding correlations.

$$\rho^2 = 1 - \frac{\hat{\sigma}_t^2}{C_{tt}} = \frac{\underline{w}_t^T W_x^T (W_x W_x^T + V_x V_x^T + \sigma_x^2 I_p)^{-1} W_x \underline{w}_t}{\sigma_t^2 + \underline{v}_t^T \underline{v}_t + \underline{w}_t^T \underline{w}_t} = \frac{C_{tx} C_{xx}^{-1} C_{xt}}{C_{tt}}$$



Missing Data

- Missing Data Strategies:
 - Ignore missing data. May not have sufficient data
 - Replace by the mean
 - Replace by the median
 - Estimate x_k from the remaining data (EM algorithm)
 - If know distribution of input data, generate x_k randomly.



Information Theoretic Co-clustering

- Most clustering algorithms seek to cluster one dimension of the matrix (e.g., documents or columns) based on similarities along the second dimension (e.g., word distribution of documents or rows).
- For sparse, noisy, and high-dimensional data, *simultaneous clustering* (“co-clustering”, “bi-clustering”) of both rows and columns is beneficial.
 - Example: given a term-document matrix, co-clustering in two dimensions simultaneously clusters terms and documents
 - Other Examples: Marketing, Dimensionality Reduction, Currency Exchange,.....
 - More robust to sparsity than traditional single dimensional (e.g., terms or documents) clustering.
 - Co-clustering can be used as a pre-processor for supervised classification or as a classifier in its own right



Key Idea of Co-clustering

- Co-clustering Problem: Find maps

$$R(X) : \{x_1, x_2, \dots, x_m\} \rightarrow \{\hat{x}_1, \hat{x}_2, \dots, \hat{x}_k\} \quad C(Y) : \{y_1, y_2, \dots, y_n\} \rightarrow \{\hat{y}_1, \hat{y}_2, \dots, \hat{y}_l\}$$

to minimize $\min_{\hat{X}, \hat{Y}} [I(X; Y) - I(\hat{X}; \hat{Y})] \Rightarrow \max_{\hat{X}, \hat{Y}} I(\hat{X}; \hat{Y})$

- $\hat{X} = R(X)$ and $\hat{Y} = C(Y) \Rightarrow H(\hat{X} | X) = H(\hat{Y} | Y) = 0.$

$$\begin{aligned} I(X; Y) - I(\hat{X}; \hat{Y}) &= [H(X) - H(\hat{X})] + [H(Y) - H(\hat{Y})] + [H(\hat{X}, \hat{Y}) - H(X, Y)] \\ &= H(X | \hat{X}) + H(Y | \hat{Y}) + H(\hat{X}, \hat{Y}) - H(X, Y) \\ &= E_{p(x, y)} \left[\log_2 \frac{p(x, y)}{p(x | \hat{x}) p(\hat{x}, \hat{y}) p(y | \hat{y})} \right] = D(p(x, y) \| q(x, y)) \end{aligned}$$

$$q(x, y) = p(x | \hat{x}) p(\hat{x}, \hat{y}) p(y | \hat{y}) \text{ where } x \in \hat{x}, y \in \hat{y}.$$

Decomposition of pmf $p(x, y)$ into a product of three matrices

Illustration of Co-clustering

$$\begin{bmatrix}
 .05 & .05 & .05 & 0 & 0 & 0 \\
 .05 & .05 & .05 & 0 & 0 & 0 \\
 0 & 0 & 0 & .05 & .05 & .05 \\
 0 & 0 & 0 & .05 & .05 & .05 \\
 .04 & .04 & 0 & .04 & .04 & .04 \\
 .04 & .04 & .04 & 0 & .04 & .04
 \end{bmatrix}$$

