

MATHEMATICS INFERENCE AND MACHINE LEARNING

PINGCHUAN MA

PARAMETER ESTIMATION

Definition.

Maximum Likelihood Estimation: A widely used approach to finding the desired parameters θ^* is maximum likelihood

$$\max_{\theta} p(\mathbf{y}|\mathbf{X}, \theta)$$

where $\mathbf{X} := [\mathbf{x}_1 | \dots | \mathbf{x}_N]^T \in \mathbb{R}^{N \times D}$, $\mathbf{y} := [y_1, \dots, y_N]^T \in \mathbb{R}^N$

Theorem.

$$\begin{aligned} \theta^* &\in \arg \min_{\theta} (-\log p(\mathbf{y}|\mathbf{X}, \theta)) \\ &= \arg \min_{\theta} (-\log \prod_{i=1}^N p(y_i|\mathbf{x}_i, \theta)) \\ &= \arg \min_{\theta} \sum_{i=1}^N -\log p(y_i|\mathbf{x}_i, \theta) \end{aligned}$$

In the linear regression model, the likelihood is Gaussian

$$\begin{aligned} -\log p(y_i|\mathbf{x}_i, \theta) &= \frac{1}{2\sigma^2} (y_i - \mathbf{x}_i^T \theta)^2 + \text{const} \\ L(\theta) &= -\log p(\mathbf{y}|\mathbf{X}, \theta) \\ &= \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\theta)^T (\mathbf{y} - \mathbf{X}\theta) \\ &= \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\theta\|^2 \end{aligned}$$

Taking the derivative w.r.t the parameter θ and set the gradient to 0, obtain

$$\begin{aligned} \frac{\partial L}{\partial \theta} &= \frac{1}{\sigma^2} (-\mathbf{y}^T \mathbf{X} + \theta \mathbf{X}^T \mathbf{X}) \\ \theta^* &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \end{aligned}$$

Remark. We consider the linear regression model:

$$y = \phi^T(\mathbf{x} + \epsilon)$$

where $\phi: \mathbb{R}^D \rightarrow \mathbb{R}^K$ is a nonlinear transformation of the input \mathbf{x} , we can simply replace \mathbf{X} by Φ , and obtain

$$\theta^* = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

Remark. In machine learning, the negative log likelihood function is also called an error function.

Definition.

The property of the polynomials fitting the noise structure is called overfitting.

Definition.

In regularization, a term is added to the log-likelihood that penalizes the amplitude of the parameter θ . Example of a loss function of the form

$$-\log p(\mathbf{y}|\mathbf{X}, \theta) + \lambda \|\theta\|_2^2$$

Definition.

Maximum A-Posterior Estimation: The parameter vector θ_{MAP} that maximizes the posterior.

Date: April 2017.

Theorem.

To find the MAP estimate. We start with the log-transform

$$\text{log}p(\theta|\mathbf{X}, \mathbf{y}) = \text{log}(\mathbf{y}|\mathbf{X}, \theta) + \text{log}p(\theta) + \text{const}$$

Remark. Choosing a Gaussian parameter prior $p(\theta) = N(0, b^2\mathbf{I})$, $b^2 = \frac{1}{2\lambda}$, the negative log-prior will be

$$-\text{log}p(\theta) = \lambda\theta^T\theta + \text{const}$$

That means for a quadratic regularization, the regularization parameter λ in corresponds to twice the precision of the Gaussian prior $p(\theta)$. The log-prior plays the role of a regularizer that penalizes implausible values.

Theorem.

For the linear regression problem where

$$y = \phi^T(\mathbf{x})\theta + \epsilon, \quad \epsilon \in N(0, \sigma^2)$$

The negative log-posterior for this model

$$-\text{log}p(\theta|\mathbf{X}, \mathbf{y}) = \frac{1}{2\sigma^2}(\mathbf{y} - \Phi\theta)^T(\mathbf{y} - \Phi\theta) + \frac{1}{2b^2}\theta^T\theta + \text{const}$$

To find θ_{MAP} , we take the derivative w.r.t the parameter θ and set the gradient to 0:

$$\frac{1}{\sigma^2}(\theta^T\Phi^T\Phi - \mathbf{y}^T\Phi) + \frac{1}{b^2}\theta^T = 0$$

$$\theta_{MAP} = (\Phi^T\Phi + \frac{\sigma^2}{b^2}\mathbf{I})^{-1}\Phi^T\mathbf{y}$$

Definition.

