SVM summarized---- Only kernels

• Minimize with respect to w, w₀

$$\sum_{n=1}^{N} \zeta_{n} + \frac{1}{2} \| \boldsymbol{w} \|^{2}$$
 (Bishop 7.21)

Solution found in dual domain with Lagrange multipliers

- a_n , $n = 1 \cdots N$ and

• This gives the support vectors S

$$\widehat{\boldsymbol{w}} = \sum_{n \in S} a_n t_n \boldsymbol{\varphi}(xn)$$
 (Bishop 7.8)

• Used for predictions

$$\hat{y} = w_0 + \boldsymbol{w}^{\mathrm{T}} \boldsymbol{\varphi}(x) = w_0 + \sum_{n \in S} a_n t_n \boldsymbol{\varphi}(x_n)^{\mathrm{T}} \boldsymbol{\varphi}(x)$$
$$= w_0 + \sum_{n \in S} a_n t_n k(x_n, x) \qquad \text{(Bishop 7.13)}$$

Finding the Decision Function

- w: maybe infinite variables
- The dual problem

$$\begin{array}{ll} \min_{\boldsymbol{\alpha}} & \frac{1}{2} \boldsymbol{\alpha}^{T} Q \boldsymbol{\alpha} - \mathbf{e}^{T} \boldsymbol{\alpha} \\ \text{subject to} & 0 \leq \alpha_{i} \leq C, i = 1, \dots, I \\ \mathbf{y}^{T} \boldsymbol{\alpha} = 0, \end{array} \qquad \begin{array}{l} \text{Corresponds to} \\ \text{(Bishop 7.32)} \\ \text{Where } Q_{ij} = y_{i} y_{j} \phi(\mathbf{x}_{i})^{T} \phi(\mathbf{x}_{j}) \text{ and } \mathbf{e} = [1, \dots, 1]^{T} \end{array} \qquad \begin{array}{l} \text{With } \mathbf{y} = \mathbf{t} \\ \text{With } \mathbf{y} = \mathbf{t} \end{array}$$

$$\mathbf{w} = \sum_{i=1}^{l} \alpha_i \mathbf{y}_i \phi(\mathbf{x}_i)$$

• A finite problem: #variables = #training data

Using these results to eliminate w, b, and $\{\xi_n\}$ from the Lagrangian, we obtain the dual Lagrangian in the form

$$\widetilde{L}(\mathbf{a}) = \sum_{n=1}^{N} a_n - \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$
(7.32)

- May 8, CODY Machine Learning for finding oil, focusing on 1) robust seismic denoising/interpolation using structured matrix approximation 2) seismic image clustering and classification, using t-SNE(t-distributed stochastic neighbor embedding) and CNN. Weichang Li, Goup Leader Aramco, Houston.
- May 10, Class HW First distribution of final projects. Ocean acoustic source tracking. Final projects. Final project is the main goal in last month. Bishop Ch 9 Mixture models
- May 15, CODY Seismology and Machine Learning, Daniel Trugman (half class), ch 8
- May 17, Class HW ch 8
- May 22, Dictionary learning, Mike Bianco (half class), Graphical models Bishop Ch 8
- May 24, Graphical models Bishop Ch 8
- May 31, No Class. Workshop, <u>Big Data and The Earth Sciences: Grand Challenges</u> <u>Workshop</u>
- June 5, Discuss workshop, ch13. Spiess Hall open for project discussion 11am-.
- June 7, Workshop report. No class
- June 12 Spiess Hall open for project discussion 9-11:30am and 2-7pm
- June 16 Final report delivered. Beer time

For final project discussion every afternoon Mark and I will be available

Chapter 13 Sequential data

- Ocean source tracking X Final Report Re-implement Source Localization in an Ocean Waveguide using Supervised Machine Learning
- X-ray spectrum absorption interpretation using NN
- Neural decoding
- Plankton
- Transfer learning and deep feature ex
- Speaker tagger
- Coral
- Amazon rainforest (Kaggle)
- **Myshake Seismic**
- Please ask questions
 - Mark and I available all afternoons. Just come or email for time slots.
 - Spiess hall 330 is open Monday 5 and 12 June. If interested I can book it at other times
- Report
 - Rather concise than long.
 - Larger group can do more.
 - Start with some very simple example. To show your idea and that it is working.
 - End with showing the advanced abilities —
 - Several figures.
 - Equations are nice.
- Delivery Zip file (Friday 16)
 - Main code (not all). It should be able to run.
 - Report (pdf preferred).



Mixtures of Gaussians (1)

Old Faithful geyser:

The time between eruptions has a <u>bimodal distribution</u>, with the mean interval being either 65 or 91 minutes, and is dependent on the length of the prior eruption. Within a margin of error of ±10 minutes, Old Faithful will erupt either 65 minutes after an eruption lasting less than 2 $\frac{1}{2}$ minutes, or 91 minutes after an eruption lasting more than 2 $\frac{1}{2}$ minutes.