$p(x, y)$

$$\begin{bmatrix}
 .5 & 0 & 0 \\
 .5 & 0 & 0 \\
 0 & .5 & 0 \\
 0 & .5 & 0 \\
 0 & 0 & .5 \\
 0 & 0 & .5
 \end{bmatrix}
 \begin{bmatrix}
 .3 & 0 \\
 0 & .3 \\
 .2 & .2
 \end{bmatrix}
 \begin{bmatrix}
 .36 & .36 & .28 & 0 & 0 & 0 \\
 0 & 0 & 0 & .28 & .36 & .36
 \end{bmatrix}
 =
 \begin{bmatrix}
 .054 & .054 & .042 & 0 & 0 & 0 \\
 .054 & .054 & .042 & 0 & 0 & 0 \\
 0 & 0 & 0 & .042 & .054 & .054 \\
 0 & 0 & 0 & .042 & .054 & .054 \\
 .036 & .036 & .028 & .028 & .036 & .036 \\
 .036 & .036 & .028 & .028 & .036 & .036
 \end{bmatrix}$$

$p(\hat{x}, \hat{y}) = \sum_{x \in \hat{x}} \sum_{y \in \hat{y}} p(x, y)$

$p(y | \hat{y}) = \frac{p(y)}{p(\hat{y})}$

$$p(x | \hat{x}) = \frac{p(x, \hat{x})}{p(\hat{x})} = \frac{p(x)}{p(\hat{x})}$$

$$q(x, y)$$

Note: $q(x, y) = p(x | \hat{x}) p(\hat{x}, \hat{y}) p(y | \hat{y}) = p(x) \underbrace{p(\hat{y} | \hat{x})}_{q(y|\hat{x})} p(y | \hat{y}) = q(x) q(y | \hat{x}) = p(y) \underbrace{p(x | \hat{x}) p(\hat{x} | \hat{y})}_{q(x|\hat{y})} = q(y) q(x | \hat{y})$

- #parameters that determine q are: $(m-k) + (kl-1) + (n-l)$



Co-clustering Algorithm

- **Step 1:** Set iteration $i=1$. Start with initial cluster maps (R_i, C_i) . Compute the pmfs

$$q^{(i,i)}(y | \hat{x}) = \sum_{\hat{y}} q^{(i,i)}(y | \hat{y}) q^{(i,i)}(\hat{y} | \hat{x}) = \sum_{\hat{y}} q^{(i,i)}(y | \hat{y}) \frac{q^{(i,i)}(\hat{x}, \hat{y})}{p(\hat{x})}$$

- **Step 2:** For every row x , assign it to the cluster that minimizes the K-L divergence $D(p(y|x) \| q^{(i,i)}(y|\hat{x}))$. The result is (R_{i+1}, C_i)

- **Step 3:** Compute the pmfs $q^{(i+1,i)}(\hat{x}, \hat{y}), q^{(i+1,i)}(x | \hat{x}), q^{(i+1,i)}(y | \hat{y}), q^{(i+1,i)}(x | \hat{y})$

$$q^{(i+1,i)}(x | \hat{y}) = \sum_{\hat{x}} q^{(i+1,i)}(x | \hat{x}) q^{(i+1,i)}(\hat{x} | \hat{y}) = \sum_{\hat{x}} q^{(i+1,i)}(x | \hat{x}) \frac{q^{(i+1,i)}(\hat{x}, \hat{y})}{p(\hat{y})}$$

- **Step 4:** For every column y , assign it to the cluster that minimizes the K-L divergence $D(p(x|y) \| q^{(i+1,i)}(x|\hat{y}))$. The result is (R_{i+1}, C_{i+1})
- **Step 5:** Compute the pmfs $q^{(i+1,i+1)}(\hat{x}, \hat{y}), q^{(i+1,i+1)}(x | \hat{x}), q^{(i+1,i+1)}(y | \hat{y}), q^{(i+1,i+1)}(y | \hat{x})$
Set $i=i+1$. Iterate Steps 2-5 until the K-L divergence converges.



1-D versus 2-D Clustering

Confusion Matrix

Co-Clustering (0.9835)			1-D Clustering (0.821)		
992	4	8	847	142	44
40	1452	7	41	954	405
1	4	1387	275	86	1099

I. Dhillon, 2003: CLASSIC 3 dataset

<ftp://ftp.cs.cornell.edu/pub/smart>



Summary

- Radial Basis Functions
- Gaussian Processes
- Relevant Vector Machines
- Feature Selection & Dimensionality Reduction
- Review of k-means and GMM-based Clustering
- Learning Vector Quantization
- Information-Theoretic Co-clustering
- Summary