To find the local minimum of a function using gradient descent, one takes steps proportional to the negative of the gradient of the function at the current point.

$$\mathbf{x}_{i+1} = \mathbf{x}_i - \gamma_i(\nabla f)(\mathbf{x}_i)$$

Definition.

Cross-validation: K-fold cross-validation effectively partitions the data into K chunks, K-1 of which form the training set \tilde{D} , and the last chunk serves as the validation set \mathcal{V} . The expected generalization error can be computed as follows.

$$\mathbb{E}_{\mathcal{V}}[G(\mathcal{V}|M)] \approx \frac{1}{K} \sum_{k=1}^K G(\mathcal{V}^{(k)}|M)$$

where $G(\mathcal{V}|M)$ is the generalization error on the validation set \mathcal{V} for model M .

Definition.

Occam's Razor: finding the simplest model that explains the data reasonably well.

FEATURE EXTRACTION

Definition.

Eigen-decomposition: Assume square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ with n linearly independent eigenvectors $\mathbf{q}_i, i = 1, \dots, n$ and n eigenvalues $\lambda_1, \dots, \lambda_n$. Then \mathbf{A} can be factorised as

$$\mathbf{A} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^{-1}$$

where $\mathbf{\Lambda}$ is a diagonal matrix whose diagonal elements are the corresponding eigenvalues, i.e., $\mathbf{\Lambda} = \text{diag}\{\lambda_1, \dots, \lambda_n\}$.

Theorem.

if \mathbf{A} is symmetric,

$$\mathbf{A} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$$

Definition.

If $\mathbf{A} \in \mathbb{R}^{m \times n}$, then there exist orthogonal matrices

$$\mathbf{U} \in \mathbb{R}^{m \times m} \text{ and } \mathbf{V} \in \mathbb{R}^{n \times n}$$

such that

$$\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T$$

where $\mathbf{\Lambda}$ is a diagnosis matrices $\mathbf{\Lambda} \in \mathbb{R}^{m \times n}$

Theorem.

The 2-norm and Frobenius norm we have the following.

$$\begin{aligned} \|\mathbf{A}\|_F^2 &= \sigma_1^2 + \dots + \sigma_p^2, p = \min\{m, n\} \\ \|\mathbf{A}\|_2 &= \sigma_1 \end{aligned}$$

Theorem.

Let the SVD of $\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^T$. If $k < r = \text{rank}(\mathbf{A})$ and

$$\mathbf{A}_k = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^T$$

then

$$\min_{\text{rank}(\mathbf{B})=k} \|\mathbf{A} - \mathbf{B}\|_2 = \|\mathbf{A} - \mathbf{A}_k\|_2 = \sigma_{k+1}$$

Proof. Since, $\mathbf{U}^T \mathbf{A}_k \mathbf{V} = \text{diag}(\sigma_1, \dots, \sigma_k, 0, \dots, 0)$ then $\text{rank}(\mathbf{A}_k)=k$ and $\mathbf{U}(\mathbf{A} - \mathbf{A}_k)\mathbf{V} = \text{diag}(0, \dots, 0, \sigma_{k+1}, \dots, \sigma_p)$. Hence, $\|\mathbf{A} - \mathbf{A}_k\|_2 = \sigma_{k+1}$. Now suppose $\text{rank}(\mathbf{B})=k$ for a $\mathbf{B} \in \mathbb{R}^{m \times n}$. It follows that

$$\text{span}\{\mathbf{q}_1, \dots, \mathbf{q}_{n-k}\} \cap \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_{k+1}\} \neq \{0\}$$

we have that

$$\|\mathbf{A} - \mathbf{B}\|_2^2 \geq \|\mathbf{A} - \mathbf{Bz}\|_2^2 = \|\mathbf{Az}\|_2^2 \geq \sigma_{k+1}^2 \sum_{i=1}^{k+1} \alpha_i^2 = \sigma_{k+1}^2$$

□

Remark. Assume we have a collection of n data samples $\mathbf{x}_1, \dots, \mathbf{x}_n$, we stack the samples as columns of matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$. The k low-rank representation of the data is

$$\mathbf{X}_k = \mathbf{U}_k \mathbf{\Lambda}_k \mathbf{V}_k^T$$

COMPONENT ANALYSIS

Theorem.