Mixtures of Gaussians (2)



Mixtures of Gaussians (3)



Mixture of Gaussians

• Mixtures of Gaussians

$$p(\mathbf{x}) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k).$$

Expressed with latent variable z

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{z}) p(\mathbf{x} | \mathbf{z}) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Posterior probability: responsibility

$$\gamma(z_k) \equiv p(z_k = 1 | \mathbf{x}) = \frac{p(z_k = 1)p(\mathbf{x} | z_k = 1)}{\sum_{j=1}^{K} p(z_j = 1)p(\mathbf{x} | z_j = 1)}$$
$$= \frac{\pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^{K} \pi_j \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}.$$
$$\mathbf{p}(\mathbf{z})\mathbf{p}(\mathbf{x} | \mathbf{z}) \quad \mathbf{N} \text{ iid } \{\mathbf{x}_n\} \text{ with latent } \{\mathbf{z}_n\}$$

EM Gauss Mix

- 1. Initialize the means μ_k , covariances Σ_k and mixing coefficients π_k , and evaluate the initial value of the log likelihood.
- 2. E step. Evaluate the responsibilities using the current parameter values

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}.$$
(9.23)

3. M step. Re-estimate the parameters using the current responsibilities

$$\boldsymbol{\mu}_{k}^{\text{new}} = \frac{1}{N_{k}} \sum_{n=1}^{N} \gamma(z_{nk}) \mathbf{x}_{n}$$
(9.24)

$$\boldsymbol{\Sigma}_{k}^{\text{new}} = \frac{1}{N_{k}} \sum_{n=1}^{N} \gamma(z_{nk}) \left(\mathbf{x}_{n} - \boldsymbol{\mu}_{k}^{\text{new}} \right) \left(\mathbf{x}_{n} - \boldsymbol{\mu}_{k}^{\text{new}} \right)^{\text{T}}$$
(9.25)

$$\pi_k^{\text{new}} = \frac{N_k}{N} \tag{9.26}$$

where

$$N_{k} = \sum_{n=1}^{N} \gamma(z_{nk}).$$
 (9.27)

4. Evaluate the log likelihood

$$\ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^{N} \ln \left\{ \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$
(9.28)

and check for convergence of either the parameters or the log likelihood. If the convergence criterion is not satisfied return to step 2.

General EM

Given a joint distribution $p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ over observed variables \mathbf{X} and latent variables \mathbf{Z} , governed by parameters $\boldsymbol{\theta}$, the goal is to maximize the likelihood function $p(\mathbf{X}|\boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$.

1. Choose an initial setting for the parameters θ^{old} .

- 2. E step Evaluate $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})$.
- 3. **M step** Evaluate θ^{new} given by

$$\boldsymbol{\theta}^{\text{new}} = \operatorname*{arg\,max}_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})$$
 (9.32)

where

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}).$$
(9.33)

4. Check for convergence of either the log likelihood or the parameter values. If the convergence criterion is not satisfied, then let

$$\boldsymbol{\theta}^{\mathrm{old}} \leftarrow \boldsymbol{\theta}^{\mathrm{new}}$$
 (9.34)

and return to step 2.

EM in general

$$p(\mathbf{X}|\boldsymbol{\theta}) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}).$$
(9.69)

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \mathcal{L}(q, \boldsymbol{\theta}) + \mathrm{KL}(q||p)$$
(9.70)

where we have defined

$$\mathcal{L}(q, \boldsymbol{\theta}) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})}{q(\mathbf{Z})} \right\}$$
(9.71)

$$\operatorname{KL}(q||p) = -\sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \right\}.$$
(9.72)

$$\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) = \ln p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}) + \ln p(\mathbf{X}|\boldsymbol{\theta})$$
(9.73)

$$\mathcal{L}(q, \theta) = \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z} | \theta) - \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}})$$
$$= \mathcal{Q}(\theta, \theta^{\text{old}}) + \text{const}$$
(9.74)

K-means

$$J = \sum_{n=1}^{N} \sum_{k=1}^{K} r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$
(9.1)

Solving for r_{nk}

$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg\min_{j} \|\mathbf{x}_{n} - \boldsymbol{\mu}_{j}\|^{2} \\ 0 & \text{otherwise.} \end{cases}$$
(9.2)

Differentiating for μ_k

$$2\sum_{n=1}^{N} r_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k) = 0$$
(9.3)

which we can easily solve for μ_k to give

$$\boldsymbol{\mu}_{k} = \frac{\sum_{n} r_{nk} \mathbf{x}_{n}}{\sum_{n} r_{nk}}.$$
(9.4)















0

2

-2

-2

Gaussian Mixtures













Mixture of Experts



Figure 11.6 (a) Some data fit with three separate regression lines. (b) Gating functions for three different "experts". (c) The conditionally weighted average of the three expert predictions. Figure generated by mixexpDemo.