Principal component analysis performs a linear mapping of the data to a lower-dimensional space in such a way that the variance of the data in low-dimensional representation is maximized.

Theorem.

The optimal features $\{y_1^o, \dots, y_n^o\} = \arg \max \frac{1}{2} \sigma_y^2$

$$\begin{aligned} \mathbf{w}_o &= \arg \max_{\mathbf{w}} \frac{1}{2n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i)^2 \\ &= \arg \max_{\mathbf{w}} \frac{1}{2n} \sum_{i=1}^n \mathbf{w}^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{w} \\ &= \arg \max_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{S}_t \mathbf{w} \end{aligned}$$

Theorem.

In order to avoid the trivial solution we incorporate the constraint $\|\mathbf{w}\|_2 = 1$. The optimisation problem can be reformulated as

$$\begin{aligned} \mathbf{w}_o &= \arg \max_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{S}_t \mathbf{w} \\ \text{s.t. } & \mathbf{w}^T \mathbf{w} = 1 \end{aligned}$$

This Lagrangian of the above optimisation problem is

$$L(\mathbf{w}, \lambda) = \frac{1}{2} \mathbf{w}^T \mathbf{S}_t \mathbf{w} - \lambda (\mathbf{w}^T \mathbf{w} - 1)$$

Forcing $\nabla_{\mathbf{w}} L(\mathbf{w}) = \mathbf{0}$, we end up to

$$\mathbf{S}_t \mathbf{w} = \lambda \mathbf{w}$$

Replacing it to the optimisation problem, we obtain

$$\lambda_o = \arg \max_{\lambda} \lambda$$

Theorem.

Assume $\mathbf{y}_i \in \mathbb{R}^d$

$$\begin{aligned} \mathbf{W}_o &= \arg \max_{\mathbf{W}} \frac{1}{2} \sum_{k=1}^d \mathbf{w}_k^T \mathbf{S}_t \mathbf{w}_k \\ &= \arg \max_{\mathbf{W}} \text{tr}(\mathbf{W}^T \mathbf{S}_t \mathbf{W}) \\ &= \arg \max_{\mathbf{W}} \text{tr}(\mathbf{\Lambda}_d) \\ &= \sum_{k=1}^d \lambda_k \end{aligned}$$

where $\mathbf{S}_t = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$, and $\mathbf{W}^T \mathbf{S}_t \mathbf{W} = \mathbf{W}^T \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T \mathbf{W}$.

Theorem.

Assume $\mathbf{B} = \mathbf{X} \mathbf{X}^T$ and $\mathbf{C} = \mathbf{X}^T \mathbf{X}$, \mathbf{B} and \mathbf{C} have the same positive eigenvalues $\mathbf{\Lambda}$. Assume $N < F$, the eigenvectors of \mathbf{U} and \mathbf{V} can be represented as

$$\mathbf{U} = \mathbf{X} \mathbf{V} \mathbf{\Lambda}^{-1/2}$$

Proof. Assume $\mathbf{X}^T \mathbf{X} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T$

$$\begin{aligned} \mathbf{S}_t &= \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T \\ \mathbf{\Lambda} &= \mathbf{U}^T \mathbf{X} \mathbf{X}^T \mathbf{U} \\ \mathbf{\Lambda} &= \mathbf{\Lambda}^{-1/2} \mathbf{V}^T \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T \mathbf{V} \mathbf{\Lambda}^{-1/2} \end{aligned}$$

In which we can obtain $\mathbf{U} = \mathbf{X} \mathbf{V} \mathbf{\Lambda}^{-1/2}$. □

- Compute $\mathbf{X}^T \mathbf{X} = [(\mathbf{x}_i - \mu)^T (\mathbf{x}_j - \mu)]$
- Perform eigenanalysis of $\mathbf{X}^T \mathbf{X} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T$
- Compute eigenvectors $\mathbf{U} = \mathbf{X} \mathbf{V} \mathbf{\Lambda}^{-1/2}$
- Compute d features $\mathbf{Y} = \mathbf{U}^T \mathbf{X}$

Definition.

Considering whitening,

$$\mathbf{Y} = \mathbf{U}^T \mathbf{X}$$

we want to $\mathbf{Y} \mathbf{Y}^T = \mathbf{U}^T \mathbf{X} \mathbf{X}^T \mathbf{U} = \mathbf{\Lambda}$, therefore, we obtain

$$\mathbf{W} = \mathbf{U} \mathbf{\Lambda}^{-1/2}$$

Definition.

Linear Discriminant Analysis is a method to find a linear combination of features that characterizes or separates two or more classes of objects or event. It can be formulated as follows.

$$\begin{aligned} \min \quad & \sigma_y^2(c_1) + \sigma_y^2(c_2) \\ \text{s.t.} \quad & (\mu_y(c_1) - \mu_y(c_2))^2 \end{aligned}$$

Definition.

The formulation is given as follows.

$$\begin{aligned} \max \quad & \mathbf{w}^T \mathbf{S}_b \mathbf{w} \\ \text{s.t.} \quad & \mathbf{w}^T \mathbf{S}_w \mathbf{w} = 1 \end{aligned}$$

Using lagrangian:

$$\begin{aligned} L(\mathbf{w}, \lambda) &= \mathbf{w}^T \mathbf{S}_b \mathbf{w} - \lambda(\mathbf{w}^T \mathbf{S}_w \mathbf{w} - 1) \\ \frac{\partial L}{\partial \mathbf{w}} &= 0 \\ \lambda \mathbf{S}_w \mathbf{w} &= \mathbf{S}_b \mathbf{w} \end{aligned}$$

Remark. PCA is an unsupervised approach for compression of data. Whitened PCA is a technique to make the input less redundant. Compared with PCA, whitened PCA can reduce the correlation among features. After whitening, each feature has the same variance. LDA is a supervised approach for compression of data. As all data is labelled, this method takes classification into account. It is different from PCA because PCA doesn't take labels into considerations. The maximal number of dimension is $C - 1$, where C is the number of labels.

SUPPORT VECTOR MACHINES

Definition.

The SVM techniques tries to find the separating hyperplane with the largest margin between two classes, measured along a line perpendicular to the hyperplane.

Remark. This means we find a line with parameters \mathbf{w} and b such that the distance between $\mathbf{w}^T \mathbf{x} + b = \pm 1$

Theorem.

Assume a point \tilde{x} on $\mathbf{w}^T \mathbf{x} + b = -1$, and we assume point $\tilde{x} + t\mathbf{w}$ touches the line $\mathbf{w}^T \mathbf{x} + b = 1$. Then, $t\mathbf{w}^T \mathbf{w} = 2$. So the length of $t\mathbf{w}$ is $\|t\mathbf{w}\|_2 = \frac{2}{\|\mathbf{w}\|_2}$. As maximizing $\frac{2}{\|\mathbf{w}\|_2}$ is equivalent to minimizing $\frac{\mathbf{w}^T \mathbf{w}}{2}$. So the equivalent problem to SVM is as follows.

$$\begin{aligned} \min_{w,b} \quad & \frac{\mathbf{w}^T \mathbf{w}}{2} \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \\ & i = 1, \dots, l \end{aligned}$$

Definition.

Constrained optimization problem is of the form

$$\begin{aligned} \min_w \quad & f(\mathbf{w}) \\ \text{s.t.} \quad & g(\mathbf{w}) \leq 0 \end{aligned}$$

We define the Lagrangian to be the original objective function added to a weighted combination of the constraints.

$$L(\mathbf{w}, a) = f(\mathbf{w}) + ag(\mathbf{w})$$

Theorem.

The original minimization problem can be written as

$$\min_{\mathbf{w}} \max_{a \geq 0} L(\mathbf{w}, a)$$

Proof. $g(\mathbf{w}) \leq 0$, we maximize $L(\mathbf{w}, a)$ by setting $a = 0$. When $g(\mathbf{w}) \geq 0$, we can get $\max_{a \geq 0} L(\mathbf{w}, a) = \infty$. Minimizing the outer loop, we obtain the minimum value of $f(\mathbf{w})$ such that $g(\mathbf{w}) \leq 0$ holds. \square

Definition.

The primal solution to the problem is given by

$$\mathbf{p}^* = \min_{\mathbf{w}} \max_{a \geq 0} L(\mathbf{w}, a)$$

The dual solution to the problem is given by

$$\mathbf{d}^* = \max_{a \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, a)$$

Definition.

$d^* \leq p^*$. Let \mathbf{w}^* be the \mathbf{w} value that corresponds to the optimal primal solution p^* . We can write for all $a \geq 0$

$$\max_{\tilde{a} \geq 0} L(\mathbf{w}^*, \tilde{a}) \geq L(\mathbf{w}^*, a) \geq \min_{\mathbf{w}} L(\mathbf{w}, a)$$

Theorem.

If certain conditions are met, namely

- $f(\mathbf{w})$ is convex
- $g(\mathbf{w})$ is affine

then $d^* = p^*$

Definition.

Support vector regression solves the following optimization problem:

$$\begin{aligned} \min_{w,b,\xi,\xi^*} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^l (\xi_i + \xi_i^*) \\ \text{s.t.} \quad & (\mathbf{w}^T \phi(\mathbf{x}_i) + b) - y_i \leq \epsilon + \xi_i \\ & y_i - (\mathbf{w}^T \phi(\mathbf{x}_i) + b) \leq \epsilon + \xi_i^* \\ & \xi_i, \xi_i^* \geq 0, \quad i = 1, \dots, l \end{aligned}$$

Slack variables ξ_i is the upper training error (ξ_i^* is the lower) subject to the ϵ -intensive tube $|y_i - (\mathbf{w}^T \phi(\mathbf{x}_i) + b)| \leq \epsilon$, $\mathbf{w}^T \mathbf{w}$ is added to smooth the function.

Remark. Linear regression finds a linear function $\mathbf{w}^T \mathbf{x} + b$ so that (\mathbf{w}, b) is an optimal solution of

$$\min_{\mathbf{w}, b} \sum_{i=1}^n (y_i - (\mathbf{w}^T \mathbf{x}_i + b))^2$$

If the data is nonlinearly distributed, a linear function is not enough, and we also want to map data to higher dimensional space by a function $\phi(\mathbf{x})$. Therefore $F \leq$ dimensionality of $\phi(\mathbf{x})$ so again overfitting happens. Support Vector Regression can remedy the overfitting problem in these two strategies.

- A threshold ϵ is given, if the point (\mathbf{x}_i, y_i) is located in the threshold, then $\xi_i, \xi_i^* = 0$.
- To reduce the model complexity and smooth the function, an additional term $\mathbf{w}^T \mathbf{w}$ is added to the objective function.

The difference between linear regression and SVM:

- loss function is different
- LR:

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right]$$

SVM:

$$\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i (y_i (w^T x_i + b) - 1)$$

LR:

$$\min_w \lambda \|w\|_2 + \sum_i \log(1 + \exp^{-y_i w^T x_i})$$

SVM:

$$\min_w \lambda \|w\|_2 + \sum_i \max\{0, 1 - y_i w^T x_i\}$$

- SVM considers the point around the optimal hyperplane while LR considers all points.
- As for solving nonlinear problems, SVM adapt kernel functions while LR don't because of the high complexity.
- SVM has a regularizer.

KERNEL PCA

Theorem.

In kernel Principal Component Analysis, $\phi(\cdot)$ may not be explicitly known or is expensive to compute and store. What we explicitly known is the dot product in H .

$$\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j)$$

Remark. All positive semi-definite functions can be used as kernels.

- Gaussian Radial Basis Function (RBF) kernel:
 $k(\mathbf{x}_i, \mathbf{x}_j) = e^{-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / r^2}$
- Polynomial kernel:
 $k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j + b)^n$
- Hyperbolic Tangent kernel:
 $k(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\mathbf{x}_i^T \mathbf{x}_j + b)$

CLUSTERING

Definition.

k-means clustering aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean.

Theorem.

Given a set of observations $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$, where each observation is a d -dimensional real vector. The objective is to find:

$$\arg \min_S \sum_{i=1}^k \sum_{x \in S_i} \|\mathbf{x} - \mu_i\|^2 = \arg \min_S \sum_{i=1}^k |S_i| \text{Var} S_i$$

The process is as follows.

Assignment step:

$$S_i^{(t)} = \{x_p : \|x_p - m_i^{(t)}\|^2 \leq \|x_p - m_j^{(t)}\|^2, 1 \leq j \leq k\}$$

Update step:

$$m_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{x_j \in S_i^{(t)}} x_j$$

Definition.

An expectation-maximization (EM) algorithm is an iterative method to find maximum likelihood or maximum a posterior estimates of parameters, where the model depends on latent variables.

Definition.

For Gaussian Mixture Models,

$$\begin{aligned} p(\mathbf{X}, \mathbf{Z}|\theta) &= p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N, \mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}|\theta) \\ &= \prod_{n=1}^N p(\mathbf{x}_n|\mathbf{z}_n, \theta_x) \prod_{n=1}^N p(\mathbf{z}_n|\theta_z) \\ &= \prod_{n=1}^N \prod_{k=1}^3 N(\mathbf{x}_n|\boldsymbol{\mu}_k, \sum_k)_{z_{nk}} \prod_{n=1}^N \prod_{k=1}^3 \pi_k^{z_{nk}} \end{aligned}$$

$$\ln p(\mathbf{X}, \mathbf{Z}|\theta) = \sum_{n=1}^N \sum_{k=1}^3 z_{nk} \{ \ln N(\mathbf{x}_n|\boldsymbol{\mu}_k, \sum_k) + \ln \pi_k \}$$

Applying operator $E_{p(\mathbf{Z}|\mathbf{X}, \theta)}$

$$\begin{aligned} &E_{p(\mathbf{Z}|\mathbf{X}, \theta)}[\ln p(\mathbf{X}, \mathbf{Z}|\theta)] \\ &= \sum_{n=1}^N \sum_{k=1}^3 E_{p(\mathbf{Z}|\mathbf{X}, \theta)}[z_{nk}] \{ \ln N(\mathbf{x}_n|\boldsymbol{\mu}_k, \sum_k) + \ln \pi_k \} \\ G(\theta) &= E_{p(\mathbf{Z}|\mathbf{X}, \theta)}[\ln p(\mathbf{X}, \mathbf{Z}|\theta)] \\ &= \sum_{n=1}^N \sum_{k=1}^3 \gamma(z_{nk}) \left\{ -\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \sum_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right. \\ &\quad \left. - \frac{1}{2} (\ln 2\pi + \ln |\sum_k|) + \ln \pi_k \right\} \end{aligned}$$

Theorem.

$$\begin{aligned}
 G(\theta) &= E_{p(Z|X,\theta)}[\ln p(X, Z|\theta)] \\
 &= \sum_{n=1}^N \sum_k E_{p(Z|X,\theta)[z_{nk}]} \ln N(x_n|\mu_k, \sum_k) + \ln \pi_k \\
 \frac{dG(\theta)}{d\mu_k} &= 0 \quad \frac{dG(\theta)}{d\sum_k} = 0 \\
 \pi_k &= \frac{\sum_{n=1}^N \gamma(z_{nk})}{N} \\
 \mu_k &= \frac{\sum_{n=1}^N \gamma(z_{nk}) x_n}{\sum_{n=1}^N \gamma(z_{nk})} \\
 \sum_k &= \frac{\sum_{n=1}^N \gamma(z_{nk}) (x_n - \mu_k)(x_n - \mu_k)^T}{\sum_{n=1}^N \gamma(z_{nk})}
 \end{aligned}$$

PROBABILISTIC PRINCIPAL COMPONENT ANALYSIS

MARKOV CHAINS

Definition.

Given a set of observations $\mathbf{D}_l = \{x_1^l, x_2^l, \dots, x_T^l\}$, $l = 1, \dots, N$, we need to find the parameters $\theta = \{\pi, \mathbf{A}\}$ that maximize $p(\mathbf{D}_1, \mathbf{D}_2, \dots, \mathbf{D}_N|\theta)$

Theorem.

$$\begin{aligned}
 p(\mathbf{D}_1, \mathbf{D}_2, \dots, \mathbf{D}_N|\theta) &= \prod_{l=1}^N p(\mathbf{D}_l|\theta) \\
 &= \prod_{l=1}^N p(\mathbf{x}_1^l) \prod_{t=2}^T p(\mathbf{x}_t^l|\mathbf{x}_{t-1}^l) \\
 &= \prod_{l=1}^N \prod_{k=1}^5 \pi_k^{x_{1k}^l} \prod_{t=2}^T \prod_{j=1}^5 \prod_{k=1}^5 a_{jk}^{x_{t-1j}^l x_{tk}^l} \\
 \ln p(\mathbf{D}_1, \mathbf{D}_2, \dots, \mathbf{D}_N|\theta) &= \sum_{l=1}^N \sum_{k=1}^5 x_{1k}^l \ln \pi_k + \\
 &\quad \sum_{l=1}^N \sum_{t=2}^T \sum_{j=1}^5 \sum_{k=1}^5 x_{1k}^l \ln a_{jk} \\
 &= \sum_{k=1}^5 \sum_{l=1}^N x_{1k}^l \ln \pi_k + \\
 &\quad \sum_{j=1}^5 \sum_{k=1}^5 \sum_{l=1}^N \sum_{t=2}^T x_{1k}^l \ln a_{jk}
 \end{aligned}$$

HIDDEN MARKOV MODEL

Definition.

Hidden Markov model is a statistical Markov model in which the system begin modeled in assumed to be a Markov process with hidden states.

Definition.

Filtering: $p(z_t|\mathbf{x}_1, \dots, \mathbf{x}_t) = \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_t, z_t)}{p(\mathbf{x}_1, \dots, \mathbf{x}_t)}$

Smoothing: $p(z_t|\mathbf{x}_1, \dots, \mathbf{x}_T)$

Decoding: $\arg \max_{\mathbf{y}_1, \dots, \mathbf{y}_t} p(\mathbf{y}_1, \dots, \mathbf{y}_t | \mathbf{x}_1, \dots, \mathbf{x}_t)$
 Evaluation: $p(\mathbf{x}_1, \dots, \mathbf{x}_T)$

Theorem.

The parameters emission probability $p(\mathbf{x}|\mathbf{z})$, probability for the first timestamp $p(\mathbf{x}_1)$ and transition matrix \mathbf{A} are obtained by maximizing the probability. After taking the expectation with regard to the posterior

$$\begin{aligned} G(\theta) &= \sum_{l=1}^N \sum_{t=1}^T \sum_{j=2}^5 \sum_{k=1}^K x_{tj}^l E[z_{tk}^l] \ln b_{jk} + \sum_{l=1}^N \sum_{k=1}^K E[z_{1k}^l] \ln \pi_k \\ &+ \sum_{l=1}^N \sum_{t=2}^T \sum_{j=1}^5 \sum_{k=1}^K E[z_{t-1j}^l z_{tk}^l] \ln a_{jk} \\ E[z_{1k}^l] &= \sum_{z_{1k}^l} z_{1k}^l p(z_{1k}^l = 1 | \mathbf{x}_1^l, \dots, \mathbf{x}_T^l) = \frac{\alpha(z_{1k}^l) \beta(z_{1k}^l)}{p(\mathbf{x}_1^l, \dots, \mathbf{x}_T^l)} \\ E[z_{tk}^l] &= \sum_{z_{tk}^l} z_{tk}^l p(z_{tk}^l = 1 | \mathbf{x}_1^l, \dots, \mathbf{x}_T^l) = \frac{\alpha(z_{tk}^l) \beta(z_{tk}^l)}{p(\mathbf{x}_1^l, \dots, \mathbf{x}_T^l)} \\ E[z_{t-1j}^l z_{tk}^l] &= \sum_{z_{t-1j}^l} \sum_{z_{tk}^l} z_{t-1j}^l z_{tk}^l p(z_{t-1j}^l = 1, z_{tk}^l = 1 | \mathbf{x}_1^l, \dots, \mathbf{x}_T^l) \\ &= \frac{a(z_{t-1j}^l) \prod_{r=1}^5 b_{kr}^{x_{tr}^l} a_{jk} \beta(z_{tk}^l)}{p(\mathbf{x}_1^l, \dots, \mathbf{x}_T^l)} \end{aligned}$$

The Lagrangian can be applied.

$$\begin{aligned} L(\theta) &= \sum_{l=1}^N \sum_{t=1}^T \sum_{j=2}^5 \sum_{k=1}^K x_{tj}^l E[z_{tk}^l] \ln b_{jk} + \sum_{l=1}^N \sum_{k=1}^K E[z_{1k}^l] \ln \pi_k \\ &+ \sum_{l=1}^N \sum_{t=2}^T \sum_{j=1}^5 \sum_{k=1}^K E[z_{t-1j}^l z_{tk}^l] \ln a_{jk} \\ &+ \lambda \left(\sum_{j=1}^5 b_{jk} - 1 \right) + \gamma \left(\sum_{k=1}^K \pi_k - 1 \right) + \delta \left(\sum_{k=1}^K a_{jk} - 1 \right) \end{aligned}$$

which gives us

$$\begin{aligned} b_{jk} &= \frac{\sum_{l=1}^N \sum_{t=1}^T E[z_{tk}^l] x_{tj}^l}{\sum_{l=1}^N \sum_{t=1}^T E[z_{tk}^l]} \\ \pi_k &= \frac{\sum_{l=1}^N E[z_{1k}^l]}{\sum_{l=1}^N \sum_{r=1}^K E[z_{1r}^l]} \\ a_{jk} &= \frac{\sum_{l=1}^N \sum_{t=2}^T E[z_{t-1j}^l z_{tk}^l]}{\sum_{r=1}^K \sum_{l=1}^N \sum_{t=2}^T E[z_{t-1j}^l z_{tr}^l]} \end{aligned}$$

MARKOV RANDOM FIELDS

Definition.

Local Markovianity: given its neighborhood, a variable is independent on the rest of the variables.

$$p(x_i | x_{V \setminus \{i\}}) = p(x_i | x_{N(i)})$$

Definition.

Global Markovianity: Let A, B, C be three disjoint subset of V . If C separates A from $B \rightarrow$

$$p(\mathbf{x}_A, \mathbf{x}_B | \mathbf{x}_C) = p(\mathbf{x}_A | \mathbf{x}_C) p(\mathbf{x}_B | \mathbf{x}_C)$$

Theorem.

If the field has the local Markov property, then $p(x)$ can be written as a Gibbs distribution.

$$p(x) = \frac{1}{Z} \prod_c \Psi_c(\mathbf{x}_c) = \frac{1}{Z} e^{-\sum_c V_c(\mathbf{x}_c)}$$

Definition.

A random vector $\mathbf{x} = (x_1, \dots, x_N)^T$ is called a GMRF wrt the graph $G = V, E$ with mean μ and a positive semi-definite precision matrix \mathbf{Q} its density has the form

$$p(\mathbf{x}) = \frac{|\mathbf{Q}|}{(2\pi)^{N/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \mathbf{Q}(\mathbf{x} - \mu)\right)$$

LINEAR DYNAMICAL SYSTEMS

$$\mathbf{x}_t = \mathbf{W}\mathbf{y}_t + \mathbf{e}_t$$

$$\mathbf{y}_1 = \mu_0 + \mathbf{u}$$

$$\mathbf{y}_t = \mathbf{A}\mathbf{y}_{t-1} + \mathbf{v}_t$$

$$\mathbf{e} \sim N(\mathbf{e}|\mathbf{0}, \Sigma)$$

$$\mathbf{u} \sim N(\mathbf{u}|\mu_0, \mathbf{P}_0)$$

$$\mathbf{v} \sim N(\mathbf{v}|\mathbf{0}, \Gamma)$$

where \mathbf{y}_i is the latent variable corresponding to \mathbf{x}_i . The parameters required to be determined is $\theta = \{\mathbf{W}, \mathbf{A}, \mu_0, \Sigma, \Gamma, \mathbf{P}_0\}$ Assuming all these variables follows Gaussian distribution, the process is described as follows.

$$p(\mathbf{y}_1) = N(\mathbf{y}_1|\mu_0, \mathbf{P}_0)$$

$$p(\mathbf{y}_t|\mathbf{y}_{t-1}) = N(\mathbf{y}_t|\mathbf{A}\mathbf{y}_{t-1}, \Gamma)$$

$$p(\mathbf{x}_t|\mathbf{y}_t) = N(\mathbf{x}_t|\mathbf{W}\mathbf{y}_t, \Sigma)$$

The difference between HMM and LDS: a Markov Chain with **discrete** latent variables ($\pi, \mathbf{A}, \mathbf{B}$) and a Markov Chain with **continuous** latent variables.